

WILEY BIOTECHNOLOGY ENCYCLOPEDIA

# THE ENCYCLOPEDIA OF MOLECULAR BIOLOGY

VOLUME 1

EDITED BY  
THOMAS E. CREIGHTON

ENCYCLOPEDIA OF

---

# MOLECULAR BIOLOGY

---

VOLUMES 1 - 4

---

Thomas E. Creighton

*European Molecular Biology Laboratory  
London, England*



A Wiley-Interscience Publication

**John Wiley & Sons, Inc.**

New York / Chichester / Weinheim / Brisbane / Singapore / Toronto

This book is printed on acid-free paper. ☺

Copyright © 1999 by John Wiley & Sons, Inc.

All rights reserved. Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (508) 750-8400, fax (508) 750-4744. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 605 Third Avenue, New York, NY 10158-0012, (212) 850-6011, fax (212) 850-6008, E-Mail: PERMREQ@WILEY.COM.  
For ordering and customer service, call 1-800-CALL-WILEY.

*Library of Congress Cataloging-in-Publication Data:*

Creighton, Thomas E., 1940–

The encyclopedia of molecular biology / Thomas E. Creighton.

p. cm.

Includes index.

ISBN 0-471-15302-8 (alk. paper)

1. Molecular biology—Encyclopedias. I. Title.

QH506.C74 1999

572.8'03—dc21

99-11575

CIP

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

# PREFACE

---

The Wiley Biotechnology Encyclopedias, composed of the *Encyclopedia of Molecular Biology*, the *Encyclopedia of Bioprocess Technology: Fermentation, Biocatalysis, and Bioseparation*, the *Encyclopedia of Cell Technology*, and the *Encyclopedia of Ethical, Legal, and Policy Issues in Biotechnology*, cover very broadly four major contemporary themes in biotechnology. The series comes at a fascinating time in that as we move into the twenty-first century, the discipline of biotechnology is undergoing striking paradigm changes.

Biotechnology is now beginning to be viewed as an informational science. In a simplistic sense, there are three types of biological information. First, there is the digital or linear information of our chromosomes and genes, with the four-letter alphabet composed of G, C, A, and T (the bases Guanine, Cytosine, Adenine, and Thymine). Variation in the order of these letters in the digital strings of our chromosomes or our expressed genes (or mRNAs) generates information of several distinct types: genes, regulatory machinery, and information that enables chromosomes to carry out their tasks as informational organelles (eg, centromeric and telomeric sequences).

Second, there is the three-dimensional information of proteins, the molecular machines of life. Proteins are strings of amino acids employing a 20-letter alphabet. Proteins pose four technical challenges: (i) Proteins are synthesized as linear strings and fold into precise three-dimensional structures as dictated by the order of amino acid residues in the string. Can we formulate the rules for protein folding to predict three-dimensional structure from primary amino acid sequence? The identification and comparative analysis of all human and model organism (bacteria, yeast, nematode, fly, mouse, etc.) genes and proteins will eventually lead to a lexicon of motifs that are the building block components of genes and proteins. These motifs will greatly constrain the shape space computational algorithms must search to successfully correlate primary amino acid sequence with the correct three-dimensional shapes. The protein-folding problem will probably be solved within the next 10 to 15 years. (ii) Can we predict protein function from knowledge of the three-dimensional structure? Once again the lexicon of motifs with their functional as well as structural correlations will play a critical role in solving this problem. (iii) How do the myriad of chemical modifications of proteins (eg, phosphorylation, acetylation) alter their structures and modify their functions? The mass spectrometer will play a key role in identifying secondary modifications. (iv) How do proteins interact with one another and/or with other macromolecules to form complex molecular machines (eg, the ribosomal subunits)? If these functional complexes can be isolated, the mass

spectrometer, coupled with a knowledge of all protein sequences that can be derived from the complete genomic sequence of the organism, will serve as a powerful tool for identifying all the components of complex molecular machines.

The third type of biological information arises from complex biological systems and networks. Systems information is four-dimensional because it varies with time. For example, the human brain has 1012 neurons making approximately 1015 connections. From this network arises systems properties such a memory, consciousness, and the ability to learn. The important point is that systems properties cannot be understood from studying the network elements (eg, neurons) one at a time; rather, the collective behavior of the elements needs to be studied together. To study most biological systems, three issues need to be stressed. First, most biological systems are too complex to study directly; therefore they must be divided into tractable subsystems whose properties in part reflect those of the system. These subsystems must be sufficiently small to analyze all their elements and connections. Second, high-throughput analytic or global tools are required for studying many systems elements at one time (see below). Finally, the systems information needs to be modeled mathematically before systems properties can be predicted and ultimately understood. This will require recruiting computer scientists and applied mathematics into biology—just as the attempts to decipher the information of complete genomes and the protein folding and structure/function problems have required the recruitment of computational scientists.

I would be remiss not to point out that there are many other molecules that generate biological information—amino acids, carbohydrates, lipids, etc. These too must be studied in the context of their specific structures and specific functions.

The deciphering and manipulation of these various types of biological information represent an enormous technical challenge for biotechnology. Yet major new and powerful tools for doing so are emerging.

One class of tools for deciphering biological information is termed high-throughput analytic or global tools. These tools can study many genes or chromosome features (genomics), many proteins (proteomics), or many cells rapidly: large-scale DNA sequencing; genome-wide genetic mapping; cDNA or oligonucleotide arrays; two-dimensional gel electrophoresis and other global protein separation technologies; mass spectrometric analysis of proteins and protein fragments; multiparameter, high-throughput cell and chromosome sorting; and high-throughput phenotypic assays.

A second approach to the deciphering and manipulation of biological information centers around combinatorial strate-

gies. The basic idea is to synthesize an informational string (DNA fragments, RNA fragments, protein fragments, antibody combining sites, etc.) using all combinations of the basic letters of the corresponding alphabet—thus creating many different shapes that can be used to activate, inhibit, or complement the biological functions of designated three-dimensional shapes (eg, a molecule in a signal transduction pathway). The power of combinatorial chemistry is just beginning to be appreciated.

A critical approach to deciphering biological information will ultimately be the ability to visualize the functioning of genes, proteins, cells, and other informational elements within living organisms (*in vivo* informational imaging).

Finally, there are the computational tools required to collect, store, analyze, model, and ultimately distribute the various types of biological information. The creation presents a challenge comparable to that of developing of new instrumentation and new chemistries. Once again, this means recruiting computer scientists and applied mathematicians to biology. The biggest challenge in this regard is the language barriers that separate different scientific disciplines. Teaching biology as an informational science has been a very effective means for breaching these language barriers.

The challenge is, of course, to decipher these various types of biological information and then be able to use this information to manipulate genes, proteins, cells, and informational pathways in living organisms to eliminate or prevent disease, produce higher yield crops, or increase the productivity of animal products and meat.

Biotechnology and its applications raise a host of social, ethical, and legal questions; for example, genetic privacy, germline genetic engineering, cloning animals, genes that influence behavior, cost of therapeutic drugs generated by biotechnology, animal rights, and the nature and control of intellectual property.

The challenge clearly is to educate society so that each citizen can thoughtfully and rationally deal with these issues, for ultimately society dictates the resources and regulations that circumscribe the development and practice of biotechnology. Ultimately, I feel enormous responsibility rests with scientists to inform and educate society about the challenges as well as the opportunities arising from biotechnology. These are critical issues for biotechnology that are developed in detail in the *Encyclopedia of Ethical, Legal, and Policy Issues in Biotechnology*.

The view that biotechnology is an informational science pervades virtually every aspect of this science—including discovery, reduction to practice, and societal concerns. These Encyclopedias of Biotechnology reinforce the emerging informational paradigm change that is powerfully positioning science as we move into the twenty-first century to more effectively decipher and manipulate for humankind's benefit the biological information of relevant living organisms.

LEROY HOOD  
University of Washington

# CONTRIBUTORS

---

- Hanna E. Abboud**, *University of Texas, Health Science Center, San Antonio, TX*
- Sankar Adhya**, *National Cancer Institute, National Institutes of Health, Bethesda, MD*
- Hiroji Aiba**, *Nagoya University, Chikusa, Nagoya, Japan*
- Philip Aisen**, *Albert Einstein College of Medicine, Bronx, NY*
- Rudolf K. Allemann**, *University of Birmingham, Birmingham, United Kingdom*
- Nicholas Allen**, *The Babraham Institute, Cambridge, United Kingdom*
- Suresh Ambudkar**, *National Cancer Institute, National Institutes of Health, Bethesda, MD*
- Vernon E. Anderson**, *Case Western Reserve University, Cleveland, OH*
- Bertil Andersson**, *Stockholm University, Stockholm, Sweden*
- Ruth Hogue Angeletti**, *Albert Einstein College of Medicine, Bronx, NY*
- Rodolfo Aramayo**, *Schering-Plough Research Institute, Kenilworth, NJ*
- Yari Argon**, *University of Chicago, Chicago, IL*
- K. Arora**, *University of California, Irvine, CA*
- Leonie K. Ashman**, *Hanson Center for Cancer Research, Adelaide, Australia*
- John F. Atkins**, *University of Utah, Salt Lake City, UT*
- William M. Atkins**, *University of Washington, Seattle, WA*
- Daniel Atkinson**, *UCLA, Los Angeles, CA*
- David Auld**, *Harvard Medical School, Boston, MA*
- Paul Babitzke**, *Pennsylvania State University, University Park, PA*
- Andrew Baird**, *Prizm Pharmaceuticals, San Diego, CA*
- Stacey J. Baker**, *Temple University School of Medicine, Philadelphia, PA*
- Tom Baldwin**, *Texas A&M University, College Station, TX*
- Michael Bamshad**, *University of Utah Health Sciences Center, Salt Lake City, UT*
- Probal Banerjee**, *College of Staten Island, New York, NY*
- Ruma Banerjee**, *University of Nebraska, Lincoln, NE*
- William R. Bauer**, *University of California, San Francisco, CA*
- Christopher Baum**, *Heinrich-Pette Institut, Hamburg, Germany*
- Edward A. Bayer**, *Weizmann Institute of Science, Rehovot, Israel*
- Miguel Beato**, *Philipps–Universitaet Marburg, Institut fuer Molekularbiologie, Marburg, Germany*
- Dorothy Beckett**, *University of Maryland, Baltimore, MD*
- Samuel Benchimol**, *University of Toronto, Toronto, Ontario, Canada*
- Steven J. Benkovic**, *Pennsylvania State University, University Park, PA*
- Tomas Bergman**, *Karolinska Institutet, Stockholm, Sweden*
- Gerald Bergtrom**, *University of Wisconsin, Milwaukee, WI*
- Alan Bernstein**, *Mount Sinai Hospital, Toronto, Ontario, Canada*
- Harris D. Bernstein**, *NIDDKD, National Institutes of Health, Bethesda, MD*
- Sophie Bertrand**, *Universiteit Gent, Gent, Belgium*
- D. M. Bethea**, *Thomas Jefferson University, Philadelphia, PA*
- Dieter Beyer**, *Rhone Poulenc Rorer Recherche Developpement, Vitry-sur-Seine, France*
- Timothy R. Billiar**, *University of Pittsburgh Medical Center, Pittsburgh, PA*
- Asgeir Bjornsson**, *University of Aarhus, Aarhus, Denmark*
- Colin C. F. Blake**, *Norfolk, United Kingdom*
- F. Bonomi**, *Universita Degli Studi Di Milano, Milano, Italy*
- Ralph A. Bradshaw**, *University of California, Irvine, CA*
- Bertram Brenig**, *Georg-August University, Göttingen, Germany*
- Kenneth J. Breslauer**, *Rutgers University, Piscataway, NJ*
- Roy J. Britten**, *California Institute of Technology, Corona del Mar, CA*
- Maurizio Brunori**, *University of Rome, "La Sapienza," Rome, Italy*
- Bernd Bukau**, *Albert-Ludwigs Universitat Freiburg, Freiburg, Germany*
- Jens R. Bundgaard**, *University of Copenhagen, Copenhagen, Denmark*
- Arsene Burny**, *Universite Libre de Bruxelles, Rhode Saint Genese, Belgium*
- Kenneth Burtis**, *University of California, Davis, CA*
- Ana Busturia**, *Universidad Aut noma de Madrid, Madrid, Spain*
- Giulio L. Cantoni**, *NIMH, National Institutes of Health, Bethesda, MD*
- M. Stella Carlomango**, *Università Degli Studi Di Napoli Federico II, Napoli, Italy*
- Gerald M. Carlson**, *University of Missouri, Kansas City, MO*
- Graham Carpenter**, *Vanderbilt University, Nashville, TN*
- Robin W. Carrell**, *University of Cambridge, Cambridge, United Kingdom*
- James Castelli-Gair**, *University of Cambridge, Cambridge, United Kingdom*
- Enrique Cerda-Olmedo**, *Universidad de Sevilla, Sevilla, Spain*
- Jonathan Chaires**, *University of Mississippi Medical Center, Jackson, MS*
- Kung-Yao Chang**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*
- Lee Chao**, *Medical University of South Carolina, Charleston, SC*
- Ansuman Chattopadhyay**, *Vanderbilt University School of Medicine, Nashville, TN*
- Walter Chazin**, *Scripps Research Institute, La Jolla, CA*
- Elizabeth H. Chen**, *University of Texas, Southwestern Medical Center, Dallas, TX*
- Donald P. Cheney**, *Northeastern University, Boston, MA*
- Andreas Chrambach**, *National Institutes of Health, Bethesda, MD*
- Jason Christiansen**, *Virginia Polytechnic Institute, Blacksburg, VA*
- Jon A. Christopher**, *Texas A&M University, College Station, TX*
- Brian F. C. Clark**, *Aarhus University, Aarhus, Denmark*
- Dennis Clegg**, *University of California, Santa Barbara, CA*
- F. Clifton**, *University of Colorado Health Sciences Center, Denver, CO*
- Georges N. Cohen**, *Institut Pasteur, Paris, France*
- Roberta Colman**, *University of Delaware, Newark, DE*
- Orla M. Conneely**, *Baylor College of Medicine, Houston, TX*
- Barry S. Cooperman**, *University of Pennsylvania, Philadelphia, PA*
- Pamela Correll**, *The Pennsylvania State University, University Park, PA*
- Pascale Cossart**, *Institut Pasteur, Paris, France*
- Nancy Craig**, *Johns Hopkins University, Baltimore, MD*
- Elliott Croke**, *Georgetown University, Washington, DC*
- Stanley T. Croke**, *ISIS Pharmaceuticals, Inc., Carlsbad, CA*
- Richard D. Cummings**, *University of Oklahoma, Oklahoma City, OK*
- James Curran**, *Wake Forest University, Winston-Salem, NC*
- Michael A. Cusanovich**, *University of Arizona, Tucson, AZ*
- Giuseppe D'Alessio**, *Universita di Napoli Federico II, Napoli, Italy*
- Antoine Danchin**, *Institut Pasteur, Paris, France*

- P. L. Davies**, *Queen's University, Kingston, Ontario, Canada*  
**Dennis R. Dean**, *Virginia Polytechnic Institute & State University, Blacksburg, VA*  
**Pieter DeHaseth**, *Case Western Reserve University, Cleveland, OH*  
**Andrew J. DeMello**, *Imperial College, London, United Kingdom*  
**David T. Denhardt**, *Rutgers University, Piscataway, NJ*  
**W. A. Denny**, *University of Auckland, Auckland, New Zealand*  
**Zygmunt Derewenda**, *University of Virginia, Charlottesville, VA*  
**Claude Desplan**, *Rockefeller University, New York, NY*  
**Raymond Devoret**, *Centre Universitaire, Orsay, France*  
**Beth Didomenico**, *Schering-Plough Research Institute, Kenilworth, NJ*  
**Patrick DiMario**, *Louisiana State University, Baton Rouge, LA*  
**Hal B. F. Dixon**, *Cambridge University, Cambridge, United Kingdom*  
**Mark Dodson**, *University of Arizona, Tucson, AZ*  
**Jan Drenth**, *Laboratory of Biophysical Chemistry, Groningen, The Netherlands*  
**Gerhart Drews**, *Institut für Biologie, Freiburg, Germany*  
**Bernard Dujon**, *Institut Pasteur, Paris, France*  
**Edward Egelman**, *University of Minnesota, Minneapolis, MN*  
**Jean Marc Egly**, *Parc d'Innovation, Illkirch, France*  
**John Ellis**, *University of Warwick, Coventry, United Kingdom*  
**Vincent Ellis**, *Thrombosis Research Institute, London, United Kingdom*  
**Walther Ellis**, *Utah State University, Logan, UT*  
**Volker Erdmann**, *Freie Universität Berlin, Berlin, Germany*  
**Gerard Evan**, *ICRF, London, United Kingdom*  
**Kenneth Ewan**, *Max-Planck Institute of Biophysical Chemistry, Göttingen, Germany*  
**Hong Fang**, *Vanderbilt University, Nashville, TN*  
**Lynette R. Ferguson**, *University of Auckland Medical School, Auckland, New Zealand*  
**Carol A. Fierke**, *Duke University, Durham, NC*  
**J. R. S. Fincham**, *University of Edinburgh, Edinburgh, Scotland*  
**Anthony L. Fink**, *University of California, Santa Cruz, CA*  
**Gunter Fischer**, *MPG Arbeitsgruppe, Halle, Germany*  
**Darrell R. Fisher**, *Pacific Northwest Laboratory, Richland, WA*  
**Catherine Florentz**, *Institut de Biologie Moléculaire et Cellulaire, Strasbourg, France*  
**Ricardo Flores**, *Instituto de Biología Molecular y Celular de Plantas, Valencia, Spain*  
**P. J. Foley**, *Thomas Jefferson University, Philadelphia, PA*  
**Josiane Fontaine-Perus**, *CNRS, Nantes Cedex, France*  
**Jefferson Foote**, *Fred Hutchinson Cancer Research Center, Seattle, WA*  
**Patricia Foster**, *Boston University, Boston, MA*  
**Michel Fougereau**, *Centre d'Immunologie Marseille Luminy, Marseille, France*  
**Francois Franceschi**, *Max-Planck Institute for Molecular Genetics, Berlin, Germany*  
**Murray J. Fraser**, *Sydney Children's Hospital, Randwick, Australia*  
**Ian Freshney**, *University of Glasgow, Glasgow, United Kingdom*  
**Terrence G. Frey**, *San Diego State University, San Diego, CA*  
**Laura Frontali**, *University of Rome, "La Sapienza," Rome, Italy*  
**Yasuo Fukami**, *Kobe University, Nada, Kobe, Japan*  
**Carl W. Fuller**, *Amersham Pharmacia Biotech, Cleveland, OH*  
**Betty Gaffney**, *Florida State University, Tallahassee, FL*  
**Fernando Garcia-Arenal**, *Ciudad Universitaria, Madrid, Spain*  
**Jean-Renaud Garel**, *Laboratoire de Biologie Structurale, CNRS, Gif-sur-Yvette, France*  
**D. R. Garrod**, *University of Manchester, Manchester, United Kingdom*  
**Frank R. Gasparro**, *Jefferson University, Philadelphia, PA*  
**C. Gaspin**, *Centre National de la Recherche Scientifique, Strasbourg, France*  
**Anthony A. Gatenby**, *E. I. duPont de Nemours & Co., Wilmington, DE*  
**Kunihiko Gekko**, *Hiroshima University, Higashi-Hiroshima, Japan*  
**Herbert Geller**, *UMDNJ—Robert Wood Johnson Medical School, Piscataway, NJ*  
**John Gerig**, *University of California, Santa Barbara, CA*  
**Norman L. Gershfeld**, *National Institutes of Health, Bethesda, MD*  
**Michael Gershon**, *Columbia University, New York, NY*  
**Jonathan M. Gershoni**, *Tel Aviv University, Tel Aviv, Israel*  
**Jean-Marie Ghuysen**, *Université de Liège, Liège, Belgium*  
**Toby Gibson**, *European Molecular Biology Laboratory, Heidelberg, Germany*  
**J. A. Girault**, *Chaire de Neuropharmacologie, Collège de France, Paris, France*  
**Adam Godzik**, *Scripps Research Institute, La Jolla, CA*  
**Jeffrey Godzik**, *Scripps Research Institute, La Jolla, CA*  
**J. Peter Gogarten**, *University of Connecticut, Storrs Mansfield, CT*  
**Takashi Gojobori**, *National Institute of Genetics, Shizuoka, Japan*  
**Alfred L. Goldberg**, *Harvard Medical School, Boston, MA*  
**Paul Gollnick**, *State University of New York at Buffalo, Buffalo, NY*  
**Pierre Goldstein**, *Centre d'Immunologie INSERM-CNRS de Marseille, Marseille, France*  
**H. Maurice Goodman**, *University of Massachusetts Medical School, Worcester, MA*  
**Horace B. Gray**, *University of Houston, Houston, TX*  
**Neil Green**, *Vanderbilt University School of Medicine, Nashville, TN*  
**William Griffiths**, *Karolinska Institutet, Stockholm, Sweden*  
**Frank Grosse**, *Institut für Molekulare Biotechnologie, Jena, Germany*  
**Marianne Grunberg-Manago**, *Institut de Biologie Physico-Chimique, Paris, France*  
**Peter Gruss**, *Max-Planck Institute for Biophysical Chemistry, Göttingen, Germany*  
**Richard Gumport**, *University of Illinois, Urbana, IL*  
**Andras Guttman**, *Genetic Biosystems, San Diego, CA*  
**William Hagan**, *College of St. Rose, Albany, NY*  
**Klaus Hahlbrock**, *Max-Planck Institut für Biochemie, Martinsried, Germany*  
**Joshua W. Hamilton**, *Dartmouth Medical School, Hanover, NH*  
**Susan Hamilton**, *Baylor College of Medicine, Houston, TX*  
**Xianlin Han**, *Washington University School of Medicine, St. Louis, MO*  
**Ryo Hanai**, *Rikkyo University, Toshima, Tokyo, Japan*  
**Yusuf Hannun**, *Duke University Medical Center, Durham, NC*  
**Stephen E. Harding**, *University of Nottingham, Sutton Bonington, United Kingdom*  
**Peter Harland**, *University of California, Berkeley, CA*  
**Ulrich Hartl**, *Max-Planck Institut für Biochemie, Martinsried, Germany*  
**R. W. Hartley**, *National Institutes of Health, Bethesda, MD*  
**Enno Hartmann**, *Georg-August Universität, Göttingen, Germany*  
**David L. Haviland**, *The University of Texas, Houston Health Science Center, Houston, TX*  
**Toshihiko Hayashi**, *University of Tokyo, Tokyo, Japan*  
**J. K. Heath**, *University of Birmingham, Birmingham, United Kingdom*  
**Robert Helling**, *University of Michigan, Ann Arbor, MI*  
**Ali Hemmati-Brivanlou**, *Rockefeller University, New York, NY*  
**M. A. Hemminga**, *Agricultural University of Wageningen, Daagenaal, Wageningen, The Netherlands*  
**Beric Henderson**, *University of Sydney, Westmead, Australia*  
**Roger W. Hendrix**, *University of Pittsburgh, Pittsburgh, PA*  
**Leonard Herzenberg**, *Stanford University, Stanford, CA*  
**George P. Hess**, *Cornell University, Ithaca, NY*  
**Christa Heyting**, *Agricultural University, Wageningen, The Netherlands*  
**D. J. Hill**, *St. Joseph's Health Center, London, Ontario, Canada*  
**James A. Hoch**, *Scripps Research Institute, La Jolla, CA*  
**Denis Hochstrasser**, *Geneva University Hospital, Geneva, Switzerland*  
**Nicholas A. Hoenich**, *The University of Newcastle upon Tyne, Newcastle, United Kingdom*  
**Wayne L. Hoffman**, *University of Texas, Southwestern Medical Center, Dallas, TX*

- J. John Holbrook**, *University of Bristol, Bristol, United Kingdom*  
**Arne Holmgren**, *Karolinska Institute, Stockholm, Sweden*  
**Thomas F. Holzman**, *Abbot Laboratories, Abbot Park, IL*  
**H. Homareda**, *Kyoin University, Tokyo, Japan*  
**Roger Hull**, *John Innes Centre, Colney, Norwich, United Kingdom*  
**Martin Hum**, *University of Texas, Southwestern Medical Center, Dallas, TX*  
**Jennifer A. Hunt**, *Duke University Medical Center, Durham, NC*  
**Atsushi Ikai**, *Tokyo Institute of Technology, Yokohama, Japan*  
**Kiyohiro Imai**, *Osaka University, Osaka, Japan*  
**Taiji Imoto**, *Kyushu University, Fukuoka, Japan*  
**Antonio Incardona**, *University of Tennessee, Memphis, TN*  
**Edward E. Ishiguro**, *University of Victoria, Victoria, B.C., Canada*  
**Akira Ishihama**, *National Institute of Genetics, Shizuoka-ken, Japan*  
**Judith A. Jaehning**, *University of Colorado, Denver, CO*  
**Yuh-Nung Jan**, *University of California, San Francisco, CA*  
**Joel Janin**, *Laboratoire de Biologie Structurale, Paris, France*  
**S. Jansen**, *Georg-August University Göttingen, Göttingen, Germany*  
**Patricia A. Jennings**, *University of California, San Diego, La Jolla, CA*  
**Randy L. Jirtle**, *Duke University, Durham, NC*  
**Jonathan Jones**, *John Innes Institute, Norwich, United Kingdom*  
**Rick Jones**, *Southern Methodist University, Dallas, TX*  
**Hans Jornvall**, *Karolinska Institutet, Stockholm, Sweden*  
**Akia Kaji**, *University of Pennsylvania, Philadelphia, PA*  
**Barton A. Kamen**, *University of Texas, Southwestern Medical Center, Dallas, TX*  
**Heather Kaminski**, *Schering-Plough Research Institute, Kenilworth, NJ*  
**Minoru Kanehisa**, *Kyoto University, Kyoto, Japan*  
**Jack Kaplan**, *Oregon Health Sciences University, Portland, OR*  
**Francois Karch**, *University of Geneva, Geneva, Switzerland*  
**Thom Kaufmann**, *Indiana University, Bloomington, IN*  
**Jack D. Keene**, *Duke University, Durham, NC*  
**Regis B. Kelly**, *University of California, San Francisco, CA*  
**Zvi Kelman**, *Memorial Sloan-Kettering Cancer Center, New York, NY*  
**Byron Kemper**, *University of Illinois, Urbana, IL*  
**James Kennison**, *National Institute of Child Health and Human Development, National Institutes of Health, Bethesda, MD*  
**Kathleen M. Kerr**, *National Institutes of Health, Bethesda, MD*  
**John R. Kirby**, *University of California, Berkeley, CA*  
**Jennifer Kitchen**, *Georgetown University, Washington, DC*  
**Horst Kleinkauf**, *Technische Universität, Berlin, Germany*  
**Jorg Klug**, *Phillips-Universität Marburg, Marburg, Germany*  
**Peter Knight**, *University of Leeds, Leeds, United Kingdom*  
**Charlotte Knudsen**, *University of Aarhus, Aarhus, Denmark*  
**Yuji Kobayashi**, *Osaka University, Osaka, Japan*  
**Ralf Koebnik**, *Biozentrum Basel, Basel, Switzerland*  
**Andrzej Kolinski**, *Scripps Research Institute, La Jolla, CA*  
**Gerald B. Koudelka**, *State University of New York at Buffalo, Buffalo, NY*  
**Joseph Kraut**, *University of California, San Diego, La Jolla, CA*  
**Toshio Kuroki**, *Showa University, Tokyo, Japan*  
**Sidney R. Kushner**, *University of Georgia, Athens, GA*  
**Kunihiko Kuwajima**, *University of Tokyo, Tokyo, Japan*  
**Harold Lahm**, *Institute of Molecular Animal Breeding, Munich, Germany*  
**Ratnesh Lal**, *University of California, Santa Barbara, CA*  
**Marc Lalonde**, *Children's Hospital, Boston, MA*  
**David Lane**, *Centre National de la Recherche Scientifique, Strasbourg, France*  
**Janos Lanyi**, *University of California, Irvine, CA*  
**Paul Lasko**, *McGill University, Montreal, Quebec, Canada*  
**Michael Laskowski Jr.**, *Purdue University, West Lafayette, IN*  
**David S. Latchman**, *University College London Medical School, London, United Kingdom*  
**Peter Lawrence**, *Medical Research Council, Cambridge, United Kingdom*  
**Judith E. Layton**, *Ludwig Institute for Cancer Research, Victoria, Australia*  
**Claude J. Lazdunski**, *Laboratoire d'Ingeniere des Systèmes Macromoléculaires, Marseille, France*  
**Robin Leake**, *University of Glasgow, Glasgow, United Kingdom*  
**Stewart H. Lecker**, *Harvard Medical School, Boston, MA*  
**Stuart LeGrice**, *Case Western Reserve School of Medicine, Cleveland, OH*  
**Ruth Lehmann**, *New York University, New York, NY*  
**Arthur M. Lesk**, *University of Cambridge, Cambridge, United Kingdom*  
**Michael Levine**, *University of California, Berkeley, Berkeley, CA*  
**Joe D. Lewis**, *Institute of Cell & Molecular Biology, Edinburgh, United Kingdom*  
**Uno Lindberg**, *Stockholm University, Stockholm, Sweden*  
**Halina Lis**, *Weizmann Institute of Science, Rehovot, Israel*  
**Uriel Littauer**, *Weizmann Institute of Science, Rehovot, Israel*  
**John Little**, *University of Arizona, Tucson, AZ*  
**Lawrence A. Loeb**, *University of Washington, Seattle, WA*  
**Peter C. Loewen**, *University of Manitoba, Winnipeg, Manitoba, Canada*  
**Douglas J. Loftus**, *National Cancer Institute, National Institutes of Health, Bethesda, MD*  
**Kenton Longenecker**, *University of Virginia Health Sciences Center, Charlottesville, VA*  
**J. Michael Lord**, *University of Warwick, Coventry, United Kingdom*  
**Martin G. Low**, *Columbia University, New York, NY*  
**Reinhard Luhrmann**, *Institut für Molekularbiologie und Tumorforschung, Marburg, Germany*  
**Paul MacDonald**, *Stanford University, Stanford, CA*  
**Robert M. Macnab**, *Yale University, New Haven, CT*  
**Neil B. Madsen**, *University of Alberta, Edmonton, Alberta, Canada*  
**Hasaji Maki**, *Ikoma City, Japan*  
**Gregory S. Makowski**, *University of Connecticut, Farmington, CT*  
**Jack Maniloff**, *University of Rochester, Rochester, NY*  
**Jeffrey R. Mann**, *Beckman Research Institute of the City of Hope, Duarte, CA*  
**Bengt Mannervik**, *Uppsala University, Uppsala, Sweden*  
**Armen Manoukian**, *Ontario Cancer Institute, Toronto, Ontario, Canada*  
**Ahmed Mansouri**, *Max-Planck Institute for Biophysical Chemistry, Göttingen, Germany*  
**Mohamed A. Marahiel**, *Philipps-Universität Marburg, Marburg, Germany*  
**Guglielmo Marin**, *Università Degli Studi di Padova, Padova, Italy*  
**Susan Marqusee**, *University of California, Berkeley, CA*  
**Mike Marsh**, *University College London, London, United Kingdom*  
**Garland Marshall**, *Washington University, St. Louis, MO*  
**Jennifer Martin**, *University of Queensland, Brisbane, Australia*  
**Nancy C. Martin**, *University of Louisville, Louisville, KY*  
**Florence Maschat**, *Institut de Genetique Humaine, Montpellier, France*  
**Christopher Mathews**, *Oregon State University, Corvallis, OR*  
**Abdul Matin**, *Stanford University, Stanford, CA*  
**Mark P. Mattson**, *University of Kentucky, Lexington, KY*  
**Russell A. Maurer**, *Case Western Reserve University, Cleveland, OH*  
**Nicola McCarthy**, *Imperial Cancer Research Fund, London, United Kingdom*  
**Richard McCarty**, *Johns Hopkins University, Baltimore, MD*  
**David R. McClay**, *Duke University, Durham, NC*  
**John R. Menninger**, *University of Iowa, Iowa City, IA*  
**Carl R. Merrill**, *National Institute of Mental Health, National Institutes of Health, Bethesda, MD*  
**David Metzler**, *Iowa State University, Ames, IA*  
**Robert H. Mitchell**, *The University of Birmingham, Birmingham, United Kingdom*  
**Marek Mlodzik**, *European Molecular Biology Laboratory, Heidelberg, Germany*  
**Evita Mohr**, *Universität Hamburg, Hamburg, Germany*  
**Ian J. Molineux**, *University of Texas, Austin, TX*



- Cesare Montecucco**, *Università degli Studi di Padova, Padova, Italy*  
**Carol Moore**, *City University of New York Medical School, NY*  
**Henning D. Mootz**, *Phillips–Universität Marburg, Marburg, Germany*  
**Gines Morata**, *University Autonoma De Madrid, Madrid, Spain*  
**Richard I. Morimoto**, *Northwestern University, Evanston, IL*  
**John F. Morrison**, *The Australian National University, Canberra, Australia*  
**Gisela Mosig**, *Vanderbilt University, Nashville, TN*  
**Jan Mous**, *F. Hoffmann–LaRoche Ltd., Basel, Switzerland*  
**Benno Müller-Hill**, *Institut für Genetik / Universität Köln, Köln, Germany*  
**J. K. Myers**, *Texas A&M University, College Station, TX*  
**Kyoshi Nagai**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*  
**Haruki Nakamura**, *Biomolecular Engineering Research Institute, Osaka, Japan*  
**Yoshikazu Nakamura**, *University of Tokyo, Tokyo, Japan*  
**Aiguo Ni**, *Medical University of South Carolina, Charleston, SC*  
**Allen W. Nicholson**, *Wayne State University, Detroit, MI*  
**Christof Niehrs**, *Deutsches Krebsforschungszentrum, Heidelberg, Germany*  
**Koich Nishigaki**, *Saitama University, Urawa, Japan*  
**Ken Nishikawa**, *National Institute of Genetics, Shizuoka, Japan*  
**Ralph A. Nixon**, *Harvard Medical School, Belmont, MA*  
**Akio Nomoto**, *University of Tokyo, Tokyo, Japan*  
**Bengt Norden**, *Chalmers University of Technology, Gothenburg, Switzerland*  
**Shiao Li Oei**, *Institute für Biochemie, Berlin, Germany*  
**Naotake Ogasawara**, *Ikoma City, Japan*  
**Yoshio Okada**, *Osaka University, Osaka, Japan*  
**Charles Ordahl**, *University of California, San Francisco, CA*  
**George Ordal**, *University of Illinois, Urbana, IL*  
**George Oster**, *University of California, Berkeley, CA*  
**Malcolm G. P. Page**, *F. Hoffmann–LaRoche Ltd., Basel, Switzerland*  
**Renato Paro**, *University of Heidelberg, Heidelberg, Germany*  
**D. A. Parry**, *Massey University, Palmerston North, New Zealand*  
**Premal H. Patel**, *University of Washington School of Medicine, Seattle, WA*  
**Barbara Pearce**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*  
**Israel Pecht**, *Weizmann Institute of Science, Rehovot, Israel*  
**Iain K. Pemberton**, *Institut Pasteur, Paris, France*  
**Bernard Perbal**, *Centre Universitaire, Orsay, France*  
**Richard N. Perham**, *University of Cambridge, Cambridge, United Kingdom*  
**Francine Perler**, *New England Biolabs, Beverly, MA*  
**Charles L. Perrin**, *University of California, San Diego, La Jolla, CA*  
**Steven Perrin**, *ARIAD Pharmaceuticals, Cambridge, MA*  
**Bengt Persson**, *Karolinska Institute, Stockholm, Sweden*  
**Donald W. Pettigrew**, *Texas A&M University, College Station, TX*  
**Larry Pfeffer**, *University of Tennessee, Memphis, TN*  
**K. Kevin Pfister**, *University of Virginia, Charlottesville, VA*  
**Margaret A. Phillips**, *University of Texas, Dallas, TX*  
**Don Phillips**, *LaTrobe University, Victoria, Australia*  
**Thomas L. Poulos**, *University of California, Irvine, CA*  
**Linda S. Powers**, *Utah State University, Logan, UT*  
**Peter E. Prevelige**, *University of Alabama, Birmingham, AL*  
**Peter L. Privalov**, *Johns Hopkins University, Baltimore, MD*  
**Gudrun Rappold**, *Institut für Humangenetik Universitätsklinikum, Heidelberg, Germany*  
**Stephen W. Raso**, *Texas A&M University, College Station, TX*  
**P. D. Rathjen**, *University of Adelaide, Adelaide, Australia*  
**E. Prekumar Reddy**, *The Fels Institute for Cancer Research and Molecular Biology, Philadelphia, PA*  
**Colin Reese**, *King's College London, University of London, United Kingdom*  
**Peter Reichard**, *Medical Nobel Institute, Karolinska Institutet, Stockholm, Sweden*  
**Steffen Reinbothe**, *Universite Joseph Fourier, Grenoble, France*  
**Dietmar Richter**, *Institut für Zellbiochemie und klinische Neurobiologie, Hamburg, Germany*  
**P. G. Righetti**, *University of Milan, Milan, Italy*  
**Richard J. Riopelle**, *Queen's University, Kingston, Ontario, Canada*  
**James Riordan**, *Harvard University, Boston, MA*  
**Anthony J. Robertson**, *National Institutes of Health, Research Triangle Park, NC*  
**Jean-David Rochaix**, *University of Geneva, Geneva, Switzerland*  
**Mario Roederer**, *Stanford University, Stanford, CA*  
**William A. Rosche**, *Boston University School of Medicine, Boston, MA*  
**Elliott M. Ross**, *University of Texas, Southwestern Medical Center, Dallas, TX*  
**Gregory M. Ross**, *Queen's University, Kingston, Ontario, Canada*  
**Lawrence Rothfield**, *University of Connecticut, Farmington, CT*  
**Harry Rubin**, *University of California, Berkeley, CA*  
**Thomas S. Rush**, *Princeton University, Princeton, NJ*  
**M. S. Saedi**, *Hybritech Co., San Diego, CA*  
**Alan Saltiel**, *Parke-Davis Pharmaceutical Research, Ann Arbor, MI*  
**Aziz Sancar**, *University of North Carolina, Chapel Hill, NC*  
**Matti Saraste**, *European Molecular Biology Laboratory, Heidelberg, Germany*  
**Zuben E. Sauna**, *National Institutes of Health, Bethesda, MD*  
**Lindsay Sawyer**, *University of Edinburgh, Edinburgh, United Kingdom*  
**Walter Schaffner**, *Universität Zurich, Zurich, Switzerland*  
**Paul Schimmel**, *Scripps Research Institute, La Jolla, CA*  
**Franz Schmid**, *Universität Bayreuth, Bayreuth, Germany*  
**Manfred Schnarr**, *Institut de Biologie Moléculaire et Cellulaire, CNRS, Strasbourg, France*  
**J. Martin Scholtz**, *Texas A&M University, College Station, TX*  
**Trudi Schupbach**, *Princeton University, Princeton, NJ*  
**James Scott**, *Royal Postgraduate Medical School, London, United Kingdom*  
**Barbara A. Seaton**, *Boston University, Boston, MA*  
**Constantin E. Sekeris**, *University of Würzburg, Würzburg, Germany*  
**Igor L. Shamovsky**, *Queen's University, Kingston, Ontario, Canada*  
**Shmuel Shaltiel**, *Weizmann Institute of Science, Rehovot, Israel*  
**Nathan Sharon**, *Weizmann Institute of Science, Rehovot, Israel*  
**G. Shaw**, *University of Florida, Gainesville, FL*  
**William V. Shaw**, *Leicester University, Leicester, United Kingdom*  
**David Sherman**, *University of Minnesota, St. Paul, MN*  
**Ben-Zion Shilo**, *Weizmann Institute of Science, Rehovot, Israel*  
**Makoto Shimoyama**, *Shimane Medical University, Izumo, Japan*  
**Israel Silman**, *Weizmann Institute of Science, Rehovot, Israel*  
**Richard J. Simpson**, *Ludwig Institute for Cancer Research, Victoria, Australia*  
**Gary Siuzdak**, *Scripps Research Institute, La Jolla, CA*  
**James Skeath**, *Washington University, St. Louis, MO*  
**Stephen J. Small**, *New York University, New York, NY*  
**Colleen M. Smith**, *Mercer University, Macon, GA*  
**Gerald R. Smith**, *Fred Hutchinson Cancer Research Center, Seattle, WA*  
**Sarah M. Smolik**, *Oregon Health Sciences University, Portland, OR*  
**Suzanne G. Sobel**, *Yale University, New Haven, CT*  
**Kunitsugu Soda**, *Nagaoka University of Technology, Niigata, Japan*  
**Lila Solnica-Krezel**, *Vanderbilt University, Nashville, TN*  
**Mark Solomon**, *Yale University, New Haven, CT*  
**Joseph Sperling**, *The Hebrew University of Jerusalem, Jerusalem, Israel*  
**Ruth Sperling**, *The Hebrew University of Jerusalem, Jerusalem, Israel*  
**Martin Spiess**, *Biozentrum, Basel, Switzerland*  
**Tom Spiro**, *Princeton University, Princeton, NJ*  
**Mathias Sprinzl**, *Universität Bayreuth, Bayreuth, Germany*  
**R. B. Spruijt**, *Wageningen Agricultural University, Wageningen, The Netherlands*  
**F. W. Stahl**, *University of Oregon, Eugene, OR*  
**Martin F. Steiger**, *F. Hoffmann–LaRoche Ltd., Basel, Switzerland*

- Scott Strobel**, *Yale University, New Haven, CT*  
**Shintaro Sugai**, *Soka University, Tokyo, Japan*  
**Roy Sundick**, *Wayne State University, Detroit, MI*  
**H. Eldon Sutton**, *University of Texas, Austin, TX*  
**Sandra S. Szegedi**, *University of Illinois at Urbana-Champaign, Urbana, IL*  
**Sevec Szmelcman**, *Institut Pasteur, Paris, France*  
**Hatsumi Taniguchi**, *University of Occupational and Environmental Health, Fukuoka, Japan*  
**Joyce Taylor-Papadimitriou**, *Imperial Cancer Research Fund Laboratories, London, United Kingdom*  
**Mariella Tegoni**, *CNRS, Marseille, France*  
**Jean Thomas**, *Cambridge University, Cambridge, United Kingdom*  
**Carl Thummel**, *University of Utah, Salt Lake City, UT*  
**Serge Timasheff**, *Brandeis University, Waltham, MA*  
**Thea D. Tlsty**, *University of California, San Francisco, CA*  
**Andrew Travers**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*  
**Jill Trehwella**, *Los Alamos National Laboratory, Los Alamos, NM*  
**Edward Trifonov**, *National Institute of Genetics, Shizuoka-ken, Japan*  
**Toshiki Tsurimoto**, *Osaka, Japan*  
**Agnes Ullmann**, *Institut Pasteur, Paris, France*  
**Dominique van der Straeten**, *Universiteit Gent, Gent, Belgium*  
**M. van Lijsebettens**, *Universiteit Gent, Gent, Belgium*  
**Marc Van Montagu**, *Universiteit Gent, Gent, Belgium*  
**Klaas van Wijk**, *Stöckholm University, Stockholm, Sweden*  
**G. Varani**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*  
**Serge N. Vinogradov**, *Wayne State University, Detroit, MI*  
**Hans von Dohren**, *Technische Universität Berlin, Berlin, Germany*  
**Edward Voss**, *University of Illinois, Urbana, IL*  
**Richard Walden**, *Horticulture Research International, Kent, United Kingdom*  
**Melinda Wales**, *Texas A&M University, College Station, TX*  
**Valeria A. Wallace**, *University College London, London, United Kingdom*  
**Frederick G. Walz**, *Kent State University, Kent, OH*  
**Andrew H. J. Wang**, *University of Illinois, Urbana, IL*  
**Hongyun Wang**, *University of California, Berkeley, CA*  
**James C. Wang**, *Harvard University, Cambridge, MA*  
**Michael J. Waring**, *Cambridge University, Cambridge, United Kingdom*  
**Amy Warren**, *Dartmouth College, Hanover, NH*  
**Graham Warren**, *Imperial College Research Fund, London, United Kingdom*  
**Hans Warrick**, *Stanford University, Stanford, CA*  
**Arieh Warshel**, *University of Southern California, Los Angeles, CA*  
**Steven A. Wasserman**, *University of Texas, Southwestern Medical Center, Dallas, TX*  
**Hans Weber**, *Institut für Molekularbiologie, Zurich, Switzerland*  
**Patricia C. Weber**, *Schering-Plough Research Institute, Kenilworth, NJ*  
**Eric Westhof**, *Institut de Biologie Moléculaire et Cellulaire, du CNRS, Strasbourg, France*  
**Richard A. Wetsel**, *University of Texas, Houston, TX*  
**Paul Whiting**, *Merck Sharp and Dohme Research Laboratories, Essex, United Kingdom*  
**Peter A. Whittaker**, *Saint Patrick's College, Maynooth, Kildare, Ireland*  
**Pernilla Whitting-Stapshede**, *California Institute of Technology, Pasadena, CA*  
**Meir Wilchek**, *Weizmann Institute of Science, Rehovot, Israel*  
**Jim R. Wild**, *Texas A&M University, College Station, TX*  
**Allison K. Wilson**, *University of Wisconsin, Madison, WI*  
**David B. Wilson**, *Cornell University, Ithaca, NY*  
**Leslie Wilson**, *University of California, Santa Barbara, CA*  
**Sam H. Wilson**, *NIEHS, National Institutes of Health, Research Triangle Park, NC*  
**R. Wolfenden**, *University of North Carolina, Chapel Hill, NC*  
**Alan Paul Wolffe**, *NICDH, National Institutes of Health, Bethesda, MD*  
**C. J. A. M. Wolfs**, *Wageningen Agricultural University, Wageningen, The Netherlands*  
**William B. Wood**, *University of Colorado, Boulder, CO*  
**Robert Woody**, *Colorado State University, Fort Collins, CO*  
**Kenneth Yamada**, *NIDR, National Institutes of Health, Bethesda, MD*  
**Charles Yanofsky**, *Stanford University, Stanford, CA*  
**Ada Yonath**, *Weizmann Institute of Science, Rehovot, Israel*  
**Hiroshi Yoshikawa**, *Nara Institute of Science & Technology, Ikoma, Japan*  
**Jennifer Zamanian**, *University of California, San Francisco, CA*  
**Guiliana Zanetti**, *University of Milan, Milan, Italy*  
**Amanda Zeffman**, *MRC Laboratory of Molecular Biology, Cambridge, United Kingdom*  
**Suisheng Zhang**, *Institute of Molecular Biotechnology, Jena, Germany*  
**Mark Zoller**, *Ariad Pharmaceuticals, Cambridge, MA*

# COMMONLY USED ACRONYMS AND ABBREVIATIONS

---

|          |  |           |  |
|----------|--|-----------|--|
| Å        | angstrom ( $10^{-10}$ m)                   | G         | Gibbs free energy of a system  |
| Ab       | antibody                                   | G         | guanine base   |
| Ac       | acetyl group                               | G-protein | guanine-nucleotide binding regulatory protein                            |
| ADP      | adenosine diphosphate                      | GABA      | <i>g</i> -aminobutyric acid  |
| Ala      | alanine residue (also A)                   | GalNAc    | <i>N</i> -acetylgalactosamine residue                                    |
| AMP      | adenosine monophosphate                    | GdmCl     | guanidinium chloride (guanidine hydrochloride)                           |
| Arg      | arginine residue (also R)                  | GDP       | guanosine diphosphate  |
| Asn      | asparagine residue (also N)                | Glc       | glucose residue  |
| Asp      | aspartic acid residue (also D)             | GlcNAc    | <i>N</i> -acetylglucosamine residue                                      |
| ATP      | adenosine triphosphate                     | Gln       | glutamine residue (also Q)   |
| ATPase   | adenosine triphosphatase                   | Glu       | glutamic acid residue (also E)   |
| Bq       | becquerel                                  | Gly       | glycine residue (also G)   |
| bp       | base-pair                                  | GMP       | guanoine monophosphate   |
| BSA      | bovine serum albumin                       | GSH       | glutathione, thiol form  |
| C        | cytosine base                              | GSSG      | glutathione, disulfide form  |
| cal      | calorie (4.18 J)                           | GTP       | guanosine triphosphate   |
| cAMP     | 3',5'-cyclic AMP                           | GTPase    | guanosine triphosphatase   |
| CD       | circular dichroism                         | H         | enthalpy of a system   |
| CoA      | coenzyme A                                 | h         | hour   |
| cDNA     | complementary DNA                          | His       | histidine residue (also H)   |
| cGMP     | 3',5'-cyclic GMP                           | HLA       | histocompatibility locus antigen   |
| Cmc      | critical micelle concentration             | HPLC      | high pressure liquid chromatography                                      |
| Cys      | cysteine residue (also C)                  | Hz        | hertz, frequency   |
| Da       | dalton                                     | IEF       | isoelectric focusing IgG, IgA, IgM, etc.<br>immunoglobulin G, A, M. etc. |
| DMSO     | dimethyl sulfoxide                         | Ile       | isoleucine residue (also I)  |
| DNA      | deoxyribonucleic acid                      | Ins       | inositol residue   |
| DNase    | deoxyribonuclease                          | IR        | infrared   |
| DTT      | dithiothreitol                             | J         | joule  |
| EDTA     | ethylenediamine tetraacetic acid           | K         | degrees Kelvin, absolute temperature                                     |
| EG       | [ethylenebis(oxonitrilo)] tetraacetic acid | k         | rate constant for a specified reaction                                   |
| EGF      | epidermal growth factor                    | $K_a$     | association constant   |
| ELISA    | enzyme-linked immunosorbent assay          | $K_d$     | dissociation constant  |
| EM       | electron microscopy                        | $K_{eq}$  | equilibrium constant for a specified reaction                            |
| EPR, ESR | electron paramagnetic (or sign) resonance  | $K_m$     | Michaelis constant   |
| ER       | endoplasmic reticulum                      | kb        | kilobases  |
| EXAFS    | extended X-ray absorption fine structure   | kcal      | kilocalorie (4.18 kJ)  |
| FAD      | flavin-adenine dinucleotide                | kDa       | kilodalton   |
| fMet     | <i>N</i> -formyl methionine                | LDL       | low density lipoprotein  |
| FMN      | flavin mononucleotide                      | Leu       | leucine residue (also L)   |
| FPLC     | fast protein liquid chromatography         | Lys       | lysine residue (also K)  |

|                   |  |            |   |
|-------------------|--|------------|---|
| m                 | meter  | RNase      | ribonuclease                                |
| mAb               | monoclonal antibody  | RNP        | ribonucleoprotein                           |
| Met               | methionine residue (also M)  | rRNA       | ribosomal RNA                               |
| MHC               | major histocompatibility locus                                     |            |   |
| min               | minute   | S          | entropy of a system                         |
| mol               | mole   | S          | Svedberg sedimentation unit ( $10^{-13}$ s) |
| $M_r$             | relative molecular mass  | s          | second                                      |
| mRNA              | messenger RNA  | s          | sedimentation coefficient                   |
| MS                | mass spectroscopy  | SDS        | sodium dodecyl sulfate                      |
| mtDNA             | mitochondrial DNA  | snRNA      | small nuclear RNA                           |
| Mur               | muramic acid residue   | SRP        | signal recognition particle                 |
| MurNAc            | <i>N</i> -acetylmuramic acid residue                               |            |   |
|                   |  | T          | temperature                                 |
| NAD <sup>+</sup>  | nicotinamide adenine dinucleotide (oxidized)                       | T          | thymine residue                             |
| NADH              | nicotinamide adenine dinucleotide (reduced)                        | TF         | transcription factor                        |
| NADP <sup>+</sup> | nicotinamide adenine dinucleotide phosphate (oxidized)             | Thr        | threonine residue (also T)                  |
| NADPH             | nicotinamide adenine dinucleotide phosphate (reduced)              | $T_m$      | midpoint of temperature induced transition  |
| NMR               | nuclear magnetic resonance   | Tri        | tris(hydroxymethyl) aminomethane            |
|                   |  | tRNA       | transfer RNA                                |
| PAGE              | polyacrylamide gel electrophoresis                                 | Trp        | tryptophan residue (also W)                 |
| PBS               | phosphate-buffered saline  | TTP        | thymidine triphosphate                      |
| PCR               | polymerase chain reaction  | Tyr        | tyrosine residue (also Y)                   |
| PEG               | polyethylene glycol  |            |   |
| pH                | negative logarithm of hydrogen ion concentration                   | U          | uracil residue                              |
| Phe               | phenylalanine residue (also F)                                     | UDP        | uridine diphosphate                         |
| $P_i$             | inorganic phosphate  | UMP        | uridine monophosphate                       |
| $PP_i$            | inorganic pyrophosphate  | UTP        | uridine triphosphate                        |
| p.p.m.            | parts per million  | UV         | ultraviolet                                 |
| Pro               | proline residue (also P)   |            |   |
|                   |  | V          | volume                                      |
| R                 | gas constant ( $8.21451 \text{ J}/(\text{K}^1 \text{ mol}^{-1})$ ) | $V_{\max}$ | maximum enzyme catalyzed velocity           |
| RIA               | radioimmunoassay   | Val        | valine residue (also V)                     |
| RNA               | ribonucleic acid   |            |   |
|                   |  | WT         | wild-type                                   |
|                   |  | Xyl        | xylose residue                              |

# CONVERSION FACTORS

---

## SI Units (Adopted 1960)

The International System of Units (abbreviated SI), is being implemented throughout the world. This measurement system is a modernized version of the MKSA (meter, kilogram, second, ampere) system, and its details are published and controlled by an international treaty organization (The International Bureau of Weights and Measures) (1).

SI units are divided into three classes:

### BASE UNITS

|  |                        |
|--|------------------------|
| length                                 | meter <sup>†</sup> (m) |
| mass                                   | kilogram (kg)          |
| time                                   | second (s)             |
| electric current                       | ampere (A)             |
| thermodynamic temperature <sup>‡</sup> | kelvin (K)             |
| amount of substance                    | mole (mol)             |
| luminous intensity                     | candela (cd)           |

### SUPPLEMENTARY UNITS

|             |                |
|-------------|----------------|
| plane angle | radian (rad)   |
| solid angle | steradian (sr) |

### DERIVED UNITS AND OTHER ACCEPTABLE UNITS

These units are formed by combining base units, supplementary units, and other derived units (2–4). Those derived units having special names and symbols are marked with an asterisk in the list below.

| Quantity                               | Unit                     | Symbol             | Acceptable equivalent   |
|--|--------------------------|--------------------|-------------------------|
| *absorbed dose                         | gray                     | Gy                 | J/kg                    |
| acceleration                           | meter per second squared | m/s <sup>2</sup>   |                         |
| *activity (of a radionuclide)          | becquerel                | Bq                 | 1/s                     |
| area                                   | square kilometer         | km <sup>2</sup>    |                         |
|  | square hectometer        | hm <sup>2</sup>    | ha (hectare)            |
|  | square meter             | m <sup>2</sup>     |                         |
| concentration (of amount of substance) | mole per cubic meter     | mol/m <sup>3</sup> |                         |
| current density                        | ampere per square meter  | A/m <sup>2</sup>   |                         |
| density, mass density                  | kilogram per cubic meter | kg/m <sup>3</sup>  | g/L; mg/cm <sup>3</sup> |
| dipole moment (quantity)               | coulomb meter            | C·m                |                         |

<sup>†</sup>The spellings “metre” and “litre” are preferred by ASTM; however, “-er” is used in the *Encyclopedia*.

<sup>‡</sup>Wide use is made of Celsius temperature ( $t$ ) defined by

$$t = T - T_0$$

where  $T$  is the thermodynamic temperature, expressed in kelvin, and  $T_0 = 273.15$  K by definition. A temperature interval may be expressed in degrees Celsius as well as in kelvin.

| Quantity   | Unit                         | Symbol                | Acceptable equivalent |
|--|------------------------------|-----------------------|-----------------------|
| *dose equivalent   | sievert                      | Sv                    | J/kg                  |
| *electric capacitance  | farad                        | F                     | C/V                   |
| *electric charge, quantity of electricity                      | coulomb                      | C                     | A·s                   |
| electric charge density  | coulomb per cubic meter      | C/m <sup>3</sup>      |                       |
| *electric conductance  | siemens                      | S                     | A/V                   |
| electric field strength  | volt per meter               | V/m                   |                       |
| electric flux density  | coulomb per square meter     | C/m <sup>2</sup>      |                       |
| *electric potential, potential difference, electromotive force | volt                         | V                     | W/A                   |
| *electric resistance   | ohm                          | Ω                     | V/A                   |
| *energy, work, quantity of heat                                | megajoule                    | MJ                    |                       |
|  | kilojoule                    | kJ                    |                       |
|  | joule                        | J                     | N·m                   |
|  | electronvolt <sup>†</sup>    | eV <sup>†</sup>       |                       |
|  | kilowatt-hour <sup>†</sup>   | kW·h <sup>†</sup>     |                       |
| energy density   | joule per cubic meter        | J/m <sup>3</sup>      |                       |
| *force   | kilonewton                   | kN                    |                       |
|  | newton                       | N                     | kg·m/s <sup>2</sup>   |
| *frequency   | megahertz                    | MHz                   |                       |
|  | hertz                        | Hz                    | 1/s                   |
| heat capacity, entropy   | joule per kelvin             | J/K                   |                       |
| heat capacity (specific), specific entropy                     | joule per kilogram kelvin    | J/(kg·K)              |                       |
| heat-transfer coefficient                                      | watt per square meter kelvin | W/(m <sup>2</sup> ·K) |                       |
| *illuminance   | lux                          | lx                    | lm/m <sup>2</sup>     |
| *inductance  | henry                        | H                     | Wb/A                  |
| linear density   | kilogram per meter           | kg/m                  |                       |
| luminance  | candela per square meter     | cd/m <sup>2</sup>     |                       |
| *luminous flux   | lumen                        | lm                    | cd·sr                 |
| magnetic field strength  | ampere per meter             | A/m                   |                       |
| *magnetic flux   | weber                        | Wb                    | V·s                   |
| *magnetic flux density   | tesla                        | T                     | Wb/m <sup>2</sup>     |
| molar energy   | joule per mole               | J/mol                 |                       |
| molar entropy, molar heat capacity                             | joule per mole kelvin        | J/(mol·K)             |                       |
| moment of force, torque  | newton meter                 | N·m                   |                       |
| momentum   | kilogram meter per second    | kg·m/s                |                       |
| permeability   | henry per meter              | H/m                   |                       |
| permittivity   | farad per meter              | F/m                   |                       |
| *power, heat flow rate, radiant flux                           | kilowatt                     | kW                    |                       |
|  | watt                         | W                     | J/s                   |
| power density, heat flux density, irradiance                   | watt per square meter        | W/m <sup>2</sup>      |                       |
| *pressure, stress  | megapascal                   | MPa                   |                       |
|  | kilopascal                   | kPa                   |                       |
|  | pascal                       | Pa                    | N/m <sup>2</sup>      |
| sound level  | decibel                      | dB                    |                       |
| specific energy  | joule per kilogram           | J/kg                  |                       |
| specific volume  | cubic meter per kilogram     | m <sup>3</sup> /kg    |                       |
| surface tension  | newton per meter             | N/m                   |                       |
| thermal conductivity   | watt per meter kelvin        | W/(m·K)               |                       |
| velocity   | meter per second             | m/s                   |                       |
|  | kilometer per hour           | km/h                  |                       |
| viscosity, dynamic   | pascal second                | Pa·s                  |                       |
|  | millipascal second           | mPa·s                 |                       |
| viscosity, kinematic   | square meter per second      | m <sup>2</sup> /s     |                       |
|  | square millimeter per second | mm <sup>2</sup> /s    |                       |
| volume   | cubic meter                  | m <sup>3</sup>        |                       |
|  | cubic diameter               | dm <sup>3</sup>       | L (liter) (5)         |
|  | cubic centimeter             | cm <sup>3</sup>       | mL                    |
| wave number  | 1 per meter                  | m <sup>-1</sup>       |                       |
|  | 1 per centimeter             | cm <sup>-1</sup>      |                       |

<sup>†</sup>This non-SI unit is recognized by the CIPM as having to be retained because of practical importance or use in specialized fields (1).

In addition, there are 16 prefixes used to indicate order of magnitude, as follows:

| Multiplication factor | Prefix | Symbol          | Note   |
|-----------------------|--------|-----------------|--|
| $10^{18}$             | exa    | E               |  |
| $10^{15}$             | peta   | P               |  |
| $10^{12}$             | tera   | T               |  |
| $10^9$                | giga   | G               |  |
| $10^6$                | mega   | M               |  |
| $10^3$                | kilo   | k               |  |
| $10^2$                | hecto  | h <sup>a</sup>  | <sup>a</sup> Although hecto, deka, deci, and centi are SI prefixes, their use should be avoided except for SI unit-multiples for area and volume and nontechnical use of centimeter, as for body and clothing measurement. |
| 10                    | deka   | da <sup>a</sup> |  |
| $10^{-1}$             | deci   | d <sup>a</sup>  |  |
| $10^{-2}$             | centi  | c <sup>a</sup>  |  |
| $10^{-3}$             | milli  | m               |  |
| $10^{-6}$             | micro  | $\mu$           |  |
| $10^{-9}$             | nano   | n               |  |
| $10^{-12}$            | pico   | p               |  |
| $10^{-15}$            | femto  | f               |  |
| $10^{-18}$            | atto   | a               |  |

For a complete description of SI and its use the reader is referred to ASTM E380 (4).

A representative list of conversion factors from non-SI to SI units is presented herewith. Factors are given to four significant figures. Exact relationships are followed by a dagger. A more complete list is given in the latest editions of ASTM E380 (4) and ANSI Z210.1 (6).

### Conversion Factors to SI Units

| To convert from                 | To  | Multiply by                  |
|---------------------------------|---|------------------------------|
| acre                            | square meter (m <sup>2</sup> )                    | $4.047 \times 10^3$          |
| angstrom                        | meter (m)   | $1.0 \times 10^{-10\dagger}$ |
| are                             | square meter (m <sup>2</sup> )                    | $1.0 \times 10^{2\dagger}$   |
| astronomical unit               | meter (m)   | $1.496 \times 10^{11}$       |
| atmosphere, standard            | pascal (Pa)                                       | $1.013 \times 10^5$          |
| bar                             | pascal (Pa)                                       | $1.0 \times 10^{5\dagger}$   |
| barn                            | square meter (m <sup>2</sup> )                    | $1.0 \times 10^{-28\dagger}$ |
| barrel (42 U.S. liquid gallons) | cubic meter (m <sup>3</sup> )                     | 0.1590                       |
| Bohr magneton ( $\mu_B$ )       | J/T   | $9.274 \times 10^{-24}$      |
| Btu (International Table)       | joule (J)   | $1.055 \times 10^3$          |
| Btu (mean)                      | joule (J)   | $1.056 \times 10^3$          |
| Btu (thermochemical)            | joule (J)   | $1.054 \times 10^3$          |
| bushel                          | cubic meter (m <sup>3</sup> )                     | $3.524 \times 10^{-2}$       |
| calorie (International Table)   | joule (J)   | 4.187                        |
| calorie (mean)                  | joule (J)   | 4.190                        |
| calorie (thermochemical)        | joule (J)   | 4.184 <sup>†</sup>           |
| centipoise                      | pascal second (Pa·s)                              | $1.0 \times 10^{-3\dagger}$  |
| centistokes                     | square millimeter per second (mm <sup>2</sup> /s) | 1.0 <sup>†</sup>             |
| cfm (cubic foot per minute)     | cubic meter per second (m <sup>3</sup> /s)        | $4.72 \times 10^{-4}$        |
| cubic inch                      | cubic meter (m <sup>3</sup> )                     | $1.639 \times 10^{-5}$       |
| cubic foot                      | cubic meter (m <sup>3</sup> )                     | $2.832 \times 10^{-2}$       |
| cubic yard                      | cubic meter (m <sup>3</sup> )                     | 0.7646                       |
| curie                           | becquerel (Bq)                                    | $3.70 \times 10^{10\dagger}$ |
| debye                           | coulomb meter (C·m)                               | $3.336 \times 10^{-30}$      |
| degree (angle)                  | radian (rad)                                      | $1.745 \times 10^{-2}$       |
| denier (international)          | kilogram per meter (kg/m)                         | $1.111 \times 10^{-7}$       |
|                                 | tex <sup>‡</sup>                                  | 0.1111                       |
| dram (apothecaries')            | kilogram (kg)                                     | $3.888 \times 10^{-3}$       |
| dram (avoirdupois)              | kilogram (kg)                                     | $1.772 \times 10^{-3}$       |
| dram (U.S. fluid)               | cubic meter (m <sup>3</sup> )                     | $3.697 \times 10^{-6}$       |
| dyne                            | newton (N)  | $1.0 \times 10^{-5\dagger}$  |
| dyne/cm                         | newton per meter (N/m)                            | $1.0 \times 10^{-3\dagger}$  |
| electronvolt                    | joule (J)   | $1.602 \times 10^{-19}$      |
| erg                             | joule (J)   | $1.0 \times 10^{-7\dagger}$  |

<sup>†</sup>Exact.

<sup>‡</sup>See footnote on p. xvi.

| To convert from                   | To  | Multiply by            |
|-----------------------------------|---|------------------------|
| fathom                            | meter (m)                                     | 1.829                  |
| fluid ounce (U.S.)                | cubic meter (m <sup>3</sup> )                 | $2.957 \times 10^{-5}$ |
| foot                              | meter (m)                                     | 0.3048 <sup>†</sup>    |
| footcandle                        | lux (lx)                                      | 10.76                  |
| furlong                           | meter (m)                                     | $2.012 \times 10^{-2}$ |
| gal                               | meter per second squared (m/s <sup>2</sup> )  | $1.0 \times 10^{-2†}$  |
| gallon (U.S. dry)                 | cubic meter (m <sup>3</sup> )                 | $4.405 \times 10^{-3}$ |
| gallon (U.S. liquid)              | cubic meter (m <sup>3</sup> )                 | $3.785 \times 10^{-3}$ |
| gallon per minute (gpm)           | cubic meter per second (m <sup>3</sup> /s)    | $6.309 \times 10^{-5}$ |
|                                   | cubic meter per hour (m <sup>3</sup> /h)      | 0.2271                 |
| gauss                             | tesla (T)                                     | $1.0 \times 10^{-4}$   |
| gilbert                           | ampere (A)                                    | 0.7958                 |
| gill (U.S.)                       | cubic meter (m <sup>3</sup> )                 | $1.183 \times 10^{-4}$ |
| grade                             | radian  | $1.571 \times 10^{-2}$ |
| grain                             | kilogram (kg)                                 | $6.480 \times 10^{-5}$ |
| gram force per denier             | newton per tex (N/tex)                        | $8.826 \times 10^{-2}$ |
| hectare                           | square meter (m <sup>2</sup> )                | $1.0 \times 10^{4†}$   |
| horsepower (550 ft·lbf/s)         | watt (W)                                      | $7.457 \times 10^2$    |
| horsepower (boiler)               | watt (W)                                      | $9.810 \times 10^3$    |
| horsepower (electric)             | watt (W)                                      | $7.46 \times 10^{2†}$  |
| hundredweight (long)              | kilogram (kg)                                 | 50.80                  |
| hundredweight (short)             | kilogram (kg)                                 | 45.36                  |
| inch                              | meter (m)                                     | $2.54 \times 10^{-2†}$ |
| inch of mercury (32°F)            | pascal (Pa)                                   | $3.386 \times 10^3$    |
| inch of water (39.2°F)            | pascal (Pa)                                   | $2.491 \times 10^2$    |
| kilogram-force                    | newton (N)                                    | 9.807                  |
| kilowatt hour                     | megajoule (MJ)                                | 3.6 <sup>†</sup>       |
| kip                               | newton (N)                                    | $4.448 \times 10^3$    |
| knot (international)              | meter per second (m/S)                        | 0.5144                 |
| lambert                           | candela per square meter (cd/m <sup>2</sup> ) | $3.183 \times 10^3$    |
| league (British nautical)         | meter (m)                                     | $5.559 \times 10^3$    |
| league (statute)                  | meter (m)                                     | $4.828 \times 10^3$    |
| light year                        | meter (m)                                     | $9.461 \times 10^{15}$ |
| liter (for fluids only)           | cubic meter (m <sup>3</sup> )                 | $1.0 \times 10^{-3†}$  |
| maxwell                           | weber (Wb)                                    | $1.0 \times 10^{-8†}$  |
| micron                            | meter (m)                                     | $1.0 \times 10^{-6†}$  |
| mil                               | meter (m)                                     | $2.54 \times 10^{-5†}$ |
| mile (statute)                    | meter (m)                                     | $1.609 \times 10^3$    |
| mile (U.S. nautical)              | meter (m)                                     | $1.852 \times 10^3†$   |
| mile per hour                     | meter per second (m/s)                        | 0.4470                 |
| millibar                          | pascal (Pa)                                   | $1.0 \times 10^2$      |
| millimeter of mercury (0°C)       | pascal (Pa)                                   | $1.333 \times 10^{2†}$ |
| minute (angular)                  | radian  | $2.909 \times 10^{-4}$ |
| myriagram                         | kilogram (kg)                                 | 10                     |
| myriameter                        | kilometer (km)                                | 10                     |
| oersted                           | ampere per meter (A/m)                        | 79.58                  |
| ounce (avoirdupois)               | kilogram (kg)                                 | $2.835 \times 10^{-2}$ |
| ounce (troy)                      | kilogram (kg)                                 | $3.110 \times 10^{-2}$ |
| ounce (U.S. fluid)                | cubic meter (m <sup>3</sup> )                 | $2.957 \times 10^{-5}$ |
| ounce-force                       | newton (N)                                    | 0.2780                 |
| peck (U.S.)                       | cubic meter (m <sup>3</sup> )                 | $8.810 \times 10^{-3}$ |
| pennyweight                       | kilogram (kg)                                 | $1.555 \times 10^{-3}$ |
| pint (U.S. dry)                   | cubic meter (m <sup>3</sup> )                 | $5.506 \times 10^{-4}$ |
| pint (U.S. liquid)                | cubic meter (m <sup>3</sup> )                 | $4.732 \times 10^{-4}$ |
| poise (absolute viscosity)        | pascal second (Pa·s)                          | 0.10 <sup>†</sup>      |
| pound (avoirdupois)               | kilogram (kg)                                 | 0.4536                 |
| pound (troy)                      | kilogram (kg)                                 | 0.3732                 |
| poundal                           | newton (N)                                    | 0.1383                 |
| pound-force                       | newton (N)                                    | 4.448                  |
| pound force per square inch (psi) | pascal (Pa)                                   | $6.895 \times 10^3$    |
| quart (U.S. dry)                  | cubic meter (m <sup>3</sup> )                 | $1.101 \times 10^{-3}$ |
| quart (U.S. liquid)               | cubic meter (m <sup>3</sup> )                 | $9.464 \times 10^{-4}$ |



| To convert from              | To  | Multiply by                   |
|------------------------------|---|-------------------------------|
| quintal                      | kilogram (kg)                               | $1.0 \times 10^{2\dagger}$    |
| rad                          | gray (Gy)                                   | $1.0 \times 10^{-2\dagger}$   |
| rod                          | meter (m)                                   | 5.029                         |
| roentgen                     | coulomb per kilogram (C/kg)                 | $2.58 \times 10^{-4}$         |
| second (angle)               | radian (rad)                                | $4.848 \times 10^{-6\dagger}$ |
| section                      | square meter (m <sup>2</sup> )              | $2.590 \times 10^6$           |
| slug                         | kilogram (kg)                               | 14.59                         |
| spherical candle power       | lumen (lm)                                  | 12.57                         |
| square inch                  | square meter (m <sup>2</sup> )              | $6.452 \times 10^{-4}$        |
| square foot                  | square meter (m <sup>2</sup> )              | $9.290 \times 10^{-2}$        |
| square mile                  | square meter (m <sup>2</sup> )              | $2.590 \times 10^6$           |
| square yard                  | square meter (m <sup>2</sup> )              | 0.8361                        |
| stere                        | cubic meter (m <sup>3</sup> )               | 1.0 <sup>†</sup>              |
| stokes (kinematic viscosity) | square meter per second (m <sup>2</sup> /s) | $1.0 \times 10^{-4\dagger}$   |
| tex                          | kilogram per meter (kg/m)                   | $1.0 \times 10^{-6\dagger}$   |
| ton (long, 2240 pounds)      | kilogram (kg)                               | $1.016 \times 10^3$           |
| ton (metric) (tonne)         | kilogram (kg)                               | $1.0 \times 10^3\dagger$      |
| ton (short, 2000 pounds)     | kilogram (kg)                               | $9.072 \times 10^2$           |
| torr                         | pascal (Pa)                                 | $1.333 \times 10^2$           |
| unit pole                    | weber (Wb)                                  | $1.257 \times 10^{-7}$        |
| yard                         | meter (m)                                   | 0.9144 <sup>†</sup>           |

<sup>†</sup>Exact.

## BIBLIOGRAPHY

1. The International Bureau of Weights and Measures, BIPM (Parc de Saint-Cloud, France) is described in Appendix X2 of Ref. 4. This bureau operates under the exclusive supervision of the International Committee for Weights and Measures (CIPM).
2. *Metric Editorial Guide (ANMC-78-1)*, latest ed., American National Metric Council, 5410 Grosvenor Lane, Bethesda, Md. 20814, 1981.
3. *SI Units and Recommendations for the Use of Their Multiples and of Certain Other Units (ISO 1000-1981)*, American National Standards Institute, 1430 Broadway, New York, 10018, 1981.
4. Based on *ASTM E380-89a (Standard Practice for Use of the International System of Units (SI))*, American Society for Testing and Materials, 1916 Race Street, Philadelphia, Pa. 19103, 1989.
5. *Fed. Reg.*, Dec. 10, 1976 (41 FR 36414).
6. For ANSI address, see Ref. 3.

R. P. LUKENS  
ASTM Committee E-43 on SI Practice

## A-DNA

A-DNA was first recognized as a [DNA structure](#) using fiber X-ray diffraction analysis. [B-DNA](#) can be converted to A-DNA under conditions of low hydration, and the process is reversible. The A-DNA double helix is short and fat, with the base pairs and backbone wrapped farther away from the helix axis (see Fig. 2 of [DNA Structure](#)). The [base pairs](#) are significantly tilted ( $\sim 19^\circ$ ) with respect to the helix axis. The major groove has a very narrow width of  $\sim 3 \text{ \AA}$  and a depth of  $\sim 13 \text{ \AA}$ , whereas the minor groove has a broad width of  $11 \text{ \AA}$  and a shallow depth of  $3 \text{ \AA}$ . The base pairs also display a minor propeller twist.

There is increasing evidence that A-DNA may play an important role in biological processes such as protein recognition and [transcription](#) regulation, as in the [TATA-box](#) sequence bound with the TATA-box binding proteins ([1](#), [2](#)). There is a sequence-dependent propensity to form A-DNA. Guanine-rich regions readily form A-DNA, whereas stretches of adenine resist it. The crystal structures of DNA oligonucleotides having guanine-rich sequences showed a characteristic intrastrand guanine–guanine stacking interaction in the A-DNA double helix, which may explain the propensity of these sequences to adopt the A-DNA conformation. Moreover, the packing of the helices in the crystal lattice revealed a characteristic pattern of the terminal base pairs from one helix abutting the minor groove surface of the neighboring helix, thus minimizing the accessibility of solvent to the wide minor groove. Since a low humidity environment favors formation of A-type helix, the displacement of surface solvent molecules by **hydrophobic** base pairs provides a driving force in stabilizing short oligonucleotides in the A-DNA conformation. Recently, it has been demonstrated that complex ions such as spermine, cobalt(III)hexamine, and [neomycin](#) can facilitate the B-DNA to A-DNA transition for DNA containing stretches of  $(dG)_n \cdot (dC)_n$  sequences.

Finally, the crystal structures of a number of DNA:RNA hybrids, such as the self-complementary r(GCG)d(TATACGC), showed that DNA–RNA hybrid helices are of the A-DNA type. All the ribose and 2'-deoxyribose sugars are in the C3'-*endo* conformation. The 2'-hydroxyl groups of the ribose are involved in different types of hydrogen bonding to the adjacent nucleotides in the chain.

### Bibliography

1. J. L. Kim, D. B. Nikolov, and S. K. Burley (1993) *Nature* **365**, 520–527.
2. Y. Kim, J. H. Geiger, S. Hahn, and P. B. Sigler (1993) *Nature* **365**, 512–520.

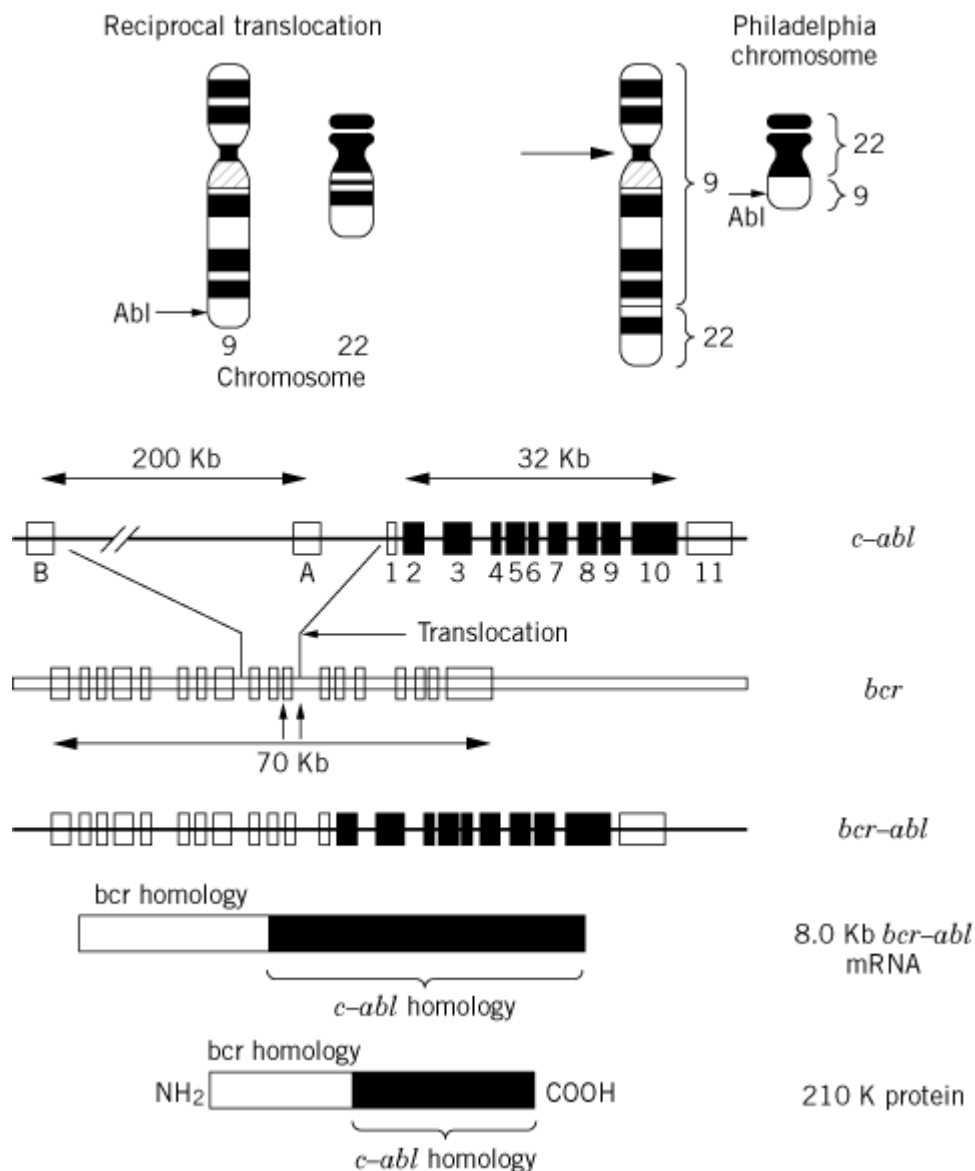
## abl Oncogenes

The *abl* gene was first identified as the transforming element of Abelson murine leukemia virus (A-MuLV), a replication-defective [retrovirus](#) that was isolated after inoculating Moloney murine leukemia virus (M-MuLV) into prednisolone-treated BALB/c mice ([1](#)). A-MuLV induces **B-cell** lymphomas *in vivo* and transforms both lymphoid and fibroblastic cells *in vitro*. The proviral genome of A-MuLV encodes a single polypeptide chain that is a fusion product of the virally-derived gag and cell-derived *abl* sequences ([2](#)). Sequence analysis of c-*abl* sequences revealed that, like the *src* gene, the *abl* gene codes for a **tyrosine kinase** that also contains the unique, SH3, SH2, and tyrosine kinase domains ([3](#)). Unlike the *src* gene, however, the *abl* gene product contains an additional C-terminal domain whose function is not entirely clear. In addition, the gene product of c-

*abl* does not contain the negative-regulatory **tyrosine** residue at its C-terminus. The tyrosine kinase activity of the c-Abl protein is negatively regulated by its SH3 domain, which is deleted from the *v-abl* gene product. Interestingly, it was also found that the *v-abl* gene product contains a point mutation in its C-terminal sequence, which enhances its tyrosine kinase and transforming activities (4). Thus, both the *v-src* and *v-abl* genes have alterations in their regulatory sequences that result in the constitutive activation of their tyrosine kinase activities, which correlates with their transforming function.

Gene mapping studies have established that the *c-abl* oncogene is located on human chromosome 9q34, the location where the break point occurs in the Philadelphia chromosome. The Philadelphia chromosome is generated when a portion of the *c-abl* gene is **translocated** to chromosome 22 and is fused to a portion of the gene called *bcr*, which itself is disrupted during the translocation process (Fig. 1). This process results in generating a new gene, called *bcr-abl*, that has enhanced oncogenic activity and whose expression leads to the development of leukemia (5, 6). It is interesting to note that the BCR-ABL gene is structurally similar to the *gag-abl* gene encoded by the Abelson murine leukemia virus. Both the *gag-abl* and *bcr-abl* genes exhibit high levels of tyrosine kinase activity, which is essential for their transforming activity.

**Figure 1.** Generation of the *bcr-abl* oncoprotein by chromosomal translocation. The *abl* gene in a normal cell, is located on chromosome 9 and encodes a tyrosine kinase. During malignant transformation of myeloid cells, a portion of chromosome 9 that contains the *abl* locus translocates to chromosome 22 at the breakpoint cluster region (*bcr*) locus and generates the chimeric *bcr-abl* oncoprotein. Because the translocation results in deleting the sequences that negatively regulate *abl* tyrosine kinase activity, the fusion protein has constitutive and increased levels of enzymatic activity.



The function of *c-abl* in normal cell growth is not fully established. The *c-abl* gene codes for two 145-kDa proteins as a result of [alternative splicing](#) of the two first exons (see [RNA Splicing](#) and [Introns, Exons](#)). This results in synthesizing two proteins that differ in their amino-terminal sequences. Both forms of *c-Abl* are found in the cytoplasm and in the nucleus. The *c-Abl* protein can bind to DNA and to [cell-cycle](#) regulators like the **retinoblastoma** protein. Recent studies show that *c-abl* gene expression is induced during cellular stress caused by agents, such as ionizing radiation and certain other genotoxic agents. These agents induce the formation of a complex involving *c-Abl*, DNA-dependent protein kinase, and Ku antigen (7). The DNA-dependent protein kinase in this complex is activated by **DNA damage** and in turn phosphorylates and activates *c-Abl*. Recent studies have also shown that *c-Abl* associates with the product of the ATM gene, which in turn activates *c-Abl* in response to ionizing radiation. DNA damage also induces binding of *c-Abl* to [p53](#) and contributes to cell-cycle arrest at the G1-phase, which is mediated by p53. Recent studies also show that the *c-Abl* protein binds to **protein kinase C- $\delta$**  and phosphorylates the latter, resulting in its activation and translocation to the nucleus, where it participates in inducing [apoptosis](#) (8). Thus, a substantial amount of evidence gathered in the past few years indicates that *c-Abl* protein has a pivotal role in mediating cellular growth arrest and the apoptotic effects that occur during exposure to ionizing radiation and genotoxic stress.

## Bibliography

1. H. T. Abelson and L. S. Rabson (1970) *Cancer Res.* **30**, 2213–2222.
2. E. P. Reddy, M. J. Smith, and A. Srinivasan (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3623–3627.
3. C. Oppi, S. K. Shore, and E. P. Reddy (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8200–8204.
4. S. K. Shore, S. L. Bogart, and E. P. Reddy (1990) *Proc. Natl. Acad. Sci. USA* **87**, 6502–6506.
5. J. Groffen, J. R. Stephenson, N. Heisterkamp, A. De Klein, C. B. Bartam, and G. Grosveld (1984) *Cell* **36**, 93–99.
6. E. Shtivelman, R. P. Lifshitz, R. P. Gale, and E. Canaani (1985) *Nature* **315**, 550–553.
7. R. Baskaran, L. D. Wood, L. L. Whittaker, Y. Xu, C. Barlow, C. E. Canman, S. E. Morgan, D. Baltimore, A. Wynshaw-Boris, M. B. Kastan, and J. Y. J. Wang (1997) *Nature* **387**, 516–519.
8. Z-M. Yuan, T. Utsugisawa, T. Ishiko, S. Nakada, Y. Huang, S. Kharbanda, R. Weichselbaum, and D. Kufe (1998) *Oncogene* **16**, 1643–1648.

## Abscisic Acid

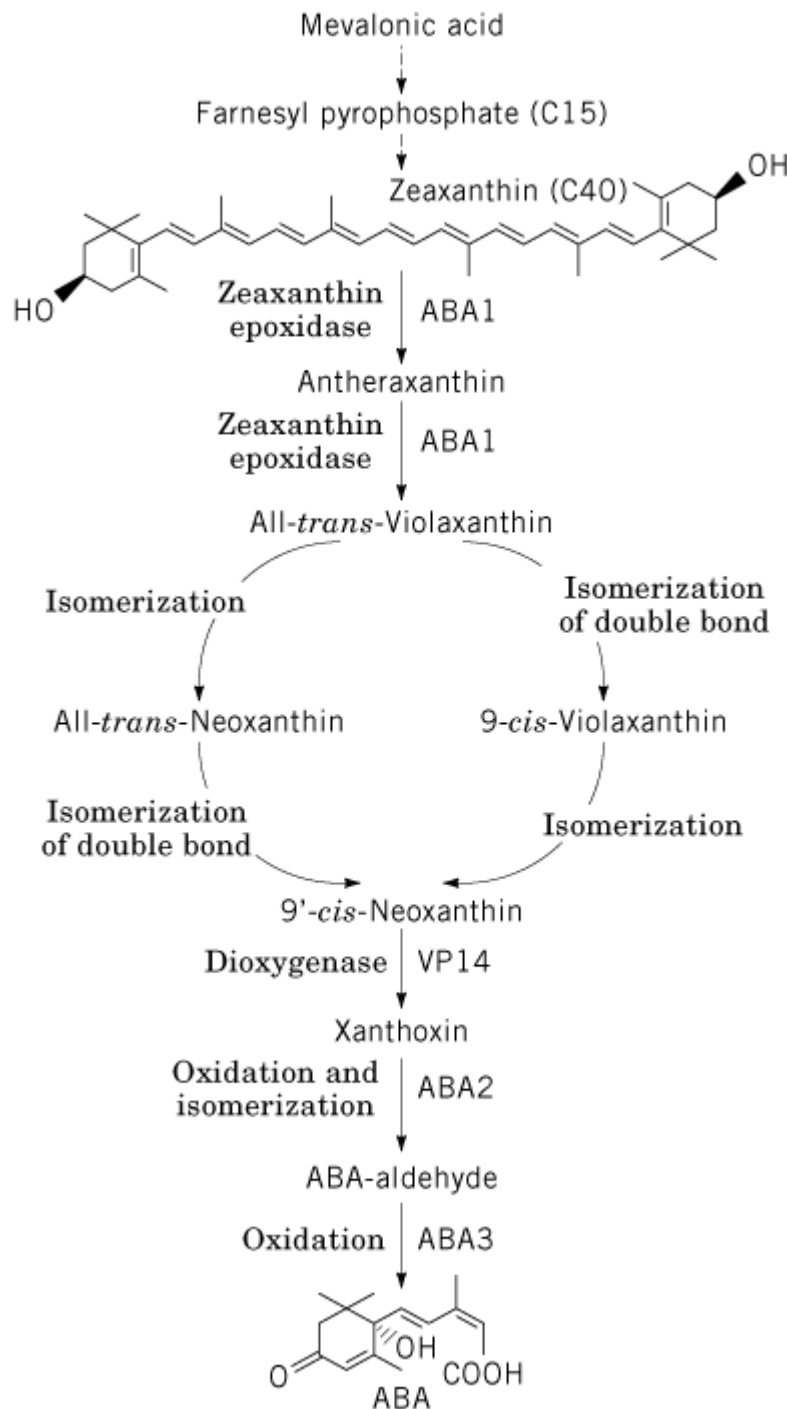
### 1. History

The first evidence for the existence of an acidic inhibitor of coleoptile growth that promoted abscission and seed maturation dates back to the early 1950s, but it was not until 1963 that abscisic acid (ABA) was identified by Frederick Addicott and coworkers (1, 2). Addicott's team was studying compounds that stimulated abscission in cotton fruits and had named the active substances *abscisin I* and *abscisin II*; the latter proved to be ABA. Two other independent efforts also culminated in the discovery of ABA. A British group headed by Wareing (3) was investigating bud dormancy of woody plants and called the most active molecule *dormin*. In New Zealand, van Stevenick (4) studied compounds that accelerated abscission of flowers and fruits of *Lupinus luteus*. In 1964, it became evident that the three groups had discovered the same [plant hormone](#), which was renamed *abscisic acid* 3 years later.

### 2. Biosynthesis and Metabolism

ABA is a universal compound in vascular plants; it is not found in bacteria, but has been reported in green algae, certain fungi, and mosses (5). ABA is a sesquiterpenoid, with mevalonic acid as a precursor. In certain fungi, ABA is produced by a direct, C15 pathway from farnesyl pyrophosphate. In higher plants, ABA is derived from xanthophylls, via an indirect C40 pathway (Fig. 1). Substantial progress in understanding the ABA biosynthetic pathway has been achieved by a combination of molecular and genetic techniques, primarily on mutants in *Arabidopsis thaliana*, *Zea mays*, *Nicotiana glauca*, and tomato. [Arabidopsis](#) mutants impaired in ABA biosynthesis were isolated on the basis of their lack of seed dormancy due to ABA deficiency and their ability to overcome a [gibberellin](#) requirement for germination, which allowed them to germinate in the presence of inhibitors of gibberellin biosynthesis or in a gibberellin-deficient background (6, 7). ABA-deficient mutants show a wilted **phenotype** when subjected to water stress.

**Figure 1.** ABA biosynthesis in higher plants. Steps mediated by ABA1, ABA2, and ABA3 in *Arabidopsis* and by VP14 in *Zea mays* are indicated. (Adapted from Ref. 19.)



Mevalonic acid is converted to farnesyl pyrophosphate, a C15 compound, via several intermediates (Fig. 1). The subsequent conversion of farnesyl pyrophosphate to zeaxanthin, a C40 carotenoid, again involves multiple steps (8). A number of *viviparous* (*vp*) mutants of maize identify loci corresponding to these early conversions. The transformation of zeaxanthin to all-*trans*-violaxanthin consists of two epoxidations at the double bonds in both cyclohexenyl rings, with antheraxanthin as an intermediate. The *aba1* mutant impaired in zeaxanthin epoxidase activity has been characterized biochemically in *Arabidopsis* (6, 9, 10). The *N. plumbaginifolia* *ABA2* gene encoding zeaxanthin epoxidase was recently cloned by [transposon](#) tagging (11). The gene encodes a **chloroplast**-imported [polypeptide chain](#) with similarity to bacterial monooxygenases and oxidases. When the *Nicotiana ABA2* gene was expressed in *Escherichia coli* heterologously, the protein was shown to catalyze both epoxidations *in vitro*. Moreover, the [complementary DNA \(cDNA\)](#) complemented both the *N.*

*plumbaginifolia aba2* mutant and the *Arabidopsis aba1* mutant. The ortholog in *Arabidopsis* (*ABA1*) was cloned by [homology](#) with *Nicotiana ABA2* (11). Although presumably a null **allele**, the *Nicotiana aba2* mutant retains up to half the ABA content of wild type in the absence of water stress. This characteristic might indicate the existence of a secondary biosynthetic pathway for ABA. Downstream of all-*trans*-violaxanthin, two isomerizations yield 9'-*cis*-neoxanthin.

The subsequent oxidative cleavage with xanthoxin as a product is thought to be the rate-limiting step in ABA biosynthesis (12, 13). A gene that most probably encodes the [enzyme](#) performing this reaction was cloned from maize, and by homology also from *Arabidopsis* (*VP14*) (13). Its overproduction in *E. coli* indicated that VP14 is probably a member of a novel class of dioxygenases.

Xanthoxin is subsequently converted to ABA-aldehyde by oxidation and isomerization. Biochemical analysis indicates that the *Arabidopsis aba2* mutant is blocked at this step in the pathway (7, 14). A final oxidation, catalyzed by ABA-aldehyde oxidase, produces ABA. Mutants in this step have been identified in several species (8). The *Arabidopsis aba3* mutant lacks several aldehyde oxidase activities that require a molybdenum cofactor (14). Yet the mutation blocks modification of the molybdenum cofactor, rather than impairing the activity of the apoprotein. Cloning of the *ABA2* and *ABA3* genes should provide further clues as to the regulation of ABA biosynthesis in these final steps.

Metabolism of ABA is mainly to phaseic acid, dihydrophaseic acid, and their respective conjugates (15, 16). Direct conjugation of ABA to an ABA glucose ester, or an ABA- $\beta$ -glucopyranoside, can also occur. Conjugation of ABA does not lead to its storage, probably because it is irreversible and yields unstable compounds. With the exception of phaseic acid, none of the metabolites carries biological activity (16).

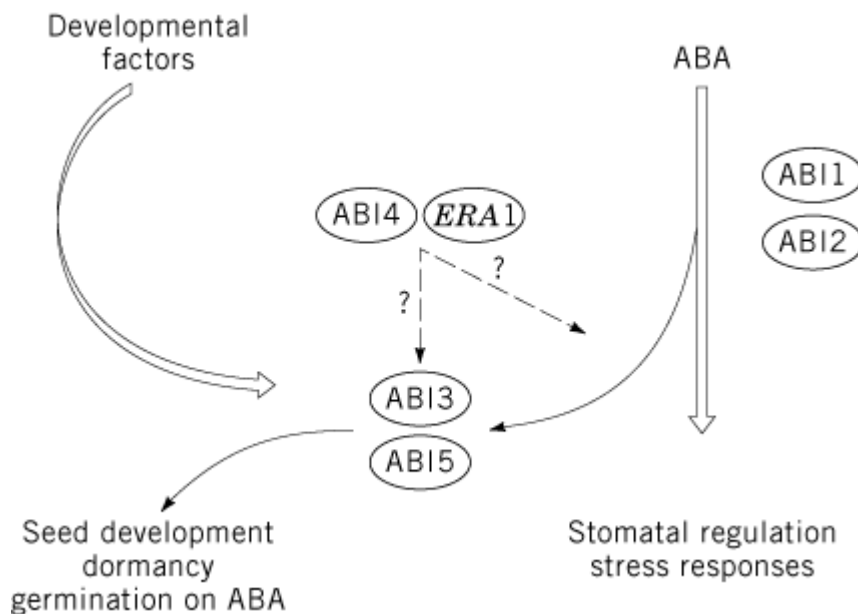
### 3. Signal Perception and Transduction

Insight into the molecular basis of ABA responses was gathered by a combination of molecular-genetic, biochemical, and electrophysiological studies (12, 17-20). Currently, little is known about the perception of ABA, despite numerous efforts to isolate ABA receptors using different approaches (20). Two groups of mutants, ABA-insensitive and ABA-hypersensitive, have largely contributed to the understanding of ABA [signal transduction](#). Whereas the five *abscisic acid-insensitive* (*abi*) mutants of *Arabidopsis* were identified by their ability to germinate on inhibiting concentrations of ABA (21, 22), the three *enhanced response to ABA* (*era*) mutants showed an increased dormancy in the presence of low ABA concentrations, compared to wild type (23). Besides their effects on seed development and dormancy, these mutants show altered responses to drought adaptation. The *abi* mutants differ from the above-mentioned *aba* mutants, because they do not have reduced endogenous ABA content and their phenotype cannot be reverted by exogenously supplied ABA. The characterization of *abi* mutants revealed that at least two ABA-response pathways exist: one primarily active in vegetative tissues, involving ABI1 and ABI2, and a second one operating predominantly during seed development and involving ABI3, ABI4, and ABI5 (19). Genes corresponding to four loci have been **cloned**. *ABI1* and *ABI2* encode highly homologous [serine/threonine kinases and phosphatases](#) 2C (24-26). These phosphatases might act in a phosphorylation/dephosphorylation cascade, as the protein phosphatase activity of recombinant ABI1 has been demonstrated *in vitro* (27). The functions of ABI1 and ABI2 are partially redundant (25). *ABI3* encodes a putative [transcription factor](#) that acts mainly during seed development and is the ortholog of *VP1* in maize (28-30). The suggested function of ABI3 consists of either activating the maturation program or preventing germination (30, 31). Finally, the predicted amino acid sequence of ERA1 shares similarity with the  $\beta$ -subunit of farnesyl transferases. In yeast and mammalian systems, this enzyme is known to modify signal transduction proteins for [membrane anchoring](#) (23). In wild-type *Arabidopsis*, ERA1 is proposed to modify a negative regulator of ABA signaling.

The current model of ABA signal transduction integrates ABI1/ABI2-dependent cascades in both

seeds and vegetative tissues (Fig. 2). These cascades interact with the ABI3 protein in seeds only, because ABI3 is expressed exclusively in seeds (32). The role of ABI3 is not confined to ABA signaling (19); it is also thought to mediate developmental signals. This role is supported by the seed phenotype of *abi3*, which is more complex than that of *abi1*, *abi2*, or *aba1* mutants (19). The evidence that VP1 can also regulate [transcription](#) independently of ABA further supports the possibility of a less strict limitation of ABI3 to ABA signaling (33, 34). The same holds true for ABI4 and ABI5 (2).

**Figure 2.** A model for ABA signal transduction.



Intermediates in the ABA transduction chain were revealed by microinjection of putative [second messengers](#) into hypocotyl cells of the *aurea* mutant of tomato (35). Following injection of ABA-inducible *KIN2* and *RD29A* **promoters** linked to b-glucuronidase, coinjection of either  $\text{Ca}^{2+}$ , cyclic ADP-ribose (cADP), ADP-ribosyl cyclase, or inositol (1,4,5)-triphosphate ( $\text{IP}_3$ ) (see [Calcium Signaling](#)) resulted in expression of the ABA-inducible genes in the absence of ABA, implying that cADP ribose, a  $\text{Ca}^{2+}$ -mobilizing second messenger, is involved in ABA responses. At present, it is difficult to integrate these data into the current model for the ABA signal transduction pathway (19). Electrophysiological studies (voltage clamp or patch clamp) on the ABA-insensitive mutants *abi1* and *abi2* have clearly shown an altered **ion channel** behavior in their stomatal guard cells and provided further evidence that phosphorylation/dephosphorylation cascades are involved in ABA signaling (36, 37).

Given the power of mutational approaches to identify signal transduction components, several groups are currently aiming at the identification of additional ABA signal transduction factors, exploiting novel screening procedures in *Arabidopsis* (12). At least eight *growth control via ABA* (*gca*) loci have been identified on the basis of their insensitivity to ABA inhibition of root growth (38). Furthermore, a *freezing sensitive* (*frs1*) mutant with a wilted phenotype appeared to be nonallelic to other ABA-deficient or -insensitive mutants of *Arabidopsis* (39). Other screens employ **transgenic** lines carrying chimeric promoter-**luciferase** constructs with promoters derived from ABA-responsive genes. Putative mutants are identified by aberrant expression of these genes in the mutagenized progeny (40). The identification of intragenic or extragenic [enhancers](#) and **suppressors** of known signaling components is yet another alternative.



Whereas part of the ABA signaling pathway is well established in seeds and stomatal guard cells, the implications of ABA signaling in the function of the vegetative meristem remain to be clarified. Furthermore, it will be most interesting to characterize the regulatory influence of light on ABA signaling, as well as the crosstalk with gibberellins, which are known to counteract ABA in many physiological processes (5).

#### 4. Downstream Targets

Of all the plant hormones, ABA is probably best known from the point of view of responsiveness at the target gene level. ABA-responsive genes encode proteins with diverse functions and are often induced under water stress (drought, salinity, low temperature). It is important to mention, however, that certain water stress-related genes are independent of ABA (17, 18). ABA-responsive genes can be classified into two groups. Primary ABA response genes are rapidly induced and are independent of **protein biosynthesis**, implying that *trans*-acting factors controlling these genes are under post-translational control. A second class of ABA response genes consists of those that require the expression of other genes.

[Cis-acting](#) and [trans-acting](#) elements involved in ABA induction of **gene expression** have been studied extensively, often in relation to water stress (17, 18). A bipartite model for regulation of transcription of the primary response genes is proposed, as for the barley *HVA1* and *HVA22* genes (41, 42). The specificity of ABA responsiveness relies on the combination of an ABA [response element](#) (ABRE) and a coupling element. The ABRE has a **consensus sequence** (PyACGTGGC) that contains a G-box core sequence and has been shown to interact with bZIP factors (eg, wheat *Em* gene-binding protein EmBP-1) (43). A set of unique coupling elements would provide specificity in ABA response at the level of individual genes. In addition, in the case of *Em* induction, it has been demonstrated that a conserved domain of 18 amino acid residues of the VP1 protein enhances DNA-binding activity of the EmBP-1 bZIP factor, thereby acting as a DNA chaperone (44). A yeast [two-hybrid system](#) screen was used to identify proteins that interact specifically with VP1 on the promoter of the *Em* gene. One such protein was related to 14–3–3 proteins (45) and thus may help to stabilize and/or activate the regulatory complex at the *Em* promoter (20).

ABREs are not the only *cis* elements that confer ABA responsiveness. Both in the case of seed desiccation and in water stress conditions, other promoter elements can be involved. For example, ABA responses can be mediated in conjunction with the Sph box in germination-specific promoters. A direct interaction between the Sph box and a 140-residue conserved domain in the C-terminal region of VP1 has been demonstrated (46).

The second pathway for ABA-dependent gene expression requires protein biosynthesis, exemplified by the *Arabidopsis* ABA-inducible *RD22* gene. A 67-bp region in the promoter of this gene contains conserved motifs for MYC and MYB factor binding (47) (see [Oncogenes, Oncoproteins](#)). Factors binding to these promoter elements (RD22BP-1 and ATMYB2) were shown to activate transcription of *RD22* in transient transactivation assays (48). Cooperative binding of these factors has been proposed to control ABA-dependent expression of *RD22*.

One of the challenges for the future consists in making the link between the ABA-induced proteins and the tissue- and development-specific physiological responses, summarized below.

#### 5. Effects

Although abscission was originally considered to be a process regulated by ABA, and the hormone was named accordingly, it is now clear that [ethylene](#) is the primary factor controlling abscission. Nevertheless, ABA also plays a key role in such diverse plant growth and developmental processes as seed maturation, germination, root growth, and stomatal closure (5). Stomatal movements involve changes in ion fluxes that occur very rapidly in response to ABA and do not require gene expression.

Furthermore, ABA is a major stress signal orchestrating responses to dehydration stress, including water stress induced by drought and high-salt concentration as well as by low-temperature conditions.

## Bibliography

1. T. A. Bennet-Clark and N. P. Kefford (1953) *Nature* **171**, 645–647.
2. K. Ohkuma, O. E. Smith, J. L. Lyon, and F. T. Addicott (1963) *Science* **142**, 1592–1593.
3. P. F. Wareing (1956) *Annu. Rev. Plant Physiol.* **7**, 191–214.
4. R. F. M. van Stevenick (1957) *J. Exp. Bot.* **8**, 373–381.
5. P. J. Davies (1995) *Plant Hormones: Physiology, Biochemistry and Molecular Biology*, Kluwer, Dordrecht, The Netherlands.
6. M. Koornneef, M. L. Jorna, D. L. C. Brinkhorst-van der Swan, and C. M. Karssen (1982) *Theor. Appl. Genet.* **61**, 385–393.
7. K. M. Léon-Kloosterziel, M. Alvarez Gil, G. J. Ruijs, S. E. Jacobsen, N. E. Olszewski, S. H. Schwartz, J. A. D. Zeevaart, and M. Koornneef (1996) *Plant J.* **10**, 655–661.
8. I. B. Taylor (1991) In *Abscisic Acid, Physiology and Biochemistry* (W. J. Davies and H. G. Jones, eds.), Bios Scientific, Oxford, UK, pp. 23–37.
9. S. C. Duckham, R. S. T. Linforth, and I. B. Taylor (1991) *Plant Cell Environ.* **14**, 631–636.
10. C. D. Rock and J. A. D. Zeevaart (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7496–7499.
11. E. Marin, L. Nussaume, A. Quesada, M. Gonneau, B. Sotta, P. Huguency, A. Frey, and A. Marion-Poll (1996) *EMBO J.* **15**, 2331–2342.
12. M. Koornneef, K. M. Léon-Kloosterziel, S. H. Schwartz, and J. A. D. Zeevaart (1998) *Plant Physiol. Biochem.* **36**, 83–89.
13. S. H. Schwartz, K. M. Léon-Kloosterziel, M. Koornneef, and J. A. D. Zeevaart (1997) *Plant Physiol.* **114**, 161–166.
14. S. H. Schwartz, B. C. Cai, D. A. Gage, J. A. D. Zeevaart, and D. R. McCarty (1997) *Science* **276**, 1872–1874.
15. J. A. D. Zeevaart and R. A. Creelman (1988) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **39**, 439–473.
16. G. Sembdner, R. Atzorn, and G. Schneider (1994) *Plant Mol. Biol.* **26**, 1459–1481.
17. J. Ingram and D. Bartels (1996) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **47**, 377–403.
18. K. Shinozaki and K. Yamaguchi-Shinozaki (1996) *Curr. Opin. Biotechnol.* **7**, 161–167.
19. S. Merlot and J. Giraudat (1997) *Plant Physiol.* **114**, 751–757.
20. R. S. Quatrano, D. Bartels, T. H. D. Ho, and M. Pagés (1997) *Plant Cell* **9**, 470–475.
21. M. Koornneef, G. Reuling, and C. M. Karssen (1984) *Physiol. Plant.* **61**, 377–383.
22. R. R. Finkelstein (1994) *Plant J.* **5**, 765–771.
23. S. Cutler, M. Ghassemian, D. Bonetta, S. Cooney, and P. McCourt (1996) *Science* **273**, 1239–1241.
24. J. Leung, M. Bouvier-Durand, P.-C. Morris, D. Guerrier, F. Cheddor, and J. Giraudat (1994) *Science* **264**, 1448–1452.
25. J. Leung, S. Merlot, and J. Giraudat (1997) *Plant Cell* **9**, 759–771.
26. K. Meyer, M. P. Leube, and E. Grill (1994) *Science* **264**, 1452–1455.
27. N. Bertauche, J. Leung, and J. Giraudat (1996) *Eur. J. Biochem.* **241**, 193–200.
28. D. R. McCarty, T. Hattori, C. B. Carson, V. Vasil, M. Lazar, and I. K. Vasil (1991) *Cell* **66**, 895–905.
29. J. Giraudat, B. M. Hauge, C. Valon, J. Smalle, F. Parcy, and H. M. Goodman (1992) *Plant Cell* **4**, 1251–1261.
30. F. Parcy, C. Valon, M. Raynal, P. Gaubier-Comella, M. Delseny, and J. Giraudat (1994) *Plant*

Cell **6**, 1567–1582.

31. E. Nambara, K. Keith, P. McCourt, and S. Naito (1995) *Development* **121**, 629–636.
32. F. Parcy and J. Giraudat (1997) *Plant J.* **11**, 693–702.
33. V. Vasil, W. R. Marcotte Jr., L. Rosenkrans, S. M. Cocciolone, I. K. Vasil, R. S. Quatrano, and D. R. McCarty (1995) *Plant Cell* **7**, 1511–1518.
34. C.-Y. Kao, S. M. Cocciolone, I. K. Vasil, and D. R. McCarty (1996) *Plant Cell* **8**, 1171–1179.
35. Y. Wu, J. Kuzma, E. Maréchal, R. Graeff, H. C. Lee, R. Foster, and C.-H. Hua (1997) *Science* **278**, 2126–2130.
36. F. Armstrong, J. Leung, A. Grabov, J. Brearley, J. Giraudat, and M. R. Blatt (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9520–9524.
37. Z.-M. Pei, K. Kuchitsu, J. M. Ward, M. Schwarz, and J. I. Schroeder (1997) *Plant Cell* **9**, 409–423.
38. G. Benning, T. Ehrlér, K. Meyer, M. Leube, P. Rodriguez, and E. Grill (1996) In *Abscisic Acid Signal Transduction in Plants*. Centro de Reuniones Internacionales sobre Biología Workshop 60 (R. S. Quatrano and M. Pagès, eds.), Instituto Juan March de Estudios e Investigaciones, Madrid, p. 34.
39. J. Salinas, F. Llorente, and J. M. Martínez-Zapater (1996) In *Abscisic Acid Signal Transduction in Plants*. Centro de Reuniones Internacionales sobre Biología Workshop 60 (R. S. Quatrano and M. Pagès, eds.), Instituto Juan March de Estudios e Investigaciones, Madrid, p. 35.
40. N.-H. Chua, Y. Wu, and R. Foster (1996) In *Abscisic Acid Signal Transduction in Plants*. Centro de Reuniones Internacionales sobre Biología Workshop 60 (R. S. Quatrano and M. Pagès, eds.), Instituto Juan March de Estudios e Investigaciones, Madrid, p. 55.
41. Q. Shen and T.-H. D. Ho (1995) *Plant Cell* **7**, 295–307.
42. Q. Shen, P. Zhang, and T.-H. D. Ho (1996) *Plant Cell* **8**, 1107–1119.
43. M. J. Guiltinan, W. R. Marcotte Jr., and R. S. Quatrano (1990) *Science* **250**, 267–271.
44. A. Hill, A. Nantel, C. D. Rock, and R. S. Quatrano (1996) *J. Biol. Chem.* **271**, 3366–3374.
45. T. F. Schultz, J. Medina, A. Hill, and R. S. Quatrano (1998) *Plant Cell* **10**, 837–847.
46. M. Suzuki, C. Y. Kao, and D. R. McCarty (1997) *Plant Cell* **9**, 799–807.
47. T. Iwasaki, K. Yamaguchi-Shinozaki, and K. Shinozaki (1995) *Mol. Gen. Genet.* **247**, 391–398.
48. H. Abe, K. Yamaguchi-Shinozaki, T. Urao, T. Iwasaki, D. Hosokawa, and K. Shinozaki (1997) *Plant Cell* **9**, 1859–1868.

## Absorption Spectroscopy

Light in the ultraviolet (UV) and visible (vis) range of the electromagnetic spectrum shows an energy that is equivalent to about 150 to 400 kJ/mol. Light with the appropriate energy is used to promote electrons from the ground state to an excited state. The absorption of energy from the incident light as a function of its wavelength is measured in absorption spectroscopy. Molecules with electrons that participate in delocalized aromatic systems often absorb light in the near-UV or visible region.

Absorption spectroscopy is usually performed on solutions of molecules in a transparent solvent. The absorbance of a solute depends linearly on its concentration, and therefore absorption spectroscopy is ideally suited for quantitative measurements. The spectral properties of a molecule depend on the

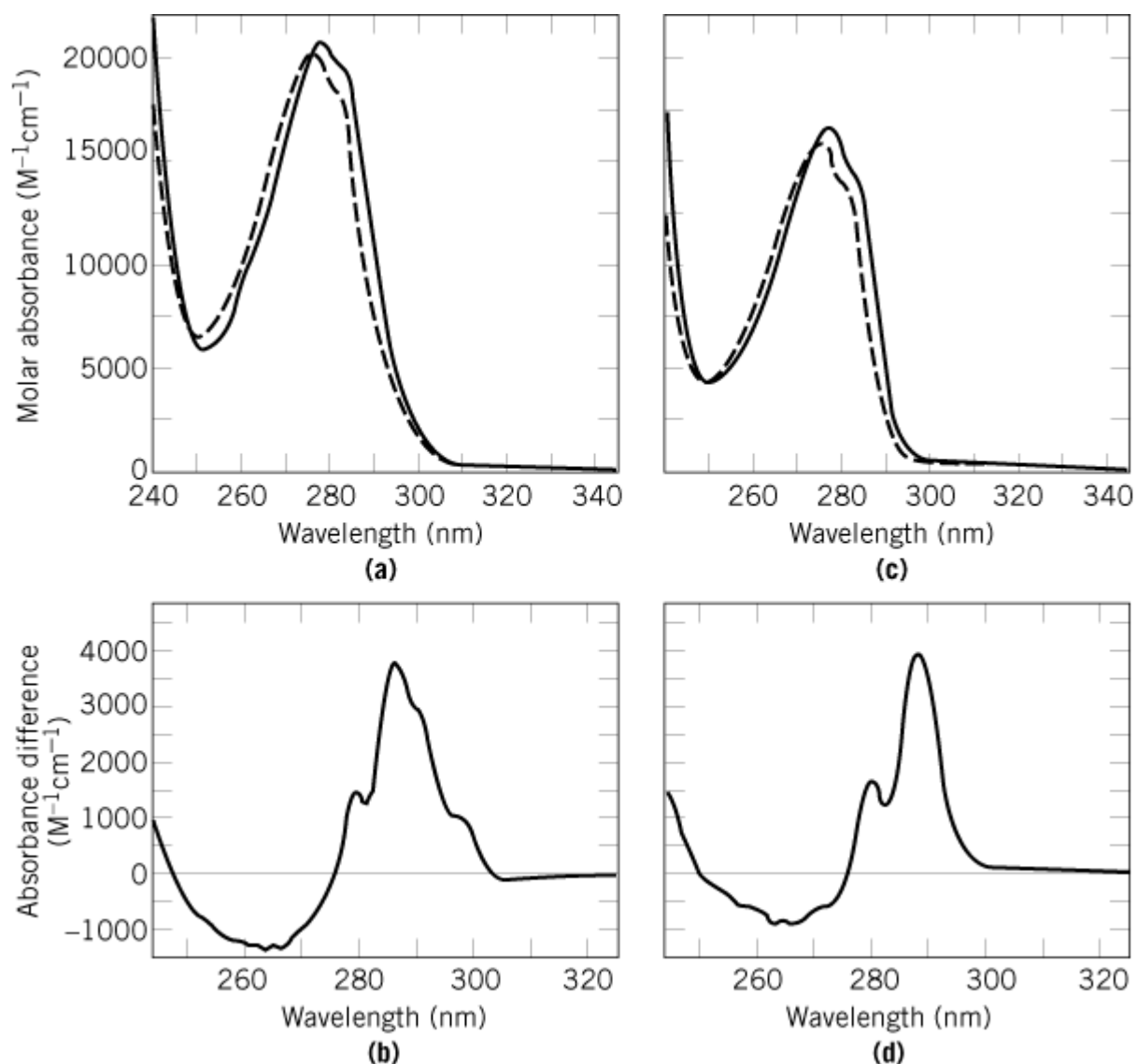
molecular environment and the mobility of its chromophores. Absorption spectroscopy and difference spectroscopy (see [Difference Spectroscopy](#)) are therefore well-suited to follow **enzyme-catalyzed reactions**, **ligand binding**, and conformational transitions in proteins and nucleic acids. Spectroscopic measurements are very sensitive, nondestructive, and require only small amounts of material for analysis. Spectrophotometers are standard laboratory equipment, and the measurement of absorbance is technically simple.

## 1. Absorbance of Proteins

In proteins, the [peptide bond](#) absorbs light in the range of 180 to 230 nm (which is called the “far-UV” range). The aromatic residues, **tyrosine** (Tyr), **tryptophan** (Trp), and **phenylalanine** (Phe), also absorb light in this region and, in addition, show bands near 260 to 280 nm (in the “near-UV”). [Disulfide bonds](#) absorb weakly near 260 nm. Some protein cofactors, such as the heme group, show absorbance in the visible range. When the peptide groups and aromatic residues are part of an asymmetric structure, or when they are immobilized within an asymmetric environment (as in folded proteins), left- and right-handed circularly polarized light is absorbed to different extents. This phenomenon is called [circular dichroism](#) (CD).

Spectra of the aromatic amino acids and model proteins are found in ([1](#)). In the near-UV, the molar absorbance of a phenylalanine residue is much smaller than that of a tyrosine or a tryptophan, and the spectrum of a protein between 240 and 300 nm is therefore dominated by the contributions from the Tyr and Trp side-chains. Phe residues contribute fine structure (“wiggles”) to the spectrum between 250 and 260 nm. The aromatic amino acids do not absorb above 310 nm, and therefore protein absorbance should be zero at wavelengths greater than 310 nm. Proteins without Trp residues do not absorb above 300 nm (see Fig. [1c](#)).

**Figure 1.** Ultraviolet absorption spectra of (a) the wild-type form and (c) the Trp59Tyr variant of RNase T<sub>1</sub>. The spectra of the native proteins (in 0.1 M sodium acetate, pH 5.0) are shown by the continuous lines. The spectra of the unfolded proteins (in 6.0 M GdmCl in the same buffer) are shown by broken lines. The difference spectra between the native and unfolded forms are shown in (b) and (d). Spectra of 15 μM protein were measured at 25°C in 1-cm cuvettes in a double-beam instrument with a band width of 1 nm at 25°C. The spectra of the native and unfolded proteins were recorded successively, stored, and subtracted. From Ref. [1](#).



A fraction of the aromatic residues are buried in the **hydrophobic** core of a native, folded protein molecule (see [Protein Structure](#)). When these residues become exposed to the aqueous solvent upon unfolding, their absorption is shifted slightly to shorter wavelengths. This is clearly seen for two forms of **ribonuclease T<sub>1</sub>** (Figs. [1a](#) and [1c](#)). The difference spectra in Figures [1b](#) and [1d](#) show that the maximal differences in absorbance occur in the 285 to 295 nm region. The **difference spectrum** for the form without a Trp residue (Fig. [1d](#)) shows a prominent maximum at 287 nm. It is typical for proteins that contain Tyr residues only. The form with a single Trp residue (Fig. [1b](#)) shows additional shoulders between 290 and 300 nm in the difference spectrum. They originate from the single Trp residue of this protein. The differences in the absorption spectra between the native and unfolded state of a protein are generally small, but they can be determined with good accuracy and are extremely useful for monitoring conformational changes of a protein.

Spectroscopic methods are, in general, the methods of choice (i) to investigate changes in the behavior of a protein under different solvent conditions, and (ii) to compare the properties of related molecules, such as homologous or mutated forms of a protein. In addition, they are widely used to measure [protein stability](#) and to follow structural transitions such as unfolding and refolding under a variety of conditions. Absorbance changes during fast reactions that occur in the milliseconds range can be followed by using rapid mixing techniques, such as in stopped-flow spectrometry (see [Kinetics](#)).

## 2. Absorbance of Nucleic Acids

Nucleic acids show a strong absorbance in the region of 240 to 275 nm. It originates from the  $\pi \rightarrow \pi^*$  transitions of the **pyrimidine** and **purine** ring systems of the nucleobases. The bases can be protonated, and therefore the spectra of DNA and RNA are sensitive to pH. At neutral pH, the absorption maxima range from 253 nm (for guanosine) to 271 nm (for cytidine), and therefore polymeric DNA and RNA show a broad and strong absorbance near 260 nm.

In native DNA, the bases are stacked in the hydrophobic core of the double helix, and therefore their absorbance is considerably decreased relative to the absorbance of single-stranded DNA, and even more so relative to oligonucleotides. This phenomenon is called hypochromism. It is widely used to follow the melting of DNA double helices.

## 3. Absorbance to Determine Concentrations

Absorbance measurements are the methods of choice to determine the concentration of proteins or nucleic acids in solution. Spectrophotometers are standard laboratory equipment, and absorbance can be measured quickly and accurately. The absorbance  $A$  is related with the intensity of the light before  $I_0$  and after  $I$  passage through the protein solution by Equation 1, and the absorbance depends linearly on concentration, according to the Lambert–Beer relationship (Eq. 2):

$$A = -\log_{10}(I/I_0) \quad (1)$$

$$A = \epsilon \times c \times l \quad (2)$$

where  $c$  is the molar concentration,  $l$  the pathlength in cm, and  $\epsilon$  the molar absorption coefficient. The concentration of a substance in solution can be determined very rapidly and accurately from its absorbance by using Equation 2. The measurement of absorbance values greater than 2 should be avoided, because only 1% of the incident light is transmitted through a solution with an absorbance of 2 (and is quantified by the photomultiplier). The value of  $\epsilon$  can be determined by a number of experimental techniques (2) or can be calculated for proteins by adding up the contributions of the constituent aromatic amino acids of a protein (2, 3).

The absorbance of nucleic acids does not vary much. The concentrations of nucleic acids in solution are routinely determined from the absorbance at 260 nm. In fact, amounts of nucleic acid are often given as “ $A_{260}$  units.” For double-stranded DNA, one  $A_{260}$  unit is equivalent to 50  $\mu\text{g}$  DNA; for single-stranded DNA, it is equivalent to 33  $\mu\text{g}$  DNA; for single-stranded RNA, it is equivalent to 40  $\mu\text{g}$  RNA. All these amounts would produce an  $A_{260}$  of 1 when dissolved in 1 mL and measured in a 1-cm cuvette. Proteins absorb much more weakly than nucleic acids. In a 1:1 mixture of nucleic acids and proteins, the proteins contribute only about 2% to the total absorbance at 260 nm. Consequently, these quantities of contaminating protein hardly affect the concentrations of nucleic acids measured by  $A_{260}$ .

## 4. Absorbance Spectrophotometers

Absorbance is measured by a spectrophotometer. Spectrophotometers consist usually of two light sources: a deuterium lamp, which emits light in the UV region, and a tungsten/halogen lamp for the visible region. After passing through a monochromator (or through optical filters), the light is focused into the cuvette, and the amount of light that passed through the sample is detected by a photomultiplier or photodiode. In diode array spectrophotometers, the sample is illuminated by the full lamp light; after passage through the cuvette, the transmitted light is spectrally decomposed by a

prism into the individual components and quantified by an array of diodes, often in intervals of 2 nm. In diode array spectrometers, the entire spectrum is recorded at the same time and not by a time-dependent scan as in conventional instruments. This is of advantage for measuring time-dependent changes at several wavelengths simultaneously (see [Kinetics](#)).

## 5. Buffers for Absorbance Spectroscopy

Good [buffers](#) for measuring difference spectra ideally should not absorb light in the wavelength range of the experiment. For work in the near-UV, buffer absorbance should be small above 230 nm, and indeed most of the solvents commonly used in biochemical experiments do not absorb in this spectral region ([1](#)). Buffer absorbance is a major problem, however, in the far-UV region below 220 nm, because buffers that contain carboxyl and/or amino groups absorb light in this wavelength range. Buffers with negligible absorbance in the far-UV include phosphate, cacodylate, and borate. Detailed procedures for the measurement of difference spectra are found in Ref. [1](#).

## Bibliography

1. F. X. Schmid (1997) "Optical spectroscopy to characterize protein conformation and conformational changes". In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 261–297.
2. S. C. Gill and P. H. von Hippel (1989) *Anal. Biochem.* **182**, 319–326.
3. C. N. Pace and F. X. Schmid (1997) "How to determine the molar absorption coefficient of a protein". In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 253–259.

## Suggestions for Further Reading

4. S. B. Brown (1980) "Ultraviolet and visible spectroscopy". In *An Introduction to Spectroscopy for Biochemists* (S. B. Brown, ed.), Academic Press, London, pp. 14–69.
5. D. B. Gordon (1994) "Spectroscopic techniques". In *Principles and Techniques of Practical Biochemistry* (K. Wilson and J. Walker, eds.), Cambridge University Press, Cambridge, UK, Chap. "7", pp. 324–380.
6. D. A. Harris and C. L. Bashford (1987) *Spectrophotometry and Spectrofluorimetry: A Practical Approach*, IRL Press, Oxford, UK

## Acceptor Stem

The acceptor stem is the site of attachment of [amino acids](#) to [transfer RNA](#) (tRNA). It is formed by 7 base pairs and has 4 single-stranded nucleotides. Nucleotides 1 to 7 from the 5' end of the tRNA base pair with nucleotides 72–66, respectively, from the 3' end of the molecule. Whereas the 5' end of the RNA has a monophosphate group, the 3' end contains a 3'-hydroxyl group, which is the site of esterification to the amino acid. The four 3'-end single-stranded nucleotides include residue 73, the discriminator base and the well-conserved nucleotides C74, C75, A76, the 3'-CCA sequence. In the three-dimensional structures of tRNA, the acceptor stem is stacked on the T arm, forming the acceptor "branch" of the L-shaped RNA fold (see Fig. 1 of [Transfer RNA](#)). The acceptor branch usually forms a regular A-type double helix, with the four single-stranded nucleotides extending in a regular helical continuity. In tRNA-like domains of some **viruses**, the acceptor stem is formed by a single RNA chain that folds into a helix due to the presence of a **pseudoknot** ([1](#)).

Specific aminoacylation of tRNA by their cognate aminoacyl-tRNA synthetases is dependent on the presence of a series of identity elements. Limited in number, these elements are preferentially located in the [anticodon](#) loop and in the acceptor stem (2-4). Residue 73, next to the CCA end, is called the “discriminator” base. The hypothesis that it contributes to discrimination of tRNA by cognate [aminoacyl tRNA synthetases](#) (5) has been largely confirmed. Residue 73 contributes strongly to specific aminoacylation of 17 different *Escherichia coli* tRNAs (6). It is the element making the major thermodynamic contribution toward aminoacylation of *E. coli* tRNA<sup>Cys</sup>, *E. coli* tRNA<sup>His</sup>, yeast tRNA<sup>His</sup> and *E. coli* tRNA<sup>Leu</sup>. Replacement of A73 to G73 in human tRNA<sup>Ser</sup> converts it to a tRNA for isoleucine (7). The crystallographic structures of two tRNA-cognate synthetase complexes suggests two possible mechanisms by which the discriminator base contributes to aminoacylation specificity. One is a direct mechanism, involving direct hydrogen interactions with the synthetase. The other is an indirect mechanism that confers a conformational change to the acceptor end of the tRNA to facilitate aminoacylation. In the complex of yeast tRNA<sup>Asp</sup>/aspartyl-tRNA synthetase, nucleotide G73 fits into the [active site](#), forming [hydrogen bond](#) interactions with side chains of the synthetase. It does not cause a conformational change in the acceptor stem (8). Alternatively, the 2-amino group of G73 of *E. coli* tRNA<sup>Gln</sup> hydrogen bonds with the phosphate oxygen of the previous nucleotide, folding the backbone of G73 back toward the 3' end of the tRNA. The formation of this fold-back hairpin enables the synthetase to open the first base pair of the acceptor stem and reach the second and third base pairs for specific interactions (9). In both complexes, additional contacts exist between the enzyme and the acceptor stems of the tRNA. The aspartyl-tRNA synthetase interacts with the tRNA<sup>Asp</sup> acceptor stem via the major groove of the RNA helix, whereas glutaminyl-tRNA synthetase interacts via the minor groove.

Acceptor stem identity-element nucleotides, as well as additional structural features, are important for synthetase recognition. Alanyl-tRNA synthetase is sensitive to both the exocyclic NH<sub>2</sub> group of G3 and local conformation of the helix due to structural characteristics of the G-U **wobble** pair (10-13). Alanine acceptance by alanyl-tRNA is modulated by additional signals within the acceptor stem, namely, A73, G1-C72, G2-C71 and G4-C69. The 5' end of histidine-specific tRNA has an additional nucleotide, residue number -1, which is a guanosine. This nucleotide, opposite the discriminator base, does form a base pair with the discriminator in many, but not all, histidyl-tRNA. This -1 guanosine is the major histidine identity element. Its influence is complemented by base pairs U2-A71 and G3-C70. Base pairs 1-72, 2-71, 3-70, and/or 4-69 contribute to aminoacylation in several cases. Glutamine identity requires a weak 1-72 base pair, G2-C71, G3-C70, in addition to signals elsewhere in the tRNA. Glycine identity is dependent on C2-G71, G3-C70 sequences, and serine identity involves G2-C71, among other signals. Methionine identity is based on the presence of A73, G2-C71, C3-G70, in addition to the CAU anticodon.

At the three-dimensional level, tRNA molecules fold into a two-domain L-shaped structure with the amino acid acceptor terminus and the anticodon at opposite ends. Minihelices (small RNA molecules containing only part of a tRNA) that mimic the amino acid acceptor domain, specifically, the acceptor stem stacked on top of the TY stem, have been shown to be useful tools for determining the contribution of the acceptor stem to tRNA function (14, 15). Minihelices containing the identity elements from a tRNA are efficient substrates for its cognate aminoacyl-tRNA synthetases. Thus, alanyl-tRNA synthetase, glycyl-tRNA synthetase, and histidyl-tRNA synthetase aminoacylate minihelices efficiently. Minihelices containing partial identity sets of six additional tRNAs are specific substrates for their cognate synthetases. Since identity elements are located very close to the CCA end, minihelices may be reduced in size to microhelices, consisting only of the amino acceptor stem closed by a loop and remain active. The smallest substrate for an aminoacyl-tRNA synthetase is derived from yeast tRNA<sup>Asp</sup> and consists of only 14 nucleotides. There are three base pairs closed by a tetraloop, plus the discriminator base and the CCA sequence (16). Double-stranded RNA duplexes and RNA/DNA heteroduplexes may also be aminoacylated (15).



The aminoacylation efficiency of minihelices is generally reduced significantly compared to that of the corresponding full-length tRNA. The amino acid charging is, however, very specific and depends on just a few nucleotides. These minimal RNA substrates are devoid of the [anticodons](#) that read the [genetic code](#) and provide the link between amino acid and **codon**. Since the relationship between the sequences and structures of the minisubstrates and the specific amino acids is maintained, there appears to be an *operational RNA code* for primitive aminoacylation. This operational code may be the primitive code from which the contemporary genetic code evolved ([17-20](#))

Acceptor stem properties are important for tRNA recognition by proteins other than synthetases. Aminoacylated initiator tRNA in **prokaryotes** is a substrate for a methionyl-tRNA transformylase that converts methionyl-tRNA<sup>Met</sup> to formylmethionine tRNA<sup>Met</sup>. Recognition of initiator tRNA by the formylase depends on the presence of methionine. Sequence and/or structural elements in the tRNA that are important for formylation by methionyl-tRNA transformylase are clustered at the end of the acceptor stem ([21](#)). The key determinants appear to be a mismatch or a weak base pair between nucleotides 1 and 72, a G-C base pair between nucleotides 2 and 71, and a C-G or, less preferably, a G-C base-pair between nucleotides 3 and 70. Mutations at G4-C69 also affect formylation kinetics slightly. In addition to the positive elements A73, G2-C71, C3-G70, and G4-C69, the occurrence of a G-C or a C-G base pair between positions 1 and 72 acts as a major negative determinant for the formylase. Formylation is a prerequisite for interaction with initiation factor IF2, which delivers the initiator tRNA to the P site of the [ribosome](#). Special structural features within the acceptor stem of eubacterial initiator tRNA contribute to their discrimination from elongator tRNA. They lack a Watson–Crick base pair between nucleotides 1 and 72 at the end of the acceptor stem. Eukaryotic initiator tRNA almost always has an A1-U72 pair, a feature not found in eukaryotic elongator tRNA.

A high-resolution [X-ray crystallography](#) structure of the ternary complex, formed by aminoacylated yeast tRNA<sup>Phe</sup>, **elongation factor** Tu from *Thermus aquaticus*, and an analog of GTP, revealed numerous contacts between the protein and the acceptor stem of the tRNA ([22](#)). These contacts include specific interactions with residue A76, with phosphates 74, 75, 67, 64, 3, and with riboses 2, 3, 63, 64, 65 along the acceptor stem and the T stem.

The stability of the 1–72 base pair also governs the degree of sensitivity of a peptidyl-tRNA to peptidyl-tRNA hydrolase. This enzyme converts the peptidyl-tRNA of *N*-acetylaminoacyl-tRNA into free tRNA plus peptides or *N*-acetyl amino acids. It is believed to play a role in the translational apparatus through the recycling of free tRNA from immature peptidyl-tRNA created by abortive protein synthesis ([23](#)).

The acceptor stem also contains specific information required during synthesis of the tRNA. [Ribonuclease P](#), an enzyme that removes extra nucleotides from the 5' end of tRNA during [tRNA biosynthesis](#), recognizes the -74CCA76- sequence in addition to nucleotides within the TY loop. Moreover, this enzyme locates its cleavage site by “measuring” the length of the helix formed by the amino acid acceptor stem fused to the TY stem. Processing at the 3' end of the tRNA, as well as the addition and repair of the conserved CCA sequence by ATP(CTP)-tRNA nucleotidyltransferase, requires information within the acceptor stem. This enzyme interacts with tRNA at the corner of the structure, where the T and D loops interact, and extends that interaction across the aminoacyl stem to the 3' end ([24, 25](#)).

## Bibliography

1. K. Rietveld, C. W. A. Pleij, and L. Bosch (1983) EMBO J. **2**, 1079–1085.
2. R. Giegé, J. D. Puglisi, and C. Florentz (1993) Prog. Nucleic Acid Res. Mol. Biol. **45**, 129–206.
3. W. H. McClain (1993) J. Mol. Biol. **234**, 257–280.
4. M. E. Saks and J. R. Sampson (1995) J. Mol. Evol. **40**, 509–518.
5. D. M. Crothers, T. Seno, and D. G. Söll (1972) Proc. Natl. Acad. Sci. USA **69**, 3063–3067.

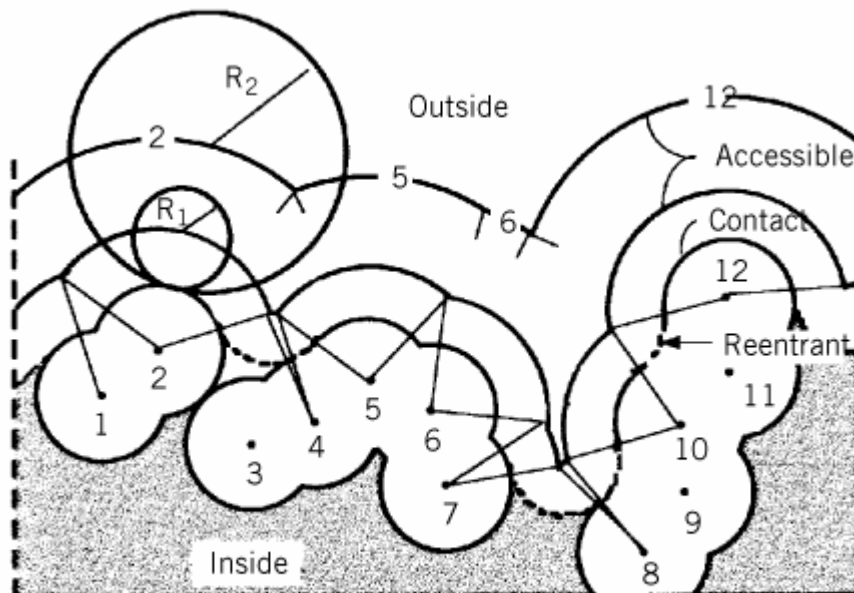
6. Y. M. Hou (1997) *Chem. Biol.* **4**, 93–96.
7. K. Breitshopf and H. J. Gross (1994) *EMBO J.* **13**, 3166–3169.
8. J. Cavarelli, B. Rees, M. Ruff, J.-C. Thierry, and D. Moras (1993) *Nature* **362**, 181–184.
9. M. A. Rould, J. J. Perona, D. Säll, and T. A. Steitz (1989) *Science* **246**, 1135–1142.
10. Y. M. Hou and P. Schimmel (1988) *Nature* **333**, 140–145.
11. Y. M. Hou and P. Schimmel (1989) *Biochemistry* **28**(17), 6800–6804.
12. W. H. McClain and K. Foss (1988) *Science* **240**, 793–796.
13. K. Musier-Forsyth and P. Schimmel (1992) *Nature* **357**, 513–515.
14. C. Francklyn, K. Musier-Forsyth, and P. Schimmel (1992) *Eur. J. Biochem.* **206**, 315–321.
15. S. A. Martinis, and P. Schimmel (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll, and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, DC, pp. 349–370.
16. M. Frugier, C. Florentz, and R. Giegé (1994) *EMBO J.* **13**, 2218–2226.
17. P. Schimmel, R. Giegé, D. Moras, and S. Yokoyama (1993) *Proc. Natl. Acad. Sci. USA* **90**, 8763–8768.
18. P. Schimmeland and L. Ribas de Pouplana (1995) *Cell* **81**, 983–986.
19. P. Schimmel (1995) *J. Mol. Evol.* **40**, 531–536.
20. P. Schimmel (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4521–4522.
21. J. M. Guillon, T. Meinel, Y. Mechulam, C. Lazennec, S. Blanquet, and G. Fayat (1992) *J. Mol. Biol.* **224**, 359–367.
22. P. Nissen, M. Kjeldgaard, S. Thirup, G. Polekhina, L. Reshetnikova, B. F. C. Clark, and J. Nyborg (1995) *Science* **270**, 1464–1472.
23. S. Dutka, T. Meinel, C. Lazennec, Y. Mechulam, and S. Blanquet (1993) *Nucleic Acids Res.* **21**, 4025–4030.
24. P. Spacciapoli, L. Doviken, J. J. Mulero, and D. L. Thurlow (1989) *J. Biol. Chem.* **264**, 3799–3805.
25. L. A. Hegg, and D. L. Thurlow (1990) *Nucleic Acids Res.* **18**, 5975–5979.

## Accessible Surface

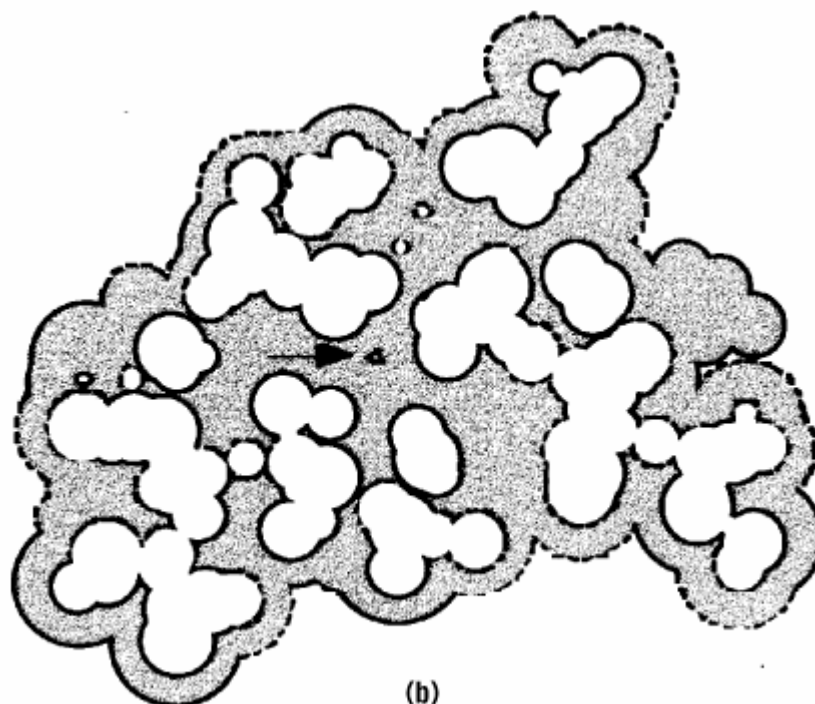
The surfaces of folded macromolecules, especially proteins, and the internal packing of their atoms have generally been analyzed using the procedure of Lee and Richards (1). The cross section of part of the surface of a native protein is depicted in Figure 1, which demonstrates a number of different surfaces and volumes (see [Protein Structure](#)). The **van der Waals surface** is defined by the spherical atoms that comprise the structure, but it is not very relevant to a folded macromolecule like a protein, where internal atoms and cavities are normally inaccessible to the solvent.

**Figure 1.** Schematic representation of possible molecular surface definitions. (a) A two-dimensional section through

part of the van der Waals envelope of a hypothetical protein including 12 atoms with the centers numbered 1–12.  $R_1$  and  $R_2$  show the radius of the probes. (b) Superposition of sections through the van der Waals and accessible surfaces of ribonuclease S. In places, the accessible surface is controlled by atoms above or below the section shown. The solid outline is the surface of carbon and sulfur atoms; the dashed outline nitrogen and oxygen. The arrow indicates a cavity inside the molecule large enough to accommodate a solvent molecule with a radius of 1.4 Å. [Taken from F. M. Richards (1977) *Ann. Rev. Biophys. Bioeng.* **6**, 151–176; B. Lee and F. M. Richards (1971) *J. Mol. Biol.* **55**, 279–400.]



(a)



(b)

The most relevant surface is the accessible surface, which is defined by its contact with molecules of the solvent. This surface is defined by rolling a spherical probe of appropriate radius  $R_1$  on the outside of the molecule, while maintaining contact with the van der Waals surface. For most

purposes, the appropriate probe is a [water](#) molecule. The accessible surface is that depicted by the center of the probe as it moves over the surface of the protein. In [Figure 1](#) the probe does not contact atoms 3, 9, or 11, and they have no accessible surface area. Such atoms are considered to be interior atoms, not part of the surface of the molecule. Those parts of the van der Waals surface in contact with the surface of the probe are designated the contact surface; they comprise a series of disconnected patches. When the probe is simultaneously in contact with more than one protein atom, its interior surface defines the reentrant surface. The contact surface and reentrant surface together make a continuous surface, which is defined as the **molecular surface**.

The accessible surface area depends on the size of the probe. When the radius of the probe increases from  $R_1$  to  $R_2$ , the number of noncontact or interior atoms in [Figure 1](#) increases from three to eight, and the accessible surface is much smoother. Thus, the smaller the probe, the larger the number of features that will be revealed. The probe is frequently taken to be a water molecule and approximated as a sphere with a radius of 1.4 Å. The accessible surface areas of individual amino acid residues are given in [Table 1](#).

**Table 1. Accessible Surface Areas of Amino Acid Residues**

| Residue Accessible surface area <sup>a</sup> (Å <sup>2</sup> ) |     |
|--|-----|
| Ala  | 113 |
| Arg  | 241 |
| Asn  | 158 |
| Asp  | 151 |
| Cys  | 140 |
| Gln  | 189 |
| Glu  | 183 |
| Gly  | 85  |
| His  | 194 |
| Ile  | 182 |
| Leu  | 180 |
| Lys  | 211 |
| Met  | 204 |
| Phe  | 218 |
| Pro  | 143 |
| Ser  | 122 |
| Thr  | 146 |
| Trp  | 259 |
| Tyr  | 229 |
| Val  | 160 |

<sup>a</sup> Values estimated for a Gly-X-Gly tripeptide in an extended conformation ([3](#)).

The structures of protein determined by [X-ray crystallography](#) indicate that the total accessible

surface area  $A_s$  of a protein is approximately proportional to the two-third power of its molecular weight. The  $A_s$  (in  $\text{\AA}^2$ ) of a typical small monomeric protein is usually related to its molecular weight  $M_w$  by the approximate relationship (2)

$$A_s = 6.3M_w^{0.73} \quad (1)$$

For oligomeric proteins

$$A_s = 5.3M_w^{0.76} \quad (2)$$

These equations are usually accurate to within  $\pm 4\%$  on average for monomers and  $5\%$  for oligomers. Equations 1 and 2 imply that an [oligomeric protein](#) has a larger accessible surface area (by 7 to 13%) than a monomeric protein of the same molecular weight in the molecular weight range up to 35,000.

For an unfolded polypeptide chain, the total accessible surface area  $A_t$  in an extended conformation is directly proportional to its molecular weight, within  $\pm 3\%$

$$A_t = 1.45M_w \quad (3)$$

Thus, for those proteins whose accessible surface area is given by Equation 1, the potential surface area buried by the protein's folding  $A_b$  will be approximately given by

$$A_b = 1.45M_w - 6.3M_w^{0.73} \quad (4)$$

This indicates that 55 to 75% of the surface area of the unfolded polypeptide chains is buried in the folding process.

The methodology of Lee and Richards (1) is also applicable to the calculation of the accessible surface area of **nucleic acids** (4). Two-thirds of the water-accessible surface area become buried on **double-helix** formation of **DNA** and **RNA**. When a probe corresponding to a single water molecule is used, the total accessible surface area is similar for [A-DNA](#) and [B-DNA](#), although marked differences appear in the major and minor groove exposures. For the larger probes, there exist considerable differences in accessible surface area between the two conformations.

## Bibliography

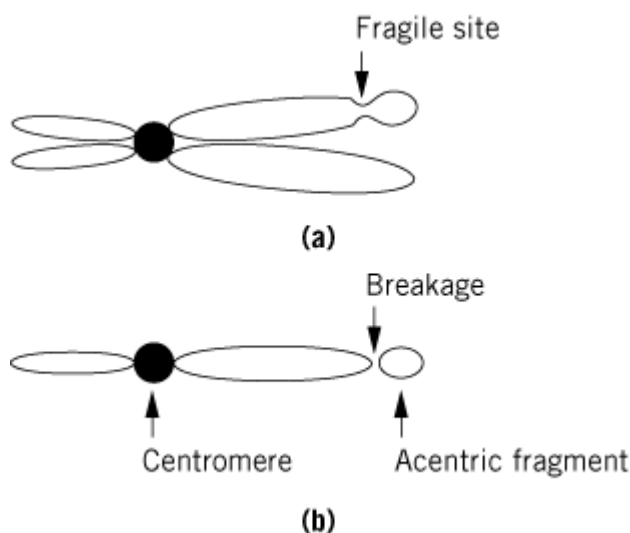
1. B. Lee and F. M. Richards (1971) *J. Mol. Biol.* **55**, 379–400.
2. J. Janin, S. Miller, and C. Chothia (1988) *J. Mol. Biol.* **204**, 155–164.
3. S. Miller et al. (1987) *J. Mol. Biol.* **196**, 641–656.
4. C. J. Alden and S. Kim (1979) *J. Mol. Biol.* **132**, 411–434.

## Acentric Fragment

Fragments of chromosomes that lack centromeres are described as acentric fragments (Fig. 1). They are formed as a result of chromosomal damage. Lacking a centromere, there is no means of

attachment to the mitotic spindle, which leads to a failure to segregate replicated chromatids in the acentric fragments to the daughter cells during cell division. Therefore, the acentric fragments are normally lost from most cells with progressive cell cycles. The breakage of the chromosome that separates the acentric fragment from the remainder of the chromosome containing the centromere can be a result of either inherent chromosomal fragility, radiation damage, or defective DNA repair mechanisms.

**Figure 1.** (a) A mitotic chromosome is shown with a fragile site, indicated. (b) Mitosis can create forces that break the chromosome at the fragile site, creating a residual chromosome that contains the centromere and an acentric fragment as indicated.



Chromosome fragile sites lead to acentric fragments in humans. The majority of recurrent cancer chromosome breakpoints are found at fragile sites (1). These breakpoints may promote oncogenic rearrangements and translocations. Chromosomal fragile sites can be induced by growing cells under abnormal conditions, by the addition of drugs or chemicals, and by infection with viruses. Fragility seems to be a consequence of incomplete compaction of the chromatin. A more open and accessible chromatin structure might explain the tendency of these sites to break, to recombine with other chromosomal regions, and to be sites of viral integration (2). A fragile site associated with a form of inherited mental retardation known as fragile X mental retardation has been characterized in some detail (5). A gene called FMR-1 (fragile X mental retardation gene 1) is present in the fragile site. This gene contains within it a trinucleotide repeat of CGG. A normal individual has between 6 and 50 of these trinucleotide repeats. Individuals with fragile X syndrome contain expansions of these repeat sequences such that they have several hundred CGG triplets. This expansion leads to chromosome fragility, to inactivation of FMR-1 gene expression, and to methylation of regulatory DNA within the gene. The expansion of the CGG trinucleotide sequences is associated with the assembly of an aberrant chromatin structure (3).

Defective chromosomal repair mechanisms lead to a much higher number of chromosomal breaks and translocations than found in normal cells. This high degree of chromosomal damage can eventually lead to malignancy. For example, in Bloom's syndrome there is a very high rate of sister chromatid exchange, and in Fanconi's anemia and ataxia telangiectasia there is a high incidence of chromosomal breakage and rejoining, leading to multicentric chromosomes and acentric fragments. Individuals with these syndromes are especially sensitive to environmental insults, such as chemical mutagens and ionizing radiation, which further increase chromosomal damage (4). Detection of acentric fragments is diagnostic of serious abnormalities in chromosomal metabolism.

## Bibliography

1. J.J. Yunis and A.L. Soreng, *Science* **226**, 1199–1204 (1984).
2. A.D. Bailey, Z. Li, T. Pavelitz, and A.M. Weiner, *Mol. Cell. Biol.* **15**, 6246–6255 (1995).
3. R.S. Hansen et al., *Cell* **73**, 1403–1409 (1993).
4. C.J. Vessey, C.J. Norbury, I.D. Hickson, *Prog Nucleic Acid Res Mol Biol* **63**, 189–221 (1999).
5. P. Jin and S.T. Warren, *Hum. Mol. Genet.* **9**, 901–908 (2000).

## Additional Reading

6. Clark M.S., and Wall W.J., *Chromosomes. The Complex Code*, Chapman and Hall, London, U.K., 1996.

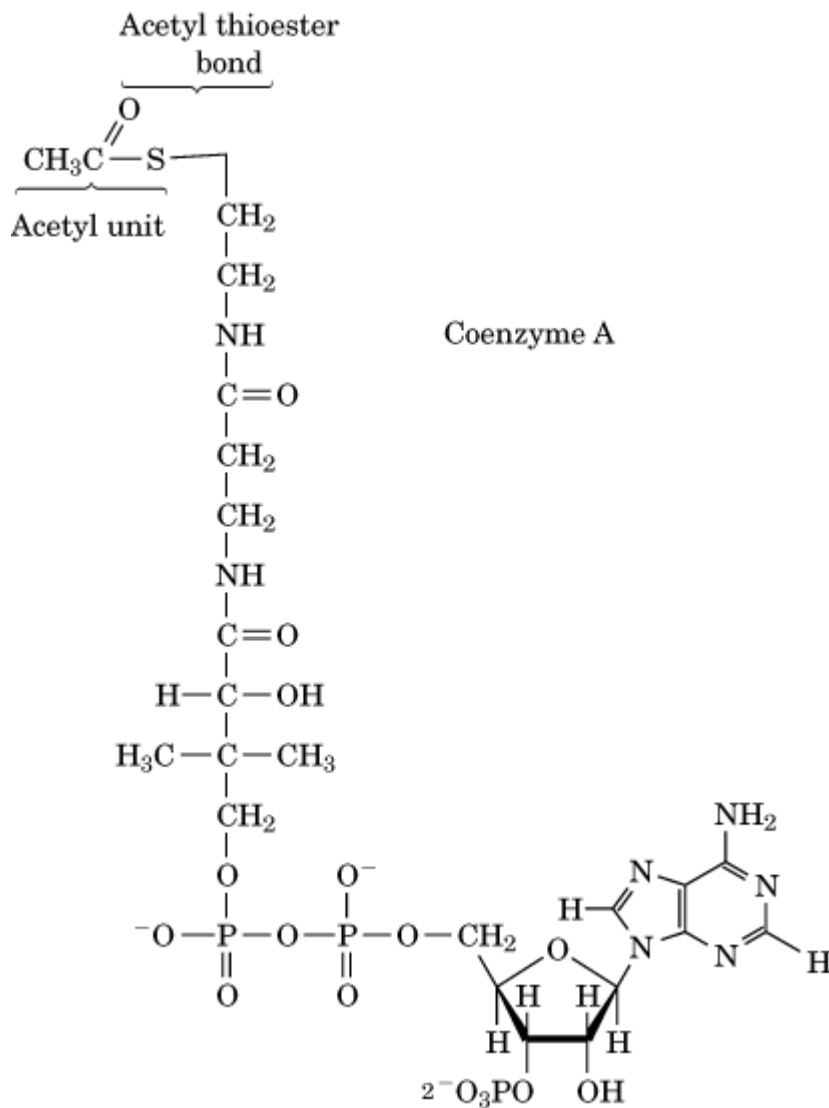
## Acetyl Coenzyme A

Acetyl coenzyme A (acetylCoA) consists of a two-carbon activated acetyl unit attached to coenzyme A in thioester linkage. AcetylCoA is central to energy generation from the degradative pathways of oxidative fuel metabolism and to a number of biosynthetic pathways that utilize the activated two-carbon acetyl unit. In aerobic cells, it is the product of all the major catabolic pathways of fuel metabolism, including  $\beta$ -oxidation of fatty acids, ketone body degradation, glycolysis and pyruvate oxidation, ethanol oxidation, and the oxidative degradation of many amino acids. The two-carbon acetyl unit of acetylCoA formed from these pathways can be completely oxidized to  $\text{CO}_2$  in the tricarboxylic acid cycle (TCA cycle), thus providing aerobic cells with energy from the complete oxidation of fuels. The acetyl unit of acetylCoA is also the basic building block of fatty acids, cholesterol, and other compounds, and it can be transferred to other molecules in acetylation reactions (eg, synthesis of N-acetylated sugars).

### 1. Structure

The ability of acetylCoA to participate in these diverse metabolic pathways is derived from the thioester bond formed between the acyl carbon of the acetyl unit and the sulfhydryl group of coenzyme A (CoASH) (Fig. 1). Because sulfur does not share its electrons, the carbonyl carbon of the thioester bond carries a more positive partial charge than that of an oxygen ester, and electrons are pulled away from the C-H bonds of the terminal methyl group. This distribution of charge facilitates nucleophilic attack at the carbonyl carbon and enhances the ability of the terminal methyl carbon to act as an electrophilic agent in condensation reactions. The acetylCoA thioester bond is a high energy bond, with a  $\Delta G^{\circ}$  for hydrolysis of  $-32.2$  kJ/mol. Transfer of the acetyl unit to other molecules, therefore, usually occurs with the release of energy.

**Figure 1.** Structure of acetyl coenzyme A (acetyl CoA).



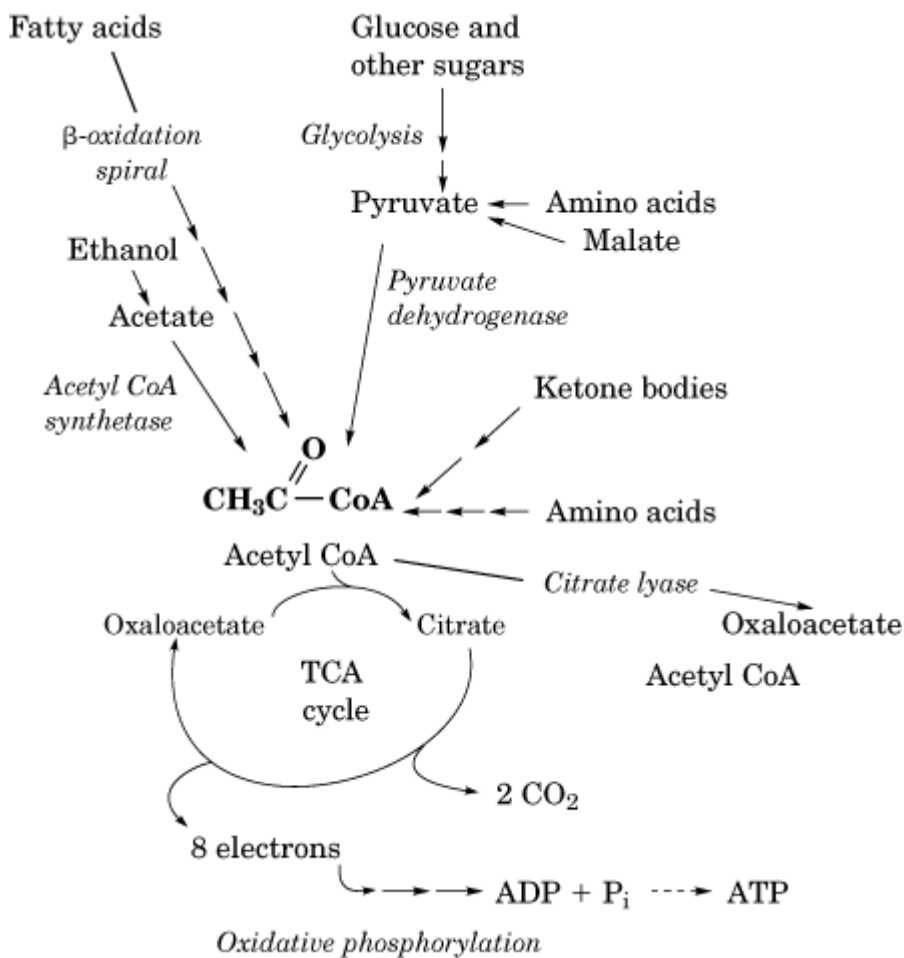
## 2. Formation of AcetylCoA in Oxidative Pathways

Three basic types of reactions exist in the oxidative pathways of fuel metabolism that generate acetylCoA: the activation of acetate, the thiolytic cleavage of  $\beta$ -ketoacyl CoAs and  $\beta$ -hydroxy acids, and the oxidative decarboxylation of pyruvate (Figure 2). In mammalian cells, acetate is the end product of ethanol metabolism and of threonine degradation. It is activated to acetylCoA in a single ATP- requiring step by the enzyme acetylCoA synthetase (Table 1). In an alternate route, bacterial cells can synthesize acetyl phosphate from acetate and then transfer the activated acetyl group to CoASH. Thiolytic cleavage of  $\beta$ -ketoacyl or  $\beta$ -hydroxy acylCoA derivatives to acetylCoA occurs in the pathways for oxidation of fatty acids, synthesis of the ketone bodies acetoacetate and  $\beta$ -hydroxybutyrate, and oxidative degradation of the amino acids isoleucine, leucine, lysine, tryptophan, phenylalanine, and tyrosine. (Most biochemistry textbooks contain general outlines of these pathways.) This type of reaction is illustrated by the  $\beta$ -oxidation spiral for fatty acids, in which one molecule of the  $\text{C}_{16}$ -fatty acid palmitate is converted to eight molecules of acetylCoA by enzymes that sequentially oxidize the molecule to a  $\beta$ -ketoacyl compound and then cleave acetylCoA from the carboxylic acid end (Table 1). The third type of reaction, the oxidative decarboxylation of pyruvate by the pyruvate dehydrogenase complex, provides the connecting link between pathways that produce pyruvate and the TCA cycle. The inhibition of the pyruvate dehydrogenase complex by acetylCoA has a regulatory role in controlling the flow of carbon into the



various pathways of intermediary metabolism.

**Figure 2.** The role of acetyl CoA in oxidative fuel metabolism.



**Table 1. Enzymes that Form Acetyl CoA**

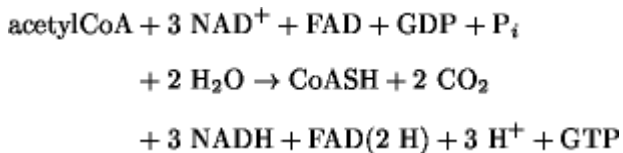
---

|  |
|--|
| <i>AcetylCoA synthetase</i>  |
| (1) acetate + ATP + CoASH $\longrightarrow$ acetyl CoA + AMP + PP <sub>i</sub>   |
| <i>β-Ketothiolase</i>  |
| (2) $\text{CH}_3(\text{CH}_2)_n\text{CH}_2\overset{\text{O}}{\parallel}\text{C}-\overset{\text{O}}{\parallel}\text{C}-\text{CH}_2-\text{SCoA} \longrightarrow \text{CH}_3(\text{CH}_2)_n\text{CH}_2\overset{\text{O}}{\parallel}\text{C}-\text{SCoA} + \text{CH}_3\overset{\text{O}}{\parallel}\text{C}-\text{SCoA}$ |
| <i>Pyruvate dehydrogenase complex (thiamine-PP, lipoate, FAD)</i>  |
| (3) pyruvate + CoASH + NAD <sup>+</sup> $\longrightarrow$ Acetyl CoA + NADH + H <sup>+</sup> + CO <sub>2</sub>   |
| <i>Citrate lyase</i>   |
| (4) citrate + CoASH + ATP $\longrightarrow$ acetyl CoA + oxaloacetate + ADP + P <sub>i</sub>   |

---

### 3. Oxidation of the Acetyl Unit of AcetylCoA in the TCA Cycle

It is estimated that about two thirds of the energy requirements of aerobic cells is met by the complete oxidation of the acetyl group of acetylCoA to CO<sub>2</sub> in the TCA cycle (also called the citric acid cycle or the Krebs cycle.) In this cyclical sequence of reactions, acetylCoA condenses with oxaloacetate to form citrate in a reaction catalyzed by the enzyme citrate synthase. Subsequent reactions of the cycle donate electrons to the **coenzymes NAD<sup>+</sup>** and **FAD** for ATP generation from oxidative phosphorylation, regenerate oxaloacetate, and release two carbons as CO<sub>2</sub>. The net reaction of the TCA cycle is:

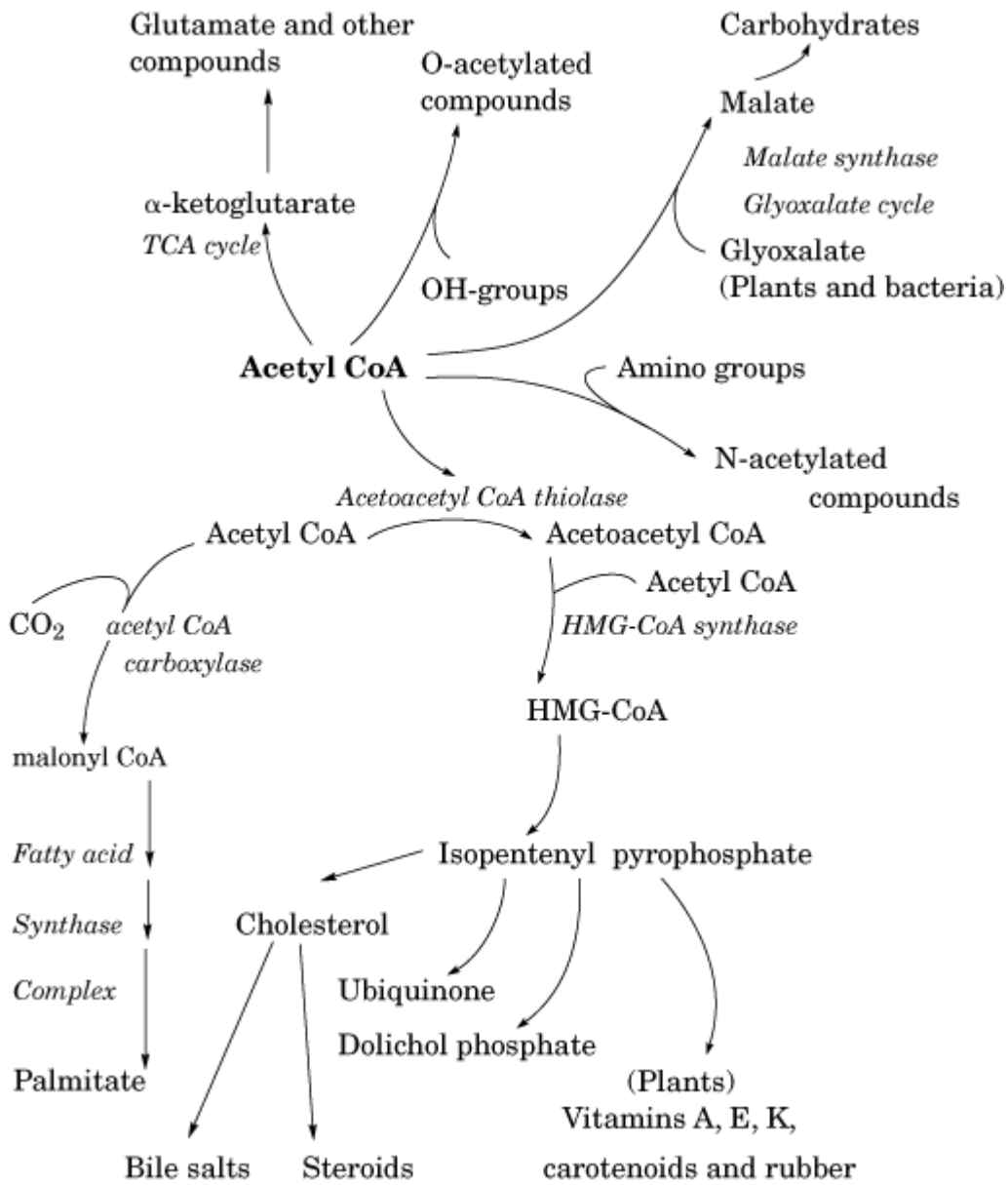


Most of the pathways that produce acetylCoA in aerobic cells of **eukaryotic** organisms are located in the **mitochondrial** matrix, and most of the biosynthetic pathways that utilize acetylCoA are outside of the mitochondrion. AcetylCoA is not directly transported through the inner mitochondrial membrane but is “transferred” from the mitochondrial matrix to the **cytosol** as citrate. It is then regenerated in the cytosol by the enzyme citrate lyase.

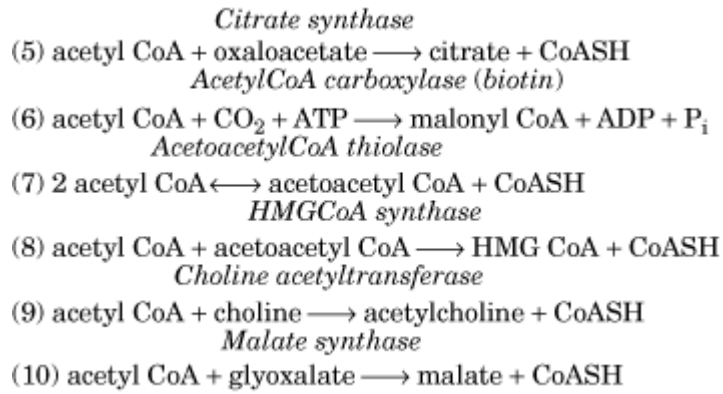
### 4. AcetylCoA as a Precursor in Biosynthetic Reactions

The two-carbon acetyl unit of acetylCoA is the precursor of a number of compounds synthesized in cells. It is the basic building block of fatty acids, cholesterol, and other compounds derived from the five-carbon isoprenoid unit (Figure 3). In the synthesis of fatty acids, acetylCoA is carboxylated to malonylCoA by the **biotin**-requiring enzyme, **acetylCoA carboxylase** (Table 2). Subsequent reactions build the C<sub>16</sub> fatty acid palmitate and other fatty acids from successive additions of the portion of malonyl CoA derived from acetylCoA, thereby providing the cell with the diverse fatty acids required for membrane lipids. In the synthesis of cholesterol, acetyl units from three acetylCoA molecules condense to form 3-hydroxy 3-methylglutaryl CoA (HMG CoA), which is subsequently decarboxylated and converted to the five-carbon isoprenoid unit isopentenyl pyrophosphate. This isoprenoid unit is one of the most common structural units of a number of compounds in mammalian cells, bacteria, and plants. AcetylCoA also contributes one carbon to compounds synthesized from the five-carbon intermediate  $\alpha$ -ketoglutarate in the TCA cycle. The acetyl unit of acetylCoA can also be transferred either to hydroxyl groups of compounds to form an acetyl oxygen ester (eg, the neurotransmitter acetylcholine; see [Acetylcholine Receptor](#)) or to an amino group to form an amide in an N-acetylation reaction (eg, acetylated amino sugars such as N-acetylglucosamine).

**Figure 3.** The central role of acetyl CoA in biosynthetic pathways.



**Table 2. Enzymes That Utilize AcetylCoA**



In mammalian cells, acetylCoA cannot provide a net source of carbon for the synthesis of glucose or other sugars. However, plants, yeast, and bacteria contain the glyoxylate cycle, which serves as a bypass of the TCA cycle. The net result of the glyoxylate cycle is the conversion of two molecules of acetylCoA to succinate, utilizing TCA cycle enzymes plus isocitrate lyase and malate synthase. As a result, *Escherichia coli* and many other bacteria are able to convert acetate to carbohydrates and amino acids and can thus utilize and grow on acetate as their sole carbon source.

#### Suggestions for Further Reading

- A. L. Lehninger, D. L. Nelson, and M. M. Cox (1990) *Principles of Biochemistry*, 2nd ed., Worth Publishers, New York. Although the pathways mentioned above can be found in more detail in most textbooks, this book is particularly good for outlines of major pathways found in mammals, bacteria, and plants.
- D. B. Marks, A. D. Marks, and C. M. Smith (1996) *Basic Medical Biochemistry: A Clinical Approach*, Williams & Wilkins, Baltimore, MD. Provides a general outline of pathways found in humans and their relationship to human physiologic and pathologic conditions.
- R. H. Bethal, D. B. Buxtion, J. G. Robertson, and M. S. Olson (1993) Regulation of the pyruvate dehydrogenase multienzyme complex. *Ann. Rev. Nutr.* **13**, 497–520. This article and the following one provide further information about the pyruvate dehydrogenase complex and acetylCoA carboxylase, two of the key enzymes in acetylCoA metabolism, which are highly regulated and the subject of current research.
- G. M. Mabrouk, I. M. Helmy, K. G. Thampy, and S. J. Wakil (1990) Acute hormonal control of acetylCoA carboxylase: the roles of insulin, glucagon, and epinephrine. *J. Biol. Chem.* **265**, 6330–6338.
- J. L. Goldstein, H. H. Hobbs, and M. S. Brown (1995) "Familial hypercholesterolemia". In *The Metabolic and Molecular Bases of Inherited Disease*, 7th ed. (C. R. Scriver, A. L. Beaudet, W. S. Sly, and D. Valle, eds.), McGraw Hill, New York: pp. 1981–2030. Covers certain aspects of pathway regulation for the synthesis of cholesterol from acetylCoA, which is of great interest for protection from heart disease.

#### Acetylcholine Receptor

Acetylcholine (ACh) is a widely distributed neurotransmitter in both the peripheral and central nervous systems (1). It is synthesized by the [enzyme](#) choline acetyltransferase within the cholinergic nerve terminal, where it is packaged into cholinergic vesicles. Arrival of an action potential at a cholinergic nerve terminal triggers quantal release of ACh by fusion of these vesicles with the presynaptic membrane and concomitant [exocytosis](#) of their contents into the synaptic cleft. The ACh released **diffuses** across the synapse to the postsynaptic membrane, where it activates its target, the acetylcholine receptor (AChR). Termination of transmission at cholinergic synapses involves rapid hydrolysis of ACh, by the enzyme acetylcholinesterase (AChE), which hydrolyzes it to choline and acetic acid (2).

The early observations of Dale (3) showed in various pharmacological preparations that ACh evoked responses similar to those evoked by either nicotine or muscarine. This provided the basis for the grouping of all AChRs into two families, nicotinic (nAChRs) and muscarinic (mAChRs). Dale's pioneering classification, based on the action of these two plant alkaloids, is still valid and used, although subtypes have been recognized and defined within each of the two families, both on the basis of specificity with respect to binding of agonists and antagonists and, in recent years, on the basis of gene [cloning](#). Two additional plant alkaloids, *d*-tubocurarine and atropine, also provide useful tools by serving as specific antagonists of nAChRs and mAChRs, respectively (4).

nAChRs are ligand-gated **ion channels**, and their activation, by ACh or other agonists, causes a rapid change in ion permeability of the membrane in which they are embedded. mAChRs are members of the family of [G-protein-coupled receptors](#) and produce much slower responses, either excitatory or inhibitory, via their corresponding [second messengers](#) (5). Initial molecular cloning studies provided the [primary structures](#) of both nAChRs (6) and mAChRs (7), revealing that they belong to distinct families of proteins that share neither sequence identity nor a similar fold. As will be discussed in detail below, nAChRs are composed of pentamers of one or more subunits that display substantial sequence [homology](#) and, most likely, similar overall folds and transmembrane topologies. Muscarinic receptors are glycoproteins of molecular weight ~80,000 which, as already mentioned, belong to the family of G-protein-coupled receptors (8). Just as for other members of this family, [hydrophobicity](#) plots predict seven **transmembrane** sequences (8). The nAChR belongs to a larger family of ligand-gated ion channels, with which it too appears to share an overall fold, a common pattern of transmembrane sequences and other structural similarities (9, 10; see text below).

The recognition that the large electric organs of electric fish such as the electric eel, *Electrophorus electricus*, and the electric ray, *Torpedo* sp., provide a highly enriched and homogeneous preparation of cholinergic synapses (11) led to their use as an experimental model and also as the source of choice for purification and characterization of the nAChR, of AChE, and of homogeneous cholinergic synaptosomes and synaptic vesicles (2). The fact that the cholinergic synapse at the neuromuscular junction was the parallel system of choice for studying the electrophysiological characteristics of synaptic transmission (12), as well as many ultrastructural and developmental parameters (for recent reviews see Refs. 13-15), meant that many of the structural data garnered by use of the electric organ model system could be directly correlated with the functional data obtained with the neuromuscular preparation (2). Thus, the nAChR of electric organ closely resembles that of muscle (16), and the molecular forms of AChE present in electric organ tissue are homologous to those found at muscle endplates (17, 18).

The finding that injection into rabbits of purified nAChR elicited an **autoimmune** condition closely resembling the human muscle disease myasthenia gravis (19) was followed by detection of circulating [antibodies](#) to the nAChR in patients suffering from the disease. This opened up a fertile area of research in which clinical aspects were closely coupled to advances in basic research on the nAChR (20).

The cholinergic synapse has thus served as the prototypic synapse for understanding many fundamental aspects of synaptic transmission at chemical synapses. Even now, when improved

electrophysiological techniques, and approaches such as  $\text{Ca}^{2+}$  imaging, allow better access to the CNS, and when genetic engineering allows expression of the receptor of choice in a system like the *Xenopus* oocyte that is amenable to user-friendly patch-clamp techniques, as well as in sufficient amounts to carry out structural characterization, the neuromuscular junction and electric organ tissue still occupy a position at the center of the stage.

## 1. The Nicotinic Acetylcholine Receptor

### 1.1. Identification and Characterization

As early as 1955, it was suggested by Nachmansohn (21) that the nAChR might be a protein in which a conformational change elicited by the neurotransmitter ACh could induce a change in permeability that would trigger the postsynaptic response, and it was in his laboratory, in the 1960s, that the first evidence was presented that the nAChR is indeed a protein. Thus, Karlin and Bartels (22), using the electric eel electroplaque preparation, showed that the response of the nAChR to carbamylcholine could be modified reversibly by reagents that either modified [thiol groups](#) or reduced [disulfide bonds](#). Subsequently, Karlin and co-workers (23, 24) and Changeux et al. (25), also in the Nachmansohn laboratory, demonstrated *in situ* [affinity labeling](#) of the nAChR in the intact electroplaque. O'Brien and co-workers (26, 27), using specific binding of the reversible ligand muscarone and homogenates of electric organ tissue of *Torpedo*, were able to obtain a good estimate (~2 nmol/g tissue) of the number of nicotinic binding sites.

A major breakthrough in the identification of the nAChR, its localization, and its molecular characterization came about by use of (a) the neurotoxin **a-bungarotoxin** (aBgt), a 74-residue polypeptide chain purified by Lee and co-workers from the venom of the banded krait, *Bungarus multicinctus* (28), and (b) the homologous **a-neurotoxins** purified from cobra venom (29). Lee's group presented electrophysiological and pharmacological evidence that these neurotoxins acted as very high-affinity (dissociation constant in the pico-to-femtomolar range) antagonists of the nAChR (30).  $^{125}\text{I}$ -aBgt was developed as a powerful tool for localization and quantification of the nAChR by high-resolution [autoradiography](#) (31). It was shown to be present as a dense quasi-crystalline array at the top of the folds of the postsynaptic membrane of skeletal muscle, opposite the putative ACh-release sites on the presynaptic membrane, with a density very similar to that in purified preparations of postsynaptic membranes obtained from *Torpedo* electroplax. A further step forward was taken when it was shown that binding of  $^{125}\text{I}$ -aBgt was retained after solubilization in nonionic [detergents](#), such as Triton X-100 and cholate, and that a single toxin-binding component, migrating at ~9 S, could be identified on sucrose gradient [sedimentation velocity centrifugation](#) (32, 33). Furthermore, the receptor so solubilized could be purified by [affinity chromatography](#) using resins to which either the a-neurotoxin (34) or a quaternary nitrogen ligand (35) had been attached. This permitted characterization of the nAChR as a multisubunit protein. Thus it was shown that it was a pentamer, of molecular weight ~250,000 (36), that contains four polypeptide subunits,  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ , with apparent molecular weights of 39,000, 48,000, 58,000, and 64,000 and a stoichiometry of  $\alpha_2\beta\gamma\delta$ . In the [electron microscope](#), such purified preparations displayed a rosette-like appearance (37), closely resembling similar, highly organized arrays of rosettes observed in purified postsynaptic membrane preparations (38). This led to the working hypothesis, still currently accepted, that the five subunits lie in the plane of the plasma membrane, surrounding the ion channel running through the center. A fifth polypeptide, the 43-kDa polypeptide (now known as rapsyn), which was found to be weakly associated with the pentamer (39), is believed to be involved in nAChR-clustering and anchoring to the [cytoskeleton](#) at the synapse (13). In skeletal muscle the  $\epsilon$  subunit is present in the nAChR of embryonic muscle, but is replaced by the  $\delta$  subunit in adult muscle (40).

Studies on model systems were important in establishing the role of the nAChR in [signal transduction](#) at the cholinergic synapse. Thus the microsac preparation (41) was used to show that ACh could induce a permeability change in such nAChR-enriched vesicles. Furthermore, prolonged exposure to high concentrations of ACh was shown to produce a closed state, with high affinity for ACh (42), equivalent to the desensitized state first described in skeletal muscle by Katz and Thesleff

(43). The availability of purified native receptor permitted reconstitution studies into liposomes and into lipid bilayers, permitting recording of single channels induced in the presence of ACh (44, 45). The purified pentamer thus contained not only the ACh-binding site, but also the ion channel and the transducing elements involved in activation and desensitization (16).

Affinity labeling studies by Karlin and co-workers (46), using a radioactive reagent, established that the disulfide bond that was susceptible to reducing reagents, as originally demonstrated by Karlin and Bartels (22), was located on the  $\alpha$  subunit, in proximity to the binding site for ACh and various quaternary ligands.

An important step forward was the finding of Raftery et al. (47) that the NH<sub>2</sub>-terminal sequences of all four subunits display substantial **sequence homology**. When, not long after, the subunits were cloned (6), it was found that such homology extends throughout the whole polypeptide chain, and the receptor pentamer can be viewed as displaying pseudo-fivefold symmetry. One issue that still remains open is the arrangement of the subunits around the lumen. This has been approached primarily by electron microscopy using subunit-specific **antibodies**, with supplementary information coming from use of affinity labels directed toward the ACh-binding site, which label both the  $\alpha$  subunit and an adjacent subunit (see text below). As many as 12 permutations are possible, but it is generally accepted that the two  $\alpha$ -subunits are not adjacent to each other (48), and discussion focused on which of the  $\beta$ -,  $\gamma$ -, and  $\delta$ -subunits is flanked on both sides by the  $\alpha$ -subunit. Karlin et al. (49) proposed that it is the  $\gamma$ -subunit that lies between the two  $\alpha$ -subunits. Although Kubalek et al. (50) proposed that the  $\beta$ -subunit might occupy this position, the assignment of Karlin and co-workers is generally favored (for detailed discussions see Refs. 10 and 51).

Advances in cloning and expression, taken together with development of the patch-clamp technique, had a dramatic influence on research on the nAChR, just as they did on research on other receptors and on ion channels. In the case of the nAChR, however, the availability of large amounts of highly purified receptor from *Torpedo* electric organ resulted in fruitful synergy between these new techniques and the techniques of protein chemistry and structural biology.

In general, research has proceeded on two fronts, one concerned with the structure of the ACh-binding site and the other with that of the ion channel, with the eventual objective of understanding the physical mechanism by which ligand-binding causes channel opening. In the following, the overall topology of the receptor and of the individual subunits will first be reviewed briefly. The current status of our knowledge of the ACh-binding site, and then of the ion channel, will be summarized, followed by a discussion of the recent structural work from the laboratory of Unwin, which is beginning to give us a first glimpse of how the receptor may be functioning. For a number of recent reviews that cover these issues in more detail than is possible here, see Refs. 9, 10, 16, and 52, as well as the recent paper of Unwin (53), which summarizes the current status of the structural work.

## 1.2. Conformation and Topology

Analysis of the sequences of the four subunits, obtained on the basis of their **complementary DNA** sequences, showed that they shared a very similar topology: A long extracellular NH<sub>2</sub>-terminal domain is followed by three putative transmembrane  **$\alpha$ -helices**, identified on the basis of hydrophobicity, M1–M3, by a cytoplasmic loop containing ~100–150 residues, and finally, by a fourth transmembrane helix, M4, so that the COOH-terminus is believed to be on the extracellular side of the postsynaptic membrane (54). Other ligand-gated ion channels that belong to the same superfamily, including receptors for  $\gamma$ -aminobutyric acid, glycine, and serotonin, display a similar topology (9, 10). A number of consensus sites for **N-Glycosylation** are present in the extracellular domains, giving rise to a sugar content of ~7% (for literature see Ref. 55); and a number of consensus sites for **phosphorylation**, present on the cytoplasmic loops, presumably fulfill regulatory functions (56).

Photoaffinity labeling, using lipid-soluble probes, has been used to assign the arrangement of the transmembrane sequences in relation to the lipid bilayer and the putative central ion channel (57-59). From such studies, it has been concluded that the M4 sequence, which is also the least conserved, is the most exposed to the lipid environment, with the M1 and M3 sequences displaying more limited exposure to the lipid environment, and M2 facing the lumen of the ion channel. This assignment of M2 as making the principal contribution to the lining of the ion channel is also supported by labeling by channel blockers and by use of the SCAM technique (see text below).

Various spectroscopic techniques show that the nAChR contains substantial amounts of both  $\alpha$ -helices and  $\beta$ -sheets (for literature see Ref. 10). In particular, Naumann et al. (60), using Fourier transform infrared (FTIR) spectroscopy on native receptor-rich membranes, have shown a predominance of  $\beta$  structure, 36 to 43%, and an  $\alpha$ -helical content of 32 to 33%. Of particular interest are the putative conformations of the membrane-spanning sequences. Görne-Tschelnokow et al. (61) attacked this problem directly, by performing FTIR spectroscopy on membrane preparations from which the extramembrane domains had been shaved by **proteolysis**. Their data suggested a **beta-sheet** content of ~40% for the transmembrane sequences. Modeling studies also suggested a substantial  $\beta$ -sheet component (62). Various labeling studies, which attempt to correlate the degree of labeling of residues with orientation toward either the ion-channel lumen or the lipid bilayer, and thus with either an  $\alpha$ -helical or a  $\beta$ -strand pitch, yield complex results (see, for example, Refs. 9, 10, 59, and 63); and precise assignments will, most likely, have to await a high-resolution 3D protein structure.

### 1.3. ACh-Binding Site

Identification of the residues comprising the ACh-binding site of the nAChR was approached by affinity labeling. Karlin and co-workers used the same radioactive affinity label that they had used to label the  $\alpha$  subunit selectively, a quaternary derivative of maleimide, to identify the residues labeled. Because they had demonstrated that labeling occurred subsequent to reduction of the receptor, it was not surprising that they found that the residues labeled were the adjacent cysteines, Cys192 and Cys193, which apparently form an intrachain disulfide bond in the native nAChR, and which are present only in the  $\alpha$ -subunit (64). As predicted, these residues are localized in the sequence that had been assigned to the extracellular domain, and would thus be the natural candidate for the ACh-binding site. The importance of this region for agonist and antagonist binding was confirmed by the observation that short synthetic peptides containing Cys192 and Cys193 (eg, from residue 185 to 196) displayed specific binding to  $\alpha$ Bgt, albeit with much lower affinity than the native receptor (65). Moreover, the  $\alpha$ -subunits of the nAChR of both snakes and of the mongoose, which are resistant to  $\alpha$ Bgt, display mutations in this region that can provide a structural basis for resistance (66).

Changeux and co-workers conducted an extensive study of the residues labeled by the tertiary photoaffinity label, DDF (67). In addition to Cys192 and Cys193, this probe labeled primarily aromatic residues, including Tyr93, Trp149, Tyr190 and Tyr198, all, again, in the putative cytoplasmic domain. Use of other labeling agents, by other laboratories, broadly confirmed these assignments and identified additional aromatic residues (for literature see Ref. 68). Thus the ACh-binding pocket of the nAChR can be viewed as an “aromatic basket” (16) in which the quaternary group of ACh interacts with these aromatics via the  $p$  electron–cation interactions (69) which X-ray **crystallographic** studies have shown to play a prominent role within the “aromatic gorge” of AChE (70, 71).

Although the studies discussed above indicate a prominent role for the  $\alpha$ -subunit in binding of both agonists and antagonists, numerous labeling studies have revealed contributions of adjacent subunits to ligand-binding. These studies have been critically reviewed by Hucho et al. (10). Taken together, the data clearly indicate that the two ACh-binding sites are at the interfaces between the two  $\alpha$ -subunits and the  $g$ - and  $d$ -subunits. This, in turn, implies structural difference between the two sites that can serve to explain nonequivalence of the two binding sites in terms of affinity for both  $\alpha$ -neurotoxins (72) and antagonists (73).



#### 1.4. Ion Channel

Pharmacological studies on the nAChR revealed a number of noncompetitive antagonists, both natural and synthetic, that blocked the change in permeability elicited by ACh without inhibiting the binding of ACh itself. It was, therefore, suggested that at least some of these compounds act by entering the ion channel and sterically blocking ion movement through it (for literature see Refs. [9](#), [10](#), and [16](#)). Based on the currently held view that the ion channel is the hole in the rosette seen in the electron microscope, such “channel blockers” are believed to exert their action by plugging this hole. Some channel blockers, such as chlorpromazine ([74](#)) and trimethylphosphonium ([75](#)), bind irreversibly upon UV irradiation; and their labeled derivatives can, accordingly, be used to identify putative elements of the channel lining. Rapid-mixing studies, using nAChR-enriched microsacs, showed that exposure to ACh leads, initially, to greatly increased labeling, followed by diminished incorporation which is consistent with desensitization as measured by ion flux measurements ([76](#)). Such labeling studies showed dominant labeling of the residues in the M2 transmembrane sequences of all four subunits, suggesting that they make the principal contribution to the lining of the lumen of the ion channel, in the open state of the channel. One such noncompetitive channel blocker, quinacrine azide, was, however, shown to label amino acid residues in M1 ([77](#)).

In parallel to such labeling studies on the purified receptor, use of [site-directed mutagenesis](#), combined with expression in *Xenopus* oocytes ([78](#)) and use of the patch-clamp technique, permitted assessment of the contributions of individual amino acid residues to the conductance properties of the ion channel, an approach pioneered by the joint efforts of the Numa and Sakmann laboratories (see, for example, Refs. [79](#) and [80](#)).

A third approach, developed more recently, is the substituted-cysteine-accessibility method, abbreviated as SCAM ([81](#)), which combines site-directed mutagenesis with chemical modification, so as to identify the amino acid residues lining the ion channel and their involvement in function. Thus, residues believed to line the ion channel (primarily from the M2 transmembrane sequence, but also from M1) are replaced one-by-one by cysteine residues, using site-directed mutagenesis, and expressed in *Xenopus* oocytes. The reactivity of their thiol groups with thiol reagents bearing negative or positive charges is examined *in situ*, both in the presence and absence of ACh, thus permitting assessment of their accessibility in both the closed and open states of the channel. Furthermore, these sulfhydryl reagents can be added both intra- and extracellularly, yielding information concerning accessibility of channel-lining residues from both surfaces ([63](#)).

These various experimental approaches have led to the identification of residues within the M2 sequences of all five subunits that determine the conductance and selectivity of the nAChR channel. In particular, using nAChR mutants expressed in *Xenopus* oocytes, three rings of negatively charged residues which may be referred to as the extracellular (or outer), the intermediate, and the cytoplasmic (or inner) ring, have been shown to play important roles in determining channel conductance ([80](#)). A ring of polar (serine and threonine) residues, named the central ring, is located above the intermediate ring, forming a constriction that may serve as part of the selectivity filter ([82](#), [83](#)). The bulk of the transmembrane sequence of M2 lies between this central ring and the outer ring.

A more detailed topographical picture has been obtained by the SCAM technique, as summarized by Karlin and Akabas ([9](#)). Thus, broadly speaking, the residues in M2 of the  $\alpha$ -subunit can be fitted to an [Amphipathic](#) helix so that, with one exception, all the residues on one face of the helix are labeled in the SCAM protocol; some are equally accessible in the presence and absence of ACh, some are more accessible in its presence, and some are more accessible in its absence. Furthermore, certain residues at the extracellular end of the M1 helix of the  $\alpha$ -subunit are also labeled. The pattern of labeling of M2 in the absence of ACh is consistent with an  $\alpha$ -helical conformation, except for a short stretch in the middle (aLeu250 to aSer252), while in the presence of ACh it is consistent with an uninterrupted  $\alpha$ -helix. The pattern of labeling of M1 cannot be ascribed to a specific conformational motif, although it too is changed in the presence of ACh. Karlin and Akabas ([9](#)) hypothesize that a movement of M1 and M2 relative to each other, with a concomitant change in **secondary structure**, may flip open the gate of the channel. More recently, the SCAM technique has also located the gate,

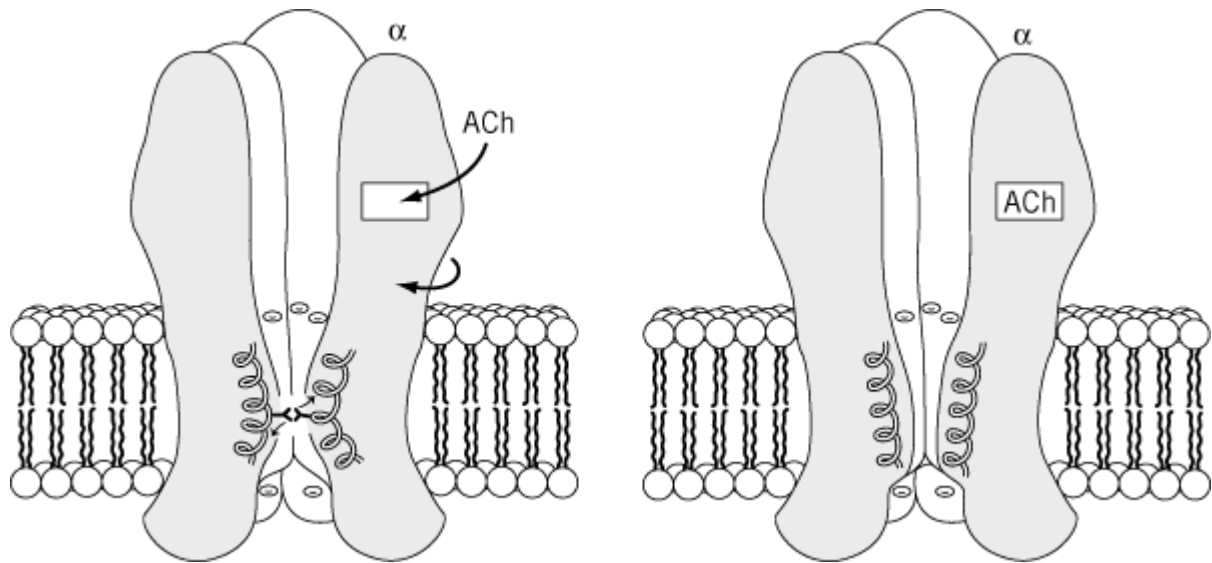
using comparative studies in which labeling was carried out from either the extracellular or cytoplasmic sides (63). In the absence of ACh—that is, in the closed state of the channel—there is a barrier to the sulfhydryl agents, when added to either side, reacting with residues aGly240 and to aThr244. ACh binding removes this barrier, which serves as an activation gate. Residues aGly240, aGlu241, aLys242, and aThr244 line a narrow part of the channel in which this gate is located. It should be noted that the second of these residues, aGlu241, belongs to the inner anionic ring, and aThr244 to the central ring, as defined above (80, 82).

### 1.5. Structural Studies

Despite the extensive studies performed on the characterization of the ACh-binding site and of the ion channel, the question with which all structural and functional studies on the nAChR are ultimately concerned is how binding of ACh is transduced into a permeability change in the ion channel. This is ultimately a problem best addressed by a structural biology approach. In spite of attempts from many laboratories to obtain nAChR crystals that would diffract X-rays (for a published example, see Ref. 84), this has not yet been achieved. Our current structural knowledge is, therefore, derived from the work of Unwin (for a recent summary see Ref. 53). As mentioned above, the postsynaptic membrane of *Torpedo* electric organ contains densely packed, partially crystalline arrays of nAChR (38). When isolated, such membranes have a natural propensity to recrystallize in tubular form (85, 86). Since the early 1980s, Unwin has been using electron microscopy to probe the structure of the nAChR from *Torpedo marmorata* and is now able to compare the receptor in its open and closed states at 9 Å resolution (53). The receptor is revealed as an elongated structure, ~125 Å long; the extracellular portion of each subunit extends ~60 Å above the membrane surface, and the intracellular portions extend ~20 Å from the cytoplasmic surface. A density underlying each receptor most likely represents the cytoskeletal protein rapsyn (see text above), normally present in 1:1 stoichiometry with the receptor (87). A view from above shows the five subunits arranged around a fivefold axis of pseudosymmetry, as predicted. The opening of the putative ion channel in the center is about 20 Å at its entrance, narrowing sharply at the membrane surface. At the current resolution, it is possible to begin to glimpse some elements of secondary structure. Thus in each subunit, in a region about 30 Å away from the extracellular surface, a group of three short rods can be detected, presumably  $\alpha$ -helices. The two subunits identified as  $\alpha$ -subunits (50) contain cavities shaped by these rods, which may correspond to the ACh-binding pockets. The narrowest portion, corresponding to the pore, appears to be shaped by five bent  $\alpha$ -helical rods, each contributed by one of the subunits (Fig. 1, left). The bends in the rods, near the middle of the membrane, are the parts closest to the axis of the pore and may represent the gate. In order to visualize the open state of the nAChR, Berriman and Unwin (88) devised a rapid-freezing device, designed to trap the open state, after application of ACh, by preventing the transition to the desensitized state. Comparison of the ACh-activated receptor with the nonactivated form revealed only slight changes (Fig. 1, right). Concerted twisting motions appeared to occur within the rods lining the putative ACh-binding pockets, which were consistent with the requirement that two ACh molecules, binding simultaneously to both  $\alpha$ -subunits, are required to open the channel. These localized disturbances were associated with small rotations extending along the axis of all the subunits, and with a switching of the helices within the pore from the bent shape seen in the nonactivated receptor to a tapered, more open configuration. If it is assumed that the threonine side chains of the central ring, which are believed to serve as the pore (see text above, along with Ref. 83) face into the pore, a minimum diameter of 10 Å can be estimated, comparable to the diameter obtained by use of various organic cations (89). Such a concerted rotation of the transmembrane sequences surrounding the lumen of the putative ion channel is consistent with the data accumulated by the various studies using site-directed mutagenesis and labeling, some of which were discussed above (see, in particular, Refs. 59 and 63).

**Figure 1.** Simplified diagram of the nAChR in its closed (left) and open (right) states, as suggested by the structural results. An ACh molecule first enters a binding site in one of the  $\alpha$ -subunits (rectangle), but significant displacements are blocked by the neighboring subunit, lying between the two  $\alpha$ -subunits. Interaction of a second ACh with the other

site then attempts to draw the neighboring subunit out of the way. A concerted localized displacement thereby takes place, initiating small rotations of the subunits along the shaft to the membrane. The rotations disrupt the gate by disrupting the association of  $\alpha$ -helices around the pore in the closed state (left), switching over to an alternate configuration in which a widened polar pathway for ion movement is created. (From Ref. [53](#), with permission.)



Obviously, when higher-resolution images become available, a more detailed mechanistic description will be feasible, but the conceptual stage appears to have been set.

### 1.6. Neuronal Nicotinic Receptors

Already in the late 1970s, it was proposed, on the basis of  $^{125}\text{I}$ - $\alpha\text{Bgt}$ -binding studies, that nAChRs were present not only in skeletal muscle, but also in the central nervous system ([90](#), [91](#)). These studies were not the subject of widespread attention, primarily due to the fact that functional correlates were lacking. When, in the mid-1980s, cloning techniques revealed the presence in brain of an nAChR [gene family](#) homologous to, but clearly distinct from, those of electric organ and muscle, and widely distributed in different brain areas ([92](#), [93](#)), it became obvious that nAChRs, like muscarinic receptors, must have important functions in brain. For many years, however, the lack of specific agonists and antagonists for a given neuronal nAChR subtype severely hampered the identification of functional nAChRs in various brain areas. These difficulties were further aggravated by the unusually fast kinetics of inactivation of some nAChR subtypes, compared to their peripheral counterpart (for literature, see Ref. [94](#)). Development of techniques permitting rapid application and removal of agonists was necessary to overcome this limitation (see, for example, Refs. [95](#) and [96](#)).

The cloning approach revealed the existence of a large number of AChR receptor subunits in the brain. In contrast to muscle and electric organ, however, these appear to fall into only two categories, a and b. By now nine a-subunits ([97](#)) and nine b-subunits have been described in vertebrate brain ([98](#)), and multiple nAChR subunits have also been described in invertebrates such as *Caenorhabditis elegans* ([99](#)) and *Drosophila* ([100](#)).

As discussed by Sargent ([98](#)), the various subtypes of subunits are expressed differentially in one brain region or another. The most recently discovered  $\alpha 9$ -subunit, for example, has a pattern of expression restricted to cochlear hair cells ([97](#)). The pharmacology of neuronal nicotinic receptors in relation to such topics as nicotine addiction ([101](#)) and Alzheimer's disease ([102](#)) are currently topics of intensive investigation; so too is their involvement in brain plasticity ([103](#)) and in behavior, learning, and memory ([104](#), [105](#)). A recent development is the observation that choline serves as a selective agonist for  $\alpha 7$  nAChRs in rat hippocampal neurons ([106](#)), in line with earlier reports for  $\alpha 7$  ectopically expressed in oocytes ([107](#), [108](#)). Because  $\alpha 7$  belongs to what appears to be the

evolutionarily oldest group of nAChRs (109), it is feasible that choline, rather than ACh, was the primeval transmitter for cholinergic receptors.

Neuronal nAChRs have not yet been purified or expressed in large amounts. Thus, essentially all our knowledge of their functional properties comes from studies in which subunits were expressed singly or in combination, most frequently by injection into *Xenopus* oocytes. Such expression studies were augmented, as for *Torpedo* and muscle nAChRs, by site-directed mutagenesis. These studies have been reviewed in detail by Sargent (98), and only a few points will be discussed here briefly.

Presence of an  $\alpha$ -subunit is necessary to obtain a functional receptor; and in the case of the  $\alpha 7$ -subunit, injection of it alone into oocytes is sufficient for it to form functional oligomers (110). Changeux, Bertrand, and co-workers took advantage of this experimental system to perform expression, combined with site-directed mutagenesis, which showed that the ligand-gated ion channel produced by the  $\alpha 7$ -subunit closely resembles that of the pentameric muscle and *Torpedo* receptors (see, for example, Ref. 111 and the review in Ref. 16). Direct experimental evidence for a pentameric stoichiometry was provided for the chick nAChR produced in oocytes by expression of the  $\alpha 4$  and  $\beta 2$  subunits with radioactive label incorporated into them. The stoichiometry of the subunits in the purified receptor was consistent with an  $\alpha_4\beta_2\beta_3$  pentamer (112). This does not mean, however, that expression of the mRNAs for any mixture of  $\alpha$ -subunits, with or without a  $\beta$ -subunit, will produce a functional oligomer. Thus far, only  $\alpha 7$  has been shown to produce functional oligomers containing only  $\alpha$  subunits. Furthermore, a recent study of Yu and Role (113) has shown that  $\alpha 5$  can participate in the function of an ACh-gated ion channel only if it is coexpressed with another  $\alpha$ - and a  $\beta$ -subunit. The structural basis for such restrictions with respect to assembly and/or function are, at this stage, unknown. Patrick and co-workers (114) observed, however, that a [Cyclophilin](#) is required for expression of functional homo-oligomeric, but not hetero-oligomeric, nAChRs. They raise the possibility that the  $\alpha$ -subunits in the homo-pentamer may not assume identical folds. Cyclophilins may thus play a critical role in the maturation of such homo-oligomers, acting directly or indirectly as prolyl *cis/trans* [Isomerases](#) or as **molecular chaperones** (115).

## Bibliography

1. P. Taylor (1990) In *The Pharmacological Basis of Therapeutics* (A. G. Gilman, T. W. Rall, A. S. Nies, and P. Taylor, eds.), 8th Edition, Pergamon Press, New York, pp. 122–130.
2. V. P. Whittaker (1992) *The Cholinergic Neuron and its Target: The Electromotor Innervation of the Electric Ray Torpedo as a Model*, Birkhäuser, Boston.
3. H. H. Dale (1914) *J. Pharmacol. Exp. Ther.* **6**, 146–190.
4. R. J. Lefkowitz, B. B. Hoffman, and P. Taylor (1990) In *The Pharmacological Basis of Therapeutics* (A. G. Gilman, T. W. Rall, A. S. Nies, and P. Taylor, eds.), 8th ed., Pergamon Press, New York, pp. 84–121.
5. D. A. Brown et al. (1997) *Life Sci.* **60**, 1137–1144.
6. S. Numa et al. (1983) *Cold Spring Harbor Symp. Quant. Biol.* **48**, 57–69.
7. T. Kubo et al. (1986) *Nature* **323**, 411–416.
8. J. Wess et al. (1997) *Life Sci.* **60**, 1007–1014.
9. A. Karlin and M. H. Akabas (1995) *Neuron* **15**, 1231–1244.
10. F. Hucho, V. I. Tsetlin, and J. Machold (1996) *Eur. J. Biochem.* **239**, 539–557.
11. D. Nachmansohn (1959) *Chemical and Molecular Basis of Nerve Activity*, Academic Press, New York.
12. B. Katz (1966) *Nerve, Muscle and Synapse*, McGraw-Hill, New York.
13. S. C. Froehner (1993) *Annu. Rev. Neurosci.* **16**, 347–368.
14. Z.-Z. Wang et al. (1996) *Cold Spring Harbor Symp. Quant. Biol.* **61**, 363–371.
15. J. R. Sanes (1997) *Curr. Opin. Neurobiol.* **7**, 93–100.
16. J.-P. Changeux (1995) *Biochem. Soc. Trans.* **23**, 195–205.

17. I. Silman and A. H. Futerman (1987) *Eur. J. Biochem.* **170**, 11–22.
18. J. Massoulié et al. (1993) *Prog. Neurobiol.* **41**, 31–91.
19. J. Patrick and J. Lindstrom (1973) *Science* **180**, 871–872.
20. A. Vincent et al. (1995) *J. Physiol. Paris* **89**, 129–136.
21. D. Nachmansohn (1955) *Harvey Lect.* **39**, 57–99.
22. A. Karlin and E. Bartels (1966) *Biochim. Biophys. Acta* **126**, 525–535.
23. A. Karlin and M. Winnik (1968) *Proc. Natl. Acad. Sci. USA* **60**, 668–674.
24. I. Silman and A. Karlin (1968) *Science* **164**, 1420–1421.
25. J.-P. Changeux, T. Podleski, and L. Wofsy (1967) *Proc. Natl. Acad. Sci. USA* **58**, 2063–2070.
26. R. D. O'Brien and L. P. Gilmour (1969) *Proc. Natl. Acad. Sci. USA* **63**, 496–503.
27. R. D. O'Brien, L. P. Gilmour, and M. E. Eldefrawi (1970) *Proc. Natl. Acad. Sci. USA* **65**, 438–445.
28. C. C. Chang and C. Y. Lee (1963) *Arch. Int. Pharmacodyn.* **144**, 241–257.
29. C. C. Chang and C. Y. Lee (1966) *Brit. J. Pharmacol.* **28**, 172–181.
30. C. Y. Lee (1972) *Annu. Rev. Pharmacol.* **12**, 265–286.
31. H. C. Fertuck and M. M. Salpeter (1976) *J. Cell Biol.* **69**, 144–158.
32. J.-P. Changeux, M. Kasai, and C. Y. Lee (1970) *Proc. Natl. Acad. Sci. USA* **67**, 1241–1247.
33. R. Miledi, P. Molinoff, and L. T. Potter (1971) *Nature* **229**, 554–557.
34. E. Karlsson, E. Heilbronn, and L. Widlund (1972) *FEBS Lett.* **28**, 107–111.
35. A. Karlin and D. A. Cowburn (1973) *Proc. Natl. Acad. Sci. USA* **70**, 3636–3640.
36. J. A. Reynolds and A. Karlin (1978) *Biochemistry* **17**, 2035–2038.
37. J. Cartaud et al. (1973) *FEBS Lett.* **33**, 109–113.
38. J. E. Heuser and S. R. Salpeter (1979) *J. Cell Biol.* **82**, 150–173.
39. R. R. Neubig, E. K. Krodel, N. D. Boyd, and J. B. Cohen (1979) *Proc. Natl. Acad. Sci. USA* **76**, 690–694.
40. M. Mishina et al. (1986) *Nature* **321**, 406–411.
41. M. Kasai and J. P. Changeux (1971) *J. Membr. Biol.* **6**, 1–80.
42. M. Weber, T. David-Pfeuty, and J. P. Changeux (1975) *Proc. Natl. Acad. Sci. USA* **72**, 3443–3447.
43. B. Katz and S. Thesleff (1957) *J. Physiol.* **138**, 63–80.
44. N. Nelson, R. Anholt, J. Lindstrom, and M. Montal (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3057–3061.
45. G. Boheim et al. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 3586–3590.
46. C. L. Weill, M. G. McNamee, and A. Karlin (1974) *Biochem. Biophys. Res. Commun.* **61**, 997–1003.
47. M. A. Rafferty, M. W. Hunkapiller, C. D. Strader, and L. E. Hood (1980) *Science* **208**, 1454–1456.
48. D. S. Wise, J. Wall, and A. Karlin (1981) *J. Biol. Chem.* **256**, 12624–12627.
49. A. Karlin et al. (1983) *J. Biol. Chem.* **258**, 6678–6681.
50. E. Kubalek, S. Ralston, J. Lindstrom, and N. Unwin (1987) *J. Cell Biol.* **105**, 9–18.
51. P. C. Kearney et al. (1996) *Neuron* **17**, 1221–1229.
52. H. A. Lester (1997) *Harvey Lect.*, **91**, 79–98.
53. N. Unwin (1998) *J. Struct. Biol.* **121**, 181–190.
54. J.-L. Popot and J.-P. Changeux (1984) *Physiol. Rev.* **64**, 1162–1239.
55. H. Shoji et al. (1992) *Eur. J. Biochem.* **207**, 631–641.
56. R. L. Haganir and P. Greengard (1987) *Trends Pharmacol. Sci.* **8**, 472–477.

57. J. Giraudat, C. Montecucco, R. Bisson, and J.-P. Changeux (1985) *Biochemistry* **24**, 3121–3127.
58. M. P. Blanton and J. B. Cohen (1994) *Biochemistry* **33**, 2859–2872.
59. M. P. Blanton et al. (1998) *J. Biol. Chem.* **273**, 8659–8668.
60. D. Naumann, C. Schultz, U. Görne-Tschelnokow, and F. Hucho (1993) *Biochemistry* **32**, 3162–3168.
61. U. Görne-Tschelnokow et al. (1994) *EMBO J.* **13**, 338–341.
62. M. O. Ortells and G. G. Lunt (1996) *Prot. Eng.* **9**, 51–59.
63. G. G. Wilson and A. Karlin (1998) *Neuron* **20**, 1269–1281.
64. P. N. Kao et al. (1984) *J. Biol. Chem.* **259**, 11662–11665.
65. D. Neumann, D. Barchan, M. Fridkin, and S. Fuchs (1986) *Proc. Natl. Acad. Sci. USA* **83**, 9250–9253.
66. D. Barchan, M. Ovadia, E. Kochva, and S. Fuchs (1995) *Biochemistry* **34**, 9172–9176.
67. J. Langenbuch-Cachat et al. (1988) *Biochemistry* **27**, 2337–2345.
68. J.-P. Changeux, J.-L. Galzi, A. Devillers-Thiéry, and D. Bertrand (1992) *Q. Rev. Biophys.* **25**, 395–432.
69. D. A. Dougherty and D. A. Stauffer (1990) *Science* **250**, 1558–1560.
70. M. Harel et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 9031–9035.
71. M. Harel et al. (1996) *J. Am. Chem. Soc.* **118**, 2340–2346.
72. A. Maelicke and E. Reich (1976) *Cold Spring Harbor Symp. Quant. Biol.* **40**, 231–235.
73. S. E. Pedersen and J. B. Cohen (1990) *Proc. Natl. Acad. Sci. USA* **87**, 2785–2789.
74. T. Heidmann, R. E. Oswald, and J. P. Changeux (1983) *Biochemistry* **22**, 3112–3127.
75. F. Hucho, W. Oberthür, and F. Lottspeich (1986) *FEBS Lett.* **205**, 137–142.
76. T. Heidmann, J. Bernhardt, E. Neumann, and J. P. Changeux (1983b) *Biochemistry* **22**, 5452–5459.
77. M. DiPaola, P. N. Kao, and A. Karlin (1990) *J. Biol. Chem.* **265**, 11017–11029.
78. E. A. Barnard, R. Miledi, and K. Sumikawa (1982) *Proc. R. Soc. Lond. B* **215**, 241–246.
79. K. Imoto et al. (1986) *Nature* **324**, 670–674.
80. K. Imoto et al. (1988) *Nature* **335**, 645–648.
81. M. H. Akabas, D. A. Stauffer, M. Xu, and A. Karlin (1992) *Science* **258**, 307–310.
82. K. Imoto et al. (1991) *FEBS Lett.* **289**, 193–200.
83. A. Villarroel, S. Herlitze, M. Koenen, and B. Sakmann (1991) *Proc. R. Soc. B* **243**, 69–74.
84. S. Hertling-Jaweed et al. (1988) *FEBS Lett.* **241**, 29–32.
85. J. Kistler and R. M. Stroud (1981) *Proc. Natl. Acad. Sci. USA* **78**, 3678–3682.
86. A. Brisson and P. N. T. Unwin (1984) *J. Cell Biol.* **99**, 1202–1211.
87. W. J. LaRochelle and S. C. Froehner (1986) *J. Biol. Chem.* **261**, 5270–5274.
88. J. A. Berriman and N. Unwin (1994) *Ultramicroscopy* **56**, 241–252.
89. B. N. Cohen, C. Labarca, N. Davidson, and H. A. Lester (1992) *J. Gen. Physiol.* **100**, 373–400.
90. Z. Vogel and M. Nirenberg (1976) *Proc. Natl. Acad. Sci. USA* **73**, 1806–1810.
91. J. Patrick and W. B. Stallcup (1977) *Proc. Natl. Acad. Sci. USA* **74**, 4689–4692.
92. D. Goldman et al. (1987) *Cell* **48**, 965–973.
93. S. Heinemann et al. (1990) *Prog. Brain Res.* **86**, 195–203.
94. E. X. Albuquerque et al. (1997) *J. Pharmacol. Exp. Ther.* **280**, 1117–1136.
95. N. G. Castro and E. X. Albuquerque (1993) *Neurosci. Lett.* **164**, 137–140.
96. Z.-W. Zhang, S. Vijayaraghavan, and D. K. Berg (1994) *Neuron* **12**, 167–177.

97. A. B. Elgoyhen et al. (1994) *Cell* **79**, 705–715.
98. P. B. Sargent (1993) *Annu. Rev. Neurosci.* **16**, 403–443.
99. J. T. Fleming et al. (1997) *J. Neurosci.* **17**, 5843–5857.
100. E. D. Gundelfinger (1992) *Trends Neurosci.* **15**, 206–211.
101. J. A. Dani and S. Heinemann (1996) *Neuron* **16**, 905–908.
102. H. Schröder et al. (1991) *Neurobiol. Aging* **12**, 259–262.
103. J.-P. Changeux et al. (1996) *Cold Spring Harbor Symp. Quant. Biol.* **61**, 343–362.
104. N. Le Novère, M. Zoli, and J. P. Changeux (1996) *Eur. J. Neurosci.* **8**, 2428–2439.
105. E. D. Levin et al. (1997) *Brain Res. Bull.* **43**, 295–304.
106. M. Alkondon et al. (1997) *Eur. J. Neurosci.* **9**, 2734–2742.
107. A. Mandelzys, P. De Koninck, and E. Cooper (1995) *J. Neurophysiol.* **74**, 1212–1221.
108. R. L. Papke, M. Bencherif, and P. Lippiello (1996) *Neurosci. Lett.* **213**, 201–204.
109. N. Le Novère and J. P. Changeux (1995) *J. Mol. Evol.* **40**, 155–172.
110. S. Couturier et al. (1990) *Neuron* **5**, 847–856.
111. D. Bertrand et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6971–6975.
112. R. Anand et al. (1991) *J. Biol. Chem.* **266**, 11192–11198.
113. C. R. Yu and L. W. Role (1998) *J. Physiol.* **509**, 667–681.
114. S. A. Helekar, D. Char, S. Neff, and J. Patrick (1994) *Neuron* **12**, 179–189.
115. S. A. Helekar and J. Patrick (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5432–5437.

### **Suggestions for Further Reading**

116. E. X. Albuquerque et al. (1997) Properties of neuronal acetylcholine receptors: pharmacological characterization and modulation of synaptic function. *J. Pharmacol Exp. Ther.* **280**, 1117–1136.
117. J.-P. Changeux (1995) The acetylcholine receptor: a model for allosteric membrane proteins. *Biochem. Soc. Trans.* **123**, 195–205.
118. F. Hucho, V. I. Tsetlin, and J. Machold (1996) The emerging three-dimensional structure of a receptor: the nicotinic acetylcholine receptor. *Eur. J. Biochem.* **239**, 539–557.
119. A. Karlin and M. H. Akabas (1995) Towards a structural basis for the function of nicotinic acetylcholine receptors and their cousins. *Neuron* **15**, 1231–1244.
120. P. B. Sargent (1993) The diversity of neuronal nicotinic acetylcholine receptors. *Ann. Rev. Neurosci.* **16**, 403–443.
121. N. Unwin (1998) The nicotinic acetylcholine receptor of the Torpedo electric ray. *J. Struct. Biol.* **121**, 181–190.
122. V. P. Whittaker (1992) "The cholinergic neuron and its target: the electromotor innervation of the electric ray *Torpedo* as a model", Birkhäuser, Boston.

### **Achaete–Scute Complex**

A central theme in developmental biology is to understand how the mechanisms that specify particular cell types integrate with those that govern where and when cell type-specific structures form during animal [development](#). Over the better part of the past century, the development of the *Drosophila* nervous system has served as a model with which to unravel the interconnected processes of cell fate specification and pattern formation. The **genes** of the *achaete–scute* complex (AS-C) bridge these two processes. The AS-C is composed of four related genes found within ~90 kbp of DNA in the distal tip of the [X-chromosome](#). The four genes—*achaete* (*ac*), *scute* (*sc*), *lethal of scute* (*l'sc*) and *asense* (*ase*)—all encode for basic **helix-loop-helix** transcriptional activator proteins and were initially identified by mutational analysis. Loss-of-function mutations in the AS-C remove sensory organs in the peripheral nervous system (PNS) and neural structures in the central nervous system (CNS), whereas gain-of-function mutations induce the formation of ectopic neural structures. The “proneural” *ac*, *sc*, and *l'sc* genes are expressed prior to neural precursor formation in precise patterns of cell clusters that forecast where neural precursors will form. During both embryonic and adult development, the primary patterning genes in *Drosophila* act through a large array of [cis-acting](#) regulatory regions found within the AS-C to create the stereotyped pattern of AS-C-expressing proneural cell clusters. Within each cluster (equivalence group), a cell communication process mediated by the [Notch signaling](#) pathway (lateral inhibition) restricts AS-C gene expression to a single cell. Within this cell, *ac*, *sc* and *l'sc*, either alone or in combination, are thought to activate a cassette of genes that directs this cell to acquire the neural precursor fate. All other cells within the cluster extinguish AS-C gene expression and are directed towards epidermal development. One target of the AS-C proneural genes in neural precursors is *ase*, which promotes the differentiation of neural structures. The AS-C thus links the process of cell fate specification to that of pattern formation: the AS-C genes receive and interpret global positional information through their regulatory regions and translate this information through their function into the formation (and differentiation) of a two-dimensional array of neural structures.

## 1. Genetic and Molecular Characterization of the *Achaete–Scute* Complex

The embryonic and adult *Drosophila* nervous systems are composed of a variety of different neural structures, each organized in reproducible but distinct patterns. For example, the neural precursors of the embryonic CNS form in orthogonal rows, whereas in the adult PNS large innervated bristles arise in an irregular, yet stereotyped pattern. Each bristle, or sensory organ, forms from a single neural precursor cell that occupies a fixed location in the developing fly. Cell migration is limited in the *Drosophila* PNS and CNS. Thus, the initial neural precursor pattern largely determines the later pattern of neural structures.

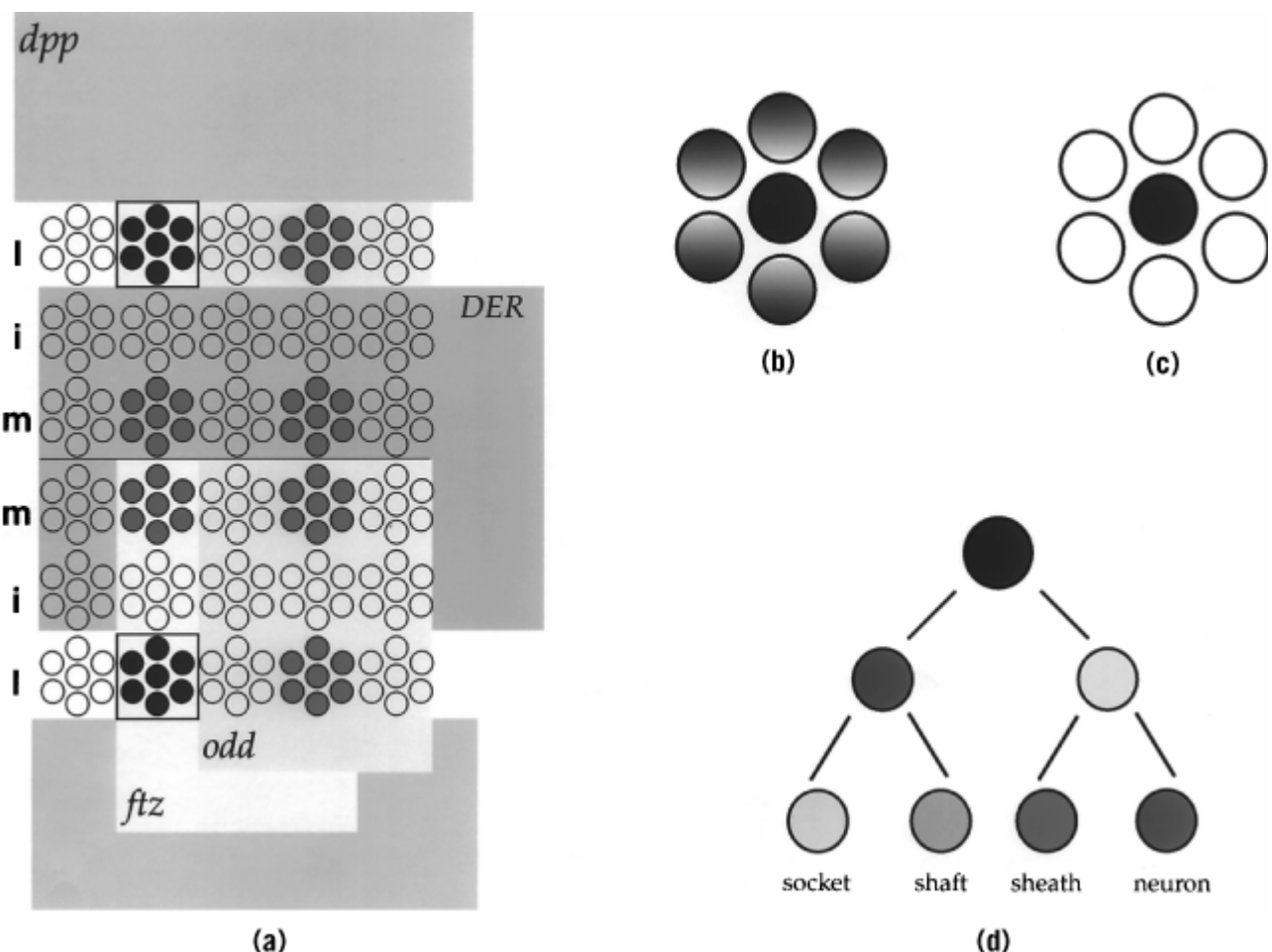
Genetic analyses identified a set of four genetic activities, designated the AS-C, located near the distal tip of the X-chromosome critical for the formation, patterning, and differentiation of neural structures (1). Loss-of-function mutations in any one of these genes remove neural structures, although the precise set of structures removed differs for each gene. For example, loss of *l'sc* activity primarily affects the CNS, whereas loss of *ac* or *sc* activity primarily results in PNS defects. Loss of *ase* function removes sensory organs along the wing blade but also causes differentiation defects in other sensory organs. Conversely, mis-expression of any AS-C gene promotes the formation of ectopic neural structures. These analyses suggested that the AS-C promotes the initial commitment of a cell to become a neural precursor and that their gene products are at least partially redundant. Detailed genetic mosaic analyses in the *Drosophila* wing [imaginal disc](#) (the larval structure that generates the wing and notum of the adult fly) showed that the activities of *ac* and *sc* define zones of neural competency from which neural precursors arise (2). These zones of neural competency are known as proneural cell clusters and are “equivalence groups” within which all cells can (although only one or a few will) become a neural precursor.

The genes *ac*, *sc*, and *l'sc* are each expressed in a complex and dynamic pattern of cell clusters, each of which quickly resolves to a single cell, the putative neural precursor (Fig. 1). The gene expression patterns of *ac* and *sc* are identical throughout neurogenesis and partially overlap with the expression



pattern of *l'sc*. On the other hand, *asense* is expressed exclusively in neural precursors and their progeny. Based on mutant phenotypes and expression patterns, *ac*, *sc*, and *l'sc* are called proneural genes whereas *ase* is called a neural precursor gene.

**Figure 1.** Patterning and specification of neural precursors. **(a)** Patterning of AS-C gene expression. Schematic representation of the ectoderm of a *Drosophila* embryo. The expression of the *ac* and *sc* genes in four proneural clusters per hemisegment is shown in red and blue. The activity of the early patterning genes *ftz*, *odd*, *DER*, and *dpp* subdivide the embryo and lead to the activation of *ac/sc* gene expression in reproducible quadrants. For example, the combined action of *ftz* activation and *odd* repression leads to *ac/sc* gene expression in domains that express *ftz* but not *odd* (yellow). Along the D/V axis, *dpp* (light blue) and *DER* (green) restrict *ac/sc* expression to the lateral column. Together, the action of these four genes delimit *ac/sc* gene expression within one proneural cluster, shown in blue. Note that other factors can overcome *DER* mediated repression in the medial column. **(b)** Singling out of a neural precursor. Within each proneural cluster, a single cell (in this example, the central cell) comes to express the AS-C proneural genes to the highest level, and this cell becomes the presumptive neural precursor (dark blue cell). **(c)** Resolution of cell fate within a proneural cluster. The presumptive neural precursor retains proneural gene expression and acts through the *Notch* signaling pathway to remove proneural gene expression in the other cells of the cluster (white). These cells are directed toward epidermal development. **(d)** Neural precursor lineage. Each neural precursor goes through an invariant cell lineage to produce a particular neural structure. The lineage of a sensory bristle of the PNS is shown. The neural precursor divides to yield two differently fated daughter cells, which in turn divide to produce a total of four cells. Each cell acquires a different fate, and together the four cells make up a functional sensory organ—an innervated bristle.



The molecular [cloning](#) of the AS-C in the 1980s identified all four genes as encoding putative transcriptional activator proteins that belong to the class B type of basic helix-loop-helix (bHLH) proteins. These four genes are contained within ~90 kbp of DNA, separated from one another by an average of ~30 kbp of DNA replete with *cis*-regulatory regions required to activate different

members of the AS-C in complex and partially overlapping patterns (3). Transcriptional activation requires the formation of a heterodimer between A and B class bHLH proteins (4); these heterodimers recognize the core DNA sequence CANNTG, where N is not specified (an 'E' box). An A-class bHLH protein with which AS-C gene products heterodimerize and activate transcription is the product of the *daughterless* (*da*) gene. Heterodimers of one AS-C gene product and Da mediate transcriptional activation of genes expressed throughout proneural clusters and in neural precursors. The *da* gene is expressed ubiquitously throughout *Drosophila* development, and AS-C genes are expressed in complex yet invariant patterns of cell clusters that predict where neural precursors form. Thus, the spatiotemporal specificity of where and when neural precursors develop is a function of the pattern of AS-C gene expression.

## 2. Cis- and Trans-Regulation of Proneural Gene Expression

Activation of the proneural genes of the AS-C in precise and reproducible patterns of proneural clusters is the first step in the formation of the stereotyped two-dimensional pattern of neural precursors in the *Drosophila* nervous system. The genetic regulatory mechanisms that dictate where and when AS-C expressing proneural clusters form are best understood for the development of the embryonic CNS (5). The anterior/posterior (A/P) and dorsal/ventral (D/V) register of AS-C-expressing cell clusters in the early embryo suggests that the segmentation genes, which segment the embryo along the A/P axis, and the dorsoventral genes, which specify pattern along the D/V axis, function to lay down directly the initial pattern of AS-C proneural clusters. A similar process mediated by *wingless* and *decapentaplegic* (*dpp*), as well as the genes of the *iroquois* complex (6), probably acts directly on the AS-C to create the pattern of AS-C-expressing cell clusters in adult structures.

During development of the embryonic *Drosophila* CNS, the AS-C proneural genes are expressed in an invariant orthogonal pattern of proneural cell clusters. For example, *ac* and *sc* are co-expressed in two rows of cell clusters in the medial and lateral—but not intermediate—columns of each segment. Prior to AS-C gene activation, the activities of different A/P and D/V patterning genes have subdivided the early embryo into an orthogonal pattern of squares, reminiscent of a checkerboard. Each square expresses a unique combination of these patterning genes, and the borders of the AS-C proneural clusters match precisely the limits of these squares (Fig. 1). This suggests a model whereby the combined activities of particular combinations of patterning genes within a square either does or does not activate a member of the AS-C. The composite pattern of AS-C-positive cell clusters would then result from the integration of the activities of each square.

Support for the checkerboard model comes from genetic analyses that assayed the expression of AS-C genes in embryos mutant for various known patterning genes. For example, the anterior border of every fourth transverse row of *ac/sc*-positive cell clusters coincides with the anterior edge of the expression of the *fushi-tarazu* (*ftz*) pair-rule gene, and the posterior border of these proneural clusters abuts the anterior border of the *odd-skipped* (*odd*) pair-rule gene (Fig. 1). In embryos mutant for *ftz*, this row of *ac/sc*-expressing proneural clusters disappears, whereas in embryos mutant for *odd* it expands posteriorly to fill the entire *ftz* domain. Thus, the combined action of *ftz* activation and *odd* repression defines the anterior and posterior boundaries respectively, of proneural clusters in this row. The D/V patterning genes control the mediolateral extent of AS-C proneural clusters. For example, the medial boundary of the lateral column of AS-C proneural clusters abuts the lateral border of the activity of the *Drosophila* epidermal growth factor (EGF) receptor (*DER*), and the lateral boundary of these clusters abuts the medial boundary of *dpp* activity (Fig. 1). In embryos mutant for *DER*, the lateral column of AS-C-expressing proneural clusters expands medially, whereas in *dpp* mutant embryos these clusters expand laterally. Thus, *DER* and *dpp* together restrict *ac/sc* gene expression to the lateral column in the developing CNS. The concerted action of *ftz* and *odd* along the A/P axis, and *DER* and *dpp* along the D/V axis, can account for the activation of *ac/sc* in one eighth of its proneural clusters. Combinations of other patterning genes probably act similarly to initiate AS-C proneural gene expression in the other regions of the developing CNS and PNS.

Scattered throughout the AS-C are regulatory regions that decode the positional information contained within the checkerboard pattern of spatial regulators and translate it into the transcriptional activation of *ac*, *sc*, and *l'sc*. The locations of many of these regulatory regions have been identified through the use of DNA rearrangements within the AS-C and of [reporter gene](#) constructs that contain specific genomic regions from the AS-C region (7). For example, *ac* and *sc* share a number of regulatory regions located between the two genes that function to activate *ac* and *sc*, but not *l'sc*, in identical patterns throughout development. It is thought that the gene products of many of the patterning genes act directly through these regulatory regions to create the precise and invariant pattern of AS-C-expressing proneural clusters.

### 3. Resolution of Proneural Cell Clusters: The Singling out of the Neural Precursor

Within each proneural cluster, AS-C expression quickly becomes restricted to one (or a few) cells (Fig. 1). This cell initiates *ase* expression, enlarges, and segregates as a neural precursor to a position below the proneural cluster from which it arose. All clusters exhibit identical dynamics of proneural gene expression, which directly reflect the cell-fate decisions made by the cells of a cluster. Initially, all cells express one or more of the proneural genes; proneural gene expression confers on all cells the ability to become a neural precursor. Then, via a cell-cell communication pathway mediated by the *Notch* signaling pathway, one cell comes to express the proneural genes to the highest level. This cell is the presumptive neural precursor; it then acts, again through the *Notch* pathway, to inhibit and eventually to extinguish proneural gene expression from all other cells of the cluster (8). The cells that lose proneural gene expression are directed toward the epidermal fate. Thus the fate of cells within a proneural cluster correlates—not with the initial expression of the proneural genes—but rather with the fate of proneural gene expression in that cell: cells that retain proneural gene expression become neural precursors, whereas cells that lose proneural gene expression are directed toward the epidermal fate.

The analysis of the AS-C provides a paradigm for how specific cell fates (and cellular structures) arise in reproducible and invariant patterns, but many questions still remain with respect to the AS-C. For example, roughly half of all CNS neural precursors still form in embryos devoid of AS-C function. Thus other “proneural genes” must exist that promote neural precursor formation in the CNS. In addition, AS-C function is not restricted to the nervous system. The gene *l'sc* mediates the initial selection of muscle progenitor cells in the mesoderm in much the same way that it (as well as *ac* and *sc*) single out neural precursors in the ectoderm. Furthermore, Galant et al (9) recently demonstrated that a butterfly AS-C **homologue** is expressed in, and thus may promote the formation of, the precursors to the pigment-producing scale cells that cover the butterfly wing. Thus, AS-C genes may mediate additional developmental events in *Drosophila* and other animals. Finally, in *Drosophila* the AS-C is an integral unit composed of four structurally related genes. Has the overall genomic structure of the AS-C and the relative position and expression of each AS-C gene been strongly conserved throughout evolution, as is observed for the genes of the **homeobox** clusters? Or does the AS-C and its constituent genes display significant plasticity over evolutionary time? Research on the AS-C has led to an understanding of some of the molecular mechanisms that specify particular cell types and those that determine where and when these cell types form. It is likely that further insights will derive from continued analysis of AS-C expression, regulation, and function.

### 4. Acknowledgments

I thank Michelle Hresko, Grace Panganiban, Lisa Nagy, and Ellen Ward for their comments. J.B.S is supported by a grant from NINDS (RO1-NS36570), the Cancer Research Fund of the Damon Runyon-Walter Winchell Foundation Award, DRS-9 and by a HHMI Research Resources Program for Medical Schools Junior Faculty Award #76296-538202.

### Bibliography

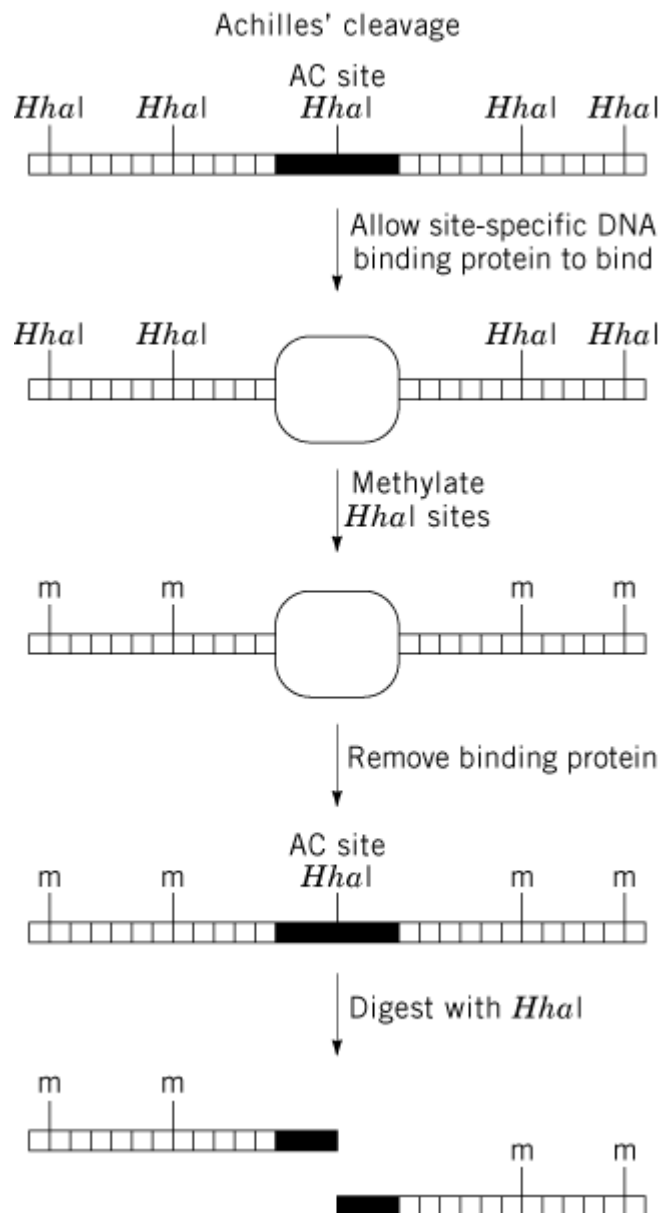
1. A. Garcia-Bellido (1979) *Genetics* **91**, 491–520.

2. C. Stern (1954) *American Scientist* **42**, 213–247.
3. S. Campuzano, L. Carramolino, C. V. Cabrera, M. Ruiz-Gomez, R. Villares, A. Boronat, and J. Modolell (1985) *Cell* **40**, 327–338.
4. C. Murre, P. Schonleber-McCaw, H. Vaessin, M. Caudy, L. Y. Jan, Y. N. Jan, C. V. Cabrera, J. N. Buskin, S. D. Hauschka, A. B. Lassar, H. Weintraub, and D. Baltimore (1989) *Cell* **58**, 537–544.
5. J. B. Skeath, G. Panganiban, J. Selegue, and S. B. Carroll (1992) *Genes Dev.* **6**, 2606–2619.
6. J. L. Gomez-Skarmeta, R. Diez del Corral, E. De La Calle-Mustienes, D. Ferres-Marco, and J. Modolell (1996) *Cell*, **85**, 95–105.
7. J. L. Gomez-Skarmeta, I. Rodriguez, C. Martinez, J. Culi, D. Ferres-Marco, D. Beamonte, and J. Modolell (1995) *Genes Dev.* **9**, 1869–1882.
8. J. A. Campos-Ortega (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez-Arias, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 1091–1130.
9. R. Galant, J. B. Skeath, S. Paddock, D. L. Lewis, and S. B. Carroll (1998) *Current Biology* **8**, 807–813.

## Achilles' Cleavage

Achilles' cleavage (AC) is a procedure developed to permit cutting double-stranded **DNA** at a single site (the AC site) in a complex [genome](#) (Fig. 1). The AC site is first protected from **methylation** by a sequence-specific [DNA-binding protein](#), and the methylated genome is subsequently cut by a methylation-sensitive [restriction enzyme](#). Other sites in the genome containing the same restriction/methylation site as the AC site are methylated and thus are not cut. The two key components of any AC system are (i) a restriction site that is recognized by both a methylation-sensitive restriction enzyme and its corresponding [methyltransferase](#) and (ii) the ability to block methylation of the chosen AC site by sequence-specific binding of a protein to that site. The overlap of the AC site and the recognition sequence for the DNA binding protein can occur naturally, or it can be engineered using [recombination](#) or [site-directed mutagenesis](#). These requirements somewhat limit the applications of the AC technique.

**Figure 1.** Achilles' Cleavage. Diagram of the method used to prevent restriction enzyme digestion at all sites except for the Achilles' cleavage (AC) site. Horizontal lines connected by short vertical lines represent double-stranded DNA. *HhaI* restriction sites are marked by bold vertical lines above the double-stranded DNA. The binding site for the DNA-binding protein is represented by a solid rectangle. The DNA-binding protein is represented by an empty oval. Sites of methylation are indicated by a vertical line and the letter 'm.'



Koob et al. (1) originally developed this procedure using two different **repressor/operator** systems. In one instance, **plasmids** carrying the operator of the [lac operon](#) were incubated with the [Lac repressor](#) and then methylated by the *Hha* I methyltransferase. Repressor binding prevented methylation of the *Hha* I site located within the *lac* operator. The methyltransferase and the operator protein were subsequently removed, leaving only the unmethylated AC site free to be cut by *Hha* I (1). Variations of the original AC method were also used to demonstrate the ability to make only one or a few cuts in the genomes of **yeast** and *Escherichia coli* (2), to turn frequent-cutting restriction enzymes into rare cutters (3), and to analyze the strength of protein-DNA interactions (4).

## 1. RecA-AC and RARE

The original AC method was made more universally applicable by using **RecA** protein and a sequence-specific **oligonucleotide** to form the protective complex, obviating the need for a sequence-specific binding protein. Two groups simultaneously developed this procedure, and their protocols were designated RecA-AC (5) and RARE (6). The RecA-AC/RARE technique still requires selection of a cleavage site that is recognized by both a restriction enzyme and its cognate methyltransferase. An oligonucleotide of 40 to 60 nucleotides in length is designed that matches the

DNA sequence containing the cleavage site. In the presence of  $Mg^{2+}$  and a non-hydrolyzable ATP analogue [eg, ATP(g-S)], RecA protein will bind to the oligonucleotide, forming a complex. When double-stranded DNA is added, the complex will bind to the matching genomic sequence to form a so-called 'triple-stranded' synaptic complex. This complex protects the cleavage site from methylation. As in the AC procedure, the methyltransferase and the RecA protein are removed after methylation, before the methylation-protected site is cleaved with the appropriate restriction enzyme.

As an alternative, the RecA complex can be used to protect specific restriction sites from restriction enzyme digestion, rather than from methyltransferase. This technique can be useful in some [cloning](#) procedures.

The RecA-AC/RARE technique has been used in various genome mapping experiments, for example to map gaps in human contiguous (**contig**) DNA and to map [telomeres](#) (7). The ability to cut genomic DNA at a few specific sites allows mapping large pieces of DNA by using [pulsed-field gel electrophoresis](#) and **Southern blotting** techniques to size DNA between two cleavage sites, or between a marker and the end of a telomere (7).

The RecA-AC/RARE technique is a useful technique for mapping and manipulating DNA and should continue to prove helpful in large-scale genome analysis studies of the future. One limitation to the technique is the number of sequences recognized by both a restriction enzyme and a methyltransferase. As more restriction enzyme/methyltransferase pairs are isolated or developed, the technique will become even more useful.

#### Bibliography

1. M. Koob, E. Grimes, and W. Szybalski (1988) *Science* **241**, 1084–1086.
2. M. Koob and W. Szybalski (1990) *Science* **250**, 271–273.
3. J. Kur, M. Koob, A. Burkiewicz, and W. Szybalski (1992) *Gene* **110**, 1–7.
4. P. M. Skowron, R. Harasimowicz, and S. M. Rutkowska (1996) *Gene* **170**, 1–8.
5. M. Koob, A. Burkiewicz, J. Kur, and W. Szybalski (1992) *NAR* **20**, 5831–5856.
6. L. J. Ferrin and R. D. Camerini-Otero (1991) *Science* **254**, 1494–1497.
7. L. J. Ferrin and R. D. Camerini-Otero (1994) *Nature Genet.* **6**, 379–383.

#### Suggestions for Further Reading

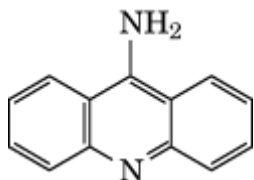
8. W. Szybalski (1997) RecA-mediated Achilles' heel cleavage. *Curr. Opin. Biotechnol.* **8**, 75–81.
9. L. J. Ferrin (1995) Manipulating and mapping DNA with RecA-assisted restriction endonuclease (RARE) cleavage, In *Genetic Engineering, Principles and Methods* (J. K. Setlow, ed.), Plenum Press, New York, Vol. **17**, pp. 21–30.

#### Acridine Dyes

Acridines are fused linear tricyclic aromatic molecules of planar geometry. Aminoacridine (Fig. 1) was originally developed as a topical antibacterial agent and is one of the most widely used and studied acridines. It is a strong base ( $pK_a$  9.9) due to the resonant [amino group](#) and is largely ionized at physiological pH. Acridines are typical DNA-**intercalating** ligands, binding tightly but reversibly to double-stranded DNA by insertion of the aromatic chromophore between adjacent base pairs, with

the long axis of the acridine parallel to the base-pair axis, ensuring maximum overlap in the binding site (1, 2). It was originally hypothesized by Brenner et al. (3) that acridine-induced mutations arise through the addition or deletion of a single base pair, an event described as **frameshift mutagenesis**.

**Figure 1.** Structure of the mutagen 9-aminoacridine.

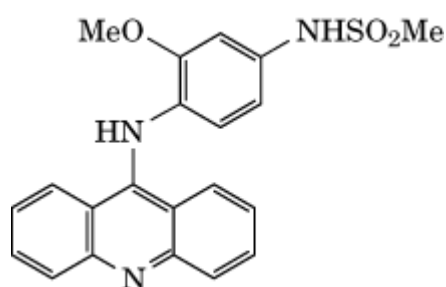


9-Aminoacridine is a strong mutagen in the rII region of **bacteriophage** T4 (4), in the CI gene of **lambda phage** (5), at various loci in *Escherichia coli* (6), and at the hisC3076 locus in *Salmonella typhimurium* (7). This and related simple acridines cause frameshift mutagenesis, especially in monotonous runs of repeated sequences. In mammalian cells their activity is somewhat weak and variable. They appear not to be causing frameshifts (at least in the most commonly used mammalian assays), but are mutagenic through weak ability to break chromosomes (8).

Substitutions on the acridine ring profoundly affect the mutagenic characteristics of these compounds, which have been divided into seven subgroups according to their distinctive mutagenic characteristics (2). 9-Aminoacridine and related simple compounds form the first group, while proflavine and other 3,6-diaminoacridines form the second group. The 3,6-diamino groups increase DNA binding affinity, and these compounds resemble 9-aminoacridine in having frameshift mutagenic activity, but are also substrates for metabolic activation (9). They are also readily activated by visible light into other products that are mutagenic in the *Salmonella* mutagenicity test (**Ames test**) (10). Some of these compounds and also the third class of quaternized acridines have considerable affinity for mitochondrial DNA. The fourth group of acridines retain DNA intercalation activity, but also act as topoisomerase II (topo II) poisons (see **DNA Topology**).

Type II topoisomerases produce double-strand breaks in DNA during replication and **transcription**, permitting strand passage through the transient break, which is then resealed by the enzyme. The potent antitumor activity of 9-anilinoacridines of the fourth group, such as amsacrine (Fig. 2), is thought to relate to their ability to stabilize the cleavable complex formed between topo II enzymes and DNA, thereby leading to an accumulation of double-strand breaks (11). These DNA breaks also result in these compounds showing strong clastogenic and recombinogenic properties. A clastogen is a physical, chemical or viral agent that produces chromosome breaks and other chromosomal mutagenesis. The fifth group of benzacridines are characterized by a higher susceptibility to metabolic activation, primarily through oxidative mechanisms (12). Benzacridines are primarily base-pair substitution mutagens, producing a spectrum of mutagenic events that is fundamentally different from that produced by other acridines and in which the acridine ring plays only a minor role. The sixth group of acridines, epitomized by those originally developed in the Institute of Cancer Research, Philadelphia, and widely known as ICR compounds (13), all carry an alkylating moiety, usually a nitrogen mustard. These compounds are not only potent frameshift mutagens in various organisms, but also cause base-pair substitution events in mammalian cells (2, 14). Although the nitroacridines can be distinguished as a seventh group, they may alternatively be considered as a subgroup of the previous class, because they can act as either simple acridines or as alkylators, depending upon the position of the nitro substitution and the nitroreductase capability of the organism concerned.

**Figure 2.** Structure of the 9-anilinoacridine amsacrine.



## Bibliography

1. L. R. Ferguson and W. A. Denny (1990) *Mutagenesis* **5**, 529–540.
2. L. R. Ferguson and W. A. Denny (1991) *Mutat. Res.* **258**, 123–160.
3. S. Brenner, L. Barnett, F. H. C. Crick, and A. Orgel (1961) *J. Mol. Biol.* **3**, 121–124.
4. A. Orgel and S. Brenner (1961) *J. Mol. Biol.* **3**, 762–768.
5. T. A. Skopek and F. Hutchinson (1984) *Mol. Gen. Genet.* **195**, 418–423.
6. S. M. Thomas and D. G. MacPhee (1985) *Mutat. Res.* **151**, 49–56.
7. L. R. Ferguson and D. G. MacPhee (1983) *Mutat. Res.* **116**, 289–296.
8. W. R. Wilson, N. M. Harris, and L. R. Ferguson (1994) *Cancer Res.* **44**, 4420–4431.
9. M. Fukenaga, Y. Mitzuguchi, and L. W. Yielding (1987) *Chem. Pharm. Bull. Japan* **35**, 792–797.
10. Y. Iwamoto, L. R. Ferguson, A. Pearson, and B. C. Baguley (1992) *Mutat. Res.* **268**, 35–41.
11. L. R. Ferguson and B. C. Baguley (1996) *Mutat. Res.* **355**, 91–101.
12. R. E. Lehr, A. W. Wood, W. Levin, A. H. Conney, and D. M. Jerina (1988) In *Polycyclic Aromatic Hydrocarbon Carcinogenesis: Structure–Activity Relationships* (S. K. Yang and B. D. Silverman, eds.), CRC Press, Boca Raton, FL, pp. 31–58.
13. H. J. Creech, R. K. Preston, R. M. Peck, and A. P. O'Connell (1972) *J. Med. Chem.* **15**, 739–746.
14. L. F. Stankowski, K. R. Tindall, and A. W. Hsie (1986) *Mutat. Res.* **160**, 133–147.

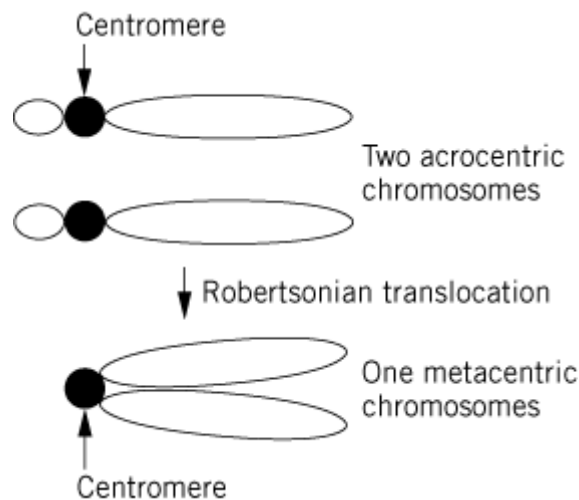
## Acrocentric Chromosome

An acrocentric chromosome has a single centromere that is localized at or near one end of the chromosome (acro = extremity). One of the most common forms of chromosomal translocation occurs when two acrocentric chromosomes fuse at the ends containing the centromeres (1). This is called a Robertsonian event (Fig. 1). This generally leads to the formation of a metacentric chromosome, where the new centromeric region appears at or near the middle of the new chromosome. Robertsonian events lead to a decrease in the apparent number of chromosomes. Much of the variation in the number of chromosomes in the members of a family or genus is related to this type of acrocentric chromosomal fusion. Whereas the number of chromosomal arms remain relatively constant, counting two for a metacentric chromosome and one for an acrocentric, the



number of chromosomes can vary widely. For example, most dogs have either 32 to 38 or 18 to 21 autosomes (any chromosome other than the X- or Y-chromosomes), but 34 to 38 chromosomal arms. Those with the greater number of autosomes have larger numbers of acrocentric chromosomes (2, 3). In genomes whose karyotypes contain exclusively acrocentric chromosomes (e.g., that of the house mouse, *Mus musculus*), rapid chromosome identification can be achieved via whole-chromosome painting techniques, e.g., spectral karyotyping (4).

**Figure 1.** Robertsonian translocation occurs when two acrocentric chromosomes fuse to create a single metacentric chromosome.



#### Bibliography

1. W. Robertson, *J. Morph.* **27**, 179–197 (1916).
2. R.K. Wayne, W.G. Nash, and S.J. O'Brian, *Cytogenet. Cell Genet.* **44**, 123–133 (1987).
3. R.K. Wayne, W.G. Nash, and S.J. O'Brian, *Cytogenet. Cell Genet.* **44**, 134–141 (1987).
4. T. Ried, E. Schrock, Y. Ning, and J. Wienberg, *Hum. Mol. Genet.* **7**, 1619–26 (1998).

#### Additional Reading

5. Wagner R.P., Maguire M.P., and Stallings R.L., *Chromosomes. A Synthesis*, Wiley-Liss, New York, 1993.

#### Acrosome

The acrosome is an organelle of the spermatozoon that covers the apical part of the [sperm](#) head. This organelle comprises an outer acrosomal [membrane](#), an inner acrosomal membrane, and plasma membranes. The acrosome plays an important role during the early stages of [fertilization](#), in that the acrosome and its associated molecules mediate recognition of the egg to be fertilized. Although acrosome morphology varies with the species, the function of the acrosome and its contents as regulating elements during the early stages of fertilization are fairly similar among species. The

acrosome contains **proteolytic** enzymes (see [Proteinases](#)) as well as proteins for binding to the zona pellucida and plasma membrane of the [Egg](#) (1). Therefore, the associated molecules of the acrosome have been proposed as targets for nonhormonal antifertilization agents.

## 1. Primary Ligands

The primary ligands are proteins that are located on or close to the acrosome; their function is to ensure that recognition between the gametes is species-specific. This is especially necessary in animals with extracorporal fertilization (eg, frog and sea urchin). On the surface of the acrosomal cap, there is a plethora of adhesion molecules (2), but their origin is not in all cases clear. Sperm adhesins become associated only upon ejaculation (3). Many of these primary ligands, such as galactosyltransferase (4), are [lectins](#) and bind to [O-linked oligosaccharides](#) and glycans of the ZP3 receptor of the oocyte's [extracellular matrix](#), the zona pellucida (5, 6). Other proteins were discovered that induce metabolic events; a 95-kDa protein was detected that induces **phosphorylation** events (7).

One of the most fascinating events during the early stages of fertilization is the acrosome reaction . After the initial contact of the sperm with the zona pellucida of the oocyte, and under the influence of progesterone (8), a unique event called the acrosome reaction takes place, which is calcium-dependent (9). Although in mammals the acrosome reaction can be induced without zona pellucida glycoproteins, they act in a catalytic manner (10). The acrosome reaction involves [exocytosis](#) of the outer membrane, causes the formation of hybrid vesicles, and leads to exposure of the inner contents of the acrosome. In many invertebrates, such as starfish, sea cucumbers, or sea urchin, a polymerization event takes place that results in the formation of [actin](#) filaments. In mollusks, the actin filaments are already present in the unreacted spermatozoon and exposed on activation.

Recent investigations report the presence of the **inositol phosphate** system in mammalian sperm. It involves the activation of a sperm **receptor** that stimulates a **G-protein**. This activates [phospholipase C](#), which cleaves subsequently phosphatidyl inositol diphosphate (PIP<sub>2</sub>) into **diacylglycerol** and inositol triphosphate (IP<sub>3</sub>). IP<sub>3</sub> triggers the release of calcium from intercellular stores, and diacylglycerol activates protein kinase [kinase C](#), which is calcium-dependent. This causes proteins to be phosphorylated and leads subsequently to the acrosome reaction (11). Upon their exposure, the molecules within the acrosome become important for the fertilization process.

## 2. Secondary Ligands: Acrosin/Proacrosin

A major component of the inner acrosome is the protein proacrosin . This is the best-characterized molecule of those involved in the early stages of fertilization. It is important for local lysis of the zona pellucida and therefore plays a central role at the stage of sperm penetration. Proacrosin is stored as its inactive form, a [zymogen](#). On raising the pH, or on contact with zona pellucida glycoproteins, the zymogen converts into its active proteolytic form, acrosin . This is accomplished by several intramolecular reorganizations (12). Proacrosin is a two-polypeptide chain molecule and, typical of [serine proteinases](#), the heavy chain starts at residue Val24. Major alterations of the molecule involve proteolytic cleavage of the proline-rich C-terminal tail, cleavage of the light chain, and linkage to the heavy chain via [disulfide bonds](#).

Proacrosin is a multifunctional protein; it does not act solely as a proteinase, but is also a fucose-binding protein (13). The molecule shows high binding affinity for fucoidan, heparin, and zona pellucida glycoprotein preparations (14). The primary mechanism of proacrosin binding to the zona pellucida is not only an interaction with carbohydrates, but involves positively charged amino acid residues of the proacrosin and negatively charged sulfate groups of the zona pellucida glycoproteins. **Homology modeling** of proacrosin revealed that the groove for binding to the zona pellucida contains loops in which the positively charged amino acids are located. The active site for proteolytic activity is located in the center, surrounded by the positively charged amino acids. The binding

groove and proteolytic center appear to form a single reactive unit, even though they are located on different regions on the [polypeptide chain](#). The proteolytic activity could be blocked without affecting proacrosin's binding characteristics.

The arrangement of the binding site and proteolytic center suggests two major tasks of the enzyme during fertilization. The molecule recognizes the zona pellucida specifically, binds to it, and digests it locally to promote sperm entry through the extracellular matrix. However, there must be further binding and proteolytic proteins involved in fertilization, as mouse **knockout strains** lacking the proacrosin gene did not produce infertile males ([15](#)); other molecules must be able to take over the role of proacrosin. On the other hand, male mice with a functionally active proacrosin gene demonstrated a significantly greater fitness in terms of reproductive success in a competition-mating experiment with the proacrosin knockout mice ([16](#)). Although proacrosin is the major compound of the acrosomal vesicle, there are further adhesion proteins. Fertilin , formerly designated PH-30, is a molecule that is involved in interactions between the sperm and egg ([17](#), [18](#)).

### Bibliography

1. A. P. Aguas and P. Pinto da Silva (1985) *J. Cell. Biol.* **100**, 528–534.
2. R. Jones (1990) *J. Reprod. Fert. Suppl.* **42**, 89–105.
3. J. J. Calvete, D. Solis, L. Sanz, T. Diaz-Maurino, W. Schäfer, K. Mann, and E. Töpfer-Petersen (1993) *Eur. J. Biochem.* **218**, 719–725.
4. X. Gong, D. H. Dubois, D. J. Miller, and B. D. Shur (1995) *Science* **269**, 1718–1721.
5. S. Kitazume-Kawaguchi, S. Inoue, Y. Inoue, and W. J. Lennarz. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3650–3655.
6. R. A. Kinloch, Y. Sakai, and P. M. Wassarman (1995) *Proc. Natl. Acad. Sci. USA* **92**, 263–267.
7. L. Leyton, P. LeGuen, D. Bunch, and P. M. Saling (1992) *Proc. Natl. Acad. Sci. USA* **89**, 11692–11695.
8. E. R. Roldan, T. Murase, and Q. X. Shi (1994) *Science* **266**, 1578–1581.
9. C. Arnoult, R. A. Cardullo, J. R. Lemos, and H. M. Florman (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13004–13009.
10. C. N. Tomes, C. R. McMaster, and P. M. Saling (1996) *Mol. Reprod. Dev.* **43**, 196–204.
11. R. Yanagimachi (1994) Mammalian fertilization. In *The Physiology of Reproduction*, 2nd ed. (E. Knobil and J. D. Neill, eds.), Raven Press, New York, pp. 189–317.
12. T. Baba, W. Kashiwabara, K. Watanabe, H. Itoh, Y. Michikawa, K. Kimura, M. Takada, A. Fukamizu, and Y. Arai (1989) *J. Biol. Chem.* **264**, 11920–11927.
13. E. Töpfer-Petersen and A. Henschen (1987) *FEBS Lett.* **226**, 68–42.
14. R. Jones and C. R. Brown (1987) *Exp. Cell Res.* **171**, 505–508.
15. T. Baba, S. Azuma, S. Kashiwabara, and Y. Toyoda (1994) *J. Biol. Chem.* **269**, 31845–31849.
16. I. M. Adham, K. Nayernia, and W. Engel (1997) *Mol. Reprod. Dev.* **46**, 370–376.
17. E. A. Almeida, A. P. Huovila, A. E. Sutherland, L. E. Stephens, P. G. Calarco, L. M. Shaw, A. M. Mercurio, A. Sonnenberg, P. Primakoff, and D. G. Myles (1995) *Cell* **81**, 1095–1104.
18. J. P. Evans, R. M. Schultz, and G. S. Kopf (1997) *Dev. Biol.* **187**, 94–106.

### Suggestions for Further Reading

19. B. T. Storey (1995) Interactions between gametes leading to fertilization: The sperm's eye view. *Reprod. Fert. Dev.* **7**, 927–942. Describes the interactions between the gametes leading to fertilization.
20. E. Töpfer-Petersen and J. J. Calvete (1995) Molecular mechanisms of the interaction between sperm and the zona pellucida in mammals: Studies on the pig. *Int. J. Androl. Suppl.* **2** **18**, 20–26. Summarizes established similarities and differences between sperm oocyte interactions of mammalian species.

## Actin

Actin is one of the most ubiquitous and conserved eukaryotic proteins. While actin was originally isolated and studied as part of the contractile apparatus of muscle (see [Thin Filament](#)), we now know that actin is found in virtually all eukaryotic cells and can be the most abundant [protein](#) present. Actin forms part of the [cytoskeleton](#) in nonmuscle cells and is involved in the maintenance of cell shape as well as cell dynamics and motility. The active form of actin is a helical polymer called F-actin (F for filamentous), assembled from monomeric subunits of G-actin (G for globular). Actin exists in a number of very similar isoforms (88% [amino acid](#) identity between yeast cytoplasmic actin and human muscle actin), and these isoforms display tissue, rather than species, specificity ([1](#), [2](#)). Thus, human cytoplasmic actin is much closer in sequence to yeast cytoplasmic actin than it is to human muscle actin. The actin monomer contains 374–376 residues, varying with the isoform, and is about 42k MW.

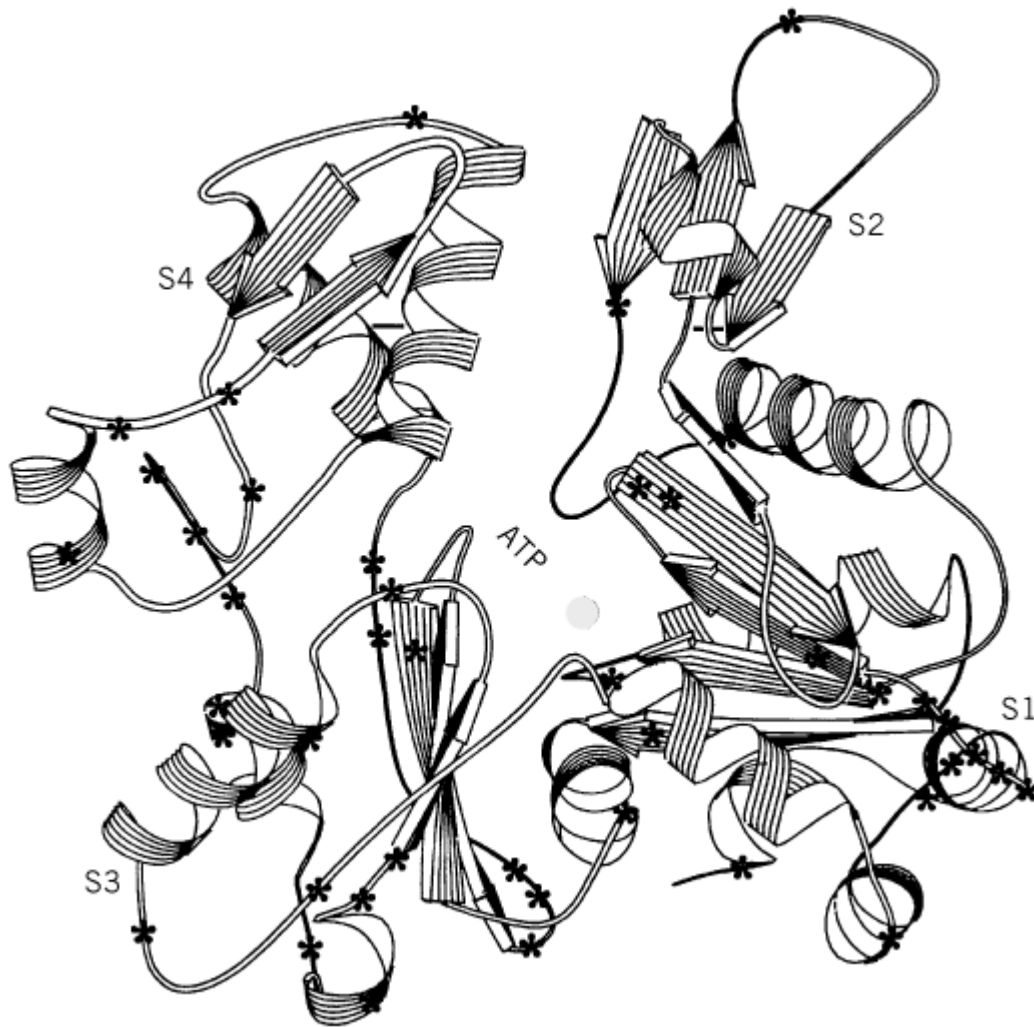
Actin interacts specifically with over 40 other proteins, and many of these interactions fall into separate classes. There are proteins that control the polymerization state of actin by nucleating, capping, or severing filaments [such as [gelsolin](#), [profilin](#), cofilin, and the b-thymosins ([3-7](#))], proteins that crosslink actin filaments together into gels and bundles [such as  $\alpha$ -**actinin**, villin, and fimbrin ([8-10](#))], and proteins that regulate the interactions of **myosin** with actin [such as **tropomyosin**, **troponin**, calponin, and caldesmon ([11](#)), (see [Thin Filament](#))]. In addition, many glycolytic enzymes bind actin, and it has been suggested that this provides spatial organization to metabolic processes via an association with the cytoskeleton ([12](#)). Other high-affinity interactions, such as between G-actin and the endonuclease DNase I (see [DNase I Sensitivity](#)) ([13](#)), do not lend themselves to simple explanations at this time.

Extensive biochemical, structural, and genetic studies have explored the properties of both G- and F-actin. The most detailed structural information exists for G-actin, since X-ray crystal structures have been obtained for complexes of G-actin with DNase I ([13](#)), gelsolin segment 1 ([14](#)), and profilin ([15](#)). The G-actin structure consists of two domains, each with two subdomains. These domains were originally called the “large” and “small” domains, but this terminology is not appropriate, as we now know that the “large” domain only contains 51% of the mass of the subunit. The N- and C-terminii are both located in the largest subdomain, subdomain-1. Between the two domains is a nucleotide-binding cleft, normally occupied by ATP in G-actin and ADP in F-actin, as well as the high-affinity metal-binding site. The subunit structure of G-actin is homologous to that of yeast hexokinase and the ATP-binding domain of HSC-70, a [chaperonin](#), suggesting a common evolutionary origin for all of these proteins ([16](#)). It has been proposed that the bacterial protein FtsA has the same fold ([17](#)), which would indicate that FtsA is a prokaryotic ancestor of the eukaryotic actins. It has already been established that the bacterial FtsZ protein has a common structure with eukaryotic [tubulins](#). Since neither hexokinase nor HSC-70 polymerize, it is evident that actin's ability to polymerize does not reside in any unique features of the secondary or tertiary structure, but rather evolved later. Thus, establishing that FtsA has a common subunit structure with actin will not, in and of itself, reveal any information about possible quaternary structures for FtsA.

Electron microscopic observations of F-actin (18-20) have been very useful in defining the orientation of the subunit in the filament, as well as revealing aspects of structural dynamics. X-ray fiber diffraction from oriented gels of F-actin has led to an atomic model for F-actin (21, 22), which provides a starting point for understanding the molecular details of the interactions of actin with many other proteins and ligands. However, many of the atomic details are still underdetermined by the method of fitting the monomeric structure into the filament, since the X-ray fiber diffraction pattern does not introduce all of the constraints needed. We thus have a general picture of the orientation of the subunit in the filament, but still lack information about many of the atomic interactions.

The subunit is oriented in the filament so that the two domains are nearly perpendicular to a plane containing the helix axis (Fig. 1). The resulting filament is about 100 Å in diameter, and subunits are spaced axially by a 27.3 Å rise along the filament axis. The change in the axial rise per subunit is less than .08 Å per subunit (<0.3%) when actin filaments in muscle are placed under full tension (23, 24), establishing that these filaments are relatively inextensible. Nevertheless, this amount of extensibility is enough to complicate the interpretation of muscle mechanics. Subunits are related to their axial neighbors by an  $\sim 167^\circ$  rotation about the filament axis, giving rise to a 59 Å pitch left-handed helix (sometimes called the genetic helix). An interesting property of the actin filament is the ability of subunits to rotate within the helix (25), and it has recently been shown that an actin-binding protein, cofilin, can change the twist of the actin filament by  $5^\circ$  per subunit, without significantly changing the axial rise per subunit (7). The helical geometry of actin also gives rise to two right-handed long-pitch strands containing actin subunits related by a 54.6 Å axial separation. These strands typically have a mean pitch of  $\sim 700\text{--}760$  Å (with the differences in the mean being due to the isoform, preparative conditions, associated proteins and ligands, etc.). The double-stranded character of these helices give rise to “crossovers” in projection at half of the pitch, or about 350–380 Å, but these crossover spacings can be variable due to the angular freedom within the actin filament (25, 26).

**Figure 1.** The G-actin subunit is shown by a ribbon representation (39) based upon the structure determined by X-ray crystallography (13). The degree of sequence conservation in actin is indicated by the fact that the only residues (45 out of 375) that differ between vertebrate striated skeletal muscle and yeast cytoplasmic actin are indicated by the asterisks. The four subdomains of actin are labeled S1–S4, and the most conserved part of the molecule is subdomain 2, with only 2 substitutions out of 38 residues between these two isoforms. The cleft between the two domains is occupied by ATP and a bound metal (indicated by the sphere). Adapted from Orlova et al. (40).



Every subunit in the actin filament is oriented with the same polarity, which gives rise to the overall polarity of the filament. This polarity has been described as the filament having a “barbed” end and a “pointed” end, based upon the morphological polarity that appears in electron micrographs when actin filaments are extensively decorated by the binding of proteolytic fragments of myosin, such as S1 or HMM (heavy meromyosin). This decoration generates a chevron appearance caused by all myosin molecules binding to the actin filament with an angle of about  $45^\circ$ . The morphological polarity has been related to a kinetic polarity apparent during *in vitro* polymerization, since the barbed end is the fast-growing end and the pointed end is the slow-growing end.

The helical symmetry of the actin filament is often described as 13/6 or 28/13, meaning that it has a helical repeat of 13 subunits in 6 turns of the  $59 \text{ \AA}$  helix, or 28 subunits in 13 turns. The helical repeat is the distance that a subunit must be translated axially to bring it into register with another subunit. An ideal 13/6 helix would have crossovers every  $355 \text{ \AA}$ , with each crossover containing a helical repeat, while an ideal 28/13 helix would have crossovers every  $382 \text{ \AA}$ , with two crossovers per helical repeat. However, there is no reason to believe that the helical symmetry needs to be described as a ratio of small integers, as there is an infinitesimally small change in structure from a symmetry of 13/6 to one of 1301/600 (a rotation of  $\sim 0.1^\circ$  per subunit), but the helical repeat changes from  $355 \text{ \AA}$  to  $35,517 \text{ \AA}$ ! The actual symmetry of the actin filament is entirely defined by the interactions between one actin subunit and four nearest neighbors ( $-2$  and  $2$  along the same long-pitch strand, and  $-1$  and  $+1$  on the opposite strand). In the absence of accessory proteins [such as **tropomyosin**, which binds to seven actin subunits along the same long-pitch strand (see [Thin Filament](#))], all changes of state must be propagated by such local interactions. Nevertheless, many *in*

*in vitro* experimental observations have shown that pure actin (in the absence of accessory proteins) can have long-range cooperative interactions within a filament (20, 27-29).

The flexibility of the actin filament has been studied by many methods, including spectroscopy (30), light microscopy (31), and electron microscopy (32, 33). Most studies have suggested a persistence, or correlation, length of about 6–7  $\mu\text{m}$ . This indicates that an actin filament much shorter than 6–7  $\mu\text{m}$  can be treated as a rigid rod, while a filament much longer than 6–7  $\mu\text{m}$  can be treated as a random coil. The persistence length does not suggest that a 6–7  $\mu\text{m}$  filament is approximated well by a rigid rod, since an actin filament only 0.6  $\mu\text{m}$  in length (one-tenth of the persistence length) will have a characteristic fluctuation of tangents at the two ends of about 25°.

Actin polymerization has been studied *in vitro* in great detail and is induced by raising the salt concentration. As for almost all other protein polymers in the cell, the assembly process involves the noncovalent binding of subunits to each other, such that every subunit (with the exception of subunits at the filament ends) is in the same environment. Probably more is known about the *in vitro* polymerization of actin than for any other protein polymer, but new findings suggest that the *in vivo* regulation of actin assembly involves the interactions with many other proteins and may be greatly different from what occurs *in vitro* (34). Polymerization can occur *in vitro* when the subunit concentration is above the critical concentration, which is the concentration of the monomer that would be at equilibrium with a population of polymers (see Treadmilling). Polymerization *in vitro* also requires that a monovalent salt, such as NaCl, be present at a concentration of ~100 mM, or a divalent salt such as MgCl<sub>2</sub> be present at ~1 mM.

Actin has a single high-affinity metal binding site, which has a greater affinity for Ca<sup>2+</sup> than for Mg<sup>2+</sup>. However, due to the much higher concentrations of Mg<sup>2+</sup> than Ca<sup>2+</sup> present *in vitro*, it is expected that this site will be occupied by Mg<sup>2+</sup>. This has been shown experimentally to be true in muscle (35), where the Mg<sup>2+</sup> concentration is three orders of magnitude higher than the Ca<sup>2+</sup> concentration. Actin also has 5–10 lower affinity metal binding sites (36). The spontaneous nucleation of actin filaments is often the rate-limiting step in actin polymerization *in vitro*, since the subsequent growth of existing filaments can occur rapidly. Within the cell, however, there is good reason to believe that actin polymerization occurs under the control of other proteins which serve to nucleate and cap growing filaments, and it is possible that spontaneous self-nucleation of actin filaments may never occur.

Actin polymerization is also loosely coupled to the hydrolysis of the bound nucleotide. While almost all attention has been focused on the primary nucleotide-binding site, actin also has a second nucleotide-binding site with mM affinity (37). Since this site is likely to be occupied in cells such as muscle where there exists a millimolar concentration of ATP, it remains to be seen whether this second site plays a physiological role. It is known that the nucleotide hydrolysis is not required for polymerization, since G-ADP actin can polymerize, as can G-actin subunits containing nonhydrolyzable analogs of ATP. The role of ATP hydrolysis is to actually destabilize the filament (38), allowing for depolymerization to occur more readily. ATP hydrolysis lags behind the initial polymerization event and occurs once a subunit is part of the polymer. The cell therefore uses the energy of ATP hydrolysis to allow for the dynamic assembly and disassembly of actin filaments, in contrast to less dynamic components of the cytoskeleton, such as intermediate filaments. The energy of ATP hydrolysis can also be used to drive treadmilling, where there is a unidirectional flux of subunits through a filament.

## Bibliography

1. P. A. Rubenstein (1990) *BioEssays* **12**, 309–315.
2. E. S. Hennessey, D. R. Drummond, and J. C. Sparrow (1993) *Biochem. J.* **291**, 657–671.
3. E. Ballweber, E. Hannappel, T. Huff, and H. G. Mannherz (1997) *Biochem. J.* **327**, 787–793.
4. M. F. Carlier and D. Pantaloni (1997) *J. Mol. Biol.* **269**, 459–467.

5. D. A. Schafer and J. A. Cooper (1995) *Ann. Rev. Cell Dev. Biol.* **11**, 497–518.
6. L. D. Burtnick, E. K. Koepf, J. Grimes, E. Y. Jones, D. I. Stuart, P. J. McLaughlin, and R. C. Robinson (1997) *Cell* **90**, 661–670.
7. A. McGough, B. Pope, W. Chiu, and A. Weeds (1997) *J. Cell Biol.* **138**, 771–781.
8. P. Matsudaira (1991) *Trends Biochem. Sci.* **16**, 87–92.
9. D. Hanein et al. (1998) *Nature Struct. Biol.* **5**, 787–792.
10. K. A. Taylor and D. W. Taylor (1994) *Biophys. J.* **67**, 1976–1983.
11. L. S. Tobacman (1996) *Ann. Rev. Physiol.* **58**, 447–481.
12. H. R. Knull and J. L. Walsh (1992) *Curr. Top. Cell. Reg.* **33**, 15–30.
13. W. Kabsch, H. G. Mannherz, D. Suck, E. F. Pai, and H. C. Holmes (1990) *Nature* **347**, 37–44.
14. P. J. McLaughlin, J. T. Gooch, H. G. Mannherz, and A. G. Weeds (1993) *Nature* **364**, 685–692.
15. C. E. Schutt, J. C. Myslik, M. D. Rozycki, N. C. W. Goonesekere, and U. Lindberg (1993) *Nature* **365**, 810–816.
16. K. M. Flaherty, D. B. McKay, W. Kabsch, and K. C. Holmes (1991) *Proc. Natl. Acad. Sci.* **88**, 5041–5045.
17. P. Bork, C. Sander, and A. Valencia (1992) *Proc. Natl. Acad. Sci.* **89**, 7290–7294.
18. R. A. Milligan, M. Whittaker, and D. Safer (1990) *Nature* **348**, 217–221.
19. A. Orlova and E. H. Egelman (1995) *J. Mol. Biol.* **245**, 582–597.
20. A. Orlova, E. Prochniewicz, and E. H. Egelman (1995) *J. Mol. Biol.* **245**, 598–607.
21. K. C. Holmes, D. Popp, W. Gebhard, and W. Kabsch (1990) *Nature* **347**, 44–49.
22. M. Lorenz, D. Popp, and K. C. Holmes (1993) *J. Mol. Biol.* **234**, 826–836.
23. K. Wakabayashi, Y. Sugimoto, H. Tanaka, Y. Ueno, Y. Takezawa, and Y. Amemiya (1994) *Biophys. J.* **67**, 2422–2435.
24. H. E. Huxley, A. Stewart, H. Sosa, and T. Irving (1994) *Biophys. J.* **67**, 2411–2421.
25. E. H. Egelman, N. Francis, and D. J. DeRosier (1982) *Nature* **298**, 131–135.
26. J. Hanson (1967) *Nature* **213**, 353–356.
27. F. Oosawa (1983) In *Muscle and Non-Muscle Motility* (A. Stracher, ed.) Academic Press, New York, p. 151.
28. G. Drewes and H. Faulstich (1993) *Eur. J. Biochem.* **212**, 247–253.
29. A. Muhrad, P. Cheung, B. Phan, C. Miller, and E. Reisler (1994) *J. Biol. Chem.* **269**, 11852–11858.
30. T. Yanagida and F. Oosawa (1978) *J. Mol. Biol.* **126**, 507–524.
31. H. Isambert, P. Venier, A. C. Maggs, A. Fattoum, R. Kassab, D. Pantaloni, and M. F. Carlier (1995) *J. Biol. Chem.* **270**, 11437–11444.
32. T. Takebayashi, Y. Morita, and F. Oosawa (1977) *Biochim. Biophys. Acta* **492**, 357–363.
33. A. Orlova and E. H. Egelman (1993) *J. Mol. Biol.* **232**, 334–341.
34. M. F. Carlier (1998) *Curr. Opin. Cell Biol.* **10**, 45–51.
35. T. Kitazawa, H. Shuman, and A. P. Somlyo (1982) *J. Musc. Res. Cell Motil.* **3**, 437–454.
36. M. F. Carlier, D. Pantaloni, and E. D. Korn (1986) *J. Biol. Chem.* **261**, 10778–10784.
37. P. Kiessling, B. Polzar, and H. G. Mannherz (1993) *Biol. Chem. Hoppe-Seyler* **374**, 183–192.
38. C. Combeau and M. F. Carlier (1988) *J. Biol. Chem.* **263**, 17429–17436.
39. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
40. A. Orlova, X. Chen, P. A. Rubenstein, and E. H. Egelman (1997) *J. Mol. Biol.* **271**, 235–243.

### Suggestions for Further Reading

41. P. Sheterline, J. Clayton, and J. Sparrow (1995) *Actin, Protein Profile* **2**, 1–103.



42. E. Reisler (1993) Actin molecular structure and function, *Curr. Opin. Cell Biol.* **5**, 41–47.
43. K. R. Ayscough and D. G. Drubin (1996) Actin: General principles from studies in yeast, *Ann. Rev. Cell Dev. Biol.* **12**, 129–160.

## Actin-Binding Proteins

### 1. The Actin Microfilament System

Actin was identified more than 50 years ago as an essential element of the force generating apparatus of muscle cells (1, 2). Together with myosin, tropomyosin, and a large number of other proteins, actin filaments constitute sarcomeres, which are ordered serially along contractile myofibrils (3, 4). In the mid 1960s, actin (5) and myosin (6, 7) were recognized as major structural components of the eukaryotic cytoplasm (8-10); for an early review of the field see (11). The unexpected observation that actin was the inhibitor of DNase I (12), and that the DNase inhibitor could be crystallized (13), led to the identification of the actin monomer-binding protein, [profilin](#) (14). The properties of profilin seemed to explain how unpolymerized actin could exist at high concentrations in cells without spontaneously polymerizing (15). It is now known that there are several other proteins controlling actin polymerization, and today the family of actin-binding proteins contains over 100 members including proteins that sequester actin monomers (profilins, thymosins), sever ([gelsolins](#)), depolymerize (cofilins), and cross-link (villin, fimbrin) actin filaments. Further characterization of the nonmuscle myosins (16), and an increased appreciation of the role of polymerization of actin suggests how cells can extend pseudopods and drag themselves along extracellular matrices (17, 18), and exhibit the high degree of motility seen on the surface of lymphocytes and in malignant cancer cells. This entry will describe general characteristics of actin-binding proteins, the analysis of which in many cases has reached the atomic level.

[Actin](#) is an obligatory component of all eukaryotes, and is one of the most highly conserved proteins during [evolution](#). Its sequence has varied less than 10% during the last 1200 million years (19). In mammalian cells, six different isoforms are known, of which the muscle specific  $\alpha$ -isoform and the nonmuscle  $\beta$ -isoform are the most studied [for review, see (20)]. On the basis of structural [homology](#) an actin-related protein, FtsA, which is involved in cell division, has been recognized in prokaryotes (21). Also, in eukaryotes, there are several actin-related proteins, ARPs (22).

Actin can be kept in a monomeric form only under low salt conditions or in the presence of actin **sequestering** proteins (review 20). Under physiological salt conditions, actin polymerizes into long filaments. In this form, actin participates in force generation and motile activity. The actin filament is polar, with a rapidly growing (+)-end (barbed end) and a slow growing (–)-end (pointed end). In the cell, actin filaments grow by addition of monomers at the (+)-end, which is located at the outer edge of cell surface protrusions: lamellipodia and filopodia. In the submembraneous zone all round the cell, actin filaments are organized into sheets and bundles held together along their lengths by a variety of [crosslinking](#) proteins, which differ in length and flexural rigidity.

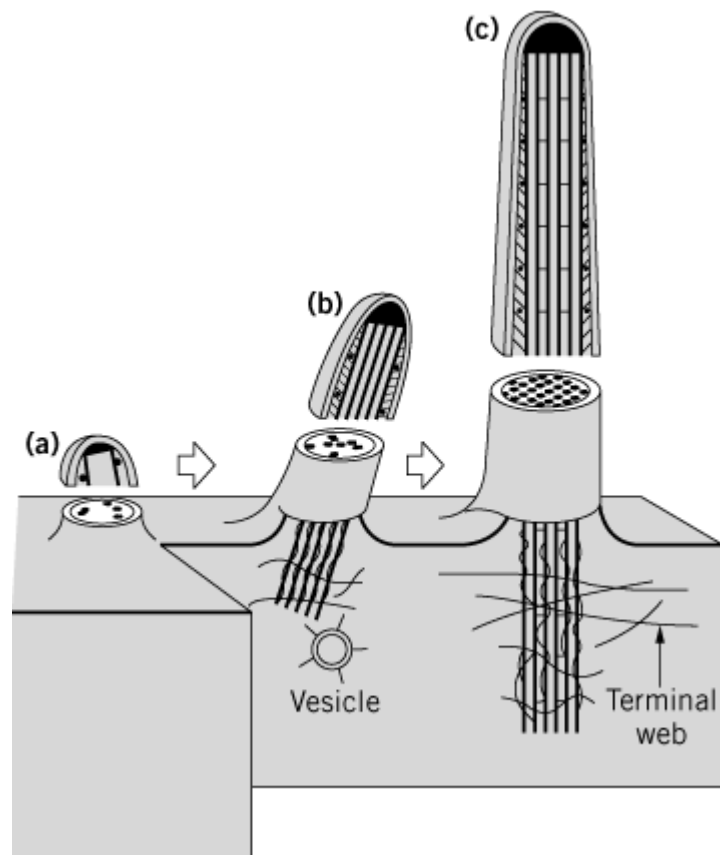
The myosins vary in structure; some are globular in shape, like myosin I, and do not polymerize, whereas others have a long  $\alpha$ -helical, supercoiled, rod domain and can form filamentous superstructures in the cytoplasm (16). Together, actin and myosin transduce chemical free energy of adenosine triphosphate (ATP) to mechanical work, not just in muscle cells, but in all eukaryotic cells. In different states of organization, the actomyosin system, is responsible for muscle contraction, cell locomotion, and cell division, and for many different local cellular transport processes (4). In nonmuscle cells, the actomyosin system is referred to as the [microfilament](#) system.

It is based on the organization of actin filaments, whose formation, organization and activity depend on the activity of a large number of actin-binding proteins, including myosin.

Microfilaments are present in particularly high concentrations in close apposition to the inner leaflet of the lipid bilayer (23-27), where they are directly or indirectly linked to cell surface proteins (receptors and adhesion proteins) continuously probing the extracellular environment. The lipid bilayer of the plasma membrane has been likened with “light machine oil” (28). It is an excellent electrochemical barrier, but it is the dense cortical weave of microfilaments which confers mechanical stability and deformability to the cell surface. Thus, the microfilament system can be viewed as an integral part of the physical barrier that separates the cellular contents from the outside world.

One of the first submembraneous actin-containing superstructures to be characterized was the red blood cell cortical network consisting of actin and the actin-binding proteins **tropomyosin**, **spectrin** and **band 4.1** (29). Linkage to the plasma membrane is provided by the spectrin-binding protein ankyrin, which binds to the transmembrane anion channel band III, and band 4.1 (an ERM protein, see below), which interacts directly with the transmembrane protein glycophorin. The adaptability of the shape of the red blood cell to the compressive forces present in the low-bore blood capillaries reflects the elastic properties of the spectrin:actin network. In the intestinal epithelium (30, 31, and refs. therein), microvilli on the apical surface of the cell are formed by tightly packed actin filament bundles stabilized by the actin crosslinking proteins villin and fimbrin (Fig. 1). The presence of nonmuscle myosin in these structures suggests that active movement may be involved in the resorptive process.

**Figure 1.** Structural organization of a microvillus. Among the simplest and best studied cell surface structures are the microvilli on polarized epithelial cells, especially those on cells of the intestinal epithelium. The majority of the proteins involved in the architecture of the microvillus have now been identified and characterized. Most of them are essential for the survival of the organism. In some cases their three-dimensional structure has been determined. It is still unclear how the assembly of the microvillus takes place and how it is controlled, how dynamic microvilli are, and what their contribution to uptake of nutrients is. The construction of a microvillus on apical surface of an intestinal epithelial cell is shown here. Many of the microvillar proteins also occur in other cell surface structures. An exception is villin, which is confined to the microvilli of a few cell types. In addition to villin, the major actin-binding proteins of the microvillar core are ezrin, fimbrin, and myosin-I. Drawing from Fath and Burgess 1995, *Curr. Biol.* **5**, 591–593.



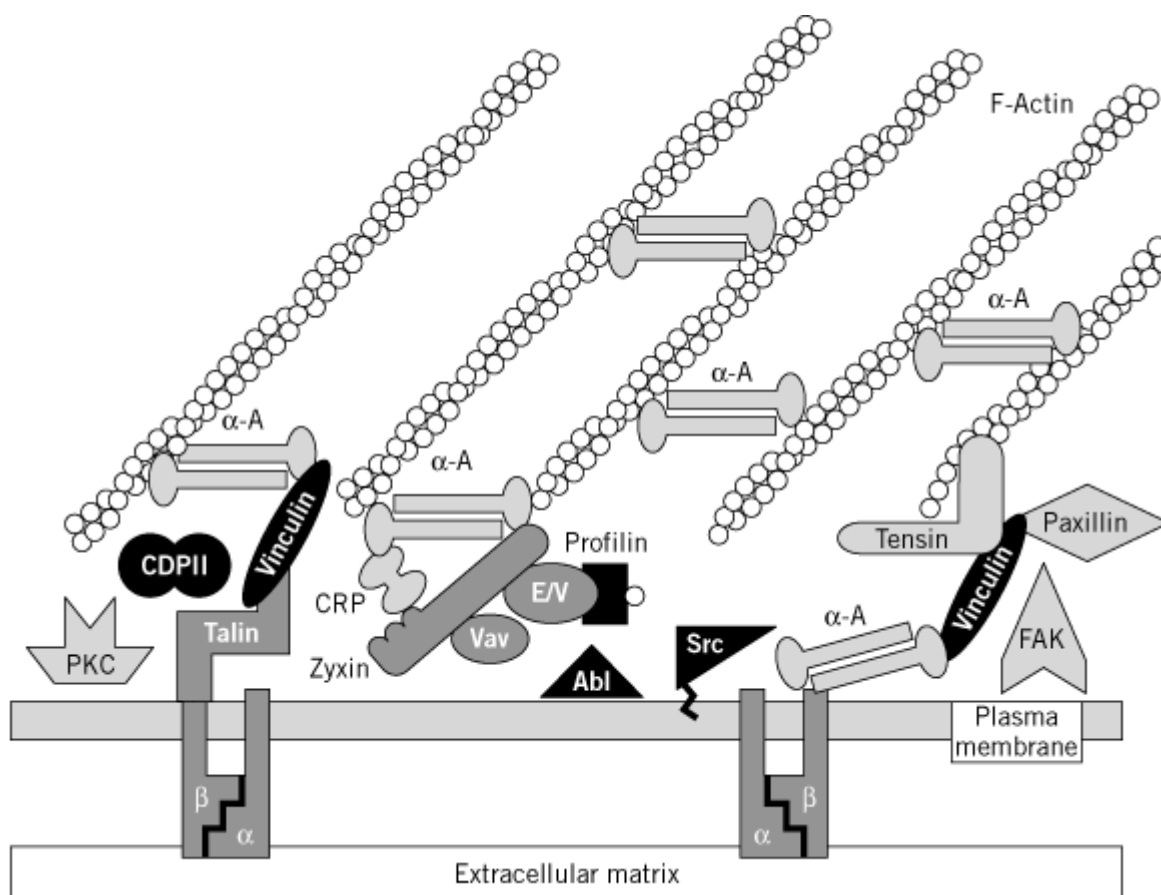
Specializations of the membrane-associated actomyosin system are found in all sensory organs, i.e. chemomechanical transduction in the actomyosin system is essential to the ability of the organism to perceive perturbations in the environment. For instance, the cellular basis for the sense of hearing lies in the microfilament-rich stereocilia in the hair cells of the inner ear, which convert subtle changes in the ambient air pressure to electrical signals transmitted to the brain via nerve impulses (32). Similarly, the mechanosensory cells in the skin confer the sense of touch (33).

Macrophages, lymphocytes, and neutrophils, which are highly motile cells, exploit the full potential of the cortical microfilament system as they forcibly move through and along extracellular matrices. When macrophages contact foreign invaders, the versatility of the cell cortex allows them to adapt their own surface structures as they engulf the pathogen.

The process of cell migration depends on the coordination of four major events involving microfilament reorganization: extension at the leading edge of membrane lamellae and filopodia, attachment to extracellular structures via transmembrane adhesion molecules, development of intracellular tension, and release of trailing end attachments with retraction of trailing ends (18, 27, 34, 35). The protrusion of lamellae at the advancing edge is driven by polymerization of actin filaments at focal complexes containing their fast growing ends at the plasma membrane, and crosslinking of formed filaments (24, 25, 36, 37). Clustering of transmembrane adhesion proteins, integrins, gives rise to focal adhesion sites (Fig. 2) linking extracellular matrix proteins to bundled actin filaments on the inside of the plasma membrane (38). The importance of the focal adhesions in transmembrane signalling is emphasized by their association with kinases regulating signal transducing protein:protein interactions and enzymes generating second messengers controlling not only the microfilament system, but also gene transcription (39). Several actin-binding proteins are targeted to these sites, contributing to formation of the integrin:actin filament linkage. The association of linkage proteins like talin, tensin, and vinculin appears to be controlled by phosphorylation. Myosin is activated through  $\text{Ca}^{2+}$ -dependent processes (40, 41), and the regulated

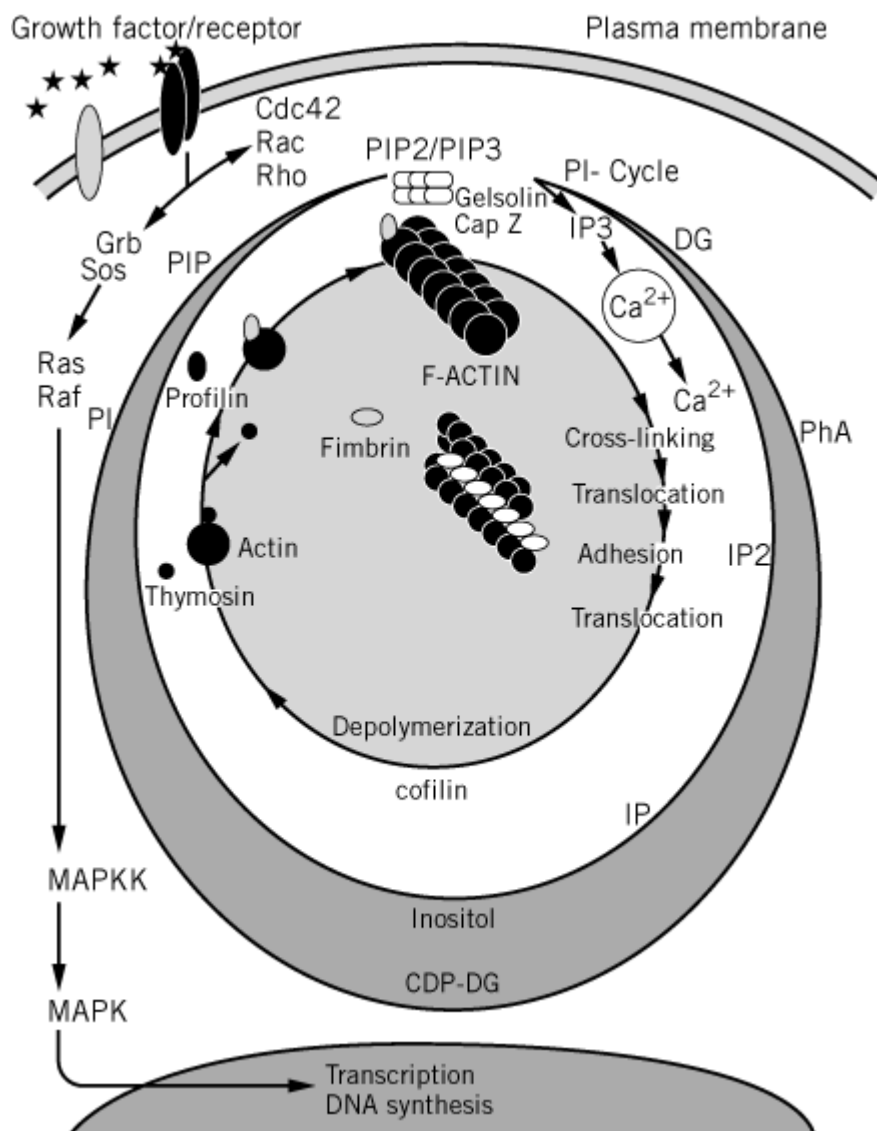
deployment of myosin molecules leads to the generation of the forces needed for cell motility and other translocation processes. Release of adhesion at the trailing end of a migrating cell is controlled by  $\text{Ca}^{2+}$  and is followed by [endocytosis](#) of the integrins and their recycling to advancing cell edges (42). Used filaments are depolymerized, a process which involves the actin monomer- and filament-binding protein ADF (or cofilin) (17), and sequestered forms of unpolymerized actin are reformed (see Fig. 3).

**Figure 2.** Proteins of focal adhesion sites. Adherens junctions comprise both cell:extracellular matrix contacts, exemplified by focal adherens in tissue cultured cells, and intercellular junctions such as the zonula adherens between epithelial cells. The regulated assembly of adherens junctions is essential to a diversity of cellular functions including wound healing, anchorage-dependent cell proliferation, adhesion, and migration. Assembly of an adherens junction may be initiated when an integrin of the b1 subfamily binds to its extracellular matrix ligand, when a cadherin recognizes an identical molecule on a neighboring cell, or when a growth factor activates its receptor. These ligand–receptor interactions result in establishment of a physical connection between the extracellular ligand and the actin cytoskeleton. Efforts to elucidate these linkages have implicated vinculin as one of the connecting molecules in a chain that involves binding of the b1 subunit of integrin to talin, talin to vinculin, vinculin to a-actinin, and a-actinin to filamentous actin. It is now known that vinculin also can bind directly to actin filaments. How the initial ligand–receptor interactions trigger recruitment and organization of the multiple proteins found at an adherens junction is unknown, but recent evidence suggest that the activities of tyrosine kinases, tyrosine phosphatases, protein kinases and members of the Rho family of small GTPases are involved. (Drawing from M. C. Beckerle (1997) *Bioassays* 18, 949–957)



**Figure 3.** Links between the phosphatidylinositol and the cell motility cycles. Binding of growth factors like platelet-derived growth factor and epidermal growth factor to their cognate receptors activate the kinase activities of the receptors resulting in phosphorylation of a number of tyrosines on the cytoplasmic part of the receptor. This leads to activation of the signal transduction elements Cde42, Rac and Rho which control outgrowth of filopodia, membrane lamellae and the formation of stress fibers, respectively (64). There is only little information so far regarding the

mechanism of activation of these small GTPases and how they relay the signal to the phosphatidylinositol (PI) cycle. Receptor activation is followed by an immediate increase in the activity of kinases that form PtdIns 4,5-bisphosphate (PIP2) and PtdIns 3,4,5-trisphosphate (PIP3). Profilin, gelsolin, CapZ, a-actinin, vinculin, the ERM family of proteins all bind these phosphoinositides, which influence their activity in the microfilament-based cell motility cycle (CM-cycle). Phosphoinositides binding to CapZ or gelsolin unblock the (+)-end of actin filaments allowing incorporation of profilin:actin at that site. Profilin:actin is then dissociated and the actin monomer becomes stably incorporated into the growing filament. Cross-linking of filaments by proteins like fimbrin gives rise to the ensembles of filaments seen in membrane lamellae and filopodia which function together with myosin in the contractile phase of cell motility. Calcium ions play important roles in regulating both the actomyosin interaction and the depolymerization of actin filaments. The depolymerization is facilitated by the dephosphorylated form of cofilin and the actin monomer sequestering proteins thymosin and profilin (rev. 163). Another branch of signal transduction from the growth factor receptor activates the gene program that leads to DNA synthesis and cell division. Grb and Sos are small proteins acting as signal transduction elements activating the small GTPases Ras, which in turn activates the serine/threonine kinase Raf. The further signalling pathway is complex involving mitogen activated protein kinases, MAPKK and MAPK (for review, see *Cell* 1998, **95**, 447–450). Other abbreviations: IP3, inositol trisphosphate; IP2, inositol bisphosphate; IP, inositol monophosphate; DG, diacylglycerol; Pha, phosphatidic acid; CDP-DG, cytidyldiphosphate-diacylglycerol; PI, phosphatidylinositol.



Thus, the combined action of protein factors controlling polymerization, crosslinking and tension development in the cell cortex is initiated by transmembrane signalling. Controlling the timing and position of sites of polymerization at the cell surface is a potent means of regulating shape change. An astonishing variety of well-ordered, yet dynamic, structures are possible, reflecting the

mechanical properties and specificities of the different kinds of actin bundling proteins and junctional connectors. Their timed and coordinated formation prevents the appearance of random rigid networks, which would be incompatible with the well-orchestrated movements seen.

## 2. Signal Transduction and Restructuring of the Actin Microfilament System

Cell surface receptors continuously monitor the extracellular milieu and pass signals to the cytoplasm by generating second messengers of various kinds which either stimulate or inhibit the activity in the microfilament system. Addition of [growth factors](#) to serum-starved cultured cells causes the immediate outgrowth of membrane lamellae and filopodia, events that depend on the polymerization and organization of actin. About 60 seconds later, there is a generalized increase in motile activity of the cells, including restructuring of the actin-containing stress fiber system and translocation of the cell. Observations of this kind have been made both with epidermal growth factor (EGF) ([43](#)) and platelet-derived growth factor ([26](#)). In the case of EGF, a significant fraction of the receptors was found in direct association with microfilaments, and a binding site for actin on the receptor has been identified ([44](#), [45](#)).

The phosphatidylinositol-cycle (PI cycle) is coupled to activation of the microfilament system ([46](#)). Ligand-induced, receptor-mediated activation of kinases generate phosphatidyl inositol 4,5-bisphosphate (PtdIns 4,5-P<sub>2</sub>) and phosphatidyl inositol 3,4,5-trisphosphate (PtdIns 3,4,5-P<sub>3</sub>) ([47](#), [48](#)), both of which affect the activity of actin-binding proteins. [Profilin](#), [gelsolin](#), [a-actinin](#), [vinculin](#), [ERM](#) family of proteins, and [myosin](#) all bind PtdIns 4,5-P<sub>2</sub>, and/or PtdIns 4,5-P<sub>3</sub>. Stimulation of platelets with low concentrations of thrombin causes a transient increase in polyphosphoinositides during the first 10 sec ([49](#)). This is accompanied by an equally transient transformation of unpolymerized actin into filaments. This initial polymerization is over in about 10–20 sec ([49](#)). The immediate precursor for this polymerization appears to be profilin:actin, since increased amounts of free profilin appears in cell extracts after stimulation. Within 2 min, these filaments become organized into supramolecular structures by the action of crosslinking proteins stabilizing the protrusions seen on the platelet surface ([50–52](#)). Simultaneous activation of phospholipase C $\alpha$ -1 results in the hydrolysis of PtdIns 4,5-P<sub>2</sub>, releasing the second messenger inositol trisphosphate, which causes release of Ca<sup>2+</sup> ions from intracellular Ca<sup>2+</sup> stores ([53](#)). Calcium ions are involved in activation of the actomyosin interaction and thereby force generation.

In addition to connections to the phosphatidyl inositol signal transduction cycle, the microfilament system is linked to the cAMP-dependent signalling system. Agonists causing increasing levels of cAMP often inhibit actin polymerization ([54](#)). A connection between adenylate kinase pathways and the actin monomer-binding protein profilin was suggested by genetic studies showing that increased levels of profilin could compensate for the deletion of the C-terminal domain CAP of the cyclase ([55](#), [56](#)). In platelets, the actin-binding protein, vasodilator-stimulated protein VASP is phosphorylated by cAMP- and cGMP-dependent kinases ([57](#)), which then correlates with the inhibition of microfilament reorganization.

The proteins involved in controlling the organization and function of the microfilament system are part of multiprotein signal transduction complexes. Several of these proteins have been identified as products of protooncogenes, ie proteins which in mutated form may cause cancer ([58](#), [59](#)). In these signal transduction complexes the proteins are linked together by adaptor modules in linker proteins such as SH2, SH3, and WW domains which recognize specific features of short polypeptide stretches, especially phosphotyrosine-containing and polyproline sequences ([60](#)). An example is the oncogene product **Vav** that has several such binding motifs suggesting how it links signal transduction at the cell surface to restructure the cortical cytoplasm ([61](#)). Vav possesses many of the homology domains commonly found in signal transduction complexes, comprising a phosphotyrosine-binding SH2 domain, two proline-rich peptide recognizing SH3 domains, a phosphatidyl inositol 4,5-bisphosphate recognizing PH domain, a guanine nucleotide exchange factor DH-domain, two cysteine-rich nuclear transcriptionally important LIM domains, and an actin-

binding calponin homology domain. The clearest demonstration that these functional linkages actually lead to nuclear transcription is that the production of Interleukin-2 (IL2) requires Vav-dependent actin cap formation in activated T-cells (62, 63).

Outgrowth of filopodia and membrane lamellae, and the formation of stress fibers, which are manifestations of the activation of the microfilament system by different kinds of ligand:receptor interactions, all depend on the activation of members of the Rho subfamily of small GTPases for their control. The small GTPase Cdc42 appears to control formation of filopodia, Rac controls membrane lamellae and Rho stress fibers (64). There is only little information so far as to how the receptor-mediated transmembrane signal is relayed to the GTPases and what their downstream targets are. It has been shown with permeabilized platelets that ligation of the thrombin receptor, and activated Rac, uncap actin filament (+)-ends through phosphoinositide synthesis (65). Some of the microfilament-controlling, signal transduction elements mentioned below interact with small GTPases. Interestingly, moesin (see below), an actin-binding protein linking actin filaments to transmembrane receptors, is required for Rho and Rac effects on the actin filament system (66).

VASP belongs to a growing collection of related proteins that are involved in stimulation of assembly of actin filaments *in vivo* (67-69). In moving cells, it is found in highly dynamic membrane regions and in integrin-rich adhesion plaques. Molecular cloning of VASP has established that it has both a central proline-rich core, which interacts directly with profilin (68, 70), and an actin-binding site. In cells, the primary site of growth of actin filaments is thought to be at advancing cell edges, and profilin:actin heterodimers are believed to be guided into this site by VASP (71). High concentrations of VASP appear to be targeted to focal adhesion sites by zyxin (72). Zyxin also binds a-actinin and Vav suggesting that orderly assembly of crosslinked actin networks can proceed at these locales. In addition to binding profilin and zyxin, VASP binds to the SH3 domains of the src-family of protein kinases (69), thus integrating signal transduction via phosphorylation with the ordered assembly of cortical structures in locally specialized sites (63). VASP and vinculin complexes can be immunoprecipitated out of cell extracts with either antivinculin or antiVASP antibodies, and recent experiments have demonstrated that the interaction between the two proteins is enhanced by PtdIns 4,5-bisphosphate binding to vinculin (73).

VASP has come to prominence owing to its relevance to the understanding of the movement of the bacterium *Listeria* through the host cell cytoplasm (74). *Listeria* expresses a proline-rich protein on its surface called ActA (zyxin-like), which binds to the modular proline-binding SH3-domain of VASP. Apparently, VASP can attract profilin:actin complexes to sites of polymerization on the surface of the bacterium, generating a propulsive force. There are many examples of pathogenic microorganisms that have acquired host cell genes that code for proteins that modulate the behavior of the microfilament system for the purpose of invasion, locomotion in the cytoplasm, and cell-to-cell spreading.

WASP is a multidomain, actin-binding protein, which is defective in patients afflicted with Wiskott-Aldrich syndrome (75, 76), a disease characterized by severe thrombocytopenia and immunodeficiency. WASP is specific to hematopoietic cells, which in case of WAS have a paucity of surface filopodia and other shape abnormalities (77). The small GTPases Cdc42 and Rac bind directly to WASP (78-80). WASP is also a binding partner for Src family of protein tyrosine kinases (81).

A protein homologous to WASP, N-WASP, is more widely distributed. N-WASP is concentrated at nerve terminal regions in the post synaptic density, and coexpression of Cdc42 and N-WASP induces extremely long actin-based filopodia in cells (82, 83). Its involvement in the formation of membrane lamellae and filopodia may be controlled by phosphoinositides in the plasma membrane. It contains a pleckstrin homology PH domain and a cofilin homology domain. Furthermore, N-WASP has been shown to be essential for the actin polymerization-dependent movement of the bacterium *Shigella flexneri* (84).

WIP is a recently discovered WASP interacting protein that is important for cortical actin assembly (85). WIP contains both profilin- and actin-binding regions, which like in VASP, are near each other in the protein sequence suggesting that they constitute the binding site for a profilin:actin heterodimer to be added to a growing actin filament. The small signal transduction GTPase Cdc42 binds to WASP (78, 79), but not to WIP. Thus, WIP function may be regulated by Cdc42 via WASP. Stimuli that activate Cdc42 may target the WASP-WIP complex to the actin filament system via interactions between the WH1 domain of WASP and the proline-rich ABM1 motifs of structural proteins such as zyxin and vinculin. The presence of SH3 binding motifs in both WIP and WASP, and the capacity to bind to the adaptor protein Nck (86), suggest that the WASP-WIP complex couples additional signaling pathways to the actin filament system.

Formins are profilin and actin-binding proteins from *S. cerevisiae*. They direct the determination of cell polarization and cytokinesis (87 for refs.). In response to a-mating factor, the formin Bni1p forms complexes with the active form of Cdc42, actin, profilin, and actin associated protein Aip3p (88-90). These proteins localize to the tips of mating projections suggesting that the formin is a target for Cdc42, linking the pheromone response pathway to activation of the actin system.

The formin family of proteins is also important in limb and kidney development in vertebrates (91). Mutations in formins lead to defects in the contractile ring formation during cytokinesis, causing the appearance of multinucleated cells. Formins are localized both in the cytoplasm and the nucleus of the cells. FH-proteins, defined by the presence of 'formin homology' regions, are important for a number of actin-dependent processes, including polarized cell growth and cytokinesis. They are large, probably multi-domain, proteins and their function may in part be mediated by an interaction with profilin.

IQGAP is yet another actin-binding, multidomain, regulatory protein (92). It was originally described as a 190-kD protein with extensive sequence similarity to the catalytic domain of RasGaps (91). These proteins control the activity of small regulatory GTPases by stimulating their GTPase activity. Subsequently an IQGAP2 was described as a liver-specific protein with a 62% homology to IQGAP1 (93). However, so far no GAP activity has been found in relation to the GTPases examined. Both IQGAPs have several copies of a ca 50 amino acid long internal repeat (IR), a single WW (SH3-like) domain presumably binding to proline-rich sequences, a number of calmodulin-binding IQ motifs, and a putative actin-interacting calponin homology CH domain. In transient transfection of COS-7 cells with epitope-tagged Cdc42, it was demonstrated that IQGAP is present in the advancing lamellae of motile cells, and a major *in vivo* target of activated Cdc42 (94). The IQGAP1 binds to filamentous actin *in vitro* via its N-terminal domain and cross-links actin filaments in a Cdc42-dependent manner (95). Recently IQGAP1 was implicated in the regulation of E-cadherin-mediated cell-cell adhesion (96).

In fission yeast, *Schizosaccharomyces pombe*, the *rng2* gene codes for an IQGAP-related protein, *rng2p*. The *rng2p* is located in the actomyosin ring and the spindle pole body and is required for the assembly of the actomyosin ring at cytokinesis (97).

Paxillin is not known as an actin-binding protein, but it is involved in the control of the microfilament system at the adhesion contacts. It was detected as a 68-kD phosphotyrosine-containing protein in RSV-transformed cells, and was later purified from smooth muscle cells (98, 99). Many features of paxillin suggest that it functions in signalling events at the adhesive membrane. It is concentrated at focal adhesions, and it is tyrosine phosphorylated in response to extracellular matrix binding and cell transformation. It binds *in vitro* to the elongated tail domain of vinculin, and thereby links vinculin to the SH3 domain of the pp60<sup>csrc</sup>, a component of focal contacts. Integrins, which are the principal transmembrane proteins at these sites, do not seem to bind directly to paxillin, and paxillin does not seem to bind either  $\alpha$ -actinin or actin, but there is evidence for direct association of paxillin to growth factor receptors. Adhesion of cells to an extracellular matrix like fibronectin or laminin via transmembrane integrins results in increased



tyrosine phosphorylation of paxillin, and a similar increase in phosphorylation takes place in connection with clustering of integrins in the plane of the membrane with antibodies. This phosphorylation appears to be caused by activation of the focal adhesion-associated tyrosine kinase, pp125<sup>FAK</sup>, and tyrosine phosphorylation of paxillin and pp125<sup>FAK</sup> initiated by cell adhesion is accompanied by the appearance of organized actin-containing stress fibres. Mutagenesis of paxillin has identified one of the paxillin LIM domains as being responsible for targeting this protein to focal contacts ([100](#))

Zyxin is a low-abundance phosphoprotein localized at sites of cell-substratum adhesion in fibroblasts ([72](#)). It has an architecture of an intracellular signal transducer with a proline-rich domain, a nuclear export signal, and three copies of the LIM motif. LIM is a double **zinc finger** domain found in many proteins that play central roles in regulation of cell differentiation. Zyxin interacts with  $\alpha$ -actinin, members of the cysteine rich protein CRP family, proteins that display Src homology 3 (SH3 domains) and Vasp/Ena family members, which in turn bind to both actin and profilin. Zyxin and its partners have been implicated in the spatial control of actin filament assembly as well as in pathways important for cell differentiation. Based on its repertoire of binding partners and its behaviour, zyxin is thought to serve an organizing centre for the assembly of multimeric protein machines that function both at sites of cell adhesion and in the nucleus.

In summary, the major signal transduction pathways and mechanisms of activation have been shown to be linked to the changes in the dynamic restructuring of the actin cortex: the phosphatidyl inositol cycle and the cAMP second messenger systems and the small GTPase coupled signal transduction mechanisms. All of this points to the growing realization that cytoplasmic signal transduction and transcriptional events in the nucleus are so interwoven that discussions of sequential control by one system or the other are oversimplifications. The idea that actin cortical movements are “downstream” events of “upstream” effectors is similarly unjustified, since actively motile cells are required for the generation of the signal itself.

There are perplexing aspects of polymerization-based propulsion of bacterial pathogens, and by implication filopodial extension. WASP and N-WASP have actin-binding sequences. It is not clear what the role of these proteins is in the regulation of actin polymerization at focal growth sites. It is known that the polar actin filaments are oriented such that their fast growing ends are “pushing” against the moving bacterial wall, or against the membranes of the advancing edges of lamellipodia and filopodia. Actin-binding proteins such as WASP are anchored in these nascent growth complexes (to ActA in the case of *Listeria*). The problem is to explain what happens when an incoming actin monomer is presented as a profilin:actin precursor to the end of the growing filament, since it would appear to be blocked by WASP. Another problem is to account for the free energy of ATP hydrolysis that is made available by the actin polymerization reaction.

### 3. Further Links of Actin Filaments to the Extracellular Matrix

Cortactin is the name of two related proteins p80 and p85, which cross-link actin filaments close to the plasma membrane. Cortactin was first recognized as a substrate for tyrosine-phosphorylation in pp60<sup>src</sup>-transformed cells, but its presence has also been demonstrated in a variety of nontransformed cells ([101](#)). Cortactins have internal repeats of 37 amino acid residues in the N-terminal end of the molecule, which are necessary for their binding to actin filaments. In the C-terminus there is an SH3-domain, and preceding that in the sequence there are proline and serine/threonine rich regions. These latter parts of the molecule are not needed for actin-binding, and their functions remain to be elucidated.

In non-transformed cells, cortactin is normally phosphorylated on serine and threonine, but becomes transiently phosphorylated on tyrosine in response to growth factor stimulation, and on transformation by activated cSrc. In platelets, cortactin becomes transiently tyrosine-phosphorylated in response to stimulation with thrombin. However, tyrosine phosphorylation of cortactin does not

appear to affect its ability to bind to actin filaments, although it does influence its localization to the cortical actin filament system.

Stimulation of cultured cells with various growth factors and platelets with thrombin results in translocation of cortactin into the actin filament system of advancing lamellae and filopodia ([102](#), [103](#)). Activation of fibroblast growth factor receptor-1 results in an association of the activated receptor with cSrc, and a correlated association of cortactin with SH2 domain of cSrc. The redistribution from the cytoplasm into advancing lamellae (ruffles and lamellipodia) appears to depend on the Rac-1-induced activation of the serine/threonine kinase PAK1, a downstream effector of Rac1 and Cdc42 ([104](#)). These GTPases are involved in the control of kinases in the phosphatidylinositol signalling pathway, and there is evidence that PtdIns 4,5-bisphosphate could be involved in the association of cortactin with the activated microfilament system ([105](#)).

Eps8 is a putative actin-binding protein which resembles cortactin in its distribution in the cell ([106](#)). It has an unusual Src homology-3 domain in the carboxy terminus, formed by an intertwined dimer of Eps8 molecules which may affect its peptide specificity. Eps8 is mostly localized in the perinuclear region. However, on stimulation of cells with growth factors, it accumulates in the peripheral cell extensions, where it is found colocalized with cortactin and filamentous actin. The Eps8 pool associated with newly formed lamellipodia and membrane lamellae appears to be mostly detergent-insoluble, as is the tyrosine phosphorylated population of Eps8 molecules generated by activation of vSrc kinase, suggesting Eps8 recruitment to specific sites during cell remodeling.

Ezrin, radixin, moesin (ERM) belong to the band 4.1 superfamily on the basis of their homology to the erythrocyte band 4.1, a protein that connects the actin/spectrin network to the erythrocyte membrane protein glycophorin C ([31](#), [107](#), [108](#)). Like cortactin, the ERM proteins appear to play structural and regulatory roles in stabilizing specialized organizations of actin filaments in connection with the plasma membrane both during development and in adult tissues. They are generally present in microvilli, filopodia, and membrane ruffles, sometimes together, but more often differentially distributed between different types of cells. They link transmembrane proteins with the cortical actin filaments, a process governed by signal transduction involving the Rho-GTPase ([109](#)). Ezrin has been shown to bind specifically to nonmuscle  $\beta$ -actin with high affinity ([110](#)).

The N-terminal domain of ezrin associates with plasma membrane components, and the C-terminal half connects these structures with the submembranous actin filaments. Purified ezrin exists in an inactive conformation that requires activation to expose sites that allow it to associate with the membrane components and with actin filament ([107](#)). It has been shown that binding of ezrin to the transmembrane hyaluronan receptor CD44 ([111](#)) requires activation of ezrin by PtdIns 4,5-bisphosphate binding to the N-terminal domain of ezrin ([112-114](#)). These are not the only transmembrane proteins binding ERM, since ERM proteins have a membrane location also in cells that do not express CD44. It is now known that ERM proteins, in addition to CD44, also bind to intercellular adhesion molecules, ICAM-1 and ICAM-2. A positively charged amino acid cluster of CD44, CD43, and ICAM2 positioned close to the lipid bilayer has been implicated in the binding ([115](#)).

A group of phosphorylated, 50–55 kDa, ezrin-binding protein has been purified from human placenta and bovine brain by affinity chromatography on immobilized N-terminal domains of ezrin or moesin ([31](#)). Isolated EBP50 has two peptide-binding PDZ domains in the N-terminal half of the molecule with capacity to bind to cytoplasmic tails of transmembrane proteins. The EBP50 can be coprecipitated with ERM and may serve to bind ERM members to one or more integral membrane proteins.

Cells opened up by the use of the detergent digitonin are still amenable to activation with non-hydrolyzable analogues of GTP. Addition of guanosine 5'-O-(3-thiotriphosphate), GTP $\gamma$ S, to such permeabilized cell models induces assembly of both actin filaments and focal adhesion complexes through activation of endogenous Rho and Rac ([66](#)). The sensitivity to GTP $\gamma$ S is lost after a few

minutes of incubation, but can be restored by the addition of the actin-binding protein moesin indicating that moesin is required for the effect of Rho and Rac on actin filament system. This is an important step towards the resolution of the signal relay which initiates the reorganization of the microfilament system.

Talin is an elongated, flexible actin-binding protein first recognized in smooth muscle (116, 117) and platelets (118, 119). It is present in focal contacts in tissue cultured cells, in adhesion plaques formed by activated platelets, in myotendinous junctions, and in dense plaques of smooth muscle, but is absent from cell:cell interaction sites (for refs, see 38). In motile cells talin is present in membrane ruffles. It is a major protein in platelets constituting more than 3% of the total protein. In response to activation, it is redistributed to the actin-rich platelet cortex (120). Thus talin is clearly of central importance for the association of actin filaments to membrane components to sites where cells make close contact with extracellular matrix proteins, and it has been shown to interact with the cytoplasmic domains of transmembrane b1-integrin (116, 121).

Talin has a molecular mass of of 270 kD as estimated from the gene sequence (122). It is about 60 nm in length, and appears to be composed of a series of globular domains as judged by electron microscopy (123). In solution, it forms homodimers by antiparallel association of the monomers. It binds to integrin, vinculin, and actin *in vitro*. There are reports that talin can nucleate, cap and cross-link actin filaments, and that its activities may be controlled by phosphorylation by both serine/threonine and tyrosine kinases (124, 125). An N-terminal domain of talin is homologous to ezrin (122), and may be the part of the molecule which associates with membrane lipids. The larger C-terminal domain contains the binding site for vinculin (see (38) and (125) for further refs.). Disruption of the talin gene in embryonic stem cells leads to a disturbance in the formation of vinculin and paxillin-containing focal contacts and cell spreading (126). The talin (-/-) cells do form embryoid bodies, and notably two morphologically distinct cell types have been observed that were able to spread and form focal adhesion-like structures with vinculin and paxillin on fibronectin. Thus, although talin was essential for the expression of b1-integrin, a subset of differentiated cells managed without it. What substitutes for talin in these cells is unclear.

Vinculin is a 117-kD, actin-binding protein found concentrated at all sites of attachment of filamentous actin to transmembrane components of adherens junctions (38, 127-131). At these sites, it enters into multiple interactions with other proteins at the cytoplasmic face of the contacts. It is critically required as a structural component in the formation of the junctions. On ligand-induced clustering of integrins, vinculin is recruited from the cytoplasm to sites of integrin-ligand interaction, where it is a substrate for tyrosine phosphorylation by pp60<sup>src</sup>.

There is electron microscopic evidence that the vinculin molecule can attain different conformational states. It is found either as a compact, globular molecule, or it is unfolded into a globular head with an extended tail (132-134). In the compact globular state, its actin-binding capacity is masked by an intramolecular association of the 95-kD head and 30-kD tail domains. The actin-binding site can be exposed by binding of acidic phospholipids to the protein suggesting that the affinity of vinculin for target molecules is regulated by modulation of the head-to-tail interaction, and that this may be used in the control of the assembly of adherens junctions (73, 135). The phospholipid PtdIns 4,5-bisphosphate binds to two discrete regions of the vinculin tail, disrupts the intramolecular head-tail interaction and induces vinculin oligomerization. The binding of PtdIns 4,5-bisphosphate to vinculin also enhances the association of vinculin to VASP the actin- and profilin-binding protein present in focal adhesion sites. The C-terminal tail region contains a binding site for the focal adhesion protein paxillin (136). Furthermore, vinculin has been identified as part of the cadherin-catenin junctional complex (137-139).

**Tensin** is a dimer of two 200 kD polypeptide chains. It is found in different types of cell contacts, including focal adhesions, zonula adherens, intercalated discs, and myotendinous and neuromuscular junctions (140-143). It contains three distinct, actin-binding sites per monomer and appears to cap as

well as cross-link actin filaments. Tensin also binds vinculin and possibly tyrosine kinase receptors through a Src homology 2 (SH2) domain ([144](#)). It has been suggested that a tensin dimer with its 6 actin-binding sites might wrap around the end of an actin filament, and bind it to a tyrosine kinase receptor through its SH2 domain. Induction of the formation of a focal contact, by the binding of extracellular ligands to integrins, causes tensin to be phosphorylated on tyrosine. Insertin, an actin-binding protein suggested to be involved in the actual incorporation of monomers into filaments, has been recognized as a proteolytic breakdown product of tensin ([145](#)).

#### 4. Proteins Controlling Polymerization of Actin

**DNase I** (mw 31 000) was the first actin monomer-binding protein found ([12](#)). However it is doubtful whether DNase I plays a role as an actin monomer-binding protein in the cell. Since the enzyme is produced by acinar cells in the pancreas and secreted into the duodenum, it has been thought of as a digestive enzyme. However, intracellular forms of the enzyme have been observed, and there is evidence that it might be involved in programmed cell death, apoptosis. ([146](#)).

The crystallization of DNase I was reported in 1948 by Kunitz ([147](#)). The same year Laskowsky and collaborators described the presence of a proteinaceous inhibitor in tissues from vertebrates ([148](#), [149](#)). Two DNase I inhibitors (I and II) were isolated and crystallized in 1966 ([13](#), [150](#)). They were shown to form gelifying high molecular weight aggregates. Searching for the identity of the inhibitors eventually led to the finding that actin was the ubiquitous inhibitor of DNase I ([12](#)). Analysis of the crystals of the DNase I inhibitor II revealed that they contained actin together with an equimolar amount of a small protein, now known as profilin ([14](#), [15](#)). Later, the inhibitor I and II were identified as profilin in complex with g-actin and b-actin, respectively ([151](#)).

DNase I is a glycoprotein. It forms a stable 1:1 complex with monomeric actin, and when mixed with equimolar amounts filamentous actin, it causes depolymerization of the filaments and formation of the 1:1 complex with actin ([152-154](#)). Kabsch and coworkers provided the first insights into the structure of actin by solving the structure of the DNase I:actin complex ([155](#); see special entry on actin). The high resolution structure of DNase I alone and complexed to an oligonucleotide has been reported as well ([156](#), [157](#)). The current model of the actin filament was arrived at using the structure of the actin monomer, as it occurs in the DNase I:actin crystal, together with diffraction data obtained with oriented gels of F-actin ([158](#)). This model is claimed to fit well into reconstructions of actin filaments from electron microscopic images ([159](#), [160](#)), and reconstructions of actin filaments with bound actin-binding proteins are interpreted in terms of this model (see below). However, doubts about this model being a relevant representation of the actin filament have been expressed, and an alternative way of deriving a model of F-actin has been suggested from the analysis of the organization of actin monomers in crystals of profilin and actin ([161](#), [162](#)). To ascertain the validity of the current model, independent information about the orientation of the actin monomer in the actin filament needs to be acquired.

Profilin (mw 12 000–15 000) is an abundant actin monomer-binding protein in eukaryotic cells ([163](#), review). (See separate article [Profilin](#).) Particularly high concentrations of profilin are found in lymphoid cells and brain cells. In lymphoid tissues the concentration of unpolymerized actin is about 100  $\mu\text{M}$  and in platelets it may be as high as 200  $\mu\text{M}$ . In both cases about 50% of the unpolymerized actin can be accounted for by profilin:b-actin and profilin:g-actin isoforms. The remaining unpolymerized actin appears to be sequestered by b-thymosins (see below). For references on specific aspects of profilin, see ([164-189](#)).

Observations point at [profilin](#) being an important factor in the control of the release of inositol trisphosphate and thus of the generation of  $\text{Ca}^{2+}$  pulses in the cell. The primary targets for  $\text{Ca}^{2+}$  regulation in cells is the actomyosin system with global changes in structure and activity as the result. It is now known that many of the microfilament-associated proteins interact with components formed in the phosphatidylinositol-cycle as a result of receptor-mediated activation of the cell, and

that these interactions modulate the activity of these proteins *vis-à-vis* actin. Although the exact physiological roles of these interactions remain to be elucidated, it suggests that the phosphatidylinositol-cycle is directly involved in controlling the microfilament-based motility cycle (see Fig. 3).

Beta thymosins constitute a conserved family of polypeptides (mw 5 000), the members of which were originally isolated from mammalian thymus and thought to be important for the immunoregulatory activities of this organ (190). Homologues of Tb4 have been found also in distantly related organisms (191). The thymosin b4 isoform is abundant in mammalian tissues, and has attracted much attention after the discovery of its association with unpolymerized actin in platelet extracts. It has a poorly defined structure as determined by NMR in solution (192). Binding studies suggested that it binds to the (+)-end of the actin monomer in competition with profilin. However, crosslinking studies imply a more complex interaction involving sites at both ends of the actin molecule (193).

Binding of thymosin b4 to actin monomers inhibits incorporation of actin monomers at both ends of the actin filament. Thymosin b4 has a rather low affinity for actin ( $K_d = 5 \mu\text{M}$ ), but since it is present in high concentrations in cells (in platelets ca 500  $\mu\text{M}$ ) it is considered to be the most important actin sequestering protein. In the presence of polymerization-inhibiting concentrations of thymosine b4, addition of profilin overcomes the thymosin effect, due to the higher affinity of profilin for actin ( $K_d = 0.25 \mu\text{M}$ ), and the fact that profilin:actin can add actin monomers onto the fast growing ends of actin filaments (see [Profilin](#)).

*In vitro*, thymosin can be shown to bind also to actin filaments and decrease the critical concentration for actin polymerization (194), and experiments on cultured cells indicate that  $\beta$ -thymosins may not be just actin sequestering factors (195). In the NIH3T3 cell line thymosin b10 is the predominant beta thymosin isoform, being present in about equimolar concentration to profilin and ADF/cofilin. Overexpression of thymosin b10 in these cells led to an increase in the cellular content of polymerized actin without any change in the total actin content of the cells. An intimate relationship with the regulation of the dynamic functioning of the microfilament system is also indicated by the increase in motility caused by overexpressing thymosin b10. Consonant with this an isoform of beta thymosine, thymosin b15, is upregulated in some forms of malignant cells (196).

The ADF/cofilin family of proteins include actin depolymerizing factor, cofilin, destrin, depactin, actophorin (197-199). These proteins are found in all kinds of eukaryotes, and when tested genetically they have proved to be essential. *In vitro*, they can form 1:1 complexes with actin monomers as well as with actin in filamentous form. The structure of the prototype ADF/cofilin, actophorin from *Acanthamoeba* have been solved by crystallography (200) and there are NMR structures of yeast cofilin and destrin (201, 202). It has recently been recognized that an ADF homology domain is present in each member of a newly identified protein family consisting of the ADF/cofilins, the twinfilins, and the drebrin/Abp1s (203).

In the cell, most of the ADF/cofilin is found in the perinuclear area, but there are also significant amounts in highly motile advancing lamella and filopodia, where the reorganization of actin through polymerization/depolymerization is most intensive. The interaction between ADF/cofilin and actin is sensitive to the pH of the medium, and as in the case of profilin, polyphosphoinositides interfere with the ADF:actin interaction. Physiologically, however, the most important mechanism controlling the activity of ADF/cofilin depends on phosphorylation/dephosphorylation of the protein. Phosphorylation of ADF/cofilin at a site near the N-terminus (S3) blocks its binding to actin. In cell extracts about 60% of the protein is phosphorylated, and in connection to receptor-mediated stimulation of cells dephosphorylation of ADF/cofilin takes place. The phosphatase involved in this reaction has not been identified, but a specific protein kinase, LIM kinase, has been found responsible for the *in vivo* phosphorylation of ADF/cofilin (204, 205).

The turnover of actin filaments inside living cells appears to be up to 100-fold faster than in *in vitro* experiments with actin alone. It has now been demonstrated in *in vitro* as well as in *in vivo* experiments that ADF/cofilin can accelerate the turnover of actin filaments (17, 199, 206-208). ADF/cofilin binds with higher affinity to ADP-containing actin, monomeric as well as filamentous, and facilitation of filament growth by active ADF/cofilin is likely to be brought about by increased depolymerisation of the actin filaments at their proximal ADP-containing ends increasing treadmilling activity. Nucleotide exchange, possibly enhanced by profilin, would in this model precede the reallocation of the ATP-containing precursor to the fast growing, distal (+)-end of the population of growing filaments. The presence of ADF in lamellipodia, where filaments are long, argues against a severing function *in vivo*, and the end-specific effects of ADF on filament assembly seen *in vitro* also indicate that the severing activity of ADF/cofilin is less important *in vivo*.

Twinfilin is a newly discovered protein composed of two cofilin-like regions (209). It was identified and characterized as an actin monomer-binding protein in budding yeast. Genes coding for homologous proteins have been recognized also in *Caenorhabditis elegans*, humans and mice. The two halves of the molecule twinfilin I and II resemble the corresponding domain from other species more than they resemble each other. Thus, twinfilins have evolved from a common ancestor and the twinfilins represent a single protein family (199). Twinfilin does not appear to act as an actin filament depolymerizing factor, but there are *in vivo* observations that suggest its involvement in the control of the dynamics of the microfilament system.

#### 4.1. Arp

The first reports describing a family of actin-related proteins, the Arps, appeared a few years ago (2, 209). These are highly conserved proteins that share a 30–60% homology with actin, but are functionally distinct from actin. Two of these actin-related proteins, Arp2 (44 kD) and Arp 3 (47 kD) have attracted a great deal of attention, since they were discovered to exist in a large complex together with five other proteins (40 kD, 35 kD, 19 kD, 18 kD, and 15 kD) that could be isolated by affinity chromatography on immobilized profilin (210). The Arp2/Arp3 complex has been reported to nucleate actin filament formation, to bind along the sides of actin filaments and to express filament crosslinking activity. A similar complex was identified in searching for factors that initiate actin assembly at the surface of *Listeria* (211). The Arp2/Arp3 complex is highly conserved among eukaryotes. Null mutations are lethal. It is localized in the cortex of amoebae and yeast and in the lamellipodia of higher eukaryotes (212, 213). In yeast, the Arp2/Arp3 complex is required for the integrity and motility of actin patches and for endocytosis (214). Thus the complex appears to play a major role in actin-based motility.

It has now been demonstrated that the Arp2/Arp3 complex isolated from *Acanthamoeba* binds profilin (215) and to the (–)-ends of actin filaments (168). It is not an efficient nucleator of actin polymerization. Capped filaments grow by the addition of monomers to the (+)-end. The complex can be seen by electron microscopy to attach the (–)-end of one filament to the side of another filament. This suggests that the Arp2/Arp3 complex might control the assembly of a branching network of actin filaments in the advancing lamellae of motile cells. Similar branching of actin filaments have been recognized earlier in electron microscopic pictures of the weave of microfilaments in advancing lamellae in moving keratinocytes (36). The cellular concentrations of the Arp2/Arp3 complex has been estimated to be high enough to cap the (–)-ends of all filaments in a cell. These observations imply that models assuming that actin filaments *in vivo* undergo an ADF/cofilin-aided depolymerization from the (–)-end during motile activity may be too simplistic, since there must also be mechanisms for the controlled uncapping of the (–)-ends of the filaments.

Further studies of the *Listeria* system have shown that the bacterial protein ActA, which is essential for the actin polymerization-dependent propulsion of the bacterium through the cytoplasm, acts synergistically with the Arp2/Arp3 complex *in vitro*. Working together, they eliminate the lag phase in the polymerization of actin and increase the initial rate of assembly 50-fold (216). This activity has been shown to reside in the N-terminal domain of the ActA protein. A cellular protein that corresponds to ActA has not yet been identified, but zyxin has been pointed out as a plausible

candidate.

Capping protein Z (CapZ) is a heterodimeric protein, which was first detected in skeletal muscle sarcomers, where it is associated with actin filaments at the Z-line end of sarcomers, i.e. at the (+)-end of the filament (217). *In vitro*, it binds with high affinity ( $K_d \approx \text{nM}$ ) to the (+)-end of actin filaments. A homologous protein that exists in non-muscle cells, capping protein  $b_2$  (CPb<sub>2</sub>), is enriched at the edge of the advancing lamellum of spreading platelets. By binding to the (+)-end of actin filaments, it blocks the addition of actin monomers to that end, and since actin filaments *in vivo* appear to grow by assembly primarily at that end, these capping proteins may be important regulators of actin polymerization. The fact that CPb<sub>2</sub> is necessary for the proper organization of actin bundles in the morphogenesis of bristles in *Drosophila*, the view is strengthened that CPb<sub>2</sub>, gelsolin in some cases, and Cap G may play active roles in actin-based motility in response to signaling. There is evidence also in this case that polyphosphoinositides are involved in the regulation of the protein.

#### 4.2. Adducin

In erythrocytes there is, in addition to CapZ, a second actin filament (+)-end capping protein, adducin (218, 225). This protein is associated with the short erythrocyte actin filaments (one a,b heterodimer per filament). Adducin has been shown to bind to spectrin:actin complexes in the cell cortex, and promote the binding between these two proteins. It also appears to be able to bind directly to actin filaments and to bundle them. Adducin has a relatively low affinity for the filament (+)-end, but even so it appears to be the capping protein bound to the erythrocyte actin filaments, and not CapZ. Its binding to actin filaments is downregulated by calcium and calmodulin, an effect possibly modulated by reversible phosphorylation of the protein. The complex relationship between adducin and the spectrin:actin network indicates that adducin represents a new type of regulated actin filament-binding and barbed end capping protein. The entire molecule is required for capping activity, and the association of adducin with actin filaments seems to be regulated by Rho-associated kinase and myosin phosphatase.

Gelsolin is an abundant and ubiquitously expressed actin-binding protein, which confers calcium sensitivity to the regulation of the dynamics of the microfilament system (219). (See separate article [Gelsolin](#).) For references on specific aspects of gelsolin, see (219-223).

Tropomodulin is a 43 kD protein that caps the (–)-ends of actin filaments in erythrocyte cortical network and in striated muscle sarcomeres (225). The interaction is particularly tight when tropomyosin is bound simultaneously. It is thought that tropomodulin is important in controlling the length of the actin filaments in these situations.

### 5. Architectural Elements of the Actin Cortex

The dynamic changes in cell surface morphology and motile behaviour, seen in response to growth hormones and other cell stimulating agents, reflect changes in the formation, organization, and activity of the microfilament system. Actin-binding proteins that cross-link actin filaments into tightly packed bundles and more loosely organized networks are necessary for the formation of the weave of microfilaments seen in membrane lamellae and filopodia and in the rest of the submembraneous zone around the cell. Actin-based surface projections on cells take many forms, such as the elaborate stereocilia of cochlear hair cells, microvilli of intestinal epithelium, and the membrane ruffles and filopodia on moving fibroblasts. The formation of these structures depends on the activity of actin-crosslinking proteins. Important members of this class of proteins belong to the spectrin superfamily which are reviewed in ref. 237.

Villin is closely related to gelsolin, an actin severing protein (described above), but distinguished from gelsolin by the presence of an additional domain comprised of 76-amino acid residues, called

the “headpiece” (224-229). The headpiece contains an additional actin-binding site, which confers actin filament bundling activity to the protein. Villin is found in the microvilli of certain types of epithelia. Crosslinking of actin filaments occurs at low  $\text{Ca}^{2+}$  concentrations. Like gelsolin, villin severs actin filaments in the presence of micromolar concentrations of  $\text{Ca}^{2+}$ , and after breakage of the actin filament, villin remains bound to the (+)-end of one of the fragments. The actin-binding properties, tissue distribution and expression during cell differentiation suggest that villin is important in organizing the actin filament core of microvilli in epithelial cells forming brush borders. In support of this, villin expressed transiently in transfected fibroblasts, results in loss of stress fibers and appearance of large numbers of microvillar protrusions on the dorsal surface of the cells. Microinjection of high concentrations of villin into cells that normally lack this protein, lose their stress fibres, and incorporate the villin into cell surface microspikes and large microvilli, an activity which was shown to depend on the presence of the C-terminal head piece.

The N-terminal domain of villin (14T) consists of 126 amino acid residues. It binds calcium and actin. The structure of this domain has been determined by NMR (230). The 76-amino acid residue long, C-terminal domain of villin is an actin-binding module that is also used in another actin-bundling protein **dematin** (band 4.9), but whose core domain is unrelated to villin. The 3-dimensional structure of the stably-folded 35 residue long subdomain of the villin headpiece has been determined by NMR (231). It is the smallest autonomously folding protein unit found.

Fimbrin belongs to the plastin family of actin-binding proteins (31, 232, 237). Members of this family are present in all types of eukaryotic cells. Fimbrin is the smallest of the actin crosslinking proteins. It consists merely of an N-terminal EF-domain and two actin-binding domains, each composed of two CH domains. Fimbrin is present together with villin in the microvilli of the intestinal brush border, and in the membrane lamellae and microspikes of most other eukaryotic cell types. A fimbrin-binding site on actin has been suggested from analysis of actin mutations suppressing fimbrin mutations (233). The structure of the N-terminal actin-crosslinking domain of fimbrin has been solved by crystallography (234), and electron cryomicroscopic images of fimbrin-decorated actin filaments has been used to localize the fimbrin-binding site on the surface of the actin filament (235).

In a suggested scenario for the formation of a microvillus (236), the first phase consists of the “streaming” of actin filaments from electron-dense plaques on the apical plasma membrane. These filaments then gathered into loose bundles that extend into rudimentary microvilli. Villin may contribute to an initial bundling of the filaments. Ezrin (see section 3), a protein controlled by tyrosine phosphorylation, links the microvillar actin filaments to the membrane components. During a second phase, the microvilli elongate and become highly organized into a full core of hexagonally-packed actin filaments. This phase correlates in time with the localization of fimbrin to the core of the microvillus. In this way a stable surface protrusion is thought to take form. Similar schemes may apply to the development of surface protrusions on the surface of other types of cells.

ABP-280, filamin, and ABP-120, are all rod-like, dimeric, actin filament crosslinking proteins, where each monomer has an N-terminal actin-binding CH-domain and a long repeat region consisting of tandemly repeated, 100 amino acid residue long, homologous segments (237). The length of the repeat region is different in the different ABPs. The ABP120 was found in *Dictyostelium discoideum*. The vertebrate variants, ABP-280 and filamin, are major cellular proteins, and products of different genes. Filamin is located specifically in the Z-lines structure of the sarcomeres in muscles. ABP-280 is a major cellular protein in non-muscle cells, where it is present in the cell cortex. The function of ABP-280 appears to be to cross-link actin filaments into orthogonal networks. The 3D-structure of the repeat region of ABP-120 (gelation factor, 6 repeats) has been determined by NMR spectroscopy, which shows it to have an immunoglobulin-like fold (226). Both filamin and ABP-280 (24 repeats) share conserved residues that form the core of the gelation factor repetitive segment structure as determined by NMR. This distinguishes these actin crosslinking proteins from the spectrins and  $\alpha$ -actinins, which have tandem repeats of an  $\alpha$ -helical



domain.

On the basis of their actin filament crosslinking activity, depending on the presence of a calponin homology CH domain, and their elongated shape, as seen by electron microscopy, these proteins have been grouped together with the spectrin and fimbrin families into a superfamily. It is true that the actin binding domains of all these proteins are homologous. However, the presence or absence of repeat domains, and the widely different structures of the repeat domains of the spectrin and filamin families, rather suggests that the three families should be considered as separate protein families employing a common actin-binding domain, since it is possible that the different families will turn out to have quite different functions in relation to the actin filament system.

Spectrin (fodrin) was first isolated as part of the matrix proteins of the erythrocyte plasma membrane (237-240). Underneath the lipid bilayer of these cells, and linked to transmembrane proteins glycophorin and the ion channel band 3, there is a regular organization of units consisting of five or six spectrin molecules attached to a short actin filament. The spectrin molecules are attached to spectrins of adjacent units to form a sheet of five- and six-sided polygons. The membrane skeletons, and isolated junctional complexes, reproducibly contain four proteins in the molar ratios of 1 spectrin dimer:2-3 actin monomers:1 band 4.1 protein:0.1–0.5 protein 4.9 (protein designations refer to mobility on gel electrophoresis). In addition tropomyosin, tropomyosin-binding protein and adducin are thought to be part of the junctional complex.

Spectrin consists of two subunits, an  $\alpha$  and a  $\beta$  subunit, which form heterodimers, which in turn associate to form tetramers in the membrane lattice. The  $\beta$ -chain has an N-terminal actin binding CH domain and 12  $\alpha$ -helical repeat domains, whose structure has been determined by NMR spectroscopy (241). The  $\alpha$ -chain associates in an antiparallel fashion with the  $\beta$ -chain. It has a C-terminal EF-hand domain, which is close to the actin-binding site of the  $\beta$ -chain followed towards the N-terminus by a binding site for protein 4.1, and a Src-like sequence in the middle of the somewhat longer repeat region. Between repeat 11 and 12, there is a calmodulin-binding sequence inserted. The structures of the spectrin repeat, the phosphoinositide-binding pleckstrin (PH) domain, the SH3 domain, and the CH domain have all been determined (241-246).

The spectrin-actin based network of proteins endows the erythrocyte with an exquisite elasticity which makes it possible for it to go through quite drastic shape changes. However, it is clear that the structural role played by the spectrin:actin system is not its sole function. The close relationship between the spectrin-actin system and transmembrane ion channels implies that the range of controlling functions is wider. The spectrin-actin system is not unique to the erythrocytes, but exists also in other tissues and in cells of as widely separated organisms as *Drosophila*, *Acanthamoeba*, *Dictyostelium*, echinoderms, and possibly also in higher plants.

Alfa-actinin is a classical skeletal muscle protein identified a long time ago (237, 247). It is highly enriched at the Z-disc-end of the actin filaments in the sarcomeres of the myofibrils. It is now known also to be an integral part of the microfilament system in nonmuscle cells, where an interaction between transmembrane  $\beta 1$  integrins and  $\alpha$ -actinin has been demonstrated (see Fig. 2). There are several isoforms of  $\alpha$ -actinin, which are products of three different genes. *In vitro*, it crosslinks actin filaments into regular ladder-like structures. It is therefore thought to contribute to the organization of actin in cells. Alfa-actinin is an elongated homodimer (200-215 kD) with antiparallel arrangement of the two subunits. Each subunit consists of a CH-domain, 4  $\alpha$ -helical domains (“spectrin repeats”), and an EF-hand domain. The binding of nonmuscle  $\alpha$ -actinin to actin filaments is modulated by  $\text{Ca}^{2+}$  binding to the EF-domain, whereas the muscle isoform of  $\alpha$ -actinin is insensitive to  $\text{Ca}^{2+}$ . The 3D-structure of the CH-domain is known by analogy to the structure of the corresponding domain from calponin solved by x-ray crystallography, and the structure of the repeated domain is known by analogy to the known  $\alpha$ -helical repeat domain in spectrin, and the EF-hand can be compared with calmodulin.

Alfa-actinin can be extracted from myofibrils without major destruction of the sarcomeric organization of the actin filaments in the myofibril, and when added back to the myofibril preparation it rebinds at the Z-disc. This implies that a-actinin may serve some function, other than a structural, in relation to the actin filaments. In nonmuscle cells, a-actinin is found enriched in spots along stress fibres and in cell adhesion plaques.

Dystrophin is the product of the gene responsible for the X-linked myopathies Duchenne and Becker muscular dystrophy. It is an elongated protein of high molecular weight (430 kD) present in low amounts in muscle and nerve cells. A closely related protein, utrophin, is more widely distributed, and there are many dystrophin/utrophin isoforms produced, due to the operation of alternative promoters and alternative splicing. The structure of these proteins place them in the spectrin/a-actinin family of proteins. They are multidomain proteins; the largest isoform is comprised of 3 N-terminal and 2 internal actin-binding sites, a peptide-binding WW domain, a calcium-binding EF domain, a zinc finger domain. The N-terminal actin binding region has sequence homology with the calponin homology CH domains of the spectrin/a-actinin family of proteins.

Dystrophin/utrophin connect cortical actin filaments with transmembrane proteins, dystroglycans and sarcoglycans, which in turn associate with extracellular matrix proteins, laminin and merosin, respectively. Skeletal muscle dystrophin can be purified from muscle cells as a large multiprotein dystrophin-glycoprotein complex, which stabilizes actin filaments *in vitro* through lateral associations. Both dystrophin and utrophin bind with higher affinity to b- than to a-actin. Comprehensive reviews on the structure and function of the dystrophin/utrophins are found in ref [237](#), [248](#).

## 6. Ion Channel-associated Actin-binding Proteins

The spectrin-actin network in erythrocytes is involved in anchoring ion exchange proteins in the membrane via the linker protein [ankyrin](#), a critical link in maintaining the characteristic biconcave shape of the cell ([249](#)). In axons, an interaction between actin, [fodrin](#), and ankyrin appears to be responsible for concentrating sodium channels at the nodes of Ranvier ([249](#)). At the neuromuscular junction, the clustering of acetylcholine receptors appears to be mediated by spectrin, actin, and the **rapsyn/43 kD actin-binding protein** ([250](#)). Various lines of evidence suggest that intact actin filaments, and perhaps actin polymerization itself, are required for calcium regulation of the postsynaptic response of NMDA receptors in the central nervous system ([251](#)). Postsynaptic excitatory synapses are most often found on the small actin-rich budlike structures, called spines, that protrude from the dendrites of highly arborized neurons such as cerebellar Purkinje cells ([252](#)). Dendritic spines have pronounced electron-dense (as seen by electron microscopy) post-synaptic densities, called PSD's, that are the likely sites of proteins linking receptors to actin. Candidate proteins include PSD-95/SAP90, chapsyn/PSD-93, SAP102, and **alpha-actinin-2**. Only alpha-actinin has a demonstrated actin binding affinity, and it has been suggested that clustering for some classes of receptors might be mediated by a looser, more "corral-like", entrapment mechanism ([253](#)). The PSD contains brain [dystrophin](#), an isoform of the muscle protein originally identified as the mutated gene product involved in muscular dystrophy that is known to link transmembrane ion channels to actin filaments ([254](#)).

The precise role for actin in the functioning of ion channels is unclear. The provision of clustered anchorage sites may confer cooperativity in the opening of channels, or support the formation of multiprotein complexes capable of integrating different signalling or environmental factors. It is not unlikely that the actin network might passively or actively transmit forces to the channels providing the impetus behind stretch-activated gating phenomena.

## 7. Myosin as an Actin-binding Protein

The myosins belong to a large superfamily of proteins, which together with actin transduce chemical energy to force generating tension and movements in the eukaryotic cells. Myosin II is the prime

partner of actin in the generation of force in muscle, as well as in nonmuscle cells. Generally myosins are thought of as the force generators, and therefore referred to as motor molecules. However, the actual mechanism of force generation is still unknown. Although, the ATP-binding head domain is largely conserved, there are many different molecular forms of myosin serving in as diverse functions as muscle contraction, cell motility, membrane traffic, and sensory perception. Myosin hydrolysis ATP to ADP.Pi still bound to the head domain. Interaction with actin filaments brings about product release, and in the presence of actin filaments, there is a rapid ATP-dependent cycling of heads on and off the actin filaments as ATP is hydrolysed. A current review of the myosin superfamily and their functional involvements is found in ref. [255](#). For a more detailed account of the structure and function of myosin see the entry on myosin in this encyclopedia.

## 8. Regulators of the Actomyosin Interaction

The N-terminus of actin is involved in the binding of a large number of actin binding proteins, including myosin (S1), tropomyosin, troponin I,  $\alpha$ -actinin, caldesmon, gelsolin, cofilin, actobindin. This actin interface is important for muscle contraction, since it is implicated not only in the activation of the myosin ATPase, but also in the regulation of actomyosin interaction by interacting with troponin I.

Tropomyosin is an elongated, dimeric, coiled-coil  $\alpha$ -helical protein that binds along the actin filament (see special entry). The tropomyosin molecule has a sevenfold repeat of nonpolar and polar amino acid residues that bind seven actin monomers in the filament. Multiple genes for tropomyosin and alternative splicing can generate a number of tropomyosins that are differentially expressed during development and in different cell types. Muscle cells express 1–2 and most vertebrate nonmuscle cells 3–8 tropomyosin isoforms. Tropomyosin in complex with the heterotrimeric protein troponin makes muscle contraction  $\text{Ca}^{2+}$  dependent, but the function of tropomyosin in nonmuscle cells has not been clarified. However, the involvement of tropomyosin in the control of chemomechanical transduction also in nonmuscle cells is indicated by the occurrence of tropomyosin in tightly bundled and contractile organizations of actin filaments (stress fibres). For further information see special entry on tropomyosin.

Troponin is a heterotrimeric protein complex that interacts directly with actin through one of its subunits. Together with tropomyosin it confers calcium sensitivity to muscle contraction in striated muscles (see special entry).

Caldesmon is involved in the regulation of the actomyosin interaction in smooth muscle ([256](#), [257](#)). It inhibits actomyosin ATPase and filament sliding *in vitro* has been shown to bind to subdomain 2 of actin. Three-dimensional image reconstruction of reconstituted thin filaments consisting of actin and smooth muscle tropomyosin, both with and without a caldesmon derivative added has been reported. In filaments containing the caldesmon derivative, tropomyosin was found in a different position as compared to the situation in the absence of the caldesmon. The observations suggest that caldesmon causes changes in the relationship between tropomyosin and the actin filament, which are different from those seen with troponin.

The giant modular protein nebulin (mw 600- 900 kD) spans the whole length of the thin filament of the striated muscle sarcomeres in vertebrates, and has been proposed to function as a “ruler” for control of the length of actin filaments in sarcomeres ([258-260](#)). It is thought to bind and stabilize F-actin. It comprises 2-3 % of the myofibrillar protein mass of skeletal muscles. There are tissue and development-specific isoforms. The C-terminal part of human nebulin is anchored in the sarcomere Z-disk and contains an SH3 domain. The nebulin SH3 sequence from several species has been determined and found strikingly conserved. Its 3D-structure has been determined in solution by NMR spectroscopy, and its interaction with poly(L-proline) has been modelled. Nebulin consists of nearly 200 tandem repeats of homologous modules about 35 residues long. These are organized into about 20 tandem superrepeats consisting of 7 different modules, and this superrepeat segment is

flanked near the N- and C-termini by single-repeat regions containing 8 modules of the same type. It has been proposed that the 35 residue module is the basic actin-binding domain and that the superrepeats reflect tropomyosin/troponin binding sites along the nebulin molecule. Exactly how the nebulin molecules are linked to actin filaments in the sarcomere of muscle is not yet known. Nebulin promotes actin nucleation and stabilizes actin filaments. Crosslinking experiments have identified the first two residues in actin to be involved in binding nebulin.

## 9. Future Directions

One of the most remarkable aspects of the contractile apparatus in muscle cells is the capacity of the system to ramp up its power output in response to increased demands for mechanical work. This phenomenon is known as the Fenn Effect, and it suggests that the acto-myosin system can increase its ATPase activity in response to higher imposed loads. It is apparently the lattice organization of filaments into sarcomeres that somehow enables the effect of external loads to be transmitted directly to the proteins producing the biochemical reactions that produce work from ATP hydrolysis. It remains to be seen if the bundles and meshworks found in non-muscle cells display similar, or even more unusual, chemomechanical feedback mechanisms. Muscle cells also perform work at nearly 100% thermodynamic efficiency, and it is interesting to ask whether the free energy transduction pathways in the cytoplasm are equally optimized for their tasks.

The newly discovered connections between the microfilament system and signal transduction raise a host of questions that will be the subject of future research. What roles do mechanical forces play in transmitting signals from the surfaces of cells to their nuclei? What regulates the assembly of these focal sites of signalling potential and what leads to their breakup after a signalling cascade has outlived its usefulness?

The large number of actin-binding proteins obscures the possible branching points that must have occurred during evolution to establish this incredibly diverse dynamical system. The analysis of evolutionary relationships in actin-binding proteins from amoeba, plants, and mammals will be of extraordinary interest for understanding how the immune system and the mammalian nervous system developed as specializations of more primitive motile mechanisms.

## Bibliography

1. F. B. Straub (1942) Studies from the Institute of Medical Chemistry, Univ. Szeged **2**, 3–15.
2. F. B. Straub (1942) Studies from the Institute of Medical Chemistry, Univ. Szeged **3**, 23–37.
3. A. G. Engel and C. Franzini-Armstrong (1994) *Myology*, vol. **1**, 2nd edition. McGraw-Hill, Inc. New York.
4. B. Alberts, D. Bray, A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walter (1998) *Essential Cell Biology*, Garland Publishing Inc., New York & London.
5. S. Hatano, T. Totsuka, and F. Oosawa (1966) *Biochim. Biophys. Acta* **127**, 488–498.
6. S. Hatano and M. Tazawa (1968) *Biochim. Biophys. Acta* **154**, 507–519.
7. M. R. Adelman and E. W. Taylor (1969) *Biochemistry* **8**, 4964–4975.
8. H. Ishikawa, R. Bishoff, and H. Holzter (1969) *J. Cell Biol.* **43**, 312–328.
9. E. Lazarides and K. Weber (1974) *Proc. Natl. Acad. Sci* **71**, 2268–2272.
10. K. Weber and U. Groeschel-Stewart (1974) *Proc. Natl. Acad. Sci* **71**, 4561–4564.
11. T. D. Pollard and R. R. Weihing (1974) *Critical Reviews in Biochemistry* **2**, 1–65.
12. E. Lazarides and U. Lindberg (1974) *Proc. Natl. Acad. Sci* **71**, 4742–4746.
13. U. Lindberg (1966) *J. Biol. Chem.* **241**, 1246–1248.
14. L. Carlsson, L.-E. Nyström, U. Lindberg, K. K. Kannan, H. Cid-Dresdner, S. Lövgren, and H. Jörnvall (1976) *J. Mol. Biol.* **105**, 353–366.
15. L. Carlsson, L.-E. Nyström, I. Sundkvist, F. Markey, and U. Lindberg (1977) *J. Mol. Biol.*

115, 465–483.

16. V. Mermall, P. L. Post, and M. S. Mooseker (1998) *Science* **279**, 527–533.
17. M.-F. Carlier and D. Pantaloni (1997) *J. Mol. Biol.* **269**, 459–467.
18. M. D. Welch, A. Mallavarapu, J. Rosenblatt, and T. J. Mitchison (1997) *Curr. Opin. Cell Biol.* **9**, 54–61.
19. R. F. Doolittle (1995) *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **349**, 235–240.
20. P. Shterline, J. Clayton, and J. Sparrow (1996) *Protein profile; Actin*, 3rd edition. Acad. Press, London.
21. M. Sanchez, A. Valencia, M. J. Ferrandiz, C. Sander, and M. Vicente (1994) *EMBO J* **13**, 4919–4925.
22. S. Frankel and M. S. Mooseker (1996) *Curr. Opin. Cell Biol.* **8**, 30–37.
23. J. V. Small and J. E. Celis (1978) *Cytobiologie* **16**, 308–325.
24. A.-S. Höglund, R. Karlsson, E. Arro, B.-A. Fredriksson, and U. Lindberg (1980) *J. Muscle Res. Cell Motil.* **1**, 127–146.
25. J. V. Small, G. Rinnerthaler, and H. Hinsen (1981) *Cold Spring Harb. Symp. Quant. Biol.* **46**, 599–611.
26. K. Mellström, A.-S. Höglund, M. Nistér, C.-H. Heldin, B. Westermark, and U. Lindberg (1983) *J. Muscle Res. Cell Motil.* **4**, 589–609.
27. J. V. Small, M. Herzog, and K. Anderson (1995) *J. Cell Biol.* **129**, 1275–1286.
28. G. L. Nicolson, G. Poste, and T.-H. Ji (1977) *Cell Surface Reviews*, vol. **3**., Dynamic aspects of cell surface organization (G. Poste and G. L. Nicolson, eds.), North-Holland, Amsterdam, p. 148.
29. A. Viel and D. Branton (1996) *Curr. Opin. Cell Biol.* **8**, 49–55
30. M. Arpin, M. Algrain, and D. Louvard (1994) *Curr. Opin. Cell Biol.* **6**, 136–141.
31. A. Bretscher, D. Reczek, and M. Berryman (1997) *J. Cell Sci.* **110**, 3011–3018.
32. A. J. Hudspeth (1997) *Curr. Opin. Neurobiol.* **7**, 480–486.
33. M. Kernan and C. Zuker (1995) *Curr. Opin. Neurobiol.* **5**, 443–448.
34. M. Abercrombie, J. E. Heysman, and S. M. Pegrum (1970) *Exp. Cell Res.* **59**, 393–398.
35. G. Albrecht-Buehler and R. M. Lancaster (1976) *J. Cell Biol.* **71**, 370–382.
36. T. M. Svitkina, A. B. Verkhovskiy, K. M. McQuade, and G. G. Borisy (1997) *J. Cell Biol.* **139**, 397–415.
37. A. Y. Chan, S. Raft, M. Bailly, J. B. Wyckoff, J. E. Segall, and J. S. Condeelis (1998) *J. Cell Sci.* **111** 199–211.
38. B. M. Jockusch, P. Bubeck, K. Giehl, M. Kroemker, J. Moschner, M. Rothkegel, M. Rüdiger, K. Schlüter, G. Stanke, and J. Winkler (1995) *Ann. Rev. Cell Dev. Biol.* **11**, 379–416. Review.
39. A. Howe, A. E. Aplin, S. K. Alahari, and R. L. Juliano (1998) *Curr. Opin. Cell Biol.* **10**, 220–231.
40. J. Kendrick-Jones, R. C. Smith, R. Craig, and S. Citi (1987) *J. Mol. Biol.* **198**, 241–252.
41. K. M. Trybus (1991) *Curr. Opin. Cell Biol.* **3**, 105–111.
42. M. A. Lawson and F. R. Maxfield (1995) *Nature* **377**, 75–79.
43. M. Chinkers, J. A. McKanna, and S. Cohen (1979) *J. Cell Biol.* **83**, 260–265
44. J. C. den Hartigh, P. M. van Bergen en Henegouwen, A. J. Verkleij, and J. Boonstra (1992) *J. Cell Biol.* **119**, 349–355
45. M. A. Van der Heyden, P. A. Oude Weernink, B. A. Van Oirschot, P. M. Van Bergen en Henegouwen, J. Boonstra, and G. Rijksen (1997) *Biochim. Biophys. Acta* **1359** (3), 211–221.
46. I. Lassing and U. Lindberg (1985) *Nature* **314**, 472–474.
47. X. D. Ren and M. A. Schwartz (1998) *Curr. Opin. Genet. Dev.* **8**, 63–67.

48. C. L. Carpenter, K. F. Tolia, A. C. Couvillon, and J. H. Hartwig (1997) *Adv. Enzyme Regul.* **37**, 377–390.
49. I. Lassing and U. Lindberg (1990) *FEBS Lett.* **262**, 231–233.
50. L. Carlsson, F. Markey, I. Blikstad, T. Persson, and U. Lindberg (1979) *Proc. Natl. Acad. Sci. USA* **76**, 6376–6380.
51. F. Markey, T. Persson, and U. Lindberg (1981) *Cell* **23**, 145–153.
52. J-H. Hartwig and K. Barkalow (1997) *Curr. Opin. Hematol.* **4**, 351–356.
53. M. J. Berridge (1993) *Nature* **361**, 315–325.
54. M. Halbrügge and U. Walter (1993) In *Protein Kinases in Blood Cell Function* (C.-K. Huang and R. I. Sha"afi, eds.), CRC Press, Boca Raton, FL, pp. 245–298.
55. J. Field, A. Vojtek, R. Ballester, G. Bolger, J. Colicelli, K. Ferguson, J. Gerst, T. Kataoka, T. Michaeli, and S. Powers (1990) *Cell* **61**, 319–327.
56. M. Fedor-Chaikin, R. J. Deschenes, and J. R. Broach (1990) *Cell* **61**, 329–340.
57. C. Haffner, T. Jarchau, M. Reinhard, J. Hoppe, S. M. Lohman, and U. Walter (1995) *EMBO J.* **14**, 19–27.
58. A. Ben-Ze"ev (1997) *Curr. Opin. Cell Biol.* **9**, 99–108.
59. E. Weisberg, M Sattler, D. S. Ewaniuk, and R Salgia (1997) *Crit. Rev. Oncog.* **8**, 343–358.
60. G. B. Cohen, R. Ren, and D. Baltimore (1995) *Cell* **80**, 237–248.
61. R. L. Collins, M Deckert, and A. Altman (1997) *Immunology Today* **18**, 221–224.
62. K. D. Fischer, Y. Y. Kong, H. Nishina, K. Tedford, L. E. Marengere, I. Kozieradzki, T. Sasaki, M. Starr, G. Chan, S. Gardener, M. P. Nghiem, D. Bouchard, M. Barbacid, A. Bernstein, and J. M. Penninger (1998) *Curr. Biol.* **8**, 554–562
63. L. J. Holsinger, I. A. Graef, W. Swat, T. Chi, D. M. Bautista, L. Davidson, R. S. Lewis, F. W. Alt, and G. R. Crabtree (1998) *Curr. Biol.* **8**, 563–572.
64. A. Hall (1998) *Science* **279**, 509–514.
65. J. H. Hartwig, G. M. Bokoch, C. L. Carpenter, P. A. Janmey, L. A. Taylor, A. Toker, and T. Stossel (1995) *Cell*, **82**, 643–653.
66. D. J. G. Mackay, F. Esch, H. Furthmayr, and A. Hall (1997) *J. Cell Biol.* **138**, 927–938.
67. M. Reinhard, M. Halbrügge, U. Scheer, C. Wiegand, B. M. Jockusch, and U. Walter (1992) *EMBO J.* **11**, 2063–2070.
68. C. Haffner, T. Jarchau, M. Reinhard, J. Hoppe, S. M. Lohmann, and U. Walter (1995) *EMBO J.* **14**, 19–27.
69. F. B. Gertler, K. Niebuhr, M. Reinhard, J. Wehland, and P. Soriano (1996) *Cell* **87**, 227–239.
70. M. Reinhard, K. Giehl, K. Abel, C. Haffner, T. Jarchau, V. Hoppe, B. M. Jockusch, and U. Walter (1995) *EMBO J.* **14**, 1583–1589.
71. L. G. Cao, G. G. Babcock, P. A. Rubenstein, and Y. L. Wang (1992) *J. Cell Biol.* **117**, 1023–1029
72. M. C. Beckerle (1997) *BioEssays* **19**, 949–956.
73. S. Hüttelmaier, O. Mayboroda, B. Harbeck, T. Jarchau, B. M. Jockusch, and M. Rüdiger (1998) *Curr. Biol.* **8**, 479–488.
74. I. Lasa and P. Cossart (1996) *Trends Cell Biol.* **6**, 109–114. Review
75. J. M. J. Derry, H. E. Ochs, and U. Francke (1994) *Cell* **78**, 635–644.
76. T. Kirchhausen and F. S. Rosen (1996) *Curr. Biol.* **6**, 676–678.
77. U. Molina, D. M. Kenney, F. D. Rosen, and E. Remold-O'Donnell (1993) *J. Exp. Med.* **151**, 4383–4390.
78. P. Aspenström, U. Lindberg, and A. Hall (1997) *Curr. Biol.* **6**, 70–75.
79. M. Symons, J. M. J. Derry, B. Kartak, S. Jiang, V. Lemahieu, F. McCormick, U. Francke, and

- A. Abo (1996) *Cell* **84**, 723–734.
80. R. Kolluri, K. F. Toliyas, C. L. Carpenter, F. S. Rosen, and T. Kirchhausen (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5615–5618.
  81. S. Banin, O. Truong, D. R. Katz, M. Waterfield, P. M. Brickell, and I. Gout (1996) *Curr. Biol.* **6**, 981–988.
  82. H. Miki, K. Miura, and T. Takenawa (1996) *EMBO J.* **14**, 5326–5335.
  83. H. Miki, T. Sasaki, Y. Takai, and T. Takenawa (1998) *Nature* **391**, 93–96.
  84. S. T. Suzuki, H. Miki, T. Takenawa, and C. Sasakawa (1998) *EMBO J.* **17**, 2767–2776.
  85. N. Ramesh, I. M. Anton, J. H. Hartwig, and R. S. Geha (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14671–14676.
  86. I. M. Anton, W. Lu, B. J. Mayer, N. Ramesh, and R. S. Geha (1998) *J. Biol. Chem.* **273**, 20992–20995.
  87. J. A. Frazier and C. M. Field (1997) *Curr. Biol.* **7**, 414–417.
  88. M. Evangelista, K. Blundell, M. S. Longtine, C. J. Chow, N. Adames, J. R. Pringle, M. Peter, and C. Boone (1997) *Science* **276**, 118–122.
  89. J. Petersen, O. Nielsen, R. Egel, and I. M. Hagan (1998) *J. Cell Biol.* **141**, 1217–1228.
  90. N. Watanabe, P. Madaule, T. Reid, T. Ishizaki, G. Watanabe, A. Kakizuka, Y. Saito, K. Nakao, B. M. Jockusch, and S. Narumiya (1997) *EMBO J.* **16**, 3044–3056.
  91. L. M. Machesky (1998) *Curr. Biol.* **8**, 202–205.
  92. L. Weissbach, J. Settleman, M. F. Kalady, A. J. Snijders, A. E. Murthy, Y.-X. Yan, and A. Bernards (1994) *J. Biol. Chem.* **269**, 20517–20521.
  93. S. Brill, S. Li, C. W. Lyman, D. M. Church, J. J. Wasmuth, L. Weissbach, A. Bernards, and A. J. Snijders (1996) *Mol. Cell Biol.* **16**, 4869–4878.
  94. J. W. Erickson, R. A. Cerione, and M. J. Hart (1997) *J. Biol. Chem.* **272**, 24443–24447.
  95. M. Fukata, S. Kuroda, K. Fujii, T. Nakamura, I. Shoji, Y. Matsuura, K. Kawa, A. Iwamatsu, A. Kikuchi, and K. Kaibuchi (1997) *J. Biol. Chem.* **272**, 29579–29583.
  96. S. Kuroda, M. Fukata, M. Nakagawa, K. Fujii, T. Nakamura, T. Ookubo, I. Izawa, T. Nagase, N. Nomura, H. Tani, I. Shoji, Y. Matsuura, S. Yonehara, and K. Kaibuchi (1998) *Science* **281**, 832–835.
  97. K. Eng, N. I. Naqvi, K. C. Wong, and M. K. Balasubramanian (1998) *Curr. Biol.* **8**, 611–621.
  98. C. E. Turner, J. R. Glenney Jr., and K. Burridge (1990) *J. Cell Biol.* **111**, 1059–1068.
  99. C. E. Turner (1994) *BioEssays* **16**, 47–52.
  100. M. C. Brown, J. A. Perrotta, and C. E. Turner (1996) *J. Cell Biol.* **135**, 1109–1123.
  101. H. Wu and J. T. Parsons (1993) *J. Cell Biol.* **120**, 1417–1426.
  102. X. Zhan, C. Plourde, X. Hu, R. Friesel, and T. Maciag (1994) *J. Biol. Chem.* **269**, 20221–20224.
  103. K. Ozawa, K. Kashiwada, M. Takahashi, and K. Sobue (1995) *Exp. Cell Res.* **221**, 197–204.
  104. S. A. Weed, Y. Du, and J. T. Parsons (1998) *J. Cell Sci.* **111**, 2433–2443.
  105. H. He, T. Watanabe, X. Ahan, C. Huang, E. Schuurin, K. Fukami, T. Takenawa, C. C. Kumar, R. J. Simpson, and H. Maruta (1998) *Mol. Cell Biol.* **18**, 3829–3837.
  106. C. Provenzano, R. Gallo, R. Carbone, P. P. Di Fiore, G. Falcone, L. Castellani, and S. Alema (1998) *Exp. Cell Res.* **242**, 186–200.
  107. K. V. Kishan, G. Scita, W. T. Wong, P. P. Di Fiore, and M. E. Newcomer (1997) *Nat. Struct. Biol.* **4**, 739–743.
  108. A. Vaheri, O. Carpén, L. Heiska, T. S. Helander, J. Jääskeläinen, P. Majander-Nordenswan, M. Sainio, T. Timonen, and O. Turunen (1997) *Curr. Opin. Cell Biol.* **9**, 659–666.
  109. R. J. Shaw, M. Henry, F. Solomon, and T. Jacks (1998) *Mol. Biol. Cell* **9**, 403–419.

110. C. B. Shuster and I. M. Herman (1995) *J. Cell Biol.* **128**, 837–848.
111. J. W. Legg and C. M. Isacke (1998) *Curr. Biol.* **8**, 705–708.
112. V. Niggli, C. Andreoli, C. Roy, and P. Mangeat (1995) *FEBS Lett.* **376**, 172–176.
113. S. Tsukita, S. Yonemura, and S. Tsukita (1997) *Trends Biochem. Sci.* **22**, 53–58.
114. M. Hirao, N. Sato, T. Kondo, S. Yonemura, M. Monden, T. Sasaki, Y. Takai, S. Tsukita, and S. Tsukita (1996) *J. Cell Biol.* **135**, 37–51.
115. S. Yonemura, M. Hirao, Y. Doi, N. Takahashi, T. Kondo, and S. Tsukita (1998) *J. Cell Biol.* **140**, 885–895.
116. K. Burridge and L. Connell (1983) *Cell Motil.* **3**, 405–417.
117. L. Molony, D. McCaslin, J. Abernethy, B. Paschal, and K. Burridge (1987) *J. Biol. Chem.* **262**, 7790–7795.
118. N. C. Collier and K. Wang (1982) *FEBS Lett.* **143**, 205–210.
119. T. O'Halloran, M. C. Beckerle, and K. Burridge (1985) *Nature* **317**, 449–451.
120. M. C. Beckerle, K. Burridge, G. N. De Martino, and D. E. Croall (1987) *Cell* **51**, 569–577.
121. A. Horwitz, K. Duggan, C. Buck, M. C. Beckerle, and K. Burridge (1986) *Nature* **320**, 531–533.
122. D. J. G. Rees, S. E. Ades, S. J. Singer, and R. O. Hynes (1990) *Nature* **347**, 685–689.
123. J. Winkler, H. Lunsdorf, and B. M. Jockusch (1997) *Eur. J. Biochem.* **243**, 430–436.
124. V. Niggli, S. Kaufmann, W. H. Goldmann, T. Weber, and G. Isenberg (1994) *Eur. J. Biochem.* **226**, 951–957.
125. W. H. Goldmann, R. M. Ezzel, E. D. Adamson, V. Niggli, and G. Isenberg (1996) *J. Muscle Res. Cell Motil.* **17**, 1–5.
126. H. Priddle, L. Hemmings, S. Monkley, A. Woods, B. Patel, D. Sutton, G. A. Dunn, D. Zicha, and D. R. Critchley (1998) *J. Cell Biol.* **142**, 1121–1133.
127. B. Geiger (1979) *Cell* **18**, 193–205.
128. B. Geiger, K. T. Tokuyasu, A. H. Dutton, and S. J. Singer (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4127–4131.
129. B. Geiger (1982) *J. Mol. Biol.* **159**, 685–701.
130. K. M. Yamada and B. Geiger (1997) *Curr. Opin. Cell Biol.* **9**, 76–85.
131. S. W. Craig and R. P. Johnson (1996) *Curr. Opin. Cell Biol.* **8**, 74–85.
132. R. P. Johnson and S. W. Craig (1994) *J. Biol. Chem.* **269**, 12611–12619.
133. L. Molony and K. Burridge (1985) *J. Cell. Biochem.* **29**, 31–36.
134. L. M. Milam (1985) *J. Mol. Biol.* **184**, 543–545.
135. J. Week, S. T. Barry, and D. R. Critchley (1996) *Biochem. J.* **314**, 827–832.
136. C. K. Wood, C. E. Turner, P. Jackson, and D. R. Critchley (1994) *J. Cell Sci.* **107**, 709–717.
137. R. B. Hazan, L. Kang, S. Roe, P. I. Borgen, and D. L. Rimm (1997) *J. Biol. Chem.* **272**, 32448–32453.
138. E. E. Weiss, M. Kroemker, A.-H. Rüdiger, B. M. Jockusch, and M. Rüdiger (1998) *J. Cell Biol.* **141**, 755–764.
139. M. Watabe-Uchida, N. Uchida, Y. Imamura, A. Nagafuchi, K. Fujimoto, T. Uemura, S. Vermeulen, F. van Roy, E. D. Adamson, and M. Takeichi (1998) *J. Cell Biol.* **142**, 847–857.
140. J. A. Wilkins, M. A. Risinger, and S. Lin (1986) *J. Cell Biol.* **103**, 1483–1494.
141. S. M. Bockholt, C. A. Otey, J. Glenney Jr., and K. Burridge (1992) *Exp. Cell Res.* **203**, 39–46.
142. S. H. Lo, P. A. Janmey, J. J. Hartwig, and L. B. Chen (1994a) *J. Cell Biol.* **125**, 1067–1075.
143. S. H. Lo, E. Weisberg, and L. B. Chen (1994b) *BioEssays* **16**, 817–823.
144. S. Davis, M. L. Lu, S. H. Lo, S. Lin, J. A. Butler, B. J. Druker, T. M. Roberts, Q. An, and L. B. Chen (1991) *Science* **252**, 712–715.



145. C. Weight, A. Gaertner, A. Wegner, H. Korte, and H. E. Meyer (1992) *J. Mol. Biol.* **227**, 593–595.
146. H. G. Mannherz, M. C. Peitsch, S. Zanotti, R. Paddenberg, and B. Polzar (1995) *Curr. Top. Microbiol. Immunol.* **198**, 161–174.
147. M. Kunitz (1948) *Science* **108**, 19–20.
148. W. Dabrowska, E. J. Cooper, and M. Laskowski (1945) *J. Biol. Chem.* **177**, 991–992.
149. E. J. Cooper, M. J. Trautman, and M. Laskowski (1950) *Proc. Soc. Exptl. Biol. Med.* **73**, 219–222.
150. U. Lindberg (1967) *Biochemistry* **6**, 323–335.
151. M. Segura and U. Lindberg (1984) *J. Biol. Chem.* **259**, 3949–3954.
152. S. E. Hitchcock, L. Carlsson, and U. Lindberg (1975) *Cell Motility Cold Spring Harbor Symp.* pp. 545–559.
153. S. E. Hitchcock, L. Carlsson, and U. Lindberg (1976) *Cell* **7**, 531–542.
154. H. G. Mannherz, H. Brehme, and U. Lamp (1975) *Eur. J. Biochem.* **60**, 109–116.
155. W. Kabsch, H. G. Mannherz, D. Suck, E. F. Pai, and K. C. Holmes (1990) *Nature* **347**, 37–44.
156. C. Oefner and D. Suck (1986) *J. Mol. Biol.* **192**, 605–632.
157. S. A. Weston, A. Lahm, and D. Suck (1992) *J. Mol. Biol.* **226**, 1237–1256.
158. K. C. Holmes, D. Popp, W. Gebhard, and W. Kabsch (1990) *Nature* **347**, 44–49.
159. K. C. Holmes, M. Tirion, D. Popp, M. Lorenz, W. Kabsch, and R. A. Milligan (1993) *Adv Exp. Med. Biol.* **332**, 15–22.
160. I. Rayment, H. M. Holden, M. Whittaker, C. B. Yohn, M. Lorenz, K. C. Holmes, and R. A. Milligan (1993) *Science* **261**, 58–65.
161. C. E. Schutt, M. D. Rozycki, J. C. Myslik, and U. Lindberg (1995) *J. Struct. Biol.* **115**, 186–198.
162. R. Page, U. Lindberg, and C. E. Schutt (1998) *J. Mol. Biol.* **280**, 463–474.
163. K. Schlüter, B. M. Jockusch, and M. Rothkegel (1997) *Biochim. Biophys. Acta* **1359**, 97–109.
164. F. Buss, C. Temm-Grove, S. Henning, and B. M. Jockusch (1992) *Cell Motil. Cytoskeleton* **22**, 51–61.
165. M. Evangelista, K. Blundell, M. S. Longtine, C. J. Chow, N. Adames, J. R. Pringle, M. Peter, and C. Boone (1997) *Science* **276**, 118–122.
166. N. Watanabe, P. Madaule, T. Reid, T. Ishizaki, G. Watanabe, A. Kakizuka, Y. Saito, K. Nakao, B. M. Jockusch, and S. Narumiya (1997) *EMBO J.* **16**, 3044–3056.
167. J. Petersen, O. Nielsen, R. Egel, and I. M. Hagan (1998) *J. Cell Sci.* **111**, 867–876.
168. R. D. Mullins, J. F. Kelleher, J. Xu, and T. D. Pollard (1998) *Mol. Biol. Cell* **9**, 841–852.
169. V. Magdolen, U. Oechsner, G. Muller, and W. Bandlow (1988) *Mol. Cell Biol.* **8**, 5108–5115.
170. M. Haugwitz, A. A. Noegel, J. Karakesisoglou, and M. Schleicher (1994) *Cell* **79**, 303–314.
171. D. A. Kaiser, M. Sato, R. F. Ebert, and T. D. Pollard (1986) *J. Cell Biol.* **102**, 221–226.
172. A. Lambrechts, J. van Damme, M. Goethals, J. Vandekerckhove, and C. Ampe (1995) *Eur. J. Biochem.* **230**, 281–286.
173. T. Fujiwara, K. Tanaka, A. Mino, M. Kikyo, K. Takahashi, K. Shimizu, and Y. Takai (1998) *Mol. Biol. Cell* **9**, 1221–1233.
174. W. Witke, A. H. Sharpe, and D. J. Kwiatkowski (1993) *Mol. Biol. Cell* **4** 149a.
175. C. Björkegren, M. Rozycki, C. E. Schutt, U. Lindberg, and R. Karlsson (1993) *FEBS Lett.* **333**, 123–126.
176. C. E. Schutt, J. C. Myslik, M. D. Rozycki, N. C. Goonesekere, and U. Lindberg (1993) *Nature* **365**, 810–816.
177. E. S. Cedergren-Zeppezauer, N. C. W. Goonesekere, M. D. Rozycki, J. C. Myslik, Z. Dauter,

- U. Lindberg, and C. E. Schutt (1994) *J. Mol. Biol.* **240**, 459–475.
178. K. S. Thorn, H. E. Christensen, R. Shigeta, D. Huddler, L. Shalaby, U. Lindberg, N. H. Chua, and C. E. Schutt (1997) *Structure* **5**, 19–32.
179. N. M. Mahoney, P. A. Janmey, and S. C. Almo (1997) *Nature Struct. Biol.* **4**, 953–960.
180. W. J. Metzler, A. J. Bell, E. Ernst, T. B. Lavoie, and L. Mueller (1994) *J. Biol. Chem.* **269**, 4620–4625.
181. H. Larsson and U. Lindberg (1988) *Biochim. Biophys. Acta* **953**, 95–105.
182. E. Korenbaum, P. Nordberg, C. Björkegren-Sjögren, C. E. Schutt, U. Lindberg, and R. Karlsson (1998) *Biochemistry* **37**, 9274–9283.
183. V. K. Vinson, E. M. De La Cruz, H. N. Higgs, and T. D. Pollard (1998) *Biochemistry* **37**, 10871–10880.
184. D. Pantaloni and M.-F. Carlier (1993) *Cell* **75**, 1007–1014.
185. M.-F. Carlier and D. Pantaloni (1993) *Sem. Cell Biol.* **5**, 183–191.
186. F. Markey, H. Larsson, K. Weber, and U. Lindberg (1982) *Biochim. Biophys. Acta* **704**, 43–51.
187. S. C. Mockrin and E. D. Korn (1980) *Biochemistry* **19**, 5359–5362.
188. P. J. Goldschmidt-Clermont, L. M. Machesky, S. K. Doberstein, and T. D. Pollard (1991) *J. Cell Biol.* **113**, 1081–1089.
189. M. Tanaka and H. Shibata (1985) *Eur. J. Biochem.* **151**, 291–297.
190. U. Lindberg, C. E. Schutt, E. Hellsten, A.-C. Tjäder, and T. Hult (1988) *Biochim. Biophys. Acta* **967**, 391–400.
191. D. Safer and V. T. Nachmias (1994) *BioEssays* **16**, 473–479.
192. V. T. Nachmias (1993) *Curr. Opin. Cell Biol.* **5**, 56–62.
193. M. Czisch, M. Schleicher, S. Horger, W. Voelter, and T. A. Holak (1993) *Eur. J. Biochem.* **218**, 335–344.
194. D. Safer, T. E. Sosnick, and M. Elzinga (1997) *Biochemistry* **36**, 5806–5816.
195. M. F. Carlier, D. Didry, I. Erk, J. Lepault, M. L. Van Troys, J. Vandekerckhove, I. Perelroizen, H. Yin, Y. Doi, and D. Pantaloni (1996) *J. Biol. Chem.* **271**, 9231–9239.
196. H. Q. Sun, K. Kwiatkowska, and H. L. Yin (1996) *J. Biol. Chem.* **271**, 9223–9230.
197. L. Bao, M. Loda, P. A. Janmey, R. Stewart, P. Anand-Apte, and B. R. Zetter (1996) *Nat. Med.* **12**, 1322–1328.
198. M. F. Carlier and D. Pantaloni (1997) *J. Mol. Biol.* **269**, 459–467.
199. J. A. Theriot (1997) *J. Cell. Biol.* **136**, 1165–1168.
200. A. Moon and D. G. Drubin (1995) *Mol. Biol. Cell* **6**, 1423–1431.
201. S. A. Leonard, A. G. Gittis, E. C. Petrella, T. D. Pollard, and E. E. Lattman (1997) *Nat. Struct. Biol.* **4**, 369–373.
202. A. A. Fedorov, P. Lappalainen, E. V. Fedorov, D. G. Drubin, and S. C. Almo (1997) *Nat. Struct. Biol.* **4**, 366–369.
203. H. Hatanaka, K. Ogura, K. Moriyama, S. Ichikawa, I. Yahara, and F. Inagaki (1996) *Cell* **85**, 1047–1055.
204. P. Lappalainen, M. M. Kessels, J. M. T. V. Cope, and D. G. Drubin (1998) *Mol. Cell Biol.* **9**, 1951–1959.
205. S. Arber, F. A. Barbayannis, H. Hanser, C. Schneider, C. A. Stanyon, O. Bernard, and P. Caroni (1998) *Nature* **393**, 805–808.
206. N. Yang, O. Higuchi, K. Ohashi, K. Nagata, A. Wada, K. Kangawa, E. Nishida, and K. Mizuno (1998) *Nature* **393**, 809–812.
207. M.-F. Carlier, V. Laurent, J. Santolini, R. Melki, D. Didry, G. X. Xia, Y. Hong, N.-H. Chua, and D. Pantaloni (1997) *J. Cell Biol.* **136**, 1307–1323.

208. J. Rosenblatt, B. J. Agnew, H. Abe, J. R. Bamburg, and T. J. Mitchison (1997) *J. Cell Biol.* **136**, 1323–1332.
209. B. L. Goode, D. G. Drubin, and P. Lappalainen (1998) *J. Cell Biol.* **142**, 723–733.
210. T. A. Schroer (1994) *J. Cell Biol.* **127**, 1–4.
211. L. M. Machesky (1997) *Curr. Biol.* **7**, 164–167.
212. M. D. Welch, A. Iwamatsu, and T. J. Mitchison (1997) *Nature* **385**, 265–269.
213. L. M. Machesky, E. Reeves, F. Wientjes, F. J. Mattheyse, A. Grogan, N. F. Totty, A. L. Burlingame, J. J. Hsuan, and A. W. Segal (1997) *Biochem. J.* **328**, 105–112.
214. M. D. Welch, A. H. DePace, S. Verma, A. Iwamatsy, and T. J. Mitchison (1997) *J. Cell Biol.* **138**, 375–384.
215. D. Winter, A. V. Podtelejnikov, M. Mann, and R. Li (1997) *Curr. Biol.* **7**, 519–529.
216. R. D. Mullins, J. A. Heuser, and T. D. Pollard (1998) *Proc. Natl. Acad. Sci. USA* **95**, 6181–6186.
217. M. D. Welch, J. Rosenblatt, J. Skoble, D. A. Portnoy, and T. J. Mitchison (1998) *Science* **281**, 105–108.
218. J. F. Casella, S. W. Craig, D. J. Maack, and A. E. Brown (1987) *J. Cell Biol.* **105**, 371–379.
219. P. A. Kuhlman, C. A. Hughes, V. Bennet, and V. M. Fowler (1996) *J. Biol. Chem.* **271**, 7986–7991.
220. H. L. Yin (1987) *Bioessays* **7**, 176–179.
221. T. Azuma, W. Witke, T. P. Stossel, J. H. Hartwig, and D. J. Kwiatkowski (1998) *EMBO J.* **17**, 1362–1370.
222. P. J. McLaughlin, J. T. Gooch, H.-G. Mannherz, and A. G. Weeds (1993) *Nature* **364**, 685–692.
223. L. D. Burtnick, E. K. Koepf, J. Grimes, E. Y. Jones, D. I. Stuart, P. J. McLaughlin, and R. C. Robinson (1997) *Cell* **90**, 661–670.
224. A. McGough, W. Chiu, and M. Way (1998) *Biophys. J.* **74**, 764–772.
225. V. M. Fowler (1996) *Curr. Opin. Cell Biol.* **8**, 86–96.
226. J. Hartwig (1994) Actin-binding proteins 1. Spectrin Superfamily. *Protein Profile.* **1**, 711–778.
227. A. Bretscher and K. Weber (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2321–2325.
228. M. S. Mooseker (1985) *Annu. Rev. Cell Biol.* **1**, 209–241.
229. D. Louvard (1989) *Curr. Opin. Cell Biol.* **1**, 51–57.
230. J. J. Otto (1994) *Curr. Opin. Cell Biol.* **6**, 105–109.
231. M. A. Markus, P. Matsudaira, and G. Wagner (1997) *Protein Sci.* **6**, 1197–1209.
232. C. J. McKnight, P. T. Matsudaira, and P. S. Kim (1997) *Nat. Struct. Biol.* **4**, 180–184.
233. A. Bretscher and K. Weber (1980) *J. Cell Biol.* **86**, 335–340.
234. J. E. Honts, T. S. Sandrock, S. M. Brower, J. L. O'Dell, and A. E. M. Adams (1994) *J. Cell Biol.* **126**, 413–422.
235. S. C. Goldsmith, N. Pokala, W. Shen, A. A. Fedorov, P. Matsudaira, and S. C. Almo (1997) *Nat. Struct. Biol.* **4**, 708–712.
236. D. Hanein, P. Matsudaira, and D. J. DeRosier (1997) *J. Cell Biol.* **139**, 387–396.
237. K. R. Fath and D. R. Burgess 1995, *Curr. Biol.* **5**, 591–593.
238. P. Fucini, C. Renner, C. Herberhold, A. A. Noegel, and T. A. Holak (1997) *Nat Struct Biol* **4**, 223–230.
239. T. W. Tillack, S. L. Marchesi, V. T. Marchesi, and E. Steers Jr. (1970) *Biochim. Biophys. Acta* **200**, 125–131.
240. G. L. Nicolson, V. T. Marchesi, and S. J. Singer (1971) *J. Cell Biol.* **51**, 265–272.
241. V. Bennett (1990) *Curr. Opin. Cell Biol.* **2**, 51–56.

242. J. Pascual, M. Pfuhl, G. Rivas, A. Pastore, and M. Saraste (1996) FEBS Lett **383**, 201–207.
243. P. Fucini, C. Renner, C. Herberhold, A. A. Noegel, and T. A. Holak (1997) Nat. Struct. Biol. **4**, 223–230.
244. J. Pascual, M. Pfuhl, D. Walther, M. Saraste, and M. Nilges (1997) J. Mol. Biol. **273**, 740–751.
245. M. Nilges, M. J. Macias, S. I. O'Donoghue, and H. Oschkinat (1997) J. Mol. Biol. **269**, 408–422.
246. F. J. Blanco, A. R. Ortiz, and L. Serrano (1997) J. Biomol. NMR **9**, 347–357.
247. K. D. Carugo, S. Banuelos, and M. Saraste (1997) Nat. Struct. Biol. **4**, 175–179.
248. S. Ebashi and F. Ebashi (1965) J. Biochem. **58**, 7–13.
249. S. J. Winder (1997) J. Muscle Res. Cell Motil. **18**, 617–629.
250. V. Bennett (1990) Physiol. Rev. **70**, 1029–1065.
251. Y. Srinivasan, L. Elmer, J. Davis, V. Bennett, and K. Angelides (1988) Nature **333**, 177–180.
252. M. Colledge and S. C. Froehner (1998) Curr. Opin. Neurobiol. **8**, 357–363.
253. C. Rosenmund and G. L. Westbrook (1993) Neuron **10**, 805–814.
254. K. M. Harris and S. B. Kater (1994) Annu. Rev. Neurosci. **17**, 341–371.
255. D. W. Allison, V. I. Gelfand, I. Spector, and A. M. Craig (1889) J. Neurosci. **18**, 2423–2436.
256. T. Kim, K. Wu, J. Xu, and I. B. Black (1992) Proc. Natl. Acad. Sci. USA **89**, 11642–11644.
257. P. A. Huber (1997) Int. J. Biochem. Cell Biol. **29**, 1047–1051.
258. J. M. Squire (1997) Curr. Opin. Struct. Biol. **7**, 247–257.
259. J. Trinick (1992) FEBS Lett. **307**, 44–48.
260. R. Littlefield and V. M. Fowler (1998) Annu. Rev. Cell Dev. Biol. **14**, 487–525.

## Actin Polymerization Toxins

Several bacteria release [toxins](#) that modify [actin](#) or proteins controlling its polymerization ([1](#)). A group of *Clostridia spp.* secrete binary toxins, whose most studied member is the C2 toxin released by *Cl. botulinum*. Botulinum C2 toxin is not neurospecific, as are botulinum [neurotoxins](#), but it affects many nonneuronal cells ([2](#)). In the intestinal loop model, it induces an acute inflammatory reaction, characterized by the alteration of endothelia and a large increase in vascular permeability. At the same time, epithelial degeneration, exfoliation, and necrosis are also observed. In cultured cells, C2 causes rounding up, with formation of blebs, followed by cell death.

Two polypeptide chains, I (55 kDa) and II (100 kDa), are needed for cell intoxication. Chain I catalyzes the **ADP-ribosylation** of Arg177 of soluble G-actin, a residue located in an area involved in protein-protein contact within the polymerized form, F-actin. The ADP-ribosylated G-actin binds to the barbed end of F-actin and prevents polymerization, whereas depolymerization at the opposite end is unaffected ([3](#)). This leads *in vitro* and *in vivo* to the disassembly of actin microfilaments, with cell rounding and release of focal adhesion plaques.

Single-chain enzymes that catalyze the ADP-ribosylation of small G proteins are released by several bacteria ([4](#)). They cannot intoxicate cells because they lack the second polypeptide chain, but they are active after cell permeabilization or injection. C3 toxin is released by some strains of *Cl. botulinum*; it catalyzes the specific ADP-ribosylation of Asn-41 of Rho, a small [GTP-binding protein](#) involved in the regulation of actin polymerization. C3 induces the depolarization of F-actin, with

rounding up and binucleation of injected cells. A different type of Rho modification is caused by cytotoxic necrotizing factors released by *E. coli* strains associated with gastroenteritis, urogenital infections, and septicemia (5). These factors cause cell ruffling, stress fiber formation, and multinucleation, by a covalent modification of Rho protein to lock it in its GTP-bound active form (6).

*Clostridia spp.* involved in the induction of diarrhea, associated with pseudomembraneous colitis, release in the intestine enterotoxins of very large dimensions (7). They are termed large clostridial toxins (LCTs) because of their size, which is in the range of 250 to 308 kDa. They are organized as A-B toxins, which are cleaved proteolytically into two polypeptide chains, A and B (see [Toxins](#)). The carboxyl-terminal part includes segments of 20 to 50 residues repeated 14 to 30 times, which are believed to mediate LCT binding to cell surface **receptors**. Such an organization in the absence of an oligomer of type B gives rise to a multivalent type of cell binding, as in the case of cholera toxin. Binding is followed by internalization of LCTs into coated **vesicles** and then into [endosomes](#). By an unknown mechanism, the amino-terminal catalytic domain of the LCTs is released from the rest of the molecule and translocates into the cytosol, where it catalyzes the transfer of a sugar residue from the corresponding UDP derivative to an actin polymerization controlling protein (7, 8). All LCTs induce rounding of different cells in culture, but the effects of the various LCT toxins can be differentiated by staining actin filaments: *Cl. difficile* and *Cl. novyi* LCTs cause a breakdown of the F-actin microfilaments, whereas *Cl. sordellii* LCT induces the formation of filopodialike structures on the cell surface, with some membrane blebbing 7. This is related to the different protein(s) targeted by the LCTs. All of them modify a threonine residue of small [GTP-binding protein\(s\)](#) of the **Ras** superfamily, in such a way that its GTPase activity is unaltered but it can no longer interact with its effector molecule. This provides a further example of the ability of bacterial toxins to “choose” essential cell targets and to modify essential cell functions.

#### Bibliography

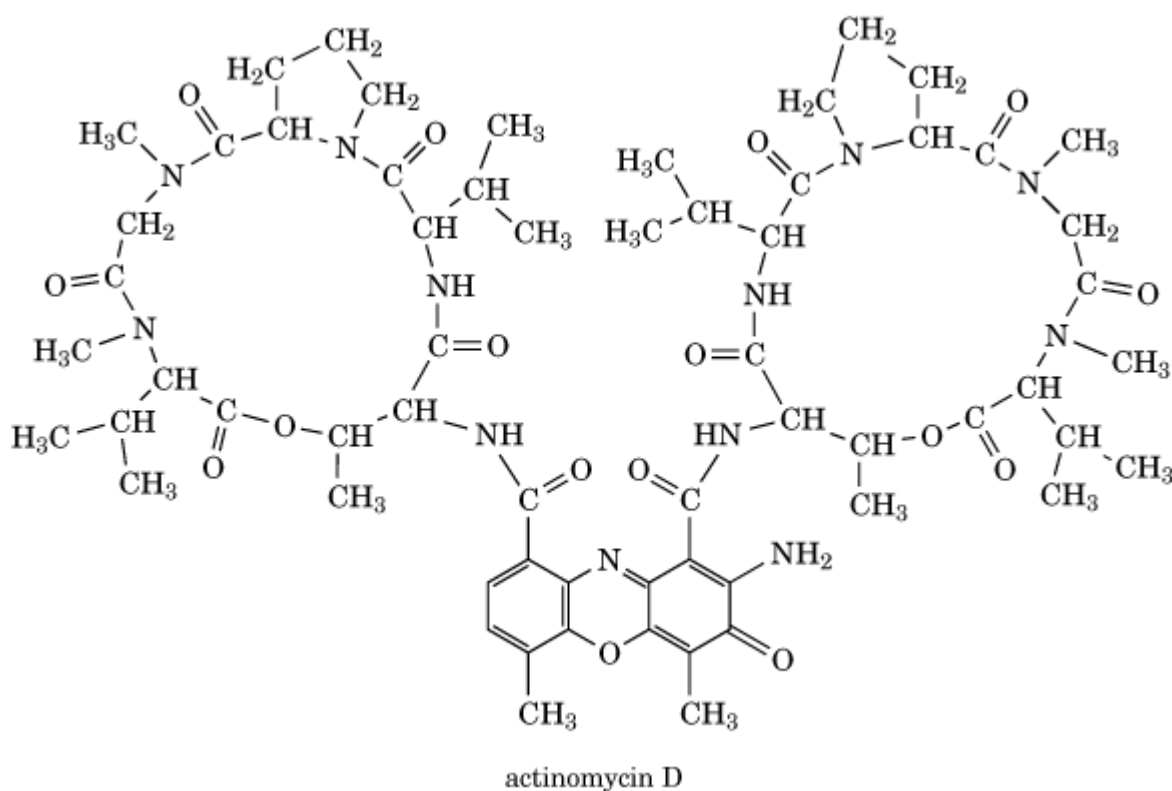
1. K. Aktories (1997) In *Guidebook to Protein Toxins and Their Use in Cell Biology* (R. Rappuoli and C. Montecucco, eds.), Sambrook and Tooze, Oxford University Press, Oxford.
2. L. L. Simpson (1989) *J. Pharmacol. Exp. Ther.* **251**, 1223–1228.
3. M. Wille, I. Just, A. Wegner, and K. Aktories (1992) *J. Biol. Chem.* **267**, 50–55.
4. K. Aktories and A. Wegner (1992) *Curr. Top. Microbiol. Immunol.* **175**, 115–131.
5. A. Caprioli et al. (1984) *Biochem. Biophys. Res. Commun.* **118**, 587–593.
6. E. Oswald et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 3814–3818.
7. C. von Eichel-Streiber, P. Boquet, M. Sauerborn, and M. Thelestam (1996) *Trends Microbiol.* **4**, 375–382.
8. I. Just et al. (1995) *Nature* **375**, 500–503.

#### Actinomycin D

Actinomycins are chromopeptide antibiotics produced by various species of streptomycetes (1). They differ only in their content of certain [amino acids](#) in their peptide rings (Fig. 1). Actinomycin D, which is active against several forms of cancer, has the distinction of being the first antibiotic discovered (in 1943) to possess useful antitumor activity. It is still used in the clinic for that purpose. It occupies a unique place in the origins of molecular biology for quite a different reason, however, namely its astounding specificity as an inhibitor of DNA-dependent RNA synthesis, which enabled it to play a vital part in the discovery of [messenger RNA](#) and helped uncover the early evidence that **transcriptional** regulation of gene expression is fundamental to cell growth and development (2).

The result of its near-absolute specificity for binding tightly to double helical B-form DNA is a total blockade of transcription by **RNA polymerase**, leading to cessation of synthesis of all forms of RNA, stable as well as unstable. Elongation of growing RNA chains is promptly terminated; there is no particular effect on the process of initiation of transcription. Thus actinomycin was used in the 1960s to measure the kinetics of breakdown of pulse-labeled RNA, thought to reflect at least partly the turnover of mRNA, and to establish the requirement for new RNA transcription in such processes as induced enzyme synthesis, mechanisms of steroid [hormone](#) and peptide hormone action, and early embryonic [development](#). Replication of DNA **viruses** is usually strongly inhibited by actinomycin, but many RNA viruses are completely insensitive to the antibiotic; this is clear evidence that their life cycles do not require the participation of DNA at any point. Because actinomycin acts essentially identically to block transcription in all cells, once it has gained access to the nucleus, it is equally applicable to the study of gene activity in eukaryotes as well as prokaryotes (it is generally ineffective against [gram-negative bacteria](#) simply because of a permeability problem, for if the walls of such cells are digested away, or otherwise weakened, the [protoplast](#) is revealed as fully sensitive to the antibiotic).

**Figure 1.** The structures of the actinomycin antibiotics.



Actinomycin does not bind to RNA or to single-stranded DNA, and its affinity for double-helical DNA is related to the latter's content of guanine nucleotides; synthetic polynucleotides composed entirely of A · T base pairs do not interact with the antibiotic at all. It took some time to establish the nature of the actinomycin–DNA complex, and for many years the matter was controversial. Eventually the issue was resolved in favor of intercalation as a result of careful structure–activity comparisons (1), binding measurements with a wide variety of synthetic as well as naturally occurring DNAs (3, 4), hydrodynamic experiments with circular DNA (5), [X-ray crystallography](#) (6), and high-resolution nuclear magnetic resonance (**NMR**) (7). A landmark was the determination of the structure of a crystalline 1:2 actinomycin:deoxyguanosine complex, which revealed for the first time that the antibiotic has an axis of near-perfect twofold rotational symmetry, which allows it

to react with two guanine nucleosides in a symmetry-related fashion (6). This observation immediately suggested plausible intercalation of the phenoxazinone chromophore between two G · C base pairs of DNA at a rotationally symmetrical site centered around a 5'-GpC-3' step. The strong preference of actinomycin for such sites in DNA was elegantly confirmed by [Footprinting](#) experiments a decade later (8, 9). With its tricyclic aromatic chromophore firmly embedded between the base pairs, the cyclic pentapeptide rings of the antibiotic are left neatly filling the minor groove of the distorted B-form helix, where they form numerous additional **van der Waals** contacts that help to stabilize the complex; positioned this way, their intrinsic right-hand twisted disposition with respect to the intercalated chromophore makes perfect sense, and the whole antibiotic molecule occupies a site covering about six base pairs in the minor groove.

Detailed [kinetic](#) studies have revealed that both the association and dissociation reactions are complicated and require several rate constants to fit the data (3, 10). The slowest processes are characterized by time constants in the range of minutes, conspicuously slower than those measured for most other DNA-binding drugs, and there is a correlation between the slowest rate constant for the dissociation reaction and the efficiency of inhibition of chain elongation by RNA polymerase. This has led to the notion that reversibly bound actinomycin molecules serve as relatively long-lived blocks to the progression of the transcribing enzyme along its template, providing a partial explanation for the extreme effectiveness of actinomycin as an inhibitor of RNA synthesis. What is not so obvious is why the functioning of DNA as a template in replication should remain unaffected even at much higher concentrations of actinomycin (2).

The critical determinant that enables the antibiotic to recognize its preferred binding sites in DNA seems to be the 2-amino group of the guanine nucleotide, as is also true for numerous other ligands (11). If the 2-amino group is removed from guanines (leaving inosine–cytosine base pairs) and transferred to adenines (forming 2,6-diaminopurine–thymine base pairs), the sites to which actinomycin binds are relocated accordingly (11, 12). The molecular recognition process includes the formation of [hydrogen bonds](#) between the purine 2-amino groups and the carbonyl substituents of the L-**threonine** residues in the peptide rings of the antibiotic; there are also hydrogen bonds from the same threonine residues to the N(3) atoms of the purine nucleotides at the binding site (6, 12). Some of the kinetic complexity of the association/dissociation reactions no doubt reflects a “shuffling” process, whereby actinomycin initially interacts with a variety of potential binding sites on the DNA lattice, then migrates one-dimensionally to locate its preferred binding sites marked by a pair of purine 2-amino groups suitably disposed in the minor groove of the double helix (13).

## Bibliography

1. J. Meienhofer and E. Atherton (1977) In *Structure–Activity Relationships Among the Semisynthetic Antibiotics* (D. Perlman, ed.), Academic Press, New York, pp. 427–529.
2. E. Reich and I. H. Goldberg (1964) *Prog. Nucleic Acid Res. Mol. Biol.* **3**, 183–234.
3. W. Müller and D. M. Crothers (1968) *J. Mol. Biol.* **35**, 251–290.
4. R. D. Wells and J. E. Larson (1970) *J. Mol. Biol.* **49**, 319–342.
5. M. J. Waring (1970) *J. Mol. Biol.* **54**, 247–279.
6. H. M. Sobell, S. C. Jain, T. D. Sakore, and C. E. Nordman (1971) *Nature New Biol.* **231**, 200–205.
7. D. J. Patel (1974) *Biochemistry* **13**, 2396–2402.
8. M. J. Lane, J. C. Dabrowiak, and J. N. Vournakis (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3260–3264.
9. M. W. Van Dyke, R. P. Hertzberg, and P. B. Dervan (1983) *Proc. Natl. Acad. Sci. USA* **79**, 5470–5474.
10. R. Bittman and L. Blau (1975) *Biochemistry* **14**, 2138–2145.
11. C. Bailly and M. J. Waring (1995) *Nucleic Acids Res.* **23**, 885–892.
12. S. Jennewein and M. J. Waring (1997) *Nucleic Acids Res.* **25**, 1502–1509.

13. K. R. Fox and M. J. Waring (1985) *Nucleic Acids Res.* **14**, 2001–2014.

### Suggestions for Further Reading

14. E. F. Gale, E. Cundliffe, P. E. Reynolds, M. H. Richmond, and M. J. Waring (1981) *The Molecular Basis of Antibiotic Action*, 2nd ed., Wiley, London, pp. 314–333.

15. H. M. Sobell (1973) The stereochemistry of actinomycin binding to DNA and its implications in molecular biology. *Prog. Nucleic Acid Res. Mol. Biol.* **13**, 153–190.

## Activation Energy

The [kinetic](#) rates of chemical and enzyme-catalyzed reactions depend on temperature. The relationship for this dependence is known as the Arrhenius law:

$$k = Ae^{E/RT} \quad (1)$$

where  $A$  and  $E$  are constants,  $R$  is the gas constant, and  $T$  is the absolute temperature in Kelvin units. The activation energy  $E$  is the height of the energy barrier that the reaction must exceed to pass from reactants to products. It is usually expressed in kilocalories per mole or in joules per mole. Equation [1](#) can be expressed as:

$$\log k = \log A - E/(2.303RT) \quad (2)$$

Therefore the magnitude of the activation energy can be obtained from the slope of a plot of the log of a rate constant for a reaction as a function of  $1/T$ . Many chemical and enzyme-catalyzed reactions increase the rate of reaction by two- to threefold for each 10°C increase in temperature. Although this relationship is useful in explaining the temperature dependence of reactions, it does not explain the rate in the thermodynamic terms of **enthalpy**  $H$  **entropy**  $S$ , or **free energy**  $G$ . This analysis comes from [transition state](#) theory. The Arrhenius activation energy does correspond to the standard enthalpy of the reaction in the van't Hoff equation ([1](#)).

### Bibliography

1. A. Cornish-Bowden (1979) *Fundamentals of Enzyme Kinetics*, Butterworths, London, Boston, pp. 9–12.

### Suggestions for Further Reading

2. D. Piszkiwicz (1977) *Kinetics of Chemical and Enzyme-Catalyzed Reactions*, Oxford University Press, New York, pp. 27–52.

## Active Site

The folding of a [polypeptide chain](#) that produces the final [protein structure](#) of an [enzyme](#) also leads to formation of the active site. From [X-ray crystallography](#) studies, it is apparent that the active site



of an enzyme is a groove, cleft, or pocket that has access to the solvent and forms only a small part of the total solvent-**accessible surface** of the protein. The relatively large sizes of enzymes are undoubtedly due to the need to obtain, at the active site, the correct spatial relationships of the amino acid residues that are involved in binding of substrates, [catalysis](#), and the release of products, as well as for any conformational changes associated with these steps. The binding of substrates or inhibitors at the active site pocket of an enzyme involves matching up of the [nonpolar](#) groups of the substrate with the nonpolar side-chains of amino acid residues, **hydrogen bonding** between the [polar](#) appropriate groups on the substrate with the backbone NH and CO groups within the active site, and even **salt bridge** formation. For substrates, these initial interactions are followed by the conformational changes that lead to the formation of the transition-state complex (see [Transition State Analogue](#)) and the chemistry for catalyzing the reaction brought about by reactive groups with the correct alignments. These may be the acidic, basic, and nucleophilic groups of the protein component of the enzyme, or the electrophilic groups of a prosthetic group (see [Coenzyme, Cofactor](#)).

### Active Site-Directed Irreversible Inhibitors

Active site-directed irreversible inhibitors of [enzymes](#) are also known as *active site-directed inactivating reagents*, *affinity labels*, and *photoaffinity labels*. They combine the features of a substrate, or substrate analogue, with those of a group-specific reagent, as in [affinity labeling](#), and have been used to determine the amino acid residues that are present in the [active site](#) and involved in enzymic [catalysis](#). They are capable of binding specifically and reversibly at the active site of an enzyme and then causing inactivation through time-dependent covalent modification of an adjacent amino acid residue. The functional group of the inhibitor is usually an electrophile that can interact with an appropriately positioned nucleophile of the enzyme, to generate a covalent bond between them. The electrophilic groups include epoxides and  $\alpha$ -haloketones. Since these compounds are reactive in solution, they could also cause some nonspecific enzyme modifications.

The action of an active-site directed irreversible inhibitor I can be illustrated by



where EI represents a Michaelis complex (see [Michaelis–Menten Kinetics](#)) and  $E - I$  denotes the covalently cross-linked form of the complex. On the basis of this formulation, it would be expected that at the early stages of the interaction, I would behave as a **competitive inhibitor** with respect to the substrate. Examples of the action of an active-site directed irreversible inhibitors are the acetylation of amino acid residues at the active site of prostaglandin synthase by aspirin (acetyl salicylate) (1) and the alkylation by L-TPCK (*N*-tosylphenylalanine chloromethyl ketone) of a **histidine** residue at the active site of  $\alpha$ -**chymotrypsin** (2).

Photoaffinity labels, such as diazoketones and aryl azides, introduce a greater degree of specificity to the modification of amino acid residues, as they are not reactive in solution. It is only after the reversible interaction of the affinity label at the active site of an enzyme, and exposure of the resulting complex to light of the correct wavelength, that a highly reactive group is formed. Such treatment with diazoketones and aryl azides leads to the formation of carbenes and azines that are extremely reactive and can add across O—H bonds of unionized carboxyl groups or unsaturated carbon–hydrogen bonds (3).

Bibliography

1. G. J. Roth, N. Stanford, J. W. Jacobs, and P. W. Majerus (1977) *Biochemistry* **16**, 4244–4248.
2. C. Walsh (1979) *Enzymatic Reaction Mechanisms*, W. H. Freeman and Company, San Francisco, Calif., p. 86.
3. V. Chowdry and F. H. Westheimer (1979) *Ann. Rev. Biochem.* **48**, 293–325.

### Active-Site Titrants

For determination of values for the [kinetic](#) rate constants and **ligand-binding** stoichiometries associated with an **enzyme**-catalyzed reaction, it is necessary to know the concentration of functional [active sites](#) of the enzyme. The latter value cannot be calculated simply from the total concentration of protein and molecular weight of the enzyme; even when the enzyme preparation has been shown to be homogeneous by a variety of techniques, it is possible that inactive enzyme is present. The inactivity could be due to the inability of the enzyme to bind substrate or perform the chemistry of the reaction, or both. Therefore, before undertaking detailed kinetic investigations on any enzyme, it is important to determine the concentrations of the active sites that can both bind the ligands and are active.

The concentration of binding sites can be obtained from data for the reversible inhibition of an enzyme by a tight-binding substrate analogue that gives rise to [competitive inhibition](#). Tight-binding inhibition occurs under conditions where the total inhibitor concentration  $I_t$  is comparable to the total enzyme concentration  $E_t$ , and a substantial fraction of the inhibitor is bound, so that allowance has to be made for the reduction of free inhibitor concentration as a result of formation of the enzyme-inhibitor complex (1). The variation of the steady-state velocity as a function of the concentrations of  $I_t$  and  $E_t$  is described by Equation 1:

$$v = \frac{kA}{2(K_a + A)} \left[ \{(K'_i + I_t - E_t)^2 + 4K'_i E_t\}^{1/2} - (K'_i + I_t - E_t) \right] \quad (1)$$

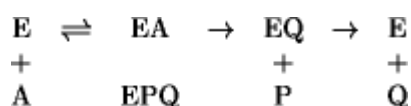
where  $k$  denotes the maximum rate of product formation in terms of moles per mole of enzyme per second,  $K_a$  is the Michaelis constant for the substrate present at concentration  $A$ , and  $K'_i$  represents an apparent inhibition constant whose relationship to the true inhibition constant  $K_i$  for the interaction of  $I$  with  $E$  is given by

$$K'_i = K_i(1 + A/K_a)$$

Equation 1 also applies to multisubstrate reactions, provided that the nonvaried substrates are present at saturating concentrations. When no assumptions are made about the purity of the enzyme,  $E_t$  of Equation 1 is replaced by  $aEx$ , where  $a$  represents the degree of purity of the enzyme and  $Ex$  denotes the total protein concentration. Fitting to the modified form of Equation 1 of steady-state velocity data obtained at varying concentrations of  $I_t$  and  $E_t$  would yield values for  $K'_i$  as well as  $a$ , which would yield a measure of the concentration of binding sites (2, 3).

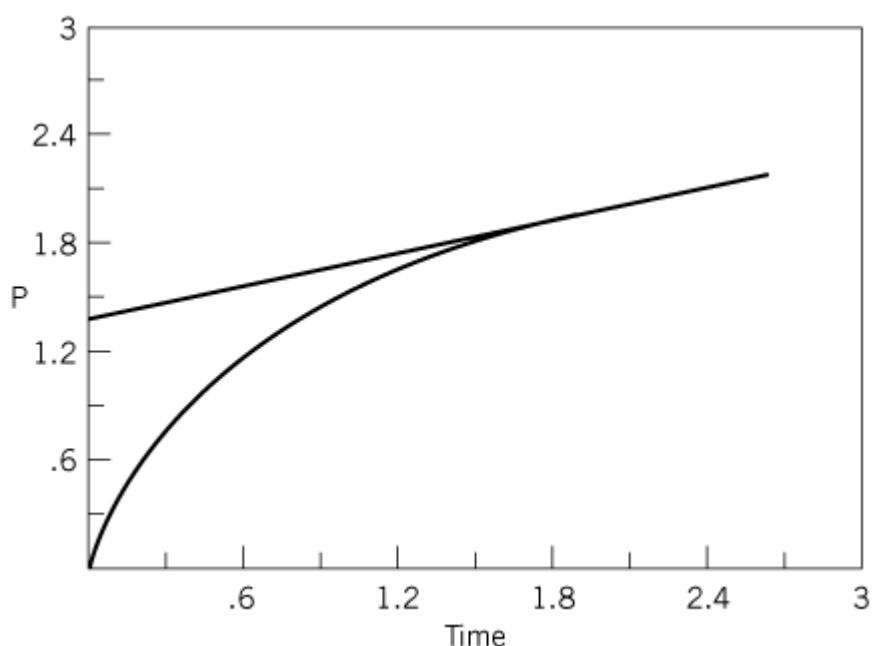
The technique of active-site titration is used to determine the concentration of catalytically active enzyme (4, 5). It requires that on mixing enzyme and the titrating substrate, there is an initial rapid burst of product formation because of the accumulation of an enzyme-bound intermediate whose rate

of breakdown is much slower than its rate of formation. The procedure can be illustrated by reference to the Uni-Bi reaction:



for which P is released before Q, and at a very much faster rate. Specific measurement of the release of P as a function of time would yield a plot (Fig. 1) that consists of curved (burst) and linear sections. The intersection of the extrapolated linear curve with the vertical ordinate would give the concentration of P that is equal to the concentration of functional active sites.

**Figure 1.** Active-site titration curve showing burst formation of a covalent enzyme-intermediate complex.



Active-site titrations may be performed with chromogenic or radioactive substrates, and rapid reaction techniques are often required. This is not the case, however, for determination of the active-site concentration of *hexokinase* (6). This can be done by incubating the enzyme with  $^{14}\text{C}$ -glucose and ATP complexed with  $\text{Cr}^{3+}$ , rather than the usual  $\text{Mg}^{2+}$ , and separating the enzyme-CrADP-glucose-6-phosphate complex on a Sepharose column. In this case, the enzyme undergoes only a single turnover, and the product complex has a half-life of over 10 min, so the amount of  $^{14}\text{C}$  associated with the enzyme gives a measure of the active-site concentration (6).

Comparison of the results obtained from tight-binding inhibition and active-site titration studies could yield information about the presence in an enzyme preparation of substantial proportion of sites that can bind substrate, but cannot catalyze the reaction.

#### Bibliography

1. J. W. Williams and J. F. Morrison (1979) *Meth. Enzymol.* **63**, 437–467.
2. J. F. Morrison and S. R. Stone (1985) *Comments Mol. Cell. Biophys.* **2**, 347–368.
3. M. J. Sculley and J. F. Morrison (1986) *Biochem. Biophys. Acta* **874**, 44–53.

4. A. Fersht (1977) *Enzyme Structure and Mechanism*, W. H. Freeman and Company, San Francisco, Calif., pp. 122–126.
5. F. J. Kezdy and E. T. Kaiser (1970) *Meth. Enzymol.* **19**, 3–20.
6. K. D. Danenberg and W. W. Cleland (1975) *Biochemistry* **14**, 28–39.

## Adenovirus

The adenoviruses constitute a large group of DNA **viruses**, the *Adenoviridae* family. Genera include *Mastadenovirus* (human, simian, bovine, equine, porcine, ovine, canine, and opossum) and *Aviadenovirus* (birds). Adenoviruses are nonenveloped icosahedral particles of 70 to 100 nm in diameter, with 252 capsomeres, of which 240 are hexons and 12 are penton bases. A projecting fiber attaches to each penton base to form a penton.

Human adenoviruses (Ads) include 47 serotypes, which can be classified into six subgroups, A to F, based on their ability for red-blood-cell agglutination and oncogenicity in rodents, plus their DNA [homology](#). Ads cause respiratory tract infections, conjunctivitis, hemorrhagic cystitis, and gastroenteritis. Highly oncogenic group A (eg, Ad12 and Ad18), weakly oncogenic group B (eg, Ad3 and Ad7), and Ad9 of group D can induce tumors in rodents, but no conclusive evidence has been reported linking adenoviruses with malignant diseases in the human. All Ads that have been tested can transform cultured rodent cells. No infectious virus is present, but viral transforming genes (E1A and E1B) are present and expressed in tumors and in transformed cells induced by adenoviruses.

The viral genome consists of a single linear molecule of double-stranded DNA [MW  $\sim 23 \times 10^6$ , size  $\sim 36$  kilobase pairs (kbp)], varying somewhat with the type. Human adenovirus type 2 (Ad2), the first to be sequenced completely, has a total of 35,937 bp. The Ad genome has [inverted terminal repeats](#) of 103 bp to 163 bp, depending on the type, and two identical [replication origins](#), one in each terminal repeat. A 55-kDa terminal protein (TP) covalently bound to each 5' end of the viral DNA molecule serves as a primer for protein-primed viral [DNA replication](#). The viral genome carries five early [transcription](#) units (E1A, E1B, E2, E3, E4), two delayed early transcription units (IX and IVa2), and one late transcription unit (major late), all of which are transcribed by **RNA polymerase II**. The viral genome also carries one or two (depending on the type) VA genes transcribed by RNA polymerase III. The *E1A*, *E1B*, *IX*, *major late*, *VA*, and *E3* are on the rightward reading DNA strands (r strand), and others are on the leftward reading DNA strand (l strand).

The E1A is the first transcription unit to be expressed shortly after infection, using cellular [transcription factors](#), encoding two major mRNAs of 12S and 13S. The encoded E1A proteins are required for productive viral replication and play a key role in cell transformation. E1A proteins transactivate early and late viral genes and cellular transcription units, and they bind cellular proteins, including p300, a family member of CBP (CREB binding protein)/p300, plus the RB family members, RB (**retinoblastoma** susceptibility gene product, tumor suppressor), p107, and p130. Binding to these cellular proteins is required for cell transformation. Binding of E1A to RB disrupts the RB · E2F complex to activate a cellular transcription factor, E2F, resulting in the activation of E2F-dependent **cell-cycle-related** genes.

Two major proteins of 19 kDa and 55 kDa, encoded by 13S and 22S [messenger RNAs](#), respectively, are generated from the E1B region. The E1B proteins are required for efficient viral growth and

cooperate with E1A products to transform rodent cells, preventing E1A-induced [apoptosis](#). The E1B 19-kDa protein has a functional similarity to the cellular anti-apoptotic Bcl-2, and the E1B 55-kDa protein interacts with **tumor suppressor** and apoptosis-related [p53](#).

Region E2 encodes proteins involved in viral replication, including the viral terminal protein precursor, pTP (coded by *E2B*), viral DNA polymerase (coded by *E2B*), and the DNA-binding protein (coded by *E2A*). None of the *E3* proteins are required for productive infections of adenoviruses in cultured cells, but are important for *in vivo* infections in humans, suppressing host defense mechanisms by [cytotoxic T lymphocytes](#) or **tumor necrosis factor- $\alpha$** . The E4 proteins are involved in transcriptional regulation, preferential transport of viral mRNA, and efficient viral DNA replication. Other functions of the E4 gene, such as transformation with E1A, tumor-suppressor p53 binding, and apoptosis-inhibiting activities, have been recently reported. The Ad9 E4 gene can oncogenically transform rat cells in the absence of the E1A and E1B genes and is required for mammary tumorigenesis.

The transcription of adenovirus late genes is triggered by the onset of viral DNA replication, yielding at least 18 distinct mRNAs by differential [poly\(A\)](#) site utilization and [alternative splicing](#). Viral DNA replication activates a major late **promoter** (MLP) located at 16.4 map units of the viral genome, and one large primary transcript is generated, terminating at 99 map units at the right end of the genome. This transcript is processed to five families of late mRNAs, L1 to L5, based on the use of common poly(A) addition sites. The late mRNAs encode structural polypeptides of the viral particle and proteins involved in polypeptide processing, capsomere assembly, and packaging of the viral genomic DNA.

Adenoviruses are being used as **vectors** for delivery of therapeutic genes to target organs *in vivo*. Recombinant adenoviruses can be constructed by replacing the E1A and E1B genes with foreign ones. The E3 region also can be deleted without significant changes in virus growth. Recombinant adenoviruses can propagate efficiently in 293 cells, which complement the defect. The advantages are that the growth of recombinant viruses does not require host cell division, and a high titer of recombinant viruses can be easily obtained.

#### Suggestions for Further Reading

T. Shenk (1996) "*Adenoviridae: The Viruses and Their Replication*". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2111–2148.

M. S. Horwitz (1996) "Adenoviruses". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2149–2171.

J. Tooze (1981) *DNA Tumor Viruses*, 2nd ed., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 943–1054.

#### Adenylate Charge

The adenylate charge, or adenylate energy charge, is a linear measure of the energy stored in the adenylate system ( $[ATP] + [ADP] + [AMP]$ ) in a living cell. It is analogous to the charge of a storage battery. Its value is defined by the expression:

$$\text{Adenylate charge} = \frac{[ATP] + 0.5[ADP]}{[ATP] + [ADP] + [AMP]}$$

The status of the adenylate system is best described by the adenylate charge and appears to be the most ubiquitous regulatory signal in metabolism; it affects the rates of nearly all metabolic conversions and the partitioning of metabolically available substances between oxidation, synthesis of cell substance, and production of storage compounds. In living organisms, the same controls that regulate metabolic partitioning also maintain the value of the adenylate charge near 0.9.

## 1. Background

Organisms are exquisitely regulated material- and energy-transducing systems. Material from the environment must be converted to the many compounds of which the organism is made, and energy, usually obtained either by absorption of sunlight or from oxidation of foodstuffs, must be converted for use in biosynthesis, movement, and membrane activities. Partitioning of material and energy between those functions must be adjusted continuously to meet the changing needs of the cell or organism. Transduction of energy through the adenylate system is an integral part of each of those functions.

Any system that is not at chemical and physical equilibrium can, in principle, supply energy. The two types of nonequilibrium situations that are mainly used in energy storage and transduction by organisms are (1) difference in electric charge or chemical potential across a membrane (see [Membrane Potentials](#)) and (2) a ratio of [ATP] to [ADP] that is far from equilibrium. All cells contain enzyme systems that catalyze transfer of energy between the two.

The energy status of a chemical storage system may be defined by either of two parameters: the molar Gibbs **free energy** change of the relevant chemical reaction or the mole fraction of the higher-energy state of the system. The two are computationally interconvertible but are not linearly related. The molar free energy change is a function of the ratio of activities of the products and reactants of the reaction; the mole fraction is a linear measure of the extent of reaction. In the familiar case of a lead storage battery, the free energy is measured by the voltage and the mole fraction by the charge (measured by a hydrometer in this case).

If a generic energy-transducing reaction is indicated by the type of reaction  $A \rightleftharpoons B$ , the molar free energy change is given by

$$\Delta G = \Delta G^\circ + RT \ln \left( \frac{[B]}{[A]} \right) \quad (1)$$

and the fractional charge (mole fraction of the more energetic state) by

$$\text{Fractional charge} = \frac{[B]}{[A] + [B]} \quad (2)$$

In the case of a lead storage battery,  $A$  is  $(2\text{Pb}(\text{SO}_4) + 2\text{H}_2\text{O})$  and  $B$  is  $(\text{Pb} + \text{PbO}_2 + 2\text{H}_2\text{SO}_4)$ . For the metabolic adenylate system,  $A$  is  $(\text{ADP} + \text{P}_i)$  and  $B$  is  $(\text{ATP})$ .

A system is fully discharged when only  $A$  is present and fully charged when it consists of  $B$  alone. For intermediate states, the stoichiometrically linear charge function indicates how much work is available as a fraction of that of a totally charged system. [This metabolic use of the term “work” is slightly nonstandard. In thermodynamic terms, the maximum work available is the integrated product of the charge multiplied by the changing value of the Gibbs free energy as the system goes to equilibrium. But the energy of the adenylate system is used stoichiometrically—for example, to affect chemical change or mechanical movement—and charge is therefore proportional to available work, defined in terms of biological effect.] It is because of its stoichiometric nature that the charge function is more relevant in most contexts than the molar free energy change for both storage

batteries and the metabolic adenylate system.

If the adenylate system merely alternated between ADP and ATP, the charge function would be the simple ATP mole fraction,  $[ATP]/([ATP] + [ADP])$ . Some enzymes, however, couple the use of ATP to synthetic reactions by converting ATP to AMP and pyrophosphate, and so the concentration of AMP must also be taken into account.



The pyrophosphate is rapidly hydrolyzed enzymically by pyrophosphatase, and this removal of product pulls the reaction forward (increases the numerical value of the negative free energy change).

Accumulation of AMP as a consequence of such reactions would deplete the cell's stores of ATP and ADP and so would be rapidly lethal. That outcome is prevented by the action of adenylate kinase, which catalyzes the phosphorylation of AMP:



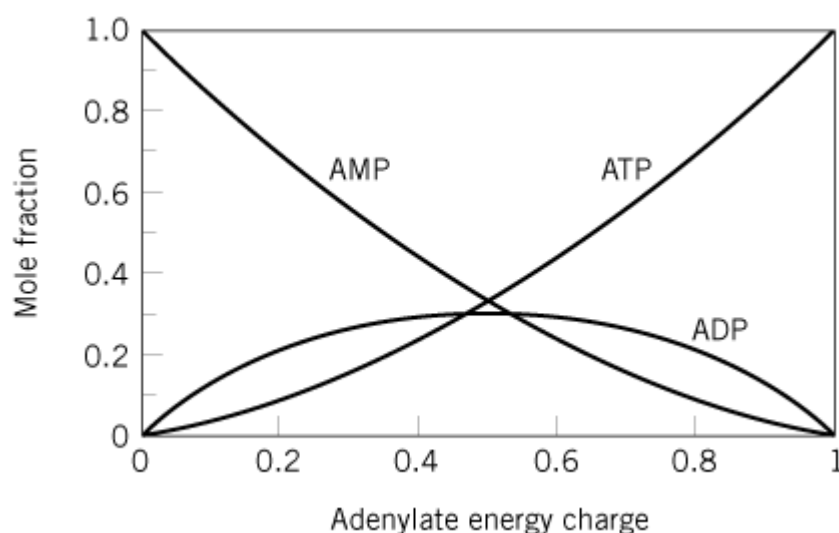
The sum of reactions (2) and (3) is conversion of two molecules of ATP to two of ADP. In such cases, a single molecule of ATP, in effect, supplies twice as much metabolic energy to a metabolic reaction as is obtained from action of an ordinary kinase that converts ATP to ADP, but a second molecule of ATP is required to make good the energy balance.

Because of the participation of AMP in metabolic energy transduction, the simple mole fraction of ATP is not an adequate measure of the energy status of the adenylate system. Reaction (3) shows that two molecules of ADP are energetically equivalent to one of ATP. Thus, the linear charge function for the adenylate system (the effective mole fraction of ATP) is seen to be

$$\text{Adenylate charge} = \frac{[ATP] + 0.5[ADP]}{[ATP] + [ADP] + [AMP]} \quad (5)$$

The adenylate charge is the linear measure of the metabolic work available. A system containing only ATP is fully charged, with an adenylate charge value of 1.0, and one containing only AMP is fully discharged, with an adenylate charge of 0. The charge value would be 0.5 if only ADP were present. If reaction (3) catalyzed by adenylate kinase is at equilibrium, the concentrations of ATP, ADP, and AMP are fixed for any particular value of the adenylate charge (Figure 1).

**Figure 1.** Mole fractions of the components of the adenylate nucleotide system as a function of the adenylate charge. The reaction catalyzed by adenylate kinase is assumed to be near equilibrium, with an apparent equilibrium constant (irrespective of differential ionization and magnesium binding) of 0.8.



## 2. Regeneration and utilization of ATP

In aerobic organisms, ATP is regenerated mostly by oxidation of foods or storage compounds to carbon dioxide (see [ATP Synthase](#)). In anaerobic organisms, substrates undergo other energy-yielding conversions—for example, fermentation of glucose to ethanol and carbon dioxide—rather than oxidation. Carbohydrates are converted to pyruvate by way of the glycolytic pathway. The pyruvate is oxidized to a derivative of acetic acid, acetyl coenzyme A, which is oxidized to carbon dioxide in the reactions of the citrate cycle, or Krebs cycle. The electrons lost in those oxidations are passed on to oxygen by mediation of a series of membrane-associated enzymes. Those electron transfers are coupled to the conversion of ADP to ATP, thus supplying metabolically available energy to the adenylate system. The routes by which other classes of foods are metabolized feed into this central pathway. Fats are broken down to produce [acetyl coenzyme A](#), which joins the carbohydrate pathway at that point. [Amino acids](#) derived from **protein degradation** are metabolized by individual pathways to produce various intermediates of glycolysis or the citrate cycle. Thus, the same central pathways are taken in the utilization of all foods.

Glycolysis and the citrate cycle are also centrally involved in biosynthesis. The biosynthetic pathways leading to the many components of a cell all begin with intermediates of glycolysis or the citrate cycle. That is, 10 or 12 intermediates of these degradative pathways are also the starting points for all synthetic sequences. Each such intermediate occupying a metabolic branchpoint must be partitioned between two competing pathways, one leading to oxidation to carbon dioxide and the other to synthesis of one or more components of the cell. The adenylate system provides energy for the chemical activities of the cell, including biosynthesis, for **active transport** of nutrients and ions across [membranes](#) against chemical potential gradients, and for most other biological requirements, including mechanical movement. All those functions are catalyzed or affected by [proteins](#). Proteins involved in those functions have evolved affinities for ATP and ADP (or AMP) that maximize their functional usefulness to the organism. Thus, it is essential that the ratio of ATP to ADP remain virtually constant. The problem would be equivalent to the regulation of voltage by the power supply of a complex electronic device, if it were not so much more complicated.

## 3. Regulation of the Regeneration and Utilization of ATP

The rate of regeneration of ATP from ADP is regulated in large part by controlling the rate at which substrate is made available to the electron transport phosphorylation system. At least five enzymes that catalyze reactions in glycolysis or the citrate cycle respond sensitively to the status of the adenylate pool and adjust the properties of the catalytic site accordingly. An increase in the energy



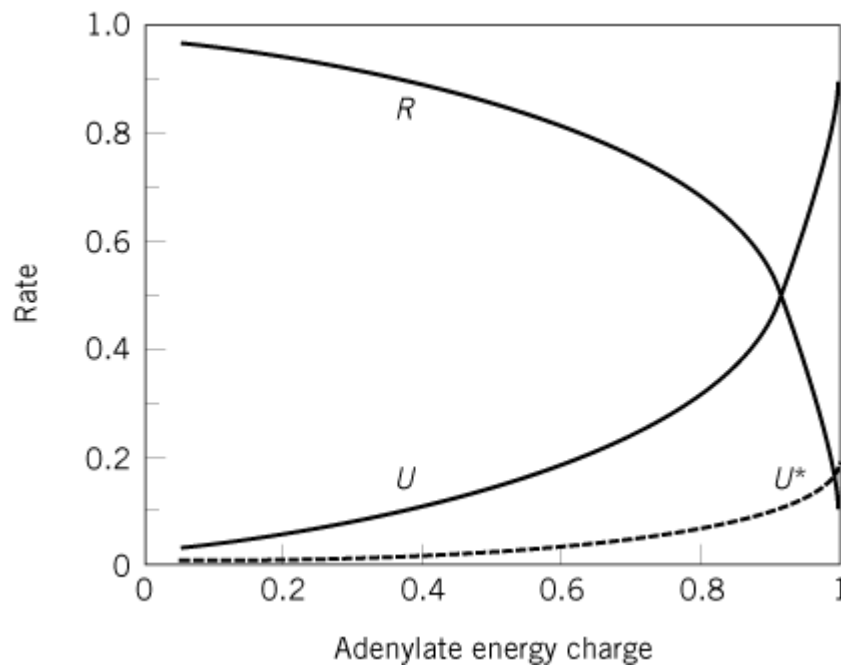
charge causes a decrease in the rate of the reaction catalyzed by the enzyme, and a decrease in charge causes an increase in rate. This **feedback inhibition** system acts to adjust the rate of regeneration of ATP to meet momentary requirements and thus to stabilize the value of the charge. The multiplicity of control sites may be surprising; a single throttle point would seem to be sufficient to regulate the rate of supply of substrate to the electron transport system. The regulatory requirements are, however, much more complex because the central pathways also supply starting materials for synthesis.

Core metabolism consists of the central pathways by which foodstuffs or storage materials of all types are prepared for oxidation and of branches from those pathways that lead to synthesis of one or more products. The primary function of oxidative pathways is regeneration of ATP from ADP, which puts energy into the adenylate system. Biosynthesis is powered by conversion of ATP to ADP, which removes energy from the adenylate system. Adjustment of the partitioning of resources between these oppositely directed pathways to meet the momentary metabolic needs of the cell or organism is the most central of the regulatory requirements that underlie cellular function and survival. At each branchpoint, the partitioning of resources must respond at least to the energy status of the cell as reflected in the ATP/ADP system and to the momentary need for the end products of the synthetic branch.

At such branchpoints, the next enzyme in the degradative pathway and the first enzyme in the biosynthetic branch compete for the branchpoint metabolite, their common substrate. Maintenance of an appropriate balance between oxidation and biosynthesis requires that both enzymes must be regulated. As well as decreasing the rate of the degradative reaction, an increase in the value of the charge leads to an increase in the rate of the competing reaction that channels the branchpoint metabolite into biosynthesis. Thus, enough substrate is oxidized to maintain the normal status of the adenylate system, and biosynthesis is allowed to the extent that the products are needed and the supply of resources allows.

This partitioning is not affected by turning enzymes on and off but by changing the affinities of the competing enzymes for the substrate, which allows for much more sensitive control. Affinity is usually expressed in terms of the Michaelis constant,  $K_m$ , the concentration of substrate at which half the catalytic sites bind substrate, leading to a reaction velocity half the maximal rate. Thus, competition between the enzymes is regulated in the most direct and effective way—by modulating their relative abilities to capture the common substrate. The competition is shown generically in Figure 2. Reactions in pathways that lead to degradation of substrate and the regeneration of ATP respond to variation in the energy status of the adenylate pool as indicated by curve  $R$ , and reactions that direct substrate into biosynthetic pathways that utilize ATP respond as shown by curve  $U$ . The result is that the value of the energy charge is maintained within a narrow range of values near the intersection of the curves. If the charge drifts upward slightly, the rate of use of ATP in biosynthesis tends to increase and the rate of regeneration of ATP decreases, counteracting the drift. A downward drift has the opposite effect.

**Figure 2.** Generalized illustration of the effects of the adenylate charge on the rates of reactions in which ATP is regenerated ( $R$ ) and in which ATP is utilized ( $U$ ). Curve  $U$  was calculated for an enzyme for which the  $K_m$  for ATP at the catalytic site is six times that of ADP. Curve  $U^*$  represents 80% depression of rate as a consequence of feedback inhibition of the enzyme by the end product of the biosynthetic sequence.



The regulatory interactions illustrated by curves *R* and *U* are necessary but not adequate. By themselves, they would adjust the rate of all biosynthetic sequences similarly and only on the basis of the availability of energy, without regard to the need for the individual products. Those controls are supplemented by product feedback inhibition. The first enzyme in nearly every biosynthetic sequence that has been studied bears a regulatory site that binds the end product of the sequence. When that site is occupied, the conformation of the enzyme changes so as to decrease the affinity for substrate at the catalytic site (increase the  $K_m$ ) (see [Allostery](#)). Thus, as the concentration of end product rises, the fraction of enzyme molecules that bind it increases, the affinity of enzyme for substrate decreases, and the enzyme competes less vigorously for substrate. When the concentration of end product falls—for example, when an amino acid is being used more rapidly for protein synthesis—a smaller fraction of the regulatory sites of the first enzyme involved in its synthesis is occupied by end product. The resulting decrease in  $K_m$  causes the substrate to be more successful in its competition with the enzyme that catalyzes the next step in the degradative pathway, and the rate of synthesis increases. Such interactions adjust synthetic rates to meet changing metabolic needs. Curve  $U^*$  in Figure 2 illustrates the response of an enzyme when the end product of its sequence is available in adequate amount from external sources. Any response between curve  $U^*$  and somewhat above curve *U* is possible. Under ordinary conditions, the response will fluctuate in the vicinity of curve *U* to adjust production to meet metabolic demand. These interactions stabilize the pools of amino acids, assuring that they will neither be depleted when demand for protein is high nor build up to unnecessary or injurious levels when demand is small. Pool levels of other metabolites are regulated similarly.

The curves of Figure 2 thus provide a general overview of the regulation of the central pathways in their roles of regenerating ATP and providing starting materials for biosynthesis (1). The interaction of curves *R* and *U* adjust degradative metabolism and overall biosynthetic rates so as to maintain a nearly constant value of energy charge and, in the process, necessarily cause an appropriate partitioning of resources between degradation and synthesis. Superposition of product feedback inhibition, illustrated by curve  $U^*$ , adjusts the rate of each individual synthetic sequence to meet the momentary needs of the organism or cell. Even when the overall rate of biosynthesis is high, the synthetic pathway leading to any given metabolite will be suppressed when that product is in good supply. Thus, the components of this control system act together to assure that the adenylate system,

the immediate energy source for nearly all cell activities, is maintained at a high and constant charge, near 0.9 (2), and that the rate of production of each biosynthetic product is determined by interaction between the overall availability of resources and the cell's need for the individual product. Regulation of biosynthetic rates is roughly equivalent to the decisions made by a human consumer who must balance how strongly an item is desired and how readily it can be afforded.

#### 4. Regulatory Properties of Enzymes

Enzymes that show *R*-type responses (see [Allostery](#)) catalyze reactions in the pathways that supply substrates for electron transport-linked regeneration of ATP but usually do not involve ATP or ADP directly. Such enzymes have, therefore, evolved regulatory, allosteric sites, distinct from the catalytic site, where ADP or AMP can bind, causing conformational changes that increase the binding affinity for the substrate at the catalytic site. Some such enzymes decrease substrate affinity when ATP binds at a regulatory site. Typically, the nucleotide binds cooperatively to two or more sites, and the effect on  $K_m$  is related to the square or higher power of the nucleotide concentration. Such interactions underlie the sensitive response of *R*-type enzymes to variation in the value of the charge (3).

In contrast, ATP and ADP are directly involved in most of the reactions at branchpoints that direct substrate into synthetic sequences. Thus, separate nucleotide-binding regulatory sites are not required. The *U*-type response is a consequence of higher affinity at the catalytic site for ADP, a product of the reaction, than for ATP, a reactant. This reversal of the usual pattern of higher affinity for reactants than for products results in pronounced inhibition by ADP across most of the energy charge range. At a charge of 0.5, for example, when the concentrations of ATP and ADP are approximately equal, ATP would be excluded from most of the catalytic sites because of competition by ADP. The result would be that most of the kinetic response to variation in the ATP/ADP ratio would occur in a rather narrow range near the high end of the energy charge scale. Curve *U* of Figure 2 is calculated for an enzyme for which the value of  $K_m$  for ATP is six times that for ADP.

The regulatory stability illustrated by Figure 2 depends on (1) precise evolutionary adjustment of the affinities for adenylates at regulatory sites of *R*-type enzymes and at catalytic sites of *U*-type enzymes and of the conformational links by which binding at those sites causes conformational changes that modulate affinity for substrate at catalytic sites and (2) the relative affinity at catalytic sites of *U*-type enzymes for ATP and ADP. As a consequence of such interactions, the ATP/ADP ratio, or the energy charge, affects reaction rates more strongly than do the absolute concentrations of the adenylate nucleotides. A mutant strain of *Escherichia coli* unable to synthesize adenylate nucleotides grew on adenine-limiting media at essentially normal rates when the intracellular concentrations of ATP, ADP, and AMP were half the normal levels. The energy charge and the ATP/ADP ratio retained their normal values. That mutant, like other organisms that have been studied, did not grow if the energy charge fell slightly below its usual value of about 0.9 (4).

The controls illustrated in Figure 2 should not be confused with thermodynamic or mass-action effects; they are strictly kinetic. The value of the Gibbs free energy change is high and negative for *U*-type reactions under all physiological conditions, as is also true for *R*-type reactions. Thus, degradation of fuels and synthesis of products are both thermodynamically favorable at all times, and evolved kinetic control mechanisms determine which conversions actually occur. The adenylate energy transduction system, which links metabolic sequences and nearly all other cell activities energetically, is ideally placed for the additional role of mediating the regulatory interactions that underlie the integrated activities of cells and organisms.

#### Bibliography

1. D. E. Atkinson (1972) In *Horizons of Biochemistry* (A. San Pietro and H. Gest, eds.), Academic Press, New York, pp. 83–96.
2. A. G. Chapman and D. E. Atkinson (1977) In *Advances in Microbiological Physiology*, vol. 15

- (A. H. Rose and D. W. Tempest, eds.), Academic Press, New York, pp. 253–308.
3. D. E. Atkinson (1970) In *The Enzymes*, 3rd ed., vol. **1** (P. D. Boyer, ed.), Academic Press, New York, pp. 461–489.
  4. J. S. Swedes, R. J. Sedo, and D. E. Atkinson (1975) *J. Biol. Chem.* **250**, 6930–6938.

### Suggestions for Further Reading

5. D. E. Atkinson (1977) *Cellular Energy Metabolism and Its Regulation*, Academic Press, New York. A general discussion of metabolic interactions and regulation. Chapters 4 (pp. 85–107) and 7 (pp. 201–224) are most relevant to the subject of this article.
6. R. H. Garrett and C. M. Grisham (1995) *Biochemistry*, Saunders/Harcourt Brace, New York, Chapters "17" and "25". An excellent general biochemistry textbook with strong coverage of metabolism.

## Adenylate Cyclases

The ubiquity of [cyclic AMP](#) (cAMP) in regulating **enzymatic** activity and/or **genetic expression** in all kingdoms of life, except for the **archaea** (but see later), accounts for the interest displayed in its mode of synthesis and the vast amount of literature devoted to the enzymes that produce it, the adenylate cyclases. These enzymes, which catalyze synthesis of cAMP from ATP and yield pyrophosphate as a by-product, can be classified into four different classes according to their common features: (1) cyclases related to enterobacterial adenylate cyclases; (2) toxic adenylate cyclases isolated from bacterial pathogens; (3) a large and probably ancient class that comprises cyclases from both eukaryotes and prokaryotes and is strongly related to [guanylate cyclases](#); and (4) one example, presently known only from the eubacteria *Aeromonas hydrophila* and *Yersinia pestis*, that differs entirely from all other classes.

### 0.1. Class I: The Enterobacterial Type

The first complete adenylate cyclase gene, *cya*, was cloned and sequenced from *Escherichia coli*. Work on other enterobacteria, such as *Erwinia chrysanthemi*, *Proteus mirabilis*, *Salmonella typhimurium*, *Yersinia intermedia*, and *Yersinia pestis* demonstrates that both the environment of the genes and the proteins specified are similar in size and overall organization to those of *E. coli* at the corresponding locus ([1](#)). Analysis of the *cya* gene from other bacterial species, related to enterobacteria but distinct from them, using genetic [complementation](#) of appropriate *cya* defective strains of *E. coli* (and more recently by direct sequencing of whole [genomes](#)) reveals that the genes from many other bacteria, in particular *Aeromonas caviae*, *Aeromonas hydrophila*, *Haemophilus influenzae*, *Pasteurella multocida*, and *Vibrio cholerae*, directly synthesize a protein structurally and phylogenetically related to the *E. coli* cyclase ([1](#), [2](#)) (see also the database at <http://www.tigr.org>).

No long stretch of **hydrophobic** amino acid residues is present to explain the membrane-bound localization of the adenylate cyclases. In all cases, the proteins are very rich in [cysteine](#) residues, an uncommon feature for proteins located in the cytoplasm or at the cytoplasmic border of the membrane. This might account for the extreme difficulty in purifying the enzymes. In addition, they are also rich in [histidine](#) residues, which could indicate that metal ions take part in the folding and/or activity of the polypeptide chain, but no experimental data support this hypothesis. Finally, the protein is made of two functionally well-defined **domains**. The catalytic domain is NH<sub>2</sub>-terminal, whereas the glucose-sensitive regulatory domain is COOH-terminal. Comparison of the polypeptide

sequence of the catalytic domain of the *E. coli* enzyme with sequences in the protein data libraries do not reveal significant identities with other known proteins. The catalytic domain sequence has been experimentally identified to be about 420 residues. Differences in the amino acid sequence of the *E. chrysanthemi* enzyme often result from the presence of complementary charged residues in place of neutral ones. This suggests that there are more [electrostatic interactions](#) (including [salt bridges](#)) that stabilize the protein at the lower growth temperature of this bacterium. This observation might be helpful when trying to understand the tertiary structure of the protein, which is still not known.

The carboxy-terminal domain of the protein is involved in regulating of the enzymatic activity, in particular its inhibition by glucose. A component of the phosphorylation cascade that mediates import of glucose in the cell, enzyme IIAGlc, is involved in this regulation, but in a way not yet understood. An aspartate residue (Asp414 in the *E. coli* enzyme) is involved in the process in an unknown way. Tonic inhibition of the catalytic domain by the regulatory domain could be relieved by phosphorylation of this residue, although such phosphorylation has never been demonstrated (3).

## 0.2. Class II: The Calmodulin-Activated Toxic Class

Whooping cough is caused by the **gram-negative** bacterium *Bordetella pertussis*, which secretes many toxic proteins into the medium, including an adenylate cyclase. In 1980 it was discovered that this enzyme is activated by a host protein, [calmodulin](#), which does not occur in bacteria (4). Two years later, Leppla (5) demonstrated that another toxic adenylate cyclase, secreted by a gram-positive bacterium, *Bacillus anthracis*, the etiological agent of anthrax, is also activated by host calmodulin. These observations stimulated intense efforts, but several years were required before the *cya* genes from either organisms could be cloned. However, in 1988 a simple idea, predating its generalization under the name "[two-hybrid system](#)," *in vivo* complementation by a **plasmid** encoding an activator of the function (in this case, calmodulin), permitted the cloning of adenylate cyclase genes coding for the calmodulin-dependent cyclases (6).

*Bordetella pertussis* adenylate cyclase is synthesized as a large bifunctional polypeptide chain of 1706 amino acid residues. This contrasts with the various low values reported for the molecular weight of the purified protein (from 43 to 70 kDa). The explanation became apparent when it was demonstrated that the N-terminal segment of the protein (400 residues) alone displays calmodulin-activated adenylate cyclase activity, whereas the rest of the molecule is responsible for hemolytic activity and for transporting the toxin. Sequence, molecular genetic, and physiological studies indicate that the adenylate cyclase domain is fused to a polypeptide chain similar to that of the *E. coli* hemolysin toxin. Therefore the name cyclolysin was coined for the toxic adenylate cyclase from *B. pertussis*.

The adenylate cyclase of *B. anthracis* has been named after the symptom it triggers in the infected host, edema factor. It is encoded in a plasmid, together with another toxin, the lethal factor and a carrier protein, the protective antigen, necessary to internalize both the edema factor and the lethal factor into host target cells. The adenylate cyclase (edema factor) protein, 800 amino acid residues long, comprises four regions of different function. The first region is a [signal peptide](#), permits secretion of the protein. The second region corresponds to the domain that binds with the protective antigen. The third region encodes the adenylate cyclase function. It is followed by the fourth region of unknown function. These toxic adenylate cyclases have been subjected to a most thorough biochemical analysis, but they have not yet been crystallized.

In spite of several attempts to isolate other members of this class, until 1998, we knew only three examples of such proteins, isolated from extremely distant bacteria, one gram-positive and two gram-negative (the adenylate cyclase from *B. bronchiseptica* is very similar to the *B. pertussis* enzyme) (7). Several examples of similar proteins have now been discovered in *Pseudomonas aeruginosa*, and in *Yersinia* species. Comparison of the catalytic regions of the *B. pertussis* and *B. anthracis* adenylate cyclases identified four conserved regions that are involved in catalysis, calmodulin binding and activation. The first region comprises a sequence, Gly-XXXX-Gly(Ala)-Lys-Ser, similar to the nucleotide-binding motif found in many ATP- or [GTP-binding proteins](#).

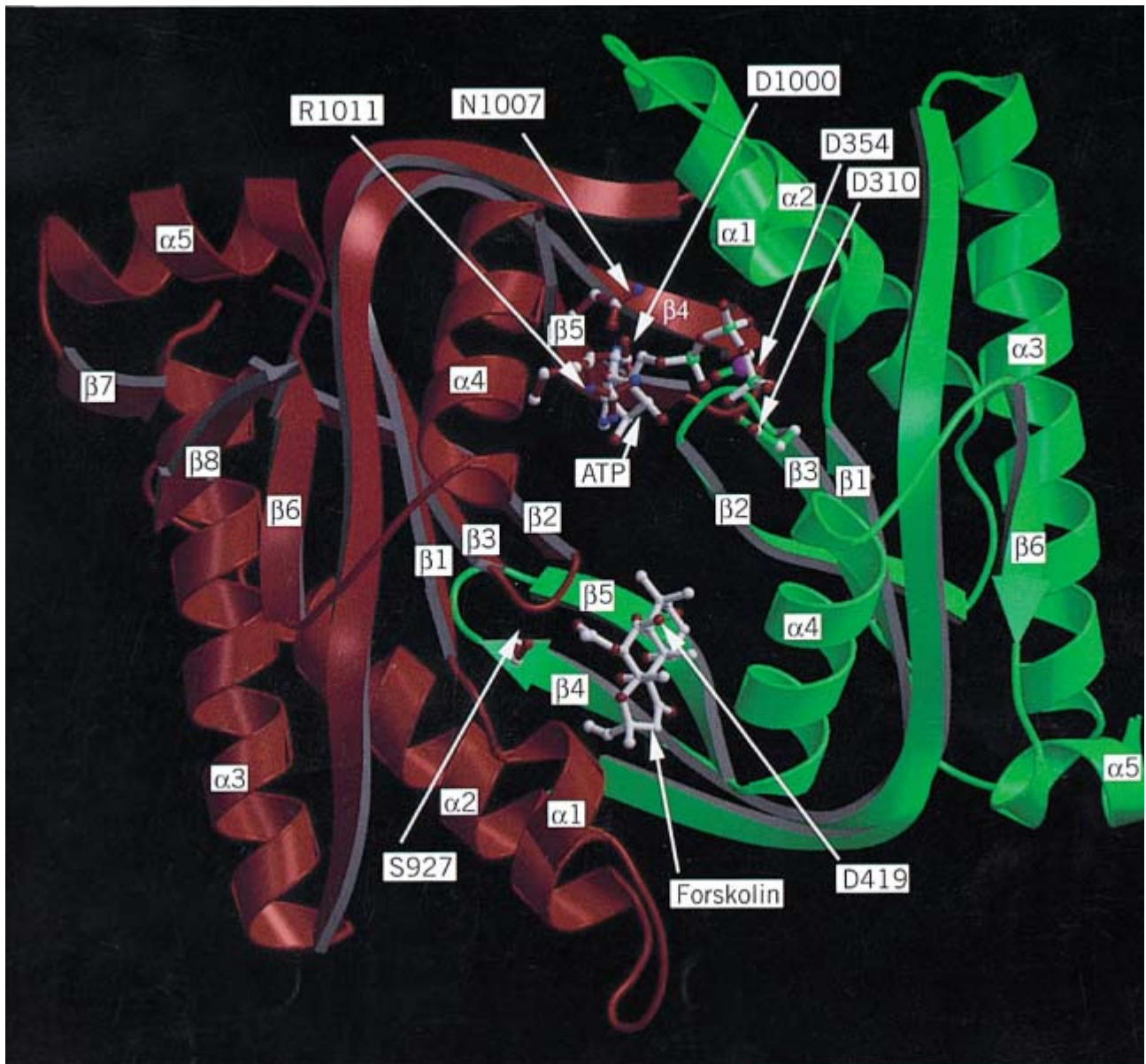
Therefore it was proposed as part of the catalytic site, and *in vitro* mutagenesis substantiated this interpretation. A second region, with the sequence Pro-Leu-Thr-Ala-Asp-Ile-Asp having some similarity in **6-phosphofructokinase**, is also involved in catalysis, and it was proposed that the aspartate residues present in this region are involved in binding ribose and magnesium-phosphate. Although it is strongly conserved, however, the first proline residue does not seem very important because it could be replaced by a leucine residue without any measurable influence on the activity or calmodulin activation of the wild-type enzyme.

Although many calmodulin-dependent enzymes have been identified, the mechanism of activation by calmodulin is still poorly understood (see [Calmodulin](#)). In several cases, limited **proteolysis** releases active, calmodulin-independent forms of the enzymes. Accordingly, it was proposed that the calmodulin-binding domain of these enzymes blocks access of substrates to the [active site](#) and that activation results because an inhibitory domain is removed upon binding calmodulin. The most original feature of the *B. pertussis* protein is that it can be split into two separate domains which recover most of the initial activity when combined. This observation, together with analysis of mutants in the region conserved between the *B. anthracis* and *B. pertussis* enzymes, indicates that these proteins may form a catalytic center from the cooperation of two halves. The function of calmodulin may be to trigger the appropriate conformational movement necessary to form an active catalytic center (3).

### 0.3. Class III: The “Universal” Class

Adenylate cyclases from multicellular eukaryotes have long remained elusive because purifying the corresponding catalytic subunit is extremely difficult. Following intense work all over the world, however, they have been the first adenylate cyclases to be crystallized and analyzed by [X-ray crystallography](#) (see Fig. 1). The activity of these enzymes is subject to a complex regulatory pattern, in particular by [GTP-binding proteins](#). Class III enzymes were first discovered in yeast, but, for convenience, we start with the eubacterial enzymes.

**Figure 1.** Three-dimensional model of the C1 and C2 domains of the catalytic core of type I adenylate cyclase derived from the X-ray crystallographic structure of the C2 homodimer of the type II enzyme (19). The polypeptide backbone of the C1 domain is green and that of C2 is red. The C-terminal nonhomologous segment of C1 is truncated after the  $\alpha$ -helix 5 (low right). Only a few side chains in the active site are shown, designated by one-letter abbreviations. When they differed from the residues in the crystal structure, they were placed according to Liu *et al.* [*Proc. Natl. Acad. Sci.* (1997) **94**, 13414–13419.] ATP was placed as described by Liu *et al.* and modified according to the interactions reported by Tesmer *et al.* [*Science* (1997) **278**, 1907–1916.] The ligand atoms are indicated by C, white; N, blue; O, red; P, green, and Mg, purple. figure is kindly provided by James Hurley. See color insert.



### 0.3.1. Eubacteria

Class III adenylate cyclases form a very diverse collection in eubacteria, both in length and in regulation. Gram-positive bacteria such as *Corynebacterium liquefaciens* secrete large amounts of cAMP because they have a very active class III adenylate cyclase, which is activated by pyruvate. *S. coelicolor* synthesizes a much less active adenylate cyclase involved in aeromycelium formation. Gram-negative bacteria, such as *Rhizobium meliloti*, synthesize at least two different adenylate cyclases. Disruption of both genes simultaneously does not alter cAMP production and suggests the presence of further enzymes. The gram-negative sliding myxobacterium *Stigmatella aurantiaca*, that exhibits an elaborated differentiation pattern, harbors at least two adenylate cyclase genes, each of them corresponding to class III enzymes. They have been partially purified and are inhibited by adenosine, as are the mammalian enzymes (see later). They comprise two domains. The catalytic domain is carboxy-terminal and the regulatory domain is a likely ion transporter in one case and the phosphorylated moiety of a two-component regulatory system in the other. Many other bacteria possess class III cyclases, in particular **cyanobacteria** (8, 9). These enzymes generally comprise two domains. The catalytic domain is carboxy-terminal. There are no indications that they must oligomerize to be active.

### 0.3.2. Lower Eukaryotes

*Saccharomyces cerevisiae* was the first organism from which class III adenylate cyclase genes were cloned and sequenced. The enzyme is activated by the *RAS* gene product (10, 11). Two forms of the enzyme may exist. A long form contains repetitions of a leucine-rich motif that plays a regulatory role and whose significance was recently substantiated and extended. Then it became clear that this eukaryotic adenylate cyclase is completely different from the enterobacterial class because of sequence differences in the catalytic center and also because the organization of the gene is different. The catalytic domain is located at the COOH-terminus in *S. cerevisiae* cyclase, whereas it is found at the NH<sub>2</sub>-terminus in *E. coli*. The yeast enzyme remained the only example of its class until Garbers, Goeddel and co-workers (12) recognized that the genes encoding [guanylate cyclases](#) that had been cloned from several metazoans were derived from an ancestor common to the yeast adenylate cyclase. Another member of class III was subsequently discovered by Young *et al.* (13), who cloned the adenylate cyclase gene from *Schizosaccharomyces pombe* by hybridization using the catalytic domain gene sequence from *S. cerevisiae* as a probe. Finally, the first higher eukaryote adenylate cyclase gene isolated in Gilman's laboratory (14) from bovine brain displayed features clearly reminiscent of this class. Since then, many other genes or [cDNA](#) for adenylate cyclases belonging to this class have been isolated and sequenced from lower eukaryotes: *Saccharomyces kluyveri*, *Trypanosoma brucei* and *T. equiperdum*, *Plasmodium falciparum*, *Neurospora crassa*, and *Dictyostelium discoideum* (3).

### 0.3.3. Higher Eukaryotes: Nine Types

Many class III adenylate cyclases have been identified in higher eukaryotes, in particular in vertebrates, but the most thorough study is in mammals, where several types differing in their regulatory properties have been identified (15). All are regulated in more or less complex ways by **G-proteins** (16). Mammalian adenylate cyclases are informally grouped into nine types according to their tissue location and activity regulation. All but type 9 are activated by the diterpene forskolin, and some are activated by protein kinase C and/or other regulators. Type 1 enzymes were described as calmodulin-activated enzymes from brain. Type 2 proteins are found in brain, lung and other tissues. Type 3 are abundant in olfactory tissue, and the smaller type 4 enzymes are present in testicular tissue. Adenylate cyclase 1, 2, and 8 are positively regulated by calcium/calmodulin, whereas types 5 and 6 are directly inhibited by calcium. Adenylate cyclase 2 and 4 are sensitive to multiple regulatory effects from diverse receptors (3, 15). Adenylate cyclase 9 mRNA, found in rat brain, is particularly abundant in the hippocampus, cerebellum, and neocortex (17). However, the classification into types is somewhat arbitrary (for example, type 4 can also be calmodulin-activated). They all have overall similar structures. Two phylogenetically related cytoplasmic domains are required for catalysis. All types are connected by an integral membrane domain and have variable integral membrane domains at the NH<sub>2</sub> terminus of the protein. Among their many functions, their role in synaptic plasticity and memory is particularly interesting (18) and will certainly make adenylate cyclases extremely fashionable again.

### 0.3.4. Three-Dimensional Structure

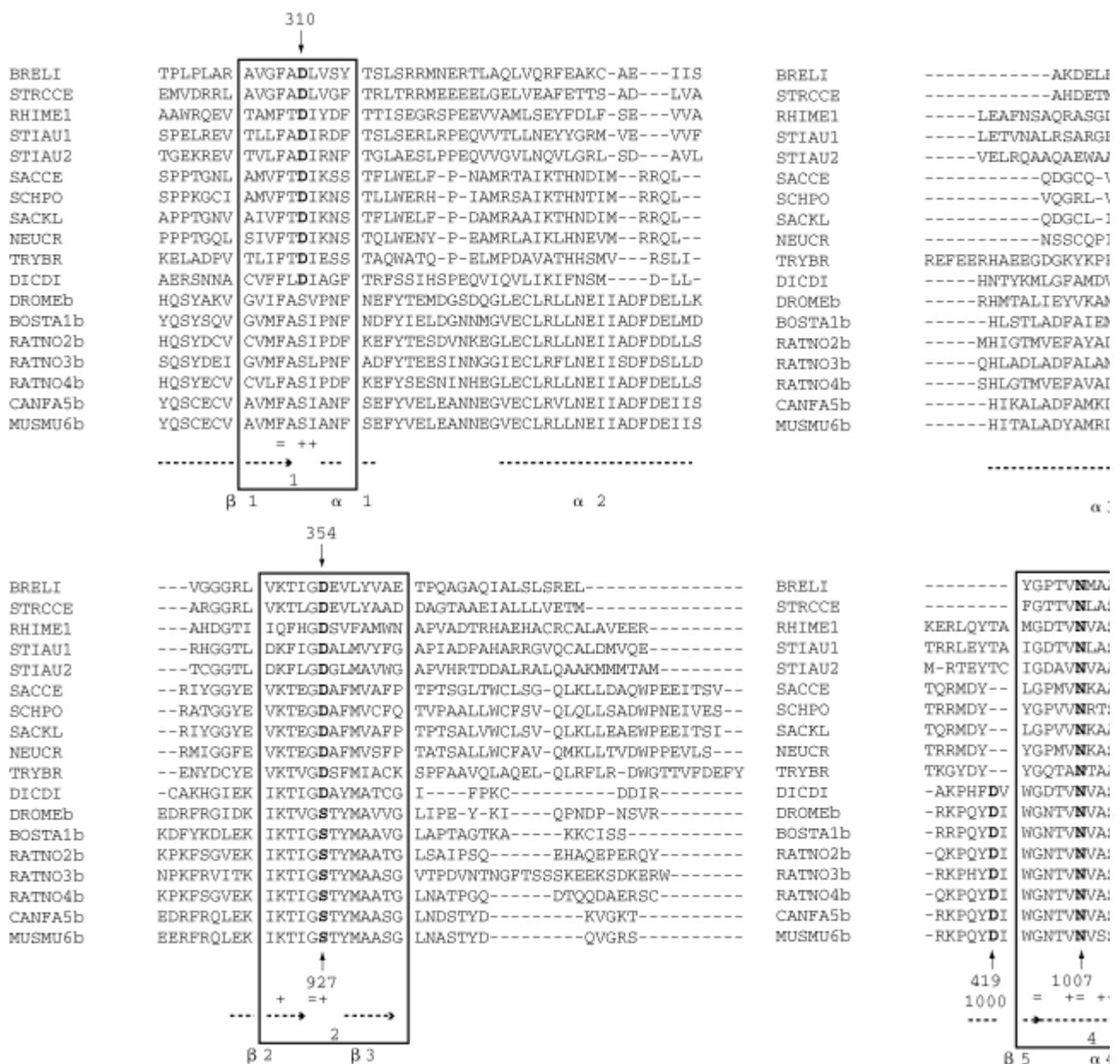
The diverse origins of class III adenylate cyclases are reflected in the wide variation in their general organizations and molecular weights. The smallest protein is the *R. meliloti* enzyme, which contains only a catalytic domain, although some data suggest that an upstream sequence may yield a much longer protein that has a complex regulatory pattern. The next shortest protein, also bacterial, is the enzyme from *C. liquefaciens*. The yeasts produce long proteins, as do higher eukaryotes (except in the case of the testicular enzyme). In all cases, the catalytic domain is located at the COOH-terminus. The mammalian enzymes consist of twelve hydrophobic membrane-spanning regions that form two distinct domains, two cytoplasmic regions that contain both variable and conserved regions, and, in particular, two well-conserved domains that are responsible for catalysis.

Comparison of the catalytic domain sequences of the class III proteins shows that four amino acid stretches are strongly conserved: (1) (Met/Leu/Ile/Val)-(Met/Leu/Ile/Val)-Phe-(Ala/Thr)-(Asp/Ser)-(Leu/Ile)-X-(Asn/Asp)-(Phe/Ser); (2) (Ile/Val)-Lys-Thr-X-Gly-(Ser/Asp)-(Ala/Ser/Thr)-(Tyr/Phe)-



Met; (3) (Met/Leu/Ile/Val)-(Arg/Lys)-(Met/Leu/Ile/Val)-Gly-(Met/Leu/Ile/Val)-(His/Asn)-X-Gly-X-(Val/Ala)-(Val/Leu)-(Ala/Ser)-Gly; and (4) (Trp/Tyr/Phe)-Gly-(Asn/Asp/Pro)-Thr-Val-Asn-X-Ala-Ser-Arg-(Met/Leu/Ile/Val) (see Fig. 2). Crystallization of the catalytic core of one protein and determination of its structure by X-ray diffraction (19) showed that these regions are part of the organization of the structure (Fig. 1). The crystal structure contains the forskolin-binding site but unfortunately not the nucleotide binding site. As a step toward understanding the evolution and function of class III cyclases, enzymes displaying significant guanylyl cyclase activity that have evolved from an adenylyl cyclase ancestor were isolated. A single amino acid residue change (Gly-Asp-Thr-Val-Asn to Gly-Asp-Thr-Ile-Asn in the region of the fourth  $\alpha$ -helix of the catalytic core) alters the nucleotide specificity of the enzyme (20). This corresponds to a pocket situated in a region of the protein that might accommodate the heterocyclic base (19).

**Figure 2.** Alignment of the amino acid sequences of Class III adenylate cyclases from various organisms. The four regions are enclosed in the four boxes. Secondary structural elements found in the crystal structure of Fig. 1 are displayed with corresponding residues of the active site depicted in Fig. 1 are indicated.



#### 0.4. Class IV Adenylate Cyclases

The preceding three classes of structurally unrelated adenylate cyclases already pose a challenging problem. So it was a surprise that *Aeromonas hydrophila* synthesizes another enzyme, a very small cyclase of 193 residues, which has an optimal temperature for activity of 65°C and is at least ten times more active than the class I adenylate cyclase in the same organism (21). No function has yet been discovered for this protein. As yet, it has been found only in various isolates of *A. hydrophila* and in *Y. pestis* (unpublished). There was one report of the presence of cAMP in Archaea, but this was later proven to be the result of an artifact of the growth culture. Therefore it was interesting to see that the sequence of adenylate cyclase from *A. hydrophila* is significantly similar to a gene product of the archaeobacterium *Methanococcus jannaschii*. The gene is expressed in *E. coli*, where it is toxic, but it does not restore cAMP synthesis. Therefore, nothing is yet known about the nature of adenylate cyclases, if they exist, in Archaea, but we may expect that, at some point, they might be discovered and that they would belong to this new class.

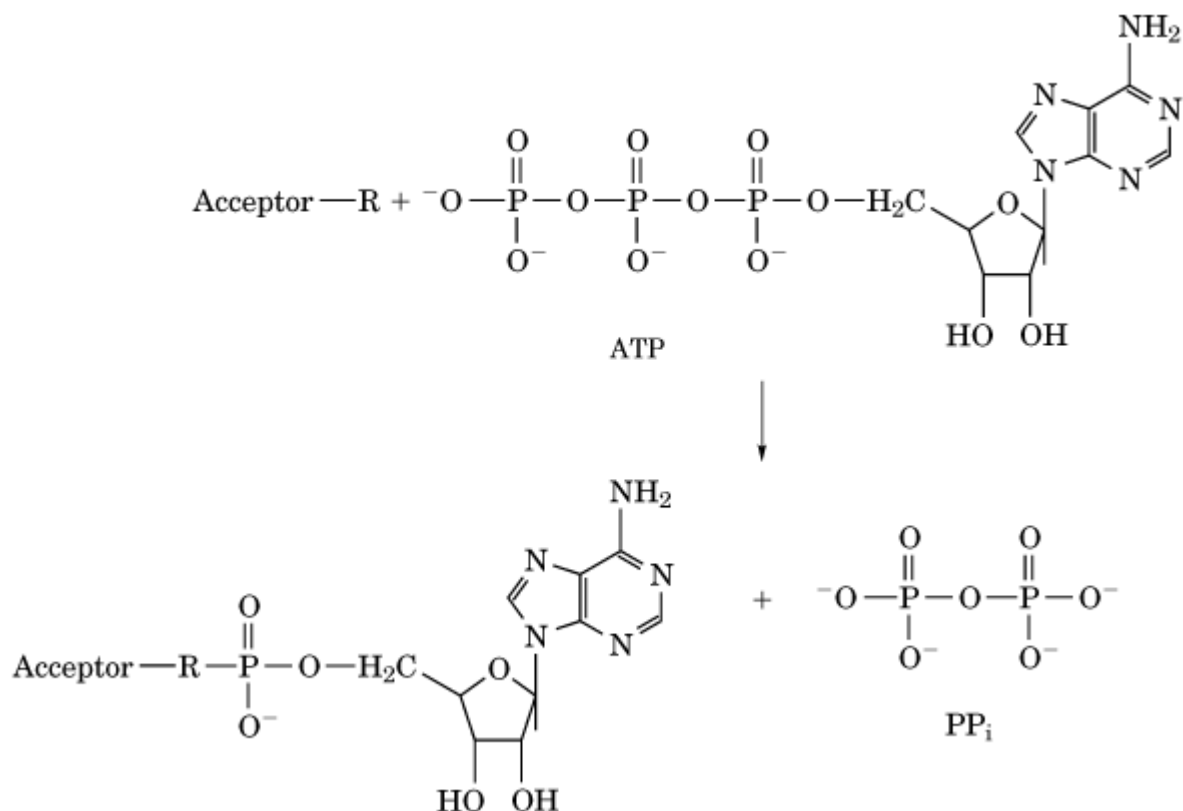
#### Bibliography

1. A. Danchin (1993) Adv. Sec. Mess. Phosphoprot. Res., **27**, 109–161.
2. P. Trotot, O. Sismeiro, C. Vivarès, P. Glaser, A. Bresson-Roy, and A. Danchin (1996) Biochimie **78**, 277–287.
3. O. Barzu and A. Danchin (1994) Prog. Nucleic Acid Res. Mol. Biol. **49**, 241–283.
4. J. Wolff, G. H. Cook, A. R. Goldhammer, and S. A. Berkowitz (1980) Proc. Natl Acad. Sci. USA **77**, 3840–3844.
5. S. H. Leppla (1982) Proc. Natl Acad. Sci. USA **79**, 3162–3166.
6. P. Glaser, D. Ladant, O. Sezer, F. Pichot, A. Ullmann, and A. Danchin (1988) Mol. Microbiol. **2**, 19–30.
7. F. Betsou, O. Sismeiro, A. Danchin, and N. Guiso (1995) Gene **162**, 165–166.
8. M. Kasahara, K. Yashiro, T. Sakamoto, and M. Ohmori (1997) Plant Cell. Physiol. **38**, 828–836.
9. M. Katayama and M. Ohmori (1997) J. Bacteriol. **179**, 3588–3593.
10. T. Kataoka, D. Broek, and M. Wigler (1985) Cell **43**, 493–505.
11. P. Masson, G. Lenzen, J. M. Jacquemin, and A. Danchin (1986) Curr. Genet. **10**, 343–352.
12. M. Chinkers, D. L. Garbers, M. S. Chang, D. G. Lowe, H. Chin, D. V. Goeddel, and S. Schulz (1989) Nature **338**, 78–83.
13. D. Young, M. Riggs, J. Field, A. Vojtek, D. Broek, and M. Wigler (1989) Proc. Natl. Acad. Sci. USA **86**, 7989–7993.
14. J. Krupinski, F. Coussen, H. A. Bakalyar, W.-J. Tang, P. G. Feinstein, K. Orth, C. Slaughter, R. R. Reed, and A. G. Gilman (1989) Science **244**, 1558–1564.
15. J. Hanoune, Y. Pouille, E. Tzavara, T. Shen, L. Lipskaya, N. Miyamoto, Y. Suzuki, and N. Defer (1997) Mol. Cell. Endocrinol. **128**, 179–194.
16. A. Marjamaki, M. Sato, R. Bouet-Alard, Q. Yang, I. Limon-Boulez, C. Legrand, and S. M. Lanier (1997) J. Biol. Chem. **272**, 16466–16473.
17. R. T. Premont, I. Matsuoka, M. G. Mattei, Y. Pouille, N. Defer, and J. Hanoune (1996) J. Biol. Chem. **271**, 13900–13907.
18. M. D. Nielsen, G. C. K. Chan, S. W. Poser, and D. R. Storm (1996) J. Biol. Chem. **271**, 33308–33316.
19. G. Zhang, Y. Liu, A. E. Ruoho, and J. H. Hurley (1997) Nature **386**, 247–253.
20. A. Beuve, E. Krin, and A. Danchin (1993) Compt. Rend. Acad. Sci. Paris **316**, 553–559.
21. O. Sismeiro, P. Trotot, F. Biville, C. Vivarès, and A. Danchin (1998) J. Bacteriol. **180**, 3339–3344.

## Adenylylation

Adenylylation is the process in which adenosine-5'-monophosphate (AMP) is covalently attached to a [protein](#), nucleic acid, or small molecule via a phosphodiester or phosphoramidate linkage. Most often, the AMP is derived from ATP, but in some bacterial adenylylation reactions  $\text{NADP}^+$  is the source. Similarly, deadenylylation is the process in which AMP is removed from the adenylylated molecule. The adenylylation/deadenylylation processes may provide regulatory control of [enzyme](#) activity, contribute to intermediate steps in individual enzymatic reaction mechanisms, or occur as intermediate steps along the biosynthetic pathway of cofactors. In this sense, adenylylation is analogous to **phosphorylation**, [sulfation](#), **methylation**, and other intracellular covalent modification reactions for which multiple functions exist. Adenylylation occurs in a wide range of organisms, including bacteria, yeast, and mammals, although it is less common than many other [post-translational modification](#) reactions as a source of enzyme regulation. The general chemical reaction for ATP-dependent adenylylation is shown in Figure 1.

**Figure 1.** General reaction scheme for ATP-dependent adenylylation. R = O or N. The a, b, and g phosphorous atoms are labeled.



It is important to distinguish adenylylation from the related processes of phosphorylation, adenylation, and **ADP-ribosylation**:

1. Phosphorylation is readily distinguished from each of the other processes by the lack of

incorporation of sugar or adenine in the acceptor.

2. Adenylation results in covalent attachment of ADP via the b-phosphoryl group. Relatively few examples of adenylation are known, and they appear to be limited to adenylation of carbohydrates to yield, for example, glucose-1-ADP.
3. ADP-ribosylation results in a covalent bond between the ribose moiety of NADPH and an acceptor.

These covalent modification reactions also contribute to multiple biological functions, including regulation of enzymatic activity. Analytically, adenylation may be distinguished from adenylylation with the appropriate radiolabeled substrates. Specifically, a-<sup>32</sup>P-ATP, but not b-<sup>32</sup>P-ATP or g-<sup>32</sup>P-ATP, will yield radiolabeled acceptor if adenylylation occurs. Typically, ATP that is radiolabeled with <sup>3</sup>H or <sup>13</sup>C in the adenine moiety is also used, in separate experiments, to confirm that label incorporated into the acceptor is not a result of phosphorylation, without adenylylation. An additional criterion that is often applied to distinguish adenylylation from phosphorylation is its sensitivity to phosphodiesterases. Phosphodiesterases will cleave the AMP from an adenylylated substrate, but they will not hydrolyze phosphate from phosphorylated protein. Specific examples of molecules that are adenylylated are discussed separately below.

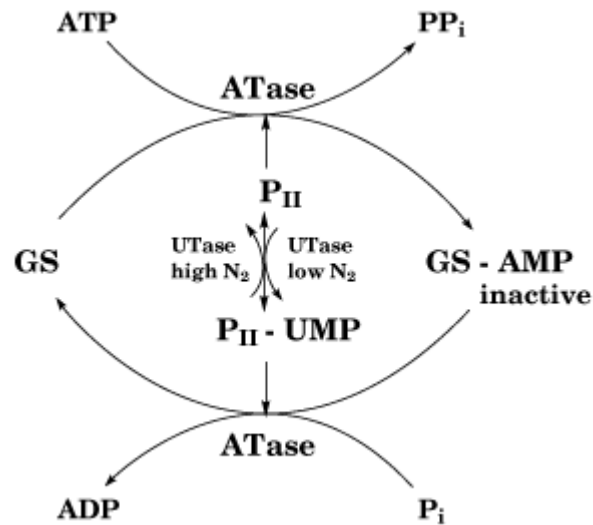
### 1. Glutamine Synthetase

In all organisms, glutamine synthetase (GS) plays a critical role in intermediary metabolism by catalyzing the ATP-dependent condensation of ammonia with glutamate, to yield glutamine. It was discovered by Stadtman and co-workers (1, 2) that *Escherichia coli* glutamine synthetase exhibited dramatically different kinetic properties when isolated from cultures grown in nitrogen-rich vs. nitrogen-starved media. They demonstrated that the differences in enzymatic activity of GS corresponded to the presence or absence of covalently attached adenylyl groups. Since their initial findings, bacterial GS and associated regulatory enzymes have provided a paradigm for understanding biological regulation of nitrogen assimilation in some prokaryotic organisms. Comparison of the molecular details of adenylylation-dependent regulation of prokaryotic GS remains an active research area. Mammalian and plant GS, in contrast, are not regulated by adenylylation.

Regulation of glutamine biosynthesis is characterized most thoroughly for *E. coli* and includes a complex bicyclic cascade that controls the adenylylation of Tyr397 on the surface of GS. Bacterial GS are dodecameric oligomers with two face-to-face hexameric rings. In the subset of bacterial strains that regulate GS activity via adenylylation, the adenylylation state of GS may vary from 0 to 12 AMPs/GS dodecamer, and the enzymatic activity decreases with increasing extent of adenylylation. The regulatory cascade is summarized in Figure 2. It includes the adenylyl transferase (ATase), a [signal transduction](#) enzyme (P<sub>II</sub>), and a uridylyl transferase (UTase) enzyme that responds directly to nitrogen levels. When the nitrogen levels are low, UTase uridylylates P<sub>II</sub> at Tyr51 to form P<sub>II</sub>-UMP. P<sub>II</sub>-UMP stimulates the deadenylylation activity of ATase, to decrease the proportion of GS in the adenylylated form and thus to increase the rate at which nitrogen is “fixed” in the amino acid glutamine, derived from glutamate by GS. When nitrogen levels are high, the UTase cleaves the UMP from P<sub>II</sub>-UMP to generate P<sub>II</sub>. Unmodified P<sub>II</sub> stimulates the adenylylation activity of ATase, to increase the adenylylation state of GS and reduce its efficiency in converting glutamate to glutamine. The complexity of these regulatory cascades underscores the importance of adenylylation as a mechanism for controlling nitrogen metabolism in some prokaryotes.

**Figure 2.** Bicyclic cascade for regulation of *E. coli* glutamine synthetase (GS). ATase catalyzes the

adenylylation/deadenylylation of GS;  $P_{II}$  stimulates the adenylylation activity, and uridylylated  $P_{II}$  ( $P_{II}$ -UMP) stimulates the deadenylylation activity. Uridylylation of  $P_{II}$  is catalyzed by UTase. UTase uridylylates  $P_{II}$  when nitrogen levels are low, and it hydrolyzes the uridylyl group from  $P_{II}$ -UMP when nitrogen levels are high.



The structural basis for the decrease in enzymatic activity of *E. coli* GS upon adenylylation is not known. Although [X-ray crystallography](#) structures are available for the unadenylylated GS from *Salmonella typhimurium* (3, 4), no structure for an adenylylated GS has been determined. [Spectroscopic](#) and [kinetic](#) analyses indicate that adenylylation causes an increase in the [K<sub>m</sub>](#) ([Michaelis constant](#)) for substrates and for the metal cofactors ( $Mg^{2+}$  or  $Mn^{2+}$ ) that are required at each GS [active site](#), although no specific interactions between the adenylyl group and protein residues have been identified (5-7). On the basis of many biophysical criteria and [hydrodynamic volume](#) properties, the adenylylated and unadenylylated forms of GS do not differ greatly in conformation (8, 9). Presumably, adenylyl groups attached at Tyr397 of each subunit in dodecameric GS induce subtle changes in the local conformation of the active sites to which they are adjacent. Spectroscopic methods have suggested that the adenine moiety of AMP attached to each subunit within the hexameric ring structure is highly dynamic, and it collides by [diffusion](#) with the subunit adjacent to it (10, 11). Thus, specific [hydrogen bonds](#) or [electrostatic interactions](#) between the adenylyl groups and protein residues may be limited to the phosphate and ribose moieties of AMP. Alternatively, AMP may simply block the enzyme active sites sterically.

Based on sequence comparisons of GS from numerous prokaryotes, and on mutational analysis of the *E. coli* and *Anabaena 7120* GS, a minimal **consensus sequence** including the Tyr 397 to be adenylylated and a properly positioned proline residue are required for efficient adenylylation by ATase. These two residues are indicated in bold in the local sequence of *E. coli* GS: Met-Asp-Lys-Asn-Leu-**Tyr**-Asp-Leu-**Pro**-Pro-Glu-Glu-Ser-Lys. Individual mutation of several residues in this sequence has negligible or modest effects on the rate of adenylylation by ATase. In contrast, replacement of the bold proline nearly abolishes the ATase-catalyzed adenylylation, and substitution of serine by proline at the analogous position in the *Anabaena 7120* is sufficient to make this GS a substrate for ATase from *E. coli* (12, 13).

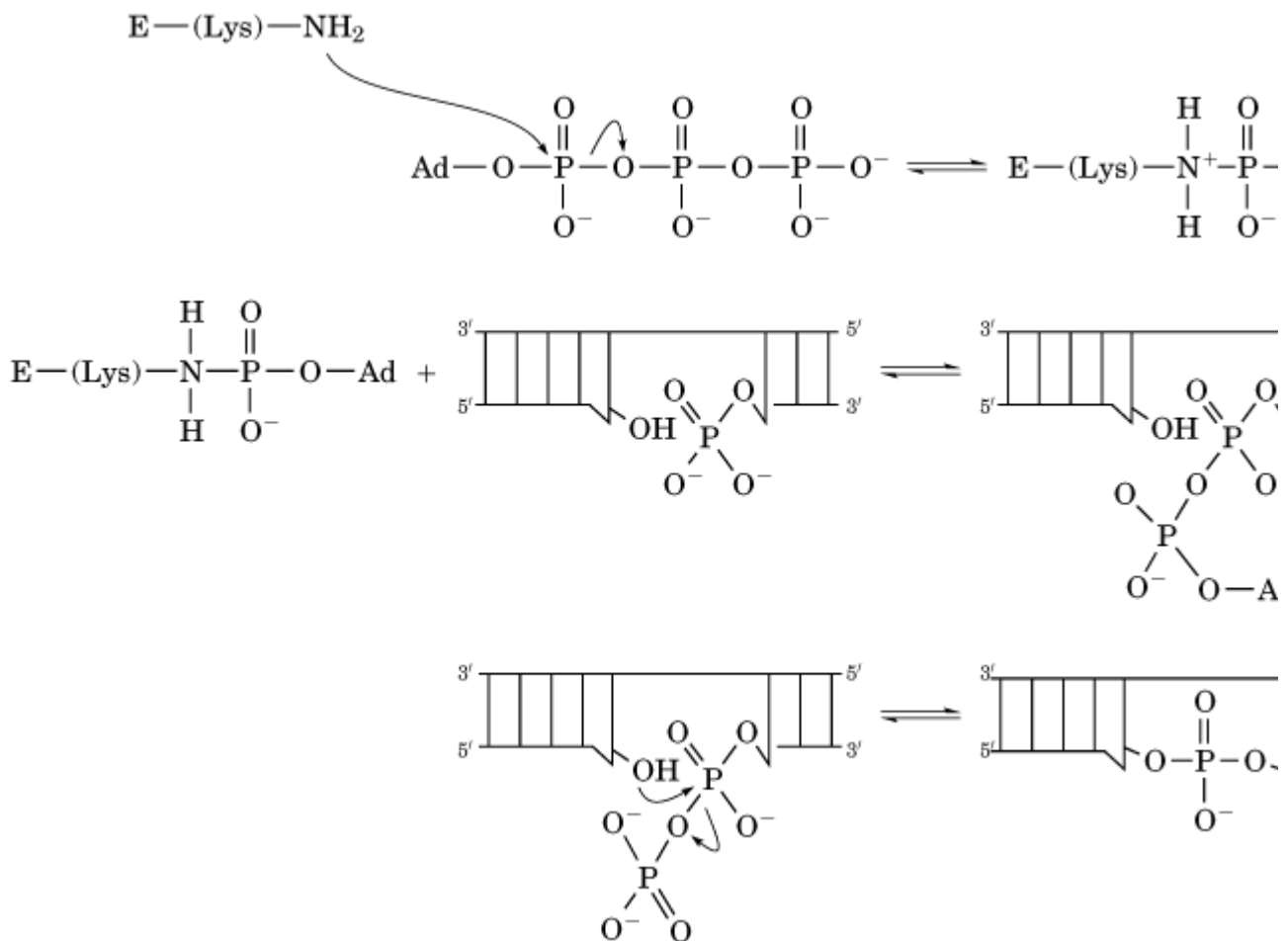
The *E. coli* Atase has been characterized genetically and biochemically; it is a constitutively expressed 945 amino acid residue protein (~115 kDa) that contains two highly homologous **domains** (8, 14). Genetically engineered N- and 3C-terminal domains that have been expressed separately and purified are both catalytically active. However, the two nonoverlapping constructs exhibit different catalytic activities. The N-terminal domain (residues 1 to 423) catalyzes deadenylylation of GS,

requires  $P_{II}$ -UMP, and is inhibited by  $P_{II}$ . In contrast, the C-terminal domain (residues 425 to 945) catalyzes the adenylation of GS and is not regulated by  $P_{II}$  or  $P_{II}$ -UMP. It is striking that  $P_{II}$  is required for the adenylation activity of the full-length ATase, whereas the C-terminal domain by itself is unresponsive to  $P_{II}$  or  $P_{II}$ -UMP. Presumably,  $P_{II}$  and  $P_{II}$ -UMP regulate the adenylation reaction of the C-terminal domain by interacting with the N-terminal domain, or the full-length construct is required to maintain the structural integrity of this regulatory site. In addition, glutamine activates the adenylation reaction of the C-terminal domain. It also is noteworthy that the adenylation and deadenylation reactions catalyzed by ATase are not the reverse of each other. Deadenylation requires inorganic phosphate ( $P_i$ ) rather than  $PP_i$ , and the reaction product is ADP rather than ATP. In effect, this deadenylation reaction may be considered formally as transfer of an adenylyl group from GS to  $P_i$ , or the adenylation of phosphate anion. Although the two functional domains of ATase exhibit significant homology at the amino acid level, they have clearly distinct functional properties. The structural basis for these differences remains unknown.

## 2. DNA Ligases, RNA Ligases, and Related Enzymes

[DNA ligases](#) and [RNA ligases](#) catalyze the formation of phosphodiester bonds at single strand breaks with adjacent 3'-hydroxyl and 5'-phosphate termini in DNA or RNA, respectively. For bacterial, mammalian, and virus ligases, the first step in the catalytic cycle is the adenylation of an active-site e-amino group of a [lysine](#) side chain (Fig. 3). Mammalian and virally encoded ligases utilize ATP to adenylylate the lysine, whereas bacterial ligases exploit  $NADP^+$ . A consensus sequence found in ligases from each of these sources that includes the adenylylated Lys, in bold, is **-Lys-X-(Asp/Asn)-Gly-** (15, 16). The phosphoramidate bond formed upon adenylation of the lysine residue is more stable toward hydrolysis than is the phosphoester bond formed upon adenylation of tyrosine residues or hydroxyl groups of aminoglycosides. Interestingly, a similar consensus sequence is found in [messenger RNA](#)-capping enzymes that transfer guanylate to the 5'-terminus of mRNA, followed by methylation to generate the mature RNA message.

**Figure 3.** Mechanism for DNA and RNA ligases. The catalytic reaction begins with ATP-dependent adenylation of the adenine; E, a ligase.



Functionally related RNA 3'-phospho-cyclases have been identified, and they appear to be present in mammals, yeast, enteric bacteria, and archaeobacteria. This wide distribution suggests that these enzymes have a critical role, but no physiological function has been determined. As part of the catalytic cycle, the 3'-phosphoryl group on the terminal nucleotide of [transfer RNA](#) or **small nuclear RNA** (snRNA) undergoes ATP-dependent adenylation, to generate a phosphodiester bond. This intermediate is attacked subsequently by the 2'-hydroxyl group to release AMP and produce a terminal cyclic phosphate ([17](#)). As with the ligase-catalyzed reactions, this cyclization proceeds via an initial adenylation reaction.

### 3. Other Mammalian Proteins

On the basis of radiolabeling experiments, it has been suggested that plasma [membrane](#) proteins from liver or parotid glands of rats are adenylylated ([18](#), [19](#)). No function has been ascribed to any of these proteins. The hepatic proteins are glycosylated, and their adenylylation is inhibited by [lectins](#) ([20](#)). An interesting possibility is that the adenylyl group is attached to these proteins via their carbohydrates. The same proteins are phosphorylated, and it has been suggested that some protein kinases may be capable of catalyzing adenylylation, in addition to phosphorylation. Further studies that demonstrate covalent attachment of adenylyl groups unambiguously are required, and the physiological function for adenylylation of these proteins remains unclear.

### 4. Aminoglycoside Antibiotics

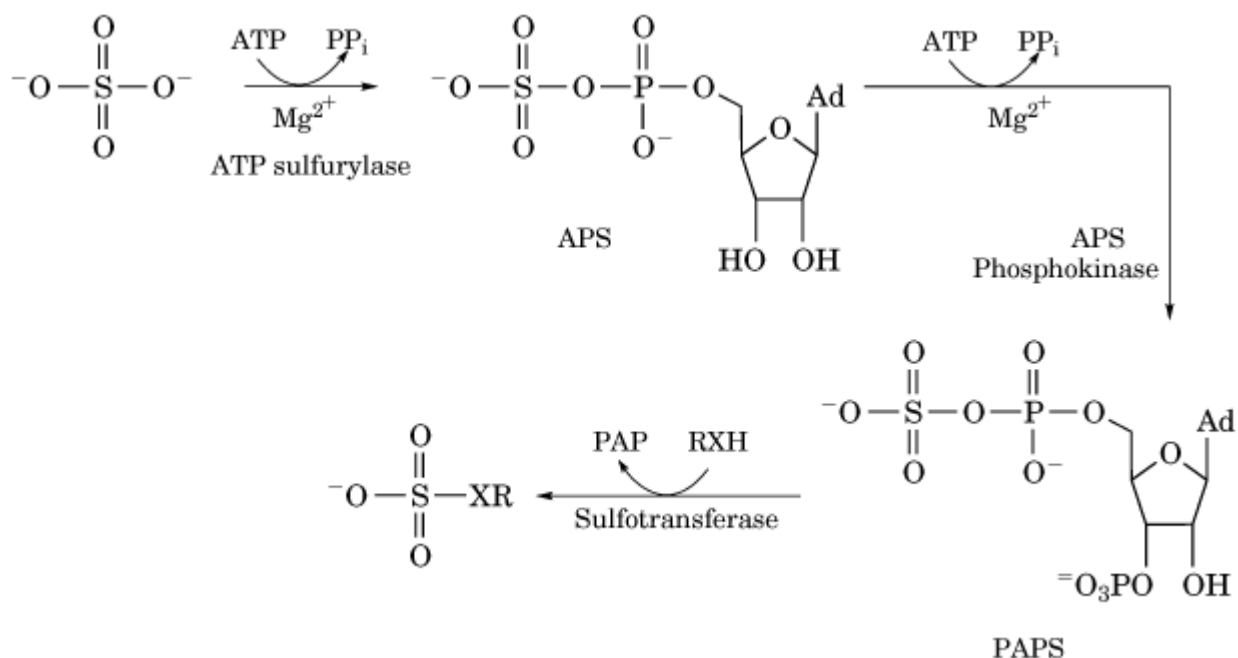
Several aminoglycoside antibiotics, including tobramycin, gentamycin, and [kanamycin](#), are adenylylated by bacterial nucleotidyl transferases, and this contributes to antibiotic resistance in

some strains. The enzymes responsible are not specific for ATP, and they also utilize GTP or UTP, thus resulting in guanylylation or uridylylation as well. These nucleotidyl transferases are **plasmid**-encoded. The adenylylation (nucleotidylylation) site on various aminoglycoside antibiotics differs, but it is most frequently either (a) the 2'- or the 3'-hydroxyl of the 3-aminoglucose ring or (b) the 4'-hydroxyl of the 6-aminoglucose ring. Based on the X-ray crystallography structure of kanamycin nucleotidyl transferase (21), the catalytic mechanism has been proposed to be an in-line displacement of  $PP_i$  by the hydroxyl group of the aminoglycoside, with an active-site glutamate residue acting as a general base. The presence of only a few contacts between the adenine ring of ATP and the nucleotidyl transferase active site is consistent with the lack of specificity for the nucleotide substrate.

## 5. Adenosine-5-Phosphosulfate

An additional example of adenylylation is provided by the biosynthetic pathway of the cofactor 3'-phospho-adenosine-5'-phosphosulfate (PAPS), which is summarized in Figure 4. This cofactor is used by several sulfotransferases in the sulfation of catechols, phenols, and other alcohols. An intermediate formed en route to PAPS is adenosine-5'-phosphosulfate (APS). APS is formed by ATP sulfurylase, in an adenylylation reaction that requires ATP and sulfate ion,  $SO_4^{2-}$ , to yield APS and  $PP_i$ . The formation of APS is partially rate-limiting in PAPS biosynthesis, and deficiency of ATP-sulfurylase may cause impaired sulfation of proteoglycans required for maintenance of extracellular matrix, cartilage, and connective tissues. Thus, the adenylylation reaction catalyzed by this enzyme may have direct clinical impact, manifested as chondrodysplasias, if it operates inadequately (22).

**Figure 4.** Biosynthesis of 3'-phospho-adenosine-5'-phosphosulfate (PAPS). The first step in the biosynthetic pathway is adenylylation of sulfate to generate adenosine-5'-phosphosulfate (APS).



## Bibliography

1. B. M. Shapiro, H. S. Kingdon, and E. R. Stadtman (1967) Proc. Natl. Acad. Sci. USA **58**, 642–649.
2. B. M. Shapiro and E. R. Stadtman (1968) J. Biol. Chem. **243**, 3769–3771.



3. R. J. Almassy, C. A. Janson, R. Hamlin, N. H. Xuong, and D. Eisenberg (1986) *Nature* **323**, 304–309.
4. M. M. Yamashita, R. J. Almassy, C. A. Janson, D. Lasek, and D. Eisenberg (1989) *J. Biol. Chem.* **264**, 17681–17689.
5. A. Ginsburg, J. Yeh, S. B. Hennig, and M. D. Denton (1970) *Biochemistry* **9**, 633–648.
6. L. M. Abell and J. J. Villafranca (1991) *Biochemistry* **30**, 1413–1418.
7. L. P. Reynaldo, J. J. Villafranca, and W. DeW. Horrocks Jr. (1996) *Prot. Sci.* **5**, 2532–2544.
8. S. B. Hennig and A. Ginsburg (1971) *Arch. Biochem. Biophys.* **144**, 611–627.
9. W. M. Atkins (1994) *Biochemistry* **33**, 14965–14973.
10. J. J. Villafranca, S. G. Rhee, and P. B. Chock (1978) *Proc. Natl. Acad. Sci. USA* **75**, 1255–1259.
11. W. M. Atkins, B. M. Cader, J. Hemmingsen, and J. J. Villafranca (1993) *Protein Sci.* **2**, 800–813.
12. E. D. Longton, B. M. Cader, J. Hemmingsen, and J. J. Villafranca (1992) In *Biosynthesis and Molecular Regulation of Amino Acids in Plants* (B. K. Singh, H. E. Flores, and J. C. Shannon, eds.), American Society of Plant Physiologists, Rockville, MD, pp. 59–68.
13. J. Hemmingsen (1991) *Doctoral dissertation*, The Pennsylvania State University, University Park, PA.
14. R. Jaggi, W. C. van Heeswijk, H. V. Westerhoff, D. Ollis, and S. G. Vasudevan (1997) *EMBO J.* **16**, 5562–5571.
15. T. Lindahl and D. E. Barnes (1992) *Annu. Rev. Biochem.* **61**, 251–281.
16. A. E. Tomkinson, N. F. Totty, M. Ginsburg, and T. Lindahl (1991) *Proc. Natl. Acad. Sci. USA* **88**, 400–404.
17. P. Genschik, E. Billy, M. Swianiewicz, and W. Filipowicz (1997) *EMBO J.* **16**, 2955–2967.
18. E. San Jose, A. Benguria, and A. Villalobo (1990) *J. Biol. Chem.* **265**, 20653–20661.
19. M. Hara-Yokoyama, H. Sugiya, and S. Furuyama. (1994) *Int. J. Biochem.* **26**, 1103–1109.
20. E. San Jose, H. J. E. Villalobo, H-J. Gabius, and A. Villalobo (1993) *Biol. Chem. Hoppe-Seyler* **374**, 133–141.
21. L. C. Pedersen, M. M. Benning, and H. M. Holden (1995) *Biochemistry* **34**, 13305–13311.
22. A. Superti-Furga (1994) *Am. J. Hum. Genet.* **55**, 1137–1145.

### **Suggestions for Further Reading**

23.

### **Reviews of Nitrogen Regulation of GS in Bacteria and GS Evolution**

24. J. R. Brown, F. T. Robb, and W. F. Doolittle (1994) *J. Mol. Evolution* **38**, 566–576.
25. B. Magasanik (1993) *J. Cell. Biochem.* **51**, 34–40.

### **Genetic Relationships of Aminoglycoside-Modifying Enzymes, Including Adenylyl Transferases**

26. K. J. Shaw, P. N. Rather, R. S. Hare, and G. H. Miller (1993) *Microbiol. Rev.* **57**, 138–163.

### **Overview of PAPS Biosynthesis**

27. C. D. Klaasen and J. W. Boles (1997) *FASEB J.* **11**, 404–418.

## Adjuvants

The term *adjuvant* designates substances that enhance the [immune response](#) without affecting the specificity of recognition. “Adjuvanticity” was first described in the 1920 by Ramon, who made the observation that mineral substances (such as metal salts and aluminum) or crude materials (such as tapioca) considerably augment the immune response to various vaccines. He then invented the name of “substances stimulantes et adjuvantes de l'immunité.” When an immune response is monitored by the kinetics of occurrence of circulating antibodies, it can easily be shown that, in the presence of adjuvants, the antibody titer is considerably higher and is maintained for a much longer period of time. Although first discovered over 70 years ago, very little progress has been made in this area, and the mode of action of adjuvants remains somewhat elusive. It is generally accepted that they act essentially in two ways: (1) retain the antigen in emulsion or aggregates, depending on the nature of the adjuvant, ensuring a slow but relatively constant release of antigen that may continuously restimulate the immune system; and (2) behave as a nonspecific activator of some partners of the immune response, like mobilizing [macrophages](#) that are acting as **antigen-presenting** cells or by exerting a polyclonal activation of lymphocytes, which is the case for the bacterial lipopolysaccharide (LPS), a potent polyclonal activator of [B cells](#).

To date, the best and universally used substance is the so-called Freund's complete adjuvant, which is an emulsion prepared with a suspension of killed mycobacteria in mineral oil. Unfortunately, it cannot be used for human purposes and must be strictly confined to laboratory animals. Attempts have been made to isolate active molecules from mycobacteria (and also from many other microorganisms). This resulted in a long list of molecules, from which the muramyl dipeptide (*N*-acetylmuramyl-L-alanyl-D-isoglutamine, or MDP) was the most extensively studied. Endowed with good adjuvant properties, it is still too toxic for human use; and many attempts have been made, and are still being made, to define nontoxic homologues. For vaccination purposes in humans thus far, the old recipes, among which are aluminum hydroxide, aluminum phosphate, or alum precipitate, are still the most widely used.

An interesting advance was made with ISCOM (which stand for immuno-stimulating complex), proposed by Morein et al. (1) ISCOM is a cage-like matrix made up of cholesterol and Quil A, which is a substance extracted from the bark of a tree, *Quillaja saponaria*, that contains, after partial purification, five main components with triterpenoid structures. The matrix spontaneously organizes in the presence of [antigen](#), most often protein isolated from a viral coat envelope. The construction ensures high immunogenicity and is used in a number of animal vaccines.

It is likely that one key event played by adjuvants is at the steps of antigen uptake and processing by antigen-presenting cells. Supporting this idea is the fact that bacterial, or more generally particulate antigens, are far better immunogens than the purified molecules isolated from cells. This constitutes an obvious difficulty for designing sophisticated “pure” vaccines, which would at first sight appear advantageous over using complete bacteria that usually contain toxic components, but that also serve as carrier with an adjuvant effect for the desired antigen. To date, an ideal vaccine would have to reconcile the purity of a [recombinant protein](#) with a well-characterized potent adjuvant. We know much about making recombinant proteins, but have at present no convincing pure adjuvant. One interesting approach is to use a living vector, such as [vaccinia virus](#), that has been genetically modified to express a given antigen at the surface of the microorganism. One may also couple a gene expressing the antigen with a gene encoding a **cytokine** that would specifically favor the amplification of a helper [T cell](#) compartment. Another possibility is to embed the antigen in liposomes that could be eventually targeted to well-defined cells of the immune system.

All these approaches have been attempted with variable success. To date, one must unfortunately admit that, despite the fantastic progress made in recent years in the understanding of the basic

mechanisms that operate in the immune system, no really significant advances have been made in defining more efficient and more rational vaccines.

See also the entries [Immune Response](#), [Immunization](#), and [Immunogen](#).

### Bibliography

1. B. Morein, B. Sundquist, S. Höglund, K. Dalsgaard, and A. Osterhaus (1984) Iscom, a novel structure for antigen presentation of membrane proteins from envelopped viruses. *Nature*, **308**, 457–460.

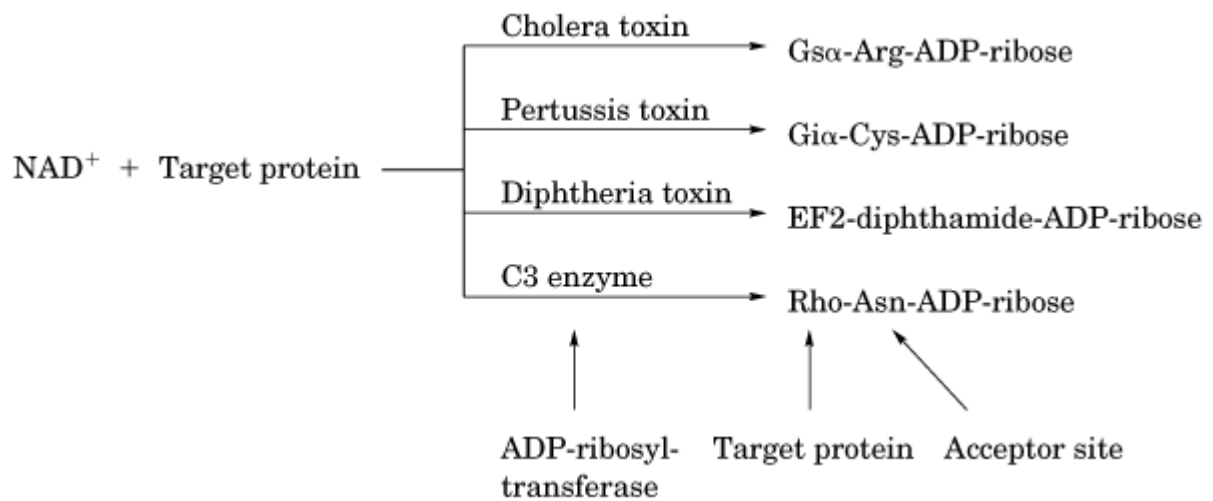
### Suggestion for Further Reading

2. F. R. Vogel (1995) Immunologic adjuvants for modern vaccine formulations. *Ann. N.Y. Acad. Sci.*, **754**, 153–160.

## ADP-Ribosylation, Mono

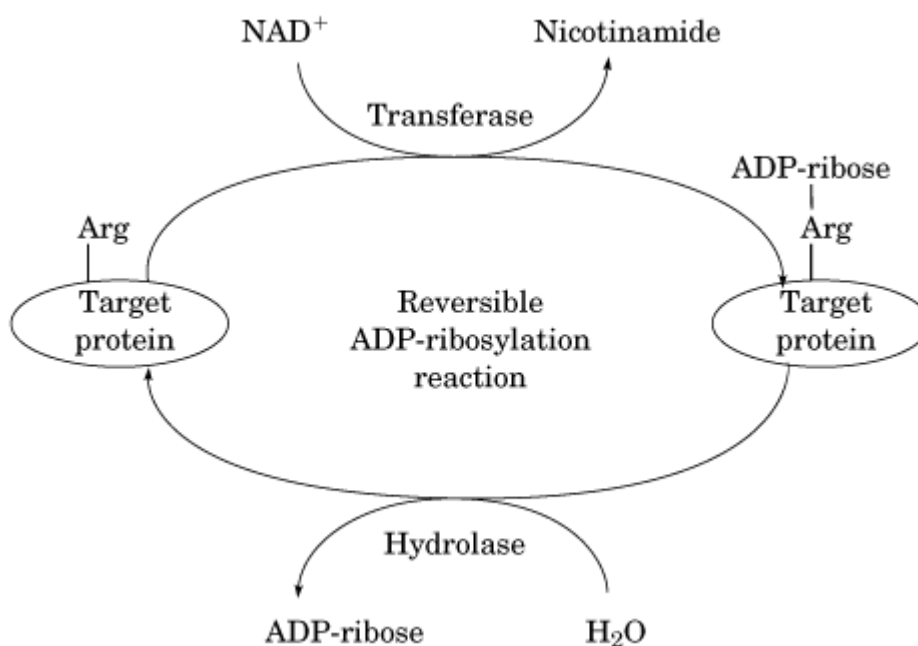
MonoADP-ribosylation is a [post-translational modification](#), catalyzed by ADP-ribosyltransferase (ADPRT), that transfers the ADP-ribose moiety from  $\text{NAD}^+$  to a specific [amino acid](#) residue of target [protein](#) and releases the nicotinamide moiety. To date, four amino acid-specific ADPRT have been reported, specific for [arginine](#), [cysteine](#), diphthamide (a modified form of [histidine](#)), and [asparagine](#) residues. These transferases were originally discovered as the **cholera**, **pertussis**, and [diphtheria toxins](#) and the *Clostridium botulinum* C3 enzyme ([1](#), [2](#)), in the order of amino acid-specificity described above. The native form of such ADP-ribosylating toxins is a heteromultimer, and one of the subunits has ADP-ribosyltransferase activity ([1](#)). The catalytic subunit penetrates the host cell with the assistance of the other subunits. The prokaryotic ADPRT modify target proteins, including a variety of [GTP-binding proteins](#), in eukaryotic cells (Fig. [1](#)). This modification leads to alterations in the target protein, and consequently in cell functions ([1](#)). The bacterial toxins and C3 enzyme have proven extremely useful in studies on [signal transduction](#) pathways in various cell types of eukaryotes. MonoADP-ribosylation must be distinguished from [poly-ADP ribosylation](#).

**Figure 1.** Target protein and acceptor site for ADP-ribosylation in eukaryotic cells by prokaryotic toxins and C3 enzyme.



The reversible Arg-specific ADP-ribosylation (Fig. 2), like protein **phosphorylation**, was initially documented as a regulatory mechanism for control of [nitrogen fixation](#) in the **photosynthetic** bacteria *Rhodospirillum rubrum* and *Azospirillum braziliense* (3). Endogenous dinitrogenase reductase-ADPRT modifies Arg101 on the target protein, dinitrogenase reductase, resulting in inactivation of the reductase. Recovery from the inactive form is achieved through cleavage of the Arg-ADP-ribose linkage by dinitrogenase reductase-activating glycohydrolase. This is the only example of metabolic regulation through endogenous ADP-ribosylation. Purification, characterization, and molecular [cloning](#) of eukaryotic Arg-specific ADP-ribosyltransferase (4-7) and of ADP-ribosyl-Arg hydrolase (8) were reported. By analogy to prokaryotes, metabolic regulation through reversible ADP-ribosylation in eukaryotes has been postulated.

**Figure 2.** Arg-specific ADP-ribosylation and de-ADP-ribosylation reactions.



Among the eukaryotic ADPRTs detected, the best-defined one is Arg-specific. Purification and

characterization of the ADPRT from turkey erythrocytes (4), rabbit skeletal muscle (9), and chicken peripheral heterophils (polymorphonuclear leukocytes) (5) revealed that these enzymes exhibit different physical and regulatory properties, kinetics, and intracellular localization. Guanidino compounds, such as arginine and agmatine, function *in vitro* as acceptors for Arg-specific ADPRT (4). With b-NAD<sup>+</sup> and arginine as substrate, a-ADP-ribosylated arginine is formed by Arg-specific ADP-ribosyltransferase; and the a anomer, but not the b anomer, is utilized as substrate by ADP-ribosyl-Arg hydrolase (10).

Some of the ADPRTs catalyze NAD glycohydrolisis or auto-Arg-specific ADP-ribosylation, when either water or the enzyme itself serve as acceptor for ADP-ribose (1, 2). Upon incubation of intact cells or cell lysates with [<sup>32</sup>P]NAD<sup>+</sup>, it would appear that the ADP-ribose released from NAD<sup>+</sup> by cellular NAD glycohydrolase is attached to some proteins nonenzymatically (11, 12). Thus, enzymatic and non-enzymatic ADP-ribosylations should be differentiated.

Molecular cloning and expression studies of [complementary DNA](#) for Arg-specific ADPRT revealed that the rabbit and human skeletal muscle transferases are glycosylphosphatidylinositol-anchored ([GPI-anchored](#)) (6) and that two forms of transferases from chicken bone marrow cells are secreted (7).

The previously known rat and mouse [T-cell](#) marker RT6 have significant sequence [Homology](#) to the ADPRT. It is now accepted that these proteins are GPI-anchored and possess transferase and/or NAD glycohydrolase activities (13-15).

A common feature of the vertebrate Arg-specific ADP-ribosyltransferase is their resemblance to the transferase from cholera toxin and to *Escherichia coli* heat-labile enterotoxin. Strictly conserved regions around the [active site](#) containing two Glu residues (E207 and E209, indicated with an asterisk in Fig. 3) are seen in all Arg-specific ADPRT detected in vertebrates, bacteria, and viruses. On the other hand, RT6.1 and RT6.2 (rat), which have NAD glycohydrolase activity, but not Arg-specific ADPRT, contain Gln in place of Glu207. The [recombinant protein](#) of a mutant Gln207Glu-RT6.1 possesses Arg-specific ADPRT activity (16, 17), but alterations in the kinetic parameters of the NAD glycohydrolase reaction are slight (17). Furthermore, the mouse homologues of rat RT6, Rt6.1, and Rt6.2 have Glu207 and ADPRT activity. When the Glu is replaced with Gln, Rt6.1 loses the activity (17).

**Figure 3.** Comparison of the amino acid sequences of a Glu-rich motif in eukaryotic, prokaryotic, and viral Arg-specific ADP-ribosyltransferases and T-cell antigen RT6. References are cited by the numbers in parentheses. ADPRT: arginine-specific ADP-ribosyltransferase.

|  |                     |      |
|--|---------------------|------|
| RT6.1                                  | 202: SFYPDQEEVLIPG  | (20) |
| RT6.2                                  | 202: SFRPDQEEVLIPG  | (13) |
| Mouse homologue of the rat RT6, Rt6.1  | 202: SYTHEEEVLIPG   | (21) |
| Rt6.2                                  | 202: SSFPREEEVLIPG  | (21) |
| Chicken bone marrow ADPRT1             | 217: SFFPSEDEVLIPP  | (7)  |
| ADPRT2                                 | 217: SFYPSEDEVLIPP  | (7)  |
| Chicken erythroblast ADPRT             | 217: STFPGEDEVLIPP  | (22) |
| Rabbit skeletal muscle ADPRT           | 233: SFFPGEEVLIPP   | (6)  |
| Human skeletal muscle ADPRT            | 233: SFFPGEEVLIPP   | (23) |
| Cholera toxin                          | 105: SPHPDEQEVLSALG | (24) |
| <i>E. coli</i> heat-labile enterotoxin | 105: SPHPYEQEVLSALG | (25) |
| <i>C. perfringens</i> iota toxin       | 414: PGYAGEYEVLLNH  | (26) |
| Bacteriophage T2                       | 583: LGIATEAEVILPR  | (27) |

It has been reported that defects in RT6 expression are associated in various animal models with the pathogenesis of autoimmune insulin-dependent diabetes and systemic lupus erythematosus (18, 19), although a possible mechanism for RT6-mediated protection from these autoimmune diseases is unknown. Definite physiological functions of monoADP-ribosylation in eukaryotes must be established by further investigation.

## Bibliography

1. S. F. Carroll and R. J. Collier (1994) *Methods Enzymol.* **235**, 631–639.
2. F. R. Althaus and C. Richiter (1987) *ADP-Ribosylation of Proteins*, Springer-Verlag, Berlin, pp. 131–194.
3. P. W. Ludden (1994) *Mol. Cell. Biochem.* **138**, 123–129.
4. J. Moss, S. J. Stanley, and P. A. Watkins (1980) *J. Biol. Chem.* **255**, 5838–5840.
5. K. Mishima, M. Terashima, S. Obara, K. Yamada, K. Imai, and M. Shimoyama (1991) *J. Biochem.* **110**, 388–394.
6. A. Zolkiewska, M. S. Nightingale, and J. Moss (1992) *Proc. Natl. Acad. Sci. USA* **89**, 11352–11356.
7. M. Tsuchiya, N. Hara, K. Yamada, H. Osago, and M. Shimoyama (1994) *J. Biol. Chem.* **269**, 27451–27457.
8. J. Moss, S. J. Stanley, M. S. Nightingale, J. J. Murtagh Jr., L. Monaco, K. Mishima, H.-C. Chen, K. C. Williamson, and S.-C. Tsai (1992) *J. Biol. Chem.* **267**, 10481–10488.
9. J. E. Peterson, S.-A. L. Jacqueline, and D. J. Graves (1990) *J. Biol. Chem.* **265**, 17602–17609.
10. J. Moss, N. J. Oppenheimer, R. E. West Jr., and S. J. Stanley (1986) *Biochemistry* **25**, 5408–5414.
11. H. Hilz, R. Bredehost, P. Adamietz, and Wielckens (1982) *ADP-Ribosylation Reactions. Biology and Medicine*, Academic Press, San Diego, CA, 207–219.
12. D. Cervantes-Laurean, D. E. Mintaer, E. L. Jacobson, and M. K. Jakobson (1993) *Biochemistry* **32**, 1528–1534.
13. F. Koch, F. Haag, A. Kashan, and H.-G. Thiele (1990) *Proc. Natl. Acad. Sci. USA* **87**, 964–967.
14. T. Takada, K. Iida, and J. Moss (1994) *J. Biol. Chem.* **269**, 9420–9423.
15. F. Koch-Nolte, D. Petersen, S. Balasubramanian, F. Haag, D. Kahlke, T. Willer, R. Kastelein, F. Bazan, and H. G. Thiele (1996) *J. Biol. Chem.* **271**, 7686–7693.
16. T. Maehama, S. Hoshino, and T. Katada (1996) *FEBS Lett.* **388**, 189–191.
17. N. Hara, M. Tsuchiya, and M. Shimoyama (1996) *J. Biol. Chem.* **271**, 29552–29555.
18. D. L. Greiner, E. S. Handler, K. Nakano, J. P. Mordes, and A. A. Rossini (1986) *J. Immunol.* **136**, 148–151.
19. F. Koch-Nolte, J. Klein, C. Hollmann, M. Kuhl, F. Haag, H. R. Gaskins, E. Leiter, and H.-G. Thiele (1995) *Int. Immunol.* **7**, 883–890.
20. F. Haag and H.-G. Thiele (1990) *Nucl. Acids Res.* **18**, 1047.
21. C. Hollmann, F. Haag, M. Schlott, A. Damaske, H. Bertuleit, M. Matthes, M. Kuhl, H. G. Thiele, and F. Koch-Nolte (1996) *Mol. Immunol.* **33**, 807–817.
22. T. Davis and S. Shall (1995) *Gene* **164**, 371–372.
23. I. J. Okazaki, A. Zolkiewska, M. S. Nightingale, and J. Moss (1994) *Biochemistry* **33**, 12828–12836.
24. J. J. Mekalanos, D. J. Swartz, G. D. Pearson, N. Harford, F. Groyne, and M. de Wilde (1983) *Nature* **306**, 551–557.
25. T. Yamamoto, T. Gojobori, and T. Yokota (1987) *J. Bacteriol.* **169**, 1352–1357.
26. S. Perelle, M. Gibert, P. Boquet, and M. R. Popoff (1995) *Infect. Immun.* **63**, 4967.

27. T. Koch and W. Ruger (1994) *Virology* **203**, 294–298.

### Suggestions for Further Reading

28. F. Haag and F. Koch-Nolte (1997) "ADP-Ribosylation in Animal Tissues; Structure, Function, and Biology of Mono(ADP-ribosyl)transferases and Related Enzymes". *Advances in Experimental Medicine and Biology*, Vol. **419**, Plenum Press, New York.
29. J. Moss and P. Zahradka (1994) "ADP-Ribosylation: Metabolic Effects and Regulatory Functions". *Molecular and Cellular Biochemistry*, Vol. **138**, Kluwer, Dordrecht, The Netherlands.

## Adriamycin

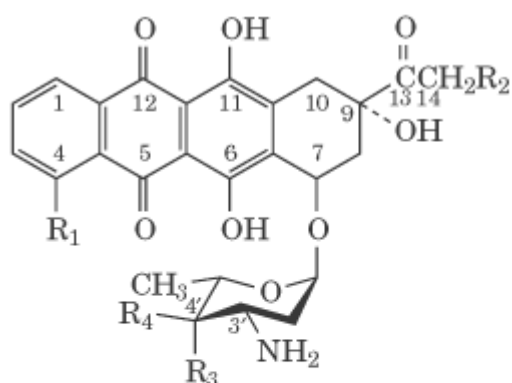
### 1. Introduction

Adriamycin is the trade name for an anthraquinone-containing antibiotic which also has the international nonproprietary name (INN) of doxorubicin (1) and the NCI internal identification number NSC 123127. It is a member of the anthracycline group of compounds that contain an anthraquinone chromophore and a polycyclic ring system (1). It has good anticancer activity against a wide spectrum of tumors (and still remains as the antitumor agent that has the widest spectrum of anticancer activity) and is one of the most extensively used of the fifty or so chemotherapeutic compounds currently in clinical use (2-4). Its full potential as an anticancer agent has not been reached because of a dose-limiting cardiotoxicity (1-4). Despite enormous effort during the last 30 years to develop more effective and less cardiotoxic derivatives, it remains one of the best, proven anticancer drugs. It also exhibits good antibacterial activity.

### 2. Structure and Chemistry

Adriamycin consists of a planar tetracyclic ring system linked by a glycosidic bond to daunosamine, an amino sugar (Fig. 1). It is a bright red compound that is normally isolated as the hydrochloride salt, has the chemical composition  $C_{27}H_{29}O_{11} \cdot HCl$ , a molecular weight of 579.98, and is soluble in water, physiological saline, and methanol (5, 6). The pKa of the amino group is 8.2 (5, 6). In neutral aqueous buffers, it self-associates (7), has absorbance maxima at 233, 480, and 530 nm, and at concentrations less than 25  $\mu M$  has an extinction coefficient of  $11,500 M^{-1}cm^{-1}$  at 480 nm (7). It is fluorescent and can be quantitated in aqueous solutions using excitation and emission wavelengths of 480 and 595 nm, respectively. The solid compound is quite stable if kept in the dark at 4°C and free of moisture. It is photosensitive and therefore solutions must also be kept in the dark at 4°C; even under these conditions, some decomposition will occur in a week or so. It chelates  $Fe^{3+}$  with exceedingly high affinity and readily forms the  $Fe(Adriamycin)_3$  complex (8).

**Figure 1.** Structure of Adriamycin and clinically relevant derivatives.



|            | R <sub>1</sub>   | R <sub>2</sub> | R <sub>3</sub> | R <sub>4</sub> |
|------------|------------------|----------------|----------------|----------------|
| Adriamycin | OCH <sub>3</sub> | OH             | OH             | H              |
| Daunomycin | OCH <sub>3</sub> | H              | OH             | H              |
| Epirubicin | OCH <sub>3</sub> | OH             | H              | OH             |
| Idarubicin | H                | H              | OH             | H              |

### 3. History

In 1958 a new species of *Streptomyces* was identified in a soil sample taken from southern Italy. A red antibiotic, daunomycin was isolated from this microorganism and exhibited good activity against a range of murine tumors (1). In a search for other potentially more active derivatives of daunomycin, the soil microorganism was subjected to the mutagen *N*-nitroso-*N*-methyl urethane and produced a mutant strain that yielded a modified form of the antibiotic. This new antibiotic (named Adriamycin because of the closeness of the original soil sample to the Adriatic Sea) exhibited both a wider spectrum of activity and an improved anticancer response against animal tumors than daunomycin (1) and was also remarkably successful in treating human tumors. Adriamycin was rapidly introduced into clinical trials in Italy and the United States and was approved for clinical use in the United States in 1974 (1).

### 4. Clinical Use

#### 4.1. Tumors

Adriamycin is particularly useful for treating solid tumors such as breast, lung, ovarian, and thyroid carcinomas, as well as soft tissue sarcomas (2-4). It is also active against lymphoid and myelogenous leukemia (2-4).

#### 4.2. Dose and Administration

It is normally administered as a bolus at a dosage of 45–75 mg/m<sup>2</sup> every 3–4 weeks, or more recently as a continuous intravenous infusion during 4–5 days (2, 3). It is inactive if administered orally because the glycosidic bond is hydrolyzed in the gastrointestinal tract. For most treatment regimes it is administered as one component of a combination of drugs (2).

#### 4.3. Side Effects

The most significant problem is a cumulative, dose-dependent cardiomyopathy that can lead to heart failure in up to 10% of patients who receive the maximum recommended dose of 550 mg/m<sup>2</sup>(2, 3). This cardiotoxicity can be acute (hours to days following treatment) or more commonly, can be delayed by months or even years. The cardiotoxicity appears to be due to the redox activity of the drug (9): Adriamycin can undergo a one electron reduction [catalyzed by a range of enzymes such as microsomal cytochrome P450 reductase, mitochondrial NADH dehydrogenase, cytochrome b5 reductase, and xanthine reductase (10)] leading to production of the semiquinone that generates



superoxide in the presence of molecular oxygen, then hydrogen peroxide, and ultimately the extremely reactive and highly toxic hydroxyl radical. These hydroxyl radicals react with any nearby molecule, and hence lead to oxidative damage of critical targets such as DNA and membrane lipids. In most tissues this lethal process is minimized by superoxide dismutase (which converts superoxide to hydrogen peroxide) and catalase (which subsequently converts hydrogen peroxide to water), as well as glutathione peroxidase (which also converts hydrogen peroxide to water). Heart tissues are particularly sensitive to Adriamycin because of their greatly compromised capacity for protection against hydroxyl radicals due to the lack of the bulk hydrogen peroxide detoxifying enzyme catalase and also because Adriamycin inhibits glutathione peroxidase activity (10). Other side effects of Adriamycin therapy are hair loss (alopecia) which usually starts after the first dose of Adriamycin, nausea and vomiting (usually overcome nowadays by appropriate anti-nausea drugs), and reduction in bone marrow function (with maximum suppression 10–14 days following treatment).

#### 4.4. Pharmacokinetics

The disappearance of Adriamycin in plasma is adequately described by a three-compartment model whose half-lives are approximately 10 min, 3 h, and 30 h (11). The plasma level is rapidly depleted and only 0.2 mM remains (average of nine patients) several hours after administration of a 30 mg/m<sup>2</sup> bolus dose of Adriamycin (12); 0.02 μM remains several hours following a 75 mg dose (either as a bolus or 4 h infusion) to a single patient (13). In contrast, the intracellular concentration of drug was approximately 4 μM and 6 μM, respectively, for these two studies (12, 13). The striking features of these results are that the intracellular concentration of drug is typically 100–1000 times greater than that in plasma (14), and that the intracellular drug pool has an extremely long half-life of 4–5 days (13). The uptake into solid tumors is similar to that of normal tissues (12).

**Pharmacology.** The major metabolite of Adriamycin involves metabolism of the C9 side chain alcohol to adriamycinol, where the conversion is catalyzed by the ubiquitous cytoplasmic NADPH-dependent aldo-keto reductase (15). Subsequent microsomal glycosidases that are present in most tissues convert adriamycinol into the inactive deoxyadriamycinol aglycon and daunosamine. The aglycon is then demethylated and conjugated to polar groups to yield more hydrophilic metabolites that are excreted mainly in the bile (15).

The drug distribution has been examined in 20 cancer patients where tissue uptake decreased in the following order: liver > lymph nodes > muscle and bone marrow > fat and skin (12). It has long been known that Adriamycin localizes in the nucleus at the subcellular level (16). This has been confirmed recently in single squamous carcinoma cells by quantitative confocal laser scanning microscopy, where the nuclear level was twice that of the cytoplasm 30 minutes after exposure to 1 μg/mL of Adriamycin, and increased even more with time (17). Interestingly, the cytoplasmic drug level was also higher than that of the medium (17), and this is surprising because Adriamycin is taken up into cells by passive diffusion (18). Nuclear localization is readily attributed to the known high affinity of Adriamycin for DNA (see below). The nature of cytoplasmic localization is less clear, but the distribution as a multitude of small patches of fluorescence throughout the cytoplasm (19) is consistent with some degree of mitochondrial localization (because the drug binds to mitochondrial membranes and, to a lesser extent, also to mitochondrial DNA); it is known that there are as many as a thousand mitochondria per cell (20).

#### 4.5. Resistance

When cells in culture are exposed to Adriamycin for an extended period of time, the cells gradually become resistant to the drug. This type of resistance is due primarily to amplification of the multidrug transporter gene *mdr1* that leads to overexpression of a 170-kDa glycoprotein now known as P-glycoprotein (21). This protein forms a pore in the cell membrane and actively pumps Adriamycin (and many other drugs) out of the cell, using ATP as an energy source. In patients treated with Adriamycin, much of the resistance that develops over several months (and diminishes the antitumor effect of the drug) and appears to derive from overexpression of P-glycoprotein. To enhance the usefulness of Adriamycin, there have been some attempts to develop specific inhibitors of this protein. This approach is no longer considered viable because of the multitude of other drug

efflux pumps that have been identified in recent times (22).

In addition to the phenomenon of multidrug resistance, other mechanisms of resistance have been identified for Adriamycin, including elevated levels of glutathione and decreased levels of topoisomerase II in cells in culture (although not yet demonstrated in tumors) (21).

## 5. Interactions with DNA

Numerous studies have shown that Adriamycin binds to DNA with high affinity by intercalating between adjacent base pairs of DNA. In this process the DNA unwinds to accommodate the drug. This is reflected by an increase in the length of small linear fragments of DNA (hence an increase in the viscosity of the solution) or in the loss of a number of negative supercoils for supercoiled DNA such as a plasmid (that results in an initial increase of viscosity and a decrease of the sedimentation coefficient). These processes have been summarized previously (3, 23). The structure of the intercalated species has been fully characterized by X-ray crystal studies of the drug–oligonucleotide complex. The dominant features are that the drug chromophore lies virtually at right angles to the adjoining base pairs and the amino sugar fits snugly into the minor groove of DNA (24, 25).

Earlier binding studies (before about 1980) misinterpreted curved Scatchard plots as indicating more than one class of DNA binding site. It is now known that this curvature which results from the extreme form of negative cooperativity is displayed when the drug binds to DNA, and is fully described by the neighbor-exclusion principle; values for the intrinsic association constant and the number of occluded base pairs are approximately  $2 \cdot 10^6 \text{ M}^{-1}$  and 3.0 bp, respectively, at near physiological ionic strengths (26, 27). Although this interpretation assumes that the drug intercalates randomly between all possible base pair combinations, this is not strictly correct because there is a slight selectivity of sequence for some sites. Theoretical quantum mechanical calculations indicated the requirement for a three base pair site, and ACG is likely to be the preferred site (28, 29). *In vitro* transcription footprinting studies confirmed the requirement for a triplet site, but showed that the preferred consensus sequence is TCA (where the drug intercalates between C and A) (30), although there was clearly only small energetic differences between this and other triplet sites.

Adriamycin is in rapid equilibrium with DNA. The on-rate is diffusion controlled, and the off-rate is also fast and usually measured by stopped-flow detergent sequestration. This has revealed that Adriamycin has an overall half-life on DNA of approximately 0.5 s at room temperature (31), although by analogy with the structurally similar drug daunomycin, the kinetic processes are likely to be much more complex, especially for longer dissociation times (32).

## 6. Mechanism of Action

The exact molecular events involved in the mechanism of action of Adriamycin have not yet been resolved. Although many possible mechanisms have been proposed and the problems in identifying the critical factors have been well summarized (33, 34), many potential mechanisms have been identified, including impairment of topoisomerase II activity, bioreductive action of the drug, free radical effects, membrane related effects (33, 34). There is a convincing body of evidence to show that the primary target is DNA (33-35). Some of this evidence is that almost all of the drug in the nucleus (>99.8% of single, living cells is associated with DNA (36); more than 80% of all of the drug present in human tumor biopsies is associated with DNA (37); and increasing drug activity (over three orders of magnitude) correlates with increasing DNA binding, damage, or impairment of DNA template activity (38). Therefore DNA binding appears to be related in some way to the anticancer properties, and at present there are two types of interactions with DNA that appear to relate to the drug action: impairment of topoisomerase II activity and formation of drug–DNA adducts.

### 6.1. Impairment of Topoisomerase II

When Adriamycin is added to tumor cells in culture, protein-associated double- and single-strand

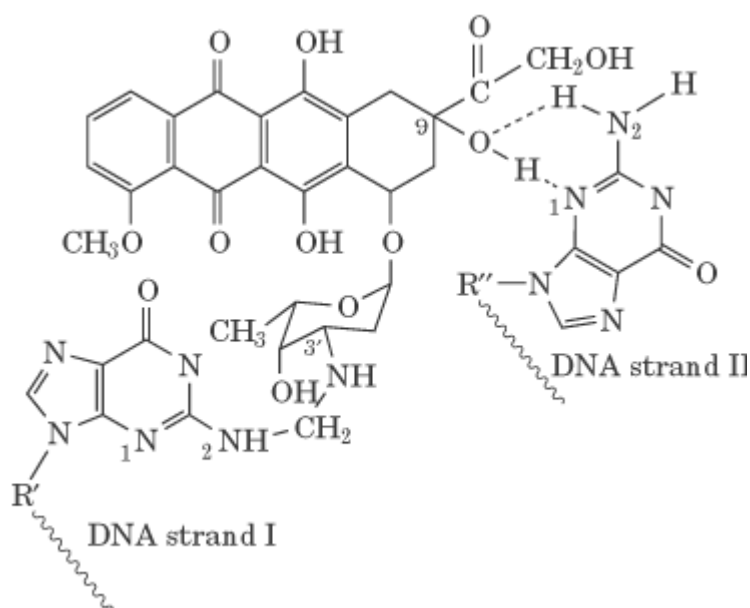
DNA breaks occur close to nuclear matrix attachment sites where topoisomerase II is localized. This enzyme regulates the topological state of DNA, and the DNA breaks appear to arise from intercalation of the drug at that site, which results in structural distortions that interfere with the religation step of topoisomerase II (2, 39, 40). There are several reasons why this appears to be a major contribution to the mechanism of action of Adriamycin: the activity of the enzyme is reduced in tumor cells which are resistant to Adriamycin (39, 40); a correlation has been shown between the induction of double-strand DNA breaks and the cytotoxicity of Adriamycin in P388 leukemia cells (41); and this type of DNA damage occurs at clinical levels of the drug (39, 40). However, there is also evidence that other factors contribute to the mechanism of action of Adriamycin and that the activity cannot be accounted for solely by impairment of topoisomerase II: some studies failed to detect DNA strand breaks in cell cultures at Adriamycin concentrations where cytotoxic responses were observed (42); topoisomerase II mediated damage is rapidly reversible following removal of the drug (39), whereas it is known that DNA double-strand breaks increase long after removal of the drug (43); some derivatives that induce a high level of double-strand breaks exhibit a low level of cytotoxicity (44); and there is little evidence that topoisomerase II is involved in some tumors (10).

## 6.2. Adriamycin–DNA Adducts

There have been many reports of Adriamycin–DNA adducts formed by enzymatic, microsomal, or cellular activation of the drug (33, 45, 46), but it was not until 1990 that it became clear from *in vitro* transcription footprinting studies that these adducts formed predominantly at 5'-GC-3' sequences (47). It was subsequently shown that adducts at these sites behave as an interstrand cross-link (48, 49), although with limited stability [half-life of 5–40 h depending on the DNA sequence and fragment length (49, 50)]. It is now clear that adducts at GC sequences and interstrand cross-links at GC sequences are one and the same lesion and that the adducts stabilize DNA sufficiently so that they function as “virtual interstrand cross-links” (51–53). The adducts have been well characterized in the solid state by X-ray crystal diffraction (54) and in solution by 2-D NMR (55).

Surprisingly, it has been shown that the adducts are mediated by formaldehyde and involve a characteristically unstable aminal linkage to the N-2 of guanine on one strand of DNA (51–55) (Fig. 2). This unusual adduct exhibits some of the characteristics of an interstrand cross-link because of the additional stability that arises from the intercalated chromophore and from additional hydrogen bonds that are formed to the second DNA strand (51–55). The cross-linkage is unstable to heat and to alkali, and this explains why it has proven so difficult to find this lesion in cells in the past (50). Conditions have been established to isolate these lesions from tumor cells (56–59), and this led to a resurgence in cellular studies of this drug. There is now good evidence that these adducts contribute to the mechanism of action of Adriamycin: the capacity of different derivatives of Adriamycin to form adducts (virtual cross-links) correlates well with their cytotoxicity (57), and gene-specific cross-linking assays have detected sufficient lesions in both the nuclear and mitochondrial genomes to show that this lesion is cytotoxic at clinical levels of the drug (58). Total cellular DNA adducts have also been quantitated directly using [<sup>14</sup>C]Adriamycin (58). Because of the formaldehyde-dependent activation of Adriamycin, formaldehyde-releasing prodrugs that have been examined in combination with Adriamycin lead to as much as a 20-fold increase of Adriamycin–DNA adduct levels and a synergistic cytotoxic response (59).

**Figure 2.** Structure of Adriamycin–DNA adducts (52, 54, 55). The exact location of H bonds has not been defined and will depend upon the DNA sequence flanking the 5'-[GC-drug-N]3' binding site (55).



### 6.3. Other Possible DNA-Related Effects

Two additional DNA-related effects also occur at low drug levels *in vitro*, but at this stage the contribution to the mechanism of action of Adriamycin under clinical conditions is unknown (34): impairment of DNA ligase activity (60): impairment of helicase activity (61).

### 6.4. Current Status

The evidence at present indicates that there is no single mechanism of action of Adriamycin. Both impairment of topoisomerase II activity (resulting in protein associated double-strand breaks) and formation of interstrand cross-links occur at clinical levels of the drug. It is known that Adriamycin activates the general cellular response to DNA damage in which the tumor suppressor protein p53 is induced, resulting in arrest of the cell cycle and the induction of apoptosis (62). Although it is known that apoptosis occurs at submicromolar concentrations of the drug (63), this response is yet to be fully understood in the context of drug-induced cytotoxicity (64). The development of new derivatives of Adriamycin in the future is likely to rely heavily on an even more detailed understanding of the molecular consequences of the DNA damage induced by this drug.

## 7. The Search for New Derivatives

Because of the dose-limiting toxicity of Adriamycin, there has been an intense international effort in the past three decades to find new derivatives that have improved anticancer activity and/or reduced cardiotoxicity. This endeavor has been extremely well summarized by Weiss (1). Unfortunately, the outcome from this search has been disappointing. From more than 2000 derivatives tested by 1992 (and perhaps as many as 2500 tested to date), none is substantially superior to Adriamycin for proven clinical anticancer activity. However, several offer advantages for clinical use: idarubicin (lacking the methoxy group at the C4 position) is more amenable to oral administration than Adriamycin; epirubicin is less cardiotoxic and results in less nausea and vomiting than Adriamycin at equimolar doses—it has anticancer activity similar to that of Adriamycin, but reduced side effects have facilitated its wide use as an alternative to Adriamycin. High-dose clinical trials are still in progress with this derivative. One new and encouraging approach that is emerging is the development of “preactivated” forms of the anthracyclines; these formaldehyde-activated derivatives (doxoform, epidoxoform) exhibit increased toxicity to tumor cells, especially to anthracycline-resistant cells (65).

## Bibliography

1. R. B. Weiss (1992) *Semin. Oncol.* **19**, 670–686.
2. V. T. DeVita, S. Hellman, and S. A. Rosenberg (2001) *Cancer: Principles and Practice of Oncology*, 6 ed., Lippincott Williams and Wilkins, Philadelphia.
3. W. B. Pratt, R. W. Ruddon, W. D. Ensminger, and J. Maybaum (1994) *The Anticancer Drugs*, 2 ed, Oxford University Press, New York, pp. 155–165.
4. T. W. Sweatman and M. Israel (1997) In *Cancer Therapeutics, Experimental and Clinical Agents* (B.A. Teicher, ed.), Humana Press, Totowa, NJ, pp. 113–135.
5. F. Arcamone (1982) *Doxorubicin: Anticancer Antibiotics*, Academic Press.
6. R. T. Dorr and D. S. Alberts (1982) In *Current Concepts in the Use of Doxorubicin Chemotherapy* (S.E. Jones, ed.), Farmitalia Carla Erba S.p.A., Milano, Italy, pp. 1–20.
7. J. B. Chaires, N. Dattagupta, and D. M. Crothers (1982) *Biochemistry* **21**, 3927–3932.
8. A. Garnier-Suillerot (1988) In *Anthracycline and Anthracenedione-Based Anticancer Agents* (J.W. Lown, ed.), Elsevier, Amsterdam, pp. 129–161.
9. J. H. Doroshow (1995) In *Anthracycline Antibiotics: New Analogues, Methods of Delivery, and Mechanisms of Action* (W. Priebe, ed.), American Chemical Society, Washington, DC, pp. 259–267.
10. pp. 160–164 of Ref. [3](#).
11. p. 378 of Ref. [2](#).
12. Y. N. Lee, K. K. Chan, P. A. Harris, and J. L. Cohen (1980) *Cancer* **45**, 2231–2239.
13. P. A. J. Speth, P. C. M. Linssen, J. B. M. Boezeman, J. M. C. Wessels, and C. Haanen (1986) *Cancer* **377**, 414–422.
14. C. E. Myers, E. G. Mimnaugh, G. C. Yeh, and B. K. Sinha (1988) In *Anthracycline and Anthracenedione-Based Anticancer Agents* (J.W. Lown, ed.), Elsevier, Amsterdam, p. 529.
15. pp. 158–160 of Ref. [3](#).
16. A. DiMarco (1975) *Cancer Chemother. Rep.* **6**, 91–106.
17. K. Kawai, Y. Minamiya, M. Kitamura, I. Matsuzaki, M. Hashimoto, H. Suzuki, and S. Abo (1997) *Cancer* **79**, 214–219.
18. M. Dalmark and H. H. Storm (1981) *Gen. Physiol.* **78**, 349–364.
19. V. Pillay, R. D. Martinus, J. S. Hill, and D. R. Phillips (1998) *J. Cell. Biochem.* **69**, 1–7.
20. J. W. Shay and H. Werbin (1987) *Mutat. Res.* **186**, 149–160.
21. pp. 50–66 of Ref. [3](#).
22. J. S. Lee, S. Scala, Y. Matsumoto, B. Dickstein, R. Robey, Z. Zhan, G. Altenberg, and S. E. Bates (1997) *J. Cell Biochem.* 513–526.
23. S. Neidle and M. R. Sanderson (1983) In *Molecular Aspects of Anti-Cancer Drug Action* (S. Neidle and M.J. Waring, eds.), Macmillan, London, pp. 35–56.
24. A. H. Wang, G. Ughetto, G. J. Quigley, and A. Rich (1987) *Biochemistry* **26**, 1152–1163.
25. G. Ughetto (1988) In *Anthracycline and Anthracenedione-Based Anticancer Agents* (J.W. Lown, ed.), Elsevier, Amsterdam, pp. 295–334.
26. F. Barcelo, J. Martorell, F. Gavilanes, and J. M. Gonzalez-Ros (1988) *Biochem. Pharmacol.* **37**, 2133–2138.
27. E. Stutter, H. Schuetz, and H. Berg (1988) In *Anthracycline and Anthracenedione-Based Anticancer Agents* (J.W. Lown, ed.), Elsevier, Amsterdam, pp. 245–293.
28. K. Chen, N. Gresh, and B. Pullman (1986) *Nucleic Acids Res.* **14**, 2251–2267.
29. B. Pullman (1991) *Anti-Cancer Drug Design* **7**, 95–105.
30. H. Trist and D. R. Phillips (1989) *Nucleic Acids Res.* **17**, 3673–3688.
31. B. Gandecha and J. R. Brown (1985) *Biochem. Pharmacol.* **34**, 733–736.
32. D. R. Phillips, P. Greif, and R. C. Boston (1988) *Mol. Pharmacol.* **33**, 225–230.

33. pp. 528–569 of Ref. [14](#).
34. D. A. Gewirtz (1999) *Biochem. Pharmacol.* **57**, 727–741.
35. C. Cullinane, A. van Rosmalen, and D. R. Phillips (1994) *Biochemistry* **33**, 4632–4638.
36. M. Gigli, S. M. Doglia, J. M. Millot, L. Valantini, and M. Manfait (1988) *Biochim. Biophys. Acta* **950**, 13–20.
37. J. Cummings and C. S. McArdle (1986) *Br. J. Cancer* **53**, 835–838.
38. L. Valentini, V. Nicoletta, E. Vannini, M. Menozzi, S. Penco, and F. Arcamone (1985) *II. Farmaco Ed. Sci.* **40**, 377–389.
39. C. Holm, J. Covey, D. Kerrigan, K. W. Kohn, and Y. Pommier (1991) In *DNA Topoisomerases in Cancer* (M. Pomesil and K. Kohn, eds.), Oxford University Press, pp. 161–171.
40. Y. Pommier (1995) In *Anthracycline Antibiotics: New Analogues, Methods of Delivery and Mechanisms of Action*. ACS Symposium Series No. 574, pp. 183–203.
41. G. J. Goldenberg, H. Wang, and G. W. Blair (1986) *Cancer Res.* **46**, 2978–2983.
42. F. A. Fornari, W. D. Jarvis, M. S. Orr, J. K. Randolph, S. Grant, and D. A. Gerwitz (1996) *Biochem. Pharmacol.* **51**, 931–940.
43. M. Binaschi, G. Capranico, P. De Isabella, M. Marini, R. Supino, and S. Tinelli (1990) *Int. J. Cancer* **45**, 347–352.
44. M. Binaschi, G. Capranico, L. Dal Bo, and F. Zunino (1997) *Mol. Pharmacol.* **51**, 1053–1059.
45. D. R. Phillips (1990) In *Molecular Basis of Specificity in Nucleic Acid-Drug Interactions* (B. Pullman and J. Jortner, eds.), Kluwer Academic, Dordrecht, The Netherlands, pp. 137–155.
46. J. Cummings, L. Anderson, N. Willmott, and J. F. Smyth (1991) *Eur. J. Cancer* **27**, 532–535.
47. C. Cullinane and D. R. Phillips (1990) *Biochemistry* **29**, 5638–5646.
48. C. Cullinane, A. van Rosmalen, and D. R. Phillips (1994) *Biochemistry* **33**, 4632–4638.
49. S. M. Cutts and D. R. Phillips (1995) *Nucleic Acids Res.* **23**, 2450–2456.
50. A. van Rosmalen, C. Cullinane, S. M. Cutts, and D. R. Phillips (1995) *Nucleic Acids Res.* **23**, 42–50.
51. D. J. Taatjes, G. Guadiano, K. Resing, and T. H. Koch (1996) *J. Med. Chem.* **39**, 4135–4138.
52. D. J. Taatjes, G. Guadiano, K. Resing, and T. H. Koch (1997) *J. Med. Chem.* **40**, 1276–1286.
53. R. A. Luce, S. Th. Sigurdsson, and P. B. Hopkins (1999) *Biochemistry* **38**, 8682–8690.
54. A. H. Wang, Y. Gao, Y. Liaw, and Y. Li (1991) *Biochemistry* **30**, 3812–3815.
55. S. M. Zeman, D. R. Phillips, and D. M. Crothers (1998) *Proc. Natl. Acad. Sci. USA* **95**, 11561–11565.
56. A. Skladanowski and J. Konopa (1994) *Biochem. Pharmacol.* **47**, 2269–2278.
57. A. Skladanowski and J. Konopa (1994) *Biochem. Pharmacol.* **47**, 2279–2287.
58. C. Cullinane, S. M. Cutts, C. Panousis, and D. R. Phillips (2000) *Nucleic Acids Res.* **28**, 1019–1025.
59. S. M. Cutts, A. Rephaeli, A. Nudelman, I. Hmelnitsky, and D. R. Phillips (2001) *Cancer Res.*, in press.
60. G. Ciarrochi, M. Lestingi, M. Fontana, S. Spadari, and A. Montecucco (1992) *Biochem J.* **279**, 141–146.
61. N. R. Bachur, R. Johnson, F. Yu, R. Hickey, N. Appelgren, and L. Malkas (1993) *Mol. Pharmacol.* **44**, 1064–1069.
62. G. Zaleskis, E. Berleth, S. Verstovek, M. J. Ehrke, and E. Mihich (1994) *Mol. Pharmacol.* **46**, 901–908.
63. A. Skladanowski and J. Konopa (1993) *Biochem. Pharmacol.* **46**, 375–382.
64. R. B. Lock and L. Stribinskiene (1996) *Cancer Res.* **56**, 4006–4012.
65. D. J. Taatjes and T. H. Koch (2001) *Curr. Med. Chem.* **8**, 15–29.

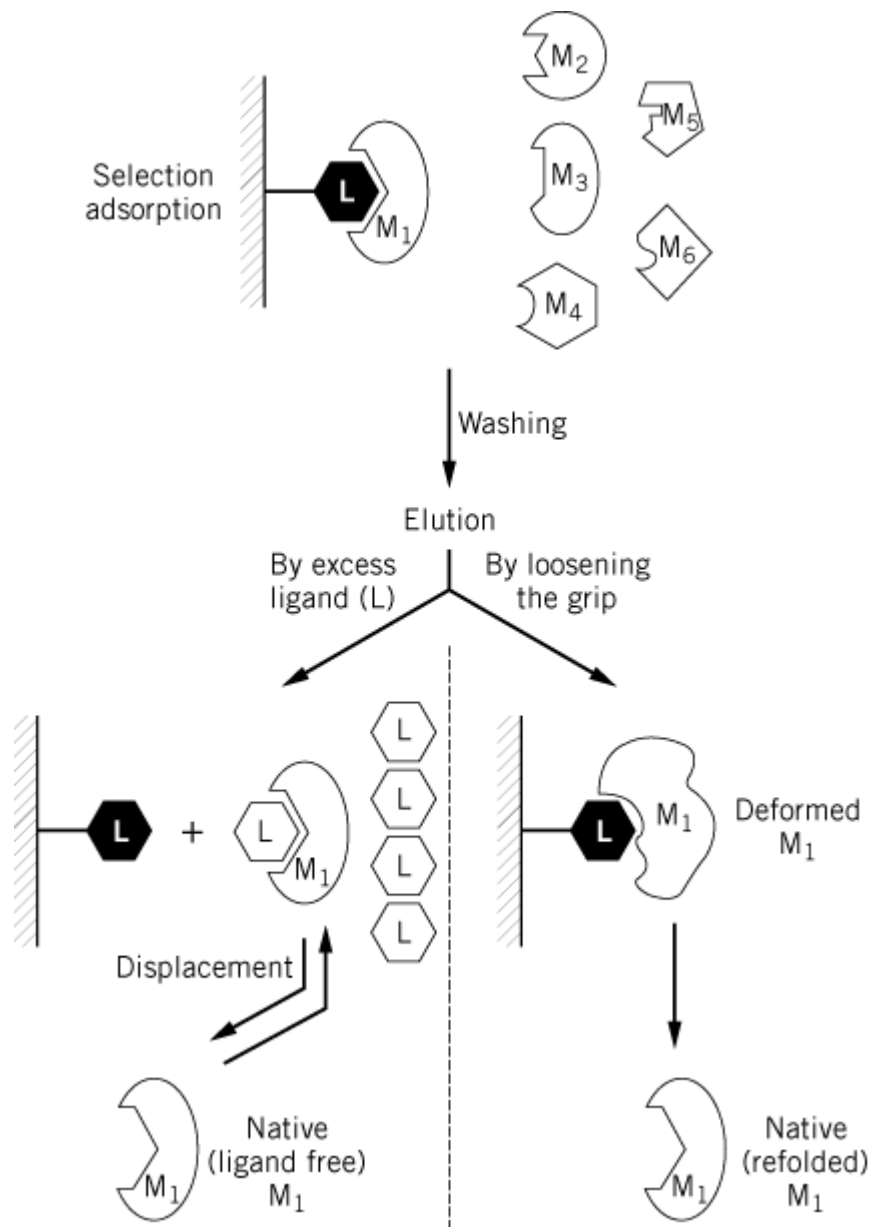
### Suggestions for Further Reading

66. R. B. Weiss (1992) The anthracyclines: Will we ever find a better doxorubicin? *Semin. Oncol.* **19**, 670–686. A detailed review of the history of Adriamycin and of the search for new derivatives.
67. D. A. Gewirtz (1999) A critical evaluation of the mechanism of action proposed for the antitumor effects of the anthracycline antibiotics, Adriamycin and daunomycin. *Biochem. Pharmacol.* **57**, 727–741. An excellent and comprehensive review of possible molecular mechanisms of action of the anthracycline class of anticancer agents.
68. V. T. DeVita, S. Hellman, and S. A. Rosenberg (2001) *Cancer: Principles and Practice of Oncology*, 6 ed., Lippincott Williams and Wilkins, Philadelphia. Contains a concise and current review of Adriamycin and emphasizes cellular and clinical aspects.

### Affinity Chromatography

#### 1. The Basic Principle

Affinity chromatography (AC) is a general **chromatographic** method for the selective extraction and purification of biological macromolecules on the basis of their biorecognition ([1-3](#)). The method makes use of the specific physiological affinity between a desired macromolecule (M) and one of its physiological ligands (L). The ligand, or its analogue (L'), actually acts as a “bait” and is used to extract or “fish out” a desired macromolecule ( $M_1$ ) (Fig. [1](#)) from a mixture of macromolecules ( $M_1$ ;  $M_2$ ;  $M_3$ ;  $M_4$ ;  $M_5$ ;  $M_6^{1/4}$ ). The other macromolecules have a very low (if any) affinity for L, presumably because they are designed to refrain from interfering *in vivo* with the physiological recognition of L by M.



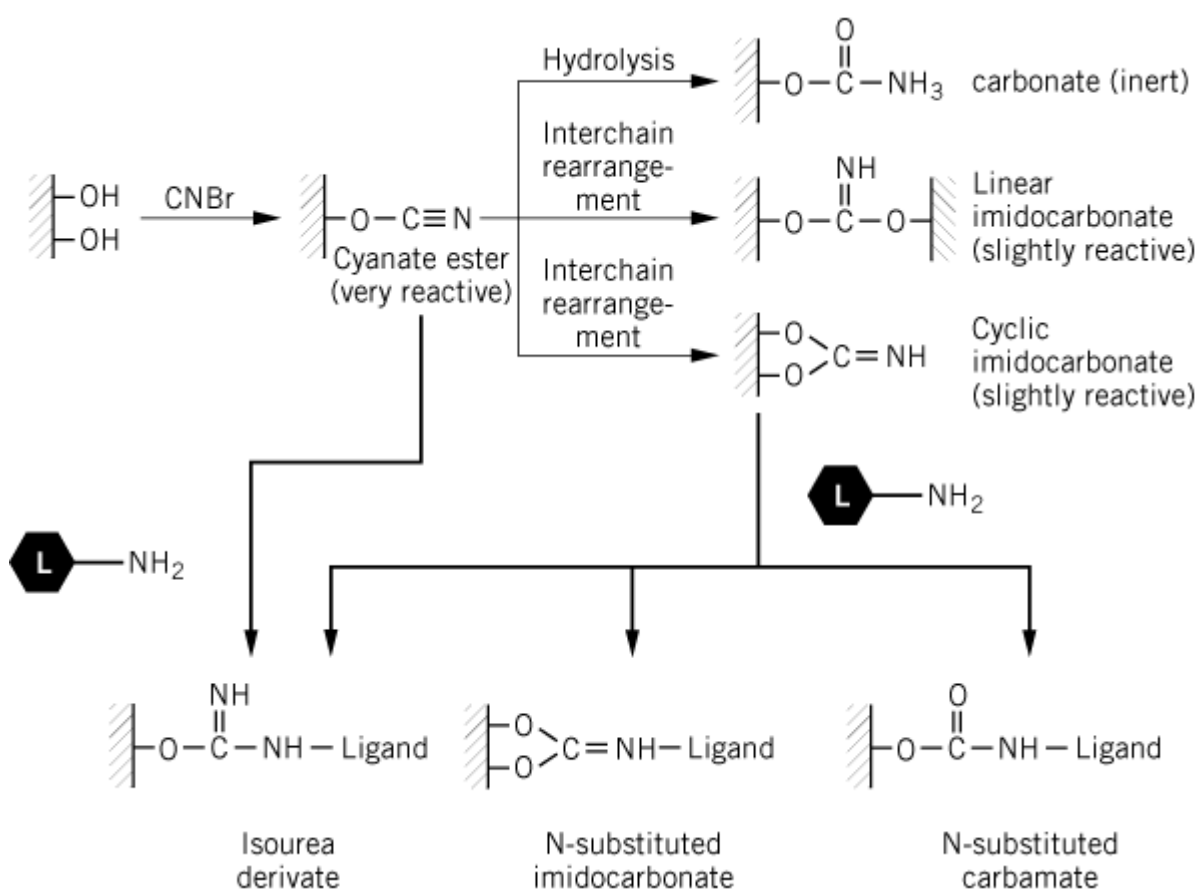
## 2. General Procedure for an AC Purification: Key Steps

1. *Immobilization (anchoring) of the ligand on an inert carrier:* L is anchored on a carrier to yield an insoluble material, usually in a beaded form. This carrier should be as inert as possible (eg, beaded [agarose](#)) to achieve true **active-site-mediated AC**. Also, the attachment point of the ligand should not involve groups that are involved in binding to the macromolecule. Over the years, different carriers and various methods for ligand immobilization were developed. These were reviewed and evaluated by Wilchek et al. (3). In general, the anchoring of L onto an inert carrier involves (i) the introduction of chemically reactive groups to the inert carrier, (ii) the covalent attachment of L to the activated carrier, and (iii) inactivation of the excess of reactive groups (if any) that may remain on the activated carrier after completion of the ligand anchoring step. The immobilized ligand can be used either batchwise or as a column. It may also find other uses—for example, to detect or demonstrate specific [protein–protein interactions](#) by the binding of a specific protein. Furthermore, it may use the resulting column material to bind and fish out another protein that interacts with it. Such immobilized ligands have been used also for labeling of cells, for the localization of proteins on cell surfaces, for the demonstration of leakage of enzymes or specific proteins from damaged tissues, and so on.



Historically, the pioneering work of Axén et al. (4) on the CNBr activation of beaded agarose had a great influence on the development of AC and the conversion of this methodology into a most widely used tool in separation science. To this day, beaded agarose continues to be the inert carrier of choice, and its activation with CNBr for ligand binding is still an activation method of choice. A thorough analytical study of the mechanism of activation of agarose by CNBr (5) showed that three major products are formed: a carbamate (chemically inert), a linear or a cyclic imidocarbonate (slightly reactive), and a cyanate ester (chemically very reactive). Analysis of freshly activated agarose showed that 60% to 85% of the total coupling capacity of the agarose is due to the formation of the cyanate esters. They are the ones that actually react and immobilize the ligand (Fig. 2). On the basis of this mechanism of activation by CNBr, it became possible to develop more efficient activation procedures, which are reviewed in Refs. 3 and 6.

**Figure 2.** The mechanism of the CNBr activation of agarose. (Modified from Ref. 3.)



2. *Selective adsorption.* The key selective step in AC is obviously the extraction of the desired macromolecule M, which is singled out and removed by the immobilized L, from the mixture in which it is present. The macromolecule—be it an [enzyme](#), an [antibody](#), a receptor, a [hormone](#), a [growth factor](#), or the like—is selectively bound by the biospecific ligand L, which can be another protein, a peptide, a polynucleotide or a nucleotide, a polysaccharide or a carbohydrate, a lipid, a vitamin, or just a metal ion. Functionally, L may be a substrate, a substrate analogue, an inhibitor, an [antigen](#), a coenzyme, a cofactor, or a regulatory metabolite. In many cases, the biospecific ligand used for the immobilization is a structural analogue of the physiological ligand (L'). It is imperative, however, to ensure that it still retains the property of selective binding to M, and ideally to M only. In choosing the ligand for an affinity chromatography column, it is often possible to adjust the grip

of M onto the anchored L, and thus to optimize both the adsorption and the elution steps. It should be noted that the adsorption conditions used (buffer, pH, ionic strength, temperature) should also be carefully chosen to secure an optimal and selective adsorption.

3. *Washing out nonspecifically bound impurities.* This is usually carried out with an excess of the buffer used for selective adsorption.

4. *Elution of the desired macromolecule.* The detachment of M from the column (elution) is one of the most important steps in purification by AC. Obviously, the ideal elution is by a specific displacement of M with an excess of its biospecific ligand (Fig. 1). This procedure preserves the native structure of M by forming the more stable complex of M with its biospecific ligand L. When such elution is achieved, it strongly suggests that true **active-site-mediated** AC is involved. However, very often biospecific ligands fail to elute the desired protein, and nonspecific means have to be applied. These usually include a change in solvent or buffer composition, a change in pH or in ionic strength, the addition of a **chaotropic** or a “deforming” buffer, a change in temperature, or a change in the electric field (**electrophoretic** desorption) (3). All these bring about a deformation of the protein (7, 8), a concomitant loosening of the grip of M for L, and consequently elution. In some cases, the binding of M to the L column is so tight that it is not possible to recover M in a fully active form. If M is an enzyme, this may yield a less active preparation (part of the M molecules may be totally inactive, or all molecules may have a lower affinity for the substrate or a lower [turnover number](#)). In some instances, the purified enzyme is fully active, but it may lose its ability to be regulated—for example, if the regulatory domain of M loses its affinity for a regulatory metabolite. Under such circumstances, immobilized ligands with lower affinity for M must be tried. Among the remedies that can be used to improve the elution step, one should note the possibility of binding the ligand to the matrix by means of an easily cleavable form—for example, through an ester bond (9, 10), which can be readily hydrolyzed with a mild base; through a link that includes vicinal hydroxyl groups, which can be readily cleaved with periodate; or through diazo bonds, which can be readily reduced with dithionates (11). It should be remembered, however, that such columns are of limited value, because they can be used only once. Electrophoresis has also been used for elution (12). Because proteins are charged, they will detach from the column and migrate toward the appropriate electrode, if the column with the adsorbed M is exposed to a strong enough electric field. This mild method of elution was successfully applied with high yields in immunoaffinity chromatography and in some AC systems.

### 3. Interposing an “Arm” Between the Ligand and the Matrix Backbone

While developing the basic principles of AC, it was observed that the purification of M is often improved by interposing a hydrocarbon chain (an “arm” or a “spacer”) between L and the matrix backbone (1). It was presumed that such an arm relieves the steric restrictions imposed by the backbone on the ligand, thereby increasing its flexibility and its availability to the protein (13). Such arms were found to improve significantly the extraction of proteins and the efficacy of the purification by AC. Initially, it was assumed that such hydrocarbon arms do not alter the inert nature of the matrix, a condition that obviously has to be ensured to preserve an active-site mediated adsorption of the extracted protein. This assumption seemed reasonable at the time because it had just been shown that at least some water-soluble proteins are quite well described as “an oil drop with a polar coat” (14), implying that the surface of water-soluble proteins is polar and thus not attracted to lipophilic “baits.” We now know that such arms, in and of themselves, may bind proteins. In fact, this observation led to the discovery of [hydrophobic chromatography](#).

### 4. The Limitations of Biospecificity—Interactions that are not Active Site-Mediated

Proteins and their physiological ligands are multifunctional molecules whose functions involve a variety of physical interactions: hydrophobic, **electrostatic**, ion-dipole, and so on. Therefore, it is reasonable to assume that a protein might interact with a column coated with a ligand (very often anchored to the beads at a local concentration much higher than its concentration *in vivo*) not only by

means of its active site. While it is sometimes possible to minimize these nonspecific interactions, it is not always possible to avoid such interfering effects, because they may be an intrinsic property of the system. For example, if ATP is linked to a matrix through its amino group or its ribose moiety; the column thus obtained may retain an enzyme having a biospecific site for ATP; but at the same time, this very column would be negatively charged due to its triphosphate groups, and it would have hydrophobic loci due to its adenine residues. Other proteins, in addition to the desired one, may therefore “regard” the column material as an **ion exchanger** by virtue of its triphosphate groups, or as a hydrophobic column by virtue of its adenine moieties. The efficiency of resolution will then depend on the magnitude of the affinity produced by charge–charge or hydrophobic interactions, as compared to the affinity between the active site of the desired macromolecule and its immobilized substrate or effector analogue. With columns of macromolecular ligands (eg, enzyme subunits, antibodies, lectins), the probability of encountering such built-in interfering effects is considerably higher, because their immobilization usually involves different anchoring points. This leads to a heterogeneous presentation of the various regions of the ligand macromolecule. In some of these presentations, the biospecific active site is available for interaction, while in other presentations the active site itself is inaccessible or sterically hindered. Hydrophobic patches in such ligands may be available for interaction not only in the biospecifically functional presentation, but also in other presentations. In fact, the tendency of a [lectin](#) such as concanavalin A to adsorb onto hydrophobic substances, in addition to its binding to sugars of the mannosyl configuration, was observed in several laboratories.

#### 5. The Relativity of Biological Recognition: Different Proteins may Share a Taste for a Biorecognition Elements

The occurrence of common biorecognition sites in different enzymes is obvious when they are functionally similar, acting on the same substrate (eg, ATP), or utilizing the same cofactor (eg, NAD). This actually forms the basis for general ligand-affinity chromatography (15). However, common biorecognition elements may also be found with proteins having no apparent functional similarity. For example, the free catalytic subunit of cAMP-dependent protein kinase (**protein kinase A**) is preferentially retarded on immobilized soybean [trypsin inhibitor](#) (16). Though initially unexpected, this is actually not surprising; in spite of the fact that **trypsin** and this [kinase](#) catalyze two different chemical reactions (hydrolysis of [peptide bonds](#) versus a phosphotransferase reaction), these two enzymes do have similar biorecognition elements (or subsites) at their active site: trypsin cleaves peptide bonds adjacent to positively charged amino acid residues ([arginine](#) and [lysine](#)), while cAMP-dependent protein kinase phosphorylates [serine](#) residues that are vicinal (in the sequence of amino acids) to the same positively charged arginine and lysine residues (17-20). Similarly, it was shown (21) that [TLCK](#) (a-*N*-tosyl-L-lysine chloromethyl ketone), an [affinity labeling](#) reagent originally designed for labeling the active site of trypsin, specifically attacks a [thiol group](#) at the active site of the catalytic subunit of cAMP-dependent protein kinase. It seems, therefore, that the retardation of the free catalytic subunit on the immobilized inhibitor is due (at least in part) to an affinity between the inhibitor and recognition subsites at the active site of the enzyme.

#### Bibliography

1. P. Cuatrecasas, M. Wilchek, and C. B. Anfinsen (1968) Proc. Natl. Acad. Sci. USA **61**, 636.
2. P. Cuatrecasas and C. B. Anfinsen (1971) Annu. Rev. Biochem. **40**, 259.
3. M. Wilchek, T. Miron, and J. Kohn (1984) Methods Enzymol. **104**, 3.
4. R. Axén, J. Porath, and S. Ernback (1967) Nature (London) **214**, 1302.
5. J. Kohn and M. Wilchek (1981) Anal. Biochem. **115**, 375.
6. S. B. Mohan and A. Lyddiatt (1997) In *Affinity Separations* (P. Matejtschuk, ed.), IRL Press, Oxford University Press, New York, p. 1.
7. S. Shaltiel, J. L. Hedrick, and E. H. Fischer (1966) Biochemistry **5**, 2108.
8. J. L. Hedrick, S. Shaltiel, and E. H. Fischer (1969) Biochemistry **8**, 2422.
9. R. J. Brown, N. E. Swaisgood, and H. R. Horton (1979) Biochemistry **18**, 4901.

10. P. Singh, S. D. Lewis, and J. A. Shafer (1979) *Arch. Biochem. Biophys.* **193**, 284.
11. P. Singh, S. D. Lewis, and J. A. Shafer (1980) *Arch. Biochem. Biophys.* **203**, 776.
12. M. R. Morgan, P. J. Brown, M. J. Lieiand, and P. D. Ocan (1978) *FEBS Lett.* **87**, 239.
13. P. Cuatrecasas (1970) *J. Biol. Chem.* **245**, 3059.
14. D. C. Phillips (1966) *Sci. Am.* November, 78.
15. K. Mosbach (1978) *Adv. Enzymol.* **46**, 205.
16. E. Alhanaty, N. Bashan, S. Moses, and S. Shaltiel (1979) *Eur. J. Biochem.* **101**, 283.
17. H. G. Nimmo and P. Cohen (1977) *Adv. Cyclic Nucl. Res.* **8**, 145.
18. O. Zetterquist et al. (1976) *Biochem. Biophys. Res. Commun.* **70**, 696.
19. B. E. Kemp, E. Benjamin, and E. G. Krebs (1976) *Proc. Natl. Acad. Sci. USA* **73**, 1038.
20. P. Daile, P. R. Carnegie, and J. D. Young (1975) *Nature* **257**, 416.
21. A. Kupfer, V. Gani, J. S. Jimenez, and S. Shaltiel (1979) *Proc. Natl. Acad. Sci. USA* **76**, 3073.

## Affinity Electrophoresis

By analogy to [affinity chromatography](#), it is possible to introduce specific **ligands** for a [macromolecule](#) into the gels of [gel electrophoresis](#) and to measure the specific retardation of the macromolecule due to its interaction with such a reagent. The advantage of such affinity methods lies in the augmented resolving power conferred by the specificity of the binding interaction.

The procedures used to introduce affinity reagents into gels have varied. In cross electrophoresis, a ligand with a net charge opposite to the species of interest migrates electrophoretically into the gel in the opposite direction. Alternatively, uncharged ligands can simply be added to the gelation mixture. Macromolecular substrates within a gel may serve as immobilized affinity reagents, either by themselves or as carriers of covalently attached affinity groups. The magnitude of the electrophoretic retardation depends on the concentration of the affinity reagent in the gel; quantitative determination of this relationship makes it possible to estimate the apparent association constant for binding of the ligand to the sample. Further information concerning the interaction can be gained from affinity electrophoresis by variation of the buffer composition (eg, the addition of metal ions to the buffer), the pH, or the temperature.

### Suggestions for Further Reading

- T. C. Bog-Hansen and J. J. Hau (1981) Glycoproteins and glycopeptides (affinity electrophoresis). In *Electrophoresis: A Survey of Techniques and Applications*, Vol. **18B** (Z. Deyl, A. Chrambach, F. M. Everaerts, and Z. Prusik, eds.), Elsevier, Amsterdam, pp. 219–252.
- T. C. Bog-Hansen and K. Takeo, eds. (1989) Symposium on affinity electrophoresis. *Electrophoresis* **10**, 811–870.
- K. Takeo (1987) Affinity electrophoresis. *Adv. Electrophoresis* **1**, 229–280.

## Affinity Labeling

Affinity labeling is a strategy to modify chemically an [amino acid](#) residue within a specific **ligand-binding** site of an [enzyme](#), either at the [active site](#) or at a regulatory, **allosteric** site. In this approach, a reagent is designed that resembles structurally the natural ligand of the enzyme but features in addition a functional group capable of reacting covalently and indiscriminately with many different amino acid residues. The designed reagent is intended to mimic the natural ligand in forming a reversible enzyme-reagent complex analogous to the enzyme-substrate complex and, once directed to that specific site, to react irreversibly with an amino acid residue accessible from that site. In the case of a purified enzyme, such a reagent allows identification of a particular ligand binding site or **domain**, which can be an experimental evaluation of a predicted binding site location based on recognition of a [protein motif](#) from its amino acid sequence. Affinity labeling constitutes a valuable starting point for selecting appropriate target sites for subsequent [site-directed mutagenesis](#) experiments and is an important tool in probing structure–function relationships in enzymes. If the synthesized reagent has a characteristic absorbance or fluorescence spectrum, affinity labeling offers a means of introducing a chromophore at a specific substrate site to report on the enzyme conformation or on distances between designated sites in the enzyme. For medicinal chemistry, affinity labeling permits mapping of the substrate binding site and establishment of its size, so that inhibitory drugs directed toward that site can be designed more rationally. In the case of a heterogeneous cell preparation or of a complex protein mixture, a specific affinity label can be used to identify a **receptor** protein or to tag a class of macromolecules, such as **nucleotide binding** proteins.

One characteristic of affinity labeling is the initial formation of a reversible enzyme–reagent complex ( $ER$ ), as indicated below.



where  $E$  and  $R$  represent the free enzyme and reagent, respectively, and  $ER'$  is the covalently modified enzyme. The formation of a reversible  $ER$  complex is often indicated by a “rate saturation effect” in which the rate constant for modification ( $k_{\text{obs}}$ ) increases as the reagent concentration is elevated until the enzyme site is saturated with reagent; subsequently, the rate constant ( $k_{\text{max}}$ ) is not changed by further increases in reagent concentration (see [Kinetics](#)). This kinetic pattern contrasts with the linear dependence on reagent concentration of the rate of a direct bimolecular chemical modification. For an affinity label,  $k_{\text{obs}}$  can be described by the equation

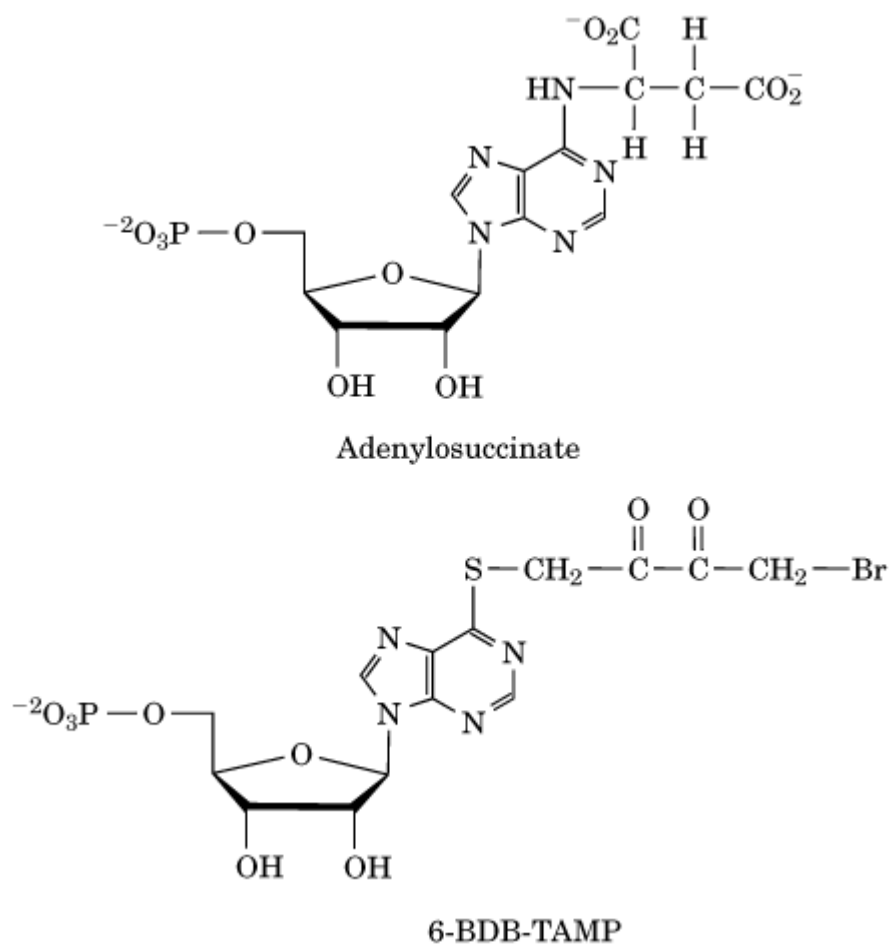
$$k_{\text{obs}} = \frac{k_{\max}}{1 + \frac{K_R}{[R]}} \quad (2)$$

where the apparent dissociation constant for the enzyme–reagent complex  $K_R$ , is given by  $(k_{-1} + k_{\max})/k_1$  and where  $k_{\max}$  is the maximum rate of modification at saturating concentrations of reagent. This type of kinetic behavior is often presented as a double reciprocal plot of  $1/k_{\text{obs}}$  versus  $1/[R]$ , based on the equation

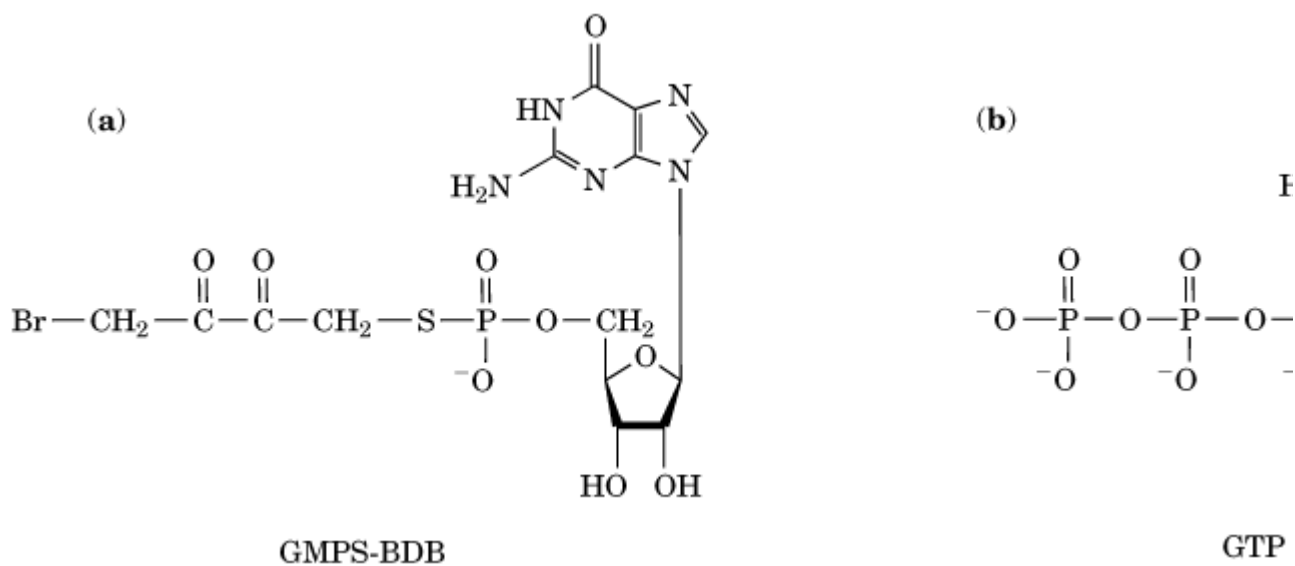
$$\frac{1}{k_{\text{obs}}} = \frac{1}{k_{\max}} + \left( \frac{K_R}{k_{\max}} \right) \left( \frac{1}{[R]} \right) \quad (3)$$

Examples of **purine** nucleotide-based affinity labels are shown in [Figures 1 and 2](#). Experiences with their use illustrate the possibilities of affinity labeling.

**Figure 1.** Comparison of the structures of the natural ligand adenylosuccinate and of the affinity label 6-(4-bromo-2,3-di 5'-monophosphate (6-BDB-TAMP).



**Figure 2.** Comparison of the structures of the affinity label (a) guanosine 5'-0-[S-(4-bromo-2,3- dioxobutyl)thio]phosphate (b) the natural ligand GTP.



## 1. BDB-TAMP

The compound 6-(4-bromo-2,3-dioxobutyl)thioadenosine 5'-monophosphate (6-BDB-TAMP) shown in Figure 1 was synthesized as a reactive nucleotide analogue to target nucleotide binding sites in enzymes (1). The bromoketo group, adjacent to the 6-position of the purine ring, can potentially react with several nucleophiles found in proteins, including those of the side chains of [cysteine](#), [histidine](#), [tyrosine](#), [lysine](#), [methionine](#), **glutamate**, and **aspartate** residues. In addition, the dioxo group provides the possibility of reaction with arginine residues.

As shown in Figure 1, the structure of 6-BDB-TAMP also exhibits a remarkable resemblance to that of adenylosuccinate, a key metabolic intermediate in the conversion of inosine monophosphate to adenosine monophosphate. The last step in that pathway is catalyzed by adenylosuccinate lyase, an enzyme proposed to initiate the cleavage reaction to AMP by attack of an enzymic general base on the b hydrogen of adenylosuccinate. Elimination of the amino group is then facilitated by protonation of the leaving group by an enzymic general acid. Despite this proposal, the key general base and acid of the enzyme had not been identified. The similarity between the structures of adenylosuccinate and 6-BDB-TAMP suggested that the latter would bind irreversibly to the adenylosuccinate site of adenylosuccinate lyase, where the bromodioxobutyl group would be likely to occupy the succinyl subsite at which the critical catalytic steps should occur.

6-BDB-TAMP has recently been confirmed to function as an affinity label of *Bacillus subtilis* adenylosuccinate lyase (2). The initial inactivation rate constant exhibits nonlinear dependence on the concentration of 6-BDB-TAMP, with an apparent reversible  $K_R = 30 \mu\text{M}$  prior to irreversible inactivation at pH 7.0 and 25°C. The tetrameric enzyme incorporates about 1 mol of 6-BDB-[<sup>32</sup>P]TAMP per mol of enzyme subunit on complete inactivation. The substrate adenylosuccinate or the products AMP plus fumarate protect against inactivation and incorporation of radioactive reagent, indicating that 6-BDB-TAMP targets the adenylosuccinate binding site. Purification of the only radioactive peptide labeled by 6-BDB-TAMP led to the identification of His<sup>141</sup> as the modified amino acid (2). These results indicated that 6-BDB-TAMP is an affinity label of His<sup>141</sup> in the substrate binding site of adenylosuccinate lyase, where it may serve as a general base accepting a proton from the succinyl group during catalysis (2) (see [Histidine \(His, H\)](#) residues). This study illustrates many of the desirable characteristics of affinity labeling: binding of the reagent by the enzyme prior to modification of a single site under mild conditions, competition between the reagent and the natural ligand (adenylosuccinate), and identification of the modified residue within a protein of known amino acid sequence.

## 2. GMPS-BDB

A reactive guanine derivative is shown in Figure 2, guanosine 5'-0-[S-4-bromo-2,3-dioxobutylthio]phosphate (GMPS-BDB), which represents a novel class of compounds containing a bromodioxobutyl group linked to the sulfur of a purine nucleotide thiophosphate (3, 4). The reactive moiety of GMPS-BDB is located at a position equivalent to that of the pyrophosphate region of GTP, as illustrated in Figure 2. Thus, it might be expected that GMPS-BDB would act as an affinity label for GTP sites in proteins (see [GTP-Binding Proteins](#)).

Adenylosuccinate synthetase catalyzes the first of the two enzymatic reactions in the conversion of IMP: the condensation of IMP and aspartate to form adenylosuccinate, as GTP is hydrolyzed to GDP and inorganic phosphate. The phosphoryl group has been proposed to be transferred to IMP from the terminal phosphate of GTP to form a 6-phosphoryl-IMP intermediate.

GMPS-BDB has been used as an affinity label of adenylosuccinate synthetase from *Escherichia coli*, and it has been demonstrated that inactivation occurs concomitantly with the modification of Arg<sup>143</sup> of each subunit of the dimeric enzyme (5). The modification of Arg<sup>143</sup> and the inactivation by GMPS-BDB is prevented by adenylosuccinate or by IMP plus GTP, implying that the reaction target is in the region of the active site. This result pointed to Arg<sup>143</sup> as a logical target for [site-directed mutagenesis](#); when the positively charged arginine was replaced by the neutral leucine, the expressed mutant enzyme exhibited a significant decrease in its affinity for nucleotides (5).

The crystal structure of *E. coli* adenylosuccinate synthetase has been determined, allowing a comparison to be made between the substrate sites identified by affinity labeling of the enzyme in solution and those assigned within the crystalline form by [X-ray crystallography](#). Arg<sup>143</sup> from one subunit projects into the putative active site of the second subunit, indicating that both subunits of dimeric adenylosuccinate synthetase contribute to each active site and that Arg<sup>143</sup> plays an important role in nucleotide binding.

### 3. AMPS-BDB

The adenosine analogue of GMPS-BDB has also been synthesized: adenosine 5'-0-[S-(4-bromo-2,3-dioxobutyl)thiophosphate] (AMPS-BDB) (4) and can be considered as a reasonable mimic of ADP or ATP. Bovine liver glutamate dehydrogenase is an allosteric enzyme that is reversibly activated by ADP. On incubation of AMPS-BDB with glutamate dehydrogenase, the enzyme reacts covalently, resulting in an irreversibly *activated* enzyme that is no longer responsive to externally added ADP (6). AMPS-BDB appears to function as an ADP substitute that is covalently bound to Arg<sup>459</sup> within the activator site of the allosteric bovine liver glutamate dehydrogenase (6).

The above are representative examples of enzymes that have been studied using the strategy of affinity labeling. They illustrate how affinity labeling, X-ray crystallography and site-directed mutagenesis can be used as complementary approaches in evaluating the functional role of particular amino acids in a protein.

### Bibliography

1. R. F. Colman, Y.-C. Huang, M. M. King, and M. Erb (1984) *Biochemistry* **23**, 3281–3286.
2. T. T. Lee, C. Worby, J. E. Dixon, and R. F. Colman (1997) *J. Biol. Chem.* **272**, 458–465.
3. D. H. Ozturk, I. Park, and R. F. Colman (1992) *Biochemistry* **31**, 10544–10555.
4. S. H. Vollmer, M. B. Walner, K. V. Tarbell, and R. F. Colman (1994) *J. Biol. Chem.* **269**, 8082–8090.
5. O. A. Moe et al. (1996) *Biochemistry* **35**, 9024–9033.
6. K. O. Wrzeszczynski and R. F. Colman (1994) *Biochemistry* **33**, 11544–11553.

### Suggestions for Further Reading

7. R. F. Colman (1983) Affinity labeling of purine nucleotide sites in proteins. *Ann. Rev. Biochem.* **52**, 67–91.
8. R. F. Colman (1990) Site-specific modification of enzyme sites. In *The Enzymes*, 3rd ed., Vol. **19** (D. S. Sigman and P. D. Boyer, eds.), Academic Press, San Diego, Calif., pp. 285–321.
9. R. F. Colman (1997) Affinity labels for NAD(P)-specific sites. *Methods Enzymol.* **280**, 186–203.
10. R. F. Colman (1997) Affinity labelling. In *Protein Function: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.) Oxford University Press, Oxford UK, pp. 155–183.
11. W. B. Jakoby and M. Wilchek, eds. (1977) Affinity labeling. *Methods Enzymol.* **46**.



## Affinity Maturation

It was observed long ago that the affinities of [antibodies](#) increase with time during [immunization](#) with a T-dependent [antigen](#). When such an antigen is injected for the first time, it induces the occurrence of low-affinity [IgM](#) antibodies, which are rapidly replaced by [IgG](#). This takes place in the first two weeks of the primary response. At the same time, an increase of the average affinity is usually observed. If the same antigen is given a second time, [IgG](#) antibodies are subsequently produced at a higher yield, and with an average affinity that is still increasing. This is characteristic of the secondary response, a phenomenon that may be repeated and amplified by multiple administrations of the antigen. This is why long-term immunization schedules are used whenever both high titers and high affinity are wanted. Affinity maturation, which is a unique property of [B cells](#) stimulated by T-dependent antigens, is the consequence of (1) the clonal organization of lymphocytes, (2) the expansion of clones that have been stimulated specifically by the antigen, and (3) the presence of a very peculiar mechanism that generates [somatic mutations](#) at each cellular division of stimulated B cells.

1. Antibody specificities are spread randomly in the huge collection of clones that represent at any given time the [repertoire](#) expressed by B lymphocytes that have rearranged their immunoglobulin genes and thus have [IgM](#) at their surface. Clones that are potentially reactive to the antigen proliferate, as the result of the cooperation between [T cells](#) and [B cells](#) and the phenomenon of **antigen processing** and presentation. This ensures that the first wave of antibodies are generally of low or moderate affinity.
2. As the immunization proceeds, expansion of the B-reactive cell population takes place, giving rapidly growing colonies, termed *germinal centers*, in secondary lymphoid organs (lymph nodes, tonsils, spleen). Germinal centers are highly organized structures that favor cell–cell interactions that will amplify local B-cell division and expansion.
3. At the same time, somatic mutations will be triggered by a mechanism that is **probably** enzyme-driven, but still not completely elucidated. The rate of the mutations that are introduced is of the order of  $1 \times 10^3$  per base pair and per generation, which is excessively high. As a result of these mutations, some clones will gain affinity, whereas others will lose. One therefore needs mechanisms to select clones of high affinity and eliminate those of lower affinity; otherwise the system would expand indefinitely and explode. Stimulation by antigen will favor the saving and expansion of clones with the highest affinity, whereas [apoptosis](#) will operate on cells that did not receive stimuli as they lost recognition ability.

**PCR** and [DNA sequencing](#) at the level of single cells allowed the group of Rajewsky to trace the successive mutations that occurred at each cell division, so that it was possible to reconstruct a genealogical tree of these mutations within the germinal center. Interestingly, hypermutated cells can ultimately follow two types of [differentiation](#): One leads to the terminal plasma cells, which will secrete circulating antibodies, and the other leads to the so-called [memory cells](#) that may persist for long periods of time and may be restimulated with great efficiency in the course of another run of immunization.

The extraordinary plasticity and adaptability of the immune system are thus quite remarkable and allow it to respond efficiently to pathogens. Raising memory cells endowed with the ability to produce rapidly high-affinity antibodies is an obvious goal of vaccination.

See also entries [Immune Response](#), [Class Switching](#), and [Somatic Hypermutations](#).

### Suggestions for Further Reading

- D. M. Tarlinton, A. Light, G. J. Nossal, and K. G. Smith (1998) Affinity maturation of the primary response by V gene diversification. *Curr. Top. Microbiol Immunol.* **229**, 71–83.
- C. J. Jolly, S. D. Wagner, C. Rada, N. Klix, C. Milstein, and M. S. Neuberger (1996) The targeting of somatic hypermutation. *Semin. Immunol.*, **8**, 159–168.
- A. Ehlich, V. Martin, W. Muller, and K. Rajewsky (1994) Analysis of the B-cell compartment at the level of single cells. *Curr. Biol.* **4**, 573–583.

## Affinity Selection

A key requirement of [combinatorial library](#) approaches is that desirable molecules must be segregated from the remaining library population. Segregation is most often accomplished by affinity partitioning, in which an immobilized target is used to capture interacting molecules from a solution-phase library. Target proteins can be immobilized either using [antibodies](#) bound to *Staphylococcus* [protein A](#), using absorption of targets onto wells of plastic microtiter plates, or covalent attachment to a variety of polymeric supports. Alternatively, combinatorial libraries attached to solid supports—for example, one-bead, one-compound (OBOC) libraries—can be used to bind soluble targets. Solid-phase libraries have also been used to screen highly complex targets, such as living cells on which a variety of receptors are expressed ([1](#)). In all cases, unbound material is removed by washing, and the bound material is recovered for amplification or identification, depending on the nature of the library.

Plastic pins were used to develop the first peptide-based combinatorial libraries. The pins are arranged such that they fit neatly into a single well of a microtiter plate. Thus, the pin system provides a convenient format for both synthesis and screening. The power of this method derives largely in the ease of handling large numbers of discrete syntheses in parallel. Sequential steps of the synthesis can be carried out by transferring the pin arrays through various microtiter reaction chambers. The peptides are then directly available for affinity screening, either on the solid phase or following cleavage from the solid support.

Lam et al. ([2](#)) developed another approach for peptide affinity selection involving split synthesis of peptides on solid support beads (OBOC libraries: see [Combinatorial Synthesis](#)). From a pool of millions of beads, a binding reaction is performed using a target molecule such as an [antibody](#), receptor, [enzyme](#), or even whole cells. Beads displaying affinity for the target are isolated, and the peptide on the bead can be microsequenced.

Fodor et al. ([3](#)) have developed immobilized combinatorial libraries on silicon microchips. Chip-based addressable libraries of peptides, oligonucleotides, and small organic molecules can be readily prepared. These methods use photolithography to control regions accessible for subsequent chemical modification. This method enables a miniaturized, fully addressable library to be generated on the surface of a silicon chip. The resulting arrays can be screened using standard affinity methods.

The affinity selection methods described above are extremely broad in their applications. Although classical nucleic acid libraries have been screened by **hybridization**, which is a highly specialized type of affinity partitioning, modern DNA and RNA aptamer libraries are screened in identical fashion to other combinatorial libraries by affinity selection over immobilized targets. Although

selection strategies have been introduced that are based on properties other than affinity, such methods are highly specialized and are not likely to displace the current reliance on affinity methods for combinatorial screening.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

### Bibliography

1. M. E. Pennington, K. S. Lam, and A. E. Cress (1996) *Mol. Diversity* **2**, 19–28.
2. K. S. Lam, S. E. Salmon, E. M. Hersh, V. J. Hruby, W. M. Kazmierski, and R. J. Knapp (1991) *Nature* **354**, 82–84.
3. S. P. Fodor, J. L. Read, M. C. Pirrung, L. Stryer, A. T. Lu, and D. Solas (1991) *Science* **251**, 767–773.

## Agarose

Agarose is one of the two most popular materials used to prepare gels for use in [gel electrophoresis](#), the other being [polyacrylamide](#). Agarose is primarily a polymer of molecular weight approximately 120 kDa of agarobiose (an anhydrogalactose-galactose disaccharide). At sufficiently high concentrations and low temperatures, the polymer forms b-helical strands, which interact by **hydrogen bonding** to form the agarose gel; they dissociate at higher temperatures to form agarose solutions. Such thermally controlled gelation is the unique feature of agarose, which provides both a high reproducibility and simplicity of preparing the gel. Since gelation requires noncovalent interactions between the copolymer strands, denaturing agents, such as solutions of [urea](#), prevent gelation.

Agarose is derived from a natural product of oceanic algae, agar. It is processed to reduce the content of acidic groups, primarily sulfate, so as to reduce the degree of [electroendosmosis](#) that occurs during electrophoresis. Further industrial processing produces a large number of commercially available agarose fractions distinguished by (1) the degree of electroendosmosis in an electric field, determined by the number of acidic groups; (2) the degree of covalent substitution of the acidic groups with chemical groupings such as hydroxyethyl or vinyl (allyl) groups, which tends to lower the melting point in proportion to the degree of substitution; (3) the addition of linear carbohydrate polymers, such as clarified locust bean gum, to increase the viscosity of the gel and thereby eliminate measurable electroendosmosis.

Most agaroses at concentrations greater than about 0.4% (w/v) form gels at room temperature. Agaroses with unusually high gel strengths can gel at concentrations as low as about 0.05%. The pore sizes of such agarose gels make them ideal for separating **DNA** molecules and **oligonucleotides**; their mobility is determined primarily by their size, as nucleic acids of the same class have the same charge density per nucleotide. The least concentrated gels exhibit mean pore sizes of up to 500 to 1,000 nm, sufficiently large for the penetration of large particles or small viruses (see [Particle Electrophoresis](#)). Highly soluble agarose species substituted to a maximum extent by hydroxyethyl groups can be used to prepare gels with concentrations of up to 9 to 10%. Such gels have pores of similar size to those of 3 to 6% polyacrylamide gels, and with greater resolving power for native proteins. Agarose can also be used as a copolymer with polyacrylamide, which provides larger pore sizes than polyacrylamide by itself. The copolymer is produced by adding agarose to the mixture of acrylamide monomers before its polymerization. For the separation

by size of DNA fragments, or of other large particles, by [capillary zone electrophoresis](#), solutions of agarose can be used.

The adherence of agarose gels to glass or plastic apparatus walls is very weak, and only horizontally oriented gel electrophoresis apparatus can normally be used. The adherence of agarose to vertical glass surfaces may, however, be strengthened by drying a thin film of agarose onto glass walls; this can permit electrophoresis in vertical glass slabs or tubes. Horizontal thin-layer agarose gels with high Joule heat dissipation capacity, which are therefore amenable to electrophoresis at high field strength without melting, can be formed on hydrophilic surfaces of thin plastic sheets (eg, “Gel-Bond”).

For preparative purposes, bands of a macromolecular sample in agarose gels may be recovered by solubilizing the agarose at increased temperature or adding the enzyme agarase. The macromolecule of interest is subsequently separated from the low-molecular-weight agarose fragments by filtration or precipitation methods.

#### Suggestions for Further Reading

- I. C. M. Dea, A. A. McKinnon, and D. A. Rees (1972) Tertiary and quaternary structure in aqueous polyacrylamide systems which model cell wall cohesion: Reversible changes in conformation and association of agarose, carrageenan and galactomannans. *J. Mol. Biol.* **68**, 153–172.
- T. G. L. Hickson and A. Polson (1968) Some physical characteristics of the agarose molecule. *Biochim. Biophys. Acta* **165**, 43–58.
- FMC Corporation (1982) *The Agarose Monograph*, Bioproducts Dept., Rockland, ME.
- J. O. Jeppson, C. B. Laurell, and B. Franzen (1979) Agarose gel electrophoresis. *Clin. Chem.* **25**, 629–638.
- P. Serwer (1983) Agarose gels: Properties and use for electrophoresis. *Electrophoresis* **4**, 375–382.
- D. Tietz (1987) Gel electrophoresis of intact subcellular particles. *J. Chromatogr.* **418**, 305–344.

## Agglutination

“Agglutinin” is an obsolete term for an [antibody](#) capable of causing cells to come together in macroscopic clumps. Agglutination (clumping) of bacteria by serum was perhaps the earliest observed manifestation of *in vitro* immune reactions. The phenomenon is due to crosslinking of cells through the reaction of multivalent antibody molecules with cell-surface [antigens](#). The quantity of antibody necessary for crosslinking is generally very small, yet the cell masses involved in agglutination generate a response visible to the naked eye. Because of these characteristics, this oldest of techniques remains one of the most sensitive and widely used types of [immunoassay](#).

Agglutination immunoassays can employ cells or synthetic particles. The two broad categories of cell-based agglutination assay are termed “direct” and “passive”. In the direct assay, antibodies cause clumping of microbes or blood cells by reaction with surface antigens endogenous to those cells. Typically, the immune status of serum in a test sample will be assessed, based on the highest dilution that will still agglutinate target cells. Direct agglutination assays are widely used to diagnose infectious disease by the presence of specific antibodies in acute or convalescent serum. The “monospot” assay for infectious mononucleosis is an example of this application ([1](#)). In passive agglutination immunoassays, the target antigen is exogenous, and is coated covalently or by

adsorption onto the cells used for the assay. Test samples are then assayed as for the direct agglutination method. Many variants of direct and passive agglutination are in use for research or diagnostic procedures. For example, the test for Rh blood type is a two-step method in which antibodies in the test sample react with endogenous Rh antigen on erythrocytes, then addition of anti-human [IgG](#) induces agglutination (2).

Availability of uniformly sized latex particles allowed development of a chemically defined substrate on which to attach antigen or antibody (3). The role of the particles was precisely analogous to that of cells, in that their aggregation in response to antibody–antigen crosslinking generated an optically detectable signal. The advantage of inert particles was in their stability, compared to cells, and in interassay reproducibility. Measurement of agglutination was at first by turbidometry. Development of laser [light-scattering](#) techniques have greatly improved the sensitivity of detection (4) and have led to development of automated instruments for quantitation of agglutination reactions.

### Bibliography

1. C. L. Lee, I. Davidsohn, and R. Slaby (1968) *Am. J. Clin. Pathol.* **49**, 3–11.
2. R. R. A. Coombs, A. E. Mourant, and R. R. Race (1945) *Br. J. Exp. Pathol.* **26**, 255–266.
3. J. M. Singer and C. M. Plotz (1956) *Am. J. Med.* **21**, 888–892.
4. V. Schulthess, R. J. Cohen, N. Sakato, and G. B. Benedek (1974) *Immunochemistry* **13**, 955–962.

## Agrobacterium

*Agrobacterium tumefaciens* and *Agrobacterium rhizogenes* are members of the *Agrobacterium* genus of soil **bacteria** responsible, respectively, for tumor formation and hairy root disease in dicotyledonous **plants**. The molecular basis of inducing neoplastic growth in plants has been one of the most intriguing areas of research in plant pathology for more than 50 years. Advances in understanding the molecular basis of tumor formation has led to the development of **vector** systems based on *Agrobacterium* to carry out [plant genetic engineering](#).

*Agrobacterium* belongs to the *Rhizobiaceae* family. Classification of the bacterium has been based on physiological and biochemical criteria, although increasingly the taxonomic structure of the genus has been based on **DNA** analysis of **chromosomal** groups of the differing biovars (1). The best studied members of the genus are *A. tumefaciens* and *A. rhizogenes*. *A. tumefaciens* is the causative agent of crown gall disease. It has a wide, obvious host range, including the majority of dicotyledonous plants and some monocotyledonous plants (2). Practically, crown gall can be a serious disease of fruit crops, including grape vines. The physical basis of crown gall disease is the transfer of a defined region of DNA, transferred or T-DNA, from a **plasmid** maintained in the bacteria, known as the Ti or tumor-inducing plasmid (see [Ti Plasmid](#)) to the [genome](#) of the infected plant cell. *A. rhizogenes* induces hairy root disease analogously. A T-DNA is transferred from the Ri, or root-inducing, plasmid to the genome of the infected plant cell (see [Ri Plasmid](#)). This natural form of plant cell [transformation](#) has been adapted and used extensively in creating **transgenic** plants (see [Plant Genetic Engineering](#)).

### Bibliography

1. H. Bouzar, J. B. Jones, and A. Bishop (1995) In *Agrobacterium Protocols* (K. M. A. Gartland and M. R. Davey, eds.), Humana Press, Totowa, NJ, pp. 9–13.

2. M. DeCleene and J. DeLey (1976) *Bot. Rev.* **113**, 81–89.

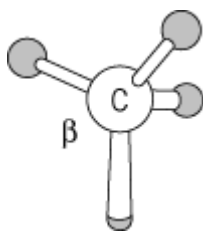
### Suggestions for Further Reading

3. K. Kersters and J. De Ley (1984) "Genus III *Agrobacterium* Conn 1943", In *Bergey's Manual of Systematic Bacteriology*, Vol 1 (N. R. Krieg and J. G. Holt, eds.), Williams and Wilkins, Baltimore, pp. 244–254.
4. G. Kahl and J. Schell (1982) *Molecular Biology of Plant Tumors*, Academic Press, London.

## Alanine (Ala, A)

The [amino acid](#) alanine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to four **codons**—GCU, GCC, GCA, and GCG—and represents approximately 8.3% of the residues of the proteins that have been characterized. The alanyl residue incorporated has a mass of 71.09 Da, a **van der Waals volume** of  $67 \text{ \AA}^3$ , and an [accessible surface area](#) of  $113 \text{ \AA}^2$ . Ala residues are usually relatively variable during [divergent evolution](#), as they are frequently interchanged in **homologous** proteins with [serine](#), [threonine](#), [valine](#), [glutamic acid](#) and [proline](#) residues.

The Ala side chain is simply a methyl group:



This [nonpolar](#) side chain makes Ala residues unreactive chemically, relatively **hydrophobic**, and not very [hydrophilic](#); consequently, 38% of the Ala residues are fully buried in the folded conformations of proteins. There, the methyl side chain undergoes rapid rotations about the  $C_a-C_b$  single bond. Ala has the greatest tendency of all the normal amino acid residues to adopt the **alpha-helical** conformation in model **peptides**. It occurs frequently in that conformation in folded [protein structures](#), but at about only half that frequency in [beta-sheets](#) and in reverse [turns](#).

### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Albumins

Albumins are operationally defined as [proteins](#) that (1) remain soluble in pure [water](#) after dialysis of protein samples such as egg white or blood serum against distilled water and (2) are not precipitated in 50% saturated ammonium **sulfate**. This contrasts with another group of proteins, called [globulins](#), which are precipitated in both distilled water and 50% saturated ammonium sulfate. This operational classification is rather obsolete today, but names like [serum albumin](#), [ovalbumin](#), and lactalbumin have remained. Other albumins include muscle albumin and plant albumins (1).

## 1. Blood Albumins

[Serum albumin](#) and other plasma proteins have been most closely studied, mainly for their medical interest. After the cells have been removed from blood by light [centrifugation](#) in the presence of chelating agents such as [EDTA](#), a yellowish fluid called plasma remains. By the addition of  $\text{Ca}^{++}$ , [fibrinogen](#) is converted to fibrin and the clotting process produces a soft gel. Removal of this gel by centrifugation or other methods will leave a clear fluid called serum. Serum thus obtained contains 70 to 80 g/L of total protein, consisting of more than 150 different kinds of proteins. Dialysis of serum against distilled water will precipitate the globulins, while albumins stay in solution. The albumin fraction contains serum albumin as the major protein (2, 3). (see [Serum Albumin](#)).

## 2. Egg Albumins

Egg proteins are first classified into yolk proteins and egg white proteins. Yolk contains proteins called vitellogenin, phosvitin, and lipovitellin, all in association with yolk lipid. Egg white is a reservoir of several kinds of proteins as a protective and nutritious environment for the fetus. Egg white proteins are further divided into albumins and globulins, as for serum proteins. The following is a list of egg white proteins that are not globulins. Descriptions of egg white globulins are found under [Globulins](#).

### 2.1. Ovalbumin

This is the major **glycoprotein** in the white of eggs, comprising 65% of the total egg white protein, with a molecular weight of 43,000 and [isoelectric point](#) (pI) of 4.7. The N-terminal glycine residue is acetylated. Two species containing one or two sites of **phosphorylation** can be separated by [electrophoresis](#). Oligosaccharide containing three moles of N-acetylglucosamine and five moles of mannose is linked to an [asparagine](#) residue through an **N-glycosidic** linkage. Treatment of ovalbumin with subtilisin yields a readily crystallizable plakalbumin, named after the platelike appearance of the resulting crystals. Ovalbumin is often used as an effective antigen in immunological studies. The protein has recently been found to have a similar amino acid sequence and three-dimensional structure to [a1-antitrypsin](#) of the [serpin](#) family (4) (see [Ovalbumin](#)).

### 2.2. Ovotransferrin (or *conalbumin*)

This protein binds ferric ions and is the same as apo-[transferrin](#), but differing in the carbohydrate moiety. In egg white, the protein is present in an almost iron-free form. This protein constitutes about 10% of the egg protein.

### 2.3. Ovomuroid

Ovomucoid is a glycoprotein (carbohydrate content about 25% by weight) with a molecular weight of 28,000 and constitutes 1.5% of the total egg white protein. It inhibits **trypsin** and **chymotrypsin** but not **plasmin**, [thrombin](#), [elastase](#), or collagenase. Its pI falls in the range 3.9 to 4.5. It remains active in the supernatant after egg white has been coagulated by heat. It has three **homologous** domains connected by linking peptides, each domain being homologous to [bovine pancreatic trypsin inhibitor](#) (BPTI). Thus, it is speculated that ovomucoid has evolved by two tandem repeats of the BPTI gene. Carbohydrates are linked to asparagine residues at residue numbers 10, 53, 69, 75, and 175. Similar inhibitors can be purified from the egg white of the quail, goose, and turkey.

## 2.4. Ovoinhibitor

This is a 48-kDa multiheaded [proteinase inhibitor](#) that can simultaneously inhibit trypsin, chymotrypsin, and elastase. Unlike ovomucoid, it also inhibits the proteinases of bacterial origins. It is a glycoprotein with 5 to 10% content of carbohydrate.

## 2.5. Avidin

This protein occupies a special position in biochemistry in that it is widely used as a probe based on its strong affinity for [biotin](#). Their binding constant reaches  $10^{14} M^{-1}$  under optimal conditions. [Avidin-biotin systems](#) are widely used for labeling macromolecular interacting systems, not only [protein-protein interactions](#), but also systems based on DNA and other macromolecules. [Avidin](#) can be substituted by [streptavidin](#), which is of bacterial origin but with a similar activity. Avidin in the egg white is largely free of biotin, and feeding rats with purified avidin as the sole protein source will cause various symptoms known to accompany biotin (vitamin H) deficiency.

## 2.6. Ovomucin

This is a glycoprotein of high molecular weight responsible for the mucous character of egg white. The treatment of ovomucin with reagents that react with [thiol groups](#) reduces the viscous nature of egg white.

## 3. Plant Albumin

**Ricin** is an example of plant albumins (from *Ricinus communis*) and has an interesting toxic function. It consists of an A subunit, of 32 kDa, pI of 7.5, and 2.4% carbohydrate, and a B subunit, of 34 kDa, pI of 4.8, and 6.5% carbohydrate. It inactivates the 60S subunit of [ribosomes](#) after being internalized into the cell and thus inhibits a peptide elongation step (see [Translation](#)). The A subunit carries out the inactivation reaction, while the B subunit, which binds to cell-surface galactose residues, helps the A subunit to be internalized into the cell. Plant albumins are more readily precipitable in 50% ammonium sulfate than animal albumins.

## 4. Milk Albumins

### 4.1. $\alpha$ -Lactalbumin

This is a single-chain 14-kDa protein that is the same as the B chain of lactose synthetase in lactating granules. The A chain of lactose synthetase alone catalyzes the synthesis of N-acetylgalactosamine from N-acetylglucosamine and UDP-galactose, but when in association with  $\alpha$ -lactalbumin, it uses glucose as a substrate and synthesizes lactose. The amino acid sequence of  $\alpha$ -lactalbumin with 143 amino acids (human, bovine, guinea pig) is homologous to that of [lysozyme](#), and its three-dimensional structure is also similar. A recent review of the relationship between lysozyme and lactalbumin is recommended ([5](#)) (see also [Alpha-Lactalbumin](#) and [Lysozymes](#)).

### 4.2. Lactoferrin

This is an iron-binding protein of 88,000 in humans and 86,000 in cows. It is similar in functional properties to serum transferrin, but immunologically and structurally distinguishable ([6](#)).

## Bibliography

1. B. Blombäck and L. A. Hanson (1979) *Plasma Proteins*, John Wiley & Sons, New York.
2. M. Perutz (1992) *Protein Structure: New Approach to Disease and Therapy*, W. H. Freeman and Co., New York.
3. T. Peters Jr. (1985) *Adv. Prot. Chem.* **37**, 161–245.
4. P. E. Stein, G. W. Leslie, J. T. Finch, and R. W. Carrell (1991) Crystal structure of uncleaved ovalbumin at 1.95 Å resolution, *J. Mol. Biol.* **221**, 942–959.
5. M. A. McKenzie and F. H. White Jr. (1991) Lysozyme and  $\alpha$ -Lactalbumin: Structure, function, and interrelationships, *Adv. Prot. Chem.* **41**, 173–315.



6. B. F. Anderson, H. M. Baker, E. J. Dodson, G. E. Harris, S. V. Rumball, J. M. Waters, and E. M. Baker (1987) Structure of lactoferrin at 3.2 Å resolution, Proc. Nat. Acad. Sci. USA **84**, 1769–1773.

## Alcohol Dehydrogenase (ADH)

ADH is a type of [enzyme](#) that catalyzes the reversible interconversion of an alcohol to an aldehyde/ketone [EC 1.1.1.1]. The substrate specificity is often wide, but it typically includes at least some activity with ethanol. The **coenzyme** is often  $\text{NAD}^+$ . ADH was initially purified from yeast (YADH) and from horse liver (LADH), sources from which it is commercially available. It is, however, now known to be widespread, with at least some forms present in **eukaryotes**, **prokaryotes** and **archaeobacteria** in general. It has been well studied at the protein and **DNA** levels from many sources. The “classical” forms of ADH are Zn-containing [metalloproteins](#) with subunits of ~40 kDa (~350 to 390 residues) (1). They are now considered part of a larger protein family (see [Protein Evolution](#)), known as medium-chain dehydrogenases/reductases (**MDR**) (2).

Although ADH is frequently regarded as well understood, much is still unknown, and novel findings of general interest are constantly obtained:

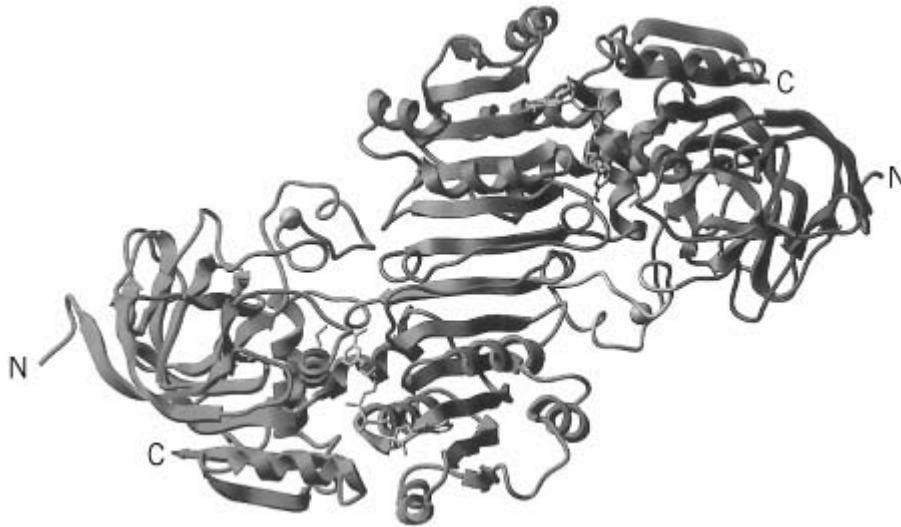
1. ADH is no longer just an enzyme, but a whole system of great complexity, involving other protein families and a multiplicity of MDR forms.
2. another such family is (short-chain dehydrogenases/reductases) [SDR](#) (3). The major ADH of insects is this type, and SDR enzymes are numerous, including **hormone**-converting and regulatory enzymes in humans and all living forms.
3. MDR-ADH occur as different classes (4) and **isozymes**. They exhibit distinct evolutionary patterns, repeated evolutionary origins, and varying properties. At least six classes and eight genes have been discerned in mammals, but these numbers are likely to grow; they differ among species because of recent [gene duplications](#). Expressions differ with organ, tissue, and age, including special fetal forms.
4. Functionally, the ADH system has long been an enigma, but recent results suggest that it has important roles. Metabolically, ADH is associated with aldehyde dehydrogenase, which is part of another large family. Medically, both activities are important in alcohol metabolism and in several disease states but also in vision and in cellular regulation and differentiation (5, 6).

### 1. MDR-ADH

#### 1.1. Classes in Vertebrates

Mammals and most vertebrates contain different classes of MDR Zn-containing ADH enzymes. These classes differ substantially in substrate specificity and in their immunological, chromatographic, and electrophoretic properties (4). They constitute separate forms that are intermediate in their properties between isozymes and distinctly different enzymes, and they differ in primary structures by about 30% but have closely related conformations (Fig. 1). The classes reflect a series of gene duplications, which occurred largely at early vertebrate times (8). Later duplications in individual lines have given rise to intraclass isozymes. The subunits of isozymes can hybridize into mixed dimers that reflect the subunit mixture, but those of different classes do not.

**Figure 1.** Structure of human ADH class I (7). The diagram and the modeling of ethanol were made with the program ICM (Molsoft, Metuchen, NJ) from coordinates in the data bank.



Class I contains the classical ADH enzyme, which is present in large amounts in liver but in lower amounts in many other organs. It is involved in ethanol metabolism and in liver function. It has good enzymatic activity with primary alcohols, cyclohexanol, retinols, and many other alcohols, plus their corresponding aldehydes. Isozymes occur and are largely species specific. The horse ADH subunits E and S derive from two separate genes and the human subunits a, b, and g from three genes. The a subunit is expressed in the fetus and b and g in the adult. The b and g types also exhibit further variability from alleles that have different population distributions in Caucasians and Orientals. This population variation, together with other alleles of the next enzyme of alcohol metabolism (aldehyde dehydrogenase), explains the different sensitivity to ethanol consumption of Caucasians and Orientals.

Class II ADH is highly variable. It is also expressed in the liver but in lower amounts. Thus far it is without a recognized functional role, but in humans it has a higher  $K_m$  for ethanol, and hence is of little importance in ethanol metabolism.

Class III ADH is expressed universally in cells and organs, and is identical to [glutathione](#) (GSH)-dependent formaldehyde dehydrogenase (EC 1.2.1.1). In the absence of GSH, it is a dehydrogenase toward long-chain alcohols and is largely inactive with ethanol. The  $K_m$  for ethanol is  $>3 M$ , so there is activity only at high ethanol concentrations. In the presence of GSH, class III ADH is active with the GSH/formaldehyde adduct, hydroxymethyl-GSH, when it functions as a formaldehyde dehydrogenase, producing formic acid.

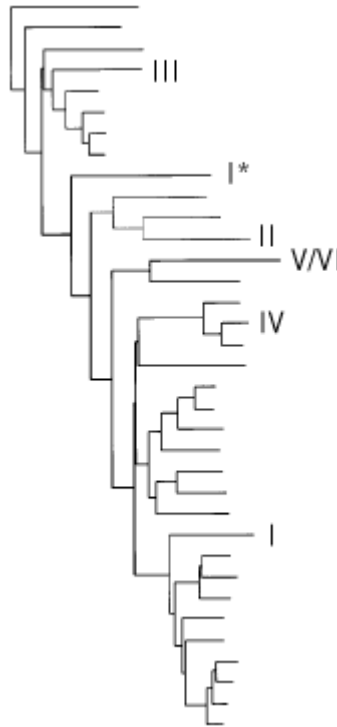
Class IV ADH is present in the stomach and in skin (epithelia). It is the most ethanol-active form and has been discussed in relation to first-pass ethanol metabolism and to retinol dehydrogenase function. Both roles, however, are far from established.

Little is known about the remaining enzyme classes in vertebrates.

Class III is the evolutionary ancestor, which is apparent from [phylogenetic trees](#) (Fig. 2), the rates of change in present forms, and in apparently being the only MDR-ADH in invertebrates. However, final judgment of the phylogeny must await characterization of further forms. Additional gene duplications probably remain to be discovered, and the true evolutionary relationships may be

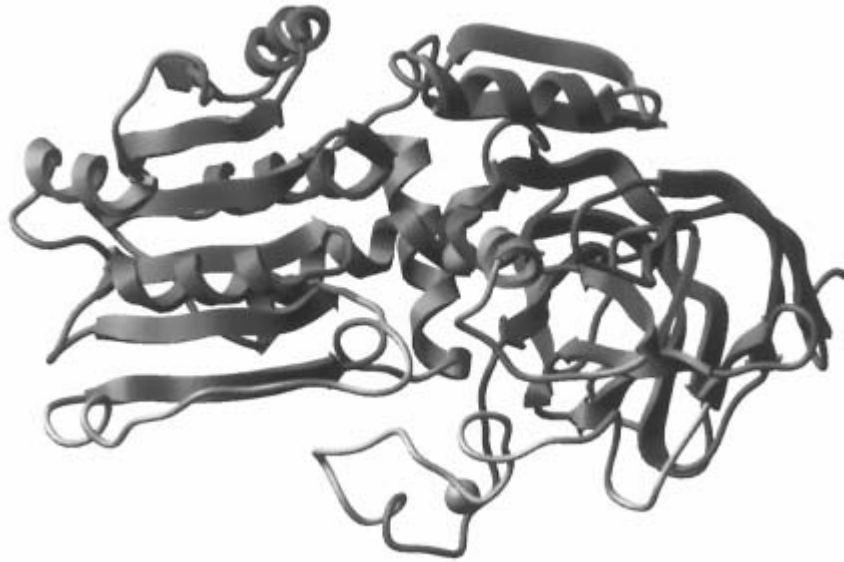
complex, involving both extant and extinct forms.

**Figure 2.** Evolutionary tree of animal ADH. The deviating class I form of bony fish is indicated by a star. The tree was calculated using ClustalW (2) with corrections for multiple substitutions.



In molecular terms, the classes exhibit different evolutionary patterns. Class III is like an enzyme of basic metabolism, and has relatively constant functional properties among species and variation primarily in the nonfunctional parts of the protein (Fig. 3). Class I is the result of an early gene duplication from class III, whereas class IV in turn diverged from the class I line. Special hybrid forms of class I/III (in fish), II/I (in birds) and IV/I (in amphibians) have been traced and may reflect either the phylogeny or simply additional gene duplications. In short, the ADH system has many interesting evolutionary features (8).

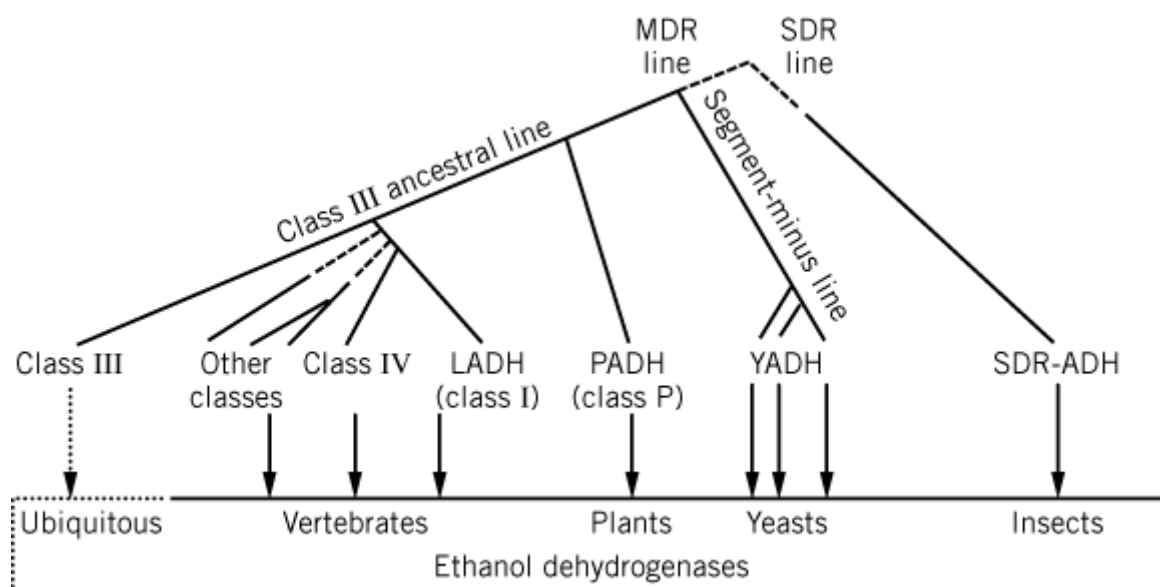
**Figure 3.** Variable segments in subunits of ADH. Those subunits that vary between class I enzymes affect functional segments, typical of a functionally variable protein. Those subunits variable between class III enzymes affect superficial segments, typical of a functionally constant protein (3).



### 1.2. Yeast, Plants, Other Nonvertebrate Forms

Nonvertebrate MDR-ADH are also frequent and have at least two classes. One is the class III enzyme, which has GSH-dependent formaldehyde dehydrogenase activity, as expected from the ancestral properties noted previously. The other is ethanol-active MDR-ADH: yeast ADH (YADH, from at least three genes) and plant ADH (PADH). YADH has greater ADH activity than LADH because of weaker **coenzyme** binding and consequently a greater turnover rate, because coenzyme dissociation is the rate-limiting step. Although YADH and PADH are related to LADH structurally and functionally, they are products of gene duplications separate from that leading to class I ADH. Combined, all of these ethanol-active forms suggest that ADHs active with alcohol have arisen several times during evolution, probably by functional **convergence**. Presumably all or many of these ethanol-active forms are derived from the ancestral class III line (Fig. 4), but they can differ in [quaternary structure](#). Ethanol-active YADH is a tetramer, whereas yeast class III ADH is a dimer.

**Figure 4.** Functional convergence toward ethanol dehydrogenases in different lines of living organisms, emphasizing the repeated appearance of ethanol-activity, many separate gene duplications, and an apparent MDR-ancestral nature of class III. Dashed lines indicate unknown relationships, and dotted lines low activity with ethanol. Segment-minus forms lack internal segments (Fig. 1) and are tetrameric. The remaining forms are dimers. Apart from the gene duplications shown, which lead to ethanol-active dehydrogenases, additional branchings lead to cinnamyl ADH (in plants), sorbitol dehydrogenase (in animals), threonine dehydrogenase (in prokaryotes), several reductases (in prokaryotes and eukaryotes), and additional enzymes.



### 1.3. Function

Liver ADH (class I) is the major metabolic enzyme for ethanol, either ingested or produced in the intestine, and for other alcohols. Overall, the formaldehyde specificity of the ancestral form (class III), the multiplicity of all ADH systems, and the gradual changes in substrate specificity suggest that the whole ADH system is a basic enzyme, together with aldehyde dehydrogenase, in cellular defense reactions against alcohols and aldehydes, converting them to acidic end products. The validity of this general role of ADH is supported by the fact that the only animals apparently lacking ADH (except for formaldehyde-active class III) are marine invertebrates (10), whose environment is low in alcohol and aldehyde levels. The formaldehyde dehydrogenase activity of class III ADH is universal and constant throughout the living world, although the  $K_m$  and enzymatic efficiencies are increased in microorganisms. This permits life at higher aldehyde concentrations (11) and suggests a general defense function for the ADH system.

Additional functions of individual classes are likely. Class IV has been discussed relative to its possible role in retinal formation important for [rhodopsin](#) and vision (6) and in [retinoic acid](#) formation, which important for vertebrate growth and [differentiation](#) (5). Similarly, special forms of ADH and aldehyde dehydrogenase are involved in the fatty acid cycle in dermal wax and plasmalogen formation in the central nervous system and are associated with an inborn error of metabolism (Sjögren-Larsson syndrome) affecting the corresponding genes (12).

### 1.4. Other MDR Enzymes

Many enzymes active on other alcohols or polyols are related to MDR-ADH in structure and origin. Among these are sorbitol dehydrogenase, threonine dehydrogenase, quinone oxidoreductase, the enoyl reductase component of the fatty acid synthesis machinery, and enzymes with undefined functions, like VAT-1 of synaptic membranes (2). These enzymes are similar conformationally and in their overall active-site relationships. But they also differ markedly, including functionally, in that all do not have the active-site zinc typical of MDR-ADH. This suggests that they have different enzymatic mechanisms. Zinc is also absent in quinone oxidoreductase, which is a crystallin in ocular lenses of some animals.

## 2. SDR-ADH

Insect ADH is the only well known SDR-ADH, but hundreds of other SDR enzymes are known. These enzymes work on hydroxyl groups of a multitude of sugar, steroid, prostaglandin, and other

compounds. They also include reductases, with other activities (like double-bond saturation), and enzymes of no less than three of the enzyme classes in general (oxidoreductases, lyases, and isomerases (3)).

### 3. Other ADH Forms

ADH activity is also known in enzymes other than the MDR and SDR forms, including iron-activated forms (13). Prokaryotes have several different forms, including long-chain alcohol dehydrogenases (14) and methanol dehydrogenases. They may even have MDR relationships, and this definitely applies to other special prokaryotic forms, such as factor-dependent formaldehyde DH (15), which is apparently the equivalent of class III MDR-ADH in **Gram-positive** bacteria.

### Bibliography

1. C.-I. Brändén, H. Jörnvall, H. Eklund, and B. Furugren (1975) *The Enzymes*, 3rd ed., chap. "11", pp. 103–190.
2. B. Persson, J. S. Zigler Jr., and H. Jörnvall (1994) *Eur. J. Biochem.* **226**, 15–22.
3. H. Jörnvall et al. (1995) *Biochemistry* **34**, 6003–6013.
4. B. L. Vallee and T. J. Bazzone (1983) *Curr. Top. Biol. Med. Res.* **8**, 219–244.
5. G. Duyster (1996) *Biochemistry* **35**, 12221–12227.
6. A. Simon et al. (1995) *J. Biol. Chem.* **270**, 1107–1112.
7. T. D. Hurley et al. (1994) *J. Mol. Biol.* **239**, 415–429.
8. O. Danielsson et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 4980–4984.
9. J. D. Thompson, D. G. Higgins, and T. J. Gibson (1994) *Nucleic Acids Res.* **22**, 4673–4680.
10. M. R. Fernández et al. (1993) *FEBS Lett.* **328**, 235–238.
11. M. R. Fernández et al. (1995) *FEBS Lett.* **370**, 23–26.
12. V. De Laurenzi et al. (1996) *Nature Genet.*, **12**, 52–57.
13. R. K. Scopes (1983) *FEBS Lett.* **156**, 303–306.
14. T. Inoue et al. (1989) *J. Bacteriol.* **171**, 3115–3122.
15. P. W. van Ophem, J. Van Beeumen, and J. A. Duine (1992) *Eur. J. Biochem.* **206**, 511–518.

### Suggestions for Further Reading

16. H. Weiner et al., eds. (1987 to 1997) *Enzymology and Molecular Biology of Carbonyl Metabolism*, Vols. 1–6, Plenum, New York. Up-to-date summaries on ADH and the metabolically-related enzymes. Published biannually from the conference series with the same name, Vol. 6 (1997) now in press.
17. H. Eklund and C.-I. Brändén (1987) Alcohol dehydrogenase. *Biol. Macromol. Assembl.* **3**, 74–142.
18. W. F. Bosron, T. Ehrig, and T.-K. Li (1993) Genetic factors in alcohol metabolism and alcoholism. *Seminars Liver Dis.* **13**, 126–135.
19. J. Shafqat et al. (1996) Pea formaldehyde-active class III alcohol dehydrogenase: Common derivation of the plant and animal forms but not of the corresponding ethanol-active forms (classes I and P). *Proc. Natl. Acad. Sci. USA* **93**, 5595–5599.

### Aligning Sequences

The most basic activity in [sequence analysis](#) involves aligning protein or nucleotide sequences together. This need arises due to the processes of molecular [evolution](#): [gene duplication](#) followed by continual **divergence** of the sequences through the accumulation of [mutations](#) over time. Comparative biological analysis, which has long been such a powerful tool for biologists (as exemplified by Linnaeus and Darwin), is arguably even more applicable in sequence analysis than in any other branch of biology, because it can be applied to an enormous number of character states at the level of individual residues in nucleic acid or protein sequences. First, however, related sequences must be correctly aligned before the power of comparative analysis can be brought to bear. Because of the difficulty of aligning highly diverged sequences, and the many applications of sequence alignment, this is one of the most active areas for method development in computational biology.

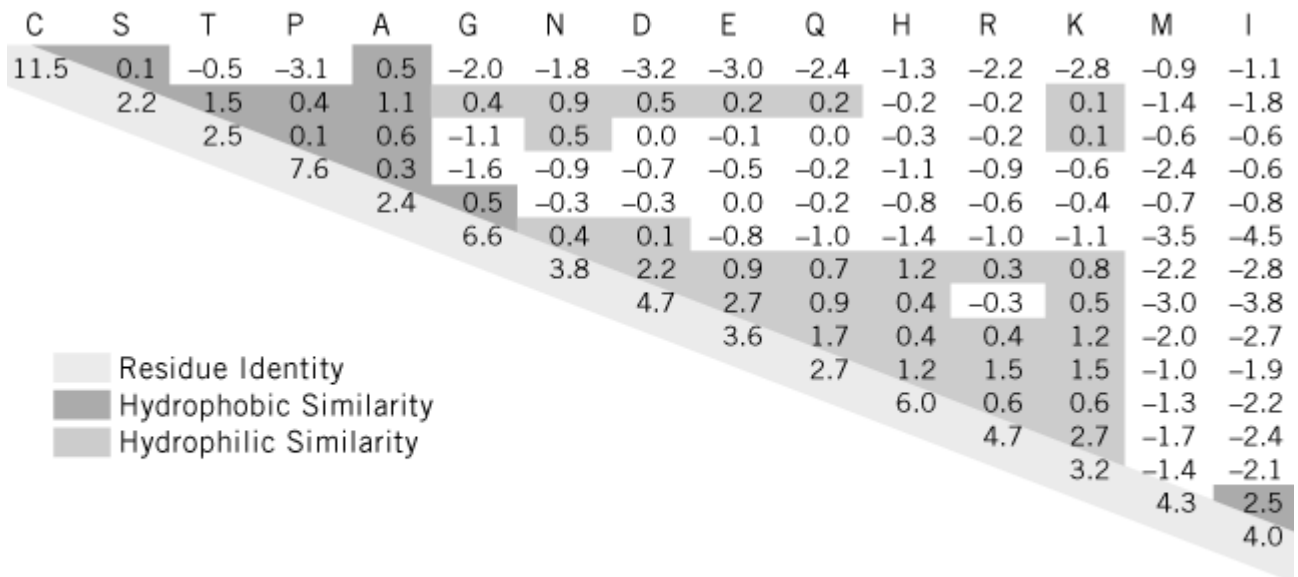
Alignment tasks generally divide into pairwise sequence alignment and multiple sequence alignment, although the underlying algorithms may share many details. The most sensitive methods for aligning sequences belong to the class of algorithms known as dynamic programming (or minimum string edit) that were initially developed for applications in text comparison. Two of the dynamic programming algorithms most used in biology are usually known as Needleman–Wunsch (1) and Smith–Waterman (2), after the researchers who first applied them to biological sequences. Because these algorithms allow gaps to be inserted at any position in the sequences, they are computationally slow. By contrast, word comparison algorithms, which do not allow gaps, are much faster, but at the expense of accuracy and sensitivity in aligning sequences.

## 1. Pairwise Sequence Alignment Algorithms

### 1.1. Dynamic Programming

The basic algorithm works through a two-dimensional matrix in which every residue in one sequence is scored against every residue in the other sequence (1). The algorithm begins in one corner (eg, top left) of the matrix and ends in the opposite corner (bottom right). At each point in the matrix, the algorithm iterates the same set of choices. Typically, it chooses which of three existing paths scores best when extended into the current point: (i) match the residues and continue aligning from the previously aligned residue pair, or (ii) pay the penalty and insert a one-residue gap into sequence X, or (iii) pay the penalty and insert a one-residue gap into sequence Y (see [Gap Penalty](#)). The algorithm is guaranteed to find the best path through the matrix, allowing for gaps at any position in either sequence. Scores for matching the residues are taken from residue exchange, or mutation, matrices. For nucleotide sequences, these are usually quite simple: for example, +1 for an *identity*, 0 for a *transition*, and –1 for a *transversion*. For proteins, typical exchange matrices are the more complex 20 \* 20 *PAM* ([point accepted mutation](#)) matrices introduced by Margaret Dayhoff (3), or subsequently derived *PAM* series derived from larger alignment datasets (4, 5). A *PAM* 250 matrix is shown in Figure 1. Gap penalties are used to control the frequency and length of gaps inserted in the sequences. Where appropriate, varying gap penalties in a position-specific manner can improve the alignment.

**Figure 1.** *PAM*250 amino acid exchange matrix developed by Gonnet and colleagues (4) and superseding the original D acids have positive log-odds exchange values while dissimilar pairs have negative values. All positive scores are colored purple highlights similar pairs of hydrophobic residues, and green indicates similar pairs of hydrophilic residues. The high aromatic residues (Phe, Tyr, Trp) and are stronger than exact matches of highly mutable residues such as Ser. The strong hydrophobic residues and small or negatively charged residues.



Dynamic programming algorithms work through a two-dimensional matrix of area  $M \times N$  in aligning sequences of lengths  $M$  and  $N$ . Therefore the computational requirement [usually symbolized as  $O(MN)$ ] has a constant factor for the calculation, multiplied by the two sequence lengths. To obtain an alignment, the algorithm makes a first pass to determine the end of the highest scoring matched segment and then a second pass working back to obtain the alignment. If only the highest score, but not the actual path, is needed (as in a database search), then only the first pass need be done. Naive implementations that first plot the whole matrix to an array will also use  $O(MN)$  memory. Because the algorithm works through the array systematically, however, it is unnecessary to store the whole array in memory. Memory-efficient implementations of the first pass are straightforward, because the alignment is not being kept (6). The second pass, to obtain the alignment, is more complicated, but memory-efficient recursive methods have been developed that have allowed large alignment tasks to be ported to small personal computers, with a small but acceptable loss in calculation speed (7-9).

## 1.2. Global Alignment

The standard Needleman–Wunsch algorithm (1) finds the optimal full-length alignment for a pair of sequences. Global alignment is appropriate where sequences are known to be both **homologous** and collinear and is therefore often used for multiple alignment of sequence families.

## 1.3. Best Local Alignment

Variants of the Smith–Waterman algorithm (2) find the optimal alignment that has a positive value for the path for a given pair of sequences. Except for highly related sequences, the best local alignment is a partial match between the sequences. Residue exchanging mutation matrices, such as *PAM250*, provide log-odds scores for the likelihood that a pair of residues will exchange as a result of mutation: Similar residues that exchange easily have positive log scores, while dissimilar residues have negative log scores. The best local alignment is taken as the highest-scoring continuously positive path. This algorithm is appropriate under conditions where sequences are not known to be both fully homologous and collinear—for example, *multidomain* proteins, or DNA regions containing rearrangements. Smith–Waterman type algorithms underlie the most sensitive methods for database searching by sequence homology yet to be devised. Because of the computational cost, they are not applied as often as search methods using ungapped alignment algorithms.



The Waterman–Eggert (10) extension of the algorithm will return sets of suboptimal paths that do not intersect with the optimal path, and in this way it can find repeats in sequences.

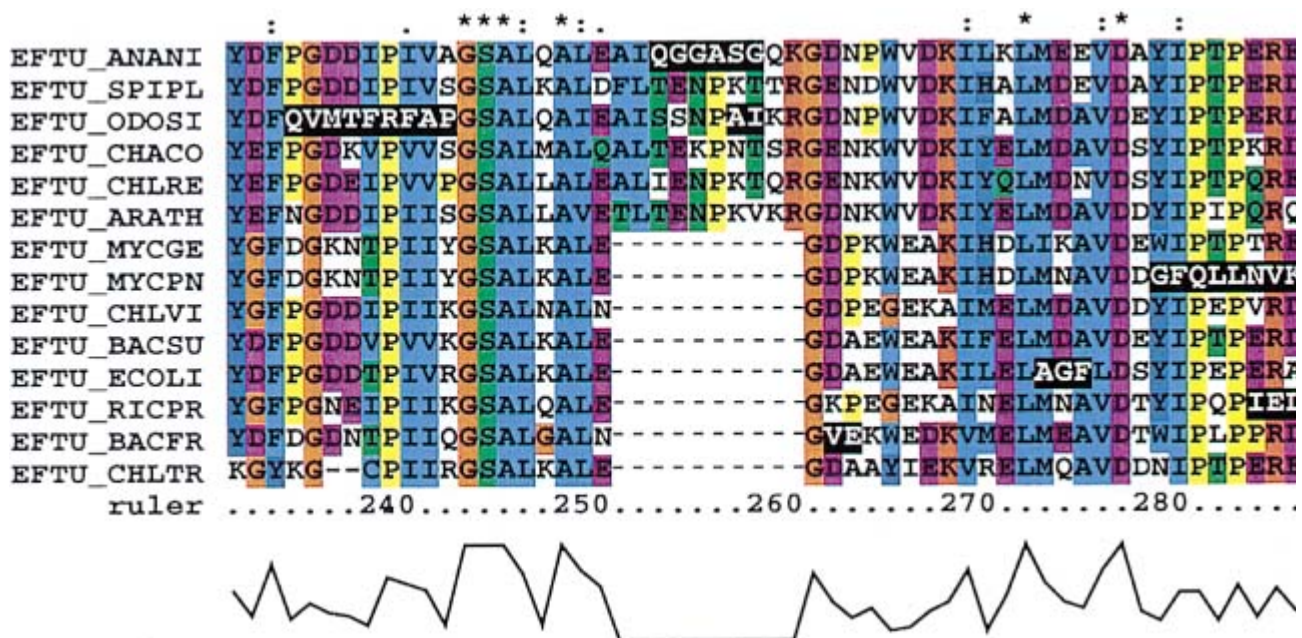
#### 1.4. Best Local Ungapped Alignment

Widely used database search tools, such as BLAST (11) and FASTA (12), search for the highest scoring matched regions without allowing for gaps. Word searches and other ungapped alignment methods are much faster than dynamic programming approaches, but at the expense of sensitivity. Thus these methods are likely to miss homologous but divergent matches. To improve the results, FASTA does a second dynamic pass on the set of top hits. For a small reduction in search speed, BLAST2 examines the gap cost between the set of ungapped positive matches between two sequences and returns composite best locally aligned regions, including gaps whenever the score is still positive. The latter algorithm is likely to approach Smith–Waterman sensitivity except for the most unusual alignment circumstances.

## 2. Multiple-Sequence Alignment

This is a set of homologous protein or nucleotide sequences that have been correctly aligned, allowing for the presence of [indels](#). Figure 2 shows an aligned region for some [elongation factor](#) TU sequences.

**Figure 2.** Part of a multiple alignment of 14 prokaryotic EFTU sequences using the one-letter code for amino acids. Gap by asterisks are completely conserved, columns marked by colons are strongly conserved, and columns marked by period bottom shows the conservation in the columns. Color is an essential aid to sequence analysis and is used here to highlight properties. Inverted characters indicate poorly matching regions of sequence. Some of these are due to natural sequence (sequence determination causing frameshifted regions: EFTUSPIPL 203–206, EFTUODOSI 234–243, EFTUMYCPN 27 errors may lead to false inferences: If R-285 were completely conserved, it would be a candidate functional residue, while incorrectly suggests a surface loop. See color insert.



### 2.1. Uses of Multiple Alignments

Multiple alignments are indispensable in computational biology. They are the basic dataset used to construct [phylogenetic trees](#), which are themselves important computational tools (eg, for weighting sequences by divergence) as well as providing insight into past evolution. They reveal conserved

residues that are likely to be structurally or functionally (eg, in catalysis) important and unconserved positions that either are unimportant or have acquired a change of function (Fig. 2). They improve the accuracy of many sequence analysis functions as compared to single sequences, such as **secondary structure prediction** (13), **coiled-coil** prediction (14), and **transmembrane helix** prediction (15). They are useful as the query input for the most sensitive homology searches (alignment profile (16) and hidden Markov model searches (17)) and can be used to detect divergent homologues that single-sequence queries cannot pick up. They are useful in the detection of **domains** and in defining their boundaries in modular, **mosaic proteins** (18). Multiple alignments of folded RNAs are a prime resource used in determining their secondary and tertiary structures, by mapping residue conservation and long-distance-coupled mutations (19, 20). DNA multiple alignments are used for identifying conserved signals, such as **promoter** elements and **RNA splice** sites (21).

## 2.2. Multiple-Alignment Algorithms

So far it has been necessary to adopt heuristic strategies to generate multiple alignments, because formally correct methods have been computationally impractical to implement. The ideal method to align  $N$  sequences would be  $N$ -dimensional dynamic programming, as this would be guaranteed to find the optimal path (ie, the optimal multiple alignment) in an  $N$ -dimensional matrix. Unfortunately the computer time required to align  $N$  sequences of length  $l$  is  $O(l^N)$  and is impractical for more than three or four sequences, although by limiting the search space to likely regions, the MSA program can align up to eight sequences (22). Another broad class of methods, those that iterate toward an optimized score for the alignment, are still computationally intensive but are becoming practical with increasing computer power. These include a number of approaches, some of which may be used in combination, such as global minimization, genetic algorithms, and trained **neural networks** (23, 24). The goal is to harness a good model description of a multiple-sequence alignment with an effective iteration strategy, so as to get high-quality alignments in a practical timescale. In the meantime, widely used alignment programs such as Clustal W (25) follow the heuristical clustered alignment strategy.

## 2.3. Progressive Clustered Alignment

This two-step approach was introduced by Feng and Doolittle (26) in an attempt to minimize errors in the final multiple alignment, by aligning the most similar sequences first and the most divergent ones last. The set of unaligned sequences are first aligned in pairwise fashion to each other, so that a matrix of the approximate pairwise similarities may be obtained. A matrix-based tree construction method, such as Neighbor-Joining (27), is used to construct a dendrogram linking the sequences according to their observed similarities. Guided by the dendrogram branching order, the sequences are then sequentially aligned together using dynamic programming, beginning with the most closely related sequences and ending by merging the most divergent groups by profile alignment. This procedure minimizes alignment errors, which become more likely with increasing sequence divergence, becoming a problem for proteins less than 25% to 30% identical in sequence. The final alignment should always be examined for misalignment, especially of the more divergent sequences, because such errors are likely to be present in any but the most straightforward multiple alignment task.

The sensitivity of the basic clustered alignment strategy can be improved by several modifications, such as weighting the sequences by divergence and position-specific **gap penalties** (25). Where **tertiary structure** information is available, gap penalty masks can be employed to guide the gaps into regions of sequence that are expected to be tolerant of **indels** (28).

## 2.4. Profile Alignment

A set of aligned sequences can be aligned to one or more new sequences by treating the group much as a single sequence. The score for one alignment column can be obtained by summing the log-odds residue exchange scores for the observed set of amino acids (16), correcting for sequence relatedness by downweighting similar sequences (25). Conserved alignment columns score more highly than unconserved columns, where the log-odds scores tend to cancel each other. Gap penalties can be

lowered at existing indels, because gaps in new sequences are more likely at these positions than at ungapped positions. The improvement in signal to noise provided by the extra information in the alignment means that a profile alignment is more accurate than independent pairwise alignment of the same set of sequences. As well as being used to merge aligned groups in the clustered alignment strategy, profiles provide a sensitive search strategy to find highly divergent homologues and are one of the main sequence analysis tools for identifying protein domain families (18).

### 2.5. Hidden Markov Model (HMM) Alignment

HMMs are a class of probabilistic models, applicable when the components of a complex linked system behave independently (a so-called Markov chain). Up to a point, this is valid for residue mutations in globular protein sequences, and HMMs are being applied increasingly in multiple-alignment algorithms and profile-type database searches (17, 29). The models are more complex than the widely used PAM model for protein evolution introduced by Margaret Dayhoff (3), providing both advantages and disadvantages. On the plus side, the models provide direct probabilities for evaluating database search matches, which can include multiple matches in a repeated sequence and can formally treat biological complexities such as splice junctions in genomic sequence. On the debit side, the extra parameters lead to more complex optimization problems at several levels. Thus, inexperienced users may set up poorly optimized HMMs while, for program developers, there is a problem as to whether the most appropriate HMMs are being applied to sequence evolution. It is important for the user not to be seduced by technical jargon, but to justify the use of newer methods such as these on the basis of convincing results.

## 3. Error in Sequence Alignments

Errors are very common and invalidate any conclusions obtained by methods based on sequence alignments (see Fig. 2). The three main sources of error are: the input sequences, mistakes by the user, or alignment algorithm failure. Causes and effects of error are manifold; some major ones are discussed below.

### 3.1. Errors in the Input Sequences

Experimental errors in determining sequences include double-insert cloning errors, truncated cDNA clones, base insertion or deletion causing translation **frameshifts**, and translation with inappropriate **genetic code** (eg, for plastid-encoded proteins). Figure 3 shows an example of a frameshift error in a database entry. Algorithm failure during alignment may be induced by sequence error, such as frameshifted regions, in which the sequences are no longer similar, and may induce incorrect gap placement (Fig. 2). Errors in sequences are very common and likely to be present in any sequence family (30, 31).

**Figure 3.** Frameshifted segment in the *Mycoplasma pneumoniae* EFTU sequence revealed by comparing the DNA to the closely related *Mycoplasma genitalium* EFTU using a dynamic programming three-frame comparison allowing shifts between translation frames (30). Exclamation points mark the frameshift sites. The first frameshift is caused by a base being dropped and the second frameshift is caused by a base being added, thereby returning to the original translation frame. See Figure multiple alignment spanning this region.

```

***** **
M.genitalium EFTA 151 AEEVRDLLTSYGFDPKNTPIIYGSALKALEGDPKWEAKIHDLIKAV
Translation AEEVRDLLTSYGFDPKNTPIIYGSALKALEGDPKWEAKIHDLMNAV
M.pneumoniae DNA 619 ggggcggttattgtggaacaatgtgcagcgggcatggaacgtaagg
caatgattccagtagaaccttagcctactagacagacataattact
aagatcaattccttcgccttttttatagtagtatttggatgcttagtat

**** *! *****
M.genitalium EFTU 200 IPTPTREVDG+++PFLLAIEDTMTITGRGTVVTGRVERGELKVGQE
Translation IPTPEREVD^^^PFLLAIEDTMTITGRGTVVTGRVERGELKVGQE
M.pneumoniae DNA 765 acacgcgggAAACctttgaggaaaaagcggaggagcggcggtagggc
tccagata ctttctaactctcgggcttcggtaggatatgaa
tattatagc gcggacaccgggtcttcgtctgtattagaataa

```

Further errors can arise in preparing sequence database entries. These can affect any part of the entry and can have particularly unpredictable consequences. However, the most common annotation errors are undoubtedly in predicted translation products and are a consequence of the limited accuracy of current gene prediction algorithms; that is, despite high-quality DNA sequence, the predicted protein sequence might contain artificially translated **introns**, missed exons, or terminal truncations or might be an artificial fusion product of two independent genes. Protein sequence alignments can often help in identifying [translation](#) problems.

### 3.2. Mistakes by the User

Generally these are due to inadequate attention to detail. Erroneous inclusion of nonhomologous sequences in the input set is quite common, particularly if keyword searches are used to extract a set of sequences; there are many examples of unrelated proteins performing equivalent functions, as well as sequences that have functions incorrectly ascribed to them. Another source of error occurs during the alignment process, when parameters may be set poorly. Trial and error is usually needed to find the parameters that best suit a particular group of sequences. For example, insertion of too few or too many gaps may suggest that the gap penalties are not optimal. It is important to take the time to get a basic understanding of how a given program works, or it is unlikely to do the best job.

### 3.3. Algorithm Failure

Clustered alignment is a heuristic strategy that is not guaranteed to find the optimal multiple alignment. The alignment process is likely to compound errors introduced by the user or in the input sequences. Alignment mistakes will also be made in difficult alignment cases, even when there are no input errors. There are many instances of homologous proteins having diverged over long periods of time, so that they have apparently very little sequence similarity. The “twilight zone” where sequence similarity merges with sequence dissimilarity is in the range 20% to 25% identity. (Five percent identity would be a random match for protein sequence, neglecting residue biases. In practice, given residue biases and gap insertions to maximize the similarity, the expectation for random sequence matches approaches ~15% identity, but higher for short or biased sequences. There is, however, no *a priori* reason why correctly aligned but extremely divergent proteins should not be found with 0% pairwise identity.) Divergent nucleotide sequences are even harder to align, because random similarity is reached at ~25% identity before gaps are added. Wherever possible, alignments need to be checked against additional data available for a sequence family, such as known [tertiary structures](#) and whether invariant catalytic residues, or any other known conserved motifs, are correctly aligned.

### 3.4. Consequences of Errors in Aligned Sequences

Errors in multiple alignments can have disastrous consequences for **phylogenetic** inference on the basis of sequence trees. Sequences with misaligned segments or translation frameshifts are apparently more divergent than they should be from the other sequences. The branch leading to that

sequence will then have a longer length (which erroneously implies a more rapid [molecular clock](#)) and the branch point may migrate toward the centre of the tree, giving a false order of divergence. Such incorrect phylogenies may be quite exciting when they seem to refute established viewpoints. Two rules of thumb are useful: (i) Infinitely more wrong tree topologies can be generated than right ones, and (ii) wrong phylogenies are more interesting than right phylogenies. Incautious advocacy of a wrong phylogeny can waste many peoples' time.

Errors also disrupt evaluation of conserved sites in sequences (Fig. 2). Catalytic residues are often absolutely conserved, so a single misaligned sequence may lead to rejection of the correct site. Most conserved residues have a structural role: Structure prediction for protein uses conserved **hydrophobic** residues and, for RNA, conserved base-pairing residues. Misalignments disrupt the conservation periodicities, leading to rejection of the correct structures. Terminally truncated sequences of multidomain proteins can lead to false inferences for the domain boundaries, which are very usefully defined by coincidence with the protein *N*- or *C*-termini.

### Bibliography

1. S. B. Needleman and C. D. Wunsch (1970) *J. Mol. Biol.* **48**, 443–453.
2. T. F. Smith and M. S. Waterman (1981) *Adv. Appl. Math.* **2**, 482–489.
3. M. O. Dayhoff, R. M. Schwartz, and B. C. Orcutt (1978) In *Atlas of Protein Sequence and Structure*, Vol. **5**, Suppl. 3 (M. O. Dayhoff, ed.), NBRF, Washington, pp. 345–352.
4. S. A. Benner, M. A. Cohen, and G. H. Gonnet (1994) *Protein Eng.* **7**, 1323–1332.
5. G. Vogt, T. Etzold, and P. Argos (1995) *J. Mol. Biol.* **249**, 816–831.
6. O. Gotoh (1982) *J. Mol. Biol.* **162**, 705–708.
7. E. W. Myers and W. Miller (1988) *CABIOS* **4**, 11–17.
8. J. D. Thompson (1995) *CABIOS* **11**, 181–186.
9. J. A. Grice, R. Hughey, and D. Speck (1997) *CABIOS*, **13**, 45–53.
10. M. S. Waterman and M. Eggert (1987) *J. Mol. Biol.* **197**, 723–728.
11. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman (1990) *J. Mol. Biol.* **215**, 403–410.
12. W. R. Pearson and D. J. Lipman (1988) *Proc. Natl. Acad. Sci. USA* **85**, 2444–2448.
13. B. Rost and C. Sander (1993) *J. Mol. Biol.* **232**, 584–599.
14. A. Lupas, M. Van Dyke, and J. Stock (1991) *Science*, **252**, 1162–1164.
15. B. Persson and P. Argos (1994) *J. Mol. Biol.* **237**, 182–192.
16. M. Gribskov, A. D. McLachlan, and D. Eisenberg (1987) *Proc. Natl. Acad. Sci. USA*, **84**, 4355–4358.
17. S. Eddy (1996) *Curr. Opin. Struct. Biol.* **6**, 361–365.
18. P. Bork and T. J. Gibson (1996) *Methods Enzymol.* **266**, 162–184.
19. C. R. Woese, R. R. Gutell, R. Gupta, and H. F. Noller (1983) *Microbiol. Rev.* **47**, 621–669.
20. F. Michel and E. Westhof (1990) *J. Mol. Biol.* **216**, 585–610.
21. R. F. Doolittle (1990) *Molecular Evolution: Computer Analysis of Protein and Nucleic Acid Sequences, Methods in Enzymology*, Vol. 183, Academic Press, San Diego, CA.
22. D. J. Lipman, S. F. Altschul, and J. D. Kececioglu (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4412–4415.
23. O. Gotoh (1996) *J. Mol. Biol.* **264**, 823–838.
24. C. Notredame and D. G. Higgins (1996) *Nucleic Acids Res.* **24**, 1515–1524.
25. J. D. Thompson, D. G. Higgins, and T. J. Gibson (1994) *Nucleic Acids Res.* **22**, 4673–4680.
26. D.-F. Feng and R. F. Doolittle (1987) *J. Mol. Evol.* **25**, 351–360.
27. N. Saitou and M. Nei (1987) *Mol. Biol. Evol.* **4**, 406–425.

28. A. M. Lesk, M. Levitt, and C. Chothia (1986) *Protein Eng.* **1**, 77–78.
29. A. Krogh, M. Brown, S. Mian, K. Sjölander, and D. Haussler (1994) *J. Mol. Biol.* **235**, 1501–1531.
30. E. Birney, J. D. Thompson, and T. J. Gibson (1996) *Nucleic Acids Res.* **24**, 2730–2739.
31. J. M. Claverie (1993) *J. Mol. Biol.* **234**, 1140–1157.

### Suggestions for Further Reading

32. R. F. Doolittle (1990) *Molecular Evolution: Computer Analysis of Protein and Nucleic Acid Sequences*, *Methods in Enzymology*, Vol. 183, Academic Press, San Diego, CA.
33. R. F. Doolittle (1996) *Computer Methods for Macromolecular Sequence Analysis*, *Methods in Enzymology*, Vol. 266, Academic Press, San Diego, CA.
34. D. Sankoff and J. B. Kruskal (1984) *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Addison-Wesley, Reading, MA.
35. M. S. Waterman (1989) *Mathematical Methods for DNA Sequences*, CRC Press, Boca Raton, FL.

## Alkaline Phosphatase

Alkaline phosphatases (E.C. 3.1.3.1) belong to a family of orthophosphoric monoester phosphohydrolases that have an alkaline pH optimum. Their genes are very frequently used as a **reporter** gene. Alkaline phosphatase activity is most commonly detected by the hydrolysis of 5-bromo-4-chloro-3-indolyl phosphate (BCIP), which, when coupled to the reduction of nitro blue tetrazolium (NBT), forms a formazan and an indigo dye that together form a strong black/purple precipitate (1). In addition, a number of fluorogenic substrates are also available (1). 4-Methylumbelliferyl phosphate (MUP) gives a blue **fluorescent** product upon hydrolysis, and 2-hydroxy-3-naphthoic acid-2-phenylamide phosphate (HNPP/Fast Red TR) fluoresces with a broad emission peak between 540 and 590 nm and can be observed using either fluorescein or rhodamine filter sets. Molecular Probes Inc. have also developed a proprietary substrate called ELF-97 (Enzyme Linked Fluorescence-97), in which cleavage of the molecule converts it from a soluble phosphate to an insoluble alcohol, with an accompanying shift from weak blue fluorescence to a bright yellow fluorescence (2).

Human placental alkaline phosphatase (hpAP) is most commonly used as a reporter enzyme in nonradioactive detection systems, where the enzyme is linked to other molecular probes (eg, specific antibodies)—for example, for the detection of proteins and nucleic acids by **Western blot**, **Southern blot**, and **Northern blot** analysis, and most commonly in **in situ hybridization**. Elegant studies have also been performed in which hpAP is fused to soluble extracellular domains of receptor molecules, which are then used as probes to detect the sites of ligand production *in vivo* and facilitate the subsequent **cloning** of the ligand genes (3).

Although animals express a number of alkaline phosphatase genes, the human placental isoform has been developed as a reporter gene, because it can be distinguished from other endogenous isoforms through its high thermostability (4). Thus all background endogenous alkaline phosphatase activities, from embryonic, intestinal, and nonspecific genes, can be minimized by preheating tissue preparations up to 80°C for prolonged periods. hpAP also retains its activity following histological processing for wax imbedding and sectioning tissues. In addition, background from endogenous

alkaline phosphatases can be further inhibited by the amino acids L-phenylalanine or L-homoarginine.

hpAP has been used in a wide range of applications, including *in vitro* transfection studies and transgenic studies *in vivo*. The sensitivity of hpAP in transient expression assays is equivalent to that of chloramphenicol acetyltransferase (CAT) (see [Reporter Genes](#)). A particularly useful variant of hpAP is a [cDNA](#) encoding a secreted form of the protein ([5](#), [6](#)) that allows hpAP activity to be assayed by sampling tissue culture medium, giving the benefit of monitoring changes in gene expression with time. In this system, background activities are further eliminated, because the endogenous isoforms are anchored to the cell membranes and do not contribute to the activity in the culture supernatants.

hpAP is also an effective reporter gene to analyze gene expression *in situ* in tissue preparations. hpAP was first used in retroviral vectors to infect small numbers of cells in developing embryos as a tool to study cell fate and lineage analysis; subsequently, hpAP has been used as a robust reporter gene in transgenic mice ([7](#)). Indeed, mice that express high levels of hpAP from a ubiquitously expressed **promoter** thrive with no adverse effects. hpAP is particularly useful to use in combination with a second reporter gene, such as *lacZ* (**beta-galactosidase**), in dual labeling studies, as the common substrates are quite distinct and give different colored products ([8](#)).

#### Bibliography

1. D. A. Knecht and R. L. Dimond (1984) *Anal. Biochem.* **136**, 180–184.
2. K. D. Larison et al. (1995) *J. Histochem. Cytochem.* **43**, 77–83.
3. J. G. Flanagan and P. Leder (1990) *Cell* **63**, 185–194.
4. P. Henthorn (1988) *Proc. Natl. Acad. Sci.* **85**, 6342–6346.
5. T. T. Yang, P. Sinai, P. A. Kitts, and S. R. Kain (1997) *Biotechniques*, **23**, 1110–1114.
6. B. R. Cullen and M. H. Malim (1992) *Methods Enzymol.* **216**, 362–368.
7. S. E. DePrimo, P. J. Stambrook, and J. R. Stringer (1996) *Transgenic Res.* **5**, 459–466.
8. X. Li, W. Wang, and T. Lufkin (1997) *Biotechniques*, **23**, 874–878.

#### Alkylation

Alkyl groups have the general formula  $C_n H_{2n+1}$ -. Alkylation is a reaction that introduces an alkyl group into a compound. Commonly used alkylating agents are olefins, alkyl halides, and alcohols. Biologically important alkylation occurs on nitrogen and oxygen atoms in polynucleotides. Alkylating agents that cause these reactions are [mutagens](#) and [carcinogens](#) and include methyl methanesulfonate, dimethyl sulfate, *N*-methyl-*N*-nitrosourea, dimethylnitrosourea, 1-methyl-3-nitro-1-nitrosoguanidine, etc. Some alkylating agents are used as anticancer drugs, including cyclophosphamide, nitrogen mustard and its oxide, and triethylenephosphoramidate.

Alkylating agents easily alkylate nucleophiles in proteins such as those on [cysteine](#) ([1](#)), [histidine](#) ([2](#)), and [methionine](#) ([3](#)) residues, [amino groups](#) ([4](#)), and sometimes [tyrosine](#) residues and [carboxyl groups](#) ([5](#)). The most well-known alkylation is that for blocking free [thiol groups](#) of cysteine residues to prevent their oxidation, especially in the course of determining the [primary structure](#) of a protein (see [Protein Sequencing](#)).

## 1. Reduction and Alkylation of Proteins

The protein at a concentration of 2% (w/v) is dissolved in 0.5 M **Tris** buffer, pH 8.1, containing 6 M **guanidinium chloride** and 0.002 M **EDTA** and reduced with **dithiothreitol** (50 mol/mol of **disulfide bond**) at 30° to 40°C for 4 h under nitrogen (6). Then the solution is cooled to room temperature, and iodoacetamide (100 mol/mol of the original disulfide bond) is added, with the addition of NH<sub>4</sub>OH as necessary to maintain a constant pH. Other alkylating agents that can be employed instead are iodoacetic acid, 4-vinylpyridine, N-ethylmaleimide, and ethyleneimine.

### Bibliography

1. J. Bridgen (1972) *Biochem. J.* **126**, 21–25.
2. A. M. Crestfield, W. H. Stein, and S. Moore (1963) *J. Biol. Chem.* **238**, 2413–2420.
3. J. M. Gleisner and R. L. Blaklay (1975) *J. Biol. Chem.* **250**, 1580–1587.
4. H. J. Goren and E. A. Barnard (1970) *Biochemistry* **9**, 974–983.
5. K. Takahashi, W. H. Stein, and S. Moore (1967) *J. Biol. Chem.* **242**, 4682–4690.
6. M. J. Waxdal et al. (1968) *Biochemistry* **7**, 1959–1966.

## Allelic Exclusion

The **clonal selection theory** put forward by Burnet (1) and Jerne (2) proposes that each single lymphocyte expresses only one single type of **antibody**. This theory received its first structural confirmation when it was shown that patients with multiple myeloma, a proliferative disorder of plasma cells, had in their serum a highly elevated level of only one molecular species of **immunoglobulin**, originating from the corresponding malignant clone. Single-cell experiments performed by the Nossal group in Australia (3) indicated that each cell isolated from mouse **immunized** with *Salmonella* **antigens** produced only one antibody of one given specificity. These observations were in agreement with the clonal theory, but did not explain by which mechanism this was achieved. An elegant approach to that problem was made with the analysis of allotypes in heterozygous animals. Allotypes, as defined by Oudin (4), are antigenic specificities that are shared by a group of animals within a given species. These genetic markers are **allelic** variants of immunoglobulins that are found on both the heavy and light chains. They have been studied extensively in the rabbits, mice, and humans. In most cases, they are the result of a few amino acid substitutions that confer specific **epitopes** (or allotypes) that can be recognized by specific antibodies. In a rabbit heterozygous for H- and L-chain allotypes, Oudin showed that any given immunoglobulin molecule always had two identical heavy chains and two identical light chains with respect to the corresponding allotypes. This ensured that the Ig molecule was symmetrical and therefore had two identical antibody combining sites, an obvious prerequisite of the clonal theory. The next step was made at the intracellular level, when Pernis et al. (5) showed with specific anti-allotype antibodies that, in a heterozygous rabbit, any given cell always expressed only one allotype of each heavy or light chain. This suggested that a **B cell**, although **diploid**, had only one of its two alleles functional, and the phenomenon was termed “allelic exclusion.”

The molecular mechanisms that account for this phenomenon could be understood only when the complex mosaic structures of immunoglobulin (Ig) genes were deciphered. Each Ig heavy or light chain contains an NH<sub>2</sub>-terminal variable and a COOH-terminal constant region, each encoded by a multiplicity of gene segments. Diversity of the variable regions, which is directly linked to the



necessity to produce a very large number of different antibody molecules with a limited number of genes, is generated by specific mechanisms of [gene rearrangements](#) that take place exclusively in lymphocytes. The heavy-chain **variable region** is encoded after the rearrangement of three discrete sets of germline gene segments ( $V_H$ , D,  $J_H$ ), and two for the light chain ( $V_L$ ,  $J_L$ ). Gene rearrangements take place sequentially during B-cell differentiation in the bone marrow, in a strictly regulated manner, as initially proposed by Alt and Baltimore. The heavy-chain gene segments rearrange D to  $J_H$  and  $V_H$  to D- $J_H$ . The gene-segment joinings are random, so only one-third of the rearranged genes have an open reading frame. As soon as one rearrangement is in frame, the corresponding m chain is expressed. Whenever the first rearrangement is successful, the second allele remains in the germline configuration, which indicates that a regulatory mechanism has taken place. It was shown by transfection experiments and by making m transgenic mice that the negative feedback inhibition of the gene rearrangement of the second allele was controlled by the m chain itself, when expressed at the cell surface. It was also shown that, at this “preB stage,” the m chain was associated with a surrogate light chain (YL), which is monomorphic and composed of two polypeptide chains encoded by the I5 and VpreB genes. Expression of m-YL has been shown to down-regulate the heavy-chain gene rearrangement, whereas rearrangement of the light chain locus is switched on, in the same random reading frames as above. The final result of this complex sequence of events is that, at each locus, only one allele is functional; the other one either remains in the germline configuration or has an out-of-frame rearranged gene, thereby accounting for the allelic exclusion phenomenon.

See also entries [Clonal Selection Theory](#), [Gene Rearrangement](#), and [Antibody](#).

#### Bibliography

1. M. F. Burnet (1959) *The Clonal Selection Theory of Acquired Immunity*, Vanderbilt University Press, Nashville, TN.
2. N. K. Jerne (1955) The natural selection theory of antibody formation. *Proc. Natl. Acad. Sci. USA* **41**, 849–857.
3. G. J. V. Nossal, A. Szenberg, G. L. Ada, and C. M. Austin (1964) Single cell studies on 19 S antibody production. *J. Exp. Med.* **119**, 485–502.
4. J. Oudin (1960) Allotypy of rabbit serum proteins. I. Immunochemical analysis leading to the in vivo utilization of seven main allotypes. *J. Exp. Med.* **112**, 107–124.
5. B. Pernis, G. Chiappino, A. S. Kelus, and P. G. H. Gell (1965) Cellular localisation of immunoglobulins with different allotypic specificities in rabbit lymphoid tissues. *J. Exp. Med.* **122**, 853–876.
6. F. Alt (1984) Exclusive immunoglobulin genes. *Nature* **312**, 502–503.

#### Suggestions for Further Reading

7. E. ten Boeckel, F. Melchers, and A. G. Rolink (1998) Precursor B cells showing H chain allelic inclusion display allelic exclusion at the level of pre-B cell receptor surface expression. *Immunity* **8**, 199–207.
8. D. Löffert, A. Ehlich, W. Müller, and K. Rajewsky (1996) Surrogate light chain expression is required to establish immunoglobulin heavy chain allelic exclusion during early B cell development. *Immunity* **4**, 133–144.

#### Alloantibody, Alloantigen

Alloantigens are **allelic** variants that can induce, when injected into animals of the same species but having a distinct genetic background, the production of the corresponding alloantibodies. The most commonly known example of alloantigens are the blood group substances, which define the basic ABO blood groups in humans. Blood group substances have long been identified as being derived from a polysaccharide core that is found in pneumococcus group XIV, onto which three genes encoding for glycosyl transferases will serially add units that will confer a specific antigenicity to the newly derived molecule. Three genes operate this system in humans: A, B, and H. The H gene is present in all individuals and will add one residue of fucose, leading to the so-called H substance. In individuals who possess the A gene, an additional GalNAC (*N*-acetylgalactosamine) will be added, providing the A substance, which is expressed at the red blood cell surface and confers the A group. In those who have the B gene, a galactose residue is added to H, making the B substance, which is also expressed at the cell surface of red cells in individuals of group B. The A and B genes are codominant, so an individual will express the A, the B, or both the A and B molecules at the red blood cell surface, depending upon what A and B genes are present. If neither the A nor B allele is present, the group will be O. There are many other blood groups in humans, with special reference to the Lewis groups, which are also derived from the same polysaccharide backbone by addition of different units. What is peculiar regarding the ABO blood group system is the fact that individuals who lack the A or the B specificity spontaneously produce alloantibodies directed against the missing substance. So an individual of group A will make anti-B alloantibodies, an individual of group B will produce anti-A, and one of group O will have both. Very severe accidents in blood transfusion would result from the agglutination of the donor red blood cells by the alloantibodies of the recipient. The reverse situation, agglutination of the host red blood cells by antibodies of an incompatible donor, should also be avoided, although accidents are less severe. Antigen compatibility between donor and recipient is therefore a must. ABO compatibility should also be observed in allografts.

Even if the genetic determinism of these alloantigens is straightforward, it still is not clear why the alloantibodies corresponding to the nonexpressed substance(s) are produced. The prevalent explanation is that blood groups cross-react with bacteria normally present in the intestinal flora, which would stimulate the immune system to produce the corresponding crossreacting antibodies. It should be stressed that these “natural” antibodies are of the [IgM](#) type, which may be directly related to the fact that they are the result of the stimulation by a T-independent polysaccharide antigen. The titer of alloantibodies varies greatly between individuals, but may be considerably elevated upon an incompatible transfusion.

Another famous case of alloantigen in the blood groups is the Rhesus factor, which was responsible for the dramatic hemolytic disease of the newborn, due to the immunization of an Rh mother against red blood cells of an Rh<sup>+</sup> fetus. This occurs during the delivery because some fetal red cells may enter the maternal circulation and induce the formation of [IgG](#) antibodies that will actively cross the placental barrier in a subsequent pregnancy and then provoke lysis of Rh<sup>+</sup> fetus red cells. Treatment involves complete transfusion of the newborn with Rh blood. It has now been generalized to prevent the immunization by injecting the Rh mother with anti-D (anti-Rhesus) antibodies immediately after delivery, to trap the red blood cells from the newborn that would have penetrated the maternal circulation at birth.

Many other alloantigens are known, but alloantibodies are generally not produced unless the antigen is given. This is the case of the major histocompatibility complex (MHC) molecules, which constitute a major problem in transplantation. The name of [major histocompatibility complex](#) indicates by itself how these molecules were first discovered as a major target for graft rejection, as the result of a very severe alloimmune response, characterized primarily by the production of [cytotoxic T lymphocytes](#). Many other systems may behave as potential alloantigens, which simply reflects the existence of allelic variants. The case of allotypes of [immunoglobulins](#) has been studied particularly by immunologists and has provided remarkable genetic markers for the study of

[immunoglobulin biosynthesis](#) and diversity at the time.

#### Suggestion for Further Reading

P. L. Mollison, C. P. Engelfriet, and M. Contreras (1987) *Blood Transfusion in Clinical Medicine*, 8th ed., Blackwell, Oxford, UK.

## Allophenic

Allophenic individuals are composed of cells of two different genotypes (often called mosaics or [chimeras](#)). Allophenic mice are the basis of the revolution in [mouse](#) genetics, allowing mice with specific gene replacements to be generated and studied. Allophenic mice are made by mixing cells from two [embryos](#) of different genotypes (1). These new composite embryos are implanted into foster mothers and allowed to develop. The early embryos are able to compensate and form a single individual from the mixture of embryonic cells. The mosaic progeny that result from these manipulations have tissues that contain cells of the two different genotypes. The allophenic mice have also been called tetraparental mice. When the germ cells are also of mixed origin, progeny can be recovered from both genotypes.

The ability to make allophenic mice in the laboratory was first used to study the contributions of cells to individual tissues. The number of cells that give rise to a particular tissue can be estimated from the proportions of the two genotypes in a large number of allophenic mice. The cellular autonomy of mutant phenotypes can be determined by generating allophenic mice between mutant and wild-type mouse embryos. The two cell types usually differ both in the mutant of interest and in some marker genotype, such as an [enzyme](#) polymorphism.

Two advances in the technology of allophenic mice have contributed to a revolution in mouse genetics in the last few years. The first advance was the ability to use embryonic teratocarcinoma cells from cell culture as one of the two cell types used to make the allophenic mice (2). The cultured cells are injected into the interior of a normal mouse embryo and are incorporated into the embryo, to form an allophenic mouse. The teratocarcinoma cells are totipotent and can even form germ cells that give rise to the **gametes** for the next generation. The second major advance was the development of techniques for [site-directed mutagenesis](#) and gene replacement in mammalian cells in culture (3) (see [Gene Targeting](#)). A gene of interest is altered in a specific way in teratocarcinoma cells in culture. A clone of cells with the altered genotype is produced, and cells from the clone are injected into early mouse embryos to produce allophenic mice. These allophenic mice can be studied as mosaic individuals or can be allowed to develop to fertile adults. If some of the germ cells in the adult allophenic mice are derived from the genetically altered teratocarcinoma cells, progeny can be recovered from the gametes produced by the teratocarcinoma cells and used to found a mutant strain of the specific genotype desired. This has led to the production of hundreds of new mouse strains mutant for developmentally important genes. It has also made the construction of mouse models for human disease far easier than in the past. These technologies have been extended from mice to a variety of other mammals; and they could be used to alter the human [genome](#) genetically, with far-reaching consequences.

#### Bibliography

1. B. Mintz (1962) *Science* **138**, 594.
2. B. Mintz and K. Illmensee (1975) *Proc. Natl. Acad. Sci. USA* **72**, 3585–3589.

3. M. R. Capecchi (1989) *Science* **244**, 1288–1292.

### Suggestions for Further Reading

4. B. Mintz (1974) *Annu. Rev. Genetics* **8**, 411–470.

5. R. Jaenisch (1988) *Science* **240**, 1468–1474.

## Allostery

Allostery is used to describe regulatory phenomena in biological systems. The term “allosteric” (other shape) was coined by Monod and coworkers to describe a particular type of regulatory behavior and marks an important synthesis in the development of biochemistry and molecular biology. Since its inception, the term has evolved to describe several related concepts and is used today to describe a variety of phenomena. To some, it is associated with a particular type of regulatory behavior, while to many it is used to describe several different aspects of regulation. The development of the concept, the meanings it has come to have, its application in describing regulatory properties of [proteins](#) and [enzymes](#), and recent developments that indicate the need to use it in its original sense are described here.

The allosteric concept was conceived in the early 1960s by Monod and coworkers to describe the properties of regulatory enzymes. In the mid-1950s, several investigators described enzymes whose catalytic properties were dependent on effector molecules other than the substrate. The first well-characterized example of this behavior was the finding by Cori and colleagues that 5'-AMP is required for the *in vitro* catalytic activity of [glycogen phosphorylase b](#) from resting skeletal muscle (1). Studies of the regulation of biosynthetic pathways in bacteria showed that the catalytic properties of the first enzyme in a pathway are often modulated by the end product of the pathway, which has a structure different from that of the substrate—for example, isoleucine inhibition of threonine deaminase (see [Threonine Operon](#)) (2) and CTP inhibition of [aspartate transcarbamoylase](#) (3). The differences in the structures of the substrates and effectors were first noted by Monod and Jacob, who postulated in a seminal paper in 1961 that the effectors bind to separate and distinct sites (4). They coined the term *allosteric site* for the effector binding site to distinguish it from the [active site](#) and postulated that the effectors act indirectly by causing changes in the conformation of the enzyme that alter its catalytic site and kinetic properties. In 1963, Monod and coworkers compiled the available data on regulatory enzymes (5). These enzymes all contained more than one subunit per enzyme molecule—that is, dimers, tetramers, and dodecamers. For several of these enzymes, changes in [quaternary structure](#) (association–dissociation of the subunits) had been shown to occur upon addition of the effectors. Monod was also aware of unpublished crystallographic results showing that the oxygen-binding sites on [hemoglobin](#) are located far from one another and that the distances between the amino acid residues labeled by heavy atoms change upon binding of oxygen (Ref. 6, pp. 577–578). On this basis, the key elements of a general model for the functional structures of regulatory enzymes were formulated:

1. Each molecule of the native enzyme contains multiple subunits with binding sites for substrates and ligands, and the regulatory behavior depends on the relationships between the subunits.
2. No direct interactions between substrate(s) or effectors are required, that is, the effectors act indirectly on the catalytic site.
3. The actions of the effectors are due entirely to a reversible conformational change in the protein that is induced by effector binding.

Monod et al. postulated that protein conformational change is the basis for control and coordination of chemical events in living cells, and Monod considered this to be the “second secret of life” (Ref. [6](#), p. 576). Structural studies using [X-ray crystallography](#) methods support the concepts of this general model.

These concepts were used to develop the concerted model, a specific molecular model for describing regulatory behavior ([7](#)). This model was developed as a plausible explanation for the positive cooperativity in the binding of oxygen to hemoglobins. Because the apparent affinity for binding of oxygen changes as its concentration is varied, these are called *homotropic interactions*. The concerted model retains the aspects of the general model with respect to oligomeric structure of the protein, multiple binding sites located far from one another, and conformational changes upon binding of one ligand that are transmitted by protein conformational changes to the other ligand binding sites, thereby altering the binding affinities. This model defines the relations between the conformational changes and ligand binding. The word *concerted* refers to the all-or-none nature of the conformational transition. That is, only two conformations are allowed for the subunits in an oligomer and all of the subunits in a given oligomer have the same conformation. The conformations are denoted by the terms *R-state* and *T-state*, and they differ in affinity for binding ligands. The R-state has the higher affinity for ligands. Positive cooperativity in oxygen binding is explained by a large predominance of the unliganded protein in the T-state with a shift to the R-state upon binding of oxygen. The model was extended to enzymes by assuming that the interactions between the substrate and enzyme are in rapid equilibrium and that all forms of the enzyme have the same  $V_{\max}$ . At this juncture, the term *allosteric* takes a different meaning. The conformational change between the two states is called an *allosteric transition*. So, in addition to referring to sites for molecules other than the substrate, it refers to interactions between substrate binding sites on different subunits. The concept of the “other shape” is extended to the enzymes, which are termed “allosteric enzymes.”

This initial formulation of the concerted model describes homotropic interactions. The original considerations of the allosteric concept were developed to treat cases in which the behavior of one molecule, the substrate, depends on the concentration of another molecule. These are called *heterotropic interactions*. The regulatory enzymes whose properties were catalogued by Monod et al. showed both homotropic and heterotropic effects ([5](#)). The phenomenon of desensitization—that is, treatments that result in loss of substrate homotropic interactions and heterotropic interactions, but do not affect catalytic activity, had been described for some of these enzymes. Monod et al. postulated that desensitization means that the molecular basis might be the same for both homotropic and heterotropic interactions. They stated that if this were the case, the concerted model could also be used to describe heterotropic effects—that is, the original allostery. In the model, activating effectors have higher affinity for the R-state, while inhibitory effectors have higher affinity for the T-state. This issue has resulted in great confusion in the use of the model. All subsequent treatments presume that a single allosteric transition is the basis for both homotropic and heterotropic effects. This presumption is particularly evident in structural studies of regulatory proteins, where the terms R-state and T-state dominate discussions.

In the three cases of regulatory proteins that have been examined with care, it is clear that the concerted model fails. More than two states are involved in oxygen binding to hemoglobin ([8](#), [9](#)). Homotropic and heterotropic interactions are not due to the same transition in aspartate transcarbamoylase ([10](#), [11](#)). Although the same T-state is predicted for [phosphofructokinase](#) with different inhibitors, experimental results show that the properties of the inhibited forms are different ([12](#)). These results argue strongly that the term *allosteric* should be used in a stricter sense to refer to effects of one ligand on the binding of another, independently of assumptions about the mechanism of the interaction.

## Bibliography

1. E. H. Fischer, A. Pocker, and J. C. Saari (1971) *Essays Biochem.* **6**, 23–68.

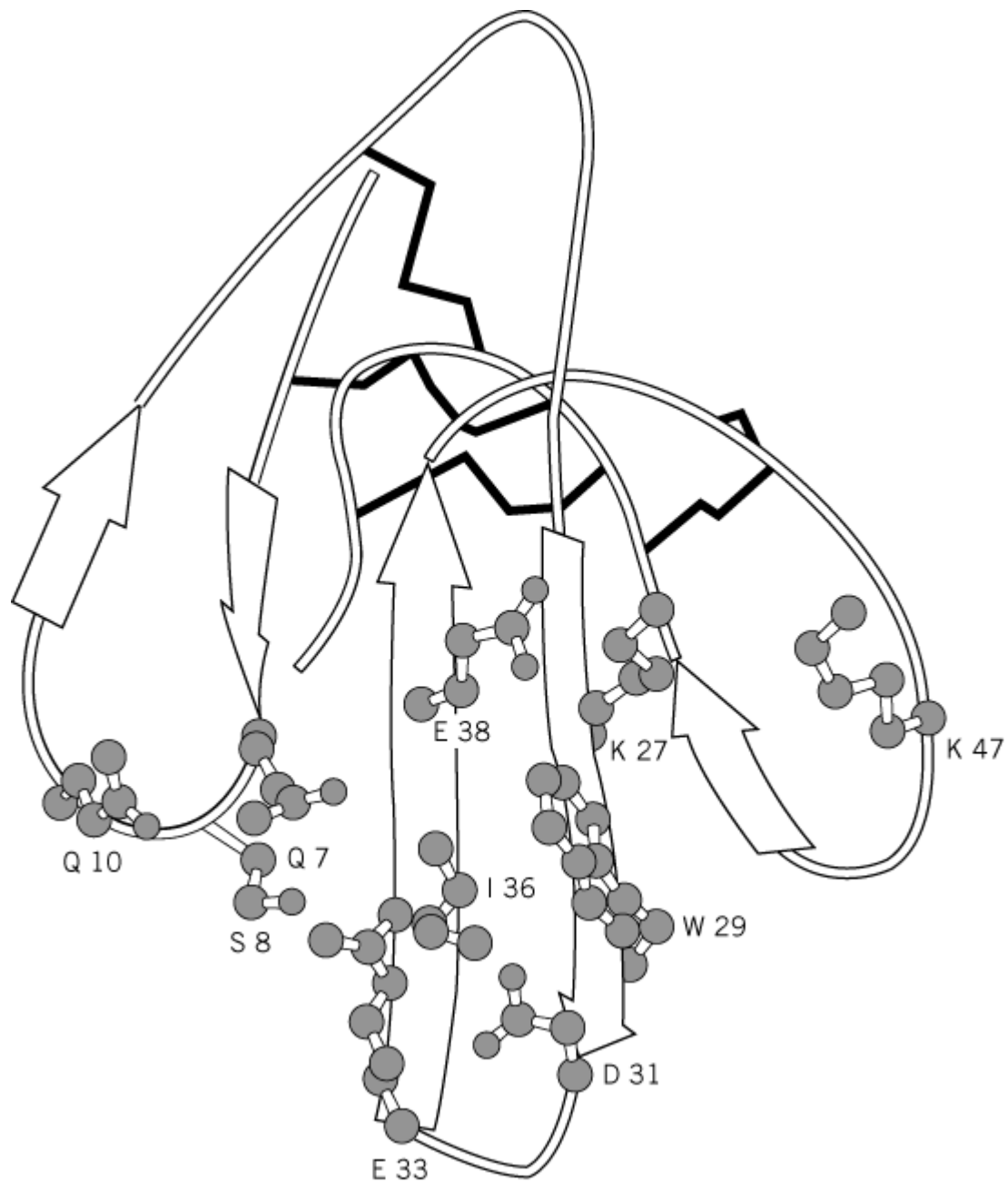
2. H. E. Umbarger (1956) *Science* **123**, 848.
3. R. A. Yates and A. B. Pardee (1956) *J. Biol. Chem.* **221**, 757–780.
4. J. Monod and F. Jacob (1961) *CSHSQB* **26**, 389–401.
5. J. Monod, J.-P. Changeaux, and F. Jacob (1963) *J. Mol. Biol.* **6**, 306–329.
6. H. F. Judson (1979) *The Eighth Day of Creation: The Makers of the Revolution in Biology*, Simon & Shuster, New York.
7. J. Monod, J. Wyman, and J.-P. Changeaux (1965) *J. Mol. Biol.* **12**, 88–118.
8. G. K. Ackers and J. H. Hazzard (1993) *Trends Biochem. Sci.* **18**, 385–390.
9. J. M. Holt and G. K. Ackers (1995) *FASEB J.* **9**, 210–218.
10. W. N. Lipscomb (1994) *Adv. Enzymol. Rel. Areas Mol. Biol.* **68**, 67–151.
11. R. C. Stevens, Y. M. Chook, C. Y. W. N. Cho, Lipscomb, and E. R. Kantrowitz (1991) *Protein. Eng.* **4**, 391–408.
12. V. L. Tlapak-Simmons and G. D. Reinhart (1994) *Arch. Biochem. Biophys.* **308**, 226–230.

## Alpha-Bungarotoxin and Curare-Mimetic Toxins

The venom of the banded krait snake (*Bungarus multicinctus*) contains a variety of protein [toxins](#) that are collectively called bungarotoxins. As is often the case, an animal venom is, in fact, a mixture of different toxins specific for different target molecules. The  $\alpha$ -neurotoxins, or curare-mimetic toxins, bind specifically to the [acetylcholine receptor](#) with high affinities ( $K_d$  in the  $10^{-9}$  to  $10^{-12}$  M range for neuronal and muscular receptors of different species) and prevent the opening of the ion **channel** caused by acetylcholine binding ([1](#), [2](#)). Such inhibition of the transmission of the nerve-muscle impulse results in a flaccid paralysis, which may end in respiratory failure and death.

$\alpha$ -bungarotoxin is the prototype of a large group of monomeric toxins, produced by many snakes, including cobra and sea snakes, and composed of 66 to 74 amino acid residues with four or five [disulfide bonds](#). Their three-dimensional structure is rather flat, with three adjacent  $\beta$ -sheet loops ([3](#)) that expose residues essential for receptor binding: Lys-27, Trp-29, Arg-39, and Lys-55 (numbering of erabutoxin) (Fig. [1](#)) ([4](#)).

**Figure 1.** Erabutoxin, a snake toxin that binds to the acetylcholine receptor. Folding of the polypeptide chain shows the typical three-finger folding of curare-mimetic snake toxins. Several residues determine the specific binding and inhibition of the acetylcholine receptor with Lys-27, Trp-29, Arg-33, and Lys 47 playing a major role ([4](#)). Reproduced from ([2](#)) with permission.



Many species of snakes, including *B. multicinctus*, also produce  $\alpha$ -neurotoxins of smaller size. Their polypeptide chains are 60 to 62 residues long but adopt the same three-finger  $\beta$ -sheet fold, stabilized by four disulfide bonds. At variance from the longer  $\alpha$ -neurotoxins, these toxins are dimers, but the same three positively charged residues are essential for receptor binding. The corresponding area of the acetylcholine receptor involved in  $\alpha$ -neurotoxin binding has not been mapped in detail, but the segment of residues 185 to 199 plays a major role, together with residues 128 to 142.

The various  $\alpha$ -neurotoxins bind with a range of affinities and specificities to the multitude of muscular and neuronal acetylcholine receptors; fluorescent and **radiolabeled** toxin derivatives are invaluable tools for the study of the structure and function of such receptors (5).

#### Bibliography

1. A. L. Harvey, ed. (1991) *Snake Toxins*, Pergamon Press, New York.
2. R. Rappuoli and C. Montecucco (1997) *Guidebook to Protein Toxins and Their Use in Cell*

*Biology*, Oxford University Press, Oxford, UK.

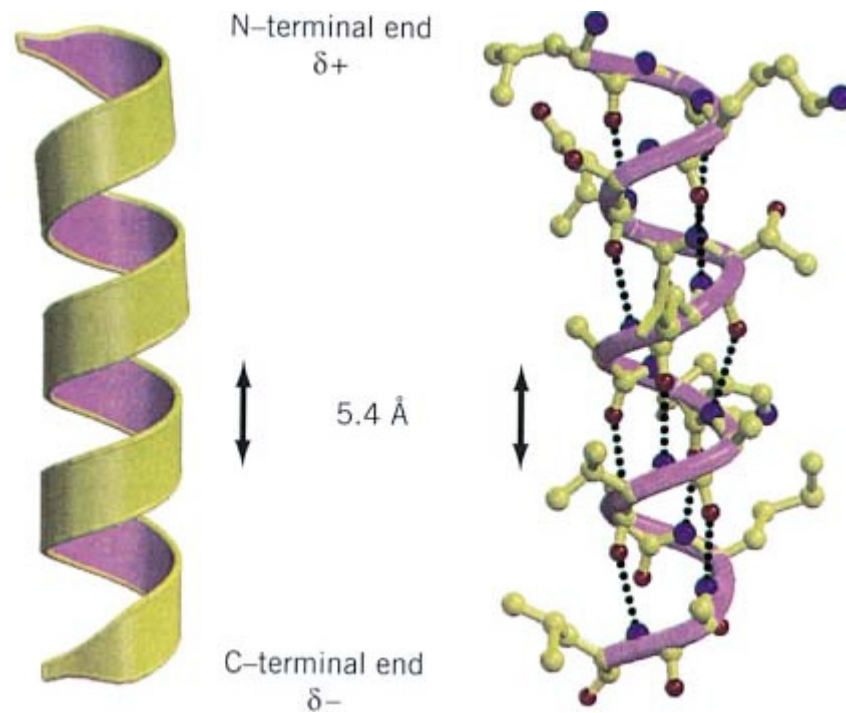
3. M. D. Walkinshaw, W. Saenger, and A. Maelicke (1980) *Proc. Natl. Acad. Sci. USA* **77**, 2400–2404.
4. O. Tremeau et al. (1995) *J. Biol. Chem.* **270**, 9362–9369.
5. A. Devillers-Thiery et al. (1993) *J. Membr. Biol.* **136**, 97–112.

### Alpha-Helix ( $3_{10}$ -Helix and Pi-Helix)

The right-handed  $\alpha$ -helix is one of two regular types of protein **secondary structure**, the other being the **b-strand** that forms **b-sheets**. Helices are formed from consecutive stretches of **residues** of the **polypeptide chain** and are characterized by a right-handed coiled **backbone** and a regular repeating pattern of backbone **hydrogen bonds**.  $\alpha$ -Helices are usually depicted as coils or cylinders in protein structure diagrams (Fig. 1). There are 3.6 residues in every turn of  $\alpha$ -helix, and for each turn the backbone is translated by 5.4 Å (or 1.5 Å per residue). In **protein structures**, the average length of an  $\alpha$ -helix is 10 residues, although much shorter and longer examples have been observed. The backbone angles are approximately  $-60^\circ$  and  $-50^\circ$  for  $\phi$  and  $\psi$ , respectively, corresponding to the allowed region in the lower left of the **Ramachandran Plot**. The atoms of the polypeptide chain pack closely together in the  $\alpha$ -helical conformation, making favorable **van der Waals interactions**. The **side chains** of each residue are oriented outward from the axis of the helix, with the Ca–Cb bond pointing toward the *N*-terminal end of the helix.

**Figure 1.** A typical protein  $\alpha$ -helix. **(Left)** The helix is depicted schematically as a coil, and the *N*- and *C*-terminal ends, with their respective partial positive and negative charges, are indicated. **(Right)** The atomic detail of the  $\alpha$ -helical structure is shown, with hydrogen bonds between backbone carbonyl oxygens (*i*) and backbone amide nitrogens of residue (*i* + 4) indicated by dotted lines. Nitrogen atoms are shown in blue and oxygen atoms are shown in red. This figure was generated using Molscrip (3) and Raster3D (4, 5). See color insert.





A regular hydrogen bond pattern is formed in the  $\alpha$ -helix, corresponding to bonds formed between the backbone carbonyl oxygen of residue ( $i$ ) and the backbone amide of residue ( $i + 4$ ) in the polypeptide chain, so that there are 13 atoms between each acceptor and donor pair. Thus, apart from the first few amide groups and last few carbonyl oxygen atoms, all the peptide bonds in an  $\alpha$ -helix form hydrogen bond interactions. The hydrogen bonds are all oriented in the same direction, as are the dipoles of each of the peptide bonds. Therefore, the helix itself also has a significant dipole moment (a partial positive charge at the  $N$ -terminus and a partial negative charge at the  $C$ -terminus), and charged residues that interact with the dipole are often found in the sequence at the appropriate end of a helix (1). For example, the negatively charged Asp residue is highly favored at the  $N$ -terminus. The helix dipole has also been implicated in protein functions, such as substrate binding and [catalysis](#).

Different amino acids have different tendencies for forming  $\alpha$ -helices (2). The amide nitrogen of [proline](#) is cyclized with its backbone and thus cannot act as a hydrogen bond donor. Proline residues therefore are not found frequently in  $\alpha$ -helices; when they are, they cause irregularities such as kinks and bends in the middle of helices. On the other hand, proline is a preferred residue at the  $N$ -terminus of the helix where other residues would have an unpaired backbone amide.

The staggered arrangement of side chains around the helix can be visualized in two dimensions through the use of a **helical wheel** (a projection down the axis of the helix). Helices in [membrane proteins](#) are often **amphipathic**, having **polar** or charged side chains arranged on one side and **hydrophobic** side chains on the opposite side of the helix.

There are several variations to the  $\alpha$ -helix conformation of proteins. The  $3_{10}$ -helix is more tightly wound than the  $\alpha$ -helix with backbone hydrogen bonds formed between residues ( $i$ ) and ( $i + 3$ ). Its name is based on its structure, with three residues per turn and 10 atoms between hydrogen bond donor and acceptor (although in real protein structures this geometry may be somewhat distorted). The  $3_{10}$ -helix is usually only observed at the ends of  $\alpha$ -helices or in short stretches of 4 to 5 residues. The P-helix is more loosely wound than the  $\alpha$ -helix with backbone hydrogen bonds formed between the ( $i$ ) and ( $i + 5$ ) residues and is very rarely observed in protein structures. The left-handed helix, having the same hydrogen bonding pattern as  $\alpha$ -helices but backbone  $\phi$  and  $\psi$  angles of  $+60^\circ$  and

+50° (rather than the -60° and -50° angles of the right-handed  $\alpha$ -helix), is not often observed, because this conformation results in steric clashes between backbone and side chain atoms. A [glycine](#) residue, however, having only a hydrogen atom as its side chain, can adopt this conformation. The poly(Pro)II helix, with no backbone hydrogen bonds, is also left-handed, with three residues per turn and a translation of 3.12 Å per residue. This helical structure is observed in [polyproline](#) and in proline-rich regions of proteins.

[See also [Protein Structure](#) and [Secondary Structure, Protein](#).]

### Bibliography

1. L. Serrano and A. R. Fersht (1989) *Nature* **342**, 296–299.
2. J. S. Richardson and D. C. Richardson (1988) *Science* **240**, 1648–1652.
3. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
4. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
5. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

6. D. J. Barlow and J. M. Thornton (1988) Helix geometry in proteins. *J. Mol. Biol.* **201**, 601–619.
7. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
8. W. G. J. Hol, P. T. Van Duijnen, and H. J. C. Berendsen (1978) The  $\alpha$ -helix dipole and the properties of proteins. *Nature* **273**, 443–446.
9. J. S. Richardson (1981) The anatomy and taxonomy of protein structure. *Adv. Protein Chem.* **34**, 167–339.
10. C. Toniolo and E. Benedetti (1991) The polypeptide  $3_{10}$  helix. *TIBS* **16**, 350–353.

## Alpha-Helix Formation

$\alpha$ -Helices are the most common type of **secondary structure** in globular proteins (see [Alpha-Helix \(310-Helix and Pi-Helix\)](#)). As such, there is intense interest in understanding the factors that contribute to  $\alpha$ -helix formation in peptides and proteins ([1-4](#)).

All peptide helices unfold with increasing temperature. This indicates that helix formation is an enthalpically favorable process (it proceeds with the release of heat). The main feature of helical structure is the repeating pattern of intramolecular [hydrogen bonds](#) formed between atoms in the peptide backbone ( $>C=O \cdot H-N<$ ). These hydrogen bonds probably provide the main enthalpic stabilization to the helical conformation ([5](#)). Opposing this **enthalpy** is the conformational **entropy** of freezing the backbone conformation, because the process of change from a disordered conformation (unfolded peptide or protein) to an ordered one (helix) is unfavorable entropically. This unfavorable entropy is about the same magnitude as the favorable enthalpy; that is, in isolation the helices are at best marginally stable.

Because the  $\alpha$ -helix has an ( $i, i + 4$ ) backbone hydrogen bonding pattern, the first four  $>N-H$  and the last four  $>C=O$  groups of the helix will not be hydrogen-bonded to a backbone partner. Polar groups that do not form hydrogen bonds to [water](#) or other protein groups are destabilizing (see [Protein Stability](#)). Therefore it is important that these groups either be hydrated or have other intramolecular partners. The side chain of the first residue of the helix (the “ $N$ -cap”) often fills this

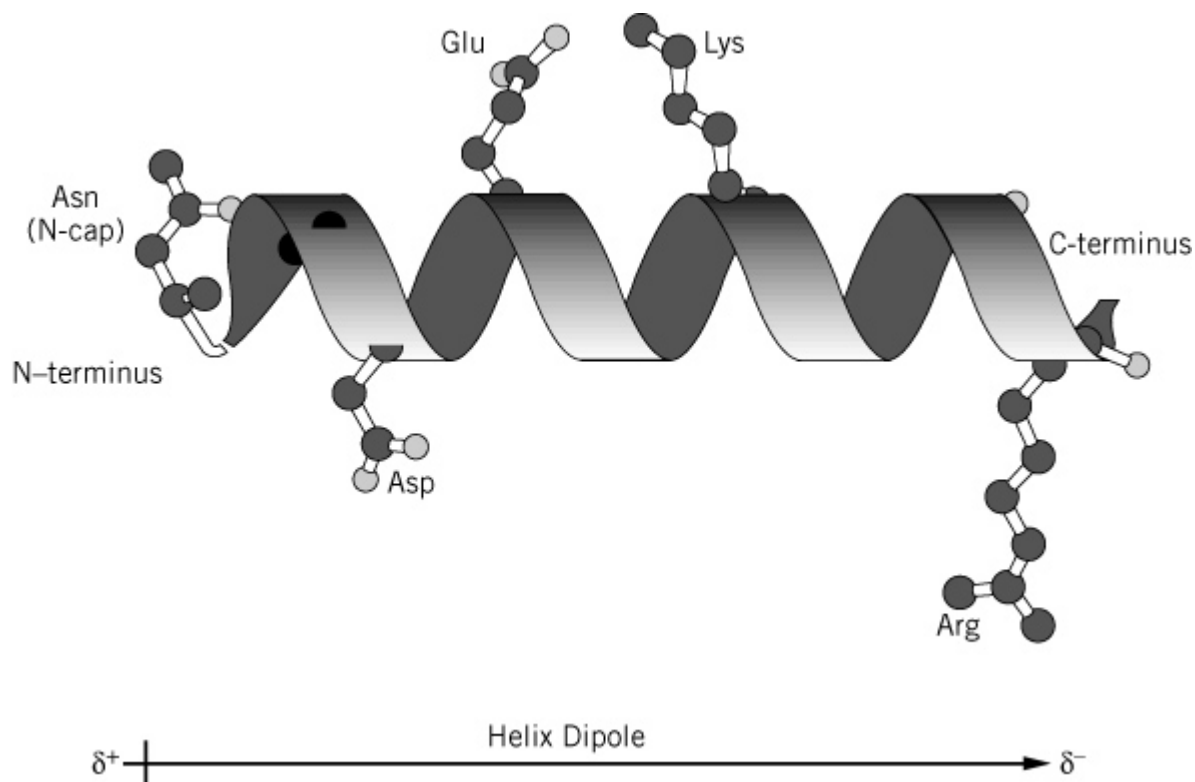
role, if it is able to hydrogen-bond to one or two otherwise unsatisfied backbone  $>N-H$  groups. These “capping” hydrogen bonds have been proposed to be important factors in determining where helices begin and end in proteins (6, 7).

The helical conformation has the peptide groups and their associated dipoles aligned in the same direction, whereas in nonhelical forms the dipoles are distributed randomly. As the helix forms, the dipoles on adjacent peptide groups align in an unfavorable fashion (see [Electrostatic Interactions](#)). However, when a full turn of the helix is completed, the first hydrogen bond is formed and the peptide groups involved in the hydrogen bond align favorably. This adds to the **cooperativity** of helix formation: Once the first hydrogen bond is formed, it is easier to continue forming the helix. The alignment of the individual peptide group dipoles in a single direction gives rise to a helix *macro-dipole* (simply the additive effects of the individual dipoles). The *N*-terminus of the helix has a partial positive charge, while the *C*-terminus carries a partial negative charge. Residues with charged side chains at the ends of a helix can interact favorably or unfavorably with the helix macrodipole, giving rise to an effect known as the *charge–dipole interaction*.

Different amino acid residues have different tendencies to form helices (see [Helix–Coil Theory](#)). For example, alanine is found in protein helices much more often than is glycine (8). The intrinsic preference to be in a helical conformation has been termed *helix propensity*. Since this is an intrinsic property, it should be characteristic of a given amino acid and independent of context. However, quantitative measurements of helix propensity in different systems have provided markedly different results (1, 2, 4). This reflects the difficulty of measuring intrinsic properties in complex molecules like proteins and peptides (however, see Ref. 9).

Interactions between side chains can also affect the stability of helices. Side chains three and four residues away from each other in the amino acid sequence will be close in space in a helical conformation. Therefore, electrostatic interactions, hydrogen bonds, or **hydrophobic** contacts between these side chains can contribute to helix stability. Likewise, if interactions are possible in the nonhelical form of the peptide or protein, this can alter the amount of  $\alpha$ -helix observed in a peptide or protein. Figure 1 shows examples of various types of side-chain interactions possible in the helix, an estimate of the magnitudes of the contributions of the main factors contributing to helix stability, and comparison of these estimates with the **free energy** provided by typical interactions in proteins.

**Figure 1.** Interactions involving amino acid side chains that affect the stability of isolated  $\alpha$ -helices. (**Top**) Ribbon model of an  $\alpha$ -helix, with stabilizing interactions involving (i) hydrogen bonding of the *N*-cap side-chain at the *N*-terminus, indicated here as Asn, (ii) favorable electrostatic interactions of ionized side chains (Asp and Arg shown here) with the helix dipole, and (iii) electrostatic interactions between oppositely charged side chains on adjacent turns of the helix. (**Bottom**) Quantitative estimates of the net stabilizing interactions in an isolated  $\alpha$ -helix and in folded globular proteins.



In addition to the stability of helices, their rates of folding and unfolding are of interest as well. The transition from helix to random coil occurs very rapidly, on the order of nanoseconds, two to three orders of magnitude faster than the fastest folding proteins (10).

#### Bibliography

1. J. M. Scholtz and R. L. Baldwin (1992) *Annu. Rev. Biophys. Biomol. Struct.* **21**, 95–118.
2. A. Chakrabarty and R. L. Baldwin (1995) *Adv. Protein Chem.* **46**, 141–176.
3. J. W. Bryson, S. F. Betz, H. S. Lu, D. J. Suich, H. X. Zhou, K. T. O'Neil, and W. F. DeGrado (1995) *Science* **270**, 935–941.
4. N. R. Kallenbach, P. Lyu, and H. Zhou (1996) In *Circular Dichroism and the Conformational Analysis of Biopolymers* (G. D. Fasman, ed.), Plenum Press, New York.
5. J. M. Scholtz, S. Marqusee, R. L. Baldwin, E. J. York, J. M. Stewart, M. Santoro, and D. W. Bolen (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2854–2858.
6. L. G. Presta and G. D. Rose (1988) *Science* **240**, 1632–1641.
7. E. T. Harper and G. D. Rose (1993) *Biochemistry* **32**, 7605–7609.
8. P. Y. Chou and G. D. Fasman (1978) *Adv. Enzymol.* **47**, 45–148.
9. J. K. Myers, C. N. Pace, and J. M. Scholtz (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2833–2837.
10. S. Williams, T. P. Causgrove, R. Gilmanishin, K. S. Fang, R. H. Callender, W. H. Woodruff, and R. B. Dyer (1996) *Biochemistry* **35**, 691–697.

#### Suggestion for Further Reading

11. The first four references cited above provide excellent reviews of helix formation in peptides and the relationship between helix formation in peptides and proteins.

## Alpha-Lactalbumin

It requires patience to achieve quiet research. The investigation of properties of  $\alpha$ -lactalbumin ( $\alpha$ -LA) was not worthy of much notice until about 1965, especially from the standpoint of molecular biology. It is only a protein component of the complex enzyme “lactose synthase” and elucidation of its function had made little progress. However, the hidden connection to molecular biology was ferreted out in about 1967 and initiated the great explosion of activity in studying the [molten globule](#) (MG) conformation of proteins. These studies of the MG proteins opened a new area of structural biology connected to various physiological functions. Modern molecular biology sometimes requires the MG protein model to understand such physiological processes. The best-characterized MG is that formed by  $\alpha$ -LA, and it can easily be prepared.  $\alpha$ -LA is now one of the star proteins in molecular biology. The development of research into  $\alpha$ -LA has been deeply affected by changing times.

One (a) of three peaks in the **sedimentation velocity** patterns of proteins in the non-casein fraction of skim milk was found responsible for a lactalbumin isolated from milk, which has subsequently been called  $\alpha$ -lactalbumin. In the middle-1960s, it was confirmed that lactose is synthesized in the mammary gland UDP-galactose (UDP-Gal) and glucose (Glc) by an enzyme “lactose synthase,” which can be resolved into two fractions. One of them is  $\alpha$ -LA, and the other is galactosyltransferase (GT). Alone in the [Golgi apparatus](#), but with metal ions such as  $Mn^{2+}$ , the latter catalyzes the transfer of Gal to GlcNAc on **glycoproteins**:



where NAc is the NAc is the N-acetyl group.  $\alpha$ -LA is synthesized in the mammary gland during the lactation. It passes through the Golgi apparatus where it combines with GT and alters its substrate specificity to inhibit reaction 1. A decrease of three of magnitude in the  $K_m$  for Glc facilitates the synthesis of lactose instead:

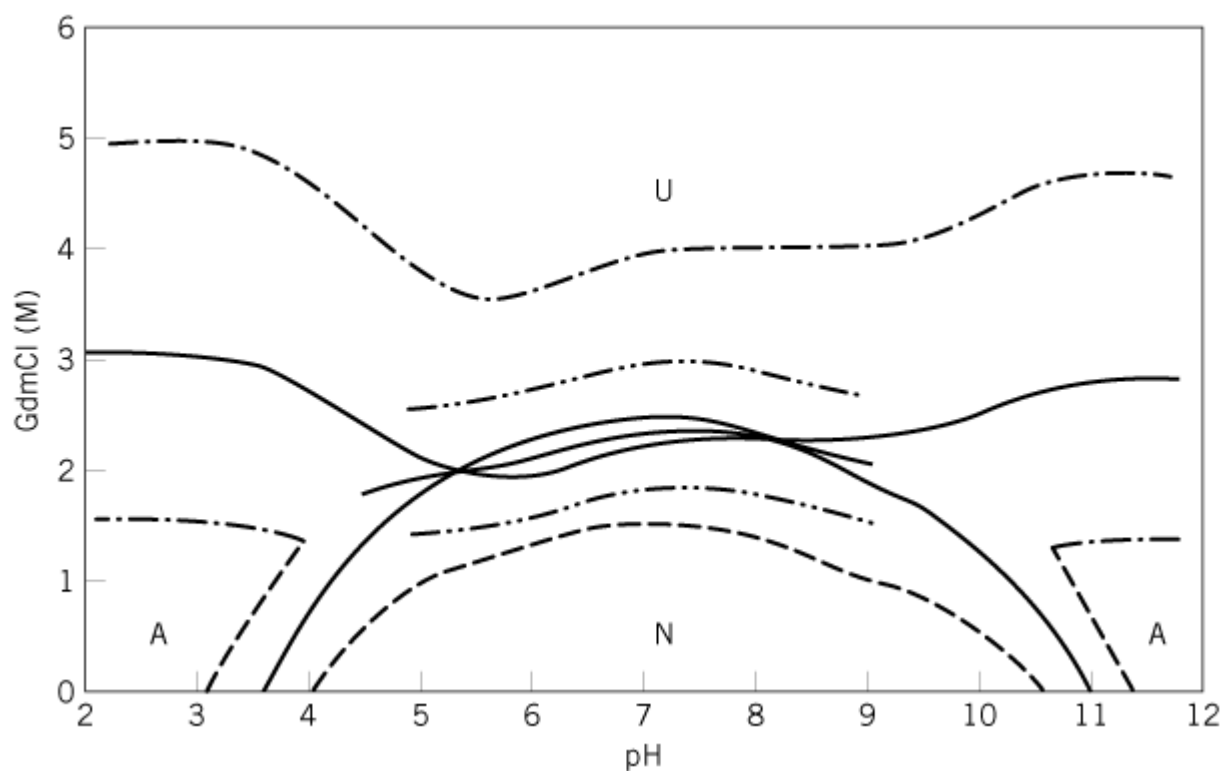


### 1. Structures

$\alpha$ -LA is normally a single polypeptide chain of 123 residues, four [disulfide bonds](#), and a molecular weight of about 14,000 except for rat  $\alpha$ -LA, which has 140 residues. About 1970, similarities were noted in the molecular weight and amino acid sequence between  $\alpha$ -LA and hen egg white [lysozyme](#) (HEWL), a lytic enzyme. Subsequently, it was found that the exon-intron organization of the genes for both proteins are the same [1](#), [2](#). These suggested **divergent** evolution of the two proteins from a common ancestral gene, and  $\alpha$ -LA was assumed to have a three-dimensional structure similar to that of HEWL, in spite of marked differences in their functions. **X-ray crystallographic** determination of the  $\alpha$ -LA structure was not successful for a long time, but the conformations of both proteins in aqueous media were compared. Until 1980, significant differences in their physical behavior were observed: 1) An intermediate conformation appears in the unfolding equilibrium of  $\alpha$ -LA, but not of HEWL, at intermediate concentrations of **guanidinium chloride** (GdmCl) at a neutral pH, which is also stable at acid pH and is designated the “A-State”; 2) One  $Ca^{2+}$  ion is tightly bound to  $\alpha$ -LA, but not to HEWL, and stabilizes the native (N) form. Removal of  $Ca^{2+}$  from  $\alpha$ -LA induces the A-state near room temperature. Figure [1](#) shows a phase diagram of bovine  $\alpha$ -LA for the N-, A-, and U- (unfolded) states in of GdmCl at 25°C, and the unfolding equilibrium has been explained in terms of

the three states (3-5). Subsequently, it was shown that the kinetic intermediate of bovine a-LA captured at the early stage of refolding from U has properties similar to the equilibrium A-state. Later, it was reported that a number of proteins including HEWL assume the equilibrium and kinetic intermediate(s) similar to the A-state of a-LA in the molten globule (MG) form, and the MG was considered to be a third physical state of proteins (see [Molten Globule](#)). However, the conformational transition of the A to U of a-LA is nonco-operative, and then the folding of a-LA along the multipathway, which is apparently in the two-state type, is proposed (6-8). Also, recently, the kinetic intermediates of folding of a-LA with non-native secondary structures were shown, especially in the apo-form (9).

**Figure 1.** Phase diagram for the three states (N, A, and U) of bovine a-lactalbumin in a solution of GdmCl at any concentration and pH at 25°C.\* The solid curves represent the pH and [GdmCl] where  $K_{eq}^{NU} K_{eq}^{AU}$ , or  $K_{eq}^{NU}$  is unit.  $K_{eq}^{ij}$  is for the conformational transition between  $i$  and  $j$ . The other curves represent the pH and [GdmCl] where  $K_{eq}^{NA} = 10$  or  $0.1$  (---),  $K_{eq}^{AU} = 10$  or  $0.1$  (-.-), and  $K_{eq}^{NU} = 10$  or  $0.1$  (-. .-). \*K. Kuwajima, Y. Ogawa, and S. Sugai (1983) *J. Biochemistry* (Tokyo) **89**, 759–770.



Since 1987, the crystal structures of a-LA from various species (baboon, human, goat, guinea pig, and bovine) and their recombinant and mutant forms have been determined at resolutions of 1.15 to 2.4 Å (10-15). The overall features of the authentic forms are similar to those of HEWL. Both structures are divided into two (a and b) subdomains by a deep cleft, the a-domain comprises residues 1 to 34 and 86 to 123, and the b-domain includes residues 35 to 85. In baboon a-domain contains four a-helices, A (residues 5 to 11), B (23 to 34), C (86 to 99), and D (105 to 109), plus single turns of  $3_{10}$ -helices (12 to 16, 101 to 104, and 115 to 119). The b-domain consists of a three-stranded antiparallel b-sheet (41 to 56), a  $3_{10}$ -helix (76 to 82), and loop regions (Fig. 2). In a-LA from other species, however, the D-helix is replaced by a loop. The C-terminal region is flexible. The conformation and the ligand coordination of a  $Ca^{2+}$ -binding site ( $K_a$  of  $10^9$  to  $10^{10} M^{-1}$  in the N-state) are observed at the interface of both the domains in all of a-LA. The  $Ca^{2+}$ -binding site is

composed of three Asp residues at 82, 87, and 88 and peptide carbonyl oxygens at residue 79 and 84 (the other two ligands are water molecules) and is in a slightly disordered, rigid pentagonal bipyramidal form (designated the “a-LA elbow”). The Asp residues are not found in the corresponding region of most lysozymes, but some have the Ca<sup>2+</sup>-binding site in an a-LA elbow (Table 1), which bind a Ca<sup>2+</sup> ion in the N-state and are important for studying the [evolution](#) of the a-LA family. A Zn<sup>2+</sup>-binding site was also found in the human a-LA crystal, located at the entrance of the cleft between the two subdomains and coordinated with Glu 49, Glu 116, and two water molecules. Also, the Mn<sup>2+</sup>-and weak Ca<sup>2+</sup>-binding sites were also indicated. a-LA binds numerous other metal ions: Na<sup>+</sup>, K<sup>+</sup>, Ba<sup>2+</sup>, Mg<sup>2+</sup>, Co<sup>2+</sup>, Cu<sup>2+</sup>, Pb<sup>2+</sup>, Hg<sup>2+</sup>, Cd<sup>2+</sup>, Sr<sup>2+</sup>, Al<sup>3+</sup>, Tb<sup>3+</sup>, Eu<sup>3+</sup>, Sc<sup>3+</sup>, and Y<sup>3+</sup>. However, elucidation of their locations and roles in lactose synthesis are awaited. Recently, binding of fatty acids to human a-LA is noted in connection with its function. The structures of a-LA determined by X-ray crystallography and by NMR in solution exhibit two **hydrophobic** clusters of aromatic residues: cluster I of Phe 31, His 32, Tyr 36, and Trp 118 and cluster II of Trp 26, Phe 53, Trp 60, and Trp 104 ([16-19](#)). These make up the cores in a-LA by combining with other hydrophobic parts in the a-domain.

**Table 1. Comparison of “Equivalent Amino Acid Residues” in Some a-Lactalbumins and Lysozymes at the Calcium-Binding Sites**

| Protein       | Source     | Residue <sup>a</sup> |     |     |     |     |
|---------------|------------|----------------------|-----|-----|-----|-----|
|               |            | 79                   | 82  | 84  | 87  | 88  |
| a-Lactalbumin | Baboon     | Lys                  | Asp | Asp | Asp | Asp |
|               | Human      | Lys                  | Asp | Asp | Asp | Asp |
|               | Bovine     | Lys                  | Asp | Asp | Asp | Asp |
|               | Equine     | Lys                  | Asp | Asp | Asp | Asp |
|               | Wallaby    | Lys                  | Asp | Asp | Asp | Asp |
|               | Guinea pig | Lys                  | Asp | Asp | Asp | Asp |
|               | Goat       | Lys                  | Asp | Asp | Asp | Asp |
|               | Camel      | Lys                  | Asp | Asp | Asp | Asp |
|               | Rabbit     | Asn                  | Asp | Asp | Asp | Asp |
|               | Rat        | Lys                  | Asp | Gly | Asp | Asp |
| Lysozyme      | Human      | Ala                  | Gln | Asn | Asp | Ala |
|               | Bovine     | Glu                  | Glu | Asp | Lys | Ala |
|               | Hen egg    | Ala                  | Ser | Asp | Ala | Ser |
|               | Equine     | Lys                  | Asp | Asn | Asp | Asp |
|               | Pigeon     | Lys                  | Asp | Asn | Asp | Asp |
|               | Canine     | Lys                  | Asp | Asn | Asp | Asp |

<sup>a</sup> Residue numbers for lysozyme are equivalent numbers for human a-lactalbumin.

## 2. Molten Globule Intermediate

The [molten globule](#) of a-LA is compact (a radius of gyration only about 10% larger than that of the N-form) and has N-like secondary structures, but without the specific **tertiary structural** packing

interactions.  $\alpha$ -LA assumes the following stable MG forms in aqueous media: 1) An equilibrium unfolding state at intermediate concentrations of GdMCl or [urea](#); 2) The acid state; 3) A partially unfolded state obtained by removing of the bound  $\text{Ca}^{2+}$  at neutral pH and low salt concentrations, and 4) After partial reduction of the disulfide bonds. NMR and amide [hydrogen exchange](#) show heterogeneous structures of the MG, an  $\alpha$ -domain relatively structured by hydrophobic interactions, and a more unfolded  $\beta$ -domain. Hydrophobic clusters I and II in the native form are somewhat rearranged locally, and some form of the hydrophobic core persists in the MG form and play as nuclei for the folding. The disulfide bond between Cys 6 and Cys 120 in native  $\alpha$ -LA can be reduced extremely quickly by thiol agents. The free [thiol groups](#) in the three-disulfide intermediate (3SS) of  $\alpha$ -LA are easily modified by **iodoacetic acid** or iodoacetoamide and generate carboxymethylated or carboxyamidomethylated 3SS  $\alpha$ -LA, respectively. These trapped derivatives of bovine  $\alpha$ -LA retain the N-form, with slight changes in the local conformation very near Cys 6 and Cys 120, but they assume the MG form in the absence of  $\text{Ca}^{2+}$  or at acid pH ([20-25](#)). A two-disulfide species and its trapped derivatives can also be obtained by subsequently reducing the 28–111 disulfide bond. The 2SS derivatives have partial MG characteristics. These disulfide intermediates and their derivatives are frequently used as models of partly unfolded proteins in the molecular biology because they are well characterized and easily prepared ([26-28](#)).

### 3. Functions

Lactose synthase (reaction 2) has long been studied to determine the binding sites of the saccharides and of metal ions and the interaction site between GT and  $\alpha$ -LA. The structure of  $\alpha$ -LA-GT complex has not been determined, but it has been shown by indirect methods such as [chemical modification](#) and [site-directed mutagenesis](#) that  $\text{Ca}^{2+}$  and some residues of  $\alpha$ -LA adjacent to the cleft including Phe 31, His 32, Ala 106, and Leu 110 are crucial for the function of lactose synthase. The flexible C-terminal region, part of the hydrophobic clusters, and the cleft are essential for the lactose synthase function. The binding sites of  $\alpha$ -LA on GT have been tentatively identified. The difference in function between  $\alpha$ -LA and lysozyme has been explained by roles of some residues such Tyr 103 in  $\alpha$ -LA. Mutations of six residues in the cleft of  $\alpha$ -LA to the corresponding ones of lysozyme created a catalytic site of lysozyme in  $\alpha$ -LA and hydrolyzed the glycoside bond, as chicken lysozyme ([29](#)).

It has been confirmed that **molecular chaperones** bind the protein and assist it folding in the cell. Physicochemical studies *in vitro* indicate that the **chaperonin** GroEL binds apo- or disulfide-reduced  $\alpha$ -LA in the MG form, although it scarcely interacts with the N-state of  $\alpha$ -LA ([30-34](#)). GroEL recognizes the hydrophobic surface exposed on the MG-form.  $\alpha$ -LA bound to GroEL is also in the MG-form.  $\alpha$ -LA and its disulfide-reduced forms in the MG interact with other chaperones and are frequently used as model substances for studies of protein folding in the cell. The insertion of soluble proteins into [membranes](#) and the conformations of the membrane-bound protein have been topics of interest ([35-38](#)). The insertion of  $\alpha$ -LA into model membranes occurs under conditions that favor formation of the MG with its hydrophobic surface. The inserted  $\alpha$ -LA is also in the MG-form, and its association with a lipid bilayer affects the chain mobility of the lipids.

Recently, some multimers of human  $\alpha$ -LA prepared from its milk casein fraction were indicated to induce **apoptosis** in cancer cells, leukoma (L210), and lung carcinoma (A549) but not in healthy cells, although the monomeric human  $\alpha$ -LA did not induce apoptosis in any cell ([39-42](#)). The active **high oligomers** are folding variants of human  $\alpha$ -LA, to which fatty acids are bound, and assume the MG form. They were shown to accumulate in the nuclei of sensitive cells and to induce DNA fragmentation. The latter process requires  $\text{Ca}^{2+}$  ions.  $\alpha$ -LA has important biological functions in addition to lactose synthesis.

$\alpha$ -LA is a protein of interesting structure, properties, and functions, and it will be noted in molecular biology in the future.

### Bibliography



1. D. C. Phillips et al (1986) *Biochem. Soc. Trans.* 15, 737–744.
2. P. K. Qasba and S. Kumar (1997) *Crit. Rev. Biochem. Mol. Biol.* 32, 255–306.
3. K. Kuwajima et al (1976) *J. Mol. Biol.* 106, 359–373.
4. K. Kuwajima et al (1977) *J. Mol. Biol.* 114, 241–258.
5. M. Mizuguchi et al (2000) *Proteins Struct. Funct. Genet.* 38, 407–413.
6. W. Pfeil (1998) *Proteins Struct. Funct. Genet.* 30, 43–48.
7. M. Nozaka et al (1978) *Biochemistry* 17, 3753–3758.
8. M. Ikeguchi et al (1998) *Protein Sci.* 7, 1564–1574.
9. A. Troullier et al (2000) *Nat. Struct. Biol.* 7, 78–86.
10. R. Acharya et al (1989) *J. Mol. Biol.* 208, 99–127.
11. R. Acharya et al (1991) *J. Mol. Biol.* 221, 571–581.
12. R. Acharya et al (1996) *Structure* 4, 691–703.
13. R. Acharya et al (1998) *Biochemistry* 37, 4767–4772.
14. K. Harata et al (1999) *J. Mol. Biol.* 287, 347–358.
15. K. Harata et al (2000) *J. Biol. Chem.* 275, 37021–37029.
16. A. T. Alexandrescu et al (1992) *Eur. J. Biochem.* 210, 699–709.
17. A. T. Alexandrescu et al (1993) *Biochemistry* 32, 1707–1718.
18. C. M. Dobson et al (1999) *J. Mol. Biol.* 286, 1567–1580.
19. C. M. Dobson et al (1999) *J. Mol. Biol.* 288, 673–688.
20. J. J. Ewbank and T. E. Creighton (1991) *Nature* 350, 518–520.
21. J. J. Ewbank and T. E. Creighton (1993) *Biochemistry* 32, 3677–3707.
22. J. J. Ewbank and T. E. Creighton (1994) *Biochemistry* 33, 1534–1538.
23. P. S. Kim et al (1995) *Trans. R. Soc. London B* 348, 43–47.
24. P. S. Kim et al (1996) *Biochemistry* 35, 859–863.
25. D. F. Moriarty et al (2000) *Biochim. Biophys. Acta* 1476, 9–19.
26. K. Brew et al (1991) *J. Biol. Chem.* 266, 698–703.
27. K. Brew et al (1996) *Biochemistry* 35, 9710–9715.
28. K. Brew et al (1999) *Protein Eng.* 12, 581–587.
29. Y. Xue et al (2001) *Proteins Struct. Funct. Genet.* 42, 17–22.
30. M. K. Hayer et al (1994) *EMBO J.* 13, 3192–3202.
31. C. V. Robinson et al (1994) *Nature* 372, 646–651.
32. K. Kuwajima et al (1996) *J. Mol. Biol.* 258, 827–838.
33. K. Kuwajima et al (1996) *J. Mol. Biol.* 264, 643–649.
34. K. Kuwajima et al (1999) *J. Mol. Biol.* 293, 125–137.
35. S. Bunuelos and A. Muga (1995) *J. Biol. Chem.* 270, 29910–29915.
36. S. Bunuelos and A. Muga (1996) *FEBS Lett.* 386, 21–25.
37. S. Bunuelos and A. Muga (1996) *Biochemistry* 35, 3892–3898.
38. K. M. Cauthern et al (1996) *Protein Sci.* 5, 1349–1405.
39. M. Svanborg et al (1995) *Proc. Natl. Acad. Sci. U.S.A.* 92, 8064–8068.
40. M. Svanborg et al (1999) *J. Biol. Chem.* 274, 6386–6396.
41. M. Svanborg et al (2000) *Proc. Natl. Acad. Sci. U.S.A.* 97, 4221–4226.
42. M. Svanborg et al (2001) *Eur. J. Biochem.* 268, 186–191.

### **Suggestions for Further Reading**

43. L. J. Berliner and D. Johnson (1988) -Lactalbumin and calmodulin, In *Calcium Binding*

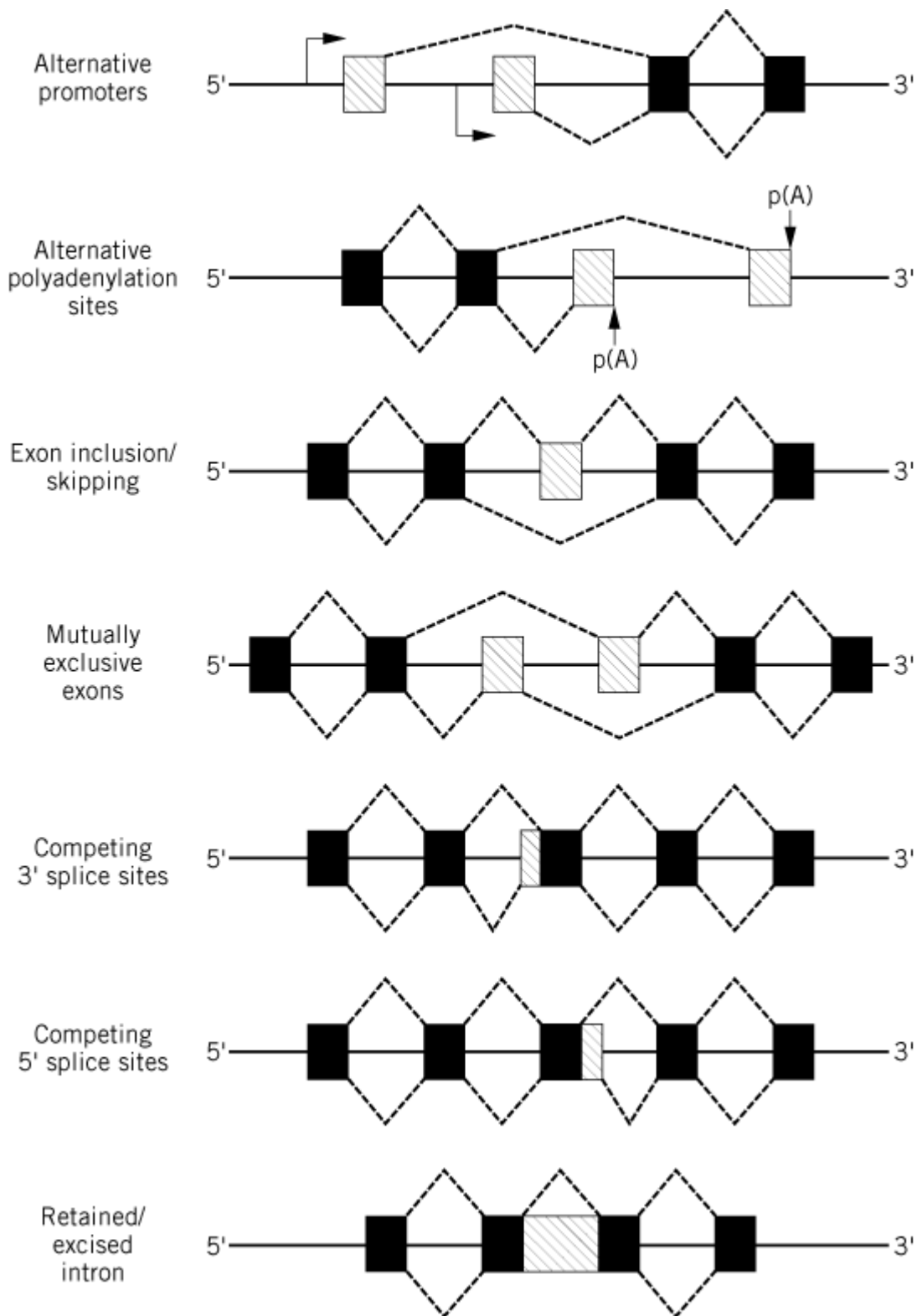
Proteins, Vol. II ( M. P. Thompson, ed.), CRC Press, Boca Raton, FL, pp. 79–116.

44. M. J. Kronman (1989) Metal-ion binding and the molecular conformational properties of -lactalbumin. *Crit. Rev. Biochem. Mol. Biol.* **24**, 565–667.
45. K. Kuwajima (1989) The molten globule state as a clue for understanding the folding and cooperatively of globular protein structure. *Proteins: Struct. Funct. Genet.* **6**, 87–103.
46. H. A. McKenzie and F. H. White Jr. (1991) Lysozyme and -lactalbumin: Structure, function and interrelationship *Adv. Protein Chem.* **41**, 171–315.
47. K. Brew and J. A. Grobler (1992) -Lactalbumin *Advanced Dairy Chemistry Proteins, Vol. 1* ( P. F. Fox, ed.), Elsevier Applied Science London and New York, pp. 191–229.
48. S. Sugai and M. Ikeguchi (1994) Conformational comparison between -lactalbumin and lysozyme. *Adv. Biophys.* **30**, 37–84.
49. K. Kuwajima (1996) The molten globule state of -lactalbumin. *FASEB J.* **10**, 102–109.
50. E. A. Permyakov and L. J. Berliner (2000) -Lactalbumin: structure and function. *FEBS Lett.* **473**, 269–274.

## Alternative Splicing

Most eukaryotic **genes** that code for proteins contain noncoding sequences (**introns**) that are interspersed among the coding regions (exons). [Transcription](#) of these genes generates so-called pre-mRNA molecules, which are converted to mature mRNAs by a process termed RNA **splicing**. During this process, the introns are precisely excised and the exons are ligated together (see [RNA Splicing](#)). The majority of nuclear pre-mRNAs are spliced constitutively; that is, only one mature mRNA species is generated from a single pre-mRNA in all tissues. In some cases however, alternative 5' and/or 3' splice sites are used during splicing, resulting in the production of more than one mRNA species from a single pre-mRNA. Alternative splicing has been documented for many eukaryotic genes, and a variety of alternative splicing patterns have been observed, as depicted schematically in [Figure 1](#). The utilization of alternative 5' and/or 3' [splice sites](#) (also referred to as donor and acceptor sites, respectively) can result in structurally distinct mRNAs by either excluding potential exon sequences or incorporating otherwise noncoding intron sequences. For some pre-mRNAs, alternative splicing is a nonregulated event such that two or more alternatively spliced mRNAs are produced at a given ratio to one another in all cell types. For others, the choice of alternative splice sites is regulated in a tissue-specific or developmental manner. This type of regulation is mediated by [trans-acting](#) factors that are differentially expressed in a particular tissue or at a specific time during development (see text below). These *trans*-acting factors may be positive or negative regulators that activate or repress the use of an alternative splice site either directly or indirectly, for example by modulating the affinity of general splicing factors.

**Figure 1.** Patterns of alternative splicing. Constitutively spliced exons are shown as black boxes, alternatively spliced exons are shown as shaded boxes, and introns are shown as a solid lines between the exons. Transcription start sites are indicated by an arrow, and polyadenylation sites are indicated by p(A). Splicing events are depicted by a dashed line.



Alternative splicing can lead to both quantitative and qualitative changes in gene expression. Quantitative changes can arise if the alternatively spliced mRNA contains a prematurely truncated open reading frame (ie, due to the presence of a [stop codon](#)) or exhibits an altered stability or [translation](#) efficiency. In many cases, alternative splicing leads to the production of so-called protein isoforms that are structurally identical in all but a specific region or domain (1). Such structural variants of a given protein often exhibit significant functional differences. Thus, through the generation of multiple protein isoforms, alternative splicing can enhance the phenotypic variability

of a single gene.

A central question in both constitutive and alternative splicing that has yet to be clearly resolved is how the correct pairs of 5' and 3' splice sites are selected for cleavage and subsequent ligation (see also [Splice Sites](#)). Pre-mRNAs contain multiple authentic 5' and 3' splice sites, which must be properly paired in order to prevent the random skipping of one or more exon. Furthermore, these sites must be distinguished from other nonauthentic sites that are also repeatedly present in most pre-mRNAs. In the case of alternative splicing, several different factors appear to be responsible for the preferential selection of one splice site over another. Firstly, features of the pre-mRNA itself (so-called [cis-acting](#) elements) can influence splice-site utilization. For example, the conserved sequences found at the 5' and 3' splice sites and the branch site contribute to splice site selection (see [Splice Sites](#)). In recent years it has become clear that the relative strengths of competing splice sites is often a deciding factor in determining which site is used preferentially. The strength of a particular 5' or 3' splice site is generally a measure of how efficiently it binds spliceosomal components, such as the U1 and U2 snRNPs or the splicing factor U2AF, which play important roles during the early stages of [spliceosome](#) assembly. This in turn is often, but not always, a function of how closely its sequences match the 5' and 3' splice site consensus sequences, or a function of the length and uridine content of the polypyrimidine tract, which determines its affinity for U2AF. The selection of weak 5' and 3' splice sites, on the other hand, can be enhanced by splicing factors that promote U1 or U2 snRNP binding (eg, SR proteins; see text below). Some 5' splice sites also activate usage of an upstream 3' splice site (eg, in the preprotachykinin pre-mRNA); and, vice versa, some 3' splice sites can promote the use of a downstream 5' splice site (eg, in the B-tropomyosin pre-mRNA) ([2](#), [3](#)). In the former case, factors bound at the 5' splice site (the U1 snRNP) enhance the interaction of U2AF with the upstream 3' splice site via a network of molecular interactions across the exon ([4](#), [5](#))

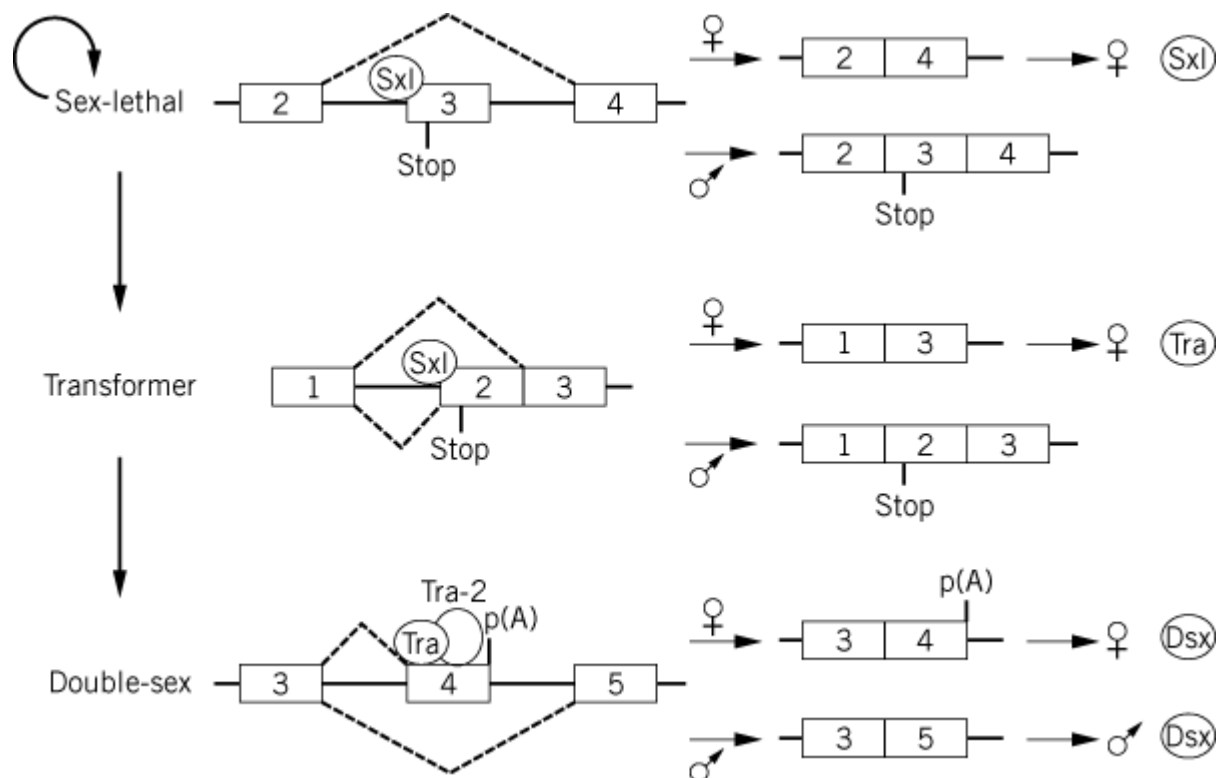
In some cases, regions of the pre-mRNA other than the 5' and 3' splice sites also contribute to alternative splice-site selection. Purine-rich sequences (so-called splicing enhancers), which are typically located in exons and often bind SR proteins, are known to enhance the use of adjacent 5' or 3' splice sites by stabilizing the interactions of spliceosomal components with them (reviewed in Ref. [6](#)). In some instances, intron sequences or sequences in the 3'-untranslated region that regulate [polyadenylation](#) have also been shown to influence alternative splicing. Finally, pre-mRNA secondary structure (eg, stem-loop structures) can affect splice-site selection, for example by blocking the interaction of spliceosomal components with a particular splice site ([7](#), [8](#)).

Splice-site selection during alternative splicing is also regulated by *trans*-acting factors. Many of the currently identified *trans*-acting factors also function in constitutive pre-mRNA splicing. In this case, variations in their concentration, or the concentration of factors that compete with them, are thought to lead to tissue-specific modulation of splice-site usage. Foremost in this category of factors are members of the evolutionarily conserved SR protein family, which are characterized by an amino-terminal **RNA-binding domain** and a C-terminal domain rich in arginine-serine (RS) dipeptides (reviewed in Ref. [9](#)). SR proteins (eg, SF2/ASF and SC35) play essential roles in constitutive nuclear pre-mRNA splicing, particularly at the earliest stages of spliceosome assembly (see [Spliceosome](#)). Moreover, in pre-mRNAs that contain multiple 5' or 3' splice sites (see Fig. [1](#)), high concentrations of SR proteins generally enhance the use of the more proximal (downstream) 5' splice site (so-called switching activity) or more proximal 3' splice site ([9](#), [10](#)). However, apparently due to differences in the affinities of individual SR proteins for different pre-mRNAs, the effect of a particular SR protein on 5' splice site selection can vary from one pre-mRNA to the next ([11](#), [12](#)). SR proteins appear to act by facilitating the interaction of the U1 snRNP with the 5' splice site, an initial step for 5' splice-site recognition, or by promoting the association of U2AF with the 3' splice site (see [Spliceosome](#)). High concentrations of SR proteins can also inhibit exon skipping (Fig. [1](#)), which likewise involves promoting the use of a proximal 5' splice site. Interestingly, the activity of SR proteins such as SF2/ASF or SC35 in 5' splice-site selection (but not in constitutive splicing) can be antagonized by

the hnRNP proteins A/B (13, 14). High concentrations of these proteins generally favor the use of more distal 5' splice sites, and the ratio of 5' splice site usage is determined by the relative amounts of hnRNP A/B and SR protein. Thus, the amount of some alternatively spliced pre-mRNAs can be modulated by varying the cellular concentrations of these proteins (15). SR proteins are also involved in activating weak 5' and 3' splice sites that are located adjacent to purine-rich splicing enhancers (reviewed in Refs. 9 and 16). The binding of specific SR proteins to these exon enhancer sequences promotes the interaction of U2AF, and thus the U2 snRNP, with the 3' splice site or the U1 snRNP with the 5' splice site.

*Trans*-acting alternative splicing factors that are expressed in a tissue-, sex- or developmental-specific manner have also been identified. In mammals, concrete, well-understood examples of cell-specific regulators of alternative splicing are currently limited. The best-characterized factors are those responsible for alternative splicing events in the fruit fly *Drosophila melanogaster*. For example, sex determination in fruit flies involves a cascade of alternative splicing events that are regulated by sex-specific proteins (17) (see Fig. 2). In males, splicing of the Sex-lethal (Sxl), transformer (tra) and doublesex (dsx) pre-mRNAs are nonregulated events that appear to require only the general splicing machinery (ie, they represent the so-called default splicing patterns). In females, alternative splice sites are activated by female-specific factors, either directly or indirectly, through the inactivation of the male-specific site. The first of these factors known to act in this cascade is the female-specific Sxl protein, which regulates its own synthesis by an alternative splicing event. In females, Sxl is thought to block (by a currently unknown mechanism) use of the 3' splice site of the third exon of the Sxl pre-mRNA, which leads to exon 3 skipping (18, 19). Because exon 3 contains a stop codon, functional Sxl protein is produced only from the female-specific Sxl mRNA (Fig. 2). In the next step of this cascade, the Sxl protein inhibits the use of the upstream 3' splice site of exon 2 of the *tra* pre-mRNA, which in turn activates splicing at a weaker downstream 3' splice site. Specifically, Sxl has been shown to bind to the stronger polypyrimidine tract of the upstream 3' splice site and thereby to inhibit the binding of U2AF (20). This results in U2AF binding to the downstream site, for which it has a lower affinity. The resulting exclusion of exon 2, which also contains a stop codon, leads to the production of functional *tra* protein in females. In the last step, the *tra* protein, in conjunction with *tra-2*, directly activates use of the weak 3' splice site of exon 4 of the doublesex pre-mRNA. *Tra* and *Tra-2*, both of which contain RS domains, interact with a purine-rich splicing enhancer present in exon 4 and recruit SR proteins, as well as U2AF, to the upstream 3' splice site (reviewed in Ref. 6). As a result, exon 4 is included in female-specific *dsx* mRNA, and [polyadenylation](#) occurs at the end of this exon. The resulting female *dsx* protein represses male differentiation, whereas the male protein represses female differentiation. Although in this particular case much has been learned about the molecular mechanisms responsible for splice-site selection, in most cases a clear understanding of the complex processes of alternative splicing awaits further investigation.

**Figure 2.** Cascade of alternative splicing events in the sex determination pathway of *Drosophila*. For simplicity, only those exons and introns involved in alternative splicing events are shown. In female flies, the sex lethal protein regulates its own synthesis and that of the *tra* protein by blocking the use of the 3' splice site of the third and second exons of the sex-lethal and *tra* pre-mRNAs, respectively. Both of these exons contain a stop codon, as indicated. The *tra* protein, together with *tra-2*, interacts with a splicing enhancer in the fourth exon of the double-sex pre-mRNA and activates use of the upstream 3' splice site. Exons are depicted as numbered boxes, and introns are depicted as solid lines. Dashed lines above the introns represent splicing events in females, and those below indicate splicing events in males. The sex-lethal (sxl), transformer (tra), and *tra-2* proteins are indicated by circles. Polyadenylation is indicated by p(A). (Adapted from Ref. 21).



## Bibliography

1. R. E. Breitbart, A. Andreadis, and B. Nadal-Ginard (1987) *Annu. Rev. Biochem.* **56**, 467–495.
2. F. H. Nasim, P. A. Spears, H. M. Hoffmann, H.-C. Kuo, and P. J. Grabowski (1990) *Genes Dev.* **4**, 1172–1184.
3. T. Tsukahara, C. Casciato, and D. M. Helfman (1994) *Nucleic Acids Res.* **22**, 2318–2325.
4. B. E. Hoffman and P. J. Grabowski (1992) *Genes Dev.* **6**, 2554–2568.
5. S. M. Berget (1995) *J. Biol. Chem.* **267**, 14902–14908.
6. K. J. Hertel, K. W. Lynch, and T. Maniatis (1997) *Curr. Opin. Cell Biol.* **9**, 350–357.
7. P. A. Estes, N. E. Cooke, and S. A. Liebhaber (1992) *J. Biol. Chem.* **267**, 14902–14908.
8. P. Sirand-Pugnet, P. Durosay, B. Clouet d'Orval, E. Brody, and J. Marie (1995) *J. Mol. Biol.* **251**, 591–602.
9. X.-D. Fu (1995) *RNA* **1**, 663–680.
10. D. S. Horowitz and A. R. Krainer (1994) *Trends Gen.* **10**, 100–106.
11. A. M. Zahler, K. M. Neugebauer, W. S. Lane, and M. B. Roth (1993) *Science* **260**, 219–222.
12. A. M. Zahler and M. B. Roth (1995) *Proc. Natl. Acad. Sci. USA* **92**, 2642–2646.
13. A. Mayeda and A. R. Krainer (1992) *Cell* **68**, 365–375.
14. X. Yang, M. R. Bani, S. J. Lu, S. J. Rowan, Y. Ben-David, and B. Chabot (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6924–6928.
15. J. F. Caceres, S. Stamm, D. M. Helfman, and A. R. Krainer (1994) *Science* **265**, 1706–1709.
16. J. L. Manley and R. Tacke (1996) *Genes Dev.* **10**, 1569–1579.
17. B. S. Baker (1989) *Nature* **340**, 521–524.
18. J. I. Horabin and P. Schedl (1993) *Mol. Cell. Biol.* **13**, 1408–1414.
19. B. Granadino, L. O. F. Penalva, M. R. Green, J. Valcárcel, and L. Sánchez (1998) *Proc. Natl. Acad. Sci. USA* **94**, 7343–7348.
20. J. Valcárcel, R. Singh, P. D. Zamore, and M. R. Green (1993) *Nature* **362**, 171–175.

21. M. J. Moore, C. C. Query, and P. A. Sharp (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 303–358.

### Suggestions for Further Reading

22. Y.-C. Wang, M. Selvakumar, and D. M. Helfman (1997) "Alternative Pre-mRNA Splicing". In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, UK, pp. 242–279.
23. B. Chabot (1996) Directing alternative splicing: cast and scenarios. *Trends Gen.* **12**, 472–478.
24. J. L. Manley and R. Tacke (1996) SR proteins and splicing control. *Genes Dev.* **10**, 1569–1579.
25. D. C. Rio (1993) Splicing of pre-mRNA: mechanism, regulation and role in development. *Curr. Opin. Genet. Dev.* **3**, 574–584.

## Alu Sequences

The [genomes](#) of almost all higher **eukaryotes** contain highly [repetitive DNA](#) sequences that are not clustered together. They are distributed throughout the genome, interspersed with longer stretches of DNA with unique (or moderately repetitive) sequences. In the human genome, the majority of such sequences belong to a single family of [Sines](#) called the **Alu family**. Each sequence is about 300 base pairs long. Although the many copies present are recognizably related, they are not precisely conserved in sequence. Their name derives from the fact that most contain a single site of cleavage for the **restriction enzyme AluI** near their middle. More than 500,000 Alu sequences are present in the human genome, accounting for 3 to 6% of the total DNA. Any particular segment of DNA of 5000 bp or longer has a high probability of containing at least one Alu sequence. Most Alu sequences are flanked by tandem [direct repeats](#) of DNA and move like [transposable elements](#) creating target-site duplications when they insert.

On average Alu DNA sequences contain about 80% identity between members of the family, but certain internal regions are more conserved: an internal 40-bp region and two sets of sequences, one near the 5' end and another one farther down in the transcriptional direction, that are homologous to sequences found in the **promoter** for **RNA polymerase III**.

One end of the Alu DNA segment is defined precisely by comparing several Alu sequences. The other end occurs at, or is adjacent to, a run of A bases of variable length that may or may not be interrupted occasionally by other bases. The internal structure of an Alu sequence is dimeric and may consist of an ancestral duplication of a segment of approximately 150 bp. In some rodents, a major SINE is 130 bp long and has sequence similarities with half of the primate Alu sequences. As in Alu, it is bound on one side by a poly(dA) sequence.

### 1. Origin

The Alu sequence derives from an internally deleted host cell 7SL RNA gene that encodes the RNA component of the [signal-recognition particle](#) (SRP) that functions in **protein biosynthesis** (1, 2). Consequently, an Alu sequence can be considered to be a transposable element or an unusually mobile [pseudogene](#). Alu sequences are transcribed from the 7SL RNA promoter, a polymerase III promoter internal to the transcript, so that it carries the information necessary for its own [transcription](#) wherever it moves. However, it needs to borrow a reverse transcriptase to transpose.

### 2. Evolution

The Alu sequences may be grouped into discrete subfamilies on the basis of their sequences. Distinct families have amplified within the human genome in recent evolutionary history (3). The Human Specific or Predicted Variant subfamily, one of the most recently formed group of Alu sequences, amplified to 500 copies within the human genome sometime after the human/great ape divergence, which is thought to have occurred 4 to 6 million years ago. Comparisons of the sequence and locations of the Alu sequences in different mammals suggest that they have multiplied only recently.

Polymorphism of the Alu family member differs from other types of **polymorphism**, such as Variable Number of Tandem Repeat (VNTR, or [minisatellite DNA](#)) or Restriction Fragment Length Polymorphism ([RFLP](#)), because individuals share Alu insertions based upon identity by descent from a common ancestor as a result of a single event that occurred one time within the human population (4). In contrast the VNTR and RFLP polymorphisms have arisen multiple times within a population. Alu sequences represent a unique source of human genetic variation and a molecular fossil record of genomic evolutionary history. These sequences are natural landmarks for physical gene mapping and for reconstructing the evolutionary history/expansion of tandemly arrayed [gene families](#) (4).

### 3. Possible Functions

The physiological role of Alu elements is unknown, although it has been proposed that they are involved in [DNA replication](#), regulation of transcription, and transport of signal recognition particle RNA to the nucleus. For example, Alu RNA and proteins that bind to Alu elements have been identified in human cells. In particular, it has been demonstrated that some Alu sequences in human gene regions have been altered in sequence so that they are now important in controlling and enhancing transcription (5). The consensus sequence of one of the major Alu families contains a functional [retinoic acid](#) binding element (see [Response Element](#)). The random insertion throughout the primate genome of thousands of Alu repeats containing a retinoic acid response element might have altered the expression of numerous genes, thereby contributing to evolutionary potential (6).

### Bibliography

1. A. M. Weiner (1980) *Cell* **22**, 209–218.
2. E. Ullu, S. Murphy, and P. M. Melli (1982) *Cell* **29**, 195–202.
3. M. A. Batzer et al. (1996) *J. Mol. Evol.* **42**, 22–29.
4. R. J. Britten (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6148–6150.
5. R. J. Britten (1996) *Proc. Natl. Acad. Sci. USA* **93**, 9374–9377.
6. G. Vansant and W. F. Reynolds (1995) *Proc. Natl. Acad. Sci. USA* **92**, 8229–8233.

### Suggestion for Further Reading

7. P. Jagadeeswaran, B.G. Forgeti, and S. M. Weissman (1981) Short interspersed repetitive DNA elements in eucaryotes: Transposable DNA elements generated by reverse transcription of RNA pol III transcripts? *Cell* **26**, 141–142.

### Amber Mutation

An amber [mutation](#) is a [nonsense mutation](#) that changes a sense **codon** (one specifying an amino acid) into the translational [stop codon](#) UAG, causing premature termination of the [polypeptide chain](#)



during [translation](#). The mutation, the codon, and the mutant are all called amber. Amber mutations arise by single base changes in the codons for eight amino acids (and in the UAA stop codon, although this is not a nonsense mutation). Mutations in the anticodons of the [transfer RNAs](#) that read those eight codons, in principle, could give rise to [amber suppressors](#), but suppressors are recovered only if another tRNA exists that reads the codon. In *Escherichia coli*, five amber suppressors that arise by a single base change have been identified. In addition, amber mutations are suppressed by [ochre suppressors](#) because of [wobble pairing](#) in the third position (5') of the [anticodon](#). Amber mutations in *E. coli* and its **bacteriophages** are easily identified by their pattern of suppression by known suppressors. In **bacteria**, amber suppressors have relatively mild effects. Many laboratory strains and even natural isolates of *E. coli* carry amber suppressors. This might be surprising, because amber suppressors are expected to prevent the proper termination of many proteins, but amber codons are used relatively infrequently in *E. coli* and related bacteria.

The name amber was originally given to mutants of bacteriophage T4 that grow on *E. coli* strain K12 (1) but not on *E. coli* strain B (1). It turned out that the K12 strain used has an amber suppressor, whereas the B strain does not. The word amber was inspired by Harris Bernstein who participated in the original experiment (Bernstein means amber in German), although published versions of the story disagree on whether the mutants were named after Harris Bernstein himself or his mother (2, 3). It also could be significant that at nearly the same time that amber mutants were being discovered, Seymour Benzer was also analyzing nonsense mutations in the *rII* genes of phage T4 and calling them “ambivalent” (4).

#### Bibliography

1. R. H. Epstein, A. Bolle, C. Steinberg, E. Kellenberger, E. Boy de la Tour, R. Chevalley, R. Edgar, M. Susman, C. Denhardt, and I. Lielausis (1964) Cold Spring Harbor Symp. Quant. Biol. **28**, 375–392.
2. R. S. Edgar (1966) In *Phage and the Origins of Molecular Biology* (J. Cairns, G. S. Stent, and J. D. Watson, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 166–170.
3. F. W. Stahl (1995) *Genetics* **141**, 439–442.
4. S. Benzer and S. P. Champe (1962) *Proc. Natl. Acad. Sci. USA* **48**, 1114–1121.

#### Amber Suppressor

Amber suppressors are mutant tRNAs that translate the UAG (amber) termination codon as a sense codon. Amber mutations cause protein synthesis to terminate prematurely, resulting in inactive, truncated polypeptides. Amber suppressors allow for protein synthesis beyond the translational block resulting in active protein. Hence the term “suppressor;” these mutant tRNAs “suppress” the phenotypes of amber mutations. These suppressors have been extensively used in prokaryotic genetic studies, and in studies of the translational apparatus and mechanisms. For complete discussions of these and other suppressors, see [Nonsense Suppression](#), [Suppressor tRNA](#), and [Genetic Suppression](#).

#### Ames Test

The *Salmonella*/mammalian **microsome** test for [mutagens](#) was originally developed in the laboratory of Bruce Ames ([1](#)) and has become sufficiently used and well-recognized to be familiarly described by his name. The assay utilizes several specially constructed strains of *Salmonella typhimurium* that normally require histidine for growth and can be reverted to **prototrophy** by a wide range of different mutagens. The assay requires that test chemicals and bacteria be plated onto a minimal agar petri dish, incorporating trace amounts of histidine and [biotin](#), which are required for growth, to allow all the bacteria to grow through a small number of generations. In the absence of mutagen, a small number of colonies will grow on these plates, whereas mutagenic chemicals may increase this number very considerably. Mutations are scored as the number of revertant colonies per dish, usually as a function of applied dose. The test protocol incorporates homogenates of (usually) rat liver directly into the petri dish, thereby permitting mammalian metabolism of many compounds that require activation before they will interact with cellular DNA.

The DNA sequence around the original mutation has been determined in those strains most commonly used for mutagenicity testing (Table [1](#)). The bacteria have been made more sensitive to mutagens by the introduction of several additional characteristics. Many of the strains carry a deletion in the *uvrB* gene and are defective in the ability for [DNA repair](#). The bacterial cell wall has increased permeability to bulky chemicals because of the *rfa* mutation, and certain introduced **plasmids** may increase the sensitivity of the bacteria to mutation by some types of chemicals. The Ames test was originally developed as a screen for chemical carcinogens ([1](#)), but this has only proved appropriate to certain chemical classes (eg, Ref. [2](#)). Nevertheless, because of the enormous number of chemicals tested in this assay, it must still occupy a premier position in testing for mutagenic properties of chemicals.

**Table 1. Genotype and Reversion Characteristics of Some *Salmonella typhimurium* Strains Commonly Used for Mutagenicity Testing**

| Histidine Mutation | Strain Number | Additional Mutations |              |         | Nature of Mutation                                       |
|--------------------|---------------|----------------------|--------------|---------|--|
|                    |               | Permeability         | Repair       | RFactor |  |
| hisC3076           | TA1537        | <i>rfa</i>           | <i>DuvrB</i> | -       | WT sequence unknown.<br>Mutant thought to be +1 near CCC |
| hisD3052           | TA1538        | <i>rfa</i>           | <i>DuvrB</i> | -       | WT:GAC-ACC-GCC-CGG-CAG <sup>1/4</sup>                    |
|                    | TA98          |                      |              | pKM101  | Mutant:GAC-ACC-GCC-GGC-AGG <sup>1/4</sup>                |
| HisD6610           | TA97          | <i>rfa</i>           | <i>DuvrB</i> | PKM101  | WT:GTC-ACC-CCT-GAA-GAG-A*TC-GCC<br>Mutant:GTC-           |

|         |        |            |              |        |   |  |
|---------|--------|------------|--------------|--------|---|--|
| hisG46  | TA1535 | <i>rfa</i> | <i>DuvrB</i> | -      | ACA-CCC-<br>CCC-TGA<br>(opal)<br>WT:.....-<br>CTC-¼ | Some base-<br>pair<br>substitution<br>events |
|         | TA100  |            |              | pKM101 | Mutant:¼-<br>CCC-¼                                  | Extragenic<br>suppressors                    |
| HisG428 | TA102  | <i>rfa</i> |              | PAQ1   | WT:CAG-<br>AGC-AAG-<br>CAA-GAG¼                     | Transitions<br>and<br>transversions          |
|         | TA104  |            |              |        | Mutant:CAG-<br>AGC-AAG-<br>TAA (ochre)              | Extragenic<br>suppressors                    |
|         |        |            |              |        |   | Small<br>deletions (-<br>3, -6)              |

---

## Bibliography

1. B. N. Ames, J. McCann, and E. Yamasaki (1975) *Mutat. Res.* **31**, 347–364.
2. J. Ashby and R. W. Tennant (1994) *Mutagenesis* **9**, 7–16.

## Amidination

Amidination takes place when an [amino group](#) is treated with an imide ester.

The reaction proceeds with reasonable yield when conducted under alkaline conditions (pH around 10). Side reactions, such as [cross-linking](#) take place at lower pH. Both  $\alpha$ - and  $\epsilon$ -amino groups can be amidinated, and the basicity of the modified residue increases. The amidinyl groups incorporated are stable in acidic media. The role of amino groups in protein function can be investigated by amidination. Proteins can be readily **radiolabeled** by amidinating with  $^{13}\text{C}$ - or  $^3\text{H}$ -labeled imide ester. Amidination of the  $\epsilon$ -amino group on [lysine](#) residues is particularly useful in [peptide mapping](#) and in determining the [primary structure](#) of a protein (see [Protein Sequencing](#)). The proteolytic enzyme **trypsin** does not cleave at the modified lysine residues, thereby limiting cleavage to [arginine](#) residues. Moreover, amidinyl groups are removed by aminolysis, and the resulting deprotected **peptides** are cleaved further with trypsin.

### 1. Acetamidination of Proteins

After reduction and [alkylation](#) of the protein (100–700 nmol), it is dissolved in a few mL of 0.2 M triethylamine-HCl [buffer](#), pH 10.3, containing 5.0 M guanidinium chloride (GdmCl) ([1](#)). Ethyl (or methyl) acetamide hydrochloride is dissolved in an equivalent amount of NaOH solution to maintain

a final acetamide concentration of 0.1 to 0.15 M (100-fold molar excess of acetamide over amino groups). The reaction mixture is incubated for 1 hr at 25°C, **dialyzed** against 0.05 M  $\text{NH}_4\text{HCO}_3$  containing 2.5 M GdmCl, and then against 0.05 M  $\text{NH}_4\text{HCO}_3$ . The protein is finally lyophilized.

## 2. Deamidation by Methylaminolysis

Acetamidated protein or peptide (2.7 mg) is dissolved in 1.6 mL of 6 to 9 M [urea](#). Then 0.9 mL of methylamine-formic acid buffer (9.6 M methylamine adjusted with HCOOH to pH 11.5) is added, and the reaction mixture is held for 4 h at 25°C. The final concentration of methylamine is 3.5 M. The reaction mixture is exhaustively dialyzed against deionized water at 4°C or, for peptides, isolated by gel filtration on Sephadex G-10, equilibrated and eluted with 0.1 M  $\text{NH}_4\text{HCO}_3$ .

## Bibliography

1. G. C. DuBois et al. (1981) *Biochem. J.* **199**, 335–340.
2. J. K. Inman et al. (1983) *Methods Enzymol.* **91**, 559–569.

## Amino Acid Analysis

The [amino acid](#) compositions of [proteins](#) are routinely determined by completely hydrolyzing the [peptide bonds](#) of the [polypeptide chain](#), and then determining quantitatively the constituent amino acids that were released.

### 1. Peptide Bond Hydrolysis

The traditional method of hydrolyzing polypeptide chains has been to incubate them anaerobically in 6 M HCl at approximately 110°C for 24–72 h (1). More modern methods use other acids, higher temperatures, and shorter periods of time. Most peptide bonds hydrolyze at similar rates, but those between the large [nonpolar](#) amino acid residues, particularly **Val**, **Leu**, and **Ile**, are hydrolyzed more slowly and require longer hydrolysis times or the addition of organic acids such as trifluoroacetic acid. Hydrolysis is presumably hindered sterically by the bulky side chains.

Any chemical procedure that hydrolyzes the peptide bonds of the backbone will also hydrolyze the chemically similar amide side chains of **Asn** and **Gln** residues, to produce the amino acids [aspartic acid](#) and [glutamic acid](#), respectively. It is feasible to measure the total number of Asn and Gln residues by measuring the amount of ammonia released during the hydrolysis, but otherwise it is not possible to distinguish between Asp and Asn and between Glu and Gln after hydrolysis of the polypeptide chain. In this case, it is common practice to designate such uncertain residues by the three-letter abbreviations Asx and Glx, and by the one-letter abbreviations B and Z, respectively.

**Trp** residues are usually destroyed completely by acid hydrolysis, probably as a result of reaction with chlorine produced by oxidation of the HCl. They can be protected by the addition of thiol or sulfonic acid compounds or of phenol to scavenge the chlorine (2). **Tyr** residues are also susceptible to chlorination, but they are usually lost only partially. The [thiol groups](#) of **Cys** residues are oxidized and the amino acid partially destroyed by acid hydrolysis; this residue is best analyzed after performic acid oxidation of the protein to convert all the Cys residues to cysteic acid.

Some of the problems with acid hydrolysis can be overcome by using other procedures, such as

hydrolysis by alkali or by [proteinases](#). Other amino acids, notably Ser and Thr, are destroyed by alkaline hydrolysis, however, and total proteinase digestion to amino acids is not straightforward. Consequently, acid hydrolysis remains in common use.

## 2. Quantifying the Amino Acids

The identities and quantities of the various amino acids present in a protein hydrolyzates are normally determined by automated amino acid analyzers. The amino acids are separated chromatographically and quantified as they emerge from the column. Traditional methods used ion-exchange chromatography of the free amino acids, followed by detection with [ninhydrin](#) or fluorescent reagents such as [fluorescamine](#). [Proline](#) does not react in the usual manner with such reagents, due to absence of an amino group, so special procedures are required to measure it. More rapid and sensitive methods now predominate, in which the amino acids are reacted with suitable reagents prior to the chromatographic separation, rather than after. The favored method at present is to react the amino acids with phenylisothiocyanate (see [Edman Degradation](#)) and then to separate the colored derivatives by [reverse-phase chromatography](#). With this procedure, a complete quantitative amino acid analysis can be carried out in just a few minutes with only picomole quantities of amino acids (3).

The relative numbers of aromatic residues (**Phe**, **Tyr**, and **Trp**) in intact proteins and peptides can usually be determined from the UV absorbance spectrum under conditions in which the polypeptide chain is fully unfolded so that its spectrum is the sum of its constituent residues (4).

Amino acid analysis does not give directly the number of residues of each amino acid per polypeptide chain. The most accurate result is the molar ratios of the various amino acids. The true molecular weight of the polypeptide chain, in the absence of any non-amino acid moieties, must be known for the amino acid analysis results to be converted to the number of residues of each amino acid per chain. Only with very accurate results, or with very small proteins, are such values usually close to the actual integer values. An alternative procedure is to use progressive chemical modification of one type of amino acid side chain for [counting residues](#).

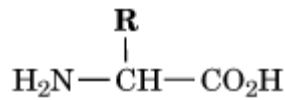
### Bibliography

1. R. L. Hill (1965) *Adv. Protein Chem.* **20**, 37–107.
2. L. T. Ng et al. (1987) *Anal. Biochem.* **167**, 47–52.
3. S. A. Cohen and D. J. Strydom (1988) *Anal. Biochem.* **174**, 1–16.
4. H. Edelhoch (1967) *Biochemistry* **6**, 1948–1954.

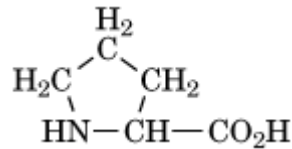
## Amino Acids

Twenty amino acids are the building blocks of [proteins](#). They are linked together in a linear [polypeptide chain](#) by forming [peptide bonds](#) between them, in an order ordained by the nucleotide sequence of the corresponding **gene** for the protein, **translated** from the corresponding [messenger RNA](#).

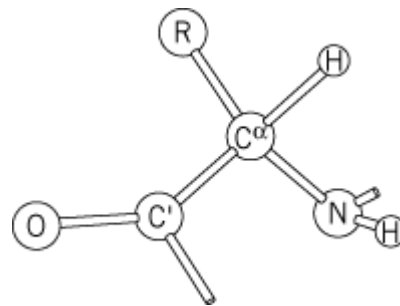
Nineteen of the amino acids have the general structure



and differ only in the chemical structures of the side chain R. The [amino group](#) and the [carboxyl group](#) give this class of compounds its name. At physiological pH values, both groups are ionized, and the zwitterion is the common form of the amino acid. The exceptional amino acid, [proline](#), differs in that its side chain is bonded to the nitrogen atom of the amino group, which is then a secondary amine, and proline is an imino acid:



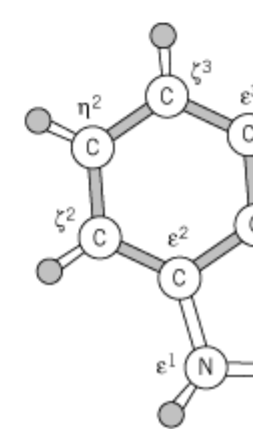
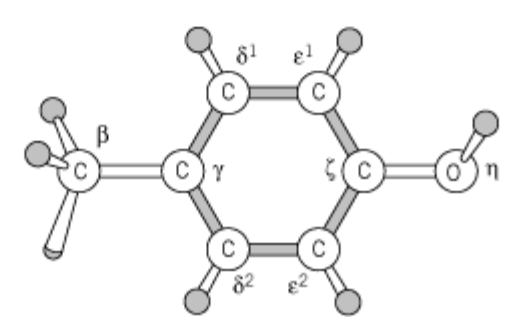
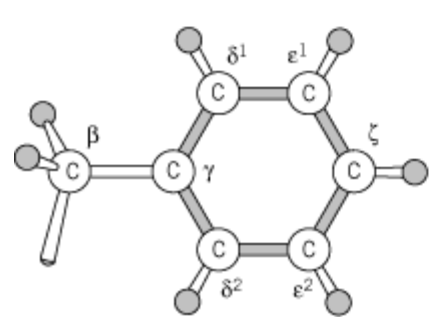
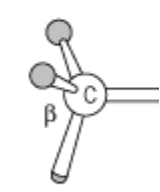
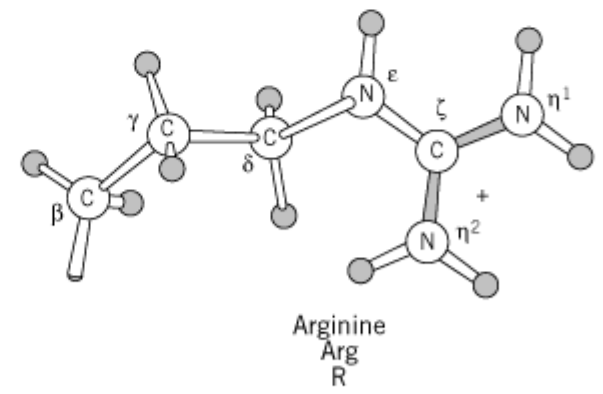
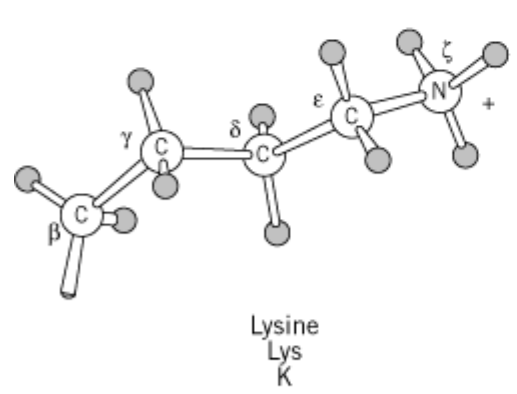
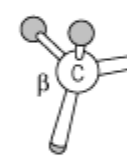
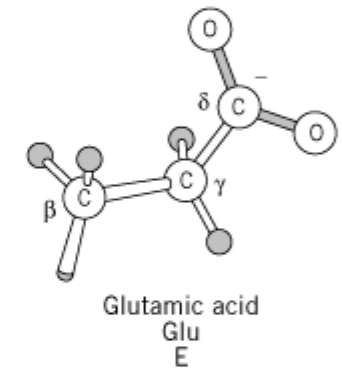
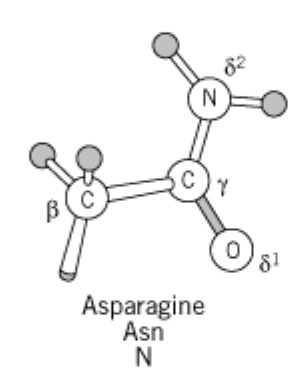
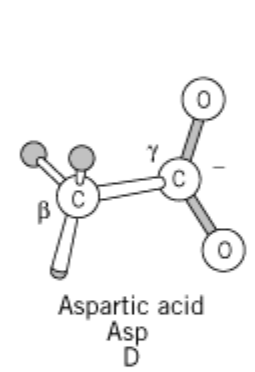
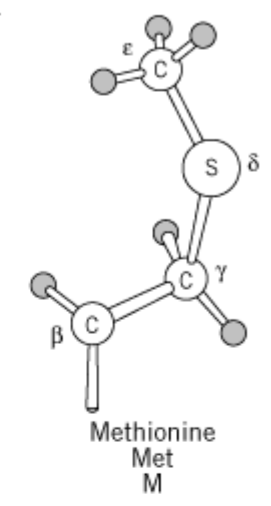
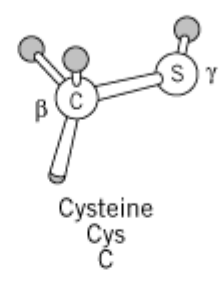
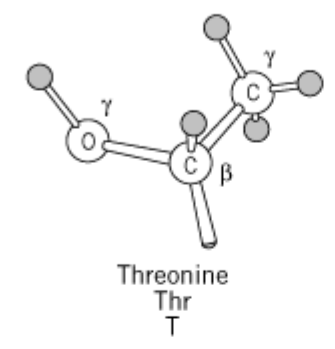
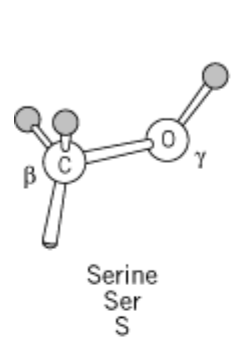
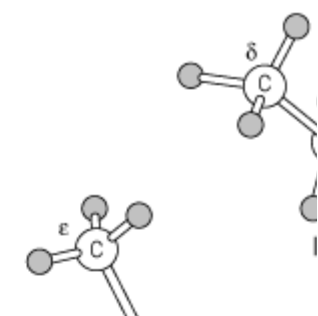
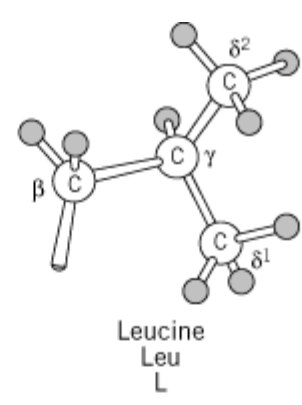
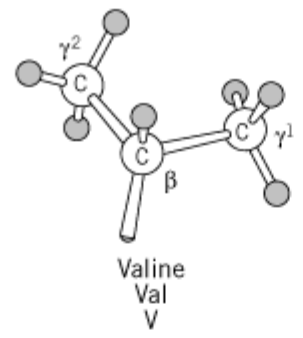
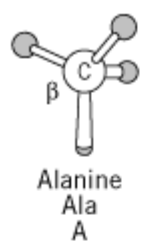
The central  $\alpha$ -carbon atom is asymmetric in 19 of the amino acids and is always the L isomer:



The exception is [glycine](#), in which the side chain is simply another hydrogen atom, so the  $C^\alpha$  atom is no longer asymmetric.

The structures of the side chains of the 20 normal amino acids used in protein biosynthesis are described in [Figure 1](#), and their properties are described in individual entries. The central, asymmetric carbon atom is designated as  $\alpha$ , and the atoms of the side chains are commonly designated  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$  and  $\zeta$  in order away from the  $C^\alpha$  atom. Chemical groups are, however, usually designated by the carbon atom to which they are bonded; hence, the  $N_\zeta$  atom of a [lysine](#) residue is part of the  $\epsilon$ -amino group. A 21st amino acid that is used in protein biosynthesis in only a few instances is [selenocysteine](#). Many variations of these 21 amino acids can be found in proteins as a result of [post-translational modifications](#) after synthesis of the polypeptide chain.

**Figure 1.** The side chains of the 20 amino acids that occur naturally in proteins. Small unlabeled spheres are hydrogen at labeled. Double bonds are black, and partial double bonds are shaded. In the case of proline, the bonds of the polypeptide included and are black. The three- and one-letter abbreviations commonly used appear below the name of the amino acid and threonine have asymmetric centers in their side chains, and only the isomer illustrated is used biologically.



Linking the amino acids into a polypeptide chain involves the condensation of the  $\alpha$ -carboxyl group of one with the  $\alpha$ -amino group of the next, with the elimination of one water molecule. The remaining amino acid within the polypeptide chain is then known as a *residue*. Each type of residue is frequently designated with either three- or one-letter abbreviations, which are given in Figure 1. The three-letter abbreviations are obvious, but the one-letter are preferred with long sequences, as they save space and are less likely to be confused; for example, Gln, Glu, and Gly can easily be confused, but not Q, E, and G. The sequences of amino acids in proteins are usually written with either abbreviation, starting at the left with the *N*-terminal residue, with the free  $\alpha$ -amino group, which is considered the first residue of the polypeptide chain. Amino acid residues in polypeptide chains are properly referred to by changing the ending of their amino acid names, adding or replacing the frequent “-ine” ending with “-yl” (eg, glycyl or alanyl residues), but the amino acid names are commonly used for the residues also. This complication can be minimized by using the one- or three-letter abbreviations. Note that these abbreviations should be used for residues in proteins only, not for the free amino acids.

Depending on the organism and its circumstances, the amino acids are derived from breaking down proteins ingested in feeding and recycling them, and also in some cases by synthesis *de novo*. In the well-known case of humans, the following amino acids cannot be synthesized and must be obtained from the diet: [histidine](#), [isoleucine](#), [leucine](#), [lysine](#), [threonine](#), [tryptophan](#), and [valine](#). [Arginine](#) is synthesized, but not at a rate sufficient during growth. [Methionine](#) is required in large amounts to produce [cysteine](#) if that amino acid is not supplied adequately; similarly, [phenylalanine](#) is required in the absence of [tyrosine](#).

#### Suggestions for Further Reading

E. J. Cohn and J. T. Edsall (1943) *Proteins, Amino Acids and Peptides*, Van Nostrand-Reinhold, Princeton, NJ.

R. E. Marsh and J. Donohue (1967) Crystal structures of amino acids and peptides, *Adv. Protein Chem.* **22**, 235–256.

## Amino Groups

Amino groups are widely distributed in biological substances such as proteins, **polynucleotides**, polysaccharides, and [lipids](#). They play important roles in [electrostatic interactions](#) because of their nucleophilicity and positive charge. The amino groups in the **nucleotides** of nucleic acids are involved in pairing bases. Modification of these groups causes serious errors in nucleic acid replication, [translation](#), and **gene expression**. Amino groups in proteins are important for maintaining their structure and solubility and, sometimes, for manifesting their biological function. In particular, an amino group plays a pivotal role in the enzymes that utilize [pyridoxal phosphate](#) derivatives as **coenzymes**.

There are two kinds of amino groups in proteins. One is the *N*-terminal  $\alpha$ -amino group, that has a  $pK_a$  of 6 to 8, and the other is the  $\epsilon$ -amino group of [lysine](#) residues that have  $pK_a$  values generally between 8 and 10.5. The *N*-terminal amino group is very important for elucidating the [primary structure](#) of a protein because a free *N*-terminal amino group is indispensable for [Edman](#)



**Degradation.** N-terminal amino groups in proteins are often protected by acylation as a [posttranslational modification](#). Attaching [ubiquitin](#) molecules to the amino groups of proteins induces their degradation. Amino groups are reactive nucleophiles and are widely utilized to immobilize proteins on solid supports for [affinity chromatography](#).

## 1. Chemical Modification of Amino Groups in Proteins

Many modification methods for amino groups have been invented based on their excellent nucleophilicities (1). An amino group is a strong nucleophile only in its nonprotonated form, so it is most reactive at high pH. Because of the differences in  $pK_a$  values between  $\alpha$ - and  $\epsilon$ -amino groups, the former may be selectively modified by controlling the pH of the reaction medium. Because many amino groups are exposed and not involved in the protein function, modification of amino groups is suitable for introducing [reporter groups](#), such as chromophores, into proteins and for **radiolabeling** them. The number of amino groups present and subsequently the extent of their modification with some modifying reagents, is determined by trinitrophenylation with 2,4,6-**trinitrobenzene sulfonic acid**. The integral number of amino groups present in a protein can be counted. Representative modification methods for amino groups are shown in Table 1.

**Table 1. Chemical Modification of Amino Groups of Proteins**

| Reaction                 | Reagent                               | pH      |
|--------------------------|---------------------------------------|---------|
| Acylation                | Acetic anhydride                      | 5.5–8   |
|                          | Acetylimidazole                       | >5      |
|                          | <i>N</i> -Acetylsuccinimide           | >4      |
|                          | Citraconic anhydride                  | 8.2     |
|                          | <i>N</i> -Hydroxysuccinimide acetate  | 6.9–8.5 |
|                          | Maleic anhydride                      | 6–10    |
|                          | Succinic anhydride                    | 7–10    |
| Alkylation and arylation | 1-Fluoro-2,4-dinitrobenzene           | 7–11    |
|                          | Iodoacetic acid                       | 7.5–9   |
|                          | 2,4,6-Trinitrobenzene sulfonic acid   | 9.5     |
| Amidination              | Methyl acetamidate                    | 7–10.5  |
| Carbamylation            | Potassium cyanate                     | >7      |
| Guanidination            | 1-Guanyl-3,5-dimethylpyrazole nitrate | 9.5     |
|                          | <i>O</i> -Methylisourea               | 10–11   |
| Reductive alkylation     | Formaldehyde + sodium borohydride     | 8–10    |

## Bibliography

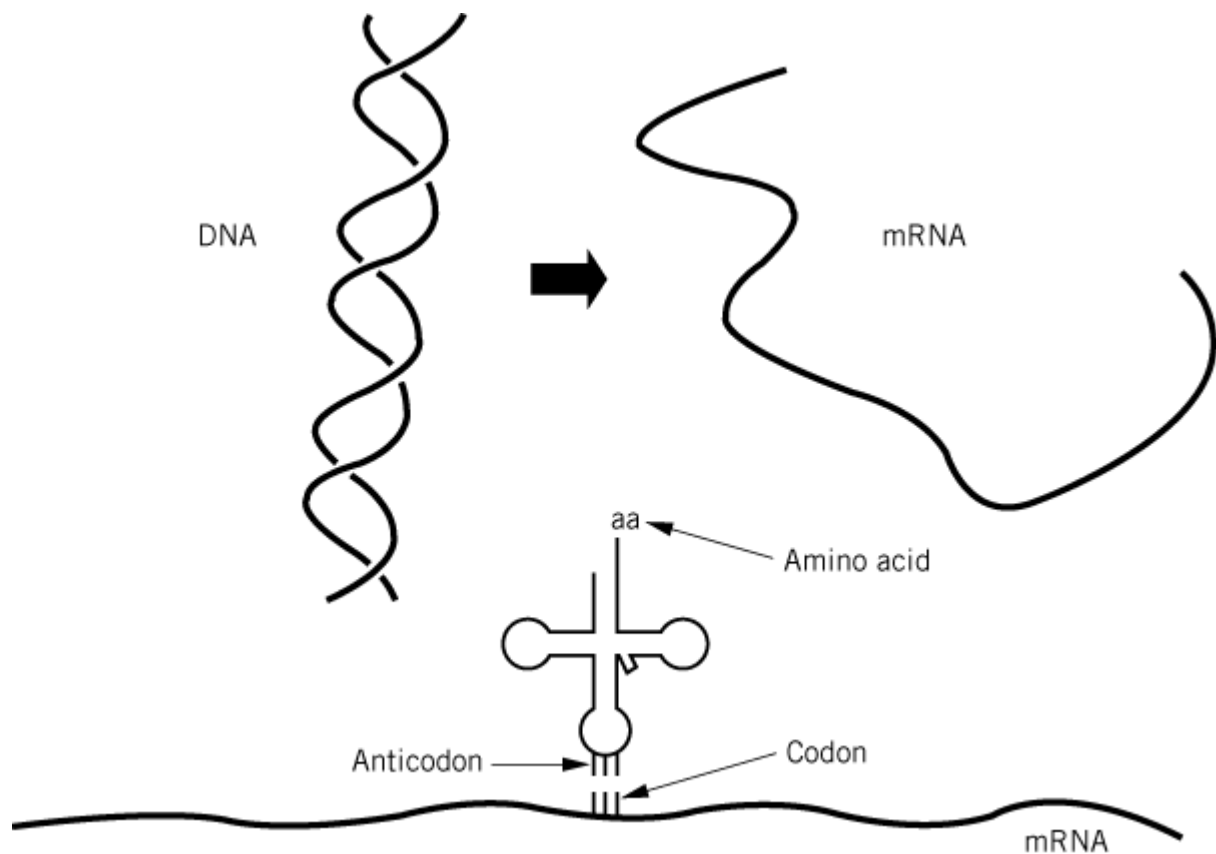
1. T. Imoto and H. Yamada (1989) In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 247–277.

## Aminoacyl tRNA Synthetases

The aminoacyl tRNA synthetases catalyze reactions that establish the rules of the [genetic code](#). For this reason, there is great interest in these [enzymes](#) and their **evolutionary** development, which is thought to be closely connected to the establishment of the code. Research on the synthetases has led to the concept of an operational RNA code for amino acids that is imbedded in the [acceptor stems](#) of [transfer RNA \(tRNA\)](#) (1). The operational RNA code may have played an important role in the assembly of the genetic code and in the overall design of tRNA synthetases.

In the flow of genetic information, [messenger RNA](#) (mRNA) is **transcribed** from DNA, and the mRNA, in turn, is the template for protein synthesis (Fig. 1). The triplet **codons** of mRNA interact with the [anticodons](#) of tRNA through **complementary base pairing**. Amino acids joined to tRNA are incorporated into the growing [polypeptide chain](#). The algorithm of the genetic code relates each amino acid to a specific trinucleotide codon. The triplet associated with a particular amino acid is determined in the aminoacylation reaction, where a given amino acid is linked to a tRNA bearing the anticodon trinucleotide that corresponds to that amino acid. These aminoacylation reactions are catalyzed by aminoacyl tRNA synthetases.

**Figure 1.** Flow of genetic information. Messenger RNA is synthesized from DNA. The mRNA has a string of trinucleotide codons that are translated into a polypeptide whose amino acid sequence is determined by the codons, according to the rules of the genetic code. The amino acid that is inserted into the polypeptide is determined by the codon–anticodon interaction with the tRNA that bears the amino acid corresponding to the particular anticodon. Therefore, the genetic code is determined by the linking of a particular amino acid with a particular anticodon triplet within a tRNA. The joining of amino acids to tRNA is catalyzed by aminoacyl tRNA synthetases. Thus, the genetic code is determined at the biochemical level in the aminoacylation reaction. (This figure was provided by Arturo Morales.)



## 1. Aminoacylation Reaction and the Genetic Code

For most tRNA synthetases, aminoacylation is carried out in a two-step reaction:



In the first reaction, the enzyme, E, uses ATP to activate an amino acid, AA, to yield the firmly bound *aminoacyladenylate* (AA-AMP). In the second step, the activated amino acid is transferred to the 3'-end of the tRNA, where it is connected by an ester linkage to the 2'- or 3'-hydroxyl group (after initial attachment, the amino acid can migrate back and forth between the 2' and the 3' positions). While most synthetases can carry out amino acid activation in the absence of tRNA, there are a few exceptions (such as glutaminyl-, glutamyl-, and argininyl-tRNA synthetases) that require the presence of the cognate tRNA for amino acid activation.

Because each of the 20 natural amino acids used in protein synthesis has a corresponding, or **cognate**, aminoacyl tRNA synthetase, there are 20 of these enzymes in each cellular compartment where proteins are synthesized. Each of them must distinguish its amino acid from all others and, at the same time, recognize the cognate tRNA that bears the anticodon corresponding to that amino acid. In **prokaryotes** and in the cytoplasm of **eukaryotes**, there is typically one tRNA synthetase for each amino acid (eukaryotic [mitochondria](#) have an additional set of synthetases that are essential for mitochondrial protein synthesis). However, the degeneracy of the genetic code means that there are 61 trinucleotides coding for the 20 amino acids. Reading these 61 triplets requires more than just 20 tRNAs. As a consequence, there are multiple tRNA isoacceptors for many of the synthetases. The

synthetases for a particular amino acid must, therefore, recognize and aminoacylate all tRNA isoacceptors for that amino acid. This consideration in itself has important implications.

For example, the codons for serine are sixfold degenerate. In order to read these six codons, the serine tRNA isoacceptors must collectively permute all three anticodon nucleotides. Thus, for a single seryl-tRNA synthetase to aminoacylate all these tRNA<sup>Ser</sup> isoacceptors, the anticodon is not suitable for discrimination. Direct experiments *in vitro* and the [X-ray crystallography](#) structure of the seryl-tRNA synthetase-tRNA<sup>Ser</sup> complex have demonstrated that, in fact, seryl-tRNA synthetase does not contact the anticodon trinucleotide (2). This observation, and others described below, showed that, for at least some amino acids, the relationship between an amino acid and the triplet of the genetic code is not direct.

## 2. Classes of Aminoacyl tRNA Synthetases

The synthetases are heterogeneous in [quaternary structures](#) and subunit sizes, and this heterogeneity obscured more fundamental relationships between these enzymes. For example, in *Escherichia coli*, the quaternary structures of synthetases include a, a<sub>2</sub>, a<sub>4</sub>, and a<sub>2</sub>b<sub>2</sub> (3). Subunit sizes vary from 303 to 951 amino acid residues (4). The 20 aminoacyl tRNA synthetases are now known to be divided into two classes of 10 enzymes each (Table 1) (8, 9). These classes are based on conserved sequence motifs and the structural architecture of the catalytic domains (8, 9). The classification is also based on the fact that the site of initial amino acid attachment on the tRNA differs between the two classes (8). The classes appear to be fixed in evolution, because there is no example of an enzyme switching classes depending on the organism to which it belongs. Thus, the two classes may have developed early in evolution.

**Table 1. Classes of Aminoacyl tRNA Synthetases**

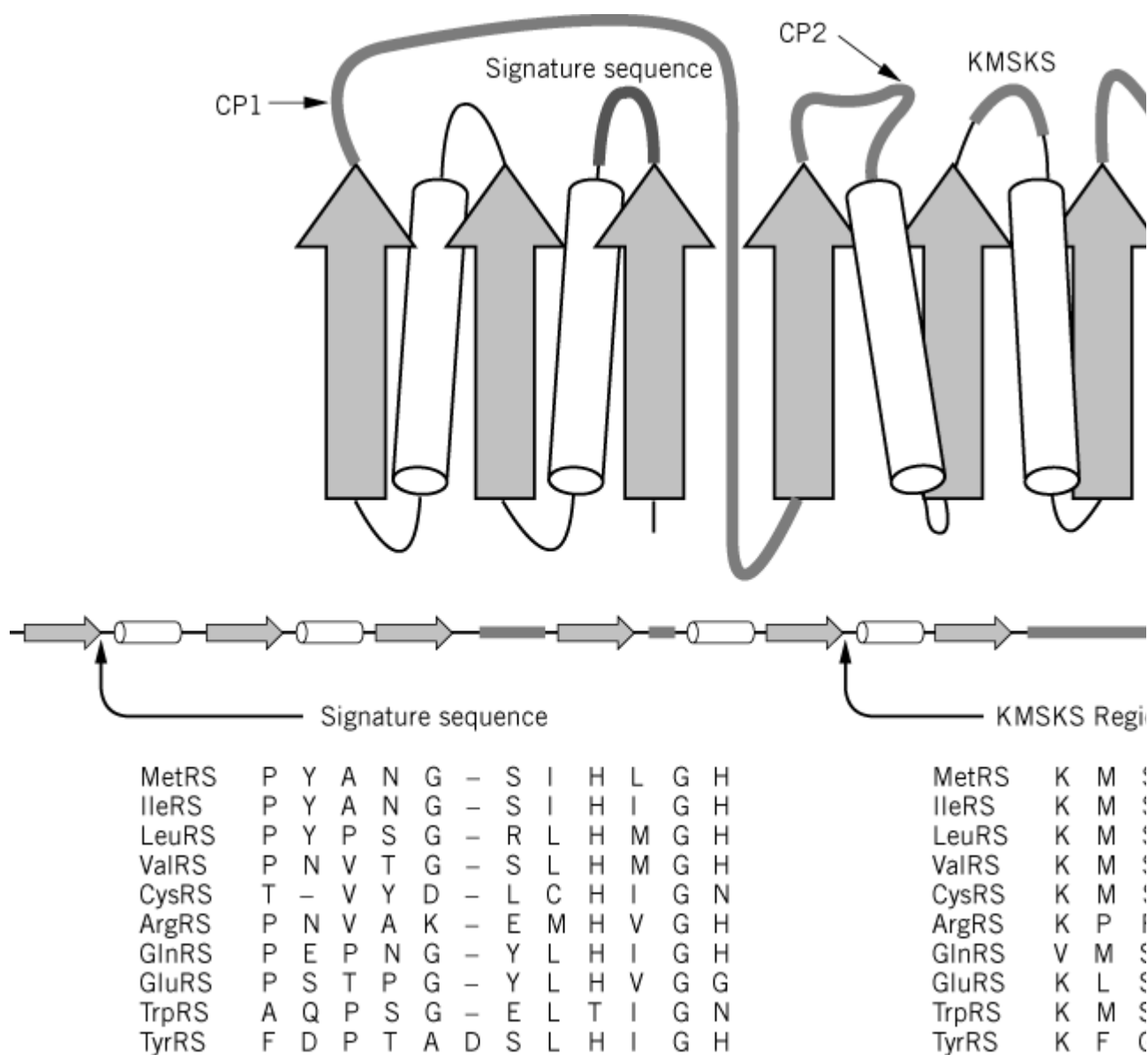
| Class I                | Class II                |
|------------------------|-------------------------|
| Arginine               | Alanine                 |
| Cysteine               | Asparagine <sup>o</sup> |
| Glutamic               | Aspartate               |
| Glutamine <sup>l</sup> | Glycine                 |
| Isoleucine             | Histidine               |
| Leucine                | Lysine                  |
| Methionine             | Phenylalanine           |
| Tryptophan             | Proline                 |
| Tyrosine               | Serine                  |
| Valine                 | Threonine               |

Gram-positive bacteria, plant chloroplasts, and animal mitochondria have been shown to have less than 20 tRNA synthetases. Instead, glutamyl-tRNA synthetase catalyzes attachment of glutamic acid to both tRNA<sup>Glu</sup> and tRNA<sup>Gln</sup> and, similarly, aspartyl-tRNA synthetases catalyzes attachment of aspartate to both tRNA<sup>Asp</sup> and tRNA<sup>Asn</sup>. An amidotransferase then catalyzes the amidation of Glu-tRNA<sup>Gln</sup> to give Gln-tRNA<sup>Gln</sup>, and, likewise, amidation of Asp-tRNA<sup>Asn</sup> gives Asn-tRNA<sup>Asn</sup> (5-7).

### 2.1. Class I

These enzymes are usually monomers and are characterized by an architecture that is similar to that seen in dehydrogenases and other **nucleotide-binding** proteins. This structural motif is a **Rossmann** nucleotide-binding fold, which consists of alternating **b-strands** and **a-helices** (Fig. 2) (10-12). In the case of class I tRNA synthetases, the fold is divided into two  $b_3a_2$  halves to give an overall  $b_6a_4$  structure. In this structure, the b-strands are arranged in parallel. A polypeptide of variable length, designated as connective polypeptide 1 (CP1), links together the two halves of the active site (13). In some class I enzymes, this insertion plays a role in **translational editing**. It also contains some of the residues for binding the synthetase to the tRNA acceptor helix (14).

**Figure 2.** Design of a class I tRNA synthetase. The nucleotide-binding fold of class I tRNA synthetases consists of alternating  $\alpha$ -helices (cylinders) and  $\beta$ -strands (arrows) that form a  $b_6a_4$  structure. A two-dimensional spatial arrangement of these elements is shown at the top, and a linear schematic is shown below. A second domain of variable size occurs after the nucleotide-binding fold. The fold is split into two  $b_3a_2$  halves by *polypeptide 1* (CP1). A second, smaller insertion (CP2) splits the second half of the fold. Two sequence elements were used to identify class I enzymes. These are known as the *12-residue signature sequence*, which ends in the HIGH tetrapeptides (10, 11) and as the *KMSKS* region. Locations in the schematic structure are shown near the label signature sequence and KMSKS. By way of example, an alignment of amino acid sequences of the 10 class I *E. coli* enzymes is shown beneath the schematic figures. Similar alignments can be made for other class I enzymes throughout evolution. (This figure was provided by Arturo Morales.)

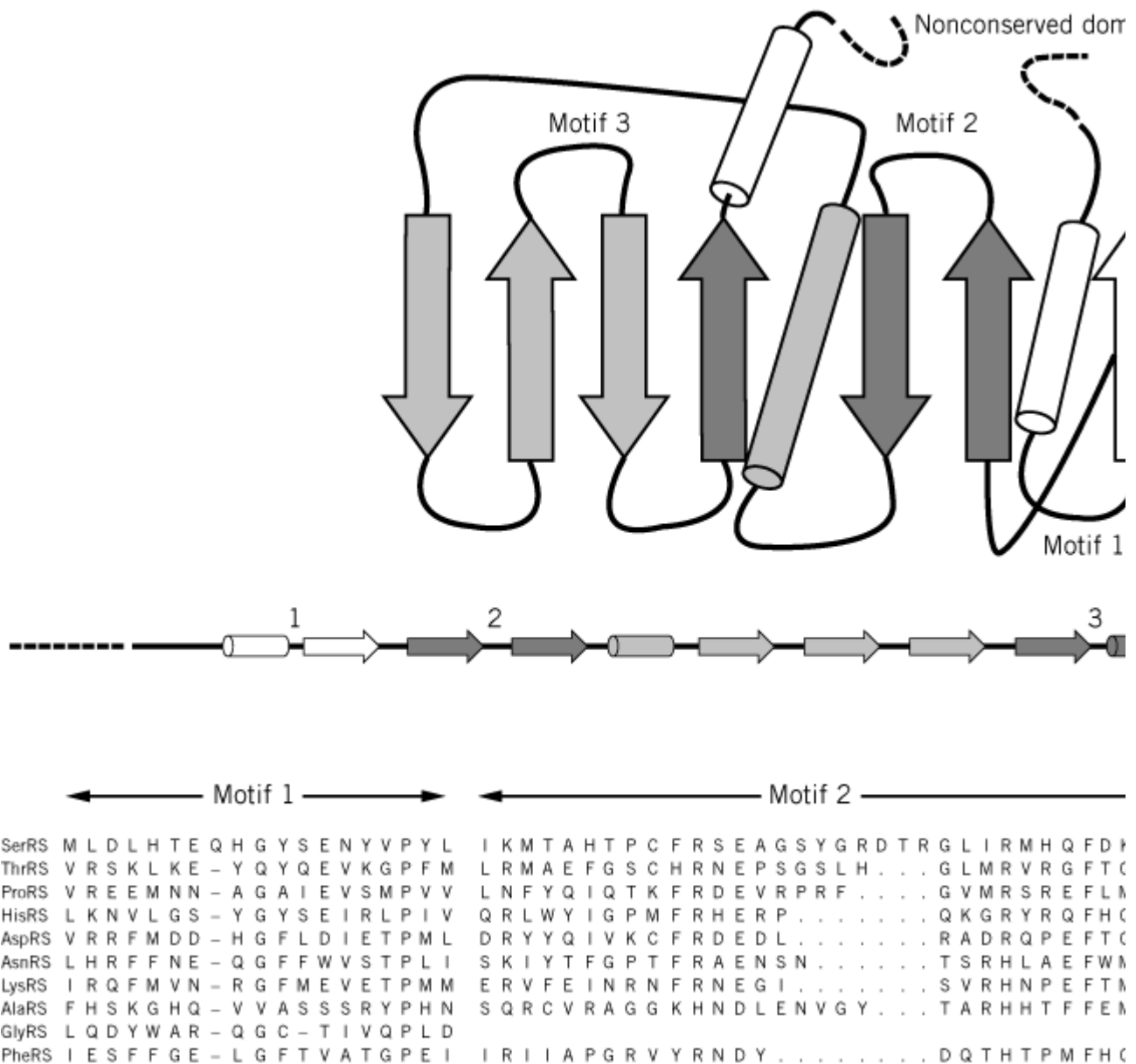


The nucleotide-binding fold contains the site for adenylate synthesis. This catalytic **domain** may be identified by two characteristic sequence motifs, without any knowledge of three-dimensional structure. One motif is the 11-amino acid element known as the signature sequence, which ends in the sequence–His–Ile–Gly–His, or HIGH in one-letter code ([10](#), [11](#)). This element is located in the first half of the nucleotide-binding fold at the end of the first b-strand and the beginning of the first a-helix. It was designated as a signature sequence because it served as a clear signature for a subgroup of related synthetases, before many crystal structures were determined. The second element is the KMSKS motif, located in the second half of the nucleotide-binding fold ([15](#)). These elements are critical parts of the [active site](#).

## 2.2. Class II

The class II enzymes are mostly  $\alpha_2$  dimers. The active sites of class II enzymes have a completely different architecture that harbors three characteristic sequence motifs. The structure consists of a seven-stranded antiparallel b-sheet with three a-helices ([9](#), [16-18](#)) (Fig. [3](#)) ([8](#), [9](#), [19](#)). The three characteristic sequence motifs are known as motifs 1, 2, and 3. The sequences of these motifs are highly degenerate ([8](#), [9](#)). They consist of a helix–loop–strand, strand–loop–strand, and strand–helix, respectively. All three of these motifs form part of the active site.

**Figure 3.** Design of a class II tRNA synthetase. The seven-stranded b-sheet with three a-helices is shown. The variable-s line that may occur on either the *N*- or the *C*-terminal side of the class-defining domain. Three characteristic sequence motifs are distinguished in this illustration by their different shadings. These motifs are highly degenerate in sequence and consist of a helix–loop–strand (motif 1), a strand–loop–strand (motif 2), and a strand–helix (motif 3). The locations of these motifs in the class-defining domain are shown. An example of the sequence motifs for *E. coli* tRNA synthetases is also shown ([8](#), [19](#)). Note the high degeneracy of these sequence motifs, especially the sequence elements in class I enzymes (Fig. [2](#)). (This figure was provided by Dr. Arturo Morales.)



### 2.3. Amino Acid Attachment

The site of initial amino acid attachment for class I enzymes is the 2'-hydroxyl, whereas the 3'-hydroxyl is used by class II enzymes (20). This distinction is now understood to result from a difference in the ways that the two enzymes approach the end of the tRNA. In particular, class I enzymes approach the end of the tRNA acceptor helix from the minor groove side, while class II enzymes approach from the major groove side (21).

### 3. Overall Structural Design and the Synthetase-tRNA Complex

The class-defining active-site domain is only a part of the tRNA synthetase structure. Inserted into the active-site domain are sequence motifs that enable the tRNA acceptor helix to bind with its 3'-end near the aminoacyl adenylate. These insertions are typically idiosyncratic to the synthetase. Two examples are the CP1 insertions of class I enzymes and the variable loops of motif 2 of class II enzymes, both of which have a role in docking the acceptor helices to the enzymes. However, in addition to insertions into the active-site domain, all synthetases have a second major domain. This domain, also typically idiosyncratic to the enzyme, provides for contacts with parts of the tRNA that

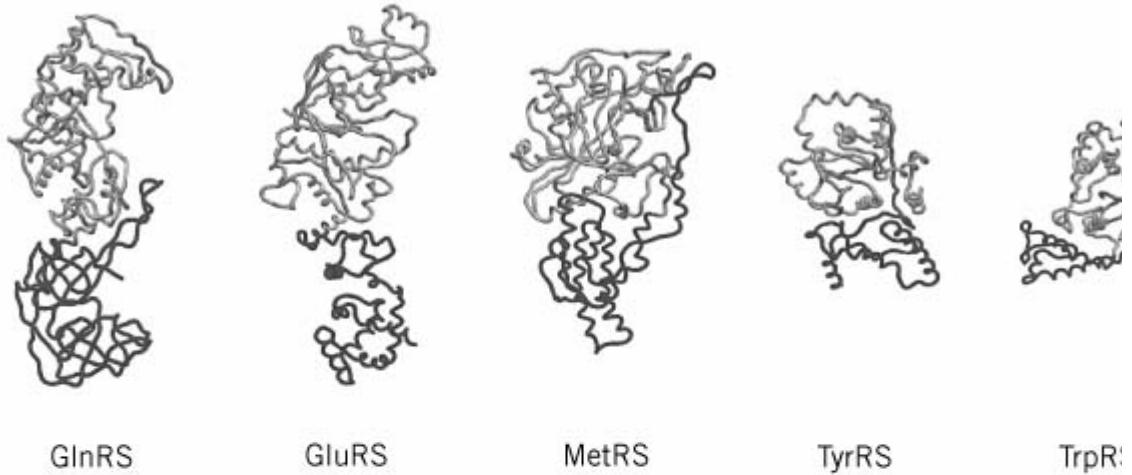
are distal to the amino acid attachment site. For many (but not all) synthetases, this includes contacts with the anticodon. For the class I methionyl- and glutaminyl-tRNA synthetases, this second, anticodon-binding, domain is largely  $\alpha$ -helical and largely  $\beta$ -structure, respectively (22, 23). This difference demonstrates that, even for enzymes in the same class, their second domains are completely unrelated. In the case of the class II seryl-tRNA synthetase, an unusual [coiled-coil](#) protrudes from the N-terminus of the enzyme. This structure, which is not found in many other class II enzymes, provides for contacts with the variable loop of tRNA<sup>Ser</sup> (2, 24).

Thus, to a rough approximation, the synthetases comprise two major domains (9, 14, 16, 23, 25-31). The tRNA molecule also comprises two major domains, which consist of the four arms of the cloverleaf secondary structure (Fig. 4). One domain is the acceptor-TyC minihelix, where the amino acid acceptor end and the TyC stem stack together to make a helix of 12 base pairs. The second domain is formed by stacking of the dihydrouridine stem with the anticodon stem. The result is an L-shaped three-dimensional structure where the amino acid acceptor end and the anticodon-containing template-reading-head are segregated into separate structural units (see [Transfer RNA](#)).

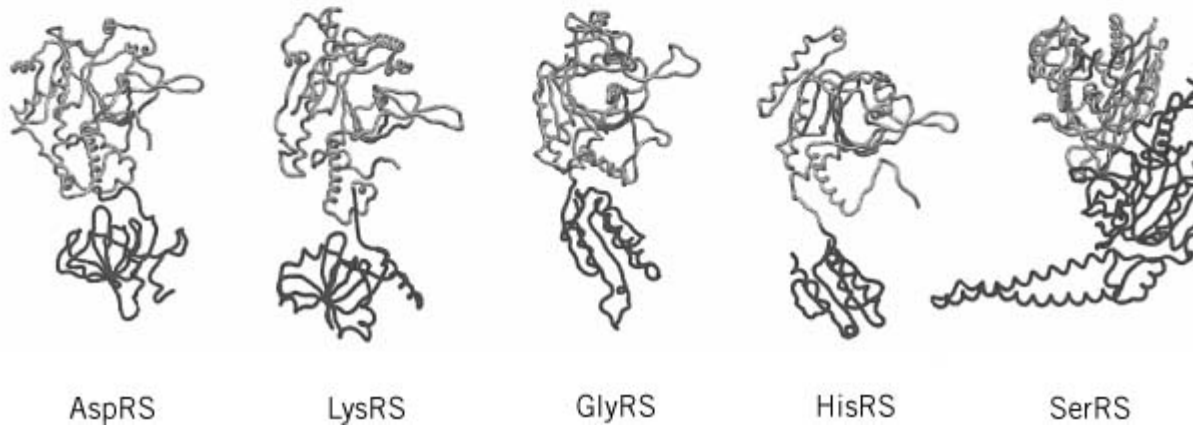
**Figure 4.** Examples of crystal structures of class I and class II tRNA synthetases. Regardless of the class to which an enzyme is assigned, its structure can be approximated as comprising two major domains. One is the class-defining active site domain which is shared by all members of the same class. This domain (gray) is thought to be the historical tRNA synthetase. The second domain (dark gray) is typically idiosyncratic to the synthetase and is not shared by all members of the same class. This domain was probably added later to the synthetase structure. In these illustrations, the domains have been defined by the obvious visual divisions in the structures and, for that reason, the catalytic domain may extend somewhat beyond the region which contains the class-defining motifs. Some of these enzymes are homodimers; in all cases, only a single subunit is shown. In the case of the tyrosyl tRNA synthetase, TyrRS, the structure of the second domain is not complete and therefore is truncated in this structural representation. Not shown is PheRS (25), which has an  $\alpha_2\beta_2$  quaternary structure. These structures were determined for GlnRS (14), GluRS (26), MetRS (23), TyrRS (27), TrpRS (28), AspRS (16), LysRS (29), GlyRS (30), HisRS (31), and SerRS (9). (This figure was provided by Arturo Morales.)



Class I

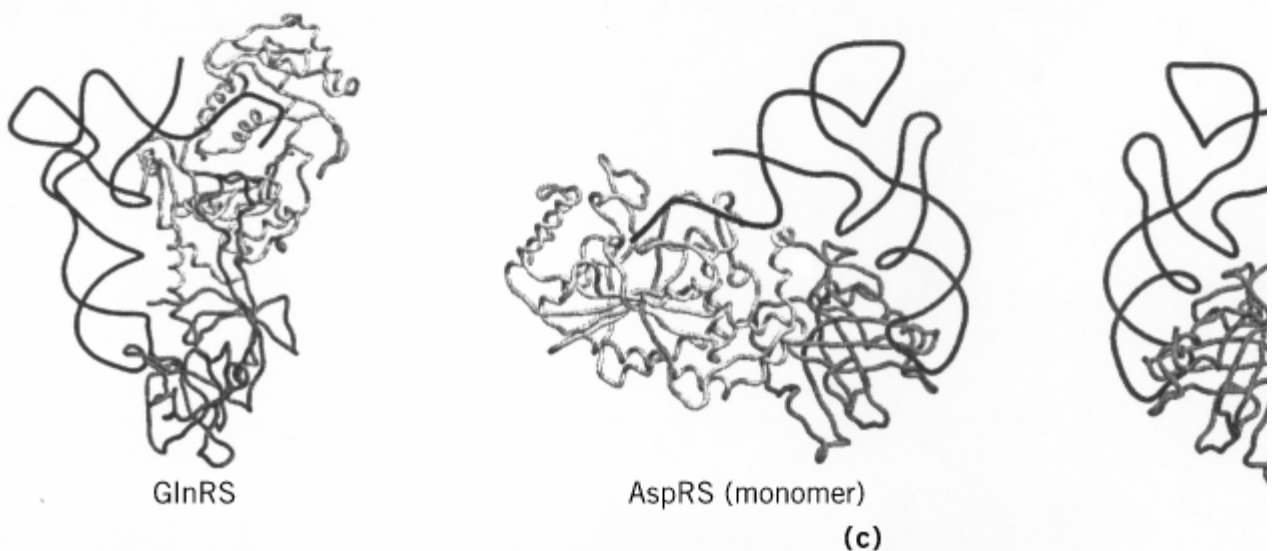
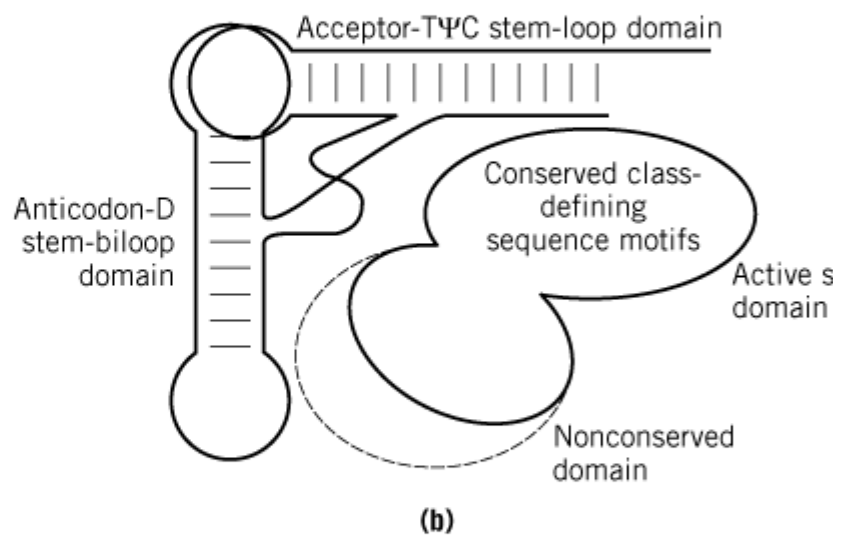
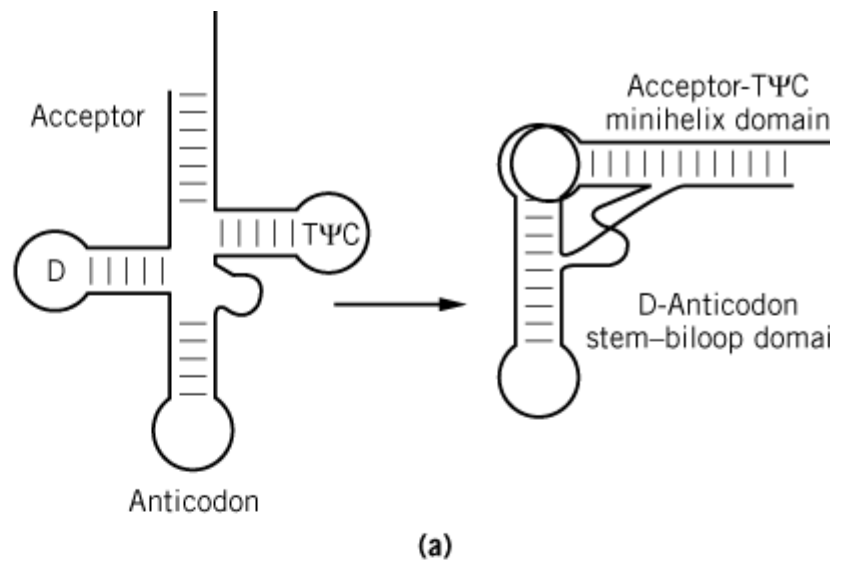


Class II



The two major domains of the synthetase make contact with the two domains of the L-shaped tRNA molecule (Fig. 5). The catalytic domain, with insertions containing RNA-binding determinants, interacts with the acceptor-TyC minihelix. The second major domain interacts with the second domain of the tRNA, where interactions may extend as far as the anticodon trinucleotide or may involve contacts with special structures, such as the large variable loop of tRNA<sup>Ser</sup>. The details of the interactions between tRNA and synthetase are described in [RNA-binding proteins](#).

**Figure 5.** Two major domains of a synthetase interacting with two domains of a tRNA. (a) Schematic representation of tRNA domains that segregate the amino acid attachment site into a 12-bp minihelix stem-loop and the anticodon triplet into a 3-bp TyC loops are indicated as common landmarks found in most all tRNA. (Adapted from Ref. 32.) (b) Domains of a synthetase interacting with the two domains of a tRNA. The second domain of the synthetase is shown with a dotted line to indicate that it varies in size and often extends further. (c) Examples of synthetase-tRNA complexes, with the active-site class-defining domain interacting with the acceptor stem, an idiosyncratic domain of the synthetase interacting with the second domain of the tRNA (5, 16). Separate shades indicate separate domains. In the case of aspartyl tRNA synthetase, the protein is a dimer. Binding of a single tRNA to the monomer is shown on the left. The dimeric complex with two bound tRNA is also shown at the right. (This figure was provided by



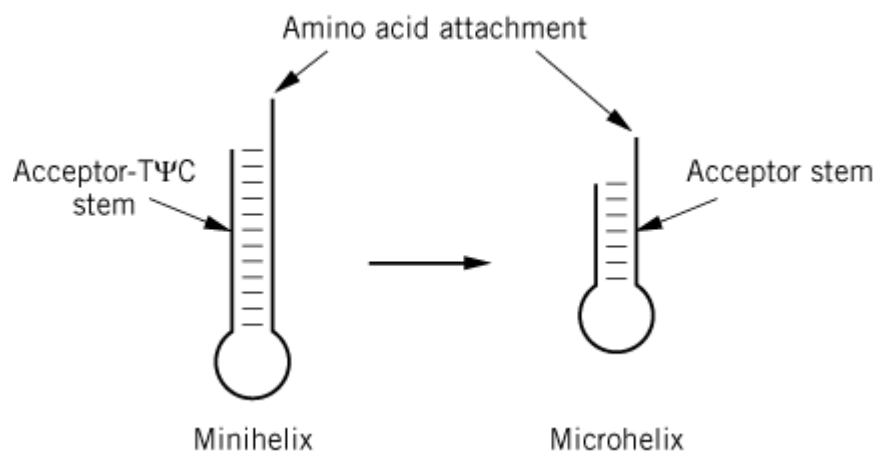
#### 4. Operational RNA Code for Amino Acids

In addition to seryl-tRNA synthetase, alanyl-tRNA synthetase is an example of a synthetase that makes no contact with the anticodon (34). Instead, an acceptor helix G3:U70 **wobble** base pair is a

major determinant of the identity of an alanine tRNA (35, 36). Alteration of this base pair to G:C, A:U, I:U, or U:G abolishes aminoacylation with alanine (37-39). Transfer of this base pair into other, nonalanine tRNA confers alanine acceptance on them. Thus, the G3:U70 pair marks a tRNA for charging with alanine.

Because the G3:U70 base pair is located in the acceptor helix, further experiments tested whether the 12-bp acceptor-TyC minihelix by itself would be a substrate for aminoacylation (Fig. 6) (40-42). Not only the minihelix, but also a 7-base-pair(bp) microhelix consisting of just the acceptor stem, is efficiently charged with alanine, provided that it contains the G3:U70 base pair (40, 43). Transfer of this base pair into other microhelices confers alanine acceptance on them. Thus, the charging behavior of the mini- and microhelices reproduces that seen with the full tRNA.

**Figure 6.** Minihelix and microhelix substrates for aminoacylation. The minihelix is derived from the 12-bp acceptor-TyC domain of the tRNA, while the microhelix is a hairpin helix whose base pairs consist of the 7-bp acceptor stem portion (40). These substrates for aminoacylation are devoid of the anticodon trinucleotides of the genetic code. About 11 examples of aminoacylation of minihelix or microhelix structures have been demonstrated. These substrates have been used to delineate the operational RNA code for amino acids (41). (Adapted from Ref. 42.)



The charging of specific RNA helices has now been demonstrated with at least 11 different tRNA synthetases, even for cases where the anticodon is known to play a significant role in the recognition of the related tRNA (41, 44-47). In the case of histidine, for example, an extra nucleotide at the 5' end of the acceptor helix is characteristic of and unique to histidine tRNA throughout evolution. RNA microhelices that contain the extra base are charged with histidine (48). Thus, the 5'-appended nucleotide marks a molecule for charging with histidine. The smallest substrates seen to be charged with specific amino acids are stem-loop hairpins with as few as four base pairs stabilized by an RNA **tetraloop** motif (Fig. 7) (49).

**Figure 7.** RNA tetraloop substrates for aminoacylation. The amino acid that can be charged onto the designated structures is indicated (49). These short helices are stabilized by an RNA tetraloop motif that confers unusual stability to short RNA helices. (Adapted from Ref. 49.)

|                   |                   |                   |                   |
|-------------------|-------------------|-------------------|-------------------|
| A                 | A                 | A                 | A                 |
| C                 | C                 | C                 | C                 |
| C                 | C                 | C                 | C                 |
| A                 | U                 | <sup>-1</sup> G-C | A                 |
| <sup>1</sup> G-C  | <sup>1</sup> G-C  | <sup>1</sup> G-C  | <sup>1</sup> C-A  |
| G-C               | C-G               | U-A               | G-C               |
| G-U <sup>70</sup> | G-C <sup>70</sup> | G-C <sup>70</sup> | C-G <sup>70</sup> |
| C-G               | C-G               | C-G               | G-C               |
| U G               | U G               | U G               | U G               |
| U C               | U C               | U C               | U C               |
| Ala               | Gly               | His               | Met               |

The charging of microhelix substrates is sometimes considerably less efficient than that for the corresponding tRNA. In each case, however, charging is sequence-specific and depends on two to four nucleotides near the amino acid attachment site. In a fine-structure mapping of the efficiently charged alanine microhelix, a constellation of atoms was identified as critical for the aminoacylation signal. Prominent among these atoms was the exocyclic 2-amino group of G of the G3:U70 base pair. Also important were specific 2'-hydroxyl groups that fell within a 5-Å (Å = 10<sup>-10</sup> m) radius of the critical 2-amino group (39).

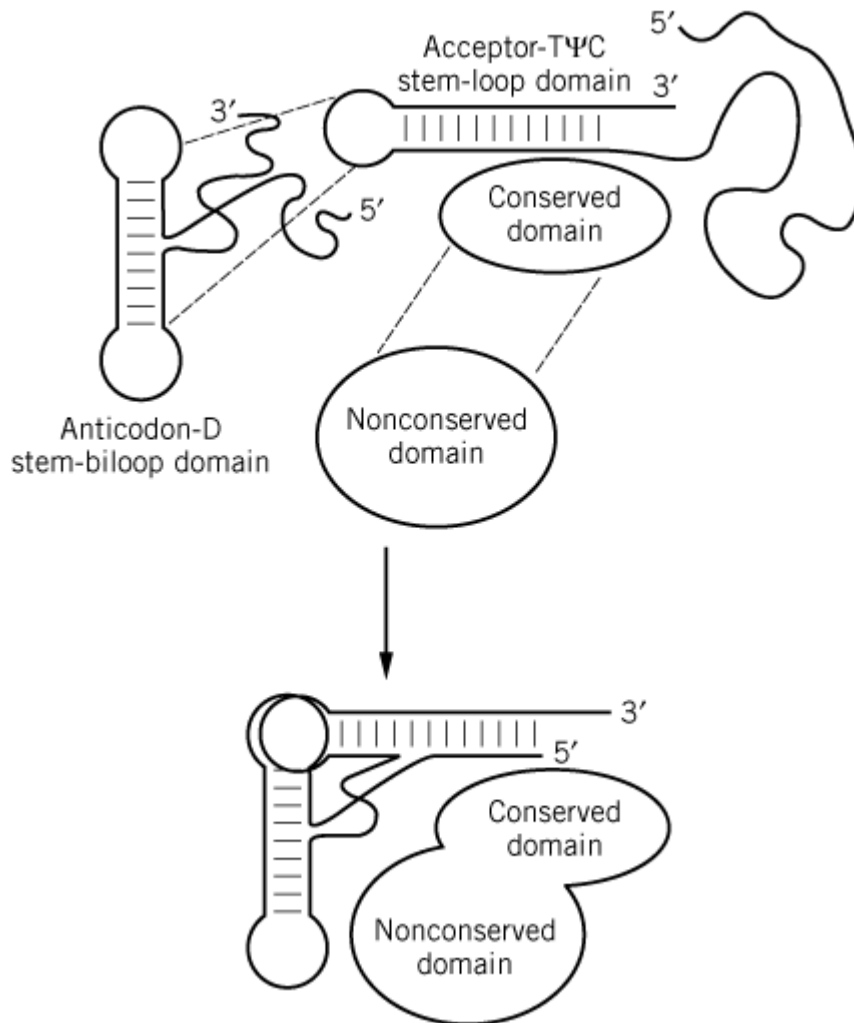
These anticodon-independent aminoacylations of oligonucleotide substrates demonstrate that specific RNA sequences and structures per se, rather than the trinucleotides of the genetic code, correspond to specific amino acids. The relationship between these RNA sequences and structures and particular amino acids constitutes an operational RNA code for amino acids that is distinct from, but related to, the genetic code. This operational RNA code may have predated the genetic code (1).

##### 5. Assembly of the Synthetase-tRNA Complex and Relationship of the Operational RNA Code to the Genetic Code

Several considerations led to the proposal that the minihelix and anticodon-containing second domain of tRNA had distinct origins (1, 50-52). The minihelix, with its amino acid attachment site, is viewed as the historical, or earliest, part of the tRNA. Similarly, the class-defining catalytic domain of tRNA synthetases is thought to be the historical enzyme, with the idiosyncratic second domain added later. Experiments with alanyl-tRNA synthetase have demonstrated that a relatively small piece of the enzyme (containing the active site) can by itself charge an RNA microhelix. This result directly demonstrated a domain-domain interaction between discrete units of the tRNA and the synthetase that may somewhat resemble an evolutionarily earlier system (33, 54).

Thus, the early synthetase may have consisted solely of a domain for adenylate synthesis. Insertions into this domain allowed the docking of RNA substrates near the activated amino acid so that aminoacylation could occur (Fig. 8). Addition of the second domain of the tRNA, with its anticodon-containing template reading head, and of the second domain of the synthetases was a second, later event. This event led to the joining of the operational RNA code to the genetic code. This scheme also suggests that the relationship between a particular amino acid and the triplet of the code is random, and simply depends on which anticodon-containing domain happened to be fused to the minihelix domain of the tRNA.

**Figure 8.** Assembly of a tRNA synthetase in evolution. A primordial tRNA synthetase is envisioned as interacting with a minihelix-like structure. As the tRNA structure developed, a template reading head (anticodon domain) was added, along with the second domain for the synthetase. (Adapted from Ref. 1.)



## 6. Conclusions

The tRNA synthetases may be among the earliest proteins, arising in evolution contemporaneously with the development of the genetic code. The first synthetases may have been **ribozymes** that catalyzed aminoacylation reactions with a specificity that depended on the sequences and structures of the RNA substrates (54, 55). The proteins that replaced these ribozymes were probably small. As a result, they could not extend much beyond the amino acid attachment site and, for that reason, gave rise to a system of interactions that based specificity of aminoacylation on interactions near the end of the acceptor helix. How these charged RNA substrates were used to synthesize specific proteins is a question of great interest.

## Bibliography

1. P. Schimmel, R. Giegé, D. Moras, and S. Yokoyama (1993) *Proc. Natl. Acad. Sci. USA* **90**, 8763–8768.
2. V. Biou, A. Yaremchuk, M. Tukalo, and S. Cusack (1994) *Science* **263**, 1404–1410.
3. P. Schimmel (1987) *Ann. Rev. Biochem.* **56**, 125–158.
4. S. A. Martinis and P. Schimmel (1996) in *Escherichia coli and Salmonella* (F. C. Neidhardt Jr.,

- ed.), ASM Press, Washington, DC, pp. 887–901.
5. A. Schon, C. G. Kannangara, S. Gough, and D. Söll (1988) *Nature* **331**, 187–190.
  6. A. W. Curnow, M. Ibba, and D. Söll (1996) *Nature* **382**, 589–590.
  7. Y. Gagnon, L. Lacoste, N. Champagne, and J. Lapointe (1996) *J. Biol. Chem.* **271**, 14856–14863.
  8. G. Eriani, M. Delarue, O. Poch, J. Gangloff, and D. Moras (1990) *Nature* **347**, 203–206.
  9. S. Cusack, C. Berthet-Colominas, M. Hartlein, N. Nassar, and R. Leberman (1990) *Nature* **347**, 249–255.
  10. T. A. Webster, H. Tsai, M. Kula, G. A. Mackie, and P. Schimmel (1984) *Science* **226**, 1315–1317.
  11. S. W. Ludmerer and P. Schimmel (1987) *J. Biol. Chem.* **262**, 10801–10806.
  12. C. Hountondji, F. Lederer, P. Dessen, and S. Blanquet (1986) *Biochemistry* **25**, 16–21.
  13. R. M. Starzyk, T. A. Webster, and P. Schimmel (1987) *Science* **237**, 1614–1618.
  14. M. A. Rould, J. J. Perona, D. Söll, and T. A. Steitz (1989) *Science* **246**, 1135–1142.
  15. C. Hountondji, P. Dessen, and S. Blanquet (1986) *Biochimie* **68**, 1071.
  16. M. Ruff, S. Krishnaswamy, M. Boeglin, A. Poterszman, A. Mitschler, A. Podjarny, B. Rees, J. C. Thierry, and D. Moras (1991) *Science* **252**, 1682–1689.
  17. S. Cusack, M. Hartlein, and R. Leberman (1991) *Nucl. Acid Res.* **19**, 3489–3498.
  18. D. Moras (1992) *Trends Biochem. Sci.* **17**, 159–164.
  19. S. Cusack (1993) *Biochimie* **75**, 1077–1081.
  20. P. R. Schimmel and D. Söll (1979) *Annu. Rev. Biochem.* **48**, 601–648.
  21. J. Cavarelli and D. Moras (1993) *FASEB J.* **7**, 79–86.
  22. M. A. Rould, J. J. Perona, D. Söll, and T. A. Steitz (1991) *Nature* **352**, 213–218.
  23. S. Brunie, C. Zelwer, and J. Risler (1990) *J. Mol. Biol.* **216**, 411–424.
  24. S. Cusack, A. Yaremchuk, and M. Tukalo (1996) *EMBO J.* **15**, 2834–2842.
  25. L. Mosyak, L. Reshetnikova, Y. Goldgur, M. Delarue, and M. G. Safro (1995) *Nature Struct. Biol.* **2**, 537–547.
  26. O. Nureki, D. G. Vassylyev, K. Katayanagi, T. Shimizu, S. Sekine, T. Kigawa, T. Miyazawa, S. Yokoyama, and K. Morikawa (1995) *Science* **267**, 1958–1965.
  27. P. Brick, T. N. Bhat, and D. M. Blow (1988) *J. Mol. Biol.* **208**, 83–98.
  28. S. Doublié, G. Bricogne, C. Gilmore, and C. W. Carter Jr. (1995) *Structure* **3**, 17–31.
  29. S. Onesti, A. D. Miller, and P. Brick (1995) *Structure* **3**, 163–176.
  30. D. T. Logan, M.-H. Mazaauric, D. Kern, and D. Moras (1995) *EMBO J.* **14**, 4156–4167.
  31. J. G. Arnez, D. C. Harris, A. Mitschler, B. Rees, C. S. Francklyn, and D. Moras (1995) *EMBO J.* **14**, 4143–4155.
  32. J. J. Burbaum and P. Schimmel (1991) *J. Biol. Chem.* **266**, 16965–16968.
  33. D. D. Buechter and P. Schimmel (1993) *Crit. Rev. Biochem. Mol. Biol.* **28**, 309–322.
  34. S. J. Park and P. Schimmel (1988) *J. Biol. Chem.* **263**, 16527–16530.
  35. Y.-M. Hou and P. Schimmel (1988) *Nature* **333**, 140–145.
  36. W. H. McClain and K. Foss (1988) *Science* **240**, 793–796.
  37. S. J. Park, Y.-M. Hou, and P. Schimmel (1989) *Biochemistry* **28**, 2740–2746.
  38. K. Musier-Forsyth, N. Usman, S. Scaringe, J. Doudna, R. Green, and P. Schimmel (1991) *Science* **253**, 784–786.
  39. K. Musier-Forsyth and P. Schimmel (1992) *Nature* **357**, 513–515.
  40. C. Francklyn and P. Schimmel (1989) *Nature* **337**, 478–481.
  41. S. A. Martinis and P. Schimmel (1995) in *tRNA: Structure, Biosynthesis and Function* (D. Söll

- and U. L. RajBhandary, eds.), American Society for Microbiology, Washington, DC, pp. 349–370.
42. P. Schimmel (1993) in *The Translation Apparatus* (K. H. Nierhaus, F. Franceschi, A. R. Subramanian, V. A. Erdmann, and B. Wittmann-Liebold, eds.), Plenum Press, New York, pp. 13–21.
  43. C. Francklyn, J.-P. Shi, and P. Schimmel (1992) *Science* **255**, 1121–1125.
  44. M. Frugier, C. Florentz, and R. Giegé (1994) *EMBO J.* **13**, 2218–2226.
  45. Y.-M. Hou, T. Sterner, and R. Bhalla (1995) *RNA* **1**, 707–713.
  46. C. L. Quinn, N. Tao, and P. Schimmel (1995) *Biochemistry* **34**, 12489–12495.
  47. M. E. Saks and J. R. Sampson (1996) *EMBO J.* **15**, 2843–2849.
  48. C. Francklyn and P. Schimmel (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8655–8659.
  49. J.-P. Shi, S. A. Martinis, and P. Schimmel (1992) *Biochemistry* **31**, 4931–4936.
  50. A. M. Weiner and N. Maizels (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7383–7387.
  51. H. F. Noller (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 137–156.
  52. N. Maizels and A. Weiner (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6729–6734.
  53. D. D. Buechter and P. Schimmel (1995) *Biochemistry* **34**, 6014–6019; Correction, p. 16352.
  54. J. A. Piccirilli, T. S. McConnell, A. J. Zaugg, H. F. Noller, and T. R. Cech (1992) *Science* **256**, 1420–1424.
  55. M. Illangasekare, G. Sanchez, T. Nickles, and M. Yarus (1995) *Science* **267**, 643–647.

## Aminopeptidases

[Enzymes](#) that catalyze the hydrolytic cleavage of the [peptide bond](#) that connects the N-terminal residue to the rest of a [peptide](#), [polypeptide](#), or [protein](#) are referred to as aminopeptidases (E.C. 3.4.11) (see [Peptidases](#)). The products of hydrolysis are, therefore, the released N-terminal [amino acid](#) and the remainder of the peptide chain. The latter can, in turn, also be an aminopeptidase substrate; hence these enzymes can release amino acids sequentially, ultimately resulting in complete hydrolysis of the polypeptide (though in practice this is not always observed). The principal determinant of specificity appears to be the free **a-amino group** of the N-terminal residue; hence these enzymes can release most, though not all, of the known amino acids, albeit at different rates. There are some aminopeptidases (eg, methionine aminopeptidase) that are quite limited in specificity. Usually, aminopeptidases can also act on amino acid amides and esters, which provide convenient substrates for routine assays. Some aminopeptidases sequentially remove dipeptides from the N-terminus of substrates, and these are referred to as dipeptidyl-peptidases. Others remove tripeptides and, hence, are tripeptidyl-peptidases.

Many aminopeptidases are zinc **metalloenzymes**, but some are [serine proteinases](#) or **thiol proteases**. Of those that require zinc, some have an [active site](#) containing a single ion (see [Thermolysin](#)), whereas others have a co-catalytic site that involves two closely spaced zinc ions.

Aminopeptidases are widely distributed in various tissues and cells. They can be monomeric (a single polypeptide chain) or have up to 12 subunits. Many are integral components of cell [membranes](#), but they are also found in the **cytosol**. They have a broad range of biological functions, including regulation of [hormone](#) concentration, control of the [cell cycle](#), and recovery of amino acids

from dietary peptides and proteins (1). All proteins synthesized by **eukaryotic** cells begin at their N-terminus with **methionine**, and its removal by methionine aminopeptidase is often crucial, not only for biological function of the protein but even for cell survival. Aminopeptidases also play important roles in the food industry—for example, ripening of cheese (2) and production of soy sauce (3).

Aminopeptidases are inhibited by the antitumor antibiotic bestatin [(2*S*,3*R*)-3-amino-2-hydroxy-4-phenylbutanol]-L-leucine, isolated from culture filtrates of *Streptomyces olivoreticuli*. It is a potent inhibitor of bovine lens aminopeptidases, with a  $K_i$  of 1.3 nM, and has numerous biological activities when administered to laboratory animals.

Membrane-bound aminopeptidases have often been identified on the basis of some other property, and then, once their amino acid sequence has been established, they are recognized to be aminopeptidases. Thus, the B-lymphocyte differentiation factor BP-1/6C3, whose expression correlates with proliferation and transformation of immature **B cells** (antibody-producing lymphocytes), has been shown to be identical to glutamyl aminopeptidase, also known as aminopeptidase A (4). Also the myeloid leukemia antigen CD-13 has been identified as aminopeptidase N (5), and the amino acid sequence of leukotriene A4 hydrolase revealed an aminopeptidase-like structure that led to the recognition of its aminopeptidase activity (6). Aminopeptidase N has also been shown to have a function unrelated to its enzymatic activity; that is, it serves as a cell-surface **receptor** for certain coronaviruses that cause upper respiratory infections (7).

The amino acid sequence and three-dimensional structure of leucine aminopeptidase from bovine lens have been determined (8). This is a broad-specificity cytosolic enzyme found in tissues of all organisms. It has a high degree of sequence similarity to several other aminopeptidases, which suggests that they all share similar structures and catalytic mechanisms.

#### Bibliography

1. A. Taylor (1993) *FASEB J.* **7**, 290–298.
2. J. Meyer, D. Howald, R. Jordi, and M. Fuerst (1989) *Milchwissenschaft* **44**, 678–681.
3. T. Nakadai (1988) *Nippon Shoyu Kenkyusho Zasshi* **14**, 50–56.
4. Q. Wu et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 993–997.
5. T. Inoue et al. (1994) *J. Clin. Endocrinol. Metab.* **79**, 171–175.
6. J. Z. Haeggstrom, A. Wetterholm, B. L. Vallee, and B. Samuelsson (1990) *Biochem. Biophys. Res. Commun.* **173**, 431–437.
7. B. Delmas et al. (1992) *Nature* **357**, 417–420; Yeager et al. *ibid.* 420–422.
8. S. K. Burley, P. R. David, A. Taylor, and W. N. Lipscomb (1990) *Proc. Natl. Acad. Sci. USA* **87**, 6878–6882.

#### Suggestions for Further Reading

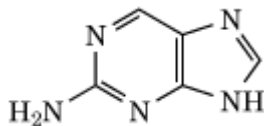
9. A. Taylor (1993) Aminopeptidases: towards a mechanism of action. *Trends Biochem. Sci.* **18**, 167–172.
10. H. Kim and W. N. Lipscomb (1994) Aspartate transcarbamylase from *Escherichia coli*: activity and regulation. *Adv. Enzymol.* **68**, 153–213.

## 2-Aminopurine (AP)



AP is a base analogue that is a [mutagen](#) in a wide range of systems and species ([1](#), [2](#)) (Fig. [1](#)). Freese ([3](#), [4](#)) originally described A.T → G.C and G.C → A.T transition [mutations](#) following AP treatment of **bacteria** and **bacteriophage**. The results were interpreted as resulting from **tautomeric** shifts and changes in base pairing behavior following the incorporation of AP into DNA. AP is readily metabolised to form deoxy-2-aminopurine triphosphate (dAPTP) ([5](#)), which may be incorporated opposite thymine during [DNA replication](#), to form an AP.T base pair. During subsequent rounds of replication, incorporation of dCMP opposite the AP would lead directly to an A.T → G.C transition mutation. Law et al. ([6](#)) found that a DNA duplex containing AP.C is thermodynamically more stable than a DNA duplex containing A.C. This is in agreement with the suggestion that the rate of insertion of an improper base by a **DNA polymerase** is determined by differences in stability between the newly formed mismatch site and the corresponding normal **Watson–Crick base pair** ([7](#)). The probability that transition mutations will occur results from a combination of the likelihood of tautomeric shifting, the availability of the analogues to the DNA replicating machinery, and the discriminatory behavior and fidelity of the DNA polymerases.

**Figure 1.** Structure of the base analogue 2-aminopurine.



Although best known as a base-pair substitution mutagen, the earliest reports of AP activity focused on the production of unequal, multipolar **mitosis** in mouse [tissue culture](#) cells, producing **aneuploid** daughter cells also bearing chromosome mutations ([1](#)). AP also causes [frameshift mutations](#)—for example, in a set of lacZ mutant variants of *Escherichia coli* ([2](#)). *E. coli dam* mutants have reduced ability for methylation of adenine. Many of these mutants also have enhanced sensitivity to AP and increased mutability by this base analogue ([8](#)). It appears that AP is partially able to saturate or inactivate the methylation-directed [mismatch repair](#) system, allowing the escape from repair of replication errors that lead to frameshift mutation. This results in indirect mutations that can be detected at certain sites ([2](#)).

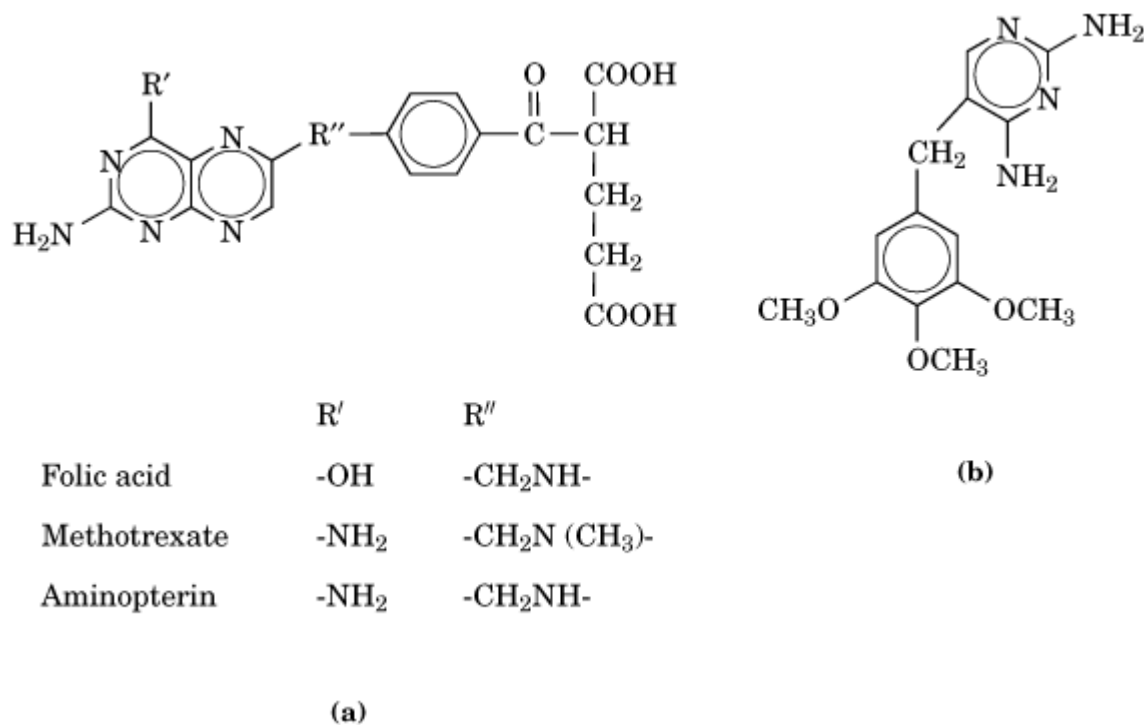
#### Bibliography

1. A. Ronen (1979) *Mutat. Res.* **75**, 1–47.
2. C. G. Cupples, M. Cabrera, C. Cruz, and J. H. Miller (1990) *Genetics* **125**, 275–280.
3. E. Freese (1959a) *Proc. Natl. Acad. Sci. USA* **45**, 622–633.
4. E. Freese (1959b) *J. Mol. Biol.* **1**, 87–105.
5. E. G. Rogan and M. J. Bessman (1970) *J. Bacteriol.* **103**, 622–633.
6. S. M. Law, R. Eritja, M. F. Goodman, and K. J. Breslauer (1996) *Biochemistry* **35**, 12329–12337.
7. M. F. Goodman, R. L. Hoskins, R. Lasker, and D. N. Mhaskar (1993) *Basic Life Sci.* **31**, 409–423.
8. B. W. Glickman, P. van den Elsen, and M. Radman (1978) *Mol. Gen. Genet.* **163**, 307–312.

## Aminopterin, Methotrexate, Trimethoprim, and Folic Acid

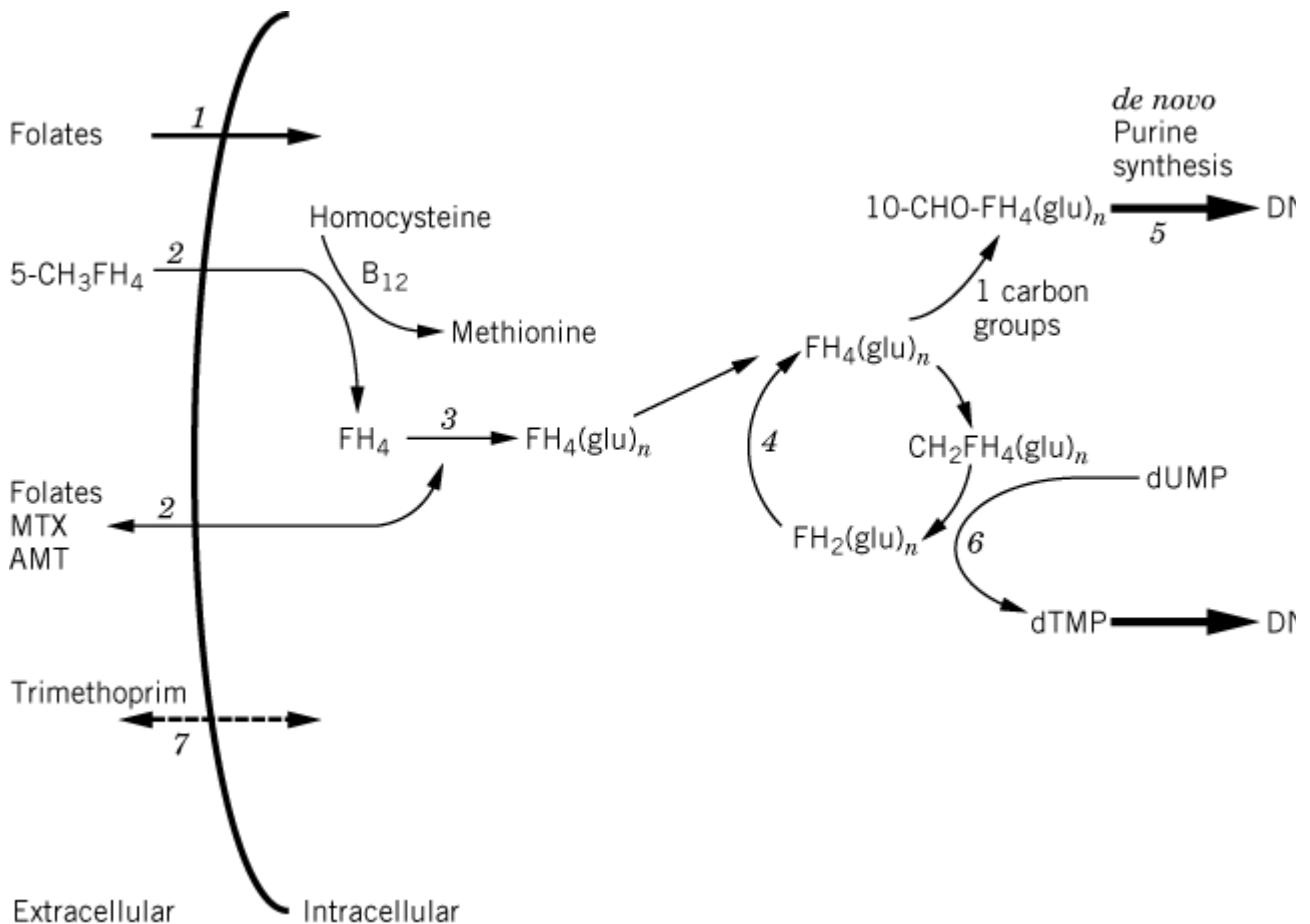
Aminopterin, methotrexate, and trimethoprim are all analogs of *folic acid* (Fig. 1) that antagonize folate-dependent metabolic pathways (1, 2). Folate is a water-soluble vitamin. Although some bacteria synthesize folate, mammalian cells cannot, and consequently it is an absolute dietary requirement. Reduced folates function as **cofactors** in many metabolic processes common to nearly all cells, such as **thymidylate** and **purine** synthesis and donation of methyl groups (Fig. 2) (1). Folates gain entry into mammalian cells by either the reduced folate carrier (RFC) and/or a hydrophobic membrane-associated folate receptor (FR) found in placental, choroid plexus and kidney cells (3). The FR is unidirectional and, once inside, cellular retention of folate is enhanced by its polyglutamation, catalyzed by the **enzyme** foylpolylglutamyl synthetase (FPGS) (4). FPGS adds up to six or seven **glutamic acid** residues via an unusual **peptide bond** through their g-carboxyl group, rather than the normal a. The number of glutamic acid residues added may play a role in regulating and distributing reduced folates, and in guaranteeing their availability as cofactors. Polyglutamation also increases the affinity of folate for folate-dependent enzymes, such as **thymidylate synthase** (TS), aminoimidazole carboxamide ribonucleotide and glycinamide ribonucleotide transformylases. The latter two are involved in purine synthesis. In cells synthesizing DNA, 5,10-methylenetetrahydrofolate serves in thymidylate synthesis as a methyl group donor for converting dUMP to dTMP. This is the only reaction in which tetrahydrofolate is partially oxidized to dihydrofolate (1, 4). **Dihydrofolate reductase** (DHFR) is the critical enzyme involved in converting dihydrofolate back to tetrahydrofolate, thus maintaining reduced folate pools to serve as one-carbon group carriers (see Fig. 2).

**Figure 1.** (a) The structures of folic acid and folate analogs: folic acid and analogs; (b) the structure of trimethoprim.



**Figure 2.** Folate transport, accumulation, and target enzymes. (1) folate receptor (FR); (2) reduced folate carrier (RFC), which is bidirectional; (3) foylpolylglutamyl synthetase (FPGS); (4) dihydrofolate reductase (DHFR); (5) aminoimidazo-

carboxamide ribonucleotide and glycinamide ribonucleotide tranformylases; (6) thymidylate synthase (TS); (7) passive diffusion of lipid-soluble compounds.  $\text{FH}_4$  and  $\text{FH}_4(\text{glu})_n$ , tetrahydrofolate with and without the added glutamyl residue;  $\text{FH}_2$ , dihydrofolate;  $\text{CH}_2\text{FH}_4$ , 5,10 methylene- $\text{FH}_4$ ; 10- $\text{CHO-FH}_4$ , 10-formyl  $\text{FH}_4$ ;  $\text{B}_{12}$ , vitamin [B12](#) [Adapted from Hum and Kamen (1996) *Investigational New Drugs* **14**, 110–111.]



Folate homeostasis is recognized as important. The structure of DHFR has been determined by [X-ray crystallography](#). One proposed mechanism of resistance to methotrexate involves DHFR [gene amplification](#) (5-7). The FR is overexpressed in some carcinomas, and its gene has been localized to the 11q13 region (3). Human FPGS has been cloned and mapped to chromosome 9q (4). The [complementary DNA](#) for RFC has been **cloned** and the RFC gene localized to the long arm of chromosome 21 (8-11).

## 1. Methotrexate (MTX)

Folate analogs entered cancer therapy in the 1940s, when aminopterin was successfully used to induce temporary remissions in children with acute lymphoblastic leukemia (ALL) (12). Because of a better therapeutic index, MTX eventually emerged as the antifolate used clinically in treating cancers, such as leukemias, lymphomas, osteosarcoma, breast cancer, choriocarcinoma, head and neck cancers, and nonmalignant disorders, such as arthritis and asthma (4, 13).

Like folates, MTX enters mammalian cells via the RFC. MTX has an apparent **dissociation constant** in the micromolar range, and via the FR has higher affinity in the nanomolar range (3). Cellular retention of MTX is enhanced by polyglutamation, which also enhances the affinity of the drug for enzymes. Once inside the cell, MTX acts as a tight-binding **competitive inhibitor** of DHFR. This leads to an accumulation of dihydrofolate and depletion of the reduced folate pools in cells actively

making dTMP via the *de novo* pathway. Accumulated dihydrofolate polyglutamates are also inhibitors of TS and the enzymes involved in the *de novo* synthesis of purines (1). The resulting imbalance of nucleotides causes base substitutions, which lead to errors in DNA synthesis and ultimately to cell death.

## 2. Aminopterin (AMT)

Aminopterin (AMT) was the first antifolate drug used clinically in childhood ALL (12) (see previous). Although more potent than MTX, in preclinical studies the toxicity of AMT was more severe and more unpredictable (13-15). Transport and metabolic studies *in vitro* have shown that AMT is the preferred substrate (16, 17). This results in greater AMT accumulation at lower concentrations and more complete polyglutamation, leading to improved cellular retention for the cytotoxic effect.

The greater potency of AMT compared to MTX has led to renewed clinical interest. AMT may find a role in treating patients with resistant or refractory malignancies or those patients in whom *in vitro* studies indicate AMT is the better choice on the basis of metabolism and accumulation (16).

## 3. Trimethoprim

Trimethoprim is an antibacterial agent developed in the 1950s. Studies by Hitchings (2) showed that its mechanism of action is the competitive inhibition of DHFR. Unique to trimethoprim is its much greater affinity (50,000- to 100,000-fold) for bacterial DHFR than for mammalian DHFR (18, 19). Trimethoprim is lipid-soluble and enters cells rapidly without requiring specific transport mechanisms. Its selective toxicity is further enhanced by the ability of folinic acid to reverse even the slight effects of trimethoprim on mammalian cells, whereas bacterial cells, unable to transport folinic acid, are not rescued by folinic acid administration (19, 20). In addition to its antibacterial activity, the drug has activity against *Pneumocystis carinii*, an opportunistic infection of the lungs encountered in severely immunocompromised patients (especially HIV patients).

## Bibliography

1. M. C. Hum and B. A. Kamen (1996) Investigational New Drugs **14**, 110–111.
2. G. H. Hitchings (1973) J. Infect. Dis. **128**(suppl), S433–S436.
3. S. Weitman, R. G. W. Anderson, and B. A. Kamen (1994) In *Vitamin Receptors: Vitamins as Ligands in Cell Communities* (K. Dakshinamurti, ed.), Cambridge University Press, Cambridge, UK, pp. 106–136.
4. T. A. Garrow, A. Admon, and B. Shane (1992) Proc. Natl. Acad. Sci. USA **89**, 9151–9155.
5. F. W. Alt et al. (1978) J. Biol. Chem. **253**, 1357–1370.
6. J. R. Bertino (1993) Ode to Methotrexate, J. Clin. Oncology **11**(1), 5–14.
7. E. Chu and C. H. Takimota (1993) In *Principles & Practice of Oncology* (V. T. De Vita, S. Hellman, and S. A. Rosenberg, eds.), Lippincott, Philadelphia, pp. 358–374.
8. J. A. Moscow et al. (1995) Cancer Res. **55**, 3790–3794.
9. P. D. Prasad et al. (1995) Biochem. Biophys. Res. Commun. **206**, 681–687.
10. S. C. Wong et al. (1995) J. Biol. Chem. **270**, 17468–17475.
11. T. L. Yang-Feng et al. (1995) Biochem. Biophys. Res. Commun. **210**(3), 874–879.
12. S. Farber et al. (1948) N. Engl. J. Med. **238**, 787–793.
13. A. Goldin et al. (1955) J. Natl. Cancer Inst. **15**, 1657–1664.
14. F. M. Sirotnak and R. C. Donsbach (1975) Biochem. Pharmacol. **24**, 156–158.
15. F. S. Philips et al. (1973) Cancer Res. **33**, 153–158.
16. A. Smith et al. (1996) Clin. Cancer Res. **2**(1), 69–73.
17. B. G. Rumberger, J. R. Barrueco, and F. M. Sirotnak (1990) Cancer Res. **50**, 4639–4643.

18. Anonymous (1971) Trimethoprim-sulfamethoxazole. *Drugs*, **1**(1), 8–53.
19. J. J. Burchall (1973) *J. Infect. Dis.* **128**(suppl), S437–S441.
20. W. Brumfitt and J. M. T. Hamilton-Miller (1980) *Brit. J. Hosp. Med.* **23**(3), 281–288.

### Suggestions for Further Reading

21. J. R. Bertino, B. A. Kamen, and A. Romanini (1997) "Folate antagonists". In *Cancer Medicine* (J. F. Holland, R. C. Bast, D. L. Morton, E. Frei, D. W. Kufe, and R. R. Weicheslbaum, eds.), Williams and Wilkins, Baltimore, pp. 907–922.
22. R. L. Blakley and S. J. Benkovic (1984) *Folates and Pterins: Chemistry and Biochemistry of Folates*, Wiley, New York, Vol. 1.
23. C. H. Takimoto and C. J. Allegra (1995) New antifolates in clinical development, *Oncology* **9**, 649–659.

### Amphipathic

Compounds that contain both highly [polar](#) and [nonpolar](#), **hydrophobic** moieties are *amphipathic*. In the presence of [water](#), these compounds tend to aggregate, forming structures that are strongly influenced by the relative extents of the polar and hydrophobic characters of the molecule, with the polar groups acting to shield the hydrophobic moieties from the solvent. In [lipids](#), the hydrophobic group is the hydrocarbon domain of the aliphatic chain (see [Lipids](#) and [Fatty Acids](#)). Typical examples are (a) the sodium salts of fatty acids, which form micelles in water, and (b) phospholipids, which form bilayers in membranes. The presence of significant amounts of nonpolar [amino acids](#), such as [alanine](#), [valine](#), [leucine](#), [isoleucine](#), [proline](#), [methionine](#), [phenylalanine](#), and [tryptophan](#), make certain [proteins](#) amphipathic; these amino acids are believed to increase the accessibility of the proteins to the hydrophobic domain of the membrane bilayer (see [Fluid Mosaic Model](#)).

### Amphoteric

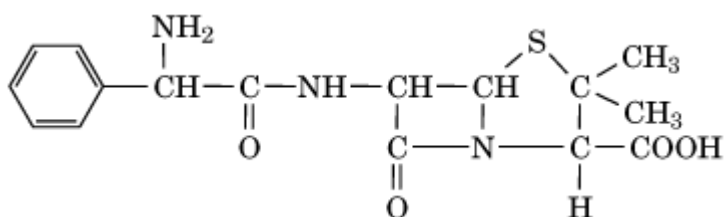
An amphoteric molecule is one that has both acidic and basic chemical features. For example, proteins are amphoteric electrolytes because they generally have both acidic and basic groups on the side chains of their amino acid residues. Extreme examples of amphoteric molecules are the ampholytes used to establish pH gradients in [isoelectric focusing](#).

### Ampicillin

Ampicillin (D[–]-a-aminobenzylpenicillin; Fig. [1](#)) is a member of a growing family of antimicrobial agents known as the semisynthetic [penicillins](#). The semisynthetic penicillins are derivatives of the natural product, 6-aminopenicillanic acid, that have been deliberately modified chemically ([1](#)). The chemical modifications are introduced to create new compounds with specific desirable properties.

In the case of ampicillin, the addition of the aminobenzyl side chain results in a product with an increased resistance to acidic pH. Thus, ampicillin is medically significant because it was the first penicillin to be administered orally in chemotherapeutic practice. Furthermore, ampicillin exhibits a broader antibacterial spectrum than do the naturally occurring penicillins and is effective against many [Gram-negative bacterial](#) species. The outer membrane component of the Gram-negative bacterial cell wall represents a permeability barrier to many antibiotics, including the natural penicillins (2). The broad activity spectrum and the relative low cost of ampicillin have made it an invaluable tool in molecular biology and genetics for studying Gram-negative model organisms, such as *Escherichia coli* and *Haemophilus influenzae*. Applications that employ ampicillin are summarized below.

**Figure 1.** Structure of ampicillin.



The first application of penicillin in genetics was for the selection of auxotrophic mutants of *E. coli* (3, 4). This technique, penicillin selection, was based on the fact that penicillin kills only actively growing bacteria, whereas nongrowing bacteria are penicillin-tolerant (see [Penicillin](#)). Although benzylpenicillin (penicillin G) was employed in the original studies describing this technique, ampicillin would be far more effective for this purpose because of its broad spectrum.

Molecular biologists have extensively exploited genetic elements encoding [b-lactamase](#) (see [Penicillin-Binding Proteins](#)), especially the enzyme designated TEM-1. TEM-1 is a broad spectrum b-lactamase that confers effective high level resistance to penicillins (5). The gene encoding TEM-1 was incorporated into the first plasmid [cloning](#) vectors (6) and has been widely used in cloning vectors since. Ampicillin has been used routinely for the selection and maintenance of bacteria carrying recombinant plasmids encoding b-lactamase, and this is undoubtedly its most common application in molecular biology. For this purpose, ampicillin is incorporated into bacteriological media at final concentrations ranging from about 50–100 µg/mL. For *E. coli*, these levels are about 10–20 times higher than the minimum inhibitory concentration (MIC) of ampicillin. The MIC is defined as the minimum concentration of the antibiotic that is necessary to inhibit bacterial growth.

When a mixture of ampicillin-sensitive and ampicillin-resistant bacteria are plated on solid media containing ampicillin, such as for selection of transformants carrying a b-lactamase-encoded plasmid, it is not uncommon to find large colonies formed by ampicillin-resistant bacteria surrounded by a zone of smaller colonies. The ampicillin-resistant bacteria in the large central colonies produce and secrete b-lactamase. The activity of the secreted b-lactamase creates a zone of reduced ampicillin concentration around the resistant colonies, and this permits the ampicillin-sensitive bacteria in the vicinity to grow. The subsequent growth of the ampicillin-sensitive bacteria results in the formation of the smaller so-called satellite colonies. The satellite colonies normally do not represent a major hindrance in these procedures, because the desired ampicillin-resistant bacteria can be readily purified by streaking on an ampicillin-containing medium. However, the problem of satellite colony formation can be minimized by substituting carbenicillin, another semisynthetic penicillin (see Fig. 1 in [Penicillin](#)) for ampicillin in the selection medium at a concentration of 50–100 µg/mL (7). Carbenicillin is less susceptible to hydrolysis by the b-lactamase and is therefore less likely to promote satellite colony formation. It is therefore often used in place of ampicillin.

The b-lactamase gene has also been introduced into plaque-forming and defective derivatives of the *E. coli* bacteriophage Mu (8). These ampicillin-selectable phages have been used for mutagenesis and for the construction of *lac* fusions in applications that take advantage of the ability of Mu to **transpose** randomly.

The concept of **transposon** mutagenesis has been applied to transposable genetic elements encoding b-lactamases for the generation of random **gene fusions** that are directly selectable with ampicillin (or carbenicillin). For example, a derivative of Tn3 designated Tn3-HoHo1 (9) is a *lacZ*-containing transposon capable of producing both transcriptional and translational **beta-galactosidase** fusions. Although it was originally developed for studies on *Agrobacterium tumefaciens*, it has been adapted for use in other bacterial genera.

b-Lactamase is a periplasmic enzyme. It has served as an important model for studying protein export to the bacterial periplasm. Urbain et al. (10) have recently developed a technique for quantifying b-lactamase activity in cultures of *E. coli* carrying recombinant plasmids that confer ampicillin resistance. Their assay involved determining the conversion of ampicillin to aminobenzylpenicilloic acid in periplasmic extracts of cells by quantitative high performance liquid chromatography (**HPLC**). The procedure may be useful for investigating the mechanism of b-lactamase translocation. It may also prove useful in studies on protein expression. For example, since b-lactamase is expressed constitutively from ampicillin-selectable recombinant plasmids, its activity could serve as a useful internal standard in protein coexpression studies.

b-lactamase has also been used as a genetic tool for studying **membrane proteins**, and these applications are based on the fact that b-lactamase is an exported protein (11, 12). A plasmid vector that permits the *in vitro* construction of translational fusions between a gene of interest, encoding either a membrane protein or an exported protein, and the mature form (ie, the exported form) of the TEM b-lactamase has been described (13). Transformants carrying recombinant plasmids with in-frame fusions are ampicillin-selectable. This technique may be used for the analysis of protein export signals or for determining the topological organization of membrane proteins. A strategy for maximizing yields of membrane and exported proteins based on this vector has also been described (14). It is notable that alkaline phosphatase has been used widely for studying protein export signals and for topological mapping of membrane proteins (15). The b-lactamase system is an attractive alternative to **alkaline phosphatase** for these purposes (13, 14) (see **Reporter Genes**). For example, b-lactamase fusions are directly selectable (with ampicillin), whereas the identification of alkaline phosphatase fusions is based on phenotypic screening. Moreover, only the periplasmic form of alkaline phosphatase is enzymatically active, whereas both the cytoplasmic and periplasmic forms of b-lactamase are active. Consequently, the cellular location of the b-lactamase fusions can be determined on the basis of the levels of ampicillin resistance; cytoplasmic b-lactamase confers ampicillin resistance only at high cell density, whereas periplasmic b-lactamase confers ampicillin resistance at low cell density. As an extension of these studies, Broome-Smith et al. (16) have constructed a transposable b-lactamase element, designated *TnblaM*, that is equivalent to the transposable alkaline phosphatase element, *TnphoA*. The b-lactamase fusions constructed with *TnblaM* are directly selectable with ampicillin, and this is a major advantage over alkaline phosphatase system, which as already noted is based on phenotypic screening.

## Bibliography

1. J. H. C. Naylor (1991) in *50 Years of Penicillin Application. History and Trends*, Public Ltd., Czech Republic, 64–74.
2. H. Nikaido and M. Vaara (1985) *Microbiol. Rev.* **49**, 1–32.
3. J. Lederberg and N. Zinder (1948) *J. Am. Chem. Soc.* **70**, 467–468.
4. B. D. Davis (1949) *Proc. Natl. Acad. Sci. USA* **35**, 1–10.
5. N. Datta and M. D. Richmond (1966) *Biochem. J.* **98**, 204–210.

6. F. Bolivar, R. L. Rodriguez, M. C. Betlach, and H. W. Boyer (1977) *Gene* **2**, 75–93.
7. F. M. Ausubel, R. Brent, R. E. Kingston, D. D. Moore, S. D. Seidman, J. A. Smith, and K. Struhl (1994) *Current Protocols in Molecular Biology*, Wiley, New York, p.1.8.7.
8. E. A. Groisman (1991) *Meth. Enzymol.* **204**, 180–212.
9. S. E. Stachel, G. An, C. Flores, and E. W. Nester (1985) *EMBO J.* **4**, 891–898.
10. J. L. Urbain, C. M. Wittich, and S. R. Campion (1998) *Anal. Biochem.* **260**, 160–165.
11. J. K. Broome-Smith, M. Tadayyon, and Y. Zhang (1990) **4**, 1637–1644.
12. M. Tadayyon, Y. Zhang, S. Gnaneshan, L. Hunt, F. Mehraein-Ghomi, and J. K. Broome-Smith (1992) *Biochem. Soc. Trans.* **20**, 598–601.
13. J. K. Broome-Smith and B. G. Spratt (1986) *Gene* **49**, 341–349.
14. J. K. Broome-Smith, L. D. Bowler, and B. G. Spratt (1989) *Mol. Microbiol.* **3**, 1813–1817.
15. C. Manoil, J. J. Mekalanos, and J. Beckwith (1989) *J. Bacteriol.* **172**, 515–518.
16. M. Tadayyon and J. K. Broome-Smith (1992) *Gene* **111**, 21–26.

## Amyloid

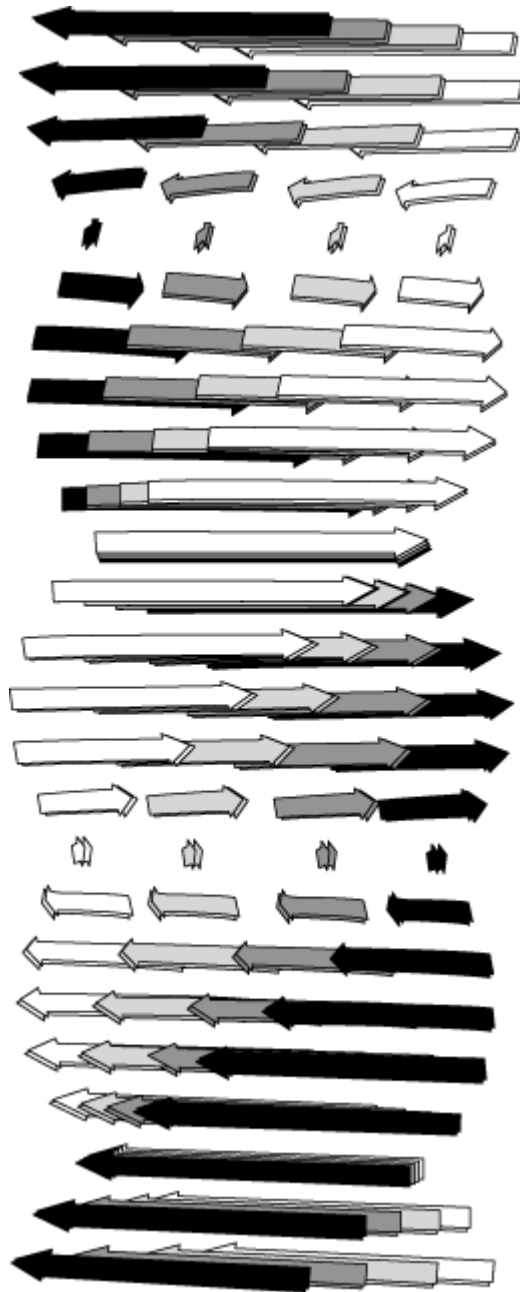
Amyloid is an insoluble, proteinaceous, fibrous material associated with a number of prominent disease states (1), which can also form spontaneously *in vitro* from [oligopeptides](#) and **denatured** proteins. The disease states in which amyloid has been implicated as an important, or even causal, factor include Alzheimer's disease (see [Amyloid Precursor Protein](#)), the transmissible spongiform encephalopathies (see [Scrapie](#)), non-insulin-dependent (type II) diabetes, and a number of polyneuropathies. It is thought that the extreme stability of amyloid fibrils permits their progressive accumulation in the extracellular spaces of vital organs whose functioning is thereby inhibited, leading to organ failure and death. The *in vivo* development of amyloid deposits from globular precursor proteins is linked to either genetic [mutation](#), incorrect processing, or the abnormal accumulation of wild-type proteins. Probably the most remarkable feature of amyloid is that its molecular structure appears to be constant and independent of the protein precursors. For example, all amyloids, regardless of the disease involved or the source of the fibrils, share similar morphological (2), tinctorial (3), and structural (4, 5) properties. This implies that amyloid formation is not a simple aggregation process but a structural conversion of globular proteins to a state that can be incorporated into a particular type of [fibrous protein](#) structure. The amyloid diseases may therefore be defined as diseases of protein misfolding.

At least 15 different proteins can form amyloid fibrils *in vivo* (1), and recently the involvement of amyloid fibrils in Huntington's disease has also been demonstrated (6). In the **electron microscope**, amyloid fibers are about 100 Å (10 nm) in diameter and usually straight or only slightly curved. Closer examination shows the fibrils to be composed of a number of smaller diameter protofilaments, arranged in a more or less parallel array. In those fibrils so far examined, the protofilaments give the fibrils the appearance of hollow cylinders or ribbons when observed in cross section. Although the number, size, and arrangement of these constituent protofilaments appear to differ somewhat, the ribbons may be merely unrolled or unformed cylinders, and hence the variation may not be so great as appears. The molecular structure of the amyloid fibrils has been established by X-ray [fiber diffraction](#). The first diffraction patterns (7, 8) showed the intense 4.7 Å (0.47 nm) meridional and 10 Å (1 nm) equatorial reflections that have been subsequently shown to characterize all amyloid X-ray patterns. These reflections indicate that the molecular structure of amyloid is composed of [beta-sheets](#) arranged parallel to the fiber axis, with their constituent [beta-strands](#) at right



angles to the axis of the fibril (4, 9). This so-called “cross-b” structure is quite different from the more common insect silk **chorions**, whose fibers are formed from b-sheets having their b-strands parallel to the fiber axis. The first use of intense synchrotron X-ray sources on amyloid (10) extended the observable X-ray pattern to 2 Å (0.2 nm) and revealed how this “classic” amyloid model could be reconciled with low-energy twisted b-sheets (11). These new data revealed a previously unobserved repeat distance of 115 Å (11.5 nm) along the fiber axis of the transthyretin amyloid of familial amyloidotic polyneuropathy (FAP), which was proposed (10) to correspond to the repeat distance of a complete helical turn of a twisted b-sheet, whose helix axis was parallel to the fiber axis. The basic helical unit therefore consists of a segment of b-sheet whose 24 b-strands complete one helical turn. This **beta-helix** model of the molecular structure of the amyloid protofilament is shown in Figure 1. Given that all amyloid fibrils share common morphological, tinctorial, and structural properties, it is reasonable to see the b-helix model as characteristic of all amyloid protofilaments. The b-helix model is distinct from other fibrous proteins, showing that the amyloid structure is unique. Its features enable low energy twisted b-sheets to be incorporated in a linear fibril in such a way that continuous b-type **hydrogen bonding** can be extended over the total length of the fibril. This, together with the opportunity for a continuous **hydrophobic** core to stabilize the b-sheet interactions along the length of the fibril, can be reasonably held to account for the known extreme stability of amyloid fibrils, which is central to their role in disease.

**Figure 1.** Drawing of one complete turn of the proposed b-helix structure of the amyloid fibril. The arrows represent b-strands and are shaded to represent the different b-sheets. The direction of the arrows has no significance, and no b-strand connections are shown.

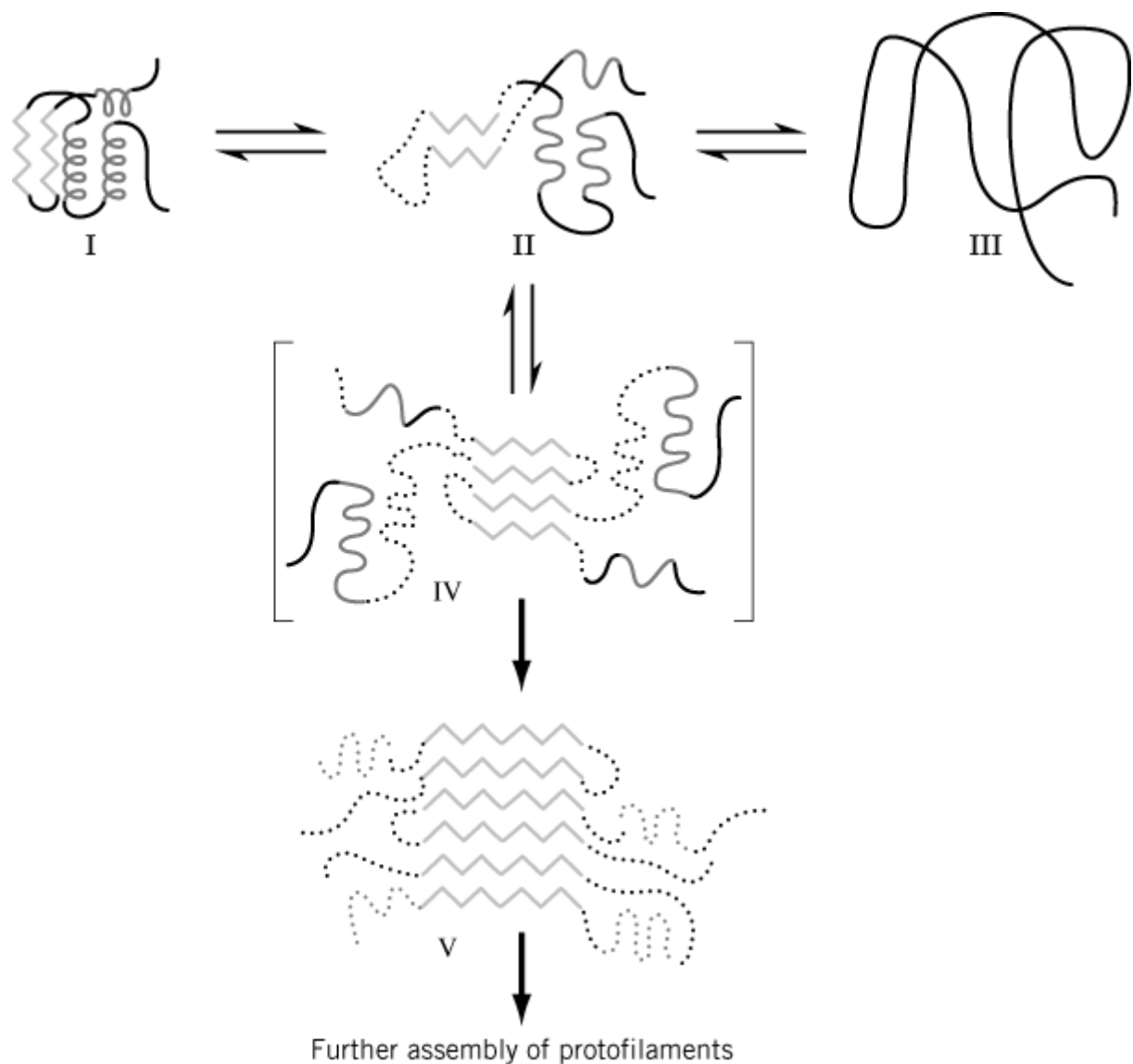


A remarkable property of amyloid fibrils is that they appear to be very similar whatever their source. This has recently been critically tested (5) by showing that the high resolution synchrotron X-ray patterns from eight different amyloids, of both disease and synthetic origins, have the same structural characteristics. This observation suggests that amyloid is a highly stable generic structure capable of accommodating proteins and peptides with a wide range of chain lengths and structures within a common fibrillar form. The repetition of b-strands in the fibril model allows polypeptide chains of different lengths to occupy the same structural framework. For example, a 10-residue peptide could form a single b-strand of the amyloid structure, consistent with the known lower limit on the length of peptides capable of forming amyloid, while longer polypeptide chains could fold back and forward on themselves, thereby forming a number of consecutive b-strands of the amyloid b-sheets. The need for loops of chain of different lengths or conformation between successive b-strands to maximize the b-strand propensity, and to generate a stable hydrophobic core between the b-sheets, is likely to give each fibril a characteristic structure within the core amyloid framework. Thus it is probably realistic to consider amyloids as a closely related family of fibrils, each with its own detailed structural differences, while maintaining a similar overall core structure.

The nature of the processes that convert the globular, soluble protein precursors into amyloid fibrils, known as either *amyloidogenesis* or *fibrillogenesis*, are important, as they underpin the development of amyloid diseases and represent a novel field of scientific investigation. It is evident that some of the proteins forming the predominantly  $\beta$ -structured amyloid fibrils are themselves largely  $\alpha$ -helical in their normal soluble states, for example, [lysozyme](#), the [prion](#) protein, and the Ab peptide in Alzheimer's disease (see [Amyloid Precursor Protein](#)). The nature of this structural rearrangement has been particularly studied in two amyloidogenic proteins: transthyretin (TTR) and lysozyme. Transthyretin amyloidosis is associated with familial amyloidotic polyneuropathy (FAP) and senile systemic amyloidosis (SAA). Both the wild-type protein (SAA) and more than 50 different genetic variants (FAP) ([13](#)) give rise to amyloid. That so many variants of transthyretin are amyloidogenic (nearly 40% of all its residues), suggests that the structure of this protein allows a particularly facile transformation into the fibrillar form. As transthyretin is a 55-kDa homotetramer, with the four monomers associating in such a way as to generate a stack of four 8-stranded  $\beta$ -sheets ([14](#)), there is a superficial structural similarity with amyloid fibril ([10](#)). Nevertheless there is evidence that the tetrameric form of TTR is not a building block for the fibril ([15](#)), and the discovery of a molecular hotspot in the pattern of amyloidogenic variants ([16](#)) suggests that structural change to the monomer is also required prior to its incorporation into the fibril. It is therefore evident that even all- $\beta$  proteins, such as TTR, must undergo a significant degree of structural rearrangement to form amyloid.

The human lysozyme molecule (see [Lysozymes](#)) has about 35%  $\alpha$ -helix and only 10%  $\beta$ -sheet structure, segregated into an  $\alpha$ - and a  $\beta$ -domain. Hence the structural changes that must accompany its incorporation into the predominantly  $\beta$ -structured amyloid fibril will need to be extensive. Two genetic variants of human lysozyme are known to form amyloid *in vivo*, Ile56Thr and Asp67His ([17](#)). A thorough biophysical study of these variants has highlighted major differences in their stability and folding behavior as compared with the wild type ([18](#)). It has generated a possible mechanism of the  $\alpha$ - to  $\beta$ -structural transformation of fibril formation of lysozyme, shown in [Figure 2](#), which may be applicable *mutatis mutandis* to other amyloidogenic proteins. Thermal denaturation studies show that the variants unfold at temperatures at least 10°C lower than wild-type lysozyme and, unlike the wild type, do not regain activity when returned to physiological conditions. Fourier transform infrared spectroscopy (see [Vibrational Spectroscopy](#)) of Asp67His lysozyme shows a significant gain in  $\beta$ -structure and loss of  $\alpha$ -structure in the soluble material, demonstrating that an  $\alpha$ -to- $\beta$  structural interconversion is associated with fibril formation. Detection of an unfolding intermediate displaying the characteristic binding to the hydrophobic dye 1-anilinonaphthalene sulfonic acid (ANS) is consistent with a **molten globule**-like intermediate, as shown in [Figure 2](#). [Hydrogen exchange](#) shows that the flexibilities of the native folds of the variants have been dramatically increased by the mutations. Inspection of the [X-ray crystallography](#) structures of the two variants ([18](#)) suggests that the key to both amyloidogenic mutations in human lysozyme lies in the effect they produce at the interface between the  $\alpha$ - and  $\beta$ -domains, with the result that domain adhesion may be weakened. Both variant lysozymes unfold dramatically faster than the wild-type protein, because the docking of the two domains, required for achieving the final rigid protein fold ([19](#)), is compromised by the presence of a [threonine](#) side-chain in the place of the wild-type **hydrophobic** anchoring residue, [isoleucine](#).

**Figure 2.** Schematic drawing of the possible mechanism of fibril formation in human lysozyme. (—)  $\beta$ -sheet structure; (—)  $\alpha$ -helical structure. A partly-folded, molten-globule form (II), distinct from the native (I), and denatured (III) forms, self-associates through the  $\beta$ -domain (IV) to initiate fibril formation. This provides a template for further deposition and the development of a  $\beta$ -sheet core structure of the amyloid fibril (V). Undefined structures are shown as broken lines.



Detailed study of the structure and behavior of variants of lysozyme has led to the proposal that transient populations of amyloidogenic proteins in a molten-globule state that lacks global cooperativity are an important feature of the conversion from a soluble to the fibrillar state (18). These observations lead to a model of the mechanism for amyloid formation for amyloidogenic lysozyme shown in Figure 2. Such a mechanism may also operate for other amyloidogenic proteins. For example, the structure of the 121–231 domain of the prion protein (20) also demonstrates that mutations associated with prion disease are involved in the maintenance of the hydrophobic core, which may be related to the earlier suggestion of a molten globule intermediate in the formation of aggregates in scrapie (21) (see [Prion](#)). Recent evidence supports the hypothesis that conformational plasticity is a key feature in prion fibril formation. A similar observation has been made for the Ab protein forming the amyloid in Alzheimer's disease (22) (see [Amyloid Precursor Protein](#)). The mechanism for helix-to-sheet conversion in lysozyme, proceeding from soluble forms of the amyloidogenic precursor proteins through transient populations of intermediates with the characteristics of molten globules, and on to intermolecular  $\beta$ -sheet association, seems to parallel the processes in the aggregation of prions and the Ab amyloid peptides, and it may occur generally in the amyloidoses. Studies of these processes is certain to extend our knowledge of **protein folding** behavior into new areas and coincidentally address some of the more intractable diseases of our times.

## Bibliography

1. M. B. Pepys (1994) in *Santer's Immunologic Diseases* (M. M. Frank, K. F. Austen, H. N. Claman, and E. R. Unanue, eds.), Little, Brown & Co., Boston, pp. 637–655.
2. A. S. Cohen, T. Shirahama, and M. Skinner (1981) in *Electron Microscopy of Protein*, Vol. 3 (I. Harriss, ed.), Academic Press, London & New York, pp. 165–205.
3. G. G. Glenner, E. D. Eanes, and D. L. Page (1972) *J. Histochem. Cytochem.* **20**, 821–826.
4. G. G. Glenner (1980) *N. Eng. J. Med.* **303**, 1283–1292.
5. M. Sunde, L. C. Serpell, M. Bartlam, P. E. Fraser, M. B. Pepys, and C. C. F. Blake (1997) *J. Mol. Biol.* **273**, 729–739.
6. E. Scherzinger et al. (1997) *Cell* **90**, 549–558
7. E. D. Eanes and G. G. Glenner (1968) *J. Histochem. Cytochem.* **16**, 673–677.
8. L. Bonar, A. S. Cohen, and M. Skinner (1967) *Proc. Soc. Exp. Biol. Med.* **131**, 1373–1375.
9. J. H. Cooper (1976) in *Amyloidosis* (O. Weheliuss and A. Paternak, eds.), Academic Press, London & New York, pp. 61–68.
10. C. C. F. Blake and L. C. Serpell (1996) *Structure* **4**, 989–998.
11. C. Chothia (1973) *J. Mol. Biol.* **75**, 295–302.
12. L. C. Serpell, M. Sunde, P. E. Fraser, P. Luther, E. Morris, E. Sandgren, E. Lundgren, and C. C. F. Blake (1995) *J. Mol. Biol.* **254**, 113–118.
13. M. J. M. Saraiva (1995) *Hum. Mutat.* **5**, 191–196.
14. C. C. F. Blake, M. J. Geisow, S. J. Oatley, B. Rerat, and C. Rerat (1978) *J. Mol. Biol.* **121**, 339–356.
15. W. Colon and J. W. Kelly (1992) *Biochemistry* **31**, 8654–8660.
16. L. C. Serpell, G. Goldsteins, I. Dacklin, E. Lundgren, and C. C. F. Blake (1996) *Amyloid: Int. J. Exp. Clin. Invest.* **3**, 75–85.
17. M. B. Pepys et al. (1993) *Nature* **362**, 553–557.
18. D. R. Booth et al. (1997) *Nature* **385**, 187–793.
19. S. D. Hooke, S. E. Radford, and C. M. Dobson (1994) *Biochemistry* **33**, 5867–5876.
20. R. Riek, S. Hornemann, G. Wider, M. Billeter, R. Glockshuber, and K. Wüthrich (1996) *Nature* **382**, 180–182.
21. J. Safar, P. Roller, D. Gajdusek, and C. Gibbs (1994) *Biochemistry* **33**, 8375–8383
22. C. Soto, E. Castano, B. Frangione, and N. Inestrosa (1995) *J. Biol. Chem.* **270**, 3063–3067.

## Suggestions for Further Reading

23. G. R. Bock and J. A. Goode, eds. (1996) *The Nature and Origin of Amyloid Fibrils*, Ciba Foundation, John Wiley and Sons Ltd., Chichester, UK.
24. L. C. Serpell, M. Sunde, and C. C. F. Blake (1997) The molecular basis of amyloidosis. *Cell. Mol. Life Sci.* **53**, 871–887.

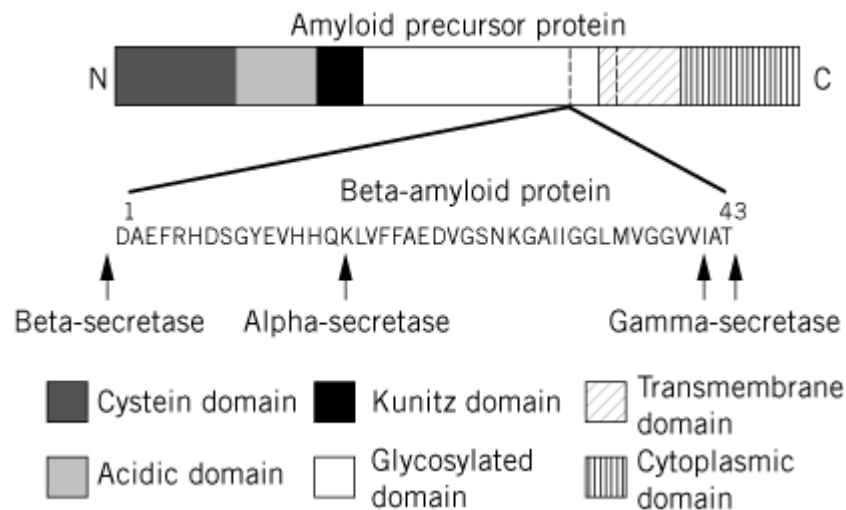
## Amyloid Precursor Protein

At the molecular level, the necessary, and probably sufficient, cause for the brain dysfunction in Alzheimer's disease is the deposition of aggregated forms of the [amyloid](#) b-protein (Ab or bA4), a **proteolytic** fragment of the amyloid precursor protein (APP), in the brain parenchyma and vascular

system (1). In the Western world, dementia is the most common neurological diagnosis and is the third leading cause of natural death, and Alzheimer's disease is by far the dominant cause of dementia (2). It has been estimated that there are 3 million patients in the United States alone with Alzheimer's disease (3). In addition, older patients with Down's syndrome ([trisomy 21](#)) tend to develop a neuropathology that is very similar to that seen in Alzheimer's patients (4). The disease is chronic, progressive, and at present untreatable. The progression of the disease begins with loss of short-term memory and disorientation, followed by complete impairment of memory, judgment, and reasoning. Firm diagnosis of the disease can be given only after postmortem examination of the brain, which is characterized by extensive loss of neurons and particular microscopic lesions, which include what are now called "neurofibrillar tangles" and "senile plaques." The neurofibrillar tangles are composed of [paired helical filaments](#) consisting of phosphorylated **tau protein**, which is normally associated with [microtubules](#). The senile plaques are composed mainly of Ab deposited in the brain parenchyma. Also present in the brain is congophilic amyloid angiopathy from the accumulation of the Ab peptide in the walls of blood vessels. The Ab protein in the senile plaques and in the vascular deposits is present in the form of [amyloid](#) (5, 6), long, stable fibrils resulting from the accumulation of precursor peptides that adopt a continuous [beta-sheet](#) structure. In this regard, Alzheimer's disease is related to Creutzfeldt–Jakob disease (CJD) and other spongiform encephalopathies, which are also characterized by the presence of amyloid deposits, in that case composed of the [prion](#) protein. All known risk factors for Alzheimer's disease appear to influence one or more of the following: (1) the concentration of Ab (7); (2) the amount of the longer, more amyloidogenic Ab chains (8, 9); or (3) the initiation of amyloid formation (10). Increases in any of these factors appear to increase the risk of Alzheimer's disease. Study of the origins of the b-amyloid protein from APP, its aggregation into amyloid fibrils, and the identification of genes involved in its inherited susceptibility are therefore central to understanding of, and devising therapies against, the scourge of Alzheimer's disease.

The amyloid precursor protein (APP), is a 110–130-kDa protein with the features of a transmembrane cell-surface **glycoprotein**. It is encoded by a gene localized on [chromosome 21](#) encoded by 18 exons (see [Introns, Exons](#)). A family of eight transmembrane glycoproteins is generated by [alternative splicing](#) of some of these exons (11). The amino acid sequence predicts that the single transmembrane domain of APP is near the C-terminus (Fig. 1) (12). The Ab sequence consists of a maximum 43-residue segment straddling the transmembrane and extracellular portions of the APP chain, and encoded by parts of exons 16 and 17. APP is processed in two distinct pathways: a major nonamyloidogenic route involving **proteolytic** cleavage within the region carrying the Ab segment to separate the extracellular and membrane-bound domains (13): and an amyloidogenic route, shown in Figure 1, leading to the release of the Ab segment (7). The cleaving of the Ab sequence from an internal site in the APP chain implies that two distinct proteolytic events are required to generate the N-terminus and the C-terminus of the Ab peptide. The proteolytic enzymes involved in processing APP are termed *secretases*: a-secretase for cleavage of APP within the Ab segment, and b- and g-secretases for cleavage of APP at the N- and C-terminal sides of the Ab segment, respectively. Soluble, C-truncated forms of APP are generated by two pathways: (1) by a-secretases cutting within the Ab sequence, thereby precluding the release of Ab (13); and (2) by b-secretases cleaving near the N-terminus of Ab producing C-terminal fragments containing the complete Ab (14). The generation of intact Ab chains suggests that g-secretases generate the C-terminus of b-amyloid from these C-terminal fragments of APP after release of the transmembrane domain from the lipid bilayer (7). The identity and location of the secretases have proved elusive, but clearly they are of prime importance in Alzheimer's disease, and targets for therapies.

**Figure 1.** Schematic diagram of the derivation of the 4-kDa b-amyloid protein (Ab) from the 110–130-kDa amyloid precursor protein (APP) by the action of the a-, b-, and g-secretase enzymes. The probable domain structure of APP has been simplified.



**Electron microscope** and X-ray [fiber diffraction](#) analyses, often using synthetic analogs or fragments of the Ab peptide, have revealed the molecular characteristics of the Ab amyloid fibrils. In the electron microscope, different forms of Ab amyloids are seen at different pH values, but the physiological form appears to be represented by a fibril about 90 Å (9 nm) diameter, composed of five or six parallel protofilaments 25–30 Å (2.5–3 nm) in diameter arranged around a hollow core (15). The X-ray fiber diffraction patterns of Ab amyloid (16, 17) show the usual 4.7 Å (0.47 nm) meridional reflection and 10 Å (1 nm) equatorial reflection that demonstrate that the molecular structure of amyloid consists of paired b-sheets running parallel to the fiber axis, with their constituent b-strands perpendicular to the fiber axis. This structure is very similar to that observed in other amyloids produced from a variety of proteins and peptides (see [Amyloid](#)). Some shorter Ab chains, such as residues 11–28, particularly when aligned in a strong magnetic field, form pseudocrystals (18), which give more detailed X-ray diffraction patterns that may be capable of yielding greater structural detail and hence begin to define the features that dispose them to form amyloid.

The discovery that Ab is produced and secreted by cells continuously under normal metabolic conditions, and is present in a soluble form in biological fluids (14, 19), suggests that other factors are involved in Ab amyloidosis. It has been proposed that the isolated Ab protein exists in two forms: a soluble form that may represent a normal host protein and a modified amyloidogenic form (20). These two forms have identical sequences and hence are likely to represent conformation isomers. The soluble form is easily degraded and appears to have an **alpha-helical/random coil** structure: [NMR](#) analysis of synthetic Ab 1–40 (21) has shown that residues 15–23 and 31–35 form a-helices, while the rest of the peptide is in a random-coil conformation, and no stable [tertiary structure](#) is present. The amyloidogenic form is more resistant to degradation, has a high b-sheet content, and forms the fibrillar aggregates found in brains with Alzheimer's disease. This behavior is very reminiscent of the normal and scrapie forms of the prion protein, which is also an agent for amyloid brain diseases (see [Prion](#), [Scrapie](#), and [Amyloid](#)). The existence of two structural forms for the Ab protein implies the possibility that other proteins may be involved in regulating their interconversion. It is known, for example, that in the rare early-onset (30–60 years of age) form of Alzheimer's disease, the disease has been linked to mutations on chromosomes 1 and 14, in addition to chromosome 21, which carries the APP gene. The genes on chromosomes 14 and 1 have been identified with presenilin 1 (22) and presenilin 2 (23, 24), respectively. The two presenilins have 65% amino acid sequence identity, and [hydrophobicity](#) plots predict seven transmembrane regions resembling a structural [membrane protein](#). The genetic variants of presenilin 1 allow the secretion of Ab peptides of longer length, and hence more amyloidogenic (25), possibly through a mechanism involving protein sorting and trafficking of APP.

The more common sporadic and familial late onset (>65 years) form of the disease has been linked with mutations on chromosome 19, subsequently narrowed to the  $\epsilon$ 4 allele of the **apolipoprotein E** (apoE) gene (26, 27). These studies have shown that carriers of the  $\epsilon$ 4 allele have an increased risk of developing Alzheimer's disease and that the inheritance of the  $\epsilon$ 4 allele correlates with an increased deposition of Ab amyloid in blood vessels and plaques. The apoE4 **isoform** differs from apoE3, the most common isoform, by an Arg/Cys substitution at position 112. The apoE protein is known to be one of the proteins associated with amyloid deposits (28). It has been found that apoE can form complexes with synthetic Ab analogues and to enhance amyloid fibril formation by Ab *in vitro* (29), possibly by a direct physical interaction. Frangione and his colleagues (30) have made the intriguing proposal that apoE can itself possibly form an amyloid-like structure (there is experimental evidence for this in its C-terminal domain), which may be able to induce other proteins such as Ab to misfold into a  $\beta$ -sheet structure that could allow them to be incorporated into a growing amyloid fibril. They term this process “conformational mimicry,” and certainly the most recent analyses of *in vivo* amyloid formation characterize it as a disease of protein misfolding, which is possibly **cooperative** (see [Amyloid](#)).

### Bibliography

1. A. I. Bush, K. Beyreuther, and C. L. Masters (1992) *Pharmacol. Ther.* **56**, 97–117.
2. A. Ott, M. M. B. Breteler, and F. van Harskamp (1995) *Br. Med. J.* **310**, 970–973.
3. R. E. Scully et al. (1993) *New Engl. J. Med.* **324**, 1255–1263.
4. K. E. Wisniewski, A. J. Dalton, D. R. Crapper-McLachlan, G. Y. Wen, and H. M. Wisniewski (1985) *Neurology* **35**, 957–961.
5. G. G. Glenner and C. W. Wong (1984) *Biochem. Biophys. Res. Comm.* **120**, 885–890.
6. C. L. Masters, G. Simms, N. A. Weinman, G. Multhaup, B. L. McDonald, and K. Beyreuther (1985) *Proc. Natl. Acad. Sci. USA* **82**, 4245–4249.
7. C. Haas and D. J. Selkoe (1993) *Cell* **75**, 1039–1042.
8. C. Hilbich, B. Kisters-Woike, J. Reed, C. L. Masters, and K. Beyreuther (1992) *J. Mol. Biol.* **228**, 460–473.
9. N. Suzuki et al. (1994) *Science* **264**, 1336–1340.
10. B. Hymen et al. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3586–3590.
11. R. Sandbrink, C. L. Masters, and K. Beyreuther (1994) *J. Biol. Chem.* **269**, 1510–1517.
12. J. Kang et al. (1987) *Nature* **325**, 733–736.
13. F. S. Esch et al. (1990) *Science* **248**, 1122–1124.
14. P. Seubert et al. (1993) *Nature* **361**, 260–263.
15. P. E. Fraser, J. Nguyen, W. Surewicz, and D. A. Kirschner (1991) *Biophys. J.* **60**, 1190–1201.
16. D. A. Kirschner, C. Abraham, and D. A. Selkoe (1986) *Proc. Natl. Acad. Sci. USA* **83**, 503–507.
17. P. E. Fraser et al. (1992) *Biochemistry* **31**, 10716–10723.
18. H. Inouye, P. E. Fraser, and D. A. Kirschner (1993) *Biophys. J.* **64**, 502–519.
19. M. Shoji et al. (1992) *Science* **258**, 126–129.
20. C. Soto, E. Castano, B. Frangione, and N. Inestrosa (1995) *J. Biol. Chem.* **270**, 3063–3067.
21. H. Sticht, P. Bayer, D. Willbold, S. Dames, C. Hilbich, K. Beyreuther, R. W. Frank, and P. Rosch (1995) *Eur. J. Biochem.* **233**, 293–298.
22. R. Sherrington et al. (1995) *Nature* **375**, 754–760.
23. E. Levy-Lahad et al. (1995) *Science* **269**, 973–977.
24. E. Rogaev et al. (1995) *Nature* **376**, 775–778.
25. M. Barinaga (1995) *Science* **268**, 1845–1846.
26. E. H. Corder et al. (1993) *Science* **261**, 921–923.



27. W. J. Strittmatter (1993) Proc. Natl. Acad. Sci. USA **90**, 1977–1980.
28. T. Wisniewski and B. Frangione (1992) Neurosci. Lett. **135**, 235–238.
29. J. Ma, A. Yee, H. B. Brewer, S. Das, and H. Potter (1994) Nature **372**, 92–94.
30. B. Frangione, E. M. Castaño, T. Wisniewski, J. Ghiso, F. Prelli, and R. Vidal (1996) in *The Nature and Origin of Amyloid Fibrils*, Ciba Foundation Symposium 199 (G. R. Bock and J. A. Goode, eds.), John Wiley and Sons Ltd., Chichester, UK, pp. 132–145.

### Suggestions for Further Reading

31. D. J. Selkoe (1994) Normal and abnormal biology of the beta-amyloid precursor protein, Annu. Rev. Neurosci. **17**, 489–517.
32. G. Evin, K. Beyreuther, and C. L. Masters (1994) Alzheimer's disease amyloid precursor protein, Amyloid: Int. J. Exp. Clin. Invest. **1**, 263–280.

## Analogy

Analogy is a similarity caused by factors other than a common genetic ancestry. In general, analogy is a similarity arising by [convergent evolution](#). Therefore, in contrast to [homology](#), analogy between two traits of interest is not based on the sharing of a common ancestor. Originally, the term “analogy” was used mostly for morphological characters. For example, the wings of bats and insects are said to be analogous organs, as they perform the same function in different species, even though they do not show a common underlying plan of structure. It is obvious that the wings of bats and insects were not derived evolutionarily from the same organ in a common ancestor.

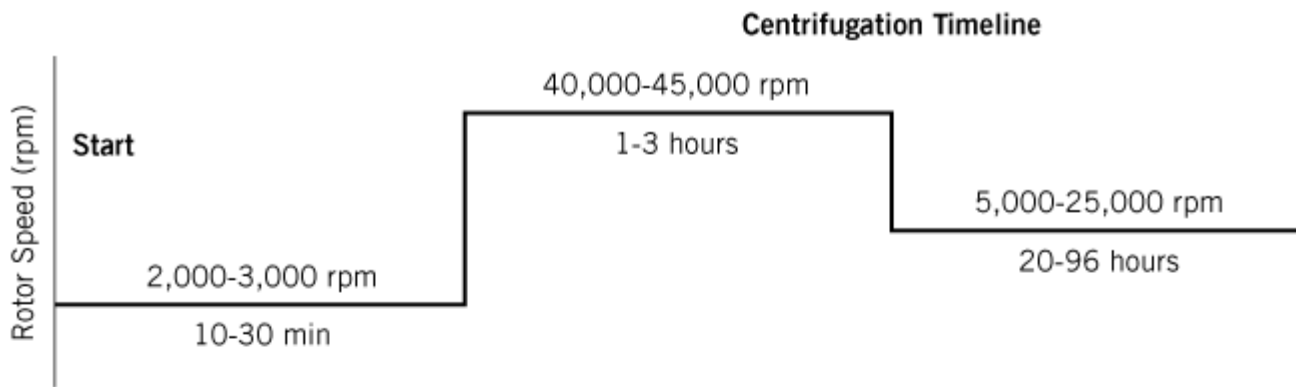
When the similarity between two different molecules is found at the level of **nucleotide sequences**, amino acid sequences, and the [tertiary structure](#) of [proteins](#), it is important to know whether the similarity is homologous or analogous. In particular, characteristic patterns of amino acids (“amino acid sequence motifs”) and the tertiary structures of protein molecules may exhibit analogy rather than homology. In contrast, any sequence similarities of nucleotides and amino acids are explained mostly by homology, because the probability that the two sequences of interest are similar by chance can be computed easily and is usually extremely small.

## Analytical Ultracentrifugation

Analytical ultracentrifugation experiments fall into two different categories, termed **sedimentation velocity** and **sedimentation equilibrium**. However, they are often employed sequentially on a given macromolecule under analysis (Fig. 1), and they yield related, complementary information about the solution behavior of a macromolecule that is of great interest for a variety of reasons. The need may be simple, for example, to determine whether a macromolecule aggregates when prepared or stored in various ways, e.g., refrigerated, frozen, or lyophilized. On a much more complex level and much more important in the context of modern protein biochemistry and molecular biology is to detect and measure the interactions of macromolecules with themselves or with other molecules and to relate this to their functional roles. For these studies, the precise analysis of their behavior in solution is

necessary. Finally, analytical ultracentrifugation plays an important role in medicine and industry. It is the reference method for analyzing the concentration-dependent solution behavior of [recombinant proteins](#) and other macromolecules employed in pharmaceutical/pharmacological research and in structure-based drug design ([1](#)).

**Figure 1.** Typical execution scheme for sedimentation velocity (SV) and sedimentation equilibrium (SE) runs.



### Run Goals

Initial phase  
Processing

SV run

SE run

**1)** Look for really big particles (aggregates)

**Yes?** Do a low speed velocity run for ~1 hour and proceed to post-run analysis

**No?** Proceed to high speed velocity run

**2)** Establish baseline absorbance scans at desired wavelengths for mass balance relationships

**1)** Obtain SV scans during run at desired wavelengths

**2)** Down-load raw data over network as run proceeds. Analyze for  $s$  and  $D$  values in near-real-time

**3)** Check for evidence of small or large molecule binding, using wavelength scans or interference optics

**4)** Determine  $M_{s,D}$  and alter final SE phase run speed for final equilibrium

**1)** Repeat cell scans every hr. and check for equilibrium

**2)** Setup SV and SE data for post-run processing averaging

**3)** Check for long-term stability to aggregation

**4)** Check for evidence of small or large molecule binding using wavelength scans or interference optics

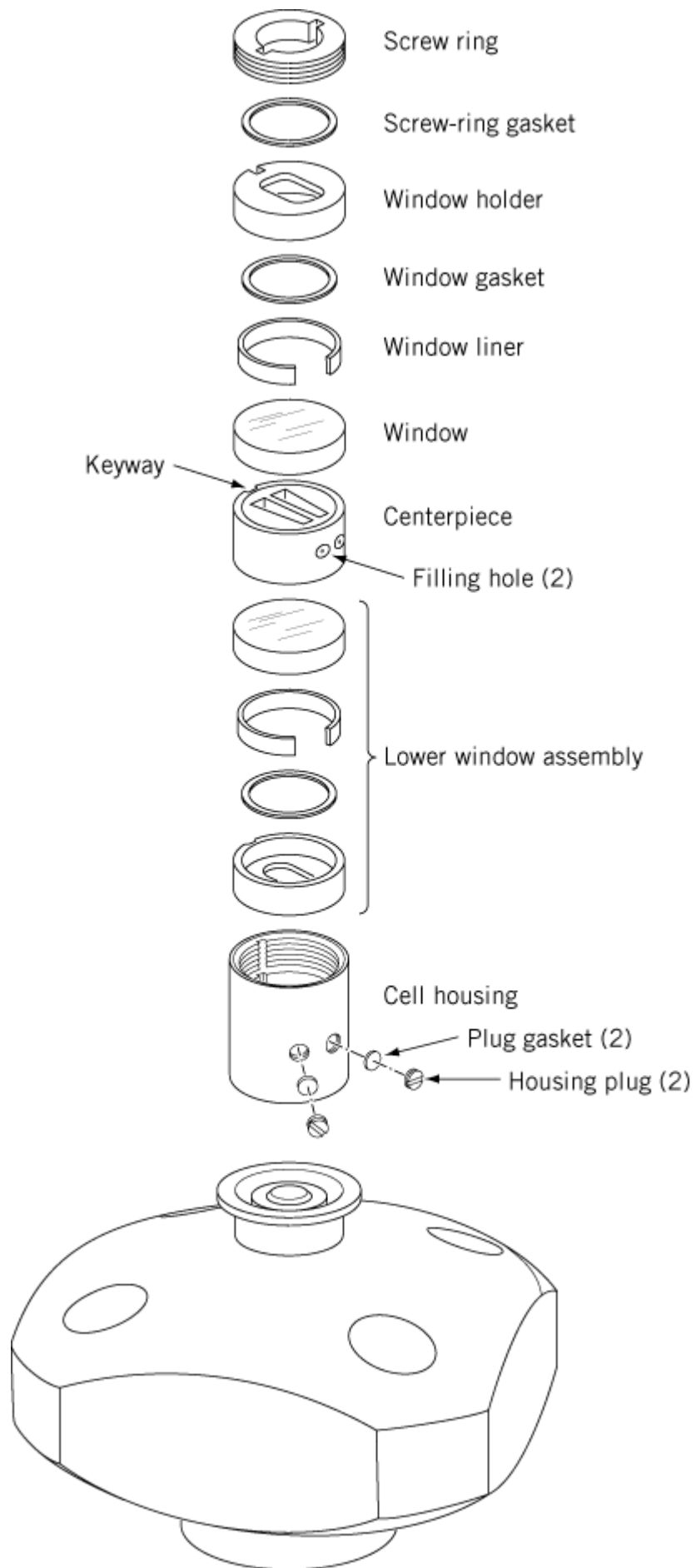
Analytical ultracentrifugation differs from preparative [centrifugation](#) in several important respects. First and foremost as the name implies, it is an *analytical* technique. If carefully applied, it can give hydrodynamic information that is precise, accurate, and unobtainable by any other single method. The hydrodynamic parameters (**molecular weights** and [sedimentation coefficients](#)) obtained from

analytical ultracentrifugation are absolute measurements because the technique is an absolute physical method in the strictest sense. It measures either the intensive or extensive physical properties of solutes directly and requires only knowledge of the rotational rate of the centrifuge rotor, the sample temperature, and the solution density and viscosity (which are extensive properties of the solute–solvent system).

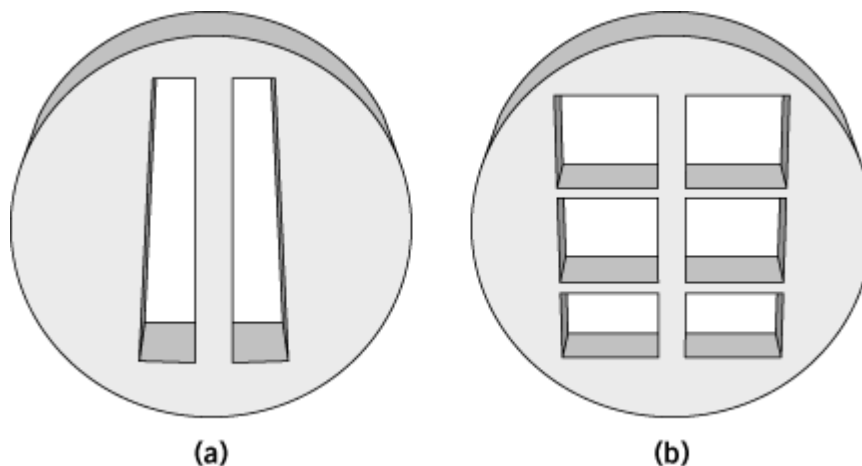
Similarly, calculations of molecular properties from ultracentrifugation data are based on direct, first-principle analyses of solute behavior in solution. Unlike indirect methods, such as static [light-scattering](#), the calculations in analyzing data from the analytical ultracentrifuge do not rely on comparison to a known standard. For example, in static light scattering the scattering behavior of a low molecular weight organic molecule (toluene, for example) is required to calibrate the instrument. [Dynamic light scattering](#) provides absolute [diffusion](#) coefficients, but the molecular weights estimated from the diffusion coefficients are indirect and depend upon the solution density, viscosity, etc. and on assumptions about molecular dimensions.

Another major difference between preparative centrifugation and analytical ultracentrifugation is that the latter is usually applied to highly purified solute samples. Although samples must be pure, analytical analyses usually require only small amounts of sample material, in the range of 10 to 1000 µg per sample. Finally, although analytical ultracentrifugation is inherently a nondestructive technique, it is not routine to recover or reuse samples that have been analyzed because analytical ultracentrifuge cells and rotors (Figs. [2](#) and [3](#)), especially those that are most useful for analyzing numerous samples, are not designed for facile sample recovery. Although analytical ultracentrifugation was at one time a useful tool for evaluating the purity of certain biological macromolecules, its use as a simple tool for purity analysis has been superseded by a variety of other bioanalytical techniques. For example, inexpensive [SDS-PAGE](#) is routinely the technique of choice for analyzing protein purity during a protein separation protocol and for estimating the molecular weights of polypeptide chains present.

**Figure 2.** A four-hole analytical ultracentrifuge rotor and a vertically exploded view of the cell assembly with a double-sector, long-channel centerpiece. Sample and reference solutions for analysis (~400 to 500 µl) are introduced into the sealed centerpiece with an appropriate syringe system. Then the cell is placed in the rotor and subjected to centrifugation. Depending upon the analyses performed, it is possible after centrifugation to recover a portion of the sample using the same syringe system (figure courtesy of Beckman Instruments, by permission).



**Figure 3.** Illustrations of centerpieces used in analytical ultracentrifugation. Left, a long-column, double-sector centerpiece used for sedimentation velocity and equilibrium studies that require large data sets. Right, the multichannel centerpiece (or “Yphantis” centerpiece after its inventor, Dr. David Yphantis) used primarily for sedimentation equilibrium studies where analysis of a larger number of samples is desired.



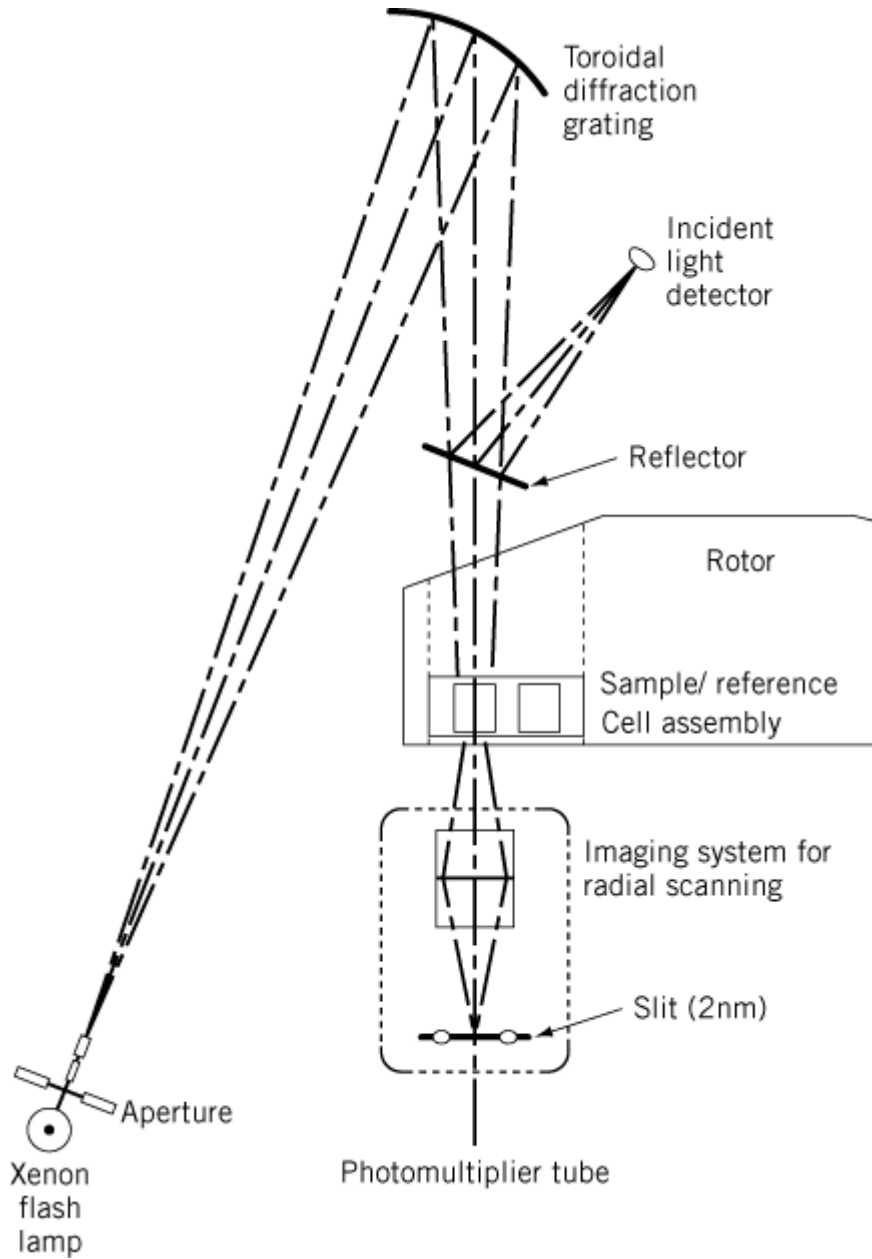
The present (and likely future) value of analytical ultracentrifugation lies in its unmatched ability to provide direct information in two primary areas: (1) the extent to which an otherwise chemically pure sample of a macromolecular solute exhibits monodispersity or polydispersity and (2) the extent to which different chemically pure molecules interact. Because these two questions are fundamental to all areas of modern cellular and molecular biology and biochemistry, analytical ultracentrifugation remains a unique and extremely powerful biophysical tool.

### 1. The Modern Analytical Ultracentrifuge

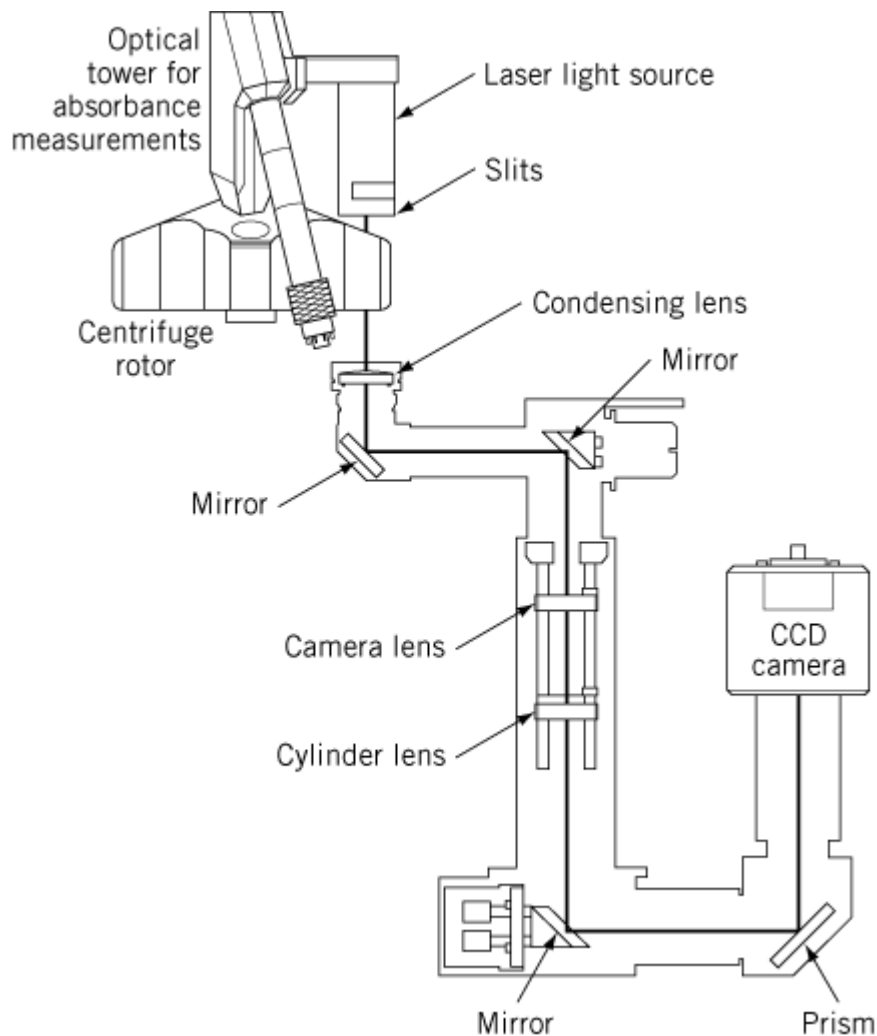
The technique requires only small sample volumes,  $\sim 20 \mu\text{l}$  to  $\sim 500 \mu\text{l}$  depending upon the type of centerpiece employed (Figs. 2 and 3). A computer-controlled optical system detects the exact extent and position of the radial displacement of solute molecules at various times during the centrifugation run. Optical detection is accomplished by directly measuring the concentration of the solute as a function of radial position in the sample cell, either by its absorbance or by the change in refractive index, using Rayleigh interference optics. Like dual-beam optical [spectroscopy](#), the technique is also differential in that it employs a sample solution to be measured and a reference solution used to “blank out” or “subtract” a signal background. The only modern, commercially available analytical ultracentrifuges are the XLA/XLI models manufactured by Beckman Instruments, which operate at speeds ranging from  $\sim 2,000$  to  $\sim 50,000$  rpm. The precise speed or speed range of operation for an experiment is determined by the type of measurement(s) to be made and depends on the size of the molecules being studied and on the type of sample cell centerpieces, centrifuge rotor, and the data collection system used.

The XLA/XLI analytical ultracentrifuges offer either absorbance-based optical detection (XLA, Fig. 4) or both absorbance and interference-based optical detection (XLI, Fig. 5). In addition to measuring the molecular weights and sedimentation coefficients of individual molecules, their distributions in samples can also be measured. The analytical ultracentrifuge can also be used for direct determination of **ligand-binding** constants and stoichiometries. Further, when combined with other physical measurements, these instruments are useful for determining the [diffusion](#) coefficients and molecular shapes of macromolecules. Thus, changes in shape related to ligand binding or solvent perturbations can also be evaluated.

**Figure 4.** Absorbance-based optical system of the XLA analytical ultracentrifuge from Beckman Instruments. Light from a high intensity flash lamp passes along an optical tube to a computer-controlled toroidal diffraction grating, which selects the wavelength(s) of observation. Light at the selected wavelength, within the range ~200–800 nm, is focused on the sample cell in the rotor. Light passing through the sample is recorded by the radially tracking photomultiplier. Then the data are analyzed using specialized software on an associated computer (figure courtesy of Beckman Instruments, by permission).



**Figure 5.** Interference and/or absorbance-based optical system of the XLI ultracentrifuge from Beckman Instruments. The absorbance-based optical system is outlined in Fig. 4. For interference measurements, a laser light source is attached to the absorbance optical tower, and an additional optical path and detector system is employed (lower right quadrant of figure). The interference system uses Rayleigh interference optics and a computer-controlled CCD camera detection system. Data are analyzed using specialized software on an associated computer. Using an eight-hole centrifuge rotor and centerpieces that hold three sample and reference solutions each, it is possible to analyze 21 samples in a single ultracentrifuge run (figure courtesy of Beckman Instruments, by permission).



### 1.1. Choosing an Optical System for Measurements

Although it is often useful to employ both absorbance and optical interference measurements when examining a sample, some considerations make one type of measurement preferable to the other. The possible choices for various types of samples are outlined in Table 1. In addition to the general preferences for detecting various analytes, other important factors influence the choice of detection methods. The relative sensitivities and dynamic ranges of absorbance versus interference detection can be a consideration. The typical useful range for absorbance measurements in a [sedimentation equilibrium centrifugation](#) experiment for an “average” protein with an absorbance at 280 nm of  $1.0 \text{ (mg/ml)}^{-1} \text{ (cm)}^{-1}$  is in the range of  $\sim 10 \text{ }\mu\text{g/ml}$  to  $\sim 4 \text{ mg/ml}$ , when the measurements employ centerpieces with both standard (12 mm) and reduced (3 mm) path lengths. In **sedimentation velocity experiments** using absorbance detection at 280 nm, the practical low end of the concentration range is about 50 to 100  $\mu\text{g/ml}$ . Sensitivity for proteins can be increased by using light with wavelengths in the range of 215 to 240 nm, where the protein peptide backbone absorbs and the absorbance of a protein is typically 3- to 10-fold greater, although this is technically more demanding. Alternatively, the protein may have a natural cofactor chromophore with a large extinction coefficient, or [chemical modification](#) can be used to introduce such a chromophore. Optical interference measurements have roughly the same absolute dynamic range as absorbance measurements, but they are routinely more useful in the mid-to-high concentration range, i.e.,  $\sim 0.1$  to  $\sim 10 \text{ mg/ml}$ .

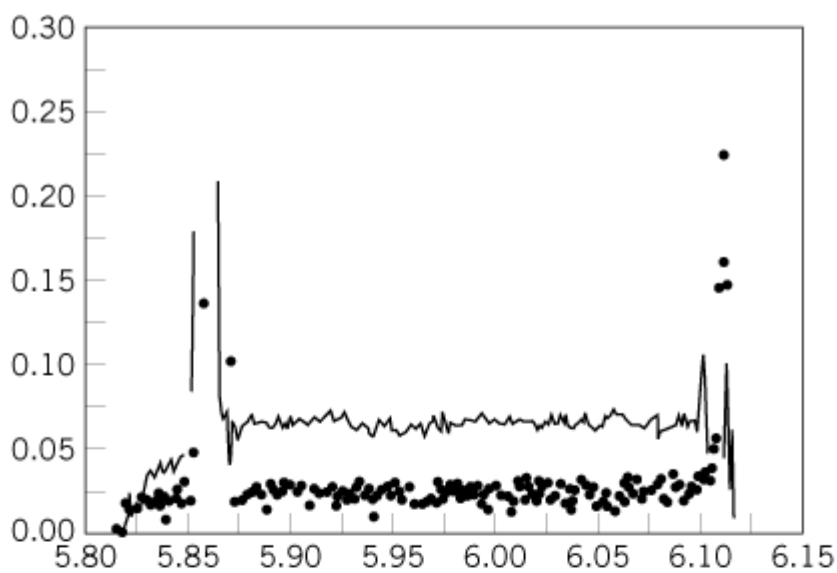


**Table 1. General Preferences in Analytical Ultracentrifugation for Detecting Different Analytes**

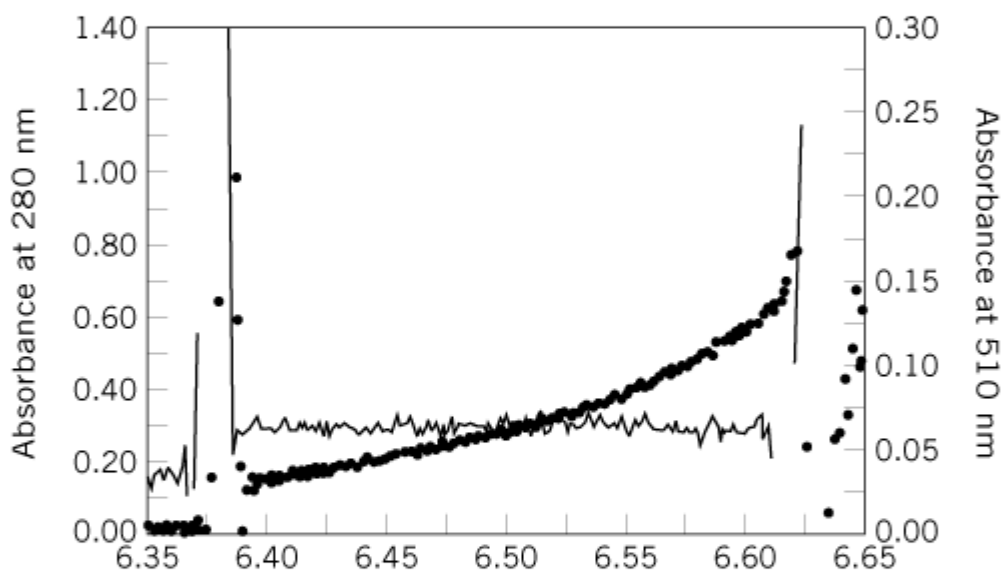
| Analyte          | Absorbance Detection   | Interference Detection                               |
|------------------|--|--|
| Carbohydrates    | Usually not useful unless solute has appreciable UV or visible absorbance properties | Usually the method of choice                         |
| Nucleic acids    | Follow absorbance at ~260 nm   | Possible, can be used at high analyte concentrations |
| Organic polymers | Usually not useful unless solute has appreciable UV or visible absorbance properties | Usually the method of choice                         |
| Polysaccharides  | Usually not useful unless solute has appreciable UV or visible absorbance properties | Usually the method of choice                         |
| Proteins         | Useful if protein contains Trp and/or Tyr residues, so that it absorbs at 280 nm     | Useful if protein lacks Trp and Tyr residues         |

The selectivity of the detection technique can also be a factor in performing an experiment. The use of this approach is illustrated in Fig. 6 for examining the potential for interactions between the Alzheimer's disease beta-amyloid peptide and the serum C1q component of the [complement system](#), which other studies had suggested is a trigger for activating the complement cascade and for potential neuronal cell loss. In this work, a 5-kDa peptide was labeled with the chromophore *fluorescein*. Binding of the labeled peptide to a ~150 kDa anti-fluorescein [antibody](#) was used as a positive control to examine the potential for b-amyloid peptide to bind to C1q. The labeled peptide binds to the antibody as expected, but it does not bind to C1q.

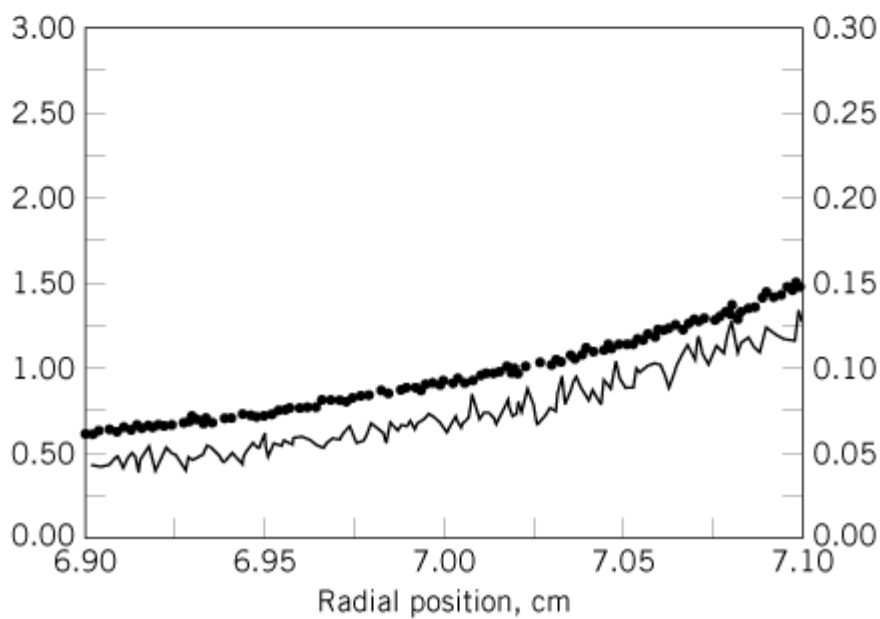
**Figure 6.** Binding of a xanthene-labeled Alzheimer's b-amyloid peptide to an anti-fluorescein antibody but not to serum complement C1q, as analyzed at sedimentation equilibrium at low speed (5000 rpm). The fluorescein analog, xanthene, absorbs light at 510 nm (solid line), and the proteins absorb light at 280 nm (dotted line). Panel (a) shows the radial distributions of labeled peptide alone at 280 and 510 nm (no sedimentation). Panel (b) shows the labeled peptide sedimented in the presence of C1q, C1q tended to sediment, but the peptide did not. Panel (c) shows the control experiment where the labeled peptide was sedimented in the presence of an anti-fluorescein antibody. The peptide binds to the antibody and sediments with a molecular weight equivalent to that of the antibody. These results demonstrate that the monomeric form of the amyloid peptide implicated in Alzheimer's disease is unlikely to bind to complement C1q and thus induce complement-mediated neural cell death (from Ref. 8, with permission).



(a)



(b)



(c)

In analytical ultracentrifugation, another important consideration is sample and reference solution buffer matching. If the solutions are poorly matched, the resulting data are usually poor, even with the extensive signal averaging employed with advanced detectors and software.

Interference detection measurements are based on the change in the refractive index of the solution caused by a solute, which is a function of its concentration. For these measurements, the system records the differential change in the refractive index of the sample versus the reference cell, as a function of radial position. The refractive index of a solution is a function of the mass of all solutes, not of their absorbance or molarity, and it is not possible to distinguish between different solutes. In fact, buffers and salts routinely contribute most of the change in refractive index because they are present in much higher concentrations than the protein. Therefore, interference measurements are extremely sensitive to small mismatches in refractive index between the buffers of the sample and reference solutions. As a result, samples analyzed by interference measurements routinely require exhaustive [dialysis](#). With appropriate sample/reference solution preparation, it is possible to achieve protein concentrations of >100 mg/ml routinely. The only effective upper limit is the intrinsic solubility of the protein itself under the conditions of analysis (rotor velocity, temperature, pH, etc., which can be easily changed). For both absorbance and interference measurements, the lower limit on size is a molecular weight of ~200 Da, and the upper limit is effectively about  $30 \times 10^6$  Da.

Absorbance optics permit the differential absorbance of two molecules to be employed to distinguish their sedimentation properties (as in Fig. [6](#)). This is not possible with interference optics, however, so it is not possible to distinguish between two different solute molecules with similar masses.

Another feature of interference measurements is their use in a differential mode. Because the interference optical system is double-beam, differential centrifugation experiments can be performed where the protein is present in both cells and a small ligand only in one. Consequently, the signal from the protein is blanked out, leaving only the signal from the ligand. A ligand that binds to the protein (even weakly) is displaced along the radial axis of centrifugal force in the centrifuge cell. Because the protein signal is blanked out, the only change in refractive index observed along the radial (force) axis is caused by the change in radial distribution of the ligand ([2](#), [3](#)).

Absorbance measurements generally have less stringent referencing requirements. In practice it is necessary only that the reference buffer has an absorbance at the wavelength of analysis at which it is electronically possible for the instrument to perform the desired subtraction. Therefore, absorbance measurements are relatively more forgiving for buffer matching in sample/reference preparation. They are not without other intrinsic problems, however, in particular those related to the absorbance of buffers and solute extinction coefficients. Most biologically useful buffers and common solution additives have some characteristic absorbance in the wavelengths used to analyze proteins in solution. Because of the high total sample/reference solution absorbance, it may be difficult or impossible to make the requisite difference absorbance measurements. However, under these conditions protein samples that have been carefully prepared by dialysis but are optically opaque at UV wavelengths, may be routinely analyzed by interference rather than absorbance detection (as above). Finally, one interesting drawback to absorbance measurements, which is routinely neglected in sedimentation equilibrium analyses (in particular), is that the analyte solution extinction coefficient(s) may vary with the composition of the solution, so that its absorbance can change during a centrifugation run. Homo- or heteroassociation processes may alter the protein extinction coefficients, which contributes nonlinear perturbations to the total

absorbance. Thus careful analyses require that the radial distributions of molecules undergoing such associations be measured at several wavelengths. Interference measurements, in comparison, are insensitive to these alterations in solute extinction coefficient.

## 1.2. Monodispersity, Polydispersity, Paucidispersity, and Nonideality

A solute is said to be monodisperse if it exhibits behavior characteristic of a single species in solution. If the solute undergoes indefinite self-association under the solution conditions examined, it is said to exhibit *polydispersity*. If a sample of solute molecules is composed of a mixture of a limited number of species with different chemical compositions, its solution behavior is *polydisperse*, for example, with respect to its molecular weight distribution, but the sample is termed *paucidisperse*. The differences between various types of solute behavior in solution are outlined in Table 2. It is often useful to define molecular weight distributions in terms of the molecular weight averages  $M_n$ ,  $M_w$ , and  $M_z$ , which are the number-, weight-, and z-averages, respectively (Table 3). They are valuable for estimating sample homogeneity. A pure solute gives identical values of  $M_n$ ,  $M_w$ , and  $M_z$ . Different methods of measurement give different averages.

**Table 2. Summary of Types of Solute Hydrodynamic Behavior**

| Hydrodynamic Characteristic | Example   | Sedimentation Velocity Behavior   | Sedimentation Equilibrium Behavior  | Consequences   |
|-----------------------------|---|---|---|--|
| <b>Monodisperse</b>         | Solute is pure, a single spherelike chemical entity that does not undergo self-interaction    | Velocity profiles fit a single-term Lamm equation (see <a href="#">Sedimentation Velocity Centrifugation</a> ), and the single sedimenting boundary gives unique values of $s$ and $D$ during run | Data fit a single exponential term at all rotor speeds and solute concentrations        | Molecular weight, $s$ , and $D$ have minimal or no dependencies on concentration; solute solution behavior is “ideal” or nearly so |
| <b>Paucidisperse</b>        | Sample comprises a small number of chemically distinct solutes, which may or may not interact | Velocity profiles fit a multiple-term Lamm equation, and sedimenting boundaries give multiple values of $s$ and $D$ during run<br><br>It may be possible to resolve separate sedimenting $g^*$    | Data do not fit a single exponential term at all rotor speeds and solute concentrations | Molecular weight, $s$ , and $D$ depend on concentration<br><br>Chances for extracting binding constants for a                      |

|                     |  |  |  |  |
|---------------------|--|--|--|--|
|                     |  | (s) distributions  | macromolecular assembly are usually severely diminished  |  |
| <b>Polydisperse</b> | Sample comprises multiple solutes, which may or may not interact                             | Velocity profiles fit a multiple-term Lamm equation, and sedimenting boundaries give multiple values of $s$ and $D$ during run   | Data do not fit a single exponential term at all rotor speeds and solute concentrations  | Molecular weight, $s$ , and $D$ depend on concentration, solute solution behavior is usually termed “nonideal” |
|                     | Sample is pure, a single spherelike chemical entity that aggregates                          | It may be possible to resolve separate sedimenting $g^*$ ( $s$ ) distributions   | Depending upon nature of self-association, it may be possible to extract binding constants for a macromolecular assembly process |  |
| <b>Nonideal</b>     | Sample comprises a single, highly asymmetric chemical entity that undergoes self-interaction | Velocity profiles fit a multiple-term Lamm equation, and sedimenting boundaries give multiple values of $s$ and $D$ during run<br><br>It may be possible to resolve separate sedimenting $g^*$ ( $s$ ) distributions | Data exhibit skewed behavior and cannot be fit by conventional analyses  | $s$ and $D$ depend on concentration, solute solution behavior is “nonideal”                                    |

**Table 3. Molecular Weight Averages**

| Type                   | Equation <sup>a</sup>                 | Methods to Determine                        |
|------------------------|---------------------------------------|---|
| $M_n$ , number-average | $\frac{\sum_i C_i}{\sum_i (C_i/M_i)}$ | Osmotic pressure, sedimentation equilibrium |

$$M_w, \text{ weight-average } \frac{\sum_i C_i M_i}{\sum_i C_i} \quad \text{Sedimentation equilibrium; light scattering}$$

$$M_z, \text{ z-average } \frac{\sum_i C_i M_i^2}{\sum_i C_i M_i} \quad \text{Sedimentation equilibrium}$$

<sup>a</sup>  $C_i$  is the concentration, and  $M_i$  is the molecular weight of the  $i$ th species in weight per unit volume.

Nonideality is typically observed at high concentrations of a solute and is caused by the effects of [excluded volume](#) (versus actual molecular volume) on solute hydrodynamic behavior. The effects of nonideality are small for spherical molecules, because their excluded and actual volumes are essentially equivalent. With such molecules, their hydrodynamic behavior varies in direct proportion to the size of the molecular sphere. When molecules are rodlike, nonideality becomes more pronounced. As the concentration increases, the effect is nonlinear and it is significantly worse at greater ratios of molecular length to width. Intermolecular interactions are another cause of nonideality. For example, electrostatic charge–charge repulsion can be significant at high solute concentrations and low solution ionic strengths. Experimentally, such charge–charge effects are usually reduced by increasing the solution ionic strength. The role of nonideality in sample behavior is identified by performing a series of studies at different solute concentrations and extrapolating the measured hydrodynamic parameters to infinite dilution.

### 1.3. Other Physical Measurements Necessary to Analyze Ultracentrifuge Data

Calculating results from sedimentation equilibrium and velocity experiments requires several additional physical measurements of the solute and of the solutions employed. To calculate the anhydrous weight-average molecular weight, the density of the solvent and the [partial specific volume](#) of the solute must be known.

There are typically two ways to obtain these values. They can be measured with an ultrasonic density meter or pycnometer or estimated from various tables if the specific buffer and solute chemical compositions are known. For a protein, the approximate partial specific volume is calculated from its amino acid composition ([4](#), [5](#)). The density of a buffer of exact composition can be estimated from an appropriate reference table (for example, from the *CRC Handbook of Chemistry and Physics*). Of the two direct methods for determining solvent densities and solute partial specific volumes, ultrasonic densitometry is preferred because it requires less solute for measurements. It is good experimental practice to measure both the solvent density and the partial specific volume of a solute sample of interest. In the case of proteins, this is particularly true if the studies analyze the tendencies of the proteins to undergo association or conformational changes when bound to a ligand. Such interactions are often accompanied by significant changes in the volume of the protein that are unaccounted for if the specific volumes of the individual amino acids comprising the protein are used in the calculation ([6](#), [7](#)). Additionally, if the protein contains carbohydrate (a **glycoprotein**, for example) or bound lipid (a **lipoprotein**), estimates from reference tables can, at best, be charitably characterized as a guess.

The final physical measurement necessary is the measurement of the solution viscosity. This value can be estimated from tables ([5](#)) (as for buffer density), but again it is good experimental practice to measure the solution viscosity directly with a simple, commercially available viscometer.

### Bibliography

1. S. W. Snyder, R. P. Edalji, F. G. Lindh, K. A. Walter, L. Solomon, S. Pratt, K. Steffy, and T. F. Holzman (1996) *J. Protein Chem.* **15**, 763–774.
2. I. Z. Steinberg and H. K. Schachman (1966) *Biochemistry* **12**, 3728–3747.

3. T. M. Lohman, C. G. Wensley, J. Cina, J. R. R. Burgess, and M. T. Record (1980) *Biochemistry* **19**, 3516–3522.
4. E. J. Cohn and J. T. Edsall (1943) In *Proteins, Amino Acids, and Peptides as Ions and Dipolar Ions* (ACS Monograph series), Reinhold Publishing, New York, pp. 155–176.
5. H. Durchschlag (1986) In *Thermodynamic Data for Biochemistry and Biotechnology* (H. J. Hinz, ed.), Springer-Verlag, New York, pp. 45–128.
6. H. Durchschlag and R. Jaenicke (1982) *Int. J. Biol. Macromol.* **5**, 143–148.
7. H. Durchschlag and R. Jaenicke (1982) *Biochem. Biophys. Res. Commun.* **108**, 1074–1079.
8. S. W. Snyder, G. T. Wang, L. Barrett, U. S. Lador, D. Casuto, C. M. Lee, G. A. Krafft, R. B. Holzman, and T. F. Holzman (1994) *Exp. Neurol.* **128**, 136–142.

### Suggestions for Further Reading

9. K. E. Van Holde (1971) *Physical Biochemistry*, Prentice–Hall, Inc., Englewood Cliffs, NJ, pp. 70–121. A short, highly readable and concisely illustrated treatment of the mechanical model of centrifugation for the biologist.
10. D. Eisenberg and D. Crothers (1979) "Physical Chemistry with Applications to the Life Sciences", Benjamin–Cummings, Menlo Park, CA, pp. 701–745. An excellent introductory physical chemistry text with a good review of fundamental flow equations and a number of useful examples.
11. C. Tanford (1961) *Physical Chemistry of Macromolecules*, Wiley, New York, pp. 317–456. A classic and still relevant text on the nature of transport processes as applied to biological systems.
12. T. M. Schuster and T. M. Laue (1994) "Modern Analytical Ultracentrifugation: Acquisition and Interpretation of Data for Biological and Synthetic Polymer Systems", Birkhäuser, Boston. The most recent compendium of modern data analytical methods and applications pertaining to the new XLA analytical ultracentrifuge.

## Androgenesis

Androgenesis is the mode of **reproduction** of **eukaryotes** where both sets of **chromosomes** of the offspring are of paternal origin. In higher **plants**, androgenesis can be obtained either *in vitro* with anther cultures or with isolated pollen grains. In animals, androgenesis can be obtained only experimentally. A technique similar to that described in the entry **Digynism** can be used, that of pronuclear transplantation: the male **pronucleus** is removed from a fertilized egg and injected into another egg, from which the female pronucleus is then removed, so that now the zygote contains two paternally derived sets of chromosomes.

Rarely, diandric embryos arise spontaneously when a fertilized egg loses its female pronucleus. These embryos are not viable. In humans, their implantation results in a uterine tumor called a hydatiform mole. In spontaneous abortions, diandric embryos are a minority among triploid conceptuses (1).

Diandric embryos are useful for studies of genetic **imprinting**.

## Bibliography

1. D. E. McFadden and J. T. Pantzar (1996) *Hum. Pathol.* **27**, 1018–1020.

## Aneuploidy

An aneuploid cell has one or more complete [chromosomes](#) either in excess or less than the normal [haploid](#), [diploid](#), or [polyploid](#) number characteristic of the species from which the cell derives. In other words, the chromosome number is not a multiple of the haploid cell number of chromosome.

Fluorescence [in situ hybridization](#) (FISH) with chromosome-specific probes has been successfully applied to rapidly detect numerical aberrations in metaphase and interphase amniotic cells ([1](#)). Aneuploidy arises primarily by the process of **nondisjunction**, probably caused by nonrandom premature [centromere](#) division in the first meiotic division of maternal **meiosis**. This varies among chromosomes, however, and a significant proportion of paternal and/or meiosis II errors have been described.

The presence of [kinetochores](#) or kinetochore proteins (detected immunochemically) in the **micronuclei** of binucleated cells indicates a cell with a high probability for aneuploidy following **cytokinesis** ([2](#), [3](#)). The relationship between [Y-Chromosome](#) aneuploidy in male humans, the number of micronuclei, the status of kinetochore proteins, and aging has been reviewed ([4](#)).

## Bibliography

1. W. L. Kuo et al. (1991) *Am. J. Hum. Genet.* **49**, 112–119.
2. D. A. Eastmond and J. D. Tucker (1989) *Mutat. Res.* **224**, 517–525.
3. B. K. Vig, H. J. Yoo, and D. Schiffmann (1991) *Mutagenesis* **5**, 361–367.
4. J. Nath, J. D. Tucker, and J. C. Hando (1995) *Chromosoma* **103**, 725–731.

## Suggestion for Further Reading

5. D. K. Griffin (1996) The incidence, origin, and etiology of aneuploidy. *Int. Rev. Cytol.* **167**, 263–296. An exhaustive review.

## Angiogenin

Angiogenin (Ang) is one of a group of proteins that are potent inducers of angiogenesis, the process by which new blood vessels are formed (see [Epidermal Growth Factor](#), [Fibroblast Growth Factors](#), [Transforming Growth Factors](#), [Tumor Necrosis Factor](#)). It was first isolated from medium conditioned by human adenocarcinoma (HT-29) cells based on the premise that tumor cells must secrete angiogenic factors in order to attract blood vessels so that they can grow. Ang was identified by its ability to stimulate angiogenesis on the chorioallantois, the outer membrane that surrounds the embryo in fertilized chicken eggs. Subsequently, it was shown to be present in normal human plasma, bovine milk, and bovine and mouse serum. It is a 14-kDa basic protein that has 33%

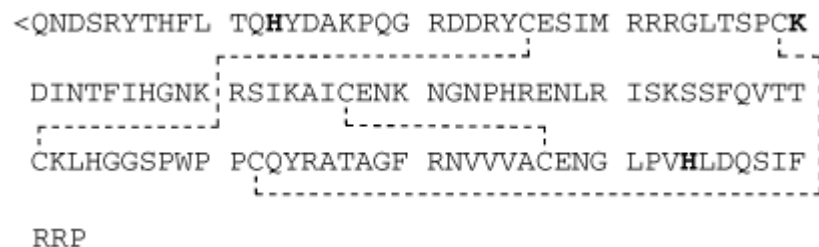


sequence identity to bovine pancreatic ribonuclease A and is one of several members of the RNase A superfamily of proteins that have unusual biological activities. It is the only angiogenic protein that is a ribonuclease, and the only ribonuclease that is angiogenic. It is an important protein that may be involved in wound healing, is critical for the growth of solid tumors, and has become a target for the treatment of metastatic cancer and other angiogenesis-related diseases. For recent reviews on Ang see Refs. [1](#) and [2](#).

## 1. Protein Chemistry

Human Ang is a single [polypeptide chain](#) of 123 [amino-acid](#) residues and three [disulfide bonds](#) (Fig. [1](#)). All of the main components of the catalytic site (see [Active Site](#)) of RNase A are present in Ang, and thus it is not surprising that it has enzymatic activity toward RNA and dinucleotides and which, like that of RNase A, is limited to cleavage after pyrimidines. What is surprising is that this activity is four to six orders of magnitude less than that of RNase A. Nevertheless, it is essential for angiogenic activity. Part of the reason for such low activity was revealed by [X-ray crystallography](#) (Fig. [2](#)). The structure of Ang is closely similar to that of RNase A, but with two major differences. On the one hand, the pyrimidine-binding ( $B_1$ ) component of the catalytic active site was found to be occluded by the side chain of Gln117. On the other, the second ( $B_2$ ) component of this site is structurally quite different from that in RNase A. In the latter, it is part of an eight-residue disulfide loop that is found in all of the 55 known mammalian members of the ribonuclease superfamily but missing in all of the nine known Ang. Instead, it is replaced by a sequence of residues that interacts with a cell-surface binding protein (see text below). Despite these likely impediments to substrate binding, Ang does catalyze RNA hydrolysis, albeit very weakly.

**Figure 1.** Amino acid sequence of human angiogenin. The active site histidine (H) and lysine (K) residues are in bold, and the disulfide bridges between cysteine (C) residues are indicated by dotted lines.



**Figure 2.** The polypeptide fold for human angiogenin drawn with the program MOLSCRIPT. (From Ref. [9](#), with permission.)



It has been suggested that through [evolution](#) Ang was endowed with a unique structure that attenuates its catalytic potency toward RNA in general, but that undergoes a conformational change when Ang is brought together with its specific, yet unidentified, substrate ([3](#)). Whether this structural rearrangement activates Ang or merely ensures its specificity is unknown. It may be that the very low activity of Ang is perfectly adequate for carrying out its unique biological function. It should be noted, however, that Ang is an effective inhibitor of cell-free protein synthesis by virtue of its ability to cleave 18S rRNA when this substrate is present in the intact [ribosome](#), and in fact is an even more effective inhibitor than RNase A. When isolated 18S RNA is used as substrate, the activity of Ang is markedly decreased relative to RNase A. This suggests that the activity of Ang depends on the environment of its substrate.

## 2. Molecular Biology

The human *Ang* gene, localized to [chromosome](#) band 14q11, is present as a single copy per haploid [genome](#), with no **introns** in the protein-coding region ([1](#), [2](#)). The gene has been cloned and expressed in both *Escherichia coli* and BHK cells, and well over 40 Ang variants have been prepared in order to explore various aspects of its structure–function relationships. Thus, [site-directed mutagenesis](#) of His13, His114, or Lys40 (the equivalents of the catalytic residues His12, His119, and Lys41 in RNase A) abolishes *both* the ribonucleolytic and the angiogenic activity of Ang. Although angiogenically inactive, these variants block the angiogenic activity of native angiogenin. By contrast, an Ang variant in which residues 58 to 70 were replaced by residues 59 to 73 of RNase A had 300- to 600-fold more activity toward RNA substrates but was angiogenically inactive. Thus, ribonuclease activity is essential but not sufficient for angiogenic activity. Moreover, this variant,

unlike the active site variants, was unable to block the angiogenic activity of native Ang, suggesting that the mutated region of the protein is involved in cell binding, probably interacting with a cell-surface binding protein. Recall that in RNase A this region of the molecule constitutes the B<sub>2</sub> component of the substrate binding site. In Ang it apparently has evolved into a cell binding site, with loss of catalytic activity but concomitant acquisition of a new biological potential.

### 3. Cell Biology

Blood vessels are composed of specialized cells known as endothelial cells. It is not surprising, given its function, that <sup>125</sup>I-labeled Ang binds to endothelial cells, and it does so in a manner that is characteristic of receptor binding—that is, time- and concentration-dependent, reversible, saturable, and competitive with unlabeled Ang. Yet Ang is a plasma protein (200 to 400 mg/L), and obviously it does not continuously stimulate new blood vessel formation as it circulates through the vasculature. This is explained by the fact that the angiogenin **receptor** is only expressed in sparsely cultured endothelial cells, not in confluent cells (as exist in blood vessels) (4). The receptor, not yet fully characterized, is a 170-kDa transmembrane protein whose presence on the cell surface correlates with the mitogenic activity of Ang. When the receptor is present, the cells respond to Ang by both increased thymidine uptake and cell proliferation. Confluent cells lack the receptor and, hence, do not respond this way to Ang.

Experiments with protein [cross-linking](#) reagents have demonstrated that Ang also interacts with a 42-kDa endothelial cell-surface protein that is a member of the [actin](#) family (1, 2). On binding Ang, it dissociates from the cell surface as an Ang–actin complex. Remarkably, this complex activates tissue plasminogen activator (tPA) and thus generates plasmin (see [Plasminogen](#)). This, in turn, stimulates cell-associated proteolytic activity to degrade the [extracellular matrix](#) and thereby facilitates cell migration, an essential feature of the angiogenesis process. Ang also acts as a [cell adhesion molecule](#) and, when coated on a plastic surface, mediates the binding of endothelial and tumor cells. *In vivo* it may help direct migrating blood vessels toward Ang-secreting tumor cells.

Addition of Ang to cultured endothelial cells stimulates a transient increase in cellular **diacylglycerol**, seemingly the result of an Ang-induced activation of [phospholipase C](#), as well as an increase in prostacyclin secretion owing to activation of phospholipase A<sub>2</sub>. Interpretation of these [second messenger](#) responses has not been possible, thus far, largely because the systems are complex and may only be stimulated indirectly by Ang.

### 4. Nuclear Translocation

Angiogenin undergoes endocytosis by sparsely cultured endothelial cells and is rapidly translocated from the cell surface to the [nucleus](#), where it accumulates in the [nucleolus](#) (1, 2). This process is receptor-mediated and regulated by a nuclear localization signal (see [Nuclear Import, Export](#)) that involves three basic residues, Arg31-Arg32-Arg33, of Ang. This signal is essential for angiogenesis and suggests that the substrate for the ribonucleolytic activity of Ang is located in the nucleolus, a highly specialized region where biogenesis of ribosomes takes place (see [Nucleolus](#) and [Ribosomes](#)). One possibility is that the ribonucleolytic activity of Ang enhances the [transcription](#) and processing of rRNA. Two enzymatically inactive variants of angiogenin whose cell-binding site is intact also accumulate in the nucleolus when added to endothelial cells, but they are angiogenically inactive. This suggests that not only is nuclear localization essential for the biological activity of Ang but, in addition, the translocated protein must be enzymatically active. This has been confirmed by recent studies in which a DNA aptamer (an oligonucleotide that binds fairly tightly to Ang) that inhibits both the enzymatic and angiogenic activities of Ang is also translocated to nucleus, but only when added to endothelial cells together with Ang. The inhibitor and the protein accumulate in the nucleus in a 1:1 stoichiometric ratio (5).

### 5. Mechanism of Action

Although Ang may well be involved in a variety of angiogenesis-related situations, its mechanism of action is likely to be quite similar for each of them; hence it is only summarized here in the context of tumor-induced angiogenesis. All tumor cells that have been examined have been found to secrete Ang, and in at least one case, pancreatic cancer, the aggressiveness of the tumor is related to the amount of angiogenin produced (6). Ang secreted by tumor cells migrates through the extracellular matrix until it reaches an endothelial cell. There it combines with cell-surface actin and dissociates as a complex that stimulates the activity of tPA and the formation of plasmin. Degradation of the extracellular matrix—for example, by plasmin-activated matrix [metalloproteinases](#)—may stimulate a few endothelial cells to migrate, and this could trigger expression of the Ang receptor. Binding of Ang to the receptor would then, on the one hand, initiate a second-messenger response, perhaps through activation of the above-mentioned phospholipases, and, on the other, promote endocytosis and nuclear localization of Ang. The ribonucleolytic action of Ang within the nucleolus would, together with signals generated via the second messengers, activate processes leading to cell proliferation. The proliferating endothelial cells would migrate through the now degraded extracellular matrix toward the tumor cell from which the Ang was released. The cell-adhesion properties of Ang may be important for ensuring cell migration in the proper direction. At present, this view of the mode of action of Ang is largely speculative, particularly because nothing is known about the events that occur within the nucleus and lead to cell division. Nevertheless, it summarizes current thinking and is a useful basis for further investigation.

## 6. Tumor Biology

The levels of Ang and its [messenger RNA](#) are increased in the tissues and cells of patients with various types of cancers, indicative of an *in vivo* role for Ang in the process of tumor angiogenesis. This being the case, anti-angiogenin agents could have potential importance in the treatment of cancer and other angiogenesis-related diseases. The most potent inhibitor of Ang is a 50-kDa protein originally isolated from placenta as an inhibitor of RNase A, hence known as placental ribonuclease inhibitor (RI) (1, 2). It inhibits Ang with a  $K_i$  of 0.7 fM, one of the strongest [protein–protein interactions](#) known, and is 60-fold stronger than its inhibition of RNase A. Unfortunately, RI is evidently unstable in extracellular fluids and has not been found to be effective in treating tumors in laboratory animals.

[Monoclonal antibodies](#) raised against Ang have been used to treat athymic mice injected with HT-29 colon carcinoma cells, and they reduced the incidence of tumors by up to 65%. Similar results have been observed with other types of tumors, and when used in conjunction with conventional therapeutic agents, a synergistic effect was seen (2). Actin was also tested for its anti-tumor activity, and it too was capable of preventing tumor growth in more than 60% of treated animals. It should be noted that anti-angiogenin agents do not affect the growth of already-established tumors. They are only effective when administered to animals at the same time as the tumor cells. They are thought to act by specific extracellular inactivation of tumor-secreted Ang and the consequent inhibition of tumor angiogenesis. These experiments are important in that they provide clear evidence of a crucial role for angiogenin in the early stages of development of these tumors.

## 7. Other Biological Properties

One of the major complications associated with regular hemodialysis is the increased morbidity and mortality arising from infections. This has been attributed to dysfunction of polymorphonuclear leukocytes. Indeed, a number of compounds have been isolated from uremic serum and shown to inhibit the biological activity of these white cells. One of these is an inhibitor of leukocyte degranulation that turns out to be Ang (7). Nanomolar concentrations of Ang inhibit both spontaneous and peptide-stimulated degranulation by 60% and 30%, respectively. However, Ang has no other effect on the cellular responses of polymorphonuclear leukocytes, such as [chemotaxis](#), [phagocytosis](#) or their peptide-stimulated oxidative respiratory burst.

Ang has also been reported to suppress significantly the proliferation of human lymphocytes stimulated by a mixed-lymphocyte culture or by phytohemagglutinin or concanavalin A (8). A maximal immunosuppressive effect was seen at an Ang dose of 50 to 100 mg/mL. It was thought that this effect might synergize with the effect of Ang on neovascularization of tumors and thereby contribute to tumor development.

## 8. Conclusions

Ang is now recognized as a pleiotropic molecule capable of inducing several intra- and extracellular activities. It induces most of the individual events in the process of angiogenesis including binding to endothelial cells, stimulating second messengers, mediating cell adhesion, activating cell-associated proteinases, inducing cell invasion, stimulating DNA synthesis and cell proliferation, and organizing the formation of tubular structures from cultured endothelial cells. Characterization of its cellular receptors, elucidation of its mechanism of nuclear translocation, identification of its intranucleolar target substrate, and understanding how these events result in cellular proliferation and blood vessel formation will provide important and novel information that can be utilized for either promoting or inhibiting angiogenesis for therapeutic purposes.

## Bibliography

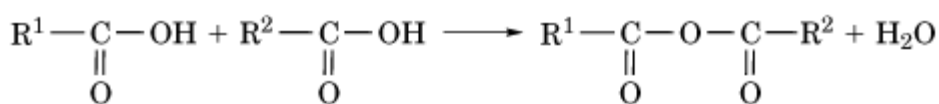
1. J. F. Riordan (1997) "Structure and Function of Angiogenin" In *Ribonucleases: Structures and Functions* (G. D'Alessio and J. F. Riordan, eds.), Academic Press, New York, pp. 445–489.
2. B. L. Vallee and J. F. Riordan (1997) *Cell. Mol. Life Sci.* **53**, 803–815.
3. R. Shapiro (1998) *Biochemistry* **37**, 6847–6856.
4. G.-f. Hu, J. F. Riordan, and B. L. Vallee (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2204–2209.
5. V. Nobile, N. Russo, G.-f. Hu, and J. F. Riordan (1998) *Biochemistry* **37**, 6857–6863.
6. S. Shimoyama et al. (1996) *Cancer Res.* **56**, 2703–2706.
7. H. Tschesche, C. Kopp, W. H. Horl, and U. Hempelmann (1994) *J. Biol. Chem.* **269**, 30274–30280.
8. J. Matousek et al. (1995) *Comp. Biochem. Physiol.* **112B**, 235–241.
9. P. J. Kraulis (1991) *J. Appl. Crystallogr. A* **24**, 946–950.

## Suggestions for Further Reading

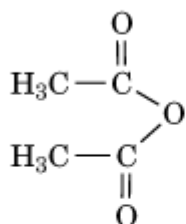
10. F. Bussolino, A. Mantovani, and G. Persico (1997) Molecular mechanisms of blood vessel formation. *Trends Biochem. Sci.* **22**, 251–256.
11. J. Folkman and P. A. D'Amore (1996) Blood vessel formation: What is its molecular basis? *Cell* **87**, 1153–1155.

## Anhydrides

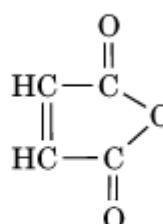
“Anhydride” means “without water,” and many anhydrides are present in nature. Anhydrides are easily hydrated in the presence of water to acids or bases. Inorganic anhydrides include  $\text{SO}_3$ ,  $\text{P}_2\text{O}_5$ ,  $\text{CaO}$ ,  $\text{Na}_2\text{O}$ , etc. Acid anhydrides are the most important anhydrides in biochemistry. An acid anhydride [I] is formed by eliminating a molecule of water from two acids (Scheme 1):



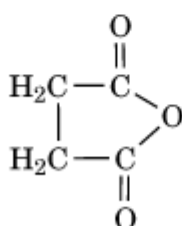
[I] Carboxylic acid anhydride



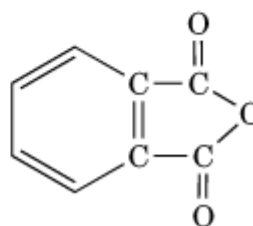
[II] Acetic anhydride



[III] Maleic anhydride



[IV] Succinic anhydride



[V] Phthalic anhydride

Adding a molecule of water reverses the reaction. Practically, a carboxylic acid anhydride is prepared by the nucleophilic displacement of chloride ion from acyl chlorides by carboxylate ion.

Acid anhydrides are reactive and are good acylating agents. They react with water to yield acids, with [amino groups](#) to yield amides, and with alcohols to yield esters. Acid anhydrides are utilized as reactive intermediates in organic syntheses, such as [peptide synthesis](#), and sometimes as **enzyme**-substrate intermediates. Acid anhydrides are unstable and therefore have large standard **free energies** of hydrolysis. The two terminal phosphate linkages in ATP are anhydride linkages, and acid anhydrides play important roles in bioenergetics (see [Adenylate Charge](#)). Acid anhydrides are widely employed for [chemical modification](#) of proteins, particularly of amino groups. The acetylation reaction by an acid anhydride with amino groups in a protein proceeds so easily that this reaction is suitable for the **radiolabeling** of a protein with a radioisotope-labeled acid anhydride (1). Acetic anhydride [II], maleic anhydride [III], and succinic anhydride [IV] are frequently employed in chemical reactions and chemical modifications of proteins. Phthalic anhydride [V] is an aromatic anhydride and is a good leaving group.

## 1. Acetic Anhydride

Acetic anhydride is the most important acid anhydride. Its boiling point is 139.6°C, and it has a characteristic penetrating odor. It is widely employed in organic syntheses. It is used to introduce the acetate group into organic compounds and is called an acetylating agent. To acetylate amino groups in a protein (2), prepare a protein solution (2 to 10%, w/v) in half-saturated sodium acetate. Over a period of 1 h at 0°C, add, in five equal portions, a weight of acetic anhydride equal to that of the protein, and continue stirring for an additional hour.

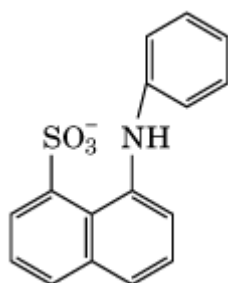
## Bibliography

1. T. N. M. Schumacher, and T. J. Tsomides (1995) In *Current Protocols in Protein Science* (J. E. Coligan et al., eds.), Wiley, pp. 3.3.11–3.3.12.
2. T. Imoto and H. Yamada (1989) In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 247–277.

## 8-Anilidonaphthalene-1-Sulfonic Acid

8-Anilidonaphthalene-1-sulfonic acid (ANS) (Fig. 1) and its dimer, 1,1'-bis(4-anilino-5-naphthalenesulfonic acid) (bis ANS) are two widely used “**hydrophobic**” probes in **fluorescence** studies of [proteins](#) and [membranes](#) (1). These dyes are minimally fluorescent in [polar](#) environments, such as aqueous solutions, but their fluorescent emission is dramatically increased in [nonpolar](#) environments. An increase in quantum yield (intensity) and a blue shift in the wavelength of maximum emission ( $I_{\max}$ ) is observed when binding to macromolecules.

**Figure 1.** The structure of 8-anilidonaphthalene-1-sulfonic acid (ANS).



The most common uses of ANS (and bis-ANS) are (1) to monitor conformational changes in proteins (2, 3); (2) to follow the kinetics of **protein folding** and unfolding (4-6); (3) to detect and characterize partially folded intermediates of proteins (both transient and equilibrium) (7, 8); (4) to measure changes in membrane properties (9, 10); (5) to detect the presence of exposed hydrophobic [accessible surfaces](#) (11); and (6) to study **ligand binding** through displacement assays (12, 13). Photoincorporation of bis-ANS has been used to locate solvent-exposed hydrophobic regions in proteins (14) and more specifically (1) to investigate **virus** capsid structure and assembly (15, 16); (2) to monitor [tubulin](#) assembly and interaction with drugs (17) because bis-ANS specifically inhibits assembly of tubulin (18); (3) to measure distances via fluorescent energy transfer from [tryptophan](#) residues to bound ANS (17, 19); (4) to study **enzyme**–substrate interactions (20), including **nucleotide** binding (21); and (5) as a sensitive staining method for proteins in [SDS-PAGE](#) (22). ANS has been used to detect conformational changes in a wide variety of functional proteins during ligand binding (ranging from small molecules to DNA) (13, 23).

The advantages of ANS as a probe reflect the simplicity of fluorescence experiments and the sensitivity of ANS to its immediate environment. Although most commonly used in the fluorescent steady-state mode, it has also been used in more complex experiments (6, 25). The **absorbance** spectrum of ANS has maxima at ~270 nm and in the 345 to 380 nm range. The fluorescent emission spectrum is very sensitive to the environment and has a maximum of ~520 nm in [water](#), which becomes progressively blue-shifted with decreasing polarity and reaches 464 nm in very nonpolar environments. A corresponding change occurs in the fluorescent lifetime from ~0.25 ns in water to ~15 ns in hydrophobic environments. The overlap between ANS excitation and tryptophan emission in the vicinity of 350 nm provides a useful donor/acceptor energy transfer system. Although the properties of ANS and bis-ANS are similar, the dimer has substantially higher affinity for many protein sites (26).

It is commonly assumed that ANS and bis-ANS bind to proteins and membranes by hydrophobic interactions mediated by the two aromatic rings, although there is little direct data (such as [X-ray crystallography](#) or [NMR](#) structures) to support this hypothesis. These dyes are anions over most of the commonly used pH range and [electrostatic interactions](#) are also important (possibly even critical) in ANS binding to proteins and membranes ([27](#)).

The applications of ANS to biomembrane studies have been summarized ([24](#)) and are not presented in any depth here. ANS binds to membranes in the vicinity of the phospholipid polar headgroups. The nonpolar part of the ANS molecule is directed toward the region of the fatty acid side chains, but the ANS does not penetrate very far into the nonpolar region of the bilayer. The binding is very sensitive to the nature of the phospholipids and to the presence of other membrane components, such as cholesterol.

For proteins, it is likely that only those ANS molecules that bind in hydrophobic regions exhibit the characteristic strong fluorescence. The remaining ANS molecules are bound to cationic sites exposed to aqueous solvent and hence are quenched ([28](#)). The thermodynamics of ANS binding to a protein have been determined: the binding is enthalpically driven at all temperatures and entropically opposed at temperatures greater than 14°C ([27](#)). The self-association of ANS and bis-ANS in water correlates with their tendency to form complexes with proteins ([29](#)). Thus care must be taken, especially when working at low pH and high ANS concentrations, to avoid ANS aggregation.

It has been recognized for some time that ANS and bis-ANS are excellent probes for partially folded intermediate states of proteins, such as the [molten globule](#) ([7-30](#)). Because the dyes do not normally bind significantly to the native states of most proteins nor to the fully **unfolded** states, they are widely used as a diagnostic for the presence of partially folded intermediates. However, ANS and bis-ANS do bind to some native proteins, usually to a specific site, e.g., the heme-binding site in apo-**myoglobin** ([31](#)). The predilection of ANS for nonnative conformations has been useful in investigating designed proteins to determine whether they have a rigid, tightly packed core, or are more like a molten globule ([32](#)), and in studies to determine the conformational state of substrate proteins interacting with molecular chaperones ([33](#)).

Fluorescent changes that correspond to binding and release of ANS during protein folding are frequently used to characterize the kinetics of folding, especially the buildup and decay of transient intermediates ([34-36](#)). Differences in partially folded intermediate conformations of proteins are distinguished by differences in the properties of bound ANS ([5](#), [37](#)).

The interpretation of results from ANS probe experiments is not without hazards, in part because the dyes perturb the system under investigation. For example, in experiments to monitor the kinetics of protein folding, the presence of ANS or bis-ANS and their strong affinity for transient intermediates significantly perturb the kinetics of folding ([35](#)). Similarly, it is customarily assumed that if ANS or bis-ANS is added to the native protein and binding is observed, it is to the native state. However, when the native state is only nominally predominant, preferential binding of the dye to a partially folded intermediate effectively shifts the equilibrium to favor the intermediate conformation ([19](#)). The presence of tight-binding ligands (eg, substrates) to the native state shifts the equilibrium back to the native conformation ([19](#)). Another source of potential complications in interpreting results arises from the fact that the dye concentration used in many studies drastically exceeds the protein concentration. It is important to use purified ANS. Most commercial preparations contain some bis-ANS, which produces misleading results because bis-ANS often binds much more tightly than ANS. The bis-ANS is removed by Sephadex LH-20 chromatography ([38](#)).

## Bibliography

1. L. Stryer (1965) *J. Mol. Biol.* **13**, 482–495.
2. T. Korte and A. Herrmann (1994) *Eur. Biophys. J.* **23**, 105–113.



3. H. Tsuruta and T. Sano (1990) *Biophys. Chem.* **35**, 75–84.
4. G. de Prat Gay et al. (1995) *J. Mol. Biol.* **254**, 968–979.
5. M. C. Shastry and J. B. Udgaonkar (1995) *J. Mol. Biol.* **247**, 1013–1027.
6. B. E. Jones, J. M. Beechem, and C. R. Matthews (1995) *Biochemistry* **34**, 1867–1877.
7. G. V. Semisotnov et al. (1991) *Biopolymers* **31**, 119–128.
8. Y. Goto and A. L. Fink (1989) *Biochemistry* **28**, 945–952.
9. A. Azzi (1975) *Q Rev. Biophys.* **8**, 237–316.
10. M. G. Tozzi-Ciancarelli, C. Di Massimo, and A. Mascioli (1992) *Cell. Mol. Biol.* **38**, 303–310.
11. M. Cardamone and N. K. Puri (1992) *Biochem. J.* **282**, 589–593.
12. C. D. Kane and D. A. Bernlohr (1996) *Anal. Biochem.* **233**, 197–204.
13. I. Taylor and G. G. Kneale (1994) *Methods Mol. Biol.* **30**, 327–337.
14. J. W. Seale, J. L. Martinez, and P. M. Horowitz (1995) *Biochemistry* **34**, 7443–7449.
15. A. T. Da Poian, J. E. Johnson, and J. L. Silva (1994) *Biochemistry* **33**, 8339–8346.
16. J. Secnik et al. (1990) *Biochemistry* **29**, 7991–7997.
17. A. Bhattacharya, B. Bhattacharyya, and S. Roy (1996) *Protein Sci.* **5**, 2029–2036.
18. P. Horowitz, V. Prasad, and R. F. Luduena (1984) *J. Biol. Chem.* **259**, 14647–14650.
19. L. Shi, D. R. Palleros, and A. L. Fink (1994) *Biochemistry* **33**, 7536–7546.
20. P. M. Horowitz and N. L. Criscimagna (1985) *Biochemistry* **24**, 2587–2593.
21. R. Takashi, Y. Tonomura, and M. F. Morales (1977) *Proc. Natl. Acad. Sci. USA* **74**, 2334–2338.
22. P. M. Horowitz and S. Bowman (1987) *Anal. Biochem.* **165**, 430–434.
23. A. R. Walmsley, G. E. Martin, and P. J. Henderson (1994) *J. Biol. Chem.* **269**, 17009–17019.
24. J. Slavik (1982) *Biochim. Biophys. Acta* **694**, 1–25.
25. E. E. Gussakovsky and E. Haas (1995) *Protein Sci.* **4**, 2319–2326.
26. C. G. Rosen and G. Weber (1969) *Biochemistry* **8**, 3915–3920.
27. W. R. Kirk, E. Kurian, and F. G. Prendergast (1996) *Biophys. J.* **70**, 69–83.
28. D. Matulis, R. Lovrien, and T. I. Richardson (1996) *J. Mol. Recognition* **9**, 433–443.
29. B. Stopa et al. (1997) *Biochimie.* **79**, 23–26.
30. Y. Goto, T. Azuma, and K. Hamaguchi (1979) *J. Biochem.* **85**, 1427–1438.
31. E. Bismuto et al. (1996) *Protein Sci.* **5**, 121–126.
32. S. F. Betz et al. (1996) *Fold Des* **1**, 57–64.
33. M. Gross et al. (1996) *Protein Sci.* **5**, 2506–2513.
34. A. F. Chaffotte et al. (1992) *Biochemistry* **31**, 4303–4308.
35. M. Engelhard and P. A. Evans (1995) *Protein Sci.* **4**, 1553–1562.
36. O. B. Ptitsyn et al. (1990) *FEBS Lett.* **262**, 20–24.
37. E. Zerovnik et al. (1997) *Eur. J. Biochem.* **245**, 364–372.
38. S. S. York, R. C. Lawson Jr., and D. M. Worah (1978) *Biochemistry* **17**, 4480–4486.

### **Suggestions for Further Reading**

39. C. A. Royer (1995) *Fluorescence spectroscopy. Methods Mol. Biol.* **40**, 65–89.
40. J. Slavik (1982) Anilinonaphthalene sulfonate as a probe of membrane composition and function. *Biochim. Biophys. Acta* **694**, 1–25.

## Animal Pole, Vegetal Pole

The cell mass of an [egg](#) is not uniformly distributed, but it exhibits significant differences in terms of morphology and at the molecular level. In order to describe this polarity, the terms animal pole and *vegetal pole* were invented to describe the two opposite poles of the egg.

Unlike the eggs of insects, which are elliptically shaped, most **oocytes** of amphibians and mammals exhibit a less pronounced asymmetry. One element of asymmetry is the location of the cell [nucleus](#), which is normally not right in the center of the oocyte, but is more peripheral, sometimes even adjacent to the egg membrane. Due to gravitational forces, the yolk within settles to the bottom of the egg and forms the vegetal pole; consequently, the polarity becomes more apparent when more yolk is present in the egg. Another component causing asymmetry is that **polar bodies** extruded during **meiosis** are frequently located in the region of the animal pole. In ascidians, [sperm](#) enters the egg somewhere in the animal hemisphere (1), causing cytoplasmic movements and rotation of the egg cortex. As a consequence, a further distinct region, the *gray crescent*, becomes visible in the fertilized egg. The gray crescent is also visible in amphibians, but is not very apparent in sea urchins. In the egg of *Unio elongatulus*, sperm entry occurs only at the vegetal pole. This is attributed to the presence of a 220-kDa binding protein that is concentrated in a restricted region of the crater region within the vegetal pole (2).

The naming of the two poles, animal and vegetal, is not based on a precise function; rather, the names have arisen from the idea that the “higher” organs evolve in the animal polar region, whereas the vegetal pole was assumed to be destined to form the “lower” organs necessary for reproduction and providing nutrition. The two poles form one of three possible coordinates, and further developmental changes in a number of amphibians correlate with the subsequent dorsal-ventral body axis of the animal. In mammals, the mechanism by which the inner cell mass settles in certain places is not fully understood. However, it was shown in the mouse that the bilateral symmetry of the early blastocyst is normally aligned with the animal-vegetal axis of the zygote. The embryonic-abembryonic axis is oriented orthogonally to the animal-vegetal axis (3).

After a sperm activates the egg, [karyogamy](#) of the male and female **pronuclei** occurs, and the egg starts to divide by **mitosis**. Most important, the cell mass does not increase during the first cell divisions; starting from the one cell, two cells are formed, after another round of divisions four cells, then eight cells, 16 cells, and so forth, until a great number of smaller cells are formed at the [morula](#) stage. These cells do not all have the same size; the smaller ones are called micromeres, whereas the larger ones are named *macromeres*. During these cleavage steps, there is relatively little **gene expression** from the nucleus of the new individual cells. In other words, the [genome](#) of the new cell does not determine its own development after [fertilization](#) and karyogamy at these very early stages of development. With regard to the new cells, the regulating elements are external, maternally derived gene products. These substances, mostly **RNA** molecules, are already present in the unfertilized egg, and they are the important factors that determine the fate of the divided cells. This has been shown by a number of experiments. Chemical inactivation or enucleation in embryos has shown that the nucleus is not necessary for the initial rounds of cleavage. Even enucleated egg fragments are able to perform developmental changes, and cross-fertilization experiments revealed that cells follow the maternal pattern of development. These maternally derived gene products are not evenly distributed in the egg, indicating that the animal and vegetal pole are not only a matter of morphological appearance, but are also related to the presence of a concentration gradient of different gene products. This was demonstrated in experiments in which sections of the animal or vegetal pole were excised and recultivated. When micromeres of the vegetal pole at the 16-cell stage were implanted into the animal pole of a donor embryo, a complete second gut developed; the micromeres are capable of changing the fate of neighboring cells (4). Cutting sea urchin eggs at the eight-cell stage into two halves, to produce two embryos each with an animal and a vegetal cell,

generates two pluteus capable of normal [development](#). In contrast, cutting the cell into an animal and an vegetal hemisphere causes considerable aberrations from the normal development.

One of the most quoted experiments to show that animal and vegetal pole cells differ significantly involved excision in which the fate of sections at the 64-cell stage were investigated. The individual cells alone were capable of forming only a **blastula**. However, adding at least four micromeres from the vegetal pole could compensate for this defect and lead to a pluteus. Fewer numbers of micromeres resulted in forms intermediate between blastula and normal pluteus. Interestingly, the vegetalization occurs not only after the addition of micromeres, but also when  $\text{Li}^+$  ions were provided. At the animal pole, [messenger RNA](#) were found that code for cell-surface proteins. These mRNA were found only in the macromeres and mesomeres, not in the micromeres (5). In the vegetal pole region are localized dorsal determinants that are necessary for dorsal axis development in [Xenopus](#). The dorsal determinants move from the vegetal pole to a subequatorial region, where they are incorporated into gastrulating cells (6). However, dorsal development may be also activated by the contact between the cortical dorsal determinant and the equatorial core cytoplasm that are brought together by cortical rotation upon fertilization (7).

**Gastrulation** is the next morphogenic event, and it starts in the region of the vegetal pole or near to the gray crescent, depending on the species. At least five different epigenic movements are observable in the further process of differentiation: invagination of cells, immigration, delamination, proliferation, and epiboly. This means that, although the polarity of the egg determines the fate of the cells at different positions, further movements and reorganization take place that are necessary for organogenesis in higher mammals. One predominantly epigenic movement is invagination during gastrulation. It was shown that vegetal egg cytoplasm is responsible for the specification of vegetal [blastomeres](#) and promotes gastrulation. Vegetal-deficient embryos of *Halocynthia roretzi* fail to enter gastrulation. They are animalized and arrested at the blastula stage. Reimplantation of the vegetal pole cytoplasm into vegetal-deficient embryos caused gastrulation at the site where implantation occurred (8).

### Bibliography

1. J. E. Speksnijder, L. F. Jaffe, and C. Sardet (1989) *Dev. Biol.* **133**, 180–184.
2. R. Focarelli and F. Rosati (1995) *Dev. Biol.* **171**, 606–614.
3. R. L. Gardner (1997) *Development* **124**, 289–301.
4. A. Ransick and E. H. Davidson (1993) *Science* **259**, 1134–1138.
5. M. Di Carlo, D. P. Romancino, G. Montana, and G. Ghersi (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5622–5626.
6. M. Sakai (1996) *Development* **122**, 2207–2214.
7. H. Kageura (1997) *Development* **124**, 1543–1551.
8. H. Nishida (1996) *Development* **122**, 1271–1279.

### Suggestion for Further Reading

9. H. Eyal-Giladi (1997) Establishment of the axis in chordates: Facts and speculations. *Development* **124**, 2285–2296. A review of the establishment of the axis in chordates.

### Ankyrins

Ankyrins are a family of conserved proteins whose prime function is to act as a link between [integral membrane proteins](#) and the **spectrin-** or fodrin-based framework lying on the cytoplasmic side of the [plasmalemma](#). The ankyrins were first characterized as those proteins that displayed a strong affinity for spectrin in **erythrocyte** membranes. They are expressed in particularly high levels in vertebrate brain but are also found in tissues such as skeletal muscle, lymphocytes, neutrophils and many epithelial tissues (1); they are also found in such primitive organisms as worms ([nematodes](#)) and fruit flies (*Drosophila*). Associations between ankyrins and a diverse selection of membrane proteins such as [cell adhesion molecules](#) and **ion channels** (the **anion exchanger**, or band 3 protein, Na<sup>+</sup>/K<sup>+</sup> [ATPase](#) and the voltage-sensitive **sodium channel**) have been characterized and confirm that a key role of the ankyrins is to mediate interactions between spectrin in erythrocytes (known as fodrin in other cells) and proteins constituting the cell membrane. Deficiencies in ankyrin (as well as spectrin and some other membrane-associated proteins) have been correlated with structurally weak erythrocytes.

Ankyrins from erythrocytes (Ank1) and brain (Ank2) are monomeric proteins with a simple **domain** substructure. Each protein has an *N*-terminal domain (about 89 to 95 kDa) with membrane-binding abilities that can interact with the anion exchanger and [tubulin](#) (amongst others). Most of this globular region comprises 24 tandem quasi-repeats of a 33-residue motif. Because six such repeats appear to represent the smallest structural entity, it is possible that there are four subdomains within the *N*-terminal domain. It has also been suggested that the 33-residue structures, while folding up in a common manner, could position nonconserved residues on surface sites that would enable them to interact with a wide range of different protein ligands, rather than with a single ligand multiple times. The sequence of the 33-residue motif is very similar to those found in [transcription factors](#), in [cell-cycle](#) control proteins, and in proteins regulating tissue [differentiation](#). The second domain of ankyrin (about 62 kDa) can be subdivided into two parts; the first is acidic, proline-rich, and about 80 residues in length. In erythrocytes it links the cytoplasmic domain of the transmembrane anion exchanger to b-spectrin at about the midpoint of the tetrameric structure that b-spectrin forms with a-spectrin (see [Spectrin](#)). The second portion is basic and much larger (almost 500 residues). It contains the consensus sequence (Arg–Arg–Arg–Lys–Phe–His–Lys/Arg) that is also required for spectrin- or fodrin-binding. It is also much more highly conserved in sequence than the acidic 80-residue subdomain. The *C*-terminal domain is about 70 kDa in both Ank1 and one variant of Ank2, but about 290 kDa in an **isoform** of Ank2 that is generated by [alternative splicing](#). The smaller isoform of Ank2 (molecular weight about 220 kDa) is common in adult rats, while the larger one (about 440 kDa) is expressed highly in neonatal animals. Furthermore, the 220-kDa protein is found widely in the brain (neuron cell bodies, dendrites and glia), whereas the 440-kDa protein seems to be found specifically in unmyelinated axons and dendrites (2). In both Ank1 and Ank2, the *C*-terminal domain has a regulatory role arising from its ability to modulate the binding affinities of both the membrane-binding and spectrin-binding domains. The sequences of the *C*-terminal domains differ in size, and the overall homology is not high, although there are some regions of similarity.

Ank3, which is the major ankyrin in kidney, is found in most epithelial cells, axons and muscle cells and has many structural similarities to Ank1 and Ank2. It too has an *N*-terminal domain (89 kDa) containing a 33-residue repeat with the **consensus sequence** (Asp/Asn–Gly–apolar–Thr–Pro/Ala–Leu–His–apolar–Ala–Ala–X–X–Gly–His/Asn–apolar–X–Val/Ile–Val/Ala–X–apolar–Leu–Leu–X–X–Gly–Ala–X–Apolar/Pro–Asn/Asp–Ala–X–Thr–Basic), a 65-kDa domain with spectrin-binding ability, and a *C*-terminal domain (56 kDa) with regulatory attributes (2). Interestingly, however, multiple transcripts are expressed, some of which lack the repeat domain at the *N*-terminal end of the molecule. It is possible that these ankyrins are involved in intracellular vesicle stabilization, sorting, or targeting (2).

## Bibliography

1. V. Bennett and D. M. Gilligan (1993) The spectrin-based membrane skeleton and micron-scale organization of the plasma membrane. *Annu. Rev. Cell Biol.* **9**, 27–66.
2. L. L. Peters, K. M. John, F. M. Lu, E. M. Eicher, A. Higgins, M. Yialamas, L. C. Turtzo, A. J.

Otsuka, and S. E. Lux (1995) Ank3 (epithelial ankyrin), a widely distributed new member of the ankyrin gene family and the major ankyrin in kidney, is expressed in alternatively spliced forms, including forms that lack the repeat domain. *J. Cell Biol.* **130**, 313–330.

## Annealing, Nucleic Acids

When two complementary nucleic acid strands are mixed under conditions favoring complex formation, the fully associated duplex is the lowest free-energy state of the system. The annealing process involves formation and disruption of partially **base-paired** regions and expansion of those regions as the system samples progressively more extensively base-paired intermediate states on the path to the lowest energy fully base-paired state. This lowest energy state may not be achieved due to kinetically trapped structures.

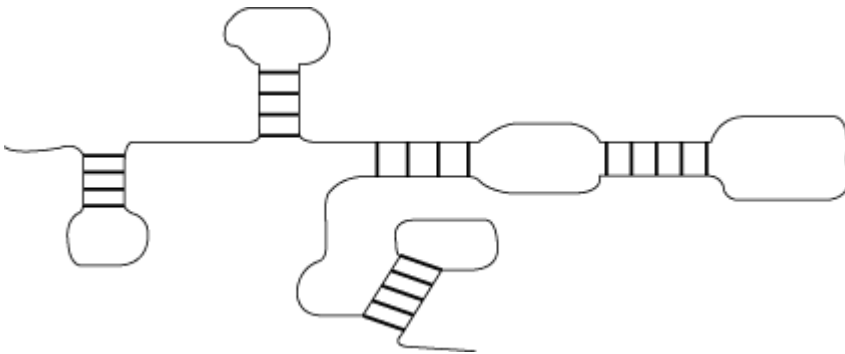
**Denaturation** of short oligonucleotide duplexes is a thermodynamically and kinetically reversible process. When the solution is cooled after melting of the duplex, the optical and hydrodynamic properties of oligonucleotide duplexes return to the initial values. The denaturation of synthetic polymers of repeating sequence are optically reversible, which indicates that most bases are in intact base pairs. However, the hydrodynamic properties display hysteresis. The change in hydrodynamic properties is due to the formation of branched structures comprised of hairpin loops and concatenated aggregates. Natural DNAs display hysteresis in both optical and hydrodynamic properties, indicating failure of a significant fraction of the bases to find a partner base and the presence of a nonlinear structure.

When two nucleic acid polymers encounter each other during hybridization or renaturation, the probability of an exact match of base pairs in the initial complex is very small. Instead, small stretches of duplex are nucleated. These complexes may grow if the sequences of the two strands permit. If there is sufficient thermal energy, ie, sufficiently high temperature, the nucleated regions will form and disassociate as the system seeks the lowest free-energy state. If there is not sufficient thermal energy, unproductive complexes will dominate, and the equilibrium configuration will not be achieved. Conversely, if there is too much thermal energy, formation of the duplex regions will be disfavored, and the strands will separate.

This leads to a common strategy of duplex formation. The annealing process is optimized by incubation at the “melting” temperature,  $T_m$ , followed by slow cooling to permit the facile reorganization of the duplex regions, to optimize sequence alignment.

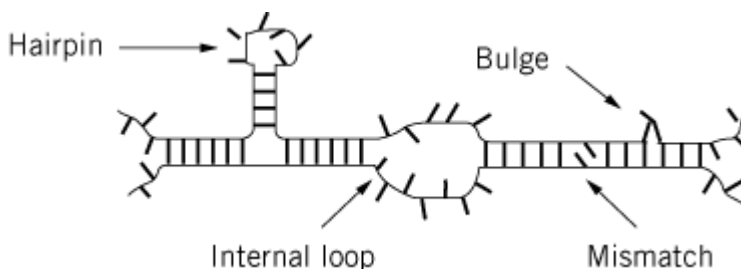
Intramolecular hairpins and internal loops in single-stranded nucleic acid polymers present a significant impediment to annealing (Fig. 1). These structures form rapidly compared to duplex structures. The rate constant for hairpin formation is  $10^5$ – $10^7$   $s^{-1}$ , meaning the structures form on the microsecond timescale. The second-order rate constant for oligonucleotide duplex formation is  $10^5$ – $10^7$   $M^{-1} s^{-1}$ . Therefore, at the low concentration of annealing experiments, duplexes form on a time scale of seconds or longer. Intramolecular structures are at a lower energy than is the fully unfolded single-stranded structure. These competing structures must be unfolded for the annealing process to proceed toward the formation of the complementary duplex. The [activation energy](#) associated with this unfolding process is large, so that at temperatures well below the  $T_m$  of these single-stranded structures, the rate of duplex formation is very slow.

**Figure 1.** Intramolecular structures formed by nucleic acid single strands.



Defects in the alignment of base sequences (Fig. 2), including hairpins and internal loops, as well as mismatches and bulges, can occur in duplexes. Here the activation energy, and thus the time required for conversion, is even greater than for single-stranded structures. Because the alphabet of nucleic acids contains only four letters, there are regions of nucleotide sequence complementarity even in random sequences. These regions may form nucleation complexes that cannot be extended, resulting in large numbers of nonproductive intermediate states. In highly repetitive sequences, large numbers of such structures may be formed. Under permissive conditions, they are readily interconvertible, and the proper sequence alignment can be achieved.

**Figure 2.** Misalignment structures in nucleic acid duplexes.

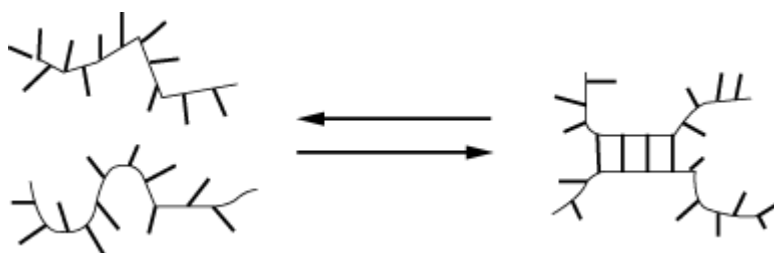


In contrast to polymers where sequence alignment dominates the kinetics of annealing, formation of short oligonucleotide duplexes is dominated by the frequency of encounter of the complementary strands. Once nucleation occurs, the pairing of the remaining bases is rapid.

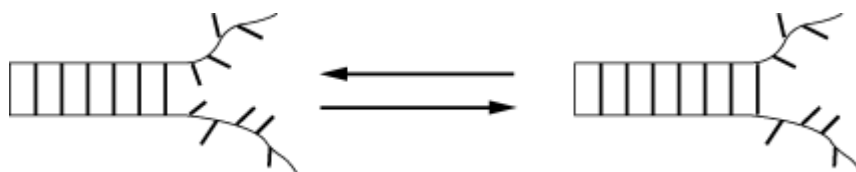
Short oligonucleotide duplexes can form without appreciable competing structures. However, competition between the duplex and single-stranded hairpin structures can occur with oligonucleotides of appropriate sequence. The *unfolded single-strand-to-hairpin* and *unfolded single-strand-to-duplex* equilibria are in competition. The two equilibria differ in molecularity, so they can be distinguished experimentally by the concentration dependence. The hairpin structure will be favored at low concentration; and the duplex, at high strand concentration. The salt concentration also affects these equilibria. Because of its lower charge density, the hairpin is favored at low salt concentrations. Therefore, a combination of salt and oligonucleotide concentration can be used to isolate one of the structures.

Because the formation of a complex involves more than one nucleic acid strand, the concentration and the length of the constituent strands also influence the stability of the nucleic acid complex. There is a complex interplay between the length and concentration in determining complex stability. Two events contribute to the formation and stabilization of the complex. The nucleation event (Fig. 3), in which the strands make initial productive contacts, is a multimolecular (bimolecular for duplex formation) process that depends on strand concentration. The propagation process (Fig. 4), also called zipping, in which the pre-nucleated base-paired chain is extended, is a pseudomonomolecular process. Formally more than one strand is involved, but the formation of a closed base pair adjacent to a preformed base pair is effectively a monomolecular process analogous to extension of a hairpin.

**Figure 3.** The nucleation event. The number of strands that are interacting changes. The process is bimolecular and depends on strand concentration.



**Figure 4.** The propagation (zipping) event. The number of strands in the complex does not change; therefore, the process is pseudomonomolecular.



For short oligonucleotides, the nucleation event dominates complex formation. The concentration dependence of the stability is simple and described adequately by the conventional mass-action model. For polymeric complexes, nucleation accounts for a relatively small part of the total stabilizing free energy. The concentration dependence of complex formation disappears for polymeric nucleic acids. For complexes of intermediate length, both nucleation and propagation make nontrivial contributions to complex formation and stability. The concentration dependence of intermediate length complex stability is reduced relative to that for short complexes of short oligonucleotides, but not vanishing as for polymers.

#### Suggestion for Further Reading

V. A. Bloomfield, D. M. Crothers, and I. Tinoco (1974) *Physical Chemistry of Nucleic Acids*, Harper & Row, New York.

## Annexins

Annexins are a family of structurally homologous, calcium- and membrane-binding proteins from eukaryotic sources including animals, plants, protists, and fungi. More than twenty distinct annexins have been characterized. Annexin homologues, notably 14-3-3 proteins, have also been identified through sequence similarities. Intracellularly, annexins vary in their organ, tissue, cell, and subcellular distribution and localization, and some annexins also occur extracellularly. The mechanism of release of annexins into the extracellular milieu is unclear because annexins lack **signal sequences**. Depending on the source, some annexins constitute a significant proportion (>1%) of total cell protein. Though their *in vivo* functions are not yet established, annexins have been implicated in many processes, including cell proliferation and [differentiation](#), [signal transduction](#), membrane trafficking and secretion, **vesicle** aggregation, cell adhesion, interactions with [cytoskeletal](#) elements, [apoptosis](#), anticoagulation, [phospholipase A<sub>2</sub>](#) inhibition, and **ion channel** activity or regulation. **Phosphorylation** may regulate some annexin-mediated processes. At least two annexins are major substrates for the tyrosine kinases, epidermal growth factor receptor, and retrovirus-encoded protein tyrosine kinase pp60v-src. Phosphorylations through other **tyrosine kinases** or **serine/threonine kinases**, for example, protein kinase C and cAMP-dependent protein kinase A, also have been observed.

Before annexins were recognized as a family, the proteins were identified in many laboratories by their calcium-dependent binding to particulate fractions or [hydrophobic chromatographic](#) column matrices. Initially these proteins were given diverse names, which reflected either their source or putative function, including lipocortins, calpactins, endonexins, placental anticoagulant proteins (PAPs), chromobindins, calcimedins, and calelectrins. Subsequent sequence analysis revealed that many of these proteins have a common identity. The general term *annexin* subsequently was adopted in reference to their common ability to *annex to* membranes. In common nomenclature, individual annexins are designated by Roman numerals for example, annexin VII.

## 1. Macromolecular Interactions

The binding of calcium ions to annexins is relatively weak, with **dissociation constants**  $K_d$  in the millimolar-to-micromolar range. Complexation with membranes or other proteins increases their calcium-binding affinities. *In vitro*, annexins bind lanthanides and divalent cations, such as strontium, barium, and zinc. Calcium-dependent binding to membranes is high-affinity, and  $K_d$  is in the nanomolar range, whereas phospholipid monomers are bound poorly. In many systems, the annexin-membrane association is reversible upon addition of calcium chelators. In other systems, annexins behave more like integral [membrane proteins](#). Typically, acidic phospholipids, such as phosphatidylserine, are strongly preferred by annexins, whereas binding to pure phosphatidylcholine membranes has not been observed. Such phospholipid binding preferences may target annexins to specific locations. The strong response of annexin V to extracellular exposure of phosphatidylserine is widely used as the basis of flow cytometric assays to detect membrane phospholipid asymmetry in apoptosis and other cellular processes.

In addition to their membrane lipid-binding properties, some annexins interact with other proteins. Some annexins interact with members of the **EF-hand** family, which suggests structural complementarity between the two protein families.

Annexins II and XI each form heterocomplexes in which the heavy chain is an annexin and the light chain resembles S-100 protein. Extracellular annexins have been described as cell surface **receptors** for some **viruses** ([1](#), [2](#)) and matrix proteins, such as [collagen](#) (reviewed in Ref. [3](#)) and tenascin C ([4](#)). Annexin II is identified as a co-receptor for **plasminogen**/tissue plasminogen activator (reviewed in Ref. [5](#)). Calcium-dependent [lectin](#) activity has been identified in some annexins that bind to specific sialoglycoproteins and glycosaminoglycans, such as heparin ([6](#), [7](#)). The biological importance of this interaction is not yet understood, but it may be related to cell surface function.



## 2. Primary Structures and Molecular Evolution

Annexin sequences are comprised of two regions, a variable N-terminal region and a conserved C-terminal core region. The core region consists of four or eight canonical repeats of approximately 70 amino acid residues. Sequence homology exists between repeats within an individual annexin and between annexins. The sequence data suggest that [gene duplication](#) and subsequent [gene fusion](#) produced a four-**domain** ancestral precursor. Further gene fusion of two four-domain annexins produced annexin VI, which is unique in its possession of eight repeats. The core region is associated with the calcium-dependent phospholipid-binding properties of annexins. Each repeat contains a highly conserved stretch of amino acid residues with a **consensus** calcium-binding sequence (Lys/Arg)-Gly-X-Gly-Thr-X<sub>38</sub>-(Asp/Glu).

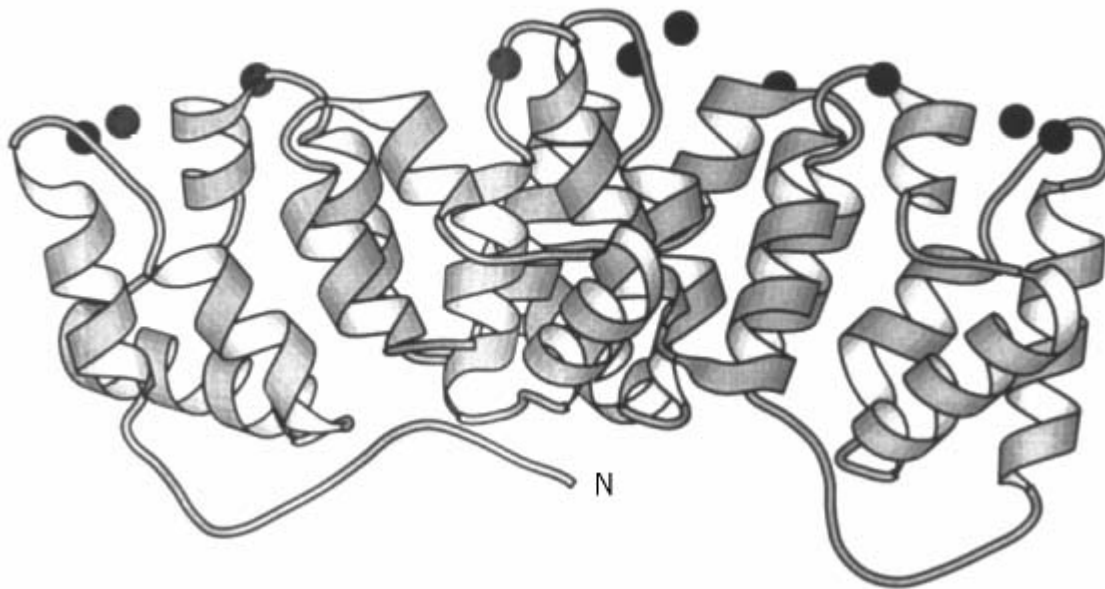
The variable N-terminal regions of annexins exhibit little homology and confer the distinct functional properties of individual annexins. Here the **intron**-exon structures are not highly conserved, whereas they are in the core region. The N-terminal region in annexins may consist of only a few amino acid residues or hundreds. Annexins VII (synexin) and XI have the largest and most hydrophobic N-terminal regions. These may considerably modify the properties of these proteins that arise from the common annexin core. In many annexins, phosphorylation sites occur within the N-terminal regions, as do some [protein-protein interaction](#) sites.

## 3. Crystal Structures of Annexins

Annexin crystal structures reveal a highly conserved **alpha-helical** fold in the C-terminal core region characteristic of this family. Each basic repeating unit is made up of five  $\alpha$ -helices, four of which run antiparallel, in a four-helix bundle and the fifth helix is perpendicular. The four-domain annexin structure is nearly planar, and the domains are arranged as a symmetrical array. The eight-domain annexin VI resembles two four-domain annexin structures that are approximately perpendicular to each other and are linked by a long,  $\alpha$ -helical segment (8, 9). Most annexin crystal structures lack the variable N-terminal domains, which are intentionally truncated or lost through proteolysis. The only exception to date, that of full-length annexin I, reveals an unexpected juxtaposition between helices in the N-terminus and core domains, a finding that may suggest a mechanism for membrane aggregation (10).

In annexin crystal structures, numerous bound calcium ions are observed in loops along one surface of the protein molecule. These loops have structural motifs that are distinct from those in EF-hand proteins or C2-domain proteins. There are several structural motifs for these calcium-binding sites. The site corresponding to the consensus calcium-binding sequence has the highest affinity for calcium and structurally resembles a type of site observed in phospholipase A<sub>2</sub>. Other sites bind lanthanides more strongly than calcium. Annexin V (Fig. 1) adopts two distinct molecular conformations, depending on whether a [tryptophan](#) side chain in the third domain extends outward from the protein surface or lies buried within the protein core. [Spectroscopic](#) and biochemical evidence indicates that the two states are interrelated and may describe a calcium-dependent conformational change that influences membrane binding. The crystal structures of annexin V complexes with calcium and phospholipid head group analogs indicate that calcium ions are involved directly in the attachment to membranes (11).

**Figure 1.** Ribbon diagram of annexin V (10). N-terminus as labeled, calcium ions as filled spheres.



#### 4. Membrane-Bound Annexins

The membrane-bound structures of two annexins have been investigated by [electron microscopy](#), and there is little evidence that the protein penetrates into the lipid bilayer. The molecular structures of membrane-bound annexins are essentially the same as in their crystals, except that the molecules reorient themselves so that their calcium-binding sites are in contact with the membrane ([8](#), [12](#), [13](#)). Several annexins exhibit calcium-dependent self-association and/or form two-dimensional arrays on membrane surfaces, a property that may underlie their biological functions ([14](#), [15](#)).

Structure-based mechanisms for annexin function have been proposed, but are not fully established. Two hypotheses have been presented to explain the annexin-induced calcium channel activity observed *in vitro*. In the “microscopic electroporation” model described for annexin V, the peripheral binding of the protein changes the electrostatic properties of the membrane ([16](#)). The calcium ion is translocated to the **cytosol** through a putative central pore in the annexin molecule. In an alternative model of calcium channel activity, based on hydra annexin XII, the protein hexamer inserts into the bilayer and creates a transmembrane structure resembling an inverted micelle ([17](#), [18](#)). Experimental data have been offered to support both hypotheses, but the mechanism and physiological relevance of the channel activity remains controversial. Mechanistic models in which annexins restrict access to membrane phospholipids and/or surfaces have been proposed to explain their inhibition of phospholipase A<sub>2</sub> ([19](#)) and thrombin ([14](#)), and various models have been suggested for processes, such as vesicle aggregation and membrane fusion. The *in vivo* functions of annexins remain under intensive investigation.

#### Bibliography

1. K. Hertogs et al. (1993) *Virology* **197**, 549–557.
2. J. F. Wright, A. Kurosky, and S. Wasi (1994) *Biochem. Biophys. Res. Commun.* **198**, 983–989.
3. K. von der Mark and J. Mollenhauer (1997) *Cell. Mol. Life Sci.* **53**, 539–545.
4. C. Y. Chung and H. P. Erickson (1994) *J. Cell. Biol.* **126**, 539–548.
5. K. A. Hajjar and J. S. Menell (1997) *Ann. N.Y. Acad. Sci.* **811**, 337–349.

6. K. Kojima et al. (1996) *J. Biol. Chem.* **271**, 7679–7685.
7. G. Kassam et al. (1997) *J. Biol. Chem.* **272**, 15093–15100.
8. J. Benz et al. (1996) *J. Mol. Biol.* **260**, 638–643.
9. H. Kawasaki, A. Avila-Sakar, C. E. Creutz, and R. H. Kretsinger (1996) *Biochim. Biophys. Acta* **1313**, 277–282.
10. A. Rosengarh, V. Gerke, and H. Luecke (2001) *J. Mol. Biol.* **306**, 489–498.
11. M. A. Swairjo et al. (1995) *Nat. Struct. Biol.* **2**, 968–974.
12. A. Oloffson, V. Mallouh, and A. Brisson (1994) *J. Struct. Biol.* **113**, 199–205.
13. D. Voges et al. (1994) *J. Mol. Biol.* **238**, 199–213.
14. H. A. M. Andree et al. (1992) *J. Biol. Chem.* **267**, 17907–17912.
15. C. Pigault et al. (1994) *J. Mol. Biol.* **236**, 199–208.
16. P. Demange et al. (1994) *Trends Biochem. Sci.* **19**, 272–276.
17. H. Luecke et al. (1995) *Nature* **378**, 512–515.
18. S. E. Moss (1995) *Nature* **378**, 446–447.
19. F. F. Davidson and E. A. Dennis (1989) *Biochem. Pharmacol.* **38**, 3645–3651.

### Suggestions for Further Reading

20. J. Mollenhauer and others (1997) *Annexins*, *Cell. Mol. Life Sci.* **53**, 506–555. A multiauthor review with nine concise articles on selected topics, including some material not reviewed elsewhere (eg, plant annexins). Up-to-date reviews on annexin genetics and molecular structure; phosphorylation; collagen-binding; and the roles of annexins in apoptosis, secretion, cancer, and autoimmune diseases.
21. B. A. Seaton, ed. (1996) *Annexins: Molecular Structure to Cellular Function*, R. G. Landes, Austin TX. Sixteen chapters covering a wide range of topics, including comprehensive reviews on annexin gene and molecular structure; biology of annexins I and XI; annexins in phagocytic leukocytes; nematode annexins; annexin binding to phospholipids and lipid assemblies; roles in membrane trafficking; calcium-independent annexin functions, including annexin/protein interactions; annexin functions involving phospholipid membrane asymmetry and clinical applications; and current experimental approaches to understanding annexin function.
22. P. Raynal and H. B. Pollard (1994) *Annexins: A novel family of calcium- and membrane-binding proteins in search of a function*, *Biochim. Biophys. Acta* **1197**, 63–93. A comprehensive survey covering most aspects of putative annexin function.
23. M. A. Swairjo and B. A. Seaton (1994) *Annexin structure and membrane interactions: a molecular perspective*, *Ann. Rev. Biophys. Biomol. Struct.* **23**, 193–213. Review with emphasis on structural and biophysical aspects of annexin-membrane associations and effects on lipid bilayer.

### Anomalous Dispersion

If an atom is hit by an X-ray beam, as in [X-ray crystallography](#), it scatters the beam in all directions. The scattered radiation can have the same wavelength as the primary beam (Rayleigh or coherent or elastic scattering) or a longer wavelength (Compton or incoherent or inelastic scattering). For diffraction, only the coherent part of the scattering is of interest. Incoherent scattering simply increases the background.

Electrons in an atom are bound by the nucleus and are, in principle, not free electrons. However, they can be regarded as such if the frequency of the incident radiation  $\omega$  is large compared with the natural absorption frequencies  $\omega_n$  of the scattering atom, or if the wavelength of the incident radiation is short compared with the absorption edge wavelength. This is normally true for the light atoms but not for the heavy atoms (Table 1).

**Table 1. The Position of the Ka-Edge for Some Elements**

| Atomic number | 6     | 16    | 26    | 34    | 78    |
|---------------|-------|-------|-------|-------|-------|
| Element       | C     | S     | Fe    | Se    | Pt    |
| Ka edge (Å)   | 43.68 | 5.018 | 1.743 | 0.980 | 0.158 |

If the electrons in an atom are regarded as free electrons, the atomic scattering amplitude (atomic scattering factor in units of electron scattering) is a real quantity  $f$  because the electron cloud is centrosymmetric. If they are not free electrons, the atomic scattering factor becomes an imaginary quantity, the scattering amplitude per electron:

$$f(\text{per electron}) = \frac{E_0 e^2}{mc^2} \times \frac{\omega^2}{\omega^2 - \omega_n^2 - i\kappa_n \omega} \quad (1)$$

where  $E_0$  is the amplitude of the electric vector of the incident beam, and  $\kappa_n$  is a damping factor for the  $n$ th orbit (K or L or  $1/4$ ). For  $\omega > \omega_n$ , Eq. (1) approaches  $f = E_0 e^2 / mc^2$ , the scattering amplitude of a free electron (1-3).

In practice, the complex atomic scattering factor, called  $f_{\text{anomalous}}$ , is separated into three parts:  $f_{\text{anomalous}} = f + f' + if''$ .  $f$  is the contribution to the scattering if the electrons were free electrons,  $f'$  is the real part of the correction to be applied for non-free electrons, and  $f''$  is the imaginary correction.  $f + f'$  is the total real part of the atomic scattering factor. Values for  $f$ ,  $f'$ , and  $f''$  are always given in units equal to the scattering by one free electron and are listed in Ref. 3. Because the anomalous contribution to the atomic scattering factor is mainly due to the electrons close to the nucleus, the value of the correction factors diminishes slowly as a function of the scattering angle, slower than for  $f$ .

Anomalous scattering causes a violation of Friedel's law:  $I(h k \ell)$  is no longer equal to  $I(h k \bar{\ell})$ . This fact can be used profitably for determining the absolute configuration (4). Moreover, it can assist in the structural determination of proteins.

#### Bibliography

1. R. W. James (1965) *The Optical Principles of the Diffraction of X-rays*, G. Bell and Sons, London, p. 135.
2. H. Hönl (1933) *Ann. der Physik*, 5. Folge **18**, 625–655.
3. International Union of Crystallography (1995) *International Tables for Crystallography*, Vol. C (A. J. C. Wilson, ed.), Kluwer Academic Dordrecht, Boston, London.

4. J. M. Bijvoet, A. F. Peerdeman, and A. J. van Bommel (1951) *Nature* **168**, 271–271.

### Suggestions for Further Reading

5. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Antennapedia Complex

The antennapedia complex (ANT-C) is a group of tightly linked **genes** that is found in the proximal portion of the right arm of the third chromosome of *Drosophila melanogaster*. The most prominent members of the complex produce striking **homeotic** transformations in adult flies carrying [mutations](#) in these genes. Thus mutations in the *Antennapedia* (*Antp*) locus cause a transformation of the antenna of the adult fly into a leg, while lesions in the *proboscipedia* (*pb*) gene cause the adult mouth parts to develop into legs rather than the normal palps used in feeding. The existence of the homeotic ANT-C was originally proposed based on the tight linkage of the *proboscipedia* (*pb*), *Sex combs reduced* (*Scr*), and *Antennapedia* (*Antp*) loci. Subsequent genetic analyses have shown that two other homeotic loci, *labial* (*lab*) and *Deformed* (*Dfd*), are also members of the complex. The homeotic loci of the ANT-C are involved in the specification of segmental identity in the posterior head (gnathocephalic) and anterior thoracic regions of the embryo and adult. Moreover, the linear order of the homeotic loci in the complex, *lab*, *pb*, *Dfd*, *Scr*, and *Antp*, corresponds to the anterior–posterior order of altered segments (intercalary, mandibular, maxillary, labial, and thoracic) found in animals bearing mutations in each of the resident loci. Taken together, the results of mutational analyses indicate that members of the complex are necessary to repress head development in the thorax (*Antp*) and elicit normal segmental identity in the anterior thorax (*Scr*) and posterior head (*Scr*, *Dfd*, *pb*, and *lab*).

A similar group of homeotic genes called the [Bithorax Complex](#) (BX-C) is found more distally on the third chromosome. This set of three homeotic genes (*Ultrabithorax*, *abdominal-A*, and *Abdominal-B*) acts in a similar fashion to the ANT-C, but in the posterior of the thorax and in the abdomen. The ANT-C is distinguished from the BX-C not only by virtue of the domain of action of its homeotic loci (anterior versus posterior) but also by the presence of loci that are not overtly homeotic in character. Two of these, *fushi tarazu* (*ftz*) and *zerknüllt* (*zen*), have been shown to affect segment enumeration (*ftz*) and the formation of dorsal structures (*zen*) in the early embryo. A third nonhomeotic gene is [bicoid](#) (*bcd*). Mutations in this locus result in female sterility and maternal effect lethality. Eggs laid by *bcd* females fail to develop normal anterior ends and instead produce mirror-image duplications of structures normally produced at the posterior terminus of the embryo.

In addition to these genetically defined loci, several other genes have been found in the ANT-C by molecular mapping. The first of these is a cluster of cuticle-protein-related genes that map between the *lab* and *pb* loci. Eight small (about 1 kbp) [transcription](#) units make up the cluster, and all have sequence similarities to known cuticle protein genes. These genes (*cc1* through *cc8*) are also apparently regulated by [ecdysone](#) in [imaginal discs](#). Deletion of the entire cluster has no apparent effect on the development or cuticle morphology of embryos, larvae, or adults. The second molecularly identified gene is the *Amalgam* (*Ama*) transcription unit. The encoded protein places the gene in the [immunoglobulin](#) superfamily and, like the cuticle cluster, the locus can be deleted from the [genome](#) with no discernible effect on the organism. Finally, there is the *zen2* or *z2* transcription unit that resides immediately adjacent to the *zen* gene. This locus is similar in structure and sequence to *zen*, but like *Ama* and the *cc* genes has no discernible function.

The entire complex has been **cloned** and sequenced and shown to cover 335 kbp of genomic DNA. The most distal transcription unit is *Antp*, which covers the distalmost 100 kbp of the complex and is made up of eight exons. Proximally, the next 75 kbp contain the *Scr* and *ftz* loci. The distal 50 kbp of this interval contain sequences necessary for *Scr* expression, as well as the two exons of the *ftz* locus and its associated regulatory elements. The proximal 25 kbp contain the three identified exons of the *Scr* transcription unit. The five exons of the *Dfd* gene are found in the central portion of the next-most-proximal 55-kbp interval. The *Dfd* transcription unit covers only 11 kbp of this region, and the 20-kbp interval flanking the gene proximally is the location of [cis-acting](#) regulatory elements for the locus. The next 25-kbp interval contains four of the nonhomeotic transcription units that help distinguish the ANT-C and BX-C. The distalmost is *Ama*, next *bcd*, and finally *zen* and *z2*. The *z2*, *zen*, and *Ama* transcription units are all relatively small (1 to 2 kbp) and comprise two exons each. The *bcd* gene is somewhat larger (3.6 kbp) and is made up of four exons. Immediately proximal to the *z2* transcription unit (about 1 kbp from its 3' end) is the 5' end of *pb*, which extends over the next 35 kbp of genomic DNA and contains nine exons. The next 25 kbp of the complex contain the cuticle cluster and its eight identified transcription units. The final 25 kbp are the sites of the *lab* gene, which is made up of three exons. Despite the nonhomeotic nature of three of the smaller transcription units (*zen*, *bcd*, and *ftz*) resident in the complex, these loci are tied to the larger homeotic genes of the region by the nature of their protein products. All five of the large homeotics (*Antp*, *Scr*, *Dfd*, *pb*, and *lab*), and the three small genes, have a **homeobox** motif, and their protein products are found in the **nuclei** of the cells in which they are expressed. Thus eight of the genes in the ANT-C encode regulatory proteins that act as [transcription factors](#). The *z2* gene also contains a homeobox; however, the biological significance of the gene is not known, as deletions of this transcription unit have no discernible effect. The cuticle-like genes and *Ama* do not contain a homeobox.

The reasons for the clustering of these developmentally significant loci of similar function are not known. The existence of common or overlapping regulatory elements, the need to insulate regulatory sequences from chromosomal position effect, and the possibility of higher-order [chromatin](#) structures for proper expression have all been proposed. Whatever the reason, the homeotic complex structure has a long evolutionary standing. Similar clusters are found in vertebrates, an observation consistent with a very early origin of these genes, probably predating the separation of protostomes and deuterostomes.

#### Suggestions for Further Reading

- M. Affolter, A. Schier, and W. J. Gehring (1990) Homeodomain proteins and the regulation of gene expression, *Curr. Opin. Cell Biol.* **2**, 485–495.
- M. D. Biggin and McGinnis W. (1997) Regulation of segmentation and segmental identity by *Drosophila* homeoproteins: the role of DNA binding in functional activity and specificity, *Development* **124**, 4425–4433.
- A. Dorn, M. Affolter, W. J. Gehring, and W. Leupin (1994) Homeodomain proteins in development and therapy, *Pharmacol. & Therap.* **61**, 155–183.
- D. Duboule and G. Morata (1994) Colinearity and functional hierarchy among genes of the homeotic complexes, *Trends Genet.* **10**, 358–364.
- R. Finkelstein and N. Perrimon (1991) The molecular genetics of head development in *Drosophila melanogaster*, *Development* **112**, 899–912.
- W. J. Gehring, M. Affolter, and T. Burglin (1994) Homeodomain proteins, *Ann. Rev. Biochem.* **63**, 487–526.
- G. Gellon and W. McGinnis (1998) Shaping animal body plans in development and evolution by modulation of Hox expression patterns, *BioEssays* **20**, 116–125.
- G. Jurgens and V. Hartenstein (1993) "The terminal regions of the body pattern". In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, pp. 687–746.

- T. C. Kaufman, M. A. Seeger, and G. Olsen (1990) Molecular and genetic organization of the Antennapedia gene complex of *Drosophila melanogaster*, *Adv. Genet.* **27**, 309–362.
- P. A. Lawrence and G. Morata (1994) Homeobox genes: their function in *Drosophila* segmentation and pattern formation, *Cell* **78**, 181–189.
- R. S. Mann (1997) Why are Hox genes clustered? *BioEssays* **19**, 661–664.
- W. McGinnis and R. Krumlauf (1992) Homeobox genes and axial patterning, *Cell* **68**, 283–302.
- G. Morata (1993) Homeotic genes of *Drosophila*, *Curr. Opin. Genet. Dev.* **3**, 606–614.
- A. Popadic, A. Abzhanov, D. Rusch, and T. C. Kaufman (1998) Understanding the genetic basis of morphological evolution: the role of homeotic genes in the diversification of the arthropod bauplan, *Int. J. Dev. Biol.* **42**, 453–461.
- B. T. Rogers and T. C. Kaufman (1997) Structure of the insect head in ontogeny and phylogeny: a view from *Drosophila*, *Int. Rev. Cytol.* **174**, 1–84.

## Antibiotic Resistance

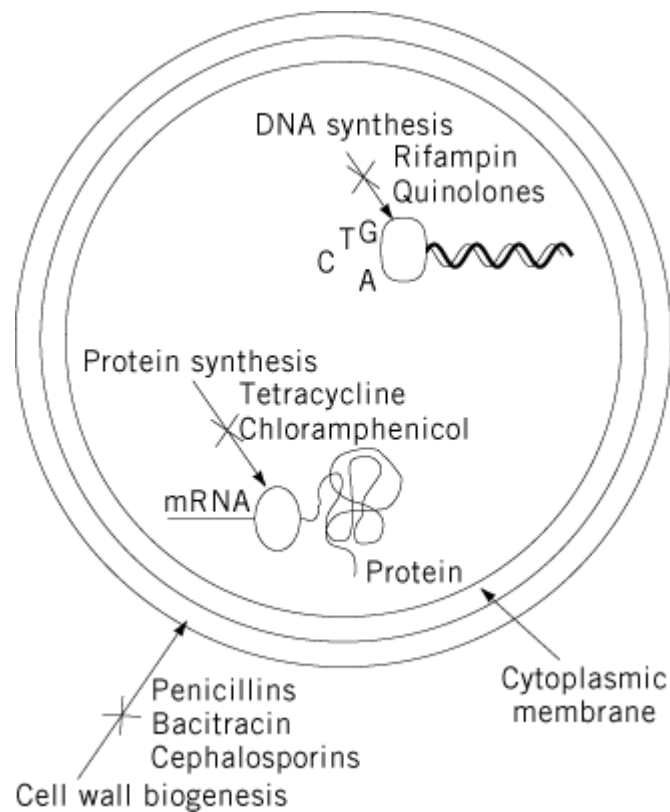
In the 1940s, the clinical use of antibiotics first curbed the widespread threat of deadly bacterial infection. These drugs effectively inhibited bacterial growth that had gone unchecked for decades. Antibiotics did not, however, eradicate the threat of bacterial infection. In fact, the widespread use of antibiotics gave a selective advantage to bacteria that had antibiotic resistance. Many strains of bacteria have developed antibiotic resistance, or insensitivity to antibiotic drugs, in response to antibiotic selection pressures. Now, bacteria employ a myriad of resistance mechanisms to circumvent the best efforts of antibiotic researchers and clinicians. Only with the development of new and potent antibiotics and the appropriate use of existing antibiotics will researchers regain control over this resilient lifeform, bacteria.

### 1. A Historical Perspective

The development of antibiotics as therapeutic agents began in the late 1930s to combat the most common cause of death, infectious disease. In the preantibiotic era, any infection could prove mortal. Subsequently, over 150 different antibiotics have been synthesized or discovered, and these drugs are used to treat bacterial infections, such as pneumonia, malaria, and tuberculosis (1, 2).

Antibiotics are a collection of natural products and synthetic compounds that kill bacteria. Naturally occurring antibiotics are isolated from molds, yeasts, and bacteria. These organisms use antibiotics as defense mechanisms to kill other bacteria. Alternatively, synthetic antibiotics are developed by understanding the architecture and function of bacteria (see Fig. 1). Some bacteria have cell walls, and many effective antibiotics, such as [penicillin](#), bacitracin, and cephalosporin, inhibit the synthesis of this cell wall. The bacterial machinery for **protein biosynthesis** differs from that of many host organisms and, therefore, is a good target for antibiotics, such as tetracycline and [chloramphenicol](#). Additionally, antibiotics, such as **rifampin** and quinolones, specifically inhibit [DNA replication](#) in bacteria (1).

**Figure 1.** Mechanisms of action of some common antibiotics. A schematic of some basic functions of bacterial cells to which antibiotics are targeted. The X indicates an inhibition of that function by the antibiotics noted (see text for details).



Antibiotics were considered the “wonder drugs” of their time, and in retrospect, this highly favored opinion resulted in their overuse. Antibiotics were commonly prescribed by physicians to cure and to appease their patients. Some patients would have genuine bacterial infections, for which antibiotic treatment is appropriate, whereas others would request antibiotics for viral infections that are not susceptible to these drugs. In addition, individuals who have suppressed immune systems, such as AIDS patients or organ transplants patients, would harbor bacteria that acquire resistance more easily. Antibiotics were also used prophylactically in agriculture and aquaculture industries to keep livestock healthy (3).

By the late 1960s, infectious disease appeared to be under control by the use of a variety of antibiotics. However, antibiotics simply depressed the propagation of bacteria. They did not eradicate it. Nonetheless, research in human health changed its focus from infectious diseases to chronic diseases, and new antibiotics were no longer being developed (2, 4). Many microbiologists warned the human health community that bacteria were potent, infectious pathogens that should not be underestimated. Bacteria have survived for more than three billion years despite numerous environmental changes on earth, and industrial wastes, insecticides, and herbicides. Obviously, the mechanisms bacteria use to survive and adapt to these adversities were very effective.

## 2. The Origins of Antibiotic Resistance

Even before the first clinical application of antibiotics, antibiotic resistance, the ability of bacteria to evade the deleterious effects of antibiotics, was postulated. In 1940, Abraham and Chain identified a bacterial enzyme that inactivates one of the first antibiotics, penicillin. Then, any bacterium that produces this enzyme would be resistant to penicillin (5). Moreover, microorganisms that use an antibiotic as defense mechanisms would require immunity to that antibiotic. This inherent resistance to a particular antibiotic was defined as a naturally-occurring trait called intrinsic resistance. A few strains of bacteria with intrinsic resistance to a particular antibiotic would not constitute a clinical threat because many diverse antibiotics are available. However, the ability of bacteria to propagate



this antibiotic resistance to other strains of bacteria had been underestimated.

The exchange of genetic information between bacteria of the same strain is a common, yet typically slow process. Mechanisms of exchange include (1) **conjugation**—a single DNA strand from one bacterium enters another bacterium and is replicated as a part of that genome, (2) **transduction**—foreign DNA is introduced into bacteria by transducing bacteriophages, and (3) transformation—autonomously replicating circular DNA plasmids are obtained by bacteria (1). Originally, it was believed that these genetic exchange mechanisms were restricted to bacteria of the same strain. However, a new method of gene exchange, using integrons, was recently identified (6, 7).

Integrons are independent, [mobile elements](#) that encode **genes** for protein functions, and encode additional DNA to guarantee the integron's expression and integration into the bacterial [genome](#). Integrons effectively generate widespread antibiotic resistance by donating antibiotic resistance genes to any strain of bacteria. Ironically, it is widely believed that integrons evolved only recently in response to antibiotic selection pressure. In other words, the use of antibiotics advanced the widespread occurrence of antibiotic resistance.

In addition to the acquisition of genes by the exchange of genetic information, bacteria also have a high mutation rate that allows them to respond to the selective pressure of antibiotics by using their own genome. For example, if bacteria were subjected to tetracycline, a random mutation in the 30S [ribosome](#) to weaken tetracycline binding would be advantageous and, therefore, would be perpetuated by the survival of the tetracycline-resistant bacteria (1). It has also been postulated that **housekeeping genes**, like acyltransferases, may have mutated to gain the ability to modify and inactivate aminoglycoside antibiotics (8-10).

The widespread phenomenon of antibiotic resistance has developed from the promiscuity of bacteria and their genomes. Initially, intrinsic antibiotic resistances were isolated incidents, but the threat of antibiotics has been readily circumvented using the acquisition of antibiotic resistance genes and the high mutational frequency of individual bacteria.

### 3. Mechanisms of Antibiotic Resistance

Using both newly acquired genes and their own mutated genes, bacteria utilize three basic mechanisms to support antibiotic resistance. Enzymes that degrade antibiotics inside the cell are key players in antibiotic resistance. Bacteria can also alter their permeability barriers to keep antibiotic concentrations below toxic levels inside the cell. Furthermore, the cellular targets of antibiotics can be modified to evade the effects of the antibiotics. These mechanisms of antibiotic resistance are distinct, but all are obtained through the acquisition or mutation of genes.

#### 3.1. Enzymatic Inactivation of Antibiotics

Degradative enzymes are a common mechanism by which bacteria become resistant to antibiotics. Such enzymes chemically modify antibiotics so that they no longer function. The genes for these degradative enzymes are obtained by acquisition of exogenous genes or mutation of endogenous genes. The expression of these genes also governs the level of antibiotic resistance in bacteria.

[b-Lactamases](#) are common examples of degradative enzymes that generate antibiotic resistance. b-Lactamases inactivate b-lactam antibiotics, a structurally similar group of penicillin-like antibiotics, all of which have a b-lactam ring structure. b-Lactamases are typically grouped into two major classes, penicillinases and cephalosporinases, based on their substrate affinity (11, 12). In addition to b-lactamases, other enzymes also degrade different antibiotics, such as the aminoglycosides, gentamycin, tobramycin, and amikacin.

All b-lactam antibiotics function similarly. Their b-lactam ring structure inhibits the final step of bacterial cell wall synthesis. Bacterial cell walls are constructed of alternating *N*-acetylglucosamine and *N*-acetyl-muramic acid residues that form long peptidoglycan chains, and the final step in cell

wall synthesis involves the enzymatic crosslinking of these peptidoglycan chains by a transpeptidase. Because the b-lactam bond resembles a portion of the peptidoglycan chains, this transpeptidase can mistake a b-lactam antibiotic for its natural substrate and hydrolyze the b-lactam bond. This hydrolysis covalently links the b-lactam drug to the transpeptidase and renders it nonfunctional (13) (see also [Penicillin-Binding Proteins](#)).

Although b-lactamases are effective at degrading some antibiotics, their mere presence is not sufficient to cause clinically relevant antibiotic resistance. In fact, these enzymes are found ubiquitously in almost all bacteria, and in some blue-green algae and mammalian tissues. b-lactamases must be present in sufficient quantities to degrade the b-lactam antibiotics effectively before they inhibit cell wall synthesis. The cellular concentration of a b-lactamase depends on its gene expression, and b-lactam antibiotics are inducers of b-lactamase expression. Furthermore, particular b-lactamases have variable affinities for b-lactam antibiotics. Therefore, the degree of antibiotic resistance due to b-lactamases is based on a combination of the ability of the b-lactam antibiotic to induce the expression of b-lactamase, and its ability to be a substrate for b-lactamase (14).

As researchers began to understand this mechanism of antibiotic resistance, more effective b-lactam antibiotics were developed. The early cephalosporins, like penicillin and amoxicillin, are extremely sensitive to b-lactamases, because these b-lactam antibiotics are potent inducers of b-lactamase expression and good substrates for the b-lactamase. In contrast, the more recently developed antibiotic, imipenem, is a strong inducer of b-lactamase expression but maintains its antibiotic activity because it is a poor substrate for most b-lactamases (15). In addition to these new antibiotics, combination therapies are also being implemented to combat antibiotic resistance. Such therapies are comprised of a b-lactam antibiotic together with b-lactamase inhibitors, like clavulanic acid, sulbactam, and tazobactam (16).

### 3.2. Altered Permeability Barriers: Pore Proteins and Efflux Systems.

The bacterial cell [membrane](#) is the major permeable barrier separating the outside of the cell from the inside. The fluidity of the membrane is generally balanced to include most nutrients, while excluding many [toxins](#). Adjusting this fluidity impedes the function of the membrane. Therefore, bacteria cannot protect themselves by changing the fluidity of their membrane. Instead, bacteria have additional structures that surround the cytoplasmic membrane or form pores through it. The alteration of these structures to exclude antibiotics is another mechanism of antibiotic resistance.

[Gram-positive](#) and [Gram-Negative Bacteria](#) have distinct structures that surround their cytoplasmic membranes. Most Gram-positive bacteria have thick cell walls that are mechanically quite strong, although very porous. Although the cell wall helps Gram-positive bacteria retain their shape, it does not exclude most antibiotics and, therefore, is not a good barrier. Thus, Gram-positive bacteria are relatively susceptible to the influx of antibiotics. Alternatively, a more effective barrier, a second lipid bilayer or membrane, surrounds Gram-negative bacteria. This outer membrane is partially composed of a lipid, lipopolysaccharide (LPS), that is not commonly found in cytoplasmic membranes. The distinguishing feature of LPS is its decreased fluidity, which makes the LPS bilayer an efficient barrier that prevents the permeation of most **hydrophobic** antibiotics into Gram-negative bacteria (17).

Enveloped by effective barriers, bacteria use pore-forming proteins, called [porins](#), to obtain nutrients from outside the cell. Porins are transmembrane proteins that function as nonspecific, aqueous channels, and allows nutrients to diffuse across the membrane. Porins generally exclude antibiotics because they are narrow and restrictive. Most antibiotics are large, uncharged molecules that cannot easily traverse the narrow porin channels that are lined with charged amino acid residues. However, some antibiotics enter the bacteria through porins, and the deletion or alteration of these porins to exclude particular antibiotics is linked to antibiotic resistance.

Because bacteria cannot develop barriers that are impermeable to all molecules, some toxins do

diffuse into bacteria along with nutrients. Therefore, bacterial cell membranes also contain transport proteins that cross the membranes and use energy to remove toxins. They are called active efflux systems, and some are directly identified as another significant cause of antibiotic resistance.

Many active efflux systems resemble other transport proteins that catalyze the efflux of common, small molecules, like glucose or cations, and it is likely that mutation has modified them to transport antibiotics. Based on their overall structure, mechanism, and sequence homologies, these transport proteins are classified into four families: (1) the major facilitator family; (2) the resistance nodulation division family; (3) the staphylococcal multidrug resistance family; and (4) the ATP-binding cassette (ABC) transporters. Of these four families, only the ABC transporters use the chemical energy generated from the hydrolysis of ATP to drive molecules across the membrane. Members of the three other families use an electrochemical proton gradient, or [proton-motive force](#), as the source of energy (18, 19).

Some active efflux systems exclude a variety of unrelated toxins from the cell. These multidrug resistance (MDR) efflux systems in bacteria are comparable to those found in mammalian cells (see [Drug Resistance](#)). An example of a bacterial MDR efflux system is the Bmr transporter that transports drugs which have diverse chemical structures and physical properties and include cationic dyes, rhodamine-6G, [ethidium bromide](#), and the antibiotics netropsin, puromycin, and fluoroquinolone (20). Other MDR efflux systems characterized in bacteria include NorA in *Staphylococcus aureus*, MexB in *Pseudomonas aeruginosa*, and EmrB in *Escherichia coli*. If MDR efflux systems occur extensively in bacteria as the source of many antibiotic resistances, they pose a far more formidable challenge than more specific mechanisms of resistance.

### 3.3. Modification of the Antibiotic's Target

Antibiotics inhibit bacterial growth by inactivating different key proteins that are essential for bacterial survival (see Fig. 1). However, the antibiotic sensitivity of these target proteins can be altered. Typically, antibiotic targets are altered by reducing their affinity for the antibiotic. Bacteria accomplish this change in affinity several ways. Bacteria acquire exogenous DNA for a mutated target protein that no longer interacts with the antibiotic, yet retains the original target protein's function. Alternatively, bacteria's endogenous genes can be mutated to achieve the same end. In contrast, DNA for novel modifying enzymes can be acquired to alter the antibiotic target posttranslationally, reducing its affinity for the antibiotic.

Altering an antibiotic's target protein directly at the DNA level is a common mechanism of target modification. An example of this modification is the mutation of genes for [penicillin-binding proteins](#) (PBPs). PBPs are transpeptidases, previously discussed, that catalyze the final step in bacterial cell wall synthesis. These PBPs have high affinity for penicillin and its derivatives, and the binding of penicillin permanently inactivates PBPs. Originating from both endogenous and exogenous DNA sources, mutated PBPs can have a lower affinity for penicillin. Therefore, PBPs are resistant to the antibiotic, yet still provide a crucial function in bacterial cell wall synthesis (21, 22). Another example of target modification via mutated DNA is a single amino acid mutation in the quinolone resistance-determining region of the **DNA gyrase** gene, *gyrA*, that can provide up to a 20-fold increase in quinolone resistance (23).

In addition to using mutated antibiotic targets, bacteria can acquire new genes that produce proteins that, in turn, alter antibiotic targets. A well-studied example is the resistance of Staphylococci to **erythromycin**. It is known that Staphylococci have acquired genes to produce a protein that methylates a residue on the 23S ribosome. The 23S ribosome is the target of erythromycin, but methylated 23S ribosome has a low affinity for erythromycin. This exogenous gene is expressed and prevents the binding of erythromycin to the ribosomes, making the bacteria erythromycin-resistant (24).

## Bibliography

## Bibliography

1. H. C. Neu (1992) *Science* **257**, 1064–1073.
2. J. Travis (1994) *Science* **264**, 360–362.
3. M. Castiglia and R. A. J. Smego (1997) *J. Am. Pharm. Assoc.* **NS37**, 383–387.
4. G. H. Cassell (1997) *FEMS Immunol. Med. Microbiol.* **18**, 271–274.
5. E. P. Abraham and E. Chain (1940) *Nature* **146**, 837.
6. H. W. Stokes and R. M. Hall (1989) *Mol. Microbiol.* **3**, 1669–1683.
7. C. M. Collis, G. Grammaticopoulos, J. Briton, H. W. Stokes, and R. M. Hall (1993) *Mol. Microbiol.* **9**, 41–52.
8. T. Udou, Y. Mizuguchi, and R. J. J. Wallace (1989) *FEMS Microbiol. Lett.* **48**, 227–230.
9. K. J. Shaw et al. (1992) *Antimicrob. Agents Chemother.* **36**, 1447–1455.
10. P. N. Rather, E. Orosz, K. J. Shaw, R. Hare, and G. Miller (1993) *J. Bacteriol.* **175**, 6492–6498.
11. M. H. Richmond and R. B. Sykes (1973) *Adv. Microb. Physiol.* **9**, 31–88.
12. K. Bush (1989) *Antimicrob. Agents Chemother.* **33**, 259–276.
13. A. Tomasz (1979) *Annu. Rev. Microbiol.* **33**, 113–137.
14. D. M. Livermore (1993) *J. Antimicrob. Chemother.* **31** (suppl. A), 9–21.
15. J. Y. Jacobs, D. M. Livermore, and K. W. M. Davy (1984) *J. Antimicrob. Chemother.* **14**, 221–229.
16. K. Coleman et al. (1994) *J. Antimicrob. Chemother.* **33**, 1091–1116.
17. P. R. Cullis and M. J. Hope (1985) In *Biochemistry of Lipids and Membranes* (D. E. Vance and J. E. Vance, eds.), Benjamin and Cummings, New York, Chap. "2".
18. S. B. Levy (1992) *Antimicrob. Agents Chemother.* **36**, 695–703.
19. K. Lewis, D. C. Hooper, and M. Ouellette (1997) *ASM News* **63**, 605–610.
20. A. A. Neyfakh, V. E. Bidnenko, and L. B. Chen (1991) *Proc. Natl. Acad. Sci. USA* **88**, 4781–4785.
21. B. G. Spratt and K. D. Cromie (1988) *Rev. Infect. Dis.* **10**, 699–711.
22. J. M. Ghuyssen (1991) *Annu. Rev. Microbiol.* **45**, 37–67.
23. G. A. Jacoby and A. A. Medeiros (1991) *Antimicrob. Agents Chemother.* **35**, 1697–1704.
24. R. Leclercq and P. Courvalin (1991) *Antimicrob. Agents Chemother.* **35**, 1267–1272.

## Suggestions for Further Reading

25. A short news article: J. Davies (1996) Bacteria on the rampage, *Nature* **383**, 219–220.
26. Two well written reviews in an excellent anthology: 1. A. Bauernfeind and N. H. Georgopapadalou (1995) In *Drug Transport in Antimicrobial and Anticancer Chemotherapy* (N. H. Georgopapadalou, ed.), Dekker, New York, Vol. 17, pp. 1–19. 2. R. E. W. Hancock (1995) In *Drug Transport in Antimicrobial and Anticancer Chemotherapy* (N. H. Georgopapadalou, ed.), Dekker, New York, Vol. 17, pp. 289–306.
27. Three exhaustive articles in a single issue of *Science* devoted to antibiotic resistance: 1. J. Davies (1994) Inactivation of antibiotics and the dissemination of resistance Genes, *Science* **264**, 375–382. 2. H. Nikaido (1994) Prevention of drug access to bacterial targets: Permeability barriers and active efflux systems, *Science* **264**, 382–388. 3. B. G. Spratt (1994) Resistance to antibiotics mediated by target alterations, *Science* **264**, 388–393.

## Antibody

Antibodies are specific [proteins](#), termed [immunoglobulins](#), that are produced by [B cells](#) upon stimulation with [antigens](#), which may be proteins, polysaccharides, nucleic acids, and so on, either in soluble form or as part of complex cellular organisms, such as bacteria, parasites, **viruses**, eukaryotic cells of animal or plant tissues, and pollens. In fact, the immune system is so built as to have the potential to make antibodies against any [macromolecule](#) of the living world. It was reported at the end of the nineteenth century by von Behring and Kitasato that a group of guinea pigs immunized against a sublethal dose of [diphtheria toxin](#) became resistant to diphtheria, whereas another group that received tetanus toxin resisted further challenge with a normally lethal dose of tetanus toxin. This key experiment showed that the protection was specific toward the immunizing agent. Thus, one master word in immunology is specificity, and this immediately raises the problem of what immunologists call the [repertoire](#). How is it possible to make the billions of different molecules that are potentially necessary to assume specific recognition of an obviously astronomical number of potential antigens? What structures are recognized on an antigen? What are the structural bases for antibody specificity? All these are basic questions that have been solved progressively during the past 40 years.

The protein nature of antibodies was established by M. Heidelberger in 1928, by showing that protein nitrogen was present in a specific antigen–antibody precipitate in which the antigen was a polysaccharide, which must have originated from the antibody part of the precipitate. Progress in the knowledge of the structure of antibodies paralleled the emergence of new technological tools in biochemistry and in molecular biology. [Analytical ultracentrifugation](#) performed by Kabat and Tiselius at the end of the 1930 indicated that antibodies were found in several classes of sizes, mostly with [sedimentation coefficients](#) of between 19S and 7S. **Electrophoretic** characterization of serum, also studied by Kabat in the same period, revealed that antibodies were located in the [globulin](#) part of the spectrum, mostly in the slowest moving g-globulin fraction. This name remained in use for years until it was replaced by a more generic one, that of [immunoglobulins](#). Immunoglobulins are thus defined as the immune proteins synthesized specifically by B cells. They are expressed in two forms: (1) as integral [membrane proteins](#) at the surface of B lymphocytes, where they represent the B-cell receptor (BCR), by analogy with the [T-cell receptor](#) (TCR) expressed at the surface of [T cells](#); and (2) as a soluble form, secreted in the bloodstream by the plasma cells, which represent the terminal stage of differentiation of the B lineage. This soluble form is what immunologists call circulating antibodies, or simply antibodies.

A main difficulty encountered in studying antibody structure was due to their extraordinary heterogeneity, which deserves some comments. As a rule, when an antigen such as [hemoglobin](#) or [serum albumin](#), or in fact any protein, is used as an [immunogen](#), a large array of antibodies with discrete fine specificities is synthesized. At the surface of a protein, many discrete regions are recognized by as many distinct antibodies. These regions are called *antigenic determinants*, or [epitopes](#). For a protein of molecular weight 50 kDa, there are 10, 20, or more epitopes that could be identified by distinct antibodies. But antibody heterogeneity is not limited to the mosaic of epitopes. It may be shown that in fact one given epitope may be recognized by slightly different antibodies, and this contributes an extraordinary multiplication of the heterogeneity. This results from the clonal organization of lymphocytes, as originally proposed by Burnet (1). B lymphocytes (and T cells as well) are organized as discrete clones, each of which makes one and only one type of antibody molecule. In humans, the average population of lymphocytes at any given time is of the order of  $10^{12}$ , of which approximately one-fifth are B cells and four-fifths are T cells. These B cells make about  $10^{20}$  molecules of circulating immunoglobulins, a number that can be calculated from the Ig concentration in serum (10 mg/mL). The size of one clone is largely dependent on the immunization state and can vary from one single resting cell to several thousand or more. What happens when an animal is being immunized is that all clones that happen to interact somewhat with any potential

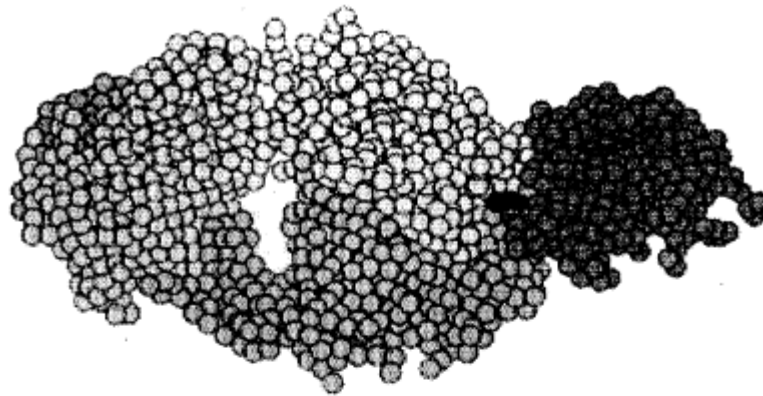
epitope of the administered antigen will proliferate, expand, and synthesize as many different antibodies, thereby contributing to a very great heterogeneity. This of course made it very difficult, if not impossible, to study the fine structure of antibody molecules, especially when considering their amino acid sequences. A bias was brought to this by the existence of multiple myeloma in humans, or its equivalent as experimentally induced plasmacytoma in the mouse. These are lymphoproliferative disorders that affect plasma cells. As a result of the malignancy, one clone will expand specifically so that, very rapidly, it will become the major B-cell clone expressed in the body. The corresponding immunoglobulin will be secreted in excessively large amounts, and this provides a homogeneous material available for structural studies. An obvious drawback with this material is that the antibody specificity is generally unknown, which prevents its use for the study of **antigen–antibody interactions**. This difficulty was turned around when Kohler and Milstein (2) succeeded in making somatic hybrids between myeloma cells and normal B lymphocytes from hyperimmunized mice. These **hybridomas** combined the properties of myeloma cells, which can grow indefinitely in culture or upon transplantation in syngeneic strains of mice, with the clones of immunized lymphocytes that could be selected for a known antigenic specificity, leading to the fantastic expansion of “**monoclonal antibodies**” (mAbs). Monoclonal antibodies were not only suitable for detailed structural analysis of antibodies, including determination of the three-dimensional **protein structure** of antibody–antigen complexes. They also allowed access, through the hybridoma cells, to isolation of **messenger RNA** encoding Ig polypeptide chains and the corresponding **complementary DNA**, and thus to the complete gene organization of the three main Ig gene loci.

The basic structure of antibodies, elucidated from the pioneer work of Porter in England (3) and Edelman in the United States (4), is given by that of the **IgG** molecule—known initially as 7S-globulin. It is a symmetrical molecule, containing two identical heavy chains (H, 52 kDa) and two identical light chains (L, 23 kDa), which are either **kappa (k) or lambda (l) chains**. Each H–L pair contains one combining site for the antigen, so the conventional  $H_2L_2$  IgG molecule is bivalent. The symmetry of the molecule is in agreement with the **clonal selection theory**, which postulated that each B cell expressed one and only one antibody specificity. This basic model could be extrapolated to other antibody classes, which differ from each other in the nature of their heavy chains and the degree of polymerization of the basic  $H_2L_2$  unit. Five classes have thus been described in higher vertebrates: **IgG**, **IgM**, **IgA**, **IgD**, and **IgE**, having the corresponding g, m, a, d, and e heavy chains, respectively. IgM are expressed as  $H_2L_2$  monomers at the cell surface of B cells, but are pentamers in the serum. IgA is mostly a dimer, whether other classes remain as monomers. The diversification of classes allows the antibodies to assume two types of function: (1) antigen recognition, common to all classes, and (2) biological or effector functions, such as **complement** fixation, active transplacental transfer, or fixation to various cell types, which amplify the action of antibodies and generally favor the elimination of a pathogen.

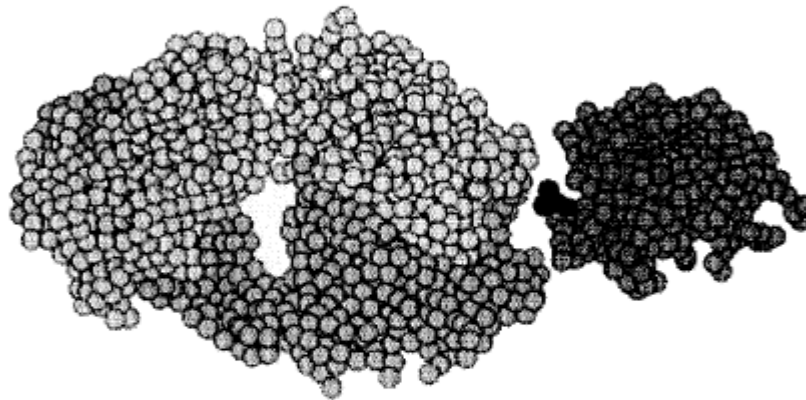
The functional duality of antibody molecules was clearly established in the late 1950 by Porter (5), who succeeded in cleaving the IgG antibody by **papain**, a proteinase, which led to the isolation of two identical fragments that were antigen binding, or Fab, and one Fc fragment that did not recognize the antigen but fixed the first component of the complement cascade. Structural support for this organization was clearly demonstrated by Hilschmann in 1965 (6), who reported the first amino acid sequence of two monoclonal human light chains. This is another example of the importance of monoclonal materials in this saga of the determination of the antibody structure. The work was performed on **Bence-Jones proteins**, which are free light chains isolated from the urine of patients with multiple myeloma. They are the result of a dysbalanced synthesis between the monoclonal heavy and light chains by the malignant plasma cells. Because of their smaller size, they pass the renal barrier and are thus easily isolated in large amounts from patient's urine. The major discovery was that the Bence-Jones light chains, or Ig light chains, were composed of one  $NH_2$ -terminal half of 110 amino acid residues that profoundly differed between two patients, whereas the  $COOH$ -terminal half, also of 110 residues, was identical. This gave clear evidence for the existence

of a huge structural diversity in antibodies, providing a unique basis for antibody specificity. Further extensive work indicated that heavy chains had also a variable region of similar size and a much longer constant region. The antibody combining site could thus be visualized as resulting from a contribution of both the  $V_H$  and the  $V_L$  regions. A more detailed analysis, computed by Kabat (7), indicated that within the **V variable regions**, subregions of hypervariability could be identified around positions 30, 50, and 100, which were later proven to participate directly in making up the antibody combining site. Three hypervariable regions, also called *complementarity determining regions* (CDR), were identified on each heavy and light chain, and their fine analysis revealed an extraordinary diversity that certainly could account for the expected huge repertoire of antibodies necessary to accommodate the potential repertoire of epitopes of the living world. Direct proof that the CDRs were implicated in antigen recognition came first from [affinity labeling](#) experiments, until a final confirmation was provided by [X-ray crystallography](#) analysis of crystals of antigen–antibody complexes. In fact, crystals were prepared from the Fab fragment of an antilysozyme antibody, combined with [lysozyme](#), because the complete IgG antibody molecule contains a floppy “hinge” region that prevents crystallization. This structure, obtained by the group of Poljak (8), presented in [Figure 1](#), indicates that about 20 amino acid residues of the antibody participate in binding the antigen, plus a similar number of residues from the antigen. On average, the area of interaction of the two partners is of the order of  $600 \text{ \AA}^2$ . Care should be taken in attempts to generalize this model to all other antigen-antibody pairs, because the size and shape of the antibody combining site varies greatly from one system to another. This variation is reflected in the wide variation of the association constants  $K_A$ , between  $10^5$  and  $10^{12} \text{ M}^{-1}$ . The interactions between antigen and antibody are exclusively non covalent and involve primarily [salt bridges](#), [van der Waals interactions](#), and [hydrogen bonds](#). The contributions of **enthalpy** and **entropy** vary immensely from one system to another, stressing again, if needed, the fantastic diversity potential of antibody molecules.

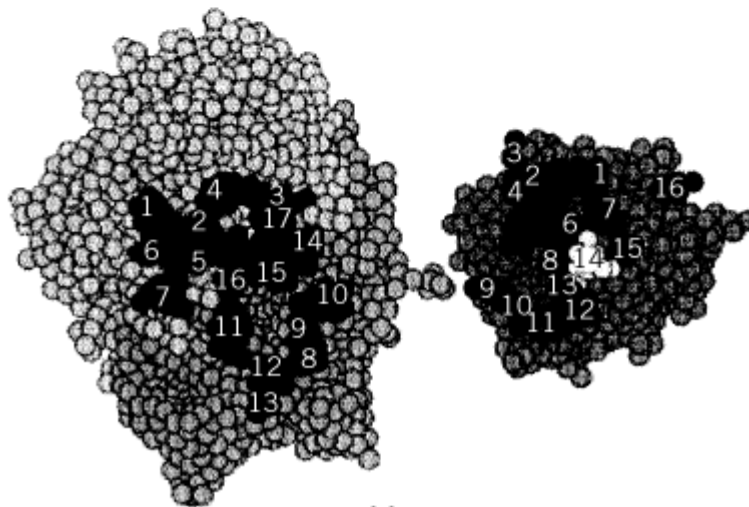
**Figure 1.** Three-dimensional structure of a lysozyme–Fab antilysozyme, showing the amino acid contributions from the heavy chain (upper part of Fab) and from the light chain (lower part). The antigen is on the right, and the Fab is on the left. (a, b) Side views. (c) Front view of the interacting residues of the antibody combining site (left) and epitope (right). (From Ref. 1, with permission from R. Poljak.)



(a)



(b)



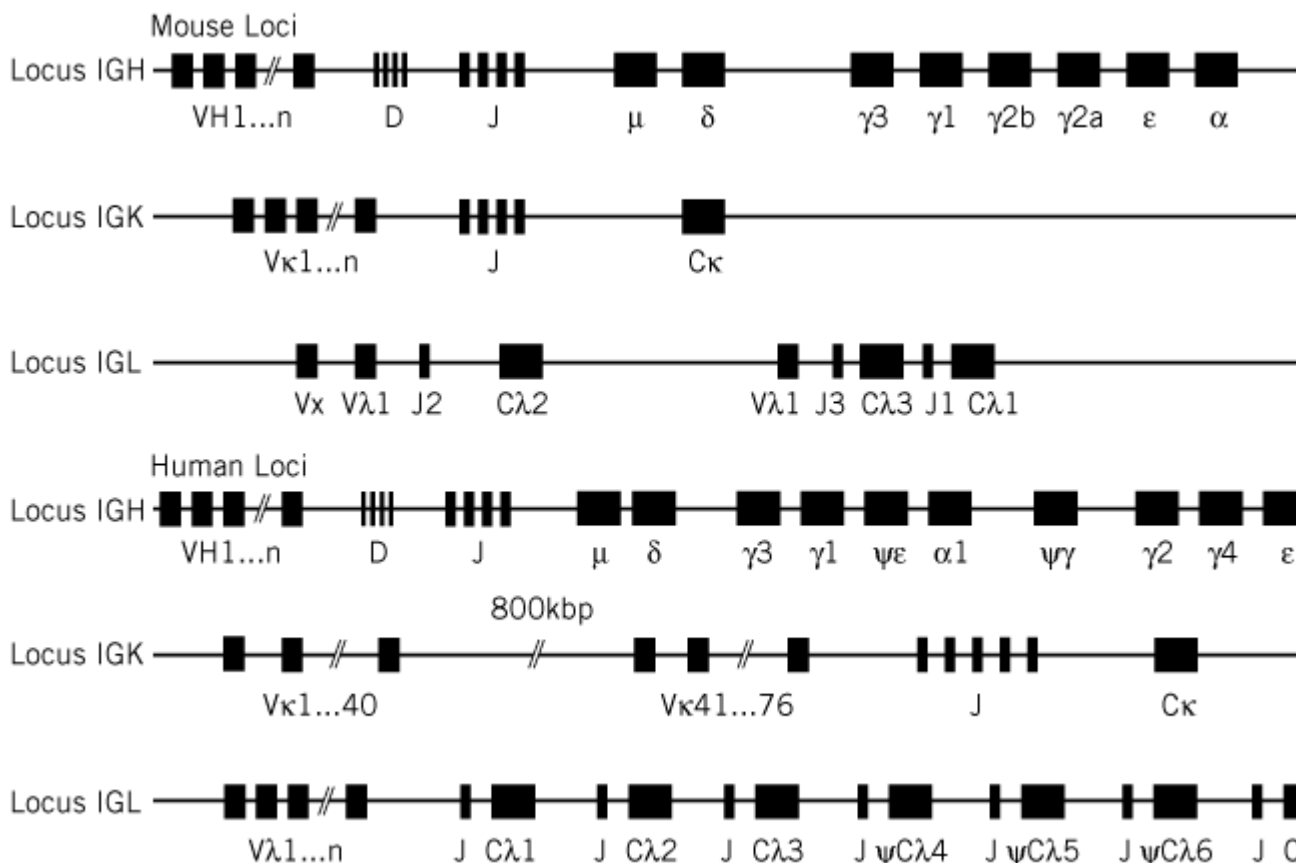
(c)

Ultimately, the main problem that the immune system has to face is how to generate such a huge diversity with a number of genes that must necessarily be limited. At the time of the first data concerning the amino acid sequences of the light chains, many hypotheses were put forward to account for this diversity. There were two extremes: One considered that the diversity was exclusively the result of [somatic hypermutations](#), whereas the other claimed that everything was encoded at the germline level, arguing that, provided that any light chain might pair with any heavy chain, 10,000 L-chain genes and 10,000 H-chain genes might generate  $10^8$  antibodies, a number that was already considered reasonable. Besides the fact that 20,000 genes would represent a large



portion of the [genome](#), this theory did not account for the conservation of the constant regions, a very serious objection that led Dreyer and Bennet (9) to propose that genes encoding the V and the C regions were separate in the germline. This was shown to be the case when the basic principles of the Ig gene organization were elucidated by the elegant experiments of Tonegawa in 1978 (10). In brief, V and C regions are encoded by separate regions within each of the three H, K, and L Ig loci, with a small number of **C genes** and a large number of **V genes** (Fig. 2). In addition, the V regions are in fact encoded by a mosaic of two gene segment clusters for the light chains ( $V_L$  and  $J_L$ ) and three for the heavy chains ( $V_H$ , D, and  $J_H$ ). Random combination of these elements takes place exclusively during B-cell differentiation and leads to a large collection of clones expressing various combinations from the basic gene mosaic. Diversity is further amplified greatly by other mechanisms, including somatic hypermutation. As a result, the number of distinct B-cell clones present at any time certainly far exceeds that necessary, especially in view of antigen–antibody recognition being partly degenerate.

**Figure 2.** Schematic organization of the three Ig gene loci in (a) mice and (b) humans. The IGK, IGL, and IGH loci, coding the k light chain, the l light chain, and the heavy chain, respectively, are located on mouse chromosomes 6, 16, and 12 and human chromosomes 2, 22, and 14, respectively.



See also entries [B Cell](#), [Immunoglobulin](#), [Clonal Selection Theory](#), [Gene Rearrangement](#), and [Repertoire](#).

#### Bibliography

1. M. F. Burnet (1959) *The Clonal Selection Theory of Acquired Immunity*. Vanderbilt University Press, Nashville, TN.
2. B. Köhler and C. Milstein (1975) Continuous culture of fused cells secreting antibody of

- predefined specificity. *Nature* **256**, 495–499.
3. J. B. Fleischman, J. B. Pain, and R. R. Porter (1962) Reduction of gammaglobulins. *Arch. Biochem. Biophys. Suppl.* **1**, 174–180.
  4. G. M. Edelman and M. D. Poulik (1961) Studies on structural units of the  $\gamma$ -globulins. *J. Exp. Med.* **113**, 861–884.
  5. R. R. Porter (1959) The hydrolysis of rabbit gammaglobulin and antibodies by crystalline papain. *Biochem. J.* **73**, 119–126.
  6. N. Hilschmann and L. Craig (1965) Amino acid sequence studies with Bence-Jones proteins. *Proc. Natl. Acad. Sci.* **53**, 1403–1409.
  7. T. T. Wu and E. A. Kabat (1970) An analysis of the sequences of the variable regions of the Bence-Jones proteins and myeloma light chains and their implications for antibody complementarity. *J. Exp. Med.* **132**, 211–250.
  8. A. G. Amit, R. A. Mariuzza, S. E. V. Phillips, and R. J. Poljak (1986) Three-dimensional structure of an antigen–antibody complex at 2.8 Å resolution. *Science* **233**, 747–753.
  9. W. J. Dreyer and J. C. Bennett (1965) The molecular basis of antibody formation: a paradox. *Proc. Natl. Acad. Sci. USA* **54**, 864–869.
  10. S. Tonegawa (1983) Somatic generation of antibody diversity. *Nature*, **302**, 575–581.

### Suggestions for Further Reading

11. G. M. Edelman, B. A. Cunningham, W. E. Gall, P. D. Gottlieb, U. Rutishauser, and M. J. Waxdal (1969) The covalent structure of an entire gamma G immunoglobulin molecule. *Proc. Natl. Acad. Sci. USA* **63**, 78–85.
12. F. Alt, T. K. Blackwell, and G. D. Yancopoulos (1987) Development of the primary antibody repertoire. *Science* **238**, 1079–1087.

## Antibody–Antigen Interactions

Binding of [antigens](#) has long been a central paradigm for molecular recognition. In addition, the biological nuances of antigen recognition have such profound ramifications for medicine that the study of antibody–antigen interactions remains a key branch of molecular immunology. In this entry we first describe the structural chemistry of antigen binding, then follow with aspects of antibody–antigen interaction that lead to unique biological phenomena.

### 1. Chemical Aspects of Antibody–Antigen Interaction

#### 1.1. General Properties

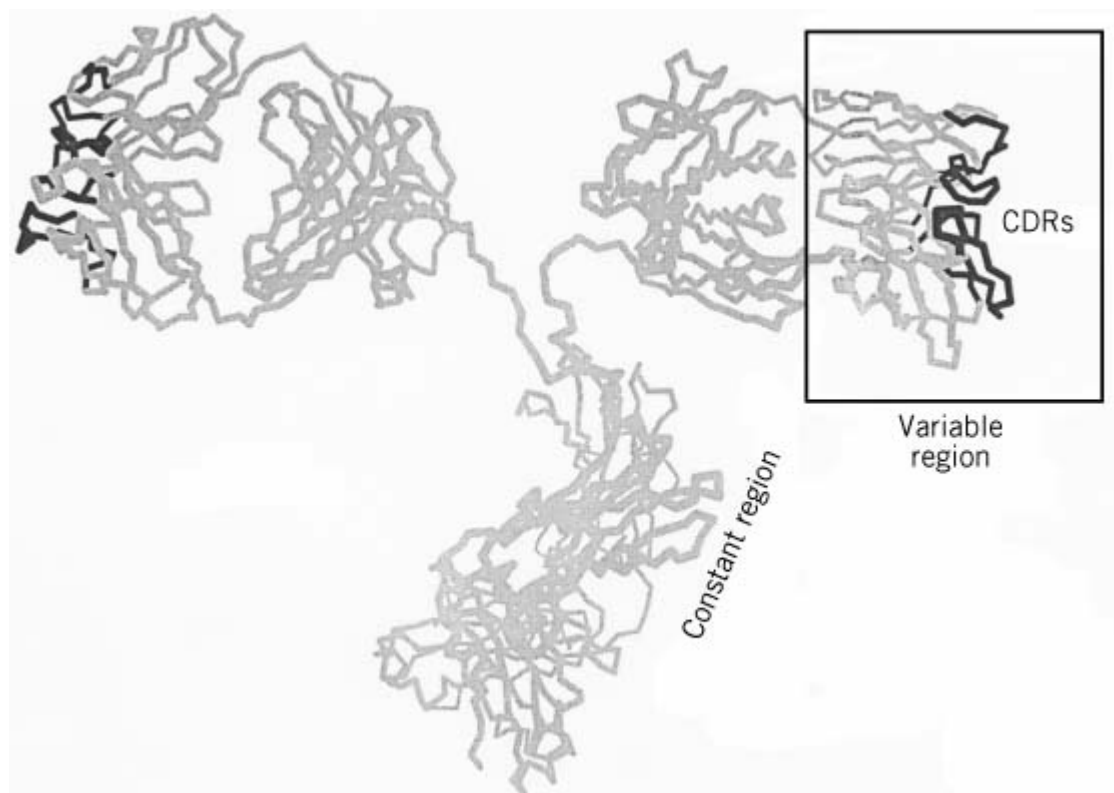
The chemical interactions between [antibody](#) and antigen do not differ substantially from other protein–ligand interactions. [Hydrogen bonds](#), [van der Waals interactions](#), and sometimes [salt bridges](#) are used to form the antibody–antigen contact. Extremely close steric complementarity between antibody and antigen surfaces seems to be a common characteristic of interfaces (1). Gaps between the opposing antibody and antigen surfaces are sometimes filled by [water](#) molecules (2). The association constant for antibody–antigen interactions ranges from  $10^5$  to  $10^{12} M^{-1}$ , with typical values for protein antigens around  $10^8$ – $10^9 M^{-1}$  (3). Rate constants for binding low molecular weight [haptens](#) can be as high as  $10^8 M^{-1}s^{-1}$ . Reactions with macromolecular antigens are slower,  $\leq 10^6 M^{-1}s^{-1}$ , except for highly charged antigens, which sometimes enhance the rate through

[electrostatic interactions](#) (4). In most cases, antigen binding does not cause easily observed changes in the binding site of an antibody. In some cases, conformational changes in the antibody are intrinsic to the binding mechanism. These changes can be caused subsequent to antigen binding, termed an “[induced fit](#)” mechanism (5); alternatively, antigen can bind selectively to one of several preexisting antibody conformations (6). Generalizations about the nature of the antibody–antigen interaction are hazardous, as is clear from the observation that the same antibody can bind one ligand through a few strong contacts and another ligand through many weak interactions (7).

### 1.2. Combining Site Structure and Function

Comparative studies showed that antibody variability was concentrated in three short stretches of sequence within each variable domain (8). These regions were postulated to confer antigenic specificity on an antibody molecule, and were termed “complementarity-determining residues” (CDRs). X-ray [crystallographic](#) studies confirmed that the CDRs, although noncontiguous in the antibody [primary structure](#), formed loops that were adjacent to each other in three-dimensional space (9). The position of CDRs within an IgG molecule is shown in Figure 1. Structural studies on antibody–antigen complexes have confirmed that the CDRs form the vast majority of intermolecular contacts. Non-CDR residues can also form contacts, however, and in some cases not all CDRs participate in the interaction with antigen (10). Of the six CDRs, heavy chain CDR3, formed at the genetic level by joining of a D segment to  $V_H$  and  $J_H$ , is the most structurally diverse, and usually the most energetically significant in binding antigen. (See [Immunoglobulin Structure](#).)

**Figure 1.** Position of CDRs (dark lines) within an IgG molecule.



### 1.3. Antigenicity

Antibodies bind to structurally precise surfaces on protein antigens (11). These surface areas can be composed of a contiguous length of polypeptide chain or from several parts of a chain that are separate in primary structure but adjacent in three-dimensional space. The only true prerequisite for a

portion of a protein to be recognized as an antigen is surface **accessibility**, although other factors, such as **hydrophilicity** and mobility, are often considered in predicting antigenicity from protein sequence (12). In the case of peptide–antibody complexes, about 7–10 peptide residues fit in the binding site of the antibody and form an ordered structure, even if the free peptide itself is not strongly ordered in solution (13). Sometimes a bound peptide structure resembles the conformation of the same sequence in an intact protein, a phenomenon that underlies the utility of peptide vaccines (14).

## 2. Biological Aspects of Antibody–Antigen Interaction

### 2.1. T-Dependent and T-Independent Immune Responses

“Normal” immune responses depend on the participation of **T cells**. In outline, antibody-producing B cells capture antigen (protein or protein–hapten conjugate) on their surface, internalize and fragment the antigen, and represent short peptide fragments on the cell surface, embedded in the binding groove of **major histocompatibility complex (MHC)** molecules. T-cells recognize the peptide–MHC combination and produce signals that cause **B cells** to enter pathways of proliferation and differentiation (15). Molecular genetic processes activated at the immunoglobulin loci include heavy-chain **class switching** and **somatic hypermutation**, leading to production of soluble **IgG**, **IgE**, and **IgA** of high affinity. T-independent antigens include virtually all nonprotein macromolecules. B-cell proliferation and differentiation also occur in a T-independent response, but chain switching and somatic mutation are difficult to detect. Structural attributes of antibodies in a T-independent response are that they are of the **IgM** isotype, contain germline (unmutated) **variable region** sequences, and show low antigen affinity.

### 2.2. Maturation of the Immune Response

Antibodies isolated soon after an initial antigen exposure, termed “primary response antibodies”, differ from those obtained later in the response or after a second administration of antigen. This transformation of the antibody **repertoire**, which leads to progressive increases in the affinity for antigen (16), is termed “maturation of the immune response.” The structural basis of this phenomenon has been determined from studies of immune responses to haptens, and can be outlined as follows. Contact with antigen induces a process within lymphocytes that introduces point mutations in the variable regions of antibody genes (17). Most mutations have a neutral or deleterious effect on antigen affinity (18). However, some mutations improve the interaction with antigen. Higher affinity confers a selective advantage on lymphocytes that express this mutation, which may be competing with nearby lymphocytes for a limited amount of antigen. Selected lymphocytes proliferate and can undergo further rounds of mutation and selection. Maturation of the immune response leads to improved affinity universally in anti-hapten responses. The same progressive affinity increase is presumed to occur with protein antigens, but unequivocal evidence for this is lacking.

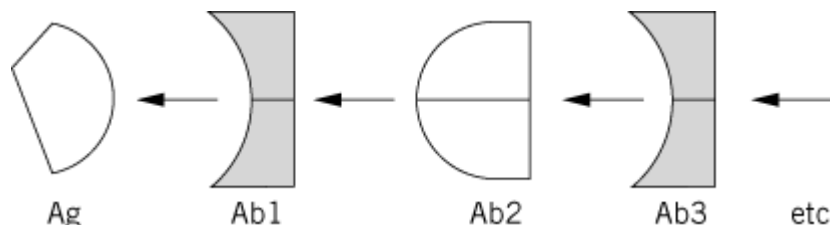
The diversity of the initial repertoire is determined by germline variable gene diversity and processes involved in **gene rearrangement**. At a structural level, antibody diversity is concentrated around the center of the antigen combining region. Somatic mutation acts to form the mature repertoire, which for protein antigens will show association constants in excess of  $10^8 M^{-1}$ . Somatic point mutations can occur anywhere in the variable region, but those selected during maturation are in general located in a band immediately peripheral to the area of diversity expressed in the early repertoire (19).

### 2.3. Anti-idiotypic Antibodies

“**Idiotypic**” is an immunological word that corresponds in structural terms to the unique CDRs of an individual antibody. If an antibody (Ab1) is used to **immunize** an animal, its CDRs can be recognized as a foreign structure, even in an animal of the same species. An antibody that reacts with the CDRs of Ab1 will be made, termed an **anti-idiotypic** antibody (Ab2). The anti-idiotypic can be used in the same way to generate an anti-anti-idiotypic (Ab3). In a theory of immune regulation (20), if Ab1 is complementary to an antigen and Ab2 is complementary to Ab1, then the binding site of

Ab2 should resemble the initiating epitope on the antigen. *In vivo*, Ab2 is then said to carry an “internal image” of the antigen. This chain of alternating complementary recognition properties is diagrammed in Figure 2).

**Figure 2.** Chain of complementary binding activities within an idiotype network.



Structural studies have investigated the extent to which the internal image model is accurate. No gross structural similarity between an antigen and Ab2 is necessary for both molecules to bind to Ab1 with high affinity. However, the structure of an antilysozyme (Ab1) complexed with an Ab2 showed that many of the Ab1 residues used for binding lysozyme were also used to bind Ab2 (21). Furthermore, many Ab2 atoms contacting Ab1 occupied positions analogous to contact atoms in lysozyme, and intermolecular hydrogen bonds and even water molecules at the interface were conserved between the Ab1:lysozyme and Ab1:Ab2 complexes. Thus Ab2 can mimic an antigen through close chemical complementarity with Ab1, even in the absence of sequence or structural homology. In the case of Ab3 molecules, structural similarity to Ab1 is unequivocal, although not surprising. CDR sequences are often nearly identical between Ab1 and Ab3, and the structure of an antigen:Ab3 complex shows chemical complementarity as stringent as would be expected at the interface between antigen and an antibody raised directly (22).

#### 2.4. Heterophile Antibodies and Molecular Mimicry

An antibody that reacts with an antigen to which the host has not been exposed is termed a “heterophile antibody” (23). The origin of heterophile antibodies is obscure; no antigen has been unambiguously identified that in a natural situation causes heterophile antibodies to appear. Structurally, heterophile antibodies tend to be low affinity IgM that react with carbohydrate antigens. Some heterophile antibodies, termed “natural antibodies,” appear to arise spontaneously and persist for the lifetime of the organism. The heterophile antibodies that determine ABO blood group compatibility are of this type. In other cases, appearance of heterophile antibodies clearly results from infection. For example, infectious mononucleosis in humans reproducibly induces antibodies specific for an antigen on sheep red blood cells (24). The supposition in such cases is of a process of “molecular mimicry”: a pathogen-specified product is made during infection that possesses sufficient structural similarity to the heterophile antigen that heterophile antibodies are induced (25). When the reactivity induced is to a self-antigen normally subject to immune **tolerance** mechanisms, [autoimmune disease](#) may result. For example, streptococcal respiratory infection is the immediate antecedent of rheumatic fever, an autoimmune attack on the heart and vascular tissue. A carbohydrate structure on the bacterial surface is thought to mimic structural glycoproteins in heart valves and vessels (26). Experimental induction of autoimmune diseases by immunization supports this type of antibody–antigen interaction as an initiating mechanism for autoimmunity (27).

#### 2.5. Polyreactive or Polyspecific Antibodies

Polyreactive antibodies do not fit the standard paradigm of one antibody recognizing a single antigen. Instead, a single polyreactive antibody can recognize a panel of many antigens that bear no obvious chemical similarity to each other. Many polyreactive antibodies are IgM, and others are IgG and other isotypes and can bind antigens with high affinity. For example, a set of [hybridomas](#) screened for simultaneous reactivity with [tubulin](#), [actin](#), myosin, and single- and double-stranded

DNA showed a median antibody affinity for all five antigens of  $10^7 M^{-1}$  (28). The molecular basis of polyreactivity is unknown. Numerous [V genes](#) have appeared in polyreactive antibodies, and there is as yet no characteristic sequence motif that predicts polyreactivity. Nevertheless, sequence studies implicated heavy-chain CDR3 as a major determinant of polyreactivity (29). Structure–function studies comparing a polyreactive and monoreactive anti-insulin with similar sequences confirmed that polyreactivity was associated with heavy-chain CDR3 (30). This CDR loop is presumed to adopt multiple conformations to accommodate different antigens, but crystallographic data on the three-dimensional structure of polyreactive antibodies are lacking.

## Bibliography

1. A. G. Amit, R. A. Mariuzza, S. E. V. Phillips, and R. J. Poljak (1986) *Science* **233**, 747–753.
2. T. N. Bhat, G. A. Bentley, G. Boulot, M. I. Green, D. Tello, W. Dall'Acqua, H. Souchon, F. P. Schwarz, R. A. Mariuzza, and R. J. Poljak (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1089–1093.
3. I. Pecht (1982) "Dynamic aspects of antibody function", in *The Antigens* (M. Sela, ed.), Academic Press, New York, Vol. **6**, pp. 1–68.
4. C. S. Raman, R. Jemmerson, B. T. Nall, and M. J. Allen (1992) *Biochemistry* **31**, 10370–10379.
5. J. M. Rini, U. Schulze-Gahmen, and I. A. Wilson (1992) *Science* **255**, 959–965.
6. D. Lancet and I. Pecht (1976) *Proc. Natl. Acad. Sci. USA* **73**, 3548–3553.
7. W. Dall'Acqua, E. R. Goldman, E. Eisenstein, and R. A. Mariuzza (1996) *Biochemistry* **35**, 9667–9676.
8. T. T. Wu and E. A. Kabat (1970) *J. Exp. Med.* **132**, 211–250.
9. R. J. Poljak, L. M. Amzel, H. P. Avey, B. L. Chen, R. P. Phizackerly, and F. Saul (1973) *Proc. Natl. Acad. Sci. USA* **70**, 3305–3310.
10. D. R. Davies, E. A. Padlan, and S. Sheriff (1990) *Annu. Rev. Biochem.* **59**, 439–473.
11. D. C. Benjamin, J. A. Berzofsky, I. J. East, F. R. N. Gurd, C. Hannum, S. J. Leach, E. Margoliash, J. G. Michael, A. Miller, E. M. Prager, M. Reichlin, E. E. Sercarz, S. J. Smith-Gill, P. E. Todd, and A. C. Wilson (1984) *Annu. Rev. Immunol.* **2**, 67–101.
12. J. A. Berzofsky (1985) *Science* **229**, 932–940.
13. R. Stanfield and I. A. Wilson (1995) *Curr. Opin. Struct. Biol.* **5**, 103–113.
14. R. A. Lerner (1982) *Nature* **299**, 592–596.
15. A. Lanzavecchia (1985) *Nature* **314**, 537–539.
16. H. N. Eisen and G. W. Siskind (1964) *Biochemistry* **3**, 996–1008.
17. M. G. Weigert, I. M. Cesari, S. J. Yonkovich, and M. Cohn (1970) *Nature* **228**, 1045–1047.
18. C. Chen, V. A. Roberts, and M. B. Rittenberg (1992) *J. Exp. Med.* **176**, 855–866.
19. I. M. Tomlinson, G. Walter, P. T. Jones, P. H. Dear, E. L. L. Sonnhammer, and G. Winter (1996) *J. Mol. Biol.* **256**, 813–817.
20. N. Jerne (1974) *Ann. Inst. Pasteur (Immunol)* **125C**, 373–389.
21. B. A. Fields, F. A. Goldbaum, X. Ysern, R. J. Poljak, and R. A. Mariuzza (1995) *Nature* **374**, 739–742.
22. K. C. Garcia, S. V. Desiderio, P. M. Ronco, P. J. Verroust, and L. M. Amzel (1992) *Science* **257**, 528–531.
23. J. Forssmann (1911) *Biochem. Z.* **37**, 78–115.
24. J. R. Paul and W. W. Bunnell (1932) *Am. J. Med. Sci.* **183**, 90–104.
25. R. T. Damian (1964) *Am. Naturalist* **98**, 129–149.
26. J. B. Zabriskie, K. C. Hsu, and B. C. Seegal (1970) *Clin. Exp. Immunol.* **7**, 147–159.
27. R. S. Fujinami and M. B. A. Oldstone (1985) *Science* **230**, 1043–1045.
28. L. Diaw, C. Magnac, O. Pritsch, M. Buckle, P. M. Alzari, and G. Dighero (1997) *J. Immunol.* **158**, 968–976.

29. C. Chen, M. P. Stenzel-Poore, and M. B. Rittenberg (1991) *J. Immunol.* **147**, 2359–2367.  
30. Y. Ichiyoshi and P. Casali (1994) *J. Exp. Med.* **180**, 885–895.

### Suggestions for Further Reading

31. E. A. Padlan (1994) Anatomy of the antibody molecule, *Mol. Immunol.* **31**, 169–217. (A superb, comprehensive review of antibody structure and function.)  
32. *FASEB J.* **9**, 1–147 (1994). An entire issue devoted to structural immunology, including reviews and research papers.  
33. *Immunol. Rev.* **163**, “Molecular anatomy of the immune response” (1998). (An entire volume of reviews of structural immunology.)

### Antibody-Conjugated Toxins

Soon after the discovery of the high specificity of [antibody](#) binding, Paul Ehrlich proposed to exploit such specificity to redirect and restrict the activity of [toxins](#) to pathological tissues such as tumors ([1-3](#)). Such “magic bullets” could be developed only long after his proposal because of the need of antibodies of well-defined specificity, such as [monoclonal antibodies](#), and to improve the understanding of the biochemistry and mechanism of the action of toxins. Such conjugates have been prepared by linking one of the many available cell-killing plant or bacterial toxins via a **disulfide** or thioether bond to an antibody specific for a cell surface antigen of cancer and noncancer cells. Chimeric toxins have also been prepared by fusing a toxin gene to a gene encoding a [hormone](#), [growth factor](#), or [cytokine](#) with the intention of depleting receptor-rich cell populations ([1-3](#)). In general, it is important that the linkage between the two immunotoxin parts be easily cleaved on or inside cells. Moreover, to restrict immunotoxin binding only to the cells to be intoxicated, the portion of the toxin determining its intrinsic specificity (the B subunit) must be either deleted or inactivated mutationally, and the  $F_c$  portion of the antibody must be removed. In this respect, **ribosome**-inactivating plant toxins lacking the B subunit have been frequently used. To avoid unspecific binding to  $F_c$ -receptor bearing cells, the toxin can be linked to the  $F(ab)_2$  portion of a monoclonal antibody. An additional problem may be presented by the presence of circulating antitoxin antibodies elicited in previous vaccination protocols (eg, **antidiphtheria**) or during the treatment with the immunotoxin. In the latter case, one can use another immunotoxin made with the same antibody and an immunologically unrelated toxin.

Several antibody-conjugated toxins are very effective in killing target cells in culture or in isolated dispersed tissues, ie, bone marrow, but are less effective *in vivo*, because of the lack of accessibility of the cancer cells in solid tumors. However, protocols are being developed to overcome these problems and exploit the potentials of immunotoxins to reach small cancer populations that are undetected by physical and surgical methods ([3](#)).

### Bibliography

1. S. Olsnes, K. Sandvig, O. W. Petersen, and B. van Deurs (1989) *Immunol. Today* **10**, 291–295.
2. D. A. Vallera (1994) *Blood* **83**, 309–317.
3. G. R. Trush, L. R. Lark, and E. S. Vitetta (1996) *Ann. Rev. Immunol.* **14**, 49–71.

## Anticodon

The anticodon of [transfer RNA](#) (tRNA) is in the central part of the linear RNA sequence, at positions 34, 35, and 36. This triplet of nucleotides forms transitory base pairs with three nucleotide codons of the **messenger RNA** (mRNA) during **protein biosynthesis**. These codon/anticodon interactions enable the mRNA to direct the order of incorporation of [amino acids](#) into the [polypeptide chain](#). These interactions occur on the small subunit of the ribosome. The presence of the anticodon in the seven-membered central loop of tRNA led to the designation of this subdomain as the anticodon loop. The anticodon stem plus the loop constitute the anticodon arm (stem and loop), which contributes to one branch of the L-shaped 3D structure of tRNA. In this structure, the anticodon is about 80 Å from the amino acid acceptor end of the tRNA (see Fig. 1 of [Transfer RNA](#)). Besides its fundamental contribution to translation of the genetic information through complementary pairing to the codon, the anticodon of tRNA may contain important identity elements, recognized by cognate [aminoacyl-tRNA synthetases](#), which lead to specific aminoacylation of the tRNA. Besides the anticodon, other members of the loop and of the nearby stem, including modified nucleotides, also contribute both to the efficiency and specificity of aminoacylation.

Codon–anticodon interaction occurs basically through classic **Watson–Crick base pairs**. However, additional types of interactions allow **wobble** to occur between the third base of the codon and the first base of the anticodon (1). The original wobble rules suggested that the first nucleoside of the anticodon can pair with more than one nucleoside at the third position of the codon. Thus, anticodons with a U at the first position could interact with codons having either A or G at the third position. Those presenting a G at position 34 could interact with codons terminating with U or C. More interestingly, tRNA presenting an inosine (deaminated adenosine) at position 34 could recognize codons terminating with either C, U, or A. For example, yeast tRNA<sup>Ala</sup>, anticodon 5'-IGC-3', interacts with three codons: 5'-GCC-3', 5'-GCU-3', and 5'-GCA-3'.

Numerous data accumulated over the years led to revised wobble rules (2, 3), which reflect new possible interactions between classic bases and take into account the vast occurrence of modified nucleosides in the anticodons of tRNA, especially those at the wobble position (residue 34, first position of the anticodon). One minor tRNA, *Escherichia coli*, tRNA<sup>Ile2</sup>, has a methionine anticodon CAU. However, in the mature tRNA, C34 is modified by covalent attachment of lysine on position 2 of the pyrimidine ring, converting it to lysidine or a “k2C” base. Lysidine at the wobble position leads to recognition of the isoleucine-specific codon 5'-AUA-3' instead of the methionine codon AUG. Interestingly, this modification is also required for specific recognition of the tRNA by isoleucyl-tRNA synthetase (see discussion below). Other examples of modification at position 34 that contribute to the specificity of codon–anticodon recognition have been reported. This is also true for nucleotide 37, at the 3' side of the anticodon triplet. Many different modified purine nucleosides have been identified at this position, and these modifications contribute to the fidelity of protein synthesis. For example, hypermodification of A to i6A (*N*-6-isopentenyladenosine), t6A (*N*-6-threonylcarbamoyladenosine), and their derivatives, stabilizes the relatively weak A-U and U-A base pairing between the third position of the anticodon and the first position of the codon. The presence of the hypermodified nucleoside ms2i6A (*N*-6-(*D*-2-isopentenyl)-2-methylthioadenosine), in contrast, prevents codon misreading by *E. coli* tRNA<sup>Phe</sup>. Modifications introduce conformational flexibility or rigidity that restrict or enlarge the number of potential base pairs. Thus the molecular mechanism by which modified bases alter codon recognition are largely structural in nature.

UAG (**amber**), UAA (**ocher**), and UGA (opal) codons do not code for an amino acid, because there are no tRNA with corresponding anticodons. Known as **nonsense codons**, they normally signal termination of [translation](#), but such codons can be created by [mutation](#), when they cause premature



termination of protein synthesis. Interestingly, nonsense **suppressor** tRNA, with mutations in the anticodon, can recognize and “suppress” nonsense mutations by inserting a specific amino acid during translation (4). The first suppressor tRNA identified was a tyrosine-specific tRNA in which a single base substitution converted the anticodon from GUA (Tyr) to CUA (amber). Conversion of a classic tRNA to an efficient suppressor RNA may also involve mutations outside the anticodon. Suppressor tRNA have been used as a tool in the search of tRNA identity elements *in vivo* and as a way of inserting desired amino acids at specific sites in proteins.

The anticodon triplet specifies which amino acid the tRNA will insert in response to a codon and is, of course, directly correlated with the amino acid bound to the -CCA end. A simple analysis of the [genetic code](#), however, reveals that the anticodon is not a common signal for synthetase recognition. For example, there are six different serine codons and, consequently, the potential for six different anticodon sequences in serine-specific tRNA. This great variability in isoacceptor tRNA precludes a common signal in the anticodon for recognition by seryl-tRNA synthetase, but it is clear that anticodon nucleotides, as well as other nucleotides within the anticodon loop, are critical for which amino acids are charged by aminoacylation (5-10).

Single-point mutations at any anticodon nucleotide of yeast tRNA<sup>Phe</sup> or yeast tRNA<sup>Asp</sup> lead to severe losses in aminoacylation efficiency (11, 12). A simple anticodon switch changes the aminoacylation specificity of methionine and valine tRNA (13). Depending on the tRNA, one, two, or three anticodon nucleotides may be involved in tRNA specificity. Position 35 is required for charging of arginine, both 35 and 36 are important for valine or threonine, and 34, 35, and 36 are involved in specificity for Asn, Asp, Cys, Gln, Ile, Lys, Met, Phe, Trp, and Tyr. In *E. coli* tRNA<sup>Gln</sup>, residue 35 is of greater importance than residues 34 and 36. In yeast tRNA<sup>Asp</sup>, nucleotides 34 and 35 contribute more than do nucleotide 36. Transplantation of anticodons into noncognate tRNA results in increases of four to six orders of magnitude in aminoacylation of the chimeric tRNA with the amino acid corresponding to the transplanted anticodon. Since it does not preclude aminoacylation of the chimeric tRNA with the original amino acid, however, it is clear that additional identity elements elsewhere in the tRNA structure dictate amino acid specificity (see [Transfer RNA](#)). The extent and the relative contribution of the different nucleotides of the anticodon to this specificity varies. The losses in specificity for yeast tRNA<sup>Asp</sup> are about 500-fold, while those for *E. coli* tRNA<sup>Val</sup> are 100,000-fold.

Other nucleotides within the anticodon domain often contribute to amino acid specificity. Residue 38 is involved in yeast tRNA<sup>Asp</sup>, and both residues A37 and U38 are involved in *E. coli* tRNA<sup>Gln</sup>. An original approach to anticodon domain function uses the *in vitro* synthesis of minihelices (shortened forms of tRNA) mimicking isolated anticodon stem and loop domains. These minihelices possess identity elements and were tested for stimulation of aminoacylation rates of the acceptor stem as substrate and/or for inhibition of full-length tRNA charging (14). Yeast valyl-tRNA (15) and *E. coli* isoleucine-tRNA (16) anticodon minihelices stimulated, up to threefold, the aminoacylation of an acceptor arm minihelix by cognate aminoacyl-tRNA synthetase. The isolated anticodon domain of tRNA<sup>fMet</sup> binds to *E. coli* methionyl-tRNA synthetase (17).

The anticodon domain is further implicated in discrimination between initiator and elongator forms of tRNA (18) and in alternate functions of tRNA, such as initiation of **reverse transcription in retroviruses** (19) (see [Transfer RNA](#)).

## Bibliography

1. F. H. C. Crick (1966) *J. Mol. Biol.* **19**, 548–555.
2. S. Yokoyama and S. Nishimura (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology, Washington, DC, pp. 207–223.
3. K. Watanabe and S. Osawa (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and

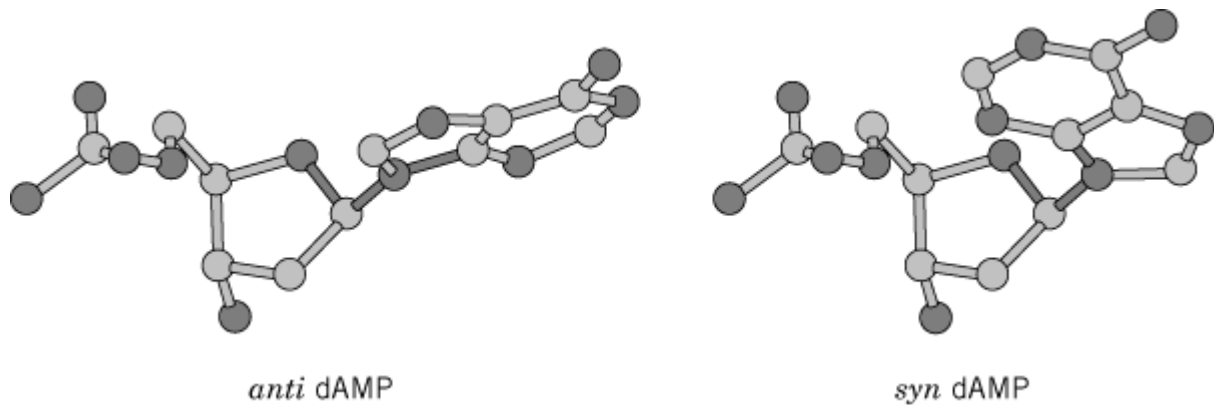
- U. L. RajBhandary, eds.), American Society for Microbiology, Washington, DC, pp. 225–250.
4. H. Ozeki, H. Inokuchi, F. K. M. Yamao, H. Sakano, T. Ikemura, and Y. Shimura (1980) in *Transfer RNA: Biological Aspects* (D. Söll, J. Abelson, and P. Schimmel, eds.), Cold Spring Harbor Laboratory, New York, pp. 341–362.
  5. L. Kisselev (1985) *Prog. Nucleic Acid Res. Mol. Biol.* **32**, 237–266.
  6. L. Pallanck, M. Pak, and L.-H. Schulma (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology, Washington DC, pp. 371–394.
  7. R. Giegé, J. D. Puglisi, and C. Florentz (1993) *Prog. Nucleic Acid Res. Mol. Biol.* **45**, 129–206.
  8. W. H. McClain (1993) *J. Mol. Biol.* **234**, 257–280.
  9. M. E. Saks, J. R. Sampson, and J. N. Abelson (1994) *Science* **263**, 191–197.
  10. W. H. McClain (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, DC, pp. 335–347.
  11. J. R. Sampson, A. B. DiRenzo, L. S. Behlen, and O. C. Uhlenbeck (1989) *Science* **243**, 1363–1366.
  12. J. Putz, J. D. Puglisi, C. Florentz, and R. Giegé (1991) *Science* **252**, 1696–1699.
  13. L. H. Schulman and H. Pelka (1988) *Science* **242**, 765–768.
  14. S. A. Martinis and P. Schimmel (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, DC, pp. 349–370.
  15. M. Frugier, C. Florentz, and R. Giegé (1992) *Proc. Natl. Acad. Sci. USA* **89**(9), 3990–3994.
  16. O. Nureki, T. Niimi, T. Muramatsu, H. Kanno, T. Kohno, C. Florentz, R. Giegé, and S. Yokoyama (1994) *J. Mol. Biol.* **236**, 710–724.
  17. T. Meinnel, Y. Mechulam, S. Blanquet, and G. Fayat (1991) *J. Mol. Biol.* **220**, 205–208.
  18. N. Mandal, D. Mangroo, J. J. Dalluge, J. A. McCloskey, and U. L. RajBhandary (1996) *RNA* **2**, 473–482.
  19. C. Isel, J.-M. Lanchy, S. F. J. Le Grice, C. Ehresmann, B. Ehresmann, and R. Marquet (1996) *EMBO J.* **15**, 917–924.

## Anti Conformation

*Anti* is used as a prefix to describe the relative orientation of two substituents on adjacent atoms of a molecule when their [torsion angle](#) is about 180° (1) (see [Conformation](#)). In addition to its formal usage, *anti* has become broadly used to describe the orientation of the **nucleic acid** bases relative to the ribose ring when describing the torsion angle about the glycosidic bond (2). The glycosidic bond torsion is defined by the relative positions of the four atoms surrounding this bond, O4'—C1—N9 (1—C4(2)) for the purine (pyrimidine) nucleosides (3). When C4 is not over the ribose ring, the conformation is *anti*, corresponding to torsion angles of –80 to –180 and 100 to 180°. Conformations where the purine ring is over the sugar are *syn* and correspond to torsion angles of –80 to 100°. These conformations for adenine are shown in Figure 1.

**Figure 1.** The *anti* and *syn* conformations of dAMP are shown. The bonds defining the glycosidic torsion angle are darkened for emphasis in both conformers. In the *anti* conformer, the torsion is 110°, and in the *syn* conformer the

torsion angle is +60.



### Bibliography

1. W. Klyne and V. Prelog (1960) *Experientia* **16**, 521.
2. A. E. V. Haschemeyer and A. Rich (1967) *J. Mol. Biol.* **27**, 369–384.
3. IUPAC-IUB Commission on Biochemical Nomenclature (1983) *Eur. J. Biochem.* **17**, 193–201.

### Suggestions for Further Reading

4. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry, Part I: The Conformation of Biological Macromolecules*, W. H. Freeman, San Francisco, CA.
5. E. L. Eliel et al. (1967) *Conformational Analysis*, Wiley-Interscience, New York.
6. W. Saenger (1984) *Principles of Nucleic Acid Structure*. "Springer Advanced Texts in Chemistry" (C. R. Cantor, ed.), Springer-Verlag, New York.

### Antifreeze Proteins

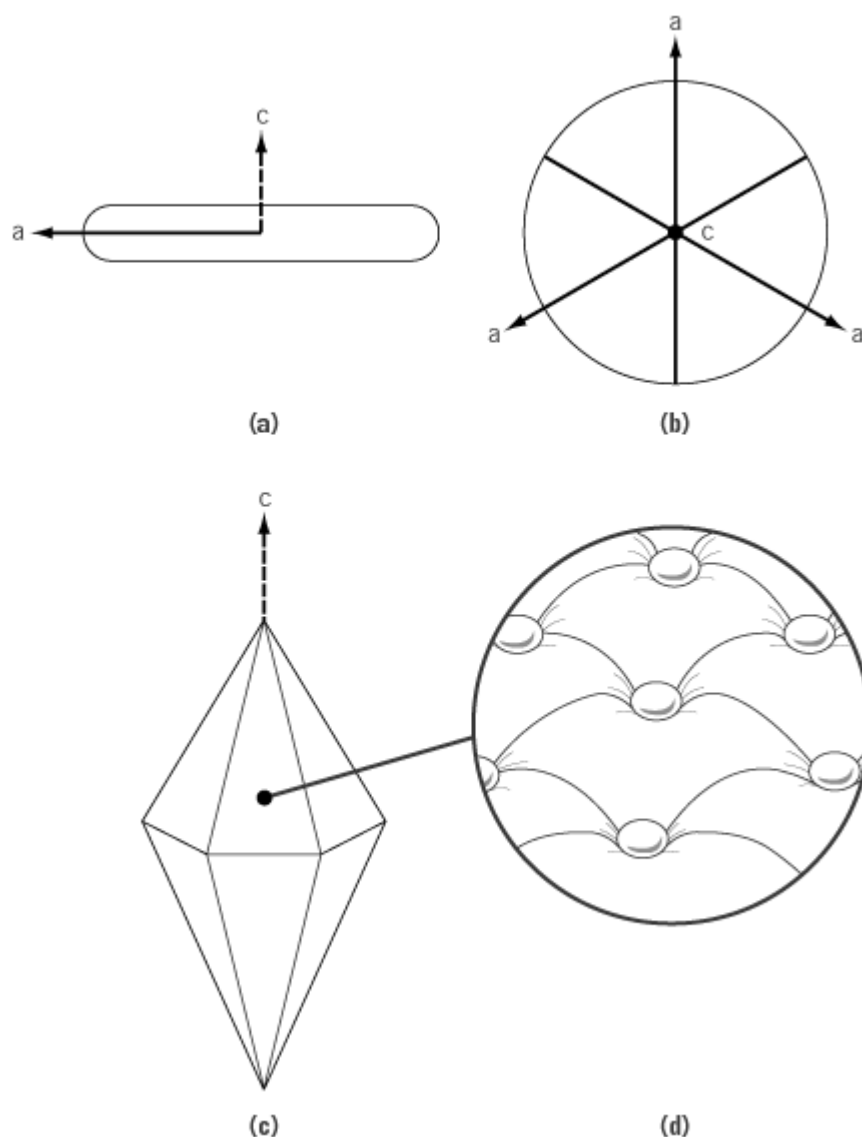
Antifreeze proteins (AFPs) adsorb specifically to the surface of ice. By doing so, they inhibit its growth and cause a nonequilibrium depression of the freezing point below the colligative melting/freezing point (1). This depression, which is referred to as thermal hysteresis ( $^{\circ}\text{C}$ ), helps organisms to avoid or resist freezing when they encounter ice at temperatures below their colligative freezing points. For example, in polar and cold temperate marine environments bony fishes (teleosts) can encounter icy seawater at its freezing point of  $\approx -1.9^{\circ}\text{C}$ , which is over  $1^{\circ}\text{C}$  colder than their colligative freezing point of  $\approx -0.7^{\circ}\text{C}$ . Those fishes that have adapted to this environment produce high concentrations of AFP in their blood (10–25 mg/mL) that can depress their freezing point by the crucial  $1^{\circ}\text{C}$  needed for survival. For practical reasons, AFPs (sometimes known as thermal hysteresis proteins) have been best characterized in fish, but have also been found in insects, plants, and bacteria that encounter subzero conditions. Not all of these organisms resist freezing. In those that do freeze, AFPs may enhance their freeze tolerance by keeping ice crystals from growing (recrystallizing) in the frozen state (2). Recrystallization of ice occurs as the melting temperature is approached and ice becomes more dynamic. It is effectively inhibited by very low AFP concentrations ( $\mu\text{g/mL}$ ). Although there are reports of AFPs neutralizing ice nucleators (3) and of having protective effects on cells above  $0^{\circ}\text{C}$  (4), the two well-established functions of

macromolecular antifreezes, thermal hysteresis and inhibition of ice recrystallization, are related activities that depend on AFP interaction with ice by an absorption-inhibition mechanism.

## 1. Inhibition of Ice Growth

In ice, each [water](#) molecule makes four tetrahedrally oriented [hydrogen bonds](#) to neighboring molecules, and the resulting ice lattice has an underlying hexagonal symmetry. When an aqueous solution is cooled in the presence of an ice crystal, water adds preferably to the atomically rough prism surfaces of the crystal much faster than to the two smoother basal planes. At slight undercooling below the freezing point, the prism surfaces are curved, and the ice grows radially as a disk (Fig. [1](#)). Most solutes are swept aside by the expanding ice front. However, AFPs with a specific affinity for an ice surface plane become adsorbed to that surface and cause growth inhibition that further defines the surface plane and typically shapes the ice crystal into a hexagonal bipyramid. The basis for growth inhibition at the surface of ice is that water is forced to add to the lattice between the bound AFP ([1](#), [5](#)). These submicroscopic ice fronts are constrained to grow with a surface curvature that makes it thermodynamically unfavorable for water to join the lattice, which in turn leads to a local depression of the freezing point below that of the bulk solvent. Rather than completely covering the surface, AFPs have been likened to buttons on a mattress (Fig. [1d](#)).

**Figure 1.** i) In the absence of AFP, a single ice crystal from the melt will grow fastest along the crystallographic *a* axes to form a disk with curved prism surfaces. ii) View of disk perpendicular to i) with *c* axis perpendicular to the page. iii) AFP adsorption to any aspect of the prism surface will force the crystal to become a hexagonal bipyramid. iv) The inset shows a section of the surface illustrating the submicroscopic ice front curvature between bound AFPs.



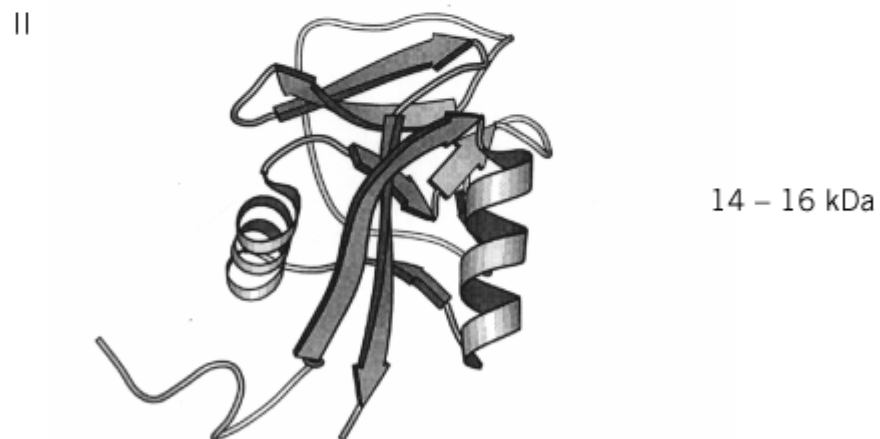
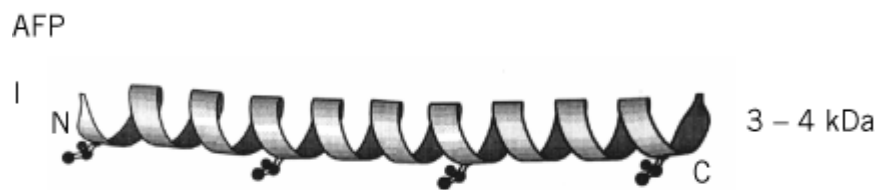
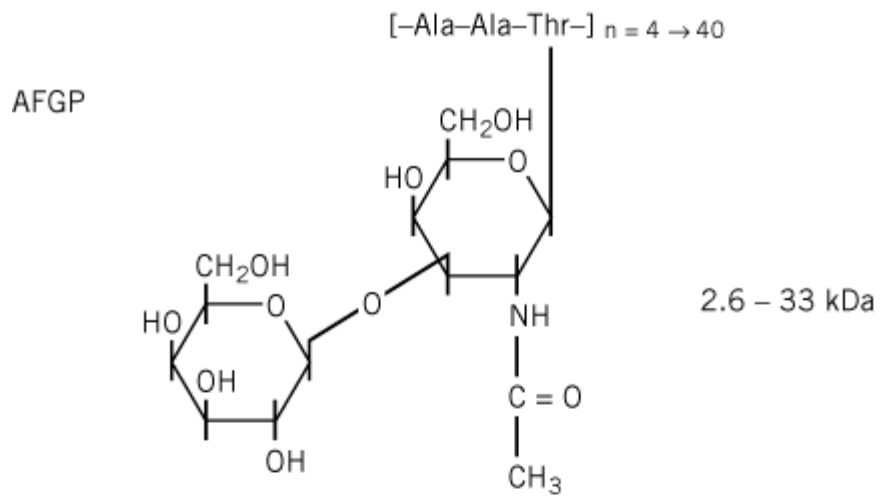
The relationship between thermal hysteresis and AFP concentration is nonlinear. Thermal hysteresis increases hyperbolically with increasing AFP concentration and approaches a plateau value of between 1 and 1.5°C for fish AFPs and 5 to 6°C for insect AFPs. If the temperature falls below the nonequilibrium freezing point at moderate to high AFP concentrations, the ice crystal will grow uncontrollably. At low AFP concentrations and moderate undercooling, the ice front will overgrow bound AFPs. This phenomenon has been used to deduce the ice surface to which AFPs bind by growing ice in a dilute AFP solution and then allowing it to sublime (ice etching) (5). The position of the AFP residue etch on the ice indicates its binding plane and, in some cases, its orientation of binding on that plane.

## 2. Adsorption of AFPs to Ice

One of the most intriguing features of AFPs is their structural diversity. In fish, there are at least five unrelated types, which are also different from insect AFPs (Fig. 2). It is likely that additional types will be discovered. They are so different that it has not been possible to identify a shared sequence or structural motif that could explain their common ice-binding activity. An added complication is that they bind to different ice surface planes (6). Early models for binding of AFPs to ice relied on a hydrogen-bonding match to the ice lattice (5, 7). However, it has been argued that hydrogen bonding of AFPs to liquid water is energetically favored over binding to ice (8). Also, correct identification of

the ice-binding sites now shows them to be more hydrophobic than originally thought (9). One structural feature that AFPs have in common is an extensive region of intimate surface complementarity between the AFP and ice. Thus the energetic contributions to binding could include [van der Waals interactions](#) over the AFP–ice interface, some hydrogen bonds, and an entropic component from not having to solvate the AFP and ice surfaces in contact (10)

**Figure 2.** AFP structures. The 3D structures of Type I AFP from winter flounder (1WFA); Type II AFP from sea raven (2AFP); Type III AFP from ocean pout (1MSI); beetle AFP (1EZG); and moth AFP were solved by X-ray crystallography or nuclear magnetic resonance. The approximate masses of the proteins are indicated in kDa.



### 3. Evolution of AFPs

Fish AFPs have been classified into five different types. Antifreeze **glycoproteins** (AFGPs) are

polymers of the tripeptide repeat Ala-Ala-Thr with a disaccharide (b-D-galactopyranosyl-(1→3)-2-acetamido-2-deoxy-a-D-galactopyranose) attached to the Thr. Type I AFP is an alanine-rich single  $\alpha$ -helix. Type II AFP is a globular protein with sequence homology and structural similarity to the carbohydrate recognition domain of calcium-dependent [lectins](#). Type III AFP is a smaller protein with a unique  $\beta$ -sandwich fold that is homologous to a portion of sialic acid synthase. The most recently discovered (Type IV) AFP has homology to apolipoprotein E and is thought to be a helix-bundle protein.

It is clear that the different AFP types have distinct origins, which suggests that they have been independently co-opted as antifreezes relatively recently in the 175-million-year radiation of the bony fishes ([11](#)). Indeed, there are now several examples of very closely related fishes producing different AFP types and at least one clear example of distantly related fishes producing the same type through **convergent** evolution ([12](#)). This pattern of recent evolution is reinforced by the observation that AFP [gene families](#) show evidence of extensive **amplification**, often with the AFP genes organized as tandem repeats and with substantial variation in gene dosage ([13](#)). The diversity of AFP types is also reflected in their expression and processing. AFP genes are either constitutively expressed or are regulated by a variety of cues, including temperature, photoperiod, and hormonal cycles ([14](#)). Fish AFPs are typically produced in the liver as pre-proteins for secretion into the serum with or without a pro-sequence. The small antifreeze glycoproteins are produced from a large [polyprotein](#) by both **proteolysis** and glycosylation ([12](#)). In winter flounder, the external tissues, skin, scales, fins, and gills produce **isoforms** representing a subclass of Type I AFP that appears to lack a signal sequence ([15](#)).

Not unexpectedly, plant and insect AFPs are also different. AFPs from winter rye appear to be homologs of three types of extracellular plant pathogenesis-related sequences, namely endochitinases, endoglucanases, and thaumatin-like proteins ([16](#)), and carrot AFP is related to polygalacturonase inhibitor ([17](#)). These similarities reinforce the links between AFPs and sugar/polysaccharide-binding proteins. Plant AFPs appear to be more active in ice recrystallization inhibition than thermal hysteresis, whereas insect AFPs show considerably more thermal hysteresis activity than fish AFPs. AFPs from moths and beetles have distinct origins, and yet they both form  $\beta$ -helices with similar threonine-rich ice-binding sites ([18](#), [19](#)). In some terrestrial insects that can resist freezing at temperatures down to  $-30^{\circ}\text{C}$ , AFPs may be part of a suite of cryoprotective mechanisms that include elevated polyol levels and the elimination of ice nucleators that would interfere with supercooling.

#### 4. Conclusion

Although adsorption-inhibition is the generally accepted mechanism of AFP action, the energetics of tight binding to ice are not yet fully understood. Other areas of current research include the basis for the hyperactivity of insect AFPs, activities of AFPs other than thermal hysteresis, and applications for AFPs in conferring freeze tolerance/resistance to unprotected organisms.

#### Bibliography

1. P. W. Wilson (1993) *Cryo-Letters* **14**, 31–36.
2. C. A. Knight, A. L. DeVries, and L. D. Oolman (1984) *Nature* **308**, 295–296.
3. A. Parody-Morreale, K. P. Murphy, E. DiCera, R. Fall, A. L. DeVries and S. J. Gill (1988) *Nature* **333**, 782–783.
4. B. Rubinsky, A. Arav, and G. L. Fletcher (1991) *Biochem. Biophys. Res. Commun.* **180**, 566–571.
5. C. A. Knight, C. C. Cheng and A. L. DeVries (1991) *Biophys. J.* **59**, 409–418.
6. C. C. Cheng and A. L. DeVries (1991) *Life Under Extreme Conditions* (G. di Prisco, ed.), Springer-Verlag, Berlin, pp. 1–14.
7. D. Wen and R. A. Laursen (1992) *Biophys. J.* **63**, 1659–1662.



8. A. D. J. Haymet, L. G. Ward, and M. M. Harding (1999) *J. Am. Chem. Soc.* **121**, 941–948.
9. J. Baardsnes, M. E. Houston, Jr., H. Chao, R. S. Hodges, C. M. Kay, and P. L. Davies (1999) *FEBS Letters* **463**, 87–91.
10. F. D. Sönnichsen, C. I. DeLuca, P. L. Davies and B. D. Sykes (1996) *Structure* **4**, 1325–1337.
11. G. K. Scott, G. L. Fletcher, and P. L. Davies (1986) *Can. J. Fisheries Aquatic Sci.* **43**, 1028–1034.
12. L. Chen, A. L. DeVries, and C. H. C. Cheng (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3817–3822.
13. C. L. Hew, N. C. Wang, S. Joshi, G. L. Fletcher, G. K. Scott, P. H. Hayes, B. Buettner, and P. L. Davies (1988) *J. Biol. Chem.* **263**, 12049–12055.
14. P. L. Davies, C. L. Hew, and G. L. Fletcher (1988) *Can. J. Zool.* **66**, 2611–2617.
15. Z. Gong, K. V. Ewart, Z. Hu, G. L. Fletcher, and C. L. Hew (1996) *J. Biol. Chem.* **271**, 4106–4112.
16. W. C. Hon, M. Griffiths, A. Mlynraz, Y. C. Kwok, and D. S. C. Yang (1995) *Plant Physiol.* **109**, 879–989.
17. D. Worrall, L. Elias, D. Ashford, M. Smallwood, C. Sidebottom, P. Lillford, J. Telford, C. Holt, and D. Bowles (1998) *Science* **282**, 115–117.
18. Y.-C. Liou, A. Tocilj, P. L. Davies, and Z. Jia (2000) *Nature* **406**, 322–324.
19. S. P. Graether, M. J. Kuiper, S. M. Gagne, V. K. Walker, Z. Jia, B. D. Sykes, and P. L. Davies (2000) *Nature* **406**, 325–328.

### Suggestions for Further Reading

20. M. Griffith and K. V. Ewart (1995) Antifreeze proteins and their potential use in frozen foods. *Biotechnol. Adv.* **13**, 375–402.
21. Y. Yeh and R. E. Feeney (1996) Antifreeze proteins: structures and mechanisms of function. *Chem. Rev.* **96**, 601–617.
22. P. L. Davies and B. D. Sykes (1997) Antifreeze proteins. *Curr. Opin. Struct. Biol.* **7**, 828–834.
23. K. V. Ewart, Q. Lin, C. L. Hew (1999) Structure, function and evolution of antifreeze proteins. *Cell. Mol. Life Sci.* **55**, 271–283.
24. J. G. Duman (2001) Antifreeze and ice nucleator proteins in terrestrial arthropods. *Annu. Rev. Physiol.* **63**, 327–357.
25. G. L. Fletcher, C. L. Hew, and P. L. Davies (2001) Antifreeze proteins of teleost fishes. *Annu. Rev. Physiol.* **63**, 359–390.

## Antigen

The historical definition of an antigen, given for many years in immunology textbooks, is a rather circular one, because it describes the antigen as a foreign substance that induces upon penetration in an animal (or a human being) the production of an [antibody](#) with which it will combine specifically, *in vivo* or *in vitro*. This classical antibody response is a characteristic of vertebrates. Although not basically wrong, this definition necessitates some comments addressing three points: (1) antibody, (2) induction, and (3) foreignness.

### 1. Antibody

When this definition was proposed, the [immune response](#) was characterized solely by the production of antibodies. When it was realized, in the 1960, that the specific functions of the immune system relied on two distinct types of lymphocytes, [B cells](#) and [T cells](#), it became clear that the immune response could no longer be limited to the production of antibodies. B and T cells “recognize” the antigen, but in very different ways. B cells express [immunoglobulins](#) at their surface, whereas T cells express the so-called [T-cell receptor](#) (TCR). By extension, surface immunoglobulins are also given the name of B-cell receptor (BCR), by analogy with the TCR. Once stimulated by an antigen, in combination with other signals (**cytokines**), the B cells differentiate into plasma cells that will express a soluble form of immunoglobulins, also known as *antibodies* or *circulating antibodies*, because they are free in the bloodstream. Immunoglobulins and TCR are constructed on the same general pattern, but are encoded by discrete sets of **genes**. They interact with the antigen in very different manners. The immunoglobulins—or antibodies—interact directly with a native form of the antigen, through a small portion of this antigen called the *antigenic determinant* or [epitope](#). The TCR does not recognize the antigen as a whole, but interacts with a fragment of it, which results from **antigen processing** and presentation. This has been well studied for **protein** antigens, which are cleaved into [peptides](#) (the processing step) in the so-called antigen-presenting cells ([macrophages](#), dendritic cells, etc.), which will then reexpress these peptides at their surface, in close association with molecules encoded by genes of the [major histocompatibility complex](#) (MHC) (presentation step). Finally, the TCR will bind to both the peptide and the MHC-presenting molecule. The peptide derived from the antigen and that interacts with the corresponding TCR is described as the “T-epitope.”

## 2. Induction

For a long period of time, in fact until it became possible to work on the chemistry of antibodies in the beginning of the 1960, immunochemists centered their interest on the most accessible partner of the antigen–antibody complex—that is, the antigen, for which simple models of well-defined structure could be worked out. Although somewhat frustrating for a physiologist (this would be like an enzymologist exclusively studying substrates), this approach shed light on the exquisite specificity of recognition by the immune system, thanks mostly to the pioneering work of Landsteiner ([1](#)) and the discovery of haptens. Haptens are small molecules that are unable to stimulate the immune system, but may bind to an antibody. To produce this antibody, the hapten had to be conjugated to a protein that served as a carrier, conferring the property of stimulating the immune system. This observation clearly indicated that induction and recognition were two separate properties of an antigen. To clarify this, it was proposed to distinguish antigenicity from immunogenicity and, consequently, an antigen from an [immunogen](#). Antigenicity describes the chemical structures that condition interaction with an antibody or with a TCR, whereas immunogenicity defines the properties of a molecule to induce an immune response. This distinction is far from pure semantics. A protein, considered as a typical “natural” immunogen, has the dual characteristics of a carrier, which is involved in processing and presentation, and of a hapten (in fact a mosaic of haptens) represented by the B epitope. This duality is of central importance for the immune system, in that it implies a close cooperation between B, T, and antigen-presenting cells. In most cases, an antibody response requires cooperation with the T-cell compartment, leading to the distinction of T-dependent antigens, as opposed to T-independent ones, which can stimulate B cells directly without this T-cell “help.”

## 3. Foreignness

At the end of the nineteenth century, the first antigens that were identified were pathogens: bacteria, **viruses**, parasites, or [toxins](#) of various origins. It was soon realized that pathogenicity and immunogenicity were not linked and that a huge variety of cells or molecules could induce an immune response, providing that they were foreign to the animal it was injected into. Furthermore, this requirement of foreignness was strengthened by the famous aphorism “horror autotoxicus” put forward in 1901 by Ehrlich, to indicate that the immune system could not be stimulated by self-components, leading to the concept that it was a master of self–nonself discrimination. Due to

independent observations of Grubb, Oudin, and Kunkel (2-4), it became apparent that this distinction between self and nonself was not as clear as initially thought. Oudin described three types of antigenic specificities, termed *isotypy*, *allotypy*, and *idiotypy*. Isotypy refers to those specificities that are characteristic of one molecule of one given animal species, for instance the mouse serum albumin. An antibody raised against albumin of one given mouse will react with the albumin of any mouse. Allotypy defines antigenic specificities that are shared by a group of individuals within a given species, as initially shown by Oudin for the rabbit immunoglobulins. This is the case of the human blood groups, for which people of, say, group A share this specificity. Allotypy simply reflects the existence of epitopes (or allotopes) that are encoded by allelic variants (see [Alloantibody](#), [Alloantigen](#)). The last type of specificity, idiotype, is more subtle, and is inherent to the molecules of the immune system itself. It was initially described by Oudin as the characteristics of one antibody molecule, synthesized by one given animal and specific for one given antigen. So, if a rabbit is given antigen X, it will produce an anti-X antibody (or Ab1). If a second rabbit, expressing the same immunoglobulin allotypes, is immunized against the anti-X antibody, it will produce an anti-anti-X antibody (or Ab2). Ab1 is an [idiotype](#), Ab2 an **anti-idiotype**, and the antigenic specificities recognized on Ab1 by Ab2 are called the *idiotopes*. Because it was later shown that Ab1 and Ab2 could be produced as discrete waves within the same animal, it was clearly realized that the self–nonself distinction was not so obvious. The horror autotoxicus dogma had also to be revisited when it was shown by Avrameas, and extended by Coutinho (5), that a low level of natural autoantibodies was consistently present in every individual, as well as autoreactive T-cell clones. Such natural antibodies have no pathogenic effect, as opposed to those identified in certain [autoimmune diseases](#), for reasons that are still not completely understood.

The term *antigen* covers a huge number of structures, from cells and macromolecular complexes, to relatively small molecules, the lower limit in molecular weight being of the order of a few thousand. Ultimately, and whatever the source of the materials considered, proteins and, to a lesser extent, polysaccharides are the main antigens. Proteins are representative of the T-dependent antigens, which implies processing and presentation. Depending upon their size, the potential number of discrete epitopes of one given antigen may vary from a very small number to several tens, in roughly linear proportion with the exposed area of the molecule. Polysaccharides behave most frequently as T-independent antigens. They are frequently constituents of the bacterial surface and are thus important for the preparation of vaccines against these microorganisms. Another example of polysaccharide antigens is given by the major blood groups in humans. Some macromolecular complexes, such as bacterial lipopolysaccharides, have been studied extensively because they also behave as polyclonal potent mitogenic activators of B cells. Nucleic acids are considered poor immunogens, although they clearly constitute a target for autoantibodies in systemic lupus erythematosus, a severe autoimmune disease. More recently, DNA has been tentatively used as a vaccine, in the form of expression vectors encoding a protein that acted secondarily as an immunizing antigen.

Detailed analysis of the antigenic determinants or epitopes of natural antigens is a difficult task and has necessitated the use of models, among which the most extensively studied in the 1960s and 1970s were polysaccharides by the group of Kabat at Columbia University (6). Based on the binding inhibition of dextran (a polymer of glucose) to anti-dextran antibodies by oligosaccharides of various lengths, the range of size of an epitope, and hence that of the corresponding antibody combining site, was estimated to be between 3 and 6 monosaccharide units. A similar approach was taken by the group of Sela at the Weizmann Institute, with synthetic peptides that mimicked the structure of natural proteins. Direct analysis of natural peptide epitopes with the corresponding antibody was made much later by [X-ray crystallography](#) and gave a clear picture of the organization of the antibody combining site. The central message is that there is a huge diversity in size and form of natural epitopes. At the surface of a protein, the number of amino acid residues that interact with the antibody combining site varies from a few to about 20. The variation in both size and number of epitopes at the antigen surface contribute a large part of the heterogeneity of the immune response.

See also entries [Antigen Processing, Presentation](#), [Epitope](#), [Immunogen](#), [Idiotypes](#), and

## [Superantigens, Xenogeneic.](#)

### Bibliography

1. K. Landsteiner and J. van der Scheer (1936) On cross-reactions of immune sera to azoproteins. *J. Exp. Med.* **63**, 325–339.
2. R. Grubb (1956) Agglutination of erythrocytes coated with incomplete anti-Rh by certain rheumatoid arthritic sera and some other sera. *Acta Path. Microbiol. Scand.* **39**, 195–197.
3. J. Oudin (1960) Allotypy of rabbit serum proteins. I. Immunochemical analysis leading to the individualization of seven main allotypes. *J. Exp. Med.* **112**, 107–124.
4. M. Harboe, C. K. Osterland, and H. G. Kunkel (1962) Genetic characters of human  $\gamma$ -globulins in myeloma proteins. *J. Exp. Med.* **116**, 719–738.
5. A. Coutinho, M. D. Kazatchkine, and S. Avrameas (1995) Natural autoantibodies. *Curr. Opin. Immunol.* **7**, 812–818.
6. E. A. Kabat (1956) Heterogeneity in extent of the combining regions of human antidextran. *J. Immunol.* **77**, 377–380.

### Suggestions for Further Reading

7. Y. Paterson (1991) The structural basis of antigenicity. *Intern. Rev. Immunol.* **7**, 121–218.
8. M. van Regenmortel (1989) Structural and functional approaches to the study of protein antigenicity. *Immunol. Today* **10**, 266–271.

## Antigen Processing, Presentation

Antigen processing and presentation are crucial steps for recognition of protein [antigens](#) by [T cells](#). The [T-cell receptor](#) (TCR) will not bind a protein directly, as do [immunoglobulins](#), but only as peptides derived from the original protein (processing step) and bound to class I or class II molecules of the [major histocompatibility complex](#), or MHC (presentation step). This applies to helper as well as to effector [cytotoxic T lymphocytes](#) and therefore condition severely every aspect of the T-dependent [immune response](#). These functions strictly condition immunogenicity and correlate well with the dichotomy between carrier and hapten that had indicated long ago that the hapten part of the complex was recognized as such by the B cell, but was unable to initiate an immune response unless conjugated to a carrier protein (see [Antigen](#)). Later it was shown that carrier recognition involved the T helper cell but was also strictly dependent on MHC molecules (MHC restriction). Antigen processing and presentation take place in specialized cells and involve a number of molecules, most of which are encoded by genes of the MHC.

The antigen-presenting cells that are most efficient and most extensively studied are the dendritic cells and the epidermal Langerhans cells. They originate from CD34<sup>+</sup> bone marrow precursors and migrate to different tissues. They express FcR e I and FcRII, as well as class I and class II molecules of MHC (in humans, HLA-DR). They internalize the antigen by the **endosomal** pathway and, after processing, reexpress at the cell surface antigen-derived peptides associated with class I molecules. Alternatively, they may present peptides associated with class I MHC after macropinocytosis. They may become interdigitated cells and migrate to secondary lymphoid organs, where they present peptides specifically to CD4 T cells. B cells also behave as antigen-presenting cells, and this is of major importance for the antibody response induced by a T-dependent antigen, because it conditions the necessary direct interaction between T and B cells. Finally, [macrophages](#) and monocytes also can

present antigens, in addition to their other functions.

There are two main pathways for antigen processing and presentation. One involves association with MHC class I molecules and presentation to CD8 T lymphocytes, while the other involves MHC class II molecules and interaction with CD4 lymphocytes.

### 1. Class I Pathway

MHC Class I molecules are composed of one [polypeptide chain](#) of 43 kDa that forms a heterodimer with the  $\beta_2$  microglobulin. A peptide binding site, which accommodates nonapeptides, is shared by the first two **domains** of the  $\alpha$  chain. There is a huge degeneracy in peptide recognition by one such site, as expected from the low number of different MHC molecules expressed in any individual. Class I molecules are specialized in the presentation of endogenous proteins, synthesized *in situ* by a virus or any intracellular bacteria or parasite. Processing starts as a fraction of these cytosolic proteins are hydrolyzed in the [proteasome](#), a **proteinase-rich** complex involved in **protein degradation**. The resulting peptides will be transported by the specialized TAP gene products to class I molecules in the [endoplasmic reticulum](#). The class I peptide complex will pass through the [Golgi apparatus](#), where the  $\alpha$  chain becomes glycosylated before being included in a **secretory vesicle** and finally expressed at the plasma membrane. Interaction with the CD8 T cell involves the TCR, which makes contact with amino acid residues of both the peptide and the class I molecule, whereas CD8 binds solely to the MHC molecule. In addition to this specific interaction, many other [cell adhesion molecules](#) contribute to this “immunological synapse,” as immunologists sometimes call it.

### 2. Class II Pathway

MHC class II molecules are composed of two chains,  $\alpha$  and  $\beta$ , each containing two domains. The peptide binding site is shared by the  $\alpha 1$  and  $\beta 1$  domains. The  $\alpha\beta$  heterodimer associates with the Ii (CD74) invariant chain that prevents any association of an endogenous peptide with the class II molecule. After glycosylation during passage through the Golgi, the complex is packed in the class II vesicles that fuse with [endosomes](#), where the invariant Ii chain is hydrolyzed, leaving only the CLIP peptide still associated with the  $\alpha\beta$  heterodimer. The last step, controlled by HLA-DM molecules, involves an exchange between CLIP and a peptide derived from an exogenous protein antigen. The complex of peptide and class II molecules may now be expressed at the cell surface. A fraction of the class II molecules is reinternalized and may bind new peptides, before being reexpressed at the cell surface. The size of peptides that bind to the class II molecule may vary from between 10 and 30 amino acid residues. Interaction with T cells is somewhat similar to the case previously described, except that it involves essentially CD4 T cells.

CD4 T cells are helper cells that will either (a) interact with the T effector compartment and promote the emergence of cytotoxic T cells or (b) interact with B cells that will ultimately produce antibodies. The latter interaction is of particular interest, because the B cell will act as both a presenting cell and an antibody-producing cell. The T-cell–B-cell interaction is initiated because the B cell presents peptides derived from the specific antigen that had been internalized after binding to the surface immunoglobulin. The antigenic peptide–class II complex will trigger the interacting TCR and activate the T-cell CD3 signaling module. Critical molecules are then produced by the T cell, amongst which are soluble cytokines and one membrane ligand, CD40L, which will bind to its receptor, the CD40 molecule, which is constitutively expressed on the B cell. This is the major signal that will ultimately activate the final phase of differentiation of the B cell and provide the key to clonal expansion and antibody production.

See also entries [B Cell](#), [Antibody](#), [T Cell](#), [T-Cell Receptor \(TCR\)](#), [Antigen](#), [Immunogen](#), and [Immune Response](#).

Suggestions for Further Reading

J. L. Whitton (1998) An overview of antigen presentation and its central role in the immune response. *Curr. Top. Microbiol. Immunol.* **232**, 1–13.

R. M. Steinman and J. Banchereau (1998) Dendritic cells and the control of immunity. *Nature* **392**, 245–252.

## Antigenic Variation

Antigenic variation is a sophisticated molecular mechanism that allows parasites, **viruses**, and some bacteria to escape the [immune response](#) of the invaded host, by changing the nature of their expressed surface [antigens](#). Classical examples are (a) **trypanosomes** amongst parasites and (b) **influenza virus**.

Trypanosomes are infectious agents transmitted by glossinas and are responsible for severe tropical diseases, such as sleeping sickness in humans. Trypanosomes express surface **glycoproteins** that induce an immune response in the host. Although at any given time one trypanosome expresses one surface glycoprotein, it may switch to the expression of another one, having a slightly different structure. There are large numbers of these different forms (up to 1000), called **variable surface glycoproteins** (VSG), each being encoded by a distinct **gene**. There are two known mechanisms that account for the successive expression of different VSG genes. One is [transposition](#) that brings a gene that was present initially as part of one of several gene clusters, scattered on different chromosomes, to one [telomere](#) end. As a consequence of the transposition, that gene is expressed, until another one becomes functional. An alternative mechanism is a change in the control of [transcription](#) of expression sites. As the structure varies from one VSG to another, there is a permanent change in antigenicity of the parasite population, thereby ensuring escape from the immune mechanisms of the host. Although a large number of genes are potentially active, only a limited fraction is used, so that the [repertoire](#) of different VSGs expressed in one host remains limited; this may explain the appearance of chronic immunity in tropical populations.

Viruses use a different mechanism to escape the immune system. One example is that of influenza A, which expresses two surface antigens, hemagglutinin and neuraminidase. It is well known that new variants always arise, sometimes being highly dangerous as they spread rapidly. New variants are the result of [mutations](#), which induce *antigenic drift*, and exchange of genetic material between different virus strains, leading to *antigenic shift*, which has a more drastic impact on spreading of the virus. Antigenic variation of viruses is a major drawback in preventing infection (see [Virus Infection, Animal](#)). One approach is to prepare modified strains by [site-directed mutagenesis](#) and thus produce vaccines that might anticipate their natural variation.

A more recent problem of antigenic variation is that by [HIV](#), which mutates rapidly, leading to the rapid accumulation of numerous variants within the same patient and providing a particularly efficient way for the virus to escape the immune system. Added to the fact that lymphocytes themselves are a major target for this virus, this stresses the unusual danger of this virus.

### Suggestions for Further Reading

G. A. Cross, L. E. Wirtz, and M. Navarro (1998) Regulation of vsg expression site transcription and switching in *Trypanosoma brucei*. *Mol. Biochem. Parasitol.* **91**, 77–91.

S. Subbarao and G. Schochetman (1996) Genetic variability of HIV-1. *AIDS* **10** (Suppl. A), S13–S23.

## Anti-Idiotype Immunoglobulins

Every biologist is familiar with the fact that injection of an [antigen](#) into an animal induces an [immune response](#) that results, among other things, in the production of specific [antibodies](#). So, if a mouse is given antigen X, it will produce the corresponding anti-X antibody. If we now purify this anti-X antibody and inject it into a mouse of the same strain, (ie, expressing the same **allotypes** as the first mouse), the second mouse will produce an antibody that will react with the anti-X antibody and that represents an anti-anti-X antibody. This defines the idiotypic specificities, which were first described independently in the 1960s by Oudin (1) in the rabbit and by Kunkel (2) in humans. These specificities are located on the **variable regions** of both the heavy and the light chains and are the result of epitope(s) [in this case we use the term of idiotope(s)], related at least in part to the antibody combining site of anti-X. It was shown subsequently that this game could be played one step further: If we now inject a third mouse with the anti-anti-X molecule, the mouse will make anti-anti-anti-X antibodies. Immunologists proposed a simplified nomenclature that reads

$X \rightarrow Ab1 \rightarrow Ab2 \rightarrow Ab3 \dots$

where Ab1 is the idiotype, Ab2 is the anti-idiotype, and Ab3 is the anti-anti-idiotype. This constitutes the so-called idiotypic cascade. One observation of particular interest, independently made by the groups of Oudin in Paris and Urbain (3) in Bruxelles, was that some antibodies in the Ab3 population behaved like Ab1, in that they could interact with antigen X, which the animal had never “seen.” This led them to consider that Ab2 had structures resembling X, because it could induce the formation of antibodies (Ab3) that behaved as anti-X. Ab2 was thus considered an “internal image” of the antigen. The observation was generalized by Jerne (4), once it was shown that this cascade could take place spontaneously within the same animal and that the level of expression of each antibody could be modulated (up or down, depending upon the experimental conditions). Jerne proposed that the entire population of antibodies of a given individual were organized as a huge network of interactions, leading to a dynamic equilibrium of B lymphocytes. This was called the *idiotypic network*.

Modulation of the [repertoire](#) of [B cells](#) expressed at any given time by injection of either partner of the cascade was thoroughly investigated during a number of years. Particular attention was paid to the diverse roles exerted by the Ab2 antibodies, because it was demonstrated that two subpopulations, termed Ab2b and Ab2a, had completely antagonistic regulatory effects, resulting either in the production of Ab3 antibodies that were Ab1-like, or in the down regulation of Ab1. Ab2b are anti-idiotype antibodies that contain a typical internal image and are therefore able to stimulate the production of anti-X antibodies, without disposing of the original antigen. This led to the development of numerous attempts to make “idiotypic vaccines.” The second type of population, Ab2a, recognizes idiotopes that are outside the antibody combining site of Ab1, and they appeared quite suitable to repress undesirable antibodies (ie, [autoantibodies](#) that have a pathogenic effect in [autoimmune diseases](#)).

The theoretical interest of idiotypic vaccines is to function as surrogate for antigens that are difficult to use for [immunization](#), either because they are difficult to isolate, are poorly immunogenic, or present a potential hazard. Many attempts have been made, especially for **viruses**, including [HIV](#), and parasites such as schistosomiasis. Antibodies were generally produced, but it turned out that in most (if not all) cases they were poorly neutralizing, and the protective effect *in vivo* was too low relative to the high efficiency that is expected from an acceptable vaccine.

Down-regulation of pathogenic autoantibodies encountered in autoimmune disease has also been extensively investigated, and it is still currently used in therapeutics. In fact, one uses immunoglobulin preparations assembled from very large pools of donors that are administered intravenously (IVIg). Although experimental arguments are compatible with the presence of anti-idiotypes, many other factors may account for the down regulation of the patient antibodies.

Finally, an interesting application of the internal image is the isolation of cellular receptors by Ab2b that mimic the corresponding ligand. Numerous examples have been reported and successfully used; the first was the Ab2b of the [insulin](#) cascade, which was shown to activate the insulin receptor of the b islet of pancreas. Adrenergic receptors or T3 and T4 [thyroid hormone](#) receptors are other examples. This could facilitate receptor isolation by the use of immunoabsorbents prepared with these anti-idiotypic internal image antibodies.

See also entries [Autoantibody](#), [Autoimmunity](#), [Idiotypes](#).

### Bibliography

1. J. Oudin and M. Michel (1969) Idiotype of rabbit antibodies. Comparison of idiotype of various kinds of antibodies formed in the same rabbit against *S. typhi*. *J. Exp. Med.* **130**, 619–629.
2. H. G. Kunkel, M. Mannick, and R. C. Williams (1963) Individual antigenic specificity of isolated antibodies. *Science* **140**, 1218–1220.
3. J. Urbain, M. Wilker, J. D. Frannsen, and C. Collignon (1977) Idiotypic regulation of the immune system by the induction of antibodies by anti-idiotypic antibodies. *Proc. Natl. Acad. Sci. USA* **74**, 5126–5129.
4. N. K. Jerne (1974) Towards a network theory of the immune system. *Ann. Immunol.* **125 C**, 373–389.

### Suggestions for Further Reading

5. C. A. Bona (1996) Internal image concept revisited. *Proc. Soc. Exp. Biol. Med.* **213**, 32–42.
6. R. J. Poljak (1994) An idiotope-anti-idiotope complex and the structural basis of molecular mimicking. *Proc. Natl. Acad. Sci. USA* **91**, 1599–1600.

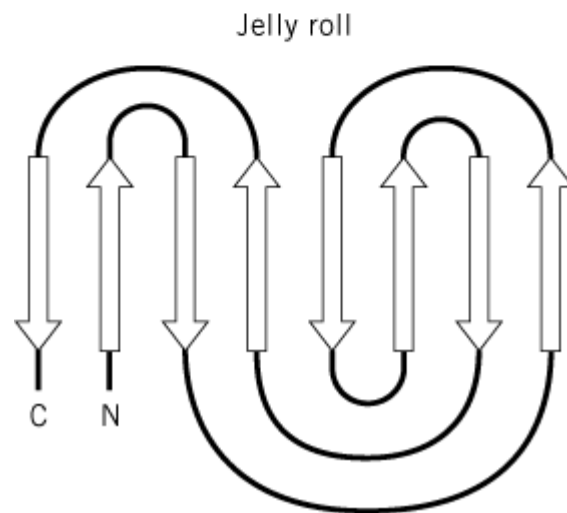
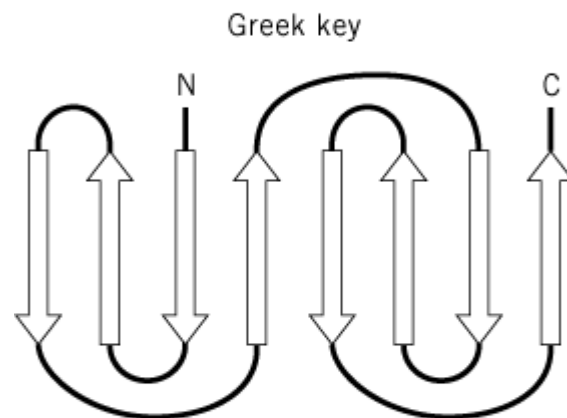
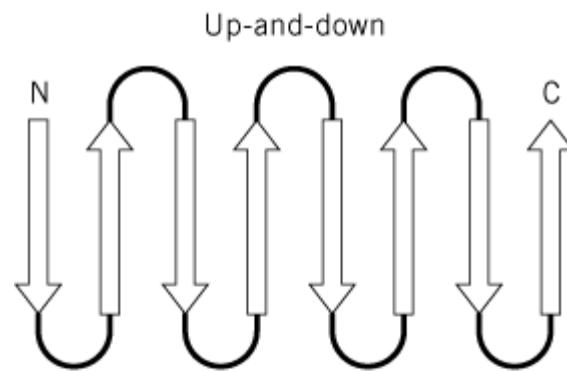
## Antiparallel Beta-Barrel Motifs

Antiparallel b-barrel motif is a general term that describes several [protein motifs](#) observed in [protein structures](#) that have in common a barrel-shaped [b-sheet](#) structure formed from antiparallel [beta-strands](#). These include the [b-barrel](#), the **Greek key** barrel (formed from two Greek key motifs), and the **jelly roll** barrel (Fig. 1). The different motifs have different connections between the strands in the barrel. Other types of antiparallel b-barrels include the [interleukin-1 motif](#) and the **neuraminidase** superbarrel. All these protein barrels comprise only antiparallel b-strands and are therefore classified as all-b structural **domains**. They therefore differ significantly from the [TIM barrel](#), which is classified as an a/b domain and is formed from a parallel b-barrel surrounded on the exterior by a-helices.

**Figure 1.** Schematic representation of the up and down (or b-meander), Greek key, and jelly roll topologies.



Antiparallel  $\beta$ -barrels can be formed from each of these three protein motifs, which represent different ways of linking consecutive  $\beta$ -strands together into a  $\beta$ -sheet.



[See also [b-Barrel](#), [Greek Key Motif](#), [Jelly Roll Motif](#), and [TIM Barrel](#).]

## Antisense Oligonucleotides

Antisense oligonucleotides are used to **hybridize** to a specific **RNA** molecule *in vivo* and thereby to inhibit its subsequent use. In most cases, the target RNA is a [messenger RNA](#) (mRNA), which cannot then be translated into a [protein](#). All that is required is knowledge of the gene's nucleotide sequence, so that an antisense oligonucleotide with the complementary sequence can be synthesized. In this way, the expression of one **gene** can be blocked, and the consequent observed effects help to elucidate the physiological role of the gene and its product.

Selection of the sites in a RNA molecule at which optimal antisense activity may be induced is complex, dependent on the terminating mechanism and influenced by the chemical class of the oligonucleotide. Each RNA displays a unique pattern of sites of sensitivity. Within the phosphorothioate oligodeoxynucleotide class, studies have shown that antisense activity can vary from undetectable to 100% by shifting an oligonucleotide by just a few bases in the RNA target ([1](#), [2](#)). Significant progress has been made in developing general rules that help to define potentially optimal sites in RNA species, but to a large extent this must be determined empirically for each RNA target and every new chemical class of oligonucleotides.

## 1. Therapeutic Uses

Besides being useful for characterizing the roles of genes, antisense oligonucleotides have the potential for therapeutic use to regulate the expression of certain genes or to block an infection by a pathogenic **bacterium** or **virus**. The therapeutic use of oligonucleotides represents a new paradigm for drug discovery, because oligonucleotides have not been studied before as potential drugs and they are being used to intervene in processes that have not been considered sites at which drugs might act. The affinity and the specificity of binding derive from the hybridization of two oligonucleotides and therefore are theoretically much greater than can be achieved with small molecules. Furthermore, the rational design of the nucleotide sequence of antisense oligonucleotide is much more straightforward than the design of small molecules interacting with proteins. Finally, it is possible to consider the design of antisense drugs to treat a very broad range of disorders, including those not amenable to other types of treatment.

This use of antisense oligonucleotides is still in its infancy, and the initial enthusiasm must be tempered by appropriate reservations concerning practical aspects. To be useful as a drug, an antisense oligonucleotide must be much more stable than ordinary nucleic acids, and it must be able to reach its desired site of action, the interior of a cell. Questions about the technology reduce to, can oligonucleotide analogs be created that have appropriate properties to be drugs? Specifically, what are the pharmacokinetic, pharmacological, and toxicological properties of these compounds and what are the scope and potential of the medicinal chemistry of oligonucleotides? Answers to many of these questions are available now.

## 2. Phosphorothioate oligodeoxynucleotides

To address the problem of stability, nonconventional oligonucleotides not susceptible to degradation or to hydrolysis by nucleases have been designed. Of the first generation of oligonucleotide analogs, the phosphorothioate class is best understood and has produced the broadest range of activities ([3](#)). Phosphorothioate oligonucleotides differ from normal in that one of the nonbridging oxygens in the phosphate group is replaced by a sulfur. The resulting compound is negatively charged, chiral at each phosphorothioate, and much more resistant to nucleases than a phosphodiester bridge.

### 2.1. Hybridization

The hybridization of phosphorothioate oligonucleotides to DNA and RNA has been thoroughly characterized. The melting temperature ( $T_m$ ) of a phosphorothioate oligodeoxynucleotide bound to RNA is lower than for a corresponding phosphodiester oligodeoxynucleotide by approximately 0.5° C per nucleotide. Compared to RNA duplex formation, a phosphorothioate oligodeoxynucleotide has a  $T_m$  approximately 2.2°C lower per nucleotide. This means that, to be effective *in vitro*,

phosphorothioate oligodeoxynucleotides must be relatively long, typically at least 17 to 20 nucleotides, and invasion of double-stranded regions of the target RNA is difficult (4-7).

### 2.1.1. Interactions with Proteins

Phosphorothioate oligonucleotides bind to proteins. These interactions can be (3) nonspecific, sequence-specific, or structure-specific, each of which may have different characteristics and effects. Nonspecific binding to a wide variety of proteins has been demonstrated, most thoroughly with [serum albumin](#). The affinity of such interactions is low. The **dissociation constant** ( $K_d$ ) for albumin is approximately 200  $\mu\text{M}$ , about the same as for its binding of aspirin or penicillin (8-10). Phosphorothioates also interact with **nucleases** and **DNA polymerases**; they are slowly metabolized by both endo and **exonucleases** and are **competitive inhibitors** of these [enzymes](#) (11). In an RNA–DNA duplex, phosphorothioates are substrates for [ribonuclease H](#) (RNaseH) (12). At higher concentrations, phosphorothioates inhibit the enzyme, presumably by binding as a single strand (11). Again, the oligonucleotides are competitive antagonists for the DNA–RNA substrate.

Phosphorothioates are competitive inhibitors of DNA polymerase  $\alpha$  and  $\beta$  with respect to the DNA template and **noncompetitive inhibitors** of DNA polymerases  $\gamma$  and  $\delta$  (12). They are also competitive inhibitors for the **reverse transcriptase** of [HIV](#) and inhibit its associated RNase H activity (13, 14). They bind to the cell surface protein CD4 and to protein kinase C (15). Phosphorothioates inhibit various viral polymerases (16), and they also cause potent, nonsequence-specific inhibition of [RNA splicing](#) (17).

### 2.1.2. In Vivo Pharmacokinetics

Binding of phosphorothioate oligonucleotides to serum albumin and alpha-2 [macroglobulin](#) provides a repository for these drugs in the serum and prevents their rapid renal excretion. Because serum protein binding is saturable, however, intact oligomer may be found in urine with high doses (18, 19), eg, 15 to 20  $\text{mg kg}^{-1}$  administered intravenously to rats. Phosphorothioate oligonucleotides are rapidly and extensively absorbed after parenteral administration, as much as 70% within 4 h (20, 21). Distribution of phosphorothioate oligonucleotides from blood after absorption or intravenous administration is extremely rapid. Distribution **half-lives** are less than one hour (18-20, 22). Clearance from the blood and plasma exhibits complex kinetics, with a terminal elimination half-life of 40 to 60 h in all species except man, where it may be somewhat longer (21). Phosphorothioates distribute broadly to all peripheral tissues, although no evidence for significant penetration of the blood brain barrier has been reported. Liver, kidney, bone marrow, skeletal muscle, and skin accumulate the highest amounts (20, 22). Liver accumulates the drug most rapidly (20% of a dose within 1 to 2 hours) and also eliminates it most rapidly (eg, the terminal half-life from liver is 62 h and from renal medulla 156 h). Within the kidney (23), oligonucleotides are probably filtered by the glomerulus, then reabsorbed by the proximal convoluted tubule epithelial cells, perhaps mediated by interactions with specific proteins in the brush border membranes. At relatively low doses, clearance of phosphorothioate oligonucleotides is caused primarily by metabolism (19, 20, 22), mediated by exo- and endonucleases.

### 2.1.3. Pharmacological Activities

Phosphorothioates also have effects inconsistent with the antisense mechanism for which they were designed. Some of these effects are sequence- or structure-specific. Others result from nonspecific interactions with proteins. These effects are particularly prominent in *in vitro* tests for antiviral activity, when high concentrations of cells, viruses, and oligonucleotides are often incubated together (24, 25). Human immune deficiency virus (HIV) is particularly problematic, because many oligonucleotides bind to the gp120 protein (26). Moreover, uncertainty as to the mode of action of antisense oligonucleotides is certainly not limited to antiviral or just *in vitro* tests (27-29). These observations indicate that, before drawing conclusions, careful analysis of dose-response curves, direct analysis of target protein or RNA, and inclusion of appropriate controls are required. In addition to interactions with proteins, other factors can contribute to unexpected results, such as overrepresented sequences of RNA and unusual structures that may be adopted by oligonucleotides

(26).

A relatively large number of reports of *in vivo* activities of phosphorothioate oligonucleotides have now appeared documenting activities after both local and systemic administration (30). However, for only a few of these reports have sufficient studies been performed to draw relatively firm conclusions concerning the mechanism of action by directly examining target RNA levels, target protein levels, and pharmacological effects, using a wide range of control oligonucleotides and examining the effects on closely related isotypes (31-34). Thus, there is a growing body of evidence that phosphorothioate oligonucleotides induce potent systemic and local effects *in vivo*, suggesting highly specific effects difficult to explain via any mechanism other than antisense.

#### 2.1.4. Toxicological Properties

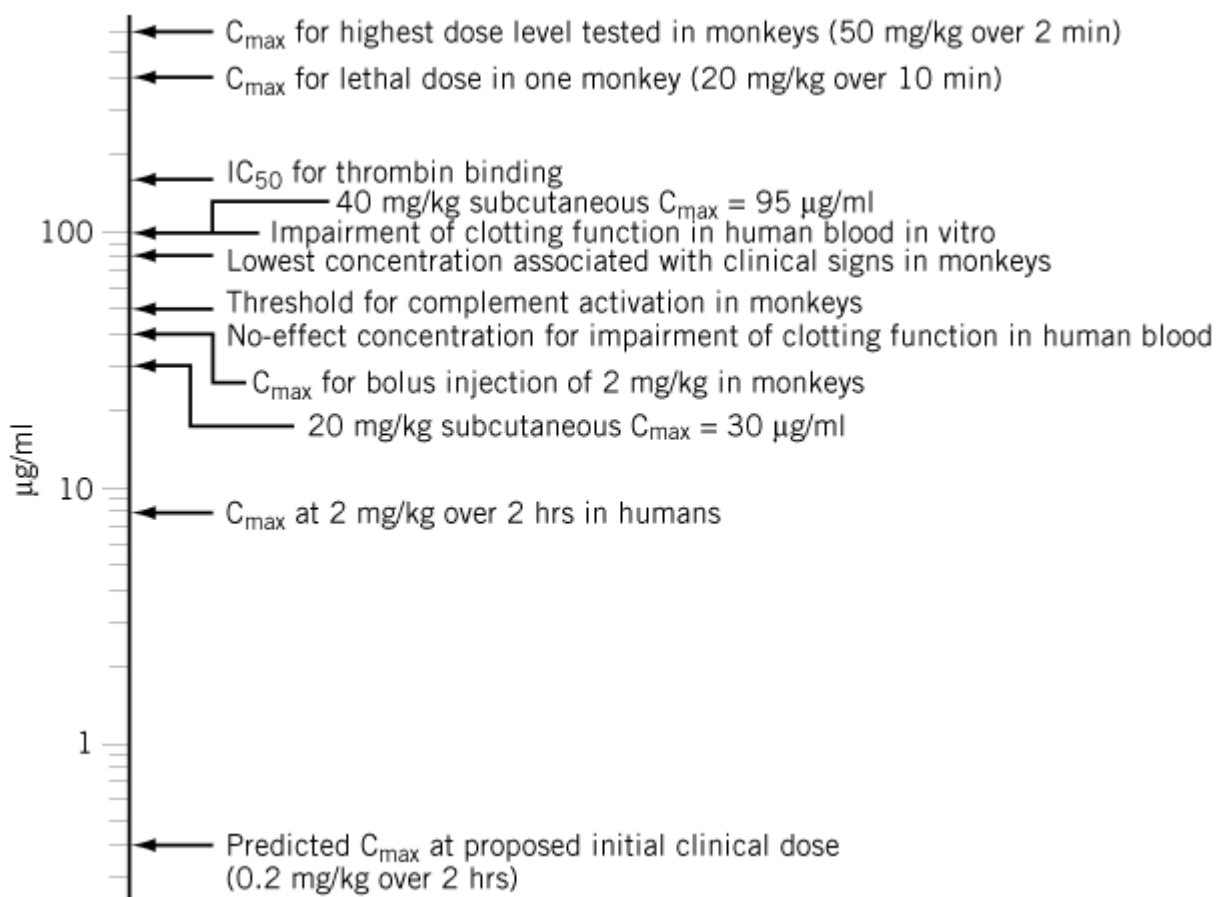
Phosphorothioate oligonucleotides are not toxic in themselves, but their use is limited by side effects. In rodents, this is immune stimulation (35, 36). In monkeys, it is sporadic reductions in blood pressure associated with bradycardia, which is often associated with activation of C-5 **complement** involving activation of the alternative complement pathway (37). All phosphorothioate oligonucleotides tested to date induce these effects, although there are slight variations in potency depending on their sequence and/or length (36, 38, 39). A second prominent toxicological effect in the monkey is on [blood clotting](#). The mechanisms responsible for these effects are probably very complex, but preliminary data suggest that direct interactions with [thrombin](#) are at least partially responsible (40).

In man, the toxicological profile differs. When ISIS 2922 is administered intravitreally to patients with [cytomegalovirus](#) retinitis, the most common adverse event is anterior chamber inflammation, which is easily managed with steroids. A relatively rare and dose-related adverse event is morphological changes in the retina associated with loss in peripheral vision (41). ISIS 2105, a 20-mer phosphorothioate designed to inhibit the replication of human papilloma viruses that cause genital warts, has been administered intradermally at doses as high as 3 mg/wart weekly for three weeks. Essentially no toxicities have been observed (42).

#### 2.1.5. Therapeutic Index

An attempt to put the toxicities and their dose-response relationships into a therapeutic context is shown in Fig. 1. This is particularly important because considerable confusion has arisen concerning the potential utility of phosphorothioate oligonucleotides for selected therapeutic purposes as a result of unsophisticated interpretation of toxicological data. As can be readily seen, the immune stimulation induced by these compounds is particularly prominent in rodents and unlikely to be dose-limiting in humans. Nor have hypotensive events in humans been observed to date. This toxicity occurs at lower doses in monkeys than in man and certainly is not dose-limiting in man. On the basis of present experience, the dose-limiting toxicity in man is likely to result from blood clotting abnormalities, associated with peak plasma concentrations well in excess of 10 µg/ml. Thus, it phosphorothioate oligonucleotides have a therapeutic index that supports their evaluation for a number of therapies.

**Figure 1.** Plasma concentrations of ISIS 2302 at which various activities are observed. These concentrations are those of intact ISIS 23902 and were determined by extracting and analyzing plasma by [capillary zone electrophoresis](#).



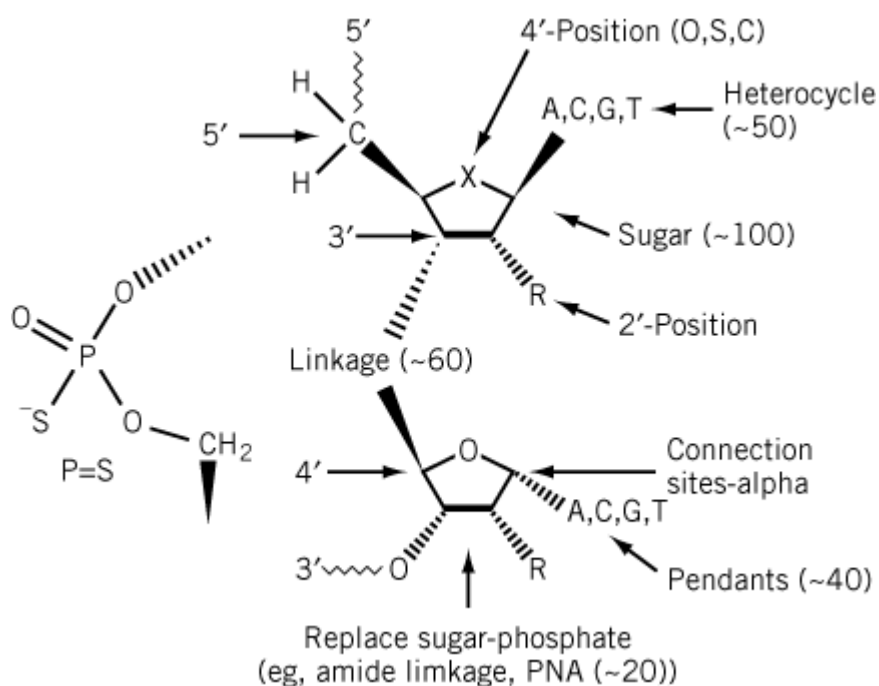
### 2.1.6. Clinical Activities

Significant therapeutic benefit has been reported in patients with cytomegalovirus retinitis treated locally with fomivirsen (43). ISIS 2302 administered every other day for one month also resulted in statistically significant improvement for five to six months in patients with steroid-dependent Crohn's disease in a randomized, double-blind placebo, controlled trial (44).

### 2.2. Medicinal Chemistry of Oligonucleotides

The core of any rational drug discovery program is medicinal chemistry. Although the synthesis of modified **nucleic acids** has been a subject of interest for some time (see [DNA Synthesis](#)), the intense focus on the medicinal chemistry of oligonucleotides dates perhaps no more than the past five years. Modifications have been made to the base, sugar, and phosphate moieties of oligonucleotides (Fig. 2). The subjects of medicinal chemical programs include approaches to (3) create enhanced and more selective affinities for RNA or duplex structures, the ability to cleave nucleic acid targets, enhanced nuclease stability, (4) cellular uptake and distribution, and *in vivo* tissue distribution, metabolism, and clearance. Arguably, the most interesting modifications to date are those that alter the sugar moiety (45) and the backbone. Modifications such as 2' methoxyethoxy enhance affinity for RNA, potency *in vivo*, provide a dramatic increase in stability, and reduce the potency for blood clotting and inflammatory effects. Also of interest are a number of modifications that replace the phosphate or the entire phosphate sugar backbone. Several novel chemical classes are being evaluated in animals and will shortly be studied in man, so it seems likely that a variety of chemical classes with differing properties will be available in the near future.

**Figure 2.** Isis oligonucleotide modifications.



## Bibliography

1. M. Y. Chiang et al. (1991) *J. Biol. Chem.* **266**, 18162–18171.
2. C. F. Bennett and S. T. Crooke (1996) "Oligonucleotide-based inhibitors of cytokine expression and function". In *Therapeutic Modulation of Cytokines* (B. Henderson, and M. W. Bodmer eds.), CRC Press, Boca Raton, pp. 171–193.
3. E. De Clercq, F. Eckstein, and T. C. Merigan (1969) *Science* **165**, 1137–1140.
4. W. F. Lima et al. (1992) *Biochemistry* **31**, 12055–12061.
5. T. Vickers et al. (1991) *Nucleic Acids Res.* **19**, 3359–3368.
6. B. P. Monia et al. (1993) *J. Biol. Chem.* **268**, 14514–14522.
7. B. P. Monia et al. (1992) *J. Biol. Chem.* **267**, 19954–19962.
8. S. T. Crooke et al. (1996) *J. Pharmacol. Exp. Ther.* **277**(2), 923–937.
9. R. W. Joos and W. H. Hall (1969) *J. Pharmacol. Exp. Ther.* **166**, 113–118.
10. S. K. Srinivasan, H. K. Tewary, and P. L. Iversen (1995) *Antisense Res. Dev.* **5**(2), 131–139.
11. S. T. Crooke et al. (1995) *Biochem. J.* **312**(2), 599–608.
12. W. Y. Gao et al. (1992) *Mol. Pharmacol.* **41**, 223–229.
13. C. Majumdar et al. (1989) *Biochemistry* **28**, 1340–1346.
14. Y. Cheng, W. Gao, and F. Han (1991) *Nucleosides Nucleotides* **10**, 155–166.
15. C. A. Stein et al. (1991) *Acquired Immune Defic. Syndr.* **4**, 686–693.
16. C. A. Stein and Y. C. Cheng (1993) *Science*, **261**, 1004–1012.
17. D. Hodges and S. T. Crooke (1995) *Mol. Pharmacol.*, **48**, 905–918.
18. S. Agrawal, J. Temsamani, and J. Y. Tang (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7595–7599.
19. P. Iversen (1991) *Anticancer Drug Des.* **6**(6), 531–8.
20. P. A. Cossum et al. (1994) *J. Pharmacol. Exp. Ther.* **269**, 89–94.
21. S. T. Crooke et al. (1994) *Clin. Pharm. Ther.* **56**, 641–646.
22. P. A. Cossum et al. (1993) *J. Pharmacol. Exp. Ther.* **267**, 1181–1190.
23. J. Rappaport et al. (1995) *Kidney Int.* **47**, 1462–1469.

24. R. F. Azad et al. (1993) *Antimicrob. Agents Chemother.* **37**(9), 1945–1954.
25. R. W. Wagner et al. (1993) *Science* **260**, 1510–1513.
26. J. R. Wyatt et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1356–1360.
27. C. M. Barton and N. R. Lemoine (1995) *Br. J. Cancer* **71**, 429–437.
28. T. L. Burgess et al. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 4051–4055.
29. M. Hertl, L. M. Neckers, and S. I. Katz (1995) *J. Invest. Dermatol.* **104**, 813–818.
30. S. T. Crooke (1995) *Therapeutic Applications of Oligonucleotides*, R. G. Landes, Austin, TX.
31. N. M. Dean and R. McKay (1994) *Proc. Natl. Acad. Sci. USA* **91**, 11762–11766.
32. N. M. Dean et al. (1996) *Cancer Res.* **56**(15), 3499–3507.
33. B. P. Monia et al. (1995) *Nature Med.* **2**(6), 668–675.
34. B. Monia et al. (1996) *J. Biol. Chem.* **271**(14), 14533–14540.
35. S. P. Henry et al. (1997) *Toxicology* **116**(1–3), p. 77–88.
36. S. P. Henry et al. *Antisense Nucleic Acid Drug Dev.*, In Press.
37. S. Henry et al. *J. Pharmacol. Exp. Ther.*, In Press.
38. K. G. Cornish et al. (1993) *Pharmacol. Commun.*, **3**, 239–247.
39. W. M. Galbraith et al. (1994) *Antisense Res. Dev.* **4**(3), 201–206.
40. S. Henry, W. Novotny, and J. Leeds. Submitted.
41. S. L. Hutcherson et al. (1995) Abstracts of the SFO, CA, *35th ICAAC*, p. 204.
42. J. M. Glover et al. (1996) *Pharmacol. Exp. Ther.*, In Press.
43. D. S. Boyer et al. (1997) In *AAO Annual Meeting*, San Francisco.
44. B. R. Yacyshyn et al. *N. Engl. J. Med.*, Submitted.
45. Y. S. Sanghvi, and P. D. Cook (1994) Symposium Series No. 580 American Chemical Society, Washington, DC, p. 232.

### Suggestions for Further Reading

46. S. T. Crooke and B. Lebleu (1993) *Antisense Research and Applications*, CRC Press, Boca Raton.
47. S. T. Crooke (1992) *Ann. Rev. Pharmacol. Toxicol.* **32**, 329–376.
48. S. T. Crooke (1993) *FASEB J.* **7**, 533–539.
49. S. M. Freier (1993) "Hybridization considerations affecting antisense drugs", In *Antisense Research and Applications* (S. T. Crooke and B. Lebleu, eds.), CRC Press, Boca Raton, pp. 67–82.
50. S. T. Crooke (1995) *Therapeutic Applications of Oligonucleotides*, R. G. Landes, Austin.

### Antisera

Antisera have long been the major reagent in immunology, and this field was for years known as serology. An antiserum is the serum of an animal that has been **immunized** against an [antigen](#), or more commonly hyperimmunized, to get the highest possible titer of the desired [antibodies](#). Very clearly, vaccination and serotherapy were the first great miracle in fighting against pathogens, long before the discovery of antibiotics. From a medical standpoint, antisera were used to fight quickly against a pathogen or a [toxin](#), because vaccination takes a long time to generate a host [immune](#)

[response](#). There are several concerns regarding the use of antisera for therapy in humans or in animals. First, one has to be certain that protection is indeed ensured by antibodies and not exclusively by cell-mediated immunity. This is not a trivial matter, because there are only a limited number of instances in which protection is due primarily to antibodies. Second, the use of antisera from a foreign animal species (heterologous antisera) will induce a strong immunization of the host against all the injected antigens. This results in the appearance, within a week or two, of a local and generalized syndrome, called *serum sickness*, that may include urticaria, local edema, rashes, arthralgia, fever, lymphadenopathy, and, ultimately, severe glomerulonephritis resulting from the formation of immune complexes. This was observed when serotherapy with horse serum was extensively used in humans, especially against **diphtheria** or tetanus. Diphtheria, a bacterial disease, can now be cured with antibiotics, and children undergo a regular schedule of vaccination that actively prevents the disease. Tetanus remains a very severe disease, with an elevated rate of mortality. Serotherapy is still used because antibodies are extremely effective against the toxin; nevertheless, everyone should be revaccinated every 10 years, because this provides the best and safest protection.

Heterologous antisera should be used only in very exceptional occasions—for example, after bites from highly dangerous snakes. Heterologous antisera are being systematically replaced, when appropriate, by purified [immunoglobulins](#), preferentially of human origin. Intravenous immunoglobulins (IVIgs) are used to correct severe immunodeficiencies of the [B-cell](#) compartment. For example, this is the case with X-linked agammaglobulinemia (Bruton disease), which is transmitted genetically through the mother and occurs in young boys. This disease is due to a variety of mutations of the BTK gene (for Bruton tyrosine kinase) that result in an early blockage of B-cell differentiation. IVIgs are also used to help patients with chronic or transient hypoglobulinemia. Whenever the deficiency is severe and would necessitate a lifelong treatment, an allogenic bone marrow graft is performed, which restores the immune system of the patient. IVIgs are also used in therapy of [autoimmune diseases](#), although the mechanisms are not absolutely understood, but may be due in part to a reequilibration of the [idiotype](#) network. Specific purified antibody of human origin can be prepared in some selected cases. The most popular preparation is the anti-D (Rhesus) immunoglobulins that are used worldwide for prevention of the hemolytic disease of the newborn (see [Alloantibody](#), [Alloantigen](#)).

With the start of organ transplantations, many attempts were made to prevent rejection by blocking the immune system of the recipient. Antilymphocyte serum, prepared in rabbits or in horses, has been used but had the usual limitation linked to the induction of serum sickness mentioned above. It could, however, help for a transient difficult acute period of rejection. It has been replaced by a [monoclonal antibody](#), prepared in the mouse and directed against the CD3 (signaling module) of the [T-cell receptor](#). Used for a short period of time, it proved efficient in down-modulating the immune response of the host, with limited risk of immunization against the heterologous immunoglobulin. Genetic engineering has also been proposed to minimize immunization against heterologous determinants. One classical approach is to insert the six specific hypervariable regions of a murine monoclonal antibody in place of the corresponding regions of a human antibody framework, using [protein engineering](#). This engineered antibody is certainly less immunogenic, but it still is, because of the [idiotype](#) determinants that cannot be avoided. Fully engineered human monoclonal antibodies would certainly be ideal for human use. This goal is being worked on actively by pharmaceutical companies. It is not yet practical, and no one has yet succeeded in making stable “natural” human monoclonal antibodies.

Besides the roles in human therapy mentioned above, antisera remain of course reagents of choice for experimental purposes. They may be used to identify new molecules and to purify them with techniques like [immunoaffinity chromatography](#) or **immunoprecipitation**, which are powerful tools to isolate rare components, such as [membrane proteins](#), receptors, [hormones](#), or diverse ligands. Monoclonal antibodies have tended to replace antisera, which is partly unfortunate and sometimes a mistake, because conventional antisera can be exquisitely specific and are very efficient, potent tools, with a mosaic of specificities that may of great help to the molecular biologist.



See also entries [Immunization](#) and [Immunogen](#).

#### Suggestion for Further Reading

E. Harlow and D. Lane, eds. (1988) *Antibodies: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

### Antitermination Control of Gene Expression

Bacteria have evolved many different complex mechanisms to control both [transcription](#) and [translation](#) of **genes** in response to environmental changes. In many cases, transcription is controlled at the level of initiation by [DNA-binding proteins](#) that either inhibit ([repressors](#)) or stimulate (activators) initiation. In addition, transcription can be regulated at the level of elongation. In some cases, transcription of a gene or [operon](#) will terminate prematurely in the absence of the action of a positive regulatory molecule. In these cases, antitermination factors allow transcription to read through termination signals and to generate full-length transcripts.

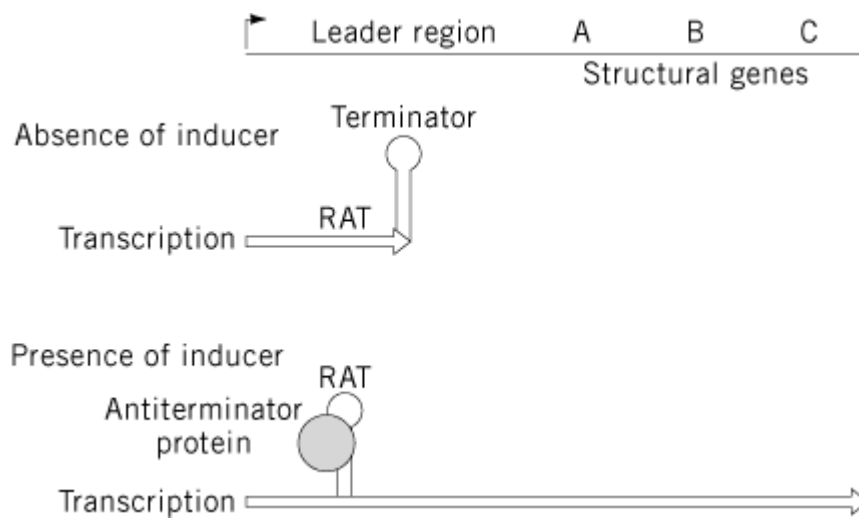
Two fundamentally different mechanisms for antitermination have been described. In one case, RNA polymerase is modified so as to allow it to read through transcription terminators. This type of mechanism controls phage development (1) and expression of rRNA operons (2). The second mechanism, covered in this chapter, involves **trans-acting factors** that interact with **RNA** and prevent formation of the terminator structure. This mechanism is very similar to **attenuation**, but antitermination can be distinguished from attenuation in that the action of the regulatory molecule results in transcription readthrough, with the default pathway being premature termination. In attenuation, the action of the regulatory molecule induces transcription termination, and the default pathway is readthrough.

Three distinct mechanisms that regulate gene expression by antitermination will be reviewed here. These mechanisms differ primarily in the type of biomolecule used as the regulator. The first mechanism uses antiterminator proteins that are activated to bind RNA targets in response to environmental stimuli. In the second mechanism, [transfer RNA](#) is used as the regulator. In this case, the degree of aminoacylation (or charging) of the tRNA is used to sense the availability of the cognate amino acid within the cell, to induce expression of genes involved in metabolism of this amino acid. Finally, in the case of the **Escherichia coli** tryptophanase operon, it appears that [ribosomes](#) are used as the regulatory molecule. Thus, bacteria have evolved a large number of mechanisms to use different biomolecules to all perform the same task—that is, to alter the conformation of the nascent mRNA to signal RNA polymerase whether it should terminate prematurely or continue transcription of the particular structural gene(s).

#### 1. RNA-Binding Protein-Mediated Antitermination: The Sac/Bgl Family of Antiterminator Proteins

Expression of several catabolic operons in bacteria is regulated by antitermination involving [RNA-binding proteins](#). These proteins prevent formation of Rho-independent transcription terminators in the nascent mRNA upstream of the regulated gene(s) (3). One such system in *E. coli* and several in *Bacillus subtilis*, appear to be highly related based on similarities of their antiterminator proteins as well as their RNA targets. In addition, several other systems function similarly but appear to have arisen independently. A general model for this mechanism is shown in Figure 1.

**Figure 1.** A general model for antitermination control by the Sac/Bgl family of antiterminator proteins. Under noninducing conditions, transcription starts at the promoter (designated by the arrow) and terminates prematurely, often in a leader region prior to the structural genes. In the presence of inducer, the antiterminator protein is activated to bind to the RAT (ribonucleic antiterminator) RNA. This binding stabilizes an RNA secondary structure involving the RAT, which prevents formation of the overlapping terminator, and transcription continues into the structural genes.



### 1.1. The *E. coli bgl* Operon

The *E. coli bglGFB* operon encodes all the functions necessary for the regulated uptake and utilization of aromatic  $\beta$ -glucosides. The operon is cryptic in wild-type strains but can become functional through spontaneous mutations. When functional, expression of this operon is regulated by antitermination mediated by the BglG protein in response to the levels of  $\beta$ -glucosides (4). In the absence of inducer, BglG does not bind RNA, and most transcripts terminate at one of two Rho-independent transcription terminators present in the leader region upstream of *bglG* and between *bglG* and *bglF*. When  $\beta$ -glucoside levels are high, BglG binds to an RNA target, named RAT for ribonucleic antiterminator, just upstream of the terminators. This binding stabilizes an alternative antiterminator RNA structure, which prevents formation of the terminator, thus allowing transcription to continue and the operon to be expressed.

The RNA-binding activity of BglG is regulated by **phosphorylation** mediated by BglF. In the absence of  $\beta$ -glucosides, BglF phosphorylates BglG, which prevents it from dimerizing and binding to the RAT (5). In the presence of  $\beta$ -glucosides, BglF dephosphorylates BglG, which now dimerizes and binds to the RAT. Phosphorylation of both  $\beta$ -glucosides and BglG is accomplished by transfer of the phosphate group from the same phosphorylated residue, Cys24, in BglF (6). These results suggest that, under conditions in which  $\beta$ -glucoside levels are high, the phosphate group can be transferred from BglG back to Cys24 in BglF. A model has been proposed in which unliganded BglF phosphorylates BglG, and  $\beta$ -glucoside binding induces BglF to undergo a conformational change that activates it to dephosphorylate BglG.

A similar system for  $\beta$ -glucoside utilization exists in the related **Gram-negative** enteric bacterium *Erwinia chrysanthemii*, although in this case the *arb* operon is not cryptic (7). ArbG shows high sequence similarity to BglG, suggesting that it functions analogously in antitermination control of the *E. chrysanthemii arb* genes. Antitermination also appears to control  $\beta$ -glucoside operons in several **Gram-Positive Bacteria** as well. A putative  $\beta$ -glucoside (*bgl*) operon has also been identified in *B. subtilis* and may be regulated by a similar antitermination mechanism (8). In addition, a protein, BglR, with homology to BglG also controls  $\beta$ -glucoside usage in *Lactococcus lactis* (9).

### 1.2. The *B. subtilis sac* Genes

Expression of two sucrose utilization operons in *B. subtilis*, *sacPA* (10) and *sacB* (11), is induced by sucrose via transcription antitermination mediated by the RNA-binding proteins SacT and SacY, respectively. SacT and SacY show extensive sequence similarity to each other, as well as to BglG from *E. coli*. The antitermination mechanisms that control these genes also appear to be quite similar to that described above for the *E. coli* *bgl* operon. Rho-independent transcription terminators exist in leader regions upstream of both *sacPA* and *sacB* and prevent transcription of the structural genes in the absence of the inducer, which is sucrose. In the presence of sucrose, SacT and SacY are activated to bind RAT sequences in the *sacPA* and *sacB* leader transcripts, respectively, and allow transcription to read through into the structural genes (12). Like BglG in *E. coli*, both of these antiterminator proteins are phosphorylated. In the case of SacY, phosphorylation negatively regulates RNA-binding activity and appears to be mediated by SacX (13). SacT is phosphorylated by HPr, a component of the phosphoenolpyruvate **phosphotransferase system**, but the role of this phosphorylation in sucrose-mediated antitermination is less clear (12).

Recently, the [protein structure](#) of the RNA-binding domain of SacY has been determined by both [NMR](#) (14) and [X-ray crystallography](#) (15). The domain exists as a dimer, with each monomer consisting of a four-stranded antiparallel [beta-sheet](#). Several amino acid residues have been identified through genetic, biochemical, and preliminary NMR studies as being important for RNA binding. These residues are clustered on the surface of one side of the protein structure (15).

### 1.3. Other Examples of Bgl/Sac Type Antiterminators

In addition to the *bgl* and *sac* systems described above, several other operons are regulated by RNA-binding antiterminator proteins with homology to BglG, SacY, and SacT. LicT regulates the *licS* gene, which is involved in  $\beta$ -glucan utilization in *B. subtilis* (16). There is also a RAT sequence overlapping a potential Rho-independent terminator upstream of *licS*.

In *Lactobacillus casei*, the lactose (*lac*) operon is regulated in response to lactose levels by LacT, which shows sequence homology to the other members of the Bgl/Sac family of antiterminators (17). The 5'-leader region of the *lac* mRNA contains a region with sequence similarity to the RAT sequence, as well as a potential stem-loop structure resembling a Rho-independent terminator.

### 1.4. Antiterminators with No Similarity to the Bgl/Sac Family

Several other systems are regulated by RNA-binding antiterminator proteins that are unrelated to those of the Bgl/Sac family; furthermore, these proteins do not appear to be related to each other. These regulatory systems thus appear to have arisen independently.

In *B. subtilis*, both the *glp* regulon, which is involved in usage of glycerol-3-phosphate, and a histidine-utilization (*hut*) operon are regulated by RNA-binding antiterminator proteins; GlpP (18) and HutP (19), respectively. The amino acid sequences of these antiterminator proteins are not similar to any other antiterminator proteins. Further, the mechanisms by which these antiterminator proteins function appear to be different from those described above, because there are no clear antiterminator RNA secondary structures near the terminators in these operons.

The amidase (*ami*) operon of *Pseudomonas aeruginosa* is regulated by antitermination in response to short-chain aliphatic amides, such as acetamide. The *amiR* gene encodes an antiterminator protein (AmiR), which is negatively regulated by AmiC, apparently through formation of an AmiC-AmiR complex (20). Acetamide destabilizes the AmiC-AmiR complex, leading to antitermination and expression of the operon. AmiR interacts with an RNA target in the 5-leader region of the *ami* mRNA that contains a Rho-independent terminator. However, no clear antiterminator RNA secondary structure is predicted. AmiR binding has been suggested to function in antitermination by interfering directly with formation of the terminator stem-loop structure (20).

In addition to all the catabolic operons described above, one anabolic operon has been shown to be regulated by antitermination. Expression of the *nas* operon of *Klebsiella pneumoniae*, which encodes enzymes required for nitrate assimilation in this bacterium, is induced by nitrate or nitrite. The NasR

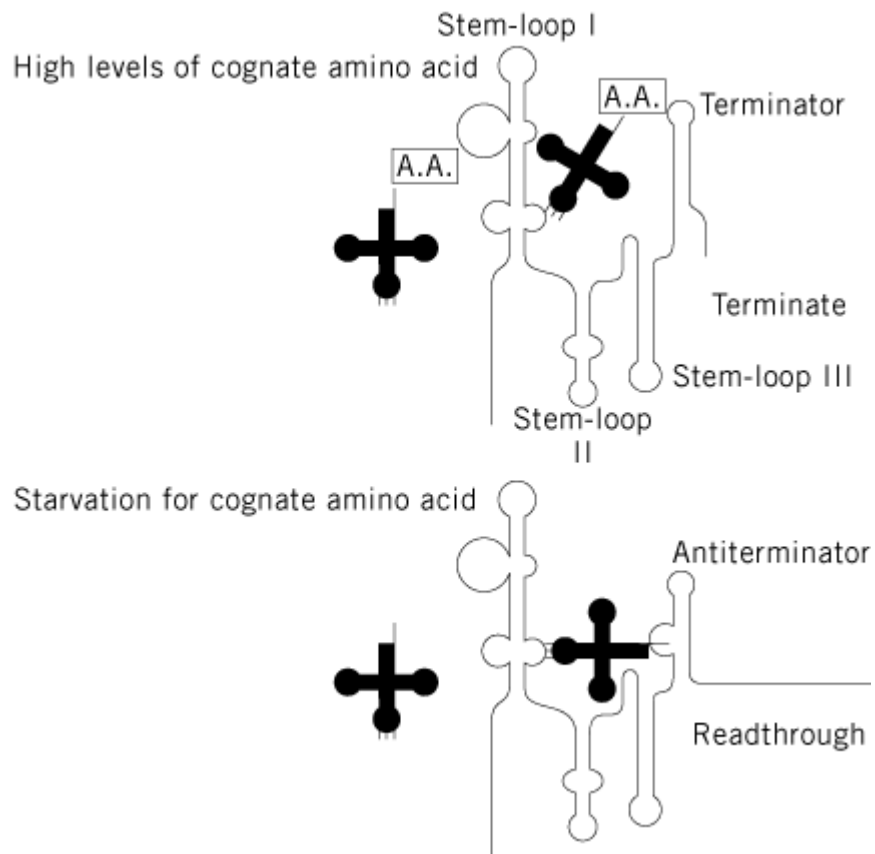
protein mediates transcription antitermination through a terminator in the leader region of the operon (21). This protein shows weak homology with AmiR in the carboxyl-terminal region.

## 2. Transfer RNA-Mediated Antitermination

An interesting variation on the antitermination mechanism involves the use of tRNA as the regulatory molecule. This mechanism regulates a large number of [aminoacyl-tRNA synthetase](#) genes in [Gram-Positive Bacteria](#) (22, 23) and several amino acid biosynthetic operons, including the *ilv-leu*, in *B. subtilis* (24, 25), and the *his* and *trp* operons in *Lactococcus lactis* (26). Expression of these genes is induced specifically by starvation for the corresponding amino acid. In the case of the amino acid operons, insufficient levels of the amino acid leads to increased expression of the corresponding biosynthetic operon. For the aminoacyl-tRNA synthetase genes, increasing the level of the synthetase is thought to allow more efficient charging of the cognate tRNA when the corresponding amino acid pool is low.

A long (approx. 300-nucleotide) untranslated leader region exists upstream of the structural gene(s) of these operons that contains several conserved features, including three stem-loop structures preceding a Rho-independent transcription terminator. Hence, in the absence of the inducing signal, transcription terminates prematurely in the leader region prior to the coding sequences. In addition to the conserved secondary structures, there is an important conserved 14-nucleotide sequence known as the T-box present in each leader region; hence these genes are known as the T-box family. An alternate arrangement of the leader region involving base-pairing between a portion of the T-box and a conserved sequence in the 5' side of the terminator stem has been proposed to form an antiterminator structure that allows transcription to read through into the structural genes (Fig. 2) (27).

**Figure 2.** Model for antitermination control by tRNA. Under conditions with adequate levels of the cognate amino acid (aa), the charged tRNA does not interact with the leader region, and the terminator forms. Under conditions of starvation for the appropriate amino acid, the uncharged tRNA interacts with the leader region via base-pairing between the anticodon and the specifier sequence, and by base-pairing between the CCA sequence at the acceptor end of the tRNA with the side bulge of the antiterminator in the leader. These interactions stabilize formation of the antiterminator conformation of the leader transcript, resulting in induction of expression of the gene. The tRNA is shown as the shaded cloverleaf structure, and a boxed “A.A.” attached to the tRNA indicates it is aminoacylated. Adapted from T. M. Henkin (1996) *Annu. Rev. Gen.* **30**, 35–57.



Another important conserved feature of the leader region of these genes is the presence of a triplet sequence corresponding to a codon for the appropriate amino acid for each operon. For example, in *tyrS*, which encodes tyrosyl-tRNA synthetase, the leader contains a UAC tyrosine codon, while the *ilv-leu* operon leader contains a CUC leucine codon. This triplet is always present in a bulged sequence in Stem-loop I (Fig. 2) and has been shown to be critical for induction in several systems. It was the presence of these triplets that led to the hypothesis that tRNAs play a role in this regulatory mechanism. This triplet was designated the “specifier sequence” because, in the case of the *B. subtilis tyrS* gene, altering the sequence to correspond to a codon for another amino acid switched induction to respond to starvation for the new amino acid (27). Other experiments demonstrated that translation of this codon was not involved in induction and that uncharged tRNA was the inducer (27). In addition, a second interaction between the CCA sequence at the 3' end of the uncharged tRNA and the complementary UGG sequence in the T-box have been shown to be important (28).

A model for tRNA-regulated antitermination regulation is presented in Fig. 2. Under starvation conditions for the corresponding amino acid, the cognate uncharged tRNA interacts with two sites in the leader region, to induce formation of the antiterminator structure and allow transcription to read through into the coding region. Aminoacylation of this tRNA is predicted to interfere with the interaction at the CCA end and prevent the charged tRNA from binding; the leader transcript then folds into the conformation with the terminator, halting transcription. It is not known if factors in addition to tRNA are required for antitermination. To date, however, it has not been possible to reconstitute tRNA-mediated antitermination in vitro, and several other lines of evidence also suggest that other factors may be involved in this mechanism (23).

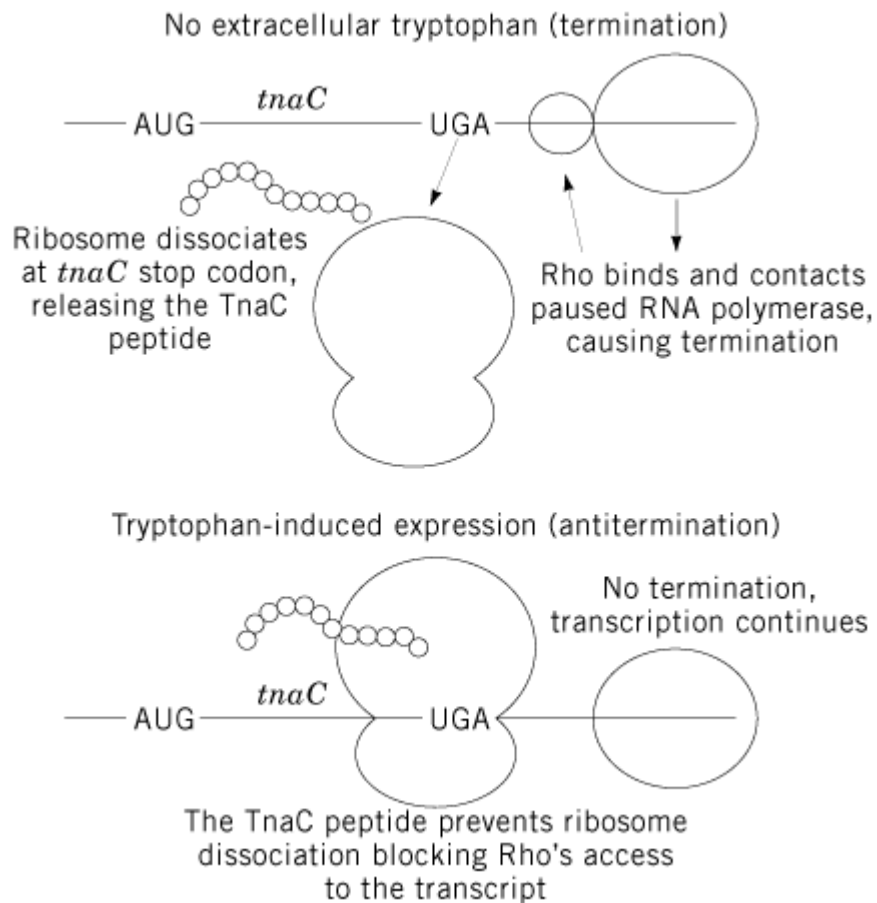
In addition to the antitermination mechanism described above, processing of the leader RNA has been shown to play a role in regulating expression of the *B. subtilis thrS* gene (29). Cleavage occurs in the loop of the antiterminator near the T-box sequence and is more efficient under threonine starvation conditions, suggesting that bound tRNA induces both antitermination and RNA

processing. This processing increases the stability of the mRNA, which would allow for increased translation and production of the threonyl-tRNA synthetase. Thus induction of expression of this gene in response to threonine starvation occurs at both the level of transcription antitermination and mRNA stability.

### 3. The *E. coli* Tryptophanase Operon

*E. coli* and several other microorganisms have the capacity to degrade tryptophan as a source of carbon, nitrogen, and/or energy (30). The degradative tryptophanase operon (*tnaCAB*) of *E. coli* is regulated by catabolite repression (31) and by an antitermination mechanism. Antitermination involves translation of a *cis-acting* 24-residue leader peptide (*tnaC*) containing a critical Trp codon (32, 33), one or more RNA polymerase pause sites between *tnaC* and *tnaA* (34), and Rho termination factor (34). While the precise antitermination mechanism responsible for controlling the *tna* operon is not firmly established, all of the data are consistent with the following model (Fig. 3) (35). During growth in a medium lacking both tryptophan and a catabolite-repressing carbon source, transcription initiation is efficient. As transcription proceeds, translation of the leader peptide occurs as soon as the coding sequence becomes available. Once the translating ribosome reaches the UGA stop codon, ribosome release exposes a *rut* (Rho utilization) site that immediately follows the stop codon. Rho then binds to the *rut* site and begins to translocate in the 3'-direction, until it encounters paused RNA polymerase, ultimately leading to transcription termination upstream of *tnaA*. When cells are growing with inducing levels of tryptophan, TnaC, or a complex of TnaC with an unidentified protein, prevents ribosome release at the *tnaC* stop codon, thereby masking the *rut* site and, hence, blocking Rho interaction with the transcript. Eventually RNA polymerase would overcome the pause signal and transcribe the structural genes encoding tryptophanase (*tnaA*) and a tryptophan permease (*tnaB*). This model assumes that there is a fundamental difference between the TnaC peptide, or the TnaC peptide-protein complex, in cells growing with or without tryptophan. It was proposed that such a complex under inducing conditions would prevent ribosome release (35), reminiscent of characterized translation attenuation mechanisms (36). The tryptophanase operon of *Proteus vulgaris* is thought to be regulated by a mechanism essentially identical to that of *E. coli* (37).

**Figure 3.** Model of *E. coli tna* operon regulation. Under noninducing conditions (no extracellular tryptophan), ribosome dissociation at the *tnaC* stop codon exposes a *rut* site, allowing Rho binding. Rho translocates to the paused RNA polymerase, leading to transcription termination. Under inducing conditions (extracellular tryptophan), ribosome stalling at the *tnaC* stop codon prevents Rho association, leading to transcription readthrough. See text for details. Adapted from C. Yanofsky, K. V. Konan and J. P. Sarsero (30).



## Bibliography

1. D. I. Freidman and D. L. Court (1995) *Mol. Microbiol.* **18**, 191–200.
2. C. Condon, C. Squires, and C. L. Squires (1995) *Microbiol. Rev.* **59**, 623–645.
3. B. Rutberg (1997) *Mol. Microbiol.* **23**, 413–421.
4. S. Mahadevan and A. Wright (1987) *Cell* **50**, 485–494.
5. O. Amster-Choder and A. Wright (1990) *Science* **249**, 540–542.
6. Q. Chen, J. C. Arents, R. Bader, P. W. Pestma, and O. Amster-Choder (1997) *EMBO J.* **16**, 4617–4627.
7. M. El Hassount, B. Henrissat, M. Chippaux, and F. Barras (1992) *J. Bacteriol.* **174**, 765–777.
8. D. Le Coq, C. Lindner, S. Krüger, M. Steinmetz, and J. Stulke (1995) *J. Bacteriol.* **177**, 1527–1535.
9. J. Bardowski, D. S. Ehrlich, and A. Chopin (1994) *J. Bacteriol.* **176**, 5681–5685.
10. M. Debarbouille, M. Arnaud, A. Fouet, A. Klier, and G. Rapoport (1990) *J. Bacteriol.* **172**, 3966–3973.
11. A. M. Crutz, M. Steinmetz, S. Aymerich, R. Richter, and D. Le Coq (1990) *J. Bacteriol.* **172**, 1043–1050.
12. M. Arnaud, M. Debarbouille, G. Rapoport, M. H. Saier Jr., and J. Reizer (1996) *J. Biol. Chem.* **271**, 18966–18972.
13. M. Idelson and O. Amster-Choder (1998) *J. Bacteriol.* **180**, 1043–1050.
14. X. Manival, Y. Yang, M. P. Strub, M. Kochoyan, M. Steinmetz, and S. Aymerich (1997) *EMBO J.* **16**, 5019–5029.
15. H. van Tilbeurgh, S. Manival, S. Aymerich, J. M. Lhoste, C. Dumas, and M. Kochoyan (1997)

EMBO J. **16**, 5030–5036.

16. K. Schnetz, J. Stulke, S. Gertz, S. Krüger, M. Krieg, M. Hecker, and B. Rak (1996) *J. Bacteriol.* **178**, 1971–1979.
17. C. A. Alpert and U. Siebers (1997) *J. Bacteriol.* **179**, 1555–1562.
18. E. Glatz, R. P. Rutberg, and B. Rutberg (1996) *Mol. Microbiol.* **19**, 319–328.
19. L. Wray Jr. and S. Fisher (1994) *J. Bacteriol.* **176**, 5466–5473.
20. S. A. Wilson, S. J. M. Wachira, R. A. Norman, L. H. Pearl, and R. E. Drew (1996) *EMBO J.* **15**, 5907–5916.
21. J. T. Lin and V. Stewart (1996) *J. Mol. Biol.* **256**, 423–435.
22. T. M. Henkin (1994) *Mol. Microbiol.* **13**, 381–387.
23. C. Condon, M. Grunberg-Manago, and H. Putzer (1996) *Biochimie* **78**, 381–389.
24. J. A. Grandoni, S. A. Zahler, and J. M. Calvo (1992) *J. Bacteriol.* **174**, 3212–3219.
25. J. A. Grandoni, S. B. Fulmer, V. Brizio, S. A. Zahler, and J. M. Calvo (1993) *J. Bacteriol.* **175**, 7581–7593.
26. P. Renault, J. J. Godon, C. Delorme, G. Corthier, and S. D. Erlich (1995) *Dev. Biol. Stand.* **85**, 431–441.
27. F. J. Grundy and T. M. Henkin (1993) *Cell* **74**, 475–482.
28. F. J. Grundy S. M. Rollins, and T. M. Henkin (1994) *J. Bacteriol.* **176**, 4518–4526.
29. C. Condon, H. Putzer, and M. Grunberg-Manago (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6992–6997.
30. C. Yanofsky, K. V. Konan, and J. P. Sarsero (1996) *Biochimie* **78**, 1017–1024.
31. J. L. Botsford and R. D. DeMoss (1971) *J. Bacteriol.* **105**, 303–312.
32. P. Gollnick and C. Yanofsky (1990) *J. Bacteriol.* **172**, 3100–3107.
33. K. Gish and C. Yanofsky (1995) *J. Bacteriol.* **177**, 7245–7254.
34. V. Stewart, R. Landick, and C. Yanofsky (1986) *J. Bacteriol.* **166**, 217–223.
35. K. V. Konan and C. Yanofsky (1997) *J. Bacteriol.* **179**, 1774–1779.
36. P. S. Lovett and E. J. Rogers (1996) *Microbiol. Rev.* **60**, 366–385.
37. A. V. Kamath and C. Yanofsky (1997) *J. Bacteriol.* **179**, 1780–1786.

### **Suggestions for Further Reading**

38. T. M. Henkin (1996) Control of transcription termination in prokaryotes. *Ann. Rev. Gen.* **30**, 35–57. A comprehensive review of the mechanisms of gene regulation by controlling transcription termination in bacteria.
39. T. Platt (1998) "RNA structure transcription elongation, termination and antitermination". In *RNA Structure and Function* (R. W. Simons and M. Grunberg-Manago, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 541–574. A more in-depth look at the processes involved in transcription and regulation of transcription termination.
40. M. Steinmetz (1993) "Carbohydrate catabolism: pathways, enzymes, genetic regulation, and evolution". In *Bacillus subtilis and Other Gram-positive Bacteria: biochemistry, physiology, and molecular genetics*. (A. L. Sonenshein, J. A. Hoch, and R. Losick, eds.), American Society for Microbiology, Washington, DC, pp. 157–170. A more in-depth look at the Sac systems in *B. subtilis*.

### **Antithrombin**



Antithrombin is a [proteinase inhibitor](#) of the [serpin](#) family that is the principal anticoagulant in human plasma; it is also known as antithrombin-III. It is composed of a single [polypeptide chain](#) of 432 amino acid residues, with four carbohydrate side chains (see [N-Glycosylation](#)). It has a reactive-center [peptide bond](#) between Arg393 and Ser394 that provides a specific target for [thrombin](#), factor Xa, and other [proteinases](#) of the [blood clotting](#) system. Antithrombin is present in plasma at a concentration of nearly 100 mg/L and circulates in a relatively inactive form, being activated conformationally by the binding of heparin. Genetic deficiency or dysfunction of antithrombin is a significant cause of familial thromboembolic disease.

#### Suggestions for Further Reading

S. T. Olson and I. Björk (1992) "Regulation of thrombin by antithrombin and heparin cofactor II", in *Thrombin, Structure and Function* (L. J. Berliner, ed.), Plenum Press, New York, pp. 159–217.

L. Jin, J. P. Abrahams, R. Skinner, M. Petitou, R. N. Pike, and R. W. Carrell (1997) The anticoagulation activation of antithrombin by heparin, *Proc. Natl. Acad. Sci. USA* **94**, 14683–14688.

### **$\alpha_1$ -Antitrypsin**

$\alpha_1$ -Antitrypsin is a [proteinase inhibitor](#) of the **serpin family** that consists of a single [polypeptide chain](#) of 294 amino acid residues and four carbohydrate side chains (see [N-Glycosylation](#)). It is also known as  $\alpha_1$ -proteinase inhibitor. Present at a concentration of 1 g/L, it is the predominant proteinase inhibitor in human plasma in terms of concentration. Its reactive-center [peptide bond](#) between Met358 and Ser359 specifically targets the [elastases](#) that are involved in connective tissue remodeling. Genetic deficiency of  $\alpha_1$ -antitrypsin is relatively common and predisposes to the development of the lung disease emphysema, with an associated risk of liver cirrhosis (see [Serpins](#)).

#### Suggestions for Further Reading

R. W. Carrell, J.-O. Jeppsson, C.-B. Laurell, S. O. Brennan, M. C. Owen, L. Vaughan, and D. R. Boswell (1982) Structure and variation of human  $\alpha_1$ -antitrypsin, *Nature* **298**, 329–334.

R. Mahadeva and D. A. Lomas (1998)  $\alpha_1$ -Antitrypsin deficiency, cirrhosis and emphysema, *Thorax* **53**, 501–505.

### **AP Endonucleases**

AP endonucleases are enzymes that cleave the phosphodiester bond on either side of an [AP site](#) in DNA and are involved in [base excision repair](#). Enzymes that cleave the phosphodiester bond to the 3'

side are referred to as type I AP endonucleases, and those that cleave the phosphodiester bond 5' to the AP site are type II AP endonucleases. Whether or not a bona fide type I AP endonuclease exists that is not associated with [DNA glycosylase](#) activity is not known. However, type II AP endonucleases, which cleave AP sites by hydrolysis, have been found in many organisms and have been characterized.

## 1. AP Lyases (Class I AP Endonucleases)

Type I AP endonucleases are also DNA glycosylases, and it appears that the same [active site](#) is involved in both the glycosylase and AP endonuclease activities. However, it has also been demonstrated that cleavage of AP site occurs by b-elimination, which involves abstraction of a hydrogen from the C2' position of the deoxyribose. Because this is a lyase reaction, and not a hydrolyase reaction, the type I AP endonucleases are referred to as AP lyases (1).

Endonuclease III and 8-oxoguanine glycosylase are examples of this class of enzymes (1, 2). The  $\alpha$ -NH<sub>2</sub> group of the *N*-terminal amino acid residue forms a [Schiff base](#) with the ring-opened form of the deoxyribose, eventually leading to cleavage of the deoxyribose–phosphate bond without water addition. Following elimination of the base by glycosylase action, the enzyme remains attached to DNA through a protonated Schiff base intermediate where the C1' of the sugar is covalently bound to the nitrogen of the *N*-terminal  $\alpha$ -amino group of the protein (3). b-Elimination of this intermediate, followed by hydrolysis of the Schiff base, results in cleavage of the phosphodiester bond and generation of a *trans*  $\alpha$ - $\beta$ -unsaturated aldehyde on the 5' side, and a 5'-phosphate on the 3' side, of the cleavage, which is concomitant with the release of the enzyme. When the AP lyase activity is associated with a glycosylase activity, the two reactions are, for the most part, coupled (2). In these instances, it is natural to assign a physiological role to the AP lyase. However, cleavage of an AP site on the 3' side occurs readily even in the absence of enzymes and can be accelerated by alkaline pH and by basic proteins such as [histones](#) and [cytochrome c](#). Hence the assignment of AP lyase activity to a protein in the absence of an associated glycosylase activity is virtually impossible.

## 2. AP Endonucleases (Class II AP Endonucleases)

These enzymes hydrolyze the phosphodiester bond 5' to an AP site. AP endonuclease IV in *E. coli*, which has a **homologue** in yeast and humans, is the only example of a “pure” AP endonuclease—that is, an endonuclease with only hydrolysis activity 5' to an AP site and with no significant activity on any other substrate. However, the best-characterized AP endonuclease is the *E. coli* exonuclease III and its homologues in eukaryotes, including yeast, *Drosophila*, and mammals. In contrast to the pure AP nuclease endonuclease IV, the exonuclease III, as the name implies, has exonuclease activity as well.

### 2.1. Exonuclease III/Rrp1/APE

The prototype of this group of enzymes, exonuclease III, has several activities, including **ribonuclease H**, 3' to 5' exonuclease, and type II AP endonuclease. The enzyme is specific for double-stranded DNA. It is a simple polypeptide chain of 30 kDa with no cofactor. The [X-ray crystallography](#) structure of an enzyme–substrate complex demonstrates that the base is flipped out of the double helix, into the active site of the enzyme (4). The *E. coli* enzyme is a potent 3' to 5' exonuclease. The exonuclease activity has a preference for blunt double-stranded termini, but it can initiate exonuclease action from a nick as well. An activity related to the 3' to 5' exonuclease function is the ability of the enzyme to remove deoxyribose fragments, including 3'-phosphoglycolate esters and 3'-phosphate, from the 3' termini of DNA strand breaks generated by attack of reactive oxygen species on DNA. Exonuclease III-defective mutants are extremely sensitive to H<sub>2</sub>O<sub>2</sub> and ionizing radiation, which most probably kills cells through the strand breaks introduced by reactive oxygen species. As a rule, these are not clean breaks with 3'-OH and 5'-P termini. Instead, they most often contain either 3'-phosphate or fragmented deoxyribose, which must

be processed further before they can be utilized by **DNA polymerases** to fill in the gap created by the DNA damage. These structures are removed efficiently by exonuclease III, but not by endonuclease IV.

The *Drosophila* exonuclease III homologue is a monomer of 75 kDa. The 250-amino-acid-residue C-terminal region of the protein is homologous to the exonuclease III family of AP endonucleases. However, the N-terminal 400 residues are not related to any known sequence, and the function of this putative **domain** is unknown at present. The *Drosophila* enzyme was initially isolated as an activity that catalyzes strand transfer and was thought to perform a **RecA**-like function in this organism, hence the name Rrp1 (Recombination related protein 1). However, later work has not confirmed a role for this protein in [recombination](#). The enzyme has a potent type II AP endonuclease activity, but its 3' to 5' exonuclease activity is rather modest compared to the *E. coli* enzyme and is only  $10^{-2}$  to  $10^{-3}$  of the AP endonuclease activity. Whether the N-terminal 400 residues play any role in recombination or in DNA repair is not known at present.

The human AP endonuclease (APE/HAP1/APEX) is highly homologous (90% to 95% sequence identity) to the enzymes from bovine and rodent sources. These enzymes are 35-kDa monomers and are clearly related to *E. coli* exonuclease III. The mammalian enzymes have potent type II AP endonuclease and 3'-phosphoglycol aldehyde esterase activities necessary for denuding a "jagged" 3' terminus caused by direct attack of reactive oxygen species or by direct hit by ionizing radiation. In contrast to exonuclease III, however, the 3' to 5' exonuclease activity of mammalian type II AP endonucleases is rather weak and undetectable with certain 3' sequences. The structure of human AP endonuclease indicates that it flips out the deoxyribose moiety into the active site ([5](#)), as is observed also in the *E. coli* exonuclease III.

A unique property of mammalian APE proteins is their capacity to reduce critical [cysteine](#) residues in certain [transcription factors](#) and by doing so activate these transcription factors. In fact, the human APE was also purified as a factor that activates the AP-1(Jun/Fos) transcription factor and was named Ref (for reducing factor 1) before the realization that it is identical to APE. The APE cysteine residue involved in this reduction reaction is located in the NH<sub>2</sub>-terminal region and plays no role in the AP endonuclease activity of this protein. In addition to Jun and Fos, other transcription factors, including USF, NF-kB, c-Myb, and v-Rel, are activated by reduction of an oxidized cysteine by APE *in vitro* ([5](#)). However, the physiological role of these *in vitro* activities is currently not known.

Mouse APE **knockout** mutants are lethal at the embryonic stage. This could be due to either (a) the importance of base excision repair for survival of the organism or (b) disruption of an activating mechanism for important transcription factors.

## 2.2. Endonuclease IV

The *E. coli* endonuclease IV, Nfo, is a monomer of 30 kDa that cleaves 5' to AP sites in a reaction that is independent of divalent metal ions such as Mg<sup>2+</sup>. It also has potent 3'-PGA diesterase and 3'-phosphatase activities but does not have exonuclease function. Although Nfo<sup>-</sup> mutants are not especially sensitive to H<sub>2</sub>O<sub>2</sub>, they are hypersensitive to killing by the oxidative agents bleomycin and *t*-butylhydroperoxide, indicating that endonuclease IV acts on a subclass of oxidative DNA lesions that are not processed efficiently by exonuclease III. The importance of Nfo in cellular defense against oxidative stress is underscored by the fact that this is the only DNA repair enzyme that is induced by oxidative stress as part of the SoxRS regulon.

The Nfo homologue in yeast, Apn1, is the major AP endonuclease in this organism. Apn1 is a monomer of 40 kDa with 31% sequence identity with the *E. coli* Nfo endonuclease. The biochemical properties of Apn1 are very similar to those of Nfo, and indeed the expression of Apn1 in *E. coli* corrects the oxidant-sensitive phenotype of *nfo* mutants, but not the H<sub>2</sub>O<sub>2</sub> sensitivity.

## Bibliography

1. V. Bailly and W. G. Verly (1987) *Biochem. J.* **242**, 565–572.
2. M. L. Dodson, R. D. Schrock, and R. S. Lloyd (1993) *Biochemistry* **32**, 8284–8290.
3. R. P. Cunningham, H. Asahara, J. F. Bank, C. P. Scholes, J. C. Salerno, K. Surerus, E. Munck, J. McCracken, J. Peisach, and M. H. Emptage (1988) *Biochemistry* **28**, 4450–4455.
4. C. D. Mol, C. F. Kuo, M. M. Thayer, R. P. Cunningham, and J. A. Tainer (1995) *Nature* **374**, 381–386.
5. M. A. Gorman, S. Morera, D. G. Rothwell, E. de la Fortelle, C. D. Mol, J. A. Tainer, I. D. Hickson, and P. S. Freemont (1997) *EMBO J.* **16**, 6548–6558.

## Suggestions for Further Reading

6. B. Demple and L. Harrison (1994) Repair of oxidative damage to DNA: enzymology and biology. *Annu. Rev. Biochem.* **63**, 915–948.
7. P. W. Doetsch and R. P. Cunningham (1990) The enzymology of apurinic/apyrimidinic endonucleases. *Mutation Res.* **236**, 173–201.
8. E. Seeberg, L. Eide, and M. Bjoras (1995) The base excision repair pathway. *Trends Biochem. Sci.* **20**, 391–397.
9. R. J. Roberts (1995) On base flipping. *Cell* **82**, 9–12.

## AP Sites (Apurinic/Apyrimidinic Sites)

AP (apurinic/apyrimidinic) sites are deoxyribose residues in the DNA that have lost the purine or pyrimidine bases.

### 1. Occurrence

AP sites are generated by a variety of mechanisms:

#### 1.1. Spontaneous Base Loss

The glycosylic bond joining the base to the deoxyribose in the DNA backbone is relatively unstable and is cleaved, even under physiological conditions, at a rate incompatible with life. It is estimated that at normal pH and temperature there are base losses at rates of about  $10^{-12} \text{ s}^{-1}$  for double-stranded DNA and  $10^{-10} \text{ s}^{-1}$  for single-stranded DNA. This means that humans (with about  $1.2 \times 10^{10}$  bp per diploid cell) lose about 1 base per minute per cell, or, taking into account that there are about  $10^{13}$  cells in the human body, lose  $1.5 \times 10^{16}$  bases from their DNA daily. The rate of base loss is greatly accelerated by heat and low pH. Because of higher susceptibility of their N9 atom to nucleophilic attack, the purine bases are hydrolyzed at a rate about 100-fold faster than are pyrimidines (1). In fact, because of this unique property of purines, one of the procedures for chemical [DNA sequencing](#) reactions consists of heating the DNA at acidic pH, followed by cleavage of the resulting AP site by alkali treatment.

#### 1.2. Ionizing Radiation

Ionizing radiation most often causes base reduction, oxidation, or fragmentation. In particular, a [urea](#) residue attached to the deoxyribose is a relatively common product of ionizing radiation and, for all practical purposes, may be considered an AP site. Such radiation causes base loss by either (a) generating reactive oxygen species that attack and destroy bases or (b) direct hits by ionizing

radiation.

### 1.3. Glycosylases

There are about a dozen [DNA glycosylases](#) with varying degrees of specificity for bases with minor modifications: uracil glycosylase, 3-methyladenine DNA glycosylase, 8-oxoguanine DNA glycosylase, and glycosylases that act on mismatches, such as T-G glycosylase, specific for the T residue, and A-G glycosylase, specific for the A residue. These enzymes release the modified, abnormal, or mismatched bases by hydrolyzing the glycosylic bond linking the base to the deoxyribose. Most of the glycosylases have dual enzymatic activities in that, after cleaving the glycosylic bond, they also cleave the phosphodiester bond by  $\beta$ -elimination. However, some, such as uracil glycosylase, are “pure” glycosylases with no such associated AP lyase activity.

## 2. Consequences

AP sites in DNA are, in fact, very unstable. Even at physiological pH and temperature, the abasic sugar can be eliminated by  $\beta$ -d elimination. The presence of the divalent metal ions that would be expected in biological fluids greatly accelerates the cleavage reaction. The reaction may be further accelerated by basic proteins, such as [cytochrome c](#). Consequently, AP sites as such do not constitute an important lesion interfering with cell survival. In contrast to  $\beta$ -elimination, however, which occurs readily under a variety of conditions, the  $\delta$ -elimination reaction, which in combination with  $\beta$ -elimination would release the abasic deoxyribose from DNA, does not occur at physiologically relevant rates. [AP endonucleases](#) hydrolyze the phosphodiester bond 5' to the AP site, and *Escherichia coli* cells that lack this enzyme are sensitive to agents that generate AP sites or fragmented deoxyribose residues.

### Bibliography

1. T. Lindahl (1976) *Nature* **259**, 64–66.

### Suggestions for Further Reading

2. B. Demple and L. Harrison (1994) Repair of oxidative damage to DNA: enzymology and biology. *Annu. Rev. Biochem.* **63**, 915–948.
3. E. Seeberg, L. Eide, and M. Bjoras (1995) The base excision repair pathway. *Trends Biochem. Sci.* **20**, 391–397.

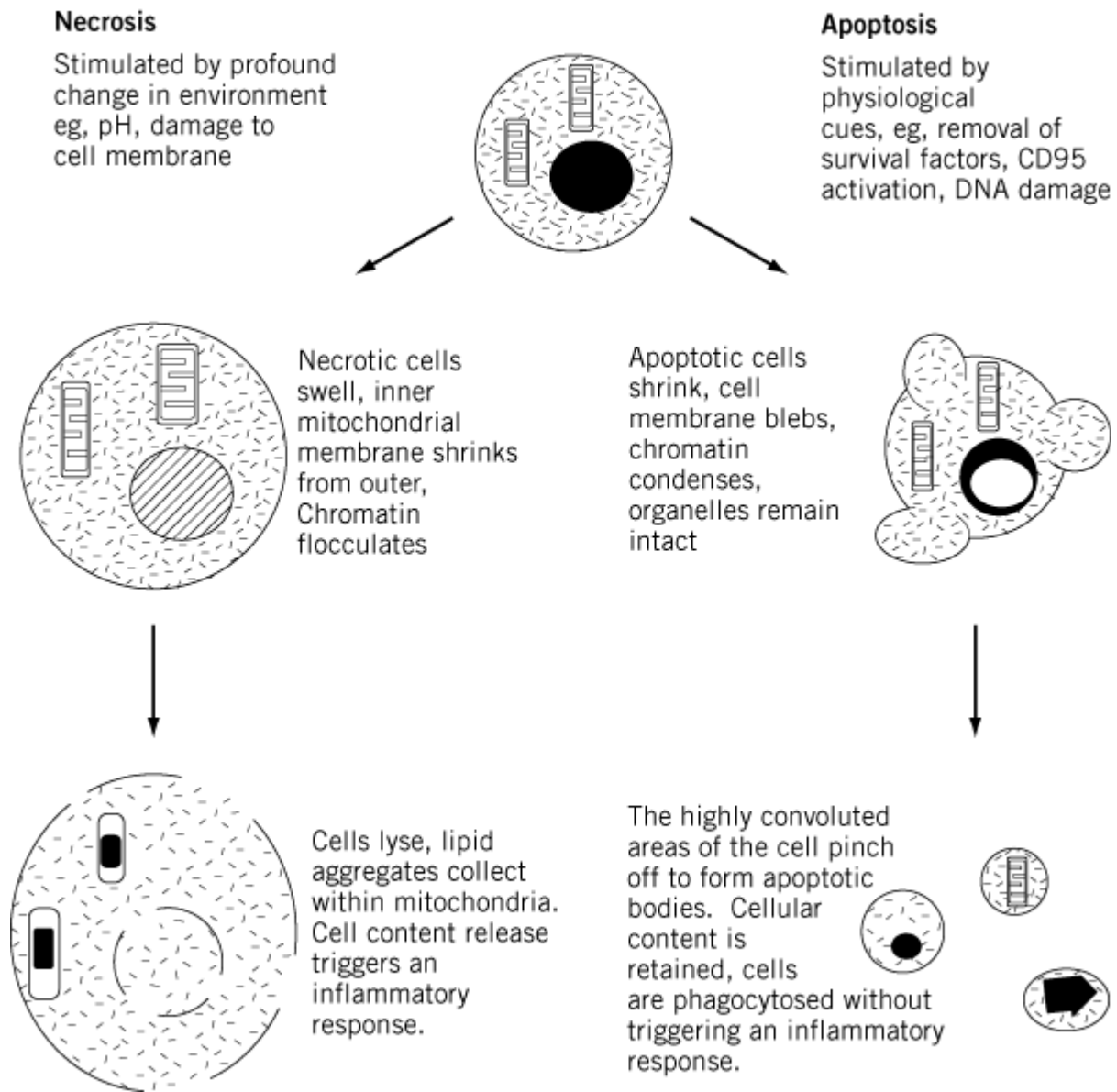
## Apoptosis

Wyllie et al. initially used the term “apoptosis” in a paper in 1972 to describe a newly observed form of [cell death](#) (1). The discovery of apoptosis has revolutionized many areas of modern biology, especially those areas concerned with disease development (2). This is primarily because apoptosis is an active form of cell death that, unlike [necrosis](#), is critical for the maintenance of tissue homeostasis (3, 4). Thus, cell populations are numerically controlled through not only [differentiation](#) and proliferation, but also the physiological loss of cells through apoptosis. Critically, this implies that a cell which loses the capacity to apoptose, or a cell which becomes inappropriately sensitive to death stimuli, can lead to perturbations in cell number regulation and, ultimately, disease development.

Cells undergoing apoptosis are evident in all tissues in response to diverse stimuli. In healthy tissues, apoptosis accounts for all the cell deaths occurring in response to normal physiological signals. Apoptosis (often referred to as [programmed cell death](#) in a developmental context) occurs at many

stages of embryonic [development](#), eg, during the formation of digits from a solid limb paddle, where interdigital cells between the digits die (5). Single apoptotic cells are evident in healthy proliferating tissues such as small gut crypts and the dermis of the skin (4). Moreover, apoptosis is also seen during endocrine-induced atrophy of tissues, in cell-mediated [immunity](#), in development of the nervous system, and in the “necrotic” oxygen-starved cores of solid tumors (6).

Morphologically, an apoptotic cell is very different from a necrotic cell (1, 6, 7); see Figure 1. One of the most prominent and easily identifiable stages of apoptosis involves the [nucleus](#). [Chromatin](#) condenses and forms aggregates near the nuclear [membrane](#), which becomes convoluted, whereas the [nucleolus](#) becomes enlarged and appears abnormally granular. The chromatin is also subject to the actions of an activated **endonuclease** that cleaves the DNA into 300- to 50-kbp fragments initially, and 180-bp fragments subsequently (8). Changes in the cytoplasm are also evident at this time; the cell shrinks visibly, adherent cells round up, and distinct protuberances or membrane blebs are discernible. **Organelles** within the shrunken cytoplasm still look normal, except for dilation of the [endoplasmic reticulum](#). These “blebbing” cells exclude vital dyes, indicating no structural failure in the cell membrane. It is at this point in the apoptotic process, which may only take 10 to 30 minutes, that apoptotic cells *in vivo* are **phagocytosed**, by either their nearest neighbors or professional [macrophages](#) (9). Thus, apoptosis occurring in single cells *in vivo* is very easily overlooked due to its rapid progression. Cells in the later stages of apoptosis, especially those *in vitro*, form apoptotic bodies as a result of pinching off of the highly convoluted blebbing areas of the cell. These apoptotic bodies are phagocytosed or, in situations where there is much cell death, the unphagocytosed apoptotic cells undergo secondary necrosis, characteristically swelling and losing membrane integrity.



## 1. Removal of Apoptotic Cells

A critical part of the apoptotic pathway is the efficient recognition and phagocytosis of apoptotic cells. The rapid disposal of apoptotic cells does not elicit an [immune response](#), consistent with apoptosis being the “no-nonsense” pathway for the disposal of unwanted cells. A breakdown in either the communication between macrophage and apoptotic cell or in the pathway of apoptosis itself may be two of the mechanisms through which chronic inflammatory disorders occur.

Recognition of apoptotic cells by professional macrophages is mediated partly by  $\alpha_v\beta_3$  [integrin](#) (the *vitronectin* receptor) and the CD36 ligand. Between these two molecules on the external surface of the cell sits a **glycoprotein** containing the sequence -R-G-D- (-Arg-Gly-Asp-), *thrombospondin*, which acts as a molecular bridge to bind to an unknown anionic site on the apoptotic cell ([10-12](#)). Apoptotic cells also express surface phosphatidylserine, which binds to an as yet uncharacterized macrophage **receptor** ([13, 14](#)). The recognition of apoptotic cells mediated by the vitronectin receptor, CD36 ligand, and thrombospondin is confined to monocyte-derived macrophages. Recognition of apoptotic cells through the expression of phosphatidylserine, which is flipped onto the cellular surface due to the inactivation of an ATP-dependent flippase, is utilized primarily by inflammatory macrophages ([9, 15](#)). Four other molecules have been implicated in apoptotic cell

recognition systems: (1) the 61D3-antigen found on human macrophages, (2) ICAM-3, found on the surface of human [B cells](#) (16), (3) the ABC1 transporter (17), and (4) the macrophage scavenger receptor (18). As mentioned previously, the interaction between apoptotic cells and macrophages does not elicit an inflammatory response, unless the apoptotic cells start to undergo secondary necrosis. Under these circumstances, other macrophage receptors are employed, and inflammatory **cytokines** are released (19). This change from “silent” death to one that activates the immune system may be a method of recruiting more phagocytes to the site of cell death to cope with the increasing number of corpses.

## 2. Regulation of Apoptosis and Cell Viability

The discovery that apoptosis is an active form of cell death which is transiently suppressable by both inhibitors of protein and RNA synthesis suggested that the cell needed to synthesize new proteins prior to its death [reviewed in (20)]. This sparked an enthusiastic hunt for a critical cell death gene that had to be translated prior to death. No such gene has been found, but many genes that regulate apoptosis have been identified as a result.

A cell can be triggered to undergo apoptosis in response to diverse stimuli, with each stimulus activating different pathways involving **gene expression**, as well as [post-translational modification](#). The discovery that apoptosis can be suppressed by the presence of specific survival factors, either cytokines or proteins expressed within the cell, suggested that regulation of cell viability was paramount to the function of the cell (21). Indeed, it appears that apoptosis represents a default pathway in all cells (22). If a cell does not receive the correct survival stimulus, it will die; consequently, cells within multicellular organisms are maintained in a viable state by a constant supply of survival factors.

Once a cell is triggered to apoptose, activation of a common set of destruction proteins, irrespective of the stimulus involved, precipitates the morphological changes that we call apoptosis. Overall, the apoptotic pathway can be divided into three phases:

1. The *decision phase*, in which the cell receives a stimulus and, depending on both its internal and external environments, may or may not be triggered to die.
2. The *commitment phase*, in which the cell is committed to death and cannot recover.
3. The *execution phase*, in which the decision to die has been made and the activated apoptotic machinery leads to the morphological changes that define an apoptotic cell.

## 3. The Decision Phase of Apoptosis

The decision phase can be compared to that of a judge hearing the evidence during a trial. The evidence presented comes from many different sources: genes, cytokines, toxic chemicals, DNA damage, and the presence of viral and bacterial infection. Only once all the evidence is heard will the sentence be passed. The genes that modulate the decision phase are many and varied. Some are more commonly associated with the regulation of the cell cycle, such as *c-myc*, *cdc25*, *c-fos*, [p53](#) and *rb* (**retinoblastoma**), whereas others are members of recently identified gene families, such as the *bcl-2* family. **Cytokines** can either trigger or delay apoptosis. [Tumor necrosis factor](#)  $\alpha$  (TNF $\alpha$ ) and CD95 ligand (CD95L, FasL, Apo-1L) are two examples of cytokines that trigger apoptosis upon binding to their receptors (23), whereas the addition of factors such as **interleukin-3** (IL-3) (24) or *insulin-like growth factor 1* (IGF-1) suppresses death (25). Conversely, removal of anti-apoptotic cytokines also induces cell death (26, 27). Both viruses and bacteria trigger apoptosis upon infection of a cell, suggesting that the cell commits suicide to prevent the invading pathogen from spreading. However, both viruses and bacteria have evolved mechanisms that manipulate this response. [Adenovirus](#), eg, has three genes required for productive infection: E1A, E1Bp19, and p55. E1A stimulates the infected cell to proliferate; however, this stimulus concomitantly triggers apoptosis. E1Bp19 and p55 delay the host cell from triggering cell death, thereby facilitating virus production (28, 29).



Alternatively, some bacteria such as *Shigella flexneri* actively promote the death of the cell by triggering the execution machinery (30-32). This allows the early release of the bacteria and triggers an immune response, which damages the surrounding tissue and aids the passage of the bacteria into its target cells.

### 3.1. Oncogenes, Tumor Suppressor Genes, and Apoptosis

In terms of suppressing tumorigenesis, long-lived organisms like humans need to limit the chances of their component cells acquiring mutations that lead to increased proliferative capacity and eventual clonogenic outgrowth. This is, in part, limited by restricting the types of cells with proliferative capacity. Many adult cell lineages exist in a senescent or postreplicative state. Moreover, for those cells that must retain the capacity to divide, apoptosis provides a “safety net.” Genes that are involved in proliferation are also involved in cell death; in fact, the two processes appear inexorably linked. This paradoxical coupling of two opposing biological processes was first realized through work on the pervasive oncogene *c-myc* (33) but has since been confirmed for several mitotic genes (34).

### 3.2. *c-myc* and Apoptosis

The protein product of the *c-myc* gene, c-Myc, is both necessary and sufficient to ensure fibroblast proliferation; in the presence of **mitogen**, all cells in cycle express c-Myc. Fibroblasts remain in cycle as long as mitogens and c-Myc are present, but only remain viable in the presence of additional survival factors. In the absence of such survival factors, cycling cells expressing deregulated Myc undergo apoptosis. This response to oncogene deregulation is thought to protect the organism against cells that acquire mutations in c-Myc. Mutant cells that continuously express c-Myc will proliferate, but since the survival factor *in vivo* is thought to be limiting, the mutant cells will die once the supply of factor is exhausted (25, 35). Thus, small proto-tumors may well arise *in vivo*, but they should regress due to lack of survival signals. If, however, another mutation has occurred producing a cooperating anti-apoptotic lesion, the mutated cells will survive. Several genes that cooperate with *myc* during tumorigenesis are anti-apoptotic. For example, v-Abl, which is a very good suppressor of apoptosis in response to many stimuli, will render haemopoietic cells cytokine-independent for growth and survival (36, 37). *bcl-2*, a gene initially identified through its translocation in the t(14:18) mutation found in follicular lymphoma (38), will cooperate with *myc*, suppressing the apoptotic signal induced through Myc expression (39, 40). The outcome of this cooperation *in vivo* is lymphoma, as seen in both double *bcl-2/Eμ-myc* **transgenic** mice (41, 42) and in human patients with the t(14:18) translocation (43, 44).

### 3.3. Bcl-2, an Anti-apoptotic Oncogene

The potent ability of Bcl-2 to suppress apoptosis triggered in response to a number of diverse stimuli prompted much research into its biochemical function. This function is still relatively obscure, but several members of the Bcl-2 family of proteins have been identified as a result (45-47). Bcl-2's form and functionality have been conserved throughout multicellular evolution, as exemplified by Ced-9, its counterpart in the nematode *C. elegans*, and the Adenovirus E1Bp19 protein (48). Ten or so members of this family now exist, and all share defined regions of **homology**, although not all act to suppress apoptosis (Table 1). The pro-apoptotic members of the family (Bax, Bak, and Bcl-x<sub>S</sub> (49-53)) share three regions of homology, BH-1, -2, -3, whereas the anti-apoptotic members (Bcl-2 and Bcl-x<sub>L</sub>) have a fourth region of homology, known as BH-4 (54, 55). Removal of the BH-4 **domain** from either Bcl-2 or Bcl-x<sub>L</sub> results in a loss of protective function. Family members can dimerize with themselves and with one another, and these interactions are important for their function (56, 57). For example, Bcl-2 preferentially binds to the pro-apoptotic family member Bax, whereas Bcl-x<sub>L</sub> binds to Bak. If the number of Bax homodimers in the cell exceeds the number of Bcl-2 homodimers or Bax/Bcl-2 heterodimers, then the cell is more likely to undergo apoptosis (47). However, it is not clear which forms of these proteins are dominant for determining life or death. A second set of Bcl-2-associated proteins has been identified that influences cell survival through interaction with Bcl-2 family members, but whose members do not have all, if any, of the BH regions (Table 2) (45, 46). Bag-1, eg, has no homology to Bcl-2, but binds to it and augments the

protective function of Bcl-2 (58). Bad, on the other hand, has homology to the family and binds Bcl-2, an interaction that disrupts Bcl-2's binding with Bax, leading to cell death (59, 60).

**Table 1. Pro-apoptotic and Anti-apoptotic Members of the Bcl-2 Family<sup>a</sup>**

| Bcl-2 Family Inhibitors | Bcl-2 Family Promoters |
|-------------------------|------------------------|
| Bcl-2 (mammalian)       | Bax (mammalian)        |
| Bcl-xL (mammalian)      | Bak (mammalian)        |
| Bcl-w (mammalian)       | Bcl-xS (mammalian)     |
| Mcl-1 (mammalian)       |                        |
| A1 (mammalian)          |                        |
| NR-13 (mammalian)       |                        |
| E1B p19 (viral)         |                        |
| BHRF1 (viral)           |                        |

<sup>a</sup> Anti-apoptotic members contain up to 4 Bcl-2 homology (BH) domains and the transmembrane (TM) domain, whereas pro-apoptotic members contain only BH 1-3 and the TM domain, except for Bcl-x<sub>S</sub> that contains only BH4, BH3, and the TM domain. Members of this gene family have been conserved throughout evolution with homologues existing in nematodes, mammals, and viruses (46, 47, and references therein].

**Table 2. Proteins That Bind to Bcl-2 Family Members and Regulate Their Function<sup>a</sup>**

| Bcl-2-like Inhibitors of Apoptosis | Bcl-2-like Promoters of Apoptosis     |
|------------------------------------|---------------------------------------|
| Bag-1 (no BH domains)              | Hrk (BH 3, TM domain)                 |
| Bra (BH 2 and BH 3)                | Bik (BH 3, TM domain)                 |
| Raf-1 (no BH domains)              | Bid (BH 3)                            |
|                                    | Bim (BH 3, TM domain)                 |
|                                    | Bad (BH 1, 2, 3, no hydrophobic tail) |

<sup>a</sup> A variety of proteins interact with Bcl-2 proteins. Some contain the BH domains that influence their function, such as Bid, Bim, Bik, whereas others have no homology with Bcl-2, but still bind and influence its function (46, 47, and references therein].

Both Bcl-x<sub>L</sub> and Bcl-2 have effects on cell-cycle transition (61). Cells expressing either gene exhibit a longer-than-average time to reenter the cell cycle after arresting in G<sub>0</sub>, suggesting that Bcl-x<sub>L</sub> and Bcl-2 interact with components of the cell-cycle machinery.

### 3.4. Tumor Suppressor Genes and Apoptosis

[Tumor suppressor gene](#) products suppress unrestrained cell proliferation through their specific

inhibitory effects on the cell cycle and are implicated in the development of neoplasia following their loss or functional inactivation. One tumor suppressor gene product, [p53](#), induces apoptosis in some tumor cell lines. p53 is a short-lived protein that is stabilized in the presence of DNA damage and triggers cell-cycle arrest, presumably to facilitate [DNA repair](#) (62). If the DNA damage is too great to repair, apoptosis is triggered, earning p53 the title “guardian of the genome” (63). At present, it is unclear if p53's induction of apoptosis is effected through its upregulation of p21<sup>Waf-1/Cip-1</sup>, leading to cell-cycle arrest (64), through some other p53-regulated gene, or is independent of p53 transcriptional activity (65). Bax, the pro-apoptotic Bcl-2 homologue, is transcriptionally regulated by p53, but not all apoptotic cell deaths induced by p53 stabilization require Bax expression (66). The role of p53 as a sensor of DNA damage was neatly illustrated by examining apoptosis in thymocytes from p53-knockout mice (67, 68). p53 null thymocytes undergo death by apoptosis as normal, except when treated with agents that damage DNA, such as etoposide or irradiation. Therefore, in the absence of p53, the cell is neither instructed to leave the cell cycle and repair the damaged DNA, nor to die. Cells therefore progress through subsequent cycles with damaged and mutated DNA. The role of p53 in inducing apoptosis in circumstances in which DNA damage is not evident is at present unclear, although in p53 knockout mice there is no evidence of a defect in tissue homeostasis due to ineffective cell death when cytokines become limiting (67, 68).

A second tumor suppressor gene implicated in control of apoptosis is the **retinoblastoma** (*rb*) gene (69-71). The absence of the functional gene product, Rb, results in massive apoptosis in the haemopoietic and nervous system possibly due to the cells being constantly in cycle and unable to arrest or terminally differentiate. This makes it impossible for cells to establish appropriate survival signals and cell:cell contacts. Loss of both p53 and Rb function results in rapid tumor progression (72), suggesting that cells which have overcome restraints upon cell-cycle progression can also escape apoptosis, underlining the link between cell proliferation and cell death.

### 3.5. Cytokine Signaling Pathways and Apoptosis

Although a wealth of information is available about the [signal transduction](#) pathways activated by cytokines and [growth factors](#) in many cell types, surprisingly little is known about how apoptosis (and conversely cell viability) is modulated. Initially, many of the factors now referred to as *survival factors* were considered solely to stimulate cells to proliferate. Thus, many of the identified signal transduction pathways are involved in proliferation but not necessarily survival. For survival signals, two pathways can exist. Cells can be stimulated to survive by one specific pathway. Loss of this signal upon removal of the cytokine may simply trigger cell death due to negation of the first signal. However, removal of the cytokine may trigger another independent signaling pathway. Only now are these different possibilities being investigated.

#### 3.5.1. Survival Signals

Cytokines with very different actions activate substantially overlapping signaling circuits, many of which appear to be required, but are not solely responsible, for regulating cell survival. However, even with these complex signaling networks, it is possible to gather some information about survival signals. For example, apoptosis stimulated by the withdrawal of **nerve growth factor** (NGF) in phaeochromocytoma (PC12) cells involves the induction of the AP-1 [transcription factor](#) c-Jun proximal to the time when neurons become committed to apoptosis. Activation of c-Jun requires phosphorylation by the protein kinase JNK (MKK4) that is, in turn, activated by MAP kinase (73) (see [Phosphorylation, Protein](#)). In contrast, the suppression of apoptosis by NGF in the same cells may involve a discrete survival signaling pathway routed through **Ras** (74), Raf, and MAP kinase (73). Hence, in this case, the induction of apoptosis and promotion of survival may be independent informational processes, rather than mere negation of one other.

In some documented cases, the ability of anti-apoptotic cytokines to suppress apoptosis does not depend on the synthesis of new genes or proteins, indicating that survival is mediated by posttranslational mechanisms. IGF-1 (25), a potent inhibitor of apoptosis in many cell types, and [epidermal growth factor](#) (EGF) (75) are reported to suppress apoptosis effectively in cells treated with inhibitors of RNA and protein synthesis. The suppression of apoptosis in such instances must be

effected by means of preexisting molecular machinery. One possible candidate signal is the activation of PI-3 kinase (PI3-K) through Ras, whose inhibition blocks the ability of NGF to mitigate apoptosis (76). However, specificity is absent, since PI3-K is activated in many signaling pathways in response to signals that do not suppress apoptosis. IGF-1 also signals via PI3-K, and the pro- and anti-apoptotic signaling pathways for this cytokine have been identified. PI3-K can be activated by its upstream effector Ras, and it can activate several downstream pathways, including the ribosomal protein p70<sup>S6K</sup>, the **Rho** family polypeptide Rac, and the serine/threonine protein kinase PKB/Akt. Of these downstream effectors, PKB/Akt is anti-apoptotic in fibroblasts expressing c-Myc and in other cell types in the absence of survival factors, whereas the Ras-mediated Raf signaling pathway is pro-apoptotic (77-79). These findings underscore the pleiotropic nature of intracellular signaling. Signals emanating from GTP-Ras trigger a plethora of potential biological outcomes, some of which promote apoptosis and some of which suppress it. The net outcome for Ras activation is presumably dictated by downstream interactions that potentiate or mitigate other signals.

### 3.5.2. Killer Cytokines

Apoptosis can also be triggered by certain **cytokines**, such as TNF $\alpha$  and CD95 ligand. These two pathways represent the best understood triggers of apoptosis, since their activation pathways are essentially mapped. CD95 ligand and TNF $\alpha$  bind and activate their specific receptors on the surface of the cell (80). The TNF type 1 and CD95 (Fas, Apo-1) receptors are functionally similar; both are transmembrane receptors that contain a region within their cytoplasmic tail which is necessary for the initiation of the apoptotic response (81-83). This region of homology is called the *death domain* (DD). Several adaptor proteins bind to this domain and can trigger apoptosis by directly linking to the downstream effector **caspase** machinery (23). Alternatively, in the case of TNFR1, the adaptor proteins can also activate NF $\kappa$ B signaling pathways (84). CD95-induced cell death, eg, results in the association of the *death-induced signaling c omplex* (DISC) (80). Upon binding of its ligand, the CD95 receptor trimerizes and binds to its death domain the protein *Fas-associated death domain* (FADD), also known as MORT1 (85, 86). FADD contains a C-terminal DD and, at its N-terminus, a *death effector domain* (DED), so called because the expression of this domain is required to trigger CD-95-induced apoptosis. The DED domain of FADD binds to the DED in caspase-8, a protein also called FLICE (FADD-like ICE, where ICE is “interleukin converting enzyme”) (87, 88). The **thiol proteinase** domain of FLICE, located at its C-terminus, is activated upon binding FADD, and FLICE then activates the *caspase cascade* by cleaving downstream caspases (89) (see **Caspases**). CD95-triggered apoptosis can be suppressed by expression of the anti-apoptotic proteins Bcl-2 and Bcl-x<sub>L</sub> (90) or by proteins that inhibit the interaction of FLICE with FADD (91). Thus, the binding of CD95 ligand to the receptor and recruitment of FADD occur within the decision phase of apoptosis. The commitment of cells to CD95-induced apoptosis does not appear to occur until FLICE is recruited to the DISC, cleaved, and therefore activated. TNFR1 exhibits a similar pathway involving an overlapping set of adaptor proteins (23).

## 4. The Commitment Phase of Apoptosis

As yet, it is unclear what commits a cell to undergo apoptosis. To a certain degree, the final trigger must be cell autonomous, because cells that have received the same stimulus do not all enter apoptosis at the same time; cells may die within 10 minutes of the signal or within 3 days. Initially, the **caspases** were proposed to act at this point, with their activation leading to certain death (92-94). Although exogenous expression of pro-caspases leads to cell death, this death is again stochastic and can be suppressed by the addition of IGF-1 (95) or thiol proteinase inhibitors such as the poxvirus protein CrmA and **baculovirus** inhibitor p35 (96, 97). Surprisingly though, not all caspase inhibitors suppress cell death. In the case of c-Myc-induced or Bax-induced apoptosis, the addition of tetrapeptide caspase inhibitors based on the pro-IL-1 $\beta$  cleavage site does not prevent cells from becoming committed to die (98, 99). Thus, even when caspases are inhibited, the cell still undergoes commitment, but it takes much longer to progress through the morphological changes of apoptosis (99). This result suggests caspases are required for the rapid execution phase of apoptosis, but not necessarily for commitment.

The unknown entity in mammalian cells is the putative homologue of the *cell death gene 4* (Ced-4) (see [Programmed Cell Death](#)). Ced-4 is required for cell death in *C. Elegans* and appears to function upstream of Ced-3, suggesting that Ced-4 may be involved in commitment (100). Recent data have shown that Ced-9, Ced-3, and Ced-4 bind one another, with Ced-4 being the critical intermediate (101-104). Without Ced-4 expression, Ced-9 does not interact with Ced-3. The expression of Ced-4 in mammalian cells produces a similar reaction between Bcl-x<sub>L</sub> and the Ced-3 homologue, FLICE (caspase-8) (101). Whether Ced-4 is a killer of mammalian cells is also not clear, but Ced-4 expression is toxic to yeast (104). In the fission yeast *Schizosaccharomyces pombe*, Ced-4 localizes to condensed [chromatin](#) and induces caspase-independent death in this apoptotically naïve cell model. Cell death is suppressed upon the expression of Ced-9, which causes relocation of Ced-4 to the membrane of the nucleus and to the [endoplasmic reticulum](#). These data strongly suggest that Ced-9 suppresses Ced-4-induced death by binding to it. This, in turn, implies that the disruption of the Bcl-x<sub>L</sub>/Ced-4-“like protein” and FLICE interaction in mammalian cells may commit the cell to death.

## 5. The Execution Phase of Apoptosis

The execution phase describes the stage where the cells morphologically change and exhibit all the characteristics of apoptosis. Not all the agents responsible for these changes have been identified, but activation of the endonuclease occurs at this stage, as does activation of the *ced-3* homologues, the [caspases](#). Overall, the caspase pathway can be thought of as a packaging process that carries out the rapid dismantling of the cell. Cells still die in the absence of caspase activity, but death is significantly prolonged (99). The cleavage of caspase substrates such as Poly-(ADP ribose) polymerase (PARP), the catalytic subunit of DNA protein kinase, and nuclear lamins [reviewed in (105)], in concert with an endonuclease, allow the nucleus and chromatin to be degraded efficiently. This saves the dying cell energy and facilitates the maintenance of membrane integrity, which in turn prevents release of the cell contents, which would stimulate an immune response.

Other effectors of the execution phase are not so well documented. Both ceramide and reactive oxygen species are found in apoptotic cells, but their exact roles are not clear (110-114): They may well be a consequence of death, rather than an effector. Changes in [mitochondria](#) also occur within apoptotic cells (115-117). Selective release of **cytochrome c** from the mitochondria appears to trigger apoptosis through activation of the caspases. A membrane permeability transition also occurs where the pores that regulate traffic between the mitochondrial outer and inner membranes become disrupted (117). However, it is unclear exactly when this change occurs within an individual cell undergoing apoptosis, making it difficult to assign these changes to a particular phase within the apoptotic pathway.

Overall, the execution phase is one where caspases and endonucleases are active and the cell is dismantled, progressing through the series of morphological changes that originally defined the process of apoptosis. Execution ends with the efficient and rapid phagocytosis of the corpse.

## 6. Summary

Apoptosis is an active form of cell death indistinguishable from programmed cell death that occurs during invertebrate and vertebrate development. Apoptosis is thought to be the default state of all cells; hence, cells must receive a survival signal in order to remain viable. Apoptosis is triggered in response to a diverse set of stimuli and plays a critical role in the maintenance of tissue homeostasis. Mutations that disrupt the apoptotic pathway are pivotal in the development of many diseases, including degenerative brain disorders and cancer.

## Bibliography

1. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1972) Cellular events in the adrenal cortex

- following ACTH deprivation. *J. Pathol.* **106**, 1.
2. C. B. Thompson (1995) Apoptosis in the pathogenesis and treatment of disease. *Science* **267**, 1456–1462.
  3. G. Evan, E. Harrington, A. Fanidi, H. Land, B. Amati, and M. Bennett (1994) Integrated control of cell proliferation and cell death by the *c-myc* oncogene. *Phil. Trans. R Soc. Lond. B* **345**, 269–275.
  4. M. D. Jacobson, M. Weil, and M. C. Raff (1997) Programmed cell death in animal development. *Cell* **88**, 347–354.
  5. J. R. Hinchliffe and D. A. Ede (1973) Cell death and the development of limb form and skeletal pattern in normal and wingless (*ws*) chick embryos. *J. Embryol. Exp. Morphol.* **30**, 753–772.
  6. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1980) Cell death: The significance of apoptosis. *Int. Rev. Cytol.* **68**, 251–306.
  7. J. F. Kerr, A. H. Wyllie, and A. R. Currie (1972) Apoptosis: A basic biological phenomenon with wide-ranging implications in tissue kinetics. *Br. J. Cancer* **26**, 239–257.
  8. A. H. Wyllie (1980) Glucocorticoid-induced thymocyte apoptosis is associated with endogenous endonuclease activation. *Nature* **284**, 555–556.
  9. J. Savill, V. Fadok, P. Henson, and C. Haslett (1993) Phagocyte recognition of cells undergoing apoptosis. *Immunol. Today* **14**, 131–136.
  10. J. Savill, I. Dransfield, N. Hogg, and C. Haslett (1990) Vitronectin receptor-mediated phagocytosis of cells undergoing apoptosis. *Nature* **343**, 170–173.
  11. J. Savill, N. Hogg, Y. Ren, and C. Haslett (1992) Thrombospondin cooperates with CD36 and the vitronectin receptor in macrophage recognition of neutrophils undergoing apoptosis. *J. Clin. Invest.* **90**, 1513–1522.
  12. Y. Ren, R. L. Silverstein, J. Allen, and J. Savill (1995) CD36 gene transfer confers capacity for phagocytosis of cells undergoing apoptosis. *J. Exp. Med.* **181**, 1857–1862.
  13. V. A. Fadok, D. R. Voelker, P. A. Campbell, D. L. Bratton, J. J. Cohen, P. W. Noble, D. W. Riches, and P. M. Henson (1993) The ability to recognize phosphatidylserine on apoptotic cells is an inducible function in murine bone marrow-derived macrophages. *Chest* **103**, 102s.
  14. V. A. Fadok, D. R. Voelker, P. A. Campbell, J. J. Cohen, D. L. Bratton, and P. M. Henson (1992) Exposure of phosphatidylserine on the surface of apoptotic lymphocytes triggers specific recognition and removal by macrophages. *J. Immunol.* **148**, 2207–2216.
  15. V. A. Fadok et al. (1992) Different populations of macrophages use either the vitronectin receptor or the phosphatidylserine receptor to recognize and remove apoptotic cells. *J. Immunol.* **149**, 4029–4035.
  16. P. K. Flora and C. D. Gregory (1994) Recognition of apoptotic cells by human macrophages: Inhibition by a monocyte/macrophage-specific monoclonal antibody. *Eur. J. Immunol.* **24**, 2625–2632.
  17. M. F. Luciani and G. Chimini (1996) The ATP binding cassette transporter ABC1 is required for the engulfment of corpses generated by apoptotic cell death. *Embo. J.* **15**, 226–235.
  18. N. Platt and S. Gordon (1995) Role of the murine scavenger receptor in the recognition of apoptotic thymocytes by macrophages. *J. Cell Biochem. Suppl.* **19B**, 300.
  19. M. Stern, J. Savill, and C. Haslett (1996) Human monocyte-derived macrophage phagocytosis of senescent eosinophils undergoing apoptosis. Mediation by alpha v beta 3/CD36/thrombospondin recognition mechanism and lack of phagocytic response. *Am. J. Pathol.* **149**, 911–921.
  20. N. J. McCarthy, C. A. Smith, and G. T. Williams (1992) Apoptosis in the development of the immune system: Growth factors, clonal selection and *bcl-2*. *Canc. Metas. Rev.* **11**, 157–178.
  21. G. I. Evan, L. Brown, M. Whyte and E. Harrington (1995) Apoptosis and the cell-cycle. *Curr. Opin. Cell Biol.* **7**, 825–834.

22. M. Raff, B. Barres, J. Burne, H. Coles, Y. Ishizaki, and M. Jacobson (1993) Programmed cell death and the control of cell survival: Lessons from the nervous system. *Science* **262**, 695–700.
23. S. Nagata (1997) Apoptosis by death factor. *Cell* **88**, 355–365.
24. M. K. Collins, J. Marvel, P. Malde, and A. Lopez-Rivas (1992) Interleukin 3 protects murine bone marrow cells from apoptosis induced by DNA damaging agents. *J. Exp. Med.* **176**, 1043–1051.
25. E. Harrington, A. Fanidi, M. Bennett, and G. Evan (1994) Modulation of Myc-induced apoptosis by specific cytokines. *Embo. J.* **13**, 3286–3295.
26. G. T. Williams, C. A. Smith, E. Spooncer, T. M. Dexter, and D. R. Taylor (1990) Haemopoietic colony stimulating factors promote cell survival by suppressing apoptosis. *Nature* **343**, 76–79.
27. R. C. Duke and J. J. Cohen (1986) IL-2 addiction: Withdrawal of growth factor activates a suicide program in dependent T cells. *Lymphokine Res.* **5**, 289–299.
28. E. White, P. Sabbatini, M. Debbas, W. Wold, D. Kusher, and L. Gooding (1992) The 19-kilodalton Adenovirus E1B transforming protein inhibits programmed cell death and prevents cytolysis by tumour necrosis factor. *Mol. Cell. Biol.* **12**, 2570–2580.
29. E. White, R. Cipriani, P. Sabbatini, and A. Denton (1991) Adenovirus E1B 19-kilodalton protein overcomes the cytotoxicity of E1A proteins. *J. Virol.* **65**, 2968–2678.
30. A. Zychlinsky, B. Kenny, R. Menard, M. C. Prevost, I. B. Holland, and P. J. Sansonetti (1994) IpaB mediates macrophage apoptosis induced by *Shigella flexneri*. *Mol. Microbiol.* **11**, 619–627.
31. K. Thirumalai, K. S. Kim, and A. Zychlinsky (1997) IpaB, a *Shigella flexneri* invasin, colocalizes with interleukin-1 beta-converting enzyme in the cytoplasm of macrophages. *Infect. Immun.* **65**, 787–793.
32. A. Guichon and A. Zychlinsky (1996) Apoptosis as a trigger of inflammation in *Shigella*-induced cell death. *Biochem. Soc. Trans.* **24**, 1051–1054.
33. G. Evan et al. (1992) Induction of apoptosis in fibroblasts by *c-myc* protein. *Cell* **63**, 119–125.
34. M. Zonig and G. I. Evan (1996) Cell cycle: On target with Myc. *Curr. Biol.* **6**, 1553–1556.
35. E. Harrington, A. Fanidi, and G. Evan (1994) Oncogenes and cell death. *Curr. Opin. Genet. Dev.* **4**, 120–129.
36. C. A. Evans, L. P. Owen, A. D. Whetton, and C. Dive (1993) Activation of the Abelson tyrosine kinase activity is associated with suppression of apoptosis in hemopoietic cells. *Cancer Res.* **53**, 1735–1738.
37. J. L. Cleveland, M. Dean, N. Rosenberg, J. Y. Wang, and U. R. Rapp (1989) Tyrosine kinase oncogenes abrogate interleukin-3 dependence of murine myeloid cells through signaling pathways involving c-myc: Conditional regulation of c-myc transcription by temperature-sensitive v-abl. *Mol. Cell. Biol.* **9**, 5685–5695.
38. Y. Tsujimoto, L. R. Finger, J. Yunis, P. C. Nowell, and C. M. Croce (1984) Cloning of the chromosome breakpoint of neoplastic B cells with the t(14;18) chromosome translocation. *Science* **226**, 1097–1099.
39. A. Fanidi, E. Harrington, and G. Evan (1992) Cooperative interaction between *c-myc* and *bcl-2* proto-oncogenes. *Nature* **359**, 554–556.
40. R. Bissonnette, F. Echeverri, A. Mahboubi, and D. Green (1992) Apoptotic cell death induced by *c-myc* is inhibited by *bcl-2*. *Nature* **359**, 552–554.
41. A. Strasser, A. W. Harris, M. L. Bath, and S. Cory (1990) Novel primitive lymphoid tumours induced in transgenic mice by cooperation between *myc* and *bcl-2*. *Nature* **348**, 331–333.
42. T. J. McDonnell and S. J. Korsmeyer (1991) Progression from lymphoid hyperplasia to high-grade malignant lymphoma in mice transgenic for the t(14;18). *Nature* **349**, 254–256.
43. S. J. Korsmeyer, T. J. McDonnell, G. Nunez, D. Hockenbery, and R. Young (1990) Bcl-2: B cell life, death and neoplasia. *Curr. Top. Microbiol. Immunol.* **166**, 203–207.

44. Y. Tsujimoto, J. Cossman, E. Jaffe, and C. M. Croce (1985) Involvement of the *bcl-2* gene in human follicular lymphoma. *Science* **228**, 1440–1443.
45. J. C. Reed (1997) Double identity for proteins of the Bcl-2 family. *Nature* **387**, 773–776.
46. M. D. Jacobson (1997) Apoptosis: Bcl-2-related proteins get connected. *Curr. Biol.* **7**, R227–R281.
47. E. Yang and S. J. Korsmeyer (1996) Molecular thanatopsis: A discourse on the BCL2 family and cell death. *Blood* **88**, 386–401.
48. G. Williams and C. Smith (1993) Molecular regulation of apoptosis—genetic-controls on cell-death. *Cell* **74**, 777–779.
49. Z. Oltvai, C. Milliman, and S. Korsmeyer (1993) Bcl-2 heterodimerizes in vivo with a conserved homolog, Bax, that accelerates programmed cell death. *Cell* **74**, 609–619.
50. L. Boise et al. (1993) *bcl-x*, a *bcl-2*-related gene that functions as a dominant regulator of apoptotic cell death. *Cell* **74**, 597–608.
51. T. Chittenden, E. Harrington, R. O'Connor, G. Evan, and B. Guild (1995) Induction of apoptosis by the Bcl-2 homologue Bak. *Nature* **374**, 733–736.
52. M. Kiefer, M. Brauer, V. C. Powers, J. Wu, S. Umansky, L. Tomei, and P. Barr (1995) Modulation of apoptosis by the widely distributed Bcl-2 homologue Bak. *Nature* **374**, 736–739.
53. S. Farrow et al. (1995) Cloning of a novel *bcl-2* homologue by interaction with adenovirus E1B 19K. *Nature* **374**, 731–733.
54. B. S. Chang, A. J. Minn, S. E. Muchmore, S. W. Fesik, and C. B. Thompson (1997) Identification of a novel regulatory domain in Bcl-xL and Bcl-2. *Embo J.* **16**, 968–977.
55. S. W. Muchmore, M. Sattler, H. Liang, R. P. Meadows, J. E. Harlan, H. S. Yoon, D. Nettlesheim, B. S. Chang, C. B. Thompson, S. L. Wong, S. L. Ng, and S. W. Fesik (1996) X-ray and NMR structure of human Bcl-xL, an inhibitor of programmed cell death. *Nature* **381**, 335–341.
56. X.-M. Yin, Z. Oltvai, and S. Korsmeyer (1994) Bh1 and bh2 domains of Bcl-2 are required for inhibition of apoptosis and heterodimerization with Bax. *Nature* **369**, 321–323.
57. T. Sato, S. Irie, S. Krajewski, and J. C. Reed (1994) Cloning and sequencing of a cDNA encoding the rat Bcl-2 protein. *Gene* **140**, 291–292.
58. S. Takayama et al. (1995) Cloning and functional-analysis of *bag-1*—a novel *bcl-2*-binding. *Cell* **80**, 279–284.
59. E. Yang, J. P. Zha, J. Jockel, L. H. Boise, C. B. Thompson, and S. J. Korsmeyer (1995) Bad, a heterodimeric partner for Bcl-xL and Bcl-2, displaces Bax and promotes cell death. *Cell* **80**, 285–291.
60. J. Zha, H. Harada, E. Yang, J. Jockel, and S. J. Korsmeyer (1996) Serine phosphorylation of death agonist BAD in response to survival factor results in binding to 14-3-3 not BCL-X(L) *Cell* **87**, 619–628 (see comments).
61. L. O'Reilly, D. Huang, and A. Strasser (1996) The cell death inhibitor Bcl-2 and its homologues influence control of cell cycle entry. *Embo. J.* **15**, 6979–6990.
62. M. B. Kastan, O. Onyekwere, D. Sidransky, B. Vogelstein, and R. W. Craig (1991) Participation of p53 protein in the cellular response to DNA damage. *Cancer Res.* **51**, 6304–6311.
63. D. P. Lane (1992) Cancer. p53, guardian of the genome. *Nature* **358**, 15–16 (news; comment).
64. W. El-Deiry et al. (1993) *WAF1*, a potential mediator of p53 tumor suppression. *Cell* **76**, 817–825.
65. C. Caelles, A. Helmberg, and M. Karin (1994) p53-dependent apoptosis in the absence of transcriptional activation of p53-target genes. *Nature* **370**, 220–223 (see comments).
66. C. Yin, C. M. Knudson, S. J. Korsmeyer, and T. Van Dyke (1997) Bax suppresses



- tumorigenesis and stimulates apoptosis *in vivo*. *Nature* **385**, 637–640.
67. S. W. Lowe, E. M. Schmitt, S. W. Smith, B. A. Osborne, and T. Jacks (1993) p53 is required for radiation-induced apoptosis in mouse thymocytes. *Nature* **362**, 847–849 (see comments).
  68. A. R. Clarke et al. (1993) Thymocyte apoptosis induced by p53-dependent and independent pathways. *Nature* **362**, 849–852 (see comments).
  69. A. Clarke et al. (1992) Requirement for a functional *Rb-1* gene in murine development. *Nature* **359**, 328–330.
  70. E.-H. Lee et al. (1992) Mice deficient for *Rb* are nonviable and show defects in neurogenesis and haematopoiesis. *Nature* **359**, 288–294.
  71. T. Jacks, A. Fazeli, E. Schmitt, R. Bronson, M. Goodell, and R. Weinberg (1992) Effects of an *Rb* mutation in the mouse. *Nature* **359**, 295–300.
  72. H. Symonds et al. (1994) p53-dependent apoptosis suppresses tumour growth and progression *in vivo*. *Cell* **76**, 703–711.
  73. Z. Xia, M. Dickens, J. Raingeaud, R. Davis, and M. Greenberg (1995) Opposing effects of ERK and JNK-p38 MAP kinases on apoptosis. *Science* **270**, 1326–1331.
  74. C. D. Nobes and A. M. Tolkovsky (1995) Neutralizing anti-p21ras Fabs suppress rat sympathetic neuron survival induced by NGF, LIF, CNTF and cAMP. *Eur. J. Neurosci.* **7**, 344–350.
  75. A. Geier, R. Beery, M. Haimsohn, R. Hemi, Z. Malik, and A. Karasik (1994) Epidermal growth factor, phorbol esters, and aurintricarboxylic acid are survival factors for MDA-231 cells exposed to adriamycin. *In Vitro Cell Dev. Biol. Anim.* **30a**, 867–874.
  76. R. Yao and G. M. Cooper (1995) Requirement for phosphatidylinositol-3 kinase in the prevention of apoptosis by nerve growth factor. *Science* **267**, 2003–2006.
  77. S. Kennedy, A. Wagner, S. Conzen, J. Jordan, A. Bellacosa, P. Tsichlis, and N. Hay (1997) The PI 3-kinase/akt signaling pathway delivers an anti-apoptotic signal. *Genes & Devel.* **11**, 701–713.
  78. H. Dudek et al. (1997) Regulation of neuronal survival by the serine-threonine protein-kinase akt. *Science* **275**, 661–665.
  79. A. Kauffmann-Zeh, P. Rodriguez-Viciana, E. Ulrich, C. Gilbert, P. Coffey, and G. Evan (1997) Suppression of c-Myc-induced apoptosis by Ras signalling through PI 3-kinase and PKB. *Nature* **385**, 544–548.
  80. P. Krammer, I. Behrmann, P. Daniel, J. Dhein, and K.-M. Debatin (1994) Regulation of apoptosis in the immune system. *Curr. Opin. Immunol.* **6**, 279–289.
  81. S. Nagata and P. Golstein (1995) The Fas death factor. *Science* **267**, 1449–1456.
  82. L. Tartaglia, T. Ayres, G. Wong, and D. Goeddel (1993) A novel domain within the 55 kd TNF receptor signals cell death. *Cell* **74**, 845–853.
  83. N. Itoh and S. Nagata (1993) A novel protein domain required for apoptosis. Mutational analysis of human Fas antigen. *J. Biol. Chem.* **268**, 10932–10937.
  84. M. Rothe, V. Sarma, V. W. Dixit, and D. V. Goeddel (1995) Traf2-mediated activation of NF-kappa-b by TNF receptor-2 and CD40. *Science* **269**, 1424–1427.
  85. M. Boldin, E. Varfolomeev, Z. Pancer, I. Mett, J. Camonis, and D. Wallach (1995) A novel protein that interacts with the death domain of Fas/Apo1 contains a sequence related to the death domain. *J. Biol. Chem.* **270**, 7795–7798.
  86. A. M. Chinnaiyan, K. O'Rourke, M. Tewari, and V. M. Dixit (1995) FADD, a novel death domain-containing protein, interacts with the death domain of Fas and initiates apoptosis. *Cell* **81**, 505–512.
  87. M. Muzio et al. (1996) FLICE, a novel FADD homologous ICE/CED-3-like protease, is recruited to the CD95 (Fas/Apo-1) death-inducing signaling complex. *Cell* **85**, 817–827.
  88. M. Boldin, T. Goncharov, Y. Goltsev, and D. Wallach (1996) Involvement of MACH, a novel

- MORT1/FADD-interacting protease, in Fas/APO-1- and TNF receptor-induced cell death. *Cell* **85**, 803–815.
89. A. Chinnaiyan and V. Dixit (1996) The cell-death machine. *Curr. Biol.* **6**, 555–562.
  90. X. Zhang et al. (1996) Up-regulation of Bcl-xL expression protects CD40-activated human B cells from Fas-mediated apoptosis. *Cell Immunol* **173**, 149–154.
  91. M. Thome et al. (1997) Viral FLICE inhibitory proteins (FLIPs) prevent apoptosis induced by death receptors. *Nature* **386**, 517–521.
  92. E. Alnemri, D. Livingston, D. Nicholson, G. Salvesan, N. Thornberry, W. Wong, and J. Yuan (1996) Human ICE/CED-3 protease nomenclature. *Cell* **87**, 171.
  93. Y. Lazebnik, A. Takahashi, G. Poirier, S. H. Kaufmann, and W. Earnshaw (1995) Characterization of the execution phase of apoptosis *in vitro* using extracts from condemned-phase cells. *J. Cell Sci.* **19**, 41–49.
  94. A. Takahashi and W. Earnshaw (1996) Ice-related proteases in apoptosis. *Curr. Opin. Gen. & Dev.* **6**, 50–55.
  95. Y. Jung, M. Miura, and J. Yuan (1996) Suppression of interleukin-1 beta-converting enzyme-mediated cell death by insulin-like growth factor. *J. Biol. Chem.* **271**, 5112–5117.
  96. V. Gagliardini, P. A. Fernandez, R. K. Lee, H. C. Drexler, R. J. Rotello, M. C. Fishman, and J. Yuan (1994) Prevention of vertebrate neuronal death by the crmA gene. *Science* **263**, 826–828.
  97. N. Bump et al. (1995) Inhibition of ICE family proteases by baculovirus antiapoptotic protein p35. *Science* **269**, 1885–1888.
  98. J. Xiang, D. Chao, and S. Korsmeyer (1996) Bax-induced cell death may not require interleukin-1 $\beta$ -converting enzyme-like proteases. *Proc. Natl. Acad. Sci. USA* **93**, 14559–14563.
  99. N. McCarthy, M. Whyte, C. Gilbert, and G. Evan (1997) Inhibition of Ced-3/ICE-related proteases does not prevent cell death induced by oncogenes, DNA damage, or the Bcl-2 homologue Bak. *J. Cell. Biol.* **136**, 215–227.
  100. J. Yuan and H. R. Horvitz (1992) The *Caenorhabditis elegans* cell death gene ced-4 encodes a novel protein and is expressed during the period of extensive programmed cell death. *Development* **116**, 309–320.
  101. A. M. Chinnaiyan, K. O'Rourke, B. R. Lane, and V. M. Dixit (1997) Interaction of CED-4 with CED-3 and CED-9: A molecular framework for cell death. *Science* **275**, 1122–1126 (see comments).
  102. M. S. Spector, S. Desnoyers, D. J. Hoepfner, and M. O. Hengartner (1997) Interaction between the *C. elegans* cell-death regulators CED-9 and CED-4. *Nature* **385**, 653–656.
  103. D. Wu, H. D. Wallen, and G. Nunez (1997) Interaction and regulation of subcellular localization of CED-4 by CED-9. *Science* **275**, 1126–1129 (see comments).
  104. C. James, S. Gschmeissner, A. Fraser, and G. Evan (1997) Ced-4 induces chromatin condensation in *S. pombe* and is inhibited by direct physical association with Ced-9. *Curr. Biol.* **7**, 246–252.
  105. M. Whyte (1996) ICE/Ced-3 proteases in apoptosis. *Trends Cell Biol.* **6**, 245–248.
  106. A. Fraser, N. McCarthy, and G. I. Evan (1996) Biochemistry of cell-death. *Curr. Opin. Neurobiol.* **6**, 71–80.
  107. T. Fernandes Alnemri et al. (1996) In-vitro activation of CPP32 and Mch3 by Mch4, a novel human apoptotic cysteine protease containing 2 FADD-like domains. *Proc. Natl. Acad. Sci. USA* **93**, 7464–7469.
  108. Y. A. Lazebnik, A. Takahashi, R. D. Moir, R. D. Goldman, G. G. Poirier, S. H. Kaufmann, and W. C. Earnshaw (1995) Studies of the lamin proteinase reveal multiple parallel biochemical pathways during apoptotic execution. *Proc. Natl. Acad. Sci. USA* **92**, 9042–9046.
  109. A. Fraser and G. Evans (1996) A license to kill. *Cell* **85**, 781–784.

110. P. J. Hartfield, G. C. Mayne, and A. W. Murray (1997) Ceramide induces apoptosis in PC12 cells. *FEBS Lett.* **401**, 148–152.
111. S. C. Wright, H. Zheng, and J. Zhong (1996) Tumor cell resistance to apoptosis due to a defect in the activation of sphingomyelinase and the 24 kDa apoptotic protease (AP24). *Faseb. J.* **10**, 325–332.
112. P. Santana et al. (1996) Acid sphingomyelinase-deficient human lymphoblasts and mice are defective in radiation-induced apoptosis. *Cell* **86**, 189–199.
113. G. J. Pronk, K. Ramer, P. Amiri, and L. T. Williams (1996) Requirement of an ICE-like protease for induction of apoptosis and ceramide generation by REAPER. *Science* **271**, 808–810.
114. M. Jacobson and M. Raff (1995) Programmed cell-death and *bcl-2* protection in very-low oxygen. *Nature* **374**, 814–816.
115. R. M. Kluck, E. Bossy Wetzels, D. R. Green, and D. D. Newmeyer (1997) The release of cytochrome c from mitochondria: A primary site for Bcl-2 regulation of apoptosis. *Science* **275**, 1132–1136 (see comments).
116. J. Yang, X. Liu, K. Bhalla, C. N. Kim, A. M. Ibrado, J. Cai, T. I. Peng, D. P. Jones, and X. Wang (1997) Prevention of apoptosis by Bcl-2: Release of cytochrome c from mitochondria blocked. *Science* **275**, 1129–1132 (see comments).
117. G. Kroemer, N. Zamzami, and S. A. Susin (1997) Mitochondrial control of apoptosis. *Immunol. Today* **18**, 44–51.

### Suggestions for Further Reading

118. M. D. Jacobson (1997) Apoptosis: Bcl-2-related protein get connected. *Curr. Biol.* **7**(5), R277–R281. An up-to-date review on the pathway of apoptosis.
119. M. D. Jacobson, M. Weil, and M. C. Raff (1997) Programmed cell death in animal development. *Cell* **88**, 347–354.
120. M. Sluysers, ed. (1996) *Apoptosis in Normal Development and Cancer*, Taylor and Francis, London.
121. E. Yang and S. J. Korsmeyer (1996) Molecular thanatopsis: A discourse on the Bcl-2 family and cell death. *Blood* **88**, 386–401.

### *ara* Operon

The pentose sugar L-arabinose is distributed widely as a component of **plant** polysaccharides in pectins, gums, and cell walls. It is also present in smaller amounts in the cell walls of some **bacteria**. Thus it is likely that, in its natural environment in the vertebrate intestine, the bacterium *Escherichia coli* periodically encounters arabinose released through degradation of these polymers (1). Not surprisingly, *E. coli*, like many other bacteria and [fungi](#), has ***ara* genes** that encode the [enzymes](#) and **transport** proteins necessary to utilize L-arabinose as a source of carbon and energy.

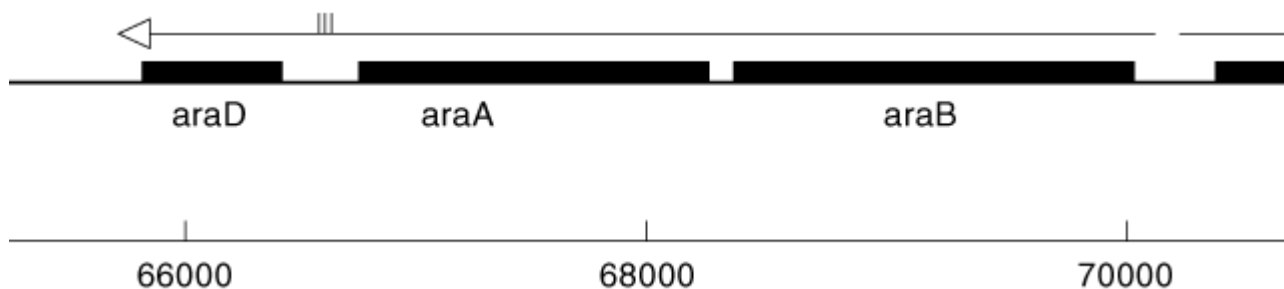
The *ara* genes of *E. coli* have been of exceptional importance to biology for several reasons. If the [lac operon](#) is the exemplar of regulation of **gene expression by repression of expression**, the *ara* operon is the counterpart as the model for positive control of gene expression, even though more complete analysis has revealed that each has both positive and negative features. (An [operon](#) is one gene or a linear sequence of related genes that are transcribed as one unit from a common initiation

point [see [Transcription](#)]. The transcript is then **translated** into protein products corresponding to the separate genes.) The first clear example of DNA looping as a feature of transcription control was the *ara* operon. The hypothesis of positive control of the *ara* operon was a challenge to establishment scientists, and its validation was the end of a fascinating story in the history of science. Moreover, the details of the operon are complex and interesting in themselves.

### 1. *ara* Operon of *Escherichia coli*

The classical *ara* operon of the bacterium *Escherichia coli* comprises three genes, *araBAD* (Fig. 1). This and four other operons, *araC*, *araE*, *araFGH*, and *araJ*, are uniquely associated with metabolism of L-arabinose (Table 1). Each set of genes is under control of the activator protein AraC, the product of the *araC* gene. The function of each gene is known, except for *araJ*, which may encode a protein that processes or transports an arabinose-containing polymer, or pumps potentially toxic arabinosides from the cell.

**Figure 1.** The positions of the *araBAD* operon and the *araC* regulatory gene on the *E. coli* chromosome. Transcription of nucleotide-pair region between *araB* and *araC*. Nucleotide location on the *E. coli* chromosome is shown beneath the gene the three [inverted repeat](#) REP sequence pairs that are assumed to produce three self-paired hairpin structures in the mRNA. The entire chromosome contains approximately 4,639,221 nucleotide pairs. The primary sequence information is in Refs. [24](#)



**Table 1. *ara* Genes and Gene Products**

| Gene        | Location on the <i>E. coli</i> chromosome (minutes) <sup>a</sup> | Size of gene (amino acid codons) | Activity of product                  |
|-------------|--|----------------------------------|--------------------------------------|
| <i>araA</i> | 1.4  | 500                              | L-arabinose isomerase                |
| <i>araB</i> | 1.5  | 566                              | L-ribulokinase                       |
| <i>araC</i> | 1.5  | 292                              | AraC regulatory protein              |
| <i>araD</i> | 1.4  | 231                              | L-ribulose-5-phosphate-4-epimerase   |
| <i>araE</i> | 64.2   | 472                              | L-arabinose transport, low affinity  |
| <i>araF</i> | 42.8   | 329                              | L-arabinose transport, high affinity |
| <i>araG</i> | 42.7   | 504                              | L-arabinose transport,               |

|             |      |     |   |
|-------------|------|-----|---|
| <i>araH</i> | 42.6 | 329 | high affinity<br>L-arabinose transport,<br>high affinity                      |
| <i>araJ</i> | 8.5  | 394 | Transport or<br>processing of<br>polymer?<br>Efflux of toxic<br>arabinosides? |

---

<sup>a</sup> Location is based on a 100 minute circular chromosome.

### 1.1. Uptake and Utilization of Arabinose

Two independent systems deliver arabinose from the environment across the cell [membrane](#) into the cell. *araE* encodes a [membrane protein](#) that mediates arabinose uptake via proton **symport** and is the lower affinity transporter ( $K_M = 50 \mu\text{M}$ ) (2). The *araFGH* operon encodes a **periplasmic** arabinose-binding protein (*araF*), a probable [ATPase](#) subunit (*araG*), and a membrane protein (*araH*), which together mediate ATP-driven arabinose transport. This transporter shows higher affinity for arabinose ( $K_M = 1 \mu\text{M}$ ) than AraE, but lower capacity (2, 3).

Internal arabinose is converted in three steps to D-xylulose-5-phosphate, a metabolite in the pentose-phosphate shunt pathway and one that is not unique to arabinose metabolism. The enzyme mediating the first step, L-arabinose isomerase, has a low affinity for arabinose with a [KM \(Michaelis constant\)](#) of 60 mM (4); this suggests that cells growing on arabinose have a very high internal arabinose concentration (5). The glucose-specific phosphotransferase enzyme IIA<sup>Glc</sup> when unphosphorylated inhibits the isomerase (unpublished results cited in 6) This inhibition may be one of the causes of the preferential use of glucose when both arabinose and glucose are in the environment.

The product of isomerase activity, L-ribulose, is converted to L-ribulose-5-phosphate by L-ribulokinase, and the phosphorylated compound is converted in turn to xylulose phosphate by the epimerase encoded by *araD*. Arabinose inhibits growth of *araD* mutants on other nutrients, presumably because accumulation of a high concentration of ribulose phosphate is toxic. Thus secondary mutants lacking isomerase or kinase activity as the result of *araA*, *araB*, or *araC* mutation, and therefore not forming ribulose phosphate, can be selected by plating an *araD* population on broth plates containing arabinose (7).

### 1.2. Regulation of *ara* Operon Expression

In the absence of arabinose, the *ara* genes are essentially not expressed, except for the regulatory gene, *araC*. On exposure to L-arabinose, all of the *ara* genes are activated, transcription of *araBAD* begins within five seconds, and the Ara proteins appear within several minutes, allowing growth on the sugar (5). The true inducer has been shown by [X-ray crystallography](#) of the inducer-AraC protein complex to be -L-arabinose (8). L-lyxose also induces the *ara* genes (9). However induction of the *ara* genes by lyxose is unlikely to occur in nature because L-lyxose is rare. *araC* mutants have been obtained that are inducible by D-fucose or by -methyl-L-arabinside, normally inhibitors of induction (1, 10, 11).

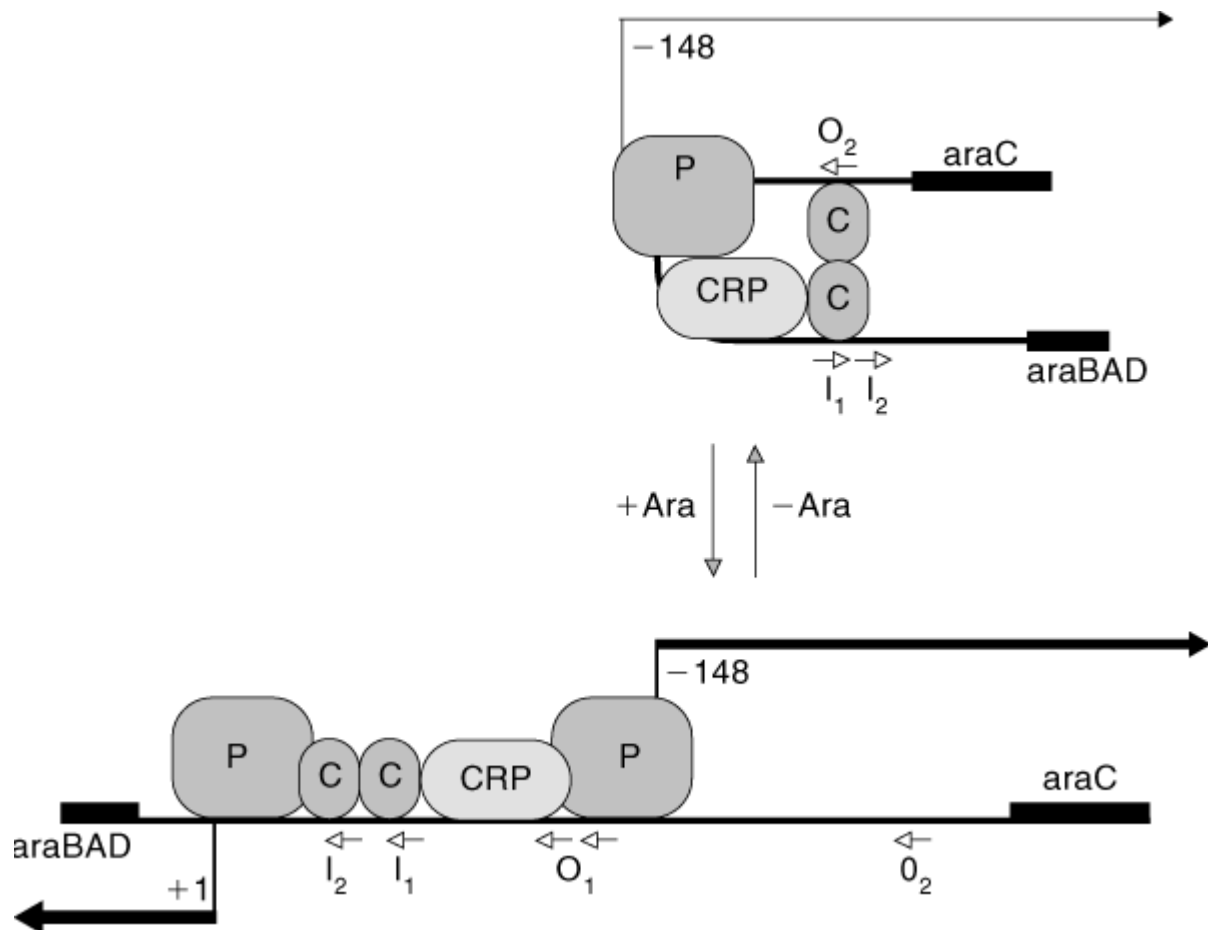
*In vivo* studies using cells lacking both arabinose transport systems showed the  $K_M$  for induction by arabinose to be unexpectedly high, nearly 10 mM (11), although, as noted above, the internal arabinose concentration is likely to be higher than this during growth on arabinose because of the low affinity of the isomerase for the pentose. The high concentration-dependence of induction suggests that there might be [natural selection](#) against a mutant isomerase with a better affinity for arabinose.

[Complementation](#) studies showed that the *presence* of AraC protein is necessary for *araBAD* gene expression (12, 13), a discovery unanticipated by those who believed all gene activation resulted from *removal* of repressors. Later, AraC was shown to be required for expression of the other *ara* operons as well (except for *araC*). *araC* is the prototype for a large family of regulatory genes that share sequence similarity in the DNA binding region and are homologous and presumably related through [evolution](#) (14).

[Cyclic AMP \(3',5'-cyclic AMP, cAMP\)](#) bound to [cyclic AMP receptor protein \(CRP\)](#) is also a positive regulator for all the *ara* operons. Expression of many genes is controlled by availability of CRP-cAMP; lack of expression due to low CRP-cAMP is referred to as [catabolite repression](#). *In vitro* transcription of *araBAD* mimics that *in vivo* in that it requires both the AraC and CRP regulatory proteins with their bound ligands. Analysis of transcription *in vitro*, and further *in vivo* studies, have given a broad understanding at the molecular level of control of *araBAD* transcription although details remain to be determined (5).

The AraC [protein structure](#) has two domains connected by a flexible polypeptide linker. The *N*-terminal domain binds arabinose and is responsible for formation of the active dimeric form of the protein. The *C*-terminal domain binds DNA at specific sites, with similar sequences, upstream of each *ara* operon. In the regulatory region between the divergently transcribed *araBAD* and *araC* operons (Fig. 1), there are five sites at which AraC can bind (Fig. 2).

**Figure 2.** Activation of transcription of the *araBAD* operon by arabinose. CRP refers to CRP-cAMP, and P refers to RNA polymerase. AraC protein is shown as a dimer bound to I<sub>1</sub> and O<sub>2</sub> in the absence of arabinose and to I<sub>1</sub> and I<sub>2</sub> in the presence of arabinose. I<sub>1</sub>, I<sub>2</sub>, O<sub>1</sub> (a double site), and O<sub>2</sub> are all potential binding sites for the AraC protein and conform to the consensus sequence 5'-TAGCN<sub>7</sub>TCCATA-3' (reading in the direction of the arrows) (16) although there is considerable variation among sites. Note how the mode of pairing of the dimers differs depending on the presence of arabinose (5, 8). On addition of arabinose, transcription of *araBAD* is initiated and expression of *araC* is increased (bottom). After a few minutes AraC-arabinose dimers are thought to bind to O<sub>1</sub> and O<sub>2</sub> so as to form a new DNA loop and reduce *araC* transcription to that characteristic of cells in the absence of arabinose (not shown).



In the absence of arabinose, the two DNA-binding domains of the AraC dimer are oriented so that they can not readily bind both of the adjacent I sites at the same time. Instead, the AraC dimer contacts  $I_1$  and  $O_2$ , thereby forming a DNA loop within the region between *araB* and *araC* (Fig. 2) (15). On addition of arabinose, the dimer undergoes a conformational shift such that the two DNA-binding domains preferentially bind  $I_2I_1$  (Fig. 2) (16). The presence of AraC at  $I_2$  stimulates addition of **RNA polymerase** and open complex formation, and transcription of *araBAD* commences, if CRP-cAMP is present (5). Although this model was proposed and refined before detailed structural information was available, X-ray crystallographic studies of the AraC N-terminal domain and linker support the model (8).

Addition of arabinose also affects expression of *araC*. When the  $I_1$ - $O_2$  loop is opened, transcription of *araC* is accelerated (Fig. 2). After a few minutes, AraC-arabinose dimers are thought to reform DNA loops, this time by bridging  $O_1$  and  $O_2$  (17, 18).  $O_1$ - $O_2$  looping does not regulate *araBAD* transcription, but interferes with RNA polymerase binding and initiation at the *araC* promoter; this reduces the rate of *araC* transcription to that characteristic of cells in the absence of arabinose (5).

*araC* is controlled by CRP-cAMP as well. CRP-cAMP binding increases, but is not essential for, *araC* transcription; it is necessary for substantial expression of the other *ara* operons.

The mechanism by which AraC-arabinose, CRP-cAMP, and RNA polymerase interact to trigger transcription is not clear. Bound CRP-cAMP helps to open the  $I_1I_2$  repression loop on addition of arabinose, but it probably activates RNA polymerase binding or initiation as well, either by direct contact or indirectly through contact with AraC (5). (CRP-cAMP does not aid AraC binding.)

Several kinds of CRP-cAMP-independent mutants have been obtained. One class (*rpoD*) has an altered RNA polymerase **sigma factor**-70 subunit (19), while another (*araC*) is altered in AraC itself (20). Other mutants (*araI*) are transcribed by RNA polymerase independently of the AraC or CRP proteins, and they result from changes at the RNA polymerase binding site in the *araBAD* promoter (21).

## 2. Other Mechanisms for Arabinose Degradation

The pathway for arabinose degradation that is common to *E. coli* and many other bacteria is not the sole means for utilizing L-arabinose in the biological world, nor is the regulatory model described above the only means for control of *ara* enzyme synthesis. At least five different pathways exist (1), and different regulatory systems are found even among organisms using the same pathway. Although the work on other organisms is still at an early stage, it is known that the **gram-positive** bacterium, *Bacillus subtilis*, utilizes the same degradatory pathway for arabinose utilization as *E. coli* but appears to use a simple repressor as the means of control of *ara* gene expression (22). The fungus *Aspergillus* uses a completely different pathway for arabinose utilization and, not surprisingly, appears to have different control systems as well (23). It will be remarkable if the *ara* genes in these other organisms are regulated by the simple repressors once hypothesized for the *ara* genes in *E. coli*.

## Bibliography

“*ara* Operon” in , Vol. 1, pp. 187–190, by Robert B. Helling, University of Michigan, Department of Biology, Ann Arbor, Michigan; “*ara* Operon” in (online), posting date: January 15, 2002, by Robert B. Helling, University of Michigan, Department of Biology, Ann Arbor, Michigan.

1. M.E. Doyle (1974) Ph.D thesis, The University of Michigan, University Microfilms, Ann Arbor, MI, 71 pp.
2. R. Schleif (1969) *J. Mol. Biol.* **46**, 185–196.
3. R.W. Hogg (1977) *J. Supramol. Struct.* **6**, 411–417.
4. J.W. Patrick and N. Lee (1968) *J. Biol. Chem.* **243**, 4312–4318.
5. R. Schleif (2000) *Trends in Genetics* **16**, 559–565.
6. C. Hoischen et al. (1996) *J. Bacteriol.* **178**, 6082–6086.
7. H. Boyer, E. Englesberg, and R. Weinberg (1962) *Genetics* **47**, 417–425.
8. S.M. Soisson et al. (1997) *Science* **276**, 421–425.
9. E. Ibanez et al. (2000) *J. Bacteriol.* **182**, 4625–4627.
10. S. Beverin, D.E. Sheppard, and S.S. Park (1971) *J. Bacteriol.* **106**, 107–112.
11. M.E. Doyle et al. (1972) *J. Bacteriol.* **110**, 56–65.
12. R.B. Helling and R. Weinberg (1963) *Genetics* **48**, 1397–1410.
13. E. Englesberg et al. (1965) *J. Bacteriol.* **90**, 946–957.
14. M.-T. Gallegos et al. (1997) *Microbiol. Mol. Biol. Rev.* **61**, 393–410.
15. K. Martin, L. Huo, and R. Schleif (1986) *Proc. Nat. Acad. Sci. USA* **83**, 3654–3658.
16. N. Lee, C. Francklyn, and E. P. Hamilton (1987) *Proc. Nat. Acad. Sci. USA* **84**, 8814–8818.
17. E. P. Hamilton and N. Lee (1988) *Proc. Nat. Acad. Sci. USA* **85**, 1749–1753.
18. L. Huo, K. J. Martin, and R. Schleif (1988) *Proc. Nat. Acad. Sci. USA* **85**, 5444–5448.
19. J.C. Hu and C.A. Gross (1985) *Mol. Gen. Genet.* **199**, 7–13.
20. L. Heffernan, R. Bass, and E. Englesberg (1976) *J. Bacteriol.* **126**, 1111–1131.
21. A.H. Horwitz, C. Morandi, and G. Wilcox (1980) *J. Bacteriol.* **142**, 659–667.
22. L. J. Mota, L. M. Sarmiento, and I. Sa-Nogueira (2001) *J. Bacteriol.* **183**, 4190–4201.
23. G.J.G. Ruijter et al. (1997) *Microbiol.* **143**, 2991–2998.



24. R.G. Wallace, N. Lee, and A.V. Fowler (1980) *Gene* **12**, 179–190.
25. N. Lee et al. (1986) *Gene* **47**, 231–244.

### Suggestions for Further Reading

26. J. Beckwith (1987) "The operon: an historical account". In *Escherichia coli and Salmonella typhimurium. Cellular and Molecular Biology* (F.C. Neidhardt, J.L. Ingraham, K.B. Low, B. Magasanik, M. Schaechter, and H.E. Umbarger, ed.), American Society for Microbiology, Washington, DC, pp. 1439–1443. (Describes the difficulty in achieving acceptance of the model of positive control.)
27. J. Gross and E. Englesberg (1959) Determination of the order of mutational sites governing L-arabinose utilization in *Escherichia coli* B/r by transduction with phage P1bt. *Virology* **9**, 314–331. (This is the research paper in which the initial studies with the *ara* operon are described. It has served as a model guiding analysis of many other genes.)
28. R. Schleif (1993) "Induction, repression and the *araBAD* operon". In *Genetics and Molecular Biology*, 2nd ed., Johns Hopkins University Press, pp. 359–383. (Systematic description of the arabinose genes and enzymes, and how it has been studied. Thoughtful and stimulating, with problems (and answers at the end of the book).)
29. R. Schleif (2000) Regulation of the L-arabinose operon of *Escherichia coli*. *Trends in Genetics* **16**, 559–565. (Recent review of control of transcription of the *ara* genes focusing on the AraC protein and its structural changes.)

## Arabidopsis

*Arabidopsis thaliana* (L.) Heynh is a dicotyledonous plant that is ideally suited for molecular-genetic studies. It is generally accepted as a model for unraveling the molecular mechanisms involved in plant growth and development, biochemical pathways, cell biology, physiology, and pathogenic interactions. Minimal genomic DNA content, few repetitive DNA sequences, and small gene families account for its technical and biological simplicity. *Arabidopsis* is a natural diploid amenable to laboratory-scale genetic experiments. It is small (15 cm to 30 cm) and is grown at high density (100 plants per 0.5 m<sup>2</sup>) without seed contamination. It has a short life cycle (6 weeks to 3 months), high seed production by self-fertilization (up to 10,000 seeds per plant), is mutagenized easily, and with its 200 different ecotypes is a natural source of genetic variation. International coordination of *Arabidopsis* research has resulted in fast exchange of material, methodology, and information and in the creation of large-scale research programs, such as the genome-sequencing and expressed-sequence tag (EST) projects. The completion of the entire genome sequence was a milestone in plant biology because for the first time a molecular overview has been obtained of common and different pathways between plants and other eukaryotes. Moreover, every plant gene is now accessible for functional analysis. In the post-genomic era *Arabidopsis* research maintains its pioneering position in the field of plant science.

### 1. History

The first botanical description of *Arabidopsis* by Johannes Thal goes back to 1577. In 1907, Laibach studied the continuity of the chromosomes (5 per haploid genome) by using the plant, and he was the first to emphasize the advantages of the species for genetic analyses (1). In 1976, Bennett and Smith showed that *Arabidopsis* has the smallest nuclear DNA content of the angiosperms analyzed. In that period, extensive chemical and irradiation mutagenesis of the plant was performed (2). Koornneef *et*

al. (3) published the first genetic map containing 76 morphological markers. In the mid eighties, Meyerowitz and coworkers demonstrated the small size and low complexity of the genome, which provoked a general interest in using the plant as an experimental model (4). In 1988, efficient transformation methods opened the potential of transgenic research in the species (5). At the same time, saturation mutagenesis of the genome by insertion of heterologous DNA was initiated and resulted in large collections that became available to the scientific community (6). In 1989, the US National Science Foundation launched the 'Multinational Coordinated Long Range Plan for Arabidopsis Genome Research', steered by an international board of scientists, with the aim of promoting *Arabidopsis* as a model system for plants, in analogy to other models such as *Drosophila melanogaster* and *Caenorhabditis elegans*. The major achievements of this initiative were a seed and DNA stock centers (The Arabidopsis Biological Resource Center [Ohio State University, Columbus, OH, USA] and the Nottingham Arabidopsis Stock Centre [University of Nottingham, Loughborough, UK]), a database, and joint efforts for physical mapping and sequencing of the genome and for gene identification. The *Arabidopsis* genome sequencing project was initiated in Europe at the end of 1993, followed by an American and a Japanese initiative. This international consortium, the Arabidopsis Genome Initiative completed the entire sequence by the end of 2000 (7).

## 2. Classification, Geographical Distribution, and Ecotypes

*Arabidopsis thaliana* is an annual herb of the mustard family (Brassicaceae, previously named Cruciferae), has bisexual flowers, and is typified by a cross-shaped corolla, tetradynamous stamen (four long and two short ones), and capsular fruit (siliques). The genus *Arabidopsis* consists of 27 species and has been classified under a new tribe, the Arabidae, based on classical morphological and molecular phylogenetic studies. *Arabidopsis* is a facultative long-day plant, meaning that long days accelerate the initiation of flowering. It originated from Eurasia and North Africa, but is now a common weed in the temperate regions of the Northern Hemisphere. Its broad geographical distribution resulted in natural variation. Approximately 200 ecotypes (wild populations) have been registered. These ecotypes represent a natural source of heritable variation: for instance, differences are observed in flowering time, response to cold, fresh weight production, and pathogen resistance. Among ecotypes, a DNA sequence polymorphism of up to 1.4% in low-copy DNA has been measured. Deletions, insertions, and substitutions have occurred. Frequently used ecotypes are Columbia, Niederzenz, Wassilewskija, and the laboratory strain Landsberg *erecta*. The Columbia ecotype was used as standard for the genome sequence, whereas other ecotypes are preferred for mutagenization. DNA polymorphisms are exploited in F2 populations for mapping and quantitative trait locus (QTL) analyses.

## 3. Genome Structure and Organization

The size of the nuclear genome of *Arabidopsis* was estimated to be between 50 to 150 Mb as determined by reassociation kinetics (see [C-Value](#)), flow cytometry, or electron microscopy. From recent physical mapping data, this size was refined to 100 to 140 Mb, which is 3- to 400-fold smaller than that of other members of the angiosperms.

In 2000, the genome sequence of *Arabidopsis* was published, covering 115.4 megabases of the 125-megabases genome (7). A whole genome duplication as well as lateral gene transfer from a cyanobacterial-like ancestor appears to be part of the evolution of the *Arabidopsis* genome. The *Arabidopsis* genome consists for 80% of single- and low-copy DNA. Four classes of highly repeated DNA have been identified, which represent only 10% of the genome and are located mainly at the telomeres and around the centromeres. Ribosomal DNA accounts for 6% of the genome and is localized at the top of chromosomes 2 and 4. In addition to the genomic sequencing, more than 30,000 redundant *Arabidopsis* ESTs have been sequenced and submitted to public databases. Comparison of EST data with data of the genome sequencing program indicates that approximately 60% of the genes are represented by an EST.

From the sequence data, several conclusions can be drawn on gene organization and DNA

composition: the gene density is one gene every 4.5 kb, the average gene length is 2 kb, genes have between 0 to 30 introns, 25,498 predicted transcripts encode proteins of approximately 11,000 families, and roughly 30% of the genes could not be assigned to functional categories. The gene families are either dispersed or in tandem arrays; long stretches (approximately 120 kb) of unique or low-copy DNA are interspersed with short stretches of a few kb of moderately repeated DNA; the GC content of the DNA is 35%, and methylation occurs in 6% of the cytosine bases.

The availability of the genome sequences of *Arabidopsis* and several other organisms allows comparisons at the genome level. For example, *Arabidopsis*, *Drosophila* and *C. elegans* have a similar amount of approximately 11,000 to 15,000 different protein types, indicating that this is the minimal number of proteins needed for a functional multicellular organism. Basic processes, such as translation, appear to be conserved across kingdoms, whereas more specialized processes use proteins that differ between plants and animals. These include membrane channels and transporters, components of signal transduction pathways, and transcription factors (8). Plants contain roughly 150 unique protein families (7).

Information on the genome project and on research in general on *Arabidopsis* is centralized in “The Arabidopsis Information Resource” (TAIR) database. A home-page on the Internet (<http://www.arabidopsis.org/home.html>) provides a lot of information on ongoing research (genetic and physical maps, genome sequence information, links to public databases, seed and DNA stock centers, etc.).

#### 4. Genetics and Physical Mapping

Based on metaphase staining and genetic linkage group analyses, the haploid chromosome number is 5, ranging in size from 13.4 Mb to 25.4 Mb. The genetic map comprises more than 460 loci. The ratio of physical to genetic distance between markers on average is 200 kb per cM. However, from the physical map construction of chromosome 4, recombination hot spots (30–50 kb/cM) and low spots ( $\geq 550$  kb/cM) were found (9). Many tools are available to map a mutation, for example, recessive visible markers, codominant embryo-lethal markers, dominant selectable markers on the located T-DNA and *Activation/Dissociation* insertions, restriction fragment length polymorphism (RFLP)-derived or PCR-based molecular markers, such as microsatellites, insertion/deletions (INDEL), and single nucleotide polymorphisms (SNP). Several molecular marker maps have been constructed based on RFLP, rapid amplified polymorphic DNA, or amplified fragment length polymorphism (AFLP) as well as on different mapping populations. A combined map, made by statistical integration, gives an approximate position and order for the markers. Recombinant inbred (RI) lines, derived from a cross between Columbia and Landsberg *erecta* (10), have been used to locate 1265 molecular markers to date. The physical map consists of contigs of DNA clones that are correlated with the mapped markers. Currently, YAC, bacterial artificial chromosome (BAC), and phage P1 artificial chromosome (PAC) contig maps are available that cover the entire genome.

#### 5. Scientific Advances and Applications

The molecular-genetic approach in *Arabidopsis* research together with cell biology tools, such as confocal microscopy, laser cell ablation, and dye-loading have led to major breakthroughs in plant developmental biology (11-13). Tremendous progress has been made in understanding the molecular control of meristem identity in vegetative meristems, during flower initiation and flower organ formation, embryo development, and pattern formation during embryogenesis, root development, epidermal cell fate specification in root hair, and trichome formation. Genes have been identified that are involved in hormone perception, biosynthesis, and signal transduction. The first hormone receptor for plants has been characterized in *Arabidopsis* (14). Much of the molecular insights into light perception and signal transduction, cell cycle regulation, and in disease resistance in higher plants comes from studies in *Arabidopsis* (15-17). Biochemical items, such as cell wall biology, are being explored. Within the next 10 years, the function of every gene will be determined. This will be achieved by reverse genetics techniques, such as knockouts and targeting-induced local lesions in

genomes (TILLING) (18), to create allelic diversity and by implementation of genome-wide methods, such as transcript profiling and microarrays. In general, biological processes will be analyzed by genomics and bioinformatics tools in order to identify all the molecular components involved.

The *Arabidopsis* genes and mutants are resources that are exploited to isolate orthologs from other species and to test their functional conservation (19), or that are used for the genetic modification of even distantly related crop plants (20). The molecular markers within contigs in *Arabidopsis* have been used for comparative mapping with *Brassica* spp. Colinearity in 5- to 10-cM regions has been demonstrated between the *Arabidopsis* genome and that of *Brassica nigra* (21). This implies that information and markers obtained from the physical mapping in *Arabidopsis* can be applied to syntenic genomic regions in mustard crops to analyze important traits in breeding programs. Microarrays of *Arabidopsis* have proven to be useful for the analysis of transcription profiles of developing seeds in related species, such as *Brassica* (22), and will be a great help in the study of seed biology in engineered plants. It is to be predicted that *Arabidopsis* will be used as a reference system to accelerate knowledge accumulation in crop species. DNA sequence analysis showed the presence of genes in the *Arabidopsis* genome that suggest the existence of secondary metabolism pathway in this plant (7). Hence, *Arabidopsis* might be a helpful tool to unravel these pathways in medicinal plants.

## Bibliography

“Arabidopsis” in , Vol. 1, pp. 190–192, by Mieke Van Lijsebettens, Nancy Terryn, and Marc Van Montagu, Ghent University, Departments of Molecular Genetics and Plant Genetics, Flanders Interuniversity Institute for Biotechnology (VIB), K.L. Ledeganckstraat 35, Gent, Belgium, B-9000, milij@gengenp.rug.ac.be; “Arabidopsis” in (online), posting date: January 15, 2002, by Mieke Van Lijsebettens, Nancy Terryn, and Marc Van Montagu, Ghent University, Departments of Molecular Genetics and Plant Genetics, Flanders Interuniversity Institute for Biotechnology (VIB), K.L. Ledeganckstraat 35, Gent, Belgium, B-9000, milij@gengenp.rug.ac.be.

1. F. Laibach (1943) Bot. Archiv **44**, 439–455.
2. G.P. Rédei (1975) Ann. Rev. Genet. **9**, 111–127.
3. M. Koornneef et al. (1983) J. Hered. **74**, 265–272.
4. R.E. Pruitt and E.M. Meyerowitz (1986) J. Mol. Biol. **187**, 169–183.
5. D. Valvekens, M. Van Montagu, and M. Van Lijsebettens (1988) Proc. Natl. Acad. Sci. USA **85**, 5536–5540.
6. K.A. Feldmann (1991) Plant J. **1**, 71–82.
7. Arabidopsis Genome Initiative (2000) Nature **408**, 796–815.
8. J.L. Riechmann et al. (2000) Science **290**, 2105–2110.
9. R. Schmidt et al. (1995) Science **270**, 480–483.
10. C. Lister and C. Dean (1993) Plant J. **4**, 745–750.
11. C. van den Berg et al. (1997) Nature **390**, 287–289.
12. M.F. Yanofsky et al. (1990) Nature **346**, 35–38.
13. J.A. Long et al. (1996) Nature **379**, 66–69.
14. G.E. Schaller and A.B. Blecker (1995) Science **270**, 1809–1811.
15. J.L. Dangl (1995) Cell **80**, 363–366.
16. C. Lin et al. (1995) Science **269**, 968–970.
17. A. Hemerly et al. (1992) Proc. Natl. Acad. Sci. USA **89**, 3295–3299.
18. C.M. McCallum, L. Comai, E.A. Greena, and S. Henikoff (2000) Nature Biotechnol. **18**, 455–457.
19. V.F. Irish and Y.T. Yamamoto (1995) Plant Cell **7**, 1635–1644.

20. D. Weigel and O. Nilsson (1995) *Nature* **377**, 495–500.
21. U. Lagercrantz, J. Putterill, G. Coupland, and D. Lydiate (1996) *Plant J.* **9**, 13–20.
22. T. Girke et al. (2000) *Plant Physiol.* **124**, 1570–1581.

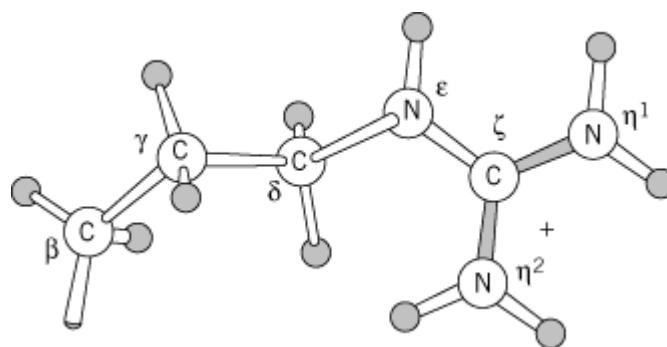
### Suggestions for Further Reading

23. J. Bowman (1994). *Arabidopsis: an Atlas of Morphology and Development*, Springer-Verlag, New York, NY.
24. C. Koncz, N.-H. Chua and J. Schell (1992). *Methods in Arabidopsis Research*, World Scientific, Singapore.
25. J.M. Martínez-Zapater and J. Salinas (1998). *Arabidopsis Protocols* (Methods in Molecular Biology Series), The Humana Press, Totowa, NJ.
26. E.M. Meyerowitz and C.R. Somerville (1994). *Arabidopsis* (Cold Spring Harbor Monograph Series, Vol. **27**), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

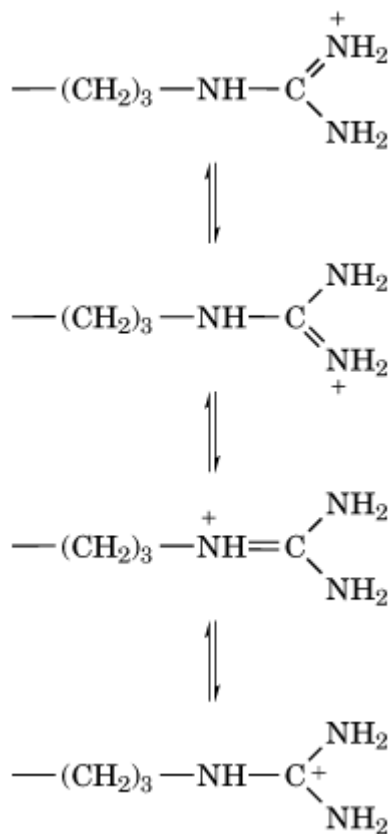
### Arginine (Arg, R)

The [amino acid](#) arginine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to six **codons**—CGU, CGC, CGA, CGG, AGA, and AGG—and represents approximately 5.7% of the residues of the proteins that have been characterized. The arginyl residue incorporated has a mass of 156.19 Da, a **van der Waals volume** of 148 Å<sup>3</sup>, and an [accessible surface](#) area of 241 Å<sup>2</sup>. Arg residues have average conservation during [divergent evolution](#); they are interchanged most frequently in **homologous** proteins with [lysine](#), the other basic residue.

The Arg side chain consists of three nonpolar methylene groups and the strongly basic d-guanido group:



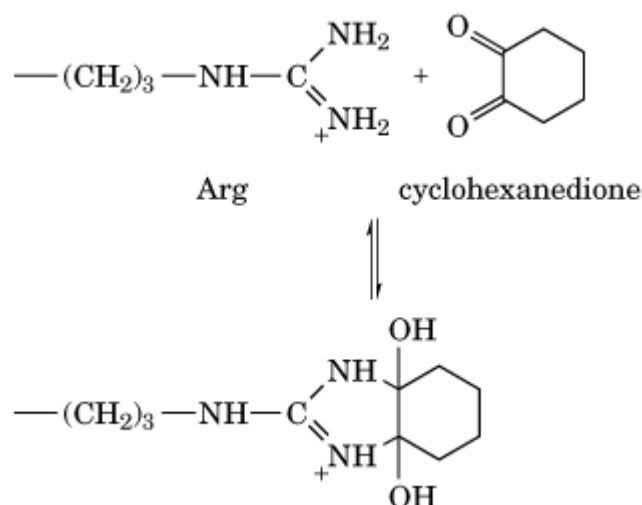
With a  $pK_a$  value usually of about 12, the guanido group is ionized over the entire pH range in which proteins exist naturally. The ionized guanido group is planar as a result of resonance:



and the positive charge is effectively distributed over the entire group. In the protonated form, the guanido group is unreactive, and only very small fractions of the nonionized form are present at physiological pH values. The guanido groups of Arg residues are almost invariably at the surfaces of native protein structures, and virtually no Arg residues are fully buried, but the nonpolar part of the side chain, and the adjoining polypeptide backbone, are frequently buried within the interior. Arg residues favor the *alpha*-helical conformation in model peptides and also occur most frequently in that **secondary structure** in folded [protein structures](#).

[Proteinases](#) frequently cleave polypeptide chains adjacent to Arg residues, as in the processing of pro-hormones, such as pro-**insulin**, at pairs of basic residues.

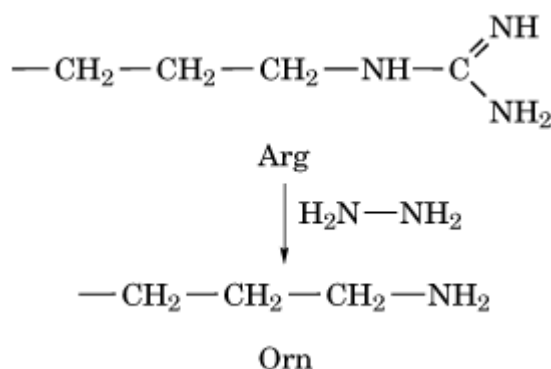
The guanido group can form heterocyclic condensation products with 1,2- and 1,3-dicarbonyl compounds, such as phenylglyoxal, 2,3-butanedione, and 1,2-cyclohexanedione:



These reactions occur readily because the distance between the two carbonyl groups of the reagents closely matches that between the two unsubstituted nitrogen atoms of the guanido group. The adduct

formed can be stabilized further by the presence of borate, which complexes with the adjacent hydroxyl groups.

The guanido group can be cleaved by hydrazine ( $\text{H}_2\text{NNH}_2$ ) to produce the side chain of ornithine:



This reaction is, however, often accompanied by cleavage of the polypeptide backbone.

#### Suggestions for Further Reading

E. L. Smith (1977) Reversible blocking at arginine by cyclohexanedione, *Meth. Enzymol.* **47**, 156–161.

A. Honegger et al. (1981) Chemical modification of peptides by hydrazine, *Biochem. J.* **199**, 53–59.

R. B. Yamasaki et al. (1980) Modification of available arginine residues in proteins by *p*-hydroxyphenylglyoxal, *Anal. Biochem.* **109**, 32–40.

## Ascaris

*Ascaris* is a genus of large animal parasitic [nematodes](#). Common species include *Ascaris equorum*, a horse parasite, *Ascaris suum*, a pig parasite, and *Ascaris lumbricoides*, which parasitizes humans. Adult ascarid nematodes, which reside in the host intestine, range from several inches to several feet in length, and females can produce more than 2 million eggs per day. Fertilized eggs are released into the environment by defecation, where they remain dormant until ingested by another host. The acidic environment of the stomach and the increase in temperature trigger resumption of embryonic development, resulting in first-stage (L1) juveniles (often termed *larvae*) inside the egg shell. After molting once, the larvae hatch from the egg as second-stage (L2) juveniles. Like all nematodes (*q.v.*), *Ascaris* molts through three more juvenile stages (L3–L4) and then to a sexually mature adult. Stages in the parasitic life cycle are tissue-specific: The L2 burrows through the wall of the small intestine and is carried via the bloodstream or lymphatic system to the lungs, bronchi, and trachea. It molts to L3 and migrates up the trachea and down the esophagus, returning to the small intestine where it molts to L4 and then, about two to three months after the initial infection, to a mature adult that begins producing eggs. About 25% of the human population is infected by *Ascaris lumbricoides*, primarily in developing countries. *Ascaris* infections are generally debilitating but not life-threatening, although heavy infections can cause severe complications resulting from intestinal obstruction, peritonitis, or allergic reaction.

As an experimental organism, *Ascaris* was important to early embryologists and cytologists such as Theodor Boveri and Otto zur Strassen in the 1890s, who took advantage of several favorable

properties of the early [embryos](#) that are still exploited by present-day researchers. The embryos are obtainable in large numbers from gravid females; they are transparent, facilitating observation by light [microscopy](#); they can be stored in their original dormant state and their development initiated synchronously by high temperature or acid shock; and the pattern and timing of their embryonic cleavages, like those of most nematodes (*q.v.*), are invariant from embryo to embryo, with many cell fates apparently determined as the cells are born. The species studied most extensively by Boveri, *Ascaris megalocephala*, has the additional advantage of having only one large [chromosome](#), which was convenient for cytological studies. *Ascaris* also has the peculiarity that there is chromosome fragmentation and diminution in all of the **somatic** precursor cells as they separate from the **germ line** in the early embryo. Although later work has shown that most of the discarded DNA does not include coding sequences, this feature allowed Boveri to demonstrate that cytoplasmic components specific to the germ line were responsible for maintaining chromosomal integrity in these cells, which supported the view that nonrandomly segregating cytoplasmic determinants were responsible for dictating early cell fates.

In the modern era, *Ascaris* has also become an important system for research in neurobiology. Part of the rationale behind Sydney Brenner's choice of the small free-living soil nematode *Caenorhabditis elegans* (see [Caenorhabditis](#)) for genetic and ultrastructural analysis of a simple nervous system and its development was the hope that *Ascaris*, only distantly related but much more suitable for electrophysiology because of its larger size, would allow functional analysis of a homologous nervous system. This hope has been largely realized: The two nervous systems appear to be very similar in both structure and function, and information from each has been valuable in understanding the other.

#### Suggestion for Further Reading

G. O. Poinar (1983) *The Natural History of Nematodes*, Prentice-Hall, Englewood Cliffs, NJ, p. 323.

## Asexual

The term “asexual” means without [sex](#); it is an adjective applied to multicellular organisms with missing or atrophied sexual organs and to forms of reproduction that do not involve the fusion of male and female gametes (see [Asexual Reproduction](#)).

The initial concept of sex referred to diploid, multicellular species in which new individuals are produced by the fusion of two haploid gametes, a sperm and an ovum (sexual reproduction). Individuals of many of these species may be classified into two sexes, *male* and *female*, depending on whether they are able to produce [sperm](#) or ova ([eggs](#)), respectively. This is the situation in most animals and in dioecious plants. The sex of an individual may be determined genetically or by the environment. Individuals are *asexual* if the organs that produce gametes are missing or atrophied; *intersexual* if their sexual organs are intermediate or ambiguous; or *hermaphrodite* if they have both kinds of sexual organs, whether they are functional or not. Hermaphroditism is prevalent in many species of plants, whether both pollen and ova are produced in the same *hermaphrodite* flowers, or in separate flowers of the same plant (*monoecious* plants).

These concepts cannot be transferred to unicellular organisms without some modifications. The [haploid](#) cells that predominate in the life cycles of many unicellular eukaryotes can often be classified into different **mating types**, sometimes more than two, with the criterion that only those of



different mating types are able to fuse and form diploids. This situation is called *heterothallism*. The mating types of the lower eukaryotes are sometimes called sexes, even though cells of different mating types are not homologous to sperm and ova of the animals.

In the situation called *homothallism*, diploid cells may be formed by the fusion of two genetically identical haploid cells. This would happen if there were a single mating type or if haploid cells change their mating type when they find no partner of the opposite one, as occurs in the homothallic strains of the yeast *Saccharomyces*.

Asexual bacteria are unable to transfer or receive DNA via conjugation (see [Hfr'S And F-Primes; F Plasmid](#)).

## Asexual Reproduction

Asexual reproduction includes all forms of reproduction that do not involve the fusion of male and female gametes. The asexual reproduction of a cell consists in the formation of genetically identical daughter cells; in the eukaryotes, this implies nuclear division via mitosis. Multicellular organisms are formed through repeated cycles of asexual reproduction from an initial zygote (see [Development and Embryology](#)). The daughter cells need not be completely identical, because cells with the same genotypes may differ in the quantitative distribution of [transcription factors](#) (see [Differentiation](#)). The genotypes of the daughter cells may differ because of mutation, [mitotic recombination](#), chromosomal nondisjunction, and other infrequent events.

There are two kinds of *asexual* reproduction of multicellular organisms, namely, [parthenogenesis](#) and vegetative reproduction. Parthenogenesis imitates sexual reproduction, in that embryogenesis is apparently normal, but without participation of a male gamete. [Haploid](#) individuals are produced through the parthenogenetic development of a normal female gamete. **Diploid** individuals are produced through the fusion of the haploid nuclei of a female gamete and a neighboring haploid cell, such as the synergids of plants and the polar bodies of animals. Other diploid individuals are produced through the parthenogenetic development of a diploid cell after failure of meiosis. In vegetative reproduction, individuals are derived from one or many somatic cells. It is a common natural process in many plants and lower animals, and it is the only multiplication mechanism in some of them.

Artificial interventions permit new forms of asexual reproduction. These include (a) the production of haploid plants from the culture of anthers and (b) the production of diploid animals through the replacement of the nucleus of the ovum by a somatic nucleus. Vegetative reproduction and asexual parthenogenesis lead to the formation of *clones*, which are sets of individuals with identical genotype. Clonal propagation is usually carried out through vegetative reproduction from a single cell or a group of identical cells.

The [parasexual cycle](#) may be seen as a form of asexual reproduction in the lower eukaryotes and in cultured cells of multicellular organisms. Haploid cells, not necessarily of different **mating type**, fuse to form [heterokaryons](#), and the nuclei of the heterokaryon fuse to produce diploid nuclei; these revert to haploidy, not by meiosis, but by random chromosome losses during mitosis, accompanied by [mitotic recombination](#).

The sexual processes of bacteria imply the conjugation of two compatible cells and the transfer of part of the genetic information of one to the other (see [Hfr'S And F-Primes; F Plasmid](#)). Formally similar are asexual processes in which some genetic information from a cell is transferred to another via a **virus** or via naked DNA (see [Transformation](#)). These processes are not exclusive to the

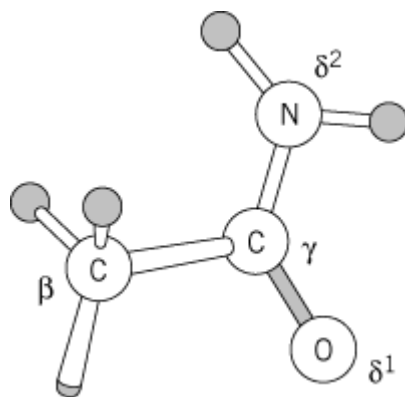
prokaryotes.

In evolutionary terms and in comparison with the various sexual processes, asexual reproduction has advantages of speed, economy, and the conservation of certain genotypes that may be lost or made less frequent after meiosis (see [Heterosis](#), and [Aneuploidy](#)). The main drawback of asexual reproduction is the limited creation of new genetic variation (see [Genetic Diversity](#)).

## Asparagine (Asn, N)

The [amino acid](#) asparagine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to six **codons**—CGU, CGC, CGA, CGG, AGA, and AGG—and represents approximately 4.4% of the residues of the proteins that have been characterized. The asparaginyl residue incorporated has a mass of 114.11 Da, a **van der Waals volume** of  $96 \text{ \AA}^3$ , and an [accessible surface](#) area of  $158 \text{ \AA}^2$ . Asn residues are those most variable during [divergent evolution](#); they are interchanged frequently in **homologous** proteins with [serine](#), [aspartic acid](#), [lysine](#), and [histidine](#) residues.

The side-chain of Asn residues is the same as that of [aspartic acid](#), except that the [carboxyl group](#) has been converted to the amide:



Both amino acids occur naturally and are incorporated directly into proteins during their biosynthesis; Asn residues do not arise from amidation of Asp in proteins. The amide side chain does not ionize and is not very reactive chemically. It is [polar](#), however, because it is both a [hydrogen bond](#) donor and acceptor. The amide group is labile at extremes of pH and at high temperatures, and Asn residues can **deamidate** to Asp. The Asn residue is especially labile at alkaline pH because its side chain is sterically suited to interact with the  $\text{—NH—}$  group of the following residue in the polypeptide chain, to form transiently a cyclic succinimidyl derivative (see [Deamidation](#)). This derivative can undergo racemization and hydrolysis to cleave the polypeptide chain or to produce a mixture of D and L isomers of Asp and isoAsp residues. In an isoAsp residue, the normal backbone peptide bond is through the side-chain carboxyl group rather than the usual  $\alpha$ -carboxyl. The deamidation reaction of Asn residues occurs 30–50 times more rapidly if the following residue is **Gly**, because the absence of a side chain there favors succinimide formation, and  $\text{—Asn—Gly—}$  sequences are especially prone to deamidation. The rate depends on the polypeptide conformation, however, because only some conformations permit succinimide formation. If the succinimide ring reacts with hydroxylamine instead of water, peptide cleavage results. Therefore, Asn—Gly peptide bonds are readily cleaved by incubation with hydroxylamine.

Asn residues are the site of [N-glycosylation](#) of proteins, one of the most prevalent [post-translational modifications](#). The Asn residue that is so glycosylated always occurs in a characteristic sequence: —Asn—Xaa—Ser—, —Asn—Xaa—, Thr—, or —Asn—Xaa—Cys—, where Xaa can be any residue except **Pro**, which also cannot immediately follow the tripeptide sequence.

#### Suggestions for Further Reading

S. J. Wearne and T. E. Creighton (1989) Effect of protein conformation on rate of deamidation: ribonuclease A, *Proteins Struct. Funct. Genet.* **5**, 8–12.

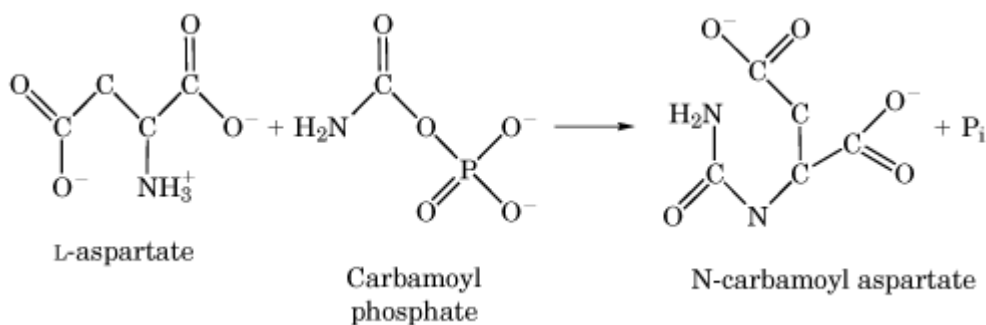
W. J. Chazin et al. (1989) Identification of an isoaspartyl linkage formed upon deamidation of bovine calbindin  $D_{9k}$  and structural characterization by 2D  $^1\text{H}$  NMR, *Biochemistry* **28**, 8646–8653.

P. Bornstein and G. Balian (1977) Cleavage at Asn-Gly bonds with hydroxylamine, *Meth. Enzymol.* **47**, 132–144.

F. Wold (1985) Reactions of the amide side-chains of glutamine and asparagine *in vivo*, *Trends Biochem. Sci.* **10**, 4–6.

### Aspartate Transcarbamoylase

Aspartate transcarbamoylase (carbamoylphosphate:L-aspartate carbamoyl-transferase, ATCase, EC 2.1.3.2) is a ubiquitous [enzyme](#) of [pyrimidine](#) biosynthesis in which it catalyzes the first unique step in the *de novo* synthetic pathway:

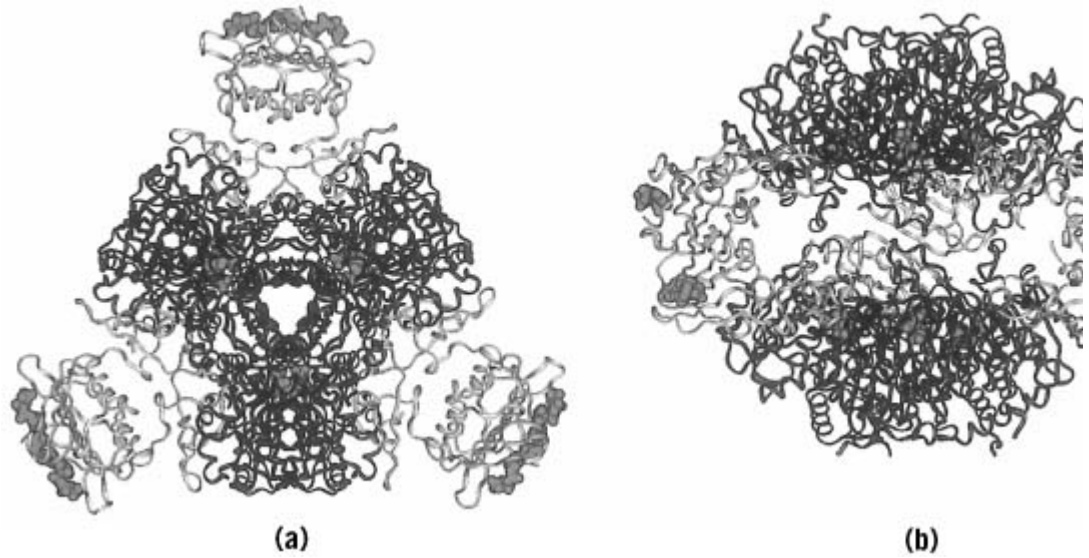


#### 1. The *Escherichia coli* Paradigm

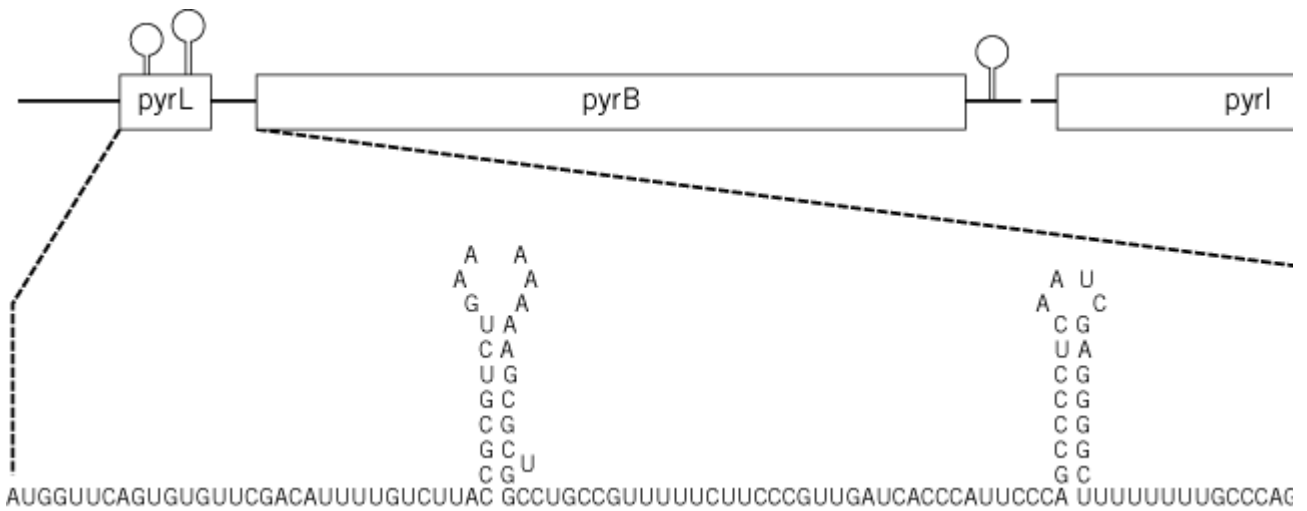
Although the enzyme exists in a variety of oligomeric and multifunctional complexes in various organisms, the most studied ATCase is that of *Escherichia coli*. It is a structurally complex, **allosterically** regulated enzyme containing two catalytic trimers associated with three regulatory dimers, 2(c<sub>3</sub>):3(r<sub>2</sub>) (Fig. 1). The catalytic polypeptide chain is encoded by the *pyrB* gene, whereas *pyrI* encodes the regulatory polypeptide. These genes are preceded by two tandem **promoters**, designated P<sub>1</sub> and P<sub>2</sub>, and an open reading frame encoding a 44-residue leader polypeptide, *pyrL*. Promoter P<sub>2</sub> has been identified as the physiologically significant promoter (see Fig. 2).

**Figure 1.** Structure of the intact *E. coli* ATCase viewed along the 3-fold axis (a) and along the approximate 2-fold axis (

located within the catalytic trimer. The allosteric sites are located within the regulatory dimers. The PDB file 8at1 was us



**Figure 2.** Genetic organization of the *E. coli pyrLBI* operon. The position of the three cistrons on the chromosome is shown. The positions of three hairpin structures adopted by the messenger RNA are indicated. The sequence around the first two is shown below. The first is the *pause hairpin*, flanked on each side by uridine-rich sequences. The second is the *terminator hairpin*, encoded by the atten



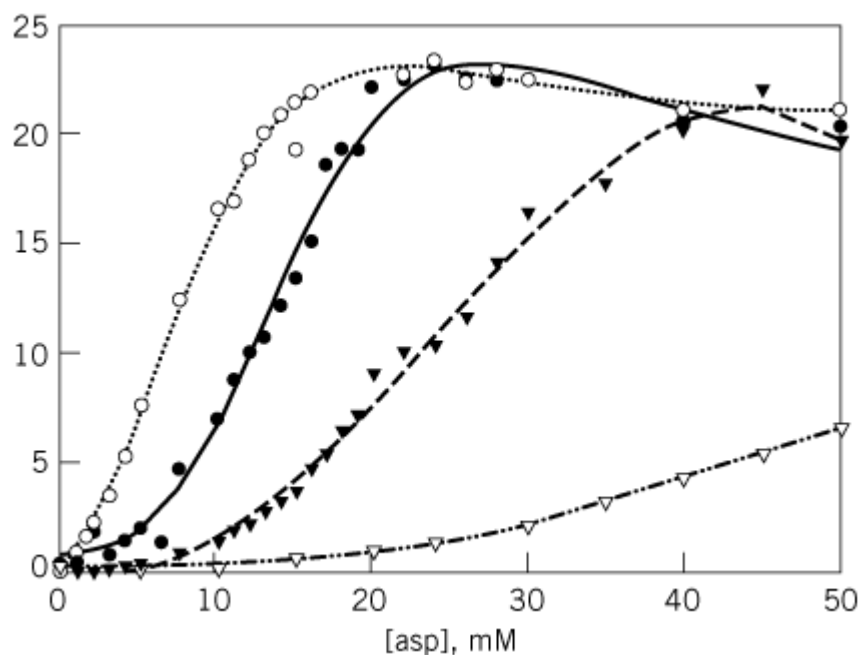
The *pyrLBI* operon is located on the linkage map of *E. coli* at 97 minutes and unlinked to any of the genes or small operons of the other six enzymes involved in pyrimidine biosynthesis. The expression of these genes is noncoordinately regulated by the intracellular levels of uridine or cytidine nucleotides, with the expression of the *pyrBI* operon negatively regulated over 300-fold. Most of this regulation (50-fold) occurs through a UTP-sensitive **attenuation** control mechanism, whereas attenuator-independent mechanisms are responsible for approximately a 6.5-fold range of regulation, including a pyrimidine-sensitive **transcriptional** initiation mechanism, and **stringent** control by ppGpp (1). According to the current model for attenuation control, transcriptional termination at the *pyrBI* attenuator (a **rho**-independent transcriptional terminator) is regulated by the relative rates of transcription and **translation** within the *pyrBI* leader region. Low intracellular levels of UTP cause **RNA polymerase** to pause at the uridine-rich region in the leader transcript (this pausing is enhanced by NusA, a general **transcription factor** that increases the efficiency of termination). This

allows time for a [ribosome](#) to initiate translation and catch up to the stalled RNA polymerase before it transcribes the attenuator region. As RNA polymerase eventually makes its way through the attenuator, the adjacent translating ribosome blocks formation of the *terminator hairpin*. This permits RNA polymerase to read through into the *pyrBI* structural genes, with translation terminating before the *pyrB* initiation codon. When intracellular levels of UTP are high, RNA polymerase does not pause during the transcription of *pyrL*. Without this pausing, the ribosomes cannot catch up to RNA polymerase before it transcribes the attenuator. This allows formation of the RNA hairpin, thus terminating transcription prior to the structural genes.

In addition to the genetic regulation of *pyrLBI*, the gene product ATCase also exerts significant control over the rate of pyrimidine biosynthesis (2). This is achieved by allosteric modification of the enzymatic activity in response to substrate concentration (**homotropic cooperativity**) and the nucleotide end products of both the purine and pyrimidine pathways (**heterotropic regulation**; see Tables 1 and 2). In the oligomeric complete enzyme, catalysis proceeds by a preferred order mechanism, with carbamoyl phosphate binding before aspartate and N-carbamoyl-L-aspartate leaving before inorganic phosphate. Homotropic cooperativity is induced by aspartate in the presence of a saturating concentration of carbamoyl phosphate and involves a structural transition of the enzyme consistent with a two-state, **concerted model**. As aspartate is bound, the enzyme shifts from the **T-state**, characterized by low activity and low affinity for substrates, to the **R-state** with high activity and high affinity for substrates. The kinetic consequence of this positive cooperativity is a sigmoidal substrate saturation curve (Fig. 3). Heterotropic activation is caused by the purine ATP, whereas the pyrimidine end products CTP and UTP **feedback inhibit** the enzyme. The pattern of inhibition exhibited by pyrimidine nucleotides is synergistic: CTP inhibits ATCase activity approximately 60% and UTP has minimal effect, while CTP in combination with UTP inhibits ATCase activity >95% (3). CTP, ATP, and UTP bind competitively, with different affinities, to a common allosteric site. CTP binds the most tightly with a pattern consistent with two classes of three sites each, whose **dissociation constants** differ by a factor of 20 ( $K_{d-CTP}$ : 5 to 20  $\mu\text{M}$ ). The binding of ATP follows a pattern similar to that of CTP with two classes of affinity sites, except that ATP binding is an order of magnitude weaker than that of CTP ( $K_{d-ATP}$ : 60 to 100  $\mu\text{M}$ ). The binding of UTP appears to be limited to three sites ( $K_{d-UTP}$ : 800  $\mu\text{M}$ ), although there may be a second class of sites too weak to be measured. Interestingly, the binding of CTP to three sites appears to enhance the binding of UTP to the remaining three sites almost 100-fold, resulting in a predicted  $K_d$  for UTP in the presence of CTP of 10  $\mu\text{M}$  (4). These biochemical binding characteristics are well suited to the physiological requirements of the pathway: Although the intracellular concentration of CTP (500  $\mu\text{M}$ ) and UTP (900  $\mu\text{M}$ ) is 3- to 6-fold lower than that of ATP (3 to 5  $\text{mM}$ ), their stronger binding can effectively displace ATP at the allosteric sites of the enzyme.

**Figure 3.** Substrate saturation profile of the *E. coli* ATCase. The velocity of the enzyme at varying concentrations of asp

presence of no allosteric effectors •, the activator ATP °, the inhibitor CTP ▼, and both pyrimidine end-products CT



**Physiological Significance of Catalytic Characteristics of the *E. coli* Aspartate Transcarbamoylase**

| Catalytic Characteristic  | Functional Consequence  |
|---|---|
| <ul style="list-style-type: none"> <li>• Ordered substrate binding (CP→Asp)</li> </ul>                    | <ul style="list-style-type: none"> <li>• Places catalytic emphasis on aspartate</li> </ul>  |
| <ul style="list-style-type: none"> <li>• Cooperative aspartate binding</li> </ul>                         | <ul style="list-style-type: none"> <li>• Allows large changes in catalytic activity with only small changes in physiological substrate concentrations</li> </ul>        |
| <ul style="list-style-type: none"> <li>• Distinct regulatory and catalytic subunits</li> </ul>            | <ul style="list-style-type: none"> <li>• Increased sophistication in the allosteric modulatory activity</li> </ul>  |
| <ul style="list-style-type: none"> <li>• Presence of shared active sites and catalytic trimers</li> </ul> | <ul style="list-style-type: none"> <li>• Facilitates modulation of catalytic activity with small conformational movements</li> </ul>                                    |
| <ul style="list-style-type: none"> <li>• T→R structural transition</li> </ul>                             | <ul style="list-style-type: none"> <li>• Provides for dramatic differences in aspartate binding and modulation of catalytic capacity at physiological levels</li> </ul> |

**Table 2. Physiological Significance of Allosteric Characteristics of the *E. coli* Aspartate Transcarbamoylase**

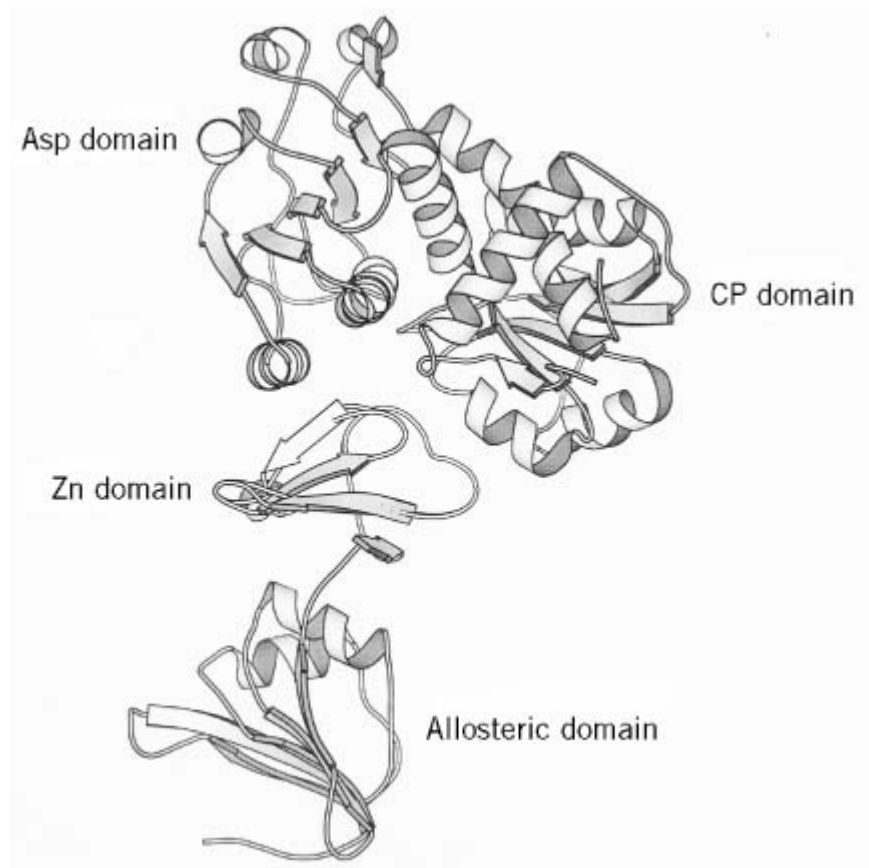
| Allosteric Characteristic  | Functional Consequence   |
|--|--|
| <ul style="list-style-type: none"> <li>• Allosteric inhibition by CTP</li> </ul>             | <ul style="list-style-type: none"> <li>• 50–70% inhibition of activity by nucleotide CTP balances new synthesis with utilization</li> </ul>                                |
| <ul style="list-style-type: none"> <li>• Synergistic inhibition by UTP and CTP</li> </ul>    | <ul style="list-style-type: none"> <li>• UTP and CTP are both end-products of the biosynthetic pathway and the synergism provides 90–95% inhibition of activity</li> </ul> |
| <ul style="list-style-type: none"> <li>• Competitive allosteric activation by ATP</li> </ul> | <ul style="list-style-type: none"> <li>• ATP competes with CTP and UTP for the same binding site, thus balancing intracellular purine and pyrimidine pools</li> </ul>      |
| <ul style="list-style-type: none"> <li>• Competitive allosteric inhibition by CTP</li> </ul> | <ul style="list-style-type: none"> <li>• Allows for displacement of ATP (activation)</li> </ul>  |

| CTP and CTP + UTP   | (CTP or CTP + UTP)   |
|---|--|
| <ul style="list-style-type: none"> <li>• Negative cooperativity in binding CTP</li> <li>• Communication of allosteric signals across protein–protein interfaces</li> <li>• Asymmetry of the holoenzyme</li> </ul> | <ul style="list-style-type: none"> <li>• Permits modulation of enzymatic activity of inhibitor concentration</li> <li>• Facilitates modulation of allosteric signal</li> <li>• Provides two distinct classes of allosteric binding sites on a single enzyme</li> </ul> |

---

Comparison of the [X-ray crystallography](#) structures of the enzyme has revealed that each of the catalytically independently folding structural **domains**: the carbamoyl phosphate (CP) domain and aspartate (Asp) domain. Likewise, the regulatory chains are composed of two structural domains: the zinc-binding (Zn) domain. Differences between the two states of the enzyme have been identified by comparison of the inhibitor (CTP) and the R-state structure binding a [bisubstrate analogue](#) [N-(phosphonacetyl)-L-aspartate]. Change in [quaternary structure](#) involves substantial conformational rearrangements as the catalytic trimers mutually reorient 10° around the 3-fold axis, while each regulatory dimer rotates 15° around the 2-fold axis. Upon CTP binding, the two domains of each catalytic chain (Asp and CP) undergo domain closure, whereas the regulatory chain (Allo and Zn) undergo domain separation. [Site-directed mutagenesis](#) studies have shown that the Asp domain is important for the formation of the high-affinity, high-activity R-state, which is required for the conformation needed for catalysis, and for homotropic cooperativity.

**Figure 4.** Domain organization of the *E. coli* ATCase illustrated with one catalytic chain (C1) and its associated regulatory chain composed of two independently folding domains: the aspartate-binding (Asp) domain and carbamoyl phosphate-binding (CP) domain. The regulatory chain is composed of the zinc-binding (Zn) domain and nucleotide-binding (allosteric) domain. This R1:C1 structure was generated using MOLSCRIPT programs.



In addition to these studies, well over 100 site- and region-specific mutations have been created in the studies, along with the many structural and biochemical analyses, have contributed to our understanding of the active-site and allosteric site locations, identifying interactions important to the stabilization of the transition state and the catalytically significant residues. More recently, site-directed mutagenesis studies have been directed at understanding the mechanism of allosteric regulation, separating the homotropic from the heterotropic effects, and CTP inhibition from CTP + UTP synergism. Although the mechanism of allosteric regulation is still unclear, current mutagenesis and structural studies are attempting to identify and differentiate between possible discrete pathways for each nucleotide signal, (2) complex and multiple interlocking pathways, or (3) energy changes.

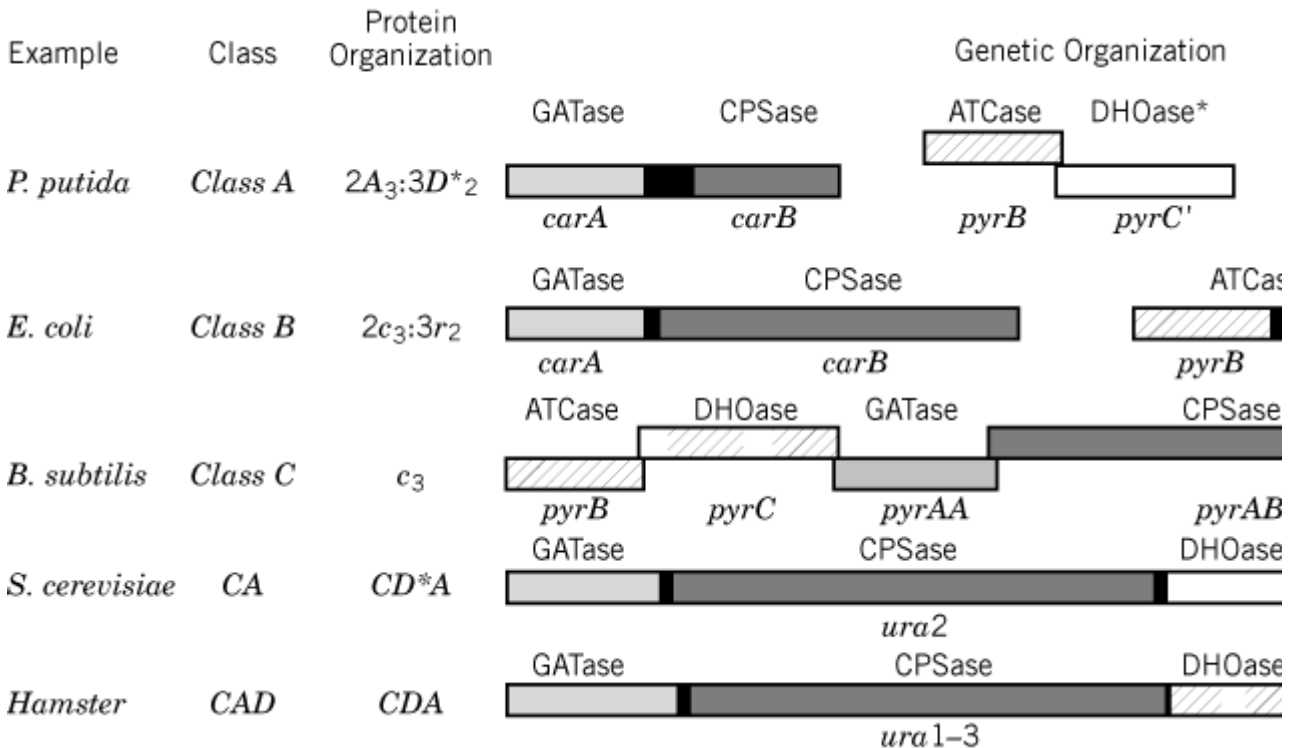
## 2. Diversity of Aspartate Transcarbamoylase

### 2.1. Prokaryotes

With the exception of some anaerobic protozoan parasites, all examined organisms are capable of *de novo* pyrimidine biosynthesis and have been found to possess the enzyme aspartate transcarbamoylase. The earliest classification of ATCases was by Jones and Jones in 1969 (16), based on enzyme size and response to nucleotides, partitioned the enzyme into three classes. Class A ATCases are the largest ATCases found in bacteria and were initially described for the pseudomonas ATCases of *E. coli* and other members of the *Enterobacteriaceae*; whereas (3) class C is the class which until recently was restricted to **Gram-positive** bacteria such as *Bacillus* (Fig. 5). The class A ATCases are the largest, and the intact enzyme exists as two catalytic trimers of 35 kDa in dodecameric association with two regulatory chains (2A<sub>3</sub>:3D<sub>2</sub>). In *P. putida*, this 45-kDa chain has been shown to be a nonfunctional homologue of the third enzyme in *de novo* pyrimidine biosynthesis immediately following ATCase in the biosynthetic pathway.



**Figure 5.** Organization of the first three enzymes and genes of pyrimidine biosynthesis: glutamine amidotransferase (GATase), DHOase. Although multiple examples are known for each category, the examples selected to represent each class are from *S. subtilis* (13), *S. cerevisiae* (9), and hamster (14). The genetic organization is represented according to the following: Enzyme designations above and gene designations below. Noncoding DNA that joins genes transcribed into polycistron segments. Overlapping genes are represented by overlapping boxes. Although all seven cistrons of *pyr* metabolism are present, the organization of only the first three is shown here.



The *Bacillus subtilis* enzyme at 100 kDa is a typical class C enzyme and corresponds in size, architecture, and catalytic trimer of the *E. coli* enzyme. Class C enzymes are composed of three identical polypeptide subunits of 34 kDa each and have no associations with other enzymes in the pathway. Furthermore, these enzymes lack a regulatory subunit nor the attendant allosteric regulation. As a consequence of the active site being shared among the three catalytic chains, the class C homotrimer may be considered the catalytic core of the bacterial ATCase. The *pyr* genes are scattered around the chromosome and are not coordinately regulated, all seven *B. subtilis* *pyr* genes are transcribed on a single [messenger RNA](#). The 3' ends of the reading frames overlap with downstream open reading frames for all **cistrons** in the cluster except *pyrB* and the preceding ORF1. The expression of *pyr* genes is repressed by pyrimidines, and expression of aspartate transcarbamoylase has been shown to decline during **sporulation**. The mechanisms for the nutritional and developmental regulation of *pyr* gene expression have not yet been elucidated.

Although the class B enzymes appeared to be restricted to the *Enterobacteriaceae*, recent studies suggest that other *Enterobacteriaceae* (7) may also contain very similar  $c_6:r_6$  enzymes. The *E. coli* ATCase provides the textbook example of allosteric inhibition by pathway end-products CTP and UTP, and activation by the end-product of the parallel pathway, ATP. However, the various ATCases from different tribes of the *Enterobacteriaceae* display diverse regulatory mechanisms of inhibition or activation. For example, while CTP serves as an allosteric feedback inhibitor in the *E. coli* ATCase, in other *Enterobacteriaceae* (Table 3). Some of these diverge from the basic regulatory paradigm of inhibition by pathway end-products as a mechanism to conserve energy and avoid the production of unneeded products. Nonetheless, the discovery of CTP + UTP synergistic inhibition provided a new paradigm for allosterically regulated enzymes: Since the concentration of ATP in actively growing cells is 2 orders of magnitude higher than that of CTP or UTP, the combination of CTP and UTP always serves as an effective antagonist to the activation by ATP unless CTP or CTP + UTP are present to reduce the activated enzyme.

**Table 3. Classification and Allosteric Characteristics of Bacterial Aspartate Tra**

| ATCase Class Characteristics | Bacterial Species <sup>a</sup>  | Alloster                                      |
|------------------------------|---|---|
| ATCase A                     | <i>Pseudomonas fluorescens</i><br><i>Acinetobacter calcoaceticus</i><br><i>Azomonas agilis</i><br><i>Azotobacter vinelandii</i> | No activation<br>UMP Inhibition               |
| ATCase B1 (I) <sup>b</sup>   | <i>Escherichia coli</i>   | ATP activation                                |
| ATCase B2 (IV)               | <i>Salmonella typhimurium</i><br><i>Yersinia intermedia</i>   | CTP, CTP + UTP<br>ATP activation              |
| ATCase B3 (V)                | <i>Erwinia carnegiana</i><br><i>Erwinia herbicola</i>   | No activation<br>sCTP <sup>c</sup> , CTP + UT |
| ATCase B4 (IV)               | <i>Yersinia enterocolitica</i>  | ATP activation<br>No inhibition               |
| ATCase B5 (IV)               | <i>Yersinia kristensenii</i><br><i>Yersinia frederiksenii</i>   | No activation<br>No inhibition                |
| ATCase B6 (II)               | <i>Aeromonas hydrophila</i><br><i>Serratia marcescens</i>   | CTP, ATP activati<br>CTP + UTP antagi         |
| ATCase B7 (III)              | <i>Proteus vulgaris</i>   | CTP, ATP activati<br>CTP + UTP Inhibi         |
| ATCase C                     | <i>Bacillus subtilis</i><br><i>Streptococcus faecalis</i><br><i>Staphylococcus epidermidis</i>                                  | None  |

<sup>a</sup> Class A and C examples are from (15).

<sup>b</sup> Tribal classifications (given in parentheses) are according to Bergey's *Manual of Determinative Ba* subgroups of the ATCase class B enzymes.

<sup>c</sup> sCTP indicates only slight inhibition by CTP (>20%).

Due to the oligomeric nature of the class B enzymes, it has been possible to form hybrids by combin enzyme with the regulatory subunits of different enzymes that have diverged allosteric patterns. The have demonstrated that the regulatory dimer determines the nature of the allosteric control of ATCas chimeric proteins have been constructed by intragenic fusion of the CP domain with the Asp domain Allo domain) forming several novel protein structures. In one instance, this type of [protein engineeri](#) chimeric enzyme by fusion of the domains of the hamster ATCase cDNA and the bacterial *pyrB* gen active ATCase trimers that, although unstable, could marginally satisfy physiological requirements.

**Table 4. The Regulatory Chain of Class B ATCases Dictates the Allosteric**

| Source of subunit    |                      | Response to Effectors <sup>a</sup> |     |
|----------------------|----------------------|------------------------------------|-----|
| Catalytic Subunit    | Regulatory Subunit   | ATP                                | CTP |
| <i>E. coli</i>       | <i>E. coli</i>       | +                                  | -   |
|                      | <i>S. marcescens</i> | +                                  | +   |
|                      | <i>P. vulgaris</i>   | +                                  | +   |
| <i>S. marcescens</i> | <i>E. coli</i>       | +                                  | -   |
|                      | <i>S. marcescens</i> | +                                  | +   |
|                      | <i>P. vulgaris</i>   | +                                  | +   |
| <i>P. vulgaris</i>   | <i>E. coli</i>       | +                                  | -   |
|                      | <i>S. marcescens</i> | +                                  | +   |
|                      | <i>P. vulgaris</i>   | +                                  | +   |

<sup>a</sup> + = activation; - = inhibition.

## 2.2. Eukaryotes

The ability to form chimeric proteins opens the possibility that intragenic fusions could provide a means for the development of new proteins by the fusion of domain modules. Among the biosynthetic pathways of examples of single polypeptides that carry multiple enzymatic activities. Eukaryotic ATCases provide examples of a multienzymatic protein. In lower eukaryotes, such as **yeast**, the first enzyme in the pathway, *carbamoyltransferase* (CPSase), and ATCase are physically linked, forming the CD<sup>\*</sup>A protein fusion complex. In *S. cerevisiae* and *S. pombe*, the enzyme architecture includes four domains; three functional domains carry out *carbamoyltransferase* (GLNase), CPSase, ATCase activities, and one dihydrorotase-like (D<sup>\*</sup>; DHOase) carries out the third enzymatic activity of the pyrimidine biosynthetic pathway just following ATCase. This structure is different from the architecture previously discussed for the class A prokaryotic enzymes, which possess a domain for DHOase. The regulation of this complex includes the feedback inhibition of both CPSase and ATCase, although the carbamoylphosphate produced by GLNase/CPSase is tightly channeled to ATCase and a domain for the GLNase/CPSase alone should be sufficient to regulate pyrimidine metabolism.

In higher eukaryotes, the CA complex is associated with a functional DHOase domain producing a domain for DHOase (physically arranged in a CDA sequence). This complex provides the central metabolic control for pyrimidine biosynthesis. CPSase subjected to allosteric inhibition by UTP and activation by 5-phosphoribosyl 1-pyrophosphate. In higher eukaryotes, regulatory and catalytic functions involve a single polypeptide chain, the multifunctional complex monomeric structure. A series of **proteolytic** and genetic truncation studies over the last 15 years have provided the structure of CAD is well-defined and simple: Each enzymatic domain is separated from the others by a polypeptide linker, and each domain can function in the absence of the other activities.

There are a number of proposals regarding the evolutionary role of the multienzymatic architecture. One proposal is that channeling occurs between CPSase and ATCase of the bifunctional CD<sup>\*</sup>A complex. Channeling of successive enzymatic activities are carried out on the same protein complex and could conceivably provide an advantage by limiting the loss of intermediate products. However, this is not always the case, as the complex freely releases the products of CPSase and ATCase, whereas ATCase and DHOase can readily utilize their substrates outside the enzymatic complex. In the case of CAD and similar enzymatic complexes, coordinate gene expression provides an alternative regulatory advantage for the evolution and maintenance of large enzymatic complexes. These complexes would provide for a smaller number of independent genes to regulate and simplify the coordination of subcellular localization of multiple enzymatic activities.

In summary, the *de novo* biosynthetic pathway involves the set of reactions that supplies UMP from metabolic pathways: aspartate, glutamine, ATP, PRPP and carbon dioxide. The *de novo* pyrimidine enzymatic steps from carbamoyl phosphate synthetase (CPSase) to orotidylate decarboxylase (OMP) evolution has developed a variety of regulatory controls and genetic organizations. Independent of it ATCase provides an important regulatory component of *de novo* pyrimidine biosynthesis in all free-living organisms. Maintaining homeostasis in intracellular nucleotide pools is an essential consideration for metabolism. Enzymological regulation mechanisms are critical for balancing the nucleotide precursors of DNA/RNA.

### Bibliography

1. C. Liu, J. P. Donahue, L. S. Heath, and C. L. Turnbough Jr. (1993) *J. Bacteriol.* **175**, 2363–2369
2. J. C. Gerhart and A. B. Pardee (1962) *J. Biol. Chem.* **237**, 891–896.
3. J. R. Wild, S. J. Loughrey, and T. C. Corder (1989) *Proc. Natl. Acad. Sci. USA* **86**, 52–56.
4. P. England and G. Hervé (1993) *Biochemistry* **31**, 9725–9732.
5. R. P. Kosman, J. E. Gouaux, and W. N. Lipscomb (1993) *Proteins* **15**, 147–176.
6. H. Ke, W. N. Lipscomb, Y. Cho, and R. B. Honzatko (1988) *J. Mol. Biol.* **204**, 725–747.
7. C. Purcarea, G. Hervé, M. M. Ladjimi, and R. Cunin (1997) *J. Bacteriol.* **179**, 4143–4157.
8. M. Lollier et al. (1995) *Curr. Genet.* **28**, 138–149.
9. L. Jaquet et al. (1995) *J. Mol. Biol.* **248**, 639–652.
10. M. J. Schurr et al. (1995) *J. Bacteriol.* **177**, 1751–1759.
11. T. A. Hoover et al. (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2462–2466.
12. C. D. Pauza, M. J. Karels, M. Navre, and H. K. Schachman (1987) *Proc. Natl. Acad. Sci. USA* **77**, 1113–1117.
13. C. L. Quinn, B. T. Stephenson, and R. L. Switzer (1991) *J. Biol. Chem.* **266**, 9113–9127.
14. J. P. Simmer et al. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4382–4386.
15. M. J. Kenny, D. McPhail, and M. Shepherdson (1996) *Microbiology* **142**, 1873–1879.
16. M. R. Bethell and M. E. Jones (1969) *Arch. Biochem. Biophys.* **134**, 352–365.

### Suggestions for Further Reading

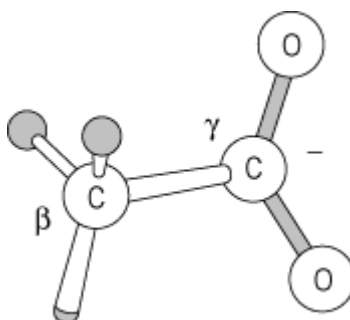
17. J. N. Davidson et al. (1993) The evolutionary history of the first three enzymes in pyrimidine biosynthesis. *J. Mol. Biol.* **231**, 157–164.
18. E. R. Kantrowitz and W. N. Lipscomb (1990) *Escherichia coli* aspartate transcarbamoylase: The concerted allosteric transition. *Trends Biochem. Sci.* **15**, 53–59.
19. W. N. Lipscomb (1996) Aspartate transcarbamoylase from *Escherichia coli*: Activity and regulation. *Enzymes* **19**, 131–151.
20. H. K. Schachman (1993) Aspartate transcarbamoylase. *Curr. Opin. Struct. Biol.* **3**, 960–967.
21. J. R. Wild and M. E. Wales (1990) Molecular evolution and genetic engineering of protein domains. Aspartate transcarbamoylase. *Ann. Rev. Microbiol.* **44**, 93–118.

## Aspartic Acid (Asp, D)

The [amino acid](#) aspartic acid is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to only two **codons**—GAU and GAG—and represents approximately 5.3% of the residues of the proteins that have been characterized. The aspartyl residue incorporated has a

mass of 115.09 Da, a **van der Waals volume** of  $91 \text{ \AA}^3$ , and an [accessible surface](#) area of  $151 \text{ \AA}^2$ . Asp residues are frequently changed during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [asparagine](#) and [glutamic acid](#) residues.

The side chain of Asp residues is dominated by its [carboxyl group](#):



This carboxyl group is normally no more reactive than are those of corresponding organic molecules, such as acetic acid. Its intrinsic  $pK_a$  value is close to 3.9, so Asp residues are ionized and very polar under physiological conditions; consequently, very few Asp residues are buried in folded protein structures, and nearly all have at least the carboxyl group on the surface. The  $pK_a$  can be shifted in folded proteins, however, and either the ionized or nonionized form can be used in the protein's function. For example, [carboxyl proteinases](#) have one active-site carboxyl group function in the ionized form, and another when nonionized. Asp carboxyl groups have a weak intrinsic affinity for  $\text{Ca}^{2+}$  ions, and they are used in many [calcium-binding proteins](#).

Asp residues differ from Glu only in having one methylene group, rather than two, so it might be thought that they would be very similar chemically and functionally in proteins, but this is not so. The slight difference in length of the side chains causes them to have different tendencies in their chemical interactions with the peptide backbone, so they have markedly different effects on the conformation and chemical reactivity of the peptide backbone. For example, Asp residues favor the **alpha-helical** conformation much less than Glu residues. In folded [protein structures](#), Asp residues occur most frequently in reverse [turns](#), whereas Glu residues are most frequently found in  $\alpha$ -helices. The polypeptide chain can be cleaved relatively easily at Asp residues, because the side-chain carboxyl group participates in the reaction. —Asp—Pro— peptide bonds are especially labile in acid, because the carboxyl group interacts with the unique tertiary N atom of the Pro residue.

#### Suggestions for Further Reading

M. Landon (1977) Cleavage at aspartyl-prolyl bonds, *Meth. Enzymol.* **47**, 145–149.

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

### Aspartyl Proteinase Inhibitors, Protein

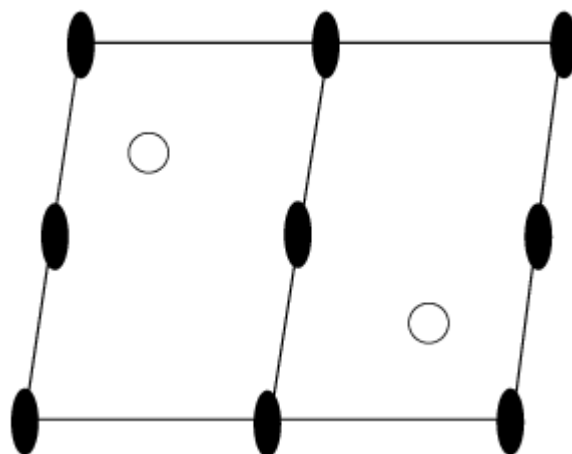
The synthetic peptide aspartyl proteinase inhibitors [see [Proteinase Inhibitors, Protein](#)] are among the most intensively designed molecules. The principal recent reason for this is that HIV proteinase inhibitors turned out to be successful drugs for fighting AIDS. In contrast, the protein inhibitors of aspartyl proteinases are very little studied. It is difficult to believe that this neglect is a result of the

rare occurrence of such inhibitors in nature. The activation of zymogens of aspartyl proteinases (pepsinogen) has been a problem of long-standing interest. The propeptide is clearly a pepsin inhibitor. Much more potent and typical is the pepsin inhibitor from the roundworm (*Ascaris lumbricoides*). The determination of the three-dimensional structure of its complex with pepsin is now in progress and may awaken this long-dormant field. A large number of reports deal with cathepsin D inhibitors classified as members of the soybean trypsin inhibitor (Kunitz), or STI family, [see [Soybean Trypsin Inhibitor \(Kunitz\), STI](#), (Kunitz)]. Many of the cathepsin D inhibitors are also said to inhibit trypsin, but the details of their interaction either with cathepsin D or with trypsin have been little studied.

## Asymmetric Unit

The asymmetric unit is the fundamental structure that is repeated in crystals of macromolecules (see [Crystallography](#); [X-Ray Crystallography](#)). Depending on the symmetry of the crystallographic lattice, each particle is present in multiple copies in the [unit cell](#) of the crystal. Only if symmetry is absent ([space group P1](#)) is a single copy present in the unit cell. In this case, the entire unit cell is an asymmetric unit. But if the unit cell has symmetry elements, it can be divided into identical parts related by this symmetry. Such a part, called the asymmetric unit, does not contain any crystallographic symmetry element. In Fig. 1, a projection of the unit cell in space group P2 is drawn. The twofold axes are perpendicular to the plane of the drawing. This unit cell has two asymmetric units. If the structure of the asymmetric unit is known, the structure of the entire crystal can be reproduced by applying the crystallographic symmetry.

**Figure 1.** Projection of a unit cell in space group P2 that has twofold axes (black ellipses) perpendicular to the plane of the drawing. The circles are particles related by the symmetry. The symmetry divides the unit cell into two exactly equal parts. Therefore, this cell has two asymmetric units.



## Bibliography

1. International Union of Crystallography (1992) *International Tables for Crystallography*, Vol. A (T. Hahn, ed.), Kluwer Academic, Dordrecht, Boston, London.

## Suggestions for Further Reading

2. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists*, VCH, New York, Weinheim, Cambridge.
3. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Atomic Force Microscopy

Binnig, Quate, and Gerber developed a [scanning probe technique](#) known as atomic force microscopy (AFM) and as *scanning force microscopy*. Unlike its predecessor, scanning tunneling microscope (STM), AFM also images nonconducting samples, such as biological specimens, in a liquid environment at molecular and even atomic resolution. Currently, structural information about biological materials at the molecular level is obtained from other **microscopic** techniques, including [electron microscopy](#), [electron crystallography](#), [X-ray crystallography](#), nuclear magnetic resonance (NMR) spectroscopy, and [vibrational spectroscopy](#). These techniques require extensive sample preparation and unfavorable operating environments, and they are unsuitable for providing real-time functional information. Molecular function is studied by various molecular biological, biochemical, and electrophysiological techniques, but it is difficult to combine both structural and functional studies in one technique. Moreover, these techniques provide incomplete information about the surfaces of biological macromolecules, the very sites of their interactions with other molecules. In contrast, AFM images the surfaces of biological specimens, where most of the regulatory biochemical and other signals are directed. Other microscopic techniques also image surfaces, for example, the **scanning electron microscope** (SEM), but AFM images living cells and molecules in a liquid environment at comparable and often greater resolution.

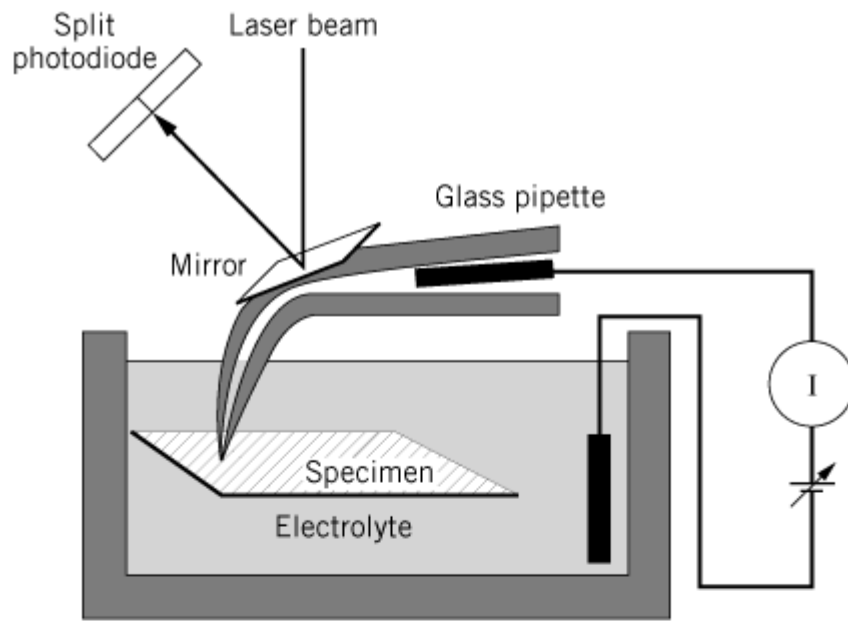
### 1. Principle of Operation

AFM is based on the general physical principle that the interactive force between two bodies is inversely proportional to some power of the distance separating them and on the physicochemical natures of the interacting bodies. A tip that is sharp on the molecular scale is attached to a cantilevered spring: As it is moved across the surface of the specimen, it is deflected by the interactive forces between the atoms of the tip and those of the specimen (see Fig. 1 of [Scanning Probe Techniques](#)). Because the spring constant of the commonly used cantilevers ( $10^{-1}$  to  $10^{-2}$  N/m) is much smaller than the intermolecular vibrational spring constant of the atoms in the specimen (10 N/m), the cantilever senses exquisitely small forces exerted by the individual sample atoms. The tip's deflection is a measure of the forces sensed by the cantilever, which are transduced to generate molecular images.

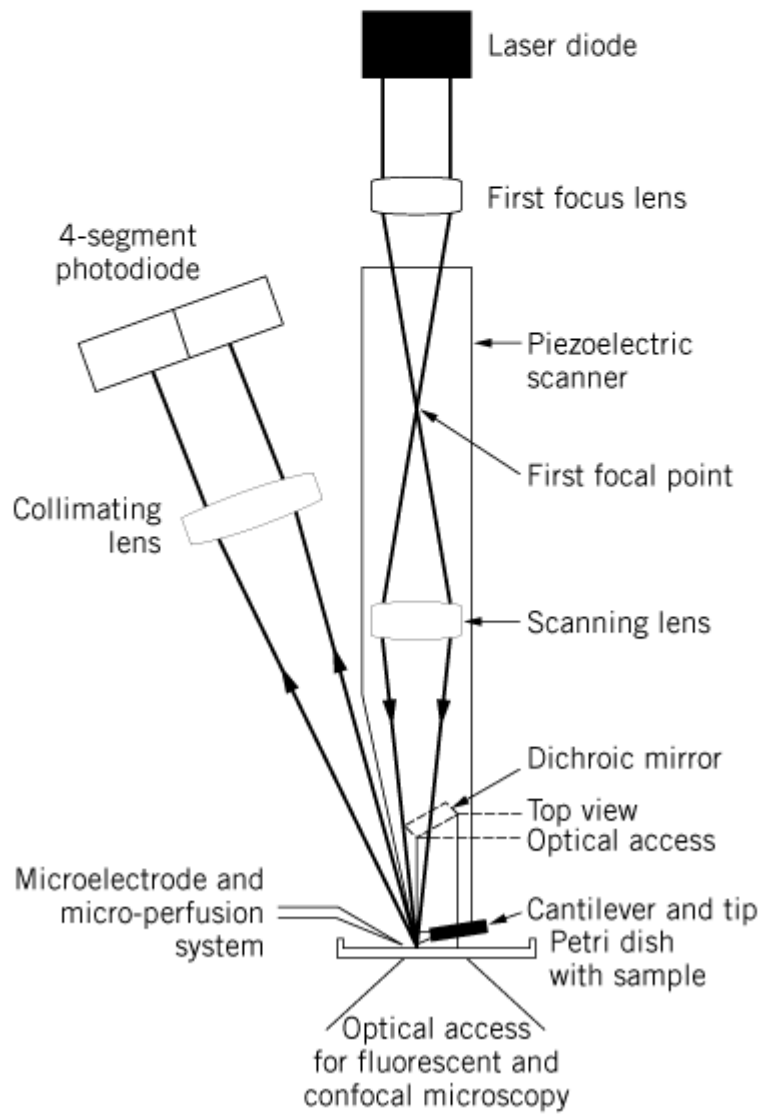
In practice, a microfabricated cantilevered tip is pressed against a sample surface by a small tracking (loading) force. The tip is raster-scanned in the x-y plane over the specimen by moving the sample beneath the tip or by moving the tip over the sample. The sample's vertical position (z) is also monitored. The three planes of movement are controlled by a piezoelectric xyz scanner, and the information about the three coordinates is used to create the image (Fig. 1). The cantilevered tip is brought sufficiently close to touch the sample (known as “contact” mode), or it oscillates at a finite distance (>a few nanometers) from the surface (“noncontact” mode). A noncontact mode microscope has the advantage that it does not perturb the sample, but the lateral resolution in these microscopes is poor, and hence they are not commonly used for biological imaging.

**Figure 1.** Schematic illustration of the operating principle of multimodal atomic force microscopes. **(a)** Schematic of the combined scanning ion conductance and atomic force microscope. A pipette serves as the probe: A laser beam reflecting from a mirror glued to the back of the pipette provides the deflection signal for the topographic image as the pipette is moved across the surface. The electrode has a nanometer-sized hole. Electrodes within the pipette and in the bath measure electrical currents (R. A. Proksch, R. Lal, P. K. Hansma, G. Morse, and G. Stucky (1996) *Biophys. J.* **71**, 2155–2157). **(b)** A combined light and atomic force microscope. The first focal point is located inside the upper portion of the piezoelectric scanner. After the positions of the lenses are adjusted, the scanning focused spot accurately tracks the cantilever, and the zero-deflection signal from the four-segment photodiode is independent of the position within the scan area. One of the key advantages of this new AFM is that there is optical access to the sample from above and below. Thus, the new AFM can be combined with an optical microscope of high numerical aperture. The other key advantage is that the sample is stationary during scanning and can be large, so that techniques for on-line perturbations and recordings are easily incorporated. For details, see P. K. Hansma, B. Drake, D. Grigg, C. B. Prater, F. Yasher, G. Gurley, V. Elings, S. Feinstein, and R. Lal (1994) *J. Appl. Phys.* **76**, 796–799; R. Lal and S. A. John (1994) *Am. J. Physiol.* **256**, C1–C21.





(a)



(b)

The deflective force is translated into a detectable signal in several ways. The most common is by an optical lever system (Fig. 1). A monochromatic laser beam reflects from the upper face of the cantilever, the angular direction of which changes as the cantilevered tip undergoes deflections. The reflected beams are captured and converted into electrical signals by position-sensitive multisegmented photodetectors. Such an optical lever amplifies the cantilever's deflection as much as a thousandfold, so deflections even less than a nanometer are measured.

### 1.1. Modes of Operation

Using an appropriate feedback system, the cantilever's deflection is kept constant or left to respond freely to the sensed forces. In the constant deflection mode (also called "constant force mode"), the feedback loop changes the height of the sample (to maintain the constant deflection) by adjusting the voltage applied to the z portion of the xyz piezoelectric scanner. The amount of z change corresponds to the sample's topological height at each point in the x-y raster. Combining the information from the three coordinates generates the 3-D image.

In the variable deflection mode ("constant height mode"), the feedback loop is open so that the cantilever deflection is proportional to the change in the tip-sample interaction, that is the force sensed by the cantilever. The surface image is constructed from the deflection information. It is called the constant height mode because the z component of the piezoelectric scanner does not change appreciably. This is usually unsuitable for a sample with large surface corrugation (e.g., cells), because the force fluctuations, and thus the cantilever's deflections, are enormous and often result in disengagement of the tip from the sample.

"Error mode" imaging relies on the imperfection in the feedback loop to operate in the constant-deflection mode. The error signal is amplified to yield contour information in the z plane. In the error mode, the feedback loop gathers high-frequency information that is normally not acquired in the constant deflection mode. This high-frequency information provides details of sharp contour changes (edges) in the sample. Measurements of actual height in error mode imaging are not accurate, however, in contrast to the other modes of operation. The main advantage of error mode operation is that imaging occurs without exerting high forces on the sample.

In "tapping mode" imaging, the cantilever is oscillated at very high frequency, normally near its resonance frequency, as it scans the sample. As the tip approaches the sample surface, its oscillatory amplitude decreases because of energy loss when the tip "taps" the surface. The amplitude of the cantilever oscillation is detected and used by the feedback system to adjust the tip-sample distance for constant amplitude. This ensures a much shorter tip-sample contact time, and smaller lateral forces are exerted on the cantilever. In this way, this mode has been successfully used for imaging delicate and individual macromolecules. The disadvantage of this mode of operation is that the vertical imaging force can be large, thus increasing the possibility of sample damage.

### 1.2. Sample Preparation

AFM is used to image specimens in aqueous, semiaqueous, or dry conditions. The imaging condition is normally chosen to maintain the specimen in as near a lifelike condition as possible. Where resolution takes priority over the physiological condition, however, the investigator is not as constrained. Imaging conditions also influence the choice of substrate, the stability of the specimen with respect to the tip interaction, and the preservation of the specimen with respect to its physiological or biochemical functions. At present, investigators rely on an empirical approach to find a suitable method, and probably will for some time to come. When it "works", the search stops for the "ideal support" or buffer.

The physicochemical characteristics of the sample determine or suggest ways under which it can be imaged. Problems encountered are as simple as getting the sample to attach to the support. Techniques for sample support include drying down of samples and adsorption to specially prepared surfaces. For example, imaging of **plasmid** DNA is vastly improved in to both resolution and

consistency, under propanol, which increases the humidity and produces a more hydrated condition. This allows reducing the tip-tracking force exerted on the sample to  $<1$  nN, thus decreasing sample deformation.

The interaction of the sample with the support determines the magnitude of the imaging force. If the sample does not adhere tightly to the support, low tracking forces must be used for imaging, or the tip literally sweeps the sample from the support. Originally, graphite (**hydrophobic** and uncharged), mica ([hydrophilic](#) and negatively charged), and glass (usually negatively-charged) were the supports most routinely used. These supports can also be modified chemically to adjust their [hydrophathy](#), charge density, and polarity. Today the repertoire has expanded greatly, and examples include gold treated with a variety of agents for DNA imaging. Gold supports maintained under potential (voltage) control have been used for DNA imaging with the **scanning transmission microscope**, and they may also prove useful for AFM.

The specimen support can also be modified or coated chemically so that it acts as a **ligand** for the sample and thus orients the specimen in a defined way. It is also possible to use artificial systems to generate constraints where there were none before. Examples include imaging isolated **cholera toxin** molecules incorporated into synthetic phospholipid bilayers, followed by covalent [cross-linking](#) or imaging the [vaccinia virus](#) protrusion out of living cells held by a suction pipette.

### 1.3. Forces in AFM

Interactive forces that deflect the cantilevered-tip are attractive or repulsive, and they vary depending on the mode of operation and the conditions used for imaging.

In contact mode imaging, the tip is deflected mainly by the repulsive forces from the overlapping electron orbitals of the atoms of the tip and sample. The dominant attractive force is a [van der Waals interaction](#) due primarily to the nonlocalized dipole–dipole interactions among atoms of the tip and specimen. When imaging in air, (attractive) surface tension is also present because of adsorbed [water](#) layers. For imaging in fluids, [electrostatic interactions](#) between charges on the specimen and the tip (occurring either naturally or induced by polarization), osmotic pressure due to charge movements and rearrangements, and structural forces due to [hydration](#), solvation, or adhesion enable a reduction in the net imaging force although both the meniscus and surface tension forces are abolished.

The interplay of local forces determines the stability of the specimen and the resolution.

Theoretically the force should be  $\leq 10^{-10}$  N for nonperturbed biological imaging. The sensitivity of AFM is sufficient to record small interactive forces, including the breaking of [hydrogen bonds](#). Imaging in contact mode under liquid, but with a net attractive rather than repulsive force, increases the resolution significantly and has produced true atomic resolution, even with an imaging force of  $10^{-11}$  N. As explained below, however, successful imaging of cells, membranes, and isolated proteins has been obtained with forces as large as  $10^{-7}$  N.

In principle, any movement of the tip caused by its interaction with the sample in the x, y, or z directions provides information about the specimen's topography. To date, most information has been obtained from z deflections. Improvements in hardware and software have, however, allowed recording movements in the zy or zx planes and measurement of lateral forces for image generation. The contribution of lateral forces to image contrast generation can be substantial.

### 1.4. Resolution in AFM

#### 1.4.1. Spatial Resolution

The limit of spatial resolution for AFM is not well defined because, unlike conventional microscopies, the images are formed by reconstructing the contours of interacting forces between the specimen and tip. The operating resolution in AFM is defined as the minimum size of two adjacent features that can be distinguished clearly. By selecting a small scan size and suitable operating conditions, one can distinguish two structures that are less than a nanometer apart. Image processing

tools used to define resolution in [X-ray crystallography](#) and [electron microscopy](#) studies may not give correct results for AFM. The operating resolution can be divided into three categories:

1. Instrumental resolution: The lateral resolution is about 1 Å and is determined by the limitations of the hardware. The vertical resolution is 0.1 Å, and hence molecular perturbations on a sample surface can be imaged.
2. Target resolution: The lateral resolution achieved depends on the characteristics of the tip, the operating environment, and the nature of the specimen. For crystalline solid specimens and many inorganic materials, atomic resolution of 1 to 2 Å has been achieved.
3. Resolution in biological specimens: The nature of the biological samples and their preparation play a key role in determining the resolution limits. For the surface of a living cell, the resolution is relatively poor (~10 nm) but greater than that by light [microscopy](#) and comparable to that by [scanning electron microscopy](#). In a biological specimen whose the density of particles is high and mobility is limited (e.g., proteins in a membrane), the resolution is comparable to that of a crystalline specimen.

#### 1.4.2. Temporal Resolution

Temporal resolution is limited by the maximum speed at which a specimen can be scanned and still have the tip accurately track surface features. Preliminary studies suggest that the scan speed should be  $\leq 2.2 \mu\text{m/s}$  for 1 nm spatial resolution on soft and deformable biological materials imaged in aqueous solution. Thus, membrane macromolecules whose dimensions are  $10 \text{ nm} \times 10 \text{ nm}$  (such as **channels** and **receptors**) divided into  $10 \times 10$  pixels (with pixel size  $\sim 1 \text{ nm}$ ) require about 45 to 50 ms to image. However, if only a single line is scanned, the image can be repeated every 4 to 5 ms. Measuring at a single point, rather than scanning, increases the temporal resolution significantly, and hence it is possible to obtain spatial information at very short time intervals. The temporal resolution also depends on the mode of operation (constant deflection or constant height mode), operating environment (solvents, pH, viscosity, elasticity), and the nature of the interactions between tip and sample. In the constant-height mode, the scan speed is limited by the speed with which the deflection of the cantilever changes in reacting to surface features. In the constant deflection mode, the scan speed is limited by how the speed with which the piezo scanner changes its z component. There is ultimately a limit to the temporal resolution imposed by the low-pass filter used to eliminate sampling noise. These filters typically have a cutoff frequency of  $\sim 15 \text{ kHz}$ , corresponding to a time resolution of  $67 \mu\text{s}$ .

Temporal resolution also depends on the material being imaged. Individual molecules at molecular resolution require faster scan speeds than cells at lower resolution. Molecular movements of biological macromolecules can be correlated with their lateral [diffusion](#) constant. Lipids in a bilayer have a typical diffusion constant of  $10^{-8} \text{ cm}^2/\text{s}$ , corresponding to a mean velocity of  $\sim 2 \mu\text{m/s}$ . The proteins embedded in natural biomembranes have diffusion constants many orders of magnitude lower (e.g., the [acetylcholine receptor](#) in myoblast patches has a diffusion constant  $< 3.0 \times 10^{-12} \text{ cm}^2/\text{s}$ ). Thus it is quite possible to obtain images at molecular resolution of proteins and other macromolecules that are properly anchored in a lipid bilayer or immobilized on a substrate.

#### 1.5. Identity of Imaged Structure

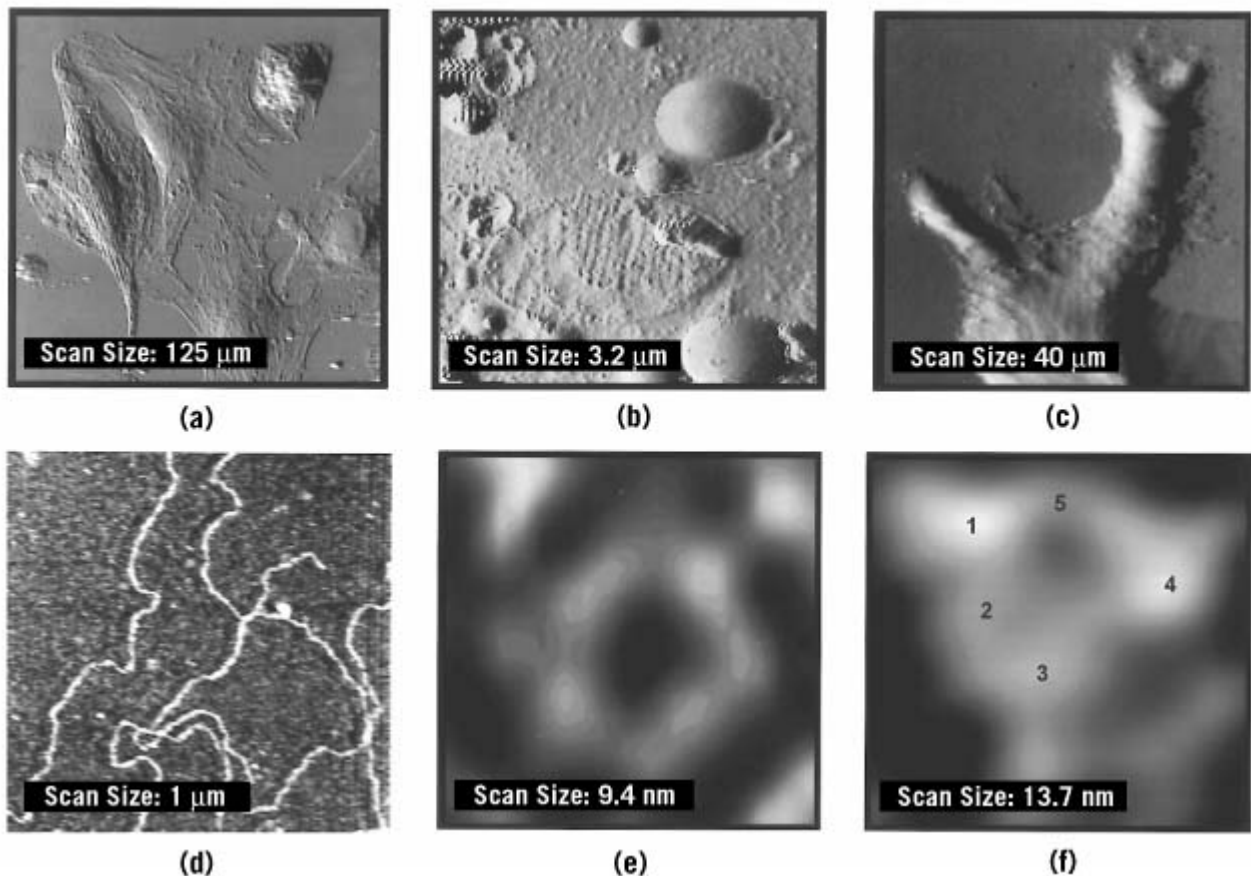
Although AFM provides molecular-resolution surface information for crystalline and amorphous materials, it is often difficult to define the nature of individual components, especially if the specimen contains a heterogeneous population of structures. This is the case with most biological systems, except in favorable systems like membranes that contain a crystalline patch of similar protein molecules. For mixed macromolecules, it is essential to compare the information obtained from AFM with that from alternative or complementary techniques, such as structural probes of electron microscopy and X-ray crystallography, biochemical and immunological binding assays, pharmacological labeling, and electrophysiological measurements.

## 1.6. Examples of AFM Imaging

### 1.6.1. Cells and Cellular Processes

AFM images cellular and subcellular structures in physiological conditions at a resolution far exceeding that of optical microscopes. Living cells have been imaged in aqueous conditions with a resolution as small as 10 nm. By applying a larger imaging force, intracellular **organelles** and **cytoskeleton** networks have also been examined (Fig. 2). The ability to view structures beneath the plasma membrane is puzzling. Two possible mechanisms are (1) the tip penetrates the bilayer and images the substructure or (2) the plasma membrane drapes around the cytoskeletal fibers and the tip images the contour of the plasma membrane. If the “drape” model is correct, these observations give wonderful demonstrations of the flexibility of biological membranes and a potential tool for measuring the “drape characteristics” of natural and synthetic membranes.

**Figure 2.** Images obtained with AFM. **(a)** Error-mode AFM image of a fixed atrial cell. The cytoskeletal network and nucleus are visible. [S. Shroff, D. Saner, and R. Lal (1995) *Am. J. Physiol.* **269**, C286–292.] **(b)** AFM error mode image of a **freeze-fracture** replica of atrial tissue. Details can be identified in the mitochondrion and in the atrial granules and vesicles. Scale Bar: 1  $\mu\text{m}$ . [L. Kordylewsky, D. Saner, and R. Lal (1994) *J. Microsc.* **173**, 173–181.] **(c)** Error-mode AFM image of a neurite outgrowth in PC-12 cells treated with **nerve growth factor** and dibutyryl **cyclic AMP**. For details, see R. Lal, B. Drake, D. Blumberg, D. Saner, P. K. Hansma, S. Feinstein (1995) *Am. J. Physiol.* **269**, C275–C285. **(d)** Lambda DNA under propanol using a regular silicon nitride tip. Note the sharp bends in the strands and a fairly regular series of lumps along the strands, 6 to 8 nm apart. The strand width is 7 to 9 nm, greater than expected, which may result from the relatively large size of the tip. Scan size is 1000 nm  $\times$  1000 nm (courtesy H. Hansma, H. G. Hansma, R. L. Sinsheimer, M. Q. Li, and P. K. Hansma (1992) *Nucleic Acids Res.* **20**, 3585–3590). **(e)** AFM height mode image of a single connexon (i.e. gap junction hemichannel) imaged on its cytoplasmic face. The subunit structure, a central pore, and the spacing between connections are apparent [S. A. John, D. Saner, J. Pitts, M. Finbow, and R. Lal (1997) *J. Struct. Biol.* **120**, 22–31.] **(f)** High-resolution AFM imaging of a single acetylcholine receptor expressed in *Xenopus* oocytes. The channel has a diameter of  $\sim 10.5$  nm, and the five subunits ( $\sim 1$  to 1.5 nm in diameter) that have a central pore-like structure are shown. The protrusion of one unit is not apparent in the image [R. Lal (1998) *Scanning Microsc.* **10**, 81–96].



In AFM imaging, the scan area can be varied from the micron to nanometer range, and hence it is possible to image features ranging from whole cells to individual macromolecules, such as **ion channels** and receptors (Fig. 2). In air-dried, hydrated *Xenopus* oocytes in which **acetylcholine receptor** proteins are expressed, the characteristic pentameric subunit structure of the expressed receptor has been observed after removing of the follicle layer. This technique of imaging expressed proteins on or in the surface of the plasma membrane of oocytes will shortly enable the characterization of a myriad of ion channels and receptors expressed in an appropriate expression system.

The major factor limiting resolution in imaging a cell surface is the mobility of the upper plasma membrane, plus the mobility of the macromolecules within the plasma membrane: the lower membrane is anchored to the substrate. Improvements in resolution may be made by (1) increasing the surface rigidity (e.g., suction of cells onto patch pipettes and thus reducing the lateral mobility of proteins); (2) by imaging with low forces (e.g., attractive force mode imaging or magnetic tapping imaging).

### 1.6.2. Membranes and Membrane Proteins

Imaging membranes, both native and reconstituted, has received wide attention because of their flattened 2-D sheet-like structure and the ease in preparing them. Purified membrane proteins, such as **bacteriorhodopsin**, **gap junctions**, **acetylcholine receptor**, and the hexagonally packed intermediate (HPI) layer from *Drosophila radiodurans*, have been imaged in aqueous conditions without fixation (Fig. 2). These membrane proteins have been characterized extensively by alternative techniques, such as electron microscopy and electron and X-ray crystallography. The results from AFM studies agree remarkably with those from other techniques. In addition, AFM images provide a direct observation of membrane polarity. For example, extracellular and cytoplasmic surfaces of gap junctions are distinguished unambiguously. The thickness measurement in AFM study is also direct and often very precise.

Synthetic membranes (Langmuir–Blodgett) and reconstituted vesicles have been imaged at molecular resolution. Images of Langmuir–Blodgett films provide direct measurement of lipid membrane thickness, obtained previously by indirect methods and theoretical extrapolations, because the height resolution in AFM is subnanometer. AFM images of Langmuir–Blodgett films show individual polar head groups and their molecular arrangement, including their long-range packing. An advantage of studying these membranes with AFM is that one can change the lipid composition on-line and study lipid–lipid interactions, lipid fluidity, and lipid–protein interactions. AFM images of proteins that are naturally embedded within membranes and form 2-D crystalline arrays and images of LB films provide some of the best evidence that the imaging of biological specimens generally agrees with that by electron microscopy, with the advantage, of course, that AFM imaging occurs in nearly physiological environments. It is also worth noting that such a correlation adds weight to the interpretation of images gathered by electron microscopy.

Proteins immobilized by synthetic membranes have also been imaged. Bacterial **porins** are one of the best-studied channel-forming **membrane proteins**. Porins reconstituted as 2-D crystals in lipid vesicles have been imaged in a liquid environment. AFM images at molecular resolution show the trimeric structure of porins illustrated by X-ray crystallography and electron microscopic **single-particle reconstruction**. In addition, recent studies show that molecular resolution can be obtained on noncrystalline specimens in liquid media. This opens a new avenue for studying the molecular structure of biological macromolecules (e.g., ion channels, receptors) that are easily expressed in an appropriate expression systems, such as the *Xenopus* oocyte, or simply isolated and anchored properly on suitable substrates. For example, purified **cholera toxin** molecules incorporated into synthetic phospholipid bilayers by covalent **cross-linking**, observed at molecular resolution, have the expected pentameric structure. Other membrane proteins imaged by AFM include the hexagonally packed intermediate (HPI) layer of *Drosophila radiodurans*,  $\text{Na}^+$ ,  $\text{K}^+$ -**ATPase**, vacuolar **proton**

**pumps** (V-H<sup>+</sup>-ATPase), and Ca<sup>2+</sup>-ATPase.

As mentioned before, AFM imaging of whole cell membranes has also been achieved. Acetylcholine receptor expressed in *Xenopus* oocytes has been observed. **Calcium channels** have been localized on the calyx-type nerve terminal of fixed chick ciliary ganglion in culture by imaging **avidin**-coated 30-nm gold particles incubated with w-**conotoxin** GVIA linked to **biotin**, although the molecular structure of individual calcium channels was not reported. The interchannel spacing of 40 nm was noticed, which may reflect the spatial limitation due to tagging with 30-nm gold particles. Individual calcium channels have much smaller diameters. Other membrane channels, such as **sodium channels**, **potassium channels**, and gap junctions, are 6 to 10 nm in diameter. Isolated cellular organelles have been imaged with AFM, including **nuclear pore complexes**, ~134 nm in outer diameter, with a central pore-like trough.

#### 1.6.3. Isolated Macromolecules, such as DNA, Amino Acids, and Proteins

Imaging isolated macromolecules is challenging because it is difficult to find suitable surfaces to which to anchor the molecules for repeatable and reliable imaging. The recent development of cryo-AFM shows good promise for obtaining high-resolution images of isolated macromolecules. Images at molecular resolution have been obtained of DNA at the **plasmid** and **chromosomal** levels, **polyamino acids**, isolated proteins, and ligand–receptor complexes. Large protein fibers, such as **actin** and **microtubules**, have also been imaged at molecular resolution, and it was possible to discern individual actin molecules. Isolated protein molecules show dynamic changes while imaging with the AFM. For example, when **glycogen phosphorylase b** binds to **phosphorylase kinase**, the dimensions and shapes of the proteins change noticeably.

Imaging DNA and nucleic acids has been appealing for many reasons. Given their well-known geometry and easy availability, they are readily identified and hence used for calibration and for studies of the interaction between tip and sample. Intriguingly, this may also open a door for structure-based sequencing and mapping of DNA AFM. However, **DNA sequencing** AFM will require an order-of-magnitude improvement in resolution (to ~2 to 3 Å). This increase may come from improvements in hardware and software, but methods to prepare DNA in extended conformations will probably be just as important in improving resolution.

Images of double-stranded DNA at molecular resolution (2 to 3 nm), in which the helical pitch and turns could be deciphered, have been obtained in air and liquid. Occasional images at higher resolution showing individual base pairs have also been obtained. Single-stranded DNA, though, has proven more intractable to image at any molecular resolution.

Images of complexes of DNA and **protein A** deposited onto mica show single proteins bound to the end of the DNA strands. In addition some single protein molecules bind to up to four DNA strands per protein molecule. When **RNA polymerase** binds to DNA, AFM images show that the modified DNA is bent at marked angles where the polymerase binds. One appeal of these approaches is in searching for **DNA-binding proteins** and, intriguingly, perhaps to image the effects of **topoisomerase** on DNA.

#### 1.6.4. Imaging Dynamic Processes

AFM, unlike other molecular level imaging systems, allows imaging in an aqueous environment. In an elegant set of studies, AFM was used to visualize real-time surface processes on **vaccinia virus pox viridae**-infected monkey kidney cells. Real-time changes in surface morphology and the **exocytosis** of enveloped virus and proteins were observed over a period of 19 hours. In contrast, cells not infected with virus showed no appreciable change in surface morphology. The real-time contractile activity of cultured atrial myocytes was also imaged. As the concentration of calcium increased, cells underwent rapid contraction, and a corresponding shortening in cytoplasmic fibers (perhaps cross-bridges) was observed.

Dynamic studies have been conducted on isolated proteins, such as formation of glycogen

phosphorylase-phosphorylase kinase complexes, [antibody–antigen interactions](#), dynamics of [immunoglobulin](#) adsorption, and binding of [streptavidin](#) to a biotinylated lipid bilayer. Real-time polymerization of **fibrin**, a protein important in blood clotting, shows that the polymer chain grows by the fusion of many short chains, rather than by successive addition of monomers to a few long chains. A change in Langmuir–Blodgett film morphology has been observed as trace amounts of a fluorescent dye are added, suggesting that the perturbation of molecular conformation by the tracer molecules may not be as insignificant as is commonly believed.

AFM can be used for *in situ* studies of the growth of protein 3-D crystals in their native solution environment and of the role of nucleation centers, lattice defects, and saturation level. These studies may provide clues for growing the 3-D crystals essential for high-resolution [X-ray crystallography](#).

### 1.7. Structure-Function Studies

AFM can be combined with other techniques, which opens the possibility, as various biochemical, pharmacological, and other perturbations are introduced on-line, of real-time dynamic studies for direct structure–function correlations at the molecular level. For example, “single cell” experiments have been reported where electrical activity and AFM images were obtained from *Xenopus* oocytes expressing acetylcholine receptor. Electrical recording of acetylcholine-sensitive current and labeling by specific binding of **a-bungarotoxin** were also conducted in parallel. The receptor density calculated from the AFM studies correlates well with that from electrical measurements and toxin binding, but the clustering of acetylcholine receptor differs from the uniform distribution of the expressed receptors that is commonly assumed in electrophysiological studies. A correlation between patch-clamp electrical recording and AFM of [transcription factor IID \(TFIID\)](#) interactions with the [nuclear pore complex](#) has also been reported, showing unplugging of the nuclear pore complex accompanied by prolonged electrical current through the channel, perhaps reflecting the reopening of the channels. A combined AFM and patch-clamp study measured the electrical current in the membrane patches excised from *Xenopus* oocytes and attached to the pipette tip, while imaging the surface topology with AFM. Although the resolution of such a study is limited, it nevertheless shows the promise of direct structure-function studies of membrane macromolecules. Another study simultaneously measured images of [bacteriorhodopsin](#) in purple membranes adsorbed onto a lipid monolayer, ion transport through the membrane, and the electrical properties of the membrane. A recent study simultaneously measured the surface structure of nuclear pore filters and electrical current passing across the filter through pores of different diameters. For such a study, a scanning ion-conductance microscope was developed that records electrical activity, while imaging the 3D-structures of various membranes (Fig. 1).

Imaging force can be varied considerably during experiments. Such a feature has been used to nanomanipulate protein and membrane structures, and two different conformations of membrane proteins have been reported.

As already mentioned, action is at biological surfaces. The potential for rapidly characterizing and cataloguing the structures of synthetic peptides opens vistas for synthesizing drugs that can interact with “similar surfaces” within the body. Because AFM measures the force of an interaction between substrate and cells, interactions of a ligand or agonist should be measurable if well-defined molecules are placed on the tip of the AFM. Then it may be possible to use this ligand as a probe to determine the presence or absence of a receptor in impure, natural preparations. Perhaps more importantly, it may be possible to measure the interactive forces between agonist and receptor. This information will prove useful in designing of drug inhibitors or mimics.

### 1.8. Analysis of Micro-mechanical Properties

Contrast mechanism and image formation in AFM reflect a sum of many local forces and the micromechanical properties of the specimen. As the choice of specimen is shifted from rigid and hard materials (such as mica, graphite, tungsten) to soft and deformable biological materials, the dominating micro-mechanical properties shift from pure frictional to viscoelastic. Frictional forces on the atomic scale have been measured between two silica surfaces, two thin films, and between



tungsten and graphite. By using local frictional force contours, fluorinated and hydrocarbon regions were distinguished in a Langmuir–Blodgett film. The hydrocarbons and fluorocarbons had separate domain structures: hydrocarbons as circular domains, fluorocarbons as the surrounding flat films. The frictional force was higher in the fluorinated region than in the hydrocarbon region. Such compositional studies provide a mechanism to identify individual components in a multicomponent sample like a cell membrane.

The viscoelastic properties of several soft, deformable biological materials, including cells, has also been obtained by a technique called *force mapping*. Such a study can be undertaken during on-line pharmacological and biochemical perturbations. In each imaging pixel, the tip is brought into proximity with the surface until a preset deflection is reached. The tip then retracts to its original position, and this process is repeated in every pixel. A height image is obtained from the amount of vertical piezo movement necessary at each point to obtain the preset deflection of the cantilever. For each pixel, a deflection versus distance curve is stored, which can be fitted to different models to obtain properties, such as the elastic modulus of the sample surface. Usually, the tip apex is approximated by a semisphere or a cone and the specimen by a spherical or planar model, depending on the shape of the features on the surface.

Thus, AFM is an imaging tool and also a system for analyzing micromechanical properties of cells, subcellular organelles, and macromolecules. It may be possible to study localized viscoelastic properties of molecular motor units, the distribution and propagation of contraction waves in a muscle cell, and the correlation between the calcium concentration wave and electrical propagation. One can also induce local shearing (frictional) force or pressure to assess the effects on the vascular system (mimicking the role of blood flow-induced shearing on vasorelaxation) or to distinguish pressure- or shear-sensitive ion channels in plasma membranes.

### 1.9. Simultaneous Multimodal Imaging

The simple design of AFM allows integrating it with other techniques, such as light [fluorescence microscopy](#), laser [confocal microscopy](#), and [near-field scanning optical microscopies](#). Such integrated systems permit simultaneous multimodal imaging and provide independent verification with appropriately labeled markers. For example, using appropriately labeled **fluorescent** signals, one can identify specific areas and then use AFM to obtain high-resolution details.

#### 1.9.1. Combined AFM and Light Fluorescence Microscope

Although conventional AFMs are ideal for high-resolution imaging, they could not be combined with large-aperture optical microscopes. In a few AFMs, however, the cantilever moves and the sample is stationary, permitting the addition of optical microscopes that have high numerical apertures. The most promising of these AFMs has the scanned-cantilever mode in which the cantilever position is accurately tracked by a scanned focused spot (Fig. 1) and is incorporated into an inverted fluorescence microscope. This combined fluorescence and force microscope has been used to image **immunolabeled** membranes and whole cells (Fig. 2). Fluorescent labels show remarkable correspondence among AFM images and the specificity of the molecules: such correspondence that one can obtain structural information at molecular resolution on biological macromolecules present individually or in small clusters, long as they have detectable fluorescent signals.

#### 1.9.2. Combined Atomic Force Microscope and Confocal Microscope

Early combined AFM and laser-scanning confocal microscopes (LSCM) included features like a stationary sample stage, an AFM with an optical tracking system for the scanned cantilever, and either a scanned-beam or tandem design confocal microscope. The limitations of such systems include a limited scan range of both the AFM and confocal images. The scan range of the independently scanned confocal spot is limited to the size of the field of view of the objective and off-axis optical aberrations, and the scanned-cantilever AFMs with optical level detection require optical tracking of the cantilever for a large scan range. The latter constraint was overcome by the scanned-cantilever (tip) design with an optical tracking feature (the features explained previously in

the combined AFM-fluorescence microscope), where the sample is scanned by a piezo system and the AFM tip and objective remain stationary. The AFM registers the topography of the sample surface, and the LSCM laser scans the surface to obtain fluorescent data on the same scan area. Although the AFM scan size is increased in this improved design, the confocal scan size is still limited by the objective. Moreover, although the AFM images are of the sample surface, the confocal image plane may not be the sample surface, but anywhere within the confocal slice, which could be no more than 100 nm thick.

A new combined AFM and LSCM allows simultaneous imaging of the sample surface in both modes, in addition to the conventional confocal imaging through the sample thickness. The salient features of such a combined microscope include a scanned-sample approach wherein the specimen is scanned above an inverted microscope objective with a fixed optical path for fluorescent LSCM imaging. An AFM positioned directly above the sample simultaneously measures the surface topography. Therefore, in this design the confocal spot and AFM cantilever remain stationary. Optical cantilever tracking is not required, and the confocal spot can be centered in the microscope's objective. The AFM feedback system ensures that the focal point is on or near the surface of the specimen, so that when the cantilever is positioned at the confocal spot, the LSCM and AFM images are acquired in direct registration, allowing image features to be easily correlated.

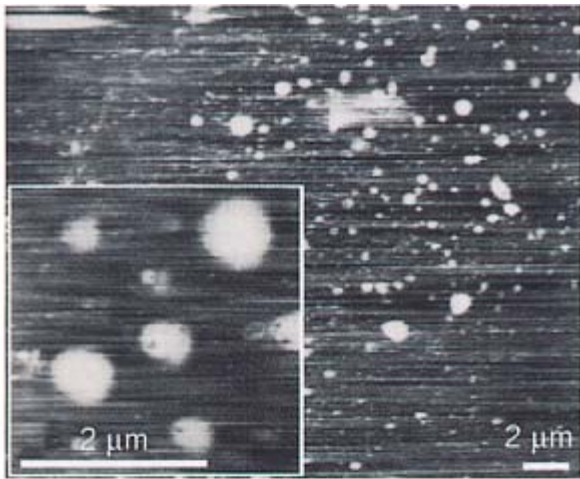
In this combined multimodal system, the confocal plane can be selected from the topmost region on a specimen surface or anywhere through its depth, so it is quite possible to follow, for example, cytoplasmic [signal transduction](#) processes leading to changes in the cellular plasma membrane surface conformations.

### 1.9.3. Combined Atomic Force Microscopy and Electrophysiological Recording

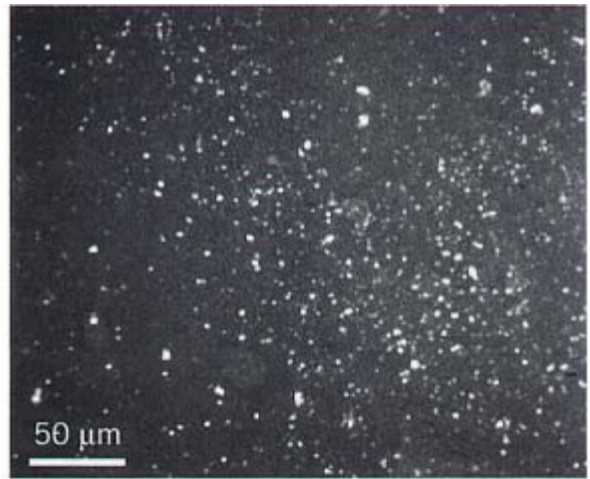
A combined tapping-mode AFM and a scanning ion-conductance microscope have been developed recently (Fig. 1). One of the salient features of this combined microscope is a bent glass pipette used as both the force sensor and the conductance probe. The force-sensing capability allows measuring of the pipette deflection, which then is used to create surface images in both regular contact mode and tapping mode. The conductance-measuring capability allows recording the electrical current flow across pores in a suitable specimen. Using such a microscope, it is possible to image the structures of channels and receptors and to measure their functional states (conducting vs. non-conducting) (Fig. 3).

**Figure 3.** Examples of AFM multimodal imaging. (a) and (b) Simultaneous immunofluorescence and atomic force microscopy of amyloid b peptide (AbP) reconstituted into liposomes. All liposomes, with or without AbP, were imaged with AFM (a), a few are shown at higher magnification in the inset. The liposomes were treated with anti-AbP antibody and subsequently identified with fluorescein-conjugated second antibody. The AbP-carrying liposomes showed strong fluorescence signals (b). [S. K. Rhee, A. P. Quist, and R. Lal (1998) *J. Biol. Chem.* **273**, 13379–13382]. (c), (d) and (e) Adhesion sites between a *Xenopus* retinal glial cell (XR1 cell line) and [extracellular matrix](#) material in a cell culture. The fluorescent images show the location of b-integrin (c) and f-actin (d) fibers detected by immunofluorescence, and the tapping mode AFM image (e) reveals the 3-D architecture of the focal point after removing of the cell body [R. Lal and R. Proksch (1997) *Int. J. Imaging Syst. Technol.* **8**, 293–300]. (f) and (g): Simultaneously combined AFM and fluorescence-confocal microscopic images. The sample was a suspension of fluorescently labeled latex beads that were dried into a gel on a plastic diffraction grating. The lines of the grating are visible in the topographic AFM image (f) but not in the confocal fluorescent image (g). The two images allow distinguishing a nonfluorescent particle (left arrow) from a fluorescent particle (right arrow), although both appear as raised bumps in the AFM image. [P. E. Hillner, D. A. Walters, R. Lal, H. G. Hansma, and P. K. Hansma (1995) *J. Micro. Soc. Am.* **1**, 123–126]. (h) and (i): Simultaneously combined AFM and scanning ion-conductance microscope (SICM) electrophysiology. (h) shows a tapping mode AFM image of a nucleopore membrane, and (i) shows the associated ionic conductivity image obtained by tapping mode SICM. Note that there are some differences in the pores detected by the two procedures. For example, the area circled in white shows a groups of pores that appear to be deep in the AFM image and highly conductive in the SICM image. The area circled in grey contains a large pore that appears deep in the AFM image but is nonconductive in the SICM image. The scale bars at the bottom are intensity-coded. Brighter is a greater height in the AFM image and a greater conductance in the SICM image. [R. A. Proksch, R. Lal, P. K. Hansma, G. Morse, and G. Stucky (1996) *Biophys. J.* **71**, 2155–2157.]

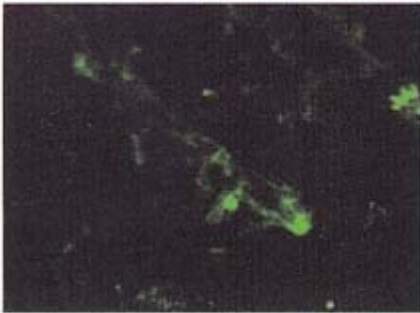




(a)



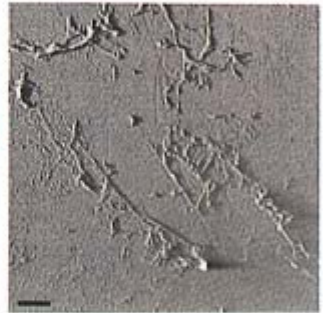
(b)



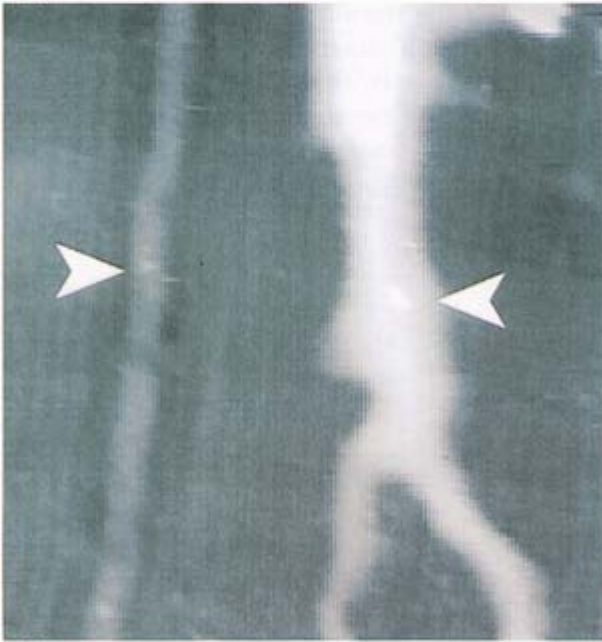
(c)



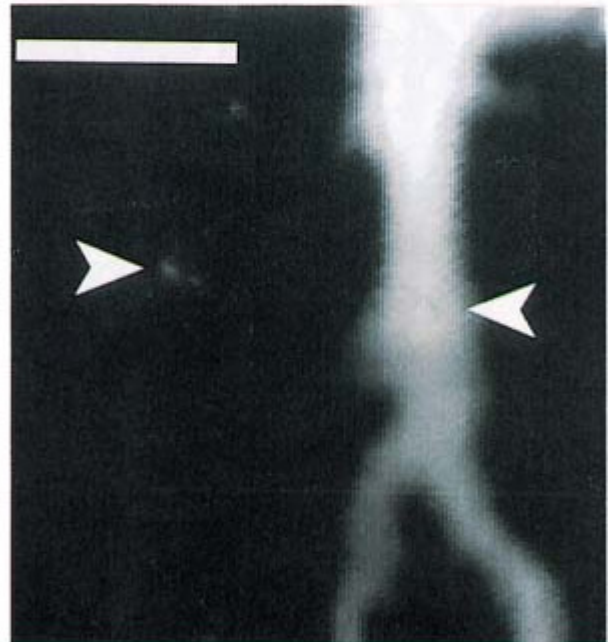
(d)



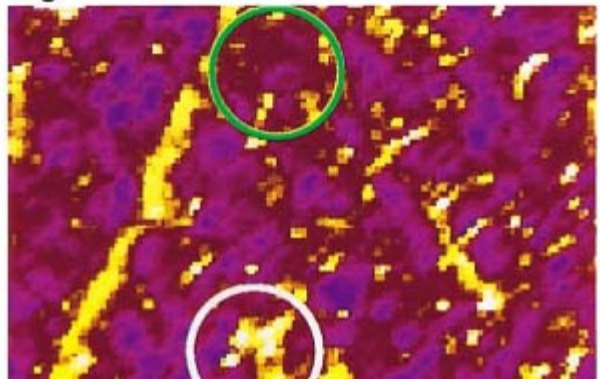
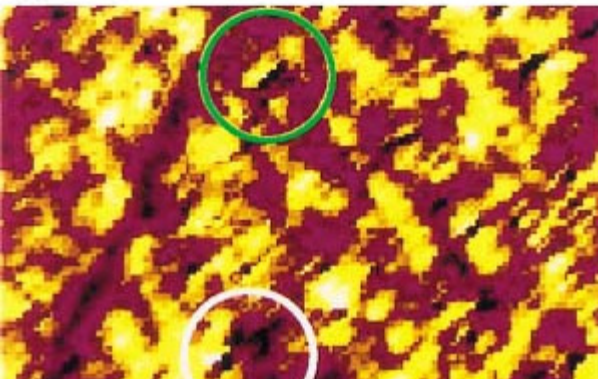
(e)



(f)



(g)



Another approach is combining AFM with the patch-clamp technique in the same experiment. Such a combined technique records electrical current in the excised membrane patches from *Xenopus* oocytes that are attached to the patch pipette tip, while simultaneously observing the surface topology with AFM. Also, The membrane surface is also deformed by applying pressure through the patch pipette and observing the lateral displacement of features. However, the resolution is limited to about 10 nm.

#### Suggestions for Further Reading

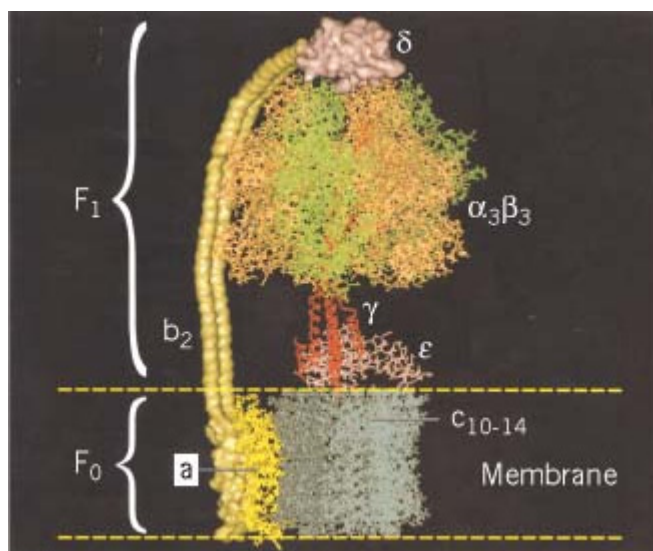
- G. Binnig, C. F. Quate, and C. Gerber (1986) Atomic force microscope, *Phys. Rev. Lett.* **56**, 930–933.
- B. Drake, S. A. C. Gould, A. L. Weisenhorn, H. G. Hansma, P. K. Hansma, C. F. Quate, C. B. Prater, T. R. Albrecht, and D. S. Cannell (1989) Imaging crystals, polymers, and processes in water with the atomic force microscope, *Science* **243**, 1586–1589.
- P. K. Hansma, V. B. Elings, C. E. Bracker, and O. Marti (1988) Scanning tunneling microscopy and atomic force microscopy—application to biology and technology, *Science* **242**, 209–216.
- E. Henderson, P. G. Haydon, and D. S. Sakaguchi (1992) Actin filament dynamics in living glial cells imaged by atomic force microscopy, *Science* **257**, 1944–1946.
- P. E. Hillner, D. A. Walters, R. Lal, H. G. Hansma, and P. K. Hansma (1995) Combined atomic force and confocal laser scanning microscope, *J. Micro. Soc. Am.* **1**, 123–126.
- J. H. Hoh, R. Lal, S. A. John, J. P. Revel, and M. F. Arnsdorf (1991) Atomic force microscopy and dissection of gap junctions, *Science* **253**, 1405–1408.
- R. Lal (1998) Imaging molecular structure of channels and receptors with an atomic force microscope, *Scanning Microsc.* **10**, 81–96.
- R. Lal and S. A. John (1994) Biological applications of atomic force microscopy, *Am. J. Physiol. Cell Physiol.* **266**, C1–C21.
- R. Lal and R. Proksch (1997) Multimodal imaging with atomic force microscopy: Combined atomic force, light fluorescence, and laser confocal microscopy and electrophysiological recordings of biological membranes, *Int. J. Imaging Syst. Technol.* **8**, 293–300.
- F. Ohnesorge and G. Binnig (1993) True atomic resolution by atomic force microscopy through repulsive and attractive forces, *Science* **260**, 1451–1456.
- B. N. J. Persson (1987) The atomic force microscope—can it be used to study biological molecules, *Chem. Phys. Lett.* **141**, 366–368.
- R. A. Proksch, R. Lal, P. K. Hansma, D. Morse, and G. Stucky (1986) Imaging the internal and external pore structures of membrane in fluid: Tapping mode scanning ion conductance microscopy, *Biophys. J.* **71**, 2155–2157.
- D. Sarid (1991) *Scanning Force Microscopy: With Applications to Electric, Magnetic and Atomic Forces*, Oxford University Press, New York.
- Y. Z. Zhang, S. Sheng, and Z. Shao (1996) Imaging biological structures with the cryoatomic force microscope, *Biophys. J.* **71**, 2168–76.

#### ATP Synthase

ATP synthase, also called  $F_0F_1$  [ATPase](#), or simply F-ATPase, is the universal [protein](#) that terminates oxidative phosphorylation by synthesizing ATP from ADP and phosphate. Nearly identical proteins are found in eukaryotic [mitochondria](#) and bacteria, and they all operate on the same principle. Electron-driven ion **pumps** set up concentration and electrical gradients across a membrane (see [Chemiosmotic Coupling](#) and [Proton Motive Force](#)). ATP synthase utilizes the energy stored in this electrochemical gradient to drive nucleotide synthesis. It does this in a surprising way by converting the electromotive force into a rotary torque that promotes phosphate binding and liberates ATP from the catalytic site where it was formed. Remarkably, this process can be reversed in certain circumstances. ATP hydrolysis can drive the engine in reverse, so that F-ATPase functions as a proton pump. Indeed, the vacuolar V-ATPases, the most ubiquitous intracellular proton pumps, are structurally similar to ATP synthase and operate according to the same principles.

ATP synthase is composed of at least eight subunit types, whose stoichiometries are denoted by the subscripts ( $a_3$ ,  $b_3$ ,  $g$ ,  $d$ ,  $\epsilon$ ,  $a$ ,  $b_2$ ,  $c_{12}$ ), that combine into two distinct regions. The geometric arrangement of the subunits is shown schematically in Fig. 1(a). The soluble  $F_1$  portion consists of a hexamer,  $a_3b_3$ . This hexamer is arranged in an annulus about a central shaft consisting of the [coiled-coil](#)  $g$  subunit. Subunits  $d$  and  $\epsilon$  are also generally isolated with  $F_1$ . The  $F_0$  portion consists of three transmembrane subunits,  $a$ ,  $b_2$ , and  $c_{12}$ . The 12 copies of the  $c$ -subunit form a disk into which the  $g$  and  $\epsilon$  subunits insert. The remainder of  $F_0$  consists of the transmembrane subunits  $a$  and  $b_2$ . The latter is attached by the  $d$  subunit to an  $a$  subunit, so that it anchors the  $a$  subunit to  $F_1$ . Thus there are two “stalks” connecting  $F_0$  to  $F_1$ ,  $g\epsilon$  and  $b_2d$ .

**Figure 1.** Schematic diagram of the subunit organization of ATP synthase (1). (a) The  $a_3b_3$  hexamer and a portion of the  $g$  shaft. The lower part of  $g$  has not been resolved. The  $c$ -subunit consists of 12 pairs of transmembrane  $\alpha$ -helices, and the  $a$  subunit consists of six transmembrane  $\alpha$ -helices. The  $\epsilon$  subunit abuts  $c$  and  $g$  and interacts with the DELSEED region of  $b$ . The  $a$  subunit is attached to an  $a$  subunit via the  $b$  and  $d$  subunits. (b) The proton motive force at the  $a$ - $c$  interface leads to the functional subdivision into two counterrotating assemblies, usually denoted as the “rotor” and “stator.” The rotor consists of subunits  $c_{12}$ - $g$ - $\epsilon$ , and the stator consists of subunits  $a$ - $b_2$ - $d$ - $a_3b_3$ .



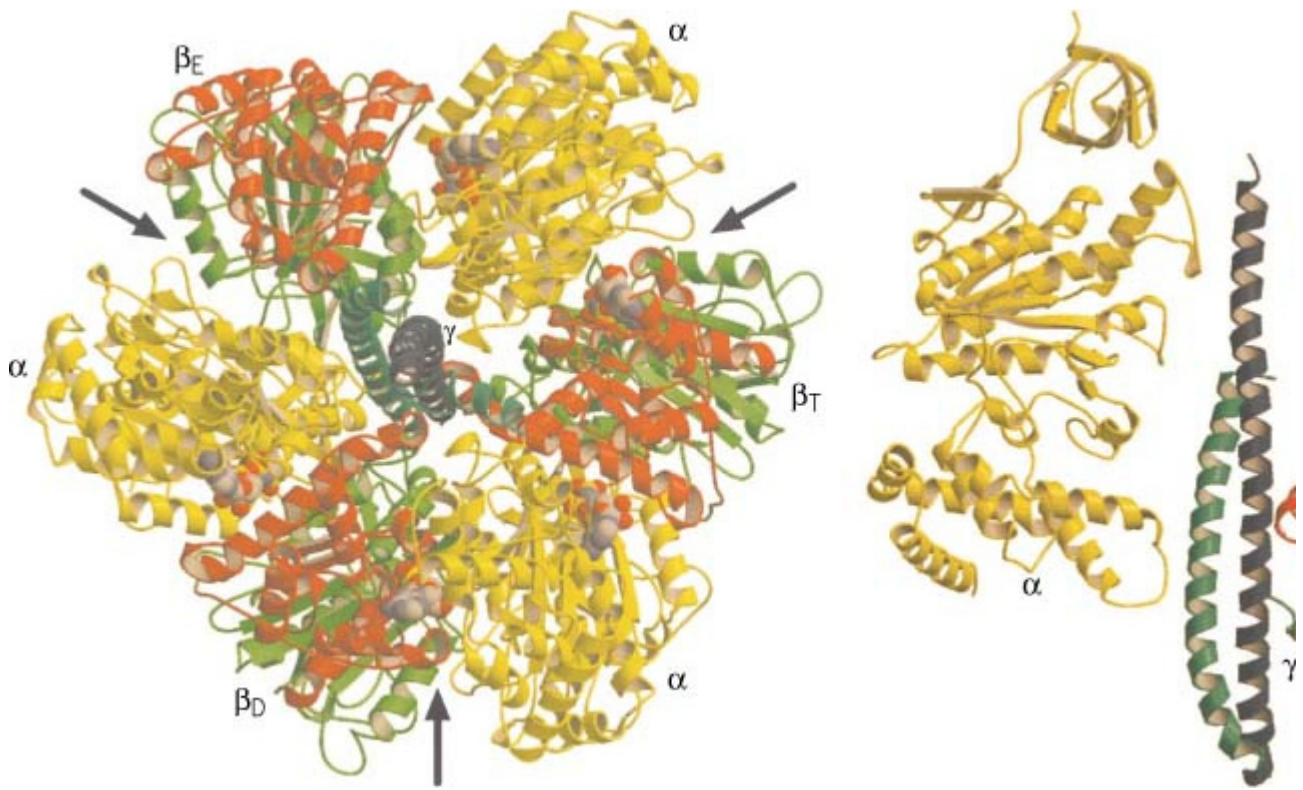
The key to understanding how ATP synthase carries out its catalytic and synthetic roles lies in this geometric organization. The entire protein can be divided into two operational regions denoted suggestively as the “rotor” and the “stator” for reasons that derive from the rotary mechanism by

which the protein operates (Fig. 1b). Indeed, it turns out that ATP synthase has two rotary engines in one. The  $F_0$  motor converts transmembrane electrochemical energy into a rotary torque on the g shaft, and  $F_1$  uses ATP hydrolysis to turn the g shaft in the opposite direction. Because they are connected, one drives the other in reverse. When  $F_0$  dominates, the rotor turns clockwise (looking upward in Fig. 1), so that  $F_1$  synthesizes ATP. When  $F_1$  dominates, so that  $F_0$  is driven counterclockwise, it can pump protons against an electrochemical gradient. Deciphering how this remarkable dual energy transduction works is one of the great triumphs of modern chemistry.

### 1. ATP Synthesis in $F_1$ is Driven by Rotation of the g-Shaft

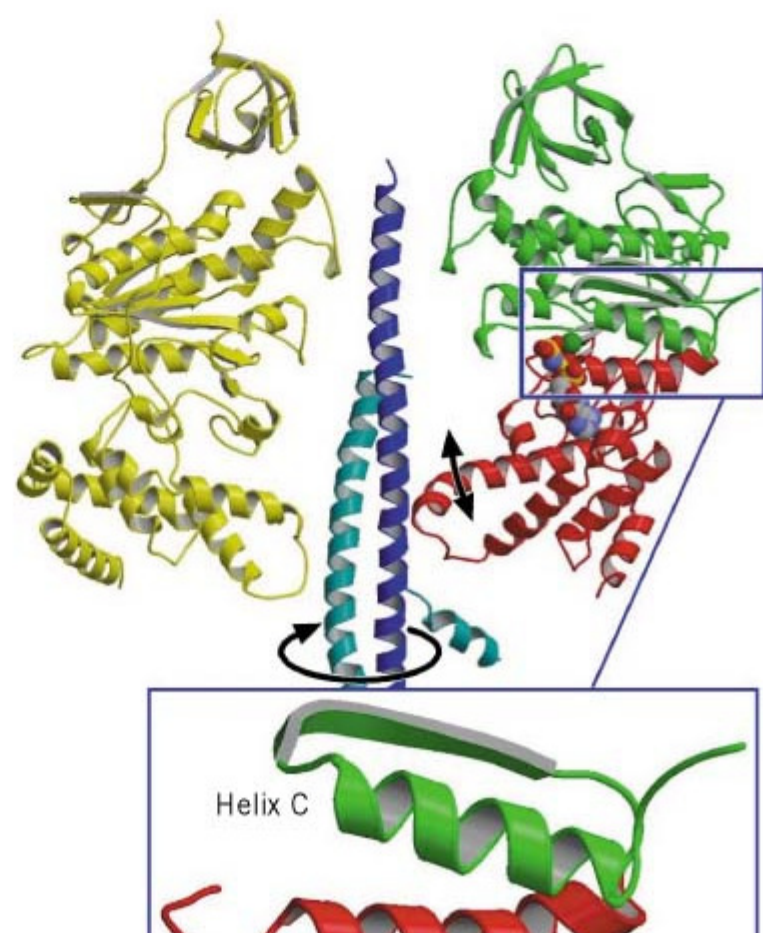
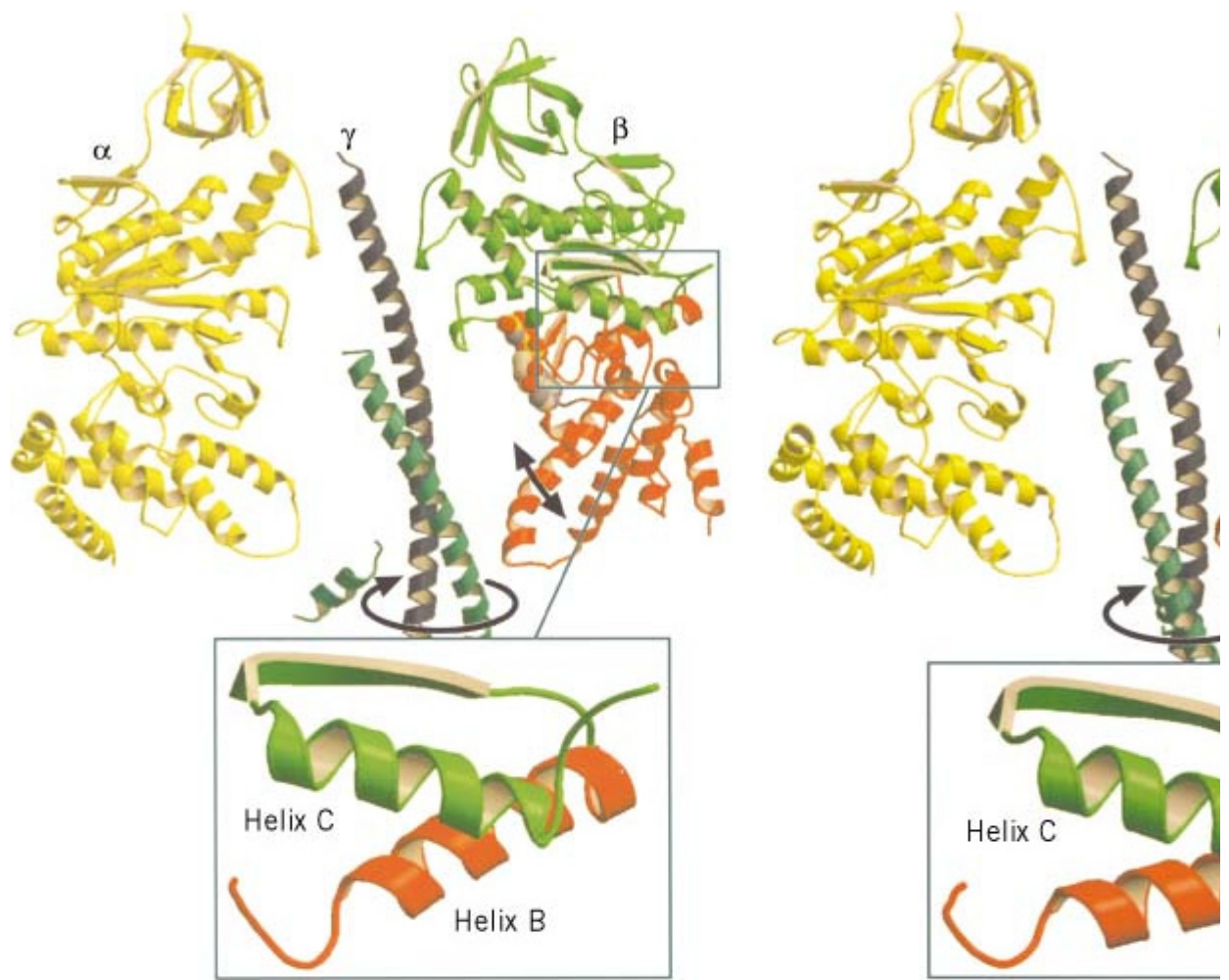
We begin with the  $F_1$  motor, because we now know precisely what it looks like. This is due to John Walker and the [X-ray crystallography](#) group at Cambridge, who worked out the exact structure of the  $a_3b_3$  hexamer and most of the g shaft (1). A stereo view of the structure is shown in Fig. 2. Walker was awarded the Nobel Prize in 1997 because his structure revealed essential asymmetries in the molecule's structure that were the key to understanding how it worked. In the early 1980's, Paul Boyer at UCLA proposed the surprising theory that, in the catalytic sites of  $F_1$ , ATP is in chemical equilibrium with its reactants, ADP and phosphate (2). So the formation of ATP is essentially without cost energetically. However, because each ATP hydrolyzed under cellular conditions liberates about 12 kcal/mol, this energetic price must be paid at some point. Boyer proposed that  $F_1$  pays this price in the mechanical work necessary to liberate the nucleotide from the catalytic site. Further, release of product (ATP) proceeds sequentially and cyclically around the  $a_3b_3$  hexamer because the synthetic reactions are synchronized in a fixed-phase relationship by rotation of the g shaft and cooperative coupling between the three catalytic sites. Boyer's "binding change" mechanism neatly fits the Walker structure, and Boyer shared the 1997 Nobel Prize. A schematic diagram of the binding change mechanism is shown in Fig. 3.

**Figure 2.** Stereo pair showing the molecular structure of the  $F_1$  subunit run through center, coils extending out from mid in light gray, and the two coils of the g subunit medium gray. Its asymmetrical structure is evident. In the b subunit: the s is in dark gray, and the lower hinge segment is medium gray. The structure of the  $c$  subunit (not shown) is known, but it the g and c subunits.

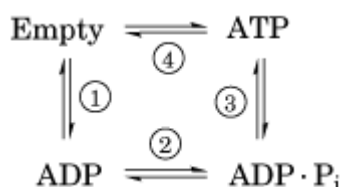


**Figure 3.** The binding change mechanism. Notation for site occupancies: T = ATP bound, DP = ADP • P<sub>i</sub> bound, D = ADP bound. The subunits are numbered clockwise. The lengths of the arrows indicate the relative binding affinities. (a) The system starts with either (b<sub>1</sub>, b<sub>2</sub>, b<sub>3</sub>) = (D, T • D • P, D). (b) Clockwise rotation of g increases the binding affinity of ADP in b<sub>1</sub>, traps ATP in b<sub>2</sub>, and (c) Further rotation of g traps ADP and allows P<sub>i</sub> binding in b<sub>1</sub>, releases the tightly bound ATP and allows ADP binding





There are actually six nucleotide-binding sites on the  $a_3b_3$  hexamer. All lie at the interfaces between the a and b subunits. The catalytic sites lie mostly in the b subunit, whereas the noncatalytic sites lie mostly in the a subunit. The role of the noncatalytic sites is uncertain, but they may help to hold the hexamer together. Each catalytic site traverses the synthetic cycle sequentially:

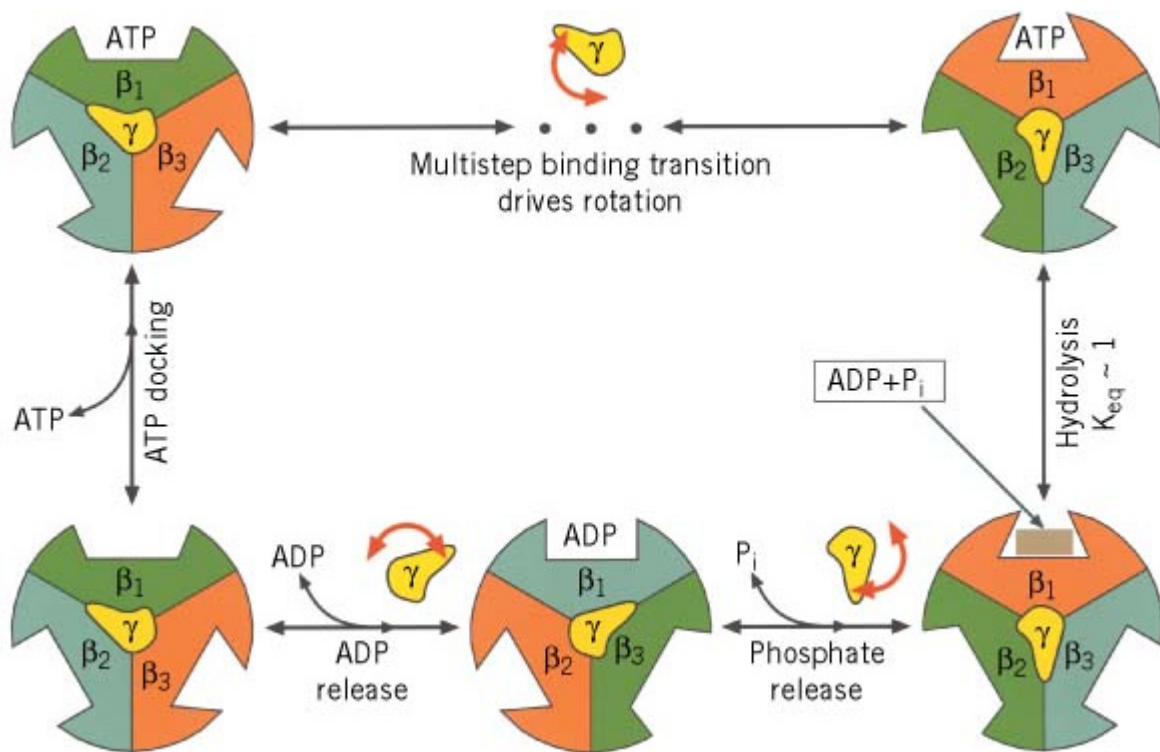


In steps 1 and 2, a site binds ADP and phosphate (not necessarily in that order). While trapped in the catalytic site in step 3, the reactants (ADP and  $P_i$ ) and product (ATP) are in chemical equilibrium.

Step 4 requires the input of mechanical torque from  $F_0$  on g to trap the reactants in the ATP state and to pry open the site, releasing the tightly bound ATP. Most of the 12 kcal/mol price of synthesis is paid in step 4. The way in which this works is found in the shape of the  $a_3b_3$  hexamer and the g shaft.

The g subunit is asymmetric and bowed. It fits into a central annulus in  $a_3b_3$ , which is itself asymmetric (Fig. 4). At the top of the  $a_3b_3$  hexamer is a **hydrophobic** “sleeve” in which the g shaft rotates. Further down, however, the annulus is offset from the center, so that as g rotates clockwise, it sequentially pushes outward on each catalytic site. In addition, the  $\epsilon$  subunit is located eccentrically and attached to the g and c subunits, so that as g rotates, it comes into contact sequentially with each b subunit in a conserved region called the DELSEED sequence (named for the single-letter abbreviation of its constituent amino acid residues). Together, this asymmetrical rotation exerts stress on the catalytic site, loosening its grip on ATP, so that thermal fluctuations can free it into solution.

**Figure 4.** Cross section of  $F_1$  showing the conformational changes in the b subunits that drive rotation of g. The a subunit is on the left and the stationary barrel region of b is on the right. During the hydrolytic cycle, the lower segment of b undergoes a hinge-bending motion that rotates it about  $30^\circ$  inward. This motion pushes on the eccentric g coiled coil, causing it to rotate within the barrel bearing. During synthesis, the rotation of g pushes on each catalytic site. The panels show three snapshots of the motion during a  $180^\circ$  rotation. Movies of the rotational sequence can be downloaded from the authors' World Wide Web site: (site currently unavailable)



(3). These interactions may mediate phosphate- and nucleotide-binding, the necessary precursors to synthesis. In addition, the catalytic sites are elastically coupled, so that the occupancy of one site affects the affinity of the other two sites. The consequence of this coupling is that when ATP concentrations are low enough so that only one site is occupied, hydrolysis proceeds much more slowly than when more than one site is occupied.

Together with the F<sub>1</sub> molecular structure, the binding change model strongly supports the idea that catalysis involves rotation of the g subunit. However, dramatic visual confirmation was provided by *in vitro* experiments in which the a<sub>3</sub>b<sub>3</sub>g subunits were isolated and attached to a bead. A fluorescently tagged [actin](#) filament was attached to the g shaft, and when ATP was supplied, the filament could be clearly seen rotating. In fact, a complete revolution takes place in three steps and consumes a single ATP per step (4).

The viscous drag on the actin filament was estimated, which allowed computing the torque developed by the F<sub>1</sub> motor and comparing it with the **free energy** available from ATP hydrolysis. The startling result was that the motor generates an average torque of more than 40 piconewton nanometers ( $40 \times 10^{-12} \text{ N} \times 10^{-9} \text{ m}$ ), more than six times the maximum force developed by [kinesin](#) or myosin. More impressively, the motor operates near 100% mechanical efficiency, this precludes any sort of heat engine, which would be limited by the Carnot efficiency (4). Several models have been proposed that address the issue of torque generation and efficiency (4-6).

The energy to drive this motion derives from the hydrolytic cycle of ATP at the catalytic site. Moreover, the conformational change that drives the hydrolysis motor must be nearly the reverse of the motion that frees ATP from the catalytic site during synthesis. Examination of the structure reveals that the major conformational change is a hinge-bending motion in the b subunits. The bottom portion of each b subunit below the nucleotide-binding site rotates inward approximately 30°, during which it pushes on the bowed g subunit, turning it much like one cranks an automobile jack (Fig. 4).

## 2. F<sub>0</sub> Converts Proton Motive Force into Rotary Torque

Currently, there have been no direct observations of rotation in the  $F_0$  portion of ATP synthase (7). However, current thinking is that the  $F_0$  assembly converts the energy contained in the transmembrane proton motive force into a rotary torque at the interface of the a and c subunits (Fig. 1). This torque turns the rotor (the c, g, and  $\epsilon$  subunits), which couples to the  $F_1$  synthetic machine.

The c assembly consists of 12 subunits, each consisting of two transmembrane [alpha-helices](#) (8). There is one essential acidic amino acid (Asp61 in the *Escherichia coli* ATP synthase) which binds protons. Because variants of ATP synthase can operate on sodium rather than protons, the interaction between the c subunit and the translocated ion has the property of an electrostatic carrier mechanism (9).

The a subunit consists of six transmembrane a-helices that contain at least one essential basic residue (Arg210 in *E. coli*) (8, 10). The [electrostatic interaction](#) of these rotor and stator charges is essential for torque generation, and several proposals have been put forward for the way this could work (11-14). Whatever the mechanism, the  $F_0$  motor must generate a torque sufficient to liberate three ATP's from the three catalytic sites in  $F_1$  for each revolution. If the proton flux through the stator is tightly coupled to the rotation of the c subunit, then a rotation of  $2\pi/3$  carries four protons down the electromotive potential of 230 mV typical of the mitochondrial inner membrane (14). This is sufficient to account for the mechanical energy required for synthesis of one ATP.

Under anaerobic conditions, the ATP synthase of the bacteria *E. coli* can reverse its operation, hydrolyzing ATP and turning the c subunit backward so that it functions as a proton pump. This is not surprising because the F-ATPases are structurally similar to the most common proton pumps, the vacuolar, or V-ATPases (7). These pumps may have been the **evolutionary** precursors of ATP synthase (15). A striking difference between the two is that the F-ATPases have 12 acidic rotor charges, whereas the V-ATPases have six. It can be shown that this enables the V-ATPases to function more efficiently as ion pumps at the expense of relinquishing their capability to synthesize ATP.

### 3. Summary

Both the  $F_1$  and  $F_0$  motors can operate in both directions.  $F_1$  is a hydrolytically driven three-piston engine that can be driven in reverse to synthesize ATP from ADP and phosphate.  $F_0$  is an ion-driven rotary engine that can be driven in reverse to function as an ion pump. The F-ATPases are structurally similar to and presumably evolutionarily related to the V-ATPase ion pumps (15). It is thought that most ion pumps function by an "alternating access" mechanism, whereby an ion is first bound strongly on the dilute side, then energy is supplied to move the ion so that it communicates with the concentrated side and weakens its binding affinity (16). In contrast with other ion pumps, however, the F and V-ATPases accomplish this by a rotary mechanism that is driven indirectly by nucleotide hydrolysis, rather than by direct phosphorylation (17). It is thought that the  $F_0$  motor is also related to the bacterial **flagellar** motor. Both can operate on sodium, although the flagellar motor has eight or more "stators" and develops far more torque than  $F_0$  (18, 19).

The mechanism driving the  $F_1$  hydrolytic motor may carry hints for other nucleotide hydrolytically fueled motors, such as kinesin, myosin, and [dynein](#). However, important structural differences may make the comparison difficult (4). For example, the motors previously mentioned all "walk" along a polymer track to which they bind tightly during a portion of their mechanochemical cycle. The power stroke of the  $F_1$  motor is driven by the b subunit, which pushes on the g shaft, but does not bind tightly to it, that is, it does not "walk" around the g shaft. Moreover, no other motor operates

with nearly the efficiency of the  $F_1$  motor, implying that there are important entropic steps in other motors that are absent in the  $F_1$  motor.

### Bibliography

1. J. Abrahams, A. Leslie, R. Lutter, and J. Walker (1994) *Nature* **370**, 621–628.
2. P. Boyer (1993) *Biochim. Biophys. Acta* **1140**, 215–250.
3. M. Al-Shawi, C. Ketchum, and R. Nakamoto (1997) *J. Biol. Chem.* **272**, 2300–2306.
4. K. Kinoshita, R. Yasuda, H. Noji, S. Ishiwata, and M. Yoshida (1998) *Cell* **93**, 21–24.
5. F. Oosawa and S. Hayashi (1986) *Adv. Biophys.* **22**, 151–183.
6. H. Wang and G. Oster (1998) *Nature* **396**, 279–282.
7. M. Finbow and M. Harrison (1997) *Biochem. J.* **324**, 697–712.
8. R. H. Fillingame (1997) *J. Exp. Biol.* **200**, 217–224.
9. P. Dimroth (1997) *Biochim. Biophys. Acta* **1318**, 11–51.
10. R. H. Fillingame (1996) *Curr. Opin. Struct. Biol.* **6**, 491–498.
11. S. B. Vik and B. J. Antonio (1994) *J. Biol. Chem.* **269**, 30364–30369.
12. W. Junge, H. Lill, and S. Engelbrecht (1997) *Trends Biochem. Sci.* **22**, 420–423.
13. G. Kaim, U. Matthey, and P. Dimroth (1998) *EMBO J.* **17**, 688–695.
14. T. Elston, H. Wang, and G. Oster (1998) *Nature* **391**, 510–514.
15. R. Cross and L. Taiz (1990) *FEBS Lett.* **259**, 227–229.
16. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. Watson (1994) *Molecular Biology of the Cell*, Garland, New York.
17. S. Khan (1997) *Biochim. Biophys. Acta* **1322**, 86–105.
18. H. Berg (1995) *Biophys. J.* **68**, 163s–166s.
19. K. Muramoto, I. Kawagishi, S. Kudo, Y. Magariyama, Y. Imae, and M. Homma (1995) *J. Mol. Biol.* **251**, 50–58.

### Suggestions for Further Reading

20. P. Boyer (1993) The binding change mechanism for ATP synthase—some probabilities and possibilities. *Biochim. Biophys. Acta* **1140**, 215–250.
21. J. Weber and A. E. Senior (1997) Catalytic mechanism of  $F_1$ -ATPase. *Biochim. Biophys. Acta* **1319** (1), 19–58.

### ATP-Binding Motif

The ATP-binding motif, also called the *Walker motif*, is a structural [protein motif](#) that is frequently found in proteins that bind ATP or GTP. It can usually be identified from just the [primary structure](#) of a protein using a fingerprint sequence that characterizes proteins as ATP or GTP-binding (see [GTP-Binding Proteins](#)). Walker et al. (1) identified two regions of sequence conservation: the A region has a stretch of small **hydrophobic** residues followed by [Gly/Ala]—X—X—Gly—X—Gly—Lys—Thr/Ser, where X is any residue. The second region of sequence conservation also has a stretch of small hydrophobic residues, this time ending in a conserved **aspartate** residue. In the

three-dimensional structures of proteins having the ATP-binding motif, such as adenylate kinase, the hydrophobic residues of the A-region form a buried  $\beta$ -strand, and the glycine-rich region forms a loop, called the **P-loop**, that interacts with the phosphate of the bound nucleotide. The second conserved region codes for another hydrophobic  $\beta$ -strand, and the conserved aspartate is required for binding the magnesium ion that usually accompanies nucleotides bound to proteins.

The ATP-binding motif represents a common, but not the only, mode of interaction between proteins and ATP/GTP. For example, it is structurally distinct from the actin fold of proteins such as [actin](#), hsp70, and hexokinase, which bind and hydrolyze ATP or GTP and where a conformational change is implicated in nucleotide binding (2).

[See also [Nucleotide-Binding Motif](#) and [P-Loop](#).]

#### Bibliography

1. J. E. Walker, M. Saraste, M. Runswick, and N. Gay (1982) EMBO J. **1**, 945–951.
2. W. Kabsch and K. C. Holmes (1995) FASEB J. **9**, 167–174.

#### Suggestions for Further Reading

3. G. E. Schulz (1992) Binding of nucleotides by proteins. Curr. Opin. Struct. Biol. **2**, 61–67 (Excellent review of nucleotide binding motifs.)

## ATPase

Any [enzyme](#) that catalyzes the hydrolysis of ATP to ADP and inorganic phosphate ( $P_i$ ) is classified as an ATPase. In some cases, the phosphate is transiently transferred to the protein before its release as a product. The hydrolysis of ATP liberates much energy, so this reaction is usually coupled to another, energetically unfavorable reaction. Three major types of ATPase are associated with membranes that couple ATP hydrolysis to the translocation of specific ions across a membrane: P-, V-, and F-ATPases.

P-ATPases have a relatively simple polypeptide composition (one or two subunits) and are phosphorylated as part of their catalytic cycle. Examples of P-ATPases are (1) the  $Na^+$ ,  $K^+$  – ATPase of the plasma membrane of animal cells; (2) the  $H^+$  – ATPase of the plasma membranes of **yeast**, **fungi**, and plants; and (3) the  $Ca^{2+}$  – ATPase of the sarcoplasmic reticulum (the [endoplasmic reticulum](#) of muscle). The plasma membrane  $Na^+$ ,  $K^+$  – ATPase and  $H^+$  – ATPase function to generate and maintain the plasma membrane electrical potential difference (inside negative), as well as the ionic disequilibria across that membrane.  $Ca^{2+}$  – ATPases function in  $Ca^{2+}$  homeostasis (see [Calcium Signaling](#)).

V-ATPases have a much more complicated polypeptide composition than do P-ATPases and are not phosphorylated during catalysis. V-ATPases are found on the membranes within the interiors of **eukaryotic** cells (endomembranes), including membranes from vacuoles (from which the “V” is derived), **lysosomes**, **Golgi**, **secretory vesicles**, **clathrin**-coated vesicles, and, in some instances, plasma membranes. V-ATPases are  $H^+$  – ATPases that show some similarity to the proton-linked [ATP synthases](#). The function of V-ATPases is to catalyze proton transport into the endomembrane interior compartments at the expense of ATP hydrolysis. Acidification of the interior of

endomembranes is required for some of their functions. In **archaebacteria**, an enzyme with similarity to V-ATPases functions as an ATP synthase.

F-ATPases (also known as “F<sub>1</sub>-F<sub>0</sub>”) have a complex polypeptide composition and, as in V-ATPases, there is no phosphorylated intermediate in the reaction mechanism. In non**photosynthetic** eukaryotes, the F-ATPase is found exclusively on the inner membrane of [mitochondria](#), whereas in green plants and algae there are two distinct F-ATPases; one in mitochondria and the other on the thylakoid membrane of [chloroplasts](#). In bacteria, F-ATPase is present in the plasma membrane.

ATPases couple the flow of protons (or Na<sup>+</sup> in some cases) to ATP hydrolysis and synthesis. The activity of the F-ATPases is very tightly regulated. Several mechanisms combine to prevent wasteful ATP hydrolysis by F-ATPase, but to allow rapid ATP synthesis at the expense of electrochemical proton (or Na<sup>+</sup>) potentials generated by proton translocation linked to electron transport (see [Chemiosmotic Coupling](#)). F-ATPases operate *in vivo* as ATP synthases; consequently, the term *ATP synthase* is preferred to *F-ATPase* when referring to the entire enzyme. The catalytic portion of the ATP synthase may be readily removed from coupling membranes and purified as a large 400-kDa water-soluble protein. Such preparations may have high ATPase activity and are called “F<sub>1</sub>” or “F<sub>1</sub>-ATPase” (see [ATP Synthase](#)).

## Attenuation Of Transcription

The ability to modulate **gene** expression in response to changing environmental signals is crucial for the survival of all organisms. Virtually every stage involved in the synthesis, function, and degradation of [macromolecules](#) is a potential target for one or more regulatory events ([1](#)). Regulatory mechanisms have been identified for all three stages of [transcription](#) (initiation, elongation, and termination). Several postinitiation regulatory mechanisms have been categorized as transcription attenuation mechanisms. Transcription attenuation can be defined as any mechanism that utilizes transcription pausing or transcription termination to modulate expression of downstream genes. For the purpose of this article, however, the definition will be restricted to situations in which the action of the regulatory molecule promotes transcription termination, with the default situation being transcriptional readthrough. There are also related **antitermination** mechanisms in which the action of the regulatory molecule promotes transcriptional readthrough.

Once transcription of a gene is initiated, the transcription elongation complex and its nascent transcript are potential targets for regulation. As transcription proceeds, the nascent transcript may fold into specific secondary and tertiary structures that signal the transcribing **RNA polymerase** to pause or terminate transcription before reaching the structural genes ([1](#)). Transcription attenuation mechanisms allow the organism to modulate the extent of transcriptional readthrough past the terminator structure in response to changing environmental signals, thereby regulating expression of the downstream genes. As will be seen, several different transcription attenuation mechanisms have been identified.

### 1. Transcription Attenuation of Biosynthetic Operons of Enteric Bacteria

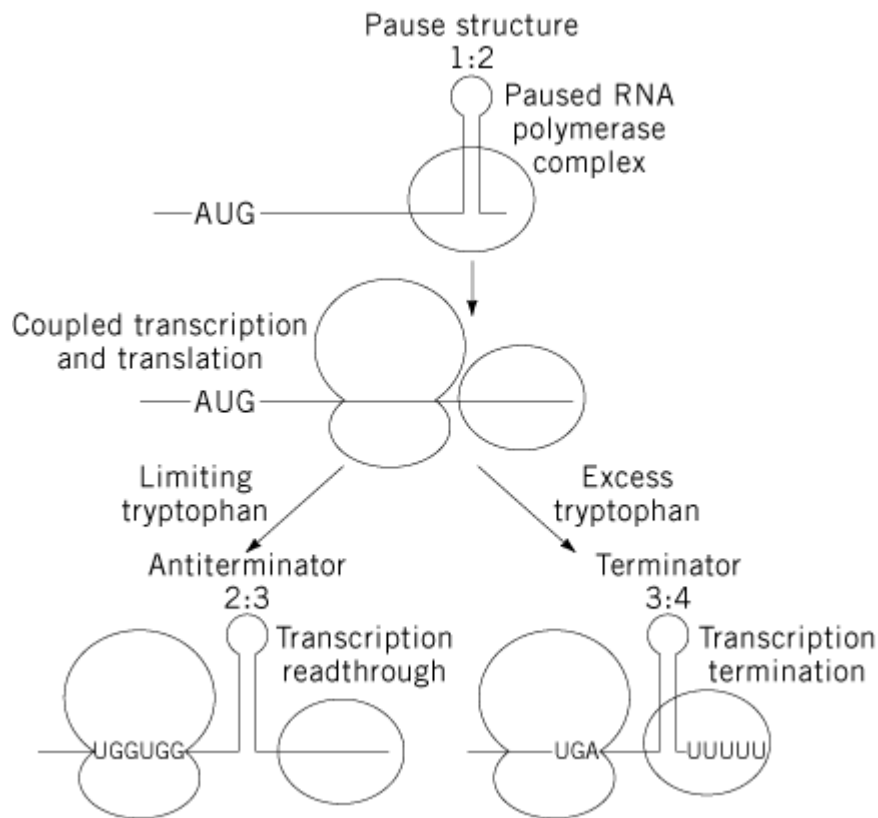
Transcription attenuation was the first demonstration that organisms can exploit [RNA structure](#) to modulate gene expression. The first attenuation mechanism was elucidated by Charles Yanofsky and his coworkers for the *Escherichia coli* tryptophan biosynthetic operon (*trpEDCBA*; see [TRP Operon](#)) ([2](#)). In addition, many other amino-acid biosynthetic [operons](#) in enteric bacteria are regulated by transcription attenuation (eg, *his*, *leu*, *ilv*, *pheA*). In each case, the genetic information required for

transcription attenuation is encoded within a 150–300-bp leader region located between the **promoter** and the first structural gene of the operon (1). Because the salient features of transcription attenuation are conserved in each system, the *E. coli trp* operon will be discussed, and the key differences with respect to other operons will be pointed out where appropriate.

Transcription initiation of the *E. coli trp* operon is regulated by TrpR, a **DNA-binding** repressor protein. Once transcription starts, the elongating transcription complex is subject to control by transcription attenuation (1). The combined actions of repression (80-fold) and transcription attenuation (eight-fold) result in approximately 600-fold regulation in response to changing concentrations of intracellular tryptophan (3). A simplified model of the *E. coli trp* operon transcription attenuation mechanism is depicted in Figure 1 (4). The 141-nucleotide leader transcript can form three overlapping RNA secondary structures, referred to as the pause structure, the antiterminator, and a Rho-independent terminator. In addition, the nascent *trp* leader transcript contains a small open reading frame that encodes a 14-amino acid residue leader peptide. Soon after transcription of the *trp* operon is initiated, a secondary structure (structure 1:2) forms in the nascent transcript that signals RNA polymerase to pause. The paused RNA polymerase complex allows sufficient time for a **ribosome** to initiate **translation** of the leader peptide. The translating ribosome then disrupts the paused RNA polymerase complex, and transcription resumes, with the ribosome closely following the molecule of RNA polymerase, thereby coupling transcription and translation. At this point, two different outcomes can occur, depending on the level of tryptophan in the cell. Under conditions of limiting tryptophan, the level of charged **transfer RNA** tRNA<sup>Trp</sup> is low. As a result of the low tryptophanyl-tRNA<sup>Trp</sup> concentration, the translating ribosome stalls at one of two tandem Trp **codons** strategically placed within the leader peptide coding sequence. Ribosome stalling at the Trp codons effectively uncouples transcription and translation. As transcription proceeds, therefore, the antiterminator structure (structure 2:3) forms and prevents formation of the overlapping Rho-independent terminator (structure 3:4), resulting in transcriptional readthrough into the *trp* structural genes. Under conditions of tryptophan excess, the level of charged tRNA<sup>Trp</sup> is sufficiently high to allow efficient translation of the tandem Trp codons, and the ribosome continues to the end of the leader peptide. When the ribosome reaches the leader peptide **stop codon**, it physically blocks formation of the antiterminator structure, thereby promoting terminator formation and, hence, termination of transcription, before RNA polymerase reaches the *trp* structural genes. Thus, expression of the *trp* operon is decreased when the cell has an adequate supply of tryptophan. As can be seen by the above example, the regulatory signal is charged tRNA<sup>Trp</sup>, and the sensory event is the capacity to translate a short peptide-coding sequence (1). The transcription attenuation mechanisms for several other amino acid biosynthetic operons, such as the *his*, *phe*, and *leu* operons, are essentially identical to that for the *trp* operon, except that the leader peptides contain seven His (5), seven Phe, and four Leu codons (6), respectively.

**Figure 1.** Model of transcription attenuation for the *E. coli trp* operon. RNA polymerase pauses following formation of the pause structure provides time for a ribosome to initiate translation of the leader peptide. Under tryptophan-limiting conditions, the ribosome stalls at the tandem Trp codons, resulting in transcription readthrough. Under conditions of tryptophan excess, the ribosome reaches the leader peptide stop codon. This ribosome position blocks formation of the antiterminator, leading to terminator formation and transcription termination. (See text for details.) Adapted from Landick and Yanofsky (4).





Expression of the *E. coli* pyrimidine biosynthetic operon, *pyrBI*, is also regulated by transcription attenuation (7, 8). In this case, the concentration of UTP serves as the regulatory signal, in conjunction with a UTP-dependent pause signal consisting of an RNA hairpin and several U residues just after the hairpin. A transcription terminator exists approximately 60 nucleotides downstream of the pause structure, but the leader transcript does not have the potential to form an antiterminator structure. Finally, there is a 44-residue leader peptide encoded by an open reading frame beginning prior to the pause signal and extending past the terminator. The model for this attenuation mechanism is as follows: When there is a deficiency of UTP, the transcribing RNA polymerase pauses at the leader pause site. This provides time for a ribosome to initiate translation of the leader peptide, which results in coupled transcription and translation. As transcription proceeds, the translating ribosome prevents formation of the transcription terminator, allowing transcription of the structural genes. However, pausing is inefficient when the cell contains an adequate supply of UTP. In this case, RNA polymerase transcribes and recognizes the terminator before the ribosome reaches this segment of the leader transcript, thus halting transcription in the leader region prior to the structural genes.

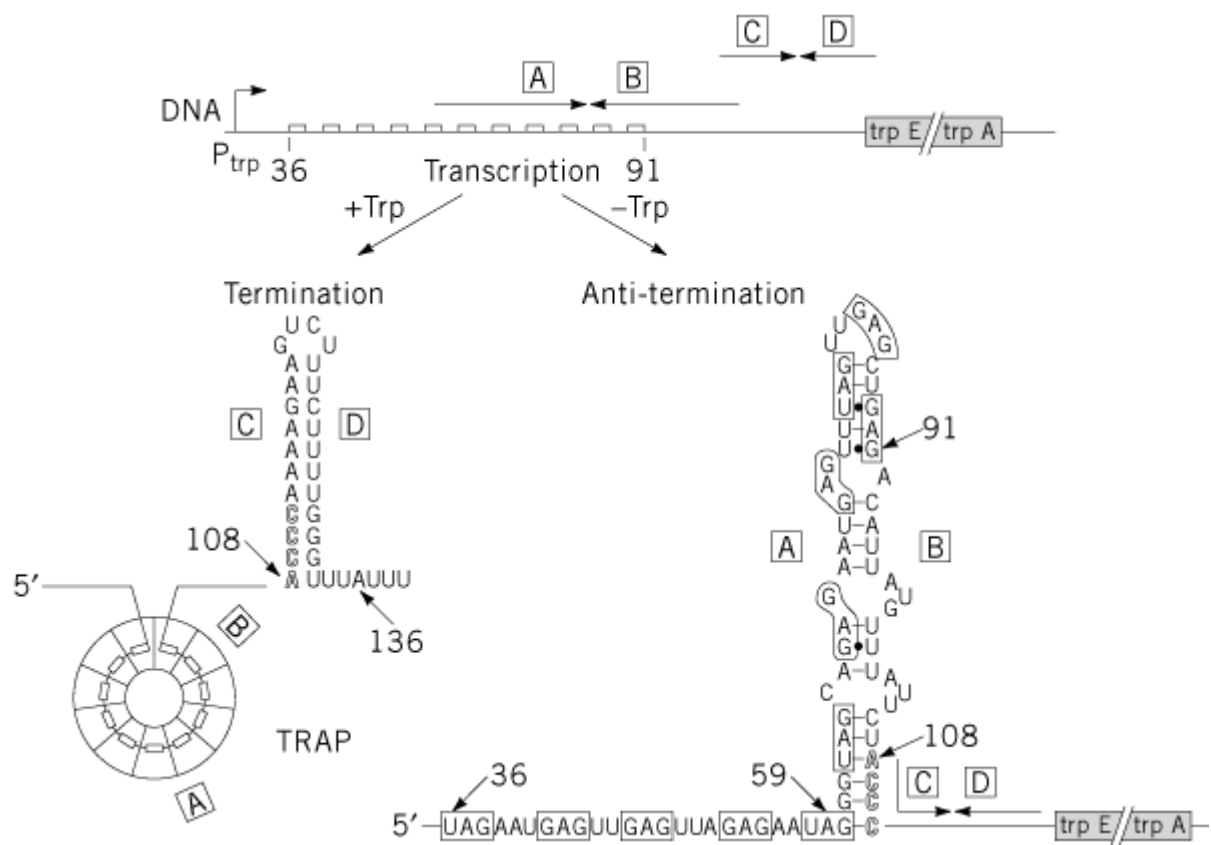
## 2. Transcription Attenuation of Gram-Positive Biosynthetic Operons

The transcription attenuation mechanisms that have been identified for the *trp* and *pyr* operons in the **Gram-positive bacterium** *Bacillus subtilis* differ dramatically from those described for the enteric bacteria. Most notably, ribosomes and tRNA molecules are not involved. Instead, sequence-specific [RNA-binding proteins](#) are responsible both for sensing the level of tryptophan or UMP in the cell and, ultimately, for the decision to terminate transcription or to readthrough into the structural genes.

Expression of the *trpEDCFBA* operon of *B. subtilis* is regulated by TRAP, the *trp* RNA-binding attenuation protein (9, 10). TRAP is composed of 11 identical subunits arranged in a single ring (11). Tryptophan binding between each adjacent subunit in a cooperative manner activates TRAP to bind RNA (11). Transcription initiation of the *trp* operon appears to be constitutive, occurring 203

nucleotides upstream of the first structural gene. A transcription attenuation model for the *B. subtilis trp* operon is depicted in Figure 2. The *B. subtilis trp* leader transcript contains inverted repeats that can form mutually exclusive antiterminator and Rho-independent terminator structures (12), although there is no apparent transcription pause signal in this case. When cells are growing in excess tryptophan, tryptophan-activated TRAP binds to 11 closely spaced (G/U)AG repeats present in the *trp* leader transcript (11, 13). Recent studies have shown that TRAP binds to these repeats by wrapping the RNA around the protein ring, with the bases of the (G/U)AG repeats interacting with several amino acid residues on adjacent subunits in the protein (14). TRAP binding blocks formation of the antiterminator, which allows formation of the overlapping terminator structure. Thus, transcription halts in the leader region prior to the *trp* structural genes. Under conditions of limited tryptophan, TRAP is not activated and does not bind to the *trp* leader transcript. Thus, as transcription proceeds, the antiterminator forms, allowing transcription readthrough into the *trp* structural genes (9, 12).

**Figure 2.** Model of transcription attenuation of the *B. subtilis trp* operon. The large boxed letters designate the complementary strands of the terminator and antiterminator RNA structures. Small rectangles represent the GAG and UAG repeats involved in TRAP binding; these triplet repeats are also outlined in the sequence of the antiterminator structure. Numbers indicate the residue positions relative to the start of transcription. Nucleotides 108 to 111 overlap between the antiterminator and terminator structures and are shown as outlined letters. The TRAP protein is represented as an 11-subunit ring, and the bound RNA is shown forming a matching circle on binding to TRAP, with each triplet repeat interacting with one subunit. From Antson et al. (11).



Expression of the *B. subtilis pyr* operon is also controlled by transcription attenuation mediated by an RNA-binding protein. This operon encodes 10 polypeptide chains involved in de novo synthesis of pyrimidines (15, 16). Transcription attenuation occurs in response to UMP levels through the action of PyrR, which is encoded by the first gene of the operon (17). There are several novel features of the attenuation mechanism that controls this operon. In contrast to the systems described previously,

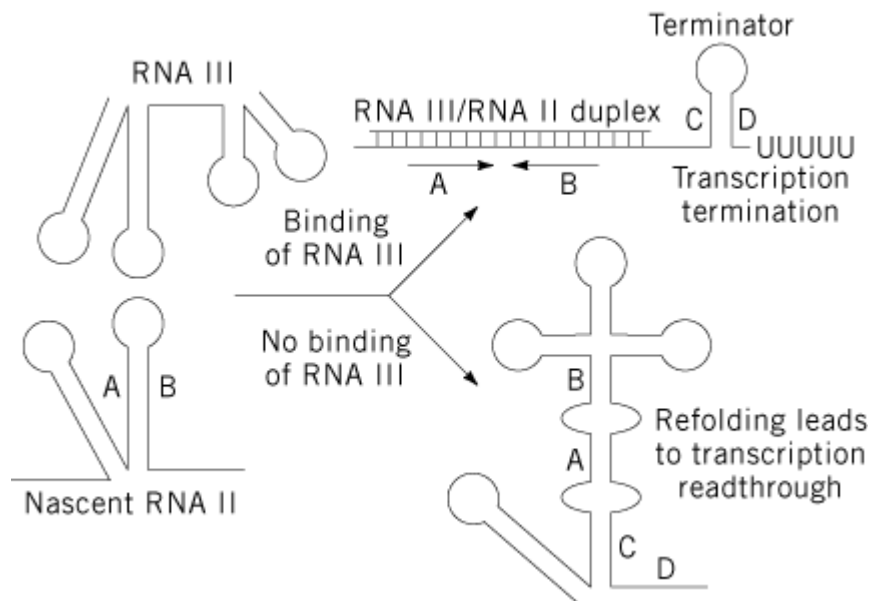
the *pyr* operon contains three attenuators: one located in the 5'-untranslated leader region, the second between the first and second genes, and the third between the second and third genes of the operon. Each attenuator can form three alternative RNA secondary structures. In addition to terminator and antiterminator structures similar to those described previously, a structure called the anti-antiterminator can form upstream of, and overlapping, the antiterminator. In the presence of UMP, PyrR is activated to act at each of the three attenuators to promote transcription termination and downregulate expression of the operon. In contrast to the mechanism by which TRAP controls attenuation in the *trp* operon, PyrR binding does not interfere directly with formation of the antiterminator structure, but rather functions by stabilizing the anti-antiterminator, which thereby indirectly stabilizes the terminator (18). The nature of the binding sites for PyrR and TRAP reflects the different effects these proteins have on RNA structure when they bind. TRAP binds to an entirely single-stranded site (19, 20), whereas PyrR binds to a stem-loop structure (21).

Another novel feature of this system is that PyrR is not only an RNA-binding regulatory protein, it is also an [enzyme](#), uracil phosphoribosyltransferase (UPRTase), which catalyzes formation of UMP from uracil and 5-phosphoribosyl-1-pyrophosphate. The physiological role of this enzyme is not clear, as *B. subtilis* has an additional UPRTase that has been shown to be more important for UMP synthesis. The structure of PyrR in the absence of UMP shows two oligomeric forms of the protein, one as a dimer and the other as a hexamer (21). Both forms appear to exist in solution, but PyrR is thought to bind RNA as a dimer. Neither the amino acid sequence nor the structure of PyrR show significant similarity to TRAP. Thus it appears that these two similar attenuation mechanisms evolved independently.

### 3. Antisense RNA-Mediated Transcription Attenuation

Antisense RNA control of gene expression has been documented for many prokaryotic genes, some of which involve an interesting form of transcription attenuation (22-24). For example, the copy number of the *Streptococcus agalactiae* plasmid pIP501 is regulated by antisense RNA-mediated transcription attenuation (23, 24). The antisense RNA (RNA III) inhibits expression of *repR*, the gene encoding the essential RepR initiator protein, by binding to the nascent *repR* leader transcript (RNA II). This interaction promotes formation of a Rho-independent terminator upstream of the *repR* coding sequence. Interaction of RNA III with RNA II is initiated by the formation of a “kissing complex” between single-stranded loops of both molecules, followed by propagation of the RNA helix. In the absence of RNA III interaction, formation of the transcriptional terminator is prevented, and expression of *repR* can proceed normally. A model illustrating this mechanism is depicted in Figure 3.

**Figure 3.** Model of plasmid-encoded *repR* transcription attenuation. Interaction of RNA III with the nascent RNA II (*repR*) transcript promotes terminator formation and transcription termination. In the absence of RNA III interaction, refolding of the nascent transcript blocks formation of the terminator, leading to transcription readthrough. (See text for details.) Adapted from Brantl and Wagner (24).



#### 4. Transcription Attenuation of Eukaryotic Genes

Expression of many eukaryotic genes is also regulated at the level of elongation of transcription by processes similar to attenuation (for reviews see [25](#), [26](#)). In many cases, the action of an inducer protein is required to release the stalled RNA polymerase II elongation complex. In several systems, including *c-myc* ([27-31](#)), *N-myc* ([32](#)), *c-fos* ([33](#)), and adenosine deaminase ([34](#)), potential stem-loop structures, similar to prokaryotic transcription terminators, have been identified near the sites of attenuation. However, the precise transcription attenuation mechanisms of these systems are not known.

The best-characterized eukaryotic system involving control of transcription elongation is transcriptional activation of HIV-1 by the transactivator protein Tat ([35](#)). In this system, transcription initiates at the HIV-1 promoter in the viral [long terminal repeat](#) (LTR). In the absence of Tat, transcription terminates prematurely prior to the structural genes. Tat recognizes an RNA target called TAR, located near the 5'-end of the viral transcript ([36](#), [37](#)). After binding to TAR, Tat interacts with several host cell proteins to enhance the processivity of RNA polymerase II and allow expression of the HIV-1 genes. Transcriptional pausing of RNA polymerase II just downstream of the TAR RNA structure was recently demonstrated. This could allow time for Tat to interact with TAR before RNA polymerase terminates prematurely ([38](#)). Interestingly, the pause signal is an RNA stem-loop similar to those seen in attenuation control in enteric bacterial systems.

#### Bibliography

1. C. Yanofsky (1988) *J. Biol. Chem.* **263**, 609–612.
2. C. Yanofsky (1981) *Nature* **289**, 751–758.
3. C. Yanofsky, R. L. Kelley, and V. Horn (1984) *J. Bacteriol.* **158**, 1018–1024.
4. R. Landick and C. Yanofsky (1987) In *Escherichia coli and Salmonella typhimurium: cellular and molecular biology* (F. C. Neidhardt, J. L. Ingraham, K. B. Low, B. Magasanik, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology, Washington, DC, pp. 1276–1301.
5. S. W. Artz and D. Holzschu (1983) In *Amino Acids: Biosynthesis and Genetic Regulation* (K. M. Herrmann and R. L. Summerville, eds.), Addison-Wesley Publishing Co., Reading, Mass, pp. 379–404.
6. P. W. Carter, J. M. Bartkus, and J. M. Calvo (1986) *Proc. Natl. Acad. Sci USA* **83**, 8127–8131.

7. K. L. Roland, F. E. Powell, and C. E. Turnbough (1985) *J. Bacteriol.* **163**, 991–999.
8. S. P. Lynn et al. (1987) *J. Mol. Biol.* **194**, 59–69.
9. P. Babitzke and C. Yanofsky (1993) *Proc. Natl. Acad. Sci. USA* **90**, 133–137.
10. J. Otridge and P. Gollnick (1993) *Proc. Natl. Acad. Sci. USA* **90**, 128–132.
11. A. A. Antson et al. (1995) *Nature* **374**, 693–700.
12. H. Shimotsu, M. I. Kuroda, C. Yanofsky, and D. J. Henner (1986) *J. Bacteriol.* **166**, 461–471.
13. P. Babitzke, J. T. Stults, S. J. Shire and C. Yanofsky (1994) *J. Biol. Chem.* **269**, 16597–16604.
14. M. Yang et al. (1997) *J. Mol. Biol.* **270**, 696–710.
15. C. G. Lerner, B. T. Stephenson, and R. L. Switzer (1987) *J. Bacteriol.* **169**, 2202–2206.
16. C. L. Quinn, B. T. Stephenson, and R. L. Switzer (1991) *J. Biol. Chem.* **266**, 9113–9127.
17. R. J. Turner, Y. Lu, and R. L. Switzer (1994) *J. Bacteriol.* **176**, 3708–3722.
18. Y. Lu, R. J. Turner, and R. L. Switzer (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14462–14467.
19. P. Babitzke, J. Yealy, and D. Campanelli (1996) *J. Bacteriol.* **178**, 5159–5163.
20. S. Xirasagar, M. B. Elliott, W. Bartolini, P. Gollnick, and P. Gottlieb (1998) *J. Biol. Chem.* **272**, 19863–19869.
21. D. R. Tomchick, R. J. Turner, R. L. Switzer, and J. L. Smith (1998) *Structure* **6**, 337–350.
22. R. P. Novick, S. Iordanescu, S. J. Projan, J. Kornblum, and I. Edelman (1989) *Cell* **59**, 395–404.
23. S. Brantl, E. Birch-Hirschfeld, and D. Behnke (1993) *J. Bacteriol.* **175**, 4052–4061.
24. S. Brantl and E. G. H. Wagner (1994) *EMBO J.* **13**, 3599–3607.
25. D. L. Bently (1995) *Curr. Opin. Genet. Dev.* **5**, 210–216.
26. S. Wright (1993) *Mol. Biol. Cell.* **4**, 661–668.
27. D. Eick and G. W. Bornkamm (1986) *Nucl. Acids. Res.* **14**, 8331–8346.
28. T. K. Kerpola and C. M. Kane (1988) *Mol. Cell Biol.* **8**, 4389–4394.
29. L. London, R. G. Keene, and R. Landick (1991) *Mol. Cell Biol.* **11**, 4599–4615.
30. N. Mechti et al. (1986) *Nucl. Acids. Res.* **14**, 9653–9666.
31. S. Wright, L. F. Mirels, M. Clara, B. Calayag, and J. M. Bishop (1991) *Proc. Natl. Acad. Sci. USA* **88**, 11383–11387.
32. L. Xu, Y. Meng, R. Wallen, and R. A. DePhino (1995) *Oncogene* **11**, 1865–1872.
33. R. Treisman (1986) *Cell* **46**, 567–574.
34. V. Ramamurthy et al. (1990) *Mol. Cell. Biol.* **10**, 1484–1491.
35. K. A. Jones (1997) *Genes and Dev.* **11**, 2593–2599.
36. B. Berkhout, R. H. Silverman, and K. T. Jeang (1989) *Cell* **59**, 273–282.
37. C. Dingwall et al. (1989) *Proc. Natl. Acad. Sci.* **86**, 6925–6929.
38. M. Palangat, T. I. Meier, R. G. Keene, and R. Landick (1998) *Mol. Cell.* **1**, 1033–1042.

### Suggestions for Further Reading

39. P. Babitzke (1997) Regulation of tryptophan biosynthesis: Trp-ing the TRAP or how *Bacillus subtilis* reinvented the wheel. *Mol. Microbiol.* **26**, 1–9. (A recent review covering TRAP-mediated regulation of the *B. subtilis trp* genes.)
40. R. Landick and C. L. Turnbough (1992) "Transcriptional attenuation". In *Transcriptional Regulation*, (S. L. McKnight and K. R. Yamamoto, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 407–446. (An in-depth review of attenuation, particularly amino acid operons in enteric bacteria.)
41. T. M. Henkin (1996) Control of transcription termination in prokaryotes. *Ann. Rev. Gen.* **30**, 35–57. (A review covering control of transcription termination in prokaryotes.)

## Autoantibody

Autoantibody is a very misleading term for certain [antibodies](#), because it covers very different situations. This probably is largely due to Ehrlich's initial proposals, formulated at the beginning of the twentieth century, of the *horror autotoxicus* dogma, leading later to the discrimination between self and nonself (1). If one sticks to this view, one must consider that autoantibodies escape from the dogma and produce a pathological situation, as is the case for [autoimmune diseases](#). Two main observations have led us to revisit this view: one is the existence of [idiotypes](#) and **anti-idiotypes**, and the other is the presence of the so-called natural autoantibodies.

It was proposed by Jerne (2) that [immunoglobulins](#) are organized as a large idiotypic network in which Ig molecules regulate each other by their interactions. Another, more pragmatic approach would be to consider that the millions of different Ig molecules present at any given time in one individual provide a sufficient number of different specificities to recognize any potential epitope, including those expressed by the immunoglobulins themselves, and more generally by any self-**epitope**, unless subject to strict negative selection during differentiation, leading to an absolute tolerance state, which would be in agreement with Ehrlich's concept. Is this the case? During B-cell differentiation, before leaving bone marrow as immature B lymphocytes, newly formed cells are screened against the local environment, and self-aggressive cells are eliminated, providing the major basis for [B-cell](#) tolerance. But this is not an all-or-none phenomenon. It is believed to be relatively flexible, based on the average affinity of [IgM](#) molecules expressed at the surface of immature B cells. This view leaves open the possibility that immature B cells circulating in the periphery have retained some autoreactive potential against a possible self antigen.

The next question is, Are there natural immunoglobulins endowed with properties of recognition of self molecules? The answer is undoubtedly yes and is supported by many observations, especially those from the groups of Avrameas and Coutinho (3). Natural antibodies in whole serum were first described as having the essential characteristic of being polyreactive, a conclusion that was not *a priori* surprising because of their polyclonality. When analysis at the **monoclonal** level was made possible with the [hybridoma](#) technology, it was realized that most natural antibodies had a wide recognition spectrum, as observed from their patterns of binding to many different molecules used for systematic screening. Most natural antibodies were IgM, and [complementary DNA](#) sequencing of their  $V_H$  and  $V_L$  regions revealed that they contained very few mutations, indicating that they reflected [transcription](#) of germline **genes**. Affinity measurements revealed moderate but significant association constants. Natural monoclonal autoantibodies thus seem best defined as primarily germline-encoded and polyreactive. Polyreactivity contains in itself the possibility of having both anti-self and anti-nonsel self specificities. Presumably due to their relatively modest affinities, low concentration, and, of course, the fact that they passed the screen for potentially aggressive anti-self specificities, these antibodies qualify as harmless and may be regarded as the pool of basic circulating immunoglobulins of wide spectrum, endowed with a high connectivity, and constituting the primary germline idiotypic network. By contrast, aggressive autoantibodies found in the serum of patients suffering autoimmune diseases generally contain mutated sequences and are often of isotypes other than IgM.

One may thus consider active immunization as bringing a transient perturbation of the background of the natural equilibrium, leading to the emergence of clones that acquire a greater affinity and an increased specificity through the occurrence of [somatic hypermutations](#).

See also entries [Autoimmunity](#), [Autoimmune Diseases](#), [Idiotypes](#).

## Bibliography

1. P. Ehrlich (1900) The Croonian lecture: On immunity. Proc. Roy. Soc. Lond. (Biol.) **66**, 424.
2. N. K. Jerne (1974) Towards a network theory of the immune system. Ann. Immunol. **125C**, 373–389.
3. A. Coutinho, M. D. Kazatchkine, and S. Avrameas (1995) Natural autoantibodies. Curr. Opin. Immunol. **7**, 812–818.

## Suggestion for Further Reading

4. A. Coutinho, M. D. Kazatchkine, and S. Avrameas (1995) Natural autoantibodies. Curr. Opin. Immunol. **7**, 812–818.

## Autoimmune Diseases

[Autoimmunity](#) is not pathologic *per se* because natural [autoantibodies](#), natural autoreactive [T-cell](#) clones, and the connected [idiotypes](#) and **anti-idiotypes** are part of the basic organization of the immune system. Autoimmune diseases (AIDS) are the consequence of a major failure in the regulation of the immune system, as the emergence of self-reactivity becomes aggressive and harmful to the organism. They represent a major cause of morbidity, and their incidence increases with aging.

The primary causes of autoimmune diseases (AIDS) are still very poorly understood, and the precise target of the autoimmune reactions is not always known. In some cases, one organ, and eventually one target molecule, can be identified. Examples of these are myasthenia gravis, in which antibodies directed against the [acetylcholine receptor](#) have been identified, or autoimmune thyroiditis, where the presence of autoantibodies directed against receptors for thyroid stimulating hormone have been documented. Such antibodies may have very different impacts upon binding. Some may act as a mimic of the hormone and therefore activate thyrocytes, thereby inducing hyperthyroidism, which is characteristic of Graves' disease. Conversely, others will block receptor activation, leading to a hypothyroidism state, as seen in Hashimoto's disease. Another example of an organ-specific AID is insulin-dependent diabetes mellitus, in which the cellular target is the  $\beta$  islets of the endocrine pancreas. Some AIDs appear to be not organ-specific, such as systemic lupus erythematosus, when a number of autoantibodies are directed against a large variety of antigens, including DNA and nucleoproteins.

Several characteristics generally underline the occurrence of AIDs: (a) In most cases there is a genetic predisposition, with a frequent linkage to certain [major histocompatibility complex](#) (MHC) haplotypes; (b) they are chronic diseases that persist for years; and (c) although the molecular target of autoantibodies or autoreactive T-cell clones may be identified, this does not clarify the nature of the initial antigen, if any.

Linkage to an MHC haplotype is indicated by the increase of the relative risk for a given AID to occur. The most spectacular example is that of ankylosing spondylitis, a very severe arthritis that most frequently affects hips, for which 95% of patients are HLA-B27. This indicates a relative risk of about 70 (ie a B27<sup>+</sup> individual has 70 times greater chance to have the disease than the average population). This AID is also particularly interesting because it has been observed that microorganisms such as *Klebsiella* share an identical sequence of six amino acid residues, Gln—Thr

—Asp—Arg—Glu—Asp (QTDRED), with an [epitope](#) of the HLA27 molecule. It was postulated that this AID might be the consequence of an epitope mimicry, an infection with *Klebsiella* inducing an immune response that subsequently would cross-react with the HLA epitope of the host. A similar hypothesis was proposed for rheumatoid arthritis, a very common rheumatological disorder that may develop for years, because it was found to be associated with certain DR4 alleles expressing the Gln—Lys—Arg—Ala—Ala (QKRAA) sequence. A similar sequence was found in Hsp 70, which might provide another example of epitope mimicry as the starter for the initial immunization.

The V<sub>H</sub> and V<sub>L</sub> regions of autoantibodies isolated from AID patients have been sequenced extensively, with the hope of detecting significant differences from natural, nonaggressive autoantibodies. Frequently, dominant expression of some V<sub>H</sub> and/or V<sub>L</sub> is observed, as in the case of anti-DNA autoantibodies isolated from lupus patients. In most occasions, but not always, these antibodies have mutated **variable regions**, as opposed to the germline sequences generally found in natural antibodies. The significance of this remains unclear, however, primarily because the autoantibodies developed by a rather large number of AID patients seem to be only passive witnesses of an autoimmune response, rather than being involved in pathogenesis, which is frequently attributed to [cytotoxic T lymphocytes](#).

Numerous animal models have been defined and extensively studied. Many have stressed the importance of the genetic background, such as the obese strain of chickens that develop a hereditary spontaneous autoimmune thyroiditis that is quite similar to Hashimoto's disease. Another example is the nonobese diabetic (NOD) mice that develop an autoimmune diabetes, which spontaneously develops an insulinitis that becomes full insulin-dependent diabetes mellitus at 7 months of age. New Zealand black (NZB) mice and other strains that spontaneously develop a systemic lupus erythematosus-like syndrome are another example. These models are interesting because they provide a possible key to understanding the genetic basis for the equivalent human diseases. Other approaches are centered more on attempts to isolate [antigens](#) that might induce an AID-like syndrome. An example of this is experimental autoimmune encephalomyelitis, which is a possible model for multiple sclerosis; it can be induced upon injection of myelin basic protein, from which encephalitogenic peptides and precise epitopes have been described. This disease can also be transferred with lymphocytes but not by serum, pointing to the role of cell-mediated immunity in this AID. Interestingly, [T cells](#) with a suppressive activity have also been identified, offering the possibility of modulating the immune system negatively by this approach. Such a treatment is eagerly awaited for most AIDs, particularly multiple sclerosis. Thus far, primarily nonspecific immunosuppressive agents are used, with the obvious problem of generating severe secondary effects.

See also entries [Autoantibody](#) and [Idiotypes](#).

#### Suggestions for Further Reading

- G. Del Prete (1998) The concept of type-1 and type-2 helper T cells and their cytokines in humans. *Int. Rev. Immunol.* **16**, 427–455.
- C. C. Goodnow (1997) Balancing immunity, autoimmunity and self-tolerance. *Ann. NY. Acad. Sci.* **815**, 55–66.
- E. G. Spack (1997) Treatment of autoimmune diseases through manipulation of antigen presentation. *Crit. Rev. Immunol.* **17**, 529–536.

## Autoimmunity



The immune system is able to discriminate between self and nonself, but there are often ambiguities with this formulation. First, it certainly is intended to stress the fact that the immune system is designed for the important role of fighting against external pathogens, whereas it should be harmless for the components of the organism. There are two ways to be harmless. One is simply not to “show off,” the second is to be present, but silent. Immunologists long lived with the idea originally put forward by Ehrlich in 1901 of the *horror autotoxicus*, thus considering that [autoantibodies](#), and more generally [autoimmunity](#), could not exist. This was even strengthened by Burnet, when he proposed the [clonal selection theory](#), based on the elimination of autoreactive lymphocytes (the so-called “forbidden clones”), leading to immune **tolerance**. On these bases, autoimmunity may only be considered to be a pathological deviation, with the emergence of [autoimmune diseases](#). Discovery of the interactions between [idiotype](#) and **anti-idiotype** and of the existence of natural antibodies endowed with auto reactivity made it necessary to reconsider the problem.

A long debate has taken place for years between immunologists, first on the existence, and then on the nature, of natural antibodies and their possible autoimmune properties. Identification of “natural” monoclonal antibodies that were polyspecific clarified the situation, because it was demonstrated that these antibodies contained **variable (V) regions** of both the heavy and the light chains that were mostly not mutated, and like the germline gene from which they were derived. As a result, these antibodies could be considered the first circle of antibody defense, with some polyspecificity that is the best witness of a certain degeneracy of recognition by the immune system. Deliberate immunization will stimulate clones that have the best affinity for the antigen, and this affinity will be amplified by the mechanism of [somatic hypermutation](#), as the specificity becomes more precise. To be polyspecific immediately implies that both auto- and xenospecificities may be found, as was shown using a large panel of random antigens for screening. Autoimmunity has also been described in the [T-cell](#) population and is subject to the same comments, with the exception that T cells do not mutate upon antigenic stimulation.

The entire [repertoire](#) expressed at a given time by the immune system is probably best visualized as a large network of interacting molecules that regulate each other. Deliberate immunization would result in transient perturbation of this equilibrium. Deviation to the production of aggressive B- and T-cell clones would thus be a central alteration of this equilibrium, ensuring some leakiness for the negative selection of potentially aggressive clones.

See also entries [Autoantibody](#), [Autoimmune Diseases](#), [Idiotypes](#).

#### Suggestions for Further Reading

C. Goodnow (1996) Balancing immunity and tolerance: deleting and tuning lymphocyte repertoires. *Proc. Natl Acad. Sci. USA* **93**, 2264–2271.

S. Lacroix-Desmazes, L. Mouthon, S. H. Spalter, S. Kaveri, and M. Kazatchkine (1996) Immunoglobulins and the regulation of autoimmunity through the immune network. *Clin. Exp. Rheumatol.* **14** (Suppl. 15), S9–S15.

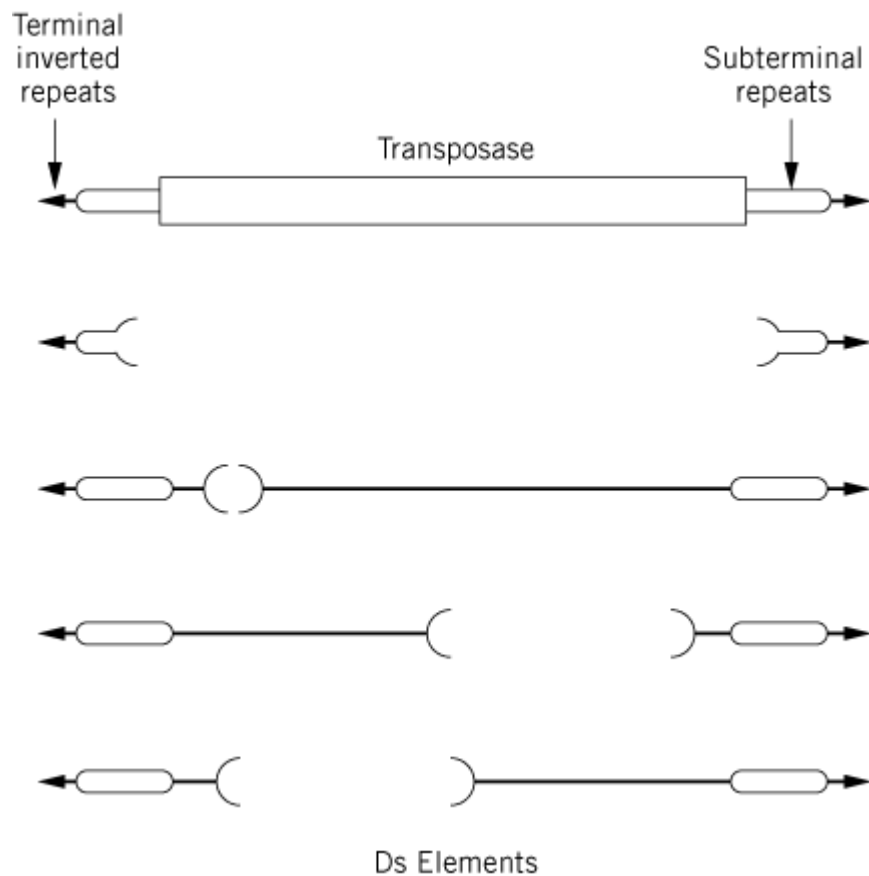
R. J. Cornall, C. C. Goodnow, and J. G. Cyster (1995) The regulation of self-reactive B cells. *Curr. Opin. Immunol.* **7**, 804–811.

#### Autonomous Controlling Elements

[Transposable elements](#) are discrete pieces of DNA that can move between nonhomologous positions in the [genome](#). When Barbara McClintock discovered transposable elements in maize, she called them [controlling elements](#), emphasizing that they could affect gene expression, in addition to being able to move from place to place within the genome (1). She discovered two classes of elements; one she termed *activator* (Ac), which could itself translocate within the genome. She also determined that Ac encoded a product that could promote the movement of a second class of elements called *dissociation* (Ds). In the presence of Ac, Ds could move and cause [chromosome](#) breaks that altered the expression of nearby genes, but was unable to move in the absence of Ac. Ds elements are so-named because their translocation could be associated with chromosome breakage.

We now know that Ac is an intact transposable element encoding a [recombinase](#), a [transposase](#), that promotes element translocation and has special sequences at the termini of the element upon which the transposase acts to promote translocation. Such intact elements are said to be “autonomous” because they can translocate without the assistance of another element. Ds, by contrast, is a nonautonomous element. It doesn't encode a transposase itself and is thus dependent on the presence of other Ac elements to supply a transposase to promote its transposition. Ds is actually an Ac element that has undergone internal deletions that have removed the transposase gene but have left intact the special terminal sequences on which transposase acts intact (Fig. 1). Thus the terminal sequences of Ds are the same as the terminal sequences of Ac, so that the Ac transposase can promote the translocation of Ds (2); there are several different Ds elements that contain different extents of the internal, nonterminal sequences. The deletions are thought to result from abortive repair of the gapped donor site after element excision (2-4).

**Figure 1.** Autonomous and nonautonomous elements. A schematic representation on an intact (autonomous) Ac element is shown (top line). It encodes a transposase formed from several exons (not shown) and special recombination sequences at the ends of the element. The arrows represent short perfect inverted repeats and the oval other repeated sequences; these are not shown to scale. Only some of the subterminal repeats are necessary for transposition. Various deletion derivatives of Ac are shown (lines 2–5); these internal deletions make the element dependent on an intact Ac for transposase and hence for translocation.



Virtually all bacterial elements that have been observed are autonomous; that is, encode a transposase and have cognate terminal sequences. Nonautonomous and autonomous elements are frequently observed in other organisms such as *Drosophila*, [Caenorhabditis elegans](#), and plants (5).

#### Bibliography

1. B. McClintock (1956) Cold Spring Harbor Symp. Quant. Biol. **21**, 197–216.
2. E. Rubin and A. A. Levy (1997) Mol. Cell. Biol. **17**, 6294–6302.
3. E. S. Coen, T. P. Robbins, J. Almeida, A. Hudson, and R. Carpenter (1989) In *Mobile DNA* (D. A. Berg and M. M. Howe, ed.), American Society for Microbiology, Washington, DC, pp. 413–436.
4. W. R. Engels, D. M. Johnson-Schlitz, W. B. Eggleston, and J. Sved (1990) Cell **62**, 515–525.
5. H. Saedler and A. Gierl (1996) Curr. Top. Microbiol. Immunol. **204**, 27–48.

#### Autonomously Replicating Sequences

The [chromosomes](#) of **yeast** contain many genomic segments at which [DNA replication](#) is ordinarily initiated. On average, one is found every 40 kb in the *Saccharomyces cerevisiae* [genome](#). These regions can act as the **origin of replication** for a circular **plasmid**. When such a segment is linked to a nonreplicative DNA fragment, it can confer on that fragment the ability to replicate autonomously

as a plasmid in a yeast cell in the S-phase of the cell cycle. This construction is called an *autonomously replicating fragment* (ARS).

An ARS-containing DNA fragment can **transform** yeast cells at much higher frequencies than a fragment without an ARS, because it can replicate autonomously without undergoing a rare homologous [recombination](#) before its **genes** can be expressed. In contrast to the latter transformants, however, those formed with an ARS-containing DNA are extremely unstable and lose the transforming DNA even under selective conditions. For unknown reasons, they tend to remain associated with the mother cell without being evenly distributed to the daughter cell upon division.

ARS have been found in *Saccharomyces cerevisiae* and but also in other yeasts, such as *Candida maltosa*, *Hansenula polymorpha*, *Kluyveromyces lactis*, and *Schizosaccharomyces pombe*. They are not restricted to yeasts and have also been described in other lower **eukaryotes**, such as *Aspergillus nidulans* and *Entamoeba histolytica*, and even in at or near potential replicative origins from *Drosophila*, Chinese hamster, and humans.

## Autoradiography

Autoradiography is the technique of recording an image of a preparation containing beta-particle emitting [radioactivity](#), using photographic film, X-ray-sensitive film, an emulsion, or other radiation-sensitive medium. The samples are placed directly against the film for a period of time to allow radioactive emissions from the sample to interact with the film emulsion and create an image. A photographic emulsion is a suspension of crystals of silver bromide embedded in gelatin. When crystals of silver bromide are struck by charged-particle or photon radiation, the silver atoms are ionized and form an invisible latent image. After exposure to the sample, the photographic grains in the emulsion are fixed using standard photographic developing, which removes silver bromide that has not been ionized. After the emulsion is developed, each small aggregate of reduced silver atoms becomes a visible dark spot on the emulsion; collectively, they make up the photographic image.

Production of a visible silver grain requires a number of ionization events, so the photographic response is not exactly linear to the amount of radiation present. Preflashing the film with a uniform low intensity of light “primes” each grain of silver to become reduced and visible after absorbing just one or a very few additional beta particles from the sample. This increases substantially the sensitivity of the film and also makes the photographic response more directly proportional to the amount of radiation in the sample. The signal-to-noise ratio is often increased by exposing the film to low temperatures. The sensitivity can also be enhanced by using scintillation screens that emit visible light on encountering a beta particle; the light is recorded by the film (see [Fluorography](#)).

Autoradiography has a large number of practical applications in the biological, chemical, and physical sciences, because it provides both qualitative and quantitative information (eg., images and amounts present). It may be used to image large, small, and microscopic specimens, including sectioned whole organisms, organs, tissues, cellular structures, and nucleic acids that contain some radiolabeled compound. An example of a whole-animal autoradiograph is shown in Figure 1. Cells may be autoradiographed either in culture as a monolayer, on a glass slide, or on thinly sectioned living tissues from an animal organ or tumor. Microautoradiography involves coating the sample directly with a radiation-sensitive emulsion; cellular constituents that have incorporated the radiolabel can be clearly identified. Autoradiography is used with [electrophoresis](#) or [chromatography](#) to image radiolabeled macromolecules and other separated chemicals for quantitative analysis. For example, autoradiography is useful for indicating the position of hybridized nucleic acids on [Southern blots](#) and **Northern blots**, and of proteins on **Western blots**.

**Figure 1.** Whole-body autoradiograph of a rat that had been injected with indium-111 chloride. Courtesy of B. Anders Jönsson, Lund University, Sweden.



## 1. Historical Development

The first autoradiograph, or **autoradiogram**, was made in 1859 by Crookes, who placed uranium rocks on photographic plates (1). Crookes would not understand the process by which the images were created until 37 years later when he visited the laboratory of Henri Becquerel. Becquerel discovered gamma rays from naturally radioactive uranium salts in 1896 when he placed uranium together with photographic plates that had been placed inside black paper to protect them from sunlight (1). This discovery closely followed the discovery of X rays by Wilhelm Roentgen in 1895 using a Crookes electron tube.

The French photographer Nicéphore Niepce had observed the darkening of silver grains on photographic plates by uranium in 1867, but he did not recognize any practical applications of this phenomenon (2). Although George de Hevesy, the pioneer of radioactive tracers, had used bismuth isotopes in animals in 1923, the first autoradiographs of radioactively labeled animal tissues were made in 1924 by Lacassagne in France (1-3). Lacassagne fed polonium-210 to rabbits and later placed thin sections of rabbit kidney tissues in paraffin blocks against photographic plates. Gettler made the first autoradiographs of human tissues from deceased persons who had ingested radium chloride in 1933 (1).

The earliest images were of poor quality. Autoradiographic techniques were improved by Leblond, who, in 1943, showed the microdistribution of iodine-131 in cells (1, 2). In 1946, Bélanger and Leblond (4) developed a method for locating radioactive elements in tissues by covering histological sections with a photographic emulsion. In 1954, Gabriel introduced techniques for identifying radioiodine using human [serum albumin](#) zone paper [electrophoresis](#) (1). In 1953, techniques were developed for labeling nucleic acids and studying kinetics of cell mitosis and division (5). High resolution techniques have more recently been employed for study of other metabolic and pharmacokinetic processes in molecular biology.

## 2. Film-Less Autoradiography

Modern trends in autoradiography involve replacing high speed X-ray film with radiation detector systems, laser scanners, and computer-based imaging systems. A variety of radiation-detecting crystals and phosphors have been developed for this purpose. Storage phosphor screens are more sensitive, by a factor of about 20–100 for beta-emitting radionuclides (2), and they are reusable. The exposure time is also much less, by a factor of about 10, over conventional X-ray film, and samples may be processed at room temperature and without a darkroom or chemicals for film developing. Applications of filmless autoradiography include **two-dimensional gels**, [Southern blots](#), **Northern**

**blots**, immunoblots, and quantitative polymerase chain reaction (**PCR**) (2).

Microchannel array detectors have been introduced to replace both X-ray films and phosphor screens (8). The new instruments are faster (by a factor of ~10) and have greater image resolution than do phosphor screens for detecting latent images from hybridization studies using macromolecules labeled with carbon-14, sulfur-35, phosphorus-32, and iodine-125 from flat gels, blots, membranes, tissue slices, and other flat specimens.

#### Bibliography

1. M. Brucer (1990) *A Chronology of Nuclear Medicine*, Heritage Publications Inc., St. Louis, MO.
2. R. Wegmann, N. Balmain, S. Ricard-Blum, and S. Guha (1995) *Cell. Mol. Biol.* **41**, 1–20.
3. A. Lacassagne, J. Lattes, and J. Lavedan (1925) *J. Radiol. Electr.* **9**, 1–14.
4. L.-F. Bélanger and C. P. Leblond (1946) *Endocrinology* **39**, 386–400.
5. A. Howard and S. R. Pelc (1970) *Heredity* **6** (suppl.), 261–273.
6. J. G. Gall and M. Pardue (1969) *Proc. Natl. Acad. Sci. USA* **63**, 378–383.
7. E. M. Southern (1975) *J. Mol. Biol.* **98**, 503–517.
8. Packard Instrument Company (1993) *Enter the New Era of Instant Autoradiography: The InstantImager™*, promotional literature from Packard Instrument Company, Meriden, Conn.

#### Suggestions for Further Reading

9. G. A. Boyd (1955) *Autoradiography in Biology and Medicine*, Academic Press, New York.
10. E. J. Hall (1994) *Radiobiology for the Radiologist*, 4th ed., J. B. Lippincott Company, Philadelphia, pp. 92–95, 397–402.
11. J. C. Kaplan and M. Delpech (1989) ““Biologie Moléculaire et Médecine””, in *Médecine et Sciences Flammarion*, Paris, p. 610.
12. D. LeGuellec and J. -Y. Sire (1991) in *Microscopie électronique. Cryométhodes, Immunocytologie, Autoradiographie, Hybridation in situ* (G. Morel, ed.), INSERM, pp. 299–304.
13. R. Wegmann, N. Balmain, S. Ricard-Blum, and S. Guha (1995) *Cell. Mol. Biol.* **41**, 1–20.

#### Autosome

An autosome is any [chromosome](#) other than the [sex](#) chromosomes. For example, the human [genome](#) has 24 chromosomes, including 22 autosomes plus the [X-](#) and the [Y-chromosomes](#). This is described as the [haploid](#) set. If a cell contains a complete haploid set of chromosomes, it is described as [euploid](#). If the chromosome set is altered by duplication or deletion then the cell is **aneuploid**. If a cell contains two sets of genomes per nucleus, it is **diploid**, and if it contains more than two, it is described as [polyploid](#).

Organisms that have a largely vegetative existence, such as some **fungi** and **algae**, have mainly haploid cells. In metazoans the **gametes**, such as the **spermatozoa** and **oocytes** in humans, are haploid. Sexual **reproduction** involves the fusion of two haploid gametes to form a diploid [egg](#) in an animal. The vast majority of the cells in an animal that develop from an egg are diploid. Thus there are two active copies of each autosome in most animal cells. In **plants**, a polyploid state is common

where the number of genomic copies in somatic cells may range from 3 to 10 (1). The number of genomic copies can actually be **polymorphic** within a particular species (2). Thus plants are relatively insensitive to the number of copies of a particular **gene** in a given cell. This is not the case with animals, where somatic cells are almost always diploid. The only exceptions are those animals with **parthenogenetic** life styles in which either males or females develop as a default state from an unfertilized egg. In general, animal cells cannot accommodate more than two copies of a particular chromosomal region.

Aneuploidy occurs when parts of chromosomes or whole chromosomes are absent from the genome. In humans the nondisjunction of homologous chromosomes during **meiosis** or of sister **chromatids** during **mitosis** can lead to aneuploidy. When a chromosome is in excess of the normal euploid number, the condition is called hyperploidy. For example, a hyperploid state in humans occurs when there are three copies of a particular chromosome, a condition known as **trisomy**. All human autosomes except chromosome 1 have been found as trisomies (3). Most trisomic embryos are spontaneously aborted, but live births occur with trisomies of chromosomes 13, 18, and 21 (4). The most common human trisomy to reach term is that of chromosome 21, resulting in Down's syndrome. These children are mentally retarded, predisposed to cancer and have a shortened life span. Presumably the deleterious effects of three copies of chromosome 21 follow from an imbalance in the expression of one or several sets of genes on chromosome 21 relative to the other genes on the other chromosomes in the cell.

Experiments in *Drosophila melanogaster* indicate that changes in the dosage of a single gene might be tolerated, but a deletion that includes several genes or a large section of a chromosome is likely to cause sterility (5). Surprisingly in trisomies of *D. melanogaster* involving whole chromosome arms, many of the duplicated autosomal genes express only diploid levels of gene product. The expression of genes outside the duplicated region is also affected, however, indicating the pleiotropic effects of altering gene copy number. An interesting aspect of sexual reproduction requiring the use of two distinct haploid gametes is that there are two copies of each autosome in a diploid cell, whereas there are one or two copies of each sex chromosome. **Dosage compensation** is the phenomenon whereby the **transcriptional** activity of genes on the sex chromosome is adjusted to similar levels independent of the number of copies, so that it is constant relative to the activity of a gene on an autosome (see **X-Chromosome Inactivation**).

Humans, again have two X-chromosomes in one sex (female) and one X-chromosome in the other (male). To achieve dosage compensation, an entire X-chromosome is silenced in female diploid cells. This process is called X-chromosome inactivation and results in the formation of the **Barr body**. The capacity to monitor the activity of genes on autosomes relative to the sex chromosomes has facilitated the discovery of novel mechanisms of regulating genes at the level of whole chromosomes.

#### Bibliography

1. J. E. Averett (1980) "Polyploidy in plant taxa: summary". In *Polyploidy* (W. H. Lewis, ed.), Plenum, New York.
2. K. Irifune (1990) *J. Sci. Hiroshima Univ.* **23**, 163–181.
3. T. J. Hassold (1986) *Trends Genet.* **2**, 105–109.
4. E. B. Hook, B. B. Topol, and P. K. Cross (1989) *Am. J. Hum. Genet.* **45**, 855–861.
5. D. L. Lindsley et al. (1972) *Genetics* **71**, 157–184.

#### Suggestion for Further Reading

6. M. S. Clark and W. J. Wall (1996) *Chromosomes. The Complex Code*, Chapman and Hall, London.

## Auxins

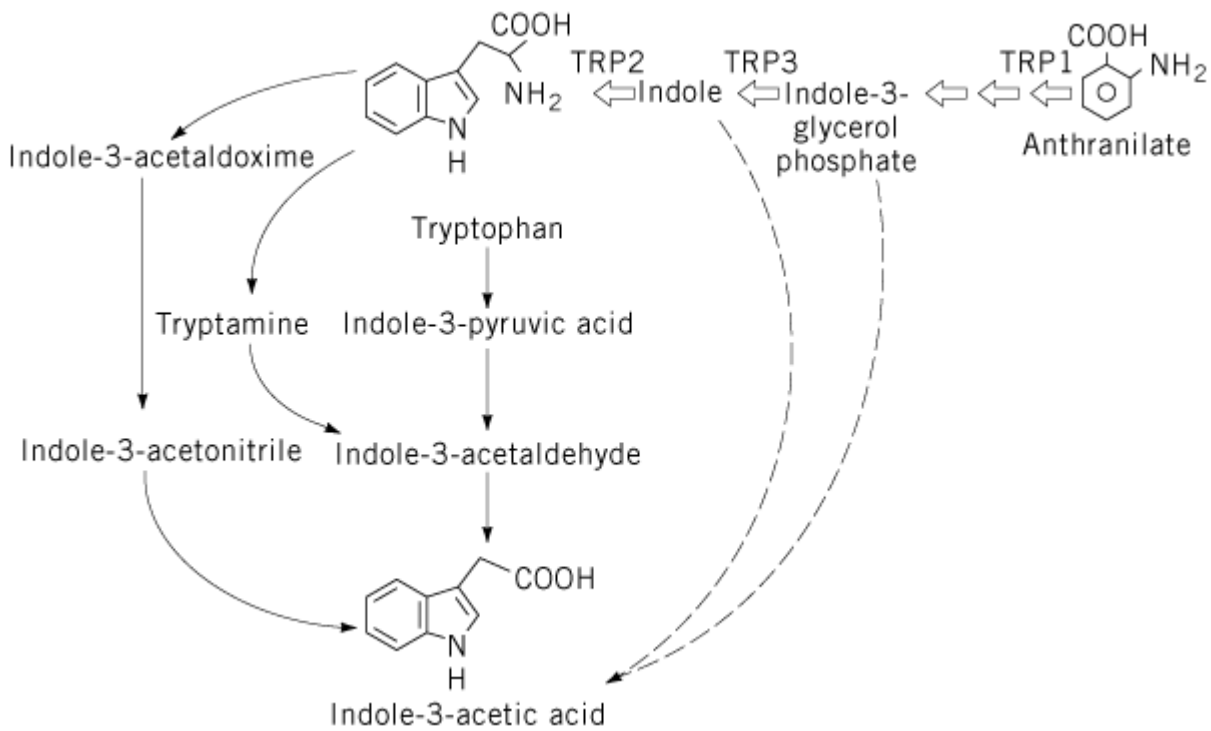
### 1. History

Auxins were the first class of [plant hormones](#) to be discovered. Initial observations by Charles Darwin and his son Francis on the phenomenon of bending of canary grass coleoptiles toward a unilateral source of light (known as phototropism) led the foundation for a series of experiments that culminated in the discovery of the first phytohormone (1). The active substance was isolated from oat coleoptiles by Went in 1926 (2) and later identified and named auxin (from the Greek “to increase”). By analogy with the animal hormone concept from Bayliss and Starling (3), plant physiologists interpreted auxin as a phytohormone.

### 2. Biosynthesis and Metabolism

The most abundant natural auxin is indole-3-acetic acid (IAA). It is mainly synthesized in the shoot apex and in young leaves. Other substances with auxin activity that are found in plants are indole-3-butyric acid, 4-chloro-indole-acetic acid, and phenyl acetic acid (4, 5); these compounds are active only at higher concentrations, and their roles in development remain largely unknown. Multiple pathways exist for auxin biosynthesis; some of these are dependent on tryptophan, while others are independent (Fig. 1). In certain species, different routes are operational (6), whereas in others, essentially all of the IAA is derived from tryptophan (7). Three possible routes for tryptophan-dependent IAA synthesis have been proposed (4, 5), which include indole-3-acetonitrile, tryptamine, and indole-3-pyruvic acid as the respective intermediates. Genes encoding tryptophan decarboxylase (which catalyzes the conversion of tryptophan to tryptamine) and nitrilase (forming IAA from indole-3-acetonitrile) have been cloned (8-10). The route via indole-3-acetonitrile has been found primarily in Brassicaceae. Likewise, the tryptamine pathway is not supposed to be of general importance, because tryptamine is not universally present in plants. A unique tryptophan-dependent pathway via indole-3-acetamide is used by bacteria and by plant cells that are transformed with *Agrobacterium tumefaciens*. The corresponding genes have been used to alter the IAA levels in **transgenic** plants (11).





Tryptophan-auxotroph mutants have revealed the existence of a tryptophan-independent pathway for IAA biosynthesis. Both maize (12) and *Arabidopsis* (13) tryptophan mutants overproduce IAA or IAA conjugates, respectively. On the other hand, mutants that block tryptophan biosynthesis before indole-3-glycerol phosphate have a **phenotype** suggestive of auxin deficiency (14). Based on these observations, IAA was proposed to be synthesized through a tryptophan precursor, either indole or indole-glycerol phosphate, thereby bypassing tryptophan (4). An auxin-overproducing mutant that accumulates both free and conjugated auxins has been isolated (15, 16). The wild-type gene product might inhibit auxin biosynthesis and shows similarity to aminotransferases (17).

Of all the auxin present in plant tissues, 95% resides in conjugated form, coupled to [amino acids](#), [peptides](#), or [proteins](#) via amide bonds, or to sugars via ester linkages (18). The role of these conjugates is diverse: the storage of free IAA (released by hydrolysis), its transport, and eventually its catabolism. A genetic approach was used to characterize a hydrolase that catalyzes the formation of free IAA from IAA-Leu (10). The *ILR1* (*IAA-leucine resistant 1*) **gene** encodes a protein with similarity to a bacterial aminoacylase that catalyzes amide cleavage.

### 3. Transport

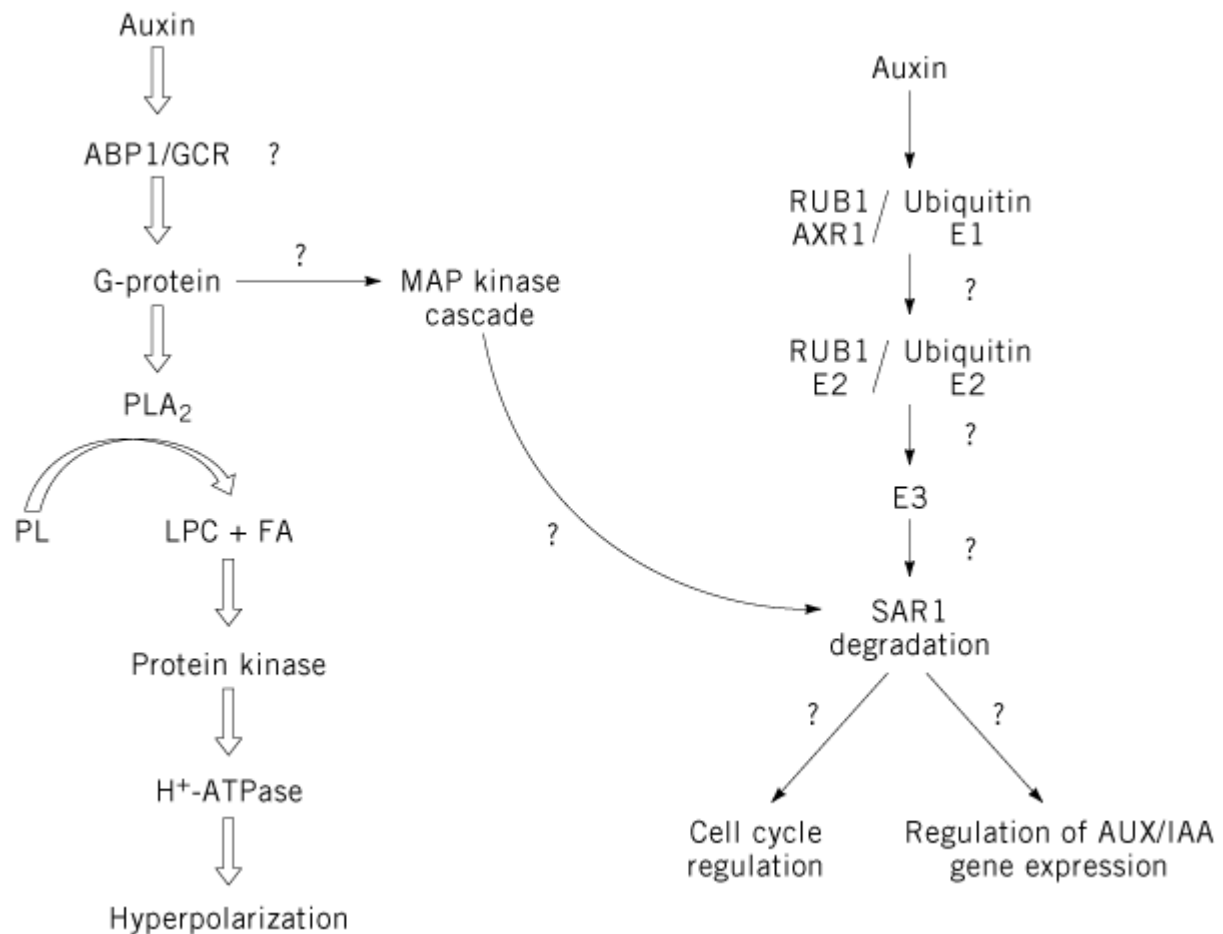
Another important aspect determining auxin levels is transport of the hormone. Auxin transport can be either (a) polar, basipetally directed, and active or (b) nonpolar and inactive (19). The latter occurs through phloem tissue. Polar transport is probably involved specifically in processes of apical dominance, vascular differentiation, and shoot–root communication. Evidence exists for the presence of influx and efflux carriers, although protonated auxin can enter cells by simple [diffusion](#). Proton–auxin **symport**, energized by the plasma membrane **proton pump**, allows a more rapid uptake than by diffusion. The AUX1 protein has been suggested to represent an auxin uptake carrier (20), based on its similarity to amino acid uptake carriers, which also function as symporters. Although several mutants have been characterized that show growth defects similar to those upon treatment with polar transport inhibitors (21), only one of them appears to be a candidate for the auxin efflux carrier. The PIN1 (*pin formed 1*) gene product shows the characteristics of a transmembrane protein (22). The corresponding mutant displays abnormal cotyledon development, together with twisted, narrow leaves with branching midveins. Auxin transport is reduced at all stages of development (23).

Recently, factors that may modulate auxin transport have been characterized (24, 25). Both the loci *RCN1* (*ROOT CURLING on NPA1*) (coding for a protein phosphatase A regulatory subunit) and *MPI* (*MONOPTEROS1*) (coding for a homologue of the auxin response factor ARF1) (26) (see text below) have been isolated based on their altered response to auxin transport inhibitors. Mutations in the *MPI* gene interfere with vascular strand formation and with initiation of the body axis in the embryo (27).

#### 4. Signal Perception and Transduction

The perception and transduction of the auxin signal remains largely enigmatic (Fig. 2). Two different approaches have been followed toward the identification of putative auxin receptors: (a) the characterization of auxin-resistant mutants and (b) the isolation of auxin-binding proteins (28, 29). Because auxin is present both intra- and extracellularly, a site of perception may be needed on both sides of the membrane. Although many auxin-binding proteins have been identified, none has been proven to function as an internal receptor at the molecular level. In contrast, an external receptor has been characterized. There is evidence that auxin causes changes in the plasma membrane potential via stimulation of the proton *ATPase* and modulation of **ion channels**. These rapid responses result from auxin binding to a protein named ABP1 (auxin-binding protein 1). Thus far, it is unknown whether ABP1 mediates cell growth. Although ABP1 possesses a C-terminal signal for retention in the endoplasmic reticulum (ER) and largely resides in the ER, it does not seem to bind auxin in this cellular compartment. Instead, ABP1 can reach the cell surface by a yet unknown mechanism, as demonstrated by immunolocalization experiments (30). It is hypothesized that ABP1 is retained at the membrane periphery by interaction with a docking protein, because ABP1 does not contain a potential membrane-spanning domain. The nature of the docking protein is undefined, but biochemical evidence supports the concept of a **G-protein-coupled receptor (GCR)** (31). These receptors contain seven transmembrane domains and are directly coupled to heterotrimeric G proteins. A possible effector could be phospholipase  $A_2$ , which produces lysophospholipids and free **fatty acids**. The former have been proposed to act as second messengers that can activate the plasma membrane proton *ATPase* (32, 33). However, G-protein-coupled receptors may activate different G proteins in different cell types, as has been shown in animal cells (34), and therefore may control a downstream branched transduction pathway. At least one report suggests the involvement of a **MAP kinase** cascade in the initiation of cell division by auxin in cell cultures (35).

**Figure 2.** A model for auxin signal transduction. PL, phospholipids;  $PLA_2$ , phospholipase  $A_2$ ; LPC, lysophospholipids; FA, fatty acids. (Modified from Refs. 29 and 39.)



The second approach to molecular characterization of auxin signal transduction involves the isolation and characterization of auxin-resistant mutants. At least seven loci were identified by screening for resistance to auxin-induced root growth inhibition (21). AXR1 (auxin resistant 1) probably plays a role at an early stage in the signaling chain, because the recessive *axr1* mutant cannot respond to auxin in all the assays tested. The *AXR1* gene encodes a protein with homology to the amino-terminal half of the **ubiquitin**-activating enzyme E1 (36). Thus, a central part of the mechanism of auxin action appears to play a regulatory role in protein degradation. Recently, two genes of *Saccharomyces cerevisiae* were found to be similar to *AXR1* (37). One of these regulates the stability of key components of **cell cycle** regulation, such as the **cyclin**-dependent kinase inhibitor Sic1p, through conjugation of the ubiquitin-like protein RUB1 to CDC53p, a component of a E3 ubiquitin—protein ligase complex (37). The speculation that *AXR1* plays a role in cell cycle control is considerably strengthened by the characterization of the *tir1* (*transport inhibitor response 1*) locus, a semidominant mutation that also confers auxin resistance (38). Sequence analysis revealed similarity between TIR1 and F-box proteins found in E3 ubiquitin-protein ligase complexes in yeast. The product of the *SAR1* (*suppressor of auxin resistance*) gene has been proposed as a candidate target for AXR1-dependent ubiquitin degradation (39). The *sar1* mutant was identified as a suppressor of the *Axr1* phenotype (40). SAR1 may function as a repressor of *AUX/IAA* genes, or it may be involved in cell cycle regulation (40). The mostly root-specific *axr4* mutation probably defines an AXR1-independent path. Double-mutant analysis indicates that at least two interacting pathways exist, in which AXR4 and AUX1 define a first route while AXR1 defines a second one (41).

## 5. Downstream Targets

Ultimately, the auxin signal leads to induction or repression of a series of target genes, some of

which are responsive in a few minutes (42, 43). The products of these primary response genes are likely candidates to play a role in auxin-stimulated growth. The induction of the early genes is independent of *de novo* **protein biosynthesis**, implying that **post-translational processes** control the activity of factors that affect their expression. The transcriptional activation of early genes is based on the presence of auxin-responsive, **cis-acting, response elements** (AuxREs). Two auxin response-element domains were defined in the *Pisum sativum* *IAA4* and *IAA5* genes, one of which acts as an auxin switch while the other acts as an **enhancer** element. The former AuxRE is present in the promoter of *AUX/IAA*, *SAUR* (*small auxin-upregulated mRNAs*), and auxin-inducible 1-aminocyclopropane-1-carboxylate synthase genes. An alternative AuxRE is found in the promoter of the *GH2* and *GH4* genes. Essentially, this element consists of a G-box motif, binds bZIP factors, and is generally responsive to stress signals. The biological role of most of the auxin early-response gene products remains to be elucidated. Nevertheless, based on similarities with proteins of known functions, it is proposed that *GH2* and *GH4* genes are **glutathione-S-transferases**. *AUX/IAA* polypeptide chains are short-lived nuclear proteins that probably function as an auxin switch in a secondary phase of signal transduction (42), implying that their target gene products probably mediate auxin responses. Thus, the *AUX/IAA* polypeptides might be involved in homo- and heterodimerization with auxin response factors, such as ARF1, that bind to AuxREs (26, 44). The *IAA17* and *AXR3* genes are identical (44). *Axr3* is a dominant auxin-resistant mutant with extremely resistant roots, although the general phenotype of the mutant resembles auxin hypersensitivity (45).

## 6. Effects

Depending on temporal and spatial factors, a particular subset of downstream target genes is activated and results in one of many described effects. Auxin has been implicated in a variety of growth and developmental processes, such as cell division, stem and root elongation, leaf expansion, formation of adventitious and lateral roots, induction of vascular tissue, apical dominance, tropisms, flowering, and fertility (46).

## Bibliography

1. C. Darwin and F. Darwin (1881) *The Power of Movement in Plants*, Appleton-Century-Crofts, New York.
2. F. W. Went (1927) Proc. Kon. Nederl. Akad. Wetensch. **30**, 10–19.
3. W. M. Bayliss and E. H. Starling (1902) Proc. R. Soc. **69**, 352–353.
4. J. Normanly, J. P. Slovin, and J. D. Cohen (1995) Plant Physiol. **107**, 323–329.
5. B. Bartel (1997) Annu. Rev. Plant Physiol. Plant Mol. Biol. **48**, 51–66.
6. L. Michalczyk, T. J. Cooke, and J. D. Cohen (1992) Phytochemistry **31**, 1097–1103.
7. K. Bialek, L. Michalczyk, and J. D. Cohen (1992) Plant Physiol. **100**, 509–517.
8. V. De Luca, C. Marineau, and N. Brisson (1989) Proc. Natl. Acad. Sci. USA **86**, 2582–2586.
9. D. Bartling, M. Seedorf, A. Mithöfer, and E. W. Weiler (1992) Eur. J. Biochem. **205**, 417–424.
10. B. Bartel and G. R. Fink (1995) Science **268**, 1745–1748.
11. H. J. Klee and C. P. Romano (1994) Crit. Rev. Plant Sci. **13**, 311–324.
12. A. D. Wright, M. B. Sampson, M. G. Neuffer, L. Michalczyk, J. P. Slovin, and J. D. Cohen (1991) Science **254**, 998–1000.
13. R. L. Last, P. H. Bissinger, D. J. Mahoney, E. R. Radwanski, and G. R. Fink (1991) Plant Cell **3**, 345–358.
14. R. L. Last and G. R. Fink (1988) Science **240**, 305–310.
15. W. Boerjan, M.-T. Cervera, M. Delarue, T. Beeckman, W. Dewitte, C. Bellini, M. Caboche, H. Van Onckelen, M. Van Montagu, and D. Inzé (1995) Plant Cell **7**, 1405–1419.
16. J. J. King, D. P. Stimart, R. H. Fisher, and A. B. Bleecker (1995) Plant Cell **7**, 2023–2037.
17. M. Gopalraj, T.-S. Tseng, and N. Olszewski (1996) Plant Physiol. **111** (Suppl.), 114 (469).
18. J. Normanly (1997) Physiol. Plant. **100**, 431–442.

19. T. L. Lomax, G. K. Muday, and P. H. Rubery (1995) In *Plant Hormones* (P. J. Davies, ed.), Kluwer Dordrecht, pp. 509–530.
20. M. J. Bennett, A. Marchant, S. T. May, H. G. Green, S. P. Ward, P. A. Millner, A. R. Walker, B. Schulz, and K. A. Feldmann (1996) *Science* **273**, 948–950.
21. O. Leyser (1997) *Physiol. Plant.* **100**, 407–414.
22. L. Galweiler, E. Wisman, A. Yephremov, and K. Palme (1996) P307
23. K. Okada, J. Ueda, M. K. Komaki, C. J. Bell, and Y. Shimura (1991) *Plant Cell* **3**, 677–684.
24. C. Garbers, A. DeLong, J. Deruère, P. Bernasconi, and D. Söll (1996) *EMBO J.* **15**, 2115–2124.
25. C. S. Hardtke and T. Berleth (1998) *EMBO J.* **17**, 1405–1411.
26. T. Ulmasov, G. Hagen, and T. J. Guilfoyle (1997) *Science* **276**, 1865–1868.
27. G. K. H. Przemeck, J. Mattsson, C. S. Hardtke, Z. R. Sung, and T. Berleth (1996) *Planta* **200**, 229–237.
28. H. Barbier-Brygoo (1995) *Crit. Rev. Plant Sci.* **14**, 1–25.
29. H. Macdonald (1997) *Physiol. Plant.* **100**, 423–430.
30. J. Deikman and M. Ulrich (1995) *Planta* **195**, 440–449.
31. P. A. Millner, D. A. Groarke, and I. R. White (1996) *Plant Growth Regul.* **18**, 143–147.
32. G. F. E. Scherer and B. André (1993) *Planta* **191**, 515–523.
33. H. Yi, D. Park, and Y. Lee (1996) *Physiol. Plant.* **96**, 359–368.
34. E. J. Neer (1995) *Cell* **80**, 249–257.
35. T. Mizoguchi, Y. Gotoh, E. Nishida, K. Yamaguchi-Shinozaki, N. Hayashida, T. Iwasaki, H. Kamada, and K. Shinozaki (1994) *Plant J.* **5**, 111–122.
36. H. M. O. Leyser, C. A. Lincoln, C. Timpte, D. Lammer, J. Turner, and M. Estelle (1993) *Nature* **364**, 161–164.
37. D. Lammer, N. Mathias, J. M. Laplaza, W. Jiang, Y. Liu, J. Callis, M. Goebel, and M. Estelle (1998) *Genes Dev.* **12**, 914–926.
38. M. Ruegger, E. Dewey, W. M. Gray, L. Hobbie, J. Turner, and M. Estelle (1998) *Genes Dev.* **12**, 198–207.
39. W. M. Gray and M. Estelle (1998) *Curr. Opin. Biotechnol.* **9**, 196–201.
40. A. Cernac, C. Lincoln, D. Lammer, and M. Estelle (1997) *Development* **124**, 1583–1591.
41. C. Timpte, C. Lincoln, F. B. Pickett, J. Turner, and M. Estelle (1995) *Plant J.* **8**, 561–569.
42. S. Abel and A. Theologis (1996) *Plant Physiol.* **111**, 9–17.
43. Y. Takahashi, S. Ishida, and T. Nagata (1995) *Plant Cell Physiol.* **36**, 383–390.
44. D. Rouse, P. Mackay, P. Stirnberg, M. Estelle, and O. Leyser (1998) *Science* **279**, 1371–1373.
45. H. M. O. Leyser, F. B. Pickett, S. Dharmasiri, and M. Estelle (1996) *Plant J.* **10**, 403–413.
46. P. J. Davies (1995) *Plant Hormones: Physiology, Biochemistry and Molecular Biology*, Kluwer, Dordrecht, The Netherlands.

## Auxotroph

An auxotroph is a [mutant](#) that requires a factor for growth that is not required by the **wild-type**. The opposite of an auxotroph is a **prototroph**, an organism that grows in minimal medium containing only inorganic salts and a carbon source (plus an additional energy source if the organism is an

**autotroph**). Typical auxotrophic mutations are in genes encoding biosynthetic [enzymes](#) for [amino acids](#), vitamins, and **nucleotides**. Some organisms are natural auxotrophs. For example, humans require that all vitamins and several amino acids be supplied in their food. In contrast many, but not all, free-living **microorganisms**, such as *Escherichia coli*, are prototrophs. A specific auxotrophic **phenotype** is called an auxotrophy, eg, an arginine auxotrophy, and the mutant is called an arginine auxotroph.

## Avidin

Avidin is a minor **glycoprotein** component of egg white that binds the vitamin [biotin](#) with the largest **association constant** known ( $K_a \sim 10^{15} \text{ M}^{-1}$ ) for a **ligand**–protein interaction ([1](#)). Its function in egg white is not known precisely, although protective or antibiotic activities have been suggested. It gained importance not because of its natural functionality, but because of its utility as a tool in molecular biology [see [Avidin-Biotin System](#)]. Together with its bacterial relative [streptavidin](#) and in complex with biotin, avidin has become one of the most important tools in molecular biology because of its myriad of applications in many fields of biology, biotechnology, and even chemistry ([2](#)). Originally, the avidin-biotin complex was introduced to localize, detect, and isolate biologically active molecules ([3](#), [4](#)), but over the years it has become a standard system for replacing radioactivity in [immunoassays](#), for DNA probes, and for various clinical and medical applications ([5](#), [6](#)).

The avidin molecule forms a very stable tetramer that consists of four identical monomers. The monomer is glycosylated at asparagine-17 and is strongly basic ([isoelectric point](#) of about 10.5) due to a surplus of [arginine](#) and [lysine](#) residues. The relatively high molar **absorption** ( $A^{1\%}_{[280 \text{ nm}]} = 1.54$ ) reflects the four [tryptophan](#) and one [tyrosine](#) residue per monomer. The molecular mass of the glycosylated tetramer is about 62,400 Da. When the oligosaccharide residue is removed, the estimated mass is 57,120 Da.

Each avidin monomer binds one molecule of biotin. The binding is considered to be kinetically noncooperative among the subunits ([7](#)), although one of the tryptophan residues from a neighboring monomer is inserted physically into the binding site and thus contributes to binding biotin and to the enhanced stability of the tetrameric structure ([8](#)). Avidin also binds to the dye HABA (4'-hydroxyazobenzene-2-carboxylic acid), but the association constant (ca.  $10^6 \text{ M}^{-1}$ ) is considerably lower than that of the avidin-biotin complex.

The unusually strong binding and the many uses of the system caused scientists to try to understand the molecular basis for this interaction. Moreover, since its inception, many groups have sought to improve various components of the system or the system as a whole. Concurrently, structure-function studies have been implemented, in which early results show the importance of aromatic amino acids in the biotin-binding sites of both avidin and streptavidin. The final structural proof is the determination of the crystal structure of both proteins (Fig. [1](#)), which shows that their overall fold, organization into the tetramer, and content and position of amino acids in the binding sites are all similar ([8-10](#)).

**Figure 1.** Ribbon diagrams of the avidin–and streptavidin–biotin monomers, showing the eight strands of the b-barrel. The biotin molecule is shown in a ball-and-stick model.



Avidin Monomer



Streptavidin Monomer

The avidin monomer forms a  $\beta$ -barrel, composed of eight antiparallel [b-strands](#), connected sequentially by loops. The biotin-binding site is a **hydrophobic** pocket that, in the absence of biotin, contains five molecules of structured [water](#). The combined structure of these water molecules emulates that of biotin itself, and as a consequence, they prevent the binding pocket from collapsing. The water molecules are expelled upon binding of biotin. Because no conformational change occurs during complex formation, the binding has a gain of **free energy**. In fact, the avidin-biotin system can be considered a natural host-guest complex. It is clear that any change of this “ideal system” leads to reduced affinity. Indeed, [site-directed mutagenesis](#) or [chemical modification](#) of residues in the binding site have reduced the affinity constant. An example is the [nitration](#) of the essential binding-site tyrosine ([11](#)), which decreases the association constant to  $\sim 10^{10} \text{ M}^{-1}$ , and renders the binding reversible at alkaline pH. Another example is the site-directed modification of tryptophan 120 in streptavidin to phenylalanine, which is also accompanied by a loss in affinity ([12](#)).

Different modifications of avidin were introduced over the years, including Extravidin, Neutralite avidin, and Lite avidin ([13](#)). Extravidin and Neutralite avidin are avidins in which the arginine residues were modified, whereas the oligosaccharide moieties were removed from Lite and Neutralite avidins. These derivatives are useful in many applications, because avidin, being a basic glycoprotein, is subject to a variety of extraneous nonspecific or undesired binding, particularly to DNA. Similar qualities of genetically engineered forms of avidin are possible because avidin was successfully expressed in a eukaryotic host ([14](#), [15](#)). This should enable preparation of different derivatives of avidin with high affinities, provided that the binding site residues are preserved. It seems that the nonbinding-site residues can be altered extensively.

### Bibliography

1. N. M. Green (1975) *Adv. Protein Chem.* **29**, 85–133.
2. M. Wilchek and E. A. Bayer (1989) *Trends Biochem. Sci.* **14**, 408–412.
3. E. A. Bayer and M. Wilchek (1978) *Trends Biochem. Sci.* **3**, N237–N239.
4. E. A. Bayer and M. Wilchek (1980) *Methods Biochem. Anal.* **26**, 1–45.
5. M. Wilchek and E. A. Bayer (1984) *Immunol. Today* **5**, 39–43.
6. M. Wilchek and E. A. Bayer (1988) *Anal. Biochem.* **171**, 1–32.
7. M. L. Jones and G. P. Kurzban (1995) *Biochemistry* **34**, 11750–11756.
8. O. Livnah, E. A. Bayer, M. Wilchek, and J. L. Sussman (1993) *Proc. Natl. Acad. Sci. USA* **90**,

5076–5080.

9. W. A. Hendrickson, A. Pähler, J. L. Smith, Y. Satow, E. A. Merritt, and R. P. Phizackerley (1989) *Proc. Natl. Acad. Sci. USA* **86**, 2190–2194.
10. P. C. Weber, D. H. Ohlendorf, J. L. Wendoloski, and F. R. Salemme (1989) *Science* **243**, 85–88.
11. E. Morag, E. A. Bayer, and M. Wilchek (1996) *Biochem. J.* **316**, 193–199.
12. T. Sano and C. R. Cantor (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3180–3184.
13. E. A. Bayer and M. Wilchek (1994) In *Egg Uses and Processing Technologies* (J. S. Sim, and S. Nakai, eds.), CAB International, Wallingford, UK, pp. 158–176.
14. K. J. Airene, P. Sarkkinen, E.-L. Punnonen, and M. S. Kulomaa (1994) *Gene* **144**, 75–80.
15. K. J. Airene, C. Oker-Blom, V. S. Marjomäki, E. A. Bayer, M. Wilchek, and M. S. Kulomaa (1997) *Protein Express. Purif.* **9**, 100–108.

### Suggestions for Further Reading

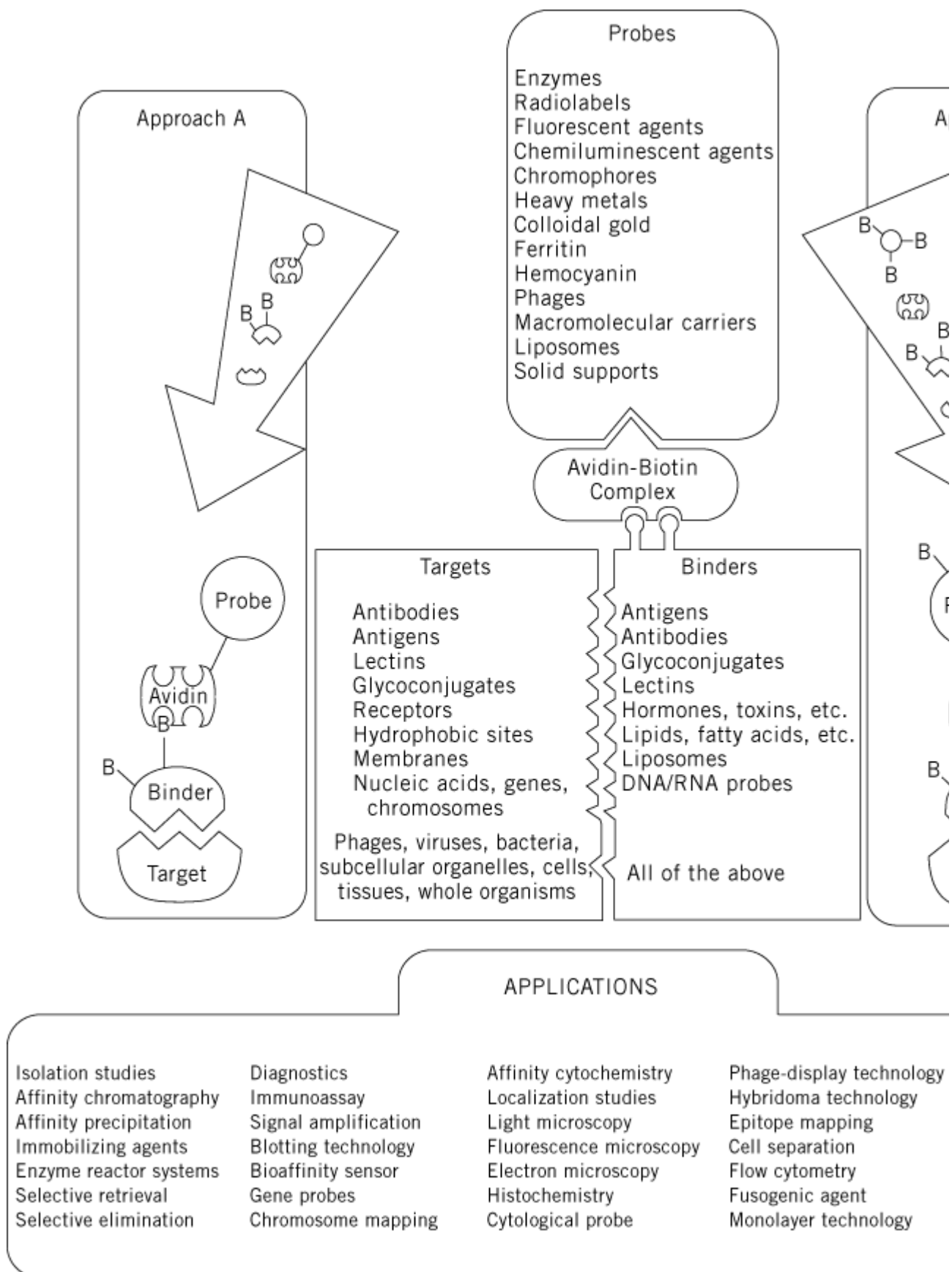
16. M. Wilchek and E. A. Bayer, eds. (1990): *Avidin-Biotin Technology*, *Methods Enzymol.*, Vol. **184**, Academic Press, San Diego, p. 746.
17. N. M. Green (1990) Avidin and streptavidin, *Methods Enzymol.* **184**, 51–67.

## Avidin-Biotin System

The avidin-biotin system has become one of the major mainstays in biochemical analysis and has far-reaching application in biotechnology, industry, and clinical medicine ([1](#)). The general idea of the avidin-biotin system is that [biotin](#), a low molecular weight vitamin, can be chemically coupled to other low or high molecular weight molecules (e.g., [proteins](#), [hormones](#), **DNA** molecules, etc.). The biotin moiety is still recognized by [avidin](#) or [streptavidin](#), either as the native protein or in derivatized form containing any one of a number of [reporter groups](#) or probes. The principle of this system is illustrated in [Fig. 1](#).

**Figure 1.** Overview of the avidin-biotin system and the two major strategies for the various applications. In both approaches, the desired experimental system is combined with a biotinylated binder molecule. Approach A involves direct interaction with the probe. In approach B, avidin is a sandwich between the biotinylated binder and the biotinylated probe. Various targets, and the applications are listed.





The avidin-biotin system has become a “universal” tool in most of the fields of the biological sciences, because of studies that commenced in the mid 1970s and the constant development that has

continued until today (2-6). The system has been applied for a wide variety of purposes (Fig. 1) and has recently been adapted for clinical use for the localization, imaging, and therapy of cancer (7). The binder and the target can be any of the components listed in Figure 1. What is required for this system is the capacity to biotinylate a binding entity so that the specificity and activity of the binding function is retained. For different approaches to biotinylation and the reagents used for binding to different functional groups, see [Biotin](#). As can be seen, most of the functional groups on biological molecules can be modified with biotin.

Because of its charge neutrality and lack of **glycosylation**, streptavidin is generally preferred over egg-white avidin in many applications, although new derivatives of avidin (e.g., Neutravidin) may prove advantageous and less expensive [see [Avidin](#) and [Streptavidin](#)]. The final component added to the system is the probe. The various probes and their potential uses are shown in Fig. 1. The probes are prepared in two ways. They are chemically conjugated directly to avidin or streptavidin (Fig. 1, Approach A) and are **fluorescent**, **radioactive**, or other types of macromolecules (proteins, polysaccharides, etc.). A second approach is to biotinylate the probes and to interact them with streptavidin under subsaturating ratios, thus leaving extra binding sites vacant (Fig. 1, Approach B). More recently, **fusion proteins** have been prepared of streptavidin with different [enzymes](#) and native fluorescent proteins.

In many cases, such as **affinity chromatographic** applications, more reversible binding of biotin to avidin would be a distinct advantage. In this context, streptavidin has a critical disadvantage in that the interaction between its subunits is very strong, and it cannot be used when reversibility of the interaction is desired. This is in contrast to avidin, from which an immobilized monovalent form can be produced (8), and the immobilized avidin monomer binds biotinylated compounds reversibly (9).

[Site-directed mutagenesis](#) and [chemical modification](#) studies are currently being performed on both avidin and streptavidin to understand better their interaction with biotin, the interaction between their subunits, and, perhaps eventually, for better application in the previously-mentioned systems (10, 11). In this regard, the single, critical [tyrosine](#) residue of the binding sites of avidin and streptavidin was **nitrated**, and the tetrameric structures of the resultant “nitro avidin” and “nitro streptavidin” were retained (12). The biotin-binding property also had sufficiently high affinity for a variety of applications, including affinity chromatography, enzyme immobilization, and **phage-display** technology (13, 14). The major difference between the nitro avidins and the native molecules is that biotinylated compounds are released by competition with free biotin, and then the latter is liberated by treating the column with basic solutions (pH 10), thereby regenerating the original biotin-binding capacity of the nitro avidin affinity column. Reduced or altered binding characteristics are also conferred on the binding sites of avidin or streptavidin by site-directed mutagenesis of selected binding site residues, such as [tryptophans](#).

Today, the avidin/streptavidin-biotin system is a real necessity in most fields of biological study.

## Bibliography

1. M. Wilchek and E. A. Bayer, eds. (1990) *Avidin-Biotin Technology Methods in Enzymology*, Vol. **184**, Academic Press, San Diego.
2. E. A. Bayer and M. Wilchek (1978) *Trends Biochem. Sci.* **3**, N237–N239.
3. E. A. Bayer and M. Wilchek (1980) *Methods Biochem. Anal.* **26**, 1–45.
4. M. Wilchek and E. A. Bayer (1984) *Immunol. Today* **5**, 39–43.
5. M. Wilchek and E. A. Bayer (1989) In *Protein Recognition of Immobilized Ligands* (T. W. Hutchens, ed.), Alan R. Liss, New York, pp. 83–90.
6. E. A. Bayer and M. Wilchek (1996) In *Immunoassay* (E. P. Diamandis and T. K. Christopoulos, eds.), Academic Press, San Diego, pp. 237–267.
7. G. Paganelli, P. Magnani, A. G. Siccardi, and F. Fazio (1995) In *Cancer Therapy with Radiolabeled Antibodies* (D. M. Goldenberg, ed.), CRC Press, Boca Raton, FL, pp. 239–254.

8. N. M. Green and E. J. Toms (1973) *Biochem. J.* **133**, 687–700.
9. K. P. Henrikson, S. H. G. Allen, and W. L. Maloy (1979) *Anal. Biochem.* **94**, 366–370.
10. A. Chilkoti, P. H. Tan, and P. S. Stayton (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1754–1758.
11. A. Chilkoti, B. L. Schwartz, R. D. Smith, C. J. Long, and P. S. Stayton (1995) *Bio/Technology* **13**, 1198–1204.
12. E. Morag, E. A. Bayer, and M. Wilchek (1996) *Biochem. J.* **316**, 193–199.
13. E. Morag, E. A. Bayer, and M. Wilchek (1996) *Anal. Biochem.* **243**, 257–263.
14. M. E. M. Balass, E. A. Bayer, S. Fuchs, M. Wilchek, and E. Katchalski-Katzir (1996) *Anal. Biochem.* **243**, 264–269.

### Suggestion for Further Reading

15. M. D. Savage, G. Mattson, S. Desai, G. Nielander, S. Morgensen, and E. J. Conklin (1992) *Avidin-Biotin Chemistry: A Handbook* Pierce Chemical Co., Rockford, Ill.

## Azurin

Azurins are [proteins](#) of bacterial origin that belong to the family of cupredoxins (ie, blue single-copper proteins) as is [plastocyanin](#) (1). Azurins consist of a single [polypeptide chain](#) that contains ca. 128 amino acid residues (~14 kDa molecular weight) (2). Several [X-ray crystallography](#) studies of azurins isolated from various bacteria, as well as of single-site mutated forms, have yielded high-resolution, three-dimensional [protein structures](#) (Fig. 1) (1, 3-6). The protein consists of eight [b-strands](#) that form a b-sandwich (barrel) structure. A segment between b-strands 4 and 5 lies outside this barrel and contains a short **a-helix** (residues 55 to 67). The copper ion is bound about 7 Å below the protein's surface and is coordinated by five ligands: three equatorial ones [a [cysteine](#)(112) thiolate and two imidazoles (His46 and His117)] and two weak ligands in axial positions (the carbonyl O atoms of Gly45 and Met121). This configuration (the carbonyl O atom in Gly45 excepted), is conserved in all known structures of “blue” type 1 (T1) copper sites determined thus far (1).

**Figure 1.** The three-dimensional structure of azurin. The copper center with its ligands is seen near the top of the molecule. The disulfide bond is found at the opposite end of the molecule.



These protein–metal interactions yield a copper site that has a distorted, trigonal bipyramidal geometry, which may explain the stabilization of the Cu(I) state relative to Cu(II) (7) (see [Redox Enzymes](#)). Structure-imposed Cu → L p-backbonding has been proposed to account further for stabilization of the cuprous state, because strong p-interaction with the  $d_p$  orbitals results in an increase of the ligand field strength. The [oxidation/reduction potentials](#) of the Cu(II)/Cu(I) in azurins (~310 mV) are higher than is generally observed for copper complexes. Studies of several mutant forms in which Met121 is replaced, demonstrated that large **hydrophobic** residues raise the Cu(II)/(I) potential, whereas negatively charged ones lower it (8): The span in reduction potential going from Glu121 to Leu121 azurin is 300 mV, or about half the range of the potentials for the naturally occurring T1 sites. The large increase in reduction potentials caused by substitutions involving bulky hydrophobic side chains at position 121 is rationalized by exclusion by such side chains of [water](#) from the metal site or simply by providing a low dielectric environment for the copper ion. In contrast, negatively charged hydrophilic residues or water are expected to stabilize the copper ion in the more positively charged +2 state and thus lower the potential, as is indeed observed. Smaller effects on the reduction potential were observed when mutations were made in other positions. It should be emphasized, however, that a thiolate copper ligand *per se* is not a prerequisite for high reduction potentials (8).

As with all other cupredoxins, azurins are also spectroscopically characterized by an intense blue color and a uniquely narrow hyperfine coupling in the electron paramagnetic resonance (EPR) spectrum (see [Electron Paramagnetic Resonance](#)). The intense blue color is due to a

$\pi S \rightarrow Cu(d_{x^2 - y^2})$  ligand-to-metal charge transfer involving the thiolate ligand as the electron donor; it has a molar extinction coefficient of 4000 to 6000  $M^{-1}cm^{-1}$  (9). This is more than 100 times larger than what is found for simple Cu(II) complexes. A second characteristic property associated with the blue (T1) Cu (II) is an EPR spectrum displaying a hyperfine splitting in the  $g_{\parallel}$  region, due to interaction of the Cu(II) nuclear and electron spins, which is unusually narrow ( $\sim 0.008\text{ cm}^{-1}$ ), or approximately 50% smaller than those of inorganic copper complexes. This is attributed to delocalization of the unpaired  $Cu(d_{x^2 - y^2})$  electron onto the Cys(S)  $pp$  orbital, thus reducing the nuclear–electron interaction (10). The unique spectroscopic properties of the T1 site are most likely related to the electron-mediating function of cupredoxins, and this relationship is a central issue in studies of biological [electron transfer proteins](#) (11).

Recent experimental investigations of the electron transfer mechanism by azurin have in part been directed toward the intramolecular process in single-site mutated forms of the proteins (12-14). This was motivated by the finding that azurin provides a useful model system for studies of long-range, intraprotein electron transfer (LRET). One way by which this reaction can be induced is from the half-reduced [disulfide bond](#) to the Cu(II) ion (12, 14). Another one is by attaching an extraneous redox center to the protein (11, 13). Results of these studies yielded insights into the mechanisms of intramolecular electron transfer in proteins and the parameters that determine the degree of electronic coupling between electron donor and acceptor. Effective electron transfer in proteins requires minimal reorganization of the redox centers as a result of the process. As described above, the metal ion in the blue copper proteins is coordinated to a site that is intermediate between the preferred geometries for Cu(II) (tetragonal planar) and Cu(I) (tetrahedral) (6). Furthermore, the relatively high reduction potentials provides an increased driving force. Finally, the tightly knit, antiparallel [b-sheet](#) structure of azurin probably improves the electronic coupling between the reaction partners through the protein matrix.

## Bibliography

1. E. T. Adman (1991) *Adv. Protein Chem.* **42**, 145–198.
2. L. Ryden and J.-O. Lundgren (1976) *Nature* **261**, 344–346.
3. E. T. Adman and L. H. Jensen (1981) *Isr. J. Chem.* **21**, 8–12.
4. H. Nar, A. Messerschmidt, and R. Huber (1991) *J. Mol. Biol.* **218**, 427–447.
5. W. E. B. Shephard, B. F. Anderson, D. A. Lewandoski, G. E. Norris, and E. N. Baker (1990) *J. Am. Chem. Soc.* **112**, 7817–7819.
6. A. Messerschmidt, L. Prade, S. J. Kroes, J. Sanders-Loehr, R. Huber, and G. W. Canters (1998) *Proc. Natl. Acad. Sci. USA* **95**, 3443–3448.
7. C. S. St. Clair, W. R. Ellis Jr., and H. B. Gray (1992) *Inorg. Chim. Acta* **191**, 149–155.
8. T. Pascher, B. G. Karlsson, M. Nordling, B. G. Malmström, and T. Vanngård (1993) *Eur. J. Biochem.* **212**, 289–296.
9. H. B. Gray and E. I. Solomon (1981) In *Copper Proteins*, Vol. **3** (T. G. Spiro, ed.), Wiley, New York, pp. 1–39.
10. W. E. Antholine, P. M. Hanna, and D. R. McMillin (1993) *Biophys. J.* **64**, 267–272.
11. A. J. Di Bilio et al. (1997) *J. Am. Chem. Soc.* **119**, 9921–9922.
12. O. Farver, L. K. Skov, G. Gilardi, G. van Pouderoyen, G. W. Canters, S. Wherland, and I. Pecht (1996) *Chem. Phys.* **204**, 271–277.
13. J. J. Regan et al. (1995) *Chem. Biol.* **2**, 489–496.
14. O. Farver and I. Pecht (1997) *J. Biol. Inorg. Chem.* **2**, 387–392.
15. T. J. Mizoguchi, A. J. Di Bilio, H. B. Gray, and J. H. Richards (1992) *J. Am. Chem. Soc.* **114**, 10076–10078.

## B Cell

Specific recognition by the immune system is mediated fundamentally by B and T lymphocytes. These cells derive from the hematopoietic [stem cells](#) through discrete differentiation pathways, although both originate in the bone marrow (and in the liver during embryonic and fetal life) (see [Hematopoiesis](#)). [T cells](#) mature in the thymus, whereas B cells continue to differentiate in the bone marrow itself. B cells express [immunoglobulins](#) that directly interact with native [epitopes](#)—that is, subregions of an [antigen](#) with a well-defined three-dimensional structure. Surface immunoglobulins are associated with a signaling module made of a Iga-Igb heterodimer to form the B-cell receptor (BCR). In contrast to the BCR, [T-cell receptors](#) (TCRs) do not interact directly with native epitopes, but they identify peptides derived from the original antigen as presented by molecules encoded by genes of the [major histocompatibility complex](#) (MHC).

Both the B and the T cells have to face the [repertoire](#) problem—that is, how to generate extremely large number of different immunoglobulins and TCRs to meet the requirements of recognizing an extraordinary large number of discrete epitopes. A theoretical approach suggests that this number might be as large as  $10^{17}$ . Because the total number of lymphocytes in a human adult averages  $10^{12}$  (roughly divided into one-fifth B cells and four-fifths T cells), and taking into account that they are clonally organized, the number of BCRs and/or TCRs expressed at one given time appears much less than  $10^{17}$ , which implies some basic degeneracy in the immune recognition system. Nevertheless, the number of different structures must still be quite large, which therefore raises the problem of how to generate such a large repertoire with, necessarily, a limited number of genes. This is the main challenge for B cells (and T cells as well), to which the way they differentiate brings the answer.

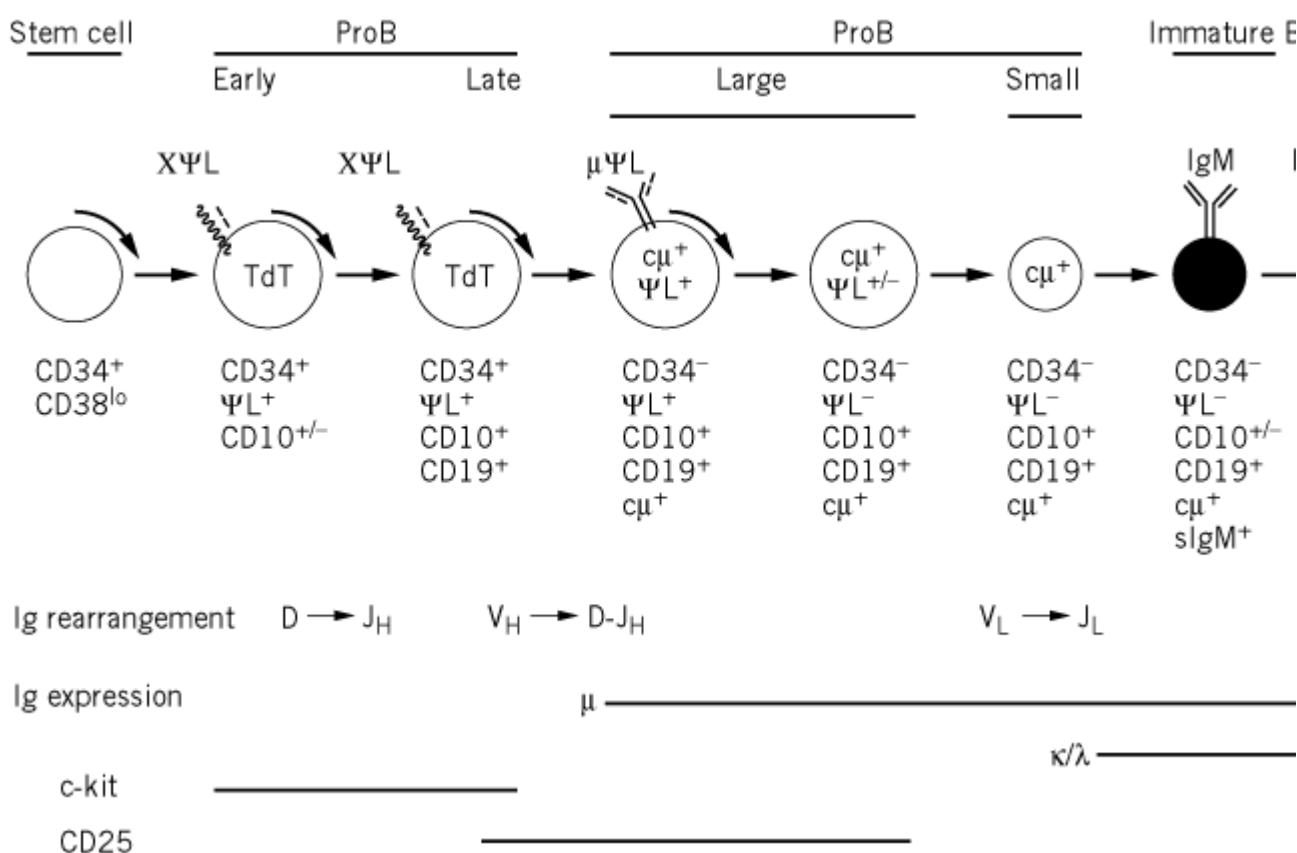
During fetal life, B lymphocytes differentiate in the liver and then in the bone marrow, which remains active throughout life. The first steps of B-cell differentiation take place in the bone marrow, which is a primary lymphoid organ, and drives precursors derived from the hematopoietic stem cells to the immature B lymphocytes. This period of differentiation is antigen-independent and is essentially devoted to generation of the basic [immunoglobulin](#) repertoire, which is the result of a complex sequence of events that involve multiple gene rearrangements. Immature B lymphocytes migrate to secondary lymphoid organs, namely spleen, lymph nodes, tonsils, and gut-associated lymphoid tissues, including Peyer's patches. Further differentiation necessitates antigen encounter and a number of complex cellular interactions, among which cooperation between the B and T cells plays a major role. A second round of diversity is then generated, which results mostly from somatic mutations, and the final steps of differentiation, including isotype switching, lead to the emergence of plasma cells and memory B cells. Two discrete lineages of B cells have been well-characterized in the mouse and are designated B1 and B2. B1 cells express the CD5 marker, have a low level of surface [IgD](#), and are mostly encountered in peritoneal and pleural cavities. They seem preferentially responsible for autoantibody synthesis. B2 cells represent the major and conventional B-cell population.

### 1. From the Hematopoietic Stem Cell to the Immature B Cell

All differentiation events that drive the emergence of the various lineages (granulocytes, erythrocytes, monocytes, megakaryocytes, **lymphocytes**) derive from a common hematopoietic stem cell that originates in the bone marrow. The B lineage presumably initiates from a precursor common to the B and the T lymphocytes. The main features of the molecular and cellular events that take place in the bone marrow and lead to B cell commitment are indicated in [Figure 1](#). Three major subpopulations define the B lineage: proB, preB, and immature B cells. The successive steps of

differentiation proceed from the periphery toward the center of the bone marrow and may be followed in several ways: (a) acquisition and/or loss of surface antigens (especially CD markers) identified by [monoclonal antibodies](#), (b) identification of cytoplasmic proteins and/or [messenger RNA](#) transcripts, and (c) analysis of the gene rearrangement status for each of the three Ig gene loci: H, L, and K.

**Figure 1.** B-cell differentiation in humans. Early steps of B-cell differentiation take place in the bone marrow and are an The major discrete steps that lead from the hematopoietic stem cell (HSC) to the immature B lymphocyte are identified by gene rearrangements, by the presence of cytoplasmic markers, and by the expression of characteristic markers at the cell early stages that go from HSC to the early proB cell encompass several precursors that are not yet completely defined. T deoxynucleotidyl transferase, responsible for generating N diversity;  $\Psi L$ , surrogate light chain that combines with the  $\mu$  heavy chain to form the preB receptor. Late events of B-cell differentiation, which take place in the periphery, are antigen-dependent. The gene organization plays a key role in mounting a T-dependent B-cell response, during which affinity is tuned upon the acquisitions mutations and the biological functions of the antibodies are adapted by isotype switching. Ultimately, this differentiation cells that secrete antibodies and in memory B cells that are involved in secondary responses. Each step of differentiation by various CD surface markers, as indicated.



The first steps of differentiation involve a sequence of direct interactions of the precursors with the stromal cells, ensured by various sets of **cellular adhesion molecules** (CAMs). The VLA4 [integrin](#) expressed at the cell surface of the precursors interacts with stromal VCAM1, and the resulting early proB cells now express Kit, a receptor that binds the stem cell factor of stromal cells. This interaction triggers proliferation of the proB cells, which will continue to differentiate through other stimuli. Late proB cells are stimulated by a soluble [growth factor](#) also of stromal origin, **interleukin-7** (IL-7), which drives proliferation of preB cells. As these steps of differentiation proceed, the cells are dividing, and Ig gene rearrangements that generate the basic repertoire take place in a sequential fashion. CD marker characterization indicates that CD34, which is already expressed in the hematopoietic stem cells, remains present at the surface of precursors and proB cells. Late proB cells can be characterized by the coexpression of CD34 with CD19, which is a very specific marker of the

B lineage because it is found on all subsequent stages, with the exception of the plasma cells. CD10, also known as the common acute lymphocytic leukemia antigen (CALLA), is expressed up to the immature B-cell stage a seems expressed slightly before CD19. Other markers are expressed as differentiation proceeds, such as CD20, CD21, CD22, CD23, and CD24.

## 2. The Basic Ig Repertoire Develops in the Bone Marrow as B Cells Differentiate

As already stated, the main feature of B-cell differentiation is the acquisition of the basic repertoire of immunoglobulins that becomes expressed as IgM on immature B cells. IgM has the classical  $H_2L_2$  organization, and both the H and L chains have variable and constant regions. The variable regions of the heavy and light chains interact in the antibody combining site, which is responsible for specific antigen recognition. Heavy and light chains are encoded by genes that are localized on three discrete gene clusters, H, k and l, located on chromosomes 14, 2, and 22, respectively. The unique feature of Ig gene organization is that they must be rearranged by the specific [recombinases](#) RAG1 and RAG2 (for recombinase activating genes) before becoming functional. The first recombination event, DH to JH, takes place in the early proB cells (see Fig. 1) and is rapidly followed by the rearrangement of one of the  $V_H$  segments to D-J. Another enzyme, terminal deoxynucleotidyl-transferase, is also active in proB cells, for which it represents an additional cytoplasmic marker. This enzyme adds nucleotides in a random fashion during the joining of D to J and of V to D-J rearrangements. ProB cells that have rearranged the  $IgV_H$  locus must “make” two decisions. One is to remain monoclonal with respect to H-chain production, the second is to activate rearrangement of the light chains. Both events are regulated by the m chain itself, which must be in its membrane form. Once the first allele of the  $IgH$  locus has completed the rearrangement process, the resulting gene may be functional or not, depending upon whether the recombination has generated a sequence of nucleotides with an open reading frame. Because of the triplet organization of the [genetic code](#), this happens only once every three rearrangements. If the first rearrangement is out of frame, the second allele will recombine, with the same probability of success. A second failure will lead to [cell death](#). Conversely, once a functional gene has been generated, the resulting heavy chain will exert a negative feedback on a further rearrangement of the  $IgH$  locus, ensuring monoclonal expression of the m chain. Much evidence suggests that the heavy chain must be expressed at the cell surface to regulate these events. Expression of the m chain at the cell surface requires that it associate with another partner, which resembles the light chain, and for this reason it is named YL or surrogate light chain. First described in the mouse and then identified in humans, the surrogate light chain is composed of two polypeptide chains, encoded by the l-like (or l5 in the mouse) and the VpreB genes, part of the regular  $IgL$  locus. The m-YL complex becomes expressed at the surface of what is now a large preB cell, which also expresses CD10 and CD19, but no longer CD34. The cell is able to undertake light-chain gene rearrangements, which occur in the order  $k \rightarrow l$ . The recombination process is regulated in the same manner as that of the heavy chains, so there is only one light chain that is expressed in any given cell, either k or l (see [1\) Light Chains](#)). The negative feedback on further rearrangements of the light chain genes is exerted by the complete IgM, which is now expressed at the surface of the immature B cell and has replaced the mYL complex. At this point, the  $Iga-Igb$  module is associated with the preB cell complex, strongly suggesting that the mYL complex might be considered a “preB” receptor.

## 3. Immature B Cells are Selected Before Leaving the Bone Marrow

Before leaving the bone marrow, immature B cells that express surface IgM are confronted with the local “self” antigenic environment. They are particularly sensitive to triggering by multivalent antigens, resulting in their clonal deletion by [apoptosis](#). Alternatively, soluble self antigens do not cause the immature B cells to die, but instead induce an anergic state that does not seem to prevent the cells migrating to the periphery. It should be noted, however, that this negative selection has a threshold that leaves a fraction of autoreactive cells going to the periphery. These cells are responsible for the presence in the bloodstream of natural [autoantibodies](#) that must be of some physiological relevance. Once validated by this “quality control,” immature B cells circulate to the



periphery through blood vessels and lymphatic system, colonize the secondary lymphoid organs, and actively recirculate. They now have become mature B cells, expressing both IgM and IgD isotypes at their surface.

#### 4. Final Steps of B-Cell Differentiation Take Place in the Periphery and are Antigen-Dependent

The B lymphocytes that are now circulating in the periphery are nondividing, short-lived cells. To achieve their ultimate differentiation, they must be triggered by an antigen, most often T-cell-dependent (Fig. 1). Within the first days after antigenic stimulation and T-cell help, activated B cells may mature to plasma cells that have an abundant [endoplasmic reticulum](#) and secrete immunoglobulins. Activated B cells may also evolve to colonize the primary follicle of a lymph node and generate a germinal center, where they interact with follicular dendritic cells, divide rapidly, switch to another [isotype](#), most frequently [IgG](#), and start to accumulate [somatic hypermutations](#) that considerably amplify diversity and affinity. Clones of high affinity will be positively selected by antigen and can now progress toward plasma cell differentiation or become long-lived B [memory cells](#).

#### 5. Conclusions

Differentiation of B cells is a highly sophisticated process that takes place continuously in the bone marrow and results in the constant emergence of a very large repertoire of immunoglobulins, expressed both as membrane B-cell receptors (BCRs) and soluble antibodies. The numerous molecular events that lead from the hematopoietic stem cell to the mature B cell are under constant selective pressure, implying a high level of cell death. Therefore, the available repertoire appears, at any given time, to be a compromise between the necessary economy in the gene number, compensated by the recombination processes, and an unavoidable wastage due to the stochastic aspects of gene rearrangements and to the negative selection of clones having a high affinity for self components.

See also entries [Antibody](#), [Clonal Selection Theory](#), [Gene Rearrangement](#), [Immunoglobulin](#), and [Repertoire](#).

#### Suggestions for Further Reading

P. D. Burrows and M. D. Cooper (1997) B cell development and differentiation. *Curr. Opin. Immunol.*, **9**, 239–244.

Y. J. Liu and C. Arpin (1997) Germinal center development. *Immunol. Rev.*, **156**, 111–126.

H. Karasuyama, A. Rolink, and F. Melchers (1996) Surrogate light chain on B cell development. *Adv. Immunol.*, **63**, 1–41.

A. Galy, M. Travis, D. Cen, and B. Chen (1995) Human T, B, natural killer and dendritic cells arise from a common bone marrow progenitor cell subset. *Immunity* **3**, 459–473.

J. Banchereau and F. Rousset (1992) Human B lymphocytes: phenotype, proliferation, and differentiation. *Adv. Immunol.* **52**, 125–262.

L. A. Herzenberg and A. B. Kantor (1993) B-cell lineages exist in the mouse. *Immunol. Today* **14**, 79–83.

## **B<sub>12</sub> (Cobalamin)**

*Cobalamins* are complex organometallic cofactors containing a central cobalt atom that is coordinated equatorially to four nitrogen atoms provided by the corrin ring. The two cofactor forms of this vitamin are *adenosylcobalamin* (AdoCbl) and *methylcobalamin* (MeCbl). A variety of other forms have been described in which the upper ligand is a water, *glutathione*, or cyano group, but the physiological significance of these forms is unknown. In nature, this cofactor is found in association with two [enzyme](#) subfamilies, the AdoCbl-dependent *isomerases* that catalyze 1,2 rearrangement reactions and the MeCbl-dependent *methyltransferases* that catalyze transmethylation reactions. Members of both subfamilies are fairly prevalent in the bacterial world, where the isomerases are involved in fermentative pathways, and the methyltransferases are involved in pathways leading to methionine, acetate, or methane. In mammals, only two B<sub>12</sub>-dependent enzymes are known, the cytoplasmic *methionine synthase* and the mitochondrial *methylmalonyl-CoA mutase*. In this review, the molecular biological aspects of B<sub>12</sub> biosynthesis, enzymology and diseases will be discussed briefly.

## 1. B<sub>12</sub> Biosynthesis

The structural cousins, vitamin B<sub>12</sub>, heme, and chlorophyll, are constructed from a common template that is derived from the precursor, 5-aminolevulinic acid. Methylation at C2 converts the common intermediate, uroporphyrinogen III, to precorrin-1 and commits it to B<sub>12</sub> biosynthesis. Thereafter, a biosynthetic assembly line, which includes a series of methylations, ring contraction (between rings A and D), cobalt insertion, alkylation, amidations, and nucleotide loop assembly reactions, takes the cofactor to its final form. Two separate pathways lead to corrin ring biosynthesis in the aerobic and anaerobic worlds (for reviews, see Refs. [1-4](#)). They run parallel at some steps and diverge at others. A unique structural feature of B<sub>12</sub> is the presence of a nucleotide loop that terminates in dimethylbenzimidazole, which can serve as the lower axial ligand to cobalt. The pathways leading to the assembly of dimethylbenzimidazole in aerobic and anaerobic organisms are distinct, but their details remain to be elucidated ([5](#)). The **genes** encoding corrin biosynthesis functions have been **cloned** from *Pseudomonas denitrificans*, where they are organized in clusters scattered along the [genome](#) (reviewed in Ref. [3](#)) and from *Salmonella typhimurium*, where they are clustered at 41 min (reviewed in Ref. [2](#)). The [transcription factor](#) PocR regulates the transcription of the cobalamin biosynthetic and propanediol utilization (dependent on a B<sub>12</sub> enzyme) genes in *S. typhimurium* ([6](#), [7](#)).

## 2. B<sub>12</sub> Enzymes

The cofactor role of AdoCbl was first described by Barker and co-workers in glutamate mutase ([8](#)), followed a few years later by the discovery of MeCbl in methionine synthase ([9](#)). B<sub>12</sub>-dependent enzymes catalyze chemically difficult reactions and manipulate the reactive cobalt–carbon bond in radically different ways (reviewed in Ref. [10](#)). In the methyltransferase subfamily, cobalamin is the initial acceptor of the methyl group donated from substrates such as methyltetrahydrofolate and methanol, and alkyltransfer to and from the cobalamin occurs *via* heterolytic cleavage of the cobalt–carbon bond. In the isomerases, the cobalt–carbon bond is broken homolytically, as the enzymes catalyze 1,2 rearrangement reactions. The migration of a diversity of groups, ranging from carbon (in methylmalonyl-CoA mutase, glutamate mutase, and methylene glutarate mutase) to nitrogen (in ethanolamine ammonia lyase and  $\epsilon$ -lysine mutase) and oxygen (in diol dehydrase), is catalyzed by the isomerases. The genes encoding a number of methyltransferases and isomerases have been cloned and the structures of the B<sub>12</sub>-binding domain of the *Escherichia coli* methionine synthase ([11](#)) and of the *Propionibacterium shermanii* methylmalonyl-CoA mutase ([12](#)) have been determined. In both enzymes, the cofactor is bound in an extended conformation in which the intramolecular base, dimethylbenzimidazole, is removed from the cobalt and replaced by a **histidine** residue donated by the protein. Another important member of the B<sub>12</sub> family of enzymes is

[ribonucleotide reductase](#), which converts ribonucleotides to deoxyribonucleotides (13). The reaction mechanisms of B<sub>12</sub>-dependent enzymes are discussed in a number of reviews (10, 14-18).

### 3. B<sub>12</sub>-Related Diseases

Functional B<sub>12</sub> deficiency can result from either nutritional insufficiency or from genetic defects (19, 20). It is manifested clinically by a combination of symptoms including hematological and neurological abnormalities, methylmalonic aciduria and homocystinuria, depending on whether one or both B<sub>12</sub>-dependent enzymes is affected. The genetic defects or inborn errors of cobalamin metabolism can result from impairments in uptake, transport, or enzymatic function and are inherited as autosomal recessive traits. Cobalamin absorption and transport abnormalities resulting from impairments in intrinsic factor, its receptor, or in transcobalamin (TC) II have been described. The cDNA encoding human intrinsic factor (21) and TC II (22) have been cloned and mapped, and genetic defects in TC-II in patient cell lines have been identified (23).

Intracellular cobalamin metabolism is complex, compartmentalized, and dependent on several enzymes. Defects in the early steps following internalization of TC II/cobalamin affect both MeCbl and AdoCbl syntheses and fall into the *cbl* C, D, and F genetic complementation groups. The identities of the proteins encoded by these loci and their precise functions are unknown. MeCbl synthesis is specifically compromised in *cblG* and *cblE* patients. Cloning of the cDNA encoding methionine synthase (24-26) and identification of mutations correlated with the *cblG* phenotype (25, 27) indicate that this locus represents methionine synthase. AdoCbl synthesis is specifically impaired in *cblA*, B, and *mut* patients. The defect in *cblB* and *mut* patients is in cobalamin adenosyltransferase and methylmalonyl-CoA mutase, respectively. The cDNA encoding methylmalonyl-CoA mutase has been cloned (28), and a number of pathogenic mutations have been described (29).

### Bibliography

1. A. Battersby (1994) *Science* **264**, 1551–1557.
2. A. I. Scott (1994) *Tetrahedron* **50**, 13315–13333.
3. F. Blanche, B. Cameron, J. Crouzet, L. Debussche, D. Thibaut, M. Vuilhorgne, F. J. Leeper, and A. R. Battersby (1995) *Angew. Chem. Int. Ed. Engl.* **34**, 383–411.
4. A. I. Scott (1996) *Pure Appl. Chem.* **68**, 2057–2063.
5. P. Renz (1998) In *Vitamin B<sub>12</sub> and B<sub>12</sub> Proteins* (B. Kraeutler, D. Arigoni, and B. T. Golding, eds.), Wiley-VCH, Weinheim, Germany.
6. T. A. Bobik, M. Ailion, and J. R. Roth (1992) *J. Bacteriol.* **174**, 2253–2266.
7. M. R. Rondon and J. C. Escalante-Semerena (1992) *J. Bacteriol.* **174**, 2267–2272.
8. H. A. Barker, H. Weissbach, and R. D. Smyth (1958) *Proc. Natl. Acad. Sci. USA* **44**, 1093–1097.
9. J. R. Guest, S. Friedman, D. D. Woods, and E. L. Smith (1962) *Nature* **195**, 340–342.
10. R. Banerjee (1997) *Chem. Biol.* **4**, 175–186.
11. C. L. Drennan, S. Huang, J. T. Drummond, R. Matthews, and M. L. Ludwig (1994) *Science* **266**, 1669–1674.
12. F. Mancina, N. H. Keep, A. Nakagawa, P. F. Leadlay, S. McSweeney, B. Rasmussen, P. Bösecke, O. Diat, and P. R. Evans (1996) *Structure* **4**, 339–350.
13. S. Licht, G. J. Gerfen, and J. Stubbe (1996) *Science* **271**, 477–481.
14. C. L. Drennan, R. G. Matthews, and M. L. Ludwig (1994) *Curr. Opin. Struct. Biol.* **4**, 919–929.
15. B. T. Golding and W. Buckel (1997) In *Comprehensive Biological Catalysis* (M. L. Sinnott, ed.), Vol. **4**, Academic Press, London.
16. B. Babior (1988) *Biofactors* **1**, 21–26.

17. J. Stubbe (1988) *Biochemistry* **27**, 3893–3900.
18. R. G. Matthews, R. V. Banerjee, and S. W. Ragsdale (1990) *Biofactors* **2**, 147–152.
19. W. A. Fenton and L. E. Rosenberg (1995) *Inherited Disorders of Cobalamin Transport and Metabolism*, McGraw-Hill, New York, pp. 3111–3128.
20. M. I. Shevell and D. S. Rosenblatt (1992) *Can. J. Neurol. Sci.* **19**, 472–486.
21. B. K. Dieckgraefe, B. Seetharam, L. Banaszak, J. F. Leykam, and D. H. Alpers (1988) *Proc. Natl. Acad. Sci. USA* **85**, 46–50.
22. O. Platica, R. Janeczko, E. V. Quadros, A. Regec, R. Romain, and S. P. Rothenberg (1991) *J. Biol. Chem.* **266**, 7860–7863.
23. L. Qian, E. V. Qaudros, and S. P. Rothenberg (1997) *Meth. Enzymol.* **281**, 269–281.
24. Y. N. Li, S. Gulati, P. J. Baker, L. C. Brody, R. Banerjee, and W. D. Kruger (1996) *Human Mol. Genet.* **5**, 1851–1858.
25. D. Leclerc, E. Campeau, P. Goyette, C. E. Adjalla, B. Christensen, M. Ross, P. Eydoux, D. S. Rosenblatt, R. Rozen, and R. A. Gravel (1996) *Human. Mol. Genet.* **5**, 1867–1874.
26. L. H. Chen, M.-L. Liu, H.-Y. Hwang, L.-S. Chen, J. Korenberg, and B. Shane (1997) *J. Biol. Chem.* **272**, 3628–3634.
27. S. G. Gulati, P. Baker, B. Fowler, Y. Li, W. Kruger, L. C. Brody, and R. Banerjee (1996) *Human Mol. Genet.* **5**, 1859–1866.
28. R. Jansen, F. Kalousek, W. A. Fenton, L. E. Rosenberg, and F. D. Ledley (1989) *Genomics* **4**, 198–205.
29. F. D. Ledlay and D. S. Rosenblatt (1997) *Human Mutation* **9**, 1–6.

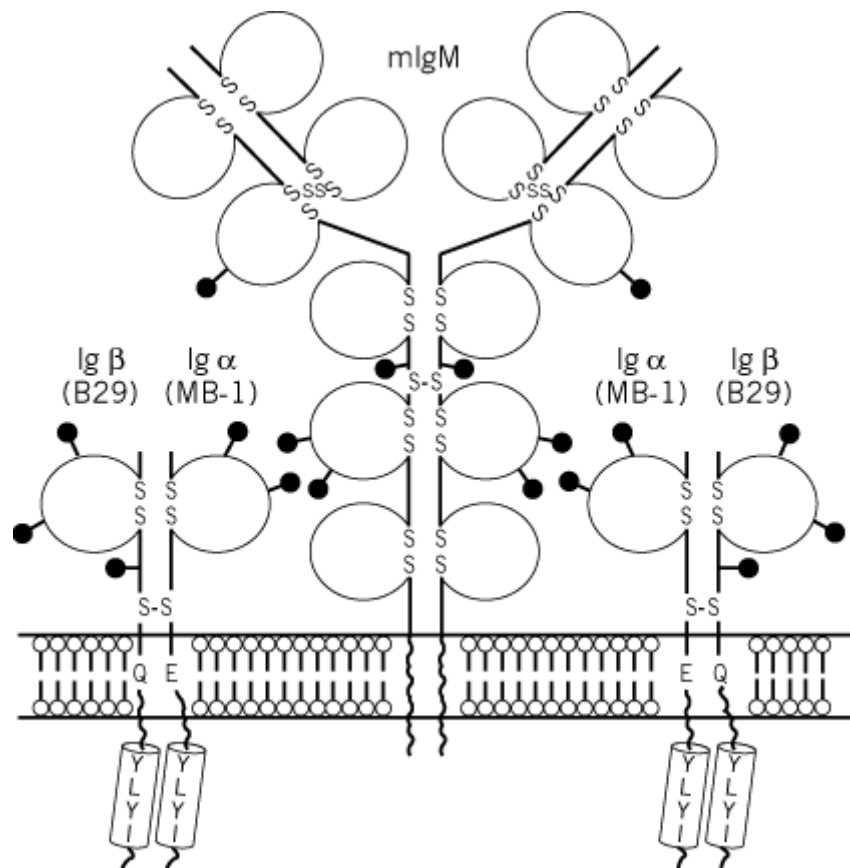
### Suggestion for Further Reading

30. R. Banerjee, ed. (1999) *Chemistry and Biochemistry of B<sub>12</sub>*, Wiley, New York.

## B-Cell Receptor (BCR)

[Immunoglobulins](#) can exist as [membrane proteins](#), expressed at the surface of [B cells](#). When B cells differentiate in the bone marrow from the initial hematopoietic [stem cells](#), they reach a stage, called *immature B lymphocytes*, that is characterized by the presence of [IgM](#) at the cell surface (Fig. [1](#)). This IgM, which is a basic m2k2 or m2l2 monomer, is anchored into the membrane bilayer by an extra portion at the COOH-terminus of the heavy chains that is not present in the soluble form. Differences in length result simply from the [alternative splicing](#) of membrane exons (see [Introns](#), [Exons](#)) present in the germline. The **transmembrane** region is terminated by only three amino acid residues located in the cytoplasm of the B cell.

**Figure 1.** Schematic representation of the B-cell receptor (BCR). The BCR is composed of one antigen binding module, the surface IgM, connected to the signaling module, a heterodimer containing the Iga and Igb, encoded by the MB-1 and B-29 genes, respectively. Note the Ig domain organization of the ectodomains of the Iga and Igb chains that contain in their cytoplasmic part the ITAM (immunoreceptor tyrosine-based activation motif) motifs that play a central role in signaling by contacting intracellular protein tyrosine kinases.



Because the surface Ig clearly acts as a B-cell receptor (BCR) upon triggering by the antigen, the necessity of a signaling module, similar to CD3 in [T cells](#), is apparent. This was indeed the case. The BCR-activating module is a heterodimer composed of one Iga and one Igb chain (also designated as CD79 a and b), each of which has an external Ig-like domain, a transmembrane portion, and an cytoplasmic region that contains several signaling motifs, called *immunoreceptor tyrosine-based activation motif* (ITAM). ITAM were first described on polypeptide chains of CD3 and consist of two Tyr—X—X—Leu/Ile sequences separated by a stretch of seven amino acid residues. They are found in a number of [signal transduction](#) modules associated with most crucial receptors of the immune system, such as the [T-cell receptor \(TCR\)](#), BCR, and FcR. Upon stimulation of BCR with cross-linking antigens, ITAM interact with protein **tyrosine kinases** of the **src** family (p59<sup>fyn</sup>, p53<sup>lyn</sup>, p56<sup>lyn</sup>, p55<sup>lck</sup>), leading to the translocation of p72<sup>syk</sup> to the membrane and activation of the **ras** pathway and **phospholipase** Cg2. As a consequence, the BCR is internalized, which is the first step of processing of antigens by B cells. Full activation of B cells require other signals, involving molecules other than the BCR.

#### Suggestion for Further Reading

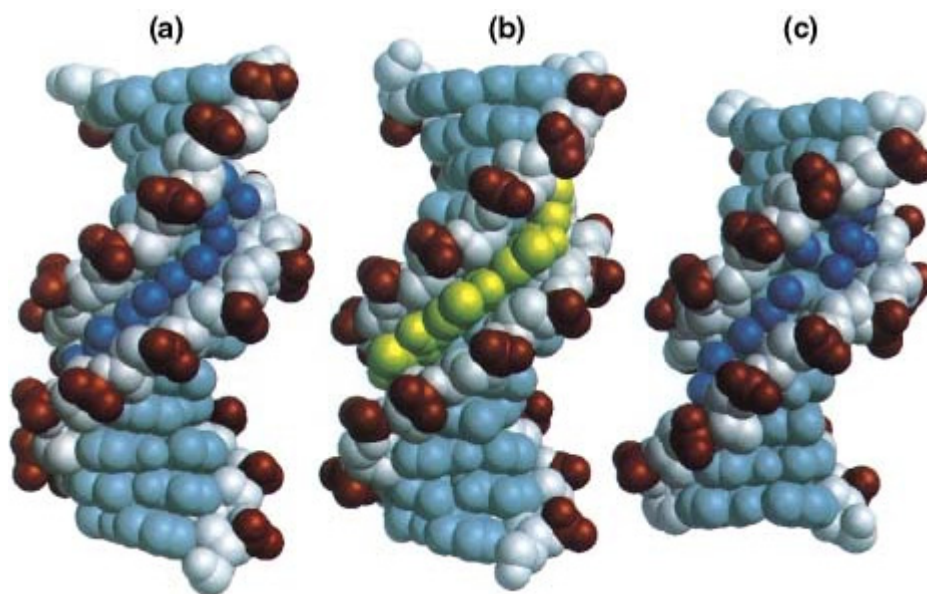
M. Reth (1995) The B cell receptor complex and co receptors. *Immunol. Today* **16**, 310–313.

#### B-DNA

B-DNA is the predominant type of [DNA structure](#) under physiological conditions. The B-DNA double helix is right-handed, with its [base pairs](#) perpendicular to the helix axis that passes through the center of the base pairs. The major groove and minor groove are roughly equivalent in depth, 8.5 and 7.5 Å, respectively, whereas the width of the major groove is ~12 Å and that of the minor groove is ~6 Å.

The first [X-ray crystallography](#) structure of B-DNA (Fig. 1A) was that of a dodecamer sequence d(CGCGAATTCGCG)<sub>2</sub> that contains the EcoRI [restriction](#) sequence GAATTC (1). There are a number of interesting structural features associated with this B-DNA structure. A distinctive feature is that the minor groove is narrower at the AATT region of the double helix than the CGCG ends. The narrow minor groove at the AATT region is filled by a spine of [water](#) molecules that form [hydrogen bonds](#) to both the O2 of thymines and the N3 of adenines, although this spine of waters has been reinterpreted recently using a higher resolution structural analysis as a spine of waters on sodium ions (2). Additionally the base pairs in the central AATT region of the helix have high propeller twist angles. The propeller twisting enhances the stacking of the bases along each strand of the double helix, and the AT base pairs are in general less restrictive to propeller twisting because of fewer hydrogen bonds than the GC base pairs (3). Lastly, the sugar pucker of the deoxyribose ring favors the C2'-endo conformation, although a range of conformations from C1'-exo to O4'-endo is also present. This may reflect the greater flexibility associated with B-DNA structures.

**Figure 1.** Examples of the structure of B-form DNA. (a) Crystal structure of the double-stranded dodecamer d(CGCGAATTCGCG) (Protein Data Base BDL084), with the water molecules in the minor groove shown as spheres. (b) Crystal structure of the double-stranded d(CGCAAATTTGCG)-distamycin A complex (Protein Database GDL003). The drug displaces the waters in the minor groove. (c) Crystal structure of the double-stranded decamer d(CCAGGCCTGG) (Protein Database BDJS30), with the waters bound to the minor groove shown as spheres. See color insert.



The narrow minor groove region associated with the A-T sequences in B-DNA affords an excellent binding site for the minor groove-binding drugs, such as distamycin A, netropsin, Hoechst 33258, Hoechst 33342, and DAPI. The crystal structure and solution structures of the complexes of several DNA oligonucleotides with those minor groove-binding drugs, most of them in 1–1 drug-to-duplex complexes (Fig. 1B), have been determined (4, 5). Those structures revealed that the drug replaces the spine of hydration in the narrow minor groove and stabilizes the DNA structure, without perturbing the overall conformation significantly. The narrow minor groove associated with the

central AATT base pairs is essential for the 1–1 binding mode of drugs. The binding energy derives in part from the gain in **entropy** associated with the displacement of the water molecules. The sequence preference for A-T regions by minor groove-binding drugs is due to the greater negative electrostatic potential at the bottom of the minor groove at A-T regions. Finally the presence of N2 amino groups of guanines provides both a charge and a steric hindrance to drug binding.

NMR study of the binding of distamycin A (a pyrrole-containing anti-tumor antibiotic) to DNA containing A-T sequences revealed a more complex pattern. Distamycin A can bind DNA not only in a 1–1 drug/duplex complex but also in a 2–1 mode, with two distamycins bound to the minor groove in an antiparallel, side-by-side manner (6). The latter binding mode requires that the minor groove width at the binding site be expanded so as to accommodate two drug molecules, reflecting the flexibility of B-DNA. This new finding has stimulated an active study in the design of new minor groove-binding compounds that can bind to all four base pairs, (A-T, T-A, C-G, and G-C), with high specificity. It was found that compounds with imidazole-containing units do not have a great discriminating power regarding the recognition toward G/C versus A/T base pairs. However, compounds that combine the imidazole units with the pyrrole units can be designed to possess excellent sequence-specific binding properties. Rules for such designs have been proposed recently (7).

Additional work on the high-resolution crystal structures of several DNA decamer oligonucleotides with “mixed” sequences, including d(CCAGGCCTGG) (Fig. 1C), showed that the minor groove width of those structures in general is wider, but with some variations. Those DNA decamers have a slightly narrow or wide minor groove, depending on the sequence. In general, the A-T regions have a narrower minor groove than the G-C regions. The [hydration](#) structure in the groove is dependent on the groove width. It was noted that a single spine of water molecules along the floor of the minor groove is associated with a narrow minor groove, whereas a ribbon of double water molecules, bridging the base edge N or O atom to the O4' atoms of the sugar ring, is associated with the wide minor groove.

### Bibliography

1. R. M. Wing, H. R. Drew, T. Takano, C. Broka, S. Tanaka, K. Itakura, and R. E. Dickerson, (1980) *Nature* **287**, 755–758.
2. X. Shui, L. McFail-Isom, G. G. Hu, and L. D. Williams (1998) *Biochemistry* **37**, 834–8355.
3. C. R. Calladine (1982) *J. Mol. Biol.* **161**, 343–352.
4. M. Coll, C. A. Frederick, A. H.-J. Wang, and A. Rich (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8385–8389.
5. H. Robinson, Y.-G. Gao, C. Bauer, C. Roberts, C. Switzer, and A. H.-J. Wang (in press) *Biochemistry*, and references cited therein.
6. J. G. Pelton and D. E. Wemmer (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5723–5727.
7. S. White, J. W. Szewczyk, J. M. Turner, E. E. Baird, and P. E. Dervan (1998) *Nature* **391**, 468–471.

### Backbone

Backbone, or main chain, is the general term used to describe the connecting chain in [polymers](#). Different kinds of polymers have different chemical backbones. For example, in [proteins](#) the backbone is a [polypeptide chain](#), but **nucleic acids** have a sugar phosphate backbone. The backbone

of a polymer, which may adopt a regular structural [conformation](#), is the constant or repeating part of the polymer, as opposed to the attached [side chains](#), which can be variable.

[See also [Polymer](#) and [Side Chain](#).]

#### Suggestions for Further Reading

L. Mandelkern (1983) *An Introduction to Macromolecules*, Springer-Verlag, New York.

P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

## Bacteriocins

Broadly defined, bacteriocins are substances produced by one **bacterium** that adversely affect another. Most of them are **peptide antibiotics** that are synthesized on [ribosomes](#), as in normal **protein biosynthesis**. Others, such as bacitracin and **gramicidin**, are synthesized in bacteria by multienzyme complexes or sequential enzyme reactions (1). In a strict sense, bacteriocin is an abbreviation of bacterial toxin (with antibacterial activity). Many bacterial toxins act against **eukaryotic** cells and are not called bacteriocins. Some of them, however, like **diphtheria**, tetanus or **cholera** toxins, resemble bacteriocins, such as [colicins](#) or pyocins, in having domain structures that convey the binding and toxic activities specific to each molecule (2, 3). The current basis for allocating a name to the agent responsible for bacteriocinlike activity produced by **Gram-positive** or -negative bacteria is to adopt some derivation of either the genus or species name of the producer strain, together with an alphabetical and/or numerical code designation specifying that strain. To avoid confusion, bacteriocins that have the same amino acid sequence, irrespective of the species of origin, should have the same name, that first published. There is a wide variety of locations of the genetic determinants of bacteriocins. They can be encoded by conjugative or nonconjugative **plasmids**, by conjugative [transposons](#), or be encoded on the [chromosome](#) (4).

Because of the different cell-wall structures of gram-positive or -negative bacteria, their bacteriocins have evolved differently in size and specificity. In gram-negative bacteria, the outer membrane necessitates receptor-mediated antagonistic activities, and very specific proteins are produced with domains for receptor binding, translocation, and activity (see [Colicins](#)). In contrast, gram-positive bacteria possess a multilayered peptidoglycan wall without an outer membrane. This favors peptides of small size that penetrate the murein network without receptor binding and specific translocation. Consequently, bacteriocins produced by gram-positive bacteria have a broad host range and are sometimes active on taxonomically unrelated genera. Thus, many of the bacteriocinlike agents produced by gram-positive bacteria kill species other than those likely to have the same ecological niche. Production of antibiotic peptides is the rule in gram-positive bacteria. Proteins the size of colicins are hardly ever encountered. In general the structures of these peptides must be stabilized by [posttranslational modifications](#).

Presently there is a burgeoning interest in bacteriocins because of the possible applications in the food industry (particularly lactic acid bacteria) and as a strategy for preventing certain infectious diseases. This review gives the general features of the main classes of bacteriocins and emphasizes the common features that are beginning to emerge.

An organism that produces a peptide antibiotic must be able to (1) synthesize the antibiotic, (2) export the antibiotic into the extracellular medium, (3) protect itself from the action of the antibiotic, and (4) interact with the sensitive cell and interfere with its growth or survival. This interaction may



or may not require entry of the antibiotic into the cell. Such entry is characteristic of some of the better-characterized ribosomally synthesized low molecular weight bacteriocins. For a detailed description of these bacteriocins, the reader is referred to the Suggestions for further reading.

## 1. Classes of Bacteriocins

A list of representative examples of lantibiotic and non-lantibiotic bacteriocins, together with their organisms of origin and mode of action is presented in Table 1. Four distinct classes of bacteriocins from lactic acid bacteria have been defined (5).

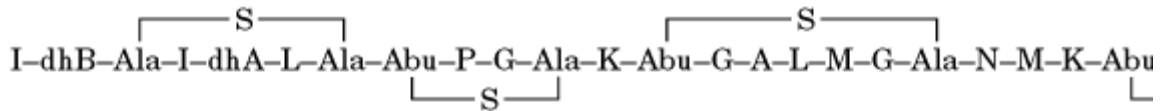
**Table 1. Representative Examples of Lanthionine-Containing Bacteriocins**

| Bacteriocin   | Type | Mass | Organism of Origin                    | Mode of Action (References)                   |
|---------------|------|------|---------------------------------------|---|
| Nisin A       | L-A  | 3353 | <i>Lactococcus lactis</i>             | Pore formation (8)                            |
| Pep5          | L-A  | 3488 | <i>Staphylococcus epidermis</i>       | Pore formation (44)                           |
| Subtilin      | L-A  | 3317 | <i>Bacillus subtilis</i><br>ATCC 6633 | Pore formation (45)                           |
| Epilancin K7  | L-A  | 3032 | <i>Staphylococcus epidermis</i>       | Pore formation (46)                           |
| Epidermin     | L-A  | 2164 | <i>Staphylococcus epidermis</i>       | Pore formation (47)                           |
| Gallidermin   | L-A  | 2164 | <i>Staphylococcus gallinarum</i>      | Pore formation (48)                           |
| Lacticin 481  | L-A  | 2901 | <i>Lactococcus lactis</i>             | Pore formation (49)                           |
| Streptococcin | L-A  | 2795 | <i>Streptococcus pyogenes</i>         | Pore formation (50)                           |
| A-FF22        |      |      |                                       |   |
| Salivaricin A | L-A  | 2315 | <i>Streptococcus salivarius</i>       | Pore formation (51)                           |
| Mutacin       | L-A  | 3245 | <i>Streptococcus mutans</i>           | Pore formation (52)                           |
| Lactocin S    | L-A  | 3764 | <i>Lactobacillus sake</i>             | Pore formation (36)                           |
| Carnocin U149 | L-A  | 4635 | <i>Carnobacterium piscicola</i>       | Pore formation (53)                           |
| Cytolysin L1  | L-A  | 4164 | <i>Enterococcus faecalis</i>          | Pore formation (54)                           |
| Cinnamycin    | L-B  | 2042 | <i>Streptomyces cinnamoneus</i>       | Membrane disorganization (55)                 |
| Duramycin     | L-B  | 2014 | <i>Streptomyces cinnamoneus</i>       | Membrane disorganization; pore formation (55) |
| Mersacidin    | L-B  | 1825 | <i>Bacillus</i> sp.                   | Cell-wall synthesis (56)                      |
| Actagardine   | L-B  | 1890 | <i>Actinoplanes</i> sp.               | Cell-wall synthesis (57)<br>inhibition        |

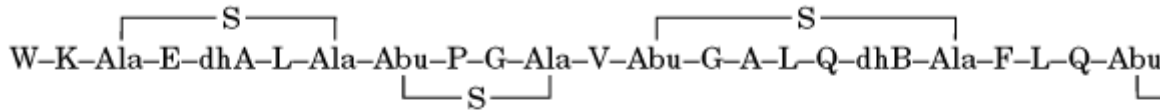
**Class** The lantibiotics (lanthionine-containing peptides) are small antibiotic peptides that are distinguished by their content of dehydro and thioether amino acids (lanthionine and 3-methyl-lanthionine) (6). Two types are distinguished by their distinctive ring structure. Type A comprises screw-shaped, amphipathic molecules that have two to seven net positive charges. Type B consists of more globular molecules that have a net positive charge or a net negative charge. The lantibiotics include a growing list of modified peptides that are active against Gram-positive bacteria: bacilli, lactococci, lactobacilli, staphylococci, streptococci, and streptomyces (Fig. 1). The function is twofold. First, these small peptide molecules are less likely to achieve stable conformations by themselves; they are generally extremely heat stable. Unusual cross-links like the thioether bridge from lanthionine help them retain their proper folded structure despite their small size (7). Secondly, unusual sidechains, such as the presence of a hydroxyaspartic acid, provide a repertoire of chemical reactivity that is very important for the biological activity of these peptides.

**Figure 1.** Primary structure of representative lantibiotics. Abbreviations used: dhA, didehydroalanine; dhB, didehydroβ-methyl lanthionine; Ala-NH-Lys, lysinoalanine; Asp-OH, hydroxyaspartic acid.

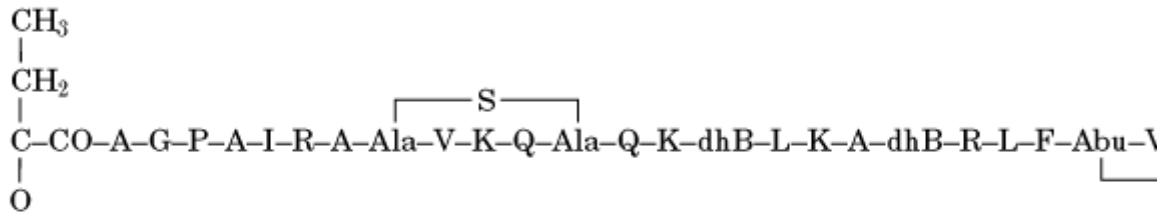
### Nisin A



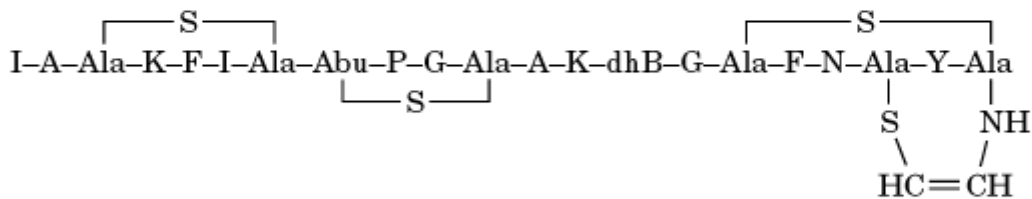
### Subtilin



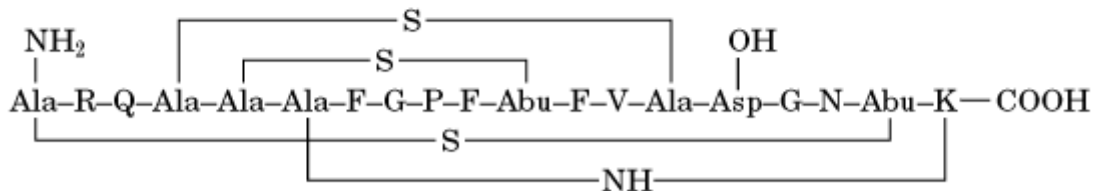
### Pep5



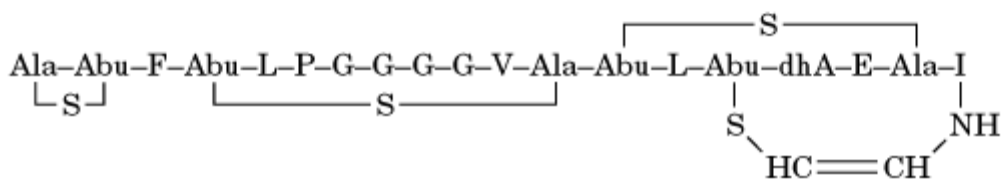
### Epidermin



### Cinnamycin



### Mersacidin



**Class II** These bacteriocins are small (<10 kDa), relatively heat stable, non-lanthionine-containing peptides. They have been subdivided into Class IIa, *Listeria*-active peptides with the N-terminal consensus sequence IIb, poration complexes requiring two different peptides for activity; Class IIc, thiol-activated peptides for activity. Unmodified peptide antibiotics (not containing lanthionine) share several characteristics: a leader peptide that has an N-terminal extension (**leader peptide**). Leader peptides are amphiphilic and contain a hydrophobic region of the **sec**-dependent exported proteins (see **protein targeting**). The leader peptide is particularly important at the proteolytic processing site, where two glycine residues are found at position (4, 9). The processed, active bacteriocins contain a substantial amount of **hydrophobic** amino acids. Accordingly, the plasma membrane is the target of most of these bacteriocins. Flanking regions of non-lanthionine peptides contain additional **open reading frames** for which a function has been demonstrated or postulated on the basis of sequence homologies. Such genes code for immunity proteins consisting of a histidine kinase and the respective regulator, and for proteins with signal transduction family. Some of these genes have an operon-like organization.

**Class III** These bacteriocins are large (>30 kDa) heat-labile proteins that include many bacteriolytic extracellular enzymes (e.g., **III** muramidases) that mimic the physiological action of bacteriocins.

**Class IV** These are complex bacteriocins that contain essential lipid or carbohydrate moieties in addition to protein. Purification of purified antibacterial components is still necessary to confirm that the presence of the additional biological activity of these molecules, thereby justifying the establishment of this class IV.

### 1.1. Microcins

An additional class of bacteriocins comprises a family of antibiotic substances called “microcins” produced by diverse members of the *Enterobacteriaceae* (10). They are distinguished from the majority of colicins by much lower molecular weight (<10 kDa) and because their synthesis is not induced by conditions that lead to induction of the **SOS repair** pathway. Instead, they are synthesized during the stationary phase, similar to that observed with most conventional antibiotics (11). The microcins were operationally defined as substances produced by gram-negative bacteria that pass through a cellophane membrane and inhibit the growth of an indicator *Escherichia coli* strain (12). Several substances known to be ribosomally synthesized fall within this group. The best known of these antibiotics peptides are microcin C7, microcin B17, and colicin V. Microcin C7 is a heptapeptide containing modifications at both the N- and C-termini that block protein synthesis *in vivo* and *in vitro* (13). Microcin B17 (Mcc B17) is a 43-residue peptide containing posttranslational modifications at the peptide backbone including serine, cysteine, and glycine residues, that result in thiazole and oxazole rings (14). Mcc B17 has much in common with lantibiotics because it is derived from a precursor peptide that is posttranslationally modified; 26 out of the 43 residues are glycine. Mcc B17 inhibits DNA replication and induces the SOS response (15). The target of this antibiotic is a DNA gyrase (16) (see [DNA Topology](#)). Colicin V (ColV) was first described in 1925 in the first report of an antibiotic substance produced by *E. coli* (17). It has a molecular weight of only 6000. The structure of mature colicin V is not totally known, even regarding side-chain modifications or N-terminal processing. Colicin V kills sensitive cells by disrupting their membrane potential (18).

## 2. Biosynthesis and Posttranslational Modifications

Generally the low molecular weight bacteriocins of gram-positive bacteria are first formed in an inactive precursor form that has a leader peptide. Following a variety of posttranslational modification reactions, the C-terminal propeptide domain is cleaved from the N-terminal leader sequence to yield the mature antimicrobial molecule. These prepeptides range from 18 to 30 amino acid residues and are not highly homologous except near the cleavage site. The prepeptides of all non-lanthionine-containing bacteriocins characterized thus far have two Gly residues at positions -2 and -1 relative to the processing site and **b-turn** promoting residues near the cleavage site. Possible functions of the prepeptide include (1) stabilizing the propeptide during [translation](#); (2) keeping the bacteriocin inactive, particularly after completion of modifications and thus helping to protect the producing strain; (3) allowing recognition of the ABC transporter system; (4) directing the precursor through a specific recognition motif toward biosynthetic enzymes (modification enzymes); (5) interacting with the propeptide region to stabilize a conformation that is essential for correct modification and thioether formation. There is evidence that the modifications of lantibiotics are made at the prepeptide stage. Cross-similarities between leaders of the modified and unmodified bacteriocins indicate that the overall features of leader peptides are important during regulation, synthesis, and generation of the immunity of peptide bacteriocins, regardless of whether or not they are modified.

**Table 2. Non-Lanthionine-Containing Bacteriocins**

---

**Mode of Action**

| Bacteriocin        | Type | Mass | Organism of Origin                         | (References)                          |
|--------------------|------|------|--|---------------------------------------|
| Pediocin PA-1      | NLC  | 4600 | <i>Pedicoccus acido lantici</i><br>PAC 1.0 | Pore formation ( <a href="#">58</a> ) |
| Mesentericin Y-105 | NLC  | 3666 | <i>Leuconostoc mesenterides</i> Y105       | Pore formation ( <a href="#">59</a> ) |
| Carnobacteriocin A | NLC  | 5100 | Carnobacterium piscicola LV17A             | ? ( <a href="#">60</a> )              |
| Sakacin A          | NLC  | 4308 | <i>Lactobacillus sake</i><br>LB706         | ? ( <a href="#">61</a> )              |
| Lactacin F         | NLC  | 5600 | <i>Lactobacillus</i> sp.                   | Pore formation ( <a href="#">62</a> ) |
| Lactococcin A      | NLC  | 5800 | <i>Lactococcus lactis</i><br>(cremoris)    | Pore formation ( <a href="#">63</a> ) |

The primary amino acid sequences of the propeptide components of many of the low molecular weight bacteriocins produced by gram-positive bacteria are known now. One important feature of the bacteriocins of gram-positive bacteria is their cysteine content. Those in which one or more cysteine residues are linked to dehydrated serine and threonine residues to form the thioether-linked amino acids lanthionine and methyl-lanthionine are called lantibiotics as previously mentioned (Fig. 1). Alternatively, bacteriocins in which pairs of cysteine residues undergo modification to form [disulfide bonds](#) are called cystibiotics (cysteine-containing antibiotics). A third subgroup of the bacteriocins, of which lactococcin B is an example, are designated thiolbiotics, because they contain only a single cysteine residue that must be present in the reduced **thiol** form to be active ([19](#)). At neutral pH, many of the low molecular weight bacteriocins are cationic. This is a unifying feature of lantibiotic and non-lanthionine containing bacteriocins and may have some significance for their activity.

Little is known about the posttranslational modification reactions and the biosynthetic enzymes involved in the maturation of lantibiotics. The lantibiotic B and C genes (designated *LanB* and *LanC*, as stipulated by the 1st and 2nd International Workshops on Lantibiotics) are probably involved in the maturation pathways, because mutation studies indicate that the gene products are essential for producing functional lantibiotics. The products of these genes, designated LanB and LanC, but also known as NisB and NisC, respectively, may be involved in the dehydration and thioether bond formation. Limited similarity was reported between LanB and *E. coli* IlvA, a threonine dehydratase, thus suggesting a similar function ([20](#)). This leads to the hypothesis that LanC and other similar proteins are involved in the enzyme-catalyzed formation of thioether bonds from the dehydrated residues and cysteine. Another modification enzyme that has been identified is EpiD, which is encoded by the epidermin gene cluster located on the 54-kb plasmid pTu32 ([21](#)). EpiD has no homologue in other lantibiotic gene clusters, and it catalyzes the biosynthesis of the C-terminal aminovinylcysteine residue of epidermin ([22](#)).

### 3. Extracellular and Maturation Release

After biosynthesis, class I and II bacteriocins are released to the extracellular medium. For each bacteriocin there is a relatively specific membrane protein system whose function is to translocate a precursor form of the antibiotic across the cytoplasmic membrane to the outside of the cell. The secretion of several peptide bacteriocins is mediated by dedicated transmembrane translocators belonging to the ATP-binding cassette (ABC) transporter superfamily ([23-26](#)). The transporters are encoded in the same operons as the bacteriocin structural gene or on a neighboring operon. The C-terminal ATP-binding domain and the N-terminal hydrophobic integral membrane domain are

expressed as a single polypeptide or as separate polypeptides. Recently, strong evidence has been reported that, very likely, all precursor peptides of the lantibiotic and non-lantibiotic bacteriocins that have leader peptides of the double glycine type (see above) are processed by their dedicated ABC-transporters, concomitant with export (27). It has also been proposed that the N-terminal domains of bacteriocin transporters are essential for initial recognition and subsequent proteolytic processing. These new proteolytic enzymes are [thiol proteinases](#). In fact, genes for potential ABC transporters have been found in all the known lantibiotic gene clusters. The lantibiotic transporters range in size from 535 to 714 amino acid residues and generally have the ATP-binding **domain** and the membrane-spanning domain in the same single protein. Sometimes a second ABC transporter system has been identified, which is not essential for production but participates in producer self-protection, thus providing some type of additional immunity (28). Although the leader sequence of the double-glycine type may be cleaved concomitant with export by the ABC-transporter itself, this is not the general case.

The ABC transporters have thiol proteinase activity, and putative proteinase genes (designated P) have been described for the leader peptides of other lantibiotics. These gene products share similarities with **subtilisin**-like [serine proteinases](#). For example, the *NisP* gene encodes an extracellular serine proteinase with a C-terminal extension anchoring it on the outer side of the cytoplasmic membrane (29). With nisin, it has been shown that cleavage of the leader peptide is not a prerequisite for export by the ABC-transporter, and this conclusion may be extended to other peptides with antimicrobial activity (30).

A detailed maturational pathway has been proposed for nisin (31). First the inducing signal (which could be nisin itself) activates the transcription at the *nisA* promoter (which has features also found in other positively regulated **promoter** sequences) via a two-component response-regulator system (32). This results in the production of pre-nisin containing free cysteines and no dehydrated residues. The pre-nisin is directed, presumably by virtue of the leader peptide, to a membrane-located complex containing the modifying enzymes LanB (probably involved in dehydration) and LanC (probably involved in thioether bonds, as previously mentioned). At this stage, the leader helps to maintain the peptide in an inactive form. Subsequently, the precursor nisin is translocated via the ABC exporter NisT at the expense of ATP hydrolysis. Finally, precursor nisin is activated by proteolytic cleavage by the extracellular protease NisP attached to the outside of the cell envelope. The role of the leader is still not fully understood, apart from its function in producing an inactive conformation. In particular, it is not known how it functions in targeting the pre-lantibiotic to the maturation and export proteins, what its fate is after cleavage, and how it contributes to processes, such as self protection, after it has been cleaved off.

This proposed sequence of events may also apply to other lantibiotics. However, additional intracellular conversions are required in specific modification reactions, such as those involving the N- and C-termini of epidermin (33), lantibiotic Pep5, epilamin K7 (34, 35), and lactocin S (36). In these cases it is likely that the leader peptide cleavage occurs intracellularly.

#### 4. Modes of Action

The various structural properties of the subgroups of bacteriocins are reflected in their three different modes of action. The primary activity of Class I type A lantibiotics is based on forming voltage-dependent, short-lived pores in the cytoplasmic membrane (37, 38). The peptides rapidly induce leakage of ions and small metabolites from bacterial cells and a collapse of the electrochemical proton gradient, leading to cessation of biosynthetic processes and eventually to cell death. Type A lantibiotics require a membrane potential of between 50 and 100 mV for pore formation, depending on the individual peptide. It is assumed that pores are formed by a transiently associated peptide oligomer in a transmembrane orientation (barrel stave model), as suggested for alamethicin. The susceptibility toward a particular peptide of different bacterial species, or even of strains within one species, varies much more than one would expect on the basis of the pore-formation model. *In vivo* pore formation or pore stability may be positively or negatively influenced by such factors as

phospholipid composition of the membrane, interactions of the peptides with integral membrane components, or the presence of surface layers. The models may need to be refined, and pore formation may depend on local perturbation of the bilayer (39). In addition, the peptides could exert secondary effects that contribute to bactericidal activity, such as autolysis of cells by activating cell-wall hydrolyzing enzymes (40).

The antibacterial effects of class I type B lantibiotics are rather weak. They bind to the head groups of phospholipids, preferentially to phosphoethanolamine, which may cause membrane permeabilization and allow the release of cations and other small solutes. Duramycin also induces the formation of complex pores (41). Mersacidin and actagardine are distinguished from other type B lantibiotics by their mode of action. Both interfere with cell-wall biosynthesis in gram-positive bacteria, which may eventually offer new possibilities in antimicrobial therapy (42).

Generally, the bactericidal action of class II bacteriocin against sensitive cells is produced principally by destabilizing membrane function, such as energy transduction, rather than disrupting the structural integrity of the membrane. This effect results from the energy-independent dissipation of the [proton motive force](#) and loss of the permeability barrier of the cytoplasmic membrane. It contrasts with the energy-dependent bactericidal action of the lantibiotics. In addition, before pores are formed, all of the non-lanthionine containing bacteriocins interact with membrane-associated receptor proteins, in contrast to class I lantibiotic bacteriocins.

## 5. Immunity of Bacteriocin-Producing cells

One of the definitive features of bacteriocin-producing cells is their ability to resist the action of their own inhibitory substances through a specific immunity mechanism. Such a mechanism is based on dedicated peptides or proteins called immunity proteins, which specifically antagonize the bacteriocin. For nisin and subtilisin, their antagonist proteins of 245 and 165 amino acid residues, respectively, display the features of bacterial lipoproteins and are similar in amino acid sequence. Despite the high degree of similarity of these lantibiotics, there is no cross-immunity between producing cells. In other cases, as with pepS, the immunity is much more closely related to the immunity systems of unmodified bacteriocins, which are short, have only 69 residues, and have a hydrophobic N-terminal segment and a strongly hydrophilic C-terminal part.

Currently, there are no clues to the mechanism of the immunity phenomenon. The proposed location of the peptides outside the cell and the observation that the cytoplasmic membrane of immune clones is not depolarized by externally added bacteriocin suggest that the bacteriocins are directly antagonized by the immunity peptide, like colicin A, for example (43).

## Bibliography

1. H. Kleinhauf and H. von Döhren (1990) *Eur. J. Biochem.* **192**, 1–15.
2. M. Kageyama, M. Kobayashi, Y. Sano, and H. Masaki (1996) *J. Bacteriol.* **178**, 103–110.
3. H. Bénédicti and V. Géli (1996) In *Handbook Biological Physics* (W. N. Konings, H. R. Kaback, and J. S. Lolkema eds.), Elsevier, Amsterdam, The Netherlands, Vol. **2**, pp. 665–691.
4. H. G. Sahl (1994) In *Antimicrobial Peptides* (Ciba Foundation Symposium 186), Wiley, Chichester, pp. 27–53.
5. T. R. Klaenhammer (1988) *Biochimie* **70**, 337–349.
6. G. Jung (1991) *Angew. Chem. Int. Ed. Engl.* **30**, 1051–1068.
7. E. Gross and J. L. Morell (1967) *J. Amer. Chem. Soc.* **53**, 2791–2792.
8. E. Gross and J. L. Morell (1971) *J. Amer. Chem. Soc.* **93**, 4634–4635.
9. W. M. de Vos, O. P. Knipers, J. R. van der Meer, and R. J. Siezen (1995) *Mol. Microbiol.* **17**, 427–437.
10. F. Baquero and F. Moreno (1984) *FEMS Microbiol. Lett.* **23**, 117–124.

11. J. F. Martin and A. L. Demain (1980) *Microbiol. Rev.* **44**, 230–251.
12. C. Asensio, C. Perez-Diaz, M. C. Martinez, and F. Baquero (1976) *Biochem. Biophys. Res. Commun.* **69**, 7–14.
13. R. Kolter and F. Moreno (1992) *Ann. Rev. Microbiol.* **46**, 141–163.
14. A. Bayer, S. Freund, C. Nicholson, and C. Jung (1993) *Angew. Chem. Int. Ed. Engl.* **32**, 1336–1339.
15. M. Herrero and F. Moreno (1986) *J. Gen. Microbiol.* **132**, 393–402.
16. J. L. Vizan, C. Hernandez-Chico, I. del Castillo, and F. Moreno (1991) *EMBO J.* **10**, 467–476.
17. A. Gratia (1925) *C.R. Soc. Biol.* **93**, 1040–1041.
18. C. Yang and J. Konisky (1984) *J. Bacteriol.* **158**, 757–759.
19. G. Bierbaum and H. Sahl (1991) In *Nisin and Novel Antibiotics* (G. Jung and H. G. Sahl, eds.), Escom Publishers, Leiden, The Netherlands, pp. 386–396.
20. Z. Gutowski-Eckel, C. Klein, K. Siegers, K. Bohm, M. Hammelmann, and K. D. Entian (1994) *Appl. Environ. Microbiol.* **60**, 1–11.
21. N. Schnell, G. Engelke, R. Augustin, F. Rosenstein, F. Götz, and K. D. Entian (1991) In *Nisin and Novel Antibiotics* (G. Jung and H. G. Sahl, eds.), Escom Publishers, Leiden, The Netherlands, pp. 269–276.
22. T. Kupke, S. Stefanovic, H. G. Sahl, and F. Götz (1992) *J. Bacteriol.* **174**, 5354–5361.
23. L. Gilson, H. K. Mahanty, and R. Kolter (1990) *EMBO J.* **9**, 3875–3884.
24. M. S. Gilmore, R. A. Segarra, and M. C. Booth (1990) *Infect. Immunol.* **58**, 3914–3923.
25. J. D. Marugg, C. F. Gonzalez, B. S. Kunka, A. M. Ledeboer, M. J. Pucci, M. Y. Toonen, S. A. Walker, L. C. Zoetmulder, and P. A. Vandebergh (1992) *Appl. Environ. Microbiol.* **58**, 2360–2367.
26. G. W. Stoddard, J. P. Petzel, M. J. van Belkum, J. Kok, and L. L. McKay (1992) *Appl. Environ. Microbiol.* **58**, 1952–1961.
27. L. S. Havarstein, H. Holo, and I. F. Nes (1994) *Microbiol.* **140**, 2393–2389.
28. K. Venema, G. Venema, and J. Kok (1995) *Trends Biochem. Sci.* **3**, 299–304.
29. J. R. van der Meer, J. Polman, M. M. Beerthuyzen, R. J. Siezen, O. P. Kuipers, and W. M. de Vos (1993) *J. Bacteriol.* **175**, 2578–2588.
30. J. R. van der Meer, H. S. Rollema, R. J. Siezen, M. M. Bethuyzen, O. P. Kuipers, and W. M. de Vos (1994) *J. Biol. Chem.* **269**, 3555–3562.
31. W. M. de Vos, O. P. Kuipers, J. R. van der Meer, and R. Siezen (1995) *Mol. Microbiol.* **17**, 427–437.
32. W. M. de Vos and G. F. Simons (1994) In *Genetics and Biotechnology of Lactic Acid Bacteria* (M. J. Gasson and W. M. de Vos, eds.), Backie Academic, Glasgow, pp. 52–105.
33. T. Kupke, C. Kempter, V. Gnan, G. Jung, and F. Götz (1994) *J. Biol. Chem.* **269**, 5653–5659.
34. M. Reis, M. Eschbach-Bludau, M. Iglesias-Wind, T. Kupke, and H. G. Sahl (1994) *Appl. Environ. Microbiol.* **60**, 2876–2883.
35. M. van de Kamp, H. W. van den Hooven, R. N. Konings, C. W. Hilbers, C. W. van de Ven, G. Bierbaum, H. G. Sahl, O. P. Kuipers, R. J. Seizen, and W. M. de Vos (1995) *Eur. J. Biochem.* **230**, 587–600.
36. M. Skangen, J. Nissen-Meyer, G. Jung, S. Stefanovic, K. Sletten, C. I. Mortreveltd-Abilgaard, and I. F. Nes (1994) *J. Biol. Chem.* **269**, 27183–27185.
37. H. G. Sahl (1991) In *Nisin and novel antibiotics Proceedings of the First International Workshop on Lantibiotics* (G. Jung and H. G. Sahl, eds.), Leiden, Escom Publishers, pp.347–359.
38. R. Benz, G. Jung, and H. G. Sahl (1991) In *Nisin and novel antibiotics Proceeding of the First International Workshop on Lantibiotics* (G. Jung and H. G. Sahl, eds.), Leiden, Escom



Publishers, pp. 359–372.

39. A. J. Driessen, H. W. van den Hooven, W. Kuiper, M. van de Kamp, H. G. Sahl, R. N. Konings, and W. Konings (1995) *Biochemistry* **34**, 1606–1614.
40. G. Bierbaum and H. G. Sahl (1987) *J. Bacteriol.* **169**, 5452–5458.
41. T. Sheth, R. M. Henderson, S. B. Hlady, and W. A. Cuthbert (1992) *Biochim. Biophys. Acta* **1107**, 179–185.
42. S. Chatterjee, D. Chatterjee, K. Jani, H. Blumbach, B. Ganguli, N. Klesel, M. Limbert, and G. Seibert (1992) *J. Antibiot.* **45**, 839–845.
43. D. Espeset, D. Duché, D. Baty, and V. Géli (1996) *EMBO J.* **15**, 2356–2364.
44. R. Kellner, G. Jung, and H. G. Sahl (1991) In *Nisin and Novel Lantibiotics* (G. Jung and H. G. Sahl, eds.), Escom Publishers, Leiden, The Netherlands, pp. 141–158.
45. E. Gross, H. H. Kitz, and E. Nebelin (1973) *Hoppe-Seyler's Z. Physiol. Chem.* **354**, 810–812.
46. M. van de Kamp, L. Horstink, H. W. van den Hooven, R. N. Konings, C. W. Hilbers, A. Frey, H. G. Sahl, J. W. Metzger, and F. J. van de Ven (1995) *Eur. J. Biochem.* **227**, 757–771.
47. H. Allgaier, G. Jung, R. Werner, U. Schneider, and H. Zähler (1986) *Eur. J. Biochem.* **160**, 9–22.
48. R. Kellner, G. Jung, T. Hörner, H. Zähler, N. Schnell, K. Entian, and F. Götz (1988) *Eur. J. Biochem.* **177**, 53–59.
49. J. C. Piard, C. Delorme, M. Novel, M. Desmageaud, and G. Novel (1993) *FEMS Microbiol. Lett.* **112**, 313–318.
50. R. W. Jack, A. Carne, J. Metzger, S. Stefanovic, H. G. Sahl, G. Jung, and J. R. Tagg (1994) *Eur. J. Biochem.* **220**, 455–462.
51. K. F. Ross, C. Ronson, and J. R. Tagg (1993) *Appl. Environ. Microbiol.* **59**, 2014–2021.
52. J. Novak, P. W. Canfield, and E. J. Miller (1994) *J. Bacteriol.* **176**, 4316–4320.
53. G. Stoffels, J. Niessen-Meyer, A. Gudmundsdottir, K. Sletten, H. Halo, and I. F. Nes (1992) *Appl. Environ. Microbiol.* **58**, 1417–1422.
54. M. S. Gilmore, R. A. Serraga, M. C. Booth, C. P. Bogie, L. R. Hall, and D. B. Clewell (1994) *J. Bacteriol.* **176**, 7335–7344.
55. A. Fredenhagen, F. Märki, G. Fendich, W. Märki, J. Gruner, J. van Oostrum, F. Raschdorf, and H. H. Peter (1991) In *Nisin and Novel Lantibiotics* (C. Jung and H. G. Sahl eds), Escom Publishers, Leiden, The Netherlands, pp. 131–140.
56. H. Kogler, H. Bauch, H. W. Fehlhaber, C. Griesinger, W. Schubert, and V. Teetz (1991) In *Nisin and Novel Lantibiotics* (G. Jung and H. G. Sahl, eds.), Escom Publishers, Leiden, The Netherlands, pp. 159–170.
57. N. Zimmermann, N. S. Freud, A. Fredenhagen, and G. Jung (1993) *Eur. J. Biochem.* **216**, 419–428.
58. M. L. Chikindas, M. J. Garcia-Garcera, A. J. Driessen, A. M. Ledeboer, J. Niessen-Meyer, I. F. Nes, T. Abee, W. N. Konings, and G. Venema (1993) *Microbiology* **59**, 3577–3584.
59. A. Maftah, T. Renault, C. Viguoles, Y. Hechard, P. Bressolier, M. Ratineaud, Y. Cenatiempo, and R. Julian (1993) *J. Bacteriol.* **175**, 3232–3235.
60. R. W. Woboro, T. Henkel, M. Sailer, K. L. Roy, J. C. Vederas, and M. E. Skiles (1994) *Microbiology* **140**, 517–526.
61. A. Holck, L. Axelson, S. S. E. Birkeland, T. Aukrust, and H. Bloom (1992) *J. Gen. Microbiol.* **138**, 2715–2720.
62. P. Muriana and T. R. Klaenhammer (1991) *J. Bacteriol.* **173**, 1779–1788.
63. H. Holo, O. Niessen, and I. F. Ness (1991) *J. Bacteriol.* **173**, 3879–3887.

### Suggestions for Further Reading

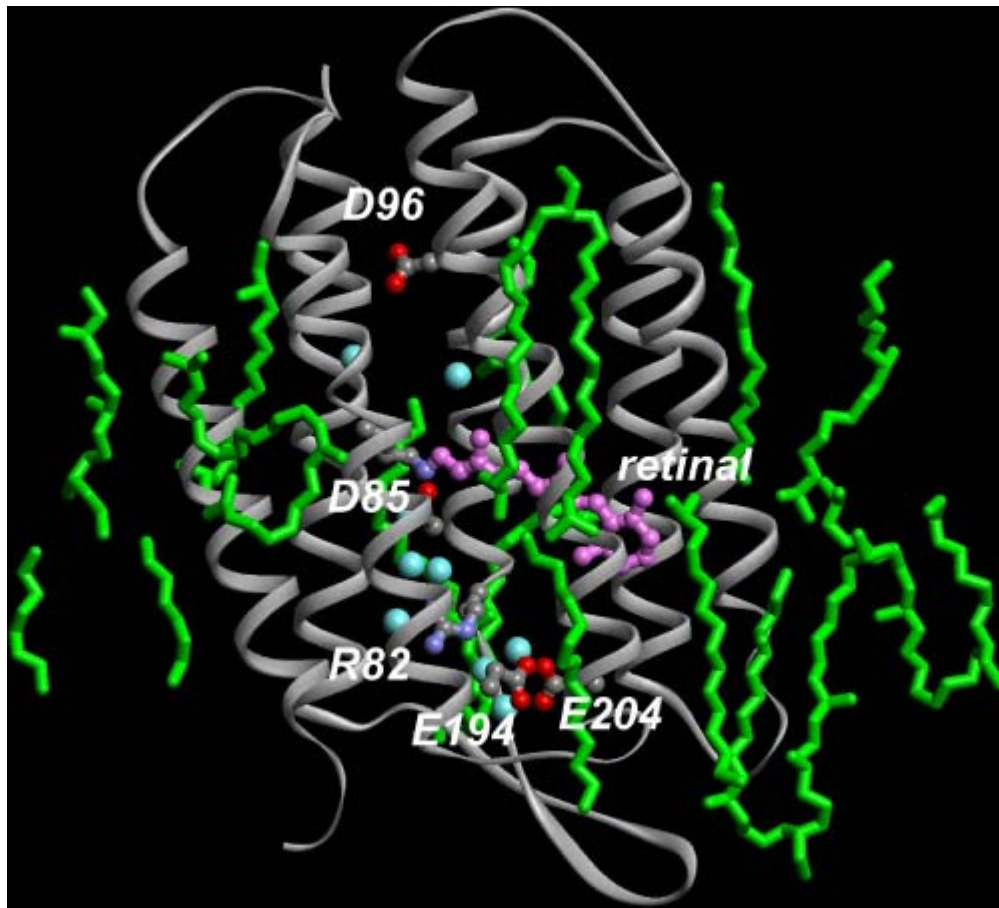
64. R. James, C. Lazdunski, and F. Pattus, eds. (1992) *Bacteriocins, Microcins and Lantibiotics*, NATO ASI Series, Springer-Verlag, Berlin, Vol. **H65**.
65. R. W. Jack, J. R. Tagg, and B. Ray (1995) Bacteriocins of gram-positive bacteria, *Microbiol. Rev.* **59**, 171–200.
66. T. R. Klaenhammer (1993) Genetics of bacteriocins produced by lactic acid bacteria, *FEMS Microbiol. Rev.* **12**, 39–86.
67. F. Moreno, J. L. San Millan, C. Hernandez-Chico, and R. Kolter (1995) Microcins, *Biotechnology*, **28** 307–321.
68. H-G. Sahl, R. W. Jack, and G. Bierbaum (1995) Biosynthesis and biological activities of lantibiotics with unique post-translational modifications, *Eur. J. Biochem.* **230**, 827–853.

## Bacteriorhodopsin

Bacteriorhodopsin is a small (26 kDa) integral [membrane protein](#), the prototype of seven-helical **G-protein**–linked **receptors**, that upon illumination transports protons across the membrane. It forms extended two-dimensional hexagonal arrays in the cytoplasmic membrane of **halobacteria**. Its transmembrane  [\$\alpha\$ -helices](#) surround the **prosthetic group** retinal, which is linked via a [Schiff Base](#) to Lys216 near the middle of helix G and lies at a small angle to the membrane plane. Photoisomerization of the retinal from all-*trans* to 13-*cis* sets off a sequence of thermal reactions (the “photocycle”) in which the interaction of the retinal and the protein causes proton transfers between various donor and acceptor groups. Together, these transfers result in the complete translocation of a proton from the cytoplasmic to the extracellular surface, thus generating a transmembrane electrochemical gradient for protons. The [proton gradient](#) is utilized in the way usual in **bacteria** for the synthesis of ATP, the uptake of nutrients (amino acids) and  $K^+$ , and the transport of  $Na^+$  out of the cells (see [Proton Motive Force](#)).

Bacteriorhodopsin forms trimers that assemble in the two-dimensional hexagonal array that constitutes the patches termed “**purple membrane**.” The purple membrane contains only bacteriorhodopsin and [lipids](#), and the regular crystalline lattice made it possible to determine its structure by cryo-**electron crystallography** at 7 Å (1) and then 3 Å (2) resolution. The protein was also crystallized from a lipid cubic phase, and its structure has now been determined by [X-ray crystallography](#) at 2.5 Å resolution (3). The protein consists of seven transmembranous  $\alpha$ -helices, with short interhelical loops and short N and C termini. Three of the helices, B, C, and D, are normal to the plane of the membrane and the other four, A, E, F, and G, are inclined at various small angles to the perpendicular (Fig. 1). The retinal is bound to the  $\epsilon$ -**amino group** of Lys216, forming a protonated Schiff base near the middle of helix G, and its polyene chain lies at about 23° from the membrane plane. Thus the Schiff base divides the protein into extracellular and cytoplasmic halves.

**Figure 1.** Structure of bacteriorhodopsin (1), and the pathway of proton transport. The seven transmembranous  $\alpha$ -helices are shown, along with only the all-*trans* retinal and the most important residues. The curved arrows identify the proton transfers that occur at different times in the photocycle (see text) and add up to the complete transport of a proton from the cytoplasmic (upper) to the extracellular surface of the membrane.



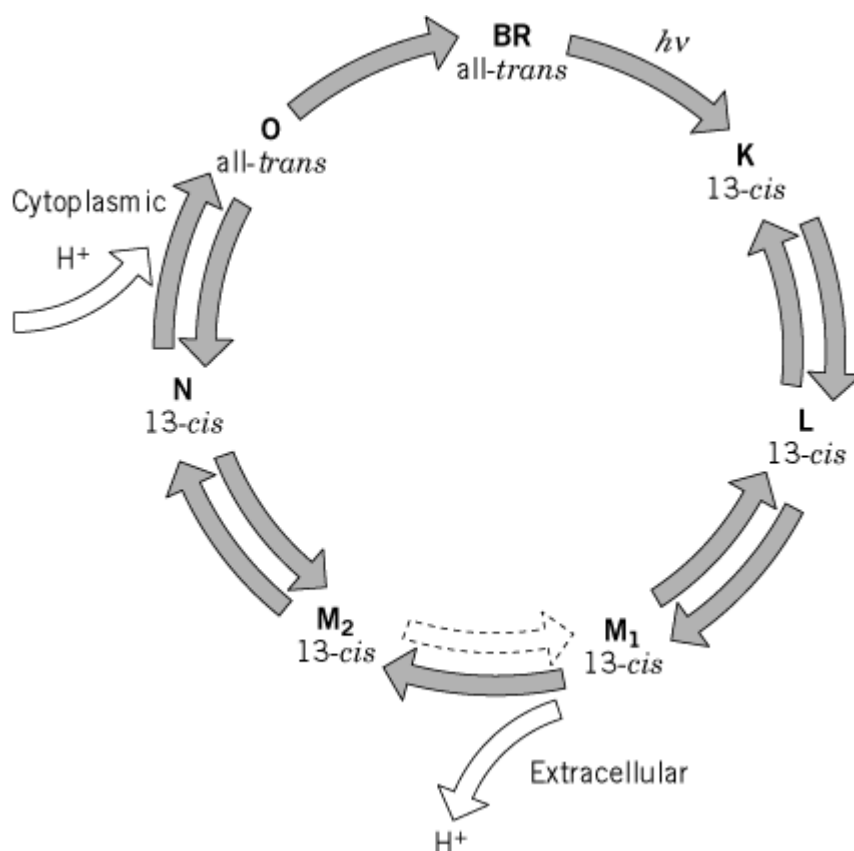
The trajectory of the transported proton from one surface to the other is through proton-conductive “half-channels” through each of the two halves of the molecule. Identification of the residues that participate in the two half-channels has been the objective of much research in the past few years. The extracellular half-channel is quite complex and contains [polar](#) and [hydrogen-bonding](#) residues. Many of them play roles in the release of protons to the extracellular surface. Asp85 and Asp212 are anionic residues located near the Schiff base. A possible pathway for protons leads from Asp85, via Arg82, Glu204, and Glu194, to the extracellular surface and is consistent with functional studies. Another plausible pathway is through Glu9 or Glu74, but the results of [mutagenesis](#) indicate that these two residues are dispensable. The cytoplasmic half-channel is simpler. It contains mostly **hydrophobic** residues between the retinal and Asp96, which is the proton donor in the reprotonation of the Schiff base. Thr46 is close enough to interact with Asp96, and mutagenesis confirms its significant role in the proton donor and acceptor function of Asp96. Acidic residues at the cytoplasmic aqueous interface include residues 36, 38, 102, and 104. It has been suggested that they form a “funnel” at the surface that directs protons into the cytoplasmic half-channel (2). Although immobilization of the protein in the purple membrane is responsible for the well-known, extraordinary thermal and photostability of bacteriorhodopsin, neither the trimeric arrangement of the monomers nor the rigid lattice structure have functional roles in the transport activity (4). Thus, availability of the structure of bacteriorhodopsin, even at the current limited resolution, has generated considerable insights into the proton transport mechanism and has provided the framework for mechanistically interpreting a wealth of spectroscopic information.

The chromophore with all-*trans*, 15-*anti* retinal has a broad **absorption** band with a maximum at 568 nm. This is considerably red-shifted from that of retinal in isolation (380 nm) or a retinal analog with a protonated Schiff base (440 nm). This red-shift results primarily from the diffuse counterion to the charged Schiff base, which comprises principally Asp85 with a contribution from Asp212,

which is several angstroms removed and providing only partial compensation of the charge (5). The all-*trans* chromophore exists in thermal equilibrium with the 13-*cis*, 15-*syn* configuration, which absorbs at 555 nm. Sustained illumination converts the retinal to 100% all-*trans*, 15-*anti*, known as “light-adaptation” (6). The photoreaction of this isomer is normally active in transport.

The photocycle of the all-*trans* chromophore is described by the intermediate states J, K, L, M, N, and O, their substates, and the sequence of their interconversions (Fig. 2). Each intermediate is characterized by a distinct absorption maximum in the visible wavelength region, numerous vibrational bands of the retinal in the infrared and Raman region and of the protein in the infrared region (see [Vibrational Spectroscopy](#)). The kinetics of the photocycle describe the sequence and the energetics of the chemical reactions that translocate protons across the membrane. In the K intermediate, which arises on the nanosecond timescale after decay of the excited state, the retinal assumes a twisted 13-*cis*, 15-*anti* configuration, as indicated by high-amplitude hydrogen-out-of-plane vibrations. Its absorption maximum is red-shifted by 30 to 40 nm from the initial state. It is converted in about one microsecond into the L state, which absorbs at 540 to 550 nm. In L, the polyene of the retinal is more relaxed, but other changes begin to appear in the protein in both the extracellular and cytoplasmic regions (7). The Schiff base forms stronger hydrogen bonds that are correlated with structural changes of bound [water](#) detected in the infrared wavelength region. One of these water molecules is bound to the Schiff base, Asp85, and Asp212 and may play a role in the  $pK_a$  shifts that result in proton transfer from the Schiff base to Asp85. The hydrogen bonding of two other water molecules near Asp96 are also affected in L, suggesting that the structural changes near the Schiff base are transmitted all the way across the protein to the cytoplasmic region.

**Figure 2.** The photocycle of bacteriorhodopsin. The intermediate states are shown, and the isomeric configuration of the retinal is indicated. The Schiff base is protonated in all but the M states.



The M state is formed by transfer of a proton from the retinal Schiff base to Asp85 (8). It has a strongly blue-shifted absorption maximum at 410 nm. The kinetics of this conversion, measured at visible and infrared wavelengths, suggest that L is in equilibrium with an early M state,  $M_1$ , and that the mixture of L and  $M_1$  decay together to form the late M state,  $M_2$  (9-11). At  $\text{pH} > 6$  this occurs in a unidirectional reaction, and thus L disappears as  $M_2$  is formed. At  $\text{pH} < 6$ , however, the  $M_1 \rightarrow M_2$  reaction is not unidirectional, and both L and  $M_1$  remain present and coexist with  $M_2$  because at the higher pH a proton dissociates from a site in the extracellular region that interacts with Asp85. This site is either Glu204 or depends on the carboxyl group of Glu204. Its proton is passed to Glu194 and is released from there to the extracellular surface. The anomalous titration properties of Asp85 indicate that the nature of the interaction between the proton affinities of this aspartate residue and the proton release site is such that either may be protonated but not both (12, 13). When Asp85 becomes protonated by the Schiff base, the  $\text{p}K_a$  of the proton release group is lowered, and the proton dissociates. When this proton is released to the bulk solution at a pH greater than the  $\text{p}K_a$  for the release, the  $\text{p}K_a$  of Asp85, in turn, is driven higher, and deprotonation of the Schiff base becomes complete. This prevents reprotonation of the Schiff base from the extracellular direction.

Reprotonation of the Schiff base is from the cytoplasmic direction by Asp96 (8, 14), which produces the N intermediate that absorbs near 560 nm but with a lower extinction than the initial state. Large-scale protein conformational changes in N are evident from a pair of negative and positive difference features in the infrared spectrum that originate from a shift of the amide I band (15). Electron and X-ray diffraction studies of the M and the N states, measured either at various times after flash illumination and freezing or in a photostationary state at ambient temperature, indicate considerable changes of conformation at the cytoplasmic surface. The most conspicuous of these is an outward tilt of the cytoplasmic end of helix F (16). Its occurrence in the M intermediate may be transitory during the  $M \rightarrow N$  conversion of the unperturbed wild-type photocycle, but this feature is clearly observable when M is stabilized in the wild-type or in mutant forms of the protein. The movement of helix F as a rigid body is confirmed by distance measurements using pairs of spin labels (17). The effects of osmotic agents, humidity, and in-plane cooperativity in the purple membrane lattice on the  $M \rightarrow N$  reaction and on the protein conformation change suggest that the rationale of the helical tilt is to increase the hydration of the cytoplasmic region and thereby to decrease the  $\text{p}K_a$  of Asp96. Thus, Asp96 becomes a proton donor to the Schiff base. The tilt of helix F is recovered during decay of the N state, presumably recovering the initial high  $\text{p}K_a$  of Asp96 and causing its reprotonation from the cytoplasmic surface.

Reisomerization of the retinal to all-*trans* occurs in the  $N \rightarrow O$  transition (18). This is made possible by the lowered barrier-to-bond rotation in the polyene chain upon protonation of the Schiff base. Residues that contact the chain near the 9-methyl and 13-methyl groups, such as Trp182 (19) and Leu93 (20), facilitate the reisomerization, probably through steric interactions that transmit residue displacements in the protein to the retinal and vice versa. The O state has a strongly red-shifted maximum at visible wavelengths, at least partly because Asp85 is still protonated. Consequently, the main component of the counterion to the protonated Schiff base is lacking. Large-amplitude hydrogen-out-of-plane vibrations indicate that, as in the K state, the retinal chain is twisted. These features disappear in the final  $O \rightarrow \text{BR}$  reaction, which appears to be limited by the rate of proton transfer from Asp85 to the still unprotonated proton release site (21). As expected from the recovery of the low initial  $\text{p}K_a$  of Asp85, this reaction is unidirectional under all conditions, and it ensures the full repopulation of the initial state and also the functioning of the **proton pump** against large transmembranous proton gradients.

Bacteriorhodopsin is one of three types of similar retinal proteins in halobacterial membranes. Their functions are all based on the photoisomerization of all-*trans* retinal to 13-*cis*, 15-*anti* and the protein reactions that accompany the thermal reisomerization. Halorhodopsin is an inwardly-directed, light-

driven chloride ion pump. It lacks Asp85 and Asp96, and the retinal Schiff base does not deprotonate during the photocycle (22). Sensory [rhodopsins](#) I and II are receptors for phototactic behavior (23). A profound similarity in the mechanisms of these proteins with different functions is indicated by the fact that their activities are interconvertible with minimal perturbations. Thus, the Asp85Thr mutant of bacteriorhodopsin binds chloride, exhibits a photocycle similar to that of halorhodopsin, and transports chloride from the extracellular to the cytoplasmic direction (24). Halorhodopsin, in turn, transports protons when the weak acid, azide, is added, by binding near the Schiff base and functioning as a proton acceptor (25). Sensory rhodopsin I transports protons like bacteriorhodopsin when the transducing protein that is normally tightly bound to it is genetically deleted (26, 27).

## Bibliography

1. N. Grigorieff et al. (1996) *J. Mol. Biol.* **259**, 393–421.
2. Y. Kimura et al. (1997) *Nature* **389**, 206–211.
3. E. Pebay-Peyroula, G. Rummel, J. P. Rosenbusch, and E. M. Landau (1997) *Science* **277**, 1676–1681.
4. N. A. Dencher and M. P. Heyn (1979) *FEBS. Lett.* **108**, 307–310.
5. K. Nakanishi et al. (1980) *J. Am. Chem. Soc.* **102**, 7945–7947.
6. G. S. Harbison et al. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1706–1709.
7. A. Maeda et al. (1997) *J. Biochem. (Tokyo)* **121**, 399–406.
8. M. S. Braiman et al. (1988) *Biochemistry* **27**, 8516–8520.
9. L. Zimányi et al. (1992) *Biochemistry* **31**, 8535–8543.
10. S. Dickopf and M. P. Heyn (1997) *Biophys. J.* **73**, 3171–3181.
11. B. Hessling, J. Herbst, R. Rammelsberg, and K. Gerwert (1997) *Biophys. J.* **73**, 2071–2080.
12. S. P. Balashov, E. S. Imasheva, R. Govindjee, and T. G. Ebrey (1996) *Biophys. J.* **70**, 473–481.
13. H. T. Richter, L. S. Brown, R. Needleman, and J. K. Lanyi (1996) *Biochemistry* **35**, 4054–4062.
14. K. Gerwert, B. Hess, J. Soppa, and D. Oesterhelt (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4943–4947.
15. M. S. Braiman, O. Bousché, and K. J. Rothschild (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2388–2392.
16. S. Subramaniam, M. Gerstein, D. Oesterhelt, and R. Henderson (1993) *EMBO J.* **12**, 1–8.
17. T. E. Thorgerisson et al. (1997) *J. Mol. Biol.* **273**, 951–957.
18. S. O. Smith et al. (1983) *Biochemistry* **22**, 6141–6148.
19. O. Weidlich et al. (1996) *Biochemistry* **35**, 10807–10814.
20. J. K. Delaney, G. Yahalom, M. Sheves, and S. Subramaniam (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5028–5033.
21. H. T. Richter et al. (1996) *Biochemistry* **35**, 15461–15466.
22. D. Oesterhelt (1995) *Israel J. Chem.* **35**, 475–494.
23. W. D. Hoff, K. H. Jung, and J. L. Spudich (1997) *Annu. Rev. Biophys. Biomol. Struct.* **26**, 223–258.
24. J. Sasaki et al. (1995) *Science* **269**, 73–75.
25. G. Váró, L. S. Brown, R. Needleman, and J. K. Lanyi (1996) *Biochemistry* **35**, 6604–6611.
26. R. A. Bogomolni et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10188–10192.
27. U. Haupts, C. Haupts, and D. Oesterhelt (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3834–3838.

## Suggestions for Further Reading

28. T. G. Ebrey (1993) "Light energy transduction in bacteriorhodopsin". In *Thermodynamics of Membranes, Receptors and Channels*, (M. Jackson, ed.), CRC Press, Boca Raton, FL, pp. 353–387.

29. J. K. Lanyi (1993) Proton translocation mechanism and energetics in the light-driven pump bacteriorhodopsin. *Biochim. Biophys. Acta* **1183**, 241–261.
30. J. K. Lanyi and G. Váró (1995) The photocycles of bacteriorhodopsin. *Israel J. Chem.* **35**, 365–386.
31. R. A. Mathies, S. W. Lin, J. B. Ames, and W. T. Pollard (1991) From femtoseconds to biology: Mechanism of bacteriorhodopsin's light-driven proton pump. *Ann. Rev. Biophys. Biophys. Chem.* **20**, 491–518.
32. D. Oesterhelt, J. Tittor, and E. Bamberg (1992) A unifying concept for ion translocation by retinal proteins. *J. Bioenerg. Biomembr.* **24**, 181–191.

## BAL 31 Nuclease

The [enzymes](#) commonly referred to as BAL 31 nuclease have been widely used for manipulation of nucleic acids, most frequently through their ability to shorten both strands of double-stranded DNA from the ends in a controllable manner. Other activities include endonucleolytic attack on double-stranded DNA in response to the presence of a variety of covalent and noncovalent alterations in DNA structure (eg, DNA-carcinogen adducts and junctions between right- and left-handed helical regions) and highly processive unidirectional exonucleolytic degradation of single-stranded DNA.

### 1. Source and Some Physical Properties

BAL 31 is the strain designation given the original isolates of the marine bacterium that produces the nucleases, which was originally determined to be a species of *Pseudomonas* (1) but was reclassified as *Alteromonas espejiana* (2). The **nuclease** activity was discovered adventitiously as a contaminant in preparations of the *Alteromonas* bacteriophage PM2, but these enzymes proved to be bacterial products secreted into the culture medium (3, 4). The common name also appears as BAL31 and BAL-31; all should be used when searching electronic databases.

The bulk of the nuclease activity in culture supernatants is owing to two molecularly and kinetically distinct forms, both of which are active as single [polypeptide chains](#) (5). The smaller of these, designated the “slow” (S) form, is produced by **proteolysis** of the larger, “fast” (F) form, which in turn derives from an even larger precursor that has not been characterized (6). *Alteromonas espejiana* copiously produces extracellular [proteinase](#) activity (6), which accounts for the progressive conversion of larger to smaller species in the culture fluid as growth proceeds. At least the S species is not a fully homogenous protein, because its N-terminal amino acid is not unique; either the presumed endoproteolytic event that generates this species from the F form does not occur at a unique site or there is exoproteolytic activity as well. Conversion of the F form to a species indistinguishable in molecular size and catalytic properties from the S form can be done by proteolysis *in vitro* (6).

The American Type Culture Collection strain of this organism (ATCC 29659) is suitable for production of the nucleases. No overproducing strains, either of *Alteromonas* itself or of heterologous hosts containing **cloned** nuclease **genes**, have been reported. Purification procedures that effect separation of the S and F forms have been described (5) and modified (6). BAL 31 nuclease is available commercially from several sources, but these products are mixtures of the S and F species. A partial amino acid sequence of an internal fragment produced by cleavage with cyanogen bromide has been obtained (7). This new sequence did not have significant [homology](#) with any reported sequences. Some physical properties of the S and F forms are presented in Table 1 (5).

**Table 1. Some Physical Properties of F and S Forms of BAL 31 Nuclease**

| Molecular Weight (kDa) | Isoelectric pH | Molar Absorption                  |  |
|------------------------|----------------|-----------------------------------|--|
|                        |                | Coefficient( $10^5$ liter/mol-cm) | Weight Absorption Coefficient(dl/g-cm) |
| F 109.0                | 4.2            | $1.10 \pm 0.09$                   | $10.2 \pm 0.8$                         |
| S 85.0                 | 4.2            | $1.00 \pm 0.03$                   | $11.7 \pm 0.3$                         |

## 2. Reactions Catalyzed

Five activities appear to account for the degradation—ultimately to mononucleotides—of DNA: (1) a 5' → 3' exonuclease activity that acts on single-stranded DNA in a highly processive manner (ie, many nucleotides are removed in a single productive enzyme-substrate encounter) (8); (2) a 3' → 5' directed exonuclease activity that acts only on duplex structures, leaving a 5'-terminated single-stranded “tail” (9); (3) an endonuclease activity against single-stranded DNA, (much slower than the exonuclease activity on single-stranded DNA [8]), that is also elicited by a variety of covalent or noncovalent distortions or lesions in duplex DNA (3-5, 10-14); (4) an exonuclease activity that can excise a small number of nucleotides starting from the site of a strand break (nick) in duplex DNA (15), which is inferred to be 5' → 3'-directed; and (5) an activity that can remove short 3'-terminated tails from otherwise duplex DNA. The F and S forms differ greatly in the rates of removal of nucleotides, at given molar concentrations of DNA ends and enzyme, via the 3' → 5' exonuclease action on double-stranded DNA, which led to their “fast” and “slow” designations (5). The two forms have comparable kinetic behavior toward single-stranded DNA (5, 8).

The BAL 31 nucleases also catalyze the terminally directed hydrolysis of double-stranded RNA and degrade RNA that contains nonduplex structure (16). Whether an analog of the endonuclease activity on altered duplex DNA exists has not been determined, except that cleavage probably occurs in response to a strand break in duplex RNA.

The 5' → 3' exonuclease activity on single-stranded DNA is highly processive, as the nucleases are able to degrade DNA polymers of nearly 500 nucleotides in length completely without dissociation of enzyme from substrate (8). The much larger (5,400-nucleotide) **ϕX174** DNA did not appear to be degraded processively, however, so there may be limits to the number of nucleotides removed in a single productive enzyme-substrate encounter (5). The values of  $V_{\max}$  and the turnover number ( $k_{\text{cat}}$ ), on the basis of the rates of internucleotide bond hydrolysis, are similar for the 500-nucleotide polymer described above and the minimal DNA substrate, a dinucleotide diphosphate (8). The kinetic data appear to require a **facilitated diffusion** mechanism for the activity on macromolecular substrates, in which nuclease molecules bind randomly and diffuse along the contour of the chain until a 5' end is located, at which time catalysis can begin (8). This apparently was the first example of a requirement for facilitated diffusion in an enzymatic reaction.

The 3' → 5' exonuclease activity on duplex DNA and the 5' → 3' exonuclease activity on single-stranded DNA can account for the progressive reduction in the length of duplex DNA. Above certain nuclease concentrations, short 5'-terminated single-stranded tails from the 3' → 5' exonuclease activity are evident in partial digests (9). If 3'-terminated tails are present, they cannot be more than a



few nucleotides in length. The 5' tails have a limiting length, at the higher enzyme concentrations examined, of about 7 nucleotide residues when about 100 nucleotides are removed by the 3' → 5' exonuclease activity. Up to 50% of the ends could be joined in [DNA Ligase](#) reactions under conditions favoring the joining of fully base-paired ends. This indicates that the 5' → 3' exonuclease activity on single-stranded DNA can terminate at the junction between single- and double-stranded regions, to leave fully base-paired ends on a significant fraction of molecules in a partially digested population. **DNA polymerase**-mediated repair, to render ends with 5'-terminated tails fully base-paired, markedly increased the fraction of ends joinable by DNA ligase, as expected (9).

As the nuclease concentrations were decreased below those mentioned above, the average length of the single-stranded tails per 100 nucleotides removed increased dramatically, and the ligase-mediated joining of ends in the absence of polymerase-mediated repair became undetectable. This apparent dependence on the concentration of the enzyme of the relative velocities of two exonuclease reactions that it catalyzes has not been explained. The velocity of the overall reaction to shorten DNA, as measured by the hyperchromicity associated with the hydrolysis of internucleotide bonds, is proportional to enzyme concentration, as expected (17). Repair of ends generated at the low nuclease concentrations noted above with DNA polymerase rendered a high percentage of the molecules ligatable (9).

There must be an activity that can remove short 3'-terminated protruding single-stranded tails from duplex DNA, because the nuclease readily degrades such starting substrates (9). No evidence exists, however, for 3' terminally directed hydrolysis of single-stranded DNAs (8).

The value of the [Km\(Michaelis constant\)](#), in terms of the molar concentration of DNA ends, for the overall exonucleolytic shortening of duplex DNA by the S form is so large that it is not practical for the substrate concentration to approach it in actual reactions. Consequently, the enzyme velocity is directly proportional to substrate concentration over the accessible range (18). A value for this  $K_m$  was reported in earlier work (5), but it appears to have been determined under conditions for which the substrate concentration did not exceed that of the enzyme, violating a premise of Michaelis-Menton-based kinetics. Kinetic parameters ( $V_{max}$  per unit concentration of nuclease and  $K_m$ ) for the length reduction of duplex DNA by the F form can be measured accurately, because the  $K_m$  value is much lower than for the S form and lies in the range of experimentally accessible concentrations of DNA ends. More recent values (18) are in reasonable agreement with the earlier data (5).

These kinetic parameters for the F nuclease depend on the length of the DNA substrate (7) (not determinable for the S form as noted). This has been shown (7) to be consistent with a model in which nonspecific binding of the nuclease away from the ends is followed by a "search" process to form a productive enzyme-substrate complex, with the enzyme bound to a terminus, as for the degradation of single-stranded DNA. However, such a mechanism is not required by the kinetics as in the case of single-stranded DNA (noted in the text above).

The kinetics of length reduction of duplex DNA are also dependent on its guanine + cytosine(G + C) content (7, 18). This is significant in the case of the S form of the nuclease, for which the rate of nucleotide removal from DNA ends decreases over fourfold over the range of 37–66 mole% G + C residues. This dependence is somewhat less for the F enzyme. Data are available to predict these effects of DNA base composition on the kinetics (18).

In contrast to the high processivity for the 5' → 3' exonuclease action on single-stranded DNA noted earlier, the 3' → 5' exonuclease is quasi-processive, removing only about 18 and 28 nucleotide residues per productive binding event for the S and F enzymes, respectively (9).

The length reduction of duplex RNA (16) presumably proceeds by a similar mechanism but has not been characterized. The BAL 31 nucleases are the most efficacious enzymes known for the

controlled length reduction of duplex DNA, and they are apparently the only enzymes that can catalyze this reaction for duplex RNA.

Duplex DNA is not attacked endonucleolytically (away from an end) at a significant rate unless there is some alteration of the duplex structure. Hence, nonsupercoiled, closed circular duplex DNA (form I° DNA) is extremely resistant to attack by the BAL 31 nucleases (11, 13, 14). The very limited attack on form I° DNA at high enzyme concentrations and long incubations (14), plus the kinetic parameters for exonucleolytic degradation of linear duplex DNA, lead to estimates of the relative rates of introduction of endonucleolytic breaks to exonucleolytic scissions of  $8 \times 10^{-11}$  and  $7 \times 10^{-12}$  for the S and F enzymes, respectively.

Negative supercoiling in closed circular DNA can elicit the endonuclease activity (12). In the majority of molecules of such supercoiled DNA, an endonucleolytic event (cleavage in one strand) is followed by the removal, in a processive manner, of several nucleotides (6.5 and 2.8 nucleotides for the F and S forms, respectively) from the initially nicked strand, to yield a gapped circular DNA intermediate (15). A fraction contain only a strand break, and no nucleotides are excised. The percentage of molecules with no nucleotides removed is significantly lower for the F nuclease than for the S species. The nicks and gaps are bounded by 5'-phosphoryl and 3'-hydroxyl termini. It is assumed that other alterations that result in endonucleolytic attack give rise to such gapped circular intermediates.

The removal of a small number of nucleotides could be accomplished exonucleolytically, starting from the site of the initial nick. Or, there could be a second endonucleolytic cut a few nucleotides away from the first one; the short oligonucleotide between endonucleolytic breaks would dissociate at room temperature to leave a gap. Available evidence supports the exonucleolytic mechanism. It was reasoned that the presumed exonucleolytic activity producing the gaps should be 3' → 5'-directed, as the exonucleolytic activity attacking base-paired ends has this characteristic. However, this proved to be inconsistent with the results of further experiments carried out assuming the 3' → 5' mode of attack; the 5' → 3' attack is thus inferred. Because this activity operates on nominally duplex DNA, it might be expected that it could remove nucleotides in a 5' → 3' direction from fully base-paired ends, to leave 3'-terminated single-stranded ends. It was noted that, if these are present, they could not be more than a few nucleotides in length, but this does not rule out very limited activity comparable to that producing the gaps.

The nicked and gapped circular DNA intermediates are then converted to linear duplex DNA by a second endonucleolytic event in the other strand, which requires a second encounter with a nuclease molecule (15). The resulting linear duplex DNA is further degraded by the mechanism noted in the text above. As expected, circular duplex DNA containing a nick introduced by means other than BAL 31 nuclease action is converted to linear duplex DNA, with the occurrence of gaps of the average sizes noted above in most of the circular molecules before their linearization (15).

### 3. Effects of Solvent Variables and Denaturing Agents

Both  $\text{Ca}^{2+}$  and  $\text{Mg}^{2+}$  at concentrations above 10–12 mM are required for maximum activity on both single- and double-stranded substrates (9). All activities are optimal near neutral pH (4, 9). The nucleases are remarkably resistant to inactivation at elevated salt concentrations (4, 10). The reaction buffer in which most of the work described above was done contained 0.6 M NaCl). This salt tolerance is not surprising, considering that sea water is the natural milieu of these extracellular enzymes.

High concentrations of agents that normally **denature** proteins fail to eliminate the nuclease activities. Single-stranded DNA is degraded at 40% of the maximum rate in 6.5 M **urea** (4), and substantial activity on this substrate was found in the presence of 6 M **guanidinium** chloride. Crude

preparations maintained activity against both single- and double-stranded substrates in the presence of 5% (w/v) SDS (3). At least a portion of the structure of the S nuclease does not become disrupted under the stringent denaturation conditions used for denaturing SDS-PAGE (6).

The nucleases are not remarkably resistant to thermal inactivation, as the half-life for disappearance of the activity on single-stranded DNA is only 3–5 minutes at 50°C (7). However, tests on preparations stored at 4°C imply that most of the activity should be retained for years in a buffer containing 5 mM Mg<sup>2+</sup> and Ca<sup>2+</sup> (4). Commercial preparations are often supplied in 50% (v/v) glycerol for storage at –20°C, under which conditions they should retain full activity indefinitely.

#### 4. Applications

By far the most extensive use of the BAL 31 nucleases, represented by hundreds of literature citations, has been the controlled length reduction of linear duplex DNA. The bulk of these reports describe the production of deletion mutants and/or the elucidation of sequences required for a particular biological activity of a (usually) cloned DNA segment. Cloned sequences can be deleted unidirectionally by cleavage of the vector containing the cloned insert at a unique restriction site on one side of the insert (see [Restriction Enzymes](#)), carrying out a partial BAL 31 degradation, releasing the shortened insert from the shortened vector by use of a restriction site on the other side of the insert, and ligation to an intact linearized vector DNA (19).

Numerous reports have appeared of BAL 31 nuclease-mediated identification of **telomeric** sequences, which occur at the ends of linear eukaryotic chromosomes and hence are degraded first by the exonuclease in intact DNA; Yao and Yao (20) and De Lange and Borst (21) apparently represent the earliest work. Where these DNAs are too large to isolate as intact molecules in solution, *in situ* lysis in [agarose](#) has been used (22) so that the initial substrate DNA is largely intact.

The progressive removal of sequences from duplex ends allows the determination of the [restriction map](#) of a DNA, either naturally linear or linearized at a unique site, by noting the order in which fragments from subsequent digestion with the restriction enzyme in question disappear from [gel electrophoresis](#) patterns of progressively shortened aliquots (19, 23). This technique is greatly enhanced if unidirectional deletions can be done, as noted above (19), as ambiguity arising from loss of sequences from both ends of the fragment is eliminated. The sites of bound proteins, including [nucleosomes](#), and interstrand crosslinks that block the exonuclease can be elucidated because intact duplex sequences between such sites will not be attacked (24, 25).

The endonuclease activity that cleaves in response to alterations in duplex structure has been used in several laboratories to detect such alterations. This attack will be followed by cleavage of the other strand and exonucleolytic attack from the ends thus generated, as noted in the text above. Where chemical modification is done, form I° DNA is used (this requires that the modification does not introduce nicks), and the rate of its loss on incubation with nuclease is monitored. Some of the chemical lesions that give rise to endonucleolytic attack are pyrimidine dimers and possibly other photoproducts of ultraviolet irradiation, adducts with [carcinogens](#) such as *N*-acetoxy-*N*-2-acetylaminofluorene and *N*-methyl-*N*-nitrosourea, interstrand cross-links produced by reaction with nitrous acid, apurinic sites, and adducts with Hg<sup>2+</sup> and Ag<sup>+</sup> ions (11, 13, 14). Strand breaks were noted in the text above as eliciting cleavage of the opposite strand. Noncovalent alterations eliciting endonucleolytic attack include moderate degrees of negative supercoiling and very high degrees of positive supercoiling (12), junctions between right-handed B-DNA and left-handed Z-DNA regions (10), the presence of unpaired nucleotides in one strand of an otherwise duplex DNA (26), and [cruciform](#) structures resulting from the extrusion of inverted repeated sequences under supercoiling stress (27, 28). **Matrix attachment sites** for **nucleolar** DNA are apparently sensitive to the nuclease (29). Finally, the lack of sensitivity of form I° DNAs has been used to help identify such species in DNA populations (30).

## Bibliography

1. R. T. Espejo and E. S. Canelo (1968) *J. Bacteriol.* **95**, 1887–1891.
2. K. Y. Chan, L. Baumann, M. M. Garza, and P. Baumann (1978) *J. Syst. Bacteriol.* **28**, 217–222.
3. H. B. Gray, Jr., D. A. Ostrander, J. L. Hodnett, R. J. Legerski, and D. L. Robberson (1975) *Nucleic Acids Res.* **2**, 1459–1492.
4. H. B. Gray, Jr., T. P. Winston, J. L. Hodnett, R. J. Legerski, D. W. Nees, C.-F. Wei, and D. L. Robberson (1981) In *Gene Amplification and Analysis* (J. G. Chirikjian and T. S. Papas, eds.), Elsevier North-Holland, New York, pp. 169–203.
5. C.-F. Wei, G. A. Alianell, G. H. Bencen, and H. B. Gray, Jr. (1983) *J. Biol. Chem.* **258**, 13506–13512.
6. C. R. Hauser and H. B. Gray, Jr. (1990) *Arch. Biochem. Biophys.* **276**, 451–459.
7. T. Lu (1992) *Some physical and catalytic properties of BAL 31 nuclease*, *Ph.D. dissertation*, University of Houston, Houston, Texas.
8. T. Lu and H. B. Gray, Jr. (1995) *Biochim. Biophys. Acta* **1251**, 125–138.
9. X.-G. Zhou and H. B. Gray, Jr. (1990) *Biochim. Biophys. Acta* **1049**, 83–91.
10. M. W. Kilpatrick, C.-F. Wei, H. B. Gray, Jr., and R. D. Wells (1983) *Nucleic Acids Res.* **11**, 3811–3822.
11. R. J. Legerski, H. B. Gray, Jr., and D. L. Robberson (1977) *J. Biol. Chem.* **252**, 8740–8746.
12. P. P. Lau and H. B. Gray, Jr. (1979) *Nucleic Acids Res.* **6**, 331–357.
13. C.-F. Wei, G. A. Alianell, H. B. Gray, Jr., R. J. Legerski, and D. L. Robberson (1983) In *DNA Repair: A Laboratory Manual of Research Procedures* (E. C. Friedberg and P. C. Hanawalt, eds.), vol. **2**, pp. 13–40.
14. C.-F. Wei, R. J. Legerski, G. A. Alianell, D. L. Robberson, and H. B. Gray, Jr. (1984) *Biochim. Biophys. Acta* **782**, 404–414.
15. A. Przykorska, C. R. Hauser, and H. B. Gray, Jr. (1988) *Biochim. Biophys. Acta* **949**, 16–26.
16. G. H. Bencen, C.-F. Wei, D. L. Robberson, and H. B. Gray, Jr. (1984) *J. Biol. Chem.* **259**, 13584–13589.
17. X.-G. Zhou (1989) *Ethidium bromide-mediated renaturation of denatured closed circular DNAs in alkaline solution: mechanistic aspects and fractionation of closed circular DNAs on a molecular weight basis. Some catalytic properties and mechanism of exonuclease action of BAL 31 nuclease*, *Ph.D. dissertation*, University of Houston, Houston, Texas.
18. H. B. Gray, Jr. and T. Lu (1993) In *Enzymes of Molecular Biology* (M. M. Burrell, ed.), Humana Press, Totowa, New Jersey, pp. 231–251.
19. C. R. Hauser and H. B. Gray, Jr. (1991) *Gene Anal. Tech. Appl.* **8**, 139–147.
20. M.-C. Yao and C. H. Yao (1981) *Proc. Natl. Acad. Sci U.S.A.* **78**, 7436–7439.
21. T. De Lange and P. Borst (1982) *Nature* **299**, 451–453.
22. R. F. Wintle, T. G. Nygaard, J. A. Herbrick, K. Kvaloy, and D. W. Cox (1997) *Genomics* **40**, 409–414.
23. R. J. Legerski, J. L. Hodnett, and H. B. Gray, Jr. (1978) *Nucleic Acids Res.* **5**, 1445–1464.
24. W. A. Scott, C. F. Walter, and B. L. Cryer (1984) *Mol. Cell. Biol.* **4**, 604–610.
25. W.-P. Zhen, O. Buchardt, H. Nielsen, and P. E. Nielsen (1986) *Biochemistry* **25**, 6598–6603.
26. D. H. Evans and A. R. Morgan (1982) *J. Mol. Biol.* **160**, 117–122.
27. L. H. Naylor, H. A. Yee, and J. H. van de Sande (1988) *J. Biomol. Struct. Dyn.* **5**, 895–912.
28. N. M. Morales, S. D. Coburn, and U. R. Muller (1990) *Nucleic Acids Res.* **18**, 2777–2782.
29. O. V. Iarovaia, M. A. Lagarkova, and S. V. Razin (1995) *Biochemistry* **34**, 4133–4138.
30. E. Cuzzoni, L. Ferretti, C. Giordani, S. Castiglione, and F. Sala (1990) *Mol. Gen. Genet.* **222**, 58–64.

## Suggestions for Further Reading

31. H. B. Gray, Jr., T. P. Winston, J. L. Hodnett, R. J. Legerski, D. W. Nees, C.-F. Wei, and D. L. Robberson (1981) "The extracellular nuclease from *Alteromonas espejiana*: an enzyme highly specific for nonduplex structure in nominally duplex DNA", In *Gene Amplification and Analysis* (J. G. Chirikjian and T. S. Papas, eds.), Elsevier North-Holland, New York, pp. 169–203.
32. H. B. Gray, Jr. and T. Lu (1993) "The BAL 31 nucleases (EC 3.1.11)", In *Enzymes of Molecular Biology* (M. M. Burrell, ed.), Humana Press, Totowa, New Jersey, pp. 231–251.

## Balbiani Ring

The [polytene chromosomes](#) of insects have extraordinary utility in the investigation of **chromosomal** structure and function. Polytene chromosomes display the phenomenon of puffing in which the [chromatin](#) associated with chromosomal bands is decondensed. In the larval salivary glands of the midge *Chironomus*, several puffs become very large and are described as *Balbiani rings*. Balbiani is the cytologist who discovered polytene chromosomes in 1881. The visible aspect of a puff is essentially the complex of **RNA** and protein that accumulates as a result of vigorous **transcriptional** activity. Among the gene products encoded by the Balbiani rings are the glue proteins required to attach the midge pupa to its substrate. All of the Balbiani ring genes are expressed coordinately and maximally during the second larval instar. Then they show differential expression in prepupae ([1](#)). The most detailed ultrastructural analysis has been performed on Balbiani ring 2 in chromosome IV of *Chironomus tentans*. Although the polytene chromosome band itself contains more than 470 kbp of DNA, the major gene is 37 kbp long and contains two very well-defined transcription units that are differentially expressed at the prepupal stage. [Electron microscopy](#) detects the first engaged **RNA polymerase** molecule at the approximate start site of transcription and the last RNA polymerase at the site of termination ([2](#), [3](#)). Within each puff, the Balbiani ring gene forms a loop of transcriptionally active chromatin with attached [ribonucleoprotein](#) complexes. As synthesis proceeds, the nascent transcripts fold up into compact ribonucleoprotein complexes that head toward the **nuclear pores** when they are exported to the cytoplasm. Each 37-kb pre-mRNA transcript is processed in the nucleus by staged association with specific proteins ([4](#)).

Contour measurements of chromosomal structure in the puff indicate that the chromatin is fully extended into an array of [nucleosomes](#) 10 nm in diameter at the site of transcription, whereas once transcription is completed, chromatin coils back up into a 30-nm diameter fiber, which is then finally packaged into a supercoiled loop. Upstream of the start site of transcription is a region free of nucleosomes, presumably corresponding to the **promoter**, and compacted chromatin fibers are further upstream and downstream of the gene. Immunologic analysis allows defining the structural components of chromatin on active and inactive segments of the Balbiani ring. Surprisingly, proteins, such as [histone](#) H1 and the core histones, remain on chromatin even while it is actively transcribed ([5](#)). This keen observation allows a rather complete ultrastructural picture of the transcription process.

## Bibliography

1. U. Lendahl and L. Weislander (1987) *Develop. Biol.* **121**, 130–138.
2. B. Bjorkroth, C. Ericsson, M. M. Lamb, and B. Daneholt (1988) *Chromosoma* **96**, 333–340.
3. C. Ericsson et al. (1989) *Cell* **56**, 631–639.
4. H. Mehlin, B. Daneholt, and U. Skoglund (1992) *Cell* **69**, 605–613.
5. C. Ericsson, U. Grossbach, B. Bjorkroth, and B. Daneholt (1990) *Cell* **60**, 73–83.

### Suggestion for Further Reading

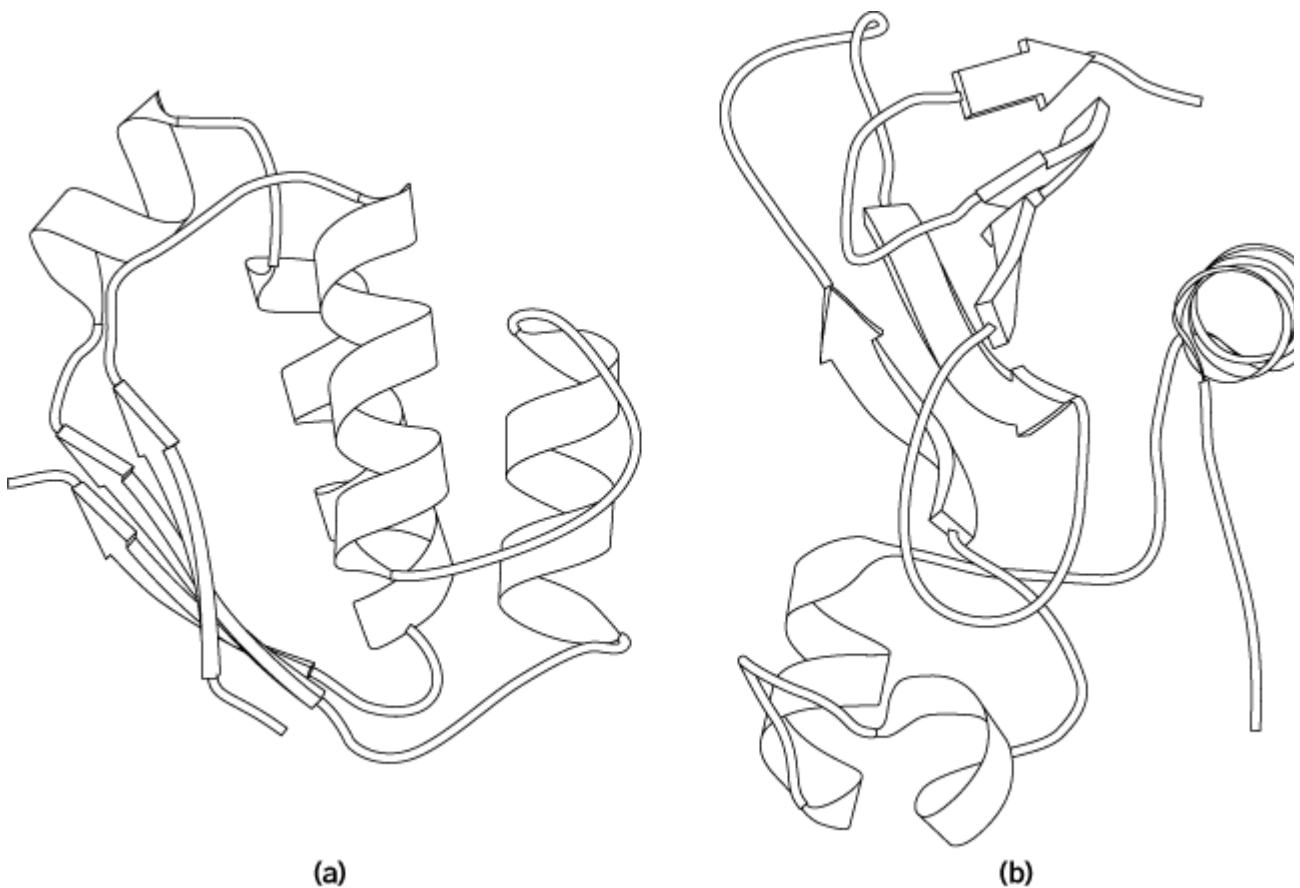
6. B. Daneholt (1997) A look at messenger RNP moving through the nuclear pore. *Cell* **88**, 585–588.

### Barnase and Barstar

Barnase is a **ribonuclease** secreted by the bacterium *Bacillus amyloliquefaciens*. It is a small protein of 110 [amino acid](#) residues with no [disulfide bonds](#) or nonpeptide components (see [Protein Structure](#)). As such, it was recognized in the early 1960s to be an ideal subject for the then-emerging study of **protein folding**. Its specific inhibitor, barstar, is produced intracellularly by the same organism and is an equally simple and even smaller protein of only 89 residues. Together they form a two-subunit complex in which the [active site](#) of barnase is buried, providing a model system for the study of [protein–protein interactions](#). The system has contributed substantially to our understanding of how proteins fold and interact and has potential for further uses.

The three-dimensional structures of both proteins and their complex have been determined (Fig. [1](#)) ([1](#)). The genes for both proteins have been **cloned** into **plasmid vectors** in *Escherichia coli* and the expressed products obtained in high yield. Expression of barnase alone is lethal, so simultaneous expression of barstar is necessary for production of barnase.

**Figure 1.** Molscript ([1](#)) representation of the structures (a) barstar and (b) barnase in their complex.



The catalytic mechanism of barnase and its relatives is essentially the same as that of the pancreatic ribonuclease family (see **Ribonuclease A**), with a histidine and a glutamic acid residue acting respectively as proton donor and acceptor, in place of the two histidines of the latter. The barnase [active site](#) lies in a broad groove on the side of the [b-sheet](#) opposite the major [a-helix](#) (see Fig. 1). In the complex, which has a **dissociation constant** on the order of  $10^{-14}$  M, barstar completely blocks the active site, with most of its contacts involving its second helix and an adjacent loop. An aspartic acid carbonyl group of barstar occupies the position of the attacked phosphate group of substrate RNA. The rate at which barnase and barstar associate is controlled largely by charged groups, positive on barnase, negative on barstar. Mutation to alanine of any of these charged residues reduces binding but increases the stability of each protein, indicating the evolutionary balance between function and stability.

Both barnase and barstar may be unfolded reversibly by heat or by **denaturants** such as [urea](#) or [guanidinium chloride](#) (see **Protein unfolding**). Both fold and unfold in a highly **cooperative** two-state manner; that is, only the native and unfolded states are significantly occupied at **equilibrium** under most conditions. During the kinetics of folding, however, barnase passes through an observable intermediate state. In this it exemplifies the behavior of most larger proteins or protein **domains** (see **Protein folding**). Barstar, on the other hand, appears to collapse more directly to the folded state, as do most other smaller proteins. Studies of barnase folding, using [kinetic](#) and equilibrium techniques and [site-directed mutagenesis](#), have traced the order in which its parts come together to form the native fold and have estimated the free energies of various interactions in several states that occur during folding.

Barnase is a member of a large family of microbial ribonucleases that share the same basic [tertiary structure](#) but have quite divergent sequences; this provides a wider system for approaching the central problem of how the fold is determined by the sequence (see [Protein Structure Prediction](#)).

The homologous ribonucleases of the *Streptomyces* share less than 25% sequence identity with barnase, but the structure of the active-site region is so well conserved that these enzymes are also inhibited by barstar, with dissociation constants as low as  $2 \times 10^{-12} M$ ; consequently, coexpression of barstar permits the high level expression of the cloned genes for these enzymes as well. As the *Streptomyces* also produce barstar homologues, it seems clear that the extreme structural conservatism of the active site of the ribonucleases is based on the strict requirement of strong binding to the inhibitor.

Separate from the continuing application of barnase and barstar as small model proteins to general problems of protein structure and chemistry are more recent uses of their genes that take advantage of the toxic effect of barnase expression in tissues of heterologous organisms. Under control of a specific **promoter**, expression of the barnase gene can ablate cells under conditions or in tissues where the promoter is turned on. This property is being applied, for example, in studies of [development](#), strategies against **viruses**, design of a conditionally lethal selective **cloning vector**, and, most spectacularly, in the development of male-sterile **plants**. In the last instance, male sterility can be reversed by inclusion of the barstar gene.

There have been several efforts at **computer simulation** of aspects of barnase folding and unfolding, and more can be expected. The continued application of directed mutagenesis, physical chemistry, and computer modeling to the folding mechanisms of barnase and barstar and their homologues will provide insight into the manner in which the sequences of these two interdependent families of proteins determine their native folds. *In vivo* use of the barnase gene, with its lethal effect limited to specific conditions or tissues, will be useful to developmental biologists and should have many practical applications in agriculture and possibly in medicine as well.

#### Bibliography

1. P. Kraulis (1991) *J. Appl. Crystallog.* **24**, 946–950.

#### Suggestions for Further Reading

2. R. W. Hartley (1997) "Barnase and Barstar", in *Ribonucleases, Structure and Functions* (G. D'Alessio and J. F. Riordan, eds.), Academic Press, New York, pp. 51–100. (Contains an exhaustive bibliography through 1995.)
3. A. Matouschek, L. Serrano, and A. R. Fersht (1994) "Analysis of protein folding by protein engineering", in *Mechanisms of Protein Folding* (R. H. Pain, ed.), IRL Press, Oxford, pp. 137–159.

#### Barr Body

In eutherian (placental) mammals, **dosage compensation** mechanisms operate in female cells to silence one of the two **X-chromosomes**, so that female cells have as many X-chromosome-derived transcripts as male cells containing a single X-chromosome. This **X-chromosome inactivation** process was first proposed by Mary Lyon (1961) and is known as the **Lyon hypothesis** ([1](#)). Female mammalian embryos begin development with two active X-chromosomes, but very early in embryogenesis almost all of the genes on one of the two X-chromosomes become inactivated ([Random X-Inactivation](#)). Although the initial choice between inactivation of the maternal or paternal X-chromosome is random, once established in a repressed state the same X-chromosome are inactivated after every cell division. The inactivation process occurs over the entire chromosome,



and practically all of the genes on the chromosome are silenced. This transcriptional inactivation is concomitant with the chromosome taking on the appearance of **heterochromatin** and also becoming late replicating during **S phase** (see [Facultative Heterochromatin](#)). Only a small fraction of the inactive X-chromosome, including the genes located in the pseudoautosomal region at Xp22.3, escapes the global silencing process (2). The inactive chromosome remains in the nucleus and can be detected cytologically as a *Barr body* in somatic cells (3). The Barr body is found only in cells containing more than one X-chromosome, if cells are trisomic for the X-chromosome, two Barr bodies will be detected. The staining procedures used to detect Barr bodies argue for a global difference in chromatin condensation (4). Differences in staining of heterochromatin compared to transcriptionally competent **euchromatin** may be caused by differences in compaction or differences in the association of many more accessory proteins in heterochromatin (5).

The use of advanced microscopy and molecular **cytogenetics**, in which fluorescent *in situ* **hybridization** is used to “paint” chromosomal territories, has recently allowed detailed dissection of chromosomal organization (see [Denaturation Mapping](#)). Light microscopic optical serial sectioning of the active and inactive X-chromosome territories reveal that they occupy similar volumes. However, reconstructed active X-chromosomes have a flatter shape and a more extended, folded surface area than the inactive X-chromosome (6). The conclusion is that the differential staining properties of the Barr body are caused mainly by the association of a distinct group of accessory proteins and RNA with this inactive chromosome (see [X-Chromosome Inactivation](#)).

Even on the active X-chromosome, most of the chromatin is not transcriptionally active. Thus, on both the inactive and the active X-chromosomes, the great majority of the chromatin is maintained in a folded state typical of a transcriptionally repressed state. There are, however, global differences in chromatin organization between active and inactive X-chromosomes. The inactive X chromosome has high overall nucleosomal density—i.e., it carries a higher “concentration” of all histones, including the unusual variant macroH2A, than all autosomes or the active X (7). This likely contributes to the optical density of the Barr Body as visualized microscopically. The Barr body contains methylated DNA, hypoacetylated histones, a specialized structural RNA (Xist), and is late replicating during the S phase. Association with Xist RNA is an important causal factor in X-inactivation, and somehow leads to the selective hypermethylation of the inactive X, which results in deacetylation of chromatin assembled over it and transcriptional quiescence. The three-dimensional distribution in the nucleus of Xist RNA coincides with the chromosomal territory occupied by the inactive X-chromosome (8). Even after removing bulk chromatin in the preparation of a nuclear matrix (which is responsible for the overall morphology of the nucleus), the Xist RNA remains in the matrix. This is consistent with the hypothesis that Xist RNA has a structural role in establishing the inactive X-chromosome territory. The Barr body still has a great deal to teach scientists about the molecular mechanisms that establish and maintain nuclear compartments.

#### Bibliography

1. M.F. Lyon, *Nature* **190**, 372–373 (1961).
2. C.M. Disteche, *Trends Genet.* **11**, 17–22 (1995).
3. M.L. Barr and E.G. Bertram, *Nature* **163**, 676–677 (1949).
4. S.M. Gartler and A.D. Riggs, *Annu. Rev. Genet.* **17**, 155–190 (1983).
5. S.W. Brown, *Science* **151**, 417–425 (1966).
6. R. Eils et al., *J. Cell Biol.* **135**, 1427–1440 (1996).
7. Perche P.Y. et al., *Curr. Biol.* **10**, 1531–1534 (2000).
8. C.M. Clemson, J.A. McNeil, H.F. Willard, and J.B. Lawrence *J. Cell Biol.* **132**, 259–275 (1996).

#### Additional Reading

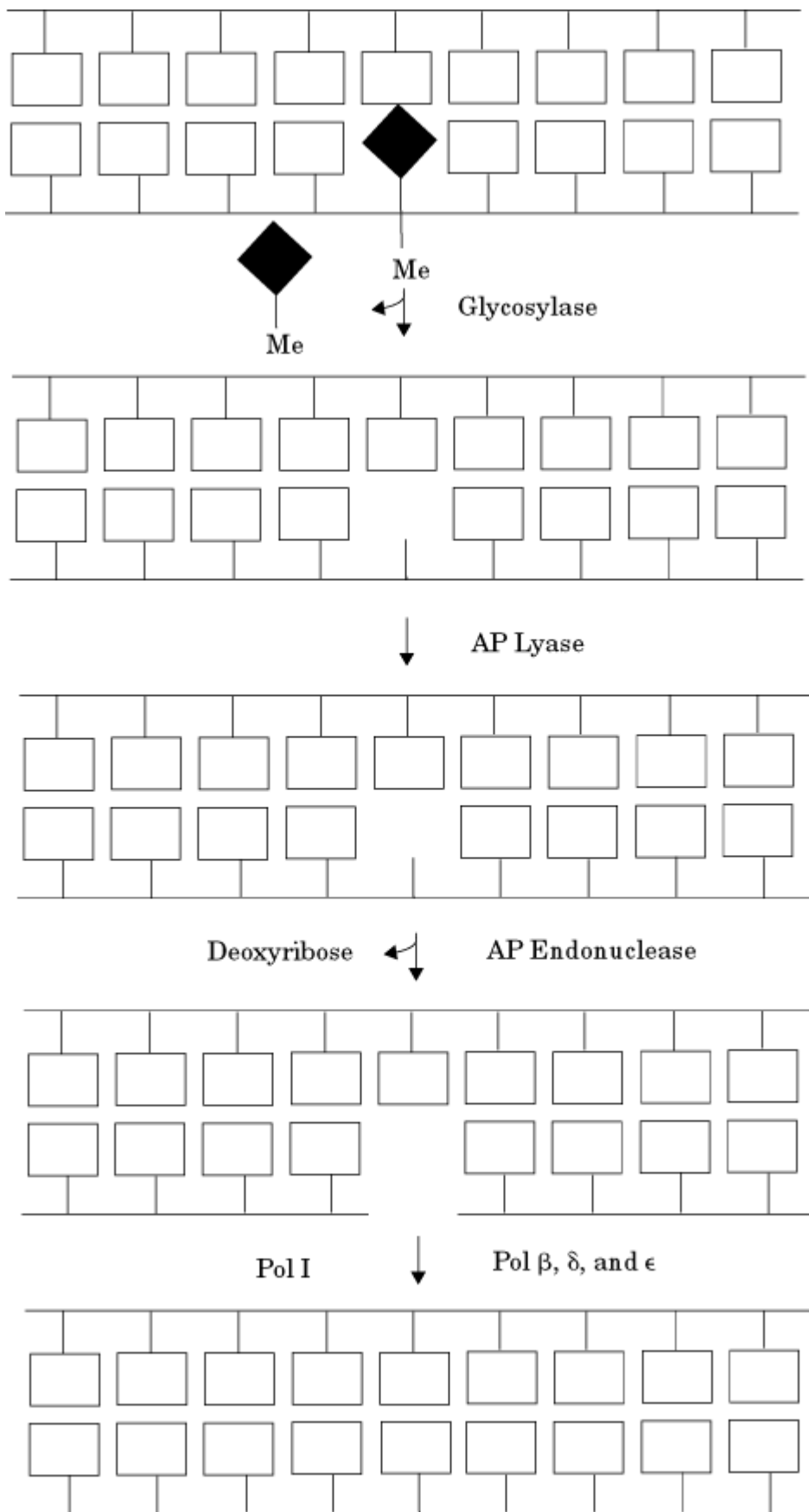
9. Cremer T. et al., Role of chromosome territories in the functional compartmentalization of the

- cell nucleus, Cold Spring Harbor Symp. Quant. Biol. **58**, 777–792 (1993).
10. Bird A.P. and Wolffe A.P., Cell **99**, 451–454 (1999).
  11. Panning B. and Jaenisch R., Cell **93**, 305–308 (1998).
  12. Csankovszky G., Nagy A., and Jaenisch R., J. Cell Biol. **153**, 773–783 (2001).

## Base Excision Repair

Base excision repair is the [DNA repair](#) system for abnormal bases in DNA, such as uracil, hydroxymethyluracil, and hypoxanthine, and for nonbulky base modifications, such as thymine glycol and hydrates, 8-hydroxyguanine, *O*<sup>4</sup>-methylthymine, and 3-methyladenine (Fig. 1). The abnormal bases and simple base lesions are removed from DNA by the combined actions of [DNA glycosylases](#) (1-3) and AP (apurinic/aprimidinic) endonucleases. The basic reaction carried out by glycosylases is the cleavage of the glycosylic bond joining the base to deoxyribose. Some of the glycosylases cleave both the glycosylic bond and the phosphodiester bond 3' to the resulting AP site in a more or less concerted reaction. Hence, DNA glycosylases have been classified either as “pure” glycosylases or as glycosylase/AP lyases. AP endonucleases are also classified as either 3'- or 5'-AP endonucleases. All 3'-AP endonucleases are lyases that cleave the phosphodiester bond by β-elimination. The 5'-AP endonucleases carry out true hydrolysis.

**Figure 1.** Repair of 3-methyl adenine by base excision repair. Methyladenine DNA glycosylase cleaves the glycosylic bond of the alkylated base; this releases the alkylated base and leaves an AP site. An AP lyase and an AP endonuclease cleave 3' and 5', respectively, of the AP site, releasing the abasic sugar. The resulting gap is filled in by Pol I in *E. coli* and by Pol β or Pol δ and ε in humans, to produce repair patches of 1 to 10 nucleotides.



The combined actions of glycosylase, AP lyase, and AP endonuclease/exonuclease creates a one- to two-nucleotide gap in the damaged strand. This gap is filled in by **DNA polymerase I** in *Escherichia*

*coli*, and mostly by DNA polymerase  $\beta$  (4) and to a lesser degree by DNA polymerases  $\delta$  and  $\epsilon$  in yeast and humans (5). The resulting nick is ligated by [DNA Ligase](#) in *E. coli*, and in eukaryotes by the ligase I– (6) or ligase III–XRCC1 complex (7). Depending on the polymerase/ligase combination involved in filling in and sealing the gap, a repair patch of 1 to 10 nucleotides may be produced.

### Bibliography

1. T. Lindahl (1974) Proc. Natl. Acad. Sci. USA **71**, 3649–3653.
2. T. Lindahl (1976) Nature **259**, 64–66.
3. J. Laval (1977) Nature **269**, 828–832.
4. R. W. Sobol, J. K. Horton, R. Kohn, H. Gu, R. K. Singhal, R. Prasad, K. Rajewsky, and S. H. Wilson (1996) Nature **379**, 183–186.
5. Z. Wang, X. Wu, and E. C. Friedberg (1993) Mol. Cell. Biol. **13**, 1051–1058.
6. R. Prasad, R. K. Singhal, D. K. Srivastava, J. T. Molina, A. E. Tomkinson, and S. H. Wilson (1996) J. Biol. Chem. **271**, 16000–16007.
7. K. W. Caldecott, C. K. McKeown, J. D. Tucker, S. Ljungquist, and L. H. Thompson (1994) Mol. Cell. Biol. **14**, 68–76.

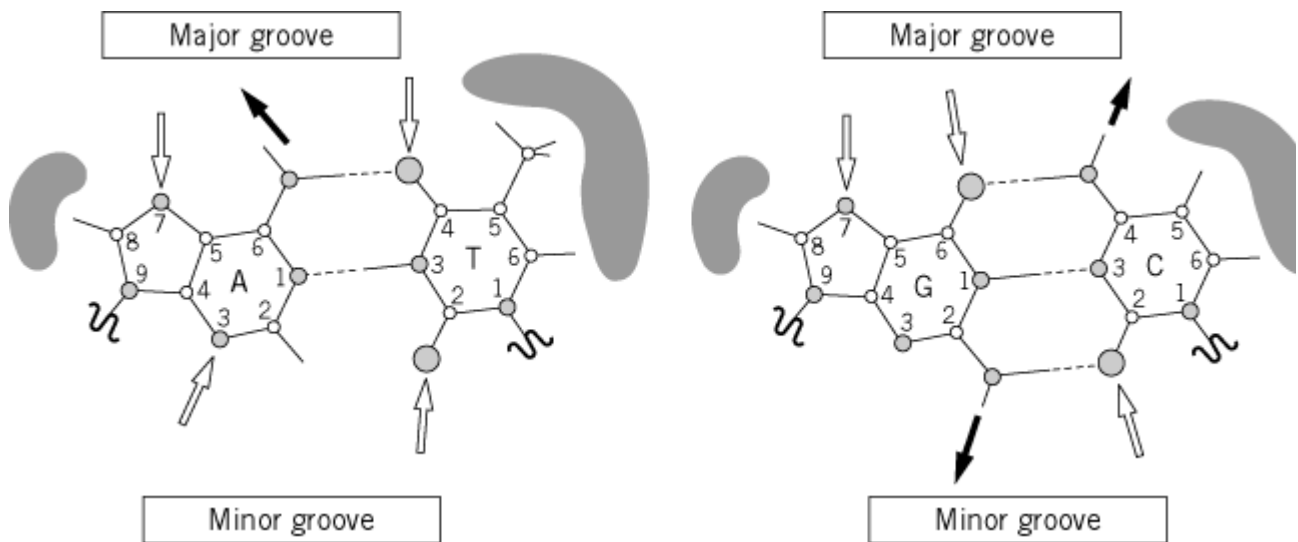
### Suggestions for Further Reading

8. B. Demple and L. Harrison (1994) Repair of oxidative damage to DNA: enzymology and biology. Annu. Rev. Biochem. **63**, 915–948.
9. E. Seeberg, L. Eide, and M. Bjoras (1995) The base excision repair pathway. Trends Biochem. Sci. **20**, 391–397.

### Base Pairs

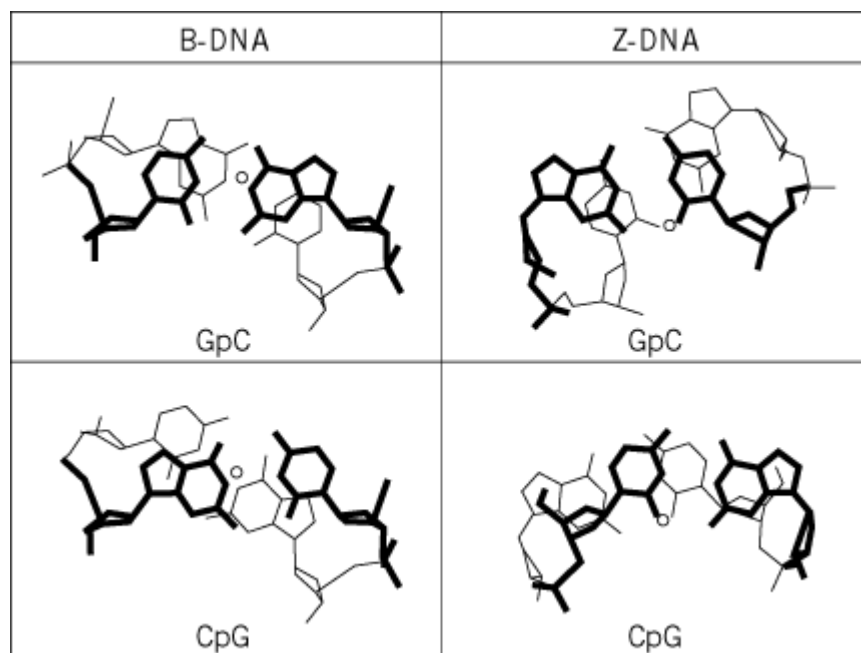
DNA exists predominantly as an antiparallel double-stranded helix. The two strands are helically coiled, which maximizes the exposure of the negatively charged sugar–phosphate backbone to water and shields the **hydrophobic** aromatic bases in the middle from [water](#). Figure 1 shows the Watson–Crick base pairs in which the specificity of base pairing is provided by [hydrogen bonds](#). The complementary arrangements of hydrogen bond donors and acceptors allow the Watson–Crick base pairing of guanine with cytosine (G:C) and adenine with thymine (A:T). The two C1' atoms within a G:C base pair and an A:T base pair are equidistant ( $\sim 10.5$  Å). It should be noted that other types of base pairs are also known to exist, and they often are involved in unusual DNA structures.

**Figure 1.** Schematic illustration of the A:T and G:C Watson–Crick base pairs. Hydrogen bonds are shown as dashed line. Hydrogen bond donors (gray arrow) and acceptors (white arrow) of the bases are shown. Note that the bases can pair in more than one way. Hydrophobic edges of the base pairs are emphasized by gray shade. The major and minor groove edges of the base pair are shown.



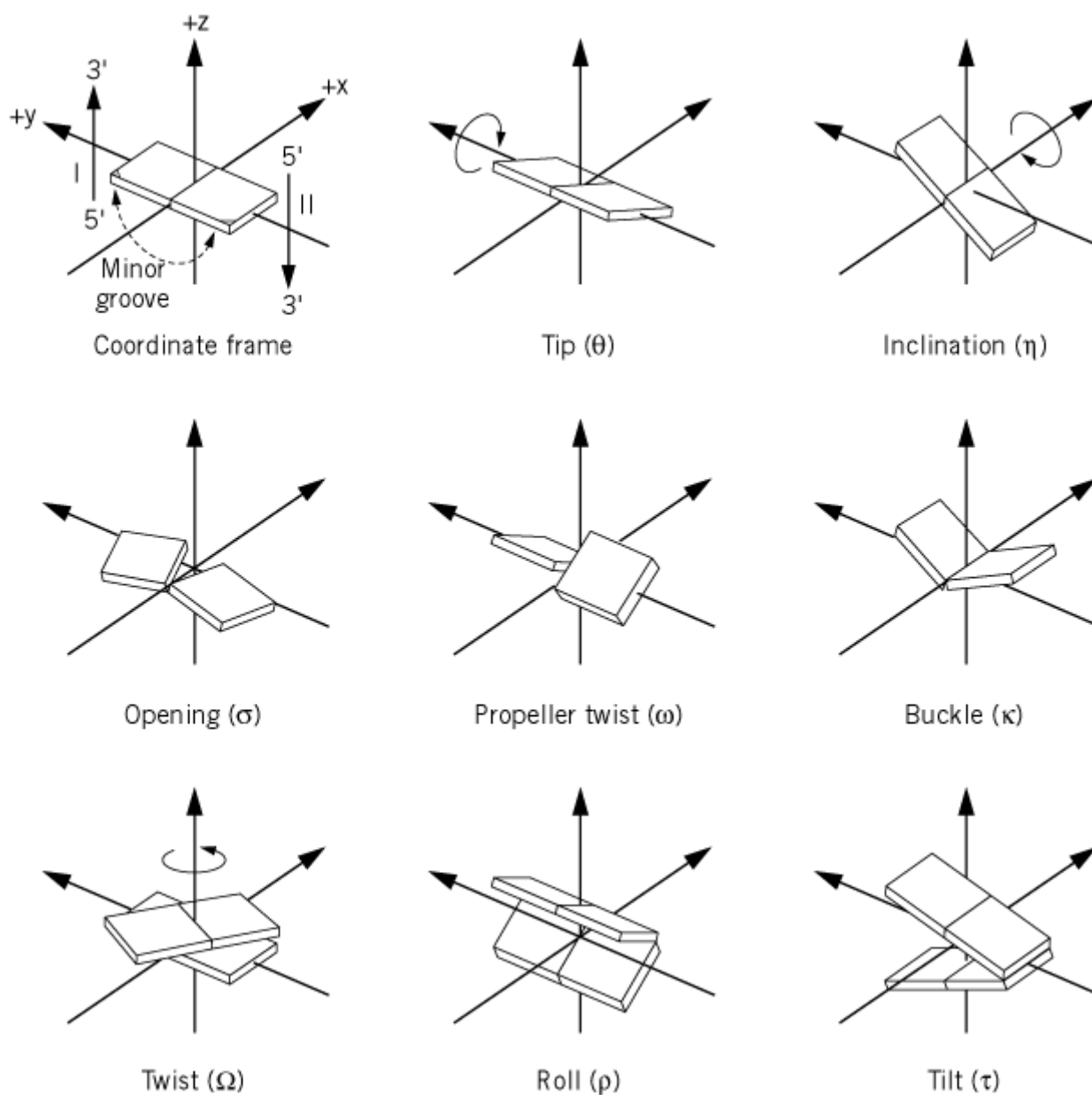
Bases in a DNA double helix are stacked on top of each other. The stacking interaction is stabilized by the electronic  $p-p$  interaction of the aromatic rings, and it also minimizes the exposure of the bases to the solvent. The base stacking has preferred interactions, depending on the characteristic dipole-dipole interactions between adjacent base pairs. Thus, the stacking patterns vary with the helix type and the particular base-base step. In B DNA, the 5'-purine-*p*-pyrimidine-3' (Pu-Py) base-pair step has a good stacking overlap, and is thus more stable, whereas the 5'-pyrimidine-*p*-purine-3' (Py-Pu) base-pair step has a poor stacking overlap, and is thus less stable (Figure 2). Such differences in stacking energy have strong implications in the bendability of DNA associated with protein binding. It has been noted that DNA sequences with Py-Pu steps in B-DNA are deformed more easily; therefore steps such as CpG, TpA, or TpG (=CpA) are often associated with kinked DNA structures.

**Figure 2.** Stacking interactions in B-DNA (left) and Z-DNA (right) showing Pu-Py step stacking (top) and Py-Pu step stacking (bottom). The circle shows the position of the helix axis. Note that the Py-Pu steps (including CpG, TpA, and CpA (=TpG)) in B-DNA have very small stacking interactions, causing those steps to be easily deformed.



Through the high resolution X-ray structures of DNA oligonucleotides, it has been found that base pairs within a DNA duplex (or other structures) have significant departures from a strict coplanar geometry. The conformational parameters associated with the DNA base pairs are defined in Figure 3. For a more complete definition, the reader is referred to the Cambridge Convention (1).

**Figure 3.** Definitions of various base pair conformations. (upper two rows) or two successive base pairs (bottom row). In the top row the motions of the bases are coordinated, while in the middle row their motions are opposed. The left, center, and right columns describe rotations about the  $z$ ,  $y$ , and  $x$  axes, respectively. The standard coordinate frame is defined at the upper left.



#### Bibliography

1. R. E. Dickerson et al. (1989) EMBO J. **8**, 1–4.

## Base-Pair Substitution

These [mutations](#) of DNA affect only one base pair at a time, possibly producing either amino acid substitutions or chain termination codons in a gene (see [Genetic Code](#)). Amber (UAG), ochre (UAA), and UGA codons will all terminate polypeptide synthesis (see [Stop Codons](#)). Amino acid substitutions often produce leaky and/or **temperature-sensitive** mutations, whereas chain-terminating mutants are usually nonleaky. It has been suggested that there is an important difference between the consequences of frameshift and base-pair substitution mutations in diploid organisms. Frameshifts are often **recessive**, whereas base-pair substitution mutations may exert partial **dominance** (1).

Freese (2) divided base-pair substitutions into two classes, *transitions* and *transversions*. Transitions are mutations in which the purine–pyrimidine orientation is conserved; that is, a purine is replaced by another purine, a pyrimidine by another pyrimidine. The eight possible transversion mutations in double-stranded DNA are: A.T ↔ T.A; T.A ↔ G.C; G.C ↔ C.G; C.G ↔ A.T. Transition mutations have the purine–pyrimidine orientation conserved, as in A.T ↔ G.C or C.G ↔ T.A.

Transitional mutations have been traditionally considered to be due to **tautomerism** of nucleotides. All of the four common DNA bases can exist in tautomeric forms, of which the biologically most interesting involve keto–enol pairs for thymine and guanine and amino–imino pairs for cytosine and adenine. The tautomer, although obeying the usual purine–pyrimidine pairing rules, is capable of hydrogen bonding with a base that is normally noncomplementary. If this occurs and is not repaired, it will lead directly to base-pair substitution mutagenesis in subsequent cell generations. Some base analogues are highly mutagenic, presumably through direct incorporation into DNA through mispairing with a normal base in a tautomeric form (see [2-Aminopurine \(AP\)](#), [5-Bromouracil](#)). Some DNA alkylating chemicals may be mutagenic through adding an alkyl group onto a normal nucleotide base, thereby affecting its base-pairing properties in subsequent rounds of replication (see [Dimethyl Sulfate \(DMS\)](#)). Such premutagenic events will be fixed at a low level during subsequent rounds of replication, typically leading to transition mutations of various types.

### Bibliography

1. C. Wills (1968) Proc. Natl. Acad. Sci. USA **61**, 937–944.
2. E. Freese (1959) Brookhaven Symp. Biol. **12**, 63–73.

## Bence-Jones Proteins

By the middle of the nineteenth century, a British physician made a curious observation on the urine of a patient suffering from a bone disease (in fact, a multiple myeloma). Perhaps because his name was Dr. Watson, this practitioner was curious. So, once back home he heated the urine and observed

the formation of a precipitate at 60°C, which redissolved upon boiling. Because he was also clever, he realized that this was not normal, but could not provide any explanation. He thus sent a sample to a member of the faculty in London, Dr. Henry Bence-Jones, who not only confirmed the fact but also gave the answer. The mysterious substance was tagged as “hydrated albumin deutoxide.” Dr. Watson was certainly pleased with the answer, but was forgotten. Dr. Bence-Jones became immortal, although the nature of the Bence-Jones protein (BJP) remained a puzzle until it was solved by Edelman and Gally (1), who showed, more than a century later, that BJPs were the free form of the monoclonal [immunoglobulin](#) light chain that is associated with a heavy chain in the corresponding myeloma protein.

BJPs have been of great interest for quite a while because they are present in large amounts in urine of most patients with multiple myeloma, thus providing a simple source of monoclonal materials, perfectly suitable for structural analysis of light chains, including amino acid sequencing. The first two sequences were obtained by Hilschmann and Craig in (2) in 1965 for two BJPs of the k [isotype](#), giving the first evidence that immunoglobulin chains contained **variable** and **constant** regions. In normal physiological situations, the synthesis rates of the heavy and light chains are similar, whereas they are most often imbalanced in malignant conditions. In most cases of multiple myeloma, an excess of light chain is produced and secreted, quite often as **disulfide-bonded** dimers, causing an associated kidney disease. The presence of BJP is still a valid sign for the diagnosis and prognosis of multiple myeloma.

In the mouse, an equivalent disease may be induced, at a frequency dependent upon the strain, with a high response in BALB/c, as shown by Potter (3). In the mouse too, BJPs are present in the urine of animals that develop a plasmacytoma. Because the tumor may be transplanted for generations and produce unchanged myeloma proteins and BJP, this provided most of the materials that allowed the immunologists to solve the basic features of [immunoglobulin structure](#) and to open the way to [monoclonal antibodies](#), as well as the key to Ig genes and thus to the elucidation of the genetic basis for antibody diversity.

See also entries [Antibody](#), [Immunoglobulin](#), and [1\) Light Chains](#).

### Bibliography

1. G. M. Edelman and J. A. Gally (1962) The nature of Bence-Jones proteins. Chemical similarities to chains of myeloma globulins and normal -globulins. *J. Exp. Med.* **116**, 207–227.
2. N. Hilschmann and L. Craig (1965) Amino acid sequence studies with Bence-Jones proteins. *Proc. Natl. Acad. Sci. USA* **53**, 1403–1409.
3. J. L. Fahey and M. Potter (1959). Bence-Jones proteinemia associated with a transplantable mouse plasma cell neoplasm. *Nature* **184**, 654–655.

### Suggestion for Further Reading

4. R. A. Kyle (1994) The monoclonal gammopathies. *Clin. Chem.* **40**, 2154–2161.

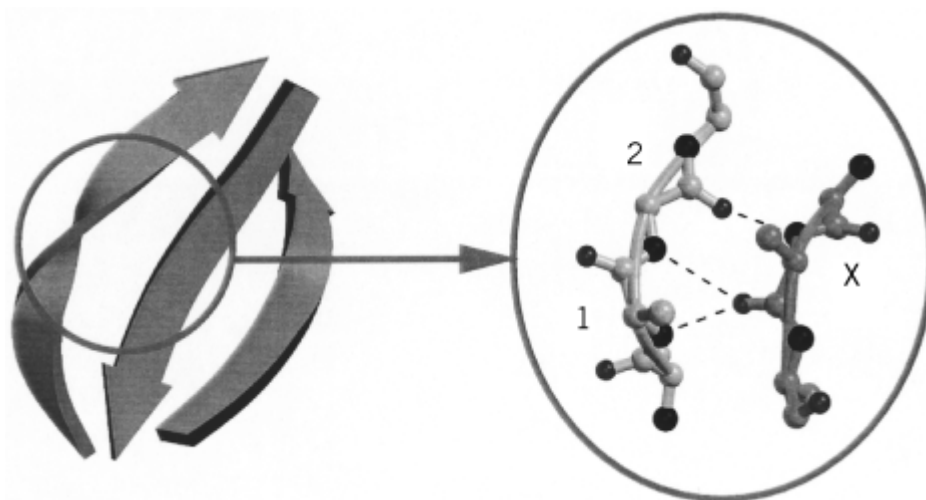
### Beta-Bulge

The b-bulge is a term used to define a specific **secondary structure** feature of [protein structures](#), where the regular [hydrogen bond](#) interactions and [backbone](#) conformation of the [b-sheet](#) are disrupted by the presence of an extra residue. The additional residue is usually located on a [b-strand](#)



at one edge of the b-sheet, where the bulge is more easily accommodated within the structure of the protein (Fig. 1). b-Bulges occur frequently, and often there are two or more per protein (1). In some cases, b-bulges have been found to be conserved in the structures of related proteins, where they may play a functional role. The more general role of the b-bulge is to alter the direction of the b-strand, so that it may be considered a type of [turn](#). However, the change in direction of the polypeptide chain is not as pronounced as that caused by other types of turns.

**Figure 1.** Example of a classic b-bulge in an antiparallel b-sheet. Three b-strands of the antiparallel b-sheet are shown. The normal twist of the b-sheet is clearly accentuated by the b-bulge, caused by an additional residue in the middle b-strand at the edge of the b-sheet. The close-up view of the b-bulge (**right**) shows residue X from the b-strand forming backbone hydrogen bonds (shown as dashed lines) with two residues (1 and 2) from the outer b-strand. For clarity, side-chain atoms have been removed, backbone nitrogen atoms and backbone oxygen atoms are shown as dark spheres. This figure was generated using Molscript (3) and Raster3D (4, 5).



A b-bulge is that region between two consecutive b-type hydrogen bonds that includes two residues (called positions 1 and 2) on one b-strand opposite a single residue (position X) on the other b-strand (2). There are several different classes of b-bulge, most (90%) occurring between antiparallel b-strands. The two most common are the classic and the G1, accounting for about 80% of b-bulges. In the classic b-bulge, the residue at position 1 has backbone dihedral angles ( $\phi = -100$ ,  $\psi = -25$ ) (see [Ramachandran Plot](#)), closer to those of an  $\alpha$ -helical than to a b-strand conformation, but residue 2 has angles closer to the b-strand conformation ( $\phi = -180$  and  $\psi = 160$ ). In the G1 b-bulge, the residue at position 1 has a positive  $\phi$  value ( $\phi = 85$  and  $\psi = 0$ ) and is therefore almost always glycine (thus the name G1). The residue at position 2 of the G1 class has dihedral angles corresponding to b-strand ( $\psi = -90$  and  $\psi = 150$ ). The G1 b-bulge often occurs in combination with a type II b-turn (which requires a glycine at position 3). Compared to the usual b-sheet structure, a b-bulge disrupts the alternating side chain placement on one of the b-strands and increases the right-handed twist of the b-strand from the usual  $10^\circ$  to  $35^\circ$  to  $45^\circ$ .

[See also [Beta-Sheet](#).]

#### Bibliography

1. A. W. E. Chan, E. G. Hutchinson, D. Harris, and J. M. Thornton (1993) *Protein Sci.* **2**, 1575–1590.
2. J. S. Richardson, E. D. Getzoff, and D. C. Richardson (1978) *Proc. Natl. Acad. Sci. USA* **75**, 2574–2578.
3. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.

4. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
5. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

## Beta-Galactosidase of *Escherichia Coli*

Exploration of the lactose (*lac*) system of *Escherichia coli* was started in the late 40s at the Pasteur Institute in Paris by Jacques Monod and his collaborators. As a result, most of the major concepts of **gene expression** and regulation were established with this system (see [Lac Operon](#)). The major gene product of the *lac* operon is b-galactosidase, the enzyme that splits lactose into glucose and galactose and the product of the *lacZ* gene. In addition to its importance for studies of enzyme induction, this protein has become very useful as a [reporter gene](#) in very many studies of **gene expression**.

### 1. Characterization

Because of its importance in genetic studies of *lac* operon, *E. coli* b-galactosidase was purified and characterized by several groups in the early 50s ([9](#), [10](#)). It could be easily measured quantitatively because of a [chromogenic substrate](#), *o*-nitrophenyl-b-D-galactoside (ONPG). The enzyme is a hydrolytic transglucosidase and accounts for up to 5% of the total protein in haploid, fully induced or constitutive strains of *E. coli*. The purified enzyme was crystallized, although its three-dimensional structure had to wait 40 years to be solved (see later). Several procedures for producing large amounts of enzyme have been described (reviewed in Ref. [6](#)). b-galactosidase is a particularly stable enzyme. It can be stored in 40% ammonium sulfate for several years at 4 °C without significant loss of activity. In buffered solution, the enzyme is stable at 37 °C for several weeks. At 57 °C the half-life of heat inactivation is 10 min ([11](#)). Active enzyme is recovered with 100% yield after **denaturation** with 8 M [urea](#), and fairly good recovery is possible after treatment by 6 M **guanidinium chloride** ([12](#)).

b-Galactosidase is specific for the b-D-galactopiranoside configuration and cleaves the bond between the anomeric carbon (C1) and the glycosyl oxygen ([10](#)). The [enzyme](#) acts as a hydrolase, but also as a transferase. The galactosyl moiety can be transferred to monosaccharides, oligosaccharides, alkyl alcohols, and phenols and also to **mercaptoethanol**. b-Galactosidase also transfers the galactosyl moiety from the C4 position to the C6 position of glucose to form *allolactose*, which is the natural **inducer** of the *lac* operon ([13](#)). b-Galactosidase requires Mg<sup>2+</sup> and Na<sup>+</sup> ions for maximal activity ([9](#)), but the reason for needing these cations is still unclear.

That b-galactosidase is a tetramer of four identical subunits, possessing one [active site](#) per subunit, is well documented, and the tetramer is the only active form of the enzyme. Because of its high turnover number (about 6000 moles of ONPG hydrolyzed/sec/mole of enzyme, one molecule of b-galactosidase per bacterial cell can be accurately measured. Several residues (Glu 461, Met 502, Tyr 503, Glu 537) have been identified as important for catalytic function or to be near the active site ([14](#), [15](#)). These residues are found in the three-dimensional structure in close proximity to one another and form a pocket likely to be the substrate-binding site ([16](#)).

The amino acid sequence of b-galactosidase, determined before DNA **sequencing** techniques had been developed, was the longest protein sequence ever determined ([17](#)). According to this sequence, the protein contains 1021 amino acid residues and a subunit molecular weight of 116,248. Subsequently, the nucleotide sequence of the *lacZ* gene was also determined ([18](#)), from which it was predicted that b-galactosidase consists of 1023 residues, with a molecular weight of 116,353 per

subunit. It is quite remarkable that the amino acid sequence deduced from the DNA sequence differs only in ten amino acid residues from that obtained by amino acid sequencing of such a large protein.

There is also a **homologue** of b-galactosidase, evolved b-galactosidase (EBG). Normally *E. coli* does not grow on lactose when the *lacZ* gene is deleted. However, cells with *lacZ* deletions were selected that grow on lactose and produce a new b-galactosidase, called EBG (19). The *ebgA* gene maps at 66 min on the *E. coli* chromosome, whereas that for *lacZ* maps at 8 min. The active EBG protein is a hexamer containing 964 residues per monomer. The sequence identity between the *ebgA* and *lacZ* genes at the DNA nucleotide level (50%) exceeds that at the amino acid level (30%). This led to the conclusion that the two genes descended by **divergence** from a common ancestral gene. b-galactosidase and EBG do not cross-react with [antisera](#) prepared against each type of protein.

## 2. b-Galactosidase Complementation

The term [complementation](#) has been generally used in genetics for the phenomenon by which a biological function lost or altered by a mutation is restored through mutual compensation by two differently altered mutant genes. In the case of *intracistronic complementation*, where the two mutations occur within different copies of the same **cistron**, it is generally accepted that the repair of function occurs at the level of the protein product of the gene and involves the interaction of differently altered polypeptide chains (see [Interallelic Complementation](#)). Such complementation involves noncovalent interaction of polypeptide chains and may occur by reassociation of differently altered subunits of an oligomeric protein or by interaction between fragments of a single polypeptide chain. Studies of b-galactosidase complementation contributed to elucidating the structure of the enzyme and to understanding the mechanism of the specific recognition, reassociation, and folding of proteins (for reviews, see Refs. 20, 21).

### 2.1. Complementation Between Point Mutants

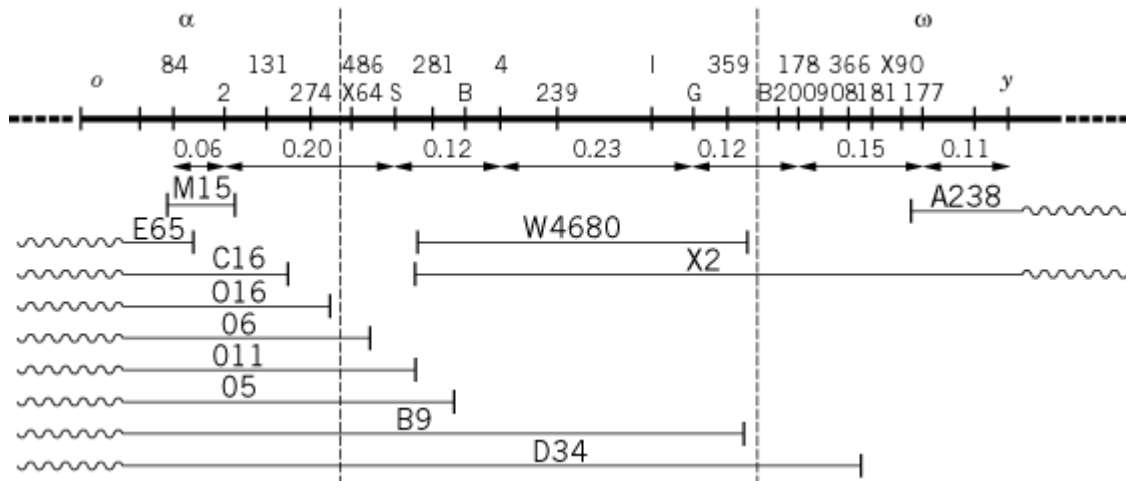
Jacob and Monod (8) showed that heterogenotes carrying different *lacZ*<sup>-</sup> point mutants may become *lacZ*<sup>+</sup> by complementation. *In vitro* studies carried out by Perrin (22) showed that many of the inactive proteins containing these point mutations exhibit monomeric structures, but cross-react immunologically with b-galactosidase. In contrast, the active, complemented, b-galactosidase had a sedimentation coefficient like tetrameric wild-type 16 S, although it was more heat-labile. The restoration of enzyme activity in this type of complementation can be accounted for by a mechanism suggested by Crick and Orgel (23), repair of lesions by reassociation of differently altered subunits.

Some inactive b-galactosidase proteins with point mutations exist as dimers. A specific class of these mutants can be activated up to 1000 times by specific antibodies raised against wild-type b-galactosidase (24, 25). The activating antibody directs the formation of active tetramers. In a strict sense this is not a complementation reaction, but the kinetics of activation are similar to those of complementation between point mutant proteins *in vitro*: In both cases, the appearance of b-galactosidase activity is a relatively slow process, requiring about 300 min at 37 °C.

### 2.2. Complementation Between Deletion Mutants

The isolation of deletion mutants of the *lacZ* gene involving quite large genetic segments (26) led to the discovery of a new type of intracistronic complementation (for a review, see Ref. 27). Inactive *lacZ* deletion mutants, lacking part of the gene corresponding to either the amino- or the carboxy-terminal region of the b-galactosidase polypeptide chain, would complement another inactive deletion mutant containing the region missing in the first. The N- and C-terminal regions involved in the two distinct classes of complementation are known as the a and w regions, respectively (Fig. 1). Both a- and w-complementation represent noncovalent reassociation of complementary fragments of the b-galactosidase subunit polypeptide chain, which then can reassemble into an enzymatically active tetrameric structure. a- and w-complementation takes place *in vivo* in partial diploids and *in vitro*, with extracts of relevant mutant strains.

**Figure 1.** Diagrammatic representation of the *lacZ* gene. The figures and letters above the solid line indicate the position of point mutations in the *lacZ* gene: *o*, operator; *y*, structural gene of Lac permease. The figures below the line represent recombination frequencies, as determined by crosses between two point mutants. The lines below indicate the extent of various deletions: M 15 isolated by Beckwith (28); A 238 and W 4680, isolated by Cook and Lederberg (29); all others by Jacob, Ullmann and Monod (26). Reprinted from (20) by permission of Cold Spring Harbor Laboratory.



### 2.3. a-Complementation

All mutants that have their promoter-proximal segment (a-segment) intact are a-donors. All genes that contain partial deletions of the N-terminus not extending beyond a barrier (which is located by the dashed line in Figure 1 between mutants 274 and X64) are a-acceptors when the proteins they produce are mixed with that of an a-donor (30). The time course of *in vitro* a-complementation is relatively slow and reaches a plateau in about 3 h. The a-complemented enzyme, like wild-type b-galactosidase, is a tetramer but is less stable to heat or urea treatment. The equilibrium constant for the complementation reaction was estimated to be approximately  $1 \text{ to } 2 \times 10^9 \text{ M}^{-1}$  (21). Given the slow association rate constant, one can predict a very low dissociation rate, indicating that the complementation process is virtually irreversible under physiological conditions.

The most unusual property of the a-peptide is its high temperature resistance in 6 M guanidinium chloride. Bacterial extracts of *lacZ* strains containing a-donors liberate soluble a-peptides after boiling in this denaturant (30) or after autoclaving in its absence (31). The smallest a-peptide still capable of complementation has a size of about 7400 Da. Upon cleavage of b-galactosidase with **cyanogen bromide**, Langley et al. (32) isolated and sequenced a peptide that has a-donor activity corresponding to residues 3 to 92 of the polypeptide chain.

The best studied a-acceptor is the M15 protein, produced by the *lacZ* mutant DM15 isolated by Beckwith (28) and shown genetically to be a small deletion in the early part of the gene (Fig. 1). The M15 protein lacks residues 11 through 41 of b-galactosidase and is a dimer (33). It has a trace level of enzyme activity and can be significantly activated by addition of anti-b-galactosidase antibodies (34), up to 15% of that obtained by a-complementation.

The nature of the a-complementation process and its probable pathway have been studied extensively (21). It is believed that the a-peptide modifies the [quaternary structure](#) of the a-acceptor by causing its tetramerization and consequently the recovery of enzymatic activity. The three-dimensional structure of b-galactosidase (16) fully confirms this prediction by showing that the a-peptide segment of the polypeptide chain participates in forming specific subunit contacts (see later).

### 2.4. w-Complementation

All mutants that have their promoter-distal region (w-region) intact are w-donors. They complement w-acceptors that are products of deletion or point mutants in the w-region (27) (Fig. 1), which represents about one-third of the total genetic length of *lacZ*. The purified w-peptide has a molecular weight of 40,000 Da (35). The kinetics of w-complementation are faster than those of a-complementation and reach a maximum extent in about 1 h. Although w-peptide is reversibly **renatured** from 8 M urea or 6 M guanidinium chloride, it is highly heat-labile.

w-Acceptor proteins active in complementation have not been obtained in pure form. In crude bacterial extracts (36, 37), w-acceptor exists in different forms of aggregation, depending on the ionic strength. At low ionic strength it is probably a tetramer that dissociates at high ionic strength, most probably into monomers. Both forms are active in complementation. The w-complemented enzyme is a tetramer, like the wild type, and exhibits very similar catalytic properties. The main difference appears after denaturation in 8 M urea. Recovery of its enzyme activity after [renaturation](#) does not exceed 5%, whereas it is practically 100% with wild-type enzyme. Moreover, the w-peptide is recovered quantitatively from the complemented enzyme after urea treatment (11).

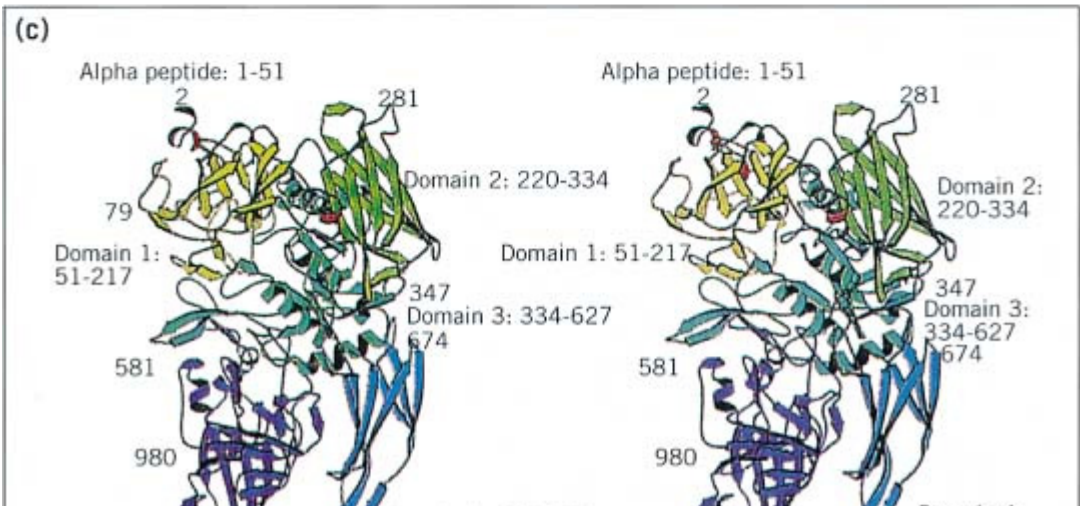
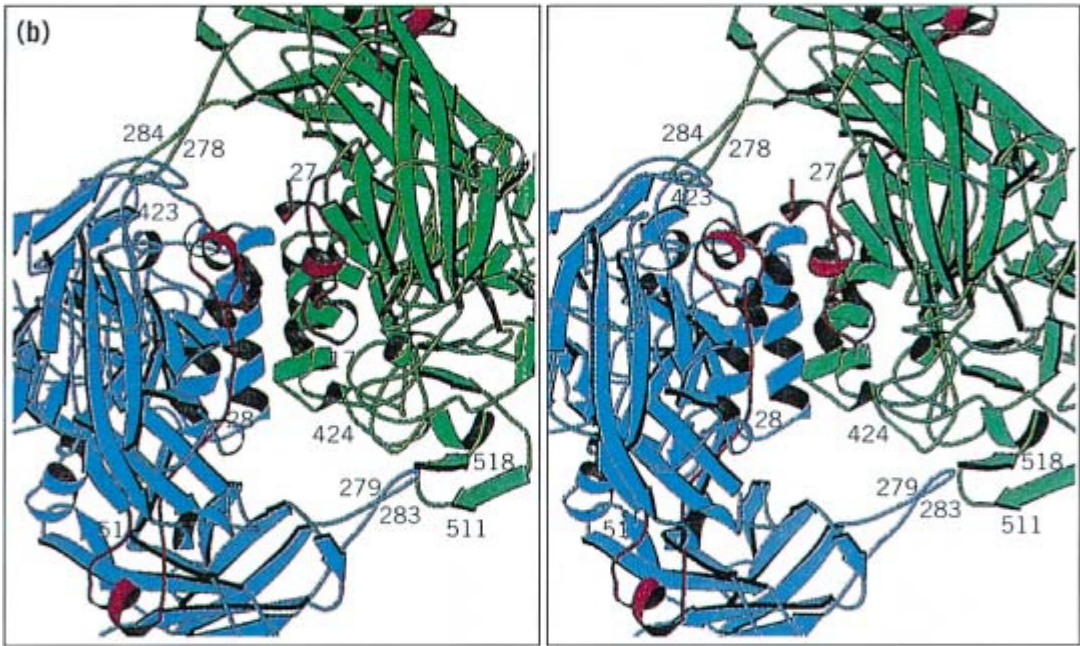
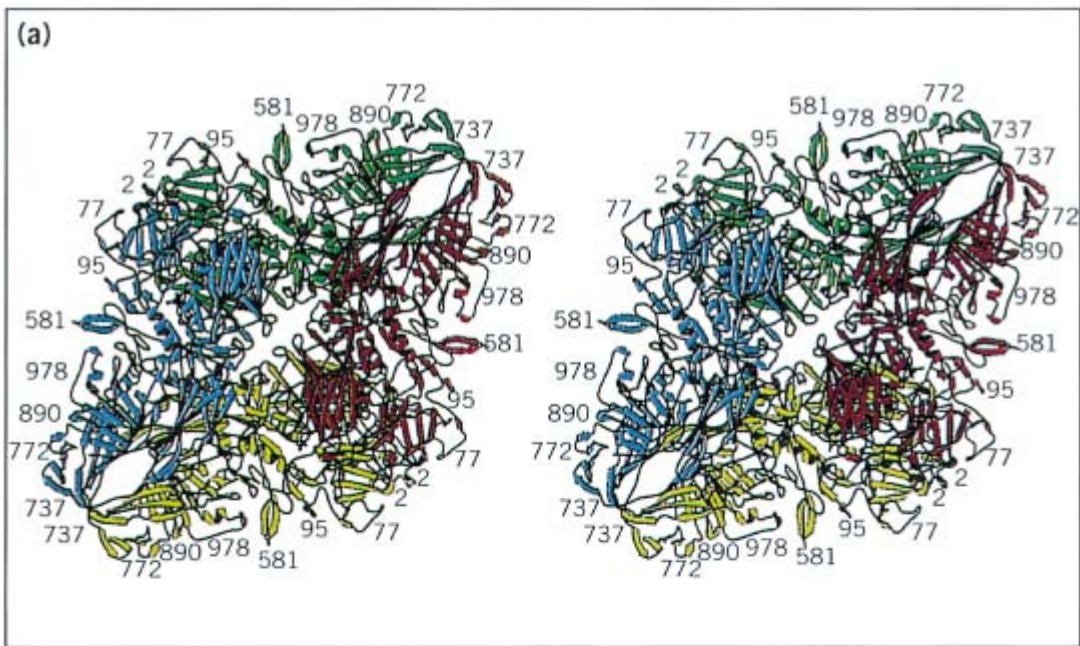
One of the most remarkable features of w-complementation is the high affinity and specificity that the peptide fragments exhibit. They reassociate even at high dilution and in the presence of a large excess of foreign proteins. The yield of complementation is considerably increased in the presence of a b-galactosidase substrate analog or in the presence of anti-b-galactosidase antibodies (36, 37). Comparing wild-type and w-complemented b-galactosidase, Goldberg (38) proposed a model for the tertiary structure of the two monomers. After having shown that the active, complemented enzyme involves reassociation into a tetrameric structure of four w-acceptor polypeptides with four w-peptides, he showed that both w and acceptor are folded to form compact globular structures, close to those of the corresponding segments of the complete polypeptide within the native enzyme. This “globule model” involves the notion that large polypeptide chains contain several distinct **domains** that assemble specifically to give the polypeptide chain its final tertiary structure. The well-defined compact w-domain is apparent within the native monomer in the three-dimensional structure of b-galactosidase (16), which also shows an unstructured 10-residue exposed segment that joins the w-domain to the rest of the chain. This explains, *a posteriori*, why mild proteolytic treatment of the wild-type b-galactosidase liberates significant amounts of w-peptide (11).

### 3. Three-Dimensional Structure

The structure of tetrameric b-galactosidase was solved by [X-ray crystallography](#) at 2.5 Å resolution by B.W. Matthews et al. (16). It is the longest polypeptide chain for which an atomic structure has been determined. The crystal structure shows that b-galactosidase is a tetramer with dimensions of roughly  $175 \times 135 \times 90$  Å along the respective twofold axes. The constituent monomers interact through two different monomer–monomer contacts. The b-galactosidase monomer consists of five compact domains and a relatively extended 50-residue N-terminal segment, corresponding to the a-peptide, which contributes to the activating interface (Fig. 2). The monomer is relatively flat and elongated and has approximate dimensions of  $50 \times 100 \times 40$  Å. The first domain consists of residues 52 to 217, the second of residues 220 to 334, the third of residues 334 to 627, the fourth of residues 627 to 736, and the final fifth domain, corresponding to the w peptide, of residues 737 to 1023. Topologically, the second domain is identical to the fourth domain. The modular structure of the monomer corroborates earlier findings (37, 38), which suggested that the folding of this rather long polypeptide chain is facilitated by each domain acting as an independent unit.

**Figure 2.** Three-dimensional structure of b-galactosidase. (a) Ribbon representation of the b-galactosidase tetramer showing the largest face of the molecule. Contacts between red/green and blue/yellow dimers form the long interface. Contacts between the red/yellow and blue/green dimers form the activating interface. Formation of the tetrameric particle results in two deep clefts that run across opposite faces of the molecule. Each contains two active sites. (b) Ribbon diagram of the blue/green dimer viewed down the molecular two-fold axis, showing the composition of the activating interface. Residues 1 to 50 from each chain, which form the a-complementation region (see text), are shown

in red. The interface includes contacts between the respective complementation peptides, between two [a-helices](#) from the respective monomers that pack together to form a four-helix bundle and between an extended loop (residues 272 to 288) from each monomer that reaches across the interface and extends into the active-site region of the neighboring monomer, stabilizing the active site structure. (c) Stereo ribbon diagram of the b-galactosidase monomer showing the domain organization of the chain. Residues corresponding to successive domains are colored in successive spectral colors. See color insert.



At present, no structure of a complex of b-galactosidase with a substrate or analog is available, so the specific interaction sites between enzyme and substrate are not known. Nevertheless, earlier biochemical studies (14, 15) concluded that residues important for catalytic function or near the **active site** (Glu461, Met502, Tyr503, Glu537) are found close to one another and located around a deep pit within the third domain, which is postulated to be the substrate-binding site. Residues from distant regions of the sequence (Trp568, Trp999) also contribute to the active site.

The X-ray structure of b-galactosidase provides a structural rationale for both a- and w-complementation. The N-terminal segment, representing the a peptide, participates in forming a subunit interface, which in turn allows formation of the active site, made up of elements from two different subunits. The w-fragment, corresponding to the fifth domain of the subunit, exhibits a specific topology. Rotation of this domain by more than 90° positions Trp 999 at the active site in the vicinity of the third domain.

#### 4. b-Galactosidase as a Tool

The b-galactosidase of *E. coli* is one of the most commonly used enzymes in molecular biology for several reasons: (1) numerous genetic and biochemical aspects of the *lac* system are known; (2) several indicator media are available for detecting lactose metabolism and monitoring Lac<sup>+</sup> and Lac<sup>-</sup> bacteria; (3) the assay procedure based on o-nitrophenyl-b-D-galactoside (ONPG) hydrolysis is exquisitely sensitive; (4) the protein is highly stable and easy to purify; and (5) the availability of substituted galactoside derivatives permits selecting and screening *lac* mutants exhibiting a variety of phenotypes. For example, the chromogenic noninducing and colorless substrate, 5-bromo-4-chloro-3-indolyl-b-D-galactoside (X-Gal), forms an insoluble blue dye upon hydrolysis, thus providing a sensitive test for b-galactosidase in solid media. X-Gal also discriminates between strains that produce high and low levels of b-galactosidase, allowing detection of plaque-forming **bacteriophages** that carry the *lacZ* gene.

#### 5. b-Galactosidase Fusions

The observation that the N-terminal 23 residues of b-galactosidase can be replaced with other amino acids without affecting enzymatic activity enabled Müller-Hill and Kania (39) to obtain fusions between *lacI* and *lacZ* genes that produce a hybrid protein that has both **Lac repressor** and b-galactosidase activities. Several genetic methods have been developed to construct *lac* fusions that produce hybrid genes which specify hybrid proteins. Subsequently, *in vitro* construction of hybrid genes facilitated obtaining a variety of fusions (40). A number of commercially available vectors exist for constructing *lacZ* fusions *in vitro*. They contain a *lacZ* gene truncated at the 5'-end and preceded by a synthetic oligonucleotide containing multiple **restriction enzyme** cleavage sites, a *polycloning* site. Thus, if a 5'-coding sequence is cloned into one of these sites so that **transcription** and **translation** are restored across *lacZ*, a hybrid protein with b-galactosidase activity is produced. The sensitive detection of this activity with X-Gal has made b-galactosidase one of the most commonly used *in vivo* **reporter enzymes**. Using the *lacZ* gene as a reporter provides a sensitive method to detect genes subject to specific regulatory signals or to study the mechanism localizing a protein to a given cellular compartment. The regulation, identification, and localization of many proteins of various origins have been uncovered by using the b-galactosidase activity of fusion proteins as a marker.

During the last few years, a number of new eukaryotic expression vectors that produce b-galactosidase fusion proteins have been designed. They are currently being used to analyze developmentally regulated genes, transgenic expression, tissue-specific expression, and a number of other aspects involved in embryogenesis or developmental biology. This approach has been



facilitated by using replication-defective retroviral vectors that encode the *lacZ* gene (41), which has become a widely used cell lineage marking system. The labeled cells are easily identified by histochemical staining, using X-Gal or **immunofluorescence** with antibodies raised against b-galactosidase.

b-galactosidase fusions provide a unique and versatile experimental tool. The hybrid proteins are easy to purify (42) and can be used as **antigen** to raise antibodies against the target gene product. In addition, because there is no restriction on the size of the target gene DNA fused to *lacZ* and because many hybrid proteins retain all or a portion of the activity of the target gene product, these fusions provide a means to identify functional domains within the target gene product.

### 5.1. a-Complementation

Perhaps the most widely exploited property of b-galactosidase is a-complementation, described previously. Whereas the region of the foreign protein that can replace the N-terminus of b-galactosidase is more or less random in hybrid fusion proteins, replacements in the a-peptide for a-complementation are more restrictive. That insertion of a few amino acid residues does not always interfere with complementation capacity (43) was widely used during the late 70s to develop new cloning vectors, now commercially available. Indeed, a-complementation has become one of the most commonly used methods for identifying bacterial colonies containing recombinant DNA.

The first vector based on a-complementation was developed by Messing et al. (44). Within the major intergenic region of the filamentous bacteriophage M13, they inserted the *lac* regulatory region and the DNA sequence that codes for the first 146 amino acids of b-galactosidase that have a-complementation capacity. When infected with these phages, bacteria that produce the M15 a-acceptor protein produce active b-galactosidase by complementation and form blue plaques in the presence of the chromogenic substrate X-Gal. Insertion of foreign DNA into the a region of M13 interferes with a-complementation and gives rise to recombinants that form pale blue or colorless plaques. This simple test has made **cloning** in filamentous bacteriophages a routine procedure. The insertion of a series of synthetic restriction enzyme sites into the a-region provides a variety of targets for cloning foreign DNA fragments. Insertion of these polycloning sites, representing up to 18 extra amino acid residues, into the a-region has practically no effect on a-complementation capacity. Insertion of additional DNA into a cloning site, however, generally abolishes a-complementation and creates recombinant bacteriophages that produce colorless plaques when grown on X-Gal plates (*blue/white screen*). The great advantage of using vectors generated from single-stranded DNA phage, like M13, is that one obtains both the double-stranded **replicative form** from the infected bacteria and the single-stranded DNA, from the progeny phage. The latter is frequently used as a **template** for sequencing with the dideoxy-mediated chain-terminating **Sanger method** (45) and also for *in vitro* **mutagenesis**. Many other vectors that permit identification of recombinant clones based on a-complementation have been designed subsequently. The most commonly used are the pUC vectors (46).

Recently, Mohler and Blau (47) adapted intracistronic complementation of the *lacZ* gene for use in eukaryotic cells. By constructing replication-defective **retroviral** vectors that express a-donors and a-acceptors and w-donors and w-acceptors, complementation of the relevant *lacZ* mutants in mammalian cells permits analyzing cell fusion and detecting colocalized interacting proteins within single intact cells. a- and w-complementation can be exploited for a wide range of studies, including **transgenic** animals that express complementing *lacZ* mutants from two **promoters** of interest, which then could reveal cell lineages in which the products of both genes coincide spatially and temporally. In line with these results, Rossi et al. (48) applied b-galactosidase complementation to monitor protein-protein interactions in intact eukaryotic cells.

Given the wide range of applications for b-galactosidase during the last decades, one can anticipate that further development of the b-galactosidase fusion and complementation systems will contribute to our understanding of many features of cellular regulation and organization.

## Bibliography

1. M. Cohn, J. Monod, H. R. Pollock, S. Spiegelman, and R. Stanier (1953) *Nature* **172**, 1096–1097.
2. J. Monod, G. Cohen-Bazire, and M. Cohn (1951) *Biochem. Biophys. Acta* **7**, 585–599.
3. J. Lederberg and E. L. Tatum (1946) *Cold Spring Harbor Symp. Quant. Biol.* **11**, 113–114.
4. E. L. Wollman, F. Jacob, and W. Hayes (1956) *Cold Spring Harbor Symp. Quant. Biol.* **21**, 141–162.
5. H. W. Rickenberg, G. N. Cohen, G. Buttin, and J. Monod (1956) *Ann. Inst. Pasteur* **91**, 829–857.
6. I. Zabin and A. V. Fowler (1978) In *The Operon* (J. H. Miller and W. S. Reznikoff, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, pp. 89–121.
7. A. B. Pardee, F. Jacob, and J. Monod (1959) *J. Mol. Biol.* **1**, 165–178.
8. F. Jacob and J. Monod (1961) *J. Mol. Biol.* **3**, 318–356.
9. M. Cohn (1957) *Bact. Rev.* **21**, 140–168.
10. K. Wallenfels and O. P. Malhotra (1961) *Adv. Carbohydrate Chem.* **16**, 239–298.
11. A. Ullmann, F. Jacob, and J. Monod (1968) *J. Mol. Biol.* **32**, 1–13.
12. A. Ullmann and J. Monod (1969) *Biochem. Biophys. Res. Commun.* **35**, 35–42.
13. C. Burstein, M. Cohn, A. Kepes, and J. Monod (1965) *Biochem. Biophys. Acta* **95**, 634–639.
14. M. Ring and R. E. Huber (1990) *Arch. Biochem. Biophys.* **283**, 342–350.
15. C. G. Cupples, J. H. Miller, and R. E. Huber (1990) *J. Biol. Chem.* **265**, 5512–5518.
16. R. H. Jacobson, X-J. Zhang, R. F. DuBose, and B. W. Matthews (1994) *Nature* **369**, 761–766.
17. A. V. Fowler and I. Zabin (1977) *Proc. Natl. Acad. Sci. USA* **74**, 1507–1510.
18. A. Kalnins, K. Otto, U. Rütger, and B. Müller-Hill (1983) *EMBO J.* **2**, 593–597.
19. H. W. Stokes, P. W. Betts, and B. G. Hall (1985) *Mol. Biol. Evol.* **2**, 469–477.
20. A. Ullmann and D. Perrin (1970) In *The Lactose Operon* (J. R. Beckwith and D. Zipser, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, pp. 143–172.
21. I. Zabin (1982) *Mol. Cell. Biochem.* **49**, 87–96.
22. D. Perrin (1963) *Cold Spring Harbor Symp. Quant. Biol.* **28**, 529–532.
23. F. H. C. Crick and L. E. Orgel (1964) *J. Mol. Biol.* **8**, 161–165.
24. B. Rotman and F. Celada (1968) *Proc. Natl. Acad. Sci. USA* **60**, 660–667.
25. E. Conway de Macario, J. Ellis, R. Guzman, and B. Rotman (1978) *Proc. Natl. Acad. Sci. USA* **75**, 720–724.
26. F. Jacob, A. Ullmann, and J. Monod (1965) *J. Mol. Biol.* **31**, 704–719.
27. A. Ullmann, D. Perrin, F. Jacob, and J. Monod (1965) *J. Mol. Biol.* **12**, 918–923.
28. J. R. Beckwith (1964) *J. Mol. Biol.* **8**, 427–430.
29. A. Cook and J. Lederberg (1962) *Genetics* **47**, 1335–1353.
30. A. Ullmann, F. Jacob, and J. Monod (1968) *J. Mol. Biol.* **24**, 339–343.
31. S. L. Morrison and D. Zipser (1970) *J. Mol. Biol.* **50**, 359–371.
32. K. E. Langley, A. V. Fowler, and I. Zabin (1975) *J. Biol. Chem.* **250**, 2587–2592.
33. K. E. Langley, M. R. Villarejo, A. V. Fowler, P. J. Zamenhof, and I. Zabin (1975) *Proc. Natl. Acad. Sci. USA* **72**, 1254–1257.
34. R. S. Accolla and F. Celada (1976) *FEBS Lett.* **67**, 299–302.
35. M. E. Goldberg and S. J. Edelstein (1969) *J. Mol. Biol.* **46**, 431–440.
36. A. Ullmann and J. Monod (1970) in *The Lactose Operon* (J. R. Beckwith and D. Zipser, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, pp. 265–272.
37. F. Celada, A. Ullmann, and J. Monod (1974) *Biochemistry* **13**, 5543–5547.

38. M. E. Goldberg (1969) *J. Mol. Biol.* **46**, 441–446.
39. B. Müller-Hill and J. Kania (1974) *Nature* **249**, 561–563.
40. T. J. Silhavy and J. Beckwith (1985) *Microbiol. Rev.* **49**, 398–418.
41. J. R. Sanes, J. L. R. Rubenstein, and J-F. Nicolas (1986) *EMBO J.* **5**, 3133–3142.
42. A. Ullmann (1984) *Gene* **29**, 27–31.
43. U. Rütger, M. Koenen, K. Otto, and B. Müller-Hill (1981) *Nucleic Acids Res.* **9**, 4087–4098.
44. J. Messing, B. Gronenborn, B. Müller-Hill, and P.H. Hofschneider (1977) *Proc. Natl. Acad. Sci. USA* **74**, 3642–3646.
45. F. Sanger, S. Nicklen, and A. R. Coulson (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463–5467.
46. J. Vieira and J. Messing (1982) *Gene* **19**, 259–268.
47. W. A. Mohler and H. M. Blau (1996) *Proc. Natl. Acad. Sci. USA* **93**, 12423–12427.
48. F. Rossi, C. A. Charlton, and H. M. Blau (1997) *Proc. Natl. Acad. Sci. USA* **94**, 8405–8410.

### Suggestions for Further Reading

49. *The Lactose Operon* (1970). (J. R. Beckwith and D. Zipser, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, 437 pages.
50. A. Ullmann (1992) Complementation in  $\beta$ -galactosidase: From protein structure to genetic engineering, *BioEssays* **14**, 201–205.
51. B. Müller-Hill (1996) *The lac Operon. A Short History of a Genetic Paradigm*, Walter de Gruyter. Berlin, New York.

## Beta-Glucuronidase

$\beta$ -Glucuronidase (E.C. 3.2.1.31; GUS) is the most widely used [reporter gene](#) in plant molecular biology (1). The *gusA* gene cloned from *Escherichia coli* encodes a 68-kDa  $\beta$ -glucuronidase that forms a stable homotetramer and catalyzes the hydrolysis of a large number of glucuronides, in which D-glucuronic acid is conjugated through a  $\beta$ -O-glycosidic linkage to any aglycone.  $\beta$ -Glucuronidase has emerged as the reporter gene of choice to be used in plants, because there is very little or no endogenous glucuronidase activity across the phyla. This minimizes the detection of background activity and allows very small quantities of GUS to be detected. GUS activity may be followed in cell lysates or *in situ*. The major use of GUS has been to study gene expression patterns in transgenic plants by expressing GUS under the control of regulatory sequences of interest. A second use has been to monitor the intracellular fate of chimeric proteins by generating GUS fusion genes. For example, GUS has been fused to gene leader sequences that target the fusion gene to different organelles and allow intracellular trafficking of proteins to be studied (2).

GUS activity may be detected using a variety of substrates. As in the case of **beta-galactosidase** and X-gal, GUS expression is most commonly detected using histochemical substrates, such as X-GlcU (5-bromo-4-chloro-3-indolyl  $\beta$ -D-glucuronide), which gives a dark blue precipitate upon hydrolysis. Alternative substrates include 5-bromo-6-chloro-3-indolyl  $\beta$ -D-glucuronide, 6-chloro-3-indolyl  $\beta$ -D-glucuronide, and indoxyl  $\beta$ -D-glucuronide that give magenta, pink, and blue precipitates, respectively (3). The variety of substrates also makes GUS a valuable component of dual (or multiple) reporter gene systems; for example, GUS and *lacZ* expression could be detected in the same tissues by using substrates that give different colored histochemical precipitates. In addition to histochemical substrates, a number of fluorescent and chemiluminescent substrates have been

developed that are analogous to the substrates for b-galactosidase, such as CFDG-GlcU (Molecular Probes Inc.) and 4-methylumbelliferyl b-D-glucuronide (MUGlcU) (4).

In the case of *lacZ*, problems may be encountered in loading substrates into cells; to load FDG, for example, cells must undergo a moderate osmotic shock that can affect cell viability. In the case of GUS, however, a second gene, *gusB*, may be expressed that encodes a **permease** that actively takes up and transports glucuronide substrates into the cell. In addition to substrates that allow reporter gene visualization, a number of other bioactive molecules can be conjugated to glucuronides, which could then be released by hydrolysis in GUS-expressing cells. Thus, combined use of *gusA* and *gusB* increases the use of the reporter gene, from merely indicating gene expression to controlling a specific cell manipulation (1).

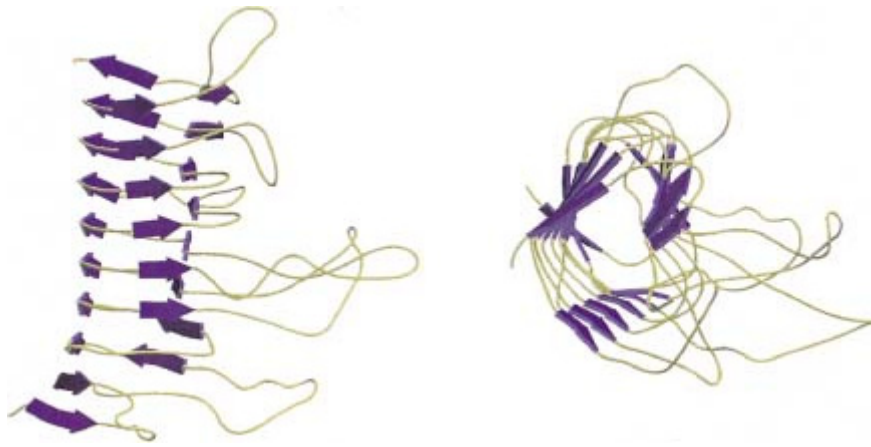
## Bibliography

1. R. A. Jefferson (1989) *Nature* **342**, 837–838.
2. R. A. Jefferson, T. A. Kavanagh, and M. W. Bevan (1987) *EMBO J.* **6**, 3901–3907.
3. G. A. Hull and M. Devic (1995) *Methods. Mol. Biol.* **49**, 125–141.
4. C. E. Olesen, J. J. Fortin, J. C. Voyta, and I. Browstein (1997) *Methods Mol. Biol.* **63**, 61–70.

## Beta-Helix

A b-helix is a type of [protein motif](#) or **domain** found in a very few [protein structures](#). It is characterized by an unusual parallel [b-sheet](#) topology formed from three parallel b-sheets wound together into a right-handed helical structure (Fig. 1). The [b-strands](#) in each b-sheet are short, having only 2 to 5 residues. Each coil of the b-helix has the same three-dimensional arrangement of a group of **secondary structure** elements and is thus similar to other coiled repeating structures, such as the b-roll, composed of two parallel b-sheets forming a b-helix, and the [leucine-rich repeat](#) of **ribonuclease inhibitor** that is composed of repeating [a-helices](#) and b-strands forming a horseshoe-shaped structure. In the b-helix, the side chains of repeating residues are packed into the center of the helix and interact with one another to form, for example, “asparagine ladders,” “serine stacks,” and/or “aromatic stacks.”

**Figure 1.** Orthogonal views of the b-helix of pectate lyase E (1), showing the three parallel b-sheets. For clarity, only the b-strand secondary structure (as purple arrows) and loops connecting the strands (yellow) are shown. Short helical regions and the *N*- and *C*-terminal loops have been removed. This figure was generated using Molscript (2) and Raster3D (3, 4). See color insert.



[See also [Beta-Sheet](#), [Protein Motif](#), [Domain](#), [Protein](#), and [Protein Structure](#).]

#### Bibliography

1. M. D. Yoder, S. E. Lietzke, and F. Jurnak (1993) *Structure* **1**, 241–251.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

5. M. D. Yoder and F. Jurnak (1995) The parallel  $\beta$ -helix and other coiled folds. *FASEB J.* **9**, 335–342.
6. C. Kisker et al. (1996) A left-handed  $\beta$ -helix revealed by the crystal structure of carbonic anhydrase from the archaeon bacteria *Methanosarcina thermophila*. *EMBO J.* **15**, 2323–2330.

## Beta-Lactamases

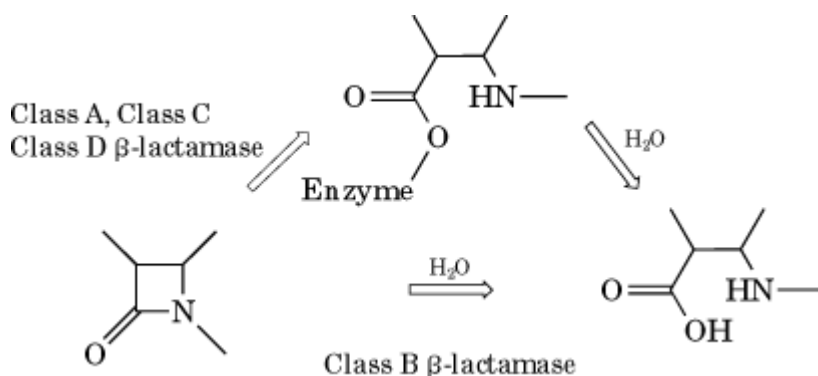
$\beta$ -Lactamases constitute one of the oldest known mechanisms of bacterial [antibiotic resistance](#) and one of the most widely distributed. The [enzymes](#) hydrolyze  $\beta$ -lactam antibiotics, such as [penicillin](#), rendering them inactive (Fig. 1). Most of these enzymes also react with other cyclic structures, such as  $\gamma$ -lactams and isatoic anhydride derivatives (1), and some have low reactivity with simple esters, [peptides](#), and depsipeptides. Richmond and Sykes (2) and Bush (3) classified  $\beta$ -lactamases according to their catalytic activity and apparent substrate specificity, which remains the most useful classification for the clinical microbiologist. Today, the ready availability of  $\beta$ -lactamase amino acid sequences obtained from their **genes** provides the basis for a structural classification first proposed by Ambler (4). There are two types of  $\beta$ -lactamase:

1. those with a [serine](#) residue at their [active site](#) that is transiently acylated by the  $\beta$ -lactam substrate during hydrolysis (Richmond and Sykes groups I and II; Ambler class A, plus the more recently distinguished classes C and D); for an example, see [Serine Proteinase](#).
2. those with metal ion cofactors, which are Zn(II) in the physiological state, and with apparently no acyl-enzyme intermediate (Richmond and Sykes group III; Ambler class B) (see

## [Metalloproteinases](#)).

The serine b-lactamases include enzymes that have both narrow and broad substrate specificities and are sensitive, to varying degrees, to b-lactam **enzyme-activated inhibitors**. This group is the more abundant and currently poses the greater clinical problem. The metallo-b-lactamases are active on a broad spectrum of b-lactam substrates and are much less sensitive to the mechanism-based inhibitors, although they can be inhibited by metal-ion chelating agents. This group has been relatively scarce in nature, but its prevalence is rising because of the increased use of b-lactam antibiotics that resist hydrolysis by serine b-lactamases.

**Figure 1.** The reaction catalyzed by b-lactamases.



### 1. Serine b-lactamases

Sequence comparisons of the active serine enzymes indicate six major groups:

1. The class A b-lactamases: Richmond and Sykes Group II enzymes; Bush groups 2a (penicillinases), 2b (broad spectrum b-lactamases), 2b' (extended broad spectrum b-lactamases), 2c (carbenicillinases), 2e (cephalosporinases). This group also includes a few proteins identified as D-Ala-D-Ala carboxypeptidases (DAC).
2. The class C b-lactamases: Richmond and Sykes and Bush Group 1, together with three proteins from *Streptomyces*, *Nocardia*, and *Bacillus subtilis*, identified as carboxypeptidase/transpeptidases.
3. The Class D b-lactamases from [gram-negative bacteria](#) (Bush Group 2d oxacillinases), together with b-lactam receptor proteins from gram-positive bacteria.
4. The bifunctional transglycosylase/transpeptidase, class A penicillin-binding proteins (PBP).
5. The monofunctional transpeptidase, class B PBP.
6. D-Ala-D-Ala carboxypeptidases (DAC).

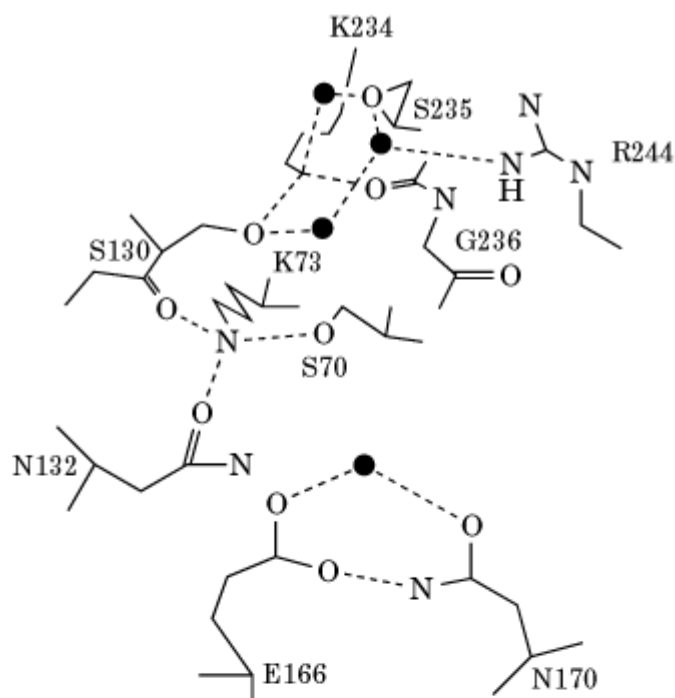
The last three groups are biosynthetic enzymes of the bacterial cell wall that react with D-Ala-D-Ala.

All six groups have three sets of conserved residues in common, which comprise the active site of the enzyme (see later). The first comprises the sequence Ser-X-X-Lys, where X is any residue and the serine residue is the one that is transiently acylated by substrates and the only residue that is absolutely conserved in all of the proteins. The second has either a serine (class A) or a [tyrosine](#) residue (in class C and class D b-lactamase groups), followed after one residue by an [asparagine](#) (except in one class D b-lactamase). The third conserved segment is usually Lys-Thr-Gly, but only the [glycine](#) residue is absolutely conserved. Sequence analyses have suggested that b-lactamase activity may have arisen several times by a process of evolutionary **convergence**.

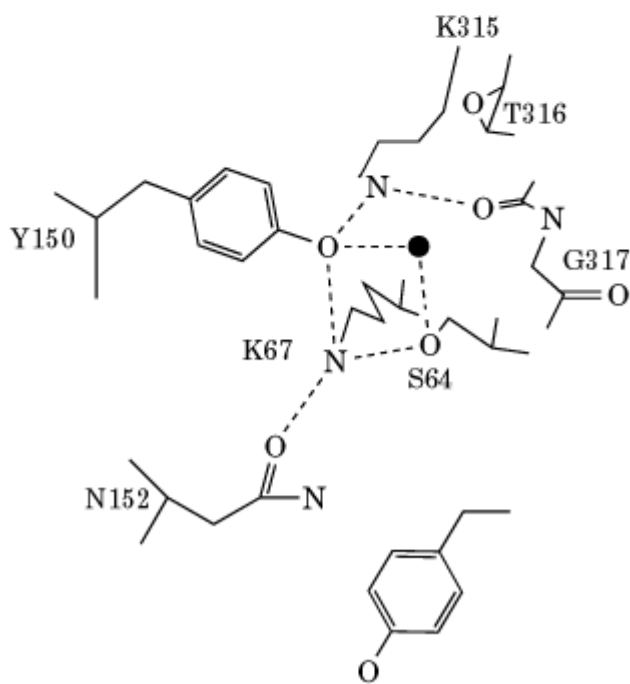
Three-dimensional structures of class A enzymes from *Staphylococcal aureus* (5), *Bacillus cereus*

(6), *B. licheniformis* (7), *Streptomyces albus* (8), and *Escherichia coli* TEM-1 (9, 10) have been published, which include only two of the five activity groups of class A enzymes. Structures of the class C  $\beta$ -lactamases from *Citrobacter freundii* (11) and *Enterobacter cloacae* (12) have been published, and this group is extended by the homologous structure of the transpeptidase domain of the *Streptomyces* carboxypeptidase/transpeptidase (13). The overall structures of all of these enzymes and a DAC are similar. The core of each molecule is a five-stranded, antiparallel  $\beta$ -sheet flanked on one side by three  $\alpha$ -helices and on the other by a larger, more diverse  $\alpha/\beta$  domain (see [Protein Structure](#)). The active site lies between the two domains and is bound by one edge of the  $\beta$ -sheet. As shown in Figure 2, the conserved residues listed previously make up the active site regions and occupy very similar positions in all of the structures (Figure 2). Even the alternative serine and tyrosine residues in the second conserved segment of class A (DAC) and class C (DAC), respectively, have their hydroxyl groups in the same positions; see Ser 130 and Tyr150 in Fig. 2 a and b, respectively. Major differences between the two classes of enzymes that are believed to have functional consequences occur in the area that serves as a recognition pocket for the side chain of the substrate. The class A enzymes have a loop (the W loop) that forms the base of the acyl amino side-chain binding pocket. This loop includes residues Glu166 and Asn170, which localize and activate a water molecule for attack on the ester bond of the substrate (14, 15). The projection of this loop into the active site also limits the scope for binding bulky acyl side chains, which was exploited in developing  $\beta$ -lactamase-resistant antibiotics, such as methicillin.

**Figure 2.** Comparison of the similar active-site regions of  $\beta$ -lactamases of class A (a) and class C (b). The residues are identified with the one-letter code. The lack spheres are water molecules observed crystallographically that are probably hydrogen-bonded (dashed lines) to various groups of the protein. The active-site serine residue that is acylated in the catalytic reaction is S70 in (a) and S64 in (b).



(a)



(b)

Evolutionary accumulation of peripheral mutations that modify the conformation of the W loop and of neighboring surface loops that define the edges of the side-chain recognition pocket has led to the appearance of enzymes that belong to groups 2b', 2c, and 2e that have extended substrate recognition profiles (16). In the class C enzymes, the water molecule that attacks the ester of the acyl intermediate is located on the opposite side of the plane of the ester and is activated by direct interaction with components of a hydrogen-bonded relay formed by the conserved residues. Identification of the individual residues involved in activating of the water has not yet been achieved,



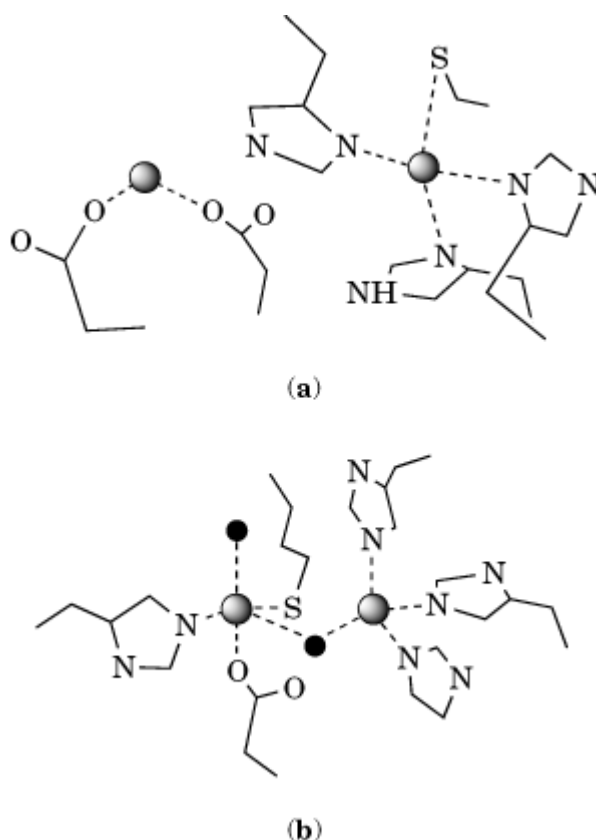
although Tyr150 has been implicated (11). The class C enzymes have a deeper, more **hydrophobic** side-chain binding pocket, lined by Tyr122, that enables them to bind substrates with large 7-acyl amino side chains.

Conformational changes during the b-lactamase reaction have been proposed for the class A, C, and D b-lactamases on a variety of grounds. With simple substrates, [NMR](#) and [circular dichroism](#) spectroscopy have suggested changes in the structure accompanying acyl-enzyme formation (17, 18). A number of substrates show nonstoichiometric bursts of hydrolysis (19), and a conformation change leading to a substrate-induced inactivation has been invoked (1, 20). These observations in solution contrast with the finding of very few differences in the structures of the free enzyme and the acyl-enzyme complex determined by [X-ray crystallography](#) (9, 11). An unambiguous description of the reaction mechanism must wait until these observations can be reconciled.

## 2. Metallo-b-lactamases

Sequence comparisons suggest that there are four groups of this type of b-lactamase that are not closely related, but share similarities with a protein from the actinorhodin biosynthetic gene cluster of *Streptomyces*. A few residues involved in binding the metal ion cofactors are absolutely conserved within this diverse group. Three structures are available, that representing two of the major subgroups of metallo-b-lactamases (Fig. 3). All three structures are similar overall to each other, but not to the serine b-lactamases, and the core of each molecule is formed from two antiparallel b-sheets, each with additional parallel strands and flanked by  $\alpha$ -helices. The structure of *Bacillus cereus* b-lactamase II was solved with Cd(II) in place of the natural Zn(II) metal ion (21). One Cd(II) ion is bound tightly, with a dissociation constant of about 1  $\mu$ M, by the side chains of three [histidine](#) and one [cysteine](#) residue (22-24). In the homologous enzyme from *Bacteroides fragilis*, one Zn(II) ion is similarly coordinated by three histidine side chains, as in the *B. cereus* enzyme, but the fourth ligand is a water molecule that is shared with a second Zn(II) ion (25). The second zinc ion is further coordinated by three other side chains, those of His, Cys, and Asp residues. Two more Cd(II) ions are located in the *B. cereus* structure, but the coordination of one of them is much weaker (23), so it may not be physiological. Early stoichiometries determined for the binding of Zn(II) and Cd(II) suggested that two ions per molecule are bound (22). The second Cd(II) ion is bound by two Asp residues and lies only 10 Å from the tightly bound ion (Fig. 3), where it might modulate the catalytic activity. It is believed that the water molecule complexed with the two Zn(II) ions in the *B. fragilis* enzyme is activated to form hydroxide for attack on the b-lactam ring, and a corresponding water molecule is probably activated by the tightly bound ion in the *B. cereus* enzyme.

**Figure 3.** Comparison of the similar active-site regions of metallo-b-lactamases from *Bacillus cereus* (a) and *Bacteroides fragilis* (b). The shaded spheres are the two metal ions, cadmium (II) in (a) and zinc (II) in (b), and the smaller solid spheres are water molecules observed crystallographically.

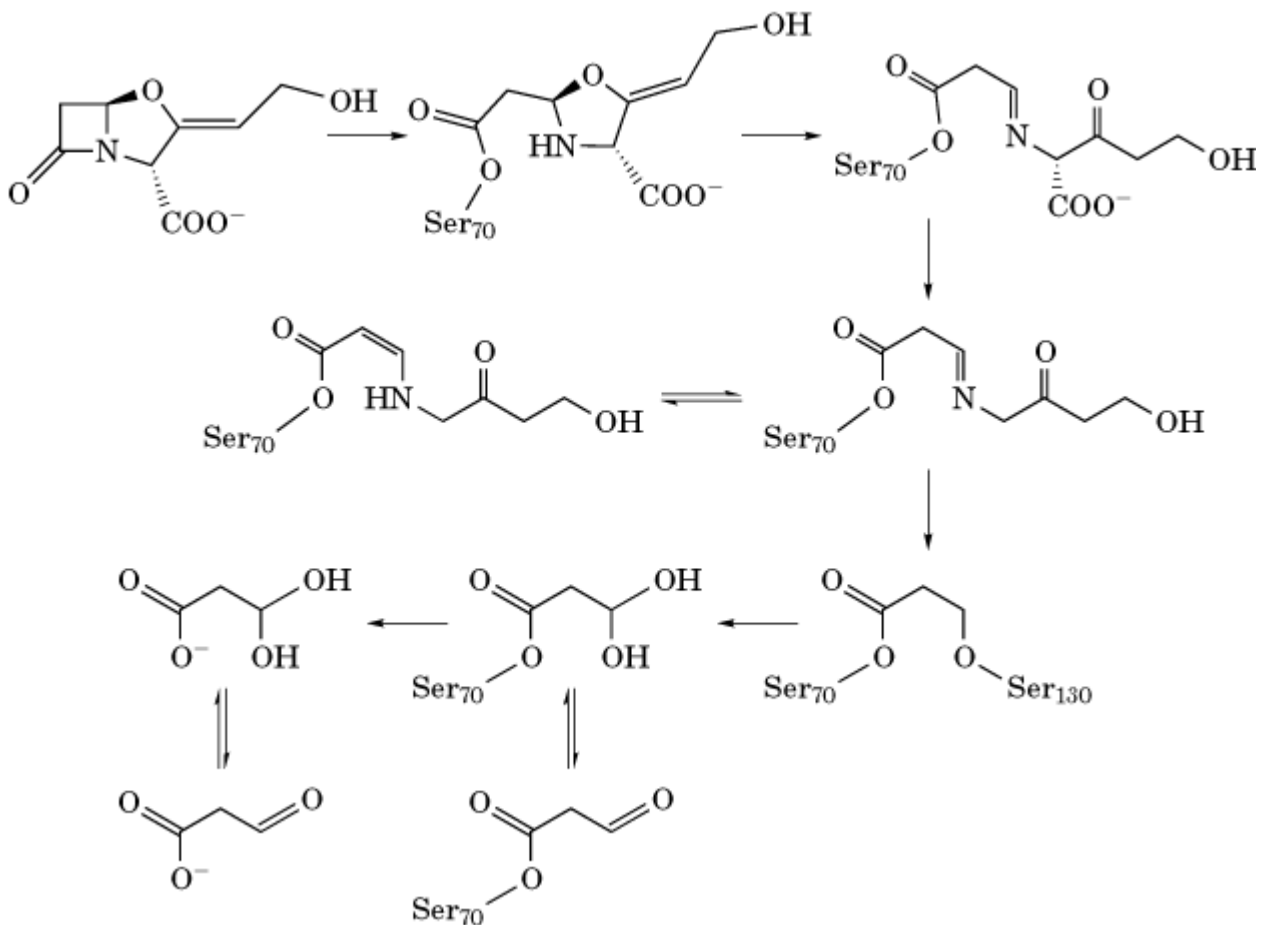


### 3. Inhibition

The important role of  $\beta$ -lactamases in antibiotic resistance has led to the development of specific  $\beta$ -lactamase inhibitors that can be used to protect  $\beta$ -lactamase-sensitive antibiotics. For the serine  $\beta$ -lactamases, several classes of mechanism-based inhibitors have been discovered, and several inhibitor proteins are known (26). In many cases, the mechanism of action of the small-molecule inhibitors is not clear, and only two classes have been used in combination with antibiotics.

*Clavulanic acid* has been widely used clinical application. It is selective for the class A and class D  $\beta$ -lactamases. Opening of the  $\beta$ -lactam ring of clavulanic acid by the initial attack of the  $\beta$ -lactamase triggers a series of chemical rearrangements in the inhibitor moiety (Fig. 4) that results in an acyl-enzyme that is resistant to attack by water and may even be cross-linked by the inhibitor moiety (27).

**Figure 4.** Suggested mechanism for the reaction of TEM-2 class A  $\beta$ -lactamases with clavulanic acid (27).



The penam sulfone acids *sulbactam* and *tazobactam* are also used clinically as  $\beta$ -lactamase inhibitors. As with clavulanic acid, a series of chemical rearrangements are provoked by the reaction with  $\beta$ -lactamase. Although sulbactam and tazobactam are relatively selective for class A and class D  $\beta$ -lactamases, analogs that have increased activity against class C  $\beta$ -lactamases have been reported (28).

#### 4. Regulation of Biosynthesis

Expression of chromosomal  $\beta$ -lactamase in *Bacillus* and of **plasmid**-encoded  $\beta$ -lactamases in staphylococci is under the control of a typical **repressor** protein (BlaI and MecI, respectively) and a second regulatory protein (BlaR and MecR, respectively). The BlaR and MecR regulatory proteins are similar. Both are integral **membrane proteins** that have a  $\beta$ -lactam-binding **domain** homologous to the class D serine  $\beta$ -lactamases (29). The predicted topology of the polypeptide in the membrane is three to five membrane-spanning  **$\alpha$ -helices**, and the N-terminus in the cytoplasm and the  $\beta$ -lactam-binding domain is on the outer surface of the membrane, where it could act as a **receptor** for  $\beta$ -lactams.

Expression of chromosomal class C  $\beta$ -lactamases in Enterobacteriaceae from the *ampC* gene is under the control of the *ampD*, *ampE*, *ampG*, and *ampR* gene products. AmpR protein is a typical **transcription** activator that is lost in organisms, such as *E. coli*, that have constitutive AmpC production. AmpR responds to the binding of 1,6-anhydro-*N*-acetylmuramyl-[L]-alanyl-[D]-glutamyl-*meso*-diaminopimelic acid (MurNAc tripeptide) produced by breakdown of the peptidoglycan of the cell wall (30). The three other proteins are involved in transmembrane signaling for induction of the  $\beta$ -lactamase. Inactivation of AmpD results in massive overproduction of  $\beta$ -lactamase, whereas loss of either AmpE or AmpG activity results in a total block of induction.

AmpD is an amidase that cleaves the tripeptide from the MurNAc tripeptide, thus inactivating it as an inducer. It is thought that AmpG is an integral membrane protein that acts as a transporter for the MurNAc tripeptide, whereas it has been suggested that AmpE, also an integral membrane protein, provides energy for its uptake (31).

## Bibliography

1. M. G. P. Page (1993) *Biochem. J.* **295**, 295–304.
2. M. H. Richmond and R. B. Sykes (1973) *Adv. Microb. Physiol.* **9**, 31–88.
3. K. Bush (1989) *Antimicrob. Agents Chemother.* **33**, 264–270, 271–276.
4. R. P. Ambler (1980) *Phil. Trans. R. Soc. Lond.* **B239**, 321–331.
5. O. Herzberg and J. Moulton (1987) *Science* **236**, 694–701.
6. B. Samaraoui, B. Sutton, R. Todd, P. Artymyuk, S. G. Waley, and D. Phillips (1986) *Nature* **320**, 378–380.
7. P. C. Moews, J. R. Knox, O. Dideberg, P. Charlier, and J.-M. Frère (1990) *Proteins: Struct. Function Genet.* **7**, 156–171.
8. O. Dideberg, P. Charlier, J. P. Wery, P. Dehottay, J. Dusart, T. Erpicum, J.-M. Frère, and J.-M. Ghuysen (1987) *Biochem. J.* **245**, 911–913.
9. N. C. J. Strynadka, H. Adachi, S. E. Jensen, K. Johns, A. Sielecki, C. Betzel, K. Sutoh, and M. N. G. James (1992) *Nature* **359**, 700–705.
10. C. Jelsch, L. Mourey, J. M. Masson, and J.-P. Samama (1992) *Proteins: Struct. Function Genet.* **16**, 364–383.
11. C. Oefner, A. D'Arcy, J. J. Daly, K. Gubernator, R. L. Charnas, I. Heinze, C. Hubschwerlen, and F. K. Winkler (1990) *Nature* **343**, 284–288.
12. E. Lobkovsky, P. C. Moews, J. Liu, H. Zhao, J.-M. Frère, and J. R. Knox (1993) *Proc. Natl. Acad. Sci. USA* **90**, 11257–11261.
13. J. A. Kelly and A. P. Kuzin (1995) *J. Mol. Biol.* **254**, 223–236.
14. R. M. Gibson, H. Christensen, and S. G. Waley (1990) *Biochem. J.* **272**, 613–619.
15. W. A. Escobar, A. K. Tan, and A. L. Fink (1991) *Biochemistry* **30**, 10783–10787.
16. G. A. Jacoby and A. A. Medeiros (1991) *Antimicrob. Agents Chemother.* **35**, 1697–1704.
17. M. Jamin, C. Damblon, A. M. Baudin-Misselyn, F. Durant, G. C. K. Roberts, P. Charlier, G. Llabres, J.-M. Frère. (1994) *Biochem. J.* **301**, 199–203.
18. P. Taibi-Tonche, I. Massova, S. B. Vakulenko, S. A. Lerner, and S. Mobashery (1996) *J. Am. Chem. Soc.* **118**, 7441–7448.
19. A. Samuni and N. Citri (1975) *Biochem. Biophys. Res. Commun.* **62**, 7–11.
20. S. G. Waley (1991) *Biochem. J.* **279**, 87–94.
21. B. J. Sutton, P. J. Artymyuk, A. E. Cordero-Borboa, C. Little, D. C. Phillips, and S. G. Waley (1987) *Biochem. J.* **248**, 181–188.
22. R. B. Davies and E. P. Abraham (1974) *Biochem. J.* **143**, 129–135.
23. G. S. Baldwin, A. Galdes, A. O. Hill, B. E. Smith, S. G. Waley, and E. P. Abraham (1978) *Biochem. J.* **175**, 441–447.
24. A. Carfi and O. Dideberg (1995) *EMBO J.* **14**, 4919–4921.
25. N. O. Concha, B. A. Rasmussen, K. Bush, and O. Herzberg (1996) *Structure* **4**, 823–836.
26. N. C. J. Strynadka, S. E. Jensen, K. Johns, H. Blanchard, M. Page, A. Matagne, J.-M. Frère, and M. N. G. James (1994) *Nature* **368**, 657–660.
27. R. P. A. Brown, R. T. Aplin, and C. J. Schofield (1996) *Biochemistry* **35**, 12421–12432.
28. H. G. F. Richter, P. Angehrn, C. Hubschwerlen, M. Kania, M. G. P. Page, J.-L. Specklin, and F. K. Winkler (1996) *J. Med. Chem.* **39**, 3712–3722.
29. T. Kobayashi, Y. F. Zhu, N. J. Nicholls, and J. O. Lampen (1987) *J. Bacteriol.* **169**, 3873–3878.

30. J. T. Park (1995) *Mol. Microbiol.* **17**, 4521–426.
31. S. Lindquist, M. Galleni, F. Lindberg, and S. Normark (1989) *Mol. Microbiol.* **3**, 1091–1102.

### Suggestion for Further Reading

32. M. I. Page, ed. (1992) *The Chemistry of  $\beta$ -Lactams* Blackie Academic & Professional, Glasgow. An excellent coverage of all aspects of the interaction of  $\beta$ -lactamases with substrates and inhibitors. Particularly useful chapters are " $\beta$ -Lactamase: mechanism of action" by S. G. Waley (pp. 198–228) and " $\beta$ -Lactamase: Inhibition" by R.F. Pratt (pp. 229–271).

## Beta-Lactoglobulin

$\beta$ -Lactoglobulin is the major whey [protein](#) in the milk of ruminants, and it has also been reported in milk from many other species, although not human, lagomorph, or rodent. Where found, genetic variants have also been reported in most cases. Isolation of whey protein from bovine milk is straightforward ([1](#)), and its concentration is 2 to 3 g/L. Consequently, the protein has been used since its first isolation in 1934 as a convenient, small, soluble protein with which to develop and calibrate new techniques. For example, it is used as a component of standard mixtures for [isoelectric focusing](#) and [SDS-PAGE](#) and a calibration sequence for automatic sequencers (see [Protein Sequencing](#)). The extensive literature on  $\beta$ -lactoglobulin, largely that from the domestic cow, can be broadly divided into studies related to the molecule as a protein and those related to its importance to the dairy industry. Reviews of both abound ([2-6](#)).

The polypeptide chain of the protein contains about 160 residues. That of ruminant species is dimeric, and in other species it is monomeric. [Figure 1](#) shows some of the known sequences from the SWISSPROT database ([7](#)). [Disulfide bonds](#) link cysteines 66 to 160 and 106 to 119. The asterisks show structurally conserved regions that were recognized following the structural determination first of the homologous plasma *retinol-binding protein* and then of bovine  $\beta$ -lactoglobulin ([8](#), [9](#)). These motifs define a widely distributed family called the *lipocalins* ([10](#), [11](#)), whose structure is an eight-stranded antiparallel  $\beta$ -barrel with (+1)<sub>8</sub> topology (see [Beta-Sheet](#)) and with an  **$\alpha$ -helix** on the outer surface. The structure of  $\beta$ -lactoglobulin is shown in [Fig. 2](#) ([12](#)). The family relationship has also been observed in the gene sequences, which show similar arrangements of **introns** and **exons** ([13](#)). The 400 bp upstream of the **open reading frame** controls the high level of expression and its restriction to the mammary gland ([14](#)). This has led to successful use of the  $\beta$ -lactoglobulin operon in [transgenic technology](#) ([15](#)).

**Figure 1.** Sequences of mature  $\beta$ -lactoglobulins with the totally conserved residues in bold. The sites of possible genetic shown in the bovine sequence in italics, and changes in the bovine B variant are shown above. The N-terminal sequence ([34](#)). The asterisks show the structurally conserved regions indicative of the lipocalin family.

\*\*\*\*\*

|          |                             |                                       |                              |                     |       |
|----------|-----------------------------|---------------------------------------|------------------------------|---------------------|-------|
| BOVINE A | <b>L</b> IVTQ <b>T</b> MKGL | <b>D</b> IQKV <b>A</b> GT <b>W</b> Y  | SLAMAAS <b>D</b> IS          | <b>L</b> LLDAQSAPLR | VYVEI |
| PIG I    | VEVTPIMTEL                  | DTQKV <b>A</b> GT <b>W</b> H          | TVAMAV <b>S</b> D <b>V</b> S | <b>L</b> LLDAKSSPLK | AYVEI |
| HORSE II | TDIPQ <b>T</b> MQDL         | DLQEV <b>A</b> GR <b>W</b> H          | SVAMV <b>A</b> S <b>D</b> IS | <b>L</b> LLDSESVPLR | VYVEI |
| DONKEY   | TNIPQ <b>T</b> MQDL         | DLQEV <b>A</b> G <b>K</b> W <b>H</b>  | SVAMAAS <b>D</b> IS          | <b>L</b> LLDSEEAPLR | VYIEI |
| DOG I    | IVVPR <b>T</b> MEDL         | DLQKV <b>A</b> GT <b>W</b> H          | SMAMAAS <b>D</b> IS          | <b>L</b> LLDSETAPLR | VYIQI |
| CAT I    | ATVPL <b>T</b> MDGL         | DLQKV <b>A</b> G <b>M</b> W <b>H</b>  | SMAMAAS <b>D</b> IS          | <b>L</b> LLDSETAPLR | VYVQI |
| DOLPHIN  | VSVIR <b>T</b> MEDL         | DIQ <b>R</b> V <b>A</b> GT <b>W</b> H | SVAMAAS <b>D</b> IS          | <b>L</b> LLDTEEAPLR | VNVEI |
| KANGAROO | VENIRSKNDL                  | GVEK <b>F</b> V <b>G</b> SWY          | LREAAKT---                   | -MEF-SIPLFDMDIKI    |       |
| MONKEY   | IDSPQ <b>T</b> MQDV         | DLPK <b>L</b> A <b>G</b> T <b>W</b> H | SMAMAA...                    |                     |       |

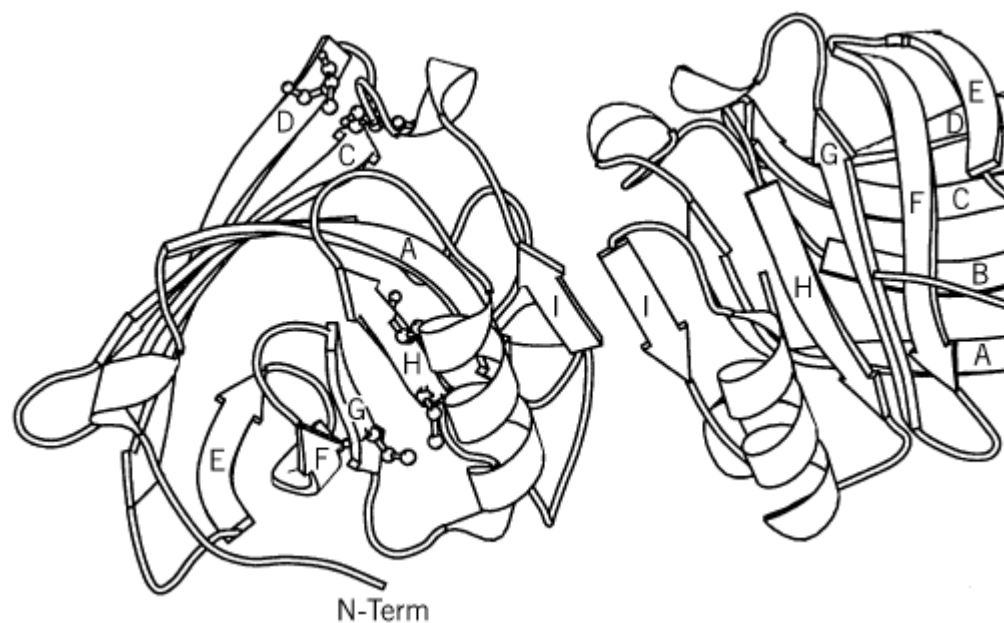
G

|          |  |  |  |                              |                                |
|----------|--|--|--|------------------------------|--------------------------------|
| BOVINE A | WEN <b>D</b> E <b>C</b> AQ <b>K</b> K                  | I <b>I</b> A <b>E</b> K <b>T</b> K <b>I</b> PA | <b>V</b> FK <b>I</b> D <b>A</b> L <b>N</b> E <b>N</b>  | --KVLVLD <b>T</b> D          | <b>Y</b> KK <b>Y</b> I         |
| PIG I    | REN <b>D</b> K <b>C</b> AQ <b>E</b> V                  | LLAK <b>K</b> T <b>D</b> I <b>P</b> A          | <b>V</b> FK <b>I</b> N <b>A</b> L <b>D</b> E <b>N</b>  | --QLFLLD <b>T</b> D          | <b>Y</b> D <b>S</b> H <b>I</b> |
| HORSE II | G <b>A</b> N <b>H</b> A <b>C</b> VER <b>N</b>          | I <b>V</b> AQ <b>K</b> T <b>E</b> D <b>P</b> A | <b>V</b> F <b>T</b> V <b>N</b> Y <b>Q</b> G <b>E</b> R | --KISVLD <b>T</b> D          | <b>Y</b> A <b>H</b> Y <b>M</b> |
| DONKEY   | GEN <b>K</b> G <b>C</b> A <b>E</b> KK                  | I <b>F</b> A <b>E</b> K <b>T</b> E <b>S</b> PA | <b>E</b> FK <b>I</b> N <b>Y</b> L <b>D</b> E <b>D</b>  | --TVFALD <b>S</b> D          | <b>Y</b> K <b>N</b> Y <b>I</b> |
| DOG I    | W <b>E</b> DGR <b>C</b> A <b>E</b> Q <b>K</b>          | V <b>L</b> A <b>E</b> K <b>T</b> E <b>V</b> PA | <b>E</b> FK <b>I</b> N <b>Y</b> V <b>E</b> E <b>N</b>  | --QIFLLD <b>T</b> D          | <b>Y</b> D <b>N</b> Y <b>I</b> |
| CAT I    | W <b>E</b> DNR <b>C</b> V <b>E</b> KK                  | V <b>L</b> A <b>E</b> K <b>T</b> E <b>C</b> AA | <b>K</b> F <b>N</b> I <b>N</b> Y <b>L</b> D <b>E</b> N | --ELIVLD <b>T</b> D          | <b>Y</b> E <b>N</b> Y <b>I</b> |
| DOLPHIN  | R <b>D</b> K <b>N</b> G <b>C</b> V <b>K</b> E <b>K</b> | I <b>I</b> A <b>E</b> K <b>T</b> E <b>I</b> PA | <b>V</b> FK <b>I</b> N <b>F</b> L <b>N</b> E <b>N</b>  | --KIFVLD <b>S</b> D          | <b>Y</b> T <b>N</b> Y <b>I</b> |
| KANGAROO | -KTDR <b>C</b> V <b>E</b> KK                           | LLL <b>K</b> K <b>T</b> K <b>K</b> P <b>T</b>  | <b>E</b> F <b>E</b> I <b>Y</b> I <b>S</b> S <b>E</b> S | SYTFCV <b>M</b> E <b>T</b> D | <b>Y</b> D <b>S</b> Y <b>I</b> |

A \*\*\*\*\*

|          |  |  |  |  |                                |
|----------|--|--|--|--|--------------------------------|
| BOVINE A | <b>V</b> CQCL <b>V</b> R <b>T</b> PE                   | V <b>D</b> E <b>A</b> L <b>E</b> K <b>F</b> D          | K <b>A</b> L <b>K</b> - <b>A</b> L <b>P</b> M <b>H</b> | IRL <b>S</b> F <b>N</b> P <b>T</b> Q <b>L</b>          | <b>E</b> E <b>Q</b> C <b>H</b> |
| PIG I    | <b>V</b> CQSL <b>A</b> R <b>T</b> LE                   | V <b>D</b> DQ <b>I</b> R <b>E</b> K <b>F</b> E         | D <b>A</b> L <b>K</b> - <b>T</b> L <b>S</b> V <b>P</b> | M <b>R</b> I--L <b>P</b> A <b>Q</b> L                  | <b>E</b> E <b>Q</b> C <b>H</b> |
| HORSE II | <b>V</b> CQYL <b>A</b> R <b>T</b> Q <b>K</b>           | V <b>D</b> E <b>E</b> V <b>M</b> E <b>K</b> F <b>S</b> | R <b>A</b> L <b>Q</b> - <b>P</b> L <b>P</b> G <b>R</b> | VQ <b>I</b> VQ <b>D</b> P <b>S</b> G <b>G</b>          | Q <b>E</b> R <b>C</b> H        |
| DONKEY   | <b>V</b> CQYL <b>A</b> R <b>T</b> Q <b>M</b>           | V <b>D</b> E <b>E</b> I <b>M</b> E <b>K</b> F <b>R</b> | R <b>A</b> L <b>Q</b> - <b>P</b> L <b>P</b> G <b>R</b> | VQ <b>I</b> V <b>P</b> D <b>L</b> T <b>R</b> M         | A <b>E</b> R <b>C</b> H        |
| DOG I    | <b>M</b> CQCL <b>A</b> R <b>T</b> LE                   | V <b>D</b> E <b>N</b> V <b>M</b> E <b>K</b> F <b>N</b> | R <b>A</b> L <b>K</b> - <b>T</b> L <b>P</b> V <b>H</b> | M <b>Q</b> L <b>L</b> N- <b>P</b> T <b>Q</b> A         | <b>E</b> E <b>Q</b> C <b>I</b> |
| CAT I    | <b>V</b> CQCL <b>T</b> R <b>T</b> L <b>K</b>           | A <b>D</b> N <b>E</b> V <b>M</b> E <b>K</b> F <b>D</b> | R <b>A</b> L <b>Q</b> - <b>T</b> L <b>P</b> V <b>H</b> | V <b>R</b> L <b>F</b> F <b>D</b> P <b>T</b> Q <b>V</b> | A <b>E</b> Q <b>C</b> H        |
| DOLPHIN  | <b>T</b> C <b>A</b> Y <b>L</b> A <b>R</b> T <b>L</b> Q | V <b>D</b> D <b>G</b> V <b>M</b> E <b>K</b> F <b>N</b> | K <b>A</b> I <b>K</b> P <b>A</b> L <b>P</b> M <b>H</b> | IRL- <b>F</b> S <b>P</b> T <b>Q</b> L                  | <b>E</b> E <b>Q</b> C <b>H</b> |
| KANGAROO | <b>A</b> C <b>A</b> H <b>Y</b> V <b>R</b> R <b>I</b> E | - <b>E</b> N <b>K</b> G <b>M</b> N <b>E</b> F <b>K</b> | K <b>I</b> L <b>R</b> -- <b>T</b> - <b>L</b> A         | M <b>P</b> Y <b>T</b> V <b>I</b> E <b>V</b> R <b>T</b> | R <b>D</b> M <b>C</b> H        |

Figure 2. A MOLSCRIPT (35) diagram of the dimer of b-lactoglobulin showing the fold of the main chain.



No definite function has yet been ascribed to b-lactoglobulin, but its presence in a family containing mostly small, secreted proteins that transport **hydrophobic** molecules perhaps implies that its function, other than the obvious nutritional one, is associated with the transport or uptake of fatty acids or retinol, both of which it binds (16, 17). Indeed, over 20 different **ligands** have been reported that have **association constants** varying between  $5 \times 10^7$  and  $4 \times 10^2 \text{M}$  (3). By analogy with other lipocalins, the binding site should be in the central pocket formed by the b-strands, but there is no definitive evidence that this is the case. Recently it has been shown that palmitate and retinol can bind independently and simultaneously (18).

Some physical parameters of the protein are summarized in Table 1. The bovine protein is stable down to pH 2, and conformational studies have shown that there are several pH-dependent changes between pH 2 and pH 9 (19, 20). At the extremes of pH, the dimer dissociates, and it aggregates above pH 8.5 to 9.0. The conformational change between pH 6 and 8 leads to increased reactivity of the free **cysteine** residue, perhaps triggered by the titration of a **carboxyl group** with a  $\text{p}K_a$  of 7.3. An association of four dimers forms an octamer at 5°C and pH 4.5 mainly in the A variant. This probably results from a carboxyl–carboxylate interaction involving the Gly64Asp change between the most common B and A genetic variants of b-lactoglobulin (21).

**Table 1. Selected Fundamental Parameters of Bovine b-Lactoglobulin**

|   |   |
|---|---|
| Number of amino acids                         | 162   |
| A variant                                     | Asp64, Val118                                     |
| B variant                                     | Gly64, Ala118                                     |
| Relative molecular mass                       | 18,400  |
| Sedimentation coefficient ( $S_{20w}^\circ$ ) | $2.83 \times 10^{13} \text{ s}^{-1}$              |
| Diffusion coefficient                         | $7.70 \times 10^{-7} \text{ cm}^2 \text{ s}^{-1}$ |
| Stokes' radius                                | 2.68 nm   |
| Radius of gyration                            | 2.17 nm   |

|                      |                                    |
|----------------------|------------------------------------|
| Axial ratio          | 2:1                                |
| Dipole moment        | 730 Debye                          |
| Absorbance at 280 nm | 0.96 (1.0 g/l solution, 1 cm cell) |

---

The complex behavior during refolding from [urea](#) between 2° and 25°C is ascribed to **disulfide interchange**. Cys106 pairs with either Cys119 or Cys121, and the former is the native ([22](#), [23](#)). It is not yet clear whether this also explains the partially folded form observed at pH 2, where the protein is predominantly a monomer ([24](#)). Refolding from **guanidinium chloride** at low pH shows a transient increase in helical content ([25](#)), an effect that is stabilized in ethanolic solution ([26](#)).

b-Lactoglobulin has been cloned and expressed both in *Escherichia coli* and in **yeast**, allowing creation of **site-directed mutations** ([27](#)) for probing the solution properties ([28](#)) or modifying the thermal denaturing behavior ([29](#)). Thermal denaturation is of particular interest to the dairy industry because heating milk is the basis of much processing and also the genetic variants exhibit different behavior ([30](#)). In addition, protein concentration, pH, temperature, ionic strength, dielectric constant, and buffer type all contribute to the denaturing behavior. Differential scanning **calorimetry** and thermal aggregation studies show that the A variant is less stable than the B variant ([31](#)), but studies at lower protein concentration show that the reverse is true ([32](#)). Thermal aggregation also involves intermolecular disulfide interchange, presumably initiated by the free [thiol group](#) ([33](#)).

#### Bibliography

1. R. Aschaffenburg and R. Drewry (1957) *Biochem. J.* **65**, 273–277.
2. J. M. A. Tilley (1960) *Dairy Sci. Abstr.* **22**, 111–125.
3. R. L. J. Lyster (1972) *J. Dairy Res.* **39**, 279–318.
4. S. G. Hambling, A. S. McAlpine, and L. Sawyer (1992) In *Advanced Dairy Chemistry I*. (P. F. Fox, ed.), Elsevier, Amsterdam, pp. 141–190.
5. M. McSwiney, H. Singh, O. Campanella, and L. K. Creamer (1994) *J. Dairy Res.* **61**, 221–232.
6. J. E. Kinsella and D. M. Whitehead (1987) *Adv. Food Nutr. Res.* **33**, 343–438.
7. A. Bairoch and B. Boeckmann (1994) *Nucleic Acid Res.* **22**, 3578–
8. M. Newcomer et al. (1984) *EMBO J.* **3**, 1451–1454.
9. M. Z. Papiz et al. (1986) *Nature* **324**, 383–385.
10. S. Pervaiz and K. Brew (1987) *FASEB J.* **1**, 209–214.
11. D. R. Flower (1996) *Biochem. J.* **318**, 1–14.
12. S. Brownlow et al. (1997) *Structure*, **5**, 481–495.
13. S. Ali and A. J. Clark (1988) *J. Mol. Biol.* **189**, 415–426.
14. J. P. Simons, M. McClenaghan, and A. J. Clark (1987) *Nature* **328**, 530–532.
15. A. L. Archibald et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 5178–5182.
16. A. A. Spector and J. E. Fletcher (1970) *Lipids* **5**, 403–411.
17. F. D. Fugate and P. S. Song (1980) *Biochim. Biophys. Acta* **652**, 28–42.
18. M. Narayan and L. J. Berliner (1997) *Protein Sci.* **7**, 150–157.
19. C. Tanford, L. G. Bunville, and Y. Nozaki (1959) *J. Am. Chem. Soc.* **81**, 4032–4036.
20. S. N. Timasheff, L. Mescanti, J. J. Basch, and R. Townend (1966) *J. Biol. Chem.* **241**, 2496–2501.
21. J. Witz, S. N. Timasheff, and V. Luzzati (1964) *J. Am. Chem. Soc.* **86**, 168–173.
22. H. A. McKenzie, G. B. Ralston, and D. C. Shaw (1972) *Biochemistry* **11**, 4539–4547.



23. T. E. Creighton (1980) *J. Mol. Biol.* **137**, 61–80.
24. L. Ragona et al. (1997) *Folding and Design* **2**, 281–290.
25. D. Hamada, S. Segawa, and Y. Goto (1996) *Nature Struct. Biol.* **3**, 868–873.
26. E. Dufour, H. C. Bertrand, and T. Haertlé (1993) *Biopolymers* **33**, 589–598.
27. C. A. Batt, L. D. Rabson, D. W. S. Wong, and J. E. Kinsella (1990) *Agric. Biol. Chem.* **54**, 949–955.
28. Y. Katakura, M. Totsuka, A. Ametani, and S. Kaminogawa (1994) *Biochim. Biophys. Acta* **1207**, 58–67.
29. C. A. Batt, J. Brady, and L. Sawyer (1994) *Tr. Food Sci.* **5**, 261–265.
30. E. Jakob and Z. Puhani (1992) *Int. Dairy J.* **2**, 157–178.
31. X. L. Huang, G. L. Catignani, and H. E. Swaisgood (1994) *J. Agric. Food Chem.* **42**, 1276–1280.
32. G. I. Imafadon, K. F. Ng-Kwai-Hang, V. R. Harwalkar, and C.-Y. Ma (1991) *J. Dairy Res.* **74**, 2416–2422.
33. S. P. F. M. Roefs and K. G. Dekruif (1994) *Eur. J. Biochem.* **226**, 883–889.
34. N. Azuma and K. Yamauchi (1991) *Comp. Biochem. Physiol.* **99B**, 917–921.
35. P. J. Kraulis (1991) *J. Appl. Cryst.* **24**, 946–950.

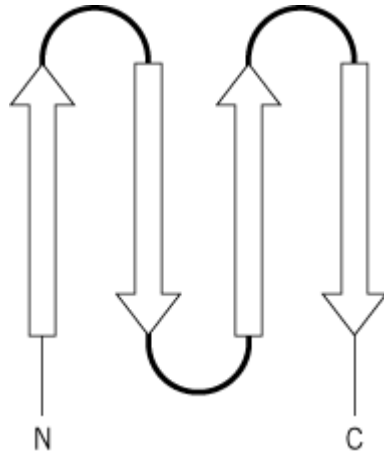
### Suggestions for Further Reading

36. L. Banaszak, N. Winter, Zh. Xu, D. A. Bernlohr, S. Cowan, and T. A. Jones (1994) Lipid-binding proteins—a family of fatty-acid and retinol transport proteins, *Adv. Protein Chem.* **45**, 89–151.
37. L. K. Creamer (1995) Effect of sodium dodecyl sulphate and palmitic acid on the equilibrium unfolding of bovine  $\beta$ -lactoglobulin *Biochemistry* **34**, 7170–7176.
38. P. F. Fox (1995) *Heat Induced Changes in Milk*, 2nd ed., International Dairy Federation, Brussels.
39. T.-R. Kim et al. (1997) High level expression of bovine  $\beta$ -lactoglobulin in *P. pastoris* and characterization of its physical properties, *Prot. Eng.* **10**, 133.
40. H. A. McKenzie (1971) " $\beta$ -Lactoglobulins". In *Milk Proteins—II*. Academic Press, New York. pp. 257–330.
41. B. Y. Qin et al. (1998) Structural basis of the Tanford transition of bovine  $\beta$ -lactoglobulin, *Biochemistry*. **37**, in press.
42. R. Townend, T. T. Herskovits, S. N. Timasheff, and M. J. Gorbunoff (1969) The state of amino acid residues in  $\beta$ -lactoglobulin *Arch. Biochem. Biophys.* **129**, 567–580.

### Beta-Meander

The  $\beta$ -meander is a [protein motif](#) frequently observed in [protein structures](#). It has a particular type of antiparallel  [\$\beta\$ -sheet](#) structure, with a very simple topology in which two or more  [\$\beta\$ -strands](#) that are consecutive in sequence are also adjacent to one another in the three-dimensional structure (Fig. [1](#)).  $\beta$ -Meanders are equivalent to multiply linked **hairpins**.

**Figure 1.** Schematic representation of a b-meander consisting of four b-strands (shown as arrows).



[See also [Beta-Sheet](#).]

### Beta-Pleated Sheet

b-Pleated sheet is an alternative name for [b-sheet](#), a type of **secondary structure** occurring in [protein structures](#). The name arises from the pleated nature of the sheet and the individual [b-strands](#) when viewed from the side.

[See also [Beta-Sheet](#).]

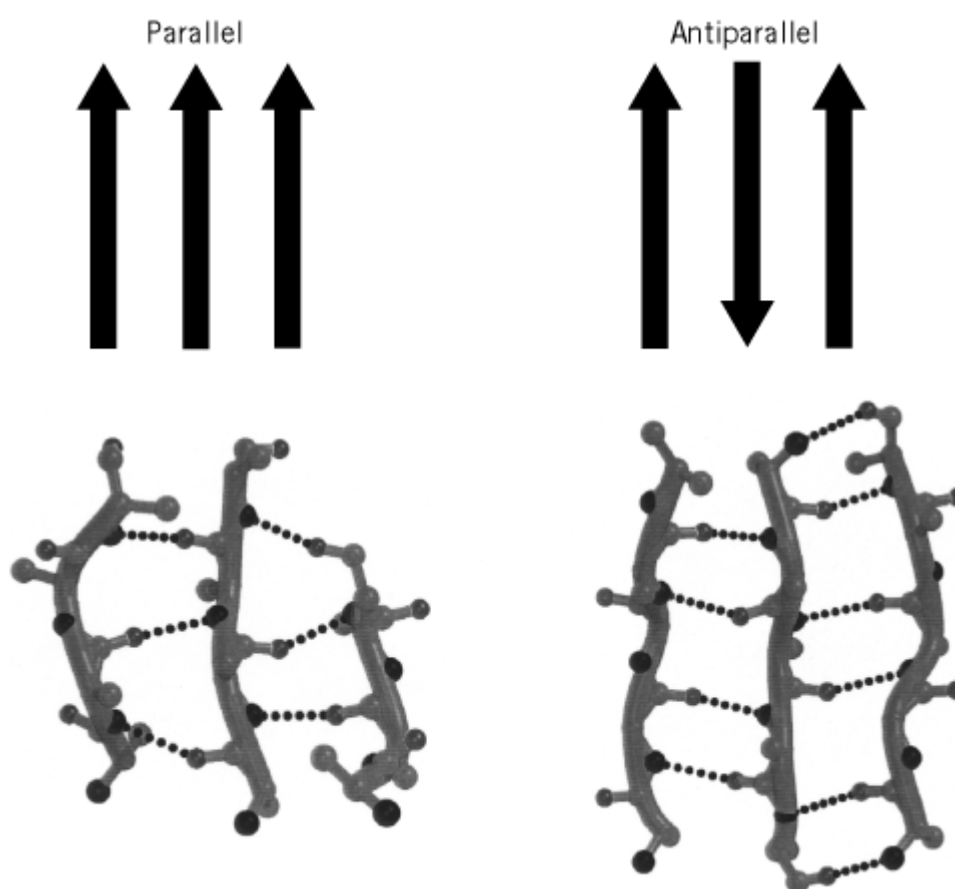
### Beta-Sheet

The b-sheet is a **secondary structure** component of [protein structure](#), comprising two or more adjacent [b-strands](#) linked by [hydrogen bonds](#). In contrast to the coiled [backbone](#) of the **a-helix**, the structure of a b-strand is characterized by an extended backbone [conformation](#), corresponding to the favorable upper left region of the [Ramachandran Plot](#). b-Sheets are formed by a side-by-side arrangement of b-strands, with backbone hydrogen bonds between adjacent strands linking the sheet together. The side chains of consecutive residues on a b-strand do not interact with each other, but lie alternately above and below the sheet, giving a pleated appearance to the b-sheet backbone (the b-sheet is sometimes referred to as a [b-pleated sheet](#)).

b-Strands are often depicted as arrows in protein structure diagrams, with the direction of the arrow running from the [N-terminal](#) end to the [C-terminal](#) end of the b-strand. When packed together in a b-sheet, the b-strands may all be oriented in the same direction (parallel b-sheet) or in alternating directions (antiparallel b-sheet). In some cases, the b-strands will be a mixture of parallel and antiparallel (mixed b-sheet). The backbone hydrogen bonding patterns of parallel and antiparallel b-

sheet structures differ (Fig. 1). In the parallel b-sheet, hydrogen bonds are evenly spaced along the b-strand and point at an angle to the other strand. The antiparallel b-sheet has alternately wide and narrow spacings between adjacent hydrogen bonds. The average backbone angles for parallel and antiparallel b-sheets also vary, with  $\phi$  of  $-119^\circ$  and  $\psi$  of  $113^\circ$  for parallel structures and  $\phi$  of  $-139^\circ$  and  $\psi$  of  $135^\circ$  for antiparallel structures (see [Ramachandran Plot](#)). Most b-sheets have a right-handed twist of about  $10^\circ$ . The regular hydrogen bonding pattern and backbone angles of b-sheets are sometimes disrupted at the edges by the formation of a [b-bulge](#), accentuating the twist to  $35^\circ$  to  $45^\circ$ .

**Figure 1.** b-Sheet structure in proteins. Schematic representation of parallel and antiparallel b-sheets, with each strand of the sheet depicted as an arrow having direction from the *N*-terminal to the *C*-terminal end. In the lower half of the figure, the atomic detail of the b-sheet structure is shown, with hydrogen bonds between backbone carbonyl oxygen atoms and amide nitrogen atoms indicated by dotted lines. Note the approximately equal spacing of hydrogen bonds in the parallel b-sheet compared with the alternating wide and narrow spacing in the antiparallel b-sheet. For clarity, only C $\beta$  atoms of side chains are shown. Nitrogen atoms and oxygen atoms are shown as dark spheres. This figure was generated using Molscript (1) and Raster3D (2, 3).



Because a b-sheet can be formed from sections of polypeptide chain that are distant from each other in the primary structure, many different types of b-sheet have been identified in protein structures. The [b-meander](#) is a simple antiparallel b-sheet topology where b-strands that are consecutive in sequence are also located adjacent to one another in the three-dimensional structure. The adjacent b-strands in a b-meander are often connected by **hairpin** turns. Other types of antiparallel b-sheet structure include the [jelly roll motif](#), the [Greek key motif](#), and the [b-barrel](#). The [b-helix](#) is formed from parallel b-sheets.

[See also [Protein Structure](#) and [Secondary Structure, Protein](#).]

## Bibliography

1. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
2. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
3. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

## Suggestions for Further Reading

4. C. Branden and J. Tooze (1991) *Introduction to protein structure*, Garland, New York.
5. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, New York.
6. J. S. Richardson (1981) The anatomy and taxonomy of protein structure. *Adv. Protein Chem.* **34**, 167–339.

## Beta-Strand

A b-strand is one of the two regular types of **secondary structure** (the other being the [a-helix](#)) that is commonly observed in [protein structures](#). The b-strand is characterized by an extended [backbone](#) conformation. The b-strand is not stable as a discrete structural unit by itself, but is stable only when it associates through [hydrogen bond](#) interactions with adjacent b-strands to form a [b-sheet](#).

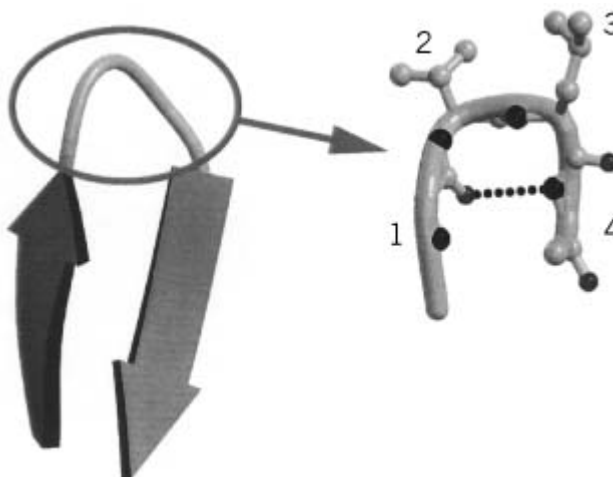
[See also [Beta-Sheet](#) and [Secondary Structure, Protein](#).]

## Beta-Turns

A b-turn is a type of nonregular **secondary structure** in [proteins](#) that causes a change in direction of the [polypeptide chain](#). In the b-turn, a [hydrogen bond](#) is formed between the [backbone](#) carbonyl oxygen of one residue ( $i$ ) and the backbone amide NH of the residue three positions further along the chain ( $i + 3$ ) (Fig. [1](#)) ([1](#)). This ( $i$ ) to ( $i + 3$ ) interaction distinguishes b-turns from [g-turns](#), which have an ( $i$ ) to ( $i + 2$ ) hydrogen bond. There are several different types of b-turns, classified according to the backbone dihedral angles of the two intervening residues, ( $i + 1$ ) and ( $i + 2$ ). The type I and II turns are the most common, together accounting for about two-thirds of all b-turns, but they have very different amino acid preferences ([2](#)). Type I turns prefer Asn, Asp, or Ser in the ( $i$ ) position; Asp, Ser, Thr or Pro in ( $i + 1$ ); Asp, Ser, Asn, or Arg in ( $i + 2$ ); and Gly, Trp, or Met in ( $i + 3$ ). In contrast, type II turns have a preference for Pro at position ( $i + 1$ ), Asn or Gly at position ( $i + 2$ ) and Gln or Arg at ( $i + 3$ ) ([2](#)). Other b-turn types include type I' (related to type I by inversion of the sign of the  $\phi$  and  $\psi$  angles of the [Ramachandran Plot](#)), II' (related to type II in the same way), IV, VIa, VIb, and VIII. Types I' and II' are almost exclusively found in **hairpins**; that is, they usually connect two adjacent antiparallel [b-strands](#). The type IV b-turn is a miscellaneous type that includes any ( $i$ ) to ( $i + 3$ ) hydrogen-bonded turn in a protein structure where the  $\phi$  and  $\psi$  angles are different (by  $>40^\circ$ ) from the values that define the other b-turn types (Table [1](#)). Types VIa and VIb have a proline *cis*-peptide bond at the ( $i + 2$ ) position; in type VIa the proline adopts a backbone conformation close to that of **a-helical** residues, whereas in type VIb the proline has a conformation similar to that of

residues in a [b-strand](#).

**Figure 1.** Schematic representations of a b-turn in a protein structure. **(Left)** The b-turn is shown connecting two b-strands in a hairpin motif. **(Right)** The detailed atomic structure of the b-turn is shown. Residues are numbered 1 to 4 for (*i*) to (*i* + 3). The hydrogen bond between the backbone carbonyl oxygen of residue 1 and the backbone amide nitrogen of residue 4 is shown as a dotted line. This example is a type I b-turn, because the backbone angles for the (*i* + 1) residue are  $\phi = -63^\circ$  and  $\psi = -33^\circ$ ; and for the (*i* + 2) residue,  $\phi$  is  $-86^\circ$  and  $\psi = -3^\circ$ . Oxygen atoms and nitrogen atoms are shown as dark spheres. This figure was generated by Molscript (3) and Raster3D (4, 5).



**Table 1. Classification of Beta-Turn Types**

| Beta-Turn Type | Backbone Dihedral Angles of Central Residues  |              |              |              |
|----------------|---|--------------|--------------|--------------|
|                | $\phi_{i+1}$  | $\psi_{i+1}$ | $\phi_{i+2}$ | $\psi_{i+2}$ |
| I              | $-60^\circ$   | $-30^\circ$  | $-90^\circ$  | $0^\circ$    |
| I'             | $60^\circ$  | $30^\circ$   | $90^\circ$   | $0^\circ$    |
| II             | $-60^\circ$   | $120^\circ$  | $80^\circ$   | $0^\circ$    |
| II'            | $60^\circ$  | $-120^\circ$ | $-80^\circ$  | $0^\circ$    |
| IV             | Any ( <i>i</i> ) to ( <i>i</i> + 3) hydrogen-bonded turn having angles that differ by more than $40^\circ$ from those of other b-turn types |              |              |              |
| VIa            | $-60^\circ$   | $120^\circ$  | $-90^\circ$  | $0^\circ$    |
| VIb            | $-120^\circ$  | $120^\circ$  | $-60^\circ$  | $0^\circ$    |
| VIII           | $-60^\circ$   | $-30^\circ$  | $-120^\circ$ | $120^\circ$  |

[See also [Secondary Structure, Protein](#), [Turns](#), [Gamma-Turns](#), and [Omega Loop](#).]

## Bibliography

1. J. S. Richardson (1981) *Adv. Protein Chem.* **34**, 167–339.
2. C. M. Wilmot and J. M. Thornton (1988) *J. Mol. Biol.* **203**, 221–232.
3. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
4. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
5. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

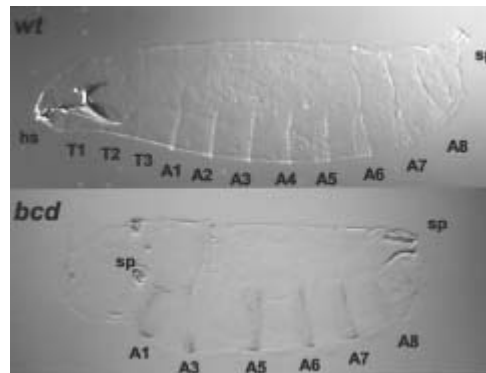
## Bicoid

Mutational analysis in the fruitfly *Drosophila melanogaster* has demonstrated that many aspects of anteroposterior patterning are already programmed in the oocyte by maternal gene activity (1, 2). During embryonic development, cells within an embryonic field receive positional information and are instructed by the concentration gradient of [morphogens](#) (3). Historically, the graded distribution of the Bicoid protein (Bcd) has provided the first evidence for the existence of such a morphogenetic gradient. Maternally derived *bcd* [messenger RNA](#) is localized to the anterior pole of the oocyte, where it serves as the source for the Bcd concentration gradient. Then, Bcd acts as a concentration-dependent activator of [transcription](#) and as a **translational repressor**. However, Bcd is only part of a wider regulatory network of maternal determinants, and its presence might actually be restricted to dipteran flies.

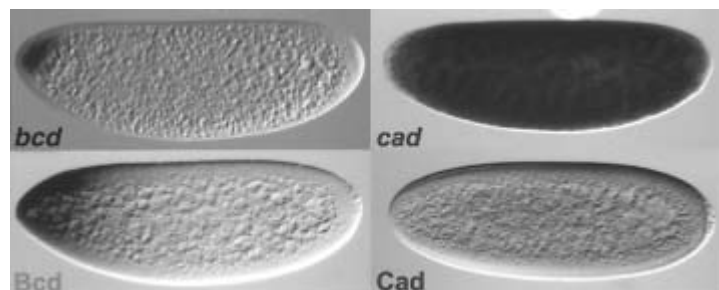
### 0.1. Anterior Localization of the Maternal Determinant *Bicoid*

Evidence for an anterior organizing center has been obtained by experimental embryology, for example, the removal and transplantation of polar cytoplasm (4) (see [Organizer](#)). The nature of the patterning agents has been identified as RNA-containing particles (5). Genetic experiments have shown that the anterior morphogenetic function requires the activity of the *bicoid* maternal gene (*bcd*). Mothers homozygous mutant for *bcd* are normal, but all of their progeny lack head, thorax, and some abdominal structures. Instead posterior terminal structures are duplicated at the anterior (Fig. 1; (6)). [Cloning](#) of the *bcd* gene revealed that *bcd* is expressed during oogenesis and that its mRNA is localized to the anterior tip of the oocyte and early embryo (Fig. 2) (7). This localization depends on the presence of the *bcd* 3'-UTR, which contains a specific structure recognized by factors that carry it to the anterior pole along the polarized [microtubule](#) array of the oocyte (8).

**Figure 1.** Phenotype of *bcd* embryos. Mothers homozygous mutant for *bcd* lay embryos that develop cuticles lacking all head structures (hs), thorax (T1, T2, T3), and two abdominal segments (A2, A4). Furthermore, the posterior structures like the spiracles (sp) are duplicated at the anterior.



**Figure 2.** mRNA and protein expression of *bcd* and *cad*. The *bcd* mRNA is localized at the most anterior tip of the embryo. Its translation and diffusion of its protein product lead to the formation of the Bcd morphogenetic gradient. Maternal *cad* mRNA is distributed throughout the embryo, but its translation is blocked by the Bcd protein. This leads to the formation of a posterior to anterior gradient of the Cad protein.



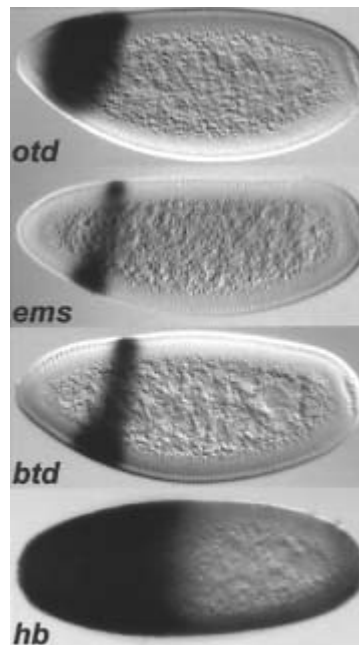
## 0.2. The Morphogenetic Gradient of the Bicoid Transcription Factor

After egg deposition, the anteriorly localized *bcd* mRNA starts to be **translated**, producing a source of the Bcd protein. Then, Bcd **diffuses** passively in the syncytial environment of the blastoderm embryo and generates an anteroposterior concentration gradient whose high point is at the anterior pole of the embryo and whose low point is near the posterior pole (Fig. 2; (9, 10)). Decreasing or increasing the copy number of the *bcd* gene changes the slope of the Bcd gradient and, consequently, changes the [fate map](#) of the early embryo. This suggests that anterior positional values are specified by Bcd in a concentration-dependent manner (9, 10).

The Bcd protein contains a homeodomain (7) that does not belong to any specific class (see [Homeobox Genes](#)). It is characterized by a **lysine** residue at residue 50 (Lys50) of the homeodomain that determines its **DNA-binding** specificity (14). Bcd functions as a transcriptional activator in cell culture and in yeast (11-13). The morphogenetic gradient of Bcd differentially activates the first zygotically active segmentation genes, the gap genes. At high concentrations of Bcd, the head gap genes *orthodenticle* (*otd*), *empty spiracles* (*ems*), and *buttonhead* (*btd*) are activated (15-18), whereas lower levels of Bcd are required to activate *hunchback* (*hb*), which is expressed in a broad anterior domain (Fig. 3; (13)). A detailed study of the proximal *hb* **promoter** suggested that the affinity of Bcd binding sites in the promoters of the different target genes determines at which threshold level, and therefore at which position along the anterior-posterior axis, these genes are activated (19, 20). The ability of Bcd to activate target genes differentially provided the first molecular explanation of how a morphogen might act.

**Figure 3.** Expression pattern of *bcd* target genes. *bcd* is a maternal gene that controls the expression of the zygotic gap

target genes. The morphogenetic gradient of Bcd protein differentially activates head (*otd*, *ems*, and *btd*) or thoracic (*hb*) gap genes.



Although Bcd is a DNA-binding protein, it also controls translation of the *caudal* gene (*cad*). *cad* is a gene required for posterior development, and its function is evolutionarily conserved in diverse organisms, including vertebrates (21). Maternal *cad* mRNA is ubiquitously distributed (22-24). Its translation is blocked at the anterior of the embryo by *bcd*, leading to a posterior-to-anterior gradient of the Cad protein (Fig. 3). *cad* is expressed ectopically in *bcd* mutant embryos throughout the embryo. Bcd binds to a sequence in the 3'-untranslated region of the *cad* mRNA (25-27). Like DNA binding, Bcd's RNA binding is conferred by the homeodomain and depends on the presence of Lys50 in the homeodomain. The Bcd homeodomain is the only one that binds RNA (25, 27).

#### 0.4. Bicoid is Integrated into a Network of Maternal Determinants

Bcd is not the only maternal morphogen involved in anteroposterior patterning. A morphogenetic gradient of the Hunchback protein (Hb) creates the correct polarity in the embryo in the absence of *bcd*. The anteroposterior Hb gradient is established by translational repression of ubiquitous maternal mRNA by the *nanos* posterior determinant (28-30). Bcd and Hb are redundant for abdominal segmentation because high levels of maternal Hb rescue the abdominal phenotype of *bcd* mutant embryos (28, 29) and the activation of the central gap gene *Krüppel*. At the anterior of the embryo, Bcd and Hb synergize to activate the anterior gap genes (31). Because each morphogen has instructive functions on its own, both must be removed from the embryo to disrupt the polarity of the embryo completely. In more posterior regions, Bcd is also redundant with *cad* for activating the gap gene *knirps* (32). This shows that *bcd* function is partially redundant with other maternal functions. It should be noted that maternal *hb* and *cad* functions are not essential for viability in the presence of *bcd*. It is likely that they represent remains of ancestral systems that patterned the embryo before the appearance of *bcd*.

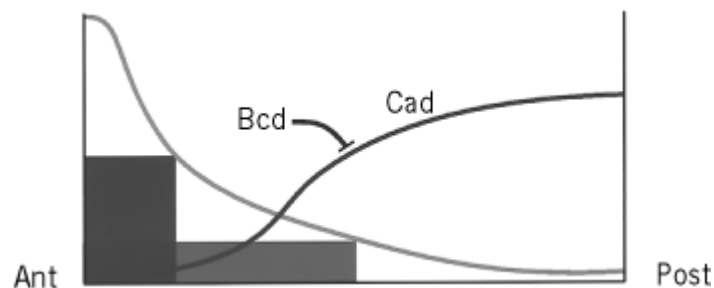
At the anterior pole of the embryo, Bcd function is influenced by the activity of the terminal system, which involves locally activating the Torso (Tor) **receptor tyrosine kinase** (33). Tor activity leads to two opposite effects on Bcd function. Tor activity increases the activator function of Bcd by lowering the threshold level of Bcd where target genes are activated (18, 34, 35). At the anterior pole, however, Bcd function is blocked in a Tor-dependent manner, and most *bcd* target genes retract from the anterior pole where Bcd concentration is the highest (36). The effect of *tor* activity on Bcd function might be mediated by *tor*-dependent phosphorylation of Bcd (36), or it may be indirect and may involve other factors (34).



### 0.5. Evolutionary Considerations on Bicoid

There are several reasons to believe that the critical functions of *bcd* in patterning the anteroposterior axis of the *Drosophila* embryo (Fig. 4) may not be conserved throughout evolution. *bcd* is unlikely to represent an ancestral morphogen. No *bcd* homologous genes have been identified outside higher dipterans (38), despite its homeobox and its specific location in the Hox cluster (7). Moreover, *bcd* shows an unusually high rate of **divergence** for a homeobox gene (39), and its function is not even conserved within higher dipterans (38, 40). Therefore, it is reasonable to assume that *bcd* is the result of recent **gene duplication** in the Hox cluster of dipterans and that it evolves relatively freely. Therefore, although the role of Bcd in *Drosophila* represents a striking example of a powerful morphogenetic **transcriptional factor**, it is likely that this function is not conserved in other insects and may exist only to allow the *Drosophila* embryo to develop quickly.

**Figure 4.** Morphogenetic functions of the Bcd gradient. The Bcd gradient differentially activates zygotic target genes (dark boxes) in head and thorax anlagen. The Bcd gradient also leads to the formation of the Cad protein gradient by blocking translation of its mRNA.



### Bibliography

1. C. Nüsslein-Volhard and E. Wieschaus (1980) *Nature* **287**, 795–801.
2. C. Nüsslein-Volhard, H. G. Frohnhof, and R. Lehmann (1987) *Science* **238**, 1675–1681.
3. L. Wolpert (1969) *J. Theor. Biol.* **25**, 1–47.
4. K. Sander (1969) *Advances in Insect Physiology* (J. E. Treherne, M. J. Berridge, and V. B. Wigglesworth, eds.), Academic Press. Vol. **12**, pp. 125–235.
5. K. Kalthoff (1979) *Determinants of Spatial Organization* (S. Subtelney and I. R. Konisberg, eds.), Academic Press, N.Y., pp. 97–126.
6. H. G. Fröhnhof and C. Nüsslein-Volhard (1986) *Nature* **324**, 120–125.
7. T. Berleth, M. Burri, G. Thoma, D. Bopp, S. Richstein, G. Frigerio, M. Noll, and C. Nüsslein-Volhard (1988) *EMBO J.* **7**, 1749–1756.
8. D. Ferrandon, L. Elphick, C. Nüsslein-Volhard, and D. St Johnston (1994) *Cell* **79**, 1221–1232.
9. W. Driever and C. Nüsslein-Volhard (1988) *Cell* **54**, 95–104.
10. W. Driever and C. Nüsslein-Volhard (1988) *Cell* **54**, 83–93.
11. S. D. Hanes and R. Brent (1989) *Cell* **57**, 1275–1283.
12. W. Driever, J. Ma, C. Nüsslein-Volhard, and M. Ptashne (1989) *Nature* **342**, 149–154.
13. G. Struhl, K. Struhl, and P. M. Macdonald (1989) *Cell* **57**, 1259–1273.
14. J. Treisman and C. Desplan (1989) *Nature* **341**, 335–337.
15. D. Dalton, R. Chadwick, and W. McGinnis (1989) *Genes and Devel.* **3**, 1940–1956.
16. R. Finkelstein and N. Perrimon (1990) *Nature* **346**, 485–488.
17. U. Walldorf and W. J. Gehring (1992) *EMBO J.* **11**, 2247–2259.

18. E. A. Wimmer, M. Simpson-Brose, S. M. Cohen, C. Desplan, and H. Jäckle (1995) *Mechanisms of Development* **53**, 235–245.
19. W. Driever, G. Thoma, and C. Nüsslein-Volhard (1989) *Nature* **340**, 363–367.
20. G. Struhl (1989) *Nature* **338**, 741–744.
21. V. Subramanian, B. I. Meyer, and P. Gruss (1995) *Cell* **83**, 641–653.
22. M. Mlodzik and W. J. Gehring (1987) *Cell* **48**, 465–478.
23. M. Mlodzik, A. Fjose, and W. J. Gehring (1985) *EMBO J.* **4**, 2961–2969.
24. P. M. MacDonald and G. Struhl (1986) *Nature* **324**, 537–545.
25. R. Rivera-Pomar, D. Niessing, U. Schmidt-Ott, W. J. Gehring, and H. Jäckle (1996) *Nature* **379**, 746–749.
26. J. Dubnau and G. Struhl (1996) *Nature* **379**, 694–699.
27. S. K. Chan and G. Struhl (1997) *Nature* **388**, 634.
28. G. Struhl, P. Johnston, and P. A. Lawrence (1992) *Cell* **69**, 237–249.
29. M. Hülskamp, C. Pfeifle, and D. Tautz (1990) *Nature* **346**, 577–580.
30. C. Schülz and D. Tautz (1995) *Development* **121**, 1023–1028.
31. M. Simpson-Brose, J. Treisman, and C. Desplan (1994) *Cell* **78**, 855–865.
32. R. Rivera-Pomar, X. Lu, N. Perrimon, H. Taubert, and H. Jackle (1995) *Nature* **376**, 253–256.
33. N. Perrimon and C. Desplan (1994) *Trends in Biochemical Sciences* **19**, 509–513.
34. Q. Gao, Y. Wang, and R. Finkelstein (1996) *Mechanisms of Development* **56**, 3–15.
35. U. Grossniklaus, K. M. Cadigan, and W. J. Gehring (1994) *Development* **120**, 3155–3171.
36. E. Ronchi, J. Treisman, N. Dostatni, G. Struhl, and C. Desplan (1993) *Cell* **74**, 347–355.
37. Y. Bellaïche, R. Bandyopadhyay, C. Desplan, and N. Dostatni (1996) *Development* **122**, 3499–3508.
38. R. Schöder and K. Sander (1993) *Roux's Arch. Dev. Biol.* **203**, 34–43.
39. R. Sommer and D. Tautz (1991) *Development* **113**, 419–430.
40. F. Bonneton, P. J. Shaw, C. Fazakerley, M. Shi, and G. A. Dover (1997) *Mechanisms of Development* **66**, 143–156.

### **Suggestions for Further Reading**

41. P. A. Lawrence (1992) *The Making of a Fly, The Genetics of Animal Design*, Blackwell Scientific Publications.
42. P. A. Lawrence and G. Struhl (1996) *Cell* **85**, 951–961.
43. D. St Johnston (1995) *Cell* **81**, 161–170.
44. D. St Johnston and C. Nüsslein-Volhard (1992) *Cell* **68**, 201–219.

## **Bifunctional Crosslinking Reagents**

### 1. Terminology

If a bifunctional chemical [crosslinking](#) agent is designed so that its two reactive groups are identical,

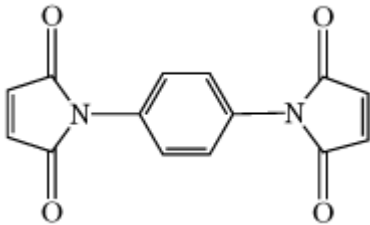
it is referred to as a *homobifunctional* reagent; if its two reactive groups are different, it is a *heterobifunctional* reagent. If one or both reactive groups of the crosslinker become so only as the result of a photochemical reaction caused by exposing the reagent to light of an appropriate wavelength, then the crosslinking reagent is photoactivatable, photoreactive, photosensitive, or **light-activated (caged)**. The portion of the crosslinking reagent other than, and usually between, its two reactive groups is referred to as the *spacer arm*, or *crossbridge*. If the spacer arm contains a covalent bond that can be easily broken (by an oxidant, reductant, base, etc), then the crosslinker is termed *cleavable*; otherwise, it is classified as *noncleavable*. With most chemical crosslinkers, a portion of the reagent is incorporated into the final crosslinked complex; however, there is a special class of reagents that facilitate crosslinking without being incorporated into the final covalent complex. These reagents are referred to as *zero-length* crosslinkers, because the groups they crosslink are not separated by a spacer arm.

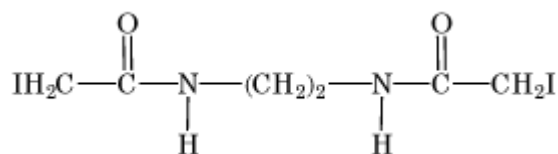
## 2. Variables in the Design of Crosslinking Reagents

By varying the chemical and physical properties of the reactive groups and spacer arm of a crosslinking agent, one can obtain crosslinkers with different properties and uses. A heterobifunctional reagent, with its two different types of reactive groups, allows for different selectivities in the functional groups that become crosslinked, compared to a homobifunctional reagent containing a pair of the same reactive groups. The selectivity shown by the reactive groups of crosslinkers for the side chains of particular amino acids mirrors the selectivity of those reactive groups when used as monofunctional reagents in the chemical modification of proteins, from which crosslinking is an outgrowth.

A number of variables pertain to spacer arms: (1) the length of the crossbridge, as in the case of bisimidoesters (Table 1) with different numbers of methylene groups in the spacer arm; (2) its geometry, which in the case of the phenylenebismaleimides (Table 1) also causes an increase in length, progressing from *ortho* to *para*; (3) its chemical nature, especially the number of charged or polar groups on the spacer arm, which influences the crosslinker's solubility properties and uses. For example, [membrane proteins](#) would more likely be crosslinked by bifunctional reagents that are **hydrophobic**, rather than **hydrophilic**. [Reporter groups](#) (chromophores, radioactive tracers, etc) can also be incorporated into spacer arms to facilitate quantification of the extent of protein modification, studies on the crosslinked complex, or identification of the crosslinked components. Incorporation into the spacer arm of cleavable bonds, such as [disulfide bonds](#), is also useful in the identification of crosslinked components, especially by [diagonal methods](#) (1).

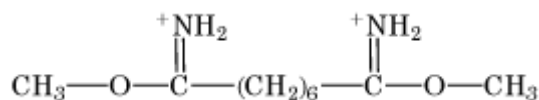
**Table 1. Examples of Bifunctional Crosslinkers**

| Compound   | Structure  | Class            |
|--|--|------------------|
| Homobifunctional                                 |  |                  |
| 1. <i>N,N'</i> - <i>p</i> -Phenylenebismaleimide |  | Bismaleimide     |
| 2. <i>N,N'</i> -Ethylene-bis(iodoacetamide)      |  | Bishaloacetamide |



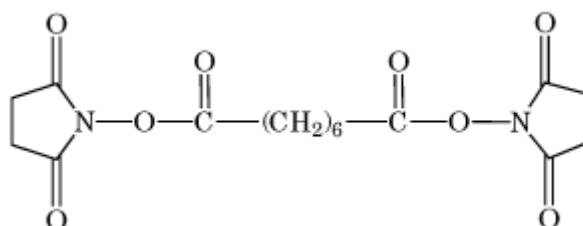
3. Dimethylsuberimidate

Bisimidoester



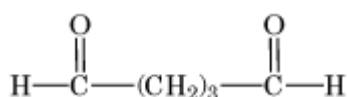
4. *N*-Hydroxysuccinimidylsuberate

*N*-Hydroxysuccinimide



5. Glutaraldehyde

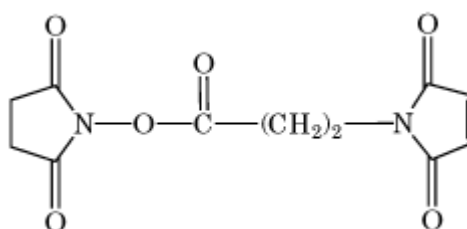
Dialdehyde



Heterobifunctional

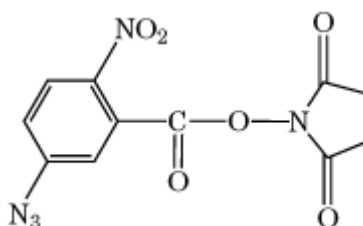
6. *N*-Succinimidyl-3-maleimidopropionate

*N*-Hydroxysuccinimide



7. *N*-5-Azido-2-nitrobenzoyloxysuccinimide

Arylazide/*N*-hydroxy ester



### 3. Examples of Bifunctional Crosslinkers

The compounds listed in Table 1 are representative of some of the most frequently used chemical classes of crosslinking reagents, but they represent only a small fraction of the literally hundreds of bifunctional crosslinking reagents that have been described in the scientific literature. The widely used *N*-substituted bismaleimides are selective for [thiol groups](#), especially at pH values near neutrality, where [amino groups](#) are predominantly protonated. This class of crosslinkers is available with many different spacer arms, having large variations in length and solubility, or having chromophores or cleavable bonds incorporated. The alkylating bishaloacetyl crosslinkers are also selective at neutral pH for sulfhydryl groups, but they will also react with imidazole side chains of some **histidine** residues under this condition. It should be remembered that few reagents used to

modify proteins chemically are absolutely specific for any given amino acid side chain; on the other hand, high selectivity can often be achieved empirically through varying the pH, reaction time, and concentrations of reactants.

The bisimidoesters were among the earliest used crosslinkers and are also available with a variety of spacer arms. These compounds preferentially react at alkaline pH with amino groups, with which they form amidine bonds, retaining the positive charge of the original amine. A more frequently used class of amine-selective reagents is the [N-hydroxysuccinimide](#) esters, which produce an amide-crosslinked product that is relatively stable in aqueous solution. The *N*-hydroxysuccinimide ester group is frequently combined with maleimide, such as in compound 6 of Table 1, to produce heterobifunctional crosslinkers with selectivity toward crosslinking amino and sulfhydryl groups. Another common class of heterobifunctional crosslinkers contains at one end a photoactivatable group, which is generally assumed to be nonselective for the functional groups it modifies. For this class of reagents, the two-step reaction protocol is particularly useful, with the chemical reaction being carried out first in the absence of activating radiation, followed by irradiation of the photosensitive group to bring about crosslinking. The arylazide shown in Table 1 (compound 7) is an example of a photosensitive crosslinker that is widely used. The bis-homobifunctional aldehyde [glutaraldehyde](#) (Table 1, compound 5) is perhaps the most frequently used chemical crosslinker.

### Bibliography

1. J. A. Cover, J. M. Lambert, C. M. Norman, and R. R. Traut (1981) *Biochemistry* **20**, 2843–2852.

### Suggestions for Further Reading

2. S. S. Wong (1993) *Chemistry of Protein Conjugation and Cross-Linking*, CRC Press, Boca Raton, FL. (An outstanding book describing all aspects of the chemistry and uses of bifunctional reagents in crosslinking and conjugation, with an extensive bibliography for each chapter.)
3. Annual products catalog of Pierce Chemical Co., Rockford, IL. [Contains structures, helpful information, and specific uses (with references) for a large number of bifunctional reagents.]

## Binding

Total binding  $B_L$  is the number of ligand molecules that at any moment are in contact with (occupy sites on) the [accessible surface](#) of a [protein](#) molecule. It is related to [preferential binding](#) (as measured by [equilibrium dialysis](#)) by

$$B_L = \left( \frac{\partial m_L}{\partial m_{pr}} \right)_{T, \mu_w, \mu_L} + \frac{m_L}{m_w} B_w \quad (1)$$

where  $B_L$  and  $B_w$  are the numbers of ligand and water molecules, respectively, that are in contact with the protein surface.

When the interactions are weak, the ligand must be used at high concentration (eg, 1 *M* sucrose, 3 *M* glycerol, 8 *M* [urea](#)). As a consequence, the total binding is *not* given by equilibrium dialysis, as the last term of Eq. (1) (binding of water) becomes significant. Total binding can be measured by nonthermodynamic techniques that respond to contacts between ligand molecules and protein [(1)]. Such techniques include **calorimetric** titration, which detects the heat of protein–ligand contact, or spectroscopic techniques (eg, **fluorescence** or **uv absorbance**) that detect spectral perturbations each

time a ligand–protein contact is made. If the total [hydration](#) is known from other techniques, then total binding can be derived from dialysis equilibrium results by the application of Eq. (1).

## Bibliography

1. S. N. Timasheff (1995) In *Protein-Solvent Interactions* (R. B. Gregory, ed.), Marcel Dekker, New York, Chap. "11".

## Bioinformatics

Bioinformatics is a new branch of molecular biology, also known as computational biology. Computational physics and computational chemistry have emerged as third branches in their respective disciplines, after the experimental and theoretical branches, due largely to the advance of computational capabilities in modern computers. In contrast, bioinformatics has emerged in molecular biology as the result of advances in experimental technology, especially in the high-throughput [DNA sequencing](#), which is generating a vast amount of **gene** and [protein](#) sequence data. Since initiation of the Human Genome Project in the late 1980s, bioinformatics has become an integral part of the coordinated efforts to sequence the entire [genomes](#) of a number of organisms from **bacteria** to human. Consequently, bioinformatics is a data-driven discipline requiring large-scale [databases](#) and associated technologies for data management and interpretation. This contrasts with the model-driven, theory-based approaches in computational physics and computational chemistry.

Bioinformatics covers a diverse range of topics, which broadly have two roles. First, bioinformatics is widely used in experimental projects, such as in determining the optimal **genetic mapping** strategy and how best to assemble raw sequence data, and in the storing and handling information and materials. Secondly and more important, the major task of bioinformatics is to develop new databases and new computational technologies that will help to understand the biological meaning encoded in the sequence data. Interpretation of sequence data cannot be achieved by numerical calculations based on first-principle equations, which are virtually nonexistent in biology, but it requires using and processing empirical knowledge acquired from experimental data. This is what human experts would do if the amount of data and knowledge were manageable in size. Thus, bioinformatics was born from the marriage of molecular biology, in the age of massive data production, with artificial intelligence, which is a mature branch of computer science, to automate knowledge processing.

Table 1 summarizes topics in bioinformatics for interpreting of sequence data. The methods of computer science have been used traditionally in artificial intelligence applications, such as speech recognition and natural language processing, but they are also quite effective in solving problems in molecular biology. There is a rough distinction between how data are organized and consequently what kinds of computational methods are used in the three categories of Table 1. In the first category of similarity searches, the database is a collection of all the known primary data, such as the [sequence database](#) and the [structure database](#), which are the repositories of all reported sequences and three-dimensional structures, respectively, of **nucleic acids** and [proteins](#). A basic operation in this category is comparing of individual sequences or structures to detect any similarity. For example, to understand the functional implications of a newly determined gene or protein sequence, it is customary to perform a sequence similarity, or [homology](#), search, comparing the query sequence with each of the sequences in the database. If any similar sequence is found in the database, the query sequence is assumed to have a similar or related function. This reasoning is based on the

empirical observation that homologous sequences generally share similar functions and similar 3-D structures, because they arose from a common evolutionary ancestor and because of functional constraints.

**Table 1. Bioinformatics for Interpretation of Sequence Data**

| <b>Problems in Biological Science</b> |  | <b>Methods in Computer Science</b>                |
|---------------------------------------|--|---|
| Similarity search                     | Pairwise sequence alignment            | Dynamic programming                               |
|                                       | Sequence similarity search             | Simulated annealing                               |
|                                       | Multiple sequence alignment            | Genetic algorithms                                |
|                                       | RNA secondary structure prediction     | Other optimization algorithms                     |
|                                       | 3-D protein structure alignment        |   |
| Structure/function prediction         | Motif extraction                       | Discriminant analysis                             |
|                                       | Functional site prediction             | Neural network                                    |
|                                       | Cellular localization prediction       | Hidden Markov model                               |
|                                       | Coding region prediction               | Formal grammar                                    |
|                                       | Transmembrane segment prediction       | Other pattern recognition and learning algorithms |
|                                       | Protein secondary structure prediction |   |
|                                       | 3-D protein structure prediction       |   |
|                                       | 3-D RNA structure prediction           |   |
| Molecular classification              | Superfamily classification             | Clustering algorithms                             |
|                                       | Fold classification                    |   |
|                                       | Ortholog/paralog grouping of genes     |   |
|                                       |  |   |

The procedure of comparing two sequences is a problem of sequence alignment, which is a major topic in [sequence analysis](#). Given a similarity measure of amino acids or nucleotides, the problem of obtaining the best sequence alignment is equivalent to optimizing a given score function that represents the overall similarity. There are variations in how the optimal alignment is made, either globally for the entire sequences or locally to detect localized regions, how many sequences are aligned at a time, and which similarity measure is to be used. By allowing insertions and deletions

([indels](#)), the number of possible alignments grows exponentially as the number and length of the sequences to be compared increases. Thus, the sequence alignment problem is a typical combinatorial optimization problem in computer science. Although the alignment of two or three sequences can be solved rigorously by using dynamic programming algorithms ([1](#)), which effectively evaluate all possibilities, multiple alignment of many sequences requires heuristic algorithms to obtain approximate solutions. Furthermore, to search large databases effectively for similar sequences, heuristic algorithms such as FASTA ([2](#)) and BLAST ([3](#)) have been developed for rapidly identifying regions of two sequences that are similar locally. Other extensions of the sequence alignment algorithms include analysis of a single RNA sequence to predict any **secondary structure** (see [RNA Structure Prediction](#)) and compare two 3-D [protein structures](#).

When humans try to interpret a sentence written in a foreign language, they use dictionaries and a knowledge of grammar. However, the process of using a homology search to interpret a sentence written in the DNA language is like comparing it with all the sentences written in the past and checking precedents as to how they were interpreted. Sequence interpretation would become more efficient if our knowledge of sequence-functional relationships were more advanced. The second and third categories in [Table 1](#) involve such reorganizations of the primary data. In the second category of structure/function prediction, the primary data are classified into a group based on a given criterion, such as functional identity, and higher level knowledge is abstracted from this group. This is based on our empirical observation that the function of a protein molecule is usually exerted at a functional site in its three-dimensional structure, which is usually a **ligand-binding** site for other molecules, and that the site is formed by one or more linear polypeptide segments that have conserved sequence patterns. These conserved patterns are called sequence motifs, and they can be extracted, for example, by multiple alignment of a set of sequences that have the same function. A motif search against a [library](#) of known motifs is widely used as an alternative method for interpreting sequence data, when no similar sequences have been found by a homology search. However, the bioinformatics problem here is not how to search the known motifs, but how best to extract and define new motifs effectively. Thus, the problem is related to machine learning in computer science, where a higher level concept or a generalized pattern is abstracted from a set of given examples.

Although the second category in [Table 1](#) is targeted at a specific function and a specific group of sequences, the third category of molecular classification represents a more global analysis, where the grouping is attempted with all the available data, and the relationship among groups, which is often hierarchical, is examined. For example, all known protein sequences can be hierarchically classified into families and [superfamilies](#) according to their level of sequence similarity. Because complete genome sequences are available, it has become possible to establish similarity relationships of all **genes** within species ([paralogous genes](#)) and across species ([orthologous genes](#)). These analyses are carried out by clustering algorithms, and the similarity score is used as a measure of the closeness of two sequences. In practice, however, there are biological difficulties, notably due to the abundance of **multidomain**, [multifunctional proteins](#), which make automatic clustering a complicated task.

Molecular biology has been dominated by a reductionist approach where, starting from a selected functional aspect of life, such as metabolism, [signal transduction](#), or [development](#), molecular components and their interactions are identified experimentally. In contrast, sequencing an entire genome identifies a complete set of genes and proteins in an organism, but it does not tell how they should be interrelated to form a functioning, living system. This is where a synthetic approach based on bioinformatics will play an increasingly more dominant role in the basic research into understanding life and in the applied research into biomedical relevance. [Table 2](#) summarizes the role of bioinformatics relative to the level of abstraction in nature. At the molecular level, an amino acid residue is represented by a symbol, for example, C for [cysteine](#), which is a compound consisting of carbon, hydrogen, nitrogen, oxygen, and sulfur atoms. When comparing two amino acid sequences, only the one-dimensional connection of symbols is considered biologically relevant, neglecting any atomic details. At a higher level of abstraction, which may be called a network level, a gene or a protein is represented by a symbol, for example, [Ras](#) for an **oncogene** product. A primary



concern here is the Ras signal transduction pathway, which is generated by a network of interacting molecules. Bioinformatics is relatively well established at the molecular level, such as in the sequence and structure databases and in the computational methods to analyze them. To step up to a higher level, however, it is necessary to develop new databases and new computational technologies, such as a network database containing all the wiring-diagram information about genes and molecules, plus the algorithms to analyze them. This direction would eventually lead to *in silico* reconstruction of a biological organism.

**Table 2. Level of Abstraction**

| Level           | Constituents                   | Entity                                | Informatics             |
|-----------------|--------------------------------|---------------------------------------|-------------------------|
| Subatomic level | Elementary particles           | Atom                                  | Computational physics   |
| Atomic level    | Atoms                          | Amino acid, etc.                      | Computational chemistry |
| Molecular level | Amino acids, nucleotides, etc. | Protein, gene, etc.                   | Bioinformatics          |
| Network level   | Proteins, genes, etc.          | Molecular pathway, molecular assembly |                         |

#### Bibliography

1. S. B. Needleman and C. D. Wunsch (1970) *J. Mol. Biol.* **48**, 443–453.
2. W. J. Wilbur and D. J. Lipman (1983) *Proc. Natl. Acad. Sci. USA* **80**, 726–730.
3. S. F. Altschul et al. (1990) *J. Mol. Biol.* **215**, 403–410.

#### Biopolymer

Biopolymers, sometimes referred to as biological [macromolecules](#), are [polymers](#) of biological importance such as [proteins](#), [nucleic acids](#), and carbohydrates. Like all polymers, biopolymers are formed by joining together many copies of repeating units to form long chains that have a repeating or constant [backbone](#) and variable [side chains](#). The repeating units of proteins are the 20 naturally occurring [amino acids](#) that form a [polypeptide chain](#). Nucleic acids are formed from [nucleotides](#) and have a sugar-phosphate backbone with [purine](#) or [pyrimidine](#) side chains.

[See also [Polymer](#) and [Macromolecule](#).]

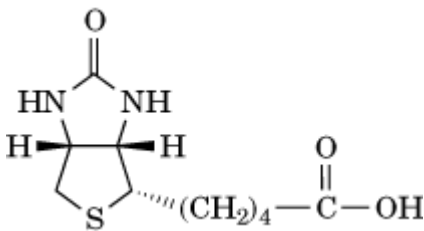
#### Suggestions for Further Reading

- A. G. Walton and J. Blackwell (1973) *Biopolymers*, Academic Press, New York.  
 P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

## Biotin

Biotin is a vitamin (vitamin H) and a [growth factor](#), first isolated by Kögl in 1935 (1). Its structure (Fig. 1) was established by du Vigneaud in 1942 (2, 3).

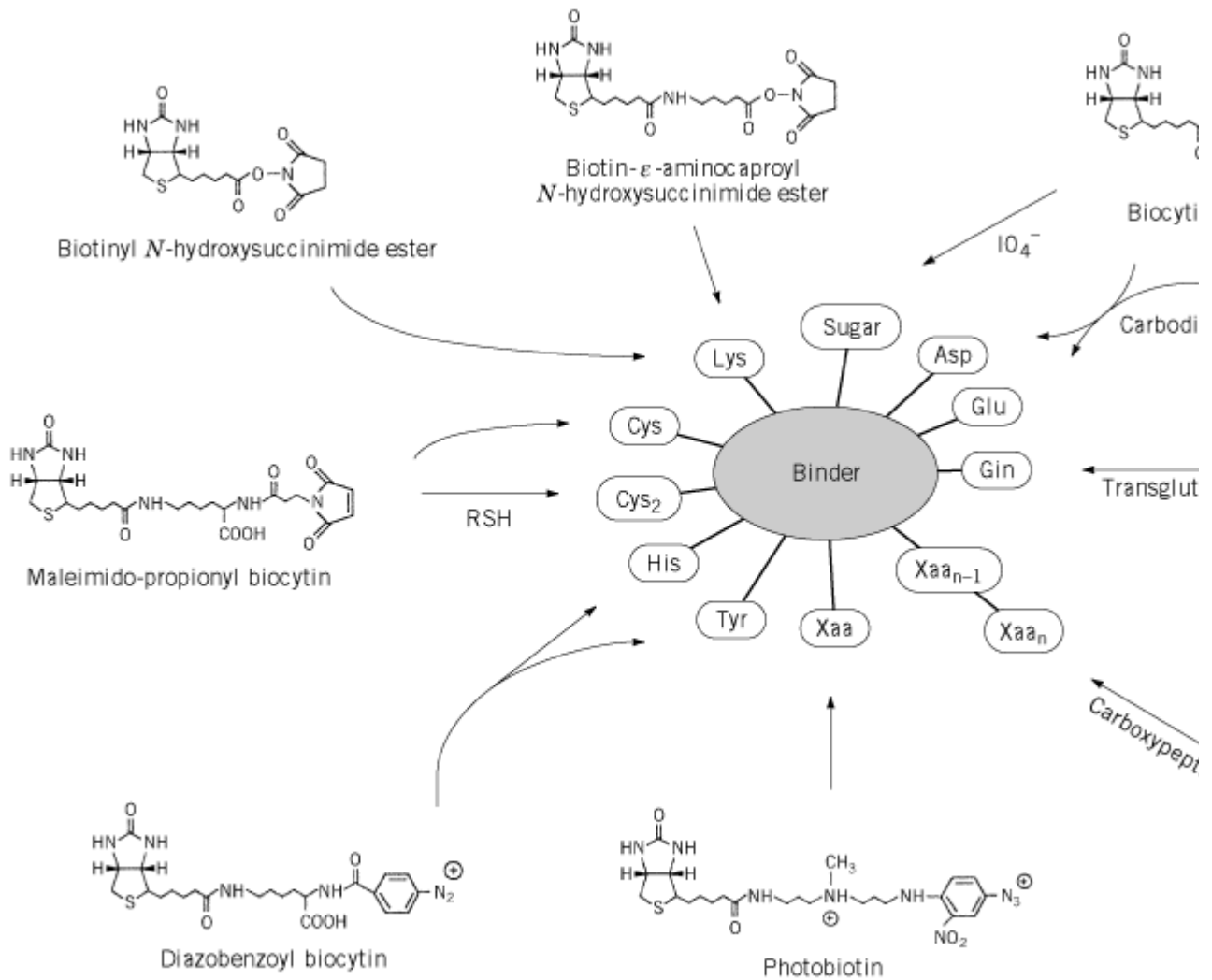
**Figure 1.** The structure of biotin.



Biotin is a prosthetic group for a family of “biotin-requiring” [enzymes](#) (carboxylases, decarboxylases and transcarboxylases) (4). These multisubunit enzymes contain two types of [active sites](#), one wherein a molecule containing a carboxylic acid group serves as a donor of the carboxyl group and the second wherein another molecule acts as an acceptor. The biotin moiety is attached through a [lysine](#) residue in a particular sequence (Ala-Met-Lys-Met) that is exposed on the surface of a special type of biotin-binding subunit. The biotin group participates as a transient carrier to which the carboxyl group is bound covalently in the process of being transferred from the donor to the acceptor molecule. The biotin moiety of the enzyme is readily available for binding to [avidin](#), which inactivates the enzyme. After serving their respective function, the biotin-requiring enzymes undergo degradation, whereby e-N-biotinyl lysine (termed biocytin) is released into the circulation. There the enzyme biotinidase hydrolyzes biocytin to regenerate biotin and lysine. A rare mutation in the gene in humans causes a deficiency in the level of this enzyme in the blood stream, which, in turn, causes a disease due to a deficiency of biotin (5).

Biotin attained true fame when it was recognized that the high-affinity interaction between avidin and biotin can be applied for a variety of biotechnological and biomedical purposes (as described in [Avidin-biotin system](#))(6, 7). The varied application of this system is usually achieved upon coupling biotin covalently to different biologically active binder molecules, [antibodies](#), [hormones](#), [DNA](#), etc. The binding of biotin to either avidin or [streptavidin](#) greatly enhances the stability of these proteins, such that the avidin–biotin or streptavidin–biotin complex is essentially irreversible. Because both proteins primarily recognize the ureido ring of the biotin molecule, the carboxylic acid side-chain can be modified almost at will, and the resultant derivatives can be used to incorporate biotin covalently into virtually any other molecule (Fig. 2).

**Figure 2.** Selected approaches using different biotin-containing reagents for biotinylating functional groups on proteins.



Both polyclonal and [monoclonal antibodies](#) have been produced against biotin. Comparing the sequence of the V<sub>H</sub> region of a monoclonal anti-biotin preparation (8) with the sequences of avidin and streptavidin revealed an astonishing similarity in their binding sites, indicating that such functional elements are conserved in nature for a given purpose, irrespective of the protein type.

### Bibliography

1. F. Kögl and B. Z. Tönnes (1936) *Physiol. Chem.* **242**, 43–73.
2. V. du Vigneaud, D. B. Melville, K. Folkers, D. E. Wolf, R. Mazingo, J. C. Kereszteolf, and S. A. Harris (1942) *J. Biol. Chem.* **146**, 475–485.
3. V. du Vigneaud, K. Hoffman, and D. B. Melville (1942) *J. Amer. Chem. Soc.* **64**, 188–189.
4. J. Moss and M. D. Lane (1971) *Adv. Enzymol.* **35**, 321–372.
5. B. Wolf, G. S. Heard, J. R. McVoy, and R. E. Grier (1985) *Ann. N.Y. Acad. Sci.* **447**, 252–62.
6. E. A. Bayer and M. Wilchek (1978) *Trends Biochem. Sci.* **3**, N237–N239.
7. E. A. Bayer and M. Wilchek (1980) *Methods Biochem. Anal.* **26**, 1–45.
8. H. Bagci, F. Kohen, U. Kuscuoglu, E. A. Bayer, and M. Wilchek (1993) *FEBS Lett.* **322**, 47–50.

## Suggestion for Further Reading

9. M. Wilchek and E. A. Bayer, eds. (1990) *Avidin-Biotin Technology, Methods Enzymol.* Vol. **184**, Academic Press, San Diego.

## Biotin Repressor

All genes for the biosynthesis of [biotin](#) of *Escherichia coli* are grouped in a single [operon](#) located at 17 min on the [chromosome](#), except for *bioH*, which is located at 75 min. The genes of the operon *bioA, BFC* are transcribed divergently from a single regulatory region located between the *bioA* and the *bioB* genes. [Transcription](#) in both directions is corepressed by biotinyl-5'-adenylate and the biotin [repressor](#). The following table shows the respective function of these genes:

---

|                             |   |
|-----------------------------|---|
| <i>bioC</i> and <i>bioH</i> | Unidentified steps in pimeloyl CoA biosynthesis |
| <i>bioF</i>                 | 7-Keto-8-aminopelargonic acid synthetase        |
| <i>bioA</i>                 | 7,8-Diaminopelargonic acid synthetase           |
| <i>bioD</i>                 | Dethiobiotin synthetase                         |
| <i>bioB</i>                 | Biotin synthase (introduction of the S atom)    |

---

No precise information is available on the function of the *bioC* and *bioH* genes, other than that a mutation in either gene results in no excretion of a known intermediate in the pathway. Therefore the products of these genes have been assigned to some early steps before 7-keto-8-aminopelargonic acid synthesis. All biotin genes are coordinately repressed when biotin is added to the growth medium in excess of 1 ng/ml.

The biotin repressor is a very interesting bifunctional protein of 321 amino acid residues, which acts at two different levels. In addition to its repressor function, it is endowed with acetyl CoA carboxylase biotin holoenzyme synthetase activity, which activates biotin to biotinyl-5'-adenylate and transfers the biotin to acceptor proteins. As soon as these proteins are totally biotinylated, biotinyl-5'-adenylate accumulates and serves as the corepressor of the biotin operon. This is a case of an enzyme synthesizing its own repressor, a unique property thus far among [DNA-binding proteins](#). Mutations in the corresponding gene (*birA*) inactivate the repressor function partially or totally and also alter the enzymatic function. This repressor binds to a 40-bp symmetrical biotin **operator** site to prevent [transcription](#) of the biotin biosynthetic genes. The structure of the repressor is highly asymmetrical and consists of three **domains**. The N-terminal domain is mostly **alpha-helical**, contains a [helix–turn–helix motif](#) and is loosely connected to the remainder of the molecule. The central domain consists of a seven-stranded mixed [beta-sheet](#), with  $\alpha$ -helices covering one face. The other side of the sheet is largely exposed to the solvent and contains the enzyme's [active site](#). The C-terminal domain comprises a six-stranded antiparallel  $\beta$ -sheet sandwich. The location of biotin is consistent with mutations affecting enzymatic activity.

Two molecules of the monomeric repressor bind cooperatively to one molecule of operator in the nanomolar concentration range. The data suggest that one molecule of repressor monomer binds to each of the two operator half sites and that they form a dimer only after they bind. Because the complex between repressor and DNA has not yet been crystallized, details of its structure remain an

open question.

## BiP (Hsp70)

The acronym BiP stands for binding protein, and it is the only member of the hsc,hsp70 family known to occur in the lumen of the [endoplasmic reticulum](#) (ER) of mammalian cells. It was identified originally as an ER protein that increases in amount when cells are starved of glucose (1) and subsequently as a protein that binds noncovalently to the heavy (H) chains of [immunoglobulins](#) as they enter the ER lumen (2). The alternative term of “Grp78” stands for **glucose-regulated protein** with an apparent subunit mass of about 78 kDa. BiP is now regarded as having a general **molecular chaperone** role for transport in the ER lumen and for [protein folding in vivo](#). BiP proteins are highly conserved in mammals, and there is 67% identity between mouse BiP and a homologue in the ER lumen of *Saccharomyces cerevisiae* called “Kar2p” (3). Like all hsc,hsp70 proteins, BiP exhibits a weak [ATPase](#) activity.

### 1. Structure

BiP is encoded in a nuclear gene located on human [chromosome](#) 9q34, and cDNA from mammalian species indicate a primary translation product of 635 amino acid residues, including an 18-residue *N*-terminal signal sequence for ER targeting and the *C*-terminal tetrapeptide **KDEL sequence** for ER retention (GenBank accession numbers M19645 and M17169). The *N* and *C* termini contain clusters of acidic residues thought to be involved in  $\text{Ca}^{2+}$  binding (4). No crystal structure is available for mammalian BiP, but a structure is known for the *N*-terminal 45-kDa fragment of bovine hsc70, which retains the ATPase activity (5), and for the *C*-terminal polypeptide-binding domain of the *Escherichia coli* hsp70 homologue called **DnaK** (6). Purification protocols for BiP utilize [affinity chromatography](#) on ATP columns (7). Recombinant hamster BiP has been purified from *E. coli* cells (8) and is marketed by StressGen Biotechnologies Corporation, as are polyclonal [antibodies](#) to rodent BiP. *In vivo* BiP exists in interconvertible [oligomeric](#) and monomeric forms and is subject to **phosphorylation** on serine and threonine residues, as well as to **ADP-ribosylation**. However, only monomeric, unmodified species of BiP are found in complexes with unfolded or unassembled polypeptides (9).

### 2. Function

BiP binds transiently to a range of newly synthesized secretory proteins, as they traverse the ER membrane and enter the ER lumen, and more permanently to misfolded, underglycosylated, or unassembled proteins whose transport from the ER is blocked (10); it does not bind to native, folded proteins. The binding is reversed by the addition of ATP and is believed to exert a molecular chaperone function by preventing premature folding and/or aggregation; this function is achieved by the shielding of potentially interactive **hydrophobic** surfaces during the time when BiP is bound. There is also evidence that the BiP homologue in yeast functions as a molecular [motor protein](#) to promote the transport of proteins across the ER membrane (11). This transport function requires the binding of the BiP homologue to the *J* domain of the yeast ER [membrane protein](#) Sec63p (12) (see also text later in this article). Studies of the binding of synthetic [peptides](#) and bacteriophage [peptide libraries](#) show that the optimum peptide length for binding is seven to eight residues, with extensive sequence diversity but a marked preference for hydrophobic residues (13, 14). These observations support the idea that BiP binds to a wide range of sequences that normally occur inside fully folded proteins. A computer program is available that scores potential BiP binding sites in protein sequences (14); it has been used to map such sequences in immunoglobulin heavy chains to the

regions that interact with light chains (15).

### 3. Interaction with Other Chaperones

Cytosolic hsc,hsp70 proteins in the bacterial and eukaryotic cytosol interact with other chaperones of the DnaJ (or hsp40) family that contain *J* domains (see [DnaK/DnaJ Proteins](#)). The yeast BiP homologue, Kar2p, interacts with two other ER proteins that contain *J* domains: Sec63p, a membrane protein involved in protein translocation across the ER membrane (16); and Scj1p, a luminal protein (17). BiP also binds either sequentially or simultaneously with chaperones such as **calnexin** and Grp94 during the folding in the ER of proteins such as immunoglobulin light (L) chains (18), thyroglobulin (19), vesicular stomatitis virus G protein (20), and [major histocompatibility complex](#) class II chains (21). BiP is thus one component in a complex set of interactions in the ER lumen between different chaperones and polypeptide chains that are folding.

### 4. Induction of BiP

BiP is an abundant protein under normal growth conditions, constituting about 5% of the ER luminal proteins, but its amount increases greatly under conditions that result in the accumulation of proteins within the ER lumen that are unable to fold correctly. These conditions include the biosynthesis of mutant chains, glucose starvation, and treatment with amino acid analogs, drugs that inhibit glycosylation, and calcium **ionophores** (22). The **promoters** of BiP genes in mammals contain several [cis-acting](#) regulatory elements required for high basal-level expression and for inducibility (22, 23).

### Bibliography

1. J. Pouyssegur, R. P. C. Shiu, and I. Pastan (1977) *Cell* **11**, 941–947.
2. I. G. Haas and M. Wabl (1983) *Nature* **306**, 387–389.
3. K. Normington, K. Kohno, Y. Kozutsumi, M. J. Gething, and J. Sambrook (1989) *Cell* **57**, 1223–1236.
4. D. R. J. Macer and G. L. E. Koch (1988) *J. Cell. Sci.* **91**, 61–70.
5. K. M. Flaherty, C. DeLuca-Flaherty, and D. B. McKay (1990) *Nature* **346**, 623–628.
6. X. Zhu, X. Zhao, W. F. Burkholder, A. Gragerov, C. M. Ogata, M. E. Gottesman, and W. A. Hendrickson (1996) *Science* **272**, 1606–1614.
7. P. J. Rowling, S. H. McLaughlin, G. S. Pollock, and R. B. Freedman (1994) *Protein Exp. Purif.* **5**, 331–336.
8. J. Wei and L. M. Hendershot (1995) *J. Biol.Chem.* **270**, 26670–26676.
9. P. J. Freiden, J. R. Gaut, and L. M. Hendershot (1992) *EMBO J.* **11**, 63–70.
10. M. J. Gething, S. Blond-Elguindi, K. Mori, and J. F. Sambrook (1994) in *The Biology of Heat Shock Proteins and Molecular Chaperones* (R. I. Morimoto, A. Tissieres and C. Georgopoulos, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, pp. 111–135.
11. S. Panzner et al. (1994) *Cell* **81**, 561–570.
12. J. L. Brodsky and R. Schekman (1993) *J. Cell Biol.* **123**, 1355–1363.
13. G. C. Flynn, J. Pohl, M. T. Flocco, and J. E. Rothman (1991) *Nature* **353**, 726–730.
14. S. Blond-Elguindi, S. E. Cwirla, W. J. Dower, R. J. Lipshutz, S. R. Sprang, J. F. Sambrook, and M. J. Gething (1993) *Cell* **75**, 717–728.
15. G. Knarr, M. J. Gething, S. Modrow, and J. Buchner (1995) *J. Biol. Chem.* **270**, 27589–27594.
16. D. Feldheim, J. Rothblatt, and R. Schekman (1992) *Mol. Cell. Biol.* **12**, 3288–3296.
17. G. Schlenstedt, S. Harris, B. Risse, R. Lill, and P. A. Silver (1995) *J. Cell Biol.* **129**, 979–988.
18. J. Melnick, J. L. Dul, and Y. Argon (1994) *Nature* **370**, 373–375.
19. P. S. Kim and P. Arvan (1995) *J. Cell Biol.* **128**, 29–38.

20. C. Hammond and A. Helenius (1994) *Science* **266**, 456–458.
21. M. S. Marks, R. N. Germain, and J. S. Bonifacino (1995) *J. Biol. Chem.* **270**, 10475–10481.
22. A. S. Lee (1992) *Curr. Opin. Cell Biol.* **4**, 267–273.
23. W. W. Li, L. Sistonen, R. I. Morimoto, and A. S. Lee (1994) *Mol. Cell. Biol.* **14**, 5533–5546.

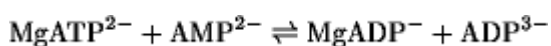
### Suggestion for Further Reading

24. D. N. Herbert, J. F. Simons, J. R. Peterson, and A. Helenius (1995) Calnexin, calreticulin and BiP/Kar2p in protein folding, *Cold Spring Harbor Symp. Quant. Biol.* **60**, 405–415.
25. F.-U. Hartl (1996) Molecular chaperones in cellular protein folding, *Nature* **381**, 571–580.
26. J. L. Brodsky (1996) Post-translational protein translocation: not all hsc70s are equal, *Trends Biochem. Sci.* **21**, 122–126.
27. M.-J. Gething, ed. (1997) *Molecular Chaperones and Protein Folding Catalysts*, Oxford University Press, Oxford (this volume contains much detailed information about BiP).

### Bisubstrate Analogue

Bisubstrate analogues were developed originally for mechanistic studies on [enzymes](#) that catalyze reactions with two substrates or products. However, they have also proved to be of value in studies on enzymes by [X-ray crystallography](#). Bisubstrate analogues are characterized by the fact that they embody in a single molecule the structural features of each of the two substrates (or products). Hence, it was expected that they would bind simultaneously at the binding sites for the two substrates, and therefore more tightly than the individual substrate molecules, and act as potent enzyme inhibitors. They have also been referred to as [transition state analogues](#), but such a classification seems inappropriate.

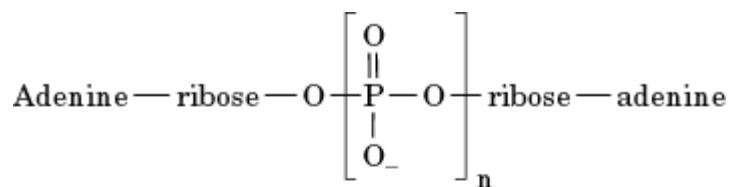
Two early, and now classical, examples of bisubstrate enzyme inhibitors are P<sup>1</sup>, P<sup>5</sup>-di-(adenosine-5') pentaphosphate (AP<sub>5</sub>A) and *N*-phosphonoacetyl-L-aspartate (PALA), whose structures are given in [Figure 1](#). AP<sub>5</sub>A was developed as an inhibitor of adenylate kinase, which is a monomeric enzyme that catalyzes the reaction:



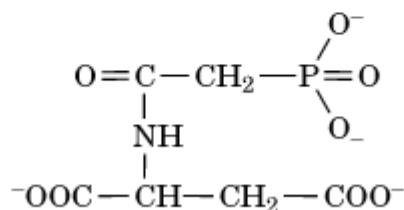
The reaction conforms to a rapid equilibrium, random **kinetic mechanism**, which implies that the enzyme possesses two distinct nucleotide-binding sites within its active site. One is for either MgATP<sup>2-</sup> or MgADP<sup>-</sup> and the other is for ADP<sup>3-</sup> or AMP<sup>2-</sup>. Irrespective of whether they are substrates or products, the moieties bound by the enzyme include two adenosine moieties and four phosphate moieties, whereas AP<sub>5</sub>A differs only in having five phosphate groups linked covalently. AP<sub>5</sub>A is a potent inhibitor of the adenylate kinase reaction, with the apparent inhibition constant  $K_i$  value being in the nanomolar region; the inhibition is competitive with respect to both MgATP and AMP, as would be expected ([1](#), [2](#)). Binding studies indicate that the stoichiometry of binding is 1:1 and the **dissociation constant** is 15 nM ([3](#)). The binding affinity is reduced seven-fold in the absence of Mg<sup>2+</sup>. It is of interest that the binding to adenylate kinase of AP<sub>4</sub>A, which is the equivalent of covalently linking ATP to AMP, or ADP to ADP, is almost 3,000-fold weaker. An increase in the number of phosphoryl groups to six also reduces the binding affinity by 400-fold. The crystal

structures of three adenylate kinases and enzyme-AP<sub>5</sub>A complexes were solved well ahead of the spatial assignment of the substrate binding sites. A review of the problems associated with the determination of the structural relationship between the two nucleotide binding sites of adenylate kinase has been presented (4).

**Figure 1.** Examples of bisubstrate analogues. (a), P<sup>1</sup>, P<sup>n</sup>-di(adenosine-5') n-phosphate, where n = 4 (AP<sub>4</sub>A), 5 (AP<sub>5</sub>A), or 6 (AP<sub>6</sub>A). (b), N-phosphonoacetyl-L-aspartate (PALA).



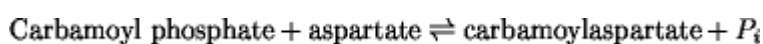
(a)



(b)

AP<sub>5</sub>A also acts as a strong inhibitor of ATP:NMP phosphotransferase from *Dictyostelium discoideum*, which utilizes either UMP or CMP as the acceptor of a phosphoryl group from ATP (5). Binding studies indicate that the dissociation constant of the enzyme-AP<sub>5</sub>A complex is 160 μM. However, the corresponding bisubstrate analogue with uridine replacing one adenine nucleotide (UP<sub>5</sub>A) is a more potent inhibitor and binds with a dissociation constant of 3 nM. The enzyme has been co-crystallized with UP<sub>5</sub>A.

PALA has been used extensively for kinetic and structural studies on [aspartate transcarbamoylase \(ATCase\)](#), which catalyzes the reaction:



PALA is considered an analogue of the two substrates linked covalently. The intact ATCase enzyme consists of both catalytic and regulatory units, does not show Michaelis-Menten kinetics, and is subject to **allosteric** control by nucleoside triphosphates. However, it is possible to prepare an active trimer of just the catalytic subunits that shows [Michaelis-Menten kinetics](#) and is not subject to allosteric control. This form of enzyme has been very useful for elucidating fundamental information about the catalytic sites and mechanism of action of ATCase. Investigations with the catalytic trimer have shown that one molecule of PALA binds to each subunit of the trimer (6) and the inhibition is linear competitive with respect to carbamoyl phosphate and linear noncompetitive relative to aspartate (7). Additional kinetic data have indicated that the kinetic mechanism for the aspartate transcarbamoylase reaction is essentially ordered, with carbamoyl phosphate being the first substrate to add; this explains the inhibition pattern. The pH-independent  $K_i$  value for the enzyme-PALA complex is 7.2 nM, which is three orders of magnitude lower than the dissociation constant for the



corresponding enzyme-carbamoyl phosphate complex (7).

PALA has also played an important role in the demonstration that the co-operativity observed with intact ATCase can be explained in terms of a two-state allosteric **concerted model** (8). Thus, the binding of up to three molecules of PALA to the intact enzyme, which has six [active sites](#), leads to activation, even though PALA is occupying active sites, as a result of displacing the equilibrium between the T (inactive) and R (active) forms of the enzyme toward R. Only with greater occupancy is PALA inhibitory. Structural studies support the idea that PALA is a true bisubstrate analogue of the substrates for the enzyme (9).

### Bibliography

1. G. E. Lienhard and I. I. Secemski (1973) *J. Biol. Chem.* **248**, 1121–1123.
2. P. Feldhaus, T. Frohlich, R. S. Goody, M. Isakov, and R. H. Schirmer (1975) *Eur. J. Biochem.* **57**, 197–204.
3. J. Reinstein, I. R. Vetter, I. Schlichting, P. Rosch, A. Wittinghofer, and R. S. Goody (1990) *Biochemistry* **29**, 7440–7450.
4. M.-D. Tsai and H. Yan (1991) *Biochemistry* **30**, 6806–6818.
5. L. Wiesmuller, K. Scheffzek, W. Kliche, R. S. Goody, A. Wittinghofer, and J. Reinstein (1995) *FEBS Lett.* **363**, 22–24.
6. G. R. Jacobson and G. R. Stark (1973) *The Enzymes* **9**, 225–308.
7. J. L. Turnbull, G. L. Waldrop, and H. K. Schachman (1992) *Biochemistry* **31**, 6562–6569.
8. L. E. Parmentier, M. H. O'Leary, H. K. Schachman, and W. W. Cleland (1992) *Biochemistry* **31**, 6598–6602.
9. W. N. Lipscomb (1994) *Adv. Enzymol.* **68**, 67–151.

### Bithorax Complex

The [homeotic genes](#) of the bithorax complex (BX-C) have been the subject of many landmark discoveries in the field of **developmental** biology. Homeotic genes were first identified in *Drosophila* by mutations that affect their expression. These mutations lead to spectacular effects on the morphology of the fly; they cause the development of a specific body structure (ie, segment, antenna, leg, wing) at the place where another structure normally develops. The first homeotic mutations, *bithorax* (*bx*) and *bithoraxoid* (*bxd*), were described in 1923 by Bridges and Morgan (1). For example, in flies homozygous for *bxd* mutations, their first abdominal segment (A1) develops like a copy of the segment immediately adjacent anteriorly—namely, the third thoracic segment (T3). Thus, the role of the *bxd*<sup>+</sup> function is to assign the identity of A1. In the late 1940s, Ed Lewis pioneered the field of developmental genetics by discovering the existence of additional homeotic mutations near *bx* and *bxd*. Although these mutations seem to affect different body segments of the thorax and abdomen, their [complementation](#) patterns turned out to be rather complex. Because they all mapped at the same genetic and cytological position, Lewis decided to name the locus the bithorax complex (BX-C). The work of Lewis was compiled in a milestone article in 1978 (2) in which he describes a series of homeotic mutations of BX-C that affect the identities of the third thoracic segment and all of the abdominal segments. Mutations in any of them transform the considered segment into a copy of the segment immediately anterior. Remarkably, genetic mapping revealed that these mutations are arranged on the [chromosome](#) in the same order along the anteroposterior axis as the segments they affect. By combining three mutations, Lewis produced flies

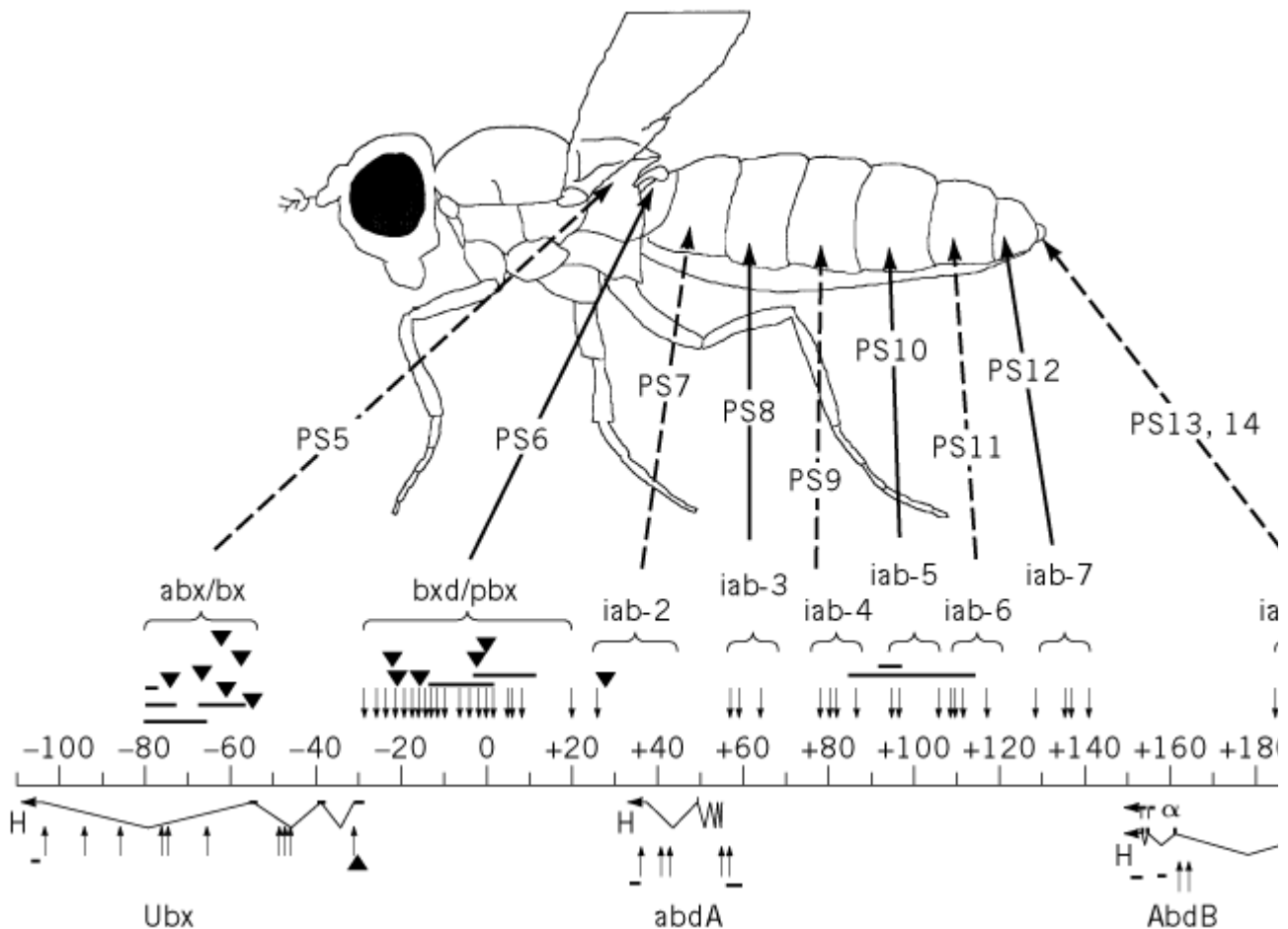
with four wings instead of two (transformation of T3 into T2). Because such animals look like more ancestor forms of insect, Lewis proposed that homeotic genes played an instrumental role in [evolution](#). The correspondence between the order of the genes on the chromosome and the order of the segments on the body of the fly has now attained almost mystical status and has turned out to be true also of the vertebrate homologues of BX-C. The Nobel Prize Committee has recognized this pioneering works by awarding its 1995 Nobel Prize of Medicine to Ed Lewis and two other *Drosophila* geneticists (Y. Nusslein-Volhard and E. Wieschaus).

In 1978, BX-C was the first *Drosophila* gene **cloned** from the chromosome without any prior knowledge of the products. To clone what turned out to be a large complex, Bender, Spierer, and Hogness developed the method of **chromosome walking** (3). Finally, in 1983, the molecular characterization of the BX-C led to the discovery of the homeobox in the laboratories of W. Gehring and of M. Scott (4, 5).

## 1. Molecular Genetics of BX-C

Figure 1 summarizes the molecular genetics of the BX-C. The complex covers 300 kbp of DNA, which are represented by the thin horizontal line marked off in kb (6, 7). Above the DNA line are represented the sites of the homeotic mutations that affect the identities of each of the segments under the control of the BX-C. The vertical arrows represent the sites of chromosomal rearrangement breaks, the triangles the sites of insertion of [transposons](#), and the horizontal lines the extent of deletions. Expression studies and analysis of the homeotic **phenotypes** in embryos have revealed that the unit transformed in each of these nine classes of mutations does not correspond to body segments; instead, it is composed of the posterior part of one segment and the anterior compartment of the next segment. These units are named parasegments (PSs) (8). For example, *bxd* mutations cause the transformation of the posterior part of T3 (pT3) and the anterior part of A1 (aA1) into p (T2) and aT3. This corresponds to the transformation of parasegment 6 (PS6) into PS5. The mutations affecting parasegment identity form nine discrete entities that, as predicted by Lewis, are aligned on the chromosome in the same order as the parasegments in which they act on the body of the fly (*abx/bx*, *bxd/pbx*, *iab-2* through *iab-8*). These mutations define nine PS-specific functions, and the arrows point toward the parasegments of the adult fly that are most affected in each class of mutations (Fig. 1). For reasons that will become clear later, it is worthwhile noting that all the mutations affecting the PS-specific functions are due to chromosomal rearrangements (more than 100 have been mapped). Thus, it seems impossible to affect these functions by point mutations, and it is unlikely that they correspond to individual **genes** coding for distinct [proteins](#).

**Figure 1.** Genetic map of the bithorax complex (BX-C). The 300 kilobases of DNA that comprise the complex are indicated by the bottom line, which is marked off in  $10^4$  base pairs. Above the DNA line are represented the sites of the homeotic mutations that affect the identities of each of the segments under control of BX-C. The vertical arrows represent the sites of chromosomal rearrangement breaks, the triangles the sites of insertion of transposons, and the horizontal lines the extent of deletions. The 14 parasegments (PS5 to PS14) that each part of the complex controls are indicated above, in the adult fly. The three transcription units (*Ubx*, *abdA*, and *AbdB*) are indicated below the DNA; in each case, transcription occurs from right to left. The exons are indicated by the thick horizontal lines, the introns by the thin V's connecting them. H indicates the homeobox domain. T1 and T2 alternative promoters a and g are indicated for *AbdB*. Mutations that affect the transcription units are indicated below the



**Northern blots** and the isolation and sequencing of [complementary DNA](#) revealed the existence of only three [transcription](#) units *Ubx*, *abd-A*, and *Abd-B*. These three transcription units are all transcribed from right to left (of Fig. 1) and contain a **homeobox** sequence at their 3' end. The *Ubx* transcription unit covers 70 kb of DNA and can generate 12 different transcripts by [alternative splicing](#) and [polyadenylation](#). [Translation](#) of these [messenger RNAs](#) yields a family of six proteins characterized by constant amino- and carboxy-proximal regions of 247 and 99 amino acid residues, respectively. The latter homeobox sequence is encoded by the 3'-terminal common exon. The members of this family are distinguishable by a short variable region that links the constant regions and consists of different combinations of three optional elements of 9, 17, and 17 residues (9-11). The spectrum of RNA products changes with time and tissue. It has been demonstrated that the *Ubx* protein containing the portion coded by the second microexon is not expressed in the nervous system (12). However, the different **isoforms** of the *Ubx* product are not all essential, since a mutation that eliminates the second microexon (and thus four isoforms) has no effect on fly viability and development (13).

The *abd-A* transcription unit is spread over a 20-kb region of DNA, and the mature *abd-A* transcript is composed of at least eight exons (14). As in *Ubx*, the *abd-A* transcription unit contains a microexon, but no alternative splicing generating multiple forms of the *abd-A* product has been detected thus far. The homeobox homology is found near the middle of the 330-residue protein. Analysis of the sequence of the whole BX-C (15) indicates that the open reading frame extends further upstream from the published ATG [initiation codon](#), raising the possibility that *abd-A* may be slightly larger than originally thought. In as much as this upstream open reading frame is conserved in *Tribolium*, it is likely that *abd-A* consists of 590 residues (16).

Finally, the *Abd-B* transcription unit consists of three classes of transcripts that are generated by the use of three alternate **promoters** and differential splicing (Fig. 1). While the alpha class of transcripts produces a protein of 55 kDa, the beta and gamma forms produce a product of 30 kDa that is truncated at the *N*-terminus (17, 18).

Mutations that affect the *Ubx*, *abd-A*, and *Abd-B* transcription units (shown below the DNA line in Fig. 1) are all lethal at the embryo stage of development, and cuticle analysis of the dead embryos detects homeotic transformations. In *Ubx* mutants, PS5 and PS6 are transformed into PS4. If such an embryo could survive, it would give rise to a fly with T3 and A1 transformed into T2 (ie, a fly with three pairs of wings). Homozygous *abd-A* embryos have PS7, 8, and 9 transformed into PS6 (in the adult this would correspond to a transformation of A2, A3, and A4 into A1). Finally, *Abd-B* mutations have PS10, 11, 12, and 13 transformed into PS9 [A5–A8 transformed into A4 (2, 7, 19, 20)].

## 2. Expression

The genetic and molecular data that have been described thus far appear to conflict. On the one hand, genetic analysis reveals the existence of nine parasegment-specific functions that are responsible for the identity of PS5 to 13, which will form the posterior thorax and the abdominal segments of an adult fly. On the other hand, molecular studies indicate that BX-C encodes only three protein products. This apparent discrepancy was solved when **antibodies** directed against the *Ubx*, *abd-A*, and *Abd-B* products became available, allowing determination of which part of the embryo these genes are expressed. A first observation derived from these studies is that each of the *Ubx*, *abd-A*, and *Abd-B* genes are expressed in domains composed of several parasegments. *Ubx* is expressed from PS5 to 13 (9, 21), *abd-A* from PS7 to 12 (14, 22), and *Abd-b* from PS10 to 14 (23, 24). Second, the expression patterns are complex, intricate, and dynamic. Comparisons of the expression patterns between wild-type embryos and those carrying mutations in the PS-specific functions revealed that the latter correspond to large *cis*-regulatory regions that are responsible for construction of the complex expression patterns of *Ubx*, *abd-A*, and *Abd-B*. The *abx/bx* and *bxl/pbx* *cis*-regulatory regions are responsible for UBX expression in PS5 and 6, respectively (21, 25, 26). The *iab-2*, *iab-4*, and *iab-4* *cis*-regulatory regions control expression of *Abd-A* in PS7, 8, and 9, respectively, while *iab-5*, *iab-6*, *iab-7*, and *iab-8* are responsible for the pattern of ABD-B expression (initiated from the a promoter) in PS10 to 13 (14, 24, 27). PS14 expresses a truncated form of ABD-B, resulting from transcription initiated from the b and g promoters.

## 3. BX-C Regulation

BX-C gene regulation can be divided into two phases, initiation and maintenance. During the early phases of embryogenesis, when parasegment identity is initially selected, the PS-specific *cis*-regulatory regions are the targets of the **gap gene** and **pair-rule gene** products (28-31). These gap and pair-rule proteins activate the *cis*-regulatory regions in successively more posterior parasegments.

The gap and pair-rule gene products are present only transiently during early development. The fact that homeotic genes are expressed throughout development implies the existence of a mechanism that maintains the activity state of each of the *cis*-regulatory regions. This maintenance system requires the **Polycomb group** (*Pc-G*) and the **trithorax group** (*trx-G*) genes (32-34). While the products of the *Pc-G* function as negative regulators, the products of the *trx-G* act as positive regulators. The products of the *Pc-G* exert their regulatory effects by interacting with specific elements in each of the *cis*-regulatory domains called **polycomb-response elements** (35-39). There may be equivalent or overlapping *trx* response elements for the *trx-G* proteins (36, 40). Though their precise mode of action is unknown, the products of the *Pc-g* and *trx-g* are thought to stabilize the expression patterns in each parasegment by **imprinting** an inactive or active **chromatin** conformation of the PS-specific *cis*-regulatory subregions (33, 41, 42).

## 4. Regulatory Elements of BX-C

Molecular studies using [reporter gene](#) constructs have revealed the existence of elements within the PS-*cis*-regulatory units that seem to be responsible for the initiation and maintenance phases of BX-C regulation. Some DNA fragments are able to initiate expression of a *Ubx-lacZ* reporter gene in the proper parasegments during early embryonic development ([28](#), [30](#), [43](#), [44](#)). In most cases, however, these patterns are not maintained, and expression expands into more anterior parasegments around the time when BX-C regulation would switch to the maintenance mode. Other BX-C DNA fragments are capable of retaining the appropriate parasegmental restrictions in *lacZ* expression after the gap and pair-rule gene products disappear. These fragments contain “maintenance elements,” also known as *Pc-g* response elements because their activity depends on *Pc-g* gene products ([35-39](#), [43](#), [44](#)). Finally, a third type of regulatory elements that has been identified in experiments with *Ubx-lacZ* reporter constructs are tissue- or cell-type-specific enhancers. They induce *lacZ* expression in specific tissue or cell types, with no restriction along the anteroposterior axis.

Many observations suggest that the PS-specific *cis*-regulatory units are organized into functionally independent domains. This is best illustrated by the expression patterns of “enhancer trap” transposons integrated in different domains of the complex ([45](#), [46](#)). These enhancer traps are subject to regulatory elements located within the same domain, but they are insensitive to regulatory elements in adjacent domains. The autonomy of each domain is ensured by elements that are believed to function as boundaries. Two such regulatory elements, *Mcp* and *Fab-7*, have been identified. *Mcp* is located between the *iab-4* and *iab-5* *cis*-regulatory units or domains, while *Fab-7* is located between *iab-6* and *iab-7* ([46-50](#)).

## 5. Concluding Remarks

Molecular analysis of the BX-C has confirmed most of the predictions that Lewis had foreseen in his 1978 model. He had initially envisioned activation of a new gene product in each parasegment. It is now clearly established that there are only three major groups of related protein products encoded by BX-C (UBX, ABD-A, and ABD-B). Discrete genetic units exist, however, that are sequentially activated in each parasegment. These units function as transcription regulatory regions (PS-specific *cis*-regulatory regions). The complex *cis*-regulation that they mediate results in a very intricate pattern of expression, both between and within parasegments. Each parasegment is a mosaic of cells expressing different homeotic products. Under the direction of these proteins, different cells adopt different fates, yielding the complex array of pattern elements that characterizes a given parasegment (or segment). The PS-specific *cis*-regulatory regions are large (*bxd* is spread over more than 30 kb; see [Fig. 1](#)) and can act from remote distances on their target promoters (*iab-5* is localized 50 kb away from its target *Abd-B* promoter). These properties suggest that the structure of the chromatin plays an important role to allow such long-distance interactions. Chromatin structure is also evoked by the properties of the *Pc-G* gene. The products of these *Pc-G* genes function as cellular memory to maintain the repressed state of the homeotic genes in body regions where they have not been activated during early development. There are analogies between the *Pc-G* repression and **mating-type** silencing in yeast or heterochromatic [position-effect](#) variegation in *Drosophila*. Though little is known at the molecular level, these analogies suggest that *Pc-G* repression involves the formation of a complex of *Pc-G* proteins leading to a chromatin structure that is refractory to transcription. The finding of boundary elements insulating adjacent PS-specific *cis*-regulatory regions has led to a model in which the sequential activation the *cis*-regulatory regions would be due to the stepwise opening of chromosomal domains ([45](#), [46](#), [49](#), [51](#)). Although no molecular clues exist to support such a model, it provides a rationale for the remarkable correspondence between the genomic organization of the BX-C and the anteroposterior axis of the fly. A similar model has been discussed recently in the case of the clusters of homeotic genes in mice, the Hox clusters ([52](#)).

## Bibliography

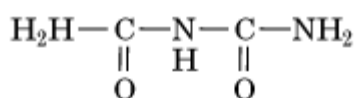
1. C. B. Bridges and T. Morgan (1923) Carnegie Inst. Washington Publ. **327**, 137.
2. E. B. Lewis (1978) Nature **276**, 565–570.

3. W. Bender, P. Spierer, and D. S. Hogness (1983) *J. Mol. Biol.* **168**, 17–33.
4. W. McGinnis, M. S. Levine, E. Hafen, A. Kuroiwa, and W. J. Gehring (1984) *Nature* **308**, 428–433.
5. M. P. Scott and A. J. Weiner (1984) *Proc. Natl. Acad. Sci. USA* **81**, 4115–4119.
6. W. Bender, M. Akam, F. Karch, P. A. Beachy, M. Peifer, P. Spierer, E. B. Lewis, and D. S. Hogness (1983) *Science* **221**, 23–29.
7. F. Karch, B. Weiffenbach, M. Peifer, W. Bender, I. Duncan, S. Celniker, M. Crosby, and E. B. Lewis (1985) *Cell* **43**, 81–96.
8. A. Martinez-Arias and P. Lawrence (1985) *Nature* **313**, 639–642.
9. P. A. Beachy, S. L. Helfand, and D. S. Hogness (1985) *Nature* **313**, 545–551.
10. M. B. O'Connors, R. Binari, L. A. Perkins, and W. Bender (1988) *EMBO J.* **7**, 435–445.
11. K. Kornfeld, R. B. Saint, P. A. Beachy, P. J. Harte, D. A. Peattie, and D. S. Hogness (1989) *Genes & Dev.* **3**, 243–258.
12. R. Weinzierl, J. M. Axton, A. Ghysen, and M. Akam (1987) *Genes & Dev.* **1**, 386–397.
13. A. Busturia, I. Vernos, J. Casanova, and G. Morata (1990) *EMBO J.* **9**, 3551–3555.
14. F. Karch, W. Bender, and B. Weiffenbach (1990) *Genes Dev.* **4**, 1573–1587.
15. C. H. Martin, C. A. Mayeda, C. A. Davis, C. L. Ericsson, J. D. Knafels, D. R. Mathog, S. E. Celniker, E. B. Lewis, and M. J. Palazzolo (1995) *Proc. Natl. Acad. Sci. USA* **92**, 8398–8402.
16. T. D. Shippy, S. J. Brown, and R. E. Denell (1998) *Dev. Genes Evol.* **207**, 446–452.
17. M. Zavortink and S. Sakonju (1989) *Genes Dev.* **3**, 1969–1981.
18. S. E. Celniker, D. J. Keelan, and E. B. Lewis (1989) *Genes Dev.* **3**, 1424–1436.
19. E. Sanchez-Herrero, I. Vernos, R. Marco, and G. Morata (1985) *Nature* **313**, 108–113.
20. J. Casanova, E. Sanchez-Herrero, A. Busturia, and G. Morata (1987) *EMBO J.* **6**, 3103–3109.
21. R. A. H. White and M. Wilcox (1985) *Nature* **318**, 563–567.
22. A. Macias, J. Casanova, and G. Morata (1990) *Development* **110**, 1197–1207.
23. M. DeLorenzi, N. Ali, G. Saari, C. Henry, M. Wilcox, and M. Bienz (1988) *EMBO J.* **7**, 3223–3231.
24. S. E. Celniker, S. Sharma, D. J. Keelan, and E. B. Lewis (1990) *EMBO J.* **9**, 4277–4286.
25. C. Cabrera, J. Botas, and A. Garcia-Bellido (1985) *Nature* **318**, 569–571.
26. R. A. H. White and M. E. Akam (1985) *Nature* **318**, 567–569.
27. E. Sanchez-Herrero (1991) *Development* **111**, 437–449.
28. S. Qian, M. Capovilla, and V. Pirrotta (1991) *EMBO J.* **10**, 1415–1425.
29. J. Muller and M. Bienz (1992) *EMBO J.* **11**, 3653–3661.
30. M. J. Shimell, J. Simon, W. Bender, and M. B. O'Connor (1994) *Science* **264**, 968–971.
31. F. Casares and E. Sanchez Herrero (1995) *Development* **121**, 1855–1866.
32. A. Shearn (1989) *Genetics* **121**, 517–525.
33. R. Paro (1990) *Trends Genet.* **6**, 416–421.
34. J. Simon, A. Chiang, and W. Bender (1992) *Development* **114**, 493–505.
35. J. Simon, A. Chiang, W. Bender, M. J. Shimell, and M. O'Connor (1993) *Dev. Biol.* **158**, 131–144.
36. C. S. Chan, L. Rastelli, and V. Pirrotta (1994) *EMBO J.* **13**, 2553–2564.
37. B. Christen and M. Bienz (1994) *Mech. Dev.* **48**, 255–266.
38. A. Chiang, M. B. O'Connor, R. Paro, J. Simon, and W. Bender (1995) *Development* **121**, 1681–1689.
39. S. Poux, C. Kostic, and V. Pirrotta (1996) *EMBO J.* **15**, 4713–4722.
40. H. Strutt, G. Cavalli, and R. Paro (1997) *EMBO J.* **16**, 3621–3632.

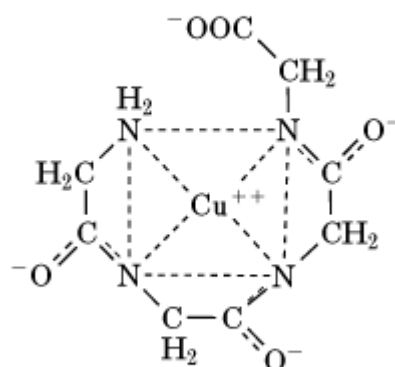
41. V. Pirrotta and L. Rastelli (1994) *Bioessays* **16**, 549–556.
42. J. Simon (1995) *Curr. Opin. Cell Biol.* **7**, 376–385.
43. J. Simon, M. Peifer, W. Bender, and M. O'Connor (1990) *EMBO J.* **9**, 3945–3956.
44. J. Muller and M. Bienz (1991) *EMBO J.* **10**, 3147–3155.
45. K. McCall, M. B. O'Connor, and W. Bender (1994) *Genetics* **138**, 387–399.
46. M. Galloni, H. Gyurkovics, P. Schedl, and F. Karch (1993) *EMBO J.* **12**, 1087–1097.
47. H. Gyurkovics, J. Gausz, J. Kummer, and F. Karch (1990) *EMBO J.* **9**, 2579–2585.
48. F. Karch, M. Galloni, L. Sipos, J. Gausz, H. Gyurkovics, and P. Schedl (1994) *Nucleic Acids Res.* **22**, 3138–3131.
49. J. Mihaly, I. Hogga, J. Gausz, H. Gyurkovics, and F. Karch (1997) *Development* **124**, 1809–1820.
50. J. Mihaly, I. Hogga, S. Barges, M. Galloni, R. Mishra, K. Hagstrom, M. Müller, P. Schedl, L. Sipos, J. Gausz, H. Gyurkovics, and F. Karch. (1998) *CMLS*, **54**, 60–70.
51. M. Peifer, F. Karch, and W. Bender (1987) *Genes Dev.* **1**, 891–898.
52. T. Kondo, J. Zákány, and D. Duboule (1998) *Molec. Cell* **1**, 289–300.

## Biuret Reaction

Proteins form biuret-like compounds (structure I) in alkaline solution by losing a proton from the nitrogen atoms of the peptide bonds. In strong alkaline solution and  $\text{Cu}^{2+}$ , this compound forms a blue colored complex (structure II). This color reaction is called the biuret reaction. The sensitivity of the reaction is not very great, but the color yield is similar from protein to protein. Thus, this reaction has been employed to quantify proteins in widely varying samples.



[I] Biuret; carbamylurea



[II] A kind of biuret compound;  
copper glycine peptide complex  
under alkaline conditions

According to the method of Gornall et al. (1), 1 to 10 mg of protein is determined by **absorbance** at 540 nm after incubation with biuret reagent (CuSO<sub>4</sub>, potassium sodium tartrate, NaOH) for 30 min at 20° to 25°C. An improved biuret method by Westley and Lambeth (2), which removes the free Cu<sup>2+</sup> ion with an ion-exchange resin and adds Na-diethyldithiocarbonate as a coloring reagent, detects 0.05 to 1 mg of protein. A simple and sensitive micro-biuret method was developed by Itzhaki and Gill (3), where UV absorbance at 310 nm is employed to quantify protein.

### Bibliography

1. A. G. Gornall, C. S. Bardawill, and M. M. David (1949) *J. Biol. Chem.* **177**, 751–766.
2. J. Westley and J. Lambeth (1960) *Biochim. Biophys. Acta* **40**, 364–366.
3. R. F. Itzhaki and D. M. Gill (1964) *Anal. Biochem.* **9**, 401–410.

## Blastoderm

The blastoderm is an [embryo](#) at the stage that consists of a single cell layer. This cell layer is formed by the cleavage divisions of the fertilized [zygote](#), largely under the control of maternally deposited gene products.

### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 5.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
- J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, U.K.

## Blastomere

After [fertilization](#), the [zygote](#) begins the cleavage divisions. During the cleavage divisions, the cells of the embryo are called *blastomeres*.

### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 5.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
- J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.



## Blastopore

The blastopore is a line of invagination on amphibian [embryos](#) for gastrulation (1). After fertilization, the amphibian [zygote](#) divides to form a mass of cells called blastomeres. The upper surface is called the *animal pole*, and the lower surface is called the *vegetal pole*. During gastrulation, the cells at the vegetal pole move inside the embryo, and the cells of the animal pole are left to form the outer surface of the embryo. The animal cells left at the surface will form the ectoderm, including the nervous system. The cells that move inside during gastrulation form the mesoderm and endoderm.

The process of gastrulation begins with the formation of a small indentation between the animal and vegetal poles on one side of the embryo. This small indentation is the beginning of formation of the blastopore. As gastrulation proceeds, the blastopore expands laterally to form a crescent, and eventually it encircles the entire embryo as gastrulation proceeds. The vegetal cells and the marginal cells between the vegetal and animal cells invaginate through the blastopore and form a second layer inside the embryo. As gastrulation continues, the blastopore shrinks to cover the remaining vegetal cells, called the *yolk plug*. When the blastopore almost completely covers the yolk plug, organogenesis begins. Interactions between various groups of cells during gastrulation cause cells to become determined to form various organs and structures. The determination of cells caused by interactions with their neighboring cells is called *developmental induction*. One of the most well-known examples is the [organizer](#) region of the amphibian embryo that specifies the determination of the nervous system. This is the dorsal lip of the blastopore and is often called the *Spemann organizer*, because it was first characterized by Hans Spemann (2).

### Bibliography

1. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 161–164.
2. H. Spemann (1938) *Embryonic Development and Induction*, Yale University Press, New Haven.

### Suggestions for Further Reading

3. S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
4. S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
5. L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
6. J. A. Moore (1972) *Heredity and Development*, 2nd ed., Oxford University Press, New York.

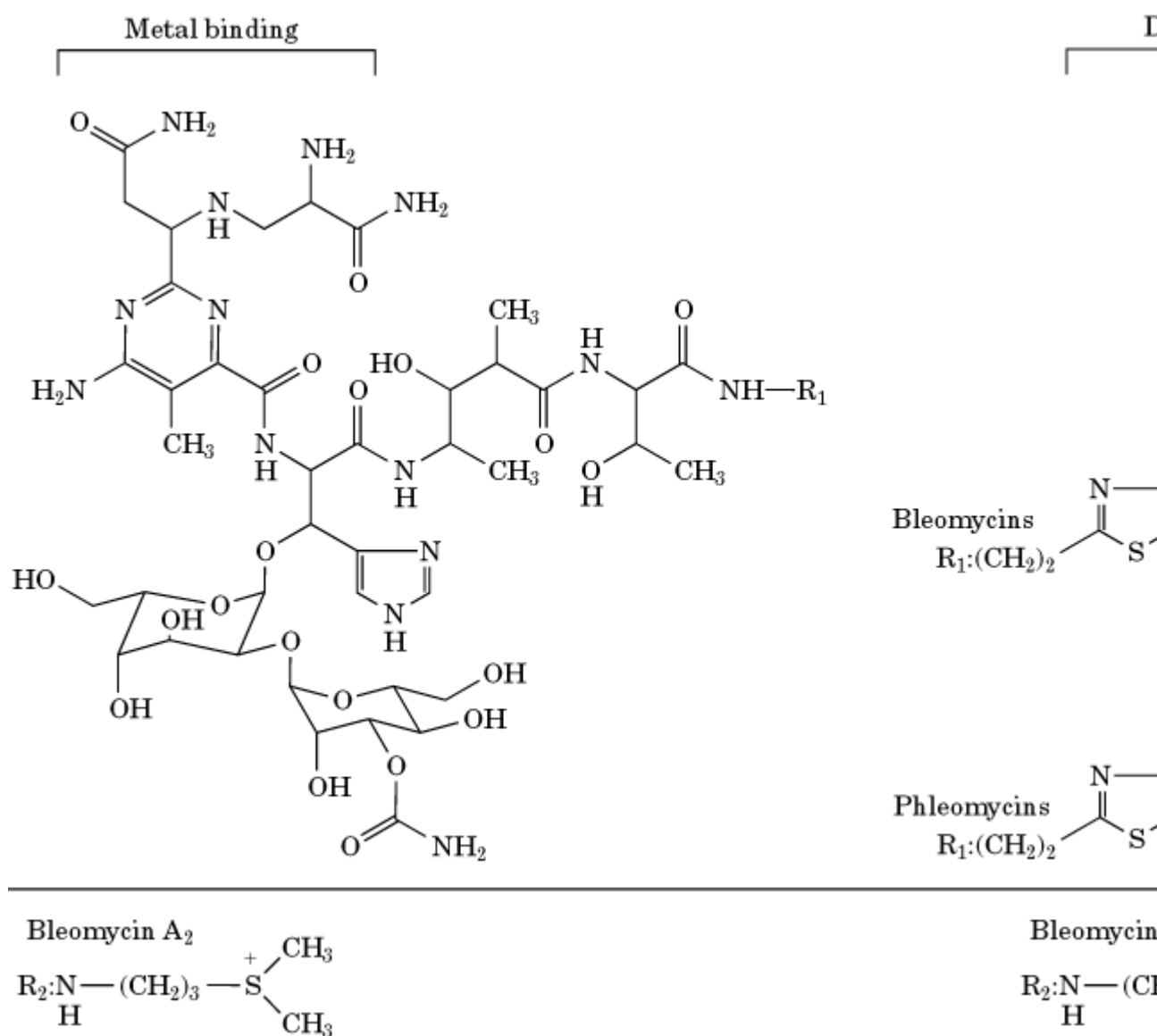
## Bleomycin

Bleomycin (BLM), discovered in 1962, is a family of low-molecular-weight metalloglycopeptides

( $M_r$  approximately 1500) produced by *Streptomyces verticillus* and isolated from the culture media of *S. verticillus* as copper complexes (1-3). It is widely used as an anticancer antibiotic in patients. BLM acts as a limited endonuclease and produces a variety of lesions in DNAs by a mechanism involving free radical attack on deoxyribose in both DNA strands. BLM damages DNAs in ways that mimic ionizing radiation. The detailed knowledge of the complicated chemistry of BLM and its interactions with DNA *in vivo* and *in vitro* drives much of its use in applications to molecular biology. BLM and structurally related analogues are considered radiomimetic and oxidative DNA-cleaving reagents, and as such they are utilized as chemical tools to study and understand the activities of this class of agents. Tools of molecular biology are also employed to understand the chemical, biological, and clinical aspects of the mechanism of action of BLMs. In addition, broadly used cloning strategies in several organisms use a gene conferring resistance to BLM and structurally related analogues as a selectable marker. Phleomycin (PLM) is used as the selective agent and is available in commercial formulations. BLM causes multiple changes to cells and is cytotoxic, and its cytotoxicity is high where there is no or only a limited barrier to BLM reaching cellular targets. BLM effectively kills all types of cells tested. Cellular resistance is conferred in several ways, including protection by nucleosomes in chromatin, DNA repair, metabolic inactivation of BLM, BLM-resistance proteins, and restricted entry of BLM. An overview of resistance mechanisms was recently published (4).

### 1. Structure and Activated Complex

The structurally complex BLMs contain a metal-binding domain and a **DNA-binding** domain (Fig. 1). BLMs require metals and oxygen species for their activity (5-9). The metal-binding domain is also the site of oxygen activation and is attached to a disaccharide group; it binds redox-active transition metals, such as Fe(II), Co(II), Cu(II), Ni(II), and Zn(II). The most stable BLM-O<sub>2</sub>-metal complex is formed with cobalt (8). When a BLM-Fe(II)-O<sub>2</sub> complex binds to DNA and the Fe(II) oxidizes to Fe(III), the complex attacks the C4' position of DNA deoxyribose (10-13). The complex thereby behaves as a limited endonuclease.



The coplanar bithiazole moiety partially intercalates into the minor groove between bases of DNA. The cationic C-terminal amines are also involved in interactions with nucleic acids. The terminal amines in BLM A<sub>2</sub> and BLM B<sub>2</sub> are dimethylsulfonium propylamine and agmatine, respectively, and are similar in length and bear one positive charge (Fig. 1). BLM B<sub>2</sub> produces considerably more DNA breaks and killing than BLM A<sub>2</sub> over a wide range of chemical concentrations (14). Without the terminal amine, the BLM molecule no longer cleaves DNA or possesses antitumor activity.

BLM and structurally related PLM (15) differ in the oxidation state of their sulfur heterocycles (Fig. 1). One of the two conjugated thiazole rings of the BLM bithiazole is modified by hydrogenation to 4,5-dihydrothiazole (thiazoline) in PLM (16, 17). In addition, the C-terminal amines in clinical preparations of BLMs differ chemically and quantitatively from prepared mixtures of PLMs. Tallysomycin is closely related to BLM (18, 19).

## 2. Anticancer Use

BLM is an important therapeutic agent that is useful as a single agent in treating several human cancers and is widely used in combination chemotherapy and radiotherapy. The water-soluble

product used in cancer treatment, Bleomycin, is a family of 11 metal-free congeners differing in their terminal amines. The clinical mixture (20) is comprised mainly of BLM A<sub>2</sub> [approximately 55% to 70% (usually 68% to 69%)] and BLM B<sub>2</sub> (approximately 25% to 32%). The effectiveness of BLM as an anticancer agent is associated with its ability to produce lesions in DNA (21, 22).

BLM is principally used in patients with lymphoma or a variety of solid tumors. It has been included in regimens for treating malignant and peripheral T-cell lymphomas, as well as in combination BEP chemotherapy (BLM, etoposide, cisplatin) for metastatic testicular teratoma and Hodgkins and non-Hodgkins lymphoma. BLM is also used in the chemotherapy and management of Kaposi's sarcoma (eg, pulmonary, gastrointestinal, epidemic, disseminated), as well as with vincristine in combination chemotherapy for epidemic Kaposi's sarcoma. BLM does not cause bone-marrow, hepatic, or renal toxicities, nor does it cause cardiotoxicity. Pulmonary fibrosis is a side effect of BLM treatment that limits its use. The molecular mechanism of lung fibrosis is being investigated, but is not fully known. Patients with genetic susceptibility are particularly susceptible to lung fibrosis (23, 24). In animal models, taurine (25) or taurine and niacin (26) counter BLM-induced lung fibrosis. The newer BLM derivatives, peplomycin and liblomycin, were developed because of their lower pulmonary toxicity and broader antitumor spectrum in animal studies (27).

### 3. Mechanism of Action on Nucleic Acids

#### 3.1. Chemical Action on DNA: DNA Damage

The unique chemical action of BLMs in the presence of oxygen and Fe(II) catalytically cleaves double-stranded DNA *in vivo* and *in vitro*. Single- and double-stranded breaks are produced, leaving 5'-phosphate- and 3'-phosphoglycolate-termini (11, 28-32). The most genotoxic and lethal lesions for cells are double-stranded breaks. BLM recognizes 5'-phosphoguananylyl(3',5')thymidine or 5'-phosphoguananylyl(3',5')cytosine sequences most frequently, releasing the pyrimidines when they are located to the 3' side of guanosine (28-30, 33-36) and leaving DNA alkali-labile (32, 37). While cleavage at the first site is G-Py-specific, the second nucleophilic attack by a BLM-Fe(II)-O<sub>2</sub> complex on the opposite DNA strand is not a sequence-specific cleavage and instead is probably targeted by the structural perturbation of the DNA at the first cleavage site (38, 39). Deoxyribose degradation also produces 3-(pyrimidin-1'-yl)-2-propenal and 3-(purin-9'-yl)-2-propenal (11, 31, 40) and oligonucleotide 3'-(phosphoro-2-O-glycolic acid) derivatives (11, 31, 41). DNA strand breaks are stoichiometric with the production of base propenals (42).

#### 3.2. Preferential Cleavage between Nucleosomes

BLM preferentially cleaves **linker regions** of **chromatin** between **nucleosomes** (43-47). The enhanced resistance of DNA bound to nucleosomes seems to be related to the necessity of a conformational change for BLM binding and intercalation. In yeast, internucleosomal cleavage and DNA degradation (46) and killing (48) are less pronounced in logarithmic phase cells than in cells that are in stationary phase. The cellular and molecular basis for this may relate to the highly effective use of BLM on particular solid tumors.

#### 3.3. Bleomycin and Phleomycin

PLM was discovered before bleomycin, but was too toxic in patients to be used as an anticancer drug. PLM is also substantially more effective than BLM on a per mole basis in yeast in producing cell killing (48), DNA breaks in intracellular DNAs (49), release of nucleosomes from chromatin (45), and genetic changes (50). Thus, the DNA lesions produced by the BLM and PLM could differ in their nature or frequency, or they could be processed differently by the cells. The mechanisms of BLM and PLM interaction with DNA *in vitro* also appear to differ (51-54). PLM exhibits a higher requirement than BLM for ferrous ions (49). Bithiazole intercalation in DNA is thought to be necessary for producing double-stranded, but not single-stranded, breaks *in vitro* (52). Accordingly, BLM produces more double-strand breaks than PLM in PM2 phage DNA (52, 54). Cu(II)BLM, but not Cu(II) PLM, intercalates (51). BLM, but not PLM, degrades relaxed DNA to a greater extent

than either positively or negatively **superhelical** DNA (54). On the other hand, BLM and PLM cleave DNA *in vitro* at similar preferred sites at similar frequencies (55, 56) and produce comparable numbers of DNA breaks under some conditions (57, 58).

### 3.4. RNA

BLM cleaves a variety of RNAs. The most studied have been [transfer RNA](#) and tRNA precursors. In contrast to DNA cleavage, BLM-induced RNA cleavage is usually wholly or predominately at a single site, often at a junction between single- and double-stranded RNA regions (59-62).

## 4. DNA Repair

### 4.1. Biochemical and Genetic Evidence for Radiomimetic Properties

The cellular processes that repair BLM-induced DNA lesions are not entirely known. The most important mechanisms for the repair of BLM-induced DNA damage are [recombinational repair](#), [base-excision repair](#), and post-replication repair. BLM and ionizing radiation produce similar lesions in DNA, although BLM produces a narrower spectrum of products than ionizing radiation. DNA breaks are introduced approximately linearly with increasing concentrations of BLM (14) and ionizing radiation (63). Some of the chromosomal lesions produced after either BLM treatments or ionizing irradiation can be ligated immediately and lead to a quick component of DNA rejoining, but other lesions require more time for processing before the termini of the DNA molecules become substrates for ligation (14, 63-65). For example, the unusual phosphoglycolate must be removed by the DNA 3'-repair diesterase activity of apurinic and apyrimidinic endonuclease (66, 67). After extension of the remaining DNA strand by one nucleotide, DNA ligase resynthesizes intact DNA molecules by forming a phosphodiester bond between adjacent 3'-hydroxyl and 5'-phosphoryl termini (14). Rapid and slower components of DNA rejoining have been identified in several laboratories. An ultrarapid phase of cellular recovery accompanies rapid repair in human cells (64, 68), suggesting that some of the cytotoxic treatment effects of BLM could be counteracted in the clinic and reduce chemotherapeutic effectiveness.

The importance of DNA repair is shown by the many studies in various organisms of mutants hypersusceptible to killing by BLM and defective in repair of DNA lesions. All *rad* mutations of *Saccharomyces cerevisiae* (69-72) that confer hypersensitivity to killing by BLM analogues also confer cross-sensitivity to ionizing radiation (73-78), so pathways are shared for the repair of chromosomal damage by bleomycins and ionizing radiation. Moreover, mutant strains with altered resistance to lethal effects of BLM have been isolated and characterized in several laboratories, and direct selection for mutations conferring hypersensitivities to lethal effects of BLM resulted in mutants exhibiting cross-hypersensitivities to ionizing radiation.

### 4.2. Genetic Recombination and Mutation

Mechanisms of recombination are important for DNA repair and rejoining [recombinant DNA](#). BLM is recombinogenic and mutagenic (38, 50, 79-85). The amount of recombination or mutation caused by BLM depends upon the assay system and treatment conditions. BLM was found to be weakly mutagenic to mitochondrial DNA (86).

## 5. Additional Cellular Targets

### 5.1. Membrane

The plasma membrane restricts BLM internalization in some mammalian cells (87, 88). This could be due to the polar and charged groups on the BLM molecule. Membrane damage by BLM (68, 89) or cell electropermealization (88) circumvents this restriction.

### 5.2. Cell Wall

BLM molecules are readily taken up into yeast cells, but are not equally distributed from cell to cell (89). BLM initially localizes to cell walls, causes cell wall and membrane damage, and aids the enzymatic conversion of cells to spheroplasts (75, 89-91). BLM alters the anchorage of several

mannoproteins in the cell wall matrix of intact cells or isolated cell walls, and it disrupts essential cell wall polymers. These activities facilitate the entry of BLM into yeast cells.

### 5.3. BLM Hydrolase

BLM hydrolase hydrolyzes and inactivates BLM. The enzyme is a [thiol proteinase](#) that has DNA-binding and peptide-cleavage domains. It binds DNA and RNA. BLM hydrolase activity protects cells from BLM toxicity, but limits the use of BLM in cancer chemotherapy. Although its normal cellular function is unknown, the enzyme is present in diverse organisms, including humans ([92-96](#)) and yeast ([97-100](#)). Expression of the yeast BLM hydrolase in mammalian cells results in resistance to BLM ([101](#)). A member of a galactose regulatory system ([102, 103](#)), yeast BLM hydrolase binds to nicked double-stranded DNA, single-stranded DNA, and RNA, without sequence specificity ([104](#)). This enzyme associates with plasma membranes and is in the cytosol ([105](#)). Human BLM hydrolase was recently shown to exhibit endopeptidase activity ([106](#)). This enzyme is thought to play a role in the development of resistance to BLM during chemotherapy ([95, 107-109](#)).

### 5.4. BLM Resistance Proteins

Several proteins in microorganisms that produce BLM or structural analogues actually confer resistance to these products. Vectors bearing genes encoding these proteins confer high levels of resistance to BLM and related antibiotics and are used as cloning vehicles. The proteins bind and form stable complexes with BLM with high specificity, thereby preventing BLM from complexing with DNA. One of these proteins is encoded by the *Streptoalloteichus hindustanus* (*Sh*) *ble* gene ([110-112](#)) found on the Tn5 bacterial [transposon](#), along with additional genes encoding resistance to other antibiotics. Expression of the *Sh ble* gene in transgenic mice reduced BLM toxicity in the mice and protected against lung fibrosis ([113, 114](#)). High levels of the protein were detected in lungs, kidney, and spleen.

### 5.5. Multiple Targets and Cytotoxicity

Multiple cellular targets of BLM are expected because of the chemical mechanism of action of the molecule and because some of the BLM-hypersensitive mutants isolated in different organisms do not exhibit hypersensitivity to lethal effects of radiation. Multiple cellular enzymes are important for surviving the toxicities of BLM, and deficiencies in their function could reduce chances of survival. The relationship and significance of each cellular target of BLM to cellular toxicity are not known. Why BLM causes [apoptosis](#) in some types of cells and not in others is also unknown. BLM is not considered apoptotic at low concentrations. Far less is known about targets of BLM in cells than about the chemical mechanism of action of BLM *in vitro* on defined substrates.

Additional factors modulate BLM activity. Generally, BLM preferentially affects mitotic cells in the G2-M phase of the [cell cycle](#). BLM is unlikely to be a useful agent to study meiosis because of the abundant lesions the molecule produces in DNA. Intracellular metal ion concentrations and pH also modulate BLM activities ([48, 115, 116](#)). Studies elucidating the roles of the multiple targets of BLM and factors that modulate BLM activities will improve our understanding and efficacious uses of the widely studied BLM family of related compounds.

## Bibliography

1. H. Umezawa, K. Maeda, T. Takeuchi, and Y. Okami (1966) *J. Antibiotics* **19**, 200–209.
2. H. Umezawa, Y. Suhara, T. Takita, and K. Maeda (1966) *J. Antibiotics* **19**, 210–215.
3. H. Umezawa (1976) *GANN* **19**, 3–36.
4. W. S. El-Deiry (1997) *Curr. Opin. Oncol.* **9**, 79–87.
5. E. A. Sausville, J. Peisach, and S. B. Horwitz (1976) *Biochem. Biophys. Res. Commun.* **91**, 871–877.
6. R. M. Burger, J. Peisach, and S. B. Horwitz (1981) *J. Biol. Chem.* **256**, 11636–11639.
7. R. M. Burger, S. J. Projan, S. B. Horwitz, and J. Peisach (1986) *J. Biol. Chem.* **261**, 15955–15959.

8. J. Stubbe and J. Kozarich (1987) *Chem. Rev.* **87**, 1107–1136.
9. P. C. Dedon and I. H. Goldberg (1992) *Chem. Res. Toxicol.* **5**, 311–332.
10. R. M. Burger, J. Peisach, and S. B. Horwitz (1981) *Life Sci.* **28**, 715–727.
11. L. Giloni et al. (1981) *J. Biol. Chem.* **256**, 8608–8615.
12. J. W. Sam and J. Peisach (1993) *Biochemistry* **43**, 1488–1491.
13. L. F. Povirk, Y. H. Han, and R. J. Steighner (1989) *Biochemistry* **28**, 5808–5814.
14. C. W. Moore (1990) *Biochemistry* **29**, 1342–1347.
15. K. Maeda, H. Kosaka, K. Yagishita, and H. Umezawa (1956) *J. Antibiotics* **9**, 82–85.
16. T. Takita et al. (1972) *J. Antibiotics* **25**, 197–199.
17. D. A. McGowan, U. Jordis, D. K. Minster, and S. M. Hecht (1977) *J. Am. Chem. Soc.* **99**, 8078–8079.
18. H. Kawaguchi et al. (1977) *J. Antibiotics* **30**, 779–788.
19. M. Konishi et al. (1977) *J. Antibiotics* **30**, 789–805.
20. C. Crooke and W. Bradner (1976) *J. Med.* **7**, 333–425.
21. S. K. Carter, S. T. Crooke, and H. Umezawa (1978) *Bleomycin, Current Status and New Developments*, Academic Press, New York.
22. S. M. Hecht, ed. (1979) *Bleomycin: Chemical, Biochemical and Biological Aspects*, Springer-Verlag, New York.
23. C. K. Haston, C. I. Amos, T. M. King, and E. L. Travis (1996) *Cancer Res.* **56**, 2596–2601.
24. R. P. Marshall, R. J. McAnulty, and G. J. Laurent (1997) *Int. J. Biochem. Cell. Biol.* **29**, 107–120.
25. R. E. Gordon, R. F. Heller, and R. F. Heller (1992) *Adv. Exp. Med. Biol.* **315**, 319–328.
26. G. Gurujeyalakshmi, M. A. Hollinger, and S. N. Giri (1998) *Am. J. Respir. Cell. Mol. Biol.* **18**, 334–342.
27. T. Takita and T. Ogino (1987) *Biomed. Pharmacother.* **41**, 219–226.
28. A. D. D'Andrea and W. A. Haseltine (1978) *Proc. Natl. Acad. Sci. USA* **75**, 3608–3612.
29. M. Takeshita, P. Grollman, E. Ohtsubo, and H. Ohtsubo (1978) *Proc. Natl. Acad. Sci. USA* **75**, 5983–5987.
30. A. P. Grollman and M. Takeshita (1980) *Adv. Enzymol. Regul.* **18**, 67–83.
31. N. Murugesan et al. (1985) *Biochemistry* **24**, 5735–5744.
32. H. Sugiyama, C. Xu, N. Murugesan, and S. M. Hecht (1985) *J. Am. Chem. Soc.* **107**, 4101–4105.
33. M. Takeshita and P. Grollman (1979) In *Bleomycin: Chemical, Biochemical and Biological Aspects* (S. Hecht, ed.), Springer-Verlag, New York, pp. 207–221.
34. C. K. Mirabelli, C.-H. Huang, A. W. Prestayko, and S. T. Crooke (1982) *Cancer Chemother. Pharmacol.* **8**, 57–65.
35. C. K. Mirabelli et al. (1982) *Cancer Res.* **42**, 2779–2785.
36. V. Murray and R. F. Martin (1985) *J. Biol. Chem.* **260**, 10389–10391.
37. J. C. Wu, J. W. Kozarich, and J. Stubbe (1983) *J. Biol. Chem.* **258**, 4694–4697.
38. R. J. Steighner and L. F. Povirk (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8350–8354.
39. L. F. Povirk, Y.-H. Han, and R. J. Steighner (1989) *Biochem.* **28**, 8508–8514.
40. R. M. Burger, A. R. Berkowitz, J. Peisach, and S. B. Horwitz (1980) *J. Biol. Chem.* **255**, 11832–11838.
41. S. Uesugi et al. (1984) *Nucleic Acids Res.* **12**, 1581–1592.
42. R. M. Burger, J. Peisach, and S. B. Horwitz (1982) *J. Biol. Chem.* **257**, 8612–8614.
43. M. T. Kuo and T. C. Hsu (1978) *Nature* **271**, 83–84.

44. M. T. Kuo and T. C. Hsu (1978) *Chromosoma* **68**, 229–240.
45. C. W. Moore, (1988) *Cancer Res.* **48**, 6837–6843.
46. C. W. Moore, C. S. Jones, and L. A. Wall (1989) *Antimicrob. Agents Chemother.* **33**, 1592–1599.
47. K. Sidik and M. J. Smerdon (1990) *Biochem.* **29**, 7501–7511.
48. C. W. Moore (1982) *Cancer Res.* **42**, 929–933.
49. C. W. Moore (1989) *Cancer Res.* **49**, 6935–6940.
50. J. F. Mc.Koy et al. (1995) *Mutat. Res. DNA Repair* **336**, 19–27.
51. L. F. Povirk, M. Hogan, N. Dattagupta, and M. Buechner (1981) *Biochem.* **20**, 665–671.
52. C.-H. Huang, C. K. Mirabelli, Y. Jan, and S. T. Crooke (1981) *Biochem.* **20**, 233–238.
53. C.-H. Huang, A. W. Prestayko, and S. T. Crooke (1982) *Biochem.* **21**, 3704–3710.
54. C.-H. Huang, C. K. Mirabelli, S. Mong, and S. T. Crooke (1983) *Cancer Res.* **43**, 2849–2856.
55. M. Takeshita et al. (1981) *Biochem.* **20**, 7599–7606.
56. J. Kross, W. D. Henner, S. M. Hecht, and W. A. Haseltine (1982) *Biochem.* **21**, 4310–4318.
57. H. Suzuki et al. (1969) *J. Antibiot.* **22**, 446–448.
58. R. Stern, J. A. Rose, and R. M. Friedman (1974) *Biochemistry* **13**, 307–312.
59. B. J. Carter et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 9373–9377.
60. A. Huttenhofer, S. Hudson, H. F. Noller, and P. K. Mascharak (1992) *J. Biol. Chem.* **267**, 24471–24475.
61. L. L. Guan et al. (1993) *Biochem. Biophys. Res. Commun.* **191**, 1338–1346.
62. M. V. Keck and S. M. Hecht (1995) *Biochemistry* **34**, 12029–12037.
63. C. W. Moore (1982) *J. Bacteriol.* **150**, 1227–1233.
64. C. W. Moore and J. B. Little (1985) *Cancer Res.* **45**, 1982–1986.
65. C. W. Moore (1988) *J. Bacteriol.* **170**, 4991–4994.
66. A. W. Johnson and B. Demple (1988) *J. Biol. Chem.* **263**, 18017–18022.
67. J. Laval (1996) *Pathol. Biol.* **44**, 14–24.
68. C. W. Moore, A. W. Malcolm, K. N. Tomkinson, and J. B. Little (1985) *Cancer Res.* **45**, 1978–1981.
69. J. Game (1993) *Semin. Cancer Biol.* **4**, 73–83.
70. S. Prakash, P. Sung, and L. Prakash (1993) *Annu. Rev. Genet.* **27**, 33–70.
71. E. C. Friedberg (1995) *DNA Genes and Mutagenesis*. ASM Press, Washington, D.C.
72. D. Ramotar and J.-Y. Masson (1996) *Mol. Cell. Biochem.* **158**, 65–75.
73. C. W. Moore (1978) *Mutat. Res.* **51**, 165–180.
74. C. W. Moore (1980) *J. Antibiot.* **33**, 1369–1375.
75. C. W. Moore (1982) *Antimicrob. Agents Chemother.* **21**, 595–600.
76. C. W. Moore (1991) *J. Bacteriol.* **173**, 3605–3608.
77. D. J. Keszenman, V. A. Salvo, and E. Nunes (1992) *J. Bacteriol.* **174**, 3125–3132.
78. C. H. He, J.-Y. Masson, and D. Ramotar (1996) *Can. J. Microbiol.* **42**, 1263–1266.
79. M. A. Hannan and A. Nasim (1978) *Mutat. Res.* **53**, 309–316.
80. M. A. Hannan, A. Nasim, and T. Brychcy (1978) *Mutat. Res.* **58**, 107–110.
81. B. K. Vig and R. Lewis (1978) *Mutat. Res.* **55**, 121–145.
82. C. W. Moore (1978) *Mutat. Res.* **58**, 41–49.
83. A. Severgnini, O. Lillo, and E. Nunes (1991) *Environ. Mol. Mutagen.* **18**, 102–106.
84. L. F. Povirk and M. J. F. Austin (1991) *Mutat. Res.* **257**, 127–143.
85. L. F. Povirk et al. (1994) *J. Mol. Biol.* **243**, 216–226.



86. L. R. Ferguson and P. M. Turner (1988) *Eur. J. Cancer Clin. Oncol.* **24**, 591–596.
87. G. Pron, J. Belehradec, Jr., S. Orłowski, and L. M. Mir (1994) *Biochem. Pharmacol.* **48**, 301–310.
88. L. M. Mir, O. Tounekti, and S. Orłowski (1996) *Gen. Pharmacol.* **27**, 745–748.
89. C. W. Moore, R. Del Valle, J. F. Mc.Koy, A. Pramanik, and R. E. Gordon (1992) *Antimicrob. Agents Chemother.* **36**, 2497–2505.
90. R. Beaudouin et al. (1993) *Antimicrob. Agents Chemother.* **37**, 1264–1269.
91. S. T. Lim, C. K. Jue, C. W. Moore, and P. N. Lipke (1995) *J. Bacteriol.* **177**, 3534–3539.
92. S. Akiyama et al. (1981) *Biochem. Biophys. Res. Commun.* **101**, 55–60.
93. J. S. Lazo, S. M. Sebti, and A. E. Filderman (1987) In *Metabolism and Mechanism of Action of Anti-cancer Drugs* (G. Powis and R. A. Prough, eds.), Taylor and Francis, London, pp. 194–210.
94. C. Nishimura, H. Suzuki, N. Tanaka, and H. Yamaguchi (1989) *Biochim. Biophys. Acta* **1012**, 29–35.
95. S. M. Sebti et al. (1989) *Biochemistry* **28**, 6544–6548.
96. S. E. Montoya, R. E. Ferrell, and J. S. Lazo (1997) *Cancer Res.* **57**, 4191–4195.
97. N. G. Kambouris, D. J. Burke, and C. E. Creutz (1992) *J. Biol. Chem.* **267**, 21570–21576.
98. C. Enenkel and D. H. Wolf (1993) *J. Biol. Chem.* **268**, 7036–7043.
99. U. Magdolen, G. Müller, V. Magdolen, and W. Bandlow (1993) *Biochim. Biophys. Acta* **1171**, 299–303.
100. H. E. Xu and S. A. Johnston (1994) *J. Biol. Chem.* **269**, 21177–21183.
101. A. Pei, T. P. Calmels, C. E. Creutz, and S. M. Sebti (1995) *Mol. Pharmacol.* **48**, 676–681.
102. L. Joshua-Tor, H. E. Xu, S. A. Johnson, and D. C. Rees (1995) *Science* **269**, 945–950.
103. W. Zheng, H. E. Xu, and S. A. Johnston (1997) *J. Biol. Chem.* **272**, 30350–30355.
104. W. Zheng and S. A. Johnston (1998) *Mol. Cell. Biol.* **18**, 3580–3585.
105. I. Niemer, G. Müller, and G. Strobel (1997) *Curr. Genet.* **32**, 41–51.
106. R. P. Koldamova, I. M. Lefterov, V. G. Gadjeva, and J. S. Lazo (1998) *Biochemistry* **37**, 2282–2290.
107. D. Drocourt, T. Calmels, J. P. Reynes, M. Baron, and G. Tiraby (1990) *Nucleic Acids Res.* **18**, 4009.
108. D. Bromme et al. (1996) *Biochemistry* **35**, 6706–6714.
109. A. A. Ferrando, A. Velasco, E. Campo, and C. Lopez-Otin (1996) *Cancer Res.* **56**, 1746–1750.
110. A. Gatignol, M. Baron, and G. Tiraby (1987) *Mol. Gen. Genet.* **207**, 342–348.
111. A. Gatignol, M. Dassain, and G. Tiraby (1990) *Gene* **91**, 35–41.
112. D. Drocourt et al. (1990) *Nucl. Acids Res.* **18**, 4009.
113. P. L. Tran et al. (1997) *J. Clin. Invest.* **99**, 608–617.
114. J. Weinbach et al. (1996) *Cancer Res.* **56**, 5659–5665.
115. C. W. Moore and D. A. Vossler (1980) *Biochim. Biophys. Acta* **610**, 425–429.
116. C. W. Moore (1994) *Antimicrob. Agents Chemother.* **38**, 1615–1619.
117. T. Takita et al. (1972) *J. Antibiot.* **25**, 755–758.
118. H. Naganawa, Y. Muraoka, T. Takita, and H. Umezawa (1977) *J. Antibiot.* **30**, 388–396.
119. T. Takita et al. (1978) *J. Antibiot.* **31**, 801–804.

### Suggestions for Further Reading

120. R. M. Burger (1998) Cleavage of nucleic acids by bleomycin. *Chem. Rev.* **98**, 1153–1170. (Reviews the chemical mechanism of bleomycin action.)

121. S. M. Hecht (1994) RNA degradation by bleomycin, a naturally occurring bioconjugate. *Bioconjugate Chem.* **5**, 513–526. (Reviews and compares mechanisms of cleavage of DNAs and RNAs.)
122. J. S. Lazo and S. M. Sebt (1997) Bleomycin. *Cancer Chemother. Biol. Response Modif.* **17**, 40–45. (Together with earlier reviews in this series by the same authors, this article provides overview of clinical uses of bleomycin and mechanisms of tumor resistance to bleomycin cytotoxicity.)
123. D. H. Petering, Q. Mao, W. Li, E. DeRose, and W. E. Antholine (1996) Metallobleomycin–DNA interactions: structures and reactions related to bleomycin-induced DNA damage. *Met. Ions Biol. Syst.* **33**, 619–648. (Reviews mechanisms of interactions of metallobleomycin and DNA in cells and *in vitro*.)
124. L. F. Povirk (1996) DNA damage and mutagenesis by radiomimetic DNA-cleaving agents: bleomycin, neocarzinostatin and other enediynes *Mutat. Res.* **355**, 71–89. (Summarizes the chemical action and relationship of DNA lesions to types of mutations produced by radiomimetic antibiotics.)

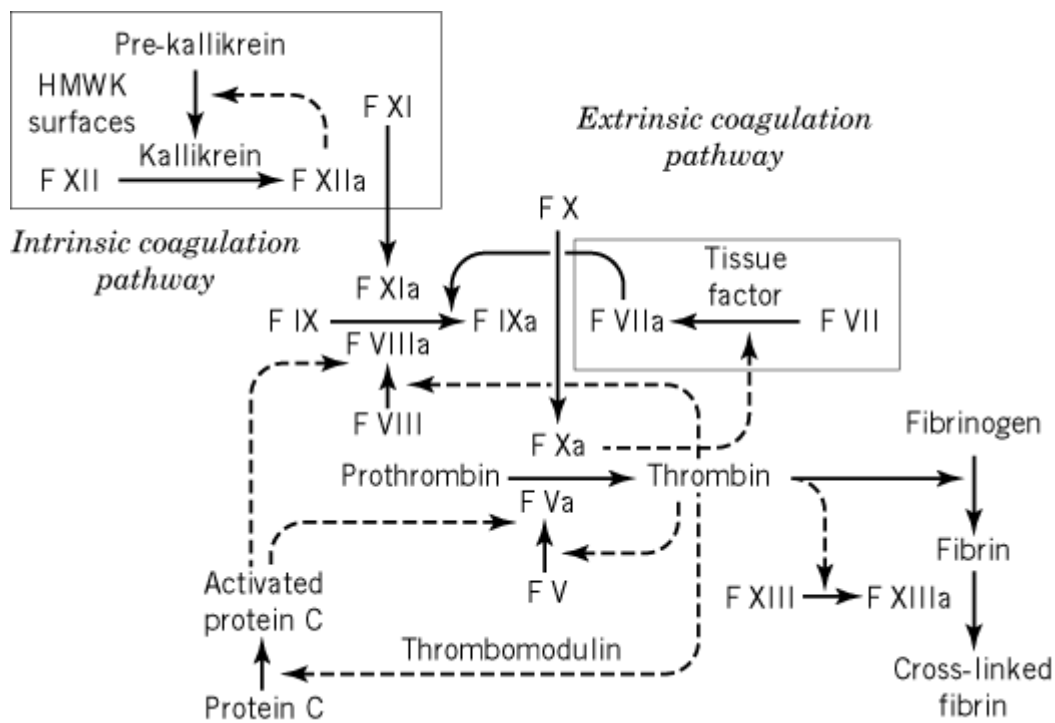
## Blood Clotting

The clotting or coagulation of blood is necessary for the maintenance of vascular integrity, and is therefore essential to survival. Cessation of bleeding can be achieved within minutes, involving both cellular and soluble protein components present in blood. The overall function of the blood coagulation system is the rapid amplification of a small initial stimulus via a series of sequential reactions that activate various [zymogens](#) and pro-cofactors of [serine proteinases](#), culminating in the generation of [thrombin](#), which proteolytically converts [fibrinogen](#) to fibrin. Fibrin monomers subsequently polymerize to form the insoluble fibrin clot. The major enzymatic components of the coagulation pathways are a series of **homologous** serine proteinases whose function is determined primarily by differences in the organization of their *N* terminal, noncatalytic **domains**. Most of the reactions are catalyzed by complexes between these serine proteinases, nonenzymatic protein cofactors, and serine proteinase zymogens that are assembled on phospholipid surfaces in the presence of calcium ions. These complexes are able to catalyze activation of the proteinase zymogen by limited **proteolysis** at rates up to  $10^6$ -fold greater than those in the absence of cofactor/phospholipid. The proteinase product of one reaction then participates in another analogous complex, thus providing the amplification, as well as the multiple levels at which both positive and negative regulation can be exerted.

Figure 1 shows the two classical pathways of blood coagulation, termed the *intrinsic* and *extrinsic* pathways. These were originally envisaged as elegant, essentially linear sequences of reactions, but **feedback inhibition** reactions and cross-reactions between the two pathways are now known to be important and to provide further amplification. The terminology used to describe the two pathways of blood coagulation is somewhat counterintuitive, tending to imply the opposite to what is now known to be the physiological clotting mechanism. Genetic deficiencies in the extrinsic pathway lead to severe bleeding tendencies, but initiation of this pathway *in vitro* requires components in the form of tissue extracts (containing tissue factor and phospholipid) that are not present in blood. By contrast, the physiological role of the “intrinsic” pathway is unclear, as genetic deficiencies in its components do not give rise to bleeding disorders, but all the components are present in blood, and it can be initiated by simply removing blood from the circulation. This is due to contact with a negatively charged surface, such as glass, giving the pathway its alternative name of “contact

activation,” and involves proteinases and cofactors quite distinct from those of the extrinsic pathway.

**Figure 1.** The proteolytic reactions of the blood clotting pathways. The zymogen/proteinase pairs are shown in black with the “forward” reactions as solid lines, the cofactors are shown in blue with the positive feedback reactions as dashed blue lines. The proteolytic negative feedback reactions are shown in red. The complexes initiating the two pathways of coagulation are boxed: the factor VII–tissue factor complex and the contact activation complex that assembles on high molecular-weight kininogen (HMWK). The Roman numeral system used for the coagulation proteinases and cofactors is based loosely on the order in which they were discovered, with the exceptions of the first four components fibrinogen, prothrombin, tissue factor–phospholipid complex, and  $\text{Ca}^{2+}$ . The proteolytically activated forms of these factors are denoted by a subscript “a,” for example, factor  $\text{X}_a$ . Additional nonproteolytic inhibitory mechanisms exist, the principal of which are the [serpin](#) antithrombin (an inhibitor of most of the Gla-proteinases) and tissue factor pathway inhibitor (TFPI), a Kunitz-type [proteinase inhibitor](#) of factors  $\text{VII}_a$  and  $\text{X}_a$ .



A complete description of the coagulation system and the many proteolytic and inhibitory reactions that lead to its exquisite functional regulation is beyond the scope of this article. We will therefore focus on two key proteolytic steps that are representative of the two types of multicomponent complexes of the extrinsic pathway and that demonstrate the fundamental interactions and mechanisms underlying catalysis of the proteolytic reactions of blood coagulation; one requiring proteolytically activated soluble cofactors, the other requiring integral [membrane protein](#) cofactors that do not require proteolytic activation.

### 1. Regulation of Blood Clotting by Complexes Involving Soluble Protein Cofactors

The only proteolytic reaction common to both pathways is the activation of prothrombin, the final serine proteinase zymogen of the coagulation cascade. Factor  $\text{X}_a$  (i.e., activated factor X) is the proteinase component of the “prothrombinase complex” that catalyzes the generation of thrombin. Both zymogen and active proteinase have an *N*-terminal “Gla domain,” which contains [g-carboxyglutamic acid](#) residues. The Gla domain is also present in prothrombin (as well as factors VII and IX, and protein C, the so-called vitamin K–dependent or Gla proteinases); it binds  $\text{Ca}^{2+}$  and participates in binding to negatively charged phospholipid surfaces, which *in vivo* are provided

principally by activated platelets. Therefore, both [enzyme](#) and substrate bind to phospholipid surfaces, with **dissociation constants** of 0.1–1  $\mu\text{M}$  (1). The cofactor of the prothrombinase complex is factor  $V_a$ , which is generated by thrombin-catalyzed limited proteolysis of factor V, providing one of the major feedback activation reactions. Factor  $V_a$  also binds to negatively-charged phospholipid surfaces, but with a much lower dissociation constant of  $\sim 5$  nM (2). It does not contain a Gla domain, and binds through a different mechanism involving partial insertion of the protein into the lipid bilayer; it is often viewed as a membrane binding site for factor  $X_a$ . The latter thus incorporates into the complex by both protein–phospholipid and [protein–protein interactions](#). Kinetic experiments suggest that the two phospholipid-bound components associate by lateral [diffusion](#), but the interactions involved are complex and to some extent **cooperative**. Although prothrombin also binds to phospholipid, it is not clear whether the relevant substrate is phospholipid-bound prothrombin or prothrombin in the solution phase but at a high local concentration in the vicinity of the membrane surface, which can be several orders of magnitude above that in bulk solution. The kinetic effect of prothrombin binding to phospholipid is a reduction in the apparent [Km \(Michaelis constant\)](#) for prothrombin activation (“apparent” as the  $K_m$  does not reflect the true concentration at the phospholipid–surface interface). This is of physiological significance, since the  $K_m$  falls from  $\sim 100$   $\mu\text{M}$  to  $< 1$   $\mu\text{M}$  (3), compared to a concentration of 1.4  $\mu\text{M}$  in blood plasma. In addition to this effect, the catalytic rate constant,  $k_{\text{cat}}$ , for prothrombin activation is greatly increased, primarily as a consequence of the  $X_a$ – $V_a$  interaction (Table 1). This suggests that the active-site environment of the proteinase is significantly affected by complex formation, although the molecular basis of this effect is not known. The prothrombinase complex also helps to achieve maximum catalytic advantage by ordering the sequence of the two proteolytic cleavages required for prothrombin activation (see [Thrombin](#)).

**Table 1. Effect of Various Components of Prothrombinase Complex on Kinetic Parameters for Prothrombin Activation<sup>o</sup>**

| Components                            | $K_m, \mu\text{M}$ | $k_{\text{cat}}, \text{s}^{-1}$ | Relative Rates |
|---------------------------------------|--------------------|---------------------------------|----------------|
| $X_a \text{ Ca}^{2+} \text{ PL } V_a$ | 0.2                | 30                              | 100            |
| $X_a \text{ Ca}^{2+} \text{ PL}$      | 0.06               | 0.04                            | 0.42           |
| $X_a \text{ Ca}^{2+} V_a$             | 34                 | 6.2                             | 0.12           |
| $X_a \text{ Ca}^{2+}$                 | 84                 | 0.01                            | 0.00009        |
| $X_a$                                 | 131                | 0.01                            | 0.00005        |

Source: Data taken from Rosing et al. (3).

The phospholipid-bound complex of factor  $IX_a$  (proteinase) and factor  $VIII_a$  (cofactor) that activates factor X has molecular and functional characteristics very similar to those of the prothrombinase complex. Factor VIII shares homology with factor V and is also activated by thrombin; factor IX has the same domain structure as factor X and is also activated by factor  $VII_a$ . There is evidence that the phospholipid-binding properties of these proteinase–zymogen components enable the proteinase

product of one complex to be transferred to the subsequent complex as its enzyme without dissociation from the membrane, thereby giving a further catalytic advantage (4).

## 2. Regulation of Blood Clotting by Complexes Involving Integral Membrane Proteins

Two integral membrane proteins with quite distinct and contrasting functions are involved in the regulation of blood coagulation: tissue factor and thrombomodulin, which are involved, respectively, in the initiation and termination of coagulation. Blood coagulation subsequent to injury of the vasculature is initiated by the exposure of tissue factor, a nonenzymatic cofactor widely expressed by cells of the subendothelium of blood vessels (although its expression can also be induced in endothelial cells and monocytes). Tissue factor binds factor VII<sub>a</sub>, a Gla-proteinase, with high affinity ( $K_d = 1-5 \text{ nM}$ ) (5). The structure of this complex has been determined by [X-ray crystallography](#) (6); the extracellular part of tissue factor is composed of two homologous C2-type **immunoglobulin**-like modules, and the protein is most closely related to the IFN $\gamma$  receptor, a member of the class 2 **cytokine** receptors (see [Interferons](#)). Factor VII<sub>a</sub> is engaged in multiple contacts with both of the tissue factor modules.

The classical substrate for factor VII<sub>a</sub> in this complex is factor X, but it is now known that factor IX is also a physiologically relevant substrate (7) (see [Hemophilia](#)). Activation of these substrates, and the subsequent incorporation of their activated forms into their cognate complexes, thus propagates the initial stimulus proteolytically. In contrast to factors V and VIII, tissue factor is constitutively active as a cofactor and does not require proteolytic processing. However, factor VII does need proteolytic activation for the complex to be active, and what initially catalyzes this cleavage and thus constitutes the initial proteolytic event in the coagulation cascade has not been a trivial problem to address, because of the extraordinary sensitivity of the system to proteolytic activation, and it remains a topic of controversy. Proposals have included (1) that factor VII is not a true zymogen and possesses a significant degree of intrinsic proteolytic activity, (2) induction of activity in factor VII by tissue factor, (3) autoactivation of factor VII, (4) and activation of factor VII by trace amounts of activated proteinases. Although not demonstrated unambiguously, mechanism 4 appears to be favored, as factor VII can be activated by factors VII<sub>a</sub>, IX<sub>a</sub>, X<sub>a</sub>, and thrombin *in vitro*. These may be generated by the contact activation system, alternative mechanisms, or “leakage” of the system, but in such low amounts that they are insignificant in the absence of tissue factor. The principal activator of factor VII once the system is initiated is considered to be factor X<sub>a</sub> (8). Tissue factor has little or no effect on factor VII activation, but assembly of the tissue factor complex increases factors IX and X activation by up to 10<sup>4</sup>-fold (9). Both the kinetic mechanisms and the molecular interactions responsible for this effect are analogous to those described for the prothrombinase complex, which involves the optimal presentation of the zymogen substrate to the catalytically efficient proteinase:cofactor complex for maximum catalytic advantage.

The other integral membrane protein of the coagulation system, thrombomodulin, is expressed by endothelial cells and acts as a cofactor for the proteolytic activation of the negative regulator of the coagulation system, protein C (10). Activated protein C (a Gla-proteinase) proteolytically inactivates the procoagulant cofactors factors V<sub>a</sub> and VIII<sub>a</sub>. In common with the principal feedback activation reactions, thrombin is again involved. But in this case, rather than proteolytically activating the cofactor, it is the enzymatic component of the thrombomodulin complex; thus thrombin can fully modulate its own generation. The interaction with thrombomodulin dramatically alters the substrate specificity of thrombin, such that protein C activation is increased 10<sup>4</sup>-fold, while its activity toward factor V and fibrinogen is diminished (1). Although protein C can bind phospholipid through its Gla domain, this interaction plays only a minor part in these effects and, because thrombin lacks a Gla domain, protein-protein interactions are the major determinants in thrombomodulin complex anticoagulant function, in contrast to the procoagulant complexes described above.

## Bibliography

1. C. T. Esmon, N. L. Esmon, and K. W. Harris (1982) *J. Biol. Chem.* **257**, 7944–7947.
2. S. Krishnaswamy and K. G. Mann (1988) *J. Biol. Chem.* **263**, 5714–5720.
3. J. Rosing, G. Tans, J. W. Govers-Riemslog, R. F. Zwaal, and H. C. Hemker (1980) *J. Biol. Chem.* **255**, 274–283.
4. K. G. Mann, M. E. Nesheim, W. R. Church, P. Haley, and S. Krishnaswamy (1990) *Blood* **76**, 1–16.
5. R. F. Kelley, K. E. Costas, M. P. O'Connell, and R. A. Lazarus (1995) *Biochemistry* **34**, 10383–10392.
6. D. W. Banner, A. D'Arcy, C. Chene, F. K. Winkler, A. Guha, W. H. Konigsberg, Y. Nemerson, and D. Kirchhofer (1996) *Nature*, **380**, 41–46.
7. K. A. Bauer, P. M. Mannucci, A. Gringeri, F. Tradati, S. Barzegar, B. L. Kass, H. ten Cate, A. S. Kestin, D. B. Brettler, and R. D. Rosenberg (1992) *Blood* **79**, 2039–2047.
8. S. Butenas and K. G. Mann, (1996) *Biochemistry* **35**, 1904–1910.
9. M. Zur and Y. Nemerson (1980) *J. Biol. Chem.* **255**, 5703–5711.
10. C. T. Esmon (1995) *FASEB J.* **9**, 946–955.

## Suggestion for Further Reading

11. K. G. Mann, R. J. Jenny, and S. Krishnaswamy (1988) Cofactor proteins in the assembly and expression of blood clotting enzyme complexes. *Annu. Rev. Biochem.* **57**, 915–956.

## Blot Overlays

Regardless of whether a blot contains immobilized **DNA**, **RNA** or **protein** (see [Blotting](#)), it is ultimately reacted with a specific probe that defines the type of **ligand-binding** interaction being studied (see [Southern Blots \(DNA Blots\)](#), [RNA Blots \(Northern Blots\)](#), and [Protein Blots \(Western Blots\)](#)). **Radiolabeled** probes are detected by [autoradiography](#), and the radioactive complexes are quantified by computer-digitized imaging or simply by excising out the bands and “counting” the radioactivity. Often the primary probe itself is not easily detectable, so a secondary probe is used, such as an [antibody](#) against the first probe conjugated to an [enzyme](#), or [avidin](#) or [streptavidin](#) when the primary probe is biotinylated (see [Avidin-Biotin System](#)). When using an enzyme-conjugated probe, one must ensure that the quencher or blocking reagent used on the blot has no interfering enzymatic activity. For example, [hemoglobin](#) should not be used when horseradish [peroxidase](#) is employed as the detection system.

### 1. DNA Probes

Traditionally, Southern and RNA blots of DNA and RNA, respectively, are probed with single-stranded DNA or double-strand DNA that has been **denatured**. The probes are usually radioactive, so that any duplex formed upon **hybridization** of the probe to a nucleic acid on the blot is detected by autoradiography. The rationale of such experiments is to identify the position of the immobilized nucleic acid in the electrophoretogram that is complementary to the sequence of the probe, so that they hybridize by **Watson-Crick base pairing**. The conditions for hybridization, the *stringency*, are regulated according to the degree of homology and complementarity between the probe and the target ([1](#), [2](#)). DNA is also used to probe protein blots, a procedure that has been named “Southwestern blotting.” As expected, this type of experiment is designed to reveal specific

interactions between defined DNA sequences and the [DNA-binding proteins](#) that bind them, such as [transcription factors](#) (3). DNA probing of *colony* or *plaque* blots is a routine approach in recombinant DNA [cloning](#) (4).

## 2. Immunoblotting

The original application of protein blots was for identifying a particular [antigen](#) in a protein gel pattern by probing the blot with its corresponding [antibody](#) (5). Detecting the immunocomplex is often accomplished by using a secondary probe, for example, a goat antimouse IgG conjugated to an enzyme, such as horseradish peroxidase, when a murine [monoclonal antibody](#) is used for the primary probe. The sensitivity of these assays is increased by combining **immunoprecipitation** prior to [gel electrophoresis](#). Thus, a crude sample of proteins is first immunoprecipitated with relevant antibodies, and the precipitated proteins are subsequently resolved by electrophoresis, blotted, and probed with a monoclonal antibody of interest. The use of antibodies to probe colony or plaque blots is very effective in screening [expression libraries](#) (6).

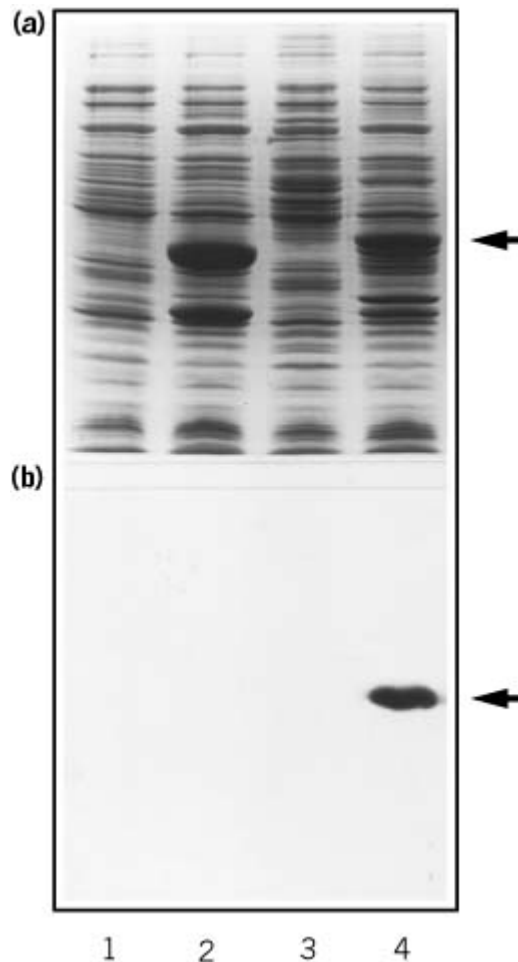
## 3. Lectin Blotting

To identify **glycoproteins**, [lectins](#) can be used to probe protein blots (7). A radioactive or enzyme-conjugated lectin is employed as the probe. A particular case of interest is the fact that the enzyme horseradish peroxidase is itself a glycoprotein. Thus, for example, a blot can be probed with the mannose-specific lectin concanavalin *A*, washed and further incubated with horseradish peroxidase directly (8). The multivalent concanavalin A binds to blotted glycoproteins containing mannose and subsequently also binds to the horseradish peroxidase without the need for chemical conjugation. When using lectins, ensure that the quencher used does not itself contain sugar.

## 4. Ligand Blotting

Blotting is usually considered for detecting [antigens](#) or nucleic acid hybrids, but it is also a very powerful way to identify all sorts of protein–protein interactions, even interactions involving nonpeptide ligands. Thus when studying any **receptor**, one should consider probing protein blots with any corresponding ligands that are detectable (9). In such experiments, it is usually advisable not to boil the protein sample or subject it to [disulfide bond](#) reduction before electrophoresis. Furthermore, the blot can be incubated in solutions that promote [renaturation](#) of the immobilized proteins. Omission of methanol from the transfer buffers used in electroblotting is also helpful for retaining functional conformations of the blotted protein. Ligand blotting has been used successfully to identify peptides that bind to [hormones](#), **cytoskeletal** components, [neurotoxins](#), nucleotides, [calmodulin](#), and even ions such as  $\text{Ca}^{2+}$  (9-11) (Fig. 1).

**Figure 1.** Neurotoxin overlays of protein blots. *Escherichia coli* transformed with pATH2 expression **vectors** containing DNA inserts corresponding to the neurotoxin binding site (lanes 3 and 4) or not (lanes 1 and 2) were induced for gene expression (lanes 2 and 4) or not (lanes 1 and 3). Then samples of all four treatments were resolved by SDS-PAGE and stained with [Coomassie brilliant blue](#) (a) or blotted onto nitrocellulose membrane filters and probed using radioactive neurotoxin. Then the washed filters were autoradiographed (b). The arrows indicate the position of the fusion protein containing the toxin binding site.



## 5. Cell Blotting

Blots have even been probed with intact cells, which has proven useful for identifying interactions between proteins and cells (12). Furthermore, bacteria can be used to probe a blot, and the interaction with a protein is detected by subsequently allowing the bacteria to grow on the surface of the blot, so that colonies are observed at the site of the immobilized protein (13). **Viruses** have also been used to probe blots to detect their corresponding receptor proteins.

In summary, one should consider probing blots with selected ligands any time bimolecular interactions are to be analyzed. Surveying literature [databases](#) using specific conjugations, such as “ligand-blot”, “calmodulin-blot” or “cell-blot,” as key words usually produces results that provide the imaginative investigator with starting points from which to proceed.

### Bibliography

1. J. Meinkoth and G. Wahl (1984) *Anal. Biochem.* **138**, 267–284.
2. L. G. Davis, M. D. Dibner, and J. F. Battey (1986) *Basic Methods in Molecular Biology*, Elsevier, New York, pp. 62–65.
3. C. N. Flytzanis (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 163–168.
4. Ref. 2 pp. 227–229.
5. H. Towbin and J. Gordon (1984) *J. Immunol. Methods* **72**, 313–340.
6. T. V. Huynh, R. A. Young, and R. W. Davis (1986) In *DNA Cloning, A Practical Approach*,



- Vol. 1 (D. M. Glover, ed.), IRL Press, Oxford, UK, pp. 49–78.
7. S. Bar-Nun and J. M. Gershoni (1994) In *Cell Biology; a laboratory handbook*, Vol 3 (J. E. Celis, ed.), Academic Press, San Diego, pp. 323–331.
  8. J. C. S. Clegg (1982) *Anal. Biochem.* **127**, 389–394.
  9. J. M. Gershoni (1988) *Methods of Biochemical Analysis* **33**, 1–58.
  10. P. Hossenlopp and M. Binoux (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 169–188.
  11. A. Vieira, R. G. Elkin, and K. Kuchler (1994) In *Cell Biology: A laboratory handbook*, Vol. 2 (J. E. Celis, ed.), Academic Press, San Diego, pp. 314–321.
  12. E. G. Hayman, E. Engvall, E. A'Hearn, D. Barnes, M. Pierschbacher, and E. Rouslahti (1982) *J. Cell Biol.* **95**, 20–23.
  13. J. M. Gershoni (1987) Protein blotting: a tool for the analytical biochemist, *Adv. Electrophoresis* **1**, 141–176.

## Blotting

Blotting is a method in which a [macromolecule](#) is immobilized on a blotting matrix and subsequently probed with a detectable ligand to determine whether the macromolecule binds that specific **ligand**. The immobilized macromolecule can be **DNA**, **RNA** or protein, in which case one generates DNA blots ([Southern blots](#)), [RNA blots \(Northern blots\)](#) (1), or [protein blots \(Western blots\)](#) (2, 3). Blots of lipids have also been produced (4). The macromolecule can be applied to the blotting matrix directly (dot blot), or it can be derived and eluted from an electrophoretic gel (gel blot) or even from a **bacterial** colony or **bacteriophage** plaque (colony blot).

### 1. Typical Blot Analysis

The most common application of blotting involves a complex mixture of DNA, RNA, or protein to be resolved by using a standard [gel electrophoresis](#) procedure, such as [agarose](#) gel separation of DNA fragments or RNA or [SDS-PAGE](#) of protein samples. After electrophoresis, the gel is dismantled from its cassette or glass plates, etc., and a sheet of an appropriate blotting matrix (eg, a [nitrocellulose](#) membrane filter; see [Blotting Matrices](#)) is cut to size and applied to the surface of the gel. The transfer of the resolved polynucleotides or peptides is accomplished via a procedure called “blotting”, and the blotted macromolecules adsorb to the surface of the matrix while retaining their relative positions, thus creating a faithful replica of the original electrophoretic pattern. The “blot” thus produced is subsequently incubated with a ligand probe, which might be **radioactive** for detection via [autoradiography](#) or conjugated to an [enzyme](#) whose activity is detectable (see [Blot Overlays](#)). Extensive washing of the blot removes the excess probe from that which is specifically associated with the immobilized macromolecule and remains bound. Subsequent detection of the retained ligand in the complex formed identifies the relevant bimolecular interaction.

### 2. Methods for Blotting

#### 2.1. Dot Blotting

This is the simplest method of applying a sample to be tested (5, 6). The sample can be an unfractionated polynucleotide or protein mixture in solution. A small volume (typically 2 to 5 µl) is applied directly to the surface of a dry blotting matrix by micropipetting the sample onto the matrix or by using commercial vacuum manifolds that enable the application of larger sample volumes

(e.g., 100 µl to 1 ml) by filtration. Such manifolds often create focused and uniform dots or thin slots of sample, thus leading to the terms “dot blots” or “slot blots”, respectively. Typically, an 8cm × 12 cm piece of blotting matrix contains as many as 96 different dots that correspond to the geometry of a standard 96-well **ELISA** plate. Once created, the dot blot is processed and probed as any other blot (see [Southern Blots \(DNA Blots\)](#), [RNA Blots \(Northern Blots\)](#), [Protein Blots \(Western Blots\)](#), and [Blot Overlays](#)). The advantages of the dot blot procedures are that they do not require any separation process and thus do not subject the sample to undue chemical modifications that could, for example, **denature** a protein sample (although some denaturation of protein occurs upon adsorption to the matrix). Moreover, dot blots are simple, cheap and quick. Where quantification is intended, direct dot blotting ensures maximal yields of sample recovery. Obviously, however, chromatography or electrophoresis is necessary to resolve a complex sample to ascribe the signal to a specific component.

## 2.2. Gel Blots

Gel electrophoresis is routinely employed to resolve macromolecules of DNA, RNA, or protein. Agarose gels and **polyacrylamide gels** can be blotted, and the common goal is to elute the “bands” efficiently from the gel to be immobilized on the surface of the blotting matrix, so as to generate a faithful replica of the electrophoretic pattern. This can be accomplished in a number of ways:

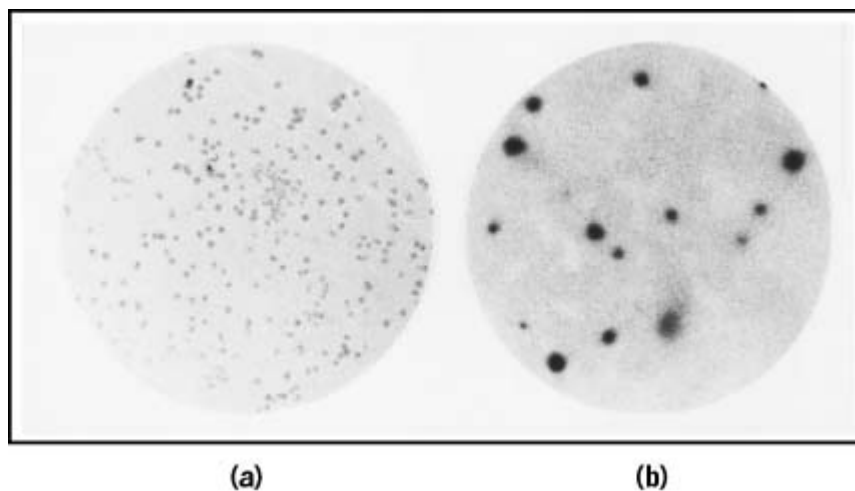
1. Diffusion blotting simply relies on the fact that the macromolecules in the gel spontaneously **diffuse** out of the gel (7). Consequently, blots are produced when a blotting matrix is simply applied to one or both sides of the gel. This approach is usually time-consuming (24 to 72 h) and of low efficiency, but it produces two equal copies simultaneously.
2. Convection blotting (also called capillary transfer) is the process of eluting the resolved macromolecules by mass flow of buffer through the gel. This is the original procedure introduced by Edwin Southern for transferring DNA restriction fragments from agarose gels onto nitrocellulose membrane filters (thus the term Southern blotting) (8). In this method, the gel is placed on top of a paper wick, which draws buffer from a reservoir. The gel is covered with a piece of blotting matrix that, in turn, is covered by a stack of paper towels or absorbent paper and a weight that ensures uniform pressure over the surface of the gel. The stack of paper towels draws a continuous flow of buffer vectorially through the gel, and with it the polynucleotides or peptides, which are eluted and deposited on the surface of the blotting matrix (9).
3. Vacuum blotting is an elaboration of blotting by convection in which the flow is accelerated by employing negative pressure (i.e., suction). Positive pressure has also been used to enhance blotting.
4. Electroblotting is achieved by applying an electric field so as to elute the proteins or polynucleotides from their corresponding gels by [electrophoresis](#). This technique has been used primarily for proteins because binding DNA and RNA to nitrocellulose requires high salt conditions that are incompatible with electroblotting. Electroblotting of nucleic acids is possible, however, with alternative blotting matrices, such as nylon membranes. A variety of commercial apparatus equipment and home-made systems are available, including those that generate gradient electric fields and provide different elution efficiencies to compensate for differences in the molecular mass of the polymers to be eluted (10, 11). Electroblotting is performed by using tank systems, which require 2 to 4 liters of transfer buffer, or semidry blotting systems, which conserve buffer and are flexible because different buffers can be employed for the anode and cathode (12).
5. Direct blotting is an elaboration of electroblotting in which a conveyor belt passes the blotting matrix by the exposed bottom of the polyacrylamide gel during electrophoresis. As the proteins or DNA fragments reach the bottom of the gel, they are deposited directly onto the slowly moving blotting matrix and generate a blot. The resolution of the bands is modulated by regulating the conveyor belt speed (13).

## 2.3. Colony Blots

At times it is necessary to identify a specific bacterial colony or phage that contains DNA of interest

or expresses a particular protein. [Libraries](#) of **recombinant DNA**-containing bacteria or [lambda phage expression libraries](#) are used to produce colony or plaque blots to be processed and screened like any other blot. The library of bacteria or phage is plated on agar after suitable incubation, and replicas are produced by simply placing a sheet of blotting matrix, such as nitrocellulose membrane or nylon membrane, onto the colony- or plaque-containing surface. The matrices pick up sufficient material from each colony or plaque at precisely the same relative position corresponding to that on the original agar plate. Then the colony blots are processed, for example, by standard DNA filter hybridization, immunoblotting, or ligand blotting ([14](#), [15](#)) (Fig. [1](#)).

**Figure 1.** Bungarotoxin overlay of bacterial colonies. *Escherichia coli* were transformed with a pATH2 expression vector containing a DNA insert coding for a fragment of the  $\alpha$  subunit of the nicotinic [acetylcholine receptor](#). This fragment is responsible for the receptor's binding of the neurotoxin  $\alpha$ -**bungarotoxin**. The bacteria were plated to a density of approximately 200 colonies per agar plate (**a**) and expression was induced. A replica of the plate was produced by using a nitrocellulose filter disc, which was processed for overlay with radio-**iodinated**  $\alpha$ -bungarotoxin. The filter was subsequently washed and autoradiographed, illustrating those colonies that contain the relevant DNA fragment (**b**) (for details see Ref [15](#)).



## Bibliography

1. J. Meinkoth and G. Wahl (1984) *Anal. Biochem.* **138**, 267–284.
2. J. M. Gershoni and G. E. Palade (1983) *Anal. Biochem.* **131**, 1–15.
3. H. Towbin and J. Gordon (1984) *J. Immunol. Methods* **72**, 313–340.
4. T. Taki, S. Handa, and D. Ishikawa (1994) *Anal. Biochem.* **221**, 312–316.
5. J. M. Gershoni (1988) *Methods Biochem. Anal.* **33**, 1–58.
6. L. G. Davis, M. D. Dibner, and J. F. Battey (1986) *Basic Methods in Molecular Biology*, Elsevier, New York, pp. 147–149.
7. B. Bowen, J. Steinberg, U. K. Laemmli, and H. Weintraub (1980) *Nucleic Acids Res.* **8**, 1–20.
8. E. M. Southern (1975) *J. Mol. Biol.* **98**, 503–517.
9. Ref. [6](#) pp. 62–65.
10. M. Bittner, P. Kupferer, and C. F. Morris (1980) *Anal. Biochem.* **102**, 459–471.
11. J. M. Gershoni, F. E. Davis, and G. E. Palade (1985) *Anal. Biochem.* **144**, 32–40.
12. G. Jacobson (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 53–72.
13. S. Beck (1993) *Methods Mol. Biol.* **23**, 219–223.
14. Ref. [6](#) pp. 185–189.

15. J. M. Gershoni (1987) Proc. Natl. Acad. Sci USA **84**, 4318–4321.

## Blotting Matrices

**Blotting** is the process of transferring macromolecules from **electrophoretic** gels to immobilizing matrices, called blotting matrices. The variety of matrices available for blotting is quite diverse, and the characteristics of each directly affects the ultimate result and quality in different blot analyses [see [Blotting](#), [Southern Blots \(DNA Blots\)](#), [Protein Blots \(Western Blots\)](#), [RNA Blots \(Northern Blots\)](#), and [Blot Overlays](#)]. Generally, two kinds of blot matrices are used: (1) chemically modified paper filters and (2) microporous membrane filters. Although filters are used, blotting depends on chemical adsorption of transferred molecules to the filter material itself, rather than filtration per se, where separation is achieved by size exclusion. Thus porosity is less important than the chemical composition, density, and thickness of the material. After binding of the desired ligand, normally the remaining binding sites on the matrix are blocked with a neutral compound, known as a quencher or blocking reagent.

### 1. Paper Filters

Initially, it was thought that transferred molecules would be best suited for blotting if they were covalently bound to the blotting matrix. This led to chemically modified paper filters that contained active moieties to bind DNA, RNA, and protein covalently. For example, **cyanogen bromide** (CNBr)-activated paper was produced (1), as was diazobenzyloxymethyl (DBM) paper, which was the most popular filter of this type (2). These filters covalently immobilize the blotted macromolecules, but they are cumbersome to handle, have the fibrous texture characteristic of blotting paper, and therefore are rarely used today.

### 2. Membrane Filters

Microporous membrane filters have become the matrices of choice for blotting. These materials are thin films of synthetic polymers with very fine, uniform surfaces. Filters with average porosities of 0.2 to 0.45  $\mu\text{m}$  are extremely suitable for all types of blotting. A variety of materials exist, and each offers unique advantages.

#### 2.1. Nitrocellulose

Nitrocellulose is by far the most commonly used blotting matrix. It binds DNA, RNA, and protein reasonably well, although the mechanism is not clear (3, 4). **Hydrophobic** interactions definitely play a role, and the salt conditions are critical, particularly for the adsorption of RNA. Nitrocellulose is often produced as a mixed ester with cellulose acetate, which reduces to some extent its binding capacity for proteins. The advantages of nitrocellulose are that it binds protein well and affords very good signal-to-background ratios in Western blotting assays. It allows staining of the immobilized protein patterns with such dyes as [Ponceau S](#) and Amido black. It becomes brittle, however, after baking at 80°C (typical for Southern and Northern blotting), which reduces the repeated use of these filters. Furthermore, nitrocellulose is less resistant to various organic solvents, such as methanol (5).

#### 2.2. Nylon Membranes

Nylon membranes were introduced initially as an alternative to nitrocellulose for protein blotting (5), and subsequently they were applied in RNA and DNA transfers (4). These membranes are usually derivatives of nylon 66 and often are modified with positive charges. They have proven exceptional

binders of DNA, RNA, and protein, but at times this presents some difficulty by producing high backgrounds, particularly in protein blotting. These filters are mechanically stable and thus can be reprobed numerous times without losing band definition or sensitivity.

### 2.3. Polyvinyl Difluoride

Polyvinyl difluoride (PVDF) membrane filters have a special application in blotting where subsequent chemical manipulation of the immobilized macromolecule is desired. A case in point is the use of these membranes in protein blots, where the individual bands are excised out of the blot and subjected to [Edman Degradation](#) for N-terminal amino acid sequencing of the resolved and blotted polypeptide chains ([6](#), [7](#)).

### Bibliography

1. L. Clarke, R. Hitzman, and J. Carbon (1979) *Methods Enzymol.* **68**, 436–442.
2. J. C. Alwine, D. J. Kemp, and G. R. Stark (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5350–5354.
3. A. De Maio (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 11–32.
4. J. Meinkoth and G. Wahl (1984) *Anal. Biochem.* **138**, 267–284.
5. J. M. Gershoni and G. E. Palade (1982) *Anal. Biochem.* **124**, 396–405.
6. M. A. Mansfield (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 33–52.
7. C. Eckerskorn and F. Lottspeich (1993) *Electrophoresis* **14**, 831–838.

### Blunt-End Ligation

Blunt-end ligation is covalently joining two or more double-stranded DNA fragments with flush ends by the enzyme [DNA Ligase](#) ([1](#)). Many [restriction enzymes](#) produce blunt ends, when they cleave both polynucleotide chains at the same site, and some **cohesive ends** are converted to blunt ends by the [fill-in reaction](#). Blunt end ligation is a much less efficient reaction than ligation of DNA fragments with cohesive ends. One technique for overcoming the low efficiency of blunt-end ligation is ligating a short oligonucleotide **linker fragment** to the DNA ends that are to be joined ([2](#)). Because a large molar concentration of the linker can be added to the reaction, its ligation to the ends of the DNA occurs more readily than joining of the ends themselves. Linkers are designed to contain a unique restriction site so that the ligation mixture is digested with the corresponding enzyme to create cohesive ends, which can be joined by conventional ligation.

The efficiency of blunt-end ligation is greatly increased by carrying out the reaction in the presence of 15% polyethylene glycol, 6,000 MW ([PEG](#)), although only linear molecules are produced in this reaction. Circular **plasmids** are required for efficient [transformation](#) of *Escherichia coli*. An efficient way of producing circular molecules from the linear product of blunt-end ligation is including a bacteriophage P1 [lox recombination](#) site in one of the molecules being ligated ([3](#)). After ligation, the DNA is treated with Cre recombinase, which by recombination between the lox sites produces a circular molecule from a linear molecule that contains two lox sites. This type of molecule is created when the molecule without a lox site is ligated to two molecules with a lox site, one at each end.

It is possible to amplify the products of a low-efficiency ligation containing the desired insert by **PCR**, using primers that bind to vector DNA on each side of the cloning site. The PCR product of the expected size is purified and recloned. This results in transformants from ligations that give no

positive colonies on direct transformation (4).

### Bibliography

1. T. Maniatis, E. F. Fritsch, and J. Sambrook (1982) *Molecular Cloning; a Laboratory Manual*, 2nd ed., Cold Spring Harbor Press, Cold Spring Harbor, NY.
2. G. J. Bhat, M. J. Lodeg, P. J. Myler, and K. D. Stuart (1991) *Nucleic Acids Res.* **19**, 398.
3. A. C. Boyd (1993) *Nucleic Acids Res.* **21**, 817–821.
4. H. W. Son and E. Lolis (1995) *BioTechniques* **18**, 644–650.

### Suggestion for Further Reading

5. A. E. T. Konson and D. S. Levin (1997) Mammalian DNA ligases, *Bioessays* **19**, 893–901; a review of DNA ligases.

## Bohr Effect

Bohr et al. (1) discovered that the oxygen dissociation curve of blood shifts to higher partial pressures of oxygen in the presence of increasing concentrations of carbon dioxide by lowering the oxygen affinity of [hemoglobin](#). This effect, it was concluded, results from the concomitant decrease in pH, and the effect of pH on the oxygen affinity of hemoglobin has subsequently been known as the Bohr effect. It was subsequently discovered that CO<sub>2</sub> also lowers the oxygen affinity of hemoglobin (Hb) directly by binding to its α-amino groups, forming carbamino compounds (2, 3):



The original observations of Bohr et al. (1) are now known as the “classical Bohr effect,” which is a composite of the specific effects of H<sup>+</sup> and CO<sub>2</sub>.

The CO<sub>2</sub> content of blood is reduced upon oxygenation of hemoglobin at a constant CO<sub>2</sub> pressure (4). This phenomenon, known as the “classical Haldane effect,” is a combination of the oxygen-linked dissociation of the carbamino groups (Eq. 1) and release of protons



followed by dehydration of bicarbonate:



The release of protons upon oxygenation or their uptake upon deoxygenation of hemoglobin is known as the “Haldane effect.”

The Bohr effect and the Haldane effect are thermodynamically equivalent (5) because they are linked functions (6). Whatever effect oxygen binding has on the affinity of hemoglobin for protons, changes in the pH must have the same effect on the affinity of hemoglobin for oxygen:

$$\left(\frac{\delta H^+}{\delta Y}\right)_{pH} = \left(\frac{\delta \log P_{O_2}}{\delta pH}\right)_Y \quad (4)$$

where  $H^+$  is the number of protons bound per heme group,  $Y$  is the fraction of hemoglobin binding sites occupied by  $O_2$ , and  $P_{O_2}$  is the partial pressure of  $O_2$ . The left-hand side of Eq. 4 is known as the “Haldane coefficient,” and the right-hand side is the “Bohr coefficient.” They give the magnitude of the Haldane and Bohr effects, which must be the same.

These effects are physiologically important for efficient and regulated transport of  $O_2$  and  $CO_2$  in opposite directions (see [Hemoglobin](#)). When the shape of the oxygen dissociation curve expressed by  $Y$  versus  $\log P_{O_2}$  is independent of changes in pH, at least within the range  $0.1 < Y < 0.9$ , Eq. 4 is simplified to

$$\delta = \frac{\Delta \log P_{50}}{\Delta pH} \quad (5)$$

where  $P_{50}$  is the  $P_{O_2}$  at  $Y = 0.5$  and  $d$  is the Bohr coefficient, which is independent of the degree of  $O_2$  binding. The Bohr coefficient expresses the number of protons (per heme group) bound to hemoglobin upon full oxygenation and is the same as  $d$  of Eq. 2.

$d$  is negative above pH 6.3, and protons are released upon oxygenation of hemoglobin, whereas  $d$  is positive at lower pH values, and protons are taken up upon oxygenation. These two phenomena are known respectively as the “alkaline” and “acid” Bohr effects. At physiological pH 7.4, human adult hemoglobin A (HbA) has  $d = -0.6$  in the presence of 0.1 M NaCl. This value is halved in the absence of  $Cl^-$ .

The Bohr effect arises from changes in the  $pK_a$  values of particular ionizable groups of hemoglobin as a result of changes in their environment upon binding oxygen. The group primarily responsible for the alkaline Bohr effect of human HbA is the imidazole side chain of the C-terminal His146 of the  $\beta$  chain. Its  $pK_a$  value changes from 8.0 to 7.1 upon oxygenation, contributing the greater part of the alkaline Bohr effect in the absence of  $Cl^-$ . In the presence of  $Cl^-$ , binding of these ions to several ionizable groups that line the central cavity of the hemoglobin molecule, including the  $\alpha$ -amino groups of the  $\alpha$  chains and the  $\epsilon$ -amino groups of Lys82 of the  $\beta$  chains, affects the protonation of those groups. Consequently, the release of  $Cl^-$  ions upon binding of oxygen produces an additional,  $Cl^-$ -dependent Bohr effect (7). It is thought that His143 of the  $\beta$  chains is one of the groups responsible for the acid Bohr effect (8), but the others are still unidentified.

The term “Bohr effect” is sometimes used for pH-dependent binding of other ligands, such as  $CO_2$  and alkylisocyanides, to hemoglobin. It is also used for pH-dependent ligand binding to [oxygen-binding proteins](#) other than hemoglobin.

## Bibliography

1. C. Bohr, K. A. Hasselbalch, and A. Krogh (1904) *Skand. Arch. Physiol.* **16**, 402–412.
2. J. K. W. Ferguson and F. J. W. Roughton (1934) *J. Physiol. (London)* **83**, 87–102.
3. J. V. Kilmartin and L. Rossi-Bernardi (1969) *Nature* **222**, 1243–1246.
4. J. Christiansen, C. C. Douglas, and J. S. Haldane (1914) *J. Physiol. (London)* **48**, 244–277.
5. I. Tyuma and Y. Ueda (1975) *Biochem. Biophys. Res. Commun.* **65**, 1278–1283.
6. J. Wyman (1964) *Adv. Protein Chem.* **19**, 223–286.
7. M. F. Perutz, G. Fermi, C. Poyart, J. Pagnier, and J. Kister (1993) *J. Mol. Biol.* **233**, 536–545.

8. M. F. Perutz, J. V. Kilmartin, K. Nishikura, J. H. Fogg, P. J. G. Butler, and H. S. Rollema (1980) *J. Mol. Biol.* **138**, 649–670.

### **Suggestions for Further Reading**

9. J. V. Kilmartin and L. Rossi-Bernardi (1973) Interaction of hemoglobin with hydrogen ions, carbon dioxide, and organic phosphates, *Physiol. Rev.* **53**, 836–890.
10. I. Tyuma (1984) The Bohr effect and the Haldane effect in human hemoglobin, *Jpn. J. Physiol.* **34**, 205–216.

### **Bowman–Birk Inhibitors**

Bowman–Birk proteinase inhibitor is a member of a widely studied family of serine proteinase inhibitors, protein. Most such inhibitors consists of a single polypeptide chain of 65–80 amino acid residues containing two homologous regions, each with a reactive site. In some Bowman–Birk inhibitors, both sites are specific for trypsin; in others, one is specific for trypsin and one for chymotrypsin; and, in still others, one is for trypsin and one for elastase. Ternary complexes of one inhibitor and two enzyme molecules are commonly found. The polypeptide chain is crosslinked by seven disulfide bridges, some within the homology regions, others between them, endowing the inhibitor with great stability.

#### **Suggestion for Further Reading**

T. Ikenaka and S. Norioka (1986) "Bowman–Birk family serine proteinase inhibitors". In *Proteinase Inhibitors* (A. Barrett and G. Salveson, eds.), Elsevier New York, 361–374.

### **BPTI (Bovine Pancreatic Trypsin Inhibitor)**

Bovine pancreatic trypsin inhibitor, BPTI (Kunitz), is also called trypsin kallikrein inactivator, TKI, aprotinin, and Trasylol™. It is the first protein inhibitor of proteinases [see [Proteinase Inhibitors, Protein](#)] to be isolated and characterized by Moses Kunitz as an inhibitor of bovine trypsin. It was independently discovered as a kallikrein inactivator. BPTI was the first protein proteinase inhibitor to be sequenced. Shortly afterward, sequencing the kallikrein inactivator revealed the identity of the two inhibitors. They consist of 58 amino acid residues, crosslinked by three disulfides. Lys<sup>15</sup> is at the P<sub>1</sub> position of the reactive site. The K<sub>i</sub> for BPTI-bovine b trypsin interaction, 5 × 10<sup>-14</sup> M, is often—probably incorrectly—referred to as the strongest of the proteinase–protein proteinase inhibitor interactions. BPTI is also extremely stable. Its T<sub>m</sub> at neutral pH is 103°C. BPTI was the first protein proteinase inhibitor to have its three-dimensional structure determined both in free form and in complex with bovine b trypsin. It has been—and still is—widely used in Europe as a drug to avoid postsurgical complications. For this use, it is isolated from bovine lungs. BPTI has been the first superbly characterized small protein that was widely available. Therefore, it became the favorite substance of many protein researchers. It was the first protein to have its three-dimensional structure



determined by nuclear magnetic resonance (NMR). It served as the object of numerous molecular dynamics and structure prediction calculations. It was the first protein for which the pattern of closure of disulfide bridges was studied in detail. BPTI (Kunitz) gave its name to a widely studied family of standard-mechanism canonical protein inhibitors of serine proteinases.

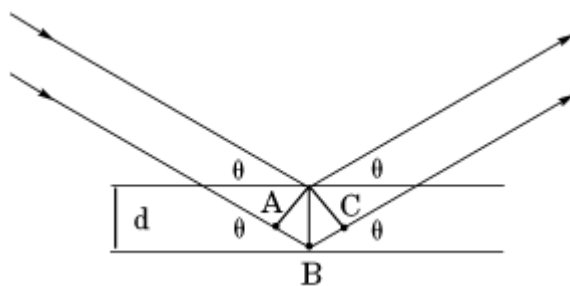
#### Suggestion for Further Reading

W. Gebhard, H. Tschesche, and H. Fritz (1986) "Bowman-Birk family serine proteinase inhibitors". In *Proteinase Inhibitors* (A. Barrett and G. Salveson, eds.), Elsevier New York, 375–388.

### Bragg Angle

The Bragg angle is important for analyzing the diffraction from crystals in [X-ray crystallography](#). W.L. Bragg noted that X-ray diffraction by a crystal can be understood without going into the details of diffraction theory (1). The process is similar to ordinary reflection of light from a mirror. In the crystal, the mirrors are the lattice planes. They are the planes constructed through the lattice points, which are the corners of the [unit cells](#). Within a set, the planes are parallel and equidistant with perpendicular distance  $d$  (Fig. 1).  $\theta$ , the reflective angle, is called the Bragg angle. The beams reflected from the upper and the lower lattice plane are in phase and reinforce each other if the difference in path lengths of the two beams, which is given by the path A–B–C, is an integral number of the wavelength  $\lambda$ . The path difference is also equal to  $2d \sin \theta$ . This results in Bragg's law:  $2d \sin \theta = n \lambda$ , where  $n$  is an integer, 1, 2, 3, etc.

**Figure 1.** Bragg's view on X-ray diffraction. The incident beam is reflected by lattice planes.  $d$  is their distance, and  $\theta$  is the reflective angle.



#### Bibliography

1. W. L. Bragg (1913) *Proc. Cambridge Phil.Soc.* **17**, 43–57.

#### Suggestions for Further Reading

2. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists*, VCH Publishers, New York, Weinheim, Cambridge.
3. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

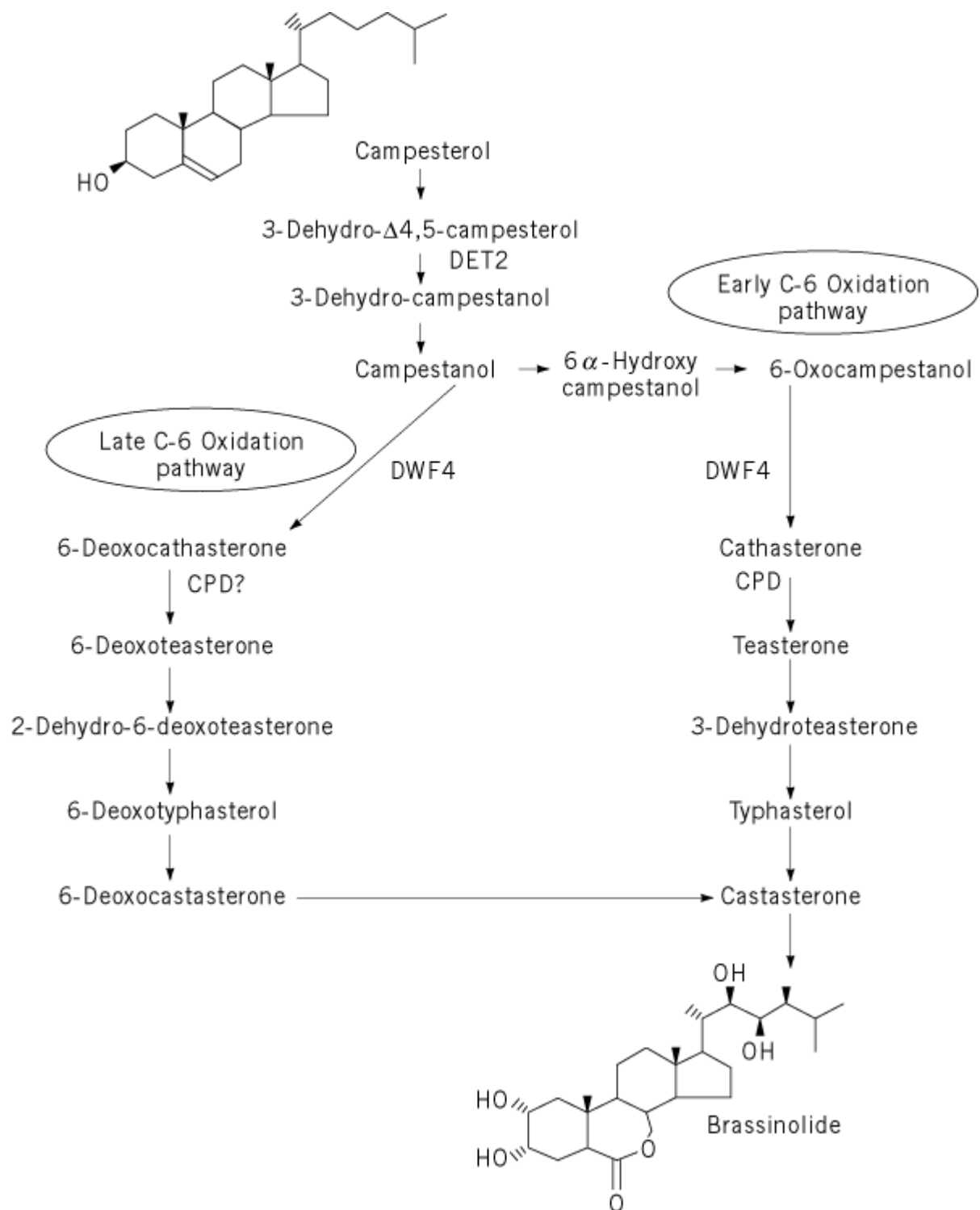
## Brassinosteroids

### 1. History

Unlike the five other “classical” [plant hormones](#), which were all discovered more than 35 years ago, brassinosteroids (BRs) have been identified only recently. In 1979, Grove et al. [\(1\)](#) reported on the first growth-promoting steroid from plants. The substance was isolated from *Brassica napus* (rapeseed) pollen and was called *brassinolide*. Brassinolide was found in minute quantities in pollen grains; about 250 kg of pollen were needed to obtain only 10 mg of the pure compound. BRs are widely present in angiosperm families, but are also found in gymnosperms, ferns, and green algae [\(2\)](#).

Recently, BRs have gained recognition as true plant hormones [\(3\)](#). Genetic experiments provided evidence that they are indispensable for normal plant growth. A second breakthrough in our knowledge of BRs owes to the recent elucidation of the entire biosynthetic pathway [\(2, 4\)](#) (Fig. [1](#)).

**Figure 1.** Brassinosteroid biosynthesis in higher plants. Steps mediated by DET2, DWF4, and CPD in *Arabidopsis* are indicated. (Adapted from Ref. [7](#).)



## 2. Biosynthesis and Metabolism

Brassinosteroids are unique as plant hormones because they are structurally related to the insect molting hormone [ecdysone](#) and to mammalian [steroid hormones](#). More than 40 BRs occur naturally, all of which share a common 5 $\alpha$ -cholestane skeleton. The relative activities of the BRs increase according to their position in the biosynthetic pathway: brassinolide is most active, whereas early precursors up to 6-oxocampestanol are virtually inactive (5). Hydroxylation at position C22 is critical for biological activity.

Brassinolide is synthesized from campesterol via two pathways (2) (Fig. 1). These pathways were established by classical labeled-precursor feeding experiments. Campesterol is derived essentially from the mevalonate pathway. Initially, the double bond in the B ring of campesterol is reduced to campestanol (6) over 3-hydro-D4,5-campesterol and 3-dehydro-campestanol. From campestanol, two possible routes lead to brassinolide. In this branch, campestanol can be oxidized on C6 to yield 6-oxocampestanol, with 6a-hydroxycampestanol as an intermediate.

The first branch is called the early C6-oxidation pathway and seems ubiquitous in higher plants. 6-Oxocampestanol is subsequently hydroxylated at position 22 on the side chain (7), resulting in cathasterone, and further hydroxylated at position 23, to yield teasterone (4). This step is followed by epimerization of the hydroxyl function on C3, yielding typhasterol, with 3-dehydroteasterone as an intermediate (8, 9). Typhasterol is further hydroxylated on C2, resulting in castasterone that is finally converted to brassinolide by lactamization of the B ring (10). Two steps in the early C6-oxidation pathway appear to be rate-limiting: (a) the formation of cathasterone from 6-oxocampestanol and (b) the final conversion of castasterone to brassinolide (2).

A second possible route from campestanol to castasterone is the so-called late C6-oxidation path, with 6-deoxo brassinosteroids as intermediates. The latter compounds are the least active BRs and were therefore considered to be dead-end products. The second biosynthetic path has been discovered only recently; it was demonstrated in cultured cells and seedlings of *Catharanthus roseus*, as well as in seedlings of rice and tobacco (11, 12). Recent evidence indicates that 6-deoxoteasterone is formed from campestanol via 6-deoxocathasterone, by subsequent hydroxylations (7). This step is followed by an epimerization to yield 6-deoxytyphasterol, comparable to the conversion from teasterone to typhasterol in the early oxidation pathway. Finally, 6-deoxytyphasterol is hydroxylated at position 2, forming castasterone.

Compelling evidence for an essential role for BRs in plant growth and development came from recent genetic studies in *Arabidopsis thaliana* (3, 13). The *det2* (*de-etiolated 2*), *dwf4* (*dwarf4*), *cpd* (*constitutive photomorphogenesis and dwarfism*), and *dim* (*diminuto*) mutants identify different steps in BR biosynthesis (7, 14-17). All of these mutants display drastic **phenotypic** effects that can be reverted by application of brassinolide. In the light, the plants are smaller than wild type because of reduced cell size, the leaves are darker green, apical dominance is reduced, and the plants are less fertile. In the dark, a de-etiolation phenotype is observed, characterized by short and thick hypocotyls, open and expanded cotyledons, and formation of primary leaves. This phenotype is accompanied by derepression of light-regulated genes, with the exception of the *dim* mutant, where a normal repression of light-regulated genes is observed in the dark (18). Collectively, these observations support the existence of cross-talk between the light and BR signaling pathways. The four **genes** corresponding to the above-mentioned loci have been **cloned** recently. The *DET2* gene product shares similarity with mammalian steroid 5 $\alpha$ -reductases. Definite proof for its function was given by two experiments that support functional conservation between mammalian and plant steroid 5 $\alpha$ -reductases (19). When expressed in human embryonic kidney cells, DET2 catalyzed 5 $\alpha$ -reduction of several animal steroid substrates. In addition, *det2* mutants could be rescued by expression of human steroid 5 $\alpha$ -reductases. The *CPD* and *DWF4* genes were cloned by T-DNA tagging, and both encode **cytochrome P450** monooxygenases that share **homology** with steroid hydroxylases (7, 15). Feeding studies indicated that DWF4 acts as a 22 $\alpha$ -hydroxylase (7), whereas CPD functions as a 23 $\alpha$ -hydroxylase (15). It is reasonable to assume that each enzyme would play a role in both branches of BR biosynthesis, although there is yet no evidence that CPD uses 6-deoxocathasterone as a substrate. Finally, the *DIM* gene encodes a protein with a putative FAD-binding domain and a nuclear localization motif (see **Nuclear Import, Export**) (18, 20). Feeding experiments indicate that *dim* mutations affect a step before typhasterol formation, so it might be involved in epimerization of teasterone (13).

As is the case for most other plant hormones, glucosylation plays a role in deactivation of BRs. Some of these conjugates may serve as a storage form, as for instance the C23-glucosylated brassinolide and esters at position C3 (21, 22). Permanent inactivation of BRs is achieved by C25

and C26 hydroxylation, as well as by cleavage of the side chain (13). It should be noted that transport of BR storage forms is facilitated after glucosylation, as a result of their increased [hydrophilicity](#). Transport of labeled BRs has been shown to occur from root to shoot, probably via xylem (23). It is not yet clear, however, whether this also plays a role in determining the endogenous levels of BRs.

### 3. Signal Perception and Transduction

The conservation of steroid-like compounds as signaling molecules across several phyla including fungi, invertebrates, and vertebrates made it tempting to speculate on the existence of soluble nuclear BR receptors in plants (24, 25). Two BR-insensitive mutants, *bri1* (*brassinosteroid insensitive 1*) and *cbb2* (*cabbage 2*), have been isolated in *Arabidopsis* (16, 26). In both mutants, hormone insensitivity is specific to BRs, and neither one can be rescued by external BR application. Expression of BR-regulated genes was abolished in BR-insensitive mutants (16). *Bri1* and *cbb2* mutations result in dramatic effects on development, including an exacerbated dwarfism, bushy phenotype, dark green and thickened leaves, and male sterility. In addition, de-etiolation of dark-grown seedlings was observed (26). *Bri1* and *cbb2* are alleles of a single locus, and the corresponding gene has been cloned using a map-based approach (27). Surprisingly, BRI1 does not share similarity with **steroid receptors**, but belongs to the class of receptor-like transmembrane [kinases](#) (28). Members of this family have an *N*-terminal extracellular [leucine-rich repeat](#) (LRR); there are 25 such repeats in BRI1, with a unique stretch of 70 amino acid residues between the 21st and 22nd LRR. These are followed by a transmembrane domain and an internal **serine/threonine kinase** domain that relays the signal. Other LRR kinases include proteins that influence meristem size and organ formation, such as *clavata1*, *erecta*, and *Xa21*, which confers resistance against *Xanthomonas oryzae* pv. *Oryzae* in rice (29-31). The *BRI1* gene is ubiquitously expressed in plant tissues in light and dark, consistent with a presumed role as a BR receptor in different cell types (27). The BRI1 protein could bind BR either directly or indirectly, via a protein intermediary. Although none of the previously identified LRR kinases can bind nonproteinaceous ligands, the island of 70 unique amino acid residues in the LRR of BRI1 could serve the purpose of directly binding BR. Mutation of this region severely affects BRI1 activity. In animal systems, LRR kinases are known to form homo- or heterodimers upon ligand binding, thereby activating the intrinsic kinase activity (28). Heterodimerization of BRI1 and XA21, or other LRR domains of disease resistance proteins, could control cross-talk between BR and pathogenic signaling (27).

### 4. Downstream Targets

To date, genes involved in the primary response to BR have not been identified. Likewise, [cis-acting](#) sequences or [trans-acting](#) factors involved in the BR response remain to be characterized. In contrast, the expression of at least 50 genes was altered upon treatment of soybean hypocotyls and epicotyls with BR for 2 to 17 h (32, 33). One of these genes was cloned, and the predicted polypeptide chain shared similarity with xyloglucan endo-*trans*-glycosylases (XETs) (34). The gene was termed *BRUI* (*BR up-regulated*) and is not responsive to either [auxin](#) or [gibberellin](#). Its recombinant product catalyzed endo-*trans*-glycosylation *in vitro* (35). The *in vivo* function of *BRUI* is thus to control expansion growth, either by cell-wall loosening or by integration of newly synthesized xyloglucans in the wall (36).

In *Arabidopsis*, the *TCH4* (*touch 4*) gene encodes an active XET, and its transcriptional activation is maximal within 2 h of BR treatment (37). XET is encoded by a [multigene family](#) in the *Arabidopsis* [genome](#). Some family members, such as *TCH4*, are responsive to BR and auxin, but not to gibberellins, whereas others, such as *meri5* (*meristem 5*), are regulated by gibberellins, less generally by BR, and not by auxin (16). Certain *XET* genes do not respond to BR. The *XET* gene family thus appears to be differentially regulated by a combination of hormonal and environmental cues (such as light, heat, cold, and touch) (38). [In Situ Hybridization](#) of soybean epicotyl cross sections with a *BRUI* probe revealed highest gene expression in vascular tissues (phloem cells and xylem parenchyma), as well as in the cortical starch sheath layer (35). In contrast to the rather restricted

action of auxin on cell elongation, BRs have been shown to have an effect both on epidermal and internal tissues (39). Nevertheless, the action is more profound on internal tissues, correlating well with the spatial expression pattern of *BRUI*. Comparative kinetic analysis of auxin and BR action indicates that auxin possibly initiates cell wall elongation at the epidermis, whereas BR continues to stimulate elongation of both outer and inner layers (39-41). BRs thus seem to have a prolonged effect on cell elongation. Brassinolide application results in a predominantly transverse orientation of cortical microtubules, allowing expansion in the longitudinal direction (42).

## 5. Effects

The final effects of BRs result from expression of an array of secondary response genes, mentioned above. A multitude of BR functions in plant development are known, including leaf bending and unrolling, stem elongation, root inhibition, xylogenesis, and pollen tube growth (5). It is clear that several of these activities may result from cross-talk with light and other hormonal signaling cascades.

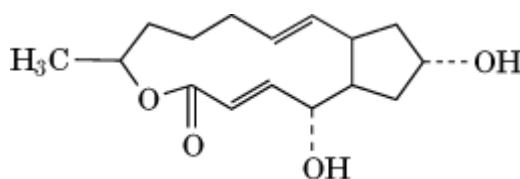
## Bibliography

1. M. D. Grove, G. F. Spencer, W. K. Rohwedder, N. Mandava, J. F. Worley, J. D. Warthen, G. I. Steffens, J. L. Flippen-Anderson, and J. C. Cook (1979) *Nature* **281**, 216–217.
2. S. Fujioka and A. Sakurai (1997) *Physiol. Plant.* **100**, 710–715.
3. S. D. Clouse (1996) *Curr. Biol.* **6**, 658–661.
4. S. Fujioka, T. Inoue, S. Takatsuto, T. Yanagisawa, T. Yokota, and A. Sakurai (1995) *Biosci. Biotech. Biochem.* **59**, 1543–1547.
5. T. Yokota (1997) *Trends Plant Sci.* **2**, 137–143.
6. H. Suzuki, T. Inoue, S. Fujioka, T. Saito, S. Takatsuto, T. Yokoto, N. Murofushi, T. Yanagisawa, and A. Sakurai (1995) *Phytochemistry* **40**, 1391–1397.
7. S. Choe, B. P. Dilkes, S. Fujioka, S. Takatsuto, A. Sakurai, and K. A. Feldmann (1998) *Plant Cell* **10**, 231–243.
8. H. Suzuki, S. Fujioka, S. Takatsuto, T. Yokota, N. Murofushi, and A. Sakurai (1994) *J. Plant Growth Regul.* **13**, 21–26.
9. H. Suzuki, T. Inoue, S. Fujioka, S. Takatsuto, T. Yanagisawa, T. Yokota, N. Murofushi, and A. Sakurai (1994) *Biosci. Biotech. Biochem.* **58**, 1186–1188.
10. T. Yokota, Y. Ogino, N. Takahashi, H. Saimoto, S. Fujioka, and A. Sakurai (1990) *Agric. Biol. Chem.* **54**, 1107–1108.
11. Y.-H. Choi, S. Fujioka, A. Harada, T. Yokota, S. Takatsuto, and A. Sakurai (1996) *Phytochemistry* **43**, 593–596.
12. Y.-H. Choi, S. Fujioka, T. Nomura, A. Harada, T. Yokota, S. Takatsuto, and A. Sakurai (1997) *Phytochemistry* **44**, 609–613.
13. M. Szekeres and C. Koncz (1998) *Plant Physiol. Biochem.* **36**, 145–155.
14. J. Li, P. Nagpal, V. Vitart, T. C. McMorris, and J. Chory (1996) *Science* **272**, 398–401.
15. M. Szekeres, K. Németh, Z. Koncz-Kálmán, J. Mathur, A. Kauschmann, T. Altmann, G. P. Rédei, F. Nagy, J. Schell, and C. Koncz (1996) *Cell* **85**, 171–182.
16. A. Kauschmann, A. Jessop, C. Koncz, M. Szekeres, L. Willmitzer, and T. Altmann (1996) *Plant J.* **9**, 701–713.
17. R. Azpiroz, Y. Wu, J. C. LoCascio, and K. A. Feldmann (1998) *Plant Cell* **10**, 219–230.
18. T. Takahashi, A. Gasch, N. Nishizawa, and N.-H. Chua (1995) *Genes Dev.* **9**, 97–107.
19. J. Li, M. G. Biswas, A. Chao, D. W. Russell, and J. Chory (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3554–3559.
20. A. R. Mushegian and E. V. Koonin (1995) *Protein Sci.* **4**, 1243–1244.
21. H. Suzuki, S. K. Kim, N. Takahashi, and T. Yokota (1993) *Phytochemistry* **33**, 1361–1367.

22. S. Asakawa, H. Abe, N. Nishikawa, M. Natsume, and M. Koshioka (1996) *Biosci. Biotech. Biochem.* **60**, 1416–1420.
23. N. Nishikawa, S. Toyama, A. Shida, and F. Futatsuya (1994) *J. Plant Res.* **107**, 125–130.
24. M. Beato, P. Herrlich, and G. Schutz (1995) *Cell* **83**, 851–857.
25. S. D. Clouse (1996) *Plant J.* **10**, 1–8.
26. S. D. Clouse, M. Langford, and T. C. McMorris (1996) *Plant Physiol.* **111**, 671–678.
27. J. Li and J. Chory (1997) *Cell* **90**, 929–938.
28. D. M. Braun and J. C. Walker (1996) *Trends Biochem. Sci.* **21**, 70–73.
29. S. E. Clark, R. W. Williams, and E. M. Meyerowitz (1997) *Cell* **89**, 575–585.
30. W.-Y. Song, G.-L. Wang, L.-L. Chen, H.-S. Kim, L.-Y. Pi, T. Holsten, J. Gardner, B. Wang, W.-X. Zhao, L.-H. Zhu, C. Fauquet, and P. Ronald (1995) *Science* **270**, 1804–1806.
31. K. U. Torii, N. Mitsukawa, T. Oosumi, Y. Matsuura, R. Yokoyama, R. F. Whitier, and Y. Komeda (1996) *Plant Cell* **8**, 735–746.
32. S. D. Clouse and D. Zurek (1991) In *Brassinosteroids: Chemistry, Bioactivity and Applications*, ACS Symposium Series, Vol. **474** (H. G. Cutler, T. Yokota, and G. Adam, eds.), American Chemical Society, Washington, D.C., pp. 122–140.
33. S. D. Clouse, D. M. Zurek, T. C. McMorris, and M. E. Baker (1992) *Plant Physiol.* **100**, 1377–1383.
34. D. M. Zurek and S. D. Clouse (1994) *Plant Physiol.* **104**, 161–170.
35. S. D. Clouse (1997) *Physiol. Plant.* **100**, 702–709.
36. D. J. Cosgrove (1997) *Plant Cell* **9**, 1031–1041.
37. W. Xu, M. M. Purugganan, D. H. Plisensky, D. M. Antosiewicz, S. C. Fry, and J. Braam (1995) *Plant Cell* **7**, 1555–1567.
38. W. Xu, P. Campbell, A. K. Vargheese, and J. Braam (1996) *Plant J.* **9**, 879–889.
39. R. Tominaga, N. Sakurai, and S. Kuraishi (1994) *Plant Cell Physiol.* **35**, 1103–1106.
40. M. Katsumi (1985) *Plant Cell Physiol.* **26**, 615–625.
41. D. M. Zurek, D. L. Rayle, T. C. McMorris, and S. D. Clouse (1994) *Plant Physiol.* **104**, 505–513.
42. K. Mayumi and H. Shibaoka (1995) *Plant Cell Physiol.* **36**, 173–181.

## Brefeldin A

Brefeldin A (BFA) is a fungal metabolite that is a powerful tool for dissecting membrane traffic and organelle dynamics. It is named after *Penicillium brefeldianum*, which produces it. BFA is a macrolide antibiotic with a molecular weight of 280.36 Da and the structure



BFA is normally used on cells in culture at concentrations of 1 to 10  $\mu\text{M}$ . As first discovered by Misumi et al. (1), BFA primarily inhibits secretion, without significantly affecting [endocytosis](#). It

does not affect protein biosynthesis significantly, does not deplete cellular ATP levels, and does not affect **cytoskeletal** elements. The unique feature of BFA, which distinguishes it from other inhibitors of protein traffic, is that it causes coalescence of organelles. Within minutes of BFA application, the [Golgi apparatus](#) fragments and fuses with the [endoplasmic reticulum](#) (ER), which does not fragment (2), leading to mixing of the two compartments.

The effects of BFA are completely reversible. Retrograde movement of Golgi proteins into the ER occurs via long, tubulovesicular processes extending out of the Golgi along [microtubules](#). This retrograde traffic pathway can be separated into two distinct phases: movement of Golgi proteins into the intermediate compartment between the ER and Golgi complex, followed by cycling between the modified ER and the intermediate compartment (3). Remarkably, many Golgi [enzymes](#) are still active in this environment and can modify proteins in the ER (4).

One of the important uses of BFA-induced redistribution has been the functional demarcation of compartments along the secretory pathway (5). It is clear that BFA affects the early Golgi compartments, the *cis* and medial cisternae. Whether BFA also affects the *trans*-most cisternae of the Golgi stack and the *trans*-Golgi network seems to be dependent on tissue and cell type. Thus, in pancreatic acinar cells, exit from the Golgi complex and formation of **secretory granules** are not BFA-sensitive (6), while in other cells the formation of secretory granules is inhibited (7). Secretion of preformed secretory granules is not inhibited by BFA. Thus, the effects of BFA point to a common mechanism in the traffic of proteins that follow either the constitutive or the regulated pathways of [exocytosis](#).

BFA's effects are not limited to the Golgi apparatus and are reiterated throughout the **endosome/lysosome** system. BFA treatment induces tubulation of these compartments (8). Similar to the mixing of the Golgi with the ER, the *trans*-Golgi network mixes with the recycling endosomal system. Remarkably, this mixed system remains partly functional, with normal cycling between plasma membrane and endosomes, but with impaired traffic between endosomes and lysosomes (8). This suggests that the vesicle-budding mechanisms (see below) in the endocytic pathway are similar to, yet distinct from, those in the exocytic pathway. Furthermore, these observations reinforce the plasticity of the endocytic pathway, which is largely functional even when some traffic steps are inhibited.

The major effects of BFA on protein traffic are explained by its ability to prevent binding of cytosolic coat proteins onto membranes (9). The target of BFA's action is a **nucleotide exchange factor** for **ADP-ribosylation factor** (ARF), a small [GTP-binding protein](#). BFA inhibits the ability of Golgi membranes to catalyze the exchange of GTP onto ARF specifically, thereby preventing ARF from interacting with the membrane. As a result, the association of the coat protein b-COP with the Golgi membrane is inhibited (10, 11). b-COP is a subunit of a cytosolic protein complex, the coatamer, that reversibly associates with Golgi membranes and controls vesicular transport.

One possible mode of ARF action is via its ability to stimulate phospholipase D in Golgi membranes (12). Phospholipase D converts phospholipids into phosphatidic acid and thus can alter the [lipid](#) content of [membranes](#). Stimulation of the Golgi-associated phospholipase D activity is BFA-sensitive, suggesting a possible link between transport events and the underlying architecture of the lipid bilayer.

The pathogenic bacterium *Staphylococcus aureus* mimics the action of BFA and disassembles the Golgi apparatus (13, 14). This effect is mediated by the secretion of EDIN (epidermal-cell differentiation inhibitor), an extracellular enzyme that ADP-ribosylates [Rho GTPase](#). As a result, all the manifestations of BFA treatment are reproduced. Thus, the regulatory circuit of coat-membrane assembly may involve more than one small GTP-binding protein. In addition to the effects of the small GTP-binding proteins, [heterotrimeric G proteins](#) of the Gi/Go subfamily also contribute to the regulation of the cycle of coatamer binding (15, 16). Activation of heterotrimeric G proteins promotes binding of b-COP to Golgi membranes and antagonizes the effect of BFA.



Several BFA-resistant [cell lines](#) exist, including the PtK1 rat kangaroo cell line, a derivative of monkey kidney Vero cells, and two derivatives of the human epidermoid carcinoma KB cell line (17). The BFA resistance is dominant and is due to a Golgi-associated factor that is **homologous** to the target of BFA in cells that are sensitive to the drug. In Golgi membranes from BFA-resistant PtK1 cells, the basal phospholipase D activity is high and insensitive to BFA (12), suggesting a mechanism for the drug resistance. Another mechanism may involve a member of the ATP-binding cassette superfamily of transport proteins, as suggested by BFA-resistant mutants in yeast (18).

The ability of BFA to affect budding of **clathrin**-coated vesicles from the *trans*-Golgi network (observed in many, but not all cells) is mediated through the rapid and reversible redistribution of g-adaptin (19). This component of the clathrin coat is specific to Golgi-derived coated vesicles and is absent from plasma membrane-derived coated vesicles. The kinetic and pharmacological similarities between BFA's effects on g-adaptin and b-COP underscore the biochemical similarities between membrane budding mechanisms that are mediated by various coat complexes.

The effects of BFA are readily reversible, and this reversibility is due partly to its detoxification. The mechanism of BFA detoxification was studied in CHO cells and is mediated by the **glutathione S-transferase** system via conversion of the antibiotic to its **glutathionyl** and **cysteinyl** derivatives, followed by secretion (20).

BFA affects cellular compartmentalization not only via protein traffic, but also through lipid metabolism (21). Its main effect on lipids is enhanced hydrolysis of sphingomyelin, a key regulator of cell proliferation and differentiation. Sphingomyelin is produced from ER-derived ceramide and is delivered to other membranes, presumably via the same trafficking system used for proteins. The effect of BFA on the level of sphingomyelin seems to be a result of the mixing of cellular membranes. BFA treatment also increases the rate of sphingomyelin biosynthesis from phosphatidylcholine, and this effect may be related to the above-mentioned action on phospholipase D. In this context, it is noteworthy that C6 ceramide, a cell-permeable ceramide analogue, partially restores BFA sensitivity in a BFA-resistant cells (22).

The various effects of BFA on a number of cellular membranous organelles indicate both the common and distinct mechanisms that operate to sort membrane components. BFA seems to inhibit a common key step in the mechanism of vesicle budding, namely, regulation of the nucleotide status of different small G proteins. The BFA-sensitive G proteins are each endowed with organelle specificity, coupled perhaps to tissue specificity, which explains the variable sensitivity of membranes to BFA. On the other hand, the use of BFA has also revealed a physiological connection between cellular [signal transduction](#) pathways and the traffic of membrane lipids and proteins through the cell.

## Bibliography

1. Y. Misumi, Y. Misumi, K. Miki, A. Takatsuki, G. Tamura, and Y. Ikehara (1986) *J. Biol. Chem.*, **261**, 11398–11403.
2. J. Lippincott-Schwartz, L. C. Yuan, J. S. Bonifacino, and R. D. Klausner (1989) *Cell* **56**, 801–813.
3. J. Lippincott-Schwartz, J. G. Donaldson, A. Schweizer, E. G. Berger, H. P. Hauri, L. C. Yuan, and R. D. Klausner (1990) *Cell* **60**, 821–836.
4. N. E. Ivessa, C. De Lemos-Chiarandini, Y. S. Tsao, A. Takatsuki, M. Adesnik, D. D. Sabatini, and G. Kreibich (1992) *J. Cell Biol.* **117**, 949–958.
5. R. D. Klausner, J. G. Donaldson, and J. Lippincott-Schwartz (1992) *Cell Biol.* **116**, 1071–1080.
6. L. C. Hendricks, S. L. McClanahan, G. E. Palade, and M. G. Farquhar (1992) *Proc. Nat. Acad. Sci. USA* **89**, 7242–7246.
7. S. G. Miller, L. Carnell, and H. H. Moore (1992) *J. Cell Biol.* **118**, 267–283.

8. J. Lippincott-Schwartz, L. Yuan, C. Tipper, M. Amherdt, L. Orci, and R. D. Klausner (1991) *Cell* **67**, 601–616.
9. L. Orci, M. Tagaya, M. Amherdt, A. Perrelet, J. G. Donaldson, J. Lippincott-Schwartz, R. D. Klausner, and J. E. Rothman (1991) *Cell* **64**, 1183–1195.
10. J. G. Donaldson, J. Lippincott-Schwartz, G. S. Bloom, T. E. Kreis, and R. D. Klausner (1990) *J. Cell Biol.* **111**, 2295–2306.
11. J. G. Donaldson, D. Finazzi, and R. D. Klausner (1992) *Nature* **360**, 350–352.
12. N. T. Ktistakis, H. A. Brown, P. C. Sternweis, and M. G. Roth (1995) *Proc. Nat. Acad. Sci. USA* **92**, 4952–4956.
13. M. Sugai, C. H. Chen, and H. C. Wu (1992) *J. Biol. Chem.* **267**, 21297–21299.
14. M. Sugai, C. H. Chen, and H. C. Wu (1992) *Proc. Nat. Acad. Sci. USA* **89**, 8903–8907.
15. J. G. Donaldson, R. A. Kahn, J. Lippincott-Schwartz, and R. D. Klausner (1991) *Science* **254**, 1197–1199.
16. N. T. Ktistakis, M. E. Linder, and M. G. Roth (1992) *Nature* **356**, 344–6.
17. N. T. Ktistakis, M. G. Roth, and G. S. Bloom (1991) *J. Cell Biol.* **113**, 1009–1023.
18. T. G. Turi and J. K. Rose (1995) *Biochem. Biophys. Res. Commun.* **213**, 410–418.
19. M. S. Robinson and T. E. Kreis (1992) *Cell* **69**, 129–138.
20. A. Bruning, T. Ishikawa, R. E. Kneusel, U. Matern, F. Lottspeich, and F. T. Wieland (1992) *J. Biol. Chem.* **267**, 7726–32.
21. C. M. Linardic, S. Jayadev, and Y. A. Hannun (1992) *J. Biol. Chem.* **267**, 14909–14911.
22. T. Oda, C. H. Chen, and H. C. Wu (1995) *J. Biol. Chem.* **270**, 4088–4092.

### Suggestions for Further Reading

23. R. D. Klausner, J. G. Donaldson, and J. Lippincott-Schwartz (1992) Brefeldin A: insights into the control of membrane traffic and organelle structure. *J. Cell Biol.* **116**, 1071–1080.
24. J. Lippincott-Schwartz, L. C. Yuan, J. S. Bonifacino, and R. D. Klausner (1989) Rapid redistribution of Golgi proteins into the ER in cells treated with brefeldin A: evidence for membrane cycling from Golgi to ER. *Cell* **56**, 801–813.
25. N. T. Ktistakis, H. A. Brown, P. C. Sternweis, and M. G. Roth (1995). Phospholipase D is present on Golgi-enriched membranes and its activation by ADP ribosylation factor is sensitive to brefeldin A. *Proc. Natl. Acad. Sci. USA* **92**, 4952–4956.

### 5-Bromouracil

Litman and Pardee (1) found that three **thymine** analogues (5-bromouracil, 5-chlorouracil and 5-iodouracil) were all incorporated fairly efficiently into the **DNA** of **bacteriophage T4**, and all were powerful **mutagens** in this organism. Further studies on 5-bromouracil (BU) mutagenesis by Freese (2, 3) were important in defining transition and transversion mutations and distinguishing these events from **frameshift mutation**. All the mutants induced by BU could be reverted by either BU or **2-aminopurine**, although spontaneous and **acridine**-induced mutants could not. The BU-induced mutations were suggested to be transitions, because it was considered that the most likely effect of incorporating this base analogue into DNA would be to induce mispairing with guanine, by a mechanism that was unknown at that time (3). Mispairings with guanine do occur and can lead to BU mutagenesis (for example, as shown in Fig. 1).

**Figure 1.** Mechanism of generation of transition mutation by 5-bromouracil (BU) after DNA replication. The presence of BU in one DNA strand in place of the original thymine (T), causes guanine (G) to be inserted in the complementary position upon DNA replication, instead of the original adenine (A).

| 1st generation | 2nd generation | 3rd generation | 4th generation |                                  |
|----------------|----------------|----------------|----------------|----------------------------------|
| A.T →          | A.BU →         | A.BU →         | A.BU           | (Mutations in later generations) |
|                | ↘              | G.BU →         | G.BU           | (Mutations in later generations) |
|                |                | ↘              | G.C            | (Transition mutations)           |

There is some evidence that ionization of bromouracil leads to a high probability of mispairing with guanine and that ionized structures are more likely to be involved in this type of mutagenesis than a neutral [wobble-pairing](#) structure containing the favored keto structures of BU (4). Goodman et al. (5) suggested that BU may exhibit substantially different base-pairing behavior depending upon whether it is present as a template base or as a deoxyribonucleoside triphosphate substrate. They interpreted differences in mutagenesis by this base analogue in different organisms and different systems as due partly to the relative size of and imbalance of the deoxynucleoside triphosphate pools. Classic interpretations of BU mutagenesis invoke direct base mispairing, and this is likely to be involved in the mutagenesis seen in bacterial and bacteriophage systems. However, studies in **yeast** (6) suggest that the majority of BU-induced mutants do not occur at the site where BU is incorporated. Szyszko et al. (7) have identified uracil as a major lesion after BU incorporation into DNA of *E. coli*, presumably because of dehalogenation of incorporated BU. They suggested that this would be followed by formation of apyrimidinic sites by the enzyme uracil-DNA glycosylase, followed by single-stranded nicks. Efficient BU mutagenesis is dependent upon various [DNA repair](#) gene functions, suggesting a major role for [SOS response](#) repair in many of the BU-induced mutations (8).

### Bibliography

1. R. M. Litman and A. B. Pardee (1956) *Nature* **178**, 529–531.
2. E. Freese (1959) *Proc. Natl. Acad. Sci. USA* **45**, 622–633.
3. E. Freese (1959) *Brookhaven Symp. Biol.* **12**, 63–73.
4. H. Yu, R. Eritja, L. B. Bloom, and M. F. Goodman (1993) *J. Biol. Chem.* **268**, 15935–15943.
5. M. F. Goodman, R. L. Hopkins, R. Lasken, and D. N. Mhaskar (1985) *Basic Life Sci.* **31**, 409–423.
6. V. Noskov, K. Negishi, A. Ono, A. Matsuda, B. Ono, and H. Hayatsu (1994) *Mutat. Res.* **308**, 43–51.
7. J. Szyszko, I. Pietrzykowska, T. Twardowski, and D. Shugar (1983) *Mutat. Res.* **108**, 13–27.
8. I. Pietrzykowska, M. Krych, and D. Shugar (1985) *Mutat. Res.* **149**, 287–296.

### Buffers

1. Basic Principles

The term *buffer* solution usually refers to a solution that minimizes changes in pH when hydrogen ions (hydrons,  $H^+$ ) are added to the solution or removed from it. (As discussed in reference [1](#), IUPAC has approved *hydron* to designate a hydrogen cation when a term independent of hydrogen isotope is desired.) Such a solution is therefore said to buffer the pH. Solutions can also be designed to buffer other species, particularly metal ions, as will be described below.

Because the concentration of hydrogen ions in biological media is low, usually about  $0.1 \mu\text{M}$ , it could be greatly changed by the production or consumption of very small quantities of hydrogen ions in chemical reactions. This is what makes buffering commonly necessary in experiments in molecular biology. The principle of a buffer is simple. Relative to the low concentration of free hydrogen ions, the buffer consists of much higher concentrations of both a base,  $A^-$  (ie, a substance that can combine reversibly with the hydrons), and its conjugate acid, HA, which is formed when the base combines with a hydron. Hence the equilibrium



is established. Subsequently, any addition of  $H^+$  causes a much smaller increase in its free concentration,  $[H^+]$ , because some of the hydrons are used up in making more HA to displace this equilibrium to the left. Likewise, any fall in  $[H^+]$  is diminished because it leads to the dissociation of HA.

## 2. Simple Theory

Users of buffers need to be aware of many features of their action if they are not to make mistakes that can easily ruin their experiments. For this, a little theory of buffer solutions is required.

### 2.1. Dissociation of Water

[Water](#) spontaneously dissociates into one hydrogen ion and one hydroxide ion:



Because of the equilibrium, the product  $[H^+][OH^-]$  is constant; its value, about  $10^{-14} \text{ M}^2$ , changes only slightly with temperature. Both the hydrons and the hydroxide ions will be largely hydrated, reversibly combined with one or more water molecules. In this article, the term *acid* is used in the sense of a Brønsted acid, a substance that can donate a hydron, and *base* is used to mean a Brønsted base, a substance that can combine with one. If no such acid or base is added to pure water, the values of  $[H^+]$  and  $[OH^-]$  are equal, namely,  $10^{-7} \text{ M}$ , ie,  $0.1 \mu\text{M}$ . Given that the pH is defined as

$$\text{pH} = -\log_{10}[H^+] \quad (3)$$

the pH is 7. The pH scale is useful in water over about the range 0 to 14, over which  $[H^+]$  varies from 1 M to  $10^{-14} \text{ M}$ , and  $[OH^-]$  from  $10^{-14} \text{ M}$  to 1 M.

### 2.2. The $pK_a$ and the pH

The acid HA will have a dissociation constant,  $K_a$ , the equilibrium constant for reaction [\(1\)](#), which is defined as

$$K_a = \frac{[H^+][A^-]}{[HA]} \quad (4)$$

Rearranging this to find  $[H^+]$  gives

$$[H^+] = \frac{K_a[A^-]}{[HA]} \quad (5)$$

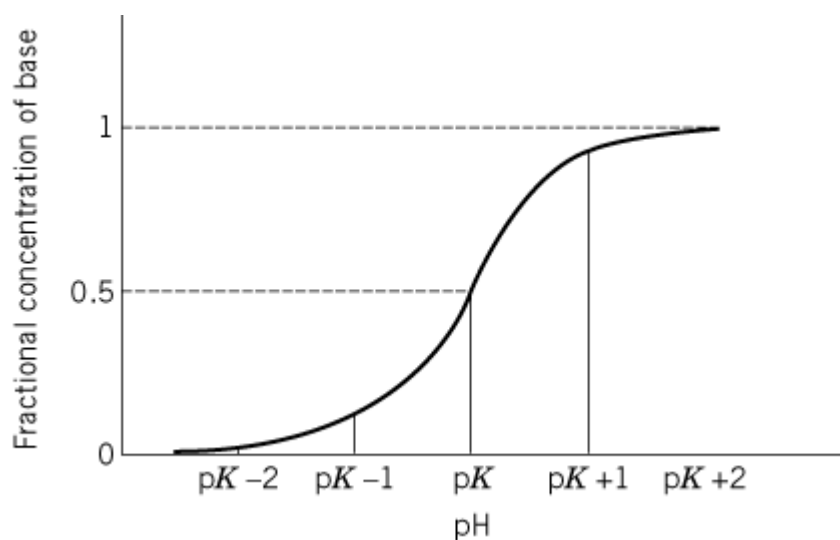
and taking the negative logarithm of each side to find the pH [eq. (3)] gives

$$pH = pK_a - \log \frac{[A^-]}{[HA]} \quad (6)$$

if we define  $pK_a$  as  $-\log K_a$ . This is known as the Henderson–Hasselbalch equation, and it is fundamental for understanding buffer action and how single groups titrate.

Figure 1 shows the fractional concentration of the base, ie,  $[A^-]/([HA] + [A^-])$ , as a function of pH. The first point to note is that, when the  $pH = pK_a$ ,  $\log([A^-]/[HA]) = 0$ , and so  $[A^-]/[HA] = 1$ , and  $[A^-]$  is 0.5 of its maximal value. Hence, the  $pK_a$  is the pH at which the acid is half-dissociated. When the pH is one unit above the  $pK_a$ ,  $\log([A^-]/[HA]) = 1$ , and so  $[A^-]/[HA] = 10$ , and hence  $[A^-]$  is 10/11, ie, 0.91, of its maximal value as shown; when the pH is two units above the  $pK_a$ ,  $\log([A^-]/[HA]) = 2$ , and so  $[A^-]/[HA] = 100$  and, hence,  $[A^-]$  is 100/101, ie, 0.99 of its maximal value. Similarly, when the pH is one unit below the  $pK_a$ ,  $\log([A^-]/[HA]) = -1$ , and so  $[A^-]/[HA] = 0.1$ , and  $[A^-]$  is 0.091 of its maximal value; when the pH is two units below the  $pK_a$ ,  $[A^-]$  is 0.01 of its maximal value. These figures illustrate that comparing the pH with the  $pK$  of a substance instantly indicates how much of the substance or group is in its hydronated and unhydronated forms. All simple acids follow the curve of Figure 1; they differ only in their  $pK_a$  values.

**Figure 1.** The dependence of the fractional concentration of the basic form of a monobasic acid on the pH. The vertical axis shows the concentration of  $A^-$  expressed as a fraction of the sum of the concentrations of  $A^-$  and HA. Vertical lines mark: (1)  $pH = pK_a - 1$ , when  $[A^-] = 0.1[HA]$ , so the fractional concentration of  $A^-$  is 1/11; (2)  $pH = pK_a$ , when  $[A^-] = [HA]$ , so the fractional concentration of  $A^-$  is 1/2; (3)  $pH = pK_a + 1$ , when  $[A^-] = 10[HA]$ , so the fractional concentration of  $A^-$  is 10/11.



### 2.3. Theory of Buffering

With properly buffered solutions, both  $[A^-]$  and  $[HA]$  will be much greater than  $[H^+]$  and  $[OH^-]$ . In this case, any  $H^+$  produced in a reaction converts  $A^-$  into  $HA$ , and any  $H^+$  taken up dissociates from  $HA$  to convert it into  $A^-$ . To see the consequences of this, it is convenient to modify the Henderson–Hasselbalch equation as follows:

$$pH = pK_a - \log \frac{[A^-]}{[HA]} = pK_a - \log[A^-] + \log[HA] \quad (7)$$

For good buffering, it is necessary for *both*  $[A^-]$  and  $[HA^-]$  to be large enough so the logarithm of neither is changed enough to produce an unacceptable change in pH when some of one is converted into the other.

This means that there are two basic requirements for a buffer: (1) It must be sufficiently concentrated (in terms of the total of  $[A^-] + [HA^-]$ ), and (2) the pH must not be too far from the  $pK_a$  of the buffer so that the concentrations of *both*  $A^-$  and  $HA$  will be large. It is the second of these that is most often forgotten. If the pH and  $pK_a$  differ by one unit, only 1/11 of the total buffer concentration will be in the minor form.

This point may be illustrated by the example of a worker who wondered why mercaptoethanol and **iodoacetate** inhibited his enzyme when neither one by itself did so. He was using more than 10 mM of each reagent in a 10-mM Tris buffer of pH 7. Hence, the 10-mM  $H^+$  released when the mercaptoethanol and iodoacetate reacted exceeded by tenfold the concentration of free Tris, which was less than 1 mM in a solution far below its  $pK_a$  of 8.1.

To a first approximation, dilution does not affect the pH of a buffer solution because it does not affect the ratio  $[A^-]/[HA^-]$ . There are, however, two qualifications to this conclusion. First, it is assumed that the dissociation constant of the buffering acid is unaltered whereas, in fact, it is likely to change with ionic strength (see below). It is also assumed that dilution is not so extreme that either  $[A^-]$  or  $[HA^-]$  become comparable to  $[H^+]$  or  $[OH^-]$ . In any case, dilution decreases the concentrations of both  $A^-$  and  $HA$ ; hence dilution makes the solution a worse buffer in that a given amount of conversion of one form into the other will produce a larger change in the ratio.

## 3. Practical Points

### 3.1. pH Measurement

Methods of measuring pH are outside the scope of this article but the process is important because reproducibility of buffer preparation often depends on it (see below). The most common method uses a glass electrode. This method depends on dipping the electrode into the solution to be tested in such a way that an electrical cell is created and that the potential given by this cell depends on the pH of the solution. Normally, a layer of thin glass separates the solution under test from one containing one of the electrodes of the cell. This glass is selectively permeable to hydrons, so that the potential across it depends on the pH of the solution. Provided that the other junction potential between this solution and the other electrode is negligible—a condition approached by a bridge containing concentrated KCl—the potential should reliably indicate the pH.

Such pH meters require calibration, and standard buffers are used for this. The makers normally give reliable instructions for calibrating the electrode, but some of these are easily overlooked. The first concerns the temperature. The pH standards have pH values that may themselves vary with temperature. Standardization is valid at only one particular temperature, and it cannot be assumed that an electrode standardized at one temperature gives a meaningful reading at a different one. The pH range specified by the makers of the apparatus should also be noted. Glass membranes may be

slightly permeable to  $\text{Na}^+$ , so that when  $[\text{H}^+]$  is low, as in alkaline solution, even a slight sensitivity to  $\text{Na}^+$  may cause overestimation of  $[\text{H}^+]$  and hence too low an estimate of a high pH.

### 3.2. Specifying a Buffer Solution

Buffer solutions are often specified in a conventional way, such as 0.2-M sodium acetate buffer, pH 4.8. The concentration given should refer to the sum of the two forms of the buffering species, designated here as  $[\text{AcO}^-]$  and  $[\text{HOAc}]$  in the case of acetate buffer, but such a specification does not make fully clear how this solution has been obtained. It might have been prepared by (1) mixing the appropriate amounts of sodium acetate and acetic acid, as calculated from the  $\text{p}K$  and the Henderson–Hasselbalch equation, (2) adding 0.2-M sodium acetate to 0.2-M acetic acid until a pH of 4.9 was reached, or (3) adding strong NaOH to 0.2-M acetic acid to reach pH 4.9, on the assumption that the volume added diluted the acetic acid negligibly. It would not have been appropriate, however, to have added strong HCl to 0.2-M sodium acetate, as this would have produced NaCl in addition to the buffering species.

Measuring pH is not straightforward, especially in view of the effects of temperature on the measurements (see above); for reproducibility, especially between different laboratories, it is much safer to specify the concentration of *each* form of the buffering substance, given that weighing material out with an accuracy of 1% to 2% is easy. Even an error of 2% in each would change the pH by only  $\log(1.04)$ , ie, by 0.017 of a pH unit. Hence, it would be better to describe the buffer just mentioned as 0.12-M sodium acetate, 0.08-M acetic acid (pH 4.8), thus adding the pH for information rather than as part of the definition.

### 3.3. Effects of Temperature

Users should remember that the pH of a buffer may change with temperature. This is because the dissociation constant of the buffering acid, like any other equilibrium constant, depends on temperature if the dissociation has an appreciable **enthalpy** change,  $DH$ . Dissociations of carboxylic acids and phosphoric acid have low values of  $DH$ , and so this effect is small with buffers based on them, but the dissociation of ammonium to ammonia and a hydron is accompanied by a large heat intake (ie,  $DH$  is positive), and increasing temperature promotes dissociation, ie, lowers the  $\text{p}K_a$ .

Hence, the pH of a buffer that uses an amino group decreases on heating. The effect is quite large, a fall of 0.028 per degree Celsius.

The importance of this needs to be considered in each experiment. If a user wished the glutamate residues of an enzyme to be hydronated to the same extent in an experiment at 0 °C as they would be in a buffer of pH 7 at 20 °C, it would be wrong to use an amino buffer with a pH of 7 at 20 °C and simply cool it. Cooling would raise the pH of the buffer, but not the  $\text{p}K_a$  of the residues, so these residues would dissociate more on cooling. But this would be a proper procedure if the experimenter was concerned with the state of the [lysine](#) residues because their  $\text{p}K_a$  would rise with that of the buffering amine so that their degree of hydronation would remain the same.

### 3.4. Effects of Concentration

The dissociation constants referred to above are those actually exhibited by the buffering species in the solution being used; ie, they are not idealized values extrapolated to standard conditions of ionic strength and temperature. They therefore change with ionic strength and the addition of other components to the solution, eg, organic solvents or deuterium oxide. Many experiments may require a solution with a particular ionic strength or particular concentrations of specified ions. Remember that these may affect the apparent value of the  $\text{p}K$  of the buffering species.

An acid and its conjugate base differ by one  $\text{H}^+$  and therefore by one charge. Typically, they will have charges 0 and  $-1$ , or  $+1$  and 0. If, however, one of the forms has more than unit charge, positive or negative, the  $\text{p}K$  will have a greater dependence on ionic strength. The most commonly used buffer to which this applies is [phosphate buffer](#), where the acid is  $\text{H}_2\text{PO}_4^-$  and its conjugate base is

the doubly charged  $\text{HPO}_4^{2-}$ . This makes its  $pK_a$  particularly sensitive to ionic strength. For example, a solution of 0.2-M phosphate may increase its pH by over 0.2 of a unit on tenfold dilution.

Many chemical reactions are catalyzed by acids or bases, so that an increased buffer concentration can accelerate them even at a fixed pH. This should be borne in mind when altering buffer concentrations. It is not very common for enzyme-catalyzed reactions because natural selection often ensures that if such catalysis accelerates the reaction, a corresponding acidic or basic group will be present in the [enzyme](#) to provide it.

### 3.5. Volatile Buffers

Separation methods such as [chromatography](#) and [electrophoresis](#) often require a buffered solution. This solution needs appreciable ionic strength (in [ion-exchange chromatography](#) to give a constant competition for the sites of the exchanger, in electrophoresis to give a constant conductivity and therefore electric field) that is negligibly affected by the substances being separated. Because it is convenient to be able to remove the buffer at the end of the procedure, a volatile buffer that can be removed by evaporation is desirable. The two forms of the buffering substance must differ in charge, however, so that one at least will have a net charge and hence be nonvolatile. Nevertheless, volatile buffers can be prepared by mixing a volatile acid with a volatile base, provided that their  $pK_a$  values are not greatly separated. An example is a mixture of pyridine ( $pK_a$  5.2) and acetic acid ( $pK_a$  4.8). On mixing, some of the acetic acid dissociates to acetate and hydronates some of the pyridine to pyridinium. Hence, the pH will be given by

$$\text{pH} = 5.2 + \log \frac{[\text{pyridine}]}{[\text{pyridinium}]} = 4.8 + \log \frac{[\text{acetate}]}{[\text{acetic acid}]} \quad (8)$$

The concentrations of pyridinium ( $\text{C}_5\text{H}_5\text{N}^+\text{-H}$ ) and acetate ( $\text{AcO}^-$ ) ions will be equal, because each acetic acid molecule that dissociates to acetate forms one pyridinium ion. As the solution is evaporated off, water, pyridine ( $\text{C}_5\text{H}_5\text{N}$ ) and acetic acid ( $\text{AcOH}$ ) evaporate, and so the equilibrium



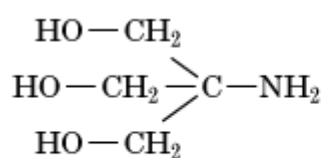
is displaced to the right to replace the un-ionized forms as they are removed.

Ammonium bicarbonate forms a useful buffer of this type. Although it is highly volatile in solution, it is not volatile when dry, as it has a stable crystal lattice. Hence to allow ammonium and bicarbonate to react to form ammonia and carbon dioxide, it may be necessary after one drying to add a little water and then dry it again. Alternatively, triethylammonium bicarbonate, prepared from triethylamine and carbon dioxide, may be used, as it does not form nonvolatile crystals. Ammonium acetate and ammonium formate have been used in this way (2), but the procedures for drying them are complicated because their  $pK_a$  values are so far apart that there is very little of the un-ionized forms at equilibrium.

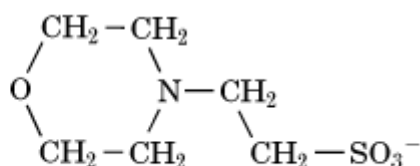
### 3.6. Appropriate $pK_a$ Values

Biological material usually has a pH near 7, and it is not always easy to find a buffering species with the desired  $pK_a$ . Carboxylic acids possess values that are too low, whereas amines have values that are too high. A large number of the commonly used buffers are amines with electron-attracting substituents, which make the lone pair of electrons on the nitrogen atom less available and so lower the atom's  $pK_a$ . Examples include:





Tris



2-Morpholinoethanesulfonic acid

(10)

with  $pK_a$  values of 8.1 and 6.2, respectively (see [Tris Buffer](#)). [Phosphate Buffers](#) are commonly used because they have a convenient  $pK_a$  of about 7 between  $\text{H}_2\text{PO}_4^-$  and  $\text{HPO}_4^{2-}$ . Their main disadvantages are their propensity to support fungal and algal growth and the sensitivity of their  $pK_a$  to ionic strength (see above).

### 3.7. The Need for Buffers

It should be remembered that it is pointless to add a buffer if the components of the reaction mixture are already highly buffering. Glycolytic intermediates, for example, are phosphate esters, and adjusting them to a pH near 7 creates a buffered solution.

### 3.8. Chemical Reactivity

It is often important to use buffering species that will not react chemically under the conditions of an experiment. Hence an experiment to acylate [amino groups](#) of proteins should use a buffer whose basic component has as low a nucleophilic reactivity as possible. Exclusion of primary amines, such as Tris, may be all that is needed. But extreme conditions may demand introduction of steric hindrance; eg, 2,6-dimethylpyridine has a  $pK_a$  close to that of pyridine but vastly less nucleophilic reactivity because the methyl groups have little steric effect on protonation of the nitrogen atom but greatly slow its reaction with larger molecules.

Another form of reactivity is the possibility of precipitating desired components in the solution, eg, cations such as  $\text{Ca}^{2+}$ , by phosphate. Occasionally, reactivity may even be of advantage. [Urea](#) is often used as a protein **denaturant** but, if kept for long in concentrated solution, especially at raised temperatures, its equilibration



can produce about 20-mM cyanate ( $\text{CNO}^-$ ) from 8 M urea (3). The harm in this is the reverse reaction, because the cyanate can carbamoylate [lysine](#) residues in proteins, forming *uncharged* ureas and destabilizing the folded protein. This process is minimized if ammonium ions are present in the buffer as they compete with the protein for the cyanate.

### 3.9. “Good” Buffers

A number of buffers were recommended by Good and colleagues (4) for a variety of reasons. The most important was that both acidic and basic forms were charged and, therefore they could not rapidly permeate biological [membranes](#). If **organelles** are suspended in a buffer without this feature, acetate, for example, molecules of acetic acid, but not acetate ions, may penetrate them and lower the internal pH below that of the external buffer. Many of the Good buffers have strongly acid groups, such as sulfonate,  $-\text{SO}_3^-$ . An example is 2-morpholinoethanesulfonate (eq. 10); its sulfonate group renders it impermeant even when its nitrilo group is un-ionized. Other desirable characteristics of buffering species are also listed (4). For many purposes, transparency at UV wavelengths is necessary.

Under other circumstances, the possession of a further charge can be a disadvantage. Ideally, buffers for ion-exchange chromatography possess only one species with a charge opposite in sign to that of

the exchanger so that the equilibration of the exchanger can be followed by pH changes; for example, a buffer cation cannot adsorb to, or desorb from, the exchanger except by exchange with  $H^+$ , which changes the pH.

#### 4. Metal-Ion Buffers

The free concentration of many metal ions in biological media is very low, and it may be necessary to buffer it. The principles are exactly the same as those that apply to buffering  $[H^+]$ . Some compound that binds the metal ion is introduced in roughly equal concentrations of the free and metal-bound forms. A substance that buffers in the correct concentration range can be chosen from a tabulation of binding constants (5).

A complication is that many of the compounds that ligate metal ions with suitable affinity also bind hydrons, resulting in competition between the cations for the ligand. This means that the affinity of the ligand for the metal ion varies with pH. The complex of the buffering species with the metal ion is often a chelate, so that two or more ligating groups are involved. Hence, the competition is not a simple 1:1 replacement, and the pH dependence of the net affinity may be complex.

#### Bibliography

1. J. F. Bunnett and R. A. Y. Jones (1988) *Pure Appl. Chem.* **60**, 1115–1116.
2. C. H. W. Hirs, S. Moore, and W. H. Stein (1952) *J. Biol. Chem.* **195**, 669–683.
3. G. R. Stark, W. H. Stein, and S. Moore (1960) *J. Biol. Chem.* **235**, 3177–3181.
4. N. E. Good and S. Izawa (1972) *Methods Enzymol.* **24**, 53–68.
5. R. M. Smith and A. E. Martell (1974-89) *Critical Stability Constants*, vols. **1–6**, Plenum Press, New York.

#### Suggestion for Further Reading

6. R. J. Beynon and J. S. Easterby (1996) *Buffer Solutions: The Basics*, IRL Press Oxford, U.K. This gives not only an excellent account but also WWW addresses for programs to calculate buffer compositions (<http://www.bi.umist.ac.uk/buffers.html>) and other help on the Internet.

#### bZip Domain

The bZip DNA binding domain and the closely related **helix-loop-helix** (HLH) domain are found in many eukaryotic [transcription](#) factors. Both bind in the major groove of DNA and have the capability of forming homo- or heterodimers with [proteins](#) of the same class. The simplest of these domains is the bZip domain which, when bound to DNA, forms a continuous  $\alpha$ -helix. The lower basic part of this helix (the b of bZip) forms sequence specific contacts with DNA whereas the upper part can dimerize with an appropriate partner through the formation of a leucine zipper. The dimer thus has the appearance of a pair of scissors binding to a palindromic DNA sequence in which each half-site is exposed on the opposite face of DNA. The helix-loop-helix domain binds to DNA in a similar manner with the difference that the  $\alpha$ -helix involved in dimerization is separated from the DNA binding  $\alpha$ -helix by a short loop.

## Buoyant Density

The buoyant density of a macromolecule is its effective density in solution. It is defined relative to the fluid medium in which it resides and/or through which it is sedimented (see [Centrifugation](#)). The dimensions of buoyant density are in  $\text{cm}^3/\text{g}$ . If a macromolecule is sedimented through a fluid column in which the density of the column is adjusted so that it spans a gradient of density *both* less than *and* greater than the macromolecule, the macromolecule moves to a fixed position in the fluid column (see [Density Gradient Centrifugation](#)). At this position the density of the macromolecule is equal to that of the solvent. This is an equilibrium position that is independent of time. This process or technique is termed density gradient sedimentation-equilibrium, equilibrium banding or isopycnic gradient centrifugation. These techniques have uses that range from the analysis of macromolecule solvation to the separation of both living cells and subcellular **organelle** systems (1). Recently the evaluation of macromolecular values of buoyant density have been employed in a variety of research situations including (1) the analysis of viral infectivity during therapy (2, 3); (2) the nature of circulating immune complexes of **viruses** (4); (3) as a preselection method for **PCR** methods (5).

### Bibliography

1. S. Daenen, W. Huiges, E. Modderman, and M. R. Halie (1993) *Leukemia Res.* **17**, 37–41.
2. A. Nagasaka, S. Hige, T. Matsushima, J. Yoshida, Y. Sasaki, I. Tsunematsu, and M. Asaka (1997) *J. Med. Virol.* **52**, 190–194.
3. T. Kanto, N. Hayashi, T. Takehara, H. Hagiwara, E. Mita, M. Naito, A. Kasahara, H. Fusamoto, and T. Kamada (1995) *J. Hepatol.* **22**, 440–448.
4. M. Hijikata, Y. K. Shimizu, H. Kato, A. Iwamoto, J. W. Shih, H. J. Alter, R. H. Purcell, and H. Yoshikura (1993) *J. Virol.* **67**, 1953–1958.
5. R. J. Carrick, G. G. Schlauder, D. A. Peterson, and I. K. Mushahwar (1992) *J. Virol. Meth.* **39**, 279–290.

### Suggestions for Further Reading

6. C. H. Chervenka (1973) *A Manual of Methods for the Analytical ultracentrifuge*, Spinco Division, Beckman Instruments, Inc. Palo Alto, CA. *The classic manual of methods for the original model E analytical ultracentrifuge*. Although the mechanical components related to model E have been superseded by the modern XLA and XLI analytical ultracentrifuges, the methods themselves are still applicable.
7. R. Hinton and M. Dobrota (1980) *Laboratory Techniques in Biochemistry and Molecular Biology: Density Gradient Centrifugation* (T. S. Work and E. Work, eds.), North-Holland, New York. A somewhat dated, but useful, compendium of techniques for preparative density gradient centrifugation.

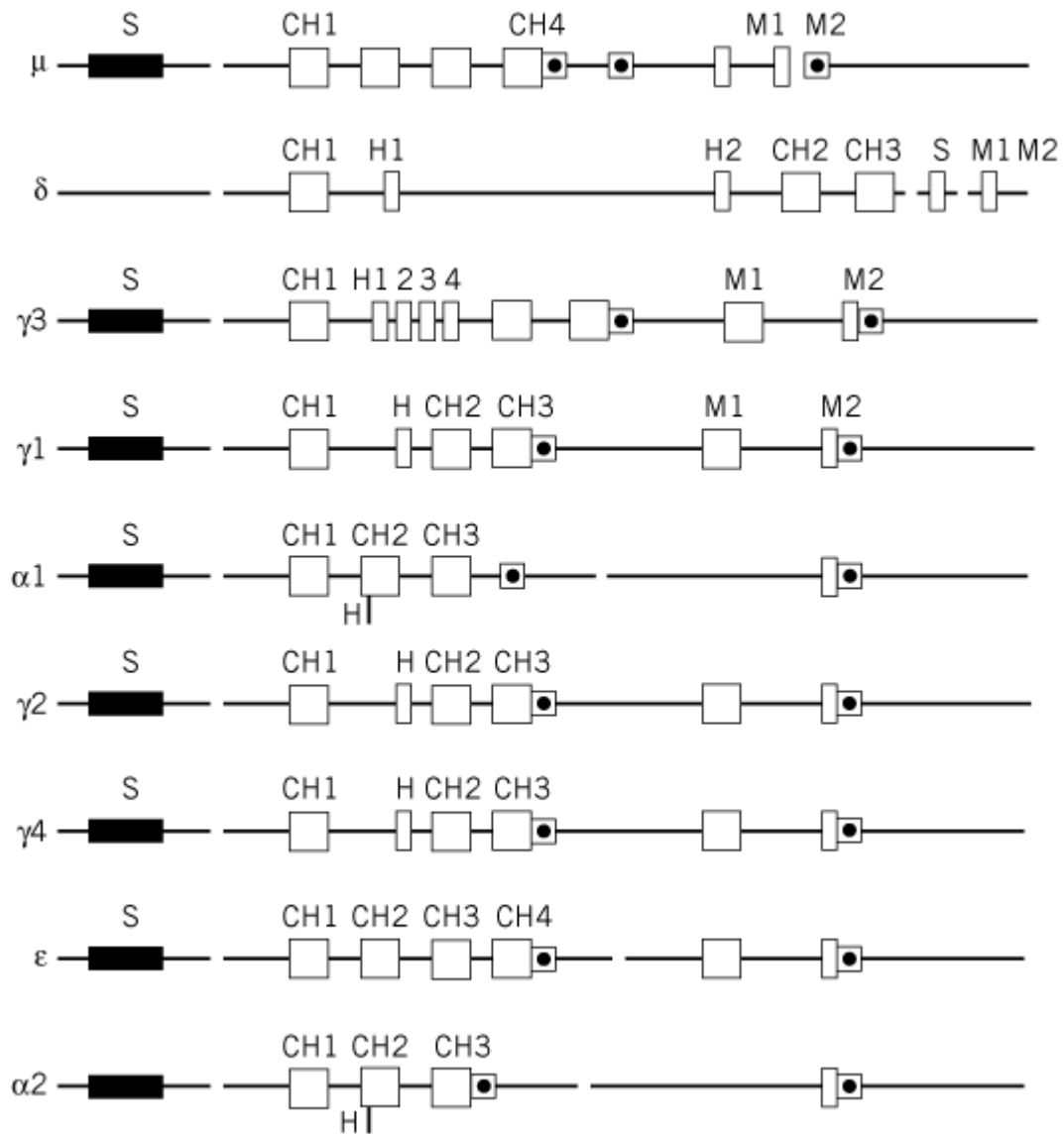
## C Genes of Immunoglobulins

C genes encode for the [constant \(C\) regions](#) of [immunoglobulin](#) (Ig) heavy and light chains and therefore define the [isotypes](#) of H (g, m, a, d and l) and L (k, l) chains. The simplest structure is that

of the Ck gene, which is unique in both the human and mouse species. It encodes for the entire constant region of the k chain. In humans, the Ck gene has three **allelic** variants based on single amino acid substitutions at two positions, which constitute the Km allotypes. The Ck gene is separated from the Jk region by a noncoding region of 1.2 kbp that contains a k [enhancer](#). A second enhancer region is present in the 3' flanking region. The organization of the Cl genes is somewhat more complicated, and there are some minor differences between humans and mouse. In both cases, they are present as tandem duplications of a JICl unit, in which each J gene is separated from the C gene by a noncoding sequence of about 1 kbp. There are four discrete tandems in the mouse, JI1C11 to JI4C14, of which only three are functional. In humans, four functional tandem JC pairs are present, in addition to three pairs of [pseudogenes](#). An enhancer region is present 3' of the coding region in both species. In humans, the IgK and IgL loci are located at **chromosomal** positions 2p12 and 22q11, respectively.

The heavy-chain constant gene locus, located at 14q32 in humans, spans over 200 kbp. The detailed organization of the individual human CH genes is given in Figure 1. Some characteristic features may be identified. First, each isotype is encoded by a mosaic of exons (see [Introns, Exons](#)) that precisely reflects the **domain** structure of the Ig. The constant region of Cm is encoded by four exons, Cm1 to Cm4, corresponding to the four constant domains of the m chain, whereas in the case of the Cg isotypes one of the exons encodes the hinge region. Second, each isotype encoding C gene has additional membrane exons (one or two, depending on the isotype) that encode the transmembrane section and the short COOH-terminal cytoplasmic tail. Coding units that make secreted or membrane heavy chains result from [alternative splicing](#) and terminate at two distinct [polyadenylation](#) sites (see Fig. 1). The overall CH region also contains two pseudogenes. It will be noticed that a switch region is located at 5' of each set of C genes, with the exception of the d locus, which accounts for the coexpression of [IgM](#) and [IgD](#) at the surface of mature B cells.

**Figure 1.** Organization of human heavy-chain constant region genes. Exons encoding discrete constant domains (CH), hinge regions (H), or transmembrane portion of surface heavy chains (M) are represented by open squares or rectangles. Polyadenylation sites are indicated by a dot. (From Ref. 1, with permission; adapted from Ref. 2.)



A final remark should be made regarding the general organization of the regulating elements that control IgH gene expression. **Promoter** regions are present at 5' of each V gene and will be discussed elsewhere. Enhancers are split into two regions; one is located at 5' of the Sm switch region, so that with each switching event it will control the newly adjacent isotype gene; the other is located at 3' of the entire IgCH locus and may thus control whatever gene is to be turned on.

See also entries [Antibody](#), [Class Switching](#), [Isotype](#), and [Switch Region](#).

#### Bibliography

1. J.-P. Revillard (1998) In *Immunologie*, De Boeck Université, Paris-Bruxelles, p. 57.
2. J. D. Capra and J. D. Duker (1989) *J. Biol. Chem.* **264**, 12745.

#### Suggestions for Further Reading

3. V. Giudicelli et al. (1997) IMGT, the international ImMunoGeneTics database. *Nucleic Acids Res.* **25**, 206–2011.
4. T. Honjo, F. W. Alt, and T. H. Rabbitts (1989) *Immunoglobulin Genes*, Academic Press, New York.

5. F. Matsuda and T. Honjo (1996) Organization of the human immunoglobulin heavy chain locus. *Adv. Immunol.* **62**, 1–29.

## C-Banding

The procedure of C-banding is used to highlight regions of constitutive [heterochromatin](#). In man, sites of C-banding include the [centromeres](#) of all of the [chromosomes](#) and large segments of heterochromatin in chromosomes 1, 3, 9, 16 and the [Y-Chromosome](#). The DNA sequences in constitutive heterochromatin are highly **repetitive** and include many [satellite DNA](#) sequences that are thought to be genetically inert. The distribution of satellite DNA sequences in chromosomes varies considerably from species to species and can even provide information that is specific to a particular individual.

C-banding involves the depurination of DNA in the chromosomes by treatment with hydrochloric acid. The sugar-phosphate backbone of DNA remains intact. Treatment with hydrochloric acid (0.2 M) is halted before depurination is complete (a treatment of several minutes). Then the chromosomes are incubated for several minutes with either sodium hydroxide (0.07 M) or barium hydroxide (0.07 M). This **denatures** the DNA, which aids solubilization. Then the chromosomes are washed overnight with  $2 \times$  SSC (0.3 M NaCl plus 0.03 M trisodium citrate) at 66°C. This warm salt wash breaks the sugar-phosphate backbone and allows DNA fragments to dissolve. Then the chromosomes are stained with **Giemsa**. Staining relies on the retention of intact nucleoprotein complexes in the highly compacted heterochromatin.

C-banding of metaphase chromosomes is greatly enhanced in resolution by two methodological improvements: first, by using chromosomes isolated early in **metaphase**, when the contraction of chromatin is incomplete; second, by the use of a sophisticated microscopic technique known as epi-illumination, in which the chromosome is viewed from the same side as it is illuminated. This provides enhanced contrast between stained and unstained bands.

### Suggestion for Further Reading

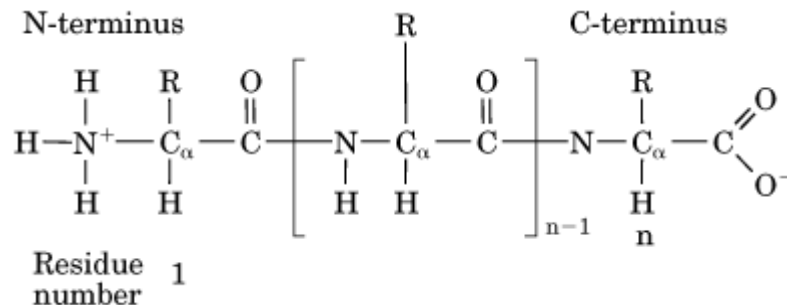
R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A Synthesis*. Wiley-Liss, New York.

## C-Terminus

The C-terminus is the term for one end of the [polypeptide chain](#) of a [protein](#). Proteins are [polymers](#) formed by condensation of the [amino groups](#) and [carboxyl groups](#) of [amino acids](#). That end of the polypeptide chain having an uncondensed or “free” carboxyl group is called the carboxyl terminus, or **C-terminus** (Fig. 1). By convention, the C-terminus is the last residue in the protein sequence or [primary structure](#). Similarly, the other end of the polypeptide chain, having an uncondensed or free amino group, is termed the amino terminus, or **N-terminus**, and by convention it is the first residue in the protein sequence. The C-terminal carboxyl group is usually negatively charged at

physiological pH.

**Figure 1.** Schematic representation of a protein, showing the C-terminus as the last residue in the polypeptide chain. The square brackets indicate the repeating part of the chain (the backbone), and the R group denotes the variable side chain of each amino acid residue. The N-terminus is the first residue in the chain.



The C-terminal carboxyl group can be changed by [post-translational modification](#) such as **amidation**. An amidated C-terminus is more common in smaller proteins, or [peptides](#), including [hormones](#) and [toxins](#), and may be important for biological function or stability. Another important covalent modification that occurs only at the C-terminal residue is the attachment of [membrane anchors](#) to secure an otherwise water-soluble protein to a hydrophobic [membrane](#). Proteins can be degraded by [carboxypeptidases](#) that specifically hydrolyze the C-terminal peptide bond.

[See also [N-Terminus](#) and [Polypeptide Chain](#).]

## C-Value

The C-value is the long-used term for the mystery, or paradox, of the wide range of [genome](#) sizes present among **eukaryotes**. They range from 13 Mbp (megabase pairs,  $1.3 \times 10^7$ ) for **yeast** to  $1.7 \times 10^{11}$  bp for *Amphiuma*, or the mudpuppy, a urodele amphibian. The units used above, the number of nucleotide pairs in one [haploid](#) set of DNA, should be used consistently. In much of the literature, the amount of DNA per cell is expressed in picograms (1 pg is  $9.8 \times 10^8$  nucleotides), and the classic tables are often in error due to the limited accuracy of the methods used. The most accurate genome sizes are those for which complete nucleotide sequences are available; for example, the yeast genome has 13,105,020 bp (base pairs), give or take a small uncertainty because not all copies of repeated sequences (see [Repetitive DNA](#)) have been sequenced. The smallest bacterial genome accurately known from sequence is that of *Mycoplasma genitalium*, with 580,073 bp. Viruses have even smaller genomes, but they are not free-living and depend on their hosts for much of the required information. A part of the range in DNA content is apparently due to the different requirements of different species for genetic information, but that is not all there is to it, by any means.

The familiar fruit fly, **Drosophila**, is a complete animal, with all of the necessary structures, and it is not surprising that it has more DNA than yeast, which has an apparently simpler structure. The genome of *Drosophila melanogaster* is about 170 Mbp. It is easy to imagine that more genes and associated regulatory systems are needed to program the development and specify the complex

structure of an animal. That is clearly true, but life is not so simple. The yeast genome has little space between its **genes** and almost no **introns**, and it is this compactness of organization that accounts for the relatively small size of its genome. The *Drosophila* genome, and plant and animal genomes in general, do not share this very compact organization, although some are more compact than others. They typically have large intergenic spaces, and most of the genome is not made up of genes and regulatory systems and has no known function. Most mammals, including humans, have about  $3 \times 10^9$  bp, and it is not obvious why we require 20 times more DNA than does *Drosophila*. Some believe it is required for the brain, but humans have no more DNA than do mice. It is not reasonable that the mudpuppy requires almost 50 times more genetic information than we do, and its large genome must have some other explanation. One might expect that the very large urodele genomes are burdened with repeated sequences, and they do have a larger fraction of repeated sequences (90%) than *Drosophila* (30%). That is not the only source of the large genome size, however, since the so-called single-copy DNA is also very large in absolute quantity in the large urodele genomes. The genome sizes of single-celled protozoa range from about  $1.0 \times 10^9$  to  $1.0 \times 10^{12}$ , but that information does nothing to reduce the mystery.

The “C-value paradox” remains as paradoxical as when it was first observed. It is known that a more compact genome is typically more compact in many features, such as smaller introns and smaller amounts of repeated sequences. [Natural selection](#) forces controlling genome size must exist, but remain unknown. How genomes have evolved is not well known, but it seems that there have been events of growth, some due to doubling of the entire set of chromosomes, others due to increases in the interspersed repeated sequences. Extraordinary reductions in DNA have also occurred in some species, while closely related species have retained their large genomes. There is a correlation between genome size and cell size, but it is not known how the two are connected, although cell size does affect metabolic rates. All that can be said is that genome size is not a parameter that is understood, although it is nevertheless intimately connected with gene structure and organization. That mammals have a quite uniform genome size is evidence for selective forces that control the C-value. Some geneticists consider the extra DNA in large genomes as a possible limitation in searching for genomic and regulatory function and have favored organisms with relatively small genomes, such as the puffer fish (*Fugu rubripes*) and a weed ([Arabidopsis thaliana](#)). Of course, the small genomes do reduce the amount of work required to determine the complete genome sequence.

#### Suggestions for Further Reading

T. Cavalier-Smith, ed. (1985) *The Evolution of Genome Size*, John Wiley & Sons, Chichester, U.K.

B. John and G. L. G. Miklos (1988) *The Eukaryote Genome in Development and Evolution*, Allen & Unwin, London.

#### CAAT Box

A CAAT box is a *cis*-acting regulatory element in [transcription](#) that contains the sequence 5'-CCAAT-3' at its core and is bound by one or more different [transcription factors](#). These include but are not limited to C/EBP and various members of the CTF/NF1 family of proteins. Because of the multiplicity of transcription factors that recognize a CAAT box, it is usually not possible to determine the identity of the particular factor that regulates expression of a gene by simple sequence inspection. CAAT boxes are found upstream of a number of genes transcribed by **RNA polymerase II** and increase the transcriptional rate of the gene above the level that could be achieved by the



transcriptional core machinery in the absence of a CAAT box. These *cis*-acting sequence elements are generally found upstream of the [TATA box](#), but their exact position relative to it is variable, and they may be found in either orientation relative to the direction of transcription from the downstream initiation site.

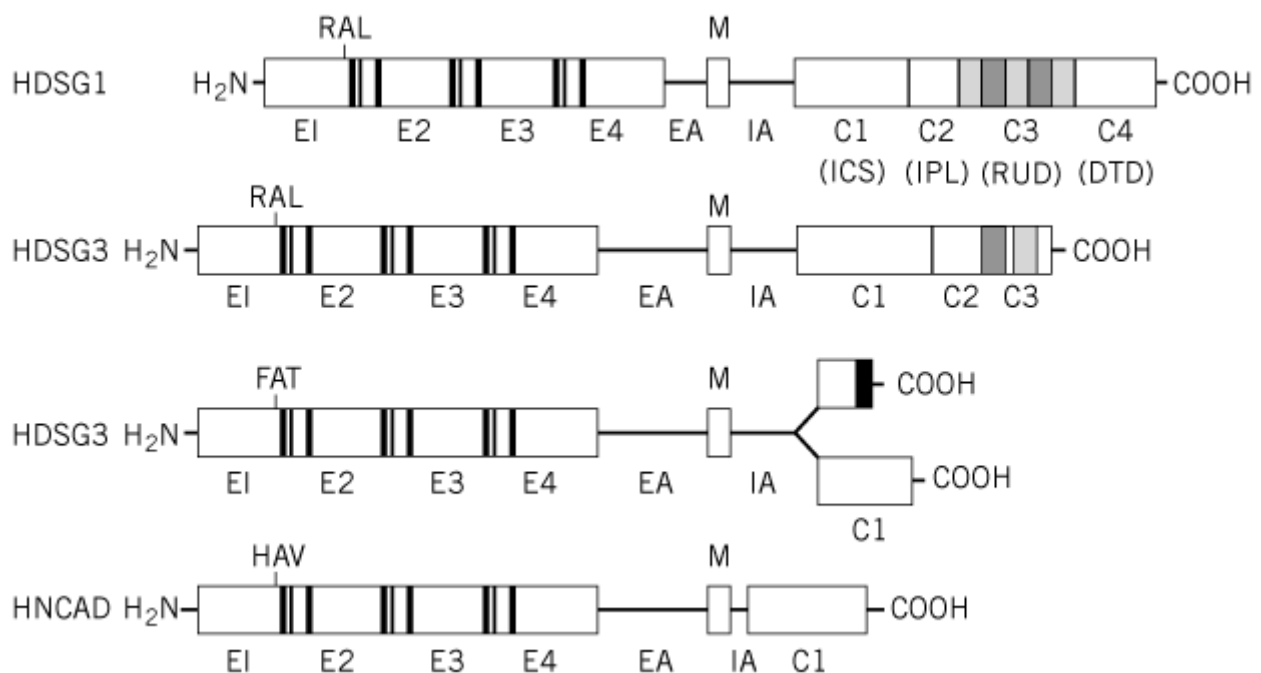
## Cadherins

The cadherins are a family of **transmembrane** proteins that regulate cell–cell adhesion in a  $\text{Ca}^{2+}$ -dependent manner during **embryogenesis** and as the animal matures ([1](#)). Three general groups of cadherins exist within the superfamily: (i) the “classical” cadherins, (ii) the desmosomal cadherins, and (iii) the protocadherins ([1](#)). The classical cadherins, of which there are more than 40 examples, display much conservation in their structures. The sequences of their cytoplasmic domains may be as much as 90% identical for proteins within the same species and up to 60% between widely divergent species. In contrast, there is more variation in sequence within the extracellular domains, and 30% to 60% identity is common. These cadherins are located at the adherens junction and at other points of contact between cells that lack a highly organized structure. The desmosomal cadherins, on the other hand, are the  $\text{Ca}^{2+}$ -dependent [cell adhesion molecules](#) found in **desmosomes**. Within the desmosomal cadherins there are two subclasses: the desmogleins and the desmocollins. **Isoforms** of both have been characterized. The desmogleins and desmocollins are structurally distinct, implying different roles *in vivo*. The third group of cadherins, the protocadherins, are more variable in their structures, and their claim for inclusion in the cadherin family lies with the repeating structures in their extracellular domains that are similar to those displayed by the classical and desmosomal cadherins. The protocadherins also have cell adhesion properties, which again link them to the cadherin family.

Cadherins are classified as Type I integral **membrane glycoproteins** (ie, their *N*-termini lie on the extracellular side of the cell membrane) and, with one exception (glycosylphosphatidylinositol-anchored T-cadherin), each molecule contains a single  **$\alpha$ -helical** transmembrane-spanning region (M) separating the cytoplasmic from the extracellular domain ([1](#)). The extracellular domain generally contains four quasi-repeating motifs (E1, E2, E3, and E4), each about 112 residues long and with one or two potential  $\text{Ca}^{2+}$ -binding sites per repeat (Fig. [1](#)). A (partial) fifth domain of quite variable length has been designated the *extracellular anchor* (EA). Some at least of the first four repeats participate in an antiparallel overlap with similar domains in cadherin molecules emanating from other cells, thus facilitating homotypic assembly under  $\text{Ca}^{2+}$  control. Other proteins are also believed to be involved in this interaction. Details of the conformation of a cadherin repeat have recently been obtained using nuclear magnetic resonance (NMR) and X-ray diffraction methods ([2](#), [3](#)). These repeats are also important in determining the specificity of adhesion between the different cadherins. The cytoplasmic domain, which is able to bind a variety of proteins, contains a substructure that shows considerable variations in size between different members of the cadherin family. Consider, for example, human desmoglein and human *N*-cadherin ([4](#)). Both display a proline-rich region (albeit of different lengths) termed the *intracellular anchor* (IA), followed by a highly charged region (C1). This represents the full extent of the human *N*-cadherin sequence, in contrast to human desmoglein, which contains a further 320 residues. The C1 domain, being highly conserved between *N*-cadherin and desmoglein, is likely to be involved in interactions with cytoskeletal [microfilaments](#) and other key proteins in the cytoplasm, such as plakoglobin. The remaining part of human desmoglein can be subdivided into a 59-residue proline-rich domain (C2), a region with novel 29-residue repeats (C3) that contain a high density of putative **phosphorylation** sites, and a *C*-terminal domain (C4) that is quite divergent in sequence between human and cow, for example. The first half of domain C4 (known as C4a) is very [nonpolar](#) and contains two partial

repeats of length 28 residues that are rich in [glycine](#) residues. They are not related to the 29-residue repeat seen in domain C3. Subdomain C4b is clearly basic in character, in contrast to all of the other cytoplasmic domains, which are acidic. A model of the spatial arrangement of desmoglein in the desmosome has been proposed (4).

**Figure 1.** Comparison of the domain structures of human desmoglein 1 (HDSG1), desmoglein 3 (HDSG3), desmocollin 3 (HDSC3), and *N*-cadherin (NCAD). This illustrates the structural similarities between the classical cadherins and those in the desmosomes. The transmembrane  $\alpha$ -helical segment (M) separates the extracellular domains E1 to E4 and the extracellular anchor EA from the intracellular anchor IA and the intracellular domains C1 to C4. Putative calcium-binding domains in the E domains are shown as shaded rectangles. The two desmocollin 3 forms result from alternative splicing. The black rectangle represents 11 amino acids at the C-terminal end that are unique to the shorter structure. (From Ref. 5, with permission.)



## Bibliography

1. J. A. MARRS and W. J. NELSON (1996) Cadherin cell adhesion molecules in differentiation and embryogenesis. *Int. Rev. Cytol.* **165**, 159–205.
2. M. OVERDUIN, T. S. HARVEY, S. BAGBY, K. I. TONG, P. YAU, M. TAKEICHI, and M. IKURA (1995) Solution structure of the epithelial cadherin domain responsible for selective cell adhesion. *Science* **267**, 386–389.
3. L. SHAPIRO, A. M. FANNON, P. D. KWONG, A. THOMPSON, M. S. LEHMANN, G. GRUBEL, J.-F. LEGRAND, J. ALS-NIELSEN, D. R. COLMAN, and W. A. HENDRICKSON, (1995) Structural basis of cell-cell adhesion by cadherins. *Nature (London)* **374**, 327–337.
4. L. A. NILLES, D. A. D. PARRY, E. E. POWERS, B. D. ANGST, R. M. WAGNER, and K. J. GREEN (1991) Structural analysis and expression of human desmoglein: a cadherin-like component of the desmosome. *J. Cell Sci.* **99**, 809–821.
5. A. P. KOWALCZYK, T. S. STAPPENBECK, D. A. D. PARRY, H. L. PALKA, M. L. A. VIRATA, E. A. BORNSLAEGER, L. A. NILLES, and K. J. GREEN (1994) Structure and function of desmosomal transmembrane core and plaque molecules. *Biophys. Chem.* **50**, 97–112.

## Suggestions for Further Reading

6. K. J. Green and J. C. R. Jones (1996) Desmosomes and hemidesmosomes: structure and function of molecular components. *FASEB J.* **10**, 871–881.
7. M. Takeichi (1995) Morphogenetic roles of classic cadherins. *Curr. Opin. Cell Biol.* **7**, 619–627.
8. O. Huber, C. Bierkamp, and A. Kemier (1996) Cadherins and catenins in development. *Curr. Opin. Cell Biol.* **8**, 685–691.

## Caenorhabditis

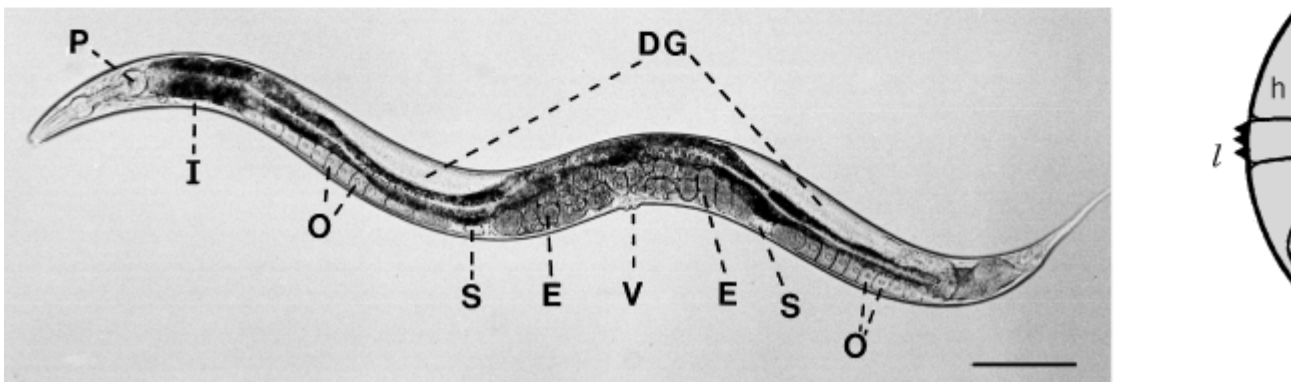
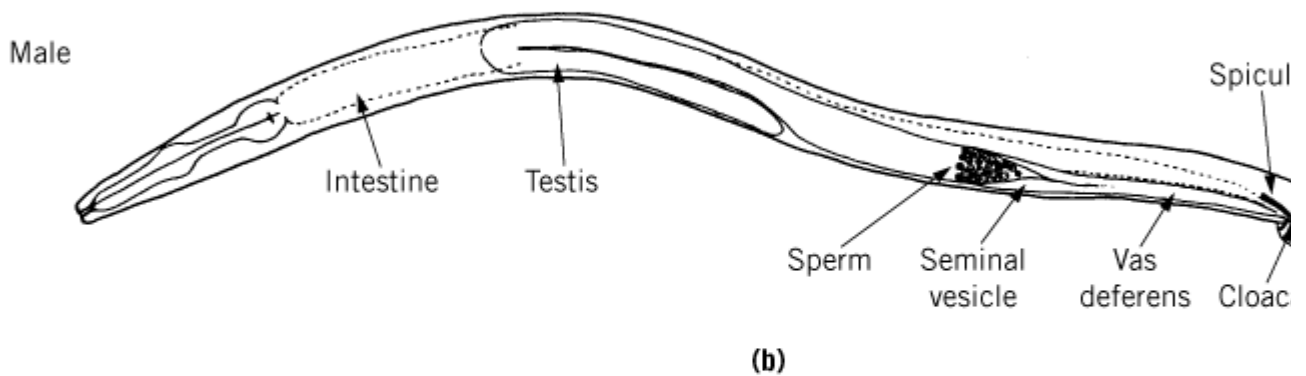
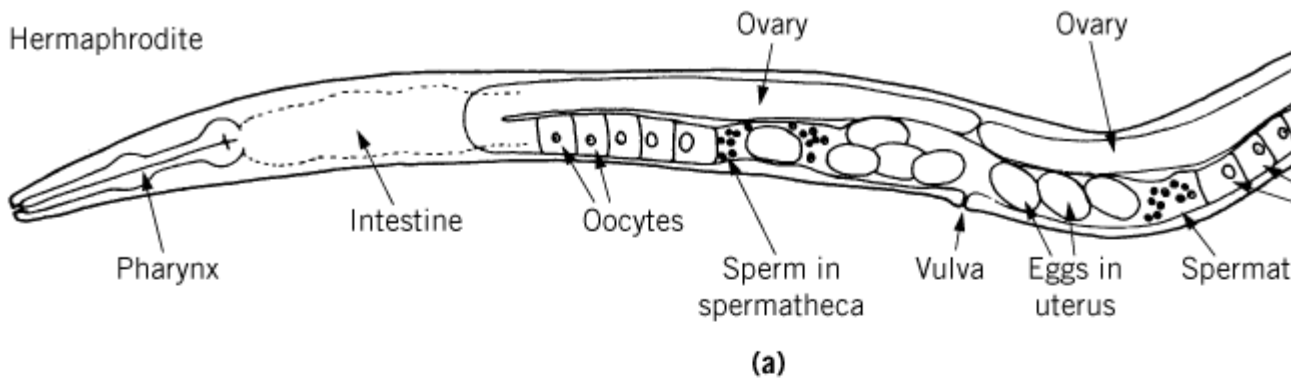
Caenorhabditis is a genus of primarily free-living soil [nematodes](#). The free-living species *Caenorhabditis elegans* has become well known after being chosen by Sydney Brenner in the 1960s for a concerted genetic, ultrastructural, and molecular analysis of animal development ([1](#)). A few other Caenorhabditis species have been characterized in less detail, largely for comparison, but the findings summarized in this article are for *C. elegans* unless otherwise indicated.

*C. elegans* was chosen for a variety of features, including its short life cycle and **hermaphroditic** mode of [development](#) (facilitating genetic analysis), its anatomical and cellular simplicity (<1000 somatic cells in the adult hermaphrodite), and its transparency and small size (facilitating light [microscopy](#) of developing animals and ultrastructural analysis, respectively). Its genome size has turned out to be relatively small as well (100 Mb), making molecular biology convenient. An entire field of *C. elegans* biology has grown out of Brenner's original investigations. As a result of this effort and the relative simplicity of “the worm,” as it is sometimes referred to, *C. elegans* is now the most completely understood metazoan in terms of genetics, ultrastructure, development, and physiology, and it has taken its place with the fruit fly *Drosophila* and the laboratory [mouse](#) as one of the most important model organisms for research biology ([2](#), [3](#)). Its entire cell lineage from zygote to adult has been determined; its ultrastructure including the complete connectivity of its nervous system has been described in detail; an extensive genetic map has been assembled; and DNA sequencing of its entire [genome](#) will be completed in 1998. It is beautifully suited for the genetic approach: identifying **genes** required for any biological process of interest by mutations that affect it, molecularly characterizing these genes after isolation by positional [cloning](#), and eventually studying the biochemical functions of the corresponding gene products. As in all organisms, the genetic control of development and physiology is still only beginning to be understood in *C. elegans*. The rapidly accumulating detailed knowledge of the worm's biology, as well as the recent realization that many developmental and physiological mechanisms are remarkably conserved throughout the animal kingdom, make *C. elegans* a powerful experimental system for investigating a variety of the most important questions in current biological research.

### 1. General Description

*C. elegans* inhabits soil in most parts of the world and feeds on soil **bacteria**. The adults are about 1 mm in length, and the two sexes, hermaphrodites and males, are morphologically distinguishable (Fig. [1](#)). The hermaphrodites produce both eggs and sperm and can self-fertilize. Males can mate with hermaphrodites to give cross progeny; hermaphrodites cannot fertilize each other. A hermaphrodite can produce about 300 self progeny, in addition to cross progeny if mated. [Embryos](#) begin development in the hermaphrodite uterus and are laid at the gastrulation stage. They hatch as first-stage juveniles (L1 larvae), which as in all nematodes grow through three subsequent larval stages punctuated by molts, before a final molt to the sexually mature adult. The adult reproductive period lasts for about 4 days, after which the animals live for an additional 2 weeks.

**Figure 1.** Anatomy of hermaphrodite and male *C. elegans* adults. Major structural features mentioned in the text are labeled. (a) Hermaphrodite. (b) Male. (c) Hermaphrodite, bright-field photomicrograph. Scale bar 0.1 mm. P, pharynx; I, intestine; C, vulva; DG, distal gonad. (d) Hermaphrodite, cross section through the anterior, viewed toward the anterior. *d*, *v*, *l*, *r*, dorsal muscle; *h*, hypodermis; *I*, intestine; *g*, gonad; *pc*, pseudocoelom; *nc*, nerve cord. Note the treads (alae) on the left and right sides.



*C. elegans* is one of the simplest metazoans and has fixed numbers of **somatic cells**: 959 in the adult hermaphrodite and 1031 in the adult male. In the laboratory it can be cultured conveniently on agar plates spread with *Escherichia coli* bacteria, or in liquid culture when larger quantities are desired for biochemical work. Individual animals can be observed and transferred from plate to plate with a platinum wire “worm pick” under a dissecting microscope, which is sufficient for scoring many

mutant phenotypes. Taking advantage of the animal's transparency throughout the life cycle, it is possible to follow individual cells during development at higher magnification using a compound microscope, preferably equipped with Nomarski differential interference-contrast optics (see [Microscopy](#)).

## 2. Genetics and the Genome

*C. elegans* is **diploid**, with five [autosomes](#) (I–V) and a sex [chromosome](#) (X). Hermaphrodites are XX and males XO; sex is determined by the X chromosome to autosome ratio. Males arise spontaneously in self-fertilizing hermaphrodite populations at a frequency of about 0.2%, as the result of X chromosome **nondisjunction**.

The [haploid](#) genome size of *C. elegans* is 100 Mb, about 10 times that of **yeast**, one-half that of *Drosophila*, and one-thirtieth that of mammals. Based on frequency of predicted coding units in the genomic DNA sequence, which is now nearly complete, the estimated total number of functional **genes** is about 13,000. Over 1000 genetic loci have been identified by [mutation](#) following chemical [mutagenesis](#) and mapped to the six linkage groups. Positional cloning of mutationally identified genes is facilitated by an extensive physical map of overlapping clones, which is anchored to the genetic map at many locations. The DNA sequence of the entire genome will soon be completely known.

## 3. Anatomy

*C. elegans* has the typical nematode body plan (see [Nematodes](#)): an outer tube of hypodermis with attached musculature and neurons, surrounding an interior space (the *pseudocoelom*), which contains the gut and, in adults, the gonad (see Fig. 1). The hypodermis secretes a three-layered, **collagenous** cuticle that is shed and replaced at each larval molt (see text below). Four bands of body-wall muscles run the length of the animal. The neuromuscular system bends the body only in the dorsal–ventral direction, so that the animal always lies on one side or the other on a solid surface. Specialized hypodermal cells on each side secrete “treads,” known as *alae*, on which the worm moves. The digestive system consists of a muscular pharynx in the head, which crushes bacteria and pumps them into the intestine, which empties at the anus near the tail. An excretory system controls the hydrostatic pressure on which the worm depends to maintain its body shape.

The entire nervous system of the hermaphrodite consists of only 302 neurons and 56 supporting and glial cells. The major ganglion is the nerve ring at the base of the pharynx, which sends processes down ventral and dorsal nerve cords to a secondary ganglion in the tail. Motor neurons from the cords enervate the body-wall muscles. Sensory neurons extend into the nerve ring from chemoreceptors and touch receptors in the head. Hermaphrodites have specialized neurons that control egg-laying (see text below); males lack these but have additional neurons that provide input from sensory structures in the tail and control male mating behavior.

The hermaphrodite reproductive system is comprised of a bilobed gonad, with one lobe extending anteriorly and one posteriorly from the uterus near the middle of the animal. The gonad consists of a somatic sheath surrounding the germ cells and has several specialized regions. The distal arm of each lobe contains **mitotically** dividing germ-cell **nuclei** in a common **cytoplasm**. The nuclei enter **meiosis** as they move away from the distal tip. The earliest nuclei to mature differentiate into [sperm](#) during the fourth larval stage; at the molt to adulthood the germ line switches sex, and subsequent meiotic nuclei are recruited to form **oocytes** as they round the bend into the proximal arm, which contains the oviduct. At the proximal end of the oviduct is the spermatheca, containing stored sperm, through which the oocytes pass and become fertilized on their way to the uterus. A muscular vulva connects the uterus to the outside and serves as the egg-laying apparatus. In the male, the gonad is single lobed and produces only [sperm](#), which are stored in the vas deferens and released through the cloaca in the tail. The fan-shaped male tail is specialized for mating, which is accomplished by deposition of sperm through the hermaphrodite vulva into the uterus, where they move to the

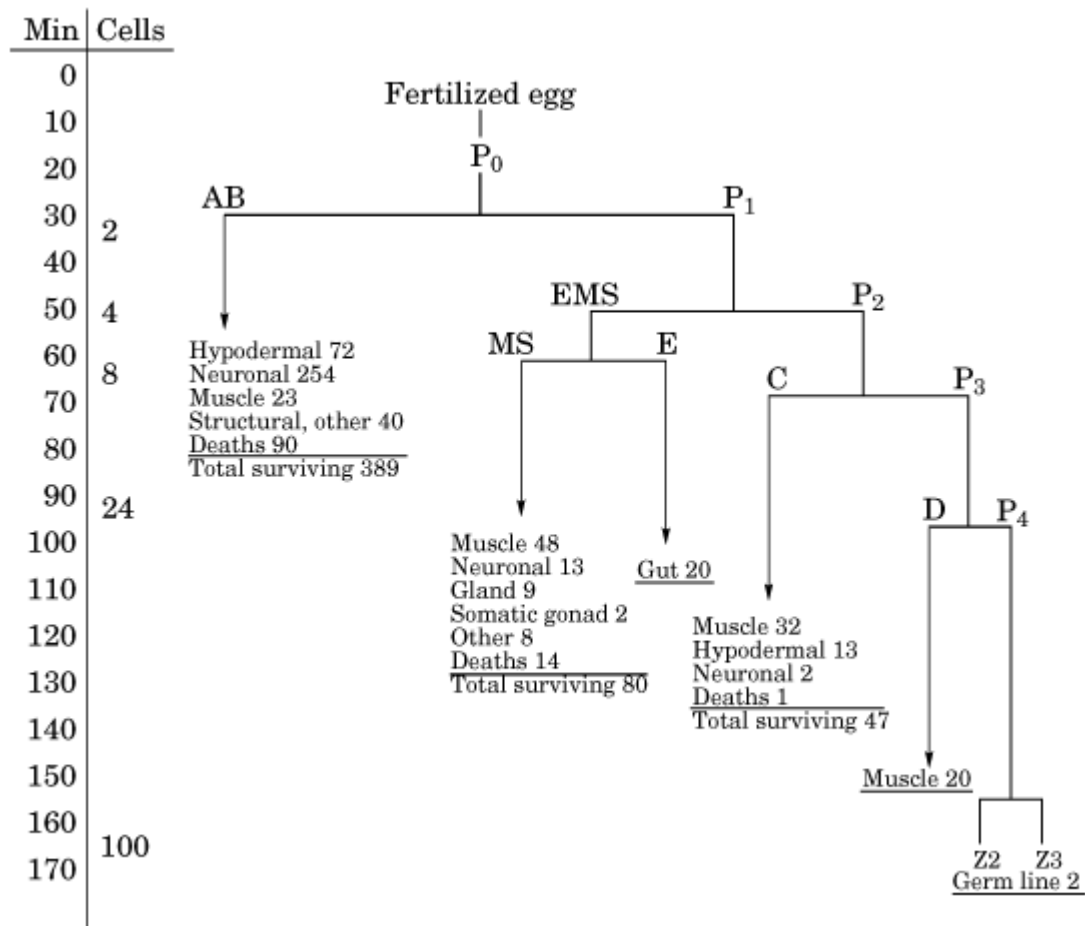
spermatheca and compete for oocytes with the resident hermaphrodite sperm.

#### 4. Fertilization and Embryonic Development

Following fertilization, a tough chitinous egg shell forms around the zygote, and embryonic cleavage begins. The spheroidal embryos, about 70  $\mu\text{m}$  in length, are viable if dissected out of the hermaphrodite at this time, but normally they remain in the uterus for about 2 h until the onset of gastrulation, when they are laid through the vulva. Gastrulation is simple; only 53 cells move from the exterior to the interior, but the result is a triploblastic embryo with outer ectodermal precursor cells that will form the hypodermis and nervous system, inner endodermal precursors that will form the gut, and in between a layer of mesodermal precursors that will give rise to muscles and the somatic gonad. By 6 h after fertilization, midway through embryogenesis, the embryo consists of about 600 cells. At this point, cell proliferation essentially ceases and a process of [morphogenesis](#) begins which literally squeezes the spheroidal embryo into the shape of a worm as organogenesis proceeds internally and cuticle is formed externally. During this time about 40 cells undergo [programmed cell death](#) and are engulfed by neighboring cells. At the end of embryogenesis the worm, about 3.5 times the length of the original embryo, digests the shell from the inside and hatches out of the egg.

The process of **embryogenesis** is essentially invariant at the cellular level. It proceeds by a stereotyped series of **cell divisions** (Fig. 2), whose timing and relative spatial orientations are the same in every embryo. This feature and the transparency of the embryo made it possible for John Sulston and his colleagues to trace out the entire embryonic cell lineage from fertilization to hatching, so that the ancestry of each of the 558 cells in the L1 is known (4). Most cells are born close to their final locations; only about 12 cells undergo long-range migrations during embryogenesis.

**Figure 2.** Early divisions in the *C. elegans* embryonic cell lineage, showing cells and tissues derived from each of the major branches. The vertical axis shows time after fertilization and total number of cells in the embryo. Horizontal lines indicate times of cell divisions.  $P_0$  through  $P_4$  are germ line cells; the cells named AB, MS, E, C, and D are somatic founder cells for the various branches of the lineage. The number of cells of each different type produced in each branch is indicated. Note that many cells are programmed to die during embryonic development, especially in the AB and MS branches.



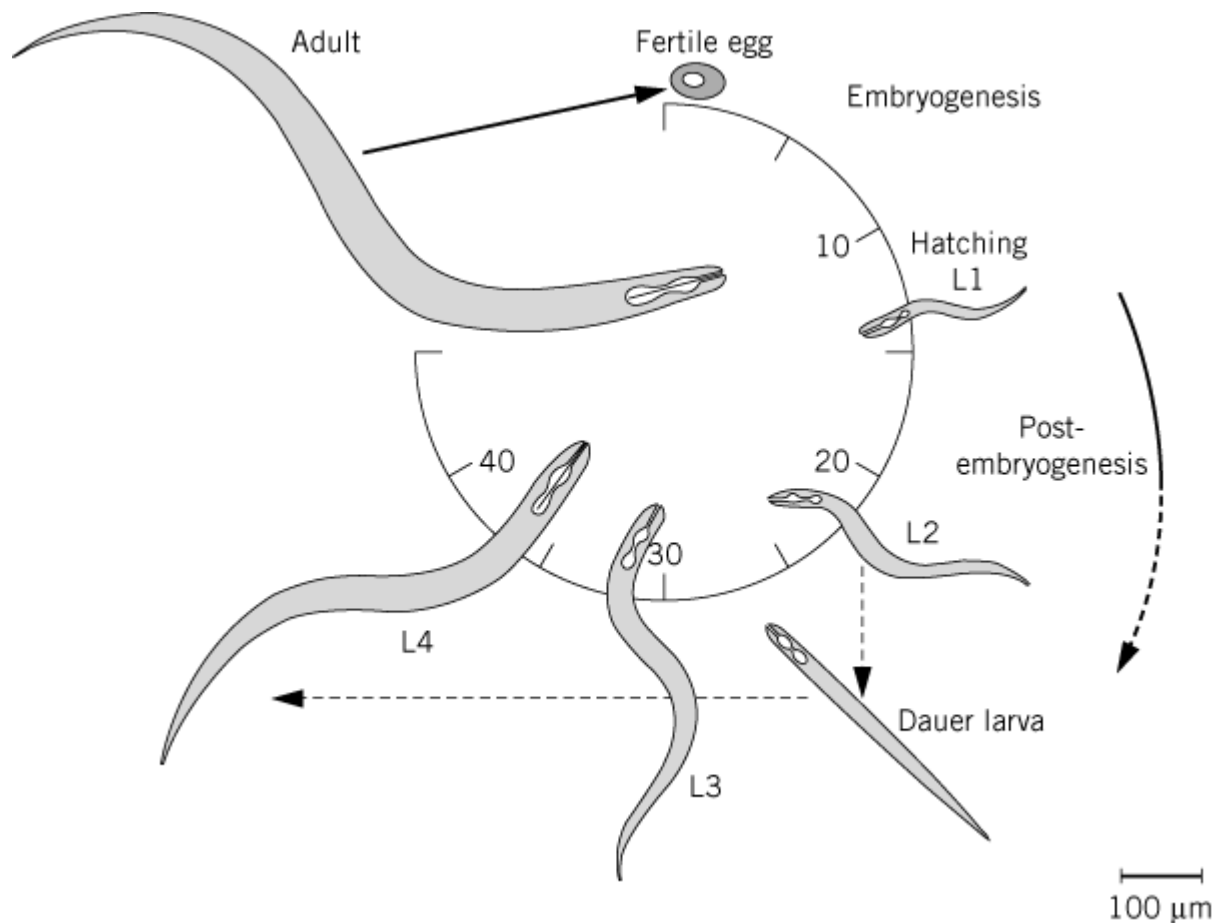
Genetic and molecular analysis of the cleavage stage has revealed that cell fates in the early embryo (Fig. 2) are determined by a combination of maternally derived factors, which segregate asymmetrically to specific cells during cleavage, and cell-signaling mechanisms by which certain cells induce new cell fates in their neighbors. This patterning process appears to proceed by evolutionarily conserved mechanisms that are common to all embryos; for example, the signaling pathways so far known to be employed several that are also important for insect and vertebrate [development](#), such as the *wingless* (Wnt), and *Notch* pathways (5).

## 5. Larval Development

### 5.1. Larval Stages

After hatching as an L1, the worm grows through three more larval stages, L2–L4 (Fig. 3), separated by molts at which a new cuticle is secreted by the underlying hypodermis and the old one shed. The newly hatched L1, 250  $\mu\text{m}$  long with 558 cells in the hermaphrodite and 560 in the male, includes functional digestive, neuromuscular, and neurosensory systems and a gonad primordium consisting of two somatic and two germ-line cells. Most of the cells in the L1 divide no further during larval development. However, 55 are *blast* cells that undergo additional divisions during larval development to produce, primarily, the adult reproductive structures and the neurons that control them: the hermaphrodite and male gonads, the vulva and egg-laying muscles in the hermaphrodite, and the sex muscles and specialized tail mating structures in the male. Additional motor neurons are also produced in the nerve cords for finer control of the body-wall muscles. As in the embryo, the invariance of these processes, the transparency of all stages, and the small number of cells involved made it possible to map the larval cell lineages in both sexes from L1 to adulthood (6). Consequently, the entire cell ancestry of *C. elegans* is known, from fertilized egg to adult.

**Figure 3.** The *C. elegans* life cycle. Clock circle indicates hours of development after fertilization at 25°C. Eggs containing developing embryos are laid at about 2 h. The L1 larvae, hatching at about 14 h, undergo four more molts at the times indicated as they grow to adulthood. See text for further explanation.



## 5.2. Cell Interactions during Larval Development

Elaboration of larval structures involves cell interactions, again occurring by evolutionarily conserved signaling pathways such as those involving homologues of Notch, Wnt, [epidermal growth factor](#) (EGF), and [transforming growth factor](#) (TGF)- $\beta$ -superfamily ligands as the signaling molecules. Application of the genetic approach to the process of vulval development, for example, has revealed details of a *ras*-based pathway by which the developing gonad signals underlying hypodermal cells to differentiate using an EGF-like ligand. The molecular components are very similar to those of *ras* pathways found in mammalian cells, in which conversion of *ras* to an **oncogene** causes a variety of cancers, and the *C. elegans* pathway has become an important model system for cancer research.

Cell interactions during larval development also include guidance cues by which several cells, including the growing axonal processes of motoneurons and the distal tip cells of the enlarging hermaphrodite gonad, find their way along the basement membrane of the pseudocoelom en route to their final destinations. Again, application of the genetic approach in *C. elegans* has shown that this process is accomplished by a conserved mechanism: The so-called *netrin* ligands and their **receptors** that guide these cells in the worm are remarkably similar to the molecules that guide axonal growth from the neural floor plate to the spinal cord in developing avian and mammalian embryos, making *C. elegans* a potentially useful model system for study of nerve regeneration.



The cell divisions and molts in larval development occur in a precisely timed sequence. Again, application of the genetic approach has identified many of the genes in the timing mechanism and is providing new information on how developmental clocks function.

### 5.3. The Dauer Larva

Like most free-living nematodes, *C. elegans* can molt to an alternative form of the L3, called a *dauer larva* (German for “enduring larva”) when conditions are unsuitable for reproduction. Dauer larvae, often simply called dauers, do not feed, are relatively resistant to drying, and can live for up to a year if desiccation is prevented. If conditions improve, they can molt to the L4 stage and resume the normal developmental pathway, with no change in subsequent lifespan. Dauer development is triggered by lack of food and overcrowding, signaled by a pheromone of unknown nature that stimulates chemosensory structures in the head. This process is of interest for at least two reasons. First, the signaling mechanism involves homologues of ligands, receptors, and downstream components of the evolutionarily conserved TGF- $\beta$  pathway, which is of widespread importance but not fully understood in mammalian development. Second, entry into the pathway of dauer development appears to turn off genes that normally limit the *C. elegans* lifespan to less than 3 weeks, allowing the dauer to live much longer. This finding, and the experimentally convenient short normal lifespan, make *C. elegans* an exciting model system for studying genetic control of aging, which is under active investigation (7).

## 6. Behavior and the Nervous System

*C. elegans* can move forward or backward and change direction, by coordinated flexing of its body-wall muscles. It is touch-sensitive, moving forward in response to a touch on the tail and backward in response to a touch on the head. Chemosensors in the head, connected by sensory neurons to the ring ganglion, can detect a variety of ions as well as volatile odorants and elicit either an attractive or repulsive response (see [Chemotaxis](#)). For example, *C. elegans* is attracted to  $\text{Na}^+$ ,  $\text{K}^+$ ,  $\text{Cl}^-$ , and several alcohols and ketones; it is repelled by  $\text{Cu}^{2+}$ , acid pH, D-tryptophan, and benzaldehyde. Males are attracted by a pheromone of unknown nature that is produced by hermaphrodites. Although some of its relatives have photoreceptors, *C. elegans* does not appear to respond to light or use it as a sensory cue. It does, however, sense and respond to temperature and will move to a preferred point in a temperature gradient.

Is an animal with only 302 neurons capable of learning? Simple conditioning experiments suggest that *C. elegans* not only can habituate to several stimuli, but also is capable of associative learning, that is, learning to use a normally neutral stimulus to predict the arrival of a second more significant stimulus. For example, if presence of food is paired with one of two equally attractive ions and absence of food with the other during conditioning, the conditioned animals will preferentially move toward a source of the paired ion when tested subsequently in the absence of food, and this preference lasts up to 7 h after training. In addition, *C. elegans* not only can sense a temperature gradient, but also can “remember” the temperature at which it has previously fed and move to the same temperature when placed in a new gradient without food. The genetic approach should allow the genes involved in learning and memory in *C. elegans* to be identified and the mechanisms of the proteins they encode to be elucidated.

*C. elegans* does not show obvious circadian rhythms, but it exhibits much higher frequency (ultradian) rhythms, such as a regular defecation cycle that is repeated about once every 45 s when the animal is feeding, regardless of the temperature. As in other organisms, the mechanisms of the temperature-compensated molecular clocks that control such cycles (like the human heartbeat) are just beginning to be understood. Application of the genetic approach to cyclical behaviors in *C. elegans* is identifying the genes and proteins that control ultradian rhythms.

## 7. Current Investigation and Online Information about *C. elegans*

Knowledge about *C. elegans* genetics, development, and behavior is accumulating rapidly in several areas of currently exciting biological research, some of which are mentioned above. These include mechanisms of pattern formation in development, origins of left–right asymmetry in animal body plans, cell fate determination, programmed cell death, cell migration and guidance, developmental timing mechanisms, physiological sensory mechanisms, organization of animal genomes, and animal [evolution](#).

As a supplement to the references listed below, much current knowledge of *C. elegans* is accessible electronically on the World Wide Web. A convenient access site can be found at <http://eatworms.swmed.edu>.

### Bibliography

1. S. Brenner (1974) *Genetics* **77**, 71–94.
2. W. B. Wood et al. (eds.) (1988) *The Nematode Caenorhabditis elegans*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
3. D. L. Riddle, T. Blumenthal, B. J. Meyer, and J. R. Priess (eds.) (1997) *C. elegans II*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
4. J. Sulston, E. Schierenberg, J. White, and J. Thomson (1983) *Dev. Biol.* **100**, 64–119.
5. R. Schnabel and J. R. Priess (1997) In *C. elegans II* (D. L. Riddle, T. Blumenthal, B. J. Meyer, and J. R. Priess, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 361–382.
6. J. Sulston and H. Horvitz (1977) *Dev. Biol.* **56**, 110–156.
7. C. Kenyon (1997) In *C. elegans II* (D. L. Riddle, T. Blumenthal, B. J. Meyer, and J. R. Priess, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 791–813.

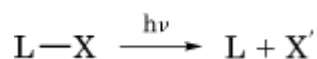
### Suggestions for Further Reading

8. S. Brenner (1974) The genetics of *Caenorhabditis elegans*. *Genetics* **77**, 71–94. (The classic paper in which Brenner first presented his rationale, descriptions of morphological and behavioral mutants, and preliminary genetic analysis of *C. elegans*.)
9. D. L. Riddle, T. Blumenthal, B. J. Meyer, and J. R. Priess (eds.) (1997) *C. elegans II*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY. (The sequel to the previous reference book. A good up-to-date summary of current knowledge about *C. elegans*.)
10. J. Sulston and H. Horvitz (1977) Post-embryonic cell lineages of the nematode *Caenorhabditis elegans*. *Dev. Biol.* **56**, 110–156. (The first extensive cell lineaging of *C. elegans*, describing all the cell lineages during larval development.)
11. J. Sulston, E. Schierenberg, J. White, and J. Thomson (1983) The embryonic cell lineage of the nematode *Caenorhabditis elegans*. *Dev. Biol.* **100**, 64–119. (The monumental paper in which Sulston and his co-workers published the complete embryonic cell lineage of *C. elegans*.)
12. W. B. Wood (1988) In *The Nematode Caenorhabditis elegans* (W. B. Wood and the community of *C. elegans* researchers, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–16. (The introductory chapter to a useful general reference on *C. elegans*.)

### Caged ATP

The synthesis of caged ATP (caged ADP and caged phosphate) and many of the properties that make these reagents useful in a variety of biological applications were first published in 1978 (1). The driving force for the development of such a reagent was the desire to be able to introduce ATP rapidly (and synchronously) at sites of biological interest at a desired time. A pulse of ATP would be released by **light activation** and would then initiate the processes being studied (Fig. 1). In contrast to [affinity labeling](#), here the ATP analogue should not bind to its site prior to activation and ATP is released freely in solution. The strategy employed was an approach that had previously been used in synthetic organic chemistry (ie, photodeprotection). In the chemical applications, a desired functionality in a molecule was modified and protected from a variety of reagents and conditions during a multistep synthesis, to be deprotected at a later step by a light-activated process. In biological applications, the caged substrates are “protected” from their receptor, [enzyme](#), or binding site by the chromophoric moiety, and following the pulse of light they are released to their biological target.

**Figure 1.** Basis of photorelease of a caged substrate. The biological ligand (L) or substrate is rendered inactive by the attachment of a photocleavable chromophore (x). Following photoactivation the ligand is released accompanied by the modified chromophore, or photofragment.



The properties of molecules appropriate for the study of biological systems are subject to more constraints than those used in synthetic organic chemistry. Much of the subsequent work on various caged ATP molecules and other caged compounds has been aimed at refining or improving their basic properties to expand the areas of application or overcome a particular limitation. The design requirements for a caged ATP are illustrative of the needs of any caged biological ligand or substrate (2).

The basic requirements for a biologically useful caged ATP (and for most other caged compounds) are as follows: (i) The activating light must be at wavelengths long enough to avoid damage to biological material (usually greater than 300 nm). (ii) The photorelease process must be as efficient as possible; that is, the quantum yield (ratio of molecules of product obtained to molecules of caged compound excited) must be as high as possible. (iii) The photorelease process should be as rapid as possible, relative to the rate of the process of interest; this has been achieved in the tens of microseconds to millisecond range, at normal temperatures and pH values. (iv) The photoreleased fragment, which usually is the modified chromophore, should not be harmful to biological materials (see text below). (v) The **absorbance** of the caged compound should be reasonably high at the exciting wavelength or wavelength range. Initially it was thought that too high an absorbance might lead to nonuniform release of substrate across the depth of an illuminated sample. Interestingly, for the recent application of multiphoton excitation, very high absorbances are an advantage. (vi) Prior to photolysis, the caged substrate should neither bind to nor interact with the biological material of interest. In many respects, the first caged ATP reported and subsequently used (but the second synthesized) fulfilled many of these requirements. In a variety of systems, one or other of these properties has been less than ideal, and this has led to various variations on the original structure in attempts to provide an improved caged ATP.

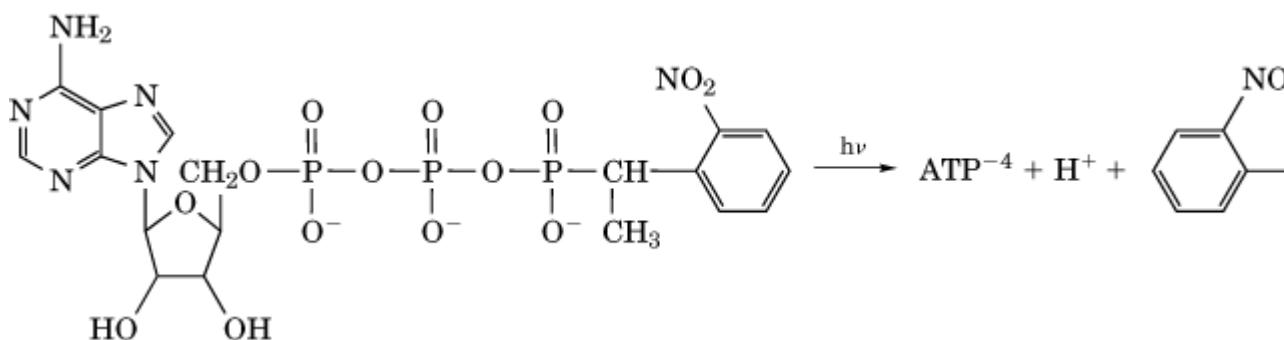
## 1. Synthesis and Photochemical Properties of Caged ATP

The initial synthesis of caged ATP was based on a strategy of first preparing caged phosphate and then coupling this to ADP. This had the advantage that the photochemical lability, quantum yield, and so on, could be first characterized, using caged phosphate and the readily assayed free substrate, inorganic phosphate (1). This route suffered from the following disadvantages: (i) It was less direct

than was desirable to achieve a gamma  $^{32}\text{P}$ -labeled caged ATP for a variety of **phosphorylation** studies, and (ii) it was not generally applicable to a wide range of phosphate-bearing biological ligands that might be desirable to cage. It would be better to have a more universal caging moiety that could be readily attached to the nucleotide or other organophosphates. These disadvantages were overcome by the development of a synthesis by Trentham and co-workers, who were able to attach the protecting 2-nitrobenzyl-based chromophore directly to the terminal phosphate of ATP, using a diazoprecursor (3). This was based on an approach that had been used by others to prepare a photosensitive **cyclic AMP** phosphotriester analogue (4). This method was used to prepare a wide array of caged phosphorylated biological ligands, including GTP, ADP, inositol trisphosphate, and so on (3).

The quantum yield for the first successful caged ATP was 0.54 and has not been bettered by subsequent analogues. The first caged ATP that was synthesized, the primary 2-nitrobenzyl analogue, yielded on photolysis some ATP, but only very low levels, even though photolysis was complete. It was hypothesized that this was due to a reaction between ATP and the released photofragment (2-nitrosobenzaldehyde), so that a chemically modified ATP resulted (1). The 2-nitrophenylethyl analogue (from the secondary benzyl alcohol) then became the molecule of choice, because photolysis yielded the less reactive 2-nitrosoacetophenone and free ATP in high yield (Fig. 2). This is the caged ATP that has been used most frequently in biological studies. The absorbance properties of caged ATP are simply the sum of the spectra of the 2-nitrobenzyl moiety and of ATP. This results in a long tail of absorbance that extends to above 350 nm. This has enabled illumination above 300 nm to be employed (using lasers or flash-lamps) so that photodamage due to absorption by most biological samples is avoided.

**Figure 2.** Photolysis of caged ATP. Photolysis of caged ATP produces a proton, free ATP, and the photofragment, 2-nitrosoacetophenone.



## 2. Mechanism and Kinetics of Photorelease

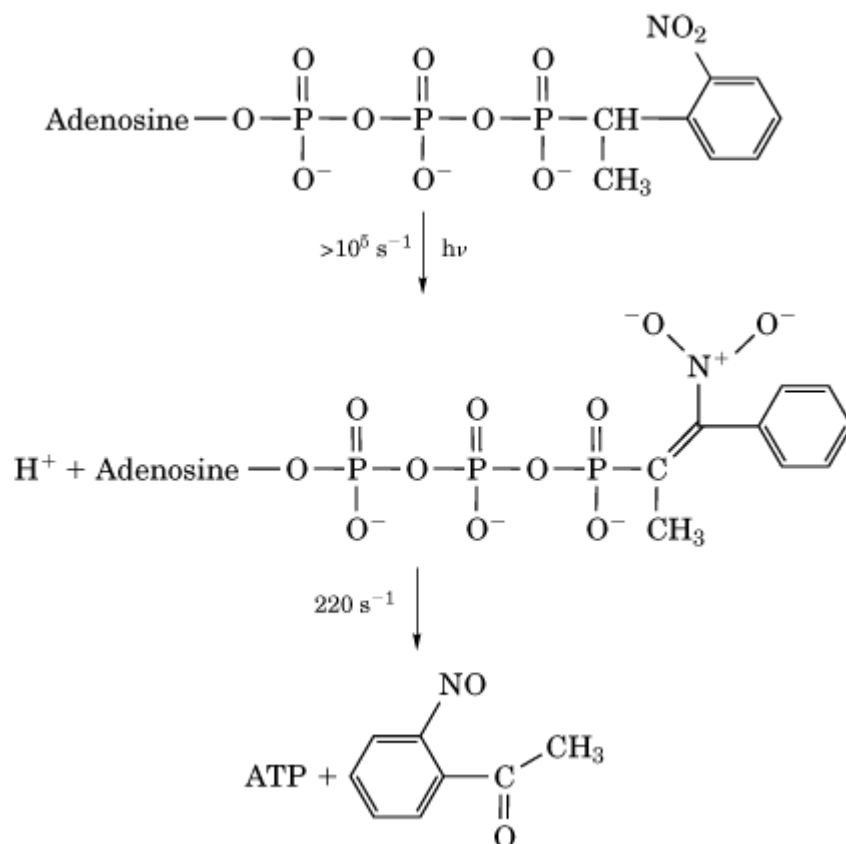
Although the earliest studies showed that ATP release could be achieved in a subsecond time frame, ideas about the kinetics of the release process and the likely mechanism were initiated by the report of McCray et al. in 1980 (5). From this work it was apparent that the rate of release of ATP from caged ATP was pH-dependent and in the millisecond time range.

Importantly, it was also pointed out in this work that there was a chromophoric intermediate on the reaction pathway, identified tentatively as an aci-nitro intermediate. The intermediacy of this type of chromophore has proven useful in later mechanistic studies of the photorelease of caged substrates, because it is central to almost all photorelease processes in which the protecting moiety is a 2-nitrobenzyl residue. Another tool that has been useful in mechanistic studies, and profitably used by Trentham and colleagues, is the reaction of the photoreleased nitrosoacetone with **thiol groups**. This

reaction had been initially employed by Kaplan et al. (1) to protect biological samples against any damaging effects of the photofragment, because it was known that nitrosoketones react readily with thiols. It was shown that the inhibition of the Na,K-ATPase enzyme due to the actions of the photofragment could be prevented by the simultaneous presence of **dithiothreitol** (DTT) or similar reagents. Furthermore, the protecting thiol had to be present in amounts at least stoichiometric with the released photofragment. In many subsequent applications, it has become routine to include reduced **glutathione** or DTT in the reaction media to mop up the released photofragment. Trentham and colleagues (6) used the rapid rate of this reaction to examine the release rates of the photofragment. This proved valuable because with caged ATP, as with many photoreleased biological substrates, it is often difficult to identify a biological process that is sufficiently fast to be used unambiguously as a bioassay to determine the photorelease rate of substrate. It should be emphasized, however, that for all new caged substrates it is essential to measure the release rate of the photoproduct of interest and not merely the rate of decay of an intermediate on the release pathway. There have been several reports of biphasic decay of putative aci-nitro intermediates when their breakdown was followed using transient ultraviolet spectroscopy (7, 8).

The most recent estimates (obtained from time-resolved infrared spectroscopy) of the formation of free ATP following the photolysis of caged ATP are about  $220 \text{ s}^{-1}$  at pH 7 and  $22^\circ\text{C}$ , a rate that is the same as the aci-nitro anion intermediate decay rate (9). An outline of the steps involved in the photolysis pathway are shown in Figure 3.

**Figure 3.** Scheme of breakdown of caged-ATP following excitation. The release rate of ATP is controlled by dark reactions that follow photoexcitation and relaxation to ground states. A detailed study of the products and intermediates has been made using time-resolved infrared spectroscopy and isotopomers of caged ATP (6).



### 3. Other Caged ATP Structures

The most familiar caged ATP is P<sup>3</sup>-(1-(2-nitrophenyl)ethyl) adenosine 5'-triphosphate. Variations on the basic theme have been made in the hope of increasing the long wavelength absorbance of the molecule or of speeding the rate of release. Unfortunately, no molecule yet reported has achieved this aim while simultaneously maintaining the relatively high quantum yield. Most attempts have involved variations in the chromophore, while maintaining the same nitrobenzyl photochemical cleavage mechanism (10). Recently potentially useful alternative photochemistries have begun to be exploited, but as yet no generally useful molecule has emerged (11-13).

#### 4. Biological Applications of Caged ATP

The two areas where caged ATP has been used most extensively are in studying the mechanism of **active transport** by **ion pumps** and in the mechanism of muscle contraction and its regulation. In single turnover studies of the Na,K-ATPase (see text below), Forbush showed that caged ATP did in fact bind with low affinity to the enzyme prior to photolysis (14). This leads to extra complications in kinetic modeling of events following photolysis. Subsequently, similar effects have been noted in muscle fibers (15). Such observations suggest that a caged ATP with a bulkier chromophore might reduce the pre-photolysis binding.

#### 5. Studies on Ion Pumps

The ion pumps are P-type **ATPases** that couple the transport of ions against their electrochemical potential gradient to the hydrolysis of ATP. This family of membrane proteins includes the Na pump or Na,K-ATPase, the Ca pump [from intracellular membranes (SERCA) or plasma membranes (PMCA)], **proton pumps** from a variety of organisms such as yeast, bacteria, and so on, and a range of heavy metal pumps that transport Cu, Cd, Ni, and so on, across cellular membranes throughout the animal and plant kingdom (16). The studies on ion pumps that have used caged ATP illustrate well several of the different advantages of this experimental strategy. These proteins carry out a series of biochemical transformations that are thought to accompany the transport of cations across the membrane. These transformations arise from the intermediate involvement of a phosphoenzyme, which is formed by transferring the terminal phosphate of ATP to the protein at an early step in the cycle, and from the protein to water in a later step. These **kinase** and **phosphatase** activities are linked to ion binding, translocation, and ion release steps (17). In order to probe the cation activation of some of these processes, it was necessary in the case of the Na pump to be able to initiate the pump cycle in a sealed system. This was required because the activating and inhibitory cations have very asymmetric effects, depending on whether they act at the extracellular or intracellular surface. The strategy employed was to trap caged ATP within resealed human erythrocyte ghosts. The extracellular medium could then be altered at will, and the Na pump process could be initiated by illumination of the ghost suspension and release of intracellular ATP from caged ATP. This enabled characterization of the side effects of activating and inhibiting cations on the transphosphorylation reactions under genuine initial rate conditions (18, 19). The essential properties of the photorelease strategy here were the stability of caged ATP during experimental procedures, until ultraviolet irradiation, when ATP could be released in good yield synchronously inside the cells in suspension. The possibility of releasing ATP from caged ATP in a rapid synchronous fashion within an ordered structure was also exploited in structural studies of the Ca pump in oriented multilayers where diffraction before and after the release of ATP from caged ATP showed movements of the protein mass relative to the membrane structure during the reaction cycle (20). Such studies would not have been possible without the photorelease approach.

Ion pumps transfer charge during their reaction cycle across cell membranes; and electrophysiological measurements of such movements, along with their analysis, have been a central area of study. In recent years, Bamberg and co-workers (21) and Apell and co-workers (22) have initiated the use of caged ATP in a novel experimental system to analyze the mechanistic consequences of these phenomena in a variety of ion pumps. These workers have used black lipid

membranes and have either attached or fused to them biomembrane fragments or vesicles containing ion pumps. Following the rapid, synchronous release of ATP from caged ATP in the medium bathing the membranes, the current-carrying and charge translocating steps can be analyzed with a conventional electrophysiological measuring system (23, 24). Such studies have probed the basis of the electrogenic nature of the Na pump and of the electrically neutral basis of the gastric proton pump, for example (25). Recent studies have dissected out the electrical properties of the bacterial active K transporter, the Kdp-ATPase of *Escherichia coli* (26).

### 5.1. Regulation and Control of Muscle Contraction

An adequate understanding of the mechanism of muscle contraction can only be achieved from experiments in a fairly intact and complex tissue. This is because the structure of the muscle fiber and its organization is an inherent feature of its function. The sliding filament model for the molecular basis of muscle contraction is the prevailing paradigm in this field, and its central feature is this relationship (27, 28). Thus it is necessary to be able to perform high-resolution kinetic studies in an ordered array of macromolecules. Recent efforts have been aimed at understanding the precise steps in the ATP hydrolysis cycle carried out by myosin that lead to the generation of force during contraction. Caged ATP has found application in detailed studies on the kinetics of ATP hydrolysis, structural changes in the fiber organization associated with contraction and relaxation (at the macroscopic and molecular levels), spectroscopic measurements on muscle fibers, and the process of relaxation associated with Ca pumping by the sarcoplasmic reticulum. Studies in all these areas have received considerable impetus by the introduction of the photorelease strategies. It has been shown that it is possible to photorelease ATP (or other substrates) in chemically skinned muscle fibers in the millisecond time range and to cause the synchronized initiation of biochemical processes that produce contraction (or relaxation) in the bundle of fibers (29, 30). Prior to this technology, such experiments were limited by the diffusional delays that were inherent in mechanically adding substrate and allowing it to diffuse into the fiber to its site of action. Initial studies that employed this approach to analyze the kinetics of the mechanical processes have now been greatly extended; by using covalently attached chromophoric reporters or paramagnetic reporters, the kinetics of conformational changes in the myosin protein can be monitored (31, 32). As well as such structural studies, it has also been possible to carry out time-resolved diffraction studies of fibers in a synchrotron beam prior to and following the release of caged ATP within the fiber (33, 34). These applications use the highly ordered skeletal muscle system to understand the basis of muscle contraction and relaxation. In smooth muscle, the contractile system has evolved to a highly complex level of cellular regulation, and many [signal transducing](#) systems and effectors play a role in regulating smooth muscle activity (35). In this area, the photorelease of caged ATP (and a variety of [second messengers](#)) continues to be a highly fruitful strategy. Recent studies by the Somlyos and their group provide a wide range of applications of this technology to the smooth muscle system (36, 37).

### 5.2. Other Systems

Since 1980 to the present time (early 1999) there have been more than 200 articles published using caged ATP, and many more using related substrates. These have extended from [difference Fourier](#) structural determinations of proteins with or without bound substrates (38, 39), to the cellular effects of occupancy of receptors (40). Structural biology of macromolecules, especially proteins, continues to provide enormous insights into molecular biological processes, and it is in this area that several interesting applications have great potential. The possibility of obtaining not merely a static picture of a protein at atomic resolution, but also one that contains dynamic information, is enticing. The accessibility of high-intensity radiation sources and a revival of interest in the [Laue Diffraction](#) method have made fast [X-ray crystallography](#) studies an attainable goal. The idea that such measurements could be made with protein crystals before and after laser-induced release of caged substrates within a crystal has received recent attention; and although many of the technical difficulties (uniformity of release, fate of photofragment, local heating effects) have not been overcome, this approach offers great promise for the future (41, 42). The number of physiological systems that have begun to be probed using caged ATP continues to increase; recent additions to this group include L-type Ca channel modulation (43), single [kinesin](#) molecules in optical traps (44),

ATP-sensitive K channels in cholinergic interneurons (45), sea-urchin sperm **flagella** motility (46), and limulus photoreceptors (47).

## 6. Summary

The introduction of the photorelease strategy that followed the description of caged ATP has now been extended to a wide array of biological substrates, including nucleoside phosphates, cyclic AMP, cyclic GMP, protons,  $\text{Ca}^{2+}$ ,  $\text{Mg}^{2+}$ , peptides, inositol phosphates, sugars, amino acids, neurotransmitters, toxins, and so on. In almost all cases, the photochemical process has the same basis as the original caged ATP strategy; along with a wide array of biological substrates, the approach can now be employed to probe systems that range in complexity from single protein crystals to brain slices.

## Bibliography

1. J. H. Kaplan, B. Forbush III, and J. F. Hoffman (1978) *Biochemistry* **17**, 1929–1935.
2. G. P. Hess (1999) (this encyclopedia).
3. J. W. Walker, G. P. Reid, J. A. McGray, and D. R. Trentham (1988) *J. Am. Chem. Soc.* **110**, 7170–7177.
4. J. Engels and E.-J. Schlaeger (1977) *J. Med. Chem.* **20**, 907–911.
5. J. A. McCray, L. Herbette, T. Kihura, and D. R. Trentham (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7237–7241.
6. A. Barth et al. (1997) *J. Am. Chem. Soc.* **119**, 4149–4159.
7. J. E. T. Corrie (1993) *J. Chem. Soc., Perkins Trans* **1993**, 2161–2166.
8. G. C. R. Ellis-Davies, J. H. Kaplan, and R. J. Barsotti (1996) *Biophys. J.* **70**, 1006–1016.
9. A. Barth et al. (1995) *J. Am. Chem. Soc.* **117**, 10311–10316.
10. J. F. Wootton and D. R. Trentham (1989) *NATO ASI Series C* **272**, 277–296.
11. J. E. Baldwin et al. (1990) *Tetrahedron* **46**, 6879–6884.
12. R. S. Givens et al. (1993) *J. Am. Chem. Soc.* **115**, 6001–6012.
13. J. E. T. Corrie and D. R. Trentham (1992) *J. Chem. Soc., Perkins Trans I*, 2409–2417.
14. B. Forbush III (1984) *Proc. Natl. Acad. Sci. USA* **81**, 5310–5314.
15. J. Sleep, C. Hermann, T. Barman, and F. Travers (1994) *Biochemistry*, **33**, 6038–6042.
16. S. Lutsenko and J. H. Kaplan (1995) *Biochemistry* **34**, 15607–15613.
17. J. H. Kaplan (1985) *Annu. Rev. Physiol.* **47**, 534–544.
18. J. H. Kaplan and R. J. Hollis (1980) *Nature* **288**, 587–589.
19. J. H. Kaplan (1982) *J. Gen. Physiol.* **80**, 915–937.
20. D. Pascolini et al. (1988) *Biophys. J.* **54**, 679–688.
21. K. Fendler, E. Grell, M. Haubs, and E. Bamberg (1985) *EMBO J.* **4**, 3079–3085.
22. R. Borlinghaus, H.-J. Apell, and P. Lauger (1987) *J. Membr. Biol.* **97**, 161–178.
23. E. Bamberg, H.-J. Butt, A. Eisenrauch, and K. Fendler (1993) *Q. Rev. Biophys.*, **26**, 1–25.
24. I. Wuddel and H.-J. Apell (1995) *Biophys. J.* **69**, 909–921.
25. M. Stengelin, K. Fendler, and E. Bamberg (1993) *J. Membr. Biol.* **132**, 211–227.
26. K. Fendler, S. Drose, K. Altendorf, and E. Bamberg (1996) *Biochemistry*, **35**, 8009–8017.
27. A. F. Huxley and R. Niedergerke (1954) *Nature* **173**, 971–973.
28. H. E. Huxley and J. Hanson (1954) *Nature* **173**, 973–976.
29. Y. E. Goldman, M. G. Hibberd, J. A. McCray, and D. R. Trentham (1982) *Nature* **300**, 701–705.
30. E. Homsher and N. C. Millar (1990) *Annu. Rev. Physiol.* **52**, 875–896.



31. J. W. Tanner, D. D. Thomas, and Y. E. Goldman (1992) *J. Mol. Biol.* **223**, 185–203.
32. C. L. Berger, F. C. Svensson, and D. D. Thomas (1989) *Proc. Natl. Acad. Sci. USA* **86**, 8753–8757.
33. K. J. V. Poole, G. Rapp, Y. Maeda, and R. S. Goody (1988) *Adv. Exp. Med. Biol.* **226**, 391–404.
34. K. Horiuti et al. (1994) *J. Biochem.* **115**, 953–957.
35. A. P. Somlyo and A. V. Somlyo (1994) *Nature* **372**, 231–236.
36. A. P. Somlyo and A. V. Somlyo (1990) *Annu. Rev. Physiol.* **52**, 857–874.
37. B. Zimmerman et al. (1995) *J. Biol. Chem.* **270**, 23966–23974.
38. A. J. Scheidig et al. (1995) *J. Mol. Biol.* **253**, 132–150.
39. C. Raimbault et al. (1997) *Eur. J. Biochem.* **250**, 773–782.
40. G. D. Housley, N. P. Raybould, and P. R. Thorne (1998) *Hear. Res.* **119**, 1–13.
41. J. Hadju and L. N. Johnson (1990) *Biochemistry* **29**, 1669–1678.
42. I. Schlichting et al. (1990) *Nature* **345**, 309–314.
43. B. O'Rourke, P. H. Backx, and E. Marban (1992) *Science* **257**, 245–248.
44. M. Higuchi, E. Muto, Y. Inoue, and T. Yanagida (1997) *Proc. Natl. Acad. Sci. USA*, **94**, 4395–4400.
45. K. Lee, A. K. Dixon, T. C. Freeman, and P. J. Richardson (1998) *J. Physiol.* **510**, 441–453.
46. T. Tani and S. Kamimura (1998) *J. Exp. Biol.* **201**, 1493–1503.
47. M. N. Faddis and J. E. Brown (1992) *J. Gen. Physiol.* **100**, 547–570.

## Calcium Signaling

Calcium ion is an important [second messenger](#) involved in cell signaling and [signal transduction](#). Most calcium in cells is sequestered in intracellular vesicles, the [endoplasmic reticulum](#) (ER) or [mitochondria](#), where it is stored for release when needed. Small, localized increases in calcium result from its regulated release from the ER, which is produced by inositol trisphosphate, IP<sub>3</sub>. IP<sub>3</sub> is produced by the **hormone**-dependent hydrolysis of phosphoinositides, to produce **inositol phosphates**, including IP<sub>3</sub>. The IP<sub>3</sub> binds to a specific ER protein, a specialized **calcium channel** with four identical subunits, each with a single membrane-spanning segment and a single IP<sub>3</sub>-binding site ([1](#)). The binding of IP<sub>3</sub> to this receptor results in a stereospecific release of calcium from the ER.

IP<sub>3</sub>-sensitive calcium channels are expressed ubiquitously in tissues. However, the process of regulated calcium release is most tightly regulated in endocrine and neuronal cells, and it is often involved in coupling stimulus to secretion, which is usually a calcium-dependent process. Calcium release from the ER is oscillatory in nature. Increases in intracellular calcium are modulated not by increases in the amount of the ion released, but instead by the frequency of the oscillations. This may explain why the release of some hormones is pulsatile.

One of the critical targets in calcium signaling is CAM kinase, a multigene family of protein [kinases](#) that are sensitive to regulation by calcium/[calmodulin](#). There are at least four members of this protein kinase family, with differences in distribution and substrate specificity. CAM kinases appear to be

especially critical to modulation of signaling in the central nervous system (2).

## Bibliography

1. M. J. Berridge (1993) *Nature* **361**, 315–325.
2. P. I. Hanson and H. Shulman (1992) *Annu. Rev. Biochem.* **61**, 229–601.

## Calcium-Binding Proteins

Calcium-binding proteins (CaBPs) are a key component linking the inorganic and organic systems that produce biological activities. Calcium is one of the most abundant inorganic elements in nature and is ubiquitous throughout biology, playing roles at the organismal, cellular, and molecular levels. Calcium is an essential component of shells and bones, where it is essentially deposited in crystals to create macroscopic support structures. At the cellular level,  $\text{Ca}^{2+}$  is one of the crucial currencies of most living organisms, acting as a [second messenger](#) in a wide range of key intracellular and extracellular systems. Calcium-binding proteins play important roles in mediating each of these effects. These proteins can be grouped into four primary categories: (1) intracellular proteins involved in  $\text{Ca}^{2+}$ -mediated [signal transduction](#), (2) [enzymes](#), (3) extracellular cell surface and [extracellular matrix proteins](#), and (4) proteins containing [g-carboxyglutamic acid](#) (Gla) residues. Although a variety of structural motifs bind  $\text{Ca}^{2+}$ , all involve oxygen atoms of the protein backbone or side chains, reflecting the intrinsic affinity of  $\text{Ca}^{2+}$  for oxygen atoms.

### 1. Intracellular Calcium-Signaling Proteins

#### 1.1. EF-Hand Calcium-Binding Proteins

This family is the most extensively studied class of CaBPs. These proteins are characterized by a highly conserved helix–loop–helix motif termed the [EF-Hand Motif](#), which consists of a 12-residue  $\text{Ca}^{2+}$ -binding loop flanked by two [alpha-helices](#). Almost all EF-hand CaBPs are composed of pairs of EF-hands. This pairing of  $\text{Ca}^{2+}$ -binding sites is presumed to stabilize the protein conformation, increase the  $\text{Ca}^{2+}$  affinity of each site over that of isolated sites, and provide a means for cooperativity in  $\text{Ca}^{2+}$  binding. The cooperativity between sites allows for an “all or nothing” response to  $\text{Ca}^{2+}$ -binding, which is crucial for the function of these CaBPs as intracellular  $\text{Ca}^{2+}$  sensors.

[Calmodulin](#) and troponin C are the best known members of this family of CaBPs. They each have two largely independent domains connected by a flexible linker. Each domain contains two EF-hands, so these proteins each bind four  $\text{Ca}^{2+}$  ions. In the resting cell, they exist in an inactive state, with either  $\text{Mg}^{2+}$  or no ion bound. When the intracellular  $\text{Ca}^{2+}$  concentration rises in response to a signal, the proteins bind  $\text{Ca}^{2+}$ . This induces a dramatic conformational change, exposing a large **hydrophobic** surface within each domain that can interact with target proteins. There are many other CaBPs thought to function in a similar manner including caltractin and the calmodulin-like domain of plant  $\text{Ca}^{2+}$ -dependent protein kinase.

Less is known about the S100 proteins, another large and important subfamily of EF-hand CaBPs. The proteins in this subfamily are composed of two EF-hands each, the first of which is a variant version with a 14-residue binding loop, termed a pseudo-EF-hand. The ligands in the pseudo EF-

hand are mainly carbonyl oxygens of the peptide backbone, whereas the canonical 12-residue loop coordinates  $\text{Ca}^{2+}$  primarily with oxygen atoms of side chains. Many S100 proteins are found as hetero- or homodimers, and they exhibit tissue-specific expression patterns. While these proteins are thought to be involved in signal transduction, the molecular mechanism of this function is not known, nor have any target proteins been positively identified. However, the expression of S100 proteins is deregulated in some diseases, including cancer, rheumatoid arthritis, and Down's syndrome, and [antibodies](#) against these proteins are commonly used as markers for screening for these diseases.

Most EF-hand CaBPs are directly involved in intracellular calcium signal transduction. However, at least three fulfill other cellular requirements. [Parvalbumin](#) and calbindin  $\text{D}_{9k}$  (an S100 protein) are thought to play roles in  $\text{Ca}^{2+}$  buffering and in  $\text{Ca}^{2+}$  uptake and transport, respectively. These and other members of the EF-hand protein family are active in various aspects of intracellular  $\text{Ca}^{2+}$  homeostasis. The diversity of the roles of EF-hand CaBPs is further illustrated by the recent structure of BM-40, a glycoprotein found in the extracellular matrix ([1](#)), which was shown to contain two EF-hands. Prior to this discovery, the EF-hand motif was thought to be unique to intracellular CaBPs. This diversity in function of the EF-hand CaBPs provides evidence of very extensive evolutionary optimization of the fit between the EF-hand CaBP fold and the calcium ion.

## 1.2. Annexins

The [annexin](#) family is a second important class of intracellular CaBPs. The exact function of the proteins in this family is unknown, but they do exhibit an intriguing  $\text{Ca}^{2+}$ -dependent high-affinity binding to phospholipids (see [Membranes](#)). Annexins bind a large number of  $\text{Ca}^{2+}$  ions with high cooperativity. They assist transport of  $\text{Ca}^{2+}$  ions across membranes *in vitro*, although the physiological relevance has not yet been established. The mechanism of membrane binding is not known. The current model involves the  $\text{Ca}^{2+}$  ions acting as “glue” by simultaneously interacting with the protein and the membrane ([2](#)).

## 2. $\text{Ca}^{2+}$ -Binding Enzymes

Two well-known enzymes exhibit  $\text{Ca}^{2+}$ -dependent translocation to the membrane fraction, presumably through a mechanism similar to that used by the annexins. Some **isoforms** of protein kinase C bind to phospholipids with high affinity in the presence of  $\text{Ca}^{2+}$  ([2](#)). The intracellular group IV [phospholipase](#)  $\text{A}_2$  also requires  $\text{Ca}^{2+}$  for membrane association ([3](#)).

Many other enzymes require  $\text{Ca}^{2+}$  ions for stability. This includes many **serine proteases**: members of both the **trypsin** and the [subtilisin](#) families have been found to bind  $\text{Ca}^{2+}$ . These proteases each bind one to three  $\text{Ca}^{2+}$  ions, which are required for structural stability but do not directly participate in catalysis. In most cases, the  $\text{Ca}^{2+}$  ions are coordinated by ligands dispersed throughout the structure. Trypsin is an exception, however, and binds  $\text{Ca}^{2+}$  using a single 12-residue surface loop. The use of  $\text{Ca}^{2+}$  ions to stabilize the protein fold is not unique to the serine proteases. Another example of an enzyme that uses  $\text{Ca}^{2+}$  in this manner is [thermolysin](#), a  $\text{Zn}^{2+}$ -dependent protease.

In other enzymes, the  $\text{Ca}^{2+}$  ions are directly involved in catalysis. This includes the secreted forms of phospholipase  $\text{A}_2$ , in which the required  $\text{Ca}^{2+}$  ion is thought to be involved in [transition state](#) stabilization ([4](#)). **Staphylococcal nuclease**, a secreted protein that catalyzes DNA and RNA hydrolysis, also relies on  $\text{Ca}^{2+}$  for catalysis. The structure of this protein shows  $\text{Ca}^{2+}$  bound in the active site, where it is thought to be used to polarize the phosphate at the scissile phosphoester bond ([5](#), [6](#)).

### 3. Extracellular Cell Surface and Extracellular Matrix Ca<sup>2+</sup>-binding Proteins

#### 3.1. Cadherins

These are tissue-specific [cell adhesion molecules](#) that are strongly Ca<sup>2+</sup>-dependent (see [Cadherins](#)). At the levels of Ca<sup>2+</sup> found in the extracellular milieu, cadherins bind several Ca<sup>2+</sup> ions. A major conformational change is seen when Ca<sup>2+</sup> is removed from these proteins. The resulting change to the apo state conformation is thought to prevent cadherins from interacting with each other.

#### 3.2. C-Type Lectins

The C-type [lectins](#) are a Ca<sup>2+</sup>-dependent class of lectins, which mediate many cell surface carbohydrate recognition events. Selectins, concanavalin A, and mammalian mannose-binding protein are examples of this type of CaBP. The structure of mannose-binding protein shows that Ca<sup>2+</sup> is directly involved in binding the carbohydrate to this protein. In concanavalin A, on the other hand, the Ca<sup>2+</sup>-binding site is  $1.0\text{--}1.4 \times 10^{-9}$  m from the carbohydrate binding site; Ca<sup>2+</sup> does not participate directly in binding the carbohydrate, but instead stabilizes the protein structure. It is not known whether selectins use Ca<sup>2+</sup> directly in carbohydrate binding or to stabilize the fold required for carbohydrate recognition.

#### 3.3. EGF Modules

A subset of [EGF motifs](#) present as **domains** in a number of proteins contain b-hydroxy-aspartic acid residues and have Ca<sup>2+</sup>-binding activity. This type of EGF motif is found in some proteins involved in the [blood clotting](#) cascade and also in the extracellular matrix protein fibrillin. Fibrillin has 54 EGF modules, 43 of which are thought to bind Ca<sup>2+</sup>. It is thought that Ca<sup>2+</sup> is important in maintaining the structure of the fibrillin monomers and also in allowing the monomers to aggregate into microfibrils. A sheath of these microfibrils covers and stabilizes the [elastin](#) fibers that give tissues such as skin, blood vessels, and lungs their needed elasticity.

### 4. Gla-Containing Proteins

The Gla-containing proteins contain [g-carboxyglutamic acid](#) residues. This is a glutamic acid residue that has been carboxylated in a vitamin K-dependent process. These proteins are found in the bones and teeth, in the kidney, and in the blood. Osteocalcin (bone Gla-protein) and matrix Gla-protein are found in bone. They are believed to be involved in the mineralization of this tissue (7). Similar proteins are thought to be involved in the mineralization of teeth. In the kidney, the Gla-containing protein nephrocalcin inhibits the nucleation, aggregation, and growth of calcium oxalate crystals. An abnormal form of this protein is found in the urine of people suffering from kidney stones. However, it is not entirely clear whether this defect is responsible for the growth of kidney stones (8).

By far the best studied of the Gla-containing proteins are those found in blood. The majority of Gla-containing proteins are [zymogen](#) forms of serine proteases involved in the [blood clotting](#) (coagulation) cascade. These proteins are: **prothrombin**, factor VII, factor IX, factor X, protein C, protein S, and protein Z. The protease domains of these proteins are very similar to the pancreatic digestive proteases **trypsin**, **chymotrypsin**, and [elastase](#). The Gla-containing coagulants require Ca<sup>2+</sup> binding in order to bind to phospholipids. This Ca<sup>2+</sup>-dependent phospholipid binding is responsible for the membrane association properties required for the function of the proteins. Ca<sup>2+</sup> binding also appears to be necessary for the formation of the native conformation of the Gla domain (9).

The Gla-containing coagulants have two classes of metal ion binding sites. There are approximately three higher affinity sites that are not metal ion-specific, and three to four lower affinity sites, which are specific for Ca<sup>2+</sup>. Only Sr<sup>2+</sup> can substitute for Ca<sup>2+</sup> in both types of binding site.

The high resolution three-dimensional structure of the Gla domain of prothrombin, with seven  $\text{Ca}^{2+}$  ions bound (10), reveals the Gla domain to involve nine to 10 turns of [alpha-helix](#), in three separate helices. Seven  $\text{Ca}^{2+}$  ions interact with 24 oxygen atoms from 16 of the 18 carboxylate groups of the nine Gla residues that are ordered in the structure. A 10th Gla residue is disordered and does not participate in  $\text{Ca}^{2+}$  binding. The coordination geometries of the  $\text{Ca}^{2+}$  ions do not correspond to any idealized polyhedra. Five of the  $\text{Ca}^{2+}$  ions are involved in a polymeric array with 18 of the liganding oxygen atoms. Four of these  $\text{Ca}^{2+}$  ions are completely buried in the protein. This complex structure is thought to nucleate the folding of the Gla domain, and is essentially electrically neutral. The complexity and irregularity of this structure explains the selectivity for  $\text{Ca}^{2+}$  ions.  $\text{Ca}^{2+}$  is able to adopt different and distorted coordination geometries.  $\text{Mg}^{2+}$ , on the other hand, is fairly rigid in its requirement for six ligands, and cannot accommodate the unusual network of ligands in the Gla domain. The remaining two metal ion sites in the Gla domain of prothrombin are solvent accessible, and carry a net charge of about +0.5 each. Because of the positive charge, these sites are thought to be involved in neutralizing the negatively-charged phospholipids, allowing the protein to associate with membranes.

### Bibliography

1. E. Hohenester, P. Maurer, C. Hohenadl, R. Timpl, J. N. Jansonius, and J. Engel (1996) *Nature Struct. Biol.* **3**, 67–73.
2. M. D. Bazzi and G. L. Nelsestuen (1993) *Cell. Signal.* **5**, 357–365.
3. J. Y. Channon and C. C. Leslie (1990) *J. Biol. Chem.* **265**, 5409–5413.
4. E. A. Dennis (1994) *J. Biol. Chem.* **269**, 13057–13060.
5. P. J. Loll and E. E. Lattmam (1989) *Prot. Struct. Funct. Genet.* **5**, 183–201.
6. F. A. Cotton, E. E. Hazen, and M. J. Legg (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2551–2555.
7. M. F. Young, J. M. Kerr, K. Ibaraki, A.-M. Heegaard, and P. G. Robey (1992) *Clin. Orthoped. Rel. Res.* **281**, 275–294.
8. F. L. Coe, Y. Nakagawa, J. Asplin, and J. H. Parks (1994) *Miner. Electrolyte Metab.* **20**, 378–384.
9. J. W. Suttie (1993) *FASEB J.* **7**, 445–452.
10. M. Soriano-Garcia, K. Padmanabhan, A. M. de Vos, and A. Tulinsky (1992) *Biochemistry* **31**, 2554–2566.

### Suggestions for Further Reading

11. M. Celio, ed. (1996) *Guidebook to the Calcium-Binding Proteins*, Oxford University Press, Oxford.
12. H. Kawasaki and R. H. Kretsinger (1995) *Protein Profile; Calcium-Binding Proteins I: EF-hands*, Vol. **2**, Academic Press, London (a good general source for information about EF-hand proteins).
13. S. Liemann and A. Lewit-Bentley (1995) Annexins: a novel family of calcium- and membrane-binding proteins in search of a function, *Structure* **3**, 233–237.
14. P. Maurer, E. Hohenester, and J. Engel (1996) Extracellular calcium-binding proteins, *Curr. Opin. Cell Biol.* **8**, 609–617.
15. C. A. McPhalen, N. C. J. Strynadka, and M. N. G. James (1991) Calcium-binding sites in proteins: a structural perspective, *Adv. Protein Chem.* **42**, 77–144.
16. L. J. Van Eldik, J. G. Zendegui, D. R. Marshak, and D. M. Watterson (1982) Calcium-binding proteins and the molecular basis of calcium action, *Intern. Rev. Cytol.* **77**, 1–61. (A review of Gla-containing proteins,  $\text{Ca}^{2+}$ -binding enzymes, and EF-hand proteins.)

## Calmodulin

Calmodulin is the quintessential member of the **EF-hand** family of [calcium-binding proteins](#) and functions as a key mediator in numerous [signal transduction](#) pathways. It is an acidic [protein](#) ([isoelectric point](#) of 4.2) of molecular weight 16.8 kDa that is found in most eukaryotic cells, from yeast to humans. It is composed of two largely independent globular **domains** connected by a flexible central [a-helix](#). The affinity of calmodulin for  $\text{Ca}^{2+}$  is fine-tuned to respond to intracellular calcium signals. Conformational changes within each of the domains induced by the binding of  $\text{Ca}^{2+}$  leads to the transduction of the  $\text{Ca}^{2+}$  signal.

### 1. Biological Function

Calmodulin is a signal transduction protein. It has been implicated in the control of a wide range of cellular functions, including cell proliferation, smooth muscle contraction, the regulation of **ion channels**, long-term potentiation and memory, and [exocytosis](#). Furthermore, it has recently been found inside the [nucleus](#), where it is thought to be involved in the regulation of [DNA replication](#), **gene expression**, and [DNA repair](#). In a resting cell with basal levels of  $\text{Ca}^{2+}$ , calmodulin exists in the inactive apo state. When a calcium signal is initiated and the intracellular levels of  $\text{Ca}^{2+}$  increase, calmodulin binds  $\text{Ca}^{2+}$  ions. This causes the protein to undergo a dramatic conformational change that exposes a large **hydrophobic** patch on the protein. This newly exposed hydrophobic surface then interacts with various cellular proteins, modulating their activity and thereby transducing the calcium signal.

Protein kinases and **phosphatases** are among the best studied of calmodulin's targets. These enzymes are activated by the release of an autoinhibitory **domain** that is bound to the [active site](#) in the resting state.  $\text{Ca}^{2+}$ -loaded calmodulin activates these proteins by binding to a site near to or overlapping with the autoinhibitory domain, causing the auto-inhibitory domain to dissociate from the active site. Myosin light chain kinase and calcium-calmodulin-dependent protein kinase II are two well-known examples of this class of calmodulin-regulated proteins. Calmodulin also regulates proteins involved in the generation of other second messengers, such as calmodulin-dependent cyclic nucleotide phosphodiesterase and nitric oxide synthase. The mechanism of regulation of these enzymes is thought to be very similar to that used to regulate the kinases and phosphatases. Calmodulin has also been shown to interact with proteins of the cytoskeleton and their regulatory proteins, such as [spectrin](#) and brush border myosin. However, the biological significance of these interactions is still unclear.

Recently, interactions between apocalmodulin and targets such as neuromodulin and unconventional myosins have been described. These interactions are thought to be mediated by an "IQ motif" (with the consensus sequence IQXXRGXXR, where X is any amino acid) in the target (see [Protein Motif](#)). This motif and its interaction with a  $\text{Ca}^{2+}$ -free EF-hand calcium-binding protein was first described in myosin (1, 2). The functional significance of the interaction of proteins containing the IQ motif with apocalmodulin is not clear. It has been suggested, however, that perhaps neuromodulin, which is associated with the plasma membrane, functions as a trap for calmodulin, sequestering it near the membrane in the absence of the activating calcium signal (3).

### 2. Ion-Binding Properties

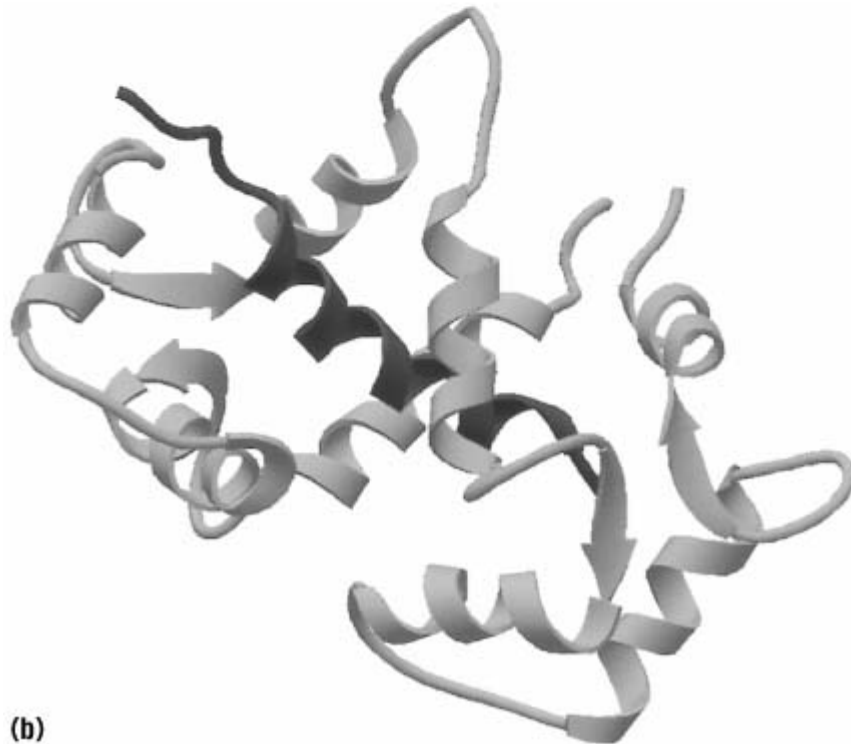
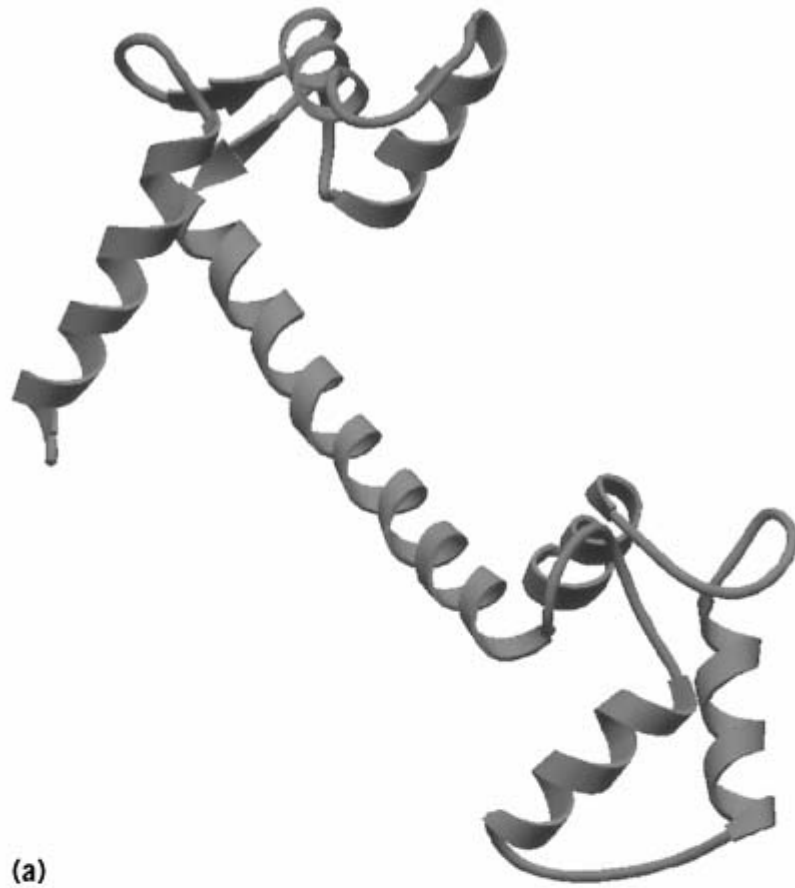
Calmodulin binds four calcium ions, with a **dissociation constant** of about  $10^{-6}$  M. Binding is

selective for  $\text{Ca}^{2+}$ ; calmodulin does not bind  $\text{Mg}^{2+}$  or monovalent ions with appreciable affinity. The binding of  $\text{Ca}^{2+}$  is also highly **cooperative**, which is very important for calmodulin's biological function as a calcium sensor, because it allows for a tightly controlled “all or nothing” response to the calcium signal: the four binding events all occur within a very narrow range of  $\text{Ca}^{2+}$  concentration. If these binding events were spread over a large range of  $\text{Ca}^{2+}$  concentrations, calmodulin would be at least partially activated over the entire range. The system would therefore lack the required sharp separation between the activate and inactive states.

### 3. Structure

Calmodulin is a largely helical protein composed of four EF-hand motifs organized into two independent globular domains, each of which contain two EF-hands. The two domains are connected by a long, flexible central  $\alpha$ -helix, giving the protein a dumbbell-like appearance (Figure 1a). This [quaternary structure](#) is relatively similar in the apo- and  $\text{Ca}^{2+}$ -loaded states of the protein. However, significant conformational changes occur within the individual domains on  $\text{Ca}^{2+}$  binding (Fig. 1b). In the apo state, each domain occupies the “closed” conformation (4-6). In this conformation, the four helices in the domain are nearly antiparallel, with interhelical angles near  $180^\circ$ . In  $\text{Ca}^{2+}$ -loaded calmodulin, each domain occupies an “open” conformation in which the four  $\alpha$ -helices are nearly perpendicular to each other (7, 8). Each domain exposes a hydrophobic surface of about  $1.25 \times 10^{-8} \text{ m}^2$  in this open conformation (5).

**Figure 1.** Ribbon representations of the three-dimensional conformations of  $(\text{Ca}^{2+})_4$ -calmodulin in the (a) absence and (b) presence of a target peptide. The peptide-free diagram was constructed using the coordinates of 1CLL (8). The two domains at the ends of the central  $\alpha$ -helix consist of two EF-hand motifs. The diagram of  $(\text{Ca}^{2+})_4$ -calmodulin and the peptide analog of the myosin light chain kinase was constructed using the coordinates of 2BBM (10). The bound peptide is depicted darker than calmodulin.



#### 4. Interactions with Target Peptides

The biophysical characterization of the interaction between calmodulin and its targets is often studied using small [peptides](#) (10–20 amino acid residues) derived from the calmodulin-binding



domain of target enzymes. Calmodulin binds to these target peptides extremely tightly, with dissociation constants ranging from  $10^{-7}$  to  $10^{-11}$  M. The peptides adopt an amphiphilic helical conformation, and tend to have a bulky hydrophobic residue at either end, often (but not always) spaced 12 residues apart. However, the sequences of the target peptides are not highly similar. Calmodulin is able to bind so tightly to such a wide array of targets because of its own plasticity. The flexible central  $\alpha$ -helix connecting its two domains can function as an “expansion joint,” allowing calmodulin to bind to peptides with different numbers of residues between the two bulky hydrophobic anchors. Furthermore, **van der Waals** interactions, which can be rather nonspecific, tend to dominate as the critical components stabilizing the interaction between calmodulin and the peptides. **Hydrogen bonds**, which are more structurally specific, are not as important in these interactions. The adaptability of the peptide binding surface of calmodulin may be further aided by its high proportion of **methionine** residues. Methionine is an unusually flexible and polarizable amino acid, which may allow calmodulin to mold its peptide-binding surface to meet the requirements of many different peptide sequences (9, 10).

High resolution three-dimensional structures of three complexes of calmodulin and peptides derived from target enzymes have been reported (11-13). These structures show that the relative disposition of the two domains of calmodulin is altered by the binding of a target peptide, but there is little change within the  $\text{Ca}^{2+}$ -activated domains themselves. The calmodulin-peptide complex forms a well-packed ellipsoid, which contrasts sharply with the dumbbell shape of calmodulin observed in the absence of target (Fig. 1). The two domains of calmodulin essentially wrap around the target peptide, forming a hydrophobic tunnel in which the peptide binds. As of yet, there are no structures of calmodulin bound to an entire target protein. It is thought, however, that the mode of binding to the target sequence on the intact enzyme will be very similar to that seen in the calmodulin-peptide complexes.

#### Bibliography

1. Xie, D. H. Harrison, I. Schlichting, R. M. Sweet, V. N. Kalabokis, and A. G. Szent-Gyorgyi (1994) *Nature* **368**, 306–312.
2. A. Houdusse and C. Cohen (1996) *Structure* **4**, 21–32.
3. Y. Liu and D. R. Storm (1990) *Trends Pharmacol. Sci.* **11**, 107–111.
4. H. Kuboniwa, N. Tjandra, S. Grzesiek, H. Ren, C. B. Klee, and A. Bax (1995) *Nature Struct. Biol.* **2**, 768–776.
5. M. Zhang, T. Tanaka, and M. Ikura (1995) *Nature Struct. Biol.* **2**, 758–767.
6. B. E. Finn, J. Evenäs, T. Drakenberg, J. P. Waltho, E. Thulin, and S. Forsén (1995) *Nature Struct. Biol.* **2**, 777–783.
7. Y. S. Babu, C. E. Bugg, and W. J. Cook (1988) *J. Mol. Biol.* **204**, 191–204.
8. R. Chattopadhyaya, W. Meador, A. Means, and F. Quijcho (1992) *J. Mol. Biol.* **228**, 1177–1192.
9. H. J. Vogel and M. Zhang (1995) *Mol. Cell. Biochem.* **149/150**, 3–15.
10. K. T. O'Neil and W. F. DeGrado (1990) *Trends Biochem. Sci.* **15**, 59–64.
11. M. Ikura, G. M. Clore, A. M. Gronenborn, G. Zhu, C. B. Klee, and A. Bax (1992) *Science* **256**, 632–638.
12. W. E. Meador, A. R. Means, and F. A. Quijcho (1992) *Science* **257**, 1251–1255.
13. W. E. Meador, A. R. Means, and F. A. Quijcho (1993) *Science* **262**, 1718–1721.

#### Suggestions for Further Reading

14. O. Bachs, N. Agell, and E. Carafoli (1994) Calmodulin and calmodulin-binding proteins in the nucleus, *Cell Calcium* **16**, 289–296 (a review of the current knowledge about calmodulin's nuclear functions).
15. A. Crivici and M. Ikura (1995) Molecular and structural basis of target recognition by

calmodulin, *Annu. Rev. Biophys. Biomol. Struct.* **24**, 85–116 (an exhaustive review of interactions between calmodulin and its targets).

16. R. D. Hinrichsen (1993) Calcium and calmodulin in the control of cellular behavior and motility, *Biochim. Biophys. Acta* **1155**, 277–293.
17. L. J. Van Eldik and D. M. Watterson, eds. (1997) *Calmodulin and Signal Transduction*, Academic Press, San Diego (a book containing articles about many aspects of calmodulin's function and structure).
18. C. B. Klee and E. Carafoli, eds. (1997) *Calcium as a Cellular Regulator*, Oxford University Press, Oxford.

## Calnexin/Calreticulin

Calnexin (also known as *p88*, *IP90*, and *CATCHER*) is a membrane-bound **molecular chaperone** present in the [endoplasmic reticulum](#) (ER) that binds selectively to many monoglucosylated **glycoproteins** that fold in this compartment during synthesis by ribosomes bound to the cytosolic face of the ER membrane. Calreticulin is a soluble homologue of calnexin, with the same glycan specificity, but occurs in the ER lumen. Together with other proteins, calnexin and calreticulin provide a means for improving the efficiency of protein folding and assembly in the ER compartment and form part of the quality-control system that ensures that only correctly folded and assembled proteins are transported from the ER compartment along the pathway followed by secretory proteins (1). Unlike some other molecular chaperones, calnexin and calreticulin are not **stress response** proteins, do not exhibit [ATPase](#) activity, and appear to recognize glucose residues in core glycans rather than **hydrophobic** surfaces in partially folded polypeptides (1). In some cell types, calreticulin also occurs in the [nucleus](#) and cytosol, where it binds to [steroid hormone](#) receptors and [integrin](#) molecules, respectively; these interactions may indicate roles for calreticulin in **gene expression** and cell signaling.

### 1. Discovery

An 88-kDa protein termed *p88* was found in transient association with class I **major histocompatibility** molecules synthesized by several murine lymphoma cell lines (2). Newly synthesized-class I heavy chains bound rapidly to *p88* before they associated with  $\text{b}_2$ -microglobulin chains. In mutant cells that lack microglobulin, the incompletely assembled class I molecules exhibited prolonged interaction with *p88* and were correspondingly impaired in their transport to the **Golgi** stacks (3). Moreover, microglobulin and peptides derived from the cytosol need to be assembled with the heavy chain before the latter can be released from *p88*. These findings led to the proposal that *p88* is a molecular chaperone that promotes assembly of class I molecules by retaining intermediates in the ER until complete complexes are formed.

Independent **pulse-chase** experiments on a human ER protein termed *IP90* showed that it transiently associates with many newly synthesized proteins, including the T-cell receptor (TCR), the B-cell antigen receptor, and class I molecules (4, 5). The binding of *IP90* to TCR complexes lacking a chains was prolonged, consistent with a retention role for this protein in the ER for incompletely assembled complexes (6). Subsequently, *p88* and *IP90* were found to be identical with the  $\text{Ca}^{2+}$ -binding phosphoprotein calnexin (4-7), a protein originally identified as one of several proteins labeled by **phosphorylation** of canine pancreatic **microsomes** with [ $\text{g-}^{32}\text{P}$ ]-GTP (8).

## 2. Structure

Calnexin is a type I nonglycosylated integral [membrane protein](#) of 573 amino acid residues, with a predicted molecular weight of 65,400 and with its substrate-binding domain in the lumen of the ER. It has a cytosolic tail of 89 residues that is phosphorylated and contains a C-terminal RKPRRE (—Arg—Lys—Pro—Arg—Arg—Glu) sequence that acts as an ER-retention signal (9). It also possesses four **domains** with high sequence similarity to calreticulin, the major [calcium-binding protein](#) of the ER lumen. Calreticulin contains 400 amino acid residues (46,000 molecular weight), which include both a KDEL ER-retrieval sequence and a strongly negatively charged region involved in Ca<sup>2+</sup> binding at its C terminus (10). Both calnexin and calreticulin appear to function as monomers.

## 3. Function

Calnexin and calreticulin are both [lectins](#) that specifically recognize glycoproteins that contain monoglucosylated core glycans; they do not recognize glycans containing two, three or no glucose residues (11, 12). When nascent chains of glycoproteins enter the ER lumen, a core glycan containing three terminal glucose residues is added en bloc to specific [asparagine](#) sidechains (see [N-Glycosylation](#)). This addition of glycan may have evolved initially to make folding intermediates more soluble under the conditions of high protein concentration that characterize the ER lumen. These glucose residues are then rapidly removed one at a time by glucosidases I and II in the ER lumen. Another luminal enzyme called UDP-Glc:glycoprotein glucosyltransferase then adds back one glucose residue. Thus glycans containing single terminal glucose residues can arise either as intermediates in the glucose-trimming process or after regeneration by the transferase. Such monoglucosylated glycans then bind to calnexin and calreticulin. In the case of pancreatic ribonuclease B, this binding occurs solely through the glycan and does not involve additional recognition of the protein moiety by the lectins (13, 14), but such additional recognition may occur in the case of other proteins (11, 15).

The function of the binding of calnexin and calreticulin to monoglucosylated glycoproteins appears to be to assist their correct folding by preventing aggregation of partially folded or misfolded chains and ensuring retention of the latter until they are degraded (see [DNA Degradation In Vivo](#)). Evidence for this conclusion comes from experiments in which proteins unable to fold correctly because of mutation remain bound to these chaperones for much longer than normal proteins and are eventually degraded (16-18). Additional possible roles include the suppression of formation of nonnative [disulfide bonds](#) (19) and the premature assembly of [oligomers](#) (18).

The distinction between correctly folded and partially folded chains is not made by these lectins but by the glucosyltransferase. This enzyme has the unusual property of distinguishing, by an unknown mechanism, between subtly different protein conformations; it adds glucose to the core glycan only if the protein attached to this glycan is not correctly folded (20, 21). Such reglucosylated glycoproteins will then bind back to calnexin and calreticulin. The removal and addition of the terminal glucose residue then continues until the protein has folded sufficiently to be no longer recognized by the transferase. The binding of calnexin to ribonuclease B *in vitro* is dynamic, indicating that the binding is readily reversible; glucosidase II cannot remove terminal glucose residues while calnexin is bound to the glycoprotein (14).

Calnexin and calreticulin associate with a wide range of proteins in the ER, including soluble secretory proteins, [extracellular matrix](#) molecules, **ion channels**, membrane **receptors**, and various glycoproteins of **viruses**. They also retain misfolded proteins that accumulate in some ER-storage diseases such as cystic fibrosis and  [\$\alpha\$ 1-antitrypsin](#) deficiency (22). Addition of glucosidase inhibitors to cells results in reduced rates of secretion of many glycoproteins, presumably due to the inhibition of their binding to the chaperones (23). All these observations support the view that the functions of calnexin and calreticulin are important for cell function. Nevertheless, viable mammalian cell lines

are known that lack either glucosidase I or II or calnexin, while mutants of *Saccharomyces cerevisiae* that lack genes for any one of these three proteins show no growth defect. However, mutants of *Saccharomyces pombe* that lack calnexin are not viable. Both species of *Saccharomyces* appear to lack calreticulin. These discrepancies may indicate some redundancy and overlap of function amongst the chaperones in the ER compartment, which include [BiP](#), hsp90, and [protein disulfide isomerase](#), as well as calnexin and calreticulin (1).

## Bibliography

1. A. Helenius, E. S. Trombetta, D. N. Hebert, and J. S. Simons (1997) *Trends Cell Biol.* **7**, 193–200.
2. E. Degen and D. B. Williams (1991) *J. Cell Biol.* **112**, 1099–1115.
3. E. Degen, M. F. Cohen-Doyle, and D. B. Williams (1992) *J. Exp. Med.* **175**, 1653–1661.
4. F. Hochenstenbach, V. David, S. Watkins, and M. B. Brenner (1992) *Proc. Natl. Acad. Sci. USA* **89**, 4734–4738.
5. K. Galvin et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8452–8456.
6. V. David, F. Hochstenbach, S. Rajagopalan, and M. B. Brenner (1993) *J. Biol. Chem.* **268**, 9585–9592.
7. N. Ahluwalia, J. J. M. Bergeron, I. Wada, E. Degen, and D. B. Williams (1992) *J. Biol. Chem.* **267**, 10914–10918.
8. I. Wada, D. Rindress, P. H. Cameron, W.-J. Ou, J. J. Doherty, D. Louvard, A. W. Bell, D. Dignard, D. Y. Thomas, and J. J. M. Bergeron (1991) *J. Biol. Chem.* **266**, 19599–19610.
9. M. Michalak, R. E. Milner, K. Burns, and M. Opas (1992) *Biochem. J.* **285**, 681–692.
10. B. Sonnichsen, J. Fullefrug, P. Nguyen, W. Diekmann, D. G. Robinson, and G. Mieskes (1994) *J. Cell Sci.* **107**, 2705–2717.
11. F. E. Ware, A. Vassilakos, P. A. Petersen, M. R. Jackson, M. A. Lehrmann, and D. B. Williams (1995) *J. Biol. Chem.* **270**, 4697–4704.
12. R. G. Spiro, Q. Zhu, V. Bhoyroo, and H.-D. Soling (1996) *J. Biol. Chem.* **271**, 11588–11594.
13. A. R. Rodan, J. F. Simons, E. S. Trombetta, and A. Helenius (1996) *EMBO J.* **15**, 6921–6930.
14. A. Zapun, S. M. Petrescu, P. M. Rudd, A. R. Dwek, D. Y. Thomas, and J. J. M. Bergeron (1997) *Cell* **88**, 29–38.
15. D. B. Williams (1995) *Biochem. Cell Biol.* **73**, 123–132.
16. C. Hammond and A. Helenius (1994) *Science* **266**, 456–458.
17. D. F. Qu, J. H. Teckman, S. Omura, and D. H. Perlmutter (1996) *J. Biol. Chem.* **271**, 22791–22795.
18. D. N. Hebert, B. Foellmer, and A. Helenius (1996) *EMBO J.* **15**, 2961–2968.
19. W. Chen, J. Helenius, I. Braakman, and A. Helenius (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6229–6233.
20. M. C. Sousa, M. A. Ferrero-Garcia, and A. J. Parodi (1992) *Biochemistry* **31**, 97–105.
21. S. E. Trombetta and A. J. Parodi (1992) *J. Biol. Chem.* **267**, 9236–9240.
22. P. J. Thomas, B. Qu, and P. L. Pederson (1995) *Trends Biochem. Sci.* **20**, 456–459.
23. V. Gross, T. Andus, T. A. Tranhl, R. T. Schwartz, K. Decker, and P. C. Heinrich (1983) *J. Biol. Chem.* **258**, 12203–12209.

## Suggestions for Further Reading

24. C. Hammond and A. Helenius (1993) A chaperone with a sweet tooth, *Curr. Biol.* **3**, 884–886.
25. J. J. M. Bergeron, M. B. Brenner, D. Y. Thomas, and D. B. Williams (1994) Calnexin: a membrane-bound chaperone of the endoplasmic reticulum, *Trends Biochem. Sci.* **19**, 124–128.
26. C. Hammond and A. Helenius (1995) Quality control in the secretory pathway, *Curr. Opin.*

Cell Biol. **7**, 523–529.

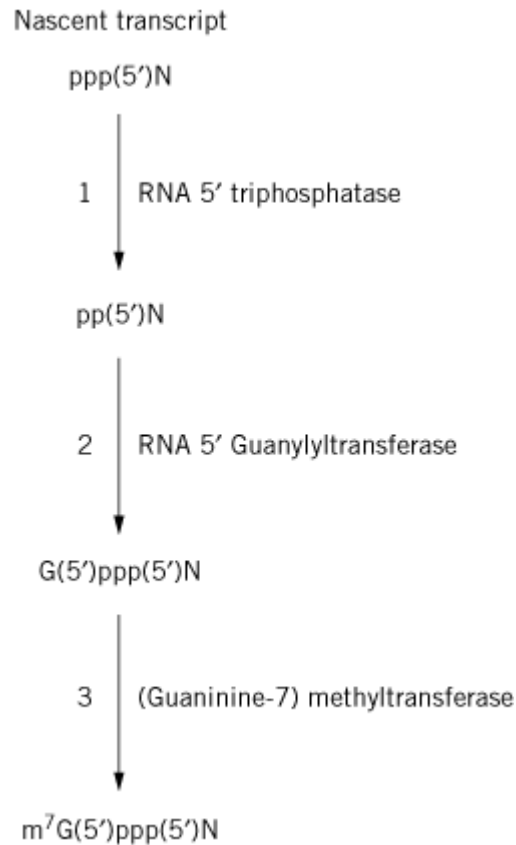
27. D. B. Williams and T. A. Watts (1995) Molecular chaperones in antigen presentation, *Curr. Opin. Immunol.* **7**, 77–84.
28. M.-J. Gething, ed. (1997) *Molecular Chaperones and Folding Catalysts*, Oxford University Press, Oxford (contains detailed articles on calnexin and calreticulin).
29. S. Dedhar (1994) Novel functions for calreticulin: interactions with integrins and modulation of gene expression? *Trends Biochem. Sci.* **19**, 269–271.
30. M. G. Coppelino, M. J. Woodside, N. Demaurex, S. Grinstein, R. St-Arnaud, and S. Dedhar (1997) Calreticulin is essential for integrin-mediated calcium signalling and cell adhesion, *Nature* **386**, 843–847.

## 5' Cap

The 5' cap is characteristic of all RNAs transcribed by [RNA polymerase II](#), and it is added to the nascent transcript soon after transcriptional initiation. Many lines of evidence suggest that the cap contributes to many aspects of RNA metabolism, including RNA stability ([1](#)), pre-mRNA **splicing** ([2, 3](#)), [polyadenylation](#) ([4](#)), U **small nuclear RNA** (snRNA) export ([5](#)), and [translation of messenger RNA](#) ([6](#)). The influence of the cap on RNA metabolism, in both the [nucleus](#) and the cytoplasm, is mediated by [cap-binding proteins](#).

The cap consists of an inverted 7-methyl guanosine nucleotide linked by a 5'-5'-triphosphate linkage to the first template-encoded nucleotide of the RNA, to give the following structure: m<sup>7</sup>G(5')ppp(5')N. The cap modification is one of the earliest covalent modifications detected on a nascent RNA transcript, and it occurs cotranscriptionally ([7](#)). The actual capping reaction requires the sequential action of several enzymatic activities and is well understood (see [Fig. 1](#)). The capping enzyme possesses an RNA 5'-triphosphatase activity, which converts the 5' terminal triphosphate of the RNA to a diphosphate. RNA guanylyltransferase then catalyzes the transfer of GMP, derived from GTP, to the diphosphate, resulting in the formation of the G(5')ppp(5')N cap. This cap structure is the substrate for methylation by RNA (guanine-7)-methyltransferase, which converts it to the monomethylated cap structure. In the absence of cap methylation, the reaction is reversible. The genes encoding the capping enzyme and methylase are essential in the yeasts *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*, underlining the importance an intact cap for viability ([8-10](#)). The cause of death in these mutants is probably due to cumulative defects in the various processes in which the cap is involved, such as RNA stability and **splicing** ([11](#)).

**Figure 1.** The capping reaction at the 5' nucleotide (N) of an RNA transcript: (1) the 5'-triphosphate of the nascent transcript is hydrolyzed by RNA 5'-triphosphatase to a diphosphate; (2) mRNA guanylyltransferase adds the cap from GMP (which is derived from GTP). The cap is then methylated on the N-7 (position 7 nitrogen) of guanosine by mRNA (guanine-7) methylase.



### Bibliography

1. Y. Furuichi, A. LaFiandra, and A. J. Shatkin (1977) *Nature* **266**, 235–239.
2. M. Konarska, R. Padgett, and P. Sharp (1984) *Cell* **38**, 731–736.
3. A. Krainer, T. Maniatis, B. Ruskin, and M. Green (1984) *Cell* **36**, 993–1005.
4. C. Cooke and J. C. Alwine (1996) *Mol. Cell. Biol.* **16**, 2579–2584.
5. J. Hamm and I. W. Mattaj (1990) *Cell* **63**, 109–118.
6. A. Shatkin (1985) *Cell* **40**, 223–224.
7. M. Salditt-Georgieff, M. Harpold, S. Chen-Kiang, and J. E. Darnell Jr. (1980) *Cell* **19**, 69–78.
8. S. Shuman and B. Schwer (1995) *Mol. Microbiol.* **17**, 405–410.
9. Y. Shibagaki, N. Itoh, H. Yamada, S. Nagata, and K. Mizumoto (1992) *J. Biol. Chem.* **267**, 9521–9528.
10. B. Schwer and S. Shuman (1994) *Proc. Natl. Acad. Sci. USA* **91**, 4328–4332.
11. B. Schwer and S. Shuman (1996) *RNA* **2**, 574–583.

### Suggestions for Further Reading

12. J. D. Lewis, S. Gunderson, and I. W. Mattaj (1995) The influence of 5' and 3' end structures on pre-mRNA metabolism, *J. Cell Sci. Suppl.* **19**, 13–19.
13. J. D. Lewis and E. Izaurralde (1997) The role of the cap structure in RNA processing and nuclear export, *Eur. J. Biochem.* **247**, 461–469.

## Cap-Binding Proteins

Most of the functions in the [nucleus](#) and cytoplasm of the [cap](#) of [messenger RNA](#) (mRNA) are mediated by cap-binding proteins (CBP), which show highly specific binding to the cap structure, m<sup>7</sup>G(5')ppp(5')N. Two major cap-binding activities have been characterized in **eukaryotes**. The first is the nuclear cap-binding complex (CBC), which plays a role in pre-mRNA **splicing** (1-3), [polyadenylation](#) (4), and U **small nuclear RNA** (snRNA) export (5). The second, eIF4E (6), is part of a trimeric complex, eIF4F, which is required for cap-dependent [translation](#) of mRNA.

CBC was initially identified in nuclear extracts from HeLa cells as an activity that bound specifically to a normal 7-methyl guanosine-5'-capped RNA. Purification of the activity showed that it was a heterodimeric complex comprising two subunits; CBP80 and CBP20 (1). Neither protein alone can bind specifically to the 5' cap—only as a heterodimer. Current evidence has shown that CBC is required for efficient (1) pre-mRNA splicing, (2) polyadenylation, and (3) the export of RNA polymerase II-transcribed U snRNA (see below).

CBC is required for efficient splicing of the cap-proximal **intron**, where its major function is to facilitate the binding of U1 snRNP to the cap-proximal 5' splice site (3) (see [Gene Splicing](#)). Experiments in the yeast *Saccharomyces cerevisiae* and in HeLa splicing extracts have demonstrated that this function of CBC is conserved in [evolution](#) (2, 7).

A role for CBC in polyadenylation has also been demonstrated *in vitro* (4): (see [Polyadenylation](#)). Depletion of CBC from HeLa cell polyadenylation extracts results in a reduction in the efficiency of polyadenylation. Further analysis of this defect demonstrated that CBC was required for efficient cleavage of the pre-mRNA, but not for the polyadenylation reaction itself.

Nucleocytoplasmic transport of some classes of RNA, specifically the 5' capped uracil-rich (U) snRNAs, is facilitated by CBC. [Antibodies](#) raised against one of the subunits, CBP20, can specifically inhibit the interaction of CBC with the cap and, as a consequence, inhibit the export of these U snRNAs from the nucleus to the cytoplasm (5).

The other CBP, eIF4E, is required for cap-dependent translation of mRNA in the cytoplasm, and it can bind directly to the cap as a monomer (6). In order to function in translation, eIF4E has to assemble with two other [polypeptide chains](#) in a heterotrimeric complex of eIF4A, eIF4G, and eIF4E, to form a complex known as eIF4F. This trimeric complex binds to the cap of mRNAs and complexes with a second translational initiation complex, eIF3, which has an [RNA helicase](#) activity. This helicase activity is thought to unwind **secondary structure** in the mRNA and to allow the cap-dependent association of the 40S **ribosomal** subunit with the mRNA. The 40S subunit then scans for the [initiation codon](#) complexes with the 60S ribosomal subunit to initiate translation.

### Bibliography

1. E. Izaurralde, J. Lewis, C. McGuigan, M. Jankowska, E. Darzynkiewicz, and I. Mattaj (1994) *Cell* **78**, 657–668.
2. J. Lewis, D. Görlich, and I. Mattaj (1996) *Nucl. Acids Res*, **24**, 3332–3336.
3. J. D. Lewis, E. Izaurralde, A. Jarmolowski, C. McGuigan, and I. W. Mattaj (1996) *Genes Devel.* **10**, 1683–1698.
4. S. M. Flaherty, P. Fortes, E. Izaurralde, I. W. Mattaj, and G. M. Gilmartin (1997) *Proc. Nat. Acad. Sci. USA* **94**, 11893–11898.
5. E. Izaurralde, J. Lewis, C. Gamberi, A. Jarmolowski, C. McGuigan, and I. W. Mattaj (1995) *Nature* **376**, 709–712.
6. N. Sonnenberg, M. Rupprecht, W. Merrick, and A. Shatkin (1979) *Proc. Natl. Acad. Sci. USA*

75, 4345–4349.

7. H. Colot, F. Stutz, and M. Rosbash (1996) *Genes Devel.* **10**, 1699–1708.

### **Suggestions for Further Reading**

8. J. D. Lewis and E. Izaurralde (1997) the role of the cap structure in RNA processing and nuclear export, *Eur. J. Biochem.* **247**, 461–469.
9. W. Merrick and J. Hershey (1996) in J. Hershey, M. Mathews, and N. Sonnenberg, eds., *Origins and Targets of Translational Control*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 1–29.



## Capillary Zone Electrophoresis

Capillary zone electrophoresis (CZE) is [electrophoresis](#) in very thin capillaries. It is an analytical separation method that is compatible with high field strength (200 to 800 V/cm) and small samples (micrograms to nanograms). The high field strength makes the separation rapid and gives high resolving power (1). The very small capillary diameter, less than 0.2 mm, counteracts dispersion of the zone, presumably through interaction of the analyte with the inner wall of the capillary. The sample is detected by its absorbance or fluorescence as it proceeds past a stationary detector at the end of the migration path.

The inner walls of the silicate capillary carry negatively charged groups, so [electroendosmosis](#) is significant; this can be used for the purpose of separation, but it is usually suppressed by coating the inner wall with a polymer such as [polyacrylamide](#) gel or dextran. Separations based predominantly on size and shape differences are achieved in CZE in the presence of soluble polymers to produce a **molecular sieve** effect. A variety of uncharged hydrophilic polymers in a wide range of molecular weights are available for that purpose (2-4). Although providing high resolving power (see [Gel Electrophoresis](#)), gels are less applicable to CZE because the gelation process may introduce air bubbles and inhomogeneities within the capillary; also, a gel-filled capillary can only be used once or a few times, while polymer solutions can be replaced easily. One drawback of CZE is that the analyte bands are enclosed within the tube and therefore are unavailable for detection by immunological, hybridization, and staining techniques. However, the eluate of CZE can be subjected to [mass spectrometry](#) to provide molecular weights of the species detected, in addition to their mobilities (5).

### Bibliography

1. F. Foret and P. Bocek (1989) *Adv. Electrophoresis* **3**, 273–347.
2. K. Ganzler, K. S. Greve, A. S. Cohen, B. L. Karger, A. Guttman, and N. C. Cooke (1992) *Anal. Chem.* **64**, 2665–2671.
3. M. C. Ruiz-Martinez, J. Berka, A. Belenkii, F. Foret, A. W. Miller, and B. L. Karger (1993) *Anal. Chem.* **65**, 2851–2858
4. D. Tietz, A. Aldroubi, H. Pulyaeva, T. Guszczynski, M. M. Garner, and A. Chrambach (1992) *Electrophoresis* **13**, 614–615.
5. D. Figeys, I. van Oostveen, A. Ducret, and R. Aebersold (1996) *Anal. Chem.* **68**, 1822–1828.

### Suggestions for Further Reading

6. P. Gebauer and P. Bocek (eds.) (1995) *Symposium on capillary electrophoresis theory*. *Electrophoresis* **16**, 1985–2174.
7. P. D. Grossman and J. C. Colburn (1992) *Capillary Electrophoresis*, Academic Press, New York, pp. 1–352.
8. B. L. Karger (ed.) (1993) *Symposium on capillary electrophoresis*. *Electrophoresis* **14**, 371–560.

## Capsids, Viral

The capsids of **viruses**—the protein part of the virus particle (virion) that surrounds and protects the nucleic acid genome—have been the subject of intense study by biochemists and structural biologists for over 50 years (see [Virus Structure](#)). Capsids are also noteworthy in that they provide one of the few examples in which the detailed properties of a biological system have been predicted successfully from “first principles.” This entry describes the principles on which we understand capsid structure—that is, what we expect capsids to be like and why we expect that, and then describes the ways in which real viruses do or do not follow those expectations.

A **gene** of double-stranded DNA can encode a [protein](#) of only about 1/20 the mass of the gene itself. As a consequence, if a virus is to encode its own capsid (and in almost all cases, viruses do so), it will only be able to make enough protein to produce a useful sized capsid if it can use multiple copies of the protein(s) encoded in its genes. Crick and Watson ([1](#)), who made the first concrete proposal for how proteins might be arranged in virus capsids, assumed that the capsids would be made of multiple identical virus-encoded protein subunits. They also made another assumption, based on a prominent property of proteins, namely that proteins are very specific in the interactions they make with other molecules, presumably including other proteins in the structure of a virus capsid. They assumed that the identical protein subunits would be packed into the capsid structure in such a way that they all made the identical set of contacts with their neighbors—so-called **equivalent packing**. Mathematically, there are only two general ways to satisfy these assumptions when packing asymmetric objects like proteins; these are to arrange the protein subunits with helical symmetry or to arrange them with one of the cubic symmetries. Helical symmetry is easily visualized: The subunits are arranged in a helical array as if they were on successive steps of a spiral staircase. Formally, each subunit is related to the preceding one by a characteristic rotation and a translation in the direction of the helix axis. “Cubic” symmetry refers to a small group of symmetries characterized by having multiple axes of rotational symmetry. These correlate with the five classical Platonic solids, which have these symmetry axes. Besides the cube, from which the entire group of symmetries takes its name, the Platonic solids include the tetrahedron, the octahedron, the dodecahedron, and the icosahedron. Of these, the icosahedron allows (or more properly, arranging protein subunits according to the five-, three-, and twofold rotational symmetry axes of an icosahedron) allows the largest structure to be made with a given size of subunit for any of the group. More to the point, of this group of symmetries, it is icosahedral symmetry that is used by real viruses.

## 1. Helical Symmetry

Many viruses can be shown to have helical symmetry. The best studied of these is the well-known [tobacco mosaic virus](#) (TMV), which has a virion with only one type of protein subunit, present in something over 2000 copies and arranged in a simple helix with 16.33 subunits per turn. There is in principle no limit on how long a helical structure such as this can be, but in TMV it is limited during assembly by the length of the genomic RNA. In the mature virion, the single-stranded RNA genome follows the helix of the proteins and lies in the groove between successive layers of the helix on the inside of the helical tube, three nucleotides per protein subunit. A domain of each subunit closes over this groove after the RNA has entered it during assembly and effectively seals off the RNA from the solvent ([2](#)).

Structurally somewhat more complex examples of helically symmetric viruses are the filamentous bacteriophages, such as phage M13. The major protein subunit of these viruses is arranged in what might best be described as a small number of helices running in parallel up a cylindrical tube. The circular single-stranded DNA genome of these viruses is stretched the length of the helical tube. The filamentous phages have small numbers of copies of additional virus-encoded proteins present on the ends of the virions. These have roles in virion assembly and in **virus infection** ([3](#)).

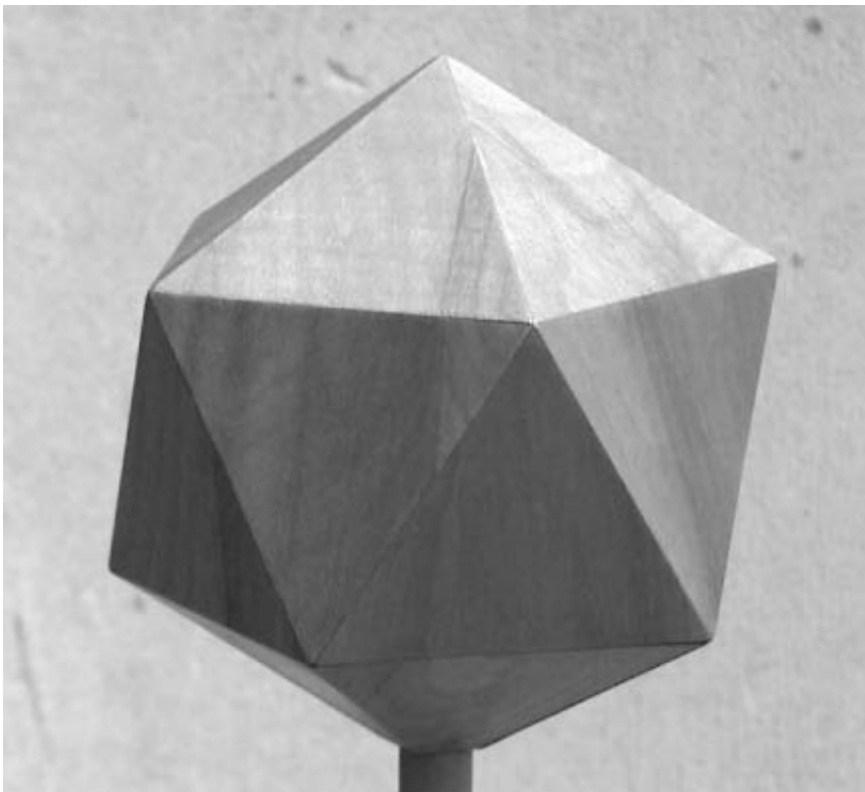
Many enveloped viruses, for example, [influenza virus](#), have their genome—single-stranded RNA in the case of Influenza—wrapped in a helical array with a virus-encoded protein. The resulting flexible nucleoprotein rods are enclosed in the viral envelope, one such structure for each of the eight

different genome segments in the case of influenza.

## 2. Icosahedral Symmetry and Quasi-Equivalent Packing

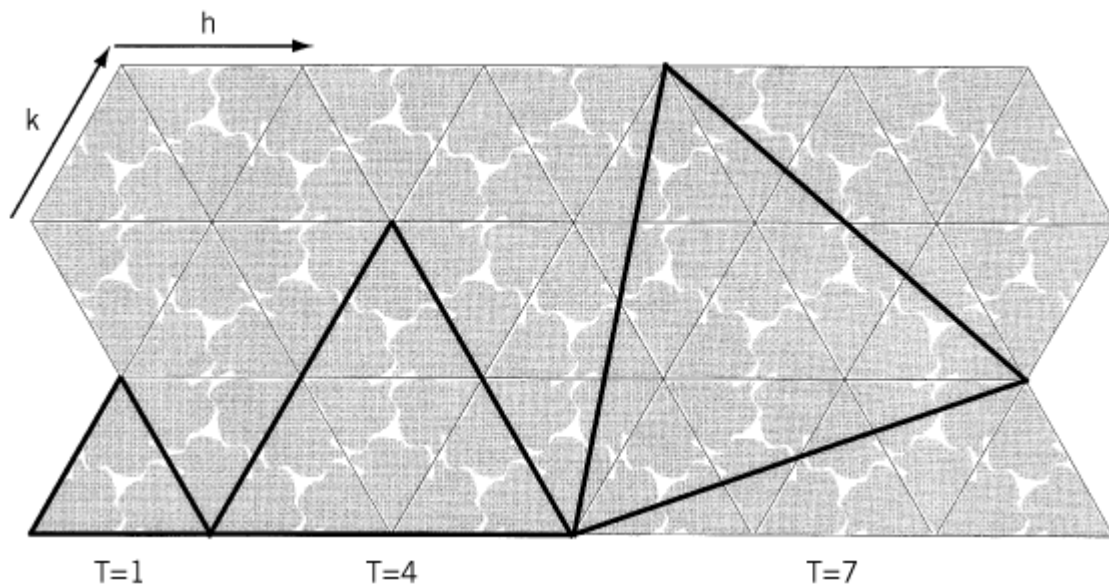
An icosahedron (Fig. 1) has 20 faces, each an equilateral triangle. Its symmetry is defined by axes of rotational symmetry, each passing through the center of the icosahedron: five-fold axes passing through the corners, three-fold axes passing through the centers of the faces, and twofold axes passing through the centers of the edges. Strictly speaking, an icosahedral virus capsid is not an icosahedron; instead, its protein subunits are related to each other by the symmetry axes of an icosahedron. However, thinking of the protein subunits as lying on the surface of an icosahedron provides a convenient way to visualize and discuss the structure of the capsid, as well as a reasonable approximation to the truth. To make an icosahedrally symmetric capsid, we place three identical protein subunits on each face of the icosahedron, symmetrically arrayed about the center. These subunits will necessarily also find themselves in symmetrical relations with their fellows around all the five-, three-, and twofold axes of the icosahedron, and inspection of a model of such a structure shows that all of the subunits lie in equivalent relationships to their neighbors.

**Figure 1.** An icosahedron.



Current knowledge of virus structure allows us to see that some small virus capsids—for example, the Parvoviruses—can be adequately described as such 60-subunit icosahedrally symmetric structures. However, it was already clear by about 1960 that the majority of “spherical” viruses are too large and have too many subunits to be so described. Because there is mathematically no way to make a larger structure while preserving strict equivalence in subunit packing, it was not clear that this way of describing capsid structure would be applicable to most viruses. However, Caspar and Klug (4) proposed a clever extension of Crick and Watson's ideas that has provided a framework for understanding and describing the great majority of spherical virus capsids. Caspar and Klug introduced the idea of “quasi-equivalence” in protein subunit packing. In a nutshell, they suggested





### 3. Icosahedral Virus Capsids in Nature

As the structures of actual virus capsids have been determined in increasing numbers and in increasing detail over the past 30 years, the general predictions of the Caspar and Klug theory have largely been confirmed. At the same time, there are now many examples of “variations on the theme”—features of the structures that are not explicitly predicted by the theory, but are not at odds with its basic concepts. In addition, there are at least two examples of structures that should not have occurred. The definitive test of how the capsid structure is organized is a high-resolution structure by [X-ray crystallography](#); however, it is often possible to get useful information from lower resolution structural measurements. For example, the hexamers and pentamers of the capsid subunit often cluster in such a way that they can be seen by [electron microscopy](#) as a distinct morphological unit, called a **capsomere**. Examining the geometrical relationships among the capsomeres on the surface of a capsid can lead to a good idea of the triangulation number of the capsid.

Parvoviruses, as mentioned above, have 60-subunit  $T = 1$  structures (5). The same is true for the small **fX174 bacteriophage**, but in this case there are two proteins present in 60 copies each, rather than only one (6). Here we can think of the heterodimer as the basic building block, with 60 copies of that heterodimer arranged with  $T = 1$  icosahedral symmetry.

$T = 3$  capsids, the first group that requires that quasi-equivalence be invoked, are especially well-populated, with examples from plant, animal, and bacterial viruses. [Tomato bushy stunt virus](#) (TBSV), the first of this group to be solved to atomic resolution, conforms quite well to the expectations of the theory, having 180 chemically identical subunits, each of which occupies one of three similar but nonidentical positions in the  $T = 3$  lattice (7). One set of subunits in TBSV, those occupying the “C” position, send an end of their polypeptide chain toward the threefold symmetry axes, where they intertwine with the corresponding parts of the two symmetrically related subunits. This was the first of what are now several examples of capsid proteins interdigitating and intertwining with their neighbors; this is presumed to provide mechanical strength to the capsid. [Poliovirus](#) (a Picornavirus) has a structure very similar to that of TBSV, including the positioning of subunits and even the fold of the polypeptide chains within those subunits, but in this case the three quasi-equivalent positions are occupied by three similar but chemically distinct proteins (8). We might imagine that two of the three coding regions encoding these three proteins were the result of [gene duplication](#) in an ancestral virus that, like TBSV, had only one capsid protein gene. A slightly different variation on this theme is found in [Cauliflower Mosaic virus](#) (a Comovirus), in which the three quasi-equivalent positions are occupied by two proteins, one of which consists of two very

similar **domains**, each occupying one of the three quasi-equivalent positions. These last two examples can be regarded either as slightly noncanonical  $T = 3$  structures or as  $T = 1$  structures made of 60 copies of the heterotrimer (Picornavirus) or heterodimer (Comovirus).

Many capsid structures with larger triangulation numbers have been identified and characterized. Among these are [lambda phage](#), **P22 phage**, bacteriophage HK97, and others ( $T = 7$ ), **T4 phage** and [Reovirus](#) ( $T = 13$ ), [Herpesvirus](#) ( $T = 16$ ), and [Adenovirus](#) ( $T = 25$ ). The largest triangulation number that has been definitively determined is that of the algal virus PBCV1 ( $T = 169$ ). A particularly instructive example is provided by bacteriophage P2 and its satellite phage P4. The P2 capsid protein assembles into a  $T = 7$  shell; but in the presence of P4, which does not encode its own capsid protein, the P2 capsid protein assembles into a  $T = 4$  structure, just big enough to enclose the P4 genome. This is accomplished through the agency of a P4-encoded protein that associates transiently with the P2 capsid protein during assembly and directs it into the  $T = 4$  geometry (9).

Although these large viruses all fit into the general picture envisioned by Caspar and Klug, they are replete with exceptions and extensions to the original picture, and these variations all expand our view of how capsids can be constructed. For bacteriophage T4 and Adenovirus, the proteins that make the pentamers at the five-fold symmetric corners of the capsid are encoded by a different gene from the proteins at the six-fold positions in the remainder of the capsid. This is analogous to the Picornavirus example above and presumably reduces the amount of conformational or bonding versatility demanded of any one protein, at the relatively minor expense of encoding an additional protein. In T4, the major component of the capsid, which occupies the hexamer positions, is, as expected, a hexamer, but in Adenovirus these positions are occupied by *trimers* of the “hexon protein” (10). However, the hexon protein is organized into two similarly folded domains, and the structure of the trimeric hexon is very close to a sixfold symmetric arrangement of those domains.

Some viruses—particularly but not exclusively some of the double-stranded DNA phages—have prolate capsids, which have a standard icosahedral arrangement of subunits, except that the shell is elongated along one of the fivefold axes of symmetry and an extra band of hexamers is inserted around the equator. Thus bacteriophage T4 has an elongated  $T = 13$  capsid, and bacteriophage f29 has an elongated  $T = 4$  capsid. These structures pose interesting questions with regard to how their length is specified and accurately achieved, but they cause no serious problems for the idea of quasi-equivalent packing.

Most viruses have other capsid subunits in addition to the main, icosahedrally packed, subunit protein. In some cases, these are present in equimolar amounts with the main subunit and packed with the same symmetry. Thus phages  $\lambda$  and T4 have well-studied examples of such “decoration proteins”; Herpesvirus has a protein clustered as trimers that fits this description, but is systematically absent from positions immediately surrounding the pentamers. Proteins in this category are known in some cases to provide additional strength and stability to the capsid. Structurally, they can be regarded simply as additional domains of the main capsid protein, albeit ones that are encoded by separate genes and generally join the structure at a different time.

Many capsids also have “minor” proteins that are not arrayed with icosahedral symmetry. The portal protein of the double-stranded DNA phages, for example, forms a grommet-like 12-subunit oligomer that replaces a pentamer at one of the five-fold corners of the icosahedral shell and provides an attachment site for the six-fold symmetric helical tail (11). Adenovirus has a trimeric spike extending out from the shell along each of its five-fold symmetry axes.

The first radical deviation from the expectations of the Caspar–Klug ideas was found in the virion of Papovavirus [SV40](#). The capsomeres of SV40 visible by electron microscopy are arranged as expected for a  $T = 7$  structure. However, all of the capsomeres are pentamers of the VP1 subunit (12), including the 60 capsomeres situated at positions of sixfold local symmetry, which would be expected to be hexamers in order to interact quasi-equivalently with their environment. The subunits adapt to this extraordinary state of affairs by having the part of the polypeptide chain that contacts

the neighboring capsomere located on a flexible arm corresponding to the C-terminus of the subunit. The C-terminal arm leaves its home subunit at dramatically different angles in different cases, allowing it to interact with its neighbor in essentially the same way in each case. This it does by invading the structure of the neighbor and forming one strand of a [b-sheet](#) structure in that subunit. This might be regarded as a form of quasi-equivalent interaction, but of a sort requiring a much more radical subunit flexibility than envisioned by Caspar and Klug.

Another surprising arrangement of subunits is found in fungal virus L-A and bacteriophage f6. These capsids have 120 identical protein subunits, not one of the “allowed” numbers. These are  $T = 1$  structures, built of 60 asymmetric homodimers. The unexpected feature of this arrangement is that the two chemically identical subunits occupy nonequivalent positions. Nonetheless, these capsids, as for those of SV40, evidently function well enough to have survived [natural selection](#) (ie, very well indeed), and an understanding of their structures enlarges our view of the capabilities of proteins.

### Bibliography

1. F. Crick and J. D. Watson (1956) *Nature* **177**, 473–475.
2. J. N. Champness et al. (1976) *Nature* **259**, 20–24.
3. L. Makowski and M. Russel (1997) In *Structural Biology of Viruses* (W. Chiu, R. M. Burnett, and R. L. Garcea, eds.), Oxford University Press, New York, pp. 352–380.
4. D. Caspar and A. Klug (1962) *Cold Spring Harbor Symp. Quant. Biol.* **27**, 1–24.
5. J. Tsao et al. (1991) *Science* **251**, 1456–1464.
6. R. McKenna, L. Ilag, and M. Rossmann (1994) *J. Mol. Biol.* **237**, 517–543.
7. A. J. Olson, G. Bricogne, and S. C. Harrison (1983) *J. Mol. Biol.* **171**, 61–93.
8. J. Hogle, M. Chow, and D. Filman (1985) *Science* **229**, 1358–1365.
9. O. Marvik et al. (1995) *J. Mol. Biol.* **245**, 59–75.
10. R. M. Burnett (1997) In *Structural Biology of Viruses* (W. Chiu, R. M. Burnett, and R. L. Garcea, eds.), Oxford University Press, New York, pp. 209–238.
11. C. Basinet and J. King (1985) *Annu. Rev. Microbiol.* **39**, 109–129.
12. R. L. Garcea and R. C. Liddington (1997) In *Structural Biology of Viruses* (W. Chiu, R. M. Burnett, and R. L. Garcea, eds.), Oxford University Press, New York, pp. 187–208.

### Suggestions for Further Reading

13. W. Chiu, R. M. Burnett, and R. L. Garcea (eds.) (1997) *Structural Biology of Viruses* Oxford University Press, New York.
14. B. N. Fields, D. M. Knipe, and P. M. Howley (eds.) *Virology* (1996) Lippincott-Raven, Philadelphia.
15. J. E. Johnson and J. A. Speir (1997) *J. Mol. Biol.* **269**, 665–675.
16. J. E. Johnson and A. J. Fisher (1994) In *Encyclopedia of Virology*, Vol. **1** (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1573–1586.

### Carbon Isotopes

Carbon is element number 6 in the periodic table and has valence states of 2, 3, or 4 ([1](#)). Thirteen isotopes of carbon have been identified ([2](#)), ranging in atomic mass number from  $^{12}\text{C}$  (half-

life =  $3 \times 10^{-22}$  s) to  $^{20}\text{C}$  (half-life = 0.01 s) (see [Radioactivity](#) and [Radioisotopes](#)). Two stable isotopes of carbon are found in nature:  $^{12}\text{C}$  at 98.9%, and  $^{13}\text{C}$  at 1.1% abundance.

The most important radioactive isotope of carbon is  $^{14}\text{C}$  (half-life = 5730 years). Carbon-14 decays by beta-minus emission to nitrogen-14, which is stable. Carbon-14 decay yields one beta particle, with an energy of 0.156 MeV maximum, 0.0495 MeV on average. Carbon-14 is produced naturally in the earth's atmosphere by cosmic ray interactions with stable nitrogen-14 according to the reaction  $^{14}\text{N}(n,p)^{14}\text{C}$ . Atmospheric  $^{14}\text{CO}_2$  is inhaled by animals and respired by plants in photosynthesis, incorporating small amounts of carbon-14 in all living organisms. Consequently, carbon-14 dating is useful for determining the age of organic matter such as wood and archeological specimens. Hundreds of different organic molecules have been labeled with carbon-14 as a tracer for studying biochemical processes and are available from commercial suppliers.

Another radioactive isotope, carbon-11 (half-life = 20.38 min), is used in nuclear medicine diagnostics with positron-emission tomography. Carbon-11 decays by beta-plus (positron) emission to boron-11, which is stable. Carbon-11 decay yields 0.98 beta particles, with an energy of 0.960 MeV maximum, 0.386 MeV on average. Positron emission is characterized by twin 0.511-MeV photons that result from the annihilation of a positron and an electron and allow it to be detected externally. Carbon-11 is produced by proton accelerators according to the reaction  $^{11}\text{B}(p,n)^{11}\text{C}$ . Carbon-11 positrons are useful for  $^{11}\text{C}$ -acetate (3) and  $^{11}\text{C}$ -palmitate (3) metabolism kinetic imaging in the assessment of heart disease (4).

Carbon-13 is very useful in studies involving nuclear magnetic resonance ([NMR](#)), because it gives an NMR signal, whereas the more abundant isotope, carbon-12, does not.

Carbon-14 was first produced artificially for metabolism studies in 1939. At present, it is produced commercially by neutron irradiation of beryllium nitride or aluminum nitride according to the reaction  $^{14}\text{N}(n,p)^{14}\text{C}$ . The low cross section of this reaction and the long half-life of carbon-14 results in low specific activities, which can make carbon-14 labeling of organic materials quite difficult.

Carbon-14 labeling of organic molecules has numerous applications in the biomedical sciences. Examples include studies of glucose metabolism and conversion to  $^{14}\text{CO}_2$ , galactose oxidation, carbohydrate metabolism, the Krebs cycle, **nucleic acid synthesis** and **hybridization**, [autoradiography](#), or [fluorography](#), and for studying many other biochemical pathways. Among the most important carbon-14 compounds used in biochemistry and metabolism studies are the carbon dioxides and monoxide, carboxyl-labeled acids, methanols, cyanides, carbides, formic acids, and cyanamides (5). Carbon-14 is detected by beta-particle liquid scintillation counting, [autoradiography](#), and [fluorography](#).

#### Bibliography

1. D. R. Lide and H. Pr. Frederikse, eds. (1995) *CRC Handbook of Chemistry and Physics*, CRC Press, Boca Raton, Fla.
2. Knolls Atomic Power Laboratory (1966) *Chart of the nuclides*, 15th ed., available from General Electric Company, San Jose, Calif.
3. H. R. Schön, H. R. Schelbert, A. Najafi et al. (1982) *Am. Heart J.* **103**, 532–547.
4. T. L. Rosamond, D. R. Abendschein, B. E. Sobel et al. (1987) *J. Nucl. Med.* **28**, 1322–1329.
5. *The Radiochemical Manual*, 2nd ed. (1966) Amersham, Bucks, England.

#### Suggestions for Further Reading

6. J. R. Catch (1961) *Carbon-14 Compounds*, Butterworths, London.



7. J. C. Harbert, W. C. Eckelman, and R. D. Neumann, eds. (1996) *Nuclear Medicine: Diagnosis and Therapy*, Thieme Medical Publishers, Inc., New York.

## Carbonic Anhydrase

The enzyme carbonic anhydrase (CA) catalyses the reversible hydration of carbon dioxide to bicarbonate and a proton, using a catalytic zinc ion bound at the active site (1, 2). Carbonic anhydrases are of interest because of their varied metabolic roles, efficient catalytic mechanisms, and high metal ion specificity. The enzyme is found in animals, plants, and bacteria, where it plays roles in respiration, [photosynthesis](#), and CO<sub>2</sub> fixation. Three genetically distinct families of CA proteins have been identified: aCA isozymes, found in all animals; bCA, commonly found in plants and bacteria; and gCA, found in a variety of bacteria (3). Of these classes of CA proteins, the a-CA family is the best characterized. This article provides a brief overview of the metabolic roles and genetic structure of the a-CA human isozymes and describes the catalytic reaction mechanism and metal-binding properties of the high activity, well-characterized human a-CA isozyme II (CAII). Recent work on the b- and g-CA enzymes, which catalyze CO<sub>2</sub> hydration using radically different protein structures, is summarized.

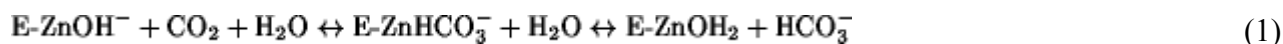
### 1. Metabolic Function and Genetic Structure of Human a-CA Isozymes

Carbonic anhydrase II (CAII) was first discovered in 1933 in human erythrocytes (4), where it facilitates CO<sub>2</sub> transport in respiration by converting CO<sub>2</sub>, released as a metabolic by-product from tissues, to bicarbonate. Bicarbonate is carried by the bloodstream to pulmonary capillaries, where its conversion back to CO<sub>2</sub> is catalyzed by CA for release by the lung. CAII is also expressed in many other tissues, including ocular epithelium, where it plays a role in maintaining intraocular pressure (5). A human genetic disease in which CAII is nonfunctional indicates that this isozyme is also crucial for bone resorption and kidney function (6). Furthermore, the inhibition of CAII by aromatic sulfonamides (RSO<sub>2</sub>NH<sub>2</sub>) is used clinically to treat glaucoma and altitude sickness (5). Additionally, at least six more CA human **isozymes** with varied activities and locations have been discovered (5). The cytosolic isozymes CAI, CAIII, and CAVI are expressed in high concentrations in red blood cells, **muscle** cells, and salivary glands, respectively. CAIV is a membrane-bound isozyme important for regulating bicarbonate levels in the kidney, whereas CAV is a [mitochondrial](#) isozyme. CAVI is secreted by the salivary glands. The exact physiological role of several of these isozymes is still under investigation.

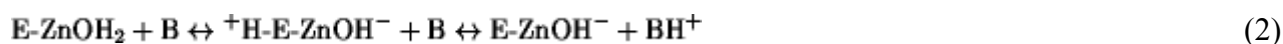
The genomic structure of the human CA isozymes lends insight into their [evolution](#) and regulation of their [transcription](#). CAI, CAII, and CAIII have been mapped to [chromosome](#) 8q22, CAV and VII to chromosome 16, CAIV to chromosome 17q, and CAVI to chromosome 1 (3). This chromosomal mapping, along with comparisons of DNA sequence homologies, indicates that [gene duplication](#) events over 450 million years ago led to distribution of the CA isozymes on four different chromosomes; further gene duplications occurred later, resulting in closely linked isozymes CAI, CAII, and CAIII on one chromosome, and CAIV and VI on another (3). The structures of the various human CA genes currently are being studied to determine the mechanism of their tissue-specific expression. CAI has an unusually long 5'-**untranslated region** containing two **promoters** that direct expression of CAI in erythroid cells and in the colon (7). Other CA genes are also likely to have complex structures to direct their expression in different locations and developmental stages.

## 2. Catalytic Mechanism and Zinc-Binding Properties of Human CAII

Carbonic anhydrase II is the most thoroughly characterized and most highly active of the CA isozymes, catalyzing the hydration of CO<sub>2</sub> with a rate constant of ~ 10<sup>6</sup> s<sup>-1</sup> (1). The crystal structure of CAII reveals a tightly bound zinc ion coordinated by the nitrogen atoms of three protein [histidine](#) residues, His94, His96, and His119, located at the bottom of a deep active site cleft (Fig. 1) (8). At physiological pH, a hydroxide ion is also bound by zinc, completing the tetrahedral geometry of the metal site. Residue Thr199 forms a [hydrogen bond](#) with the zinc hydroxide, orienting it for efficient catalysis. Catalysis occurs in two main steps (1), initiated by nucleophilic attack of the hydroxide ion on the carbonyl carbon of CO<sub>2</sub> to form bicarbonate.

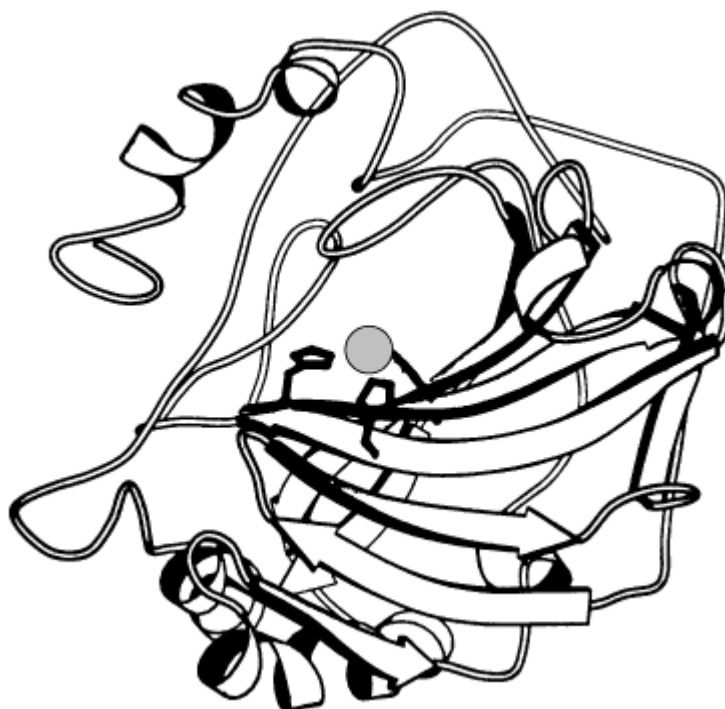


Water displaces bicarbonate, and the product is released. In the second, rate-limiting step in catalysis (Eq. (2)), the catalytically active enzyme species is regenerated by transfer of the product proton from the zinc-water to protein residue His64. This residue is exposed to solvent and transfers the proton to buffer B:



This proton shuttle is essential for rapid catalysis, as transfer of the proton to buffer is much faster than transfer to water. Carbonic anhydrase is one of the few biological systems in which proton transfer can be studied (9).

**Figure 1.** The crystal structure of human CAII determined by X-ray crystallography (8). It reveals a twisted b-sheet structure in which three histidine residues coordinate the catalytic zinc ion. Figure generated using MOLSCRIPT (13).



Just as the active site of CA has evolved for rapid catalysis, the protein structure is optimized for

tight, specific binding of the zinc ion, with a  $K_d$  of  $\sim 2$  pM (2, 10) (see [Zinc-Binding Proteins](#)). CAII binds only  $\text{Cu}^{2+}$ , and  $\text{Hg}^{2+}$  with comparable affinity, and neither of these metals confers catalytic activity to the enzyme (10). CAII has a much lower affinity, with a  $K_d$  of nanomolar, for metals such as  $\text{Co}^{2+}$ ,  $\text{Mg}^{2+}$ ,  $\text{Cd}^{2+}$ , and  $\text{Ni}^{2+}$ , and only the  $\text{Co}^{2+}$ -substituted enzyme is active. As the metal binding properties of this protein have been studied extensively, CAII is often used as a model for creating novel metal sites in proteins.

Studies in which [site-directed mutagenesis](#) has been used to alter the nature of conserved residues in the zinc site (2) reveal the importance of several factors for high metal affinity. First, in almost all CAII variants studied, the zinc ion retains tetrahedral geometry even if the surrounding protein structure must be rearranged to accommodate new side-chain positions. Second, the distance between the metal and ligand is crucial. Third, conserved residues that form hydrogen bonds with the histidine zinc ligands each contribute modestly to zinc affinity but play a large role in controlling the rates of metal equilibration (2, 11). The histidine ligands, the residues that form hydrogen bonds to these ligands, and the surrounding protein structure contribute to the reactivity and binding properties of the zinc ion.

### 3. Genetically Distinct CA Families Arose Through Functional Convergence

Examples of the b-CA and g-CA carbonic anhydrase families are only now being discovered and characterized. In an amazing example of evolutionary functional **convergence**, these protein families are completely unrelated genetically to one another and to the a-CA family, yet these enzymes catalyze the same reaction using a catalytic zinc ion (3). The structure determined by [X-ray crystallography](#) of a g-CA isolated from [archaeobacteria](#) (12) reveals a trimeric protein containing three active sites; in each active site, a zinc ion is bound at the monomer interface by three histidine residues from different polypeptide chains. As in a-CA proteins, the zinc is bound in tetrahedral geometry, the fourth ligand being a solvent molecule. In this case, a [glutamic acid](#) residue located near the zinc ion in g-CA may be the functional analog of Thr199 in human CAII. No structure of b-CA has been solved yet, but phylogenetic and spectroscopic data suggest that the zinc ion may be coordinated by glutamate or [cysteine](#) residues rather than histidine (3).

### 4. Summary

In recent years, the techniques of site-directed mutagenesis and X-ray crystallography have greatly expanded our knowledge of the relationships between protein structure and function in human CAII (2). The structure of CAII appears to have evolved for maximum catalytic activity and metal affinity. Applying these methods to the carbonic anhydrase enzymes that have evolved independently to catalyze the same reaction will give further insight into how protein structure can dictate the reactivity and affinity of metals.

### Bibliography

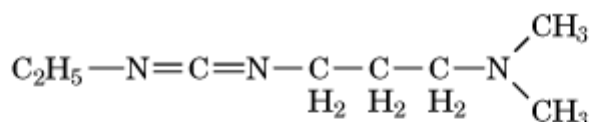
1. D. N. Silverman and S. Lindskog (1988) *Acc. Chem. Res.* **21**, 30–36.
2. D. W. Christianson and C. A. Fierke (1996) *Acc. Chem. Res.* **29**, 331–339.
3. D. Hewett-Emmett and R. E. Tashian (1996) *Mol. Phylogenet. Evol.* **5**, 50–77.
4. W. C. Stadie and H. O'Brien (1933) *J. Biochem.* **103**, 521–529.
5. S. J. Dodgson, R. E. Tashian, G. Gros, and N. D. Carter (1991) *The Carbonic Anhydrases: Cellular Physiology and Molecular Genetics*, Plenum Press, New York.
6. W. S. Sly et al. (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2752–2756.
7. H. J. M. Brady et al. (1991) *Biochem. J.* **277**, 903–905.
8. K. Hakansson, M. Carlsson, L. A. Svensson, and A. Liljas (1992) *J. Mol. Biol.* **227**, 1192–1204.
9. D. N. Silverman (1995) *Meth. Enzymol.* **249**, 479–503.

10. S. Lindskog and P. O. Nyman (1964) *Biochim. Biophys. Acta* **85**, 462–474.
11. C.-C. Huang et al. (1996) *Biochemistry* **35**, 3439–3446.
12. C. Kisker et al. (1996) *EMBO J.* **15**, 2323–2330.
13. P. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.

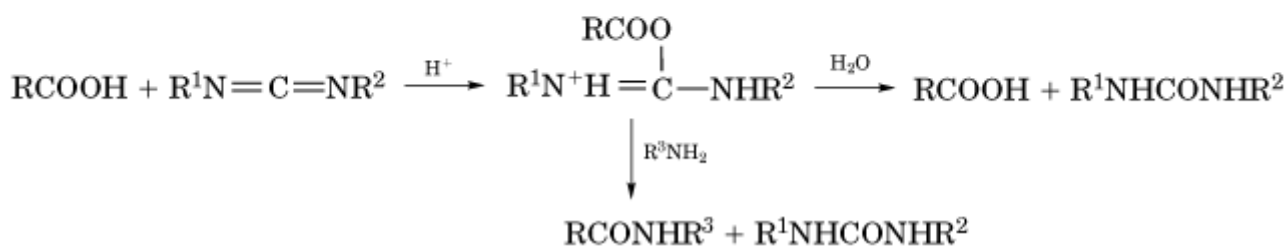
### Suggestions for Further Reading

14. D. W. Christianson (1991) Structural biology of zinc. *Adv. Prot. Chem.* **42**, 218–355.
15. S. J. Dodgson, R. E. Tashian, G. Gros, and N. D. Carter (1991) *The Carbonic Anhydrases: Cellular Physiology and Molecular Genetics*, Plenum Press, New York.
16. W. S. Sly and P. Y. Hu (1995) Human carbonic anhydrases and carbonic anhydrase deficiencies. *Ann. Rev. Biochem.* **64**, 375–401.
17. R. E. Tashian (1992) Genetics of the mammalian carbonic anhydrases. *Adv. Gen.* **30**, 321–356.

## Carboxyl Groups



[I] EDC



A carboxyl group has the formula  $\text{—COOH}$ . Among compounds that have a carboxyl group are carboxylic acids. A carboxyl group shows relatively weak acidity. Many carboxylic acids are found in nature as free acids or in the form of esters, amides, or salts. Biologically important acids are those found in the metabolic pathways of glycolysis, fatty acids, [amino acids](#), and so on. A carboxylic acid reacts with a base to form a salt and can be reduced to an aldehyde and an alcohol. A carboxylic acid forms an ester with an alcohol, an amide with an amine, and an acid anhydride with a carboxylic acid [see [Anhydrides](#)]. Carboxyl groups are not very reactive, and some catalyst or activation processes are required for their reactions.

There are three kinds of carboxyl groups in proteins, the C-terminal  $\alpha$ -carboxyl, the  $\beta$ -carboxyl of [aspartic acid](#), and the  $\gamma$ -carboxyl of [glutamic acid](#). They have intrinsic  $\text{pK}_a$  values of 3.8, 4.0, and 4.4, respectively. Carboxyl groups are often employed as catalytic groups in [enzymes](#), and they often play critical roles in **ligand binding**. Thus, carboxyl groups are important in protein chemistry.

### 1. Chemical Modification of Carboxyl Groups in Proteins

Esterification and amidation are the usual chemical modifications of carboxyl groups (1). Carboxyl groups in proteins are esterified by treating the dry protein with methanol-HCl. However, this method is drastic and not suitable for selective modification of carboxyl groups. Trialkyloxonium salts, such as triethylloxonium tetrafluoroborate, are commonly used for mild esterification of carboxyl groups in proteins (2). Diazoacetyl compounds, such as diazoacetamide, methyl diazoacetate, and *N*-diazoacetylglycinamide, react with carboxyl groups at acidic pH to form esters. Other reagents, such as acid halogen compounds and ethyleneimine have also been used to esterify the carboxyl groups in proteins.

Carboxyl groups are amidated with simple amines after activating the carboxyl groups with water soluble carbodiimides, such as 1-ethyl-3-(3-dimethylaminopropyl) carbodiimide (EDC) [1]. At pH 4.75 to 5, the reaction proceeds as in Scheme 1. Exhaustive amidation is achieved by using high concentrations of both amine (1 M) and EDC (0.1 M) at pH 4.75 under **denaturing** conditions (3). Limited amidation is possible using low concentrations of amine (0.1 M) and EDC (2.6 mM) under mild conditions at pH 5.0 (4). Isoxazolium salts, including Woodward's Reagent K (2-ethyl-5-phenylisoxazolium 3'-sulfonate), can be used in the place of EDC at pH 3 to 5. Amidation is also achieved by ammonolysis of the esters in liquid ammonia at  $-55^{\circ}\text{C}$ . In this case, selective amidation is possible starting with a selectively esterified protein (5).

### Bibliography

1. T. Imoto and H. Yamada (1989) In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, U.K., pp. 247–277.
2. S. M. Parsons et al. (1969) *Biochemistry* **8**, 700–712.
3. T.-Y. Lin and D.E. Koshland Jr. (1969) *J. Biol. Chem.* **244**, 505–508.
4. H. Yamada et al. (1981) *Biochemistry* **20**, 4836–4842.
5. R. Kuroki et al. (1986) *J. Biol. Chem.* **261**, 13571–13574.

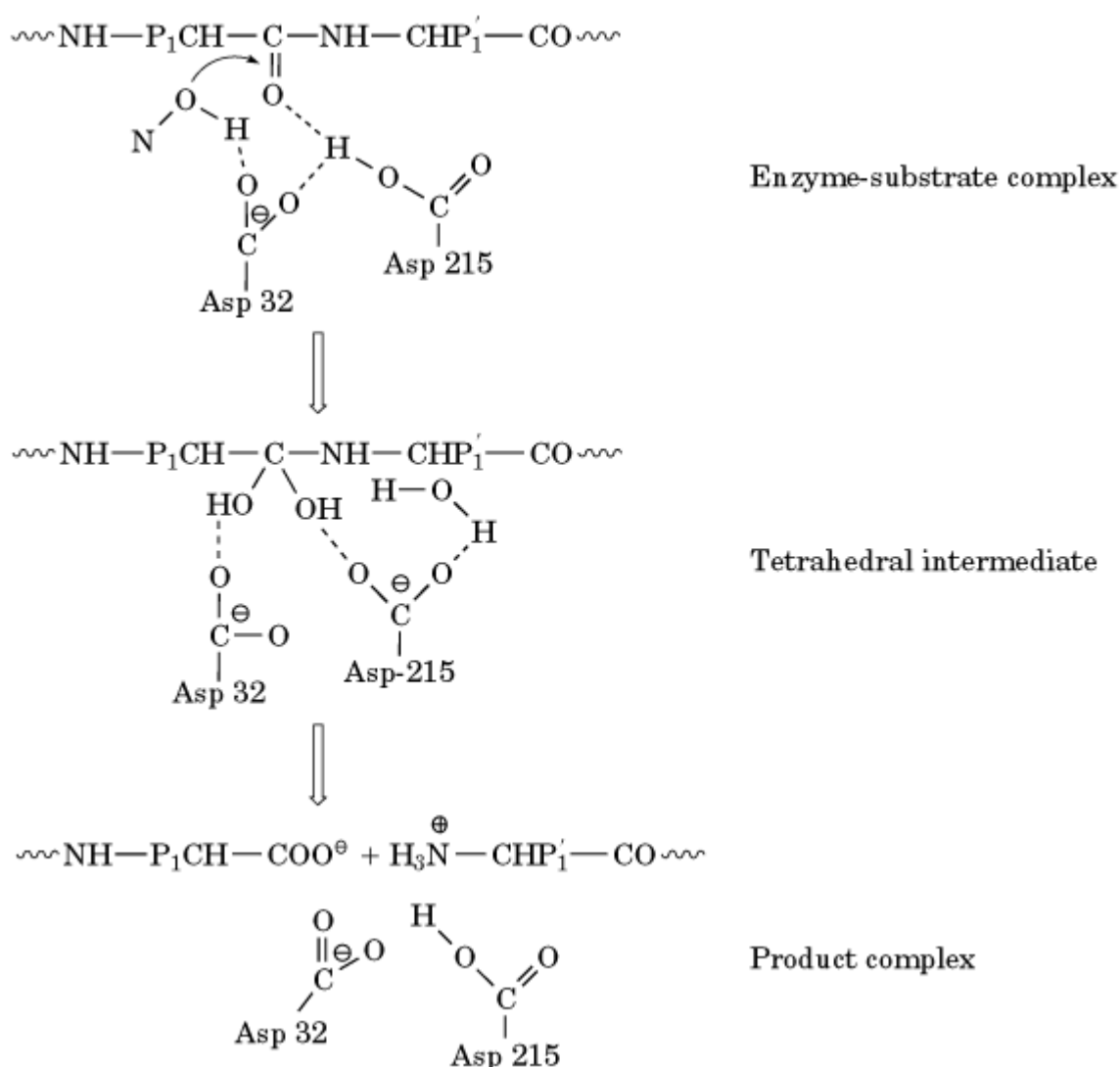
### Carboxyl Proteinase

The carboxyl proteinases (E.C. 3.4.23) are also known as *aspartyl proteinases* or *acid proteinases*, so called because the mechanism by which they catalyze the hydrolysis of [peptide bonds](#) in [proteins](#) (see [Proteinases](#)) involves the participation of two [carboxyl groups](#), usually provided by the side chains of two [aspartic acid](#) residues (1). The prototypical member of this large family of structurally-related enzymes is the gastric proteinase, **pepsin**, which is optimally active at acid pH. Other members include *rennin* (*chymosin*) from the fourth stomach of the calf; *renin*, a kidney enzyme found in blood plasma; [cathepsin D](#), found in **lysosomes**; numerous fungal enzymes, such as penicillopepsin; and a number of retroviral proteinases from [retroviruses](#), including the human immunodeficiency virus [HIV](#) proteinase that is the target of the proteinase inhibitors used in the treatment of acquired immune deficiency syndrome (AIDS). Pepsin and renin illustrate the range of specificity exhibited by proteinases in general and carboxyl proteinases in particular. Pepsin cleaves multiple bonds in virtually all proteins found in the diet, whereas renin cleaves a single bond in its only known biological substrate, *angiotensinogen*—the precursor of *angiotensin*.

Carboxyl proteinases occur in all **eukaryotic** organisms and are important for protein processing and degradation. They are generally synthesized as inactive precursors that spontaneously activate under acidic conditions. Protein and peptide substrates bind to the [active site](#) of these enzymes, and one of the carboxylate groups (in the case of pepsin, it is from Asp32) activates a water molecule to attack

the carboxyl group of the peptide bond that is to be cleaved (Figure 1). The other carboxyl group, which is protonated (in pepsin, it is Asp215), donates a proton to the peptide nitrogen atom to form an activated intermediate that dissociates into the products. It is believed that this general mechanism pertains to all of the carboxyl proteinases.

**Figure 1.** Schematic representation of the mechanism of action of carboxyl proteinases. The substrate (top line) binds to the active site of the enzyme [represented by the carboxyl groups of Asp32 and Asp215 (porcine pepsin numbering)] to form an enzyme-substrate complex. P<sub>1</sub> and P'<sub>1</sub> are the sidechains of the main specificity-determining amino acid residues that contribute the -CO- and -NH- groups to the peptide bond that will be hydrolyzed. The carboxylate group of Asp 32 promotes the attack of a water molecule on the carbonyl carbon of the peptide bond, to form a tetrahedral intermediate that is stabilized by transfer of a proton from Asp 215. The intermediate undergoes rearrangement to form a product complex, which subsequently dissociates, thereby regenerating the original enzyme.



It is interesting that in carboxyl proteinases such as pepsin and renin, both catalytic carboxyl groups are present in the same molecule. In the retroviral proteinase from HIV-1, which only contains 99 amino acids, two molecules form a dimer with a single active site containing one catalytic aspartic acid from each subunit. The larger mammalian and microbial carboxyl proteinases contain approximately 325 residues and appear to have **homologous** amino- and carboxyl-terminal **domains** that likely evolved from an early [gene duplication](#) and **fusion** event. Each domain contributes a catalytically active aspartic acid residue.

A characteristic feature of carboxyl proteinases is their susceptibility to inhibition by [pepstatin](#), an acylated pentapeptide analogue produced by **Streptomyces**. Recently a number of carboxyl proteinase inhibitors have been found useful for the treatment of AIDS, particularly when used in conjunction with other viral enzyme inhibitors. Carboxyl proteinase inhibitors have also been proposed for the treatment of malaria. Two such proteinases are believed to be essential for degradation of [hemoglobin](#) in the red blood cell phase of the life cycle of the parasite.

#### Bibliography

1. D. R. Davies (1990) *Annu. Rev. Biophys. Biophys. Chem.* **19**, 189–215.

### Carboxypeptidases

[Enzymes](#) that catalyze the hydrolytic cleavage of the carboxyl-terminal [amino acid](#) residue from an oligo- or [polypeptide chain](#) are called carboxypeptidases.

One group of these enzymes uses a catalytic mechanism analogous to that of the [serine proteinases](#) and, in fact, are inhibited by the prototypical serine proteinase inhibitor, diisopropylfluorophosphate ([DIFP](#)). These enzymes are found in the **vacuoles** of higher **plants** and **fungi** and in the **lysosomes** of animal cells, and they presumably participate in intracellular **protein degradation**. Many fungi also secrete serine carboxypeptidases. *Carboxypeptidase Y* (E.C. 3.4.16.1) is a serine carboxypeptidase present in yeast vacuoles ([1](#)). It is a very useful enzyme for determining the amino acid sequence of peptides and proteins (see [Protein Sequencing](#)), because it can release all C-terminal amino acids, including [proline](#). Moreover, it is active in the presence of **denaturants**, such as [urea](#) or sodium dodecyl sulfate (**SDS**), and in the pH range from 4 to 7.

The other group contains metallo-carboxypeptidases, which have a zinc ion at the catalytic site ([2](#)). The best known of these are the pancreatic enzymes, carboxypeptidases A and B, which are synthesized as largely inactive [zymogens](#), procarboxypeptidases A and B, and stored in the zymogen granules of the pancreatic acinar cells. On ingestion of a meal, they are released into the duodenum through the pancreatic duct, become activated by the action of trypsin, and help digest dietary proteins and peptides. The specificity of carboxypeptidase A (E.C. 3.4.17.1) is toward aromatic and bulky **hydrophobic** amino acids, nicely complementary to that of **chymotrypsin**, which generates peptides with just such amino acids at their C-termini. Similarly, the specificity of carboxypeptidase B (E.C. 3.4.17.2) is toward peptides with the C-terminal basic amino acids [arginine](#) and [lysine](#), which can be generated by the action of **trypsin**. The combined action of these endo- and exopeptidases ensures optimum formation of essential amino acids. Carboxypeptidase N circulates in plasma and is a peptidyl-L-arginine hydrolase that is thought to be responsible for the degradation of bradykinin and other [hormones](#).

An important dipeptidyl carboxypeptidase, *angiotensin converting enzyme* (ACE), cleaves C-terminal dipeptides from a variety of substrates, most notably angiotensin I and bradykinin. The former reaction generates angiotensin II, a potent vasopressor, and the latter inactivates a vasodilator. This dual effect on blood pressure has led to the widespread use of ACE inhibitors as antihypertensive agents.

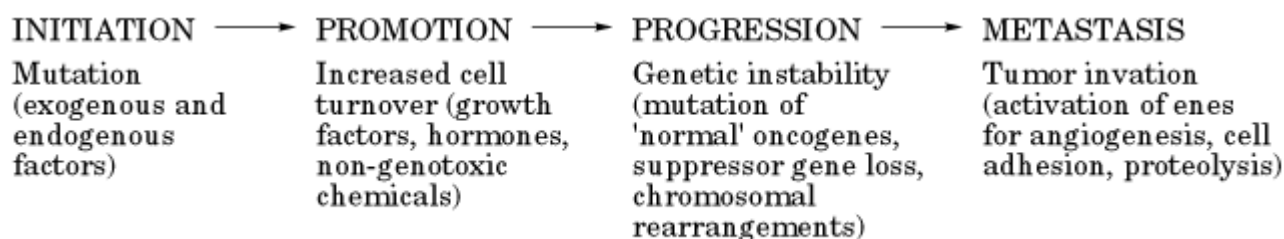
#### Bibliography

1. K. Breddam (1986) *Carlsberg Res. Commun.* **51**, 83–128.

## Carcinogen

Translated literally, the term “carcinogen” means giving rise to carcinomas, or epithelial malignancies. In practice, “carcinogen” is used to describe any physical or chemical agent that increases the incidence of tumors in animal models or in human populations. Chemical carcinogens are a diverse group of agents with a range of properties. Much of the early work on carcinogenesis in animal models identified two stages: tumor initiation, thought to involve a mutational event, and tumor promotion, which probably did not. Although this classification of events is still useful in animal studies, it is now recognized as an oversimplification. Cancer develops through a series of stages, which have been exceptionally well-characterized in certain cancers such as those of the colon (1). (An historic perspective on the development of understanding the nature of carcinogenesis is given by Harris (2)). Many of the steps in carcinogenesis involve either [mutations](#) or [epigenetic](#) changes, with cells gaining a selective advantage and undergoing clonal expansion as a result of activation of **protooncogenes** and/or inactivation of [tumor suppressor genes](#). “Epigenetic” is defined as “all processes relating to the expression (transcription and translation) and the interaction of the genetic material. Epigenetic mechanisms may act at three levels of cell organization: (1) Turn genes off or modulate **protein biosynthesis**, (2) regulate the [translation](#) of [messenger RNA](#) into proteins, and (3) regulate the topographic distribution and function of proteins.” (3). Metastatic spread of the cancer involves the activation of further oncogenes controlling various aspects of cell adhesion and movement. A scheme showing some of the events involved in various stages of cancer development is shown in Figure 1.

**Figure 1.** Factors involved in the various stages of carcinogenesis.



Mutagenicity is the major mechanism for the activation of proto-oncogenes, and most carcinogens are also [mutagens](#). Miller and Miller (4) claimed that the property in common to all the diverse carcinogens is that they either interact directly with DNA or can be metabolically activated to nucleophilic intermediates that are reactive with DNA. While this is certainly true for most carcinogens, there is increasing recognition that not all carcinogens are DNA reactive. For example, inhibitors of topoisomerase II enzymes have been shown to induce human cancers (5), and these act to cause DNA breaks indirectly by poisoning the enzyme (see [DNA Topology](#)). A number of carcinogens may act through epigenetic mechanisms. There is also a class of carcinogens that have no action at all on the genetic material (eg, **peroxisome** proliferators, which are thought to work through **receptor** binding and modification of fatty acid metabolism (6)).



While there is no doubt that mutagenic mechanisms are involved in carcinogenesis, there is considerable doubt as to the relative roles of exogenous and endogenous mutagens in the development of human cancers. Exogenous mutagens include alkylating chemicals such as nitrogen mustards and [nitrosamines](#), usually encountered in cancer therapy, and industrial chemicals such as vinyl chloride, [dimethyl sulfate](#), and bis(chloromethyl) ether. Prominent endogenous mutagenic processes include oxy-radical DNA damage, DNA depurination, **DNA polymerase** infidelity, and deamination of [5-methylcytosine](#).

### Bibliography

1. B. Vogelstein, E. R. Fearon, S. R. Hamilton, S. E. Kern et al. (1988) *N. Engl. J. Med.* **319**, 525–532.
2. C. C. Harris (1991) *Cancer Res.* **51** (Suppl.), 5023s–5044s.
3. R. Reiger, A. Michaelis, and M. M. Green (1991) *Glossary of Genetics: Classical and Molecular*. Springer-Verlag, New York.
4. E. C. Miller and J. A. Miller (1966) *Pharmacol. Rev.* **18**, 805.
5. L. R. Ferguson and B. C. Baguley (1994) *Environ. Mol. Mutagen.* **24**, 245–261.
6. S. Green (1995) *Mutat. Res.* **333**, 101–109.

### Casein

Several kinds of phosphoproteins collectively called casein constitute 80% of the [protein](#) of cow's milk and 40% of that of humans. Casein is made up of all 20 kinds of [amino acids](#), which is an important characteristic from a nutritional point of view. Skim milk, which is the milk left after the cream is removed, is brought to pH 4.7 to form a precipitate that contains the casein. The supernatant is called whey and contains 20% of the total protein. [Electrophoresis](#) at alkaline pH separates casein into the a, b, and g fractions. The a fraction is further separated into  $a_{s1}$  and  $a_{s2}$ , the common subscript *s* signifying sensitivity to precipitation by the addition of calcium and k-casein. Both  $a_s$ -casein and b-casein are precipitated with  $Ca^{++}$ , whereas k-casein is not. In milk, these caseins constitute micelles (calcium caseinate-phosphate complex) of a diameter in the range of 30 to 300 nm. The a- and b- caseins are rich in phosphate, mainly in the form of O-phosphoserine residues, and they form a calcium complex that is not soluble, but precipitated. When surrounded by the [hydrophilic](#) parts of k-casein in milk, however, it is soluble as large complexes. The g-casein fraction is a heterogeneous mixture primarily of proteolytic fragments of b-casein.

Bovine b-casein contains 209 amino acid residues (with a molecular weight of 23,600) and is very hydrophobic. Five phosphoserine residues are found in the N-terminal region. In milk, the phosphoserine residues are associated with  $Ca^{++}$  ions, forming calcium caseinate. b-casein and its partial hydrolysis products present in milk do not have definite three-dimensional structures. They are considered to have a [random coil](#) structure that is easily digestible, an advantage of casein as a nutritional protein. The relative casein composition of bulk milk is given in Table [1](#).

**Table 1. Average Casein Composition of Bovine Bulk Milk Samples**

---

### Casein Type Weight %

---

|                 |       |
|-----------------|-------|
| a <sub>s1</sub> | 38.1  |
| a <sub>s2</sub> | 10.2  |
| b               | 35.7  |
| k               | 12.8  |
| g               | 3.2   |
| Total           | 100.0 |

---

The fourth stomach of ruminating animals (abomasum) contains a proteolytic enzyme called *chymosin*, which liberates a glycopeptide from k-casein at neutral pH. The remaining molecule, called para-casein, reacts with Ca<sup>++</sup> and then forms an insoluble curd, the source of cheese. The enzyme preferentially cleaves the peptide bond between Phe105-Met106 of k-casein to produce para-k-casein. Humans do not have chymosin, but **pepsin** in the stomach converts k-casein to para-casein to start the digestion of milk in the infant stomach. Chymosin and pepsin are structurally related.

The casein micelle has a complex structure, and the relative disposition of the four subspecies of casein, and especially that of k-casein, has long been the focus of discussion. Since chymosin preferentially digests k-casein, this casein is thought to be on the periphery of the micelle. Recent **electron micrographs** of variously treated micelles show granular structures covering the surface, prompting a suggestion that a casein micelle is composed of smaller-sized submicelles. The idea of a submicellar structure has actually been a long standing one from reconstitution experiments, and by incorporating various observations, a currently accepted model of casein micelles is presented by Holt (1), as reproduced in Figure 1.



[Secretory vesicles](#) of casein micelles are found abundantly in mammary secretory cells. After the casein proteins are synthesized, they undergo [O-glycosylation](#) and **phosphorylation** before being secreted. These changes occur in concert with lactose synthesis by the enzyme known as lactose synthetase. All the enzymes involved in glycosylation, phosphorylation, and lactose synthesis can tolerate  $\text{Ca}^{++}$  concentrations in the range of 1 *mM*. Accumulation of  $\text{Ca}^{++}$  in the mammary gland cell, together with phosphorylation of some serine residues, is necessary for the formation of casein micelles.

#### Bibliography

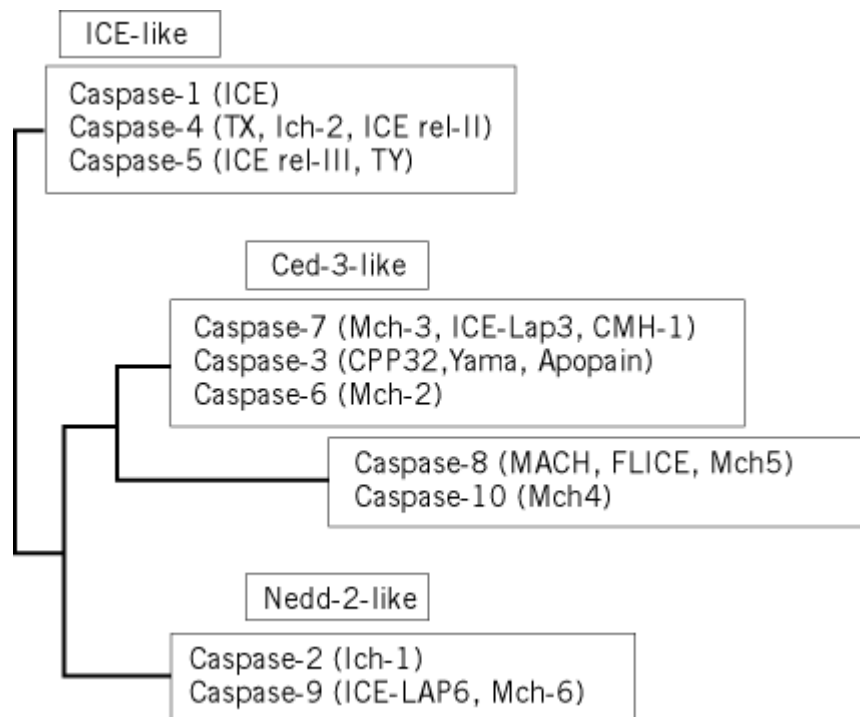
1. C. Holt (1992) Structure and stability of bovine casein micelles, *Adv. Prot. Chem.* **43**, 63–151.

#### Caspases

*Caspases* are [proteases](#) involved in [apoptosis](#) and [programmed cell death](#). They are mammalian **homologues** of the Ced-3/interleukin-1 $\beta$  converting enzyme (ICE) products of [cell death](#) genes ([1](#)). The caspase family in mammalian cells currently numbers 14, as shown in Figure [1](#). They are [thiol proteases](#) that are synthesized as inactive [pro-proteins](#) that are processed to a pro-domain and two subunits of approximately 10 and 20 kDa in size. The pro-caspases are inert until proteolytic cleavage of the pro-domain and/or the p10 subunit from the p20 subunit. Caspases share a conserved [active-site](#) sequence, -Gln-Ala-Cys-Arg-Gly- (QACRG) or -Gln-Ala-Cys-Gln-Gly- (QACQG), and

cleave their substrates at the peptide bond after an aspartate residue; hence, their name- **cysteine aspartate-specific proteinases** (1, 2). Caspases cleave specific substrates during **apoptosis** and are required for the rapid degradation of the cell (2-6); however, none is thought to be critical for apoptosis to occur.

**Figure 1.** The mammalian caspase family. The various members have been arranged in a phylogeny on the basis of their relative amino acid sequence homologies. The members of the Ced-3-like family of caspases are thought to be the downstream effectors of apoptosis and to be activated in a caspase cascade by the apical members represented by caspases 2, 8, 9, and 10 (1).



Caspases can be divided into subgroups, with each sharing a specific substrate specificity beyond the overall caspase requirement for Asp at the P1 position of the substrate (see **Proteinases**). Caspases 1, 4, and 5, for example, preferentially cleave after -Tyr-Val-Ala-Asp- (YVAD) sequences and have not been shown to be critical for apoptosis to occur. Caspases 3 and 6, however, are active during cell death and cleave primarily substrates with the sequence -Asp-Xaa-Xaa-Asp- (DxxD). Caspases 2, 8, and 10 are different again and are thought to activate primarily the downstream caspases, such as caspase 3 and 6 (7). This suggests that caspases act in a hierarchy, with activated upstream caspases cleaving downstream caspases and activating them in turn, to generate a *caspase cascade* (4, 8). The apical caspases are proposed to be those with large pro-domains such as caspases 2, 8, and 10. The pro-domains contain a protein structure that facilitates their interaction with other proteins leading to caspase activation, as seen in CD95 (Fas/Apo-1)-induced apoptosis (9).

#### Bibliography

1. E. Alnemri, D. Livingston, D. Nicholson, G. Salvesan, N. Thornberry, W. Wong, and Yuan, J. (1996) Human ICE/CED-3 protease nomenclature. *Cell* **87**, 171.
2. A. Takahashi and W. Earnshaw (1996) Ice-related proteases in apoptosis. *Curr. Opin. Gen. Dev.* **6**, 50–55.
3. A. Chinnaiyan and V. Dixit (1996) The cell-death machine. *Curr. Biol.* **6**, 555–562.
4. Y. Lazebnik, A. Takahashi, G. Poirier, S. H. Kaufmann, and W. Earnshaw (1995)

Characterization of the execution phase of apoptosis *in vitro* using extracts from condemned phase cells. *J. Cell. Sci.* **19**, 41–49.

5. M. Whyte (1996) ICE/Ced-3 proteases in apoptosis. *Trends Cell Biol.* **6**, 245–248.
6. A. Fraser, N. McCarthy, and G. I. Evan (1996) Biochemistry of cell death. *Curr. Opin. Neurobiol.* **6**, 71–80.
7. T. Fernandes Alnemri, R. C. Armstrong, J. Krebs, S. M. Srinivasula, L. Wang, F. Bullrich, L. C. Fritz, J. A. Trapani, K. J. Tomaselli, G. Litwack, and E. S. Alnemri (1996) *In vitro* activation of CPP32 and Mch-3 by Mch-4, a novel human apoptotic cysteine protease containing 2 FADD-like domains. *Proc. Natl. Acad. Sci. USA* **93**, 7464–7469.
8. Y. A. Lazebnik, A. Takahashi, R. D. Moir R. D. Goldman, G. G. Poirier, S. H. Kaufmann, and W. C. Earnshaw (1995) Studies of the lamin proteinase reveal multiple parallel biochemical pathways during apoptotic execution. *Proc. Natl. Acad. Sci. USA* **92**, 9042–9046.
9. A. Fraser and G. Evan (1996) A license to kill. *Cell* **85**, 781–784.

## Cassette Mutagenesis

Cassette mutagenesis involves replacing a wild-type **DNA** sequence with synthetic double-stranded oligonucleotides in order to introduce one or more [mutations](#) (1-6). The technique allows saturation of a target amino acid codon with [mutations](#), and can be used to probe the role of that particular residue in a protein and the effect of specific point mutations in it. In their description of the method, Wells et al. (4) synthesized single-stranded oligonucleotides containing different codons over the target in separate pools. Oligonucleotide-**site-directed mutagenesis** procedures are used to generate **restriction sites** that closely flank the target codon in the **plasmid** containing the gene of interest. Use of the appropriate restriction endonucleases permits insertion of the synthetic duplex cassettes. These are designed to restore fully the wild-type coding sequence except over the target codon, and also to eliminate one of the restriction sites. This latter point is important for selection of the clones containing the mutant cassette.

Several variations on the original technique have been developed for specific purposes. For example, Reidhaar-Olsen and Sauer (5) described a combinatorial method, in which they randomized two or three positions by oligonucleotide cassette mutagenesis, selected for functional protein and then sequenced to determine the spectrum of allowable substitutions at each position. They applied the method repeatedly in order to examine the role of different substitutions on the DNA-binding domain of **lambda repressor**. Kegler-Ebo et al. (6) described codon cassette mutagenesis as a way of depositing single codons at specific sites in double-stranded DNA. Using this method, a series of 11 cassettes is sufficient to insert all possible amino acids at any constructed target site.

## Bibliography

1. S. Green (1955) *Mutat. Res.* **333**, 101–109.
2. M. D. Matteucci and H. L. Heyneker (1983) *Nucleic Acids Res.* **11**, 3113–3121.
3. J. B. McNeil and M. Smith (1985) *Mol. Cell. Biol.* **5**, 3545–3551.
4. J. A. Wells, M. Vasser, and D. B. Powers (1985) *Gene* **34**, 315–323.
5. J. F. Reidhaar-Olsen and R. T. Sauer (1988) *Science* **241**, 53–57.
6. D. M. Kegler-Ebo, G. W. Polack, and D. DiMaio (1994) *Methods Mol. Biol.* **57**, 297–310.

## Catabolite Repression

The ability of glucose to inhibit the synthesis of certain [enzymes](#), referred to as the *glucose effect*, was recognized in **bacteria** at an early date by Monod (1). He observed that when *Escherichia coli* encounters both glucose and lactose, for example, it metabolizes the glucose first and represses the use of lactose, resulting in biphasic growth (diauxie). This phenomenon is due to the repressive effect of glucose on the synthesis of enzymes required for the metabolism of other sugars. Later, the term catabolite repression was introduced as a general name for the glucose effect because compounds closely related to glucose elicited varying degrees of repression of glucose-sensitive enzymes, and the catabolites derived from the repressing compound were assumed to cause the glucose effect (2). Studies on catabolite repression up to the early 1970s brought about the discovery of a positive control system of [transcription](#) by [cyclic AMP](#) (cAMP) and the cyclic AMP receptor protein (CRP) (also called the catabolite gene activator protein, (CAP)) and led to the concept that catabolite repression is caused by reduction of the level of intracellular cAMP (3). We now know that multiple and different mechanisms are operating, depending on the growth conditions and the target [operons](#), and that the cAMP-dependent mechanism is just one aspect of catabolite repression. Ironically, the preferential utilization of glucose over lactose, the prototype of catabolite repression, is not due to the reduction of the cAMP-CRP complex (4, 5).

Sometimes the term catabolite repression has been used to describe the glucose effect that is independent of the operon-specific regulator, such as the [Lac repressor](#). Here, we shall use it as a general name for the glucose effect as originally defined. According to this view, which is apparently generally accepted, catabolite repression includes all forms of the glucose effect, regardless of their mechanisms. We also concentrate on catabolite repression in *E. coli*. Although several mechanisms of glucose catabolite repression are now understood in *E. coli*, a common feature is that glucose causes catabolite repression ultimately by modulating the activity of [transcription factors](#) involved in the regulation of target operons.

### 1. cAMP–CRP Complex

The best characterized mechanism of catabolite repression involves the regulation of the intracellular concentration of the cAMP–CRP complex. It is well known that the presence of glucose in the growth medium lowers the intracellular cAMP level under certain conditions (3, 6). Cyclic AMP is synthesized from ATP by the enzyme [adenylate cyclase](#). Although the mechanism of regulation of the cAMP level remains elusive, glucose is thought to decrease cAMP by decreasing the level of the phosphorylated form of enzyme  $\text{IIA}^{\text{Glc}}$ , which is involved in the activation of adenylate cyclase.  $\text{IIA}^{\text{Glc}}$  is one of the enzymes of the phosphoenolpyruvate-dependent carbohydrate **phosphotransferase system** (PTS) and is directly responsible for the **active transport** and phosphorylation of glucose (7). Recently, it was discovered that the concentration of CRP is also lowered by the presence of glucose and that this is an additional factor contributing to catabolite repression (8). The decreased CRP is a consequence of the complex autoregulation of expression of the *crp* gene (9). It should be noted that the reduction in the cAMP-CRP level by glucose is usually rather moderate (in the range of several-fold).

### 2. Inducer Exclusion

The second mechanism of catabolite repression is inducer exclusion, by which glucose lowers the intracellular concentration of inducers necessary for the induction of catabolic operons (7). The target of glucose signaling in inducer exclusion is operon-specific regulators, such as the Lac

repressor. The dephosphorylated enzyme  $\text{IIA}^{\text{Glc}}$ , which accumulates in the presence of glucose, binds to and inactivates (for example) the Lac permease, resulting in an increase of the active unliganded Lac repressor (see [Lac Operon](#)). Inducer exclusion is a mechanism by which glucose inhibits more strictly the expression of target operons.

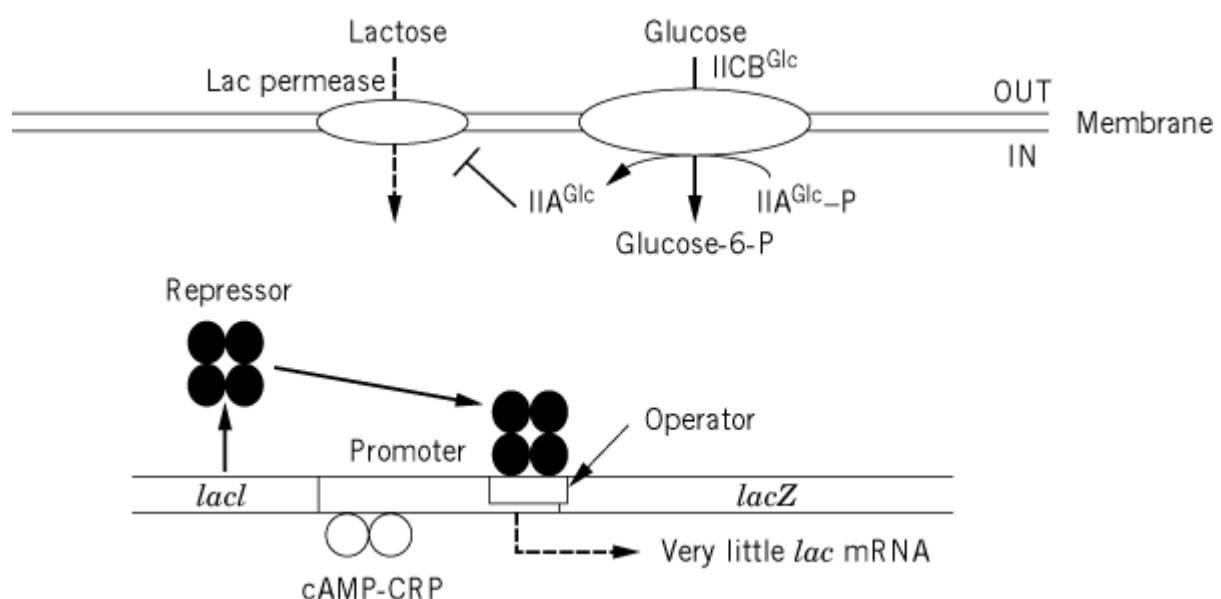
### 3. Catabolite Repressor/Activator Protein

The third mechanism of catabolite repression is mediated by the catabolite repressor/activator (Cra) protein, which acts as a global regulator of genes encoding enzymes of central carbohydrate metabolism (10). The unliganded form of Cra binds to the operator regions of target operons, causing either activation or inhibition of transcription. The presence of glucose or other PTS sugars produces glycolytic catabolites, such as fructose-1-phosphate, which bind to the Cra protein and cause it to dissociate from the target DNA, resulting in either catabolite repression or catabolite activation.

### 4. Relationships between the Various Mechanisms

While multiple mechanisms of catabolite repression have been identified in *E. coli*, their signaling pathways appear to be interrelated to each other. For example, the PTS plays a pivotal role in the regulation of the intracellular concentrations of cAMP, inducer, and glycolytic catabolites. In addition, it is particularly important to realize that the contribution of each mechanism varies, depending on growth conditions and the target genes. For example, the Cra-mediated mechanism may play no role in catabolite repression of the *lac* operon, because this operon is not under the control of Cra. An unexpected finding is that the presence of glucose in the lactose medium does not affect the intracellular cAMP level (4). This means that catabolite repression mediated by the reduction in cAMP never happens in glucose–lactose diauxie. The presence of unliganded Lac repressor through inducer exclusion is the principal mechanism for this historical phenomenon (4, 5). Figure 1 illustrates the present understanding of catabolite repression in the glucose–lactose system.

**Figure 1.** Mechanism of catabolite repression in the glucose-lactose system (4). When both lactose and glucose are present, glucose is transported and phosphorylated by the glucose PTS ( $\text{IIA}^{\text{Glc}} + \text{IICB}^{\text{Glc}}$ ), increasing the concentration of the nonphosphorylated form of  $\text{IIA}^{\text{Glc}}$ , which prevents the uptake of lactose by inhibiting the Lac permease activity. Thus, the concentration of *lac* inducer is very low in the presence of glucose, so the Lac repressor is active and represses transcription of the *lac* operon. It should be noted that glucose does not affect the binding of cAMP–CRP to the promoter, because the levels of cAMP and CRP are not reduced by the presence of glucose.



There are more stories yet to be elucidated before catabolite repression in *E. coli* is fully understood. Further diversity in the mechanism of catabolite repression is known in **Gram-positive** bacteria (11).

### Bibliography

1. J. Monod (1947) Growth **11**, 223–289.
2. B. Magasanik (1961) Cold Spring Harbor Symp. Quant. Biol. **26**, 249–256.
3. I. Pastan and R. Perlman (1970) Science **169**, 339–344.
4. T. Inada, K. Kimata, and H. Aiba (1996) Genes to Cells **1**, 293–301.
5. K. Kimata, H. Takahashi, T. Inada, P. Postma, and H. Aiba (1997) Proc. Natl. Acad. Sci. USA **94**, 12914–12919.
6. R. S. Makman and E. W. Sutherland (1965) J. Biol. Chem. **240**, 1309–1314.
7. P. W. Postma, J. W. Lengeler, and G. R. Jacobson (1993) Microbiol. Rev. **57**, 543–594.
8. H. Ishizuka, A. Hanamura, T. Kunimura, and H. Aiba (1993) Mol. Microbiol. **10**, 341–350.
9. H. Ishizuka, A. Hanamura, T. Inada, and H. Aiba (1994) EMBO J. **13**, 3077–3082.
10. M. H. Saier and T. M. Ramseier (1996) J. Bacteriol. **178**, 3411–3417.
11. C. J. Hueck and W. Hillen (1995) Mol. Microbiol. **15**, 395–401.

### Suggestions for Further Reading

12. B. Magasanik (1970) "Glucose effects: inducer exclusion and repression", in *The Lactose Operon*, J. Beckwith and D. Zipser, eds., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 189–220.
13. A. Ullmann and A. Danchin (1983) Role of cAMP in bacteria, Adv. Cyclic Nucleotide Res. **15**, 1–53.
14. N. D. Meadow, D. K. Fox, and S. Roseman (1990) The bacterial phosphoenolpyruvate:glycose phosphotransferase system. Annu. Rev. Biochem. **59**, 497–542.
15. M. H. Saier Jr., T. M. Ramseier, and J. Reizer (1996) "Regulation of carbon utilization", in *Escherichia coli and Salmonella typhimurium: Cellular and Molecular Biology*, F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, J. K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds., American Society for Microbiology, Washington, DC, pp. 1325–1343.

### Catalase

Catalase, also called *hydroperoxidase*, is a protective [enzyme](#) that has been studied for over a century; the concept of a specific protein catalyzing (hence the name *catalase*) the degradation of hydrogen peroxide to oxygen and water first appeared in 1900 (1). Removal of H<sub>2</sub>O<sub>2</sub>, as shown in the following reaction, avoids unwanted side reactions and prevents the formation of the even more reactive hydroxyl radical:





Although the removal of these reactive oxygen species is not essential for growth, catalases do enhance long-term survival in an aerobic environment (2).

A number of developments have spawned considerable interest in the enzyme and spurred extensive studies of both the enzyme and its genes. The first was the simplicity of the “drop test” assay, which involves the visual monitoring of oxygen evolution after the application of a drop of 30% H<sub>2</sub>O<sub>2</sub> to the edge of a bacterial colony. Such ease of **phenotypic** scoring has resulted in extensive use of the enzyme as a diagnostic tool for microbiological strain identification. Catalases have been found to be synthesized as part of a variety of **stress response** systems, resulting in their frequent use as indicators of stress response activation. More recently, the catalase–peroxidase family has gained notoriety from its role as the *in vivo* activator of isoniazid into an effective antibiotic in *Mycobacterium tuberculosis* (3). Loss of the enzyme through [mutation](#) results in isoniazid resistance, one of the reasons for the increasing spread of tuberculosis.

### 1. Assay

Two quantitative assay methods are commonly used for catalase activity. One involves the measurement of oxygen evolution using an oxygraph equipped with a Clark electrode (4); the second is a spectrophotometric assay of H<sub>2</sub>O<sub>2</sub> by its **absorbance** at 240 nm (5). The two assay procedures produce comparable values in reasonably clear protein solutions, but the oxygraph protocol has the advantage that it can be used for catalase determinations in whole-cell suspensions and in quite turbid extracts. Catalase activity can be visualized on [polyacrylamide](#) gels following [electrophoresis](#) under nondenaturing conditions (6), producing a clear band in a brown background. This same visualization procedure can be easily modified to visualize [peroxidase](#) activities as brown bands on a clear background.

### 2. Classification

Because of the protection afforded against active oxygen species by catalases, most aerobic organisms produce at least one catalase from among the three main classes of the enzyme. The largest class includes the “typical,” heme-containing, monofunctional catalases with either small (~60 kDa) or large (>80 kDa) subunits. The next largest class includes the bifunctional catalase–peroxidases, which are also heme-containing but with sequence similarity to plant and fungal peroxidases. The third and, to date, smallest class includes the non-heme- or Mn-containing catalases. Small-subunit, monofunctional catalases have been found in prokaryotes and eukaryotes but not in an archaeal species, and the large-subunit catalases are restricted to fungi and bacteria. The catalase-peroxidases are found in **prokaryotes** and **archaeobacteria**, and the non-heme-containing enzymes have been found only in bacteria (7).

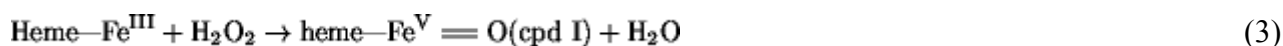
### 3. Reaction

Catalases and peroxidases both degrade H<sub>2</sub>O<sub>2</sub> but, whereas catalase employs H<sub>2</sub>O<sub>2</sub> as both electron donor and acceptor (reaction 1), a peroxidase uses organic substrates as the electron donor to reduce H<sub>2</sub>O<sub>2</sub>:

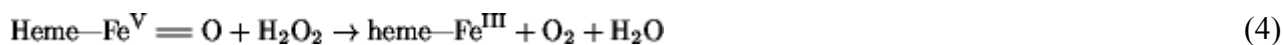


The catalase-peroxidases catalyze both reactions 1 and 2 at significant rates.

The well-characterized catalytic reaction pathway is a two-step process. First, H<sub>2</sub>O<sub>2</sub> is converted to water and an oxyferryl species, compound I, with the iron in the +5 oxidation state but with part of the charge delocalized in a heme cation radical:



Then compound I reacts with a second molecule of  $\text{H}_2\text{O}_2$  to produce water and molecular oxygen:

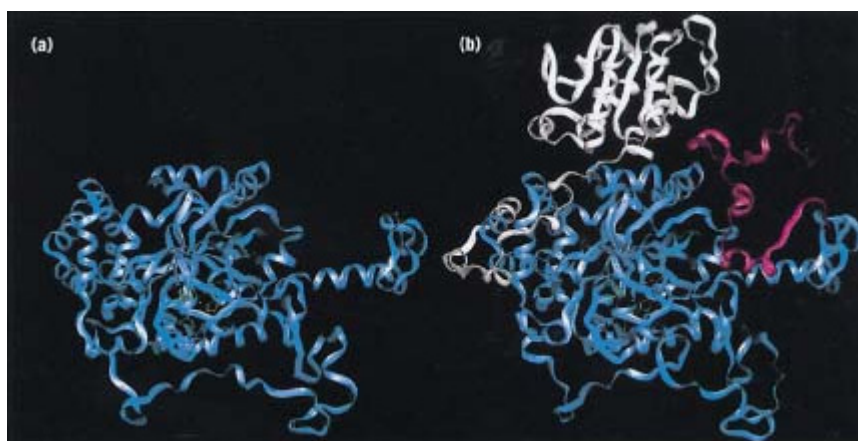


In the absence of a suitable substrate for the reduction of compound I, other intermediates can be formed such as the catalytically inactive compound II.

#### 4. Structure and Properties

Detailed structural information is now available for all three types of catalases. Small- and large-subunit catalases share a common highly conserved core sequence of about 350 residues organized in a b-barrel structure (see [Beta-Sheet](#)). This is illustrated in Figure 1 (8, 9), which compares the small subunit of *Proteus mirabilis* catalase (Fig. 1a) and the large subunit of *Escherichia coli* HP11 (Fig. 1b). An *N*-terminal extension of 10–90 residues (colored red in Fig. 1B) and a *C*-terminal extension of 150–170 residues (colored white in Fig. 1b), the latter organized in a flavodoxin-like structure, are added to the core structure in the large-subunit enzyme. The tetrameric organization of small-subunit enzymes is sufficiently stable to resist denaturation by organic reagents and to retain activity over a broad pH range from 4 to 11. In large-subunit enzymes, the extended sequences allow the amino terminus of one subunit to be overlapped or trapped by the carboxyl terminus of an adjacent subunit (10), and this interweaving further enhances the resistance of the tetramer to denaturation by heat, SDS, and [urea](#).

**Figure 1.** Comparison of the structures of the small subunit of *Proteus mirabilis* catalase (a) (8) and the large subunit of *E. coli* HP11 catalase (b) (9). The conserved b-barrel core structure is colored blue in both subunits. The 78 additional residues at the *N* terminus of HP11 in (b) are colored red, and the 195 additional residues at the *C* terminus of HP11 in (b) are colored white. See color insert.



Besides size, the small- and large-subunit catalases have two other characteristic differences. The first is that small-subunit enzymes contain heme *b*, which in large-subunit enzymes is converted to a *cis*-hydroxy *g*-spirolactone heme *d* and flipped 180°. The second is that small-subunit, but not large-subunit, catalases contain NADPH (11), which is postulated to have a role in the reduction of inactive compound II. A number of unusual modifications have been identified in catalases, including

1. A [methionine](#) sulfone in the active site of the *P. mirabilis* enzyme (12)

2. A blocked [cysteine](#) residue in HP11 from *E. coli* (13)
3. A [histidine-tyrosine](#) bond in HP11

The narrow channel leading to the deeply buried [active site](#) of catalases limits access to small substrates and inhibitors, and results in a 10-fold higher apparent  $K_m$  for  $H_2O_2$  compared to the non-heme-containing catalases and catalase–peroxidases. Within the catalase active site, a histidine residue immediately above the heme, and an asparagine residue situated to one side, orient the  $H_2O_2$  over the heme iron (12) to facilitate a very fast reaction in both small- and large-subunit enzymes [ $k_{cat}$  @  $3.5 \times 10^5 s^{-1}$  (15) and  $1.2 \times 10^5 s^{-1}$ , respectively] as compared to the catalase–peroxidases and non-heme-containing catalases [ $k_{cat}$  @  $1.6 \times 10^4 s^{-1}$  (16) and  $3.9 \times 10^5 s^{-1}$  (17), respectively].

## 5. Cloning and Expression

The number of catalase enzymes in a given organism varies: one in animals, one or two in **fungi**, one to three in plants, and zero to three in bacteria. The availability of catalase-deficient hosts, the ease of assay, and the development of protocols using oligonucleotide and **PCR** probes have facilitated the [cloning](#) and characterization of over 90 catalase genes. Most of the isolated genes are from plants and bacteria, with a small number from fungi and animals. A **phylogenetic** comparison of the core protein sequences has revealed separate branches for the small-subunit enzymes from bacterial, plant, animal, and fungal sources, and one additional branch containing the large-subunit enzymes from both fungal and bacterial origins.

The most detailed studies of catalase gene expression have been carried out in plants and bacteria. In maize, catalase expression is tied to developmental changes during kernel development and seed germination and to environmental factors such as light and growth regulators (18). A multitude of species-specific mechanisms control catalase and catalase–peroxidase expression in bacteria, but two common threads are evident among these expression patterns involving induction either by an active oxygen species such as  $H_2O_2$  or in stationary phase. In most cases, the control mechanisms can be rationalized as either protective responses to oxidative stress or as a means of enhancing long-term survival during metabolic limitation (7, 19).

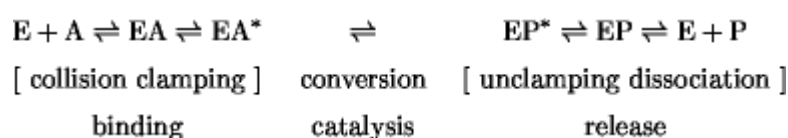
## Bibliography

1. O. Loew (1900) U.S. Dept. Agriculture Rpt. 65, Govt. Printing Office, Washington, DC.
2. M. R. Mulvey, J. Switala, A. Borys, and P. C. Loewen (1990) *J. Bacteriol.* **172**, 6713–6720.
3. J. S. Blanchard (1996) *Annu. Rev. Biochem.* **65**, 215–239.
4. M. Rorth and P. K. Jensen. (1967) *Biochim. Biophys. Acta* **139**, 171–173.
5. A. G. Hildebrandt and I. Roots (1975) *Arch. Biochem. Biophys.* **171**, 385–397.
6. E. M. Gregory and I. Fridovich. (1974) *Anal. Biochem.* **58**, 57–62.
7. P. C. Loewen (1997) In J. G. Scandalios, ed., *Oxidative Stress and the Molecular Biology of Antioxidant Defenses*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 273–308.
8. P. Gouet, H. M. Jouve, and O. Dideberg (1995) *J. Mol. Biol.* **249**, 933–954.
9. J. Bravo, N. Verdager, J. Tormo, C. Betzel, J. Switala, P. C. Loewen, and I. Fita (1995) *Structure* **3**, 491–502.
10. W. R. Melik-Adamyany, V. V. Barynin, A. A. Vagin, V. V. Borisov, B. K. Vainshtein, I. Fita, M. R. N. Murthy, and M. G. Rossmann (1986) *J. Mol. Biol.* **188**, 63–72.
11. H. N. Kirkman and G. F. Gaetani (1984) *Proc. Natl. Acad. Sci. USA* **81**, 4343–4347.
12. A. Buzy, V. Bracchi, R. Sterjiades, J. Chroboczek, P. Thibault, J. Gagnon, H.-M. Jouve, and G.

- Hudry-Clergeon (1995) *J. Protein Chem.* **14**, 59–72.
13. S. Sevinc, W. Ens, and P. C. Loewen (1995) *Eur. J. Biochem.* **230**, 127–132.
  14. I. Fita and M. G. Rossmann (1985) *J. Mol. Biol.* **185**, 21–37.
  15. A. Deisseroth and A. L. Dounce (1970) *Physiol. Rev.* **50**, 319–375.
  16. A. Claiborne and I. Fridovich (1979) *J. Biol. Chem.* **254**, 4245–4252.
  17. Y. Kono and I. Fridovich (1983) *J. Biol. Chem.* **258**, 6015–6019.
  18. J. G. Scandalios (1992) in J. G. Scandalios, ed., *Molecular Biology of Free Radical Scavenging Systems*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 117–152.
  19. H. Schellhorn (1994) *FEMS Microbiol. Lett.* **131**, 113–119.

## Catalysis

Catalysis by [enzymes](#) starts with the collision between the enzyme  $E$  and its substrate  $A$  to form an enzyme-substrate complex  $EA$ :



This is a second-order reaction, as it involves two reactants, and the collision rate is equal to  $k[E_i][A_i]$ , where  $k$  is a second-order rate constant. There is a distinct limit to the value of  $k$ , as it cannot exceed the [diffusion](#) rate (see [Diffusion-Controlled Reactions](#) and [Turnover Number](#)).

Nonproductive complexes can also form as a result of the interaction of a substrate with an enzyme. If the structure of a substrate is complementary to that of the [active site](#) of the enzyme, the interaction rate is high and the binding tight. If, however, this is not the case, some of the binding energy has to be utilized to facilitate the interaction through a conformational change or clamping reaction on the enzyme. Binding is the product of the initial interaction and the subsequent conformational change.

After being bound to the enzyme, the substrate must undergo activation. That is, a proportion of the total number of substrate molecules in the ground state must be raised to the [transition state](#). This is a short-lived, unstable species with a structure intermediate between that of the substrate and the product. These activated molecules can fall back to the ground state or be converted to  $EP^*$  at a rate that is independent of the structure of the reactants. The [activation energy](#) barrier for the conversion of ground-state molecules to transition-state molecules is very much lower for an enzyme-catalyzed reaction than for either a noncatalyzed reaction or chemically catalyzed reaction. In the activated state, the substrate is bound more tightly to the enzyme than it is in the ground state (see [Transition State Analogue](#)). If such tighter binding did not occur, there would be no catalysis. In fact, the increase in rate is proportional to the increase in binding. It is for this reason that transition state analogues are potent inhibitors of enzymes. After formation of the  $EP^*$  complex, the enzyme undergoes another conformational or unclamping reaction to release the product into the bulk medium and regenerate the free form of enzyme. Such regeneration is an essential part of enzymic catalysis.

## Bibliography

1. D. D. Hackney (1990) *The Enzymes* **19**, 1–36.

## Catalytic Antibodies

The concept of *catalytic antibodies* owes its origin to Pauling and Jencks, who proposed that an [antibody](#) normally differs from an enzyme by its inability to bind selectively and to stabilize the [transition state](#) of a chemical reaction. An antibody that did happen to be specific for a transition state should therefore function like an enzyme and promote chemical catalysis of the corresponding reaction (1). Advances in chemical synthesis and in [hybridoma](#) technology to produce [monoclonal antibodies](#) have enabled the exploitation of the diversity and affinity of the immune repertoire to develop catalytic antibodies that are elicited against [transition state analogues](#). The binding sites of these antibodies are anticipated to bind substrates structurally related to the transition state and to process them to products through a pathway lower in free energy, and therefore more rapid, than the normal one that occurs in the absence of antibody. By clever design of appropriate transition state analogues, “tailor-made catalysts” should be created to catalyze reactions with no enzymic counterparts. Challenge of the immune system with a transition state analogue to induce a catalytic antibody is a necessary, but not a sufficient, condition to generate an effective catalyst; factors other than high affinity for the transition state, such as precise orientation of catalytic residues and the effective release of reaction products at the active site, also contribute to the overall catalytic efficiency of an antibody. Strategies have evolved to design transition state analogue [immunogens](#) that would solicit catalytic functions within the antibody combining site for efficient catalysis.

### 1. Reactions Catalyzed by Antibodies

There are now approximately 100 reactions that have been catalyzed by antibodies (1-3). These reactions include (1) pericyclic processes (oxy-Cope rearrangement, Diels–Alder condensation, Claisen rearrangement), (2) elimination reactions (decarboxylation, dehydration, *syn* elimination of HF from fluoroketones, E2 (bimolecular) elimination of benisoxazole), (3) hydrolyses (carbonate esters, esters, amide, lactones, enol ethers), (4) bond-forming reactions (lactonization, [peptide synthesis](#), cationic cyclization, aldol condensation), and (5) redox reactions (ketone reduction, epoxidation, sulfoxide oxidation). Antibodies can catalyze, with high stereo- and regioselectivity, reactions that may not generally be catalyzed by enzymes. This property has appeal in the potential application of antibody catalysis to commercial synthetic and medical uses, such as pharmaceutical synthesis and prodrug activation.

### 2. The Nature of the Catalytic Site

The [active sites](#) of catalytic antibodies elicited by a given immunogen exhibit high sequence homologies and are often structurally convergent (4). These sites harbor shallow clefts complementing the structural and electronic features of the immunogen. The positioning of active site residues in the binding pocket is accomplished by somatic mutation of the germline ancestral antibody arising from an immunological response (5). The germline antibody undergoes a substantial amount of [induced-fit](#) conformational change on binding the immunogen. By the process of [affinity maturation](#), the active site residues in the mature antibody become preoriented such that a rigid binding pocket is generated for optimal binding (lock-and-key fit).

### 3. Nature of Catalysis

Catalysis by antibodies is like that by enzymes, in that it exhibits [Michaelis–Menten kinetics](#) in which substrate binding precedes a chemical transformation, followed by dissociation of the product. Kinetic characterization experiments (6) indicate that the antibody-catalyzed reaction recapitulates a number of characteristics expected of an enzyme, but the chemical transformation step is often rate-limiting. The rate enhancement is generally less than that observed for enzymes catalyzing similar reactions. Transition state analogues are only approximations of the true transition state, and they also do not demonstrate the extremely tight binding to enzymes that is predicted theoretically (see [Transition State Analogue](#)). Furthermore, antibodies catalyze reactions primarily through restrictions on the translational and rotational movement of substrates, and only to a lesser extent through the active-site acid/base or nucleophilic catalysis that occurs in enzymes. The difference in catalytic efficiency between antibody and enzyme can also be attributed in part to the latter's ability to provide more extensive [electrostatic interactions](#) and [hydrogen bonding](#) during catalysis through conformational mobility (7).

#### 4. Generation of Catalytic Antibodies

The recovery of catalytic antibodies relies on an efficient means of sampling the immune repertoire and screening for effective catalysts. The murine immune repertoire is estimated to have a diversity of  $10^8$ , and it can be further expanded by immunization. The majority of catalytic antibodies have been derived from hybridomas that generally capture only 0.01% of the immune repertoire. The development of recombinant antibody fragments as combinatorial [libraries](#) consisting of  $\geq 10^6$  members in **lambda** and in [M13 phage](#) has increased the access to the immune repertoire and possibly the recovery of catalytic antibodies. The diversity of an antibody combinatorial library, and potentially the yield of efficient catalytic antibodies, can be further expanded by polypeptide chain-shuffling experiments in which the heavy-chain fragment of a catalytic antibody is allowed to cross with a library of light-chain fragments elicited by the same immunoglobulin, and vice versa (8).

#### 5. Screening for Catalytic Antibodies

Catalytic antibodies have been identified by their high affinities for the appropriate haptens, as well as their catalytic activities for the desirable reactions. Given the often low activity of catalytic antibodies, sensitive but specific methods have been developed for identifying the catalyst. CatELISA is based on screening for the presence of antibody-catalyzed reaction products using product-specific antibodies as primary antibody indicators in [enzyme-linked immunosorbent assay](#) (ELISA) (9). An alternative method depends on the catalytic release of a reactive species that can be trapped via covalent modification of the phage-displayed antibody, which then permits the recovery of DNA encoding catalytic antibody clones from the phage particles (10). Another method involves the recovery of antibody clones that sustain growth of **microorganisms**, such as **yeast** or bacterial [auxotrophs](#), through catalysis of essential biological pathways (11). The host organisms, however, may also impose an additional selection on what catalytic antibodies can be recovered from the screen, since the host will not survive expression of antibodies that are toxic.

#### Bibliography

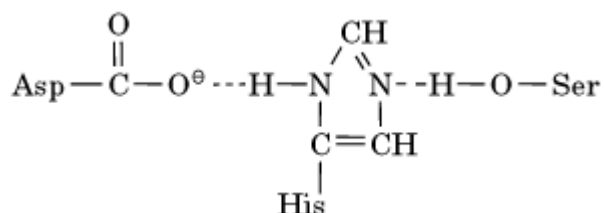
1. R. A. Lerner, S. J. Benkovic, and P. G. Schultz (1991) *Science* **252**, 659–667.
2. S. J. Benkovic (1992) *Annu. Rev. Biochemistry* **61**, 29–54.
3. J. R. Jacobsen and P. G. Schultz (1995) *Curr. Opin. Struct. Biol.* **5**, 818–824.
4. J. B. Charbonnier et al. (1997) *Science* **275**, 1140–1142.
5. G. J. Wedemayer et al. (1997) *Science* **276**, 1665–1669.
6. J. D. Stewart et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 7404–7409.
7. M. R. Haynes et al. (1994) *Science* **263**, 646–652.
8. B. Posner et al. (1994) *Trends Biochem. Sci.* **19**, 145–150.
9. D. S. Tawfik et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 373–377.

10. K. D. Janda et al. (1997) *Science* **275**, 945–948.
11. Y. Tang, J. B. Hicks, and D. Hilvert (1991) *Proc. Natl. Acad. Sci. USA* **88**, 8784–8786.
12. J. A. Smiley and S. J. Benkovic (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8319–8323.

## Catalytic Triad

The term *catalytic triad* is used to describe the arrangement of amino acid residues in the [active sites](#) of [serine proteinases](#) that underlies their mechanism of [enzyme](#) action. It was first applied to the active site of **chymotrypsin** when the structure of that enzyme was determined by [X-ray crystallography](#) (1). Earlier studies had shown that chymotrypsin could be inactivated by chemical modification of one of its serine residues with the nerve gas, diisopropylfluorophosphate (DIFP). Examination of the enzyme structure revealed that the hydroxyl group of this [serine](#), residue 195, was “activated” by the imidazole group of [histidine](#) 57, which in turn was “activated” by the carboxyl group of [aspartate](#) 102 (Fig. 1). Originally described as a *charge relay system*, these three residues are more commonly called a *catalytic triad*.

**Figure 1.** The catalytic triad, the hydrolytic apparatus of serine proteinases. It consists of the side chains of an aspartic acid, a histidine, and a serine residue (eg, residues 102, 57 and 195, respectively, in bovine chymotrypsin). In the absence of substrate, the imidazole group of the histidine is unprotonated, but it potentiates the nucleophilic properties of the hydroxyl side chain of the serine. When substrate binds to the active site, the serine proton is transferred to the imidazole group and the oxygen attacks the carbonyl group of the substrate. The resulting positively charged imidazole is stabilized by interaction with the negative charge of the aspartate carboxyl group. The activated serine is also highly reactive toward inhibitors such as PMSF and DIFP.



Similar triads have been observed in both the chymotrypsin and [subtilisin](#) evolutionary families of serine proteinases (2), as well as in other hydrolytic enzymes (3). It is remarkable that the bacterial serine proteinase [subtilisin](#) has the same geometrical arrangement of aspartate, histidine, and serine residues as chymotrypsin, but all other structural aspects of the two proteins are quite different. This is a classic example of **convergent** evolution. The activated serine interacts with the carbonyl carbon atom of the peptide bond to be hydrolyzed, forming an oxyanion intermediate, which is converted to an acylated serine with release of the amine component of the peptide. Subsequent hydrolysis of the acylenzyme occurs by the reverse reaction in which [water](#) substitutes for the activated serine.

## Bibliography

1. D. M. Blow (1976) *Acc. Chem. Res.* **9**, 145–152.
2. T. A. Steitz and R. G. Schulman (1982) *Annu. Rev. Biochem. Biophys.* **11**, 419–444.
3. D. M. Blow (1990) *Nature* **343**, 694–695.

## Cathepsins

The term *cathepsin* is derived from the Greek word meaning “to digest” and is used to describe a broad range of intracellular [proteinases](#) that serve important biological functions. Two major pathways have been identified that control the degradation of cellular proteins. One operates within the **cytoplasm** and is called the **ubiquitin**-mediated pathway. The other functions within specific cellular compartments, especially the **lysosomes**, but also in [endosomes](#) and the [endoplasmic reticulum](#). Lysosomes are the major site of intracellular protein degradation (1). They are amply endowed for this function because they contain as many as 20 different proteolytic enzymes and as many other hydrolytic enzymes (eg, **phosphatases**, [lipases](#), **nucleases**, and sulfatases) (1).

The lysosomal proteinases are called *cathepsins*. They are single-chain proteins that range in molecular mass from 20 kDa to about 40 kDa. Typically they have a broad substrate specificity, are optimally active at somewhat acidic pH, and belong to either the [thiol proteinase](#) or [carboxyl proteinase](#) classes. For no particular reason, they have been designated by letters as cathepsins B, L, H, M, N, S, and T, all of which are thiol proteinases or, more commonly, *cysteine proteinases*; and cathepsins D and E, which are carboxyl (or *aspartyl*) proteinases (2). All of these are [endopeptidases](#), enzymes that cleave [peptide bonds](#) in the internal part of a protein. In addition, there are several lysosomal *exopeptidases*, including *cathepsin C*, which removes dipeptides from the amino-terminus of polypeptide chains and is therefore a dipeptidyl aminopeptidase, and at least two carboxycathepsins that are [carboxypeptidases](#). Cathepsin G is a chymotrypsin-like serine proteinase found in neutrophils, and cathepsin R is another serine proteinase found in the [endoplasmic reticulum](#).

The lysosomal cathepsins have been considered to be the most active proteinases in the body, because they degrade four times as much protein as the pancreatic and gastric proteinases combined (3). Much of this is essentially turnover of cytosolic proteins to balance **protein synthesis** and maintain homeostasis. However, cathepsins are also important components of the **immune system**: In order for a foreign protein to generate an [antibody](#) response, it must be taken up by a specialized **antigen-presenting** cell and degraded within endosomes. The resultant peptide fragments are then “presented” on the cell surface, where they trigger antibody production (4).

### Bibliography

1. F. Authier, B. I. Posner, and J. J. M. Burgeron (1994) In *Cellular Proteolytic Systems* (A. J. Ciechanover and A. L. Schwartz, eds.), Wiley-Liss, New York, pp. 89–113.
2. J. S. Bond and P. E. Butler (1987) *Annu. Rev. Biochem.* **56**, 333–364.
3. A. J. Barrett and H. Kirschke (1981) *Methods Enzymol.* **80**, 535–561.
4. C. V. Harding (1994) In *Cellular Proteolytic Systems* (A. J. Ciechanover and A. L. Schwartz, eds.), Wiley-Liss, New York, pp. 163–180.

## Caudal Protein



Caudal protein (Cad) was identified originally as a **homeodomain** protein in *Drosophila*, but **homologues** have subsequently been found in vertebrates as well. It belongs to the superfamily of **DNA-binding** homeodomain **transcription factors**. Caudal **messenger RNA**, which is supplied by the mother in the unfertilized **egg**; is required for proper segmentation of posterior segments in insects and for the development of the posterior lineage in vertebrates. In *Drosophila*, Caudal has been shown to bind to **promoter** elements and to activate directly the **transcription** of several segmentation genes. In this process, it cooperates with other transcriptional activators of segmentation genes, eg, **Bicoid** (another homeodomain protein). At the preblastoderm stage, Caudal protein forms a gradient, which is generated by **translational repression** caused by Bicoid binding to the *cad* mRNA. Although this early function appears to be specific to *Drosophila* (and other long-germ-band insects), the zygotic function of Caudal, specification of the most posterior segments in insects or posterior cell lineage in vertebrates, appears to be conserved throughout the animal kingdom.

The *caudal* (*cad*) gene of *Drosophila* was first isolated as a **homeobox**-containing gene by cross-hybridization to the homeobox of *Ultrabithorax* (*Ubx*) and others (1, 2). The *cad* homeodomain shares 58% identity with that of *ftz*, and 53% and 52%, respectively, with those of *Antp* and *Ubx*. During *Drosophila* development, *cad* is expressed both maternally and zygotically. The maternal expression is detected in the **germ line** in nurse cells during **oogenesis**, and its mRNA accumulates in the developing **oocyte**. At the end of oogenesis, the *cad* mRNA is distributed evenly throughout the oocyte. No Cad protein is detected during oogenesis. Cad protein is first detected beginning with nuclear division cycles 6 or 7, just prior to the formation of the syncytial blastoderm. From this stage on to the 13 nuclear cycle, the distribution of Cad protein forms a concentration gradient throughout the embryo and along the anterior-posterior axis, with its maximum at the posterior end. The *cad* mRNA remains uniformly distributed throughout the cytoplasm to the 12 nuclear cleavage; subsequently, it follows the protein gradient, when it is then distributed as a concentration gradient itself. Nevertheless, the protein gradient precedes the mRNA gradient by several nuclear cycles, so that the mRNA gradient is not relevant for the generation of its protein gradient (see below). Following the final nuclear cleavage during cellularization of the blastoderm, the maternally-derived *cad* mRNA disappears, and with it the protein translated from it (1-3).

Formation of the protein gradient from the initially uniformly distributed mRNA is generated by translational repression in the anterior regions of the embryo. Another maternally provided homeodomain protein, Bicoid, binds the *cad* mRNA and prevents its translation (4, 5). This is the only reported example of a homeodomain protein controlling the expression of another homeodomain factor at the level of its mRNA by translational repression. Caudal protein also accumulates in the pole cells during their formation and is present there for several hours of embryogenesis. There is no *de novo* expression in the pole cells, however, and all the protein present was synthesized from the maternally provided mRNA. Formation of the Caudal protein gradient does not depend on zygotic gene expression. In older unfertilized eggs, an anterior-posterior gradient is detected that is similar to that normally present in embryos of the same stage (2, 3). Maternally expressed Caudal has also been reported in other insects, such as *Bombyx mori* (6).

Zygotic expression of Caudal protein in *Drosophila* begins during cellularization of the blastoderm as a stripe in the posterior region of the embryo corresponding to the anlagen of the most posterior segments (the presumptive abdominal segments A9/A10, parts of the telson) and the hindgut. At the extended-germ-band stage, and throughout the rest of embryogenesis, *cad* is expressed in posterior midgut cells, the Malpighian tubules, the anal plate and, at weaker levels, during germ-band extension in pair-rule-like stripes (1-3). In third instar larvae, *cad* is expressed in the Malpighian tubules, the posterior midgut, the gonads of both sexes, and in the anlagen of the anal plates and the hindgut, within the genital (3). Most aspects of the zygotic expression are conserved in those vertebrates that have been analyzed (7-10).

In *Drosophila*, the function of the Caudal transcription factor is required for segment specification in the thoracic and abdominal regions by regulating the transcription of several segmentation genes (11-14). Embryos lacking both maternal and zygotic Cad protein are severely shortened, with variable deletions of many of the thoracic and abdominal segments. Zygotically provided protein can, however, largely rescue these defects. Ectopic expression of Caudal protein throughout the embryo at the blastoderm stage leads to **phenotypes** that are somewhat opposite to the phenotypes observed in *cad*-defective embryos (15). Analysis of the [cis-acting](#) elements of several segmentation genes has demonstrated that Cad is directly activating transcription of *fushi tarazu* (11) and other segmentation genes (12-14). Caudal acts in concert with Bicoid, and is partially redundant, in controlling the expression of the downstream targets, the segmentation genes (12, 13).

Embryos from wild-type mothers that lack zygotic function display an absence of parts of the most posterior segments, the terminalia; in particular, the anal tuft, parts of the anal pads, and the terminal sense organs, all structures that are derived from a cryptic tenth abdominal segment, are deleted (2). The isolation and developmental expression analysis of Caudal homologues in vertebrates (7-10) has suggested that this aspect of Caudal function, specification of most posterior segments or posterior lineage in vertebrates, has been conserved throughout [evolution](#). Based on the expression patterns of some of the vertebrate Caudal homologues, additional roles for this subfamily of homeodomain transcription factors have been proposed; eg, participation in rostrocaudal axial patterning or in the development and regeneration of the liver (16, 17). However, all the assumptions of Caudal function (s) in vertebrate development are based purely on the respective expression patterns, and no loss-of-function mutants have been reported to date. Thus, all these proposed roles of vertebrate Caudal homologues await confirmation by the analysis of specific mutants, eg, **knockout** mice, or the potential identification of mutants in other model systems, such as the [zebrafish](#).

#### Bibliography

1. M. Mlodzik, A. Fjose, A. Gehring, and W. J. Gehring (1985) *EMBO J.* **4**, 2961–2969.
2. P. M. MacDonald and G. Struhl (1986) *Nature* **324**, 537–545.
3. M. Mlodzik and W. J. Gehring (1987) *Cell* **48**, 465–478.
4. J. Dubnau and G. Struhl (1996) *Nature* **379**, 694–699.
5. R. Rivera-Pomar, D. Niessing, U. Schmidt-Ott, W. J. Gehring, and H. Jäckle (1996) *Nature* **379**, 746–749.
6. X. Xu, P. X. Xu, and Y. Suzuki (1994) *Development* **120**, 277–285.
7. A. Frumkin, Z. Rangini, A. Ben-Yehuda, Y. Gruenbaum, and A. Fainsod (1991) *Development* **112**, 207–219.
8. A. Frumkin et al. (1993) *Development* **118**, 553–562.
9. J. S. Joly et al. (1992) *Differentiation* **50**, 75–87.
10. F. Beck, T. Erler, A. Russell, and R. James (1995) *Dev. Dyn.* **204**, 219–227.
11. C. Dearolf, J. Topol, and C. Parker (1989) *Nature* **341**, 3430–3433.
12. R. Rivera-Pomar, X. Lu, N. Perrimon, H. Taubert, and H. Jäckle (1995) *Nature* **376**, 253–256.
13. C. Schulz and D. Tautz (1995) *Development* **121**, 1023–1028.
14. T. Hader et al. (1998) *Mech. Dev.* **71**, 177–186.
15. M. Mlodzik, G. Gibson, and W. J. Gehring (1990) *Development* **109**, 271–277.
16. A. V. Morales, E. J. de la Rosa, and F. Pablo (1996) *Dev. Dyn.* **206**, 343–353.
17. U. Doll and J. Niessing (1996) *Eur. J. Cell Biol.* **70**, 260–268.

#### Suggestions for Further Reading

18. R. Rivera-Pomar and H. Jäckle (1996) From gradients to stripes in *Drosophila* embryogenesis: filling in the gaps. *Trends Genet.* **12**, 478–483.
19. D. Duboule (1994) *Guidebook to Homeobox Genes*. Sambrook and Tooze Publication, Oxford

## Cauliflower Mosaic Virus

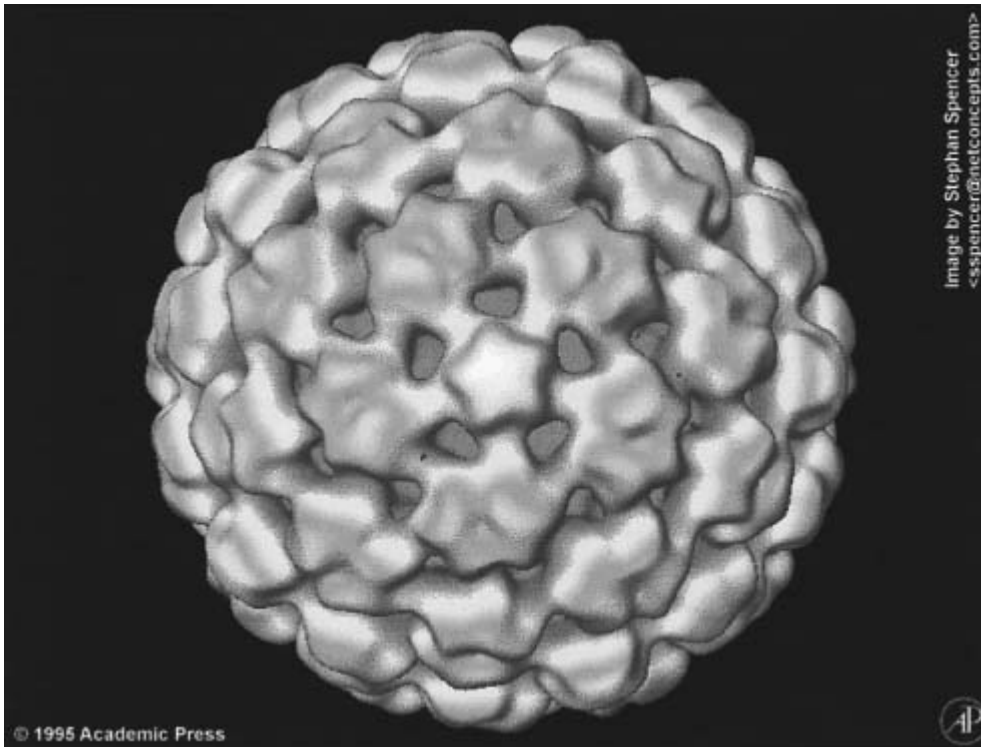
*Cauliflower mosaic virus* (CaMV) was the first plant **virus** shown to have a DNA [genome](#) and the first shown to **replicate** by **reverse transcription**. The virus is found worldwide but only causes significant losses locally. It is transmitted by aphids in the externally borne or stylet-borne manner (see [Virus Infection, Plant](#)) and encodes a helper factor that is required for transmission. It has a narrow host range, being restricted primarily to the Cruciferae. CaMV is the type member of the Caulimovirus genus, which contains 11 species and 6 possible members.

This virus has had a significant impact on plant virology and plant molecular biology. The unravelling of the reverse transcription replication mechanism led to the concept of *pararetroviruses* (1, 2). The virus is an important source of gene regulatory elements, which have been used extensively in the genetic manipulation of plants.

### 1. Virus Structure

CaMV particles are isometric, about 50 nm in diameter (**Fig. 1**). Although the virus particles have been crystallized, [X-ray crystallography](#) did not yield detailed information, most probably because of the heterogeneity of the coat protein (3). The structure has been shown by [cryoelectron microscopy](#) and three-dimensional **image reconstruction** to comprise three concentric layers of solvent-excluded density (4). The outermost layer is made up of a total of 420 coat protein subunits arranged in  $T = 7$  **icosahedral symmetry**; it is the first example of a  $T = 7$  virus that obeys the stoichiometric rules of icosahedral symmetry. The DNA genome is distributed in layers II and III, together with some of the viral coat protein. The viral coat protein is expressed from **open reading frame IV** (see below) as a 58-kDa precursor molecule that is processed **proteolytically** to several molecular sizes, the major ones being 42 and 37 kDa. The protein in virions is **glycosylated** and **phosphorylated**, and the virions also contain minor amounts of the products of open reading frames III and V.

**Figure 1.** External structure of cauliflower mosaic virus, viewed down a local five-fold axis.

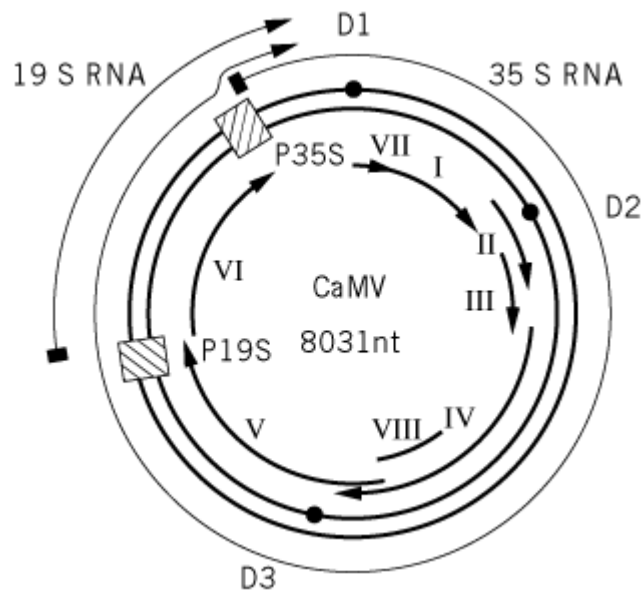


CaMV particles are very stable and are only dissociated at high pH (5) or by digestion of the coat protein with proteinases.

## 2. Viral Genome

Each virus particle contains a single molecule of circular double-stranded DNA of 8 kbp. The DNA of most isolates of CaMV has three discontinuities (D1, D2, and D3 in Fig. 2), one in one strand (the a strand) and two in the other (yielding the b and g strands); one strain has only two discontinuities, one in each strand. The discontinuities have an unusual structure, with a fixed 5' DNA nucleotide (with sometimes one or two ribonucleotides attached) and variable 3' end, overlapping the 5' end by 10 to 30 nucleotides. They are the sites of the priming of (+)- and (-)-strand [DNA replication](#), and this unusual structure results from viral replication (see below). The DNA has a twisted conformation that, because of the discontinuities, cannot be **supercoiled**, the constraining forces being unknown.

**Figure 2.** Genome organization of CaMV. The double-stranded DNA genome is represented by the thick double circle, with the discontinuities shown as ●. The promoters for the 35S and 19S transcripts are indicated by ☒, and the positions of the transcripts are shown by the outer arcs. The inner arcs are the open reading frames I to VIII.



CaMV DNA has six, or possibly eight, open reading frames (ORFs) on the complement to the a strand (**Fig. 2**). ORF I product (37 kDa) is involved in the cell-to-cell movement of the virus, and ORF II product (18 kDa) is the aphid transmission protein (see [Virus Infection, Plant](#)). The product of ORF III is a protein of 15 kDa that is processed to 11 kDa and thought to be involved in packaging the DNA genome (6). ORF IV encodes the viral coat protein, which contains the “cys” motif  $C^{ys-} X_2 C^{ys-} X_4 H^{is-} X_4 C^{ys-}$  characteristic of [retrovirus gag](#) or nucleoproteins and a highly basic lysine-rich region. The viral [polymerase](#) (reverse transcriptase +ribonucleaseH) is encoded by ORF V (79 kDa). The product of ORF VI forms the matrix of the [inclusion bodies](#) in which virus particles are found embedded in the **cytoplasm** of infected cells. It has several functions, including being the site of the reverse transcription phase of viral replication (7), involved in virus assembly (8), a transactivator during viral expression (see below), and a symptom determinant. No products have been found for ORFs VII and VIII, and it is uncertain whether they are expressed.

### 3. Transcription and Translation

CaMV DNA is transcribed asymmetrically from the a strand to give a more-than-genome length transcript (the 35S RNA), with a terminal repeat of 180 nucleotides and subgenomic transcript (the 19S RNA) that is the [messenger RNA](#) for ORF VI (**Fig. 2**); the transcripts are **polyadenylated**. There are suggestions of other transcripts, such as one for ORF V (9), but these have not been characterized. **Splicing** events between the leader sequence and an acceptor site in ORF II (10) are essential for virus infectivity. It is possible that these events might form the mRNA for expression of ORF III and, possibly, ORFs IV and V.


The **promoter** for the 35S RNA is well characterized and widely used in the expression of transgenes in plants. It directs constitutive, high-level expression in most tissues of most plant species. The promoter has a modular structure, with regions that control the tissue specificity of expression [reviewed in 11]. Nuclear factors have been identified that bind to specific regions of the promoter. The 19S promoter is much weaker than the 35S promoter and is less well characterized.



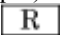
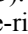
The CaMV genome also contains the signal sequence for the polyadenylation of both the 35S and 19S transcripts. Due to the terminal repeat, the transcription of the 35S RNA has to pass this signal the first time it is encountered and recognize it the second time. It is thought that the bypass is due to the proximity of the promoter (12).

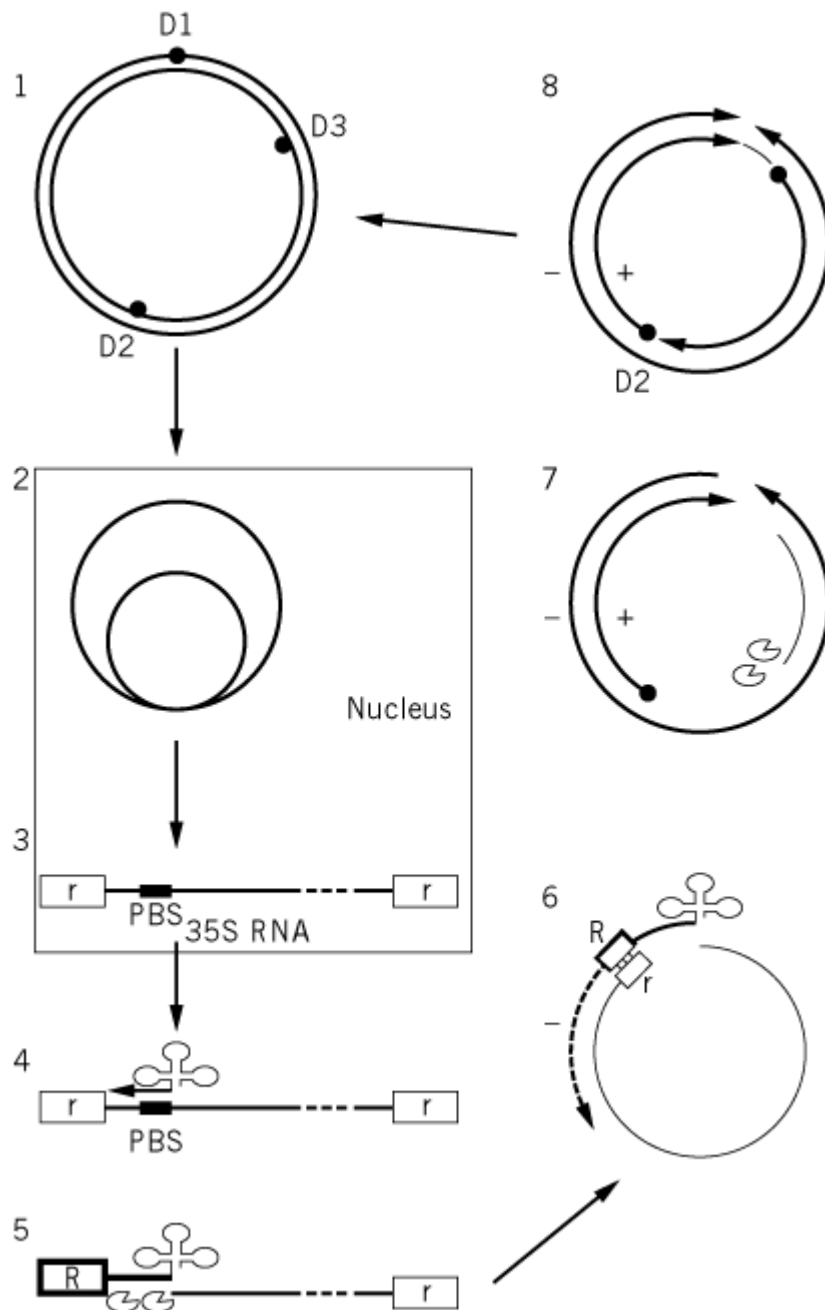
[Translation](#) of the 35S RNA is effected by several unusual mechanisms [reviewed in (11) and further refined in (10)]. The leader sequence upstream of ORFI is long (>600 nucleotides) and contains several **AUG** start codons, which could inhibit downstream translation. However, this leader sequence is folded into a complex stem-loop structure, which enables small ribosomal subunits to “shunt” from the 5′ to 3′ end of the leader sequence without scanning the sequence in between; this then opens up ORFI. Translation of downstream ORFs on this polycistronic message is dependent on the presence of a **transactivator**, which is the product of ORF VI, and possibly on the splicing mechanism described above. This process of transactivation seems to enable [ribosomes](#) that have translated one ORF to remain competent to translate the next ORF downstream.

#### 4. Viral Replication

The replication cycle of CaMV (Fig. 3) has been studied in much detail, but there are still several aspects that are not fully understood. Virus particles enter the uninfected cell and are unencapsidated by an unknown mechanism. The viral DNA genome enters the [nucleus](#), where the discontinuities are sealed and the molecule associates with [histones](#) forming a [minichromosome](#) (Fig. 3, step 2). This is the template for transcription by the host DNA-dependent **RNA polymerase**, giving the 35S and 19S RNAs, which pass to the cytoplasm (Fig. 3, step 3). The 19S RNA is translated to yield the inclusion body protein. The 35S RNA has two functions, being the mRNA for at least the products of ORFs I and II and the template for the reverse transcription phase of replication. There is some evidence that this phase of replication takes place in viruslike particles in the inclusion bodies (7). The first event in reverse transcription is the annealing of the 3′ end of tRNA<sup>met</sup><sub>init</sub> to a site, that of discontinuity 1 (D1), about 600 nucleotides from the 5′ end of the 35S RNA (Fig. 3 step 4). This acts as a primer for the synthesis of a-strand DNA toward the 5′ end of the template. RNaseH activity removes the RNA moiety of the DNA–RNA duplex thus formed and leaves DNA (Fig. 3, step 5) that, because of the terminal repeat of the 35S RNA, is complementary to the 3′ end of the template. Thus, a strand switch is effected, and synthesis of the a strand DNA continues (Fig. 3, step 6). The digestion of the RNA by RNaseH leaves purine-rich regions at the positions of discontinuities 2 and 3 (D2 and D3), which act as primers for the synthesis of the b and g strands (Fig. 3 steps 7 and 8). When the oncoming strand of (+)-strand DNA reaches the next priming site, it displaces the primer and some newly synthesized DNA, resulting in the triple-strand discontinuity described above. A second strand switch occurs at the tRNA priming site between the oncoming a and g strands, allowing the synthesis of the g strand to continue and making discontinuity 1.

**Figure 3.** Replication cycle of CaMV. Step 1: The input DNA containing discontinuities 1–3 (D1–D3) passes to the nucleus where it forms minichromosomes (step 2). These are the template for the transcription of the 35S RNA (step 3) with terminal repeats ,

which moves to the cytoplasm. Priming of the 35S RNA occurs by the annealing of the 3′ end of tRNA<sup>met</sup><sub>init</sub>  to the primer binding site (PBS) on the 35S RNA and leads to the synthesis of (–)-strand DNA (step 4). RNaseH  removes the RNA of the RNA DNA duplex, leaving the complement to the terminal repeat sequence  (step 5). This anneals to the 3′ end of the 35S RNA, and DNA synthesis continues (step 6). RNaseH activity leaves purine-rich regions , which act as primers for (+)-strand DNA synthesis (steps 7 and 8). The oncoming strand displaces the primer sequence, thus giving the characteristic discontinuities, and a second strand-switch is effected at D1 to complete the molecule.



## Bibliography

1. R. Hull and H. Will (1989) *Trends Genet.* **5**, 357–359.
2. H. Temin (1989) *Nature* **339**, 254–255.
3. Z. X. Gong, H. Wu, R. H. Cheng, R. Hull, and M. G. Rossmann (1990) *Virology* **179**, 941–945.
4. R. H. Cheng, N. H. Olson, and T. S. Baker (1992) *Virology* **186**, 655–668.
5. R. Al Ani, P. Pfeiffer, G. Lebeurier, and L. Hirth (1979) *Virology* **93**, 175–187.
6. J. L. Mougeot, T. Guidasci, T. Wurch, G. Lebeurier, and J. M. Mesnard (1993) *Proc. Natl. Acad. Sci. USA* **90**, 1470–1473.
7. C. M. Thomas, R. Hull, J. A. Bryant, and A. J. Maule (1985) *Nucl. Acids Res.* **13**, 4557–4576.
8. A. Himmelbach, Y. Chapdelaine, and T. Hohn (1996) *Virology* **217**, 147–157.
9. A. L. Plant, S. N. Covey, and D. Grierson (1985) *Nucl. Acids Res.* **13**, 8305–8321.
10. Z. Kiss-László, S. Blanc, and T. Hohn (1995) *EMBO J.* **14**, 3552–3562.

11. H. M. Rothnie, Y. Chapdelaine, and T. Hohn (1994) *Adv. Virus Res.* **44**, 1–67.
12. H. Sanfaçon and T. Hohn (1990) *Nature (Lond.)* **346**, 81–84.

### Suggestion for Further Reading

13. R. J. Shepherd (1994) "Caulimoviruses". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 223–226.

## cDNA Libraries

In contrast to [genomic libraries](#), which contain raw DNA sequences harvested from an organism's [chromosomes](#), cDNA libraries are composed of processed nucleic acid sequences harvested from an organism's RNA pools (see [Complementary DNA \(cDNA\)](#)). cDNA can be prepared from any RNA source, including in RNA and noncoding RNA, isolated from single cell organisms, cultured metazoan cells, and isolated tissues. cDNA libraries provide a powerful means of examining cell- and tissue-specific gene expression. For example, mammalian cDNA libraries contain only the fraction of sequences that are expressed in the subject tissues at the time of harvest, typically in the range of 10% of the **genes** carried by the [genome](#). Also, cDNA is typically prepared enzymatically from **polyadenylated** mRNA, a population of RNA that has previously undergone post-transcriptional editing and removal of any intervening sequences (**introns**). Whereas genes in genomic DNA are often interrupted by numerous introns, the copy of a gene in a cDNA or mRNA commonly contains an uninterrupted sequence encoding the gene product. Thus, cDNA inserts are directly expressible as both RNA and protein gene products. Vectors used in cDNA libraries often contain expression elements that mediate transcription and translation of the cloned cDNAs.

In constructing cDNA libraries, the goal is usually to represent full-length protein coding sequences for all of the mRNAs expressed in the subject cells or tissue. Thus, it is essential to preserve the integrity of cellular mRNA during extraction and purification protocols. Libraries made from fragmented RNA will not preserve the full open reading frame, rendering the resultant partial clones less useful. Unfortunately, RNA tends to be a transient species by design, as cells have developed complex post-transcriptional mechanisms to regulate gene expression through mRNA decay. The problem is worse in certain tissues, such as pancreas, that secrete large amounts of degradative enzymes as part of their normal function. Though troublesome, the RNA decay problem can be overcome through specialized procedures that seek to inactivate **ribonuclease** enzymes concomitantly with RNA extraction from intact tissue.

Once total RNA is extracted from the subject tissue, the mRNA fraction is purified by annealing to oligodeoxythymidine (oligo-(dT)) cellulose. Oligo-(dT) is a short polymer that hybridizes specifically to polyadenylate. Because the majority of cellular mRNAs are **polyadenylated**, affinity chromatography over oligo-(dT) cellulose provides a simple means to remove contaminating RNA species, such as ribosomal and transfer RNAs, which do not contain [poly A](#) tails. The purified mRNA is then annealed to oligo-(dT), which serves as a primer that can be extended by the enzyme **reverse transcriptase** (RT). RT extension thus generates **antisense** DNA copies of the purified mRNA species. Under appropriate conditions, RT can be induced to complete a second-strand DNA synthesis, using the first DNA strand as a template. The result is a collection of double-stranded DNA copies of the parental mRNA population (hence the “c” in “cDNA”). The double-stranded DNA population can then be cloned into an appropriate plasmid or bacteriophage vector to constitute a cDNA library.



cDNA libraries may be screened either by **hybridization** to the cloned inserts or by detection of RNA or protein expression products encoded by the cDNA inserts. Hybridization screening is similar in principle for both genomic and cDNA libraries and uses synthetic oligonucleotide probes or gene probes derived from existing clones. The probes are labeled with radioactivity or a fluorescent tag and then annealed to the total nucleic acids contained within a population of bacterial colonies or viral plaques. Colonies or plaques that contain nucleic acids complementary to the probe can then be detected and purified. Protein expression libraries can be screened based on affinity to [antibodies](#), proteins, nucleic acids, or small-molecule ligands. In addition, functional properties can be screened for, such as [catalysis](#) or induction of [transcription](#) of a [reporter gene](#).

Useful information regarding tissue-specific patterns of gene expression can be ascertained through screening cDNA libraries. By definition, a cDNA population is derived from the fraction of genes expressed in a given tissue through sampling the mRNA pool. These expressed sequences can be sampled and characterized in brute force fashion by high-throughput random sequencing of cDNA clones isolated from a library. It is not necessary to determine the full-length cDNA sequence in order to derive useful information from this approach. The partial sequences obtained are entered into an [Expressed Sequence Tag \(EST\) database](#) where the frequency of each partial sequence (known as a “tag”) can be scored (1). In this manner, it is possible to determine the expression level of individual genes for any tissue. Furthermore, investigators that identify a gene product of interest in the EST database are already one step ahead on their sequencing efforts, because much of the open reading frame may already be in the public database. Finally, it is possible to request many EST clones from the I.M.A.G.E. Consortium (Integrated Molecular Analysis of Genomes and their Expression; <http://www-bio.llnl.gov/bbrp/image/image.html>), providing an easy source of genetic material for subcloning projects.

The concept of ESTs as quantitative indicators of gene expression levels has been further refined through a new method termed Serial Analysis of Gene Expression (SAGE) (2). SAGE provides a means of tagging every mRNA in the cell with a small, intrinsic nine-nucleotide sequence that is essentially unique to each transcript. In addition, SAGE provides a means to eliminate artifacts arising through processing and polymerase chain reaction (PCR) amplification of the original cDNA, thereby rendering the method highly quantitative. Finally, SAGE enables high throughput sequencing of the nine-nucleotide tags on a scale sufficient to quantify several hundred thousand cDNA species for a given tissue (3). Although the sequence information obtained from a single SAGE tag is minimal, it is usually sufficient to identify the transcript via hybridization or database searching methods. SAGE is being conducted on a large scale through the Cancer Genome Anatomy Project (CGAP; (site currently unavailable)), with the goal of rendering in detail the molecular differences between various normal tissues and tumors, in order to understand better the molecular basis of cancer. Ultimately the CGAP database, among others, will provide the capability of performing “virtual” **Northern blots**, in which any sequence segment can be screened through WWW-based databases to ascertain all known sequence, genetic, and biological information associated with mRNAs containing the sequence segment.

Tissue- or cell-specific gene expression can also be studied through specialized cDNA libraries derived from the subset of genes that are differentially expressed between two preparations of mRNA (4). For example, a cell culture line can be treated with a [hormone](#), such as [dexamethasone](#), and RNA can be prepared from samples of the cell culture before and after treatment. cDNAs prepared from the pretreatment population can then be used to remove selectively any complementary mRNAs from the post treatment preparation through a process known as [subtractive hybridization](#). The majority of the mRNA species remaining in the post treatment pool should be those transcripts that were synthesized as a direct result of hormone treatment. These specific mRNAs can then be converted into cDNA and cloned into a suitable library vector, to constitute a differential cDNA library. Screening or sequence tagging a subtractive library greatly increases the efficiency of characterizing tissue-specific or stimulus-specific differences in gene expression.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA](#)

[Libraries](#), [Genomic Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

## Bibliography

1. M. D. Adams, M. Dubnick, A. R. Kerlavage, R. Moreno, J. M. Kelley, T. R. Utterback, J. W. Nagle, C. Fields, and J. C. Venter (1992) *Nature* **355**, 632–634.
2. V. E. Velculescu, L. Zhang, B. Vogelstein, and K. W. Kinzler (1995) *Science* **270**, 484–487.
3. L. Zhang, W. Zhou, V. E. Velculescu, S. E. Kern, R. H. Hruban, S. R. Hamilton, B. Vogelstein, and K. W. Kinzler (1997) *Science* **276**, 1268–1272.
4. J. S. Wan, S. J. Sharp, G. M. Poirier, P. C. Wagaman, J. Chambers, J. Pyati, Y. L. Hom, J. E. Galindo, A. Huvar, P. A. Peterson, M. R. Jackson, and M. G. Erlander (1996) *Nat Biotechnol* **14**, 1685–1691.

## Cell Adhesion Molecules

Cell adhesion molecules play central roles in embryonic development, tissue organization, and overall maintenance of the structure and form of multicellular organisms ([1](#), [2](#)) They provide the physical links between cells and other cells, and between cells and the extracellular matrix. Adhesion molecules are also involved in dynamic cell interactions, such as during tissue morphogenesis, cell migration, immunological and inflammatory responses, and possibly memory formation ([1](#), [3-7](#)). Adhesive proteins are also intimately involved in intracellular signaling processes that regulate cell growth, differentiation, gene regulation, and programmed cell death ([8-13](#)).

Cell adhesion molecules mediate adhesive interactions by forming specific protein-protein interactions or between proteins and complex carbohydrates. In addition, cell adhesion molecules frequently link directly to multimolecular protein complexes on the cytoplasmic face of the plasma membrane, which, in turn, mediate interactions with the cytoskeleton and signal transduction pathways ([3](#), [9-13](#)). Consequently, these cell adhesion and signaling complexes not only link cells with other cells and the extracellular matrix but also help to integrate extracellular physical information with the major signal transduction pathways within cells (see [Integrins](#)). For example, they appear to transduce information from inputs as diverse as transient contacts with other cells, binding of different connective tissue molecules, and tension or torsion at the cell surface. Each of these inputs connects to intracellular signaling pathways, such as phosphorylation cascades ([9-12](#)).

### 1. Types of Adhesion

The physical linkages involved in cell adhesion are provided by adhesion molecules that function as parts of general adhesion systems or in specialized adhesive structures. For example, epithelial cells can adhere to other cells along broad expanses of plasma membrane using general-purpose adhesive molecules such as cadherins ([14-16](#)). They can also adhere to each other by specialized adhesive structures, such as adherens junctions, desmosomes, and tight junctions ([1](#)). Each type of specialized junctional complex requires specific adhesion molecule components, which can include specialized cadherins and other proteins. Fibroblasts can also form complexes with other cells, but they most characteristically adhere to extracellular matrix. Cell-to-matrix adhesions can be broad and flat, as when epithelial cells adhere via integrins to basement membranes, or they can involve specialized structures that include hemidesmosomes or the focal contacts of fibroblastic and endothelial cells. The protein complexes involved in cell-to-matrix adhesion contain different complements of adhesive, cytoskeletal, and signal transduction molecules than those involved in cell-to-cell adhesion ([1](#), [13](#), [15](#)). This entry will first provide a general overview of specific types of cell adhesion

proteins, then briefly discuss how they are integrated into cytoskeletal and signaling networks.

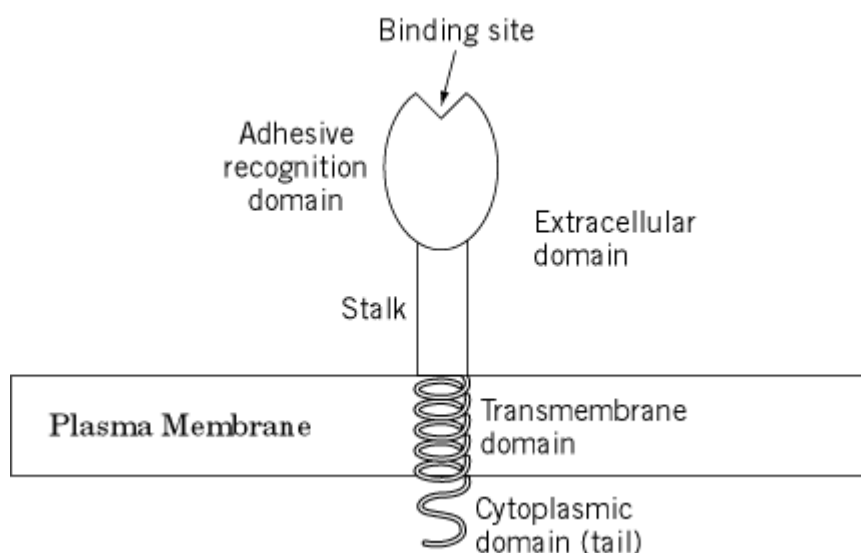
## 2. General Features of Cell Adhesion Proteins

Cell adhesiveness is generally based on the specific binding of a protein to another molecule at the cell surface. When adhesion results from binding of an adhesion molecule to the same type of protein on a neighboring cell, the interaction is termed ‘*homophilic*.’ The cadherin family is a major mediator of such homophilic interactions (14-16) (see [Cadherins](#)). Cadherins form complexes in which both intercellular and lateral binding interactions cooperate to create tightly packed adhesion complexes that mediate adhesion with high avidity.

In many systems, however, adhesiveness involves binding of a receptor to a specific ligand. In cell-to-cell adhesion, the target protein of an adhesion receptor can be a “counter-receptor” or a complex carbohydrate on a protein anchor in the plasma membrane. In cell-to-matrix interactions, a plasma membrane adhesion protein such as an integrin can bind to an extracellular matrix protein that is itself considered to be an adhesive protein. For example, the protein bound could be a fibronectin or a laminin, which are complex, multifunctional proteins involved in both cell adhesion and anchorage to structural components of the extracellular matrix (2, 15, 17).

Consequently, there are two broad classes of adhesion molecules or receptors (Fig. 1). One class is bound to the plasma membrane, often as a transmembrane protein (Fig. 1). This type of molecule is generally a receptor, a homophilic adhesion molecule, or a counterreceptor. It often consists of an extracellular domain containing one or more cell-interaction domains or sites, as well as a stalk region, a hydrophobic transmembrane domain, and (usually) a cytoplasmic domain or tail. This type of molecule is often also involved in transmembrane transmission of signals.

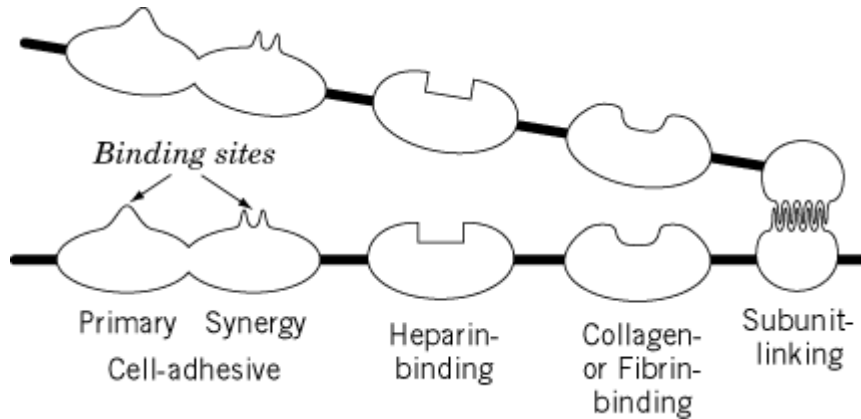
**Figure 1.** Schematic diagram of a generalized transmembrane adhesion molecule consisting of an extracellular domain, transmembrane segment, and cytoplasmic domain. See text for discussion.



The second broad class of adhesion molecules consists of proteins that are often classified as cell surface or extracellular matrix proteins (see [Extracellular Matrix](#)), but that contain domains or sites involved in cellular adhesion (Fig. 2). Molecules in this class include fibronectins, laminins, vitronectin, tenascins, thrombospondins, and the collagens. All contain one or more cell-binding domains, which consist of a primary recognition motif consisting of a short peptide sequence (eg, Arg—Gly—Asp), and sometimes a synergy site that provides a substantial increase in receptor-

binding specificity and affinity. As discussed below, many of these proteins also contain a variety of other functional domains.

**Figure 2.** Schematic summary of common features of extracellular matrix adhesion molecules. Linear arrays of structural and functional domains provide a series of specific binding sites. See text for further discussion.



In general, adhesion molecules frequently have the following properties (1-3, 6, 14-16, 18-25):

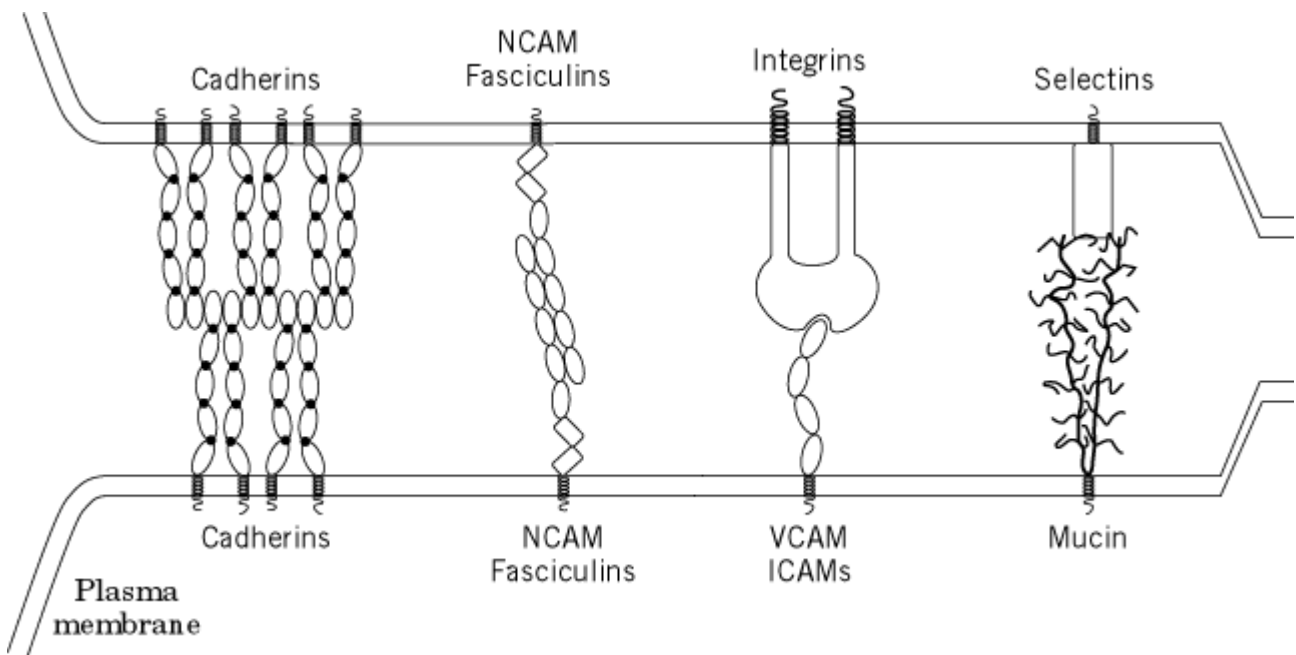
- Composed of multiple repeats of protein motifs, such as the immunoglobulin (Ig) motif, the epidermal growth factor (EGF) repeat, or the fibronectin motif.
- Specialized functional domains, including a domain for formation of dimers or higher polymers.
- Moderate affinity, for example, with dissociation constants ( $K_d$ ) in the range of  $10^{-6}$ – $10^{-7}M$  for fibronectin, and even as weak as  $10^{-4}M$  for leukocyte adhesion molecules that mediate rolling adhesion. In contrast, many well-known protein—protein interactions have affinities with  $K_d = 10^{-9}$ – $10^{-11} M$ .
- High avidity after clustering. Adhesion molecules often form functional clusters or aggregates in the plane of the plasma membrane, which causes them to develop strong total avidity due to the cooperation of the otherwise weak binding of individual molecules.
- Regulation by activation, such as the “inside-out” signaling that increases the affinity of integrins in platelet activation and leukocyte adhesion, or regulation by phosphorylation of the cadherin system.

### 3. Specific Cell Adhesion Molecules

The broad class of adhesion molecules embedded in the plasma membrane contains several large groups of proteins that share common structural motifs, especially the Ig motif (Fig. 3). Cadherins represent a large superfamily of proteins involved in homophilic cell-to-cell adhesive interactions (14-16) (see [Cadherins](#)). Cadherins bind to cadherins of the same type on other cells via cell interaction sites that can include the recognition sequence His—Ala—Val. Besides the “classic” cadherins, such as E-cadherin and N-cadherin, there are a number of other types of cadherins, as the highly specialized cadherins termed *desmocollins* and *desmogleins* found exclusively on desmosomes, which link cells together at these particularly strong attachment sites connected to intermediate filaments such as keratins or vimentin. Cadherins are quite sensitive to depletion of calcium in the surrounding medium, accounting for the ability of calcium chelators such as EDTA to dissociate tissues into component cells. Although desmosomal cadherins are present in epithelia of all ages, cadherins appear to be of particular importance during embryonic development, when they

mediate cell—cell adhesion and help define tissues by their propensity to bind primarily to cadherins of the same type, rather than to other cadherins on unrelated cell types (14-16). These specific adhesive interactions result in cell sorting to form tissue regions of closely related cells.

**Figure 3.** Examples of cell adhesion protein complexes linking plasma membranes of one cell (top) with another (bottom) are listed. Adhesion molecules comprising repeating immunoglobulin motifs include cadherins, NCAM, fasciculins, VC, the plasma membrane. Cadherins appear to exist as dimers that form lateral and intercellular connections. Selectins bind carbohydrates in certain mucins. CD2, CD48, and CD58 contain only two Ig repeats. Myelin protein zero ( $P_0$ ) forms hon spaced plasma membranes.



Integrins are major cell-surface receptors for a host of extracellular matrix proteins, as well as “counterreceptors” on other cells (8) (see [Integrins](#)). There are more than 20 types of integrin subunits. Genetic absence of most integrin subunits leads to disease or death, often in the embryo or near the time of birth. Integrins are heterodimers of one  $\alpha$  and one  $\beta$  subunit, each of which provides adhesive specificity. They have a highly distinctive appearance, with a bulbous head domain contributed by each subunit, a binding site that also appears to involve part of each subunit, and two spindly legs (Fig. 3). The legs penetrate through the plasma membrane and usually terminate in rather short cytoplasmic domains. Although short, these cytoplasmic tails can mediate a remarkable range of signaling events (see text below). Integrin functions are also often inhibited by depletion of divalent cations.

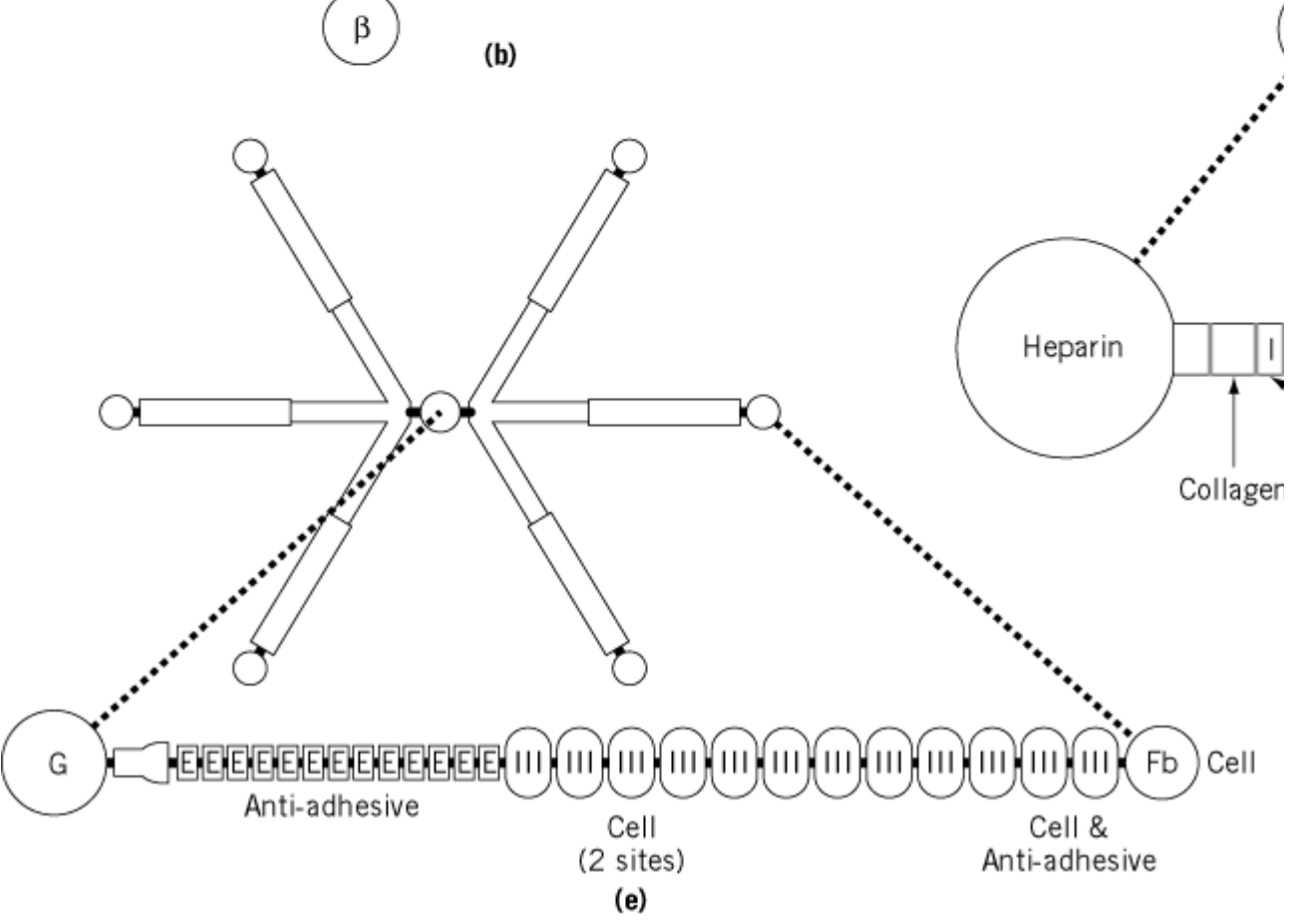
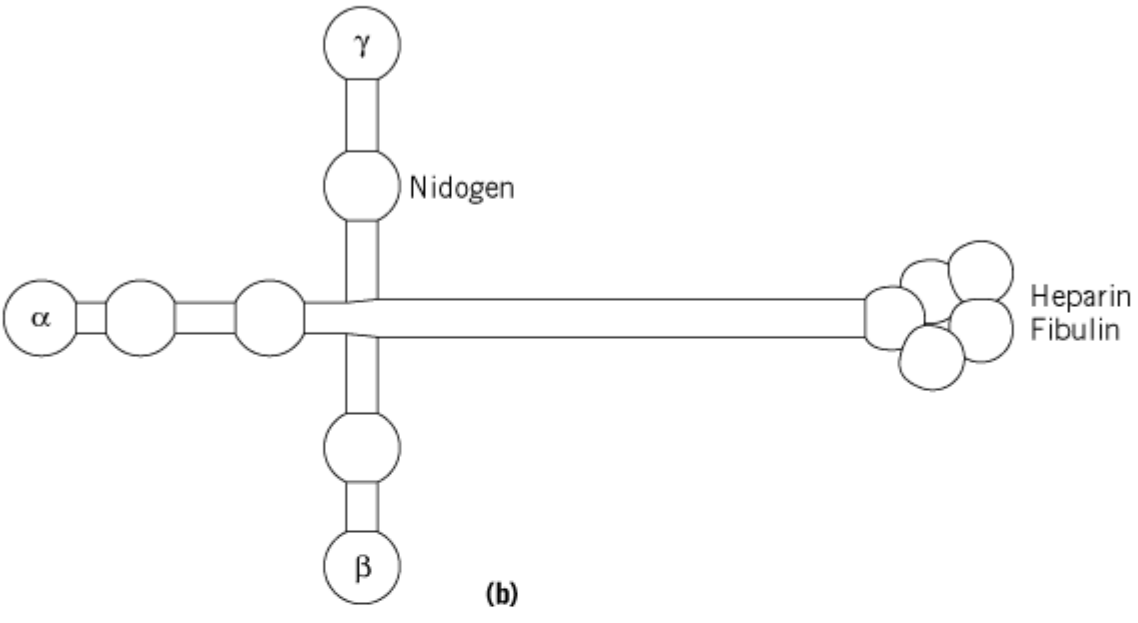
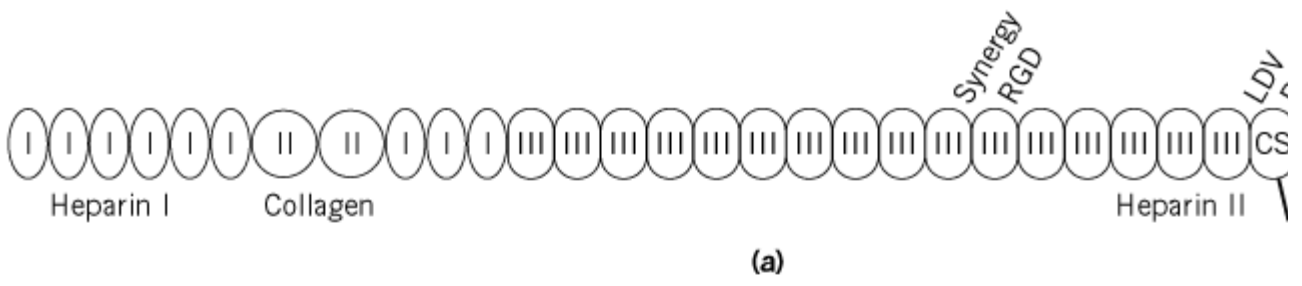
Although cell adhesion molecules (CAMs) and counterreceptors are structurally related by their use of the Ig motif, and they often even share the -CAM designation, they can differ functionally (Fig. 3). Molecules such as NCAM are homophilic adhesive molecules that bind to the same type of molecule on adhering cell surfaces (22). In contrast, counterreceptors such as the ICAMs and VCAM have specialized peptide recognition sites that are bound specifically by integrins such as LFA-1 (CD11a/CD18 or  $\alpha_1\beta_2$ ) or VLA-4 ( $\alpha_4\beta_1$ ). The functions and sites of expression of these molecules differ widely. For example, molecules such as NCAM are implicated along with cadherins in embryonic developmental events, such as axonal bundling and guidance, whereas the ICAM counterreceptors are present as targets for binding by cells in the blood circulation. Levels of counterreceptors on the cell surface can often also be regulated rapidly in response to cytokines.

A new family of transmembrane proteins termed *ADAMs*, which are membrane proteins with a disintegrin and a metalloproteinase domain (26), contain a proteinase-like domain, as well as an integrin recognition site in the disintegrin domain. A member of this family is present on sperm and may play a role in adhesive interactions, but the functions of other members of this rapidly growing family remain to be characterized. Syndecans are cell-surface heparan sulfate proteoglycans, with a protein core that crosses the plasma membrane and terminates in a cytoplasmic tail. Syndecans appear to function as “coreceptors,” mediating signaling in association with a primary adhesion molecule (27, 28). For example, syndecans bind to fibronectin at the same time that cell surface integrins bind to a cell-binding domain of fibronectin, and these binding partners cooperate during formation of the specialized adhesion sites known as focal contacts. Syndecans have several other signaling functions as well, and they can interact with both cytoskeletal proteins and protein kinase C (27, 28).

#### 4. Extracellular Adhesion Molecules

The second major class of adhesive molecules exists extracellularly, and its members are generally targets for cell-surface adhesion receptors such as the integrins. A gallery of such molecules is presented in Figure 4, showing functional binding sites for recognition by cell-surface receptors and domains for binding to a variety of ligands. Many contain one or more heparin-binding domains, which are used to bind to heparan sulfate proteoglycans in the extracellular matrix or in cell-surface syndecan molecules. Many also contain other types of binding domains for extracellular molecules, such as for collagen/gelatin, fibrin, entactin (nidogen), and fibulin (Figs. 2 and 4). Many of these molecules are both large and multimeric, with a specific domain for covalent cross-linking to other subunits. Finally, most of these proteins are either alternatively spliced or encoded by a set of closely related genes, thereby generating a number of isoforms. A number of reviews provide details about these complex proteins (see [Extracellular Matrix](#) and Refs. 2, 18-22).

**Figure 4.** Composite of six examples of cell adhesion molecules of the extracellular matrix: (a) Fibronectin, (b) Laminin, (c) Vitronectin, (d) Thrombospondin, (e) Tenascin, (f) Nidogen (entactin). These molecules vary widely in size and shape, and are organized into functional binding domains. Specific molecular or cell-binding sites are labeled.



5. Peptide Adhesive Recognition Sites

A striking characteristic of a number of cell adhesion molecules is their frequent use of short peptide recognition sites (Refs. [18](#), [23-25](#), [29-31](#) and references cited therein). Short peptides with these sequences can at least partially mimic the binding function of intact molecules to adhesion receptors. Examples include Arg–Gly–Asp (RGD) and Leu–Asp–Val (LDV), which can bind directly to certain integrin receptors. Nevertheless, binding can be 25- to 200-fold more active in larger peptides or proteins, such as in association with a “synergy” site that functions synergistically with the primary adhesive peptide site ([30](#)). Some peptide recognition sites appear to be “cryptic” in the intact molecule, and may be active only in small fragments of the protein, such as after proteolysis. Table [1](#) lists the putative peptide recognition sequences in cell adhesion molecules. It is important to stress that a number of adhesion molecules also function by adhering to other proteins over much broader intermolecular contact areas, via multiple molecular contacts and with no simple peptide motif, in analogy to classic high affinity noncovalent protein–protein interactions ([19](#)). It is nevertheless striking that so many cell adhesion molecules use simple, specific, peptide adhesive-recognition motifs.

**Table 1. Adhesive Recognition Sequences**

| <b>Protein and Peptide Name</b>        | <b>Recognition Sequence</b>   |
|--|-------------------------------|
| Fibronectin                            |                               |
| Cell-binding determinant               | <i>GRGDS</i>                  |
| Synergy site                           | <i>PHSRN</i>                  |
| Potential synergistic sites            | <i>RNS, SDV</i>               |
| CS1 site of III <sub>1</sub> CS domain | <i>DELPQLPHPNLHGPEILDVPS</i>  |
| CS5 site of III <sub>1</sub> CS domain | <i>REDV</i>                   |
| FN-C/H-I                               | <i>YEKPGSPPREVVPRPRPGV</i>    |
| FN-C/H-II                              | <i>KNNQKSEPLIGRKKK</i>        |
| FN-C/H-III                             | <i>YRYRYTPKEKTGPMKE</i>       |
| FN-C/H-IV                              | <i>SPPRRARVT</i>              |
| FN-C/H-V                               | <i>WQPPRARI</i>               |
| Type III <sub>5</sub> site             | <i>KLDAPT</i>                 |
| E1                                     | <i>TDIDAPS</i>                |
| Laminin $\alpha_1$ chain               |                               |
| PA22.2 (IKVAV)                         | <i>SRARKNAASIKVAVSADR</i>     |
| RGD site                               | <i>RGDN</i>                   |
| GD-1                                   | <i>KATPMLKMRTSFHGCIK</i>      |
| GD-2                                   | <i>KEGYKVRLDLNITLEFRITTSK</i> |
| GD-3                                   | <i>KNLEISRSTFDLLRNSYGVRK</i>  |
| GD-4                                   | <i>DGKWHTVKTEYIKRKAF</i>      |
| GD-6                                   | <i>KQNCLSSRASFRGCVRNLRLSR</i> |



|   |   |
|---|---|
| AG-10   | NRWHSIYITRFG                                |
| AG-32   | TWYKIAFQRNRK                                |
| AG-73   | RKRLQVQLSIRT                                |
| Laminin b <sub>1</sub> chain                        |   |
| YIGSR   | <i>YIGSRC</i>                               |
| PDSGR   | <i>PDSGR</i>                                |
| F9 (RYVVLPR)  | <i>RYVVLPRPVCFEKGKGMNYVR</i>                |
| LG TIPG   | <i>LG TIPG</i>                              |
| Integrin α <sub>2</sub> β <sub>1</sub> binding site | YGYYGDALR                                   |
| B-7   | AFGVLALWGTRV                                |
| B-133   | DSITKYFQMSLE                                |
| B-160   | VILQQSAADIAR                                |
| Laminin b <sub>2</sub> chain                        |   |
| LRE site  | <i>LRE</i>                                  |
| Laminin g <sub>1</sub> chain                        |   |
| P20   | RNIAEIIKDI                                  |
| C-16  | KAFDITYVRLKF                                |
| C-28  | TDIRVTLNRLNTF                               |
| C-64  | SETTVKYIFRLHE                               |
| C-68  | TSIKIRGTYSER                                |
| Fibrinogen  |   |
| RGD sites   | <i>RGDS</i> and <i>RGDF</i>                 |
| g chain peptide                                     | <i>HHLGGAKQAGDV</i>                         |
| Integrin α <sub>M</sub> β <sub>2</sub> binding site | <i>KRLDGGG</i>                              |
| von Willebrand factor                               |   |
| RGD site  | <i>RGDS</i>                                 |
| GPIb site   | CQEPGGLVVPPTDAP plus LCDLAPEAPPPTLPP        |
| Asp-514—Glu-542                                     | DLVFLLDGSSRLSEAEF- EVLKAFVVDemme            |
| Vitronectin   | <i>RGDV</i>                                 |
| Entactin/nidogen                                    | <i>SIGFRGDTC</i>                            |
| Fertilin  | <i>TDE</i>                                  |
| Circumzooite protein                                | <i>VTCG</i>                                 |
| Thrombospondin (CD36-binding motif)                 | <i>CSVTXG</i>                               |
| Heparin-binding motif                               | $\frac{1}{4}$ <i>YSXY</i>                   |
| N-Terminal domain site                              | <i>VDAVRTEKGFLLLASLRQMKKTRG-TLLALERKDHS</i> |
| GAG-independent site                                | <i>FQGVLQNVRFVF</i>                         |
| 4NIs  | <i>RFYVVMWK</i>                             |
| Collagen type I                                     |   |
| RGD sites   | <i>RGDTP</i> and <i>SRGDTG</i>              |

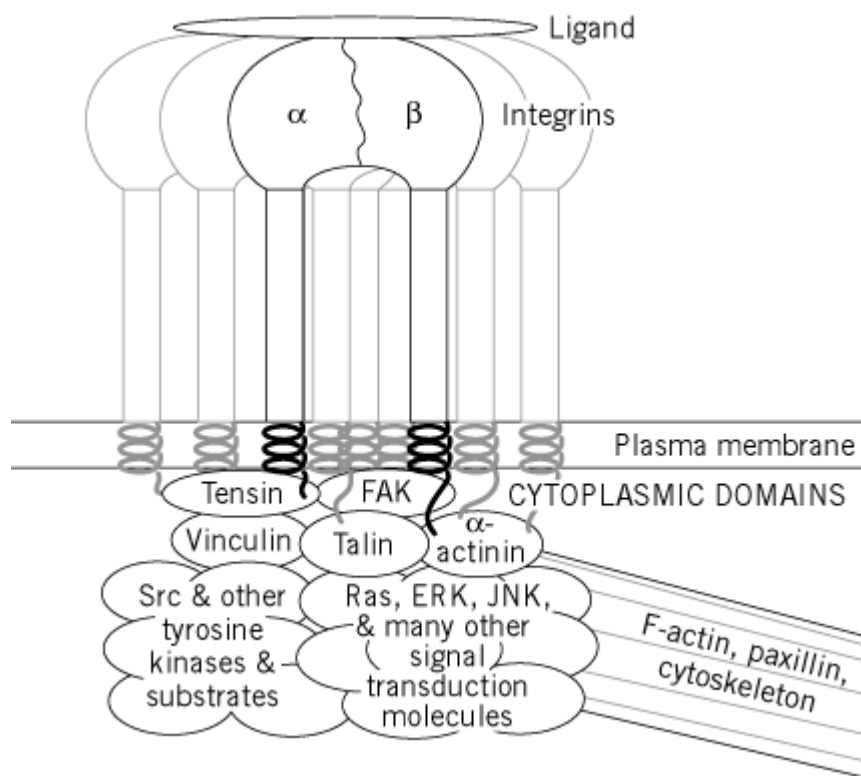
|                                       |   |
|---------------------------------------|---|
| DGEA site                             | <i>DGEA</i>   |
| Collagen type IV                      |   |
| IV-H1                                 | <i>GVKGDKGNPGWPGAP</i>  |
| Hep-I                                 | <i>TAGSCLRKFSTM</i>   |
| Hep III                               | <i>GEFYFDLRLKGDK</i>  |
| $\alpha_3$ (IV) 185–203               | <i>CNYYSNSYSFWLASLNPER</i>  |
| $\alpha_2\beta_1$ binding site        | <i>FYFDLR</i>   |
| Amyloid P component                   | <i>FTLCPR</i>   |
| Amyloid precursor protein             | <i>RHDS</i>   |
| Bone sialoprotein                     | <i>EPRGDNYR</i> (cyclic)  |
| L1                                    | <i>PSITWRGDGRDLQEL</i>  |
| ICAM-1                                |   |
| JF9                                   | <i>VLYGPRLDERDAPGNWTW- PENSQQTPMC</i>                                     |
| ICAM40-51                             | <i>KELLPGNNRKV</i>  |
| Cyclic peptide 1                      | <i>PSKVILPRGGC</i> (cyclic)   |
| $\alpha_L\beta_2$ binding site        | <i>LET, IET</i>   |
| LLG-C4                                | <i>CPCFLLGCC</i> (cyclic)   |
| ICAM-2                                | <i>GKSFTIECRVPTVEP</i>  |
| NCAM                                  | <i>KYSFNYDGSE</i>   |
| VCAM-1                                |   |
| Domain 1 C-D loop<br>(binding motifs) | <i>QIDSP</i> (cyclic peptide with terminal Cs)<br><i>IDSP, GNEH, KLEK</i> |
| <i>N</i> -Cadherin                    | <i>LRAHAVDVNG</i>   |
| Myelin protein 0 (P <sub>0</sub> )    | <i>YSDNGTF</i>  |
| Ninjurin                              | <i>PPRWGLRNRPIN</i>   |
| Leishmania pg63                       | $\frac{1}{4}$ <i>SRYD</i> $\frac{1}{4}$                                   |
| Streptavidin                          | <i>GRYDS</i> <sup>a</sup>   |

<sup>a</sup> See Refs. [14](#), [16](#), and [22](#) for additional information.

Two key sites of membrane adhesion molecules that allow them to function as signal transduction receptors are their ligand-binding domains and their cytoplasmic domains. The ligand-binding domains are of obvious importance in binding to extracellular molecules, but their roles are more intriguing in at least some cases. There appear to be separable functions for ligand occupancy, specifically, for filling the binding site with a ligand or an antibody, as opposed to receptor clustering (which can be induced by a multivalent ligand, such as fibrils of fibronectin or collagen). These two inputs can cooperate in promoting accumulation of specific cytoskeletal proteins, such as  $\alpha$ -actinin and actin, which are thought to be crucial for forming strong adhesions ([32](#)). Even though they lack intrinsic enzymatic activities, integrins appear to be able to function as signaling receptors and as regulators of actin cytoskeletal organization by recruiting other molecules to their cytoplasmic domains. Integrins can reportedly bind directly to certain cytoplasmic proteins directly, such as to talin,  $\alpha$ -actinin, and focal adhesion kinase, perhaps regulated by ligand occupancy on the outside of the cell. Integrin clustering appears to play a central role in forming large intracytoplasmic

complexes of over 30 molecules that can serve as signaling centers, such as for MAP kinase activation (Fig. 5; Refs. 9-13, 32).

**Figure 5.** Integrin adhesion and signaling complex formed after cell adhesion molecule interactions. Binding of a ligand and aggregation of adhesion receptors in the plane of the plasma membrane results in the formation of multimolecular complexes that can contain over 30 cytoskeletal and signal transduction molecules, which link to the actin-containing cytoskeleton and mediate intracellular signaling.



## 6. Major Functions of Cell Adhesion Molecules Beyond Adhesiveness

Cell adhesion molecules and receptors can be intuitively understood as mechanisms for attaching cells to other cells, for permitting cell adherence to extracellular matrix molecules, and for mediating traction during cell migration. It has become clear, however, that they also play critical roles in cellular signaling. In fact, their signaling functions may be of equal biological importance. A current view is that “adhesion” proteins are actually “cell-interaction” proteins that have multiple functions involving the bidirectional transfer of information at the cell surface. In fact, some of these proteins, such as thrombospondin and tenascin, can have antiadhesive activities in certain situations and can modulate a host of cellular functions. Even classic “adhesive” proteins, such as fibronectin, can have a bewildering range of activities, mediated by integrin binding, that range from activating or modulating nearly every known mammalian signal transduction pathway. These pathways include tyrosine phosphorylation by tyrosine kinases, mitogen-activated protein kinases (MAP kinases),  $\text{Ca}^{2+}$  and  $\text{H}^+$  fluxes, inositol phosphate pathways, and protein kinase C, which prevent apoptosis and activate specific gene transcription (9-12).

This same theme of complex formation, binding of cytoskeletal proteins such as actin, and accumulation of signaling molecules also appears to occur for cell-to-cell adhesion molecules (14-16, 22). Cadherin cytoplasmic domains interact with catenins, which then link to the actin cytoskeleton. In addition, however, the binding of  $\beta$ -catenin to cadherins can titrate its cytoplasmic

levels; b-catenin molecules that are not bound by cadherins can enter the nucleus and regulate gene expression (3). A variety of other signaling pathways exist in cell-cell and cell-matrix adhesion systems. Although much more remains to be learned about precise pathways and overall integration of the many cytoplasmic effects of cell adhesion molecules, it is clear that adhesion molecules are crucial components of the basic regulatory mechanisms of cells and provide dynamic links to the external environment.

## Bibliography

1. B. Alberts et al., in *Molecular Biology of the Cell*, 3rd, ed., Garland Publishing, New York, 1994, pp. 949–1009.
2. E.D. Hay, ed., *Cell Biology of Extracellular Matrix*, Plenum Press, New York 1991.
3. B.M. Gumbiner, *Cell* **84**, 345–357 (1996).
4. S.F. Gilbert, *Developmental Biology*, Sinauer Associates, Sunderland, Mass., 1994.
5. R.O. Hynes and Q. Zhao, *J. Cell Biol.* **150**, F89–96 (2000).
6. G.M. Edelman and J.P. Thiery, *The Cell in Contact: Adhesions and Junctions as Morphogenetic Determinants*, John Wiley & Sons, New York, 1985.
7. D.L. Benson, L.M. Schnapp, L. Shapiro and G.W. Huntley, *Trends Cell Biol.* **10**, 473–482 (2000).
8. R.O. Hynes, *Cell* **69**, 11–25 (1992).
9. E.A. Clark and J.S. Brugge., *Science* **268**, 233–239 (1995).
10. M.A. Schwartz, M.D. Schaller, and M.H. Ginsberg, *Annu. Rev. Cell Devel. Biol.* **11**, 549–599 (1995).
11. A.E. Aplin, A.K. Howe, and R.L. Juliano, *Curr. Opin. Cell Biol.* **11**, 737–744 (1999).
12. F.G. Giancotti and E. Ruoslahti, *Science* **285**, 1028–1032 (1999).
13. K.M. Yamada and B. Geiger, *Curr. Opin. Cell Biol.* **9**, 76–85 (1997).
14. M. Takeichi, *Annu. Rev. Biochem.* **59**, 237–252 (1990).
15. B.M. Gumbiner, *J. Cell Biol.* **148**, 399–404 (2000).
16. U. Tepass et al., *Nat. Rev. Mol. Cell Biol.* **1**, 91–100 (2000).
17. S.K. Sastry and K. Burridge, *Exp. Cell Res.* **261**, 25–36 (2000).
18. K.M. Yamada, *J. Biol. Chem.* **266**, 12809–12812 (1991).
19. C. Chothia and E.Y. Jones, *Annu. Rev. Biochem.* **66**, 823–862 (1997).
20. R. Gonzalez-Amaro and F. Sanchez-Madrid, *Crit. Rev. Immunol.* **19**, 389–429 (1999).
21. M.J. Humphries, *Biochem. Soc. Trans.* **28**, 311–339 (2000).
22. H. Kamiguchi and V. Leemmon, *Curr. Opin. Cell Biol.* **12**, 598–605 (2000).
23. S. Ayad, *The Extracellular Matrix Factsbook*, 2nd ed., Academic Press, San Diego 1998.
24. T. Kreis and R. Vale, *Guidebook to the Extracellular Matrix, Anchor, and Adhesion Proteins*, Oxford University Press, New York 1999.
25. C. Isacke and M. A. Horton, *The Adhesion Molecule Factsbook*, Academic Press, San Diego 2000.
26. J. Schlondorff and C.P. Blobel, *J. Cell Sci.* **112**, 3603–3617 (1999).
27. J.R. Couchman and A. Woods, *J. Cell Sci.* **112**, 3415–3420 (1999).
28. A.C. Rapraeger, *J. Cell Biol.* **149**, 995–998 (2000).
29. E. Ruoslahti, *Annu. Rev. Cell Devel. Biol.* **12**, 697–715 (1996).
30. S. Aota, M. Nomizu, and K. M. Yamada, *J. Biol. Chem.* **269**, 24756–24761 (1994).
31. M. Nomizu et al., *Arch. Biochem. Biophys.* **378**, 311–320 (2000).
32. K.M. Yamada and E.H.J. Danen, in J. S. Gutkind, ed., *Signaling Networks and Cell Cycle Control*, Humana Press, Totowa, NJ, 2000, pp. 1–25.

### Additional Reading

33. Alberts B. et al., "Cell junctions, cell adhesion, and the extracellular matrix", in *Molecular Biology of the Cell*, 3rd ed., Garland Publishing, Inc., New York, pp. 949–1009 1994.
34. Hay E.D., ed., *Cell Biology of Extracellular Matrix*, Plenum Press, New York 1991.
35. Chothia C. and Jones E.Y., *Annu. Rev. Biochem.* **66**, 823–862 (1997).
36. Gumbiner B.M., *Cell* **84**, 345–357 (1996).

Another excellent source of current information on adhesion molecules is found each year in issue number 5 of *Current Opinion in Cell Biology*. These annual issues are devoted to “cell-to-cell contact and extracellular matrix,” and provide useful updates in the field.

## Cell Cycle

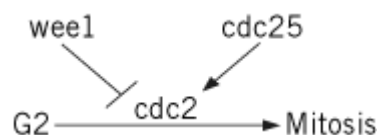
A dramatic change in how we understand the progression of the eukaryotic cell cycle, from phenomenological to biochemical, occurred in the late 1980s. The cell cycle consists of four main phases: G1, S, G2, and M. [DNA replication](#) occurs in S phase, and the [chromosomes](#) are segregated to daughter cells during mitosis (M phase). The so-called gap phases (G1 and G2) are defined simply as the periods separating DNA replication and mitosis. The length of the cell cycle ranges from just a few minutes in certain early embryos (which don't need to increase their mass between divisions), to 1.5 h in budding yeast, to approximately one day in typical mammalian [tissue culture](#) cells. In mammalian cells, the durations of S phase (~7 h) and mitosis (~1 h) are relatively inflexible, whereas the lengths of the gap phases can vary greatly. The length of G1 is particularly variable and is sensitive to the levels of nutrients and growth factors. In general, once a cell progresses past the “Restriction Point” in G1, it becomes committed to complete the cell cycle. A convergence of research in yeast, marine invertebrates, and frogs led to a phenomenal growth in our mechanistic understanding of how all eukaryotic cells regulate these events and the transitions between them. The basic machinery controlling the cell cycle consists of a subfamily of protein [kinases](#), termed cyclin-dependent kinases (Cdk's). These enzymes are in turn regulated by a number of mechanisms, including [transcription](#), inhibitory and activating [phosphorylations](#), binding to inhibitory proteins, and [ubiquitin](#)-mediated protein degradation. Checkpoints that prevent the execution of one cellular event before a prior step has been completed further modulate cell cycle progression.

### 1. Historical Threads: Yeast Genetics, MPF, and Cyclins

Genetic studies of the cell cycle in the yeasts *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe* illuminated important principles underlying the logic of cell cycle control and identified many of its key players. Screens were performed for conditional [temperature-sensitive mutations](#), causing cells to arrest quickly with uniform morphologies, indicative of defects in executing individual cell cycle steps. Analysis of a large collection of such *cdc* (cell division cycle) mutants from *S. cerevisiae* placed them into series of dependent and independent pathways (1). For example, inactivation of *CDC28* causes cells to arrest at “START” (equivalent to the restriction point of mammalian cells) and blocks pathways leading to formation of the bud, DNA replication, and duplication of the spindle pole body. Inactivation of genes acting after *CDC28* can block one of these pathways, while the other two continue. A similar screen in the distantly related yeast *Schizosaccharomyces pombe* successfully identified the most proximal regulators of the G2 to mitosis transition (Fig. 1) (2). *cdc2<sup>+</sup>* was found to be the key regulator of entry into mitosis and to encode a protein kinase. Inactivation of *cdc2<sup>+</sup>* causes cells to arrest in G2, whereas its premature

activation leads to early entry into mitosis. Genetic analysis indicated that the *cdc25*<sup>+</sup> gene product functioned in a pathway leading to activation of Cdc2 and that the product of *wee1*<sup>+</sup> functioned in a pathway leading to inhibition of Cdc2. Sequence analysis of these genes indicated that both *cdc2*<sup>+</sup> and *wee1*<sup>+</sup> encoded protein kinases, but initially shed little light on the function of Cdc25 (see text below).

**Figure 1.** Regulation of the G2-to-mitosis transition in *Schizosaccharomyces pombe*. *cdc2*<sup>+</sup> is required for this transition. In the absence of its activity, cells arrest in G2. *cdc25*<sup>+</sup> is a positive regulator of *cdc2*<sup>+</sup>, and its inactivation also gives a G2 arrest. In contrast, *wee1*<sup>+</sup> negatively regulates *cdc2*<sup>+</sup>, and its inactivation leads to premature activation of Cdc2 and entry of cells into mitosis at an unusually small size (a “wee” phenotype).



Two other lines of research provide essential historical threads leading to our current understanding of the cell cycle. The first involves [maturation promoting factor](#) (MPF), which was originally defined as a developmental factor found in the cytoplasm of metaphase-arrested frog eggs that could induce oocytes to proceed through meiosis and to “mature” into eggs following its injection into their cytoplasm (3). Subsequent research found that MPF activity is present in all eukaryotic cells during M phases (meiosis or mitosis). Thus, MPF is a universal regulator of entry into mitosis. The second thread concerns the [cyclins](#). The first cyclins were found during studies of translational control before and after fertilization of sea urchin eggs conducted as part of the Physiology course at the Marine Biological Laboratory in Woods Hole (4). These proteins were synthesized continuously, and they accumulated until their abrupt degradation during mitosis. This pattern of accumulation hinted that cyclins might play an important role during the cell cycle. Later work showed that the injection of cyclin [messenger RNA](#) led to the maturation of frog oocytes, suggesting that cyclins, like MPF, functioned as positive regulators of the G2-to-mitosis transition (5).

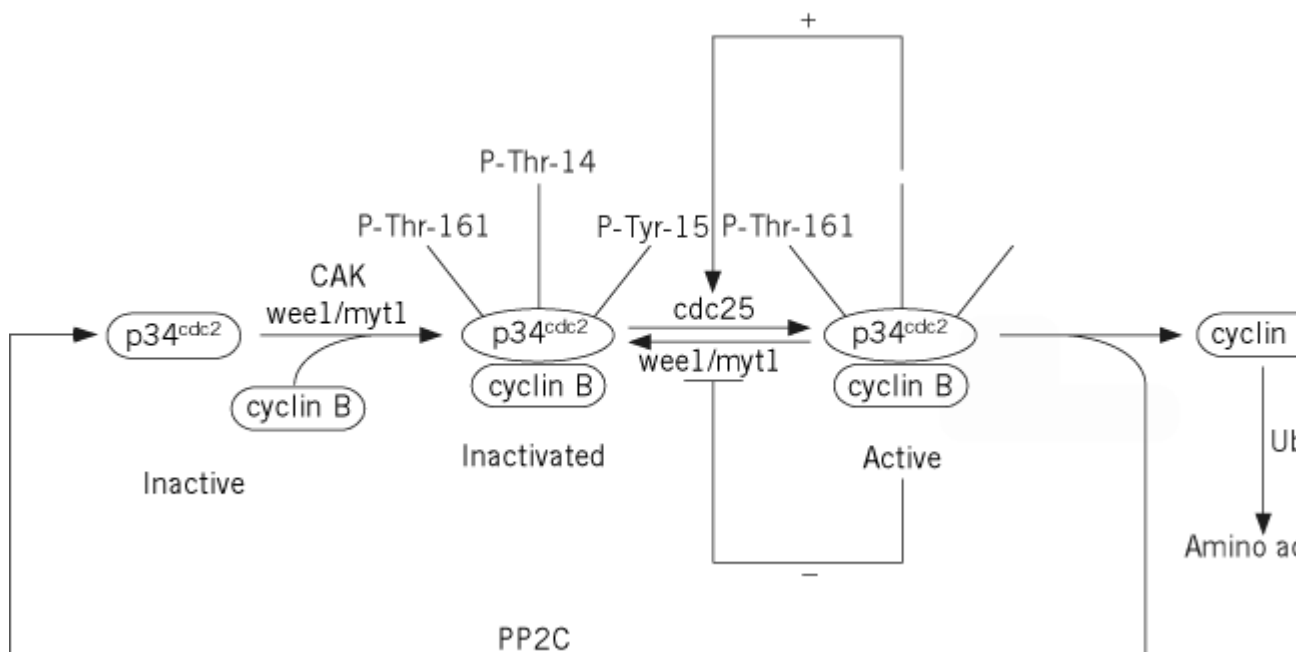
Work on cyclins, MPF, and the yeast *cdc* genes came to an explosive convergence in the late 1980s with the purification of MPF and the identification of its subunits as homologs of *cdc2*<sup>+</sup> and cyclin (6-9). Thus, this key regulator of entry into mitosis was a protein kinase (Cdc2) whose activity was controlled at least in part via the cyclic accumulation and degradation of a regulatory subunit (cyclin). The activity of Cdc2 is typically assayed by its ability to **phosphorylate** a convenient substrate, [histone](#) H1. Cdc2 is a workhorse master regulator. Rather than sitting at the apex of a large regulatory pathway, Cdc2 generally phosphorylates the downstream-most targets of mitotic regulation. For example, its direct phosphorylation of nuclear lamins causes the disassembly of the lamin-containing filaments of the nuclear lamina (10, 11).

## 2. Regulation of Cdc2 Activity: Phosphorylation and Cyclin Degradation

In addition to cyclin binding, enzymes such as Cdc2 are regulated by multiple phosphorylations (12, 13) (see Fig. 2). Cdc2 is negatively regulated via phosphorylation of Thr14 and Tyr15. (The amino acid positions in this section refer to the phosphorylation sites in human Cdc2.) These sites are phosphorylated by Wee1 (the negative regulator of Cdc2 identified genetically in *Schizosaccharomyces pombe* (see Fig. 1)) and by Wee1-like protein kinases (Mik1 in *Schizosaccharomyces pombe* and the membrane-bound Myt1 protein in vertebrates). Dephosphorylation of these inhibitory sites is carried out by Cdc25 proteins, genetically identified activators of Cdc2 that share distant similarity to protein **tyrosine phosphatases**. Finally, activating

phosphorylation of Cdc2 occurs on Thr161 within the so-called T-loop and is carried out by the Cdk-activating kinase (CAK). Dephosphorylation of Thr161 has not been as well studied. It occurs following cyclin degradation, possibly by a phosphatase termed KAP. Two feedback loops operate as cells enter mitosis and lead to the inhibition of Wee1-like enzymes and the stimulation of Cdc25, thus ensuring the abrupt and irreversible transition from G2 to mitosis. The roles of cyclin binding and activating phosphorylation have been studied using [X-ray crystallography](#) of Cdk2 (14-16), a close relative of Cdc2. Monomeric Cdk2 is inactive for two major reasons. First, residues that interact with the phosphates of ATP are out of position, resulting in the misalignment of the  $\gamma$ -phosphate of ATP for phosphoryl transfer. Second, the T-loop is positioned so that it would physically prevent protein substrates from gaining access to the site of catalysis. Binding to cyclin relieves this steric block and moves residues that interact with ATP into their proper positions. The cyclin A-Cdk2 complex is about 1% active. Full activation requires activating phosphorylation (on Thr160 in Cdk2), which alters the positions of some T-loop residues and creates an acidic patch for the binding of optimal substrates containing a key basic residue.

**Figure 2.** Activation of Cdc2. Monomeric Cdc2 is unphosphorylated and inactive. Binding to cyclin induces its phosphorylation on inhibitory sites (Thr14 and Tyr15) by Wee1-like protein kinases and its phosphorylation on an active site (Thr161) by CAK. The transition into mitosis is accompanied by two feedback loops involving stimulation of Cdc25 (the phosphatase acting on Thr14 and Tyr15) and inhibition of Wee1-like enzymes, resulting in the abrupt and irreversible activation of Cdc2 and entry into mitosis. The degradation of cyclin by the ubiquitin system toward the end of mitosis leads to the inactivation of Cdc2 and its dephosphorylation, probably by a phosphatase called KAP.



One of the irreversible ratchet steps in the cell cycle is the degradation of cyclins by the ubiquitin system (17). Ubiquitin is a 76-residue protein whose covalent attachment to proteins can target them for **proteolysis** by the [proteasome](#), a huge multiproteinase unwinding and degrading machine. Ubiquitin is activated by its ATP-dependent covalent attachment to an enzyme called E1. E1 then transfers ubiquitin to one of many E2 enzymes. The E2s can ubiquitinate substrates, often with the help of an E3. The E3s may be the most diverse and interesting components of this system. Some E3s receive the covalently bound ubiquitin; others serve as matchmakers that bring together a substrate and the appropriate E2. For the mitotic cyclins, the E3 was first termed the cyclosome, though it is now generally termed the *anaphase promoting complex* (APC), an 8- to 12-subunit complex. The APC is the regulated component of the ubiquitin system for the degradation of the

mitotic cyclins and serves as the target for checkpoint signals (see below) that can block cyclin degradation. The ubiquitin-dependent degradation of cyclins that act earlier in the cell cycle generally does not require the APC and has been best studied in *S. cerevisiae*.

### 3. Relatives: Cdks and Cyclins

Both Cdc2 and the mitotic cyclins are members of large families of related proteins, many of which function in cell cycle regulation. The Cdc2 relatives are termed Cdks (for cyclin-dependent kinases) and the cyclin relatives are designated by letter. The different cyclins accumulate at various times during the cell cycle and pair with a limited subset of the Cdks to execute their respective functions. The major cyclin–Cdk complexes and their functions are shown in Table 1. Complexes of D-type cyclins with Cdk4 and Cdk6 function in progression through G1 and are the first kinases to phosphorylate the cell-cycle inhibitor retinoblastoma protein (Rb). The cyclin E–Cdk2 complex stimulates the initiation of DNA synthesis, and the cyclin A–Cdk2 complex is generally thought to promote passage through S phase. Cyclin B partners with Cdc2 to promote the G2-to-mitosis transition. All of these complexes are regulated by phosphorylations on sites equivalent to those used in Cdc2. In addition to these enzymes, a large number of cyclin–Cdk complexes function in processes other than cell cycle regulation, including transcription and neuronal differentiation. In general, the regulatory subunits of these Cdks are stable proteins and are classified as cyclins only on the basis of sequence similarity.

**Table 1. Functions of Cyclin–Cdk Complexes Involved in Cell Cycle Control**

| Catalytic Subunit | Major Cyclin Partner | Function |
|-------------------|----------------------|----------|
| Cdc2              | Cyclin B             | G2 to M  |
| Cdk2              | Cyclin A             | S        |
| Cdk2              | Cyclin E             | G1 to S  |
| Cdk4              | Cyclin D             | G1       |
| Cdk6              | Cyclin D             | G1       |

### 4. Cdk Inhibitors

In addition to cyclin binding and phosphorylation, Cdks are also regulated by the binding of inhibitory proteins (18). These can be divided into two classes: the CIP/KIP family (consisting of p21<sup>CIP1</sup>, p27<sup>KIP1</sup>, and p27<sup>KIP2</sup>) and the INK4 family (consisting of p15<sup>INK4b</sup>, p16<sup>INK4a</sup>, p18<sup>INK4c</sup>, and p19<sup>INK4d</sup>). The CIP/KIP proteins bind preferentially to the cyclin–Cdk complex and cause its inactivation. These proteins bind to numerous Cdks, but show a preference for Cdk2. In contrast, the INK4 proteins bind preferentially to monomeric Cdks (inhibiting the binding of cyclin) and are specific for Cdk4 and Cdk6. Cdk inhibitors play very important roles both in cell cycle progression (particularly in responses to [growth factors](#)) and during development. One of the best-studied roles for a Cdk inhibitor involves the regulation of the initiation of DNA synthesis in *S. cerevisiae* by Sic1p. Entry into S phase in *S. cerevisiae* requires the activity of complexes containing an S-phase cyclin (Clb5p or Clb6p) and the budding yeast Cdk, Cdc28p. These complexes are inhibited by Sic1p. Degradation of Sic1p by the ubiquitin system is triggered following the phosphorylation of Sic1p by Cdc28p bound to G1 cyclins, thereby freeing the Clb5p and Clb6p complexes to initiation DNA replication (19).



## 5. Checkpoints

Among the key mechanisms for ensuring the orderly progression of cell cycle events are checkpoints, which monitor the completion of certain events and prevent the execution of subsequent events if the monitored event has not been completed (20, 21). For example, checkpoints can block entry into mitosis if the DNA is damaged or has not been fully replicated, or they can block exit from mitosis if the spindle has not been properly assembled. In *S. cerevisiae* a checkpoint even monitors whether a bud has formed. One of the defining features of a checkpoint sensor is that it is not part of the machinery being sensed and that it is dispensable, though only for the short term. The checkpoint that provides the paradigm for all others involves the sensing of **DNA damage** and unreplicated DNA by Rad9p in *S. cerevisiae*. A number of yeast *cdc* mutants cause G2 arrests that are dependent on this checkpoint. For example, inactivation of *CDC9*, which encodes [DNA Ligase](#), leads to G2 arrest if *RAD9* is functional. In contrast, inactivation of *CDC9* in a *rad9* mutant does not lead to an immediate arrest. Instead, cells progress into mitosis and continue through a couple of cell cycles, until they arrest as microcolonies. The G2-arrested *RAD9* cells remain viable while they attempt to repair their DNA, whereas the *rad9* cells rapidly lose viability as they proceed through the checkpoint. Unperturbed wild-type cells grow normally in the absence of this checkpoint. The importance of checkpoints to a proper cell cycle and to human health is underscored by the [p53](#) protein. p53 is the most frequently mutated protein in human cancers and functions as part of a checkpoint responding to DNA damage by transcriptionally inducing the p21 Cdk inhibitor (22), thereby blocking cell cycle progression. The absence of this checkpoint leads to genomic instability (the accumulation of further mutations, loss of chromosomes, etc.) and a greatly increased risk of tumor formation.

### Bibliography

1. J. R. Pringle and L. H. Hartwell (1981) In *The Molecular Biology of the Yeast Saccharomyces* (J. Strathern, E. Jones, and J. Broach, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 97–142.
2. P. Nurse (1990) *Nature* **344**, 503–508.
3. Y. Masui and C. L. Markert (1971) *J. Exp. Zool.* **177**, 129–146.
4. T. Evans, E. T. Rosenthal, J. Youngblom, D. Distel, and T. Hunt (1983) *Cell* **33**, 389–396.
5. K. I. Swenson, K. M. Farrell, and J. V. Ruderman (1986) *Cell* **47**, 861–870.
6. W. G. Dunphy, L. Brizuela, D. Beach, and J. Newport (1988) *Cell* **54**, 423–431.
7. J. Gautier, C. Norbury, M. Lohka, P. Nurse, and J. Maller (1988) *Cell* **54**, 433–439.
8. J. C. Labbé, J. P. Capony, D. Caput, J. C. Cavadore, J. Derancourt, M. Kaghad, J. M. Lelias, A. Picard, and M. Dorée (1989) *EMBO J.* **8**, 3053–3058.
9. J. Gautier, J. Minshull, M. Lohka, M. Glotzer, T. Hunt, and J. L. Maller (1990) *Cell* **60**, 487–494.
10. M. Peter, J. Nakagawa, M. Dorée, J. C. Labbé, and E. A. Nigg (1990) *Cell* **61**, 591–602.
11. G. E. Ward and M. W. Kirschner (1990) *Cell* **61**, 561–577.
12. M. J. Solomon (1993) *Curr. Opin. Cell Biol.* **5**, 180–186.
13. M. J. Solomon and P. Kaldis (1998) In *Results and Problems in Cell Differentiation*, Vol. **22**: Cell Cycle Control (M. Pagano, ed.), Springer, Heidelberg, pp. 79–109.
14. J. L. DeBodt, J. Rosenblatt, J. Jancarik, H. D. Jones, D. O. Morgan, and S.-H. Kim (1993) *Nature* **363**, 595–602.
15. P. D. Jeffrey, A. A. Russo, K. Polyak, E. Gibbs, J. Hurwitz, J. Massagué, and N. P. Pavletich (1995) *Nature* **376**, 313–320.
16. A. A. Russo, P. D. Jeffrey, and N. P. Pavletich (1996) *Nat. Struct. Biol.* **3**, 696–700.
17. R. W. King, R. J. Deshaies, J.-M. Peters, and M. W. Kirschner (1996) *Science* **274**, 1652–1659.
18. T. J. Soos, M. Park, H. Kiyokawa, and A. Koff (1998) In *Results and Problems in Cell*

*Differentiation*, Vol. **22**: Cell Cycle Control (M. Pagano, ed.), Springer, Heidelberg, pp. 111–131.

19. R. M. Feldman, C. C. Correll, K. B. Kaplan, and R. J. Deshaies (1997) *Cell* **91**, 221–230.
20. L. H. Hartwell and T. A. Weinert (1989) *Science* **246**, 629–634.
21. S. J. Elledge (1996) *Science* **274**, 1664–1672.
22. J. W. Harper, G. R. Adami, N. Wei, K. Keyomarsi, and S. J. Elledge (1993) *Cell* **75**, 805–816.

### **Suggestions for Further Reading**

23. W. G. Dunphy, ed. (1997) *Methods in Enzymology*, Vol. **283**: Cell Cycle Control, Academic Press, San Diego, CA.
24. D. O. Morgan (1997) Cyclin-dependent kinases: engines, clocks, and microprocessors. *Annu. Rev. Cell Dev. Biol.* **13**, 261–291.
25. A. Murray and T. Hunt (1993) *The Cell Cycle*, W. H. Freeman, New York.
26. M. Pagano, ed. (1998) *Results and Problems in Cell Differentiation*, Vol. **22**: Cell Cycle Control, Springer, Heidelberg.

## **Cell Death**

It is not initially obvious why cell death should be considered important. The death of single-celled organisms, such as **amoeba** or **bacteria**, has no apparent advantage for the individual cell. However, in multicellular organisms damage to an individual cell can have repercussions for the whole animal. It therefore makes biological sense to dispose of this cell and replace it with a healthy one. Cells become damaged in a variety of different ways, but for long-lived, multicellular animals, the greatest risk is from a damaged cell that has acquired an **oncogenic** mutation leading to unrestricted clonal growth. To minimize this risk, several pathways have evolved that limit the expansion of somatic cells, one of which is cell death.

Tissue size is restricted by several different methods, not all of which actually require the physical loss of the cell. The problem facing multicellular organisms is how to maintain the capacity to divide and replace damaged cells while minimizing the risk of mutation leading to enhanced growth potential. One way around this problem is to make tissues post-replicative, as occurs in the central nervous system (CNS). All the neurons within the CNS are produced during embryogenesis and are then maintained throughout the life span of the organism. For some tissues, however, such an approach is not practical due to persistent physical damage to the cells, such as occurs in the epithelial lining of the gut or skin. Here, cells are continuously produced by a select number of proliferating **stem cells**. The daughter cells they produce move upward, away from the basement membrane, terminally differentiate, and are eventually shed. Because these cells are continuously replaced, any damaged cell should be automatically eliminated. Tissue size is also restricted by the need for vascularization. If there is no blood supply, only a small proportion of cells can be maintained by the diffusion of solutes, which is one of the reasons why many solid tumors have central necrotic zones.

For tissues that require a significant number of cells in cycle, the risk of mutation is minimized by both **senescence** and cell death. Senescence is a pathway that is invoked once cells have undergone a specific number of doublings, causing them to arrest permanently. Cells in this senescent state are unable to reenter the cell cycle (**1**). Thus, a cell that does acquire a mutation allowing constant

proliferation is limited in its capacity to divide, restricting clonal outgrowth [see [Senescence](#)]. Finally, cells can be triggered to die when no longer required, or when damaged. This “programmed cell death” is essential for the maintenance of tissue homeostasis and is invoked in disparate situations including DNA damage, absence of survival signals, and oncogene activation (2, 3) [see [Programmed Cell Death](#) and [Apoptosis](#)]. Thus, multicellular organisms are protected against somatic mutation by a number of independent mechanisms that act in concert to limit the potential for [neoplastic transformation](#).

## 1. The Meaning of Cell Death

The concept that cell death is fundamentally important for restricting cell population expansion took some time to be accepted for several reasons. Dead cells are not obvious in healthy tissues, whereas they are abundant in areas of damaged tissue, such as ischemic heart tissue resulting from myocardial infarction, allowing dead cells to be fundamentally associated with disease. Only during embryogenesis are dying cells seen in abundance, but these were generally regarded as cells deleted as a result of overproduction during development. In addition, death is seen as a bad outcome in human terms, a misconception that has allowed human anthropomorphic confusions to outweigh scientific evidence. More recently, however, several of these objections have been overruled, allowing cell death to become an accepted, essential daily process.

Early investigators of both invertebrate and vertebrate development observed that developmental cell deaths occurred in response to several biological cues and could be suppressed by inhibitors of both protein and RNA synthesis (4, 5), suggesting a requirement for macromolecular synthesis (6, 7). Moreover, they made the important connection that these cell deaths were an essential part of the developmental program of the organism concerned; hence the term “programmed cell death” (PCD) (8). It is now possible in less complex invertebrate models, such as the [nematode](#) *Caenorhabditis elegans*, to map both the fate of all individual cells within the developing organism and the genes that dictate these fates (9) [see [Programmed Cell Death](#)].

For vertebrates, the importance of cell death was not appreciated until the detailed characterization of the form of cell death termed [apoptosis](#) (10). The programmed cell deaths observed during [development](#) are identical to apoptotic cells in mammalian tissues, suggesting that a regulated form of cell death has been conserved throughout [evolution](#) (11, 12). From these early observations, our understanding of apoptosis/PCD has grown to include a complex pathway that is both morphologically and biochemically distinct from classical necrosis or accidental cell death. [Necrosis](#) is a passive event, in which cells that become irreversibly damaged, and therefore useless, die (13, 14). This form of cell death, which requires no input from the dying cell, occurs in cells subject to physical disruption or severe toxic stress. Conversely, apoptosis/PCD describes a pathway of events in which the cell is actively involved. Due to the active nature of apoptosis, which is triggered by many physiological and toxic stimuli, this form of cell death is sometimes referred to as “cell suicide.” Overall, the importance of “active” cell death is underlined by the fact that death is the default state for all cells. Hence, all cells are programmed to die, unless signaled to survive (15).

The importance of cell death in the regulation of tissue homeostasis is paramount. The number of proliferating cells determines the cell population number, as does the number of differentiating cells and the number of dying cells. Research over the last 20 years has graphically shown how cells that have mutations which disrupt their capacity to undergo apoptosis in response to physiological stimuli are involved in the etiology of several diseases, including Alzheimer's, AIDS, and cancer (16).

Aside from the three pathways of death described above, the “death” of a cell does not always result in the loss of viability or cellular function. Within a cell population, few cells are actually proliferating, many are in a resting state outside of the [cell cycle](#), known as  $GM_0$  phase. Certain cells can reenter the cell cycle given the appropriate stimuli, but cells that enter a permanent  $G_0$  state do not respond to these proliferation signals. Cells in the latter state are termed senescent and can

complete all normal functions except division (1) [see [Senescence](#)]. Cells that escape the confines of senescence are able to proliferate continuously and, more important, the mutations that lead to the evasion of senescence are common and perhaps mandatory in tumor cells (17).

Thus, the description “cell death” encapsulates several diverse processes. This can involve the physical loss of the cells through apoptosis/PCD or, in some circumstances, necrosis. Alternatively, cell death can refer to a genetic death, where cells no longer retain the ability to replicate, but continue to survive. Which method of cell death is induced depends on several factors, including the cell's internal environment, its external environment, and its developmental history. It is not yet clear whether the genes that regulate apoptosis/PCD and senescence overlap. For example, genes such as [p53](#) and **retinoblastoma** appear to be required for the regulation of apoptosis and senescence, but it is not clear whether these genes perform similar or different roles in each pathway. Indeed, cells that have evaded the first signal to senesce enter *crisis*, which is defined as the point at which the culture exhibits both apoptosis and proliferation. This suggests that there must be some overlap in the control of different types of cell death, and this may enable one type of cell death to be employed when another either is not suitable or is not able to be induced.

### Bibliography

1. C. A. Afshari and J. C. Barrett (1996) "Molecular genetics of in vivo cellular senescence". In *Cellular Aging and Cell Death* (N. J. Holbrook, G. R. Martin, and R. A. Lockshin, eds.), Wiley-Liss, New York, pp. 109–122.
2. A. J. Hale, C. A. Smith, L. C. Sutherland, V. E. Stoneman, V. Longthorne, A. C. Culhane, and G. T. Williams (1996) Apoptosis: Molecular regulation of cell death. *Eur. J. Biochem.* **237**, 884.
3. S. J. Martin, and D. R. Green (1995) Apoptosis and cancer: The failure of controls on cell death and cell survival. *Crit. Rev. Oncol. Hematol.* **18**, 137–153.
4. J. R. Tata (1966) Requirement for RNA and protein synthesis for induced regression of the tadpole tail in organ culture. *Dev. Biol.* **13**, 77–94.
5. R. A. Lockshin (1969) Programmed cell death. Activation of lysis by a mechanism involving the synthesis of protein. *J. Insect Physiol.* **15**, 1505–1516.
6. A. Glucksmann (1965) Cell death in normal development. *Arch. Biol. Liege* **76**, 419–437.
7. P. G. Clarke and S. Clarke (1996) Nineteenth century research on naturally occurring cell death and related phenomena. *Anat. Embryol. Berl.* **193**, 81–99.
8. R. A. Lockshin and C. M. Williams (1964) Programmed cell death. II. Endocrine potentiation of the breakdown of the intersegmental muscles of silkworms. *J. Insect Physiol.*, **10**, 643.
9. R. E. Ellis, J. Y. Yuan, and H. R. Horvitz (1991) Mechanisms and functions of cell death. *Ann. Rev. Cell Biol.* **7**, 663–698.
10. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1972) Cellular events in the adrenal cortex following ACTH deprivation. *J. Pathol.* **106**, ix.
11. M. D. Jacobson, M. Weil, and M. C. Raff (1997) Programmed cell death in animal development. *Cell* **88**, 347–354.
12. A. Fraser, N. McCarthy, and G. I. Evan (1996) Biochemistry of cell-death. *Curr. Opin. Neurobiol.* **6**, 71–80.
13. B. F. Trump, J. M. Valigorsky, J. H. Dees, W. J. Mergner, K. M. Kim, R. T. Jones, R. E. Pendergrass, J. Garbus, and R. A. Cowley (1973) Cellular change in human disease. A new method of pathological analysis. *Hum. Pathol.* **4**, 89–109.
14. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1980) Cell death: The significance of apoptosis. *Int. Rev. Cytol.* **68**, 251–306.
15. M. C. Raff (1992) Social controls on cell survival and cell death. *Nature* **356**, 397–400.
16. C. B. Thompson (1995) Apoptosis in the pathogenesis and treatment of disease. *Science* **267**, 1456–1462.
17. W. E. Wright and J. W. Shay (1996) "Mechanisms of escaping senescence in human diploid

cells". In *Cellular Aging and Cell Death* (N. J. Holbrook, G. R. Martin, and R. A. Lockshin, eds.), Wiley-Liss, New York, pp. 153–166.

### Suggestions for Further Reading

18. N. J. Holbrook, G. R. Martin, and R. A. Lockshin (1996) *Cellular Aging and Cell Death*. Vol. 16, Modern Cell Biology (J. B. Harford, series ed.), Wiley-Liss, New York. A good introductory text to both senescence and cell death.
19. M. Sluysers (1996) *Apoptosis in Normal Cancer and Development*, Taylor & Francis, London. This book covers all aspects of apoptosis, giving both a physical description of the process and covering all the molecular aspects of the pathway.

### Cell Fusion, Cell Hybrids

Fusions of both the external and intracellular [membranes](#) of cells are important for [differentiation](#) and [development](#). Moreover, enveloped **viruses** infect cells via fusion of their envelopes with cell or [endosome](#) membranes.

*Cell fusion* is the process of fusion of the membranes of two or more cells and results in the formation of cells with multiple **nuclei**. It occurs at various stages of the natural development of organisms, such as in the first step of fertilization of an **oocyte** with [sperm](#) and in myotube formation by fusion of myoblasts during differentiation of skeletal muscles. Artificial cell fusion can be induced by the addition of a high concentration of [Sendai virus](#), an enveloped virus of the paramyxovirus group, as was demonstrated in 1957–63 (1-3). This finding coincided with cell biology's beginning focus on the culture of **somatic cells** of birds and mammals. In 1961, Barski et al. (4) reported the appearance of hybrid cells after a few months of mixed culture of two different mouse cancer cell lines. These hybrid cells had a single nucleus containing [chromosomes](#) from both parent cell lines, and their appearance was considered to be due to the spontaneous fusion of cells of the two cell lines.

Sendai virus proved to be an important development in this field because it has some useful characteristics for the fusion of somatic cells:

1. Its targets are sialoglycoproteins and sialolipids, which are present in the cell membranes of almost all mammalian and fowl cells. Consequently, it can induce the fusion of a wide range of cells (3).
2. Its cell fusion activity is not affected by procedures that inactivate the viral [genome](#), such as exposure to ultraviolet (UV) light. Thus, fused cells could be prepared under conditions inhibiting virus growth, using a UV-inactivated virus.
3. There is no species specificity in its fusion of cells, so interspecific [heterokaryons](#) can be induced easily (5, 6), unlike in the case of fertilization.
4. The frequency of virus-induced hybrid formation is at least 1000 times greater than that of spontaneous hybridization.

Littlefield reported in 1964 (7) further progress in methods for the selection of hybrid cell clones, by fusing two different mutants defective in the **salvage pathway** for nucleotide biosynthesis and culturing them in a medium containing **aminopterin**, which inhibits *de novo* nucleotide synthesis. One of the mutant cell lines that were fused lacked **thymidine kinase**, the other **hypoxanthine-**

**guanine phosphoribosyl transferase**. The only cells that could grow from the mixed culture were fused cells that had acquired mutual [complementation](#) of the two mutant defects. Based on these new techniques, the field of *somatic cell genetics* was established in the 1960s.

That the chemical fusogen **polyethylene glycol** is effective for the fusion of [protoplasts](#) of **plant** cells was first reported in 1974 (8). With this fusogen, hybrid plants, such as the “pomato,” can be prepared in cultures of somatic hybrid cells. Moreover, *electroporation* (see [Transfection](#)) was found to be useful as a physical method for cell fusion in the 1980s (10, 11). Cell fusion by these two methods does not require viral receptors, so they have made possible the fusion of cells from all kinds of organisms.

In the early stage of somatic cell genetics in the 1960s and 70s, many reports provided important information on some basic phenomena in cell biology, permitting a focus on molecular biology in the next stage of the 1980s and 90s. The main findings were as follows:

1. Immediately after multinucleated cell formation by artificial cell fusion, some nuclear proteins derived from the parents are rapidly transferred and mixed between the various nuclei. This tends to induce synchronization of the stage of [DNA replication](#) of the nuclei after the fusion of randomly growing cells. The degree of this synchronization is greatest in cells with only two nuclei and decreases sharply with an increasing number of nuclei in polykaryocytes (12, 13). This may be why almost all randomly isolated hybrid cells have in their nucleus one set of chromosomes derived from each parent. As a special case, the fusion of cells in **mitosis** with cells in interphase causes the rapid dissolution of the nuclear membrane of the interphase cells, followed by the condensation of its chromosomes, a process named *chromosomal pulverization* (14) or *premature chromosome condensation* (15). Reactivation of dormant nuclei from chick erythrocytes was demonstrated upon their fusion with cultured mammalian cells (16, 17). In the case of cell [neoplastic transformation](#) by tumor viruses, induction of Simian Virus 40 ([SV40](#)) production was observed after the fusion of SV40-transformed nonproducer hamster cells with monkey cells, which are a permissive host for SV40 (18). Later, a similar observation was reported on the fusion of [Rous sarcoma virus](#) (RSV)-transformed rat cells with chicken cells (a permissive host of RSV) (19). In the case of cell differentiation, the distinction between luxury and household functions of cells was observed by the formation of hybrids of cells with different phenotypes (20). The findings were important for choosing the correct combination of cells for hybridization of differentiated cells. In 1974, Köhler and Milstein (21) reported that [monoclonal antibodies](#) could be prepared from hybrid cells ([hybridomas](#)) of **B cells** and **myelomas** (tumor cells obtained from B cells). This method of preparing monoclonal antibodies is now a major technique in molecular, cell, and developmental biology, and in medicine.
2. Weiss and Green (22) observed that the chromosomal balance is unstable in interspecific man/mouse hybrids and human chromosomes disappear randomly during serial passage in culture. Based on this finding, mapping of the human chromosomes became possible (23) and was followed by the technique of direct [in situ hybridization](#) of the chromosomes with [complementary DNA](#).
3. In 1972, Bootsma and colleagues (24) demonstrated that genetic complementation groups in a hereditary disease, *xeroderma pigmentosum*, could be classified by the cell fusion technique. This was the first successful genetic analysis of a human hereditary disease.

## 1. Mechanism of Cell Fusion

Artificial cell fusion using the various fusogens mentioned above has been established as a routine laboratory method. These fusogens, viral, chemical, and physical, have different modes of action on cell membranes, but induce similar changes in the cell membrane for fusion of the [lipid](#) bilayers. In general, **glycoproteins** are distributed on cell surfaces, with their **hydrophobic** domains embedded in the lipid bilayer of the cell membrane, their [hydrophilic](#) domains exposed to the outside, while their intracytoplasmic domains are associated with the [cytoskeletal](#) system. [Water](#) molecules are also

associated with the outside of cell membranes. These structures on the exteriors of cells inhibit the close contact of lipid bilayers of neighboring cell membranes, which is essential for cell fusion. Close contact of membranes requires on the cell membrane surface the transient appearance of areas from which glycoproteins are excluded. All the fusogens can induce such areas.

Sendai virus has two glycoproteins: HN with **receptor** binding and destroying activities, and F with fusogenic activity that is essential for fusion. The first step in cell fusion is the attachment of the virus to cell surfaces and agglutination of the cells, which is produced by HN activity. The second step is insertion of the fusogenic domain of F into the lipid bilayer. This domain is located at the N-terminus of F and consists of 15 relatively hydrophobic residues (25); it has the unique characteristic of trapping cholesterol molecules in its tertiary structure at 37°C (26). The amino acid sequence of this domain is well conserved in paramyxoviruses.

Simultaneous removal of cholesterol molecules from multiple sites in the cell membrane by the attachment of several hundred virus particles perturbs the cell membrane and causes breakdown of the normal barrier to ions. At this stage, calcium ions from the medium promptly penetrate the cell and induce changes in cell structures, such as separating the connection between the cytoplasmic domains of membrane proteins and the cytoskeletal system (see [Calcium Signaling](#)). Macromolecules can also **diffuse** through the cell membrane, and membrane fluidity increases, so that membrane proteins can move more freely in the lipid bilayer. This results in clustering of intramembrane domains of the proteins (demonstrable as cold-induced clustering) and the appearance of areas with no membrane proteins. The appearance of these areas is followed by close attachment of the lipid bilayers of neighboring cells, due to strong hydrophobic interactions, and cell fusion (27).

Polyethylene glycol is reported to induce clustering of membrane proteins similar to that induced by Sendai virus (28). Electroporation by electric pulses causes pores to be formed in cell membranes that may allow  $\text{Ca}^{2+}$  ions into the cytoplasm, as is also caused by Sendai virus. Thus, the fundamental mechanisms of these three fusogens appear to be similar.

It is known that calcium ions (29) and an energy (ATP) supply (2) are required for cell fusion. In the absence of either one, cells rapidly degenerate and fusion is greatly decreased. Evidence suggests that calcium ions associate directly with phospholipid molecules to normalize the perturbed membrane structure and promote a connection between the two lipid bilayers, in addition to acting in the cytoplasm as mentioned above. Why does cell fusion require an energy supply? In simple terms, cell fusion would be expected to produce a favorable decrease in **free energy**, by changing from a number of small vesicles to one large one, in which case an energy supply would not be necessary. Energy may be required for the removal of excess calcium ions introduced into the cytoplasm during cell fusion, but not for the membrane fusion itself. On culture of fused cells, the calcium ions are rapidly sequestered in **organelles** of the cells, and their level in the cytoplasm returns to normal. If the supply of energy is delayed, the cells degenerate.

## 2. Supplement

1. In the case of Sendai virus, cell-to-cell fusion is also possible as a result of viral envelope fusion itself, if one virus envelope fuses with the membranes of two different cells. But envelope fusion is slower than that induced by a high concentration of the virus and is observed to occur after completion of cell-to-cell fusion. Finally, as a result of viral envelope fusion, many viral glycoproteins are integrated into the fused cell surfaces, which are excluded from the cell surface by internalization by coated vesicles in culture (30) (see [Clathrin](#)).

Polykaryocytes (syncytia) are often observed in pathological tissues infected by an enveloped virus, such as paramyxoviruses, [retroviruses](#) and the [Herpes virus](#). This polykaryocyte formation may occur by cell fusion via viral envelope fusion. Newly synthesized viral glycoproteins are distributed

massively on the surfaces of infected cells, which consequently are quite similar to the viral envelope, and they may fuse with neighboring noninfected cells, as in viral envelope fusion.

2. Enveloped viruses cause infection by fusion of their envelopes either (a) with cell membranes at neutral pH or (b) with endosome membranes at acidic pH after their internalization from the cell surface. Infections by paramyxoviruses, retroviruses, and the Herpes virus are of the first type and induce syncytia formation *in vivo* and *in vitro*. [Influenza virus](#) is of the second type and does not induce syncytia *in vivo*, but the viral envelope can fuse with cell membranes under acidic conditions *in vitro*. This is because the influenza fusogenic glycoprotein (HA, or hemagglutinin) is not functional at neutral pH, but becomes functional under acidic conditions by a conformational change, forming trimers ([31](#)).

The fusogenic glycoproteins of various enveloped viruses differ, but all of them contain a hydrophobic domain that can interact with the lipid bilayer of cellular membranes. In some cases, the glycoproteins of the viruses are known to be synthesized in an inactive form in which their hydrophobic domain is hidden and then activated by **proteolytic** cleavage exposing this domain. This was first demonstrated by Homma ([32](#)) with the Sendai virus. In this virus, the F glycoprotein is synthesized as an inactive form  $F_0$  and is then cleaved to  $F_1$  and  $F_2$ , with the fusogenic domain being exposed at the N-terminus of  $F_1$ . The glycoprotein (gp160) of human immunodeficiency virus ([HIV](#)) a retrovirus, is cleaved to gp120 and gp41, and the fusogenic domain is exposed at the N-terminus of gp41 ([33](#)). The HA glycoprotein of influenza virus is also activated by the cleavage of inactive  $HA_0$  to  $HA_1$  and  $HA_2$ , with the hydrophobic domain being exposed at the N-terminus of  $HA_2$  ([34](#)) and becoming functional by trimer formation under acidic conditions.

### 3. Application of Fusogenic Reactions to Cell Engineering

Various biotechniques for the reconstitution of cells or introduction of macromolecules into cells have been developed using the unique characteristics of the interactions of fusogens with cell membranes. Introduction of macromolecules such as DNA, RNA, and proteins from the medium into the cytoplasm has become possible by perturbation of the cell membrane. Electroporation is usually used for that purpose ([35](#)), but treatment of the cells with Sendai virus can also be used ([36](#)).

Another technique for the introduction of macromolecules is based on the mechanism of infection of Sendai virus with cell membranes. In 1979, Uchida et al. ([37](#)) reported on the reassembly of viral envelopes as artificial **vesicles**, with HN and F glycoproteins embedded in their surface and any macromolecules present trapped inside. These pseudovirus particles containing the macromolecules instead of viral nucleocapsids will introduce those macromolecules into a cell that they “infect.” This technique is especially useful *in vivo*. Subsequently, modifications have been made to simplify the preparation procedure, using spontaneous fusion of the UV-inactivated virus with simple, artificial [liposomes](#) containing the desired macromolecules. Using such a preparation, the human [hepatitis B virus](#) genome has been introduced into rat liver cells *in vivo* to induce hepatitis ([38](#)). This technique is useful for drug delivery *in vivo* and gene therapy ([39](#)).

Reconstruction of cells was first reported by Veomett et al. in 1974 ([40](#)). On incubation with [cytochalasin](#), cells could be separated into *nucleoplasts* (enclosed by cell membranes but lacking cytoplasm) and *cytoplasts* (without a nucleus). Viable cells could be reconstituted by the fusion of nucleoplasts with cytoplasts. Heterologous combinations could also be constituted. This technique has been expanded to the preparation of cloned animals. Introduction of nuclei from somatic cells at the early stages after cleavage of fertilized eggs into enucleated eggs, by Sendai virus-mediated fusion, is reported to result in a high frequency of development of animals ([41](#)).

### Bibliography

1. Y. Okada, T. Suzuki, and Y. Hosaka (1957) *Med. J. Osaka Univ.* **7**, 709–717.



2. Y. Okada (1962) *Exp. Cell Res.* **26**, 98–128.
3. Y. Okada and J. Tadokoro (1963) *Exp. Cell Res.* **32**, 417–430.
4. G. Barski, S. Sorieul, and F. Cornfert (1961) *J. Natl. Cancer Inst.* **26**, 1269–1277.
5. H. Harris and J. F. Watkins (1965) *Nature* **205**, 640–646.
6. Y. Okada and F. Murayama (1965) *Exp. Cell Res.* **40**, 154–158.
7. J. Littlefield (1966) *Exp. Cell Res.* **41**, 190–196.
8. K. N. Kao and M. R. Michayluk (1974) *Planta* **115**, 355–367.
9. G. Melchers, M. D. Sacristan, and A. A. Holder (1978) *Carlsberg Res. Commun.* **43**, 203–218.
10. M. Senda, J. Takeda, S. Abe, and T. Nakamura (1979) *Plant Cell Physiol.* **20**, 1441–1443.
11. U. Zimmermann and J. Vienken (1982) *J. Membrane Biol.* **67**, 165–182.
12. T. Yamanaka and Y. Okada (1966) *Biken's J.* **9**, 159–175.
13. T. Yamanaka and Y. Okada (1968) *Exp. Cell Res.* **49**, 461–469.
14. H. Kato and A. A. Sanderberg (1968) *J. Nat. Cancer Inst.* **41**, 1117–1123.
15. P. N. Rao and R. T. Johnson (1972) *J. Cell Sci.* **10**, 495–513.
16. H. Harris (1965) *Nature* **206**, 583–588.
17. N. R. Ringertz (1974) In *Somatic Cell Hybridization* (R. L. Davidson and F. de la Cruz, eds.), Raven Press, New York, pp. 239–264.
18. P. Gerver (1966) *Virology* **28**, 501–509.
19. J. Svoboda, O. Machala, and I. Holzanek (1967) *Acta Virol.* **13**, 155–157.
20. R. L. Davidson and K. Yamamoto (1968) *Proc. Nat. Acad. Sci. USA* **60**, 894–901.
21. G. Köhler and C. Milstein (1975) *Nature* **256**, 495–497.
22. M. C. Weiss and H. Green (1967) *Proc. Nat. Acad. Sci. USA* **58**, 1104–1111.
23. R. J. Klebe, T. Chen, and F. H. Ruddle (1970) *Nat. Acad. Sci. USA* **66**, 1220–1227.
24. E. A. de Weerd-Kastelein, W. Keijzer, and D. Bootsma (1972) *Nature New Biol.* **238**, 80–83.
25. M. J. Gething, J. M. White, and M. D. Waterfield (1978) *Proc. Nat. Acad. Sci. USA* **75**, 2737–2740.
26. K. Asano and A. Asano (1985) *Biochem. Int.* **10**, 115–122.
27. J. Kim and Y. Okada (1981) *Exp. Cell Res.* **132**, 125–136.
28. D. S. Roos, J. M. Robinson, and R. L. Davidson (1983) *J. Cell Biol.* **97**, 909–917.
29. Y. Okada and F. Murayama (1966) *Exp. Cell Res.* **44**, 527–551.
30. J. Kim and Y. Okada (1982) *Exp. Cell Res.* **140**, 127–136.
31. I. A. Wilson, J. J. Skehel, and D. C. Wiley (1981) *Nature (Lond.)* **289**, 366–373.
32. M. Homma (1971) *J. Virol.* **8**, 619–629.
33. J. M. McCune et al. (1988) *Cell* **53**, 55–67.
34. S. G. Lazarowitz and P. W. Choppin (1975) *Virology* **68**, 440–454.
35. U. Zimmermann, J. Vienken, and G. Pilwat (1980) *Bioelectrochem. Bioenerg.* **7**, 553–574.
36. K. Tanaka, M. Sekiguchi, and Y. Okada (1975) *Proc. Nat. Acad. Sci. USA* **72**, 4071–4075.
37. T. Uchida et al. (1979) *Biochem. Biophys. Res. Commun.* **87**, 371–379.
38. K. Kato et al. (1991) *J. Biol. Chem.* **266**, 3361–3364.
39. V. J. Dzau, M. J. Mann, R. Morishita, and Y. Kaneda (1996) *Proc. Nat. Acad. Sci. USA* **93**, 11421–11425.
40. G. Veomett, D. M. Prescott, J. Shay, and K. R. Porter (1974) *Proc. Nat. Acad. Sci. USA* **71**, 1999–2002.
41. J. McGrath and D. Solter (1983) *Science* **220**, 1300–1302.

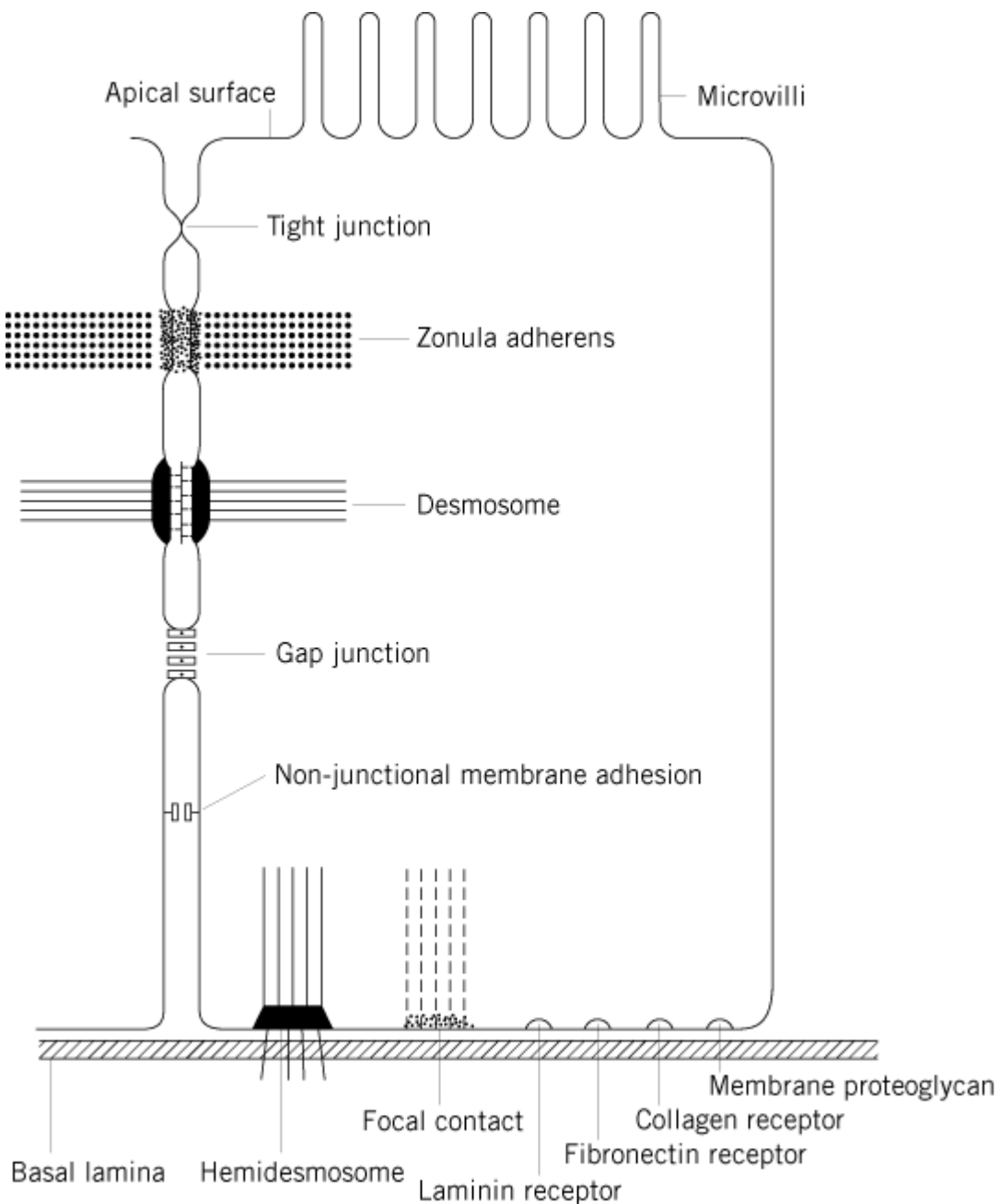
### Suggestions for Further Reading

42. N. Düzgunes, ed. (1993) *Membrane Fusion Technique*, Part B, Methods in Enzymology, Vol. **221**, Academic Press.
43. N. Düzgunes and F. Bronner, eds. (1988) *Membrane Fusion in Fertilization, Cellular Transport, and Viral Infection*, Current Topics in Membranes and Transport, Vol. **32**, Academic Press.
44. J. Wilschut and D. Hoekstra, eds. (1991) *Membrane Fusion*, Marcel Dekker, Inc.

### Cell Junctions

A cell junction is a cell surface structure observed by [electron microscopy](#) that mediates cellular interactions. Cell junctions are of two types, those that mediate cell–cell interactions and those that mediate cell–substratum interactions. In the former category, the cells of vertebrates possess four principal types, all of which are represented at the lateral surfaces of simple epithelial cells, such as those that line the intestine and kidney tubules (Fig. [1](#)). The four principal types of cell–cell junction are (i) the [tight junction](#) or zonula occludens (pl. zonulae occludentes); (ii) the [intermediate junction](#) or zonula adherens (pl. zonulae adherentes); (iii) the **desmosome** or macula adherens (pl. maculae adherentes); (iv) the [gap junction](#) or nexus. The principal type of cell–substratum junction is the [hemidesmosome](#), which is present in contact with the basement membrane in the basal cells of stratified epithelia, such as epidermis, in some simple epithelial cells, especially amnion, and mammary epithelium. Cells that are well spread in culture exhibit a second type of cell–substratum junction, the focal adhesion or [focal contact](#) (Fig. [1](#)). This is a region where the lower cell surface is most closely apposed to the culture substratum (separation 10 to 15 nm). It is not clear whether the focal contact has a strict equivalent *in vivo*, but its study has made an extensive and important contribution to the understanding of cell adhesion. The basic structures of the principal types of cell junctions were elucidated by electron microscopy in the 1960s ([1-3](#)). Since then, a substantial amount of detailed knowledge of the molecular composition of each junctional type has been revealed.

**Figure 1.** Generalized simple epithelial cell showing cell junctions and other adhesion mechanisms. The lateral surface mediates cell–cell interactions via adhesive junctions, the zonula adherens (intermediate junction) and the desmosome, and communicating gap junctions that allow small molecules to be exchanged between the cytoplasm of adjacent cells. The tight junction, or zonula occludens, restricts and regulates paracellular permeability. The basal surface mediates adhesion to the basement membrane. Hemidesmosomes are present in some simple epithelia but are most numerous in epidermis. Integrin and non-integrin receptors for several matrix components are present here. When cells are cultured, many of these become localized in focal contacts, small areas where adhesion to the substratum is strongest. (From *Journal of Cell Science*, with permission of the Company of Biologists.)



The intermediate junction, the desmosome, the hemidesmosome, and the focal contact are primarily adhesive junctions that bind cells to each other or to their substratum. Each represents a cell-surface site at which cell adhesion molecules (see [Cell Surface Adhesion Receptors](#)) are concentrated. In desmosomes and intermediate junctions, the adhesion molecules are specific types of [cadherins](#) (4, 5); in hemidesmosomes and focal contacts, they are [integrins](#) (6). These junctions also represent cell-surface attachment points for elements of the [cytoskeleton](#). In desmosomes and hemidesmosomes, the associated cytoskeletal elements are [intermediate filaments](#), also known as tonofilaments. These are important in providing structural continuity throughout a tissue. This continuity is disrupted in certain **autoimmune** and genetic diseases of desmosomes (eg, pemphigus) (7), hemidesmosomes (eg, bullous pemphigoid, junctional epidermolysis bullosa) (8), or intermediate filaments (eg, epidermolysis bullosa simplex) (9). Such diseases result in structural disruption of the epidermis and/or mucosal tissues, with effects ranging from mild to lethal. In intermediate junctions and focal contacts, the cytoskeletal elements are [actin](#) filaments, also known as [microfilaments](#). The actin filaments associated with focal contacts are organized into large bundles, called stress fibers, of

several micrometers in length. During animal development, morphogenetic events such as gastrulation and neural tube rolling involve gross changes in cell shape that appear to be caused by contraction of submembranous microfilament rings associated with zonula adherens-type junctions (10, 11). More limited contractile activity of such rings may be involved in regulation of the permeability of tight junctions, with which zonulae adherentes are frequently closely associated in the simple epithelial junctional complex, a region where junctions are concentrated at the extreme apicolateral interface between the cells.

Transduction of signals across the plasma membrane is a vital function in the regulation of normal cell behavior in both developing and adult organisms. Many adhesion receptors have now been shown to transduce signals, either from inside the cell to outside regulating the function of the receptors themselves, or from outside to inside regulating signalling pathways and gene expression in response to adhesive stimuli. Such [signal transduction](#) is an additional important function of adhesive junctions (12).

The principal function of the tight junctions is to occlude the paracellular channels of epithelia, thereby restricting intercellular leakage of molecules between the distinct biological compartments separated by the epithelium (13). Tight junctions are also important in forming tight seals between some endothelial cells, such as those of the blood–brain barrier. Their sealing properties arise because they are zonular in nature, encircling the entire apicolateral margins of simple epithelial cells, and because there is direct contact between the plasma membranes of adjacent cells. Tight junctions are also important in maintaining the composition of the different membrane domains of epithelial cells, because they prevent [diffusion](#) of molecules within the outer leaflet of the plasma membrane, thereby restricting them to either the apical or the basolateral domain (14). This property is sometimes referred to as the “fence” function of tight junctions, and their occluding properties are described as a “gate” function.

Gap junctions, so-called because they exhibit a regular plasma membrane separation of 2 nm, are punctate membrane sites that provide [hydrophilic](#) channels for direct cell-to-cell communication (15). They permit diffusion of molecules of less than 1000 daltons between cells. These ubiquitous junctions are important in both excitable and nonexcitable tissues. In the former they facilitate direct transmission of electrical impulses, while in the latter they mediate a phenomenon called metabolic cooperation, and they may also transmit signaling molecules such as **calcium** and **inositol phosphates** between cells. Their function is essential in embryonic development (16) and in the normal functioning of adult tissues—for example, in coordination of the contraction of cardiac muscle and smooth muscle in the uterus during parturition (17, 18).

Gap junctions appear to be widespread in the animal kingdom, being present even in primitive organisms such as *Hydra*. Intermediate junctions have also been clearly demonstrated in invertebrates (eg, in *Drosophila*). The occurrence of desmosomes and tight junctions in invertebrates is much less clear. Desmosome-like structures are clearly present in insects—for example, between the apposed epithelia of the upper and lower surface of the wing blade in *Drosophila*. However, these desmosome-like structures appear to be associated with [microtubules](#) rather than intermediate filaments, and their adhesion may be mediated by integrins rather than cadherins. Tight junctions are absent from insects, which instead possess septate junctions, structures that have a ladder-like ultrastructural appearance between adjacent plasma membranes and which, like tight junctions, appear to restrict paracellular permeability (19).

## Bibliography

1. M. G. Farquhar and G. E. Palade (1963) *J. Cell Biol.* **17**, 375–412.
2. D. E. Kelly (1966) *J. Cell Biol.* **28**, 51–72.
3. J-P. Revel and M. Karnovsky (1967) *J. Cell Biol.* **33**, 7–12.
4. J. L. Holton et al. (1990) *J. Cell Sci.* **97**, 239–246.

5. P. J. Koch et al. (1990) *Eur. J. Cell Biol.* **53**, 1–12.
6. A. Sonnenberg et al. (1991) *J. Cell Biol.* **113**, 907–917.
7. J. R. Stanley (1995) In “Cell adhesion and human disease.” *Ciba Found. Symp.* **189**.
8. R. E. Burgeson and A. M. Christiano (1997) *Curr. Opin. Cell Biol.* **9**, 651–658.
9. W. H. I. McLean and E. B. Lane (1995) *Curr. Opin. Cell Biol.* **7**, 118–125.
10. T. E. Schroeder (1970) *J. Embryol. Exp. Morphol.* **23**, 427–462.
11. P. Karfunkel (1971) *Dev. Biol.* **25**, 30–56.
12. C. Rosales et al. (1995) *Biochem. Biophys. Acta.* **1242**, 77–98.
13. J. L. Madara and K. Darmsathaphorn (1985) *J. Cell Biol.* **101**, 2124–2133.
14. K. Simons (1990) In *Morphoregulatory Molecules* (G. M. Edelman, B. A. Cunningham, and J.-P. Thiery, eds.), Wiley, New York, pp. 341–356.
15. N. B. Gilula, O. R. Reeves, and A. Steinbach (1972) *Nature* **235**, 262–265.
16. A. E. Warner, S. C. Guthrie, and N. B. Gilula (1984) *Nature* **311**, 127–131.
17. W. C. Cole and R. E. Garfield (1985) In *Gap Junctions* (M. V. L. Bennett and D. C. Spray, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 215–230.
18. E. C. Beyer, D. L. Paul, and D. A. Goodenough (1987) *J. Cell Biol.* **195**, 2621–2629.
19. P. J. Bryant (1994) In *Molecular Mechanisms of Epithelial Cell Junctions: From Development to Disease* (S. Citi, ed.) R. G. Landes, Austin, TX, pp. 1–21.

### Suggestions for Further Reading

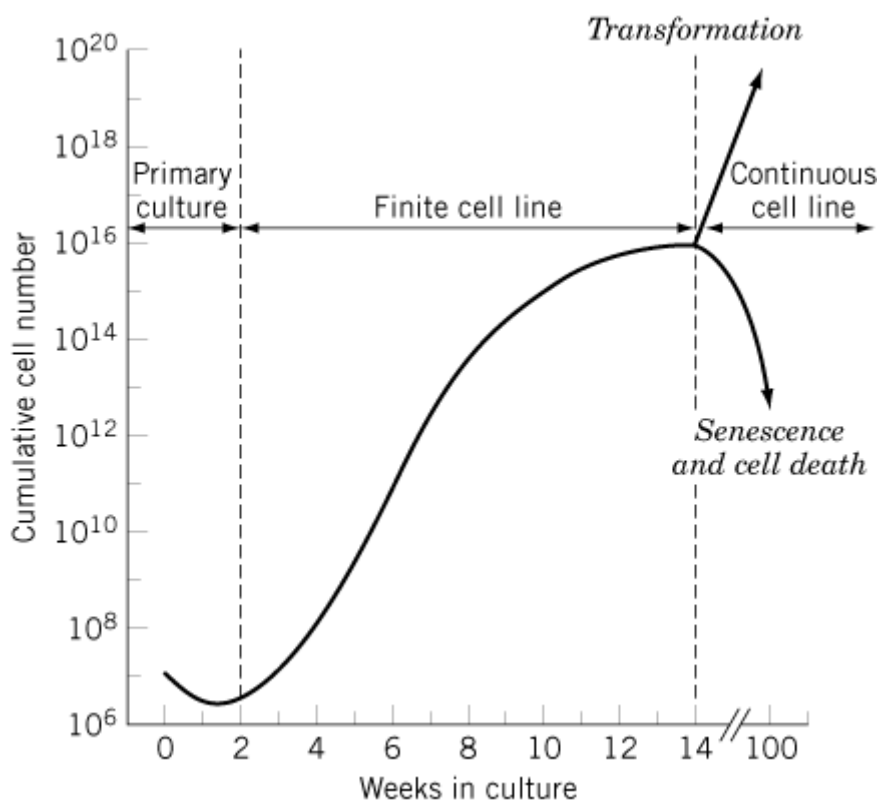
20. D. R. Garrod and J. E. Collins (1992) "Intercellular Junctions and Cell Adhesion in Epithelial Cells". In “*Epithelial Organization and Development*” (T. P. Fleming, ed.), Chapman and Hall, London, pp. 1–52. (A source of reference to the early literature.)
21. D. R. Garrod, M. A. J. Chidgey, and A. J. North (eds.) (1998) *Adhesive Interactions of Cells*, JAI Press, Greenwich, CT. (Reviews on desmosomes by the editors, tight junctions by S. Citi, and adherens junctions and focal contacts by A. Ben-Ze"ev, as well as other reviews in cell adhesion.)
22. D. A. Goodenough, J. A. Goliger, and D. L. Paul (1996) *Annu. Rev. Biochem.* **65**, 475–502. (A thorough review of gap junctions.)
23. L. Barradori and A. Sonnenberg (1996). *Curr. Opin. Cell Biol.* **8**, 647–656. (An introduction to the literature on hemidesmosomes.)

### Cell Line

When a primary cell culture (see [Tissue Culture](#)) is subcultured, it becomes a cell line. This implies that multiple lineages of cells, not necessarily distinct from each other, coexist in the culture ([1](#)). Serial subculture, while maintaining these parallel lineages, will tend to show convergence towards whatever common **phenotype** is best adapted to the culture conditions being employed, as the cell type with the greatest proliferative ability will predominate. The cell line may be finite and die out after a fixed number of population doublings, or become a continuous cell line (see [Immortalization](#)) (Fig. [1](#)). Cell lines, particularly continuous cell lines, can become a valuable resource, particularly if preserved in liquid nitrogen after appropriate characterization and validation. Many continuous cell lines are currently available, with wide-ranging properties, including [drug resistance](#) markers, inducibility for specific [enzymes](#), estrogen sensitivity, cornification, blood vessel formation,

hemoglobin synthesis, myogenesis, and transfection susceptibility (Table 1). Finite cell lines with distinct phenotypic properties are also becoming available through the development of serum-free selective media (2-4). (see Serum Dependence) and can be derived in the laboratory or purchased from commercial suppliers.

**Figure 1.** Effect of continued culture passage on cumulative cell number, assuming that no cells are discarded. The curve shows the initial decline due to selection, then exponential growth during the replicative phase, and then growth arrest and eventual deterioration following senescence, in a finite cell line, or continued proliferation, often at an enhanced rate, following transformation. Reproduced from Freshney's *Culture of Animal Cells, A Multi-Media Guide*, 1999, Wiley-Liss, N.Y.



**Table 1. Cell Lines in Common Use<sup>a</sup>**

| Cell line | Morphology | Origin     | Age       | Status | Ploidy  | Characteristics                      | Referen   |
|-----------|------------|------------|-----------|--------|---------|--------------------------------------|---|
| MRC5      | Fibroblast | Human lung | Embryonic | Normal | Diploid | Susceptible to human viral infection | Jacobs (1970) <i>Nature</i> <b>22</b> : 168     |
| WI38      | Fibroblast | Human lung | Embryonic | Normal | Diploid | Susceptible to human viral infection | Hayflick Moorhea (1961) <i>E: Cell Res.</i> 585 |

|             |             |                       |           |            |           |   |  |
|-------------|-------------|-----------------------|-----------|------------|-----------|---|--|
| IMR90       | Fibroblast  | Human lung            | Embryonic | Normal     | Diploid   | Susceptible to human viral infection                  | Nichols et al (1977) <i>Science</i> <b>1</b> 60                      |
| A2780       | Epithelial  | Human ovary           | Adult     | Neoplastic | Aneuploid | Chemosensitive with resistant variants                | Tsuruo et al (1986) <i>J. Cancer Res</i> <b>77</b> , 941             |
| A549        | Epithelial  | Human lung            | Adult     | Neoplastic | Aneuploid | Synthesizes pulmonary surfactant                      | Giard et al (1972) <i>J. Natl. Cancer Inst.</i> <b>51</b> , 1        |
| A9          | Fibroblast  | Mouse subcutaneous    | Adult     | Neoplastic | Aneuploid | HGPRT-ve: deriv. L929                                 | Littlefield (1964) <i>Nature</i> <b>201</b> 1142                     |
| BHK21-C13   | Fibroblast  | Syrian hamster kidney | NB        | Normal     | Aneuploid | Transformable by polyoma virus                        | Macpherson and Stokes (1962) <i>Virology</i> 147                     |
| BRL3A       | Epithelial  | Rat liver             | NB        | Normal     |           | Produce IGF-II  | Coon (1971) <i>J. Cell Biol.</i> <b>39</b> , 29a                     |
| Caco-2      | Epithelial  | Human colon           | Adult     | Neoplastic | Aneuploid | Transports ions and amino acids                       | Fogh (1971) <i>J. Natl. Cancer Inst.</i> <b>58</b> , 209             |
| CHANG liver | Epithelial  | Human liver           | Embryonic | Normal?    | Aneuploid | HeLa contaminated                                     | Chang (1965) <i>Proc. Soc. Exp. Biol. Med.</i> <b>87</b> ,           |
| CHO-K1      | Fibroblast  | Chinese hamster ovary | Adult     | Normal     | Diploid   | Simple karyotype                                      | Puck et al (1958) <i>J. Exp. Med.</i> <b>108</b> , 945               |
| EB-3        | Lymphocytic | Human                 | Juvenile  | Neoplastic | Diploid   | EB virus + ve   | Epstein and Barr (1964) <i>Lancet</i> <b>1</b> ,                     |
| GH1, GH3    | Epithelial  | Rat                   | Adult     | Neoplastic | Aneuploid | Produce growth hormone                                | Yasumura et al. (1966) <i>Science</i> <b>1</b> 1186                  |
| HeLa        | Epithelial  | Human cervix          | Adult     | Neoplastic | Aneuploid | G6PD type A   | Gey et al (1952) <i>Cancer Res</i> <b>12</b> , 364                   |
| HeLa-S3     | Epithelial  | Human cervix          | Adult     | Neoplastic | Aneuploid | High plating efficiency; will grow well in suspension | Puck and Marcus (1955) <i>Proc. Natl. Acad. Sci.</i> <b>41</b> , 852 |

|        |             |                               |           |  |  |  |
|--------|-------------|-------------------------------|-----------|--|--|--|
|        |             |                               |           |  |  | <i>Sci. USA</i><br>432   |
| Hep-2  | Epithelial  | Human larynx                  | Adult     | Neoplastic Aneuploid HeLa contaminated                   |  | Moore et al (1955) <i>Cancer Res</i> <b>15</b> , 598   |
| HT-29  | Epithelial  | Human colon                   | Adult     | Neoplastic Aneuploid Differentiation inducible with NaBt |  | Fogh and Trempe (1975) in <i>Human Tumor Cells in vitro</i> , J. Fogh, ed. Academic Press, New York, p.              |
| KB     | Epithelial  | Human oral                    | Adult     | Neoplastic Aneuploid HeLa contaminated                   |  | Eagle (1955) <i>Proc. Soc. Exp. Biol. (N.Y.)</i> <b>89</b> , 362   |
| L1210  | Lymphocytic | Mouse                         | Adult     | Neoplastic Aneuploid Rapidly growing; suspension         |  | Law et al (1949) <i>J. Natl. Cancer Inst.</i> <b>10</b> , 1  |
| L5178Y | Lymphocytic | Mouse                         | Adult     | Neoplastic Aneuploid Rapidly growing suspension          |  |  |
| L929   | Fibroblast  | Mouse                         | Adult     | Normal Aneuploid Clone of L cell                         |  | Sanford et al (1948) <i>J. Natl. Cancer Inst.</i> <b>9</b> , 21  |
| LS     | Fibroblast  | Mouse                         | Adult     | Neoplastic Aneuploid Grow in suspension: deriv. L929     |  | Paul and Struthers (unpublished)   |
| MCF7   | Epithelial  | Human breast pleural effusion | Adult     | Neoplastic Aneuploid Estrogen receptor +ve               |  | Soule et al (1973) <i>J. Natl. Cancer Inst.</i> <b>51</b> , 1  |
| P388D  | Lymphocytic | Mouse                         | Adult     | Neoplastic Aneuploid Grow in suspension                  |  | Dawe and Potter (1973) <i>Am. J. Pathol.</i> <b>33</b> , 603; Koren et al (1975) <i>J. Immunol.</i> <b>114</b> , 894 |
| S180   | Fibroblast  | Mouse                         | Adult     | Neoplastic Aneuploid Cancer chemotherapy screening       |  | Dunham and Stewart (1953) <i>J. Natl. Cancer Inst.</i> <b>13</b> , 1   |
| STO    | Fibroblast  | Mouse                         | Embryonic | Normal Aneuploid Used as feeder                          |  | Bernstein  |



|         |            |                             |           |            |           |   |  |
|---------|------------|-----------------------------|-----------|------------|-----------|---|--|
|         |            |                             |           |            |           | layer for embryonal stem cells                                | (1975) <i>Proc. Natl. Acad. Sci. USA</i> 1441      |
| 3T3-L1  | Fibroblast | Mouse, Swiss                | Embryonic | Normal     | Aneuploid | Adultipose differentiation                                    | Green and Kehinde (1974) <i>Cancer Res.</i> 1, 113 |
| 3T3 A31 | Fibroblast | Mouse BALB/c                | Embryonic | Normal     | Aneuploid | Contact inhibited; readily transformed                        | Aaronsor and Toda (1968) <i>J. Physiol.</i> 141    |
| NRK49F  | Fibroblast | Rat kidney                  | Adult     | Normal     | Aneuploid | Induction of suspension growth by transforming growth factors | DeLarco and Todaro (1978) <i>J. Physiol.</i> 335   |
| Vero    | Fibroblast | Monkey kidney               | Adult     | Normal     | Aneuploid | Viral substrate and assay                                     | Hopps et al (1963) <i>J. Immunol.</i> 416          |
| ZR-75-1 | Epithelial | Human breast, ascites fluid | Adult     | Neoplastic | Aneuploid | ER-ve, EGFr +ve   | Engel (1975) <i>Cancer Res.</i> 38, 3352           |

<sup>a</sup> Modified from R. I. Freshney *Culture of Animal Cells*, 1994, 3rd ed., Wiley-Liss, N.Y., p. 151.

## 1. Origin of Cell Lines

Regenerating tissues *in vivo* are made up of a small, self-repopulating, [stem cell](#) pool, an expandable pool of proliferating **progenitor cells**, and a nonproliferating differentiated cell pool. On demand, cells leave the stem cell pool and enter the progenitor compartment, where expansion is regulated to meet the current demand in the differentiated cell compartment. When cells enter the differentiated cell compartment, this may be an irreversible process, as seen with erythrocytes, keratinocytes, and neurons; or it may be reversible, as seen with hepatocytes, endothelial cells, and fibrocytes. Cell lines can be derived from any tissue with a proliferating compartment, or from cells that can re-enter a proliferating compartment. It is possible that cell lines will contain stem cells, but, except for hemopoietic cells (see [Hematopoiesis](#)), markers are unavailable to determine this. As many cell lines express lineage markers and can, under appropriate conditions, differentiate, it seems most likely that they are derived from the progenitor cells of the tissue.

As cell lines are derived from random outgrowth or from disaggregated cells, multiple lineages of cells will be present. The purity of the culture will be determined by how well these lineages are matched, that is, derived from cells with the same phenotypic fate, or are derived from cells with dissimilar fates. In practice, the selective pressure exerted by the medium will tend to limit the cell population to cells of like phenotype with similar survival capacity, and cells that proliferate less rapidly will gradually be overgrown. It is possible that cells within one phenotypic group have different proliferative capacities but interactions between them produces a uniform proliferative rate in the entire population.

## 2. Finite Cell Lines

When a cell line is generated at the first subculture of a primary culture, its lifespan is determined, initially, by the environmental conditions. If the conditions are inadequate, the cell line will die out within one or two subcultures, but if the medium, substrate, and other conditions are satisfactory, then the cell line may progress through several serial subcultures. The ultimate limit to the number of subcultures will be determined by the potential regenerative life-span of the cells (see [Senescence](#); [Immortalization](#)). Most cell lines from normal tissues will undergo a fixed number of population doublings (Fig. 1), given by the product of the number of passages and the estimated number of doublings per passage (5). These are known as finite cell lines. Normal human skin fibroblasts will usually achieve around 50–60 population doublings and then enter [senescence](#). Senescent cells no longer proliferate, may be partially differentiated, and can remain viable for several months. Cultures of finite cell lines must, therefore, be stored frozen at an early passage level, thawed when required, and used between predetermined passage levels, usually between 10 and 30, to ensure an adequate and consistent supply.

## 3. Continuous Cell Lines

Some cell lines, such as those derived from mice or from tumors of many species, are not limited by a finite lifespan. They progress either by smooth transition to an immortal cell line (see [Immortalization](#)) or undergo [neoplastic transformation](#) at some stage and, instead of dying out after a set number of population doublings, continue to proliferate unchecked. Some cultures show evidence of selection and go through a period called *crisis*, when most cells in the population die out by senescence, but a transformed subpopulation survives, usually with an enhanced growth rate, increased cloning efficiency, loss of [contact inhibition](#), acquisition of anchorage-independent growth, ability to grow to a higher saturation density, and increased tumorigenicity in animals (see: [Neoplastic Transformation](#)).

There are a large number of continuous cell lines in existence, many of which are banked in repositories such as the American Type Culture Collection (ATCC) or the European Collection of Animal Cell Cultures (ECACC). They form a valuable resource of vigorously growing cultures, capable of unlimited expansion, but are subject to several caveats.

1. Their origin must be validated before use. Acquisition from a reputable cell bank, or the originator, will usually guarantee this, but there are many instances of continuous cell lines being cross-contaminated by other, more vigorously growing cell lines such as [HeLa Cells](#) (6-8). The identity of a cell line needs to be confirmed prior to extensive use. This is done most effectively by performing a [DNA fingerprinting](#), but it is also possible by immunotyping, **isoenzyme** analysis, or [chromosome](#) analysis.

2. The cells must be shown to be free from infection. Continuous cell lines form an ideal substrate for [mycoplasma](#), which can grow undetected in the culture (9). Again, most reputable cell banks will be able to demonstrate that cell lines being distributed are mycoplasma-free, but stocks still need regular checks every 1 or 2 months to ensure that they remain uninfected.

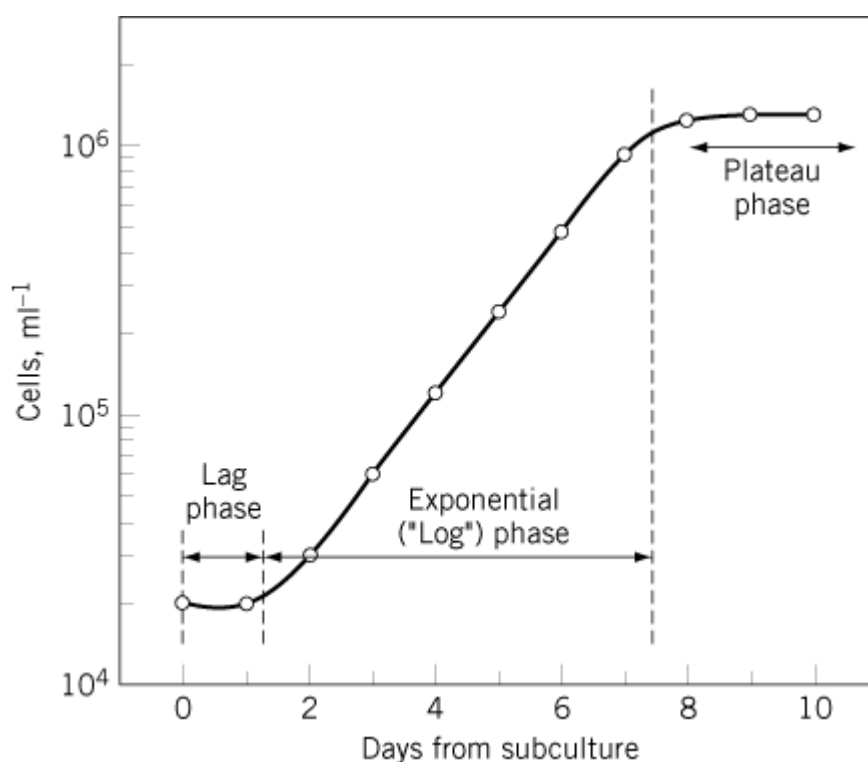
3. Continuous cell lines are genetically unstable. They are usually [aneuploid](#) (the chromosomal complement differs in number, and by chromosomal rearrangements, from the donor) and heteroploid (contain subpopulations with differing chromosomal constitutions). This is reflected in their phenotypic diversity, including variations in morphology, enzyme activity, antigenic expression, and growth characteristics. It is normal practice to clone such populations, select a clone with the required properties, and cryopreserve sufficient ampoules (12–100) for future use. Cultures are generally maintained for a limited period, usually about 3 months, and then replaced from frozen stock.

## 4. Subculture

Subculture, or passage, is the transfer of a culture from one vessel to another, usually using **trypsin** to dislodge the cells from the substrate and dissociate them from each other. If the cell line is proliferating, then subculture will also imply diluting the cell concentration achieved at the end of the culture period, to a new seeding concentration to initiate a new culture. Growth from seeding until the next subculture is known as a growth cycle, which is repeated each time the culture is passaged. A series of growth cycles constitute serial propagation, and this should not normally proceed beyond about 3 months without replacing stock from the freezer (see above).

Each growth cycle is composed of a lag phase (Fig. 2), where the cells are recovering from trypsinization, synthesizing new **extracellular matrix**, adhering to the substrate, spreading, repolymerizing the **actin** microfilament **cytoskeleton**, and responding to signaling from extracellular matrix adhesions and cytoskeletal rearrangements, which leads to expression of **cyclins** and, eventually, reentry to the cell cycle. When the cells start to divide, they enter the exponential or “log” phase of the growth cycle and will continue exponential growth until all the available growth surface is occupied or the medium is exhausted. If the medium is limiting, it is customary to replace it half way through the exponential phase of growth. When all the available growth surface is occupied, the culture is said to be *confluent* and will require to be subcultured again. If the culture is allowed to progress beyond confluence, growth will slow down, as a result of density limitation of growth (see **Contact Inhibition**), and the culture enters the plateau phase of the growth cycle. If the cells are normal and respond to normal growth control signals, there will be little evidence of cell proliferation or cell loss in plateau, specifically, minimal turnover. If, however, the cells are transformed, they will continue to proliferate for several cell generations after reaching confluence and reach a higher saturation density at plateau, due to deficiencies in contact inhibition of cell motility and density limitation of cell proliferation (10, 11). Even in plateau, there will be continued cell proliferation, although at a lower level than in exponential growth, balanced by increased cell loss from **apoptosis**, resulting in a greater turnover than occurs in normal cells in plateau.

**Figure 2.** Phases in the growth cycle of cultured cells following subculture. Modified from Freshney's *Culture of Animal Cells, A Multi-Media Guide*, 1999, Wiley-Liss, N.Y.



The numerical parameters that can be derived from the growth curve are the duration of the lag period, the doubling time in mid–log phase, and the saturation density in plateau. The last will depend on the feeding regimen and, if cell kinetics are being determined, should be measured under nonlimiting medium conditions. Where these parameters are consistent, they can be used to calculate the split ratio, the degree of dilution required to reinitiate a new growth cycle with a short lag period and to achieve the appropriate density for subculture at a convenient time in the future, usually around one week. The split ratio should be a power of 2 when handling finite cell lines. This allows an approximate calculation of the number of generations that have elapsed since the last subculture; for instance, a split ratio of 8 would imply that the cells had undergone 3 population doublings in each growth cycle. Split ratios are less useful with continuous cell lines. They are not needed to determine the stage in the life cycle of the cell line, as it is immortal, and the population doubling time is so much shorter that split ratios of the order of 1–100 are required to produce a week-long growth cycle. At this level of dilution, it is more advisable to determine the cell concentration at each subculture and to dilute to give the seeding concentration that will allow growth to the next subculture in a suitable interval, say, one week.

## 5. Selective Culture

Alterations in the choice of medium and/or substrate will determine which cells survive in primary and early passage cultures. Most selective media are serum-free (see [Serum Dependence](#)), and recipes are available for the culture of many different cell types, including fibroblasts, epithelial cells (epidermal, mammary), glial cells, melanocytes, endothelial cells, and smooth muscle Table 2. Selection of a purified cell culture, with defined characteristics, by cloning, use of a selective medium, or physical separation, allows the resultant culture to be called a *cell strain*.

**Table 2. Examples of Selective Media<sup>a</sup>**

| Cells or Cell Line          | Medium   | Reference  |
|-----------------------------|----------|--|
| Fibroblasts                 | MCDB 202 | McKeehan et al. (1977) <i>In Vitro</i> <b>13</b> , 399.                      |
| Fibroblasts                 | MCDB 110 | Bettger et al. (1981) <i>Proc. Natl. Acad. Sci. USA</i> <b>78</b> , 5588.    |
| Keratinocytes               | MCDB 153 | Tsao et al. (1982) <i>J. Cell Physiol.</i> <b>110</b> , 219                  |
| Bronchial epithelium        | LHC      | Lechner and LaVeck (1985) <i>J. Tissue Cult. Meth.</i> <b>9</b> , 43         |
| Mammary epithelium          | MCDB 170 | Hammond et al. (1984) <i>Proc. Natl. Acad. Sci. USA</i> <b>81</b> , 5435     |
| Prostate epithelium (rat)   | WAJC 401 | McKeehan et al. (1982) <i>In Vitro</i> <b>18</b> , 87.                       |
| Prostate epithelium (human) | WAJC 404 | McKeehan et al. (1984) <i>Cancer Res.</i> <b>44</b> , 1998                   |
| Glial cells                 |          | Michler-Stuke and Bottenstein (1982) <i>J. Neurosci. Res.</i> <b>7</b> , 215 |
| Melanocytes                 |          | Naeyaert et al. (1991) <i>Br. J. Dermatol.</i> <b>125</b> , 297              |

|                        |          |   |
|------------------------|----------|---|
| Small cell lung cancer | HITES    | Carney et al. (1981) <i>Proc. Natl. Acad. Sci. USA</i> <b>78</b> , 3185 |
| Adenocarcinoma of lung |          | Brower et al. (1986) <i>Cancer Res.</i> <b>46</b> , 798                 |
| Colon carcinoma        |          | Van der Bosch et al. (1981) <i>Cancer Res.</i> <b>41</b> , 611          |
| Endothelium            | MCDB 130 | Knedler and Ham (1987) <i>In Vitro</i> <b>23</b> , 481                  |

<sup>a</sup> Modified from R. I. Freshney *Culture of Animal Cells*, 1994, 3rd ed., Wiley-Liss, New York.

## 6. Characterization

Cell morphology and growth characteristics are simple criteria for characterization. Spindle-shaped cells are usually considered to be fibroblast-like, without necessarily implying that they are known to be of the fibrocyte lineage. Similarly, cells that grow in patches of polygonal, or pavement-like, cells with distinct boundaries between the cells are said to be epithelial-like, although numerous other cell types can assume a similar morphology. Characterization of a cell strain employs such criteria as (1) identification of the type of [intermediate filament](#) protein (eg, [cytokeratin](#) subtype in different epithelial cells, desmin in muscle cells, glial fibrillary acidic protein in astrocytes, and [neurofilament](#) protein in neurons and some neurendocrine cells) ([12](#)); (2) expression of cell-surface [antigens](#), such as epithelial membrane antigen ([13](#)); (3) enzymes, such as tyrosine aminotransferase and its inducibility by glucocorticoids in hepatocytes ([14](#)); (4) marker chromosomes ([15](#)); (5) **isoenzymes** ([16](#)); and (6) the **DNA fingerprint** ([17](#), [18](#)). The same criteria may be used to characterize a cell line, but would give only an average for the population in a cell strain.

## 7. Cell Banks

Most of the cell lines in common use are lodged in cell banks containing records of their identifying characteristics. Cell lines or strains that have been characterized and shown to be free of mycoplasma can be submitted to a cell bank for safe-keeping, with or without authority to distribute to other users. Cell lines may be obtained from cell banks for a set charge (Table 3) ([19](#)).

**Table 3. Cell Banks<sup>a</sup>**

|                          |   |
|--------------------------|---|
| United States and Canada | American Type Culture Collection (ATCC), 12301 Parklawn Drive, Rockville, MD 20852  |
|                          | National Institute of General Medical Sciences (NIGMS), Human Genetic Mutant Cell and National Institute on Ageing Cell Culture Repositories, Coriell Institute for Medical Research, Copewood Street, Camden, NJ 08103 |
|                          | Repository for Human Cell Strains, and Cell Repository for Neuromuscular Diseases, Children's Hospital Research Institute, McGill University, 2300 rue Tupper Street, Montreal, Quebec H3H 1P3, Canada                  |

|           |  |
|-----------|--|
| Europe    | European Collection for Animal Cell Cultures (ECACC),<br>PHLS/CAMR, Porton, Salisbury, England   |
|           | European Collection for Biomedical Research, Dept. Cell Biology<br>and Genetics, Erasmus University, P.O. Box 1783,<br>Rotterdam, Netherlands                            |
|           | Collection Nationale des Cultures des Microorganismes, Institute<br>Pasteur, 25 rue du Dr. Roux, F-75724 Paris Cédex 15, France  |
|           | Human Genetic Cell Repository, Hospices Civils de Lyon,<br>Hôpital Debrousse, 29 rue Soeur Bourrier, F-69322 Lyon<br>Cédex 05, France                                    |
|           | Tumorbank, Institut für Experimentelle Pathologie (dkfz),<br>Deutsche Krebsforschungszentrum, in Neuenheimer Feld 280,<br>Postfach 101949, D-6900 Heidelberg 1, Germany  |
|           | Centro Substrati Cellulari, Istituto Zooprofilattico Sperimentale<br>della Lombardia, e dell'Emilia, Via A. Biandri 7, I-25100 Brescia,<br>Italy                         |
|           | National Bank for Industrial Microorganisms and Cell Cultures<br>(NBIMCC), Blvd. Lenin 125 BL 2, V floor, Sofia, Bulgaria  |
|           | National Collection of Agricultural and Industrial<br>Microorganisms (NCAIM), Dept. Microbiology, University of<br>Horticulture, Somloiu 14-16, H-1118 Budapest, Hungary |
| Japan     | Japanese Cancer Research Resources Bank (JCRB), National<br>Institute of Hygienic Sciences, Kami-Yoga, Setagaya-Ku, Tokyo  |
|           | General Cell Bank, Institute of Physical and Chemical Research<br>of RIKEN, 2-1 Hirosawa, Wako, Saitama, 351-01  |
| Australia | Commonwealth Serum Laboratories, 45 Poplar Road, Parkville,<br>Victoria 3052   |

---

<sup>a</sup> Modified from R. I. Freshney, *Culture of Animal Cells*, 1994, 3rd ed., Wiley-Liss, New York.

Information on cell lines is also available through a number of data banks ([20](https://www.atcc.org)) such as ATCC ([www.atcc.org](http://www.atcc.org)) and ECACC ([www.ecacc.org.uk](http://www.ecacc.org.uk)).

### Bibliography

1. W. I. Schaeffer (1990) *In Vitro Cell. Dev. Biol.* **26**, 97–101.
2. D. W. Barnes, D. A. Sirbasku, and G. H. Sato, eds. (1984) *Methods for Serum-Free Culture of Epithelial and Fibroblastic Cells*, Alan R. Liss, New York, pp. 3–24.
3. R. Maurer (1992) in *Animal Cell Culture, a Practical Approach*, 2nd ed. R. I. Freshney, ed., IRL Press at Oxford University Press, Oxford pp. 15–46.

4. D. W. Jayme and D. F. Gruber (1994) in *Cell Biology, a Laboratory Handbook*, J. E. Celis, ed., Academic Press, New York, pp. 18–24.
5. L. Hayflick, P. S. Moorhead (1961) *Exp. Cell Res.* **25**, 585–621.
6. C. S. Stulberg, W. D. Peterson, Jr., and W. F. Simpson, (1976) *Am. J. Hematol.* **1**, 237–42.
7. W. A. Nelson Rees, D. W. Daniels, and R. R. Flandermeyer (1981) *Science* **212**, 446–452.
8. R. J. Hay (1991) *Dev. in Biol. Stand.* **75**, 193–204.
9. G. T. McGarrity, D. G. Murphy, and W. W. Nichols eds., *Mycoplasma Infection of Cell Cultures*, Plenum Press, New York, pp. 87–104.
10. R. Dulbecco and J. Elkington (1973) *Nature* **246**, 197–199.
11. B. Westermark (1974) *Int-J. Cancer* **12**, 438–451.
12. F. C. S. Ramaekers, J. J. G. Puts, A. Kant, O. Moesker, P. H. K. Jap, and G. P. Vooijs (1982) *Cold Spring Harbor Symp. Quant. Biol.* **46**, 331–339.
13. E. Heyderman, K. Steele, and M. G. Ormerod (1979) *J. Clin. Pathol.* **32**, 35–39.
14. T. Ikeda, N. Sawada, M. Satoh, and M. Mori (1998) *J. Cell. Physiol.* **175**, 41–49.
15. T. R. Chen (1988) *Cytogen. Cell Genet.* **48**, 19–24.
16. K. G. Steube, D. Grunicke, and H. G. Drexler (1995) *In Vitro Cell. Dev. Biol.* **31**, 115–119.
17. G. N. Stacey, B. J. Bolton, D. Morgan, S. A. Clark, and A. Doyle (1992) *Cytotechnology* **8**, 13–20.
18. G. Stacey, B. Bolton, A. Doyle, and B. Griffiths (1992) *Cytotechnology* **9**, 211–216.
19. A. Doyle, R. Hay, and B. E. Kirsop, eds. (1990) *Living Resources for Biotechnology*, Cambridge Univ. Press, Cambridge, U.K., pp. 1–15.
20. A. Doyle, R. Hay, and B. E. Kirsop, eds. (1990) *Living Resources for Biotechnology*, Cambridge Univer. Press, Cambridge, England, pp. 17–49.

## Cell Lineage

The daunting complexity of vertebrate organisms has led investigators to adopt a reductionist approach to their investigation. This approach is based on the assumption that an otherwise impossibly complicated system can be effectively understood through an analysis of its constituent cells and molecules. A reductionist concept of development envisions the process as the assembly of elementary units in many combinations and permutations, unfolding over time in stages that escalate progressively in their complexity. The transition from one stage to the next is influenced by the conditions that prevailed during the earlier period. Developmental biology thus seeks to infer the mechanisms of morphogenesis and histogenesis from the properties, activities, and responses of the precursors that give rise to the various lineages of the cells of the body.

That the concept of lineage should be an important one in developmental biology is something of a paradox. In a simplistic sense, it seems odd that lineage can be a meaningful issue, when all of the cells of an animal descend from a single precursor, the fertilized zygote. Distinct lineages, however, arise very early in embryogenesis. Formation of a body plan in lower vertebrates, including amphibia, is a process that begins in the oocyte, continues after fertilization and cleavage, and is finally completed at gastrulation (1, 2). Polarity in the oocyte causes cytoplasmic constituents to be distributed asymmetrically in the egg and newly formed blastula (3-6). This asymmetric distribution establishes different lineages in blastomeres, which do not each inherit identical cytoplasmic determinants from their predecessors. The resulting differential inheritance creates cellular

differences between the blastomere-derived lineages and allows cells to interact with one another (7, 8). Regional inequalities and cellular interactions may then be expressed as functional differences between cells in different locations. These processes form an ordered cascade that leads to the establishment of a pattern of small groups of like cells that follow a developmental agenda common to themselves, but different from those followed by other small groups of cells located elsewhere in the embryo (4). Functional differences between cells located in particular sites in the developing embryo are necessary for the differential cellular behaviors that underlie morphogenetic movements, like gastrulation and neurulation (9-11). It is also these different cellular loci, and the progeny to which they give rise, that provide developmental meaning to the concept of lineage.

In mammals, the establishment of lineages is delayed relative to that in amphibian embryos. Early cleavage cells of mammalian embryos appear to be identical to one another and totipotent (12). Each of the blastomeres of a two-cell embryo is capable of generating a complete adult, and will do so if the embryo is split (13, 14). Destruction of a single blastomere, moreover, may not prevent the normal development of a mature animal. **Chimeric** animals can also be constructed by combining morulae (15). It is thus unlikely that localized determinants play a role in mammalian development prior to or during the eight-cell stage. The pattern of cell diversification must be generated at subsequent stages, through the interactions of cells with one another and with their microenvironment.

## 1. Lineage as a Pedigree

Lineage, understood as a set of cells for which the ancestry can be traced back to a single progenitor, can possibly be delineated in simple transparent embryos, such as that of *Caenorhabditis elegans* by direct observation (16, 17). Cells can be recognized and followed from egg to adult in a living embryo as they divide, migrate, differentiate, and/or die. Keeping track of the cells is facilitated by the fact that nematode development is virtually invariant. Because the process is so constant, it is possible to predict the fates of descendant cells from their positions on the lineage tree. It is, of course, not possible to use similar direct observation by itself to determine the pedigrees of the cells of an adult mammal. Not only are the numbers of cells in a mammalian embryo too large to permit the progeny of unmarked individual cells to be followed, the details of cell lineage in mammals may also show random variations, even between genetically identical individuals. Lineage can possibly be followed in mammals by marking individual precursor cells in one or another location, for example, with a replication-deficient **retrovirus** that carries the gene encoding *Escherichia coli* **beta-galactosidase** (18-20). The activity of the bacterial enzyme or the protein can then be detected histochemically or immunocytochemically, respectively, thereby permitting all of the cells of the clone that ultimately develops from the single virus-infected progenitor to be identified. This type of analysis has been especially powerful in studies of the development of the central nervous system (CNS).

Studies of *C. elegans* have been particularly important in establishing that the terminally differentiated cells of an adult animal do not necessarily constitute the lineal descendants of a single **founder cell** (16, 17, 21). Different phenotypic classes of cells (with the exception of intestinal and germ cells) in *C. elegans* thus arise from several founder cells, which themselves are found on different arms of the lineage tree. The descendants of a single primitive embryonic precursor, therefore, may become cells of the hypodermis, musculature, and nervous system of the mature animal (21). Lineage in the straightforward pedigree sense thus does not predict the commitment of multipotent precursor cells to a specific phenotype. That choice, the moment when a cell's fate has become determined, may be made long before the terminally differentiated phenotype becomes evident in the appearance of the cell, and may be shared by cells originating from more than a single progenitor (22). Determination in this sense implies that a cell has undergone a lasting and autoperpetuating change that will forevermore distinguish that cell and its progeny from other kinds of cell and will commit the determined cell and its descendants to a common specialized developmental program. Cells may become determined either before, or when, they overtly differentiate (display their ultimate phenotype in a detectable manner); moreover, all of the cells that



are determined to follow a shared developmental program do not necessarily arise from a common precursor, and multipotent precursor cells may divide asymmetrically so that not all of their progeny are identically determined. The observation that cells of a particular terminally differentiated phenotype may not share an origin from a single “founder” has led the concept of lineage to acquire a meaning distinct from its linkage to cellular pedigrees.

## 2. Lineage and Determination

Cells of embryos that share a determined fate, as well as those that descend from a common founder, are often referred to as members of a common lineage. For example, the multipotent [stem cells](#) that are found in the neuroepithelium of the neural tube are self-sustaining, but they also give rise to progeny determined to be neurons or glia; following determination, the cells are said to be developing in neuronal or glial lineages (23). Lineage is used to denote their common determined phenotype, not their origin from a common precursor. The neuroepithelial stem cells proliferate in response to the actions of defined mitogens, including [epidermal growth factor](#) (EGF) and [fibroblast growth factor](#) (FGF), and differentiate as neurons or glia under the influence of a variety of microenvironmental factors. Factors that have been found under at least some circumstances (*in vitro* or *in vivo*) to promote the restriction of subsets of central nervous system stem cells to the glia lineage include [sonic hedgehog](#) (Shh) (24), ciliary neurotrophic factor (CNTF) (25), leukemia inhibitory factor (LIF) (25), and bone morphogenetic protein 4 (BMP4) (26).

The generation of the phenotypic diversity of cells in various brain regions occurs in stages, progressing from pluripotentiality to lineage restriction (a phenomenon that prevents cells from giving rise to certain types of progeny), to the determination of specific phenotypes, and finally to the manifestation of the features of the terminally differentiated cells (27). Because of their responsiveness to their microenvironment, transplanted stem cells can integrate into a variety of brain regions, and thus potentially may provide a tool for treating the diseased or injured brain (23). The progenitor cells of the vertebrate brain, however, are themselves heterogeneous and thus members of different lineages (27). These differences help to explain how and when neuroepithelial cells begin to respond to the molecular signals provided by their microenvironment. Environmental factors, be they mitogens, trophic [growth factors](#), or [morphogens](#), can only affect cells that have prepared themselves to be affected. To respond to an extracellular factor, the responding cell must have previously acquired the relevant receptors. Most of the stem cells of the vertebrate nervous system express the [intermediate filament](#) protein, nestin (23, 25). Lineages of progenitors, however, can be defined by the expression in common of molecular properties that distinguish the cells of the lineage from other progenitors in the neuroepithelium; such lineage-defining molecular traits include expression of the *achaete/scute* family of basic helix–loop–helix (bHLH) [transcription factors](#), **EGF receptors**, ventricular zone gene 1 (*vzg-1*), and the embryonic neural [cell adhesion molecule](#) (E-NCAM) (27). Progenitors in different locations may respond differently to the gliogenesis-promoting factors noted above.

## 3. Lineages in the Development of the Neural Crest

A source of precursor cells in developing vertebrates that has been extremely valuable for studies of the role(s) of cell lineage in development has been the neural crest (28, 29). This is a transitory structure that appears during embryogenesis and disappears as its component cells disperse, migrating through the embryo to give rise to a wide variety of terminally differentiated cell types (30, 31). The lineages of cells as diverse as melanocytes, fibroblasts, endocrine cells, smooth and skeletal muscle, cartilage, bone, meningeal cells, Schwann cells, satellite cells, enteric glia, autonomic, enteric, and sensory neurons can all be traced back to the neural crest. A great deal of investigation has been devoted to understanding when crest-derived cells become determined, and what factors cause them to become so. Considerable evidence indicates that at least some of the cells of the premigratory crest are pluripotent. Clones of these cells give rise in culture to several different classes of terminally differentiated cells (32-35). Lineage studies that have marked (by microinjection) individual cells of the premigratory crest (36, 37), as well as early-migrating crest-

derived cells (38), have shown that progeny of the marked precursors reach different destinations and thus give rise to a variety of cell types. Furthermore, cultures of stem cells can be prepared from the neural crest (39, 40), just as they can from the neuroepithelium of the neural tube (23, 25).

Although it is clear that the neural crest contains multipotent cells, the population of cells in the premigratory crest is actually heterogeneous. Some of the cells, even in the premigratory crest itself, are already determined (41). Cells determined to develop in the melanocytic lineage, for example, migrate selectively along the dorsolateral pathway (42). Among the uncommitted cells, there appears to be a progressive loss of developmental potential that occurs as a function of age (28, 29, 43). Whether the loss of developmental potential is simply determined by the age of the cells, or by the signals the cells receive as a result of their position in the embryo, has yet to be determined.

Some of the crest-derived cells that arrive in the organs that lie at the ends of their migratory pathways are still multipotent when the crest-derived cells arrive. The bowel has been carefully investigated in this regard and is a good model to illustrate the role of lineages in neural crest development (44). Two techniques have been employed to investigate the developmental potential of the crest-derived cells that colonize the gut. One, back-transplantation, consists of placing a segment of avian gut that has received a complement of émigrés from the neural crest in a migration pathway of a younger host embryo (45, 46). When this is done, crest-derived cells leave the grafts and remigrate in their new hosts. Where they go depends on where the grafts are placed. If the grafts are situated so as to replace the host's crest in locations that normally provide crest-derived émigrés for the bowel, the donor cells migrate to the host's gut and participate in forming the enteric nervous system (ENS). On the other hand, if the grafts are placed anywhere else in the embryo, the enteric crest-derived cells of the donor fail to migrate to the host's bowel. Instead, they migrate to peripheral nerves, sensory, and sympathetic ganglia. Donor cells that arrive in sympathetic ganglia become catecholaminergic sympathetic neurons, even though catecholaminergic neurons are not found in the enteric ganglia of the avian gut. Donor cells that find themselves in peripheral nerves develop as Schwann cells, even though the ENS contains glia rather than Schwann cells. Enteric crest-derived cells thus learn nothing from their previous journey to the bowel and evidently remember nothing of it. Crest-derived cells from the donor gut, however, never give rise to muscle, melanocytes, or connective tissue, indicating that they have lost at least some of developmental potential of their ancestors; nevertheless, they remain pluripotent to the extent that they are able to give rise to cell types that they would not have done had they been left *in situ*.

Clonal analysis of crest-derived cells isolated from the embryonic avian (47) or fetal mammalian bowel (48, 49) has also suggested that the population of crest-derived émigrés that colonizes the gut is a heterogeneous mixture of multipotent and determined cells. One or more terminally differentiated types of cells may be found in clones arising from individual precursors. None of the clones developing *in vitro* from enteric crest-derived cells contain melanocytes, endocrine cells, or skeletal muscle, although some contain smooth muscle. The heterogeneity of the crest-derived cell population that colonizes the bowel is reflected in the birthdays of the different types of enteric neuron (50). Some of these cells, containing acetylcholine or serotonin, are born very early in embryonic life, beginning in the mouse at E8.5, an age that precedes the arrival in the bowel of even the most precocious of enteric neurons. The crest-derived population that arrives in the gut, therefore, must contain a mixture of already-determined postmitotic neurons, as well as the undetermined pluripotent precursors that were revealed by the experiments involving back-transplantation and clonal analysis described above. On the other hand, other types of enteric neuron, such as those that contain calcitonin gene-related peptide (CGRP), do not even begin to leave the mitotic cycle until after the last serotonergic neurons have been born. The fact that the timing of the birthdays of enteric neurons is a function of their phenotypic class was among the earliest observations to suggest that multiple lineages of crest-derived precursors participate in the formation of the ENS. It is at least plausible that neurons in the same lineage might be born at similar times.

Confirmation of the idea that the ENS is formed by multiple lineages of crest-derived precursors came from a study that combined an *in vitro* investigation of enteric neuronal development with an

analysis of the effects of the targeted deletion of the *mash-1* transcription factor (51). It had previously been known that cells in the developing mammalian (but not avian) bowel are transiently catecholaminergic (TC) (52-55), and that these cells are derivatives of the neural crest (56-58). The TC cells were found not to be neurons, because they proliferate (56, 57, 59), although it was postulated that the TC cells might be neural precursors. That hypothesis was confirmed, and the disappearance of TC cells was found to occur not because they die, but because they lose their catecholaminergic properties as they acquire the characteristics of their terminally differentiated phenotype (56).

When it was discovered that TC cells express the same cell-surface differentiation antigens as precursors in the sympathoadrenal lineage, it was proposed that enteric and sympathetic neurons arise from a common lineage of crest-derived progenitors (60, 61). This idea was furthered by the demonstration that *mash-1* is expressed during the development of both sympathetic and enteric neurons (62, 63). On the other hand, knockout of *mash-1* is accompanied by the loss of almost the entire sympathetic nervous system (the superior cervical ganglion is spared), while the development of the ENS (except for that of the esophagus, which is lost) is merely delayed (63, 64).

The shared features expressed by sympathoadrenal and ENS precursors made it possible to test the hypothesis that sympathetic and enteric neurons are derivatives of a common precursor lineage.

**Complement**-mediated lysis of all of the cells isolated from the bowel that express **antigens** common to the sympathoadrenal lineage reduces the numbers of enteric neurons developing *in vitro*, but it does not prevent their formation (51). No serotonergic neurons develop *in vitro* following lysis of TC cells, although neurons expressing CGRP appear normally. The cells in the developing bowel that express *mash-1* were found to be TC cells, and the knockout of *mash-1* in mice was demonstrated to result in the complete loss of early-developing serotonergic neurons, while late-born CGRP-containing cells appear on schedule. The apparent delay in timing of the formation of the ENS in *mash-1* knockout mice is thus due to the selective loss of the early-developing lineage of enteric neurons. These observations established for the first time that the ENS is, in fact, derived from at least two lineages of crest-derived precursors. One of these, which may be common to the precursors of sympathoadrenal cells, obligatorily expresses *mash-1*, gives rise to neurons that are born early, and is the origin of all of the serotonergic neurons of the bowel. The other lineage is never catecholaminergic, does not express *mash-1*, gives rise to neurons that are born late, and is the origin of peptidergic enteric neurons, including those that contain CGRP.

Since the discovery of the role of *mash-1* in the development of one lineage of enteric neurons, additional lineages or sublineages of enteric neuronal precursor have been discovered. All of the crest-derived cells that colonize the bowel appear to express p75<sup>NTR</sup>, the common neurotrophin receptor (44, 56, 65). The cells (from the vagal and sacral regions of the neural crest) evidently also express Ret (66, 67), which must be stimulated by its functional ligand, glial cell line-derived neurotrophic factor (GDNF) (68-70). Actually, in order to stimulate Ret, GDNF has to form a complex with an  $\alpha$ -component, GFR $\alpha$ -1 (71-73). If Ret (74, 75), GDNF (69, 76), or GFR $\alpha$ -1 (77) are knocked out, there are no neurons or glia in the entire bowel below the esophagus and in the immediately adjacent region of the stomach. GDNF is a mitogen for early crest-derived enteric neuronal precursors and later supports their development as neurons (78-80). Because the defect in the *mash-1* knockout mice is much more limited than that seen after the elimination of Ret or the components of its ligand complex, it follows that the *mash-1*-dependent lineage is a subset of a larger lineage of GDNF-dependent precursors.

More limited sublineages have been identified by the dependence of subsets of enteric neurons on other growth factor/receptor complexes, such as a still-to-be identified neuropoietic cytokine that activates the  $\alpha$  component of the receptor for CNTF (65). Other factors, such as NT-3/TrkC (81) and laminin-1 (82, 83), have been identified that promote the development of enteric neurons, but these factors do not define particular lineages. Endothelin-3 (ET-3) and its preferred receptor, endothelin B (ET<sub>B</sub>), also appear to be necessary for a subset of enteric neurons (84, 85). This

factor/receptor complex, however, does not appear to define a particular lineage of enteric neurons. When either ET-3 (84) or ET<sub>B</sub> (85) is knocked out in mice, the terminal colon becomes aganglionic, just as it does in humans with Hirschsprung's disease, some of whom also display ET-3 (86) or ET<sub>B</sub> (87) mutations. The lesion that results from deletion of ET-3 or ET<sub>B</sub> thus is geographically restricted, but within the affected zone all lineages of crest-derived enteric neurons are affected. ET-3 has been found to inhibit enteric neuronal development, and it has been postulated to act by preventing the premature differentiation of enteric neurons before their crest-derived have completed their colonization of the gut (78, 88).

#### 4. Summary

The concept of lineage is used in two ways in describing development. In the pedigree sense, lineage is used to refer to the descendants of a common progenitor. Sharing a lineage defined as a pedigree, however, does not necessarily mean that the descendants of the common progenitor share the same developmental program or fate. Lineage as a pedigree is thus independent of the determination of its members. In a second developmental sense, the term *lineage* is used to refer to cells that do follow a common developmental program, regardless of whether these cells trace their ancestry back to a single shared precursor. Usually they do not. Identity of determination thus is now the common denominator that defines lineage.

#### Bibliography

1. J. B. L. Bard (1990) *Morphogenesis: The Cellular and molecular processes of Developmental Anatomy*, Cambridge University Press, Cambridge, U.K.
2. J. P. Trinkaus (1984) *Cells into Organs: The Forces That Shape the Embryo*, 2nd ed., Prentice-Hall, Englewood Cliffs, NJ.
3. M. R. Rebagliati et al. (1985) *Cell* **42**, 769–777.
4. J. B. Gurdon (1992) *Cell* **68**, 185–199.
5. J. P. Vincent, G. F. Oster, and J. C. Gerhart (1986) *Dev. Biol.* **113**, 484–500.
6. M. Kirchner, J. Newport, and J. Gerhart (1985) *Trends Genet.* **1**, 41–47.
7. S. C. Guthrie and N. B. Gilula (1989) *Trends Neurosci.* **12**, 12–16.
8. J. D. Hardin and L. Y. Cheng (1986) *Dev. Biol.* **115**, 490–501.
9. H. Spemann (1938) *Embryonic Development and Induction*, Yale University Press, New Haven, CT.
10. W. C. Smith et al. (1993) *Nature* **361**, 547–549.
11. K. W. Y. Cho et al. (1991) *Cell* **66**, 1111–1120.
12. R. L. Gardner (1985) *Philos. Trans. R. Soc. Lond. (Biol.)* **312**, 163–178.
13. S. J. Kelly (1977) *J. Exp. Zool.* **200**, 365–376.
14. A. K. Tarkowski (1959) *Nature* **184**, 1286–1287.
15. A. McLaren (1976) *Mammalian Chimeras*, Cambridge University Press, Cambridge, U.K.
16. C. Kenyon (1985) *Philos. Trans. R. Soc. Lond. (Biol.)* **312**, 21–38.
17. J. E. Sulston and H. R. Horvitz (1977) *Dev. Biol.* **56**, 110–156.
18. D. D. Galileo et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 458–462.
19. C. Walsh and C. Reid (1995) *Ciba Found. Symp.* **193**, 21–40.
20. J. G. Parnavelas, M. C. Mione, and A. Lavdas (1995) *Ciba Found Symp.* **193**, 41–58.
21. J. E. Sulston, E. Schierenberg, and J. G. White (1983) *Dev. Biol.* **100**, 64–119.
22. B. Christ, H. J. Jacob, and M. Jacob (1977) *Anat. Embryol.* **150**, 171–186.
23. R. McKay (1997) *Science* **276**, 66–71.
24. D. M. Orentas and R. H. Miller (1996) *Dev. Biol.* **177**, 43–53.

25. K. K. Johe et al. (1997) *Genes Dev.* **10**, 3129–3140.
26. R. E. Gross et al. (1996) *Neuron* **17**, 595–606.
27. L. Lillien (1998) *Curr. Opin. Neurobiol.* **8**(1), 37–44.
28. N. M. Le Douarin and E. Dupin (1993) *J. Neurobiol.* **24**, 146–161.
29. D. J. Anderson (1989) *Neuron* **3**, 1–12.
30. N. M. Le Douarin (1982) *The Neural Crest*, Cambridge University Press, Cambridge, U.K.
31. N. M. Le Douarin (1986) *Science* **231**, 1515–1522.
32. K. Ito, T. Morita, and M. Sieber-Blum (1993) *Dev. Biol.* **157**, 517–525.
33. M. Sieber-Blum and A. M. Cohen (1980) *Dev. Biol.* **80**, 96–106.
34. M. Sieber-Blum et al. (1993) *J. Neurobiol.* **24**(2), 173–184.
35. A. Baroffio, E. Dupin, and N. M. Le Douarin (1988) *Proc. Natl. Acad. Sci. USA* **85**, 5325–5329.
36. M. Bronner-Fraser and S. E. Fraser (1988) *Nature* **335**, 161–164.
37. M. Bronner-Fraser and S. Fraser (1989) *Neuron* **3**, 755–766.
38. S. E. Fraser and M. Bronner-Fraser (1991) *Development* **112**, 913–920.
39. D. L. Stemple and D. J. Anderson (1993) *Dev. Biol.* **159**, 12–23.
40. D. Stemple and D. J. Anderson (1992) *Cell* **71**, 973–985.
41. P. D. Henion and J. A. Weston (1997) *Development* **124**(21), 4351–4359.
42. C. A. Erickson and T. L. Goins (1995) *Development* **121**, 915–924.
43. D. J. Anderson (1993) *Annu. Rev. Neurosci.* **16**, 129–158.
44. M. D. Gershon (1997) *Curr. Opin. Neurobiol.* **7**, 101–109.
45. T. P. Rothman et al. (1990) *Development* **109**, 411–423.
46. T. P. Rothman et al. (1993) *Dev. Dyn.* **196**, 217–233.
47. F. Sextier-Sainte-Claire Deville, C. Ziller, and N. M. Le Douarin (1994) *Dev. Biol.* **163**, 141–151.
48. L. Lo and D. J. Anderson (1995) *Neuron* **15**(3), 527–539.
49. L. Lo, L. Sommer, and D. J. Anderson (1997) *Curr. Biol.* **7**, 440–450.
50. T. D. Pham, M. D. Gershon, and T. P. Rothman (1991) *J. Comp. Neurol.* **314**, 789–798.
51. E. Blaugrund et al. (1996) *Development* **122**, 309–320.
52. P. Cochard, M. Goldstein, and I. B. Black (1978) *Proc. Natl. Acad. Sci. USA* **75**, 2986–2990.
53. G. M. Jonakait et al. (1979) *Proc. Natl. Acad. Sci. USA* **76**, 4683–4686.
54. G. Teitelman, T. H. Joh, and D. J. Reis (1978) *Brain Res.* **158**, 229–234.
55. M. D. Gershon et al. (1984) *J. Neurosci.* **4**, 2269–2280.
56. G. Baetge, J. E. Pintar, and M. D. Gershon (1990) *Dev. Biol.* **141**, 353–380.
57. G. Baetge and M. D. Gershon (1989) *Dev. Biol.* **132**, 189–211.
58. G. Baetge, K. A. Schneider, and M. D. Gershon (1990) *Development* **110**, 689–701.
59. G. Teitelman et al. (1981) *Dev. Biol.* **86**, 348–355.
60. D. J. Anderson et al. (1991) *J. Neurosci.* **11**, 3507–3519.
61. J. F. Carnahan, D. J. Anderson, and P. H. Patterson (1991) *Dev. Biol.* **148**, 552–561.
62. F. Guillemot and A. L. Joyner (1993) *Mech. Dev.* **42**, 171–185.
63. F. Guillemot et al. (1993) *Cell* **75**, 463–476.
64. L. Lo et al. (1994) *Perspect. Dev. Neurobiol.* **2**, 191–201.
65. A. Chalazonitis et al. (1998) *Dev. Biol.* **198**(2), 343–365.
66. V. Pachnis, B. Mankoo, and F. Costantini (1993) *Development* **119**, 1005–1017.
67. P. L. Durbec et al. (1996) *Development* **122**, 349–358.

68. P. Durbec et al. (1996) *Nature* 381–793.
69. S. Jing et al. (1996) *Cell* **85**, 1113–1124.
70. M. Trupp et al. (1996) *Nature* **381**, 785–789.
71. M. Trupp et al. (1997) *J. Neurosci.* **17**, 3554–3567.
72. A. M. Davies et al. (1997) *Neuron* **19**, 485.
73. J. Widenfalk et al. (1997) *J. Neurosci.* **17**, 8506–8519.
74. A. Schuchardt et al. (1994) *Nature* **367**, 380–383.
75. J. J. S. Treanor et al. (1996) *Nature* **382**, 80–83.
76. M. S nchez et al. (1996) *Nature* **382**, 70–73.
77. G. Cacalano et al. (1998) *Neuron* **21**, 53–62.
78. C. J. Hearn, M. Murphy, and D. Newgreen (1998) *Dev. Biol.* **197**, 93–105.
79. R. O. Heuckeroth et al. (1998) *Dev. Biol.* **200**(1), 116–129.
80. A. Chalazonitis et al. (1998) *Dev. Biol.* **204**, 385–406.
81. A. Chalazonitis et al. (1994) *J. Neurosci.* **14**, 6571–6584.
82. T. P. Rothman et al. (1996) *Dev. Biol.* **178**, 498–513.
83. A. Chalazonitis et al. (1997) *J. Neurobiol.* **33**, 118–138.
84. A. G. Baynash et al. (1994) *Cell* **79**, 1277–1285.
85. K. Hosoda et al. (1994) *Cell* **79**, 1267–1276.
86. P. Edery et al. (1996) *Nat. Genet.* **12**(4), 442–444.
87. E. G. Puffenberger et al. (1994) *Cell* **79**, 1257–1266.
88. J. J. Wu, T. P. Rothman, and M. D. Gershon (1997) *Neurosci. Abst.* **23**, 24.

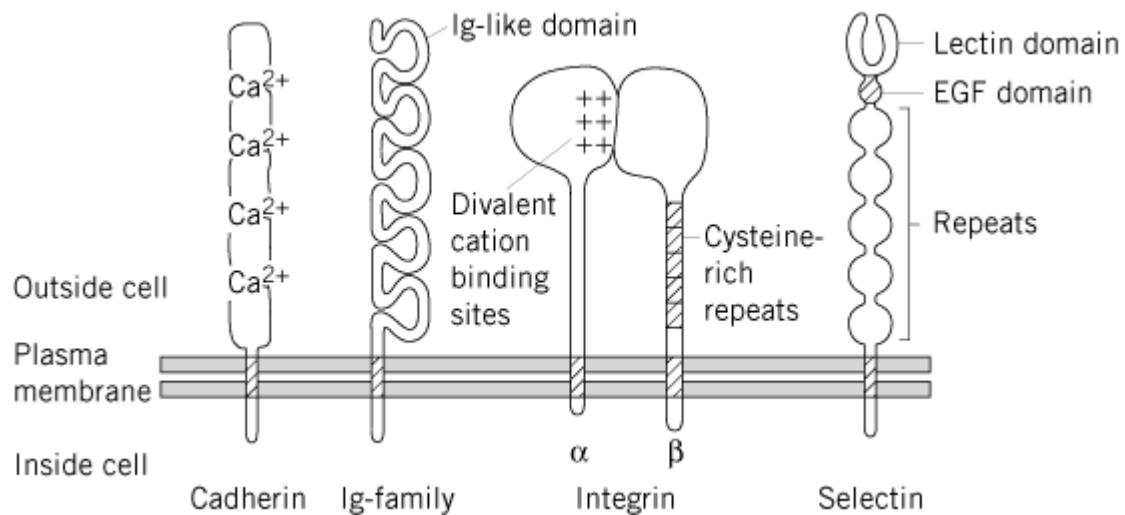
### Suggestions for Further Reading

89. B. Alberts et al. (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing, New York, pp. 1037–1107.
90. N. M. Le Douarin and E. Dupin (1993) *J. Neurobiol.* **24**, 146–161.
91. L. Lillien (1998) *Curr. Opin. Neurobiol.* **8**(1), 37–44.

## Cell Surface Adhesion Receptors

Cell surface adhesion receptors are molecules that mediate cell adhesion by binding to other molecules on the surface of an adjacent cell or to a component of the [extracellular matrix](#). Many adhesion receptors have been described. Most of them belong to one of a small number of families of related molecules in which individual members share the same basic molecular structure. The principal families are the [cadherins](#), the [immunoglobulin](#) (Ig) superfamily, the [integrins](#), and the selectins (Fig. 1) as follows.

**Figure 1.** Diagram showing the generalized molecular structures of the four principal families of cell-surface adhesion receptors. (From Ref. 7, with permission of BMJ Publishing Group.)



## 1. Cadherins

In most tissues, a major contribution to intercellular adhesion is made by calcium-dependent cell adhesion molecules known as cadherins (1). In general, these are simple transmembrane glycoproteins. The extracellular domain has an adhesion site toward the *N*-terminal region and several **calcium-binding** sites. Adhesive binding is mainly homophilic: a cadherin molecule on one cell binds to another cadherin molecule of the same type on the next cell. Linkage of the cytoplasmic **domain** to the [cytoskeleton](#) through proteins known as catenins is necessary for cadherin function. The best characterized is epithelial cadherin, E-cadherin. This appears very early in development, when it is involved in compaction of the eight-cell embryo and cell polarization. In adult epithelia—for example, intestinal epithelium—it is present on the lateral cell surfaces but is concentrated in intercellular junctions, known as the zonulae adherentes (see [Cell Junctions](#); [Intermediate Junction](#)), which ring the apicolateral margins of cells. These junctions are characterized by a cortical ring of cytoskeleton, the major component of which is [actin](#). Cadherins are also present in the nervous system, where they play an essential role in neural development.

The adhesive glycoproteins of the other major intercellular junctions of epithelia, the **desmosomes** are also members of the cadherin family (2). Their extracellular domains are very like those of cadherin, but their cytoplasmic domains differ, being specialized for forming desmosomal plaques and, thereby, attachment to the keratin [intermediate filament](#) cytoskeleton, rather than to actin.

## 2. Immunoglobulin Superfamily

Another major group of cell-to-cell adhesion molecules are members of the Ig superfamily. Their extracellular portions are characterized by the presence of at least one, and usually multiple, immunoglobulin-like domains (3, 4). Included in this group are several nervous system [cell adhesion molecules](#), such as the neural cell adhesion molecule (N-CAM), L-1, and TAG, which are involved in neuronal guidance and fasciculation. Several members of the immunoglobulin family are concerned with [antigen](#) recognition and adhesion in [T cells](#). These include the following: the [T-cell receptor](#) (CD3) and its co-receptors CD4 and CD8, which together recognize the complexes of antigen peptide and [major histocompatibility complexes](#) on other cells; the major histocompatibility complex molecules themselves; and lymphocyte function related antigen 2 (LFA-2 or CD2), a receptor for another immunoglobulin-like molecule, LFA-3, expressed on other cells. Another group of immunoglobulin-like cell adhesion molecules includes the so-called intercellular adhesion molecules, ICAM-1, -2, and -3, which are more widely expressed, for example, on epithelial and endothelial cells, and V-CAM on endothelial cells. These are involved in the inflammatory response.

The immunoglobulin superfamily is large and diverse, probably because the basic structure of the immunoglobulin domain is versatile and readily adaptable to different binding functions (see [Immunoglobulin Structure](#)). Among these molecules, however, the T cell receptor and the immunoglobulins themselves have somatically-variable domains necessary for antigen recognition. Members of the superfamily are even present in insects, where they are involved in nerve connections; the association of immunoglobulin-like domains in cellular recognition preceded the immune system in evolution.

### 3. Integrins

Both cell-to-cell and cell-to-matrix receptors are contained within this large family of adhesion molecules. Integrins are heterodimers consisting of one  $\alpha$  chain and one  $\beta$  chain, both of which are necessary for adhesive binding (5). Seventeen different  $\alpha$  chains and eight different  $\beta$  chains are now known. Integrins may be classified into subfamilies according to which  $\beta$  subunit is involved in the complex. Thus  $\beta 1$  integrin may associate with one of nine different  $\alpha$  subunits, to give a series of matrix receptors of differing specificity. The  $\beta 2$  integrins, on the other hand, are a family of cell-to-cell adhesion molecules of lymphoid cells with three alternative  $\alpha$  subunits. The classification is made more complicated because some  $\alpha$  subunits can associate with different  $\beta$  subunits (for example,  $\alpha 6\beta 1$  and  $\alpha 6\beta 4$ ).

Some integrins are apparently quite specific in their ligand-binding properties—for example,  $\alpha 5\beta 1$  for the Arg–Gly–Asp tripeptide sequence of [fibronectin](#)—whereas others are promiscuous—for example,  $\alpha v\beta 3$ , once regarded as the vitronectin receptor, also binds fibronectin, [fibrinogen](#), von Willebrand factor, thrombospondin, and osteopontin. An interesting example is  $\alpha 4\beta 1$ , which binds both the IIIICS domain of fibronectin and V-CAM on endothelial cells. To complicate matters further, individual cell types usually express multiple integrins. A good example to consider here is the blood platelets that express predominantly  $\alpha IIb\beta 3$  (GPIIb/IIIa), which binds fibrinogen, fibronectin, von Willebrand factor, and vitronectin, but also lesser amounts of  $\alpha V\beta 3$ ,  $\alpha 5\beta 1$ ,  $\alpha 2\beta 1$  ([collagen](#)), and  $\alpha 6\beta 1$  ([laminin](#)).

### 4. Selectins

Most cellular interactions seem to entail homophilic or heterophilic [protein–protein interactions](#). However, the selectins constitute a family of cell adhesion proteins that bind to carbohydrate. Selectins have **lectin**-like carbohydrate-binding domains at their extracellular *N*-terminal extremities (6). One of these, L-selectin (L = leucocyte) is a “homing receptor,” mediating regionally specific adhesion of lymphocytes to endothelium in peripheral lymph nodes. This molecule is also involved in the adhesion of neutrophils to endothelium during the inflammatory response. Two other members of this family, E-selectin (E = endothelial) and P-selectin (P = platelet), also participate in the inflammatory response. E-selectin is unregulated on endothelial cells over a period of hours after stimulation by inflammatory mediators. P-selectin is contained within Weibel–Palade bodies of endothelial cells and platelet  $\alpha$  granules, from which it is rapidly mobilized on activation, mediating adhesion to neutrophils and monocytes.

### Bibliography

1. M. Takeichi (1990) Cadherins: a molecular family important for selective cell–cell adhesion. *Annu. Rev. Biochem.* **59**, 237–252.
2. D. R. Garrod, M. A. J. Chidgey, and A. J. North (1996) *Curr. Opin. Cell Biol.* **8**, 670–678.
3. T. A. Springer (1990) *Nature* **346**, 425–434.
4. G. M. Edelman and K. L. Crossin (1991) *Annu. Rev. Biochem.* **60**, 155–190.
5. R. O. Hynes (1992). *Cell* **69**, 11–25.
6. M. P. Bevilacqua and R. M. Nelson (1993) *J. Clin. Invest.* **91**, 379–387.
7. D. R. Garrod (1997). "Cell to cell and cell to matrix adhesion". In *Basic Molecular and Cell*



### Suggestions for Further Reading

8. D. R. Garrod, M. A. J. Chidgey, and A. J. North (eds.) (1999) *Adhesive Interactions of Cells*, JAI Press, Greenwich, CT. (A series of up-to-date reviews on cell adhesion receptors and intercellular junctions.)
9. “Cell adhesion and human disease.” *Ciba Found. Symp.* **189**. (A valuable collection of reviews on cell adhesion molecules and their involvement in human disease.)
10. M. Hortsch and C. S. Goodman (1991) *Annu. Rev. Cell Biol.* **7**, 505–557. (Adhesion receptors in *Drosophila*—the invertebrate situation.)

### CEN Sequences

[Autonomously replicating sequences](#) (ARS) are not stable *in vivo* during cell division and cannot be used readily as plasmid **vectors**. But, if certain DNA sequences that function as [centromeres](#) in yeast (CEN sequences) are grafted onto a **plasmid** already containing an ARS, the resulting vector plasmid is stabilized and segregates accurately during **mitosis** and **meiosis**. The CEN sequences are necessary for attaching [chromosomes](#) to the **mitotic spindle**, presumably through some specific connector proteins that join the [microtubules](#) to the centromere. In addition, CEN sequences can be used to construct linear chromosomes by preventing circularization by the addition of **telomeric** sequences that are normally at the end of chromosomes ([1](#), [2](#)).

Conserved features of *Saccharomyces cerevisiae* CEN SEQUENCES are confined to a region of about 120 bp. The highly conserved 8 bp at the left (PuTCACPuTG, where Pu = purine) constitute the left boundary of a functional CEN sequence. The right boundary lies within or just beyond the 25 bp at the right, with the consensus sequence TGT-T-TG-TTCCGAA—AAA, where - indicates no specific base. A mutant that lacks the left conserved element can still assemble into a centromere that is partially functional in mitosis and well functioning in meiosis. The sequences between the two conserved terminal DNA elements are not essential for centromere function. Their lengths can be increased by 50% or their sequence from 6 to 12% in GC content without measurable changes in mitotic or meiotic segregation of plasmids carrying such CEN mutations ([3](#)). The right boundary sequence appears to be a binding site for a protein, as evidenced by various inactive mutations, especially in the central CCG sequence, and by an exonuclease blocking assay ([4](#)). The left element, which carries the **palindromic** sequence CACPuTG, binds the **helix–loop–helix motif** protein CPF1 ([5](#)).

### Bibliography

1. L. Clarke and J. Carbon (1980) *Nature* **287**, 504–509.
2. C–L Hsiao and J. Carbon (1981) *Gene* **15**, 157–166.
3. L. Panzeri et al. (1985) *EMBO J.* **4**, 1867–1874.
4. R. Ng., J. Ness and J. Carbon (1986) *Basic Life Sci.* **40**, 479–492.
5. R. Niedenthal, R. Stoll, and J. H. Hegemann (1991) *Mol. Cell. Biol.* **11**, 3545–3553.

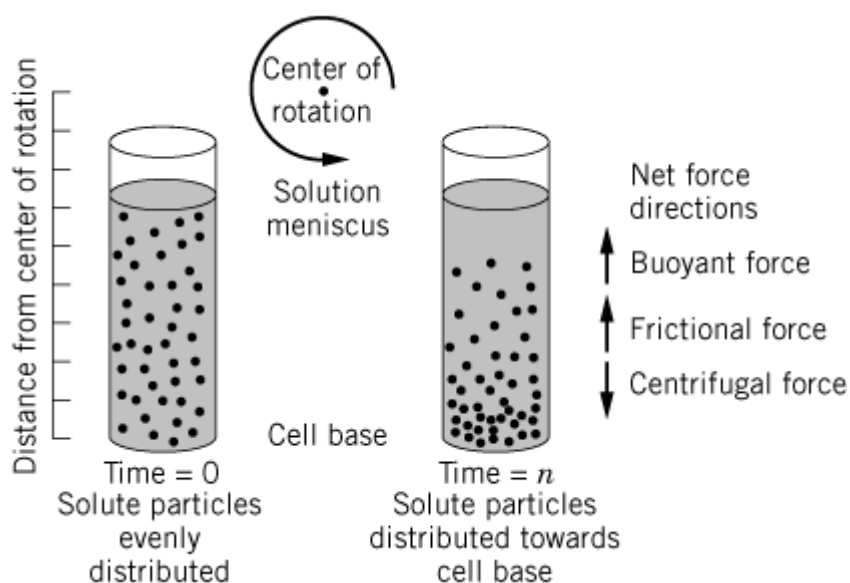
### Suggestion for Further Reading

6. L. Clarke and J. Carbon (1985) “The structure and function of yeast centromeres,” *Ann. Rev. Genet.* **19**, 29–57.

## Centrifugation

Centrifugation is a technique often employed during isolation or analysis of various cells, **organelles**, and biopolymers, including **proteins**, **nucleic acids**, lipids, and carbohydrates dissolved or dispersed in biologically relevant solvents (ie, typically aqueous buffers). It is also applicable to synthetic macromolecules dispersed or dissolved in nonaqueous, organic solvents. In this technique, the sample (comprising a liquid phase and a solute) is placed in a suitable vessel, and the vessel is spun in a centrifugal rotor. The centrifugal force created by the spinning rotor causes the solute sample to sediment out of solution (typically, though not always, toward the base of the vessel) (Fig. 1). The centrifugal force applied to the sample is akin to gravitational force (acceleration) and is measured in gravities, as in  $n \times G$  (the gravitational constant  $G$  equals  $6.6720 \times 10^{-11} \text{ N}\cdot\text{m}^2/\text{Kg}^2$ ). The extent to which the sample sediments toward the base of the vessel is a function of a series of complex interacting factors related to the properties of the solute alone and of the system as a whole (solvent and solute) (Fig. 1, Tables 1 and 2).

**Figure 1.** A schematic diagram of a centrifugation experiment.



**Table 1. The Forces Involved in Centrifugation**

| Proportionality Constant<br>or Force Equation | Cause or Definition | Dimensions |
|---|---------------------|------------|
|---|---------------------|------------|

|                             |   |                         |
|-----------------------------|---|-------------------------|
| Centrifugal force = $w^2rm$ | Centrifugation of solute of mass, $m$                                     | (g-cm)<br>(radians/sec) |
| Buoyant force = $-w^2rm_s$  | Solvent of mass $m_s$ ; displacement by sedimenting solute                | (g-cm)<br>(radians/sec) |
| Frictional force = $-fv$    | Resistance to movement of solute particle at velocity $v$ through solvent | g-cm/sec <sup>2</sup>   |

**Table 2. General Proportionality and Descriptive Constants for Centrifugation**

| Term   | Definition  | Dimensions                        |
|--|---|-----------------------------------|
| Angular velocity = $w = (2\pi r/60) \times \text{rpm}$ | $2\pi \left( \frac{\text{radians}}{\text{revolution}} \right) \left( \frac{\text{revolutions}}{\text{min}} \right) \left( \frac{1 \text{ min}}{60 \text{ sec}} \right)$ | radians/sec                       |
| Frictional coefficient = $f = 6\pi\eta r_s$            | Stokes' law frictional coefficient for a sphere of equivalent radius $r_s$ in a solvent of viscosity $\eta$ .   | g/sec                             |
| Diffusion coefficient $D$                              | Defined by Fick's first law (see Refs. 1 and 2) for an ideal solution at infinite dilution as $D = RT/Nf$ , e.g., 1 Fick = $10^7 \text{ cm}^2/\text{sec}$               | $\text{cm}^2/\text{sec}$          |
| Sedimentation coefficient $s$                          | The velocity of sedimentation per unit of centrifugal force   | sec                               |
| Solution density $\rho$                                | Mass per unit volume  | $\text{g}/\text{cm}^3$            |
| Solution viscosity $\eta$                              | The resistance to change within a fluid, a type of internal friction  | poise or dyne-sec/cm <sup>2</sup> |
| Solute partial specific volume $v$                     | The change in volume per unit mass  | $\text{cm}^3/\text{g}$            |
| Solute buoyant density                                 | Effective density determined by isopycnic centrifugation in a specific gradient medium  | $\text{g}/\text{cm}^3$            |
| Avagadro's number $N$                                  | $6.022 \times 10^{23}$  | mole <sup>-1</sup>                |
| Gas constant $R$                                       | $8.314 \times 10^7$   | ergs/K-mole                       |

The factors that define the entire sedimenting system include (1) the rate  $w$  at which the rotor spins; (2) the length of time for which the centrifugal force is applied; (3) the solvent density  $\rho$ ; and (4) the system temperature  $T$ . The factors specific to the sedimented solute include (1) the radial distance  $r$  of the solute from the central axis of rotation; (2) the solute [partial specific volume](#)  $v$ ; (3) the solute [diffusion](#) coefficient  $D$ ; (4) the solute **frictional coefficient**  $f$ ; and (5) the solute [sedimentation coefficient](#)  $s$ . During centrifugation (Fig. 1), the centrifugal force is opposed by [buoyant density](#) and frictional forces. Under these conditions it can be shown that the rate of sedimentation of the solute

molecules of weight-average molecular weight  $M_w$  and moving as a radial boundary  $r_b$  is defined by

$$\text{Boundary velocity} = v = \frac{M_w(1 - \bar{v}\rho)}{Nf} \omega^2 r_b. \quad (1)$$

Operationally, the technique can be divided into preparative and analytical modes, which differ principally in whether the actual recovery of the materials that have been centrifugally separated is desired or practical (Fig. 2).

**Figure 2.** Common preparative and analytical methods used in centrifugation.

| Preparative<br>Isolation and<br>Recovery of Solute(s) | Solute<br>Analytical<br>Information |
|---|-------------------------------------|
| Common Methods  | Common Methods                      |
| Selective pelleting                                   | Sedimentation<br>velocity           |
| Density gradient separations                          | Sedimentation<br>equilibrium        |
| Rate zonal  |                                     |
| Flotation   |                                     |
| Isopycnic   |                                     |

## 1. Preparative Centrifugation and Ultracentrifugation

Preparative centrifugation (Fig. 2) is routinely employed in isolating, separating, and purifying a variety of biological and cellular/subcellular components. Some examples include [chromatin](#), coated vesicles (see [Clathrin](#)), cytosolic proteins, DNA, **Golgi**, [inclusion bodies](#), **lipoproteins**, **lysosomes**, **Microorganisms**, **Microsomes**, [mitochondria](#), myelin, **nuclei**, **glycoproteins**, **RNA**, **Peroxisomes**, plasma [membranes](#), **polysomes**, proteoglycans, [ribosomes](#), synaptosomes, and **viruses**. Preparative centrifugation has two primary modes of separation. Solutes can be selectively pelleted from solution, as described here, or they can be subjected to separation along a solution density gradient (see [Density Gradient Centrifugation](#)).

### 1.1. Selective Pelleting or Differential Centrifugation

Clearing macromolecular components from solution is probably the most commonly employed preparative centrifugation technique. In this method the sample is subjected to centrifugation, usually, though not necessarily, at a constant rotor velocity. Over time, a pellet of sedimented material is deposited along the most radially distant wall of the tube containing the sample (Table 3). At a sufficiently high rotor speed (and therefore high centrifugal force), virtually all macromolecular solute components heavier than the solvent are pelleted out of solution over time. Because it is possible to select a variety of rotor velocities, sample tube dimensions, rotor types, centrifugation run times, etc., it is often possible to accomplish a significant level of purification by pelleting alone because the method of separation relies on the differential sedimentation rates of small and large solute components. For example, large components are typically removed by sedimentation at low speeds for short periods of time. If the supernatant from such a centrifugation is placed in a new container and centrifuged again at a higher rotor velocity for an appropriate period of time, it is possible to collect another distribution of particles from solution, this time of smaller dimensions

than the first. Therefore repeating this process yields a series of pelleted samples with distributions of solute molecules, each of which have accumulated under conditions differing only in their centrifugal forces and thus have different sedimentation rates. Most samples of biological interest are complex mixtures of interacting solutes in solution at widely different concentrations. Because the size, shape, and concentration distributions of these components are routinely broad, together they exhibit very broad distributions of sedimentation rates. Thus the technique of solute pelleting, though widely applied, is only rarely sufficient to purify a single component completely. Typically it is used to selectively enrich, often to a high degree, a preparation of cellular or subcellular components in a single component. For small biological solutes, such as soluble cytosolic proteins prepared by a protein isolation procedure, pelleting is routinely employed along with differential [precipitation](#) (**salting out**) using ammonium sulfate (see [Sulfate Salts](#)) to yield highly enriched, sedimented precipitates of the desired protein components.

**Table 3. Descriptive Constants for Preparative Centrifugation**

| Term                                | Definition  |
|-------------------------------------|---|
| Clearing factor, $k$                | time (hours) = $k/s$ , $s = \text{svedbergs}$ ; the value of $k$ varies with each rotor type and compares the relative rotor efficiencies for pelleting   |
| Clearing time, $t$                  | $t = k/s$ (see clearing factor, above); the time required to pellet a particle to the base of the tube  |
| Flotation Coefficient, $f_s$        | For particles which float in a particular solution in a centrifugal field, it is similar to the sedimentation coefficient, units are in sec.  |
| Relative centrifugal force or (RCF) | $\text{RCF} = 1.12r (\text{rpm}/1000)^2$ , where $r$ is in mm; the ratio of the centrifugal force at a specific rotation rate in rpm for a particular rotor at a particular radius $r$ from the center of rotation of that rotor to the gravitational force of the Earth at sea level.                                    |
| Square-root speed reduction law     | $\text{[Maximum allowed rotor speed]} = \frac{\text{Desired speed}}{\sqrt{\frac{\text{Max allowed solution density}}{\text{actual solution density}}}}$ Used to calculate the permissible rotor speed in rpm when using nonprecipitating gradient materials at densities greater than the manufacturer permits in a rotor |

Pelleting is also of great use when applied to harvesting viruses or cells from various growth media. When cells or viruses are used to produce various recombinant proteins of pharmaceutical interest, centrifugal harvesting is often employed to capture the biological component from the (usually large) fermentation broth. In these applications, a preparative centrifuge and rotor capable of accepting a continuous flow into and through the rotor are employed. The centrifugal force applied to the sample stream as it passes through the rotor is adjusted to permit selective pelleting of cells within the centrifuge rotor. The cell-depleted broth passes out of the rotor/centrifuge and into a receiver vessel, and, then the harvested cells are recovered from the centrifuge rotor.

#### Bibliography

1. K. E. Van Holde (1971) *Physical Biochemistry*, Prentice Hall, Englewood, Cliffs, NJ, pp. 70–121.

2. H. Fugita (1962) *Mathematical Theory of Sedimentation Analysis*. Academic Press, New York.

### Suggestions for Further Reading

3. K. E. Van Holde (1971) *Physical Biochemistry*, Prentice–Hall, Englewood Cliffs, NJ, pp. 70–121. A short, highly readable and concisely illustrated treatment of the mechanical model of centrifugation for the biologist.
4. D. Eisenberg and D. Crothers (1979) *Physical Chemistry with Applications to the Life Sciences*, Benjamin–Cummings, Menlo Park, CA, pp. 701–745. An excellent introductory physical chemistry text with a good review of fundamental flow equations and a number of useful examples.
5. C. Tanford (1961) *Physical Chemistry of Macromolecules*, Wiley, New York, pp. 317–456. A classic and still relevant text on the nature of transport processes as applied to biological systems.

## Centromeres

The centromere (kentron = center; meros = part) is the region of the mitotic chromosome that participates in chromosomal movement. The mitotic spindle attaches to a specialized structure at the centromere known as the kinetochore. The motor responsible for the movement of the chromosomes toward the spindle poles during mitosis is also located at the centromere. Morphologically, centromeres are distinguished by their appearance at metaphase and anaphase as constrictions in the chromatin and by their heterochromatin staining pattern. A chromosome without a centromere is described as acentric. It does not segregate properly during cell division and is rapidly lost in successive cell cycles.

Centromeres occupy different positions in chromosomes, and they are useful markers. Chromosomes with a single centromere are called monocentric. Chromosomes with centromeres near one end are called acrocentric, those with the centromere visible at or near the middle are metacentric, and those with the centromere truly at the end are telocentric. The monocentric chromosomes found in most metazoan plants and animals almost always have centromeres embedded in segments of heterochromatin. Some organisms, such as plants in the genus *Luzula*, have holocentric chromosomes with diffuse centromeres. In these chromosomes the spindle attaches to the centromeric heterochromatin that is distributed along the entire length of the chromosome (1). Mammalian tissue culture cells occasionally develop chromosomes with multiple centromeres (2). A chromosome with two centromeres is called dicentric. Dicentric chromosomes are often the result of chromosome breakage followed by fusion. A dicentric chromosome normally breaks at anaphase, when the centromeres are pulled in opposite directions. The holocentric chromosomes found in certain plants have developed as yet unknown mechanisms to avoid this problem.

Centromeres represent highly specialized chromosomal organelles. The simplest types of centromeres are found in the yeast *Saccharomyces cerevisiae*. The presence of these centromeric sequences, together with the autonomously replicating sequences that serve to direct DNA replication within exogenous plasmid DNAs allow these small minichromosomes to be stably maintained through cell division. The centromeric sequence allows the minichromosomes to be segregated in mitosis and meiosis with accuracy. Maintenance of minichromosomes through cell division also provides a simple assay for the definition of functional centromeric sequences in yeast.

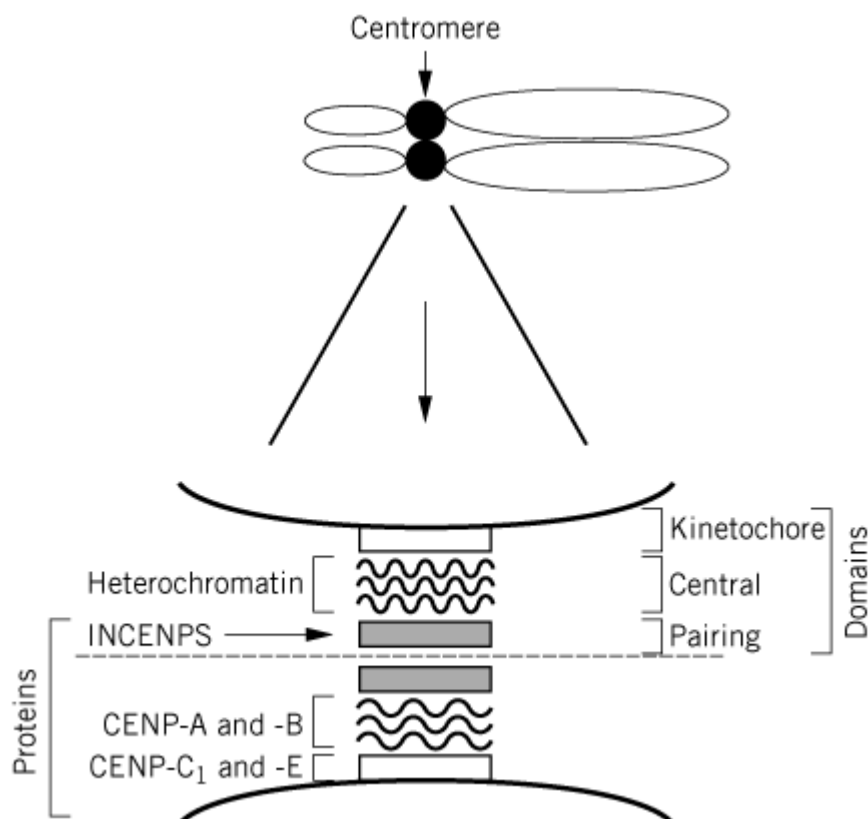
All of the yeast centromeres are functional as small segments of DNA about 1,000 base pairs (1 kbp)

or less in length. Nucleosomes are assembled in specific positions on these DNA sequences, and they contain a specialized histone H3 molecule (Cse4p) (16). Within each segment of centromeric DNA are similar nucleotide sequences. In the middle of the centromeric DNA is a 220-bp segment that constitutes the centromeric core, that contains the minimal sequence necessary for centromeric function. Only a single spindle fiber attaches to a yeast centromere (4). If this attachment is to the core, then the nucleoprotein complex assembled on the core also functions as a kinetochore.

Mammalian centromeres are considerably more complex than those in *S. cerevisiae*. Early experiments on the centromere of the human Y-chromosome established that much more DNA is required for the chromosomal segregation function than in yeast. Over 300 kbp of DNA, including 200 kbp of  $\alpha$ -satellite DNA, are required to generate a functional centromere, suggesting the epigenetic—i.e., meta-DNA—mechanisms establish centromere location in cells of metazoa (17), a notion supported by the appearance of neo-centromeres in mammalian cells over nonrepetitive DNA (18). For native centromeres, however, deletion of the  $\alpha$ -satellite DNA makes the centromere inactive (5).

The centromere can be broken up into three distinct structures, the kinetochore and the central and pairing domains (Fig. 1). The protein constituents of the centromere have been identified primarily through the use of autoantibodies from patients with rheumatic disease. These include the inner centromeric proteins (INCENP) and the CENP A, B, and C proteins. Immunological staining reveals that the INCENPs are in the pairing domain. The central domain contains dense chromatin known as constitutive heterochromatin. The kinetochore is anchored to this heterochromatin (Fig. 1). The DNA within this heterochromatin is composed primarily of various families of repetitive DNA (satellite DNA).

**Figure 1.** An expanded view of the centromere of a mammalian chromosome. The centromere is a bipartite structure. The central axis between the pairing domains is indicated by the dashed line. The relative positions of heterochromatin, the INCENPs, and CENPs are indicated within the kinetochore and the central and pairing domains.



The a-satellite family of DNA sequences (which comprises 5% of the human genome) is probably present at the centromeres of all human chromosomes. The basic repeat 171 bp long, occurs in large arrays up to  $3 \times 10^6$  bp long. The a-satellite repeats assemble a specialized chromatin structure that provides many insights into the molecular nature of constitutive heterochromatin. It has been long that nucleosomes are positioned with respect to the nucleotide sequence found in a-satellite DNA so that histone-DNA contacts would begin and end at particular sites (6). An important point is that a single nucleosome is believed to exist on every 171 bp repeat. In modified form, this DNA sequence also provides the foundation for crystallization of the nucleosome more (7). Several interesting specialized chromosomal proteins also associate with a-satellite DNA. A 10-kDa protein, high mobility group protein (HMG-I/Y) binds to a-satellite DNA specifically *in vitro* (8). The HMG-I/Y protein binds selectively in the minor groove of the double helix, associating with runs of six or more AT base pairs. In addition HMG-I/Y also probably recognizes certain secondary structural features of the DNA molecule. CENP-B specifically recognizes a 17-bp DNA sequence (5' CTTCGTTGGAA CGGGA 3') in a subset of a-satellite DNA repeats. The DNA-binding domain of CENP-B is necessary and sufficient for localization to the centromere *in vivo* (9). CENP-B contains anionic regions rich in aspartic and glutamic acid residues, which are characteristic of many proteins that interact with chromatin. CENP-B is found throughout the centromeric heterochromatin beneath the kinetochore plates (see [Kinetochore](#)). Although the exact functions of CENP-B in heterochromatin are not known, it may play a role in the higher order folding of centromeric chromatin through self-assembly mechanisms. Although CENP-B is found at all mammalian centromeres examined, the CENP-B recognition element is not found in all a-satellite DNA (10). It is possible that CENP-B can also be targeted to centromeres through interactions with other protein components.

CENP-A is also specifically associated with centromeric DNA and is especially interesting because it shares homology with core histone H3. Recent sequence analysis has revealed that CENP-A is a very specialized H3 variant (11). Each core histone and CENP-A have two domains, and amino-terminal tail domain that lies on the outside of nucleosomal DNA and a carboxyl-terminal, histone-fold domain that is involved in protein—protein interactions and in wrapping the DNA. CENP-A has a highly divergent amino-terminal tail domain and a relatively conserved carboxy-terminal, histone-fold domain. The targeting of CENP-A to centromeric DNA is directed by the histone-fold domain, and it is very likely that a highly specialized, dedicated chromatin assembly mechanism exists to enable the selective heterodimerization of CENP-A with histone H4 over centromeric DNA (19). Although the exact mechanism by which this is achieved is still unknown. Characterization of a CENP-A homologue in yeast (CSE4, Ref. 3) establishes that this specialized histone is essential for normal chromosomal segregation during mitosis. Genetic evidence suggest that CSE4 heterodimerizes with histone H4 (19). This establishes that specialized nucleosomal structures are present at the centromere, in addition to the presence of positioned nucleosomes. Clearly, the centromere represents a highly differentiated chromosomal domain even at the most fundamental level of chromatin structure.

CENP-C is a large 107-kDa protein that is highly hydrophilic and basic (pI = 9.4) (13). CENP-C is concentrated in a narrow band immediately below the inner kinetochoric plate at the interface between the chromosome and the kinetochore (Fig. 1). CENP-C binds DNA directly and is required for normal kinetochoric assembly. CENP-C is found only at the active centromere of a stable dicentric chromosome, suggesting a direct role in centromeric function (14). Other components of the centromere include CENP-E, which resembles microtubule-binding motor proteins. CENP-E is a 312-kDa polypeptide chain with a tripartite structure consisting of globular domains at the N- and C-termini, separated by a 15,000-residue  $\alpha$ -helical domain that is predicted to form coiled-coils. CENP-E colocalizes to the centromere and kinetochore during metaphase, but it is released from the centromere at the onset of anaphase, when it is degraded. CENP-E functions as a kinetochoric motor during the early part of mitosis (15) (See [Kinetochore](#)). The proper cohesion between centromeres on



sister chromatids prior to the metaphase-anaphase transition in mitosis, and centromere nondisjunction on sister chromatids during meiosis I are essential for proper chromosome behavior. Both are effected by the preferential association of the cohesion complex (120) with centromeres (16, 21); in addition to the cohesins, in metazoa, noncohesion proteins as MEI-S322 also associate with the centromere and mediate kinetochore assembly as well as sister chromatid cohesion during meiosis (22).

The assembly of the centromere and the mechanisms of kinetochoric association with chromatin offer perhaps our best opportunity to understand the construction of a specialized chromosomal domain at the biochemical level.

## Bibliography

1. S. Pimpinelli and C. Goday, *Trends Genet.* **5**, 310–313 (1989).
2. B.K. Vig and N. Paweletz, *Chromosoma* **96**, 275–283 (1988).
3. S. Stoler, K.C. Keith, K.E. Curnick, and M. Fitzgerald-Hayes, *Genes Dev.* **8**, 573–586 (1995).
4. J.B. Peterson and H. Ris, *J. Cell Sci.* **22**, 219–226 (1976).
5. C. Tyler-Smith, R. Oakley, and Z. Larin, *Nat. Genet.* **5**, 368–375 (1993).
6. R.T. Simpson, *Prog. Nucl. Acids Res. Mol. Biol.* **40**, 143–184 (1991).
7. K. Luger et al., *Nature* **389**, 251–260 (1997).
8. M.J. Solomon, F. Strauss, and A. Varshavsky *Proc. Natl. Acad. Sci. U.S.A.* **83**, 1276–1280 (1986).
9. A.F. Pluta et al., *Science* **270**, 5191–5194 (1995).
10. I.G. Golberg et al., *Mol. Cell Biol.* **16**, 5156–5168 (1996).
11. K.F. Sullivan, M. Hechenberger, and K. Masri, *J. Cell Biol* **127**, 581–592 (1994).
12. M.M. Smith et al., *Mol. Cell Biol.* **16**, 1017–1026 (1996).
13. H. Saitoh et al., *Cell* **70**, 115–125 (1992).
14. J. Tomkiel et al., *J. Cell Biol.* **125**, 531–545 (1994).
15. K.D. Brown, K.W. Wood, and D.W. Cleveland, *J. Cell Sci.* **109**, 961–969 (1996).
16. K.F. Sullivan, *Curr. Opin. Genet. Dev.* **11**, 182–188 (2001).
17. G.H. Karpen and R.C. Allshire, *Trends Genet.* **13**, 489–496 (1997).
18. A.E. Barry et al., *Hum. Mol. Genet.* **8**, 217–227 (1999).
19. L. Glowczewski et al., *Mol. Cell Biol.* **20**, 5700–5711 (2000).
20. A.V. Strunnikov, *Trends Cell Biol.* **8**, 454–459 (1998).
21. T. Tanaka, M.P. Cosma, K. Wirth, and K. Nasmyth, *Cell* **98**, 847–858 (1999).
22. J.M. Lopez, G.H. Karpen, and T.L. Orr-Weaver, *Curr. Biol.* **10**, 997–1000 (2000).

## Čerenkov Radiation

Čerenkov radiation, or the Čerenkov effect, is the visible bluish-white light surrounding a radioactive source observed in a pool of water. The phenomenon was first reported by the Russian physicist P. A. Čerenkov in 1934 (1), who was one of several scientists who observed the glow of light in close proximity to intense gamma-radiation sources. This effect is the reason why radiation sources are said to “glow in the dark.” Its importance for molecular biology is that it is routinely used to measure

the amount of the isotope  $^{32}\text{P}$  in a sample (see [Radioactivity](#) and [Phosphorous Isotopes](#)).

The Čerenkov effect can be explained by classic electromagnetic theory and principles of optic science. Analysis of Čerenkov radiation has shown fundamental relationships between the velocity of charged particles, light, the intensity of the light, and its wavelength spectrum (2). Only a small fraction (<0.1%) of the charged-particle radiation is emitted by the absorbing medium as coherent light. According to Jelley (3), who provided the first scientific explanation for the Čerenkov effect, the observed light results from charged particles traversing a transparent dielectric medium, one that does not conduct electricity. Charged particles are produced in the absorbing medium (water) when gamma radiation from the radioactive source interacts with the absorber. These charged particles produce local polarization along their path in the dielectric. Light in the visible spectrum is emitted when the polarized molecules in the medium return to their rest state soon after passage of the charged particle. If the velocity of the charged particles is less than that of light in the same medium, the light emitted from molecules in the dielectric is overridden and not observed. However, if the velocity of the charged particles is greater than that of light in the dielectric, a wavefront of light is produced from individual molecules in the dielectric, and the emission is reinforced by constructive interference. The Čerenkov effect is analogous to the bow wave from a ship that travels faster than the velocity of the surface waves, or to the shock wave trailing a supersonic aircraft passing through air. In other words, if the velocity  $v$  of a charged particle traversing a transparent dielectric material of refractive index  $n (= bc)$  exceeds the velocity of light ( $c/n$ ) in the medium, or  $v > c/n$  and  $b > 1/n$ , then Čerenkov radiation is emitted at an angle  $x$  relative to the particle direction, where  $x = \arccos(1/bc)$ . The velocity  $c$  of light in a vacuum is  $2.997 \times 10^{10}$  cm/s, and  $b = v/c$ . The Čerenkov photons form a conical wavefront of half-angle ( $90^\circ - x$ ) behind the particle.

The polarization effect actually decreases the energy lost by a charged particle that traverses a condensed medium, whereas the production of Čerenkov radiation increases the loss of energy by the particle. The visible spectrum is produced at a frequency interval of about  $3 \times 10^{14}$  Hz. Applications of Čerenkov radiation theory have been developed for detecting single high energy charged particles, measuring the energy of charged particles, and determining angles of incidence.

The light pulses emitted by charged particles traveling in a transparent medium can be collected and counted by modern scintillation counters. The amount of light produced in water is small compared to that produced in the presence of a scintillator, but it can be detected from beta-emitting radionuclides if their energy is greater than the threshold energy of 265 keV. The average energy of phosphorous-32 beta particles is 695 keV, so the majority of emitted beta particles can be detected.

#### Bibliography

1. P. A. Čerenkov (1934) *Compt. Rend. Acad. Sci. URSS* **2**, 451.
2. I. Frank and I. Tamm (1937) *Compt. Rend. Acad. Sci. URSS* **3**, 109.
3. J. V. Jelley (1953) *Atomics* **4**, 81.

#### Suggestions for Further Reading

4. F. H. Attix (1986) *Introduction to Radiological Physics and Radiation Dosimetry*, Wiley, New York.
5. W. J. Price (1958) *Nuclear Radiation Detection*, 2nd ed., McGraw-Hill, New York.

### Chain-Termination (Dideoxy) DNA Sequencing

The relative length of a DNA chain is determined easily and accurately by [gel electrophoresis](#), and the most common types of [DNA sequencing](#) use methods that map sequence information to DNA chain lengths. For chain-termination DNA sequencing, this is done by synthesizing the complement of the DNA using a **DNA polymerase** under conditions that terminate synthesis at sites where only one of the four bases occurs. Thus, all sequencing experiments are done in three steps. First, a single pure DNA segment is isolated for sequencing. Second, this DNA is used as a template for synthesis catalyzed by a **DNA polymerase** with mixtures of normal and chain-terminating nucleotides. Finally, the products of this synthesis are separated according to size by gel electrophoresis. Numerous variations on each of these steps are commonly used.

## 1. Isolation of Specific Segments of DNA

Most of the techniques of molecular biology rely on isolating specific segments of DNA, and sequencing is no exception. In fact, the first practical application of chain-termination sequencing relied on [cloning](#) specific DNA segments using vectors derived from **M13 bacteriophage**. These bacteriophages contain a single-stranded DNA **chromosome** that accommodates inserts of more than 5000 bases. Isolating the single-stranded DNA in pure form from these phages is simple and inexpensive, so these vectors are still commonly used. Similarly, virtually all plasmid vectors are commonly used for DNA sequencing, and essentially any clone of up to about 200 kb can be sequenced by cycle sequencing techniques, provided sufficient DNA can be purified. Another popular way of isolating DNA for sequencing is by using the polymerase chain reaction (**PCR**). With nested PCR primers, it is now relatively simple to amplify segments of **genomic** DNA for direct sequencing in a matter of hours.

## 2. Chain Termination Reactions

Chain termination reactions require the isolated template DNA, a suitable primer, 2'-deoxynucleoside triphosphates (dNTPs), 2',3'-dideoxynucleoside triphosphate (ddNTP) chain terminators, and a DNA polymerase. The most critical component is the DNA polymerase. Many DNA polymerases have been used for sequencing, including those from eubacteria, such as the DNA polymerase I (**Klenow fragment**) of *Escherichia coli*, **reverse transcriptases** from retroviruses, such as avian myeloblastosis virus, polymerases from bacteriophage, such as **T7 phage** DNA polymerase, and polymerases from archaea, such as *Thermococcus litoralis*. Virtually all of these have been genetically or chemically modified to eliminate exonuclease activities or to improve reaction rates with the ddNTPs.

The most recent examples of polymerases specifically engineered for DNA sequencing are enzymes derived from *Thermus aquaticus*, [Taq DNA polymerase](#). Two regions of this polymerase are modified to produce particularly effective polymerases for DNA sequencing. First, *Taq* DNA polymerase has 5'-3' exonuclease activity, which degrades sequencing primers. The first 300 amino acid residues at the N-terminus of this enzyme are required for this exonuclease activity. Portions of this **domain** can be deleted, or the activity can be eliminated by point mutation. Secondly, normally *Taq* DNA polymerase is relatively inefficient at using ddNTPs. As discovered by Tabor and Richardson, this can be improved more than 10<sup>4</sup>-fold by changing residue Phe667 to Tyr. This modification improves ddNTP usage, and it also greatly improves the quality of the sequence data obtained. Native *Taq* DNA polymerase produces sequencing bands that vary in intensity more than 15-fold, depending on the nearby sequence. In contrast, Tyr667 polymerase produces bands that vary in intensity by less than threefold. This makes interpreting the results of the electrophoretic separation much more accurate. A number of polymerases that have this modification are now commercially available, as is T7 DNA polymerase, which naturally has a tyrosine at the corresponding position.

### 3. Cycle Sequencing

Cycle sequencing is the process of using repeated cycles of thermal denaturation and polymerization to produce greater amounts of product in a DNA sequencing reaction. The amount of product DNA increases linearly with the number of cycles. (This distinguishes it from PCR, which uses two primers so that the amount of product increases exponentially with the number of cycles.) During each cycle, the thermostable DNA polymerase extends the annealed primer molecules, typically at 60° to 70°C. The mixture is heated above the melting temperature of DNA (95°C), dissociating the extended primer from the template. Then, the mixture is cooled, allowing another molecule of primer (which is present in excess) to anneal to the limited supply of template. Further cycles of extension and denaturation result in producing much more extended primer than the amount of template used. This improves the sensitivity of the sequencing experiment, and it also allows ready use of double-stranded templates for sequencing. Generally, cycle sequencing works much more reliably over a wider range of template concentrations than noncycled protocols. This accounts for its nearly universal application for large-scale DNA sequencing projects.

### 4. Methods for Labeling DNA Sequences

The products of the chain termination reactions must be labeled for all practical DNA sequencing methods. The original label was  $\alpha$ -<sup>32</sup>P dATP that was simply added to the chain-termination reaction. Newly synthesized DNA was labeled with radioactive phosphorous, and detected by simple [autoradiography](#). More recently, the lower energy isotopes <sup>33</sup>P and <sup>35</sup>S (in the form of  $\alpha$ -thio-dATP) have been used because they generate autoradiograms with higher resolution. These offer the advantage of using less total radioactivity than other methods. In addition, only specifically terminated, elongated DNA chains are labeled and therefore visualized by autoradiography, which eliminates the background bands and stops normally observed on DNA sequencing autoradiograms and results in extremely clean sequence data.

Automated, **fluorescent** DNA sequencing methods were introduced in 1987 and have become essential tools for large-scale sequencing efforts. The sequence products used by these automated systems are labeled by fluorescent primers (dye primers) or fluorescent dideoxynucleotides (dye terminators). These have been used in single-color detection instruments, and in four-color multiplex instruments in which the four bases are distinguished by color. Recent innovations include fluorescent dye-labeled DNA primers that exploit fluorescent **energy transfer** to optimize the absorption and emission properties of the label. These primers carry a fluorescein derivative at the 5'-end as a common donor and rhodamine derivatives attached to a modified thymidine within the primer sequence as acceptors. Adjustment of the donor–acceptor spacing by placing the modified thymidine in the primer sequence allows generating four primers. All have strong absorption at a common excitation wavelength (488 nm) and efficient fluorescent emission at 525, 555, 580, and 605 nm. These improve the sensitivity and accuracy of the automated sequencing system.

Fluorescent dye-labeled ddNTP terminators have also been used extensively for DNA sequencing, and those that use the energy-transfer principle are also commercially available. Like the radio-labeled terminators, they have the advantage of labeling only specifically terminated, elongated DNA chains, so that background bands are eliminated.

### 5. Electrophoresis and Automated Sequencing

The high-resolution separation of DNA fragments by size is essential for all sequencing methods. For radioactively labeled DNA sequencing experiments, this is done by using gels cast in glass plates that are 0.2 to 0.4 mm thick, 40 to 80 cm long, and wide enough to accommodate 32 to 96 samples. Typically, the gels are 4 to 8% polyacrylamide cross-linked with *N,N'*-methylene bisacrylamide (see [Polyacrylamide](#)) and contain tris borate buffer (0.089 M, pH 8.3) and 7 to 8 M urea. After electrophoresis for 2 to 18 hours, the gels must be removed from the glass plates for autoradiography

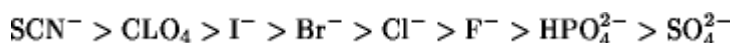
and reading of the sequence of 200 to 400 nucleotides. Because these gels are cumbersome to make and use, considerable effort has been made to improve separation methods. The most commonly used methods involve a sensitive fluorescent detection instrument that continuously monitors the migration of fluorescent-labeled DNA past a fixed position on the gel. The results are collected and evaluated directly by computer, producing finished sequence information. This saves considerable labor in “reading” the sequence from the gels and improves the resolution sufficiently to read 500 or more bases routinely from a single sequence experiment. Noncross-linked “gels” have also been introduced that run in 50 to 100 micron diameter, 40 to 70 cm long capillaries with fluorescent detection. The efficient heat transfer of these electrophoresis media allow faster, high-resolution separations.

#### Suggestions for Further Reading

- C. W. Fuller and M. A. Reeve (1995) Thermo sequenase - A novel thermostable polymerase for DNA sequencing, *Nature* **376**, 796–797.
- J. Ju et al. (1995) Design and synthesis of fluorescence energy transfer dye-labeled primers and their application for DNA sequencing and analysis, *Anal. Biochem.* **231**, 131–140.
- L. G. Lee et al. (1992) DNA sequencing with dye-labeled terminators and T7 DNA polymerase: Effect of dyes and dNTPs on incorporation of dye-terminators, and probability analysis of termination fragments, *Nucleic Acids Res.* **20**, 2471–2483.
- J. Messing (1983) Bacteriophage M13 vectors for DNA sequencing, *Methods Enzymol.* **101**, 33–38.
- G. M. Prober et al. (1987) system for rapid DNA sequencing with fluorescent chain-Terminating Dideoxynucleotides, *Science* **238**, 336–341.
- B. B. Rosenblum et al. (1997) New dye-labeled terminators for improved DNA sequencing patterns, *Nucleic Acids Res* **25**(22), 4500–4504.
- F. Sanger, S. Nicklen, and A. R. Coulson (1977) DNA sequencing with chain terminating inhibitors, *Proc. Natl. Acad. Sci. USA* **74**(12), 5463–5467.
- L. M. Smith et al. (1986) Fluorescence detection in automated DNA sequence analysis, *Nature* **321** (6071), 674–679.
- S. Tabor and C. C. Richardson (1995) A single hydroxyl moiety in Pol I-T DNA polymerase is responsible for distinguishing between deoxy- and dideoxyribonucleotides, *Proc. Natl. Acad. Sci. USA* **92**, 6339–6343.
- P. B. Vander Horn et al. (1997) Thermo sequenase DNA polymerase and *T. acidophilum* pyrophosphatase: New thermostable enzymes for DNA sequencing, *BioTechniques* **22**, 758–765.

#### Chaotropes: Kosmotropes

Compounds that increase and decrease, respectively, the aqueous solubility of proteins are classified as chaotropes and kosmotropes ([1](#)). Therefore, the [Hofmeister series](#) of ions may be divided into these two categories. For the series of anions, solubilization in water is promoted in the order



Those on the left-hand side of  $\text{Cl}^-$  generally increase the solubility; they are called *chaotropes*.

Those on the right-hand side of  $\text{Cl}^-$  act as salting-out agents; they are called *kosmotropes*. Chaotropes are **water**-structure breakers; kosmotropes, to the contrary, are water-structure makers. A similar division may be made about the  $\text{Na}^+$  ion in the cationic Hofmeister series. In the salts, the predominant effect is that of anions. An important characteristic is that the effects are additive for all the species in the solution.

Organic molecules also affect water structure and can be classified as chaotropes and kosmotropes. Thus, [urea](#), glycine, formamide, and acetamide are chaotropes. Organic kosmotropes are sucrose, polyols, the methylacetamines, methyl formamides, and methyl ureas. Except for sugar and the polyols, caution must be exercised in using organic kosmotropes, since their action may reverse itself at high concentration.

The mechanism of action of these molecules in stabilizing or destabilizing protein structure, and as a corollary in acting as **salting in** or **salting out** agents, is related to their effects on the orientation of water molecules through polar interactions with the water [hydrogen bond](#) donor and acceptor properties. By perturbing or strengthening the interactions between water molecules (maximization of hydrogen bond formation), these molecules exercise their chaotropic or kosmotropic actions. At interfaces, the water structure is already perturbed by the nonavailability at the surface of other water molecules for hydrogen bond formation. This is true for the water–air interface. Similar perturbations may occur at interfaces with proteins and other biological entities—hence, the effect on their properties. Control of reactions by the chaotropic/kosmotropic effect is very widespread in biological systems. Its principal diagnosis is that the compound which affects a particular process must be present at a high concentration (local concentration for cell compartments and organelles).

#### Bibliography

1. K. D. Collins and M. W. Washabaugh (1985) *Quart. Rev. Biophys.* **18**, 323–422.

## Chaperonin

The chaperonins are a family of **molecular chaperone** found in all types of cell whose function is to assist the correct folding of certain polypeptide chains that have been either newly synthesized or generated from native proteins by environmental stresses that cause partial unfolding.

### 1. Nomenclature

The term *chaperonin* was suggested by Sean Hemmingsen (1) to describe a family of highly sequence-related molecular chaperones found in [chloroplasts](#), [mitochondria](#), and **eubacteria** such as *Escherichia coli*. This term was proposed to simplify the existing complex nomenclature for different members of this family whose close relationship was only realized when cDNA for the chloroplast chaperonin was sequenced (1). The term was subsequently extended to distant homologues found in **archaeobacteria** and the eukaryotic cytosol (2, 3). The two subfamilies of the chaperonin family are referred to as follows:

1. The GroE or group I subfamily, found in chloroplasts and other plastids, mitochondria, and all eubacteria
2. The TCP-1 or group II subfamily, found in archaeobacteria and the eukaryotic cytosol

The GroE terminology reflects the fact that the eubacterial protein was first identified by genetic studies in four laboratories in 1972/73 as a bacterial protein required for the replication of **bacteriophage** such as lambda in *E. coli* (4). “Gro” refers to phage growth, and the suffix “E” refers

to the observation that the phage growth defect is overcome when the phage carries a mutation in the head gene E. The TCP-1 terminology is derived from the identification of a protein called the *t*-complex polypeptide encoded by the mouse T locus (5). The mitochondrial members of the GroE subfamily are sometimes referred to as “hsp60” proteins, but this terminology should not be used to describe the chaperonins as a whole, since the eukaryotic TCP-1 members and the chloroplast GroE members are not **heat shock** proteins.

Table 1 presents a suggested nomenclature and useful abbreviations for members in both subfamilies, and lists other names that are used. Some authors confine the abbreviation “cpn60” to the GroE subfamily, but this restriction is inconsistent with the fact that the TCP-1 members share with the GroE members subunit molecular masses of approximately 60 kDa.

**Table 1. Nomenclature of the Chaperonins**

| Preferred Name                 | Other Names  | Useful Abbreviations        |
|--------------------------------|--|-----------------------------|
| GroE subfamily                 |  |                             |
| Eubacterial chaperonin 60      | GroEL ( <i>E. coli</i> ), 65-kDa antigen   | Eu cpn60                    |
| Eubacterial chaperonin 10      | GroES ( <i>E. coli</i> ), cochaperonin   | Eu cpn10                    |
| Mitochondrial chaperonin 60    | hsp60, mitonin, HuCha60  | Mt cpn60                    |
| Mitochondrial chaperonin 10    | hsp10, cochaperonin  | Mt cpn10                    |
| Chloroplast chaperonin 60      | Rubisco subunit binding protein  | Ch cpn60                    |
| Chloroplast chaperonin 10      | chaperonin 21, cochaperonin  | Ch cpn10, Ch cpn21          |
| TCP-1 subfamily                |  |                             |
| Cytosolic chaperonin 60        | <i>t</i> -Complex polypeptide 1, chaperonin-containing TCP-1, TCP-1 ring complex | Cyt cpn60, TCP-1, CCT, TRiC |
| Archaeobacterial chaperonin 60 | Thermophilic factor 55, thermosome   | Ar cpn60, TF55              |

## 2. Function

The main function of the chaperonins is to assist the folding of certain newly synthesized polypeptide chains into their biologically active conformations (see [Protein Folding In Vivo](#)). They achieve this end, not by providing steric information required for each chain to fold correctly, but by sequestering each partially-folded chain in a protected compartment generated by the oligomeric structure of the chaperonin molecule. In this compartment, each chain can continue to fold to a point where aggregation with other partially folded chains is no longer a problem (6-9). This aggregation arises because some proteins fold via intermediate states that expose **hydrophobic** surfaces

transiently and are thus subject to intermolecular interactions with similar states. This aggregation effect can sometimes be observed when chemically denatured proteins refold spontaneously in dilute solution *in vitro* (see [Protein Folding In Vitro](#)), but its magnitude is expected to be enhanced in the intact cell because of the phenomenon of macromolecular crowding, or [excluded volume](#), thus increasing the effective concentrations of these states by two or three orders of magnitude ([8, 9](#)). It should be stressed that, although the chaperonins increase the yield of correctly folded proteins by minimizing aggregation, they do not increase the rate of folding above that achieved by the fastest folding fraction of spontaneously refolding chains that manage to avoid aggregation.

The best evidence for this view of chaperonin function comes from genetic and **pulse-chase** labeling studies with living cells. Several proteins imported into the mitochondria of yeast cells form aggregates when the function of the mitochondrial chaperonin is impaired by mutation ([10](#)), whereas about 30% of newly synthesized soluble cytoplasmic proteins in *E. coli* cells become either insoluble or inactive when *cpn60* function is switched off by means of a [temperature-sensitive mutation](#) ([11](#)). *In vitro* studies of the ameliorating effect of added GroE proteins on the aggregation of various pure denatured proteins during their refolding after removal of denaturant are consistent with this antiaggregation role ([12-14](#)). *In vitro* studies also suggest that a side effect of this antiaggregation role is the partial unfolding of nonaggregated chains that have become kinetically trapped in misfolded conformations, thus allowing these chains another chance to fold correctly ([15](#)); it is not clear how important such an effect may be *in vivo*. In addition, the GroE chaperonins in eubacteria and mitochondria ([16](#)) and the TCP-1 chaperonins in the archaeobacteria ([17](#)) are stress proteins that prevent and reverse the denaturation of some fully folded proteins by stresses such as high temperature.

It must be emphasized that the actual spectrum of proteins that use chaperonin function to fold correctly in intact cells is not known, but preliminary calculations and observations suggest that it is probably a minority in the case of both the GroE chaperonins ([9, 18](#)) and the TCP-1 chaperonins ([7](#)). Genetic evidence derived from the study of yeast mutants suggests that a major function of the eukaryotic TCP-1 is to assist the folding of [tubulin](#) and [actin](#) ([19](#)), consistent with the observation that newly synthesized chains of actin and tubulin can be isolated from pulse-labeled cells of CHO cells in the form of complexes with TCP-1 chaperonin oligomers ([20](#)). *In vitro* protein refolding experiments also suggest that the substrate specificities of the GroE chaperonins and the TCP-1 chaperonins are different ([6](#)).

### 3. Structure

The chaperonins occur as large oligomeric structures consisting of two stacked rings of subunits, each about 55–60 kDa in size, surrounding a central cavity or cage in which the protein substrate binds; each subunit catalyzes the slow hydrolysis of ATP to ADP. The GroE chaperonins have seven subunits per ring, whereas the TCP-1 chaperonins have eight or nine. The subunit sequence identity between members of the GroE subfamily is in the range 42–76%, whereas that between members of the TCP-1 subfamily is around 30–40%, but the identity between the subfamilies is much less (15–20%) and is confined to regions corresponding to the [ATPase](#) domain in GroE *cpn60*. The evolutionary implications of these similarities are under debate ([7, 21](#)).

A convenient abbreviation for the chaperonin(s) is *cpn(s)*. Thus the subunits of the ring can be called *cpn60 subunits*, and the oligomer can be called *cpn60*. The GroE subfamily also contains another type of oligomeric ring made of seven smaller 10-kDa subunits, called *cpn10*, because there is a slight similarity between the *cpn10* subunit sequence and part of the *cpn60* subunit sequence ([21](#)). GroE *cpn60* and *cpn10* oligomers bind to each other in 1:1 and 1:2 ratios in the presence of either ATP or ADP, and these binary complexes play an essential role in the protein folding function of these molecules ([8, 9](#)). The chloroplast *cpn10* oligomer is unusual in that each subunit consists of two copies of a *cpn10* sequence fused head to tail and thus is often referred to as *cpn21* ([22](#)).

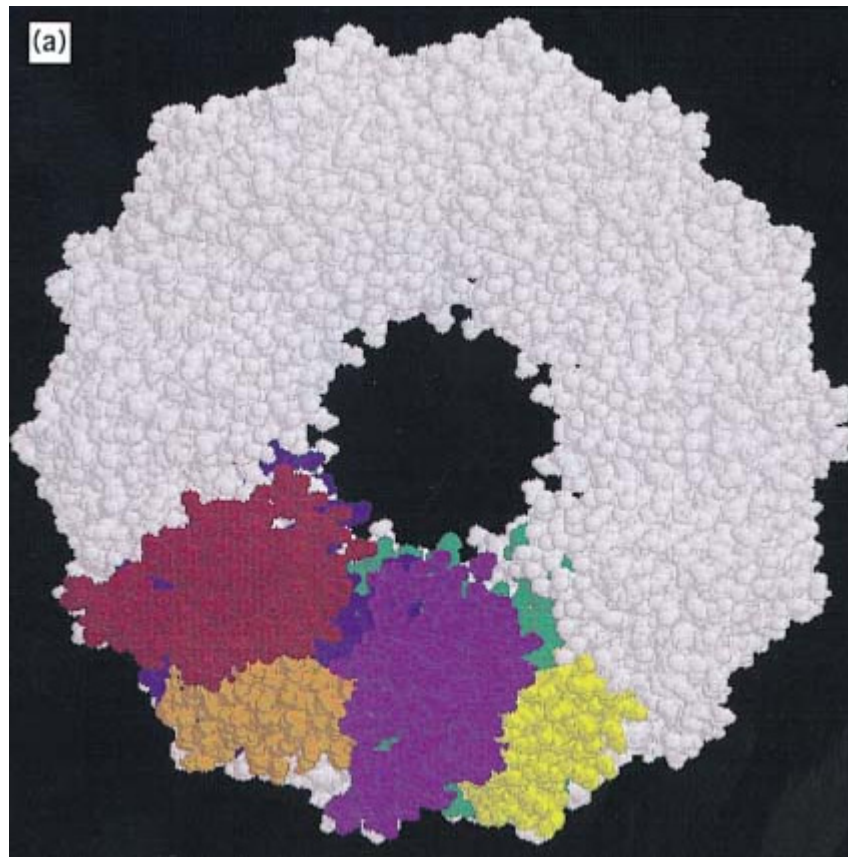
The TCP-1 oligomers found in the eukaryotic cytosol are much more variable than those found in



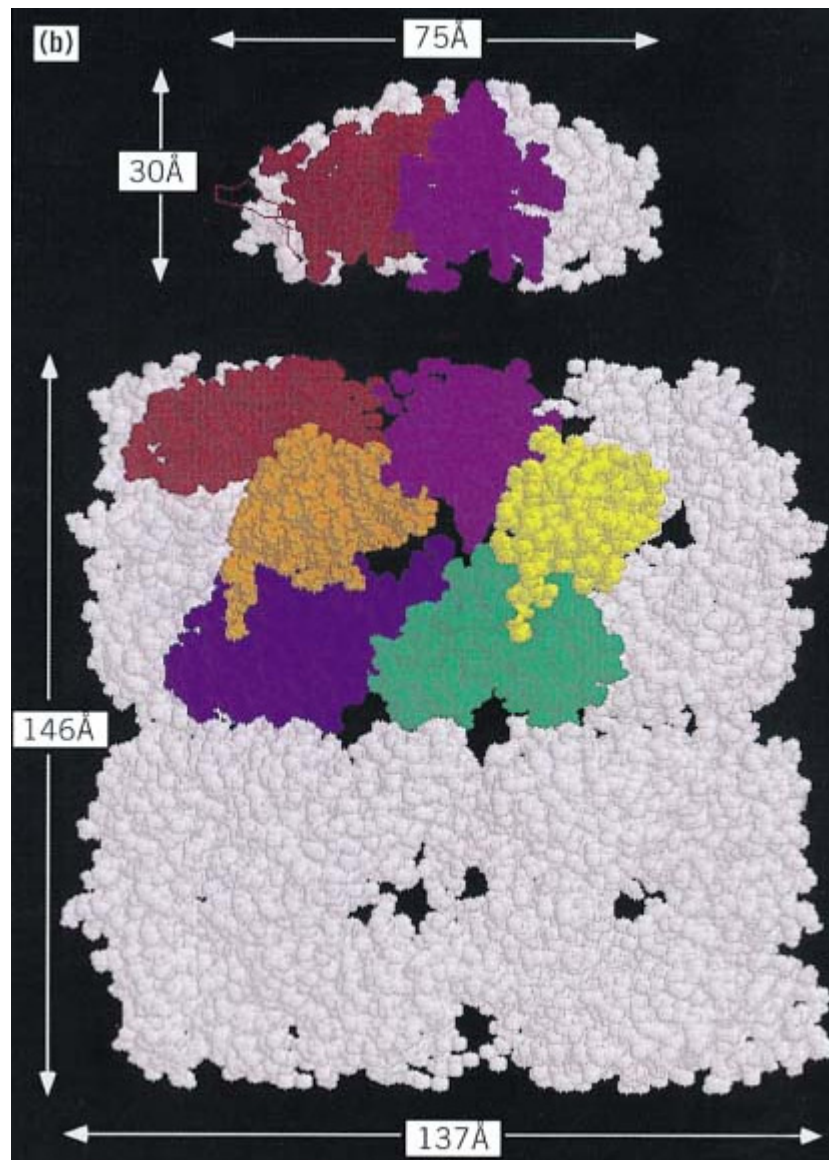
either the archaeobacteria or the GroE subfamily, since subunits of the latter consist of only one or two sequences, whereas the former have at least eight different, but related, subunit sequences in each ring (6, 7, 23). For this reason, the eukaryotic TCP-1 oligomers are also called *CCT complexes*, for chaperonin-containing TCP-1 complexes, or *TriC* for TCP-1 ring complexes (see Table 1). The TCP-1 chaperonins do not appear to contain cpn10-like members.

The best-studied chaperonins are the cpn60 and cpn10 oligomers from *E. coli*, termed GroEL and GroES, respectively. GroEL has been extensively studied by both [negative stain](#) and [cryoelectron microscopy](#) (24, 25) and by [X-ray crystallography](#) of a double mutant form (26). GroEL is a cylindrical structure containing a central cavity about 50 Å in diameter (Fig. 1a). Each subunit consists of (1) an apical domain to which the protein substrate and GroES bind; (2) an equatorial domain that protrudes into the central cavity, which contains the ATPase site and is responsible for most of the intersubunit interactions; and (3) an intermediate domain that contains potential hinge sites responsible for the considerable movements of the other two domains revealed by electron microscopy (24, 25). These conformational changes result from the binding of nucleotide and GroES to the GroEL (26) cylinder. The crystal structure of GroES shows a dome-shaped structure with a hydrophilic inner surface that fits over one end of the GroEL cylinder (24, 27, 28) (see also Fig. 1b,c), creating a large enclosed dome-shaped space about 65 × 80 Å in size, inside which the protein substrate continues to fold (29-32) (see also Fig. 1c,d).

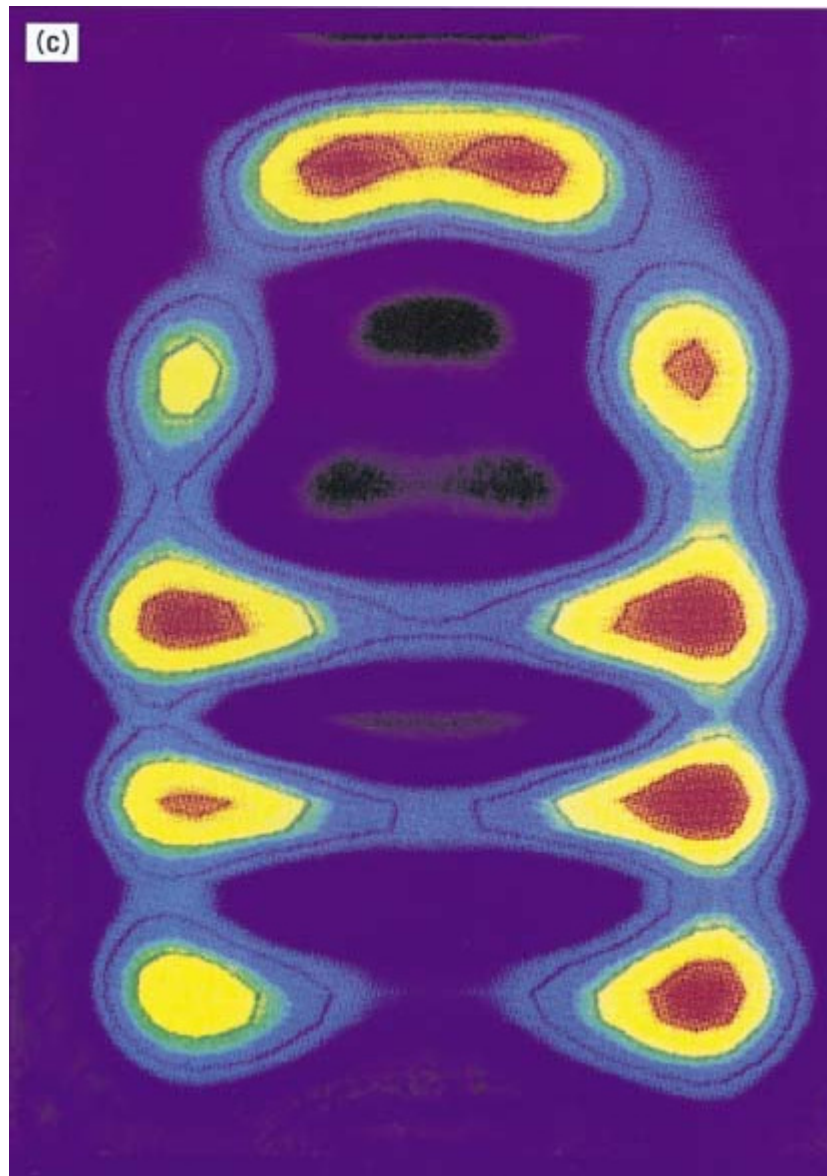
**Figure 1.** Structure and function of the chaperonin system. (a,b) Space-filling representations showing a top and side view, respectively, of the crystal structure of a double mutant form of GroEL (26). Two adjacent subunits are colored with the apical domains in red and purple, the intermediate domains in orange and yellow, and the equatorial domains in blue and green, respectively. Free passage between the GroEL rings is obstructed by *N*- and *C*-terminal residues not resolved in the crystal structure. (b) side view of GroES (28) at the top; two adjacent domains and a single mobile loop that is structured in the GroES crystal due to crystal packing are colored. The loop regions normally protrude from the base of GroES downward toward the GroEL. (c) Asymmetrical GroEL/GroES complexes as revealed by cryoelectron microscopy (45). Note the upward and outward movement of the apical GroEL domains interacting with GroES. (d) Model of the GroEL/GroES reaction cycle in assisting protein folding (30); see text for details. The term *unfolded protein* refers to a partially folded intermediate that is represented by the light pink spheres; folded protein is represented by the dark pink sphere, while the hatched spheres represent a mixture of folded and partially- folded proteins expected in a population of GroEL molecules. At step 4, GroES may associate with either the protein-containing ring of GroEL or with the opposite empty ring; the latter possibility is not shown. Reprinted from *Nature* 381, 571–580 (1996), with permission; copyright (1996) Macmillan Magazines Ltd. See color insert.



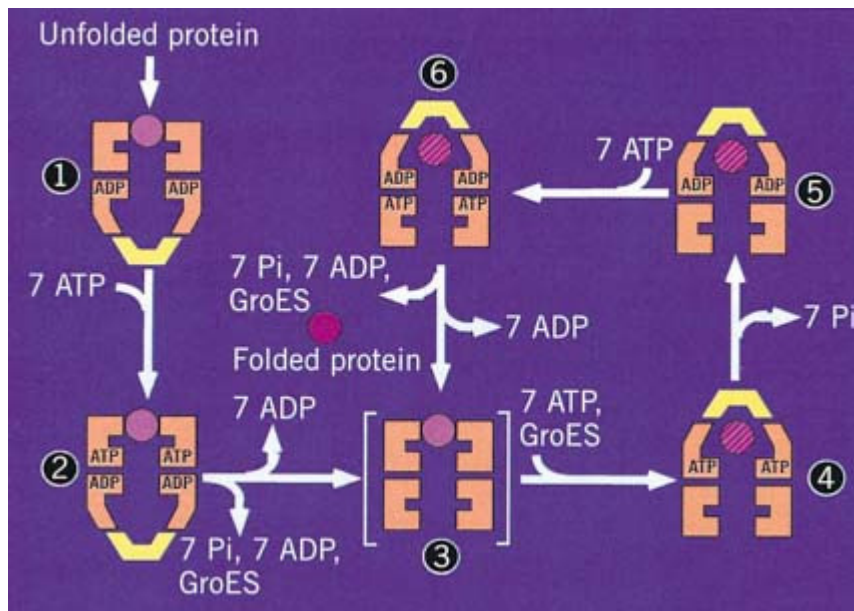
**Figure 1.** (*continued*) Structure and function of the chaperonin system. **(a,b)** Space-filling representations showing a top and side view, respectively, of the crystal structure of a double mutant form of GroEL (26). Two adjacent subunits are colored with the apical domains in red and purple, the intermediate domains in orange and yellow, and the equatorial domains in blue and green, respectively. Free passage between the GroEL rings is obstructed by *N*- and *C*-terminal residues not resolved in the crystal structure. **(b)** side view of GroES (28) at the top; two adjacent domains and a single mobile loop that is structured in the GroES crystal due to crystal packing are colored. The loop regions normally protrude from the base of GroES downward toward the GroEL. **(c)** Asymmetrical GroEL/GroES complexes as revealed by cryoelectron microscopy (45). Note the upward and outward movement of the apical GroEL domains interacting with GroES. **(d)** Model of the GroEL/GroES reaction cycle in assisting protein folding (30); see text for details. The term *unfolded protein* refers to a partially folded intermediate that is represented by the light pink spheres; folded protein is represented by the dark pink sphere, while the hatched spheres represent a mixture of folded and partially- folded proteins expected in a population of GroEL molecules. At step 4, GroES may associate with either the protein-containing ring of GroEL or with the opposite empty ring; the latter possibility is not shown. Reprinted from *Nature* 381, 571–580 (1996), with permission; copyright (1996) Macmillan Magazines Ltd. See color insert.



**Figure 1.** (*continued*) Structure and function of the chaperonin system. (**a,b**) Space-filling representations showing a top and side view, respectively, of the crystal structure of a double mutant form of GroEL (26). Two adjacent subunits are colored with the apical domains in red and purple, the intermediate domains in orange and yellow, and the equatorial domains in blue and green, respectively. Free passage between the GroEL rings is obstructed by *N*- and *C*-terminal residues not resolved in the crystal structure. (**b**) side view of GroES (28) at the top; two adjacent domains and a single mobile loop that is structured in the GroES crystal due to crystal packing are colored. The loop regions normally protrude from the base of GroES downward toward the GroEL. (**c**) Asymmetrical GroEL/GroES complexes as revealed by cryoelectron microscopy (45). Note the upward and outward movement of the apical GroEL domains interacting with GroES. (**d**) Model of the GroEL/GroES reaction cycle in assisting protein folding (30); see text for details. The term *unfolded protein* refers to a partially folded intermediate that is represented by the light pink spheres; folded protein is represented by the dark pink sphere, while the hatched spheres represent a mixture of folded and partially- folded proteins expected in a population of GroEL molecules. At step 4, GroES may associate with either the protein-containing ring of GroEL or with the opposite empty ring; the latter possibility is not shown. Reprinted from *Nature* 381, 571–580 (1996), with permission; copyright (1996) Macmillan Magazines Ltd. See color insert.



**Figure 1.** (*continued*) Structure and function of the chaperonin system. **(a,b)** Space-filling representations showing a top and side view, respectively, of the crystal structure of a double mutant form of GroEL (26). Two adjacent subunits are colored with the apical domains in red and purple, the intermediate domains in orange and yellow, and the equatorial domains in blue and green, respectively. Free passage between the GroEL rings is obstructed by *N*- and *C*-terminal residues not resolved in the crystal structure. **(b)** side view of GroES (28) at the top; two adjacent domains and a single mobile loop that is structured in the GroES crystal due to crystal packing are colored. The loop regions normally protrude from the base of GroES downward toward the GroEL. **(c)** Asymmetrical GroEL/GroES complexes as revealed by cryoelectron microscopy (45). Note the upward and outward movement of the apical GroEL domains interacting with GroES. **(d)** Model of the GroEL/GroES reaction cycle in assisting protein folding (30); see text for details. The term *unfolded protein* refers to a partially folded intermediate that is represented by the light pink spheres; folded protein is represented by the dark pink sphere, while the hatched spheres represent a mixture of folded and partially- folded proteins expected in a population of GroEL molecules. At step 4, GroES may associate with either the protein-containing ring of GroEL or with the opposite empty ring; the latter possibility is not shown. Reprinted from *Nature* 381, 571–580 (1996), with permission; copyright (1996) Macmillan Magazines Ltd. See color insert.



No crystal structure for any TCP-1 chaperonin is yet available, but electron-microscopic studies suggest that the general architecture and overall domain structure of the TCP-1 cylinder resembles that of GroEL (33, 34). The archaeobacterial chaperonin was originally termed thermophilic factor 55, because it is virtually the only protein made by thermophilic archaeobacteria under heat-shock conditions (35), but it is now termed the *thermosome* (36). Thermosomes generally contain eight subunits per ring, but there are several reports that a significant subset of particles contain nine subunits per ring, while the chaperonins of some *Sulfolobus* species appear all to contain nine subunits per ring (37). The thermosomes of *Pyrodictium*, *Thermoplasma*, and *Sulfolobus* consist of two different but related subunits of almost identical molecular mass (33). Electron-microscopic image analysis suggests that the two types of subunit in the thermosome of *Thermoplasma acidophilum* alternate within each ring, generating a four-fold symmetry (38).

Averaged electron-microscopic images of the mammalian TCP-1 particle reveal a ringlike structure with eightfold quasi-rotational symmetry with a central channel (33). Primary structures of eight distinct types of related subunit have been determined by cloning and sequencing of mouse cDNA; these types share about 30% amino acid residue identity. A striking observation is that the sequence of each type is more highly conserved among different eukaryotic species than are the sequences between types in any one species; in mammals the sequence identity is over 96% for each type of subunit, and around 60% between yeast and mouse. These observations have prompted the conclusion that each type of subunit diverged early in evolution, has changed only slowly during the evolution of eukaryotes, and thus may have a specific function in binding subsets of sequences in substrate proteins and/or other molecular chaperones (6, 7). Consistent with this suggestion is the observation that the different types of subunit in each TCP-1 chaperonin ring deviate most from one another in the sequences of the apical domains, where the protein substrate binds.

#### 4. Mechanism of Action

The molecular details of how the chaperonins function have, and continue to be, studied by means of *in vitro* protein refolding experiments using defined components, but such experiments suffer from the disadvantage that the conditions under which they are performed are different in important respects from those operating during protein folding in the intact cell (9, 18). Thus it is not surprising that this field is subject to intense debate between people supporting different views (39). Since the chaperonins evolved to fold proteins inside cells, and not inside test tubes, it is essential when evaluating the validity of these views to define those aspects of chaperonin action observed *in vitro*

that are mechanistically essential to improve the yield of correctly folded protein, when extrapolated as far as possible to *in vivo* conditions. In our opinion, the model outlined in Figure 1d best meets this criterion for the GroE chaperonins (30, 40).

In this model, newly synthesized polypeptide chains are released from ribosomes in the form of partially folded intermediates that may contain bound molecular chaperones of the hsp70/DnaJ families (see Protein Folding In Vivo). The latter chaperones prevent aggregation and premature folding before each chain is complete. Such intermediates are commonly referred to as *unfolded proteins*, but in fact they often are molten globules. These partially folded, compact intermediates lose the hsp70/DnaJ chaperones by an uncharacterized mechanism that results in each intermediate molecule binding by hydrophobic interactions to the apical domains at one end of a GroEL cylinder; the other end of the cylinder is occupied by GroES as a result of ADP binding to the GroEL ring proximal to the GroES (see Fig. 1d, step 1). The binding of polypeptide and the hydrolysis of ATP by the ring not occupied by GroES causes the GroES and the ADP to dissociate (step 2). GroES and ATP then rebind with equal probability to either ring (not both rings), resulting in 50% of the bound polypeptide being encapsulated inside the cavity capped by GroES (step 4). This rebinding of GroES results in the displacement of the bound polypeptide into the cavity, where it can continue to fold; it is for this reason that this model has been dubbed the "Anfinsen cage" model to stress the idea that the protein in the cavity folds in a similar, but not necessarily identical, manner to that by which it folds in a protein renaturation experiment of the type pioneered by Anfinsen (41, 42). Thus one essential difference between chaperonin-assisted protein folding *in vivo* and spontaneous protein refolding *in vitro* is that the former process segregates each chain into the protected compartment provided by the GroEL/GroES cylinder to avoid the problem of aggregation.

The time available for this protected folding is set by the time it takes each ring of GroEL to hydrolyze seven ATP molecules. This hydrolysis is positively cooperative in each ring, and two rounds take about 10 s at 37°C (steps 4 and 5). The GroES then dissociates, and the protein is free either to diffuse into the cytosol or to rebind to the apical domains of the same GroEL cylinder (step 6). Rebinding occurs only if sufficient hydrophobic residues are still exposed on the surface of the compact intermediate. The cycle of release into the cage then repeats until the protein is folded sufficiently to the point where aggregation with other partially folded proteins is no longer a problem for the cell. The number of reaction cycles will vary, depending on the rate at which a given protein can fold when released into the cavity. Even with a single type of protein substrate, the population of folding chains will not be in synchrony or occupy the same number of cycles, since some chains fold more rapidly than others of the same sequence (see Protein Folding In Vitro). Usually only a fraction of the polypeptide chains fold to completion in a single cycle of GroES binding and unbinding, depending on the intrinsic rate of folding of a particular polypeptide and its tendency to become trapped in misfolded states. Rebinding of these intermediates to the apical domains of GroEL results in partial unfolding in preparation for a subsequent folding trial (30). This unfolding function distinguishes the GroEL/ES machinery from a simple folding cage, which acts solely to prevent intermolecular aggregation.

Uncertainties about this model include the fate of the 50% of polypeptide bound to the GroEL ring not capped by GroES; such polypeptides would be prone to aggregation and degradation if released into the cytosol. Release of such polypeptides *in vitro* is reduced in the presence of high concentrations of synthetic polymers that mimic the macromolecular crowding, excluded volume effect present in the cytosol (40). However, some release of partially folded chains may be essential *in vivo* because of the danger that proteins that are unable to fold correctly as a result of mutation will clog up the GroEL apical domains; an escape route for such molecules may be provided by release from the ring not capped by GroES.

Another problem concerns the significance of GroEL molecules with GroES bound to each end of the cylinder that have been observed *in vitro* (43). The ratio of total GroES subunits to total GroEL subunits in extracts of *E. coli* cells is equimolar, consistent with such GroES/GroEL/GroES complexes occurring *in vivo*, where their formation may be favored by the effect of macromolecular

crowding on protein association. Nevertheless, such complexes are not mechanistically essential, at least as deduced from *in vitro* refolding experiments, nor can they be permanent if the compact intermediates are to enter the GroEL cavities.

Another limitation imposed by this model is the size range that can be accommodated inside the cage, which is likely to have an upper limit around 60 kDa. Larger proteins than this, up to about 150 kDa, exist inside bacterial cells and even larger ones in the cytosol of eukaryotic cells. A GroES equivalent to cap the TCP-1 chaperonin oligomer has not been identified, so it is possible that the latter assists protein folding in a manner somewhat different from that for the GroE chaperonins. There is evidence that, unlike GroEL, the mammalian TCP-1 chaperonin binds to chains of actin being synthesized by extracts of reticulocytes before the chains have been released from the ribosomes (44).

There are four known intracellular compartments in which proteins fold that do not appear to contain any type of chaperonin: the [endoplasmic reticulum](#) lumen, the intermembrane **mitochondrial** space, the intrathylakoidal lumen, and the **periplasmic** space of **Gram-negative** bacteria. Additionally, the eukaryotic cytosol seems to lack a general chaperonin, since the TCP-1 chaperonin appears to be restricted to a small subset of protein substrates, including actin and tubulin. Thus it is certainly not the case that the folding of all proteins in all cells involves the chaperonins. Given that the intracellular environment strongly favors aggregation, it is likely that the chaperonin-independent proteins utilize the assistance of other types of molecular chaperone for their folding. This assistance may also apply to proteins larger than 60 kDa that cannot fit into the chaperonin cage. For these proteins, the cotranslational and sequential folding of protein **domains** may circumvent the requirement for a sequestered folding compartment.

It is also possible that some chaperonins have roles additional to those discussed above with respect to protein folding, since there are sporadic reports of chaperonin-like molecules occurring on the cell surface in both prokaryotic and eukaryotic cells and in the blood serum, as well as stimulating **cytokine** production by animal cell cultures. The bacterial chaperonins are also the major immunogens in all human bacterial infections, and their possible roles in protective immunity and autoimmune disease is an active area of research (45).

## 5. Note Added in Proof

The cage mechanism for chaperonin-assisted protein folding is strongly supported by recent structural and functional data. The crystal structure of the asymmetrical GroEL:GroES complex with bound ADP shows that the cage is sufficient to accept proteins up to at least 70 kDa and that the wall of the cage provides a hydrophilic environment permissive for folding (Z. Xu et al. *Nature* **388**, 741–750, 1997). Crystal structures of the archaean thermosome suggest that in this chaperonin the folding cage is closed by loop sequences that emanate from the apical domains of the chaperonin subunits, explaining the lack of a separate GroES-like factor (M. Klumpff et al. *Cell* **91**, 263–270, 1997; L. Ditzel et al. *Cell* **93**, 125–138, 1998). A single ring mitochondrial chaperonin is fully functional in protein folding in a GroES-dependent manner *in vivo* (K. L. Nielsen et al. *Mol. Cell* **2**, 93–99, 1998). The size range of GroEL substrate proteins and their kinetics of interaction *in vivo* are consistent with the cage mechanism (K. L. Ewalt et al. *Cell* **90**, 491–500, 1997). Oligomeric GroEL with an intact central cavity is essential for the maintenance of growth of *E. coli* (F. Weber et al. *Nature Struct. Biol.* **11**, 977–985, 1998).

## Bibliography

1. S. M. Hemmingsen, C. Woolford, S. M. van der Vies, K. Tilly, D. T. Dennis, G. C. Georgopoulos, R. W. Hendrix, and R. J. Ellis (1988) *Nature* **333**, 330–334.
2. R. S. Gupta (1990) *Biochem. Int.* **20**, 833–841.
3. R. J. Ellis (1990) *Science* **250**, 954–959.
4. C. Georgopoulos, D. Ang, K. Liberek, and M. Zylicz (1990) in R. I. Morimoto, A. Tissieres,

and C. Georgopoulos, eds., *Stress Proteins in Biology and Medicine*, Cold Spring Harbor Laboratory Press, New York, pp. 191–221.

5. L. Silver, K. Artzt, and D. Bennett (1979) *Cell* **17**, 275–284.
6. H. Kubota, G. Hynes, and K. Willison (1995) *Eur. J. Biochem.* **230**, 3–16.
7. K. R. Willison and A. L. Horwich (1996) in R. J. Ellis, ed., *The Chaperonins*, Academic Press, San Diego, pp. 108–136.
8. F. U. Hartl (1996) *Nature* **381**, 571–580.
9. R. J. Ellis and F. U. Hartl (1996) *FASEB J.* **10**, 20–26.
10. E. M. Hallberg, Y. Shu, and R. L. Hallberg (1993) *Mol. Cell Biol.* **13**, 3050–3057.
11. A. L. Horwich, K. B. Low, F. A. Fenton, I. N. Hirshfield, and K. Furtak (1993) *Cell* **74**, 909–917.
12. P. Goloubinoff, J. P. Christeller, A. A. Gatenby, and G. H. Lorimer (1989) *Nature* **342**, 884–889.
13. J. Buchner, M. Schmidt, M. Fuchs, R. Jaenicke, R. Rudolph, F. X. Schmid, and T. Kiefhaber (1991) *Biochemistry* **30**, 1586–1591.
14. J. Martin, T. Langer, R. Boteva, A. Schramel, A. L. Horwich, and F.-U. Hartl (1991) *Nature* **352**, 36–42.
15. G. S. Jackson, R. A. Staniforth, D. J. Halsall, T. Atkinson, J. J. Holbrook, A. R. Clarke, and S. G. Burston (1993) *Biochemistry* **32**, 2554–2563.
16. J. Martin, A. L. Horwich, and F. U. Hartl (1992) *Science* **258**, 995–998.
17. A. Guagliardi, L. Cerchia, S. Bartolucci, and M. Rossi (1994) *Protein Sci.* **3**, 1436–1443.
18. R. J. Ellis (1996) in R. J. Ellis, ed., *The Chaperonins*, Academic Press, San Diego, pp. 1–25.
19. D. Ursic, J. C. Sedbrook, K. L. Himmel, and M. R. Culbertson (1994) *Mol. Biol. Cell* **5**, 1065–1080.
20. H. Sternlicht, G. W. Farr, M. L. Sternlicht, J. K. Driscoll, K. R. Willison, and M. B. Yaffe (1993) *Proc. Natl. Acad. Sci. USA* **90**, 9422–9426.
21. R. S. Gupta (1996) in R. J. Ellis, ed., *The Chaperonins*, Academic Press, San Diego, pp. 27–64.
22. U. Bertsch, J. Soll, R. Seetharam, and P. V. Viitanen (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8696–8700.
23. K. F. Liou and K. R. Willison (1997) *EMBO J.* **16**, 101–106.
24. A. M. Roseman, S. Chen, H. White, K. Braig, and H. R. Saibil (1996) *Cell* **87**, 241–251.
25. O. Llorca, S. Marco, J. L. Carrascosa, and J. M. Valpuesta (1997) *J. Struct. Biol.* **118**, 31–42.
26. K. Braig, Z. Otwinowski, R. Hegde, D. C. Boisvert, A. Joachimiak, A. L. Horwich, and P. B. Sigler (1994) *Nature* **371**, 578–586.
27. S. C. Mande, V. Mehra, B. Bloom, and W. G. J. Hol (1996) *Science* **271**, 203–207.
28. J. F. Hunt, A. J. Weaver, S. J. Landry, L. Gierasch, and J. Dieneshofer (1996) *Nature* **379**, 37–42.
29. J. Martin, M. Mayhew, T. Langer, and F. U. Hartl (1993) *Nature* **366**, 228–233.
30. M. Mayhew, A. C. R. da Silva, J. Martin, H. Erdjument-Bromage, P. Tempst, and F. U. Hartl (1996) *Nature* **379**, 420–426.
31. J. S. Weissman, H. S. Rye, W. A. Fenton, J. M. Beecham, and A. L. Horwich (1996) *Cell* **84**, 481–490.
32. S. Chen, A. M. Roseman, A. S. Hunter, S. P. Wood, S. G. Burston, N. A. Ranson, A. R. Clarke, and H. R. Saibil (1994) *Nature* **371**, 261–264.
33. T. Waldmann, E. Nimmesgern, M. Nitsch, J. Peters, G. Pfeifer, S. Muller, J. Kellermann, A. Engel, F. U. Hartl, and W. Baumeister (1995) *Eur. J. Biochem.* **227**, 848–856.
34. V. A. Lewis, G. M. Hynes, D. Zheng, H. Saibil, and K. Willison (1992) *Nature* **358**, 249–252.



35. B. M. Phipps, A. Hoffmann, K. O. Stetter, and W. Baumeister (1991) *EMBO J.* **10**, 1711–1722.
36. B. M. Phipps, D. Typke, R. Hegerl, S. Volker, A. Hoffmann, K. O. Stetter, and W. Baumeister (1993) *Nature* **361**, 475–477.
37. S. Marco, D. Urena, J. L. Carrascosa, T. Waldmann, J. Peters, R. Hegerl, G. Pfeifer, H. Sack-Kongehl, and W. Baumeister (1994) *FEBS Lett.* **341**, 152–155.
38. M. Nitsch, M. Klumpp, A. Lupas, and W. Baumeister (1997) *J. Mol. Biol.* **267**, 142–149.
39. A. C. Clarke and P. A. Lund (1996) R. J. Ellis, ed., in *The Chaperonins* Academic Press, San Diego, pp. 168–212.
40. J. Martin and F. U. Hartl (1997) *Proc. Natl. Acad. Sci. USA* **94**, 1107–1112.
41. R. J. Ellis (1994) *Curr. Biol.* **4**, 633–635.
42. R. J. Ellis (1996) *Fold. Des.* **1**, R9–R15.
43. Z. Torok, L. Vigh, and P. Goloubinoff (1996) *J. Biol. Chem.* **271**, 16180–16186.
44. J. Frydman and F. U. Hartl (1996) *Science* **272**, 1497–1502.
45. A. R. M. Coates (1996) in R. J. Ellis, ed., *The Chaperonins* Academic Press, San Diego, pp. 268–296.

### Suggestions for Further Reading

46. R. J. Ellis, ed. (1996) *The Chaperonins*, Academic Press, San Diego (this is the first book devoted to the chaperonins and contains 10 chapters by different authors covering, as well as the main topics discussed above, other aspects such as gene regulation, gene accession numbers, primary structures, and immunological and medical implications).
47. M.-J. Gething, ed. (1997), *Molecular Chaperones and Protein Folding Catalysts*, Oxford University Press, Oxford (this book contains detailed information about the chaperonins, as well as other molecular chaperones).

## Chemical Modification

Chemical modification is one of the most useful methods of identifying the functional groups of a [protein](#). Chemical modification is also used for labeling proteins with [reporter groups](#) to monitor their conformations for **radiolabeling** and to increase their stability. If proteins are modified chemically under mild conditions, the modifications occur to the native conformation, and the modified proteins usually retain their native conformations. However, conformational changes in the modified proteins sometimes occur, and one must be careful about checking the conformation of the modified proteins, especially when they lose their functions. One must also be careful about the side reactions that often accompany otherwise specific chemical modifications.

It is possible to modify chemically the residues of [aspartic acid](#), [glutamic acid](#), [histidine](#), [lysine](#), [arginine](#), [methionine](#), [tryptophan](#), [tyrosine](#), and [cysteine](#), and less easily [serine](#), [threonine](#), [asparagine](#), and [glutamine](#), but it is not possible to modify specifically [glycine](#), [alanine](#), [valine](#), [leucine](#), [isoleucine](#), [phenylalanine](#), and [proline](#) residues.

Residues that participate in the function are usually **accessible** to the **solvent** and consequently susceptible to reaction with a chemical reagent. The residues that show different reactivity in the presence and absence of a **ligand** are important for the function. A good method for discriminating the functional residues after modification involves analyzing the responsible residues within the

proteins retained by an appropriate affinity column [see [Affinity Chromatography](#)]. [Affinity labeling](#) is a sophisticated way of modifying active site residues selectively. Modification with a **suicide substrate** is a very good method for modifying catalytic residues specifically.

The stability of a protein can be altered by chemical modification (see [Protein Stability](#)). Proteins are stabilized by introducing additional stabilizing forces such as **hydrophobic** forces, [electrostatic interactions](#), and [hydrogen bonds](#). Intra- or intermolecular **cross-links** are also effective for stabilizing proteins, as is attaching polymers to their surfaces.

[Amino groups](#) are reactive nucleophiles that can be modified by many chemical reactions, including acylation, [amidation](#), and [guanidination](#). [Carboxyl groups](#) can be modified after activation with **carbodiimides** or by esterification. Tyrosine residues are modified by [nitration](#) with [tetranitromethane](#), **iodination**, or **acetylation** with acetylimidazole. The [thiol group](#) is the strongest nucleophile among all of the functional groups of amino acids, so there are many reagents that react specifically with it. Among these reactions are metal binding with *p*-mercurylbenzoate (PCMB); mixed-disulfide formation with disulfide reagents, such as 5,5'-dithiobis(2-nitrobenzoic acid) (DTNB, or Ellman's reagent) or dithiodipyridine; [alkylation](#) with **iodoacetate**, methyl iodide, ethyleneimine or *N*-**ethylmaleimide**, cyanylation with 2-nitro-5-thiocyanobenzoic acid (NTCB), and oxidation with many oxidants. Reduction of the [disulfide bonds](#) of proteins and alkylation of the resulting thiol groups is an important step in protein chemistry. The imidazole group of histidine residues can be modified with [diethylpyrocarbonate](#), or they can be photooxidized in the presence of photosensitizing dyes. The sulfur of methionine residues can be oxidized to the sulfoxide by air or by oxidants, or alkylated with agents like methyl iodide under acidic conditions. The latter reaction can be reversed by thiols, so an isotopic label can be introduced in 50% of the terminal methyl group of methionine residues using labeled methyl iodide. The guanido group of arginine residues forms heterocyclic condensation products with 1,2- and 1,3-dicarbonyl compounds, such as phenylglyoxal, 2,3-butanedione, and 1,2-cyclohexanedione. The indole ring of the tryptophan residue can be modified with various oxidants, such as *N*-bromosuccinimide, iodine, or ozone, or by electrophilic reagents, such as 2-hydroxy-5-nitrobenzylbromide or 2-nitrophenylsulfenylchloride.

#### Suggestions for Further Reading

- C.H.W. Hirs (ed.) (1967) Modification reactions. Specific modification reactions, *Methods Enzymol.* **11**, 481–711.
- C.H.W. Hirs (ed.) (1972) Modification reactions. Specific modification reactions, *Methods Enzymol.* **25**, 387–671.
- C.H.W. Hirs and S.N. Timasheff (eds.) (1977) Chemical modification, *Methods Enzymol.* **47**, 407–498.
- C.H.W. Hirs and S.N. Timasheff (ed.) (1983) Chemical modification. Active-site labeling, *Methods Enzymol.* **91**, 549–642.
- T. Imoto and H. Yamada (1989) "Chemical modification". In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, U.K., pp. 247–277.

#### Chemical Shift

A nuclear magnetic resonance (NMR) spectrum of the protons (hydrogen atoms),  $^{13}\text{C}$ , or other nuclei of a molecule typically features groups of signals displayed over a range of frequencies. The frequencies can ultimately be related to the Larmor (resonance) frequencies of the nuclei under

study. The frequencies of absorbed or emitted energy represented in, for example, the proton NMR spectrum are different for protons in structurally or chemically distinct environments. Chemical shift refers to the sensitivity of the Larmor frequency to the covalent structure of which the observed nucleus is part, sensitivity to the intra- and intermolecular interactions present, and sensitivity to other sample variables, such as temperature, concentration, and pressure. The chemical shift is an important aspect of the NMR experiment because it ultimately provides a unique “signature” for each spin of a molecule of interest.

Atomic nuclei such as  $^1\text{H}$ ,  $^{13}\text{C}$ , and  $^{15}\text{N}$ , which have the property referred to as spin, take up certain quantum-mechanically allowed orientations when they are in a magnetic field. In the magnetic field, these nuclei undergo a precessional motion around the direction of the magnetic field. The frequency of this motion ( $\nu$ ) is characteristic of the type of nucleus and its local chemical environment and can be expressed as

$$\nu = \gamma B_0(1 - \sigma)/2\pi$$

where  $\gamma$  is the gyromagnetic ratio of the nucleus,  $B_0$  is the magnitude of the laboratory magnetic field, and  $\sigma$  is the chemical shielding parameter, or screening constant. (The precessional frequency is also dependent on possible scalar and dipolar coupling interactions between nuclei; in terms of energy these are smaller influences on the precessional frequency and are neglected for the present discussion [see scalar coupling].) Chemical and structural information is implicit in the value of the shielding parameter; it is the purpose of NMR experiments to measure nuclear precessional frequencies and thus define the values of  $\sigma$  for nuclear spins of a sample. The shielding parameter is dimensionless and typically is in the range  $1 - 1000 \times 10^{-6}$ . Values for  $\sigma$  are conveniently discussed in terms of parts per million (ppm).

The chemical shift, as represented by the shielding parameter ( $\sigma$ ), for a given spin in a molecule depends principally on the chemical bonds that hold the spin to the molecule. The carbon atom of a methyl group has a shielding parameter that is about 180 ppm different from the shielding parameter of the carbon atom of a carbonyl group, primarily because the local electronic structures about the two carbons are so different (see [NMR \(Nuclear Magnetic Resonance\)](#)). Because local electronic structures for a given chemical group tend to be similar from molecule to molecule, the shielding parameters for the molecules' nuclei tend to fall into narrow, identifiable bands. Thus constructing correlation tables that reliably indicate the approximate expected values of the shielding parameters of each nucleus in a molecule is possible. Organic chemistry textbooks generally contain such tables.

Shielding parameters also subtly reflect the entire electronic structure of a molecule and the electronic structures of solvent molecules or other solutes that are present in a sample. The tertiary structures of biopolymers result in characteristic effects on the shielding parameters of their component nuclei ([1](#), [2](#)).

When considered in detail, it becomes apparent that the value of the shielding parameter is dependent on the orientation of the molecule in a magnetic field. That is,  $\sigma$  is anisotropic and must be represented by a tensor. In gases and most liquid samples, rotational and translational motions of the molecules are rapid, and molecules change their positions and orientations rapidly in the magnetic field. Consequently, in these samples the value of the shielding parameter detected experimentally is the average of all possible orientations of the molecule. If the motions of the molecules in a sample are restricted, as would be the case in solids or in liquid crystalline solutions that are partially ordered, a range of  $\sigma$  values will be detected for a given nucleus.

By convention, NMR spectra are displayed so that shielding parameters (or their equivalents) increase from left to right along a chemical shift axis. That is, the NMR absorption or emission signals for nuclei with the largest shielding parameter appear to the right in the spectrum, whereas nuclei with smaller shielding parameters appear progressively to the left. Also, by convention the

algebraic signs of the numbers along the shielding parameter axis are the negatives of the actual number. Finally, measuring absolute values of shielding parameters with confidence is impractical. The NMR signal from a convenient material is thus chosen as a standard or reference signal, and shielding parameters are measured relative to this signal. Therefore, a signal appearing in a proton NMR spectrum at the position labeled 7 ppm has a shielding parameter that is 7 ppm smaller than the shielding parameter that is characteristic of the reference signal, arbitrarily assigned 0 ppm.

### Bibliography

1. D. S. Wishart and B. D. Sykes (1994) *Methods Enzymol.* **239**, 363–392.
2. S. S. Wijmenga, M. Kruithof, and C. W. Hilbers (1997) *J. Biomol. NMR* **10**, 337–350.

### Suggestions for Further Reading

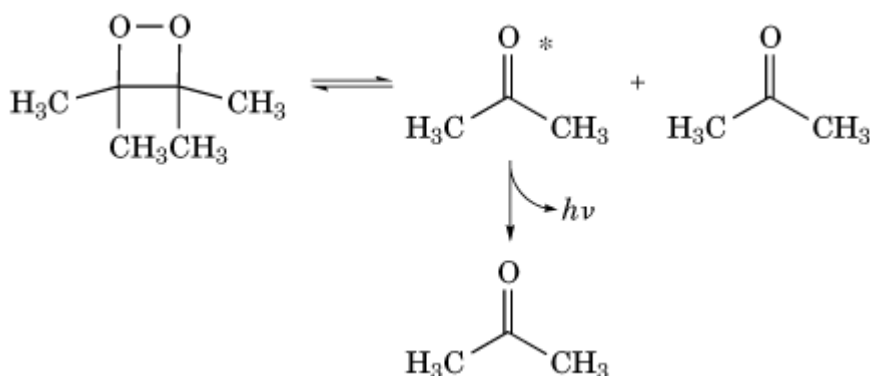
3. R. J. Abraham, J. Fisher, and P. Loftus (1988) *Introduction to NMR Spectroscopy*, Wiley, New York.
4. F. A. Bovey (1988) *Nuclear Magnetic Resonance Spectroscopy*, Academic, San Diego.
5. R. K. Harris (1983) *Nuclear Magnetic Resonance Spectroscopy: a physicochemical view*, Pitman, Marshfield, Mass.
6. S. W. Homans (1992) *A Dictionary of Concepts in NMR*, Clarendon, Oxford.
7. R. Kitamaru (1990) *Nuclear Magnetic Resonance: principles and theory*, Elsevier, New York.
8. R. S. Macomber (1998) *A Complete Introduction to Modern NMR Spectroscopy*, Wiley, New York.
9. C. H. Yoder and C. D. Schaeffer Jr. (1987) *Introduction to Multinuclear NMR*, Benjamin/Cummings, Menlo Park.
10. The NMR literature is replete with theoretical and experimental studies of chemical shielding effects. A database of proton, carbon, and nitrogen chemical shifts in proteins is maintained (<http://www.bmrb.wisc.edu>).

## Chemiluminescence

Chemiluminescence is the production of light via a chemical reaction. As in all [luminescence](#), the source of the light is the decay of an electron from a higher energy excited state to the ground state. In chemiluminescence, the reactants are generally in their ground-state configuration, but one or more of the products is formed with an electron in a high-energy orbital (the product molecule is formed in an excited state). Thus, chemiluminescence is the conversion of chemical energy (Gibbs **free energy**) to radiant energy (light). A necessary precondition for chemiluminescence is that the change in free energy for the reaction be sufficiently large to produce one of the products in an excited-state configuration.

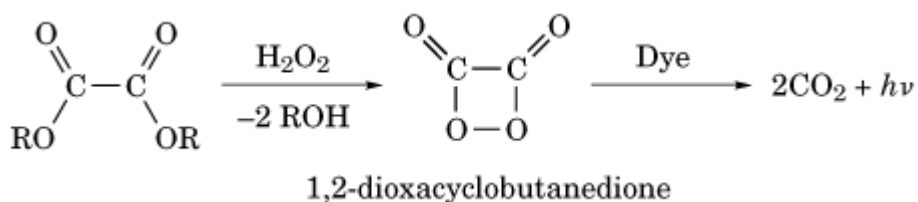
Perhaps the simplest example of a chemiluminescent reaction is the thermal dissociation of the cyclic peroxide tetramethyl 1,2-dioxetane to two molecules of acetone (Fig. [1](#)). The free energy of the [transition state](#) of the reaction is about 90 kcal/mol above the ground state for the acetone products. This is sufficient energy for one of the acetone molecules to be produced in an electronically excited state. When this molecule decays to the ground state, a photon of light is emitted.

**Figure 1.** Simple example of a chemiluminescent reaction, the thermal dissociation of the cyclic peroxide tetramethyl 1,2-dioxetane to two molecules of acetone.



In many cases, a chemiluminescent reaction is carried out in the presence of a dye. The energy from the excited state of the reaction product is transferred to the dye, so that the reaction product is in the ground state and the dye is now in an electronically excited state. The dye will then decay to the ground state and emit a photon of a different color, depending on the chemical structure of the dye. Thus, a reaction that produces UV light may be coupled with a dye that emits in the visible region in order to provide useful illumination. A second reason for coupling a chemiluminescent reaction with a dye is to increase the lifetime of the light emission. Many chemiluminescent reactions are over quickly, and the light is released in a brief, intense flash. If, however, the reaction is coupled with a phosphorescent dye (see [Luminescence](#)), the radiant energy is transferred to the dye and is then released at a much slower rate, producing a longer-lasting, but less intense glow. This is the basis of such commercial products as the glow-stick for emergency lighting, as well as the luminescent novelty necklaces or other items often seen at carnivals. An example of the of reaction used in these devices is shown in Figure 2.

**Figure 2.** Chemiluminescent reaction coupled with a phosphorescent dye. The radiant energy from the chemical reaction is transferred to the dye and is then released at a much slower rate, producing a longer-lasting, but less intense glow. This is the basis of such commercial products as the glow-stick for emergency lighting and luminescent necklaces and other novelty items.



Chemiluminescent reactions are common in living organisms, where they are known as bioluminescent reactions. In some natural systems, energy from the chemical reaction will be passed to “antenna proteins,” which serve the same functions as dyes in the strictly chemical reaction (see [Luciferases And Luciferins](#)).

#### Suggestion for Further Reading

J. D. Roberts and M. C. Caserio (1977) *Basic Principles of Organic Chemistry*, 2nd ed., W. A. Benjamin, Menlo Park, CA, pp. 1371–1418.

## Chemiosmotic Coupling

Many forms of energy are interconverted by living systems. Chemical energy may be used to drive mechanical processes, to generate osmotic and electrical gradients, and even to emit light. In [photosynthesis](#), light energy is converted to chemical energy and is used indirectly to drive the thermodynamically unfavorable synthesis of ATP from ADP and  $P_i$ . In oxidative metabolism in [mitochondria](#) and some bacteria, ATP synthesis is driven by the energy released by oxidation of metabolites. During the 1950s, it was clearly established that both oxidative phosphorylation and photosynthetic phosphorylation are dependent on electron transport. How electron transport could be linked to the formation of a phosphoanhydride bond in ATP was a major question in bioenergetics for the next two decades.

Oxidation–reduction reactions are clearly quite different in character from the removal of the elements of water from ADP and  $P_i$  to form ATP. Yet, there is no doubt that these two processes are coupled to one another. In biochemical parlance, two reactions are said to be coupled when they share a common intermediate. The nature of the intermediate common to electron transfer and to ATP synthesis was elusive. In 1961 in a short, entirely theoretical article, Peter Mitchell proposed a radically different idea for coupling (1). Rather than a chemical intermediate linking electron transport and phosphorylation, as was then generally thought, Mitchell suggested that transmembrane electrochemical proton potentials  $\Delta\tilde{\mu}_{H^+}$  could couple the two reactions. In addition, Mitchell realized that specific **active transport** systems could be linked to  $\Delta\tilde{\mu}_{H^+}$ .

### 1. The Chemiosmotic Hypothesis

The postulates of the chemiosmotic hypothesis (1-3) in brief are as follows:

1. Membranes that catalyze oxidative phosphorylation (the inner mitochondrial membrane and some bacterial plasma membranes) and photosynthetic phosphorylation ([chloroplast](#) thylakoids and plasma membranes of some bacteria) are poorly permeable to protons.
2. Electron transport generates  $\Delta\tilde{\mu}_{H^+}$  by vectorial transport of electrons and protons.
3. An ATPase is driven in reverse (ATP synthesis) by the energetically favorable flow of protons down their electrochemical potential.
4. Coupling membranes contain specific exchange metabolite transport systems that may be linked to the  $\Delta\tilde{\mu}_{H^+}$ .

Before considering these postulates in further detail, we will define  $\Delta\tilde{\mu}_{H^+}$ . The chemical potential of a substance,  $m$ , is the free energy of a system per mole. The energetics of the movement of an ion across a membrane has two components: chemical (osmotic) and electrical. Chemical work must be done to generate a concentration (actually activity) gradient and electrical work, to generate the charge imbalance (the membrane potential). The combined electrochemical potential is  $\tilde{\mu}$ . Just as is the case for Gibbs free energy,  $G$ , it is the change in electrochemical potential ( $\Delta\tilde{\mu}$ ) that is of interest. In general, at constant temperature and pressure

$$\Delta\tilde{\mu} = RT \ln [x^{z+}]_a / [x^{z+}]_b + zF\Delta\psi \quad (1)$$

where  $R$  is the gas constant ( $8.3 \text{ kJ K}^{-1} \text{ mol}^{-1}$ );  $T$ , the absolute temperature;  $\ln, \log_e$ ;  $x^{z+}$ , a cation of  $z$  positive charges;  $a$  and  $b$ , two compartments separated by a membrane;  $z$ , the charge on the cation;  $F$ , Faraday's constant ( $96.5 \text{ kJ mol}^{-1} \text{ V}^{-1}$ ); and  $Dy$ , the membrane potential. For protons,  $z = 1$ , and Equation 1 becomes

$$\Delta\tilde{\mu}_{\text{H}^+} = RT \ln[\text{H}^+]_a/[\text{H}^+]_b + F\Delta\psi \quad (2)$$

or

$$\Delta\tilde{\mu}_{\text{H}^+} = 2.30RT \log[\text{H}^+]_a/[\text{H}^+]_b + F\Delta\psi \quad (3)$$

or, since  $\text{pH} = -\log[\text{H}^+]$ ,

$$\Delta\tilde{\mu}_{\text{H}^+} = 2.30RT(\text{pH}_b - \text{pH}_a) + F\Delta\psi \quad (4)$$

or

$$\Delta\tilde{\mu}_{\text{H}^+} = 2.30RT(\Delta\text{pH}) + F\Delta\psi \quad (5)$$

Mitchell preferred to express Equation 5 in electrical units and coined the term “proton motive force” (pmf or  $Dp$ ), which is simply  $\Delta\tilde{\mu}_{\text{H}^+}$  divided by  $F$ :

$$\Delta p = \Delta\tilde{\mu}_{\text{H}^+}/F = 2.30RT/F(\Delta\text{pH}) + \Delta\psi \quad (6)$$

As forms of energy may be interconverted, so may their units. In some ways it makes sense to consider the energetics of proton gradients in electrical units. Proton fluxes across energy-coupling membranes are analogous to electric circuits. Some people prefer to use  $\Delta\tilde{\mu}_{\text{H}^+}$ , others  $Dp$ . Since the two terms are readily interconvertible, this usage, although potentially confusing, is not a problem. At  $25^\circ\text{C}$ , Equation 5 may be written

$$\Delta\tilde{\mu}_{\text{H}^+} = 5.69 \Delta\text{pH} + 96.5\Delta\psi \quad (7)$$

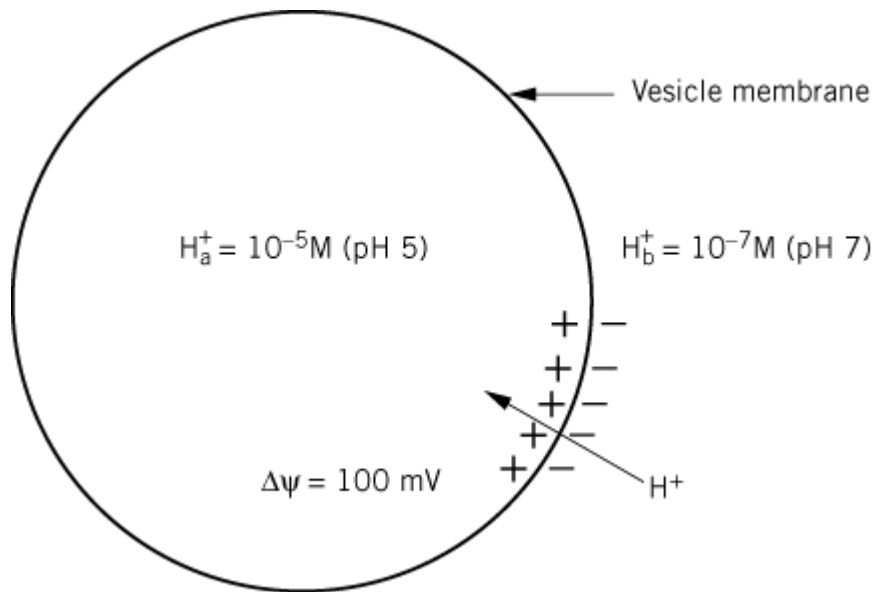
where  $\Delta\tilde{\mu}_{\text{H}^+}$  is also given in kilojoules per mole and  $Dy$  is expressed in volts. The proton motive force (in millivolts) may be expressed as

$$\Delta p = 59 \Delta\text{pH} + \Delta\psi \quad (8)$$

where  $Dy$  is also given in millivolts.

As in any thermodynamic analysis, sign conventions are very important to keep straight for  $\Delta\tilde{\mu}_{\text{H}^+}$ . Consider a topologically closed phospholipid bilayer vesicle (Fig. 1). This vesicle has an inside aqueous compartment,  $a$ , and an external compartment,  $b$ , the suspending medium.  $DpH$  is defined as  $\text{pH}_{\text{out}} - \text{pH}_{\text{in}}$ , which is 2.0 units for the case shown in Figure 1. A charge difference ( $Dy$ ) where the inside of the membrane is more positive than the outside is defined as having a positive sign. Suppose  $Dy$  is +100 mV. From the data given in Figure 1, the  $\Delta\tilde{\mu}_{\text{H}^+}$  or  $Dp$  may be calculated as  $\Delta\tilde{\mu}_{\text{H}^+} = 5.69(2) + 96.5(0.1) = 21.0 \text{ kJ mol}^{-1}$  or  $Dp = 59(2) + 100 = 218 \text{ mV}$ . Notice the sign of the answers, which are positive numbers. These values give the energy cost per mole ( $+\Delta\tilde{\mu}_{\text{H}^+}$  or  $+Dp$ ) to generate the electrochemical protein potential or  $Dp$ . The flow of protons down their electrochemical potential is exergonic, and the absolute magnitudes of  $\Delta\tilde{\mu}_{\text{H}^+}$  and  $Dp$  are the same, but the sign changes. This statement is equivalent to saying that if a chemical reaction has a positive  $DG$  value in the forward reaction, the reverse reaction must have a negative  $DG$  of the same absolute magnitude.

**Figure 1.**  $\Delta\tilde{\mu}_{\text{H}^+}$  Across a vesicular membrane.  $\Delta\text{pH}$  is defined as  $\text{pH}_{\text{out}}$  (compartment *b*)  $-$   $\text{pH}_{\text{in}}$  (compartment *a*), and  $\Delta\psi$  as positive since the inside of the membrane is more positive than the outside.



The essence of chemiosmotic coupling is that exergonic electron transport by the respiratory and photosynthetic electron chains is obligatorily linked to transmembrane proton flux, resulting in the generation of  $\Delta\tilde{\mu}_{\text{H}^+}$ . The flow of protons down their electrochemical potential provides the driving force for ATP synthesis. There is compelling evidence in support of these proposals.

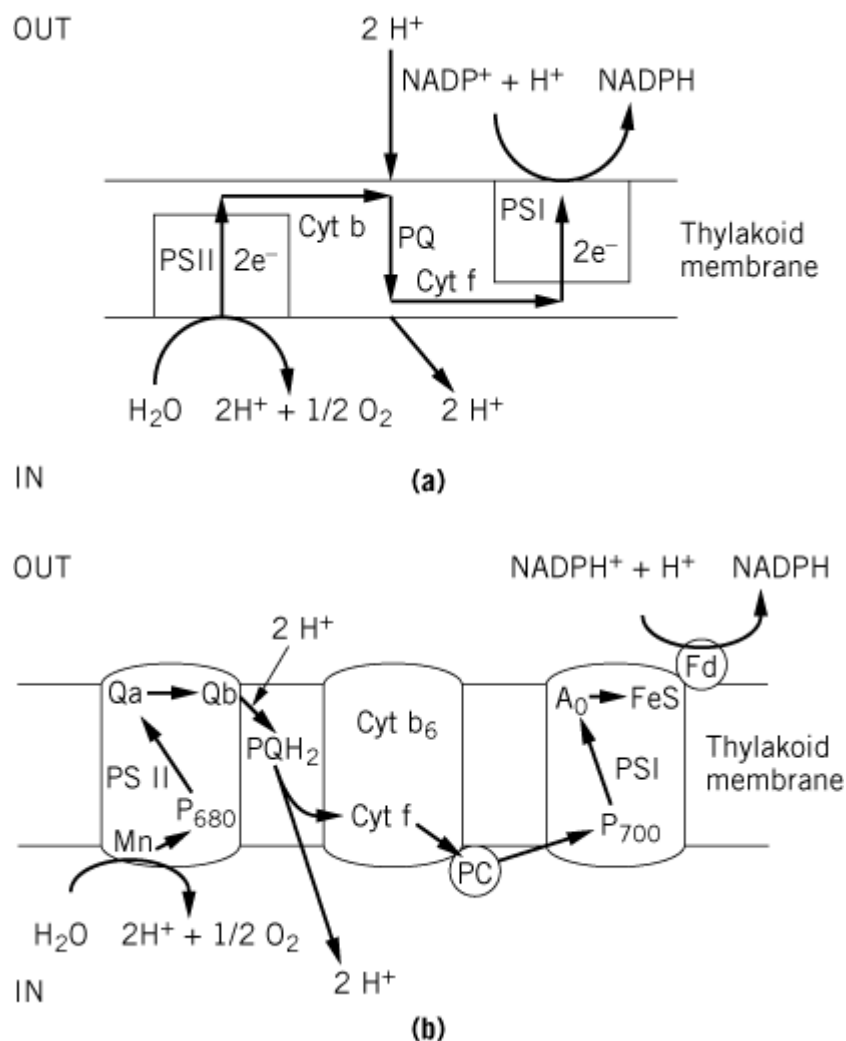
## 2. Evidence in Support of the Chemiosmotic Coupling Hypothesis

As Mitchell predicted (1-3), the mitochondrial, bacterial, and chloroplast membranes that couple ATP synthesis to electron transport are poorly permeable to protons, except when proton-linked processes, such as ATP formation, occur at high rates. Proton transport was shown to be linked to electron transport in mitochondria (4), chloroplasts (5), and bacteria (6). The measurements of the magnitudes of  $\Delta\tilde{\mu}_{\text{H}^+}$  across these membranes turned out to be difficult, but in most instances,  $\Delta\tilde{\mu}_{\text{H}^+}$  values approaching  $20\text{ kJ mol}^{-1}$  ( $\Delta p \sim 200\text{ mV}$ ) have been measured during steady state, rapid electron transport. As predicted by Mitchell, lipophilic weak acids (eg, 2,4-dinitrophenol) could collapse the  $\Delta\tilde{\mu}_{\text{H}^+}$  by shuttling protons across the membrane. ATP synthesis is also inhibited by these reagents, which are termed “uncouplers” because they uncouple electron flow and ATP synthesis.

When Mitchell proposed the chemiosmotic hypothesis, the understanding of biological membrane structure, especially how proteins may be integrated into membranes, was rudimentary. We take for granted today that membrane proteins may be plugged into the lipid bilayer in an asymmetric manner (see [Membranes](#)). But in the 1960s Mitchell's proposals for membrane protein asymmetry, transmembrane electron transport, and membrane sidedness were novel. Although not all of Mitchell's predictions turned out to be correct (cytochrome oxidase, for example, translocates protons via a mechanism that Mitchell did not foresee), many were. His model for proton and electron transport in thylakoid membranes is compared to a much more recent view in Figure 2. There is a remarkable similarity between the two schemes.



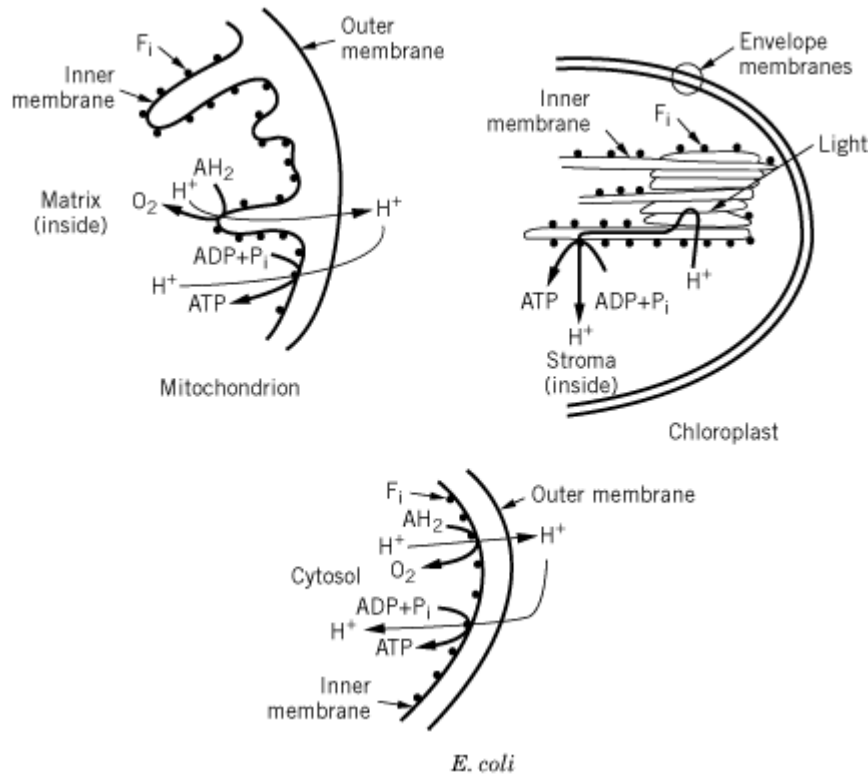
**Figure 2.** Schemes for electron and proton transport in chloroplast thylakoid membranes. In (a) Mitchell's early scheme (2) is illustrated with the membrane in the horizontal position more commonly used by others. In (b) a modern interpretation of electron and proton transport in thylakoids is given. The similarity between the two schemes is remarkable. Note in particular the transmembrane electron transport by photosystem II (PS II) and photosystem I (PS I), the oxidation of water to  $2\text{H}^+$  and  $1/2 \text{O}_2$  inside (the thylakoid lumen) and proton translocation coupled to plastoquinone (PQ) oxidation/reduction. These elements are essentially the same in the two schemes. Cytochrome *b* (Cyt *b*) was not properly placed by Mitchell. In (b) Mn stands for the manganese of the oxygen evolving complex, Qa and Qb, for the bound plastoquinone electron acceptors of PS II;  $\text{P}_{680}$  and  $\text{P}_{700}$ , for the reaction center chlorophylls; Cyt *f*, for cytochrome *f*; PC, for plastocyanin;  $\text{A}_0$ , for the primary electron acceptor in PSI; FeS, for iron–sulfur electron transport proteins; Fd, for [ferredoxin](#). Ferredoxin– $\text{NADP}^+$  oxidoreductase and the FeS protein of the cytochrome *b6f* complex were omitted for clarity.



As illustrated in Figure 3 (7), the sidedness of the [chloroplast](#) thylakoid membrane is opposite that of mitochondria and bacteria (8). As a result of photoelectron transport, thylakoids accumulate protons from the stroma. In contrast, oxidative electron transport in mitochondria and bacteria cause protons to be ejected from the mitochondria or bacteria. The steady state  $\Delta\psi$  across the thylakoid membrane is very low, and the  $\Delta\text{pH}$  is high (as much as 3.5 units). The  $\Delta\text{pH}$  component of the  $\Delta\bar{\mu}_{\text{H}^+}$  in mitochondria and bacteria is, however, usually small, and  $\Delta\psi$  is high. These observations make sense when it is realized that the mitochondrial matrix and bacterial cytoplasm are the internal compartments bounded by the coupling membranes. The pH of these compartments must be controlled to allow the numerous enzymes they contain to function. In contrast, the thylakoid lumen contains no metabolic enzymes. Electrical neutrality in thylakoids is maintained by  $\text{Cl}^-$  flux.

Essentially, thylakoids accumulate HCl in the light.

**Figure 3.** Sidedness of coupling membranes. Mitochondrial and bacterial (*E. coli*) electron transport is coupled to the ejection of protons, whereas in thylakoids, protons are accumulated. Note that, in bacteria and chloroplasts, ATP is generated within the same compartment in which it is used. In contrast, most of the ATP synthesized within mitochondria is exported to the cytoplasm. From Ref. 47, with permission.



Dramatic evidence in favor of the chemiosmotic hypothesis was obtained in André T. Jagendorf's laboratory during the mid-1960s (9, 10). Simply by adjusting the pH of thylakoid suspensions to pH 4 and rapidly raising the pH to 8.0 in the presence of ADP + P<sub>i</sub>, relatively large amounts of ATP are formed *in the dark*. ATP synthesis induced by acid to base transitions is abolished by uncouplers and reagents that block ATP synthesis. Later, it was shown (11) that artificially generated electric fields could also generate ATP. In mitochondrial vesicles, Thayer and Hinkle (12) showed that the rate of Dp-driven ATP synthesis is as fast or faster than that coupled to electron transfer.

The membrane of the halophilic bacterium, *Halobacter halobium* (*salinarium*) elaborates purple patches composed entirely of bacteriorhodopsin. Bacteriorhodopsin is a light-dependent proton pump that, when incorporated into lipid membranes, can generate substantial  $\Delta\tilde{\mu}_{H^+}$  (12). When the enzyme from bovine heart mitochondria that catalyzes ATP synthesis was coreconstituted with bacteriorhodopsin, vesicles were obtained that catalyzed light-dependent, uncoupler-sensitive ATP formation (14).

Mitchell's proposal that ATP synthesis is accomplished by driving a proton-linked ATPase in reverse was also novel. By the time Mitchell was developing the chemiosmotic hypothesis, it had been clearly established that the energy of ATP hydrolysis could be used to generate ion gradients. For example, the Na<sup>+</sup> and K<sup>+</sup> concentration gradients and Dy across mammalian plasma membranes are generated by an ATPase, the Na<sup>+</sup>,K<sup>+</sup>-ATPase. If, Mitchell reasoned, an ATPase acts as an ion pump, might an ATPase operate in reverse to utilize the  $\Delta\tilde{\mu}_{H^+}$  generated by electron transport to drive ATP

synthesis? Work on “coupling factors” during the 1960s established that mitochondria and chloroplasts contain an enzyme that is required for ATP synthesis that could under appropriate conditions also hydrolyze ATP (15). The catalytic part of what is now known as “ATP synthase” is called “F”<sub>1</sub> and is an extrinsic membrane protein. F<sub>1</sub> may be removed from the membrane and is very water-soluble. (See article on [ATP Synthase](#).)

Mitochondrial ATP synthesis and hydrolysis are strongly inhibited by the antibiotic, [oligomycin](#). The ATPase activity of F<sub>1</sub>, after its removal from the inner mitochondrial membrane, is insensitive to oligomycin—an observation that seemed incompatible with a role of F<sub>1</sub> in oxidative phosphorylation. Racker then showed that a crude detergent fraction of mitochondrial inner membranes contain a factor that confers sensitivity to oligomycin to F<sub>1</sub>. This factor was called “F<sub>0</sub>”; F<sub>0</sub> is not a coupling factor, and it is improper to denote F<sub>0</sub> as “F zero” or “F naught.” Kagawa and Racker (16) pioneered the isolation and purification of the entire ATPase complex, F<sub>1</sub> – F<sub>0</sub>, which could be incorporated into phospholipid vesicles that catalyzed energy-linked activities. Similar studies were carried out with the chloroplast (CF<sub>1</sub> – CF<sub>0</sub>) and bacterial enzymes.

Mitchell predicted that F<sub>0</sub> contains a mechanism for transmembrane proton movement. Removal of CF<sub>1</sub> was shown (17) to greatly enhance the proton permeability of the thylakoid membrane. *N,N'*-Dicyclohexylcarbodiimide (DCCD) at low concentrations inhibits ATP synthesis and ATP hydrolysis by F<sub>1</sub> – F<sub>0</sub> (17). DCCD acts by blocking proton conductance by F<sub>0</sub> and, in thylakoids, restores net light-dependent proton uptake and high  $\Delta\tilde{\mu}_{H^+}$  to membranes from which CF<sub>1</sub> had been removed. The DCCD reacts with an Asp or Glu residue in a very hydrophobic, 8-kDa polypeptide of F<sub>0</sub> (18).

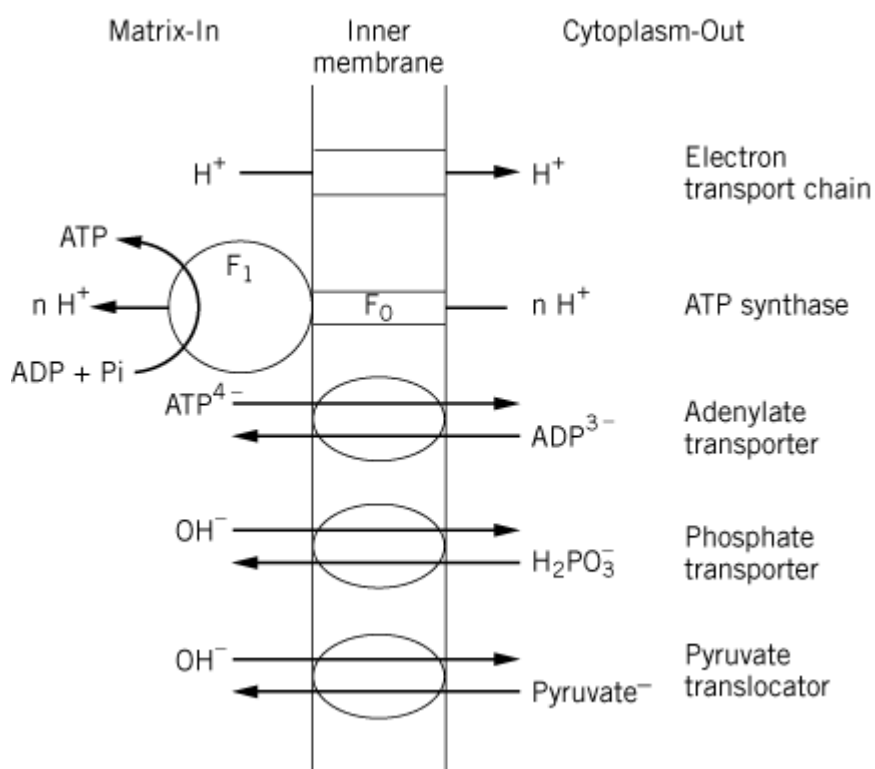
The evidence that F<sub>1</sub> – F<sub>0</sub> is a proton-translocating ATPase/ATP synthase is overwhelming. Phosphorylation driven by imposed  $\Delta pH$  or  $\Delta y$  is abolished by specific inhibitors of either F<sub>1</sub> or F<sub>0</sub> function. As described above, F<sub>0</sub> conducts protons rapidly, and blocking F<sub>0</sub> proton transport also inhibits ATP synthesis. ATP hydrolysis by F<sub>1</sub> – F<sub>0</sub> is linked to proton translocation. In thylakoids ATP hydrolysis in the dark can generate  $\Delta\tilde{\mu}_{H^+}$  values of significant magnitude to drive some ATP synthesis. This ATP-dependent ATP synthesis is very sensitive to inhibition by uncouplers. The  $\Delta\tilde{\mu}_{H^+}$  generated by ATP hydrolysis can drive electron transport in reverse (19), a process that is blocked by ATP synthase inhibitors and uncouplers. Other evidence that the electron-transport chain and ATP synthase are linked via  $\Delta\tilde{\mu}_{H^+}$  includes the facts that phosphorylation decreases the magnitude of  $\Delta\tilde{\mu}_{H^+}$  and stimulates the rate of electron transport. Specific inhibition of ATP synthesis increases the  $\Delta\tilde{\mu}_{H^+}$  and slows electron transport. Uncouplers strongly decrease  $\Delta\tilde{\mu}_{H^+}$  and inhibit ATP synthesis, but stimulate electron transport. An extensive analysis of the coupling between photosynthetic electron transport and ATP synthesis in the steady state gave results strongly indicating that  $\Delta\tilde{\mu}_{H^+}$  is the intermediate linking the two processes.

### 3. Metabolite Transport and Motility

Mitchell's fourth postulate, that coupling membranes contain specific exchange transport systems, has also received broad experimental support, especially for mitochondrial and bacterial coupling membranes. The chloroplast thylakoid membrane does not appear to contain metabolite transporters. ATP is synthesized in chloroplasts within the same compartment (the stroma) in which it is utilized in photosynthesis (see Fig. 3). The inner membrane of the envelope, not the thylakoid membrane, is the permeability barrier of the chloroplast. In contrast, the inner membrane of the mitochondrion and the plasma membrane of bacteria are both coupling membranes and permeability barriers. In these membranes, transport must coexist with ATP synthesis.

Most of the ATP generated by oxidative phosphorylation in mitochondria is exported to the cytoplasm and is hydrolyzed to drive cellular processes. Mitochondria must contain mechanisms to export ATP and import ADP and  $P_i$ . In addition, mitochondria will take up pyruvate, di- and tricarboxylic acids, and some amino acids. Cation exchangers are also present. Many of these transporters utilize the  $\Delta\tilde{\mu}_{H^+}$  established by the electron transport chain (Fig. 4). A few transporters will be briefly considered in this article, with emphasis on ADP/ATP and  $P_i$  transporters in mitochondria. (see [Mitochondria](#)).

**Figure 4.** Proton circuits in the mitochondrial inner membrane. Electron transport ejects protons from the mitochondrion, resulting in the formation of  $\Delta\tilde{\mu}_{H^+}$ . A number of processes including ATP synthesis, ADP/ATP exchange transport,  $P_i - OH^-$  exchange transport, and the uptake of certain metabolites (illustrated by pyruvate uptake in exchange for  $OH^-$ ) are linked to the  $\Delta\tilde{\mu}_{H^+}$ .



At physiological pH, ADP exists as  $ADP^{3-}$ , and ATP as  $ATP^{4-}$ . The inner membrane contains a well-characterized transporter, known as the *adenine nucleotide translocator* (20). This transporter catalyzes the counterexchange translocation of  $ADP^{3-}$  and  $ATP^{4-}$ . The one-for-one exchange transport (**antiport**), with  $ATP^{4-}$  exported from the mitochondrion and  $ADP^{3-}$  imported, is electrogenic; there would be the net transfer of one negative charge out of the mitochondrion per transport event. The  $\Delta\psi$  across the mitochondrial inner membrane is outside positive. Thus,  $ATP^{4-}/ADP^{3-}$  exchange transport would be promoted in the direction of ATP export and ADP import by the  $\Delta\psi$  generated by electron flow.

$P_i$  accumulation by mitochondria is also catalyzed by a transporter.  $P_i$  is either cotransported with  $H^+$  or counter-exchanged for  $OH^-$ . Since the form of  $P_i$  that is preferentially transported is  $H_2PO_4^-$ , one-for-one transport of  $P_i$  with either  $H^+$  (**symport**) or  $OH^-$  (**antiport**) is electrically neutral and is not

affected by  $\Delta\psi$ .  $P_i$  transport is, however, linked to  $\Delta\text{pH}$ . The pH of the matrix of respiring mitochondria is about 0.5 unit more alkaline than that of the cytoplasm. Either  $\text{OH}^-$  antiport or  $\text{H}^+$  symport would favor  $P_i$  uptake into mitochondria. This mechanism and the rapid utilization of  $P_i$  by oxidative phosphorylation assure that the  $P_i$  is delivered to the mitochondrial matrix from the cytoplasm. The transport of ATP/ADP and  $P_i$  by mitochondria costs the equivalent of one proton/ATP synthesized. Pyruvate, the end product of aerobic glycolysis, is transported into mitochondria via an electrically neutral counterexchange with  $\text{OH}^-$ , and its transport, like that of  $P_i$ , is linked to  $\Delta\text{pH}$ .

Bacteria have the potential to carry out the active transport of a large number of compounds, including sugars and amino acids (21). Transport may be linked to the  $\Delta\mu_{\text{H}^+}$  generated by electron transport, or, under anaerobic conditions, by ATP hydrolysis by  $F_1 - F_0$ . The best-studied example of  $\Delta\mu_{\text{H}^+}$ -dependent active transport in bacteria is the lactose transporter (lac **permease**) of *Escherichia coli* (22). Lactose uptake is obligatorily linked to  $\text{H}^+$  uptake into the cell, at a probable stoichiometry of one lactose to one proton. The bacterium maintains, by either electron transport or ATP hydrolysis, a negative  $\Delta\psi$  (outside positive), and, depending upon the conditions, a negative  $\Delta\text{pH}$  (inside alkaline). The flow of protons down the  $\Delta\mu_{\text{H}^+}$  provides the driving force for the accumulation of lactose to concentrations in excess of 1,000 times that of the medium. In addition to proton-linked transport systems, bacteria contain ABC (ATP binding cassette) transporters (23) that utilize ATP for transport to metabolites or ions directly, as well as the phosphotransferase system (24) that utilizes phosphoenolpyruvate as both an energy source and a source of phosphate for the transport and concurrent phosphorylation of some sugars.

The concepts developed by Mitchell for metabolite and ion transport were subsequently found to apply to membranes other than coupling membranes. The plasma membranes of higher plant cells, yeast, and fungi contain an enzyme that is structurally, functionally, and mechanistically related to the  $\text{Na}^+$ ,  $\text{K}^+$ -ATPase of animal cell plasma membranes. This enzyme, the  $\text{H}^+$ -ATPase (25), is the primary ion pump in plant, yeast, and fungal plasma membranes. ATP hydrolysis pumps protons out of the cells, generating the transplasma membrane  $\Delta\mu_{\text{H}^+}$ . The  $\Delta\text{pH}$  component in plant cells is 1.0–2.0 pH units, whereas  $\Delta\psi$  may be in excess of –200 mV. The active transport across the plasma membrane is driven by the energetically favorable flow of protons down the  $\Delta\mu_{\text{H}^+}$ . As is the case for lactose transport in *E. coli*, specific transporters that mediate the translocation of metabolites across the membrane are obligatorily coupled to  $\text{H}^+$  transport. An important example of  $\text{H}^+$ : metabolite cotransport (symport) in plants is the active loading of sucrose (26) into the phloem. As a consequence of sucrose:  $\text{H}^+$  cotransport, the sucrose concentration within the phloem may be as high as 0.5 M.

An entirely new class of  $\text{H}^+$ -ATPases was discovered during the 1980s and characterized in the 1990s (27). The members of this class are called “V-ATPases,” in which the V stands for *vacuolar*. All eukaryotic cells contain V-ATPases. In plant and yeast cells, the V-ATPase is present on the vacuolar membrane and in the [Golgi apparatus](#). In animal cells, V-ATPases are present in **clathrin**-coated vesicles, chromaffin granules (adrenal medulla), Golgi membranes, **lysosomes**, synaptic vesicles, and, in some specialized cells, the plasma membrane. The V-ATPases are structurally much more complex than plasma membrane  $\text{H}^+$ -ATPases, and, although the V-ATPases large subunits have some sequence similarity to the a and b subunits of  $F_1$ , V-ATPases and  $F_1 - F_0$  ATPases are distinct.

The function of V-ATPases is to hydrolyze ATP and to pump protons from the cytoplasm into the internal aqueous compartment of endomembranes. The interiors of vacuoles and lysosomes are more acidic than that of the cytoplasm. Metabolite and ion fluxes are linked to the  $\Delta\mu_{\text{H}^+}$  generated by the V-ATPase. A  $\text{Ca}^{2+}/\text{H}^+$  exchange transporter has, for example, been found in the vacuolar membrane

of higher plants (28). Acidification of the vesicles that are formed by receptor-mediated endocytosis is required for dissociation of the ligand from its receptor (29).

The V-ATPase of eukaryotes never functions as an ATP synthase. Even in the best of conditions *in vitro*, imposed proton gradients will not cause the pump to reverse. The H<sup>+</sup>-ATPase of **archaebacteria** (eg, *Halobacter salinarium* [*halobium*], (30)) more closely resembles V-ATPases than F<sub>1</sub> – F<sub>0</sub>. Yet, the V-ATPase (sometimes called “A-ATPase” for Archeon), of the archeon can make ATP at the expense of the  $\Delta\tilde{\mu}_{H^+}$  generated by the light-dependent bacteriorhodopsin proton pump. Eukaryotes tailored the V-ATPase to fit their needs.

In addition to its utilization in ATP synthesis and membrane transport,  $\Delta\tilde{\mu}_{H^+}$  can be converted to mechanical energy in the form of motility of bacteria with **flagella** such as *E. coli* and *Bacillus subtilis*. Flagellar rotation is driven by  $\Delta\tilde{\mu}_{H^+}$ , rather than directly by ATP. Starved *B. subtilis* is immotile, but motility may be restored in transient manner by the artificial imposition of  $\Delta\psi$ .  $\Delta\psi$  (inside alkaline) is also capable of driving flagellar motion (31).

#### 4. Sodium Electrochemical Potential

The electrochemical sodium ion potential,  $\Delta\tilde{\mu}_{Na^+}$ , is given by

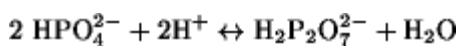
$$\Delta\tilde{\mu}_{Na^+} = RT \ln [Na^+]_a / [Na^+]_b + F \Delta\psi \quad (9)$$

where  $R$ ,  $T$ ,  $\ln$ ,  $\Delta\psi$ , and  $F$  are as defined in Equation 1. In some bacteria,  $\Delta\tilde{\mu}_{Na^+}$  is exploited for ATP synthesis, active transport, and motility. In addition, Na<sup>+</sup>/metabolite cotransport systems are present in animal cells. Transporters within the plasma membrane of these cells link the translocation of the substrates (eg, amino acids) to the flux of Na<sup>+</sup> down the  $\Delta\tilde{\mu}_{Na^+}$  generated by the Na<sup>+</sup>, K<sup>+</sup>-ATPase (32). In bacteria, part of the  $\Delta\tilde{\mu}_{H^+}$  generated by respiration may be converted to  $\Delta\tilde{\mu}_{Na^+}$  by counterexchange transport of Na<sup>+</sup> for H<sup>+</sup>. In some bacteria, Na<sup>+</sup> translocation may be linked to metabolism or electron transport without the involvement of  $\Delta\tilde{\mu}_{H^+}$ . For example, in *Propionigenium modestum*, the anaerobic decarboxylation of succinate to form propionate is coupled to the transport of Na<sup>+</sup> out of bacterium, generating a  $\Delta\tilde{\mu}_{Na^+}$  (33). In **halophiles**, Na<sup>+</sup> transport has been shown to be driven by electron transport through the NADH dehydrogenase region of the electron transport chain.

Vesicles of the *P. modestum* plasma membrane catalyze ATP-dependent Na<sup>+</sup> transport, and an imposed  $\Delta\tilde{\mu}_{Na^+}$  drives ATP synthesis (34). The membrane contains a Na<sup>+</sup>-ATPase (ATP synthase) that is analogous to the F<sub>1</sub> – F<sub>0</sub> ATP synthase of coupling membranes. It is interesting to note that F<sub>0</sub> subunits of the *P. modestum* Na<sup>+</sup>-ATPase confer to F<sub>1</sub> from *E. coli* the ability to translocate Na<sup>+</sup> (35). One archeon, *Methanosarina mazei* Gö 1, contains an F<sub>1</sub> – F<sub>0</sub> that uses  $\Delta\tilde{\mu}_{Na^+}$  and a V-ATPase (V<sub>1</sub> – V<sub>0</sub>) that uses  $\Delta\tilde{\mu}_{H^+}$  (36). It is clear, therefore, that proton transport per se is not required for ATP synthesis by F<sub>1</sub> – F<sub>0</sub>. This observation makes mechanisms of ATP synthesis in which translocated protons are directly involved in catalysis very unlikely.  $\Delta\tilde{\mu}_{Na^+}$  may also be the driving force for motility in some marine bacteria.

#### 5. Pyrophosphate and Energy Coupling

Pyrophosphate (PP<sub>i</sub>) is merely two P<sub>i</sub> molecules joined by an anhydride bond:



The equilibrium for this reaction lies far to the left; at pH 7 and 25°C, the standard free-energy change for (PP<sub>i</sub>) synthesis from 2 P<sub>i</sub> is about +16.7 kJ mol<sup>-1</sup>. PP<sub>i</sub> is a product of several biosynthetic reactions, including aminoacyl tRNA synthesis, fatty-acid activation, and DNA and RNA synthesis. These reactions occur with equilibrium constants close to 1 and are freely reversible. By coupling these reactions to PP<sub>i</sub> hydrolysis, the reactions are pulled in the synthetic direction. Cells contain soluble [pyrophosphatases](#) that hydrolyze PP<sub>i</sub> to 2 P<sub>i</sub>.

Other roles for PP<sub>i</sub> in plants have recently been elucidated. The vacuolar membrane contains a pyrophosphatase (37) that couples pyrophosphate hydrolysis to inward proton translocation. Both the V-ATPase and the pyrophosphatase contribute significantly to the  $\Delta\tilde{\mu}_{\text{H}^+}$  across the vacuolar membrane. Thus, some of the energy of PP<sub>i</sub> hydrolysis is conserved by generation of the  $\Delta\tilde{\mu}_{\text{H}^+}$ . In addition, PP<sub>i</sub> may be utilized as a phosphoryl donor, as an alternative to ATP, in glycolysis.

In the photosynthetic bacterium, *Rhodospirillum rubrum*, a pyrophosphatase, and F<sub>1</sub> – F<sub>0</sub> ATPase (ATP synthase) exist within the same membrane (38). Net PP<sub>i</sub> synthesis driven by the photoelectron transport chain of the bacterium has been observed. PP<sub>i</sub> synthesis is inhibited by proton ionophores (uncouplers), but not by reagents that block ATP synthesis by interacting with F<sub>1</sub>. Membrane vesicles (chromatophores) from *R. rubrum* are oriented so that the catalytic sites of both the ATP synthase and the pyrophosphatase face the outside. ATP hydrolysis by these vesicles in the dark is coupled to inward H<sup>+</sup> translocation, resulting in the generation of substantial  $\Delta\tilde{\mu}_{\text{H}^+}$  values. It is even possible to detect ATP-dependent pyrophosphate synthesis that is also uncoupler-sensitive. ATP hydrolysis and pyrophosphate synthesis are linked by the  $\Delta\tilde{\mu}_{\text{H}^+}$ .

## 6. Concluding Remarks

The four postulates of the chemiosmotic coupling hypothesis have received broad experimental support. Coupling membranes are poorly permeable to H<sup>+</sup>. Electron transport results in the generation of  $\Delta\tilde{\mu}_{\text{H}^+}$  at values that are consistent with the energetic demands of ATP synthesis. An ATPase is present that is reversible and can make ATP, provided the magnitude of  $\Delta\tilde{\mu}_{\text{H}^+}$  is sufficient, or hydrolyze ATP if the  $\Delta\tilde{\mu}_{\text{H}^+}$  is low. Artificially imposed  $\Delta\tilde{\mu}_{\text{H}^+}$  can drive ATP synthesis at rates as fast as in oxidative or photosynthetic phosphorylation.  $\Delta\tilde{\mu}_{\text{H}^+}$  is energetically and kinetically competent to be the intermediate common between electron transport and ATP synthesis. The fourth postulate, the existence of exchange transporters and their possible coupling to the  $\Delta\tilde{\mu}_{\text{H}^+}$ , has had a tremendous impact on the field of metabolite transport, not only in mitochondria but also in many other membranes. In part, the chemiosmotic hypothesis had its origins in Mitchell's long-standing interest in active transport in bacteria. The chemiosmotic hypothesis provided a new theoretical basis for the interpretation of transport experiments. Very early on, Mitchell proposed that membrane proteins may be organized in specific orientations within membranes and that electron transfer may in part occur across membranes. For his remarkable insights and experimental contributions, Peter Mitchell was awarded the Nobel Prize in Chemistry in 1978.

That H<sup>+</sup> (or Na<sup>+</sup>) translocation is coupled to ATP synthesis and hydrolysis is generally accepted. Some evidence suggests that under some circumstances, the protons translocated during electron transport need not equilibrate with the bulk aqueous phase prior to their passage through F<sub>1</sub> – F<sub>0</sub> to make ATP. Because of the high buffering capacity of the interior of thylakoids, 50 ms or more of illumination may be required to generate ΔpH values of magnitudes sufficient to drive ATP synthesis. When the generation of Δy by proton transport is prevented, permeant buffers should delay the onset of ATP synthesis by increasing the internal buffering capacity. The delay is not as long as it should be, if it is assumed that the added buffer is uniformly distributed throughout the internal volume of the thylakoids (39, 40). These and other observations kept alive the localized

proton theory of R. J. P. Williams (41). According to this theory, dubbed “microchemiosmosis” by Mitchell, protons are trapped in localized regions of the membrane that are not necessarily in equilibrium with each other. The transmembrane bulk phase  $\Delta\tilde{\mu}_{H^+}$  was proposed to be used only for energy storage.

It is difficult to design experiments that provide unequivocal support for the localized proton hypothesis. If protons may be delivered to the ATP synthase without equilibrating with the bulk aqueous phase, some barrier must exist to prevent equilibration. Disruption of thylakoid structure has been suggested to convert the coupling mechanism from localized to delocalized (ie, chemiosmotic). It would thus be very difficult to establish the nature of the barrier. Results of a number of other reported experiments at first glance seem at odds with the chemiosmotic hypothesis. In most cases, flaws in the methods used or in the interpretation of the results can explain the apparent discrepancy.

The chemiosmotic hypothesis has stood the test of time. During the past few years, much has been learned about the structure of some of the proton and electron translocators of coupling membranes. Yet, we know little about how cytochrome oxidase pumps protons or how the ATP synthase exploits the  $\Delta\tilde{\mu}_{H^+}$  to drive ATP synthesis. The emphasis of bioenergetics has shifted away from *what* couples electron transport to ATP synthesis to *how* the proton translocating elements work.

### Bibliography

1. P. Mitchell (1961) *Nature* **191**, 144–148.
2. P. Mitchell (1966) *Chemiosmotic Coupling in Oxidative and Photosynthetic Phosphorylation*, Glynn Research, Bodmin, Cornwall, England.
3. P. Mitchell (1979) *Les Prix Nobel en 1978, Nobel Foundation*, Stockholm, pp. 142–143.
4. P. Mitchell and J. Moyle (1965) *Nature* **298**, 147–151.
5. J. S. Neuman and A. T. Jagendorf (1964) *Arch Biochem Biophys.* **107**, 109–119.
6. P. Mitchell (1962) *J. Gen. Microbiol.* **29**, 144–148.
7. R. E. McCarty (1985) *Bioscience* **35**, 27.
8. P. C. Hinkle and R. E. McCarty (1978) *Sci. Am.* **238**, 104–123.
9. A. T. Jagendorf and E. Uribe (1966) *Proc. Natl. Acad. Sci. USA* **55**, 170–177.
10. A. T. Jagendorf (1967) *Fed. Proc.* **26**, 1361–1369.
11. H. T. Witt (1979) *Biochim. Biophys. Acta* **505**, 355–427.
12. W. S. Thayer and P. C. Hinkle (1975) *J. Biol. Chem.* **250**, 5330–5335.
13. W. Stoekenius (1985) *Trends Biochem Sci.* **10**, 483–485.
14. E. Racker and W. Stoekenius (1974) *J. Biol. Chem.* **249** 662–663.
15. E. Racker (1976) *A New Look at Mechanisms in Bioenergetics*, Academic Press, New York.
16. Y. Kagawa, A. Kandrach, and E. Racker (1973) *J. Biol. Chem.* **248**, 676–684.
17. R. E. McCarty and E. Racker (1967) *Brookhaven Symp. Biol.* **19**, 202–212.
18. K. J. Cattell, C. R. Lindop, I. G. Knight, and R. B. Beechey (1971) *Biochem. J.* **125**, 169–177.
19. B. Chance and G. Hollunger (1961) *J. Biol. Chem.* **258**, 1474–1486.
20. M. Klingenberg (1989) *Arch. Biochem Biophys* **270**, 1–14.
21. R. D. Simoni and P. W. Postman (1975) *Annu. Rev. Biochem.* **44**, 523–554.
22. H. R. Kaback (1986) *Annu. Rev. Biophys. Chem.* **15**, 279–319.
23. C. F. Higgins (1992) *Annu. Rev. Cell Biol.* **8**, 67–113.
24. N. D. Meadow, D. K. Fox, and S. Roseman (1990) *Ann. Rev. Biochem.* **59**, 497–542.
25. R. Serrano (1990) in C. Larsson and I. M. Moller, eds., *The Plant Plasma Membrane: Structure, Function and Molecular Biology*, Springer-Verlag, New York, pp. 127–153.
26. N. Sauer (1992) in D. T. Cooke and D. T. Clarkson, eds., *Transport and Receptor Proteins of Plant Membranes: Molecular Structure and Function*, Plenum Press, New York, pp. 67–75.



27. N. Nelson (1989) *J. Bioenerg. Biomembr.* **21**, 533–571.
28. K. S. Schumaker and H. Sze (1986) *J. Biol. Chem.* **261**, 12172–12178.
29. I. Melman (1992) *J. Exp. Biol.* **172**, 39–45.
30. K. Ihara, T. Abe, K.-I. Sugimura, and Y. Mukohata (1992) *J. Exp. Biol.* **172**, 475–485.
31. S. Matsuura, J. Shioi, and Y. Imae (1977) *FEBS Lett.* **82**, 187–190.
32. E. J. Collarini and D. L. Oxender (1987) *Annu. Rev. Nutr.* **7**, 75–90.
33. P. Dimroth (1991) *Bioessays* **13**, 463–468.
34. W. Hilpert, B. Schink, and P. Dimroth (1984) *EMBO J.* **3**, 1665–1670.
35. W. Laubinger, G. Dekers-Hebestreit, N. Altendorf, and P. Dimroth (1990) *Biochemistry* **25**, 5458–5463.
36. R. Wilms, C. Frieberg, E. Wegerk, I. Meier, F. Mayer, and V. Müller (1996) *J. Biol. Chem* **271**, 18843–18852.
37. P. A. Rea, Y. Kim, V. Sarafian, R. J. Poole, J. M. Davies, and D. Sanders (1992) *Trends Biochem. Sci.* **17**, 348–353.
38. P. Nyring, B. F. Nore, and A. Strid (1991) *Biochemistry* **30**, 2883–2887.
39. R. D. Horner and E. N. Moudrianakis (1983) *J. Biol. Chem.* **258**, 11643–11647.
40. R. A. Dilley, S. M. Theg, and W. A. Beard (1987) *Annu. Rev. Plant Physiol.* **38**, 347–389.
41. R. J. P. Williams (1969) *Curr. Top. Bioenerg.* **3**, 79–156.

### Suggestions for Further Reading

42. W. A. Cramer and D. B. Kraff (1990) *Energy Transduction in Biological Membranes. A Text of Bioenergetics*, Springer-Verlag, New York.
43. G. D. Greville (1969) A scrutiny of Mitchell's chemiosmotic hypothesis of respiratory chain and photosynthetic phosphorylation, *Curr. Top. Bioenerg.* **3**, 1–78 (a very clear and useful review of the state of chemiosmotic coupling in the early days).
44. D. G. Nichols and S. J. Ferguson (1992) *Bioenergetics*, Academic Press, London.
45. R. N. Robertson (1968) *Protons, Electrons, Phosphorylation and Active Transport*, Cambridge University Press, London (in this slender volume, the author describes thoughts about charge separations and ion fluxes that predated the chemiosmotic hypothesis).
46. V. P. Skulachev and P. C. Hinkle, eds. (1982) *Chemiosmotic Proton Circuits in Biological Membranes*, Addison-Wesley, London.

### Chemokines

Chemokines are small secreted [proteins](#), generally 8 to 15 kDa in mass, that were originally identified by their ability to stimulate [chemotaxis](#) and/or activate leukocytes (lymphocytes, neutrophils, eosinophils, mast cells, monocytes, and macrophages). They are now known to be multifunctional and to be secreted by many different cell types in response to various stimuli (for reviews see Refs. [1-5](#)). In inflamed tissues, they may induce effector functions such as the generation of an oxidative burst or the secretion, from storage granules, of [proteinases](#) (by neutrophils and monocytes), histamine (by basophils), and cytotoxic proteins (by eosinophils). Chemokines activate **cell-surface receptors** to promote adhesion of circulating leukocytes to the vascular endothelium and subsequently to facilitate extravasation (diapedesis) of the adherent leukocyte into adjacent tissue, the source of the chemokine. They have proinflammatory effects and can give rise to acute or

chronic inflammatory responses; some have been shown to be involved in movement and activation of cells in both angiogenesis/vasculogenesis and neural pathfinding, thus implicating them in aspects of [development](#) not directly involved in the [immune response](#).

## 1. Chemokine Structure and Chromosome Location

The chemokine family (at least 40 members have been identified in humans) consists of four classes distinguished by the positioning of four conserved [cysteine](#) residues near the N-terminus of the protein. These four cysteine residues define a distinct structural motif based on the two [disulfide bonds](#) (C16C3, C26C4) that they form. All the chemokines possess a central [b-pleated sheet](#) region consisting of three antiparallel [beta-strands](#) (a [Greek key motif](#)) followed in most cases by a C-terminal [alpha-helix](#). There is evidence that the chemokines can form dimers and higher oligomers, but whether this is important for their function is controversial (2). The N-terminal sequence of the chemokines is of particular importance in determining their specificity. Alteration or removal of one or a few [amino acid](#) residues can dramatically affect their activity or target cell specificity. This provides the opportunity for the specificity of a chemokine to be altered by local factors subsequent to its secretion.

The CXC or a-class chemokines possess a single residue between the first and second cysteines. This family can be further subdivided by the presence or absence of the amino acid sequence Glu-Leu-Arg (ELR) preceding the first [cysteine](#). This sequence is recognized by chemokine receptors on specific cell types, and thus, has functional implications. The CC or b-class chemokines have no residues between the first and second cysteines. Some members of this group have two additional cysteine residues that form a third disulfide bond (6).

Two additional chemokines have been identified that do not fit into either the CC or CXC classes. The g class is missing both the first and third cysteines. This C chemokine has lymphotactin as its prototype member. The d class possesses a CXXXC motif, and thus far is represented only by fractalkine/neurotactin. This atypical membrane-bound CXXXC chemokine has, on the C-terminal side of the b-sheet sequences, a long glycosylated mucin structure, followed by a transmembrane segment and a short C-terminal cytoplasmic **domain**. The chemokines, with their alternative names, are listed in Table 1 by class.

**Table 1. Human Chemokines and Their Receptors**

| Chemokine (and Alternative Names) | Receptors                  |
|-----------------------------------|----------------------------|
| <i>ELR + CXC</i>                  |                            |
| IL-8                              | CXCR1, R2                  |
| GROa (MGSA-a)                     | CXCR2, R1                  |
| GROb (MGSA-b, MIP-2a)             | CXCR2                      |
| GROg (MIP-2b)                     | CXCR2                      |
| ENA-78                            | CXCR2                      |
| LDGF-PBP                          | CXCR2                      |
| GCP-2                             | CXCR2                      |
| <i>Non-ELR CXC</i>                |                            |
| PF4                               | Unknown, GAGs <sup>a</sup> |
| Mig                               | CXCR3                      |
| IP-10                             | CXCR3, GAGs <sup>a</sup>   |

|                                 |                        |
|---------------------------------|------------------------|
| ITAC                            | CXCR3                  |
| SDF-1a/b                        | CXCR4                  |
| BCA-1/BLC                       | CXCR5                  |
|                                 | <i>CC</i>              |
| MIP-1a (LD78a)                  | CCR1, CCR5, CCR9       |
| MIP-1b (Act-2, HC21)            | CCR1, CCR5, CCR9       |
| MDC (STCP-1)                    | CCR4                   |
| TECK                            | Unknown                |
| TARC                            | CCR4                   |
| RANTES                          | CCR1, CCR3, CCR4, CCR5 |
| HCC-1                           | CCR9                   |
| HCC-4 (NCC-4, LEC)              | Unknown                |
| DC-CK1 (PARC, MIP-4, AMAC-1)    | Unknown                |
| MIP-3a (LARC, Exodus)           | CCR6                   |
| MIP-3b (ELC)                    | CCR7                   |
| MCP-1                           | CCR2, CCR9             |
| MCP-2                           | CCR2, CCR9             |
| MCP-3                           | CCR2, CCR9             |
| MCP-4                           | CCR2, CCR3, CCR9       |
| Eotaxin                         | CCR3, CCR9             |
| Eotaxin-2/MPIF-2                | CCR3                   |
|                                 | <i>Six CysteineCC</i>  |
| I-309                           | CCR8                   |
| MIP-5/HCC-2(Lkn-1)              | CCR1, CCR3             |
| MPIF-1 (CKb-8)                  | Unknown                |
| 6Ckine (SLC,Exodus-2, TCA-4)    | CCR7, CXCR3            |
|                                 | <i>Others</i>          |
| Lymphotactin (SCM-1)            | XCR1                   |
| CX3C (Fraktalkine, Neurotactin) | CX3CR1                 |
| Multiple Chemokines             | DARC (Duffy antigen)   |

<sup>a</sup> GAG = Glycosyl amino glycan.

Members within a class are linked both by their structural similarities and by their location in the human genome. Although there are exceptions in each class, members of the CXC group generally map to the chromosomal locus 4q13. CC chemokines are found primarily on chromosomes 7, 8, 9, and 17. Chromosomal location may provide a useful tool in identifying new members of the g and d classes.

## 2. Chemokine Expression

Chemokines are expressed by virtually all cell types. Regulation of chemokine expression is complex and highly dependent on specific signaling pathways. Cells can express chemokines constitutively or in response to stimulation by, for example, **interleukin-1** (IL-1), interleukin-4, **tumor necrosis factor** [a](#), **interferon-g** (IFN-g), or pathogenic agents such as [endotoxin](#)

(lipopolysaccharide). Low molecular-weight fragments (200 kDa) of the glycosaminoglycan hyaluronan (resulting from tissue destruction, for example) can enhance expression of some chemokines (eg, MIP-1, RANTES, or MCP-1) in specific macrophage types (eg, bone marrow-derived macrophages or elicited peritoneal macrophages) (7). Depending upon the situation, the expression of these same chemokines can be either suppressed or enhanced by inflammatory mediators, such as IL-10 or IFN-g (8).

One well-characterized example of induction of chemokines is provided by the [macrophage](#). Macrophages can be stimulated to produce chemokines during T cell-directed delayed type hypersensitive reactions in tissues by contact with activated T cells displaying CD40L. Kornbluth et al. (9) found that, although LPS was also a potent inducer of chemokines, the macrophage-stimulating [lymphokines](#) IFN-g and GM-CSF (granulocyte macrophage colony stimulating factor) produced by activated T cells were not good inducers.

### 3. Chemokine Receptors

Chemokines bind to seven-transmembrane-spanning G-protein-coupled receptors (see [Heterotrimeric G Proteins](#)) whose sequence [homology](#) is distinct from that of other such receptors. The interaction of a chemokine with its receptor can activate diverse signaling pathways, depending upon the specific signaling elements expressed by the target cell. Chemokine receptors are named on the basis of the cysteine motif that they recognize, and most will bind more than one chemokine ligand in that group. Likewise, a particular chemokine can interact with more than one receptor. Table 1 lists the chemokines and their respective receptors. Some receptors are expressed constitutively by particular cell types, whereas the expression of others may require an inducing signal. Thus CCR1 and CCR2 are constitutively expressed by monocytes, but are expressed by lymphocytes only after stimulation by IL-2 (10). Some chemokine receptors are expressed by nonhematopoietic cells, eg, endothelial cells, neurons, astrocytes, and epithelial cells. This emphasizes recent findings that chemokine receptors may have other important functions besides that of leukocyte chemotaxis. The Duffy antigen receptor, the determinant of the Duffy blood group, is found on a variety of cells, including erythrocytes and endothelial cells, as well as cells of the nervous system. It is unusual in its ability to bind both CC and CXC chemokines without eliciting any known biologic response. Possibly it serves to store or scavenge chemokines (11).

Specific chemokines generally act on specific subsets of leukocytes as a function of the receptors expressed by the cell (Table 2). Most CXC chemokines containing the ELR sequence are chemotactic for neutrophils, whereas non-ELR chemokines generally attract lymphocytes. Neutrophils express the CXCR1 receptor, which is stimulated by IL-8 and granulocyte chemotactic protein-2 (GCP-2), and the CXCR2 receptor, which responds to IL-8, GCP-2, the LPS-induced CXC chemokine (LIX), the neutrophil-activating peptide-2 (NAP-2), the epithelial cell-derived neutrophil activating peptide 78 (ENA-78), and the growth-regulated oncogenes, GRO-a, GRO-b and GRO-g. IL-2-stimulated lymphocytes express CXCR3, the receptor for MIG, IP-10, ITAC, and 6CKine.

**Table 2. Primary Expression of Chemokine Receptors<sup>a</sup>**

| Receptor (Alternative Names) | Primary Expression                               |
|------------------------------|--|
| CCR1                         | Monocytes, eosinophils, activated T cells        |
| CCR2                         | Monocytes, eosinophils, activated T cells        |
| CCR3                         | Basophils, eosinophils, subset of Th2, monocytes |
| CCR4                         | Activated T lymphocytes                          |

|                              |  |
|------------------------------|--|
| CCR5                         | Monocytes, T lymphocytes                                 |
| CCR6 (GPR-CY4, DRY6, CKR-L3) | T and B lymphocytes, bone marrow-derived dendritic cells |
| CCR7 (EBI1)                  | Activated T and B lymphocytes                            |
| CCR8 (Ter1, CemR1)           | Monocytes, T lymphocytes                                 |
| CCR9 (D6)                    | Bone marrow, T lymphocytes, monocytes                    |
| CXCR1 (IL-8RA)               | Polymorphonuclear cells (PMNs, neutrophils)              |
| CXCR2 (IL-8RB)               | Polymorphonuclear cells (PMNs, neutrophils)              |
| CXCR3                        | Activated T lymphocytes                                  |
| CXCR4 (LESTR, HMSTR, Fusin)  | Monocytes, T lymphocytes, PMNs                           |
| CXCR5 (BLR-1)                | B cells  |
| XCR1 (GPR5)                  | Unknown  |
| CX3CR1 (V28)                 | Neurons, endothelial cells                               |
| DARC (Duffy antigen)         | Red blood cells, postcapillary endothelial cells         |

<sup>a</sup> This table lists cell types that have been widely accepted as expressors of chemokine receptors. It should be noted that many other cell types express these receptors as well.

Eosinophils, basophils, monocytes, NK cells, and T cells respond to members of the family of CC chemokines. These include monocyte chemoattractant proteins (MCP-1, -2, -3, -4, and -5), macrophage inflammatory protein-1a and -1b (MIP-1a, MIP-1b), RANTES (regulated on activation, normal T expressed and secreted), eotaxin-1 and -2, and I-309. Receptors for these chemokines expressed by eosinophils are CCR1 (recognized by MCP-3 and -4; MIP-1; RANTES) and CCR3 (recognized by MCP-3 and -4; eotaxin-1 and -2; RANTES). Basophils express CCR3 and CCR2, which is activated by all the MCPs. Receptors expressed on monocytes are CCR1, 2, 5, and 8. Ligands for CCR5 include MIP-1a, MIP-1b, and RANTES, whereas the ligand for CCR8 is I-309. Activated T cells express CCR1, CCR2, CCR4 (recognized by TARC, the thymus and activation-regulated chemokine), CCR5, and CCR7, whose ligand is MIP-3 and 6Ckine (3, 4).

#### 4. Chemokine Signaling

The ability of [pertussis toxin](#) to suppress many of the chemokine-induced signals suggests that the heterotrimeric G<sub>abg</sub> proteins mediating chemokine signaling belong to the G<sub>i</sub> group, which is defined by their ability to inhibit [adenylate cyclase](#). Activation of the G<sub>i</sub> protein leads to replacement of GDP with GTP and dissociation into α<sub>i</sub>•GTP and βγ subunits, both membrane-associated.

Downstream mediators of general G<sub>abg</sub> signaling may include the *src* of nonreceptor protein **tyrosine kinases**, phosphatidylinositol 3-kinase (PI3-K), [phospholipase C](#), protein kinase C species, and the Ras-Raf-MEK-ERK kinase cascade (12, 13). Phospholipase C hydrolyzes phosphatidylinositol 4,5-bisphosphate to produce inositol trisphosphate, a mobilizer of intracellular Ca<sup>2+</sup>, and **diacylglycerol**, which, together with Ca<sup>2+</sup>, activates protein kinase C (see [Second Messengers](#) and [Calcium Signaling](#)). Activation of phospholipase D, which generates phosphatidic acid, and phospholipase A<sub>2</sub>, which releases **arachidonic acid**, may be direct or indirect downstream effects of chemokine stimulation. Immediate effects (within seconds) on cell behavior are generally

mediated by reversible changes in protein **phosphorylation** of **cytoskeletal** proteins, directed in part by the small [Gtpases](#) Rho, Rac and Cdc42, thereby affecting the polymerization and depolymerization of [actin](#), which underlies the extension and retraction of lamellipodia, implementors of leukocyte migration (3). Longer term changes are effected by **transcriptional** or **posttranscriptional** changes in gene expression.

Evidence for specific signaling intermediates includes the following: The response of neutrophils to chemoattractants can be suppressed by inhibitors of [tyrosine](#) phosphorylation, possibly acting on members of the *src* of nonreceptor protein tyrosine kinases (12). In human T lymphocytes, RANTES induces PI3-K; also, wortmannin, which is a potent PI3-K inhibitor, can inhibit RANTES-induced T-lymphocyte responses (14). In a T-cell hybrid, MCP-1 mobilized intracellular  $Ca^{2+}$  concentrations, stimulated  $Ca^{2+}$  import, and transiently increased tyrosine phosphorylation of the ERKs by a pertussis toxin-sensitive process (15). In human monocytes, Yen et al. (16) obtained evidence for a putative  $G_q$ -mediated pathway, distinct from the pertussis toxin-sensitive  $G_i$ -dependent pathway, which activated ERK2 via a protein kinase C-dependent process that probably did not require Ras. Both the  $G_i$  and the proposed  $G_q$  pathways were required for a chemotactic response to MCP-1.

MIP-1 activation of the CCR5 receptor expressed in a transfected murine pre-B lymphoma cell line led to phosphorylation and activation of a protein **tyrosine kinase**, known variously as RAFTK, Pyk2, or CAK-b, that is related to focal adhesion kinase (17). The cytoskeletal protein paxillin became phosphorylated and associated with RAFTK, and the downstream signaling kinases JNK/SAPK and p38 also appeared to be stimulated. Thus, chemokines have the ability to activate a variety of [signal transduction](#) mechanisms that connect the extracellular environment with the cytoskeleton.

## 5. Chemokine Functions

As part of the immune defense system, lymphocytes patrol the body, continually surveying it for pathogens or cells expressing abnormal surface [antigens](#) (11, 18). They move from the blood stream into tissues and then back into the blood stream via the lymphatic circulation. Granulocytes and monocytes exit the blood stream, but cannot recirculate. Extravasation is a coordinated multistep process that requires a series of complex changes in the activity of molecules on the cell's surface as it undergoes the transition from a freely floating cell to a weakly adherent rolling cell, and then to a strongly adherent stationary cell able to insinuate itself between the endothelial cells and penetrate the underlying stroma. Chemokines help drive this process and, as noted in Table 3, determine what types of cells infiltrate into the tissues in various inflammatory diseases. Imai et al. (19) report that the MDC and TARC chemokines, which are produced by dendritic cells in the thymus, act via the CCR4 receptor to recruit and activate T lymphocytes in the thymus. DC-CK1, another C-C chemokine derived from dendritic cells, is a chemoattractant for naive T cells (20). In allergic diseases such as asthma, expression of eotaxin and monocyte chemoattractant proteins appear particularly important in stimulating the accumulation and activation of eosinophils and mast cells, leading to histamine release and the allergic response (21). Similar events are thought to occur in rhinitis and atopic dermatitis.

**Table 3. Chemokine Involvement in Inflammatory Diseases**

| <b>Inflammatory Disease</b>         | <b>Infiltrate</b>     | <b>Chemokines Involved</b> |
|-------------------------------------|-----------------------|----------------------------|
| Acute respiratory distress syndrome | Neutrophils           | IL-8, GRO-a,b,g, ENA-78    |
| Asthma                              | Eosinophils, T-cells, | MCP-1,4, MIP-1a,           |

|                            |  |  |
|----------------------------|--|--|
| Bacterial pneumonia        | monocytes, basophils<br>Neutrophils          | Eotaxin, Rantes<br>IL-8, ENA-78, GRO-a,b,g |
| Sarcoidosis                | T-cells, monocytes                           | IP-10, MIP-1a, MCP-1, Rantes               |
| Glomerulonephritis         | Monocytes, T-cells, neutrophils              | MCP-1, GRO                                 |
| Rheumatoid arthritis       | Monocytes, neutrophils                       | IL-8, MCP-1, IL-8, ENA-78                  |
| Osteoarthritis             | Monocytes, neutrophils                       | MIP-1b, IL-8                               |
| Atherosclerosis            | T-cells, monocytes                           | IL-8, IP-10, Rantes, MCP-1                 |
| Inflammatory bowel disease | Monocytes, neutrophils, T-cells, eosinophils | MCP-1, IL-8, ENA-78, Eotaxin               |
| Psoriasis                  | T-cells, neutrophils                         | MCP-1, MIG, IL-8, Rantes, MIP-1            |
| Bacterial meningitis       | Neutrophils, monocytes                       | IL-8, GRO-a, MCP-1, MIP-1a,b               |
| Viral meningitis           | T-cells, monocytes                           | MCP-1, IP-10                               |
| Multiple sclerosis         | T-cells                                      | Rantes, IP-10, MIP-1a                      |
| Seasonal allergic rhinitis | Eosinophils                                  | Eotaxin, Rantes, IL-8, MIP-1, MCP-1        |
| Periodontal disease        | Monocytes, PMNs                              | IL-8, MCP-1                                |

---

Subsets of lymphocytes express different surface molecules that allow them to target different tissues in the body. Extravasation of the lymphocyte from the blood stream is initiated when a free flowing lymphocyte becomes tethered to the endothelium such that it can still roll in the direction of flow. Rolling of leukocytes is controlled in part by the interaction of selectins with specific sialylated carbohydrate determinants. L-selectin, expressed by lymphocytes, monocytes, neutrophils, and eosinophils, interacts with receptors such as CD34 on the activated endothelium. P-selectin and E-selectin molecules are expressed by activated endothelial cells and bind structures related to the tetrasaccharide sialyl Lewis antigens on the leukocyte surface. Thrombin and histamine can mobilize P-selectin from intracellular stores, whereas up-regulation of E-selectin requires induction by inflammatory agents such as IL-1, tumor necrosis factor, or LPS. This interaction appears to result in activation of leukocyte 1 and 2 integrins, enabling them to bind to ligands of the [immunoglobulin](#) superfamily (eg, VCAM-1, the ICAMs; see [Cell Adhesion Molecules](#)) either expressed constitutively by the endothelial cells or induced by inflammatory agents. Various subsets of lymphocytes rolling on an ICAM-1 monolayer can be induced to arrest within one second by SDF-1, 6Ckine, MIP-3, or MIP-3; adhesion, which can be prevented by pertussis toxin, is transient, lasting some 5 to 8 minutes (22).

In inflammatory diseases, it is likely that sentinel cells at the site of inflammation produce chemokines that induce the strong interactions required to elicit leukocyte invasion. Chemokine expression has been demonstrated in glomerulonephritis, asthma, inflammatory bowel disease, and allogenic transplant rejection (2, 4). Because multiple inflammatory mediators are generally elaborated at sites of inflammation, the task of determining which is responsible for any particular pathogenic response is difficult. Although there is considerable redundancy built into the system, some chemokines predominate in some forms of injury. For example, IL-8 is a major causative

factor associated with reperfusion injury and acid-induced pneumonitis, which is the consequence of neutrophil invasion and the release of inflammatory mediators (23).

Angiogenesis (see [Angiogenin](#)), the induction of synthesis of new blood vessels, typically with reference to the vascularization of a tumor, is a complex response affected by many factors, including chemokines. Although the specific role of chemokines in angiogenesis is unclear, evidence has suggested that CXC chemokines possessing the ELR motif generally promote angiogenesis, whereas those lacking the ELR motif tend to inhibit it (24). Mice lacking either SDF-1 or the counterpart receptor CXCR4 are defective in formation (vasculogenesis) of the large vessels supplying the intestinal tract (25). Vessels supplying the small intestine were missing major branches, and this led to hemorrhaging. Likewise, the vascular supply to the stomach in these mice was bereft of all large arteries and veins.

Many functions regulated by chemokines in the immune system, including cell migration and adhesion, as well as cytoskeletal reorganization, are important in the development and function of the nervous system. Defects in brain development in the CXCR4-deficient mouse (and its counterpart, the SDF-1 deficient mouse) are likely to be harbingers of other developmental and brain-specific functions of chemokines (26). The CXCR4-deficient mouse is embryonic-lethal, although with some animals surviving to term. The most prominent defects are found in the development of the cerebellum, where the granule neurons migrate prematurely. Since the radial glial cells are preserved in these animals, the evidence is consistent with an inhibitory role of SDF-1 on granule cell migration. This inhibition could occur by blocking the action of a migratory-promoting substance, or by altering adhesion molecules on the cell surface of the radial glial cell or the granule neuron. Since some of the guidance cues for this migratory process have been identified, it is of obvious importance to investigate if SDF-1 has any action to modulate the expression or conformation of these proteins.

Additional evidence for central nervous system actions of chemokines has been provided by the **transgenic** (see [Transgenic Technology](#)) expression of chemokines under the influence of the myelin basic protein **promoter**. Directed expression of the murine chemokine KC stimulated neutrophil infiltration of the brain at various sites (27). Likewise, overexpression of MCP-1 in both thymus and brain led to infiltration of either MAC-1, F4/80, or macrophages and monocytes, respectively (28). These studies demonstrate that each organ is capable of recruiting mononuclear cells under specific conditions.

## 6. Viral Adaptation to Chemokines

**Viruses** often exploit endogenous cell surface receptors, including chemokine receptors, in the initial stages of infection. For example, mice deficient in MIP-1 as the result of a targeted gene disruption (**knockout** mice) exhibit attenuated responses to certain viral infections (Coxsackie virus, [influenza virus](#)). Several of the chemokine receptors serve as entry mechanisms for various strains of [HIV](#). Roger (29) has reviewed the role of chemokine receptors and the HLA **haplotype** in the development of AIDS. The CCR5 receptor appears to be the preferred coreceptor for macrophage-tropic, nonsyncytium-forming HIV-1 strains, whereas the CXCR4 receptor (fusin) primarily services the syncytium-inducing, T-cell-line-tropic HIV-1 strains (30). Some HIV isolates can utilize both receptors, while a few other variants appear to have acquired the ability to co-opt other chemokine receptors as coreceptors for virus infection. Individuals with mutations in chemokine receptors may exhibit considerable resistance to infection or to progression of the disease. One mutation in particular should be noted. People homozygous for a deletion mutation in the CCR5 gene (CCR5 d32) are resistant to infection by the virus; heterozygotes for the same mutation had a significantly slower onset of disease (31). Chemokines can also affect the progression of an HIV infection. SDF-1, a ligand for CXCR4, inhibits HIV entry into cells (32, 33), and the CC chemokines MIP-1a, MIP-1b, and RANTES, all ligands for CCR5, likewise seem to retard HIV-1 infection (34).

## 7. Conclusion



The chemokine field is evolving rapidly. Although the initial focus has been on the role of chemokines in the immune system, future directions for chemokine research lie in elucidating chemokine expression and function in other organs and during [development](#). Importantly, the involvement of chemokines and their receptors in many disease processes makes them primary targets for drug development.

## 8. Acknowledgments

Research in the authors' laboratories has been supported by the National Institutes of Health.

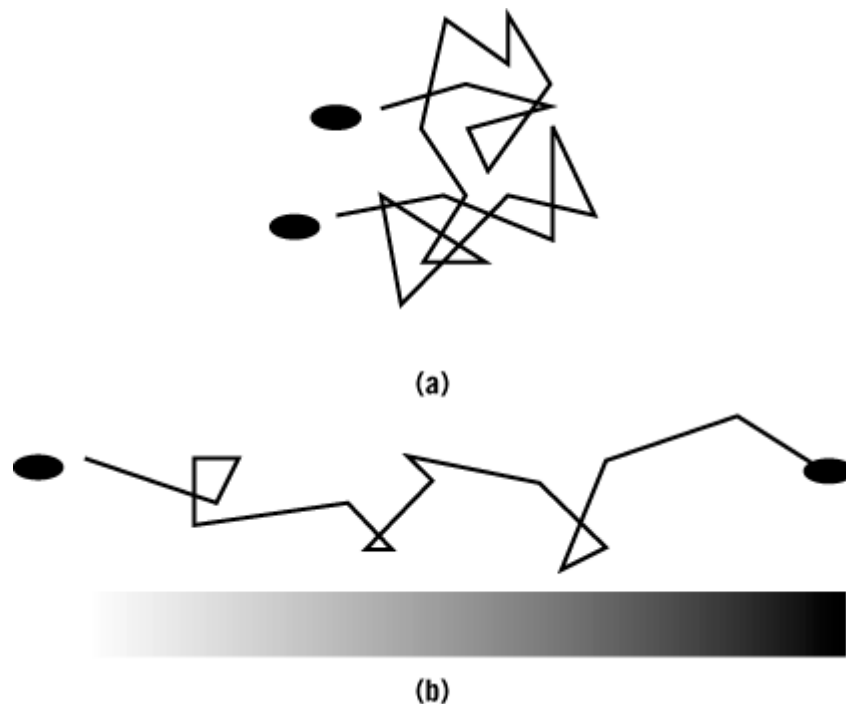
## Bibliography

1. N. W. Schluger and W. N. Rom (1997) *Current Opin. Immunol.* **9**, 504–508.
2. B. J. Rollins (1997) *Blood* **90**, 909–928.
3. M. Baggiolini (1998) *Nature* **392**, 565–568.
4. A. D. Luster (1998) *New England J. Med.* **338**, 436–445.
5. A. Zlotnick, J. Morales, and J. A. Hedrick (1998) *Adv. Immunol.* in press.
6. A. Pardigol et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 6308–6313.
7. C. McKee et al. (1996) *J. Clin. Invest.* **98**, 2403–2413.
8. M. R. Horton et al. (1998) *J. Immunol.* **160**, 3023–3030.
9. R. S. Kornbluth, K. Kee, and D. D. Richman (1998) *Proc. Natl. Acad. Sci. USA* **95**, 5205–5210.
10. P. Loetscher, M. Seitz, M. Baggiolini, and B. Moser (1996) *J. Exp. Med.* **184**, 569–577.
11. M. B. Furie and G. J. Randolph (1995) *Amer. J. Pathol.* **146**, 1287–1301.
12. G. M. Bokoch (1995) *Blood* **86**, 1649–1660.
13. J. S. Gutkind (1998) *J. Biol. Chem.* **273**, 1839–1842.
14. L. Turner, S. G. Ward, and J. Westwick (1995) *J. Immunol.* **155**, 2437–2444.
15. P. M. Dubois et al. (1996) *J. Immunol.* **156**, 1356–1361.
16. H.-h. Yen, Y. Zhang, S. Penfold, and B. J. Rollins (1997) *J. Leukoc. Biol.* **61**, 529–532.
17. R. K. Ganju et al. (1998) *Blood* **91**, 791–797.
18. T. A. Springer (1994) *Cell* **76**, 301–314.
19. T. Imai et al. (1998) *J. Biol. Chem.* **273**, 1764–1768.
20. G. J. Adema et al. (1997) *Nature* **387**, 713–717.
21. B. Lamkhieoued et al. (1997) *J. Immunol.* **159**, 4593–4601.
22. J. J. Campbell et al. (1998) *Science* **279**, 381–384.
23. H. G. Folkesson, et al. (1995) *J. Clin. Invest.* **96**, 107–116.
24. R. M. Strieter et al. (1995) *J. Leukocyte Biol.* **57**, 752–762.
25. K. Tachinaba et al. (1998) *Nature* **393**, 591–594.
26. Y-R. Zou et al. (1998) *Nature* **393**, 595–599.
27. M. E. Fuentes et al. (1995) *J. Immunol.* **155**, 5769–5776.
28. M. Tani et al. (1996) *J. Clin. Invest.* **98**, 529–539.
29. M. Roger (1998) *FASEB J.* **12**, 625–632.
30. Y. Feng, C. C. Broder, P. E. Kennedy, and E. A. Berger (1996) *Science* **272**, 872–877.
31. M. W. Smith et al. (1997) *Science* **277**, 959–965.
32. C. C. Bleul (1996) *Nature* **382**, 829–833.
33. E. Oberlin et al. (1996) *Nature* **382**, 833–835.
34. D. Zagury et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 3857–3861.

## Chemotaxis

Chemotaxis is the movement of a motile cell in response to chemical changes in the environment, and it implies a directional sense. In most cases where **bacteria** migrate up concentration gradients of attractants, direction is not sensed directly. The small size of bacteria would make spatial comparisons of concentration unreliable. Instead, bacteria usually sense temporal changes in concentration by using chemoreceptors, **transducer proteins** that bind or release attractants and repellents. This temporal measurement allows the cells to alter motion of their **flagella** briefly so that a “biased random walk” brings the bacteria to the right destination (1, 2) (Fig. 1). Eukaryotic cells carry out chemotaxis by an unrelated mechanism, and they are much larger and probably do sense direction.

**Figure 1.** Motion of peritrichous bacteria. (a) Random walk in an isotropic medium. (b) Biased random walk in a gradient of an attractant or repellent.



### 1. Bacteria

There are three main types of motile bacteria: 1) the peritrichous (flagella over the surface), 2) the polar (flagella at one or both poles), and 3) the gliding bacteria, which lack flagella.

The peritrichous bacteria, such as *Escherichia coli* and *Bacillus subtilis*, normally swim smoothly by rotating their flagella counterclockwise (looking along the flagella toward the cell body) and tumble by rotating their flagella clockwise, which discoordinates the bundle of flagella. Chemotaxis in *E. coli* occurs when the cell happens to move toward higher attractant concentrations. It decreases the probability of tumbling and, thus, increases the tendency to travel up the attractant gradient (1).

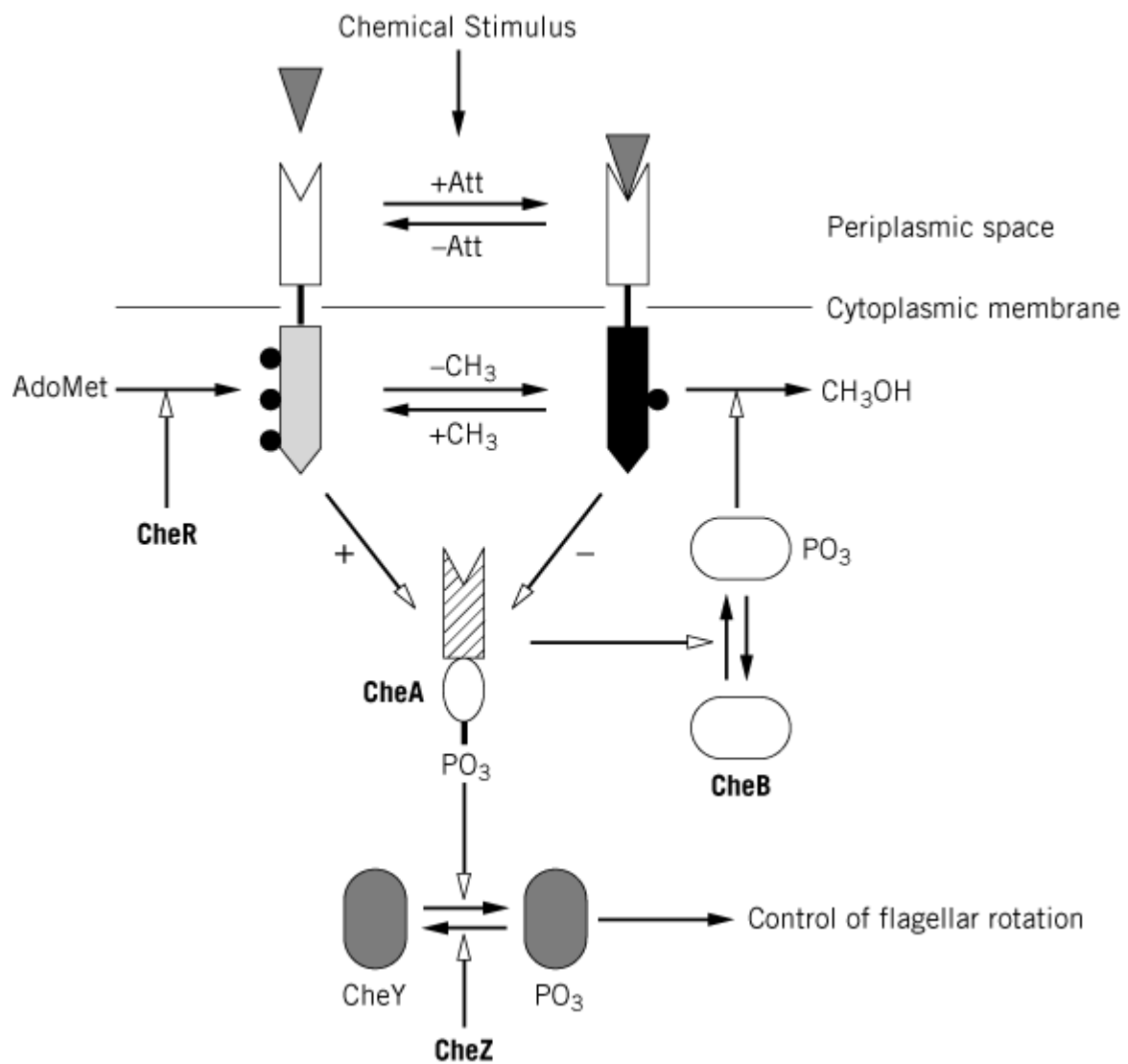
This decreased probability of tumbling results (Fig. 2) because the attractant molecules bind to methylated chemoreceptors to cause a brief inactivation of the cytoplasmic CheA **kinase** (see **Transducer proteins**). Consequently, there is less phosphorylated CheA (CheA-P) and less CheY-P (2) (see **Phosphorylation**). CheY-P dissociates from the switch that controls flagellar rotation, so that smooth swimming ensues. The enhanced smooth swimming is short-lived because CheZ hydrolyzes the phosphoryl group from CheY, and the CheR **methyltransferase** increases the degree of methylation of the chemoreceptor and restores the prestimulus activity of CheA (3) (see also Methylation, Protein). Thus, there is an inverse relation between the levels of CheY-P and termination of a smooth swimming event.

Logically, it seems plausible that bacteria moving by chance toward *lower* attractant concentrations would show an increased tendency to tumble and, hence, to begin heading in another direction. In the case of *E. coli*, however, there is no demonstrable increased tendency for these bacteria to tumble compared with bacteria in an isotropic medium (1). Interestingly, CheY-P cannot affect rotational direction of the flagella by itself. Fumarate needs to be released or synthesized for switching to occur (4) and also promotes tumbling (5). The connection between CheY-P and fumarate in controlling direction of flagellar rotation has not been elucidated.

The role of CheY-P is less clear in the case of polar bacteria, such as *Caulobacter crescentus* and *Pseudomonas aeruginosa*. When these bacteria move by chance toward lower concentrations of attractant, they reverse motion by rotating their flagella in the opposite direction, to head back up the gradient. If repellent is added, however, they show a series of reversals of motion. This result implies that CheY-P does not govern the *direction* of rotation of the flagella but instead causes repeated reversals of behavior. Such an event might ensue if binding of CheY-P to the switch reduces the energy of activation of switching the flagellar rotational direction. In the absence of binding, the energy of activation is presumed to be so high that reversals are very infrequent. In the case of *Halobacterium salinarium*, negative stimuli cause significant increases in the concentrations of free fumarate in the cytoplasm (4).

Other bacteria, like *Sinorhizobium meliloti*, a **Gram-negative** bacterium phylogenetically rather distant from *E. coli*, rotate their flagella unidirectionally. In this case, the CheA/CheY system also exists with the complication that there are two species of CheY. CheY2-P is the major species that affects flagellar movement, and it acts by slowing flagellar rotation (6). Some bacteria, such as *Rhodobacter sphaeroides*, have multiple copies of many of the *che* genes (7).

**Figure 2.** Diagram of the mechanism of chemotaxis in *Escherichia coli*. Open arrows indicate control of the designated reaction or signaling step. At the *top*, the attractant (Att) binds to its transmembrane chemoreceptor, which is methylated on Glu residues. The methyl groups are depicted as closed black circles. This binding decreases the phosphorylated form of the protein CheA (*center*), which decreases phosphorylation of the proteins CheY and CheB. Dephosphorylated CheB is much less active in removing the methyl groups from the chemoreceptor, so CheR restore them. Phosphorylated CheY dissociates from the switch that controls flagellar rotation, and smooth swimming ensues. CheZ hydrolyzes the phosphoryl group from CheY, ending the response to the attractant.



Thus, it appears that the basic system involving the two-component couple CheA and CheY is universal, but how it functions varies for different groups of organisms.

The potential role for CheY-P is still less clear in the case of the gliding bacteria. One of them, *Myxococcus xanthus*, has two motility systems, the S (social) system powered by extrusion and retraction of pili from the ends and the A (adventurous) system, powered by a mechanism of unknown structure along the sides. The former coordinates cells and, unless the agar concentration is quite low, requires cell contact for movement. The frequency and pole at which pili activity occurs is governed by the Frz proteins, which are homologous to the chemotaxis proteins in flagellated cells. However, none of the genes encode a single-domain CheY (8, 9).

One issue that has existed for many years has been whether the rotation of multiple flagella are coordinated or independent. The finding that bundles of flagella located 50 mm apart on *Spirillum volutans* were observed switching synchronously within 10 ms implies electrical coordination of direction of the flagellar rotation (1). In contrast, certain mutant strains of *Salmonella typhimurium* show random and uncorrelated switching of their flagellar direction when observed under partially de-energized conditions (10). *E. coli* filamentous cells switch their flagella asynchronously although biases of nearby, but not distant, motors were correlated (11). In *H. salinarium*, an archaeon, the rotational direction of flagellar bundles on each end is also uncoordinated, and cells were frequently observed to be immobilized because the bundles at both ends worked against each other. Thus, there

is no general rule governing flagellar coordination. Indeed, the mechanism by which an electrical connection could work, as in *S. volutans*, is obscure.

## 2. Eukaryotes

Chemotaxis in eukaryotes occurs by an unrelated mechanism. The process is best known from studies in the aggregating slime mold *Dictyostelium discoideum* and in leucocytes; it is probably quite similar in both. In the slime mold, chemotaxis plays a role in feeding and in the aggregation that precedes **differentiation** and **sporulation**. Chemotaxis in *Dictyostelium* involves apparent polarization of the cell to form pseudopodia or lamellipodia at the leading edge at the highest concentration of attractant molecules. Polarization and formation of pseudopodia are concerted events regulated by G protein-coupled/serpentine **cyclic AMP** (cAMP) receptors (cARs) along the surface of the organism. During aggregation, cAMP is produced and binds to a cAR, which activates a heterotrimeric **G-protein**. **Guanylate cyclase** is activated by operation of this G-protein, thereby producing cGMP. cGMP activates a cGMP-dependent protein kinase, which regulates myosin II heavy and light chain kinases and may also help regulate the actin cytoskeleton. Phosphorylation of the myosin heavy chain causes it to depolymerize so that pseudopod extension can only occur locally. Posteriorly, contraction occurs and pseudopod formation is blocked. At the same time cAMP binding to its receptor leads to activation of phosphatidylinositol-3 kinase. The product, phosphoinositol-3,4,5-triphosphate, attracts PH (pleckstrin homology) proteins, including the kinase Akt/PKB, to the localized region where the local amount of cAMP-bound receptor is highest. Another PH domain protein is CRAC. It is required for receptor activation of adenylyl cyclase (see below), which then propagates the cAMP signal to neighboring cells. This membrane localization lasts only 5 to 8 s; however, that is long enough to produce an “activation domain” as a focus for multiple pathways needed for chemotaxis, pseudopod extension, and cell polarization. F-actin is assembled primarily there. Actin polymerization is activated by formation of an Arp2/3 protein complex. Nucleation is induced by Arp2/3 interaction with WASp and Scar adaptor proteins. These proteins provide the structural linkage to the activated cAR through the activated form of Rac and Cdc42, which are regulated by the G protein G<sub>bg</sub>. Pseudopod formation also requires actin crosslinking and filament growth regulated by actin binding proteins ([12](#), [13](#)).

Binding of cAMP to the receptors transiently activates both adenylyl cyclase and ERK2 MAP kinase. A rise in internal cAMP activates protein kinase A such that it inhibits ERK2 and leads to loss of ligand binding by the receptor. ERK2 also phosphorylates the cAMP diesterase REG A that reduces the internal concentration of cAMP. A secreted membrane-associated phosphodiesterase reduces the external cAMP concentrations, and the cells become sensitive again. These features account for spontaneous oscillations of external cAMP production of about 6 min that permit chemotactic migration of cells toward centers during the aggregation process ([14](#)).

## Bibliography

1. H. C. Berg (1975) *Annu. Rev. Biophys. Bioeng.* **4**, 119–136.
2. J. B. Stock and M. G. Surette (1996) In *Escherichia coli and Salmonella, Cellular and Molecular Biology*, 2 ed. (F. C. Neidhardt, ed.), ASM Press, Washington, DC, Vol. **I**, pp. 1103–1129.
3. J. S. Parkinson (1993) *Cell* **73**, 857–871.
4. R. Barak and M. Eisenbach (1996) *Curr. Top. Cell. Reg.* **34**, 137–158.
5. K. Prasad, S. R. Caplan, and M. Eisenbach (1998) *J. Mol. Biol.* **280**, 821–828.
6. V. Sourjik and R. Schmitt (1996) *Mol. Microbiol.* **22**, 427–436.
7. J. P. Armitage *Adv. Microbiol. Physiol.* **41**, 229–289.
8. H. Sun, D. R. Zusman, and W. Shi (2000) *Curr. Biol.* **10**, 1143–1146.
9. M. J. Ward and D. R. Zusman (1997) *Mol. Microbiol.* **24**, 885–893.
10. R. M. Macnab and D. P. Han (1983) *Cell* **32**, 109–117.

11. A. Ishihara, J. E. Segall, S. M. Block, and H. C. Berg (1983) *J. Bacteriol.* **155**, 228–237.
12. R. A. Firtel and R. Meili (2000) *Curr. Opin. Genet. Develop.* **10**, 421–427.
13. R. A. Firtel and C. Y. Chung (2000) *BioEssays* **22**, 603–615.
14. M. T. Laub and W. F. Loomis (1998) *Mol. Biol. Cell* **12**, 3521–3532.

### Suggestions for Further Reading

15. J. J. Falke and S. H. Kim (2000) Structure of a conserved receptor domain that regulates kinase activity: the cytoplasmic domain of bacterial taxis receptors. *Curr. Opin. Struct. Biol.* **10**, 462–469.

## Chimera

Spemann (1, 2) was the first to employ the term “chimera” and to consider the great potential for surgically created chimeric embryos in the analysis of **developmental** mechanisms. The chimera method has frequently involved imaginative experimental procedures by which cells of one species are grafted into another. Any animal thus composed of different cell populations that derive from more than one fertilized egg should be considered as a chimera. This type of animal can currently be constructed in amphibians, birds, and mammals.

Many fundamental concepts of embryology have been at least partly formulated on the basis of results of cell or tissue transplants between two different embryos, usually separate species of amphibians. The most spectacular transplantation experiments, published by Spemann and Mangold in 1924 (3), demonstrated the organizing power of the dorsal lip of the blastopore during gastrulation by interspecific transplantations of this area. More recently, a model of lens induction was developed by using chimeric eyes (4). The technical advantages of producing amphibian chimeras are straightforward, owing to the independence of the embryos from their parents. They are easily accessible and receptive to foreign tissue, even across species barriers. All these qualities are not shared by higher vertebrates, such as birds and mammals. Nonetheless, bird embryos have several advantages over other vertebrate embryos, making certain interesting approaches feasible. The greatest advantage is continual accessibility within the egg throughout the developmental period. Another is the ease with which the various rudiments can be delineated, and thus removed and replaced, with extreme precision. An avian chimera obtained by combining quail and chick cells has been the most successful method, having provided a continual source of new data about developmental mechanisms for almost 30 years (5, 6). With the advent of the quail-chick nuclear marker, which is particularly simple to employ, easy to identify, and endowed with great resolving power, avian chimeras have been used to study the ontogeny of the nervous system, the development of the hematopoietic and immune systems, and the formation of muscles and skeleton. Quail **heterochromatin** in the nucleus is concentrated around the **nucleolus**. This creates a large, deeply staining mass that is easily distinguishable from the diffuse heterochromatin of chick cells. Moreover, there are some **antigens** that are quail-specific and not detectable in chick cells. These phenomena allow individual quail cells to be readily distinguished, even when most of the cell population is chick.

Although the avian **embryo** is a practical model, perfectly suitable for tissue graft experiments after the incubated egg is opened, it is difficult to undertake this type of investigation in the mammalian fetus *in utero*. Nonetheless, it has become routine to remove postimplanted mammalian embryos from the uterus, manipulate them, and return them to a foster mother for further development. Thus,

chimeric mice are the result of two or more early-cleavage (usually 4- or 8-cell) embryos that have been artificially aggregated to form a composite embryo. Since each cell is able to produce any component of the body, the construction of the chimeric mouse has very important consequences for the study of mammalian ontogeny. A very powerful application of this technique is the transfer of genes into every cell of the mouse embryo. During mouse development, there is a stage when only two cell types are present: outer cells, which will form the fetal portion of the placenta, and inner cells, which will give rise to the embryo itself. These inner cells are known as embryonic [stem cells](#) because each in isolation can generate all the cells of the embryo (7, 8). These cells can be grown in culture, where they are treated to incorporate new DNA. The new embryonic stem cells can then be injected into another early-stage mouse embryo, resulting in a chimeric mouse. Mice that derive from these animals are **transgenic** mice. Our understanding of regulatory mechanisms in mammalian development is improving increasingly rapidly as a result of the construction of these genetically modified mice. A combination of the tools of developmental genetics with those of embryology should lead to real advances in the study of such mechanisms. In this field, our group has pioneered the grafting of embryonic tissues from transgenic mice into the chick embryo (9, 10). Owing to these interspecies grafting experiments, it is possible to monitor factors that regulate the expression of a particular gene *in vivo*. The value of this technique is greatly increased when a [reporter gene](#) is used to follow the changes in gene expression of the grafted cells. The possibility of conducting grafts until late stages of *in vivo* development allows the behavior of wild and mutant mouse cells to be observed at any developmental stage and location.

## 1. Amphibian Chimera

Using amphibians, Spemann and Mangold (3) improved our understanding of the specification of the nervous system by transplanting dorsal blastopore lip tissue from an early gastrula into the ventral ectoderm of another gastrula. They used differently pigmented embryos from two species of newt: darkly pigmented *Triturus taeniatus* and nonpigmented *Triturus cristatus*. On the basis of color, it was easy to distinguish host and donor tissues. The dorsal blastopore lip tissue from early *Triturus taeniatus* gastrula, once transplanted into an early *Triturus cristatus* gastrula in the region, would normally become ventral epidermis. In fact, the donor tissue did not become belly skin but invaginated and formed a secondary embryo, face to face with its host. The more recent use of nuclear markers has allowed Spemann's results to be confirmed (11). Such chimeras elegantly demonstrate the organizing power of the dorsal lip of the blastopore in amphibian gastrula, since whole secondary embryos formed under the influence of the transplanted tissue.

Considerable advances have also taken place in the field of differentiation and organogenesis through the use of amphibian chimeras. The more common examples are those involving the interaction of epithelia with adjacent mesenchyme. After being separated, embryonic epithelium and mesenchyme can be recombined in different ways (12). In a classic experiment, Spemann and Schotté (13) transplanted flank ectoderm from an early frog gastrula into the region of newt gastrula destined to become part of the mouth. Similarly, the presumptive flank ectodermal tissue of newt gastrula was placed into the presumptive oral regions of frog embryos. The structures of the mouth region differ greatly between salamander and frog larvae. The *Triturus* salamander larva has club-shaped balancers beneath its mouth, whereas frog tadpoles produce mucus-secreting glands and suckers. Frog tadpoles also have a horny jaw without teeth, whereas the salamander has a set of calcareous teeth in its jaw. The larvae resulting from the transplants were chimeras. The salamander larvae had froglike mouths, and the frog tadpoles had salamander teeth and balancers. In other words, the mesodermal cells instructed the ectoderm to make a mouth, but the ectoderm responded by making the only mouth it “knew” how to make. Thus, instructions sent by mesenchymal tissue can cross species barriers, although the response of the epithelium is species-specific. Thus, organ-type specificity is usually controlled by the mesenchyme within a species, but species specificity is usually controlled by the responding epithelium.

The cells that form the lens are derived from a region of head ectoderm in contact with optic vesicles of the anterior neural plate. Servetnick and Grainger (14) removed animal cap ectoderm from various

gastrula stages and then transplanted them into the presumptive lens region of neural-plate-stage embryos. Ectoderm from early gastrula showed little or no competence to form lenses, but ectoderm from slightly later stages was able to. By the end of gastrulation, this ability to respond to the neural plate signal had been lost. This competence was found to be inherent within the ectoderm itself and not induced by other surrounding tissues. These observations showed that only mid- to late-gastrula ectoderm is able to respond to signals from the anterior neural plate. It has recently been demonstrated that the transcription factor Pax6 may play a role in the determination processes of eye tissue (see [Pax Genes](#)).

## 2. Avian Chimera

When portions of quail embryo are grafted into a similar region of chick embryo, the cells become integrated into the host and participate in the construction of the appropriate organs. This grafting is done while the embryo is still inside the egg, and the chick that hatches is a “chimera,” having a portion of its body composed of quail cells (15). The quail is usually chosen as a donor because it is easier to identify quail cells among chick cells than the reverse.

The ontogeny of the peripheral nervous system is one of the fields in which the use of avian chimeras is the most fruitful. This system arises almost entirely from the [neural crest](#), a transient structure that develops at the neural tube apex. To follow the fate of chick neural tube-bearing neural crests, a segment is replaced by a homologous quail fragment. This method has permitted the diverse range of neural crest potentialities to be well-described, and the list of neural crest derivatives well-defined (6).

In relation to its origin, neural crest not only participates in the formation of spinal and cranial sensory ganglia, sympathetic and parasympathetic ganglia and plexuses, and Schwann cells of the peripheral nerves, but it also gives rise to endocrine and paraendocrine cells (calcitonin-producing cells, carotid body type-I cells, adrenomedullary cells) and pigment cells. The construction of an avian chimera has demonstrated the contribution of the cephalic neural crest to mesectodermal derivatives. For example, nasal and maxillary processes are built up partly by crest cells of mesencephalic origin. The mesencephalic crest cells also form the cephalic skeleton, the upper and lower jaws, the palate, and the tongue. The rhombencephalic crest participates in formation of the pre-optic region and the hyoid arch skeleton. In addition to cartilage and bone, the cephalic neural crest takes part in other derivatives of the head and neck, such as dermis and connective tissue.

Brain development has also been characterized by the use of the quail-chick chimera technique (16). Quail-chick brain chimeras can hatch and survive without showing impaired movement or locomotion, which indicates that functional synapses have been established between host and donor neurones, as well as between donor neurones and host muscles. Moreover, the normal behavior of the chimeras demonstrates that proper neuronal connections develop in the brain, which means that quail axons recognize local signals for growth and directionality in the chick environment, as they do in normal development. An interesting application of the chimeras is illustrated by chicks with transplanted quail mesencephalic-diencephalic brain areas, which then exhibit quail vocalization (17). Another important issue is the elucidation of when and where groups of cells in the brain make commitments to particular development pathways. Genes with spatially restricted expression are of particular interest, since they may indicate the existence of committed groups of cells that are important in pattern formation, but are not discernible on the basis of morphology. An example is provided by the **zinc-finger** gene *Krox-20*, which has restricted domains of expression in the early neural plate. *Krox-20* is expressed in the early neural epithelium, first in one stripe, and then in two stripes in the hindbrain. Subsequently, *Krox-20* is expressed in two alternating segments (rhombomeres) in the hindbrain. The generation of stripes of *Krox-20* expression in the early neural plate suggests that rhombomeric precursor cells are committed prior to the morphological appearance of the rhombomere.

Hox [homeotic gene](#) expression is seen along the dorsal axis, from the anterior boundary through to



the tail. Direct evidence for hindbrain plasticity comes from quail-chick chimera experiments showing that anterior-to-posterior transpositions can reprogram Hox expression and induce a transformation in cell fate (18). The quail-chick system has also shown that environmental cues play a significant role in maintaining the Hox code in the neuroepithelium.

Chimeric experiments involving the muscles of the body, limbs, and skeleton have also been carried out in birds. For development, somites are the primitive metamereric structures of the vertebrate body from which arise the vertebrae that surround the spinal cord, the muscles and connective tissue holding the vertebrae, the dermal layer of the skin of the back, and the limb musculature. All these data were obtained by constructing a chimeric embryo after interspecific grafts of somites between quail and chick embryos (19, 20). The somites appear as pairs of epithelial spheres that bud off from the unsegmented paraxial mesoderm in a craniocaudal direction. They become polarized into a ventral mesenchymal compartment, the sclerotome, which yields the dorsal skeleton, and a dorsal epithelial component, the dermomyotome, from which striated muscle and dermis arise. By constructing a quail-chick chimera after interspecific exchanges of medial or lateral halves of newly-formed somites, Ordhal and Le Douarin (21) determined that, regardless of the somites, cells farthest from the neural tube (ie, lateral) migrate to form the body wall and limb musculature.

The interspecies graft method in birds has contributed considerably to understanding the control mechanisms of somite patterning. This method demonstrates that specification of the somite is accomplished by the interaction of different tissues that form its environment. In fact, the newly formed somite is composed mostly of unspecified cells, and the determination of somite compartments toward the different lineages is regulated by environmental cues. The ventral-medial portion of the somite is induced to become the sclerotome by factors, especially *Sonic hedgehog* protein, that are secreted from the notochord and neural tube floor plate. If portions of the notochord, the source of **Sonic hedgehog**, are transplanted next to other regions of the somite, those regions also become sclerotome cells. These sclerotome cells express a new transcription factor, Pax1, which activates cartilage-specific genes and is necessary for the formation of vertebrae (22). In similar ways, the myotome is induced by two distinct signals. The muscles surrounding the body axis, which arise from the medial portion of the somite, are induced by factors from the dorsal neural tube, probably members of the *Wnt* family (23, 24). The muscles derived from the lateral portion of the somite, which form the musculature of the limbs and body wall, are probably induced through the combination of Wnt proteins from the epidermis with bone morphogenetic protein-4 (BMP4) protein from the lateral plate mesoderm (25). These factors cause the myotome cells to express particular transcription factors, MyoD and Myf5, that activate muscle-specific genes. The dermatome differentiates in response to another factor secreted by the neural tube, neurotrophin 3 (NT-3) (26).

### 3. Mammalian Chimera

The laboratory mouse is by far the most popular mammal used to construct chimeras for developmental studies. Mice offer a large variety of genetically well-characterized strains that provide scope for using genetic markers. Moreover, they breed throughout the year and thus continually supply embryos. The most convenient time for manipulation of a mammalian embryo is before its implantation. During the first three to four days of gestation, mouse embryos have not yet formed an attachment in the mother's reproductive tract and can thus be explanted, manipulated, and then retransplanted into another mouse to continue their development. Chimeric mice are generally formed by the fusion of two early embryos that organize to produce a single mouse with two distinct cell populations.

A skeletal muscle cell (myotube) is an elongated cell containing many nuclei. It has been widely debated whether this cell is derived from the fusion of several mononucleated precursor cells (myoblasts) or from a single cell that undergoes nuclear division without cell division. Evidence for the fusion process has been provided by mouse chimeric constructions. Mintz and Baker (27) fused mouse embryos from two strains that produce different types of the dimeric enzyme isocitrate dehydrogenase: one strain makes the A subunit and the other B. If myotubes are formed from one

cell whose nuclei divide without cytokinesis, the dimeric enzyme will be purely AA or BB. If, however, myotubes are formed by fusion between cells, some might code for the B subunit and others for A, in which case the molecules of the enzyme will be hybrid (AB). [Electrophoresis](#) can separate these three types (see [Isozyme, Isoenzyme](#)). The presence of the hybrid AB enzyme in extracts of skeletal muscle tissue confirms the fusion model.

The study of mutations that impair development or function has long been recognized as a valuable means of elucidating the normal role of genes in such processes. Genetically chimeric animals consisting of mixtures of mutant and wild type cells can be of considerable value for studying the genes that are primarily implicated in developmental mechanisms. Methods of manipulating mouse embryos and transferring genes into every embryonic cell are now standard procedures. In recent years, an extraordinary increase in the use of these methods has led to numerous exciting prospects.

The most widely used method to produce transgenic mice is to inject **cloned** DNA into a pronucleus of a fertilized egg. Although this method has been used primarily for studies of **gene expression**, an unexpected benefit is the frequent integration of the injected DNA into genes, causing insertional mutations. Experimental infection of embryos by [retroviruses](#) has also been developed as a method for gene transfer. As in the case of spontaneous infection, this approach has resulted in insertional mutations. An additional approach to producing transgenic mice involves genetic modification of tissue culture lines of embryonic stem cells, followed by their re-incorporation into growing embryos. By applying somatic cell genetic techniques to these cells while in culture, both random and selected insertional mutations have been produced. Another method for assaying the function of a gene is to eliminate it from the genome of the whole organism. Very recent studies, using strategies to inactivate genes by homologous [recombination](#) in embryonic stem cells, with the subsequent generation of germ-line chimeras, have made this scenario possible and thereby initiated a new era in mammalian developmental biology.

The power of this technology can be illustrated by the *int-1* gene, which is expressed during central nervous system development in a temporally and spatially restricted fashion. Following inactivation of the gene by homologous recombination in embryonic stem cells, germ-line chimeras were obtained and bred to derive homozygous *int-1*-deficient mice. These homozygous mice were not viable, because specific regions of the brain were missing. This experiment indicates that *int-1* is an essential gene for early development and suggests that its expression might determine the fate of a cell. The *Hox a-3* gene has been found to control segment-specific gene expression in *Drosophila*, but what is its role in mammals? Chisaka and Capecchi (28) used [gene targeting](#) to determine the function of *Hox a-3* in the development of the mouse. Homozygous mutants of *Hox a-3* were found to have severe anomalies in development of the hindbrain and in mesectodermal neural crest derivatives and to lack thyroid, parathyroid, and thymus glands. These studies clearly open up new avenues for the study of mammalian development at the molecular level and will certainly be instrumental in unraveling the molecular networks responsible for the functioning of a multicellular organism.

#### 4. Mammalian Avian Chimera

In chimeras, the analysis of cell lineages depends on the ability to distinguish between cells of different origins. Markers are required to trace the developmental fate of cells and to recognize them at all ages in chimeric embryos. In mammals, genetic manipulations appear to be valuable means of elucidating the normal role of genes and identifying altered cells. Furthermore, the environment of the chimera allows these cells to survive and to display their phenotype. In chimeras, identification of all the descendants of genetically modified cells has thus far provided a reliable and sensitive means of measuring alterations in embryogenesis. In conjunction with an *in situ* marker, this may help in determining the cell type and the nature of the functions affected by the particular mutant. Extremely sensitive labeling is possible with transgenic mouse lines since cells that integrate foreign DNA, such as a reporter gene encoding for *Escherichia coli*  $\beta$ -galactosidase, are detectable by simple histochemical revelation, which can be used to differentiate grafted cells from host cells. Currently,

the major means of exploring this genetic tool is by *in vitro* experiments in which mutant-induced tissue is cultured with wild-type tissue. The culture system is efficient in tissue development for only limited periods, however, and the media used may affect the tissue outcome. *In vivo* micromanipulation remains the most powerful means of studying the fates of cells, their origins, and the cell-cell interactions that govern development. Since implantation of mouse embryo is particularly unsuitable for *in vivo* manipulation, chick embryo is used as the site for developing mouse embryonic cells. Studies in such chimeras provide considerable information on the lineage and the differentiation of mouse embryonic cells, since *in ovo* grafted mouse cells are accessible throughout chick host embryogenesis. Moreover, these results are indicative of the information to be obtained through genetic manipulations in conjunction with embryonic tissue transplantation. It has been clearly demonstrated that *in ovo* transplantation provides a suitable environment for the development of mammalian cells and that the information supplied by this environment is capable of promoting mouse cell differentiation.

Chimeras were prepared by transplanting somites from nine-day postcoitum mouse embryos into two-day-old chick embryos at different axial levels. Mouse somitic cells then differentiated *in ovo* into dermis, cartilage, and skeletal muscle, as they normally do in the course of development, and were able to migrate into chick host limb. To trace the behavior of somitic myogenic stem cells more closely, somites arising from mice bearing a transgene of the desmin gene linked to a reporter gene coding for *E. coli*  $\beta$ -galactosidase were grafted *in ovo*. Interestingly, the transgene was rapidly expressed in myotomal muscles derived from implants. In the limb muscle mass, positive cells were found several days after implantation. This method facilitates investigation of the mechanisms of mammalian development, allowing the normal fate of implanted mouse cells to be studied and providing suitable conditions for identification of descendants of genetically modified cells.

Chimeras were also prepared by transplanting fragments of neural primordium mouse embryos into chick embryos at different axial levels. Mouse neuroepithelial cells differentiated and organized to form the different cellular compartments normally constituting the central nervous system. The graft entered into the development of the peripheral nervous system through migration of neural crest cells associated with mouse neuroepithelium. Depending on the graft level, mouse crest cells participated in the formation of various derivatives, such as head components, sensory ganglia, orthosympathetic ganglionic chain, nerves, and neuroendocrine glands. Knock-out mice with the tenascin genes inactivated, which express LacZ instead of tenascin and show no tenascin production (29), were used specifically to follow Schwann cells lining nerves derived from the implant. Together with the previous results on somite development, this study shows that the chick embryo constitutes a privileged environment, facilitating access to the developmental potentials of normal or defective mammalian cells. It allows study of the histogenesis and precise timing of formation of a known structure, as well as the implications of a given gene at all equivalent mammalian embryonic stages.

## Bibliography

1. J. Spemann (1901) Arch. Entwicklungsmech. Org. **12**, 224–264.
2. J. Spemann (1921) Arch. Entwicklungsmech. Org. **48**, 533–570.
3. J. Spemann and H. Mangold (1924). Arch. Mikrost. Anat. Arch. Entwicklungsmech. **100**, 599–638.
4. J. J. Henry and R. M. Grainger (1987) Dev. Biol. **124**, 200–214.
5. N. M. Le Douarin (1969) Bull. Biol. Fr. Belg. **103**, 435–452.
6. N. M. Le Douarin (1982) *The Neural Crest*, Cambridge University Press, London and New York.
7. R. L. Gardner (1968) Nature **220**, 596–597.
8. L. A. Moustafa and R. L. Brinster (1972) J. Exp. Zool. **181**, 193–202.
9. J. Fontaine-Pérus, V. Jarno, C. Fournier Le Ray, Z. L. Li, and D. Paulin (1995) Development **121**, 1705–1718.

10. J. Fontaine-Péru, P. Halgand, Y. Chéraud, T. Rouaud, M. E. Velasco, C. Cifuentes-Diaz, and F. Rieger (1997) *Development* **124**, (16), 3025–3036.
11. R. L. Gimlick and J. Cook (1983) *Nature* **306**, 471–473.
12. J. W. Saunders Jr, J. M. Cairns, and M. T. Gasseling (1957) *J. Morphol.* **101**, 57–88.
13. H. Spemann and O. Schötté (1932) *Naturwissenschaften* **20**, 463–467.
14. M. Servetnick and R. M. Grainger (1991) *Development* **112**, 177–188.
15. N. M. Le Douarin and M. A. Teillet (1973) *J. Embryol. Exp. Morphol.* **30**, 31–48.
16. N. M. Le Douarin (1993) *TINS* **16**, 64–72.
17. E. Balaban (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3657–3660.
18. A. Grappin-Botton, M. A. Bonnin, L. A. McNaughton, R. Krumlauf, and N. M. Le Douarin (1995) *Development* **121**, 2707–2721.
19. A. Chevallier, M. Kieny, A. Mauger, and P. Sengel (1977) *Vertebrate Limb and Somite Morphogenesis*, Cambridge University Press, Cambridge, pp. 421–432.
20. B. Christ, H. J. Jacob, and M. Jacob (1977) *Anat. Embryol.* **150**, 171–186.
21. C. P. Ordhal and N. Le Douarin (1992) *Development* **114**, 339–353.
22. C. A. Smith and R. S. Tuan (1996) *Teratology* **52**, 333–345.
23. A. E. Münsterberg, J. Kitajewski, D. A. Bumcrot, A. P. McMahon, and A. Lassar (1995) *Genes Dev.* **9**, 2911–2922.
24. H. M. Stern, A. M. C. Brown, and S. D. Hauschka (1995) *Development* **121**, 3675–3686.
25. O. Pourquié et al. (1996) *Cell* **86**, 461–471.
26. G. Brill, N. Kahane, C. Carmeli, D. Von Schack, Y. A. Barde, and C. Kalcheim (1995) *Development* **121**, 2583–2594.
27. B. Mintz and W. W. Baker (1967) *Proc. Natl. Acad. Sci. USA* **58**, 592–598.
28. O. Chisaka and M. R. Capecchi (1991) *Nature* **350**, 473–479.
29. Y. Saga, J. Yagi, Y. Ikawa, T. Sakakura, and S. Aizawa (1992) *Genes Dev.* **6**, 1821–1838.

### Suggestions for Further Reading

30. E. Dupin, C. Ziller, and N. M. Le Douarin (1998) “The Avian Embryo as a Model in Development Studies: Chimeras and *in vitro* Clonal Analysis,” *Curr. Top. Dev. Biol.* **36**, 1–35.15
31. N. Le Douarin, F. Dieterlen Lièvre, and M. A. Teillet (1996) “Quail-Chick Transplantations” *Methods Cell Biol.* **51**, 23–59.

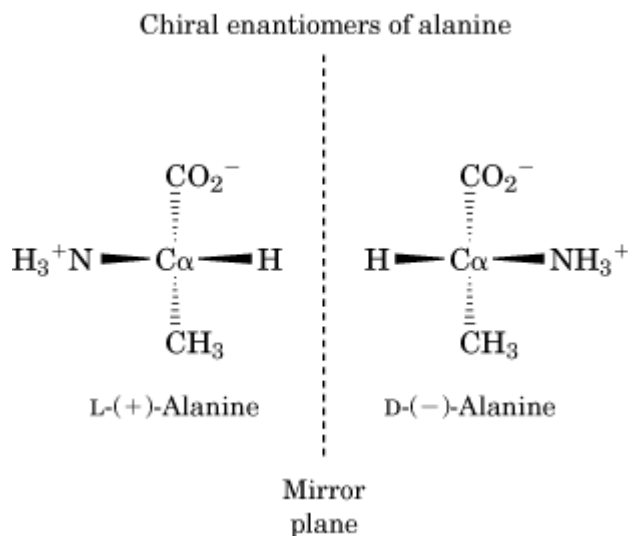
### Chiral and Chiral Center

The word chiral is derived from the Greek word “cheir” for “hand.” It was first used to describe any molecule that could not be superimposed on its mirror image, like left and right hands, by Lord Kelvin (1). This supplanted the less precise term dissymmetrical used by Pasteur (2). The nonsuperimposable mirror image isomers are [enantiomers](#). The physical property that differentiates enantiomers is the direction that they rotate plane polarized light. Thus, the two enantiomers are differentiated as either dextro(+) or levo(–) rotatory, depending on whether the rotation of the polarized light is clockwise or counterclockwise, respectively. Solutions containing an excess of one of the enantiomers are said to be optically active. A solution containing exactly equal concentrations

of both enantiomers is called racemic (see [Racemic And Racemization](#)).

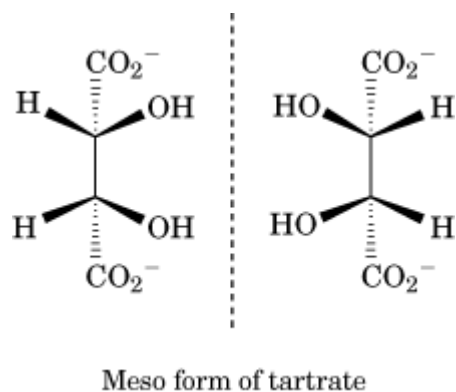
The most common structural element that generates chiral molecules is the presence of a tetrahedral atom with four different substituents, by definition a chiral center. The  $\alpha$ -carbon of the  $\alpha$ -amino acids (Fig. 1), except for [glycine](#), and the hydroxymethylene carbons of carbohydrates are all chiral centers. The mirror images of chiral molecules have different absolute [configurations](#), which are uniquely identified by the rules proposed by Cahn et al. (3) (see [Configuration](#)).

**Figure 1.** The enantiomers of alanine, which are mirror images. Note that the  $\alpha$ -carbon is chiral because of the four different substituents. The (L) and (D) designate the absolute configuration of the chiral center, whereas the (+) and (-) indicate which direction an aqueous solution of the enantiomer will rotate the plane of polarized light.



Possessing a tetrahedral chiral center is not required for chirality, nor is the presence of a chiral center proof that the molecule will be chiral. Other arrangements of atoms can generate nonsuperimposable mirror images. One example of biochemical importance is the hexacoordinate complexes of metal ions, such as the Mg ATP complex (4). Other chiral molecular structures have been surveyed by Testa (5). A molecule with an internal mirror plane may have chiral centers but not be chiral. These are meso compounds. Because of their internal symmetry, these molecules can be superimposed on their mirror image. The **diastereomeric** form of tartaric acid, which has two chiral centers, but a mirror plane between C2 and C3, provides an example of a meso compound (Fig. 2).

**Figure 2.** Meso tartrate, which contains two chiral centers C2 and C3, but is not chiral because there is an internal plane of symmetry. Note that, unlike the enantiomers of alanine, these mirror images can be superimposed by a 180° rotation.



### Bibliography

1. L. Kelvin (1894) Boyle Lect. J. Oxford Univ. Junior Sci. Club (May 25), 25.
2. L. Pasteur (1853) Pharmacol. J. **XIII**, 111.
3. R. S. Cahn, C. K. Ingold, and V. Prelog (1966) *Angew. Chem. Int. Ed.* **5**, 385–415.
4. E. A. Merritt et al. (1978) *Biochemistry* **17**(16), 3274–3278.
5. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), Vol. **3**, Elsevier, Amsterdam, pp. 1–48.

### Suggestions for Further Reading

6. J. March (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, pp. 82–100.
7. K. Mislow (1966) *Introduction to Stereochemistry*, W. A. Benjamin, New York, pp. 50–60.

## Chironomus

The subfamily chironomini includes *Chironomus* and related genera that by some estimates contain over 5000 species. Because of this diversity, many studies of chironomids focus on environmental biology and systematics. But it was the giant polytene chromosomes of salivary gland cells and the bright red color of the larvae of these ubiquitous flies that first drew the attention of experimental biologists more than 100 years ago.

The alternating dark bands (chromomeres) and light bands of the polytene chromosomes are punctuated by puffed regions, the largest of which were named Balbiani rings to honor their discovery by E. G. Balbiani in 1881. The possible relationship between chromomeres and genetic loci on the one hand, and between puffs and gene activity on the other, excited the interest of cytologists and geneticists. One of these, W. Beermann, examined the morphology of chromosomes in four cells in the “special lobe” of the salivary glands of larvae that were derived from a cross between *C. tentans* and *C. pallidivittatus*. He was able to correlate the production of prominent secretion granules by these cells with the presence of a specific chromosome puff in both the hybrid and in *C. pallidivittatus*, and showed that the puff was absent from the special cell chromosomes in *C. tentans*, which do not contain the secretion granules. C. Pelling later demonstrated that polytene chromosome puffs selectively incorporated radioactive uridine, firmly establishing that they are sites of active gene transcription. The sequential appearance of some, and the disappearance of other chromosome puffs during development could be mimicked by treating cultured late-stage larval

salivary glands with the steroid hormone, 20-hydroxyecdysone. Thus, the same hormone that signals molting and metamorphosis in insects also controls the differential transcription of genes encoded by responsive puffs (see [Ecdysone](#)).

In a prodigious feat of molecular isolation that predates molecular cloning by almost a decade, J. E. Edstrom and his colleagues used micro manipulation to dissect individual Balbiani rings. A very large (75S) mRNA (messenger RNA) encoding a polypeptide of about 850 kDa was extracted from pooled Balbiani ring 2 isolates. This and other salivary gland secretory proteins encoded by the major Balbiani rings are extruded through the larval mouth and spun together with benthic detritus to make tabular houses. The stickiness of *Chironomus* “silk” results from the aggregation properties of interacting secretory polypeptides from the major Balbiani rings, and several studies focus on those aspects of structure that could explain the adhesive properties of these unusual biopolymers. In *C. tentans*, each major Balbiani ring encodes polypeptides with internal repetitive domains. The core repeat of one of these domains has four cysteine residues in invariant positions; a subrepeat region contains tandem *basic-pro-acidic* residue motifs. This arrangement might encourage protein strand alignment and polymerization of the “silk”. In other studies, B. Daneholt and colleagues have capitalized on the large size of the Balbiani Ring gene transcripts to investigate messenger ribonucleoprotein (mRNP) assembly and transport through nuclear pores; polytene chromosome centromere and telomere evolution, replication, and function are also areas of active research.

Recent studies have used expression of a chironomid/drosophilid ecdysone receptor (ultraspiracle) gene to understand hormone receptor domain function. Still other seek to define sex-determining genes and proteins, and to characterize diverse transposable elements in chironomids. A significant number of reports have focused on *Chironomus* globin gene structure, polymorphism, and evolution. The rest of this article deals with developmental and functional studies of *Chironomus* hemoglobins and with the structure, organization, and evolution of the many globin genes in this large multigene family.

## 1. *Chironomus* Hemoglobin Structure and Function

The red pigment so visible through the larval integument was identified as hemoglobin by E. R. Lankester in 1872. In 1932, T. Svedberg et al. used their then-new technique of sedimentation velocity ultracentrifugation to estimate the molecular weight of hemoglobin from *Chironomus* (and other invertebrates); chironomid hemoglobins “weighed in” at 32 kDa. Physiological studies focused on the unusually high oxygen affinity of the hemoglobin (compared to its vertebrate counterparts) and on how it facilitated oxygen supply to larval tissues in a relatively hypoxic benthic habitat.

*Chironomus* hemoglobin is actually a mixture of many hemoglobin proteins synthesized and secreted by larval fat body. Long before the first chironomid globin gene was cloned, much was known about the structure of these proteins. In contrast to the tetrameric vertebrate hemoglobins (see [Hemoglobin](#)), chironomid hemoglobins do not show cooperative binding of oxygen, existing either as monomers or homo-dimers under physiological conditions (Svedberg’s 32 kDa “hemoglobin” was presumably a mixture of homodimers). Over the course of a decade, the amino acid sequences of 12 *C. thummi* globins were determined in the laboratory of G. Braunitzer. Since none is truly similar in sequence to its vertebrate counterparts, alignments of chironomid and vertebrate globin sequences were only possible because they share several amino acids at crucial positions known to be involved in binding the heme group. On the other hand, comparison of the three-dimensional structure of a monomeric chironomid hemoglobin (determined by X-ray crystallography) with that of vertebrate globins reveals many shared structural features, including similarly organized  $\alpha$ -helical domains and a characteristic hydrophobic heme pocket. These observations clearly indicate that similar molecular functions and three-dimensional structures can be attained in different proteins despite considerable differences in primary structure. The differences in primary structure between globins of very different organisms reflects adaptation to species-specific oxygen requirements over long evolutionary periods.

Whereas at least two *C. thummi* globin genes encode identical polypeptides (1), pairwise comparisons of amino acid sequences reveal remarkably dissimilar globins. These differences have also been ascribed to evolutionary specialization of different hemoglobins that now serve different respiratory functions (2). However, the physiological roles of individual hemoglobins are at best only poorly understood. *In vivo* determinations of larval O<sub>2</sub> binding and storage abilities under different conditions of oxygen tension necessarily represent a function of the aggregate set of hemoglobins present at the time of measurement. As expected, the O<sub>2</sub> affinities determined for purified hemoglobins that are present in different proportions in the hemolymph are not easily related to the aggregate *in vivo* data. What does emerge from such experiments is a loose correlation of the generally high oxygen affinities of the hemoglobins and the benthic, often oxygen-depleted habitat of larvae.

Larvae live in the mud at the bottom of freshwater ponds and lakes, developing in burrows constructed to detritus cemented by the sticky secreted salivary gland proteins discussed above. Moreover, chironomids can withstand the oxygen stress of polluted or eutrophic waters, often comprising more than 80% of all animal individuals. This tolerance is due in part to the high concentrations of hemolymph hemoglobin. Zebe (3) suggested that hemoglobins of *C. thummi* may not be required for oxygen delivery at higher partial pressures (oxygen needs under these conditions are presumably satisfied by diffusion). Rather, because they can bind O<sub>2</sub> at low partial pressures (50 to 10 torr), the hemoglobins allow continued respiration in times of severe oxygen deficit, when they become physiologically significant. From this, one might predict that oxygen stress will favor larval production of hemoglobins with particularly high O<sub>2</sub> binding affinities. Homodimeric hemoglobins have higher O<sub>2</sub> affinities than the monomeric hemoglobins and do indeed rise in concentration in response to circadian and seasonal declines in O<sub>2</sub> availability and water quality, where they may be more involved in O<sub>2</sub> transport under these conditions than monomeric hemoglobins. With a more pronounced Bohr effect than monomeric hemoglobins, dimeric hemoglobins may also allow response to more immediate changes in oxygenation conditions due to local fluctuations in pH.

Overall hemoglobin levels reach their highest concentrations in the hemolymph of *C. thummi* in the fourth, or last larval stage (instar). Although it is likely that newly synthesized hemoglobins function mainly in respiration, hemolymph hemoglobin is also taken up by developing oocytes and deposited in yolk granules during oogenesis in pupae (4, 5). This hemoglobin undoubtedly serves a nutrient rather than a respiratory function during embryogenesis. At the same time, developing eggs may contain hemoglobins of their own making. Cross-reactive antihemoglobin antibodies tagged with colloidal gold particles are localized over rough endoplasmic reticulum in young embryos, suggesting that hemoglobin synthesis begins at embryogenesis (6). Hemoglobins first accumulate sufficiently to become visible as a red pigment in the blood of second instar larvae. Thereafter, hemoglobin synthesis is cyclic, being high during the instar and declining at each molt. Hemoglobin synthesis occurs for the last time just after pupariation, ceasing sometime in the pupal stage, by which time most existing hemoglobins are degraded. Hemoglobins have disappeared and are virtually undetectable several hours after adult emergence, at which time the adult flies have a fully developed tracheal respiratory system. Once again, in spite of these observations, it is not possible to state with certainty the physiological role played by any given hemoglobin at any given moment. As we will see, this issue may be moot for many of the hemoglobins, since there is reason to believe that much of the structural variation between hemoglobins has not been selected to perform a unique physiological function, but instead is the result of neutral evolution.

## 2. Developmental Regulation of *Chironomus* Hemoglobin Synthesis

While cyclic hemoglobin production is tied to molting, it should also come as no surprise that the synthesis of individual members of such a large family of hemoglobins is regulated differentially during development. For example, polyacrylamide gel electrophoresis (PAGE) of *C. thummi* fourth



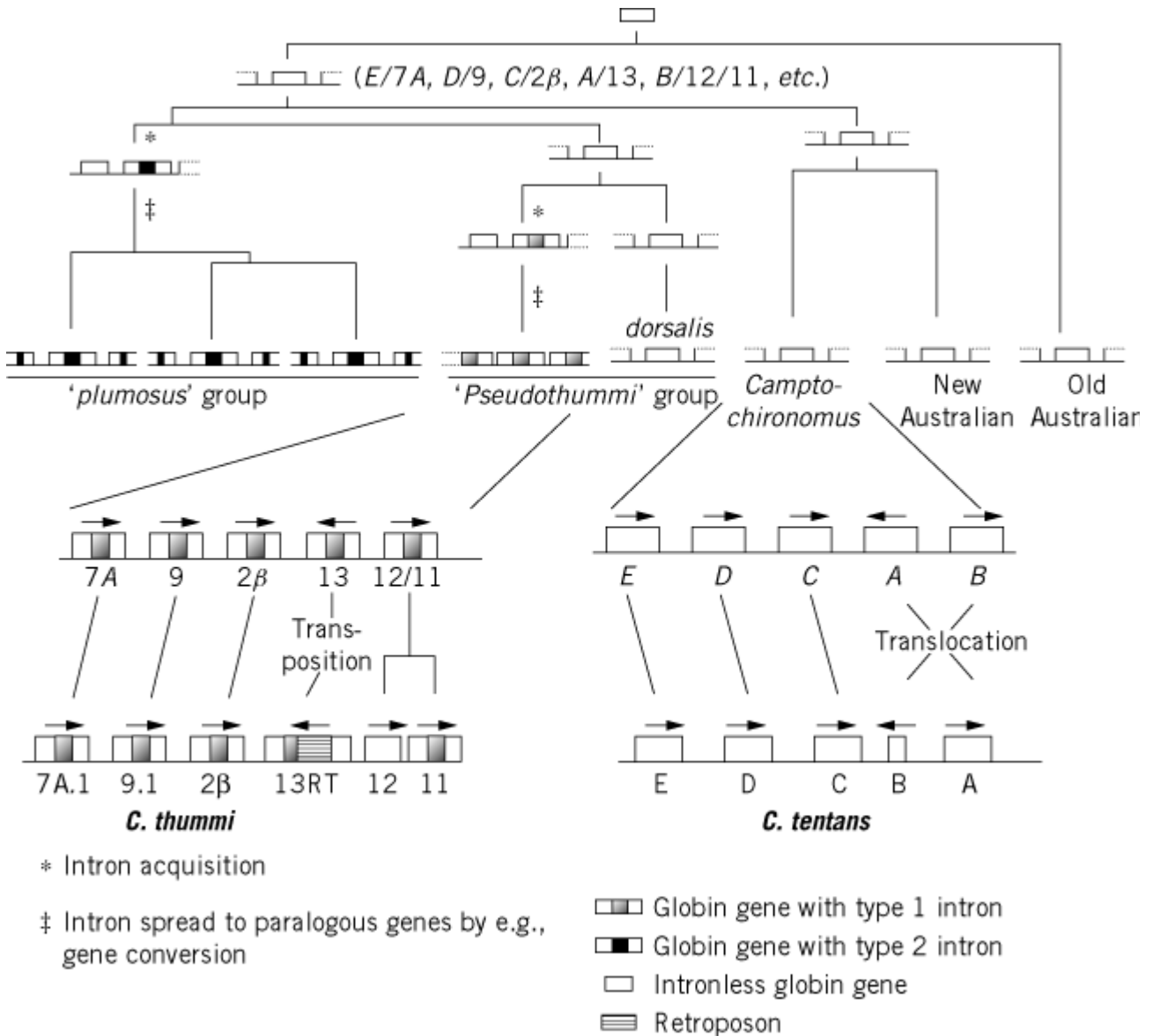
instar hemolymph reveals two hemoglobins that are absent in electrophoretic separations of third instar larval hemolymph. Another hemoglobin, detectable in embryo homogenates by electrophoresis and immunoblotting, is undetectable in last instar larval hemolymph. In S1 nuclease protection assays performed to follow specific globin gene expression during development, transcripts of the ct-1, ct-4, ct-7B4, and ct-7B5 globin genes (ct-, ctp-, ctn-, and kc- are prefixes for *C. thummi*, *C. piger*, *C. tentans*, and *K. cornishi* globins and globin genes.) are detected at high levels in 3rd instar *C. thummi* larvae (7). Transcripts of each of these globin genes decline to undetectable levels immediately after molting to the fourth instar, reappearing 24 hr later, along with two new fourth-instar-specific globin transcripts (ct-3, ct-6). The expression of each of these globin mRNAs was tracked further during fourth instar development. The results of these experiment indicate that globin transcript steady-state levels follow one of two kinetic patterns. Transcripts of the ct-3 and ct-6 genes rise from very low or undetectable levels after the molt, and then remain elevated through 8 days of the instar. Globin mRNAs expressed in both instars also rise from low or undetectable levels just after the molt. But unlike the specific globin mRNAs, they reach maximum levels after 2 to 3 days and then decline to near basal levels. The absence of detectable globin transcripts just after the third to fourth instar molt is consistent with the observation that the rate of synthesis of hemoglobins of cultured larval tissues exposed to the molting hormone 20-hydroxyecdysone (20HE) is very much reduced. Like other steroid hormones, 20HE binds to an intracellular receptor, which then binds to a hormone response element(s), directly controlling the transcription of adjacent genes. Such direct action on globin is unlikely, since a search in DNA lying between or flanking any *C. thummi* globin gene sequences published or recorded in international DNA sequence databases failed to detect 20HE response elements. As yet unidentified factors must mediate the hormone-induced repression of hemoglobin synthesis during molting. There is no evidence at all for a role for 20HE in the selective repression of “non-stage-specific” globin mRNA levels during the fourth instar, indicating the action of still other factors in the selective silencing of some but not all globin genes.

### 3. *Chironomus* Globin Gene Structure and Organization

At the same time as the differential expression of hemoglobins or globin transcripts was being described, many more globin genes were being characterized. Known *C. thummi* genes exceed the number needed to account for the 12 hemoglobins that were first isolated. The loci of individual and clustered globin genes were mapped by *in situ* hybridization of cloned DNAs to polytene salivary gland chromosomes, and as noted, the expression of many globin genes including several for which no protein product has yet been isolated, has been shown. Table 1 summarizes some key properties of cloned globin genes characterized thus far. Genes that are clearly inactive include ct-13RT (which suffered insertion of a retroposon in coding DNA; see Fig. 1), ct-V (with an inappropriately located stop codon), and ctn-ORFB (with its coding region truncated at the 5' and 3' ends). With these exceptions, all of the known chironomid globin genes possess the structural prerequisites of active genes, including appropriately located TATA (promoter) and AATAAA (polyadenylation) signals, an open reading frame, and a stop codon. In addition to the first of full-length globin genes in Table 1, many partial genes (8, 9) amplified by PCR from numerous species (eg, for globins 2b, 7, and 9) are probably active. When compared to vertebrate (eg, human) globin genes, a striking difference is that the latter are separated by thousands of nucleotides, whereas chironomid globin genes are often separated by less than 500 nucleotides. Another difference is that chironomid globin genes lack the hallmark structure of three exons separated at conserved positions by two introns that characterize vertebrate (and even plant) globin genes. Instead, different chironomid globin genes are either intronless or contain a single introns at divergent locations. The growing database of chironomid globin genes from several species, their unusual structure, and their organization have sparked considerable interest in the evolution of this unique multigene family.

**Figure 1.** Intron acquisition by chironomid globin genes. Relative times of independent acquisition of introns at different locations (“type 1” vs “type 2”) by ancestral globin genes in the “*pseudothummi*” and “*plumosus*” phylogenetic groups are suggested. Only *C. dorsalis*, the oldest member of the “*pseudothummi*” group, has intronless orthologues of the

intron-bearing genes. Phylogenetic clusters are taken from Guryev et al. 2001 (9). The evolutionary history of intron-bearing and intronless orthologous gene clusters *C. thummi* and its distant relative, *C. tentans* is also shown. ct-, ctn-, and ORF (“open reading frame”) designations are deleted to save space. Arrows above the genes indicate direction of transcription. Ancestral genes are designated in italics. Details are explained in the text.



**Table 1. Chironomid Globin Genes**

| Organism         | Protein |         | Status <sup>a</sup> | Homologues | locus <sup>b</sup> | Comments  |
|------------------|---------|---------|---------------------|------------|--------------------|---|
|                  | Gene    | Product |                     |            |                    |   |
| <i>C. thummi</i> | ct-1    | I       | active              | —          | F2b3               | Allele of ct-1A; intronless; aa sequence avail. |
|                  | ct-1A   | IA      | active              | —          | F2b3               | Allele of ct-1;                                 |

|              |          |          |   |      |  |
|--------------|----------|----------|---|------|--|
|              |          |          |   |      | intronless; aa<br>sequence avail.                    |
| ct-3-1       | III      | *        | — | A1B2 | Intronless; * 3-1<br>&/or3-2 = avail.<br>aa sequence |
| ct-3-2       | III      | *        | — | A1B2 | Intronless, *3-1<br>&/or3-2 = avail.<br>aa sequence  |
| ct-4-1       | IV       | *        | — | A1B2 | Intronless, *4-1<br>&/or4-2 = avail.<br>aa sequence  |
| ct-4-2       | IV       | *        | — | A1B2 | Intronless, *4-1<br>&/or4-2 = avail.<br>aa sequence  |
| ct-E         | Putative | Putative | — | A1B2 | Intronless   |
| ct-6         | VI       | Active   | — | F2b3 | Intronless; aa<br>sequence avail.                    |
| Not<br>found | VIIB-1   | Active   | — | —    | aa sequence<br>avail.                                |
| Not<br>found | VIIB-2   | Active   | — | —    | aa sequence<br>avail.                                |
| ct-7B3       | VIIB-3   | Active   | — | N.D. | Intronless;<br>protein product<br>assumed            |
| ct-7B4       | VIIB-4   | Active   | — | F2b3 | Intronless;<br>protein product<br>assumed            |
| ct-7B5       | VIIB-5   | Active   | — | F2b3 | Intronless;<br>protein product<br>assumed            |
| ct-7B6       | VIIB-6   | Active   | — | F2b3 | Intronless;<br>protein product<br>assumed            |
| ct-7B7       | Putative | Putative | — | F2b3 | Intronless;<br>protein product<br>assumed            |

|          |          |          |   |      |   |
|----------|----------|----------|---|------|---|
| ct-7B8   | Putative | Putative | — | F2b3 | Intronless;<br>protein product<br>assumed   |
| ct-7B9   | Putative | Putative | — | F2b3 | Intronless;<br>protein product<br>assumed   |
| ct-7B10  | Putative | Putative | — | F2b3 | Intronless;<br>protein product<br>assumed   |
| ct-7B9/5 | Putative | Putative | — | N.D. | Intronless<br>chimera; protein<br>product<br>assumed                                |
| ct-Y     | Putative | Active   | — | F2b3 | Partial DNA<br>sequence, RT-<br>PCR product   |
| ct-W     | Putative | Active   | — | F2b3 | Intronless;<br>protein product<br>assumed   |
| ct-V     | None     | Inactive | — | F2b3 | Stop codon in<br>middle of<br>coding region   |
| ct-Z     | Putative | Active   | — | F2b3 | Intronless;<br>protein product<br>assumed   |
| ct-7A.1  | VIIA     | Active   | — | F2b3 | Central intron;<br>aa sequence<br>avail.  |
| ct-9.1   | IX       | Active   | — | F2b3 | Central intron;<br>aa sequence of<br>allele avail.                                  |
| ct-2b    | II2b     | Active   | — | F2b3 | Central intron;<br>aa sequence<br>avail.  |
| ct-13RT  | None     | Inactive | — | F2b3 | Central intron;<br>retroson<br>insertion in<br>coding DNA;<br>may be<br>transcribed |

|                   |          |          |                  |      |   |
|-------------------|----------|----------|------------------|------|---|
|                   | ct-12    | Putative | Putative —       | F2b3 | Intronless  |
|                   | ct-11    | Putative | Putative —       | F2b3 | Central intron  |
| <i>C. piger</i>   | ctp-7B4  | Putative | Putative ct-7B4  | F2b3 | Intronless; protein product assumed   |
|                   | ctp-7B5  | Putative | Putative ct-7B5  | F2b3 | Intronless; protein product assumed   |
|                   | ctp-7B6  | Putative | Putative ct-7B6  | F2b3 | Intronless; protein product assumed   |
|                   | ctp-7B7  | Putative | Putative ct-7B7  | F2b3 | Intronless; protein product assumed   |
|                   | ctp-7B8  | Putative | Putative *       | F2b3 | Intronless; protein product assumed; * chimera of ct-like 7B8 and 7B10 genes; |
|                   | ctp-Y    | Putative | Putative ct-Y    | F2b3 | Intronless; protein product assumed   |
|                   | ctp-W    | Putative | Putative ct-W    | F2b3 | Intronless; protein product assumed   |
|                   | ctp-V    | Putative | Putative ct-V    | F2b3 | Intronless; protein product assumed   |
|                   | ctp-Z    | Putative | Putative ct-Z    | F2b3 | Intronless; protein product assumed   |
| <i>C. tentans</i> | ctn-ORFE | Putative | Putative ct-7A.1 | —    | Intronless; protein product assumed   |
|                   | ctn-ORFD | Putative | Putative ct-9.1  | —    | Intronless; protein product   |

|                           |          |          |                    |   |   |
|---------------------------|----------|----------|--------------------|---|---|
|                           |          |          |                    |   | assumed   |
|                           | ctn-ORFC | Putative | Putative ct-2      | — | Intronless; protein product assumed                           |
|                           | ctn-ORFB | None     | Inactive *ct-11/12 | — | Intronless; may be transcribed; * ancestor to ct-11 and ct-12 |
|                           | ctn-ORFA | Putative | Putative ct-13RT   | — | Intronless; protein product assumed                           |
| <i>C. pallidivittatus</i> | cpa-3-1  | Putative | Putative ct-3-1    | — | Intronless; protein product assumed                           |
|                           | cpa-3-2  | Putative | Putative ct-3-2    | — | Intronless; protein product assumed                           |
|                           | cpa-4-1  | Putative | Putative ct-4-1    | — | Intronless; protein product assumed                           |
|                           | cpa-4-2  | Putative | Putative ct-4-2    | — | Intronless; protein product assumed                           |
|                           | cpa-E    | Putative | Putative cpa-E     | — | Intronless; protein product assumed                           |
|                           | cpa-F    | Putative | Putative ?         | — | Intronless; protein product assumed                           |
| <i>K. cornishi</i>        | kc-7BA   | Putative | Putative ?         | — | Intronless; protein product assumed                           |
|                           | kc-7BB   | Putative | Putative ?         | — | Intronless; protein product assumed                           |
|                           | kc-7BC   | Putative | Putative ?         | — | Intronless; protein product assumed                           |

|        |          |            |   |  |
|--------|----------|------------|---|--|
| kc-7BD | Putative | Putative ? | — | Intronless;<br>protein product assumed |
| kc-7BE | Putative | Putative ? | — | Intronless;<br>protein product assumed |
| kc-7BF | Putative | Putative ? | — | Intronless;<br>protein product assumed |
| kc-7BG | Putative | Putative ? | — | Intronless;<br>protein product assumed |

<sup>a</sup> Active genes have been shown to be transcribed and/or translated, or to correspond to an identified globin polypeptide.

<sup>b</sup> Positions mapped to polytene salivary gland chromosomes of *C. thummi*; See Ref. 44 for the mapping of the same genes to the salivary gland chromosomes of other chironomid species.

#### 4. Evolution of *Chironomus* Globins

Phylogenetic analyses support an early origin of monomeric hemoglobins about 500 million years ago and the appearance of homodimeric hemoglobins about 250 to 300 million years ago, corresponding roughly with the origin of Diptera. Recently, a search of the *Drosophila* genome database using a chironomid globin gene query found a distant relative, designated a neuroglobin because it was expressed in brain tissue in *D. melanogaster*; and mouse and human neuroglobins have since been characterized (10-13). The chironomid globin gene family has experienced many of the same evolutionary events as their multigene families, such as point-mutational inactivation, translocations, gene loss, gene conversion, etc. They also support two somewhat unorthodox hypotheses. One that has received much attention in recent years, is that eukaryotic genes can acquire introns, complete with correctly placed intron/exon junctions and splicing signals. The other hypothesis is that much of the diversity of hemoglobin structure in *Chironomus* has arisen by neutral evolution rather than by functional adaptations of the proteins. There is resistance to ascribe to neutral evolution those mutations causing even minor differences in amino acid sequence between the products of homologous genes in different species or paralogous genes in the same organism. Such differences between globin genes in *Chironomus*, and between the products of multigene family members in general, are widely thought to be the result of adaptive selection. However, evidence from the different chironomid globins and globin genes suggest that many were not selected for their adaptive value, but exist instead as a neutral consequence of positive selection of the expression of high levels of hemoglobin from multiple genes.

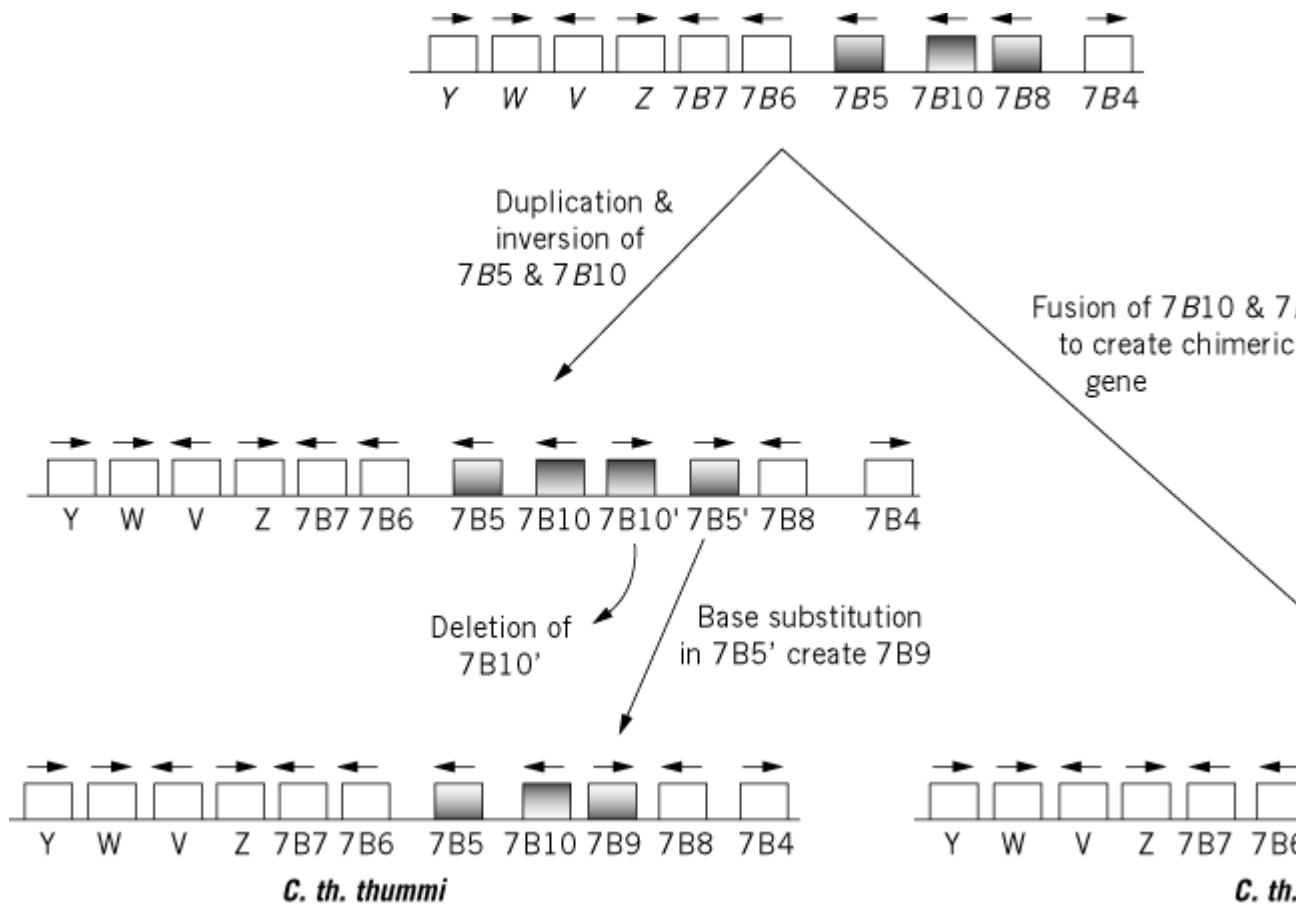
##### 4.1. Origin of A Central Intron

By the time the first chironomid globin genes were cloned in 1984, the sequencing of vertebrate globin genes had become a cottage industry. Each new report presented a globin gene with the same mosaic structure: three exons and two introns, the latter always interrupting the coding DNA at the same location. In a prophetic para (14), Go analyzed X-ray diffraction data for a vertebrate globin and concluded that in addition to the outer introns, the ancestral globin gene had a third intron in the center of the gene such that the entire gene was separated into four exons of similar length, each

coding for a separate structural domain. Go's prediction was fulfilled shortly thereafter with the report of the first plant (legume) globin gene, which had a centrally located intron in addition to two introns at the same position as those in the vertebrate globin genes. Apparently the structure of plant globin genes represents the ancestral state. Many of the first chironomid globin genes to be reported were intronless, a condition presumably due to the reverse transcription of a completely processed transcript and reinsertion of the complementary DNA into the genome. The discovery of the *C. thummi* ct-2b and ct-9.1 genes (15), and of several other invertebrate globin genes with a central intron, however, raised questions about the origin of chironomid globin gene introns. Unlike the conservation of location of the outer vertebrate globin gene introns, the location of central introns is not conserved across kingdoms or phyla, and none is at the position predicted by Go. Could they, like the outer introns in globin genes, still be descendants of an ancestral intron? Gilbert suggested that introns at nonconserved locations have moved by "intron sliding." Others have argued that the discordant location of introns is more easily explained if introns can be acquired by previously uninterrupted genes. Hankeln et al. (8) reported that the globin genes of some species have a "central" intron at a slightly different location than those of other species. A chironomid mitochondrial gene phylogeny groups flies with an intron at one location ("type 1") on the "*pseudothummi*" branch, and groups those with an intron at the other location ("type 2") on a separate "*plumosus*" evolutionary branch (8, 16, 17). The phylogeny supports the separate acquisition of an intron by the homologous genes after divergence of the two lineages (Fig. 1). Several plausible mechanisms for intron acquisition have been suggested. In considering these, Kao et al. (15) offered a mechanism whereby a chironomid globin gene containing a centrally located  $\frac{1}{4}$ AGGT $\frac{1}{4}$  tetramer (to provide donor and acceptor splice sites) could have acquired a central intron by a simple gene duplication, accompanied (or followed) by deletion of all but the DNA presently found in the intron. A mechanism proposed by Hankeln et al. (8), where introns are acquired by a reversal of the splicing reaction, might explain intron insertion in globin genes lacking these AGGT splice site precursors. An intron once acquired by a single gene might spread to nearby related genes by gene conversion (8); see Fig. 2).

**Figure 2.** Evolutionary history of orthologous gene clusters in *C. thummi* and its sibling subspecies, *C. piger*. ct-, and ct<sub>1</sub> save space. Arrows above the genes indicate direction of transcription. Ancestral genes are designated in italics. Shaded recent duplications, losses, and fusions.





#### 4.2. Evolution of a Globin Multigene Family

The DNA sequences of orthologous globin gene clusters in different chironomid species support the adaptive selection of a high globin gene copy number and the essentially neutral evolution of individual globin genes. The *ctp-Y*, *ctp-W*, *ctp-V*, and *ctp-Z* genes, and five *ctp-7B* genes in *C. piger* (18), and the *ctp-Y*, *ctp-W*, *ctp-V*, and *ctp-Z* genes, and seven *ctp-7B* genes of sibling species *C. thummi* (1) are a case in point. The phylogenetic study summarized in Figure 2 traces the evolutionary history of the two clusters, showing the seven genes that are clearly homologous, and a series of gene duplications and deletions necessary to explain the current sequence and organization of the clusters. Pairwise comparisons show that at the nucleotide level, homologous globin genes or logical gene pairs (ie, genes of indeterminate orthology, but sharing immediate common ancestry) in the two species share from 96.7% to 99.2% identity. In three cases, homologous genes in each species code for identical polypeptides, despite several millions of years of species divergence. Within *C. thummi*, the *ct-7B5*, *ct-7B6*, *ct-7B9*, and *ct-7B10* paralogous genes are also very similar, as shown by the amino acid differences in the pairwise comparisons shown below:

|      | 7B6 | 7B10 | 7B5 | 7B9 |
|------|-----|------|-----|-----|
| 7B8  | 6   | 4    | 6   | 6   |
| 7B6  |     | 2    | 1   | 1   |
| 7B10 |     |      | 2   | 2   |
| 7B5  |     |      |     | 0   |

In addition to identical paralogues (*ct-7B5*, *ct-7B9*), sequencing of an independent genomic clone revealed a gene identical *ct-7B5* (19). Designated *ct-7B9/5* because its 5' and 3' flanking DNA are

similar to the 5' and 3' regions of ct-7B9 and ct-7B5, respectively, it probably arose by intrachromosomal crossing over between oppositely oriented ct-7B9 and ct-7B5 in a looped (bent) chromosome structure (with deletion of most or all of ct-7b10). Since the recombination point is in fact external to the transcribed region of the gene, 7B9/5 must simply be another ct-7B5 gene or allele. The recombined region (without ct-7B10) may be an independent gene locus or an allelic haplotype. Clearly the 7B globin gene subfamily predates the *thummi/piger* divergence. The ability of *C. thummi* and *C. piger* to lose different descendants of an ancestral cluster by gene fusion, and the overall similarity of paralogous genes, implies that any small differences that accumulate in the genes before or after fusions and corrections are evolutionarily inconsequential. The 7B gene cluster seems to have undergone periodic expansion and contraction (gene duplication, deletion), while maintaining the ability to encode identical or nearly identical globins for millions of years. The best explanation for this phenomenon is that the few differences that accumulate between 7B genes are the result of random drift (neutral evolution), but that the cluster itself has experienced positive selection of a high number of copies of globin genes as a means to ensure the synthesis of large amounts of hemoglobin.

Whereas the orthologous globin gene clusters in *C. thummi* and *C. piger* illustrate early evolutionary events in the life of a multigene family, globin gene clusters in *C. thummi* and *C. tentans* are separated by 60 million years of evolution (16, 20). The extent paralogous and homologous genes sequences are much less similar to each other (except for a region of the ct-2b and ct-9A genes that underwent partial gene correction, homogenizing intronic and exonic DNA in the middle of the two genes). The divergence of the *C. thummi* and *C. tentans* gene clusters offers a wider window through which to test the concept of gene copy number selection. Figure 1 summarizes evolution of these genes after acquisition and spread of a "type 1" intron in one of the clusters. The ancestral *C. thummi* gene cluster must have had five intron-bearing genes before a gene duplication created ct-11 and ct-12. ct-12 has since lost its intron, and both ct-11 and ct-12 have been diverging for some time (16). In contrast, the *C. tentans* cluster contains only intronless homologues, originally ascribed to multiple intron losses from intron-bearing ancestors (16, 20). Figure 1 reflects more recent phylogenetic evidence that some or all of the genes in these clusters arose early in the genus, and that the *C. tentans* genes never acquired an intron in the first place (9). A difference in the organization of the orthologous genes between *C. thummi* and *C. tentans* is explained if a translocation is assumed between ctn-A and ctn-B (or vice-versa in *C. thummi*). Of special note, the recent inactivation of a different (non homologous) gene in each cluster suggests that not all globin genes in a cluster are indispensable, even after 60 million years of evolution. In *C. thummi*, ct-13RT suffered the insertion of a retroposon into coding DNA in exon I. This occurred very recently, because the rest of the coding region, the intron, and the flanking DNA containing promoter and polyadenylation motifs are all still intact. Recent retroposition is supported by the observation that ct-13RT is actually an allele coexisting with a normal, uninterrupted ct-13 allele (17). In *C. tentans*, the ctn-ORFB gene was truncated at both the 5' and 3' ends. Again, this event must have occurred very recently, because the gene retains all other attributes of a viable, transcribable gene. The phylogenetic analyses of the orthologous *C. thummi* and *C. tentans* gene clusters supports the recent inactivation of the ct-13RT allele and the ctn-ORFB genes, and the likelihood that accumulated amino acid differences between many (if not all) homologous and paralogous genes are the result of neutral evolution. If the greater diversity of the *tentans* and *thummi* genes cannot be completely explained as the positive adaptation of structurally diverse hemoglobins, then their maintenance, like that of the large number of 7B genes in *thummi* and *piger*, must be the result of positive gene copy number selection favoring high levels of hemoglobin synthesis.

In sum, the evolution of *Chironomus* globin genes spans more than 250 million years, in which time individual sequences within and across species have accumulated many amino acid substitutions. Some especially diverged hemoglobins may have evolved to serve unique functions, for example in environments that undergo cyclic changes in oxygenation. The specific expression of some globins during development might reflect habitat (and therefore, oxygenation) differences as larvae first descend from the waters' surface to the benthos, and later ascend to the surface at the time of eclosion. On the other hand, differences between many of the globin genes may be neutral, selection

favoring the proliferation of a large family of genes encoding proteins of physiologically similar function. Kimura (21) suggested that duplicated genes accumulating neutral change are preadapted, later becoming substrates for positive adaptation after speciation. A consequence of high gene copy number selection is that some of the many functionally redundant globin genes can serve as raw materials for Darwinian selection, which could explain the evolution and spread of more than 5000 species of Chironomous to diverse habitats.

### Bibliography

1. P. M. Trewitt, R. A. Luhm, F. Samad, S. Ramakrishnan, W-Y. Kao, and G. Bergtrom (1995) *J. Mol. Evol.* **41**, 313–328.
2. M. Goodman, J. Pedwaydon, J. Czelusniak, T. Suzuki, T. Gotoh, L. Moens, F. Shishikura, D. Walz, and S. Vinogradov (1988) *J. Mol. Evol.* **27**, 236–249.
3. E. Zebe (1991) *Comp. Biochem. Physiol.* **99A**, 525–529.
4. P. M. Trewitt, D. R. Boyer, and G. Bergtrom (1986). *J. Insect Physiol.* **32**, 963–969.
5. C. R. Myers, G. Bergtrom, S. M. Crook, D. R. Boyer, M. L. P. Collins (1986). *J. Histochem. Cytochem.* **34**, 221–226.
6. C. R. Myers, M. L. P. Collins, M. Agresti, and G. Bergtrom (1986) *J. Insect Physiol.* **32**, 845–851.
7. D. A. Saffarini, P. M. Trewitt, R. A. Luhm, and G. Bergtrom (1991) *Gene* **101**, 215–222.
8. T. Hankeln, H. Freidl, I. Ebersberger, J. Martin, and E. R. Schmidt (1997) *Gene* **205**, 151–160.
9. V. Guryev, I. Makarevitch, A. Blinov, and J. Martin. (2001) *Mol. Phyl. Evol.* **19**, 9–21.
10. T. Burmester, T. Hankeln (1999) *Mol. Biol. Evol.* **16**, 1809–1811.
11. T. Burmester, B. Weich, S. Reinhardt, T. Hankeln (2000) *Nature* **407**, 520–523.
12. M. Couture, T. Burmester, T. Hankeln, D. L. Rousseau (2001) *J. Biol. Chem.* (in press).
13. J. T. Trent III, R. A. Watts, M. S. Hargrove (2001) *J. Biol. Chem.* (in press).
14. M. Go (1981) *Nature* **291**, 90–92.
15. W-Y. Kao, P. M. Trewitt, and G. Bergtrom (1994) *J. Mol. Evol.* **38**, 241–249.
16. M. C. Gruhl, W.-Y. Kao, and G. Bergtrom (1997) *J. Mol. Evol.* **45**, 499–508.
17. M. C. Gruhl, S. V. Scherbik, K. G. Aimanova, A. Blinov, J-L. Diez, and G. Bergtrom (2000) *Gene* **252**, 153–163.
18. T. Hankeln, C. Luther, P. Rozynek, E. R. Schmidt (1991). In *Structure and Function of Invertebrate Oxygen Carriers* (S. Vinogradov and O. H. Kapp, eds.), Springer Verlag, New York, pp. 287–296.
19. W.-Y. Kao and G. Bergtrom (1995) *Gene* **153**, 241–249.
20. P. Rozynek, M. Broecker, T. Hankeln, E. R. Schmidt (1991) In *Structure and Function of Invertebrate Oxygen Carriers* (S. Vinogradov and O. H. Kapp, eds.), Springer Verlag, New York, pp. 287–296.
21. M. Kimura (1991) *Jpn. J. Genet.* **66**, 367–386.

### Suggestions for Further Reading

22. W. Beermann (1963) Cytological aspects of information transfer in cellular differentiation. *Am. Zool.* **23**–32.
23. S. T. Case and L. Wieslander (1992) Secretory proteins of Chironomous salivary glands: structural motifs and assembly characteristics of a novel biopolymer. *Results Probl. Cell Differ.* **19**, 187–226.
24. B. Daneholt (2001) Assembly and transport of a premessenger RNP particle. *Proc. Natl. Acad. Sci. U.S.A.* **98**, 7012–7017.
25. B. N. Green, A. R. Kuchumov, T. Hankeln, E. R. Schmidt, G. Bergtrom, and S. N. Vinogradov (1998) An electrospray ionization mass spectroscopic study of the extracellular hemoglobins

- from *Chironomus thummi thummi*. Biochim. Biophys. Acta. **1383**, 143–150.
26. R. Lewin R (1984) Surprise finding with insect globin genes. Science **226**, 328.
  27. P. A. Osmulski and W. Leyko (1991) "The structure and function of vhironomus hemoglobins". In *Structure and Function of Invertebrate Oxygen Carriers* (O. H. Kapp, and S. N. Vinogradov, eds.), Springer-Verlag, New York, pp. 305–312.
  28. J. H. Rogers (1989) How were introns inserted into nuclear genes? Trends. Genet. **5**, 340–343.
  29. A. Stoltzfus and W. F. Doolittle (1993) Slippery introns and globin gene evolution. Curr. Biol. **3**, 215–217.
  30. A. Stoltzfus, D. F. Spencer, M. Zucker, J. M. Logsdon, Jr., and W. F. Doolittle (1995) Introns and the origin of protein-coding genes. Science **268**, 1366–1367.

## Chloramphenicol

Chloramphenicol is a bacteriostatic agent that inhibits the growth of many species of [Gram-positive](#) and [Gram-negative bacteria](#); it was the first broad-spectrum antibiotic to be used clinically. Originally obtained from cultures of the soil bacterium *Streptomyces venezuelae*, chloramphenicol inhibits bacterial **protein biosynthesis** by interfering with the intrinsic catalytic activity of the peptidyl-transferase of the [ribosome](#) during the elongation phase of [translation](#).

Eukaryotic cells are generally not affected by chloramphenicol. Nevertheless, the clinical application of this drug for the treatment of systemic infections has been severely curtailed, primarily because of the haematotoxicity associated with its use. Chloramphenicol continues to be used topically, particularly in the treatment of eye infections. The emergence of chloramphenicol-resistant bacteria, apparently in response to the selective pressure exerted by the drug (see [Antibiotic Resistance](#)), has also restricted the use of chloramphenicol in the treatment of bacterial infections.

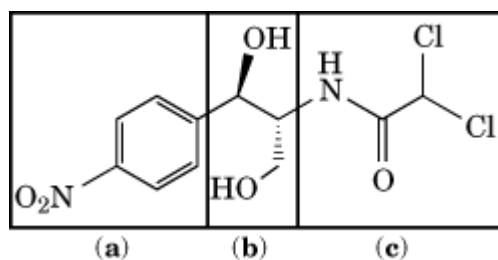
Chloramphenicol resistance is frequently encountered in many genera of bacteria. Several mechanisms of resistance have been described: (i) modification of the target, the bacterial ribosome; (ii) alteration in the permeability of the bacterial cell, leading to low intracellular antibiotic concentrations; and (iii) enzymatic modification of the antibiotic. Resistance most commonly occurs as the result of inactivation of chloramphenicol by [chloramphenicol acetyltransferase](#) (CAT). This enzyme is of value to molecular biologists, because it can be assayed with specificity and sensitivity.

Chloramphenicol is used today in molecular biology primarily as a component of selective media for genetic studies, to distinguish chloramphenicol-resistant from chloramphenicol-susceptible bacteria. Chloramphenicol is also an essential component of CAT assays, which are used to study gene regulation, generally in eukaryotic cells. Finally, chloramphenicol provides a tool with which to investigate the peptidyl-transferase center of the bacterial ribosome.

### 1. Chemistry

Chloramphenicol [D(-)-*threo*-2-dichloroacetamido-1-*p*-nitrophenyl-1,3-propanediol; Fig. [1](#)] has two asymmetric carbon atoms. Of the four stereoisomers, only one, the D-*threo* isomer, has antibacterial activity. This compound can be produced on an industrial scale by chemical synthesis or by fermentation. Chloramphenicol is poorly soluble in water, but dissolves readily in organic solvents (eg, methanol).

**Figure 1.** Chemical structure of chloramphenicol [D(-)-*threo*-2-dichloroacetamido-1-*p*-nitrophenyl-1,3-propanediol].



Derivatives having a wide range of para substituents are active; therefore, this part of the molecule is not involved in specific drug–target interactions. The propanediol moiety, the site of the D-threo configuration, is essential for activity. Acetylation of one or both hydroxyl groups inactivates the drug. Removal of the dichloroacetamido side chain virtually eliminates biological activity, although substitution by a number of groups does not render the drug inactive. Some substituents reduce activity against intact bacteria, but enhance the inhibition of cell-free protein synthesis (1). Apparently, these derivatives are not taken up by the bacterial cell as efficiently as the parent compound.

## 2. Antibacterial Spectrum

Chloramphenicol is a broad-spectrum antibiotic, active against archaeobacteria and both gram-positive and gram-negative bacteria. Organisms that fall within its spectrum include *Bacillus anthracis*, *Brucella abortus*, *Corynebacterium diphtheriae*, *Escherichia coli*, *Hemophilus influenzae*, *Klebsiella pneumoniae*, *Neisseria meningitidis*, *Salmonella* spp., *Serratia marcescens*, *Staphylococcus aureus*, *Streptococcus* spp., and so on (2).

Chloramphenicol inhibits translation in bacterial cells, intact [chloroplasts](#), and intact [mitochondria](#). Some of the toxic side effects associated with chloramphenicol may arise from the effect of the drug on mitochondrial metabolism. Translation by mammalian cells, yeast, and protozoa is not inhibited by chloramphenicol.

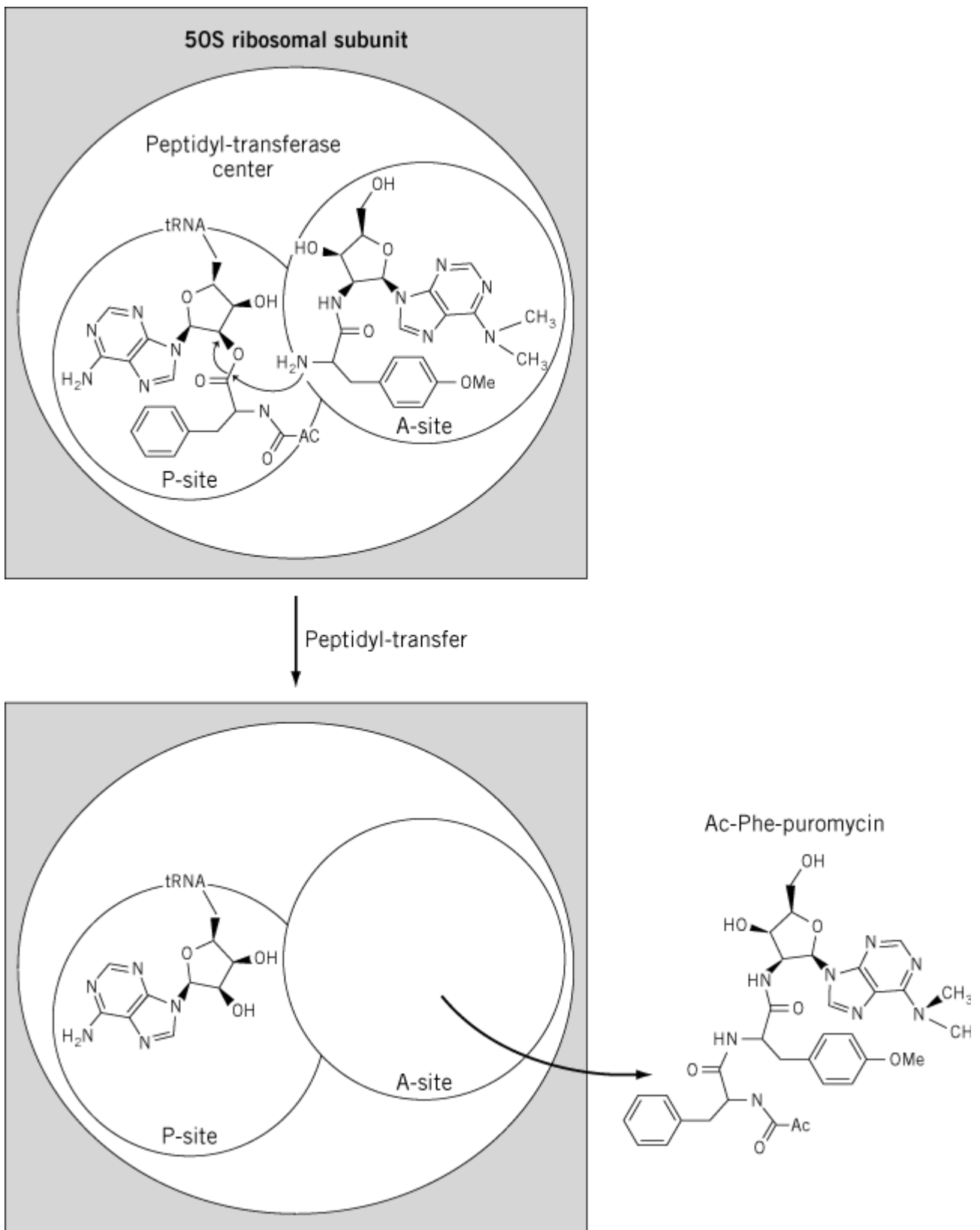
## 3. Mode of Action

Chloramphenicol inhibits bacterial peptidyl-transferase, the enzymatic reaction by which one more amino acid is added to the growing [polypeptide chain](#). During the elongation phase of protein biosynthesis in bacteria, chloramphenicol “freezes” [polysomes](#) on the [messenger RNA](#), thereby fixing the peptidyl-tRNA to the ribosomes (3). Chloramphenicol competes for binding to the 50S subunit of the bacterial ribosome with molecules that specifically bind to the A site of the peptidyl-transferase center (aminoacyl-tRNA and its analogues, eg, CACCA-Leu) and is, therefore, thought to bind next to, or within, the A site of the peptidyl-transferase center (4). In contrast, no competition is observed with substrates that bind at the P site (peptidyl-tRNA analogues, eg, CACCA-LeuAc or AcPhe-tRNA). Furthermore, neither the dinucleotide ApC nor amino acids block the binding of chloramphenicol, whereas both [puromycin](#) and the dinucleotide CpA do (5). Thus, the drug presumably binds in the same region of the A site as the universally conserved CCA part of aminoacyl-tRNA.

The peptidyl-transferase assay in which [peptide bonds](#) are formed between a peptidyl-tRNA and puromycin on the 50S subunit (the “puromycin reaction”) has been used to unravel the mechanism of action of chloramphenicol (6) (Fig. 2). Chloramphenicol blocks this reaction, provided that the

donor peptidyl-tRNA, bound at the P site, is short (eg, AcPhe-tRNA). In this case, chloramphenicol appears to disturb the correct positioning of the nascent peptide during its synthesis, thereby causing the premature release of the peptidyl-tRNA from the ribosome. Chloramphenicol does not, however, affect the puromycin reaction when a larger peptidyl-tRNA—for example, Ac(Phe)<sub>2</sub>-tRNA—is used as donor (7). It is likely, therefore, that chloramphenicol interferes with the binding of the terminal CCA of the aminoacyl-tRNA at the A site and also interferes with the entry of the nascent polypeptide into the “tunnel” that normally guides it away from the peptidyl-transferase center.

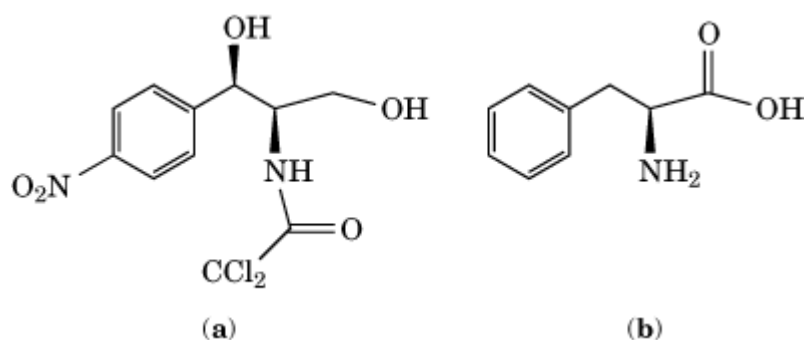
**Figure 2.** The puromycin reaction. A peptidyl-tRNA analogue [eg, Ac[<sup>14</sup>C]Phe-tRNA) bound to the ribosomal P site of 50S ribosomal subunit reacts with puromycin (an aminoacyl-tRNA analogue) bound to the A site of the peptidyl-transferase center. The resulting peptidyl-puromycin (eg, Ac[<sup>14</sup>C]Phe-puromycin) is falling off the 50S subunit and can be separated from the substrate (Ac[<sup>14</sup>C]Phe-tRNA) by extraction.



Chloramphenicol inhibits protein synthesis in bacterial extracts with varying potencies depending on the template employed. In particular, synthesis promoted by poly(U) is markedly more resistant to chloramphenicol than that promoted by poly(C) or poly(A) (8). Synthesis occurring with a natural mRNA (eg, MS2 phage RNA) is usually very sensitive to chloramphenicol. The precise reason for the template-dependence is not known. Chloramphenicol may inhibit poly(U)-dependent poly(Phe)

synthesis less than it does the synthesis of other polypeptides because the structures of D-threo chloramphenicol and L-phenylalanine are similar (Fig. 3). Competition between these molecules during poly(Phe) synthesis may involve steric hindrance between their phenyl groups and, consequently, chloramphenicol may be chased from the ribosome. Nevertheless, this does not explain why the same template dependence is also observed with several other peptidyl-transferase inhibitors that differ chemically from chloramphenicol and, in particular, that do not have a phenyl group (e.g. the group B streptogramins and the macrolides). Further experiments are needed to elucidate this phenomenon in greater detail.

**Figure 3.** Structural similarity of D(-)-*threo* chloramphenicol and L-phenylalanine. Competition between these molecules during poly(Phe) synthesis may involve steric hindrance between their phenyl groups.



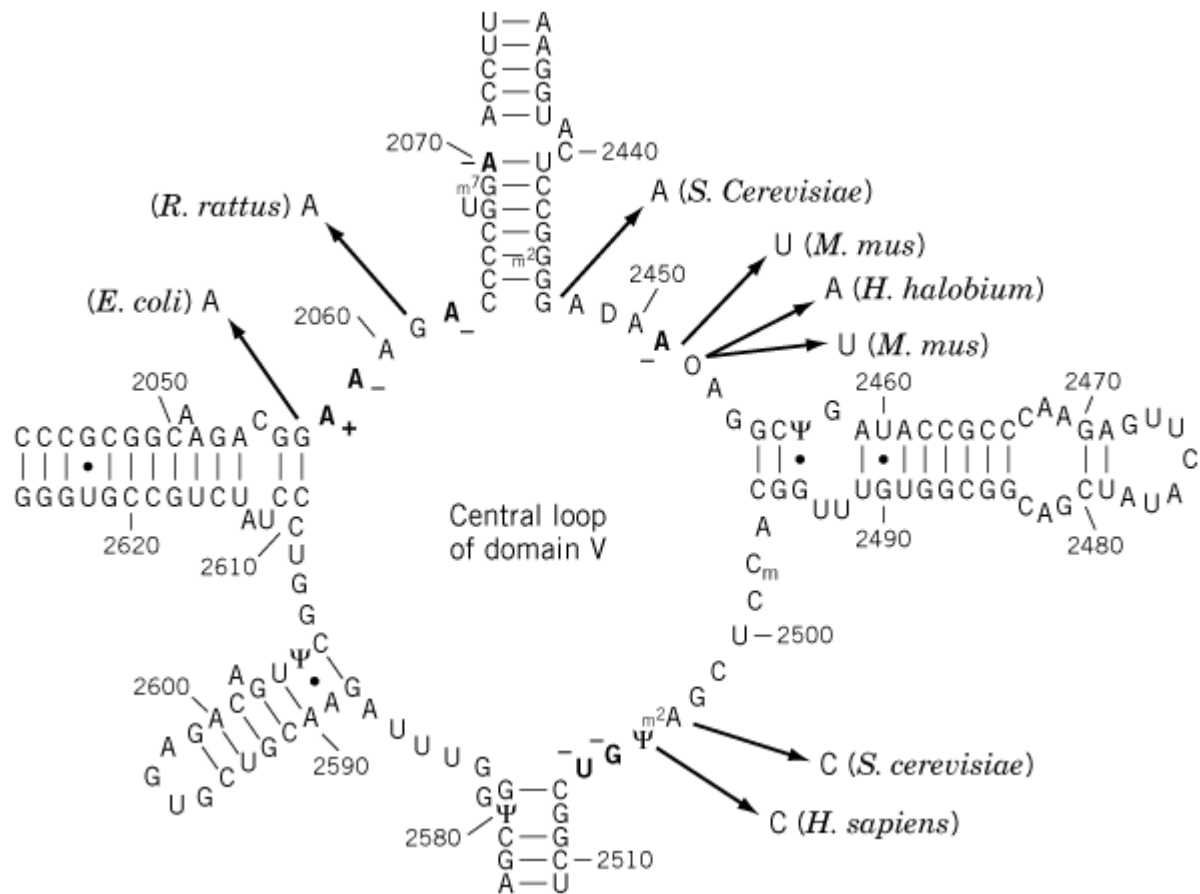
#### 4. Ribosomal Binding Site

Chloramphenicol has a single high-affinity binding site located on the 50S ribosomal subunit ( $K_d = 2 \times 10^{-6}$  M), as well as one or more low-affinity sites, at least one of which is located on the 30S subunit ( $K_d = 2 \times 10^{-4}$  M) (9). A number of 50S proteins were labeled in [affinity-labeling](#) experiments, including L1, L2, L11, L16, L19, and L27 (10). Most of these proteins are located in or next to the peptidyl-transferase center of the ribosome, and several can be affinity-labeled with other antibiotics that interact with this center (**erythromycin**, puromycin, lincosamides, streptogramins).

That the peptidyl-transferase domain is the primary binding site of chloramphenicol was confirmed by RNA [footprinting](#) experiments, in which bases of 23S rRNA were compared in antibiotic-free or chloramphenicol-bound ribosomes for their reactivity to chemical probes. Chloramphenicol protected bases A2059, A2062, A2451, and G2505 and enhanced the reactivity of A2058 (11). Protection of A2070 and U2506 has also been observed (12). All these nucleotides are located in the central loop of domain V of 23S rRNA (Fig. 4). Additionally, several mutations of RNA that confer resistance to chloramphenicol have been identified in the same domain (Fig. 4; see references in Ref. 12).

**Figure 4.** Interaction of chloramphenicol with the central loop of domain V of 23S rRNA. The secondary structure of the central loop of domain V in *E. coli* 23S rRNA is derived from phylogenetic sequence comparisons. Differences in the accessibility of bases to chemical reagents in the presence and absence of chloramphenicol (footprints) are shown in red (–, protection of base; +, enhanced activity). Mutations leading to chloramphenicol resistance are indicated by arrows, and the corresponding organisms are indicated in brackets. All of the eukaryotic mutations were characterized for mitochondrial rRNA. (Figure derived from Refs. 11 and 12.)





## 5. Cold-Shock Induction

Shifting growing *E. coli* from 37°C to 20°C inhibits cellular growth and shuts down protein biosynthesis. To compensate for the stress produced by the temperature shift, the levels of (p)ppGpp drop and the synthesis of a subset of proteins is induced; this induction is known as **cold-shock**. Certain inhibitors of protein synthesis, including chloramphenicol, specifically induce the synthesis of the cold-shock proteins in the absence of a temperature shift (13). Ribosomal inhibitors have, in fact, been classified according to their ability to elicit a cold-shock- or a heat-shock-like response. Inhibitors in whose presence the A site of the peptidyl-transferase center is occupied, either by the antibiotic or by an aminoacyl-tRNA (chloramphenicol, erythromycin, fusidic acid, spiramycin, **tetracycline**), produce a cold-shock-like response. In contrast, inhibitors in whose presence the A site of the peptidyl-transferase remains empty ([kanamycin](#), [streptomycin](#)) produce a heat-shock-like response.

## 6. Chloramphenicol Resistance

### 6.1. Modification of the Bacterial Ribosome

Several mutations in 23S rRNA confer resistance to chloramphenicol. The residues at which these mutations occur (residues G2057, G2061, G2447, A2451, C2452, A2503, U2504; *E. coli* numbering) are located in the peptidyl-transferase loop of 23S rRNA (Fig. 4), which is in accordance with their proposed binding site.

### 6.2. Alteration in Permeability to Chloramphenicol

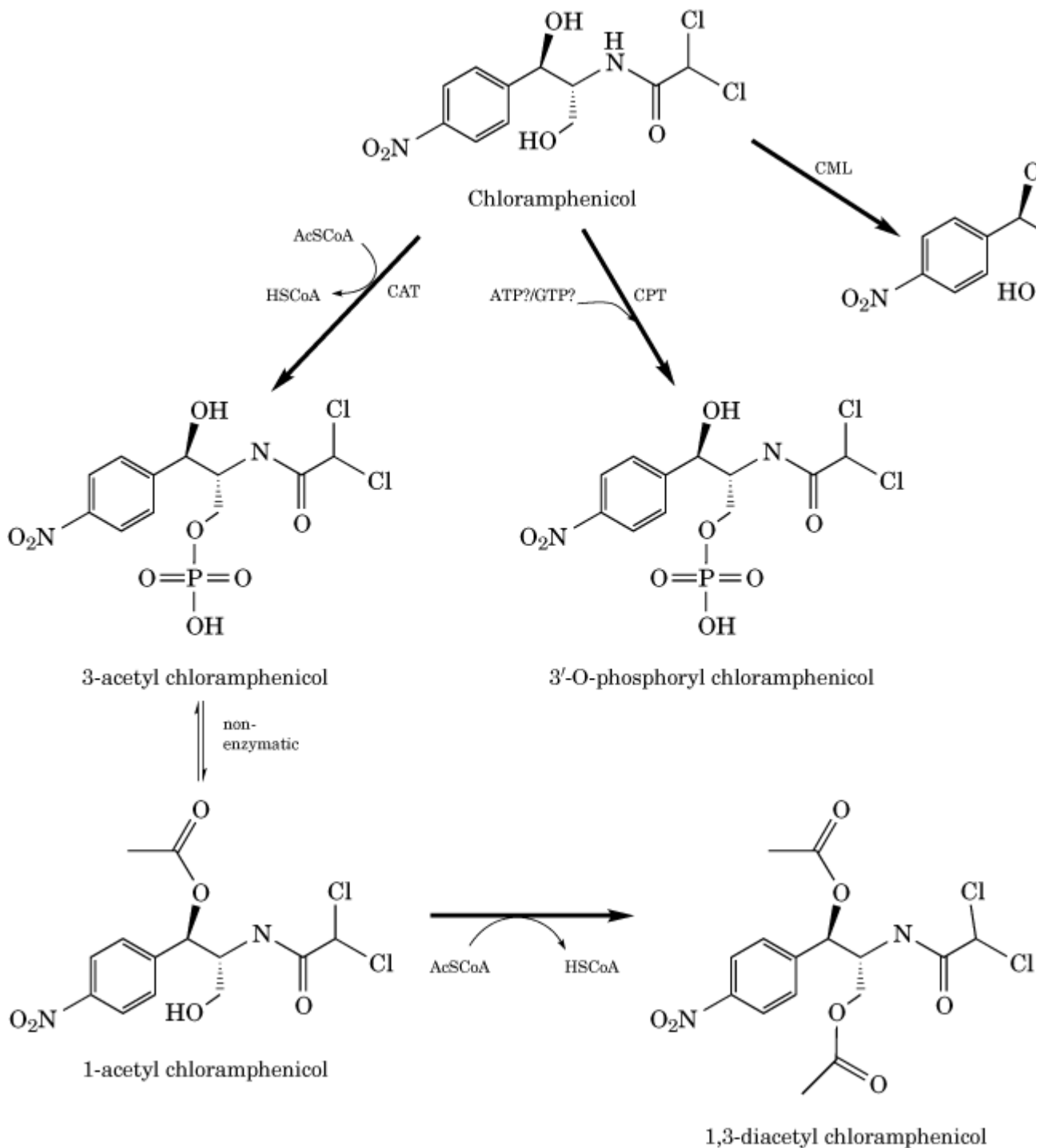
Protein synthesis in *S. venezuelae* cells that have begun to produce chloramphenicol becomes insensitive to the antibiotic. There is no modification of the ribosome in these cells; instead, the antibiotic is excluded from the cytoplasm during biosynthesis. The amino acid sequence deduced

from the gene conferring resistance (*cml*) is markedly similar to those of chloramphenicol-resistance genes from *Streptomyces lividans*, *Rhodococcus fascias*, and other bacteria. These polypeptides contain 12 **hydrophobic** regions characteristic of membrane-associated **transport** proteins and are related to a known family of efflux proteins (14).

### 6.3. Enzymatic Modification of Chloramphenicol

The most common mechanism of chloramphenicol resistance is inactivation by the enzyme chloramphenicol 3'-*O*-acetyltransferase (CAT), encoded by the *cat* gene. The biological and immunological properties of three types of CAT (I, II, III) have been described. The structure of CAT III has been determined by [X-ray crystallography](#), as well as its tertiary complex with chloramphenicol and CoA (15). The enzyme is a trimer; its [active sites](#) lie at each of the interfaces between adjacent subunits. CAT catalyzes the transfer of an acetyl moiety from acetyl-CoA to the 3'-hydroxyl of chloramphenicol. The initial product of the reaction, 3-acetyl chloramphenicol, can undergo nonenzymatic intramolecular rearrangement to 1-acetyl chloramphenicol. The latter compound is a substrate for a second round of enzymatic acetylation at the 3-hydroxyl, which yields 1,3-diacetyl chloramphenicol as the final product (Fig. 5). Neither the monoacetylated nor the diacetylated derivative binds to the bacterial ribosome; both, therefore, are devoid of antibiotic activity.

**Figure 5.** Enzymatic inactivation of chloramphenicol. CAT, chloramphenicol 3'-*O*-acetyltransferase; CPT, chloramphenicol hydrolase.



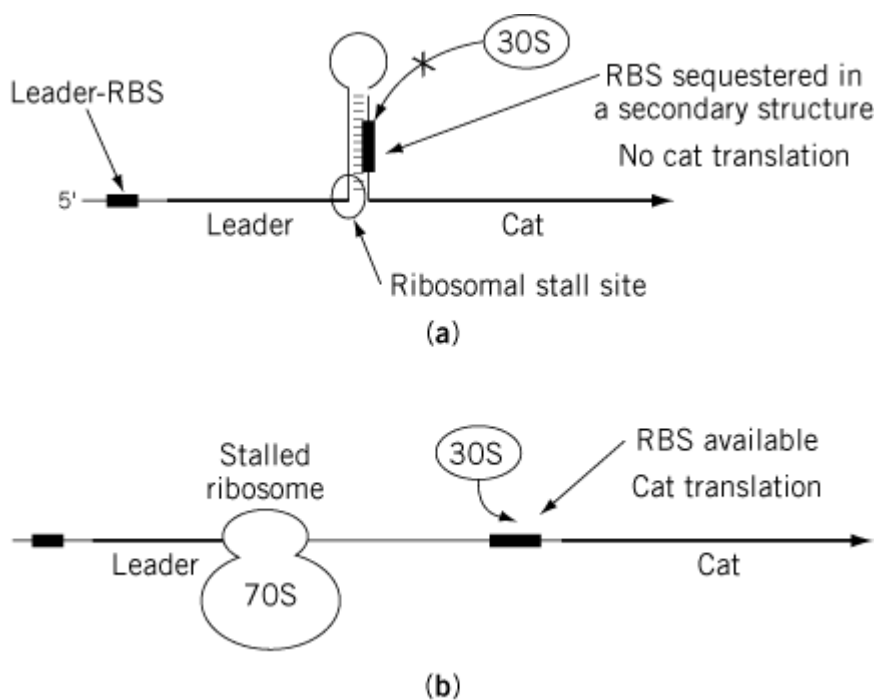
Recently, a second chloramphenicol 3'-*O*-modifying gene was discovered (16). This enzyme, chloramphenicol 3'-*O*-phosphotransferase (CPT) phosphorylates chloramphenicol, yielding 3'-*O*-phosphochloramphenicol, which binds only very weakly to ribosomes and thus does not inhibit translation.

Chloramphenicol is also inactivated by chloramphenicol hydrolase, which hydrolyzes the amide bond (17) (Fig. 5). This enzyme, encoded by the *cml* gene, was first isolated from the mycelium of a chloramphenicol-producing strain of *Streptomyces*. A similar enzyme has been obtained from several other bacteria, including *E. coli*.

#### 6.4. Induction of Chloramphenicol Resistance by Translational Attenuation

Chloramphenicol resistance by enzymatic modification in both Gram-positive and Gram-negative bacteria can be either constitutive or inducible. What is unusual about inducible chloramphenicol resistance is that it occurs by translational attenuation, rather than by transcriptional **attenuation**. The ribosome-binding site for the resistance determinant is sequestered in a secondary “hairpin” structure situated within the encoding mRNA; a short, translated open reading frame, termed the leader, lies upstream of this structure. In the absence of the antibiotic, the secondary structure is maintained, the ribosome does not translate the mRNA that constitutes the hairpin, and the resistance gene is not translated. In the presence of the antibiotic, the ribosome stalls at the leader sequence and the secondary structure relaxes, allowing the *cat* (or *cml*) determinant to be translated (18) (Fig. 6).

**Figure 6.** Chloramphenicol resistance by translational attenuation. RBS, ribosomal binding site; 30S, small ribosomal subunit; 70S, entire bacterial ribosome.



## 7. The Use of CAT in Molecular Biology

CAT is a gene product that can be assayed with specificity and sensitivity, and consequently it is of considerable value to molecular biologists. When the *cat* gene is fused to heterologous regulatory signals, gene expression can be determined by measuring the CAT activity. Eukaryotes lack endogenous CAT activity, and the amount of CAT activity correlates well with the amount of mRNA transcribed, so the CAT assay has proved to be a powerful tool for investigating and quantifying gene expression in eukaryotic cells. Conventional CAT assays involve the incubation of cell extracts with radioactively labeled chloramphenicol, chromatographic separation of the reaction products, and [autoradiography](#). The recent use of **fluorescent**, rather than radioactive, substrates has dramatically speeded up and simplified the CAT assay (19).

### Bibliography

1. O. Pongs (1979) In *Antibiotics* Vol. 1 (F. H. Hahn ed.), Springer-Verlag, Heidelberg, Germany, pp. 26–56.
2. A. A. Yunis (1988) *Annu. Rev. Pharmacol. Toxicol.* **28**, 83–100.
3. E. Cundliffe and K. McQuillen (1967) *J. Mol. Biol.* **30**, 137–143.

4. D. Nierhaus and K. H. Nierhaus (1973) Proc. Natl. Acad. Sci. USA **70**, 2224–2228.
5. M. L. Celma, R. E. Monro, and D. Vazquez (1971) FEBS Lett. **13**, 247–251.
6. R. R. Traut and R. E. Monro (1964) J. Mol. Biol. **10**, 63–72.
7. H.-J. Rheinberger and K. H. Nierhaus (1990) Eur. J. Biochem. **193**, 643–650.
8. Z. Kucan and F. Lipmann (1964) J. Biol. Chem. **239**, 516–520.
9. J. L. Lessard and S. Pestka (1972) J. Biol. Chem. **247**, 6909.
10. B. S. Cooperman, C. J. Weitzmann, and C. L. Fernandez (1990) In *The Ribosome: Structure, Function, & Evolution* (W. E. Hill et al., eds.), American Society for Microbiology, Washington, D.C., pp. 491–501.
11. D. Moazed and H. F. Noller (1987) Biochimie **69**, 879–884.
12. C. Rodriguez et al. (1995) J. Mol. Biol. **247**, 224–235.
13. P. A. Lagosky and F. N. Chang (1981) J. Biol. Chem. **256**, 11651–11656.
14. H. Nikaido and M. H. Sailer Jr. (1992) Science **258**, 936–942.
15. A. G. W. Leslie (1990) J. Mol. Biol. **213**, 167.
16. R. H. Mosher et al. (1995) J. Biol. Chem. **270**, 27000–27006.
17. V. S. Malik and L. C. Viking (1971) Can. J. Microbiol. **17**, 1287–1290.
18. P. S. Lovett (1996) Gene **179**, 157–162.
19. C. K. Lefevre et al (1995) Biotechniques **19**, 488–493.

### Suggestions for Further Reading

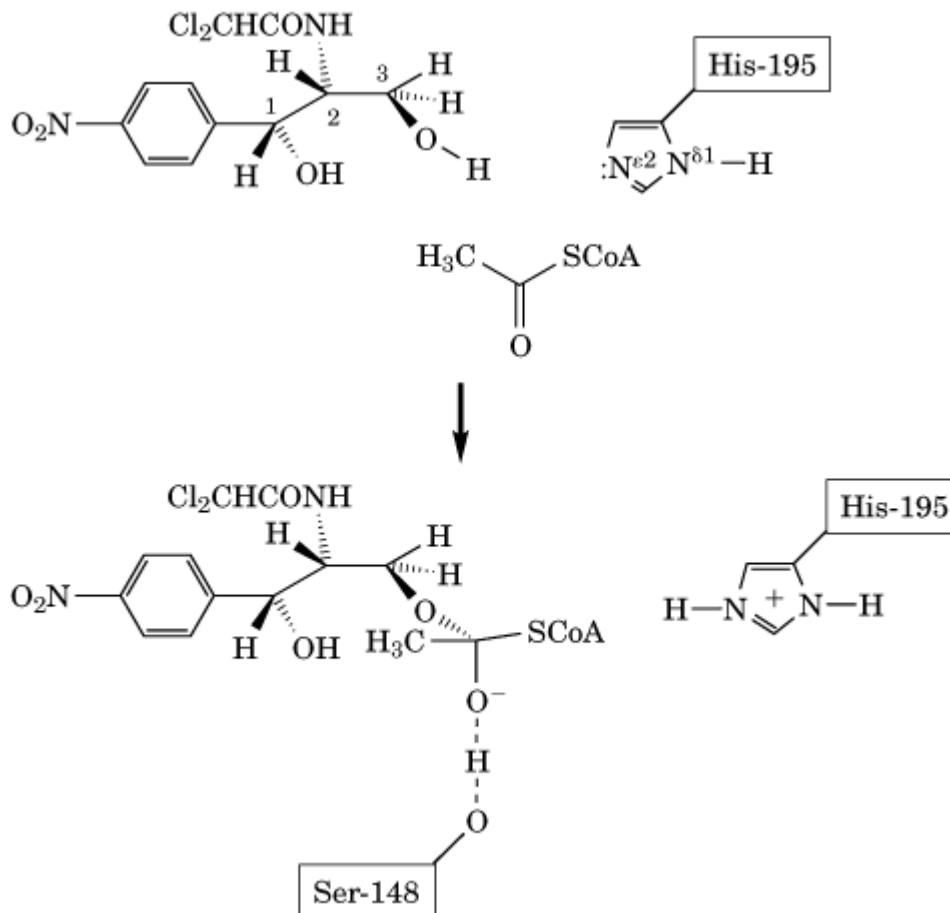
20. E. F. Gale, E. Cundliffe, P. E. Reynolds, M. H. Richmond, and M. J. Waring (1981) Chloramphenicol. In *The Molecular Basis of Antibiotic Action*, 2nd ed. Wiley, New York, pp. 460–468.
21. F. E. Hahn (1983) Chloramphenicol. In *Antibiotics*, Vol. VI (F. H. Hahn, ed.), Springer-Verlag, Berlin, pp. 34–45.
22. L. C. Vining and C. Stuttard (1995) Chloramphenicol. *Biotechnology* **28**, 505–530.

## Chloramphenicol Acetyltransferase

The **enzymic** mechanisms by which bacteria modify biologically active molecules are limited by two constraints: (i) the structure–activity correlations for each class of chemical agent and (ii) the metabolic repertoire available to the bacteria. In the case of [chloramphenicol](#) (Cml), which has a number of vulnerable functional groups, there are several possibilities. Each functional group of Cml (Fig. [1](#)) contributes to its effectiveness as an inhibitor of ribosomal peptidyltransferase activity ([1](#)), and there are examples (reviewed in Ref. [2](#)) of enzyme-mediated resistance to Cml due to dehalogenation, nitro group reduction, hydrolysis of the amide bond, and modification of the hydroxyl groups by phosphorylation ([3](#)) or acetylation ([2](#), [4](#), [5](#)). Nonetheless, after more than four decades of medical and veterinary use, the preponderant enzymic modification mechanism for Cml resistance in bacteria of clinical importance is that of *O*-acetylation of the 3-hydroxyl group, catalyzed by chloramphenicol acetyltransferase (CAT) (Fig. [1](#)).

**Figure 1.** The mechanism of acetylation of the 3-hydroxyl group of chloramphenicol by acetyl-CoA as catalyzed by

CAT. The two substrates and the His195 residue of CAT are shown (**top**). Within the transition state or tetrahedral intermediate (**bottom**), His195 has abstracted a proton from the 3-hydroxyl group, to generate an “oxyanion” intermediate that has attacked the carbonyl of acetyl-CoA. The intermediate and transition state are stabilized by hydrogen bonding with the side-chain hydroxyl group of Ser148 of CAT.



**Genes** for CAT are widespread among [gram-positive](#) and [gram-negative bacteria](#); the *cat* gene in each case is either chromosomal or carried by a mobile genetic element, such as a plasmid or [transposon](#). More important biochemically are the properties of representative variants within the CAT “family” and the pattern of conservation of the [primary structure](#), deduced from nucleotide sequences of its genes, which in turn is related to its three-dimensional [protein structure](#), substrate specificity, and catalysis.

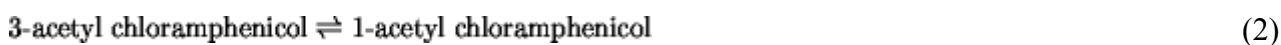
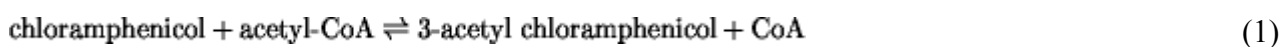
All CAT [polypeptide chains](#) are in the range 24 to 26 kDa and normally exist in solution as compact and very stable *homo*-trimers. However, some CAT variants associate *in vivo* and *in vitro* to give hybrids, functional  $\alpha_2\beta$  and  $\alpha\beta_2$  *hetero*-trimers with physical and catalytic characteristics reflecting the properties of the parental trimers (6). Because natural isolates of Cml-resistant bacteria may occasionally harbor more than one *cat* gene, it is possible in such instances that the intracellular CAT pool will include both parental and hybrid trimers.

The extent of structural variation within the CAT family can be appreciated from a comparison of the deduced primary structures for the products of known *cat* determinants, which yields a lower limit of 28% identity for the most divergent pair of known sequences. Only ~11% of the amino acid residues appear to be identical in all CAT variants, comprising not only residues with side chains that are involved in catalysis and substrate binding, but also those that contribute to critical structural elements necessary for the precise folding and stable packing of the polypeptide chains. A reference point for understanding the variety in primary structures of CAT variants is the type III enzyme

(CAT<sub>III</sub>), for which a wealth of information is available, including the [tertiary structure](#) at high resolution ([7, 8](#)) and the structural determinants for the binding of each substrate ([2, 4, 5, 7, 8](#)). However, the type I enzyme (CAT<sub>I</sub>) may be the most widely distributed variant, because it is specified by many “F-like” R plasmids of gram-negative bacteria, by transposon *Tn9* ([9](#)), and by promoter-less *cat* “cassettes,” constructed in the laboratory for insertion “downstream” of the noncoding sequences of other genes, to study the regulation of their expression (see [Reporter Genes](#)). Type II CAT (CAT<sub>II</sub>), notable among “enteric” CATs for its particular sensitivity to inhibition by reagents that react with [thiol groups](#) and by its association with *Haemophilus influenzae* ([10](#)), is less commonly encountered. CAT<sub>I</sub> has the remarkable properties of a high affinity for triphenylmethane dyes, such as crystal violet (reviewed in Ref. [2](#)), and the ability to bind a steroidal antibiotic (fusidic acid) both tightly and specifically ([11](#)). The latter property is sufficient to confer resistance to fusidate in mutant strains of *Escherichia coli* that were selected for their sensitivity to the antibiotic prior to the introduction of the gene for CAT<sub>I</sub>. A plausible mechanism is the effective sequestration of fusidate by high levels of CAT<sub>I</sub>, thereby impeding access of the antibiotic to its cellular target, ribosomal [elongation factor](#) G. The structural basis for the curious binding of fusidate, which is competitive with respect to Cml, has been deduced by [protein engineering](#) and [X-ray crystallography](#) ([11](#)). Of the eight residues in the Cml binding pocket of CAT<sub>III</sub> that differ from those of CAT<sub>I</sub>, only four appear to be responsible for the former's low affinity ( $K_i = 279 \mu\text{M}$ ) for fusidate. Replacement of each of the four with their counterparts in CAT<sub>I</sub> is sufficient to confer upon CAT<sub>III</sub> an affinity for fusidate ( $K_i = 5.4 \mu\text{M}$ ) that approaches that of wild-type CAT<sub>I</sub> ( $K_i = 1.5 \mu\text{M}$ ).

The catalytic machinery of CAT<sub>III</sub> comprises a number of amino acid side chains (as well as [backbone](#) atoms and ordered [water](#) molecules) that make precise contacts with one another or with a substrate molecule. The homo-trimers have rotational symmetry at the molecular threefold axis, and consequently have three identical active sites, so only atoms in a single monomer need be addressed. Nonetheless, each of the structurally and functionally equivalent [active sites](#) lies deep in the interfacial clefts between subunits. Central to catalysis ([12](#)) is His195, which arises from one face of each cleft to supply the general base ( $\text{N}^{\text{e}2}$  in Fig. [1](#)) to deprotonate the C3 hydroxyl of Cml, producing an “oxyanion” intermediate that in turn attacks the carbonyl (C2) carbon of acetyl CoA to yield a tetrahedral intermediate (Fig. [1](#)). Essential for the stabilization of the latter, *en route* to the [transition state](#) for the reaction, and confirmed by [site-directed mutagenesis](#) ([13](#)), is a negatively charged [hydrogen bond](#) (Fig. [1](#)) between the oxyanion and the hydroxyl of Ser148, another residue conserved in all CATs. A neighboring participant in catalysis is Thr174, also conserved, which is hydrogen-bonded to a water molecule that in turn probably makes two hydrogen bonds with the putative tetrahedral intermediate (Fig. [1](#)), one to the 1-hydroxyl of Cml and the other to the 3-oxygen of the intermediate. Two additional conserved residues (Arg18 and Asp199) facilitate catalysis via a network of hydrogen bonds with His195, anchoring the side chain of the latter in a novel conformation that allows it to fulfill its general base role ([7, 8, 13-15](#)).

Chloramphenicol becomes acetylated on its 1-hydroxyl group also, albeit at a rate much slower than the 3-acetyl derivative is generated ([16](#)). The equations below, all of which are reversible, indicate the transformations involved in the two reactions, indicating that in both cases CAT acetylates only the 3-hydroxyl:





Reaction 2 is a nonenzymic acetyl migration, slow and reversible, which yields at equilibrium a mixture of mono-acetyl products. The 1-acetyl Cml so formed is available for a second round of enzymic acetylation at the C3 position, yielding 1,3-diacetyl Cml. Reaction 3 is ~150-fold less efficient than reaction 1, almost certainly due to an unfavorable “fit” of the substrate at the active site because of the bulky 1-acetyl substituent (17). In any case, the sluggish final step (reaction 3) is of little microbiological significance, because both mono-acetyl derivatives of Cml are already devoid of significant antimicrobial activity, making the rate of reaction 1 the prime determinant of the Cml-resistance phenotype.

In summary, the precise geometry and chemical properties of both substrate binding sites (for Cml and acetyl CoA) and of the catalytic center of CAT<sub>III</sub> each contribute to its extraordinary efficiency, with a turnover number of 600 s<sup>-1</sup> (25°C) with a  $K_m$  for Cml of 12 μM (~4 μg/mL), reassuringly close to the concentrations at which it inhibits most bacteria of clinical importance. A derived kinetic parameter, the so-called specificity constant ( $k_{cat}/K_m$ ), which combines a measure of substrate affinity with one for catalytic competence, is actually the second-order rate constant for productive collisions of an enzyme with its substrate(s). The value for CAT<sub>III</sub> and Cml ( $5 \times 10^7 \text{ s}^{-1}\text{M}^{-1}$ ) approaches that of well-characterized enzyme reactions that are limited by diffusion of the reactants (typically  $10^8$  to  $10^9 \text{ s}^{-1}\text{M}^{-1}$ ) and hence at the limit of evolutionary development. By such criteria, CAT<sub>III</sub> has evolved to a state approaching “perfection” in biological catalysis (18), wherein a fine balance has been struck between rate acceleration ( $k_{cat}$ ) and specificity (and affinity) for substrate ( $K_m$ ). Although it is not clear how the specificity and catalytic efficiency of CAT<sub>III</sub> (and related variants) have evolved, the acetyltransferase (E2p) of the pyruvate dehydrogenase complex, which generates acetyl-CoA for central metabolism, has a three-dimensional structure that is virtually identical to that of CAT, as well as the same mechanism of catalysis, but the two proteins have very few identities in primary structure—only those involved with the active site (19).

It is of interest that there is a large family of “xenobiotic” *O*-acetyltransferases (XATs) with a range of specificities for natural products; several of these enzymes have a low affinity for Cml and hence were first detected as effectors of low-level resistance to the antibiotic (20). All appear to be trimeric but have no sequence homologies with members of the *bona fide* CAT family described above. One such XAT (with CAT activity) has been studied by X-ray crystallography (21) and shown to have a tertiary structure quite different from that of CAT, but the catalytic mechanism may well involve general base catalysis involving a conserved histidine residue.

## Bibliography

1. E. F. Gale et al. (1981) *The Molecular Basis of Antibiotic Action*, 2nd ed., Wiley, London, pp. 462–468.
2. W. V. Shaw (1983) *CRC Crit. Rev. Biochem.* **14**, 1–46.
3. R. H. Mosher et al. (1995) *J. Biol. Chem.* **27**, 27000–27006.
4. W. V. Shaw (1992) *Sci. Progress (Oxford)* **76**, 565–580.
5. W. V. Shaw and A. G. W. Leslie (1991) *Annu. Rev. Biophys. Biophys. Chem.* **20**, 363–386.
6. P. J. Day, I. A. Murray, and W. V. Shaw (1995) *Biochemistry* **34**, 6416–6422.
7. A. G. W. Leslie (1990) *J. Mol. Biol.* **213**, 167–186.
8. A. G. W. Leslie, P. C. E. Moody, and W. V. Shaw (1988) *Proc. Natl. Acad. Sci. USA* **85**, 4133–4137.
9. N. E. Alton and D. Vapnek (1979) *Nature* **282**, 864–869.
10. I. A. Murray, J. V. Martinez-Suarez, T. J. Close, and W. V. Shaw (1990) *Biochem. J.* **272**, 505–



11. I. A. Murray et al. (1995) *J. Mol. Biol.* **254**, 993–1005.
12. A. Lewendon et al. (1994) *Biochemistry* **33**, 1944–1950.
13. A. Lewendon, I. A. Murray, W. V. Shaw, M. R. Gibbs, and A. G. W. Leslie (1990) *Biochemistry* **9**, 2075–2080.
14. A. Lewendon et al. (1988) *Biochemistry* **27**, 7385–7390.
15. A. Lewendon and W. V. Shaw (1993) *J. Biol. Chem.* **268**, 20997–21001.
16. J. Ellis, C. R. Bagshaw, and W. V. Shaw (1995) *Biochemistry* **34**, 16852–16859.
17. I. A. Murray et al. (1991) *Biochemistry* **30**, 3763–3770.
18. W. J. Albery and J. R. Knowles (1976) *Biochemistry* **15**, 5631–5640.
19. A. Mattevi et al. (1993) *Biochemistry* **32**, 3887–3901.
20. I. A. Murray and W. V. Shaw (1997) *Antimicrob. Agents Chemother.* **41**, 1–6.
21. T. W. Beaman, M. Sugantino, and S. L. Roderick (1998) *Biochemistry* **37**, 6689–6696.

## Chloroplast

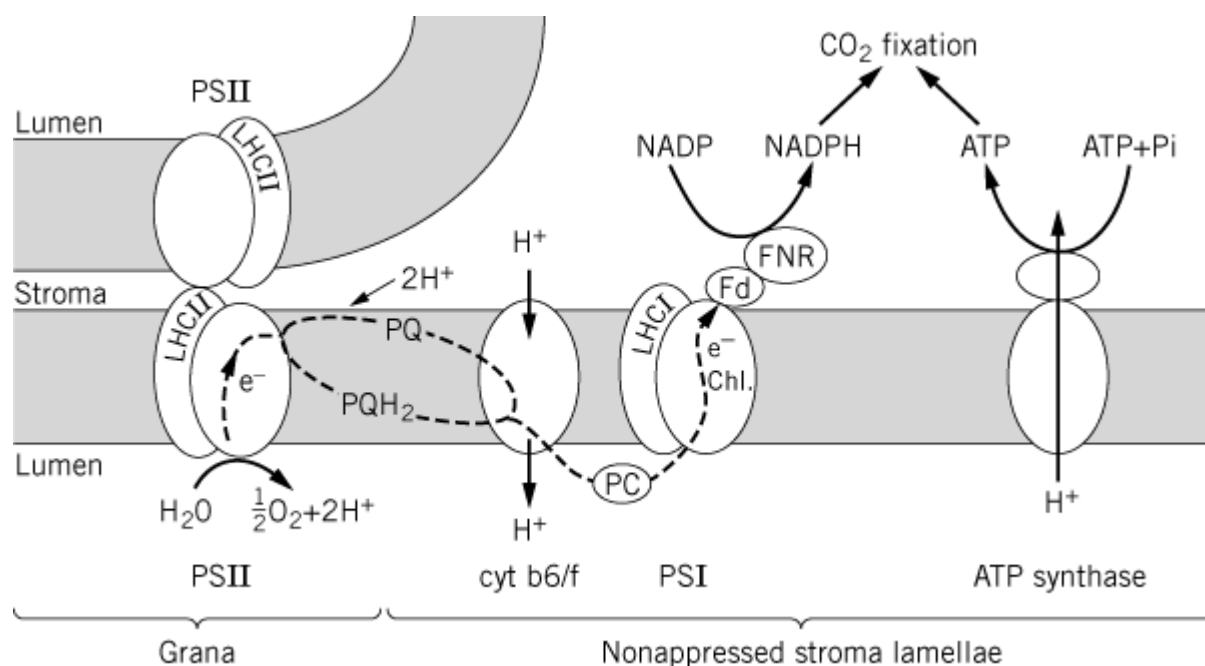
A distinctive feature of chloroplasts of **plants** and **algae** is their extensive, internal, green, chlorophyll-containing [membrane](#) system, called thylakoid membranes, where the primary reactions of [photosynthesis](#) occur. This system of [photosynthetic reaction centers](#) converts light energy into chemical energy, which is used to drive cellular metabolism. Besides their important role in photosynthesis, chloroplasts are also involved in several biochemical pathways, such as the biosynthesis of [amino acids](#), fatty acids, tetrapyrroles including chlorophyll and heme, carotenoids, isoprenoids and pyrimidines. Chloroplasts are also involved in carbon metabolism and in nitrogen and sulfur assimilation ([1](#)). Like [mitochondria](#), chloroplasts possess their own genetic system, which cooperates closely with the [nucleus](#) in biosynthesizing numerous organellar components. Chloroplasts represent one type of plastid derived from colorless proplastids in the meristematic cells of plant leaves and shoots, which have only a rudimentary internal membrane system ([1](#)). Light profoundly affects the development of proplastids. They differentiate into chloroplasts in the presence of light, whereas in its absence they differentiate into etioplasts, which lack chlorophyll and contain a prolamellar body. Upon subsequent illumination, the prolamellar body gives rise to lamellae of the thylakoid membrane. Depending on the plant tissues, the developmental stage, and the environmental conditions, proplastids also differentiate into chromoplasts in petals or fruits, into leucoplasts in roots, or into amyloplasts in tubers in which starch is accumulated. Proplastids also develop into elaioplasts in glands, certain fruits and seeds, where they are involved in synthesizing lipids, terpenoids, carotenoids, and carbohydrates. Although these various plastid forms have rather distinct morphologies, plastid differentiation is reversible to a large extent, because chloroplasts develop from leucoplasts or amyloplasts, and viceversa. During transitions from chloroplasts to the other plastid forms, the expression of most organellar genes is reduced, whereas specific nuclear genes encoding plastid proteins are activated ([1](#), [2](#)). An important point is that all plastid types contain an internal membrane system that is crucial for their interconversion.

### 1. Thylakoid Membranes and the Photosynthetic Apparatus

The internal thylakoid membrane system consists of appressed and non-appressed flattened membrane vesicles, called grana and stroma lamellae, respectively (Fig. [1](#)). The primary reactions of photosynthesis are catalyzed by four major protein-pigment complexes of the thylakoid membrane:

(i) photosystem II and (ii) photosystem I, and their associated chlorophyll antennae, (iii) the [cytochrome b6/f complex](#), and (iv) the [ATP synthase](#) (Fig. 1; see also [Photosynthesis](#)). Briefly, light energy is captured by the antennae and channeled to the reaction centers of photosystem II and photosystem I. The energy is used to energize an electron in chlorophyll and to create a stable charge separation across the membrane. This triggers a series of oxido-reductions along the photosynthetic **electron-transfer chain**. At one end of this chain, [water](#) is oxidized by photosystem II with concomitant evolution of oxygen and release of protons into the lumen. Then electrons are transferred to plastoquinone, to the cytochrome b6/f complex, which acts as a **proton pump**, and to the soluble electron carrier plastocyanin in the thylakoid lumen. At the other end of the chain, photosystem I oxidizes plastocyanin upon light absorption and transfers electrons to [ferredoxin](#) and then to NADP to form NADPH. The resulting pH gradient is used by the fourth complex, ATP synthase, to produce ATP on the stromal side. This enzyme also functions in the opposite direction by hydrolyzing ATP to pump protons into the thylakoid lumen and thus generate a pH gradient. Because the abundance of the thylakoid membrane complexes facilitates their biochemical analysis and because the state of the redox cofactors is monitored readily by **spectroscopic** techniques, the thylakoid membrane has been studied intensively and represents one of the best-studied membrane systems.

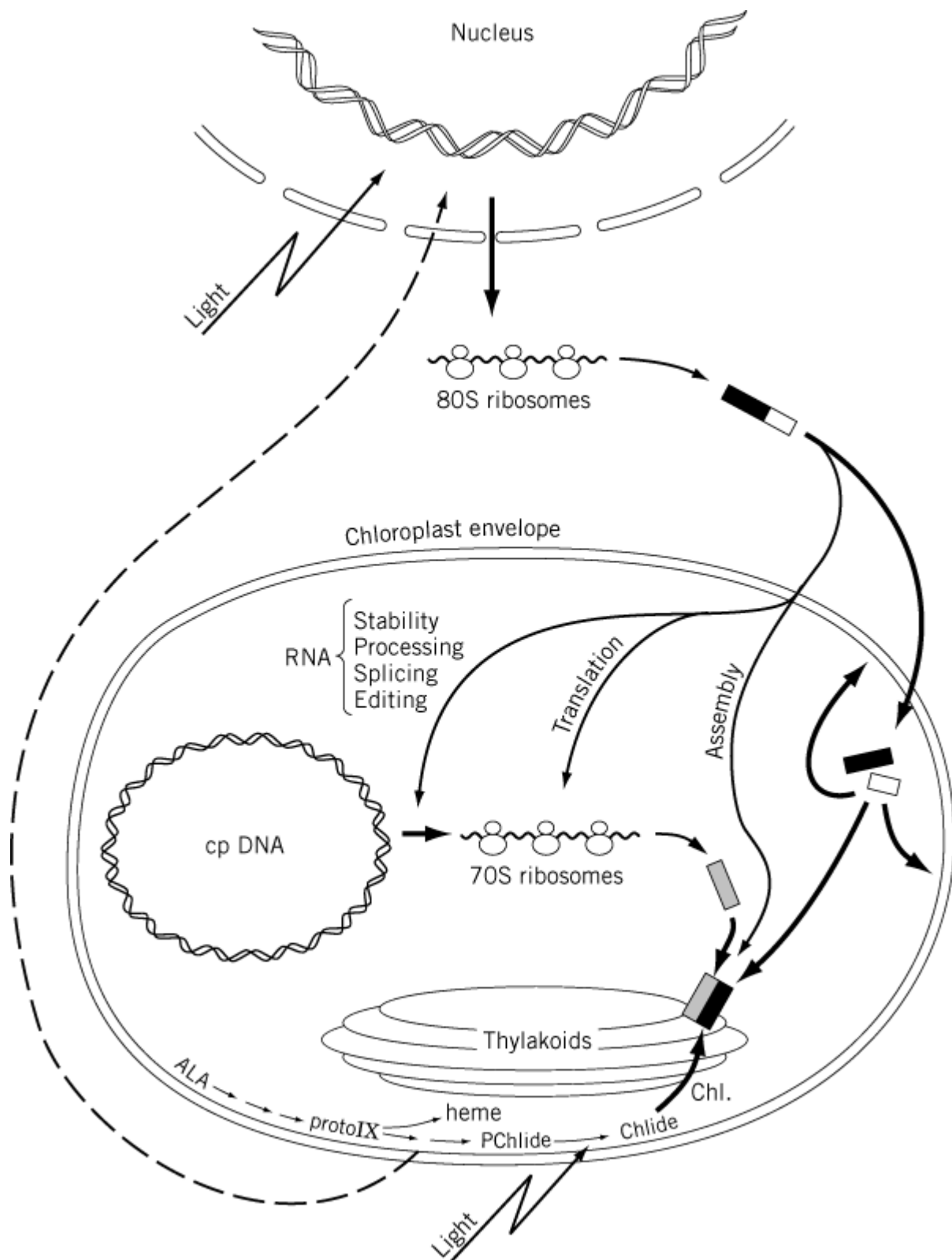
**Figure 1.** Photosynthetic complexes in the thylakoid membrane of chloroplasts. PSII (photosystem II) is located within the appressed grana region, whereas PSI (photosystem I) is located within the nonappressed stroma lamellae. The photosynthetic electron transfer chain is shown starting with water as electron donor to PSII, to plastoquinone (PQ), to the cytochrome b6/f complex (cytb6/f), to the soluble electron transfer protein plastocyanin (PC), to PSI, to ferredoxin (Fd), to ferredoxin-NADP oxidoreductase (FNR), and to NADP as final electron acceptor. Electron flow is coupled to proton translocation into the lumen. The resulting pH gradient across the thylakoid membrane drives ATP synthesis. Both ATP and NADPH are used for CO<sub>2</sub> fixation.



Each of the four photosynthetic complexes contains numerous protein subunits, some of which are encoded by the chloroplast [genome](#), whereas others are encoded by the nuclear genome (Table 1; Fig. 2). The two principal reaction center polypeptides of photosystems I and II are highly **hydrophobic** and contain 11 and 5 transmembrane [a-helices](#), respectively, to which most of the redox cofactors and several chlorophylls are bound with an asymmetrical distribution across the thylakoid membrane. This asymmetry is crucial for the vectorial electron transport in the membrane.

The distribution of these complexes is unequal between the appressed (grana) and nonappressed thylakoid membrane regions (3). Photosystem II is localized predominantly in the grana regions, whereas photosystem I and the ATP synthase complex are found exclusively in the nonappressed regions. The cytochrome b6/f complex is present in both the grana and nonappressed regions. Destacking and restacking of thylakoid membranes are induced experimentally by decreasing and then increasing again the cation concentration. Remarkably, the lateral segregation of the photosynthetic complexes between grana and stromal membranes is lost upon destacking because of random mixing, but it is restored upon restacking the membranes (4).

**Figure 2.** Biosynthesis of the photosynthetic apparatus and protein traffic in the chloroplast. Photosynthetic complexes consist of nucleus- and chloroplast-encoded subunits. The former are synthesized as precursors on cytosolic 80S ribosomes and targeted to the chloroplast. Upon import into the organelle, the N-terminal stromal transit peptide domain is cleaved, and the protein is directed to the stroma, to the envelope, or to the thylakoids. In the latter case, the protein contains an additional cleavable thylakoid targeting domain. Several posttranscriptional steps in the chloroplast, such as RNA stability, processing, splicing, editing, and translation, plus the assembly of protein complexes, require the action of numerous nucleus-encoded factors. Chlorophyll, the major pigment of the thylakoid membrane is synthesized entirely in the chloroplast. Synthesis starts from d-aminolevulinic acid (ALA) and involves several steps in common with heme biosynthesis until protoporphyrin IX (proto IX). One of the last steps of chlorophyll synthesis, conversion of protochlorophyllide (Pchl<sub>id</sub>) to chlorophyllide (Chl<sub>id</sub>), requires light in land plants. Chlorophyll synthesis is tightly coordinated with the synthesis of its apoproteins. Expression of nuclear genes of photosynthetic proteins is strongly stimulated by light. Some of the chlorophyll precursors influence, directly or indirectly, expression of nuclear genes involved in photosynthesis.



**Table 1. Informational Content of Chloroplast DNA from Land Plants and Green Algae**

**Number in Chloroplast**

| Genes Involved in                                       | (Nucleus) <sup>a</sup> |
|---|------------------------|
| Photosynthesis  |                        |
| Photosystem II  | 14 (5)                 |
| Photosystem I   | 5 (8)                  |
| Cytochrome b6/f complex                                 | 5 (2)                  |
| ATP synthase  | 6 (3)                  |
| Ribulose biphosphate<br>carboxylase/oxygenase           | 1 (1)                  |
| NADH dehydrogenase <sup>b</sup>                         | 11                     |
| Light-independent chlorophyll synthesis<br><sup>c</sup> | 3                      |
| CO <sub>2</sub> uptake <sup>d</sup>                     | 1                      |
| Protein Biosynthesis                                    |                        |
| Ribosomal RNAs  | 4                      |
| Ribosomal proteins                                      | 20                     |
| tRNAs   | 30–31                  |
| RNA polymerase subunits                                 | 4                      |
| Other Functions   |                        |
| ClpP subunit of ATP-dependent protease                  | 1                      |
| ORFs of unknown function                                | 10–20                  |

<sup>a</sup> The number of chloroplast-encoded subunits of photosynthetic complexes. The numbers in parenthesis refer to the number of nucleus-encoded subunits.

<sup>b</sup> These genes are found only in the chloroplast genomes of land plants.

<sup>c</sup> The number of chloroplast genes involved in light-independent chlorophyll synthesis and CO<sub>2</sub> uptake. These genes are found only in the chloroplast genomes of gymnosperms, liverwort, and Chlamydomonas.

<sup>d</sup> This gene is involved in CO<sub>2</sub> uptake.

Other dynamic changes in thylakoid membrane organization occur when plants and algae are subjected to light of different wavelengths that is preferentially absorbed by either photosystems II or I. Under these conditions, part of the chlorophyll antenna is displaced from one photosystem to the other, so as to achieve balanced light absorption and hence optimal functioning of the two photosystems (5). Thus when photosystem II is preferentially activated, the plastoquinone pool is reduced. This leads to the activation of a protein kinase associated with the cytochrome b6/f complex, to the phosphorylation of the [light-harvesting complex](#) (LHCII), and to the concomitant movement of part of the photosystem II antenna to photosystem I. If photosystem I is preferentially activated, LHCII is dephosphorylated, and it returns to photosystem II in the grana regions.

### 1.1. Stroma and Carbon-Fixation Cycle

The ATP and NADPH produced by the primary light reactions of photosynthesis are used as sources of energy and reducing power to drive the reactions of the carbon fixation cycle, which convert CO<sub>2</sub> into glyceraldehyde 3-phosphate, a precursor to sugars, amino acids, and fatty acids. Although these reactions are also called the “dark reactions,” the enzymes involved are inactivated in the dark and

need to be reactivated by light through the reducing power generated by photosynthesis. The key reaction, which involves converting one atom of inorganic carbon, as CO<sub>2</sub>, into organic carbon, is catalyzed by the enzyme **ribulose 1,5 biphosphate carboxylase** (Rubisco), a large stromal enzyme that works only sluggishly (6). Therefore, it is required in large amounts and is thought to be the most abundant protein on earth. This enzyme also has an oxygenase activity, which predominates if the concentration of CO<sub>2</sub> is low. Under these conditions it catalyzes the first step of a pathway called photorespiration, which ultimately liberates CO<sub>2</sub> and thereby reverses the photosynthetic reaction. In addition, the stroma includes a large number of proteins involved in several important metabolic pathways (amino acid and fatty acid synthesis, sulfur and nitrogen assimilation). The chloroplast transcription and translation systems are also contained in this compartment.

## 1.2. Chloroplast DNA and Its Informational Content

Chloroplasts together with mitochondria, are the only cellular organelles containing their own apparatus for **protein biosynthesis**. It consists of chloroplast **DNA**, **RNA polymerase**, **enzymes** involved in RNA metabolism, **ribosomes**, **transfer RNA**, and several **translation** factors. Chloroplast ribosomes resemble those of **bacteria** and have similar ribosomal RNA and proteins and sensitivity to a similar spectrum of **antibiotics**. The circular chloroplast DNA molecules range in size between 70 kb and 400 kb (7) and are present in about 100 copies per chloroplast. A typical mesophyll cell contains close to 100 plastids and thus about 10,000 chloroplast DNA circles. A dozen chloroplast genomes from several vascular plants and algae have been sequenced. These sequences have revealed the existence of about 120 chloroplast genes in plants and green algae. They include about 50 genes that encode components of the transcriptional apparatus (subunits of RNA polymerase) and of the translational apparatus (ribosomal RNA, ribosomal proteins, transfer RNA, and translation factors). About 40 genes are involved in photosynthesis, and they encode some of the subunits of photosystems I and II, the cytochrome b6/f complex, ATP synthase and Rubisco (see Table 1). The other subunits of these complexes are encoded by the nuclear genome, translated on cytosolic ribosomes, and imported posttranslationally into the chloroplast. The genes involved in the plastid protein synthesizing system and in photosynthesis have been conserved during evolutionary **divergence** of the chloroplast genomes of plants and green algae. The remaining chloroplast genes, however, have not been universally conserved. Eleven genes encoding subunits of NADH dehydrogenase are present in the chloroplasts of plants, but not in algae. Whereas the role of the mitochondrial NADH dehydrogenase in respiration is well understood, the function of the chloroplast enzyme has not yet been elucidated. It could be involved in a chlororespiratory pathway by reducing the plastoquinone pool in the dark, which is ultimately oxidized by molecular oxygen via unknown redox components (8).

In most plants, one of the last steps of the chlorophyll synthesis pathway, the conversion of protochlorophyllide into chlorophyllide, is light-dependent. In green algae and **gymnosperms**, an alternative light-independent pathway for chlorophyll synthesis is mediated by three chloroplast genes that are absent in **angiosperms**. The sequences of chloroplast genomes have revealed additional genes whose functions are still unknown.

The chloroplast genomes of nongreen algae contain twice as many genes as those of higher plants. Additional genes include those required for photosynthesis that are nucleus-encoded in plants and green algae, genes involved in the synthesis of fatty acids, amino acids, and pigments, genes required for **protein folding** and transport, and additional genes of unknown function (9). The smallest plastid genome identified, that of the white parasitic plant *Epifagus virginiana*, is only 70 kbp in size. It has lost all the genes involved in photosynthesis, and the remaining genes encode mostly components of the plastid protein synthesizing system.

It is generally admitted that plastids originated as the result of an endosymbiotic event in which a **prokaryotic** photosynthetic organism, probably similar to a **cyanobacterium**, invaded a primitive **eukaryotic** cell. Strong support for this endosymbiotic hypothesis arises from the considerable similarity between the transcriptional and translational systems of prokaryotes and plastids. It is

thought that during evolution genetic information from the intruder was gradually lost and transferred to the nucleus of the host. The question thus arises why chloroplast DNA has been maintained. One possibility is that this evolutionary plastid genome size reduction is still in progress and has not yet reached its final stage. Another possibility is that the plastid protein synthesizing apparatus is essential for synthesizing the large hydrophobic polypeptides of the photosynthetic reaction centers, which cannot be translocated across the plastid envelope membrane. A third recently advanced hypothesis is that the presence of the plastid protein synthesizing system is essential to allow a rapid response of plastid gene expression to environmental changes (10).

### 1.3. Chloroplast Gene Expression

Two distinct RNA polymerases are present in the chloroplasts of higher plants. One is similar to its bacterial homologue, and its subunits are encoded by chloroplast genes. This enzyme transcribes primarily genes involved in photosynthesis, which are expressed at a high level. The second plastid RNA polymerase is nucleus-encoded and is required for expressing the nonphotosynthetic plastid functions necessary for plant growth (11). Many chloroplast genes are organized in large transcription units. These units are transcribed into large precursor transcripts, which then are processed into individual [messenger RNA](#) (mRNA) molecules. Chloroplasts contain RNA **splicing** systems, because several plastid genes contain **introns**, mostly group II and group I, which have a characteristic secondary structure (12). These introns have also been found in mitochondrial genes, and some of them are **self-splicing**. Splicing in the chloroplast is rather complex, as in the case of the *psaA* gene encoding one of the reaction center polypeptides of photosystem I in the green alga *Chlamydomonas*. This gene consists of three coding regions (**exons**) that are widely separated on the chloroplast genome and are flanked by group II intron sequences (13). They are transcribed individually, and maturation of the *psaA* mRNA depends on two trans-splicing reactions in which the separate transcripts of the three exons are spliced together. A particularly intriguing feature is that one of the introns is split into three parts (14). This has interesting evolutionary implications because it is thought that group II introns represent the precursors of nuclear introns and their associated splicing factors. In this view, the split chloroplast intron may represent an intermediate between group II and nuclear introns. The chloroplast genetic system has evolved at a rather slow rate and could have therefore maintained some ancient gene organization.

Another unusual feature of chloroplast RNA metabolism is RNA **editing** in vascular plants (15). Editing in chloroplasts is a posttranscriptional process in which specific C residues of a primary transcript are changed to U. Editing has important implications for interpreting DNA genomic sequence data. As an example, an ACG triplet may be edited to AUG, thereby creating a new initiation codon, which could not be identified in the DNA sequence. Alternatively, an editing event may change an internal codon and thus change the corresponding amino acid predicted by the DNA sequence. Therefore, sequencing of chloroplast genomes may not allow identifying of all of the plastid genes.

Because the subunits, redox cofactors, and pigments of photosynthetic complexes are synthesized by two distinct genetic systems, the process has to occur in a coordinated way (Fig. 2). Genetic studies with *Chlamydomonas* and **maize** have indeed revealed the existence of highly complex interactions between nucleus and chloroplast (16). A large number of nuclear genes are involved in chloroplast gene expression. They encode factors targeted to the chloroplast that act at different posttranscriptional steps, such as RNA processing, RNA stability, RNA splicing, translation, and the assembly of photosynthetic complexes. Light strongly enhances some of these steps, especially translation. Several translational activators have been identified, which act at the level of initiating translation. Translation in the chloroplast occurs on chloroplast ribosomes, which are often closely associated with the thylakoid membrane. Cotranslational insertion into the thylakoid membrane has been proposed for the hydrophobic reaction center polypeptides. In addition, synthesis of chlorophyll and its apoproteins needs to be strictly coordinated, because free chlorophyll is highly photoreactive and causes serious damage to the cell.

### 1.4. Chloroplast-Nuclear Cross talk

Chloroplast function and development depend to a large extent on the nucleus. A large number of nuclear genes encode chloroplast structural components and enzymes and are involved in regulating chloroplast gene expression. Reciprocally, chloroplasts also influence nuclear gene activity. This is apparent in mutant plants with defective chloroplasts, where nuclear genes of proteins involved in photosynthesis are no longer expressed. As an example, when carotenoid synthesis is inhibited, chloroplasts rapidly bleach in strong light because chlorophyll is photooxidized in the absence of carotenoids (17). Under these conditions, expression of nuclear genes that code for several abundant chloroplast proteins involved in photosynthesis is specifically repressed. A block in chloroplast protein synthesis has a similar effect (18). These observations imply the existence of a plastid-derived factor that directly or indirectly influences nuclear gene activity. There are mutants of *Arabidopsis* in which the transduction of this plastid-derived signal to the nucleus is affected (19). The nature of the plastid factor is still unknown in plants, although studies with *Chlamydomonas* suggest that some porphyrin compounds, which act as intermediates in the chlorophyll biosynthetic pathway, are involved in this response (Fig. 2, 20).

### 1.5. Protein Sorting in the Chloroplast

Chloroplasts are bounded by an envelope that consists of the outer and inner membranes. The outer membrane is freely permeable to ions and small molecules, whereas the inner membrane is highly selective and contains specific translocators and **permeases** that allow regulated metabolic transport between cytosol and stroma. The envelope also contains the protein import system.

From just the modest size of the plastid genome, it is clear that the majority of the chloroplast proteins are encoded by nuclear genes and imported into the chloroplast. Six chloroplast compartments can be distinguished: (1) the outer envelope membrane, (2) the intermembrane space, (3) the inner envelope membrane, (4) the stroma, (5) the thylakoid membrane, and (6) the lumen. Nucleus-encoded proteins destined to the chloroplast are synthesized as precursor proteins containing, in most cases, a transient N-terminal **transit peptide** (21). Transit peptides are both necessary and sufficient to import a polypeptide into the chloroplast. Transit peptides of stromal proteins consist of 30 to 120 residues in only a poorly conserved sequence. The only distinguishing feature is that they are rich in hydroxylated amino acids and deficient in acidic residues. Recognition of the protein import **receptor** by the transit peptide is followed by translocation of the precursor protein in an extended conformation across the two envelope membranes. ATP and GTP are the sole energy sources for this process, which also requires the participation of several factors to unfold protein on the outside and to refold protein on the inside of the organelle. Several **molecular chaperones** play an important role in the proper folding of the polypeptides that enter the chloroplast (21). Translocation of the precursor of protochlorophyllide oxidoreductase, an enzyme involved in the last step of chlorophyll synthesis, also requires the presence of its substrate, protochlorophyllide, inside the plastid (22). This raises the possibility that the substrate drives the translocation by inducing or stabilizing folding of the enzyme on the stromal side of the envelope.

Thylakoid precursor proteins contain a bipartite transit peptide. The first domain targets the protein to the stroma, and the second hydrophobic domain, which resembles the **signal sequences** of **secretory proteins**, acts as the thylakoid targeting domain. Surprisingly, there are four pathways for protein translocation into or across the thylakoid membrane (21). The first corresponds to the bacterial protein secretion system and uses **Sec proteins** homologous to the bacterial SecA and SecY proteins. The second uses a system involving a **signal recognition particle**. The third pathway is rather unique because it uses only the trans-thylakoid pH gradient as an energy source (see **Chemiosmotic Coupling**). The fourth pathway involves spontaneous insertion of certain proteins into the thylakoid membrane.

Insertion of proteins into the chloroplast envelope occurs by several routes. Some nucleus-encoded polypeptide chains lack a cleavable transit peptide and are inserted directly into the outer and inner membranes. Other envelope membrane proteins containing a cleavable transit peptide use the general import pathway. At least one inner membrane envelope protein is encoded by the chloroplast genome, so it must contain an appropriate targeting signal.



## 1.6. Chloroplast Engineering

A major breakthrough in chloroplast research in 1988 was the development of an efficient method for genetically **transforming** chloroplasts of the green alga *Chlamydomonas* (23), which was subsequently adapted to higher plants (24). In this method, tungsten or gold particles are coated with DNA and bombarded into cells with a particle gun (see [Transfection](#)). Upon entry of the particles into chloroplasts, the DNA is released and integrated into the chloroplast chromosome by homologous [recombination](#). The existence of an efficient chloroplast homologous recombination system and the development of selectable markers for chloroplast transformation have opened the door to manipulating the chloroplast genome. In particular, this new technology allows directed chloroplast gene disruption, a powerful tool for elucidating the role of genes of unknown function. It has also permitted [site-directed mutagenesis](#) of specific residues of photosynthetic reaction center polypeptides so as to gain new insights into their structure-function relationship, and it has been very useful for studying chloroplast gene expression. Chloroplast transformation has important applications for plant biotechnology and [protein engineering](#). Because the chloroplast genome is present in multiple copies, up to 10,000 per cell, new genetic information introduced into plastids is amplified. In principle, this opens the possibility of expressing foreign proteins of commercial interest in large quantities. The expression of foreign genes in the chloroplast compartment offers the additional advantage of considerably reducing the risk of transfer of new genetic material to the environment because the chloroplasts from the male parent are not transmitted to the progeny in most crop plants.

## Bibliography

1. N. W. Gillham (1994) *Organelle Genes and Genomes*. Oxford University Press, New York.
2. J. K. Hooper (1984) *Chloroplasts*. Plenum Press, New York.
3. J. Olive and O. Vallon (1991) *J. Electron. Microsc. Technol.* **18**, 360–374.
4. G. Ojakian and P. Satir (1974) *Proc. Natl. Acad. Sci. USA* **21**, 2052–2056.
5. J. F. Allen (1992) *Biochim. Biophys. Acta* **1098**, 275–335.
6. R. J. Spreitzer (1993) *Ann. Rev. Plant Physiol. Plant Mol. Biol.* **44**, 411–434.
7. M. Sugiura (1996) In *Molecular Genetics of Photosynthesis*. (B. Andersson, A. H. Salter, and J. Barber, eds.), Oxford University Press, Oxford, New York, pp 58–74.
8. P. Bennoun (1982) *Proc. Natl. Acad. Sci. USA* **79**, 4352–4356.
9. M. Reith and J. Munholland (1995) *Plant Mol. Biol. Rep.* **13**, 333–342
10. J. Allen (1995) *J. Theor. Biol.* **165**, 609–631.
11. L. A. Allison, L. D. Simon, and P. Maliga (1996) *EMBO J.* **15**, 2802–2809.
12. M. Sugita and M. Sugira (1996) *Plant Mol. Biol.* **32**, 315–326.
13. U. Kück, Y. Choquet, M. Schneider, M. Dron, and P. Bennoun (1987) *EMBO J.* **6**, 2185–2195.
14. M. Goldschmidt-Clermont et al. (1991) *Cell* **65**, 135–143.
15. H. Kössel et al. (1993) In *Plant Mitochondria*, A. Brennicke and U. Kück (eds.), VCH, Weinheim, Germany, pp. 93–102.
16. J.-D. Rochaix (1992) *Ann. Rev. Cell Biol.* **8**, 1–28.
17. W. Taylor (1989) *Ann. Rev. Plant Physiol. Plant Mol. Biol.* **40**, 211–233.
18. J. Gray (1996) In *Membranes: Specialized Functions in Plants*, M. Smallwood, J. P. Knox, and D. J. Bowles, eds., Bios Scientific Oxford, pp. 441–455.
19. R. E. Susek, F. M. Ausubel, and J. Chory (1993) *Cell* **74**, 787–799.
20. U. Johannigmeier and S. H. Howell (1984) *J. Biol. Chem.* **259**, 13541–13549.
21. K. Cline and R. Henry (1996) *Annu. Rev. Cell Dev. Biol.* **12**, 1–26.
22. S. Reinbothe, S. Runge, B. Reinbothe, B. von Cleve, and K. C. Apel (1995) *Plant Cell* **7**, 161–172.
23. J. E. Boynton et al. (1988) *Science* **240**, 1534–1538.

24. P. Maliga (1993) *Trends Biotechnol.* **11**, 101–107.

### Suggestions for Further Reading

25. N. W. Gillham (1994) *Organelle Genes and Genomes*, Oxford University Press, New York.

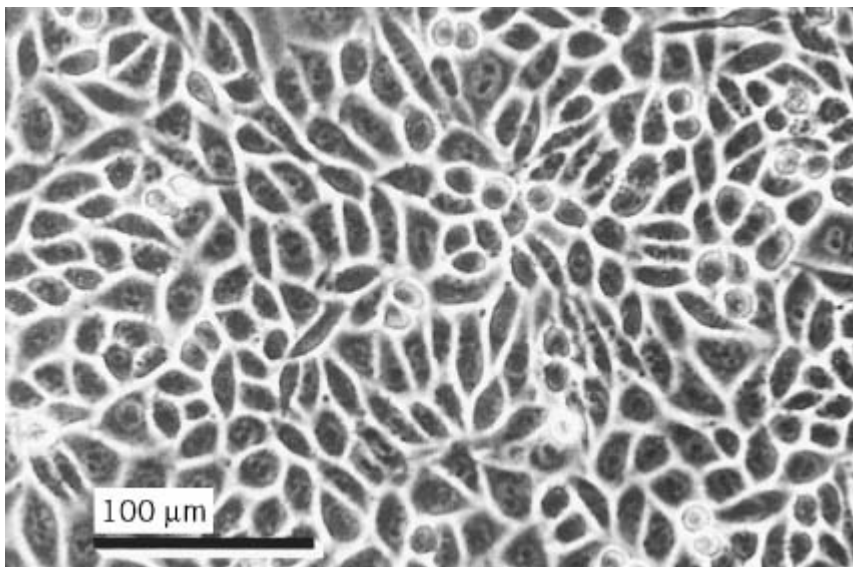
26. J. K. Hooper (1984) *Chloroplasts*, Plenum Press, New York.

## CHO Cells

### 1. Origin

CHO cells were originated as a [cell line](#) by Theodore Puck in 1958 (1) from the Chinese hamster, *Cricetulus griseus*. CHO/Pro cells requiring proline for growth were subsequently derived by nutritional selection. The parental cells were treated with bromodeoxyuridine (see [5-Bromouracil](#)) in proline-deficient medium and then exposed to the near-visible UV from a fluorescent light. This killed the growing cells, but not those that required proline. Proline-requiring clones were then grown up by feeding the surviving cells with proline-rich medium. CHO/Pro<sup>-</sup> was subcloned by dilution cloning, and cell line CHO/Pro-K1 was isolated. The current designation of this cell line in common use is CHO-K1, and it still has a requirement for proline, which is present in Ham's F12, the medium usually recommended for its propagation. It is listed in the American Type Culture Collection (ATCC) catalogue as CCL-61 (Fig. 1).

**Figure 1.** Confluent culture of CHO-K1 cells. Phase contrast, Olympus CK microscope, 20× objective.



### 2. Properties

CHO-K1 (Fig. 1) is a continuous cell line and near-**diploid** with 20 [chromosomes](#) ( $2C = 22$ ). It has a very short doubling time of around 15 h, making it popular as a host for [transfection](#) and biotechnology. The plating efficiency is also very high, and can be 100% under optimal conditions,

so these cells have always been popular for clonogenic survival studies of nutritional mutants and radiation survival. Chinese hamster cells originally became popular for genetic studies because of their relatively small number of readily distinguishable chromosomes, but the advent of chromosome banding techniques and chromosome painting by fluorescence [in situ hybridization](#) makes the distinction of individual chromosomes easier in many other species.

CHO-K1 cells, although transformed, still retain nutrient-dependent G<sub>1</sub> [cell-cycle](#) blockade. Isoleucine deprivation blocks the cells in G<sub>1</sub>, and its restoration generates a synchronous population (2).

### 3. Usage

CHO-K1 cells have been used in cell-cycle control and signaling (3) and are used extensively in biotechnology (4). They are used frequently in DNA [transfection](#) (5) and virally mediated DNA transfer (6). They can be maintained in suspension culture to generate large numbers of cells or product. Using methotrexate-induced co-amplification with cotransfected [dihydrofolate reductase](#), they have been used for production of **interferon-g** (7) and **prothrombin-2** (8). They have been adapted to grow in serum-free medium (9) (see [Serum Dependence](#)). Together with V79 cells, another Chinese hamster cell line, they have been used extensively in genotoxicity studies (10).

### Bibliography

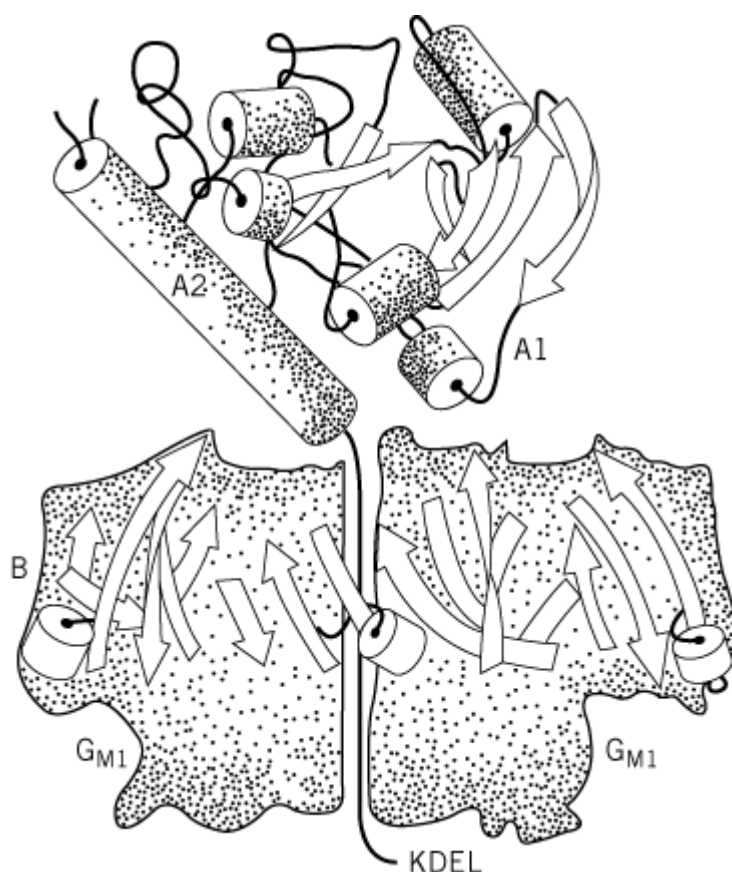
1. T. T. Puck, S. J. Cieciura, and A. Robinson (1958) *J. Exp. Med.* **108**, 945–956.
2. K. D. Ley and R. A. Tobey (1970) *J. Cell. Biol.* **47**, 453–459.
3. G. Tortora, S. Pepe, C. Bianco, V. Damiano, A. Ruggiero, G. Baldassarre, C. Corbo, Y. S. Chochung, A. R. Bianco, and F. Ciardiello (1994) *Intl. J. Cancer* **59**, 712–716.
4. K. B. Konstantinov (1996) *Biotechnol. Bioeng.* **52**, 271–289.
5. S. Subramanian and F. Srienc (1996) *J. Biotechnol.* **49**, 137–151.
6. T. Schoneberg, V. Sandig, J. Wess, T. Gudermann, and G. Schultz (1997) *J. Clin. Invest.* **100**, 1547–1556.
7. V. Leelavatcharamas, A. N. Emery, and M. Al-Rubeai (1994) *Cytotechnology* **15**, 65–71.
8. G. Russo, A. Gast, E. J. Schlaeger, A. Angiolillo, and C. Pietropaolo (1997) *Protein Express. Purif.* **10**, 214–225.
9. M. J. Keen and N. T. Rapson (1995) *Cytotechnology* **17**, 153–163.
10. H. F. L. Mark, R. Naram, T. Pham, K. Shah, L. P. Cousens, C. Wiersch, E. Airall, M. Samy, K. Zolnierz, R. Mark, K. Santoro, L. Beauregard, and P. H. Lamarche (1994) *Ann. Clin. Lab. Sci.* **24**, 387–395

### Cholera Toxin and Enterotoxins

Cholera toxin (CLT) and the closely similar heat-labile enterotoxins (LTs) are released by toxigenic strains of *Vibrio cholerae* and *Escherichia coli*, respectively, the causative agents of epidemic diarrheas (1). These bacteria attach firmly to the apical portion of the intestinal epithelium and produce several [toxins](#), including CLT or LT, encoded by genes contained in a transferable virulence cassette (2). Five B subunits (103 residues) are secreted into the bacterial **periplasm**, where they assemble into a pentamer that binds the A chain. A is then proteolytically cleaved at a single point to generate the catalytic A1 chain (192 residues) linked via a [disulfide bond](#) to the A2 peptide (3, 4)

(Fig. 1). The structural organization of these enterotoxins is shared by *Shiga* toxins (other enterotoxins, which cause bloody diarrheas and necrosis of intestinal epithelium and are produced by *Shiga spp.* and *E. coli spp.*) and [pertussis toxin](#) (3-6). The B pentamer has the shape of an irregular cylinder with a flat surface and convoluted bottom surface. The five B monomers are arranged around a 5-fold axis in the oligomer-binding fold, recently identified in a set of oligonucleotide- and oligosaccharide-binding proteins (7). Each monomer contains five b-strands in an antiparallel [beta-sheet](#) and one [alpha-helix](#). The five helices form a central pore hosting the carboxyl-terminal segment of A2, which emerges on the other side with a Lys/Arg-Asp-Glu-Leu motif that may be important in cell penetration (8). The amino-terminal half of A2 forms a long  $\alpha$ -helix, rising from the flat surface, involved in interaction with A1. Thus, very few protein–protein contacts exist between A1 and B. The convoluted oligomer B surface, distal from A1, forms five sugar-binding sites on the outside edge of the cylinder (1, 3, 4).

**Figure 1.** Structural organization of cholera toxin and related toxins, showing a cross section of the molecule of cholera toxin and the *E. coli* heat-labile enterotoxins (1, 3, 4). A is composed of two polypeptide chains: A1, endowed with ADP-ribosyltransferase activity, and A2, which consists of a long  $\alpha$ -helix, involved in interaction with A1, followed by a structureless segment that penetrates the small central hole of the B pentamer. Each of the five B subunits forming the B oligomer has a binding site for the oligosaccharide portion of the ganglioside  $G_{M1}$ . Thus, cholera toxin binds to the membrane via multiple interactions, with the catalytic A1 subunit pointing away from the membrane, with little protein–protein contact between A1 and B.



The structure of A1 shows a cleft where NAD binds and Arg-7 and Glu-112, two residues essential for activity, are located (9). Each B subunit of CLT or LTs possesses a single binding site specific for the oligosaccharide portion of a ganglioside (10), and one molecule of toxin binds five glycolipid molecules. Oligosaccharide binding affinity is rather low, but a high affinity of the toxin for the cell ( $K_d$  of the order of  $10^{-10}$ ) is obtained because of the pentavalent binding. This strategy of

multivalent binding to obtain a strong cell association is displayed by other bacterial toxins and viruses (11).

The target of CLT and LTs is localized on the basolateral membrane; therefore, these toxins have to transcytose through the cell to display their activity. Available evidence indicates that CLT and LTs are internalized inside apical [endosomes](#) that may recycle the glycolipid–toxin complex to the surface or deliver it to later endosomes. As a result of the binding of the **KDEL sequence** of the A subunit to a **KDEL receptor**, the toxin moves retrogradely through the **Golgi** cisternae, and some toxin molecules are then expected to be sorted into basolateral endosomes that fuse with the basolateral membrane (8). At some stage of this intracellular trafficking, these toxins have to be reduced and the A1 subunit has to translocate from the luminal to the cytosolic side of the membrane to interact with their target on the cytosolic face of the basolateral membrane. Structural and membrane photolabeling data indicate that oligomer B does not penetrate the lipid bilayer (1, 10, 12), but A1 inserts in the membrane upon reduction of the A1—A2 interchain disulfide bond. Hence, it is likely that as soon as the disulfide bridge is reduced, A1 “rolls over” oligomer B and inserts into the membrane. This is at variance from postulated mechanism of membrane penetration of other toxins, whose protomer B plays an active role in the insertion of the catalytic subunit. Neither the chemical nature of the reducing agent nor the intracellular stage at which this step takes place are known.

The A1 subunits of CLT and LTs catalyze the transfer of ADP-ribose from NAD to an Arg residue present in the LRXRVT conserved sequence of the  $\alpha$  subunit of the  $G_S$ ,  $G_P$ , and  $G_{olf}$  large trimeric [GTP-binding proteins](#) involved in the coupling of cell surface receptors to the [adenylate cyclase](#) (13). Such modification results in a permanent activation of this latter enzyme and a large increase in cellular [cyclic AMP](#) (cAMP) level, which initiates a cascade of [signal transduction](#) pathways. A1 ADP-ribosylating activity is enhanced by a group of cytosolic or membrane GTP-binding proteins, termed ARF (ADP-ribosylating factors), present in eukaryotic cells (14). ARF are strongly conserved from yeast to humans and are involved in the control of membrane trafficking and protein transport inside cells (15).

The increased cAMP level brought about by CLT or LTs in the enterocyte has a series of consequences, but the inhibition of a **sodium channel** and activation of a **chloride channel** localized on the apical membrane appear to be very relevant to diarrhea. In fact, a decreased sodium reabsorption and increased chloride secretion cause an osmosis-driven loss of water into the intestine. Also important in cholera is the activation of entero-chromaffin cells, which respond to the cAMP increase with release of VIP (vasointestinal peptide), which further lowers intestinal water reabsorption and inhibits muscle cells with an alteration of intestinal peristalsis.

## Bibliography

1. B. D. Spangler (1992) *Microbiol. Rev.* **56**, 622–647.
2. M. K. Waldor and J. J. Mekalanos (1996) *Science* **272**, 1910–1914.
3. T. Sixma et al. (1993) *J. Mol. Biol.* **230**, 890–918.
4. R. G. Zhang et al. (1995) *J. Mol. Biol.* **251**, 563–573.
5. M. E. Fraser, M. M. Chernaiia, Y. V. Kozlov, and M. N. G. James (1994) *Nature Struct. Biol.* **1**, 59–64.
6. P. E. Stein et al. (1994) *Structure* **2**, 45–57.
7. A. G. Murzin (1993) *EMBO J.* **12**, 861–867.
8. W. I. Lencer et al. (1995) *J. Cell Biol.* **131**, 951–962.
9. M. Domenighini, C. Magagnoli, M. Pizza, and R. Rappuoli (1994) *Mol. Microbiol.* **14**, 41–50.
10. E. A. Merritt et al. (1994) *Protein Sci.* **3**, 166–175.
11. G. Menestrina, G. Schiavo, and C. Montecucco (1994) *Mol. Aspects Med.* **15**, 81–193.
12. M. Tomasi and C. Montecucco (1981) *J. Biol. Chem.* **256**, 11177–11181.

13. D. M. Gill and M. J. Woolkalis (1991) *Methods Enzymol.* **195**, 267–280.
14. J. Moss and M. Vaugham (1991) *Mol. Microbiol.* **5**, 2621–2627.
15. J. E. Rothman and F. T. Wieland (1996) *Science* **272**, 227–234.

## Chorion Genes and Proteins

The outer eggshell of those numerous insects whose [embryos](#) develop externally is known as the *chorion*, and it must be mechanically robust in order to provide the environment in which the [egg](#) can develop successfully. The chorion must also have a structure that minimizes water loss, while still permitting the gas exchange vital for embryonic respiration. In *Drosophila* the chorion consists of an outer exochorion, an endochorion, a thin inner chorionic layer, a wax layer and the vitelline membrane. The endochorion has the most important structural and gas-transporting roles *in vivo* (1). The wax and inner chorionic layers act as the waterproofing agent for the egg. Proteins in the chorion undergo extensive [post-translational modifications](#) and consequent peroxidase-catalyzed cross-linking of two or three [tyrosine](#) residues, which results in a significant increase in the mechanical stability of the entire structure. In contrast to the chorion in the *Drosophila*, however, the gypsy moth and other Lepidoptera have a chorion with an outer “sieve” layer (30 nm thick) overlying about 50 to 60 lamellae of variable thickness (but each about 0.2 μm thick on average). The innermost region of the chorion (ie, that beneath the lamellar region) is called the *trabecular layer* and has a thickness of about 0.5 μm (2). Lamellar substructures of this general type are not found outside Diptera. An individual lamella consists of a fibrous layer in which fibrillar elements reported to be about 3 to 4 nm in diameter lie parallel to the surface. The orientation of the lamellae change in a relatively systematic manner from one lamella to the next, thus generating a liquid crystal-like, cholesteric phase ultrastructure. These and other observations reveal that the *Drosophila* and silk moths show remarkable similarities, but equally remarkable differences. For example, they share regulatory elements that direct chorion **gene expression** to the follicular epithelium at the end of **oogenesis** (2). In contrast, they differ in their structural **gene** sequences, their **chromosomal** organization, their chorionic ultrastructures, and their modes of morphogenesis.

The chorion in general is a very complex structure and, in some cases at least, contains a hundred or more different proteins. The proteins in silk moths, however, all have similarities in sequence and fall within either the a or b branches of the chorion superfamily. These proteins in turn are encoded by numerous duplicated genes, all of which contain a single **intron**. Fortunately, the chorion genes too fall into one of a small number of families - A, B, and C. Each family contains multiple genes that occurred during evolution by [gene duplication](#) and sequence **divergence**. The families themselves are related and constitute a [superfamily](#) with A genes in one branch and B and C in the other. The A and B genes, which exist in pairs in divergent orientation in the chromosome, are coordinately expressed and **transcribed** in opposite directions under the direction of a bidirectional **promoter**. The chorion genes in *Drosophila* are quite different from those in silkmoths with regard to sequence and organization. In particular, the *Drosophila* genes are tandemly orientated, and each has its own unidirectional promoter that is temporally unique.

Proteins in the mature silkmoth chorion are usually small (10 to 20 kDa) and have *N*- and *C*-terminal **domains** that are variable in sequence and structure. These regions endow the protein with specific functional and structural attributes. Quasi-repeats are not uncommon (3): In *Drosophila* the dipeptide (Gly—His) is repeated five times towards the *N*-terminal end of s36 and nine times in the *C*-terminal region of s38; a hexapeptide repeat based on tyrosine is found in the *C*-terminal region of s36. Peptide repeats of the form (Cys—Gly), (Cys—Gly—Gly), and (Gly—Tyr—Gly—Gly—Leu) are

found in the case of silkmoths. Unlike the chorion proteins in *Drosophila*, a central domain is largely conserved in the a and b families of the silkmoth proteins. This is characterized by tandem repeats. In the case of the a family, the central domain is 52 residues in length and displays a six-residue quasi-repeat with a **consensus sequence** (Gly—X—Val/Ile—Y—Val/Ile—Z), where X is generally a charged or large polar residue, Y is variable, and Z is commonly Ala, Gly or Cys. The b family also displays a similar but different six-residue quasi-repeat. The predicted **secondary structure** of both hexapeptide repeats is that of an antiparallel **b-strand** terminated by a **b-turn**. Interestingly, the occurrence of valine and isoleucine residues two apart in the a repeat places them on the same side of the b-sheet in a manner analogous to that in feather **keratin**. Depending on the precise form of the b-turns, one face of the twisted **b-sheet** would be almost entirely apolar and could give rise to a b-barrel structure in which the apolar residues would be located internally. In the case of the gypsy moth Ld15 protein (and the 292a protein of *A. polyphemus*), both penta- and hexapeptide quasi-repeats are found with consensus sequences of (Gly—Leu—X—Pro/Gly—Tyr), (Tyr—Gly—X—X—Gly/Ala), and (Gly—X—Val—X—apolar—Ala/Gly) (2). *Drosophila* chorion proteins s36 and s38, on the other hand, contain three consecutive large apolar residues in a quasi-repeat, and these are likely to form a b-strand. They are followed by three to seven residues rich in proline and basic amino acids that are likely to form a connecting b-turn or loop. These few examples illustrate that although there are significant variations in the detail of the sequence repeats, the antiparallel b-sheet conformation that is predicted to occur in all cases would seem to be a constant feature of all the chorion proteins. The paired organization of the moth chorion genes may imply that the a and b members of the chorion superfamily interact with one another to form the core of the structure, thus leaving the N- and C-terminal regions in an external location where they are able to provide other important functional properties.

### Bibliography

1. A. C. Spradling (1993) "Developmental Genetics of Oogenesis". In *The Development of Drosophila melanogaster*, Vol. 1, (M. Bate and A. M. Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–70.
2. R. F. Leclerc and J. C. Regier (1993) Choriogenesis in the Lepidoptera: morphogenesis, protein synthesis, specific mRNA accumulation, and the primary structure of a chorion cDNA from the gypsy moth. *Dev. Biol.* **160**, 28–38.
3. S. J. Hamodrakas, A. Batrinou, and T. Christophoratos, Structural and functional features of *Drosophila* chorion proteins s36 and s38 from analysis of primary structure and infrared spectroscopy. (1989) *Int. J. Biol. Macromol.* **11**, 307–313.

### Suggestions for Further Reading

4. F. C. Kafatos, G. Tzertzinis, N. A. Spoerel, and H. T. Nguyen (1995) "Chorion Genes: An Overview of Their Structure, Function and Transcriptional Regulation". In *Molecular Model Systems in the Lepidoptera* (M. R. Goldsmith and A. Wilkins, eds.) Cambridge University Press, Cambridge, England, pp. 181–215.
5. T. L. Orr-Weaver (1991) *Drosophila* genes: cracking the eggshell's secrets. *Bioessays* **13**, 97–105.

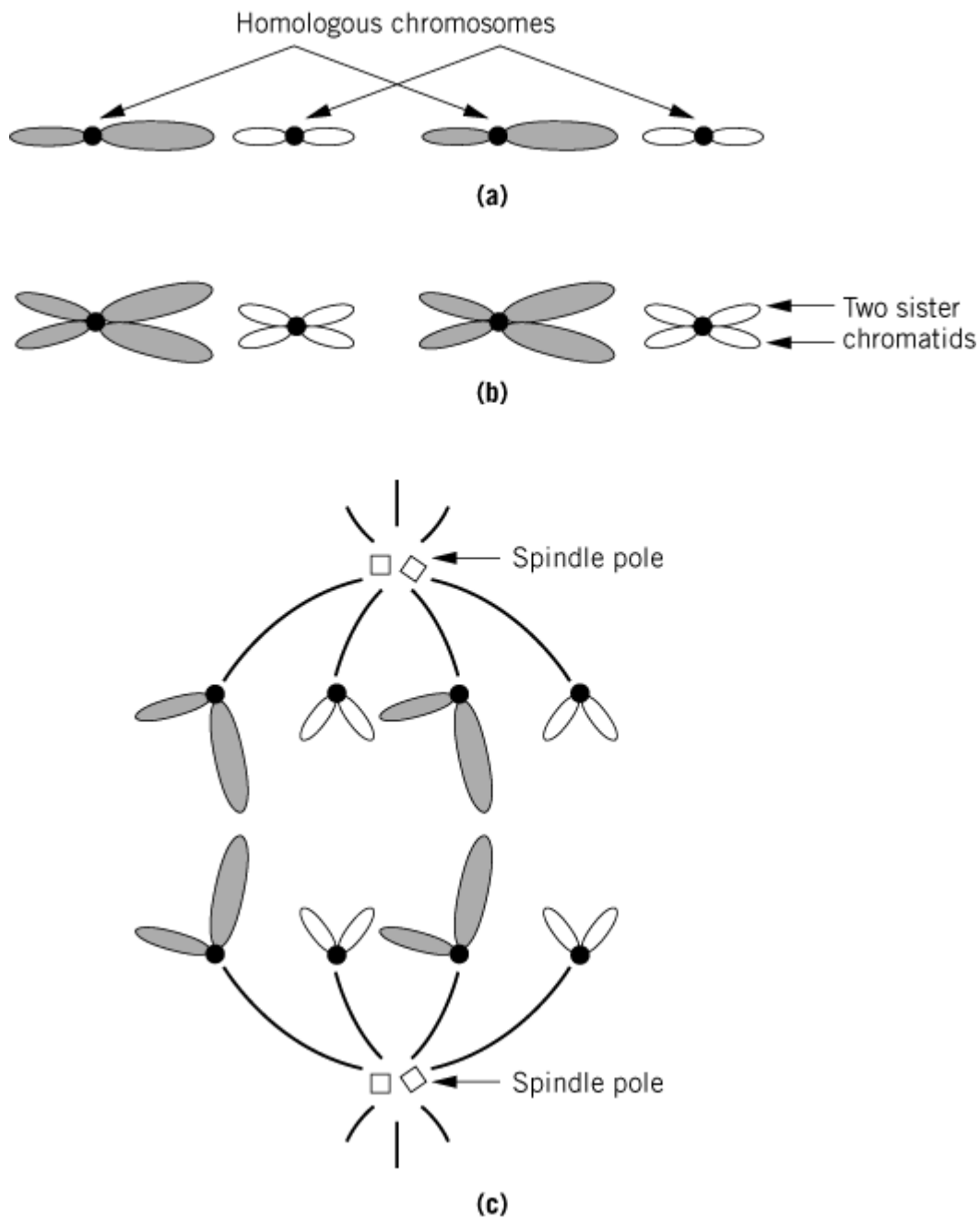
### Chromatid

A chromatid is one of the two adjacent strands of a **chromosome** that are generated as a result of chromosomal duplication. The two chromatids are held together at the **centromere** and become

visible during **mitosis** (Fig. 1, see top of next page). Chromosomal dynamics during the [cell cycle](#) are conveniently discussed in terms of the appearance and behavior of chromatids. It is during **interphase** that DNA is replicated (**S-phase**). Before the S-phase there is a gap known as G1, when the chromosomes are considered to be unineme, meaning that a single duplex DNA molecule runs the entire length of the chromosome (see [Lampbrush Chromosomes](#)). The amount of DNA in G1 per haploid genome is called the **C-value**. A **diploid** cell has a 2C DNA content. After the S-phase, in the gap known as G2, the chromosomes have duplicated and there is a 4C DNA content. Now each chromosome consists of two unineme sister chromatids. Mitosis begins at the end of G2. Whereas the chromosomes are too diffuse and decondensed to be visualized when somatic cells are in the interphase, during mitosis the chromosomes and chromatids become visible. Mitosis can be divided into four stages. During **prophase**, the [chromatin](#) in the chromosomes begins to condense so that they appear as morphologically indistinct, diffuse threads. At the end of the prophase, the chromosomes become well separated. Now each chromosome is distinctly double-stranded, consistent with two sister chromatids joined at the centromere. In the metaphase the chromosomes are maximally contracted and align on a plane, known as the metaphase plate. The metaphase plate is midway between the **spindle poles** and represents an integral component of the **mitotic spindle**. In the **anaphase**, the pairing region within the centromere dissolves (see [Centromeres](#)). This allows the mitotic spindle to pull the two sister chromatids apart. Each pole receives a complete set of chromatids, representing all of the chromosomes. Finally in **telophase** the mitotic spindle breaks down, and the nuclear membrane is reassembled around the chromosomes. The chromosomes progressively decondense and cease to be visible by light [microscopy](#).

**Figure 1.** Chromosomal dynamics during the cell cycle. **(a)** At the end of mitosis and in G1, there are two copies of each chromosome, called homologous chromosomes, in a diploid cell. **(b)** At the end of the S-phase and in G2, each chromatid in each chromosome has been duplicated, creating two sister chromatids joined at the centromere. **(c)** During the anaphase of mitosis the two sister chromatids separate and move to opposite spindle poles.





Sister chromatids also interact during **meiosis**, which represents the cell division, chromosomal and chromatid segregation events associated with **gamete** formation. There are two cell divisions during meiosis. Once again chromosomal dynamics are best described in terms of the segregation of chromatids. Initially, the chromosomes replicate, then the cell enters the first meiotic cell division. Homologous chromosomes are attracted to each other, forming stable pairs held together by the [synaptonemal complex](#) (see [Homologous Chromosomes](#)). Thus each structure consists of two chromosomes, each containing two sister chromatids. [Recombination](#) occurs between the maternal- and paternal-derived chromosomes. Sites of such genetic **crossing over** are called **chiasmata**. Chiasmata represent sites of physical breakage and reunion between nonsister chromatids. They generally encompass a few hundred to a few thousand base pairs of DNA. Finally, each pair of chromosomes separates. At the start of the second meiotic cell division, each chromosome consists of two chromatids joined by a centromere. Then the chromatids are aligned on a metaphase plate attached to a spindle. The centromere splits, and the two sister chromatids move to opposite poles. Then nuclei are reformed in each of the four haploid daughter cells. Thus for both mitosis and

meiosis the mechanics of chromosome condensation, recombination, and segregation emphasize the fate of individual chromatids.

### Suggestions for Further Reading

A. T. C. Carpenter (1987) Gene conversion, recombination nodules, and the initiation of meiotic synapsis. *BioEssays* **6**, 232–236.

G. J. Gorbsky (1992) Chromosome motion in mitosis. *BioEssays* **14**, 73–80.

G. S. Roeder (1990) Chromosome synapsis and genetic recombination: Their roles in meiotic chromosome segregation. *Trends Genet.* **6**, 385–389

## Chromatin

Chromatin is the complex of [histones](#) and DNA, associated with smaller amounts of other proteins, into which the chromosomal DNA of all eukaryotes is organized. As well as being a means of packing DNA in an orderly fashion, chromatin structure plays a crucial role in the regulation of gene expression, both activation and repression. Three books contain much useful information about chromatin structure and function ([1-3](#)); for the most up-to-date information in many areas of this rapidly moving field, the reader will be referred to recent articles and reviews.

### 1. Overview

Chromatin has a beaded appearance in the [electron microscope](#) at low ionic strength, arising from a regularly repeating structure ([4](#)). The repeating unit is the [nucleosome](#), which contains ~166 to 240 bp of DNA, of which ~166 bp is wound in two left-handed superhelical turns around an octameric complex of the four core [histones](#) (H3, H4, H2A, and H2B) and stabilized by one molecule of histone H1; the nucleosome also includes a variable length of [linker DNA](#) (~0 to 76 bp) that connects neighboring nucleosomes. The octamer is organized as a central H3<sub>2</sub>H4<sub>2</sub> tetramer, which has a rather flat, twisted horseshoe shape, flanked above and below by an H2A–H2B dimer. The details of the organization of histone and DNA in the core of the nucleosome are revealed in the high-resolution structure of the 146-bp [nucleosome](#) core particle (which contains no linker DNA and no H1) recently determined by [X-ray crystallography](#) ([5](#)). The globular domain of histone H1 binds asymmetrically, bridging a point near the dyad and an entering or exiting DNA double helix ([6](#)) (see [Nucleosome](#) for further details); the basic C-terminal domain interacts with the linker DNA, partially neutralizing its charge. The array of nucleosomes is known as a nucleosome filament, or the 10-nm filament (10 nm being the diameter of the nucleosome). The amount of DNA contained in the repeating unit, which reflects the length of linker DNA, is determined from the sizes of the DNA in the chromatin fragments liberated by incomplete digestion with micrococcal nuclease (**Staphylococcal nuclease**), which cuts in the linker. A series of fragments, usually visualized as a “ladder” of bands in [gel electrophoresis](#), occurs at multiples of the unit repeat size. This is commonly about 200 bp (when there is a linker length of 200 – 166 = 34 bp), but it may be as long as ~240 bp (in sea urchin sperm), or as short as ~166 bp (essentially no linker DNA) in yeast and in mammalian cerebral cortex neurons (the glial cells in the same tissue have a repeat length of about ~200 bp). The repeat length measured in this way is the average for the tissue or cell population under examination; there will be local variations about this mean (see [Linker DNA](#)). Packaging of the DNA in nucleosomes achieves about a sixfold length compaction of the DNA, and further packing is achieved by salt-dependent, histone H1-assisted folding of the 10-nm nucleosome filament into a 30-nm filament. Much evidence supports a simple helical coiling of the nucleosome filament into a

“solenoid” with six nucleosomes per turn, but there are other models (see text below). Most of the chromatin in the cell is in the form of a 30-nm filament throughout most of the [cell cycle](#), and at interphase appears to be looped onto the nuclear matrix; at mitosis, the 30-nm filament is subjected to further levels of folding to give the highly dense metaphase chromosome in which the overall length compaction of the DNA is about 10,000-fold.

The 30-nm filament is not a suitable template for [transcription](#) by the large eukaryotic **RNA polymerases** and clearly has to be unfolded first; formation of an “open” (“transcriptionally competent”) chromatin state is an essential prerequisite for transcription. This involves acetylation of the basic *N*-terminal regions of the core histones, which appear to play a role in chromatin folding, and some depletion of histone H1. It has recently become clear that acetylases are targeted to particular regions of the chromosome through recruitment, directly or indirectly, by gene regulatory proteins bound to specific DNA sequences at **promoters** or upstream activating sequences. Cells appear to contain many acetylase (and deacetylase) activities, which function as components of multiprotein complexes (see [Histone Acetylation](#)). Even in the 10-nm filament form, nucleosomes may act as blocks to the formation of transcription initiation complexes by occluding promoters, which may include strategically positioned nucleosomes. It is now clear that the cell has specialized energy-dependent mechanisms that are required for “remodeling” some promoters, so that they become accessible to transcription factors, as well as mechanisms for facilitating RNA chain elongation. These mechanisms also involve multiprotein complexes, which may have subunits in common with acetylases and deacetylases. Chromatin is thus both an effective means of packing and the target for an integrated network of machines that modify it. There is a wealth of genetic evidence showing that the growing list of factors that regulate transcription includes both *bona fide* chromatin components (eg, histones) and proteins that modify chromatin (eg, acetyltransferases and deacetylases). An unexpected discovery was that a structural motif found in all the core histones is also found in a number of components of the transcription initiation factor TFIID (see [Histone Fold](#)), as well as in some histone acetyltransferase complexes (see [Histone Acetylation](#)), suggesting a structural and evolutionary link between chromatin and the assemblies that act upon it. Another surprising finding has been that histone H1 may have a specific gene regulatory role, in addition to its role in packaging and as a general repressor (see [Histones](#)).

## 2. Nucleosome Positioning and Mobility

Histones package the whole of the [genome](#) and therefore appear to be largely indifferent to DNA sequence. Locally, the position of a nucleosome is determined by the underlying DNA sequence and, in particular (because the DNA, which is normally rigid, is bent around the octamer surface), by the local bendability of DNA, which is sequence-dependent. Analysis of the DNA sequences of bulk nucleosome core particles shows that A and T di- and trinucleotides are likely to occur where the minor groove faces inward, and G and C di- and trinucleotides where the minor groove faces out. This is because the minor groove of AT-rich DNA is naturally narrow and so, by having A and T where the minor groove faces in, the compression of the groove that occurs on wrapping the DNA around the octamer is readily accommodated ([7](#)). The dinucleotide periodicity is on average 10.2 bp, exactly the average structural periodicity given by the X-ray crystal structure of the nucleosome core particle (see [Nucleosome](#)). The helical periodicity in solution is 10.6 bp, so the DNA is overtwisted in the nucleosome, presumably to achieve the match between the structural and sequence periodicities that underlie the rules that govern rotational positioning of the octamer with respect to the DNA sequence (the orientation of a face of the duplex with respect to the histone octamer surface). This is thus determined by the local bendability (flexibility toward curvature) of the DNA, which is determined by the properties of the individual base steps. A rotationally positioned duplex shows a characteristic 10-nucleotide pattern of nicking by DNase I, as each strand rises from the surface about every 10 nucleotides. Translational positioning, on the other hand (the basis of which is not well understood), refers to the choice of a particular stretch of DNA by the histone octamer, rather than other stretches of the same length translated forward or backward along the DNA by about 10 bp, which would allow the same rotational setting. The abundance of AA/TT dinucleotides separated by roughly 10 bp led to the construction of fragments with alternating A/T and C/C

sequences  $[(A \text{ or } T)_3 \text{NN}(G \text{ or } C)_3 \text{NN}]_n$ —the so-called TG pentamer—that indeed formed very stable nucleosomes (8). Selection of naturally occurring 146-bp DNA sequences with a high affinity for the histone octamer revealed (in addition to the A/T, G/C periodicity) that sequences with repeated TATAAAACGCC motifs (“phased TATA” sequences) formed nucleosomes that were even more stable, suggesting that high-affinity binding is aided by flexible sequences (9); selection from longer (220 bp) synthetic DNA sequences also revealed a signal (CTAG) that favored high-affinity binding (10). Such sequences, even if they occurred infrequently in the genome, could have important implications for chromatin organization and regulation, as could sequences that are refractory to nucleosome formation (eg, TGGGA repeats identified in a negative-selection approach (11)). It seems likely that the signals discovered for rotational positioning *in vitro* (7) also function *in vivo*, as revealed by analysis of the complete genome sequences of the yeast *Saccharomyces cerevisiae* and the nematode *Caenorhabditis elegans* (12), where dinucleotide periodicities are compatible with nucleosomal constraints. Analysis of 168-bp chicken erythrocyte chromatosome sequences [ie, containing H5(H1)] shows that the statistical preferences of the octamers for particular sequences are slightly modulated (13), presumably in order to optimize the stability of the octamer–DNA–H1 ternary complex. This would be reflected in slight differences in translational positions of nucleosomes in the presence of H5(H1), which might have functional consequences (eg, in exposing or revealing short DNA sequences in an H1-dependent manner). An additional feature of chromosomal DNA is the frequent occurrence of AGGA within half a double helical turn of one terminus, imposing asymmetry on the chromatosome. Whether this is related to the asymmetric binding of H1 to the nucleosome remains to be seen. The rules for nucleosome organization on DNA sequences have recently been reviewed (14).

Although histones can, in general, package DNA essentially irrespective of sequence, but capitalize on the local bendability of DNA to promote a better fit on the octamer surface, some DNA sequences have a particularly high affinity for the octamer, leading to a “positioned nucleosome” in which a particular translational position is preferred. Nucleosome reconstitution experiments show that the H3–H4 tetramer alone is sufficient to confer positioning on a defined DNA sequence (15, 16). Positioned nucleosomes have been identified *in vivo* by nucleosome mapping techniques and *in vitro* in nucleosome reconstitution experiments on defined DNA sequences. The rules for nucleosome positioning *in vivo* are not well understood: When the TG pentamer (see above) was introduced into yeast, it appeared to exclude nucleosomes rather than provide optimal positioning (17), showing that a strong rotational setting is not sufficient for positioning *in vivo*. Positioned nucleosomes may act as a boundary (18) and determine the positions of several neighboring nucleosomes on either side. More than one positioned nucleosome separated at less than a nucleosome's length of DNA would serve to keep the region clear of histones (eg, at promoters) for the binding of sequence-specific proteins with various gene regulatory or structural roles. At the inducible yeast *PHO5* promoter, events leading to promoter remodeling and transcription of the *PHO5* gene, which encodes a phosphatase, are initiated in response to low phosphate levels by binding of the transcription factor PHO4 to a site between two positioned nucleosomes in a set of four over the promoter (19) (see text below). Positioned arrays of nucleosomes, essential for the functioning of the promoter, are also found in many other cases—for example, the mouse mammary tumor virus (MMTV) major late promoter (20) (see also Ref. 2). It has recently become apparent (21) that nucleosome positioning is the probable explanation for the differential regulation of *Xenopus* oocyte and somatic 5S genes *in vivo* (22), where the oocyte genes are repressed and the somatic genes remain active when H1 accumulates after the mid-blastula transition. Sequence differences between the two types of gene mean that the oocyte nucleosome occludes the TFIID binding site and binds H1 preferentially (so the gene is turned off), whereas the converse is broadly true for the somatic gene, which remains active.

A nucleosome core can be reconstituted on to DNA *in vitro* by dialysis of an equimolar mixture of the histone octamer (or H3–H4 tetramers and H2A–H2B dimers) and DNA from high ionic strength (eg, 2 M NaCl) to low ionic strength. Even in cases where there are strong nucleosome positioning signals in the DNA, octamers on different fragments will occupy a population of less favorable

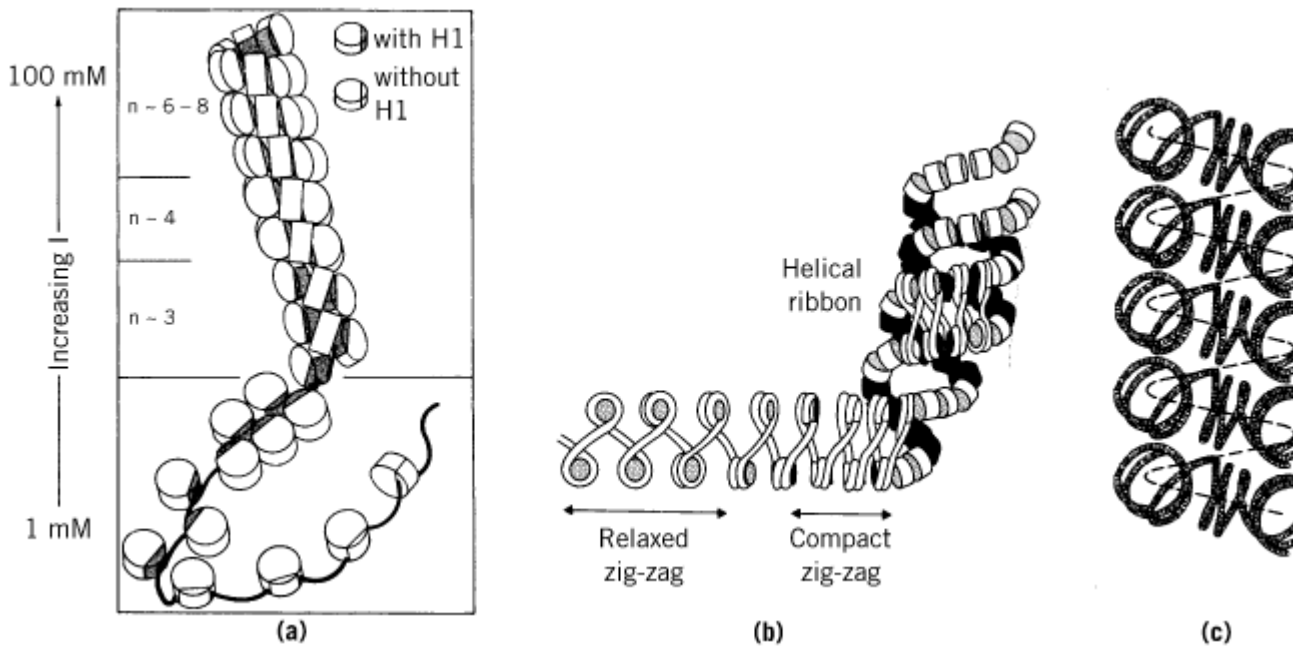
minor translational positions in addition to a major position, all with the same rotational setting of the DNA. For example, on the *Xenopus borealis* somatic 5S rRNA gene, which is regarded as a strong positioning sequence, octamers can occupy 10 different positions (six being predominant), all related to each other by the 10 bp helical periodicity of the DNA, which must be caused by a strong rotational signal in the DNA sequence (21). The occupancy of less favorable octamer positions on a DNA sequence is higher at low temperatures (eg, 4°C), when the octamers are trapped, than if the temperature is then shifted to 37° (23). The change in positions populated demonstrates that nucleosomes have an intrinsic mobility, which could be useful, for example, in remodeling of promoters. In the presence of H1, nucleosome mobility is greatly reduced (24, 25); conversely, loss of H1 might enable nucleosomes to move, for example, to expose binding sites for gene regulatory proteins.

### 3. Higher-Order Structure: The 30-nm Filament

The 10-nm nucleosome filament, which represents the first level of chromatin structure, is the form in which chromatin exists at low ionic strength (eg, >20 mM) *in vitro*. As the ionic strength is increased, the 10-nm filament condenses into a 30-nm (diameter) filament, visible by electron microscopy. Most of the chromatin in the nucleus is in this form. The folding is a consequence of ionic-strength induced screening of the residual negative charge on the linker DNA (26), much of the charge having already been neutralized by the basic tail(s) of histone H1 and, probably, the basic *N*-terminal tails of at least some of the core histones, and results in new interactions (eg, nucleosome-nucleosome) within the higher-order structure.

Various models have been proposed for the 30-nm filament (27), some of which are shown in Figure 1. All agree on a radial distribution of nucleosomes with their faces parallel to the filament axis, but they differ fundamentally in the path taken by the linker DNA, whether bent, curved, or straight. The evidence on this point is conflicting (see [Linker DNA](#)). The evidence for the various models has been extensively reviewed (28-31). Much evidence supports the original “solenoid” model (32), which proposes a simple helical coiling of the nucleosome filament, with about six nucleosomes per turn at physiological ionic strength, and necessarily bent linker DNA (Fig. 1a). [The details are easily modified to accommodate the asymmetric, rather than symmetric, location of the globular domain of H1 near the dyad (6); see text above.] The solenoid places the dyads of the nucleosomes, histone H1, and the linker DNA between nucleosomes, on the inside of the solenoid; H1 has indeed been shown to be located within the 30-nm filament (33). However, other models have been proposed, such as the helical-ribbon model (34, 35) (Fig. 1b), which invokes straight linkers and is based on the zigzag appearance of the nucleosome filament in electron micrographs, and these too would be compatible with an internal location of H1. The coiled linker model (36), like that shown in Fig. 1c, which is related to the solenoid, is not necessarily compatible, because the dyads of the nucleosomes, and hence H1, are alternately inside and outside the solenoid for some linker lengths (eg, in chicken erythrocyte chromatin).

**Figure 1.** Some models of chromatin higher-order structure. The 30 filament is shown as (a) a solenoid (32), (b) a helical ribbon (34), and (c) as a helical conformation with coiled linkers (36). In (a) and (b) the most condensed structures are shown at the top, and the extended structures present at low ionic strength at the bottom. (From Refs. 32, 34, and 36, respectively, with permission.)



Whatever the model that best describes the spatial relationships between nucleosomes in the higher-order structure, the 30-nm filament in the cell nucleus is unlikely to be completely regular, partly because of local heterogeneity in linker length along the filament within the average value, partly due to the presence of different histone variants or subtypes, and so on. The 30-nm filament is also a dynamic structure that undergoes thermal “breathing”, and this property is probably central to gene regulation through control of chromatin folding. The stability of the 30-nm filament is relevant to the ease with which it may be unpackaged for transcription and appears to be increased by special variants of H1 (such as H5 and spH1 in the condensed chromatin of nucleated erythrocytes and sea urchin sperm, respectively) (37), and it is probably decreased in regions of transcriptionally competent chromatin by (partial) loss of H1 (eg, Ref. 38) and by acetylation of the core histone *N*-terminal tails (see [Histone Acetylation](#)).

#### 4. Chromatin and Transcription

In eukaryotes, transcription works in a chromatin context (39, 40). Some genes are active in all cell types, whereas others are active only in specific cell types. Active genes in higher eukaryotes have a chromatin structure characterized by a greater sensitivity (7- to 10-fold) to the endonuclease DNase I than the bulk of the chromatin in the nucleus (1, 2). This state precedes active transcription and marks “transcriptional competence”; it is the first step in control of transcriptional activity. DNase I sensitivity is probably due to (partial) relaxation of the 30-nm filament toward the 10-nm filament state, which is probably facilitated by partial loss of linker histones and acetylation of the core histone *N*-terminal tails (41) (see [Histone Acetylation](#)). DNase I sensitivity associated with a particular gene extends well beyond the coding region and encompasses an entire chromosomal domain (eg, tens of kilobases) and includes all the regulatory elements (eg, [enhancers](#)) for the genes (2). The boundaries of these chromosomal domains have been characterized in only a few cases, and no common feature is apparent. In at least some instances, the boundaries act as “insulators”, preserving the structural and functional autonomy of the domain and also blocking the influence of regulatory elements from the outside (42). Short localized regions, termed DNase I hypersensitive sites (typically a few hundred base pairs), in the genome show a much higher sensitivity to DNase I than described above (43). These are due to the absence of a canonical nucleosome and the presence of gene regulatory (sequence-specific) proteins. They often occur within a domain of general DNase I sensitivity, coinciding with promoters or enhancers, and may also occur at the boundaries of chromosomal domains. Hypersensitive sites at promoters may be constitutive or inducible,

depending on the gene, and are often flanked by “positioned nucleosomes”.

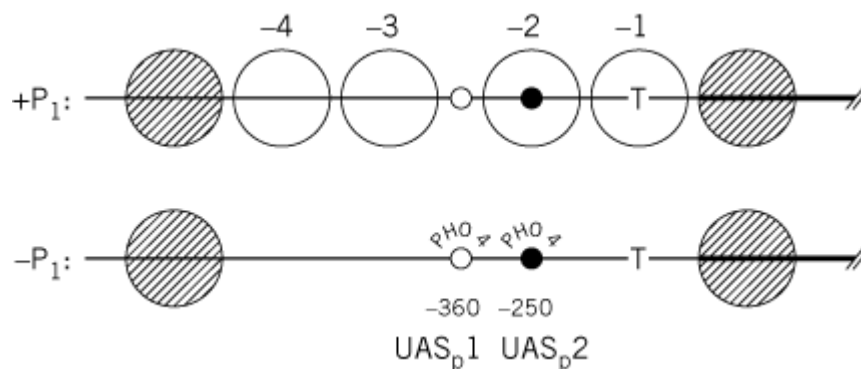
The presence of a nucleosome at a promoter generally inactivates it (39). Promoters therefore either have to be kept clear of nucleosomes, or one or more nucleosomes has to be perturbed or disrupted to permit assembly of a preinitiation complex and recruitment of RNA polymerase. Promoters may be kept clear either by flanking positioned nucleosomes or by the presence of sequence-specific proteins that bind during replication—for example, the Grf2 protein in yeast, which binds to an extended region of DNA and precludes nucleosome formation (39). The promoters that are packaged in nucleosomes have to be opened up (“remodeled”) in order to permit assembly of a preinitiation complex. The nature and extent of the disruption required to open up a promoter is unclear and may be different in different cases, perhaps depending on where the factor-binding sequences are within a nucleosome (eg, whether they can be exposed simply by loss of H1 or by loss of an H2A-H2B dimer, or whether they also require disruption of the central tetramer in the core of the nucleosome; see Ref. 2). It does not appear to result in complete histone loss. Promoter remodeling and disruption *in vivo* appears to require the action of energy-dependent “chromatin remodeling machines” (see text below) which may be recruited to particular promoters in response to the initial signal for activation, and possibly also histone acetyltransferases (see [Histone Acetylation](#)).

*In vitro* many transcription factors have been shown to be able to can bind to their cognate sites on the surface of a nucleosomes, albeit with lower affinity than to free DNA. The effects can be accounted for by transient “site exposure” as the interactions between the octamer and DNA fluctuate (44); this model predicts that binding of factors to multiple juxtaposed binding sites on the same nucleosome will be cooperative (45), as indeed already observed. Binding occurs without global displacement of the DNA, which will remain anchored at sites not bound by transcription factors. Consistent with this, a recent report shows that the erythroid transcription factor GATA-1 causes extensive disruption of histone–DNA contacts in nucleosome core particles containing GATA-1 sites, except over 50 bp around the dyad that serves to anchor the DNA (46). The complex is stable in solution, and the disruption is reversed on removal of GATA-1. The inherent flexibility of the nucleosome, and the fact that many promoters have clusters of factor-binding sites, suggests that this mechanism of gaining access to binding sites within nucleosomes at promoters could have physiological relevance.

At some inducible promoters, the binding site for the transcription factor that initiates the remodeling events is already exposed, in the linker between two nucleosomes in a short positioned array, rather than within a nucleosome. This is the situation at the *PHO5* promoter in yeast (19). Activation of the promoter in response to low phosphate relies on the accessibility in the uninduced state of a (weak) binding site for a transcription factor (Pho4) in the 70-bp gap between two of the four positioned nucleosomes at the promoter (Fig. 2). Binding of Pho4 leads to disruption of two nucleosomes on either side of the binding site and exposure of a second, stronger Pho4 binding site, as well as binding sites for another transcription factor, Pho2, which binds cooperatively with Pho4. The nature of the disruption is not clear but is reflected in the accessibility to a [restriction enzyme](#) of a site normally incorporated into a nucleosome and probably involves a “remodeling complex” (see text below). In other cases, a nucleosome at a promoter may be beneficial or even essential for activation of transcription. In the mouse mammary tumor virus (MMTV) [long terminal repeat](#) (LTR) promoter, six positioned nucleosomes place two [glucocorticoid response elements](#) in rotational positions (see text above) on one of the nucleosomes, such that they are accessible to glucocorticoid receptor binding in response to hormone; the two sites are also in close proximity on one face of the nucleosome. The binding site for the transcription factor NF1 is occluded until the nucleosomes are disrupted or perturbed in response to receptor binding (20). In a slightly different scenario, wrapping DNA around a nucleosome may bring into proximity two sites outside the nucleosome core, which are separated by about a nucleosome repeat length of DNA. There are several instances of this (2)—for example, the *Drosophila hsp26* promoter (47), where a critical nucleosome upstream of the start site of transcription performs such a scaffolding role and brings binding sites for **heat shock** transcription factors into proximity with each other and with TFIID, to create a preset promoter, with bound RNA polymerase, waiting only for the binding of heat shock transcription factor to induce it.

The positioned nucleosome is flanked by binding sites for the GAGA protein, which appears to play a crucial role in the architecture of the promoter, by helping to determine the position of the nucleosome (2).

**Figure 2.** The yeast *PHO5* promoter: chromatin structure in the uninduced (+Pi) and induced (-Pi) states. Large open circles (-1 to -4): the four positioned nucleosomes that are disrupted and become “transparent” upon induction. Small circles: Pho4 binding sites [UASp1 (open) and UASp2 (filled)]. T: TATA box. The thick horizontal line marks the beginning of the coding sequence. From Ref. 75, with permission.



In contrast to initiation, transcriptional elongation can proceed through nucleosomes *in vivo*. In a model system using SP6 polymerase and a single nucleosome, the histone octamer is displaced during elongation (39) and then rebinds upstream of its initial position; the mechanism proposed is transferred from in front of the advancing polymerase to behind it, via a DNA loop that is displaced (48). However, no such studies have yet been carried out with the large eukaryotic RNA polymerases. The indications are that a reconstituted nucleosome array presents a block to elongation by RNA polymerase II, which may be relieved by an energy-dependent mechanism (see next section).

## 5. Chromatin Remodeling Machines

Exciting discoveries in the last couple of years have revealed several multiprotein complexes (molecular masses ranging from 0.5 to 2 MDa) in *Drosophila*, yeast, and human cells that are able to remodel repressive chromatin in an ATP-dependent manner to permit transcription initiation at promoters *in vitro*. These include three complexes derived from a *Drosophila* embryo extract: NURF (nucleosome remodeling factor), CHRAC (chromatin accessibility complex) and ACF (ATP-dependent chromatin assembly and remodeling factor); the *Saccharomyces cerevisiae* SWI-SNF (switching-sucrose nonfermenting) and related RSC (remodeling the structure of chromatin) complexes; and SWI/SNF-related complexes in human cells and *Drosophila*. They have been comprehensively discussed in a number of recent reviews (49-53) that cite the original literature. They all contain an ATPase activity; in the case of the three *Drosophila* complexes, this is the same ISWI (imitation SWI2) subunit. The various complexes were purified initially using a range of functional assays to follow the purification, and they appear to be functionally and mechanistically distinct; they are likely to differ in their modes of nucleosomal perturbation. NURF (four protein subunits) was purified as a factor required for ATP-dependent remodeling of chromatin templates so that the GAGA protein is able to bind (GAGA binds to GA-rich sites in several *Drosophila* heat-shock promoters, such as the *hsp26* promoter, already mentioned); and it facilitates transcriptional activation. Interestingly, the 55-kDa subunit of NURF (which contains “WD repeats,” which are protein-interaction motifs) is also found in the *Drosophila* chromatin assembly factor dCAF-1 (54), and human homologues are components of histone acetylase and deacetylase complexes (see [Histone Acetylation](#)). CHRAC (five subunits) was purified as a factor that generated ATP-dependent



accessibility of sites in a chromatin template, reconstituted *in vitro*, to restriction enzymes. It is also able to generate a regular array from an irregular array of nucleosomes and has been described as an energy-dependent nucleosome spacing factor; it may therefore play a role in chromatin assembly (see text below) (where newly assembled nucleosomes are initially deposited in an irregular array that “matures” to give a regular spacing), although a spacing activity might, of course, also promote transcription factor binding. The third *Drosophila* complex, ACF (four subunits), was also purified as an ATP-dependent nucleosome spacing factor in a fractionated system and plays a role in nucleosome assembly; it can also mediate transcription-factor-mediated chromatin disruption in much the same way as NURF. SWI/SNF (9 to 12 subunits) facilitates transcription factor binding, and both SWI/SNF and RSC (15 subunits) disrupt the repeating 10-nucleotide cleavage pattern given by DNase I on mononucleosomes, showing disruption of the rotational setting of the DNA on the octamer surface. In yeast, the SWI/SNF complex is not abundant and not essential for viability, whereas the RSC complex is abundant and essential. SWI/SNF appears to be necessary for the transcription of only a small subset of genes, and in an *in vitro* system the dependence of transactivation on SWI/SNF activity is much reduced if there are multiple binding sites for the activator on the same nucleosome (55). It is possible that *in vivo* the SWI/SNF complex is needed only for remodeling of weak promoters. Altered forms of mononucleosomes generated by both SWI/SNF (56, 57) and RSC (58) complexes have recently been reported. They retain a full complement of histones and have altered physical properties that would be consistent with either a heavily modified mononucleosome or a dinucleosome. The SWI/SNF products have a higher affinity for transcription factors than unmodified nucleosomes. Much remains to be done to (a) elucidate the mechanisms by which the various chromatin remodeling machines act on chromatin, (b) define the roles of the various subunits, and (c) establish which mechanisms, *in vivo*, play a role in promoter disruption for transcription initiation, and whether others are primarily involved in reversing such disruption or in nucleosome assembly, and so on. The relationship between these complexes and the large acetylase, and deacetylase, complexes (see [Histone Acetylation](#)) also remains to be clarified, although the indications are that they cooperate in the remodeling of certain promoters.

The foregoing account has dealt solely with complexes that disrupt chromatin at promoters, allowing transcription initiation complexes to form. The transcriptional machinery also confronts nucleosomes during RNA chain elongation. Experiments on a reconstituted array of nucleosome cores show that the passage of eukaryotic RNA polymerase II is blocked by a nucleosome core and that a heterodimeric protein isolated from human cells will facilitate RNA chain elongation in a chromatin context. The protein has been designated FACT (*facilitates chromatin transcription*) (59) and appears to be a [DNA-binding protein](#), one subunit being an HMG-box protein. Little is yet known about it how it might work, but its action is not ATP-dependent. The *in vitro* assays show a rather modest extension of the transcript (50 to 150 nucleotides) before FACT stalls, and it is not yet clear whether FACT is designed to work only on the promoter-proximal nucleosome or needs cofactors for its normal action, or whether the modest extension is merely due to a technical limitation of the experiment. The chromatin template contains closely packed nucleosome cores (no linker DNA, no H1), and if octamer transfer via a loop of DNA is required, as demonstrated in model systems (48) (see text above), this may not be possible here.

## 6. Chromatin and Repression

There are hierarchies of repression, through chromatin structure, as there are of transcriptional activation. The most basic level, promoter occlusion by a nucleosome, was discussed above. Superimposed on this is repression through folding of the nucleosome filament into a 30-nm filament (whether this occurs in yeast is unclear). The higher-order structure is stabilized by a hypoacetylated state of the *N*-terminal tails of the core histones (see [Histone Acetylation](#)). In organisms other than yeast and *Drosophila*, DNA **methylation** is known to contribute to gene repression through chromatin structure, and a **5-methylcytosine**-binding protein (MeCP2) has recently been reported to recruit deacetylases (60, 61). This could favor a more stable 30-nm filament, which might be further stabilized by the tighter binding of H1 to methylated DNA (62). The 30-nm filament is not a template for transcription, and in some cases, 30-nm filament stabilization (eg, by specialized linker

histone variants) and packing in the nucleus may be at least major players in general repression. This may be the situation in mature nucleated erythrocytes (eg, of birds, fish, and amphibia), where H1 is largely replaced by histone H5 and where there is a global and terminal shutdown of transcription; and in sea-urchin sperm, where highly condensed chromatin is linked to the presence of sperm-specific H1. Other more complex mechanisms of repression also exist, however, resulting in [heterochromatin](#). The essence of these mechanisms—although they differ in many respects in, for example, yeast and *Drosophila*, where they have been extensively studied—is the assembly, over particular chromosomal regions, of repressive multiprotein complexes that interact with chromatin. In the case of *Drosophila* at least, the chromatin is probably in the form of the 30-nm filament, and the mechanisms add an additional layer of repression for genes whose permanent inactivation is crucial and where inappropriate expression cannot be tolerated.

Yeast chromosomes are too small to allow heterochromatin to be seen cytologically, but nonetheless some chromosomal regions show many of the features of heterochromatin in more complex eukaryotes (replication late in S phase, localization at the nuclear envelope, and [position effects](#) on gene expression that are inherited **epigenetically**). In *Saccharomyces cerevisiae*, such regions occur at the silent mating type loci, *HML* and *HMR*, and at [telomeres](#) (the ends of the chromosomes; genes placed here are repressed) and have been extensively studied by genetic analysis (63). In each case, a multiprotein complex of Sir2, Sir3, and Sir4 proteins (Sir = silent information regulator information regulator) interacts with the *N*-terminal regions of histones H3 and H4. At the telomeres, which contain the sequence  $[C_{1-3}A]_n$ , the **silencers** consist of multiple binding sites for RAP1 (repressor activator protein 1), whereas silencers at *HML* and *HMR* consist of binding sites for ABF1 (*ARS*-binding factor 1), the origin recognition complex (ORC) and RAP1. Residues 16 to 29 of H4 are required for interaction with Sir3 and Sir4 and for silencing, presumably by providing interactions necessary for stabilization of the multiprotein complex. In addition, although the core histones at the silenced loci in yeast are hypoacetylated, there is a specific requirement for Lys12 in acetylated form, possibly for interaction with Sir3 (63). In the current model for telomeric repression, RAP1 binds to the telomeric repeats, initiating polymerization of Sir3/Sir4 over the adjacent chromatin domain, the binding of the Sir proteins being mediated by interaction with the H3 and H4 *N*-termini, possibly stabilized by acetylation of H4 at Lys12. Other repressive mechanisms in yeast also involve the *N*-terminal tails of H3 and H4 (64).

Extensive genetic studies of heterochromatin have also been carried out in *Drosophila*. Genes brought into juxtaposition with heterochromatin (eg, at centromeres) may be silenced in a subset of cells that normally express the gene—an effect known as position effect variegation (PEV) (2). Chromatin structure is altered in variegating genes, as shown by the properties of the *hsp26* transgene inserted at various chromosomal locations; the heterochromatic *hsp26* gene has a noninducible promoter and is packaged in an unusually regular nucleosomal array (65). Modifiers of PEV, identified genetically, are either (a) structural proteins that are components of multimeric protein complexes that interact with chromatin or (b) proteins that could play indirect roles in regulating the formation of chromatin—for example, by modifying histones or DNA. One such protein is the histone deacetylase RPD3 (see [Histone Acetylation](#)). However, mutations in this protein that increase acetylation *increase* PEV, rather than decreasing it, as might be expected; this is probably because of the requirement for acetylation at Lys12, as in yeast heterochromatin, against a background of general histone hypoacetylation. The multiprotein complexes that affect PEV, and presumably stabilize chromatin higher-order structure, contain (a) HP1 (heterochromatin protein 1), which binds to chromatin but not to DNA, and (b) the protein termed Su(var)3–7, which interacts with HP1 and chromatin (65). The principles underlying PEV in *Drosophila* seem also to apply to the stable inactivation of **euchromatic** genes needed for maintaining differential expression of developmental regulators (eg, the [homeotic genes](#) that control the segmental identity of the insect body). Locally silenced heritable chromatin structures are generated by the assembly of multiprotein complexes that contain members of the [Polycomb group](#) (Pc-G) of proteins, which are homologues of HP1 and contain a “chromodomain” (66). The Pc-G genes are responsible for maintaining long-range repression of homeotic genes over multiple enhancers in the later stages of *Drosophila*

development. Like HP1, Pc-G proteins appear not to bind DNA. However, specific DNA sequences that recruit Polycomb group proteins and associated silencing factors [Polycomb response elements; (PREs)] have been identified and appear to be essential for stabilizing determined expression states, which are usually (but not invariably) repressive (see Ref. 66). Various models have been proposed for the mechanism of repression, including propagation of a repressive array of interacting proteins over chromosomal domains (the “spreading model” by analogy with the “spreading” effect of heterochromatin into euchromatin in PEV), an alternative model has been proposed in which clusters of Pc-G proteins interact cooperatively with each other to create tethered loop domains of chromatin which ensure repression (67). Mammalian homologues of HP1-like proteins have been found, as have interacting partners, and it seems likely that the same mechanism of silencing, by assembly into repressive heterochromatin, may be used to regulate some developmentally pivotal mammalian genes whose leaky expression would be disastrous for the cell.

Mammalian [X-Chromosome Inactivation](#) is an extreme case of heterochromatin formation over a whole chromosome. The precise molecular basis is not known, but there are some notable features of the chromatin composition: hypoacetylation of histone H4; the presence of an unusual H2A variant, macroH2A (see [Histones](#)); heavy methylation of the DNA; and the presence of an unusual large untranslated RNA, termed Xist, which is expressed from the inactive X-chromosome, and which is of prime importance in an early event in X-inactivation (68). It would not be too surprising if X-inactivation turned out to have at least some features in common with telomeric and position effect variegation in yeast and *Drosophila*.

## 7. Replication-Linked Chromatin Assembly

During S-phase of the cell cycle in eukaryotes, it is not just the DNA that undergoes replication, but the chromatin as a whole. Chromatin assembly is tightly coupled to [DNA replication](#) and the passage of the [replication fork](#). It involves assembly of nucleosomes on daughter strands, followed by conversion of an initially irregularly spaced array of nucleosomes to a regular array—a process known as maturation. There has been recent progress in understanding the participants in both of these processes. Assembly involves, in equal amounts, both old (parental) histones in preexisting nucleosomes, which are recycled, and newly synthesized histones, which are apparently randomly segregated on to the two nascent duplexes. The temporal order of deposition of newly synthesized histones on to DNA reflects the structural organization of the nucleosome; it begins with the deposition of H3 and H4 to form the central organizing tetrameric core, followed by H2A and H2B as two dimers, and finally addition of H1.

The newly synthesized H4 deposited at the replication fork is diacetylated, at lysines 5 and 12, in all species so far examined (see [Histone Acetylation](#)). Newly synthesized H3 is apparently also modified (acetylated or phosphorylated?), although perhaps not in all cell types. The role, if any, played by acetylation *per se* in the assembly process is not well understood. Acetylation is transient and is erased by deacetylases shortly after deposition of histones on to the DNA; one possibility might be that any tendency for premature formation of internucleosome contacts is avoided by the acetylation (see [Nucleosome](#)). The deposition of acetylated H3 and H4 involves the participation of chromatin assembly factor 1 (CAF-1), (69), which acts as a core histone **molecular chaperone**. It forms a stable complex with newly synthesized acetylated H3 and H4 and somehow targets them specifically to the replication fork (70); the histone chaperone NAP1 appears to handle H2A and H2B in *Drosophila* (71). Other proteins may also be involved in the assembly process. CAF-1 purified from human cells contains three subunits, indicated by their mass in kDa: p48, p60, and p150; homologues exist in yeast and *Drosophila*. The small subunit (p48) of human CAF-1 is homologous to the regulatory subunit (p46) of the human B-type histone acetyltransferase (responsible for deposition-related acetylation); both subunits bind to histone H4 and, significantly, at a site that is accessible only in the free histone and not in chromatin (as shown by the structure of the nucleosome core particle; see [Nucleosome](#)). Strikingly, similar subunits also occur in chromatin remodeling machines (eg, the p55 subunit of *Drosophila* NURF) as noted, as well as in histone deacetylases (see [Histone Acetylation](#)), leading to the suggestion that these subunits target the

various complexes to their histone substrates in a manner that is regulated by nucleosomal DNA (72). They are members of the highly conserved family of proteins that contain “WD repeats,” multiple distinctive sequence motifs that appear to be involved in [protein–protein interactions](#). Because the site of interaction of CAF-1 with H4 is accessible only in the free histone, CAF-1 would have to be displaced before nucleosome assembly could occur; the interactions appear to be disrupted by ACF (ATP-requiring chromatin assembly and remodeling factor; see text above). ACF was initially identified as a factor capable of regularly spacing an irregularly spaced nucleosome array (51) and would be well-suited to this role during the maturation step of chromatin assembly. It is also potentially able to participate in the remodeling of nucleosomes that accompanies transcriptional activation (see text above: chromatin remodeling machines).

## 8. Chromatin Reconstitution

Reconstitution *in vitro* of chromatin from its component parts, like any complex biochemical structure, would allow questions to be asked about structure–function relationships. Reconstitution of the nucleosome core is straightforward. It is achieved by gradual dialysis of a 1:1 molar ratio of histone octamer and 146 bp DNA from 2M NaCl to low salt (eg, 10 mM NaCl). Nucleosome core particles (see [Nucleosome](#)) for X-ray crystallography were made in this way. The histones can be presented either as the intact octamer or as a 1:1 molar ratio of H2A–H2B dimers and (H3)<sub>2</sub>(H4)<sub>2</sub> tetramers.

Long “chromatin” reconstituted *in vitro* by mixing histones at high ionic strength (eg, 2M NaCl), and then dialyzing to low ionic strength, shows close packing of octamers, irrespective of the presence of H1, which has no effect because it is the last histone to bind in this procedure. In attempts to determine what cellular factors determine nucleosome spacing (nucleosome repeat length), several cell-free extracts, from *Xenopus* eggs and oocytes and from *Drosophila* embryos, have been developed that will assemble plasmid DNA, irrespective of sequence, into “physiologically” spaced chromatin in an ATP-dependent fashion, independent of replication. Nucleosome spacing (and hence linker length) is increased by H1. The earliest extracts were from *Xenopus* eggs, which contain a large maternal histone pool, and led to the discovery of two acidic histone chaperone proteins, nucleoplasmin and N1/N2, which facilitated chromatin assembly in this system. Work on the *Xenopus* system is summarized in Ref. 2 (Section 3.4.2). The true role of these proteins may be storage of the large histone pool in the egg. Whether similar proteins exist in somatic cells is unclear, although nucleoplasmin-like proteins have been reported. The *Drosophila* embryo extracts are the same ones that led to the discovery of CHRAC (53) and ACF (51) (see text above: Chromatin remodeling) which themselves have nucleosome spacing activity. Spacing may also have an electrostatic component, involving neutralization of charges on linker DNA (73), and a connection between nucleosome spacing and formation of higher-order structure, which is also ionic strength dependent, has been suggested.

A well-defined *in vitro* system starting with pure components, which has recently been further explored (74), will also assemble “properly spaced” chromatin, in an H1-dependent, ATP-independent manner. Irregularly spaced octamers are first deposited on the DNA by dialysis from high to low NaCl concentration, and then a more regularly spaced array is generated by incubation at “physiological ionic strength” (150 mM NaCl) with H1(H5) in the presence of polyglutamate (which probably binds histones reversibly and allows equilibrium to be reached in the histone positions on the DNA). Interestingly, the system is sensitive to the DNA sequence, and regular nucleosome arrays were generated on only about half of the ~2-kbp cloned chicken genomic sequences tested; the sequences tested (by **Southern blotting**) behaved in essentially the same way in chicken liver nuclei. It is argued (74) that genomic chromatin is a mosaic of less well-ordered regions and regular regions that are about 10 nucleosomes long and contain a nucleosome positioning sequence (or sequences), which effectively acts as a boundary (18) against which linker-histone-dependent nucleosome alignment can occur. It is possible that egg, oocyte, and embryo extracts, which appear to be indifferent to DNA sequence as well as ATP-dependent, are designed to assemble chromatin through

rapid cycles of cell division, and not to be responsive, for example, to genomic signals that may be necessary to specify arrays of nucleosomes compatible with gene expression and regulation in somatic cells.

## Bibliography

1. K. van Holde (1988) *Chromatin*, Springer-Verlag, New York.
2. A. P. Wolffe (1995) *Chromatin; Structure and Function*, 2nd ed., Academic Press, New York.
3. S. C. R. Elgin (ed.) (1995) *Chromatin Structure and Gene Expression* [Frontiers in Molecular Biology series (B. D. Hames and D. M. Glover, eds.)], Oxford University Press, Oxford, U.K.
4. R. D. Kornberg (1977) *Annu. Rev. Biochem.* **46**, 931–954.
5. K. Luger, A. W. Mader, R. K. Richmond, D. F. Sargent, and T. J. Richmond (1997) *Nature* **389**, 251–260.
6. Y.-B. Zhou et al. (1999) *Nature* (in press).
7. A. A. Travers (1987) *Trends Biochem. Sci.* **12**, 108–112.
8. T. E. Shrader and D. M. Crothers (1989) *Proc. Natl Acad. Sci. USA* **86**, 7418–7422.
9. H. R. Widlund et al. (1997) *J. Mol. Biol.* **267**, 807–817.
10. P. T. Lowary and J. Widom (1998) *J. Mol. Biol.* **276**, 19–42.
11. H. Cao, H. R. Widlund, T. Simonsson, and M. Kubista (1998) *J. Mol. Biol.* **281**, 253–260.
12. J. Widom (1996) *J. Mol. Biol.* **259**, 579–588.
13. S. Muyldermans and A. A. Travers (1994) *J. Mol. Biol.* **235**, 855–870.
14. A. Travers and H. Drew (1999) *Biopolymers* (in press).
15. J. J. Hayes, D. J. Clark, and A. P. Wolffe (1991) *Proc. Natl. Acad. Sci. USA* **88**, 6829–6833.
16. C. Spangenberg et al. (1988) *J. Mol. Biol.* **278**, 725–739.
17. S. Tanaka, M. Zatchej, and F. Thoma (1992) *EMBO J.* **11**, 1187.
18. R. D. Kornberg and L. Stryer (1988) *Nucleic Acids Res.* **16**, 6677–6690.
19. J. Svaren and W. Horz (1997) *Trends Biochem. Sci.* **22**, 93–97.
20. T. Archer, M. G. Cordingley, R. G. Wolford, and G. L. Hager (1991) *Mol. Cell. Biol.* **11**, 688–698.
21. G. Panetta et al. (1999) *J. Mol. Biol.* (in press).
22. P. Bouvet, S. Dimitrov, and A. P. Wolffe (1994) *Genes Dev.* **8**, 1147–1159.
23. G. Meersseman, S. Pennings, and E. M. Bradbury (1992) *EMBO J.* **11**, 2951–2959.
24. S. Pennings, G. Meersseman, and E. M. Bradbury (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10275–10279.
25. K. Ura, J. J. Hayes, and A. P. Wolffe (1995) *EMBO J.* **14**, 3752–3765.
26. D. J. Clark and T. Kimura (1990) *J. Mol. Biol.* **211**, 883–896.
27. G. Felsenfeld and J. D. McGhee (1986) Structure of the 30 nm fibre. *Cell* **44**, 375–377.
28. J. Widom (1989) *Annu. Rev. Biophys. Biomol. Struct.* **18**, 365–395.
29. K. van Holde and J. Zlatanova (1996) *Prog. Nucleic Acid Res. Mol. Biol.* **52**, 217–259.
30. V. Ramakrishnan (1997) *Crit. Rev. Eukaryot. Gene Expr.* **7**, 215–230.
31. J. Widom (1998) *Annu. Rev. Biophys. Biomol. Struct.* **27**, 285–327.
32. F. Thoma, T. Koller, and A. Klug (1979) *J. Cell Biol.* **83**, 403–427.
33. V. Graziano, S. E. Gerchman, D. K. Schneider, and V. Ramakrishnan (1994) *Nature* **368**, 351–354.
34. C. L. F. Woodcock, L.-L. Frado, and J. B. Rattner (1984) *J. Cell Biol.* **99**, 42–52.
35. C. L. Woodcock and R. Horowitz (1995) *Trends Cell Biol.* **5**, 272–277.
36. J. D. McGhee, J. M. Nickol, G. Felsenfeld, and D. C. Rau (1983) *Cell* **33**, 831–841.

37. J. O. Thomas, C. Rees, and P. J. G. Butler (1986). *Eur. J. Biochem.* **154**, 343–348.
38. R. T. Kamakaka and J. O. Thomas (1990). *EMBO J.* **9**, 3997–4006.
39. R. D. Kornberg and Y. Lorch (1992) *Annu Rev. Cell Biol.* **8**, 563–589.
40. G. Felsenfeld (1992) *Nature* **335**, 219–224.
41. T. R. Hebbes, A. W. Thorne, and C. Crane-Robinson (1988) *EMBO J.* **7**, 1395–1402.
42. P. Geyer (1997) *Curr. Opin. Genet. Dev.* **7**, 242–248.
43. S. Elgin (1988) *J. Biol. Chem.* **263**, 19259–19262.
44. K. J. Polach and J. Widom (1995) *J. Mol. Biol.* **254**, 130–149.
45. K. J. Polach and J. Widom (1996) *J. Mol. Biol.* **258**, 800–812.
46. J. Boyes et al. (1998) *J. Mol. Biol.* **279**, 529–544.
47. G. H. Thomas and S. C. R. Elgin (1988) *EMBO J.* **7**, 2191–2201.
48. V. M. Studitsky, D. J. Clark, and G. Felsenfeld (1994) *Cell* **76**, 371–382.
49. T. Tsukiyama and C. Wu (1997) *Curr. Opin. Genet. Dev.* **7**, 182–191.
50. M. J. Pazin and J. T. Kadonaga (1997) *Cell* **88**, 737–740.
51. T. Ito, J. K. Tyler, and J. T. Kadonaga (1997) *Genes to Cells* **2**, 593–600.
52. B. Cairns (1998) *Trends Biochem. Sci.* **23**, 20–25.
53. P. D. Varga-Weisz and P. B. Becker (1998) *Curr. Opin. Cell Biol.* **10**, 346–353.
54. M. A. Martinez-Balbas, T. Tsukiyama, D. G. Dula, and C. Wu (1998) *Proc. Natl. Acad. Sci. USA* **95**, 132–137.
55. R. T. Utley, J. Cote, T. Owen-Hughes, and J. L. Workman (1997) *J Biol. Chem.* **272**, 12642–12649.
56. J. Coté, C. L. Peterson, and J. L. Workman (1998) *Proc. Natl. Acad. Sci USA* **95**, 4947–4952.
57. G. Schnitzler, S. Sif, and R. E. Kingston (1998) *Cell* **94** 17–27.
58. Y. Lorch, B. R. Cairns, M. Zhang, and R. D. Kornberg (1998) *Cell* **94**, 29–34.
59. G. Orphanides et al. (1998) *Cell* **92**, 105–116.
60. X. Nan et al. (1998) *Nature* **393**, 386–389.
61. P. L. Jones et al. (1998) *Nat. Genet.* **19**, 187–191.
62. M. A. McArthur and J. O. Thomas (1996) *EMBO J.* **15**, 1705–1714.
63. M. Grunstein (1998) *Cell* **93**, 325–328.
64. M. Grunstein (1997) *Nature* **389**, 349–352.
65. L. L. Wallrath (1998) *Curr. Opin. Genet. Dev.* **8**, 147–153.
66. G. Cavalli and R. Paro (1998) *Curr. Opin. Cell Biol.* **10**, 354–360.
67. V. Pirrotta (1998) *Cell* **93**, 333–336.
68. A. M. Keohane, J. S. Lavender, L. P. O'Neill, and B. M. Turner (1998) *Dev. Genet.* **22**, 65–73.
69. S. Smith and B. Stillman (1989) *Cell* **58** 15–25.
70. A. Verreault, P. D. Kaufman, R. Kobayashim, and B. Stillman (1996) *Cell* **87**, 95–104.
71. T. Ito, M. Bulger, R. Kobayashi, and J. T. Kadonaga (1996) *Mol. Cell. Biol.* **16**, 3112–3124.
72. A. Verreault, P. D. Kaufman, R. Kobayashi, and B. Stillman (1997) *Curr. Biol.* **8**, 96–108.
73. T. A. Blank and P. B. Becker (1995) *J. Mol. Biol.* **252**, 305–313.
74. K. Liu and A. Stein (1997) *J. Mol. Biol.* **270**, 559–573.
75. P. D. Gregory and W. Horz (1998) *Eur. J. Biochem.* **251**, 9–18.

### **Suggestions for Further Reading**

76. S. C. R. Elgin (1996) Heterochromatin and gene regulation in *Drosophila*. *Curr. Opin. Genet. Dev.* **6**, 193–202.

77. G. Felsenfeld (1996) Chromatin unfolds. *Cell* **86**, 13–19.
78. E. P. Geiduschek (1998) Chromatin transcription: clearing the gridlock. *Curr. Biol.* **8**, R373–R375.
79. P. D. Gregory and W. Horz (1998) Chromatin and transcription. *Eur. J. Biochem.* **251**, 9–18.
80. P. Kaufman (1996) Nucleosome assembly: the CAF and the HAT. *Curr. Opin. Cell Biol.* **8**, 369–373.
81. R. E. Kingston, C. A. Bunker, and A. N. Imbalzano (1996) Repression and activation by multiprotein complexes that alter chromatin structure. *Genes Dev.* **10**, 905–920.
82. K. Struhl (1998) Histone acetylation and transcriptional regulatory mechanisms. *Genes Dev.* **12**, 599–606.

## Chromatofocusing

Biological [macromolecules](#) differ in their compositions of **buffering** groups, such as [amino acids](#) and **purine** or **pyrimidine** bases, thus giving each biomolecule its own unique charge properties. Generally, at low or acidic pH values, most biomolecules have a net positive charge, whereas at a higher or basic pH they carry a net negative charge. Of course, at a certain pH, the isoelectric pH or [isoelectric point](#) ( $pI$ ), these biomolecules do not possess a net electric charge. This unique isoelectric pH provides the basis for the selectivity of chromatofocusing in **chromatographically** separating biomolecules.

Chromatofocusing is a special kind of technique classified within the [ion-exchange chromatography](#) of biomolecules. A weak ion-exchange matrix and a pH gradient are used in chromatofocusing, instead of a strong ion exchanger and a salt gradient as in normal ion-exchange chromatography. This technique was first described theoretically and experimentally demonstrated by Sluyterman and co-workers ([1](#), [2](#)). They proposed that a pH gradient could be produced in an ion-exchange column packed with an appropriate ion-exchange resin with good buffering capacity. Although a pH gradient in a column can be produced in a manner similar to that of a salt gradient by using two different pH buffers in a mixing chamber of a gradient maker, a pH gradient can be created internally in the column by taking advantage of the buffering capacity of the weak ion-exchange resin. In practice, a certain pH buffer is used to equilibrate a column packed with a weak ion exchanger. Then another buffer with a different pH is passed through the column, which generates a pH gradient in the column. If such a pH gradient is used to elute biomolecules bound to the ion-exchange resin, the biomolecules elute in order of their isoelectric pH.

The mechanism of chromatofocusing is based on the buffering action of the charged groups on the ion-exchange resin and on the fact that a biomolecule has a net negative charge at a pH above its isoelectric point. In a descending pH gradient, a single molecular species exists in three charged states—negative, neutral, and positive. When a positively charged column (ie, packed with an anion exchange resin) is equilibrated with a starting buffer of high pH, biomolecules that become negatively charged are initially retained on the column. When an elution buffer of lower pH is passed through the column, a pH gradient develops and the individual molecules continuously change their charged states. The molecules at the rear of the sample zone are the first to be titrated by the low-pH buffer. These molecules become positively charged when the pH is less than their  $pI$ , so that they are repelled from the column matrix and are carried rapidly to the front of the sample zone, because of the high velocity of the moving buffer. In traveling to the front of the sample zone, the molecules encounter an increase in pH and are titrated from their positive form to neutrality and back

to their negative form. Once the molecules become negatively charged again, they re-adsorb on the matrix and again fall back to the rear of the sample zone. The cycle between the front and rear of the sample zone results in “focusing” (ie, a continuous narrowing of this zone) until the molecules elute from the column. At this point, the pH of the column effluent is approximately the same as *pI* of the components eluting.

Chromatofocusing is an analytical or preparative technique for separating biomolecules according to their *pI*. The details of this technique relating to column packing, sample preparation, and sample application used in this technique are similar to those for [affinity chromatography](#). A detailed operational protocol is beyond the scope of this article, so the interested reader is directed to other reviews ([3](#), [4](#)).

There are a number of advantages in using chromatofocusing. In this technique, a biomolecule is not subjected to a pH greater than its *pI*, and the resulting focusing effects concentrate the sample into a sharp, highly resolved band. One of the great benefits of chromatofocusing is its ease of operation. No gradient-forming devices or mixers are required. The pH gradient is formed with a single isocratic eluent. Chromatofocusing is used widely in research as the method of choice for resolving **isozymes** and molecular species with very similar charge characteristics, such as [transferrin](#), [ferritin](#), and [hemoglobins](#) [see refs. in [3-5](#)].

#### Bibliography

1. L. A. A. Sluyterman and O. Elgersma (1978) *J. Chromatogr.* **150**, 17–30.
2. L. A. A. Sluyterman and J. Wijdenes (1978) *J. Chromatogr.* **150**, 31–44.
3. L. Giri (1990) in *Guide to Protein Purification* (M.P. Deutscher, ed.), *Methods in Enzymology* **182**, Academic Press, New York, pp. 380–392.
4. T. W. Hutchens (1989) in *Protein Purification: Principles, High Resolution Methods, and Applications* (J.-C. Janson, and L. Rydén, eds.), VCH, New York, pp. 149–174.
5. Pharmacia (1985) *FPLC Ion Exchange and Chromatofocusing-Principles and Methods*, Offsetcenter AB, Uppsala, Sweden.

## Chromatography

Chromatography constitutes a family of closely related methods for separating and analyzing a wide variety of chemicals. Today, despite developments in analytical chemistry that link scientists to many modern and extremely sophisticated devices, the classic methodology of chromatography still plays a very important role among analytical techniques. It is almost impossible to imagine a laboratory without chromatographic equipment. Quality control, product purification, and basic research are some of the fields in which chromatography is used.

### 1. Definition

The term chromatography, first used by Tsvet (1872–1919), a Russian botanist, derived from the two Greek words *Khromatos* (color) and *graphos* (written). He used the term chromatography to describe his studies on pigment separation using a chalk column ([1](#), [2](#)), thus defining chromatography as a method by which the components of a mixture were separated on an adsorbent column in a flowing system ([1](#), [2](#)). The International Union of Pure and Applied Chemistry (IUPAC) has further defined chromatography as



A method, used primarily for separation of the components of a sample, in which the components are distributed between two phases, one of which is stationary while the other moves. The stationary phase may be a solid, or a liquid supported on a solid, or a gel. The stationary phase may be packed in a column, spread as a layer, or distributed as a film, etc.; in these definitions chromatographic bed is used as a general term to denote any of the different forms in which the stationary phase may be used. The mobile phase may be gaseous or liquid (3).

## 2. Historical Perspective

Although some phenomena that form the basis of chromatographic methods have been known for a long time, Tsvet is generally referred to as the father of chromatography. In 1906, he described the separation of plant pigments by column liquid chromatography using over 100 adsorption media (1, 4, 5). During the next forty years after Tsvet's work, there were some important developments in the field. For example, Kuhn et al. published two papers in 1931 on separating carotenoids on a calcium carbonate column (6, 7). Kuhn, Karrer, and Ruzicka applied the chromatographic technique to their own fields of interest and were awarded the Nobel prize (1937, 1938, and 1939, respectively) for their contributions to chromatography. Tiselius (8) and Claesson (9) developed the now classical procedures involving the continuous observation of optical properties of solutions flowing out of chromatographic columns. Tiselius was awarded the Nobel prize in 1948 for his research on "Electrophoresis and Adsorption Analysis". The introduction of gradient elution in 1952 (10) was an important contribution to all column chromatographic methods.

Chromatographic developments greatly accelerated after the famous paper by Martin and Synge appeared in 1941 (11). They presented the invention of liquid-liquid (or partition) chromatography in columns and in planar form (paper chromatography). They also provided a theoretical framework for the basic chromatographic process. Martin and Synge were awarded the Nobel prize in 1952 for this work. When thin layers of supported silica gel were introduced as an alternative for paper in the late 1950s (12, 13), the field of **thin-layer chromatography** (TLC) was born and became so popular that it has largely replaced the older technique. Another main development in the progress of chromatography was the introduction of **gas-liquid chromatography** (GLC) by James and Martin (14), which had an unprecedented impact on the analytical chemistry of organic compounds. Porath and Flodin introduced **size-exclusion chromatography** in 1959 (15), allowing easy separation of macromolecules. Modern liquid chromatography (**HPLC**) was introduced in the early 1970s, permitting the efficient separation of a wide range of components.

The theory of chromatography was first studied by Wilson (16), who discussed the quantitative aspects in terms of **diffusion**, adsorption rate, and isothermal nonlinearity. The plate theory was first presented by Martin and Synge (11) and was further explored by Craig (17) and Gluechauf (18). In this theory, chromatography is described in terms closely analogous in its mode of operation to distillation and extraction fractionating columns. Lapidus and Admanson (19), followed by van Deemter and co-workers (20), developed the rate theory, an alternative to the plate theory. In this theory, column efficiency was described as a function of the mobile-phase's flow rate and diffusion properties and the stationary-phase particle size. In 1959, Giddings published another paper on this topic (21), and the rate theory has since become the backbone of chromatographic theory. In 1963, Giddings (22) pointed out that, if the efficiencies of gas chromatography were to be achieved in liquid chromatography, particle sizes of 2 to 20  $\mu\text{m}$  were required. This prediction was found to be correct with the development of HPLC systems. Numerous detailed descriptions of chromatographic theory exist in the literature. Notable examples include a monograph edited by Jönsson (23) and the excellent reviews of Snyder (24) and Horváth and Melander (25).

This is only a small glimpse of the historical development of chromatography. The book *75 Years of Chromatography—A Historical Dialogue* (26), which describes many of the individuals who deserve credit for developing this technique, can be consulted for more complete accounts.

## 3. Retention Mechanism Classification

There are three common ways to classify chromatographic methods. The first and most popular classification is based on the mechanisms of retention, the manner in which the analyte interacts with the stationary phase. In this classification, chromatography may be divided into the five following basic types:

### 3.1. Adsorption Chromatography

This technique is based on competition for neutral analytes between the mobile phase (gas or liquid) and a neutral solid adsorbent. Therefore, analytes with **polar groups** are retained longer by a polar adsorbent, and **nonpolar** analytes interact better with a nonpolar stationary phase. In this type of chromatography, the analytes are simultaneously in contact with both the stationary phase and the mobile phase.

### 3.2. Partition Chromatography

This technique is based on competition for neutral analytes between the mobile phase (gas or liquid) and a neutral liquid or liquid-like stationary phase (the latter is usually called a “bonded-phase” when long alkyl chains or their derivatives are bonded to a matrix and behave like a liquid). In partition chromatography, the analyte is transferred from the bulk of one phase into the bulk of the other, so that the analyte molecules are surrounded only by molecules of one phase. Separation in partition chromatography is achieved by differences in the **partition coefficients** of the analytes between the mobile and stationary phases.

### 3.3. Ion-Exchange Chromatography

This technique is based on the **electrostatic interaction** between a charged solute and an oppositely charged solid stationary phase. Separation in **ion-exchange chromatography** is achieved by the differing affinities of ions in solution for oppositely charged ionic groups in the stationary phase. Ion-exchange chromatography is applicable to any solute that acquires a charge in solution. Thus, even carbohydrates, which are largely uncharged below pH 12, are separated by this type of chromatography at sufficiently high pH.

### 3.4. Size-Exclusion Chromatography

This technique is based on the sieving principle and is variously known as gel chromatography, **gel filtration** and gel-permeation chromatography (see **Size Exclusion Chromatography**). In this technique, the stationary-phase particles have a wide range of pore sizes, causing the stationary phase to behave like a molecular sieve. Small molecules permeate the pores, and large bulky molecules are excluded. Thus, the solutes are separated on the basis of molecular weight and size, and the larger ones elute first.

### 3.5. Affinity Chromatography

This technique is based on the unique biological specificity of **ligand-binding** interactions (ie, the lock-and-key mechanism) (see **Affinity Chromatography**). The ligand is covalently bonded to the matrix that forms packing material for the column. Separation is achieved after the applied **macromolecule** becomes specifically, but not irreversibly, bound to the ligand. The macromolecule is eluted by altering the composition or pH of the eluent so as to weaken its interaction with the ligand, thus promoting dissociation and facilitating elution of the retained compounds.

Many other types of chromatography, such as **hydrophobic-interaction**, **chiral**, **ion-pair**, and **salting-out** chromatography, are also frequently used. Furthermore, in practical chromatography intermediate or mixed types are often used. So, although one dominant mechanism is presented, the chromatographic modes are not mutually exclusive. To know more about the relationships of different chromatographic modes, see a schematic diagram by Saunders ([27](#)).

## 4. Development Procedure Classification

The second classification is based on the development procedure, the mechanism by which the sample is removed from the column, and therefore depends on the nature of the mobile phase. This

classification was introduced by Tiselius (8) in 1940. There are three chromatographic modes in the classification: (1) elution development, (2) displacement development, and (3) frontal analysis. The principles of each mode are illustrated by Braithwaite and Smith (28). In practice, only elution development and to a lesser extent displacement development are commonly used.

## 5. Fractionation Phase Classification

The third classification is based on the phases between which the fractionation process takes place. In chromatography, one phase is held immobile or stationary, and the other one (the mobile phase) is passed over it. Therefore chromatography is mainly divided into two large groups named according to the state of aggregation of the mobile phase, liquid chromatography and gas chromatography. Further groupings can be made by naming both the mobile and stationary phases, thus liquid-liquid, liquid-solid, gas-liquid, and gas-solid chromatography have been named. More recently, supercritical fluids have been used as mobile phases, and these techniques have been named supercritical fluid chromatography, irrespective of the state of the stationary phase.

## Bibliography

1. M. S. Tsvet (1906) *Ber. Deut. Botan. Ges.* **24**, 316.
2. M. S. Tsvet (1906) *Ber. Deut. Botan. Ges.* **24**, 384.
3. *Recommendations on Nomenclature for chromatography, Rules Approved 1973*, IUPAC Analytical Chemistry Division Commission on Analytical Nomenclature (1974) *Pure Appl. Chem.* **37**, 447.
4. M. S. Tsvet (1910) *Khromofilly v Rastitel'nom i Zhivotnom Mire Tipogr. Varshavskago Uchebnago Okruga*, Warsaw.
5. H. H. Strain and J. Sherma (1967) *J. Chem. Educ.* **44**, 238–242.
6. R. Kuhn and E. Lederer (1931) *Ber.* **64**, 1349.
7. R. Kuhn, A. Winterstein, and E. Lederer (1931) *Hoppe Seyler's Z. Physiol. Chem.* **197**, 141–160.
8. A. Tiselius (1940) *Ark. Kemi. Mineral. Geol.* **14B**, 22.
9. S. Claesson (1946) *Ark. Kemi. Mineral. Geol.* **23A**, 1.
10. R. S. Alm, R. J. P. Williams, and A. Tiselius (1952) *Acta Chem. Scand.* **6**, 826–836.
11. A. J. P. Martin and R. L. M. Synge (1941) *Biochem. J.* **35**, 1358–1368.
12. E. Stahl et al. (1956) *Pharmazie* **11**, 633.
13. E. Stahl (ed.) (1962) *Thin Layer Chromatography*, Academic Press, New York.
14. A. T. James and A. J. P. Martin (1952) *Biochem. J.* **50**, 679–690.
15. J. Porath and P. Flodin (1959) *Nature* **183**, 1657–1659.
16. J. N. Wilson (1940) *J. Am. Chem. Soc.* **62**, 1583–1591.
17. L. C. Craig (1950) *Anal. Chem.* **22**, 1346–1352.
18. E. Glueckauf (1955) *Trans. Faraday Soc.* **51**, 34–44.
19. L. Lapidus and N. R. Amundson (1952) *J. Phys. Chem.* **56**, 984–988.
20. J. J. van Deemter, F. J. Zuiderweg, and A. Klinkenberg (1956) *Chem. Eng. Sci.* **5**, 271.
21. J. C. Giddings (1959) *J. Chem. Phys.* **31**, 1462–1467.
22. J. C. Giddings (1963) *Anal. Chem.* **35**, 2215–2216.
23. J. Å. Jönsson (ed.) (1987) *Chromatographic Theory and Basic Principles (Chromatographic Science Series 38)*, Marcel Dekker, New York.
24. L. R. Snyder (1992) in *Chromatography* (E. Heftmann, ed.) (*J. Chromatogr. Library*, Vol. **51A**), Elsevier, Amsterdam, pp. A1–A68.
25. C. Horváth and W. R. Melander (1983) in *Chromatography* (E. Heftmann, ed.) (*J. Chromatogr. Library*, Vol. **22A**), Elsevier, Amsterdam, pp. A27–A135.

26. L. S. Ettre and A. Zlatkis (1979) *75 Years of Chromatography-A Historical Dialogue* (L.S. Ettre and A. Zlatkis, eds.) (J. Chromatogr. Library, Vol. 17), Elsevier, Amsterdam.
27. D. L. Saunders (1975) in *Chromatography*, 3rd ed. (E. Heftmann, ed), Van Nostrand Reinhold, New York, p. 81.
28. A. Braithwaite and F. J. Smith (1985) *Chromatographic Methods*, 4th ed., Chapman and Hall, New York, p. 8.

## Chromocenter

The salivary gland nuclei of *Drosophila melanogaster* contain [polytene chromosomes](#). These specialized [chromosomes](#) essentially contain duplicated repeats of the DNA in the entire chromosome. Polytene chromosomes are fused at their [centromeres](#) to form the chromocenter. This specialized chromosomal domain in *Drosophila* consists of constitutive [heterochromatin](#). The DNA sequences within the chromocenter consist of two types of repeats. One type that is assembled into a-heterochromatin is composed of highly repeated simple DNA sequences, whereas the second, b-type is more complex. Structural components of heterochromatin at the chromocenter have been identified (1). The best-characterized protein that accumulates selectively at the chromocenter is known as heterochromatin protein 1 (HP1). [Antibodies](#) to HP1 colocalize with the type of DNA repeat found in b-heterochromatin called the Dr.D element (2). This is at the edges of the chromocenter in what is termed pericentric heterochromatin.

HP1 contains a protein motif that is found in other chromatin binding proteins and is known as the chromodomain (for chromatin organization modifier). Recent studies have established that the chromodomain family of proteins comprises more than 40 members (3) that can be subdivided into two major groups. Proteins, such as HP1, contain both an amino-terminal chromodomain and a carboxy-terminal “shadow” chromodomain (4). The amino-terminal chromodomain directly binds to heterochromatin, whereas the carboxy-terminal “shadow” chromodomain determines nuclear localization and assists in binding to chromatin (5). The second group of proteins relies on interactions with other proteins to target association with particular chromatin domains (6). The structure of the chromodomain was recently determined using [NMR](#) (7). The chromodomain has strong homology to two archaeobacterial DNA-binding proteins. However, the eukaryotic chromodomain does not interact with DNA and is involved in [protein-protein interactions](#). Each chromodomain consists of an amino-terminal, three-stranded, antiparallel [beta sheet](#) that folds against a carboxy-terminal [alpha-helix](#). The presence of both a chromodomain and a shadow chromodomain are thought to allow proteins, such as HP1, to function as adapters in assembling large multicomponent proteins.

The chromocenter is a useful chromosomal domain for identifying the structural components of heterochromatin and potentially for actually understanding how heterochromatin is organized at a molecular level. The formation of the chromocenter indicates how similar [nucleoprotein](#) complexes that share common structural components can self-associate. It provides a nice example of a specialized nuclear compartment that is assembled so that depends on protein–nucleic acid interactions.

## Bibliography

1. T. C. James and S. C. R. Elgin (1986) *Mol. Cell Biol.* **6**, 3862–3872.
2. G. L. G. Miklos, M.-T. Yamamoto, J. Davies, and V. Pirotta (1986) *Proc. Natl. Acad. Sci. USA*

85, 2051–2055.

3. E. V. Koonin, S. B. Zhou, and J. C. Lucchesi (1995) *Nucl. Acids Res.* **23**, 4229–4233.
4. R. Aasland and A. F. Stewart (1995) *Nucl. Acids Res.* **23**, 3168–3173.
5. J. S. Platero, T. Hartnett, and J. C. Eissenberg (1995) *EMBO J.* **14**, 3977–3986.
6. A. Lorentz, K. Ostermann, O. Fleck, and H. Schmidt (1994) *Gene* **143**, 1–8.
7. L. J. Ball et al. (1997) *EMBO J.* **16**, 2473–2481.

## Chromogenic Substrate

A chromophore is any light-absorbing group (see [Absorption Spectroscopy](#)), and a *chromogenic* substrate is one that is acted on by an [enzyme](#) so as to increase or decrease the absorption of light at a particular wavelength as substrate is converted to product. An example of a naturally occurring chromogenic substrate is **NADH**, which absorbs light strongly at 340 nm. The absorbance decreases as NADH is oxidized by a pyridine nucleotide-dependent dehydrogenase, because NAD does not absorb substantially at 340 nm. Artificial chromogenic substrates have been used extensively for kinetic studies on a range of enzymes. *p*-Nitrophenylphosphate is hydrolyzed by **phosphatases** to release *p*-nitrophenol, which has a yellow color in alkaline solution. The increase in absorbance with hydrolysis is measured at 400 nm. *p*-Nitrophenylesters are used in a similar manner to measure the activity of *esterases*. Phenazine methosulfate is used as a convenient means of determining the activity of flavoprotein enzymes, as the oxidized form of the electron acceptor is yellow and the reduced form is colorless. An extensive list of artificial, chromogenic substrates has been recorded ([1](#)).

## Bibliography

1. R. M. C. Dawson, D. C. Elliott, W. H. Elliott, and K. M. Jones (1986). *Data for Biochemical Research*, Clarendon Press, Oxford, pp. 350–377.

## Chromomere

Chromomeres represent discrete structures in [chromosomes](#) that are visible under the light microscope in mitotic or meiotic **prophase** (see [Chromatid](#)). [Electron microscopy](#) and staining indicate that chromomeres represent regions of more complex, higher order structure than the 30-nm [chromatin](#) fiber ([1](#)). It has been suggested that they represent sites of the preferential chromatin compaction that occurs early in the prophase, coincident with the separation of sister chromatids ([2](#)). However, electron microscopic images of native chromatin in the **interphase** nuclei at physiological ionic strengths reveal that chromomeres may be identical to globular clusters of [nucleosomes](#) [(superbeads)] that probably represent local coiling of the 30-nm fiber. It is also possible to detect bead-like discontinuities in the chromatin fiber in sectioned nuclei. Under certain conditions, these superbeads are remarkably uniform and contain between 8 and 48 nucleosomes. This suggests that chromomeres represent discrete structural units ([3](#)). The failure to find a ubiquitous organization of chromatin into chromomeres or superbeads, however, indicates that all chromatin does not adopt this

form of higher order structure in interphase nuclei *in vivo*. This irregularity may in fact represent the true state of affairs within the nucleus, because, clearly, eukaryotic DNA is not packaged into structures of a crystalline order and stability.

Chromomeres are useful for discriminating between different chromosomes. They are arranged in a specific and consistent pattern along each chromosome. In [lampbrush chromosomes](#), chromomeres occur where there are long regions of inactive chromatin that are consequently compacted into higher order structures.

### Bibliography

1. D. E. Comings (1978) *Ann. Rev. Genet.* **12**, 25–45.
2. G. F. Bahr and P. M. Larsen (1974) *Adv. Cell. Mol. Biol.* **3**, 192–215.
3. H. Zentgraf and W. W. Franke (1984) *J. Cell Biol.* **99**, 272–286.

## Chromosomes

Chromosomes are the [nucleoprotein](#) complexes that provide the structural framework for the expression of **genes** and mediate the transfer of genetic information from generation to generation. The [nuclei](#), [mitochondria](#), and [chloroplasts](#) of all cells, both **eukaryotes** and **prokaryotes**, plus **viruses**, all contain chromosomes. These various chromosomes can vary greatly in size, from very large (see [Polytene Chromosome](#), [Lampbrush Chromosomes](#)) to very small (see [Double Minute Chromosome](#); [Minichromosome](#)). All chromosomes contain **DNA**, which can be packaged by different proteins to assemble a nucleoprotein complex. The DNA in the chromosomes of a eukaryotic cell nucleus is packed with small basic proteins known as [histones](#). The most striking property of every chromosome within the eukaryotic cell nucleus is the length of each molecule of DNA incorporated and folded into it. The human [genome](#) of  $3 \times 10^9$  bp would extend over a meter if unraveled and straightened, yet it is compacted into a nucleus only  $10^{-5}$  m in diameter. It is an astonishing feat of engineering to organize such a long linear DNA molecule within ordered structures that can reversibly fold and unfold within the chromosome.

The basic architectural matrix of the chromosomes found within eukaryotic cell nuclei is [chromatin](#). All of the DNA in the nucleus of a somatic cell is packaged into chromatin by the histone proteins, together with other [DNA-binding proteins](#). In specialized **germ cells**, such as [sperm](#), [protamines](#) can replace the histones as the primary means of packaging DNA. There is considerable specialization in the type of chromatin assembled in different regions of a chromosome, depending on its functional requirements. Chromatin containing genes that are being expressed in a cell is called [euchromatin](#). A large chromosomal region that contains inactive genes is assembled into [facultative heterochromatin](#). Specialized chromosomal structures that contain very few genes but have other essential architectural roles in the chromosome are assembled into constitutive [heterochromatin](#). The [centromere](#) is an essential structure containing constitutive heterochromatin that mediates segregation of the chromosomes. The molecular motors that drive this process are found within the [kinetochore](#) that is attached to the centromeric heterochromatin. Other heterochromatic domains are at the ends of chromosomes in specialized structures known as [telomeres](#). Within the telomere are many reiterated [terminal repeat](#) sequences resulting from the activity of the enzyme [telomerase](#). The telomere protects the end of the chromosome from degradation, fusion, and progressive loss of DNA during chromosomal **duplication**.

The position of the centromere relative to the rest of the chromosome provides an important reference point for describing different types of chromosomes. Chromosomes can be **metacentric**, with a centromere near the middle of the chromosome, or **acrocentric**, when it is near the end of the chromosome, or even potentially **telocentric**, when it is at the very tip. Most eukaryotic chromosomes are **monocentric**, having a single centromere, but some are **holocentric** and have multiple centromeric domains. Holocentric chromosomes have specialized mechanisms for segregating chromosomes during **cell division** that differ from those found in most plant and animal cells.

In animal cells, the presence of multiple centromeres in a single chromosome is often the product of chromosomal rearrangements. Such events usually occur as a result of chromosomal damage and can involve **duplications**, **inversions**, and [translocations](#). A **dicentric** chromosome contains two centromeres following one of these rearrangements, whereas an **acentric** chromosome lacks a centromere entirely. [Isochromosomes](#) might also be formed, where a chromosome contains multiple identical chromosomal arms. All of these rearrangements lead to major problems for segregating chromosomes during the cell division of **meiosis** and **mitosis**. Chromosomal damage can be very useful in mapping the positions of genes relative to each other by methodologies that apply somatic cell genetics (see [Radiation Hybrid](#)).

Each **diploid** cell receives two sex chromosomes, either two [X-chromosomes](#) in females or one X- and one [Y-Chromosome](#) in males, plus two copies of each of the other chromosomes, known as [autosomes](#). Each pair of autosomes consists of two [homologous chromosomes](#), one from the father and one from the mother in humans. Immediately after mitotic division, each chromosome consists of a single DNA molecule. This single molecule of DNA is assembled into a [chromatid](#), which duplicates during the **S-phase** of the [cell cycle](#) to form two sister chromatids. Then these sister chromatids are segregated to the two daughter cells at mitosis. In meiosis, [haploid](#) germ cells are created. This involves segregating of homologous chromosomes at the first meiotic division before the sister chromatids are segregated at the second meiotic division. Because homologous chromosomal regions are aligned during meiosis, chromosomal rearrangements can lead to major problems in segregating homologous chromosomal material. This can lead to [multivalent](#) chromosomes, whereas a normal meiosis would produce only bivalent chromosomes.

For organisms that undergo sexual reproduction, special [sex](#) chromosomes exist. In humans, these are known as the X- and Y-chromosomes. Female cells contain two X-chromosomes, whereas male cells contain a single X-chromosome and a Y-chromosome. Mary Lyon proposed that **gene expression** from the two X-chromosomes in a female cell should be the same as that from the single X-chromosome in a male cell (see [Lyon Hypothesis](#)). This is accomplished by silencing one of the two X-chromosomes in a female cell (see [X-Chromosome Inactivation](#); [Random X-Inactivation](#)). The silenced X-chromosome is converted to a [Barr body](#) made up of facultative heterochromatin. X- and Y-chromosomes share some limited regions of [homology](#) that mediate their pairing and subsequent segregation during meiosis.

An important aspect of chromosomal biology involves mapping individual genes on chromosomes, which facilitates diagnosing particular diseases. A wealth of mapping procedures exist. These make up the field of [cytogenetics](#), which is based on visualizing chromosomes. Mapping procedures involve using particular stains for the DNA and protein components of the chromosome (see [C-Banding](#); [G Banding](#)). In recent times, these methodologies have been supplemented with high-resolution [denaturation mapping](#) techniques. Then all of the stained chromosomes of a particular cell that make up the [karyotype](#) can be displayed formally as an [ideogram](#).

At a more refined level, structure at the level of chromosome and chromatin also determines the functions of particular genes. It has long been known from visualizing the large polytene chromosomes that individual chromosomal **domains** exist. Individual domains like the [Balbiani rings](#) contain chromatin that reversibly changes structure to reflect different functional states. Balbiani rings appear as **puffs** when they are being **transcribed**. Other regions of the polytene

chromosome, known as [interbands](#), contain the regulatory DNA sequences that control the appearance of puffs. Chromosomal rearrangements can lead to positioning genes next to large domains of constitutive heterochromatin. Placing a gene next to the [chromocenter](#) in *Drosophila*, where the centromeres fuse together, **represses** expression of the gene. This [position effect](#) indicates that structurally specialized domains of chromosomes exist and presents a major problem for biotechnology and gene therapy in expressing foreign DNA in target cells.

**Nucleases** like **DNase I** have been very useful in mapping regulatory DNA in the chromosome. The **locus control regions**, [enhancers](#), and **promoters** that control gene expression exist as sites that are **hypersensitive** to digestion by DNase I. These elements often control the expression of clusters of [contiguous genes](#) that have related functions in the cell. Sites with extreme [DNase I sensitivity](#) are also found at the regulatory elements controlling the initiation of [DNA replication](#) at the **origin of replication**. The **terminally redundant** regulatory elements of [retroviruses](#) are also assembled into nucleoprotein complexes that retain accessibility to DNase I.

Chromosomes are fascinating structures around which much of modern medicine and biotechnology revolves. They also provide the foundation for considering basic molecular mechanisms that control gene expression and for considering the forces that facilitate [evolution](#). Genes require a chromosomal environment to be maintained throughout the generations and within which to realize their full regulatory potential.

## Chymotrypsin, Chymotrypsinogen

*Chymotrypsin* (E.C. 3.4.21.1) is a mammalian digestive [serine proteinase](#), of the **trypsin** family, that is synthesized in the acinar cells of the pancreas in the form of an inactive precursor, *chymotrypsinogen*. This [zymogen](#) is stored, along with other digestive [enzymes](#) and enzyme precursors, in pancreatic granules, whose contents are released into the duodenum when the pancreas is stimulated by the [hormones](#) cholecystokinin and acetylcholine, which are secreted in response to eating a meal.

Chymotrypsinogen is a single [polypeptide chain](#) of 245 amino acid residues that becomes enzymatically active on **proteolytic** cleavage of the peptide bond that connects Arg15 and Ile16. This cleavage is catalyzed by another pancreatic enzyme, **trypsin**, which in turn is generated from its zymogen precursor, trypsinogen, by [enterokinase](#), an intestinal serine proteinase. Cleavage of the Arg15—Ile16 bond results in a reorganization of the protein structure, alignment of the [catalytic triad](#), and generation of a substrate binding site (1). Other proteolytic cleavages accompany activation of chymotrypsinogen, but only the one cited is crucial for generating enzymatic activity.

The biological function of chymotrypsin is to digest dietary proteins in the small intestine. It catalyzes the cleavage of [peptide bonds](#) in which the carbonyl group is supplied primarily by aromatic or bulky **hydrophobic** amino acids ([tyrosine](#), [tryptophan](#), [phenylalanine](#), [leucine](#), [isoleucine](#), [methionine](#)). These residues become the C-termini of the product peptides and are subsequently removed by digestion with [carboxypeptidase A](#). The catalytic mechanism of chymotrypsin and other serine proteinases is described elsewhere (see [Serine Proteinase](#)).

Numerous chymotrypsin-like enzymes generally referred to as *chymases* have been identified in various animal tissues, particularly in mast cells, neutrophils, and lymphocytes (2). The biological functions of these enzymes are not well understood, but a chymase isolated from human heart is a major angiotensin II -forming enzyme (3). The *prostate-specific antigen*, which is used in the diagnosis of prostate cancer, is a serine proteinase with limited chymotrypsin-like activity (4).



Being a serine proteinase, chymotrypsin is inhibited by diisopropylfluorophosphate ([DIFP](#)) and also by numerous, naturally occurring **serine proteinase inhibitors** and [serpins](#) (5), as well as many small synthetic peptide analogues or [active site](#) reagents (6).

### Bibliography

1. R. M. Stroud, A. A. Kossiakoff, and J. L. Chambers (1977) *Annu. Rev. Biophys. Bioeng.* **6**, 177–193.
2. J. A. Nadel (1991) *Ann. N. Y. Acad. Sci.* **629**, 319–331.
3. Y. Liao and A. Husain (1995) *Can. J. Pharmacol.* **11**(Suppl. F), 13F–19F.
4. A. M. el-Shirbiny (1994) *Adv. Clin. Chem.* **31**, 99–133.
5. J. Potempa, E. Korzus, and J. Travis (1994) *J. Biol. Chem.* **269**, 15957–15960.
6. J. C. Powers et al. (1989) *J. Cell. Biochem.* **39**, 33–46.

### Suggestions for Further Reading

7. R. M. Stroud (1974) A family of protein-cutting enzymes. *Sci. Am.* **231**, 24–88.
8. T. A. Steitz and R. G. Shulman (1982) Crystallographic and NMR studies of the serine proteases. *Annu. Rev. Biophys. Bioeng.* **11**, 419–444.

## Cilia and Eukaryotic Flagella

Cilia and eukaryotic flagella are long (10  $\mu\text{m}$  to 40  $\mu\text{m}$ ), narrow, membrane-bounded structures that contain a highly ordered, stable [microtubule](#) array, consisting of nine fused outer doublet microtubules surrounding a central pair of singlet microtubules. These microtubules are firmly anchored at the base of the cilium or flagellum in basal bodies. The microtubules of cilia and flagella are extremely stable, and their movement is not complicated by dynamics. In addition to the tubulin backbone, these microtubules contain a large and diverse array of [microtubule-associated proteins](#) (MAPs), most of which have not been characterized, and several forms of [dynein](#), a large multicomponent ATP-transducing mechanochemical [motor protein](#). The various MAPs create and maintain the stable axoneme structure, and the dynein is responsible for creating the movement. The movements of cilia and flagella involve bending of these long organelles, which is caused by the sliding of stable double-microtubule pairs past each other. Sliding is effected by the dynein motors; the base of the motor is attached permanently to the A microtubule of one outer doublet pair, and the head, which contains the motor **domain**, transiently attaches and detaches to the B subfiber of an adjacent outer doublet pair. The dynein motor “walks” along the B subfiber toward the minus end, creating the movement. The dynein-mediated sliding of outer double-microtubules occurs in a fashion similar to filament sliding in muscle contraction, in which [actin](#) filaments and myosin filaments slide past each other through the action of the ATP-transducing myosin motor “heads.”

## Circadian Rhythms and Clocks in Fungi

Circadian clocks are ubiquitous and present in organisms as distantly related as **fungi** and mammals. Classically defined as rhythms that persist under constant environmental conditions, circadian clocks control daily rhythms that correlate with changes in the external environment. Most likely, circadian rhythms were selected early in [evolution](#) for conferring an evolutionary advantage on the organisms containing them.

To compare the properties of circadian rhythms among different organisms, the concept of circadian time (CT) was developed. The circadian day is divided into 24 equal parts, each one circadian hour. By convention, CT0 is subjective dawn, and CT12 represents subjective dusk. The normal circadian cycle in *Neurospora* can be reset in a time-dependent manner by external signals (a property called entrainment) such as light (1) or temperature (2).

Lower eukaryotes are ideal models for clock study because of their genetic and biochemical tractability. *Neurospora crassa* is such an example. Circadian rhythms in *Neurospora* are easily observed with the use of specialized growth tubes called **race tubes**. When wild-type strains of *Neurospora* are inoculated on one side of a race tube, the organism initiates a rapid vegetative growth toward the opposite side. After growth for a day in constant light, the position of the growth front is marked, and the culture is transferred to constant dark. The light-to-dark transfer synchronizes the culture on CT12 (subjective dusk). Monitoring of the growth front every 24 h reveals a typical pattern of alternate conidial and aconidial bands, whose frequency permits the determination of both period length and phase of the rhythm. The *Neurospora* circadian rhythm has a period of 22 h at 24 °C. This period does not change significantly between 18°C and 30°C, as demanded by the circadian definition of a rhythm (a phenomenon called temperature compensation).

Three general questions drive research in fungal chronobiology: What is the biochemical and/or genetic basis for the clock? How does the clock get its input from its extracellular and extraorganismal environment? How is time, as paced by the clock, transduced within the cell?

## 1. The Clock

The central oscillator of the *Neurospora* clock is the *frq* locus. Originally identified by chemical and ultraviolet mutagenesis, mutations at this locus alter the normal circadian clock of the cell but have no severe morphological defects in growth or development (3, 4). Cloning of the *frq* locus elucidated a complex pattern of gene [transcription](#) (3, 5). Two overlapping transcripts of 4 and 4.5 kb were identified, both containing long leader sequences, with several upstream open reading frames. These transcripts have been predicted to encode a putative 989-amino-acid-residue protein. Recent data, however, support the existence of two forms of FRQ (6). The long form (FRQ<sup>1-989</sup>) starts at the first in-frame ATG [start codon](#), whereas the short form (FRQ<sup>100-989</sup>) starts at the third in-frame ATG of the longest reading frame in the region (6). Temperature regulates the ratio of FRQ forms by favoring the use of different initiation codons at different temperatures. For example, the long form of FRQ (FRQ<sup>1-989</sup>) is favored at high temperatures, whereas at lower temperatures the short form of FRQ (FRQ<sup>100-989</sup>) is predominant. These similar but functionally distinct forms of the clock protein FRQ represent a novel adaptive mechanism designed to keep the *Neurospora* clock running over a wide range of temperatures (6). Despite its fine regulation, the biochemical function of FRQ remains unknown (7). However, the presence of a nuclear localization signal [required for FRQ activity (8)], a weak **helix–turn–helix** DNA-binding domain (9), and a conserved acidic region strongly suggest the involvement of FRQ in transcription. As predicted for protein regions important for function, all of these FRQ transcription-factor-like signatures are conserved in FRQ homologues isolated from other distantly related fungal species (10).

The autoregulatory feedback cycle controlled by *frq* takes one day to complete. This loop involves the transcriptional activation and repression of the *frq* locus. Starting with low levels of *frq* transcript and protein (dawn), *frq* transcripts begin to rise, reaching a peak accumulation 4 to 5 h later [just before noon—(8)]. The maximal level of FRQ protein can only be observed 4 to 6 h later (relative to

the peak in *frq* mRNA). This time lag is probably due to post-transcriptional regulation of *frq* mRNA or FRQ protein. Synthesized FRQ protein enters the nucleus (8) and acts rapidly, either directly or indirectly, to repress *frq* (11). The repressing activity of FRQ remains until the protein begins to be **phosphorylated**, linking the activity of unknown protein kinases to the clock. Phosphorylation of FRQ increases its turnover (12), until its scarcity allows positive factors like *white collar-1* (*wc-1*) and *white-collar-2* (*wc-2*) (1, 7) to restart the cycle (see text below).

## 2. Input

Light and temperature are among the most important environmental stimuli controlling circadian rhythmicity. The light-dependent [signal transduction](#) pathway is thought to involve an as yet uncharacterized flavin-dependent blue-light photoreceptor (13, 14). When wild-type strains of *Neurospora* grown in the dark receive a pulse of light, the levels of *frq* mRNA increase. Two global regulatory genes, *wc-1* and *wc-2*, have been established to be involved directly or indirectly in the induction and maintenance of *frq* mRNA levels following this light treatment. For example, the light response observed in wild-type strains (ie, the increase in *frq* mRNA levels) cannot be detected in strains mutant in *wc-1*, but it is normal in strains mutant in *wc-2*. Interestingly, the products of both loci (*wc-1* and *wc-2*) are required for the maintenance of high-level expression of *frq* mRNA in the light. Therefore, the pulse of light positively stimulates *frq* mRNA levels in an FRQ-independent, WC-dependent manner (1, 15). Moreover, the light dosage positively correlates with the amount of *frq* mRNA, and, depending when this pulse of light is perceived, it can reset the clock, thereby phase advancing or phase delaying the conidiation rhythms of the organism. Complete elucidation of the roles of WC-1 and WC-2 on *frq* gene expression will be essential for a complete understanding of the signal transduction pathway responsible for clock input.

## 3. Output

Genes whose rhythmic regulation equals the period of the strain under study are said to be part of the output component of the oscillator (ie, they execute the commands given by the master oscillator and establish the cellular rhythm). To be defined as such, loss-of-function mutations in these classes of genes must not affect the functioning of the clock. With this rationale in mind, a systematic search for clock-controlled genes (*ccg*) was carried out in *Neurospora* (16). Most, but not all, of the *ccg* genes identified were found to be induced during the asexual **sporulation** pathway. The ones that were not were postulated to be involved in other developmental pathways (eg, sexual sporulation) or other unknown output pathways (7). Although rhythmically regulated, the levels and amplitude of expression differ among different clock-controlled genes. In addition, their expression pattern in both *frq*<sup>+</sup>, and *frq* changes from gene to gene [ie, some genes are repressed, while others are derepressed or unchanged (7)]. To date, a direct link between *frq* (the central master oscillator) and *ccg* (the foot soldiers) has not been established, despite being an area of active investigation. It is not surprising to find that the *ccg* identified have nothing or very little in common. For example, *ccg-1* codes for a small polypeptide of 71 amino acid residues with no known homologues in other organisms (17), despite its elevated level of expression (18), whereas *ccg-12*, encodes copper [metallothionein](#) (18). See Refs. 18 and 19 for a complete description of *ccg* genes. While the mechanisms of circadian rhythms are incompletely understood at this time, it remains undeniable that the molecular strategies employed by fungi are an excellent model system with which to study similar systems in higher eukaryotes.

## Bibliography

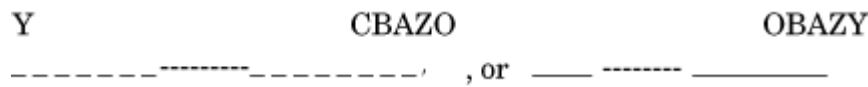
1. S. K. Crosthwaite, J. J. Loros, and J. C. Dunlap (1995) *Cell* **81**, 1003–1012.
2. Y. Liu, M. Merrow, J. J. Loros, and J. C. Dunlap (1998) *Science* **281**, 825–829.
3. J. C. Dunlap (1996) *Annu. Rev. Genet.* **30**, 579–601.
4. J. F. Feldman (1982) *Annu. Rev. Plant Physiol.* **33**, 583–608.
5. C. R. McClung, B. A. Fox, and J. C. Dunlap (1989) *Nature* **339**, 558–562.

6. Y. Liu, Y. Garceau, J. J. Loros, and J. C. Dunlap (1997) *Cell* **89**, 477–486.
7. D. Bell-Pedersen (1998) *Microbiol.* **144**, 1699–1711.
8. C. Luo, J. J. Loros, and J. C. Dunlap (1998) *EMBO J.* **17**, 1228–1235.
9. M. T. Lewis, L. Morgan, and J. F. Feldman (1996) *Mol. Gen. Genet.* **253**, 401–414.
10. M. Merrow and J. C. Dunlap (1994) *EMBO J.* **13**, 2257–2266.
11. M. W. Merrow, N. Y. Garceau, and J. C. Dunlap (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3877–3882.
12. N. Y. Garceau, Y. Liu, J. J. Loros, and J. C. Dunlap (1997) *Cell* **89**, 469–476.
13. E. Klemm and H. Ninnemann (1979) *Photochem. Photobiol. Rev.* **4**, 207–265.
14. J. Paietta and M. L. Sargent (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5573–5577.
15. S. K. Crosthwaite, J. J. Loros, and J. C. Dunlap (1997) *Science* **276**, 763–769.
16. J. J. Loros, S. A. Denome, and J. C. Dunlap (1989) *Science* **243**, 385–388.
17. J. J. Loros (1995) *Semin. Neurosci.* **7**, 3–13.
18. D. Bell-Pedersen, M. Shinohara, J. J. Loros, and J. C. Dunlap (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13096–13101.
19. D. Bell-Pedersen, N. Garceau, and J. J. Loros (1996) *J. Genet.* **75**, 387–401.

## Circular Chromosome

All **Hfr strains** (for high frequency of [recombination](#)) of *Escherichia coli* originate from  $F^+$  strains (males, possess the  $F$  sex factor). When mated with  $F^-$  recipients (females), they generate about a thousand times more recombinants than the equivalent  $F^+$  strain under identical conditions ([1](#)). The recombinants obtained inherit a particular segment of the donor chromosome exclusively ([1](#)). The way in which *Hfr* donors transfer their chromosome to recipients was elucidated by violent agitation of a mixture of mating Hfr and  $F^-$  in a high-speed mixer. The mating was interrupted at different time intervals and the mixture was plated for detection of recombinants. It was found that recombinants acquire different selected characters from the donor cell at different times. The order in which these characters appear is the order of their arrangement on the chromosome. The only reasonable interpretation of these results is that the effect of agitation is a rupturing of the chromosome during the transit from the donor to recipient bacteria. Only those genes that have already penetrated the recipient bacteria at the time of treatment can appear in recombinants. Thus, the strain Hfr used comprises a homogeneous population of donor bacteria, all of which transfer their chromosome in the same specifically oriented and linear way, so that a particular extremity, designated O (for origin), is always the first to penetrate the recipient bacteria ([2](#), [3](#)).

A number of Hfr stains of independent origin were analyzed with respect to the markers they transferred to recipients and to the orientation of the transfer. The outcome was remarkable and unforeseen. Although all of the bacteria of any one Hfr strain transfer a given set of markers in a particular sequence, different strains transfer different chromosomal segments, parts of which may overlap. Some strains transfer their genes in a direction reverse to others. The relationship of the various genes to one another however, is the same irrespective of the Hfr strain used to map them. This gives an unequivocal general linkage map for *E. coli*, consisting of a single linkage group. Each strain transfers its chromosome as a linear, oriented structure with a specific head and tail, for example,



Another strain can always be found which transfers genes A and Z as closely linked markers, for example,



Thus, in spite of this linear transfer, it is impossible to define any extremities on the chromosome. It is concluded that the chromosome of *F*<sup>+</sup> donors is originally continuous or circular and that Hfr strains arise by opening up the circle at a point characteristic for each Hfr type to yield a linear, transferable structure (3).

The circular nature of the *E. coli* chromosome has been confirmed by gentle liberation and dispersion of its DNA and visualization by [autoradiography](#) or by [electron microscopy](#). Therefore there is indisputable, direct evidence that this chromosome is indeed a continuous, circular thread of double-stranded DNA (4, 5).

#### References

1. F. Hayes (1953) Cold Spring Harbor. Symp. Quant. Biol. **18**, 75–93.
2. E. L. Wollman and F. Jacob (1955) C. R. Acad. Sc. Paris **240**, 2449–2451.
3. E. L. Wollman and F. Jacob (1958) Ann. Inst. Pasteur **95**, 641–666.
4. J. Cairns (1962) J. Mol. Biol. **4**, 407–409.
5. J. Cairns (1963) Cold Spring Harbor Symp. Quant. Biol. **28**, 43–46.

#### Suggestion for Further Reading

6. F. Jacob and E. L. Wollman (1961) *Sexuality and the Genetics of Bacteria*, Academic Press, New York, pp. 164–165.

## Circular Dichroism

Circular dichroism (CD) is a form of [spectroscopy](#) that uses circularly polarized light and is highly sensitive to the [conformations](#) of molecules. Because of this conformational sensitivity, CD is one of the most widely used techniques for characterizing the conformations of [proteins](#), **nucleic acids**, and carbohydrates, and for monitoring conformational transitions induced by temperature, solvent composition, **ligand binding**, etc.

The principal conformations of polypeptide chains—[alpha-helix](#), [beta-sheet](#), [b-turns](#), and the so-called unordered conformation—all have distinctive CD spectra that permits these **secondary structures** to be readily recognized in their pure forms or in simple combinations. Even in globular proteins, which represent complex mixtures of these and other conformations, in most cases the CD spectrum can be analyzed to provide reasonably accurate estimates of the fraction of residues in the various secondary structures. CD signals due to aromatic side chains and nonprotein chromophores provide sensitive probes of protein [tertiary structures](#) and of protein–ligand interactions. In the field of nucleic acids, CD also provides valuable information about the nature of the secondary and higher

order structures detecting the binding of small molecules and proteins. CD has more restricted application to carbohydrates and polysaccharides, but it has some important applications to studies of carbohydrate conformation.

In addition to its pronounced sensitivity to conformation, CD has other advantageous features. Only small amounts of material are required, comparable to those required for a UV-visible absorption spectrum, and it is relatively easy to survey a broad range of conditions of temperature, pH, ionic strength, and solvent composition. Moreover, the cost of the instrumentation is moderate, and it is easy to operate and maintain. The principal limitation of CD is that it provides lower resolution structural information than [X-ray crystallography](#) or [NMR](#) spectroscopy. A CD signal can rarely be assigned to a specific residue or group in a macromolecule. Whereas CD gives an estimate of the fraction of residues in a globular protein that are in  $\alpha$ -helical segments, it cannot provide information on their location in the sequence. These limitations, however, are outweighed by the conformational sensitivity and other experimental advantages of CD.

## 1. Fundamentals of Circular Dichroism

### 1.1. Definition of Circular Dichroism

Circular dichroism is a type of spectroscopy in which the difference in absorption of right- and left-circularly polarized light is measured. In circularly polarized light (cpl), the vector  $\mathbf{E}$  describing the electric field of the light, rotates about the direction of propagation once in each period or wavelength of the light. As seen by an observer looking toward the light source, the electric vector may rotate in either a clockwise sense, corresponding to right-circularly polarized light (rcpl), or in a counterclockwise sense, corresponding to left-circularly polarized light (lcpl). At a given instant, the tip of the electric vector of rcpl follows a right-handed helix in space, whereas the tip of the electric vector of lcpl follows a left-handed helix. Therefore, cpl is **chiral**, and rcpl and lcpl are related as mirror images, that is, they are **enantiomeric**. Because of their enantiomeric character, it is to be expected that rcpl and lcpl interact differently with a chiral molecule, giving rise to differences in optical properties for a chiral molecule measured with rcpl and lcpl.

CD is defined as the difference in absorbance of rcpl and lcpl:

$$CD = \Delta A = A_l - A_r = \epsilon_l c \ell - \epsilon_r c \ell = (\epsilon_l - \epsilon_r) c \ell = \Delta \epsilon_M c \ell \quad (1)$$

Here  $A_l$  and  $A_r$  are, respectively, the absorbances for lcpl and rcpl;  $c$  is the molar concentration of the chiral species;  $\ell$  is the pathlength of the sample cell (in cm);  $\epsilon_l$  and  $\epsilon_r$  are the molar extinction coefficients of the chiral species for lcpl and rcpl, respectively; and  $\Delta \epsilon_M$  is the molar CD, commonly called simply the CD.

To facilitate comparisons of the CD of proteins or nucleic acids of very different molecular masses, it is nearly universal practice to normalize the CD by dividing the molar CD by the number of amino acid or nucleotide residues, yielding the mean residue CD,

$$\Delta \epsilon = \Delta \epsilon_M / n_r = \Delta A / c_r \ell \quad (2)$$

where  $n_r$  is the number of residues and  $c_r$  is the molar concentration of residues. Near-UV CD spectra of proteins are reported as either molar CD or mean residue CD, and visible CD spectra are usually reported as molar CD of the protein or the molar CD per visible chromophore, for example, per heme in vertebrate hemoglobins. The reader must always carefully check the basis for reporting CD spectra in the literature.

Although nearly all CD instrumentation measures  $\Delta A$  and hence yields  $\Delta \epsilon$  most directly, many instruments are calibrated in units of ellipticity, an angular unit based on an earlier method for

measuring CD. The molar ellipticity is defined as

$$[\theta]_{\text{M}} = 100\theta/c\ell \quad (3)$$

where  $\theta$  is the measured ellipticity in degrees and  $c$  and  $\ell$  have the same meaning as before. Molar ellipticity is directly proportional to molar CD:

$$[\theta]_{\text{M}} = 3298\Delta\epsilon_{\text{M}} \approx 3300\Delta\epsilon_{\text{M}} \quad (4)$$

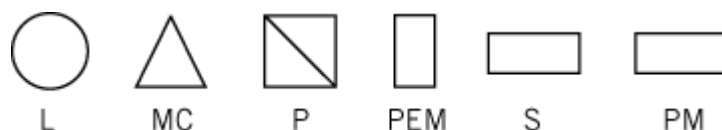
Mean residue ellipticity  $[\theta]$  is defined by analogy to mean residue CD and is commonly used for reporting protein CD data in the far UV and for CD of DNA and RNA.

## 1.2. Instrumentation (1)

Measurement of CD requires high sensitivity because the difference between  $e_l$  and  $e_r$  is small in comparison with their average value  $e = (e_l + e_r)/2$ , which corresponds to the molar absorbance coefficient measured with unpolarized light. Typically,  $\Delta e/e$  is of the order of  $10^{-3}$  to  $10^{-4}$ . To achieve adequate signal-to-noise ratios, CD spectrophotometers use a modulation technique. The light incident on a sample is switched between left- and right-circular polarization at a specific frequency. The difference in the sample response to the two forms of circularly polarized light is detected directly by selecting the part of the response that has a frequency component matching that of the modulation of the incident light.

Figure 1 shows a highly schematic representation of a CD spectrophotometer. Light from the source L is dispersed in the monochromator MC, and then a narrow band of wavelengths passes through the linear polarizer IP. The polarizer splits the unpolarized monochromatic beam into two linearly polarized beams, one polarized in the  $x$  direction, the other in the  $y$  direction ( $z$  is the direction of propagation). Then one of the two linearly polarized beams passes through the photoelastic modulator (PEM), which consists of a plate made of a transparent, optically isotropic material bonded to a piezoelectric quartz crystal. When an alternating electric field is applied, the light emerging from the PEM switches from  $l_{\text{cpl}}$  to  $r_{\text{cpl}}$  and back with the frequency of the applied electric field, typically about 50 kHz. If the sample S exhibits CD, the amount of light absorbed varies periodically with the polarization of the incident light, and so the intensity of light that reaches the photomultiplier PM exhibits sinusoidal intensity variations at the frequency of the field applied to the PEM. Thus, the photomultiplier output consists of a signal with a small alternating current (ac) component superimposed on a direct current (dc) component. The ac component is filtered out and amplified. The ratio of the ac to the dc component is directly proportional to the circular dichroism of the sample, and this quantity is recorded as a function of wavelength to provide a CD spectrum.

**Figure 1.** A block diagram of a CD spectrophotometer. L is the light source, MC the monochromator, P the polarizer, PEM the photoelastic modulator, S the sample, and PM the photomultiplier.



CD is most commonly measured in the visible and ultraviolet regions of the spectrum. Commercially available instruments are used in the 180 to 700-nm region. Some commercial instruments permit measurements to 170 nm, at the edge of the vacuum UV, or to 1000 nm, in the near infrared (IR). Xenon arc lamps, quartz optics, and photomultiplier tubes are used for the light source, the

dispersing and polarizing elements, and the detectors, respectively.

CD is also observed in IR absorption bands, associated with vibrational transitions (see [Vibrational Spectroscopy](#)). This type of CD is called vibrational circular dichroism (VCD) (2). Some VCD instruments use the same type of layout as that shown in Fig. 1, although the light sources, dispersing and polarizing optics, and detectors are necessarily different. In addition, VCD instruments have been constructed that utilize the Fourier transform principle.

In recent years, CD has been widely used to measure the kinetics (3) of conformational changes in proteins and nucleic acids, especially in investigating [protein folding in vitro](#) and unfolding. Reactions on the timescale of a minute or longer are readily measured by using manual mixing and conventional CD instrumentation. A conventional CD instrument coupled to a stopped-flow mixing device permits measurements of reactions with halftimes down to a few milliseconds. Still faster reactions, triggered by temperature or pressure jumps, or by laser or electron pulses, require specialized CD instrumentation (4) that does not utilize the modulation technique. Such instruments have been constructed for nanosecond and picosecond timescales.

### 1.3. Theoretical Background (5)

The connection between theory and experiment in CD spectroscopy is provided by a quantity called the rotational strength,  $R$ . Experimentally, the rotational strength is proportional to the area under a CD band, that is, the integral

$$R \propto \int (\Delta\epsilon/\lambda) d\lambda \approx (\lambda_{\max})^{-1} \int \Delta\epsilon d\lambda \quad (5)$$

The integral is taken over a CD band attributable to a single electronic or vibrational transition. CD bands are generally rather sharp, so it is a satisfactory approximation to replace the  $\lambda^{-1}$  factor in the integrand of Eq. 5 by  $\lambda_{\max}^{-1}$ , where  $\lambda_{\max}$  is the wavelength at which the CD band has maximal intensity.

The rotational strength is related to molecular properties by the equation

$$R = \text{Im}\{\boldsymbol{\mu}_{0a} \cdot \mathbf{m}_{a0}\}, \quad (6)$$

where  $\text{Im}$  indicates the imaginary part of a complex quantity. The first factor in the curly brackets is the electric dipole transition moment for the transition from the ground state 0 to the excited state  $a$ . The second factor is the magnetic dipole transition moment of the transition. Qualitatively,  $\mathbf{m}_{0a}$  can be interpreted as the linear displacement of charge associated with the transition, and  $\mathbf{m}_{a0}$  as the circular displacement of charge. Therefore a non zero rotational strength requires that the transition have both a linear and a circular displacement of charge and that the axis of the circulation not be perpendicular to the linear displacement. This corresponds to a helical motion of electronic charge, where the sense of the helix is determined by the relative orientation of  $\mathbf{m}_{0a}$  and  $\mathbf{m}_{a0}$  and this determines the sign of the CD band.

### 1.4. Related Phenomena

Phenomena that depend upon differences in the interaction of a sample with left- and right-circularly polarized light are called chiroptical phenomena. CD is the most widely used chiroptical phenomenon in molecular biology, but others are important for historical reasons or have specialized applications. The first chiroptical phenomenon discovered was optical rotation (OR), the rotation of the plane of polarization of plane-polarized light as it passes through a chiral medium. The angle through which the plane is rotated, commonly measured at 589 nm (the Na D line), is still widely used in organic chemistry to characterize chiral molecules. OR was also used in protein chemistry, but was supplanted in the 1950s by optical rotatory dispersion (6) (ORD), the wavelength



dependence of OR, and in the 1960s by CD. The specific rotation  $\alpha$  of a sample is defined as

$$[\alpha] = 10\alpha/c'\ell \quad (7)$$

where  $\alpha$  is the angle through which the plane of polarization is rotated by the sample,  $\ell$  is defined as in Eq. (1), and  $c'$  is the concentration of the chiral substance in g/mL. Molar rotation and mean residue rotation are defined by analogy to the corresponding CD parameters. It should be noted that refraction is a scattering phenomenon and occurs at all wavelengths, not just in absorbance bands. For this reason, OR and ORD can be measured at visible wavelengths for substances that absorb only in the ultraviolet.

## 2. Proteins

### 2.1. General Aspects

In the far UV, the peptide groups of the backbone generally dominate the CD spectrum. The peptide group has two electronic transitions in the normally accessible far UV. These are the  $np^*$  transition near 220 nm and the  $pp^*$  transition near 190 nm in secondary amides and 200 nm in tertiary amides (X-Pro peptide groups). The  $np^*$  transition is weak in absorbance, but it gives rise to strong CD bands. The  $pp^*$  transition is associated with strong absorbance and CD. Because of the strong electric dipole transition moment,  $pp^*$  transitions in neighboring peptide groups interact with each other, giving rise to two or more absorbance and CD bands. This phenomenon, called exciton splitting, is most clearly seen in the  $\alpha$ -helix CD spectrum, described later.

The aromatic side chains of **Phe**, **Tyr**, and **Trp** residues have strong absorbance bands in the far UV that contribute to the CD spectrum (7). In most cases, their contribution is small compared to those of the much more numerous peptide groups. For some proteins, however, aromatic CD bands are clearly discernible.

In the near UV, the CD spectrum of proteins is dominated by the aromatic and disulfide transitions. The near-UV CD bands of the aromatic side chains are generally relatively sharp and have a characteristic fine structure due to vibrational effects. In proteins with a small number of aromatic side chains, the near-UV CD bands can frequently be assigned to one of the three types of aromatic side chains and in some cases, through [site-directed mutagenesis](#), to specific residues in the sequence. The CD bands due to the [disulfide bond](#) are generally distinguishable from aromatic CD bands by their much greater width.

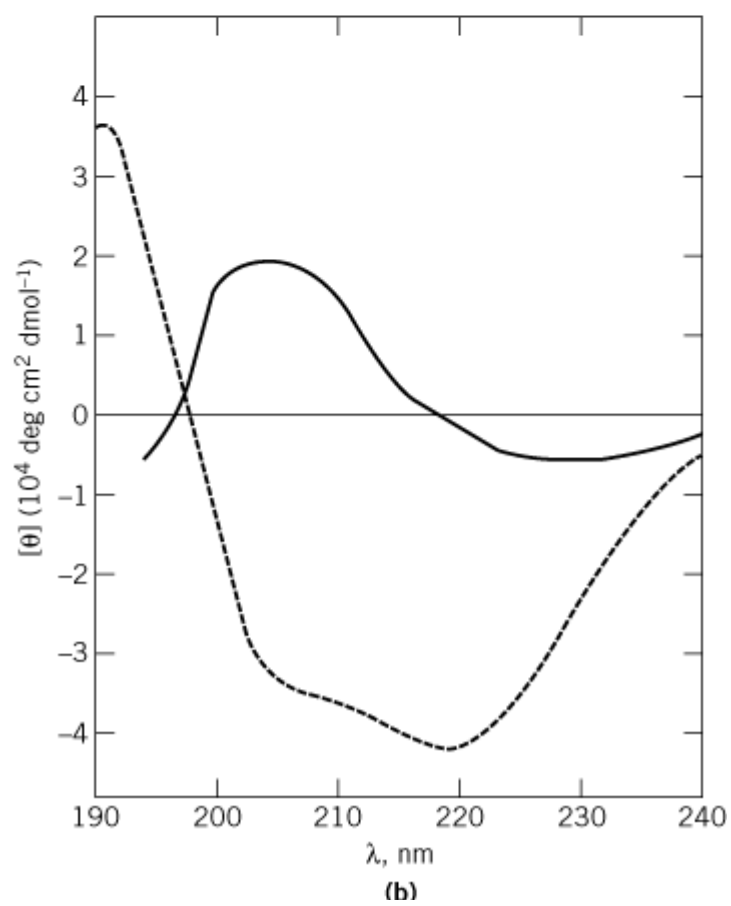
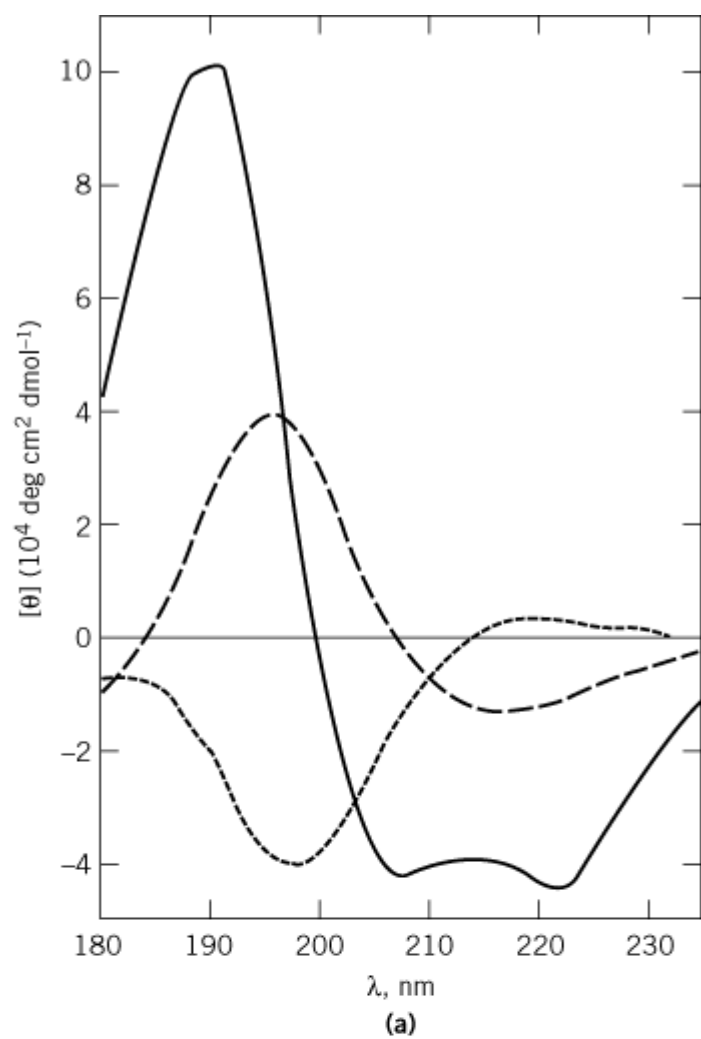
For proteins containing only the normal amino acids, there are no CD bands at wavelengths above 300 nm. Many **prosthetic groups**, **coenzymes**, transition-metal ions, and other ligands have absorbance bands in this wavelength region, and these are associated with CD bands in complexes with proteins.

### 2.2. Secondary Structure Analysis

The various types of secondary structure (see [Secondary Structure, Protein](#)) in proteins have characteristic CD spectra, as established by studies of model oligo- and polypeptides. Figure 2 shows CD spectra of  $\alpha$ -helix,  $\beta$ -sheet, unordered polypeptides, and  $\beta$ -turns. The  $\alpha$ -helix has the most distinctive and strongest CD spectrum with two negative bands of comparable magnitude at ca. 222 and 208 nm and a stronger positive band near 190 nm. The 222 nm band results from the  $np^*$  transition of the amide group, whereas the 208 and 190 nm bands both arise from the amide  $pp^*$  transition. The latter two bands are from exciton splitting of the  $pp^*$  transitions in the amide groups that are held in a well-defined helical geometry. Interactions between the transition dipole moments in a very long helical array give rise to three absorbance bands, one at 208 nm polarized parallel to the helix axis, and two bands at 190 nm, polarized in the two independent directions perpendicular to the helix axis. For a right-handed  $\alpha$ -helix, the parallel band is associated with the negative CD band

at 208 nm, and the perpendicular bands with the positive CD band at 190 nm. The CD of the  $\alpha$ -helix is largely, but not completely, independent of the solvent and of the sequence of amino acids. Aromatic residues (Phe, Tyr, Trp) modify the  $\alpha$ -helix CD spectrum, especially if they constitute a sizeable fraction of the residues. In homopolymers of aromatic amino acids, the  $\alpha$ -helix CD spectrum is modified beyond recognition.

**Figure 2.** (a) CD spectra of model polypeptides in three conformations:  $\alpha$ -helix (—), poly(Glu), pH 4.5 (33);  $\beta$ -sheet (---), poly(Lys-Leu), 0.5 M NaF, pH 7.4 (34); and unordered (·····), poly(Lys-Leu) in salt-free aqueous solution (34). (b) CD spectra of model peptides in two  $\beta$ -turn conformations: Type I  $\beta$ -turn (---), cyclo(L-Ala-L-Ala-Aca), where Aca is  $\epsilon$ -aminocaproyl; Type II  $\beta$ -turn (—), cyclo(L-Ala-D-Ala-Aca). Spectra acquired in water at 22°C (35).



The CD of b-sheets is more variable than that of a-helices. A negative band near 217 nm and a positive band in the 195 to 200-nm region are characteristic of b-sheets (Fig. 2). Theory (8) predicts that the absolute value of the ratio of the ellipticity at the positive maximum near 197 nm to that at the negative maximum near 217 nm increases with increasing twisting of the sheet, and is larger for parallel than for antiparallel twisted sheets. The relative amplitude of the 217 nm and 197 nm bands varies considerably among model systems, ranging from ~1.7 for the antiparallel b-sheet formed in poly(Lys) by heating to 52°C at pH 11.1 to 8.3 for the parallel b-sheet of Boc(Val)<sub>7</sub>OMe in trifluoroethanol.

Several types of b-turn can be distinguished, but only two fundamental types are common in proteins: the type I turn (and its variant, type III), which can accommodate any L-amino acid at either of the positions in the turn, and the type II turn, which usually has a Gly at the second residue of the turn. The CD characteristics of b-turns vary in accordance with the range of conformations. However, studies of cyclic peptides with well-defined b-turns indicate that type I turns have a-helix-like CD spectra, with negative maxima at ~220 and 210 nm and a positive band near 190 nm, whereas type II turns have a CD spectrum like that of a b-sheet, but with the maxima shifted 5 to 10 nm to longer wavelengths, that is, the negative band in the 220 to 225-nm region, and the positive band between 200 and 210 nm (Fig. 2).

All models for unordered [random coil](#) polypeptides have a strong negative band near 200 nm (Fig. 2), but some have a long-wavelength positive band and others a negative shoulder at longer wavelengths. The resemblance of the CD spectra of charged poly(Lys) and poly(Glu) to that of poly(Pro) II led to the proposal that these polypeptides are not truly unordered but contain short segments of the poly(Pro)II helix. The current view (9) of unordered polypeptides is that those peptides with a negative shoulder in the 210- to 220-nm region have amino acid residues predominantly in the a-helix and b-sheet regions of conformational space (see [Ramachandran Plot](#)) and those with a positive long-wavelength band have a substantial fraction of residues in the poly(Pro)II conformation. The latter systems undergo a noncooperative transition to the former as the temperature increases.

### 2.2.1. Protein Secondary Structure

Early applications of OR and ORD to the analysis of helix content in proteins have been reviewed (10). Initially, protein CD data were analyzed by fitting the CD spectra of the protein to a linear combination of data for model polypeptides, such as those shown in Fig. 2. However, such methods generally gave poor results because the CD spectra of globular proteins are too complex to be adequately represented by a simple linear combination of homopolymer spectra. Many methods for analyzing protein CD data to obtain secondary structure have been proposed (11, 12). Two features are essential for satisfactory results. First, a basis set of proteins of known secondary structure is needed to calibrate the method by providing information on the CD contributions of the various types of secondary structure existing in real protein structures. Second, it is necessary to allow flexibility in weighting the proteins in this basis set when analyzing each protein. Several methods incorporating these features provide useful estimates of the fraction of a-helix, b-sheet, b-turn, and unordered conformations (11-14). VCD and IR absorption spectra (See [Vibrational Spectroscopy](#)) in the amide I region also provide estimates of the secondary structure content of globular proteins by using analytical methods similar to those used for electronic CD (15). Combinations of electronic CD, VCD and IR absorption can also be used.

The methods used for analyzing the secondary structure of proteins should be applied with caution to [peptides](#), because globular proteins are used to calibrate these methods. The CD measured at 220 or 222 nm has frequently been used to determine the helix content of peptides. The difference in CD between the 217 nm negative maximum and the positive maximum near 195 nm, characteristic of b-sheets, has been utilized to quantify the b-strand conformation of peptides (16).

### 2.2.2. Tertiary Structure of Proteins

Far-UV CD spectra of proteins have been used to assign proteins to the broad classes of all- $\alpha$ , all- $\beta$ ,  $\alpha + \beta$ , and  $\alpha/\beta$  by cluster analysis (17). Aromatic side-chain bands (7) in both the near and far UV have been used as markers to identify [tertiary structures](#). The **T-state** (deoxy) conformation of [hemoglobins](#) has a strong negative CD band at 287 nm, whereas the **R-state** (liganded) conformation has weak negative or positive CD at this wavelength. The magnitude of the negative CD band at 250 to 255 nm correlates with the T  $\rightarrow$  R transition in [insulin](#). Tertiary structural changes in the **chymotrypsinogen**  $\rightarrow$  chymotrypsin conversion are associated with a reorientation of the Trp 175 - Trp 215 pair, leading to a large change in the 225 nm Trp CD band.

### 2.2.3. Protein Folding

One of the major aims in equilibrium studies of protein folding (see [Protein Folding In Vitro](#)) is detecting intermediates. If equilibrium intermediates are not detectable, the folding/unfolding process is said to be two-state because only the fully folded and fully unfolded forms are present at significant levels. If only these forms are present, the extent of unfolding is the same regardless of the physical property used to monitor it. CD is especially useful in this regard because the far-UV and near-UV CD measure primarily the secondary and tertiary structures, respectively. If the extent of unfolding measured by these two probes is identical within experimental error, the process is likely to be two-state. On the other hand, lack of coincidence of the unfolding curves measured in these two spectral regions is evidence that one or more intermediates exist.

A type of folding/unfolding intermediate observed in a number of proteins is called the [molten globule](#). This intermediate is nearly as compact as the native protein and has a secondary structure content comparable to that of the native protein. However, the side chains are quite mobile and are not locked into a well-defined tertiary structure. CD is one of the most useful techniques for identifying the molten globule form of proteins. The CD spectrum of a molten globule is similar to that of the native protein in the far UV and reflects a substantial amount of  $\alpha$ -helix and/or  $\beta$ -sheet but has only weak features in the near UV because the aromatic side chains do not have well-defined positions and conformations.

Another aspect of equilibrium protein folding studies in which CD has played a major role is investigating the effect of sequence on the stability of intact proteins and peptide segments. Here, the far UV is used to follow the **denaturation** profiles of the wild-type protein and of a series of mutant forms in which specific residues have been altered. From reversible thermal denaturation studies, changes in the free energy of unfolding,  $DDG_{unf}$  can be obtained that measure the effect of the amino acid substitution upon the net stability of the protein and thus shed light on the effects of steric bulk, **hydrophobic** character, and charge effects. There have also been numerous CD studies of model peptides that adopt the helix conformation to examine systematically the effects of amino acid substitutions on  $\alpha$ -helix stability (18). Similar studies have recently been reported for  $\beta$ -sheet-forming peptides and protein fragments (19).

Kinetic studies of protein folding by CD have been limited to stopped-flow methods and thus to the millisecond timescale. During the so-called dead time of the experiment, many proteins studied thus far (20) acquire a substantial fraction of  $\alpha$ -helix and  $\beta$ -sheet, as evidenced by  $[q]_{220}$  values that are generally closer to those of the native form than to the fully unfolded form. By contrast, the CD in the near-UV detected immediately after the dead time is essentially that of the unfolded form. These results generally have been interpreted as resulting from a so-called burst phase in which a molten globule-like intermediate is formed on the timescale of milliseconds or less. Then subsequent slower processes lead to additional formation or remodeling of the secondary structure and formation of the tertiary structure.

### 2.2.4. Membrane Proteins

The study of [membrane proteins](#) in their native environment poses problems for CD. Two kinds of potential artifacts must be avoided. Strongly scattering suspensions of membrane fragments exhibit

preferential scattering of  $rcpl$  or  $lcpl$ , which manifests itself as an apparent CD signal that distorts the true CD of the membrane protein. The particulate character of membrane fragments also distort the CD signals because of the CD analog of Duysens' flattening (21), which occurs in samples with a highly inhomogeneous distribution of absorbing chromophores that is associated with clustering of the chromophores. Methods have been devised that correct for or avoid each of these difficulties, but in the flattening effect, they require that the membrane protein be transferred from the native membrane to small unilamellar vesicles that have, on average, no more than one protein molecule per vesicle. Alternatively, these artifacts are circumvented by solubilizing the membrane protein in a nonionic detergent. This poses some risk of inducing conformational changes in the protein, although nonionic detergents are usually relatively benign.

### 2.2.5. Protein–Ligand Interactions

CD is widely used to characterize ligand binding to proteins and yield binding constants, stoichiometry, and information about conformational changes in the ligand and/or protein, in favorable cases. If the ligand absorbs at wavelengths above 300 nm, it is convenient to monitor the corresponding CD bands. This is especially so if, as in many cases, the ligand is achiral or a rapidly interconverting racemate. In such cases, the free ligand contributes no background CD, and only the bound ligand contributes to the CD signal. This is effectively so even in chiral ligands, such as NADH, FMN, FAD, in which the chromophore is conformationally mobile in the free form and therefore has only a weak CD signal. Ligand binding is also studied at shorter wavelengths, where protein CD bands are normally studied, either as induced ligand CD bands or as a perturbation of the aromatic and peptide CD bands. These perturbations arise either through direct interactions with the ligand or through changes in the tertiary or secondary structure induced upon ligand binding.

The use of CD to determine **dissociation constants** ( $K_d$ ) is straightforward, but there are limitations to the range of dissociation constants that can be determined accurately. To provide an accurate analysis, the concentration of the species being monitored must be comparable to the  $K_d$ . Given the typical magnitudes of CD, this requires  $K_d > ca. 1 \mu M$  for successful measurements by conventional CD. Fluorescence-detected CD (22) may extend this range to ca. 10 nM.

If  $K_d < ca. 1 \mu M$ , CD cannot be used to measure reversible dissociation, but the protein–ligand complex dissociates negligibly at  $\mu M$  concentrations and above. Studies of such complexes are very fruitful because the ligand serves as a tightly but noncovalently bound spectroscopic probe, capable of reporting conformational changes in the protein. Heme proteins (23) and flavoproteins are important examples. The CD signals from such proteins are highly sensitive to oxidation state, spin state (in heme proteins), and to substrate, inhibitor, and allosteric effector binding.

## 3. CD of Nucleic Acids

### 3.1. General Aspects (24)

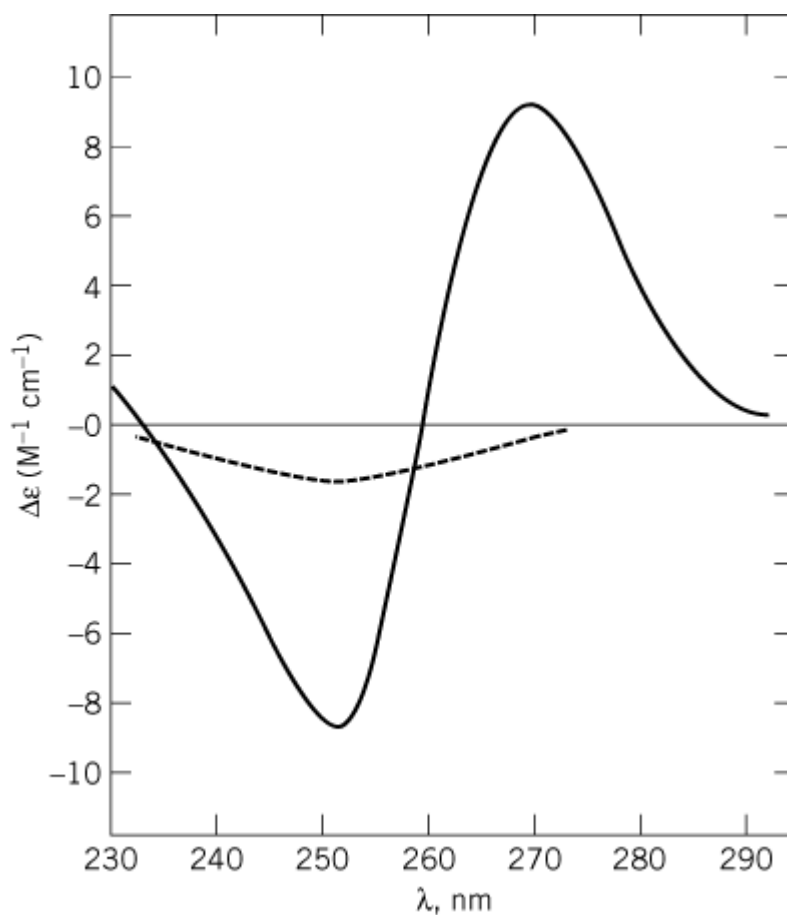
The **purine** and **pyrimidine** bases of **DNA** and **RNA** are largely responsible for the CD spectra of nucleic acids in the wavelength range normally studied (180 to 300 nm). The sugar and phosphate groups do not absorb significantly above 200 and 180 nm, respectively. From the standpoint of CD, their main function is maintaining the relative geometries of the bases. Each base has a characteristic set of  $pp^*$  transitions in the 180- to 300-nm region. The corresponding absorbance and CD bands are relatively broad. All five natural bases have one or two moderately intense  $pp^*$  bands near 260 nm and several more intense bands in the 180- to 200-nm region. In addition, each base is expected to have several  $np^*$  transitions in the 180- to 300-nm region, but these bands absorb weakly. Though potentially strong in CD, few  $np^*$  bands have been identified, and the CD spectra of nucleosides, nucleotides, and polynucleotides are dominated by  $pp^*$  contributions.

The CD spectra of nucleosides and nucleotides are relatively weak compared to those of polynucleotides. In the monomer, only base–sugar and base–phosphate interactions contribute, and

the CD is averaged over a broad range of conformations. By contrast, oligo- and polynucleotide CD spectra are dominated by base–base interactions, and the range of conformations is usually narrower.

The CD spectra of oligonucleotides differ markedly from the sum of the constituent monomer spectra. This is illustrated by the spectra of ApA and pA (Fig. 3). The large difference is attributed to coupling of the  $pp^*$  transitions in the two bases, another example of exciton coupling. This led to a model in which the adenine rings are stacked in a right-handed helix, like two successive bases in A-RNA or A-DNA. In fact, the CD of ApA closely resembles that of poly(A), except that it is only about half as large. This model is supported by theoretical calculations of the CD and more recently by NMR. Interestingly, although the CD spectra of dA, dAp, and dpA are similar to one another and to those of the riboA monomers, the CD spectrum of d(ApA) is quite different from that of ApA and reflect a difference in the geometry of base stacking.

**Figure 3.** Near-UV CD spectra of the dinucleoside monophosphate ApA (—) and the mononucleotide pA (----). Spectra obtained in 0.1 M Tris, 0.1 M NaCl, pH 7.4 at 5.5°C for ApA and at 20°C for pA (36).



The CD of nucleotide trimers can be described in most cases by the equation

$$[\theta]_{A_p B_p C} = (2[\theta]_{A_p B} + 2[\theta]_{B_p C} - [\theta]_{B_p})/3 \quad (8)$$

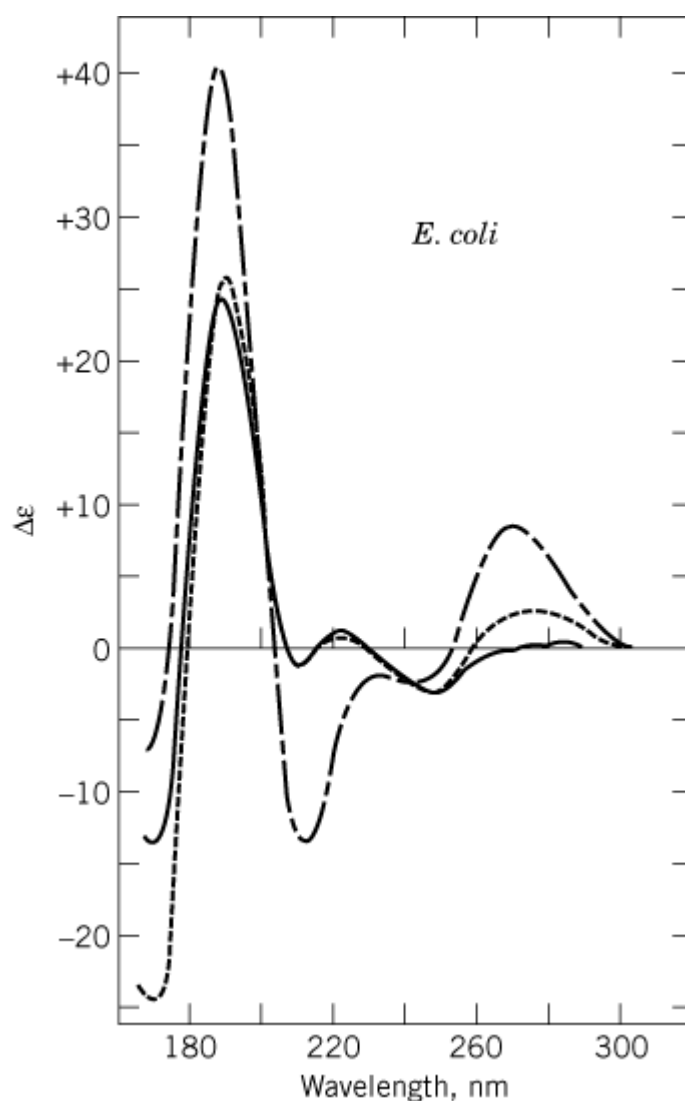
This model assumes that (1) nearest neighbor bases stack similarly in dimers and in higher oligomers and (2) only nearest neighbor interactions are significant for CD. The model works reasonably well for higher oligomers and for trimers and has been extended to double-stranded oligonucleotides of DNA and RNA. The same model has been used to derive nearest neighbor frequencies for DNA and

RNA from base-composition and CD data.

### 3.1.1. Secondary Structure

The CD of a DNA is highly diagnostic of its secondary structure. The [B-DNA](#) conformation normally found in aqueous solution has a positive CD band near 275 nm and a negative band of similar magnitude near 245 nm, as shown in Fig. 4. A pair of closely spaced bands of opposite sign is called a couplet, and it is characterized by the sign of its long-wavelength component as a positive or negative couplet. A couplet in which the two lobes are of very similar amplitude is called a conservative couplet. In the far UV, the B-DNA spectrum has a positive couplet with peaks near 190 and 175 nm that is nearly an order of magnitude more intense than the near-UV bands.

**Figure 4.** The CD spectra of *E. coli* DNA in the B form at low ionic strength (10 mM sodium phosphate buffer, - - -), the A form (0.67 mM sodium phosphate buffer, 80% trifluoroethanol, - · -), and the B form at high ionic strengths (6 M  $\text{NH}_4\text{F}$ , —) (37). (Reprinted with permission from Ref. 37, © 1985, Wiley-Liss.)



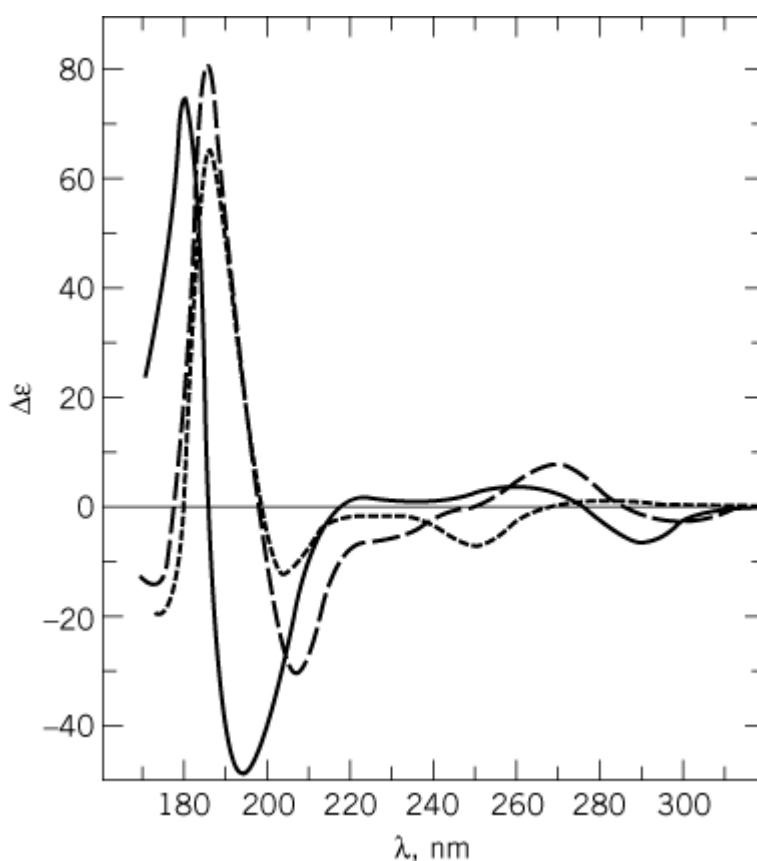
The A-form of DNA (see [A-DNA](#)) is favored by low water activity and is induced by adding alcohols. In the near UV, the B  $\rightarrow$  A transition is marked by a large increase in the positive band and a decrease in the magnitude of the 245 nm band. Therefore, the A-DNA CD spectrum is distinctly nonconservative. A strong negative band near 210 nm is also a characteristic feature of the A-DNA



spectrum in contrast to the weak negative band in the B-DNA spectrum. Below 200 nm, a positive couplet is observed as in B-DNA, but the couplet is asymmetrical.

The B- and A-forms of DNA are both right-handed double helices. As expected, the left-handed Z form of DNA (see [Z-DNA](#)) has a very different CD spectrum, as shown for Z-form poly[d(G-C)] in Fig. 5. The near-UV bands for the Z-form are opposite in sign to those for the B-form and are shifted to longer wavelengths. The strong bands in the far UV are also reversed in sign and are red-shifted relative to those for the B-form. (Note that the CD spectrum of poly[d(G-C)] in the B-form differs from that of natural or “random-sequence” DNA. This is generally true for DNAs with simple di- and trinucleotide repeating sequences, in which specific features of geometry and electronic interactions are not averaged.) The reversed signs of the near-UV CD in the Z-form of poly[d(G-C)] were the first indications of a dramatically different conformation for this copolymer at high salt concentrations, but subsequent studies have shown that the signs of the near-UV CD bands do not reliably indicate helix sense for the double helix. However, the signs of the far-UV CD bands correlate with the helix sense in all known cases: a positive band between 180 and 192 nm for right-handed A- and B-forms and a negative band between 185 and 200 nm for the left-handed Z-form ([25](#)). Near-UV CD is widely used to monitor B → Z transitions in DNA.

**Figure 5.** CD spectra of poly[d(GC)] in three conformations. B form (10 mM sodium phosphate buffer, .....); A form (0.67 M sodium phosphate buffer, 80% trifluoroethanol, - - -); Z form (10 mM sodium phosphate buffer, 2 M NaClO<sub>4</sub>, —) ([38](#)). (Reprinted with permission from Ref. [38](#), © 1985, IRL Press.)



tive to ionic strength. At ionic strengths comparable to physiological, the conservative CD pattern shown in Fig. 4 is observed. As ionic strength increases, the 275 nm band is selectively diminished, and at 6 M NH<sub>4</sub>F, it nearly vanishes (Fig. 4). By contrast, the 245 nm band undergoes little change. Studies with closed circular DNA show that changes in ionic strength from ca. 0.1 M to 6 M change the number of base pairs per turn from ca. 10.4 to 10.2. The DNA in [chromatin](#) has

a CD spectrum similar to that in 6 M  $\text{NH}_4\text{F}$ , indicating that it is a somewhat underwound form of B-DNA.

Double-stranded RNA is nearly always found in the A-form, and never in the B-form. Z-form RNA has been observed for poly[r(G-C)], but the conditions required to drive it into this form are more stringent than for the deoxy copolymer. A-form RNA has a CD spectrum similar to that of A-DNA, except that a weak negative band is commonly observed near 290 nm on the long-wavelength side of the 260 nm band. As with DNAs, the far-UV CD is diagnostic of the helix sense in double-stranded RNA, and the correlation between the sign of the far-UV couplet and the helix sense is the same.

In both RNA and DNA, the single-stranded forms have CD spectra that qualitatively resemble the double-stranded forms because base stacking has a much greater influence on the CD than does base pairing. Strand separation eliminates the latter but does not affect stacking significantly. Thermal denaturation of DNA leads to relatively small CD changes in the near UV but substantially larger far-UV CD changes. CD is also useful in characterizing **triplex** (26) forms of DNA and RNA.

### 3.1.2. Tertiary Structure

**Supercoiling** of DNA leads to readily detectable changes in the CD spectrum (27). The positive 275-nm band of B-DNA becomes more positive for negative superhelical supercoiling and conversely for decreased negative superhelical densities. For superhelical densities in the normal range, the CD changes are directly proportional to the superhelical density.

Although there is no evidence that the CD of RNA is directly sensitive to tertiary structure, it is sensitive to sequence and base pairing. Thus, CD has been used to test alternative models of RNA tertiary structure that differ in secondary structure distribution (28).

### 3.1.3. Ligand Binding to Nucleic Acids

Complexes of DNA and RNA with dyes, antibiotics, and other small ligands (See [Intercalation](#)) give rise to induced CD in the ligand absorption bands and to changes in the nucleic acid CD bands (29). These changes permit the determination of dissociation constants and stoichiometries for complexes of such small ligands with DNA and RNA. In some cases, additional structural information about the ligand or nucleic acid conformation is obtained.

CD is widely used to study DNA-protein and RNA-protein interactions (30). A very useful feature of CD in such studies is that the nucleic acid generally dominates the CD spectrum in the 260-nm region, whereas the protein dominates in the 200- to 240-nm region. Because each of these regions responds to the conformation of the dominant component, CD is of great value in monitoring conformational changes upon complex formation. Nonspecific protein–nucleic acid complexes often have dissociation constants in the micromolar range where CD is well suited for measuring  $K_d$  values. Specific protein–nucleic acid complexes usually have  $K_d$  s that are smaller by several orders of magnitude at physiological ionic strengths. In some cases, the  $K_d$  is shifted into the measurable range by increasing the ionic strength. Even if the binding is too tight to permit measuring the  $K_d$  by CD, conformational changes in the protein and/or nucleic acid are detected and characterized by CD. CD is also well suited for measuring the kinetics of protein–nucleic acid complex formation.

## 4. Carbohydrates

Circular dichroism has important applications in studying carbohydrates (31), although these are more limited than for proteins and nucleic acids. Of the chromophores common in carbohydrates, only the amide (N-acetyl sugars) and [carboxyl groups](#) (uronic acids) have CD bands above

200 nm. The hydroxyl, ether, acetal, and ketal chromophores that are most typical have their first longest wavelength CD bands near the short-wavelength limit of conventional CD instruments, near 180 nm. Higher energy transitions are studied only with vacuum UV instruments, but they are usually obscured by solvent absorption, so studies are limited to thin solid films. For these reasons, optical rotation and ORD measurements in the near UV and visible regions continue to play an important role in studying carbohydrates.

Monomeric sugars have been investigated extensively, and sector rules relating the signs of CD bands near 170 nm to the conformation of the sugar have been derived. A useful model relating the optical rotation at the Na D line to the conformation of a sugar has been developed. This model is useful in testing features of the conformational energy surface in disaccharides. In the study of polysaccharides, CD and optical rotation have been used to characterize order-disorder, single- to multistranded, and sol to gel transitions (32).

## Bibliography

1. W. C. Johnson Jr. (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 635–652.
2. T. A. Keiderling (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 555–598.
3. K. Kuwajima (1996), in *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 159–182.
4. J. W. Lewis, R. A. Goldbeck, D. S. Kliger, X. Xie, R. C. Dunn, and J. D. Simon (1992) *J. Phys. Chem.* **96**, 5243–5254.
5. R. W. Woody (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 25–67.
6. K. Imahori and N. A. Nicola (1973) In *Physical Principles and Techniques of Protein Chemistry* (S. J. Leach, ed.), Academic Press, New York, Part C, pp. 357–444.
7. R. W. Woody and A. K. Dunker (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 109–157.
8. M. C. Manning, M. Illangasekare, and R. W. Woody (1988) *Biophys. Chem.* **31**, 77–86.
9. R. W. Woody (1992) *Adv. Biophys. Chem.* **2**, 37–79.
10. P. Urnes and P. Doty (1961) *Adv. Protein Chem.* **16**, 401–544.
11. W. C. Johnson Jr. (1990) *Proteins: Struct. Funct. Genet.* **7**, 205–214.
12. N. Sreerama and R. W. Woody (1994) *J. Mol. Biol.* **242**, 497–507.
13. S. Yu. Venyaminov and J. T. Yang (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 69–107.
14. N. J. Greenfield (1996) *Anal. Biochem.* **235**, 1–10.
15. V. Baumruk, P. Pancoska, and T. A. Keiderling (1996) *J. Mol. Biol.* **259**, 744–791.
16. L. Zhong and W. C. Johnson Jr. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 4462–4465.
17. S. Yu. Venyaminov and K. S. Vassilenko (1994) *Anal. Biochem.* **222**, 176–184.
18. J. M. Scholz and R. L. Baldwin (1992) *Ann. Rev. Biophys. Biomol. Struct.* **21**, 95–118.
19. D. L. Minor Jr. and P. S. Kim (1994) *Nature* **367**, 660–663.
20. H. Roder and G. A. Elöve (1994) In *Mechanisms of Protein Folding* (R. H. Pain, ed.), IRL Press, Oxford, pp. 26–54.
21. L. M. N. Duysens (1956) *Biochim. Biophys. Acta* **19**, 1–12.
22. D. H. Turner (1978) *Meth. Enzymol.* **49**, 199–214.
23. Y. P. Myer and A. J. Pande (1978) In *The Porphyrins* (D. Dolphin, ed.), Academic Press, New York, Vol. **3**, pp. 271–322.

24. W. C. Johnson Jr. (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 433–468.
25. J. H. Riazance, W. C. Johnson Jr., L. P. McIntosh, and T. M. Jovin (1985) *Nucleic Acids Res.* **15**, 7627–7636.
26. D. M. Gray, S.-H. Hung, and K. H. Johnson (1995) *Meth. Enzymol.* **246**, 19–34.
27. S. Brahms, S. Nakasu, A. Kikuchi, and J. G. Brahms (1989) *Eur. J. Biochem.* **184**, 297–303.
28. K. H. Johnson and D. M. Gray (1991) *Biopolymers* **31**, 385–395.
29. C. Zimmer and G. Luck (1992) *Adv. DNA Sequence-Specific Agents* **1**, 51–88.
30. D. M. Gray (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 469–500.
31. E. S. Stevens (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 501–530.
32. E. R. Morris, D. A. Rees, D. Thom, and E.J. Welsh (1977) *J. Supramol. Struct.* **6**, 259–274.
33. W. C. Johnson Jr. and I. Tinoco Jr. (1972) *J. Am. Chem. Soc.* **94**, 4389–4390.
34. S. Brahms and J. Brahms (1980) *J. Mol. Biol.* **138**, 149–178.
35. J. Bandekar, D. J. Evans, S. Krimm, S. J. Leach, S. Lee, J. R. McQuie, E. Minasian, G. Nemethy, M. S. Pottle, H. A. Scheraga, E. R. Stimson, and R. W. Woody (1982) *Int. J. Peptide Protein Res.* **19**, 187–205.
36. K. E. van Holde, J. Brahms, and A. M. Michelson (1965) *J. Mol. Biol.* **12**, 726–739.
37. W. C. Johnson Jr. (1985) *Meth. Biochem. Anal.* **31**, 62–125.
38. J. H. Riazance, W. A. Baase, W. C. Johnson Jr., K. Hall, P. Cruz, and I. Tinoco Jr. (1985) *Nucleic Acids Res.* **13**, 4983–4989.

### Suggestions for Further Reading

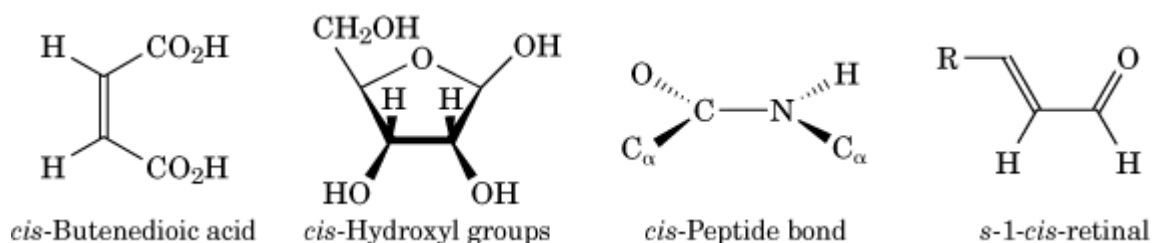
39. G.D. Fasman (1996) *Circular Dichroism and the Conformational Analysis of Biomolecules*, Plenum Press, New York.
40. W.C. Johnson Jr. (1985) Circular dichroism and its empirical application to biopolymers, *Meth. Biochem. Anal.* **31**, 62–125.
41. K. Nakanishi, N. Berova, and R.W. Woody (1994) *Circular Dichroism: Principles and Applications*, VCH, New York.
42. D.W. Sears and S. Beychok (1973) "Circular dichroism", In *Physical Principles and Techniques of Protein Chemistry* (S. J. Leach, ed.), Academic Press, New York, Part C, pp. 445–593.
43. R.W. Woody (1995) Circular dichroism, *Meth. Enzymol.* **246**, 34–71.

### Cis Configuration

The prefix *cis* has been used in chemistry to indicate “on the same side” (1). The original use came from describing [stereoisomers](#) of molecules with a ring or containing carbon–carbon double bonds. When similar substituents are on the same side of the ring or double bond, the [configuration](#) is referred to as *cis*, and when they are on opposite sides, they are referred to as *trans*. The C2 and C3 hydroxyl groups of the furanose form of **ribose** are *cis*. The [conformation](#) about a single bond may

also be noted as *cis* or *trans*, particularly when all the substituents lie in a plane. The *cis* conformation of a [peptide bond](#) has both  $\alpha$ -carbons on the same side of the  $N-C$  amide bond (see [Cis/Trans Isomerization](#)). The notation *s-cis* (2) is used to emphasize that the conformation about a single bond is being described, as in *s-1-cis* retinal. These different uses of the chemical prefix *cis* are shown in Figure 1.

**Figure 1.** The various uses of the prefix *cis* are shown. The common feature is that the prefix denotes “on the same side.”



## Bibliography

1. A. Baeyer (1888) *Liebig's Annalen der Chemie* **CCXLV**, 137.
2. R. Mullikan (1942) *Rev. Mod. Phys.* **14**, 265–267.

## Suggestions for Further Reading

3. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York, pp. 318–371.
4. J. March (1985) *Advanced Organic Chemistry*, 3rd ed., Wiley-Interscience, New York, pp. 109–115.

## *Cis*-Acting

This term is used in the context of regulation of gene [transcription](#) to denote DNA sequences that have to be physically joined to the **genes** in question to influence their activity. [Trans-acting](#) sequences, on the other hand, exert their regulatory effects regardless of whether or not they are present on the same chromosome as the regulated gene. *Trans*-acting genes generally encode proteins that bind to *cis*-acting DNA sequences.

*Cis*-acting sequences were classically identified as a result of their [mutation](#), either silencing or activating the controlled gene, depending on whether the *cis*-acting sequence provided a binding site for **RNA polymerase** or other necessary **transcriptional factor** or for a [repressor](#) protein. These two kinds of *cis*-acting sequences, respectively positive and negative in their effect, were first identified in the *Escherichia coli* [Lac operon](#), where they were called the **promoter** and **operator** sequences. With the development of [recombinant DNA](#) technology, it became possible to detect protein-binding sites in DNA without mutation (e.g., by [footprinting](#) or [gel retardation assays](#)), and then to determine their effects on transcription by “engineered” deletions.

*Cis*-acting sequences that act otherwise than by providing binding sites for regulatory proteins are comparatively uncommon, but two quite different examples should be mentioned. In bacterial [operons](#), a mutation that causes premature translational termination of an [upstream](#) gene can cause a drastic reduction in the expression of genes in the same [operon](#) further [downstream](#). This is a strictly *cis*-acting effect, acting within a single unit of transcription. The mechanism is explained under [Lac Operon](#).

The second example is from female mammals where one of the two [X-chromosomes](#) in each cell lineage is largely inactivated with the notable exception of one gene, active only in the otherwise “silent” X, that is transcribed into a long RNA molecule called Xist. This transcript, which does not encode a protein, plays an essential role in silencing the rest of the chromosome. It acts only in *cis* and remains associated with the chromosome from which it is transcribed without effect on the other, active, X-chromosome (see [X-Chromosome Inactivation](#)).

X-chromosome inactivation involves propagating a certain type of condensed and transcriptionally inactivating chromatin structure, generally called [heterochromatin](#), which, in most chromosomes, is restricted to certain segments. In *Drosophila melanogaster*, heterochromatin exerts a *cis*-effect in suppressing gene activity when, as a result of segmental chromosomal rearrangement, it is brought close to genes not normally associated with it. The multiprotein heterochromatin complex spreads along the chromosome, inactivating the genes in its path with a probability that decreases with distance (see [Position Effect](#)). Genes that are normally close to heterochromatin are apparently insulated in some way against this effect.

## ***Cis*-Dominance**

**Dominance**, in the genetic sense, is a property of one form (**allele**) of a **gene** relative to another allele of the same gene, said to be recessive. Where dominant and **recessive alleles** are present together, the dominant, not the recessive, registers its effect on the organism. Most mutant alleles have lost part or all of their activity and then are usually recessive. When mutant alleles are dominant it is usually because they are hyperactive or have been freed from some normal constraint. The term *cis*-dominant means that the mutation is in an element that has to be physically joined to the gene to exert its effect (ie, in *cis*), not separated from it on another chromosome.

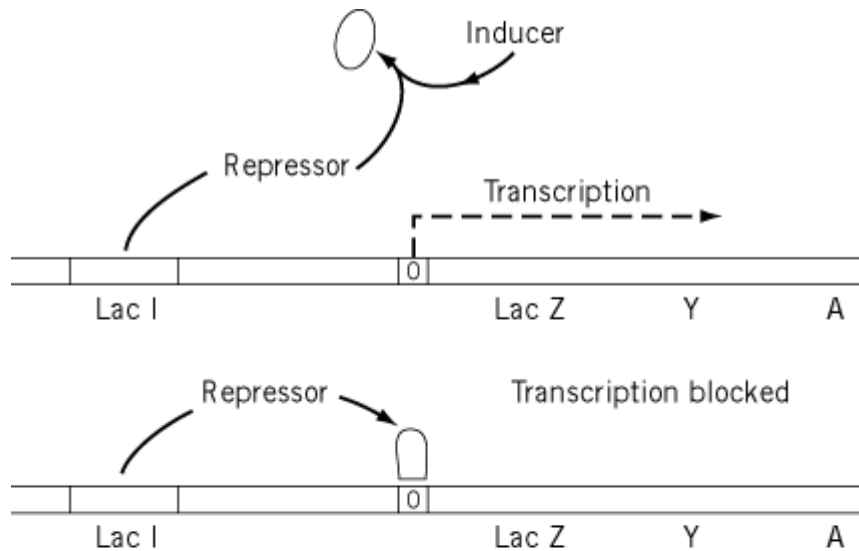
Insofar as a gene is, by definition, an integrated (*cis*-acting) unit of genetic function, any part of a dominant allele, including its coding sequence, could be said to be *cis*-dominant. But, in practice, the term (if it is used at all) is usually reserved for DNA sequences that govern gene activity in *cis* from outside the transcribed or translated sequence.

The concept of *cis*-dominance arose from the classical work of Francois Jacob and Jacques Monod at the Pasteur Institute in Paris on the utilization of the sugar lactose by the bacterium *E. coli* (see [Lac Operon](#)). Lactose utilization by the bacterium depends on two proteins, the enzyme **b-galactosidase**, which cleaves lactose to glucose and galactose, and a [membrane protein](#), b-galactoside permease, needed for uptake of lactose into the cell. These proteins are normally synthesized by *E. coli* only in the presence of lactose in the growth medium, or a lactose analogue, that acts as an **inducer**.

Two classes of mutations result in failure to grow on lactose as a carbon source: *lacZ* mutants lack b-galactosidase and *lacY* mutants lack the permease. *LacZ* and *lacY* are two genes that are closely linked (together with a third, seemingly inessential gene *LacA*) in what was later shown to be a single unit of [transcription](#), or [operon](#), transcribed in the sequence Z - Y - A.

Two other kinds of mutants differ from wild type in producing b-galactosidase and permease **constitutively**, that is, whether inducer is present or not. They were mapped respectively in a separate gene, *LacI*, closely linked to *LacZ*, and in a so-called **operator** segment (*o*) at the end of *LacZ* that overlaps the *LacZ* transcription start point. The *LacI* constitutive mutants (*LacI*<sup>-</sup>) are recessive to wild type (*LacI*<sup>+</sup>), which is interpreted as meaning that *LacI*<sup>+</sup> encodes a transcriptional [repressor](#) and that the repressor function is nullified by inducer. The wild-type operator segment (*o*<sup>+</sup>) was given the hypothetical role of binding to the *LacI* repressor, blocking *lac* operon transcription. The operator-constitutive mutants (*o*<sup>c</sup>) are interpreted as changes in the DNA sequence of the operator so that it no longer bound the repressor. The inducer was thought to act by binding to the repressor protein and changing its conformation, so that it is no longer bound to the operator (see Fig. 1). All of these elements of the model have since been confirmed by molecular analysis.

**Figure 1.** The regulation of the *Escherichia coli lac* operon and two kinds of mutations that have dominant effects: *LacI*<sup>S</sup> (noninducible), dominant both in *cis* and in *trans*; and *o*<sup>c</sup> (constitutive), dominant only in *cis*. The partial diploids for testing dominance (items 5–8) were made by introducing F' plasmids that carry a second copy of the *lac* segment of the bacterial genome. *o*<sup>c</sup> is resistant to the wild type (*LacI*<sup>+</sup>) repressor and also to the *LacI*<sup>S</sup> superrepressor. *LacI*<sup>-</sup> mutants, not shown here, are constitutive and recessive. Information from Refs. 1 and 2.



|   |   | Lac Z (β-galactosidase) activity |              |
|---|---|----------------------------------|--------------|
|   |   | No inducer                       | With inducer |
| 1 | [I <sup>+</sup> O <sup>+</sup> Z <sup>+</sup> ] | -                                | +            |
| 2 | [I <sup>+</sup> O <sup>c</sup> Z <sup>+</sup> ] | +                                | +            |
| 3 | [I <sup>S</sup> O <sup>+</sup> Z <sup>+</sup> ] | -                                | -            |
| 4 | [I <sup>+</sup> O <sup>+</sup> Z <sup>-</sup> ] | -                                | -            |
| 5 | [ $\frac{I^+ O^+ Z^-}{F' I^+ O^c Z^+}$ ]        | +                                | +            |
| 6 | [ $\frac{I^+ O^+ Z^+}{F' I^+ O^c Z^-}$ ]        | -                                | +            |
| 7 | [ $\frac{I^S O^+ Z^+}{F' I^+ O^+ Z^-}$ ]        | -                                | -            |
| 8 | [ $\frac{I^S O^+ Z^-}{F' I^+ O^c Z^+}$ ]        | +                                | +            |

The tests for dominance were made in partial diploids ([merodiploids](#)) that have a second copy of the *Lac* operon brought into the cell as part of an [plasmid](#). Some of the key results are shown in the table of Fig. 1. The formal demonstration of *cis*-dominance comes from the comparison of items 5 and 6. The *trans* merodiploid,  $o^+ LacZ^- / o^c LacZ^+$ , is constitutive, but the *cis* merodiploid,  $o^+ LacZ^+ / o^c LacZ^-$ , is inducible.

The *cis*-dominance of  $o^c$  is readily understandable. There is no way in which the binding of a protein to one piece of DNA could affect the transcription of an unconnected piece. The operator segment is an integral part of the *lac* operon and cannot function when separated from it. On the other hand, *LacI* is a separate and functionally autonomous gene that acts on *LacZ* in *trans* just as well as in *cis*.

Mutations of the rather uncommon *LacI<sup>s</sup>* type change the repressor into a superrepressor that is insensitive to inducer and confers noninducibility. Its dominance over *LacI<sup>+</sup>* (both in *cis* and in *trans* - see item 7 in Fig. 1) is due to the *LacI* repressor functioning as a **dimer**. In mixed *LacI<sup>+</sup>/LacI<sup>s</sup>* protein dimers, the noninducibility of the *LacI<sup>s</sup>* component is imposed on the whole.

The concept of *cis*-dominance was important at a time when the idea that gene sequences had regulatory as well as structurally determining functions was just emerging. Monod and Jacob found that *LacI* acts in *trans* on *LacZ* (and *LacY*) activity and the *o* segment acts only in *cis*. Today it is recognized that the activities of most genes are regulated by multiple *trans*-acting proteins, encoded by other genes, that bind to DNA sequences (**promoters**, **enhancers**, **silencers**) that act in *cis*. But it is questionable whether *cis*-dominance is an appropriate term in reference to all these *cis*-acting transcriptionally controlling sites. Dominance implies that there are recessive alleles over which the dominance can be demonstrated, which is seldom the case, and the *cis-trans* comparison, which could provide the evidence that the effect applies only in *cis*, is hardly ever available.

One dominant mutation, in *Drosophila*, for which the *cis-trans* comparison has been made in *Contrabithorax* (*Cbx*), which maps close to the *Ultrabithorax* (*Ubx*) gene and extends its activity in the [embryo](#) in an anterior direction. Here the effect of *Cbx* is strong in flies of constitution *Cbx + / + Ubx* (*Ubx* being a virtually null allele) but weak, though still significant, in *Cbx Ubx / ++*. In the latter genotype, the *Cbx* regulator is separated from the gene that it regulates. *Cbx* can be called *cis*-dominant, but the observation that it also has a small effect in *trans* is interesting. E. B. Lewis (3) showed that this and other *trans* effects to do with *Ubx* depend on the close pairing of homologous chromosomes that is a special feature of several *Drosophila* tissues.

The concept of *cis*-dominance is an important part of molecular biological history, even though the current usefulness of the term is rather limited. The principle of *cis*-acting controls of gene activity, to which the concept led, is now a commonplace and essential part of molecular biology.

### Bibliography

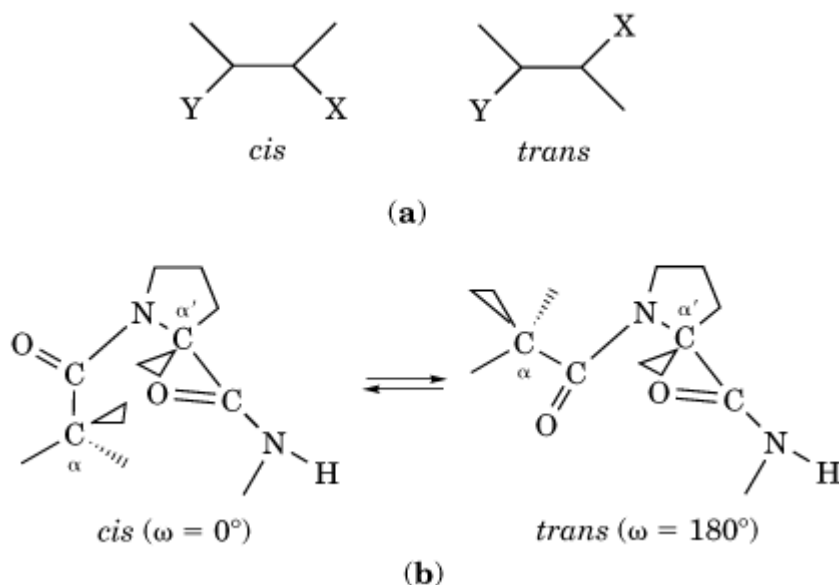
1. J. R. Beckwith (1970) In *The Lactose Operon* (J. R. Beckwith and D. Zipser, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 1–26.
2. F. Jacob and J. Monod (1961) On the regulation of gene activity. Cold Spring Harbor Symp. Quant. Biol. **26**, 193–209.
3. E. B. Lewis (1955) Am. Nat. **89**, 73–89.



## Cis/Trans Isomerization

The distinction between the *cis* and *trans* isomers of a molecule originates from a geometry-based classification of molecular structures. When the two substituents X and Y are on the same side of the structural unit [Fig. 1(a)] the isomer is *cis*. It is designated *trans* in the opposite arrangement. Ambiguities in the designation of the isomers are avoided by using the more sophisticated *E/Z* nomenclature, which often has *E* for the *cis* form and *Z* for the *trans* form. The interconversion between the isomers, the *cis/trans* isomerization, occurs by rotation about the central linkage [Fig. 1(a)]. There is a barrier to rotation quantified by the **free energy** of activation  $\Delta G^\ddagger$  that is proportional to a first-order rate constant  $k_{\text{obs}}$ , where  $k_{\text{obs}} = (k_{\text{cis to trans}} + k_{\text{trans to cis}})$  for the reversible isomerization. The bond order of the central linkage correlates with the magnitude of the energy barrier  $\Delta G^\ddagger$ . It ranges from high values of  $\Delta G^\ddagger > 100$  kJ/mol for double bonds, via an intermediate range for linkages having a partial double bond character, down to low values in the range of  $k_B T$  (where  $k_B$  is the Boltzmann constant and  $T$  the absolute temperature) for C–C single bonds. Under ambient conditions, a rotational barrier  $>90$  kJ/mol indicates that the individual isomers are not readily interconverted, and they may exhibit quite distinct chemical properties. A corresponding  $\Delta G^\ddagger$  value has been used to discriminate semantically between different [conformations](#) (barrier to rotation  $<90$  kJ/mole) and [configurations](#) of a molecule.

**Figure 1.** (a) General representation of *cis/trans* isomerism (b) *Cis/trans* isomerism of a prolyl peptide bond.



In a polypeptide chain, both the amide -CONH- and the **imide** —CON < [peptide](#) groups are intrinsically competent to undergo *cis/trans* isomerization by rotation about the torsion angle  $w$  of the peptide bond (see [Ramachandran Plot](#)). Among the gene-encoded amino acids, the imide peptide bond is exclusively formed by [proline](#) residues. Because the lone electron pair of the nitrogen atom is delocalized over the peptide bond, the C–N linkage has a partial double-bond character. Typically, this destabilizes twisted conformations but stabilizes planar arrangements (*trans*,  $w = 180^\circ$  *cis*,  $w = 0^\circ$ ) of the two  $\alpha$ -C atoms adjacent to the peptide bond [Fig. 1(b)]. The *cis/trans* isomerization of the -CONH- moiety is relatively fast (half-time  $<1$  at room temperature) and leads to a very small percentage of *cis* isomer,  $\gg 1\%$ , at equilibrium. In contrast, the peptidyl-proline moiety (in the

following referred to as a prolyl bond) often has *cis/trans* isomers in comparable amounts. Steric constraints, which favor a *trans* arrangement in the case of -CONH-, are similar in both isomers for prolyl bonds, and the *cis* isomer generally occurs in about 20% of the prolyl peptide bonds. There is a relatively high rotational barrier  $\Delta G^\ddagger$  of about 80 kJ/mol for the prolyl bond. The combination of a substantial population of both isomers and their slow interconversion implies that *cis* prolyl bonds have considerable influence on biochemical reactions of the polypeptide backbone.

In the absence of a folded conformation, polypeptide chains can theoretically form  $2^n$  *cis/trans* isomers, where  $n$  is the number of prolyl bonds in the molecule. In simple cases, as with a four-proline octapeptide derived from the prolactin **receptor**, all of the possible isomers have been detected and quantified in solution (1). Structural formation in the peptide chain, however, reduces the number of isomers substantially. A particular prolyl bond in native proteins is usually either *cis* or *trans* in all the molecules, although *cis/trans* isomerization in the native structure has been observed for a few proteins by NMR methods in solution. In the latter case, structural alterations propagate through the backbone around the isomeric proline residue, which can be accompanied by distinct biological activities of the isomeric proteins (2, 3). Depending on the polypeptide structure, the **half-time** for prolyl bond isomerization ranges from seconds to hours.

From stereochemical considerations, the peptidyl transferase center on the **ribosome** is thought to be constructed for synthesizing all peptide bonds in the *trans* conformation. However, in native, globular proteins of known three-dimensional structure, about 5 to 6% of prolyl peptide bonds are *cis* (4, 5). When proline-containing proteins are unfolded in the presence of high concentrations of **denaturants**, such as **urea** or **guanidinium chloride** (GdmCl), the **random coil** polypeptides generally equilibrate slowly to a mixture of numerous *cis/trans* isomers. For the fraction of unfolded molecules that have one or more nonnative isomers of a prolyl bond, subsequent **refolding** of the protein has to start from different conformational states. If the *cis/trans* isomerization is slower than refolding, slow kinetic phases of folding may be apparent when the time course of refolding is monitored. For that reason, *cis/trans* isomerization is the rate-limiting step in folding for some proteins. It is generally most easily detectable when there is a *cis* prolyl bond in the native state, because then a large fraction of unfolded molecules has the incorrect isomer (6, 7).

Chemical catalysis of prolyl bond isomerization is rare. Most organic solvents and micelles or phospholipid **vesicles**, cause a moderate decrease in the rotational barrier. The rate constants are independent of the pH value in the physiological range, unless dissociable groups are located adjacent to proline. An increased rotational barrier results from the O-protonation of the peptide bond in acidic solution, characterized by a  $\text{pK}_a = -1$ . In strong acids, N-protonated species become populated ( $\text{pK}_a = -7$ ) that accelerate *cis/trans* isomerization (8). Many relationships between the peptide structure and prolyl bond isomerization have been elucidated. For example, when measured in oligopeptides, the increased barrier to rotation caused by aromatic amino acids that precede the proline residue is accompanied by an increase in the *cis* population of up to 40%, whereas small aliphatic side chains in the same position lead to lower  $\Delta G^\ddagger$  values in conjunction and lower *cis* contents of 5 to 10% (9).

The conformational constraints on the polypeptide backbone induced by the prolyl *cis/trans* isomerization restrict bimolecular recognition processes. Conformational selectivity was demonstrated by the **hydrogen bond** directed preferential binding of proline peptides in the *cis* conformation to **b-cyclodextrins** (10) and to a synthetic, multidentate terephthaloyl amide (11). For the recognition of opioid peptides of the dermorphin type by m- and  $\delta$ -receptors, specificity for the *cis* conformers was suggested (12). In a biological context, it may be important that endoproteases, such as **chymotrypsin**, **trypsin**, **thrombin** and clostripain cannot readily cleave a peptide bond adjacent to a *cis* prolyl moiety, even if the isomeric bond occupies a position remote from the scissile bond (13, 14). Due to this conformational specificity, the rate of the *cis* to *trans* isomerization of the prolyl bond limits the rate of proteolysis for good substrates in the presence of high protease

concentrations, and this is a useful assay for isomerization. Even *in vivo*, the time course of bradykinin (Arg-Pro-Pro-Gly-Phe-Ser-Pro-Phe-Arg) degradation by pulmonary endothelial [peptidases](#) is controlled by the conformational specificity of the proteinases (15).

### Bibliography

1. K. D. Oneal et al. (1996) *Biochem. J.* **315**, 833–844.
2. A. P. Hinck, E. S. Eberhardt, and J. L. Markley (1993) *Biochemistry* **32**, 11810–11818.
3. J. Kordel, S. Forsen, T. Drakenberg, and W. J. Chazin (1990) *Biochemistry* **29**, 4400–4409.
4. M. W. MacArthur and J. M. Thornton (1991) *J. Mol. Biol.* **218**, 397–412.
5. D. E. Stewart, A. Sarkar, and J. E. Wampler (1990) *J. Mol. Biol.* **214**, 253–260.
6. J. F. Brandts, H. R. Halvorson, and M. Brennan (1975) *Biochemistry* **14**, 4953–4963.
7. F. X. Schmid (1986) *Methods Enzymol.* **131**, 70–82.
8. H. Sigel and B. R. Martin (1982) *Chem. Rev.* **82**, 385–426.
9. R. K. Harrison and R. L. Stein (1992) *J. Amer. Chem. Soc.* **114**, 3464–3471.
10. M. Lin et al. (1995) *Anal. Chim. Acta* **307**, 449–457.
11. C. Vicent, S. C. Hirst, F. Garciatellado, and A. D. Hamilton (1991) *J. Amer. Chem. Soc.* **113**, 5466–5467.
12. R. Schmidt et al. (1995) *Int. J. Peptide Protein Res.* **46**, 47–55.
13. G. Fischer, H. Bang, E. Berger, and A. Schellenberger (1984) *Biochim. Biophys. Acta* **791**, 87–97.
14. S. Meyer, M. Drewello, and G. Fischer (1996) *Biol. Chem.* **377**, 489–495.
15. M. P. Merker and C. A. Dawson (1995) *Biochem. Pharmacol.* **50**, 2085–2091.

### Suggestions for Further Reading

16. T.E. Creighton, ed. (1992) *Protein Folding*, W.H., New York; A book written by specialists that summarizes proline-limited protein folding in the context of other folding events.
17. B. Testa (1982) "The geometry of molecules: basic principles and nomenclature", In *Stereochemistry* (C. Tamm, ed.) Elsevier Biomedical Press, Amsterdam, New York, Oxford, pp. 1–47.

## Clamp Loaders, Processivity Complex

Chromosomal **replicases**, the **DNA polymerases** that replicate [chromosomes](#), are multiprotein complexes characterized by the high processivity of their DNA synthesis. These replicating machines polymerize thousands of nucleotides without dissociating from the [template](#). The replicases of the three well-characterized systems (*Escherichia coli*, **eukaryotes**, and **bacteriophage T4**), which span the evolutionary spectrum, are similar in function, structure, and overall organization. Their remarkable processivity is achieved by a ring-shaped processivity factor (or "**sliding clamp**") that binds the polymerase catalytic unit and tethers it to the DNA. Another complex of accessory proteins is required to load the clamp onto the DNA. The sliding clamp does not have any affinity for DNA and, therefore, other proteins are needed to assemble the clamp around the DNA. The accessory protein complex ("clamp loader"; also called a "molecular matchmaker") binds to the primer terminus and **couple**s ATP hydrolysis to the assembly of the ring around the DNA primer. Therefore, the replicase can be thought of as having three components: (1) a

polymerase, (2) a clamp loader complex that assembles the clamp around DNA and (3) a DNA sliding clamp (Table 1).

**Table 1. Three-Component Structures of Chromosomal Replicases**

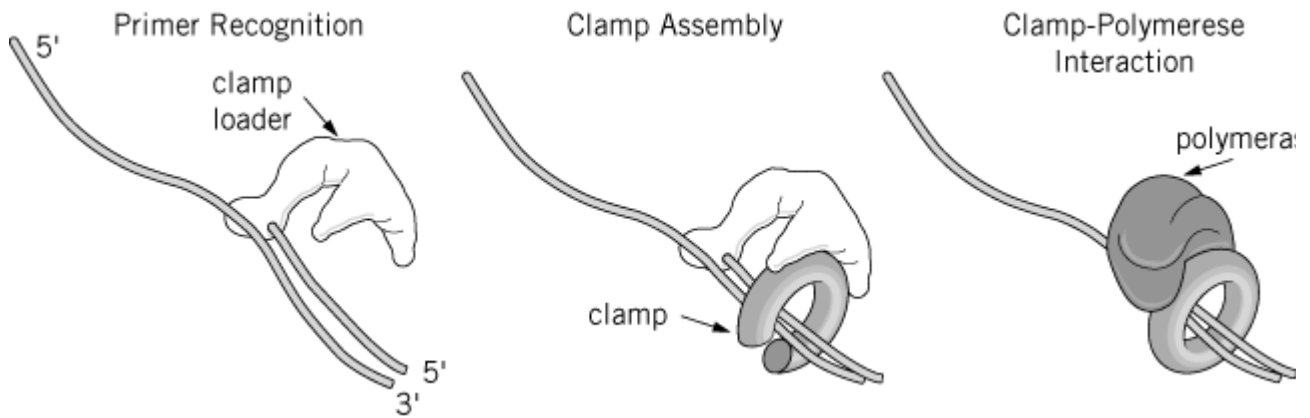
|                | DNA Polymerase               | Clamp loader              | Sliding Clamp |
|----------------|------------------------------|---------------------------|---------------|
| Eukaryotes     | pold (3 subunits)            | RF-C complex (5 subunits) | PCNA          |
| <i>E. coli</i> | Core polymerase (3 subunits) | g-complex (5 subunits)    | b-subunit     |
| Phage-T4       | gp43                         | gp44/62                   | gp45          |

## 1. The Clamp Loaders

The three best-studied replicases are the *E. coli* DNA polymerase III holoenzyme (polIII), the eukaryotic polymerase d (pold), and the replicase of phage T4. These three multisubunit complexes share structural and functional similarities (reviewed in Refs. 1-3). The clamp loaders for these systems are the g-complex of prokaryotes, the replication factor-C (RF-C) complex (also called Activator-1) of eukaryotes, and a complex of the products of gene 44 and gene 62 (gp44/62) of phage T4. The g-complex is composed of five subunits called g, d, d', c and y. RF-C is also a five-subunit complex, composed of one large and four small subunits (p140, p40, p38, p37, and p36). gp44/62 contains only two polypeptides gp44 and gp62. In the complex, a tetramer of gp44 is tightly associated with one gp62. Interestingly, the three clamp loaders share amino acid similarities among the subunits (4, 5).

A detailed mechanism by which these clamp loaders operate in assembling the sliding clamp around DNA is not yet known. The basic steps of the loading process, however, have been elucidated (Fig. 1). In general, the clamp loader recognizes the 3' end of the junction between the single strand and duplex DNA (primer/template) and uses ATP hydrolysis to assemble the clamp around the primer. The status of the clamp loader after assembling the clamp around DNA is not yet clear. It has not yet been determined whether the clamp loader leaves the clamp on the DNA, where it interacts with the polymerase to initiate processive DNA synthesis, or if the clamp loader is also needed to assist binding the polymerase to the clamp (the different features of each clamp loader are discussed later).

**Figure 1.** A model for clamp loader activity. The clamp loader recognizes the 3' terminus of a primer and assembles the sliding clamp around the DNA. Following assembly, the clamp interacts with the polymerase and tethers it to the DNA for processive DNA synthesis.



In prokaryotes and eukaryotes, the clamp loader also functions as a clamp unloader (6). Upon completion of an [Okazaki fragment](#), the polymerase rapidly dissociates from the clamp (7-9), leaving the clamp on the DNA (10). The clamps of polIII (the b subunit) and pold [proliferating cell nuclear antigen (PCNA)] are relatively stable on DNA (6, and references therein). Therefore, it was postulated that the sliding clamp remains assembled around the DNA when the synthesis of an Okazaki fragment is complete. During lagging strand synthesis, a new sliding clamp is needed for the synthesis of each Okazaki fragment. In *E. coli* about 10 times more Okazaki fragments are formed during replication than there are b-subunits within the cell. In human cells, the number of Okazaki fragments formed during one round of chromosomal replication has been estimated to be 100 times greater than the amount of PCNA. Therefore, recycling of the clamp (b and PCNA) is needed to fulfill the requirement for a constant supply of clamps for further DNA synthesis. In phage T4, the clamp readily dissociates from DNA upon completion of an Okazaki fragment. Thus, the clamp loader has a dual function during DNA replication: (1) loading of the sliding clamp onto DNA to initiate processive DNA synthesis and (2) recycling the clamps on the lagging strand.

### 1.1. g-Complex

The g-complex is composed of five subunits (g, d, d', c, y) (reviewed in Ref. 11). The genes encoding them have been identified, and the purified proteins have been used to study the function of the individual subunits and subassemblies of the g-complex (discussed later). Four of the subunits (d, d', c and y) are encoded by unique genes. The g-subunit is formed from the same gene (*dnaX*) that encodes t (another subunit of polIII) by an efficient translational **frameshift** mechanism that produces g in an amount equal to t (reviewed in Ref. 12). As a result, the g-subunit consists of the N-terminal 430 residues of t followed by a unique C-terminal Glu residue. In addition, the g- and d'-subunits show amino acid sequence similarity.

No one subunit alone can assemble the b sliding clamp around DNA (11 and references therein). At low ionic strength, a combination of g and d assemble b onto DNA, but the reaction is feeble. The d', c- and y-subunits are needed for an efficient loading reaction under physiological conditions (13).

The g-complex has only weak affinity for single-stranded and duplex DNA, and it exhibits weak DNA-dependent ATPase activity, but the ATPase activity is stimulated by the sliding clamp (b-subunit) (14). The best DNA effector for the ATPase activity is a singly primed template, indicating an interaction between the g-complex and the primer/template junction. The g-subunit is the only subunit of the g-complex with an [ATP-binding motif](#), and it binds ATP (15). Furthermore, mutation of the ATP-binding site of g within the g-complex destroys the ATPase activity and the ability to assemble the b-ring onto DNA. The g-subunit, however, lacks significant ATP hydrolysis activity even in the presence of DNA (14, 15), and the d- and d'-subunits are required for ATPase activity, implying that the gdd' complex recognizes DNA. The g-complex binds b in the presence of ATP. The d subunit interacts with b, with a strength similar to that of the entire g-complex (16), indicating

that b binds the g-complex mainly through the d-subunit. Interestingly, however, the interaction between the b-subunit and d does not depend on ATP (whereas the g-complex requires ATP to bind the clamp). These apparently contradictory observations can be explained if the d subunit is buried within the g-complex and ATP induces conformational changes that lead to the interaction of d with b (16). The g-complex also interacts with [single-stranded DNA binding protein](#) (SSB) via the c subunit (11). These interactions are important for efficient clamp assembly at physiological ionic strength.

The results obtained with the individual subunits and their subassemblies demonstrate the basic features of the g-complex action during loading of the clamp onto DNA (Fig. 1). Upon binding of ATP, the g-complex undergoes conformational changes leading to interaction of the d-subunit with b. The g-complex recognizes a primed template, aided by the interaction between c and SSB. Assembly of the ring around the DNA primer stimulates the ATPase activity of the g-complex. ATP hydrolysis results in conformational changes such that d is again buried within the g-complex leaving the clamp around the DNA primer. Now the clamp interacts with the polymerase to initiate processive DNA synthesis (Fig. 1).

### 1.2. Replication Factor-C (RF-C)

RF-C was first isolated from human 293 cells as an essential replication factor for the *in vitro* replication of simian virus-40 (SV-40) (17). The factor was purified later from different sources, including yeast and human cell lines (18 and references therein), based on its ability to stimulate processive DNA synthesis by pold in the presence of PCNA and ATP.

Purified RF-C is a five-subunit complex (p140, p40, p38, p37, and p36). The genes that encode all RF-C subunits from yeast and humans have been identified (4, 5 and references therein), and all are essential in yeast. Although encoded by different genes, they show extensive amino acid sequence similarities in the central region of the protein (reviewed in Ref. 5).

Early biochemical studies determined the role of RF-C as a clamp loader of PCNA and suggested a mode of action during the assembly reaction (Fig. 1). RF-C binds preferentially to single-strand/duplex DNA (template/primer) junction in the presence of ATP (19, 20). PCNA interacts with the RF-C/DNA complex. The DNA-dependent ATPase activity of RF-C is stimulated by the binding of PCNA. The hydrolysis of ATP leads to conformational changes within the RF-C, which locks the PCNA sliding clamp around the DNA. Then the clamp interacts with the polymerase to initiate processive DNA synthesis (Fig. 1).

Biochemical analysis of RF-C individual subunits and the protein–protein interactions within the RF-C complex is underway (eg, 21, 22). Purified proteins, subcomplexes and mutational analysis have shed light on the function of the individual subunits and subassemblies of the RF-C complex. Although each subunit contains an ATP-binding motif, only p40 binds ATP (20). Yeast p36 (23) and a subcomplex of human RF-C (p36-p37-p40) exhibit DNA-dependent ATPase activity. The DNA-binding activity resides within the large subunit (p140) (20), although the p37 subunit also binds DNA weakly. All of the subunits, except p37, bind PCNA (24-26). The intact complex, however, is needed for the assembly of PCNA around the DNA and, to date, no subunit combination analyzed performs this task. Unloading is probably a simpler mechanism, as the p40 subunit is sufficient to unload PCNA from DNA and also interacts with pold (26). The role of this interaction in DNA synthesis, however, is not yet clear. RF-C also interacts with SSB. The interaction between RF-C and SSB might be needed for the assembly of the clamp, as demonstrated for the *E. coli* g-complex.

### 1.3. gp44/62

Similar to the g-complex and RF-C, the clamp loader of phage T4, the gp44/62 complex, is also composed of five subunits (Table 1). In this case, however, only two polypeptides form the pentamer, where a tetramer of gp44 binds a single subunit of gp62 (27). Like the clamp loader of eukaryotes and prokaryotes, the gp44/62 complex exhibits ATPase activity that is stimulated by the sliding clamp (gp45) and by DNA (28, reviewed in Ref. 29). A DNA structure that resembles a

template/primer junction is the best effector for stimulating ATPase activity (27). Cross-linking and protein-DNA footprinting assays demonstrate that the clamp loader interacts with DNA in the absence of gp45, but these interactions are relatively weak, and the clamp is needed for binding to DNA (reviewed in Refs. 29, 30). Studies with individual subunits of the gp44/62 complex revealed that the gp44 exhibits the ATPase activity, whereas gp62 interacts with the clamp (31).

The exact mechanism of action of the gp44/62 complex is not yet fully understood. The overall features of its activity, however, have been determined and are similar to those of the g-complex and RF-C (Fig. 1). The gp44/62 complex interacts with gp45 in an ATP-dependent manner. The complex recognizes the primer terminus, which stimulates the ATPase activity of the clamp loader, bringing about conformational changes that lock the clamp around the primer and allow it to interact with the polymerase to initiate processive DNA synthesis.

## 2. Concluding Remarks

The polymerases responsible for replicating chromosomal DNA during cell division are multiprotein complexes. The processivity of these enzymes relies on a ring-shaped protein that encircles DNA and tethers the catalytic unit to DNA for processive DNA synthesis. The mechanisms by which the clamp loaders assemble the ring around DNA and disassemble them from DNA are only beginning to be understood. The complexity of this mechanism is inferred from the observation that it requires the coordinate activity of several polypeptides in eukaryotes, prokaryotes, and phage T4. Future studies are needed to elucidate the mechanisms by which the clamp loaders operate in the loading and unloading reactions. Purification of the clamp loader complexes and the individual subunits will enable further analysis of these mechanisms.

## Bibliography

“Clamp Loaders, Processivity Complex” in , Vol. 1, pp. 472–474, by Zvi Kelman, University of Maryland, Biotechnology Institute, Rockville, MD and Lori M. Kelman, Montgomery College, Germantown, MD; “Clamp Loaders, Processivity Complex” in (online), posting date: January 15, 2002, by Zvi Kelman, University of Maryland, Biotechnology Institute, Rockville, MD and Lori M. Kelman, Montgomery College, Germantown, MD.

1. B. Stillman (1994) *Cell* **78**, 725–728.
2. Z. Kelman and M. O'Donnell (1994) *Curr. Opin. Genet. Dev.* **4**, 185–195.
3. B. Stillman (1996) *DNA Replication in Eukaryotic Cells*. M. L. DePamphilis, ed., CSH Laboratory Press, Cold Spring Harbor, NY, pp. 435–460.
4. M. O'Donnell, R. Onrust, F. B. Dean, M. Chen, and J. Hurwitz (1993) *Nucleic Acids Res.* **21**, 1–3.
5. G. Cullmann, K. Fien, R. Kobayashi, and B. Stillman (1995) *Mol. Cell. Biol.* **15**, 4661–4671.
6. N. Yao, J. Turner, Z. Kelman, P. T. Stukenberg, F. Dean, D. Shechter, Z.-Q. Pan, J. Hurwitz, and M. O'Donnell (1996) *Genes to Cell* **1**, 101–113.
7. C. A. Wu, E. L. Zechner, and K. J. Marians (1992) *J. Biol. Chem.* **267**, 4030–4044.
8. P. T. Stukenberg, J. Turner, and M. O'Donnell (1994) *Cell* **78**, 877–887.
9. K. J. Hacker and B. M. Alberts (1994) *J. Biol. Chem.* **269**, 24221–24228.
10. A. Yuzhakov, J. Turner, and M. O'Donnell (1996) *Cell* **86**, 877–886.
11. Z. Kelman and M. O'Donnell (1995) *Ann. Rev. Biochem.* **64**, 171–200.
12. C. S. McHenry (1988) *Ann. Rev. Biochem.* **57**, 519–550.
13. M. O'Donnell and P. S. Studwell (1990) *J. Biol. Chem.* **265**, 1179–1187.
14. R. Onrust, P. T. Stukenberg, and M. O'Donnell (1991) *J. Biol. Chem.* **266**, 21681–21686.
15. Z. Tsuchihashi and A. Kornberg (1989) *J. Biol. Chem.* **264**, 17790–17795.
16. V. Naktinis, R. Onrust, L. Fang, and M. O'Donnell (1995) *J. Biol. Chem.* **270**, 13358–13365.

17. T. Tsurimoto and B. Stillman (1989) *Mol. Cell. Biol.* **9**, 609–619.
18. U. Hubscher, G. Maga, and V. N. Podust (1996) *DNA Replication in Eukaryotic Cells*, M. L. DePamphilis (ed.), CSH Laboratory Press, Cold Spring Harbor, NY, pp. 525–543.
19. S.-H. Lee, A. D. Kwong, Z.-Q. Pan, and J. Hurwitz (1991) *J. Biol. Chem.* **266**, 594–602.
20. T. Tsurimoto and B. Stillman (1991) *J. Biol. Chem.* **266**, 1950–1960.
21. F. Uhlmann, J. Chi, H. Flores-Rozas, F. B. Dean, J. Finkelstein, M. O'Donnell, and J. Hurwitz (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6521–6526.
22. V. N. Podust and E. Fanning (1997) *J. Biol. Chem.* **272**, 6303–6310.
23. X. Li and P. M. Burgers (1994) *Proc. Natl. Acad. Sci. USA* **91**, 868–872.
24. M. A. McAlear, E. A. Howell, K. K. Espenshade, and C. Holm (1994) *Mol. Cell. Biol.* **14**, 4390–4397.
25. M. Chen, Z.-Q. Pan, and J. Hurwitz (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2516–2520.
26. Z.-Q. Pan, M. Chen, and J. Hurwitz (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6–10.
27. T. C. Jarvis, L. S. Paul, and P. H. von Hippel (1989) *J. Biol. Chem.* **264**, 12709–12716.
28. J. R. Piperno and B. M. Alberts (1978) *J. Biol. Chem.* **253**, 5174–5179.
29. M. C. Young, M. K. Reddy, T. C. Jarvis, E. P. Gogol, M. K. Dolejsi, and P. H. von Hippel (1994) *Molecular Biology of Bacteriophage T4*. (J. D. Karam, ed.), American Society for Microbiology, Washington, DC, pp. 313–317.
30. N. G. Nossel (1994) *Molecular Biology of Bacteriophage T4*. (J. D. Karam, ed.), American Society for Microbiology, Washington, DC, pp. 43–53.
31. J. Rush, T.-C. Lin, M. Quinones, E. K. Spicer, I. Douglas, K. R. Williams, and W. H. Konigsberg (1989) *J. Biol. Chem.* **264**, 10943–10953.

### Suggestions for Further Reading

32. U. Hubscher, G. Maga, and V. N. Podust (1996) "DNA replication accessory proteins", In *DNA Replication in Eukaryotic Cells*. (M. L. DePamphilis, ed.), CSH Laboratory Press, Cold Spring Harbor, NY, pp. 525–543.
33. Z. Kelman and M. O'Donnell (1994) DNA replication enzymology and mechanisms, *Curr. Opinion Genet. Dev.* **4**, 185–195.
34. M. O'Donnell (1992) Accessory protein function in the DNA polymerase III holoenzyme from *E. coli*. *BioEssays* **14**, 105–111.
35. B. Stillman (1996) "Comparison of DNA replication in cells from prokarya and eukarya", In *DNA Replication in Eukaryotic Cells* (M. L. DePamphilis, ed.), CSH Laboratory Press, Cold Spring Harbor, NY, pp. 435–460.
36. M. C. Young, M. K. Reddy, and P. H. von Hippel (1992) Structure and function of the bacteriophage T4 DNA polymerase holoenzyme, *Biochemistry* **31**, 8675–8690.

### Class Switching

Upon stimulation with an [immunogen](#), the circulating [antibodies](#) belong first to the [IgM](#) isotype, before being replaced by [immunoglobulins](#) of another class (most frequently an [IgG](#) in the bloodstream). This phenomenon is known as [isotype](#) switching (or switch) and is the consequence of a gene rearrangement that takes place in the centrocytes, located in the light zone of the germinal center. As a result, the V–D–J rearranged section of the mature B cell, which was initially associated



with a Cm [constant region](#) to make a complete m chain, will now become associated with the constant region of another isotype. This gene rearrangement involves recognition sequences termed “[switch regions](#)” that are located at 5' of each constant CH gene (with the exception of the d chain gene). These regions are very similar in sequence and are presumably recognized by an [enzyme](#) system, which has not been identified to date, but does not involve [recombinase](#). As a result of this recombination event, the portion of DNA localized between the two switch regions concerned is deleted. The light-chain genes are not affected by the isotype switching, so the antibody combining site is unaffected. In other words, isotype switching has changed the Ig isotype without modifying the specificity of the cell. It has simply conferred on the antibody a distinct biological function, which will amplify the physiological action of the immune system in a number of ways, such as to fight more appropriately against pathogens, ensure fetal protection through transplacental transfer, be expressed in secretions, and so on, depending upon the selected isotype. Isotype switching is another example of a mechanism that necessitates interaction between [T cells](#) and [B cells](#) and makes use of two sets of molecules. One involves CD40 ligand (CD40-L) and CD40, expressed at the cell surfaces of T<sub>H</sub> and B lymphocytes, respectively, and ensuring direct contact between the cells. The second signal is provided by the T<sub>H</sub> cell and is a **cytokine** that is released as a soluble factor in the immediate proximity of the B cell. Depending upon the nature of the cytokine, and therefore on the type of interacting T cell, the switch mechanism will generate a discrete isotype. For example, secretion of **interleukin-4** (IL-4) will favor a switch toward the IgE class, whereas [transforming growth factor](#) b1 will induce switching to IgA.

See also entries [Antibody](#), [Immune Response](#), [Immunoglobulin](#), [Isotype](#), [Switch Region](#).

#### Suggestions for Further Reading

J. Stavnezer (1996) Immunoglobulin class switching. *Curr. Opin. Immunol.* **8**, 199–205.

T. Honjo and T. Kataoka (1978) Organisation of immunoglobulin heavy-chain genes and allelic deletion model. *Proc. Natl Acad. Sci. USA* **75**, 2140–2144.

## Clathrin

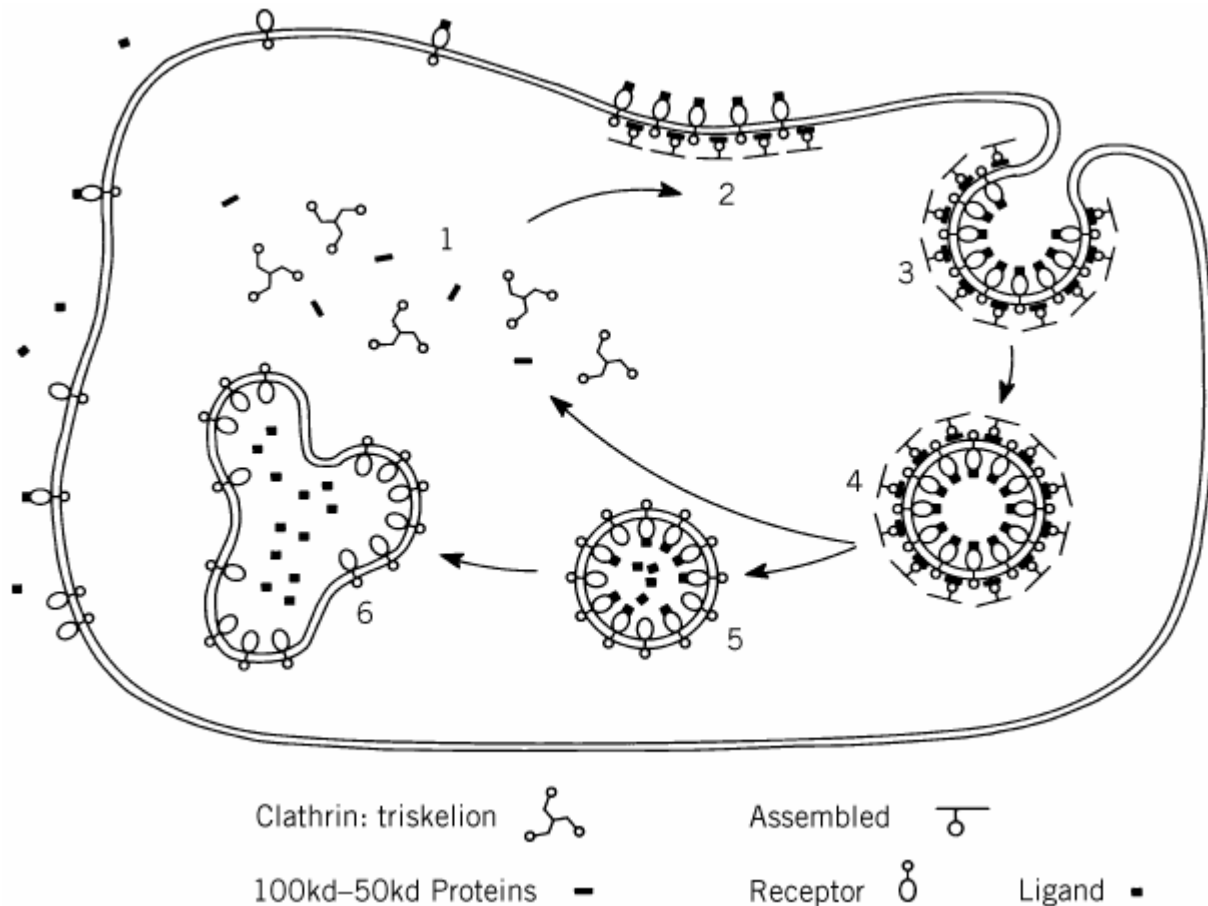
Clathrin is the major component of purified **coated vesicles** ([1](#)), intracellular structures that derive from coated membrane involved in [endocytosis](#) and **receptor** recycling.

Clathrin-mediated endocytosis ([2](#), [3](#)) is the major route for the uptake of specific macromolecules into cells for example, **low density lipoprotein** carrying cholesterol and [transferrin](#) bearing iron. Such macromolecules are concentrated from the surrounding cellular fluid by binding to the receptors exposed on the surface of cells. The receptors, in turn, have features in their cytoplasmic tails that enable them to cluster into clathrin-coated pits and thus be endocytosed efficiently in small vesicles. (Fig. [1](#)). The vesicles, in addition, carry other molecules that target them to an internal compartment, the [endosome](#), and allow them to fuse with the endosomal membrane, thus delivering the cargo of nutrients into the lumen for use by the cell. Meanwhile, the coat proteins and receptors are recycled through the endosomal compartment and back to the plasma [membrane](#). The coat proteins diffuse through the cytoplasm to reform coated pits there for a further round of endocytosis.

**Figure 1.** Endocytosis promoted by a round of clathrin assembly and recycling. [Taken from B.M.F. Pearse and R.A.

Crowther (1987) *Ann. Rev. Biophys. Biophys. Chem.* **16**, 49–68.] (1) Cytoplasmic clathrin triskelions and adaptors (containing specific 100 kDa adaptins and associated subunits) are triggered to assemble on the membrane in a reaction involving GTP-binding proteins with a selection of receptors to form a coated pit(2). The coated pit invaginates as further receptors and coat proteins assemble(3). Pinching-off the completed coated vesicle(4) requires dynamin to promote fusion in the neck region. Uncoating follows, releasing the vesicle (5), which carries molecules primed for fusion with the endosome(6), and the soluble coat components, which recycle(1).

A similar coated vesicle cycle takes place in many nerve terminals. In this case, the cycle replenishes the synaptic vesicle population, whose efficient fusion with the presynaptic membrane, on stimulation, depends on, among other molecules, synaptobrevin, syntaxin, and munc18 (29-32).



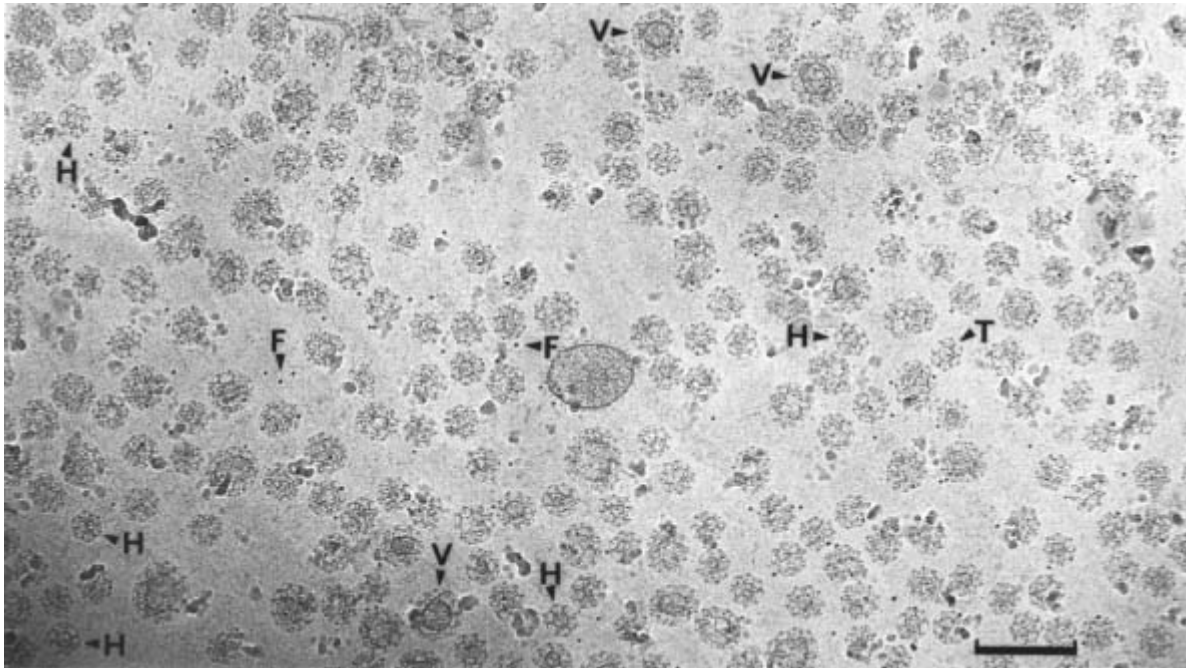
This system of budding coated vesicles is used again and again throughout cellular biology. Thus clathrin-coated vesicles are involved in **antigen presentation** in immunology (4) in delivering receptors and signaling molecules at the right place and time during [development](#) (5) and in recycling synaptic vesicle components in nerve synapses (6). In essence, clathrin and its associated proteins provide a mechanism for efficient recycling of specific receptors that are controlled in various ways. Other non-clathrin coated vesicle systems that contain a spectrum of related molecules and certain quite distinct coat components mediate sorting and recycling steps throughout the secretory pathway (7, 8). Clathrin-coated vesicles are more particularly associated with sorting specific molecules from the trans **Golgi** network into the pre-lysosomal/endosome network in addition to their endocytic function.

Unfortunately, **viruses** [eg, influenza virus (9)] exploit the endocytic system to infect cells, and various other pathogens and foreign substances do damage to cells by this route. Defects in the system are also the root cause of certain medical conditions, for example, familial hypercholesterolaemia (2), and I-cell disease (10).

## 1. Clathrin

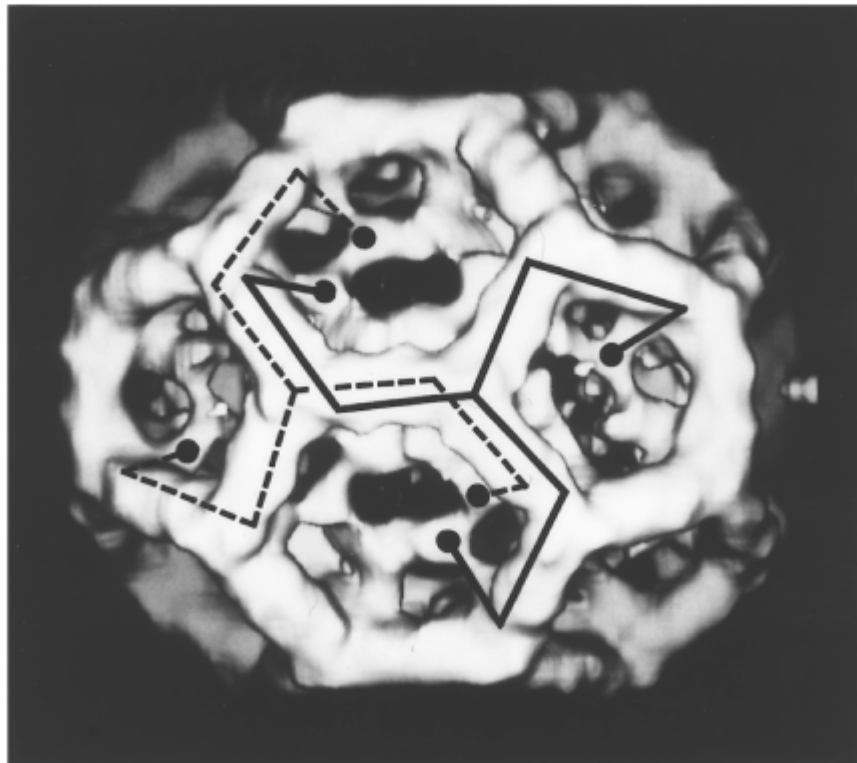
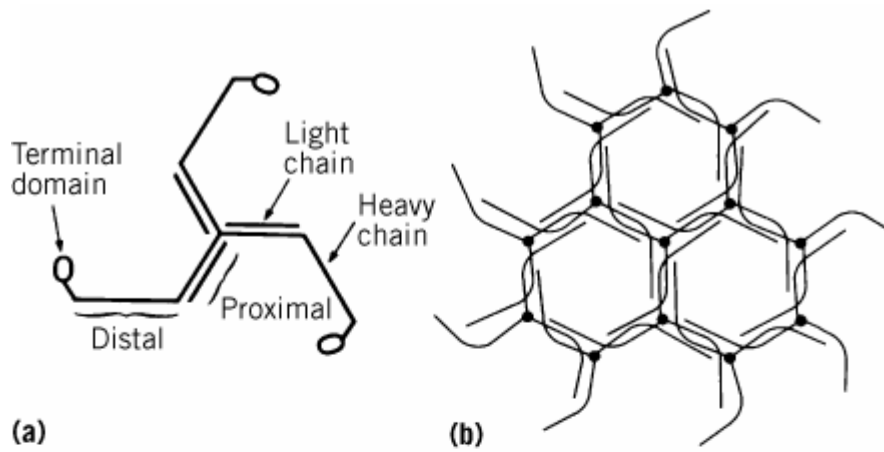
The function of clathrin is to form the strong outer, or cytoplasmic, surface of the coat, a remarkable honeycomb of hexagons and pentagons (Fig. 2). This flexible, structural network accommodates extensive areas of flattish membrane or encloses a range of vesicles, the smallest of which (250 Å diameter) is contained within a truncated icosahedron composed of 20 hexagons and 12 pentagons (for review, see Ref. 11).

**Figure 2.** Field of unstained placental coated vesicles in ice. Hexagonal barrels (H), tennis ball structures (T), larger coats containing vesicle (V), and ferritin (F) are indicated. Scale bar 200 nm. [Taken from G.P.A. Vigers, R.A. Crowther, and B.M.F. Pearse (1986) *EMBO Journal* 5, 529–534.]



The modular design of the clathrin molecule is extraordinary and unexpected [Fig. 3(a)]. The overall shape is that of a triskelion (12, 13). Three heavy chains (180 kDa) and three light chains (one LCa and 2 LCb) form the three-legged structure. The C-termini of the heavy chains come together to form the hub of the structure, and the N-termini are in the terminal domains of the extended legs. The light chains are located along the proximal leg regions, possibly contributing in some way to the geometry at the vertex. Beyond the light chains, the heavy chains exhibit a kink, and the distal ends bend yet again to form a more globular terminal domain.

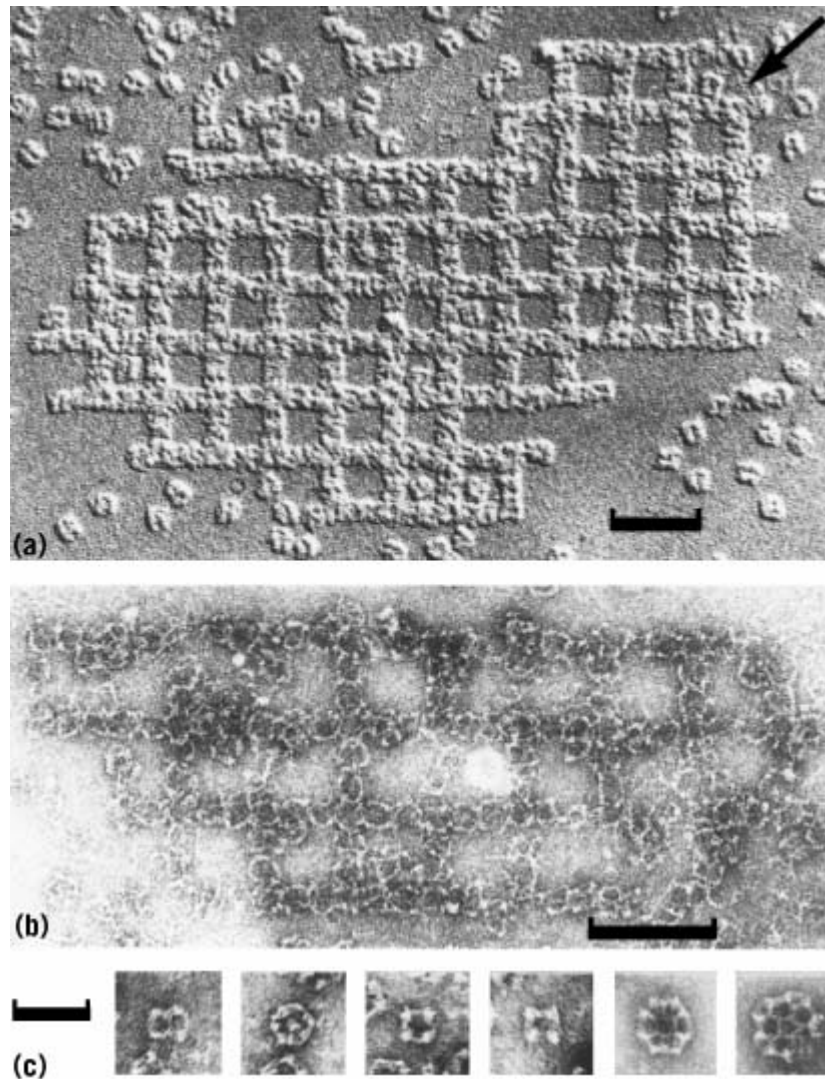
**Figure 3.** (a) Schematic drawing showing the modular structure of the triskelion. (b) Packing diagram showing how triskelions (lacking terminal domains for simplicity) form a hexagonal lattice. (c) Three-dimensional map of a clathrin cage containing 12 pentagons and 8 hexagons, computed from electron micrographs of unstained specimens embedded in vitreous ice. Each triskelion leg runs from one vertex along two neighboring polygonal edges and then turns inward. Its terminal domain forms the inner shell of density. [Taken from B.M.F. Pearse and R.A. Crowther (1987) *Ann. Rev. Biophys. Biophys. Chem.* 16, 49–68.]



The packing arrangement of the triskelions (lacking terminal domains for simplicity) to form a hexagonal lattice is shown in Fig. 3(b). The profiles of two complete triskelions are shown at adjacent vertices of a clathrin cage in the form of a hexagonal barrel in Fig. 3(c).

*In vitro* in artificial conditions, clathrin exhibits its versatility as a construction molecule (Fig. 4). Among other small particles, clathrin makes cubes [Fig. 4(c)] a variant of the normal cage (14). In turn the cubes pack together to form remarkable arrays that resemble the foundations of buildings [Fig. 4(a) and (b)].

**Figure 4.** The most astonishing type of clathrin aggregate produced *in vitro* is an open square packing of cubes in a pattern reminiscent of the foundations of an ancient building visualized by (a) unidirectional shadowing; (b) negative staining in uranyl acetate (bar 0.2  $\mu\text{m}$ ); (c) comparison of cubes with cages, showing from left to right, two-fold, three-fold, and two four-fold views of the cube plus cages of the hexagonal barrel and truncated icosahedron type (football). The edge of the cube is more than twice as long as the vertex to vertex distance in the cages. Bar, 0.1  $\mu\text{m}$ . [Taken from P.K. Sorger, R.A. Crowther, J.T. Finch, and B.M.F. Pearse (1986) *J. Biol.* **103**, 1213–1219.]



cDNAs coding for clathrin heavy chains and light chains have been **cloned** and sequenced (15, 16). The sequences suggest that part of the light chains form a **coiled-coil** structure with a corresponding part of the heavy chain to form the proximal legs. Extensive studies of the interaction between light chains and heavy chains, using **antibodies** and **mutational** analysis, have provided further evidence for the positioning of the light chains along the proximal leg, studies also indicate that the C-termini of the light chains occur at the vertex, where they might influence the structural angles involved in cage assembly (see, eg, Ref. 17). Now the prospects are good for obtaining crystals suitable for determining the high-resolution structure by **X-ray crystallography** from material produced by expressing parts of the triskelion (e.g., the hub region) in *E. coli* (17).

An interesting problem is how spurious clathrin cage formation is prevented in the cytoplasm or indeed how clathrin triskelions are disassembled from the coat structure after vesicle budding but are not prevented from forming coated pits. Recently, other attendant “**chaperone**” proteins, hsp70c and cofactor auxilin, have been found, which modulate cage assembly in the cytoplasm (18).

Cloning of clathrin genes has also led to further exploration of the role of clathrin by gene knockout and mutation. Deletion of the clathrin gene in yeast makes various strains very sick, slow growing, or dead. In those that survive, membrane organization is affected. In particular, the processing of pre-pro-a-factor by the kex2 endoprotease is disrupted, leading to secretion of the immature a-factor and failure in sexual reproduction, and endocytosis of specific molecules, for example, kex2,

is reduced, and the cells fill with abnormal vacuoles (19). In a temperature-sensitive mutant of clathrin, transient defects have been observed in sorting to the vacuole. The picture emerging (20) is not dissimilar to that observed in mammalian cells, that is, clathrin is involved in specific sorting steps both in the trans Golgi region and during endocytosis at the plasma membrane. These are the sites where the characteristic coated structures, now confirmed as clathrin-coated pits, were seen in abundance in early **electron microscope** pictures of fixed cell sections (21, 22). In a more complex creature, the **slime mold**, *Dictyostelium discoideum* (23), failure of clathrin heavy chain expression in cells impairs endocytosis and causes a lack of endosomes and contractile vacuoles, leading to defects in osmoregulation. Such cells also cannot follow the developmental program.

In *Drosophila melanogaster*, the clathrin heavy chain gene is essential (24). However, in flies, a dramatic effect is caused by the shibire mutation. This is a temperature-sensitive mutation in the molecule, now known to be **dynamin**, that is required for budding off a clathrin-coated vesicle (25, 26). At the nonpermissive temperature, the flies drop down as if dead. They cannot recycle their synaptic vesicle components in synapses and therefore cease to fly. However, when cooled down again, they start to fly as usual. These results confirm and extend the original electron microscope observations of abundant clathrin-coated vesicles in nerve synapses (27) and particularly the neuromuscular junction (28). Recent studies have identified many more components in the synaptic vesicle cycle (6) and explored their role by genetic manipulation in *Drosophila* and *Caenorhabditis elegans* [see Fig. 1; (29-32)].

In the worm, *C. elegans*, clathrin-mediated sorting has been implicated in the development of the vulval region (5). This study suggests that even apparently quite subtle perturbations in sorting by the coated vesicle system have profound effects in the development of a complex organism.

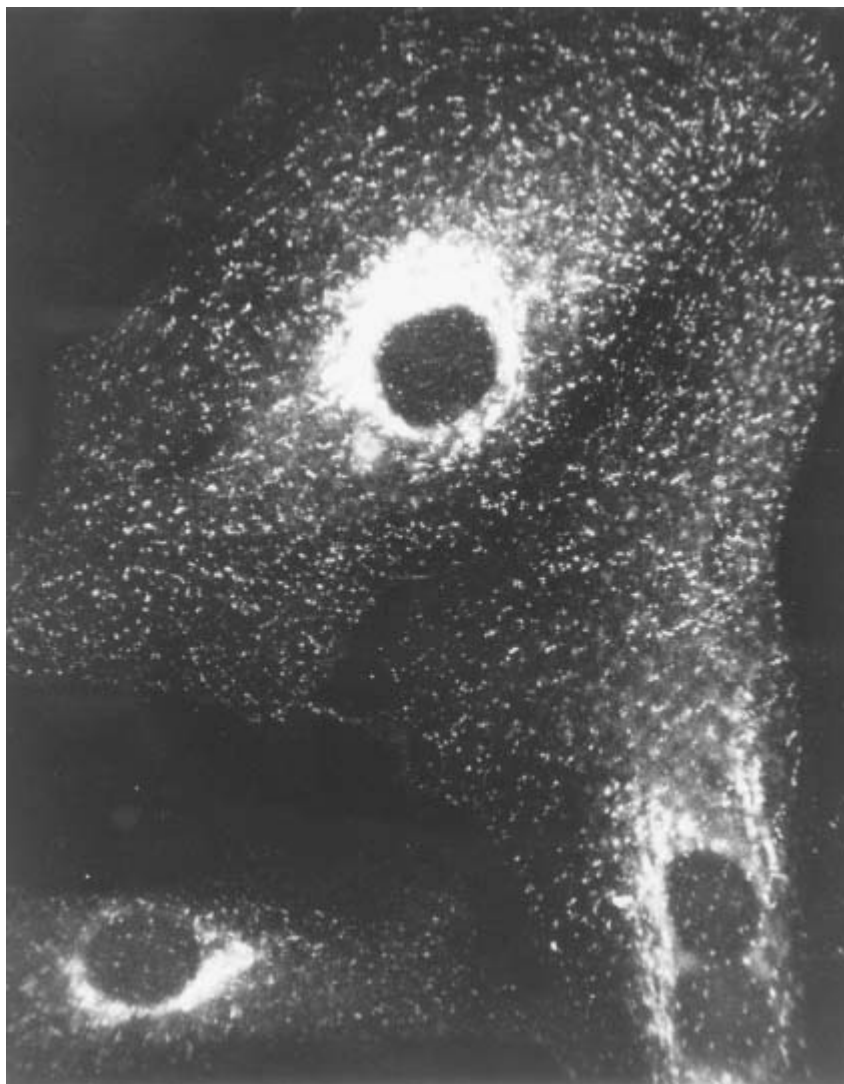
Recently, a second clathrin heavy-chain gene (CLTD) has been identified in humans, which has its maximal level of expression in skeletal muscle (33). This gene was found in the region commonly deleted in velo-cardio-facial syndrome (VCFS). Based on the location and expression pattern of CLTD, the suggestion is that **hemizygoty** at this locus plays a role in the etiology of one of the VCFS-associated **phenotypes**.

In summary, the clathrin molecule is an extraordinary building unit that has an intricate packing arrangement and forms coated structures of striking beauty. It carries out an important function.

The assembly properties of clathrin with the coated structures and vesicles it encloses have allowed purifying and identifying many of the other functional components involved. Chief among these are the clathrin adaptors, heterotetrameric complexes that coassemble with clathrin to form the vesicle coat.

### 1.1. Clathrin Adaptors

Two distinct types of clathrin adaptor complexes have been identified (3, 34). One of these (the PM-adaptor) consists of an a-adaptin (~100 kDa) combined with b-adaptin (~100 kDa) and two smaller subunits, AP50 and AP17. This adaptor is found by **immunofluorescence** with a **monoclonal antibody** against a-adaptin mainly in plasma-membrane-coated pits. In contrast, the second, the Golgi adaptor, consists of g-adaptin (~95 kDa) combined with b'-adaptin and two other subunits, AP47 and AP20, and is found by immunofluorescence using a monoclonal antibody against g-adaptin in coated pits in the trans Golgi network (Fig. 5).



Then the problem immediately arises of how these different adaptor coat complexes assemble on the others. In fact, the problem increases as new related adaptors (35), coats and coatomer complexes (7) are identified in the intracellular membrane system. The problem is still incompletely understood, although numerous ways in the assembly of particular coated vesicles, their budding, and the control of membrane trafficking have been identified in the cytoplasm and on membranes by a combination of **yeast** genetics and biochemistry. An earlier control protein identified in post-Golgi membrane transport is Sec4, a small **GTP-binding** protein. A range of GTP-binding proteins have been implicated throughout the recycling system, including Rab proteins (36), the **Rho** (38) and heterotrimeric G-proteins (39). For instance, distinct Rab proteins take up a characteristic particular membrane compartments in the cell and are found in coated pits on those membranes (40). The kinetics of fusion of vesicles from one compartment with another (41), although precisely how this is unclear. ARF, on the other hand, has been implicated in coat binding to membranes notably in the formation of clathrin-coated vesicles. Dynamins are concerned with actual budding of coated vesicles, as previously mentioned (see Fig. 1).

When clathrin-coated vesicles assemble on a membrane, they concentrate certain specific receptors and other membrane proteins behind. These specific receptors are recognized by features in their cytoplasmic tails. Each receptor finds its own characteristic **steady-state** distribution throughout the sorting system because each receptor has a specific amino acid sequence. There is no single, defined, clear-cut sequence indicating that a receptor will assemble efficiently. Nevertheless, at least in endocytosis, mutational analysis has pinpointed certain small groups of residues in the cytoplasmic tails of receptors that constitute primary sequence 'features' which allow the normal accumulation of their receptors on membranes. Typically, such a 'feature' contains four residues, the first and probably most important of which is a tyrosine (T) residue (one or both of which is often charged) followed by a bulky aliphatic residue. Such tetrapeptides are recognized by the AP50 subunit of the plasma membrane adaptor (44). However, as more examples are studied, it becomes clear that other possible signals exist and that other factors, such as phosphorylation of receptor tails, play a role in p

(46).

One of the most investigated examples of a receptor that is routed via clathrin-coated vesicles containing a Golgi network (TGN) is the cation-independent mannose 6-phosphate/IGF-II receptor (MPR). The MPR recognizes luminal sorting signals, **mannose 6-phosphate**-containing sugar chains on **lysosomal** enzymes, and also serves as cargo inside coated vesicles on their way to lysosomes (47, 48). The MPR has a 163-residue cytoplasmic tail with features involved in directing the receptor along its complex, intracellular, recycling pathway, including Tyr-Lys-Tyr-Ser-Lys-Val (49, 50). The tail region containing these residues, apparently, also acts as a sorting signal for lysosomal enzymes, separately or in combination with the C-terminal region (51).

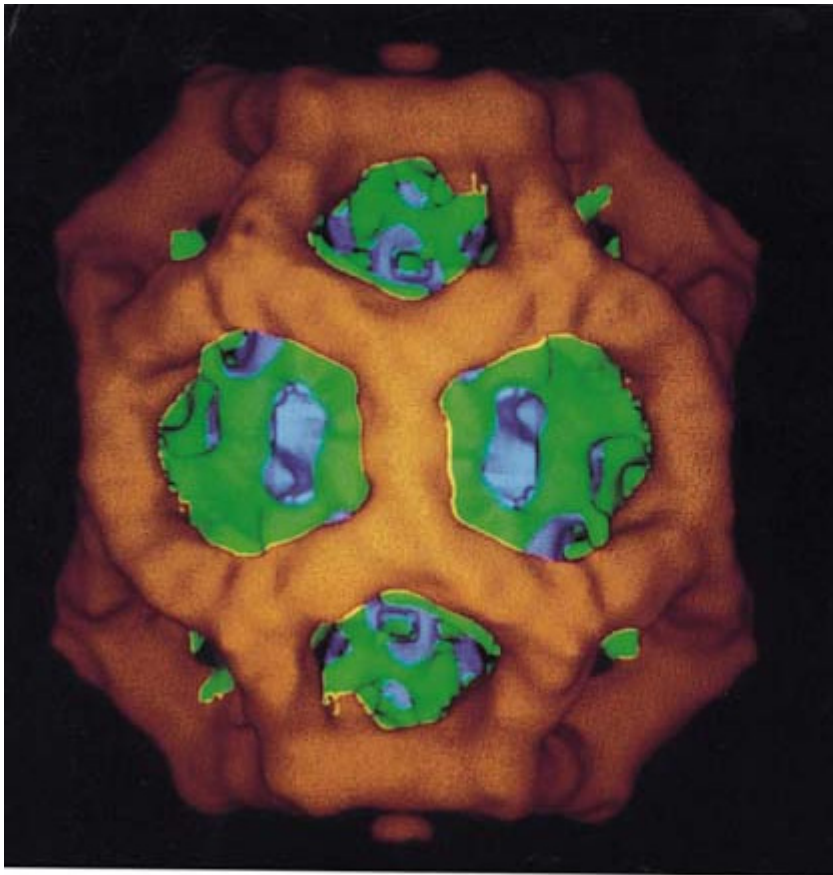
Both types of coated vesicle adaptors bind to the MPR tail. Recognition by the plasma membrane adaptor is dependent on the tyrosine residues important for the endocytosis signal, but these same mutations do not abolish the clathrin adaptor (52). There is also evidence that the Golgi adaptor preferentially binds to the MPR tail when in addition to these cytoplasmic features that determine the routing of the MPR, the extracellular domain also plays some role in the precise steady-state distribution of the receptor, at least in CV1 and COS cells (53).

### 1.2. Clathrin Coats

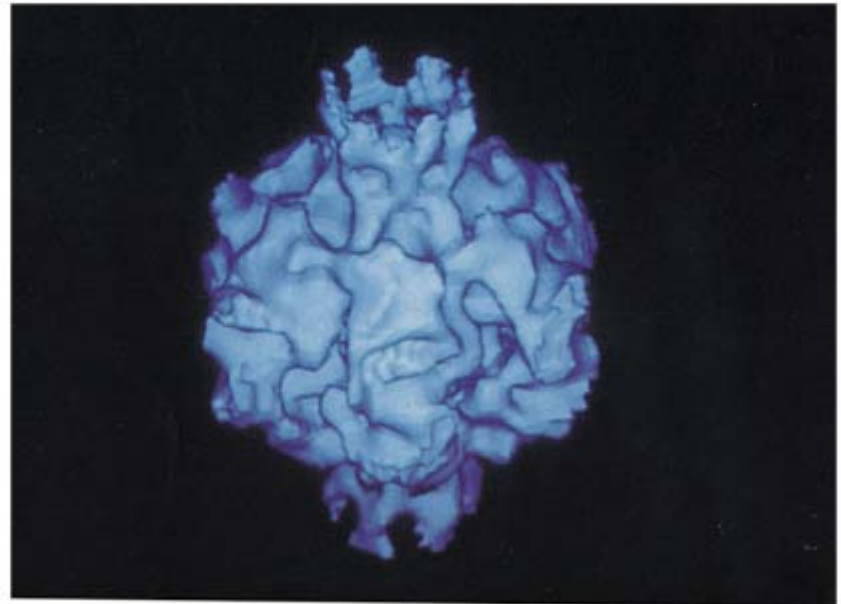
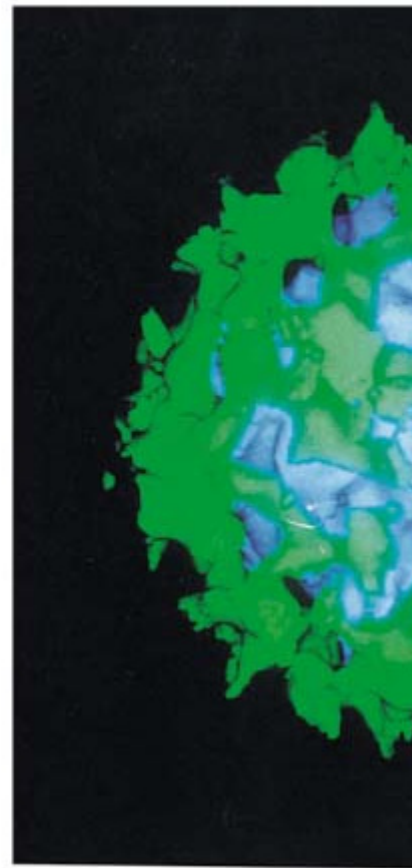
A low ( $\sim 50$  Å) resolution three-dimensional map of a clathrin coat has been generated from electron cryomicroscopy in vitreous ice (Fig. 6) (see [Electron Imaging](#)). The major components of the coat, clathrin, and adaptor were purified from purified coated vesicles and separated. Three types of particles were reassembled from these preparations: the first was the clathrin cage itself, the second was derived from the first by partial trypsin digestion to remove the triskelions, and the third was the complete coat containing both clathrin and adaptors. Maps of each type were obtained by computer imaging of tilted specimens of several examples of individual particles of the three different types. The incomplete structures from those of the complete coat, three different layers of the overall coat were highlighted in the image presented (Fig. 6). Thus the outer polyhedral clathrin lattice is highlighted in red, the **triskelion domains** extends into the structure (green), and the adaptors form an inner shell (blue). These particles, in the absence of a vesicle, and in fact are actually too small to accommodate a vesicle. Probably the smallest particles that contain a vesicle is a truncated icosahedron, which can also be reconstituted as a coat. If purified coats are imaged in projection they exhibit the same coat thickness composed of three shells of density, whereas the vesicles of high density contained within the coat. Now a higher ( $\sim 20$  Å) resolution version of the coat structure has been obtained from types of hexagonal barrel preparations and the equivalent icosahedral specimens by using improved cryomicroscopy and highly developed imaging techniques currently available. This will allow further identification of small features, especially if combined with specific antibodies to decorate those domains. It is also possible to generate a 3D map of vesicles from their cytoplasmic constituents and membranes enriched in relevant proteins in physiological conditions, an exciting prospect as it may lead to a greater understanding of the specificity underlying the assembly and function of these organelles and to locate some of the controlling elements in the structure.

**Figure 6.** See color insert section. Clathrin encapsulates coated vesicles, the organelles responsible for uptake of essential nutrients (especially across the placenta in mammals and into oocytes in chickens). (a) A three-dimensional map of a clathrin coat reassembled from their constituent protein complexes. Colors highlight the outer polyhedral clathrin lattice (red), the **triskelion domains** (green) revealed in (b) and an inner shell of adaptors (blue) exposed in (c). These three maps led to the descriptive LEGO model. (Courtesy of G.P.A. Vigers, R.A. Crowther, and B.M.F. Pearse from data presented in *EMBO J.* (1986) 5:111-120.)





(a)



(c)

### Bibliography

1. B. M. F. Pearse (1975) *J. Mol. Biol.* **97**, 93–98.
2. J. L. Goldstein, M. S. Brown, R. G. W. Anderson, D. W. Russell, and W. J. Schneider (1985) *Ar*
3. B. M. F. Pearse and M. S. Robinson (1990) *Ann. Rev. Cell Biol.* **6**, 151–171.

4. P. R. Wolf and H. L. Ploegh (1995) *Ann. Rev. Cell Dev. Biol.* **11**, 267–306.
5. J. Lee, G. D. Jongeward, and P. W. Sternberg (1994) *Genes Dev.* **8**, 60–73.
6. P. R. Maycox, E. Link, A. Reetz, S. A. Morris, and R. Jahn (1992) *J. Cell Biol.* **118**, 1379–1388.
7. R. Schekman and L. Orci (1996) *Science* **271**, 1526–1533.
8. J. E. Rothman and F. T. Wieland (1996) *Science* **272**, 227–234.
9. T. Stegmann, R. W. Doms, and A. Helenius (1989) *Ann. Rev. Biophys. Biophys. Chem.* **18**, 187
10. S. Kornfeld (1986) *J. Clin. Invest.* **77**, 1–6.
11. B. M. F. Pearse and R. A. Crowther (1987) *Ann. Rev. Biophys. Biophys. Chem.* **16**, 49–68.
12. E. Ungewickell and D. Branton (1981) *Nature* **289**, 420–422.
13. T. Kirchhausen and S. C. Harrison (1981) *Cell* **23**, 755–761.
14. P. K. Sorger, R. A. Crowther, J. T. Finch, and B. M. F. Pearse (1986) *J. Cell Biol.* **103**, 1213–12
15. A. P. Jackson, H. Seow, N. Holmes, K. Drickamer, and P. Parham (1987) *Nature* **326**, 154–159.
16. T. Kirchhausen, S. C. Harrison, E. Ping Chow, R. J. Mattaliano, K. L. Ramachandran, J. Smart, *Acad. Sci. USA* **84**, 8805–8809.
17. S. Liu, M. Wong, C. S. Craik, and F. M. Brodsky (1995) *Cell* **83**, 257–267.
18. E. Ungewickell, H. Ungewickell, S. E. H. Holstein, R. Lindner, K. Prasad, W. Barouch, B. Mart (1995) *Nature* **378**, 632–635.
19. G. S. Payne and R. Schekman (1989) *Science* **245**, 1358–1365.
20. E. Conibear and T. H. Stevens (1995) *Cell* **83**, 513–516.
21. T. F. Roth and K. R. Porter (1964) *J. Cell Biol.* **20**, 313–332.
22. D. S. Friend and M.G. Farquhar (1967) *J. Cell Biol.* **35**, 357–376.
23. T. J. O'Halloran and R. G. W. Anderson (1992) *J. Cell Biol.* **118**, 1371–1377.
24. C. Bazinet, A. L. Katzen, M. Morgan, A. P. Mahowald, and S. K. Lemmon (1993) *Genetics* **134**
25. T. Kosaka and K. Ikeda (1983) *J. Cell Biol.* **97**, 499–507.
26. H. Damke, T. Baba, A. M. van der Blik, and S. L. Schmid (1995) *J. Cell Biol.* **131**, 69–80.
27. E. G. Gray and R. A. Willis (1970) *Brain Res.* **24**, 149–168.
28. J. E. Heuser and T. S. Reese (1973) *J. Cell Biol.* **57**, 315–344.
29. A. DiAntonio and T. L. Schwarz (1994) *Neuron* **12**, 909–920.
30. K. Broadie, A. Prokop, H. J. Bellen, C. J. O'Kane, K. L. Schulze, and S. T. Sweeney (1995) *Neu*
31. Y. Fujita, T. Sasaki, K. Fukui, H. Kotani, T. Kimura, Y. Hata, T. C. Sudhof, R. H. Scheller, and **271**, 7265–7268.
32. E. M. Jorgensen, E. Hartweg, K. Schuske, M. L. Nonet, Y. Jin, and H. R. Horvitz (1995) *Nature*
33. H. Sirotkin, B. Morrow, R. DasGupta, R. Goldberg, S. R. Patanjali, G. Shi, L. Cannizzaro, R. Sh Kucherlapati (1996) *Human Molecular Genetics* **5**, 617–624.
34. S. Ahle, A. Mann, U. Eichelsbacher, and E. Ungewickell (1988) *EMBO J.* **7**, 919–929.
35. F. Simpson, N. A. Bright, M. A. West, L.S. Newman, R. B. Darnell, and M. S. Robinson (1996)
36. P. Cosson, C. Démollière, S. Hennecke, R. Duden, and F. Letourneur (1996) *EMBO J.* **15**, 1792
37. A. Salminen and P. J. Novick (1987) *Cell* **49**, 527–538.
38. C. Lamaze, T. Chuang, L. J. Terlecky, G.M. Bokoch, and S. L. Schmid (1996) *Nature* **382**, 177–
39. R. H. Kehlenbach, J. Matthey, and W. B. Huttner (1994) *Nature* **372**, 804–809.
40. P. Chavrier, R.G. Parton, H. P. Hauri, K. Simons, and M. Zerial (1990) *Cell* **62**, 317–329.
41. V. Rybin, O. Ullrich, M. Rubino, K. Alexandrov, I. Simon, M. C. Seabra, R. Goody, and M. Zer
42. A. Peyroche, S. Paris, and C. L. Jackson (1996) *Nature* **384**, 479–481.
43. P. Chardin, S. Paris, B. Antonny, S. Robineau, S. Béraud-Dufour, C. L. Jackson, and M. Chabre
44. W. Boll, H. Ohno, Z. Songyang, I. Rapoport, L. C. Cantley, J. S. Bonifacino, and T. Kirchhausen

5795.

45. I. V. Sandoval and O. Bakke (1994) *Trends Cell Biol.* **4**, 292–297.
46. A. Pelchen-Matthews, I. J. Parsons, and M. Marsh (1993) *J. Exp. Med.* **178**, 1209–1222.
47. S. Kornfeld (1992) *Ann. Rev. Biochem.* **61**, 307–330.
48. J. Klumperman, A. Hille, T. Veenendaal, V. Oorschot, W. Stoorvogel, K. von Figura, and H. J. C. (1997) *J. Biol. Chem.* **272**, 997–1010.
49. P. Lobel, K. Fujimoto, R. D. Ye, G. Griffiths, and S. Kornfeld (1989) *Cell* **57**, 787–796.
50. M. Jadot, W. M. Canfield, W. Gregory, and S. Kornfeld (1992) *J. Biol. Chem.* **267**, 11069–11077.
51. K. F. Johnson and S. Kornfeld (1992) *J. Cell Biol.* **119**, 249–257.
52. J. N. Glickman, E. Conibear, and B. M. F. Pearse (1989) *EMBO J.* **8**, 1041–1047.
53. R. Le Borgne, A. Schmidt, F. Mauxion, G. Griffiths, and B. Hoflack (1993) *J. Biol. Chem.* **268**, 11069–11077.
54. E. Conibear and B. M. F. Pearse (1994) *J. Cell Science* **107**, 923–932.

## Cleveland Map

This method of [peptide mapping](#) involves generating of [protein](#) fragments via limited **proteolysis**, usually in the presence of **SDS**, followed by electrophoretic separation in [SDS-PAGE](#). The resulting one-dimensional peptide patterns are characteristic of the protein substrate and are used to evaluate structural relationships between proteins and to aid identification of individual polypeptide chains. The technique is named after the first author of a 1977 publication describing the procedure (1), and many variants of this method have subsequently been presented (2).

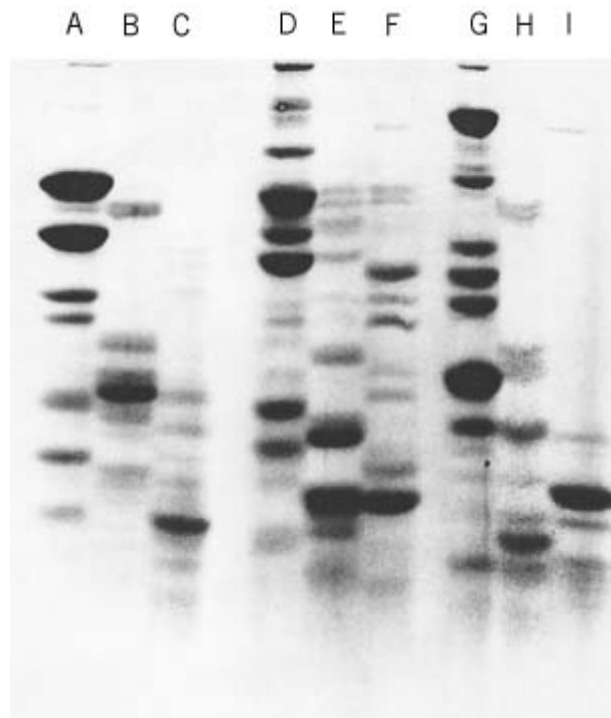
After initial separation and excision from a primary gel, the protein substrate is digested in solution or, typically, in the stacking zone of the SDS-PAGE “mapping gel.” The solution approach relies on purifying the protein by conventional chromatographic procedures or preparative SDS-PAGE followed by excision and [electroelution](#) from the gel (3). Protein preparation in SDS-containing gels is feasible because the subsequent digestion is performed in the presence of SDS. Frequently, SDS-PAGE isolation of the protein is possible at an early stage of purification, which is important when only a small amount of the protein is available. For Cleveland mapping, the efficient approach is to place the gel slice that contains the protein band of interest directly into a well of a second SDS/polyacrylamide gel and to add a suitable proteolytic enzyme. After electrophoresis just long enough for the protein substrate to migrate out of the gel piece and enter the stacking zone with the proteinase, the power is temporarily switched off and the gel left for digestion to take place. When the power is turned on again, the fragments of the protein substrate are separated from the proteinase and are mapped in the resolving gel. This procedure is preferred when handling small amounts of protein because of the direct transfer between primary and secondary gels.

The original Cleveland protocol (1) involved brief staining of the primary gel with [Coomassie Brilliant Blue](#), followed by rapid destaining and excision of the band corresponding to the protein substrate. Because fixation of the protein is undesirable due to the electrotransfer, negative staining with 1 M KCl may be an alternative and has been used for subsequent electroblotting and sequence analysis (4). The excised gel piece is equilibrated in buffer containing 0.1% (w/v) SDS to unfold the protein substrate and to obtain a suitable pH of 6.8 for stacking gel electrophoresis and digestion (1). Often, 5 to 10  $\mu\text{g}$  of protein substrate is required for application onto the primary gel to generate a

proteolytic pattern detectable by Coomassie Blue staining, but the minimum amount can be substantially lowered if **silver staining** **fluorescence** techniques or **radiolabeling** are used to visualize the peptide map.

After equilibration, the gel piece that contains the protein is placed in a well of the “mapping gel,” and, in the original protocol, overlaid with a solution of the appropriate proteinase in the same buffer. In later protocols, loading of the proteinase below the gel piece is suggested to avoid a variable depth of enzyme solution around the gel piece (2). Once a narrow stack has formed after the initial period of current, usually when the **tracking dye** (e.g., bromphenol blue) is approaching the bottom of the stacking gel, the power is turned off for 30 to 60 min. This is normally sufficient for the proteinase to cleave susceptible peptide bonds, which typically generates 10 to 20 fragments of 50- to 70-kDa proteins (1). It is crucial that the proteolytic enzyme is active in 0.1% SDS (up to 0.5% for in-solution digestions). Suitable enzymes are **trypsin**,  **$\alpha$ -chymotrypsin**, and **papain**, in addition to the popular V8 proteinase from *Staphylococcus aureus*. All are added in proportions corresponding to an enzyme to substrate ratio in the range of 1:10 to 1:50 (w:w). After digestion, electrophoresis is resumed, and the resulting fragments are resolved, typically in a 15 or 20% acrylamide gel of the discontinuous type (Fig. 1). A unique proteolytic pattern is obtained for each protein substrate and for each individual proteinase. This discriminating feature makes Cleveland mapping useful as a tool for visually judging the similarity between proteins.

**Figure 1.** One-dimensional Cleveland maps of **albumin** (A,D,G), **tubulin** (B,E,H), and **alkaline phosphatase** (C,F,I), generated by cleavage with papain (A–C), *Staphylococcus aureus* V8 proteinase (D–F), and chymotrypsin (G–I), followed by SDS-PAGE in a 20% (w/v) acrylamide gel. Each combination of protein substrate and proteinase gives a unique pattern. From Ref. 1 with permission.



In addition to enzymes, specific chemical cleavage is employed where the protein substrate within a gel is treated in a test tube. Suitable reagents are **cyanogen bromide** for cleavage after **methionine** residues, *N*-chlorosuccinimide for cleavage at **tryptophan**, **hydroxylamine** for cleavage at Asn-Gly peptide bonds, and formic acid for cleavage at Asp-Pro bonds. As for the size of the protein substrates, the method is generally applicable to polypeptide chains in the range 10 to 100 kDa.

Peptide fragments smaller than 10 kDa are difficult to resolve into informative peptide maps, although separating such small peptides is possible by using different SDS-PAGE systems. An alternative is to recover small fragments separately by [size exclusion chromatography](#) for subsequent peptide mapping by reverse-phase **HPLC**. This can be useful also for proteins large than 10 kDa in order not to waste small fragments (5).

### Bibliography

1. D. W. Cleveland et al. (1977) *J. Biol. Chem.* **252**, 1102–1106.
2. W. J. Gullick (1986) In *Practical Protein Chemistry: A Handbook* (A. Darbre, ed.), Wiley, Chichester, U.K., pp. 207–225.
3. M. W. Hunkapiller et al. (1983) *Methods Enzymol.* **91**, 227–236.
4. T. Bergman and H. Jörnvall (1987) *Eur. J. Biochem.* **169**, 9–12.
5. T. Bergman et al. (1991) *J. Protein Chem.* **10**, 25–29.

### Clonal Selection Theory

At the beginning of the twentieth century, Ehrlich was the first to propose that [antibodies](#) did exist prior to the introduction of [antigen](#) or [immunogen](#) and that antigenic stimulation only promoted amplification of molecules of the appropriate specificity. This proposal was quite pertinent, because it contained several key ideas, proven correct some 70 years later. One is the selective process of a preexisting [repertoire](#), the other is that [immunoglobulins](#) (Igs) may exist as surface or secreted molecules. The only aspect that was not verified is that Ehrlich imagined a cell to express many different antibodies. After a long period during which [hapten](#) analysis indicated the exquisite precision of [antibody–antigen interactions](#), it seemed impossible that antibody structures might preexist to recognize laboratory artifacts, like the haptens being studied. This resulted in the template theory of antibody formation, which required the presence of antigen prior to the synthesis of the corresponding antibodies. Revival of the selection theory was proposed in 1955 by Jerne, but it took another 4 years before the final form of the “clonal selection theory” was put forward by Burnet. The essence of the theory was that antibody specificities were preformed, but in a clonally distributed manner, so that one cell makes only one antibody of one given specificity. This was crucial to provide a mechanism accounting for acquisition of discrimination between self and nonself, after the observation on induction of immunological **tolerance** by the group of Medawar in 1953, when it was demonstrated that tolerance was an acquired phenomenon, resulting from a contact between the immature immune system and the self antigens during gestation or the perinatal period. Burnet most simply explained tolerance by a deletion of those clones that could interact with self components, leading to the concept of “forbidden clones.” Once the basic anti-nonself repertoire established, an antigen was seen as merely selecting and amplifying the cell population of clones expressing the corresponding specific antibodies.

Many investigators worked to define models that would prove the clonal theory, and the literature is amply documented on this point, with many very elegant approaches, such as those based on limiting dilutions, either *in vitro* or *in vivo* after transfer in irradiated mice. The ultimate proof came from molecular analysis, first at the [protein](#) level on myeloma proteins and then when the mechanism of [gene rearrangement](#) of Ig **genes** was elucidated, under the control of [allelic exclusion](#). Some exceptions to the clonal expression of Ig molecules have been periodically reported. Some were not exceptions, as in the case of lymphocytes coexpressing multiple [isotypes](#) that, in fact, share identical  $V_H$  and  $V_L$  regions, and therefore having the same specificity. Upon reactivation of the [recombinase](#)

genes upon [immunization](#), for example, a transient coexpression of different V genes may occur, leading to secondary gene rearrangements. But this rarely results in long-term production of antibodies of different specificities. The clonal theory is widely accepted and considered proven, and the fantastic use of [monoclonal antibodies](#) is (for most biologists) daily living proof of its pertinence.

See also entries [Allelic Exclusion](#), [Antibody](#), [Gene Rearrangement](#), and [Immunoglobulin](#).

#### Suggestions for Further Reading

M. F. Burnet (1959) *The Clonal Selection Theory of Acquired Immunity*, Vanderbilt University Press, Nashville, TN.

K. Rajewsky (1996) Clonal selection and learning in the antibody system. *Nature* **381**, 751–758.

M. Nussenzweig et al. (1987) Allelic exclusion in transgenic mice that express the membrane form of immunoglobulin. *Science* **236**, 816–819.

## Cloning

The word “cloning” has several different meanings. For example, it is possible to clone cells, that is, to cause cells to reproduce themselves so as to make a population of identical cells. In molecular biology, any piece of **DNA** can be cloned—inserted into a **vector** that replicates in a host to produce many copies of the same recombinant vector. It is also possible to clone **genes**. Gene cloning is the process of identifying and isolating a specific gene of interest. Cloning genes is a major focus of molecular biology, and the ability to clone genes has revolutionized biology. Once a gene is identified and cloned into an appropriate vector, it can be manipulated in many different ways. Cloning dramatically amplifies the DNA so that it can be **sequenced**. A cloned gene inserted into an [expression system](#) creates an organism that produces up to 25% of its total protein or more as the gene product, allowing large-scale production of the protein (1). A cloned gene can be mutated, and the mutated form inserted back into an organism that lacks a functional copy of the gene (2). This allows structure-function analysis of the gene product. A cloned gene or set of genes can be introduced into a new host to create a new metabolic pathway or to modify an existing pathway (3).

In many cases, gene cloning is carried out by fragmenting the **genomic** DNA from the organism containing the gene of interest and inserting the fragments into self-replicating DNA molecules (vectors) using [recombinant DNA](#) techniques; **transforming** the resultant population of recombinant molecules, which comprise a [DNA library](#), into a host organism, screening or selecting the transformants to identify cells that contain the desired DNA, and isolating the vector DNA.

### 1. Preparation of DNA for Cloning

Genomic DNA must be fragmented before cloning into a suitable vector. It is important to have fragments that contain the gene intact. One way to ensure this is to use a [restriction enzyme](#) that makes frequent cuts in the DNA under conditions of partial digestion (so that only some of the sites are cleaved) so as to obtain fragments of the desired size. The restriction enzyme *Sau3A* is often used because it has a four-base recognition sequence and therefore cuts DNA at many sites in the genome. Additionally, it produces [cohesive, sticky ends](#) that ligate with BamHI, a site often present in [polylinkers](#). Mechanical shearing also fragments the DNA, but sheared DNA cannot be cloned directly. First, the ragged ends must be filled in and polylinkers added.

Because small fragments of DNA **ligate** to a vector more readily than large fragments, the DNA to

be cloned is often fractionated by sucrose density gradient centrifugation (see [Sedimentation Velocity Centrifugation](#)) to obtain DNA of the desired size (pp. 2.85–2.86 of Ref. 4). It is important to set up the conditions for ligation so that most of the recombinant molecules produced contain only a single DNA fragment. This means that nearly equal numbers of vector molecules and fragment molecules should be present and that the total DNA concentration should not be too high (~40 µg of vector DNA per mL).

It is necessary to produce a large enough number of transformants or recombinant **bacteriophage** to have a high probability that a fragment carrying the intact gene is present in the population. The exact number depends on the size of the genome from which the gene is being cloned (the larger the genome, the more transformants need to be screened) and the size of the inserts (the larger the insert size, the fewer transformants need to be screened). A major problem in cloning is restriction enzymes in the host, which can degrade foreign DNA. Mutant strains of *Escherichia coli* are available that lack restriction enzymes, and these are best for constructing large [libraries](#) for screening or selection. Once a gene is cloned, restriction enzymes in a host are not usually a problem because large amounts of DNA overcome the restriction barrier (5).

**Prokaryotic** genes are usually cloned from genomic DNA, whereas **eukaryotic** genes are often cloned from [complementary DNA](#) (cDNA), which is prepared by copying [messenger RNA](#) with **reverse transcriptase**. Most eukaryotic genes contain **introns**, which prevent them from being expressed in prokaryotes, but the introns are not in the cDNA. The cDNA clone is often used as a probe to isolate the gene from a [genomic library](#). The genomic clone can be sequenced to identify the nature and location of the introns and the [upstream](#) and [downstream](#) sequences of the gene.

## 2. Cloning Vectors

There are four general types of vectors that are used in *E. coli*: **plasmids**, phage, **cosmids** or bacterial artificial [chromosomes](#) (BACs):

1. Plasmids have the advantage that they are smaller than the other vectors and usually give larger amounts of vector DNA from a given volume of culture because of their high **copy number**.
2. Phage vectors allow larger sized inserts of foreign DNA, so that a smaller number of **plaques** are needed to give a high probability that the entire genome of the organism is present in the library that is being screened. Furthermore, phage systems give more plaques per weight of DNA than plasmids give colonies, and plaques are easier to screen than colonies because they are smaller and more uniform. Phage T4 vectors allow cloning inserts as large as 120 kbp (6), whereas phage P1 vectors allow 80-kbp inserts (7), and [lambda phage](#) vectors allow up to 40-kbp inserts (pp. 2.2 to 2.125 of Ref. 4). Once a positive plaque is identified, the insert DNA is subcloned into a plasmid vector to take advantage of the greater amplification and ease of DNA preparation associated with plasmid vectors.
3. Cosmid vectors combine features of the preceding two classes because they can be packaged into phage but also replicate as plasmids. They also allow inserts up to 40 kbp (8).
4. Bacterial artificial chromosomes are based on the single-copy **F factor** plasmid and accept >300-kbp DNA fragments (9). Another important system for cloning very large DNA pieces is **yeast artificial chromosomes** (YAC). YACs accept fragments over 200 kbp and are used in genome sequencing and positional cloning of genes (10). Although **yeast** is not as easy to work with as *E. coli*, it is useful for studying eukaryotic genes that cannot be expressed in *E. coli*.

## 3. Identifying the Gene of Interest

Screening is usually the most difficult step, and there are many different approaches available. The simplest method is to complement a mutant gene in a host organism (see [Complementation](#)), which allows direct selection of the desired clone (11). A powerful example of direct selection is cloning a **transposon-tagged** gene where selection is for [antibiotic resistance](#) encoded in the [transposon](#).

Another approach is to screen transformant colonies for the gene product, either with specific [antibodies](#) (pp. 12.11–12.29 of Ref. 4) or by an enzymatic assay (12). This requires a sensitive, specific assay because foreign genes are often expressed at a low level, and enzymes from the host organism must not interfere with the assay. In addition, the gene to be cloned needs to be expressed in the host organism.

The most general screening approach is to transfer the colonies to a filter, lyse the cells, **denature** the DNA on the filter and **hybridize** a labeled nucleic acid probe to the filter (see [Southern Blots \(DNA Blots\)](#)). The probe pairs with complementary sequences in the DNA on the filter, and the colony to which the probe hybridizes is visualized by the label. There are various methods to design or obtain an oligonucleotide probe that specifically hybridizes to the desired gene.

If one wants to clone the gene coding for a known protein, it is possible to sequence the protein, usually at its N-terminus (see [Edman Degradation](#)), and to use [reverse translation](#) to determine the probe sequence. However, sequences from internal peptides produced by specific cleavage of the polypeptide chain (eg, by an **arginine**-specific [proteinase](#) or by **cyanogen bromide** cleavage at **methionine residues**) are also used if the N-terminus is blocked or if the N-terminal sequence is unsuitable for reverse translation (pp. 11.44–11.57 of Ref. 4).

Another method is using a segment of the DNA from the gene of a closely related organism as the hybridization probe (13). A modification of this approach is looking at **homologous** sequences from a set of organisms to identify highly conserved regions that are used to design an oligonucleotide probe specific for the desired gene. It is important to choose organisms with DNA base compositions close to that of the organism from which the gene is to be cloned.

With the large amount of genomic and cDNA sequence information (see [Expressed Sequence Tag](#)) becoming available in computer [databases](#), genes of interest are frequently identified during database searches. Hybridization probes are designed from the sequence information and used to screen a library or to amplify the gene directly.

An ingenious way of selecting for a full-length clone in a phage vector, when a partial clone is available, is inserting the sequence next to the *supF* gene in a plasmid, infect *recA*<sup>+</sup> cells carrying the plasmid with the phage library, and plate the resulting phage on an *E. coli* strain, where only phage containing *supF* grows (14). Recombination between the complementary sequences introduces the *supF* only into the desired phage.

#### 4. Positional Cloning of a Gene Identified by its Mutant Phenotype

Often genes are identified because they have an interesting **phenotype** when mutated. The mutation is mapped by using DNA-based markers. The most closely linked marker is used as a hybridization probe to screen a YAC, BAC, or cosmid library, thus initiating a **chromosome walk** to the gene (15). Once a library clone is identified that is deemed to contain the gene of interest, it is subcloned into a vector used to transform the mutant. The subclone containing the gene is identified by its ability to complement the mutant phenotype.

##### 4.1. Potential Problems with Library-Based Cloning

One problem with direct selection is potential contamination of the library with host DNA, which leads to cloning a host gene, rather than the gene from the original organism. Therefore it is necessary to show that genes isolated by complementation are from the desired organism, which can be done by a Southern blot. Other potential cloning problems are that some DNA fragments cannot be cloned because they contain a [poison sequence](#). DNA fragments containing **repeated DNA** sequences are often unstable, because they recombine easily. So *recA* mutant strains of *E. coli* are often used as a host for cloning to minimize [recombination](#) (16).

A major problem in cloning a DNA fragment that has identical ends (cohesive or blunt) is



recircularization of vector molecules that do not contain an insert because the two ends of the cut vector are also identical. This gives a high background of transformants or plaques lacking an insert. One method of minimizing this problem is treating the cut vector with calf [alkaline phosphatase](#) to remove its 5' phosphate groups (pp. 1.60–1.62 of Ref. 4). Dephosphorylated vector molecules cannot circularize, but they ligate to the 5'-phosphate groups on the DNA fragments to be cloned. The resulting recombinant molecule circularizes, producing a circular molecule with a single-stranded nick in each strand. These molecules give transformants. It is necessary to completely inactivate the alkaline phosphatase before the DNA to be cloned is added to the vector. The dephosphorylation reaction must be run under conditions that remove most 5'-phosphate groups but do not inactivate the plasmid, so the reaction has to be monitored carefully.

Another method of dealing with the problem of vector molecules lacking inserts is blue-white screening. Many *E. coli* vectors contain an N-terminal portion of the **beta-galactosidase** gene that codes for the a-fragment that complements the inactive b-fragment. A polylinker is introduced into the b-galactosidase gene so that the a-fragment is inactivated when a DNA molecule is ligated into the polylinker. The ligation mixture is transformed into an *E. coli* host that produces the b-fragment, and the transformants are plated on selective plates. The colonies that contain an insert are white, whereas colonies without an insert are blue (pp. 1.85–1.86 of Ref. 4).

When a probe that hybridizes to the gene to be cloned is available, the cloning process is simplified by using the probe to identify an appropriately sized restriction fragment that contains the gene by using Southern blots run on digests of a number of restriction enzymes. The appropriately sized fragments are eluted from a preparative agarose [gel electrophoresis](#) run on genomic DNA digested with the chosen restriction enzyme. The fragments are ligated into a vector cut with a compatible restriction enzyme, and the resulting transformants are screened with the probe. This process significantly reduces the number of colonies or plaques that need to be screened (1 of 40 positive vs. 1 of 1,000 for a *Thermomonospora fusca* gene) (17).

## 5. PCR to Clone Genes

Whenever sequence information is available for a gene, it is possible to design **PCR** primers that amplify all or part of a gene directly from DNA isolated from the desired organism. The PCR product is used as a probe to screen for the clone, or the PCR product is cloned.

## 6. Insertional Mutagenesis

Transposons (see [Transposon Tagging](#)) and [retroviruses](#) are used as insertional mutagens to “tag” genes. If a previously identified transposon or retrovirus inserts into a gene, the sequence of the transposon or retrovirus is used as a tag to identify the host DNA flanking the insertion. Techniques, such as IPCR, plasmid rescue, and library screening (using the element as a probe) are all successful in identifying flanking DNA from a “tagged” individual.

## 7. Cloning Tissue- or Treatment-Specific Genes

In some cases, researchers do not want to identify a specific gene, but instead they are interested in a certain class of gene, for example, liver-specific genes or genes induced by a [hormone](#). There are numerous methods for identifying such specific types of gene, including [subtractive hybridization](#), differential display, database analysis, and analysis of microarrays (see [Expressed Sequence Tag](#)).

## 8. Immunochemical Methods

Another method for cloning a gene, where the protein has been identified previously, is using an [antibody](#) that recognizes the protein to immunoprecipitate **polysomes** making the protein and then to isolate its mRNA. The mRNA produces a cDNA copy that is cloned. The cDNA clone is also used

as a probe to isolate a genomic clone.

## 9. Conclusion

There are a large number of ways to clone a gene or set of genes of interest. The method used depends on what materials are available, such as protein sequence, antibodies to the protein coded for by the gene, gene sequence information, useful libraries for screening, or a “tagged” or “untagged” mutation in the gene of interest. The purpose of cloning a gene is to generate large amounts of DNA for further analysis, such as sequencing or mutational analysis. A cloned gene is also useful for producing the protein or for **antisense** studies, and in some cases it is used to disrupt the endogenous gene by homologous recombination. A cloned gene is also useful as a probe in **gene expression** studies. Cloning a gene is one of the first steps to understanding its function.

## Bibliography

1. C. Gellissen and L. P. Hollenberg (1997) *Gene* **190**, 87–97.
2. S. N. Ho, H. D. Hunt, R. M. Horton, J. K. Pullen, and L. R. Pease (1989) *Gene* **77**, 51–59.
3. R. A. Dixon, C. J. Lamb, S. Masoud, V. J. H. Sewalt, and N. L. Paiva (1996) *Gene* **179**, 61–71.
4. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning; a Laboratory Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
5. G. G. Wilson and N. E. Murray (1991) *Annu. Rev. Genet.* **25**, 585–562.
6. V. B. Rao, V. Thaker, and L. W. Black (1992) *Gene* **113**, 25–33.
7. J. C. Pierce and N. L. Sternberg (1992) *Methods Enzymology* **216**, 549–574.
8. N. Fairweather (1997) *Methods Mol. Biol.* **68**, 137–148.
9. B. A. C. Shizuya, B. Birren, U.-J. Kim, V. Mancino, T. Slepak, Y. Tarhiri, and M. Simon (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8794–8797.
10. D. T. Burke, G. F. Carle, and M. V. Olson (1987) *Science* **236**, 806–812.
11. D. Mazel, E. Loic, S. Blanchard, W. Savrin, and P. Marliere (1997) *J. Mol. Biol.* **266**, 939–949.
12. R. M. Teather and P. J. Wood (1982) *Appl. Environ. Microbiol.* **43**, 777–780.
13. J. Agnan, C. Korch, and C. Selitrennikoff (1997) *Fungal Genet. Biol.* **21**, 292–301.
14. A. J. Hanzlik, M. M. Osemlak-Hanzlik, and D. M. Kurnit (1992) *Gene* **122**, 171–174.
15. S. D. Tanksley, M. W. Ganai, and G. B. Martin (1995) *Trends Genet.* **11**, 63–68.
16. A. C. M. Radding (1982) *Annu. Rev. Genet.* **16**, 405–437.
17. S. Zhang, G. Lao, and D. B. Wilson (1995) *Biochemistry* **34**, 3386–3395.

## Suggestions for Further Reading

18. D. M. Glover (1985) *DNA Cloning; a Practical Approach*, Vol. **2**, IRL Press, Washington, D.C., describes methods for cloning into microorganisms other than *E. coli*.
19. D. M. Glover and B. D. Hanes (1995) *DNA Cloning; a Practical Approach*, IRL Press, New York, describes methods for introducing genes into mammalian cells and mammals.
20. D. W. S. Wong (1997) *The ABCs of Gene Cloning*, Chapman and Hall, New York; an overview of cloning.

## ClpAP and ClpXP Proteinases

These large ATP-dependent [proteases](#) are involved in **protein degradation** in prokaryotes. Each enzyme is composed of two subcomplexes, both of which are essential for ATP-dependent **proteolysis**. (A similar enzyme system is also found in chloroplasts of plant.)

## 1. ClpP

The ClpP proteolytic component is a [serine proteinase](#) composed of 14 identical subunits, organized in two seven-membered rings. The rings enclose a central chamber where the [active sites](#) are localized. To be degraded, proteins must enter this chamber through openings in the rings. By itself, ClpP can hydrolyze small [peptides](#) but not polypeptides. To degrade proteins, this proteolytic component must associate with either of two homologous ATPase complexes (ClpA or ClpX), which determine the substrate specificity.

## 2. ClpA and ClpX

ClpA and ClpX complexes are hexameric rings that associate with each end of the ClpP to form a four-ring active enzyme. These ATPases bind different types of polypeptides, probably unfold them, and facilitate their entry into ClpP for degradation.

### Suggestion for Further Reading

S. Gottesman, M. R. Maurizi, and S. Wickner (1997) Regulatory subunits of energy-dependent proteases. *Cell* **91**, 435–438.

## Coding Strand

This is an ambiguous term, and perhaps best avoided. Because **DNA** is double-stranded and only one strand of the DNA of a **gene** is usually **transcribed** to make **RNA**, there is a distinction between the transcribed and the nontranscribed strand at any point in the [genome](#). Using the useful terminology “upstream”/“downstream” to refer to the direction in which the gene is transcribed and translated, the nontranscribed DNA strand runs upstream to downstream in the chemical direction 5' to 3', whereas the complementary transcribed strand has the opposite 3' to 5' polarity. The transcript RNA runs 5' to 3', like the nontranscribed DNA strand. Thus, the nontranscribed DNA strand matches the mRNA in its base sequence (with thymine replacing uracil), not the transcribed strand which has the back-to-front complements of the codons, as explained in Figure 1. Consequently, in predicting encoded [amino acid](#) sequence of the encoded [protein](#) from a gene DNA sequence, one uses the sequence of the nontranscribed DNA strand because this is the strand from which the codons can be read directly. For this reason, the nontranscribed strand is sometimes called the coding strand.

Some workers, however, adopt the opposite convention because, after all, the transcribed strand provides the information for the codon sequence of the mRNA. Therefore, it is best to eschew the coding versus noncoding terminology and to distinguish the two DNA strands as transcribed versus nontranscribed.

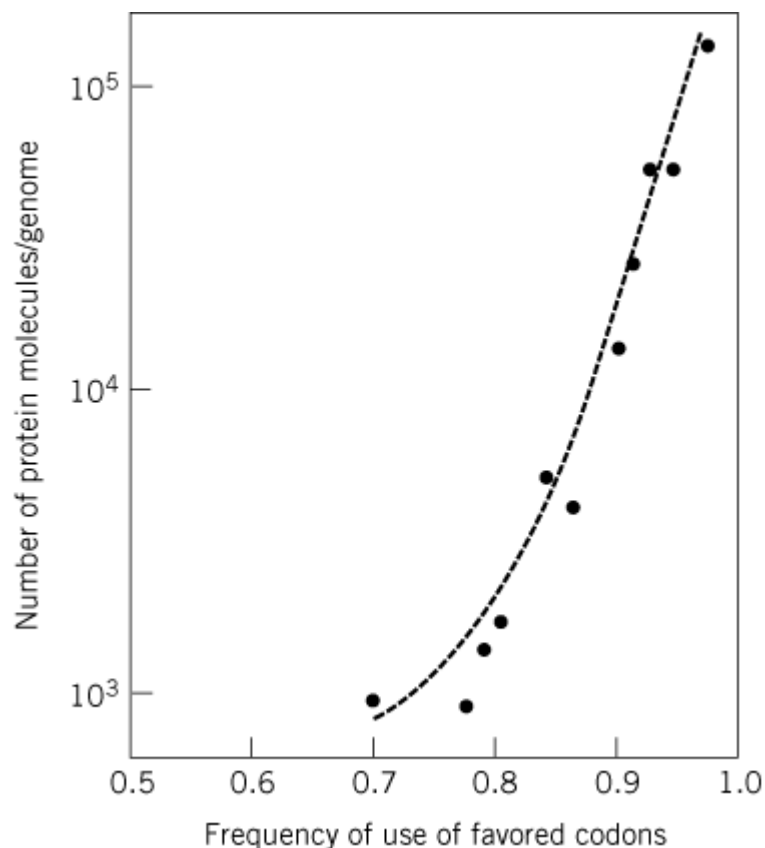
It should be remembered that both strands of a particular segment of DNA can be transcribed, though probably never in the same cell type and at the same time. This may be the case, for example, when genes are “nested” within the introns of other genes (see [Gene Structure](#)).

## Codon Usage and Bias

The [genetic code](#) is degenerate, in the sense that (except for methionine and tryptophan) a given [amino acid](#) can be specified by either two, three, four, or six alternative **codons**. When the same amino acid is specified by three, four, or six codons, they cannot all be recognized by the same [transfer RNA](#) anticodon, and so in these cases at least two or three distinct transfer RNAs must carry the same amino acid (isoaccepting) but with different anticodons (see [Transfer RNA](#) and [Aminoacyl tRNA Synthetases](#)).

Not all codons for a particular amino acid are used equally frequently. Generally, strong biases favor some codons and are against others, and different codons are favored in different organisms (Table 1). Some genes within an organism show much greater codon usage bias than others, and there is a strong correlation between the degree of bias shown by a gene and its level of expression in terms of the amount of encoded protein produced. For example, in all organisms where the matter has been examined, the genes that encode ribosomal proteins, which are always relatively abundant, are among the most highly biased. The excellent correlation in *Escherichia coli* between codon bias and level of protein product is shown in Fig. 1. This suggests that the preferred codons are more efficiently translated than the relatively rarely used codons.

**Figure 1.** Frequency of using favored (“optimal”) codons in *E. coli* genes plotted against the quantity of the gene protein products. From Ref. 1 by permission.



**Table 1. Codon Bias in Different Organisms<sup>a</sup>**

| Amino Acid    | Possible Codons            | Favored [and Disfavored] Codons |                      |                         |
|---------------|----------------------------|---------------------------------|----------------------|-------------------------|
|               |                            | <i>Drosophila</i>               | <i>Saccharomyces</i> | <i>Escherichia coli</i> |
| Leucine       | CUU CUC CUA<br>CUG UUA UUG | CUG                             | UUG                  | <b>CUG</b>              |
| Arginine      | CGU CGC CGA<br>CGG AGA AGG | CGU CGC                         | AGA                  | CGU CGC                 |
| Proline       | CCU CCC CCA<br>CCG         | CCC                             | CCA                  | CCG                     |
| Glutamine     | CAA CAG                    | <b>CAG</b>                      | <b>CAA</b>           | CAG                     |
| Lysine        | AAA AAG                    | <b>AAG</b>                      |                      |                         |
| Alanine       | GCU GCC GCA<br>GCG         | GCC                             | GCU,GCC              |                         |
| Valine        | GUU GUC GUA<br>GUG         | <i>[GUA]</i>                    | GUU GUC              | <i>[GUC]</i>            |
| Glycine       | GGU GGC GGA<br>GGG         | <i>[GGG]</i>                    | <b>GGU</b>           | GGU GGC                 |
| Serine        | UCU UCC UCA<br>UCG AGU AGC | <i>[AGU]</i>                    | UCU UCC              | UCU UCC<br>AGC          |
| Threonine     | ACU ACC ACA<br>ACG         | ACC                             | ACU ACC              | ACC,ACU                 |
| Isoleucine    | AUU AUC AUA                | AUC                             | AUU AUC              | AUC                     |
| Asparagine    | AAU AAC                    | <b>AAC</b>                      | AAC                  | <b>AAC</b>              |
| Phenylalanine | UUU UUC                    | <b>UUC</b>                      | UUC                  | UUC                     |
| Tyrosine      | UAU UAC                    | UAC                             | UAC                  | UAC                     |
| Glutamic acid | GAA GAG                    | <b>GAG</b>                      | GAA                  | GAA                     |
| Cysteine      | UGU UGC                    | <b>UGC</b>                      | UGU                  |                         |
| Histidine     | CAU CAC                    | CAC                             |                      | CAC                     |
| Aspartic acid | GAU GAC                    |                                 |                      |                         |

<sup>a</sup> *Drosophila* data from a pool of 15 high-expression genes (3); *E. coli* and *Saccharomyces* from pools of highly expressed genes (1, 2). The codons indicated show strong majority usage (heavy type indicates 90%) except those italicized in brackets where the bias is strongly against. Blanks indicate no strong bias or insufficient data.

Different translation efficiencies of different codons is largely accounted for by a number of factors, pointed out by Ikemura (“Ikemura’s rules”), the first of which is that when different isoaccepting tRNAs are present at very different concentrations, the favored codons will be those that are serviced by the most abundant tRNAs. There is evidence for this generalization from *E. coli*, *Saccharomyces* (1, 2) and *Drosophila* (3). Where a particular amino acid has both codon bias and unequal concentrations of isoaccepting tRNAs, the rule holds virtually without exception. For example, the

favored arginine codons in *E. coli* are CGU and CGC, whereas AGA is strongly favored in *Saccharomyces*. In the former organism, the tRNA that has anticodon GCI (I for inosine) recognizing CUC/U is present at several-fold higher concentration than that recognizing AGA, with anticodon UCU\* (U\* is a modified uridine), whereas in yeast the tRNA inequality is very much the opposite.

The different relative abundances of tRNAs do not, however, explain codon usage bias when it operates between codons served by the same isoaccepting tRNA. Here the bias can be explained in terms of Ikemura's other rules, which account for the fact that an anticodon may pair more comfortably with one codon than with another, even though it recognizes both. Thus an anticodon that has a thiolated uridine or 5-carboxymethyl uridine in the 5' position (the "wobble" pairing position, see [Anticodon](#)) prefers a 3' A in the codon over 3' G (rule 2); 5' inosine (I) prefers U or C to A, even though it pairs with all three (rule 3); and codon-anticodon interactions dependent on the relatively weak A–U or T–A pairings in both of the first two positions work best with a stronger G–C or C–G pairing in the third (rule 4). Examples from *E. coli* consistent with these rules are the biases in favor of GAA over GAG (rule 2), CGU or CGC over CGA (rule 3) and AUC over AUU (rule 4) (1).

### 1. Effects of Altering Codons

If the more rarely used codons really are translated relatively slowly, one would expect that their introduction into normally highly expressed genes would decrease protein yields. This prediction has been tested experimentally, and substantial effects have indeed been found, especially when multiple and contiguous rare codons have been introduced into the upstream ends of open reading frames. Hoekema et al. (4) carried out a total replacement of the normal codons by synonymous infrequently used codons in the upstream third of the normally highly expressed *Saccharomyces PGKI* (phosphoglycerate kinase) gene and found a tenfold reduction in enzyme yield. There was also a threefold reduction in the level of **messenger mRNA**. In the gene encoding glutamate dehydrogenase, in *Neurospora crassa*, a double frameshift mutant that generated three successive rare codons, numbers 54 to 56 of a 453-codon sequence, brought about a threefold reduction in enzyme without appreciable effect on mRNA level (5). And in *E. coli*, the insertion of an additional mRNA segment containing five contiguous rare AGG arginine codons 24 bases downstream of the initiation codon of *lacZ* (**b-galactosidase**) reduced the enzyme yield by 90% in the early exponential phase and virtually to zero thereafter. Another insertion, identical except that AGG was replaced by the favoured CGU codon (6), had little effect. However, when the distance between the inserted AGG-containing sequence and the initiation codon was increased by additional insertions, the yield of b-galactosidase increased almost linearly with distance, as if the rare codons have maximum effect when translation is just starting and progressively less effect as it proceeds. The mechanism of such distance-dependence is unclear. Chen and Inouye (6), reviewing the positions of rare codons in *E. coli* genes, noted that they fall within the first 25 codons after the initiating AUG and suggested that their function is to modulate the rate of translation. This would only be a regulating mechanism if there were some means of increasing or decreasing the effect, for instance, by adjusting the levels of the least abundant tRNAs. But there is as yet no evidence for this.

If rare codons have strong negative effects on the protein yields of genes, they may be an obstacle to the expression of **transgenes** in alien host cells. A codon that is abundant in the donor organism may be rare in the recipient. Therefore, to obtain a good yield of a foreign protein, it may be necessary to adjust the coding sequence of the introduced gene by *in vitro* DNA manipulation. Perhaps the best example of successful application of this idea is provided by the work of Cormack *et al.* (7), who wished to confer on the fungus *Candida albicans* the ability to produce the [green fluorescent protein](#) (GFP) native to the jellyfish *Aequoria victoria*. Introduction of the natural GFP gene into *Candida* on a replicating plasmid **vector** resulted in the formation of some GFP mRNA in the fungus, but no detectable GFP. Replacing the single CUG codon with CUU, in case the *Candida* coding idiosyncrasy of translating CUG as serine rather than leucine (8) was responsible for the failure,

brought no improvement. But when all the *Aequoria* codons that are comparatively rare in *Candida* were replaced by codons that are common in the fungus, introduction of the transgene produced an abundance of GF protein. The level of GFP mRNA was also considerably increased, which suggests, as does the yeast *PGK1* example mentioned before (4), that blocked translation, and hence failure to cover the mRNA with ribosomes, may expose the mRNA to enzymic degradation.

## Bibliography

1. T. Ikemura (1985) *Mol. Biol. Evol.* **2**, 13–34.
2. P. M. Sharp, T. M. F. Tuohy, and K. R. Mosurski (1986) *Nucleic Acids Res.* **14**, 5125–5143.
3. D. C. Shields, P. M. Sharp, D. G. Higgins, and F. Wright (1988) *Mol. Biol. Evol.* **5**, 704–716.
4. A. Hoekema, R. A. Kastelein, M. Vasser, and H. A. de Boer (1987) *Mol. Cell. Biol.* **7**, 2914–2924.
5. J. A. Kinnaird and J. R. S. Fincham (1991) *J. Mol. Biol.* **221**, 733–736.
6. G.-F. T. Chen and M. Inouye (1990) *Nucleic Acids Res.* **18**, 1465–1473.
7. B. P. Cormack et al. (1997) *Microbiology* **143**, 303–311.
8. M. A. S. Santos and M. F. Tuite (1995) *Nucleic Acids Res.* **23**, 1481–1486.

## Coenzyme, Cofactor

The term cofactor has been used as a general term to indicate that a compound was required, in addition to the [enzyme](#), for a reaction to proceed and that the compound remained unchanged at the end of the reaction. Cofactors were considered to be either activators, such as metal ions, or *coenzymes*, organic molecules that participated in enzymic reactions. **NAD** and **NADP** were named originally as Coenzyme I (CoI) and Coenzyme II (CoII), respectively, because in studies on metabolic reactions they appeared to function as electron carriers. They certainly perform that role in coupled reactions, such as those catalyzed by glyceraldehyde-3-phosphate dehydrogenase and [lactate dehydrogenase](#):



with the sum of reactions (1) and (2) being



However, this is not the role of **NAD/NADH** when functioning as substrates for a dehydrogenase reaction. Similar comments can be made about Coenzyme A (CoA), which functions as an acyl carrier, and Coenzyme Q (CoQ, ubiquinone), which acts as an electron carrier. In this connection, it is of interest that nucleoside phosphates, such as **ATP**, are not usually classified as carriers of phosphoryl groups.

Other compounds that have been considered as coenzymes are [biotin](#), flavin adenine di- and mono-nucleotide (**FAD**, **FMN**), [lipoic acid](#), [pyridoxal phosphate](#), and thiamin pyrophosphate. Biotin, lipoic acid, and pyridoxal phosphate are bound covalently to carboxylases, acyl transferases, and aminotransferases, respectively. FAD may, or may not, be bound covalently to oxidases. The interaction of thiamin pyrophosphate with decarboxylases is noncovalent, although the binding is tight.

It would seem preferable to consider a compound as a coenzyme only under conditions where it is functioning as a carrier and not when it is simply a substrate for a single enzyme. Compounds that are covalently or tightly bound to an enzyme might better be regarded as **prosthetic groups**. The features that these compounds have in common are the helping groups that they provide to enzymes. These groups possess chemical attributes that the enzymes do not have and that are essential for catalysis.

## Cohesive, Sticky Ends

Cohesive ends are short single-stranded sequences of **DNA**, usually 1 to 3 bases long, that are produced at the end of a double-stranded DNA molecule by the action of a type II [restriction enzyme](#) that makes [staggered cuts](#) in a symmetrical, **palindromic** sequence (1). For example, *Pst*1 cuts the sequence 5' ..CTGCAG.. to give the ends ..C and TGCAG.. The two ends 3' ..GACGTC.. ..GACGT.. C.. produced in this way are complementary, so they bind together by base pairing and allow [DNA Ligase](#) to repair the break. Furthermore, any fragment cut by an enzyme of this type binds to any other DNA fragment cut by the same enzyme or to one that creates the same cohesive ends. When the cohesive ends ligated together are cut by enzymes that have different recognition sequences, the new sequence is not cut by either enzyme. In the case of *Bam*HI, 5' ..GGATCC.. and *Sau*3a, ..GATC.., which produce the same sticky ends ..CCTAGG.. ..CTAG..but have different length recognition sites, the sequence created can always be cut by *Sau*3a but can be cut only by *Bam*HI if the sequence cut by *Sau*3a is a *Bam*HI site. Some restriction enzymes recognize degenerate sequences, so that only when both DNA fragments cut with this type of enzyme have the same sequence do they pair and allow ligation. Depending on the position of the cuts in the two strands, either a 5'-strand or the 3'-strand, creating a 5'-overhang or a 3'-overhang is possible, but each enzyme produces only one type.

Cohesive ends are also present on some **bacteriophages**. For example, on [lambda phage](#) they are 12 bases long and are produced by staggered cutting of a symmetrical sequence site (cos site) by a phage [enzyme](#) during DNA packaging (2).

## Bibliography

1. G. G. Wilson and N. E. Murray (1991) *Ann. Rev. Genet.* **25**, 585–562.
2. O. Yang, A. Hanagan, and C. E. Catalano (1997) *Biochem.* **36**, 2744–2752.

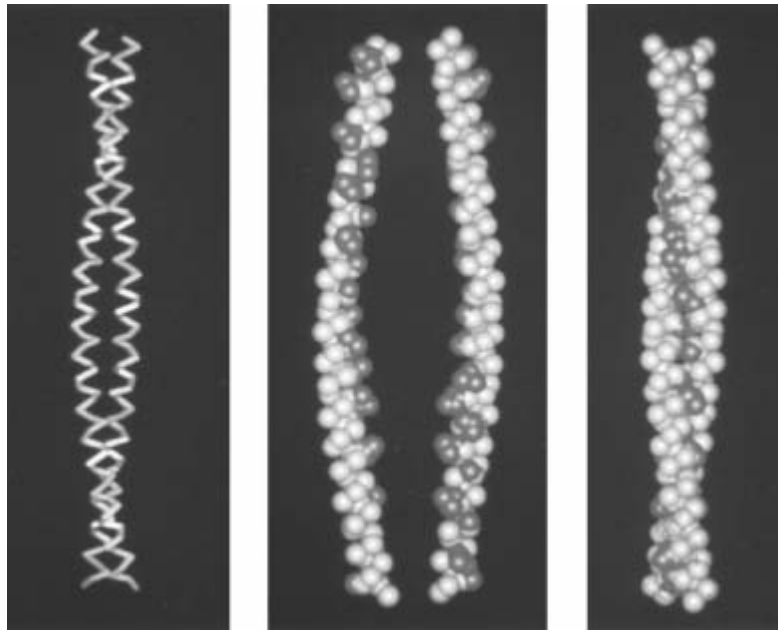
## Coiled-Coils



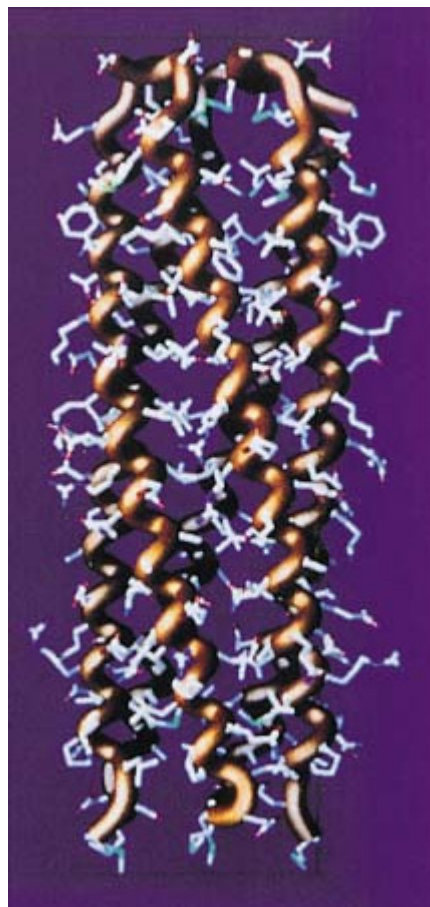
A coiled-coil is a generic name for any helical structure that has an axis that is itself helical. A general feature of all such coiled-coil conformations is that the handedness of coiling alternates at successive levels of structure, in line with the well-established practice employed by ropemakers over the centuries. This maximizes the interactions between strands and minimizes relative slippage. There are two particularly well-known classes of coiled-coil structures in biological structures: (i) the  $\alpha$ -fibrous proteins and (ii) the [collagens](#). The structure of each will be described.

$\alpha$ -Fibrous proteins display a characteristic **heptad** quasi-repeat of the form  $(a-b-c-d-e-f-g)_n$ , where about 75% of the  $a$  and  $d$  positions are occupied by **nonpolar** residues such as [leucine](#), [isoleucine](#), and [valine](#) (1, 2). These sequences adopt a right-handed  **$\alpha$ -helical** conformation with about 3.6 residues per turn. Because the apolar residues in the heptad repeat are 3.5 residues apart, on average, it follows that they will form an apolar stripe on the surface of the  $\alpha$ -helix, and this will wind around the helix in a left-handed manner (3). In an aqueous environment, apolar residues tend to pack as closely together as possible to shield each other from the [water](#) and, in doing so, provide the **hydrophobic** driving force for assembly. This can be facilitated in this case by two or more  $\alpha$ -helices coming together, optimizing the packing of the apolar residues along their interface (so-called knobs-into-holes packing) and winding around one another to generate a left-handed coiled-coil structure (Fig. 1). Favorable [electrostatic interactions](#) can also be made between the chains, which help to specify both the relative chain direction and the axial displacement between the chains (1). The interchain ionic interactions occur predominantly between oppositely charged residues in positions  $e$  and  $g$  of different chains. Calculations and experimental observations on two-stranded  $\alpha$ -helical coiled-coils have shown that the strands are parallel (rather than antiparallel) and in axial alignment. The types of apolar residues in positions  $a$  and  $d$  are important in specifying the number of chains in the coiled-coil molecule (see [Leucine Zippers](#)). Two-stranded coiled-coils occur in muscle **myosin**, [intermediate filaments](#), plectin, streptococcal M proteins, centrophilin, [kinesin](#), b-giardin, and a host of other proteins. Three-stranded structures are found in the [laminins](#), [fibrinogen](#), the [spectrin](#) superfamily of proteins, bacteriophage leg proteins such as gp17, cartilage matrix protein, mannose-binding protein, [macrophage](#) scavenger receptor protein, and many others. Four-stranded ropes are found in the silks of the bees, wasps, and ants (*Hymenoptera aculeata*), as well as in globular proteins (such as the **four-helix motif**). Five-stranded coiled-coils have recently been described in [HIV](#) capsid protein and in collagen oligomeric matrix protein (COMP) (Fig. 2). As Cohen and Parry (3) have pointed out, it is not the bending of the axis of the  $\alpha$ -helix in the supercoiled conformation that is crucial but rather the way in which it permits systematic apolar side-chain interactions to be made. It must also be pointed out that many sequences show discontinuities in their heptad substructures. A recent study by Brown et al. (4) has shown that all such (short) discontinuities can be classified as either “stutters” or “stammers”; these correspond to deletions of three and four residues, respectively, from an otherwise continuous repeat. Physically, a stutter results in a region in which the coiled-coil undergoes a degree of local underwinding to generate a longer supercoil pitch length, whereas the stammer causes a degree of local overwinding, thus giving rise to a shorter supercoil pitch length. The latter is likely to be stereochemically more difficult to achieve. The former has been observed experimentally. A coiled-coil is not confined to fibrous proteins with long rod-like domains as seen in myosin, intermediate filaments, and desmoplakin (for example), but occurs commonly in globular proteins in the form of  $\alpha$ -helical bundles. These are generally short in length (say 3 to 10 heptads) and contain anything from two or three  $\alpha$ -helices to sizable bundles containing six or even more. The same underlying heptad repeat is present, but it becomes less easy to recognize as the helix length decreases.

**Figure 1.** The structures of two-stranded coiled-coils. (Left) The backbone model of the two-stranded coiled-coil from a portion of the *N*-terminal end of **tropomyosin**. Two right-handed  $\alpha$ -helices coil around each other in a left-handed manner. (Center) A space-filling model of the same structure with the strands shown separated. The apolar residues are shown in black. (Right) The two strands brought together. The apolar residues are interlocked in a systematic way along the axis of the coiled-coil and are shielded from water as a result. (From Ref. 9 with permission.)



**Figure 2.** The five-stranded coiled-coil oligomerization domain of cartilage oligomeric matrix protein (COMP). Each of the  $\alpha$ -helices forms about one-third of a complete turn over the length of the structure. The N-terminal end is at the bottom of the page. (Courtesy of R. A. Kammerer.) See color insert.



The  $\alpha$ -chains in collagen contain a triplet repeat of the form  $(\text{Gly-X-Y})_n$ , where X and Y can be almost any amino acid residue but are commonly proline and hydroxyproline, respectively. The chain folds up into a left-handed helical structure that is very similar to that seen in polyglycine II and polyproline II. Three  $\alpha$ -chains then aggregate in parallel with a one-residue relative axial displacement between chains to generate a right-handed, triple-helical structure in which the glycine residues lie along the axis of the molecule (see [Collagen](#)). This class of conformation was originally formulated by Ramachandran and Kartha (5) and Rich and Crick (6) and has subsequently been refined by Fraser et al. (7, 8). The structure has 10 repeating units in three turns, a unit rise of 0.2894 nm, and a pitch length of 0.9647 nm. The individual  $\alpha$ -chains have a pitch length of 8.68 nm. The conformation is stabilized by a single hydrogen bond per triplet between the peptide NH group of a glycyl residue and the peptide carbonyl group of an X residue in another chain. Furthermore, water molecules bridge the glycyl NH to the C=O of a prolyl residue in another chain and the OH group of a hydroxyprolyl residue in the same chain. The diameter of a collagen molecule is about 1.4 nm. Its length varies from one type of collagen to another, but for a Type I collagen molecule its length is close to 300 nm. The molecule thus has a high axial ratio in excess of 200.

### Bibliography

1. J. F. Conway and D. A. D. Parry (1990) Structural features in the heptad substructure and longer range repeats of two-stranded  $\alpha$ -fibrous proteins. *Int. J. Biol. Macromol.* **12**, 328–334.
2. J. F. Conway and D. A. D. Parry, Three-stranded  $\alpha$ -fibrous proteins: the heptad repeat and its implications for structure. *Int. J. Biol. Macromol.* **13**, 14–16.
3. C. Cohen and D. A. D. Parry (1986)  $\alpha$ -Helical coiled-coils—a widespread motif in proteins. *Trends Biochem. Sci.* **11**, 245–248.
4. J. H. Brown, C. Cohen, and D. A. D. Parry (1996) Heptad breaks in  $\alpha$ -helical coiled coils: stutters and stammers. *Proteins Struct. Funct. Genet.* **26**, 134–145.
5. G. N. Ramachandran and G. Kartha (1955) Structure of collagen. *Nature (London)* **176**, 593–595.
6. A. Rich and F. H. C. Crick (1961) The molecular structure of collagen. *J. Mol. Biol.* **3**, 483–506.
7. R. D. B. Fraser, T. P. MacRae, and E. Suzuki (1979) Chain conformation in the collagen molecule. *J. Mol. Biol.* **129**, 463–481.
8. R. D. B. Fraser, T. P. MacRae, A. Miller, and E. Suzuki, (1983) Molecular conformation and packing in collagen fibrils. *J. Mol. Biol.* **167**, 497–521.
9. C. Cohen and D. A. D. Parry (1990)  $\alpha$ -Helical coiled-coils and bundles: how to design an  $\alpha$ -helical bundle. *Proteins Struct. Funct. Genet.* **7**, 1–15.

### Suggestions for Further Reading

10. C. Cohen and D. A. D. Parry,  $\alpha$ -Helical coiled-coils and bundles: how to design an  $\alpha$ -helical bundle. (1990) *Proteins Struct. Funct. Genet.* **7**, 1–15.
11. R. D. B. Fraser and T. P. MacRae (1973) *Conformation in Fibrous Proteins and Related Synthetic Polypeptides*, Academic Press, London.
12. A. Lupas (1996) Coiled-coils: new structures and new functions. *Trends Biochem. Sci.* **21**, 375–382.

### Cointegrative Vectors

Cointegrative vectors were the first type of **vector** developed for transferring foreign **DNA** from the bacterium [Agrobacterium](#) to **plant** cells, for use in [plant genetic engineering](#) (1). Although cointegrative vectors may be relatively difficult to work with in practice, compared with the alternative binary vectors, they offer **plasmid** stability in the **bacterial** cell.

The general concept of using cointegrative vectors relies on homologous [recombination](#) within *Agrobacterium* to introduce into a modified T-DNA the DNA to be transferred to the plant cell. First, the DNA to be transferred is introduced into an intermediate [cloning](#) vector based on an *Escherichia coli* plasmid. The intermediate vector that contains the foreign DNA is introduced by conjugation into *Agrobacterium* that contains the cointegrative vector. The cointegrative vector is a [Ti plasmid](#) from which the T-DNA genes that encode oncogenic function have been removed and/or replaced with a sequence that is also contained in the intermediate vector. The intermediate vector is not stable in *Agrobacterium*. Homologous recombination between the intermediate vector and the cointegrative vector results in transferring the foreign DNA to the cointegrative vector. These vectors are designed so that once the foreign DNA, is integrated into the cointegrative vector, it is located between its border sequences. The most generally used cointegrative vector system is the “split-end vector system” (2)

#### Bibliography

1. P. Zambryski et al. (1983) EMBO J. **2**, 2143–2150.
2. R. Fraley et al. (1985) Bio/Technology **3**, 629–635.

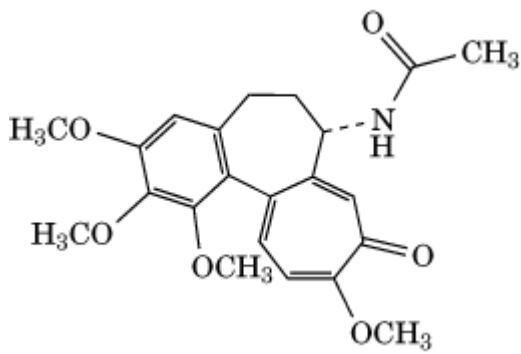
#### Suggestions for Further Reading

3. R. Walden (1988) *Genetic Transformation in Plants*, Prentice-Hall, Englewood Cliffs, NJ.
4. F. F. White (1993) "Vectors for gene transfer in higher plants". In *Transgenic Plants*, Vol. **1** (D. Kung and R. Wu, eds.), Academic Press, San Diego, pp. 15–48.

## Colchicine

Colchicine (Fig. 1), a compound obtained from the autumn crocus, *Colchicum autumnale*, is an ancient drug that has been used for centuries for the treatment of gout (1). Its use in molecular biology began in the mid-1930s when its potent ability to inhibit eukaryotic cell proliferation at mitosis was discovered. Since then, colchicine has had a remarkable history as an experimental tool for characterizing the biochemical properties of [tubulin](#), the protein subunit of [microtubules](#), for characterizing the diverse processes in eukaryotic cells that are dependent upon drug-sensitive microtubules, and for studying the polymerization and dynamics properties of microtubules. For example, colchicine has been used to determine the role of microtubules in mitosis, protein **secretion**, axonal transport, and the development and maintenance of asymmetric cell shape. It has also been used extensively as a cytogenetic tool to determine [chromosome](#) numbers in [karyotypes](#). Also, radiolabeled colchicine was used as an affinity marker to effect the first purification of tubulin from brain (2).

**Figure 1.** Structure of colchicine.



## 1. Colchicine Binding to Tubulin

The binding of colchicine to tubulin is not a simple process (3-6). The binding reaction is slow and highly temperature-dependent. Binding is extremely slow at 0°C, and 2 to 3 h is required to reach equilibrium at 37°C. The [activation energy](#) for the forward reaction is high, ~20 kcal/mol. The kinetics of binding are biphasic. The initial step is the formation of a reversible preequilibrium complex that is followed by a slow step in which conformational changes in tubulin lead to the formation of a final-state tubulin-colchicine (TC) complex. Although the binding of colchicine to tubulin is noncovalent, the final state TC complex is very poorly reversible.

Colchicine binds to tubulin at a single site with a dissociation constant in the range of 0.02 μM to 5 μM. Colchicine has very different affinities for tubulin isolated from different sources. For example, colchicine binds strongly to vertebrate brain tubulin, but it binds very weakly to plant, fungal, and yeast tubulins. Also, there are different tubulin isotypes, and colchicine binds differently to the isotypes purified from same tubulin source (7). Despite intensive investigation, the precise location in tubulin of the colchicine-binding site is not clear. The binding site is thought to reside in the b-tubulin subunit, near residues Cys 354 and Cys 241, close to the intradimer interface.

The binding of colchicine to tubulin induces conformational changes in tubulin, as well as in colchicine (6, 8, 9). For example, the binding of colchicine to tubulin quenches the intrinsic tryptophan **fluorescence** of tubulin, indicating that it induces a small change in the [tertiary structure](#) of the protein. It also perturbs the far-ultraviolet [circular dichroism](#) spectrum of tubulin, indicating that it changes the **secondary structure** of the protein. In addition, colchicine binding to tubulin strongly increases the intrinsic [GTPase](#) activity of tubulin, increases the affinity of the ab dimer association by threefold, changes the exposure of [thiol groups](#) in tubulin, and, under certain conditions, induces tubulin to self-assemble into nonmicrotubule polymeric structures. Tubulin is also subject to a time- and temperature-dependent irreversible decay of its [protein structure](#). Binding of colchicine to tubulin slows the rate of decay. The concept that colchicine undergoes conformational changes upon binding to tubulin is evident from the development of colchicine-tubulin fluorescence and the change of the colchicine circular dichroism spectrum upon binding to the protein.

## 2. Inhibition of Microtubule Polymerization by Colchicine

Colchicine inhibits microtubule polymerization at concentrations that are far below the total concentration of tubulin (10, 11), indicating that colchicine inhibits microtubule polymerization by acting at the microtubule ends. In order to produce its potent actions on microtubule polymerization, colchicine must first form a TC complex. Substoichiometric concentrations of TC-complex only partially depolymerize microtubules, and relatively low concentrations of TC-complex can stabilize microtubules against dilution-induced disassembly (12). These studies support the hypothesis that the TC-complex forms a stabilizing “cap” at the end of the microtubule. Also, the TC-complex can

form copolymers with unliganded tubulin when microtubules are assembled in the presence of TC complex (13).

### 3. Kinetic Suppression of Microtubule Dynamics

Microtubules exhibit two kinds of nonequilibrium dynamics: [treadmilling](#) and dynamic instability (see [Microtubules](#)). Recent studies have revealed that colchicine and other compounds that depolymerize microtubules (see [Vinblastine](#)) strongly suppress these dynamics at relatively low concentrations in the absence of appreciable microtubule depolymerization. Colchicine was found some years ago to suppress treadmilling *in vitro* and the rate of microtubule disassembly upon dilution of the microtubules [now called “kinetic capping” (11, 12)]. However, its stabilizing effects on dynamics were only fully appreciated with the introduction of real-time differential-interference contrast video microscopy, which enabled one to visualize directly in real time the stabilizing action of the drug on the growing and shortening dynamics of individual microtubules (14). Small numbers of incorporated TC complexes strongly suppress the rates and extents of growing and shortening and greatly increase the percentage of time that the microtubules spend in an attenuated state. In addition, the TC complex strongly suppresses the catastrophe frequency and increases the rescue frequency. At low submicromolar concentrations, TC complex suppresses the dynamics without reducing the polymer mass. Significant reduction of polymer mass requires relatively high TC complex concentrations. However, the surviving microtubules are extremely stable. Colchicine appears to suppress microtubule dynamics by binding at the microtubule ends, most probably by inducing a conformational change and/or by steric hindrance at the ends.

#### Bibliography

1. P. Dustin (1984). *Microtubels*, Springer-Verlag, Berlin, pp. 1–482.
2. R. C. Weisenberg, G. G. Borisy, and E. W. Taylor (1968) *Biochemistry* **7**, 4466–4478.
3. G. G. Borisy and E. W. Taylor (1967) *J. Cell. Biol.* **34**, 525–533.
4. L. Wilson and M. Friedkin (1967) *Biochemistry* **6**, 3126–3135.
5. B. Bhattacharyya and J. Wolff (1974) *Proc. Natl. Acad. Sci. USA* **71**, 2627–2631.
6. D. L. Garland (1978) *Biochemistry* **17**, 4266–4272.
7. R. F. Luduena (1983) *Mol. Biol. Cell* **4**, 445–457.
8. J. M. Andreu and S. N. Timasheff (1982) *Biochemistry* **21**, 6465–6476.
9. T. David-Pfeuty, C. Simon, and D. Pantaloni (1979) *J. Biol. Chem.* **254**, 11696–11702.
10. J. B. Olmsted and G. G. Borisy (1973) *Biochemistry* **12**, 4282–4289.
11. R. L. Margolis and L. Wilson (1977) *Proc. Natl. Acad. Sci. USA* **74**, 3466–3478.
12. R. L. Margolis, C. T. Rauch, and L. Wilson (1980) *Biochemistry* **19**, 5550–5557.
13. H. Sternlicht and I. Ringel (1979) *J. Biol. Chem.* **254**, 10540–10550.
14. D. Panda, J. E. Daijo, M. A. Jordan, and L. Wilson (1995) *Biochemistry* **34**, 9921–9929.

#### Suggestions for Further Reading

15. O. J. Eigsti and P. Dustin Jr. (1955) *Colchicine, in Agriculture, Medicine, Biology, and Chemistry*, The Iowa State College Press, Ames, IA.
16. L. Wilson and K. W. Farrell (1986) *Ann. NY Acad. Sci.* **466**, 690–708.
17. S. B. Hastie (1991) *Pharmacol. Ther.* **512**, 377–401.
18. L. Wilson and M. A. Jordan (1994) in *Microtubules*, J. S. Hyams and C. W. Lloyd, eds., Wiley-Liss, New York, pp. 59–83.
19. A. Vandecandelaere, S. Martin, M. Schilstra, and P. Bayley (1994) *Biochemistry* **33**, 2792–2801.

## Cold-Sensitive Mutants

Cold-sensitive [mutants](#) are a class of **conditional lethal mutants** that cannot grow at a temperature below the organism's normal optimum. The low temperature may be lethal because the mutant form of the [protein](#) loses its function or because it forms inappropriate or inhibitory interactions with other proteins. Finding and characterizing second-site revertant mutations that **suppress** the cold-sensitive phenotype identifies [protein–protein interactions](#) (see [Temperature-Sensitive Mutation](#)).

## Colicins

Colicins produced by and active against coliform **bacteria** constitute a subset of the [bacteriocins](#) generated by many groups of bacteria. They have been a subject of interest since very early in this century. In 1925, Gratia demonstrated that *Escherichia coli* strain V (virulent in experimental infection) in liquid media produces a heat-stable substance that in high dilution inhibits the growth of *E. coli* (1). This protein was designated as colicin V. It has now been shown that it fits better with the description of the so-called “microcins” (see [Bacteriocins](#)) (2). Then a whole series of colicins produced by *E. coli* and closely related members of Enterobacteriaceae were discovered. Mainly as a result of the influence and efforts of Fredericq (3), knowledge of the colicins advanced at a great rate, and more than 20 different types were identified on the basis of their action against a set of specific resistant (generally receptor-deficient) mutants (Table 1). In contrast to the bacteriocins from **Gram-positive** bacteria that kill species other than those that are likely to have the same ecological niche, colicins are active only against *E. coli* and closely related bacteria (there is a similar relationship between the cloacins and *Enterobacter cloacae*, the klebicins and *Klebsiella* species, and the pyocins and pseudomonaceae). A characteristic feature of colicinogenic bacteria is to be specifically immune to the colicin they make but not to other colicins. A large number of colicins have now been identified, and each has been characterized by the corresponding specific immunity protein. Colicins are **plasmid**-encoded, and each plasmid also bears an immunity protein, thus ensuring that plasmid carriers are protected from the colicin they themselves produce. There is probably a selective advantage for Enterobacteriaceae to produce a colicin because 30 to 40% of natural isolates of *E. coli* carry colicinogenic plasmids (4). The colicins exert their lethal effect through a single-hit mechanism (5). There was some controversy about the meaning of this terminology, but now there is a consensus that sensitive cells are killed by a single event, implying that a single molecule kills, although not every one does.

**Table 1. Characteristics of Colicins**

| Colicin     | Colicin Activity Group | Receptor | Translocation System |
|-------------|------------------------|----------|----------------------|
| E2,E7,E8,E9 | A                      | DNase    | BtuB TolA, B, Q, R   |

|            |   |   |      |                           |
|------------|---|---|------|---------------------------|
| E3,E6,DF13 | A | RNase                                     | BtuB | TolA, B, Q, R             |
| E1         | A | Pore-forming                              | BtuB | TolC, TolA, TolQ          |
| A          | A | Pore-forming                              | BtuB | OmpF, TolA, B, Q, R       |
| N          | A | Pore-forming                              | OmpF | OmpF, TolA, Q             |
| K          | A | Pore-forming                              | Tsx  | OmpF, OmpA, TolA, B, Q, R |
| Col5       | B | Pore-forming                              | Tsx  | TolC, TonB, ExbB, D       |
| Col10      | B | Pore-forming                              | Tsx  | TolC, TonB, ExbB, D       |
| Ia,Ib      | B | Pore-forming                              | Cir  | TonB, ExbB, D             |
| D          | B | Inhibition of protein synthesis           | FepA | TonB, ExbB, D             |
| M          | B | Inhibition of synthesis of murein and LPS | FhuA | TonB, ExbB, D             |

---

Colicins are soluble proteins of 29 to 70 kDa (for colicin V, see previous). Their amino acid sequence is known, and they share the property of being linearly organized in three **domains** that have specific functions. They are also highly asymmetrical molecules and have axial ratios of 8 to 10 (6, 7). Another common feature is that their production is induced by exposure of colicinogenic bacteria to agents like UV light and genotoxic chemicals, such as [mitomycin C](#), that elicit the [SOS response](#). After induction, they are produced in large amounts, so they provide useful model systems to study fundamental biological problems, such as protein–protein interactions, polypeptide translocation and insertion across and into [membranes](#), functioning of voltage-gated pores, etc.

The classification of colicins reflects bacterial rather than colicin properties. According to its activity spectrum against a variety of mutants, a particular colicin is unambiguously assigned to one of two groups, A and B (8, 9). Group B colicins are inactive on strains that have a lesion in the *tonB* gene but are active against strains mutant in *tolA* and *tolB* genes. Group A colicins show the opposite specificity. Now we know that group A colicins (A, E1 to E9, K, L, N, and cloacin DF13) and filamentous **bacteriophage** (f1, fd, and M13) need the Tol proteins to penetrate into cells (9-11). Group B colicins (B, D, Ia, Ib, M, 5, and 10) and phages T1 and F80 need TonB and its associated proteins (8, 9). In all cases studied so far, the determinants for colicinogeny are located on plasmids (Table 1) of two different types: small multicopy plasmids or large low-copy-number plasmids that generally correspond to the A and B groups of colicins. In general, group A colicins are encoded by small plasmids and are actively released to the extracellular medium, whereas group B colicins are encoded by large plasmids and are very poorly secreted.

At least three hypotheses have been proposed to explain the evolution of colicin plasmids: (1) positive selection of diversity, (2) recombinational shuffling, and (3) transposition (12). The different colicins may have evolved by DNA recombination of fragments encoding different colicin domains (see later). This is best exemplified by the common uptake route for colicins B and D, which have a highly homologous N-terminal and central polypeptide sequence (defining the translocation and receptor domains respectively; see later) but very different C-terminal domains with different types of activities (13).

## 1. Genes for Colicins and Associated Proteins



Regardless of the plasmid type, the genes for the colicin, immunity, and lysis (when there is one) are always clustered (14). The difference between immunity to the enzymatic colicins and to the channel-forming colicins and colicin M (whose target is the cytoplasmic membrane) is reflected in the regulation of their synthesis and the arrangement of the various colicin [operons](#). In enzymatic colicins, the immunity gene is transcribed in the same direction as the colicin and lysis genes (15-17), whereas for the channel-forming colicins and for colicin M, immunity is encoded on the opposite DNA strand and is thus transcribed in the direction opposite to the colicin and lysis genes (18-20). Enzymatic colicin and cognate immunity genes form an operon under the control of an SOS **promoter** (see later). Thus, induction of colicin synthesis results in increasing the amount of immunity protein synthesized, but the presence of a terminator causes partial transcription arrest upstream of the lysis gene. In contrast, inducing the synthesis of pore-forming colicins from their SOS promoter does not result in a concomitant increase in the immunity protein. The immunity is constitutively expressed from a weak promoter (19, 21-23). The situation is similar for the murein-synthesis-inhibiting colicin M and its immunity gene (24).

Colicin production is inducible and suicidal. Under normal conditions, very few cells produce colicin because the expression is repressed by the [LexA repressor](#) and is switched on only during the SOS response normally associated with repairing damaged DNA. This means that, in culture, colicin production is inducible in a growing population of cells by [mutagens](#) (UV light, mitomycin C). This induction is very efficient. In the absence of mutagens, LexA exhibits a very tight repression (25) and ensures that only a small proportion of cells go down the dead-end road of colicin expression and export. Once they have done so, it is important to the plasmid clone as a whole that the sacrificial cells produce as much colicin as possible. This maximizes selection against non-plasmid-bearing cells. Meanwhile, the silent preservation of many identical copies of the same plasmid in immune cells ensures no significant loss of the plasmids in the producer cell.

## 2. Extracellular Release of Colicins

Most of the colicins (group A, but not group B) are actively released into the growth medium. Their release mechanism differs from that of proteins secreted by Gram-negative bacteria (see [Protein Secretion](#)):

1. they do not contain an N-terminal or C-terminal signal sequence like **Sec**-dependent exported polypeptides or those depending on ABC transporters;
2. their release depends on the function of a lytic protein (variously called *kil*, *lys*, *brp*) whose gene is part of the colicin operon, downstream of the colicin gene;
3. they are released from the host cell some time after synthesis;
4. their extracellular release is not specific, and quasi-lysis proteins have various effect on cells in addition to causing the release of assorted proteins and small molecules (26).

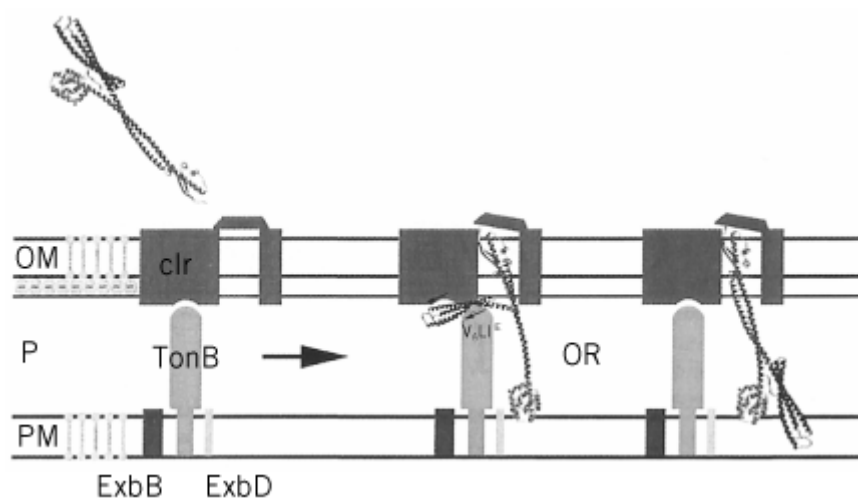
The lysis genes encode short proteins produced in a precursor form with a **signal sequence** and a typical cysteine lipid-modification consensus box, implying that the precursor undergoes the following [posttranslational modifications](#): (1) acylation of the cysteine residue, (2) cleavage of the signal sequence by signal peptidase II, and (3) fatty acylation of the [amino group](#) of the now N-terminal acylated cysteine. The maturation and processing of lysis proteins occurs slowly, so every intermediate form can be observed. Both the mature form and the signal sequence accumulate in *E. coli* after processing (27, 28). The addition of globomycin to *E. coli* cells stops processing of lysis protein precursors and inhibits colicin release (29). Acylation is essential for quasi-lysis (29, 30). The mature lysis proteins activate the normally dormant [phospholipase A](#) in the outer membrane of colicin-producing cells, thus causing production of lysophospholipids, quasi-lysis, and colicin release (31). However, lysis proteins must also affect the permeability of the inner membrane (32). The lysis protein must reach a critical concentration within the cells before quasi-lysis and colicin release, which also indicates that the lysis protein directly affects the inner membrane. Local

modifications in the structure of the bilayer are induced by interactions with lipopeptide micelles, as previously suggested in the case of Iturin A and bacillomycin L (33).

### 3. Mechanism of Entry and Domain Organization of Colicins

Colicins bind to specific receptors in the bacterial outer membrane, from where they are translocated to, and eventually through, the inner membrane to reach their targets. Consistent with these three steps, their polypeptide chains are linearly organized into three domains: the *N*-terminal domain is involved in the translocation step, the central domain is responsible for receptor binding, and the *C*-terminal domain carries the lethal activity (Fig. 1). The boundaries of these domains were defined first by limited **proteolytic** digestion (34-36), then by genetic engineering (37, 38), and from sequence homologies between colicins with the same type of lethal activity and using the same receptors (E-type, for example) or the same translocation pathways. Even colicin M, which has a molecular mass of 29 kDa, has three functional domains (24).

**Figure 1.** The mechanism of colicin Ia attachment and translocation. The outer membrane (OM) receptor for colicin Ia consists of Cir, TonB, and accessory plasma membrane (PM) proteins ExbB and ExbD. The translocation domain (blue), receptor binding domain (green), and channel-forming domain (red) separated by a pair of helices (gray), each 16 nm long, are indicated in colicin Ia. The TonB box of colicin Ia may compete with the TonB box of Cir for linkage to TonB. The channel-forming C domain reaches the plasma membrane, where it forms an ion-conducting channel by subsequent insertion and rearrangement of helices within the membrane (not shown). The translocation T domain may remain near the periplasmic surface of the outer membrane during channel activity. However, the presence of the 16 nm-long T3 helix indicates an alternative possibility in which the T domain crosses the periplasmic space (P) to participate in channel formation *in vivo*. The location of the TonB box of colicin Ia on the upper surface of the T domain of sequence Glu(E)-Ile(I)-Met(M) Ala(A)-Val(V) is indicated, reading right to left, as EIMAV. The arrows indicate the locations of the TonB box in the receptor and colicin Ia. Reproduced from (68) with permission.



To enter the cells, colicins have borrowed multiprotein systems used by sensitive cells for important biological functions. These proteins include **porins**, **vitamin B<sub>12</sub>**, siderophore, nucleoside receptors, and multiprotein systems that cooperate with these proteins. The group A and group B colicins correspond to two different pathways of entry beyond receptor binding: the TonB and the Tol pathway. The receptors for the various colicins are given in Table 1. The protein most frequently used is BtuB, the vitamin B<sub>12</sub> receptor, which defines the so-called E-type colicins (E2 to E9) (12).

High-affinity receptors or iron siderophores are also used by many colicins. The nucleoside Tsx porin and the major porin OmpF are also receptors (Table 1). Depending on the colicin class, the colicin receptors function cooperatively with either the Ton B system or the Tol system.

The *tol* locus was originally defined as comprising four loci, *tol A, B, Q, R*, on the basis of group A colicin tolerance (39). These genes are organized in two operons that comprise three additional genes called *orf1, orf2* and *pal* (40), whose gene products are not directly involved in colicin uptake. However, one of them, the Pal protein (peptidoglycan associated protein), a **lipoprotein** of the outer membrane, interacts with TolB (41). Both *tol* and *pal* mutants have an altered outer membrane (39). The only soluble Tol protein is TolB, which is principally located in the **periplasm** (42) but is linked to the outer membrane through interaction with Pal. The other Tol proteins are associated with the cytoplasmic membrane. TolQ contains three membrane-spanning segments, and TolA and TolR are anchored to this membrane by their *N*-terminal regions. All three proteins interact through membrane-spanning ***α*-helices** (43, 44). There is yet no direct evidence that TolB interacts with membrane-associated Tol proteins, but each of the Tol proteins cofractionates with a membrane fraction accounting for the presumed contact sites between the inner and outer membranes of *E. coli*. In colicin A-treated cells, the toxin was also found in this fraction, and the relative amount of Tol proteins was doubled, suggesting that the colicin itself may recruit more Tol proteins at the contact sites (45, 46).

Although both TolB and TolR are required to take up most group A colicins, neither is required for colicin E1, which in contrast requires the so-called TolC protein, a minor outer membrane protein that forms pores *in vitro* (47) and is also involved in secreting hemolysin (48) and colicin V (49). OmpF, another porin, is required for translocating many group A colicins (A, E-type, N) (50, 51). The *N*-terminal domain of group A colicins contains all the information needed for the translocation step (52) beyond receptor binding, including binding to OmpF (or TolC for colicin E1) and to Tol proteins like TolA (53) and TolB (54). There are about 1000 translocation sites for colicin A per bacterium (55). In addition, it has been shown that: (1) unfolding occurs upon receptor binding (55), (2) the *N*-terminal domain interacts with Tol proteins, (3) the polypeptide chain spans the entire cell envelope after inserting the pore-forming domain (56), and (4) a reduced constriction of the lumen of the OmpF pore prevents translocation of colicin A and N (57). A hypothetical model of colicin translocation based on these results has been presented (11, 58).

Group B colicins are imported into sensitive cells through high-affinity siderophore receptors and the Ton B pathway. The Ton B system consists of a complex of proteins TonB-ExbB-ExbD, which facilitates the flow of energy from the cytoplasmic to the outer membrane for the energy-dependent transport of ferric siderophores and of vitamin B<sub>12</sub>. ExbB and ExbD are physically and functionally homologous to TolQ and TolB (59). Ton B, like Tol A, has its *N*-terminal anchor in the cytoplasmic membrane and spans the periplasm (60). It has been proposed that an energized conformation of TonB, like FepA or FhuA, opens channels in the outer membrane, which act as receptors for the ferric siderophores. Until recently, each of the group B colicins studied uses a Ton B-dependent receptor, but it has now been reported that a new colicin (colicin 10) uses the Ton B system and binds to a Ton B-independent receptor Tsx. Colicin 10 also requires TolC for its uptake (61). Thus the interaction between the colicin receptor and Ton B is not obligatory if TolC is involved. This situation may be comparable to that of the group A colicin E1, which does not need TolB or TolR but does need TolC.

The high-affinity receptors for siderophores and vitamin B<sub>12</sub> contain a pentapeptide motif close to the *N*-terminus, called the “Ton B box” (62). Mutations in the Ton B box reduce receptor activity dramatically. They are also **suppressed** by mutations in TonB (63), indicating physical and functional interactions between TonB and the receptors. Group B colicins, such as colicins B and D, which bind to the FepA receptor, and colicin M, which binds to the FhuA receptor, also contain a TonB box sequence close to the *N*-terminus. This indicates that these colicins interact directly with TonB during translocation. Mutations in the TonB box of colicin B affect its uptake and are suppressed by secondary substitutions clustered in TonB (63). A similar situation is observed with TonB, FhuA, and colicin M. These results suggest that group B colicins need to interact with TonB to be translocated, just as group A colicins need to interact with Tol proteins. Both the *N*- and *C*-terminal domains of the protein also interact with the inner membrane and therefore are protected

from proteolytic degradation (64). This suggests that the Y-shaped structure (65) allows colicin Ia to span the entire cell envelope, as previously reported with colicin A (66). The receptor-binding domain would be the only part remaining exposed outside the cells.

A molecular model of colicin M uptake has been proposed (24). After first binding to FhuA in the outer membrane, colicin M is taken up in an energy-coupled process through the action of the TonB protein, which is anchored to the inner membrane and spans the periplasmic space. TonB binds to FhuA, inducing the release conformation of FhuA, which, in turn, results in the vectorial translocation of colicin M across the outer membrane. TonB itself has two conformations, energized and unenergized. Energization of TonB takes place in the cytoplasmic membrane. Induction of the receptor release conformation consumes energy, so TonB dissociates from the FhuA receptor and has to be reenergized to induce the next round of colicin M uptake. Because both FhuA and colicin M contain a TonB box, the TonB protein most likely interacts with both of them sequentially. The role of ExbB and D is still unclear. They are involved either in energy transduction (cycling of the TonB conformational change) from the cytoplasmic membrane to TonB or in stabilizing TonB, which is physically and functionally unstable (24, 67, 68).

The 626-residue colicin Ia is about 21 nm long and consists of three functional domains separated by a pair of  $\alpha$ -helices 16 nm long. A central domain at the bend of the hairpinlike structure mediates binding to the outer membrane receptor. A second domain mediates translocation across the outer membrane via the TonB pathway. The TonB box recognition motif of colicin Ia is on one side of three 8 nm-long helices arranged as a helical sheet. The third domain, made up of 10  $\alpha$ -helices, is the pore-forming domain (see below). The two exceptionally long 16 nm  $\alpha$ -helices enable colicin Ia to span the periplasmic space between the outer and inner membranes (Fig. 1) (68). This type of structure, which will very likely be found in each colicin, strongly suggests how these toxins exploit the machinery of the target cell to get across the periplasmic space (Fig. 1).

#### 4. Mode of Action of Colicins

The biochemical effects of colicins on sensitive bacteria fall into four classes: (1) the RNase colicins typified by colicin E3, (2) DNase colicins typified by colicins E2 and E9, (3) pore-forming colicins like colicin E1, and the still unique colicin M, which inhibits cell wall biosynthesis (24).

##### 4.1. RNase Colicins

The two most extensively studied members of this class are colicin E3 and cloacin DF13 (derived from a strain of *Enterobacter cloacae* closely related to *E. coli*). The 11 kDa and 12 kDa C-terminal segments of these bacteriocins cleave the 16S ribosomal RNA, either in isolation or in intact 30S ribosomal subunits, which explains their inhibitory effect on protein synthesis (69). Colicins E5 and E6 have the same type of activity (70). Although they are highly homologous in their RNase regions (71), colicins E3, E5, E6, and cloacin DF13, each have a specific immunity protein. No structures are yet available for any of these RNase domains.

##### 4.2. DNase Colicins

In 1970, it was shown that colicin E2 causes both single- and double-strand breaks in the DNA of sensitive cells (72), and three new types of DNase colicins were subsequently described (colicin E7, E8, and E9). The sequential identity between the various DNase domains is greater than 80% (71), yet each is inhibited by a specific immunity protein. This is similar to the situation in the RNase-type colicins. Structural and biophysical analysis of the DNase immunity protein complexes has been aided by the *E. coli* overexpression systems developed by Wallis et al. (73-75).

##### 4.3. Pore-Forming Colicins

Pore-forming colicins constitute the largest group. Colicin V must be set apart in this class because it fits the definition of microcins better than that of colicins (2). Pore-forming colicins dissipate the membrane potential (76-80), which causes a series of metabolic effects, such as inhibition of **active transport** and of protein and nucleic acid synthesis, decrease of the internal ATP concentration (see

[Adenylate Charge](#)), and leakage of potassium ions ([81-85](#)). Cell death results from the depletion of the ATP pool following efflux of phosphate through the channel and ATP hydrolysis ([86](#)). Ten pore-forming colicins have been identified thus far. Sequential homologies indicate that they should be separated into two groups: (1) E1, 5, Ia, Ib and (2) A, B, and N ([87](#)), differing mainly in a segment containing the loop between helices 8 and 9, as defined in the X-ray crystal structure of the C-domain of colicin A ([88](#)). The structures of soluble forms of the channel domains of colicins A, Ia and E1, are known to varying extents ([89-92](#)). The pore-forming domain is composed of a bundle of 10  $\alpha$ -helices arranged in three layers in a novel protein fold ([93](#)). The N-terminal layer (helices 1 and 2) is connected to another layer consisting of two pairs of [amphipathic](#) helices (helices 3 + 4 and 6 + 7), and a single helix (helix 5) connects the two helix pairs. The middle layer is composed of helices 8 and 9, which form a **hydrophobic** hairpin buried in the core of the molecule (Fig. [1](#)). The external sides of the peripheral helices are [hydrophilic](#), which explains the paradox of how the same protein domain exists either in a water-soluble form or in a membrane-inserted state (see later). Because the C-terminal domain of colicins is easily purified by proteolytic digestion of the whole colicins or by genetic engineering, many biophysical studies have been carried out to improve our understanding of this step in the action of pore-forming colicins (for a review, see [Suggestions for further reading](#)).

Schein et al. ([78](#)) were the first to demonstrate that colicin A forms well-defined channels in planar lipid bilayers that permit ions to cross the membrane ([94, 95](#)). Pore-forming colicins form channels of 10 to 30 picosiemens (pS), which correspond to  $10^7$  ions/sec/channel. These channels are characterized by their sensitivity to the electrical potential across the membrane. The gating voltage (corresponding to the activation or to the inactivation of 50% of the channels) varies for a given colicin. Values are +21, +50, and +70 mV for colicins N, A, and E1, respectively. There are multiple membranes states at the single-channel level. With colicin E1, channels of 10 pS are initially observed. With increased activation in 1 M KCl, channels of 20, 40, and 60 pS with substates appear, and flickering is observed. Similar observations with other colicins suggested the existence of many closed, open, and inactivated states. Thus, the concept of a single-membrane-inserted conformation appears to be incorrect, although a simplified model with a single or a few membrane-inserted forms has been very helpful.

The channel properties of the C-domain and of whole colicin are qualitatively similar for colicin E1, but not for colicins A ([96](#)) and Ia ([97](#)). Differences are ascribed to the influence of the other domains on the pore-forming domain in the intact molecule. The pores are permeable to cations and anions, but the rates at which they are transported are low ([98](#)). There is a relative preference for cations versus anions that is modulated by the pH and the lipid composition ([99, 100](#)). Using the relative permeabilities of large organic cations and anions and taking into account their asymmetrical shapes, the size of the channel diameter was estimated at only 0.4–0.5 nm. This implies strong interactions between permeable ions and the side chains of residues exposed in the channel lumen ([100](#)).

The transition of pore-forming colicins from a water-soluble to a membrane-inserted state involves large structural changes, which is of interest in the context of **protein folding**, protein translocation, and protein insertion into membranes. From studies with model membrane systems, a likely sequence of events for membrane insertion and pore formation has been proposed for colicins A and E1 (see [Suggestions for Further Reading](#)). Briefly, the C-domain first binds to the outer face of the cytoplasmic membrane. The negative charge density plays a part in the binding and the kinetics of insertion ([101, 102](#)). After binding, but before insertion, the pore-forming domain must have lost the tertiary structure of the water-soluble form and adopted primarily a **molten globular** conformation, which is still compact and has kept its native secondary structure ([87, 103, 104](#)). As mentioned later, the interaction of the N-terminal and central domains with the translocation machinery might trigger the appearance of a molten globular conformation of the channel domain *in vivo* ([105](#)). The topology of the membrane-inserted state has been extensively studied, mainly with the pore-forming domains of colicins E1, A, and Ia. The membrane-bound state involves insertion of the hydrophobic helical hairpin formed by H8 and H9 into the bilayer. In contrast to colicin A, which is more dependent on

[electrostatic interactions](#), the protein would lie flat on the membrane surface. The hydrophobic hairpin could not insert because the density of charged residues on the membrane surface is too high ([103](#)).

The polypeptide translocation involved in voltage-dependent gating has been studied through three main experimental approaches: (1) By using lipophilic radiolabeled probes, Merrill and Cramer ([106](#)) demonstrated that a 42-residue region of colicin E1 corresponding to helices 5 and 6 was labeled in the presence, but not in the absence, of a membrane potential in a vesicle system with **valinomycin**-induced diffusion potential. Another segment corresponding to the hydrophobic hairpin was strongly labeled in the presence or absence of membrane potential. (2) The existence of a drastic change triggered by a transnegative membrane potential *in vivo* was also demonstrated by [disulfide-bond engineering](#) ([107](#)). The same technique further showed that  $\alpha$ -helices 1, 2, 3, and 10 remain at the membrane surface after application of a membrane potential. (3) Colicin Ia with [biotin](#), conjugated to specific cysteine residues introduced by [site-directed mutagenesis](#), were inserted into a lipid bilayer, and the channels were opened or closed by varying the membrane potential. Then [streptavidin](#) was added on the *cis*- or *trans*-side of the membrane. The results demonstrated ([108](#)) that a region of the pore-forming domain of colicin Ia reversibly flips across the membrane, implying a large conformational change when potential is applied. A region of at least 68 residues is able to flip back and forth with channel opening and closing. Several open, channel structures probably exist ([109](#)).

#### 4.4. Inhibition of Murein Biosynthesis

Colicin M is unique among the colicins in that it inhibits murein biosynthesis ([110](#)) by interfering with the dephosphorylation of C<sub>55</sub>-polyisoprenyl pyrophosphate, which leads to cell lysis ([111](#)). It also has a much shorter polypeptide chain than other colicins (29 kDa as compared to 40 to 70 kDa). Yet, it has the same basic design of three functional domains (see previous). The largest part of the C-terminal domain of colicin M resides in the cytoplasmic membrane because the carrier lipid and pyrophosphatase activity were found in the membrane fraction ([112](#)). Colicin M does not penetrate deeply into the cytoplasmic membrane but acts mostly at the periplasmic side of this membrane ([24](#)). Furthermore, cytoplasmic colicin M does not kill the cells, and only external colicin M, after its import through the energy-dependent TonB system (see previous), has access to the target in the cytoplasmic membrane. A chimeric protein of colicin M linked to a signal peptide directing the C-domain to the periplasm (the outer face of inner membrane) causes cell lysis. The amount of colicin M at the target was so high that immunity broke down ([24](#)).

### 5. The Systems Immune to Colicins

The immunity of colicinogenic cells to the action of the colicin they produce was first noted by Fredericq in 1957 ([113](#)). The cell relies for its survival on the immunity protein. The individual cell is safe from its own colicin when that gene is carried on the same plasmid as the *imm* gene. Much progress with regard to the specificity determinants has been made during recent years on the basis of the immune activity, especially for nuclease-type and pore-forming colicins.

#### 5.1. Immunity Against RNases

The determinants for specificity of colicin-immunity interactions with colicins E3, E6 and cloacin DF13 are likely to reside in a limited number of amino acid residues, eight in the nuclease domains and up to nine in the corresponding immunity proteins ([114](#), [115](#)). These determinants are located in the N-terminal regions of both the RNase domain and the immunity protein. A single residue either in Im3 or in Im6 is critical for defining the specificity ([116](#)). Other residues are also likely to be involved in the interaction but have more peripheral roles in defining specificity. The 84-residue Im3 ([117](#)) is folded into a four-stranded antiparallel  $\beta$ -sheet connected by loops and a single short  $\alpha$ -helix, in marked contrast to the structure of a DNase-specific immunity protein Im9 ([118](#)), which contains four  $\alpha$ -helices. The specificity-determining residues of Im3 at the two main positions are exposed and line one face of the  $\beta$ -sheet ([117](#)). No structures are yet available for any RNase domains.

#### 5.2. Immunity Against DNases

The sequences of the DNase domains of colicins E2, E7, E8, and E9 are more than 80% identical, yet they have specific immunity proteins. Site-directed mutagenetic studies identified six Col E9 residues in the C-terminal DNase domain as possible specificity determinants (119). The sequences of Im2, Im7, Im8, and Im9 are 58% identical but are not homologous to the RNase-type immunity proteins. Two amino acid residues were identified as specificity determinants of Im8 and Im9 by using homologous recombination (74) and site-directed mutagenesis (118, 120).

Further structural and biophysical studies of DNase-immunity protein complexes have shown that the DNase E9-Im9 complex is extremely stable (like the Col E9-Im9 complex) with a  $K_d$  of  $9.3 \times 10^{-17}$  M, one of the highest affinities ever measured for a protein interaction (121). The association of the nuclease with the Im9 protein is essentially **diffusion-controlled**, involves electrostatic steering, and follows a two-step mechanism in which the proteins form an initial encounter complex before undergoing a conformational change to yield the final stable complex (121). Although there is no cross-reactivity between DNase-type colicins (E2, E7, E8, and E9) and noncognate immunity proteins under normal levels of expression, biophysical studies indicated that the latter bind to the DNase domain of colicin E9 and inhibit its activity (122). The  $K_d$  values range from  $10^{-17}$  M (IM9) to  $10^{-4}$  M. Consistent with this result, overexpressing each of the noncognate *Im* genes in bacteria results in significant levels of cross-reactivity toward the ColE9 toxin, and the order of *in vivo* cross-reactivity (Im9 > Im2 > Im8 > Im7) mirrors exactly the measured *in vivo* affinities (122). The specificity for the DNase-Im complexes is controlled through the dissociation constant, which is more than  $10^6$ -fold faster for the noncognate Im proteins. The Im9 fold consists of a distorted antiparallel four-helix bundle in which the second helix (from the N-terminus) is the main specificity determinant (118, 120). The surface regions of Im9 that interact with the DNase include the two central helices of the molecule, and the binding surface is heavily negatively charged, consistent with a positively charged partner DNase domain (118).

### 5.3. Immunity to Pore-Forming Colicins

The immunity proteins directed against nuclease-type colicins described previously are expressed at a level similar to that of their cognate colicins. In contrast, the proteins immune to pore-forming colicins are expressed constitutively at very low levels ( $10^2$  to  $10^3$  molecules per cell). Another major difference is that these proteins protect the cells against external colicin because the membrane potential has the wrong orientation for colicin activity from the inside. Thus the immunity is directed against colicins produced by other cells.

Membrane vesicles from cells immune to colicin Ia are depolarized by colicin E1, but not Ia (123), which first indicated a cytoplasmic membrane localization for these types of immunity proteins. In addition, the construction of hybrids between colicins Ia and Ib, A and E1, or A and B, demonstrated that, independent of the translocation pathway, immunity is directed specifically against the C-terminal, pore-forming domain (19, 52, 124). ImmA has four transmembrane  $\alpha$ -helices, and both the N- and C-termini are located in the cytoplasm (125). The shorter ImE1 has three transmembrane  $\alpha$ -helices. The N- and C-termini are on the cytoplasmic and periplasmic sides of the membrane, respectively (127). These orientations are in agreement with the inside-positive prediction rule (128).

As with Im proteins directed against nucleases, there is no cross-reactivity between homologous immunity proteins directed against homologous A, B, or N colicins. Similarly, overproduction of the immunity protein leads to partial cross-reactivity (124). The main determinant for specific immunity recognition is located in the hydrophobic hairpin of the channel domain of the colicin (124, 129). Either the whole bacteriocin or the C-terminal domain of pore-forming colicins produced in the cytoplasm of *E. coli* are devoid of cytotoxicity. However, when the C-terminal domain is fused to a signal sequence, the channel is inserted and functional, and cell death follows. This could be inhibited by coproduction of the cognate immunity protein (130, 131). The cytotoxicity of the hybrid protein is independent of the uptake machinery normally used, demonstrating that the C-domain alone forms the channel *in vivo*, as in *in vitro*. The interaction of ImA with colicin A requires

the immunity protein to assemble functionally but does not require the channel to be in the open state (131). In other words, the immunity protein interacts with the hydrophobic helical hairpin, as first suggested by the studies mentioned previously (124). This interaction was demonstrated directly using an epitope-tagged immunity protein (131). Site-directed mutagenetic studies of ImA indicate that a role for polar regions of ImA cannot be excluded, in addition to the intramembrane helix–helix interactions. Two roles are proposed for the hydrophilic loops in this protein: (1) they may stabilize its interactions with the channel on both sides of the membrane; (2) they may be required for the functional assembly of the transmembrane helices of ImA (132). The colicin E1 immunity protein tolerates a higher degree of substitution than ImA (127).

#### 5.4. Immunity to Colicin M

Like the protein it inhibits, the immunity to colicin M is unique when compared with other Im proteins. It prevents colicin M from inhibiting murein synthesis. This 14 kDa protein has been localized in the cytoplasmic membrane, and a substantial portion is exposed to the periplasmic space. There is indirect evidence that the colicin M-immunity interaction occurs at the periplasmic side of the cytoplasmic membrane (133).

#### Bibliography

1. A. Gratia (1925) C.R. Soc. Biol. **93**, 1041–1041.
2. R. Kolter and F. Moreno (1992) Ann. Rev. Microbiol. **46**, 141–163.
3. P. Fredericq (1963) Ann. Rev. Microbiol. **11**, 7–22.
4. M. Riley and D. Gordon (1992) J. Gen. Microbiol. **138**, 1345–1352.
5. F. Jacob, L. Siminovitch, and L. Wollman (1952) Ann. Inst. Pasteur **83**, 295–315.
6. J. Konisky (1982) Ann. Rev. Microbiol. **36**, 125–144.
7. D. Cavard, P. Sauve, F. Heitz, F. Pattus, C. Martinez, R. Dijkman, and C. Lazdunski (1988) Eur. J. Biochem. **172**, 507–512.
8. J. K. Davis and P. Reeves (1975) J. Bacteriol. **123**, 96–101.
9. J. K. Davis and P. Reeves (1975) J. Bacteriol. **123**, 102–117.
10. R. Nagel del Zwaig and S. E. Luria (1967) J. Bacteriol. **94**, 1112–1123.
11. C. Lazdunski (1995) Mol. Microbiol. **16**, 1059–1066.
12. M. A. Riley (1993) Mol. Biol. Evol. **10**, 1048–1059.
13. U. Ross, E. E. Harkness, and V. Braun (1989) Mol. Microbiol. **3**, 891–902.
14. S. Luria and J. Suit (1987) In *Escherichia coli* and *Salmonella typhimurium* (F. E. Neidhardt, ed.), ASM, Washington, DC, pp. 1615–1624.
15. H. Masaki and T. Ohta (1985) J. Mol. Biol. **82**, 217–227.
16. S. T. Cole, B. Saint-Joanis, and A. P. Pugsley (1985) Mol. Gen. Genetics **198**, 465–472.
17. K. F. Chak and R. James (1985) Nucleic Acids Res. **13**, 2519–2530.
18. P. T. Chan, H. Ohmori, J. I. Tomizawa, and J. Lebowitz (1985) J. Biol. Chem. **260**, 8925–8935.
19. J. A. Mankovich, C. H. Hsu, and J. Konisky (1986) J. Bacteriol. **168**, 228–236.
20. J. Morlon, M. Chartier, M. Bidaud, and C. Lazdunski (1988) Mol. Gen. Genetics **211**, 231–243.
21. R. Lloubes, D. Baty, and C. Lazdunski (1986) Nucleic Acids Res. **14**, 2621–2636.
22. A. P. Pugsley (1988) Mol. Gen. Genetics **211**, 335–341.
23. S. Zhang, L. Yan, and G. Zubay (1988) J. Bacteriol. **170**, 5460–5467.
24. V. Braun, S. Gaisser, C. Glaser, R. Harkness, T. Ölschager, and J. Mende (1992) In *Bacteriocins, Microcins and Lantibiotics* (R. James, C. Lazdunski, and F. Pattus, eds.), Springer-Verlag, Berlin, Heidelberg, NATO ASI Series, Vol. **65**, pp. 119–125.
25. R. Lloubes, M. Granger-Schnarr, C. Lazdunski, and M. Schnarr (1991) J. Mol. Biol. **217**, 421–



26. D. Cavard and B. Oudega (1992) In *Bacteriocins, Microcins and Lantibiotics* (R. James, C. Lazdunski, and F. Pattus, eds.), Springer-Verlag, Berlin, Heidelberg, NATO ASI Series, Vol. **65**, pp. 297–305.
27. B. Oudega, A. Ykema, F. Stegehuis, and F. de Graaf (1984) *FEMS Microbiol. Lett.* **22**, 101–108.
28. D. Cavard, R. Lloubes, J. Morlon, M. Chartier, and C. Lazdunski (1985) *Mol. Gen. Genet.* **199**, 95–100.
29. D. Cavard, D. Baty, S. P. Howard, H. Verheij, and C. Lazdunski (1987) *J. Bacteriol.* **169**, 2187–2194.
30. A. P. Pugsley and S. T. Cole (1987) *J. Gen. Microbiol.* **133**, 2411–2420.
31. A. P. Pugsley and M. Schwartz (1984) *EMBO J.* **3**, 2393–2397.
32. S. P. Howard, D. Cavard, and C. Lazdunski (1991) *J. Gen. Microbiol.* **137**, 81–89.
33. R. Maget-Dana, F. Heitz, M. Ptak, F. Peypoux, and M. Guinaud (1985) *Biochem. Biophys. Res. Commun.* **129**, 965–971.
34. Y. Ohno-Iwashita and K. Imahori (1980) *Biochemistry* **19**, 652–659.
35. F. de Graaf and B. Oudega (1986) *Curr. Top. Microbiol. Immunol.* **125**, 183–205.
36. R. Dreher, V. Braun, and B. Wittman-Liebold (1985) *Arch. Microbiol.* **140**, 343–346.
37. D. Baty, M. Frenette, R. Lloubes, V. Géli, S. P. Howard, F. Pattus, and C. Lazdunski (1988) *Mol. Microbiol.* **2**, 807–811.
38. M. Frenette, H. Bénédicti, A. Bernadac, D. Baty, and C. Lazdunski (1991) *J. Mol. Biol.* **217**, 2509–2514.
39. R. Webster (1991) *Mol. Microbiol.* **5**, 1005–1011.
40. A. Vianney, M. Michelle Muller, T. Clavel, J. C. Lazzaroni, R. Portalier, and R. E. Webster (1996) *J. Bacteriol.* **178**, 4031–4038.
41. E. Bouveret, R. Derouiche, A. Rigal, R. Lloubes, C. Lazdunski, and H. Bénédicti (1995) *J. Biol. Chem.* **270**, 11071–11077.
42. M. Isnard, A. Rigal, J. C. Lazzaroni, C. Lazdunski, and R. Lloubes (1994) *J. Bacteriol.* **176**, 6392–6396.
43. R. Derouiche, H. Bénédicti, J. C. Lazzaroni, C. Lazdunski, and R. Lloubes (1995) *J. Biol. Chem.* **270**, 11078–11084.
44. J. C. Lazzaroni, A. Vianney, J. L. Popot, H. Bénédicti, S. Samatey, C. Lazdunski, R. Portalier, and V. Géli (1995) *J. Mol. Biol.* **246**, 1–7.
45. J. P. Bourdineaud, S. P. Howard, and C. Lazdunski (1989) *J. Bacteriol.* **171**, 2458–2465.
46. G. Guiard, P. Boulanger, H. Bénédicti, R. Lloubes, M. Besnard, and L. Letellier (1993) *J. Biol. Chem.* **269**, 5874–5880.
47. R. Benz, E. Maier, and I. Gentscher (1993) *Zbl Bakt* **278**, 187–196.
48. C. Wandersman and P. Deleplaire (1990) *Proc. Natl. Acad. Sci. USA* **87**, 4776–4780.
49. L. Gilson, H. Mahanty, and R. Kolter (1990) *EMBO J.* **9**, 3875–3884.
50. H. Bénédicti, M. Frenette, D. Baty, R. Lloubes, V. Géli, and C. Lazdunski (1989) *J. Gen. Microbiol.* **135**, 3413–3420.
51. J. P. Bourdineaud, H. P. Fierobe, C. Lazdunski, and J. M. Pagès (1990) *Mol. Microbiol.* **4**, 1737–1743.
52. H. Bénédicti, M. Frénette, D. Baty, R. Lloubes, M. Knibiehler, F. Pattus, and C. Lazdunski (1991) *J. Mol. Biol.* **217**, 429–439.
53. H. Bénédicti, C. Lazdunski, and R. Lloubes (1991) *EMBO J.* **10**, 1989–1995.
54. E. Bouveret, A. Rigal, C. Lazdunski, and H. Bénédicti (1997) *Mol. Microbiol.*, **23**, 909–920.
55. D. Duché, D. Baty, M. Chartier, and L. Letellier (1994) *J. Biol. Chem.* **269**, 24820–24825.

56. H. Bénédetti, R. Llobes, C. Lazdunski, and L. Letellier (1992) *EMBO J.* **11**, 441–447.
57. D. Jeanteur, T. Schirmer, D. Fourel, V. Simonet, G. Rummel, C. Widner, J. P. Rosenbusch, F. Pattus, and J. M. Pagès (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10675–10679.
58. H. Bénédetti, L. Letellier, R. Llobes, V. Géli, D. Baty, J. M. Pagès, and C. Lazdunski (1992) In *Dynamics of Membrane Assembly* (J.A.F. Op den Kamp, ed.), NATO ASI Series, Springer-Verlag, Berlin, Vol. **63**, pp. 316–332.
59. K. Eick-Helmerich and V. Braun (1995) *J. Bacteriol.* **171**, 5117–5126.
60. K. Postle (1990) *Mol. Microbiol.* **4**, 2019–2925.
61. H. Pils and V. Braun (1995) *Mol. Microbiol.* **16**, 57–67.
62. E. Schramm, J. Mende, V. Braun, and R. M. Kemp (1987) *J. Bacteriol.* **169**, 3350–3357.
63. V. Braun (1995) *FEMS Microbiol. Rev.* **16**, 295–307.
64. S. F. Mel, A. M. Falick, A. L. Burlingame, and R. M. Stroud (1993) *Biochemistry* **312**, 9473–9479.
65. P. Ghosh, S. Mel, and R. M. Stroud (1994) *Nature Struct. Biol.* **1**, 597–604.
66. H. Bénédetti, R. Llobes, C. Lazdunski, and L. Letellier (1992) *EMBO J.* **11**, 441–447.
67. K. Postle and J. Skare (1988) *J. Biol. Chem.* **263**, 11000–11007.
68. M. Wiener, D. Freymann, P. Ghosh, and R. Stroud (1997) *Nature* **385**, 461–464.
69. K. Jakes (1982) In *Molecular Action of Toxins and Viruses* (P. Cohen and S. von Heinegen, eds.), Elsevier, Amsterdam.
70. M. Mock and A. P. Pugsley (1982) *J. Bacteriol.* **150**, 1069–1076.
71. P. C. K. Lau, M. Parsons, and T. Uchimura (1992) In *Bacteriocins, Microcins and Lantibiotics* (R. James, C. Lazdunski, and F. Pattus, eds.), Springer-Verlag, Berlin, pp. 353–378.
72. P. S. Ringrose (1970) *Biochim. Biophys. Acta* **213**, 320–334.
73. R. Wallis, A. Reilly, A. Rowe, G. Moore, R. James, and C. Kleanthous (1992) *Eur. J. Biochem.* **207**, 687–695.
74. R. Wallis, G. R. Moore, C. Kleanthous, and R. James (1992) *Eur. J. Biochem.* **210**, 925–930.
75. R. Wallis, A. Reilly, K. Barnes, C. Abell, D. Campbell, G. Moore, R. James, and C. Kleanthous (1994) *Eur. J. Biochem.* **220**, 447–454.
76. J. Weiss and S. Luria (1978) *Proc. Natl. Acad. Sci. USA* **75**, 2483–2487.
77. H. Tokuda and J. Koninsky (1978) *Proc. Natl. Acad. Sci. USA* **76**, 6167–6171.
78. S. Schein, B. Kagan, and A. Finkelstein (1978) *Nature* **276**, 159–163.
79. W. Cramer, J. Dankert, and Y. Uratami (1983) *Biochim. Biophys. Acta* **737**, 173–193.
80. J. P. Bourdineaud, P. Boulanger, C. Lazdunski, and L. Letellier (1990) *Proc. Natl. Acad. Sci. USA* **87**, 1037–1041.
81. K. Fields and S. Luria (1969) *J. Bacteriol.* **97**, 57–63.
82. K. Fields and S. Luria (1969) *J. Bacteriol.* **97**, 64–77.
83. A. Kopecky, D. Copeland, and J. Lusk (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4631–4634.
84. C. Plate, J. Suit, A. Jetten, and S. Luria (1974) *J. Biol. Chem.* **19**, 6138–6143.
85. J. Gould and W. Cramer (1977) *J. Biol. Chem.* **252**, 5491–5497.
86. G. Guihard, H. Bénédetti, M. Besnard, and L. Letellier (1994) *J. Biol. Chem.* **268**, 17775–17780.
87. C. Lazdunski, D. Baty, V. Géli, D. Cavard, J. Morlon, R. Llobes, P. Howard, M. Knibiehler, M. Chartier, S. Varenne, M. Frenette, J. L. Dasseux, and F. Pattus (1988) *Biochim. Biophys. Acta* **947**, 445–464.
88. M. Parker, F. Pattus, A. Tucker, and D. Tsernoglou (1989) *Nature* **337**, 93–96.
89. M. Parker, J. Postna, F. Pattus, A. Tucker, and D. Tsernoglou (1992) *J. Mol. Biol.* **224**, 639–657.

90. P. Ghosh, S. Mel, and R. Stroud (1994) *Nat. Struct. Biol.* **1**, 597–504.
91. P. Elkins, H. Y. Song, W. Cramer, and C. Stauffacher (1994) *Proteins Struct. Funct. Genet.* **19**, 150–157.
92. M. Wormald, A. Merril, W. Cramer, and R. Williams (1990) *Eur. J. Biochem.* **191**, 155–161.
93. L. Holm and C. Sander (1993) *FEBS Lett.* **315**, 301–306.
94. F. Pattus, D. Cavard, R. Verger, C. Lazdunski, and H. Schindler (1983) in *Physical Chemistry of Transmembrane Ion Motions* (G. Spach, ed.), Elsevier, Amsterdam, pp. 407–413.
95. F. Pattus, D. Massote, H. Wilmsen, J. Lakey, D. Tsernoglou, A. Tucker, and M. Parker (1990) *Experientia* **96**, 180–192.
96. M. Collarini, G. Amblard, C. Lazdunski, and F. Pattus (1987) *Eur. Biophys. J.* **14**, 147–153.
97. P. Gosh, S. Mel, and R. Stroud (1993) *J. Membrane Biol.* **134**, 85–92.
98. L. Raymond, S. Slatin, and A. Finkelstein (1985) *J. Membrane Biol.* **84**, 173–181.
99. J. Bullock (1992) *J. Membrane Biol.* **125**, 255–257.
100. J. Bullock, E. Kolen, and J. L. Shear (1992) *J. Membrane Biol.* **128**, 1–16.
101. F. van der Goot, N. Didat, F. Pattus, W. Dowhan, and L. Letellier (1993) *Eur. J. Biochem.* **213**, 217–221.
102. G. van der Goot, J. Gonzalez-Manas, J. Lakey, and F. Pattus (1991) *Nature* **354**, 408–410.
103. J. Lakey, G. van der Goot, and F. Pattus (1994) *Toxicology* **87**, 85–108.
104. M. Parker and F. Pattus (1993) *Trends Biochem. Sci.* **18**, 391–395.
105. D. Duché, D. Baty, M. Chartier, and L. Letellier (1994) *J. Biol. Chem.* **269**, 24820–24825.
106. A. Merrill and W. Cramer (1990) *Biochemistry* **29**, 8529–8534.
107. D. Duché, M. Parker, J. Gonzalez-Manas, F. Pattus, and D. Baty (1994) *J. Biol. Chem.* **269**, 6332–6339.
108. S. Slatin, X-Q. Qiu, K. Jakes, and A. Finkelstein (1994) *Nature* **371**, 158–161.
109. X. Q. Qiu, K. Jakes, P. Kienker, A. Finkelstein, and S. Slatin (1996) *J. Gen. Physiol.* **107**, 313–328.
110. K. Schaller, J. Höltje, and V. Braun (1982) *J. Bacteriol.* **152**, 994–1000.
111. R. Harkness and V. Braun (1989) *J. Biol. Chem.* **264**, 6177–6182.
112. G. Siewert and J. Strominger (1967) *Proc. Natl. Acad. Sci. USA* **57**, 767–773.
113. P. Fredericq (1957) *Ann. Rev. Microbiol.* **11**, 7–21.
114. H. Masaki, S. Yajima, A. Akutsur-Koide, T. Ohta, and T. Uozumi (1992) In *“Bacteriocins, Microcins and Lantibiotics”* (R. James, C. Lazdunski, and F. Patus eds.), Springer-Verlag, Berlin, pp. 379–395.
115. A. Akutso, H. Masaki, and T. Ohta (1989) *J. Bacteriol.* **171**, 6430–6436.
116. H. Masaki, A. Akutsu, T. Uozumi and T. Ohta (1991) *Gene* **107**, 133–138.
117. S. Yajima, Y. Muto, S. Yokoyama, H. Masaki, and T. Uozumi (1992) *Biochemistry* **31**, 5578–5586.
118. M. Osborne, A. L. Breez, L. Y. Lian, A. Reilly, R. James, C. Kleanthous, and G. Moore (1996) *Biochemistry* **35**, 9505–9512.
119. M. Curtis and R. James (1991) *Mol. Microbiol.* **5**, 2727–2733.
120. M. Osborne, L-Y. Lian, R. Wallis, A. Reilly, R. James, C. Kleanthous, and G. Moore (1994) *Biochemistry* **33**, 12347–12355.
121. R. Wallis, G. Moore, R. James, and C. Kleanthous (1995) *Biochemistry* **34**, 13743–13750.
122. R. Wallis, K. Y. Leung, A. Pomoner, H. Videler, G. Moore, R. James, and C. Kleanthous (1995) *Biochemistry* **34**, 13751–13759.
123. C. Weaver, A. Redborg and J. J. Konisky (1981) *J. Bacteriol.* **148**, 817–828.
124. V. Géli and C. Lazdunski (1992) *J. Bacteriol.* **174**, 6432–6437.

125. V. Géli, D. Baty, and C. Lazdunski (1988) *Proc. Natl. Acad. Sci. USA* **85**, 689–693.
126. V. Géli, D. Baty, F. Pattus, and C. Lazdunski (1989) *Mol. Microbiol.* **3**, 679–687.
127. H. Song and W. Cramer (1991) *J. Bacteriol.* **173**, 2935–2943.
128. G. von Heijne (1992) *J. Mol. Biol.* **225**, 487–494.
129. Y. Zhang and W. Cramer (1993) *J. Biol. Chem.* **268**, 1–8.
130. D. Espeset, Y. Corda, K. Cunningham, H. Bénédicti, R. Lloubes, C. Lazdunski, and V. Géli (1994) *Mol. Microbiol.* **13**, 1121–1131.
131. D. Espeset, D. Duché, D. Baty, and V. Géli (1996) *EMBO J.* **15**, 2356–2364.
132. D. Espeset, P. Piet, C. Lazdunski, and V. Géli (1994) *Mol. Microbiol.* **10**, 1111–1120.
133. T. Öschlänger, A. Turba, and V. Braun (1991) *Mol. Microbiol.* **5**, 1105–1111.

### Suggestions for Further Reading

134. H. Bénédicti and V. Géli (1996) "Colicin transport, channel formation and inhibition". In *Handbook of Biological Physics* (W. Konings, H. Kaback, and J. S. Lolkema, eds.), Elsevier, Amsterdam, Vol. **2**, pp. 665–691.
135. W. Cramer, J. Heymann, S. Schendel, B. Deriy, F. Cohen, P. Elkins, and C. Stauffacher (1995) Structure-function of the channel-forming colicins. *Ann. Rev. Biophys. Biomol. Struct.* **24**, 611–641.
136. R. James, C. Kleanthous, and G. Moore (1996) The biology of E colicins: Paradigms and paradoxes. *Microbiology* **142**, 1569–1580.
137. A. P. Pugsley (1984) The ins and outs of colicins: Part 1. Production and translocation across membranes. Part 2. Lethal action, immunity and ecological implications, *Microbiol. Sci.* **1**, 168–175 and 203–205.
138. V. Braun, H. Pilsel, and P. Grob (1994) Colicins: Structures, modes of action, transfer through membranes and evolution. *Arch. Microbiol.* **161**, 199–206.

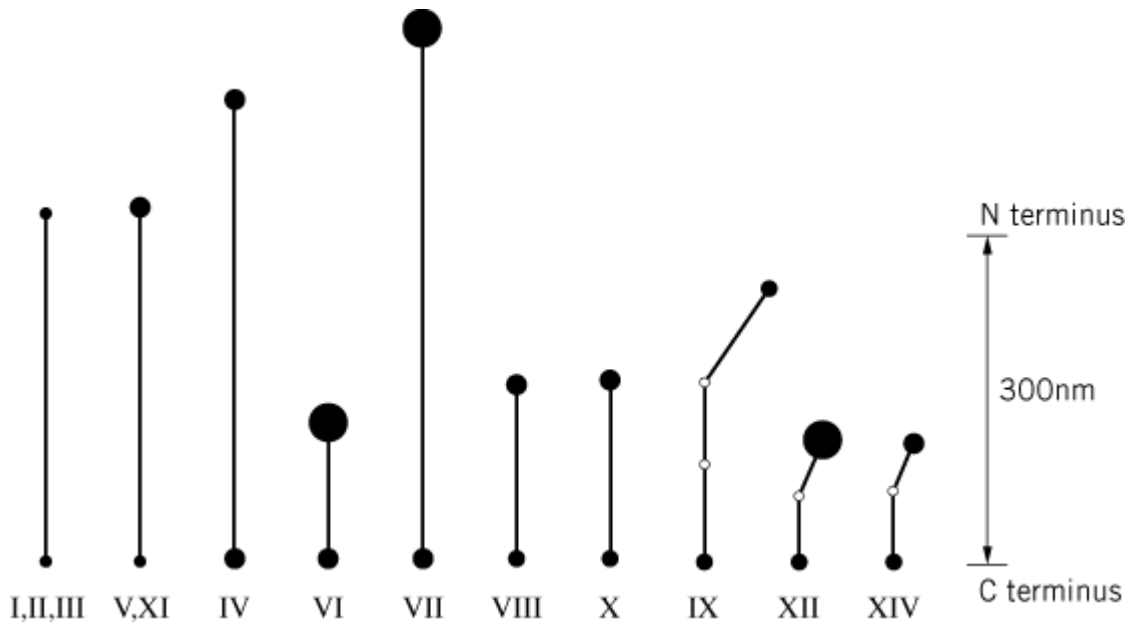
## Collagen

[Proteins](#) of the *collagen* family are the major components of the [extracellular matrix](#) and facilitate the formation and maintenance of a multicellular system. Collagens serve as solid-state regulators for cellular function and as scaffolding of the tissue architecture, particularly in large vertebrates. They contain many [proline](#) and [glycine](#) residues and a distinct **secondary structure**: the [polyproline II-like helix](#), which is distinct from the [a-helix](#), [b-sheet](#) and [turn](#) secondary structures in other [protein structures](#). This regular secondary structure arises because the sequences the collagen polypeptide chain consist largely of the repeated sequence Gly—X—Y, with abundant proline residues at the X positions and rich in *hydroxyproline* (Hyp) residues at the Y positions. Three polyproline II-like helices constitute the [supersecondary structure](#) of the collagenous triple helix, which is stabilized through [hydrogen bonds](#) nearly perpendicular to triple helical axis (see [Triple-Helical Proteins](#)). The collagen protein family includes all the structural proteins of the extracellular matrix with triple-helical collagenous **domains** in their molecular architecture.

The collagen superfamily is classified into groups (Table [1](#)) according to their molecular and/or supramolecular structures. The molecular structures of the collagenous proteins can be depicted with ball-and-stick models (Fig. [1](#)). The balls represent noncollagenous or globular domains, without abundant glycine and proline residues, while the sticks show the collagenous triple helices. Some of the triple-helical domains, including that of the type IV collagen, have interruptions in the Gly—X

—Y triplets, in that glycine residues do not always occupy every third position.

**Figure 1.** Molecular architecture of collagen superfamily depicted with ball and stick models (stick = collagenous, triple-helical domain; ball = noncollagenous domain).



**Table 1. Collagen Classification**

| Family   | Types of Molecuels or Chains                 |
|--|--|
| Fibrillar collagen   | Type I, type II, type III, type V, type XI   |
| Meshwork-forming collagen  | a1, a2, a3, a4, a5, and a6 chains of type IV |
| Fibril associated collagen with type interruptedtriple helices (FACIT) | Type IX, type XII, type XIV, type XVI        |
| Collagen with long triple helix  | Type VII                                     |
| Collagen with short triple helix                                       | Type VIII, type X, type VI                   |
| Membrane associated collagen   | Type XVII                                    |
| Others   | Type XIII, type XV, type XVIII, type XIX     |

The structure, assembly, and supramolecular aggregation of type I collagen is the prototype from which has developed our understanding of collagenous structure, particularly the fibrillar collagens. Type I collagen is one of the major components of the fibrous collagens that occur in the greatest

amounts. The structure and characteristic properties of triple-helical domains have been deduced from the study of type I collagen, together with comparative studies of other types of collagen. Unless otherwise mentioned, the description of collagen triple helices given below is based primarily on the information obtained through the studies type I collagen. These properties are generally shared with the triple-helical domains in other collagen types, especially regarding the characteristic features distinct from  $\alpha$ -helix or  $\beta$  structure in noncollagenous proteins. However, most recent studies suggest that the triple-helical regions have structures and properties specific for each type, particularly in their intermolecular interactions. The characteristic features of the various collagen types are due to differences in the distribution of charged residues along the triple helix, content of glycosylated hydroxylysine and bulky **hydrophobic** residues, as well as imperfection of the Gly—X—Y repeat, which is a prerequisite for the triple-helical conformation.

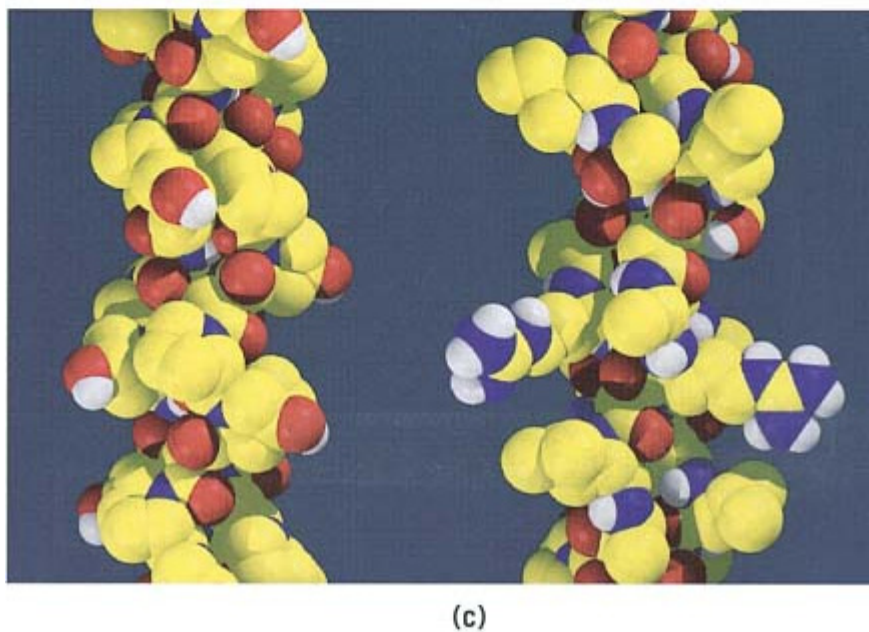
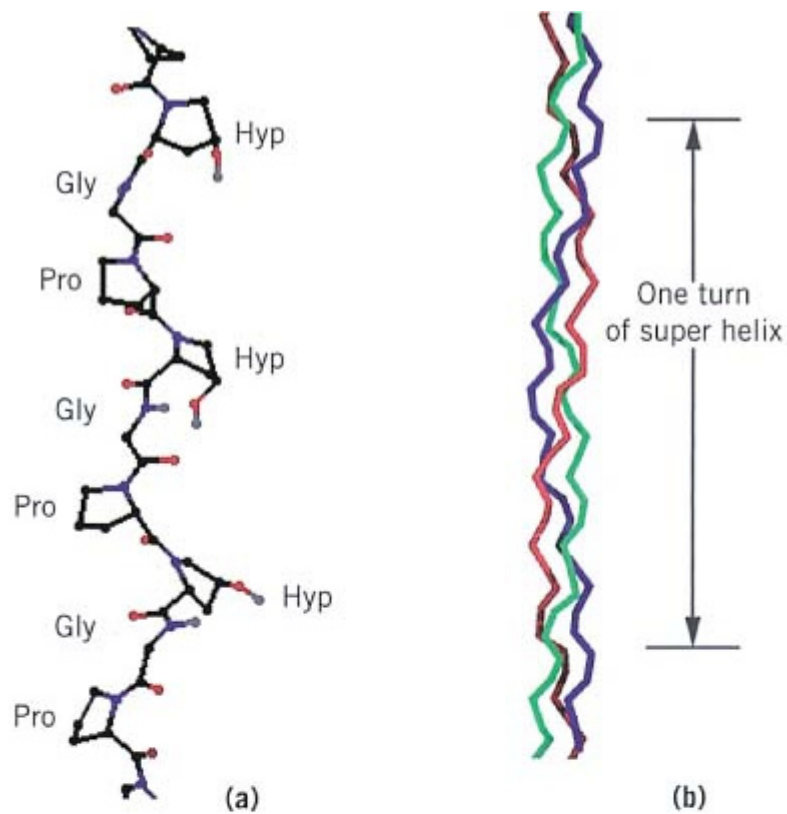
## 1. Collagenous Triple Helix; Common and Specific Features among Different Types of Collagen

### 1.1. Shape and Structure of the Triple-Helical Domains

A variety of physicochemical studies, including direct visualization of individual molecules by rotary shadowing techniques in the **electron microscope**, have demonstrated the rod-like nature of the collagen triple-helical domain. The molecules of type I collagen have a length of about 300 nm (Fig. 1) and a diameter of about 1.4 nm. The rod is neither rigid nor randomly flexible, but appears to possess an intermediate level of semiflexibility.

Details of the three-dimensional structure of the collagen triple helix were established by model building to fit X-ray [fiber diffraction](#) data. Frequent occurrence of proline and hydroxyproline residues (which together account for approximately two-ninths (~22%) of the amino acid residues) favor a polyproline II-like conformation. The axial distance between one amino acid and the next in the polyproline II-like helical structure is 0.286 nm, close to twice that in the  $\alpha$ -helix (0.15 nm). The triple-helical structure of the synthetic peptide (Pro—Pro—Gly)<sub>10</sub> was also determined by [X-ray crystallography](#). The overall helical symmetry is left-handed, with 10 residues per three turns (108°/residue) or 7 residues per two turns (103°/residue), with a pitch of 2.9 nm. The three polyproline II-like helical chains are further coiled about a central axis, to form a right-handed helix (Fig. 2). The occurrence of glycine as every third residue gives rise to a polymer of repeating tripeptide units with the formula of (—Gly—X—Y—).

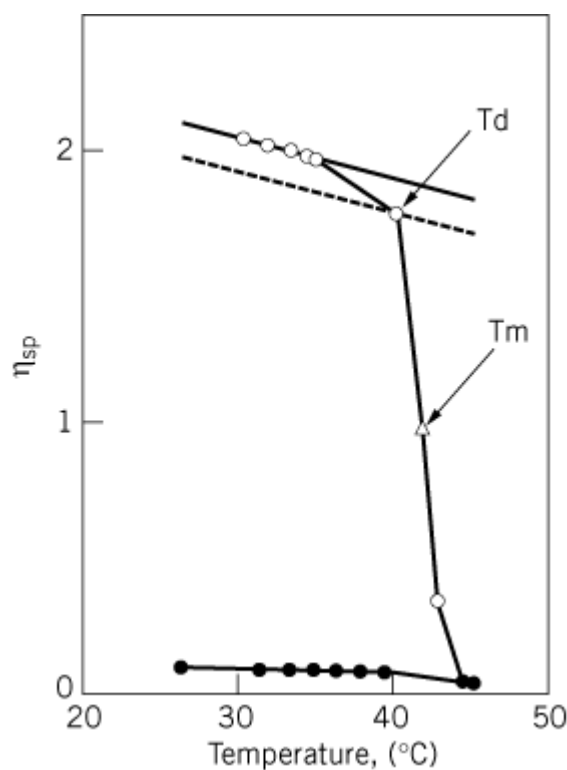
**Figure 2.** Collagenous triple helix. (a) Gly—Pro—Hyp repeated sequence of a chain in a model of part of one a chain in the collagenous triple helix. Both NH and CO groups project perpendicular to the fibrillar axis (C, black; N, blue; O, red; H, gray). (b) Backbone of the collagenous triple helix. (c) Gly—Pro—Hyp trimer (*left*). Note that there is a groove on the surface of the helix. A part of human type III collagen molecule, [a1(III)]<sub>3</sub> in the sequence of GITGARGLAGP (*right*). Note that all the residues except Gly project to the outer surface of the molecule (C, yellow; N, blue; O, red; H, white). See color insert.



## 1.2. Thermal Stability of the Triple-Helical Conformation

Flexibility of the collagenous triple helix may vary along the chain. Collagen triple helices of different types have varying flexibility, depending on what residues occupy the X and Y positions. Yet the thermal stability in terms of the helix-to-coil transition temperature (see [Helix-Coil Theory](#)) is similar, regardless of the type of collagen, in the same animal or animal tissues. The collagenous domains are heat-stable up to the upper limit of animal body temperature. Heat denaturation of collagenous domains starts around 37°C in mammalian collagens (Fig. 3). The most recent study on the denaturation temperature of the bovine type IV collagen triple-helical domain indicated that the domain also has denaturation temperature above 37°C. This suggests that interruptions in the triplet repeats do not greatly decrease the thermal stability.

**Figure 3.** Denaturation temperature of collagen triple helix. The temperature was raised stepwise by 1.5°C at intervals of 20 min, and the specific viscosity ( $\eta_{sp}$ ) was measured. Pepsin-treated acid-soluble collagen from calf skin (0.8 mg/mL) was dissolved in 0.15 M potassium phosphate buffer, pH 6.8, and 1 M glucose. Open circles indicate the specific viscosity of the native collagen solution with increasing temperature; filled circles correspond to the values of the denatured collagen solution when the temperature was lowered the same way as it was raised. The broken line is drawn through the values 5% less in the specific viscosity compared with that expected for the native collagen solution. The triangle ( $T_m$ ) corresponds to the denaturation temperature obtained by the conventional method of taking the midpoint of the curve.

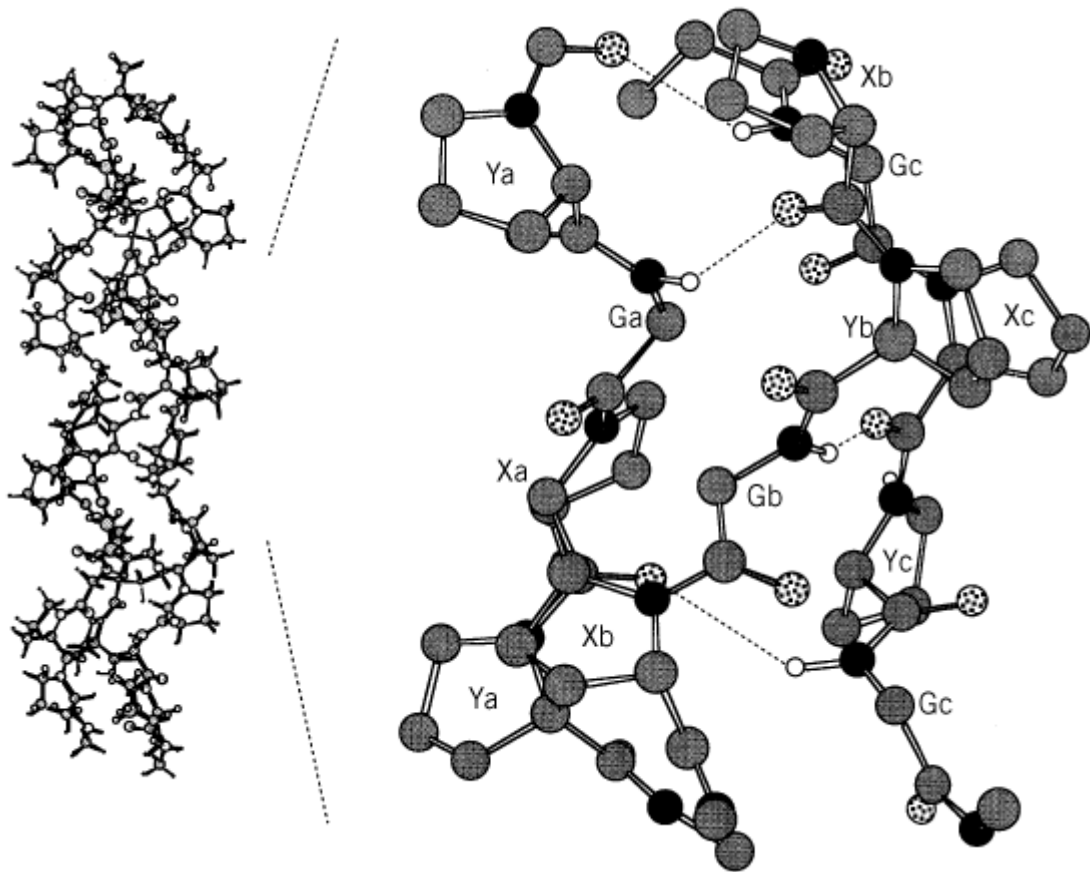


### 1.3. Primary Sequence Required for the Triple-Helical Conformation

The characteristic [primary structure](#) of polypeptides adopting collagenous triple-helical structures consists of repeats of the sequence Gly—X—Y, where Gly represents glycine and X or Y represents any other amino acid residue. Glycine is the only amino acid that can pack tightly at the center of the triple-stranded collagen monomer, where it provides HN groups for hydrogen bonding to O=C— groups in the [peptide bonds](#) of the other chains (Fig. 4). The stability of the triple-helical conformation is due in part to these hydrogen bonds, which are aligned nearly perpendicular to the helical axis. Interruptions of the Gly—X—Y repeats decrease the stability of the conformation. Substitution of a glycine residue occurring within the sequence Gly—X—Y of triple-helical domains destabilizes and disrupts the helical conformation. In the triple helix, the side chains of all the X and Y residues are exposed on the surface of the triple helix.

**Figure 4.** A schematic drawing illustrating direct interchain hydrogen bonding between (Gly)NH·CO(Pro in X position) the collagenous triple-helix. The three chains (a, b, and c) with repeated Gly-Pro(x)-Pro(y) sequence, eg, -Ga-Xa-Ya, -C Yb-.





#### 1.4. Proline Residues at Position X and Hydroxyproline Residues at Position Y

The helical conformation of individual  $\alpha$  chains arises largely as a result of steric repulsion between proline residues in the X position (approximately 120 residues per  $\alpha 1(I)$ ) and 4-hydroxyproline residues in the Y position (approximately 100 residues per  $\alpha 1(I)$  chain) and because the five-membered rings of imino acids are rigid and limit rotation about the peptide N—C bond. The proline and hydroxyproline residues also stabilize the triple helix. The contribution to helix stability from the pyrrolidine rings of proline and hydroxyproline is thought to be entropic, in that these residues may not acquire as much freedom of rotation upon **denaturation** as other residues. Another interpretation for contribution of the pyrrolidine rings to the stability of triple-helical conformation is related to the fact that these side chains are located on the surface of the triple helix. Pyrrolidine rings are surprisingly favorable in contact with water. Furthermore, [hydroxylation](#) of the proline residues before Gly or Y positions increases the thermal stability greatly although hydroxyproline residues at X positions or after Gly decrease the stability. Whether the hydroxyl group is at the 3 or 4 position of proline residues also influences greatly the thermal stability of the triple helical conformation.

The amino acid sequence responsible for the formation and stability of the collagenous triple helix is susceptible to **proteolysis** by collagenases. The sequence specifically recognized by bacterial collagenases is usually in sequences such as —GlyXYGlyProYGlyXHypGlyXY—. Cleavage occurs at the amino side of the Gly residues.

#### 1.5. Resistance of the Triple-Helical Domains to Pepsin or Other Proteinases

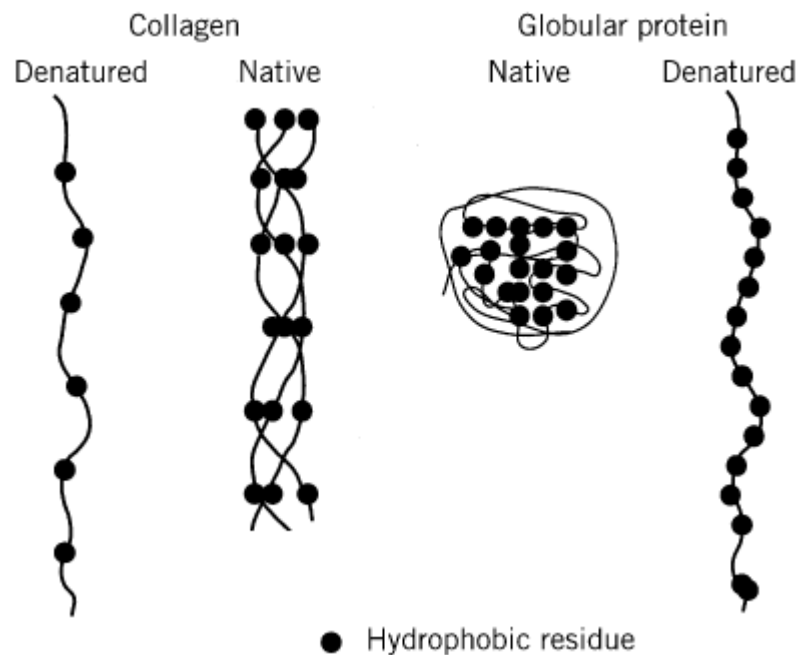
The intact triple-helical domain is generally resistant to most [proteinases](#). However, when heated above physiological temperatures, it undergoes a helix-to-coil transition and, once melted, becomes susceptible to degradative enzymes. Protein consisting primarily of collagenous domains stays in solution at acidic pH. The collagenous domains have generally been isolated by **pepsin** treatment of otherwise insoluble tissues. In a number of triple helices of recently discovered collagen family members, the occurrence of glycine every third residue is occasionally interrupted. The interrupted sites are speculated to lower the stability of the triple helix and might form kinks in the rod-like triple helix. Thus, collagen helix that is otherwise unsusceptible to many proteinases might become susceptible at the interrupted sites of Gly—X—Y tripeptides.

Thus, type IV collagen may be degraded gradually at the interrupted sites of glycine—X—Y triplets, being liberated from the aggregates with pepsin treatment.

### 1.6. Other Amino Acid Residues at Positions X and Y

All side chains at positions X and Y protrude out along the surface of the triple helix. Consequently contribute to the [hydrophilicity](#), [ionization](#), hydrophobicity (see Fig. 5), and steric roughness of the **molecular surface**, including the size of the groove of the triple helix (Fig. 2). Charged groups, together with their neighboring sequences, may also affect the stability of the triple helix, presumably due to differential contributions from [water](#) of [hydration](#) on the surface.

**Figure 5.** Hydrophobic residues in collagen polypeptide chains. The content of hydrophobic residues in collagenous do (indicated as filled circles) is relatively low compared to the globular proteins. However, they projected to the surface of triple helix, while the globular protein keeps most of the hydrophobic residues inside.



The amino acid residues other than proline or hydroxyproline residues at positions X and Y can provide another classification of the collagen protein family. In general, the collagenous triple-helical domains contain a higher content of basic (arginine + lysine) than of acidic (aspartate + glutamate) residues, resulting in a basic [isoelectric point](#). Two groups of collagens can be classified on the basis of their relative content of arginine and lysine: the high arginine group ( $\text{Arg/Lys} > 1$ ) and the low arginine group ( $\text{Arg/Lys} < 1$ ). A greater content of Lys residues may result in a greater content of hydroxylysine residues, which further provides a possibility for additional **glycosylation**. A high content of glycosylated hydroxylysine should contribute greatly to the surface roughness of the triple-helical domains. Another general feature about the amino acid composition of collagens is the low content of hydrophobic amino acid residues. The interaction within the triple helix is not stabilized by hydrophobic interactions, but they are strong between triple helices because all the hydrophobic residues are exposed on the surface of the triple helix. The content of large hydrophobic residues in the collagenous triple helices also classifies the collagenous proteins into two groups. One group contains a high ratio of Ala/hydrophobic amino acids (Val, Leu, Ile, Phe, and Met) and the other group contains a low ratio.

## 2. Classification of Collagens Based on the Primary Sequence

In the fibrillar collagens, the triple-helical conformation occurs throughout 95% of the length of the monomer (Fig. 3). Thus, of the 1057 residues in the  $\alpha 1(I)$  chain of human collagen, 1014 occur in regular Gly—X—Y triplets. The *N*-terminal 17 residues and the *C*-terminal 26 residues (referred to as *telopeptides*) do not have glycine as every third residue. The type I collagen helical molecule is a heterotrimer composed of two identical  $\alpha 1(I)$  chains and one  $\alpha 2(I)$  chain. The  $\alpha 1(I)$  and  $\alpha 2(I)$  chains are very similar, but their primary structures are coded by separate genes and are sufficiently different for the chains to be separated by [exchange chromatography](#) and by [SDS-PAGE](#). The  $\alpha$ -chains each contain over 1000 amino acid residues and have molecular weights of approximately 95,000.

The existence of a family of collagenous proteins in the connective tissues of vertebrates was first identified when cartilage collagen (type II) was found to be genetically distinct from the type I collagen of skin and tendon. A third collagen (designated type III) was detected in skin. More than 30 different collagen polypeptides have been found in the extracellular matrix in the form of at least 19 different collagen types (designated from I to XIX; see Table 2).

**Table 2. Polypeptide Chain Composition of Genetically Distinct Collagen Types**

| Type  | Known chains    | Known or Putative Composition of Molecules at Present | Length of Triple-Helical Domain (nm) | Main or Known Distribution                              | Aggregate Structure of Purified Protein or that Estimated with Immunohistochemical study, etc. |
|---|-----------------|---|--------------------------------------|---|--|
| <b>Fibrillar Collagen</b>   |                 |   |                                      |   |  |
| I   | $\alpha 1(I)$   | $[\alpha 1(I)]_2\alpha 2(I)$                          | 300                                  | Almost all connective tissues without hyaline cartilage | Fibril   |
|   | $\alpha 2(I)$   | $[\alpha 1(I)]_3?$                                    |                                      |   |  |
| II  | $\alpha 1(II)$  | $[\alpha 1(II)]_3$                                    | 300                                  | Cartilage   | Fibril   |
| III   | $\alpha 1(III)$ | $[\alpha 1(III)]_3$                                   | 300                                  | Almost similar to type I                                | Fibril   |
| V   | $\alpha 1(V)$   | $[\alpha 1(V)]_2\alpha 2(V)$                          | 300                                  | Almost similar to type I, adult cartilage               | Fibril   |
|   | $\alpha 2(V)$   | $\alpha 1(V)\alpha 2(V)\alpha 3(V)$                   |                                      |   |  |
|   | $\alpha 3(V)$   | $[\alpha 1(V)]_3?$                                    |                                      |   |  |
| XI  | $\alpha 1(XI)$  | $\alpha 1(XI)\alpha 2(XI)\alpha 3(XI)$                | 300                                  | Cartilage   | Fibril   |
|   | $\alpha 2(XI)$  |   |                                      |   |  |
|   | $\alpha 3(XI)$  |   |                                      |   |  |
| V/XI  | $\alpha 1(XI)$  | $[\alpha 1(XI)]_2\alpha 2(V)$                         | 300                                  | Vitreous body   | Fibril   |
|   | $\alpha 2(V)$   |   |                                      |   |  |
| <b>FACIT (Fibril-Associated Collagen with Interrupted Triple-Helices)</b> |                 |   |                                      |   |  |
| IX  | $\alpha 1(IX)$  | $\alpha 1(IX)\alpha 2(IX)\alpha 3(IX)$                |                                      | Surface of the  | Aggregated with fibrillar  |

|                                  |           |   |     |                                       |   |
|----------------------------------|-----------|---|-----|---------------------------------------|---|
|                                  |           | (IX)  |     | cartilage fibril                      | collagen periodically on the cartilage collagen fibrils |
|                                  | a2(IX)    |   |     |                                       |   |
|                                  | a3(IX)    |   |     |                                       |   |
| XII                              | a1(XII)   | [a1(XII)] <sub>3</sub>                                    |     | Tissues rich in type I                |   |
| XIV                              | a1(XIV)   | [a1(XIV)] <sub>3</sub>                                    |     | Tissues rich in type I                |   |
| XVI                              | a1(XVI)   | [a1(XVI)] <sub>3</sub>                                    |     |                                       |   |
| <b>Others</b>                    |           |   |     |                                       |   |
| VI                               | a1(VI)    | a1(VI)a2(VI)a3(VI)  | 100 | Almost all connective tissues         | Beaded microfibril                                      |
|                                  | a2(VI)    | [a1(VI)] <sub>2</sub> a2(VI)                              |     |                                       |   |
|                                  | a3(VI)    |   |     |                                       |   |
| VII                              | a1(VII)   | [a1(VII)] <sub>3</sub> ?                                  | 420 | Anchoring fibril                      | Short dimer   |
| VIII                             | a1(VIII)  | [a1(VIII)] <sub>3</sub>                                   | 150 | Basement membrane of endothelial cell | Hexagonal array   |
|                                  | a2(VIII)  | [a1(VIII)] <sub>2</sub> a2(VIII), [a2(VIII)] <sub>3</sub> |     |                                       |   |
| X                                | a1(X)     | [a1(X)] <sub>3</sub>                                      | 130 | Hypertrophic cartilage                | Hexagonal array   |
| XIII                             | a1(XIII)  | [a1(XIII)] <sub>3</sub> ?                                 |     |                                       |   |
| XV                               | a1(XV)    | [a1(XV)] <sub>3</sub> ?                                   | 150 |                                       |   |
| XVII                             | a1(XVII)  | [a1(XVII)] <sub>3</sub>                                   |     |                                       |   |
| XVIII                            | a1(XVIII) | [a1(XVIII)] <sub>3</sub> ?                                |     |                                       |   |
| XIX                              | a1(XIX)   | ?   |     |                                       | FACIT?  |
| <b>Meshwork-Forming Collagen</b> |           |   |     |                                       |   |
| IV                               | a1(IV)    | [a1(IV)] <sub>2</sub> a2(IV)                              | 350 | Basement membrane                     | Polygonal meshwork                                      |
|                                  | a2(IV)    | a3(IV)a4(IV)a5(IV), [a3(IV)] <sub>2</sub> a4(IV)          |     | Sinusoid                              |   |
|                                  | a3(IV)    | [a5(IV)] <sub>2</sub> a6(IV)?                             |     |                                       |   |
|                                  | a4(IV)    |   |     |                                       |   |
|                                  | a5(IV)    |   |     |                                       |   |
|                                  | a6(IV)    |   |     |                                       |   |

---

The collagen numbering system (with Roman numerals for each collagen type and Arabic numerals for individual  $\alpha$ -chains) to some extent reflects the relative abundance of the various collagens, in that the more abundant collagens were identified earliest. In addition to these collagens, there exists a number of secreted proteins that contain collagenous amino acid sequences and short triple-helical conformations, such as the **complement** component C1q, acetylcholine esterase, lung surfactant protein, conglutinin, serum mannose-binding protein, scavenger receptors (AR-I and AR-II), and MARCO. The collagenous sequences in these proteins contribute to their distinctive structures and functions. Since they have no known structural role in the extracellular matrix, however, they are not classified as collagens.

From the data derived from amino acid and gene sequencing, collagen molecules can be grouped into groups shown in Figure 1 and Table 1. Fibrillar collagen molecules are characterized by an uninterrupted triple-helical domain of approximately 300 nm. They are synthesized as procollagens comprised of three polypeptide chains that undergo processing to single chains and subsequently assemble into collagen fibrils and fibers. Fibrillar collagen molecules (ie, types I, II, III, V, and XI) exhibit several common structural features that reflect the highly conserved exon-intron structure of the genes.

Polygonal meshwork-forming collagens (type IV collagen polypeptides) have large triple-helical domains (>160 kDa) with a length of >350 nm. Their primary structures are characterized by imperfections in the  $\text{—Gly—X—Y—}$  triplet sequence. These interruptions are a particular feature of type IV collagen, in which the helical domain contains more than 20 short stretches of non-helix-forming amino acids.

Short triple-helical collagen molecules (types VI, VIII, X) contain interruptions in the helical domain (as in types IX, XII, and XIV). Collagen types VIII and X show remarkable homology and might have similar roles in tissues. Type XII and type XIV collagens have similarities to type IX collagen in their domain structures. A portion of these triple-helical domains have the potential to interact with fibrillar collagen. Thus, these three types of collagen, plus type XVI, comprises a group of fibril-associated collagens with interrupted triple helices (FACIT collagens).

An alternative approach for classifying the collagens depends on supramolecular structures that might be related to their physiological function. Individual collagen types may themselves represent a family of related collagenous structures in the extracellular matrix. Type IV collagen is a family of six homotrimeric  $\alpha$  chains ( $\alpha 1$ ,  $\alpha 2$ ,  $\alpha 3$ ,  $\alpha 4$ ,  $\alpha 5$ , and  $\alpha 6$ ), and type V/XI is a family of 6  $\alpha$  chains:  $\alpha 1(V)$ ,  $\alpha 2(V)$ ,  $\alpha 3(V)$ ,  $\alpha 1(XI)$ ,  $\alpha 2(XI)$ , and  $\alpha 3(XI)$ .

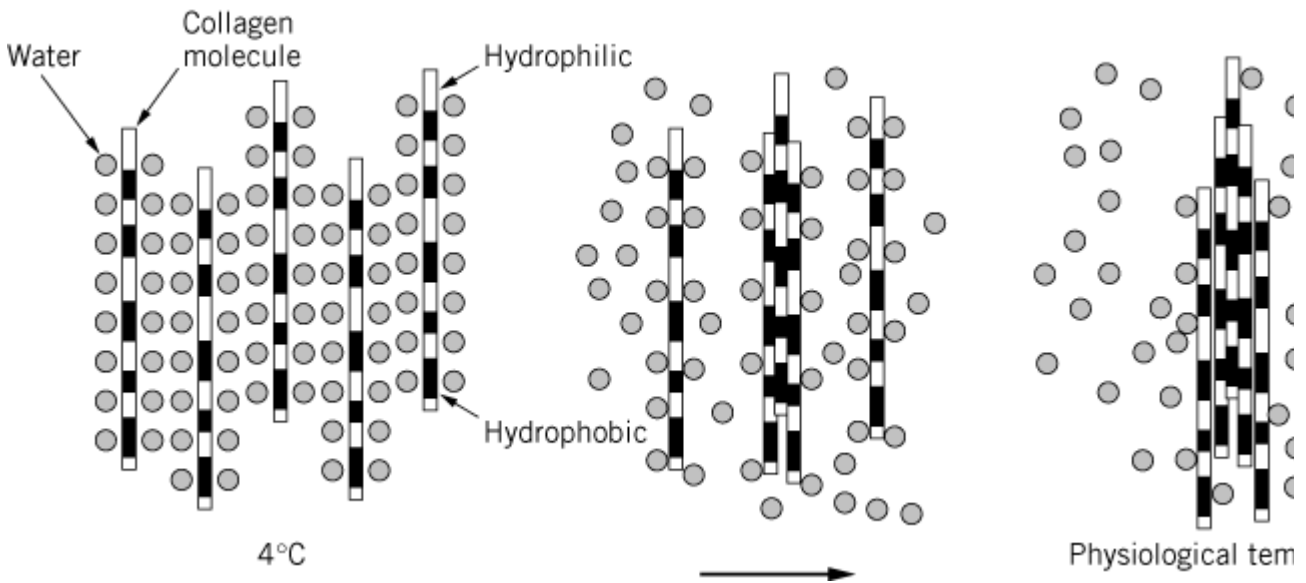
### 3. Self-assembly into Supramolecular Structure (Fibrillar Aggregates or Polygonal Meshworks)

#### 3.1. Lateral Interactions Between Triple-Helical Domains

The collagen triple helix is essentially a one-dimensional molecule, with its near-constant axial separation between amino acid residues throughout the triple-helical domain. This has permitted a direct correlation between structural data obtained by electron microscopy with chemical sequence data. Lateral association of collagenous triple-helical domains, to form fibrils, is a common feature of the collagen family. Fibrils self-assemble spontaneously from solutions of extracted type I collagen when the pH, temperature, and ionic strength are adjusted to physiological values. This lateral association appears to be promoted by elevated temperature, suggesting that **hydrophobic interactions** are the principal driving force (Fig. 6). Fibril formation is controlled to a large extent by the amino acid sequence of the collagen and, in particular, the distribution of **polar** and hydrophobic residues that are exposed on the surface of the triple-helical domain. Hydroxylysine and glycosylated hydroxylysine residues might be a most potent candidate involved in limiting the lateral growth of the collagenous domains. The residues on the surface of collagen triple helices may well be bulky enough to cause steric hindrance in the lateral association of the triple-helical domains (Fig. 15, see later). The situation would be particularly pronounced in the triple-helical domains of type I and type IV collagens, which respectively contain more than 20 and 50 glycosylated hydroxylysine residues.

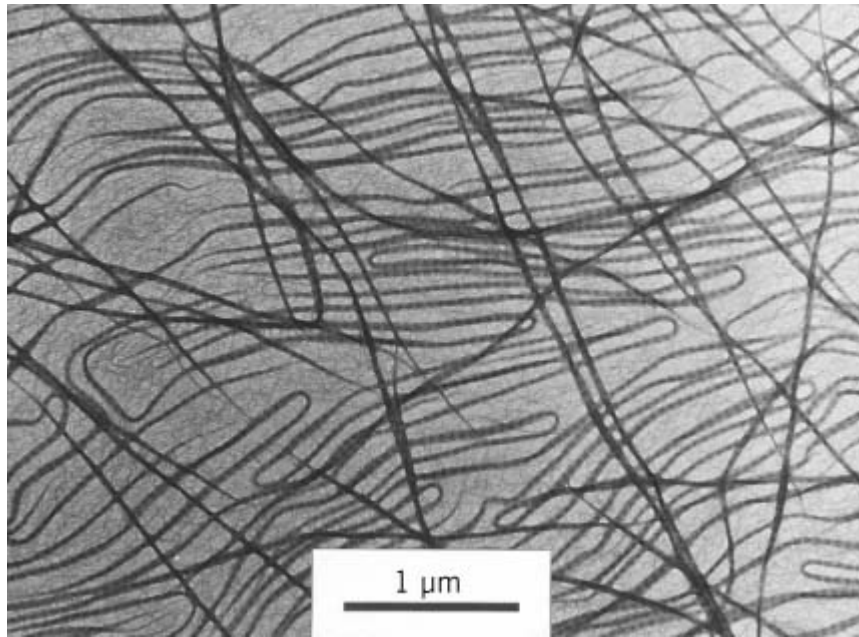
in a single polypeptide.

**Figure 6.** Lateral association of triple-helical domains. Triple-helical domain is highly hydrated at 4°C. At physiological temperature such as 37°C, water molecules on the hydrophobic domains will be dehydration. Collagen molecules accomplish lateral association of the triple-helical domains through the exposed hydrophobic regions.



There is stringent control over lateral growth of the collagen fibrils *in vivo*, in that the distributions of fibrillar diameter are sharp in a wide range of extracellular matrices of various tissues. That type I collagen can be organized in fibrils with different diameters may indicate the involvement of other regulatory elements. Heterotypic fibrils containing molecules of type III and type V, in addition to type I, collagen molecules may predominate *in vivo*. The presence of the **propeptide** of type III procollagen, for example, on the surface of collagen fibrils has led to the suggestion that this peptide may have a role in limiting the fibril diameter. Type V collagen was shown to form banded fibrils with finite diameter (Fig. 7). It has been postulated that the type V collagen can form hybrid fibrils with the type I collagen. The responsible component in the type V collagen appears to be the triple-helical domain, since pepsin-treated type V collagen forms only fine fibrils. This finding suggests that type V collagen functions as one of the regulatory elements in limiting fibril diameter, since the triple-helical domain of type V collagen forms D-period banded fibrils, with a potent ability to limit the lateral growth.

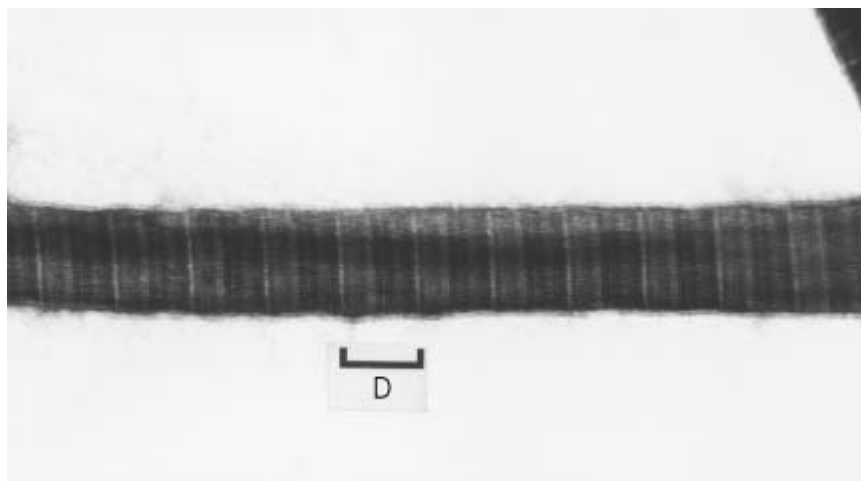
**Figure 7.** Reconstituted type V collagen fibrils. Purified pepsin-treated type V collagen from human placenta with the composition of  $[a1(V)]_2a2(V)$  forms the banded fibrils with finite thickness. The periodicity of the banding pattern is approximately 65 nm. Fibril thickness is 35 nm in average.



### 3.2. D-Staggered Lateral Association

The fibrillar structures with banded patterns are accounted for by the specific parallel and mutually staggered alignment of the triple-helical domains of the fibrillar collagens (Fig. 8). The sequence of hydrophobic and charged residues distributed along the triple helix provide maximum electrostatic and hydrophobic interactions between the neighboring molecules. The reconstituted fibrils exhibit a cross-striated banding pattern in the electron microscope, demonstrating conclusively that the aggregation of collagen molecules into axially ordered fibrillar structures is basically a [self-assembly](#) process, where the information for association is contained within the assembling molecules themselves.

**Figure 8.** D-staggered array of type I collagen in banded fibrils.



It has been recognized for over 30 years that the cross-striated periodic structure of the native collagen fibril is a consequence of the assembly of molecules in a parallel array, but mutually staggered (ie, axially displaced with respect to one another) by approximately one-quarter of their length—often, which is referred to as the *quarter-staggered array*. The periodicity of the cross-striations in the fibril is explained by the fact that each collagen monomer has eight highly charged regions 67 nm apart that appear under appropriate conditions as stained bands.

The banded period (the D period) is confirmed by low-angle [X-ray scattering](#) of rat-tail tendon fibrils, and the overall length of a fibrillar collagen monomer is 4.4 D units (when D = 234 amino acids, the length of one cross-striation period is 67 nm). Within a collagen fibril, the molecules are staggered by integral multiples of the distance D (Fig. [8](#)).

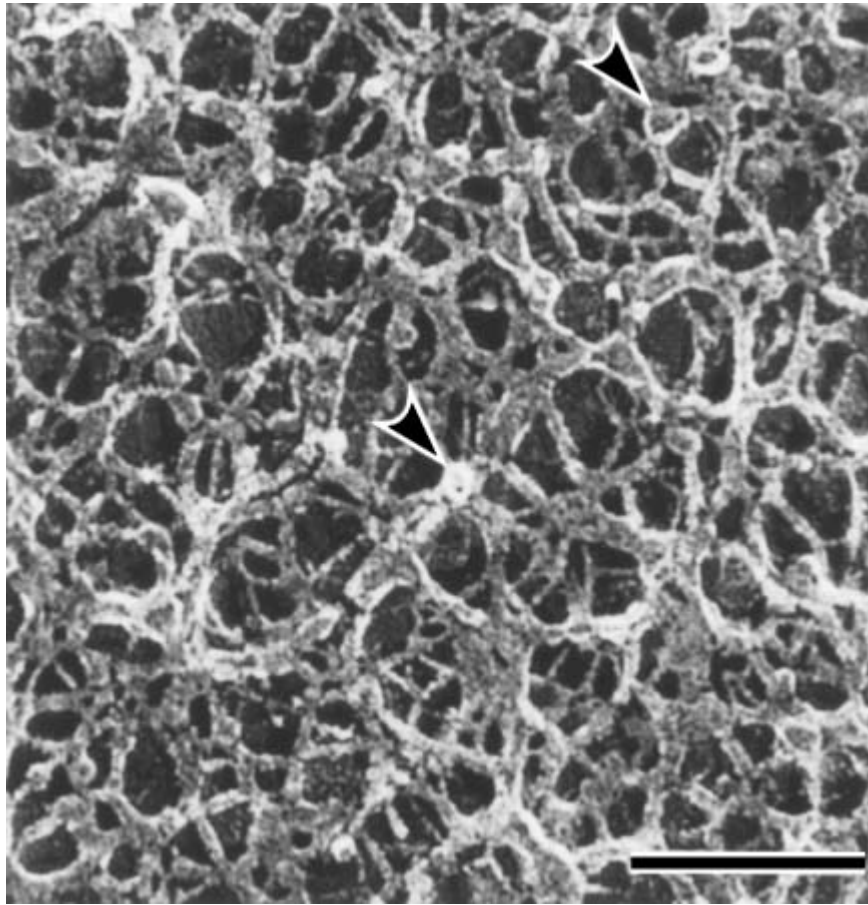
### 3.3. Collagen Supramolecular Aggregates in Tissues

The organization and arrangement of the collagen fibrils and polygonal meshwork in tissues are subject to considerable variation. Fibrillar collagen, type I collagen in particular, provides the major mechanical strength of skin, tendon, bone, dentine, cornea, sclera, and so on. The fibrillar collagen types form long, unbranched, banded fibrils with a characteristic periodicity of 60–70 nm. The fibrillar aggregate of collagen is classified as a quasi-crystal, because of its highly symmetrical insoluble structure built up of essentially identical subunits. Especially in tendon, type I collagen accounts for approximately 90–95% of the dry weight of the tissue and is present as large, highly orientated fibers. All the fibrillar units are arranged in large parallel bundles, with the average diameter of the fibrils varying between 50 and 500 nm. In dermis, where type I collagen accounts for 80–90% of the collagenous proteins, the fibrils form a coarse network partially oriented in the plane of the skin, with the average diameter of the fibrils between 40 and 100 nm. The orthogonally arranged and precisely packed fibrils in corneal stroma are of particular interest, since the arrangement may well be related to tissue transparency. In the cornea, which has a high content of type V and type I collagens, the collagen fibrils have a uniform diameter of approximately 25 nm. Concentric circles of collagen fibrils in cortical bone provides another type of organization and arrangement.

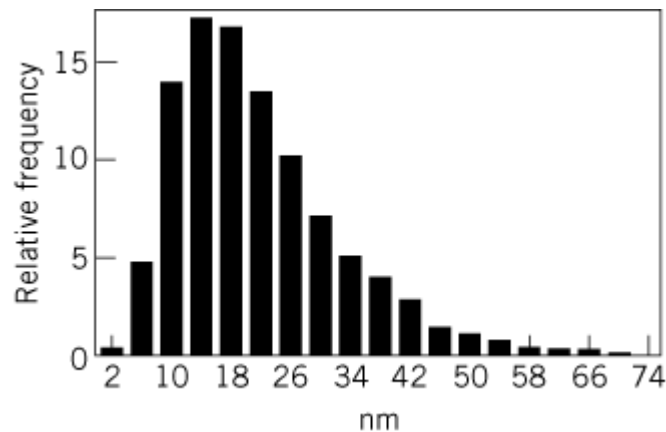
The skeletal structure of the basal lamina has a special three-dimensional meshwork with an average pore size of about 18 nm. Isolated type IV collagen can reconstitute into polygonal meshwork (Fig. [9](#)) with an average pore size also of about 18 nm (Fig. [10](#)). This suggests that the self-assembled supramolecular structure of the type IV collagen is formed primarily through the lateral interactions of the triple-helical domain, which has kinks or bending points due to interruptions of the Gly—X—Y repeat.

**Figure 9.** Polygonal meshwork structure of type IV collagen aggregates. Bar in lower right corner is equal to 100 nm. Arrowheads indicate globules presumably formed through NCI domains.





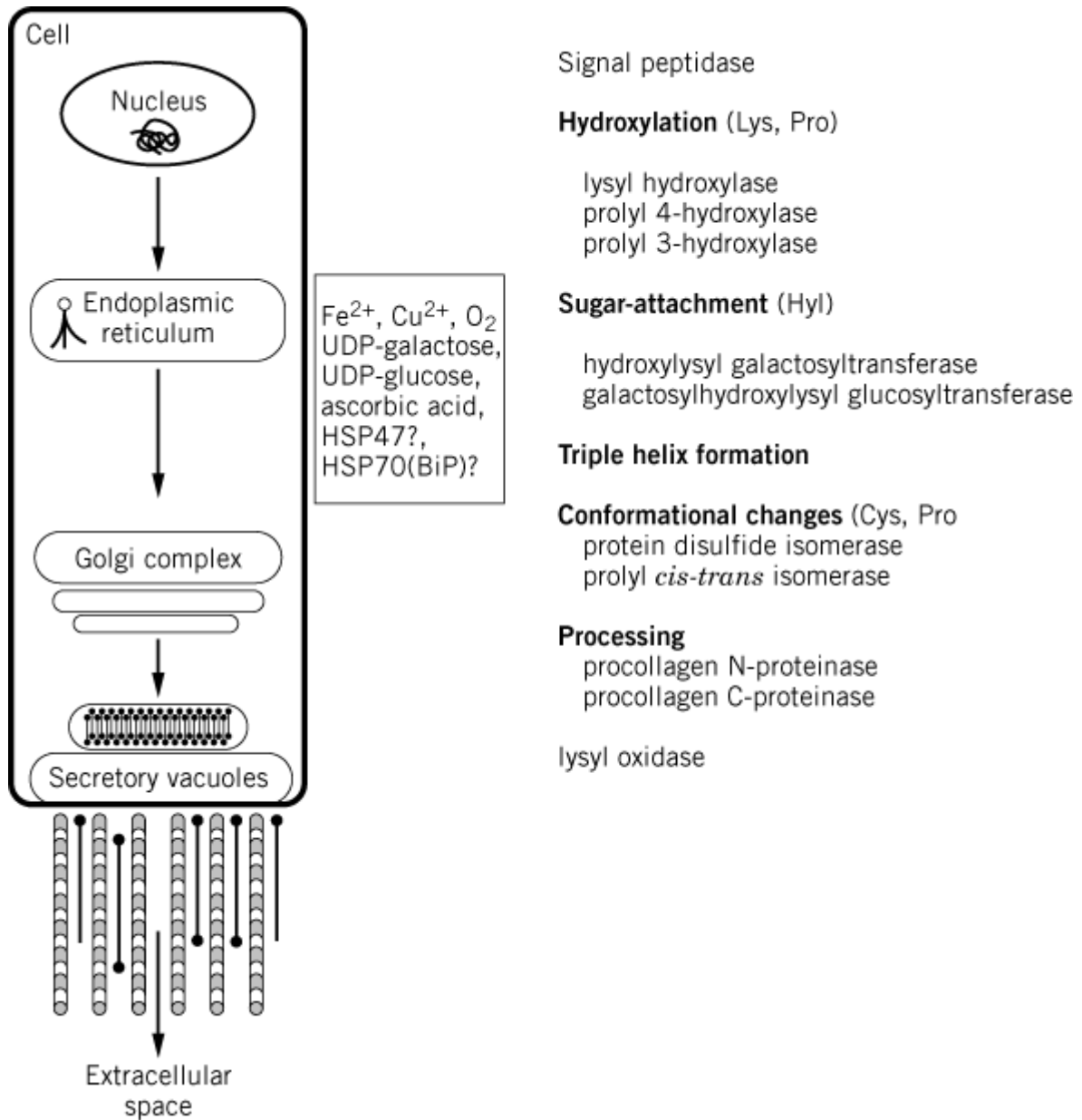
**Figure 10.** Histograms of the distances between two branching points of polygonal meshwork of type IV collagen aggregates.



#### 4. Biosynthesis and Chain Assembly of Collagen

Biosynthesis of the collagen molecule is a complex process that begins with [transcription](#) of the individual collagen **genes** and culminates in the chains assembled into collagen molecules (Fig. [11](#)). This biosynthetic process involves a number of co- and [posttranslational modifications](#) unique to collagenous sequences. Intracellular modifications of the newly synthesized polypeptide chains result in the formation of triple-helical molecules. At least 10 different enzymes have been implicated in the posttranslational processing of the collagen molecule. The biosynthesis of type I collagen can be regarded as a useful model exemplifying many of the common features of collagen biosynthesis. The nonfibrillar collagens may, however, deviate from the general scheme.

**Figure 11.** Schematic drawing of the complex process of biosynthesis of collagen. Collagen translation products are directed into the endoplasmic reticulum by the *N*-terminal signal peptide that is cleaved from the polypeptide chain by signal peptidase. Individual procollagen chains undergo complex enzyme-catalyzed post-translational modifications prior to the completion of chain assembly and the folding of the triple helix. Hydroxylation of Lys and Pro and attachment of sugars to hydroxylysine (Hyl) are characteristic post-translational modifications of collagenous polypeptides.



#### 4.1. Pre-procollagen mRNA Translation

The pre-procollagen chains that are the primary [translation](#) product contain hydrophobic *N*-terminal **signal sequences** of 22–26 amino acid residues, similar to those in most other secreted proteins and are essential for targeting of nascent polypeptides to the [endoplasmic reticulum](#) (see [Protein Secretion](#)). During or shortly after translocation, the signal sequence is removed by the [signal peptidase](#) on the luminal side of the endoplasmic reticulum membrane.

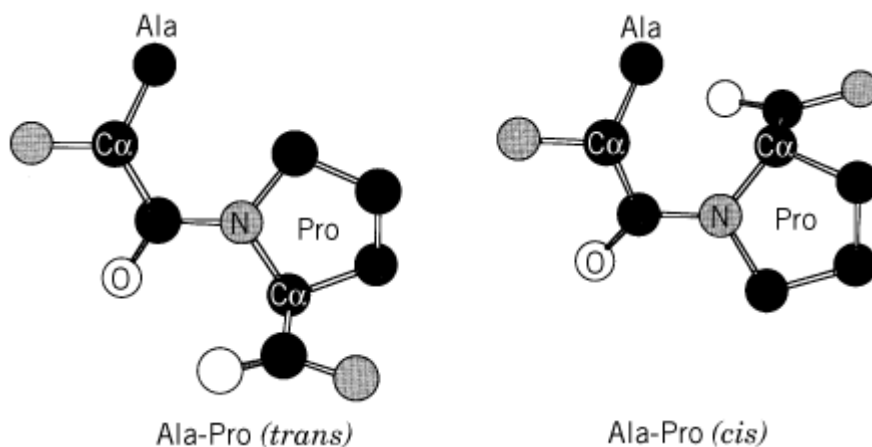
[Hydroxylation](#) of proline and lysine residues of the newly synthesized collagen polypeptide chains results in the formation of procollagen molecules containing an array of hydroxyproline and hydroxylysine residues. This hydroxylation is produced by collagen-specific enzymes within the cell

(see [Hydroxylation \(Lysine, Proline\)](#)).

#### 4.2. Helix Formation and Secretion

The assembly of three polypeptides into a triple-helical collagen molecule is a complex process initiated by association of the three *C*-terminal **propeptides** (in the case of type I collagen), which involves chain alignment, nucleation, and propagation into the triple-helical conformation. Requiring the *C* terminus of the polypeptide chain, chain association and alignment (in the case of type I collagen) takes place only as the chains near full elongation, and folding may proceed in the *C*- to *N*-terminal direction. Association of the *C* termini is stabilized by disulfide bond formation catalyzed by the enzyme [protein disulfide isomerase](#) (PDI). Assembly of the triple helical molecule might be limited by [cis–trans isomerization](#) (Fig. 12) of the peptide bonds preceding proline residues, which is catalyzed by a cytosolic enzyme, [peptidylprolyl cis–trans isomerase](#) (PPI). Formation of the triple-helical conformation appears to preclude further enzymatic post-translational modifications. The procollagen in a triple-helical conformation is ready to be secreted into the extracellular space.

**Figure 12.** *Cis–trans* isomerization of Ala-Pro peptide bonds. H atoms are not shown.



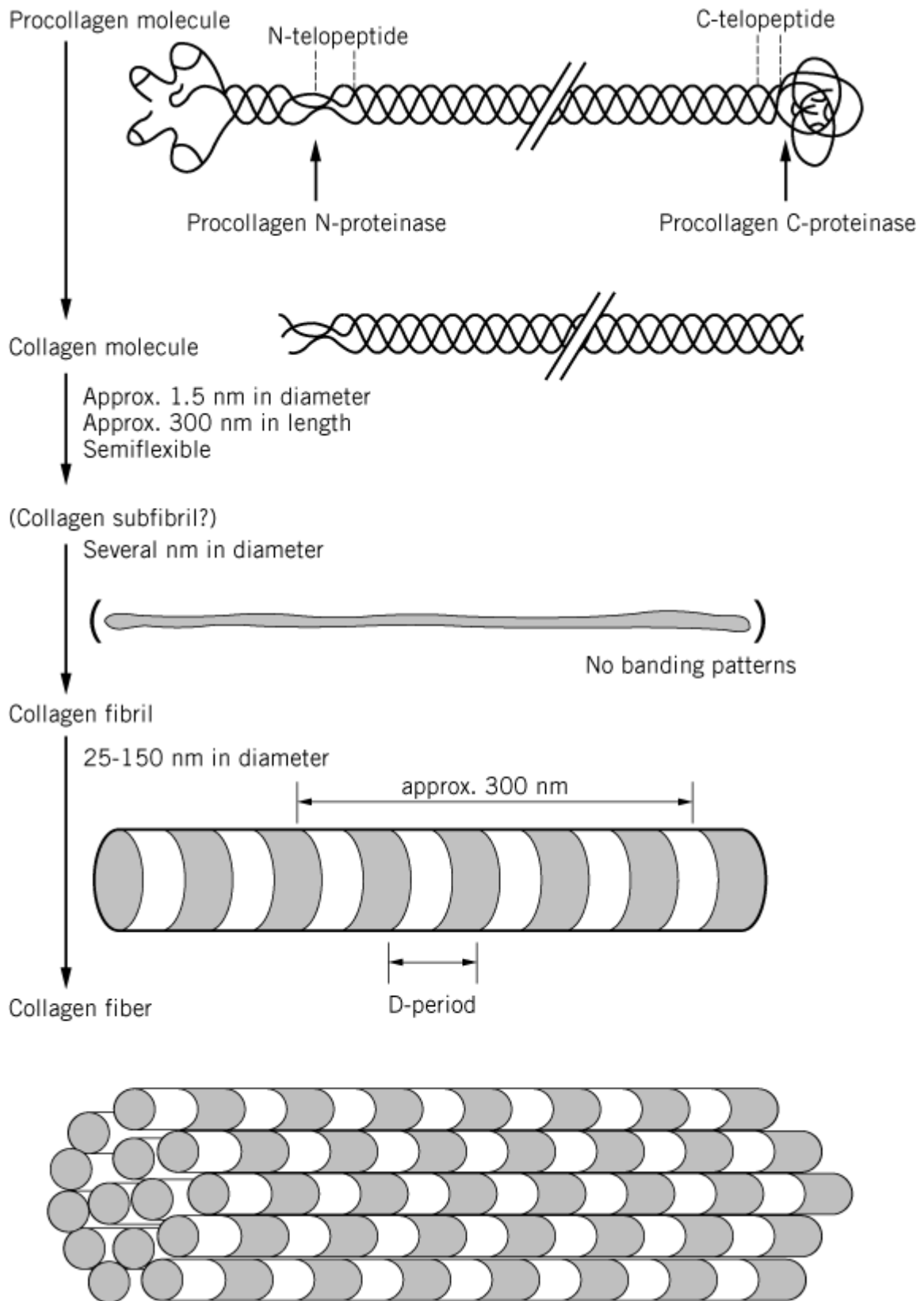
Stable triple helix formation requires 4-hydroxyproline in the Y position of a high proportion of —Gly—X—Y— triplets (at least 100 of the 1000 amino acids of the helical domain of type I collagen) (see [Hydroxylation \(Lysine, Proline\)](#)). The normal rate of procollagen secretion is apparently dependent on triple-helical conformation. If triple helix formation is prevented by inhibiting prolyl hydroxylase (by the use of *a,a'*-dipyridyl or by limiting the availability of the cofactors  $\text{Fe}^{2+}$  or  $\text{O}_2$ ), the nonhelical proa(I) chains first accumulate within endoplasmic reticulum.

#### 4.3. Extracellular Polymerization

Fibril formation involves the specific enzymatic cleavage of the procollagen *N*- and *C*-terminal propeptides in the extracellular space (Fig. 13), and these enzymes must therefore play a key role in regulating the aggregation. Such processing is characteristic of the fibrillar collagens and may not be a universal prerequisite for the assembly of all collagen types. The cleavage of the *N*- and *C*-terminal propeptides from the procollagen I molecule at specific peptide bond cleavage sites is achieved by two specific neutral [metalloproteinases](#)—procollagen *N*-proteinase and procollagen *C*-proteinase—both of which require  $\text{Ca}^{2+}$  for activity and are inhibited by metal chelators. These enzymes act essentially only on molecules in the triple-helical conformation. Type I procollagen *N*-proteinase cleaves the *N*-terminal propeptides of types I and II procollagens between a proline and a glutamine residue. A separate enzyme is involved in the processing of the *N*-propeptide of type III procollagen. The type I procollagen *C*-terminal proteinase cleaves the *C*-terminal propeptides of both the proa1(I) and proa2(I) procollagen chains at an Ala—Asp bond. Most recently, procollagen *C*-proteinase was found to be BMP-1 (bone morphogenetic protein-1), and recombinant BMP-1 was shown to act as procollagen *C*-proteinase. The enzyme is also able to cleave the *C*-terminal propeptides of type I

procollagen homotrimer, type II procollagen, type III procollagen, and other proteins, including laminin-5 and prolyl oxidase.

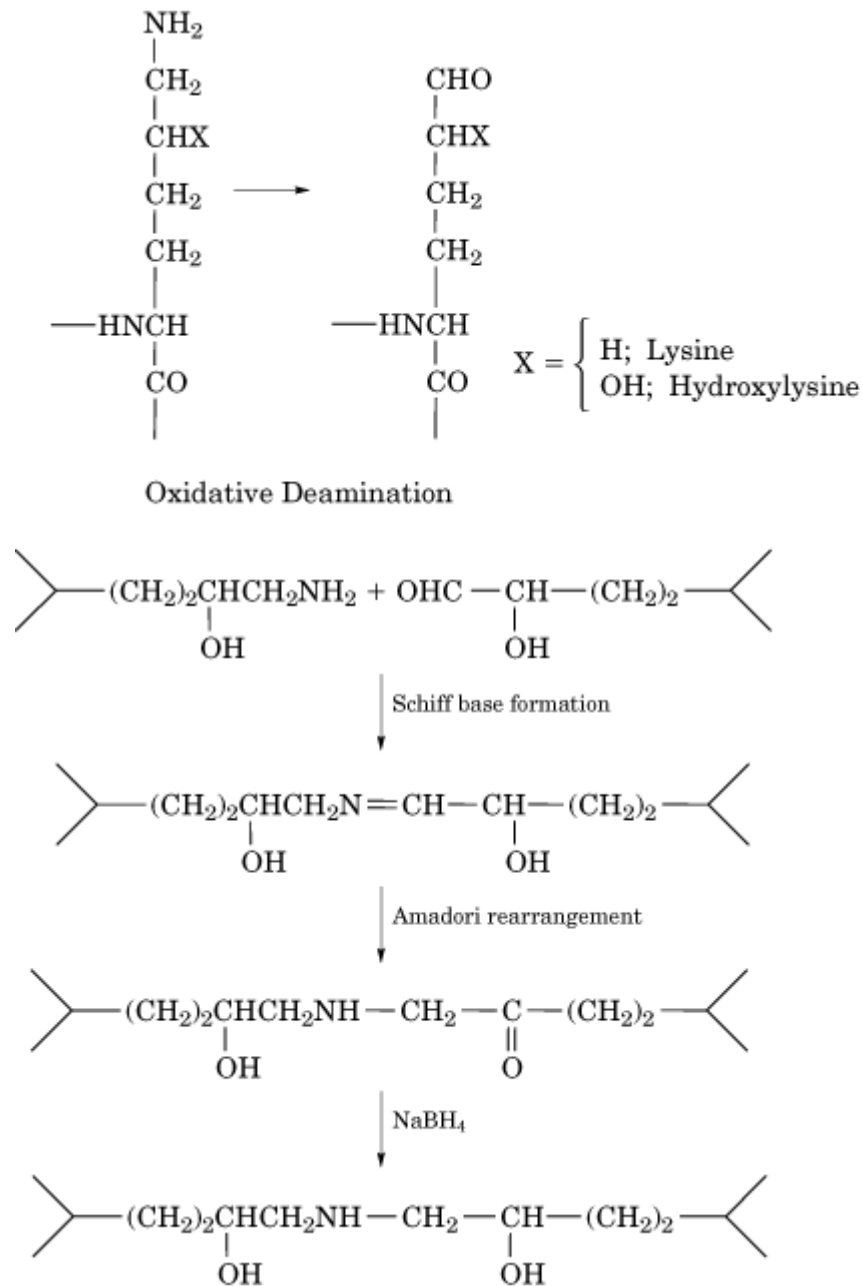
Figure 13. Processing of type I procollagen to fibril formation.



#### 4.4. Crosslinking of Type I Collagen in the Fibrils

Covalent crosslinking of type I collagen in the fibrils involves oxidative deamination of specific lysine or hydroxylysine residues in the nonhelical (telopeptide) regions at the ends of the triple helix. Thus, lysyl and hydroxylysyl aldehydes are generated by lysyl oxidase (protein lysine 6-oxidase) (Fig. 14). The catalytic activity of lysyl oxidase is dependent on strict steric requirements (in the case of type I collagen, the quarter-stagger arrangement of molecules in the fibril) and on the sequence of amino acids surrounding the “target” lysyl/hydroxylysyl residues. The recombinant enzyme was produced as a proenzyme, which could be processed for activation by BMP-1 or procollagen C-proteinase.

**Figure 14.** Oxidative deamination and crosslinking of lysine or hydroxylysine residues. The oxidative deamination, Schiff base formation, and Amadori rearrangements take place *in vivo*. The last step, treatment with sodium borohydride, is carried out in the laboratory to stabilize the crosslinks for chemical analysis.



The majority of covalent crosslinks stabilizing the fibrillar collagens involve the telopeptides and, consequently, pepsin treatment of insoluble, crosslinked fibrils tends to release the triple-helical domain, which can be recovered in its native conformation.

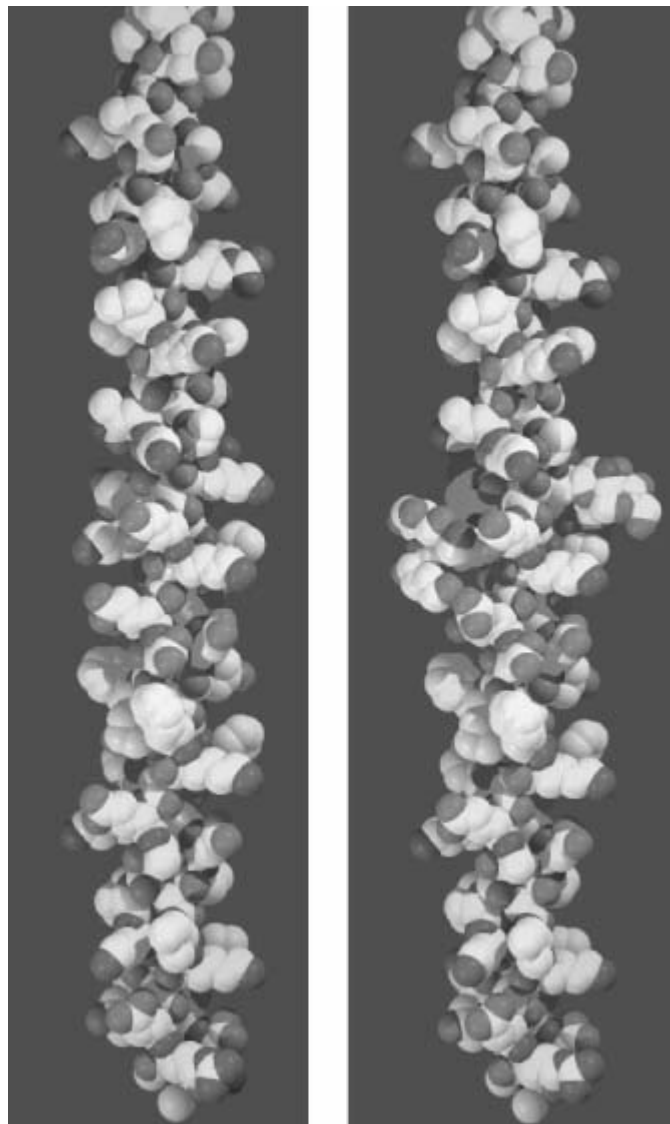
## 5. Posttranslational Modifications of Proline and Lysine Residues

Hydroxylation of proline and lysine residues is a post-translational modification that is specific to collagen (see [Hydroxylation \(Lysine, Proline\)](#)). The reason for hydroxylation of proline residues at the 3 position is not known, but the presence of an appropriate number of 4-hydroxyproline residues in the —Y— position of —Gly—X—Y— is a crucial determinant of the stability of the triple helix under physiological conditions. The hydroxyl group forms hydrogen bonds that contribute to the thermal stability. The transition temperature (thermal denaturation temperature; melting temperature,  $T_m$ ) of triple helices formed *in vitro* by nonhydroxylated collagen is about 25°C, which is 15°C lower than the fully hydroxylated collagen triple helix (Fig. 3). The importance of hydroxylysine may be inferred from the inherited connective tissue disorder, Ehlers–Danlos syndrome type VI, in which lysyl hydroxylase deficiency results in impaired crosslink formation and consequent susceptibility to mechanical disruption of tissues.

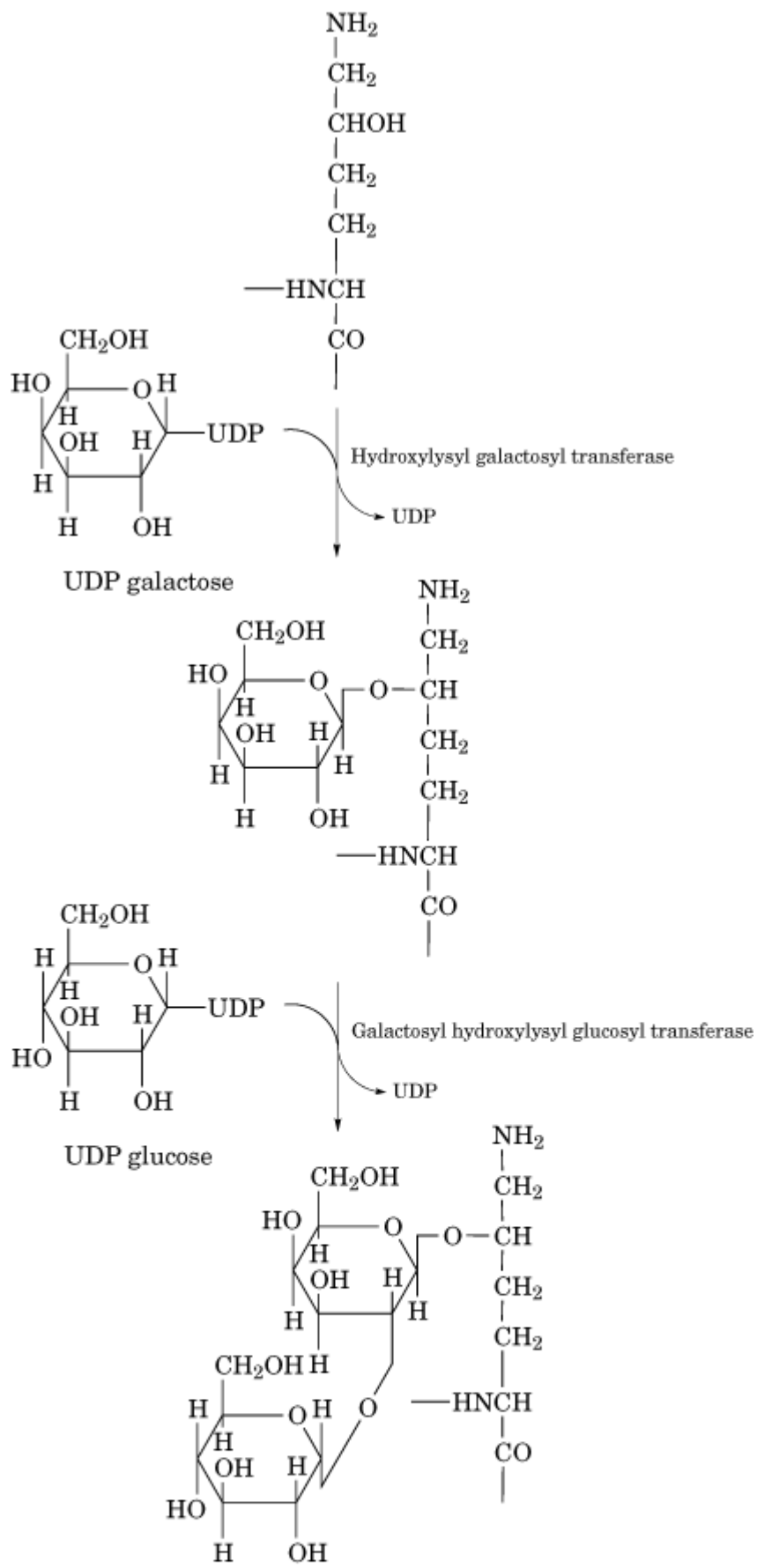
### 5.1. Glycosylation of Hydroxylysine Residues

Vertebrate collagen molecules contain the monosaccharide galactose and the disaccharide glucosyl–galactose. The glycosides are covalently linked through the hydroxyl group to hydroxylysine residues within the triple-helical domains (see Figs. 15 and 16).

**Figure 15.** Glycosylated hydroxylysine residues on the surface of the collagen triple helix. The molecular model of a portion (36 residues/chain) of human type V collagen, [a1(V)]<sub>2</sub>a2(V), was built up by the McMolw (Molecular Images Software, San Diego, CA, USA) and Chem3D (CambridgeSoft Corporation, Cambridge, MA, USA) programs, using the 2CLG file of the Protein Data Bank, which is the theoretical model of the (Gly—Pro—Hyp)<sub>4</sub> trimer. The two lysine residues of the two a1(V) chains are either left unchanged (*left*) or changed to glucosylgalactosylhydroxylysine (*right*). The glycosylated hydroxylysine residues are evidently bulkier than other residues.



**Figure 16.** Attachment of galactose and glucose groups to hydroxylysine residues of collagen. (a) Enzymatic steps in coupling the two sugars; (b) stereo view (direct) of glucosylgalactosylhydroxylysine. H atoms are not shown.





This glycosylation is unique to collagen, and its extent is variable both between collagen types and within the same collagen in different tissues and at different ages. A potential biological role of the specific carbohydrate units would be the regulation of the lateral assembly of the triple-helical domains. An inverse relationship exists between carbohydrate content and collagen fibril diameter, suggesting that the steric hindrance to the formation of highly ordered fibrils might be caused by the bulky glycosylated hydroxylysine residues. The glycosylation of hydroxylysine residues within the triple helix proceeds through the action of two isolated enzymes: hydroxylysyl galactosyl transferase (UDP galactose: 5-hydroxylysine–collagen galactosyltransferase) and galactosyl hydroxylysyl glucosyl transferase (UDP glucose: 5-hydroxylysine-collagen glucosyltransferase). Galactose is transferred to hydroxylysine residues by the former enzyme, while glucose is transferred to galactosyl-hydroxylysine residues by the latter enzyme. The carbohydrates are provided in both reactions by the appropriate UDP-glycosides. Both enzymes require a divalent cation (preferably  $Mn^{2+}$ ) for the activity. A free  $\epsilon$ -amino group in the substrate hydroxylysyl residue and a nonhelical polypeptide conformation are requirements for both transferases. The reactions are more favorable with longer peptides. Hydroxylysyl galactosyl transferase can bind at least two manganese ions per molecule of active enzyme.

#### Suggestions for Further Reading

B. Brodsky (1995) *FASEB J.* **9**, 1537–1546.

E. Adachi, I. Hopkinson, and T. Hayashi (1997) Basement-membrane stromal relationships: interactions between collagen fibrils and the lamina densa. *Intl. Rev. Cytolo.* **173**, 73–156.

P. Yurchenco, D. E. Birk, and R. P. Mecham, eds. (1994) *Extracellular Matrix Assembly and Structure*, Academic Press, San Diego.

J. F. Bateman, S. R. Lamandé, and J. A. M. Ramshaw (1996) "Collagen superfamily", in *Extracellular Matrix*, Vol. **2**, W. D. Comper, ed., Harwood Academic Publishers, Amsterdam.

C. M. Kielty, I. Hopkinson, and M. E. Grant (1992) "Collagen", in *Connective Tissue and Its Heritable Disorders*, R. M. Royce and B. Steinmann, eds., Wiley, New York.

## Colony-Stimulating Factors

Colony-stimulating factors (CSFs) are **glycoproteins** that stimulate the proliferation, differentiation, and survival of hematopoietic (blood) cells and also activate mature myeloid cell functions. CSFs were first described in the 1960s as activities in various conditioned media that stimulated the growth of colonies of bone marrow cells in semisolid cultures. These colony-forming cells are the precursors of our red and white blood cells (see [Hematopoiesis](#)). Purification and characterization of these activities were difficult with the biochemical techniques available because they were present in conditioned media in very low concentrations and also because they were heterogeneous, largely due to variable glycosylation. Eventually, in the late 1970s and 1980s, four different murine CSFs were purified, followed by their human counterparts. These were called

1. Multi-CSF or interleukin-3 (IL-3)
2. Granulocyte-macrophage (GM)-CSF
3. Granulocyte (G)-CSF
4. Macrophage (M)-CSF, or CSF-1

Since then, the development of molecular biology has led to the [cloning](#) of CSF [complementary DNAs](#) and the production of purified [recombinant proteins](#), facilitating *in vitro* and *in vivo* studies of their activities. All four CSFs have been tested in clinical trials, and both G-CSF and GM-CSF are used in clinical practice to elevate leukocyte levels in a variety of situations. In addition to these CSFs, many [interleukins](#) and [growth factors](#) have been described that stimulate or co-stimulate bone marrow cell colony growth. The most important of these are the early-acting stem cell factor (SCF) and FLT3/FLK-2 ligand (FL) (1), the erythroid lineage stimulator erythropoietin, and the megakaryocyte and platelet stimulator thrombopoietin (TPO). Recently, TPO has been found to have multilineage activities (2).

## 1. Biochemical Characteristics

The CSFs belong to a family of **cytokines** that have little sequence identity but were predicted to be similar in structure (3). The [protein structures](#) of all four have now been determined, although for IL-3 a mutant molecule with enhanced stability was used, rather than the native molecule (4). They are all **four-helix bundles** with up–up–down–down connectivity. IL-3, G-CSF, and GM-CSF are monomeric secreted proteins, but the soluble form of M-CSF is a **disulfide-bonded** homodimer produced by **proteolytic** cleavage of cell-membrane molecules. Other characteristics are given in Table 1.

**Table 1. Biochemical Characteristics of the CSFs**

| Property                            | IL-3  |      | GM-CSF |       | G-CSF |     | M-CSF                    |                    |
|-------------------------------------|-------|------|--------|-------|-------|-----|--------------------------|--------------------|
|                                     | Hu    | Mu   | Hu     | Mu    | Hu    | Mu  | Hu                       | Mu                 |
| Amino acids                         | 133   | 140  | 127    | 124   | 174   | 178 | 224/406/522 <sup>a</sup> | 519                |
| $M_r$ (kDa)                         |       |      |        |       |       |     |                          |                    |
| Predicted                           | 15.1  | 15.7 | 16.3   | 16    | 19    | 19  | 60                       | 61                 |
| Expressed                           | 14–30 | 28   | 18–33  | 18–33 | 22    | 25  | 45–90 <sup>b</sup>       | 45–90 <sup>b</sup> |
| Glycosylation                       |       |      |        |       |       |     |                          |                    |
| N-linked                            | 2     | 4    | 2      | 2     | 0     | 0   | 4                        | 3                  |
| O-linked                            | ND    | ND   | 2      | ND    | 1     | 1   | 1                        | 1                  |
| Disulfides                          | 1     | 2    | 2      | 2     | 2     | 2   | 7–9                      | 7–9                |
| a-Helix type                        | Short |      | Short  |       | Long  |     | Short                    |                    |
| Cross-species activity <sup>c</sup> | –     | –    | –      | –     | +     | +   | +                        | –                  |

<sup>a</sup> Different forms resulting from alternative splicing.

<sup>b</sup> Disulfide-bonded dimer of proteolytically cleaved transmembrane protein.

<sup>c</sup> Between human and murine. Data derived from Refs. 4, 18, and 19.

## 2. CSF Genes and Their Expression

The cloning of **genomic** DNA encoding the CSFs and analysis of the regions 5' of the genes has resulted in a better understanding of the control of CSF expression (see Table 2). A single **messenger RNA** species is transcribed for IL-3 and GM-CSF, whereas a minor **alternatively spliced** mRNA is produced for human G-CSF that results in a three-amino-acid-residue insertion in the protein after residue Leu35 and a 20-fold reduction in activity. Transcription of M-CSF mRNA is more complex, with four species resulting from alternative splicing of the 10 exons (see **Introns, Exons**). The largest two differ only in the 3' untranslated region and encode identical 522-residue **transmembrane** proteins. The smaller two result from splicing out of part of exon 6, which results in loss of the proteolytic cleavage site from the protein produced from the 1.6-kb **transcript** and expression of a stable transmembrane protein (5). The DNA regions upstream of the coding regions contain **promoter** and **enhancer** elements that have been partially characterized and shown to respond to CSF inducers such as **tumor necrosis factor** (TNF), IL-1, lipopolysaccharide (LPS), and **antigens**. The control of CSF gene expression is not understood completely, but it involves control of mRNA stability by elements in the 3' untranslated region of the mRNAs, as well as control of transcription.

**Table 2. Characteristics of the Genes Encoding the CSFs**

| Property             | IL-3   |         | GM-CSF                            |         | G-CSF                                       |        | M-CSF                                       |             |
|----------------------|--|---------|-----------------------------------|---------|---|--------|---|-------------|
|                      | Hu   | Mu      | Hu                                | Mu      | Hu  | Mu     | Hu  | Mu          |
| Size (kbp)           | 2.2  | 2.7     | 2.5                               | 2.5     | 2.5   | 2.5    | 20  | 20          |
| Chromosomal location | 5q23-31  | 11A5-B1 | 5q21-31                           | 11A5-B1 | 17q11-22                                    | 11D-E2 | 1p13-21                                     | 3F3         |
| Number of exons      | 5  | 5       | 4                                 | 4       | 5   | 5      | 10  | 10          |
| Alternative splicing | No   | No      | No                                | No      | Yes   | No     | Yes   | Yes         |
| mRNA species (kb)    | 0.9  | 0.9     | 0.78                              | 0.78    | 1.5   | 1.5    | 4, 2.3, 2.0, 1.6                            | 4, 2.6, 1.4 |
| Promoter elements    | AP-1, ELF-1, NIP <sup>a</sup> , NF-IL-3, CRE, CK-1, CK-2, GC |         | CK-1, CK-2, IgbB, GC, AP-1, ELF-1 |         | PU box, CK-1, NF-IL-6, octamer, GPE-3, CK-2 |        | SP-1, AP-1, AP-2, NF-1, NF-IL-6, CK-1, CK-2 |             |

<sup>a</sup> Repressor element. Data from Refs. 18 and 19.

### 3. CSF Production

Many cell types have been shown to produce CSFs *in vitro* in response to various stimuli, such as bacterial products (LPS) and inflammatory cytokines (TNF, IL-1), but the role of CSFs in maintaining white blood cell levels in normal healthy individuals is not entirely clear (see later). Whether unstimulated cells produce CSFs is still debated (6). GM-CSF, G-CSF, and M-CSF are produced by macrophages, endothelial cells, fibroblasts, and stromal cells. In addition, T lymphocytes produce GM-CSF, M-CSF (in humans), and IL-3. Mast cells have been shown to

produce IL-3 and GM-CSF. Other cell types have been reported to produce GM-CSF and IL-3, but these studies have yet to be confirmed.

#### 4. CSF Receptors

The CSF receptors are transmembrane proteins belonging to the cytokine receptor class I family, except for the M-CSF receptor, which is a [tyrosine kinase receptor](#). The class I receptors contain a structurally conserved module called the cytokine-binding **domain** (CBD), the hemopoietin domain, or the cytokine receptor homology domain. This module comprises two [fibronectin](#) type III domains with four conserved [cysteine](#) residues in the *N*-terminal domain and a conserved Trp—Ser—X—Trp—Ser sequence in the *C*-terminal domain. Ligand binding causes dimerization or oligomerization of the receptors, which initiates [signal transduction](#). The IL-3 and GM-CSF receptors form heterodimers of a ligand-specific  $\alpha$ -chain and a common  $\beta$ -chain that is also shared with the IL-5 receptor. In the mouse, there is an additional  $\beta$ -chain that is unique to the IL-3 receptor. The G-CSF and M-CSF receptors form homodimers. Receptor dimerization activates tyrosine kinases (intrinsic in the M-CSF receptor, cytoplasmic **JAK** kinases for the other receptors) that **phosphorylate** the receptors, providing binding sites for other signaling molecules that are in turn activated by phosphorylation of [tyrosine](#) residues and, in some cases, [serine](#) or [threonine](#). Ultimately, these signaling cascades lead to alterations in gene transcription, but the components of the biochemical pathways are still being defined (reviewed in Refs. [7-9](#)).

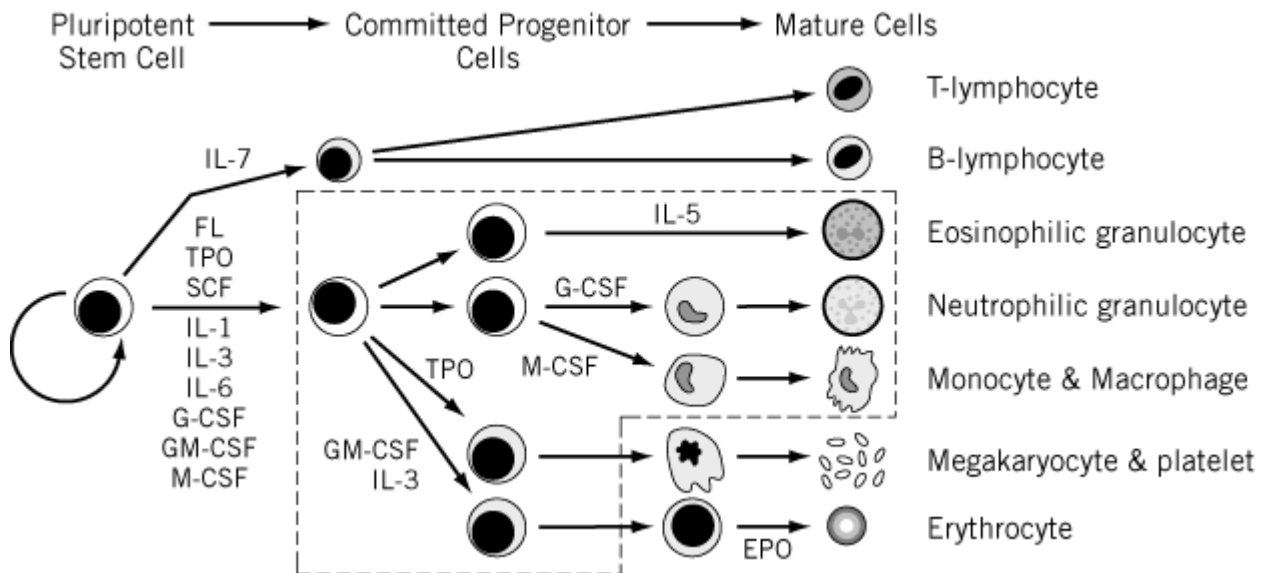
The CSF receptors are expressed at low levels on hematopoietic cells, usually a few hundred receptors per cell. IL-3 receptors are expressed on neutrophil, eosinophil, and macrophage lineages, on mast cells, and on undifferentiated blast cells. GM-CSF receptors are expressed on neutrophil, eosinophil, and macrophage lineages. G- and M-CSF receptor expression is more restricted; G-CSF receptors are present on neutrophils and at a low level on macrophages, and M-CSF receptors are found on macrophages, osteoclasts, and placental trophoblasts. Recent studies have detected expression of receptors on purified stem-cell populations. Human CD34<sup>+</sup> cells expressed SCF, FL, G-CSF, and, at a low level, GM-CSF receptors ([10](#)). Murine stem cell populations expressed IL-1 $\alpha$ , IL-3, IL-6, and G-CSF receptors ([11](#)).

#### 5. Biological Activity

##### 5.1. *In Vitro*

All four of the CSFs stimulate the proliferation and differentiation of bone-marrow progenitors to form colonies in semisolid agar. The type of mature cells in the colonies varies with the stimulator (see Fig. [1](#)). IL-3 was otherwise called multi-CSF because it stimulates the largest number of different cell lineages. Colonies of granulocytes, macrophages, eosinophils, megakaryocytes, undifferentiated blast cells, and mixed cell lineages are obtained. In liquid cultures, proliferation of mast cells is strongly stimulated by IL-3, and some types of B lymphocytes are reported to respond. GM-CSF stimulates growth of granulocytic, macrophage, and eosinophil colonies. Production of dendritic cells, which are important **antigen-presenting** cells in the immune response, is stimulated by GM-CSF. G-CSF and M-CSF are more lineage restricted, producing mainly granulocytic and macrophage colonies, respectively. The CSFs can give enhanced or synergistic responses when used in combination with each other and with other interleukins and growth factors. Some effects of CSFs on nonhematopoietic cells have been reported. G- and GM-CSF can stimulate proliferation of endothelial cells. Placenta-derived cells and oligodendrocytes proliferate in response to GM-CSF. M-CSF stimulates proliferation of trophoblast cells and may be important in pregnancy.

**Figure 1.** Schematic representation of hematopoiesis and the CSFs and other cytokines required for mature cell production. IL, interleukin; FL, FLT3/FLK-2 ligand; SCF, stem cell factor; TPO, thrombopoietin. GM-CSF and IL-3 stimulate cells inside the dotted line box.



## 5.2. *In Vivo*

The effects of CSFs *in vivo* have been determined both by injection in mice and humans and by analysis of CSF-deficient mice created by gene targeting and, in the case of M-CSF, a naturally occurring mutation. Injection of IL-3 or GM-CSF elevates circulating levels of neutrophils, monocytes, and eosinophils. G-CSF causes a greater elevation of neutrophils than does IL-3 or GM-CSF, but only a weak increase in monocytes and no effect on other lineages. M-CSF causes a selective elevation of monocytes. The magnitude of these increases depends on the dose of CSF administered. It is interesting that in response to infection, circulating IL-3 or GM-CSF are not usually detected, whereas G-CSF and, to a lesser extent, M-CSF are frequently elevated in serum. These observations suggest that IL-3 and GM-CSF may be important in local inflammatory responses in the tissues rather than in a systemic response to infection.

The importance of the CSFs in maintaining basal hematopoiesis has been assessed by analysis of CSF-deficient mice. For IL-3, GM-CSF, and G-CSF, gene “knockout” mice were created by introduction of a defective version of the gene of interest into embryonic stem cells by homologous recombination. The embryonic stem cells were injected into blastocysts to create chimeric mice that were then bred to create heterozygous and homozygous deficient mice. For M-CSF, mice with a naturally occurring point mutation (*op/op*) that resulted in a frameshift in the sequence and a loss of M-CSF production were discovered. The IL-3- and GM-CSF-deficient mice were able to produce normal levels of circulating blood cells and bone marrow progenitors. The GM-CSF-deficient mice developed pulmonary alveolar proteinosis and increased pulmonary infections, presumably as a consequence of defective pulmonary macrophage function. Recent studies in IL-3-deficient mice have shown that IL-3 is important for elevation of basophils and mast cells in response to a parasitic infection (12) and for some delayed-type hypersensitivity (T-lymphocyte-dependent) responses (13). G-CSF defective mice were neutropenic, with neutrophil numbers reduced to 20–35% of wild-type control mice. Bone marrow cellularity was normal, but some granulocytic and macrophage progenitors were reduced in number. These mice were more susceptible to a bacterial infection (*Listeria monocytogenes*) than wild-type mice. Young M-CSF-deficient mice have osteopetrosis, a severe macrophage and osteoclast deficiency resulting in failure of teeth eruption and overgrowth of bone that occludes the marrow cavity and reduces marrow cellularity. Unlike the other CSF-deficient mice, the *op/op* mice overcome their defects with age, which suggests that other cytokines may be able to compensate for the lack of M-CSF. Thus, the CSF-deficient mice have shown that each CSF has a unique function *in vivo*. GM-CSF and IL-3 are not required for steady-state hematopoiesis. G- and M-CSF are required for normal neutrophil and macrophage production, respectively. All four CSFs contribute in different ways to control of infection. They may have other nonunique functions

that are not revealed in the deficient mice because of compensatory activities of other cytokines. These may be detected in future studies of multiply deficient mice.

## 6. Myeloid Leukemia

Acute and chronic myeloid leukemias (AMLs and CMLs) are clonal diseases of bone marrow cells in which the predominant cell types are granulocytes and macrophages, although leukemic cells of other lineages are also found. There are potentially two ways in which CSFs could be involved in leukemia: They may be required for the proliferation of the leukemic cells, and they may cause differentiation and clonal extinction of the leukemia. In semisolid agar cultures, CML cells will produce colonies of normal appearance containing granulocytes, macrophages, and eosinophils. The colony growth is dependent on CSF stimulation, and all four CSFs are active. In contrast, AML cells produce small clones with little maturation, and growth is very variable from patient to patient. The cells are responsive to CSFs to a variable extent, and in a proportion of patients (<30%) there is colony growth in the absence of CSFs. The colonies do not contain clonogenic cells that are able to form further colonies in secondary cultures. The leukemic stem cells have not been identified, and their dependence on CSFs is unclear. Whether autocrine production of CSFs by leukemic cells is important is also unknown, with the exception of leukemic monocytes, which produce CSFs like their normal counterparts. Lopez and colleagues have produced a GM-CSF antagonist that induced remission in a murine model of juvenile myelomonocytic leukemia (14). Although there are some leukemic cell lines that differentiate and cease proliferation in response to CSFs *in vitro*, many cell lines are not suppressed by CSFs. There may be defects in the signal transduction pathways that prevent differentiation. In a small number of patients with congenital neutropenia who develop AML, a mutation in the G-CSF receptor expressed by the leukemic cells has been described that prevents differentiation in response to G-CSF (15). However, it is not known whether this mutation contributed to the development of AML.

## 7. Clinical Use of CSFs

The use of chemotherapy in cancer treatment is limited by its toxicity, primarily neutropenia. Both G-CSF and GM-CSF have been tested extensively in clinical trials and shown to reduce neutropenia, allowing the intensification of chemotherapy doses. Clinical trials have not yet shown, however, whether dose intensification improves survival. G-CSF has fewer side effects and is more effective at elevating neutrophil levels than GM-CSF and thus is usually the CSF of choice, but GM-CSF may be more effective in reducing infectious complications (16). In high-dose chemotherapy followed by autologous bone marrow transplantation, both CSFs reduced the time taken for recovery to  $\geq 500$  cells/mm<sup>3</sup> by about one week and, in some trials, reduced the duration of hospital stay. G-CSF (and GM-CSF) treatment increased the levels of circulating progenitor and stem cells so that G-CSF-mobilized peripheral blood progenitor cells are now used instead of bone marrow cells for reconstitution after high-dose chemotherapy. This treatment results in more rapid recovery than bone marrow transplantation.

Patients with congenital neutropenia or idiopathic chronic neutropenia respond to G-CSF treatment with elevation of neutrophil levels, resulting in fewer infections and less time in hospital. A few of these patients (3% in Ref. 17) develop myelodysplasia or leukemia, but at present there is no evidence whether G-CSF increases the frequency of these cases. Despite initial concerns that treatment of leukemia patients with CSFs might exacerbate the disease, there is no evidence that this occurs. Neutrophil recovery has occurred more rapidly in leukemia patients receiving G-CSF or GM-CSF after chemotherapy, but reduction in duration of hospital stay was not usually observed.

There have been fewer clinical trials with M-CSF and IL-3. M-CSF appears to improve recovery from fungal infections and is approved in Japan for use in allogeneic transplantation, induction therapy of AML, and dose-intensive therapy of ovarian cancer. IL-3 given after chemotherapy reduced the recovery time from leukopenia and elevated the levels of neutrophils, monocytes, and

eosinophils. In the future, it is likely that combinations of cytokines will be tested.

## Bibliography

1. Y. Yonemura, H. Ku, S. D. Lyman, and M. Ogawa (1997) *Blood* **89**, 1915–1921.
2. G. P. Solar, W. G. Kerr, F. C. Zeigler, D. Hess, C. Donahue, F. J. de Sauvage, and D. L. Eaton (1998) *Blood* **92**, 4–10.
3. J. F. Bazan (1990) *Immunol. Today* **11**, 350–354.
4. Y. Feng, B. K. Klein, and C. A. McWherter (1996) *J. Mol. Biol.* **259**, 524–541.
5. M. Baccarini and E. R. Stanley (1990) In *Growth Factors, Differentiation Factors, and Cytokines* (A. Habenicht, ed.), Springer-Verlag, Berlin, pp. 188–200.
6. F. H. M. Cluitmans, B. H. J. Esendam, J. E. Landegent, R. Willemze, and J. H. F. Falkenburg (1995) *Blood* **85**, 2038–2044.
7. T. Hara and A. Miyajima (1996) *Stem Cells* **14**, 605–618.
8. B. R. Avalos (1996) *Blood* **88**, 761–777.
9. L. R. Rohrschneider, R. P. Bourette, M. N. Lioubin, P. A. Algate, G. M. Myles, and K. Carlberg (1997) *Mol. Reprod. Dev.* **46**, 96–103 (Abstract).
10. A. W. Roberts, M. Zaiss, A. W. Boyd, and N. A. Nicola (1997) *Exp. Hematol.* **25**, 289–305.
11. W. J. McKinstry, C. Li, J. E. J. Rasko, N. A. Nicola, G. R. Johnson, and D. Metcalf (1997) *Blood* **89**, 65–71.
12. C. S. Lanz, J. Boesiger, C. H. Song, N. Mach, T. Kobayashi, R. C. Mulligan, Y. Nawa, G. Dranoff, and S. J. Galli (1998) *Nature* **392**, 90–93.
13. N. Mach, C. S. Lanz, S. J. Galli, G. Reznikoff, M. Mihm, C. Small, R. Granstein, S. Beissert, M. Sadelain, R. C. Mulligan, and G. Dranoff (1998) *Blood* **91**, 778–783.
14. P. O. Iversen, I. D. Lewis, S. Turczynowicz, H. Hasle, C. Niemeyer, K. Schmiegelow, S. Bastiras, A. Biondi, T. P. Hughes, and A. F. Lopez (1997) *Blood* **90**, 4910–4917.
15. F. Dong, R. K. Brynes, R. K. Tidow, K. Welte, B. Lowenberg, and I. P. Touw (1995) *N. Engl. J. Med.* **333**, 487–493.
16. J. Nemunaitis (1997) *Drugs* **54**, 709–729.
17. D. C. Dale, M. A. Bonilla, M. W. Davis, A. M. Nakanishi, W. P. Hammond, J. Kurtzberg, W. Wang, A. Jakubowski, E. Winton, P. Lalezari, W. Robinson, J. A. Glaspy, S. Emerson, J. Gabilove, M. Vincent, and L. A. Boxer (1993) *Blood* **81**, 2496–2502.
18. D. Metcalf and N. A. Nicola (1995) *The Hemopoietic Colony-Stimulating Factors. From Biology to Clinical Applications*, Cambridge University Press, Cambridge, England.
19. R. Callard and A. Gearing (1994) *The Cytokine Facts Book*, Academic Press, London.

## Suggestions for Further Reading

20. G. J. Lieschke (1997) CSF-deficient mice—what have they taught us? Ciba Foundation Symposium **204**, 60–77. A detailed analysis of CSF deficient mice.
21. D. Metcalf (1998) Lineage commitment and maturation in hematopoietic cells: the case for extrinsic regulation. *Blood* **92**, 345–348.
22. T. Enver, C. M. Heyworth, and T. M. Dexter (1998) Do stem cells play dice? *Blood* **92**, 348–351. Two sides of the debate about hematopoietic lineage commitment.
23. S. S. Watowitch, H. Wu, M. Socolovsky, U. Klingmuller, S. N. Constantinescu, and H. F. Lodish (1996) Cytokine receptor signal transduction and the control of hematopoietic cell development. *Annu. Rev. Cell Dev. Biol.* **12**, 91–128.
24. *References 18 and 19 above. Reference 18 gives a detailed, well referenced description of the four CSFs. Reference 19 is less detailed, but covers other cytokines in addition to the original four.*

## Combinatorial Libraries

Few fields have emerged as rapidly as the technologies collectively referred to as combinatorial chemistry. The origins of combinatorial “thinking” as applied to biological problems can be traced to our understanding of the immune system. Humoral immunity is based on the generation of billions of different [antibody](#) molecules through unique **gene recombination** events. Upon exposure to a foreign [immunogen](#), the subset of the antibody [repertoire](#) with specificity for the invader is clonally expanded. During expansion, the active antibody subpool undergoes additional [recombination](#) and [mutation](#) events, and the most active variants are further selected. This process is reiterated until mature antibody molecules are produced that are highly optimized against the foreign [antigen](#). In essence, antibody repertoires are combinatorial libraries—collections of related but distinct molecules representing different combinations of genetic elements. The recognition of these facts preceded the creation of the first combinatorial library, which was based upon the synthesis and screening of large numbers of synthetic peptides in order to identify selected antigenic determinants of a viral pathogen (1). Given the parallels between the immune system and combinatorial libraries, it is not surprising that the original application of this emerging technology was in the derivation of immunogens capable of stimulating immune protection in animals not yet exposed to certain pathogens.

The hallmarks of a combinatorial approach are the generation and use of complex compound mixtures to identify molecules with desirable properties. Four basic operations underpin combinatorial approaches: (i) generation of libraries containing related but diverse compounds, (ii) selection of molecules with desirable properties from the library, (iii) characterization of the selected molecules, and (iv) amplification of the identified products. Depending upon the chemical composition of the library, the selected molecules may require amplification before structure determination. For example, [DNA libraries](#) can be amplified using the polymerase chain reaction (PCR) before the final products are identified by sequencing. Synthetic peptides and small molecules cannot be amplified per se and must be subjected to structural elucidation, and then resynthesized for subsequent use. Regardless of the library composition, the four operations enumerated above can be applied in an iterative fashion to progressively refine the selected compounds.

Several entries in this volume deal with the application of combinatorial approaches to complex biological problems. Issues central to the implementation of combinatorial technologies in biological systems will be described below, including cognitive aspects of library design, practical aspects of library construction, basic issues in affinity selection and library screening, and discussions of different library types and their applications. The breadth of the topic Combinatorial Libraries is so expansive that it is impossible to cover any topic in great detail. The entries entitled [Libraries](#), [Combinatorial Synthesis](#), and [Affinity Selection](#) are intended to provide a generic introduction to the principles underlying combinatorial operations. The remaining entries, [DNA libraries](#), [Genomic libraries](#), [cDNA libraries](#), [Expression libraries](#), [Peptide libraries](#), and [Phage display libraries](#), provide more specific details regarding the major types of combinatorial libraries in widespread use. Small molecule synthesis or drug discovery applications are not dealt with in detail, although many basic concepts from these fields are presented. The reader is referred to the suggested reading list for additional discussion of these related topics.

See also [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), [Phage Display Libraries](#).

## Bibliography



1. H. M. Geysen, R. H. Meleon, and S. J. Barteling (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3998–4002.

### Suggestions for Further Reading

2. J. D. Watson (1987) *Molecular Biology of the Gene*, 4th ed., Benjamin-Cummings, Menlo Park, CA.
3. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning: A Laboratory Manual*, 2nd ed., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
4. D. J. Kenan, D. E. Tsai, and J. D. Keene (1994) *Trends Biochem. Sci* **19**, 57–64.
5. J. N. Abelson, ed. (1996) *Combinatorial Chemistry*, Vol. **267**, Methods in Enzymology, Academic Press, San Diego.
6. L. A. Thompson and J. A. Ellman (1996) *Chem. Rev.* **96**, 555–600.
7. B. K. Kay, J. Winter, and J. McCafferty (1996) *Phage Display of Peptides and Proteins*, Academic Press, San Diego.

## Combinatorial Synthesis

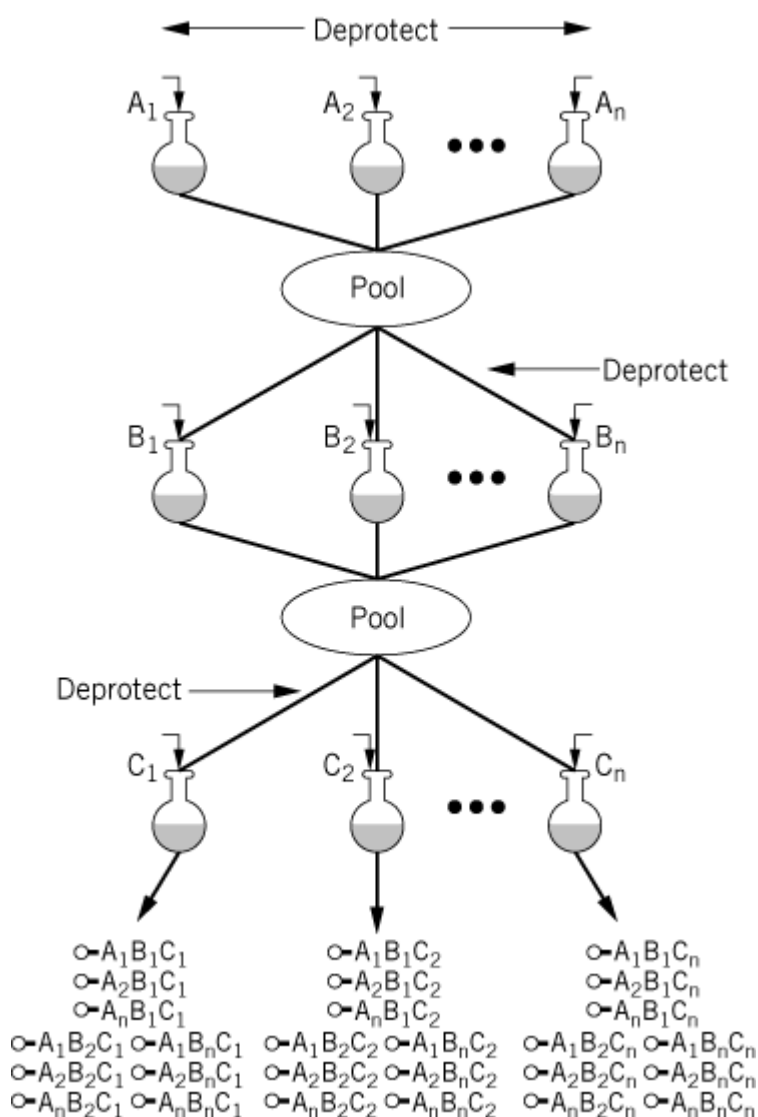
The diversity of a molecular [library](#) may be derived from natural sources, or it may be created through specialized techniques in synthetic chemistry. Traditionally, synthetic chemists focus on preparing individual compounds, paying careful attention to stereochemical control and to achieving the greatest possible yield and purity. This often requires the use of highly specialized reactions and can be a painstakingly slow process. In contrast, chemists engaged in combinatorial organic synthesis employ transformations that are more general in scope and can be applied to the simultaneous preparation of many chemically distinct compounds. In exchange for the ability to prepare many compounds in a short time frame, compromises are accepted that may result in sacrificing, within acceptable limits, yield and purity of individual reaction products. This entry deals with the basic issues faced in preparing and screening synthetic [combinatorial libraries](#) and some of the relevant technologies and methods. No attempt has been made to provide exhaustive coverage of the rapidly expanding range of chemistries now amenable to library synthesis, nor have we attempted to cover the application of combinatorial chemistry as it applies directly to the field of drug discovery. See [Combinatorial Libraries](#).

Combinatorial chemistry seeks to create all possible combinations of a set of building blocks and then extract the identity of members that exhibit a desired property. Thus, if a traditional synthesis involves incorporating building block A, followed by B, and then C to give the desired compound A–B–C, then the analogous combinatorial synthesis would employ a set of monomers, A<sub>1</sub>–A<sub>n</sub>, followed by a second set of monomers, B<sub>1</sub>–B<sub>n</sub>, and then a third set of monomers, C<sub>1</sub>–C<sub>n</sub>, giving all possible combinations of (A<sub>1 to n</sub>)–(B<sub>1 to n</sub>)–(C<sub>1 to n</sub>). Such processes can be described as  $n^k$  sets, where  $n$  is the number of building blocks used in each cycle of synthesis and  $k$  is the total number of synthetic cycles. The number of individual species generated by these approaches become vast very quickly. For example, a relatively short DNA library of 100 random nucleotide positions has a theoretical complexity of  $4^{100}$ , or  $1.6 \times 10^{60}$  different sequences. This number exceeds the number of protons in the sun. Complexities of more realistic peptide libraries using the 20 amino acids incorporated into proteins are as follows: A library of all possible tripeptides will contain 8000 ( $20^3$ ) distinct peptide sequences; a library of all possible tetrapeptides will contain 160,000 ( $20^4$ ) distinct

peptide sequences; and a library of all possible pentapeptides will contain 3,200,000 ( $20^5$ ) distinct peptide sequences. Specialized methods are required to prepare, assay, and elucidate the identity of active compounds from collections of this size. In order to facilitate product purification, compounds are typically synthesized while covalently attached to inert solid supports, including polymeric pins or resin beads, although other materials have been used. Use of solid phase synthesis enables many reactions to be carried out simultaneously using parallel arrays of individual reaction chambers. If each reaction yields a single product, the result is a library of spatially discrete compounds. This approach has the advantage that the identity of compounds can be ascertained immediately, based on the position in the synthetic array and the associated synthetic history. The initial application of this approach was reported by Geysen et al. for synthesis of peptide libraries (1); however, many variations have been reported since then, involving a wide range of chemistries. Although convenient and straightforward to implement, the number of compounds that can be synthesized and screened using this approach is limited by the physical dimensions of the arrays to approximately 10,000 different species. Miniaturization and photolithographic masking techniques have enabled arrays containing more than 100,000 compounds to be synthesized on silica chips (2). Although this has been applied to a variety of chemistries, it is perhaps most widely recognized as the [DNA chip](#) technology of Affymetrix (3). Unfortunately, this approach requires highly specialized instrumentation for the photolithographic masking, as well as for assay of the resulting arrays.

The synthesis of larger libraries requires pooling methods that overcome the size limitations of spatial arrays (4). Pooling methods employ mixtures of reactants at each synthetic step (5). For example, in random DNA oligonucleotide synthesis, two or more deoxyribonucleotide phosphoramidites are introduced during each coupling cycle. One major problem with mixture synthesis is that variation in reactivity can result in some building blocks being under- or overrepresented in the final compound pool. At least in some cases, the initial concentration of reactants can be adjusted to correct for these variations in reactivity. A more attractive approach, first introduced by Furka, is the “split synthesis” strategy (6). Briefly, a portion of synthesis resin is “split” or divided into equal portions and placed into individual reaction vessels (Fig. 1). Each vessel receives a single building block in large excess, which drives the reaction to completion. The resin from all the vessels is then recombined, mixed thoroughly, and then redistributed. A second set of building blocks is then introduced. This process is repeated the desired number of cycles, to produce all possible combinations of the desired building blocks. Because each resin bead is exposed to a single building block at each step, there will be only one compound (but many copies) attached to each bead. This has led to the term one-bead, one-compound (OBOC) libraries (7, 8).

**Figure 1.** Schematic flow of a split synthesis. The resin is first divided into  $n$  equal portions. Following deprotection, each resin batch is reacted with a different monomer from set A. The resin is pooled, mixed, and divided into equal portions. After deprotection, individual monomers from set B are coupled. The process is repeated, and monomers from set C are added. Each final reaction batch is held separately for subsequent assay.



Compound identification is straightforward for spatially discrete libraries; however, identifying active species in the large and complex mixtures that result from pooling strategies presents a formidable problem. Typically, each member of a pool is present in a quantity that, although sufficient for determining biological activity, is insufficient for isolation and structure determination by standard analytical methods. As a result, strategies to identify the active members must be built into the synthetic schemes and closely coupled to the types of assays that will be performed with a particular library.

Compounds can be identified from soluble libraries using deconvolution strategies such as those introduced by Houghten (9, 10). This approach relies on cycles of synthesis and screening to focus progressively in on active species. An initial library is prepared in which one position of diversity is held constant. Each individual pool of compounds contains a different building block at the constant position but a mixture of building blocks at the remaining positions of diversity. Screening identifies the pool and therefore the optimal building block at the defined location. A second library is then synthesized in which all members have the optimal building block at this first position, but each pool contains a different building block at a second position of diversity. This process is repeated until the preferred building block for each position is determined. The chief disadvantage of this approach is that it is a laborious and time-consuming process. Moreover, this strategy does not guarantee that the most potent library member will be identified since selection of a “preferred” building block depends on the cumulative activity of a pool. Selection of a particular building block could, therefore, result

from the presence in a pool of a relatively few number of potent compounds or a large number of relatively poor inhibitors.

Several strategies have been developed that allow the active members of a library to be determined without requiring cycles of synthesis and screening. The method of positional scanning (11) involves synthesizing a series of combinatorial libraries, each fixing a single position of diversity. Synthesis of each library is divided such that each component pool is prepared using a single building block at the defined position, while the remaining positions are synthesized using mixtures of building blocks. Although this method does not require resynthesis, it does require additional synthetic effort to prepare a complete set of scanning libraries, and the total number of syntheses is greater than for deconvolution approaches. Moreover, because each position is defined independently, there is an increased chance of missing the most active members in the library.

Orthogonal pooling strategies also avoid the resynthesis burden of deconvolution methods (12). In this approach, two separate libraries are synthesized such that any given pool of the first library has only one compound in common with any pool of the second library. Determining the common species between positive pools from the two libraries reveals the identity of active compounds. The chance for misidentifying the optimal building block at any given position is avoided by this approach, but there is a considerable synthetic burden, as well as compound tracking issues.

The second general assay format is designed to measure the interaction of a target (enzyme, receptor, cell, etc) with compounds while they remain attached to the solid support used for library synthesis. Geysen pioneered this approach for assaying libraries of peptides arrayed on pins (1); however, the full potential was realized only when it was coupled to screening of OBOC libraries (8). Because each resin bead in an OBOC library carries a single library compound and the beads are physically distinct, each bead behaves as an isolated assay system. Briefly, a target is incubated with a pool of beads. The excess target is removed by washing, leaving behind only target that associates specifically with a compound attached to a resin bead. The active compounds are identified following visualization of those beads to which the target specifically interacts.

Specialized techniques are required to identify the compounds attached to individual beads, because each bead can carry at most a few nanomoles of compound, which is insufficient for structure determination except in a few rare cases (eg, peptides). These methods seek to “tag” or “encode” the beads with the synthetic history in a manner that is easily deciphered (13, 14). The code is introduced either prior to library synthesis or following each cycle of building block addition. In the latter case, the chemistries for library synthesis and code synthesis must be mutually compatible, which is one of the major obstacles in the development and utilization of this approach. Examples of encoding methods include the use of peptides read by [Edman Degradation](#) (15), nucleic acids read by PCR amplification and sequencing (16), several small molecule chemical classes identified by various analytical methods (17, 18), nonradioactive isotopic tags read by mass spectrometry (19), and radioisotopes encapsulated within the solid-phase support (20, 21). The development of encoding methods is an active area of research because the split synthesis method is currently the only tractable approach to preparing large libraries in which all possible compounds are represented.

See also [Combinatorial Libraries](#), [Libraries](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

## Bibliography

1. H. M. Geysen, R. H. Meleon, and S. J. Barteling (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3998–4002.
2. S. P. Fodor, J. L. Read, M. C. Pirrung, L. Stryer, A. T. Lu, and D. Solas (1991) *Science* **251**, 767–773.
3. A. C. Pease, D. Solas, E. J. Sullivan, M. T. Cronin, C. P. Holmes, and S. P. Fodor (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5022–5026.

4. C. Pinilla, J. Appel, S. Blondelle, C. Dooley, B. Dörner, J. Eichler, J. Ostresh, and R. A. Houghten (1995) *Biopolymers (Peptide Sci.)* **37**, 221–240.
5. H. M. Geysen, S. J. Rodda, and T. J. Mason (1986) *Mol. Immunol.* **23**, 709–715.
6. A. Furka, F. Sebestyen, M. Asgedom, and G. Dibo (1991) *Int. J. Pept. Protein Res.* **37**, 487–493.
7. K. S. Lam, M. Lebl, and V. Krchnak (1997) *Chem. Rev.* **97**, 411–448.
8. K. S. Lam, S. E. Salmon, E. M. Hersh, V. J. Hruby, W. M. Kazmierski, and R. J. Knapp (1991) *Nature* **354**, 82–84.
9. C. T. Dooley, N. N. Chung, B. C. Wilkes, P. W. Schiller, J. M. Bidlack, G. W. Pasternak, and R. A. Houghten (1994) *Science* **266**, 2019–2022.
10. S. E. Blondelle, E. Takahashi, P. A. Weber, and R. A. Houghten (1994) *Antimicrob. Agents Chemother.* **38**, 2280–2286.
11. R. A. Houghten, C. Pinilla, S. E. Blondelle, J. R. Appel, C. T. Dooley, and J. H. Cuervo (1991) *Nature* **354**, 84–86.
12. B. Deprez, X. Williard, L. Bourel, H. Coste, F. Hyafil, and A. Tartar (1995) *J. Am. Chem. Soc.* **117**, 5405–5406.
13. S. Brenner and R. A. Lerner (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5381–5383.
14. K. D. Janda (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10779–10785.
15. J. M. Kerr, S. C. Banville, and R. N. Zuckermann (1993) *J. Am. Chem. Soc.* **115**, 2529–2531.
16. M. C. Needels, D. G. Jones, E. H. Tate, G. L. Heinkel, L. M. Kochersperger, W. J. Dower, R. W. Barrett, and M. A. Gallop (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10700–10704.
17. M. H. J. Ohlmeyer, R. N. Swanson, L. W. Dillard, J. C. Reader, G. Asouline, R. Kobayashi, M. Wigler, and W. C. Still (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10922–10926.
18. Z.-J. Ni, D. Maclean, C. P. Holmes, M. M. Murphy, B. Ruhland, J. W. Jacobs, E. M. Gordon, and M. A. Gallop (1996) *J. Med. Chem.* **39**, 1601–1608.
19. H. M. Geysen, C. D. Wagner, W. M. Bodnar, C. J. Markworth, G. J. Parke, F. J. Schoenen, D. S. Wagner, and D. S. Kinder (1996) *Chem. Biol.* **3**, 679–688.
20. K. C. Nicolaou, X.-Y. Xiao, Z. Parandoosh, A. Senyei, and M. P. Novz (1995) *Angew. Chem. Int. Ed. English* **34**, 2289–2291.
21. E. J. Moran, S. Sarshar, J. F. Cargill, M. M. Shahbaz, A. Lio, A. M. M. Mjalli, and R. W. Armstrong (1995) *J. Am. Chem. Soc.* **117**, 10787–10788.

## Competence

Competence is the condition of bacterial cells that enables them to take up naked **DNA**. When this DNA can integrate into the **chromosome** by **recombination** or, in the case of **plasmids**, establish itself in the cytoplasm, the cells are said to be transformed (see **Transformation**). Competence is manifested by several species of bacteria—representatives of 29 different genera according to a recent count (**1**)—but far from all. Even in transformable strains, competence generally depends on the establishment of a particular physiological state. Four species have served as the major experimental models for studies of competence, and we shall therefore confine discussion largely to these: *Bacillus subtilis*, *Streptococcus pneumoniae* (“pneumococcus”), *Haemophilus influenzae*, and *Neisseria gonorrhoeae*. In view of the universal utility of *Escherichia coli*, which does not acquire competence naturally, we shall also review briefly the artificial induction of competence in this species.

The basic mechanism of DNA uptake appears to be the same in all naturally transformable organisms. Duplex DNA binds to the cell surface, and is then processed so that just one of the two strands passes through a specialized transmembrane channel into the cytoplasm, where it undergoes recombination with the resident chromosome. In **gram-negative** bacteria, a means of first passing through the outer membrane is also needed. The systems that regulate the uptake process, however, as well as the environmental stimuli to which these systems respond, appear to be more idiosyncratic, presumably reflecting the very different natural habitats in which the capacity to take up DNA is induced. Furthermore, the regulatory pathways of competence intersect with those of other metabolic activities. It is with the regulation of competence that this entry is chiefly concerned, the uptake process itself being deferred largely to the entry [Transformation](#). We begin with the bacterium for which most is known at both the mechanistic and regulatory levels.

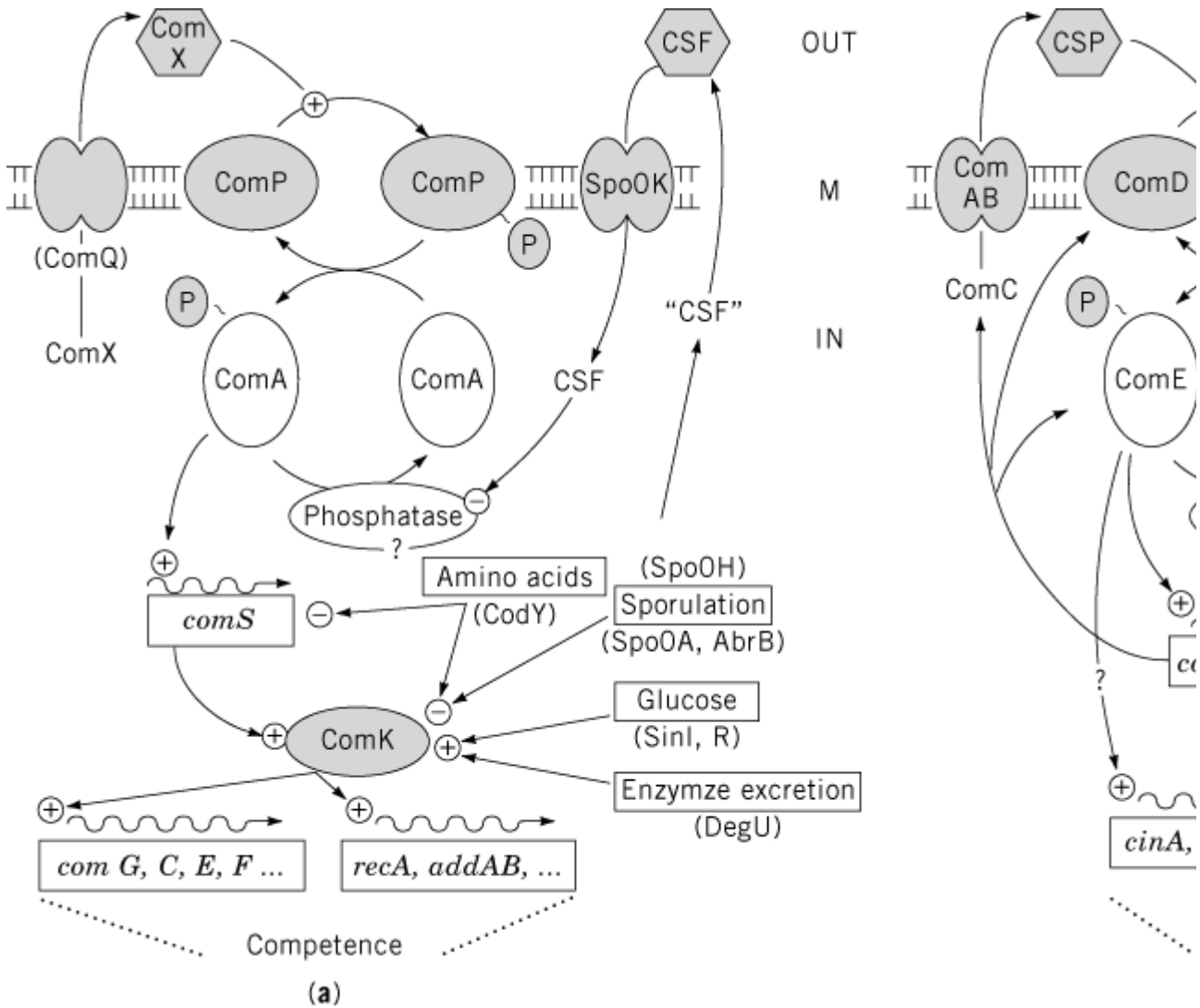
## 1. *B. subtilis*

*B. subtilis* survives and grows in many environments, but its natural reservoir is the soil. It was early appreciated that competence in this bacterium is a phenomenon induced by decelerating growth (2), such as occurs at the transition from exponential to stationary phase or upon depleting the medium of amino acids. Dependence of competence on transition to the stationary phase appears to reflect inhibition of competence by added amino acids (3). Prototrophs in minimal medium can attain competence in exponential phase (4). We now know that it is just one of a number of late-growth responses, which also include excretion of degradative [enzymes](#), antibiotic production, development of motility, and **sporulation** (5). Regardless of the many overlapping inputs from these other pathways, induction of competence is mediated by a single transcriptional activator, the ComK protein (6). This protein stimulates expression of the genes that encode the 12 components of the DNA-uptake machine identified to date (7-9). These genes are distributed among four polycistronic [transcription](#) units, *com G*, *C*, *E* and *F*, which specify the DNA-binding apparatus and its assembly, as well as the functions that translocate a single strand of DNA into the cytoplasm. The **nuclease** that degrades the other strand has not been identified. ComK also stimulates (a) the synthesis of **RecA**, which catalyzes the final step of transformation, (b) the recombination of the imported strand with the chromosome, and (c) the synthesis of AadAB, a nuclease/[DNA helicase](#) that assists the process.

Several intersecting metabolic currents determine the net level of active ComK at a given moment, and hence the degree of competence. Some of these are indicated in Figure 1. What is remarkable about this system is that positive input from all these factors is necessary for high levels of competence; the absence of any one of them diminishes activation of ComK or reduces transcription of its gene. Stochastic variability, leading to limits on one or more of these factors, may explain why in *B. subtilis* cultures, unlike those of other species, only about 10% of the cells become competent, even with optimum inducing treatments. The special condition of these cells is underscored by assorted observations of their relatively quiescent metabolic state and reduced buoyant density (4).

**Figure 1.** Competence induction in *B. subtilis* and *S. pneumoniae*. The pathways are drawn to emphasize similarities and are thought to be the main participants. M indicates the membrane, OUT indicates the medium, and IN indicates the cyto and in the text; OBL denotes oligopeptide-binding lipoprotein.; cin stands for competence induction; + indicates stimulat activity, – indicates repression or inhibition; question marks indicate factors not yet identified or presumed activities not transcripts. *B. subtilis*: The ComX pheromone is a peptide of 10 amino acid residues that is cleaved from the longer *com*. a *comQ*-dependent pore to act at a surface domain of the membrane protein, ComP, causing the latter to phosphorylate it ComP, the sensor of a two-component regulatory system, in turn phosphorylates ComA, converting it to a transcriptional Another short (4–7 residue) peptide, competence stimulating factor [CSF (12)], is cleaved from a precursor during expor through an oligopeptide **permease** pore and to protect ComA ~ P from inactivation by inhibiting a **phosphatase** (12, 13) ComK from a complex in which it is normally sequestered by the chaperone-like MecA and MecB proteins, liberating it competence genes (17, 18). *ComG* specifies the DNA-binding apparatus, the *comC* gene product is responsible for proce ComG complex in the membrane (9), *comE* encodes the DNA translocation functions, and *comF* appears to specify a me stranded DNA into the cytoplasm (8). The *recA* and *addAB* genes encode recombination functions. ComK also stimulate amino acids in the medium is mediated by the CodY protein, which represses transcription of *comS* and *comK* (3). The si

involve the freeing of an activator of *comK* transcription, SinR (19). *S. pneumoniae*: Competence stimulating peptide (CSP) peptide, probably generated by removal of the **signal sequence** from the ComC prepeptide during secretion by a specific activator stimulates the sensor-regulator protein pair, ComD-ComE, and the resulting ComE ~ P protein stimulates transcription of *comK*. ComD homologue interacts directly with CSP (25). Because CSP ultimately stimulates transcription of its own gene, it initiates wholesale induction of the culture, effectively synchronizing competence development. The *comC*, *D*, and *E* genes constitute coregulation. The importation of other peptides appears to modulate competence induction (27-29). Strains mutant for the *AliAB* develop competence at much lower cell densities than the wild type, suggesting that importation of peptides serves to repress *comCDE* transcription until exhausted at higher cell densities.



The finding that competence was attained earlier in cultures growing in medium already used for growth of cells to the competent state implied that extracellular signaling molecules promote development of competence (10), and this proved to be the case. Two such molecules are known, both generated by the cell and exported to the medium. The ComX pheromone interacts with the cell-surface protein ComP to initiate a cascade of interactions that culminate in competence (11, 12). Competence stimulating factor (CSF) is probably imported to stimulate this pathway. These interactions are shown in Figure 1 and described in the legend. It is possible that the apparent redundancy of competence-inducing function reflects only the co-opting of the CSF sporulation factor to protect the primary product of the competence-specific pathway, the phosphorylated form of ComA, ComA ~ P (13). Unlike ComX, which seems to be specific for competence, CSF production is intimately linked with the system that regulates sporulation; transcription of its gene requires the Spo0H sigma factor and is stimulated, indirectly, by the Spo0A ~ P transcription

regulator. Grossman (14) has suggested that CSF is also a signaling molecule for sporulation.

The external actions of ComX pheromone and of CSF have led to the suggestion that these factors communicate cell-density information to the system that regulates competence (14). ComA thus serves as a depot for these signals, and in its phosphorylated form it sends them on by activating transcription of *comS* (12, 15-18).

This general induction pathway is also subject to nutritional signals: amino acids via CodY (3), and glucose via SinI and SinR (19). How these, as well as signals from the sporulation and enzyme secretion pathways (AbrB, DegU ¼), are integrated to allow the cell to decide whether to commit itself to competence remains an intriguing issue.

## 2. *S. pneumoniae*

This organism is an obligate parasite of the respiratory tract. In contrast to *B. subtilis*, *S. pneumoniae* readily attains competence during exponential growth in rich medium. Moreover, all the cells become competent, albeit briefly. Competence in *S. pneumoniae* does share with that of *B. subtilis*, however, an association with a distinct physiological state. This takes the form of altered surface properties, cessation of cell-wall synthesis, increased cell-wall fragility, and release of DNA into the medium. The onset of competence is very sensitive to culture conditions: temperature, the nature of the medium buffer, the concentrations of magnesium and calcium, and pH. In low-pH medium, there is only one peak of competence, at high cell density, whereas in high-pH medium there is a succession of peaks (20).

The competence induction pathway of *S. pneumoniae* shows notable similarities to that of *B. subtilis*, as depicted in Figure 1. Induction of competence is mediated by an extracellular signaling molecule, the competence stimulating peptide (CSP) (21-23), and this initiates the induction cascade via a sensor-regulator protein pair (24-26). Induction is also modulated by imported peptides (27, 28), though in this case the effect is negative (29).

Several proteins other than those shown in Figure 1 make their appearance at competence, some of them associated with the membrane (30); but the only one demonstrated to have a clear role in transformation, indeed in competence itself, is the RecA protein, which mediates homologous recombination (31).

*S. pneumoniae* cultures exhibit a particular aspect of competence that, though detected long ago (32), has received little attention. Filtrates of cultures that have passed their peak of competence prevent development of competence when added to another culture, suggesting that a specific inhibitor of competence is secreted into the medium. Such a factor might be induced at competence just in order to turn it off again, accounting for the brevity of the competence peak.

## 3. *H. influenzae*

*H. influenzae* lives in the mucus of the respiratory tract, where it may be exposed to large amounts of naked DNA (33). Only its own DNA, however, is taken up with high efficiency. This selectivity results from the presence throughout the genome of numerous copies of a short, specific base sequence, the uptake-signal sequence (USS), which is recognized by the uptake apparatus (34, 35). The outer membranes of competent cells extrude vesicles, termed *transformasomes* (or “blebs”), which encapsulate the DNA and then release it into the **periplasm** to allow transport into the cell (36-38). As in *B. subtilis*, growth-inhibiting conditions, such as transient anaerobiosis, nutritional downshift, and approach to stationary phase, induce competence. The level of competence attained varies and can be close to 100% of the population. Nevertheless, there appears to be no cell density sensing involved, each cell developing competence independently. *H. influenzae* cells implanted in the rat peritoneum rapidly develop the capacity to import DNA added with the inoculum and manifest other changes characteristic of the competent state (39).



Several factors important to the development of competence in *H. influenzae* have been identified, but the picture remains fragmented. Of crucial importance is the intracellular level of [cyclic AMP \(cAMP\)](#). Competence was abolished by mutations in the [cyclic AMP receptor protein \(CRP\)](#) and [adenylate cyclase \(Cya\)](#) genes, but could be restored to a *cya* mutant by addition of cAMP (40, 41). The cell senses sugar availability through a carbohydrate phosphotransferase system that increases cAMP synthesis in response to declining sugar availability: Mutations in two genes specifying components of this system, *ptsI* and *crr*, caused partial loss of competence, a deficiency that can be compensated for by added cAMP (42, 43). Competence in *H. influenzae* thus responds directly to nutritional state.

Some clues as to how cAMP might act in competence induction have recently emerged. Notably, cAMP stimulates transcription from a divergent **promoter** region of two genes essential for transformation, *tfoX* and *rec-1*. *TfoX* was originally identified by a mutation, *sxy-1*, that increased transformation 100-fold (44). Null mutations in the gene eliminated DNA binding and transformation, whereas overexpression of *tfoX* caused constitutive competence (45, 46), implying that TfoX is a positive regulator of competence. The finding that TfoX is needed for expression of other competence genes (*com101A*, *dprA*, and *rec-2*) (46, 47) supports this suggestion and hints at a central role in competence similar to that played by the ComK regulator of *B. subtilis*. Whereas *tfoX* transcription normally increases only as cells approach stationary phase, addition of cAMP to a log-phase culture causes an immediate stimulation (46).

On the other hand, transcription of *rec-1*, which encodes the general recombination enzyme that is analogous to RecA of *E. coli*, is only weakly stimulated by cAMP-CRP (48). It is, however, induced by single-stranded DNA, as befits its function of integrating transforming DNA. This common property of Rec proteins has been proposed to endow Rec-1 with another quite distinct function, that of triggering competence. Induction of competence in *H. influenzae* may be associated with the conversion of up to 5% of the chromosomal DNA to single-stranded form. Observing that overproduction of Rec-1 led to early development of competence, Stuy (49) suggested that formation of Rec-1/single-stranded DNA complexes might induce competence gene expression, as the equivalent complexes induce the SOS regulon in *E. coli*. The idea remains untested.

Several new proteins appear in the outer membrane at competence. Some of them, rather than being synthesized *de novo*, are redistributed from the inner membrane. One mutant that fails to do this also fails to bind DNA. It carries a lesion in the gene *por*, which encodes a periplasmic disulfide oxidoreductase, suggesting the need to assemble or fold **disulfide bond**-containing proteins (50) (see [Protein Disulfide Isomerase](#)). Two proteins probably involved in conducting DNA into the cytoplasm are induced at competence: Rec-2, which shares [homology](#) with the presumed transmembrane channel protein ComEC of *B. subtilis*, and DprA (47).

There is as yet no indication that *H. influenzae* competence responds to a specific extracellular signal. If it does, it now seems unlikely to do so via a sensor-regulator protein pair, such as ComP-ComA in *B. subtilis*: Mutation of the genes encoding four such systems, as deduced from the genome sequence, failed to affect competence (43). What is clear is that the identity and function of most of the components of the competence regulatory system are still unknown, and that the true place of the known ones in the system likewise remains to be shown.

#### 4. *N. gonorrhoeae*

Unlike the bacteria already discussed, for which competence is a transitory response to growth conditions, transformable *Neisseria* strains are competent at all phases of growth (51). An early indication (52) that competence could be induced by a **proteinase**-sensitive factor present in conditioned medium has not been confirmed. The transformation mutants so far reported are defective in one or more of the individual functions known to be involved in transformation-[pili](#) and proteins for DNA binding, transport, and recombination (53-55), but none have suggested the

existence of competence factors or regulators. The competent state may be so integrated into the general physiology of *Neisseria* that a specific regulatory system is not needed, or the regulatory elements are essential and therefore missed in mutant hunts. Transformation of cells in culture is influenced by pH and by the concentrations of mono- and divalent cations (53), but these factors probably affect DNA-cell envelope interactions directly, rather than by regulating competence gene output.

## 5. *E. coli*

In 1970, Mandel and Higa (56) discovered that *E. coli* could be induced to take up DNA by treatment of the cells with  $\text{Ca}^{2+}$  ions. Although this artificial competence has been a key technique in molecular genetics, its basis is not understood. Suspension of *E. coli* cells in cold  $\text{CaCl}_2$  solution is known to induce the formation of transmembrane pores consisting of an outer shell of polyhydroxybutyrate, connected by  $\text{Ca}^{2+}$  to an inner sheath of polyphosphate (57). The role of these pores is unclear, because their internal diameter appears to be too small to allow the passage of double-stranded DNA. They do, however, have marked effects on temperature-induced membrane fluidity, and it is possible that, by anchoring an otherwise mobile membrane, they create tension that opens holes elsewhere in the membrane to admit the DNA (57, 58).

## 6. Summary

One of the most striking aspects of research on competence to emerge in recent years is the interrelationships between the proteins and regulatory factors involved and those of other phenomena previously considered distinct, perhaps the clearest example being *B. subtilis* sporulation (14). Indeed, acquisition of competence is usually accompanied by other major physiological changes. Because of this, and because our knowledge of competence is based largely on laboratory observation, it is fair to question whether DNA uptake itself is relevant in natural bacterial habitats. In fact there are good reasons for thinking that it is, and this issue is taken up in the entry [Transformation](#).

## Bibliography

1. M. G. Lorenz and W. Wackernagel (1994) *Microbiol. Rev.* **58**, 563–602.
2. C. Anagnostopoulos and J. Spizizen (1961) *J. Bacteriol.* **81**, 741–746.
3. P. Serror and A. L. Sonenshein (1996) *J. Bacteriol.* **178**, 5910–5915.
4. D. Dubnau (1993) In *Bacillus subtilis and Other Gram-Positive Bacteria* (A. L. Sonenshein, J. A. Hoch, and R. Losick, eds.), ASM, Washington, D.C., pp. 555–584.
5. D. C. Dooley, C. T. Hadden, and E. W. Nester (1971) *J. Bacteriol.* **108**, 668–679.
6. D. van Sinderen and G. Venema (1994) *J. Bacteriol.* **176**, 5762–5770.
7. D. van Sinderen, A. Luttinger, L. Kong, D. Dubnau, G. Venema, and L. Hamoen (1995) *Mol. Microbiol.* **15**, 455–462.
8. D. Dubnau (1997) *Gene* **192**, 191–198.
9. Y. S. Chung and D. Dubnau (1994) *Mol. Microbiol.* **15**, 543–551.
10. H. Joenje, M. Gruber, and G. Venema (1972) *Biochim. Biophys. Acta* **262**, 189–199.
11. R. Magnuson, J. Solomon, and A. D. Grossman (1994) *Cell* **77**, 207–216.
12. J. Solomon, R. Magnuson, A. Srivastava, and A. D. Grossman (1995) *Genes Dev.* **9**, 547–558.
13. J. Solomon, B. Lazazzera, and A. D. Grossman (1996) *Genes Dev.* **10**, 2014–2024.
14. A. D. Grossman (1995) *Annu. Rev. Genet.* **29**, 477–508.
15. N. M. Nakano and P. Zuber (1991) *J. Bacteriol.* **173**, 7269–7274.
16. J. Hahn and D. Dubnau (1991) *J. Bacteriol.* **173**, 7275–7282.
17. L. Kong and D. Dubnau (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5793–5797.

18. T. Msadek, F. Kunst, and G. Rapoport (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5788–5792.
19. U. Bai, I. Mandic-Mulic, and I. Smith (1993) *Genes Dev.* **7**, 139–148.
20. J. D. Chen and D. A. Morrison (1987) *J. Gen. Microbiol.* **133**, 1959–1967.
21. A. Tomasz and J. L. Mosser (1966) *Proc. Natl. Acad. Sci. USA* **55**, 58–66.
22. L. S. Havarstein, G. Coomaraswami, and D. A. Morrison (1995) *Proc. Natl. Acad. Sci. USA* **92**, 11140–11144.
23. L. S. Havarstein, D. B. Diep, and I. F. Nes (1995) *Mol. Microbiol.* **16**, 229–240.
24. E. V. Pestova, L. S. Havarstein, and D. A. Morrison (1996) *Mol. Microbiol.* **21**, 853–862.
25. L. S. Havarstein, P. Gaustad, I. F. Nes, and D. A. Morrison (1996) *Mol. Microbiol.* **21**, 965–971.
26. Q. Cheng, E. A. Campbell, A. M. Naughton, S. Johnson, and H. R. Masure (1997) *Mol. Microbiol.* **23**, 683–692.
27. G. Alloing, P. de Philip, and J.-P. Claverys (1996) *J. Mol. Biol.* **241**, 44–58.
28. B. J. Pearce, A. M. Naughton, and H. R. Masure (1994) *Mol. Microbiol.* **12**, 881–892.
29. G. Alloing, B. Martin, C. Granadel, and J.-P. Claverys (1998) *Mol. Microbiol.* **29**, 75–83.
30. M. N. Vijayakumar and D. A. Morrison (1986) *J. Bacteriol.* **165**, 689–695.
31. B. Martin, P. Garcia, M.-P. Castaniè, and J.-P. Claverys (1995) *Mol. Microbiol.* **15**, 367–379.
32. A. Tomasz and R. D. Hotchkiss (1964) *Proc. Natl. Acad. Sci. USA* **51**, 480–487.
33. L. Matthews, S. Spector, L. Lemm, and J. Potter (1963) *Am. Rev. Respir. Dis.* **88**, 199–204.
34. K. L. Sisco and H. O. Smith (1979) *Proc. Natl. Acad. Sci. USA* **76**, 972–976.
35. D. B. Danner, R. A. Deich, K. L. Sisco, and H. O. Smith (1980) *Gene* **11**, 311–318.
36. M. Kahn, M. Concino, R. Gromkova, and S. H. Goodgal (1979) *Biochem. Biophys. Res. Commun.* **87**, 764–772.
37. R. A. Deich and L. C. Hoyer (1982) *J. Bacteriol.* **152**, 855–864.
38. M. E. Kahn, G. Maul, and S. H. Goodgal (1982) *Proc. Natl. Acad. Sci. USA* **79**, 6370–6374.
39. M. Dargis, P. Gourde, D. Beauchamp, B. Foiry, M. Jaques, and F. Malouin (1992) *Infect. Immun.* **60**, 4024–4031.
40. M. S. Chandler (1992) *Proc. Natl. Acad. Sci. USA* **89**, 1626–1630.
41. I. R. Dorocicz, P. M. Williams, and R. J. Redfield (1993) *J. Bacteriol.* **175**, 7142–7149.
42. L. P. Macfadyen, I. R. Dorocicz, J. Reizer, M. H. Saier Jr., and R. J. Redfield (1996) *Mol. Microbiol.* **21**, 941–952.
43. M. L. Gwinn, D. Yi, H. O. Smith, and J. F. Tomb (1996) *J. Bacteriol.* **178**, 6366–6368.
44. R. J. Redfield (1991) *J. Bacteriol.* **173**, 5612–5618.
45. P. M. Williams, L. A. Bannister, and R. J. Redfield (1994) *J. Bacteriol.* **176**, 6789–6794.
46. J. J. Zulty and G. J. Barcak (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3616–3620.
47. S. Karudapuram and G. J. Barcak (1997) *J. Bacteriol.* **179**, 4815–4820.
48. J. J. Zulty and G. J. Barcak (1993) *J. Bacteriol.* **175**, 7269–7281.
49. J. H. Stuy (1989) In *Genetic transformation and expression* (L. O. Butler et al., eds.), Intercept Inc., Andover, pp. 85–112.
50. J.-F. Tomb (1992) *Proc. Natl. Acad. Sci. USA* **89**, 10252–10256.
51. P. F. Sparling (1966) *J. Bacteriol.* **92**, 1364–1369.
52. A. Siddiqui and I. D. Goldberg (1975) *Biochem. Biophys. Res. Commun.* **64**, 34–42.
53. G. D. Biswas, T. Sox, E. Blackman, and P. F. Sparling (1977) *J. Bacteriol.* **129**, 983–992.
54. G. D. Biswas, S. A. Lacks, and P. F. Sparling (1989) *J. Bacteriol.* **171**, 657–664.
55. I. J. Mehr and H. S. Seifert (1997) *Mol. Microbiol.* **23**, 1121–1131.
56. M. Mandel and A. Higa (1970) *J. Mol. Biol.* **53**, 159–162.

57. R. N. Reusch and H. L. Sadoff (1988) Proc. Natl. Acad. Sci. USA **85**, 4176–4180.  
 58. C. E. Castuma, R. Huang, A. Kornberg, and R. N. Reusch (1995) J. Biol. Chem. **270**, 12980–12983.

### Suggestions for Further Reading

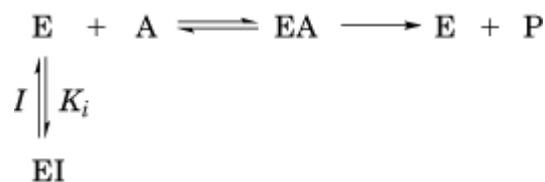
59. J. M. Solomon and A. D. Grossman (1996) Who's competent and when: regulation of natural genetic competence in bacteria. Trends Genet. **12**, 150–155. A brief but comprehensive review, providing access to the important literature.  
 60. A. D. Grossman (1995) Genetic networks controlling the initiation of sporulation and the development of genetic competence in *Bacillus subtilis*. Annu. Rev. Genet. **29**, 477–508. A review of competence development in *B. subtilis* and its interconnections with sporulation and nutritional state.  
 61. J.-P. Claverys, A. Dintilhac, I. Mortier-Barrière, B. Martin, and G. Allosing (1997) Regulation of competence for genetic transformation in *Streptococcus pneumoniae*. J. Appl. Microbiol. Symp. Suppl. **83**, 32S–42S. A review of competence regulation in this pathogen, emphasizing the influence of peptide transport and other factors.

## Competitive Inhibition

Analogs of the substrate of an [enzyme](#) can often bind at the same place within the enzyme's [active site](#) as the substrate but are not acted on by the enzyme. They are inhibitors of the enzyme, giving rise to [dead-end inhibition](#), and they usually compete with the substrate, giving rise to competitive inhibition .

### 1. Linear

The simplest type of competitive inhibition occurs when an inhibitor combines reversibly at the active site of the same form of the enzyme as the substrate. The inhibitor may be a structural analogue of the substrate or a product of that substrate. In the following scheme, the inhibitor I combines reversibly with the free enzyme to form a dead-end EI complex that can only



dissociate back to the components from which it was formed. Therefore, the interaction of E and I is at thermodynamic equilibrium, and the inhibition constant  $K_i$  is a **dissociation constant**. The general equation that describes this type of inhibition is shown in Equation [1](#):

$$v = \frac{VA}{K_a(1 + I/K_{is}) + A} \quad (1)$$

where  $v$  is the observed velocity of the reaction,  $V$  the maximum velocity,  $K_a$  the Michaelis constant for the substrate present at concentration  $A$ ,  $I$  the concentration of inhibitor, and  $K_{is}$  the apparent

inhibition constant. For a single-substrate reaction,  $K_{is}$  would be equal to  $K_i$ . The double-reciprocal form, [Lineweaver–Burk plot](#), of the equation for linear competitive inhibition (Eq. 2):

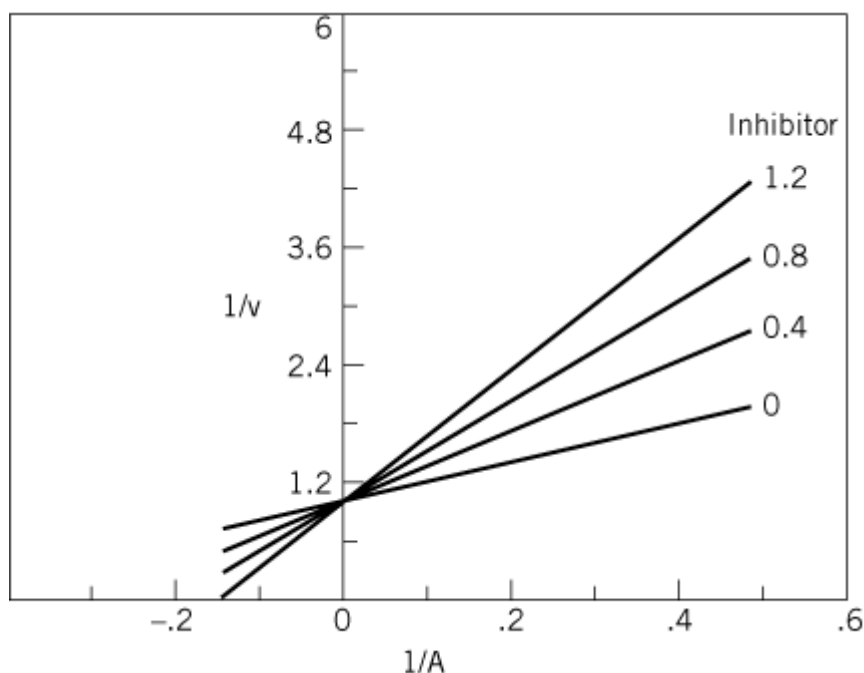
$$\frac{1}{v} = \frac{K_a}{V} \left( 1 + \frac{I}{K_{is}} \right) \frac{1}{A} + \frac{1}{V} \quad (2)$$

indicates that a primary plot of  $1/v$  against  $1/A$  at different concentrations of  $I$  would yield a family of straight lines that intersect at the same point on the vertical ordinate (Fig. 1). This is the signature of competitive inhibition; sufficiently high concentrations of substrate can compete successfully with the inhibitor. The slope of each curve is given by Equation 3:

$$\text{Slope} = \frac{K_a}{VK_{is}} I + \frac{K_a}{V} \quad (3)$$

A plot of slope against inhibitor concentration would yield a straight line that intersects the abscissa at the point where  $I = -K_{is}$ . Since this secondary plot should be linear, the inhibition would be classified as linear competitive. Although values for  $K_{is}$  can be obtained by such graphical procedures, it is preferable to obtain its value by an overall fit of the primary data to the general equation by use of a least-squares fitting method (1).

**Figure 1.** Primary double-reciprocal plot for linear competitive inhibition.

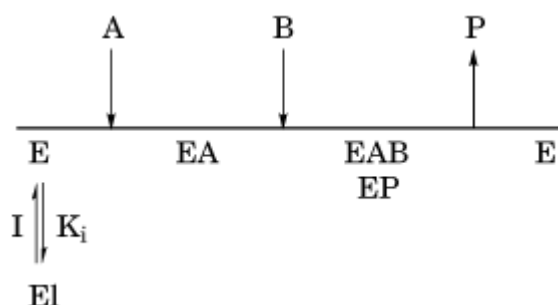


When two substrates are involved in a reaction, the value determined directly for  $K_{is}$  may be only an apparent constant, with its value dependent on the concentration of the nonvaried substrate (see [Lineweaver–Burk Plot](#)). In this case, it would also have to be stated that the inhibitor was linear competitive with respect to a particular substrate.

Linear competitive inhibition will also be observed with a two-substrate reaction when the inhibitor is an analogue of the first substrate to add in an equilibrium ordered mechanism, A, and the second

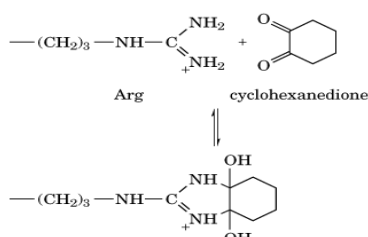
substrate to add, B, is varied (Fig. 2). The inhibition by I with respect to A will be linear competitive because both A and I compete directly for the same form of enzyme (E). The inhibition by I with respect to B will also be linear competitive, even though I and B combine with different enzyme forms, because the reaction between E and A in an equilibrium-ordered mechanism is at thermodynamic equilibrium. Thus, increasing concentrations of B will ultimately reduce the concentrations of both EA and E to zero, so that no free enzyme remains to interact with I.

**Figure 2.** Equilibrium-ordered mechanism for which inhibition by I would be linear competitive with respect to both substrates.



## 2. Hyperbolic

This type of inhibition occurs when an inhibitor binds reversibly at a site on the enzyme, other than the active site, so as to make it more difficult for the substrate to combine at the active site. Thus, the inhibitor and substrate can be present on the enzyme at the same time. The reactions involved can be illustrated as follows:



where  $K_{is}$  and  $K_{id}$  denote dissociation constants for EI and EIA, respectively. It is a thermodynamic requirement that the affinity of the substrate A for the E and EI forms differs, so that

$$K_A = K_a \frac{K_{id}}{K_{is}}$$

For the inhibition to be competitive, the presence of the inhibitor on the enzyme must not affect the rate of product formation. The general equation that describes hyperbolic competitive inhibition is given in Equation 4:

$$v = \frac{V A [1 + \frac{I}{K_{id}}]}{K_a (1 + \frac{I}{K_{id}}) + A} \quad (4)$$

which can be arranged in double-reciprocal form as

$$\frac{1}{v} = \frac{K_a}{V} \left[ \frac{1 + \frac{I}{K_{is}}}{1 + \frac{I}{K_{id}}} \right] \frac{1}{A} + \frac{1}{V} \quad (5)$$

This is an equation of a straight line for a plot of  $1/v$  against  $1/A$ , but it is apparent that the slopes of the lines will vary as a hyperbolic function of the concentration of  $I$  and the intercept of the curves with the vertical ordinate is independent of  $I$ . For a plot of slope as a function of  $I$ , the curve would be a concave-down, nonrectangular hyperbola, with limiting values of  $K_a/V$  and  $K_a K_{id}/(V K_{is})$  at zero and infinite concentrations of  $I$ , respectively. Values for  $K_{id}$  and  $K_{is}$  may be obtained from a plot of  $1/(\text{slope}_o - \text{slope}_i)$  against  $1/I$ , where  $\text{slope}_o$  and  $\text{slope}_i$  are slopes of the primary plot in the absence and presence of inhibitor, respectively. Values for  $K_{is}$  and  $K_{id}$  are best obtained by an overall least-squares fit of the data to the equation that describes hyperbolic competitive inhibition (1).

It should be noted that the same kinetic mechanism would apply if  $I$  were an activator that enhances the combination of the substrate with the enzyme. In this case, the plot of slope against  $I$  would yield a concave-up, nonrectangular hyperbola (2).

### Bibliography

1. W. W. Cleland (1979) *Meth. Enzymol.* **63**, 103–138.
2. V. L. Schramm and J. F. Morrison (1970) *Biochemistry* **9**, 671–677.

### Competitive Labeling

Competitive labeling is a method of [chemical modification](#) for determining the  $pK_a$  and the reactivity of functional groups in proteins. The extent of modification of a particular protein group at various pH values is compared to that of the same group in a standard compound. This gives the reactivity of the protein group relative to its intrinsic reactivity, and the pH-dependence gives the  $pK_a$  value.

These data reflect the **accessibility** of the group in the folded protein, which provides information about the chemical reactivity of the group and its environment in the folded protein. This indirectly gives conformational information about the protein. Conformational transitions associated with function, association, and **ligand binding** can be monitored with this method.

Competitive labeling was originally employed in a study of the [amino groups](#) in [elastase](#) (1). A mixture of elastase and phenylalanine (the standard compound) was acetylated with a small amount of  $^3\text{H}$ -labeled acetic anhydride (see [Anhydrides](#)). It is important that the amount of the labeling agent is less than that of the groups being modified. Under these conditions, the relative amounts of label in the standard and in the specific protein group will depend upon their relative ionization and reactivities. The fraction of the reactive, nonionized form of the amino group ( $\alpha_i$ ) and the ratio of reactivity ( $r = k_i/k_{\text{standard}}$ ) of the  $i$ th group are derived from the known  $pK_a$  of the standard ( $\alpha_{\text{standard}}$ ) and the measured ratio of labeling:

$$\alpha_i r = \alpha_{\text{standard}} (\text{label on } i\text{th group} / \text{label on standard}) \quad (1)$$

The label on the  $i$ th group could be determined by the radioactivity of the appropriate peptide separated by [peptide mapping](#). The values on the right side of the equation are obtained

experimentally at each pH and are plotted against pH. The values of the  $pK_a$  and of  $r$  of the protein group are obtained by theoretical curve fitting.

Besides using acetic anhydride to react with amino groups, 1-fluoro-2,4-dinitrobenzene has been used to react with amino groups and with [histidine](#) and [tyrosine](#) residues, and **iodoacetate** with [thiol groups](#). A small agent is most suitable for this purpose of modification. Double labeling with  $^3\text{H}$ -label, followed by  $^{13}\text{C}$ -labeling after **denaturation** gives a much more sensitive analysis. In this case, the same group in the denatured protein is the standard.

#### Bibliography

1. H. Kaplan, K. J. Stevenson, and B. S. Hartley (1971) *Biochem. J.* **124**, 289–299.

#### Suggestion for Further Reading

2. N.M. Young and H. Kaplan (1989) "Chemical characterization of functional groups in proteins by competitive labelling". In *Protein Structure: A Practical Approach* (T.E. Creighton, ed.), IRL Press, Oxford, U.K., pp. 225–245.

## Complement Fixation

**Complement** is a set of serum proteins that undergo a sequential series of cleavage and association reactions in response to antibody–antigen complexes. If these complexes occur on a cell surface, complement reactions can cause cell lysis (1). The lysis of erythrocytes is particularly easy to detect visually, due to the release of hemoglobin, and this lysis forms the principle for a type of immunoassay. Antibody–antigen complexes in solution cause complement reactions to occur without cell lysis. Since the complement reactions are stoichiometric and the products short-lived, these non-cell-associated complexes deplete a portion of the complement proteins, a process termed “complement fixation.” This depletion reduces the lysis observed when the complement is subsequently added to antibody-coated erythrocytes.

Key reagents for complement fixation assays, all available commercially, are complement, erythrocytes, and hemolytic antibody (2, 3). Guinea pig serum is usually used as a source of complement, since guinea pig complement reacts with antibodies of most experimentally significant mammalian species. Sheep erythrocytes, washed free of serum components, are used for lysis. Hemolytic antibody is obtained by immunizing rabbits with sheep erythrocyte membranes. The erythrocytes are sensitized in preparation for complement lysis by incubation with the hemolytic antibody. Once amounts of complement and sensitized erythrocytes appropriate for easy detection of lysis are determined, the reagents can be used to measure the interaction of an antibody–antigen pair (4). The antibody and antigen to be tested are incubated for a set time in the presence of complement. Standards containing known amounts of antibody and antigen are incubated simultaneously. Sensitized erythrocytes are next added, and lysis is measured after further incubation, by spectrophotometric determination of hemoglobin. Diminished lysis relative to a control indicates that the antibody–antigen pair in question has formed an aggregate and induced complement reactions to occur in solution, reducing the amount of complement available to lyse cells. Results observed with the standards are used to calibrate the correspondence between antibody–antigen complex formation and lysis. Because antibody–antigen aggregates are the most active species in fixing complement, bell-shaped response curves result when percent complement fixed is plotted against added antibody or antigen. Complement fixation is most efficient at roughly equivalent amounts of antibody and



antigen, whereas either component in excess will favor formation of binary and ternary complexes that do not fix complement well. Consequently, several concentrations of antigen or antibody must be tested to establish whether the observed response lies on the increasing or decreasing slope of the standard curve.

As an analytical technique, complement fixation is sensitive to nanogram quantities of antigen, but is subject to reagent variation and interference by chemical components of the sample. A notable value, however, is in study of cross-reactivity between structurally similar antigens. Proteins that differ slightly in sequence give easily measured differences in complement fixation ability when they are probed with a single antiserum. This effect is so regular that complement fixation by homologous antigens can be used to establish the evolutionary relationship between closely related species (5, 6).

## Bibliography

1. J. Bordet and O. Gengou (1901) *Ann. Inst. Pasteur* **15**, 289–302.
2. M. M. Mayer (1961) "Complement and complement fixation", in *Experimental Immunochimistry*, Thomas, Springfield, IL, pp. 133–240.
3. L. Levine and H. Van Vunakis (1967) *Meth. Enzymol.* **11**, 928–936.
4. E. Wasserman and L. Levine (1961) *J. Immunol.* **87**, 290–295.
5. E. M. Prager and A. C. Wilson (1971) *J. Biol. Chem.* **246**, 5978–5989.
6. A. B. Champion, E. M. Prager, D. Wachter, and A. C. Wilson (1974) "Microcomplement fixation", in *Biochemical and Immunological Taxonomy of Animals*, C. A. Wright, ed., Academic Press, London, pp. 397–416.

## Complement System

### 1. Introduction

After the discovery of humoral immunity, Hans Buchner found that if fresh serum possessing an antibacterial antibody was added to bacteria, the bacteria were quickly lysed. However, if the serum was heated to 56°C, the lytic capacity of the serum was lost. He also demonstrated that the loss of lytic capacity could be restored by using unheated nonimmune serum in conjunction with heated immune serum. These studies were extended by others, including Jules Bordet, who concluded that serum contained two components necessary for cellular lysis. These two components are heat-stable [antibodies](#) and a heat-labile component that “complements” the lytic function of antibodies. He reasoned that antibodies had two binding sites, one for [antigen](#) and the other for heat-labile substance that was given the name *complement*.

Complement is now known not to be a single component, but it consists of at least 20 chemically and immunologically distinct plasma proteins capable of interacting with one another in a highly regulated manner to provide at least four main biological functions. First, *cytolysis* is mediated by the association, or polymerization, of specifically activated complement components on the surface of target cells. These components form pores that disrupt the integrity of the lipid membrane, and the cell is killed by osmotic lysis. Second, antibody and antigen combine to form *immune complexes* that, unless removed, can result in damage of body tissues. The binding of complement proteins prevent the damage from immune complexes by mediating their solubilization and clearance. Third, the *opsonization* of foreign particles is mediated by the binding of complement proteins (opsonins). Phagocytic leukocytes bear receptors for these complement proteins, so that opsonized particles are cleared from the body by [phagocytosis](#). Finally, through the activation of complement, *proteolytic*

fragments are released whose function are to *mediate inflammation*. These fragments, or “anaphylatoxins,” can act on several target cells such as neutrophils, smooth muscle, and vascular endothelium, as well as organ systems of the body.

The complement proteins are normally present in the circulation and are produced primarily by the liver and other extrahepatic sources, such as [macrophages](#) and fibroblasts. Some of the components are produced as functionally inactive [proteinases](#) that are activated only when proteolytically cleaved themselves by previously activated complement proteins. Complement activation occurs only at localized sites under specific conditions. The binding of specific antibody to antigen can initiate complement activation through what is called the *classical pathway*. Additionally, in the absence of antibody, some complement components are directly activated by binding to the surfaces of infectious agents, such as bacteria, **fungi**, and **viruses**, through what is called the **alternative pathway**.

## 2. Overview of the Complement System

Individual proteins of the complement system are present in the serum as functionally inactive molecules. The individual precursor proteins are designated numerically, such as C1, C2, up to C9. Other proteins involved with the complement system have retained their common, or trivial, names such as properdin, factor B, and factor H. Individual components must be activated sequentially and specifically for the complement cascade to progress and mediate inflammation and lysis of foreign organisms. The activation of complement is a dynamic process that allows the proteins of the system to interact functionally and is not a static singular event.

An overview of the operation of the complement system is shown in [Figure 1](#), and biochemical and physical information of the proteins involved in the complement system are listed in [Table 1](#). A protein called C3 is the central component of the complement system. The alternative and classical pathways are two parallel, but entirely independent, mechanisms that lead to the formation of C3 *convertases* whose function is to cleave C3 into C3a and C3b. In the classical pathway, initial events involve binding of specific antibody to antigen. This antigen–antibody complex binds complement C1 and sequentially activates complement proteins C4 and C2, leading to the formation of the C4b2b complex that functions as the *classical pathway C3 convertase*. In contrast, activation of complement via the alternative pathway occurs independent of antigen–antibody complexes. Instead, alternative pathway activation is dependent on activators that allow the formation and deposition of the C3bBb complex (**alternative pathway C3 convertase**) on their surfaces. Activators of the alternative pathway include many bacteria, fungi, and viruses. Both classical and alternative convertases can cleave additional C3, generating more C3b (see [Fig. 1](#)). The binding of additional C3b to these C3 convertases changes them conformationally to **C5 convertases**, which specifically cleave C5. Once cleaved, both the alternative and classical pathways share the same terminal steps. These steps do not involve the proteolytic cleavage of additional components but instead involve the sequential binding of components C6, C7, C8, and C9. This terminal assembly leads to the formation of the *membrane attack complex (MAC)*, resulting in the osmotic lysis or cytolysis of bacteria or other affected cells.

**Figure 1.** Schematic representation of the alternative and classical activation of the complement pathway leading to the assembly of the membrane attack complex.

Antigen-antibody complex  
(IgG or IgM)

Classical path

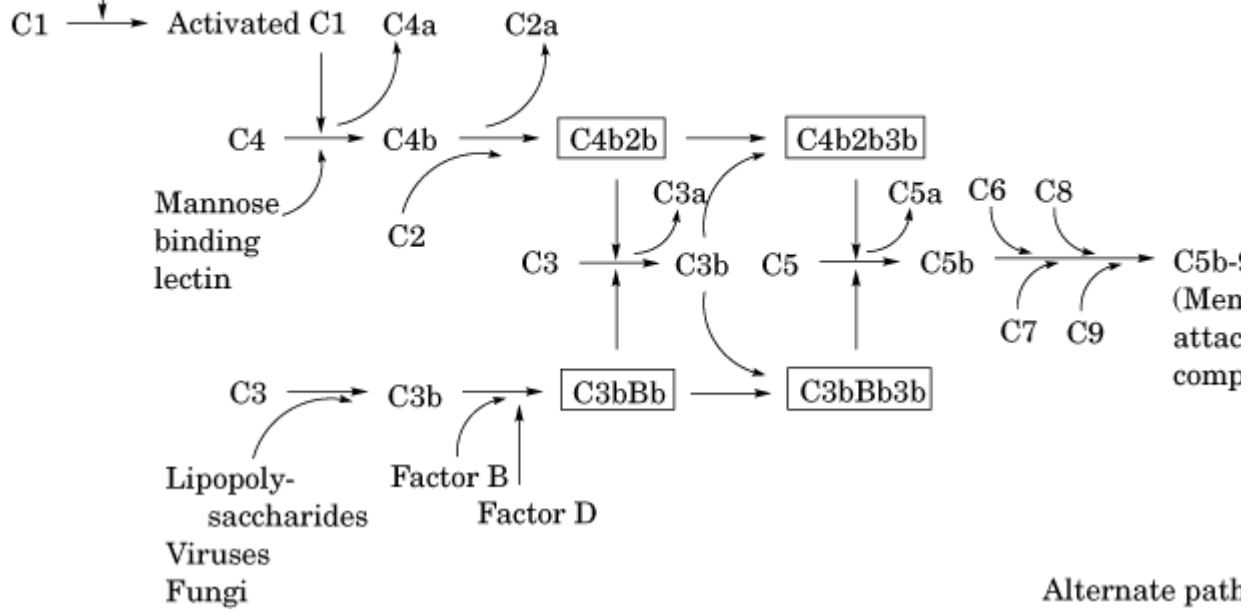


Table 1. Protein Components of the Classical and Alternative Pathways<sup>a</sup>

| Component                | Serum Concentration, $\mu\text{g/mL}$ | Molecular Weight | Number of Polypeptide Chains                             | mRNA <sup>b</sup> Size, kb | Chromosome Location | Gene <sup>b</sup> Size, kb |
|--------------------------|---------------------------------------|------------------|--|----------------------------|---------------------|----------------------------|
| <u>Classical Pathway</u> |                                       |                  |  |                            |                     |                            |
| C1q                      | 75                                    | 410,000          | 6 of A at 24,000<br>6 of B at 23,000<br>6 of C at 22,000 | 0.52<br>1<br>1.6           | 1<br>1<br>1         | 2.5<br>2.6<br>3.2          |
| C1r                      | 34                                    | 85,000           | Single chain   | 2                          | 12                  | 10.5 <sup>c</sup>          |
| C1s                      | 30                                    | 85,000           | Single chain   | 2                          | 12                  | 10.5 <sup>c</sup>          |
| C4                       | 450                                   | 210,000          | a-Chain, 93,000; b-chain, 75,000; g-chain, 33,000        | 5.3                        | 6                   | 16                         |
| C2                       | 25                                    | 95,000           | Single chain   | 2.9                        | 6                   | 18                         |
| C3                       | 1500                                  | 195,000          | a-Chain, 110,000; b-                                     | 5.2                        | 19                  | 41                         |

|  |     |                     |   |         |    |                  |
|--|-----|---------------------|---|---------|----|------------------|
|  |     |                     | chain,<br>75,000                                |         |    |                  |
| C5                                       | 75  | 180,000             | a-Chain,<br>115,000; b-<br>chain,<br>75,000     | 5.5     | 9  | 80               |
| C6                                       | 60  | 128,000             | Single chain                                    |         | 5  | 85               |
| C7                                       | 60  | 121,000             | Single chain                                    | 3.9     | 5  | 78               |
| C8                                       | 80  | 150,000             | a-Chain,<br>64,000;                             | 2.5     | 1  | 70               |
|  |     |                     | b-chain,<br>68,000;                             | 2.6     | 1  | 40               |
|  |     |                     | g-chain,<br>22,000                              | 1       | 9  | 1.8              |
| C9                                       | 58  | 79,000              | Single chain                                    | 2.4     | 5  | 80               |
| Mannose<br>binding Lectin                | 1   | ~ 600,000           | Multimer of<br>32,000                           | 3.5     | 10 | 7                |
| <u>Alternative Pathway</u>               |     |                     |   |         |    |                  |
| Properdin                                | 25  | 220,000             | 4 at 56,000                                     | 1.6     | X  | 6.4              |
| Factor B                                 | 225 | 100,000             | Single chain                                    | 2.9     | 6  | 6                |
| Factor D                                 | 1   | 25,000              | Single chain                                    | 1       | 4  | 2.5 <sup>d</sup> |
| <u>Inhibitors—Soluble</u>                |     |                     |   |         |    |                  |
| Factor I                                 | 34  | 105,000             | 46,000  | 2.4     | 4  | 63               |
| Factor H                                 | 500 | 150,000             | Single chain                                    | 4.4     | 1  | 7                |
| C1 inhibitor                             | 275 | 105,000             | Single chain                                    | 1.8     | 11 | 17               |
| C4 binding<br>protein                    | 150 | 560,000             | 6, a-Chain,<br>70,000; 1,<br>b-chain,<br>45,000 | 2.5     | 1  | 40               |
|  |     |                     |   | 1       | 1  | 10               |
| S protein                                | 500 | 83,000              | Single chain                                    | 1.6     | 17 | 5.3              |
| <u>Membrane—Inhibitors and Receptors</u> |     |                     |   |         |    |                  |
| DAF (CD55)                               |     | 70,000              | Single chain                                    | 3.1     | 1  | 40               |
| MCP (CD46)                               |     | 58,000–<br>63,000   | Single chain                                    | 4.2     | 1  | 43               |
| CR1 (CD35)                               |     | 190,000–<br>280,000 | Single chain                                    | 7.3–13  | 1  | 130–<br>160      |
| CR2 (CD21)                               |     | 140,000             | Single chain                                    | 5       | 1  | 20               |
| CR3<br>(CD11b/CD18)                      |     | 260,000             | a-Chain,<br>165,000;                            | 4.1     | 16 | 55               |
|  |     |                     | b-chain,<br>95,000                              | 3       | 21 | 32               |
| CR4<br>(CD11c/CD18)                      |     | 245,000             | a-Chain,<br>150,000;                            | 4.7     | 16 | 25               |
|  |     |                     | b-chain,<br>95,000                              | 3       | 21 | 32               |
| MIRL (CD59)                              |     | 18,000              | Single chain                                    | 0.6–6.0 | 1  | 27               |
| C5aR (CD88)                              |     | 42,000              | Single  | 3.0     | 19 | 9                |

|      |        |   |     |    |   |
|------|--------|---|-----|----|---|
| C3aR | 50,000 | Single chain <sup>e</sup><br>chain <sup>e</sup> | 3.0 | 12 | 8 |
|------|--------|---|-----|----|---|

<sup>a</sup> All data for proteins listed are for human complement components.

<sup>b</sup> Approximate size.

<sup>c</sup> C1r size estimated from proximity and homology to C1s.

<sup>d</sup> Complete factor D cDNA hybridizes within a 2.5-kb genomic fragment.

<sup>e</sup> Single chain with 7-transmembrane domains.

During activation, through either the classical or alternative pathway, various cleavage products mediate a variety of biological functions. As discussed, cytolysis is mediated through the membrane attack complex. Opsonization and subsequent phagocytosis by leukocytes is mediated primarily through C3b, a cleavage fragment of C3 that binds receptors for C3b. Cleavage products of C3, C4, and C5 (C3a, C4a, and C5a anaphylatoxins) mediate inflammatory events and the recruitment and activation of leukocytes.

### 3. Assembly and Activation of the Classical Pathway

The classical pathway is activated by antigen–antibody complexes or aggregated immunoglobulin acting on complement component C1 (Table 1). C1 consists of three distinct protein molecules, C1q, C1r, and C1s, which are held together by Ca<sup>2+</sup>-dependent bonds. C1 is present in serum as a firm C1q-C1r-C1s complex, while individual C1 components are found only in pathologic conditions. C1 contains one molecule of C1q and two molecules of both C1r and C1s.

A molecule of C1q is comprised of 18 [polypeptide chains](#) of three distinct types (A, B, and C) which are twisted together comprising one subunit. Each C1q molecule is then made up of such six subunits connected together; once assembled, its structure is similar to that of [collagen](#). The C1q molecule bears the sites that enable the entire C1 molecule to bind the Fc region of [IgM](#) and [IgG](#) immunoglobulin molecules. A single C1 molecule can bind about six IgG molecules. After binding to Ag-Ab complexes, C1q undergoes a conformational change that causes C1r to activate itself by a limited self-cleavage.

Activated C1r cleaves a single [peptide bond](#) in C1s, which then acquires enzymatic activity of its own. This newly generated enzyme is a serine esterase type and mediates the cleavage of component C4. With the activation of the C1s enzyme, the initial process is complete, and the earlier reactants, including antibody, antigen, C1q, and C1r, are no longer needed.

C4 is composed of three nonidentical polypeptide chains and is synthesized as a single-chain precursor in a beta–alpha–gamma orientation. C1 cleaves a single bond on the C4 alpha chain leading to the production of C4a and C4b. This cleavage leads to the formation of a labile binding site in the larger fragment of C4b, which enables it to bind to the activating surface (activator). The C4b alpha chain, like C3, contains an internal *thioester* bond formed between a [glutamic acid](#) and [cysteine](#) residue. Cleavage of the alpha chain of C4 is followed by stress-induced hydrolysis of the thioester bond. This permits the reactive acyl group of the glutamyl residue to form a covalent bond with a reactive hydroxyl or [amino group](#) on the surface of the activator.

C2 cleavage by C1s also generates a labile binding site of unknown chemical composition in the larger C2b fragment, which allows it to bind to C4b. Mg<sup>2+</sup> ions are required for the formation of the C4b2b complex. Formation of the C4b2b complex is not very efficient, as the majority of C2 and C4

molecules entering into this reaction lose their labile binding sites before achieving union with membranes or with each other and diffuse away as inactive reaction products. The newly formed C4b2b is a proteolytic enzyme that assumes the role of continuing the complement reaction cascade, so earlier activating components are no longer required. The C4b2b complex is also called the **classical C3 convertase**, which acts to cleave C3.

The substrate for C4b2b is C3, which is synthesized as a single chain of beta–alpha orientation that is processed after translation. The larger, alpha chain, is cleaved at a single site located near the amino terminus. The smaller resulting fragment C3a (9000  $M_r$  (relative molecular mass)) is a biologically potent peptide that mediates inflammation and will be discussed later. A labile binding site is generated in the larger fragment C3b, which enables the molecule to attach to membranes at sites near, but distinct from, those utilized by antibody and C4b2b. Often described as the “lynchpin,” C3 is the precursor of several biologically active fragments that function by association with the other proteins of the complement system leading to lysis of the target cell cell via either classical or alternative pathways. C3 has numerous distinct binding sites, one of which is the thioester domain that allows the molecule to bind covalently to target sites and particles, such as immune complexes and membrane surfaces.

The chemical site of the C3 thioester has the sequence Gly-Cys-Gly-Glu-Glu-Asn with the Cys and second Glu residues joined by a thioester bond, specifically a b-cysteinyl–g–glutamyl thioester bond. This thioester bond is present in the C3d domain of the alpha chain. With the cleavage of C3 into C3a and C3b, the thioester undergoes a stress-mediated hydrolysis, and the reactive acyl group of the glutamyl residue forms a covalent bond with a reactive hydroxyl or amino group on the activator surface. A major amount of reactive C3 fails to achieve binding with activators as most of these reactive thioesters have reacted with water.

The attachment of C3b to membranes in the vicinity of C4b2b molecules leads to the generation of the last enzyme of the classical pathway, C4b2b3b, **or the classical C5 convertase**. To this end, the classical pathway has finished the initial steps of activation. From here, the C5 convertase acts on C5 to initiate steps in the forming of the MAC.

#### 4. Activation of the Classical Pathway

The activation of the classical pathway can occur by one of two ways: activation through the use of immunoglobulin or through nonimmunoglobulin activation by cleavage of C1. Immunoglobulin activation of the classical pathway involves the use of antibodies belonging to subclasses IgG<sub>1</sub>, IgG<sub>2</sub>, and IgG<sub>3</sub>, as well as IgM, are capable of initiating the classical pathway. Immunoglobulins IgG<sub>4</sub>, IgA, IgD, and IgE are inactive in this regard. Among the three listed above, IgG<sub>3</sub> is the most active in interacting with C1, then IgG<sub>1</sub> and IgG<sub>2</sub>. Activation occurs when the first component C1 binds to a site in the Fc region of the IgG or IgM. Only one molecule of bound IgM is capable of initiating the activation of the classical pathway; however, it is estimated that six molecules of IgG are required.

Nonimmunologic activation of the classical pathway can be accomplished by diverse substances such as DNA, certain viruses, and **trypsin**-like enzymes, by direct proteolytic attack on the C1 molecule. Of particular interest are **lectins** that also act as a nonimmunological activator of the classical pathway. Mannose binding lectin (MBL) is a serum protein found in all mammals and is regulated as an acute-phase protein. MBL is capable of binding carbohydrate moieties present on many microorganisms and subsequently capable of activating the classical pathway through C4 without the need for specific antibody or C1q.

#### 5. Alternative Complement Pathway or Properdin Pathway

The alternative pathway may be activated immunologically by aggregates of human IgA and by certain complex polysaccharides, fungi, viruses, bacterial lipopolysaccharides, and trypsin-like enzymes. This system was originally described as the properdin system, a group of proteins involved in resistance to infection, but distinct from complement. The proteins of the alternative pathway, C3, factor B, factor D, and properdin (factor P) perform the functions of activation, recognition, and amplification of the pathway, resulting in the formation of the activator bound C3/C5 convertase. The alternate pathway system was found to be involved in the destruction of certain bacteria, neutralization of some viruses, and red blood cell (RBC) lysis from patients with paroxysmal nocturnal hemoglobinuria.

## 6. Activation of the Alternate Pathway

Alternate pathway activation initially proceeds in a manner different from that for the classical pathway. The exact mechanism by which the first C3b molecule is produced is still controversial. However, it is clear that antibody is not required since mixtures of only highly purified alternative pathway proteins behave as well as serum. The ultimate aim of the proteins, C3, factor B, factor D, and properdin is the initiation, recognition, and amplification of the pathway that results in the formation of the activator bound **C3/C5 convertases**. The first requirement for activation is the presence of C3, which participates in the initiation and amplification of the pathway. It is generally accepted that activated C3 (but not-yet cleaved) [C3\*] reacts with proenzyme B, and this complex is then cleaved by factor D, yielding two fragments, Ba and Bb. Bb then attaches to C3b forming the alternative pathway C3 convertase. The C3 convertase is able to cleave additional C3 into C3a and C3b. While most C3b remains in the fluid phase, some binds to various cell surfaces. Thus C3b is continuously generated and is deposited on the surface of the activator. The alternative C3 convertase is unstable and decays rapidly unless another alternative pathway member, **properdin**, binds to the C3 convertase and stabilizes it. Properdin binding of C3bBb extends its functional half-life eightfold. The function of the C3 convertase is to produce enough meta-stable C3 to deposit on the surface of surrounding particles. Surface-bound C3b then interacts with factor B and is activated by D to form more C3bBb. This enzyme is capable of cleaving large amounts of C3, and forms more C3bBb (C3bBb3b); thus, a positive feedback mechanism is formed that amplifies the initial stimulus. Many of the C3b molecules bind to the surface of the activator in close proximity to these enzymes, forming C3bBb3b or the alternative pathway C5 convertase. Like its counterpart in the classical pathway (C4b2b3b), the function of C3bBb3b is to proteolytically cleave C5 and initiate the assembly of the membrane attack complex.

## 7. The Membrane Attack Complex: C5–C9 Reaction

The terminal portion of the complement sequence is known as the MAC (complex). C5–C9 must become membrane bound for damage to occur. This complex may attach to a cell bearing the activation enzymes of the classic or alternate pathways, or it may attach to a bystander cell. The MAC is initiated by the cleavage of C5 by the alternative and classical C5 convertases.

Similar to C3, pro-C5 is secreted in a beta–alpha orientation prior to [post translational modification](#). C5 shares sequence [homology](#) with C3 and C4, including the domain corresponding to the thioester region; however, C5 lacks the essential cysteine and glutamic acid residues necessary for thioester formation. Assembly of the MAC initiates with the cleavage of C5. Cleavage at residue 74 of the alpha chain, results in the generation of C5a anaphylatoxin and C5b (Table 1), where C5b can then bind C6. C6 is a single chain peptide that, when bound to C5b, remains loosely attached to the membrane until the addition of C7. Once a single C7 has bound the C5b6, the resulting complex is highly lipophilic and inserts into the lipid bilayer. When inserted, the C5b67 complex serves as a high affinity integral membrane receptor for a single C8 molecule. C8 is composed of three nonidentical polypeptide chains and has an interesting structure in that the alpha and gamma chains are covalently linked by one or more disulfide bonds and the beta chain is joined to the alpha–gamma chains by noncovalent forces. The C5b-8 becomes stably attached to the membrane by insertion of the C8 gamma chain into the lipid bilayer. With this early “MAC” formed, cellular

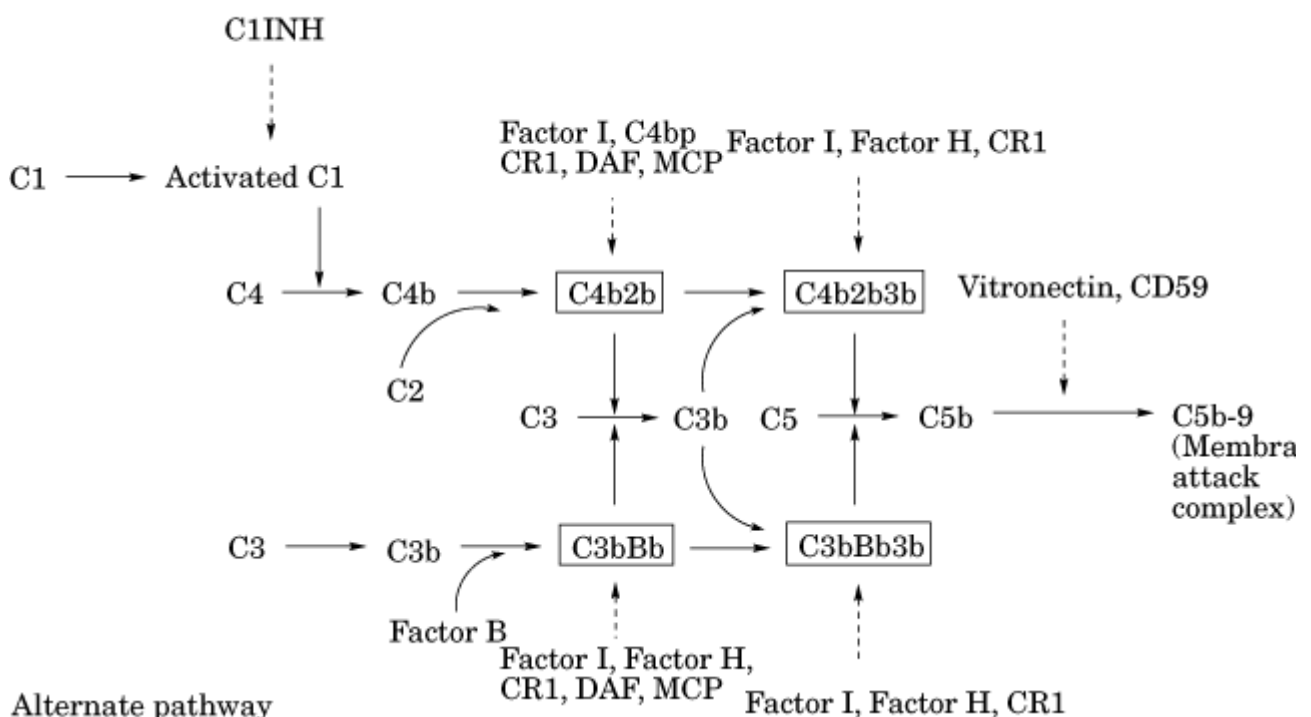
leakage can occur at this stage. The cytolytic process is greatly accelerated by the attachment of C9, the final complement component. The C9 monomer polymerizes at the site of the C5b-8 complex. As few as 4 C9 molecules can result in lysis for many microorganisms. However, when the MAC contains as many 12 to 15 C9 molecules, the “poly-C9” forms a pore in the membrane and permits passive exchange of cytoplasmic contents with the extracellular media with subsequent osmotic lysis. The pores have an internal diameter of about 110 Å and appear similar in structure to those formed by lytic T cells when the related protein **perforin** is deposited on target cells.

## 8. Control Mechanisms of the Complement System

Given the lytic potential of the complement system, uncontrolled activation can lead to the formation of the MAC on normal tissues and significant generation of C3a, C4a, and C5a anaphylatoxins. Multiple fluid phase and membrane bound proteins tightly regulate both the alternative and classical pathways by acting at specific steps during activation. In addition, the uncontrolled activation of the complement cascade is also prevented by the necessity of specific activators, such as antibody–antigen complexes, as well as the lability of the activated complement components such as C3b. Deficiencies in regulatory proteins result in unregulated complement activation and are apparent in various disease conditions. An overview of the regulation complement pathway activation by both membrane-bound and soluble inhibitors are shown in Figure 2. The biochemical and physical information of these inhibitors are listed in Table 1.

**Figure 2.** Regulation of classical and alternative complement pathway activation by membrane and soluble inhibitors. Principal components and C3/C5 convertases are retained from Figure 1. Regulators of complement activation are shown in italics.

### Classical pathway



## 9. Classical Pathway Inactivators

A protein called *C1 inhibitor* (**C1INH**) acts to inhibit the initiation of the classical pathway. C1INH



is a member of the serine protease inhibitor family and acts by blocking C1r and C1s of the C1 complex. As with other serpins, such as  $\alpha_1$ -antitrypsin, C1INH uses a “bait” sequence that mimics the normal substrates of C1r and C1s. However, when C1INH is cleaved by C1r or C1s, it forms a covalent ester linkage that inhibits cleavage of C4 and C2, thus inhibiting the progression of the classical pathway. Most of the C1 in the serum is already bound to C1INH. On specific binding of C1 to antigen–antibody complexes, C1INH is released allowing classical pathway activation to proceed.

Other means to control the classical pathway involve interfering with the formation of the C3 convertase. Such inhibition can be mediated either through interfering with the assembly of the convertase or in promoting the dissociation of the convertase. Three proteins called ***C4-binding protein (C4 bp)***, ***complement receptor type I (CRI)*** (also known as CD35), and ***decay accelerating factor (DAF)*** (also known as CD55) function as inhibitors of the classical pathway by binding to C4b. In doing so, they compete with C2 for C4b and thus inhibit the formation of the classical C3 convertase. In binding C4b, they also promote the dissociation of the C3 convertase. All three proteins are structurally similar and are homologous members of a family of proteins called the regulators of complement activity (RCAs).

Another way that the C3 convertase can be inhibited is by the cleavage of C4b by a protein called **factor I**. Factor I is a serine esterase that cleaves C4b into two additional fragments: C4c and C4d. C4d remains bound to the original activating surface; however, it is unable to contribute to the formation of the classical C3 convertase. Cleavage of C4b by factor I can only occur in the presence of a “cofactor.” The proteins that can serve a cofactor for factor I cleavage include CR1, C4b, and membrane cofactor protein (MCP), or CD46. MCP is also a member of the RCA family but does not possess decay accelerating activity.

## 10. Alternative Pathway Inactivators

Regulation of the alternative pathway is also accomplished by several circulating and membrane bound proteins, many of which also act on the classical pathway. In analogous fashion, the inhibition of convertase activity through either decay acceleration or by providing cofactor activity for factor I are the major ways in which this is accomplished.

**Factor H** is a soluble serum protein that possesses decay-accelerating activity that is specific for the Bb fragment of the alternative pathway. By competitively binding Bb, the formation of the alternative C3 convertase is impaired. Factor H is also a member of the RCA family and in binding Bb can also promote the disassembly of the alternative C3 convertase.

In addition to its role in inhibiting the classical pathway activation, factor I can also cleave C3b using factor H, MCP, and CR1 as cofactors. The resulting fragment of C3b (iC3b) is unable to participate in the formation of the alternative C3 convertase.

## 11. Inactivators of the Membrane Attack Complex

With the alternative and classical pathways activated, the potential for normal or bystander host cells to be subjected to lysis is significant. Excessive cell lysis is normally prevented by a number of proteins that act to inhibit the formation of the membrane attack complex.

A membrane bound protein that is capable of inactivating the membrane attack complex is the ***membrane inhibitor of reactive lysis (MIRL)***, more commonly known as **CD59**. CD59 is broadly expressed among different cell types and is probably the most important molecule in minimizing damage to bystander cells. It is thought that CD59 acts by binding C8 and C9 and thus inhibits the polymerization of C9 as well as the insertion of the MAC into the membrane.

As discussed above, when C7 binds C5b6, the resulting complex is lipophilic and it inserts into the membrane. This process of membrane insertion is inhibited by **S protein**, also known as **vitronectin**, a protein related to fibronectin and laminin. The S protein is thought to function by binding the C5b67 complex before it inserts into the membrane.

## 12. Complement Receptors and Biological Consequences of Complement Activation

Through the activation of the complement cascade, through either the classical or alternative pathways, numerous complement derived fragments are generated. Many of these fragments possess biological activities that are mediated through the use of specific receptors.

In addition to regulating the activation of complement, **CR1 (CD35)** serves as a high affinity receptor for C3b and C4b. CR1 functions as a receptor for C3b- and C4b-coated particles wherein macrophages ingest and remove these coated particles. Another receptor, **CR2 (CD21)**, is expressed primarily on B cells, dendritic cells, and epithelial cells. CR2 is known to specifically bind iC3b and C3dg—both factor I-generated cleavage products. In addition, CR2 is also known as the Epstein–Barr virus receptor. **Mac-1** or **CR3 (CD11b/CD18)** is related to the integrin family and is a specific receptor for the factor I cleavage product of C3b, iC3b. CR3 is expressed on many different cells derived from the bone marrow including neutrophils, mononuclear phagocytes, mast cells and natural killer (NK) cells. The primary role of CR3 is the phagocytosis and clearance of iC3b-coated particles. **CR4 (CD11c/CD18)**, like CR3, is also a member of the integrin family and also binds iC3b as well as C3dg.

The cleavage of C3, C4, and C5 results in the production of **C3a**, **C4a**, and **C5a anaphylatoxins**. These hormone-like peptides mediate a variety of cellular and biochemical responses from leukocytes, including the generation of bacteriocidal superoxide radicals, proteolytic enzyme release from intracellular granules, cellular aggregation, smooth muscle contraction, and phagocytosis. In addition, C5a induces chemotaxis, or the migration, of leukocytes into areas of complement activation. The induction of chemotaxis was thought to be unique for C5a; however, recent findings suggest that C3a may be chemotactic for eosinophils, a cell whose function is mediating allergic reactions. These anaphylatoxins mediate their effects through specific receptors that are coupled to GTP-binding proteins and whose signal transduction can be abrogated by pertussis toxin. The C3a and C4a receptors have been found to be expressed on mast cells, basophils, lymphocytes, and smooth muscle cells. The C5a receptor has been found on mast cells, basophils, neutrophils, monocytes/macrophages, endothelial cells, smooth muscle, liver parenchyma, and epithelial lining of lung as well as astrocytes and microglia of the central nervous system. Of the three anaphylatoxins, the most potent is C5a in mediating biological effects. C3a is approximately 20-fold less potent, while C4a is about 2500-fold less potent than C5a. C5a has been shown to stimulate the release of [tumor necrosis factor](#) (TNF) from mast cells as well as stimulating P-selectin (CD62P) expression on vascular endothelium to promote neutrophil binding.

## 13. Complement Genes and Relationships

On the basis of sequence homologies, many of the complement proteins can be grouped together as members of families. Members of these families share structural and/or functional similarities by which relationships can be assessed.

One of the more defined families is that of the regulators of complement activity (RCAs) as previously discussed. Factor H, CR1, CR2, DAF, C4 bp, and MCP are all members of the RCA, and all share the ability to bind both C3b and C4b. At the peptide level, all members of this family possess multiple, tandemly arranged repeated structures called short consensus repeats (SCRs). Each SCR is 65–70 amino acids long with 11–14 conserved amino acid residues. The composition and number of SCRs each of these members contain vary, but structurally and functionally they form the foundation by which this family is based. The genes encoding H, CR1, CR2, DAF, C4 bp, and MCP are found within an 800-kb genomic segment on the long arm of chromosome 1.

Another group of complement genes that are linked together are C2, factor B, and C4 that all map within the major histocompatibility complex of both humans and mice. In the human, these complement genes, often called the Class III genes, map between the Class II HLA-DR and Class I B loci of chromosome 6. Similar to the Class I and Class II genes, the Class III genes are polymorphic as multiple alleles exist for each of these genes.

Complement components C3, C4, and C5 constitute a structurally homologous family of proteins that includes  $\alpha_2$ -**macroglobulin** and pregnancy zonal protein. All except C5 are characterized by the presence of an internal thioester bond allowing these proteins to covalently bind with cell surfaces or other proteins. However, unlike the RCA family, C3, C4, and C5 reside on separate chromosomes (Table 1).

#### 14. Pathologies Related to the Complement System

The complement cascade is a potent and powerful system whose function is the destruction of target cells and infectious agents. Unwanted deposition of complement as well as the production of significant quantities of inflammatory mediators can result in substantial damage to normal healthy tissues. The regulatory mechanism of the MAC and alternative and classical pathways establishes a fine balance between activation and inhibition so as to protect autologous cells but allow the destruction of foreign agents. Deficiencies in nearly any component of the complement system compromise the health of an individual by tipping the balance one way or another. In a general sense, complement deficiencies can be due to the absence of a component due to structural defect in the gene, or the production of a component that is incapable of functioning in its specific role in the cascade. Table 2 summarizes complement deficiencies, the resulting abnormalities, and associated pathologies.

**Table 2. Complement Deficiencies and Associated Clinical Abnormalities**

| Component                      | Biological Defect  | Associated Disease                                       |
|--------------------------------|--|--|
| <u>Classical Pathway</u>       |  |  |
| C1q, C1r, C1s                  | Defective classical pathway                                | Systemic lupus erythematosus (SLE), bacterial infections |
| C4, C2                         | Defective classical pathway                                | SLE, bacterial infections, glomerulonephritis            |
| C3                             | Defective classical and alternative pathways               | Bacterial infections, glomerulonephritis                 |
| <u>Alternative Pathway</u>     |  |  |
| Properdin, factor D            | Defective alternative pathway                              | Bacterial infections                                     |
| <u>Membrane Attack Complex</u> |  |  |
| C5, C6, C7, C8                 | Defective MAC assembly                                     | Recurrent neisserial infections                          |
| <u>Inhibitors</u>              |  |  |
| Factor I, factor H             | Deregulation of complement activation<br>Consumption of C3 | Bacterial infections                                     |
| CR3 (CD11b/CD18)               | Impaired opsinization                                      | Bacterial infections                                     |

|              |                                   |                                       |
|--------------|-----------------------------------|---------------------------------------|
| C1 inhibitor | Deregulation of classical pathway | Hereditary angioneurotic edema (HANE) |
| DAF (CD55)   | Deregulated C3 convertase         | Paroxysmal nocturnal hemoglobinuria   |
| MIRL (CD59)  | Impaired MAC regulation           | Paroxysmal nocturnal hemoglobinuria   |

---

Deficiencies in alternative and classical pathway components usually manifest clinically in recurrent bacterial infections. The most serious deficiencies are those in C3. Given its role in opsonization, phagocytosis, and lysis of bacteria, homozygous deficiencies in C3 often prove fatal. In contrast, C9 deficient individuals cannot effectively generate MAC formation, yet these patients have minimal or no associated pathology as the addition of C8 to the C5b-7 complex can result in osmotic lysis of target cells. Individuals deficient in the early classical pathway components (C1, C2, and C4) are impaired in their ability to solubilize and clear immune complexes, resulting in local inflammation associated with autoimmune diseases such as *systemic lupus erythematosus* (SLE).

Patients that are deficient in the terminal complement components (C5 through C9) have impaired ability to assemble the MAC. Of particular interest is that these patients seem to have a propensity for Neisseria bacterial infections suggesting that cytolysis may be the major defense mechanism against these bacteria.

Deficiencies in the soluble and membrane bound regulatory components of complement result in abnormal activation and deposition of complement. Defects in C1INH results in a condition called *hereditary angioneurotic edema* (HANE), in which the intermittent accumulation of edema (swelling) in the skin and mucosa occurs. The exact mechanism of how the edema occurs is not known; however, without C1INH to inhibit non-specific C1 activation, C2 and C4 are readily cleaved and their by products are implicated. Deficiencies in the expression of phosphatidylinositol-linked membrane proteins result in a condition called *paroxysmal nocturnal hemoglobinuria* (PNH), where patients suffer from recurrent bouts of intravascular hemolysis. As DAF and CD59 are PI-linked, they also serve to inhibit C3 convertase formation. However, erythrocytes are especially sensitive to lysis as they do not possess membrane bound form of CD59 as other cell types do.

## 15. Summary

The proteins of the complement system consists of at least 20 chemically and immunologically distinct plasma proteins that can dynamically interact with one another in a highly regulated manner. The complement system can be activated by using either of two converging initiation pathways: the classical pathway, which is activated by antibody–antigen complexes; and the alternative pathway, which is activated by the surface of invading organisms. This interaction of proteins results in (1) the lysis of target cells and bacteria, (2) the binding and clearance of immune complexes throughout the body, (3) the binding of complement based proteins to foreign particles and subsequently cleared by phagocytosis, and (4) the generation of inflammatory mediators to enhance the humoral and cellular immune response. The activation of complement is tightly regulated so as to protect host cells from nonspecific attack deposition while allowing the destruction of foreign agents. Deficiencies in complement components of either pathway or in regulatory components are associated with recurrent bacterial infections and autoimmune diseases.

## 16. Acknowledgments

This is publication 133-IMM from the Institute of Molecular Medicine for the Prevention of Human

Diseases, University of Texas-Houston Health Science Center. This work was supported by GM56050 (DLH) and AI25011 (RAW). This article is dedicated to the memory of Dr. Hans Müller-Eberhard.

### Suggestions for Further Reading

- A. K. Abbas, A. H. Lichtman, and J. S. Pober (1994) *Cellular and Molecular Immunology*, 2nd ed., Saunders, Philadelphia.
- H. Müller-Eberhard (1986) The membrane attack complex, *Annu. Rev. Immunol.* **4**, 503–528.
- H. Müller-Eberhard (1988) Molecular organization and function of the complement system, *Annu. Rev. Biochem.* **57**, 321–347.
- R. A. Wetsel and H. R. Colten (1990) in *Inheritance of Kidney and Urinary Tract Diseases*, A. Spitzer and E. D. Avner, eds., Kluwer Academic, Chapter "18", pp. 401–429.
- G. D. Ross, ed. (1986) *Immunobiology of the Complement System*, Academic Press, San Diego, CA.
- K. Rother and G. O. Till, eds. (1988) *The Complement System*, Springer-Verlag, New York.
- K. Rother and U. Rother, eds. (1986) *Hereditary and Acquired Complement Deficiencies in Animals and Man*, *Progress in Allergy*, Vol. 39, Karger.

## Complementary DNA (cDNA)

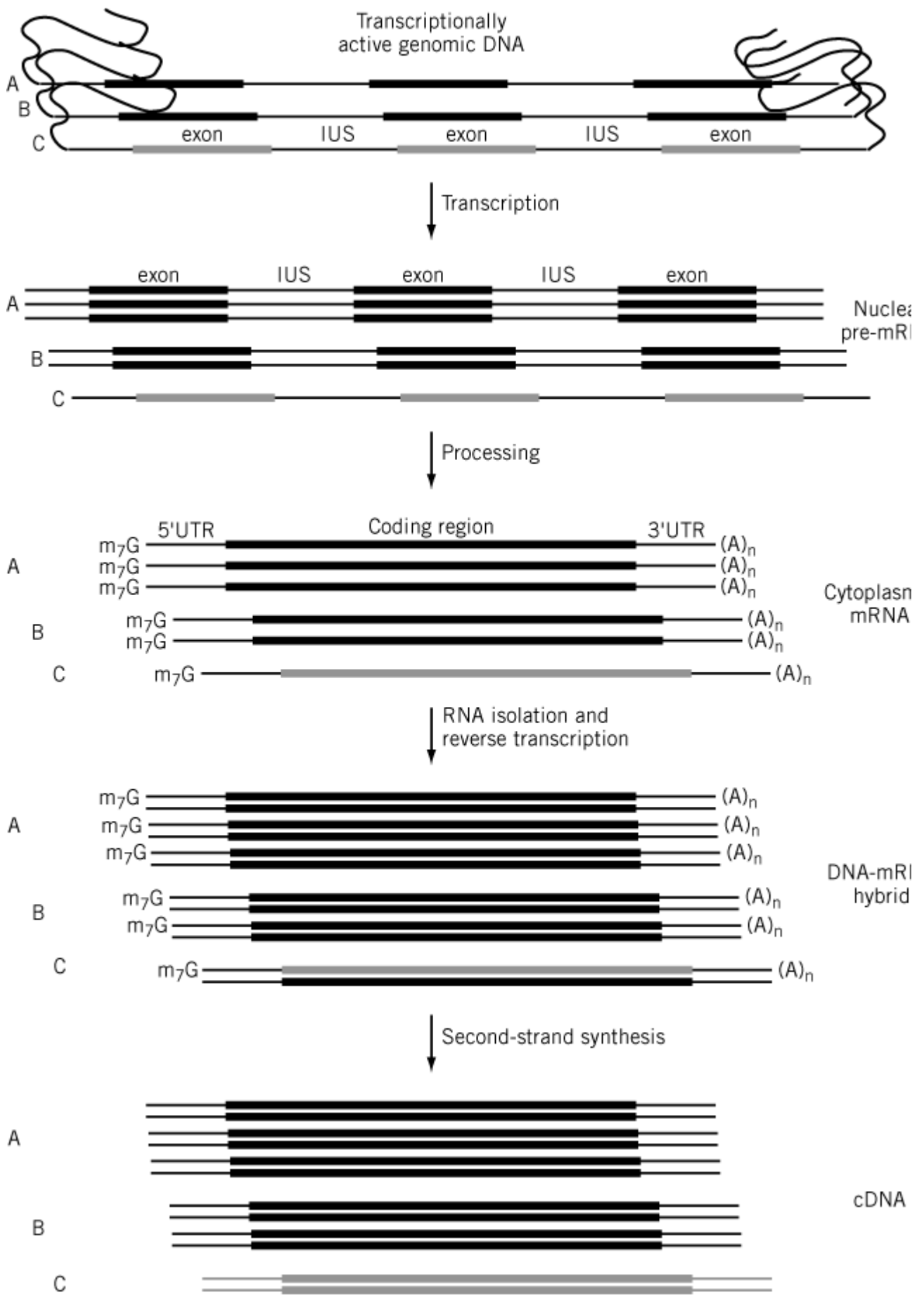
The base-pair complementarity between a **gene** and its **transcript** leads to the formation of a DNA–RNA hybrid under experimental conditions that favor specific **hybridization**. Pioneering work by Gillespie and Spiegelman that made use of the complementarity between viral genomic DNA and RNA in infected cells ([1](#)) helped establish the beginnings of a new technology for the detection of specific RNA molecules hidden among an RNA background, in this case of viral RNA among host cell RNA. The potential advantages of *complementary DNA* (cDNA) as a probe for RNA detection were recognized before techniques existed to make specific cDNA conveniently for tagging genes. The great abundance in the genome of extra, nontranscribed DNA, which is a common feature of higher organisms, complicates the use of genomic DNA as a cDNA probe. Some time after the first references to cDNA, a **DNA polymerase** activity that depends on the presence of an RNA **template** was observed in RNA **viruses** ([2](#), [3](#)). This discovery, in addition to extending our knowledge of RNA virus replication intermediates, led to the identification and isolation of a previously unknown form of DNA polymerase called *reverse transcriptase* and to the identification of buffer conditions that allow some DNA-dependent DNA polymerases to use a DNA-primed RNA template to form a DNA product ([4](#), [5](#)). Largely due to the discovery of reverse transcription and the technical advantages it provides for *in vitro* cDNA synthesis, cDNA has become a fundamental research tool for use with higher organisms.

### 1. cDNA Copies of messenger RNA

Transcriptionally active genes represent only a relatively small fraction of the total host genomic DNA in higher organisms, only 3% in human cells, with the noncoding or nonexpressed DNA representing the majority of the genome. Reverse transcribed cDNA synthesized under carefully controlled conditions is a faithful and stable double-strand DNA copy of cellular RNA and is trimmed of excess genomic sequence. A full-length cDNA molecule made from **messenger RNA** (mRNA) contains only the protein-coding region of an expressed gene, together with the adjacent untranslated regulatory sequences; it lacks any unexpressed gene or genomic sequences not

contained in the mature mRNA, such as **introns** (Fig. [1](#)).

**Figure 1.** Schematic representation of cDNA synthesis from mRNA. Transcriptionally active genes are transcribed, resulting in pre-mRNA that is processed to the mature form. Nuclear processing includes the removal of intervening sequences (IVS, introns) and addition of a 5' m<sub>7</sub>G cap and 3' poly dA tail. Different gene activities result in different mRNA abundance shown as A, B, or C. After mRNA isolation and reverse transcription, a population of DNA–mRNA hybrid molecules are formed where the DNA is anti-sense in relation to the mRNA. Second-strand synthesis produces a cDNA population in which the more active gene is represented more than the less active gene.



For research that involves mRNA, the preparation of cDNA not only provides a DNA strand complementary to RNA, but it also serves to minimize the technical problems of working directly with RNA, which is chemically unstable at alkaline pH and susceptible to notoriously stable and ubiquitous cellular ribonucleases. A further advantage of cDNA is that in many applications it is integrated into a **prokaryotic DNA vector** to form [recombinant DNA](#) capable of autonomous replication when **transformed** into a host organism with a favorable genetic background. Engineered vectors with a variety of features are available for many different methods of cDNA analysis. As recombinant DNA, cDNA may be amplified through culturing of the host cells. Because each transformed host cell receives a single recombinant vector, clonal lines of host cells differing only in the recombinant vector received are established. As a copy of total cellular mRNA, however, cDNA is a complex mixture representing the complexity of the template mRNA used to direct its synthesis. Thus, cDNA integrated into a vector, transformed into a host, and cultured represents a complex mixture or **library** of transcripts, a [cDNA library](#). The use of cDNA frequently requires the selection of a single **clone** or cDNA copy from a library, since the analysis of expression from a single gene is often the desired aim. The efficient selection of a previously uncharacterized cDNA clone from an entire cDNA library poses a major technical challenge. Selection of a new cDNA is complicated by the lack of any previous information to assist in the selection process. Often, the selection might start with little information other than observations about the encoded function of the desired cDNA clone or atypical cellular characteristics that could be associated with it.

As a stable copy of RNA, a single selected cDNA clone may be transcribed through vector signals to regenerate the RNA, which may in turn be used for synthesis of the encoded protein, either within a host cell or through the use of **protein biosynthesis *in vitro*** system. The purity of RNA or protein prepared from a single cDNA clone is absolute, thus justifying the effort to prepare a cDNA library and select a desired clone. This is an alternative to the task of biochemical purification of a protein from a tissue or cell extract. A purified cDNA in a regulated protein expression system may be used to provide a continuous renewable protein source for analysis of the structure, activity, or any characteristic of the encoded protein. Also, the effect that a cDNA's expression has on a relevant **eukaryotic** cell line may provide information about the biological role of the encoded function. Thus, cDNA-directed protein expression not only provides a source of pure protein for either *in vitro* or *in vivo* analyses, but it may also provide a method of cloning (selection) of a cDNA from a cDNA library. Many selections of new cDNAs start without protein or gene sequence information, but instead rely on the detection of the function or biological effect. This process is known as direct expression or functional cloning. Any direct expression cloning method mimics the classic genetic methods of rescue or [complementation](#) and so could be considered genetics in the absence of mating. To describe functional cloning in a general way, an altered host cell's **phenotype** or measured activity is restored to the wild-type when the genetic material encoding the host's deficiency is supplied. Although experimental approaches that rely on classic genetic methods have produced much information about cellular processes, they depend on good genetic systems to which the genetic material can be delivered. Systems that lend themselves to an imposed manipulation and selection, generally single-cell organisms with a short generation time, are best for this purpose. A myriad of eukaryotic cell cultures that have been purposefully altered or carefully isolated from the affected organs of diseased patients and that show an atypical phenotype have been widely used as host cells for functional cloning. The applications for the functional cloning of new cDNAs from higher organisms also depend on a selectable phenotype or activity to isolate the host cell that harbors the relevant genetic material. As an example, the anchorage-dependent phenotype of nearly all primary tissue cells is not exhibited by many tumor cells that will grow in suspension (6). When the genetic material used is cDNA, the details of the selection method must be designed so that the cDNA library contains the cDNA encoding the deficient host activity, and so that the cDNA is incorporated into the host cells and expressed functionally. With expression of a cDNA library in eukaryotic host cells, the selection of a desired clone depends on acquisition of a new phenotype. In the example of the anchorage-independent phenotype exhibited by some tumor cells, the induction of the anchorage dependency allows selection for a cDNA-encoded activity related to the anchorage phenotype and possibly to the cause of tumor formation (7, 8). For any selection where the desired



cDNA is of an mRNA that is in low abundance, the selection will require a large-scale transformation of host cells, in order to produce even a single positive clone.

Acquiring a cDNA opens up many types of experimental approaches for the analysis of gene expression and studies of the protein encoded by the cDNA. As an example, protocols are available to target a single codon in a cDNA for [site-directed mutagenesis](#). By using available methods to express the altered cDNA, the effects of the mutation on the encoded protein's activity may be measured, and in this way, a structure–function relation is probed.

Although the methods to prepare and make use of cDNA are advanced and commercially available, it remains a challenge to select a completely new cDNA clone from a cDNA library. This may limit the uses of cDNA technology to better characterize known cellular processes rather than uncharacterized cellular processes. The successful recovery of a novel cDNA depends on a specific selection strategy that ultimately will test each cDNA for an encoded activity. The task of individually testing each clone in a cDNA library is nearly overwhelming. If a screening method designed to test each cDNA in a collective manner cannot be designed, it forces the need for biochemical methods to purify the protein activity in order to provide enough sequence information to obtain the cDNA by more conventional methods, such as library screening with a sequence probe or selective **PCR** library amplification with degenerate oligonucleotide primers. Evolving methods of protein expression have demonstrated the potential of cDNA selection as a possible alternative to protein purification. However, functional cloning may not always be possible. An alternative approach, which can temporarily bypass the challenge presented by the need to design or apply a new isolation strategy for each new cDNA, is to randomly isolate and analyze cDNA clones first and then subsequently to identify the encoded function. The use of sequence information derived from randomly selected cDNA samples has become part of the Human Genome Project as a way of providing an aid for the physical mapping of expressed genes, as gene tags ([9](#), [10](#)). The cDNA sequences obtained are collected into the [expressed sequence tag](#) (EST) database, which is a collection of nonrepetitious human cDNA sequences of randomly selected cDNAs from 250 cDNA libraries from RNA isolated from 37 distinct human organs and tissues ([11](#)). The EST information is providing genetic markers for physical gene mapping and also information about the expression pattern of the mapped genes as the cell lineage is known for each sample analyzed. Some hints about the encoded function of the mapped genes may be revealed by the expression pattern. Furthermore, the large volume of sequence information produced by random sequencing may be compared to the available sequence information of the growing number of fully sequenced genes, in order to identify the presence of conserved protein domains with established activities in a process that has come to be known as *functional genomics* ([12](#), [13](#)). This process of extending available information about function to randomly acquired new sequence information, which necessarily requires a unified, large-scale approach, will undoubtedly contribute greatly to the identification of the encoded functions of previously unknown human genes.

## 2. Preparation of cDNA

### 2.1. Isolation of RNA Template

The preparation of cDNA begins with the isolation of the RNA that will serve as a template to direct the cDNA synthesis. Since cDNA is representative of the RNA template, an inciteful choice of a cellular source of RNA, such as a tissue that is known or at least thought to contain the specific mRNA of interest, is important. In some cases, the RNA might be tested for the presence of an RNA of interest. In a simple case, the RNA sample could be analyzed prior to use for cDNA preparation by use of a sequence probe. This sort of analysis requires, of course, some definite sequence information about the mRNA of interest. In a more complex case, when no sequence information about the mRNA is available, but detection of the encoded protein is possible, the RNA source could be analyzed for the ability to direct synthesis of the protein. Detection of a protein may be accomplished by any known function, such as **enzymatic** activity or **ligand binding** activity. The presence of a protein in a cell does not insure the presence of the protein-encoding mRNA, as the protein could be present as a stable form from an earlier expression of an mRNA that is now

degraded. On the other hand, a mRNA encoding a protein of interest could be present in a tissue, but remain untranslated. One technique that can be used to insure the presence of the mRNA of interest in an RNA preparation is to translate the total RNA in an *in vitro* translation reaction and then assay the total resulting protein product for the expected protein. The two high-activity *in vitro* translation systems in common use are those prepared from rabbit reticulocytes and wheat germ (14, 15).

Given that a tissue source can be identified that contains the RNA of interest, there may be a possibility of enriching for the desired mRNA. Only a small fraction of total cellular RNA is mRNA that is directed to ribosomes for translation into protein (Table 1). The majority of cellular RNAs function in other ways to aid translation, often with structural roles. Enrichment methods aim to remove the more abundant structural RNAs selectively from the mRNAs. Enrichment of polyA + mRNA by hybridization with immobilized oligo (dT), followed by stringent washing to remove the abundant, structural, poly A – RNAs, is a common method used to prepare the mRNA template for synthesis of cDNA representing protein-encoding mRNAs (16). To reduce contamination by mitochondrial poly A + RNA, the mitochondria may be removed by cell fractionation prior to total RNA extraction, although the extra time and handling of cells required could lead to some loss or degradation of RNA. Isolation of full-length mRNA of nuclear origin based on the unique m<sup>7</sup>G cap structure might also serve as an alternative to cell fractionation to reduce those RNA messages of mitochondrial origin often present in RNA preparations (17). Translational inhibitors have been used to stabilize ribosome-bound mRNA that is degraded during translation (18), and as selective RNA degradation pathways are increasingly well characterized, there is an increasing promise of specific inhibitors of RNA degradation for use in preparing mRNA that normally has a relatively short half-life.

**Table 1. Cellular Distribution of RNA in HeLa Cells**

| Cellular Location    | % of Total | Average $t_{1/2}$     |
|----------------------|------------|-----------------------|
| <b>Nucleus</b>       | <b>10</b>  |                       |
| hnRNA                | 4.0        | 20 min                |
| snRNA                | 2.5        | 100 hours             |
| pre-rRNA             | 3.5        | 15 min                |
| <b>Cytoplasm</b>     | <b>85</b>  |                       |
| mRNA                 | 2.0        | 8–10 hours            |
| rRNA                 | 75.0       | 107 hours             |
| tRNA                 | 8.0        | 80 hours              |
| <b>Mitochondrion</b> | <b>5</b>   |                       |
| mit-mRNA             | 0.5        | 8–10 hours            |
| mit-rRNA             | 3.5        | 107 hours             |
| mit-tRNA             | 1.0        | 80 hours <sup>a</sup> |

<sup>a</sup> Note: These values are approximate ones based on experimental observations made during repeated isolation of RNA from HeLa cell cultures.

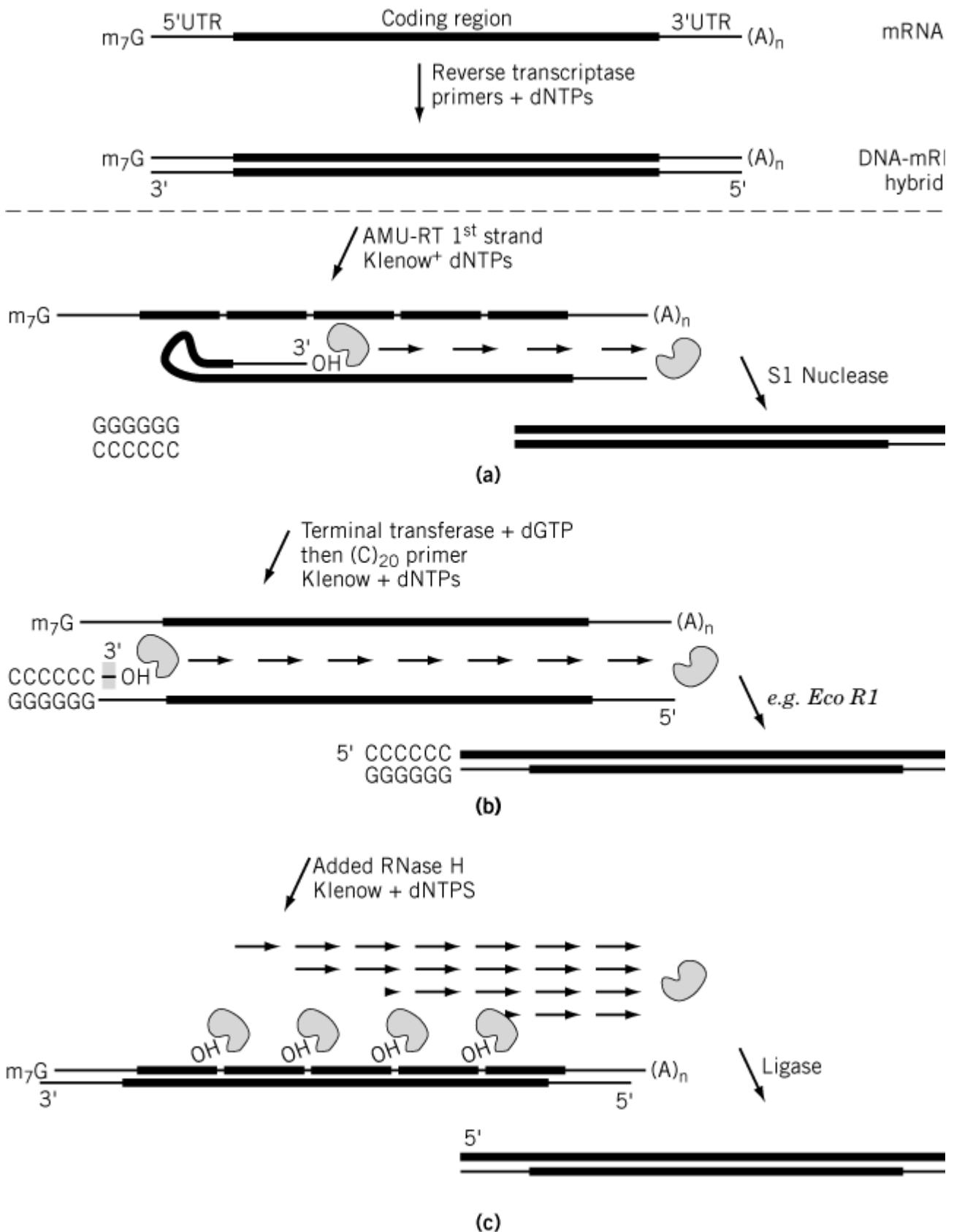
Among poly A + RNA, there are messages of low, middle, and high abundance. The process of normalizing cDNA is used to reduce the relative number of high- and middle-abundance mRNA

represented in the cDNA product. In this way, relatively fewer cDNA clones need to be screened in order to isolate a low-abundance message copy from a cDNA library. Preferred removal of cDNAs representing high- or middle-copy abundance mRNA is based on the rate of strand reannealing in a denatured cDNA sample (19). The two strands of a cDNA representing a high-abundance message have a higher probability of finding each other in a mixture and anneal before those of low-abundance cDNA. When annealed or double-stranded, the cDNA is separated from unannealed or single-strand cDNA by the use of [hydroxyapatite chromatography](#). After a short annealing time, the low-abundance cDNA will remain in the single-strand chromatography fractions. These fractions are collected and allowed to anneal fully for use as normalized cDNA. Alternatively, [subtractive hybridization](#) is a technique used for the enrichment of, eg, tissue-specific or developmental stage-specific mRNA in the preparation of subtraction cDNA libraries.

## 2.2. Synthesis of the cDNA First Strand

The two cDNA strands are prepared independently. In the first-strand synthesis, the use of an RNA-dependent DNA polymerase results in a mixture of cDNA–RNA hybrid molecules in which the DNA strand is complementary to the original RNA, but possibly of shorter length (Fig. 2). The reverse transcription systems in common use are based on the avian myeloblastosis reverse transcriptase (AMV-RT) prepared from purified avian myeloblastosis virus particles or on the Moloney murine leukemia virus reverse transcriptase (MMLV-RT) prepared from a recombinant *Escherichia coli* strain (20). The purified AMV-RT has a stronger endogenous RNase H activity than the cloned MMLV-RT, thus increasing the possibility of unwanted side reactions when using AMV-RT. However, purified AMV-RT may be more stable than MMLV-RT at the elevated temperatures used to reduce premature termination caused by **secondary structure** in the RNA template. RNase H is an endoribonuclease that specifically hydrolyzes the phosphodiester bonds of RNA in a DNA–RNA hybrid to produce products with terminal groups 3'-OH and 5'-P. Both AMV-RT and MMLV-RT have a tendency to begin second-strand synthesis prematurely through a reaction in which reverse transcriptase uses DNA rather than RNA as a template by forming a hairpin structure that serves as a template-primer substrate (21). Hairpin formation may be inhibited by the addition of sodium pyrophosphate and **spermidine**, although the RNase H activity is not inhibited with these conditions. The RNase H activity present in AMV-RT may provide a better opportunity for favorable hairpin formation by leaving single-stranded regions of DNA available for self-hybridization; this side reaction is more of a problem when AMV-RT is used. Heat-stable DNA polymerase isolated from *Thermus thermophilus* is capable of reverse transcription in the presence of MnCl<sub>2</sub>, which is a common alternative to the use of reverse transcriptases (22). Each system of reverse transcription has some advantages, as will be discussed, and should be considered with a view of how the second-strand synthesis will be performed. In addition to reverse transcriptase, first-strand synthesis requires deoxyribonucleotide triphosphates, dATP, dGTP, dTTP, and dCTP (or 5-methyl dCTP which may be substituted for dCTP in order to protect the cDNA from 5-methyl dCTP-sensitive [restriction enzymes](#) that may be used in subsequent cloning steps). The first-strand synthesis reaction also requires primers.

**Figure 2.** Three separate schemes of second-strand cDNA synthesis from the DNA–mRNA hybrid. (a) A classic method makes advantage of a looped back first strand as a primer for second-strand synthesis. After second-strand synthesis, S1 nuclease digestion of the single-strand loop structure results in an available free end for vector incorporation. (b) Terminal transferase-directed extension of the 3' end of the first strand results in a dG stretch that provides a binding site for an oligo dC primer to initiate second-strand synthesis. The undesired dG stretch on the first-strand oligo dT primer is removed by restriction endonuclease digestion targeted to a site present in the oligo dT primer or in the recipient vector if vector priming is used. (c) Added RNase H degrades the RNA in the DNA–RNA hybrid product of first-strand synthesis, resulting in a series of RNA primers that are extended during second-strand synthesis. Subsequent ligation may aid the recovery of full-length cDNAs.



Unlike **RNA polymerases**, which bind a recognition sequence in a double-strand DNA template to initiate polymerization, replicative-type DNA polymerases require primers with a free 3'-OH group and a short double-stranded sequence. Note that a consideration of the primers to be used for the

first-strand synthesis may serve to direct a more representative cDNA or cDNA better suited to the research aim. Because all mRNA, except animal [histone](#) mRNA, terminate with a 3' poly A tract, the use of a p(dT)<sub>10–15</sub> oligonucleotide as a first-strand primer ensures that the synthesis begins at the 3' end of mRNA, even when total RNA is used as a template for reverse transcription. A disadvantage of this primer is that the 5' end of longer mRNA may be underrepresented in the cDNA product, since the chance that processive synthesis will continue to the 5' end decreases with product length. The mRNA 3' terminal poly A tracts with undesirable average lengths of ≈200 nucleotides can be avoided by the use of an “anchored p(dT) oligo,” which has an equal concentration of dA, dC, or dG at the 3' end of p(dT)<sub>12</sub> dN. Such an oligo mixture initiates first-strand synthesis at the position immediately adjacent to the beginning (5') of the poly A tract, with an expected reduction in the length of the poly A tract captured in the cDNA. A random primer of six oligodeoxynucleotides used in place of p(dT)<sub>10–15</sub> may result in a more representative first-strand synthesis, since at least some of the primers should find matching hybridization sequences near the middle and 5' end of the RNA and thus are frequently extended to the 5' end of the template. A random primer is a better choice for use with poly A(–)RNA, which is not complementary to p(dT)<sub>10–15</sub>. A technical consideration for random primed first-strand synthesis is the amount or ratio of primer to template that should be used, since more than one primer for each template could lead to interference between polymerase initiating on different primers hybridized to the same RNA template molecule. This will result in premature termination and a smaller average cDNA product length in the sample. Mixtures of random-primed and oligo dT-primed first-strand synthesis products should equally represent the length of RNA.

Another choice of primer to initiate first-strand synthesis is a sequence-specific primer, a primer complementary to a known sequence within the mRNA of interest. The design of a sequence-specific primer depends, of course, on some available sequence information. A very conserved sequence common to a group of possibly related mRNAs used for the preparation of a limited cDNA library is an example of the application of a sequence-specific primer. Whichever type of primer is used for the initiation of first-strand synthesis, extra sequence at the 5' end of a primer, such as a restriction site sequence or a second primer-binding site, should not interfere with polymerase-directed extension of the 3' primer end. Thus, the 5' end of the primer may serve as a “sticky end” to incorporate the cDNA into a vector or as a primer target for subsequent amplification of the cDNA by PCR.

### 2.3. Synthesis of the Second Strand

As for first-strand synthesis, second-strand synthesis requires a polymerase, dNTPs, and a primer added to a template. In this case, the template is the product of the first-strand synthesis reaction. Since the template is DNA, a DNA-dependent DNA polymerase, most commonly *E. coli* DNA polymerase I or the **Klenow fragment** of DNA polymerase I, is used for second-strand synthesis. Additionally, reverse transcriptase added as a “chaser” reportedly facilitates second-strand synthesis through obstructions caused by secondary structure in the template ([23](#)). The product of the second-strand synthesis reaction is double-strand cDNA ready for steps leading to recombinant vector integration. The choice of primers that might serve to direct the synthesis of the second strand is limited, since the 3' terminal sequence of the first strand varies greatly among the first-strand synthesis product population; there is no common sequence at the 3' end of the first-strand product available for primer hybridization. Primer design for second-strand synthesis has conventionally been approached in one of three ways: (1) the classic method of hairpin-loop priming followed by [S1 nuclease](#) digestion, (2) RNase H treatment to generate RNA primers, or (3) homopolymeric tailing, also known as vector priming (see Fig. [2](#)). RNase H is commonly available in a pure form for the second-strand synthesis reaction (Fig. [2c](#)). Added RNase H is used to generate short RNA primers that are extended to form the second strand ([24](#), [25](#)). The relative technical simplicity, high yield, and large average product length achieved with the RNase H method have proven it to be the superior approach generally, although other methods may continue to have specific applications and certainly

have had a central role in characterization of the potential side reactions that affect the yield and average size of the cDNA product.

The classic method of second-strand synthesis (26, 27) takes advantage of the tendency for the 3' end of the first strand to fold back and form a DNA duplex with sequences located further 5' or upstream on the same DNA strand (Fig. 2a). A consequence of the duplex or hairpin formation is availability of the self-primed template for second-strand synthesis. After removal of the RNA, reagents for second-strand synthesis are added, and so the second strand is generated. The product is double-strand cDNA covalently closed at the end corresponding to the 5' end of the original mRNA. The addition of  $S_1$  nuclease, a single-strand-specific nuclease, results in cleavage of the single-strand loop to produce free 5' and 3' ends, which are necessary for subsequent integration into a cloning vector. Although the use of  $S_1$  nuclease, which has a potential of hydrolyze double-stranded DNA when used in excess, is often cited as the major disadvantage of the classical method, it is the necessary loss of sequence information that is the true technical fault in the method. Since at least some sequence is involved in the hairpin formation, the sequence of the 5' end of the original mRNA is not represented in the cDNA product when this classical method of second-strand synthesis is used.

The method of homopolymeric tailing or vector priming relies on the use of *terminal deoxynucleotidyltransferase*, an enzyme that catalyzes the sequential addition of available nucleotide triphosphates to the free 3' end of a DNA strand. Thus, a homopolymer tail made by extension of the free 3' end of the cDNA strand resulting from the first-strand synthesis reaction serves as the primer site for second-strand synthesis (Fig. 2b). The application of homopolymer tailing to second-strand synthesis preserves the cDNA sequence representing the original 5' end of the mRNA template, since priming is from sequence added by terminal transferase (28, 25). After first-strand synthesis initiated from a T-tailed vector, the 3' end of the newly formed cDNA strand is the target for homopolymer tailing. The use of dGTP for the tailing reaction consistently results in an average of 20 G nucleotides on both the 3' end of the cDNA first strand and on the free 3' end of the vector. The consistent length of the tail with use of dGTP is apparently due to structural constraints, which cause decreased tailing after 20 Gs (29, 30). A restriction site in the vector is used to remove the G-tail from the vector, without affecting the G-tail on the cDNA, which as a DNA–RNA hybrid is a poor substrate for restriction endonucleolytic cleavage (31). Finally, an oligo p(dC)<sub>20</sub> primer is added, along with other reagents required for second-strand synthesis. The p(dC)<sub>20</sub> primer may be designed to have a 5' end complementary to the overhang of the vector restricted end to aid in subsequent circularization of the recombinant DNA; otherwise, both the cDNA and vector ends may be made flush or blunt prior to circularization. Many variations of the vector priming method briefly described here have been described in detail elsewhere, but all make use of terminal transferase (32). The disadvantages of being more technically difficult and having more steps, as well as having a stretch of 20 G nucleotides at the vector-cDNA boundary (which could require special techniques for sequencing and subcloning), might be balanced by the important potential of increasing the number of full-length clones when using this method of second-strand synthesis.

Isolation of a cDNA corresponding to the intact full-length mRNA molecule often presents a challenge. Methods that increase the frequency of obtaining full-length clones in a cDNA synthesis reaction are often preferred. Losses of cDNA may occur at each step of the synthesis, from RNA isolation to vector incorporation. Losses of full-length cDNA may also occur due to incomplete replication of template by polymerase. For example, the presence of secondary structure in the RNA template may impede the progress of reverse transcriptase during first-strand synthesis, resulting in the formation of a partial cDNA beginning at the primer site and ending at the site where the secondary structure is encountered by the polymerase. Technical manuals that contain detailed protocols for cDNA synthesis often describe ways to monitor the progress of cDNA synthesis to give an indication of the reaction yields, or the incurred losses, and also the average size of the cDNA

strand. A method designed to extend a partial cDNA isolated by a selection procedure called RACE (rapid amplification of cDNA ends) may be used to obtain the corresponding full-length clone (33). The actual 5' end of an mRNA should, however, be mapped against the gene and/or correlated with the N-terminal sequence of the protein where possible. [Primer Extension](#) is frequently used to determine the precise 5' end of an mRNA, and this method only requires enough sequence information for primer design.

### 3. Uses of cDNA

The final steps in cDNA synthesis are often those leading to vector integration. The use or application of cDNA is largely determined by the type of vector that is chosen to contain the cDNA. Vector integration is, therefore, the first step toward fitting the cDNA to the type of application that will be made or toward developing a selection strategy for the isolation of a cDNA of interest from a library. Recent progress in the development of cloning vectors for specialized applications and of methods of efficiently adapting the cDNA ends to fit the vector ends have resulted in a lot of flexibility in the combination of available vector features when making a choice from the long list of commercially available cloning vectors. The immediate aim of retaining all the newly synthesized cDNA through the process of vector integration and host transformation, while keeping the desired vector features, is frequently satisfied by the use of a **phagemid** (34). A phagemid vector combines features of [lambda phage](#), filamentous phage, and plasmid vectors into one cloning vector (35, 36). The efficiency of host cell transformation through packaging of cDNA into viral particles, and subsequent viral infection as performed with the use of a l phage cloning vector, cannot be matched by even the optimal conditions for direct host-cell transformation with plasmid DNA. There is consequently less potential for cDNA loss during host cell transformation with a l phage system when compared to a plasmid system. The lower background that results from the screening of plaques rather than colonies with sequence hybridization probes or [antibody](#) probes offers another significant advantage to the use of viral vectors. Plasmid vectors offer advantages over l phage vectors for the characterization of the cDNA inserts, due to the necessarily large size of l phage vectors. Phagemid vectors contain a plasmid vector within a l phage vector and additional sequences derived from filamentous phage that contribute to the *in vivo* excision of the plasmid from the l phage vector. Thus, for characterization of individual selected cDNA clones, the original phage clones are converted to plasmid clones. Although it is technically straightforward to subclone a single cDNA clone from one vector to another, it is technically difficult to subclone an entire cDNA library from one to another vector without some loss of cDNA. Whichever type of vector is used to receive newly made cDNA, consideration of the vector features required for selection of a cDNA of interest from a library and the efficiency of the process of integration and host-cell transformation should be given prior to vector integration of cDNA as a library.

Possibly the greatest benefit of acquiring a new cDNA of interest from a cDNA library is the sequence information that can be derived from it. In addition to the encoded protein sequence, comparison of a cDNA sequence with the corresponding gene sequence reveals or confirms the location of any introns present in the gene. Analysis of a set of positives identified by rescreening of a library with the now available hybridization probe may provide some indication of the variability of the gene transcripts, eg, the presence of alternate splice sites or multicopy genes with minor sequence variations. The presence of viral RNA polymerase recognition sequences, most commonly SP6 and T7 **promoters**, in many cloning vectors are useful for the *in vitro* production of anti-sense RNA probes that may be used for the high-sensitivity detection of gene expression. Likewise, PCR primers may be designed from the sequence information to provide an alternate method of sensitive or quantitative gene expression analysis. The presence of potential post-transcriptional regulatory elements in the untranslated regions represented in the cDNA may be identified through computer-assisted [sequence analysis](#) and [database](#) comparisons. Although any number of extremely sophisticated uses for cDNA sequence information may be envisioned for a particular application, a generally useful benefit is the possibility to generate peptide antibodies to encoded peptide [epitopes](#) predicted to be on the protein surface.

Vectors are available with many combinations of features. A particularly useful vector feature for many types of applications is the inclusion of sequences that support cDNA expression. As the sequences required for correct initiation of translation differ among organisms, expression vectors, or an expression cassette within a cloning vector, are designed for use as a set with a matching host cell. The promoter sequences included to drive the transcription of the expression cassette and efficient transcriptional termination signals are likewise active in the same host. As there are no introns to be removed, expression of many encoded eukaryotic proteins is supported in strains of *E. coli*, provided that these proteins are not toxic to the *E. coli* host strain used. In order to avoid the possible toxic effects of some eukaryotic proteins, expression vectors with inducible promoters are available to postpone the cDNA expression to a point of vigorous growth so that significant protein production is accomplished before the host culture dies. The activity of eukaryotic proteins produced in prokaryotic host cells may be deficient due to the absence of factors required for proper [protein folding in vivo](#), for [post-translational modifications](#) essential for the activity, or even for correct intracellular localization. Eukaryotic expression systems may better support the production of eukaryotic proteins with natural levels of activity and are designed for maximum protein production. Such expression systems may support continued replication of the vector containing a cDNA insert, in addition to containing an **expression cassette** with an active promoter and signals for transcription termination and translation initiation. Both prokaryotic and eukaryotic expression systems with matched vector and host are designed to produce proteins in quantities required for their structural characterization from manageable volumes of cell culture. Accurate physical measurements that require purified proteins may take advantage of “built-in” purification aids available in some expression vectors. Expression of a cDNA as a **fusion protein** aids in the isolation of protein from cell cultures. Provided a cDNA is cloned in-frame in relation to a vector sequence that encodes a [polypeptide chain](#) with a known ligand affinity, the resulting fusion protein may be purified from cell lysates by passage through an [affinity chromatography](#) column containing the immobilized ligand. The presence of a [proteinase](#) peptidase recognition sequence facilitates proteolytic removal of the vector-encoded polypeptide and elution of the cDNA-encoded protein from the affinity column.

An additional type of expression vector that is adapted to cDNA selection is one designed for use with a yeast host strain. The yeast 1, 2, and 3-hybrid systems have been designed to isolate a cDNA clone from a library based on the encoded protein's recognition or specific binding to a DNA, protein, or RNA sequence, respectively (37-39) (see [Two-Hybrid Systems](#)). In this way, cDNA selection is facilitated by its expression inside a yeast host cell that has been genetically altered for the selection. In order to perform the selection, an entire cDNA library is cloned into a shuttle vector that will support replication in *E. coli*, as well as replication and expression (as a fusion protein) in the yeast host. The library is amplified for transformation into the recipient yeast host strain, which is altered in such a way that specific binding of a target sequence or “bait sequence” results in a positive selection. The yeast hybrid systems typify advances in the use of traditional genetic strains adapted for functional cloning of cDNA from higher organisms, and yeast has consequently become a tool for functional cloning through molecular recognition, extending the use of cDNA as an aid to the selection of novel factors. Thus, a factor that binds to a DNA sequence, such as a [transcription factor](#) that binds to a gene promoter element, may be difficult to isolate by biochemical means due to its low intracellular abundance, but it might be cloned with the yeast 1-hybrid system. Likewise, a factor that binds to a protein might be cloned with the 2-hybrid system, or an RNA-binding protein with the yeast 3-hybrid system.

#### 4. Summary

cDNA is a faithful double-stranded DNA copy of RNA and, as such, can represent those genes that are expressed as RNA in a source tissue. When integrated into an appropriate cloning vector to form recombinant DNA, cDNA may provide a continuous renewable source of genetic material for hybridization probes, designed reporters of biological activity, or pure encoded protein. The benefits of protein production and a myriad of specific research applications have made cDNA a universal tool for basic and applied cell and molecular biology research. The potential for using cDNA in



molecular methods designed to select for activities that contribute detectable cellular phenotypes has provided possible alternatives or improvements to genetic and biochemical approaches.

### Bibliography

1. D. Gillespie and S. Spiegelman (1965) *J. Mol. Biol.* **12**, 829–842.
2. H. M. Temin and S. Mizutani (1970) *Nature* **226**, 1211–1213.
3. D. Baltimore (1970) *Nature* **226**, 1209–1211.
4. M. G. Sarngadharan et al. (1972) *Nature New Biol.* **240**, 67–72.
5. E. M. Scolnick (1971) *Develop. Biol.* **26**, 175–176.
6. J. C. Barrett et al. (1979) *Cancer Res.* **39**, 1504–1510.
7. M. Noda (1990) *Molec. Carcinogenesis* **3**, 251–253.
8. C. P. Carstens et al. (1995) *Gene* **164**, 195–202.
9. L. Rowen et al. (1997) *Science* **278**, 605–607.
10. M. Boguski and G. D. Schuler (1995) *Nature Genet.* **10**, 36971.
11. M. D. Adams et al. (1995) *Nature* **377**, Suppl., 3–174.
12. S. Henikoff et al. (1997) *Science* **278**, 609–614.
13. R. L. Tatusov et al. (1997) *Science* **278**, 631–637.
14. B. E. Roberts and B. M. Paterson (1973) *Proc. Natl. Acad. Sci. USA* **70**, 2330.
15. C. W. Anderson et al. (1983) *Methods Enzymol.* **101**, 635.
16. H. Aviv and P. Leder (1972) *Proc. Natl. Acad. Sci. USA* **69**, 1408–1412.
17. I. Edery et al. (1995) *Mol. Cell Biol.* **15**, 3363–3371.
18. J. Ross (1997) *Bioessays* **19**, 527–529.
19. T. G. Coche (1997) *Methods Mol. Biol.* **67**, 359–369.
20. M. Roth et al. (1985) *J. Biol. Chem.* **260**, 9326–9335.
21. M. S. Krug and S. L. Berger (1987) *Methods Enzymol.* **152**, 318.
22. C. Rüttimann et al. (1985) *Eur. J. Biochem.* **149**, 41–46.
23. U. Gubler (1987) *Methods Enzymol.* **152**, 325.
24. H. Okayama and P. Berg (1982) *Mol. Cell Biol.* **2**, 161–170.
25. U. Gubler and B. J. Hoffman (1983) *Gene* **25**, 263–269.
26. A. Efstratiadis et al. (1976) *Cell* **7**, 279–288.
27. F. Rougeon and B. Mach (1976) *Proc. Natl. Acad. Sci. USA* **73**, 3418–3422.
28. H. Land et al. (1981) *Nucl. Acids Res.* **9**, 2251–2266.
29. A. Dugaiczky et al. (1980) *Biochemistry* **19**, 5869–5873.
30. A. Otsuka (1981) *Gene* **13**, 339–346.
31. W. H. Eschenfeldt et al. (1987) *Methods Enzymol.* **152**, 339.
32. P. L. Deininger (1987) *Methods Enzymol.* **152**, 376.
33. M. A. Frohman et al. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 8998–9002.
34. M. A. Alting-Mees et al. (1992) *Methods Enzymol.* **216**, 483.
35. D. Hanahan (1983) *J. Mol. Biol.* **166**, 557–580.
36. J. M. Short et al. (1988) *Nucl. Acids Res.* **16**, 7583–7600.
37. S. Fields and O. Song (1989) *Nature* **340**, 245–247.
38. D. J. SenGupta et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8496–8501.
39. Z. F. Wang et al. (1996) *Genes Develop.* **10**, 3028–3040.

### Suggestions for Further Reading

40. V. B. Chanda, ed. (1997) *Current Protocols in Molecular Biology*, John Wiley & Sons, New York.
41. I. G. Cowell and C. A. Austin, eds. (1997) *Methods in Molecular Biology—cDNA Library Protocols*, Humana Press, Totowa, NJ.
42. Suppliers of the reagents used for cDNA synthesis often supply excellent technical manuals.

## Complementation

The term complementation usually refers to the situation where two defective genomes in the same cell together support a normal or nearly normal **phenotype**, each supplying the function that the other lacks. Another meaning of the term, in the context of DNA manipulation, is the repair of a mutational deficiency in an organism or cell culture by artificially introducing a **gene** that supplies the missing function. This is a potent way of [cloning](#) genes selected according to function. In either sense, complementation provides a way of discriminating between gene functions and of defining genes as functional units.

### 1. Complementation and the definition of the functional gene

Classical genetics was based upon clear-cut heritable variants, either induced by mutagenic treatments or, in earlier days, just turning up as “sports” in wild or cultivated populations. The geneticists' favorite organisms, notably **maize** (*Zea mays*) and the fruit-fly (*Drosophila melanogaster*), were **diploid**. Most mutations were recessive to **wild type**, which is what would be expected if they mostly represented losses of function that could be more or less adequately supplied if one gene out of two was normally active.

As mutants accumulated, they were attributed to the same or different genes on the basis of two criteria. First, it was supposed that mutations in the same gene were inseparable by [recombination](#). A diploid that carries two different mutant **alleles** (ie, alternative forms) of a particular gene should never produce nonmutant or double-mutant germ cells by free reassortment or **crossing-over** at **meiosis**. According to the second criterion, mutations in the same gene should affect the organism in similar or at least related ways. Similarity or overlap of functional effects was judged by individual appearance (phenotype) and also and more rigorously by failure of the mutant chromosomes to complement each other's defects in the hybrid diploid. (In fact, the term complementation did not come into use until the late 1950s, and I use the word with hindsight.)

A classical example of using the complementation criterion, and one of the first in which it was thrown into doubt, involves the sex-linked *white* (*w*) gene of *Drosophila*, one of several genes governing eye color. A considerable number of recessive mutant *w* alleles had been identified. They diluted the normal red-brown eye pigment to various degrees when the same allele was present on both [X-chromosomes](#) (**homozygous**) in females, or on the single X in males. Females that had two different mutant *w* alleles (**heterozygous**) had dilute eye colors generally intermediate between the colors of the two homozygotes. Thus females of constitution *white/apricot* ( $w/w^a$ ) have pale apricot eyes, not wild-type red as would be expected if each mutant X-chromosome could provide the pigment-forming function lacking in the other.

In 1952, E. B. Lewis (1) found that, at very low frequency (about 1 in 10,000),  $w/w^a$  females produce eggs of recombinant types, with neither mutation or both. By the attached-X technique, which need

not be described here, he showed that these recombinants arise by crossing-over between the X-chromosomes, just as if *white* and *apricot* were mutations at distinct, though very closely linked loci. Accordingly he gave them different symbols, *w* and *a*, and rewrote the constitution of the female parent  $w^+/w^+a$ , where the + superscripts denote the respective wild-type alleles.

The remaining anomaly was what Lewis called the *cis/trans* position effect (Fig. 1). Two kinds of doubly heterozygous constitution could be compared, the original so-called *trans* arrangement  $w^+/w^+a$ , that has the two mutations on opposite chromosomes, and the *cis* arrangement,  $w^+a^+/wa$ . They had identical overall gene content, but the former had pale apricot eyes and the latter was fully wild type.  $w^+$  and  $a^+$  could act cooperatively to promote normal levels of eye pigment only when they were together on the same chromosome. Lewis called their relationship **pseudoallelic**: allelic according to their failure to complement each other but nonallelic according to their ability to recombine. A formally identical situation at the *Drosophila lozenge* locus had already been reported by M. M. Green (2).

**Figure 1.** Lewis's (1) *cis-trans* comparison using two mutations in the *Drosophila melanogaster* X-linked *white* (*w*) gene. The mutations *white* (*w*) and *apricot* ( $w^a$ , shown here as *a*) are each recessive to wild type in females but do not complement one another in *trans*. By this criterion they are allelic in spite of being separable at low frequency by recombination.

|              |   |          |          |   |          |                  |
|--------------|---|----------|----------|---|----------|------------------|
| <i>trans</i> | <table style="border-collapse: collapse; margin: 0 auto;"> <tr> <td style="padding: 0 10px;"><i>w</i></td> <td style="padding: 0 10px;">+</td> </tr> <tr> <td style="padding: 0 10px;">+</td> <td style="padding: 0 10px;"><i>a</i></td> </tr> </table> | <i>w</i> | +        | + | <i>a</i> | Wild type<br>red |
| <i>w</i>     | +   |          |          |   |          |                  |
| +            | <i>a</i>  |          |          |   |          |                  |
| <i>cis</i>   | <table style="border-collapse: collapse; margin: 0 auto;"> <tr> <td style="padding: 0 10px;"><i>w</i></td> <td style="padding: 0 10px;"><i>a</i></td> </tr> <tr> <td style="padding: 0 10px;">+</td> <td style="padding: 0 10px;">+</td> </tr> </table> | <i>w</i> | <i>a</i> | + | +        | Pale<br>apricot  |
| <i>w</i>     | <i>a</i>  |          |          |   |          |                  |
| +            | +   |          |          |   |          |                  |

For Lewis and Green at that time, the recombination criterion for allelism took priority, and the reason for the position effect was a matter for speculation. Before long, however, high-resolution recombination analysis of series of mutants in the microbial world, first in the fungus *Aspergillus nidulans* (3) and then in the bacterial virus (**bacteriophage**) T4 (4), revealed that pseudoallelism, in Lewis's sense, was much more the rule than the exception. It was rather rare for independently occurring mutations to be “truly” allelic. Almost all pairwise combinations of noncomplementing mutants showed some low frequency of recombination when crossed together. Soon afterward, the same was shown to be true in two bacterial species, *Escherichia coli* and *Salmonella typhimurium* (5), which have no regular diploidy or meiosis but other means exist for obtaining partial diploids and tests for complementation and recombination (see [Complementation Tests](#)). In one example after another, it was shown that most mutations within functional genes are at different mutually recombinable sites which, by different methods in different organisms, could be mapped in a closely spaced linear sequence that was continuous with the much longer sequence of the whole chromosome.

Following his very extensive analysis of the *rII* series of mutants in phage T4 ((4), see [Complementation Tests](#)), S. Benzer proposed replacing the word gene with a new term, **cistron**, meaning the unit of function as defined by Lewis's *cis/trans* comparison. When the *cis* and *trans* phenotypes are identical (wild type, if the mutants involved were recessive), the mutations are in different cistrons, affecting different functions. If *trans* is mutant and *cis* normal (or relatively so) the mutants are in the same cistron. The concept has been universally accepted, but the term cistron has dropped out of use. Now “gene” is generally used in the same sense, and the full *cis-trans* test is

rarely applied. The theoretical reason for the comparing *trans* ( $m1 + / + m2$ ) with *cis* ( $m1m2/ ++$ ), rather than just *trans* with the wild type, is that it controls the possibility that the cumulative effect of two heterozygous mutations in different genes could be a sub-wild phenotype, even though each is individually recessive to its wild-type allele. In practice, however, recessive alleles usually remain recessive even when several of them are present together, and so the *trans*/wild comparison, the simple complementation test, is usually deemed sufficient.

Although diploidy provides the means for complementation testing in higher plants and animals and also in budding **yeast**, alternative methods, described under [Complementation tests](#), have to be used for habitually haploid organisms, such as bacteria, viruses, and most fungi.

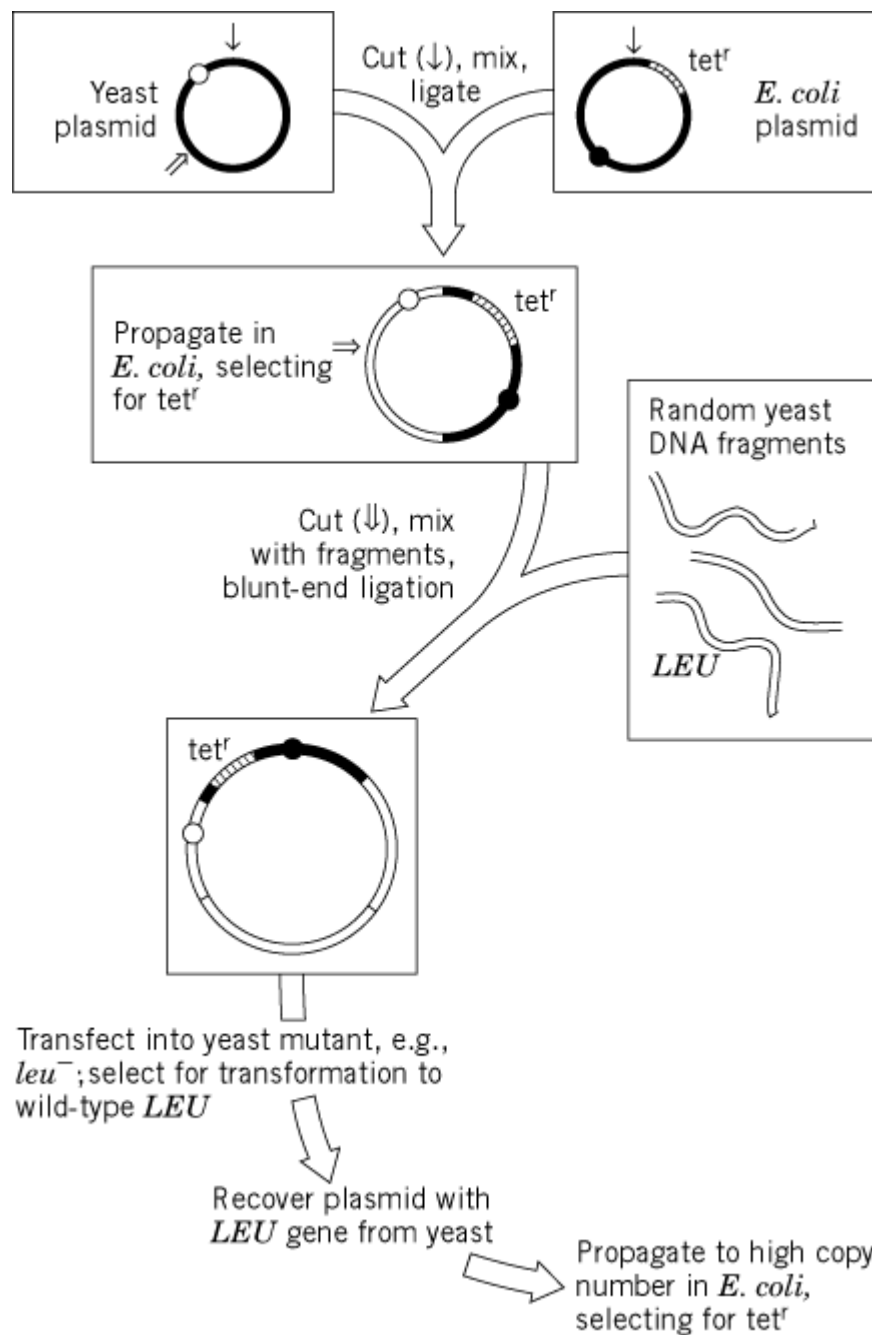
### 1.1. Complementation and gene cloning *in yeasts*

Cloning wild-type alleles of yeast genes with metabolic functions and mutating to give **auxotrophic** phenotypes, became straightforward following the development of the *Saccharomyces* two-micron (2- $\mu$ m) **plasmid** as a cloning vehicle (7). This plasmid also replicates in the other major experimental yeast species, the fission yeast *Schizosaccharomyces pombe*. It has been most useful as a component of a series of hybrid **shuttle vectors**, capable of replication either in yeast or in *Escherichia coli*. These usually incorporate the replicative origin of the *E. coli* plasmid ColE1 and a gene for [antibiotic resistance](#) (e.g., to **tetracycline**) that can be selected for in the bacterial host.

To select and clone a yeast gene, fragments of total yeast DNA, produced by digestion with a **restriction endonuclease** (or sometimes by sonication), are ligated into the closed-loop, double-stranded DNA of a shuttle vector and introduced en masse by one of the standard [transformation](#) procedures into mutant yeast cells, which then are plated on a medium on which they grow to form colonies only if their mutational deficiency has been repaired. (see [Cloning](#)) The first yeast genes to be cloned were selected because they conferred on auxotrophic (nutritionally exacting) mutants the ability to form colonies on a minimal (un-supplemented) growth medium. But any gene that mutates to give a conditional no-growth phenotype can be cloned in essentially the same way. Temperature sensitive mutants have been a particularly rich source of cloned genes. Almost any essential protein can be mutated to a temperature-sensitive form, functional at, say, 30°C but not at 36°C. Mutant cells grown at the permissive temperature are transformed with a gene “[library](#)” and selected for repair of function at the restrictive temperature.

The shuttle vector bearing the selected gene replicates autonomously either in yeast or after transfer to *E. coli*. The bacterial host is more efficient for mass propagation of the clone. When sufficiently amplified, the cloned gene is cut out of the vector, and its structure and function are analyzed in various ways, most obviously by DNA sequencing. The general method is outlined in Fig. 2.

**Figure 2.** Cloning and amplifying a yeast gene in a shuttle vector selected by its complementation of a yeast auxotrophic (leucine-requiring) mutant. Open and filled circles indicate replication origins for propagation in yeast and *E. coli*, respectively. Based on Beggs' first demonstration of the method (7).



The complementation of yeast mutants is used for cloning yeast genes and also genes from such distant organisms as *Drosophila* and humans. One good example is cloning the human equivalent of the *Schizosaccharomyces pombe* (fission yeast) *CDC2* gene (8). *S. pombe* temperature-sensitive *cdc2* mutants cannot complete their cell cycle at elevated temperature (36°C) because the mutant *cdc2* gene product, a protein [kinase](#) that provides a signal essential for initiating cell division, is inactivated at this temperature. Transformation of mutant cells with a human [cDNA library](#) cloned in a replicating plasmid led to the isolation of a few colonies that harbor hybrid plasmids carrying the desired cDNA sequence. The cloning plasmid used in this instance differs from that illustrated in Fig. 2 in that its replication in fission yeast depends on a DNA sequence from the monkey **SV40 virus**, which replicates well in *S. pombe*, not on the yeast 2- $\mu$  plasmid. To clone mammalian gene sequences in yeasts, it is generally necessary to use cDNA rather than genomic DNA. Mammalian **introns** are not effectively spliced out in yeast cells and, in any case, often make the gene too long for cloning in one piece.

### 1.2. In *Neurospora*

Although yeast species that have a single-cell, colony-forming habit and are accessible to replicating plasmids are especially convenient for cloning genes, the complementation principle can also be applied to filamentous fungi and has been particularly successful in *Neurospora crassa*. This species has the disadvantage that it has no convenient plasmid that stably replicates without integration into the chromosomes. Consequently, although auxotrophic mutants are “rescued” by transformation with fragments of total *Neurospora* DNA, it is not so easy to recover the transforming sequences from the transformed cultures. This difficulty was overcome by sib selection (9).

The method is making a library of *Neurospora* DNA sequences in some suitable vector (sufficient in number so that it is probable that the great majority of genes are represented), dividing the clones into a number of pools, and testing DNA isolated from each pool for its ability to repair the mutant of interest. Then a pool that works (and there is usually at least one) is subdivided into smaller pools for further tests, and then to still smaller pools, until the hunt is narrowed down to a small number of clones that can be tested individually. Though obviously laborious, the method has been very successful. It has been used for cloning genes of metabolic function and also for the loci (actually gene complexes) that govern fungal **mating type** by complementation of mating-deficient mutants (10, 11).

### 1.3. In *Drosophila*

Many gene functions involved with basic cellular processes are common to yeast and mammals. Then it is not surprising to find that others, more to do with tissue differentiation, are common to most animals. Because the genetic control of **development** is better understood in *Drosophila* than in any other animal, it is of interest to ask whether some of the relatively well-characterized *Drosophila* genes have mammalian equivalents. Flies cannot be used for the kind of mass screening of gene libraries that is possible with yeast, but it is not difficult to find out whether some particular mouse gene or cDNA clone complements a *Drosophila* mutant.

It is relatively easy to get exotic DNA into the *Drosophila* genome by inserting it into the transposable **P element** and injecting the construct into early embryos. Thus the mouse gene M33, cloned on the basis of its sequence similarity to the *Drosophila polycomb* (*Pc*) gene, was successfully introduced into *Pc* mutant flies (12). The mutants, which usually have sex combs on all legs of the male instead of just on the front pair as in the wild type, were partially normalized by inserting one copy of M33 into their genome and were restored almost to wild type by two copies. *Pc* protein is one of a group of chromosomal proteins that act together to prevent certain developmentally important genes from being expressed in the wrong places or at the wrong times. This experiment strongly indicates that there is a mammalian gene that has a similar function.

### Bibliography

1. E. B. Lewis (1952) Proc. Natl. Acad. Sci. USA **3**, 953–961.
2. M. M. Green and K. C. Green (1949) Proc. Natl. Acad. Sci. USA **35**, 586–591.
3. R. H. Pritchard (1953) Heredity **9**, 343–371.
4. S. Benzer (1955) Proc. Natl. Acad. Sci. USA **41**, 344–354.
5. P. E. Hartman, J. C. Loper, and D. Serman (1960) J. Gen. Microbiol. **22**, 323–368.
6. S. Benzer (1958) In *The Chemical Basis of Heredity* (W. D. McElroy and B. Glass, eds.), Johns Hopkins Press, Baltimore, pp. 70–93.
7. J. D. Beggs (1978) Nature **275**, 104–109.
8. M. G. Lee and P. M. Nurse (1987) Nature **327**, 31–33.
9. S. J. Vollmer and C. Yanofsky (1986) Proc. Natl. Acad. Sci. USA **83**, 4867–4873.
10. N. L. Glass, J. Grotelueschen, and R. L. Metz (1990) Proc. Natl. Acad. Sci. USA **87**, 4912–4916.
11. C. Staben and C. Yanofsky (1990) Proc. Natl. Acad. Sci. USA **87**, 4917–4921.

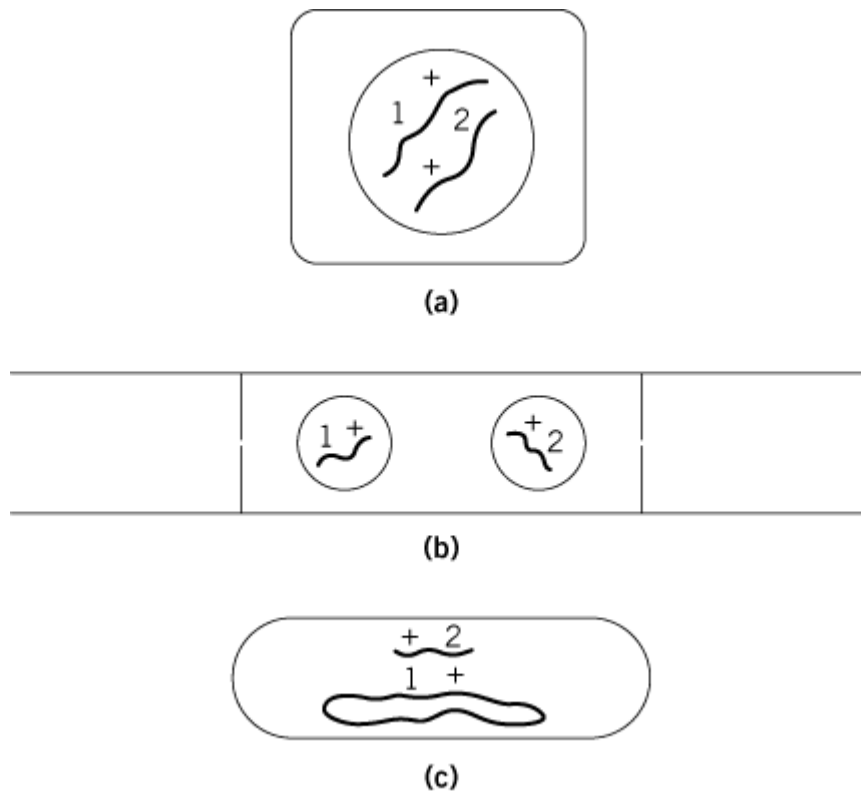
## Complementation Tests

The **gene** is best defined as an integral unit of genetic function, in the sense that none of its parts functions normally in isolation from the others. If so, two different forms (**alleles**) of the same gene, both more or less inactivated by **mutation**, should not make good each other's defects when together in the same cell. They should not, in other words, show **complementation**. On the other hand, two **genomes** that have mutations in different genes should complement each other to produce a wild-type phenotype, provided that each mutant allele is recessive to its **wild-type** counterpart. This complementation criterion is of great importance because it permits assigning mutations to the same or to different genes without a detailed analysis of gene functions. When mutations are closely linked and have fairly similar phenotypic effects, it is always likely that they are in the same gene. But it is not uncommon to find close linkages of separate genes whose function is related. A complementation test provides much stronger evidence.

The interpretation of complementation tests is itself subject to certain difficulties, which are discussed separately under the heading **Allelic complementation**. Here we make the simplifying assumption that complementation between recessive mutations means that they are in different genes.

For a complementation test, it is necessary to bring two mutant genomes, or at least relevant parts of genomes, together in the same cell. In organisms including animals, higher plants and budding **yeasts**, that have a stable **diploid** phase in their life cycles, this can be simply achieved by making a sexual cross between mutant strains. In microorganisms that have no stable diploid phase, various other means have to be employed. In **haploid** filamentous **fungi**, **heterokaryons** that have two genetically different kinds of nuclei in a common cytoplasm provide a satisfactory substitute for diploidy. In **bacteria**, where neither diploidy nor heterokaryosis is usually available, there are various ways of introducing fragments of the genome from one cell to another to make partial diploids (**merodiploids**). This makes it possible to test for complementation of genes located within the duplicated fragment (Fig. 1). In **viruses**, especially bacterial viruses (**bacteriophages**), complementation tests can be carried out simply by mixed infections.

**Figure 1.** Three general ways of making complementation tests: (a) by diploidy (higher plants, animals and the budding yeast *Saccharomyces cerevisiae*); (b) by heterokaryosis (fungi, such as *Neurospora crassa* and sometimes cultured mammalian cells); (c) by addition of a genomic fragment to a haploid genome (bacteria, such as *Escherichia coli* and *Salmonella typhimurium*). In each diagram, 1 and 2 are different mutations, and + stands for the corresponding nonmutated sites.



## 1. Complementation Tests in Diploids

The most extensive complementation testing using diploids has been carried out in the budding yeast *Saccharomyces cerevisiae* and the fruit fly *Drosophila melanogaster*. These are the diploid organisms that have most experimentally induced mutants sorted into genes.

*Saccharomyces* has the great advantage as an experimental organism of propagating itself equally well as a haploid or a diploid. Haploid cells are of two different mating types, called *a* and *a*. In wild strains, the cells are constantly interconverted by a genetic switching mechanism but are stabilized in nonswitching laboratory strains. Most yeast genetics is based on the use of **auxotrophic** mutants that grow only if given some specific nutritional supplement which they can no longer synthesize—most commonly an [amino acid](#), a **purine** or **pyrimidine** base, or a vitamin. Such mutants can be readily obtained in haploid strains. For comprehensive complementation testing, it is necessary to isolate each mutant in both mating types from a cross to wild type. Then haploid cells of pairs of mutants of opposite mating types are mixed at marked positions on plates of unsupplemented agar growth medium. Sexual fusion to form diploid cells occurs almost immediately, and the diploids grow if the two mutations complement each other's nutritional deficiencies. The result is quite clear after overnight incubation. Auxotrophs that have different nutritional requirements always complement, but mutants that have the same requirement may or may not do so because several different genes acting in sequence are usually necessary to synthesize a single end product. The results split the mutants that have a common requirement, say, for the amino acid [histidine](#), into clear-cut complementation groups - *his1*, *his2*, *his3*, etc. Every member of one group complements every members of any other group, but only sporadic and fairly infrequent complementation occurs within groups (see [Interallelic Complementation](#)). In general, complementation groups equate to genes, and each one is deficient in one enzyme of the biosynthetic pathway.

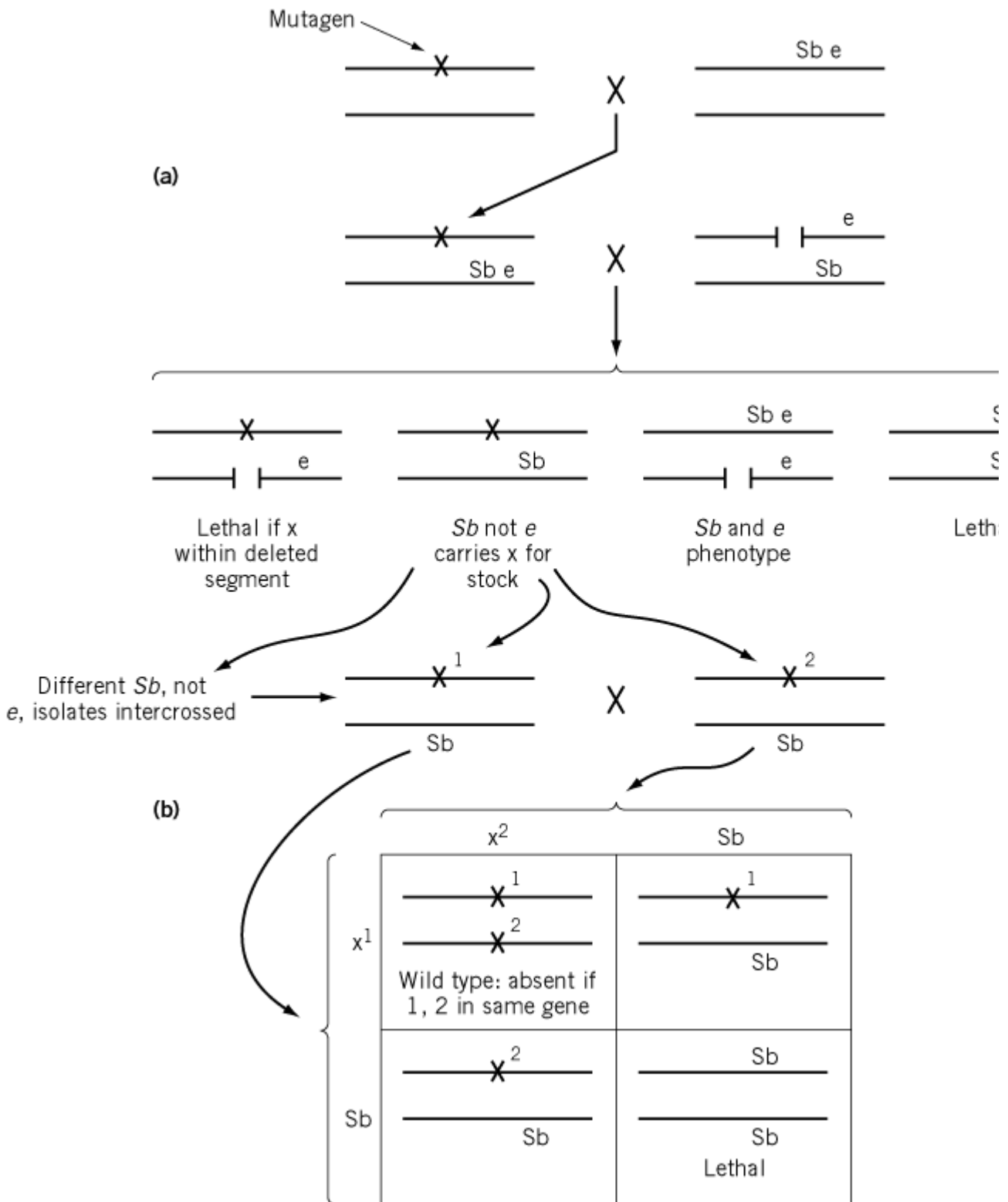
In the fruit-fly, *Drosophila melanogaster*, which is certainly the most intensively investigated multicellular organism, many mutants have been assigned to allelic series based on the mutant phenotypes of their heteroallelic combinations. The famous example of multiple alleles of the *w*



(white-eye) gene is referred to under [Complementation](#). Most of the allelic series have consisted of viable mutants, all within the same short [chromosome](#) segment, and obviously related phenotypes that made their allelic relationship are very likely even without complementation testing. However, nonvisible mutations, **recessive lethals**, have been much used in attempts to establish the total number of genes that have essential functions residing within a particular chromosome segment. Complementation is the only guide for classification of lethals.

A broadly applicable experimental scheme is shown in Fig. 2. The method is selecting, after heavy mutagenization, a large number of **recessive lethal** chromosomes based on their failure to complement a deletion of a short, defined chromosome segment, and maintaining them in breeding stocks combined with a “balancer” chromosome. A *Drosophila* balancer chromosome carries a dominant viable mutation to show that it is present and a recessive lethal mutation to ensure that any fly that carries it also has the chromosome against which the balancer is supposed to be balanced. The balancer chromosome also contains an inverted segment, or complex of inverted segments, to inhibit its recombination with the mutagenized chromosome, or else has genetic markers to ensure that recombination within the mutagenized region can be avoided. Crosses between flies that have different recessive lethals, each in combination with the same balancer, produce a predicted one-third of their progeny free of the balancer, if the two are in different genes and therefore complement one another, but no progeny without the balancer if they are in the same gene.

**Figure 2.** The isolation of recessive lethal mutations in *Drosophila* and testing them for complementation. **(a)** Male flies treated with a mutagen and mated to females carrying a balancer third chromosome that has the marker mutations *Sb* (do bristle phenotype, recessive lethal) and *e* (recessive ebony body colour). Flies in the next generation that have mutagenized chromosome 3 covered by the balancer, are crossed to flies that have a third-chromosome deletion and the *e* marker, balanced against *Sb*. If the mutagenized third chromosome has a recessive lethal mutation (x) at a site within the deletion, there are viable progeny that do not carry *Sb*. Then, the recessive lethal is present, balanced against the *Sb* chromosome in all viable progeny that are not ebony (*e/e*). **(b)** When *1/Sb* and *2/Sb* flies are crossed (*1* and *2* are two different lethals), one-third of viable progeny should be wild type (without *Sb*) if *1* and *2* are in different complementation groups (equivalent to genes) exhibit *Sb* if *1* and *2* are in the same complementation group. Based on Ref. [1](#).



In the study on which Fig. 2 is based (1), a total of 268 recessive lethals that fall within a segment of chromosome 3 comprising 26 *polytene* chromosomal bands could be assigned to 25 complementation groups. The numbers were sufficiently great to make it likely that most, if not all, of the genes that have essential functions and therefore can mutate to recessive lethality had been identified. The gene number was probably underestimated because there are certainly genes that can be deleted without lethality. Nevertheless, the agreement between chromosome band number and the apparent number of essential genes (now found in several different experiments) is too close to be

without significance.

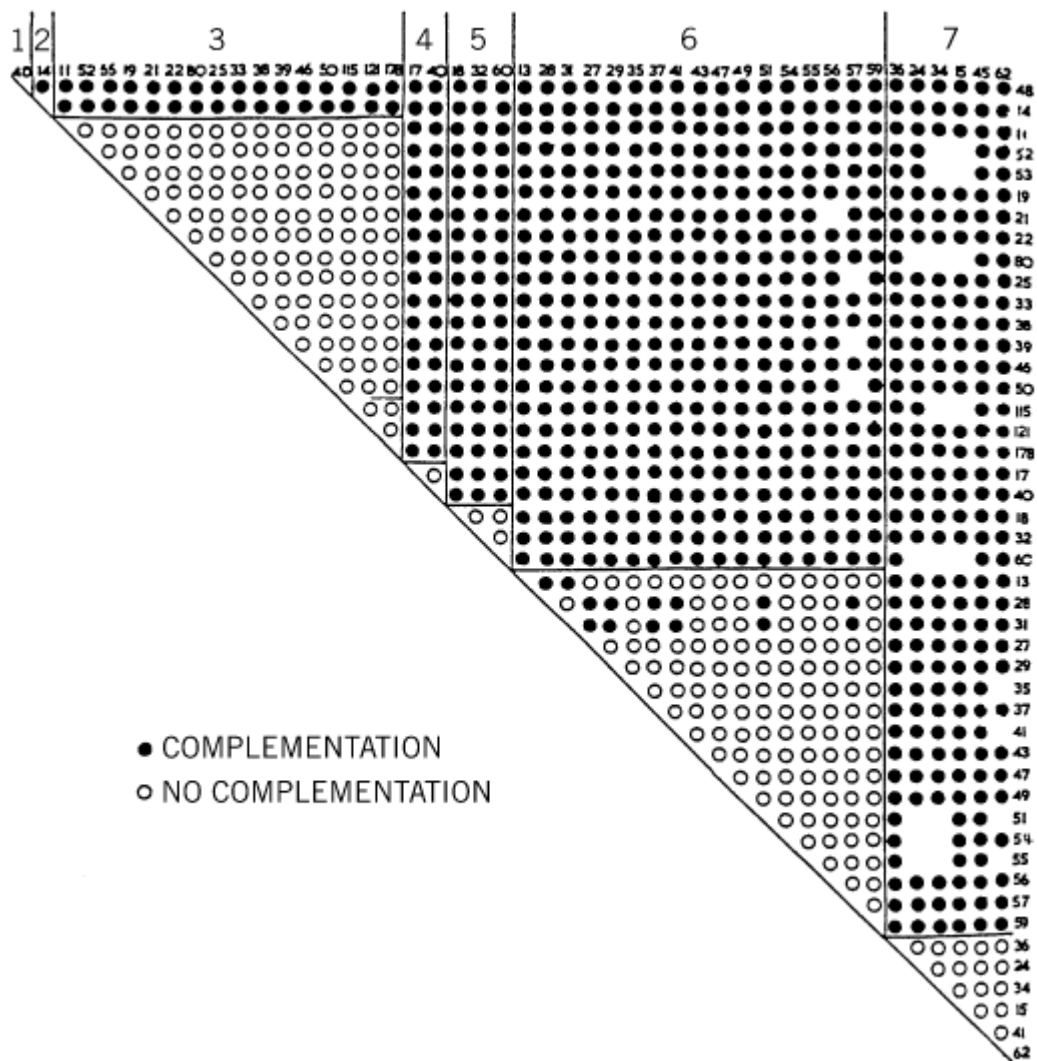
## 2. The Use of Heterokaryons

Filamentous fungi, or at least those (eg, *Neurospora crassa*, *Aspergillus nidulans*) that have been well studied genetically, are naturally haploid except for the immediate product of sexual fusion, which immediately undergoes meiosis to provide nuclei for haploid **spores**. In *Aspergillus*, however, rare vegetative diploids are obtained by selection (2). In *Neurospora*, where this is not possible, complementation tests are easily carried out with haploid heterokaryons, mixed cultures in which the filamentous hyphae have undergone fusions to bring the two different kinds of haploid nuclei together in the same cytoplasm (3).

For heterokaryon formation, *Neurospora crassa* strains have to be of the same mating type, and identical with respect to a number of other genes governing vegetative compatibility. These requirements are automatically met if the mutants under test have all been isolated in the same wild-type strain, as is usually the case. Given vegetative compatibility, fusions occur readily as soon as mixed mutant inocula of conidia (asexual spores) begin to germinate. *Neurospora* filaments (hyphae) are coenocytic. Numerous nuclei are present together in each cytoplasmic compartment, and the nuclear ratio in a heterokaryon is quite variable but is determined at least approximately by the ratio of the two kinds of conidia in the initial mixed inoculum. The nuclear ratio is, in fact, not usually critical. If, as is usually the case, the mutants under test are auxotrophic (ie, have special nutritional requirements) and the mixed inoculation is on a minimal medium on which neither mutant grows by itself, the heterokaryon grows vigorously if the mutants complement one another.

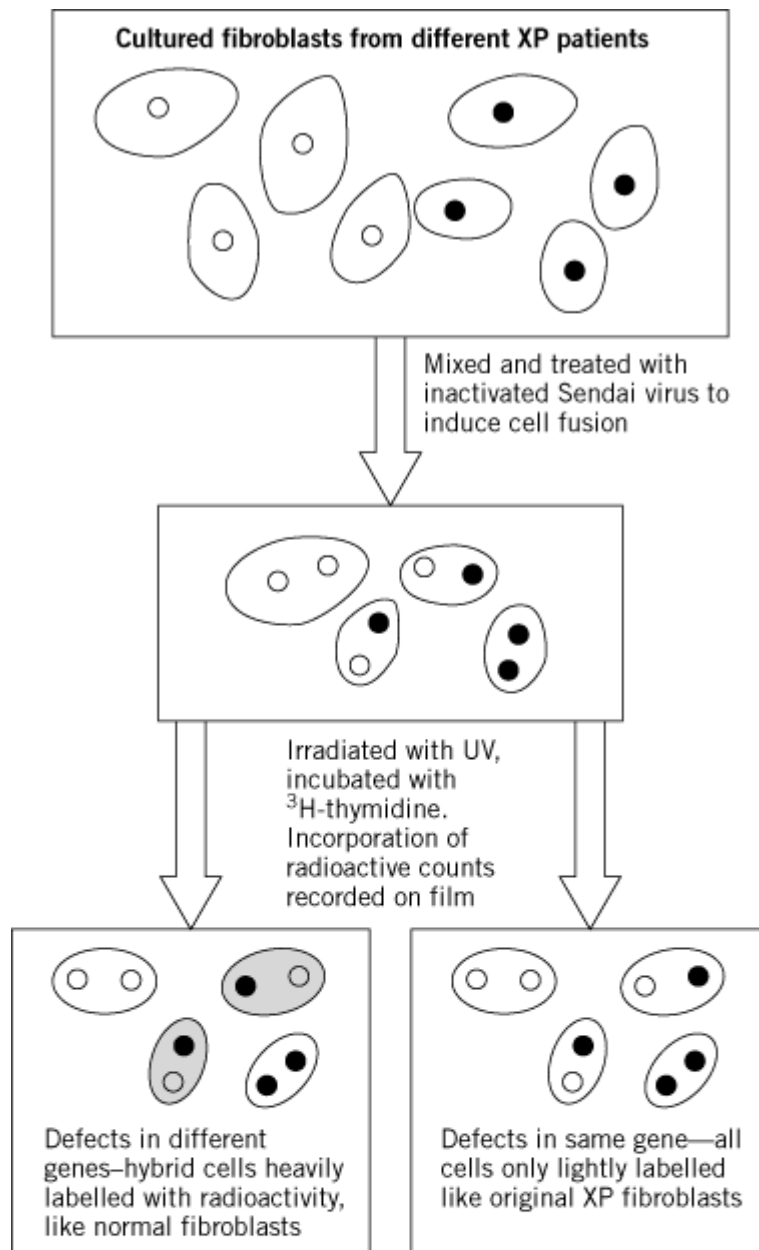
At least in *Neurospora*, different genes can complement each other perfectly efficiently in heterokaryons even though they are separated in different nuclei. This is not at all surprising. Genes act, after all, by exporting their [messenger RNA](#) to the **cytoplasm**. Reports indicate that complementation in *Aspergillus* is more efficient in diploids than in heterokaryons, but one suspects that this is because, in this fungus, the different nuclear types are not sufficiently closely intermingled. A large number of *Neurospora* auxotrophic mutants have been efficiently classified into complementation groups, equivalent to genes, by heterokaryon testing. An example is shown in Fig. 3 (4).

**Figure 3.** The result of complementation tests by heterokaryon formation on 46 *Neurospora* mutants that cannot utilize acetate. Nearly all possible pairwise mixtures of mutant conidia were inoculated into a medium containing acetate as the sole carbon source. Growth occurred only if the mutants showed complementation. In the matrix shown, the mutants, designated by their original isolation numbers, are divided into seven different complementation groups. All intergroup combinations showed complementation and grew like wild type. Most or all combinations within a group did not grow. The growth shown by a few pairs of mutants within group 6 is evidence for *inter allelic complementation*. From Ref. 4.



Humans, of course, are diploid, but complementation testing by controlled sexual crossing is obviously out of the question. Here again, heterokaryosis may provide an alternative, and cells in culture may serve as surrogates for people. One good example is an analysis of the human *xeroderma pigmentosum* syndrome, caused by genetic deficiency in [DNA repair](#) (5). Cultured cells isolated from different patients were induced to fuse by treating them with heat-inactivated [Sendai virus](#). Complementation in the resulting heterokaryotic cells was assessed by their ability or inability to incorporate a radioactive precursor into their DNA following damage by UV light. The results permitted assigning the different mutations to four different complementation groups, putatively different genes (Fig. 4).

**Figure 4.** The use of heterokaryosis for complementation testing of cultured cells from human Xeroderma Pigmentosum (XP) patients (5). The syndrome that is excessive sensitivity to uv light is a consequence of failure to repair damaged DNA by excision followed by new synthesis. (a) Cultured fibroblasts from different patients were mixed and induced to undergo fusions. The resulting cells, some now heterokaryotic, were irradiated and cultured in the presence of tritiated **thymidine** to see whether they responded to the DNA damage by incorporating this radioactive DNA precursor into their nuclei. The original XP mutant cells, or noncomplemented heterokaryons, incorporated no tritium, but certain pairs of mutant cells showed complementation. (b) The complementation matrix with the division of 12 XP mutations into four complementation groups A-D. Here the cells are diploid and homozygous for the different mutations instead of haploid as in the *Neurospora* example (Fig. 3), but the principle is the same. +, complementation; 0, no complementation; -, test not done.



(a)

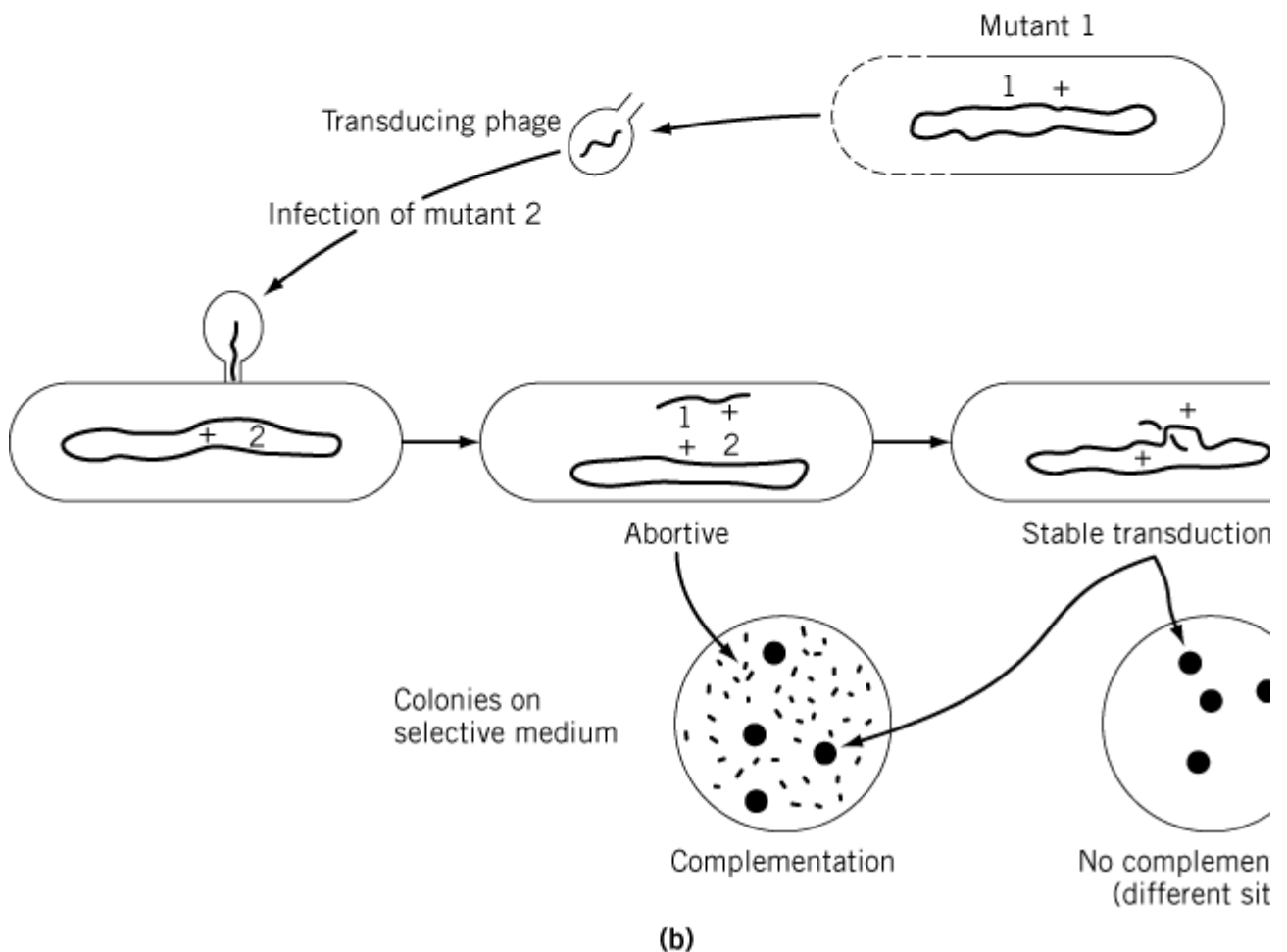
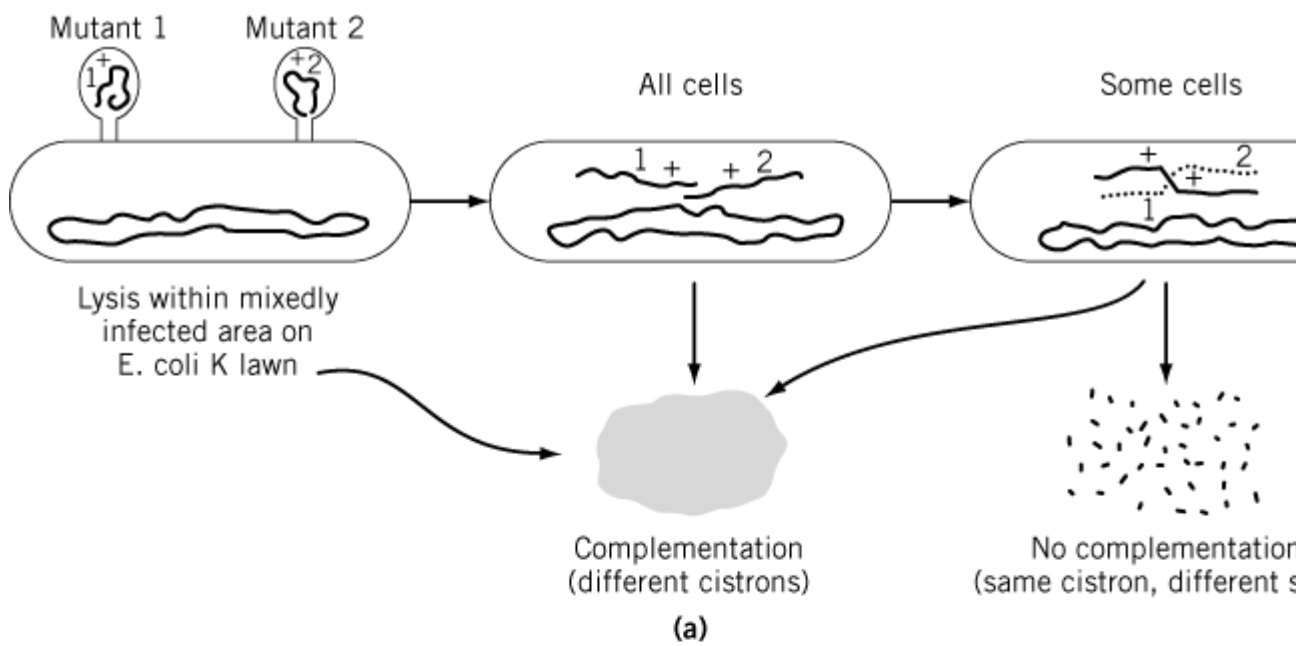
| A |   |   | B |   |   | C |   |   |   |   | D |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 1 | 2 | 3 | 1 | 2 | 3 | 4 | 5 | 1 | 2 | 3 |   |
| 0 | 0 | 0 | + | + | + | + | + | - | + | + | + | - | + | 1 |
| - | 0 | + | + | + | + | - | - | - | - | + | - | + | 2 | A |
| 0 | + | + | + | + | + | + | + | + | + | - | + | + | 3 |   |
| 0 | + | + | + | + | + | + | + | + | + | + | + | + | 1 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | + | + | + | + | 1 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | + | - | + | + | 2 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | + | - | + | + | 3 | C |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | + | - | + | + | 4 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | + | + | + | + | 5 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | D |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |   |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 |   |

(b)

## 2.1. In Viruses by Mixed Infection

S. Benzer's very extensive complementation testing of bacteriophage T4 mutants (see also [Complementation](#)) depended on infection of bacterial cells with mixtures of phage particles (6). Mixed infection is analogous to heterokaryosis because two different whole genomes are present in a somewhat indeterminate ratio, but the difference is that the simple viral genomes are not enclosed in **nuclei**. The mutants under study (class *rII*) were distinguished by growing on one *Escherichia coli* strain (B) but not on another (K). Complementation between pairs of mutants was checked by superimposing *rII* phage inocula on lawns of K cells, usually with a small admixture of B cells to allow at least some viral propagation. Three results were clearly distinguished: (1) complementation of the two mutants—total clearing of the mixedly infected patch due to lysis of all the cells in that area; (2) no complementation, but some formation of wild-type phage by recombination during mixed growth, shown by scattered spots of lysis; and (3) much less commonly, neither complementation nor recombination—no lysis at all (Fig. 5a). These results were interpreted to mean that (1) mutant sites are in different cistrons (ie, different functional genes); (2) mutants are at different sites in the same cistron; and (3) mutants are at the same site, or one is a deletion overlapping the other. The complementation results split the *rII* mutants cleanly into two mutually complementing groups or cistrons. Recombination analysis using a set of overlapping deletion mutants showed that these two cistrons correspond to two adjacent, nonoverlapping segments of the phage “chromosome”, which is not actually a chromosome in the cytological sense but a single molecule of **double-stranded DNA**.

**Figure 5.** Two examples from bacteria of complementation tests that distinguish between complementation and recombination. **(a)** Testing for complementation between *rII* mutants of bacteriophage T4. These mutants grow well in *Escherichia coli* B but not at all in strain K12(l). Mixed infections on a lawn of K cells that have a minor admixture of B cells to permit some joint growth give total clearing (lysis) if mutants *rII-1* and *rII-2* are in different genes (cistrons), and lysis only in spots of recombination, if they are at different sites in the same gene. Based on Ref. 6. **(b)** Abortive transduction as a criterion for complementation in *Salmonella typhimurium*. Phage P22 grown on histidine-requiring (*his*-) mutant 1 was used to infect mutant 2. Surviving bacteria were tested for their ability to grow on a histidine-free medium. Rather numerous tiny (abortive) colonies showed that a nonreplicating fragment of DNA from mutant 1 complements mutant 2. The absence of such tiny colonies signifies no complementation. In either case, vigorously growing colonies arise by recombination between the bacterial chromosome and the incoming fragment, provided that 1 and 2 are at different sites. Based on Ref. 7.



## 2.2. By Partial Diploids (Merodiploids) in Bacteria

It is not possible to test in bacteria for complementation between whole genomes because, at least in those species that are well-analyzed genetically, there is neither a stable diploid phase nor any sexual process that transfers a whole genome from one cell to another. However, bacteria acquire extra genomic fragments in three general ways: (1) by infection with defective bacteriophages which, “by

mistake”, have picked up bacterial DNA (**transduction**); (2) by joining a fragment of bacterial DNA to the DNA of a plasmid that brings about its own intracellular transfer (**conjugation**); and (3) by uptake of free DNA (**transformation**). The first two methods have been important in complementation analysis.

Transduction has been the principal means of classification and mapping of auxotrophic mutants of the bacterium *Salmonella typhimurium*. Bacteriophage **P22** grown on one mutant strain (the donor) is used to infect a second mutant strain, perhaps requiring supplementation with the same nutrient, such as an amino acid. Most of the infected bacteria are lysed, but some survive, are now lysogenic, and harbor the phage in latent form (see **Lysogeny**). Some survivors have acquired some fragment or other of the genome of the donor bacterium by infection from defective bacteriophage. The donor DNA is replicated in the recipient cell only if it recombines into the recipient chromosome and replaces the equivalent recipient segment. More commonly the donor DNA persists without replication, and its genes are expressed only in the primary recipient cell and in a few surrounding cells by diffusion of gene products. The appearance of tiny colonies of cells (abortive transductants) that grow transiently without the nutritional supplement, demonstrates complementation between the recipient genome and the donor genomic fragment. Usually a much smaller number of large strongly growing colonies results from chromosomal integration of a donor fragment and recombination between donor and recipient mutational sites to reconstitute a wild-type gene. The presence of large colonies without abortive ones shows that the mutations are at different sites within the same gene or cistron (Fig. **5b**). Failure to form any colonies at all means (apart from possible experimental error) that the two mutations are at the same site or that one is a deletion overlapping the other (**7**).

The plasmid that is used to bring about partial diploidy in *E. coli* is the transmissible “sex factor” F, a closed-loop DNA molecule that replicates autonomously in the bacterial cell and is occasionally integrated by **crossing-over** into the much larger loop of the bacterial chromosome. The integrated plasmid DNA undergoes occasional excision, usually precisely, but sometimes carries with it an adjacent segment of bacterial DNA, the nature of which depends on where the plasmid happened to have been inserted. This modified F is called F' and, like F itself, is transmitted from cell to cell. F'-*lac* plasmids that carry segments of DNA from the **lac operon** of the bacterial chromosome, were of great importance in elucidating the complementation relationships of *E. coli lac* mutants (a topic dealt with separately under **cis-dominance**).

## Bibliography

1. J. Gausz, H. Gyurovics, G. Beneze et al. (1981) *Genetics* **98**, 775–789.
2. J. A. Roper (1952). *Experientia* **8**, 14–15.
3. G. W. Beadle and V. L. Coonradt (1944) *Genetics* **29**, 291–308.
4. R. B. Flavell and J. R. S. Fincham (1968) *J. Bacteriol.* **95**, 1056–1062.
5. K. H. Kraemer, H. G. Coon, R. A. Petinga et al. (1975) *Proc. Natl. Acad. Sci. USA* **72**, 59–63.
6. S. Benzer (1959) *Proc. Natl. Acad. Sci. USA* **45**, 1607–1620.
7. P. E. Hartman, Z. Hartman, and D. Serman (1960) *J. Gen. Microbiol.* **22**, 354–358.

## Complex Loci

In classical **genetics**, the word locus referred to a position on a **chromosome** marked by a **gene** difference that has a **phenotypic** effect. Neither the locus nor the gene were thought to be subdivisible by **recombination**, and so locus and gene were treated almost as synonyms. The idea of



complex loci arose from examples, especially from *Drosophila* and **maize**, of different [mutations](#), apparently at the same chromosome locus, that give a range of different phenotypic effects, suggesting that the locus harbors some rather complex kind of gene.

With our modern knowledge, the word locus can take either of two very different meanings. It may mean just the genetically mapped site of a mutation, which, at the molecular level, is quite likely to be no more than a single base-pair change within the long DNA sequence of a gene. But, in the context of complex loci, it can mean a relatively short chromosomal segment (not frequently split up by genetic recombination) that contains a cluster of genes among which there is some functional connection. So complex loci, though of somewhat confused etymology, is a useful heading under which to review the evidence for chromosomal units of function at a level higher than that of the individual gene.

This article does not cover examples of loci that mutate to give various phenotype effects when the varied effects are due to mutations in a single gene that has multiple functions. Such genes that encode “multiheaded” enzymes (see [Multifunctional Proteins](#)) are sometimes called “cluster genes” in fungal genetics to distinguish them from the [gene clusters](#) with which they may initially be confused. Some examples are reviewed under **Allelic Complementation**. Nor are we concerned with single genes that have multiple modes of **intron** splicing, which yield alternative [messenger RNAs](#) (see examples under [Gene Structure](#)). And we obviously cannot include groups of functionally diverse genes that just happen to be very closely linked, for that would make virtually all genes parts of one complex locus or another. To qualify as a complex locus, in our present sense, a gene cluster must have some kind of unity in structure, in function, or in both.

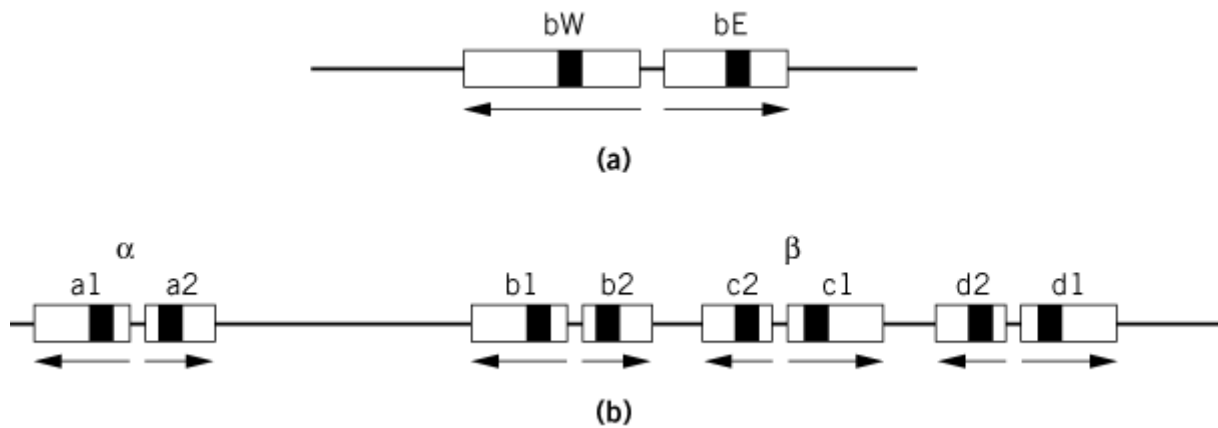
Gene clusters that fulfil this criterion have been found in **fungi**, flowering **plants**, *Drosophila* and mammals, and this article reviews examples from each of these kinds of organisms. We do not deal here with **bacteria**. [Operons](#), which might be viewed as complex loci but are not usually so termed, are dealt with under their own heading. Also excluded from this article are highly repetitious tandemly arrayed genes required by the cell in very high dosage. The preeminent example is **ribosomal** RNA genes in the **nucleolus organizer** chromosomal regions of eukaryotes. [Histone](#) genes are also highly repetitive and are concentrated in clusters in some organisms, including *Saccharomyces* yeast. This is arguably more repetition than complexity.

## 1. Fungi

In fungi, the self-incompatibility and cross-compatibility in sexual mating in many species is determined by **mating-type** loci, which virtually always consist of more than one gene. A relatively simple example is the maize smut fungus, *Ustilago maydis*, whose two loci, *a* and *b*, control the specificity of sexual fusion and the formation of an invasive dikaryotic mycelium after fusion, respectively. The *a* locus has two “**alleles**,” *a1* and *a2*, each of which consists of two genes, one that encodes a polypeptide pheromone and the other a pheromone **receptor**. The *a1* pheromone is recognized by the *a2* receptor and the *a2* pheromone by the *a1* receptor. There is no close [homology](#) between *a1* and *a2*, and they do not recombine with each other. The *b* locus exists in numerous alternative forms. Each consists of two closely spaced, divergently transcribed genes called *bE* and *bW* that encode the two components of a dimeric transcriptional activator. But dimerization occurs only between *bE* and *bW* products of different mating types (1).

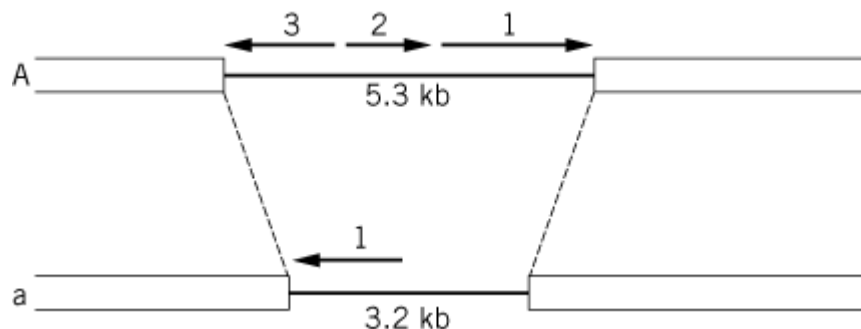
The ink-cap mushroom fungus, *Coprinus cinereus* (like *Ustilago*, a member of the Basidiomycete class), also has two mating-type loci, *A* and *B*. The *A* locus, which controls the ability to form a fruiting dikaryon, is the most thoroughly investigated. Here the situation resembles that at the *Ustilago b* locus, except that genes that encode dimerizing monomers, thought to be transcriptional activators on the basis of their **homeobox**-like sequence motifs, are distributed in multiple copies between two subloci, *Aa* and *Ab*, between which rare recombination occurs (2) (Fig. 1). Again, dimerization occurs only between the products of different mating types.

**Figure 1.** The structures of mating type loci in Basidiomycete fungi. **(a)** The *b* locus of *Ustilago maydis*. The *bW* and *bE* gene products form W-E protein dimers to activate transcription of genes necessary for sexual development. Each exists in multiple allelic forms. Effective dimers are not formed by W-E combinations from the same strain. Gene segments that encode homeodomain-like protein sequences are filled. **(b)** The archetypal *A* mating type locus of *Coprinus cinereus*. The system works in principle as in the *Ustilago b* locus, except that here there are four pairs of genes, *a*, *b*, *c*, and *d*, any one of which can function in promoting fertility. The archetype is a composite of known mating-types, all of which are missing one or another part of it. From Refs. 1 and 2 by permission.



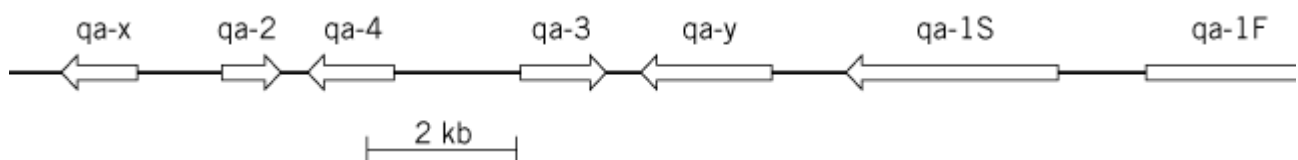
Ascomycete fungi, such as *Neurospora crassa* and *Saccharomyces cerevisiae*, simpler than Basidiomycetes, have only two mating types (*A* and *a* in *Neurospora*, *a* and *a* in *Saccharomyces*), but the genetic loci that determine them are still complex. *Neurospora A* contains three genes (Fig. 2) and *Saccharomyces a* contains two, though each of the *a* mating-type loci has only a single gene. The two alternative mating-type loci in both organisms are completely nonhomologous, and the same is true for the *Ustilago a* locus. It has been customary to call the alternative occupants of these mating-type loci “alleles,” but in the light of present knowledge this terminology is inappropriate because it implies different forms of a single gene. The new word idiomorphic has been suggested for the relationship between quite different genes or gene clusters that are alternative occupants of the same locus (3). Much the same applies to the yeast mating-type system that has the additional complication of mating-type switching.

**Figure 2.** The mating type loci of the Ascomycete fungus *Neurospora crassa*. *A1* and *a1*, quite unrelated in sequence, act in a complementary way to permit sexual reproduction. *A2* and *A3* function in postfertilization fruiting-body development. From Ref. 3 by permission.



Fungi also provide examples of another kind of locus complexity, the tight clustering of separate genes that co-operate in a particular area of metabolism. The *Neurospora crassa* *Qa* cluster is a particularly good example (Fig. 3). Here a DNA segment of only 18 kbp harbors seven different genes, the transcription of which is increased 50- to 1000-fold when the fungus is forced to grow using quinate as a sole carbon source. Three of these genes encode the three enzymes that, acting in sequence, oxidize quinic acid to protocatechuic acid. Two genes encode regulatory proteins that control the [transcription](#) of the whole cluster. One gene is a transcriptional activator and the other a repressor that cancels the function of the activator but whose own effect is nullified by binding to quinate. The remaining two genes have no identified function and are recognized solely as **open reading frames** (4). The significance of the gene clustering is not fully understood, but it seems plausible to suggest that derepression of transcription is facilitated by an induced unfolding of **nucleosomal** structure that extends across the entire complex.

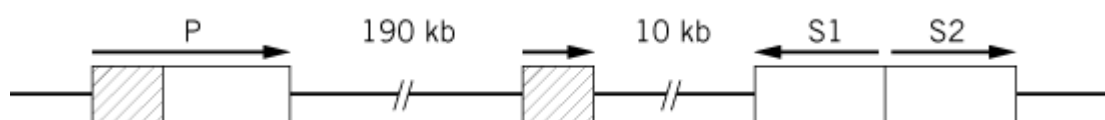
**Figure 3.** The *Qa* (quinate-utilization) gene cluster of *Neurospora crassa*. The genes *qa-3*, *qa-2*, and *qa-4*, respectively, encode quinate dehydrogenase, catabolic dehydroquinase, and dehydroshikimate dehydratase, which act sequentially to degrade quinate. *qa-1S* and *qa-1F* encode regulatory proteins, that respectively repress and activate transcription of the whole cluster. *qa-Y*, it is thought, encodes a quinate transporter protein, and *qa-X* has no known function (4).



## 2. Flowering plants

A long-standing example from what used to be, pre-*Arabidopsis*, the geneticists' favorite plant, is the *R* locus of *Zea mays* that is necessary for pigmentation and encodes transcriptional activators for genes that encode pigment-synthesizing enzymes. It was an early example of locus complexity because it has two independent functions. Some recessive loss-of-pigment alleles affect only the red color in the plant whereas others affected only the seed. Analysis at the DNA level has revealed the situation shown in Fig. 4 (5). The *R* locus contains three genes. One gene (*P*) is specific for plant pigment, and the other two (*S1*, *S2*) apparently duplicate the function in seeds. Presumably their distinct tissue specificities are due to different promoter/**enhancer** sequences. There is also a partial and presumably functionless duplication of *P*. An origin by **duplication** and subsequent partial **divergence** of function is the most obvious explanation for the clustering of these functionally related genes, although the opposite orientations of *S1* and *S2* are a puzzle.

**Figure 4.** The complex *R* locus of *Zea mays*. *P* controls red pigment synthesis in the plant generally and the duplicate *S* genes in the seed. *S1* and *S2*, transcribed divergently, are each similar to *P* except that they lack the segment shown hatched. This segment is present as a tandemly duplicated fragment, presumably nonfunctional. Sequence information suggests that this scrambled arrangement is due to the activity of a transposable element (5). “Alleles” of *R*, *R<sup>S</sup>* (red seed, green plant), and *r<sup>r</sup>* (colorless seed, red plant) are, respectively, due to loss of *P* and *S* function. The *P* and *S* components can be separated by meiotic crossing over at frequencies on the order of  $10^{-4}$  to  $10^{-3}$ .

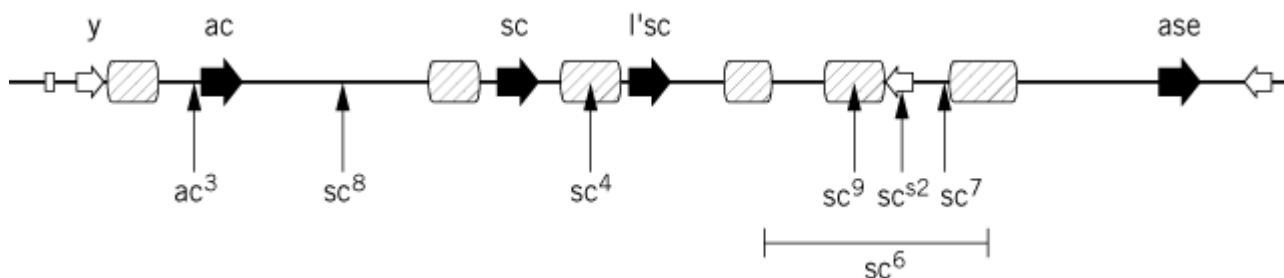


A long-standing problem in plant biology has been the various systems of self-incompatibility, often controlled by what are multiple alleles at a single locus, usually called *S*. The general principle is that pollination is prohibited if the pollen grain and stigma (or, in the *Brassica* family, the pollen mother cell and stigma) have an allele in common. The *Brassica* *S* locus, one of the best analyzed, contains two genes, each present in the population in multiple allelic forms, that control the incompatibility response of the stigma. A third *S* locus gene, expressed in the pollen and determining pollen specificity, has also been identified (6). Here, as in the analogous self-incompatibility in fungi, the functional reason for the multiple genes all in the same locus is obvious. For the system to work, the different specificities have to be mutually exclusive in the [haploid](#) genome.

### 3. *Drosophila*

In the fruit fly, *Drosophila melanogaster*, a classical example of locus complexity, perhaps originally generated by tandem duplications, is provided by the scute (*sc*) and achaete (*ac*) families of mutations (7, 8). These affect the bristles in the wings and mid-thorax of the fly, really parts of the fly's neural system, the origin of which can be traced back to the transcription of scute and achaete, usually both together but sometimes one more than the other, in rather precisely positioned bristle mother cells in the wing [imaginal disc](#) of the larva. It seems that *sc* and *ac* largely duplicate each other's functions and, before molecular analysis, the basis for recognizing two different genes was not clear. Different mutations in the *sc/ac* locus eliminate different bristles or different combinations of bristles, and, to the extent that their effects are nonoverlapping, complement one another in *trans*. Now that the *sc/ac* locus has been analyzed at the DNA level, it is seen that there are six genes within the approximately 100 kbp that comprise the *sc/ac* complex (7, 8) (Fig. 5), of which at least four are involved with the neural/bristle system. *sc* and *ac* apparently have the major roles. Most of the mutations that affect the positioning of bristles actually fall between the genes. Several are chromosomal rearrangements that have breakpoints within the *ac/sc* complex, and it is thought they exert their effects by separating *ac* and/or *sc* from enhancer sequences, several of which have been identified in the *sc/ac* region. The patterns of transcriptional activation of *sc* and *ac* and hence the positioning of bristles, it is thought, reflect the distribution of a number of enhancer-binding proteins, the identity of which remain to be established. The genes and the enhancers together may be seen as an integrated system for responding to a pattern of positional information set by the products of other genes.

**Figure 5.** The *scute-achaete* (*sc/ac*) complex at the tip of the *Drosophila melanogaster* X-chromosome. Genes that have functions in neural/bristle development are shown as filled arrows, and other genes, apparently unrelated in function, as open arrows (*y* = yellow). Thin vertical arrows indicate the breakpoints of segmental rearrangements that affect the bristle phenotype (*sc*<sup>6</sup> is a deletion). Regions that include enhancer sequences are shown hatched. The different bristle phenotypes are believed to be due to separation of *sc* and/or *ac* from enhancers. Redrawn and simplified from Ref. 8 by permission.



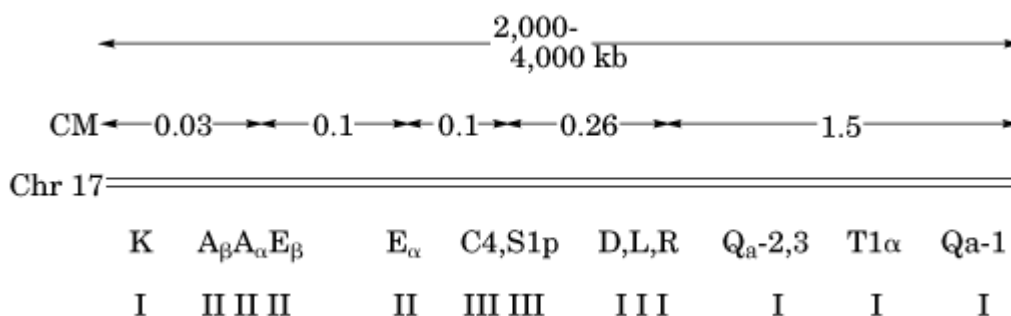
### 4. Mammals

Many clusters of related genes have originated through tandem duplication and subsequent functional divergence. Obvious examples are the globin gene clusters in animals and humans (see under [Globins](#), [Gene Structure](#)). The b-globin gene cluster has a common ancestral origin and is also functionally integrated through the locus control region.

The clustered **HOX genes**, which are of profound importance in animal development and encode homeobox-containing proteins, were discovered in mammals and now in other animal groups through their significant similarities to the genes of the linked *bithorax* and *Antennapedia* complexes of *Drosophila*. Five unlinked clusters of HOX genes in mouse and humans each consist of up to 13 tandemly arranged genes within a region on the order of 100 kbp of DNA. Through their homeobox-encoding segments, the genes of one cluster have only a very limited degree of sequence similarity to each other but have a greater resemblance to genes in corresponding positions in other clusters in the same species or even in different species. Because they were not identified through the classical genetic route of mutation mapping, the HOX clusters of mammals are not usually called loci though, confined as they are within chromosomal segments that correspond to a small fraction of a **centimorgan**, they are sufficiently tightly linked to qualify for that title, and they also fulfil the criterion of functional integration. The almost uncanny correspondence between the linear sequence of genes within each cluster and the anterior-posterior order of the body segments that the genes affect clearly indicates some function for the clustering and ordering of the genes, though this is still mysterious.

The mammalian [major histocompatibility complex](#) (MHC), intensively investigated in mouse, and its human equivalent the human leucocyte antigen (HLA) complex, is encoded by genes spread over rather too large a section of chromosome to be called a locus. Indeed, the positions of individual antigen-encoding genes within it are themselves often called loci. In humans it encompasses about four megabases of DNA and 3 to 4 centimorgans. But it may be regarded as at the upper end of the size range of clusters of functionally related genes, and so it is relevant to mention it here. A diagram showing of the structure of the complex, and its scale in relative to the whole chromosome is shown in Fig. 6. Because there is significant homology among the different genes, at least within each of the two main subclusters, the close linkage is likely to be due to duplication by short-range [transposition](#) and subsequent diversification. But it may also have been maintained by [natural selection](#), because their continued coupling is advantageous if certain combinations of alleles of different genes confer disease resistance ([10](#)).

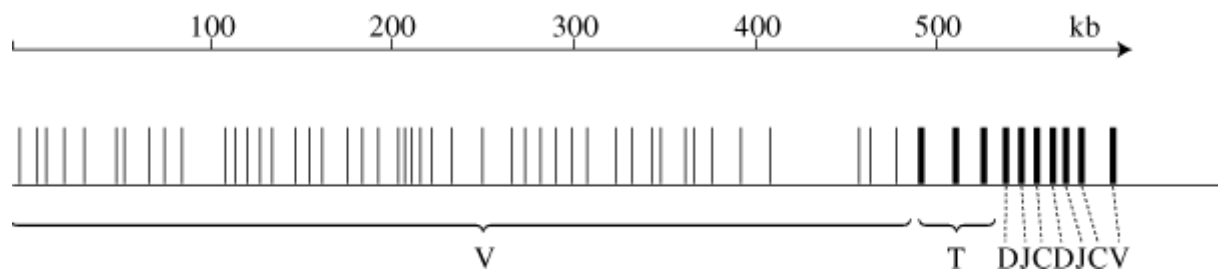
**Figure 6.** The arrangement and spacing of genes within the major histocompatibility complex (MHC) of mouse. Class I genes encode transplantation antigens mediating graft rejection (K, D, L), or other cell-surface antigens (Qa, Ti). Class II genes encode antigens involved in the immune response. Class III genes encode components of the complement system. Meiotic recombination within the complex is infrequent but not very rare (genetic map distances are shown in centimorgans, cM). After Ref. [9](#) by permission.



An entirely different kind of complex locus is involved in determining hypervariable proteins of the mammalian immune system, [immunoglobulins](#) and [T-cell receptor](#) proteins. These are not really

gene clusters but rather banks of gene components that are assembled by DNA splicing reactions in a large number of different ways to provide coding sequences for proteins of almost any required specificity. In mammals, as exemplified by mouse, there are six known rearranging complexes of the mouse immune system, three (heavy-chain and two light-chain) for immunoglobulin components and three (a, b and g polypeptides) for T-cell receptor components. A diagram to scale of the completely sequenced human b-chain locus that extends over 600 kb is shown in Fig. 7 (11). Here the functional reason for clustering is obvious. It is only surprising that the sequences destined to be spliced together are as widely separated as they are. Presumably they are brought close during the splicing process by the formation of some very extensive DNA loops.

**Figure 7.** The arrangement and spacing of the part-genes which, when appropriately spliced together, encode b-chains of the human T-cell receptor. Splicing of DNA during T-cell differentiation brings together one C “constant”, one J (“junction”), one D (“diversity”), and one V (“variable”) sequence to encode a large number of alternative b-chains. J on the map stands in each case for six or seven short alternatively used junction sequences. The other marked elements are all single units. All are transcribed/translated from left to right, except for the rightmost V which is inverted and doubtfully functional. The splicing occurs in stages, first D to J, then DJ to V, and finally VDJ to C. The diagram has been simplified by omitting the [pseudogenes](#) and incomplete V sequences, which approximately equal the functional genes in number and are interspersed among them. T = trypsinogen gene, the presence of which, in multiple copies, within the b-chain complex, has uncertain functional significance. After Ref. 11 by permission.



## Bibliography

1. F. Banuett (1992) *Trends in Genetics* **8**, 174–180.
2. E. H. Pardoe, S. F. O'Shea, and L. S. Casselton (1996) *Genetics* **144**, 87–94.
3. N. L. Glass and M. A. Nelson (1994) In *The Mycota I*. (J. G. H. Wessels and F. Meinhardt, eds.), Springer-Verlag, Berlin, pp. 295–306.
4. R. F. Geever, L. Huiet, J. A. Baum, B. M. Tyler, V. B. Patel, B. J. Rutledge, M. E. Case, and N. H. Giles (1989) *J. Mol. Biol.* **207**, 15–34.
5. E. L. Walker, T. P. Robbins, T. E. Bureau, J. Kermicle, and S. L. Dellaporte (1995) *EMBO J.* **14**, 2350–2363.
6. C. R. Schopfer, M. E. Nasrallah, and J. B. Nasrallah (1999) *Science* **286**, 1697–1700.
7. S. Campuzano and J. Mondolell (1992) *Trends in Genetics* **8**, 202–207.
8. J. L. Gomez-Skarmeta, I. Rodriguez, C. Martinez, J. Culi, D. Ferres-Marco, D. Beamonte, and J. Mondolell (1995) *Genes Dev.* **9**, 1869–1882.
9. L. Hood, M. Steinmetz, and B. Malissen (1983) *Annu. Rev. Immunol.* **1**, 529–568.
10. I. P. M. Tomlinson and W. F. Bodmer (1995) *Trends Genet.* **11**, 493–498.
11. L. Rowen, B. F. Koop, and L. Hood (1996) *Science*, **272**, 1755–1762.

## Compressibility

Compressibility is an important measure of the atomic packing and flexibility of protein molecules. Measurement of this quantity is possible, in principle, only for proteins in solution. The partial specific compressibility of a protein in solution  $b$  is defined as the change in the [partial specific volume](#)  $v_2^0$  with increasing pressure. Since the atomic [van der Waals volume](#),  $v_c$ , may be assumed to be incompressible, the experimentally determined  $b$  of a protein can be mainly attributed to its cavities  $v_{cav}$ , and to changes in [hydration](#)  $Dv_{sol}$ :

$$\bar{\beta} - (1/v_2^0)(\partial v_2^0/\partial P) = -(1/v_2^0)[(\partial v_{cav}/\partial P) + (\partial \Delta v_{sol}/\partial P)] \quad (1)$$

(see [Partial Specific \(Or Molar\) Volume](#)). The first term on the right-hand side contributes positively, and the second term negatively, to  $b$ . Thus, a positive  $b$  value can be ascribed to a large cavity effect overcoming the hydration effect. There exist two types of compressibilities, depending on the experimental conditions: adiabatic and isothermal. Most compressibility studies of proteins have been of adiabatic compressibility because of its accuracy and technical convenience.

### 1. Adiabatic Compressibility

The partial specific adiabatic compressibility  $b_s$  is calculated with the Laplace equation,  $b_s = 1/u^2 d$ , where  $u$  is the sound velocity and  $d$  the density of the protein solution. The sound velocity is usually measured with an accuracy of  $1 \text{ cm sec}^{-1}$  by using a resonance or ring-around pulse method at 3 to 6 MHz. All amino acids show large negative  $b_s$ , with values between  $-62.5 \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$  (for [glycine](#)) and  $-21 \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$  (for [tryptophan](#)) at  $25^\circ\text{C}$ , because they have a negligibly small amount of cavity. Data for  $b_s$  have been reported for about 40 different proteins; some values in water (or in dilute buffer at neutral pH) are listed in [Table 1](#). Fibrous proteins show negative  $b_s$ , indicating the dominant effects of the surface on hydration relative to internal cavity volume. On the contrary, most globular proteins have positive  $b_s$ , due to the large contribution of the internal cavities overcoming the hydration effect. This is consistent with the fact that more than 50% of the total surface area is buried in the interior of a folded protein molecule.

**Table 1. Partial specific volume ( $v_2^0$ ), compressibility ( $b_s$  and  $b_T$ ), and volume fluctuation ( $dV_{rms}$ ) of proteins in water at  $25^\circ\text{C}$  (1, 4)**

| Protein                   | $v_2^0$<br>( $\text{cm}^3 \text{ g}^{-1}$ ) | $b_s$ ( $10^{-12}$<br>$\text{cm}^2 \text{ dyn}^{-1}$ ) | $b_T^a$ ( $10^{-12}$<br>$\text{cm}^2 \text{ dyn}^{-1}$ ) | $dV_{rms}^b$<br>( $\text{cm}^3 \text{ mol}^{-1}$ ) |
|---------------------------|---|--|--|--|
| <i>Globular proteins</i>  |   |  |  |  |
| Cytochrome c              | 0.725                                       | 0.07   | 4.27   | 30.8 (0.34)  |
| Soybean trypsin inhibitor | 0.713                                       | 0.17   | 4.44   | 41.1 (0.20)  |
| Trypsin                   | 0.717                                       | 0.92   | 5.16   | 46.0 (0.28)  |
| Ribonuclease A            | 0.704                                       | 1.12   | 5.48 (5)   | 36.2 (0.38)  |

|                         |       |             |             |             |
|-------------------------|-------|-------------|-------------|-------------|
| Peroxidase              | 0.702 | 2.36        | 6.70        | 68.3 (0.24) |
| a-Chymotrypsin          | 0.717 | 4.15        | 8.32        | 62.2 (0.33) |
| Lysozyme                | 0.712 | 4.67        | 7.73 (12.3) | 44.2 (0.43) |
| Carbonic anhydrase      | 0.742 | 6.37        | 10.5        | 76.0 (0.34) |
| a-Lactalbumin           | 0.736 | 8.27        | 12.4        | 56.9 (0.54) |
| b-Lactoglobulin         | 0.751 | 8.45        | 11.8        | 63.6 (0.46) |
| Myoglobin               | 0.742 | 8.98        | 13.1        | 64.1 (0.51) |
| Ovalbumin               | 0.746 | 9.18        | 12.1        | 101 (0.30)  |
| BSA                     | 0.735 | 10.5        | 14.6 (13.4) | 135 (0.27)  |
| <i>Fibrous proteins</i> |       |             |             |             |
| Gelatin                 | 0.689 | -2.5        |             |             |
| F-actin                 | 0.720 | -6.3 (20°C) |             |             |
| Myosin                  | 0.724 | -18 (20°C)  |             |             |
| Tropomyosin             | 0.733 | -41 (20°C)  |             |             |

<sup>a</sup> The values in parentheses represent the experimental data; the other values were calculated with Equation 2 of the text.

<sup>b</sup> The values in parentheses represent the ratio (%) of  $dV_{\text{rms}}$  to the total protein volume.

Some assumptions are necessary for estimating separately the contributions of cavities and hydration to  $b_s$ . A rough estimation suggests that the intrinsic compressibility of globular proteins free from hydration is of the order of  $(10 \text{ to } 20) \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$ . This value is comparable to the adiabatic compressibility of normal ice, suggesting that an internal protein structure is as rigid as ice.

As can be seen in Table 1, the value of  $b_s$  varies over a wide range and is sensitive to the structural characteristics of individual proteins. For example, [lysozyme](#) has a considerably smaller  $b_s$  than does a-lactalbumin, in spite of their great similarities in primary and [tertiary structures](#). However, the compressibility-structure relationships of proteins have been scarcely discussed on a molecular level, because of the complicated contribution of [hydration](#). The contributions of some structural factors have been deduced by the statistical analyses of  $b_s$  data on globular proteins (Ref. 1). From the ratio of accessible surface area to volume,  $b_s$  would be expected to increase with increasing molecular weight, but there appears no definite correlation between  $b_s$  and molecular weight. Instead, a positive correlation is found between  $b_s$  and the partial specific volume. Hydrophobic proteins show large positive  $b_s$ , probably due to an enhanced imperfect packing of nonpolar residues localized in the interior of the molecule. Typical  $\alpha$ -helical proteins, such as [myoglobin](#) and bovine [serum albumin](#) have very large  $b_s$ , even though the  $\alpha$ -helix itself is rigid, predominantly nonhelical proteins, such as [trypsin](#) and [soybean trypsin inhibitor](#), have small  $b_s$ . These results suggest that the  $\alpha$ -helix could be a dynamic domain for the thermal fluctuation of proteins. Four amino acids (Leu, Glu, Phe, and His) show statistically a strong ability to increase  $b_s$ , whereas another four (Asn, Gly, Ser, and Thr) decrease it, although the meaning of this is obscure. A single amino acid substitution can bring about a noticeable change in  $b_s$ , probably due to modified atomic packing, even though it is rare to observe any visible changes in the tertiary structures of mutants by [X-ray crystallography](#).



Adiabatic compressibility data for nonnative conformations and for unfolding processes, in combination with partial specific volume, also present useful information on the principles of protein structure that cannot be obtained by spectroscopic techniques. Unfolding of a protein with strong denaturants such as **guanidinium** chloride decreases the compressibility and specific volume, but it is known that the compressibility increases and the partial specific volume decreases on thermal and pressure denaturation.

## 2. Isothermal Compressibility

The volume fluctuations and the pressure-dependent properties of proteins are theoretically related to the partial specific isothermal compressibility  $b_T$ , rather than the adiabatic one  $b_s$ . To determine  $b_T$ , the solution density or partial specific volume must be measured as a function of pressure, under hydrostatic pressure or centrifugal force. Such an experiment is very difficult, however, and few  $b_T$  data are available for proteins, because high pressure may cause protein denaturation and modify preferential solvent interactions. Approximate values of  $b_T$  may be estimated from  $b_s$  by utilizing the relationship

$$\beta_T = \beta_s + \alpha^2 T / C_p d \quad (2)$$

if the thermal expansion coefficient  $\alpha$  and the **heat capacity** at constant pressure  $C_p$ , are known or can be reasonably inferred. Such assumed  $b_T$  values are listed in Table 1, along with some experimental data. Despite many assumptions, the calculated values of  $b_T$  are not so very different from the experimentally observed ones. The value of  $b_T$  is greater than  $b_s$  by  $(3 \text{ to } 4) \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$ ; these differences are comparable to those found for the amino acids.

There are some other methods to evaluate the isothermal compressibility of proteins, although the value obtained is not a partial quantity in solution. Kundrot and Richard (2) reported  $4.7 \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$  for  $b_T$  of lysozyme by comparing its X-ray crystallographic volume at atmospheric pressure and 1000 atmospheres. Contraction of the molecule was distributed nonuniformly; one domain was essentially incompressible, but another had a large compressibility. Computer simulations, such as **molecular dynamics**, **Monte Carlo**, and normal mode analyses, are also available. Normal mode analysis of myoglobin yielded  $b_T = 9.37 \times 10^{-12} \text{ cm}^2 \text{ dyn}^{-1}$ .

According to statistical thermodynamics, the volume fluctuation  $dV_p$  of a protein with volume  $V_p$  is related to its isothermal compressibility  $b_T$  (Ref. 3):

$$\overline{\partial V_p^2} = k T V_p \beta_T \quad (3)$$

where  $k$  is the Boltzmann constant and  $T$  the absolute temperature. The root-mean-square fluctuation of the partial molar volume  $dV_{\text{rms}}$ , which was estimated by using  $b_T$  instead of  $b_s$ , is listed in the last column of Table 1. The volume fluctuation is only about 0.3% of the overall dimensions of protein, but the  $dV_{\text{rms}}$  values estimated by this method are the lower limit of the probable fluctuations, since the solvation factor is still included in the value of  $b_T$  that was used. If the intrinsic isothermal compressibility of a protein itself, which is not easy to determine, is used instead of  $b_T$ , a greater volume fluctuation (up to 50% greater) might be expected for a protein without hydration. Although many assumptions are used in estimating  $b_T$  and the fluctuation of volume, the values obtained are considered to be reasonable; if concentrated in one area at one particular moment, the fluctuation in volume could produce sufficient cavities or channels to allow the entry of solvent or probe molecules

to account for phenomena such as [hydrogen exchange](#) and fluorescence quenching observed with folded proteins.

### Bibliography

1. K. Gekko and Y. Hasegawa (1986) *Biochemistry* **25**, 6563–6571.
2. E. Kundrot and F. M. Richards (1987) *J. Mol. Biol.* **193**, 157–170.
3. A. Cooper (1976) *Proc. Natl. Acad. Sci. U.S.A.* **73**, 2740–2741.
4. K. Gekko (1991) In *Water Relationship in Food* (H. Levine and L. Slade, eds.), Plenum Press, New York, pp. 753–771.

### Suggestions for Further Reading

5. A. P. Sarvazyan (1991) Ultrasonic velocity of biological compounds. *Annu. Rev. Biophys. Chem.* **20**, 321–342.
6. K. Gekko (1996) Hard proteins and soft proteins: Structural flexibility as revealed by compressibility. *Protein, Nucleic Acid and Enzyme (Tanpakusitu, Kakusan, Koso, in Japanese)* **41**, 2025–2036.

## Computer Simulation of Biological Molecules

Biological molecules are systems of enormous complexity. This makes it exceedingly difficult to elucidate the details of biological activities using only experimental techniques. Despite enormous progress in structural studies of [proteins](#), it is hard to determine uniquely how the particular function is determined by the given structure. A mathematical description of the energy of a macromolecule in terms of its atomic coordinates is, in principle, sufficient to predict relevant biological properties. However, such a description can only be handled by powerful computers. Having the proper model of a macromolecule should allow one to simulate its properties as if one is dealing with the actual molecule. Such an approach, called *computer simulation*, offers what is probably the best way of describing the function of macromolecules.

The relationship between the energy of molecules and the position of its atoms is called a [potential function](#) or a *force field*. A reasonable potential function provides the first step for modeling a molecule. In principle, one can obtain potential functions by using quantum mechanical methods that describe the entire molecule as a supermolecule. However, such methods are not yet practical for describing macromolecules. Some progress can be made using **hybrid QM/MM approaches**, but, in general, it is essential to use approximated potential functions whose parameters are based on empirical information ([1-5](#)).

Potential functions can be based on **all-atom molecular models** that describe the interactions between all the atoms of the given molecule or on **simplified models** where the forces due to several atoms are described by a single interaction center. Examination of molecular properties using the corresponding potential functions is known as [molecular mechanics](#). This term is usually reserved for the use of **energy minimization** and **normal mode analysis** methods. The use of potential functions to solve the equation of motion of atoms is known as [molecular dynamics](#) (MD). This method gives the atomic position as a function of time and provides a powerful way of obtaining dynamical and time average properties. In some cases one can obtain useful insight by using **Brownian dynamics**. Frequently it is sufficient to know the properties of macromolecules at their minima. Generating starting configurations in the search for such minima can be done by **Monte**

**Carlo methods** [Monte Carlo calculations](#) and other sampling approaches. Apparently, the most important biological properties are reflected in most cases by thermodynamic properties and not dynamical properties. In particular the **Free Energy** provides the most important connection between structure and function, while the knowledge of the individual contributions to free energy from **enthalpy** and **entropy** is less critical for reproducing functional properties of macromolecules. The free energy of macromolecules can be evaluated by several approaches and, in particular, by the free energy perturbation (FEP) and linear response approximation (LRA) methods. [Free energy calculations](#) are frequently based on the use of **thermodynamic cycles** (3, 4, 6).

The potential functions that describe the energetics of macromolecules reflect different types of interactions. This includes **bonding** and [van der Waals interactions](#) that determine the steric properties of the given molecule and [electrostatic interactions](#) (that include [hydrogen bonding](#) and other **dipolar** and **ionic interactions**). The free energy of macromolecules is largely dependent on electrostatic energies, which reflect the effect of polar and ionized groups and its compensation by the polarization of the solvent around the protein (6, 7). These energies involve **long-range interactions**, whose accurate evaluation present a major challenge. Electrostatic energies can be modeled by using **implicit models** where the solvent molecules (and sometimes the protein dipoles) are represented implicitly. The use of a **dielectric constant** in an electrostatic model is a form of implicit model, and, in fact, the dielectric constant in many electrostatic models is not related at all to the true protein dielectric constant but to the properties that are not simulated explicitly (6, 8). *Computer simulations of biological molecules has been particularly effective in studies of the following subjects.*

## 1. Energetics (Thermodynamics and Equilibrium Properties)

Simulation of the free energies of different biological processes provide the “missing link” between the corresponding structures and functions. This is particularly true with regard to electrostatic free energies. However, obtaining reliable free energies by computer simulations is very challenging and requires validation by discriminative benchmarks. One of the most effective test in of calculations of free energies in proteins in general and electrostatic models in particular is provided by calculations of absolute **pK<sub>a</sub>s of ionizable groups** in proteins (8). The basis for calculating this quantity is the relationship (6)

$$2.3RT(pK_a^p - pK_a^w) = \Delta\Delta G_{sol}^p(AH \rightarrow A^-) - \Delta\Delta G_{sol}^w(AH \rightarrow A^-) \quad (1)$$

where the superscripts *p* and *w* denote protein and water, respectively, and  $\Delta\Delta G_{sol}(AH \rightarrow A^-)$  is the difference in solvation energy between the protonated and deprotonated states of the acid. Although early models only considered the change of pK<sub>a</sub>s due to the effect of protein ionized groups (9), it is now recognized that the most important factor is the self-energy (6) (or “solvation energy”), which is the energy of forming the charge in its protein environment when the other charges are turned off. The modulation of this effect is largely determined by the protein permanent dipoles (6). Obtaining reliable estimates of the self-energy by Free Energy Perturbation is very challenging, but significant progress has been made (3, 10). Apparently, it is essential to treat long-range effects in a consistent way in order to obtain meaningful results in pK<sub>a</sub> calculations (10). It is frequently simpler to obtain better results by simplified solvent models (6) and by models that treat the solvent as a dielectric continuum (11). However, one must bear in mind that models that treat the protein in a uniform dielectric without considering the microscopic nature of the protein permanent dipoles cannot reproduce the correct physics of ionized groups in proteins (6, 8).

**Redox energies** of biological cofactors play a major role in **electron transfer** processes. Such energies can be evaluated by simulation methods using a similar approach to that used in equation 1 and evaluating the change in solvation energy upon moving the charge from water to the protein. FEP calculations of redox energies were reported by several workers, and related studies in

[photosynthetic](#) reaction centers were also presented. Here the proper treatment of the solvent around the protein is quite crucial ([12](#)). As in the case of  $pK_a$  calculations, one finds that simplified models ([13](#)) and continuum models ([14](#), [15](#)) are quite effective.

Reliable calculations of free energies of **ligand binding** are crucial for quantitative progress in rational drug design. FEP has been used extensively in studies of binding free energies ([3](#)). Most early studies involve simple cases of small changes in the ligand ([16](#)). Calculations of absolute binding free energies are much more challenging ([3](#), [17](#)) and involve major convergence problems. Here the use of the LRA has been found to be quite effective ([17-19](#)). Extensive progress has also been made with simplified models ([17](#), [19](#), [20](#)) and with more phenomenological approaches ([21](#), [22](#)).

The permeation of ions through membrane channels plays a major role in biophysics. Here, the challenge is in evaluating the energy of an ion in the channel relative to its energy in solution and using the energetics in determining the conductance of the channel. Reliable simulation of such systems requires one to properly represent the channel and its surrounding membrane and water environment. The energy contribution of the membrane-induced dipoles must be included in consistent calculations. Consistent FEP studies of the energetics of moving an  $Na^+$  ion to the center of the **gramicidin A** channel have been reported with a reasonable result ([23](#)). Incomplete modeling of the boundaries of the simulated system and long-range interactions can lead to significant overestimates of the penetration barrier (eg, the results of ref. [24](#)). Semimacroscopic studies of ion channels are quite useful and much less challenging than full microscopic calculations ([10](#), [25](#), [26](#)). The recent emergence of the structure of a bacterial **potassium channel** ([27](#)) have highlighted the need for reliable FEP calculations, and some progress has been already made ([28](#)).

## 2. Structural Properties

Molecular simulations can be used to refine, and sometimes predict, structural properties. Energy minimization approaches were the first to be used successfully in relaxing bad steric contacts in [X-ray crystallography](#) ([29](#)). Such approaches were also used in more systematic refinements of protein structures ([30](#)). The emergence of powerful computers allows one to use more rigorous statistical mechanical approaches coupled with MD simulations in exploring structural properties ([31](#)).

Simulations are also very useful in evaluating and analyzing [temperature factors](#) and interpreting the meaning of the corresponding experimental results.

All-atom simulation approaches are still quite ineffective in predicting protein-structures, and one has to resort to simplified models. Even when one starts from the X-ray structure, it is not trivial to have the simulated structure identical to the observed one. The problems involved include the need for reliable potential function and, perhaps more importantly, the need to have a realistic description of the solvent ([32](#)) and proper treatment of long-range forces.

## 3. Dielectric Properties

Although proteins cannot and should not be described as a medium of a uniform dielectric constant, it is important to understand the dielectric properties of proteins. Experimental attempts to determine local dielectric constants are far from being unique, and simulation studies can provide great insight. In relating the simulated dielectric constant to experimental observations, one should take into account the effect of the solvent around the protein, which is also referred to as the *reaction field* ([33](#)). Consistent simulations of the overall dielectric constant of solvated proteins have been reported by several workers ([34-36](#)). It has been found that the local dielectric at protein [active sites](#) is quite large and quite different than the low value usually deduced from dielectric measurements ([37](#)). Simulations are also very effective in analyzing the fast time components of dielectric **relaxation times** ([34](#), [36](#)).

#### 4. Dynamical Properties

MD simulations and related approaches provide powerful ways of evaluating dynamical properties. Many of these properties can be casted in terms of the **corresponding time correlation function** (38).

$$C_A(t) = \frac{1}{\tau} \int_0^\tau A(t')A(t+t')dt' \quad (2)$$

where A is the property of interest and t is a sufficiently long simulation time. For example, the [diffusion](#) constant, D, can be obtained from the velocity autocorrelation,  $C_v(t)$ , by

$$D = \frac{1}{3} \int_0^\infty C_v(t)dt \quad (3)$$

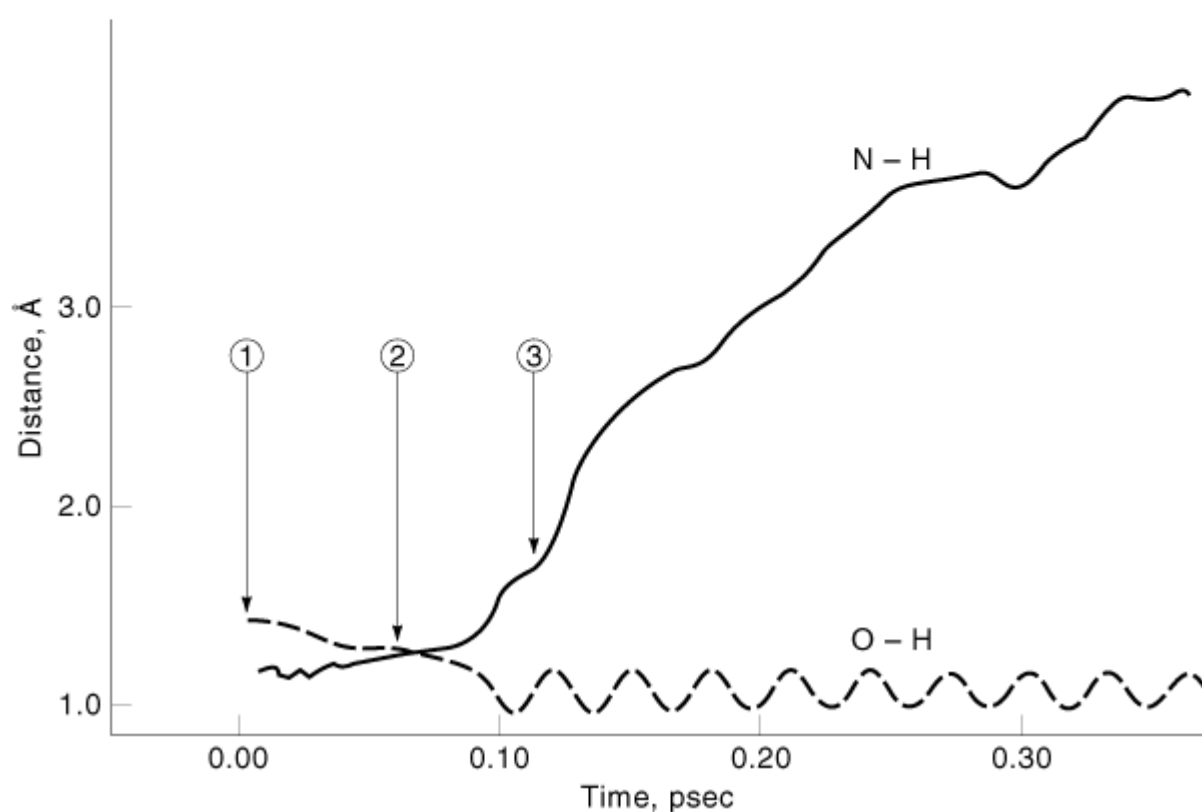
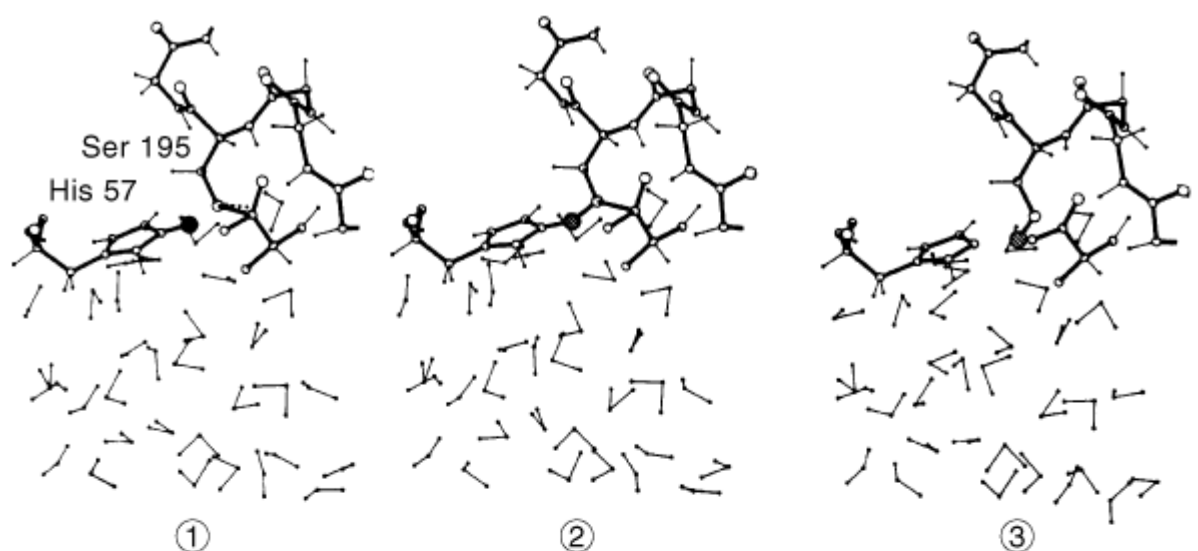
This relationship can be used in comparing simulations to experiments and in analyzing the microscopic meaning of experimentally deduced parameters. The dynamical range accessible to simulation studies can be probed by several approaches, and significant progress has been made in comparing simulations to the results of inelastic **neutron scattering** (39).

Molecular fluctuations are characterized by correlation times, which are closely related to the underlined relaxation processes. Relevant relaxation times can be evaluated by simulation studies using:

$$\tau_A^{-1} = C_A(0)^{-1} \int C_A(t)dt \quad (4)$$

The fluctuations of the macroenvironments of proteins can play an important role in fast reactions. This is true in particular with regard to the effect of electrostatic fluctuations on charge transfer processes such as electron transfer reactions. The rate of such reactions is determined by the effect of the protein fluctuation on the difference between the energies of the reactant and product state. Such an approach has been used in studies of the dynamics of the primary event in [photosynthesis 26, 40-42](#). The same approach can be used in studies of photoisomerization processes, although in such cases one needs to account for a large coupling between the electronic states. As much as rate constants of regular chemical reactions are concerned, the most important factors are [activation energies](#) (see below) and not dynamical effects. Nevertheless, the chance that a trajectory that reaches the [transition state](#) will lead to reaction of the transmission factor is a dynamical effect. This factor can be evaluated using MD simulation by running trajectories downhill from the transition state (43) (Fig. 1). However, the transition factor is usually close to unity in reactions with significant activation barrier (4, 26) (Fig. 1).

**Figure 1.** A downhill trajectory for the proton transfer step in the catalytic reaction of trypsin. The trajectory moves on the actual ground state potential, from the top of the barrier to the relaxed enzyme–substrate complex. 1, 2, and 3 designate different points along the trajectory, whose respective configurations are depicted in the upper part of the figure. The time reversal of this trajectory corresponds to a very rare fluctuation that leads to a proton transfer from Ser 195 to His 57. (adapted from ref. 4).



## 5. Rate of Biological Processes

The rate constant of processes in condensed phases can be expressed as:

$$k = F(k_B T/h) \exp\{-\Delta g^\ddagger/k_B T\} \quad (5)$$

where  $k_B$  is the Boltzmann constant,  $T$  is the absolute temperature,  $h$  is the Planck constant,  $F$  is the **transmission factor** of section 3, and  $Dg^\ddagger$  is the [activation-free energy](#). For most reactions with large  $Dg^\ddagger$ , one finds that  $F$  is close to unity and, thus, the most important factor is  $Dg^\ddagger$ . [Enzymes](#) catalyze their reactions by reducing activation-free energies (44), and one of the most important questions in structure function correlation is how this reduction of  $Dg^\ddagger$  is accomplished (4).

Molecular simulations should allow one to evaluate activation free energies and transmission factors and to probe the molecular origin of enzymatic reactions. However, such reactions involve bond making and bond breaking processes that cannot be described by simple force fields. It is also impractical to describe the complete enzyme/substrate complex by high level quantum mechanical approaches. A possible alternative is provided by **hybrid Quantum mechanics/molecular mechanics (hybrid QM/MM)** methods [4](#), [45-49](#). Such methods describe only parts of the reacting system quantum mechanically, while treating the rest of the protein classically. QM/MM approaches are very promising, but the corresponding results are not yet quantitative ([26](#)). An alternative, which is, at present, the most reliable way of simulating enzymatic reaction, is provided by the Empirical Valence Bond (EVB) approach ([4](#), [26](#), [50](#), [51](#)). The EVB method describes the reaction potential surface by a valence bond method that considers the different resonance structures of the reactant and product fragments and their interaction with the surrounding protein and solvent. This method is calibrated by considering the reference reaction in water so that it focuses on the difference between the reaction in protein and solution, which is much easier to obtain than to predict the reaction potential surface in the protein by a first principle approach.

The EVB approach can be used to calculate conveniently and efficiently activation-free energies of enzymatic reactions using a combination of FEP and umbrella sampling approaches ([4](#), [50](#)). The activation energies can be converted to rate constants of enzymatic reactions using equation [5](#).

## 6. Protein Folding

Models that use **simplified representation** of the protein amino acids ([52](#)) have provided a powerful understanding of the protein folding process ([53-56](#)). Such models are sometimes classified as *on-lattice* and *off-lattice* models. On-lattice models consider the simplified representation along a regular cubic lattice (thus allowing exact enumeration at the expense of somewhat unrealistic structure), while the off-lattice models try to provide more realistic representation ([52](#)). Simplified models are used extensively in understanding the relationship between sequence to folding probability and to the specific folded structure ([52-56](#)) and in deducing the sequence of events in the folding process ([53](#)). The gradual increase in computer power makes it tempting to use all-atom simulations in studies of folding processes. Some advances have been made ([57-59](#)), and more will follow. A promising effective approach is the use of the simplified representation as a reference potential for all-atom calculations ([60](#)).

## Bibliography

1. S. Lifson and A. Warshel (1968) *J. Chem. Phys.* **49**, 5116.
2. U. Burkert and N. L. Allinger (1982) *Molecular Mechanics*, American Chemical Society, Washington, DC.
3. P. Kollman (1993) *Chem. Rev.* **93**, 2395–2417.
4. A. Warshel (1991) *Computer Modeling of Chemical Reactions in Enzymes and Solutions* John Wiley & Sons, New York.
5. A. T. Hagler, E. Huler, and S. Lifson (1977) *J. Am. Chem. Soc.* **96**, 5319.
6. A. Warshel and J. Åqvist (1991) *Ann. Rev. Biophys. Chem.* **20**, 267–298.
7. P. E. Smith and B. M. Pettitt (1994) *J. Phys. Chem.* **98**, 9700–9711.
8. C. N. Schutz and A. Warshel (2001) *Proteins: Struct., Func., and Gen.* **44**, 400–417.
9. C. Tanford and J. G. Kirkwood (1957) *J. Am. Chem. Soc.* **79**, 5333.
10. F. S. Lee, Z. T. Chu, and A. Warshel (1993) *J. Comp. Chem.* **14**, 161–185.
11. B. Honig, K. Sharp, R. Sampogna, M. R. Gunner, and A. S. Yang (1993) *Proteins: Struct., Func., and Gen.* **15**, 252–265.
12. R. G. Alden, W. W. Parson, Z. T. Chu, and A. J. Warshel (1995) *J. Am. Chem. Soc.* **117**, 12284–12298.
13. P. J. Stephens, D. R. Jollie, and A. Warshel (1996) *Chem. Rev.* **96**, 2491.

14. M. R. Gunner and B. Honig (1991) *Proc. Natl. Acad. Sci. U.S.A.* **88**, 9151–9155.
15. H.; X. Zhou (1994) *J. Am. Chem. Soc.* **116**, 10362–10375.
16. T. P. Straatsma and J. A. McCammon (1992) *Ann. Rev. Phys. Chem.* **43**, 407–435.
17. F. S. Lee, Z. T. Chu, M. B. Bolger, and A. Warshel (1992) *Prot. Eng.* **5**, 215–228.
18. J. Åqvist, C. Medina, and J.–E. Samuelsson. (1994) *Prot. Eng.* **7**, 385–391.
19. Y. Y. Sham, Z. T. Chu, H. Tao, and A. Warshel (2000) *Proteins: Struct. Funct. Genet.* **39**, 393–407.
20. N. Froloff, A. Windemuth, and B. Honig (1997) *Protein Sci.* **6**, 1293–1301.
21. G. Verkhivker, K. Appelt, S. T. Freer, and J. E. Villarfranca (1995) *Prot. Eng.* **8**, 677–691.
22. I. Muegge and Y. C. Martin (1999) *J. Med. Chem.* **42**, 791–804.
23. J. Åqvist and A. Warshel (1989) *Comments Mol. Cell Biophys.* **6**, 91.
24. B. Roux and M. J. Karplus (1993) *Am. Chem. Soc.* **115**, 3250–3260.
25. V. Dorman, M. B. Partenskii, and P. C. Jordan (1996) *Biophys. J.* **70**, 121–134.
26. A. Warshel and W. W. Parson *Quart. Rev. Biophys.* In Press.
27. D. A. Doyle, J. M. Cabral, R. A. Pfuetzner, A. L. Kuo, J. M. Gulbis, S. L. Cohen, B. T. Chait, and R. MacKinnon (1998) *Science* **280**, 69–77.
28. V. Luzhkov and J. Aqvist (2001) *Biochimica et Biophysica Acta* **36446**, 1–9.
29. M. Levitt and S. J. Lifson (1969) *J. Mol. Biol.* **46**, 269–279.
30. A. Jack and M. Levitt (1978) *Acta Cryst. A* **34**, 931.
31. A. T. Brunger (1988) *XPLORE Manual*, 1.5 ed., Yale University, New Haven, CT.
32. M. Levitt and R. Sharon (1988) *Proc. Natl. Acad. Sci. U.S.A.* **85**, 7557–7561.
33. L. J. Onsager (1936) *Am. Chem. Soc.* **58**, 1486.
34. G. King, F. S. Lee, and A. J. Warshel (1991) *Chem. Phys.* **95**, 4366–4377.
35. T. Simonson and D. J. Perahia (1995) *Am. Chem. Soc.* **117**, 7987.
36. W. F. vanGunsteren and A. E. Mark (1992) *Eur. J. Biochem.* **204**, 947.
37. A. Warshel and A. Papazyan (1998) *Curr. Opin. Struct. Biol.* **8**, 211–217.
38. D. A. McQuarrie (1976) *Statistical Mechanics*, Harper and Row, New York.
39. S. Cusack (1989) *Chemica Scripta* **29A**, 103–107.
40. A. Warshel and W. W. Parson (1991) *Annu. Rev. Phys. Chem.* **42**, 279–309.
41. S. Creighton, J.–K. Hwang, A. Warshel, W. W. Parson, and J. Norris (1988) *Biochemistry* **27**, 774–781.
42. K. Schulten and M. Tesch (1991) *Chem. Phys.* **158**, 421–446.
43. C. H. Bennet (1997) *Algorithms for Chemical Computations*, ACS, Washington, DC.
44. L. Pauling (1946) *Chem. Eng. News* **263**, 294.
45. A. Warshel and M. Levitt (1976) *J. Mol. Biol.* **103**, 227–249.
46. V. Théry, D. Rinaldi, J.-L. Rivail, B. Maignret, and G. G. Ferenczy (1994) *J. Comp. Chem.* **15**, 269–282.
47. J. Gao (1995) *Reviews in Computational Chemistry*, Vol. **7**, VCH, New York.
48. A. J. Mulholland, G. H. Grant, and W. G. Richards (1994) *Prot. Eng.* **6**, 133–147.
49. R. B. Murphy, D. M. Philipp, and R. Friesner (2000) *J. Comput. Chem.* **21**, 1442–1457.
50. J. Åqvist and A. Warshel (1993) *A. Chem. Rev.* **93**, 2523–2544.
51. J. Lobaugh and G. A. Voth (1996) *J. Chem. Phys.* **104**, 2056–2069.
52. M. Levitt and A. Warshel (1975) *Nature* **253**, 694.
53. J. N. Onuchic, P. G. Wolynes, Z. Luthey-Schulten, and N.D. Socci (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 3626–3630.



54. A. Godzik, A. Kolinski, and J. J. Skolnick (1993) *Comp.-Aided Mol. Des.* **7**, 397–438.
55. K. Yue and K. A. Dill (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 4163–4167.
56. E. Shakhnovich, G. Farztdinov, A. M. Gutin, and M. Karplus (1991) *Phys. Rev. Lett.* **67**, 1665–1668.
57. E. M. Boczko and C. L. Brooks, III (1995) *Science* **269**, 393–396.
58. V. Daggett and M. J. Levitt (1993) *Mol. Biol.* **232**, 600–619.
59. T. Lazaridis and M. Karplus (1997) *Science* **278**, 1928–1931.
60. Z. Z. Fan, J. K. Hwang, and A. Warshel (1999) *Theor. Chem. Accounts* **103**, 77–80.

### Suggestions for Further Reading

61. M. P. Allen and D. J. Tildesley (1987) *Computer Simulation of Liquids*, Oxford University Press, Oxford.
62. A. Warshel (1991) *Computer Modeling of Chemical Reactions in Enzymes and Solutions*, John Wiley & Sons, New York.

## Concatemers

A concatemer is a multimeric molecule of **DNA** formed by identical monomers arranged linearly in the same head-to-tail orientation. Closure of a **bacteriophage** DNA molecule upon entry into the bacterial host cell is a relatively simple process. The complementary single-stranded ends anneal in a reaction that is thermodynamically favored and by conditions that restrict [diffusion](#) and keep the ends close to each other. **Annealing** is followed by covalent joining, catalyzed by the host [DNA Ligase](#), which forms phosphodiester bonds. Opening of the circular DNA molecule, coupled to its encapsidation to form mature phage particles, however, is rather complex.

**Capsid** precursors interact with the newly replicated DNA molecules in the form of concatemers ([1](#), [2](#)), which are the product of DNA replication proceeding in two stages. First, the DNA injected into the host is circularized and covalently closed, as previously described. In the case of the [lambda phage](#), the DNA undergoes several stages of DNA replication in which the circular DNA rings generate daughter rings. Then, during phage vegetative growth, **rolling-circle** replication gives rise to *linear* concatemers of 1 DNA. These concatemers, which in the absence of packaging may be up to 10 chromosomal units long, are the normal substrates for packaging.

RNA primers for DNA synthesis are excised during the process. This raises the question of how the extreme ends of linear DNA can be completed. Not only must a small RNA fragment initiate at the 3' terminal nucleotide of the template strand, but the sequence must be filled in after ribonucleotide removal. This dilemma was resolved using phage T7 DNA, which has the property of having redundant ends. The sequence of about 160 bp found at the left end is repeated exactly at the right end of the molecule. When the linear T7 molecules replicate, they do not generate unit-length progeny molecules, but very long concatemers (units linked end to end) containing the complete genome sequence repeated over and over.

### Bibliography

1. D. Kaiser and T. Masuda (1973) *Proc. Natl. Acad. Sci. USA* **70**, 260–264.
2. D. Kaiser, M. Syvanen, and T. Masuda (1974) *J. Supramol. Struct.* **2**, 318–328.

## Conditional Lethal Mutations

[Mutations](#) that cause lethality under one condition (the restrictive or [nonpermissive condition](#)) but not another (the [permissive condition](#)) are called conditional lethal mutations, or loosely, conditional lethals. Conditional lethal mutations have an honored position in molecular biology. The recognition that a **gene** is not indivisible, the determination of the triplet nature of the [genetic code](#), and the discovery of **nonsense codons**, all resulted from studies of a particular class of conditional lethal mutants of **bacteriophage T4**. T4 *rII* mutants grow on *Escherichia coli* strain B (the permissive host) but not on *E. coli* strain K-12 when it is **lysogenic** for bacteriophage **lambda** (the nonpermissive host). Interestingly, the function of the proteins encoded by the two *rII* genes remains unknown. The conditional lethality of the *rII* mutants allowed geneticists to grow the phage, to perform genetic crosses on the permissive host, and then to analyze the results on the nonpermissive host.

Conditional lethal mutations are the only genetic way to identify essential genes in [haploid](#) organisms or cells. By definition, a mutant with a [null mutation](#) in an essential gene cannot be isolated in the haploid state (in **diploid** organisms, [recessive lethal mutations](#) on [autosomes](#) are **complemented** by the other functional **allele**). Conditional lethal mutations are usually identified by screening mutant clones for growth under permissive conditions but not under restrictive conditions. Lethal in this case does not always mean conveying death. For example, [auxotrophies](#) are considered a class of conditional lethals because the auxotrophic cells grow only in the presence of the required growth factor, although they do not necessarily die in its absence. The more typical classes of conditional lethals are host-range mutants (described above), **temperature-sensitive mutants**, and [cold-sensitive mutants](#). [Nonsense mutations](#) are considered conditionally lethal because the mutation is suppressed when a suppressor [transfer RNA](#) is present.

Conditional lethal mutations are a powerful way to identify [protein–protein interactions](#). A mutation that renders a protein inactive at high or low temperature, for example, may be suppressed by a compensating mutation in the gene encoding a protein that interacts with it. To conclude that the two proteins interact, the suppression must be allele-specific, ie, not caused by general suppression (as in the case of tRNA nonsense suppressors). In addition, other events, such as [gene duplication](#), suppress conditional lethals. Interacting proteins are also identified by synthetic lethal mutations, which are mutations in each of two genes, neither of which is lethal alone but which are lethal when together in the same cell.

### Suggestions for Further Reading

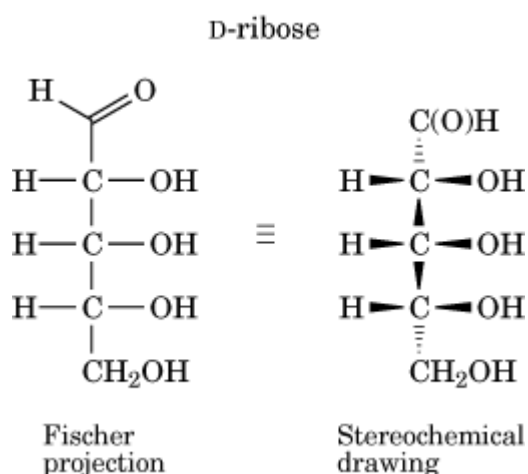
- A. Adams, D. E. Gottschling, C. Kaiser, and T. Sterns (1998) *Methods in Yeast Genetics: A Laboratory Course Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
- E. M. Phizicky and S. Fields (1995) Protein-protein interactions: methods for detection and analysis, *Microbiol. Rev.* **59**, 94–123.
- A. Griffiths, J. H. Miller, D. Suzuki, R. Lewontin, and W. Gelbart (1996) *An Introduction to Genetic Analysis*, 6th ed., W. H. Freeman, New York.
- J. Beckwith and T. Silhavy (1992) *The Power of Bacterial Genetics: A Literature-Based Course*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
- T. D. Brock (1990) *The Emergence of Bacterial Genetics*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

## Configuration

Although configuration and [conformation](#) are listed as synonyms in the *Oxford English Dictionary* (1), in chemistry and biochemistry they have specific and separate meanings differentiated by time or energy. Both terms refer to the organization of atoms in space, but configuration refers to the element that remains invariant in time in the absence of any covalent bonds being altered, whereas conformation refers to the relative orientation in space of atoms that can vary by rotation about single bonds and consequently varies in real time for biological molecules. The absolute configuration identifies the arrangement of atoms that generates **chirality** in a molecule. Following Pasteur's identification of the two chiral forms of tartaric acid, it was only known that they differed by being nonsuperimposable mirror images. This ambiguity was resolved when the crystal structure of (+) tartaric acid was determined by Bijvoet and co-workers (2). By this means, the absolute configuration of many related molecules became identifiable.

The designation of the configuration of a chiral center was problematic for nearly a century. Fischer introduced a general procedure that identified enantiomers as either D- or L- (3) based on whether the nonhydrogen substituent was on the right or left when the molecule was drawn as a "Fischer projection" (Fig. 1). This nomenclature permeates biochemistry through the now colloquial names of [amino acids](#), carbohydrates, and [lipids](#).

**Figure 1.** A Fischer projection. By convention, it has the vertical bonds directed away from the viewer and horizontal bonds directed out toward the viewer. The configuration of each chiral carbon in D-ribose is D because the nonhydrogen substituent is drawn to the right. For sugars, the enantiomer is defined by the bottom chiral carbon when the carbon chain is oriented vertically with the carbonyl carbon at the top; thus, D-ribose is pictured.

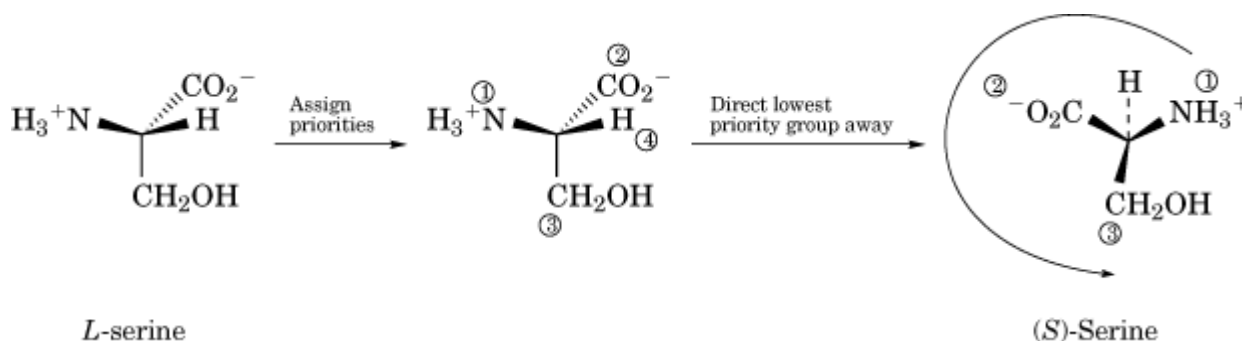


The Fischer nomenclature leads to ambiguities (4). A rigorous and unambiguous method of identifying configuration proposed by Cahn et al. has been adopted for specifying absolute configuration as either (*R*) or (*S*) for chiral tetrahedral centers (5). The procedure requires the assignment of priority to the four substituents that generate a chiral center, followed by a procedure to identify the arrangement as either (*R*) or (*S*). The four rules for assigning priority are:

1. Substituents are assigned their priority in order of decreasing atomic number of the atom directly bonded to the chiral center.
2. If two or more atoms receive the same priority in step 1, the atoms bonded to each of the equal priority atoms are examined one at a time. If the two groups are not differentiated by the atom of greatest priority, the second and third atoms are compared successively.
3. Heavier isotopes take precedence over lighter isotopes, eg,  $^2\text{H}$  over  $^1\text{H}$  and  $^{14}\text{C}$  over  $^{12}\text{C}$ .
4. Double bonds are counted as two bonds to the same atom.

Once the relative priorities of the four substituents are assigned, the bond to the lowest priority group is oriented directly away from the observer, and the remaining three groups are viewed in a plane. If the arc connecting the three groups in order of highest to lowest priority is clockwise, the chiral center is assigned the (*R*) configuration (from the Latin *rectus*), whereas if the arc is counterclockwise, the center is assigned the (*S*) configuration (from the Latin *sinister*). These assignments are shown in Figure 2.

**Figure 2.** Assignment of the relative priorities of the four substituents on the chiral  $\alpha$ -carbon of serine. The carboxylate carbon takes precedence over the hydroxymethyl carbon because it has three bonds to oxygen. The assignment of the (*S*) configuration results from the counterclockwise arc that connects the substituents in order of precedence.



. Assignment of the relative priorities of the four substituents on the chiral  $\alpha$ -carbon of serine. The carboxylate carbon takes precedence over the hydroxymethyl carbon because it has three bonds to oxygen. The assignment of the (*S*) configuration results from the counterclockwise arc that connects the substituents in order of precedence. [[Full View](#)]

## Bibliography

1. *Oxford English Dictionary*, 2nd ed. (1989) Clarendon Press, Oxford.
2. J. M. Bijvoet, A. F. Peerdeman, and A. J. van Bommel (1951) *Nature* **168**, 271.
3. E. Fischer (1891) *Ber. Dtsch. Chem. Ges.* **24**, 2683.
4. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York, Chap. "5".
5. R. S. Cahn, C. K. Ingold, and V. Prelog (1966) *Angew. Chem. Int. Ed.* **5**, 385–415.

## Suggestions for Further Reading

6. J. March (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, pp. 93–99.
7. K. Mislow (1966) *Introduction to Stereochemistry*, W. A. Benjamin, New York, pp. 86–97.

## Confocal Microscopy

Confocal microscopy is now a well-established tool for the examination of subcellular structure and function and complements light and electron microscopy. Because of its high temporal resolution, confocal microscopy allows for the visualization of living as well as fixed tissues and cells, and therefore dynamic processes can be examined and analyzed quantitatively as they actually occur. Confocal microscopy offers several advantages over conventional light microscopy, among which are an increase in contrast, resolution, and clarity (1). Conventional light microscopy provides a two-dimensional (2-D) image of the specimen in the focal plane of the objective lens, but this image is contaminated by out-of-focus images of the specimen above and below the focal plane. Confocal microscopy provides a 2-D image in the focal plane without the out-of-focus information. Furthermore, the resolution of images from a confocal microscope is improved by a factor of 1.4–1.75 (2). Through computer control of the focus and acquisition of images, modern confocal light microscopes can collect a series of 2-D images (or “optical sections”) through the specimen producing a three-dimensional (3-D) image. Four-dimensional imaging (4-D), defined as 3-D imaging over time, is a recently developed extension in which 3-D images are recorded at periodic time intervals (3).

The basic principle of the confocal microscope is to eliminate the scattered, reflected, or fluorescent light from out-of-focus planes by making the illumination, specimen, and detector all have the same focus, ie, they are confocal. In effect, this microscope will image only the very thin optical section on which the beam is focused. Matched pinholes are used, one at the light source which is imaged onto the specimen to function as a probe that is scanned over the specimen, and one at the detector to capture only a narrow plane of focus. Thus, the out-of-focus blur from areas above and below the focal plane are eliminated. The matched pinhole apertures improve the lateral resolution over conventional light microscopes by a factor of 1.4 with the use of circular apertures and 1.75 with annular apertures (2). Confocal microscopes are often designed to scan in a raster pattern over the sample in which the microscope illuminates one spot at a time, scanning the spot along parallel lines in the focal plane. Lasers are an ideal illumination source for raster scanning because they provide an intense beam of monochromatic radiation that can be condensed onto a small spot. Hence many confocal microscopes are laser scanning (LSCM).

Applications of LSCM include (i) determining the location of [proteins](#), [lipids](#) and nucleic acids, cytoskeletal structures and organelles within cells (4, 5), using fluorescent dyes, antibodies, [phalloidin](#), and [lectins](#), (ii) observing ionic fluctuations, such as calcium, magnesium, and pH, in cells and organelles (6), and, (iii) measuring [membrane potential](#) using fluorescent dyes (7). The power of 4-D imaging in molecular biology was illustrated by the noninvasive monitoring by LSCM of mitotic events, and cleavage and migration patterns of fertilized sea urchin eggs labeled with DiOC<sub>6</sub> (8). [Macromolecules](#) and subunits can be characterized by immunocytochemical fluorescence probes, which involves the use of antibodies labeled with fluorophores (9). When more than one fluorophore is used, 3-D multilabel (multicolor) imaging can be used to map the 3-D relationship of the labeled structures. For example, z-series collected from two different channels, fluorescein and rhodamine, can be merged into a single reconstruction. By rotating the rendered volume, particular nuances of the structural relationships highlighted by the bound fluorophores may be revealed (10, 11). It is not only important to know that a particular compound is present in a specific cell type or subcellular component, but it is also important to detect and quantify changes in local concentrations of such compounds. Quantitative immunocytochemistry utilizes the principles of stereology and statistical analyses (12) and can be coupled to LSCM to detect differences in immunoreactivity. This method has been used to measure the distribution of integral [membrane proteins](#) in the vertebrate retina and to determine the distribution of transport vesicles (13).

Many of the advantages of a confocal light microscope can be achieved with a standard light microscope equipped with a digital camera interfaced to a computer, which controls the microscope stage and focus control as well as run rapid deconvolution algorithms (14). With knowledge of the 3-D point spread function of the microscope, the deconvolution algorithms remove out-of-focus signal to produce images comparable in quality to LSCM images. The advantages of the deconvolution confocal microscope are that (i) the system expense is a fraction of the cost for a LSCM system, and (ii) lower illumination levels can be used, since all of the light emitted from the specimen is used in forming the 3-D image. Lower illumination minimizes the problem of photobleaching of fluorophores used to label the specimen. A disadvantage of this approach is that the full 3-D deconvolution process can take significantly more time to produce a 3-D image than does a confocal microscope.

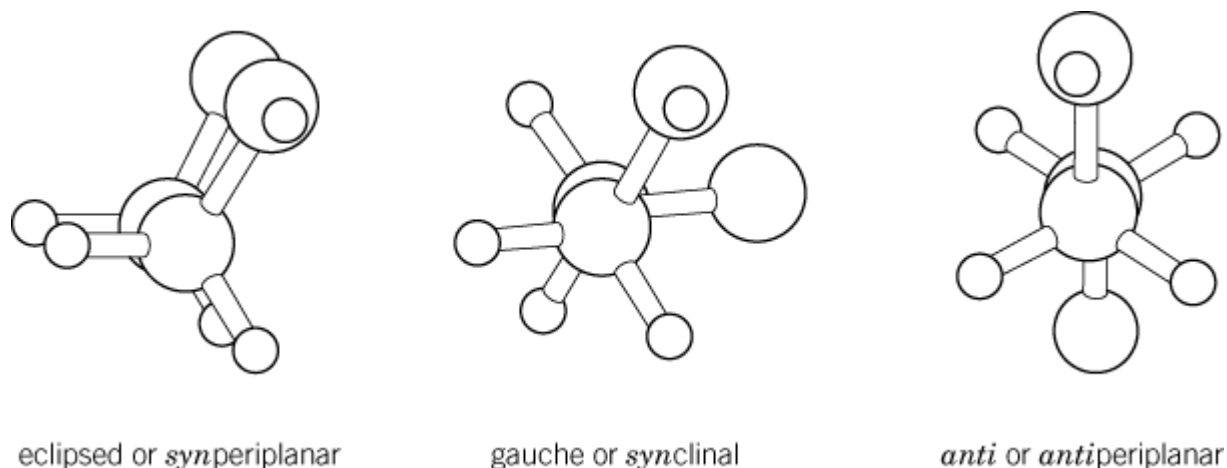
## Bibliography

1. A. Boyde (1990) "Confocal optical microscopy", In *Modern Microscopies* (P. J. Duke and A. G. Michette, eds.), Plenum Press, New York, pp. 185–204.
2. E. M. Slayter and H. S. Slayter (1992) *Light and Electron Microscopy*, Cambridge University Press, Cambridge, UK.
3. S. A. Stricker, S. Paddock, and G. Schatten (1990) *J. Cell Biol.* **111**, 113a.
4. A. H. Cornell-Bell et al. (1993) "Membrane glycolipid trafficking in living polarized pancreatic acinar cells: Assessment by confocal microscopy", In *Methods in Cell Biology*, Vol. **38**: *Cell Biological Applications of Confocal Microscopy* (B. Matsumoto, ed.), Academic Press, San Diego, pp. 222–241.
5. I. L. Hale and B. Matsumoto (1993) "Resolution of subcellular detail in thick tissue sections: Immunohistochemical preparation and fluorescence confocal microscopy", In *Methods in Cell Biology*, Vol. **38**: *Cell Biological Applications of Confocal Microscopy* (B. Matsumoto, ed.), Academic Press, San Diego, pp. 290–325.
6. A. Boyde (1995) "Confocal optical microscopy", In *Image Analysis in Histology: Conventional and Confocal Microscopy* (R. Wootton et al., eds.), Cambridge University Press, Cambridge, UK, pp. 151–196.
7. L. M. Loew (1993) "Confocal microscopy of potentiometric fluorescent dyes", In *Methods in Cell Biology*, Vol. **38**: *Cell Biological Applications of Confocal Microscopy* (B. Matsumoto, ed.), Academic Press, San Diego, pp. 195–210.
8. S. A. Stricker et al. (1992) *Dev. Biol.* **149**, 370–380.
9. T. C. Brelje, M. W. Wessendorf, and R. L. Sorenson (1993) "multicolor laser scanning confocal immunofluorescence microscopy: Practical application and limitations", In *Methods in Cell Biology*, Vol. **38**: *Cell Biological Applications of Confocal Microscopy* (B. Matsumoto, ed.), Academic Press, San Diego, pp. 98–182.
10. W. Galbraith et al. (1989) *Soc. Photo-Opt. Instrum. Eng.* **1063**, 19–20.
11. M. W. Wessendorf (1990) In *Handbook of Chemical Neuroanatomy*, Vol. **8** (A. Bjorklund et al., eds.), Elsevier, Amsterdam, pp. 1–45.
12. J. T. McBride (1995) "Quantitative immunocytochemistry", In *Image Analysis in Histology: Conventional and Confocal Microscopy* (R. Wootton et al., eds.), Cambridge University Press, Cambridge, UK, pp. 339–354.
13. B. Matsumoto and I. L. Hale (1993) In *Methods in Neuroscience*, Vol. **15** (P. C. Hargrave, ed.), Academic Press, San Diego, pp. 54–71.
14. D. Agard (1984) *Ann. Rev. Biophys. Bioeng.* **13**, 191–219.

## Conformation

Nonidentical spatial arrangements of the atoms of a molecule achieved by rotation about single bonds are referred to as different conformations (1). Two molecules differing in their conformation may be referred to as conformers. Because rotations about single bonds are usually rapid, conformers are rapidly interconverted, making it difficult to separate individual conformers. Note one exception to this rule is rotation about the C–N single bond of the [peptide bond](#) of polypeptides, which has partial double-bond character, so it is planar and the *cis* and *trans* conformations are only slowly interconverted (see [Cis/Trans Isomerization](#)). The conformation of a molecule may then be largely specified by characterizing the rotations about its single bonds. These rotations are quantified by determining the [torsion angle](#) or [dihedral angle](#). Individual torsion angles may be qualitatively described as eclipsed, **gauche**, or **anti** (Fig. 1), or more formally named *syn* periplanar, clinal, and *anti* periplanar; they are quantified, as described in [Torsion angle](#).

**Figure 1.** The common description of a conformation, illustrated by ethylene glycol. When the substituents are superimposed, the conformation is described as *eclipsed*; when they are staggered, the conformation is *gauche* if the large substituents are roughly separated by 60° and *trans* or *anti* if they are opposite.



Other common combinations of torsion angles lead to descriptive names for common groups of atoms, eg, the 2'-endo conformation of **ribose** describes all five of the torsion angles of the furanose ring. This practice extends to larger biomacromolecules. The conformation of a small repeating segment of a biomacromolecule is characterized by its torsion angles along the backbone. Thus specification of the phi ( $\phi$ ) and psi ( $\psi$ ) angles for an oligopeptide (see [Ramachandran Plot](#)) describes its conformation as alpha-**helical** or beta-**sheet** (2). The specification of the torsion angles along the phosphodiester backbone of a polynucleotide also will determine whether the conformation is of the A, B, or Z type (3) (see [DNA Structure](#)).

A stable three-dimensional structure of a biological macromolecule is referred to as its conformation. Because rotation about single bonds is energetically easy, a unique conformation does not exist, as at least one single bond will be changing rapidly. A cooperative change of several torsion angles, such as occurs in the protein unfolding or melting of double-stranded DNA, may have a significant energetic barrier. The interconversion of the two conformations is then identified as a *conformational change*, where the implication is that the macromolecule has been converted from one family of conformations to an experimentally distinguishable family of conformations.

## Bibliography

1. E. L. Eliel et al. (1967) *Conformational Analysis*, Wiley-Interscience, New York.
2. G. E. Schulz and R. H. Schirmer (1979) "Principles of Protein Structure", *Springer Advanced Texts in Chemistry* (C. R. ed.), Springer-Verlag, New York.
3. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer Advanced Texts in Chemistry (C. R. Cantor, ed.), Springer-Verlag, New York.

## Suggestions for Further Reading

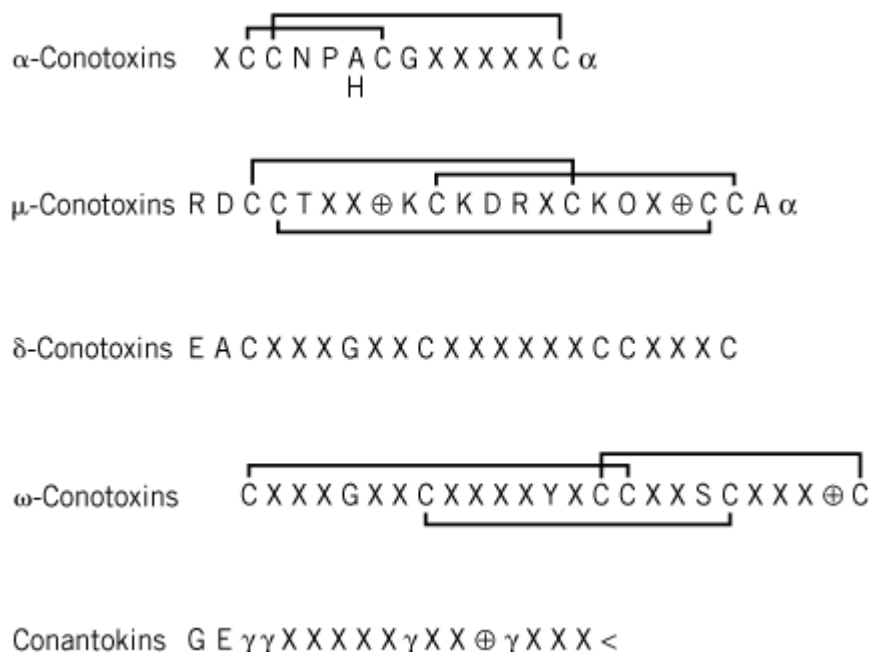
4. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry, Part I: The Conformation of Biological Macromolecules*, W. H. Freeman, San Francisco, CA.
5. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York, pp. 124–137.

## Conotoxins

Perhaps the largest array of **channel-binding toxins** is produced by marine gastropods of the genus *Conus*, which includes more than 500 species ([1](#), [2](#)). These animals appear to have developed the ability of generating a very large number of toxin variants starting from a basic peptide structure consisting of relatively few residues (Fig. [1](#)). In such a way, they produce venoms containing complex mixtures of toxins, capable of binding and blocking the activity of several ion channels. Moreover, some conotoxins are capable of binding differentially to the various **isoforms** of an ion-specific channel, thus providing precious tools to pharmacologists and neuroscientists ([1](#)).

**Figure 1.** Peptide toxins contained in *Conus* venoms. *Conus* spp. produce a highly complex mixture of toxins. Shown here are the consensus sequences and disulfide connections of toxins specific for the acetylcholine receptor (a-conotoxins), sodium channels (m- and d-conotoxins), calcium channels (w-conotoxins), and the NMDA glutamate channel (conantokins). O indicates *trans*-4-hydroxyproline, plus a positively charged residue, g a g-carboxyglutamate, and a amidated  $\alpha$ -carboxyl. The disulfide connectivity of d-conotoxins is not established.





The  $m$ -conotoxins are 22 residue-long basic peptides and include three hydroxyproline residues and six [cysteines](#) that form three [disulfide bonds](#) (Fig. 1). These conotoxins have an ellipsoid shape with the basic residues, which are essential for the binding activity, clustered on one side of the molecule (2, 3).  $m$ -conotoxins bind to the tetrodotoxin binding site of the **sodium channel**, thereby inhibiting the propagation of the action potential and causing flaccid paralysis. Another family of conotoxins specific for sodium channels is that of the  $d$ -conotoxins, whose binding causes an increased conductance of such voltage-gated channels (4).

$w$ -conotoxins act on voltage-gated **calcium channels** and inhibit calcium entry through the presynaptic membrane, thus preventing the release of neurotransmitters. Their peptide chain consists of 24 to 30 residues and three disulfide bonds (Fig. 1). A frequently used  $w$ -conotoxin is  $G_{VIA}$ , a 27 residue-long peptide that is very specific for calcium channels containing the  $\alpha 1B$  subunit, such as the N-type calcium channel present at the neuromuscular junction of vertebrates. The specificity of different  $w$ -conotoxins for different vertebrates is exploited in studying the function and anatomical distribution of calcium channels (1).

The largest family of conotoxins is that of the  $a$ -conotoxins, 13 to 18 residue-long peptides, whose folding is dominated by two disulfide bonds (Fig. 1). These paralytic neurotoxins bind specifically to the nicotinic [acetylcholine receptor](#) and different  $a$ -conotoxins are able to distinguish between different neuronal isoforms of the receptor. *Conantokins* are *Conus* peptides with yet another channel specificity: They bind to the NMDA-sensitive glutamate channels. They are peptides of 17 to 21 residues with no cysteine but four **gamma-carboxyglutamate** residues (1) (Fig. 1). Very short *Conus* peptides are conopressins, which consist of nine amino acid residues with a single disulfide bond and are specific for the vasopressin receptor.

## Bibliography

1. R. Rappuoli and C. Montecucco (1997) *Guidebook to Protein Toxins and Their Use in Cell Biology*, Sambrook and Tooze, Oxford University Press, Oxford, UK.
2. B. M. Olivera, G. Miljanich, J. Ramachandran, and M. E. Adams (1994) *Ann. Rev. Biochem.* **63**, 823–867.
3. J.-M. Lancelin et al. (1991) *Biochemistry* **30**, 6908–6916.
4. K. Shon (1995) *Biochemistry* **34**, 4913.

## Conservative Substitutions

[Missense mutations](#) cause the substitution of one [amino acid](#) for another in the [protein](#) encoded by the mutant **gene**. The effect of such a change on the protein's function will often depend on how similar the new amino acid is to the original one. Amino acids are classified by their acidity, polarity, [hydrophobicity](#), and structure. Substitutions are considered conservative if the change is between amino acids of the same class. However, a better criterion is whether the substitution has occurred during the [evolution](#) of related proteins. The mutations that give rise to conservative amino acid substitutions are considered evolutionarily [neutral mutations](#), or nearly so. Lists of likely conservative substitutions are found in many laboratory manuals, for example ([1](#)).

### Bibliography

1. A. Ellington and J. M. Cherry (1991) In *Current Protocols in Molecular Biology* (F. M. Ausubel, R. Brent, R. E. Kingston, D. D. Moore, J. G. Seidman, J. A. Smith, and K. Struhl, eds.), Wiley, New York, pp. A.1C.1–A.1C.10.

## Constant (C) Region

The heavy and light chains of [immunoglobulins](#) contain (a) **variable (V) regions** of approximately equal size (ie, about 110 amino acid residues), and (b) constant (C) regions that contain a multiple of 110 residues, one for the L chain and three or four for the H chains, depending upon the [isotype](#). Whereas the V regions are involved in [antigen](#) recognition, C regions are devoted to effector functions, such as **complement** fixation, transplacental passage, or attachment to various cell types. C regions are organized on a basic **domain** structure, corresponding to the 110-residue subunit of Ig with an autonomous three-dimensional [protein structure](#), as shown by [X-ray crystallography](#) (see [Immunoglobulin Structure](#)). This structure, known as the immunoglobulin constant domain fold, provides the basis for the general organization of a large number of proteins, known as the [Ig superfamily](#).

The existence of different C regions was first shown by immunochemistry, because specific antisera could distinguish two types of light chains, known as [L light chains](#), and five discrete isotypes of H chains, which dictate the organization into the five classes of immunoglobulins ([IgA](#), [IgD](#), [IgE](#), [IgG](#), [IgM](#)). Definitive understanding of the basis of isotype structure was obtained by [protein sequencing](#) and finally by [DNA sequencing](#). Within each Ig class, one may distinguish subclasses, as in the case of the g chain or the a chains in humans.

Most of the effector functions of immunoglobulins are supported by the last two or three COOH-terminal domains of the heavy chains, depending upon the class. The most interesting feature of constant regions is their ability to bind to a receptor family, termed **Fc receptors (FcRs)**, because binding involves the Fc region of Ig. FcRs represent adaptors that can confer on a variety of cell types any antibody specificity, depending on which Ig will bind to the receptor. There is an FcR

specific of each soluble isotype, IgG, IgA, and IgE. One distinguishes FcRs of high affinity (FcR type I) that fix Ig monomers and FcRs of low affinity (FcR type II) that bind Ig only when aggregated or as multivalent complexes. Depending on which cell type FcR is anchored, they express a cytoplasmic tail of variable size, which may result from alternative splicing and that contains various signaling motifs such as an immunoreceptor tyrosine-based activation motif and an immunoreceptor tyrosine-based inhibition motif. Fc $\epsilon$ RI, present on mast cells, basophils, and eosinophils, plays a major role in anaphylaxis, because it will fix monomeric IgE. Upon cross-linking by antigen, this FcR is triggered, leading to the cascade of events that are responsible for degranulation of the corresponding cells and thus to the liberation of pharmacologically active molecules that account for immediate hypersensitivity. Fc $\epsilon$ RII (CD23) binds aggregated IgE, or IgE already cross-linked by antigen. It is present at the surface of many cell types, including [B cells](#). It may also be cleaved and release several soluble products that have a **cytokine**-like activity. Fc $\gamma$ RI is found on monocytes, eosinophils, and neutrophils and binds to monomeric IgG, whereas Fc $\gamma$ RaI is present on macrophages, neutrophils, and T and B cells.

This demonstrates the extraordinary plasticity and amplification of molecules of the immune system and illustrates the huge diversity of functions that are linked to the constant regions of immunoglobulins.

See also entries [Antibody](#), [Immunoglobulin](#), and [Isotype](#).

#### Suggestions for Further Reading

F. Shakib (1986) *Basic and Clinical Aspects of IgG Subclasses, Monographs in Allergy*, Vol. **19**, Karger, Basel.

J. V. Ravetch and J. Kinet (1993) Fc receptors. *Annu. Rev. Immunol.* **9**, 457–492.

## Contact Inhibition

Cultured cells, particularly when grown as a monolayer attached to a substrate, will go through a reproducible growth cycle following subculture (1) (see [Cell Line](#)). The phases normally defined are (1) the lag phase, before the culture starts to proliferate; (2) the exponential, or log, phase, when the cells are dividing rapidly and the population doubles over a fixed time; and (3) the plateau phase, when the cell number remains stable with further time in culture. As a culture reaches the end of the exponential phase and enters plateau, it responds to cell contact, reduced substrate availability, a reduction in nutrients, and a number of other stimuli that inhibit a further increase in cell numbers. These have been divided into two major components: contact inhibition and density limitation of cell proliferation. Contact inhibition is the cessation of cell motility that occurs when a cell culture reaches confluence. It was first described from observations of randomly migrating fibroblasts made in time-lapse cinemicrography (2). Normal fibroblasts cultured at low cell density show a spindle-shaped morphology, with a reversible polarity along their long axis, which is also the direction of cell migration. When the cell encounters another similar cell, contact induces a reversal of the polarity; the ruffling in the leading tip of the cell ceases and transfers to the opposite tip, which then becomes the anterior end, and migration resumes in the opposite direction. When contact is made with another cell, the process of reversal of polarity is repeated, unless the cell is surrounded by other cells, when membrane ruffling ceases and cell migration is inhibited. In the cases of normal human or chick diploid fibroblasts, the cell assumes a parallel array that gives the appearance of whorls in the cell monolayer under low power [microscopy](#) or naked eye observation.

At the time that motility ceases, the cells are still capable of proliferation, but the culture becomes crowded following one or two population doublings, the cells assume a narrower spindle shape, and cell spreading is reduced beyond the point where the cell is able to re-enter the [cell cycle](#) (3). At this point, proliferation ceases in a normal cell culture, and the culture remains as a monolayer. Further growth is inhibited by density limitation of cell proliferation, and the culture enters the plateau phase of the growth cycle (see [Cell Line](#)). Density limitation of cell proliferation is distinct from contact inhibition, which affects primarily membrane ruffling, polarity, and cell migration. The mechanisms are, however, similar, although the detail remains obscure. Subconfluent cells have many more adhesions to the substrate than to other cells, favoring motility and proliferation. When the culture becomes confluent, the number of cell–cell adhesions (cadherins) increases, but they are still exceeded by cell–substrate adhesions (integrins), which inhibit cell motility but not proliferation. When the culture becomes crowded, the number of cell–cell adhesions increases, and cell–substrate adhesions are reduced, maintaining inhibition of motility but now inhibiting re-entry into the cell cycle. As contact with adjacent cells and contact with the substrate are mediated via different classes of **receptors** (4), the signaling that results from activation of these receptors is potentially different, although exactly how remains unclear.

In addition to alterations in cell spreading, adhesion receptors, and the [actin](#) cytoskeleton when cells reach a high density, the cells are also exposed to reduced nutrient levels, due to a higher rate of utilization, and increased levels of catabolites, both of which will tend to inhibit proliferation. In a static culture, this may generate a depleted layer above the cells, across which nutrients and [growth factors](#) must diffuse to reach the cells. The [diffusion](#) rate will be significantly slower with higher molecular weight components of the medium, including [growth factors](#) and peptide hormones such as [insulin](#). Support for this diffusion boundary layer hypothesis was provided by experiments where local irrigation was increased, by means of a micropump, and resulted in a localized increase in the frequency of labeling with [<sup>3</sup>H]-thymidine (5). As the range of the diffusion boundary was short, it was proposed that the limitation was principally in growth factors, depleted by binding to cell-surface receptors and [endocytosis](#), which was confirmed by the observation that addition of growth factors or serum to the medium will induce further proliferation (6).

On the other hand, one of the main observations that confirms that contact inhibition and density limitation of cell proliferation are not simply evidence of environmental deterioration is the result of the so-called “wounded monolayer” experiment (7). If a confluent, growth-arrested monolayer is scored with a sharp instrument, creating a bare patch devoid of cells, the surrounding cells respond as follows. The edge cells start to show ruffling of the cell membrane bordering the wound, spread out, and migrate into the space. When they have spread beyond a critical point, they enter the cell cycle and divide, and they continue to do so until the space of the wound is filled. They then stop migrating, but continue to proliferate until the cell density in the wound matches the crowded state of the rest of the monolayer; at this point they stop dividing. Meanwhile, no other cells distant from the wound show any sign of entering the cell cycle. This confirms that a major component of density limitation of cell proliferation is due to geometry and cell shape, rather than simply nutritional or growth factor limitation.

If a culture of normal fibroblasts that is growth arrested at high density is fed with medium containing serum or growth factors, many of the cells will reenter the cell cycle, resulting in the formation of a second cell layer over the first; if this is repeated, multilayered cultures can be produced. It has been proposed that the first layer of cells secretes a [collagen](#) overlay (which may produce a further diffusion barrier) and the second cell layer grows on top of this, technically not disobeying the rules of contact inhibition and density limitation of cell proliferation. This second layer, however, is still limited by contact and density, and when it reaches confluence, it will form a second layer of parallel arrays of cells, not colinear with the first. If perfused, these cultures may continue to increase up to 20 or 30 cell layers deep (8).

When cultures of normal epithelium reach confluence, they also become contact-inhibited. In less

dense cultures, however, they do not show the random migration seen in cultures of fibroblasts, but tend to grow in patches with ruffling of the membrane, cell spreading, and proliferation restricted to the outer edge of the patch (9). Their motility is limited, but when it occurs, it tends to involve the whole patch, which moves as a unit. Multilayering in normal epithelial cultures also occurs, but tends to imply maturation perpendicular to the substrate, rather than overgrowth. Epidermal cells, for example, will become stratified, with the upper layers of cells becoming progressively more keratinized. This is accentuated if the cells are grown on collagen, particularly in the presence of dermal fibroblasts and due to integrin down-regulation. Epithelial cells, which are derived from simple epithelium that is only one cell thick *in vivo*, tend to be obligate monolayers in culture. Likewise, normal endothelium from the lining of blood vessels will not pile up in culture, although after some time at high density it will tend to curl up and form secondary structures resembling capillaries (10), particularly if grown on an extracellular matrix, such as Matrigel.

**Transformation** allows cells to escape from the restrictions of contact inhibition and density limitation of cell proliferation (see definition of transformation in **Immortalization**). On one hand, they often lack, or have modified, the appropriate receptors or **cell-adhesion molecules** to recognize each other on contact; on the other hand, they are no longer dependent on cell spreading to allow entry into the cell cycle. Many transformed cell lines will propagate in suspension, without any substrate attachment. In addition, transformed cells often produce autocrine growth factors or have permanently active steps in the **signal transduction** cascades that promote cell proliferation (11-13). Hence, transformed cells will grow to a higher density and form multilayered cultures quite readily. The limitations imposed by density now tend to be diffusion-related, principally nutrient, catabolite, and gaseous (dissolved O<sub>2</sub> in and CO<sub>2</sub> out). One major limitation is the release of lactic acid by the cells, as they are generally more anaerobic in their metabolism than normal cells, and this depresses the pH, initially inhibiting proliferation and ultimately killing the cells. Although transformed cell cultures enter a plateau phase of culture when they reach a high density, the growth fraction in transformed cultures in plateau can be quite high (10–20%), unlike the plateau in normal cell cultures, where the growth fraction is very low (<5%). A steady state is reached when cell proliferation is matched by cell **necrosis** or **apoptosis**, and the cells deteriorate rapidly and irreversibly if the medium is not replenished.

The differences in response to high cell density between normal and transformed cells has been exploited in the isolation of transformed foci. If a monolayer of 3T3 cells is **transfected** with an **oncogene**, especially in the exponential phase of growth, when transfection is more efficient, it will show foci of transformed cells when the culture reaches confluence, because of the lack of contact inhibition of cell motility and density limitation of cell proliferation in successfully transformed cells (14). Likewise, transformed cells plated on a confluent monolayer will continue to grow and may form foci, or they may infiltrate the monolayer. Normal cells plated on a confluent monolayer will give a varying response, depending on the cell lineage of the confluent monolayer. If the plated normal cells are of the same lineage as the confluent monolayer, they will not grow; for instance, fibroblasts will not grow on a confluent monolayer of fibroblasts, and normal glial cells will not grow on a confluent monolayer of normal glial cells (15), whereas their transformed counterparts will. On the other hand, normal glial cells will grow on a confluent monolayer of fibroblasts, and normal keratinocytes will grow on a monolayer of normal 3T3 cells (16).

Cultures that are propagated in suspension will, of course, not encounter contact inhibition. Nevertheless, they are subject to density limitation of cell proliferation. Part of this is undoubtedly due to nutrient depletion, catabolite buildup, and pH depression, but even if limiting nutrients are replaced and the pH is stabilized, high density suspension cultures do not increase significantly above  $1-2 \times 10^6$  cells/ml. Quite often suspension cultures, such as **hybridomas**, will remain in the plateau phase for only a short time and then deteriorate rapidly, leaving few viable cells in the culture. The reason for this appears to be that the cells tend to enter apoptosis when they reach a high density. Some moderate success has been achieved overexpressing *bcl* in these cells to inhibit apoptosis (17).

## Bibliography

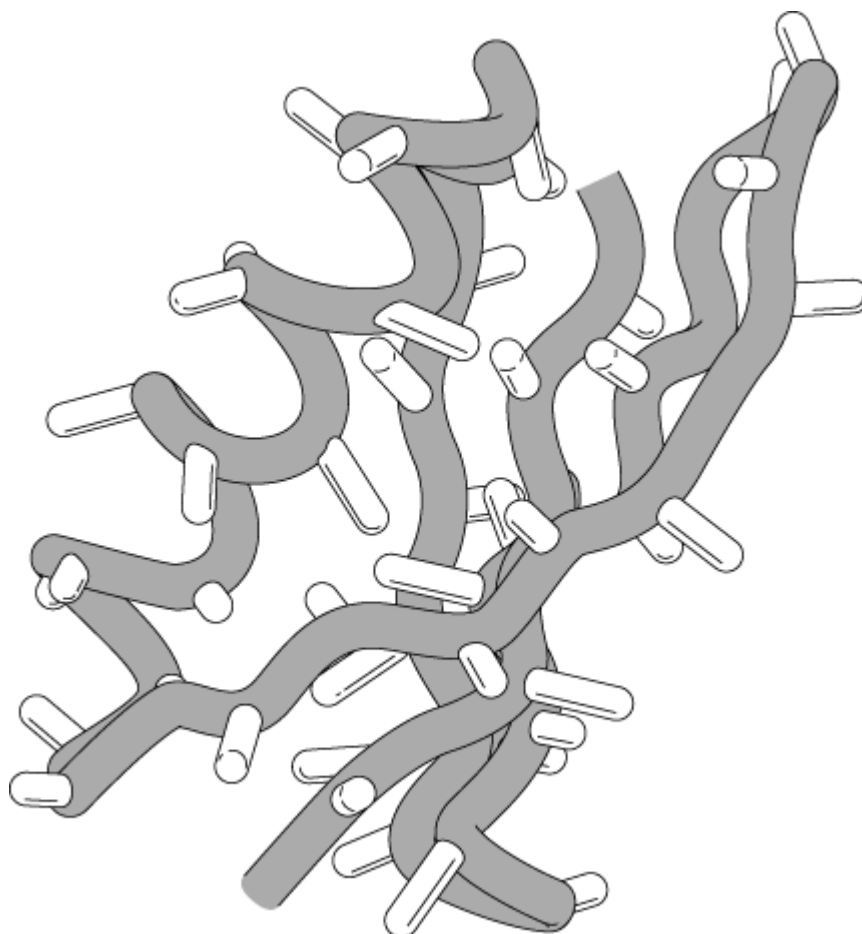
1. R. I. Freshney (1994) *Culture of Animal Cells, a Manual of Basic Technique*, Wiley-Liss, New York, pp. 153–157.
2. M. Abercrombie and J. E. M. Heaysman (1954) *Exp. Cell Res.* **6**, 293–306.
3. I. Folkman and A. Moscona, (1978) *Nature* **273**, 345–349.
4. S. Levenberg, B.-Z. Katz, K. M. Yamada, and B. Geiger (1998) *J. Cell Sci.* **111** 347–357.
5. M. G. P. Stoker (1973) *Nature* **246**, 200–203.
6. G. A. Dunn and G. W. Ireland (1984) *Nature* **312**, 63–65.
7. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1989) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing, New York, p. 898.
8. P. F. Kruse Jr. and E. Miedema (1965) *J. Cell Biol.* **27**, 273.
9. I. McKay and J. Taylor-Papadimitriou (1981) *Exp. Cell Res.* **134**, 465–470.
10. J. Folkman and C. Haudenschild (1980) *Nature* **288**, 551–556.
11. P. Kahn and T. Graf (1986) *Oncogenes and Growth Control*, Springer, Berlin.
12. K. Siegfried, Y. H. Han, M. A. A. DeMichele, J. D. Hunt, A. L. Gaither, and F. Cuttitta (1994) *J. Biol. Chem.* **269**, 8596–8603.
13. A. Balmain and K. Brown (1988) *Adv. Cancer Res.* **51**, 147–182.
14. R. I. Freshney (1994) *Culture of Animal Cells, a Manual of Basic Technique*, Wiley-Liss, New York, p. 233.
15. C. M. MacDonald, R. I. Freshney, E. Hart, and D. I. Graham (1985) *Exp. Cell Biol.* **53**, 130–137.
16. J. G. Rheinwald and H. Green (1975) *Cell* **6**, 331–344.
17. S. Terada, Y. Itoh, H. Ueda, and E. Suzuki (1997) *Cytotechnology* **24**, 135–141.

## Contact Maps

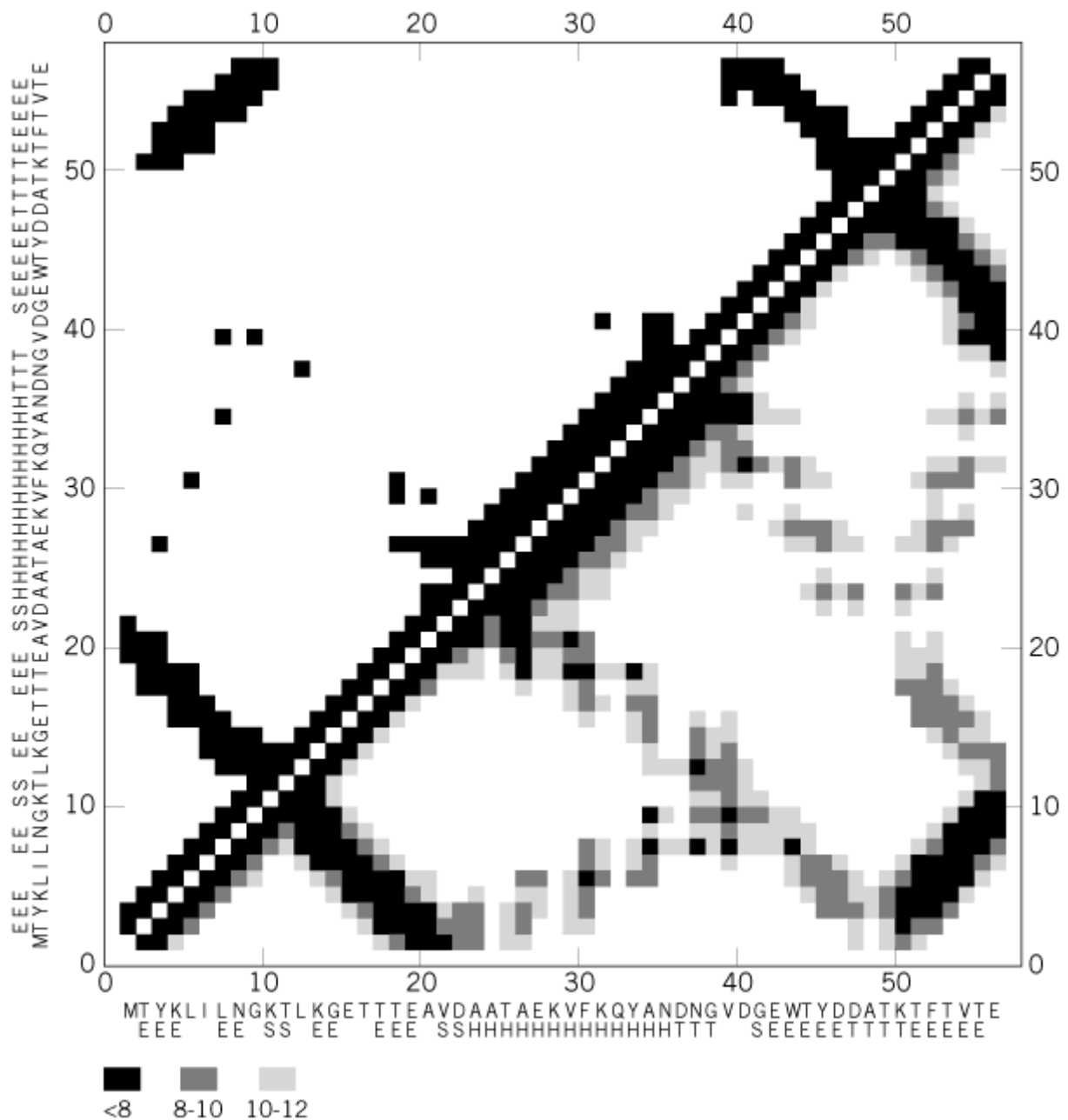
Contact maps are two-dimensional representations of three-dimensional [protein structures](#). A three-dimensional description of a protein structure composed of  $N$  structural units could be expressed as an  $N \times N$  array of the pairwise distances (see [Distance Geometry](#)). This could be done for all pairs of atoms, for selected types of atoms (eg, Ca atoms), for groups of atoms (eg, [side-chain](#) centers of mass), or for entire [amino-acid](#) residues. Contact maps are generated from such matrices by taking a certain cutoff value for the pairwise distances. For example, the  $N \times N$  matrix of the distances between protein Ca atoms ([1](#), [2](#)) can be transformed into a Ca-based contact map ([3](#)). Those Ca atoms that are closer to each other in the protein structure than the chosen cutoff distance are considered to be “in contact.” This produces a binary  $N \times N$  matrix—a so-called black and white contact map. Alternatively, one may assume a set of several critical values for the distances between Ca atoms (or for other atoms or groups of atoms) and generate an integer matrix—the equivalent of a contact map with a colored or gray scale. A map of main-chain [hydrogen bonds](#) could also be considered as a variant of a protein contact map. The choice of structural units being mapped and the choice of cutoff distances determine the quality and range of structural information being stored in a contact map (see Fig. [1](#)).

**Figure 1.** Schematic drawing of the structure of the B domain of protein G and its contact maps. (a) Three-dimensional structure of the B domain. (b) Ca-based contact maps: above the diagonal, the black and white map with a cutoff distance of 8 Å; below the diagonal is the gray-scale map where various shades of gray correspond to three values of the cutoff

distance, 8 (darkest), 10, and 12 Å. (c) Side-chain-based contact map. Above the diagonal, the dark squares correspond to the pairs of side chains for which the distance between at least one pair of heavy atoms is less than 5 Å. Below the diagonal, a 6.25-Å cutoff criterion has been applied to the centers of mass of the side chains. (d) The main chain hydrogen bond map for the same structure. The amino acid sequence of the protein is given along each axis in one-letter code, and the secondary structure is indicated as follows: E is extended b-strand, H is helix, and S and T are two types of turns.

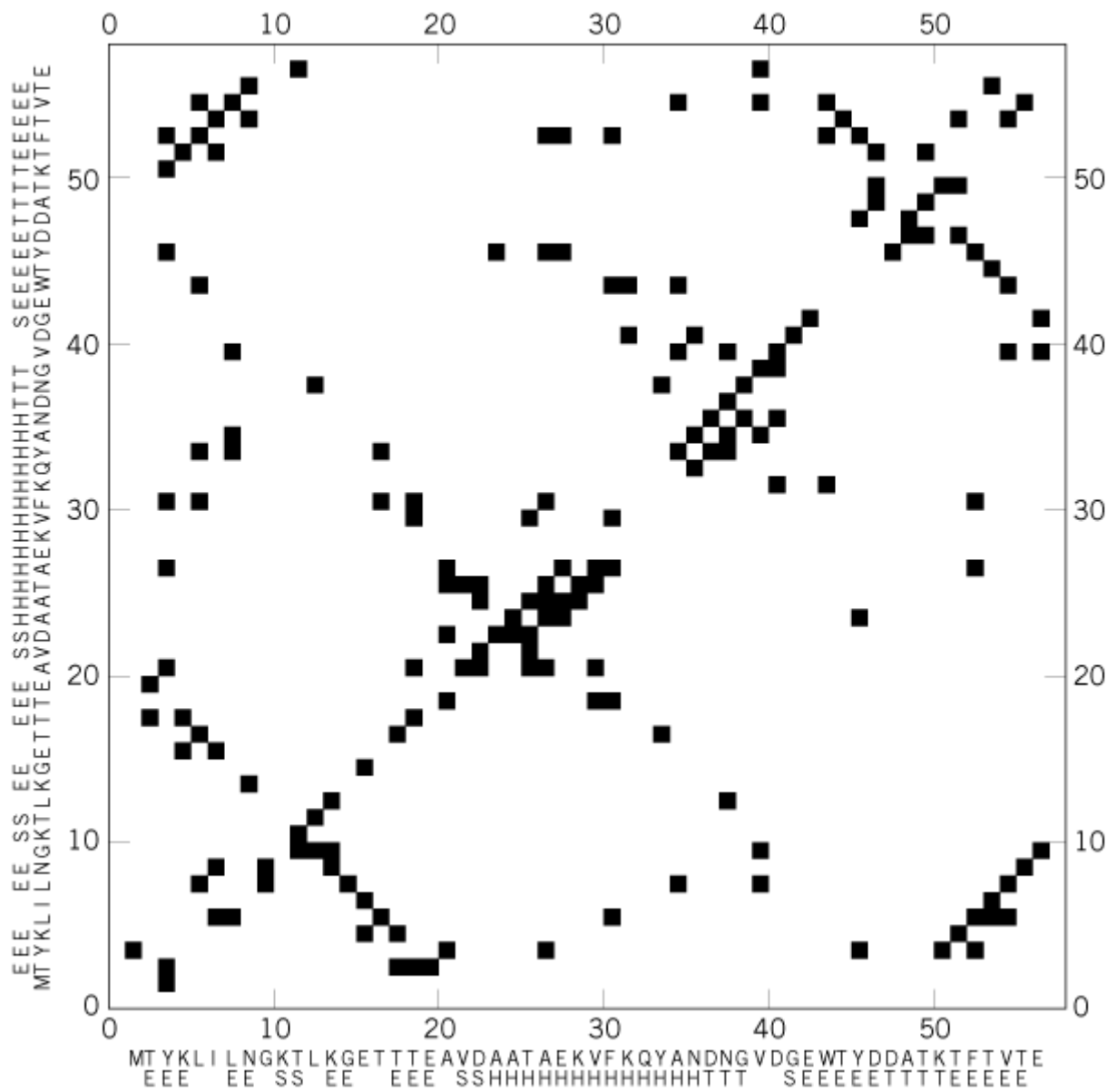


**Figure 1.** (*continued*) Schematic drawing of the structure of the B domain of protein G and its contact maps. (a) Three-dimensional structure of the B domain. (b) Ca-based contact maps: above the diagonal, the black and white map with a cutoff distance of 8 Å; below the diagonal is the gray-scale map where various shades of gray correspond to three values of the cutoff distance, 8 (darkest), 10, and 12 Å. (c) Side-chain-based contact map. Above the diagonal, the dark squares correspond to the pairs of side chains for which the distance between at least one pair of heavy atoms is less than 5 Å. Below the diagonal, a 6.25-Å cutoff criterion has been applied to the centers of mass of the side chains. (d) The main chain hydrogen bond map for the same structure. The amino acid sequence of the protein is given along each axis in one-letter code, and the secondary structure is indicated as follows: E is extended b-strand, H is helix, and S and T are two types of turns.

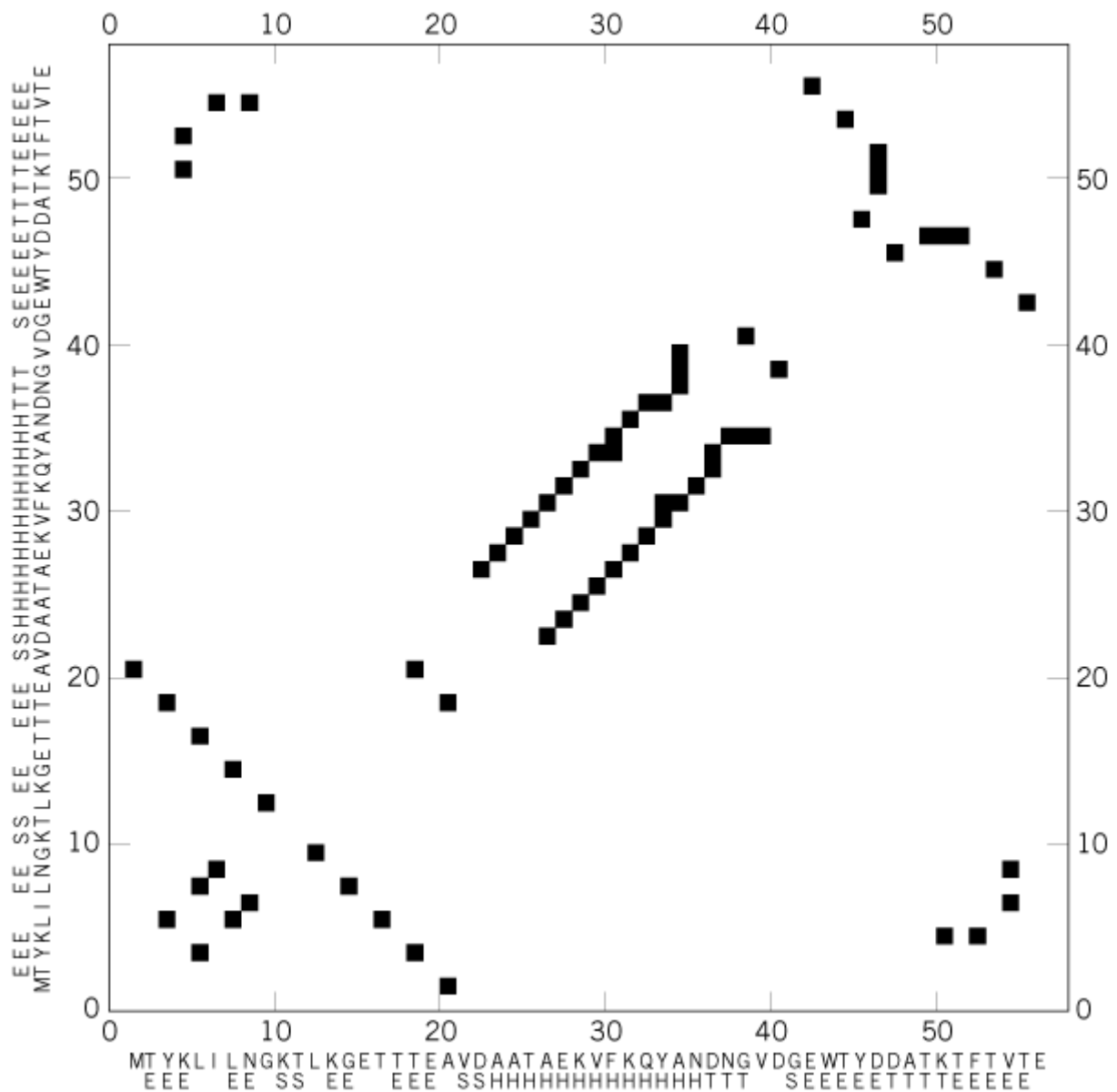


**Figure 1.** (*continued*) Schematic drawing of the structure of the B domain of protein G and its contact maps. (a) Three-dimensional structure of the B domain. (b) Ca-based contact maps: above the diagonal, the black and white map with a cutoff distance of 8 Å; below the diagonal is the gray-scale map where various shades of gray correspond to three values of the cutoff distance, 8 (darkest), 10, and 12 Å. (c) Side-chain-based contact map. Above the diagonal, the dark squares correspond to the pairs of side chains for which the distance between at least one pair of heavy atoms is less than 5 Å. Below the diagonal, a 6.25-Å cutoff criterion has been applied to the centers of mass of the side chains. (d) The main chain hydrogen bond map for the same structure. The amino acid sequence of the protein is given along each axis in one-letter code, and the secondary structure is indicated as follows: E is extended b-strand, H is helix, and S and T are two types of turns.





**Figure 1.** (*continued*) Schematic drawing of the structure of the B domain of protein G and its contact maps. (a) Three-dimensional structure of the B domain. (b) Ca-based contact maps: above the diagonal, the black and white map with a cutoff distance of 8 Å; below the diagonal is the gray-scale map where various shades of gray correspond to three values of the cutoff distance, 8 (darkest), 10, and 12 Å. (c) Side-chain-based contact map. Above the diagonal, the dark squares correspond to the pairs of side chains for which the distance between at least one pair of heavy atoms is less than 5 Å. Below the diagonal, a 6.25-Å cutoff criterion has been applied to the centers of mass of the side chains. (d) The main chain hydrogen bond map for the same structure. The amino acid sequence of the protein is given along each axis in one-letter code, and the secondary structure is indicated as follows: E is extended b-strand, H is helix, and S and T are two types of turns.



## 1. Contact Maps as a Fingerprint of Protein Three-Dimensional Structure

A contact map constitutes a structural “fingerprint” of a protein (4). Each protein can be identified based on its contact map. The **secondary structure**, fold topology, and side-chain packing patterns (for side-chain contact maps) can be visualized conveniently and read from the contact map. Furthermore, structural similarity between a pair of proteins is immediately apparent by a very pronounced similarity of their contact maps; in comparing two protein structures, there is no need to search all their possible relative orientations. The reconstruction of a protein structure from its contact map is more complex, although low-to-moderate resolution three-dimensional models can be easily built, even from a fragmentary contact map (5). The accuracy of the model depends on the type of contact map and the computational tools employed. A combination of protein nuclear magnetic resonance [NMR](#) spectra (see [NOESY Spectrum](#); [COSY Spectrum](#)) constitutes a hybrid contact map of a protein, and model building from these data is an example of a map-to-structure modeling procedure (6, 7).

## 2. Ca-Based Contact Maps

Ca-based contact maps and distance matrices were perhaps the first commonly used maps for visualization of protein structures (3, 8, 9). An example is given in Figure 1b. These contact maps reflect well the overall topology of the protein fold, but only rather coarse structural details can be

read from them. This is due to the fact that the Ca–Ca distance distributions extracted from protein structures have several convoluted peaks. These peaks correspond to various distances between pairs of various secondary structure elements ( $\alpha$ -helices, [beta-strands](#), etc.). Thus, a single cutoff distance is always inadequate: Too small a value would miss some helix-to-helix contacts, while too large a value may create some problems with identification of the secondary structure patterns. Gray scale maps (several cutoff ranges) communicate much more detailed structural information.

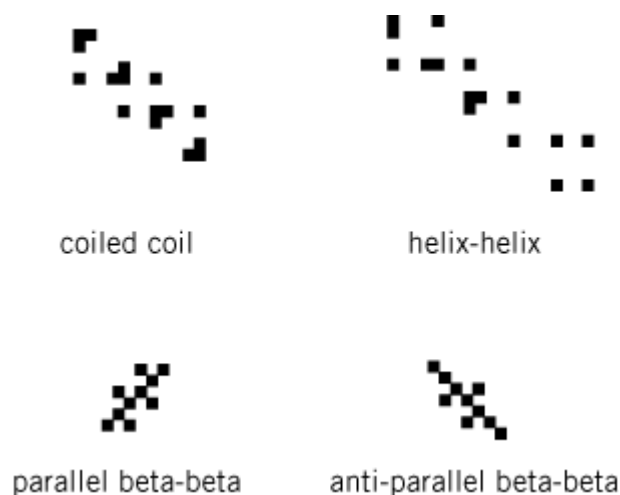
### 3. Side-Chain Contact Maps

Side-chain contact maps contain much richer information, not only about the topology of a protein fold and its secondary structure, but also many fine details about the packing patterns of the protein side-chains. Various conventions can be used to build side-chain contact maps; two are illustrated in [Figure 1](#). In one case, two residues are assumed to be in contact when any two heavy atoms (ie, all except hydrogen) are a shorter distance from each other than some assumed cutoff. Due to the comparable size of all the united atom types constituting the side chains (eg,  $\text{CH}_2$ , CH,  $\text{NH}_2$ , etc.), a good choice of cutoff distance is between 4.5 and 5.0 Å ([4](#), [10](#)). In this range, the number of detected contacts is not sensitive to the particular choice of cutoff, and the packing pattern of the side chains is always described with high fidelity. Characteristic patterns of contacts between elements of secondary structure are an important and useful feature of these contact maps ([11](#)). Alternatively, one may build a side-chain contact map using the side-chain centers of mass as a reference. In this case, a larger value of the cutoff distance needs to be used. These cutoff values for the distance between side-chain centers of mass could be made specific for certain amino acid pairs, on the basis of the different sizes of the side chains, which produce different average contact distances for various pairs. As seen in [Figure 1](#), the two approaches (atom based and center of mass based) lead to very similar protein representations. The patterns of the atom-based contact maps are slightly better defined.

### 4. Regularities of the Contact Maps Reflect Regularities of Protein Structures

Different types of contact maps reflect different aspects of the regularities seen in protein structures. The side-chain-based contact maps are a very good example ([10](#)). Near the diagonal of the map, the distinct features of the protein secondary structure can be easily read. Indeed, for extended fragments of the polypeptide chain, only residues  $i$  and  $i + 2$  can be in contact. For  $\alpha$ -helices, the  $i$ ,  $i + 3$  and  $i$ ,  $i + 4$  patterns of contacts are well pronounced. Furthermore, characteristic clusters of contacts further away from the diagonal reflect the packing between particular pairs of  $\beta$ -strands within the  $\beta$ -sheets. Parallel and antiparallel structures have very different features on the contact patterns. Very characteristic patterns could also be observed for other pairs of secondary-structure elements. It is even very easy to distinguish between the patterns for two helices in a helical or  $\alpha/\beta$  protein and in the [coiled-coil](#) structural motifs. [Figure 2](#) shows some typical side-chain contact map patterns.

**Figure 2.** Representative patterns of side-chain contact maps describing interactions between  $\alpha$ -helices (top) and between  $\beta$ -strands in parallel and antiparallel [beta-sheets](#) (bottom).



## 5. Application of Contact Maps in Modeling Protein Structure

Contact maps provide a very convenient way of identifying pairwise interactions within protein structures. This has been applied in various algorithms for [threading protein sequences](#), where contact maps are actually used as ersatz template structures (4). In [computer simulations of biological molecules](#), contact maps provide a convenient way of displaying structural changes. Also, regularities of the patterns seen in all classes of proteins can provide a guideline for designing knowledge-based multibody potentials (12, 13) and for protein modeling in a reduced (and perhaps also in all-atom) representation (13, 14). Such potentials may be necessary to reproduce the all-or-none character of protein folding transitions in model simulations (15).

One can easily recognize the distinct features of well-defined contact maps, with their characteristic patterns, after inspecting several maps of various proteins. This is an excellent example of a pattern recognition problem that could be learned by **neural network calculations**. Then such a trained network can be used for the automated recognition of good versus poor models of protein structure (16).

### Bibliography

1. I. D. Kuntz (1975) *J. Am. Chem. Soc.* **97**, 4362–4366.
2. F. M. Richards and C. E. Kundrot (1988) *Proteins* **3**, 71–84.
3. M. Levitt (1976) *J. Mol. Biol.* **104**, 59–107.
4. A. Godzik, J. Skolnick, and A. Kolinski (1992) *J. Mol. Biol.* **227**, 227–238.
5. J. Skolnick, A. Kolinski, and A. R. Ortiz (1997) *J. Mol. Biol.* **265**, 217–241.
6. W. Braun and N. Go (1985) *J. Mol. Biol.* **186**, 611–626.
7. R. Kaptein, R. Boelens, R. M. Scheek, and W. F. van Gunsteren (1988) *Biochemistry* **27**, 5389–5395.
8. M. N. Liebman, C. A. Venanzi, and H. Weinstein (1985) *Biopolymers* **24**, 1722–1758.
9. D. C. Phillips (1970) *Biochem. Soc. Symp.* **30**, 11–28.
10. A. Godzik and C. Sander (1989) *Protein Eng.* **2**, 589–596.
11. A. Godzik, J. Skolnick, and A. Kolinski (1993) *Protein Eng.* **6**, 801–810.
12. A. Kolinski, A. Godzik, and J. Skolnick (1993) *J. Chem. Phys.* **98**, 7420–7433.
13. A. Kolinski and J. Skolnick (1996) *Lattice Models of Protein Folding, Dynamics and Thermodynamics*, R. G. Landes, Austin, TX.
14. K. A. Olszewski, A. Kolinski, and J. Skolnick (1996) *Proteins* **25**, 286–299.
15. A. Kolinski, W. Galazka, and J. Skolnick (1996) *Proteins* **26**, 271–287.
16. M. Milik, A. Kolinski, and J. Skolnick (1995) *Protein Eng.* **8**, 225–236.

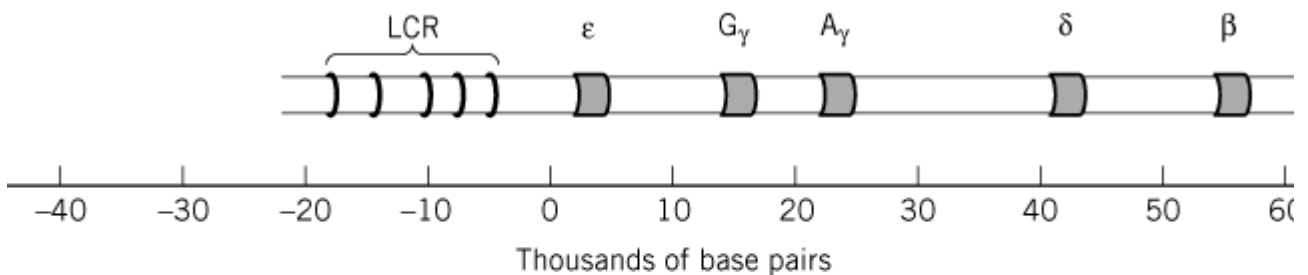
## Contiguous Genes

The order of genes on a [chromosome](#) and the existence of directly contiguous genes can be useful in understanding their individual functions and in mapping the genome to understand the function of chromosomal **domains**. In **prokaryotes** there is often a close linkage of genes involved in a common metabolic pathway. This type of linkage is uncommon in **eukaryotes**, although examples exist, including five genes in *Neurospora crassa* involved in synthesis of chorismic acid (1) and two genes in *Drosophila melanogaster* involved in purine metabolism (2).

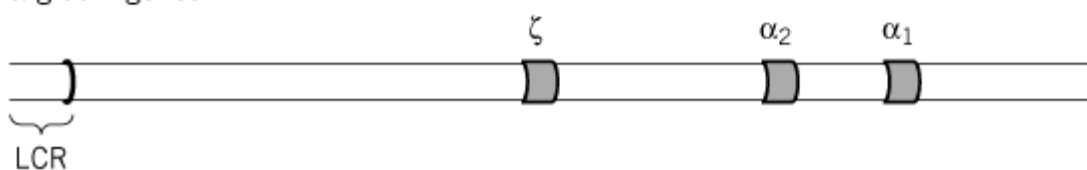
There are several families of closely linked genes in humans. These include genes for [alcohol dehydrogenase](#), [α-amylase](#), the [major histocompatibility complex](#), [α-globins](#), [β-globins](#), [chorionic gonadotropin](#), [luteinizing hormone](#), and [growth hormone](#). Every chromosome has at least one family cluster of undispersed genes. These clusters vary in size from several thousand base pairs to millions of base pairs with a large cluster like the major histocompatibility complex. The members of a single **gene family** are derived from the [gene duplication](#) of single ancestral genes. Whereas gene families are rare in prokaryotes, the more complex differentiative states in eukaryotes require similar gene products with subtle differences in function. An example of a well-studied locus, where genes with closely related functions are linked together, occurs at the α- and β-globin loci. At these chromosomal sites, globin genes specific for embryos, neonatal and adult animals are linked together and expressed in a sequential fashion during [development](#) (Fig. 1). The globin loci also present interesting examples of where gene linkage helps understanding the chromosomal function. All of the globin genes depend on a super [enhancer](#) known as a **locus control region** (LCR), which functions over an entire chromosomal **domain**.

**Figure 1.** The β-globin gene cluster, which is on the short arm of chromosome 11 in humans and includes five genes that are sequentially expressed during development. Its locus control region (LCR) consists of five short stretches of DNA extending over a span of 12 kb. Similarly, there are three α-globin genes on the short arm of chromosome 16, but that LCR comprises only one DNA sequence far upstream.

### β-globin genes



### α-globin genes



The relative order of genes is often conserved within a chromosome. This is termed linkage conservation. Linkage conservation often occurs in chromosomes isolated from distinct species. For example, in the [X-chromosome](#) there are many genes that have the same order along the chromosome in both [mouse](#) and humans. This is a useful aid in mapping new genes. There are several potential evolutionary advantages to linking certain genes together, which include common regulatory mechanisms, at the level of establishing a functional chromosomal domain, as described for the globin genes (3), or in terms of coordinate expression, as described for two genes involved in a particular metabolic pathway. There may also be regulatory mechanisms where genes compete with each other for a common regulatory element, as occurs in the globin genes at certain times in development (4).

### Bibliography

1. F. H. Gaertner and K. W. Cole (1977) *Biochem. Biophys. Res. Comm.* **75**, 259–270.
2. M. E. Johnstone (1985) *Biochem. Genet.* **23**, 539–546.
3. F. Grosveld, G. B. van Assendelft, D. R. Greaves, and G. Kollias (1987) *Cell* **51**, 975–985.
4. M. Wijgerde, F. Grosveld, and P. Fraser (1995) *Nature* **377**, 209–213.

### Contrast Variation

“Contrast” in scattering and diffraction experiments is defined as the difference between the mean scattering density of a component and its background, which can be a solvent or other components of an assembly of biological molecules. The greater the contrast, the more readily a component can be distinguished from its surroundings. “Contrast variation” involves the manipulation of the contrasts of specific components in a system in order to extract structural information on individual components. Contrast variation used in combination with [small-angle scattering](#) is a powerful method for examining the shapes and interactions of biological molecules in solution. If one has ordered samples, then contrast variation used in combination with **neutron diffraction** experiments can give important information on the location of disordered components that cannot be seen in X-ray diffraction experiments.

#### 1. Contrast Variation with Neutrons

Either contrast variation experiments take advantage of inherent differences in scattering density between components of a complex or assembly, or the experimenter introduces contrast into the system by manipulating the scattering density of a specific component. Scattering densities are calculated by summing the scattering amplitudes of each atom within a volume and dividing by that volume. Because X-rays are scattered by electrons and the X-ray scattering amplitudes of atoms increase monotonically with the number of electrons, it is difficult to change the scattering density of a biological molecule in a benign way. On the other hand, neutrons provide extremely elegant and practical means for contrast variation. Neutrons are scattered by atomic nuclei in a sample. Hence neutron scattering amplitudes depend upon the complex properties of the neutron–nucleus interaction, and they show no systematic dependence on atomic number. Furthermore, isotopes of the same element can have very different neutron scattering properties. For neutrons, one of the largest differences in neutron scattering amplitude is between the isotopes of hydrogen ( $^1\text{H}$  and  $^2\text{H}$ ). Table 1 lists the coherent, elastic neutron scattering amplitudes for the atoms commonly found in biological systems. Note that the scattering amplitudes for most nuclei are positive and

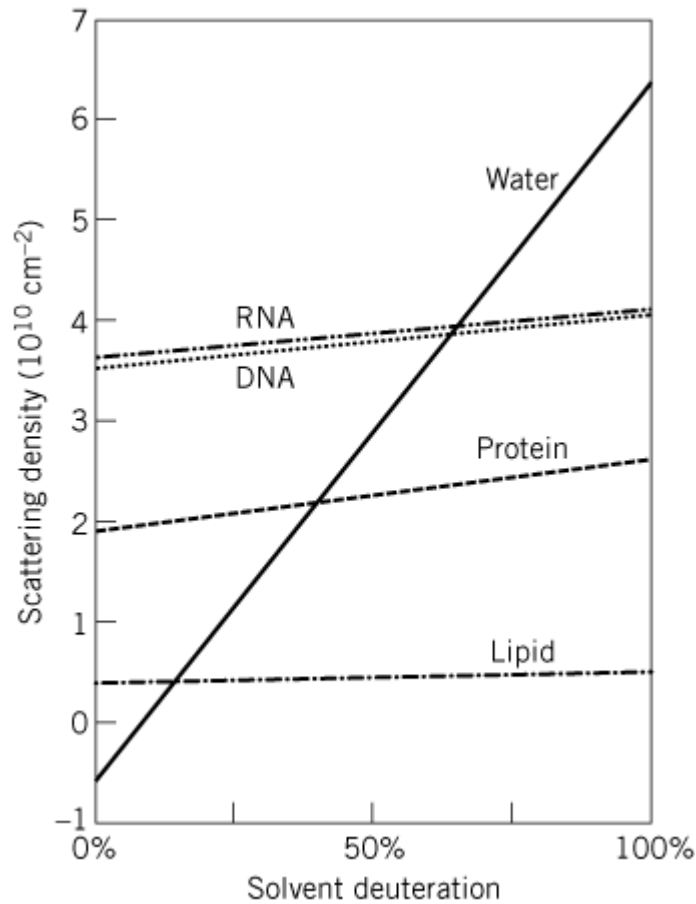
approximately equal. The exception is the scattering amplitude for  $^1\text{H}$  which is negative, resulting from a  $180^\circ$  phase shift between neutrons scattered by  $^1\text{H}$  compared to the other nuclei. As a consequence, the neutron scattering density of a particle depends strongly on its mean hydrogen content. The basic biological constituents, proteins, polynucleotides, and lipids, each have quite different mean neutron scattering densities (Fig. 1). Furthermore, the scattering densities for pure  $^1\text{H}_2\text{O}$  and  $^2\text{H}_2\text{O}$  bracket the values for biological constituents. Thus by simply varying the  $^2\text{H}_2\text{O}:^1\text{H}_2\text{O}$  ratio in the solvent, one can do a contrast variation experiment on a protein–DNA complex, or a membrane protein in a bilayer. For a complex of proteins, selective deuteration of individual proteins provides a way of altering their mean neutron scattering densities for contrast variation studies. Recently, it has been demonstrated that one can obtain dramatically increased contrast effects using a novel method of “spin contrast variation.” This method uses polarized neutron scattering from samples labeled with dynamically polarized nuclei and has been demonstrated to be effective for studying biological systems by locating RNA (1) and proteins (2, 3) in ribosomes.

**Table 1. Coherent Neutron Scattering Lengths,  $b_{\text{coh}}$ , and Corresponding X-Ray Scattering Factors,  $f_{\text{X-ray}}$ , for Biologically Relevant Nuclei**

| Atom        | Nucleus                | $(10^{-12} \text{ cm}) f_{\text{X}}$ for $q = 0 (10^{-12} \text{ cm})^a$ |      |
|-------------|------------------------|--|------|
| Hydrogen    | $^1\text{H}$           | −0.3742  | 0.28 |
| Deuterium   | $^2\text{H}$           | 0.6671   | 0.28 |
| Carbon      | $^{12}\text{C}$        | 0.6651   | 1.69 |
| Nitrogen    | $^{14}\text{N}$        | 0.940  | 1.97 |
| Oxygen      | $^{16}\text{O}$        | 0.5804   | 2.25 |
| Phosphorous | $^{31}\text{P}$        | 0.517  | 4.23 |
| Sulfur      | mostly $^{32}\text{S}$ | 0.2847   | 4.5  |

<sup>a</sup> X-ray scattering amplitudes are normally given in units of electrons, but have been converted to cm here for comparison with neutron scattering amplitudes. At very short wavelengths and low- $Q$ , the X-ray coherent scattering cross-section of an atom with  $Z$  electrons is  $4\pi(Zr_0)^2$ , where  $r_0 = e^2/m_e c^2 = 0.28 \times 10^{-12} \text{ cm}$ .

**Figure 1.** Average scattering length densities for biological molecules as a function of the fraction of  $\text{D}_2\text{O}$  in the solvent. The slope of the lines arises from the exchange of labile hydrogens.



## 2. Small-Angle Scattering with Contrast Variation

The [scattering intensity distribution](#) from a homogeneous solution of monodisperse particles in a solution can then be expressed as

$$I(Q) = \int \int \Delta\rho(\mathbf{r}_1)\Delta\rho(\mathbf{r}_2) \frac{\sin Q|\mathbf{r}_1 - \mathbf{r}_2|}{Q|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (1)$$

The integration is taken over the volume of the particle.  $Q$  is the momentum transfer or scattering vector amplitude and is equal to  $4\pi(\sin q)/\lambda$ , where  $q$  is half the scattering angle, and  $\lambda$  is the wavelength of the incident neutron radiation.  $\Delta\rho(\mathbf{r}) = \rho(\mathbf{r}) - \rho_s$  is the “contrast,” or neutron scattering density difference between the particle and the solvent. For uniform scattering density particles, if one can “match,” or make equal, the scattering density of a particle and its solvent, then the small-angle scattering from that particle will be zero; that is, the particle will become “invisible” in the scattering experiment. For a multicomponent assembly, one can use this technique to study the shape of an individual component within the assembly. Rigorous solvent matching can be difficult, however, and is made more difficult when there are significant internal density fluctuations within the component being matched.

For a two-component complex in solution, one can conduct a series of measurements at different contrast points (“contrast series”) to extract the basic scattering functions of individual components. Ignoring internal scattering density fluctuations, the scattering from a two-component complex in solution can be written as

$$I(Q, \Delta\rho_1, \Delta\rho_2) = \Delta\rho_1^2 I_1(Q) + \Delta\rho_1 \Delta\rho_2 I_{12}(Q) + \Delta\rho_2^2 I_2(Q) \quad (2)$$



The subscripts 1 and 2 refer to each component; and  $\Delta r_{1(2)} = r_{1(2)} - r_s$ , where  $r_{1(2)}$  is that mean scattering density for component 1 (or 2).  $I_1(Q)$  and  $I_2(Q)$  represent the scattering of components 1 and 2, respectively, while  $I_{12}(Q)$  is the cross term. A set of scattering measurements with different solvent scattering densities, or contrasts, gives a set of equations in the form of equation 2, which can be solved by linear least-squares regression to give the three basic scattering functions  $I_1(Q)$ ,  $I_2(Q)$ , and  $I_{12}(Q)$ . Using these functions, the individual structural parameters for each component such as molecular weight, radius of gyration,  $R_g$ , and the pair distribution function,  $P(r)$ , can be calculated, as well as information on the relative dispositions of the individual components (see [Small-Angle Scattering](#) and [Radius Of Gyration](#) for details).  $P(r)$  is the probable frequency distribution of all vector lengths between pairs of volume elements within the scattering component, weighted by the product of the scattering densities of each pair of volume elements.  $P(r)$  is calculated as the indirect Fourier transform of the basic scattering function for each component, and its second moment gives  $R_g$  for that component. The indirect Fourier transform of the cross term  $I_{12}(Q)$  gives the pair distribution function containing intercomponent vector lengths, which can give information on the distance between the two components.

Ibel and Stuhrmann (4) proposed an alternative approach to determining the relative dispositions of components in a complex. They showed that the  $R_g^2$  dependence on the scattering contrast can be written as

$$R_g^2 = R_m^2 + \frac{\alpha}{\Delta\rho} - \frac{\beta}{\Delta\rho^2} \quad (3)$$

where  $R_m$  is  $R_g$  at infinite contrast; and  $\Delta r = r - r_s$ , where  $r$  is the mean scattering density for the complex and  $r_s$  is the scattering density of the solvent. The coefficient  $\alpha$  is related to the second moment of the scattering density fluctuations about the mean value for the complex, while  $\beta$  is related to the square of the first moment of the density fluctuations about the mean. If the sign of  $\alpha$  is positive, the lower scattering density component is more toward the inside of the complex than the higher scattering density component. A negative  $\alpha$  indicates the reverse.  $\beta$  is proportional to the square of the separation of the centers of mass of the two components. If  $\beta$  is zero, then the centers of mass are coincident.

### 3. Deuteration of Biological Molecules

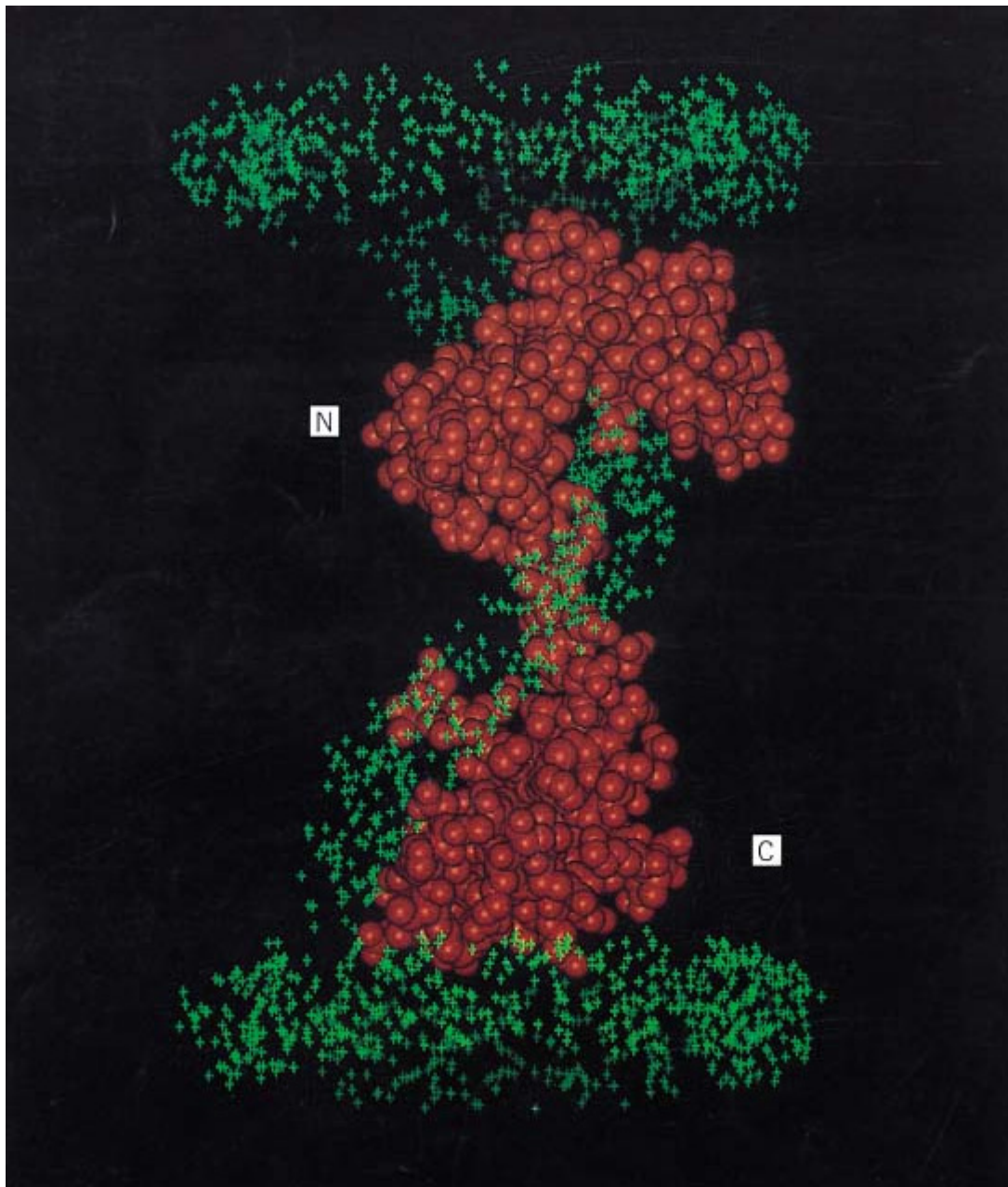
Neutron scattering and/or diffraction experiments using [contrast variation](#) frequently depend upon the ability to reconstitute complexes with specific components labeled uniformly with deuterium. Deuteration of biological molecules is achieved by growing microorganisms on deuterated media. Relatively high average levels of deuteration (50% for nucleic acids, 70 to 80% for proteins) can be achieved by growing organisms on  $^2\text{H}_2\text{O}$  with normal protiated carbon sources (5). Typically, however, the aromatic groups are not highly deuterated using this approach. Higher levels of enrichment (> 99%) are achieved using deuterated nutrients and  $^2\text{H}_2\text{O}$  (5). Deuterated algal hydrolysate provides an excellent rich growth medium for *E. coli* containing an expression system with a specific protein of interest. If the expression is sufficiently robust, deuteration can be achieved in minimal growth media with deuterated glucose, glycerol, or succinate in  $^2\text{H}_2\text{O}$ . Intermediate levels of uniform deuterium labeling can be achieved by diluting  $^2\text{H}$  with  $^1\text{H}$  in the growth medium.

### 4. Examples of the Application of Contrast Variation

#### 4.1. Small-Angle Scattering and Contrast Variation

Examples of significant advances based on small-angle neutron scattering with contrast variation include: the solution structure of the [nucleosome](#) core particle showing the DNA winding around the outside of the [histone](#) core ([6](#), [7](#)); mapping the radial distribution of protein and nucleic acid components in the [influenza virus](#) ([8](#)); determining the subunit structure of DNA-dependent **RNA polymerase** ([9](#)) and its relationship to a 130-bp DNA fragment ([10](#)); the subunit structure of a ribosome particle, showing the dispositions its components ([11](#), [12](#)) (see [Neutron Diffraction And Scattering](#)); and the determination of the conformations and dispositions calmodulin complexed with its target enzymes ([13](#)) and of the evolutionarily related muscle protein troponin C (TnC) complexed with its regulatory target troponin I (TnI) ([13](#)) which acts as a calcium-sensitive “switch” to regulate the contractile apparatus. Figure [2](#) illustrates this latter experiment, showing the model derived from a neutron contrast series on complexes of TnI with deuterated TnC in its calcium-saturated form ([14](#)). The neutron data show TnC in the complex has the unusual dumbbell shape evident in the TnC crystal structure ([15](#), [16](#)), while TnI forms a spiral structure that passes through two hydrophobic clefts in each globular domain of TnC. The diameter of the TnI central spiral is close to that expected for an  $\alpha$ -helix. The model suggests a mechanism for the action of the calcium-sensitive TnC/TnI molecular switch. The hydrophobic cleft of the *N*-terminal regulatory domain of TnC alternately opens and closes (in response to calcium binding and release) and thus binds and releases TnI, which in turn switches “off” and “on” TnI's inhibition of the interactions between the muscle fibers required for contraction.

**Figure 2.** Model structure for the muscle protein complex troponin C/troponin I derived from small-angle neutron scattering and contrast variation ([14](#)). Troponin C is represented as a space filling model based on its crystal structure, and the green crosses represent the volume occupied by the troponin I. This figure is based on work done at Los Alamos National Laboratory. See color insert.



#### 4.2. Diffraction Experiments with Contrast Variation

In crystals of large biomolecular complexes, there are frequently disordered regions, even when the diffraction data extend to relatively high resolution. The scattering from the disordered components is seen predominantly at low- $Q$ —that is, in the “small-angle” region where the contrast effects can be strong. Deuteration strategies similar to those used in small-angle neutron scattering allow for contrast manipulation. Phase information can be derived for the low- $Q$  data using contrast-variation procedures (17), beginning with starting phases from a [X-ray crystallography](#) or a model structure and using the contrast-variation relationships to produce the best scattering-density maps (18).

Examples of experiments using contrast variation and low-resolution neutron crystallography include the determination of the location of 25% of the protein and the RNA in crystals of the [Tomato Bushy Stunt Virus](#) (19) that were “missing” in the 2.8 Å resolution X-ray structure. The most important conclusion of this study is that at the virus radius where the protein and RNA interact, the quasi-equivalence of the 180-protein subunits breaks down, and the RNA interacts preferentially with 60 of the subunits. This observation favors a model for self-assembly in which the RNA interacts with trimers of protein to form an icosahedral scaffolding in which the gaps are filled by the other 120

subunits which have nonspecific contacts with the RNA (see [Capsids, Viral](#)). Another example of the application of low-resolution neutron diffraction and contrast variation is in [membrane protein](#) crystallography, where it has been used to locate and orient [detergent](#) molecules with respect to the protein ([20](#), [21](#)). These experiments take advantage of the inherent high contrast between detergent and  $^2\text{H}_2\text{O}$ ; and they provide information on the role of detergent in membrane protein crystallization, as well as the potential lipid–protein interactions.

## Bibliography

1. J. Wadzack et al. (1997) *J. Mol. Biol.* **266**, 343–356.
2. R. Willumeit et al. (1996) *J. Mol. Struct.* **383**, 201–211.
3. J. Zhao and H. B. Stuhmann (1993) *J. Phys. IV* **3**, 233–236.
4. K. Ibel and H. B. Stuhmann (1975) *J. Mol. Biol.* **93**, 255–265.
5. H. L. Crespi (1977) *Stable Isot. Life. Sci., Proc. Tech. Comm. Meet. Mod. Trends Biol. Appl. Stable Isot.*, IAEA, Vienna, pp. 111–121.
6. J. F. Pardon et al. (1975) *Nucleic Acids Res.* **2**, 2163–2176.
7. P. Suau et al. (1977) *Nucleic Acids Res.* **4**, 3769–3786.
8. S. Cusack, R. W. H. Ruigrok, P. C. J. Krygsman, and J. E. Mellema (1985) *J. Mol. Biol.* **186**, 565–582.
9. P. Stöckel et al. (1980) *Eur. J. Biochem.* **112**, 411–417 and 419–423.
10. H. Lederer, K. Mortensen, R. P. May, G. Baer, H. L. Crespi, and H. Heumann (1991) *J. Mol. Biol.* **219**, 747–755.
11. M. S. Capel, D. M. Engelman, B. R. Freeborn, M. Kjeldgaard, J. A. Langer, V. Ramakrishnan, D. G. Schindler, D. K. Scheider, B. P. Schoenborn, I. Y. Sillers, S. Yabuki, and P. B. Moore (1987) *Science* **238**, 1403–1406.
12. R. P. May, V. Nowotny, P. Nowotny, H. Voß, and K. H. Nierhaus (1992) *EMBO J.* **11**, 373–378.
13. J. K. Krueger, G. Zhi, J. T. Stull, and J. Trewhella (1998) *Biochemistry* **37**, 13997–14004.
14. G. A. Olah and J. Trewhella (1994) *Biochemistry* **33**, 12800–12806.
15. O. Herzberg and M. N. G. James (1985) *Nature (London)* **313**, 653–659.
16. M. Sundaralingham et al. (1985) *Science* **227**, 945–948.
17. M. Roth, A. Lewitt-Bentley, and G. A. Bentley (1984) *J. Appl. Crystallogr.* **17**, 77–84.
18. M. Roth (1991) In *Int. Union Crystallogr., Crystallogr. Symp.* **5** (Crystallogr. Comput. **5**) (D. M. Moras, and A. D. Podjarny, eds.), Oxford University Press, Oxford, U.K., pp. 229–248.
19. P. Timmins, D. Wild, and J. Witz (1994) *Structure* **2**, 1191.
20. M. Roth et al. (1989) *Nature* **340**, 659–662.
21. P. Timmins, E. Pebay-Peyroula, and W. Welte (1994) *Biophys. Chem.* **53**, 27–36.

## Suggestions for Further Reading

22. L. A. Feigen and D. I. Svergun (1987) *Structure Analysis by Small-Angle X-Ray and Neutron Scattering*, Plenum Press, New York.
23. B. Jacrot (1976) *Rep. Prog. Phys.* **39**, 911–953.
24. P. B. Moore (1982) "Small-Angle Scattering Techniques for the Study of Biological Macromolecules and Macromolecular Aggregates", In *Methods of Experimental Physics*, Vol. **20** (G. Ehrenstein and H. Lecar, eds.), Academic Press, New York, pp. 337–390.
25. B. P. Schoenborn and R. B. Knott (eds.) (1996) *Basic Life Sciences*, Vol. **64**: Neutrons in Biology, Plenum Press, New York.
26. H. B. Stuhmann (1987) "Molecular Biology". In *Methods of Experimental Physics*, **23**, Part C, Academic Press, New York.

## Controlling Element

*Controlling element* is the term used by Barbara McClintock (1) to describe the new genetic elements she found in the late 1940s that “can modify and control the action of the genes themselves” and “may move from location to location within the chromosome complement without losing its identity.” We now know such elements as [transposable elements](#)—that is, discrete pieces of DNA that can move from place to place within a [genome](#). Such elements are extremely widespread and have been found in virtually all organisms that have been examined. It is interesting that she discovered these elements that we now know to be mobile DNA before the structure of DNA was established in the early 1950s. The first molecular characterization of transposable elements came nearly 20 years later, after their discovery in bacteria, where DNAs containing them could be isolated and characterized (2).

McClintock was a maize geneticist and cytogeneticist. Her discovery of controlling elements resulted from the study of irregular patterns of pigment in maize kernels (ie, variegations of pigment patterns) that arose in some ears. Such variegation results in changes in pigment gene expression during kernel development. We now know that such variegation reflects the insertion of a transposable element to inactivate a gene and the subsequent excision of the element at a later time, restoring gene function. Having earlier studied chromosome breaks that resulted from unusual chromosome fusions, she was also quick to appreciate that some variegation resulted from chromosomal breakage at a particular site. She proposed that there was an element *dissociation* (Ds) at the breakage site and suggested that variegation results from the loss of genes downstream of this site following chromosome breakage.

Identification of an element that could cause chromosome breaks was novel, but even more dramatic was her finding that the element could move to another chromosomal site and again cause instability via chromosome breakage. McClintock realized that another element in addition to Ds was required to cause breaks at Ds and also to promote the translocation of Ds to another chromosomal position—that is, that another element encodes a product that promotes the movement of Ds. She named this other element *activator* (Ac) and also established that Ac could move from place to place.

We now know that Ac is an intact (autonomous) element (about 2.9 kbp) encoding a [recombinase](#), a [transposase](#), that acts on special sequences at the tips of the element to promote its translocation; Ds is a deleted version of Ac that lacks the transposase but still has the special terminal recombination sequences, so that Ac transposase can promote movement of Ds. Ds is called a **nonautonomous element** because it requires the presence of transposase provided by another element.

### Bibliography

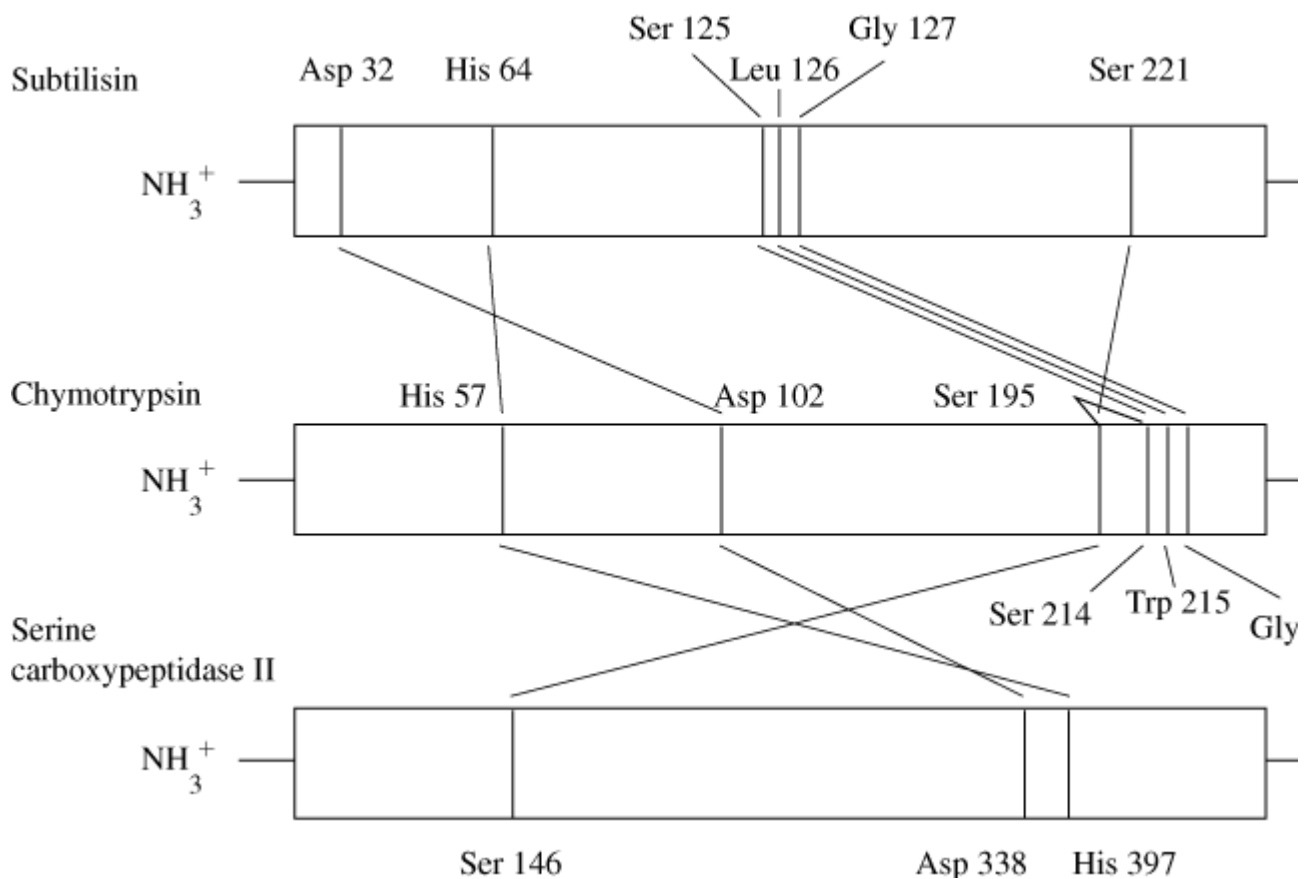
1. B. McClintock (1956) Cold Spring Harbor Symp. Quant. Biol. **21**, 197–216.
2. P. Starlinger (1977) In *DNA Insertion Elements, Plasmids, and Episomes* (A. I. Bukhari, J. A. Shapiro, and S. L. Adhya, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 25–30.

## Convergent Evolution

Convergence is defined as an evolutionary event in which two morphological or molecular traits become similar during [evolution](#) due to a similarity in environment or selection pressure, even though these traits have independent ancestors or origins. In the case of amino acid or **nucleotide sequences**, there is little known evidence for convergence; the probability of two sequences becoming similar by convergence seems vanishingly small, and any similarities in sequence are attributed to [divergent evolution](#). Yet even the absence of sequence similarity cannot be taken as evidence of convergence, as this will also result from extensive divergence between two sequences that originated from a common ancestor. In many such cases where proteins related functionally have no detectable sequence [homology](#), their three-dimensional [protein structures](#) are very similar; protein structure appears to change more slowly than does the [primary structure](#), and such cases seem to be further examples of divergent evolution.

There are, however, some examples of convergence that are apparent in the protein structure and the [active sites](#) of [enzymes](#). For example, there are several kinds of [serine proteinases](#). One family of these proteinases, represented by **chymotrypsin**, has very similar [tertiary structures](#) and active sites, containing a [catalytic triad](#) of serine, histidine, and aspartic acid residues. The serine proteinase [subtilisin](#) also contains an extremely similar triad structure, even though the rest of its [tertiary structure](#) and sequence is unrelated to the chymotrypsin family (Fig. 1) (1). Moreover, the three residues constituting the catalytic triad occur in different orders in the primary structures, making it virtually inconceivable that the two families of proteins arose from a common ancestor. Furthermore, serine carboxypeptidase II also contains a catalytic triad at its active site, even though its catalytic mechanism is entirely different. It is considered that the very similar triad structures of these different protein families emerged from their respective and independent ancestors by positive [natural selection](#). This is one of the best characterized examples of convergent evolution.

**Figure 1.** Convergent evolution of the active site residues in subtilisin, chymotrypsin, and serine carboxypeptidase II. Although relative positions of active-site residues differ in the amino acid sequences among subtilisin, chymotrypsin, and serine carboxypeptidase II, they have a common catalytic triad consisting of Ser 221, His 64, and Asp 32 in subtilisin; of Ser 19 and Asp102 in chymotrypsin; and of Ser 146, His 397, and Asp 338 in serine carboxypeptidase II. Thus, they are considered examples of convergence or convergent evolution (2).



There are other examples of enzymes with the same functions but sufficient differences to make it clear that they have not arisen from the same ancestor, but that the same function has arisen independently by a type of convergent evolution (Table 1) (2).

**Table 1. Some Enzymes that Have Evolved Independently on More Than One Occasion, by Convergent Evolution**

---

| Superoxide dismutases              |
|------------------------------------|
| Adolases                           |
| Sugar kinases                      |
| Serine proteinases                 |
| Alcohol dehydrogenases             |
| Aminoacyl tRNA synthetases         |
| Ribonucleotide reductases          |
| Topoisomerases                     |
| Phosphoenolpyruvate carboxykinases |
| Malate dehydrogenases              |

---

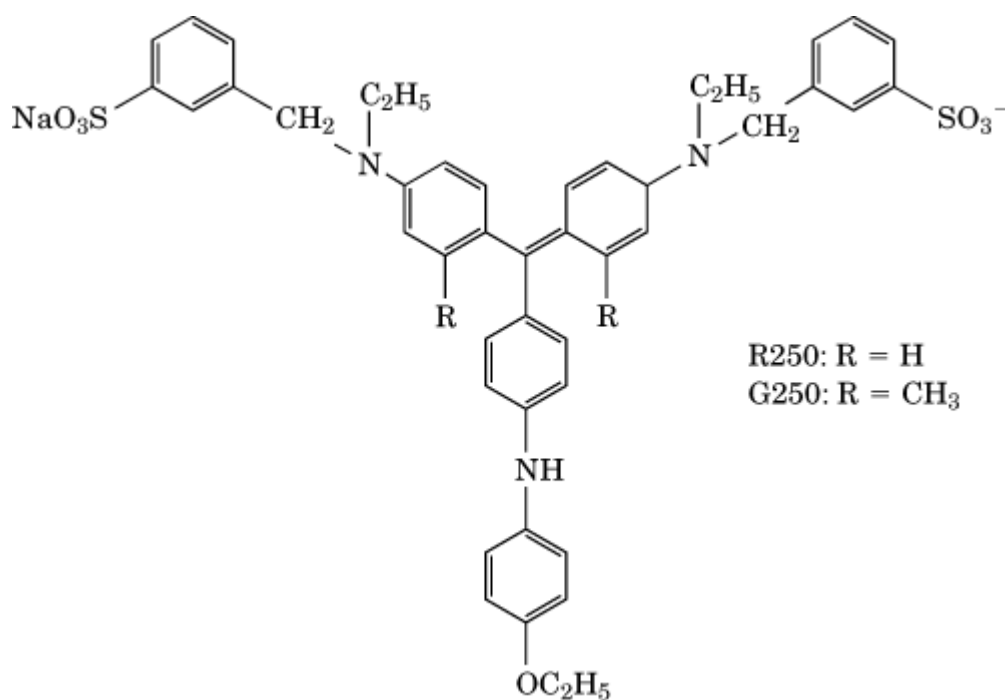
## Bibliography

1. D. Voet and J. G. Voet (1995) *Biochemistry*, 2nd ed., Wiley, New York.
2. A. Doolittle (1994) *Trends. Biochem. Sci.* **19**, 15–18.

## Coomassie Brilliant Blue

Proteins are most frequently detected after [gel electrophoresis](#) by fixing them in the gel, so that they do not diffuse further, and then staining them to produce a colored zone. The most sensitive stain is [silver stain](#), but the stain used most frequently is Coomassie Brilliant Blue. It has two frequently encountered forms, R250 and G250, whose structures are depicted in Figure 1. Although the two dyes are structurally very similar, they require different physical and staining procedures and are not interchangeable in any particular protocol. These two dyes do not react chemically with proteins but merely form noncovalent complexes. The interaction with proteins is believed to be primarily ionic and involves the acidic sulfonate groups on the dye and basic groups on the protein, but nonpolar **van der Waals** forces are probably also involved. Consequently, the dyes do not bind equally to all proteins, so two electrophoretic bands with the same blue color need not contain the same amount of protein (1).

**Figure 1.** The chemical structures of the R250 and G250 forms of Coomassie Brilliant Blue. They differ only in the nature of the group R.



Electrophoresis gels are readily stained by placing the gel in a liquid containing the Coomassie dye plus an agent to fix the protein bands, usually acetic acid plus methanol; a solution of 10% (w/v)



**trichloroacetic acid** plus 10%(w/v) sulfosalicylic acid is very effective for fixing and staining. Fixed and insoluble proteins bind the dye tightly. It is usually necessary to remove the excess dye from the staining mixture before the bands on the gel can be seen. This is usually accomplished simply by washing the gel in the fixative solution, but it can also be accomplished by adding an agent, such as polyurethane foam, that binds the excess dye tightly. It is important not to remove all of the excess dye from the solution because then the dye bound to the protein bands dissociates and the gel becomes bleached. The procedure is simplified if the Coomassie dye is only slightly soluble in the original staining mixture; the protein takes up the dye from the solution, but the background color remains weak.

A number of other dyes, such as [Ponceau S](#) and Amido black, can be used in the same way, but they are not as sensitive as the Coomassie dyes. Coomassie Brilliant Blue detects approximately 0.1 µg of protein in a band on a [polyacrylamide](#) gel.

Coomassie Brilliant Blue, particularly G250, is also used to quantify the amount of protein in solution, which is known as the Bradford assay (2). The dye complexed to protein has an altered **absorbance** spectrum. Under certain acidic conditions, the absorbance maximum shifts from 465 to 595 nm upon binding to a protein. Even though different proteins give somewhat different responses in this assay, its simplicity makes it widely used.

#### Bibliography

1. M. Tal, A. Silberstein, and E. Nusser (1980) *J. Biol. Chem.* **260**, 9976–9980.
2. M. M. Bradford (1976) *Anal. Biochem.* **72**, 248–254.

## Cordycepin

Cordycepin is 3'-deoxyadenosine, but is normally used in molecular biology as the triphosphate form. 3'-Deoxyadenosine triphosphate is a chain-terminating nucleotide analogue of ATP that can be incorporated into an RNA by [polyadenylate polymerase](#). This inhibits further elongation of the [poly A](#) tail due to the absence of the 3'-hydroxyl group. The cordycepin triphosphate has been used primarily to study the [polyadenylation](#) reaction *in vitro*, where it has allowed this reaction to be resolved into two distinct steps: cleavage of the [messenger RNA](#) precursor and the addition of the poly A tail (see [Polyadenylation](#)) (1, 2).

#### Bibliography

1. M. Sheets, P. Stephenson, and M. Wickens (1987) *Mol. Cell Biol.* **7**, 1518–1529.
2. C. Moore, H. Skolnik-David, and P. Sharp (1986) *EMBO J.* **5**, 1929–1938.

#### Suggestion for Further Reading

3. W. Keller and L. Minvielle-Sebastia (1997) A comparison of mammalian and yeast pre-mRNA 3'-end processing, *Curr. Opin. Cell Biol.*, **9**, 329–336.

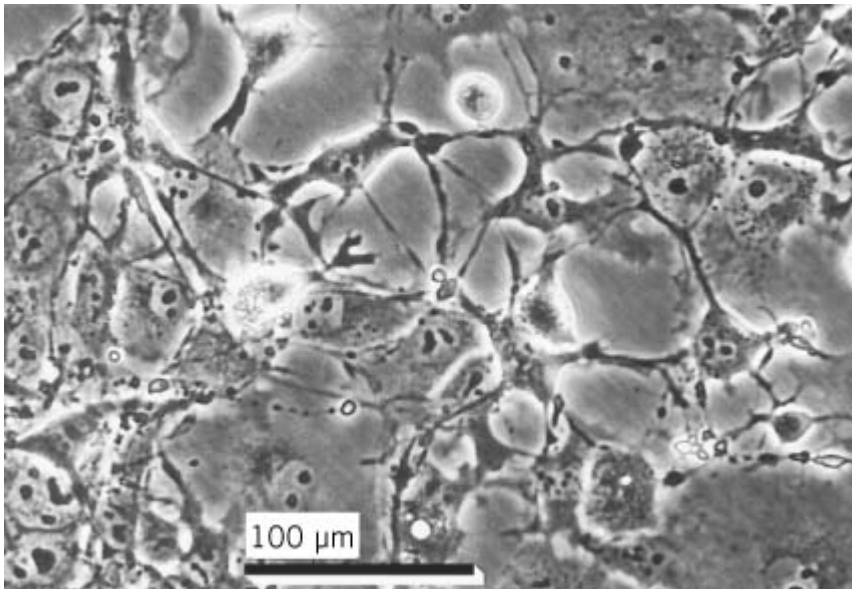
## Corepressor

A corepressor is a molecule that binds to a sequence-specific DNA binding protein and enables that protein to direct transcriptional repression. The term is used to apply to two types of mechanism. In the first case, a small molecule binds to and effects a conformational change in a DNA-binding protein, thus, enabling binding to a specific DNA sequence. One example of this mode of action is the binding of tryptophan to the tryptophan repressor. The second type of corepressor is a protein that forms a specific complex with a sequence-specific DNA binding protein and then itself inhibits transcription either directly or indirectly. One example of such an activity is the complex formed in eukaryotes between histone deacetylases and DNA-targeted regulators via an adaptor protein.

## COS Cells

COS cells are monkey kidney fibroblasts (Fig. 1); they have become popular in recent years because they support certain strains of replication-deficient animal **viruses** (1, 2) and can be **transfected** easily (3, 4). Two lines are in common use, COS-7 (ATCC CRL-1651) and COS-1 (ATCC CRL-1650).

**Figure 1.** Confluent culture of COS-7 cells. Phase contrast, Olympus CK microscope, 20× objective.



### 1. Origin

Both cell lines were derived from CV-1 cells (ATCC CCL-70), an established [cell line](#) of simian fibroblasts, by transformation with an origin-defective mutant of **SV40 virus** (1). Both are **T-antigen**-positive, permit lytic growth of SV40, and support the replication of temperature-sensitive mutant SV40 tsA209 and of SV40 mutants with deletions in the early region. COS-1 cells contain a

single integrated copy of the complete early region of the SV40 [genome](#).

## 2. Properties

COS cells are propagated as a monolayer in the alpha modification of Eagle's minimal essential medium (α-MEM) supplemented with 10% fetal bovine serum (see [Serum Dependence](#)) or in Dulbecco's modification of Eagle's medium. They grow rapidly, with an approximate population doubling time of 18 h.

## 3. Usage

COS cells have been used extensively in DNA transfer studies ([5-7](#)) and have been transfected by a number of different methods, including calcium phosphate ([8](#)), DEAE ([4](#)), and lipid transfection ([9](#)).

## Bibliography

1. Y. Gluzman (1981) SV40-transformed simian cells support the replication of early SV40 mutants. *Cell* **23**, 175–182.
2. F. Unckell, R. E. Streeck, and M. Sapp (1997) *J. Virol.* **71**, 2934–2939.
3. Y. Cao, D. M. Stafforini, G. A. Zimmerman, T. M. McIntyre, and S. M. Prescott (1997) *J. Biol. Chem.* **273**, 4012–4020.
4. J. J. Schwartz and R. Rosenberg (1998) in *DNA Transfer to Cultured Cells*, K. Ravid and R. I. Freshney, eds., Wiley-Liss, New York, pp. 179–192.
5. V. Misra, H. J. Klamut, and A. M. Rauth (1998) Transfection of COS-1 cells with DT-diaphorase cDNA: Role of base change at position 609. *Br. J. Cancer* **77**, 1236–1240.
6. S. G. Plonk, S. K. Park, and J. H. Exton (1998) *J. Biol. Chem.* **273**, 4823–4826.
7. B. Corsi, F. Perrone, M. Bourgeois, C. Beaumont, M. C. Panzeri, A. Cozzi, R. Sangregorio, P. Santambrogio, A. Albertini, P. Arosio, and S. Levi (1998) *Biochem. J.* **330**, 315–320.
8. J. V. O'Mahoney and T. E. Adams (1998) in *DNA Transfer to Cultured Cells*, K. Ravid and R. I. Freshney, eds., Wiley-Liss, New York, pp. 125–145.
9. V. V. Bickko (1998) In *DNA Transfer to Cultured Cells* (K. Ravid and R. I. Freshney, eds., Wiley-Liss, New York, pp. 193–212.

## COSY Spectrum

Correlation spectroscopy (COSY) was among the first two-dimensional (2D) nuclear magnetic resonance (NMR) experiments to be developed. In COSY-type experiments, cross peaks arise in a 2D map at the intersections of two [chemical shifts](#) if a resolved spin-coupling interaction exists between the nuclei characterized by the two shifts. The cross peaks have considerable fine structure, consisting of positive and negative signals if the lines of the spectrum are sufficiently narrow.

The appearance of the cross peaks in a COSY-type experiment relates to a quantum-mechanical phenomenon known as coherence transfer. In NMR, a coherence is characterized by a specific Larmor precessional frequency that is essentially defined by the magnetic field of the spectrometer and the chemical shifts of the spins involved in the creation of the coherence. A cross peak in COSY-type experiments develops because magnetization that was initially characterized by one precessional frequency is converted at some point to a coherence that precesses at a different

frequency, defined by another chemical shift. The COSY 2D experiment in effect measures the precessional frequency of the coherence during both parts of the experiment. Coherence transfer occurs only if spins are J-coupled to each other. In peptide and protein systems, the appearance of a cross peak in proton-proton COSY experiments generally identifies sets of protons that are adjacent to each other in the covalent structure. With sufficiently narrow lines, it is possible to analyze the structure of COSY cross peaks to obtain values for the spin-spin coupling constant involved in the coupling interaction signaled by the cross peak.

Many variations on the basic ideas of the COSY experiment are now available, including DQFCOSY (double quantum filtered COSY). The advantages of the DQFCOSY experiment include a much narrower set of diagonal peaks (so that cross peaks close to the diagonal can be detected more readily) and drastic reduction of the intensities of singlets. Such singlets tend to be very intense and can cause noise and various artifacts in a 2D spectrum.

Note the requirement that a resolved spin-coupling interaction be present for optimum detection of a COSY-type cross peak. With greater molecular weights of samples, proton NMR spectral lines increase in width. Under these conditions, the positive and negative components of COSY cross peaks can begin to cancel one another, with the detection of the cross peak relative to noise becoming more difficult. Thus, although the presence of a cross peak in a COSY-type spectrum indicates mutual spin coupling between (probably) adjacent groups of protons, the absence of a cross peak in a COSY-type spectrum cannot be taken as evidence that two nuclei are not coupled to each other.

#### Suggestions for Further Reading

R. J. Abraham, J. Fisher and P. Loftus (1988) *Introduction to NMR Spectroscopy*, Wiley, New York.

*Two-dimensional NMR Spectroscopy: applications for chemists and biochemists*, 2nd ed. (1994) (W. R. Croasmun and R. M. K. Carlson, eds.), V.C.H., New York.

A. E. Derome (1987) *Modern NMR Techniques for Chemistry Research*, Pergamon, Oxford

S. W. Homans (1992) *A Dictionary of Concepts in NMR*, Clarendon, Oxford.

R. S. Macomber (1998) *A Complete Introduction to Modern NMR Spectroscopy*, Wiley, New York.

D. L. Turner (1985) *Prog. NMR Spectrosc.* **17**, 281–358.

D. E. Wemmer (1988) In *Recent Advances in Organic NMR Spectroscopy* (J. B. Lambert and R. Rittner, eds.), Norell Press, Landisville, New Jersey.

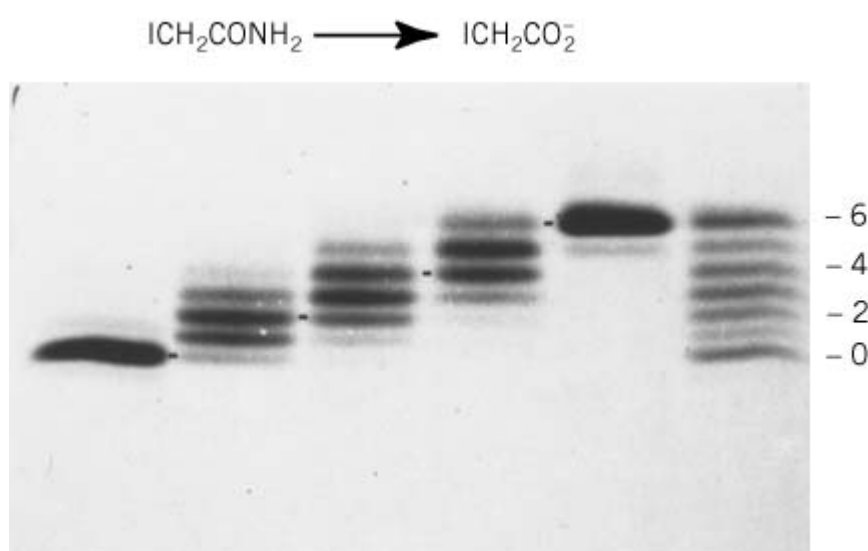
#### Counting Residues

Natural [proteins](#) have integral numbers of each of the twenty [amino acids](#) in a [polypeptide chain](#), but most currently accepted methods of determining this number by [amino acid analysis](#) yield only a nonintegral ratio of moles of amino acid per mole of protein. This value is rarely close to an integer because of experimental error and uncertainty about the **molecular weight** of the protein. The method described here determines the integral number of residues of certain amino acids, independently of any other property of the polypeptide chain, including its molecular weight. Then this information is combined with amino acid analysis to determine a more accurate value for the molecular weight and the number of other amino acids. The method also provides information about the net charge on a protein.

The general procedure is to modify the  $N$  residues of a given amino acid gradually and specifically, so as to generate a complete spectrum of molecules with 0, 1, 2, 3,  $\frac{1}{4}$ ,  $N$  of the groups modified. A single modification reaction can be used, varying the average extent of the reaction. Alternatively, one uses the competition with a related, but distinguishable reagent and carry out the modification to completion. Then the modified species generated are counted after separating them by a technique sensitive only to the number of groups modified by the particular reagent. The most useful modifications are those that alter the electric charge of a residue. Then the number of modifications introduced is determined by separating the various species by [electrophoresis](#), [isoelectric focusing](#), or [ion exchange chromatography](#). The separation must be sensitive only to the number of groups modified in charge, not which particular ones in the polypeptide chain. This is usually accomplished by carrying out the modification and the separation on unfolded molecules, where differences between the  $N$  residues are minimized. There should be  $(N + 1)$  species present.

The [thiol groups](#) of [cysteine](#) residues are reacted with a reagent like **iodoacetic acid**, which introduces an acidic group. This reaction competes with that with the closely related iodoacetamide, which is neutral and does not introduce a charged group. The reaction of the thiol groups is taken to completion, and the number of acidic groups introduced is varied by varying the ratio of iodoacetic acid to iodoacetamide (Fig. 1). The number of acidic groups introduced is determined by electrophoresis under denaturing conditions. In addition to the number of total cysteine residues or thiol groups, the number of [disulfide bonds](#) can be determined by counting the thiol groups present before and after reduction of the disulfide bonds.

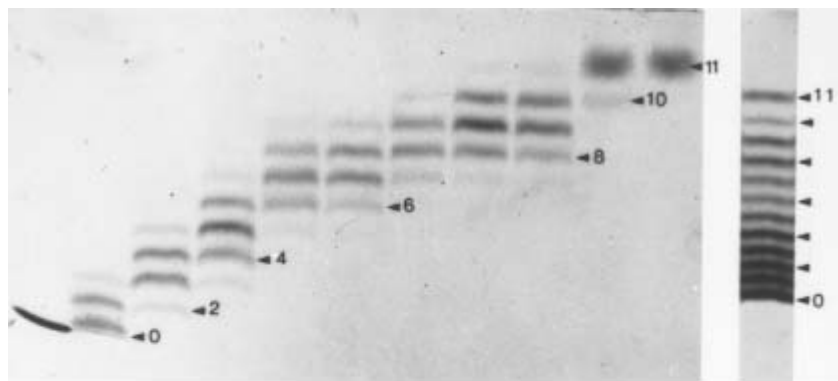
**Figure 1.** Counting the six cysteine residues of reduced [BPTI](#). The reduced protein was reacted with the neutral iodoacetamide (1st lane), acidic iodoacetic acid (5th lane), and mixtures of the two in the ratios 1:1, 1:3, and 1:9 in lanes 2, 3, and 4, respectively. Lane 6 contains a mixture of equal portions of the samples applied to lanes 1 to 5. Electrophoresis of the basic protein was in 8 M [urea](#) from top to bottom. The electrophoretic mobility is decreased in proportion to the number of acidic carboxymethyl groups on each molecule; the number for each of the electrophoretic bands is indicated at the right. The competition between the two reagents indicates that iodoacetamide reacts three times more rapidly than iodoacetic acid under the conditions used here. From Ref. 1.



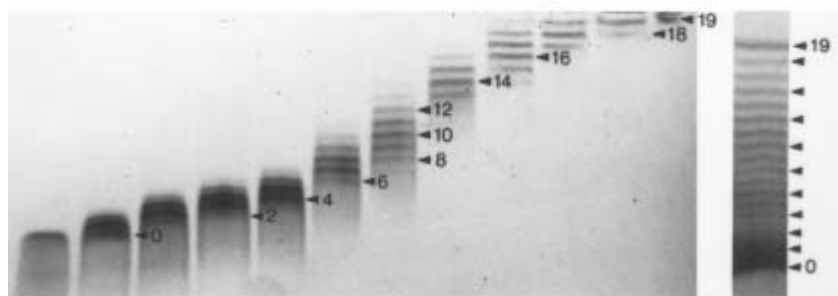
The amino groups of **lysine residues** and at the N-terminus can be modified specifically by reacting them with a reagent like succinic anhydride, which replaces a basic amino group with an acidic group. Consequently, the net charge of the protein can change by up two units for each amino group reacted. The reaction is varied in its extent by varying the amount of reagent added to the unfolded

protein (Fig. 2).

**Figure 2.** Counting amino groups of bovine **ribonuclease A** (a) and horse ferricytochrome **c** (b) by electrophoretically separating the mixtures produced by progressive succinylation. The original unmodified protein is on the left, and the degree of succinylation increases to the right. Electrophoresis was in 8 M urea at pH 3.6 from top to bottom. The separate lane on the far right is of a mixture obtained by combining all of the individual samples. In this case, electrophoresis was at the slightly lower pH of 3.45. Alternating bands are marked by arrows, and the number of succinyl groups is indicated for a few of the bands. The results confirm the presence of 11 and 19 amino groups in ribonuclease A and cytochrome *c*, respectively. From Ref. 2.



(a)



(b)

Other amino acids for which specific modifications are possible may be counted in the same way, in theory, but the specific techniques have not yet been developed.

In each of the previous procedures, the number of groups modified is counted by the number of new bands apparent by electrophoresis. Electrophoretic mobility is proportional to the net charge of the protein molecule, which is changed gradually by modifying specific residues. Therefore, the net charge on the original protein molecule can be determined by measuring the degree of modification necessary to give the protein zero electrophoretic mobility or by using stepwise modification to establish a scale for electrophoretic mobility as a function of changes in net charge.

#### Bibliography

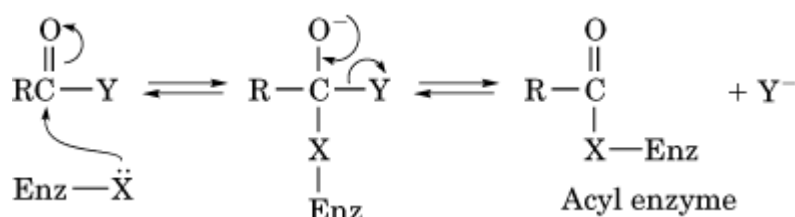
1. T. E. Creighton (1980) *Nature* **284**, 487–489.
2. M. Hollecker and T. E. Creighton (1980) *FEBS Lett.* **119**, 187–189.

### Suggestion for Further Reading

3. M. Hollecker (1997) In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.) IRL Press, Oxford, pp. 151–164.

### Covalent Catalysis

Covalent [catalysis](#) by an [enzyme](#) involves the formation, along the reaction pathway, of a covalent enzyme intermediate that can be formed and broken down at a rapid rate. Covalent bonds are formed most commonly as a result of the attack by an enzyme nucleophilic (electron-rich) group on an electrophilic (electron-deficient) moiety of the substrate that is bound at the [active site](#). The nucleophilic groups that are present on the side-chains of [amino acid](#) residues are  $\text{RCOO}^-$ ,  $\text{RNH}_2$ ,  $\text{R-OH}$ , and the nitrogen atoms of the imidazole ring of [histidine](#) residues. The electrophilic moieties of substrates may be acyl, phosphoryl, or glycosyl groups, so the covalent intermediates would be acyl-, phosphoryl-, and glycosyl-enzyme complexes, respectively. The reactions involved in the formation of an acyl-enzyme are



The second step in covalent catalysis is the attack by a low-molecular-weight nucleophile on the covalent intermediate to release the second reaction product. If the attacking nucleophile is water, the enzyme would catalyze a hydrolytic reaction that conforms to a **kinetic mechanism** with the ordered release of products. If the attacking nucleophile were a second substrate, the reaction would have a ping-pong kinetic mechanism.

Enzyme molecules are poor in electrophilic groups. But electrophilic catalysis does occur with those enzymes that contain metals (see [Metal-Requiring Enzymes](#)) or prosthetic groups (see [Coenzyme, Cofactor](#)) that act as electron sinks during catalysis.

As a large number of enzymes catalyze reactions through the formation of covalent intermediates, it appears that this type of catalysis offers some general advantages (1). It has been considered that the immobilization of a covalent intermediate should provide a significant **entropic** driving force and catalytic efficiency is increased by having multiple bond-making and bond-breaking steps occurring within a single active site.

### Bibliography

1. C. Walsh (1979) *Enzymatic Reaction Mechanisms*, W. H. Freeman and Company, San Francisco, Calif., pp. 39–41.

## CpG Islands

In some animals, cytosine bases occurring in the sequence CG are frequently methylated to [5-methylcytosine](#) (see [Methylation, DNA](#)). Yet even in species where most such residues are methylated, there are genome segments, known as *CpG islands*, which are conspicuously not methylated. These are short, dispersed regions of DNA with a high frequency of CpG dinucleotides relative to the bulk genome, but they are not methylated. In the haploid human genome, it is estimated that there are some  $30$  to  $45 \times 10^3$  CpG islands. Such regions of DNA are at least 200 bp in length and have a *G + C* content greater than 50%. Outside these islands, the frequency of CpG sequences is depleted to about only 20% of that expected on a random basis. This bias against CpG sequences that would be methylated is believed to be caused by methylcytosine bases that deaminate spontaneously 10 to 20 times more readily than normal cytosine bases. Furthermore, the product of deamination of 5-methylcytosine is thymine, a normal base, so this mutagenic event is more difficult to repair than usual (see [DNA Repair](#)).

CpG islands are associated with the 5' ends of **housekeeping genes** and of a few tissue-specific genes (highly tissue-specific genes usually lack islands). They have an open [chromatin](#) structure, and it has been postulated that they are sites of interaction between [transcription factors](#) and **promoters**. All known widely expressed genes are associated with more than one CpG island, which usually includes the [transcription](#) start site, and a few of these genes have an additional island in the 3' direction. The average length of CpG islands is about 240 bp in widely expressed genes, but there is no typical size. Most are between 200 and 1400 bp, and a majority of islands are 200 to 400 bp. The average size of islands associated with genes with limited expression is slightly smaller.

Less than half of tissue-specific expressed genes and of those with limited expression are associated with CpG islands that occur in one or more exons. This implies that CpG islands may be used to identify transcripts, since they occur in exons. Furthermore, CpG islands cover the whole or part of the promoter regions containing canonical [TATA box](#) and [CAAT box](#), and they include GC-rich promoters lacking a typical TATA box. Such promoters are considered typical of housekeeping genes.

### Suggestion for Further Reading

S. H. Cross and A. P. Bird (1995) CpG islands and genes, *Curr. Opin. Genet. Dev.* **5**, 309–314.

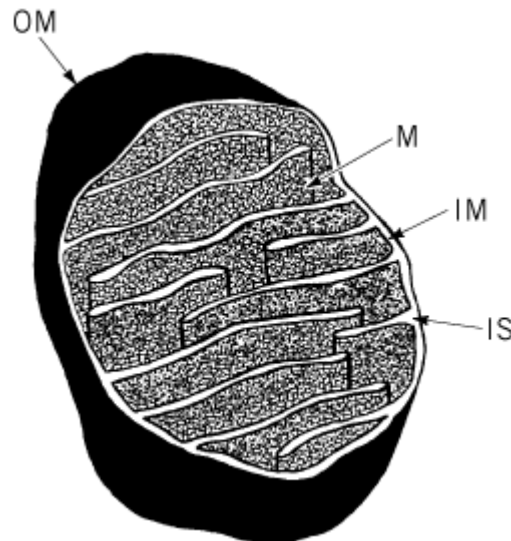
## Cristae

[Electron microscopy](#) of [mitochondria](#) demonstrates that the inner mitochondrial membrane (IM) folds towards the inner part of the mitochondrion and then folds back (Figure [1](#)), thereby forming structures that are called cristae (see [Mitochondria](#)). The same morphological evidence has shown that a relationship exists between the development of cristae and the functional activity of mitochondria, which is, in turn, related to the energy requirement of the cell. In fact, folding of the inner mitochondrial membrane is a prerequisite for the structural organization of respiratory enzyme complexes (see [ATP Synthase](#)). The IM defines the intermembrane space (IS) and the mitochondrial matrix (M), where most soluble enzymes and the **protein biosynthesis** apparatus are localized, while the respiratory enzyme complexes are localized in the IM. Recently, detailed analysis of



mitochondrial inner and outer membranes (OM), of their contacts, and of the localization of proteins in the two membranes, in the intermembrane space, and in the matrix, has led to very important developments in our understanding of the essential aspects of mitochondrial transport. Since the IM and OM separate two distinct spaces in mitochondria, proteins encoded by the cell [nucleus](#) and targeted to the mitochondria can be destined to either one of these spaces, or to be integral components of the IM or OM.

**Figure 1.** Morphology of a mitochondrion derived from electron microscopy (OM = outer membrane; IM = inner membrane; IS = intermembrane space, M = matrix).



Sorting of proteins to mitochondria is generally mediated by presequences, which are by signals composed of 15–35 amino acid residues at the amino terminus (see [Signal Peptide](#)). Most proteins to be targeted to the mitochondrion are therefore synthesized on cytoplasmic [ribosomes](#) as precursors and then transferred to **receptor** proteins on the mitochondrial surface. The pre-sequences are then inserted across both membranes at contact points, through distinct outer and inner protein-conducting channels. Many components of this complex machinery for importing proteins have recently been identified ([1](#)).

Transport requires unfolding and refolding of proteins by [chaperonins](#) and is followed by removal of the presequence by a [proteinase](#) in the mitochondrial matrix. [Signal peptides](#) are needed not only for routing proteins to mitochondria but also to reach the correct mitochondrial location. Proteins may actually remain in the matrix, assemble with other proteins in the IM, or be transported to the IS through the IM. An inner membrane-associated proteinase is then used to cleave a second signal peptide from these proteins. However, this very interesting scheme is not general, since some proteins lacking a signal peptide can also be imported into mitochondria.

#### Bibliography

1. T. Lithgow, B. S. Glick, and G. Schatz (1995) Trends Biochem. Sci. **20**, 98–101.

#### Critical Micelle Concentration

Many properties of a [detergent](#) in aqueous solution depend on its critical micelle concentration (cmc), which in turn is determined by the chemical structure of the detergent. In aqueous solution, detergents remain either in aggregated (micellar) or monomeric states, and the cmc is the maximum concentration at which the detergent molecules remain as monomers. An increase in detergent concentration above the cmc leads to the formation of micelles or mixed micelles with the solubilized lipid and protein molecules (see [Detergents](#)). Although the formation of mixed micelles of detergents with lipid and protein molecules is an important requirement in detergent solubilization of plasma [membrane proteins](#) and lipids, a low cmc may not be essential for efficient membrane solubilization because the cmc value alone does not always reflect a detergent's ability to dissociate membranes. For example, sodium dodecyl sulfate ([SDS](#)) is strongly ionic and has a cmc of ~0.23% at 25°C, whereas MEGA-10 is [polar](#) but nonionic and has virtually the same cmc of ~0.22%. Yet SDS is probably the strongest **denaturing** detergent known, whereas MEGA-10, just like MEGA-8, is much less disruptive to biomembranes ([1](#), [2](#)). A very similar picture emerges from comparing sodium cholate, which is strongly ionic and denaturing, with CHAPS, which is moderately zwitterionic and much less denaturing, even though both have a cmc of ~0.5% ([2](#)). Because the cmc values are also the result of certain more fundamental factors, such as structural features and ionic properties of the detergents (ie, nonionic, moderately ionic, or strongly ionic, and the membrane-dissociating property of a detergent is directly related to such factors), one effective way of classifying detergents is based on their structure and ionizability, rather than their cmc values (see [Detergents](#)).

## Bibliography

1. M. Hanatani, K. Nishifuji, M. Futai, and T. Tsuchiya (1984) *J. Biochem.* **95**, 1349–1353.
2. A. C. Newby, A. Chrambach, and E. M. Bailyes (1982) *Techniques in Lipid and Membrane Biochemistry*, B409, Elsevier/North-Holland, pp. 1–22.

## Crosslinking

*Crosslinking* refers to linking separate chemical functionalities covalently. Crosslinking is usually achieved through use of a [bifunctional crosslinking reagent](#), specifically, a compound having two reactive groups that react chemically with the two functional groups that become crosslinked. In this case, only two functional groups are usually crosslinked, but, with other types of reagents, higher order crosslinking is possible. Crosslinking can be either intramolecular (between separate functional groups on the same molecule) or intermolecular (between functional groups located on separate molecules). In the latter case, it is usually between two macromolecules or between a macromolecule and a relatively small molecule. Crosslinking is usually used to detect groups that are in much greater proximity than would be expected from just their bulk concentrations, due to the structures and associations of the molecules of which they are part. In biochemical applications, the chemical functions that can be crosslinked may be components of all classes of [macromolecules](#) ([proteins](#), **nucleic acids**, carbohydrates, and [lipids](#)). For simplicity, this article focuses on crosslinking involving the functional groups of proteins, predominantly **amino**, **thiol**, imidazole and [carboxyl groups](#) (ie, the nucleophilic side chains of certain amino acids); however, the concepts are the same for the other classes of macromolecules. It should be noted that crosslinking also occurs naturally, resulting from oxidation, radiation, or the action of certain [enzymes](#).

A distinction is sometimes made between the terms *crosslinking* and *conjugation* (or bioconjugation). In this case, crosslinking is used to denote the covalent linkage of only molecules that are associated naturally, such as the subunits of an [oligomeric protein](#) or a [hormone](#) and its **receptor**. Conjugation, in contrast, refers to the covalent linkage of molecules that lack affinity for one another. The resulting covalent complexes, referred to as conjugates, have a multitude of uses, especially in biotechnology (see [Enzyme Immobilization And Conjugation](#)). The distinction between crosslinking and conjugation is maintained in this article.

## 1. Uses of crosslinking

On the basis of the variety of the crosslinking reagents that are available and of the molecules that can be crosslinked, plus the types of experimental questions that can be answered, crosslinking is among the most versatile of biochemical techniques. In the case of proteins, for example, some information about their three-dimensional structure can be gained through identifying specific amino acid residues that are readily crosslinked intramolecularly (1). Crosslinking can also be used as a general conformational probe to detect structural changes in a protein, such as might be induced by **ligand binding** (2). Inasmuch as crosslinking can hinder conformational transitions and the dissociation of oligomers, it can also be used to stabilize the **tertiary** or [quaternary structures](#) of proteins. Thus, it can “lock” a protein into a particular functional state (3) or stabilize an oligomeric protein for subsequent physical studies, such as [electron microscopy](#). A careful analysis of the mass and composition of crosslinked complexes allows a minimal subunit stoichiometry to be determined; this is especially useful for insoluble oligomers, such as those found in [membranes](#) (4). Information concerning the symmetry of oligomeric proteins can also be obtained from crosslinking (5, 6). One of the most common uses of crosslinking is to identify neighboring proteins in a complex (7), and this can be extended to include determination of the specific regions of each that are crosslinked (8). By varying the length of the crosslinker, one can also obtain a limit for the maximum distance that can separate neighboring proteins (6). In addition to giving both relative and absolute structural information, crosslinking can be used for general screening. For example, one can identify the receptor for a particular ligand, such as a hormone (9).

It must be emphasized that for all the above applications, interpretations can be made only for those crosslinked complexes that are generated. The absence of crosslinking does not mean that two molecules do not interact, only that they are not crosslinked under the specific set of experimental conditions chosen, which includes the crosslinking reagent used. Because crosslinking requires the proper positioning of those functional groups having the appropriate chemical reactivities, one cannot expect all interacting molecules to be crosslinked by any given reagent.

## 2. Crosslinking Procedures

Crosslinking, like chemical modification of proteins, is an empirical procedure, and both processes are influenced by the same variables, namely temperature, pH, concentrations of reactants, and time. In crosslinking, however, the influence of the latter two variables is more complex. Because nucleophilic side chains of amino acids are considerably more reactive when deprotonated, the pH of the reaction is important in controlling the extent of crosslinking and the selectivity for particular side chains. Regarding concentrations, the amounts of both the crosslinker and the proteins can influence crosslinking, with low concentrations of protein and crosslinker generally favoring intramolecular crosslinking. Adding the crosslinker over time, instead of as a single addition, favors a greater extent of crosslinking. Differences in selectivity of the reactive groups of bifunctional reagents, especially heterobifunctional crosslinkers, for protein functional groups can be maximized using a two-step reaction. In this procedure, one protein is first derivatized with the crosslinking reagent and then purified to remove excess free reagent before being exposed to the second protein. This two-step procedure, which allows for each of the two reactions to be performed under different conditions, is particularly useful for [enzyme immobilization and conjugation](#), when the two molecules do not normally associate with each other.

Following crosslinking, the covalent complexes formed are commonly identified as new [electrophoresis](#) bands in sodium dodecyl sulfate–polyacrylamide gel electrophoresis ([SDS-PAGE](#)), with altered molecular weights, although care must be taken, because intramolecular crosslinking can alter the migration of individual polypeptides with this technique. Alternatively, in simple systems, covalently linked complexes can be identified as components with increased **Stokes radii** in [size exclusion chromatography](#). If either the crosslinker or the crosslinked protein has identifiable physical or chemical characteristics, these can aid in the identification of complexes. Once crosslinked complexes have been identified, complete analysis of their composition and stoichiometry can still be challenging, especially when multiple proteins are present. One useful method for identifying the proteins composing a complex is to use a cleavable crosslinker so that the original components can be regenerated from the fractionated complex ([7](#)). Another useful method is to perform [Western blots](#) on the fractionated complex using antibodies against potential components ([2](#)). Once the components have been identified, their stoichiometry in the complex is deduced from its overall mass. The ultimate step in analysis is to determine directly the actual amino acid residues that are crosslinked. Using conventional [protein sequencing](#) by [Edman Degradation](#), the sequences of the two crosslinked peptides are obtained simultaneously, with the crosslinked residues represented by gaps in the sequences; however, identification of crosslinked peptides by [mass spectrometry](#) is becoming increasingly common.

### Bibliography

1. F. C. Hartman and F. Wold (1967) *Biochemistry* **6**, 2439–2448.
2. O. W. Nadeau, D. B. Sacks, and G. M. Carlson (1997) *J. Biol. Chem.* **272**, 26296–26201.
3. R. E. Benesch and S. Kwong (1991) *J. Prot. Chem.* **10**, 503–510.
4. E. Heymann and R. Montlein (1980) *Biochem. Biophys. Res. Commun.* **95**, 577–582.
5. F. Hucho, H. Mullner, and H. Sund (1975) *Eur. J. Biochem.* **59**, 79–87.
6. J. Hajdu, S. R. Wyss, and H. Aebi (1977) *Eur. J. Biochem.* **80**, 199–207.
7. D. Dey, D. E. Bochkariov, G. G. Jokhadze, and R. R. Traut (1998) *J. Biol. Chem.* **273**, 1670–1676.
8. V. Rossi, C. Gaboriaud, M. Lacroix, J. Ulrich, J. C. Fontecilla-Camps, J. Gagnon, and G. J. Arlaud (1995) *Biochemistry* **34**, 7311–7321.
9. T. H. Ji (1977) *J. Biol. Chem.* **252**, 1566–1570.

### Suggestions for Further Reading

10. G. Mattson, E. Conklin, S. Desai, G. Nielander, M. D. Savage, and S. Morgenson (1993) A practical approach to crosslinking, *Mol. Biol. Rep.* **17**, 167–183.
11. S. S. Wong (1993) *Chemistry of Protein Conjugation and Cross-Linking*, CRC Press, Boca Raton, FL. (An outstanding book describing all aspects of the chemistry and uses of bifunctional reagents in crosslinking and conjugation, with an extensive bibliography for each chapter.)

### Cruciform (Holliday Junction)

A key intermediate in DNA [recombination](#) is the **Holliday junction** that involves the joining of four homologous DNA strands into a cruciform arrangement, as visualized by [electron microscopy](#) some years ago. **Supercoiling** also can induce the formation of the cruciform for [inverted repeat](#) sequences. Detailed structural analysis of the cruciform [DNA structure](#) became possible when specific sequences were designed to form a so-called immobile junction. Studies on those molecules

using [gel electrophoresis](#) and **fluorescence** indicated that the structure has a stacked X shape (ie, not tetrahedral) in the presence of metal ions (eg,  $Mg^{+}$  or  $Co(NH_3)_6^{3+}$ ). Models of the detailed structure at the junctions have been proposed. One such model involves a right-handed, antiparallel X structure, with two stacked helices crossing over each other with a crossover angle of  $\sim 60^\circ$  (or  $120^\circ$ ). Within each helix of the X structure, the two arms are collinear and the junction site is folded. A number of [B-DNA](#) structures determined by [X-ray crystallography](#) showed helix–helix packing contacts in the crystal lattice that are remarkably similar to the X structure. Recently the crystal structures of two proteins that are involved in the recombination process have been solved in the presence of their DNA substrate. The crystal structure at 2.4 Å resolution of Cre recombinase bound to a loxP DNA substrate showed that four recombinases and two loxP sites form a synapsed structure in which the bound DNA adopts a conformation similar to the four-way Holliday junction ([1](#)). The synapse assembly has a pseudofourfold symmetry, and the junction core is less compact than that proposed for the “free” Holliday junction.

## Bibliography

1. F. Guo, D. N. Gopaul, and G. D. Van Duyne (1997) *Nature* **389**, 40–46.

## Cryoelectron Microscopy

Cryoelectron microscopy can be divided into two classes, depending on whether the specimen is *hydrated* or *dry*. [Table 1](#) provides a brief description of the possible methods of specimen preparation. In each method, the specimen is initially frozen and can subsequently be observed in a cryoelectron microscope, or can be further preserved by drying or making a replica for observation in a conventional electron microscope.

**Table 1. Methods of Specimen Preparation in Cryoelectron Microscopy**

| Frozen Hydrated Specimens  | Dry Specimen Previously Frozen   |
|--|--|
| <p>1. <i>Vitrified film of ice</i><br/>A thin layer of macromolecules in aqueous suspension is rapidly frozen by plunging in a cryogen, eg, ethane, and transferred to a cryoelectron microscope.</p>  | <p>1. <i>Freeze drying</i><br/>Cells or tissues are frozen, then lyophilized. The remaining material is dried and stabilized for transfer to a conventional electron microscope. Alternatively, a replica is made.</p> |
| <p>2. <i>Frozen section</i><br/>Cells or tissues are rapidly frozen by contact with a cold metal block, sometimes under high pressure, then sliced into thin sections at <math>\sim -160\text{ C}^\circ</math>.<br/>Subsequently, sections are transferred to a cryoelectron</p> | <p>2. <i>Freeze fracture</i><br/>Cells or tissues are frozen, then fractured. A replica of the fractured surface is made.</p>  |
|  | <p>3. <i>Freeze etching</i></p>  |

microscope.

Freeze fracture followed by freeze drying.

#### 4. *Freeze substitution*

Tissues are frozen, then the ice is replaced at low temperature by an organic solvent. Conventional plastic embedment follows.

#### 5. *Cryoimmunohistochemistry*

Tissues are first chemically fixed, then suffused with a cryoprotectant and frozen for sectioning. Immunolabelling of the thawed sections follows.

---

The development of techniques for visualizing frozen hydrated tissues and [macromolecules](#) has been a boon for studying structures in a state that closely approximates their native condition (1). Because freezing takes place very rapidly, these techniques offer the ability to examine dynamic events with high time resolution. For example, the structure of the open-channel form of the [acetylcholine receptor](#) was elucidated from membrane crystals that were briefly activated (< 5 ms) with droplets containing acetylcholine (2). Small conformational changes triggered by acetylcholine were detected in the frozen hydrated crystals. The spatial resolution of periodic structures, eg, protein crystals, can also be high, extending to 3.3 Å (3; Perkins, unpublished results). Moreover, these structures can be examined unstained by utilizing the intrinsic contrast difference between the biological specimen and vitreous ice. Thus, the effects of solvent composition on the structure and aspects of interacting components can be examined (2).

Fundamental to examining frozen hydrated specimens is the preservation of the environment by rapid freezing (4). Rapid freezing causes the water of the specimen and its environment to be trapped in a vitreous (amorphous or glass-like) state. This vitreous state prevents damage to the specimen caused by the formation of ice crystals when the specimen is frozen more slowly. The specimens preserved in thin films of vitreous ice generally retain a high degree of structural integrity. However, it is crucial to maintain the specimen at below ~ -140° C because vitreous ice tends to crystallize above this temperature (4).

A wide range of frozen hydrated specimens has been examined in cryoelectron microscopes. Most studies have concentrated on macromolecules and their assemblies hydrated and frozen in thin films. To a lesser extent cryosections of biological tissues have also been studied (5-9), but the difficulty of freezing of bulk material without formation of ice crystals has hampered progress. The first published results obtained of macromolecules were from viruses, fairly robust particles (10). Early work produced 3-D reconstructions at relatively low resolution, ~ 2-3 nm (10-13), although certain viruses produced data to ~ 1 nm resolution (14). Important new information about viral structures was obtained even at lower resolution. Recent work has extended the resolution of complex viruses beyond 1 nm (15, 16). By combining cryoelectron microscopy with [X-ray crystallography](#), near atomic resolution detail of the overall architecture of large complexes and of the interactions between the components can be obtained. An instructive example is the 3-D structure of rhinovirus 14 complexed with Fab fragments determined at 0.4 nm resolution (17).

Although the contrast in frozen hydrated specimens is low compared to negatively stained specimens, the signal can be significantly boosted by averaging many identical particles together using Fourier averaging for crystals (2, 17) and correlational averaging for single particles as discussed below (15, 16, 18). Vitreous ice embedment is particularly useful for visualizing the specimen interior compared to negative stain, which accumulates around areas accessible to the

aqueous phase, thus usually showing only the molecular envelope. Hence, in general only those portions of integral [membrane proteins](#) extending beyond the lipid bilayer are observed in negative stain. However, exceptions to this limitation of negatively stained membrane proteins have been found ([19-22](#)). In certain instances, frozen hydrated crystals of membrane proteins do not provide resolution as high as the same crystals in negative stain ([23](#)). Vitreous ice embedment of delicate specimens, such as [chromosomes](#) or DNA fragments ([9, 24, 25](#)) has provided impressive results, since this type of specimen is commonly poorly preserved in other embedment media. Another example is the direct visualization by cryoelectron microscopy of A-, P-, and E-site transfer RNAs in the *E. coli* [ribosome](#) ([26](#)).

Cryofixation of tissues is an alternative to conventional chemical fixation of subcellular ultrastructure. Vitrification of ice can frequently be achieved to a depth of 10–20 µm when tissues are frozen by slamming them against a polished copper block cooled to the temperature of liquid nitrogen or liquid helium ([27](#)). High-pressure freezing can freeze samples in vitreous ice to depths as great as 100 µm, since the increase in volume going from water to ice is impaired by a high-pressure environment inhibiting the growth of ice crystals. Cryofixation is usually the choice for immunocytochemistry because it retains the antigenic properties of the specimen better than chemical fixation. With good vitrification, freeze-etching or freeze-substitution (Table [1](#)) can provide information about structures stabilized with a mechanism different from chemical fixation.

### Bibliography

1. J. Jaffe and R. M. Glaeser (1984) *Ultramicroscopy* **13**, 373–378.
2. P. N. T. Unwin (1995) *Nature* **373**, 37–43.
3. J. Brink, W. Chiu, and M. Dougherty (1992) *Ultramicroscopy* **46**, 229–240.
4. J. Dubochet et al. (1982) *J. Microscopy* **128**, 219–237.
5. J.-J. Chang et al. (1983) *J. Microscopy* **132**, 109–123.
6. J. Dubochet et al. (1983) *J. Bacteriol.* **155**, 381–390.
7. A. W. McDowell et al. (1983) *J. Microscopy* **131**, 1–9.
8. A. W. McDowell et al. (1984) *J. Mol. Biol.* **178**, 105–11.
9. A. W. McDowell, J. M. Smith, and J. Dubochet (1986) *EMBO J.* **5**, 1395–1402.
10. J. Adrian et al. (1984) *Nature* **308**, 32–36.
11. J. Lepault and K. Leonard (1985) *J. Mol. Biol.* **182**, 431–441.
12. F. P. Booy et al. (1985) *J. Mol. Biol.* **184**, 667–676.
13. R. H. Vogel et al. (1986) *Nature* **320**, 533–535.
14. J. Lepault (1985) *J. Microscopy* **140**, 73–80.
15. B. Bottcher, S. A. Wynne, and R. A. Crowther (1997) *Nature* **386**, 88–91.
16. J. F. Conway et al. (1997) *Nature* **386**, 91–94.
17. T. J. Smith (1996) *Nature* **383**, 350–354.
18. M. van Heel, G. Harauz, and E. V. Orlova (1992) *J. Struct. Biol.* **116**, 17–24.
19. B. Karlsson et al. (1983) *J. Mol. Biol.* **165**, 287–302.
20. B. Bottcher, P. Graber, and E. J. Boekema (1992) *Biochem. Biophys. Acta* **1100**, 125–136.
21. S. Karrasch et al. (1996) *J. Mol. Biol.* **262**, 336–348.
22. G. A. Perkins et al. (1997) *Biophys. J.* **72**, 533–544.
23. J. M. Valpuesta, R. Henderson, and T. G. Frey (1990) *J. Mol. Biol.* **214**, 237–251.
24. J. Dubochet et al. (1994) *Nature Struct. Biol.* **1**, 361–363.
25. J. Bednar et al. (1995) *J. Mol. Biol.* **254**, 579–594.
26. R. K. Agrawal et al. (1996) *Science* **271**, 1000–1002.
27. J. Dubochet (1995) *Trends Cell Biol.* **5**, 366–369.

## Crystallization

The determination of the three-dimensional structures of [proteins](#) by [X-ray crystallography](#) requires the availability of large, well-ordered crystals. These are difficult to obtain. The crystallization of any given protein requires long painstaking work, including much trial and error in the choice of environmental conditions (pH, temperature) and crystallizing agents.

The basic requirement is the availability of the protein in *very* pure form, free of all other molecules, such as traces of nucleic acids, in a single form (absence of **isoforms**), in a solution that is free of any dust particles. If liganded to any **coenzymes**, cofactors, or effectors, all molecules must be equally liganded; all molecules must be in the native state, and there must be no partial oligomer formation.

The basic approach is to prepare a protein solution at high concentration (10 to 100 mg/ml) in a crystallizing solvent at conditions at which the solution is close to saturation. Then the solution is gradually brought to supersaturation by, eg, slow variation in pH or temperature. All changes must be brought about slowly and the solution not subjected to shocks, such as mixing. This allows the formation of crystal nuclei. Amorphous aggregates must be totally avoided. Once nuclei are formed, the system is induced to crystallize further by the addition of protein to these nuclei, avoiding the formation of new ones, since the aim is to form large crystals. See the specialized literature for further details.

Typical crystallizing agents are (1) salts  $(\text{NH}_4)_2\text{SO}_4$ ,  $\text{Na}_2\text{SO}_4$ , Na citrate,  $\text{MgSO}_4$ ; (2) organic molecules [2-methyl-2,4-pentanediol (MPD), ethanol, isopropanol, acetone, dioxane, 1,3-propanediol]; (3) polyethylene glycols (molecular weight between 2,000 and 20,000). See also [Vapor Phase Crystallization](#).

### Suggestions for Further Reading

T. L. Blundell and L. N. Johnson (1976) *Protein Crystallography*, Academic Press, New York.

H. W. Wyckoff, C. H. W. Hirs, and S. N. Timasheff, eds., (1985) *Methods in Enzymology*, Chaps. "2"—"15", Academic Press, New York.

## Crystallography

Crystals are very important for molecular biology because they are required for detailed structure determination of macromolecules using [X-ray crystallography](#). The scientific study of crystals is very old. Its origin can be traced back to 1669 when Nicolaus Steno noticed that the angles between the faces in growing crystals remain constant. Individual crystals may differ in shape appreciably, but the angles between corresponding faces are always the same. A century later, René Just Haüy proposed that crystals are built from small units whose shape is a parallelepiped. An extensive review of the history of crystallography is given in Ref. [1](#). The physical properties of crystals are,



generally not equal in all directions. Their optical, magnetic, and electric properties vary with crystal direction, and this is the reason that crystals are widely used for physical measurements and in industrial processes, for instance, a quartz crystal in a clock or watch. The wealth of variable characteristics is a consequence of the regular packing of the ions, atoms, or molecules in the crystal. This internal regularity had been suggested since Haüy, but it could only be proven with the introduction of X-ray crystallography. In this technique, the crystal acts as a grating and diffracts X-rays in directions specific for the repeating distances in the crystal. In addition, it became possible to elucidate the internal structure of crystals, that is, the arrangement and structure of the molecular composition.

In the regular packing inside the crystal, three repeating vectors can be recognized: **a**, **b**, and **c**, with angles  $\alpha$ ,  $\beta$ , and  $\gamma$  between them. These three vectors define a [unit cell](#) in the crystal lattice. From morphological crystal symmetry considerations, that is, the rotational axes, mirror planes, and centers of symmetry observed in the shape of crystals, 32 combinations of these symmetry elements (point groups) are possible. They can be assigned to seven, and not more than seven, crystal systems (2): triclinic, monoclinic, orthorhombic, tetragonal, trigonal, hexagonal, and cubic. Internally, a crystal has more symmetry operators. Translational symmetry is added to the symmetry axes and the mirror planes. These operators can be combined in 230 different ways, leading to 230 [space groups](#), introduced by Schoenflies and Fedorow at the end of the nineteenth century (Ref. 1, p. 114; 2).

Because biological macromolecules have an asymmetrical structure, they can only form crystals without mirror planes or centers of symmetry. This restricts the number of space groups appreciably for these molecules.

#### Bibliography

1. P. Groth (1926) *Entwicklungsgeschichte der Mineralogischen Wissenschaften*, Julius Springer, Berlin.
2. International Union of Crystallography (1992) *International Tables for Crystallography*, Vol. A (T. Hahn, ed.) Kluwer Academic Dordrecht, Boston, London.

#### Suggestion for Further Reading

3. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists*, VCH, New York, Weinheim, Cambridge.

#### $C_0t$ Curve

A  $C_0t$  curve describes the sequential complexity of a DNA sample. It describes the [kinetics](#) of reassociation and **reannealing** of a sample of double-stranded DNA that has been (1) fragmented to small pieces, approximately 400 nucleotides, (2) denatured to single strands, and (3) then permitted to reanneal. The extent to which two complementary DNA strands reassociate is proportional to  $C_0t$ .  $C_0$  is the original molar concentration of the complementary strands of DNA, and  $t$  is the incubation time.  $C_0t_{1/2}$  is the value of  $C_0t$  required for completing half of the reannealing reaction. If the sequences of all initial DNA fragments were identical, complementary fragments encounter each other readily and the mixture reanneals readily at a low value of  $C_0t$ . At the other extreme, reannealing is very slow and occurs only at high  $C_0t$  values if the original DNA contains many

different sequences, so that the concentrations of complementary strands are very low and they encounter each other only infrequently. The value of  $C_{0t_{1/2}}$  is a measure of the sequential complexity of the original DNA.

Such studies on genomic DNA were one of the first indications that [genomes](#) often contain [repetitive DNA](#). Such DNA segments are present in higher concentrations than those segments present in only one copy per genome, and this fraction of the DNA reanneals at a correspondingly lower  $C_{0t}$  value.

The results indicate that different sequences are present in widely varying frequencies throughout the genome, and only a fraction occurs only once per genome.

## ***Cubitus Interruptus* Genes**

In *Drosophila* [development](#), the **phenotype** associated with a gain-of-function [mutation](#) of the *cubitus interruptus* (*ci*) gene was first described in 1971 by Hochman (1). Later, this gene was categorized as a segment polarity gene because the loss of *ci* function results in the loss of naked regions in the segments of larval cuticle (2). The *ci* gene encodes a **zinc-finger** protein, CI, that stimulates [transcription](#) in response to the *hedgehog* (*hh*) signal transduction cascade, which helps to establish the anterior-posterior positional information in the embryonic segments and [imaginal discs](#). The [morphogen](#) in vertebrates, [Sonic hedgehog](#) (*Shh*), induces and patterns many developmental processes, at least in part through regulating the expression of the vertebrate *ci* **homologues** *Gli1*, *Gli2*, and *Gli3*, which then transcribe a second level of *Shh* target genes. The importance of these genes in growth and development is illustrated by the fact that the Gli1 protein is amplified in glioblastoma cells.

### 1. Protein Structure of CI and Gli

Both the *ci* and *Gli* genes encode proteins with five highly homologous, tandemly repeated zinc fingers that act as **DNA-binding** domains and are distinguishable from other zinc-finger domains (3-6). The protein products of these genes are required for the transcription of *hh* target genes during [development](#). Because the *Gli* and *ci* genes respond to the same signaling pathway and share a common DNA-binding domain, they are considered to be homologues, even though there is very little conservation on the amino acid sequence level between CI and Gli proteins outside of the zinc-finger region. Three members of the *Gli* gene family (*Gli1*, *Gli2*, and *Gli3*) have been identified in mouse, chicken, and [zebrafish](#) (7, 8). Two *Gli* genes have been identified in *Caenorhabditis elegans* (9), four in [Xenopus](#) (10). They share various degrees of similarity with the zinc finger region, which represents the greatest degree of homology among the family. The zinc fingers are responsible for the DNA-binding activity of both CI and Gli proteins as demonstrated by both in vitro binding assays and in vivo experiments, in which expression of a transgene is driven by Gli binding sites (4, 5, 11, 12). The consensus DNA-binding site for CI and Gli proteins has also been characterized by isolating and comparing **genomic** DNA fragments that can bind to Gli1. The core sequences were further confirmed by DNase [footprinting](#) (13). The [X-ray crystallography](#) structure of the Gli1 zinc fingers bound to the DNA binding site has been solved (14). Zinc finger 1 does not contact the DNA, whereas the rest of the four fingers wrap around the DNA, with fingers 4 and 5 making most of the base contacts with the DNA. This structure might explain the higher homology shared by zinc fingers 3–5 among the members of the *Gli* gene family. **Bandshift assays** demonstrate that full-length Gli1, full-length Gli3, or a GST protein fused to the CI zinc finger domain translated in vitro can bind to the Gli DNA binding site **consensus sequence** (5, 11, 13). When the CI zinc-finger **domain** is fused to a heterologous activation domain, it functions as a transcriptional activator in fly imaginal discs and embryos (15). Furthermore, CI can transactivate a [reporter gene](#) driven by Gli

binding sites in yeast (12) and in fly cell lines (16). All of these experiments suggest that *ci* functions as a [transcription factor](#) that binds DNA through its zinc-finger domain. The amino acid residues of CI that are N-terminal to the zinc-finger domain are rich in alanine residues, a characteristic of transcriptional repressor domains (17). When a CI protein that is truncated C-terminal to the zinc-finger region is expressed in the anterior compartment of the fly embryonic segments, it causes downregulation of *ci* target genes (15). This result supports the idea that the N-terminus of CI has a [repressor](#) function. The C-terminus of the CI protein is highly acidic, a characteristic of transcriptional activators (18). When introduced into *Drosophila* embryos, a CI protein in which the N-terminus is deleted can function as a transcriptional activator to enhance *ci* target gene expression, indicating that the activation domain resides in the C-terminus (15).

## 2. Expression and Function of *ci*

The repeating segments of the fly embryo are established during blastoderm stage through a series of hierarchical interactions among the maternal gap and [pair-rule genes](#). The segment polarity genes specify the anterior-posterior polarity of each segment during gastrulation. Consequently, it is no surprise that most of the segment polarity genes describe interdependent signal transduction cascades. *Ci* is a member of the *hh* signal transduction pathway and encodes a zygotic gene that is not detected until early cellular blastoderm stage (stage 5). Before stage 11, expression of CI protein closely follows the pattern of *ci* transcripts. Both protein and message are detected initially on the dorsal surface of the embryo. As development progresses, the domain of *ci* transcription and protein expression extends ventrally and, by stage 10, encompasses the entire embryo. Coincident with the onset of ectodermal segmentation, both the message and the protein are restricted to the anterior compartment of each segment by the *engrailed* (*en*) gene, resulting in 15 broad, metameric repeating stripes. At stage 11, the protein expression pattern no longer coincides with the transcript pattern. Although transcripts are still uniformly distributed throughout the anterior compartment of each segment, the protein stripe within each segment is graded so that higher levels of protein are detected in the rows of cells that define the posterior edge of the anterior-posterior (A/P) boundary and lower levels are seen in the middle of the CI-expressing stripe. This pattern of expression also exists in imaginal discs, where *ci* is uniformly transcribed in the entire anterior compartment and higher levels of CI are detected along the A/P boundary (3, 19). It should be noted that the protein staining patterns were obtained with [antibody](#) made against the C-terminus of CI and detected only the full-length protein. The discrepancy found between transcript and protein levels in the segments and discs indicates that CI is subject to [post-transcriptional regulation](#).

The *ci* transcription pattern is regulated by various unidentified factors and is mediated through the regulatory sequences in the *ci* promoter. Using transgenic flies that carry the *LacZ* reporter gene under the control of different *ci* regulatory sequences, it has been shown that various elements in the *ci* promoter differentiate embryonic from disc expression (20). Although the sequences required for the repression of *ci* transcription in the posterior compartment of the embryonic segments and imaginal discs are different, they both contain binding sites for the *engrailed* protein, EN, consistent with the genetic data showing that *en* inhibits *ci* transcription in the posterior compartment. EN protein expressed in bacteria protects specific DNA sequences in the *ci* regulatory region that have been shown to mediate *ci* repression in the posterior compartment (20). In addition, EN binds the *ci* gene in [polytene chromosomes](#), and misexpression of *en* by the **heat-shock** promoter greatly reduces *ci* mRNA levels (20). Furthermore, Eaton and Kornberg (21) have shown that *ci* transcripts expand into the posterior compartment of both embryonic segments and imaginal discs in *en* mutant flies (21).

Genetic [epistasis](#) experiments place *ci* within the *hh* signaling cascade. It acts downstream from the *hh* receptor *patched* (*ptc*) and *smoothened* (*smo*), a transmembrane protein that is coupled to *ptc* and transduces the *hh* signal. *Ci* also acts upstream of the *hh* target genes (see [Hedgehog Signaling](#)). CI is an activator of transcription in the cells that receive an *hh* signal. These cells abut the posterior edge of the A/P boundary where high levels of CI are expressed. The cells expressing lower levels of CI are outside the range of *hh* activity and, in these cells, CI is a repressor of the *hh* target genes (19,

22). Recently, great progress has been made toward understanding the cellular mechanism involved in determining whether *ci* is an activator or a repressor. The first insight comes from the observation that CI exists in two forms: the 155-kDa full-length activator form and a 75-kDa proteolytic fragment (22). In the absence of an *hh* signal, a proteolytic event generates the 75-kDa form by cleaving the sequences C-terminal to the zinc-finger domain from the full-length protein. Because the resulting 75-kDa proteolytic product has both the zinc-finger DNA-binding domain and the N-terminal repressor domain, it functions as a repressor, preventing the activation of *ci* target genes. Antibodies made against the N-terminus of CI (AbN), which detect both full-length CI and the proteolytic fragment, stain the entire anterior compartment of the embryonic segments and imaginal discs. On the cellular level, the antibody staining is distributed evenly in both the nucleus and the cytoplasm (15, 22). On the other hand, full-length CI resides primarily in the cytoplasm. The levels are very low in the cells of the anterior compartment that do not receive an *hh* signal and are very high along the A/P boundary that contacts the *hh*-producing cells (19). This expression profile suggests that *hh* signaling regulates full-length CI levels post-transcriptionally. Genetic and cell culture experiments support this hypothesis. The *hh*<sup>Mrt</sup> mutation results in a dominant misexpression of *hh* in the wing disc. In these animals, high levels of CI protein are detected along the anterior wing margin, where ectopic *hh* is expressed (23). Suppression of *hh* expression by shifting a temperature-sensitive *hh* mutant (*hh*<sup>9k94</sup>) to nonpermissive temperatures causes decreased levels of CI along the A/P boundary in imaginal discs (15). The transmembrane protein PTC, the product of the *patched* gene, is the receptor for HH and, in the absence of HH binding, inhibits the *hh* signaling pathway. The binding of HH to PTC relieves the inhibition of *ptc* on *smo*, the positive transducer of *hh* signaling (see [Hedgehog Signaling](#)). In *ptc* mutants, *smo* is constitutively active, and *hh* target gene expression becomes independent of a *hh* signal (24-28). As expected, high levels of full-length CI are distributed uniformly in the anterior compartment of *ptc* mutants (15). Overexpression of *ptc* in wing discs results in the downregulation of full length CI along the A/P boundary without changing the *ci* messenger RNA levels, providing further evidence that *hh* regulates CI protein levels post-transcriptionally (23).

The mechanism involved in the *hh* regulation of full-length CI levels is currently under intense investigation. It has been shown that a domain C-terminal to the zinc-finger domain and N-terminal to the sites phosphorylated by protein kinase A is responsible for the retention of the 155-kDa full-length CI in the cytoplasm (22). Because the 75-kDa repressor form of CI does not contain this domain, its nuclear localization is not regulated. The 155-kDa CI protein forms a complex with *costal2* (*cos2*), a protein related to kinesin, *fused* (*fu*), a serine/threonine protein kinase, *suppressor of fused* (*su(fu)*), and other unidentified proteins. This complex is tethered to microtubules, thus anchoring CI in the cytoplasm. When cells receive an *hh* signal, the complex dissociates from microtubules, releasing the full-length CI from the cytoskeletal structure and presumably making it more accessible to nuclear translocation. *Hh* signaling also inhibits CI proteolysis, which accounts for the increased levels of full-length CI detected along the A/P boundary (22, 29, 30).

Protein kinase A (PKA) antagonizes *hh* and negatively regulates *ci* activity (see [Hedgehog Signaling](#)). Recent studies have helped to elucidate the mechanism involved in this regulation. Loss of PKA function or inhibition of PKA activity increases the levels of full-length CI in both embryos and discs (23). Cell culture experiments have provided details about the mechanism involved in the PKA regulation of CI (31). CI has consensus PKA phosphorylation sites at four serine residues in its C-terminus. Substitution of the serine residue with alanine in any of the first three PKA sites inhibits CI proteolysis, thus increasing the levels of full-length CI and stimulating CI-mediated transcription, suggesting that CI is the direct target of PKA regulation. Recently an F-box/WD40-repeat protein encoded by *slimb* has been identified in flies (32). Loss of *slimb* function causes the formation of supernumerary limbs and the accumulation of high levels of full-length CI. Because *slimb* is related to a yeast protein, *cdc4p*, that is involved in targeting cell-cycle regulators to the ubiquitin-mediated protein degradation pathway, *slimb* has been hypothesized to mediate CI proteolysis.

Evidence obtained both *in vitro* and *in vivo* suggests that *ci* directly mediates transcription of the *hh*

target genes. Consensus CI binding sites have been identified in the promoter regions of *ptc* and *wg*, and CI has been shown to bind to these sites both *in vitro* and *in vivo* (11, 12, 22). In cell culture, CI transactivates a reporter gene driven by the *wg* (*wingless*) promoter region that contains the CI binding sites (11). Both endogenous and ectopic *hh* signals transactivate a *lacZ* gene driven by the 758-bp minimal *ptc* promoter in a CI binding site-dependent manner in transgenic flies (12). These experiments support the idea that CI is a transcription factor that mediates the *hh* signal in the nucleus. *Ci* also functions as a suppressor of *hh* gene expression. *Ci* mutant clones generated in the anterior compartment of the imaginal discs ectopically express HH and induce the ectopic expression of *ci* and *hh* target genes in the surrounding wild type cells, presumably through the inductive function of the ectopic *hh* (33).

### 3. Gli Expression and Function in Vertebrates

A number of different *hh* homologues exist in vertebrates, and they also play important roles in the specification of cell fate. The transcription factors that respond to these vertebrate signaling cascades are considered homologues of *ci* because the zinc-finger motif that defines their DNA binding domains is more similar to that of CI than any other zinc-finger proteins. To date, three *ci* homologues, *Gli1*, *Gli2* and *Gli3*, have been studied in detail.

In mouse, all three *Gli* genes are detected initially during gastrulation at 7.5 days postcoitum (dpc) in broad domains in both ectoderm and mesoderm. As development progresses, the expression patterns of the three genes become more restricted and, by the completion of organogenesis, their expression is no longer detected. *Gli1* expression overlaps with *Sonic hedgehog* (*Shh*) expression initially in the ventral midline; however, by 9.5 dpc *Gli1* is excluded from *Shh*-expressing cells and is restricted to a narrow stripe of cells that lies along the caudal-rostral axis of the ventral neural tube and lateral to the *Shh*-expressing floor plate cells. The *Gli2* and *Gli3* expression patterns are widespread within the dorsal neural tube in domains that do not overlap with regions of *Gli1* expression. While *Gli2* expression is relatively homogenous and broad, *Gli3* is expressed as a gradient, with highest concentrations found in the dorso-lateral cells of the neural tube. This dorsal-ventral pattern of expression is maintained throughout the development of the brain and spinal cord (7, 34, 35). In the somites, all three *Gli* genes are expressed in the dorsal-lateral mesenchymal cells (dermomyotome), suggesting that the *Gli* genes are involved in myogenesis. The *Gli1* message is also detected in the ventral medial somite (sclerotome), whose development is patterned by *Shh*. *Gli3* is expressed throughout the newly formed somite, and its expression becomes restricted to the dorsal medial myotomal cells as development progresses (7, 36). As in the developing neural tube, the initial *Gli1* expression pattern in the limb bud overlaps that of *Shh* in the posterior mesenchyme and then becomes restricted to the posterior cells immediately adjacent to the ZPA (zone of polarization activity) that expresses *Shh*. Except in the posterior margin, *Gli2* and *Gli3* are expressed throughout the limb bud, with higher levels of *Gli3* expression detected in the autopod. During bone formation, *Gli1* is expressed strongly in the perichondrium surrounding the cartilage elements where *Indian hedgehog* (*Ihh*) is expressed, and *Gli2* and *Gli3* are expressed in the tissues surrounding the perichondrium (37). *Gli1* transcripts are also detected in developing lung, gut, gonad, and eye (7, 38).

That *Gli1* is expressed in cells adjacent to those that express *Shh* or *Ihh* is reminiscent of the relationship between *hh* and *ci* and suggests that *Gli1* is one of the *hh* target genes and may transduce the *hh* inductive pathways. Transgenic or Strong's luxoid mutant mice that misexpress *Shh* induce ectopic *Gli1* expression in the neighboring cells (37, 38). Removal of the notochord from wild-type quail embryos abolishes *Gli1* expression, and supplementing the same embryos with SHH restores *Gli1* expression (36), supporting the notion that *Gli1* is a target of *Shh*. Ectopic expression of *Gli1* induces ectopic expression of the ventral neural tube markers *HNF-3b* (hepatocyte nuclear factor-3b), *patched* (*ptch*), and *Shh*. In addition, ectopic *Gli1* expression in mouse and *Xenopus* suppresses the expression of the dorsal neural tube marker *Pax-3*. This phenotype is reminiscent of the ectopic floor plate differentiation obtained with a gain of *Shh* function. *Gli1* also induces ectopic ventral neuronal differentiation, as shown by its ability to induce the expression of the ventral neuron

markers, serotonin and dopamine. This inductive activity may be a secondary effect of *Shh*, because *Gli1* also induces *Shh* expression (34, 35, 39). The promoter region of *HNF-3b* contains Gli binding sites that can bind to *Gli1* *in vitro*. Tissue culture experiments demonstrate that *Shh* can transactivate a reporter gene that is controlled by the Gli binding sites found in the *HNF-3b* promoter. Mutating the Gli consensus binding site in the *HNF-3b* promoter abrogates the *Shh* response, showing that this site is required for *Shh*-dependent transactivation (40). The best evidence that *Gli1* is a mediator of *Shh* signaling comes from experiments in transgenic mice that express a *LacZ* reporter gene driven by the *HNF-3b* minimal promoter. In these animals, b-galactosidase expression is detected in the floor plate, and this expression depends on the presence of an intact Gli binding site (40). Taken together, these studies imply that the transcription of the *Gli1* gene is induced by *Shh*, whereupon the *Gli1* protein activates an array of genes needed to differentiate the floor plate and the ventral neurons. This pattern of regulation also exists in limb development. Marigo et al (8) have shown that misexpression of *Shh* in the anterior limb bud causes ectopic *Gli1* expression, and expression of a Gli-VP16 fusion protein activates the *Shh* target gene *ptch*. *Gli* is also a target for *Ihh* and mediates *Ihh* function in bone morphogenesis. In a developing bone, the prehypertrophic cells of cartilage express *Ihh* and the receptor for parathyroid/parathyroid related protein (PTH/PTHrP). On secretion, *Ihh* stimulates *Gli1* expression in the perichondrial cells, leading to the synthesis of PTHrP, which is then secreted and binds to its receptor on the *Ihh*-producing cells. Stimulation of the PTHrP signaling pathway ensures the proliferative state of the prehypertrophic cells and prevents their premature ossification (41, 42).

Our understanding of *Gli2* function comes mainly from experiments with *Gli2* knock-out mice. *Gli2* loss-of-function mutations result in the absence of floor plate differentiation, severe axial skeletal abnormalities, and downregulation of the *Shh* target genes, such as *ptch* and *Gli1*. These results suggest that *Gli2* functions as a transcriptional activator of *Shh* target gene expression in certain structures during development (43). Injecting frog embryos with *Shh* or *Gli1* causes ectopic expression of *Gli2*, implying that *Gli2* transcription is activated by *Gli1* in response to *Shh* signaling (35).

*Gli3* has been proposed for several reasons to function as a repressor of *Shh*. First, the expression patterns of *Gli3* and *Shh* are mutually exclusive (8, 34, 35, 44). Second, *Gli3* can suppress the activation of the *HNF-3b* promoter by *Gli1* in cotransfection assays (40). Third, *Gli3* knock-out mice (*Xt*) misexpress *Shh* in the anterior limb bud and dorsal neural tube, domains that normally express *Gli3* (35). Genes that are involved in limb outgrowth and patterning, such as *ptch* and *Hox*, and those for [fibroblast growth factors](#) (FGFs) and bone morphogenic proteins (*Bmp*) are also ectopically expressed in *Gli3* mutant mice due to the misexpression of *Shh* (45). This cascade of misexpression might explain the extra autopods associated with the heterozygous *Xt* mouse. Whereas *Gli3* is a repressor of *Shh*, *Shh* acts to repress *Gli3*. The misexpression of *Shh* by viral infection downregulates *Gli3* in the chick limb, and the ectopic expression of *Gli1* in the dorsal tube of transgenic mice suppresses the normal dorsal *Gli3* expression (8, 39). These results support the notion that CI function is accomplished by two Gli proteins in vertebrates, with *Gli1* acting as an activator and *Gli3* as a repressor, respectively, of *Shh* target gene expression.

Although Gli proteins have distinct expression patterns and cause different phenotypes when disrupted individually, they share some functional redundancies. Both *Gli1* and *Gli2* can induce motor neuron differentiation, and this redundancy may explain the fact that motor neuron induction is unaffected in homozygous *Gli2* mutant mice (43).

From an evolutionary perspective, many interesting parallels exist between the *Drosophila* and vertebrate *hh* signaling pathways. In *Drosophila*, *hh* suppresses the activity of the CI repressor by inhibiting the proteolysis of the 155-kDa CI protein. *Shh* signaling suppresses the activity of the *Gli3* repressor by inhibiting its transcription. CI suppresses *hh* expression, and *Gli3* suppresses the transcription of *Shh* in the cells where CI or Gli is active. In the presence of an *hh* signal, the 155-kDa CI protein is an activator of *hh* target genes, and this activity has evolved in vertebrates into the functions of *Gli1* and *Gli2*. Clearly, continued studies of CI and the Gli family members will

elucidate the evolution of the signaling systems that drive the developmental process.

## Bibliography

1. B. Hochman (1971) *Genetics* **67**, 235–252.
2. C. Nusslein-Volhard and E. Wieschaus (1980) *Nature* **287**, 795–801.
3. T. V. Orenic, D. C. Slusarski, K. L. Kroll, and R. A. Holmgren (1990) *Genes Dev.* **4**, 1053–1067.
4. K. W. Kinzler, J. M. Ruppert, S. H. Bigner, and B. Vogelstein (1988) *Nature* **332**, 371–374.
5. J. M. Ruppert, B. Vogelstein, K. Arheden, and K. W. Kinzler (1990) *Mol. Cell. Biol.* **10**, 5408–5415.
6. D. C. Hughes, J. Allen, G. Morley, K. Sutherland, W. Ahmed, J. Prosser, I. Lettice, G. Allan, M.-G. Mattei, M. Farrall, and R. E. Hill (1997) *Genomics* **39**, 205–215.
7. C.-C. Hui, D. Slusarski, K. A. Platt, R. Holmgren, and A. L. Joyner (1994) *Dev. Biol.* **162**, 402–413.
8. V. Marigo, R. L. Johnson, A. Vortkamp, and C. J. Tabin (1996) *Dev. Biol.* **180**, 273–283.
9. D. Zarkower and J. Hodgkin (1992) *Cell* **70**, 237–249.
10. J.-C. Marine, E. J. Bellefroid, H. Pendeville, J. A. Martial, and T. Pieler (1997) *Mech. Dev.* **63**, 211–225.
11. T. v. Ohlen, D. Lessing, R. Nusse, and J. E. Hooper (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2404–2409.
12. C. Alexandre, A. Jacinto, and P. W. Ingham (1996) *Genes Dev.* **10**, 2003–2013.
13. K. W. Kinzler and B. Vogelstein (1990) *Mol. Cell. Biol.* **10**(2), 634–642.
14. N. P. Pavletich and C. O. Pabo (1993) *Science* **261**, 1701–1707.
15. J. Hepker, Q.-T. Wang, C. K. Motzny, R. Holmgren, and T. v. Orenic (1997) *Development* **124**, 549–558.
16. H. Akimaru, Y. Chen, P. Dai, D.-X. Hou, M. Nonaka, S. M. Smolik, S. Armstrong, R. H. Goodman, and S. Ishii (1997) *Nature* **386**, 735–738.
17. K. Han and J. L. Manley (1993) *EMBO J.* **12**, 2723–2733.
18. P. J. Mitchell and R. Tjian (1989) *Science* **245**, 371–378.
19. C. K. Motzny and R. Holmgren (1995) *Mech. Dev.* **52**, 137–150.
20. C. Schwartz, J. Locke, C. Nishida, and T. B. Kornberg (1995) *Development* **121**, 1625–1635.
21. S. Eaton and T. B. Kornberg (1990) *Genes Dev.* **4**, 1068–1077.
22. P. Aza-Blanc, F.-A. Ramirez-Weber, and T. B. Kornberg (1997) *Cell* **89**, 1043–1053.
23. R. L. Johnson, J. K. Grenier, and M. P. Scott (1995) *Development* **121**, 4161–4170.
24. A. Martinez-Arias, N. E. Baker, and P. W. Ingham (1988) *Development* **103**, 157–170.
25. T. Tabata and T. B. Kornberg (1994) *Cell* **76**, 89–102.
26. J. Capdevila, M. P. Estrada, E. Sanchez-Herrero, and I. Guerrero (1994) *EMBO J.* **13**, 71–82.
27. J. Capdevila and I. Guerrero (1994) *EMBO J.* **13**, 4459–4468.
28. W. Li, J. T. Ohlmeyer, M. E. Lane, and D. Kalderon (1995) *Cell* **80**, 553–562.
29. D. J. Robbins, K. E. Nybakken, R. Kobayashi, J. C. Sisson, J. M. Bishop, and P. P. Therond (1997) *Cell* **90**, 225–234.
30. J. C. Sisson, K. S. Ho, K. Suyama, and M. P. Scott (1997) *Cell* **90**, 235–245.
31. Y. Chen, N. Gallaher, R. H. Goodman, and S. M. Smolik (1998) *Proc. Natl. Acad. Sci. USA* **95**, 2349–2354.
32. J. Jiang and G. Struhl (1998) *Nature* **391**, 493–496.
33. M. Dominguez, M. Brunner, E. Hafen, and K. Basler (1996) *Science* **272**, 1621–1625.
34. J. Lee, K. A. Platt, P. Censullo, and A. R. i. Altaba (1997) *Development* **124**, 2537–2552.

35. A. R. i. Altaba (1998) *Development* **125**, 2203–2212.
36. A.-G. Borycki, L. Mendham, and C. P. Emerson Jr. (1998) *Development* **125**, 777–790.
37. K. A. Platt, J. Michaud, and A. L. Joyner (1997) *Mech. Dev.* **62**, 121–135.
38. J. C. Grindley, S. Bellusci, D. Perkins, and B. L. M. Hogan (1997) *Dev. Biol.* **188**, 337–348.
39. M. Hynes, D. M. Stone, M. Dowd, S. Pitts-Meek, A. Goddard, A. Gurney, and A. Rosenthal (1997) *Neuron* **19**, 15–26.
40. H. Sasaki, C.-c. Hui, M. Nakafuku, and H. Kondoh (1997) *Development* **124**, 1313–1322.
41. A. Vortkamp, K. Lee, B. Lanske, G. V. Segre, H. M. Kronenberg, and C. J. Tabin (1996) *Science* **273**, 613–621.
42. B. Lanske, A. C. Karaplis, K. Lee, A. Luz, A. Vortkamp, A. Pirro, M. Karperien, L. H. K. Defize, C. Ho, R. C. Mulligan, A.-B. Abou-Samra, H. Juppner, G. V. Segre, and H. M. Kronenberg (1996) *Science* **273**, 663–666.
43. Q. Ding, J. Motoyama, S. Gasca, R. Mo, H. Sasaki, J. Rossant, and C.-c. Hui (1998) *Development* **125**, 2533–2543.
44. R. Burke and K. Basler (1997) *Curr. Opin. Neurobiol.* **7**, 55–61.
45. D. Buscher, B. Bosse, J. Heymer, and U. Ruther (1997) *Mech. Dev.* **62**, 175–182.

### **Suggestions for Further Reading**

46. A. R. Altaba (1997) Catching a Gli-mpse of hedgehog. *Cell* **90**, 193–196.
47. L. G. Biesecker (1997) Strike three for *GLI3*. *Nature Genet.* **17**, 259–260.
48. U. Radhakrishna, A. Wild, K.-H. Grzeschik, and S. E. Antonarakis (1997) Mutation in *Gli3* in postaxial polydactyly type A. *Nature Genet.* **17**, 269–271.

### **Cyclic AMP (3',5'-cyclic AMP, cAMP)**

Nucleotides play a universal role in life, as components of nucleic acids, as forms in which chemical free energy is stored, and as regulators of gene expression or enzyme activity. Cyclic adenosine 3',5'-monophosphate (cAMP) plays a universal role in the control of gene expression as well as in the integration of metabolic functions. It is present both in eucaryotes and in procaryotes (see [1](#) for early recognition of cAMP presence in bacterial genera). cAMP seems to be absent only from archaeobacteria (but see [Adenylate Cyclases](#)). Its presence was controversial in plants, but a work from the J. Schell laboratory investigating opine catabolism in plants resulted in the cloning of a gene that appeared to specify cAMP synthesis in plants. It was, however, later reported that this work was a fake, so that it is now admitted that plants do not make cAMP ([2](#)). Cyclic AMP has been reported to exist in cyanobacteria and in algae. This ubiquity explains the major interest displayed in its mode of synthesis, and the vast amount of literature devoted to the enzymes that produce cAMP from ATP, the adenylate cyclases (E. C. 4.6.1.1). Because cAMP is a regulatory molecule, it must be either excreted in the environment or inactivated in order not to accumulate. This is performed by 3',5'-cyclic-nucleotide phosphodiesterases (E. C. 3.1.4.17). These enzymes are generally specific for cyclic nucleotides (namely cAMP and cGMP) and sometimes specific for cAMP or cGMP alone. A variety of natural inhibitors modulate their activity (nucleoside triphosphates, pyrophosphate, and especially methylated xanthines, such as theophylline) by a variety of processes involving protein phosphorylation or calcium. Adenylate cyclases (see entry) form four independent classes of enzymes, and this raises the question of the origin of cyclic nucleotides as regulatory molecules, as well as their universal implication in regulatory networks. Because it is very polar and negatively



charged, cAMP does not permeate easily into cells (unless through specific transporters, generally unknown at present). More lipophilic analogs such as N<sup>6</sup>,O<sup>2'</sup>-dibutyryl adenosine-3',5' monophosphate are, therefore, used to modulate its concentration and mimic its effect in cell cultures *ex vivo*, but *in vivo* inhibitors of phosphodiesterase or specific mediators (neuromediators in particular) are used for therapeutic purposes where cAMP concentration must be altered.

Cyclic AMP was discovered in 1958 by E. Sutherland, who obtained a Nobel prize in 1971 for this and other discoveries on hormone action. As he has himself written, it is within the scope of molecular biology that cAMP was discovered: “*When I first entered the study of hormone action, some 25 years ago, there was a widespread feeling among biologists that hormone action could not be studied meaningfully in the absence of organized cell structure. However, as I reflected upon the history of biochemistry, it seemed to me there was a real possibility that hormones might act at the molecular level*”. Sutherland built up a cell-free system where well-known hormones could control glycolysis *in vitro*. Using this system he isolated a small thermostable molecule that was able to activate glycogen phosphorylase. Chemical analysis of the molecule permitted its identification as adenosine 3'-5' cyclic monophosphate. Synthesis of cAMP was shown to be the result of the action of an enzyme, adenylate cyclase, that generated cAMP and PP<sub>i</sub> from ATP, when activated by adrenaline (3). Since this pioneering work, the study of cAMP-mediated effects required the identification of the structure, function, and regulation of adenylate cyclases, the cAMP synthesizing enzymes. And, contrary to expectation, this did not yield a unifying picture of the role of cAMP, but, rather, demonstrated that this molecule has been used over and over again by living organisms for very different functions.

At the time of cAMP discovery, the aphorism of Jacques Monod, “*what is true for Escherichia coli is true for the elephant,*” induced biochemists to try bacterial systems to unravel cAMP function. After the discovery of cAMP by Sutherland in 1958, Mackman and Sutherland (4) demonstrated that glucose-starved *E. coli* cells accumulated cAMP. Ullmann and Monod (5) later established that part of the catabolite repression phenomenon (see entry) was controlled by cAMP. This discovery raised hopes that the study of this mediator in bacteria would help to understand what happens in eucaryotes (even perhaps in higher eucaryotes). However, it soon became clear that cAMP in eucaryotes was generally, as found by Sutherland, a “second messenger” that was used as an intracellular relay molecule to the action of extracellular hormones, while it acted directly on transcription via its receptor, the Catabolite Activator Protein (see entry) in *E. coli* (6). Study of the slime mold *Dictyostelium discoideum* revealed another function of cAMP, phylogenetically linked to is hormone-mediated action in higher eucaryotes, namely a pulsatile synthesis and degradation used by bacteria as a signal to control their aggregation properties as a differentiating multicellular organism (7).

The universal role of cAMP in controlling such diverse metabolic processes is puzzling because the enzymes needed for its synthesis, adenylate cyclases, are extremely varied and submitted to a wide variety of regulations (see entry). Why does this result in the synthesis of the same molecular species, cAMP? Is not all the regulatory process lost in this way? How can the cAMP signal generated by one enzyme type be distinguished from another? Compartmentalization is often invoked in this process, but, while this is relatively easily accounted for in the case of macromolecules, this is difficult to see in the case of small molecules such as cAMP. Another usual answer is to say that it is the combination of hormonal receptors of differing types and cAMP—and not cAMP alone—that is required for specificity. But would not cAMP synthesized from different sources also be recognized? Another answer is to remark that cAMP is known to be only one among many second messengers: cGMP has been added to the list as well as inositol phosphates, phosphatidyldiglycerides, calcium, etc (see corresponding entries). This certainly permits generation of a combinatorial control of activities, but would be very sensitive to accidental synthesis of cAMP. Cyclic AMP is not synthesized in a steady state way (even in bacteria). It is, therefore, important not to consider cAMP as such, with some average concentration, but to consider the shape of its time-dependent variation in concentration. In fact, observations are accumulating that strongly suggest

that cAMP does not have the same effect when it is delivered in a steady state fashion, rather than in a pulse (or a series of pulses) (7).

The motile and aggregating amoeba *D. discoideum* has been used as a paradigm for cell differentiation because undifferentiated cells start to differentiate into specific tissues some time after starvation. Secreted in the external medium, cAMP is necessary for aggregation. The genetics, biochemistry, cellular biology, and physiology of phenomena involving cAMP have been investigated in detail in this organism, where it controls, as in higher eucaryotes, a protein phosphorylation cascade, initiated after a regulatory cAMP-binding subunit of a protein kinase detaches from its target enzyme. This cascade is necessary not only for chemotaxis and aggregation but also for the triggering of genes involved in differentiation. The regulation of cAMP pulsatile concentration is mediated by two sets of enzymes adenylate cyclases and phosphodiesterases. In contrast with the situation with higher eucaryotes, however, cAMP and phosphodiesterase control operates not from the interior of the cell but from the external medium. This requires specific membrane receptors for cAMP and a process of signal transduction (see entry). In *D. discoideum* the pulses are generated by an appropriate coupling between adenylate cyclase activity, phosphodiesterase activity, and diffusion. The main observation is that variation in the cAMP pulse frequency changes the response of the cell. Many biochemical models can account for such cAMP pulses. These models require simple enzyme properties (in particular standard nonlinear features, such as self-activation and desensitization after saturating activity). They do not require the existence of many gene products, but only a specific behavior of enzymes (appropriate  $V_m$  and  $K_m$  of biosynthetic and degradative enzymes). Assuming that cAMP concentration modulation in time is the control event is, therefore, not a biochemical paradox.

In bacteria, Utsumi et al (8) have investigated cyclic AMP synthesis during the cell cycle of *E. coli* on synchronized cells, and they have given an unambiguous demonstration that there was a strong correlation of cAMP synthesis and replication or cell division, suggesting that the molecule may play some role in the cell cycle. This is also correlated with the position of the adenylate cyclase gene near the chromosome's origin of replication and with its very low level of expression, suggesting that expression is strongly coupled to DNA replication. This observation has long been overlooked because cells deficient for adenylate cyclase or CAP are viable, suggesting that cAMP is dispensible. Specific time-dependent variation of the concentration of cAMP for fine coordination of replication and division in *E. coli* is achieved by excretion of the nucleotide, rather than coupling to the activity of a phosphodiesterase (6). In this respect it is interesting that high concentrations of cAMP produced by foreign genes in *E. coli* are not toxic until they reach a very high level (at least 10-fold the normal concentration), whereas much lower concentration of cAMP produced by the endogeneous adenylate cyclase are toxic (9). Cyclic AMP has been formally linked to catabolite repression, but there are many catabolite sensitive operons that do not respond to cAMP. In addition cAMP synthesis is very strong in *E. coli* when cells enter the stationary phase of growth, suggesting that it could be a cell-to-cell signal as it is in *D. discoideum*. This may be one of its function in other bacteria (such as *Rhizobium* species), where it is clearly not linked to catabolite repression.

In the same way, intracellular and extracellular levels of cAMP vary during the cell cycle of *Saccharomyces cerevisiae*. Using centrifugal elutriation, Smith et al (10) showed that the intracellular cAMP concentration followed the stages of the cell cycle, being highest during the division cycle and lowest immediately before or just after cell separation; at the same time the external cAMP concentration did not vary. Therefore, in yeast, as in *E. coli*, it appears that the role of the external medium is to behave as a sink. These observations substantiate the demonstration that, under normal conditions, appropriate enzyme systems can generate a specific time-dependent pattern of cAMP concentration. As in the case of *E. coli*, it is known that in *S. cerevisiae* adenylate cyclase is dispensible in mutants of the cAMP receptor, and in *S. pombe* adenylate cyclase is dispensible during vegetative growth (11). But, as in this former case, cells that carry the mutation and are deficient in adenylate cyclase have several growth defects. In this respect, the function of the time-dependent cAMP pattern could be optimization of transient processes, in particular cell division and chromosome segregation.

In all these cases cAMP is recognized on the cell surface by a specific receptor. It is, therefore, interesting to identify cases where membrane targets of cAMP have been demonstrated (*D. discoideum* aside). Nerve cells typically generate and are sensitive to transient signals; they also have very involved patterns of adenylate and guanulate cyclase regulation (see entries). In this respect it is important to observe that cGMP (but also cAMP) has been shown to be involved as a central molecule in vision, taste, and olfaction. In particular, in addition to their role as second messengers in protein phosphorylation cascades, cyclic nucleotides are involved at the membrane surface, but intracellularly, in gating ion channels in olfactory and taste neurons. This certainly permits generation of a variety of time-dependent patterns for cAMP regulation, as a function of environmental inputs as well as of the fine molecular structure of the enzyme or its subunits (12, 13). Because ion channels are involved in the main functions of neurons (firing patterns), this makes cyclic nucleotides important in learning processes.

Indeed, many experiments have demonstrated that cAMP is involved as a mediator of learning and memory in invertebrates [*Aplysia* (14) and *Drosophila melanogaster* (15)], as well as vertebrates (16, 17). The study of mutants of *D. melanogaster* that are defective in learning or memory has been of major importance in our understanding of the physiology, biochemistry, and anatomy underlying conditioned behaviors. *D. melanogaster* learning mutants have been separated into two general classes: those with structural defects in the brain and those without obvious brain alterations. From studies of mutants affected in the brain structure, two areas have been found to be involved in conditioned behavior: the mushroom bodies and the central complex. Analysis of the mutants has shown that many types of molecules are involved in learning, but the cAMP-mediated phosphorylation cascade has emerged as especially important. During learning, time-dependent processes are involved in the stabilization of synapses, a general view being that they are created during growth as transient entities that can either regress or be stabilized. In this process, the evolution of the synaptic pattern is dependent on the pattern of neurotransmitter delivery. Analysis of the minimal requirements for synapses stabilization suggests that neurotransmitter release must be coupled to some other transient metabolic process in a retrograde manner in order to yield a stable geometry (18). In the cases where cAMP is involved, one can, therefore, speculate that the role of this mediator is to trigger an appropriate biochemical process when the proper time-dependent control of its synthesis is at work (19). Accordingly, once again, it is not the cAMP concentration that is important, but, rather, the time-variation of its concentration. In the process of learning, the regulation of adenylate cyclase activity would, therefore, be exquisitely tuned to permit delivery of the molecule in the proper time-dependent manner.

In *D. melanogaster*, five different genes have proven important for normal learning: *dunce* (a cAMP phosphodiesterase), *rutabaga* (an adenylyl cyclase), *amnesiac* (a product similar to adenylate cyclase activating peptides), *DCO* (protein kinase A), and *dCREB2* (a cAMP-response element binding protein). The products of many of these learning mutants are enriched in mushroom bodies. A process involving control of transcription by the cAMP response element binding protein (CREB)-responsive plays a central role in the formation of long-term memory in *D. melanogaster*, *Aplysia* and mammals. This is one of the examples where cAMP is involved in the control of transcription in eucaryotes, as it is in eubacteria, although through a different chain of events. Agents that prevent CREB activity interfere with the formation of long-term memory, whereas agents that increase the amount or activity of the transcription factor accelerate the process, thus indicating that CREB is essential for the switch from short-term memory to long-term memory (protein synthesis dependent) (20). Further work involving inbred mice strains as well as knock-out mutants affecting the hippocampal region demonstrated that both the genetic background and the temporal pattern of synaptic activity affects the cAMP-dependent synaptic plasticity (21).

## Bibliography

1. M. Ide (1971) Arch. Bioch. Biophys. **144**, 262–268.
2. T. Ichikawa, Y. Suzuki, I. Czaja, C. Schommer, A. Lessnick, J. Schell, and R. Walden (1998)

Nature 390–396.

3. E. W. Sutherland (1972) *Science* **177**, 401–408.
4. R. S. Mackman and E. W. Sutherland (1963) *Fed. Proc.* **22**, 470.
5. A. Ullmann and J. Monod (1968) *FEBS Lett.* **2**, 714–717.
6. A. Ullmann and A. Danchin (1983) *Adv. Cyclic Nucl. Res.* **15**, 1–52.
7. J. Dallon and H. Othmer (1997) *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **352**, 391–417.
8. R. Utsumi, M. Kawamukai, H. Aiba, M. Himeno, and T. Komano (1986) *J. Bacteriol.* **168**, 1408–1414.
9. A. Danchin (1993) **27**, 109–162.
10. M. E. Smith, J. R. Dickinson, and A. E. Wheals (1990) *Yeast* **6**, 53–60.
11. T. Maeda, N. Mochizuki, and M. Yamamoto (1990) *Proc. Natl. Acad. Sci. U.S.A.* **87**, 7814–7818.
12. T. Leinders-Zufall, M. Rand, G. Shepherd, C. Greer, and F. Zufall (1997) *J. Neurosci.* **17**, 4136–4148.
13. W. Zagotta and S. Siegelbaum (1996) *Annu. Rev. Neurosci.* **19**, 235–263.
14. C. Bailey, C. Alberini, M. Ghirardi, and E. Kandel (1994) *Adv. Second Messenger Phosphoprotein Res.* **29**, 529–544.
15. R. Davis (1996) *Physiol. Rev.* **76**, 299–317.
16. R. Bourtschuladze, B. Frenguelli, J. Blendy, D. Cioffi, G. Schutz, and A. Silva (1994) *Cell* **79**, 59–68.
17. Z. Xia and D. Storm (1997) *Curr. Opin. Neurobiol.* **7**, 391–396.
18. J. P. Changeux and A. Danchin (1976) *Nature* **264**, 705–712.
19. Y. Zhong and C. F. Wu (1991) *Science*, **251**, 198–201.
20. J. Yin and T. Tully (1996) *Curr. Opin. Neurobiol.* **6**, 264–268.
21. P. V. Nguyen, S. N. Duffy, and J. Z. Young (2000) *J. Neurophysiol.* **84**, 2484–2493.

## Cyclic Amp Receptor Protein (CRP)/Catabolite Gene Activator Protein (CAP)

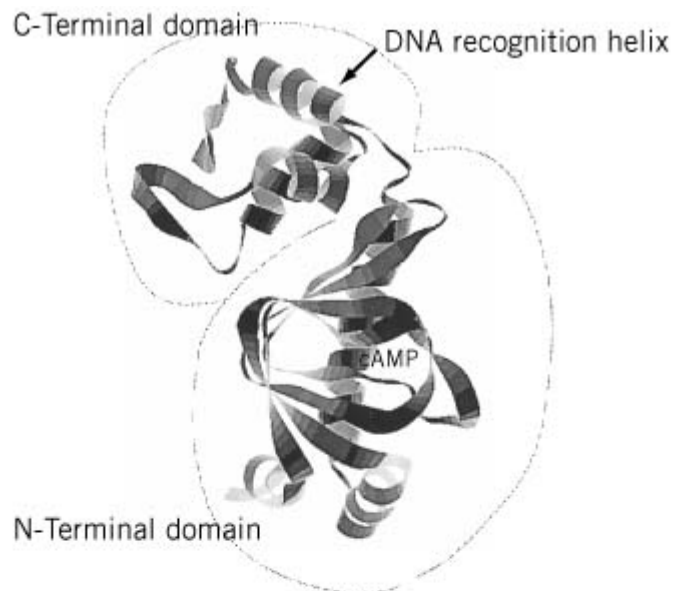
The cyclic AMP receptor protein (CRP), or catabolite gene activator protein (CAP), of *Escherichia coli* is one of the best characterized [transcription factors](#). It was discovered more than 25 years ago as a protein that binds **cyclic AMP (cAMP)** and stimulates **gene expression** of the [lac operon](#) (1, 2). Since then, CRP has attracted attention as the paradigm of gene activator proteins. It is now known that CRP, in conjunction with cAMP, participates as a global transcription factor in a wide regulatory network, both activating and repressing a large set of [operons](#). One mechanism of [catabolite repression](#) is due to the reduction in the intracellular concentrations of both cAMP and CRP.

The CRP protein is composed of two identical subunits of 209 amino acid residues each. It is active when complexed with cAMP, which behaves as an **allosteric** effector. It is a prototype of sequence-specific [DNA-binding proteins](#) containing a [helix–turn–helix motif](#). The cAMP–CRP complex binds to specific DNA sites in various operons. On binding DNA, it bends the DNA and interacts with **RNA polymerase** and/or other regulatory proteins to regulate [transcription](#) of the target operons.

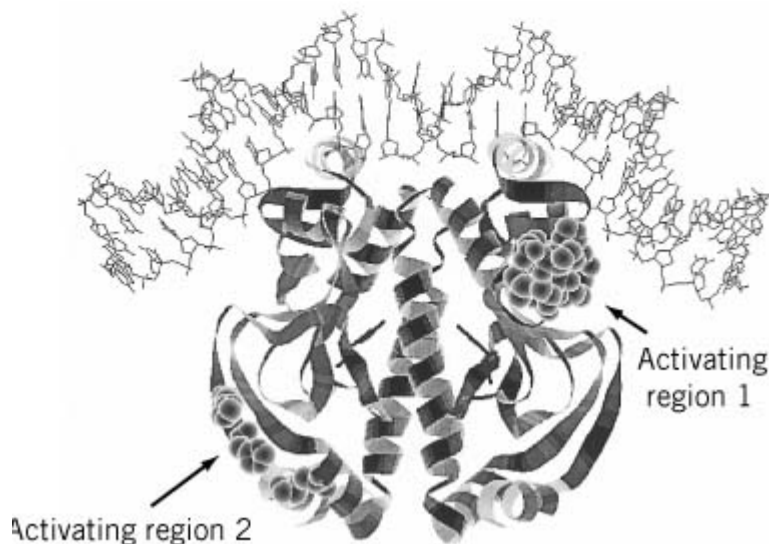
Following early studies on the properties of purified CRP, the *crp* gene was **cloned** and sequenced in

1982 (3, 4). The [X-ray crystallography](#) structure of CRP complexed with cAMP was solved in 1981 (5). This confirmed the dimeric nature and **domain** structure of CRP deduced from biochemical studies. The subunit has two domains (Fig. 1a). The larger *N*-terminal domain is responsible for cAMP binding and for dimerization. The smaller *C*-terminal domain contains a helix–turn–helix DNA-binding motif. CRP is one of the first regulatory proteins where the helix–turn–helix motif was identified. CRP alone is able to bind to DNA nonspecifically. The binding of cAMP induces a conformational state that binds to specific sequences with a dyad symmetry. The detailed nature of the structural changes in CRP caused by cAMP remains unknown.

**Figure 1.** Structures of the CRP monomer and the CRP-DNA complex. The crystallographic coordinates were obtained from the Brookhaven Protein databank (accession code 1CGP). The images were generated by Protein Adviser (FQS). **(a)** The two domains of a CRP monomer. The larger *N*-terminal domain is responsible for cAMP binding and for dimerization. The smaller *C*-terminal domain contains a helix–turn–helix motif that is involved in DNA binding. **(b)** The CRP dimer bound to a consensus DNA site. The positions that are responsible for transcription activation (activating regions) are highlighted by space-filling models. Activating region 1 contacts the *C*-terminal domain of a subunit of RNA polymerase (8), while activating region 2 contacts the *N*-terminal domain of a subunit of RNA polymerase (9).



(a)



(b)

The DNA sequences of many CRP binding sites, and a resulting **consensus sequence**, have been determined. They include variations of the 22-bp [palindromic](#) sequence, 5'-AAATGTGATCTAGATCACATTT-3', in which the two TGTGA motifs are relatively well conserved in different promoters. The binding of cAMP-CRP to target DNA induces bending of the DNA (6). The structure of cocrystals of cAMP-CRP bound to a 30-bp DNA target was determined in 1991 (7). It shows the helix-turn-helix motifs of the two subunits inserted into successive major groves of a DNA that is bent by 90°C (Fig. 1b).

In promoters where CRP alone is sufficient to activate transcription, the role of the cAMP-CRP complex is to enhance functional binding of RNA polymerase (8, 9). Binding of cAMP-CRP or RNA polymerase to these promoters stimulates the other's binding. An example is the *lac* promoter, where the CRP-binding site is centered on base pair -61, which is one of the preferable positions for CRP action. Another position where CRP activates transcription efficiently is -41, as in the *gal*

promoter (see [Gal Operon](#)). Changing these two standard distances reduces or eliminates CRP action. Only when this change is by an integral number of turns of DNA double helix is the ability of CRP to activate transcription retained to some extent, suggesting that both CRP and RNA polymerase must bind to the same face of the DNA ([10](#), [11](#)). Genetic and biochemical analysis revealed that cAMP–CRP activates transcription by directly contacting RNA polymerase at these simple promoters. A surface-exposed loop in the C-terminal domain of CRP and the C-terminal part of the  $\alpha$  subunit of RNA polymerase are primarily responsible for the contact between two proteins ([8](#)). An additional contact between the N-terminal domain of CRP and the N-terminal domain of a subunit is also involved in transcription activation at the promoters where the CRP site lies at –41 ([9](#)). The interaction between CRP and RNA polymerase is needed transiently to stimulate events leading to the formation of an open complex in these simple promoters ([12](#)). The contribution of CRP-induced DNA bending to transcription activation is still unclear. The CRP-binding site lies well upstream in several CRP-dependent promoters. In these cases, CRP acts as a coactivator of a second activator that is presumed to be responsible for the direct interaction with RNA polymerase. In the example of the *araBAD* promoter, CRP binds around –94 and the second activator AraC binds in the region between the CRP and RNA polymerase sites ([13](#)) (see [ara Operon](#)).

CRP also acts as a repressor or a corepressor in several operons. A simple example is the *cya* promoter, where cAMP–CRP inhibits the action of RNA polymerase by binding a target site located within the promoter ([14](#)). A more complex mechanism of repression is found in the operons that are coordinately regulated by CRP and CytR. Most of these promoters possess tandem CRP-binding sites that flank the CytR operator. An example is the *deo* promoter, where the CRP-binding sites are located at –41 and –94. The binding of cAMP–CRP to these sites dramatically enhances the binding of CytR to its operator, resulting in the formation of a repression complex containing both cAMP–CRP and CytR ([15](#)).

Although many studies have been identifying the regions of both transcription factors and RNA polymerase that are responsible for protein–protein and protein–DNA interactions, how these interactions lead to transcription activation is largely unknown. CRP, along with its target promoters, will continue to be a useful system for further understanding of molecular mechanisms of transcriptional regulation, including this fundamental question.

## Bibliography

1. G. Zubay, D. Schwartz, and J. R. Beckwith (1970) *Proc. Natl. Acad. Sci. USA* **66**, 104–110.
2. M. Emmer, B. deCrombrughe, I. Pastan, and R. L. Perlman (1970) *Proc. Natl. Acad. Sci. USA* **66**, 480–487.
3. H. Aiba, S. Fujimoto, and N. Ozaki (1982) *Nucleic Acids Res.* **10**, 1345–1361.
4. P. Cossart and B. Gicquel-Sanzey (1982) *Nucleic Acids Res.* **10**, 1363–1378.
5. D. B. McKay and T. A. Steitz (1981) *Nature* **290**, 744–749.
6. H.-M. Wu and D. M. Crothers (1984) *Nature* **308**, 509–513.
7. S. C. Schultz, G. C. Shields, and T. A. Steitz (1991) *Science* **253**, 1001–1007.
8. S. Busby and R. H. Ebright (1994) *Cell* **79**, 743–746.
9. S. Busby and R. H. Ebright (1997) *Mol. Microbiol.* **23**, 853–859.
10. K. Gaston, A. Bell, A. Kolb, H. Buc, and S. Busby (1990) *Cell* **62**, 733–743.
11. C. Ushida and H. Aiba (1990) *Nucleic Acids Res.* **18**, 6325–6330.
12. H. Tagami and H. Aiba (1998) *EMBO J* **17**, 1759–1767.
13. R. B. Lobell and R. F. Schleif (1991) *J. Mol. Biol.* **218**, 45–54.
14. H. Aiba (1985) *J. Biol. Chem.* **260**, 3063–3070.
15. L. Sogaard-Anderson and P. Valentin-Hansen (1993) *Cell* **75**, 557–566.

## Suggestions for Further Reading

16. G. Zubay (1980) The isolation and properties of CAP, the catabolite gene activator. *Methods Enzymol.* **65**, 856–877.
17. J. L. Botsford and J. G. Harman (1992) Cyclic AMP in prokaryotes, *Microbiol. Rev.* **56**, 100–122.
18. D. M. Crothers and T. A. Steitz (1992) "Transcriptional activation by *Escherichia coli* CAP protein", in *Transcriptional Regulation*, S. L. McKnight and K. R. Yamamoto, eds., Cold Spring Harbor Press, Cold Spring Harbor, New York, Vol. **1**, pp. 501–534.
19. A. Kolb, S. Busby, H. Buc, S. Garges, and S. Adhya (1993) Transcriptional regulation by cAMP and its receptor protein. *Annu. Rev. Biochem.* **62**, 749–795.

### Cyclic GMP (Cyclic Guanosine 3',5'-Monophosphate, cGMP)

Cyclic GMP, a structural analog of **cyclic AMP**, occurs at similarly low concentrations (*i.e.*, in the micromolar range) in many animal tissues. It is synthesized by [guanylate cyclases](#), which correspond to the [adenylate cyclases](#) that synthesize cyclic AMP, and is destroyed by specific **phosphodiesterases**.

A functional role for cyclic GMP in **bacteria** is controversial. It is present in *Escherichia coli* but at such a low intracellular concentration (nanomolar, corresponding to about one molecule per cell) that hardly makes it a significant molecule. In addition, the sequence of the *E. coli* [genome](#) does not reveal any gene product that could code for a guanylate cyclase. There are many reports suggesting the presence of cGMP in bacteria ([1](#)), but its presence is perhaps most likely in myxobacteria, where it could be involved in cell aggregation and [differentiation](#).

In contrast, cGMP is universally in **eukaryotes** (except in **plants**). It is generally involved in processes leading to activation of specific regulation cascades, which are differentially controlled by appropriate mediators, or leading to control specific processes for neuronal activation of sense organs, such as the sensitivity to light of retina receptor cells and the triggering of olfaction and taste ([2-4](#)). The organization of guanylate cyclase and the control elements is sometimes similar to but distinct from the organization of hormonally regulated adenylate cyclase. In particular, **G-protein-mediated** regulation of vision operates on the phosphodiesterase rather than on the cyclase.

Phototransduction systems in vertebrates and invertebrates share a great deal of overall strategic similarity but differ significantly in the underlying molecular machinery. In the dark, vertebrate retinal rod cells synthesize a high level of cGMP that keeps gated **sodium channels** open in the plasma membrane of the outer segment ([5](#)). Light closes these channels by activating an enzymatic cascade that leads to rapid hydrolysis of cGMP by cGMP-specific phosphodiesterase ([6](#)). This hyperpolarizes the cell and modulates transmitter release at the synaptic buttons. Photoexcited [rhodopsin](#) triggers a transducer protein ([transducin](#)), which is related to G-proteins by catalyzing the exchange of GTP for bound GDP. Subsequently, the activated GTP-form of transducin switches on the phosphodiesterase. The cascade that has an overall gain of  $10^5$ , is turned off by the [GTPase](#) activity of transducin and by the action of two proteins, rhodopsin kinase and arrestin. The [kinetics](#) of the reactions in the cGMP cascade limit the temporal resolution of the visual system, and statistical fluctuations in the reactions limit the reliability of detecting dim light. Together with **calcium** ions and **inositol phosphates**, cGMP controls visual excitation and adaptation. A light-induced fall in the internal free  $\text{Ca}^{2+}$  concentration subsequently stimulates resynthesis of cGMP, antagonizes the catalytic activity of rhodopsin, and restores the high affinity of the light-regulated



sodium channel for cGMP, allowing the cell to adapt to background light (7).

The initial events in mollusks and arthropods are probably similar to those of vertebrates. However, whereas light activation of vertebrate photoreceptors leads to activation of cGMP-phosphodiesterase and generation of a hyperpolarizing response, activation of photoreceptors of invertebrates like *Drosophila* leads to stimulation of [phospholipase C](#) and generation of a depolarizing response (8). Cyclic GMP has also been implicated in modulating behavior in insects (9).

Cyclic GMP is also a [second messenger](#) in regulation mediated by natriuretic peptide hormones, but it is probably in the recently discovered [nitric oxide](#) (NO) regulation cascade that cGMP has the most unexpected role. Nitric oxide and atrial natriuretic peptide hormones play key roles in a number of neuronal functions, including learning, memory, and in blood circulation. Most experiments suggest that they exert converging actions by elevating of intracellular cGMP levels by activating soluble and membrane-bound guanylate cyclases. Cyclic GMP is the starting point for multiple [signal transduction](#) cascades, which are now beginning to be unraveled (10). In more than a quarter of a century since the discoveries of atrial granules and volume receptors in the heart atria, the search for natriuretic hormones has led to the isolation and identification of many atrial natriuretic factors (ANF) [(11), as a specific example see (12)]. In the heart, for example, ANF peptides are synthesized and stored in the [Golgi apparatus](#) of cardiac myocytes and are released in response to atrial wall stretch following acute plasma volume expansion and increased central blood volume. The mechanisms of the renal action of these potent natriuretic hormones are not yet completely unraveled. The renal hemodynamic, tubular, and adrenal, and systemic vascular effects are related to enhanced cGMP synthesis in specific medium-sized arteria, in the glomeruli and specific tubular segments, and in adrenal tissue. Specific ANF-binding sites have been detected in these target organs. A primary action of elevated cGMP levels is stimulating cGMP-dependent protein kinase (PKG), the major intracellular receptor protein for cGMP, which phosphorylates substrate proteins to trigger a regulation cascade (10).

Cyclic GMP-dependent protein kinases (PKG) also mediate some of the neuronal effects of cGMP (13), but unfortunately few PKG substrates are known in the brain. In striatonigral nerve terminals, for example, NO mediates phosphorylation of the protein phosphatase regulator, dopamine- and cyclic AMP-regulated phosphoprotein by PKG. PKG substrates are critically placed in the protein phosphorylation network and regulate protein phosphatases, intracellular calcium levels, and the function of many ion channels and neurotransmitter receptors. Nitric oxide is a signaling molecule in the nervous system of both mammals and insects. In contrast to classical transmitters, NO permeates membranes and acts on neighboring targets normally limited by diffusion barriers. This diffuse signaling is evolutionarily highly conserved. The NO forming enzyme, NO synthase, is present mostly in the nervous system, especially the brain. A soluble form of guanylate cyclase is the major target of NO action. Usually there is cellular separation of the release site and target site of NO, although exceptions to this rule exist.

As with cAMP, cGMP is important for memory in insects. In the honeybee, for example, the NO/cGMP system in the antennal lobes is implicated in the processing of adaptive mechanisms during chemosensory processing, and experimental data support a specific role of the NO system in memory formation (14).

Signal transduction in gastric and intestinal smooth muscle is mediated by receptors coupled via distinct G-proteins to various effector enzymes. Calcium is implicated in signal transduction in different ways according to the cell type (e.g., circular and longitudinal muscle cells). The initial steps involve  $\text{Ca}^{2+}$ /calmodulin-dependent activation of myosin light-chain kinase and the interaction of [actin](#) and [myosin](#). Relaxation is mediated by cAMP-and/or cGMP-dependent protein kinases. A specific cascade involves G-protein-dependent stimulation of  $\text{Ca}^{2+}$  influx, leading to  $\text{Ca}^{2+}$ /calmodulin-dependent activation of a constitutive NO synthase in muscle cells that activates soluble guanylyl cyclase. The resulting activation of protein kinase A and PKG is jointly responsible

for muscle relaxation (15).

Therefore, cyclic GMP is a secondary messenger that acts on targets that are sometimes similar to those of cAMP but proceed through a completely different cascade, in which the diffusible NO (and sometimes carbon monoxide) play a major role.

## Bibliography

1. N. B. Bhatnagar, R. Bhatnagar, and T. A. Venkitasubramanian (1984) *Biochem. Biophys. Res. Commun.*, **121**, 634–640.
2. L. Stryer (1986) *Ann. Rev. Neurosci.* **9**, 87–119.
3. T. Misaka, Y. Kusakabe, Y. Emori, T. Gono, S. Arai, and K. Abe (1997) *J. Biol. Chem.* **272**, 22623–22629.
4. T. Nakamura and G. H. Gold (1987) *Nature* **325**, 442–444.
5. J. T. Finn, M. E. Grunwald, and K. W. Yau (1996) *Ann. Rev. Physiol.* **58**, 395–426.
6. M. Chabre, J. Bigay, F. Bruckert, F. Bornancin, P. Deterre, C. Pfister, T. Vuong, and D. Baylor (1988) *Cold Spring Harbor Symp. Quant. Biol.* **53** (*Pt 1*), 313–324.
7. D. Baylor (1996) *Proc. Natl. Acad. Sci. USA* **93**, 560–565.
8. C. Zuker (1996) *Proc. Natl. Acad. Sci. USA* **93**, 571–576.
9. K. Osborne, A. Robichon, E. Burgess, S. Butland, R. Shaw, A. Coulthard, H. Pereira, G. RJ, and M. Sokolowski (1997) *Science* **277**, 834–836.
10. S. M. Lohmann, A. B. Vaandrager, A. Smolenski, U. Walter, and H. R. DeJonge (1997) *Trends Biochem. Sci.*, **22**, 307–312.
11. H. J. Kramer and B. Lichardus (1986) *Klinische Wochenschrift* **64**, 719–731.
12. L. R. Forte, X. Fan, and F. K. Hamra (1996) *Am. J. Kidney Dis.* **28**, 296–304.
13. X. Wang and P. J. Robinson (1997) *J. Neurochem.* **68**, 443–456.
14. U. Muller (1997) *Prog. Neurobiol.* **51**, 363–381.
15. G. Makhoulouf and K. Murthy (1997) *Cell Signal* **9**, 269–276.

## Cyclins

From a simple start as a family of proteins with interesting patterns of accumulation during the cell cycle, the cyclins have grown to become key regulators of diverse cellular processes, in particular the cell cycle. Most cyclins, whether they are present only at specific times during the cell cycle or constitutively, exert their functions through their associated cyclin-dependent kinase (Cdk) binding partners. Binding to cyclins is one of the required steps in the activation of Cdks. The degradation of many cyclins by the ubiquitin system provides a means of inactivating the associated Cdk following completion of its function. A large number of cyclins have been identified including some, like the original cyclins, that have roles in cell cycle progression, and others that don't cycle and that activate Cdks involved in very different activities, such as transcription.

The first cyclins were found during studies of translational control before and after fertilization of sea urchin eggs conducted as part of the Physiology course at the Marine Biological Laboratory in Woods Hole (1). These proteins were synthesized continuously, and accumulated until their abrupt degradation during mitosis. This sawtooth pattern of accumulation hinted that cyclins might play an important role during the cell cycle, either as inducers of cell cycle transitions or, perhaps less

interestingly, as proteins that responded to cell cycle states to perform functions important for that stage. Later work showed that the injection of cyclin mRNA caused frog oocytes to mature into eggs (2), that is, to progress through meiosis and to arrest in second meiotic metaphase, ready for fertilization. This result suggested that cyclins were actually inducers of the transitions into meiosis and mitosis. Subsequent work revealed that the mitotic cyclins are the regulatory subunits of maturation promoting factor (MPF) (3, 4). MPF had been characterized as a proteinaceous activity that, when withdrawn from the cytoplasm of an egg and injected into an oocyte, caused that oocyte to mature into an egg. The catalytic subunit of *Xenopus* MPF had recently been determined to be Cdc2, a protein kinase first identified by genetic studies in the fission yeast *Schizosaccharomyces pombe* as the key regulator of the G2-M phase transition. The following years rapidly revealed large families of proteins showing sequence similarity to the original, mitotic cyclins, and to Cdc2. The cyclins are referred to by letter (cyclin A, cyclin B, 1/4) and the kinase partners have been called cyclin-dependent kinases (Cdk2, Cdk3 1/4). (For further discussion of these and other aspects of cyclin function, see [Cell Cycle](#).)

One of the irreversible ratchet steps in the cell cycle is the degradation of cyclins by the ubiquitin system (5). Ubiquitin is a 76-aa protein whose covalent attachment to proteins can target them for proteolysis by the proteasome, a huge multiprotease unwinding and degrading machine. Ubiquitin is activated by its ATP-dependent covalent attachment to a cysteine side chain on an enzyme called E1. E1 then transfers ubiquitin to a cysteine of one of many E2 enzymes. The E2s can ubiquitinate substrates, often with the help of an E3. The E3s may be the most diverse and interesting components of this system. Some E3s receive the covalently bound ubiquitin; others serve as matchmakers that bring together a substrate and the appropriate E2. For the mitotic cyclins, the E3 was first termed the cyclosome, though it is now generally termed the anaphase promoting complex (APC), which consists of 8 to 12 subunits. The APC is the regulated component of the ubiquitin system for the degradation of the mitotic cyclins and serves as the target for checkpoint signals that can block cyclin degradation. Proteolysis of cyclins that act earlier in the cell cycle has been best studied in the budding yeast *S. cerevisiae* and, though mediated by the ubiquitin system, uses an "SCF complex" as the E3 component, rather than the APC.

The cyclins now comprise a large family of proteins with diverse functions, each bound to a cyclin-dependent kinase (Cdt) catalytic partner. All cyclins resemble the first mitotic cyclins (cyclin A and B) in sequence, but not all cycle during the cell cycle. Cyclins A, B, D, and E play major roles in regulating cell cycle transitions. The cell cycle stage at which each kinase functions is largely determined by when each cyclin partner accumulates during the cell cycle. Quite a few cyclins (and their associated Cdks) function in transcription. For instance, cyclin H is one of the subunits of TFIIH, a general transcription factor for RNA polymerase II involved in the phosphorylation of the C-terminal domain (CTD) of the large subunit of the polymerase. Cyclin C is a subunit of the RNA polymerase II holoenzyme, which can also phosphorylate the CTD. As more cyclins are discovered (and the alphabet is exhausted!) the majority of all cyclins will probably be noncycling regulatory subunits of protein kinases functioning in processes other than the cell cycle. The observation that the first cyclins had cyclic patterns of accumulation and were involved in cell cycle control may represent more a historical footnote owing to their relative ease of discovery than a reflection of fundamental properties of this family of proteins.

Direct roles of cellular cyclins in diseases are extremely rare. For instance, some tumor cells contain cyclin gene amplifications. Many tumors overexpress cyclin messenger RNAs, presumably a secondary consequence of increased rates of cell proliferation. Much more important, however, is the general circumvention of normal cell cycle controls that is a hallmark of cancers. This topic is discussed in [Cell Cycle](#). Interestingly, some viruses have co-opted cyclins to subvert normal cell cycle controls (6). For instance, the Kaposi Sarcoma Associated Herpesvirus encodes a D-type cyclin. Following infection, this cyclin associates with Cdk6, which is involved in progression through the G1 phase of the cell cycle, and activates it in a manner that makes it resistant to Cdk inhibitor proteins that normally restrain G1 progression. The infected cell is thereby pushed into S phase, allowing the virus to replicate and to produce progeny virus. This situation provides yet

another example of how viruses have adapted normal cellular proteins for their own ends.

### Bibliography

1. T. Evans, E. T. Rosenthal, J. Youngblom, D. Distel, and T. Hunt (1983) *Cell* **33**, 389–396.
2. K. I. Swenson, K. M. Farrell, and J. V. Ruderman (1986) *Cell* **47**, 861–870.
3. J. C. Labbé, J. P. Capony, D. Caput, J. C. Cavadore, J. Derancourt, M. Kaghad, J. M. Lelias, A. Picard, and M. Dorée (1989) *EMBO J.* **8**, 3053–3058.
4. J. Gautier, J. Minshull, M. Lohka, M. Glotzer, T. Hunt, and J. L. Maller (1990) *Cell* **60**, 487–494.
5. A. M. Page and P. Hieter (1999) *Annu. Rev. Biochem.* **68** 583–609.
6. H. Laman, D. J. Mann, and N. C. Jones (2000) *Curr. Opin. Genet. Dev.* **10**, 70–74.

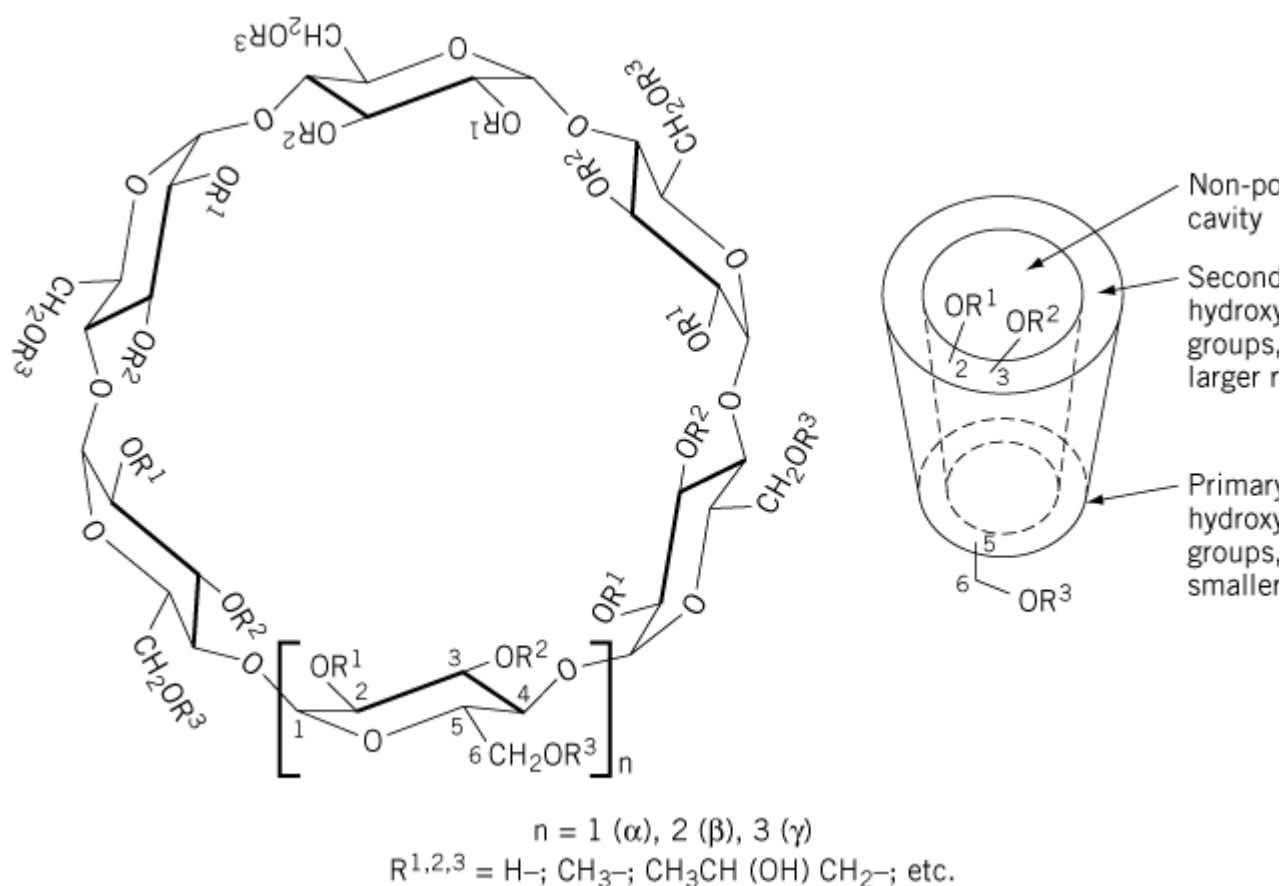
### Suggestions for Further Reading

7. D. O. Morgan (1997) Cyclin-dependent kinases: engines, clocks, and microprocessors. *Annu. Rev. Cell Dev. Biol.* **13**, 261–291.
8. A. Murray and T. Hunt (1993) *The Cell Cycle*, W.H. Freeman and Company, New York.

## Cyclodextrins

Cyclodextrins are cyclic oligosaccharides with a truncated cone shape and an axial void cavity (Fig. 1). The diameter and the volume of the cavity vary with the number of glucose units in the cyclodextrin ring (1). The most commonly used cyclodextrin is b-cyclodextrin, which has seven glucose units and a cavity with a diameter of 0.78 nm and a volume of approximately 35 nm<sup>3</sup>. Other natural cyclodextrins, such as a- and g-cyclodextrin, with six and eight glucose units and 0.57- and 0.95-nm cavity diameters, respectively, are also frequently used. The outer surface of the cyclodextrin molecule is [hydrophilic](#), because the majority of the hydroxyl groups project outward, resulting in good water solubility. The internal cavity is relatively [nonpolar](#), and it can encapsulate nonpolar solutes of appropriate dimensions, with binding occurring through various nonpolar interactions. Such binding is known as *inclusion complexation*.

**Figure 1.** Cyclodextrin structure.



The conformation of cyclodextrins in aqueous solution is believed to be that of the truncated cone of Figure 1. Molecules of **hydrophobic** compounds of appropriate size and shape penetrate into the cavity and are bound mainly through hydrophobic interactions, whose strengths depend on the efficiency of the contact. The edge of the torus of the larger circumference consists of secondary hydroxyl groups that are attached to **chiral** carbons (C2 and C3 of the glucose units). This structure results in variable binding affinities for different **enantiomers**, probably due to interactions of the chiral solute with the chiral entrance to the cavity (2). The primary hydroxyl groups of the glucose monomers make up the smaller edge of the cone. Chemical derivatization of natural cyclodextrins via modification of their hydroxyl groups is currently an area of very active research (3-6). Such modifications are yielding materials of varying complexation selectivities and physicochemical properties, such as improved solubility (7).

Hoffmann and Bock (8) have examined complex formation between different cyclodextrins and **nucleotides**. They found that adenosine, cytidine, guanosine, uridine, inositol, and deoxythymidine monophosphates (AMP, CMP, GMP, UMP, IMP, and dTMP) did not form a complex with  $\alpha$ -cyclodextrin, but  $\beta$ -cyclodextrin readily bound AMP and IMP. It was concluded that these six nucleotides are too bulky to fit into the cavity of  $\alpha$ -cyclodextrin. When a complex is formed with  $\beta$ -cyclodextrin, the ribose and phosphate groups of the nucleotides exert a stabilizing effect by establishing hydrogen bonds with the outer rim of the cyclodextrin molecules. The position of the phosphate group is not important; increasing distance of the phosphate group from the base increased the stability of the complex. Larger oligonucleotides exhibited decreasing tendencies for complex formation, but the extent of complexation depended significantly on their base composition. Interestingly, polynucleotides with double- or triple-helical structures did not show any complexation with  $\beta$ -cyclodextrin (9). **Transfer RNA**, which comprises both helical and nonhelical structures, can interact with cyclodextrins, and these complexes have been used extensively for tRNA studies.

Cyclodextrins have also been used in labeling nucleic acid molecules in various biochemical analysis applications, especially in [DNA sequencing](#). Cyclodextrin labels provide potentially high signal efficiency and versatility in label colors, while maintaining uniform chemical and physical properties. Cyclodextrin tracers are prepared by coupling the cyclodextrin to specific binding substances, such as nucleic acids, and forming inclusion complexes with the fluorophores. Then the DNA molecules are sequenced using cyclodextrin-labeled chain terminators. This method allows DNA sequencing with high sensitivity and high throughput ([10](#)).

Ikeda et al. ([11](#)) reported the use of anthryl(alkylamino)–cyclodextrin complexes as chemically switched DNA intercalators that were [allosteric](#). On adding a ligand that is tightly bound in the cyclodextrin cavity, such as 1-adamantol, the host molecule releases the anthryl unit, which then leads to strong intercalation with the double-stranded DNA molecule. This principle could be extremely useful in nucleic acid reactions of medicinal and biotechnological importance, particularly in view of the established methods to modify cyclodextrins for specific interactions with a wide range of substances, for new drug delivery systems.

Cyclodextrins and their analogs can also be used as carriers to increase cellular uptake of phosphorothioate [antisense oligonucleotides](#). Cellular uptake of phosphorothioate oligodeoxynucleotides in the presence of various cyclodextrin analogs was found to depend on the concentration and the time. In particular, 2-hydroxypropyl- $\beta$ -cyclodextrin (HP $\beta$ CD), 2-hydroxyethyl- $\beta$ -cyclodextrin (HE $\beta$ CD), and a mixture of various HP $\beta$ CDs having different degree of substitution were observed to increase the uptake of phosphorothioate oligodeoxynucleotides two- to threefold in 48 h ([12](#)).

Cyclodextrins and derivatives are used to enhance the bioavailability of many water-insoluble pharmaceuticals ([13](#)). The solubilization results in fast and quantitative *in vivo* delivery for intravenous and intramuscular dosing. A decrease in irritation at the administered site can be observed ([14](#)). The solubility of [hormones](#) such as hydrocortisone was enhanced 72-fold using randomly methylated  $\beta$ -cyclodextrin ([15](#)). Sublingual administration of  $\gamma$ -cyclodextrin complex of testosterone avoided rapid first-pass loss of the hormone and directed it effectively into the circulation; administration of the complex into the stomach resulted in much lower circulatory hormone levels ([16](#)).

The relatively good water solubility of the cyclodextrins makes these materials useful in [chromatography](#) and [electrophoresis](#) separation methods, such as in [HPLC](#) (high-performance liquid chromatography) and [capillary zone electrophoresis](#) (CZE) ([17](#), [18](#)). Buffer systems, even with organic modifiers, can be used to control the pH or modify other secondary equilibrium. Temperature also has a significant effect on selectivity in cyclodextrin-mediated separation systems.

## 1. Acknowledgment

The authors gratefully acknowledge Professor József Szejtli for his stimulating discussions.

## Bibliography

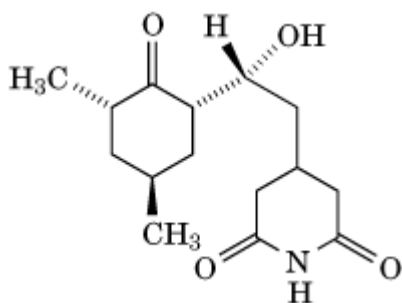
1. J. Szejtli (1988) In *Cyclodextrins and Their Inclusion Complexes*, Cyclodextrin Technology, Kluywer, Dordrecht, The Netherlands.
2. A. Guttman et al. (1988) *J. Chromatogr.* **488**, 41–53.
3. B. Sébille (1987) In *Cyclodextrin Derivatives, Cyclodextrins and their Industrial Uses*, D. Duchêne, ed., Editions de Santé, Paris, pp. 353–385.
4. K. Uekama (1985) *Pharm. Intl.* **6**, 61–65.
5. J. Pitha et al. (1986) *Intl. J. Pharm.* **29**, 73–82.
6. B. W. Müller and U. Brauns (1985) *Intl. J. Pharm.* **26**, 77–88.

7. Y. Y. Rawjee, D. E. Stark, and Gy. Vigh (1993) *J. Chromatogr.* **635**, 291–306.
8. L. Hoffmann and R. M. Bock (1970) *Biochemistry* **9**, 3542–3550.
9. J. Szejtli, ed. (1982) *Proceedings of the 1st International Symposium on Cyclodextrins*, Reidel, Dordrecht, The Netherlands.
10. K. M. Kosak, PCT Intl. Appl. (WO 9102040, 2/21/1991), 69pp.
11. T. Ikeda, K. Yoshida, and H. J. Schneider (1995) *J. Am. Chem. Soc.* **117**, 1453–1454.
12. Q. Zhao, J. Temsamani, and S. Agrawal (1995) *Antisense Res. Devel.* **5**, 185–192.
13. R. A. Rajewski and V. Stella (1996) *J. Pharm. Sci.* **85**, 1142–1169.
14. T. Loftsson and M. E. Brewster (1996) *J. Pharm. Sci.* **85**, 1017–1025.
15. J. Szejtli (1994) *Med. Res. Rev.* **14**, 353–386.
16. J. Pitha, E. J. Anaissie, and K. Uekama (1987) *J. Pharm. Sci.* **76**, 788–790.
17. S. Ahuja (1991) In *Chiral Separations by Liquid Chromatography*, American Chemical Society, Washington, DC.
18. A. Guttman (1996) In *Handbook of Capillary Electrophoresis*, 2nd ed., J. P. Landers, ed., CRC Press, Boca Raton, pp. 75–100.

## Cycloheximide

$C_{15}H_{23}NO_4$  (Fig. 1; molecular mass, 281.35). Cycloheximide is an antibiotic produced by *Streptomyces griseus*, and it inhibits protein biosynthesis in eukaryotes. Cycloheximide does not affect protein synthesis in prokaryotes or in [mitochondria](#). It binds to the 60S subunit of eukaryotic [ribosomes](#) and inhibits the [peptidyl transferase](#) activity. The inhibition of peptidyl transfer causes a stack of ribosomes on the [messenger RNA](#), generating polysome accumulation in cells. Cycloheximide is frequently used for the study of gene expression in eukaryotes as an inhibitor of protein synthesis *in vivo*, as well as *in vitro*.

**Figure 1.** Cycloheximide.



## Cyclophilin

Historically, an 18-kDa cytosolic protein from pig kidney was the first [peptidyl prolyl cis/trans isomerase](#) (abbreviated PPIase) discovered and characterized enzymatically (1). These enzymes catalyze the rotation of a peptidyl prolyl moiety in [peptides](#) and [proteins](#). In 1989 it was found that this protein is nearly identical in its [primary structure](#) to the already known cytosolic receptor of the immunosuppressive undecapeptide [cyclosporin A](#) in human T lymphocytes, previously named cyclophilin [Cyp18; for nomenclature see (2)] (3-5). This enzyme represents the archetype of cyclophilins in that it constitutes the catalytic **domain** of larger cyclophilin-like proteins. Now two additional families of such isomerases are known, those that bind [FK506](#) (FKBP) and the parvulin family.

Alignment of amino acid sequences of the cyclophilin family defines residues that are highly conserved among vertebrates, plants, fungi, and bacteria. Nearly perfect conservation is located primarily in the central segment of the protein, including the sequences -Phe-His-Arg-Ile/Val-Ile-(Xaa)<sub>5</sub>-Gln-Gly-Gly- at positions 53 to 65 and -Met-Ala-(Xaa)<sub>9-10</sub>-Gln-Phe-Phe/Tyr-Ile/Val- at positions 100 to 114, where Xaa is any amino acid and / separates alternatives. These sequences are used as typical **motifs** to search for other cyclophilins. The three-dimensional structure of the archetypal cyclophilin consists of an eight-stranded antiparallel [b-sheet](#) barrel capped by **a-helices** (6, 7).

The characterization of 11 distinct cyclophilin genes in the genome of the [nematode](#) *Caenorhabditis elegans* demonstrates that organisms use numerous cyclophilins (8). As many as 12 human homologues have been identified thus far from the initial assessment of gene diversity by [expressed sequence tag](#) analysis. The C-terminal and/or N-terminal amino acid extensions of larger cyclophilins consist of additional **domains**, directing the proteins into specific cellular compartments, mediating protein/protein and DNA/protein interactions, or generating other biochemical functions.

The cytosolic Cyp18 often represents the major cyclophilin of cells. Kidney tubules and endothelial cells, for example, contain 10 µg of Cyp18 per mg total protein (9). Despite the numerous studies reporting the three-dimensional structures of Cyp18-bound oligopeptide substrates and inhibitors (10-12), a definite catalytic mechanism for cyclophilins and other PPIases is still lacking. It is currently believed that catalysis arises from a combination of desolvation of the substrate and substrate-assisted catalysis. Electrophilic assistance, for example, by the Arg55 residue of Cyp18, can also be envisaged to account for the more effective catalysis found for cyclophilins, compared to FKBP. The positively charged side-chain of this amino acid residue is located within **hydrogen-bonding** distance and perpendicular to the plane of the substrate proline ring in Cyp18/substrate complexes. This environment may cause additional weakening of the reactive C-N linkage by immobilizing the lone pair of electrons of the nitrogen atom. Indeed, a Cyp18 variant in which Arg55 has been replaced retains only 0.1% of the wild-type enzymatic activity (13).

Because of the putative functional redundancy and overlapping of the many PPIases present in most organisms, gene deletion experiments often do not lead to a recognizable **phenotype** under normal growth conditions. Among the seven cyclophilins found in *S. cerevisiae*, deletions of only two are associated with any phenotype. Disruption of *Cpr3*, the gene encoding the **mitochondrial** isoform of yeast cyclophilin, affects growth on L-lactate medium (14), whereas *cpr7D* cells are defective in normal cell growth (15). The expression of many PPIase genes seems to be sensitive to **heat shock** and to the chemical **stress response**.

Recently, it was shown that host cell cytosolic Cyp18 is required for **HIV-1** infection before **reverse transcription** but subsequent to receptor binding and membrane fusion in T cells (16). A proline-rich segment of the capsid domain of Pr55<sup>gag</sup> mediates incorporation of Cyp18 into HIV-1 virions



(17). This conserved segment that contains four proline residues occurs in the sequence -Pro-(Xaa)<sub>4</sub>-Pro<sub>222</sub>-(Xaa)<sub>2</sub>-Pro-(Xaa)<sub>5</sub>-Pro-. Mutant proteins of HIV-1<sub>HXB2</sub> that have **site-directed mutations** in which Pro<sub>222</sub> is replaced by Ala or Gly<sub>221</sub> by Ala, fail to bind to the **fusion protein** glutathione-S-transferase/Cyp18. Virions that contain these altered proteins cannot sequester Cyp18 into the released virions, emphasizing the importance of the Gly-Pro<sup>222</sup> bond for Pr55<sup>gag</sup>/Cyp18 complex formation. The ability of many cyclosporin derivatives to dissociate the Pr55<sup>gag</sup>/Cyp18 complex reveals a quantitative relationship between Cyp18 inhibition and complex decomposition. Obviously, the antiviral effect of cyclosporin A must be caused by a pathway distinctive from immunosuppression because cyclosporin A derivatives with negligible immunosuppressive activity, but high affinities for the active site of Cyp18, retain potent anti-HIV activity (18-20).

A genetic **cDNA** screen was used to identify the bovine homologue of the retina-specific cyclophilin NinaA of *Drosophila*. The membrane-localized PPIase of the secretory pathway of the fly is required for proper folding and trafficking of Rh1-6 opsin in photoreceptor cells (21). The bovine RanBP2 cyclophilin has a domain that binds to the **GTPase** Ran, as well as to red/green opsin. A still unknown modification of opsin, possibly a prolyl bond isomerization catalyzed by the Cyp-domain of RanBP2, augments and stabilizes the interaction between the Ran-binding domain and opsin. This modification is important in membrane trafficking of long-wavelength opsin in photoreceptors (22).

### Bibliography

1. G. Fischer, H. Bang, and C. Mech (1984) *Biomed. Biochim. Acta* **43**, 1101–1111.
2. G. Fischer (1994) *Angew. Chem., Int. Ed. Engl.* **33**, 1415–1436.
3. R.E. Handschumacher et al. (1984) *Science* **226**, 544–547.
4. G. Fischer et al. (1989) *Nature* **337**, 476–478.
5. N. Takahashi, T. Hayano, and M. Suzuki (1989) *Nature* **337**, 473–475.
6. R.T. Clubb, S.B. Ferguson, C.T. Walsh, and G. Wagner (1994) *Biochemistry* **33**, 2761–2772.
7. H. Ke (1992) *J. Mol. Biol.* **228**, 539–550.
8. A.P. Page, K. Macniven, and M.O. Hengartner (1996) *Biochem. J.* **317**, 179–185.
9. B. Ryffel et al. (1991) *Immunology* **72**, 399–404.
10. G. Pflügl et al. (1993) *Nature* **361**, 91–94.
11. H.M. Ke et al. (1994) *Structure* **2**, 33–44.
12. Y.D. Zhao and H.M. Ke (1996) *Biochemistry* **35**, 7362–7368.
13. L.D. Zydowsky et al. (1992) *Protein Sci.* **1**, 1092–1099.
14. E.S. Davis et al. (1992) *Proc. Natl Acad. Sci. U.S.A.* **89**, 11169–11173.
15. A.A. Duina, J.A. Marsh, and R.F. Gaber (1996) *Yeast* **12**, 943–952.
16. D. Braaten, E.A. Franke, and J. Luban, (1996) *J. Virol.* **70**, 3551–3560.
17. E.T. Franke, H.E.H. Yuan, and J. Luban (1994) *Nature* **372**, 359–362.
18. B. Rosenwirth et al. (1994) *Antimicrob. Agents Chemother.* **38**, 1763–1772.
19. S.R. Bartz et al. (1995) *Proc. Natl Acad. Sci. USA* **92**, 5381–5385.
20. C. Aberham, S. Weber, and W. Phares (1996) *J. Virol.* **70**, 3536–3544.
21. E.K. Baker, N.J. Colley, and C.S. Zuker (1994) *EMBO J.* **13**, 4886–4895.
22. P.A. Ferreira, T.A. Nakayama, W.L. Pak, and G.H. Travis (1996) *Nature* **383**, 637–640.

### Suggestions for Further Reading

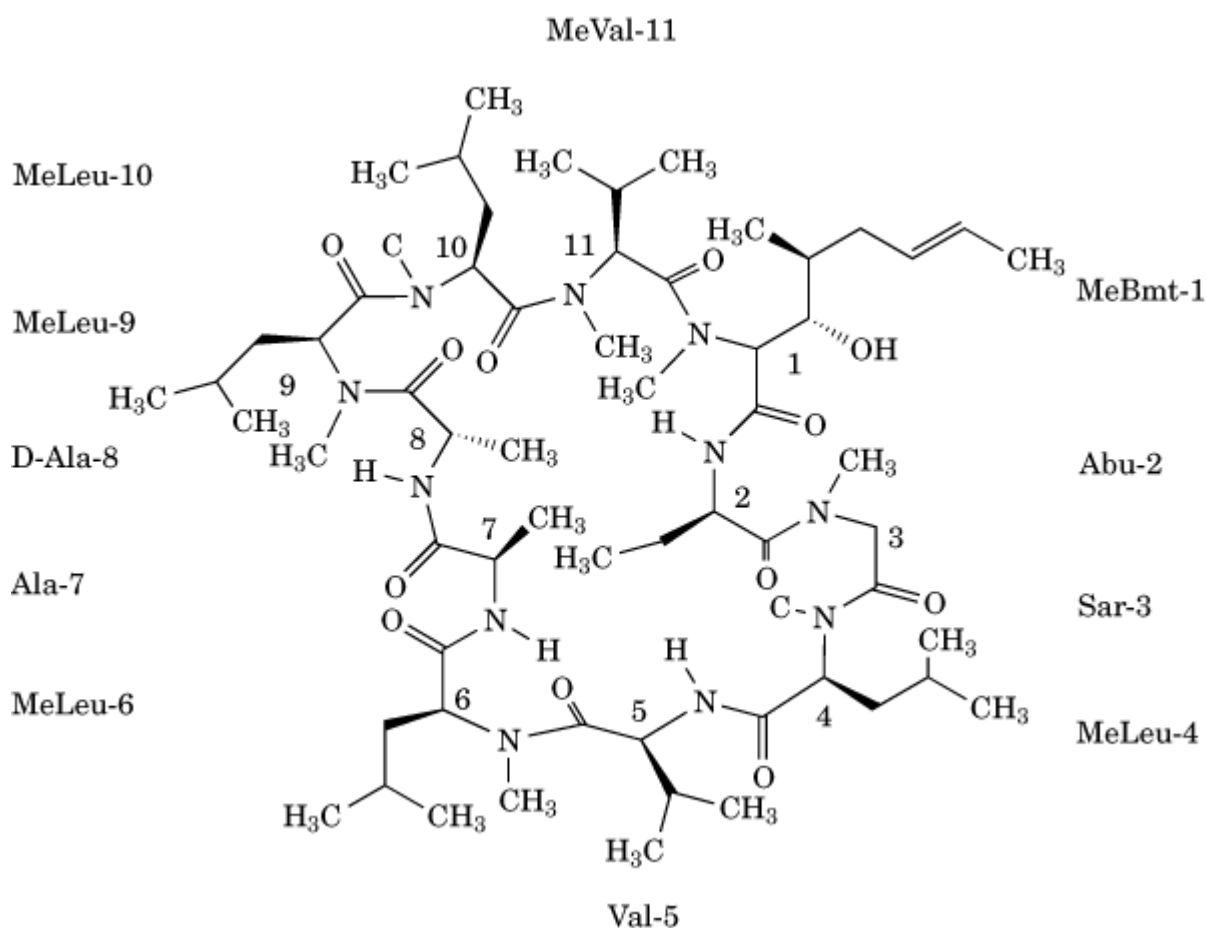
23. A. Galat and S.M. Metcalfe (1995) *Peptidylproline cis/trans isomerases*, *Prog. Biophys. Molec. Biol.* **63**, 67–118; contains an exhaustive bibliography of many aspects of PPIases.
24. C.C. Trandinh, G.M. Pao, and M.H. Saier (1992) *Structural and evolutionary relationships among the immunophilins—two ubiquitous families of peptidyl-prolyl cis/trans isomerases*,

- FASEB J. **6**, 3410–3420; contains a detailed analysis of the amino acid sequences of PPIases.
25. G. Wiederrecht and F. Etzkorn (1994) The immunophilins, *Persp. Drug Disc. Design* **2**, 57–84.
26. D.M. Armistead and M.W. Harding, (1993) Immunophilins and immunosuppressive drug action, *Ann. Rep. Medicinal Chem.* **28**, 207–213.

## Cyclosporin

The cyclosporin family of secondary metabolites from the **fungus** imperfectus *Beauveria nivea* are cyclic undecapeptides characterized by an enormous spectrum of biological effects. The most widely used derivative, cyclosporin A (CsA), is immunosuppressive, antiinflammatory, and antichemotactic, and combats multidrug-resistance and viral, fungal, and parasitic infections, without being cytotoxic to mammalian cells (1, 2). Among a number of other *N*-methyl amino acids, CsA contains an uncommon amino acid at position 1, the (4*R*)-4([(*E*)-2-butenyl]-4-*N*-dimethyl-L-threonine (Fig. 1). Several minor congeners of CsA, designated CsB, CsC ¼ CsZ, have been found as by-products of the fermentation process. The [D-MeVal<sup>11</sup>]-CsA derivative (CsH) is devoid of most biological effects, thereby serving as a control for nonspecific interactions. Other derivatives have been developed that exhibit a more restricted spectrum of biological activities compared to CsA. For example, [L-MeLeu(3-OH)<sup>1</sup>,MeAla<sup>4,6</sup>]-CsA is nonimmunosuppressive but has enhanced antiviral efficacy (3). In the fungus, non-ribosomal synthesis of cyclosporins occurs on a huge, single-chain enzyme with a molecular mass of 1,689 kDa, termed cyclosporin synthetase . It utilizes unmethylated natural amino acids, ATP/Mg<sup>2+</sup> and [S-adenosyl-L-methionine](#) as starting materials in biosynthesis (4). Beside its use as a powerful tool for unraveling the molecular basis of the cellular [immune response](#), the therapeutic application of CsA has revolutionized the likelihood of success for organ transplantation since its introduction in 1983 (5). Other disorders with immune components, like psoriasis and juvenile diabetes, are relieved with the aid of cyclosporins.

**Figure 1.** Structure of cyclosporin A.



The dramatic effects of CsA on the cellular immune response, in which the T cells remain in a nonproliferative state, results primarily from the impaired expression of a number of [lymphokine genes](#). This could include a variety of [interleukins](#) (IL-2, IL-3, IL-4), factors, such as granulocyte macrophage [colony stimulating factor](#) (GM-CSF), [tumor necrosis factor](#) (TNF- $\alpha$ ) and **g-interferon**, and the **oncogene** c-myc, all synthesized subsequent to antigen stimulation of T cells. Initial investigations identified cytosolic [cyclophilin](#) Cyp18 as the specific receptor for CsA in T cells, but minor **isoenzymes** cannot be completely ruled out yet.

The three-dimensional structure of the inhibitory CsA/Cyp18 complex was determined by [X-ray crystallographic](#) and [NMR](#) methods (6-8). The structure of Cyp18 is hardly affected by complex formation, whereas that of CsA bound to Cyp18 is dramatically different when in nonpolar organic solvents or in the solid state. In response to the altered conformation of CsA in the complex, the inhibition of the PPIase activity of Cyp18 by CsA is time-dependent, including a *cis/trans* **isomerization** of the peptide bond between MeLeu<sup>9</sup>-MeLeu<sup>10</sup>. [Site-directed mutagenesis](#) of Cyp18 indicates the involvement of a **hydrophobic** pocket of side-chains, with Trp121 as a major determinant, in the tight binding of CsA to Cyp18, with an inhibition constant,  $K_i$  of 2.3 nM.

Comparison of Cyp18/substrate complexes with Cyp18/CsA indicates that the peptide chain of CsA runs in the opposite direction. The MeVal11 side-chain of the inhibitor occupies the place of the substrate proline ring.

In some cases, binding of cyclosporins to the targeted cyclophilin and inhibition of its PPIase activity may explain the biological effect of cyclosporins. In contrast, binding of CsA to Cyp18 is necessary, but not sufficient, to arrest signals for the cellular immune response. For example, [MeAla<sup>6</sup>]-CsA is an excellent inhibitor of the PPIase activity of Cyp18 but has only 1% of the immunosuppressive

effect of CsA, whereas the weak inhibitor [MeBm<sub>2</sub>t<sup>1</sup>]-CsA (K<sub>i</sub> > 1 μM) still shows a considerable fraction of the CsA effect (9).

Thus, for immunosuppression CsA (like [FK506](#) and [rapamycin](#)) is best viewed as a pro-drug that is activated functionally when bound to Cyp18 (10). Once formed from the components, the complexes CsA/Cyp18 or FK506/FKBP12 bind to and inactivate reversibly, in a noncompetitive mechanism, the Ca<sup>2+</sup>- and calmodulin-dependent, heterodimeric protein **phosphatase 2B** (calcineurin). This enzyme catalyzes the dephosphorylation of Ser/Thr residues in phosphorylated proteins (11, 12). As a result of dephosphorylation, this enzyme induces the phosphorylated cytoplasmic subunit of NF-AT to translocate into the [nucleus](#) for association with the newly synthesized nuclear subunit of NF-AT to form a functional [transcription factor](#) (13).

### Bibliography

1. J. F. Borel, C. Feurer, H. Gubler, and H. Stahelin (1976) *Agents Actions* **6**, 468–475.
2. J. F. Borel (1989) *Pharmacol. Rev.* **41**, 259–371.
3. S. R. Bartz et al. (1995) *Proc. Natl Acad. Sci. U.S.A.* **92**, 5381–5385.
4. A. Lawen and R. Zocher (1990) *J. Biol. Chem.* **265**, 11355–11360.
5. C. R. Stiller (1996) *Transplant. Proc.* **28**, 2005–2012.
6. J. Kallen et al. (1991) *Nature* **353**, 276–279.
7. H. M. Ke et al. (1994) *Structure* **2**, 33–44.
8. S. W. Fesik et al. (1991) *Biochemistry* **30**, 6574–6583.
9. N. H. Sigal et al. (1991) *J. Exp. Med.* **173**, 619–628.
10. S. L. Schreiber (1991) *Science* **251**, 283–287.
11. J. Liu et al. (1991) *Cell* **66**, 807–815.
12. J. Liu et al. (1992) *Biochemistry* **31**, 3896–3901.
13. N. A. Clipstone, D. F. Fiorentino, and G. R. Crabtree (1994) *J. Biol. Chem.* **269**, 26431–26437.

### Suggestions for Further Reading

14. A. Lawen (1996) Biosynthesis and mechanism of action of cyclosporins, *Prog. Med. Chem.* **33**, 53–97. This nice review about cyclosporin A-mediated immunosuppression includes a summary of the biosynthesis of cyclosporin derivatives.
15. P. F. Halloran and J. Madrenas (1991) The mechanism of action of cyclosporine—a perspective for the 90s, *Clin. Biochem.* **24**, 3–7.
16. M. Thali (1995) Cyclosporins: Immunosuppressive drugs with anti-HIV-1 activity, *Mol. Med. Today*, **6**, 287–291.
17. J. Clardy (1995) The chemistry of signal transduction, *Proc. Natl Acad. Sci. U.S.A.* **92**, 56–61.
18. H. Fliri, G. Baumann, A. Enz, J. Kallen, M. Luyten, V. Mikol, R. Movva, V. Quesniaux, M. Schreier, M. Walkinshaw, R. Wenger, G. Zenke, and M. Zurini (1993) Cyclosporins—structure-activity relationships, *Immunosuppressive Antiinflammatory Drugs* **696**, 47–53.
19. G. Fischer (1994) Peptidyl-prolyl cis/trans isomerases and their effectors, *Angew. Chem., Int. Ed. Engl.* **33**, 1415–1436.

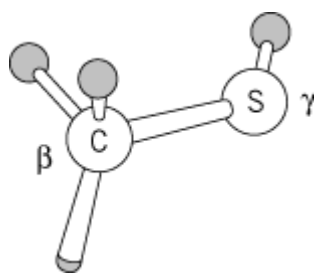
### Cystatin

Cystatin is a protein inhibitor of cysteine proteinases from chicken egg white that gave its name to a family and superfamily of such inhibitors [see [Cysteine \(Cys, C\)](#); [Proteinase Inhibitors, Protein](#)].

## Cysteine (Cys, C)

The [amino acid](#) cysteine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—UGU and UGC—and represents only 1.7% of the residues of the proteins that have been characterized. The cysteinyl residue incorporated has a mass of 103.15 Da, a **van der Waals volume** of  $86 \text{ \AA}^3$ , and an [accessible surface](#) area of  $140 \text{ \AA}^2$ . Cys residues are one of the most conserved type of residue during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [serine](#) residues.

The [thiol group](#) of the Cys side-chain



is the most reactive of any of the normal amino acids. It undergoes numerous chemical reactions with a variety of reagents (see [Thiol Groups](#)). These reactions provide means for counting and measuring the number of Cys residues present in a protein (see [Counting Residues](#) and [-Dithiobis-\(2-Nitrobenzoic Acid\), DTNB](#)). It is the ionized thiolate anion that is reactive, and only it has **absorbance** in the UV region; the  $pK_a$  value of a typical exposed Cys thiol group is 8.7, but deviations as low as 3.5 and greater than 11 are known, depending on the environment (1).

In the nonionized form, the thiol group is not very [polar](#) or reactive, and thiol groups may be buried in folded proteins without the need for any [hydrogen bonding](#). Approximately 40–50% of Cys residues are fully buried in folded [protein structures](#). They occur predominantly in [beta-sheet](#) type of **secondary structure**. Cys residues can serve as efficient “capping” residues at the *N*-termini of  $\alpha$ -helices in which their thiol group interacts with the hydrogen bonding groups of the backbone; the  $pK_a$  can then be decreased and the thiol groups made more reactive; probably for that reason, Cys residues occur at such positions in folded proteins only when important for function (2).

A notable tendency is for two Cys residues of the same or different polypeptide chains to form [disulfide bonds](#). This occurs only if the protein conformation permits, and disulfide formation is a useful probe of **protein folding** *in vitro*. *In vivo*, disulfide bonds are usually inserted primarily into proteins to be secreted by the catalyst [protein disulfide isomerase](#).

If a protein is hydrolyzed to amino acids under conditions in which the disulfide bond persists, two cysteine amino acids linked by a disulfide bond result, which historically was known as the amino acid cystine. Peptide bonds preceding Cys residues can be cleaved by cyanation with 2-nitro-5-thiocyanobenzoic acid, followed by incubation at alkaline pH (3). The reaction is specific for Cys residues but has the disadvantage that the new *N*-terminus generated is blocked and not amenable to

further **sequencing** by the [Edman Degradation](#).

Thiol groups have intrinsic affinities for many metal ions, and they are used in proteins to ligate ions of zinc, iron, and copper. In an extreme example, [metallothioneins](#) about a third of the residues are Cys. The thiol groups of Cys residues are also used as nucleophiles in catalysis in the [thiol proteinases](#).

Cys residues are the sites of several types of [post-translational modifications](#), in addition to disulfide formation. **Farnesyl groups** and geranylgeranyl groups can be added to Cys residues at the C-termini of certain proteins, after removal of a few C-terminal residues. **Palmitoyl groups** can be attached in thioester linkages to the side chains of Cys residues.

### Bibliography

1. J. W. Nelson and T. E. Creighton (1994) *Biochemistry* **33**, 5974–5983.
2. T. Kortemme and T. E. Creighton (1995) *J. Mol. Biol.* **253**, 799–812.
3. G. R. Stark (1977) *Meth. Enzymol.* **47**, 129–132.

## Cysteine Proteinase Inhibitors

Cysteine proteinase inhibitors attracted wide attention and, after serine proteinase inhibitors [see [Serine Proteinase Inhibitors, Protein](#)], have the most frequently described members. Most but not all of these belong to the [cystatin](#) superfamily, which is frequently divided into three families: the stefin, the cystatin, and the kininogen families. The stefins are low molecular weight (about 100 residues) molecules without attached carbohydrate and without any intramolecular disulfide bridges. They are extremely stable intracellular proteins. The word *stefin* comes from the name of the J. Stefin Institute in Ljubljana (Slovenia), where the first stefin was discovered. Cystatins proper are only slightly bigger than stefins (120–130 residues). They contain two disulfide bridges near their COOH terminals and typically are secretory proteins. The family and superfamily are named after chicken egg white cystatin, which was discovered long before the other cystatins became popular. The third family of cystatins is the kininogen family. Proteins are assigned to it on the basis of the high molecular weights of its members, predominantly HMW (high molecular weight) (120 kD) and LMW (low molecular weight) (68 kD) kininogens. Kininogens are glycosylated and contain numerous disulfide bridges. Each contains three cystatin like repeats.

All cystatin family members inhibit only cysteine proteinases and only those related to papain but, within this broad group, they exhibit little discrimination between individual enzymes. In strong contrast to serine proteinase inhibitors, cystatins react with enzymes whose catalytic sulfhydryl group is significantly blocked. The mode of association was elucidated in detail from the three-dimensional structure of a complex of stefin A with papain.

In avian and mammalian cells, there are present cysteine proteinases called calcium-activated neutral proteinases (CANP), or calpains. These complicated multidomain enzymes are of great importance in the metabolism of muscles. They are controlled by a closely related family of molecularly complex endogenous protein inhibitors—calpastatins.

### Suggestions for Further Reading

- A. J. Barrett (1981) Cystatin, the egg white inhibitor of cysteine proteases. *Methods Enzymol.* **80**,

771–778.

W. M. Brown and K. M. Dziegielewska (1997) Friends and relations of the cystatin superfamily—new members and their evolution. *Protein Sci.* **6**, 5–12.

H.-H. Otto and T. Schimeister (1997) Cysteine proteases and their inhibitors. *Chem. Rev.* **97**, 133–171.

## Cystine Knot

The cystine knot is a common [protein motif](#) that occurs in some [protein structures](#). It incorporates an antiparallel [b-sheet](#) and three [disulfide bonds](#), one of which passes through a ring formed by the other two disulfides (Fig. 1). Two families of cystine knot motifs have been identified. One is the [growth factor](#) cystine knot family (1, 2) that includes **nerve growth factor**, **transforming growth factor b2**, [platelet derived growth factor](#), and human chorionic gonadotrophin. In this topology, the six cysteine residues (designated CI through CVI) form three disulfide bridges (CI-IV, CII-V, CIII-VI), with CI-IV passing through the ring formed by the disulfides of CII-V and CIII-VI. The size of the ring varies from 8 to 14 residues. The b-sheet is formed from four antiparallel b-strands. These cystine knot growth factors are all dimeric, but their dimer interfaces differ.

**Figure 1.** Schematic representation of the backbone structure of the growth factor cystine knot protein, nerve growth factor (4). b-Strands are shown as arrows, and the three disulfide bonds of the cystine knot are shown in gray, with the sulfur atoms depicted as gray spheres. The *N*- and *C*-termini are labeled. This figure was generated using Molscript (5) and Raster3D (6, 7).



The second group is the inhibitor cystine knot family (3) and includes the neurotoxin [w-conotoxin GVIA](#) and the uterotonic peptide kalata B1. The inhibitor cystine knot motif has the same disulfide bond pattern as the growth factor cystine knot, but in this case the knot is formed by CIII-VI passing through the ring formed by disulfides CI-IV and CII-V. The ring of the inhibitor cystine knot varies from 8 to 12 atoms, and the b-sheet is triple stranded and antiparallel.

[See also [Disulfide Bonds](#) and [Beta-Sheet](#).]

#### Bibliography

1. N. Q. McDonald and W. A. Hendrickson (1993) *Cell* **73**, 421–424.
2. N. W. Isaacs (1995) *Curr. Opin. Struct. Biol.* **5**, 391–395.
3. P. K. Pallaghy, K. J. Nielsen, D. J. Craik, and R. S. Norton (1994) *Protein Sci.* **3**, 1833–1839.
4. N. Q. McDonald et al. (1991) *Nature* **354**, 411–414.
5. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
6. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
7. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.



## Cytochalasin

Cytochalasins are a group of fungal metabolites that can permeate cell membranes and exert very strong effects on the [cytoskeleton](#): ruffling and [translocation](#) are stopped, and cells round up (1, 2). In addition, cytochalasins A and B also inhibit glucose transport across the plasma membrane. Although it was originally believed that the main action of the cytochalasins on the cytoskeleton was binding with high affinity to the barbed, or fast-growing end of the [actin](#) filament, and inhibiting subunit addition and loss, it has become clear that the effects are more complex. The concentrations of cytochalasins used in cell biological experiments will have strong effects on actin monomers and may convert much of a cell's monomeric G-ATP actin to G-ADP actin (3). In addition, even at saturating concentrations cytochalasins do not completely block the fast-growing ends with respect to subunit addition or loss (3). Thus, the action of cytochalasin on a living cell is likely to be complicated and may affect the interactions of actin with many [actin-binding proteins](#). Nevertheless, as a result of the large effect cytochalasins have on actin, they have been widely used to study the role of the actin cytoskeleton in cellular processes.

### Bibliography

1. I. Yahara, F. Harada, S. Sekita, K. Yoshihira, and S. Natori (1982) *J. Cell Biol.* **92**, 69–78.
2. M. Schliwa (1982) *J. Cell Biol.* **92**, 79–91.
3. P. Sampath and T. D. Pollard (1991) *Biochemistry* **30**, 1973–1980.

## Cytochrome P450

Cytochromes P450 (designated here as P450) arise from a [superfamily](#) of **genes** that encode hemoproteins with molecular weights of 40 to 55 kDa and extraordinarily diverse properties and functions. P450 is nearly ubiquitous in the biotic world, with over 500 genes characterized in **plants**, animals, **fungi**, and **bacteria**, although it is not present in all organisms. Of the substrates of P450, 1000 are known, and there are perhaps as many as  $10^6$  in total; more than 40 different reactions are catalyzed (1). These enzymes play a critical role in:

1. Detoxification of bioactive molecules in vertebrates, insects, and plants
2. Metabolic activation of compounds that are both beneficial (drug precursors, [hormones](#)) and detrimental ([carcinogens](#))
3. Biogenesis and catabolism of many endogenous compounds, including [steroid hormones](#), vitamin D, biological insecticides, [fatty acids](#), insect pheromones, and plant lignins and pigments

The P450 that metabolize xenobiotics appear to have no critical physiological functions, since mice with disrupted genes for these proteins had no obvious phenotypic abnormalities (2, 3). Most P450 proteins are membrane-bound in the [endoplasmic reticulum \(microsomes\)](#) or [mitochondria](#), but soluble forms are present in bacteria. In mammals, the greatest concentrations of P450, particularly the detoxification forms, are present in the liver, but they are probably present in most tissues of the body, with relatively high concentrations in the intestines, lungs, nasal epithelia, and steroidogenic tissues. The genes for these enzymes differ in size, **intron** number, and **chromosomal** location. **Gene expression** of subsets of the P450 is regulated by a variety of substances, including

xenobiotics such as aromatic hydrocarbons, barbiturates, and **peroxisomal** proliferators, and by endogenous steroid and polypeptide hormones.

The name cytochrome P450 is based on the characteristic increase in the absorbance of the protein heme group at 450 nm upon binding of CO and reduction of its iron atom. P450 differ structurally from most [cytochromes](#) because the fifth coordination group of the iron atom of the heme is occupied by a thiolate group from a cysteine residue and the sixth is occupied by water. They differ functionally because P450 are monooxygenases and not simply electron carriers (see [Cytochromes](#)).

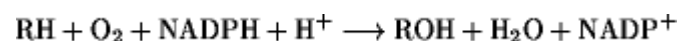
## 1. Classification

On the basis of amino acid sequence similarity, the P450 are divided into more than 100 families and subfamilies (4). The various P450 within a [gene family](#) generally have >40% sequence identity, and members within a subfamily have identities >55% within mammalian species and >46% overall. P450 with >97% identity are arbitrarily considered **allelic** variants unless there is independent evidence for separate gene loci. Genes from animals, fungi, plants, or bacteria are assigned family numbers 1–49, 51–69, 71–99, and 101 and above, respectively. Family 51 is represented in plants and mammals, as well as fungi.

The recommended nomenclature for P450 genes is an italicized *CYP* (*Cyp* for mouse and *Drosophila*) followed by an Arabic number designating the family, a letter for the subfamily, and a second number for the member of the subfamily, eg, *CYP1A1* (4). Identical terminology is used for gene products, except that the name is not italicized. Trivial names of P450 and other designations of the proteins, such as cytochrome P450 1A1, P450 1A1 or 1A1, are also considered acceptable in publications so long as the official name is specified.

## 2. Catalytic Activity

The overall reaction catalyzed by P450 is



where R is an organic substrate. One oxygen atom of the O<sub>2</sub> is incorporated into the substrate and the other is reduced to water, so P450 is both a monooxygenase and mixed-function oxidase. Further reactions and rearrangements lead to more than 40 different types of chemical reactions, including aliphatic and aromatic oxidations, N-, S-, and O-dealkylations, oxidative deamination, and peroxidation.

The P450 reaction cycle starts with the binding of the substrate to the ferric form of P450, which is followed by reduction of the heme iron atom by acceptance of an electron, binding of O<sub>2</sub>, activation of the oxygen by acceptance of a second electron, and scission of the oxygen molecule, followed by formation of water, a hydroxylated product, and ferric P450. The two electrons from NADPH are donated either to microsomal P450, through a flavoprotein, P450 reductase, or to mitochondrial and bacterial P450, sequentially through a flavoprotein (a ferredoxin reductase) and an iron-sulfur [ferredoxin](#). In rare cases, the P450 and P450 reductase are fused into a single protein (5). Cytochrome b<sub>5</sub> increases the activity of some P450 *in vitro*, either as a source of the second electron or **allosterically** (6).

## 3. Structure

The catalytic domain of membrane-bound P450 is oriented toward the cytoplasm for the microsomal forms and toward the matrix for the mitochondrial forms. The microsomal P450 are inserted into the membrane in a cotranslational, signal recognition particle-dependent manner (7). The N-terminal 20 to 25 amino acid residues function as a membrane-insertion and halt-transfer signal and are not

cleaved, so they anchor the protein to the membrane. The mitochondrial P450 contain standard cleavable mitochondrial targeting signals at their N-terminus, but the mode of insertion into the inner membrane of the mitochondria is not known.

The three-dimensional structures of soluble bacterial P450, but not any membrane-bound P450, have been determined. Although the primary sequence is dissimilar among these bacterial proteins, a common “P450 fold” is present in each (see [Protein Structure](#)), which suggests that a similar core structure will also be present in eukaryotic membrane-bound P450 (8). The P450 structure contains an **alpha-helical** region and a [beta-sheet](#) region, which constitute about 70 and 20%, respectively, of the polypeptide chain. The sequence FXXGXXXCXG, which includes the heme-binding [cysteine](#) residue, is conserved in nearly all P450. Structures of eukaryotic P450 have been predicted from amino acid sequence alignments and homology molecular modeling with bacteria P450 (see [Protein Structure Prediction](#)), and predictions of residues that interact with substrates have been generally consistent with the observed effects of mutations (9, 10).

#### 4. Gene Structure

The sizes of the eukaryotic P450 genes, from the [transcription](#) initiation site to the [polyadenylation](#) site, vary dramatically, from about 3 kb for the *CYP21* genes to >70 kb for *CYP19* genes (11). In higher eukaryotes, the number of **introns** varies from 6 to 13, and the sites of insertion of introns are almost always conserved within a family, but only rarely between families. **Polymorphism** is common in P450 genes and underlies genetic differences in drug metabolism and congenital diseases in humans. The families of P450 genes are dispersed among at least nine chromosomes in humans and mice, but occasionally families cluster together, such as families 1, 11, and 19. Subfamilies within a family also may be dispersed to separate chromosomes or may cluster together, and members of the *CYP2A*, *CYP2B*, and *CYP2F* subfamilies are intermingled (4, 12). The [gene duplication](#) leading to the *CYP21A* locus included the C4 [histocompatibility](#) gene, so the two P450 genes alternate with the C4 gene. In spite of the diversity of P450 genes, their similarity in [primary structures](#), the common three-dimensional fold, and some conserved intron sites among P450 genes support the idea that the P450 superfamily evolved by **divergence** from a common ancestral gene.

#### 5. Gene Expression

P450 genes are expressed in tissue-, developmental-, and sex-specific patterns. The mechanisms regulating these patterns of expression are best understood in the liver and adrenal gland. An **orphan receptor** containing a **zinc finger**, termed Ad4BP or steroid factor-1 (SF-1), has been shown to be critical for the developmental and tissue-specific expression of P450 genes that are required for the biogenesis of adrenal and sex steroid hormones (13). The SF-1 gene has been disrupted in mice, with major effects on differentiation of the adrenal gland and gonads, pituitary gonadotropin function, and the ventral medial hypothalamus (14), so this gene has important developmental functions as well as effects on the expression of P450 genes. *CYP19*, aromatase, contains SF-1 sites, but its tissue-specific expression is also dependent on alternate **promoters** and transcription initiation sites (15).

All the known liver-enriched regulatory factors have been implicated in the expression of one or more P450 genes. Examples in which functional binding sites for liver-enriched factors have been demonstrated include HNF-1a for *CYP2E1*, HNF-4 for *CYP2C1/2/3*, DBP for *CYP2C6* and *CYP7*, HNF-3 for *CYP2C6*, C/EBPb for *CYP2D5*, and C/EBPa and C/EBPb for *CYP2B1/2* (16, 17). A variety of ubiquitously expressed factors also contribute to basal expression of the *CYP* genes. The developmental expression of *CYP2C6* correlates with the developmental activation of DBP (16), and the diurnal expression of DBP underlies the diurnal regulation of *CYP7* (18). Activation of *CYP2D5* by C/EBPb also requires the constitutive factor Sp1, which has a binding site near that of C/EBPb (16).

##### 5.1. Sex-dependent Expression

The sex-dependent expression of several P450 genes in rodents may be mediated by more than one

mechanism. In rats, the sex-specific expression of *CYP* genes is dependent on different patterns of [growth hormone](#) secretion in the two sexes. STAT [transcription factors](#) (19) and [phospholipase](#)<sub>A2</sub> (20) may be involved in the growth hormone-regulated male-specific expression of specific *CYP* genes, but additional studies are required to establish the exact mechanism. In mice, a sex difference information (SDI) element is present in male-specific *Cyp2d9* and female-specific *Cyp2a4* (21). The SDI contains a CpG sequence that is **methyated** in a sex-dependent way, and a factor has been identified that binds in a methylation-dependent manner (22). The exact role of this factor is not yet resolved, as its tissue and developmental expression differs from that of the *CYP* genes.

## 5.2. Aromatic Hydrocarbons

The expression of subsets of the xenobiotic-metabolizing P450 genes is regulated by aromatic hydrocarbons, peroxisomal proliferators, and barbiturates. The induction of *CYP1A1* by aromatic hydrocarbons, including the halogenated aromatic hydrocarbon 2,3,7,8-tetrachlorodibenzo-*p*-dioxin, is best understood. Genetic differences in the induction of *Cyp1a1* in mice and in variants of cultured mouse hepatoma Hepa 1 cells contributed critically to the identification of a regulatory locus *AhR*, which encodes the aromatic hydrocarbon receptor, and a second gene, the Ah receptor nuclear translocator *Arnt*, which encodes a protein required for the **nuclear import** of AhR (23, 24). The AhR is complexed with the **molecular chaperone** hsp90 in the cytoplasm in uninduced cells, which is essential to maintain a functional AhR. Binding of the ligand results in the dissociation of hsp90 and localization of AhR in the nucleus. Arnt is a nuclear protein that binds to liganded AhR and is required for the binding of AhR to the DNA binding site (25). AhR/Arnt bind to multiple sites within an [enhancer](#) about 1 kb from the start site of transcription and disrupt the **nucleosomal** structure of the enhancer (23). In turn, the nucleosomal structure of the promoter region is also disrupted in a manner dependent on the **transactivation** domain of AhR, which then allows the binding of regulatory proteins to the proximal promoter. The core sequence of the binding sites, which have been designated xenobiotic-, dioxin-, or Ah-**response elements** by different investigators, is 5'-TNGCGTG-3'. Arnt binds to the core sequence, and AhR binds to the flanking sequence. Both AhR and Arnt contain (1) basic loop-helix-loop (bLHL) motifs near their N-termini; (2) PAS **domains**, which exhibit homology with the *Drosophila* proteins, Per and Sim; and (3) glutamine-rich, putative transcriptional activation domains in the C-terminal region (24). The ligand and hsp90 interact with sequences near the PAS domain of AhR, and the binding of AhR to Arnt involves the bHLH and PAS domains of each protein. Arnt is capable of transactivation independently, but AhR transactivation requires the association of AhR with Arnt.

## 5.3. Peroxisomal Proliferators

Induction of expression of the *CYP4A* genes by peroxisomal proliferators and fatty acids requires peroxisomal proliferator-activated receptor- $\alpha$  (PPAR $\alpha$ ), a member of the steroid [hormone receptor](#) family (26). CYP4A P450 catalyze the  $\omega$ -hydroxylation of fatty acids, including arachidonic acid, which may ultimately result in the degradation of the fatty acids by peroxisomes. The induction of *CYP4A* genes is defective in mice with disrupted PPAR $\alpha$  genes (27). PPAR $\alpha$  forms a heterodimer with retinoid-X-receptor (RXR) to activate the gene. The predominant sequence in *Cyp4A6* that confers the response to peroxisomal proliferators contains an imperfect repeat of the sequence AGGTCA, with a single base pair between the repeats and additional sequence to the 5' side of the repeat. The 5' sequence appears to be important for the specificity of binding of PPAR $\alpha$ /RXR $\alpha$  to the nonconsensus, imperfect repeat of the peroxisomal proliferator response element present in *CYP4A6* (26).

## 5.4. Barbiturates

Induction of P450 genes by barbiturates is not well understood. Neither a barbiturate receptor nor specific barbiturate response elements in mammalian genes have been defined, and different mechanisms may be involved for different P450 (28). The system best understood is the induction by barbiturates of Bm-1, Bm-2, and Bm-3 P450 in *Bacillus megaterium* (29). In untreated cells, a repressor, Bm3R1, binds to a palindromic 20-bp sequence upstream of a promoter that drives the expression of both the repressor and Bm-3. In the presence of barbiturate, the binding is reversed,

which may result from direct effects on the repressor, but positively acting factors are also induced that compete for the binding of the repressor at the palindromic site and at sites termed “barbie boxes” present in the 5' regions of the Bm-3 and Bm-1 genes. In the mammalian *CYP2B1/2* genes, evidence has been presented for phenobarbital responsive sequences in both the proximal promoter region, including barbie-box-like sequences (30), and in distal phenobarbital-responsive enhancers (31, 32). Identification of specific responsive elements and characterization of the cognate binding proteins will be required to establish the role of these regions in phenobarbital induction.

### 5.5. Glucocorticoids

A number of *CYP* genes are also modulated by [glucocorticoids](#). In some cases, the effect is mediated through the classical glucocorticoid receptors, such as modulation by glucocorticoids of aromatic hydrocarbon induction of *CYP1A1* and phenobarbital induction of *CYP2B* genes (33). In other cases, such as *CYP3A* genes, the response to glucocorticoids occurs by nonclassical mechanisms that remain to be elucidated (34).

### 5.6. Cyclic AMP

Genes for P450 involved in steroid metabolism are induced by polypeptide hormones via **cyclic AMP**-mediated (cAMP) mechanisms (35). In the adrenal gland, *CYP11A*, *CYP11B*, *CYP17*, and *CYP21* are regulated in this way. Interestingly, even though each of these genes probably evolved from a common ancestor, different regulatory elements in the promoter of each gene mediate the cAMP response. A classical CREB binding site is present in *CYP11B*, but binding sites for proteins other than CREB are present in other genes, and different sites are present in *orthologous genes* in different species. In bovine *CYP11A*, a binding site for Sp1, a constitutive factor, mediates the cAMP response (36), and members of the **homeodomain** PBX family of [helix-turn-helix motif](#) proteins bind to the cAMP-responsive element in *CYP17* (37). Multiple proteins bind to many of these elements, and those that predominate *in vivo* remain to be identified.

## Bibliography

1. M. J. Coon, A. D. N. Vaz, and L. L. Bestervelt (1996) *FASEB J.* **10**, 428–434.
2. H. C. L. Liang et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1671–1676.
3. S. S. T. Lee, J. T. M. Buters, T. Pineau, P. Fernandez-Salguero, and F. J. Gonzalez (1996) *J. Biol. Chem.* **271**, 12063–12067.
4. D. R. Nelson et al. (1996) *Pharmacogenetics* **6**, 1–42.
5. L. O. Narhi and A. J. Fulco (1987) *J. Biol. Chem.* **262**, 6683–6690.
6. H. Yamazaki, W. W. Johnson, Y.-F. Ueng, T. Shimada, and F. P. Guengerich (1996) *J. Biol. Chem.* **271**, 27438–27445.
7. B. Kemper and E. Szczesna-Skorupa (1989) *Drug Metabol. Rev.* **20**, 811–820.
8. C. A. Hasemann, R. G. Kurumbail, S. S. Boddupalli, J. A. Peterson, and J. Deisenhofer (1995) *Structure* **2**, 41–62.
9. O. Gotoh (1992) *J. Biol. Chem.* **267**, 83–90.
10. C. von Wachenfeldt and E. F. Johnson (1995) In *Cytochrome P450—Structure, Mechanism, and Biochemistry* (P. Ortiz de Montellano, ed.), Plenum Press, New York, pp. 183–223.
11. B. Kemper (1993) In *Frontiers in Biotransformation* (K. Ruckpaul and H. Rein, eds.), Vol. **8**, Akkademie Verlag, Berlin, pp. 1–58.
12. S. M. G. Hoffman, P. Fernandez-Salguero, F. J. Gonzalez, and H. W. Mohrenweiser (1995) *J. Mol. Evol.* **41**, 894–900.
13. D. S. Lala, D. A. Rice, and K. L. Parker (1992) *Mol. Endocrinol.* **6**, 1249–1258.
14. K. L. Parker and B. P. Schimmer (1996) *Trends Endo. Metab.* **7**, 203–207.
15. M. S. Mahendroo, C. R. Mendelson, and E. R. Simpson (1993) *J. Biol. Chem.* **268**, 19463–19471.
16. F. J. Gonzalez and Y.-H. Lee (1996) *FASEB J.* **10**, 1112–1118.

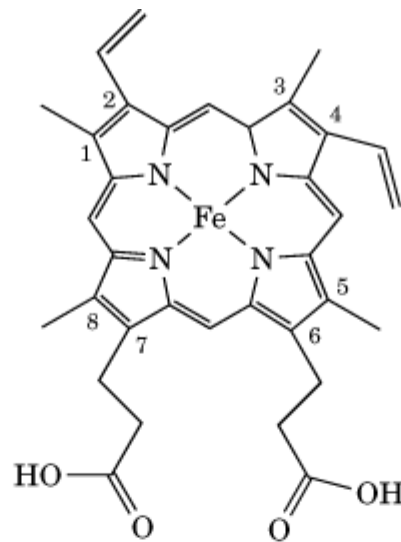
17. Y. Park and B. Kemper (1996) *DNA Cell Biol.* **15**, 693–701.
18. D. J. Lavery and U. Schibler (1993) *Genes Develop.* **7**, 1871–1884.
19. D. J. Waxman, P. A. Ram, S. H. Park, and H. K. Choi (1995) *J. Biol. Chem.* **270**, 13262–13270.
20. P. Tollet, M. Hamberg, J. Å. Gustafsson, and A. Mode (1995) *J. Biol. Chem.* **270**, 12569–12577.
21. H. Yoshioka, M. Lang, G. Wong, and M. Negishi (1990) *J. Biol. Chem.* **265**, 14612–14617.
22. N. Yokomori, R. Kobayashi, R. Moore, T. Sueyoshi, and M. Negishi (1995) *Mol. Cell Biol.* **15**, 5355–5362.
23. J. P. Whitlock Jr. et al. (1996) *FASEB J.* **10**, 809–818.
24. O. Hankinson (1995) *Ann. Rev. Pharmacol. Toxicol.* **35**, 307–340.
25. R. S. Pollenz, C. A. Sattler, and A. Poland (1993) *Mol. Pharmacol.* **45**, 428–438.
26. E. F. Johnson, C. N. A. Palmer, K. J. Griffin, and M.-H. Hsu (1996) *FASEB J.* **10**, 1241–1249.
27. S. S.-T. Lee et al. (1995) *Mol. Cell Biol.* **15**, 3012–3022.
28. D. J. Waxman and L. Azaroff (1992) *Biochem. J.* **281**, 577–592.
29. Q. Liang and A. J. Fulco (1995) *J. Biol. Chem.* **270**, 18606–18615.
30. L. Prabhu et al. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9628–9632.
31. E. Trottier, A. Belzil, C. Stoltz, and A. Anderson (1995) *Gene* **158**, 263–268.
32. Y. Park, H. Li, and B. Kemper (1996) *J. Biol. Chem.* **271**, 23725–23728.
33. R. A. Prough, M. W. Linder, J. A. Pinaire, G.-H. Xiao, and K. C. Falkner (1996) *FASEB J.* **10**, 1369–1377.
34. L. C. Quattrochi, A. S. Mills, J. L. Barwick, C. B. Yockey, and P. S. Guzelian (1995) *J. Biol. Chem.* **270**, 28917–28924.
35. M. Waterman (1994) *J. Biol. Chem.* **269**, 27783–27786.
36. P. Venepally and M. Waterman (1995) *J. Biol. Chem.* **270**, 25402–25411.
37. N. Kagawa, A. Ogo, Y. Takahashi, A. Iwamatsu, and M. R. Waterman (1994) *J. Biol. Chem.* **269**, 18716–18719.

### Suggestions for Further Reading

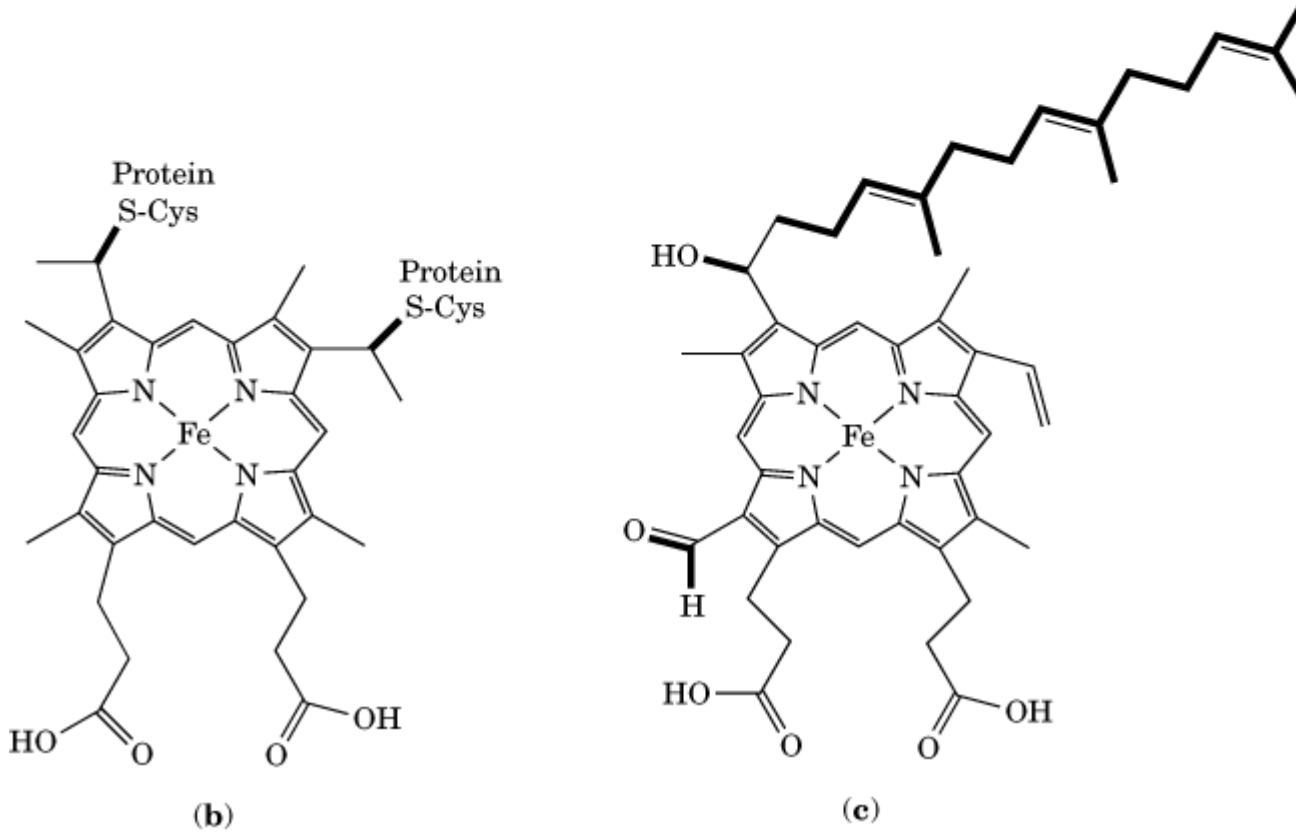
38. FASEB Journal series of 12 reviews, “*The cytochrome P450: structure, function, regulation, and genetics*” (1996/97) Initial review in *FASEB J.* (1996) **10**, 205. This series presents four reviews each on the structures of P450, regulation of xenobiotic-metabolizing P450, and regulation of P450 that metabolize endogenous compounds.
39. C. Ioannides, ed. (1996) *Cytochromes P450—Metabolic and Toxicological Aspects*, CRC Press, New York. Comprehensive reviews of the drug metabolizing P450, primarily families 1–4.
40. T. Omura, Y. Ishimura, and Y. Fujii-Kuriyama, eds. (1993) *Cytochrome P450*, VCH, New York. A concise introduction to P450.
41. P. R. Ortiz de Mantellano, ed. (1995) *Cytochrome P450—Structure Mechanism, and Biochemistry*, 2nd ed., Plenum Press, New York. Excellent comprehensive reviews of P450, including chapters on the regulation of gene expression by xenobiotics, hormones, and cAMP-mediated mechanisms.
42. M. A. Schuler (1996) Plant cytochrome P450 monooxygenases. *Crit. Rev. Plant Sci.* **15**, 235–284. A comprehensive review of plant P450.

Cytochromes are iron-containing proteins that facilitate the movement of electrons in a wide variety of metabolic processes (see [Electron Transfer Proteins](#)). Cytochromes are widely distributed throughout nature. They are found in all **eukaryotes** and in most, but not all, **prokaryotes**. The prosthetic group or chromophore of cytochromes is heme, an iron-containing porphyrin. The prototypic form of heme in cytochromes is protoheme IX, which consists of iron, two vinyl side chains, four methyl groups and two propionic acid side chains on a conjugated tetrapyrrole ring (Fig. [1](#)). Protoheme IX is the prosthetic group of *b*-type cytochromes and the family of proteins known as [cytochrome P450](#). The *c*-type cytochromes have heme *c* as the prosthetic group, in which protoheme IX is covalently bound to the protein through one or, more commonly, two thioether bonds to [cysteine](#) side-chains. Heme *a*, found in *a*-type cytochromes, has a long isoprenoid tail substituted on one of the vinyl groups and a formyl group replacing a methyl. Further variations in heme structure are found in isolated cases, for example, cytochrome *d* and cytochrome *o* ([1](#)) but will not be discussed here. Cytochrome P450 and **cytochrome *c*** are discussed elsewhere in this volume.

**Figure 1.** The structures of (a) protoheme IX, (b) heme *c*, and (c) heme *a*.



(a)



(b)

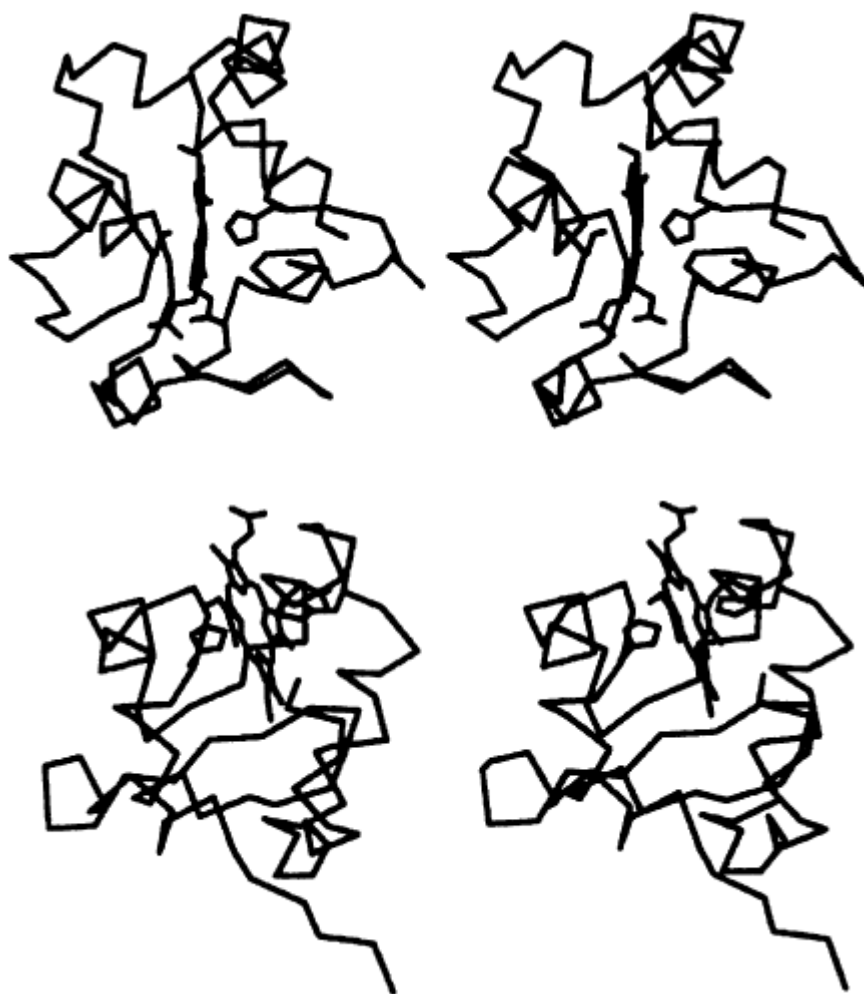
(c)

The iron in cytochromes, is either penta- or hexa-coordinate. Porphyrin nitrogen atoms fill four coordination positions (in-plane coordination, see Fig. 1) and appropriate ligands are positioned below the heme plane (5-coordinate) and, in most cases also above the heme plane (6-coordinate). The out-of-plane ligands in typical *c*-type cytochromes, are the N<sub>3</sub> atoms of [histidine](#) residues and the sulfur atom of [methionine](#) (Fig. 2). Ligation in the *b*-type cytochromes is generally by two His residues (Fig. 2) and by one or two His in the cytochromes *a*. Variations exist, however. For example, the out-of-plane ligation is by two His residues in cytochrome *c*<sub>3</sub>, and by one His and an open position in cytochrome *c*'. Cytochromes are associated with the inner membrane of [mitochondria](#), that of [microsomes](#), the [thylakoid](#) membrane in [chloroplasts](#), and the plasma



membrane in **bacteria** (see [Membranes](#)). However, soluble forms (typically *c*-type cytochromes) are found in the space between the inner and outer mitochondrial membranes and in the **periplasmic** space in bacteria.

**Figure 2.** Stereo representations of the structure of tuna cytochrome *c* (upper stereo pair) and microsomal cytochrome *b<sub>5</sub>* (lower stereo pair). In cytochrome *c*, the N-terminal  $\alpha$ -helix is at the *top*, the Met80 ligand is on the *left*, and the His18 ligand on the *right*. In cytochrome *b<sub>5</sub>*, the N-terminus is at the lower *right* the His61 ligand is on the *left*, and the His37 ligand is on the right. The heme propionates are oriented toward the solvent.



The physiological activity in all cytochromes is the reversible oxidation and reduction of the iron atom, which cycles between the ferric ( $\text{Fe}^{3+}$ ) and ferrous ( $\text{Fe}^{2+}$ ) states. From a physiological standpoint, an important property of cytochromes is their oxidation-reduction (**redox**) potential, which is  $E_0'$  at pH 7. The redox potential measures the relative stabilities of the oxidized (ferric) and reduced (ferrous) forms and therefore their tendencies to donate and accept electrons, respectively. At pH 7 in aqueous solution, the redox potential of  $\text{H}_2$  is  $-414$  mV and that of  $\text{O}_2$  is  $+816$  mV. With relatively minor exceptions, these define the biological range. Cytochromes have redox potentials somewhere between these limits. In biological electron transfer, the direction of electron flow is from lower to higher potentials in the limit from hydrogen to oxygen. Thus the redox potential of a particular cytochrome determines, to a first approximation, where it is positioned relative to other redox compounds in electron transfer pathways (see [Electron Transfer Proteins](#)). The *c*-type

cytochromes have redox potentials ranging from  $-400$  to  $+450$  mV, depending on the structural family and species, although they typically vary between only  $+200$  and  $+300$  mV in eukaryotic *c*-type cytochromes. The *b*-type cytochromes have redox potentials that range from  $-200$  to  $+150$  mV, and the *a*-type cytochromes have potentials in the  $+220$  to  $+400$  mV range.

The role of the protein moiety in cytochromes is threefold. First, the protein provides the out-of-plane ligands and, in the case of *c*-type cytochromes, the covalent attachment sites for the heme. Secondly, the protein provides the proximal environment for the heme, a nonpolar molecule that is not water-soluble, controls the exposure of the heme to the solvent, and interacts with the heme to modulate the redox potential. Thirdly, the protein surface provides the sites for the interaction with electron donors and acceptors because in all known cases the heme iron is buried in the protein's interior. Thus, the redox potential is modulated by a combination of the nature of the out-of-plane ligands, the immediate protein environment, and the solvent exposure of the heme.

The redox potentials dictate the direction of electron flow. For example, mitochondrial cytochrome *c* ( $E'_0 = +260$  mV) cannot significantly reduce oxidized pyridine nucleotide ( $\text{NAD}^+$ ,  $E'_0 = -320$  mV). They do not, however, explain the specificity and efficiency of biological electron transfer. Within a biological milieu, a large number of redox compounds exist close to each other. These include small organic compounds, such as flavins and pyridine nucleotides (for example,  $\text{NAD}^+$ ); oxygen (in aerobes); inorganics, such as nitrates, nitrites, sulfates, and sulfhydryls; and a variety of redox proteins, including non-heme-iron sulfur proteins, flavoproteins, cytochromes, and copper proteins (see [Electron Transfer Proteins](#)). Thus, from an energetic standpoint, substantial opportunities exist for electrons to flow in a great variety of directions. This does not happen, and electron flow occurs with very high efficiency (approaching 100%) only within specific pathways (for example, in respiration or [photosynthesis](#)). This incredible efficiency, the essence of life, is mediated by the surface of the protein, which provides recognition sites that determine which redox compounds interact on a physiological timescale, that is, it determines how fast electrons are transferred, and the direction of electron flow ([2](#), [3](#)). This kinetic control means that energetically equivalent reactions can take place on timescales that differ by many orders of magnitude. For example, reduced pyridine nucleotide reduces mitochondrial cytochrome *c* slowly (in seconds), whereas cytochrome *c*<sub>1</sub> ( $E'_0 = +220$  mV) reduces the same cytochrome in milliseconds, even though reduction by pyridine nucleotide is more strongly favored energetically in the first case ( $65.7$  kJ/mol versus  $3.8$  kJ/mol).

The interactive domains of the cytochromes are defined by the amino acid side chains positioned on their molecular surface that dictate their sites of interaction. The structure of the complex between two electron transfer proteins determines the proximity and orientation of the relevant prosthetic groups (the two hemes in the case of electron transfer between two cytochromes). Three factors play a role in defining the interactive domains. First, the positioning of charged amino acid side chains on the surface of interacting molecules result in electrostatic interactions, either attractive or repulsive, that define the interaction domain ([4](#)). Secondly, the distance between redox centers (heme irons in the case of cytochromes) influences the rate of electron transfer. This distance is controlled by the size and nature of the amino acid side chains in the interactive domain ([2](#)). Third, two interactive domains that interact must be complementary physically and structurally, even if defined by nonpolar [van der Waals interactions](#) ([2](#), [3](#)). Apparently, during the course of evolution, mutations have occurred that optimize the rate of electron flow through relevant metabolic pathways by modulating the amino acid side chains defining the interactive domains, thus allowing competing redox reactions to occur in close proximity and channeling the flow of electrons through the appropriate pathway.

From the standpoint of solubility, cytochromes exist in three types of environments. Soluble forms are typically *c*-type cytochromes that are free to diffuse between electron donors and acceptors. Integral [membrane proteins](#), that is, transmembrane proteins, have only part of the molecule exposed to solvent. The majority of the molecule is integrated into the membrane. Examples are *a*-type cytochromes and some *b*-type cytochromes. Peripheral membrane proteins are strongly associated

with a membrane or with integral membrane proteins and have a substantial portion of the molecule in solution. Examples are cytochrome  $b_5$  and cytochromes  $c_1$  and  $f$ .

## 1. A-type Cytochromes

Cytochromes  $a$  and  $a_3$  are components of cytochrome oxidase, or complex IV, of the mitochondrial electron transfer chain and are found in all eukaryotic organisms and many aerobic prokaryotes. Cytochrome oxidase catalyzes the four-electron reduction of  $O_2$  to  $H_2O$ , utilizing cytochrome  $c$  as the electron donor. Cytochrome  $a$ , which is hexa-coordinate, mediates the transfer of electrons from cytochrome  $c$  to the penta-coordinate cytochrome  $a_3$  (with the intermediate participation of a copper atom;), which provides the site for oxygen reduction. Variants of cytochrome oxidase are found in some aerobic bacteria that do not contain hemes other than heme  $a$ . Examples are cytochrome  $bo$ , which uses quinols as the electron donor in *Escherichia coli* and cytochrome  $cbb_3$  in *Rhodobacter sphaeroides*, which has cytochrome  $c$  as an electron donor (1). Thus, the heme type can vary, although the function, the reduction of molecular oxygen to water, remains constant. Cytochrome oxidase is an integral membrane protein and couples the transfer of electrons and the corresponding redox energy with the movement of protons across the inner mitochondrial or periplasmic membrane, driving the synthesis of ATP (see [Proton Motive Force](#)). From an evolutionary standpoint, there is structural and amino acid sequence homology between eukaryotic and bacterial cytochrome oxidase (1). In addition, a variety of other enzymes exist that are structurally related to cytochrome oxidase and generally use cytochrome  $c$  or an electron donor, such as nitrous oxide reductase, which reduces  $N_2O$  to  $N_2$  and utilizes copper as the prosthetic group, but has sequence homology with cytochrome oxidase and utilizes cytochrome  $c$  as its electron donor (1). Another example is nitric oxide reductase, which catalyzes the reduction of [nitric oxide](#) to nitrous oxide. Although it is a cytochrome  $bc$  complex, it has a transmembrane component and homology to portions of cytochrome oxidase (1).

## 2. B-type Cytochromes

The  $b$ -type cytochromes fall into two broad categories: (1) a family of relatively large cytochromes containing protoheme IX that are part of integral membrane protein complexes like cytochrome  $bc_1$  and cytochrome  $b_6f$ ; (2) a family of cytochrome  $b_5$ -like proteins that are peripheral proteins, typically anchored to a membrane, or are components of protein complexes like flavocytochrome  $b_2$  and sulfite oxidase. The cytochrome  $bc_1$  complex and its homologue in plants, cytochrome  $b_6f$ , catalyze the transfer of electrons from quinones to cytochrome  $c$  (in animals) or to plastocyanin (in plants) coupled with the movement of protons across the inner mitochondrial or the thylakoid membranes, which drives the synthesis of ATP (see [ATP Synthase](#)). The mitochondrial cytochrome  $bc_1$  complex, also known as complex III, has the cytochrome  $b$  component integrated into the membrane through a number of transmembrane  $\alpha$ -helices (5). The cytochrome  $c_1$  component is anchored to the membrane and is the electron donor to cytochrome  $c$ . Cytochrome  $b$  is also a component of succinate dehydrogenase (also known as complex II), a complex integral membrane protein that participates in the mitochondrial electron transfer pathway (6). Succinate dehydrogenase oxidizes succinate to fumarate. The electrons reduce quinone and feed into the mitochondrial electron transfer chain.

Because of their large size and strong association with membranes, cytochromes  $b$  are not as well understood as the generally small and water-soluble  $c$ -type cytochromes (cytochromes  $c_1$  and  $f$  are exceptions in terms of size and solubility). Thus the extent of structural and functional diversity among the cytochromes  $b$ , particularly in prokaryotes, remains an open question.

Proteins of the cytochrome  $b_5$  family contain a cytochrome  $b_5$  domain (Fig. 2) and, in the case of the mitochondrial and microsomal members, a hydrophobic tail or anchor that associates the cytochrome with the relevant membrane (7). The structure of cytochrome  $b_5$  is shown in Fig. 2. In nitrate reductase, sulfite oxidase, and flavocytochromes  $b_2$ , the cytochrome  $b_5$  domain is tightly associated with a flavoprotein that specifies the particular enzymatic activity. The cytochromes  $b_5$  domain acts as an electron acceptor from the molybdenum (Mo)-pterin cofactor in sulfite oxidase or from protein-bound flavin in cytochrome  $b_5$  reductase, nitrate reductase, P450 reductase, and flavocytochrome  $b_2$ . The cytochrome  $b_5$  electron acceptor is cytochrome  $c$  in the flavocytochrome  $b_2$  and sulfite oxidase systems, and a Mo-pterin cofactor in nitrate reductase. In microsomal systems, cytochrome P450 and a variety of acyl-CoA desaturases are the electron acceptors. In red blood cells, cytochrome  $b_5$  reduces methemoglobin (see [Hemoglobin](#)). Thus, cytochrome  $b_5$  has evolved the ability to function in a diverse group of electron transfer pathways and with a variety of electron donors and acceptors (7).

### 3. C-type Cytochromes

These cytochromes are discussed in more detail elsewhere in this encyclopedia, but it is important to note that this is a diverse group of cytochromes consisting of a variety of structural motifs and functions (8). Although restricted to mitochondrial cytochrome  $c$  and cytochrome  $c_1$  (discussed previously) in eukaryotes and key components in the respiratory electron transfer chain, a variety of bacterial  $c$ -type cytochromes have evolved into a diverse and complex array of structures and function. In terms of structural families, the cytochromes  $c$  include the mitochondrial cytochrome  $c$  family and their bacterial homologies, cytochrome  $c_1$  and  $f$ , and the cytochromes  $c_3$ , cytochromes  $c'$ , flavocytochromes  $c$ , cytochrome  $c_4$ , and cytochromes  $c_5$ . Now, only cytochrome  $c_3$  and  $c$ -type cytochromes that are related to mitochondrial cytochrome  $c$  are well understood in both structure and function.

The basic structural motif of the mitochondrial cytochrome  $c$  family has been maintained throughout evolution. Figure 2 presents the structure of tuna cytochrome  $c$ . In contrast to cytochrome  $b_5$ , the heme propionates are buried and not solvent accessible. Cytochromes utilizing the cytochrome  $c$  structural motif are found in almost all bacteria but differ primarily in their redox potentials, which vary from approximately +50 to +450 mV, and in the type and distribution of amino acid side chains on their molecular surfaces.

### 4. Summary

Cytochromes evolved early in the development of life on this planet, and they have retained the basic functional group protoheme IX (with relatively minor variations) throughout evolution. A number of basic structural motifs have evolved that are found in both prokaryotes and eukaryotes, including the mitochondrial cytochromes  $c$  and their bacterial precursors, the cytochrome  $b_5$  family, the cytochrome oxidases, the membrane-spanning cytochromes  $b$ , and the cytochrome  $c_1/b_6$  family. It is striking that a relatively small number of structural families have evolved to take advantage of diverse metabolic opportunities, capturing redox energy to produce ATP and adapting to available ecological niches.

### Bibliography

1. M. Saraste, J. Castresana, D. Higgins, M. Lübben, and M. Wilmanns (1996) In *Origin and Evolution of Biological Energy Conversion* (H. Baltscheffsky, ed.) VCH, New York, pp. 255–289.
2. G. Tollin, T. E. Meyer, and M. A. Cusanovich (1986) *Biochim. Biophys. Acta* **853**, 29–41.

3. T. E. Meyer, G. Tollin, and M. A. Cusanovich (1994) *Biochimie* **76**, 480–488.
4. J. A. Watkins, T. E. Meyer, G. Tollin, and M. A. Cusanovich (1994) *Protein Sci.* **3**, 2104–2114.
5. P. N. Furbacker, G.-S. Tae, and W. A. Cramer (1996) In *Origin and Evolution of Biological Energy Conversion* (H. Baltscheffsky, ed.), VCH, New York, pp. 221–253.
6. B. A. C. Ackrell, M. K. Johnson, R. P. Gunsalus, and G. Cecchini (1992) In *Chemistry and Biochemistry of Flavoenzymes*, Vol. **III** (F. Müller, ed.), CRC Press, Boca Raton, pp. 229–297.
7. F. Lederer (1994) *Biochimie* **76**, 674–692.
8. T. E. Meyer, J. J. van Beeumen, R. P. Ambler, and M. A. Cusanovich (1996) In *Origin and Evolution of Biological Energy Conversion* (H. Baltscheffsky, ed.), VCH, New York, pp. 71–108.

### Suggestions for Further Reading

9. T. E. Meyer and M. A. Cusanovich (1989) Structure, function and distribution of soluble bacterial redox proteins, *Biochim. Biophys. Acta* **975**, 1–28.
10. R. H. Scott and A. G. Mauk, eds. (1996) *Cytochrome c. A. Multidisciplinary Approach*, University Science Books, Mill Valley.
11. G. R. Moore and G. W. Pettigrew (1990) *Cytochromes c*, Springer-Verlag, New York.

## Cytogenetics

Cytogenetics is the study of **genetics** by visualization of [chromosomes](#) and chromosomal aberrations, usually through light [microscopy](#). There is a distinguished history of chromosomal visualization for analytical purposes. Longley defined the morphology of **maize** chromosomes (reviewed in Ref. [1](#)), which opened the door for subsequent work on maize cytogenetics by McClintock ([2](#)). Staining of **prophase** and early **metaphase** maize chromosomes with acetocarmine allowed McClintock to identify 10 distinct linkage groups. In turn this laid the foundation for other maize cytogeneticists to identify chromosomal [translocations](#) and **inversions** visually. For example, it was possible to demonstrate that there is a correspondence between marker genes and cytologically defined chromosomal domains in crosses between heterozygous strains ([3](#)). Cytogenetics in *Drosophila melanogaster* received an enormous stimulus from the systematic study of the [polytene chromosomes](#) in salivary glands ([4](#)). These chromosomes have visibly distinct banding patterns under a phase microscope and after staining. Each of the four *Drosophila melanogaster* chromosomes have easily differentiated band patterns. The thousands of bands served as markers that allowed cytogeneticists to discern small deletions, duplications, translocations, and transpositions. Cytological experiments in *Drosophila* also allowed the discovery of the [position effect](#), the phenomenon in which the position of a gene relative to heterochromatin influences the efficiency with which it is expressed ([5](#)).

Improvements in cytogenetics eventually allowed visualizing metaphase chromosomes in somatic cells in culture. Only in 1956 was it established by these methodologies that human cells contain 46 chromosomes ([6](#)). Likewise cytogenetic approaches allowed the discovery that children with Down's syndrome are **trisomic** for chromosome 21 (see [Autosome](#)). Medical cytogenetics has since been advanced by chromosome stains, such as Giemsa, that enable separating individual chromosomes in sufficient numbers for highly focused analysis of their biochemistry and genetics (see [G Banding](#)).

## Bibliography

1. A. E. Longley (1952) *Bot. Rev.* **18**, 399–412.
2. B. McClintock (1929) *Science* **69**, 629–632.
3. H. B. Creighton and B. McClintock (1931) *Proc. Natl. Acad. Sci. USA* **17**, 485–491.
4. E. Heitz and H. Bauer (1933) *Z. Zellforsch Mikrosoc. Anat.* **17**, 67–81.
5. A. H. Sturtevant (1925) *Genetics* **10**, 117–147.
6. J. H. Tijo and A. Lean (1956) *Hereditas* **42**, 1–6.

## Cytokeratins

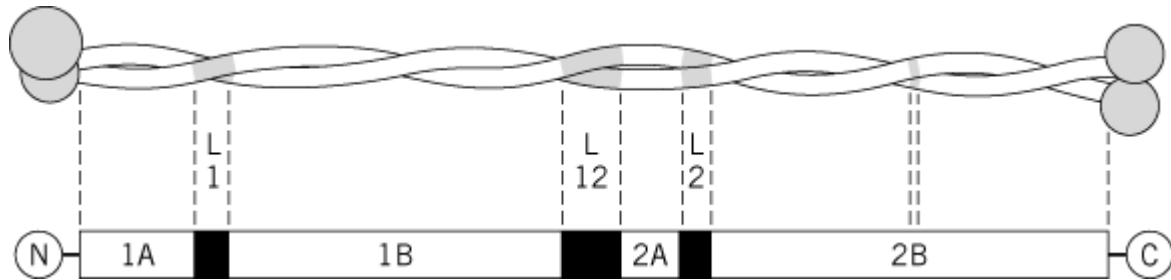
Cytokeratins form long filaments (they are one of the family of [intermediate filament](#) proteins; see also [Keratins](#)) that extend in a sinuous manner between the cytoplasmic surface of the cell [nucleus](#) and the plasma membrane of epithelial cells. The cytokeratin filaments appear to attach to the cytoplasmic side of the [nuclear envelope](#), often close to the [nuclear pore complexes](#). From there they radiate outward and link to or loop through the portions of the **desmosomes** or adhesion plaques that lie on the inner side of the cell membrane.

Cytokeratin chains can be divided into two groups: “acidic” Type I and “neutral-basic” Type II. In general, the Type II chains have higher [isoelectric points](#) and higher molecular weights than their Type I counterparts. Each chain comprises nonhelical end **domains** separated by a flexible rodlike structure. The *N*- and *C*-terminal domains in the epidermal keratins (or cytokeratins) have a characteristic substructure arranged about the rod domain with bilateral symmetry: E1, V1, H1 in the *N*-terminal domain and H2, V2, E2 in the *C*-terminal domain (1). The H1 and H2 subdomains, located respectively immediately *N*- and *C*-terminal to the rod domain, have sequences that are highly **homologous** regions within each chain type. The V1 and V2 subdomains are variable in length and sequence but always have a high content of [glycine](#) and [serine](#) residues. The E1 and E2 subdomains lie at the extreme ends of the molecule and are short, generally basic regions.

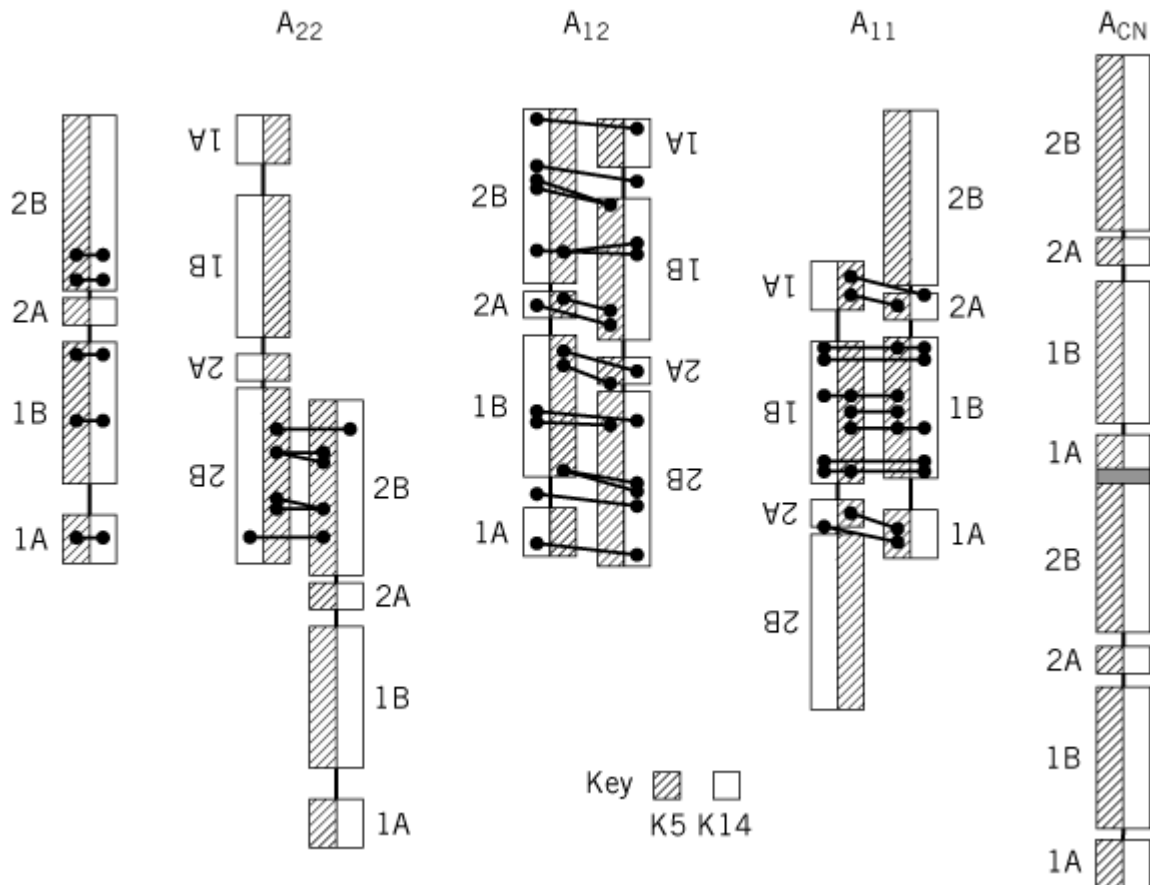
Between the terminal domains, the amino acid sequence has a [heptad repeat](#) substructure characteristic of a [coiled-coil](#) conformation (2). The latter thus consists essentially of coiled-coil rod domains (segments 1A, 1B, 2A, and 2B), joined by short linking domains (L1, L12, and L2), and results in a rod of total length 46.2 nm (Fig. 1). Each molecule contains a pair of different chains (Type I and Type II) aligned parallel to one another and in axial register (3–5). Segments 1B and 2 (= 2 A + L2 + 2B) each contain highly regular distributions of both acidic and basic residues (2). Thus each rod segment can be modeled as alternating bands of positive and negative charge, with a period of about 1.4 to 1.5 nm. Pairs of molecules aggregate through their rod domains in three antiparallel modes: (i) with their 1B segments largely overlapped, (ii) with their 2B segments largely overlapped, and (iii) with their entire lengths largely overlapped. In cytokeratins there is also a short head-to-tail overlap of about eight residues between parallel molecules (6). The antiparallel interactions originate in large part by the maximization of [electrostatic interactions](#), which are possible as a consequence of the regular charge distribution in neighboring rod domains. They are also both specified and stabilized by interactions between the H1 (and presumably H2) subdomains, the head-to-tail overlap region, and the link L2. Cross-links formed between [lysine](#) residues in cytokeratins have been characterized (Fig. 2), thus confirming the modes of aggregation originally proposed on theoretical grounds (see earlier). This work also allowed the projected (axial) lengths of the links L1, L12, and L2 to be determined, together with the relative axial staggers between various coiled-coil rod segments (6, 7). These data allow a surface lattice structure to be constructed that predicts a near-axial period in the cytokeratin of 22.6 nm (the experimentally observed value is

22.7 nm). **Scanning transmission electron microscope (STEM)** data indicate that *in vitro* assemblies may contain about 16, 24, 32, or 40 chains, but it seems that only the 32-chain variant is present *in vivo*. The STEM data, nonetheless, would indicate that the cytokeratin may have a substructure in which an eight-chain entity, a protofibril, is present. *In vitro* two, three, four, or five of these protofibrils may aggregate in a nonspecific manner.

**Figure 1.** Schematic cytokeratin intermediate filament molecule. The rod domain is of approximate length 46 nm, with coiled-coil segments 1A, 1B, 2A, and 2B and non-coiled-coil links L1, L12, and L2. A break in the phasing of the heptad repeat, known as a stutter, occurs at the center of segment 2B. The *N*- and *C*-terminal domains contain E (end), V (variable but rich in glycine and serine residues), and H (homologous) subdomains arranged about the rod domain with bilateral symmetry. (From Ref. 8, with permission.)



**Figure 2.** The cross-links induced chemically between lysine residues are shown for human K5/K14 epidermal keratin IF. The cross-links between the two polypeptide chains within a molecule (*left*) show that the chains are parallel and in axial register. Cross-links between molecules indicate that the molecules aggregate antiparallel through one of three modes, with (i) their 2B segments largely overlapped ( $A_{22}$ ), (ii) the entire molecules largely overlapped ( $A_{12}$ ), and/or (iii) their 1B segments largely overlapped ( $A_{11}$ ). Combination of  $A_{22}$  and  $A_{11}$  results in a small head-to-tail overlap of about 1 nm between parallel molecules. (From Ref. 7, with permission.)



Characterization of a large number of gene mutations leading to keratinopathies has been achieved in recent years. The resulting sequence changes lie very largely in or close to the head-to-tail overlap region comprising the conserved sequences at the *N*-terminal end of segment 1A and the *C*-terminal end of segment 2B. A reasonable conclusion is that this interaction is subtly modified by the mutated residue and becomes less able to stabilize the cytokeratin *in vivo*.

## Bibliography

1. P. M. Steinert, D. A. D. Parry, W. W. Idler, L. D. Johnson, A. C. Steven, and D. R. Roop (1985) Amino acid sequences of mouse and human epidermal type II keratins of  $M_r$  67000 provide a systematic basis for the structural and functional diversity of the end domains of keratin intermediate filament subunits. *J. Biol. Chem.* **260**, 7142–7149.
2. W. G. Crewther, L. M. Dowling, P. M. Steinert, and D. A. D. Parry (1983) Structure of intermediate filaments. *Int. J. Biol. Macromol.* **5**, 267–274.
3. P. M. Steinert (1990) The two-chain coiled-coil molecule of native epidermal keratin intermediate filaments is a type I–type II heterodimer. *J. Biol. Chem.* **265**, 8766–8774.
4. M. Hatzfeld and K. Weber (1990) The coiled-coil of *in vitro* assembled keratin filaments is a heterodimer of type I and type II keratin: use of site-specific mutagenesis and recombinant protein expression. *J. Cell Biol.* **110**, 1199–1210.
5. P. A. Coulombe and E. Fuchs (1990) Elucidating the early stages of keratin filament assembly. *J. Cell Biol.* **111**, 153–169.
6. P. M. Steinert, L. N. Marekov, R. D. B. Fraser, and D. A. D. Parry (1993) Keratin intermediate filament structure: crosslinking studies yield quantitative information on molecular dimensions and mechanism of assembly. *J. Mol. Biol.* **230**, 436–452.
7. P. M. Steinert, L. N. Marekov, and D. A. D. Parry (1993) Conservation of the structure of



keratin intermediate filaments: molecular mechanism by which different keratin molecules integrate into pre-existing keratin intermediate filaments during differentiation. *Biochemistry* **32**, 10046–10056.

8. D. A. D. Parry and R. D. B. Fraser (1985) Intermediate filament structure: 1, Analysis of IF protein sequence data. *Int. J. Biol. Macromol.* **7**, 203–213.

### Suggestions for Further Reading

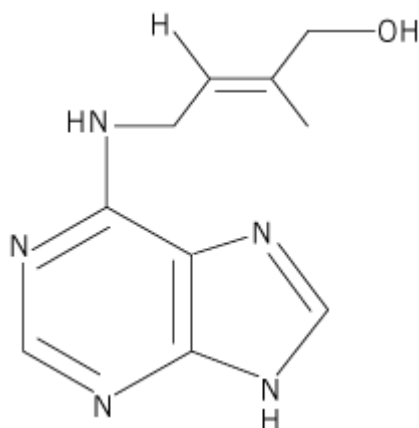
9. D. A. D. Parry and P. M. Steinert (1995) *Intermediate Filament Structure*, Springer-Verlag, Heidelberg.
10. R. D. Goldman and P. M. Steinert (1990) *Cellular and Molecular Biology of Intermediate Filaments* (eds.), Plenum Press, New York.
11. E. Fuchs and K. Weber (1994) Intermediate filaments: structure, dynamics, function and diseases. *Annu. Rev. Biochem.* **63**, 345–382.

## Cytokinins

### 1. History

The first discovery that plants make compounds that stimulate cell division dates back to 1913, when Haberlandt (1) demonstrated that an unknown substance present in vascular tissues of various plants caused cork cambium formation and wound healing in cut potato tubers. These molecules were later called *cytokinins*, from the Greek “cytokinesis” (cell division). In the 1940s, van Overbeek et al. (2) showed that similar compounds were also present in the milky endosperm of immature coconuts. In 1954, Miller et al. (3) isolated a substance, formed by partial breakdown of herring sperm DNA, that, in synergy with *auxin*, stimulated cell division in tobacco pith cells very actively. The molecule was called *kinetin* (6-furfurylamino purine). *Trans*-zeatin (6-(4-hydroxy-3-methylbut-*trans*-2-enylamino) purine) was the first cytokinin to be identified from maize endosperm (4) (Fig. 1).

**Figure 1.** Structural formula of *trans*-zeatin.



### 2. Biosynthesis and Metabolism

Natural cytokinins (CKs) are primarily purine derivatives with an  $N^6$ -substituted isopentenylated side chain (iP-type CKs) (5). However, a series of CKs with an aromatic side chain (eg, benzyladenine) have been demonstrated in certain species (6, 7). The most abundant naturally occurring CKs are CK nucleosides and CK nucleotides (5). Certain phenolic compounds also possess cell-division-promoting activity, but are not classified as true cytokinins (8).

The origin of CKs in plants remains a subject of debate. Roots are the primary sites of cytokinin biosynthesis, although limited synthesis in the shoot has also been demonstrated (9, 10). CKs have been proposed to be transported from root to shoot by the transpiration stream in xylem tissues, but conclusive evidence is lacking (11). Even though present throughout the plant, CKs have never been proven to be synthesized by endogenous, plant-borne [enzymes](#). Because many root-associated bacteria are cytokinin producers, it has been hypothesized that CKs are made exclusively by these microbial symbionts (12).

Cytokinins can be synthesized via a direct (*de novo*) or an indirect ([transfer RNA](#)) pathway. The main source of free CKs is the *de novo* pathway. The key step is the formation of  $N^6$ -(D<sup>2</sup>-isopentenyl) adenosine-5'-monophosphate from D<sup>2</sup>-isopentenyl pyrophosphate and AMP catalyzed by isopentenyltransferase. The isopentenyl side chain originates from mevalonate, which also serves as a precursor to [abscisic acid](#), [gibberellins](#), and [brassinosteroids](#). Isopentenyltransferase was first characterized from bacteria and slime molds (5, 7). Different constructs using the *Agrobacterium* isopentenyltransferase *IPT* (*tmr*) **gene** were used to generate constitutive or controlled cytokinin overproduction in **transgenic** plants (13-15). These studies mainly confirmed the classical view of cytokinin action in plants. Specific manipulation of [senescence](#) was achieved by placing the *IPT* gene under the control of a senescence-specific **promoter** (16). Recently, isopentenyltransferase activity was also demonstrated in plants (17), but the enzyme was never purified to homogeneity, nor have plant genes encoding isopentenyltransferase been **cloned**. This is the major reason why CK production by the plant itself is still questioned. An *Arabidopsis* mutant with sixfold elevated CK levels was isolated from a screen for seedlings with multiple cotyledons (18, 19). The **phenotype** of this *amp1* (*altered meristem program 1*) mutant is reminiscent of several CK effects, including enhanced lateral branching, elevated rate of leaf formation, induction of flowering and floral morphological abnormalities, and delayed senescence. In addition, *amp1* displays a de-etiolated phenotype in the dark and is an **allele** of *cop2* (*constitutive photomorphogenic 2*). It was proposed that CK induces light-regulated genes in a light-independent fashion (19). The *AMP1* gene product may be required for CK degradation, or it could be a negative regulator of CK biosynthesis. A second mutant *cri* (*crystal*) that accumulates sixfold more CK than the wild type was recently identified (20). *Cri1* mutants are phenotypically different from *amp1* and show symptoms of vitrification.

The second biosynthetic route for CKs is via tRNA. Modified bases in tRNA, including isopentenylated bases, have been found in virtually all organisms investigated (21). CK modification in tRNA is always at the adenosine residue immediately 3' to the [anticodon](#) of tRNAs that read codons starting with U (5). These CK residues in tRNA may serve as a direct, albeit not major, precursor of free CK, as a result of tRNA turnover. The rate of CK release in this way is dependent on the metabolic activity of the cell. The major CK in tRNA is *cis*-zeatin, with very low CK activity, but it can be converted to the biologically active *trans* isomer (Fig. 1) by a *cis-trans* isomerase (22) (see [Cis/Trans Isomerization](#)).

Several modifications can occur to the  $N^6$ -substituted side chain and purine ring. Formation of zeatin and its nucleoside and nucleotide derivatives, most probably by hydroxylation of the isopentenyl adenine moiety, is common in most plant tissues. The hydroxylase that mediates this conversion was isolated from cauliflower (23). [O-Glycosylation](#) of hydroxylated CKs is often observed, but leads to a considerable decrease in activity (5). CK-*O*-glucosides can be converted to free active CK by the action of  $\beta$ -glucosidases. In contrast, modification of the purine ring by *N*-glucosylation of the 3, 7,

or 9 position yields inactive, stable derivatives that cannot generate free CKs. *N*-Glucosylation is thus expected to play an important role in regulating the levels of active CKs. An additional control is exerted by CK oxidase (24, 25). This enzyme catalyzes degradation of the *N*<sup>6</sup>-(D<sup>2</sup>-isopentenyl) side chain and irreversibly leads to loss of CK activity. CK nucleotides are not accepted as a substrate. A general role for CK oxidase is to contribute to CK homeostasis in plants (26).

### 3. Signal Perception and Transduction

A number of CK-binding proteins have been described, one of which is a potential CK **receptor** (27, 28). The latter is a cytosolic protein that reversibly binds zeatin and is involved in transcriptional activation in the presence of zeatin, but not in its absence. Only the *trans* isomer elicited the observed effect.

Two candidate plasma membrane CK receptors were identified by mutational analysis (29, 30). Mutants of *Arabidopsis* that exhibit shoot formation in tissue culture in the absence of CK were obtained by an activation-tagging approach. The *CKII* (*cytokinin-independent 1*) gene was isolated and encodes a protein similar to two-component regulators of bacteria (29). As in the case of the ethylene receptor ETR1 (see [Ethylene](#)), CKII is a hybrid [kinase](#) with fused sensor and receiver domains. Based on this analogy and on the fact that constitutive overexpression of *CKII* results in typical effects of CK action, CKII is hypothesized to encode a CK receptor. Overexpression would enable cells to sense endogenous concentrations below the threshold, which normally leads to CK response. Intriguingly, phytochromes also display similarities with two-component systems. It has been speculated that they might regulate other two-component systems, thus enabling cross-talk between the ethylene and CK pathways (31). A second candidate CK receptor is a seven-transmembrane domain receptor (7TM), characteristic in **G-protein** signaling and involved in transduction of extracellular signals (30). The *Arabidopsis GCR1* (*G-coupled receptor 1*) gene was isolated by polymerase chain reaction (PCR), and **antisense** mutagenesis with a 35S *aGCR1* construct resulted in reduced cotyledon and leaf expansion and a single flowering stem, a phenotype reminiscent of the *cyr1* (*cytokinin resistant 1*) mutant (see text below). However, *GCR1* and *CYR1* do not represent the same gene (30).

The above findings confirm the involvement of both phosphorylation and G proteins in CK signaling (32). Sense and antisense expression in tobacco of a gene encoding a small GTP-binding protein from rice resulted in a dwarf phenotype and abnormal flower development, correlating with a sixfold elevation of CK levels compared to wild type (33). **Phosphorylation** plays an important role in [cell cycle](#) control by CK (32). In addition, evidence was given for involvement of intracellular Ca<sup>+</sup> in CK signaling (see [Calcium Signaling](#)) mostly from studies in mosses (32). Finally, the similarity of a CK-binding protein to *S*-adenosyl-homocysteine hydrolase raised the possibility that some CK effects might be mediated through control of **methylation** of DNA and/or proteins (34).

Further progress in understanding CK [Signal Transduction](#) is expected from mutant studies. Several CK-resistant mutants have been identified, mainly in *Arabidopsis* (35). Five loci have been studied in detail. The majority were isolated using a root elongation response in the presence of CK. In certain cases, resistance to other hormones was found. For instance, the *ckr1* (*cytokinin resistant 1*) mutant is allelic to *ein2* (*ethylene-insensitive 2*), indicating cross-talk with the ethylene pathway. It was demonstrated that the inhibition of root elongation by CK results from an induction of ethylene biosynthesis (36). This was exploited by screening for the absence of triple response in the dark, in the presence of kinetin (37). The *cin5* (*cytokinin-insensitive 5*) mutant affects the 1-aminocyclopropane-1-carboxylate synthase gene 5 in *Arabidopsis*. Disruption of the carboxyl terminus of the same isoform leads to ethylene overproduction, previously identified as an *ethylene overproducer* (*eto2*) mutant. A third class of CK-resistant mutants is *cyr1* (38). In contrast to *ckr1* and *cin5*, *cyr1* displays an abnormal shoot phenotype: cotyledons and leaves fail to expand, a limited number of leaves are formed, and a single infertile flower is made. No resistance to other hormones was observed; however, *cyr1* is hypersensitive to abscisic acid (38). *Cyr1* is probably allelic to *emf2*

(*embryonic flower 2*). *EMF* gene products might function in the maintenance of vegetative growth (39). The *stp1* (*stunted plant 1*) mutant has a reduced sensitivity to CKs, both in assays for elongation and radial swelling of the root (40). The general morphology was not affected, but a reduction in growth rate was observed. The function of the *STP1* gene product remains to be elucidated. Finally, a recent study describes identification of three loci involved in control of cell division and plant development (41). The corresponding mutants were called *pasticcino* (*pas*) and showed hypertrophy of their apical parts when grown on CK-containing medium, a phenotype reminiscent of abnormal shoots regenerated on media with unbalanced auxin/CK ratio and strikingly similar to the fasciation disease caused by *Rhodococcus fascians* (41). The altered embryo, root, and leaf development result from uncoordinated cell divisions. *PAS* genes are hypothesized to play a role in CK signaling, because no differences in CK levels were found compared to wild type.

#### 4. Downstream Targets

Primary CK response genes have not yet been isolated, nor have CK responsive [cis-acting](#) elements or [trans-acting](#) factors been identified (42). A large number of genes that are either induced or repressed upon CK treatment have been described. Controls are exerted at a transcriptional or post-transcriptional level. Changes in gene expression were observed between 2 h and 48 h after treatment. Many of these secondary CK response genes are involved in light regulation and nutrition (32, 42).

#### 5. Effects

Depending on temporal and spatial factors, a particular subset of downstream target genes is activated by CK treatment and results in one of many described effects. Cytokinins regulate many aspects of plant growth and development, including cell division, enlargement, and differentiation, as well as [chloroplast](#) development, release of apical dominance, nutrient mobilization, flowering, and delay of senescence (43, 44).

#### Bibliography

1. G. Haberlandt (1913) *Sitzungsber. K. Preuss. Akad. Wissensch.*, 318–345.
2. J. van Overbeek, M. E. Conklin, and A. F. Blakeslee (1941) *Science* **94**, 350–351.
3. C. O. Miller, F. Skoog, M. H. Von Saltza, and F. Strong (1955) *J. Am. Chem. Soc.* **77**, 1392–1393.
4. D. S. Letham (1963) *Life Sci.* **8**, 569–573.
5. C.-m. Chen (1997) *Physiol. Plant.* **101**, 665–673.
6. M. Strnad, W. Peters, E. Beck, and M. Kamínek (1992) *Plant Physiol.* **99**, 74–80.
7. E. Prinsen, M. Kamínek, and H. Van Onckelen (1997) *Plant Growth Regul.* **23**, 3–15.
8. R. A. Teutonico, M. W. Dudley, J. D. Orr, D. G. Lynn, and A. N. Binns (1991) *Plant Physiol.* **97**, 288–297.
9. A. Carmi, and J. Van Staden (1983) *Plant Physiol.* **73**, 76–78.
10. C.-M. Chen, J. R. Ertl, S. M. Leisner, and C.-C. Chang (1985) *Plant Physiol.* **78**, 510–513.
11. P. D. Hare, W. A. Cress, and J. van Staden (1997) *Plant Growth Regul.* **23**, 79–103.
12. M. A. Holland (1997) *Plant Physiol.* **115**, 865–868.
13. J. I. Medford, R. Horgan, Z. El-Sawi, and H. J. Klee (1989) *Plant Cell* **1**, 403–413.
14. T. Schmülling, S. Beinsberger, J. De Greef, J. Schell, H. A. Van Onckelen, and A. Spena (1989) *FEBS Lett.* **249**, 401–406.
15. J. J. Estruch, E. Prinsen, H. Van Onckelen, J. Schell, and A. Spena (1991) *Science* **254**, 1364–1367.
16. S. Gan and R. M. Amasino (1997) *Plant Physiol.* **113**, 313–319.
17. J. R. Blackwell and R. Horgan (1994) *Phytochemistry* **35**, 339–342.

18. A. M. Chaudhury, S. Letham, S. Craig, and E. S. Dennis (1993) *Plant J.* **4**, 907–916.
19. A. N. Chin-Atkins, S. Craig, C. H. Hocart, E. S. Dennis, and A. M. Chaudhury (1996) *Planta* **198**, 549–556.
20. V. Santoni, M. Delarue, M. Caboche, and C. Bellini (1997) *Planta* **202**, 62–69.
21. N. Murai (1981) In *Cytokinins: Chemistry, Activity and Function* (D. W. S. Mok and M. C. Mok, eds.), CRC Press, Boca Raton, FL, pp. 87–99.
22. N. V. Bassil, D. W. S. Mok, and M. C. Mok (1993) *Plant Physiol.* **102**, 867–872.
23. C.-m. Chen and S. M. Leisner (1984) *Plant Physiol.* **75**, 442–446.
24. D. J. Armstrong (1994) In *Cytokinins: Chemistry, Activity, and Function* (D. W. S. Mok and M. C. Mok, eds.), CRC Press, Boca Raton, FL, pp. 139–154.
25. R. J. Jones and B. M. N. Schreiber (1997) *Plant Growth Regul.* **23**, 123–134.
26. S. Eklöf, C. Astot, T. Moritz, J. Blackwell, O. Olsson, and G. Sandberg (1996) *Physiol. Plant.* **98**, 333–344.
27. C. Brinegar (1994) In *Cytokinins: Chemistry, Activity, and Function* (D. W. S. Mok and M. C. Mok, eds.), CRC Press, Boca Raton, FL, pp. 217–232.
28. O. N. Kulaeva, N. N. Karavaiko, S. Y. Selivankina, Y. V. Zemlyachenko, and S. V. Shipilova (1995) *FEBS Lett.* **366**, 26–28.
29. T. Kakimoto (1996) *Science* **274**, 982–985.
30. S. Plakidou-Dymock, D. Dymock, and R. Hooley (1998) *Curr. Biol.* **8**, 315–324.
31. J. W. Reed (1998) *Trends Plant Sci.* **3**, 43–44.
32. P. D. Hare, and J. van Staden (1997) *Plant Growth Regul.* **23**, 41–78.
33. H. Sano, S. Deo, E. Orudjev, S. Youssefian, K. Isizuka, and Y. Ohashi (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10556–10560.
34. S. Mitsui, T. Wakasugi, and M. Sugiura (1996) *Plant Growth Regul.* **18**, 39–43.
35. J. Deikman (1997) *Plant Growth Regul.* **23**, 33–40.
36. A. J. Cary, W. Liu, and S. H. Howell (1995) *Plant Physiol.* **107**, 1075–1082.
37. J. P. Vogel, K. E. Woeste, A. Theologis, and J. J. Kieber (1998) *Proc. Natl. Acad. Sci. USA* **95**, 4766–4771.
38. J. Deikman and M. Ulrich (1995) *Planta* **195**, 440–449.
39. C.-H. Yang, L.-J. Chen, and Z. R. Sung (1995) *Dev. Biol.* **169**, 421–435.
40. T. I. Baskin, A. Cork, R. E. Williamson, and J. R. Gorst (1995) *Plant Physiol.* **107**, 233–242.
41. J.-D. Faure, P. Vittorioso, V. Santoni, V. Fraisier, E. Prinsen, I. Barlier, H. Van Onckelen, M. Caboche, and C. Bellini (1998) *Development* **125**, 909–918.
42. T. Schmülling, S. Schäfer, and G. Romanov (1997) *Physiol. Plant.* **100**, 505–519.
43. A. N. Binns (1994) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **45**, 173–196.
44. P. J. Davies (1995) *Plant Hormones: Physiology, Biochemistry and Molecular Biology*, Kluwer, Dordrecht, The Netherlands.

## Cytomegalovirus

Cytomegaloviruses (CMVs) are a distinct subgroup of [herpesvirus](#) and are classified as subfamily *Betaherpesvirinae* in the *Herpesviridae* family. CMV forms a characteristic nuclear inclusion, which

gives an “owl's eye” appearance to infected cells. Many animal species have their own species-specific CMV.

Human CMV maintains a collinear genomic organization, but no significant nucleotide sequence [homology](#) with other animal CMVs. Human CMV is also designated as human herpes virus 5 (HHV-5). A great majority of the human population are infected with CMV, without any symptoms by adulthood. CMV is excreted in the urine of infants and intermittently in saliva of healthy seropositive adults. Human CMV may cause diseases such as congenital cytomegalic inclusion disease and infectious mononucleosis. The importance of human CMV as a pathogen has been increasingly correlated with the rising number of immunodeficient patients due to immunosuppressive therapy after transplantation, acquired immune deficiency syndrome (AIDS), and so on.

Mature virions of CMV are 150 to 200 nm in diameter and have an envelope composed of a lipid bilayer and viral glycoproteins. Inside the envelope, a tegment or matrix surrounds an icosahedral capsid, which contains a double-stranded, linear DNA [genome](#). Two types of noninfectious particles, in addition to infectious virions, are produced in cells infected with human CMV. Dense bodies do not contain nucleocapsid or viral DNA and are essentially enveloped tegment protein. Noninfectious enveloped particles have a capsid but no electron-dense DNA core.

Equine CMV has a genome of 180 kbp, while the genomes of human, simian, and murine CMVs are 230 to 240 kbp. Human CMV is unique among animal CMVs in that it has a genome structure similar to herpes simplex virus 1 (HSV-1, class E structure). Two unique sequences (UL and US) are flanked by [inverted repeats](#). Inversion of the UL and US segments mediated by inverted repeats can give rise to four different genomes. Other animal CMVs are less complicated in genome structure, and genome inversion is not known for them.

The genome of the AD169 strain of human CMV has been sequenced. Two hundred and eight open reading frames (ORFs) with a coding capacity of at least 100 amino acid residues were identified in the genome. Although only a few of these ORFs have been actually shown to encode proteins, at least 40 have been shown to be dispensable for replication. ORFs have been designated by their location with regard to the unique and repeated sequences (TRL, UL, IRL, IRS, and US) and sequential numbers. By analogy to HSV-1, roughly one-quarter of these ORFs are predicted to encode proteins required for [DNA replication](#) and metabolism functions, and three-quarters for virion maturation and virion structure.

Human CMV shows a highly restricted host range and takes 48 to 72 h to yield detectable amounts of progeny virions. Primary human differentiated cells, such as skin or lung fibroblasts, as well as chimpanzee cells, support human CMV replication, but undifferentiated, transformed, or **aneuploid** cells are usually nonpermissive for CMV replication. The genes of CMV have been classified into three sequential, kinetic classes: a(immediate early), b(delayed early) and g(late) genes. Expression of the latter two relies on the a gene products. The b gene products mainly serve DNA replication and metabolism, while g gene products encode structural proteins.

#### Suggestion for Further Reading

E. S. Mocarski Jr. (1996) "Cytomegaloviruses and Their Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2447–2492.

## Cytoplasmic Inheritance

Most of the genetic information of a eukaryotic cell resides in its [nucleus](#), and some is found in the cytoplasm. Each **mitochondrion** and [chloroplast](#) contains several copies of specific DNA molecules that carry some of the genetic information for the organelle; the rest is encoded in nuclear **genes**. Organelle genes are **transcribed** and **translated** in the organelles themselves through a [genetic code](#) that diverges in many cases from the standard one, while the [proteins](#) encoded in nuclear genes are imported into the organelle after their **protein biosynthesis** on cytoplasmic [ribosomes](#). As a consequence of this shared genetic determination, there are many cases of nuclear and organelle mutants that exhibit similar phenotypes.

Most organelle [genomes](#) consist of circular double-stranded DNA molecules (Table 1). After considerable controversy, it seems clear that many mitochondrial DNA molecules are linear rather than circular.

**Table 1. Some Genomes from Mitochondria and Chloroplasts, with Indication of Their Size, Number of Confirmed or Putative Translatable Genes (Open Reading Frames, ORFs), and number of genes for transfer RNA and ribosomal RNA<sup>a</sup>**

| Species                          | Number         | Size (kb) | ORF | tRNA | rRNA |
|----------------------------------|----------------|-----------|-----|------|------|
| <b>Mitochondrial Genomes</b>     |                |           |     |      |      |
| <i>Saccharomyces cerevisiae</i>  | L36885         | 78        | 13  | 24   | 4    |
| <i>Schizosaccharomyces pombe</i> | X54421         | 19        | 10  | 25   | 2    |
| <i>Drosophila melanogaster</i>   | U37541         | 20        | 13  | 22   | 2    |
| <i>Caenorhabditis elegans</i>    | X54252         | 14        | 12  | 22   | 2    |
| <i>Homo sapiens</i>              | X93334         | 17        | 13  | 22   | 2    |
| <i>Chlamydomonas reinhardtii</i> | U03843         | 16        | 8   | 3    | 14   |
| <i>Arabidopsis thaliana</i>      | Y08501, Y08502 | 367       | 81  | 19   | 3    |
| <b>Chloroplast Genomes</b>       |                |           |     |      |      |
| <i>Euglena gracilis</i>          | X70810         | 143       | 66  | 41   | 11   |
| <i>Nicotiana tabacum</i>         | Z00044         | 156       | 107 | 37   | 8    |
| <i>Oryza sativa</i>              | X15901         | 135       | 113 | 46   | 8    |

<sup>a</sup> Data taken from the indicated EMBL sequence accession numbers.

Organelles may contain molecules of DNA other than their own [chromosomes](#). Examples are the two linear plasmids of the mitochondria of *Zea mays* that result in the absence of functional pollen (androsterility); they are very useful for the commercial production of hybrid seed (see [Heterosis](#)). Another example is *kalilo*, a linear plasmid of *Neurospora crassa* that blocks mitochondrial functions when it integrates into the mitochondrial DNA; the result is [senescence](#), the inability for indefinite vegetative growth.

Other cytoplasmic bodies with their own genetic information range from various self-replicating particles to **viruses** and complete endosymbiotic cells. Their genomes may consist of DNA or RNA. First to be discovered were the *kappa* particles of *Paramecium aurelia* (1), but many others are known, such as the *killer* particles in *Saccharomyces cerevisiae* and the *s* particles that render

*Drosophila melanogaster* specially sensitive to carbon dioxide.

## 1. Heteroplasmons and Mosaicism

A cell contains numerous copies of organelle DNA molecules, up to several thousand, and these need not be identical (heteroplasmon). Organelle genomes are notoriously unstable and suffer frequent point mutations and large rearrangements. Heteroplasmons can be produced by these spontaneous changes or by cell fusion. Contrary to the usual behavior of **heterozygotes**, the relative proportion of two variants in a heteroplasmon is not fixed, giving rise to mosaicism in cell clones.

In some cases, one of the variants imposes itself because of its faster reproduction. Thus, most *petite* mutants of the yeast *Saccharomyces cerevisiae*, lacking aerobic metabolism, carry changes in their mitochondrial DNA. Some of these are called *suppressible*, because mutant mitochondria displace wild-type mitochondria in heteroplasmons. Wild-type mitochondria, on the other hand, impose their phenotype when mixed with mitochondria that carry *petite* mutations called *neutral*.

In other cases, the two variants may be found in various relative proportions, or even pure, in different cells, as would be expected from random [genetic drift](#). Many variegated ornamental plants exhibit patches or streaks of different colors in their stalks and leaves, ranging from normal green to whitish. This phenotype may have many causes; in the plant *Mirabilis jalapa*, the colors depend on the presence of mutant chloroplasts, mixed with the normal ones in various proportions.

[Recombination](#) can occur between the various forms of chloroplast DNA or mitochondrial DNA in heteroplasmons, as shown for example for chloroplast markers of *Chlamydomonas reinhardtii* and mitochondrial markers in *Saccharomyces cerevisiae*.

## 2. Heterokaryon Test

The fusion of two cells that differ in both nuclear and cytoplasmic markers results in cells that are at the same time [heterokaryons](#) and heteroplasmons. Vegetative multiplication of the heterokaryon produces three kinds of cells: heterokaryons and two homokaryons, one for each of the kinds of nuclei. Cytoplasmic traits may segregate in various ways, all independent of the segregation of the nuclei. Homokaryotic segregants may show new combinations of cytoplasmic and nuclear markers. This test requires that the nuclei in the heterokaryon do not fuse, but remain separate, so that they can segregate to different daughter cells. The test has been extended to the yeast *Saccharomyces* by the use of dominant *kar* mutations that block nuclear fusion after mating, thus giving rise to heterokaryons.

## 3. Reciprocal Crosses

The results of reciprocal crosses offer a test for cytoplasmic inheritance when the two parents do not contribute equally to the cytoplasm of the zygote. The most frequent case in animals and plants is *matrilinear inheritance*, in which all individuals inherit from their mothers only, because of the larger size and contribution of female gametes. Examples are the variegation in *Mirabilis*, the first report of this heredity pattern (2), androsterility in *Zea*, and several human diseases due to mitochondrial mutations, such as Leber's hereditary optic neuropathy.

The matrilinear rule is not absolute. The s particles of *Drosophila* are transmitted most often by the ova ([egg](#)), but sometimes by [sperm](#). Chloroplasts may be provided by the male, as in the conifers, or by both parents, as in *Oenothera*. Mammals have powerful mechanisms to exclude the mitochondrial DNA of the sperm from the embryo, and the failure of these mechanisms might produce exceptions to the matrilinear rule.

Uniparental inheritance is observed with the *poky* mitochondrial mutants and the *kalilo* plasmid in the fungus *Neurospora crassa*. When two haploid strains of different **mating type**, *A* or *a*, are



crossed, either can provide the large protoperithecia or the small conidia that fuse to produce the zygote. All offspring inherit the cytoplasmic markers from the protoperithecia donor, irrespective of its mating type.

Uniparental inheritance can occur even when the gametes are similar in size, as in *Chlamydomonas reinhardtii*. The zygote is formed by the fusion of seeming identical cells of the two mating types, called  $mt^+$  and  $mt^-$ . Almost all the offspring inherit the chloroplast markers of the  $mt^+$  parent only, because the DNA of the  $mt^-$  parent is usually destroyed (3).

There are cases of maternal inheritance that are not due to cytoplasmic genes.

### Bibliography

1. T. M. Sonneborn (1938) *Science* **88**, 503.
2. C. Correns (1909) *Z. Indukt. Abst. Vererbungsl.* **1**, 291–329.
3. R. Sanger (1954) *Proc. Natl. Acad. Sci. USA* **40**, 356–363.

### Suggestions for Further Reading

4. G. J. Hermann and J. M. Shaw (1998) Mitochondrial dynamics in yeast. *Annu. Rev. Cell Dev. Biol.* **14**, 265–303.
5. R. N. Lightowlers, P. F. Chinnery, D. M. Turnbull, and N. Howell (1997) Mammalian mitochondrial genetics: heredity, heteroplasmy and disease. *Trends Genet.* **13**, 450–455.
6. N. G. Larsson and D. A. Clayton (1995) Molecular genetic aspects of human mitochondrial disorders. *Annu. Rev. Genet.* **29**, 151–178.
7. J. Nosek, L. Tomaska, H. Fukuhara, Y. Suyama, and L. Kovac (1998) Linear mitochondrial genomes: 30 years down the line. *Trends Genet.* **14**, 184–188.
8. M. Sugiura (1995) The chloroplast genome. *Essays Biochem.* **30**, 49–57.

## Cytoskeleton

The cytoskeleton consists primarily of three types of filament: [microtubules](#) (about 25 nm in diameter), **actin**-containing [microfilaments](#) (initially thought to be 5 to 7 nm in diameter but now known to be about 10 nm in diameter), and [intermediate filaments](#) (about 10 nm in diameter). The latter were so named because they were intermediate in size between the long-established and well-characterized microtubules and microfilaments. Microtubules are filamentous structures often many microns in length that occur in the cytoplasm of eukaryotic cells. The actins are a highly conserved family of eukaryotic proteins, and actin filaments are perhaps best known for the highly specific interactions they make with **myosin** thick filaments that results in **muscle** contraction. Many actin-binding proteins play important roles in establishing the cytoskeleton *in vivo* and in allowing it to function appropriately. For example, fimbrin and **a-actinin** aggregate and cross-link microfilaments into bundles, severin severs actin filaments, b-actinin acts as an actin filament-capping protein, and ponticulin (a **transmembrane protein**) binds actin filaments and [spectrin](#) binds and links the microfilaments to the cell membrane. See [Actin](#), [Microfilament](#), [Microtubules](#), [Intermediate Filaments](#), [Cytokeratins](#), and [Keratins](#) for further information on various components of the cytoskeleton.

## Cytotoxic T Lymphocytes

The main interest to molecular biology of cytotoxic T lymphocytes is probably mechanistic. These cells are able to destroy other cells. How do they achieve this? The 40-odd year-old history of this field started with immunology and defense, and the cellular phenomenology that was possible at the time. It now culminates with the realization that the long-sought mechanisms of cytotoxicity convey in fact molecular signals of [cell death](#). Thus, fundamental research on cytotoxic T cells now goes beyond immunology, becoming part of the booming research effort on [programmed cell death](#).

It was shown in 1960 that white peripheral blood cells from alloimmune dogs were able to lyse *in vitro* target cells bearing the corresponding **alloantigens** (1). We now know that many different cell types have this ability to destroy other cells, in particular lymphoid cells, and among these T (i.e., *thymus-derived*) lymphocytes. There are two varieties: CD4+ and CD8+ lymphocytes. When activated, both varieties can be cytotoxic through a Fas-based mechanism, and in addition CD8+ cells can be cytotoxic via a perforin-granzyme-based mechanism. These two mechanisms will be considered in more detail below. A main role of cytotoxic T cells is to lyse syngeneic cells presenting surface determinants originating from **viruses** or other intracellular microorganisms. These, when within cells, are out of reach of other immune system effectors. Destruction of infected cells destroys the microorganism or exposes it to other effector mechanisms. Cytotoxic T cells thus fight intracellular pathogens. Interestingly, cytotoxic T cells can also destroy normal, uninfected cells, such as Fas-bearing activated lymphocytes, thus in part ensuring down-regulation of immune responses. Neither the perforin- nor the Fas-based mechanisms in isolation, nor both together, seem to account fully for graft rejection.

When an effector cell approaches a target cell, it “recognizes” it through an array of diverse molecules at the surface of the effector cell (TCR/CD3/peptide, CD4 or CD8, LFA-1, etc.) and of the target cell (Class I or Class II MHC, ICAM-1, etc.). How, then, does the cytotoxic cell destroy the target cell? First indications were obtained at a phenomenological level. Engagement of the above-mentioned molecules leads to effector cell activation, including an increase in the concentration of free cytoplasmic calcium (2) within a matter of minutes, immediately followed by reorientation within the effector cell of the secretory apparatus towards an area facing the target cell (3). This suggested the possible participation of secretory phenomena in at least one mechanism of cytotoxicity. Many cytotoxic T cells indeed contain peculiar cytoplasmic granules. Their content can be released by degranulation in the extracellular fluid on target cell recognition (see [Protein Secretion](#)).

Among the phenomena then observed in the target cell is the early disintegration of nuclear DNA (4, 5) into fragments of around 180 bp, and multiples thereof, corresponding to the size of **nucleosome**-shielded DNA, which is due to the action of one or several **endonucleases** activated early in the sequence of events leading to target cell death. Other events in this sequence include an influx of calcium, condensation of the cytoplasm, fragmentation of the cytoplasm and nucleus, and only very late secondary membrane disruption. This phenomenology of cell death is very similar, if not identical, for example in developmental circumstances. In all cases, the same program seems to govern the course of events in a cell that dies, including, as we now know, a similar cascade of molecular events and the same “apoptotic” morphological traits (6). What may be different is merely the signaling of this death program. The nature of the signals produced by cytotoxic cells has been largely unraveled by studies at the molecular level.

### 1. The Perforin/Granzyme B Mechanism of Cytotoxicity

A first approach to identify molecules involved in cytotoxicity has been to look for effector cell

molecules themselves endowed with cytotoxic activity. This approach led to the detection (7, 8), characterization, and [cloning](#) of *perforin*, a 60-kDa protein present in granules in the cytoplasm of many cytotoxic T and NK (natural killer) cells. At high concentrations, it can act on membranes as a calcium-dependent channel-former; most interestingly, it shows significant homology to the terminal components of the [complement system](#) cascade. A formal demonstration that perforin was involved in a mechanism of cytotoxicity was provided by several groups (9), all showing that cytotoxicity mediated by T cells was greatly impaired in mice made perforin-deficient through [gene targeting](#).

Other molecules present in granules of cytotoxic T cells were identified by subtractive cloning. Prominent among molecules isolated this way were a number of [serine proteinases](#), such as CTLA-1/CCP1/Granzyme B (10-12) and H Factor/CTLA-3/Granzyme A (10, 12, 13). Granzyme B was demonstrated through gene targeting to be required in a mechanism of cytotoxicity (14).

A current view on the cooperation between perforin and Granzyme B within the same mechanism of cytotoxicity is as follows. On recognition of the target cell, perforin- and granzyme B-containing granules in the effector cell reorient towards the target cell and are triggered to **exocytose**. Granzyme B may somehow enter the target cell and migrate into a target cell compartment where it would be innocuous, but from which it can be released by perforin at sublytic concentrations (15-17). Intracellularly released Granzyme B would then activate the standard programmed cell death cascade, through direct cleavage and therefore activation (18) of some of the ICE-family cysteine proteases called [caspases](#). However, while target cell death mediated by Granzyme B (or at least by purified granules) may involve caspase activation for phenomena such as DNA fragmentation, it may, interestingly, use a noncaspase pathway for cell membrane disruption (19).

## 2. The Fas-Based Mechanism of Cytotoxicity

Through a **somatic cell** genetic approach it was found (20) that another mechanism of T cell-mediated cytotoxicity required a protein called Fas/APO-1/CD95 (21, 22) at the target cell surface. This in turn led to the cloning of the Fas ligand that is expressed at the effector cell surface (23). Fas belongs to a family of **receptors** (including Fas, TNF-R1, DR3, DR4, DR5), most or all of which are very efficient at transducing a signal interpreted as a cell death signal, while Fas ligand belongs to a family of corresponding ligands (including the Fas ligand, the TNFs, TRAIL, etc.). Involvement of Fas in cytotoxicity may be a reflection *in vitro* of the main roles of the Fas system *in vivo*, such as down-regulation of the [immune response](#) (24), protection against the immune system (25, 26), or potential major physiopathological effects (27).

Target cell Fas and effector cell Fas ligand define molecularly the Fas-based mechanism of cytotoxicity. In this mechanism, when an effector cell encounters a target cell, engagement of the [T cell receptor](#) of the effector cell by target cell [major histocompatibility complex](#) (MHC) leads to expression of the Fas ligand at the effector cell surface. Fas-ligand expression can be induced on CD8+ and on Th0 and Th1 CD4+ [T cells](#) (28-30) following specific recognition of [antigen](#) through the T cell receptor. This induction seems to involve the activation of **tyrosine kinases**, requires RNA and protein synthesis and the presence of calcium, and is inhibited by [cyclosporin A](#) (31-34). Other molecules may be involved, such as Myc and Max, *9-cis-retinoic acid*, and its receptors, the **orphan receptor** for steroids, Nur77, and ALG-3.

Once expressed, effector cell Fas ligand engages target cell Fas, which leads to target cell death. In more detail, engagement of Fas leads in a matter of seconds to protein recruitment via the Fas “death domain,” the cytoplasmic segment of Fas that is necessary and sufficient to transduce a death signal (35, 36). First, FADD/MORT-1 (37, 38) directly binds the death domain of Fas via its own C-terminal death domain. FADD/MORT-1 also includes an N-terminal death effector domain through which it associates with another molecule, FLICE/MACH (39, 40). FLICE/MACH includes two death effector domains at its N terminus and, strikingly, a caspase-homology domain at its C terminus (39, 40). Thus, the Fas-FADD-FLICE complex provides a remarkably direct link from a membrane signal to caspase activation, which is required for programmed cell death.

### 3. Concluding Remarks

The perforin-based and the Fas-based pathways account for most, and perhaps all, of T cell-mediated cytotoxicity (41-44), at least as assessed in a 4-hr assay *in vitro*. Other molecules, such as TNF-R1 or the TRAIL receptors, may play a role in longer assays. The Fas pathway is used by both CD4<sup>+</sup> (mostly Th1) and CD8<sup>+</sup> effector cells, while the latter also use the perforin/granzyme pathway. Thus, in CD8<sup>+</sup> T cells, two signals may stem from the T cell receptor/CD3 complex on antigen-specific recognition, one of them leading to granule exocytosis and the other one leading to [transcription](#) of the Fas ligand gene. Neither the exact nature of these distinct signals (45), nor whether they can be triggered simultaneously in a given cytotoxic cell is known as yet.

More generally, both mechanisms of T cell-mediated cytotoxicity seem to act by signaling the caspase activation step of the evolutionarily conserved programmed cell death cascade within the target cell. In this sense, the emergence in evolution of T cell-mediated cytotoxicity has not required the invention of new mechanisms of killing, but merely of new ways of signaling a preexisting programmed cell death cascade.

### Bibliography

1. A. Govaerts (1960) *J. Immunol.* **85**, 516–522.
2. M. Poenie, R. Y. Tsien, and A.-M. Schmitt-Verhulst (1987) *EMBO J.* **6**, 2223–2232.
3. A. Kupfer, G. Dennert, and S. J. Singer (1985) *J. Mol. Cell. Immunol.* **2**, 37–49.
4. J. H. Russell, V. Masakovski, T. Rucinsky, and G. Phillips (1982) *J. Immunol.* **128**, 2087–2094.
5. J. J. Cohen et al. (1985) *Adv. Exp. Med. Biol.* **184**, 493–508.
6. J. F. R. Kerr, A. H. Wyllie, and A. R. Currie (1972) *Br. J. Cancer* **26**, 239–257.
7. P. A. Henkart (1985) *Annu. Rev. Immunol.* **3**, 31–58.
8. E. R. Podack (1985) *Immunol. Today* **6**, 21–27.
9. D. Kägi et al. (1994) *Nature* **369**, 31–37.
10. J.-F. Brunet et al. (1986) *Nature* **322**, 268–271.
11. C. G. Lobe et al. (1986) *Science* **232**, 858–861.
12. D. Masson and J. Tschopp (1987) *Cell* **49**, 679–685.
13. H. K. Gershenfeld and I. L. Weissman (1986) *Science* **232**, 854–858.
14. J. W. Heusel et al. (1994) *Cell* **76**, 977–987.
15. C. J. Froelich et al. (1996) *J. Biol. Chem.* **271**, 29073–29079.
16. D. A. Jans et al. (1996) *J. Biol. Chem.* **271**, 30781–30789.
17. L. Shi et al. (1997) *J Exp. Med.* **185**, 855–866.
18. A. J. Darmon, D. W. Nicholson, and R. C. Bleackley (1995) *Nature* **377**, 446–448.
19. A. Sarin et al. (1997) *Immunity* **6**, 209–215.
20. E. Rouvier, M.-F. Luciani, and P. Golstein (1993) *J. exp. Med.* **177**, 195–200.
21. S. Yonehara, A. Ishii, and M. Yonehara (1989) *J. Exp. Med.* **169**, 1747–1756.
22. B. C. Trauth et al. (1989) *Science* **245**, 301–305.
23. T. Suda, T. Takahashi, P. Golstein, and S. Nagata (1993) *Cell* **75**, 1169–1178.
24. J. H. Russell, B. Rush, C. Weaver, and R. Wang (1993) *Proc. Nat. Acad. Sci. USA* **90**, 4409–4413.
25. D. Bellgrau et al. (1995) *Nature* **377**, 630–632.
26. T. S. Griffith et al. (1995) *Science* **270**, 1189–1192.
27. L. E. French and J. Tschopp (1997) *Nat. Med.* **3**, 387–388.
28. T. Suda et al. (1995) *J. Immunol.* **154**, 3806–3813.

29. F. Ramsdell et al. (1994) *Internat. Immunol.* **6**, 1545–1553.
30. L. L. Carter and R. W. Dutton (1995) *J. Immunol.* **155**, 1028–1031.
31. A. Anel et al. (1994) *Eur. J. Immunol.* **24**, 2469–2476.
32. A. Anel et al. (1995) *Eur. J. Immunol.* **25**, 3381–3387.
33. F. Vignaux et al. (1995) *J. Exp. Med.* **181**, 781–786.
34. M.-F. Luciani and P. Golstein (1994) *Roy. Soc. Phil. Trans. B* **345**, 303–309.
35. N. Itoh and S. Nagata (1993) *J. Biol. Chem.* **268**, 10932–10937.
36. M. P. Boldin et al. (1995) *J. Biol. Chem.* **270**, 387–391.
37. M. P. Boldin et al. (1995) *J. Biol. Chem.* **270**, 7795–7798.
38. A. M. Chinnaiyan, K. O'Rourke, M. Tewari, and V. M. Dixit (1995) *Cell* **81**, 505–512.
39. M. P. Boldin, T. M. Goncharov, Y. V. Goltsev, and D. Wallach (1996) *Cell* **85**, 803–815.
40. M. Muzio et al. (1996) *Cell* **85**, 817–827.
41. D. Kägi et al. (1994) *Science* **265**, 528–530.
42. B. Lowin, M. Hahne, C. Mattmann, and J. Tschopp (1994) *Nature* **370**, 650–652.
43. H. Kojima et al. (1994) *Immunity* **1**, 357–364.
44. C. M. Walsh et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10854–10858.
45. M. T. Esser, B. Krishnamurthy, and V. L. Braciale (1996) *J. Exp. Med.* **183**, 1697–1706.

### **Suggestions for Further Reading**

46. G. Berke (1994). The binding and lysis of target cells by cytotoxic lymphocytes: Molecular and cellular aspects. *Annu. Rev. Immunol.* **12**, 735–773.
47. P. A. Henkart (1994). Lymphocyte-mediated cytotoxicity: two pathways and multiple effector molecules. *Immunity* **1**, 343–346.
48. S. Nagata and P. Golstein (1995). The Fas death factor. *Science* **267**, 1449–1456.
49. M. E. Peter, F. C. Kischkel, S. Hellbardt, A. M. Chinnaiyan, P. H. Krammer, and V. M. Dixit (1996). CD95(APO-1/Fas)-associating signalling proteins. *Cell Death Differ.* **3**, 161–170.

### **D Gene Segment**

The existence of D genes encoding a short region of the third complementarity determining region (CDR3) of [immunoglobulin](#) heavy chains was first postulated when it was apparent that expressed H chains between the sequences encoded by the **variable** (V) and **J (joining)** genes contained a segment variable in length that had no counterpart in either gene. It turned out that variability of this region resulted from 2 causes: One was the existence of a set of D genes (D for Diversity); the second was due to the presence of nucleotides that were not germline encoded but added enzymatically at both the VD and the DJ junctions by terminal-deoxynucleotidyl-transferase, which is specifically expressed in lymphocytes. As the result of these combined events, heavy-chain CDR3 appeared as the most hypervariable region of Ig molecules, and therefore of prime importance as a basis for antibody specificity.

D gene segments appear to be present in all heavy-chain gene loci, from the lowest vertebrate species. In humans, about 30 D gene segments have been identified, clustered between the IgVH locus and the JH genes. The various D genes differ in length, but are always limited to a small

number of potential **codons**. They are flanked on both their 5' and 3' ends by the conventional recombination signal sequences recognized by the RAG1 and RAG2 [recombinases](#). The spacers are of different length on each side, so they meet with the 12/23 alternate size rule, to recombine with VH and JH gene segments. D genes make a major contribution to heavy-chain variability of CDR3, not only because of their number and sequence diversity, but also because they are used in many alternate possibilities: (a) They may be inserted in either orientation; (b) any reading frame may be used (in humans, not in the mouse; see text below), although one is more frequently encountered; (c) they may join to each other, so that examples of as many as four fused D genes have been described, thus providing a CDR3 of unusual length; and (d) they may be partially deleted at both ends before being inserted between the V and J genes. All of this clearly can generate a huge diversity potential.

In the mouse, minor differences from the above situation in D gene use and function have been reported, due to the fact that there is essentially only one reading frame used by murine D genes. There are two reasons for this. One [gene rearrangement](#) introduces a [stop codon](#) in the 3' JC sequence. The second reason is more complicated. Some mouse D genes have a **promoter**-like region that can be activated after the first gene rearrangement making the DJ joint has been performed. In that case, a D-J-C<sub>m</sub> protein is synthesized, which becomes exposed at the surface of the preB cell and blocks any further gene rearrangement (see [B Cell](#)), but is unable to contribute a functional immunoglobulin.

The genes for [T-cell receptors](#) (TCRs) have also D genes on the b and d chains. There are only a number of D genes, which are used in either orientation and reading frame. They have 12-bp spacers of identical length on both 5' and 3' flanking recombination signal sequences, corresponding to the 23-bp spacers of the V and J genes. TCR D genes also contribute to CDR3, a region that is most variable, as for immunoglobulins.

See also entries [Gene Rearrangement](#) and [Recombinase](#).

#### Suggestions for Further Reading

F. Matsuda and T. Honjo (1994) Organization of the human immunoglobulin heavy-chain locus. *Adv. Immunol.* **62**, 1–29.

H. Sakano, Y. Kurosawa, M. Weigert, and S. Tonegawa (1982) Identification and nucleotide sequence of a diversity DNA segment (D) of immunoglobulin heavy chain genes. *Nature* **290**, 562–565.

H. Gu, D. Kitamura, and K. Rajewsky (1991) B-cell development regulated by gene rearrangement—arrest of maturation of membrane bound D<sub>m</sub> protein and selection of DH element reading frames. *Cell*, **65**, 47–54.

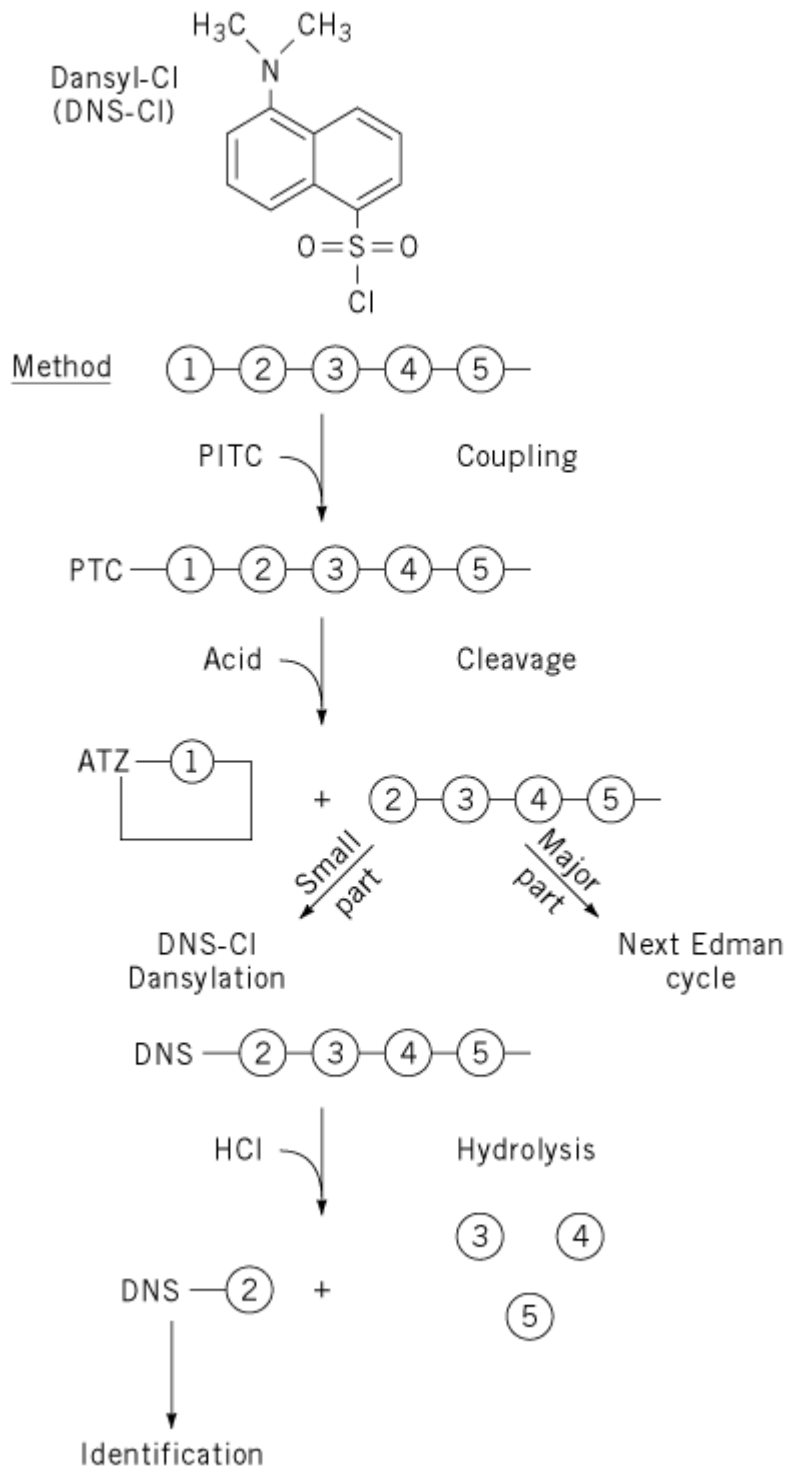
## Dansyl Chloride

The repetitive nature of the [Edman Degradation](#) reaction opened the way to [protein sequencing](#) and to automation in the form of protein sequencers, but manual Edman degradations also became popular during much of the 1960s and 1970s.

One major manual method used dansyl (1-dimethylaminonaphthalene-5-sulphonyl) chloride (Fig. [1](#), see top of next page) to detect the new N-terminus after each Edman cycle ([1](#)). Combined with purification of peptides on paper by high-voltage [electrophoresis](#) and [chromatography](#) (see [Peptide](#)

[Mapping](#)), it made protein sequence analysis accessible to many laboratories without the need for large equipment. It also avoided the difficulties in identification of the phenylthiohydantoin (PTH)-amino acids that existed early, and replaced PTH-identification with dansyl-amino acid identification. The latter became rapid and reliable by **thin-layer chromatography** (2). Hence, many of the original protein [primary structures](#) determined used this methodology, originating from the MRC Laboratory of Molecular Biology, Cambridge, England. Together with the Sanger method of DNA sequencing later developed there, the “dansyl-Edman/paper purification” technique made this laboratory the leader in sequencing technology for a long time in the 1960s and 1970s.

**Figure 1.** Principle of the “dansyl-Edman” modification of the Edman degradation reaction, with N-terminal analysis of separate samples. PITC, phenylisothiocyanate; ATZ, anilinothiazolinone; PTH, phenylthiohydantoin; Dansyl, 1-dimethylaminonaphthalene-5-sulphonyl.



Dansyl chloride reacts with protein N-termini much like the original Sanger reagent, fluorodinitrobenzene (FDNB) (3), but it is easier to use, more sensitive, and fluorescent, making it easily detectable. Dansyl chloride is still used for protein labeling in many cases where fluorescence is desirable.

#### Bibliography

1. W. R. Gray (1967) *Methods Enzymol.* **11**, 469–475 and (1972) **25**, 121–138.
2. K. R. Woods and K.-T. Wang (1967) *Biochim. Biophys. Acta* **133**, 369–370.
3. F. Sanger (1945) *Biochem. J.* **39**, 507–515.



## Databases

Molecular biology is an empirical discipline that requires observing and understanding different types of data. For example, a **gene** may imply a physical location on the [chromosome](#), a **nucleotide sequence**, an [amino acid](#) sequence, a three-dimensional [protein structure](#), a molecular component of cellular function, a regulatory mechanism of **gene expression**, or even a **phenotypic** difference caused by [mutation](#). In addition to this variability, the quantity of molecular biology data is increasing rapidly, especially for gene and protein sequences and 3-D structures, due to advances in experimental technologies. Molecular biology databases are a number of resources available over the Internet that comprise a [bioinformatics](#) infrastructure for biomedical sciences. Each database contains a specific type of data that has cross-references to other databases, which can be used for integrated information retrieval. This is possible because a database is generally organized as a collection of entries, and connections can be made at the level of entries without standardizing how data items should be organized within an entry. The advent of the World Wide Web (WWW) in the early 1990s was a boon to molecular biology databases because the concept of hyperlinks is fully compatible with the practice of cross-references. The WWW also dramatically increased the accessibility of computers to biologists. The primary resources of molecular biology databases are bibliographic databases, [sequence databases](#), and [structure databases](#).

Although bibliographic, sequence, and structural aspects of molecular biology are relatively easy to computerize, the next step is to organize their functional aspects. There are resources in that direction, such as motif libraries that contain higher level knowledge abstracted from sets of functionally related sequences and pathway databases that contain computerized knowledge of molecular interactions and biochemical pathways. Representative examples of these resources are shown in Table [1](#).

**Table 1. Selected List of Molecular Biology Databases**

| <b>Data Type</b>      | <b>Database</b> | <b>Organization</b>                                  |
|-----------------------|-----------------|--|
| Biomedical literature | Medline         | National Center for Biotechnology Information (NCBI) |
| Nucleotide sequence   | GenBank         |  |
|                       | EMBL            | European Bioinformatics Institute (EBI)              |
|                       | DDBJ            | National Institute of Genetics, Japan                |
| Amino acid sequence   | PIR             | National Biomedical Research Foundation              |
|                       | SWISS-PROT      | Swiss Institute of Bioinformatics (SIB)              |
| Sequence motif        | PROSITE         |  |
| 3-D structure         | PDB             | Brookhaven National Laboratory                       |

Abstraction of a real problem is made through a data model. For example, data are organized in two-dimensional tables in the relational data model. The relational database based on the relational model has been widely used in a number of applications, including some of the sequence databases. Although in principle all different types of molecular biology data can be stored in a single, unified, relational database, this is impossible in practice because of the varying views of how data items should be organized and related. In the current web of molecular biology databases, different types of data are integrated by a loose coupling based on links (cross-references), rather than a tight coupling based on unified schema. This approach is extended to include other types of links, especially similarity links computed by similarity search algorithms and biological links representing molecular interactions, which can also be integrated for biological reasoning (1). The major bioinformatics servers shown in Table 2 provide link-based database retrieval systems, such as Entrez at NCBI, SRS at EBI, and DBGET/LinkDB at GenomeNet in Kyoto.

**Table 2. WWW Addresses for the Major Bioinformatics Servers**

---

| Server    | Address  |
|-----------|--|
| NCBI      | <a href="http://www.ncbi.nlm.nih.gov">www.ncbi.nlm.nih.gov</a> |
| EBI       | <a href="http://www.ebi.ac.uk">www.ebi.ac.uk</a>               |
| ExPASy    | (site currently unavailable)                                   |
| GenomeNet | <a href="http://www.genome.ad.jp">www.genome.ad.jp</a>         |

---

In scientific disciplines, the merit of storing and managing information in a computer was first realized in bibliographic databases, which were designed for humans to read and understand. In the next step, factual data reported in the literature were computerized in factual databases, such as in sequence databases, which made more sophisticated retrieval available, for example, sequence similarity searches. Even in this case, however, the database is still a static resource to be retrieved, and it is up to humans to make sense out of the retrieved data. In contrast, the knowledge of links or relationships is more dynamic in nature. For example, ancestors in a family can be retrieved from a “deductive” database that contains parent–child relationships and the rules for combining them. Thus, knowledge is different from data or information, in that new knowledge can be generated dynamically from existing knowledge by logical reasoning. In the era of mass data production, molecular biology requires logical computation based on empirical knowledge rather than numerical computation based on first principles (see [Bioinformatics](#)). The web of molecular biology databases also requires a new generation of knowledge bases.

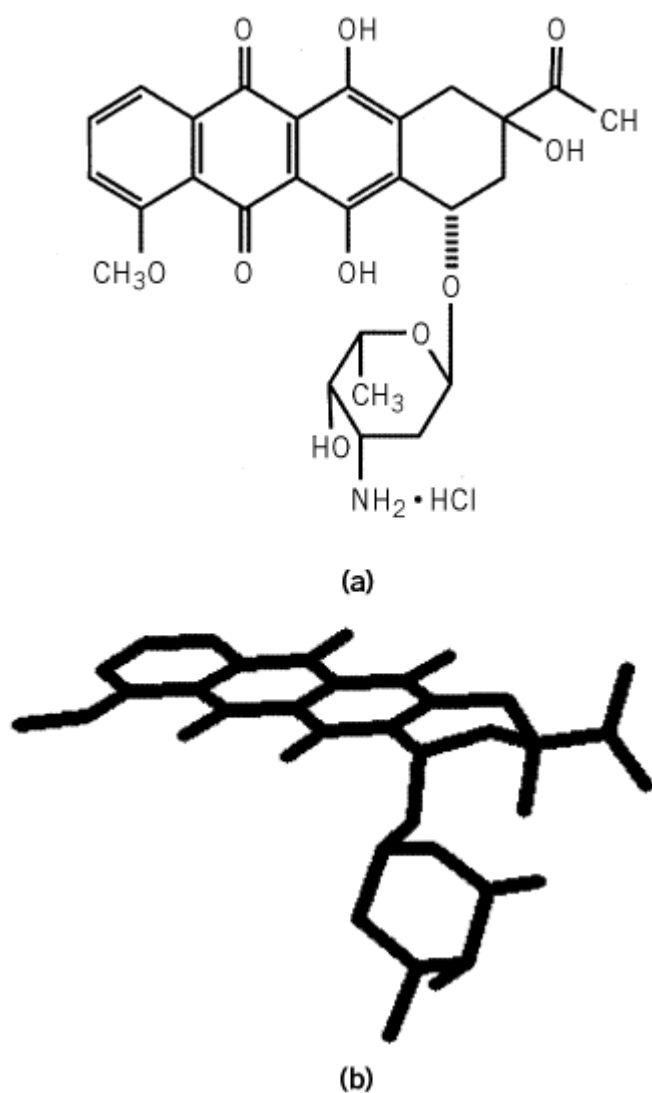
#### Bibliography

1. M. Kanehisa (1997) Trends Biochem. Sci. **22**, 442–444.

## Daunomycin

Daunomycin (synonym: daunorubicin) is the prototype of the clinically important anthracycline antibiotics that are widely used in cancer chemotherapy. Daunomycin and its close relative [adriamycin](#) are natural products isolated from various *Streptomyces* species. These two anthracyclines are probably the most widely used antitumor agents worldwide ([1](#), [2](#)). Daunomycin is primarily used in the treatment of acute leukemias. The structure of daunomycin is shown in [Figure 1](#); it is composed of two major parts: an anthraquinone ring system and the pendant daunosamine, an amino sugar.

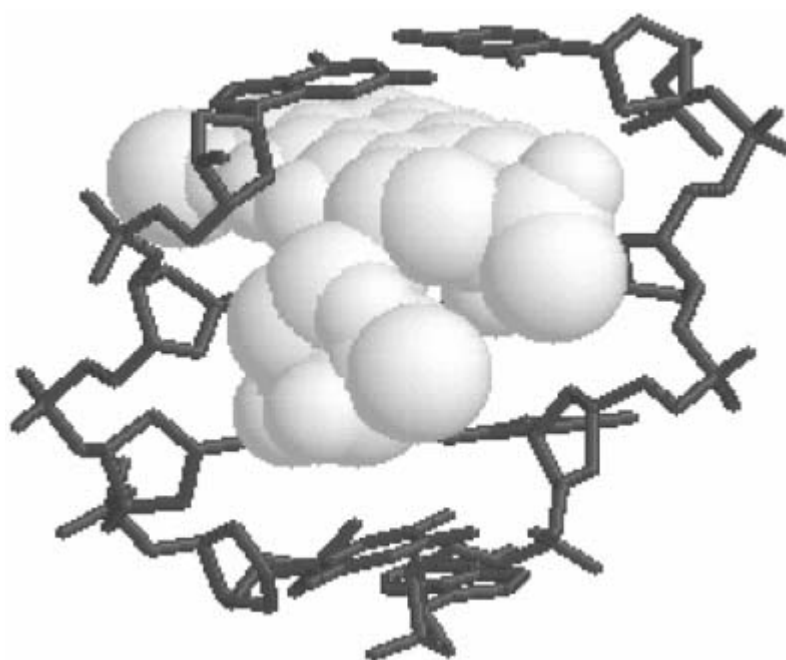
**Figure 1.** The chemical structure (a) and molecular model (b) of the anthracycline antibiotic daunomycin.



Daunomycin rapidly permeates into cells, passing through the [membrane](#) into the cytoplasm, and then accumulates in the nucleus, where it binds avidly to **DNA**. It binds to DNA by the process of

[intercalation](#), in which the anthraquinone ring inserts between adjacent DNA base pairs. The daunosamine moiety provides additional interactions that stabilize the complex by fitting into the minor groove of the DNA double helix, where it participates in [van der Waals interactions](#). These additional interactions in the minor groove distinguish daunomycin from simple intercalators like [ethidium bromide](#), which are stabilized only by the interactions within the intercalation site. Figure 2 shows the structure of a daunomycin–DNA complex determined by [X-ray crystallography](#) (3). Intercalation simultaneously lengthens and unwinds the DNA duplex, and it stabilizes DNA against thermal denaturation, increasing its **T<sub>m</sub> (melting temperature)**.

**Figure 2.** Structure of a daunomycin–DNA complex. The anthraquinone portion of daunomycin (shown in a space-filling representation) is intercalated between the CpG step in a DNA oligonucleotide (shown in a stick representation). The daunosamine portion of daunomycin lies in the DNA minor groove. This figure was prepared from coordinates (file DDF01) deposited in the Nucleic Acid Database (Rutgers University), <http://ndbserver.rutgers.edu>.



Daunomycin, despite its small size (527.5 MW), exhibits several levels of **molecular recognition** (4). It is structurally specific and strongly prefers to bind to right-handed, B-form DNA over alternate helical forms. Daunomycin will **allosterically** convert left-handed Z DNA to an intercalated, righted-handed form (5). There is pronounced sequence selectivity in the binding of daunomycin to DNA, which is revealed by DNase I [footprinting](#) experiments. Daunomycin binds preferentially to the triplet sequences 5'(A/T)GC and 5'(A/T)CG, where the notation (A/T) indicates that either A or T may occupy the sequence position (6, 7).

The mechanism by which daunomycin acts to kill cancer cells remains incompletely defined, and it may involve inhibition of several different cellular functions. Daunomycin, like most intercalators, is an effective inhibitor of [DNA replication](#) and of [transcription](#), both *in vitro* and *in vivo*. The most likely intracellular target for daunomycin, according to recent evidence, is the enzyme topoisomerase II (see [DNA Topology](#)). Daunomycin poisons this key enzyme, by trapping an intermediate enzymatic species in which cleaved DNA strands are covalently attached to topoisomerase II, by preventing the normal resealing of the duplex DNA after strand passage (8). A comprehensive study of several anthracycline derivatives, including daunomycin, showed that intercalation is necessary, but not sufficient, for inhibition of topoisomerase II (9).

The anthraquinone moiety of daunomycin readily undergoes a number of oxidation–reduction reactions, with the production of oxygen free radicals. These free radicals can damage cells, and they may represent another important mechanism in the cytotoxicity of daunomycin. A serious, complicating side effect in the clinical use of daunomycin in cancer chemotherapy is its cardiotoxicity (1, 2). Free-radical production has been implicated as the primary cause of such cardiotoxicity.

Daunomycin was one of the first effective cancer chemotherapeutic agents identified, and it remains clinically important after over 30 years of use. It is one of the best characterized and understood intercalators and serves as an important model for understanding how small molecules bind to DNA.

### Bibliography

1. R. B. Weiss (1992) *Semin. Oncol.* **19**, 670–686.
2. G. N. Hortobagyi (1997) *Drugs* **54**, 1–7.
3. A. H.-J. Wang (1992) *Curr. Opin. Struct. Biol.* **2**, 361–368.
4. J. B. Chaires (1996) in *Advances in DNA Sequence Specific Agents*, Vol. 2 (L. H. Hurley and J. B. Chaires, eds.), JAI Press, Greenwich, CT, pp. 141–167.
5. J. B. Chaires (1986) *J. Biol. Chem.* **261**, 8899–8907.
6. B. Pullman (1990) *Anti-cancer Drug Design* **7**, 96–105.
7. J. B. Chaires, J. E. Herrera, and M. J. Waring (1990) *Biochemistry* **29**, 6145–6153.
8. F. Zunnino and G. Capranico (1990) *Anti-cancer Drug Design* **5**, 307–317.
9. A. Bodley et al. (1989) *Cancer Res.* **49**, 5969–5978.

### Suggestions for Further Reading

10. F. Arcamone (1981) *Doxorubicin*, Academic Press, New York.
11. J. W. Lown, ed. (1988) *Anthraacycline and Anthracenedione Based Anticancer Agents*, Elsevier, Amsterdam.
12. W. Priebe, ed. (1995) *Anthracycline Antibiotics: New Analogues, Methods of Delivery, and Mechanism of Action*, ACS Symposium Series 574, American Chemical Society, Washington, DC.

## DEAD and DEAH Domains

DEAD and DEAH domains contain a characteristic Asp-Glu-Ala-Asp/His sequence that is usually abbreviated DEAD or DEAH, using the amino-acid one-letter code. The **domains** belong to a large [superfamily](#) of [enzymes](#) involved in nucleic acid metabolism. Both DEAD and DEAH domains alter nucleic acid secondary structure in a nucleoside triphosphate (NTP)-dependent manner. Many function as either [DNA helicases](#) or [RNA Helicases](#), which unwind DNA and RNA duplexes, respectively.

The general properties of DEAD/DEAH domains are summarized, and the functions of several well-characterized DEAD and DEAH domains are reviewed. Included are the DEAD domain of eukaryotic initiation factor 4a, eIF-4A, and the precursor RNA processing (PRP) DEAH domains found in the [spliceosome](#), a multiprotein complex responsible for processing pre-[messenger RNA](#) into mRNA (see [RNA Splicing](#)). Finally, the structural basis for sequence conservation will be

discussed using the **X-ray crystallographic** structures of the DExx DNA helicases from *Bacillus stearothermophilus* and *Escherichia coli* and the DExH RNA helicase from the hepatitis C virus (HCV).

## 1. Identification and Classification of DEAD and DEAH Domains

Helicases are enzymes that remove double-stranded regions of nucleic acid when single-stranded RNA and DNA are required. Regulated production of single-stranded nucleic acids occurs during many cellular processes, including replication, [recombination](#), and [transcription](#). The energy required for unwinding is derived from the hydrolysis of NTPs such as ATP.

Amino acid sequences can be used to identify and classify helicases. Helicase sequences contain seven regions or motifs of conserved residues (1). On the basis of the sequences of the individual motifs, helicases have been broadly grouped into two major [superfamilies](#), SFI and SFII, and several smaller families [(2), Fig. 1). With the increased availability of [genome](#) information, helicases are frequently identified from their sequences alone, prior to the demonstration of helicase activity *in vitro*.

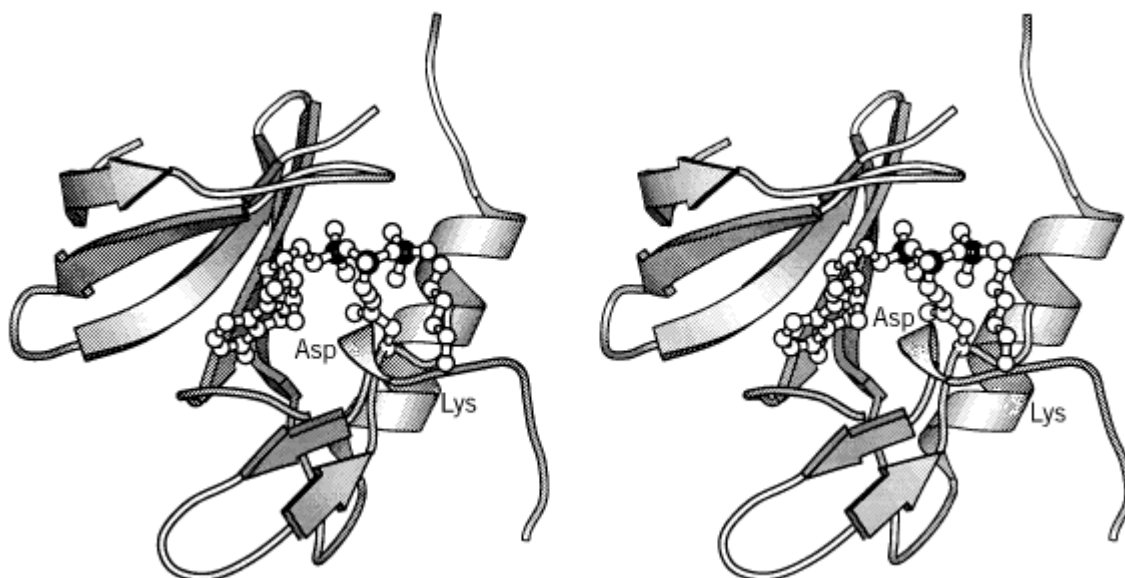
**Figure 1.** Representative examples of conserved sequence motifs. The seven common sequence motifs of DEAD and DEAH domains are shown in one-letter code. Helicases from both superfamilies I and II are given. The NTP-binding sequences within motifs I and II of the structures of the *E. coli* Rep DNA and HCV RNA helicases revealed similar spatial arrangements of all motifs except the following abbreviations: Rep, *E. coli* Rep DNA helicase; PcrA, *B. stearothermophilus* PcrA DNA helicase; eIF-4A, eIF-4A helicase; HCV, RNA helicase domain of HCV NS3 protein, and PRP22, *S. cerevisiae* precursor RNA processing protein

|         | Motif    |        |      |         |           |      |  |
|---------|----------|--------|------|---------|-----------|------|--|
|         | I        | Ia     | II   | III     | IV        | V    |  |
| Protein |          |        |      |         |           |      |  |
| SFI     |          |        |      |         |           |      |  |
| Rep     | AGAGSGKT | FTNKAA | DEYQ | VGDDDQS | IKLEQNYRS | QLMT |  |
| PCRA    | AGAGSGKT | FTNKAA | DEYQ | VGADQDS | ILLEQNYRS | MLMT |  |
| SF2     |          |        |      |         |           |      |  |
| eIf-4A  | AQSGTGKT | PTRELA | DEAD | LSATMP  |           | LITT |  |
| HCV     | APTGSQKS | PSVAAT | DECH | ATATPP  | LIFCHSKKK | VVST |  |
| PRP22   | GETGSQKT | PRRVAA | DEAH | VTSATL  |           | DGIY |  |

The energy for helicase-mediated unwinding of nucleic acid duplexes is derived from the hydrolysis of NTPs. Consequently, amino acid residues known to bind NTPs in other proteins are among the most conserved across the helicase families. By analogy with structures of ATP:protein complexes, the conserved lysine residue embedded in the Ala-x-x-Gly-x-Gly-Lys-Thr sequence of motif I probably interacts with the b and g phosphates of NTP (Fig. 2). Similarly, NTPases share an aspartic

acid residue in motif II that ligates the magnesium of the NTP:  $Mg^{+2}$  complex (Fig. 2). In the DEAD and DEAH domains, this conserved aspartic acid occurs within the title sequences Asp-Glu-Ala-Asp/His. In 1982, Walker et al. (3) identified the functional significance of these motifs through sequence comparisons of ATP-utilizing enzymes. For this reason, motifs I and II are also referred to as Walker motifs A and B. In the literature, DEAD domains are occasionally referred to as DEAD-box proteins.

**Figure 2.** Stereoscopic views of an enzyme:ATP complex. Interactions made by the conserved lysine and aspartic acid residues of the protein with the ATP phosphates and  $Mg^{2+}$ , respectively, are highlighted. The protein backbone is shown as ribbons and the bound ATP in the ball-and-stick representation. Phosphate atoms are shown as solid spheres. The  $Mg^{2+}$  ion is situated between the phosphate oxygen and aspartic acid. For clarity, the complete coordination of the  $Mg^{2+}$  is not shown. Coordinates were taken from the crystal structure of the human cyclin-dependent kinase 2:ATP complex [(28), PDB entry 1HCK].



Initially, the classification of helicases by residues within motif II grouped functionally related proteins. However, as more helicase sequences became known, many exceptions to the strict conservation of motif II residues within functional families have been found. Indeed, a new DExH domain family containing helicases from positive-strand viruses has been proposed (4). For these enzymes, a variety of residues can be accommodated at the third position. Overall, amino acid sequence comparisons suggest that the DEAD and DEAH domains belong to a larger family of DExx domains, where x indicates sequence variability at the third and fourth positions of the DEAD/DEAH motif.

## 2. General Properties of DEAD and DEAH Domains

Isolated DEAD and DEAH domains typically exhibit NTPase activity that can be stimulated by the addition of nucleic acid polymers. The domains bind NTPs with **dissociation constants** in the micromolar range.  $K_m$  values for ATP of 80, 260, and 250  $\mu M$  were measured for eIF-4A (5), PRP16 (6) and HCV helicases (7), respectively. Unlike eIF-4A, which is an obligate ATPase (8), little specificity for NTP is observed for some DEAD and DEAH domains.  $K_m$  values for ATP, CTP, GTP, and UTP ranged from 260 to 580  $\mu M$  for the PRP16 DEAH domain (6) and from 250 to 1300  $\mu M$  for the HCV DExH helicase (7). For the HCV helicase, the addition of poly U increases its ATPase activity approximately 10-fold (9). In the case of eIF-4A, the addition of poly U increases

the ATPase activity approximately 8-fold (8).

The mechanism by which NTP hydrolysis is coupled to the modification of a nucleic acid secondary structure remains one of the most interesting questions regarding the function of DEAD and DEAH domains. As outlined below, several lines of evidence support the hypothesis that the hydrolysis reaction produces conformational changes in the domain. The changes alter the stability of the bound nucleic acid to favor the single-stranded form.

DEAD and DEAH domains typically contain ~300 residues and frequently reside within a longer sequence, or as subunits of a larger protein complex. For example, each of the seven *Saccharomyces cerevisiae* PRP splicing factors contains a DEAD/H domain. One of the larger splicing factors, PRP22, also contains an amino-terminal RNA-binding domain (6). The human DNA helicase II has two double-stranded **RNA-binding** domains, a DExH domain, and a glycine-rich domain (10). The HCV nonstructural protein 3 (NS3) is a bifunctional enzyme with an amino-terminal [serine proteinase](#) activity and a helicase activity in the carboxy-terminal DExH domain (11).

### 3. DEAD Domains

One of the most well-characterized DEAD domains is eukaryotic initiation factor 4A, eIF-4A. eIF-4A is one of three initiation factors required for binding of the 40S ribosomal subunit to mRNA during the initiation of translation. The ribosome-binding reaction requires ATP hydrolysis. During initiation of translation, eIF-4A may use the energy of ATP hydrolysis to disentangle structured regions in the 5'-untranslated regions of messenger RNA. When complexed with eIF-4B *in vitro*, eIF-4A unwinds duplex RNA in an ATP-dependent manner.

The biochemical characterization of eIF-4A provided the first detailed views of RNA helicase activity, which serve as a framework for comparison of newly discovered helicases. Studies of eIF-4A are reviewed in (12) and (13).

Mutational studies on eIF-4A demonstrated the functional importance of residues in motifs I and II. For example, the conversion of residues within the DEAD sequence to either NEAD or DQAD (substitution of the first conserved aspartic acid with asparagine, or the second residue with glutamine) abolished ATPase activity (12). The same result was observed when the conserved lysine residue in motif I was changed to asparagine. Similar experiments with other DExH domains have confirmed that these amino acid residues are essential.

Studies of the eIF-4A mechanism support the hypothesis that alterations in RNA structure are induced by conformational changes in the DEAD domain. By monitoring the domain's susceptibility to **proteolysis**, ADP was found to cause conformational changes in both the free eIF-4A and eIF-4A:RNA complex (14). Conformational changes were also observed on the binding of an ATP analog to the eIF-4A:RNA complex. Similarly, the observed nucleic acid-dependent differences in eIF-4A affinity for either  $\text{ADP.Mg}^{+2}$  or  $\text{ATP.Mg}^{+2}$  were also thought to reflect structural changes in eIF-4A.

### 4. DEAH Domains

The DEAH domains are involved in RNA processing and display many features common among helicases. Originally classified as DEAH domains, new sequence information has expanded the classification to DEAD/H domains. DEAD/H helicases are associated with the large [ribonucleoprotein](#) (RNP) complex called the spliceosome. The spliceosome processes pre-messenger RNA to remove intervening untranslated sequences called **introns**, so the resultant messenger RNA contains only expressed sequences, or exons.

A number of precursor RNA processing (PRP) proteins associate with the spliceosome. The proteins



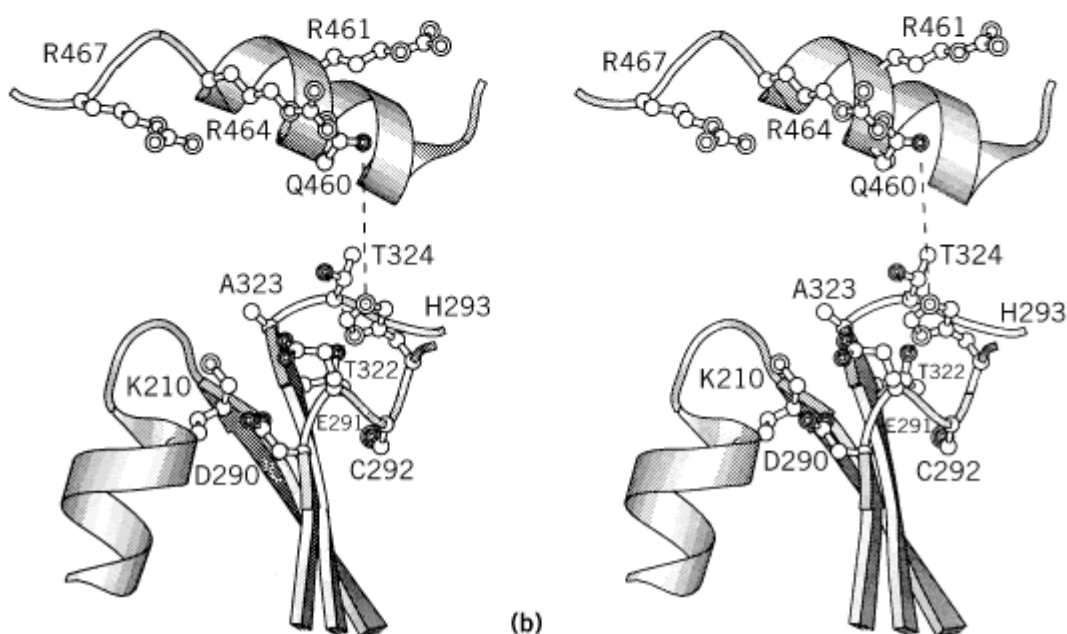
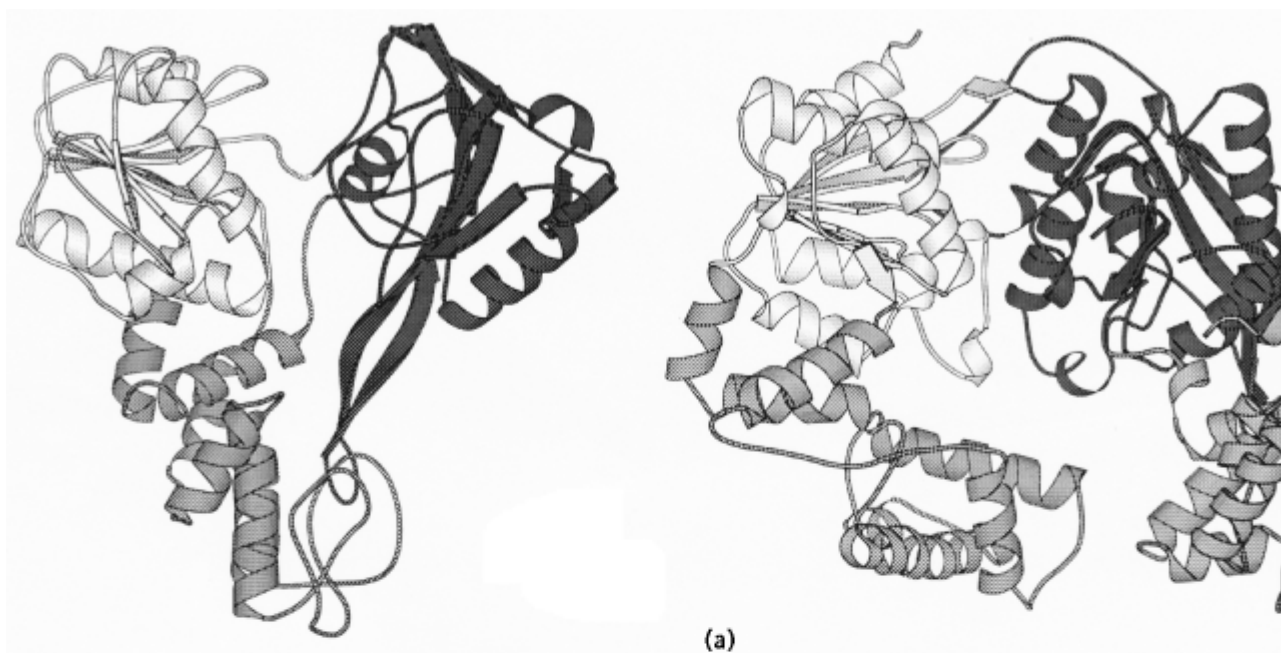
contain either DEAD or DEAH domains and are also called splicing factors. Similar to many helicases, the isolated DEAH domain protein, PRP22, exhibits ATPase activity that is stimulated by the addition of poly(U) RNA (15). Mutagenesis studies of PRP22 variants demonstrated that the ATPase activity requires the conserved lysine residue of motif I. Although helicase activity has been demonstrated for the DEAH splicing factor PRP16 *in vitro*, the functional significance of these data is unclear (6). Most evidence suggests that the PRP splicing factors interact directly with the RNA in the spliceosome, to modulate RNA structure for efficient processing (15-17).

## 5. Structures of DExx and DExH Domains

**X-ray crystallographic** structures of several helicases have been reported. The structures include two DExx DNA helicases belonging to SFI, *Bacillus stearothermophilus* PcrA helicase (18), and *E. coli* Rep helicase (19), and the HCV RNA DExH domain of SFII (20, 21).

The helicases are multidomain structures (Fig. 3). The DNA helicases fold into structurally similar halves, with two domains each. Each half is composed of a bab domain and an **alpha-helical** domain. A deep groove separates the two major helicase domains. The HCV RNA helicase domain contains three folded units: two bab domains and one a-helical domain. Comparison of the helicase structures revealed that the bab domains are structurally similar in both overall fold and their relative spatial orientation (22).

**Figure 3.** Structural features of DExx and DExH domains. The ribbons trace the polypeptide backbones of the HCV RN (left drawing, part a) and *E. coli* Rep DNA helicase (right drawing, part a). The upper left domain in each structure is the bab domain (shown in the lightest ribbon). The NTP-binding domain and second domain (shown in the darkest ribbon) share bab fold. Many of the conserved sequence motifs are located near the deep groove that separates the bab domains. A detailed view of the HCV RNA helicase is shown in stereoscopic representation (part b). Figures prepared with Molscript (29).



Both helicase motifs I and II are located within one of the bab domains. The overall fold of the bab domain and spatial arrangement of residues in motifs I and II are similar to conserved features present in other NTP-utilizing enzymes. The crystallographic structure of the PcrA:ADP complex also showed nucleotide binding to this domain (18). The structural studies taken together, including comparison of the superfamily I and II proteins, have identified the NTP-binding domain.

The structural studies also revealed conformational flexibility that is potentially important in helicase function. Comparison of the independent molecules in the HCV RNA helicase crystals revealed that one of the bab domains could rotate relative to the remainder of the molecule (21). When the helicase structures are compared to those of NTP-binding proteins known to function as molecular switches, the observed structural changes in the helicase were found to involve residues at the end of a [beta-strand](#) previously shown to constitute part of the energy transducing mechanism (18, 21). These residues in the DExH domain map to motif III, which mutagenesis studies have shown is required for coupling NTP hydrolysis to unwinding (12, 23). Some indication of the extent of

possible structural changes is found in the structure of the Rep DNA helicase complexed with a short oligonucleotide, where two orientations differing by a 130-deg rotation about a hinge region of the protein were observed for one of the helical subdomains (19).

## 6. Roles of DEAD and DEAH Residues

The mechanistic roles of the first and last residues of the DEAD and DEAH have been largely deciphered. Structural studies on the free and ADP-bound forms of DExx and DExH domains support the participation of the conserved aspartic acid in binding the NTP phosphate groups. Comparison of the structures of ATP:protein complexes suggests that the aspartic acid probably coordinates either the  $Mg^{+2}$  of the NTP. $Mg^{+2}$  complex, or it interacts with the phosphate groups via an intervening water molecule. However, the exact nature of the interaction has yet to be defined.

The last residue of the motif II tetrapeptide is involved in coupling the energy generated by NTP hydrolysis to duplex unwinding. Researchers working with eIF-4A generated a molecule in which the wild-type DEAD sequence was changed to DEAH (12). The ATPase activity of the modified eIF-4A was slightly enhanced, but the helicase activity was decreased 90% relative to the wild-type molecule. In a similar set of experiments, the DExA mutants of the DExH vaccinia virus nucleoside triphosphate phosphohydrolase II (NP-II) (5) and HCV RNA helicase exhibited NTPase activity and RNA binding comparable to the wild-type protein (24). In contrast, the RNA unwinding activity of these proteins was essentially undetectable.

Crystal structures of the DExH domains show that the histidine side chain is oriented toward the deep groove separating the two bab domains. Although a detailed structural description of the precise role played by the last residue of motif II in the helicase reaction has not been determined, the structures show that the histidine residue is close to conserved residues in motif VI. Mutagenesis experiments on eIF-4A (25) and vaccinia virus NP-II (26) demonstrated that residues in motif VI are required for helicase activity. Altered proteins exhibit decreased ATPase and helicase activities, without compromised RNA binding (27). The structural proximity of the DExH histidine and residues of motif IV offers the possibility that the residues interact directly during the helicase catalytic cycle (Fig. 3).

## 7. Conclusion

The well-characterized DEAD and DEAH domains play essential roles in nucleic acid metabolism. Structural and mutagenesis studies, and amino acid sequence comparisons, have revealed the functions of the first D residue and the terminal D and H residues, although the mechanistic details of their interactions have not been determined. The DEAD and DEAH sequences are required for nucleoside triphosphate hydrolysis and, along with residues of motifs Ia, III, and VI, are required for coupling hydrolysis to mechanical action.

DEAD and DEAH domains occur widely in nature and exhibit a variety of biological functions. The functional diversity of these domains constitutes a challenge for classification by sequence-based approaches. Overall, the DEAD, DEAH, and DExH domains appear to be required in processes involving large-scale conformational changes of domain and substrate. Thus, their utilization as molecular motors is emerging as a unifying theme (14).

## Bibliography

1. S. Schmid and P. Linder (1992) D-E-A-D protein family of putative RNA helicases. *Mol. Microbiol.* **6**, 283–291.
2. A. Gorbalenya and E. Koonin (1993) Helicase: Amino acid sequence comparisons and structure-function relationships. *Curr. Opin. Struct. Biol.* **3**, 419–429.
3. J. Walker, M. Saraste, M. Runswick, and N. Gay (1982) Distantly related sequences in the alpha- and beta-subunits of ATP synthase, myosin, kinases and other ATP-requiring enzymes

and a common nucleotide binding fold. *EMBO J.* **1**, 945–951.

4. E. Koonin (1991) Similarities in RNA helicases. *Nature* **352**, 290.
5. J. Lorsch and D. Herschlag (1998) The DEAD box protein eIF4A. 1. A minimal kinetic and thermodynamic framework reveals coupled binding of RNA and nucleotide. *Biochemistry* **37**, 2180–2193.
6. Y. Wang, J. Wagner, and C. Guthrie (1998) The DEAH-box splicing factor Prp16 unwinds RNA duplexes in vitro. *Curr. Biol.* **8**, 441–451.
7. K. Morgenstern, J. Landro, K. Hsiao, C. Lin, Y. Gu, M. Su, and J. Thomson (1997) Polynucleotide modulation of the protease, nucleoside triphosphatase, and helicase activities of a hepatitis C virus NS3-NS4A complex isolated from transfected COS cells. *J. Virol.* **71**, 3767–3775.
8. F. Rozen, I. Edery, K. Meerovitch, T. Dever, W. Merrick, and N. Sonenberg (1990) Bidirectional RNA helicase activity of eucaryotic translation initiation factors 4A and 4F. *Mol. Cell Biol.* **10**, 1134–1144.
9. F. Preugschat, D. Averett, B. Clarke, and D. Porter (1996) A steady-state and pre-steady-state kinetic analysis of the NTPase activity associated with the hepatitis C virus NS3 helicase domain. *J. Biol. Chem.* **271**, 24449–24457.
10. S. Zhang and F. Grosse (1997) Domain structure of human nuclear DNA helicase II (RNA helicase A). *J. Biol. Chem.* **272**, 11487–11494.
11. L. Jin and D. L. Peterson (1995) Expression, isolation, and characterization of the hepatitis C Virus APTase/RNA helicase. *Archives Biochem. and Biophys.* **323**, 47–53.
12. A. Pause and N. Sonenberg (1992) Mutational analysis of a DEAD box RNA helicase: The mammalian translation initiation factor eIF-4A. *EMBO J.* **11**, 2643–2654.
13. D. Wassarman and J. Steitz (1991) RNA splicing. Alive with DEAD proteins. *Nature* **349**, 463–464.
14. J. Lorsch and D. Herschlag (1998) The DEAD box protein eIF4A. 2. A cycle of nucleotide and RNA-dependent conformational changes. *Biochemistry* **37**, 2194–2206.
15. J. Wagner, E. Jankowsky, M. Company, A. Pyle, and J. Abelson (1998) The DEAH-box protein PRP22 is an ATPase that mediates ATP-dependent mRNA release from the spliceosome and unwinds RNA duplexes. *EMBO J.* **17**, 2926–2937.
16. S. Teigelkamp, M. McGarvey, M. Plumpton, and J. Beggs (1994) The splicing factor PRP2, a putative RNA helicase, interacts directly with pre-mRNA. *EMBO J.* **13**, 888–897.
17. S. Kim and R. Lin (1996) Spliceosome activation by PRP2 ATPase prior to the first transesterification reaction of pre-mRNA splicing. *Mol. Cell Biol.* **16**, 6810–6819.
18. H. Subramanya, L. Bird, J. Brannigan, and D. Wigley (1996) Crystal structure of a DExx box DNA helicase. *Nature* **384**, 379–383.
19. S. Korolev, J. Hsieh, G. Gauss, T. Lohman, and G. Waksman (1997) Major domain swiveling revealed by the crystal structures of complexes of *E. coli* Rep helicase bound to single-stranded DNA and ADP. *Cell* **90**, 635–647.
20. H. Cho, H., N. Ha, L. Kang, K. Chung, S. Back, S. Jang, and B. Oh (1998) Crystal structure of RNA helicase from genotype 1b hepatitis C virus. A feasible mechanism of unwinding duplex RNA. *J. Biol. Chem.* **273**, 15045–15052.
21. N. Yao, T. Hesson, M. Cable, Z. Hong, A. Kwong, H. V. Le, and P. C. Weber (1997) Structure of the hepatitis C virus RNA helicase domain. *Nature Struct. Biol.* **4**, 463–467.
22. S. Korolev, N. Yao, T. Lohman, P. C. Weber, and G. Waksman (1998) Comparisons between the structures of HCV and Rep helicases reveal structural similarities between SF1 and SF2 super-families of helicases. *Prot. Sci.* **7**, 605–610.
23. D. Kim, J. Kim, Y. Gwack, J. Han, and J. Choe (1997) Mutational analysis of the hepatitis C virus RNA helicase. *J. Virol.* **71**, 9400–9409.
24. G. Heilek and M. Peterson (1997) A point mutation abolishes the helicase but not the nucleoside

- triphosphatase activity of hepatitis C virus NS3 protein. *J. Virol.* **71**, 6264–6266.
25. A. Pause, N. M'ethot, and N. Sonenberg (1993) The HRIGRXXXR region of the DEAD box RNA helicase eukaryotic translation initiation factor 4A is required for RNA binding and ATP hydrolysis. *Mol. Cell Biol.* **13**, 6789–6798.
  26. C. Gross and S. Shuman (1996) The QRxGRxGRxxxG motif of the vaccinia virus DExH box RNA helicase NPH-II is required for ATP hydrolysis and RNA unwinding but not for RNA binding. *J. Virol.* **70**, 1706–1713.
  27. C. Gross and S. Shuman (1995) Mutational analysis of vaccinia virus nucleoside triphosphate phosphohydrolase II, a DExH box RNA helicase. *J. Virol.* **69**, 4727–4736.
  28. U. Schulze-Gahmen, H. L. De Bondt, and S.-H. Kim (1996) High-resolution crystal structures of human cyclin-dependent kinase 2 with and without ATP: Bound waters and natural ligand as guides for inhibitor design. *J. Med. Chem.* **39**, 4540–4546.
  29. P. J. Kraulis (1991) Molscript: A program to produce both detailed and schematic plots of protein structures. *J. Appl. Cryst.* **24**, 946–950.

### Suggestions for Further Reading

30. L. Bird, H. Subramanya, and D. Wigley (1998) Helicases: A unifying structural theme? *Curr. Opin. Struct. Biol.* **8**, 14–18.
31. C. Guthrie (1996) "The spliceosome is a dynamic ribonucleoprotein machine". The Harvey Lecture Series, Wiley-Liss, New York, Vol. **90**, pp. 59–80.
32. K. Holmes, C. Sander, and A. Valencia (1993) A new ATP-binding fold in actin, hexokinase and Hsc70. *Trends in Cell Biol.* **3**, 53–59.
33. G. Kadaré and A. Haenni (1997) Virus-encoded RNA helicases. *J. Virol.* **71**, 2583–2590.
34. T. Lohman (1992) Escherichia coli DNA helicases: mechanisms of DNA unwinding. *Mol. Microbiol.* **6**, 5–14.

### Dead-End Inhibition

Dead-end inhibition of an [enzyme](#) occurs when a compound, usually a structural analogue of a substrate, combines at the [active site](#) of the enzyme, but does not undergo reaction. The resulting complex can only dissociate back to the components from which it was formed. Dead-end inhibitors are useful for determining kinetic reaction mechanisms because dead-end inhibition patterns can be characteristic of a particular mechanism. Rules have been formulated for the qualitative prediction of such inhibition patterns ([1](#)) and these rules can be summarized as follows:

1. A dead-end inhibitor affects the *slope* of a double-reciprocal plot (see [Lineweaver–Burk Plot](#)) when
  - the inhibitor and the variable substrate combine with the same enzyme form or
  - a reversible connection exists between the points of addition of the inhibitor and the variable substrate. Such a reversible connection exists only when the variable substrate adds after the inhibitor.
2. A dead-end inhibitor affects the *intercept* of a double-reciprocal plot when the inhibitor and variable substrate combine with different forms of enzyme and saturation with the variable substrate does not overcome the inhibition. Intercept effects will always be observed unless steps between the points of addition of the variable substrate and the inhibitor are at thermodynamic equilibrium.

The results of the application of the rules to kinetic mechanisms for Bi–Bi reaction mechanisms are given in Table 1. Further examples are given in [Competitive inhibition](#), [Noncompetitive Inhibition](#), [Product Inhibition](#), [Substrate Inhibition](#), and [Uncompetitive inhibition](#).

**Table 1. Dead-end Inhibition Patterns for Bi–Bi Reaction Mechanisms<sup>a</sup>**

| Mechanism                    | Inhibitor Combines with | Variable Substrate |          |
|------------------------------|-------------------------|--------------------|----------|
|                              |                         | <i>A</i>           | <i>B</i> |
| Ordered                      | E                       | C                  | NC       |
|                              | EA                      | UC                 | C        |
|                              | EQ                      | UC                 | UC       |
|                              | EA + EQ                 | UC                 | NC       |
| Equilibrium-ordered          | E                       | C                  | C        |
|                              | EA                      | UC                 | C        |
| Rapid equilibrium,<br>random | E + EB                  | C                  | NC       |
|                              | E + EA                  | NC                 | C        |
| Ping-pong                    | E                       | C                  | UC       |
|                              | F                       | UC                 | C        |
|                              | E + F                   | NC                 | NC       |

<sup>a</sup> The mechanisms are illustrated under **Kinetic mechanisms**. C = linear [competitive inhibition](#). NC = linear [noncompetitive inhibition](#). UC = linear [uncompetitive inhibition](#).

## Bibliography

1. W. W. Cleland (1963) *Biochim. Biophys. Acta* **67**, 188–196.

## Deamidation

Deamidation is the reverse reaction of amidation (see [Carboxyl Groups](#)) and involves release of an amine from an amide. Deamidation occurs spontaneously in proteins at [asparagine](#) and [glutamine](#) residues, which causes heterogeneity, instability, and sometimes defects in the function of the proteins.

Deamidation in proteins occurs enzymatically or chemically. [Enzymes](#) that participate in deamidating of proteins are transglutaminase and peptidoglutaminase. An amide bond in an asparagine or glutamine residue is relatively labile, especially in acidic media, and is easily

hydrolyzed chemically. The rate of deamidation is sensitive to the environment of the amide. Exposed amides are more readily deamidated than those buried in the interior of a protein. Thus, denatured proteins often have high potential for deamidation.

The Asn-Gly sequence is especially notorious for rapid deamidation. The reaction shown below proceeds readily in this sequence. Additional reactions such as **racemization**, transpeptidation, and cleavage of peptide bonds take place during this reaction. These reactions also occur with the Asp-Gly sequence and pivotal components of protein deterioration.

The rates of deamidation vary with different peptides. The **half-times** for deamidation of various asparagine residues under the same conditions range from 3.3 to 277 days. Deamidation causes conformational destabilization of proteins and accelerates their degradation. Thus the deamidation of Asn or Gln residues has been regarded as a timer of protein turnover (1).

Asparaginase hydrolyzes asparagine to aspartic acid (deamidation). Because asparagine in the blood is an essential nutrient for the growth of malignant white cells, the intravenous administration of asparaginase suppresses the growth of some types of animal and human leukemias.

#### Bibliography

1. A.B. Robinson, J.H. McKerrow, and P. Cary (1970) Proc. Natl. Acad. Sci. USA **66**, 753–757.

#### Suggestion for Further Reading

2. H.T. Wright (1991) Nonenzymatic deamidation of asparagine and glutamine residues in proteins. Crit. Rev. Biochem. Mol. Biol. **26**, 1–52.

## Decapentaplegic

The *decapentaplegic* (*dpp*) gene plays a key role in patterning both embryonic and adult structures in the fruitfly *Drosophila*. The name of the **locus** underscores the multiple requirements for *dpp* activity during development and refers to the wide range of defects caused by **mutations** in the **gene** (*decapentaplegic* = 15 defects) (1). *dpp* is one of the few haplo-insufficient loci in *Drosophila*; consequently, even a 50% reduction in *dpp* activity leads to embryonic lethality. Less severe alleles of *dpp* are viable but result in adults that have defects in structures derived from one or more **imaginal discs**. The *dpp* gene encodes a **protein** that is **homologous** to secreted **growth factors** of the **transforming growth factor- $\beta$**  (TGF- $\beta$ ) **superfamily**. Ligands that belong to this group of signaling molecules have been identified in animals across the phylogenetic spectrum and are involved in cellular and developmental processes ranging from regulation of the **cell cycle** and the **extracellular matrix**, to establishment of embryonic pattern and aging. Genetic and molecular studies of *dpp* and the genes involved in *dpp* **signal transduction** have proven to be of broad general interest, because the basic mechanisms involved in generating and receiving TGF- $\beta$  signals are evolutionarily conserved.

### 1. Molecular Features

The *dpp* gene spans ~55 kb of DNA, the bulk of which consists of a large array of *cis*-regulatory elements (2). *dpp* **messenger RNA** is expressed in a complex spatial and temporal pattern that reflects the multiple roles the ligand plays during development. Three major **cis-acting** domains regulate transcription of the gene: an upstream *dpp*-shv region required for expression in the

embryonic gut and pupal wing, a central *dpp*-Hin region that directs embryonic expression, and a downstream *dpp*-disk region that regulates transcription in imaginal discs (2). These dispersed regulatory elements contribute to the genetic complexity of the locus and mediate *dpp* transcription in response to inputs from a variety of signaling pathways. For example, *dpp* expression in the blastoderm stage embryo is regulated by **Dorsal**, a [transcription factor](#) homologous to mammalian NFκB. Later in embryogenesis, *dpp* transcription in the dorsal-most cells of the germband is controlled by genes acting in the Jun kinase signaling pathway, while expression of *dpp* in the midgut is dependent on multiple factors, including **homeobox** DNA-binding proteins, the [growth factor](#) Wingless, and Dpp itself. During imaginal disc development, localized expression of *dpp* in the wing and leg discs is regulated by the *hedgehog* and *wingless* genes.

The *dpp* locus gives rise to three major and two minor transcripts through the use of alternative **promoters** that are utilized at different stages of development (2). While the mRNAs have unique 5'-untranslated sequences, they share common second and third exons that contain an open reading frame coding for a single 588-amino-acid-residue protein (see [Cell Cycle](#) and [Extracellular Matrix](#)). Dpp is most closely related to the vertebrate bone morphogenetic proteins BMP-2 and BMP-4, and it shares ~75% identity with these proteins in the carboxy-terminal region that constitutes the mature ligand **domain**. While the vertebrate BMPs were first identified by their ability to induce ectopic bone, they are now known to have important roles in embryonic development in a number of organisms, including frogs, mice, and humans. Dpp and BMP-4 can substitute for one another functionally, because a human BMP-4 transgene can rescue patterning defects in a *Drosophila* embryo lacking Dpp (3). Conversely, the fly protein can induce the formation of ectopic bone when injected subcutaneously into rats. Like other ligands belonging to the TGF-β superfamily, Dpp is processed and secreted as a **disulfide**-linked dimer of ~30kDa.

## 2. Dpp Signal Transduction

Dpp acts as a secreted ligand to influence the developmental fate of cells that receive the signal. Genes involved in Dpp signal transduction have been primarily identified using two strategies. Some genes have been recovered in genetic screens to isolate [enhancers](#) of weak *dpp* **alleles**, while others have been isolated in low stringency **hybridization** screens, based on their homology to TGF-β signaling components identified in other organisms.

According to the current paradigm for BMP signaling, the ligand binds a heteromeric complex of two structurally related transmembrane **serine-threonine kinases**, called the type I and type II receptors (4, 5). Formation of the ligand–receptor complex allows the type II kinase to **phosphorylate** and activate the type I receptor. A type II receptor, Punt, as well as two type I receptors, Thick veins (Tkv) and Saxophone (Sax), have been implicated in Dpp signaling. Although the role of Punt and Tkv as receptors for Dpp is well established, it appears that Sax may primarily mediate the response to other BMP-related ligands in *Drosophila*. Activation of Tkv results in the direct phosphorylation of a cytoplasmic protein encoded by *mothers against dpp* (*mad*). This modification triggers **Mad** to form a complex with Medea (a structurally related protein), and it enables their translocation from the cytoplasm into the [nucleus](#). Mad and Medea contain DNA-binding domains and are thought to regulate the expression of downstream target genes in association with other transcription factors (4, 6). The Mad family of proteins is evolutionarily conserved. The human homologue of Medea (DPC4) has been identified as a **tumor-suppressor** gene, a result that appears logical in light of the known antiproliferative effects of TGF-β.

## 3. Biological Role of Dpp

Dpp plays an important role in the specification of cell fate and morphogenesis in a number of tissues (1). Mutations in *dpp* that interfere with production of an active ligand affect all developmental events regulated by the gene, while mutations that disrupt specific enhancer elements affect *dpp* function in a tissue- or stage-specific manner. Among the processes that require *dpp* signaling are: oogenesis; establishment of dorsal–ventral pattern during embryogenesis;



morphogenetic movements of dorsal closure; subdivision of the mesoderm along the dorsal–ventral axis; specification of the visceral mesoderm and endoderm in the embryonic gut; development of the heart, gastric caecae, salivary glands, and the trachea; and growth and patterning of imaginal discs.

The ability of Dpp to trigger distinct responses in a single field of cells has generated a great deal of interest in understanding the mechanisms underlying *dpp* function. Dpp has been shown to specify cell fate in a concentration-dependent manner; that is cells respond to different thresholds of Dpp by following distinct pathways of differentiation. A critical issue that arises is how gradients of *dpp* activity are established, and how such gradients are interpreted. Based on recent studies, two alternative mechanisms have evolved: one is based on [diffusion](#) of Dpp from its site of synthesis to generate a protein concentration gradient, and the second involves an inhibitor that diffuses into the domain where *dpp* is expressed and interferes with Dpp signaling in a graded manner. This results in a gradient of ligand activity, rather than concentration. Both types of gradients are discussed further below.

Adult structures in the fruitfly arise from imaginal discs, small groups of epithelial cells that are set aside during embryogenesis. These discs grow and are patterned during the larval and pupal stages in response to different signals. In the wing imaginal disc, Dpp acts as a long-range [morphogen](#) to specify cell fate along the anterior–posterior axis. *dpp* is expressed in a narrow domain of cells at the anterior–posterior compartment boundary, from where it diffuses to generate a gradient of protein in the surrounding tissue. Cells up to 20 cell diameters away respond to different threshold concentrations of Dpp by activating the transcription of target genes like *spalt* and *optomotor blind* ([7](#), [8](#)). Expression of *spalt* occurs close to the source of *dpp* protein, while *optomotor blind* expression overlaps with, and extends further than, *spalt* expression. The nested domains of Dpp target gene expression further subdivide the wing disc into distinct regions that are specified by the combination of genes expressed.

During early development, a gradient of Dpp signaling is required to establish cell fates within the dorsal half of the embryo. Peak levels of Dpp signaling specify the dorsal-most amnioserosa tissue, while lower levels are required to specify the dorsal ectoderm ([9](#), [10](#)). Reduction in the level of Dpp signaling results in progressive loss of dorsal structures, while increasing concentrations of Dpp can induce dorsal cell fates. Since *dpp* mRNA is expressed uniformly in all dorsal cells, it is generally believed that Dpp activity, rather than its concentration, is graded. A number of extracellular proteins are involved in generating a gradient of Dpp signaling in the embryo. A second BMP ligand, Screw, acts synergistically with Dpp to enhance signaling in the dorsal-most cells. The activity of these ligands is antagonized by a secreted factor, Short gastrulation (Sog), that can prevent ligand binding to the receptor. The inhibitor Sog is expressed in ventral cells and diffuses dorsally to generate a ligand gradient of the opposite polarity. In addition, a [metalloproteinase](#), Tolloid (Tld), promotes signaling in dorsal cells by cleaving Sog and releasing the ligand. Thus modulation of ligand activity at multiple levels contributes to establishment of a gradient of BMP signaling in the embryo ([11-14](#)).

Similar antagonistic interactions involving homologous proteins are involved in patterning the dorsal–ventral axis in vertebrate embryos. In *Xenopus*, BMP-4 promotes ventral development, while a homologue of Sog (Chordin) promotes dorsal cell fates. Recently an amphibian homologue of Tld (Xolloid) has been shown to cleave Chordin ([12](#)). These and other studies suggest that the dorsal–ventral axes in *Drosophila* and vertebrates are specified by a similar mechanism, although they are inverted relative to one another. The extensive parallels between TGF- $\beta$ /BMP signaling in vertebrates and invertebrates allow one to extend the insights gained from studying Dpp signaling in *Drosophila* to other organisms, including humans.

## Bibliography

1. F. A. Spencer, F. M. Hoffman, and W. M. Gelbart (1982) Decapentaplegic: a gene complex affecting morphogenesis in *Drosophila melanogaster*. *Cell* **28**, 451–461.
2. R. D. St. Johnston, F. M. Hoffman, R. K. Blackman, D. Segal, R. Grimaila, R. W. Padgett, H.

- A. Irick, and W. M. Gelbart (1990) The molecular organization of the *decapentaplegic* gene in *Drosophila melanogaster*. *Genes Dev.* **4**, 1114–1127.
3. R. W. Padgett, J. M. Wozney, and W. M. Gelbart (1993) Human BMP sequences can confer normal dorsal–ventral patterning in the *Drosophila* embryo. *Proc. Natl. Acad. Sci. USA* **90**, 2905–2909.
  4. C.-H. Heldin, K. Miyazono, and P. ten Dijke (1997) TGF- signalling from cell membrane to nucleus through SMAD proteins. *Nature* **390**, 465–471.
  5. J. Massague (1996) TGF- signaling: receptors, transducers, and Mad proteins. *Cell* **85**, 947–950.
  6. K. Arora, H. Dai, S. G. Kazuko, J. Jamal, M. B. O'Connor, A. Letsou, and R. Warrior (1995) The *Drosophila schnurri* gene acts in the Dpp/TGF- signaling pathway and encodes a transcription factor homologous to the human MBP family. *Cell* **81**, 781–790.
  7. T. Lecuit, W. J. Brook, M. Ng, M. Calleja, H. Sun, and S. M. Cohen (1996) Two distinct mechanisms for long-range patterning by Decapentaplegic in the *Drosophila* wing. *Nature* **381**, 1387–1393.
  8. D. Nellen, R. Burke, G. Struhl, and K. Basler (1996) Direct and long-range action of a DPP morphogen gradient. *Cell* **85**, 357–368.
  9. E. L. Ferguson and K. V. Anderson (1992) Decapentaplegic acts as a morphogen to organize dorsal–ventral pattern in the *Drosophila* embryo. *Cell* **71**, 451–461.
  10. K. A. Wharton, R. P. Ray, and W. M. Gelbart (1993) An activity gradient of *decapentaplegic* is necessary for the specification of dorsal pattern elements in the *Drosophila* embryo. *Development* **117**, 807–22.
  11. G. Marqués, M. Musacchio, M. J. Shimell, K. Wunnenberg-Stapleton, K. W. Cho, and M. B. O'Connor (1997) Production of a DPP activity gradient in the early *Drosophila* embryo through the opposing actions of the SOG and TLD proteins. *Cell* **91**, 417–426.
  12. S. Piccolo, E. Agius, B. Lu, S. Goodman, L. Dale, and E. M. De Robertis (1997) Cleavage of Chordin by Xolloid metalloprotease suggests a role for proteolytic processing in the regulation of Spemann organizer activity. *Cell* **91**, 407–416.
  13. J. L. Neul and E. L. Ferguson (1998) Spatially restricted activation of the SAX receptor by SCW modulates DPP/TKV signaling in *Drosophila* dorsal-ventral patterning. *Cell* **95**, 483–494.
  14. M. Nguyen, S. Park, G. Marqués, and K. Arora (1998) Interpretation of a BMP activity gradient in *Drosophila* embryos depends on synergistic signaling by two type I receptors, SAX and TKV. *Cell* **95**, 495–506.

### Suggestions for Further Reading

15. Flybase Consortium (1998) FlyBase: a *Drosophila* database. *Nucleic Acids Res.* **26**, 85–88. (<http://flybase.bio.indiana.edu/>)
16. W. M. Gelbart (1989) The decapentaplegic gene: a TGF- homolog controlling pattern formation in *Drosophila*. *Development* **107**, 65–74.
17. P. A. Lawrence and G. Struhl (1996) Morphogens, compartments, and pattern: lessons from *Drosophila*? *Cell* **85**, 951–961.

### Degeneracy of the Genetic Code

The [genetic code](#) for [translation](#) of the gene sequence in **protein biosynthesis** is degenerate—there is more than one triplet **codon** for most [amino acids](#). For example, tyrosine is encoded by UAU and

UAC, and arginine is encoded by CGU, CGC, CGA, CGG, AGA, and AGG. Two codons that share the same first two nucleotides encode the same amino acid if the third nucleotide is either C or U, and they often do if the third nucleotide is A or G. Alternate synonymous codons are not used randomly. The GC content of the DNA, especially of bacteria, varies widely in different species, and there is a corresponding difference in **codon usage**. In bacteria and the yeast *Saccharomyces cerevisiae*, and even to some extent in *Drosophila* (except for the [histone](#) genes), the most frequently used codons are those that are decoded by abundant [transfer RNA](#), and they are strongly preferred in genes expressed at high levels. In genes expressed at low levels, codon usage is more uniform (1-3). In bacteria with extreme DNA base composition biases (for example, the *Mycoplasma capricolum* genome is 25% G + C, whereas that of *Micrococcus luteus* is 74% G + C), not surprisingly, all genes have similar codon usage using those codons with the corresponding G + C content (4). Codon usage also varies among genes from the same genome. This is true of mammalian genomes but also, to some extent, in *S. cerevisiae*. Regions of differing G + C content were initially recognized in mammalian [chromosomes](#) and termed [isochores](#). Overlying differences in codon usage between individual genes is the apparent higher density of genes in GC-rich isochores, and this may influence codon usage. As genome sequencing projects proceed, this relationship is expected to be clarified.

### Bibliography

1. T. Ikemura (1981) *J. Mol. Biol.* **151**, 389–409.
2. M. Gouy and C. Gautier (1982) *Nucl. Acids Res.* **10**, 7055–7074.
3. H. Dong, L. Nilsson, and C. G. Kurland (1996) *J. Mol. Biol.* **260**, 649–663.
4. T. Ohama, A. Muto, and S. Osawa (1990) *Nucl. Acids Res.* **18**, 1565–1569.

### Suggestion for Further Reading

5. P. M. Sharp and G. Matassi (1994) Codon usage and genome evolution. *Curr. Opin. Genet. Develop.* **4**, 851–860.

## Dehydrogenases

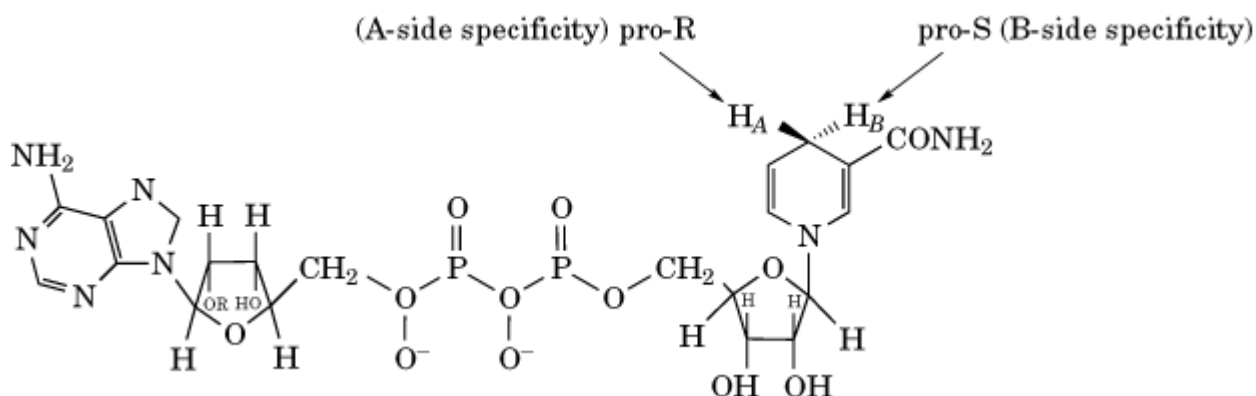
Pyridine nucleotide-linked dehydrogenases are [enzymes](#) that catalyze oxidation–reduction reactions involving a pyridine nucleotide and the transfer of a hydride (one proton and two electrons) from or to a second substrate molecule. Most dehydrogenase reactions are readily reversible, but some are not, and the enzymes responsible for essentially irreversible reactions are referred to as *reductases*; there are exceptions, however, such as *formate dehydrogenase*. The oxidized forms of the pyridine nucleotides are **NAD** (nicotinamide adenine dinucleotide) and **NADP**. The latter nucleotide has an additional phosphoryl group attached at the 2'-position of the ribose that is part of the adenosyl moiety (Fig. 1). The general form of a dehydrogenase reaction can be illustrated by reference to the reaction catalyzed by [lactate dehydrogenase](#) (Eq. 1):



Not all dehydrogenase reactions are as simple as that described by Equation 1. The oxidation reactions catalyzed by *isocitrate* and *glutamate* dehydrogenases give rise to the release of CO<sub>2</sub> and NH<sub>3</sub>, respectively. NAD and NADP are written frequently as NAD<sup>+</sup> and NADP<sup>+</sup> even though the overall charge on the molecule is negative in the neutral pH region; the plus refers to the charge on the nitrogen atom of the nicotinamide moiety of the oxidized forms of the pyridine nucleotides.

Dehydrogenases usually show a distinct preference for one or the other of the two pyridine nucleotides; examples are given in Table 1, and a more extensive list is available (1). The binding to the enzyme of the nonpreferred nucleotide is generally weaker, and the maximum velocity is usually lower.

**Figure 1.** Structure of NADH / NADPH, showing the two prochiral hydrogen atoms at the 4-position of the nicotinamide moiety. R at the 2'-position of the ribose moiety represents either a hydrogen atom (NADH) or phosphoryl group (NADPH).



**Table 1. Pyridine Nucleotide Specificity of Some Dehydrogenases**

| Dehydrogenase/Reductase    | Pyridine Nucleotide | Stereospecificity |
|----------------------------|---------------------|-------------------|
| Alcohol                    | NAD                 | A                 |
| Glyceraldehyde-3-phosphate | NAD                 | B                 |
| 3-Phosphoglycerate         | NAD                 | B                 |
| Lactate                    | NAD                 | A                 |
| Glutamate                  | NAD/NADP            | B                 |
| Glutathione                | NAD/NADP            | B                 |
| Dihydrofolate              | NADP                | A                 |
| Glucose-6-phosphate        | NADP                | B                 |
| Isocitrate                 | NADP                | A                 |

In dehydrogenase reactions (Eq. 1), the hydride is transferred to the 4-position of the nicotinamide ring of NAD or NADP (Fig. 1). The two hydrogen atoms at this position in NADH and NADPH are chemically equivalent, and the tetrahedral carbon atom is not **chiral**. However, the two hydrogen atoms at the 4-position are **prochiral**. The binding of reduced pyridine nucleotides to dehydrogenases is such that hydride transfer involves only one of the two hydrogen atoms. Thus, the enzymes exhibit regioselectivity, and hydride transfer occurs either from the pro-R hydrogen (A-side specificity) or from the pro-S hydrogen (B-side specificity). Examples of dehydrogenases that fall into each of these two classes are given in Table 1. Attempts have been made to explain why some

dehydrogenases are A-side and others are B-side, but no convincing explanation has been forthcoming (2).

The side-specificity of dehydrogenases allows the preparation of reduced pyridine nucleotides with a deuterium (D) or tritium (T) atom replacing a hydrogen atom in either the pro-R or pro-S positions (see [Hydrogen Isotopes](#)). Extensive use has been made of deuterio- and tritio-pyridine nucleotides to gain information about enzymic catalysis (3). The cleavage of a bond can be expected to be two to eight times slower than with a bond. With chemical reactions, the full isotope effect difference in rates is usually observed, because the bond-breaking step is rate-limiting. But with enzyme-catalyzed reactions, steps other than the catalytic one can be at least partly rate-limiting. This results in the observed isotope effect, with NADD or NADPD being smaller than the absolute value for the bond-breaking step. Indications as to where the rate-limiting steps lie for a particular dehydrogenase reaction come from the use of initial velocity studies to determine the apparent first-order rate constants  $k$  and the maximum velocities  $V$  with the protio- and deuterio-substrates. The observed isotope effects are then calculated from the ratio of the values for  $(V/K)_H$  and  $(V/K)_D$ , expressed as  $^D(V/K)$  and the ratio of the values for  $V_H$  and  $V_D$ , expressed as  $^D V$ . The values for these ratios usually range from 1.0 to 3.5. Low values for  $^D(V/K)$  can result from the stickiness of substrates. Values for  $^D V$  will be low whenever product release is rate-limiting (3). The tritiated forms, NADT and NADPT, can be used for determinations of  $^T(V/K)$ , but not from initial velocity studies, as the tritium is only a trace label. For this same reason, it is not possible to obtain values for  $^T V$ . The value for  $^T(V/K)$  will be higher than that for  $^D(V/K)$  as the bond is stiffer than the bond. As  $^T k = (^D k)^{1.442}$ , the general function  $[\ ^D k - 1 ] / [ \ ^D k^{1.442} - 1 ]$  can be used to determine the intrinsic isotope effect for a dehydrogenase reaction (4). It is this value that gives insights into the chemical mechanism of the reaction and the nature of the [transition state](#).

The pyridine nucleotides that function as effective substrates for most dehydrogenases have a nicotinamide-ribose linkage (Fig. 1) with a b-anomeric configuration and are referred to, for example, as b-NAD. These pyridine nucleotides also occur with an a-anomeric configuration and are found in tissue extracts, as well as commercial preparations. Both a-NADH and b-NADH undergo thermal **epimerization** reactions to produce an equilibrium mixture with 10% of the a form (5). There are no examples of a dehydrogenase that specifically uses a-NADH or a-NADPH, although [alcohol dehydrogenase](#) and lactate dehydrogenases utilize a-NADH with maximum velocities some three to four orders of magnitude lower than with b-NADH. The same side-specificity is observed with both anomers of NADH (5). By contrast, a-NADPH is a relatively good substrate for the [dihydrofolate reductase](#) that is encoded by an R-plasmid (6). The  $K_m$  value of 17  $\mu\text{M}$  is comparable to the value of 4  $\mu\text{M}$  for b-NADPH. The maximum velocity of the reaction with a-NADPH is 70% of that with b-NADPH. It has been proposed that those dehydrogenases which show activity with both the a- and b-anomers must have sufficient conformational freedom to permit specific orientation of the dihydronicotinamide ring for the reduction of other substrates.

#### Bibliography

1. G. Popjak (1970) *The Enzymes* **2**, 115–215.
2. N. J. Oppenheimer (1984) *J. Am. Chem. Soc.* **106**, 3032–3033.
3. W. W. Cleland (1982) *CRC Crit. Rev. Biochem.* **13**, 385–428.
4. D. B. Northrop (1975) *Biochemistry* **14**, 2644–2651.
5. N. J. Oppenheimer and N. O. Kaplan (1975) *Arch. Biochem. Biophys.* **166**, 526–535.
6. S. L. Smith and J. J. Burchall (1983) *Proc. Natl. Acad. Sci.* **80**, 4619–4623.

## Denaturants Stabilizers

The structures of globular [proteins](#) (and other biological macromolecules) are not very stable (see [Protein Stability](#)). Therefore, it is often necessary to increase the stability by addition of inert compounds in order to maintain a biochemical activity or to preserve an assembled structure, such as various **organelles** or multisubunit enzymes. The stabilizing compounds must characteristically be added at high concentration (>1 M), since their interactions with proteins are weak (see [Stabilization And Destabilization By Co-Solvents](#)). Strongly interacting molecules, although effective in particular cases, generally tend to alter some properties of the protein molecules (see [Binding](#)).

The most frequently used stabilizing agents are 1 M sucrose, other sugars (trehalose, glucose), 3 to 4 M glycerol, neutral [amino acids](#) at 1 M levels (glycine, proline, alanine), methyl amines (betaine, sarcosine, trimethyl amine N oxide), polyols (sorbitol, inositol, mannitol), some salts [ $\text{Na}_2\text{SO}_4$ ,  $(\text{NH}_4)_2\text{SO}_4$ ] (see [Sulfate Salts](#)) (1). All have the property of being preferentially excluded from proteins, which imparts stabilization and precludes excessive surface contacts that could alter the local protein structure (see [Preferential Hydration](#)). Most stabilizers used in the laboratory are also cryoprotectants and osmolytes, in that they are chosen by nature to maintain high osmotic pressure in, eg, amphibians, and to protect against freezing (2).

Conversely, a number of compounds are used to denature proteins when so desired. These again act at high concentration and are characterized by their [preferential binding](#) to proteins. The most common denaturants are [urea](#), **guanidinium** chloride, sodium dodecyl sulfate (**SDS**), trichloroacetic acid,  $\text{CaCl}_2$ , LiBr, and other salts on the denaturant side of the [Hofmeister series](#) (3). (See [Guanidinium Salts](#), [Hofmeister Series](#).)

### Bibliography

1. S. N. Timasheff (1992) In *Stability of Protein Pharmaceuticals* (T. J. Ahern and M. C. Manning, eds.), Plenum, New York, pp. 265, 286.
2. P. H. Yancey, M. E. Clark, S. C. Hand, R. D. Bowlus, and G. N. Somero (1982) *Science* **217**, 1214–1222.
3. P. H. von Hippel and T. Schleich (1969) In *Structure and Stability of Biological Macromolecules* (S. N. Timasheff and G. D. Fasman, eds.), Marcel Dekker, New York, Chap. "6".

## Denaturation Mapping

**Denaturation** of the **double-helical** structure of **DNA** is an essential first step in many of the most common methods for investigating **chromosomal** structure. The method relies on the unique order of bases in DNA and the capacity of two separate complementary strands to recognize each other through the process of **hybridization**. Denaturation mapping can be done using either sections of tissue or chromosomal preparations. The sample slide to be analyzed is heated to denature the double helix into two separate DNA strands. Rapid cooling of the preparation in the presence of formaldehyde keeps the DNA strands separated. The cloned gene or sequence of interest is prepared as a DNA probe in single-stranded form with a label attached. This label could be either a source of [radioactivity](#), or an [epitope](#) recognized by an [antibody](#) (often using [streptavidin](#) as an [antigen](#)). A

change in salt concentration and temperature facilitates hybridization between complementary DNA strands. The probe and chromosomal sequences compete with each other, but a sufficient excess of probe can be added so that some of the probe hybridizes to the specific site on a chromosome.

When initially developed by Pardue and Gall in 1970, the probe was labeled with the radioactive [hydrogen isotope](#) tritium (1). More recently the use of antibodies has become popular (2, 3). The antibodies are usually linked to an enzyme that acts on an appropriate substrate to generate a visible signal. This could either be a colored precipitate on the chromosome, that is detectable under a light microscope, or the antibody could be linked to a fluorochrome that is visualized by a **fluorescence** microscope. This latter technique is known as fluorescent [in situ hybridization](#) (FISH). Now modifications of this technique allow analysis under the electron microscope for increased resolution or in the **confocal microscope** that builds up three-dimensional images by optical sectioning.

A number of specialized methods have acquired their own acronyms. These include fluorescence in situ hybridization (FISH) as described; genomic in situ hybridization (GISH); whole chromosome painting (WCP); and primed in situ labeling (PRINS). GISH is a method for detecting species-specific chromosomes when the entire genomic DNA is labeled and hybridized to whole chromosomal spreads. It is useful in analyzing of cell fusion hybrids and in evolutionary studies. In WCP, DNA from a single chromosome is labeled and prehybridized with repetitive DNA to stop cross-reaction with other chromosomes from the same organism. WCP probes are available commercially for all 22 [autosomes](#) and the sex chromosomes in humans. In PRINS, the polymerase chain reaction (**PCR**) is used to label sequences on the actual chromosome or tissue. There are many other adaptations of these important methodologies.

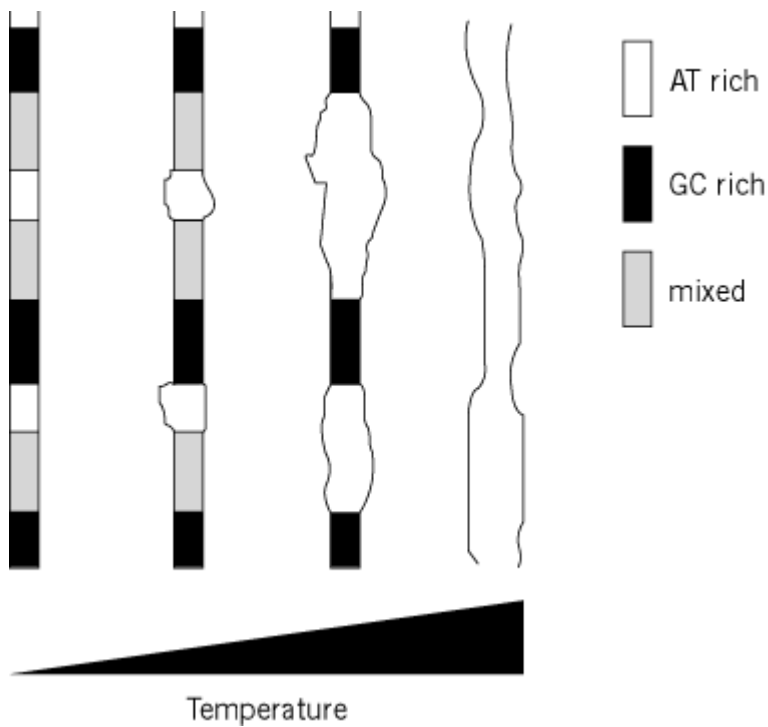
#### Bibliography

1. M. L. Pardue and J. Gall (1970) *Science* **168**, 1356–1358.
2. J. Leary, D. Brigati, and D. C. Ward (1983) *Proc. Natl. Acad. Sci. USA* **80**, 4045–4049.
3. P. Lichter et al. (1990) *Science* **247**, 64–68.

#### Denaturation, Nucleic Acids

Disruption of the secondary and tertiary structures of nucleic acid polymers is promoted by exposure to extreme conditions. Most commonly, nucleic acid denaturation refers to the formation of two free single-stranded nucleic acids from a duplex. Denaturation (loss of native structure) of the polymer duplex can be promoted by a variety of conditions or agents, including heat, chemical denaturants, and extremes of pH. The stability of the duplex varies by region, and the base pair composition (% GC) determines the stability of a given region (Fig. 1). Under most conditions, regions of the duplex that are composed predominantly of AT (for DNA) or AU (for RNA) base pairs are less stable than GC-rich regions. As the conditions shift toward those favoring single-stranded polymers, the path from fully base-paired duplex to fully denatured single strands passes through intermediate states, where first the AT/AU rich regions “melt,” forming internal loops, then mixed sequence regions, then finally the GC-rich regions. As the last GC-rich regions “melt,” the strands separate.

**Figure 1.** Mechanism of melting of polymeric DNA.



## 1. Thermal Denaturation

Because it is easy to manipulate, heat is frequently employed to denature nucleic acids. The temperature can be changed rapidly and reversibly in the small volumes typically employed in nucleic acid studies. Heating to temperatures above the  $T_m$  and rapid cooling can trap the strands of a duplex in the separated state. The melting transitions of mixed-sequence nucleic acid polymer duplexes are quite broad. To ensure strand separation, the duplex must be exposed to temperatures significantly above  $T_m$ . Failure to separate the strands results in a prenucleated duplex that can readily reanneal. This method does not work for oligonucleotide duplexes for which thermal denaturation is readily reversible.

Equilibrium thermal denaturation curves are used to study duplex (or other structure) thermodynamic properties. The extent of the thermally induced denaturation transition as a function of temperature is monitored either by some optical property, usually ultraviolet absorbance, or via heat capacity, as in differential scanning calorimetry. "Melting curves," where the observed parameter is plotted as a function of temperature, contain quantitative information about the stability of the duplex and the thermodynamic origins of that stability (1). Optically monitored thermal denaturation curves are analyzed using the van't Hoff model and can readily provide thermodynamic ( $\Delta H^\circ$ ,  $\Delta S^\circ$ , and  $\Delta G^\circ$ ) and extrathermodynamic data ( $T_m$ ). Not all melting transitions are consistent with the assumptions underlying the models used to extract thermodynamic data from equilibrium melting curves. Differential scanning calorimetry provides a direct model-independent means of generating thermodynamic data on thermally induced denaturation.

## 2. Chemical Denaturation

Chemical denaturants reduce the thermal stability of nucleic acid duplexes. Some denaturants, such as formaldehyde, form covalent bonds to the bases. Others act by altering the noncovalent interactions of the solvent with the nucleotide bases. A wide variety of reagents can induce denaturation of DNA (2). The isolated bases that, when attached to the backbone, comprise a nucleic acid are rather insoluble in aqueous solution. The **hydrophobic** character of the bases, coupled with



the [hydrophilic](#) character of the sugar–phosphate backbone, push the single-strand–duplex equilibrium toward duplex formation. The stacked configuration of the bases in a duplex is favored over the unstacked bases of the single strand (see [Base Pairs](#)). There are conditions under which the bases in a single strand will stack, but the duplex also competes successfully against these structures. Inclusion of a **denaturant** solubilizes the bases, thereby reducing the energetic cost of losing the stacking interactions of the duplex. As a consequence, the thermal stability of the duplex is reduced.

### 3. pH-Dependent Denaturation

The stability of nucleic acid duplexes is rather insensitive to changes in hydrogen ion concentration in the vicinity of neutral pH. This is because none of the functional groups present in typical nucleic acids titrate in this region. If the pH is decreased to a value below 5 or increased to a value above 9, significant destabilization of the nucleic acid duplex is observed. Because pH is simple to manipulate, high pH frequently is used to denature nucleic acid duplexes. Care must be taken to avoid degradation of the polymers under high pH conditions, which can promote depurination and strand breakage, particularly at high temperature. The presence of mismatches or modified bases may alter the pH sensitivity of the duplex because often these anomalous structures exhibit protonation equilibria in the vicinity of neutral pH. Certain nonnatural backbone structures alter radically the pH dependence of duplex stability. Low pH stabilizes cytosine-containing pyrimidine-purine-pyrimidine triple helices. Therefore, when sequence requirements for **triple helix** formation are met, care must be taken to adjust pH so as to avoid unintended triple helix formation.

#### Bibliography

1. L. A. Marky and K. J. Breslauer (1987) *Biopolymers* **26**, 1601.
2. L. Levine, J. A. Gordon, and W. P. Jenks (1963) *Biochemistry* **2**, 168.

#### Suggestion for Further Reading

3. V. A. Bloomfield, D. M. Crothers, and I. Tinoco (1974) *Physical Chemistry of Nucleic Acids*, Harper & Row, New York.

## Denaturation, Protein

The term protein denaturation usually refers to the process during which the specific properties of a native, functionally active [protein](#) disappear under some nonphysiological conditions without chemical modification of the molecule. Because the functional properties of a protein are provided by the unique disposition of all of its functional groups, the unique, native, three-dimensional [protein structure](#), denaturation of a protein assumes that the native structure is disrupted, and the [polypeptide chain](#) partially or completely unfolds, but without any change of the covalent bonds.

The native protein structure is maintained by a delicate balance of various types of interactions between the various protein groups and the surrounding media, which depend on the environmental conditions (see [Protein Stability](#)). *In vitro* this balance is usually optimal under conditions close to physiological. Departure from them can result in protein denaturation. Correspondingly, denaturation is usually classified in accordance with the external condition inducing it, for example, by heating (see Thermal denaturation, below), by cooling (Cold denaturation), by pressure (Pressure denaturation), by acids or alkali (pH denaturation), or by **denaturants** such as [urea](#) or [guanidinium salts](#), for example, guanidinium chloride (GdmCl) (Denaturant denaturation).

Denaturation can be characterized (1) by the extent to which the native protein structure is disrupted, the completeness of polypeptide chain unfolding (see: [Unfolded Proteins](#)); (2) by its reversibility upon restoring conditions optimal for the native state (see Reversibility, below); and (3) by its **cooperativity**. The completeness of disruption of the protein structure depends on the environmental conditions and also on the ability of the given protein to withstand the denaturing action of the environment (see [Protein Stability](#)). Studying all these aspects of protein denaturation is important for understanding the mechanism of **protein folding** and stabilization of the native protein structure, in other words, understanding the forces maintaining this structure. It is also important for the practical purpose of designing new proteins and modifying of natural proteins (see [Protein Engineering](#)). Correspondingly, this process has been extensively studied for many years using various proteins and their mutant forms, fragments and synthetic models in various solvents and by various experimental techniques. The results of these studies are summarized in numerous reviews (see [bibSect](#)). The following discussion assumes that the protein being studied is pure and no covalently modified variants are present in substantial quantity.

## 1. Reversibility of Denaturation

Under the appropriate conditions, the native state of a protein is usually restored when the denaturing conditions are removed. Thus the change of conformation of a protein upon denaturation, even complete unfolding of the polypeptide chain, in principle, is a reversible process determined by the environmental conditions. The folding of a polypeptide chain into a native conformation is a **thermodynamically** driven process in accordance with the genetic information included in the sequence of amino acid residues of the [polypeptide chain](#) (see [Self-Assembly](#)). This was realized first by Anson and Mirsky ([1](#), [2](#)), followed by Eisenberg and Schwert ([3](#)), and formulated as one of the most fundamental principles of protein science by Anfinsen ([4](#), [5](#)). That is why denaturation is usually reversible, so long as the newly synthesized protein has not been chemically modified *in vivo* (processed) after it has folded (see [Post-Translational Modifications](#)). However, some of the groups of proteins that are exposed by denaturation (unfolding), such as [thiol groups](#), are highly reactive and they become covalently modified. Others interact chaotically with each other and with the groups of other molecules to form aggregates. Furthermore, exposure of [nonpolar](#) groups upon unfolding decreases the solubility of proteins, thereby enhancing their aggregation. Consequently, the observation of complete reversibility of protein denaturation is not always straightforward, especially for large proteins that have high [effective molarities](#) of reactive and nonpolar groups in the unfolded state. Observation of reversibility often requires special conditions: (1) low concentrations of protein to decrease the probability of contacts between the protein molecules; (2) increase of electrostatic repulsions between the protein molecules by the appropriate pH and ionic strength of the solution; (3) exclusion from the solution of all other reactive molecules that might react irreversibly with protein groups (eg, oxygen, compounds containing thiol groups and [disulfide bonds](#), etc.); and (4) the use of co-solutes that increase the solubility of the unfolded proteins (eg, urea or GdmCl). The reversibility of unfolding diminishes at higher temperatures, especially above 70°C, and with prolonged incubation because of secondary temperature-induced chemical changes of protein groups, for example, reshuffling of disulfide bonds and [deamidation](#) of Asn and Gln residues ([6](#)). Renaturation (refolding) is also delayed by [cis-trans isomerization](#) of proline [peptide bonds](#) in the unfolded state ([7](#), [8](#)). This is especially pronounced if there are many proline residues and if in the folded state they are in the intrinsically unfavorable *cis* conformation ([9](#)). Searching for conditions under which a protein denatures reversibly is important for studying the thermodynamics of formation of protein structure, its energetic basis (see [Protein Stability](#)).

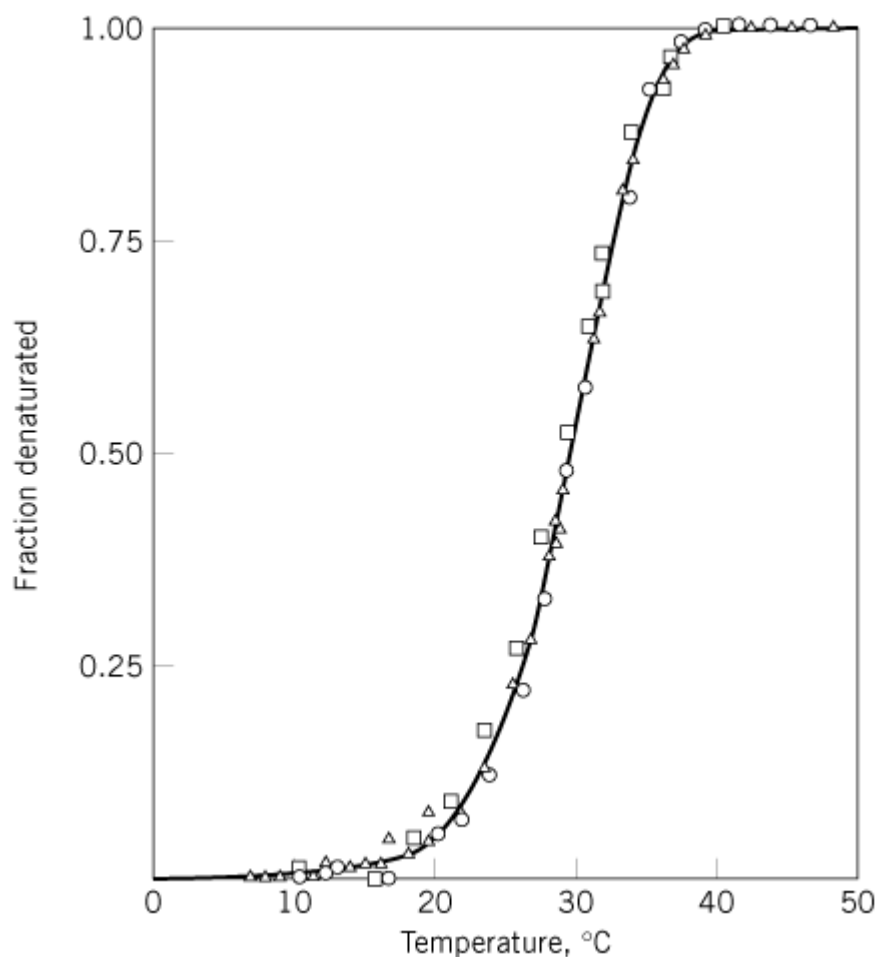
## 2. Cooperativity of Unfolding

The various properties of a protein molecule, observed by different methods upon varying the external conditions, change simultaneously in some cases, but not in others. Changes of different properties of a protein reflect different changes in its structure, so denaturation of protein structure can be either a single- or multistage process. Under denaturing conditions, the protein structure can change either as a whole or in several stages. Various parts of the molecule can be disrupted by

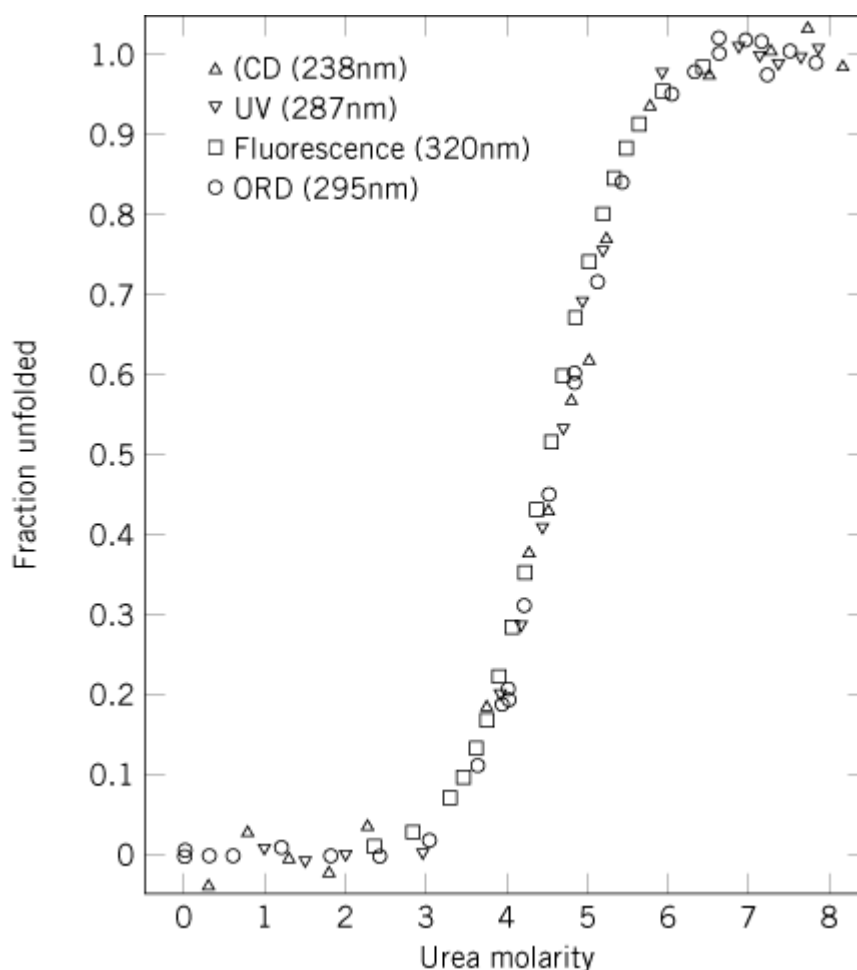
different denaturing conditions.

For small globular proteins, whose a molecular mass is less than about 20 kDa, all parameters that are sensitive to structural change usually change almost simultaneously with an increase of denaturing action (Figs. 1 and 2). Therefore, denaturation of these proteins is a single-step process in which the protein undergoes a single transition from the unique native state to the denatured state. Partly folded intermediate states probably occur during unfolding/folding, but they are highly unstable and would be present at equilibrium at only very low levels. Denaturation of such proteins is approximated rather well by a two-state transition that ignores all possible intermediate states. It is assumed that their native structure breaks down and refolds in an all-or-none manner (12, 13). Because proteins consist of many structural elements connected by covalent bonds into a polypeptide chain, this assumes that all the groups (amino acid residues) of small globular proteins change their states simultaneously, that is, cooperatively.

**Figure 1.** Thermal denaturation of ribonuclease A at pH 2.19 according to Ref. 10; (*f*) intrinsic viscosity, (○) optical rotation at 365 nm, and (△) absorption spectroscopy.

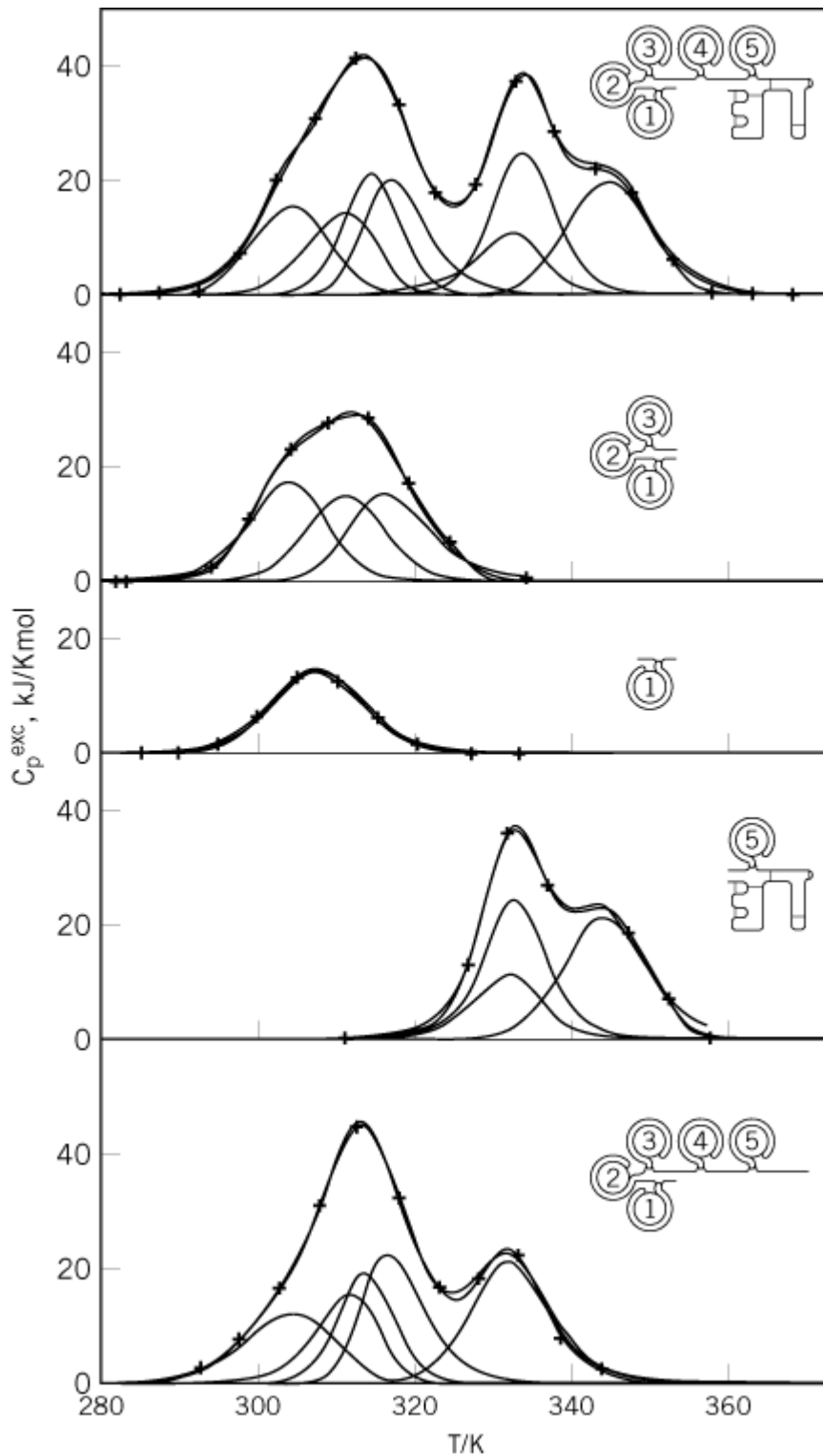


**Figure 2.** Urea denaturation of ribonuclease T<sub>1</sub> according to Thomson et al. (11). (△) circular dichroism ellipticity (CD) at 238 nm; (▽) UV absorbance at 287 nm; (○) optical rotatory dispersion (ORD) at 295 nm; and (*f*) fluorescence at 320 nm.



In contrast, multistage denaturation is normal for large proteins. Some smaller fragments of large proteins that separately preserve their native-like structure also denature reversibly and highly cooperatively (14-16). These observations have shown that the structures of large proteins are subdivided into cooperative units, called structural **domains**, which fold and unfold more or less independently ((17-19); for a review see Ref. 20). The stabilities of the different protein domains usually differ. So they unfold under different denaturing conditions, and the denaturation profile of such multidomain proteins is usually very complex (Fig. 3). Conversely, a complex denaturation profile usually means that the protein has a multidomain structure. It follows that the thermodynamic analysis of protein denaturation can provide information on the domain organization of large protein molecules (20-22).

**Figure 3.** Calorimetric characterization of the thermal unfolding of plasminogen and its proteolytically liberated domains, according to Ref. 16. The thin lines show the deconvolution of the excess heat absorption profiles into individual two-state transitions for the various domains of each of the molecules. The inserts show schematically the domain structures of the protein molecule used; complete plasminogen is at the top.



Simultaneity in changes of all properties of a protein upon denaturation is one of the tests for cooperativity of this process and, correspondingly, for the cooperativity of the protein structure. However, if the domains of the protein are similar, they denature under the same denaturing conditions, and all parameters that specify the state of the protein change simultaneously, giving the impression that the protein denatures in a single cooperative process. If such a transition is analyzed assuming that it is a two-state transition, fictitious thermodynamic parameters result that have no real meaning.

A stronger test of the cooperativity of the denaturation process consists of comparing the effective thermodynamic parameters derived from the equilibrium analysis whose real values are measured

by direct experimental methods. In thermodynamic analysis of the equilibrium between two states, the effective equilibrium constant  $K^{\text{eff}}$  can be expressed as:

$$K^{\text{eff}} = \frac{Y_x - Y_n}{Y_d - Y_n} = \exp\left(-\frac{\Delta G^{\text{eff}}}{RT}\right) \quad (1)$$

where  $Y_n$  and  $Y_d$  are values of any observed parameters that characterize the pure native and the pure denatured states, respectively,  $Y_x$  represents the value of the index under the given conditions, and  $\Delta G^{\text{eff}}$  is the effective Gibbs **free energy** difference between the two states. By following the dependence of the equilibrium constant on various fundamental intensive external parameters, such as temperature ( $T$ ), pressure ( $P$ ), or the activity of given ligand I ( $a_i$ ), one can find the effective extensive thermodynamic parameters that specify the process, the changes in **enthalpy**  $\Delta H^{\text{eff}}$ , **volume**  $\Delta V^{\text{eff}}$ , and the number of bound **ligands**  $\Delta N_i^{\text{eff}}$ .

$$\Delta H^{\text{eff}} = RT^2 \frac{\partial \ln K^{\text{eff}}}{\partial T} \quad (2)$$

$$\Delta V^{\text{eff}} = -RT \frac{\partial \ln K^{\text{eff}}}{\partial P} \quad (3)$$

$$\Delta N_i^{\text{eff}} = \frac{\partial \ln K^{\text{eff}}}{\partial a_i} \quad (4)$$

These effective parameters are derived assuming that the observed process represents a two-state transition, so if they correspond to the real values measured by direct experimental methods, this is a strong indication that the observed process indeed represents a two-state transition. The most efficient approach for this type of analysis is temperature-induced changes, particularly when they are observed by the associated heat effects measured by supersensitive scanning micro **calorimetry** (13). Comparing the calorimetrically measured enthalpy of the temperature-induced unfolding of a protein with the effective or van't Hoff enthalpy calculated from the sharpness of transition indicates how close the observed process is to a two-state transition.

More elaborate methods of analyzing of cooperativity, which reveal intermediate stages in the temperature-induced denaturation process consist of computer deconvolution of the calorimetrically measured excess **heat capacity** function of temperature (21, 22) or simulation of this function based on statistical mechanical modeling of the process (23, 24). Using these methods for thermodynamic analysis of the calorimetrically measured heat capacity function of a protein and of its proteolytic fragments corresponding to various domains under different solvent conditions yields quite detailed information about the domain organization of a protein. This is especially valuable for large proteins, the structures of which often cannot be studied by **NMR** because their large size or by **X-ray crystallography** methods because to their reluctance to crystallize (for a review, see Ref. 20). This approach has permitted investigating the domain organization of such large proteins as **fibrinogen** (15), **plasminogen** (16), **fibronectin** (25), **plasminogen activator** (26), recombinant factor XIII (27), and protein C (28) (see **Blood Clotting**).

NMR is a powerful method for analyzing the cooperativity of unfolding of small proteins and identifying any domains. Using the **chemical shifts** of individual proton resonances, one can follow the state of individual amino acid residues of a protein and their rate of **hydrogen exchange** as a function of the environmental conditions. Such analysis showed that small proteins unfold

cooperatively (29), but that some of the proteins reveal indications of subdomain organization under some conditions (30-32).

The mechanism of the extreme cooperativity of unfolding/folding of single-domain proteins is still a subject of discussion. Perhaps the most important conclusion, which is most commonly accepted, is that extreme cooperativity is a specific property of heterogeneous polypeptide chains that can fold into a unique compact structure in which all groups are tightly packed that is, it represents a lock and key system. For the theoretical models of cooperativity, see (23, 33-35).

### 2.1. Denaturant Denaturation

**Denaturants** are various agents (cosolutes) that lead to denaturation of proteins when they are added to aqueous solutions (Fig. 2). The most commonly used denaturants are **urea** and **guanidinium** chloride (GdmCl). Most proteins denature in concentrated solutions (eg, 5 to 8 M) of these compounds. Studies of the dependence of protein conformation on the concentration of these denaturants is the simplest and therefore the most popular method for quantitatively evaluating protein stability, the more so because denaturation by denaturants is also the most reversible of the various types of denaturation. Furthermore, the very significant changes in the hydrodynamic and optical properties of proteins in concentrated solutions of these denaturants led to the conclusion that they induce the most complete unfolding of polypeptide chains (see [Unfolded Proteins](#)). The denaturant-induced denaturation of small single-domain proteins usually occurs very sharply over a small change in denaturant concentration (Fig. 2), and it usually is regarded as a cooperative, two-state transition. If so, the equilibrium studies of the influence of denaturants on a protein can provide information on the energetics of protein unfolding/folding, that is, the stability of the native protein structure. Using any parameter  $Y$  sensitive to the protein conformation, the equilibrium constant and the Gibbs free energy of the transition can be estimated by, Eq. (1). Usually, the value of  $DG$  varies linearly with denaturant concentration  $[D]$ :

$$\Delta G = \Delta G^\circ - m[D] \quad (5)$$

and determines the value of  $DG$  at zero concentration  $DG^\circ$  of denaturant (36). The parameter  $m$  is a measure of the dependence of  $DG$  on the denaturant concentration. It is believed that the value of  $m$  correlates very strongly with the amount of protein surface exposed to solvent upon unfolding (37).

The main difference between the effects of urea and GdmCl on protein stability is that, although urea is likely to be only a destabilizing agent, GdmCl has a dual role. At high concentrations GdmCl destabilizes the protein native state, but at low concentrations it stabilizes it (38, 39) because, in contrast to urea, GdmCl is a salt. An increase in its concentration increases its activity as a denaturant and also the ionic strength of the solution, which reduces destabilizing [electrostatic interactions](#) within proteins. If the GdmCl concentration varies and the ionic strength of the solution is kept constant by adding of neutral salts, the stability of the protein is directly proportional to the concentration of GdmCl over an extended concentration range (40). Then, the stability of the protein can be determined from Eq. (5), and it corresponds well with the estimates of  $DG^\circ$  determined from calorimetric data (see [Protein Stability](#)). This equation is valid, however, only if the denaturant-induced unfolding of the protein is truly a two-state transition. Because this is not evident *a priori*, it needs to be verified.

Notwithstanding the extensive experimental data on the influence of denaturants on the state of proteins, the mechanism of this influence is yet to be understood. It is generally believed that one or both of two mechanisms are responsible for the ability of these compounds to destabilize the native state of proteins. In the first model, these solutes increase the solubility of the [nonpolar](#) side chains that form the compact cores of globular proteins (41-43). A sound argument for this hypothesis is the observation that denaturants increase the aqueous solubility of small nonpolar molecules. According to the other model, highly polar denaturants form [hydrogen bonds](#) with polar

groups of proteins, particularly the peptide groups, and increase the probability of their exposure. These hydrogen bonds would have to be stronger than those formed by [water](#) (44). A strong interaction of urea with the polar groups of proteins has been concluded by studies of the influence of these cosolutes on the solubility of model compounds (45). According to calorimetric studies, urea and GdmCl interact with proteins in both the folded and unfolded states and cause a strong heat effect that is proportional to the exposed surface area and depends strongly on temperature (46, 47). It became evident that the denaturing action of urea and GdmCl is caused mainly by the enthalpic binding of these compounds. With increased temperature, the bound denaturants dissociate gradually absorbing heat that results in an apparent heat capacity increase of the protein in the presence of urea and GdmCl. This explains the observed decrease of the intrinsic **viscosity** of the unfolded polypeptide chain with increasing temperature (48). Analyzing calorimetric data by the standard model for multisite **ligand binding** was determined the number of denaturant binding sites and the binding constant (47). Timasheff (49) examined these calorimetric data in terms of [preferential binding](#). He concluded that the weak overall thermodynamic interaction of the denaturants with proteins results from the strong favorable interactions of proteins with denaturants at the exchangeable sites that are compensated for by unfavorable interactions, almost as strong, at loci on the protein surface at which the denaturants cannot displace water. Schellman & Gassner (50) have analyzed calorimetric data for the interactions of urea and GdmCl with proteins based on the solvent-exchange model. According to this model, the free energy and preferential interactions for very weak binding and high concentrations of cosolvent do not depend on the high concentrations of cosolvent, but on its excess occupation over that of the bulk solution. The binding constant  $k$  and the number of binding sites  $n$  derived by using this model differ considerably from those derived using the standard binding model, but the values of the binding enthalpy  $DH$ , and of  $nkDH$  are rather close.

The specific interaction of urea with polar groups of proteins has been shown by X-ray crystallography using crystals of diketopiperazines (51) and proteins (52). The number of binding sites for urea and GdmCl were determined crystallographically for **ribonuclease A** and [dihydrofolate reductase](#), and they closely correspond with the numbers estimated from calorimetric data by the standard binding model. Furthermore, binding of the denaturant significantly decreases the crystallographic **B temperature factor** of the protein, that is, the mobility of its native state (53). Each molecule of urea or GdmCl must interact simultaneously with several groups of the protein by forming multiple hydrogen bonds (54). If so, one might expect these denaturants to form even more extensive networks of hydrogen bonds with the groups of the unfolded polypeptide chain, increasing its rigidity and consequently its intrinsic viscosity. Then, the conformation of polypeptide chains in the presence of urea or GdmCl could not be regarded as a standard [random coil](#) because it would be a stretched random coil, at least at room temperature.

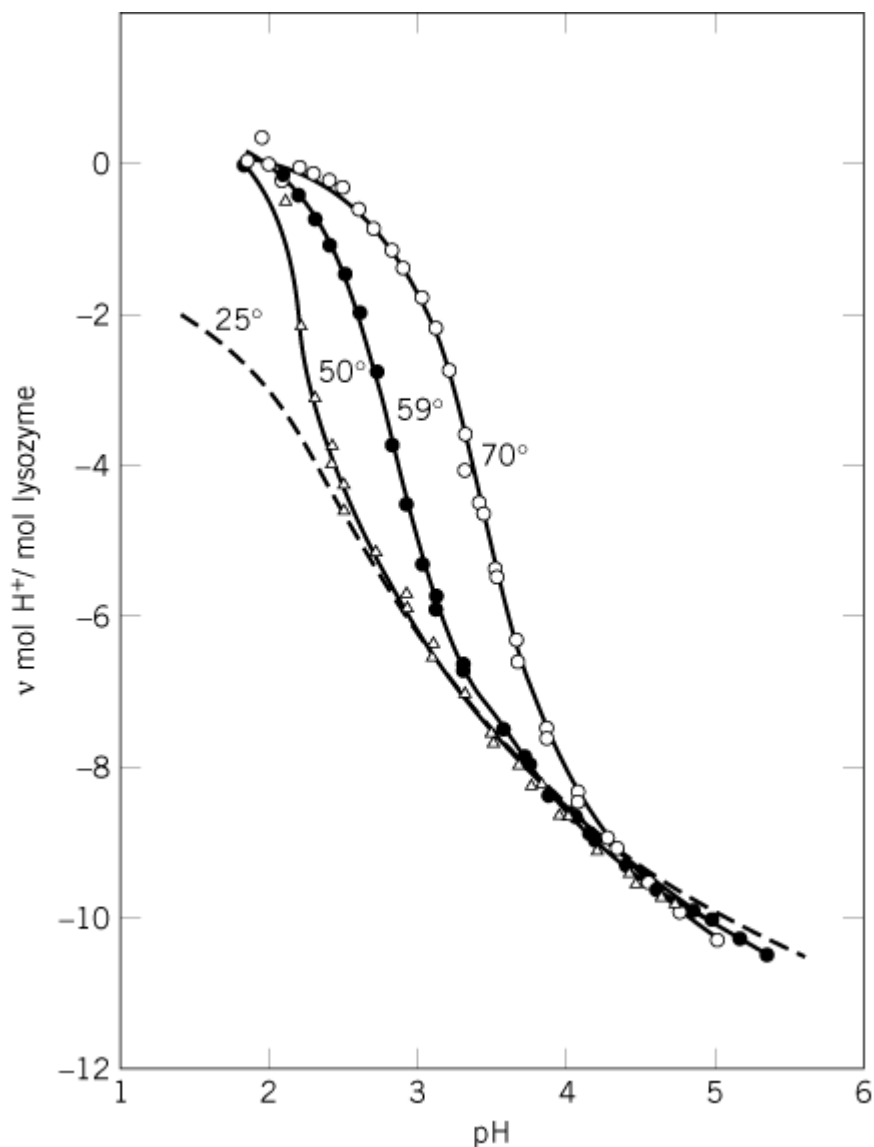
## 2.2. pH Denaturation

The pH of the solution is one of the most important factors that determines the state of a protein. The changes of protein structure with variation of pH are frequently complex. For of small, single-domain proteins, the intensive, denaturational changes in all parameters that characterize protein conformation usually occur over a narrow range of pH and are clearly distinguishable above a smooth background of gradual changes. In accordance with the Le Chatelier principle, pH-induced denaturation is associated with the release or uptake of protons (Fig. 4). Potentiometric titration of proteins revealed that smooth changes result from titrating of protein groups that have  $pK_a$  values not very different from those of the corresponding free amino acids, whereas the gross conformational changes associated with pH denaturation are accompanied primarily by the unmasking of buried groups (56, 57).

**Figure 4.** Protonation state of hen [lysozyme](#) versus pH at different temperatures, according to Ref. 55. At 25°C the protein is native at all pH values, but at higher temperatures it unfolds with decreasing pH, unmasking and



protonating the buried groups.



pH-induced denaturation depends on temperature and, therefore, temperature-induced denaturation depends on the pH (see Thermal denaturation). Because the temperature-induced denaturation of a small single-domain protein is a two-state transition, it was assumed that the transition induced by variation of pH at a fixed temperature is also two-state, so that the initial and the final characteristics of the protein are independent of the sequence of pH or temperature variation. Then, one can determine the Gibbs free energy difference between the native and denatured state from the potentiometrically measured pH-dependence of the protonation difference  $\Delta v(\text{pH})$  between these states (55):

$$\Delta G(\text{pH}) = 2.3RT \int_{\text{pH}_t}^{\text{pH}} \Delta v(\text{pH}) d\text{pH} \quad (6)$$

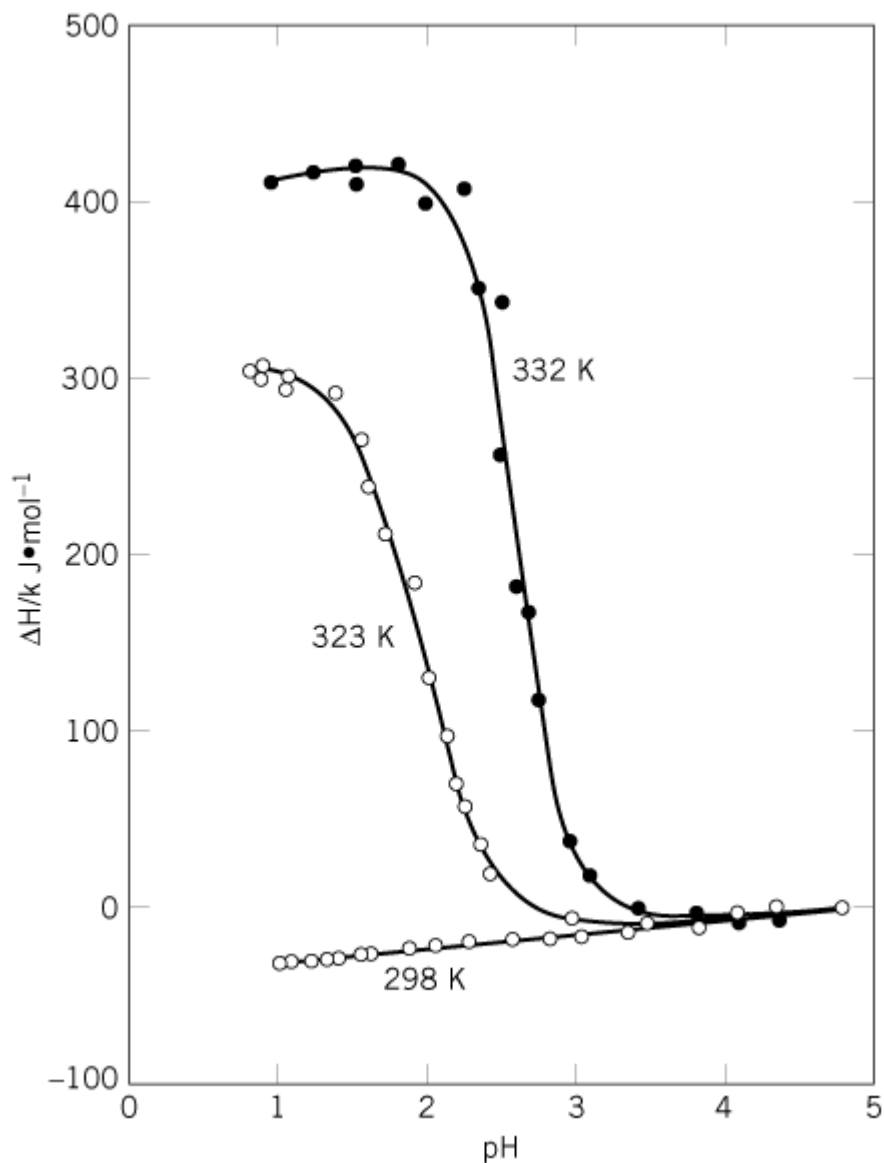
where  $\text{pH}_t$  is the transition pH at which half of the molecules are in the denatured state.

Comparison of Gibbs free energy values of pH- and temperature-induced denaturation showed that they indeed correspond very well, a strong argument that the pH-induced denaturation of small globular proteins can be regarded as a two-state transition (58).

Calorimetric titration studies have shown that pH denaturation of proteins creates significant heat

effects that depend on temperature (Fig. 5). If corrected for the heats of protonation, the enthalpy of pH-induced protein transitions is identical to the enthalpy of the temperature-induced transition and depends similarly on temperature (55). Thus, the pH-induced denaturation of protein results in a heat capacity increase similar to that observed for temperature induced-denaturation. For the theory of acid-induced denaturation, see Refs. 56, 57, and 59 and for a review, see Ref. 44.

**Figure 5.** Heat of transfer of lysozyme from pH 4.8 to a given pH, at constant temperature, according to Ref. 55. At 298 K (25°C), the protein remains folded at all pH values, but it unfolds at low pH at the higher temperatures used.

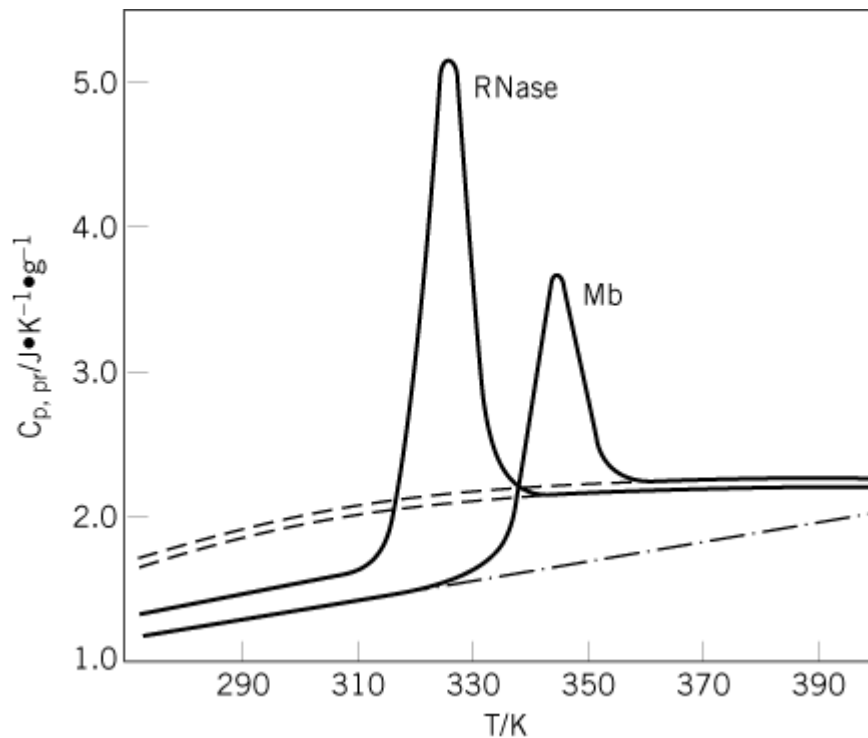


With increased temperature increasing thermal motion induces disruption of the ordered native structure of a protein. Breakdown of the native structure of small globular proteins occurs over a relatively narrow temperature range, in which all specific characteristics of the native protein change simultaneously (Fig. 1). In that temperature range the microscopic fluctuations coalesce into a global macroscopic fluctuation, that is, the macroscopic unfolding of a protein structure (60), which results in drastic changes in the rate of hydrogen exchange of the protein's amide groups (61). Therefore, the thermal denaturation of a single-domain small protein is a highly

cooperative process. Because this process is induced by increased temperature, it is obvious that it should result in an increase in the enthalpy, which is the extensive thermodynamic parameter conjugate with temperature, that is, it should proceed with heat absorption in accordance with the Le Chatelier's principle. Assuming that this process does indeed represent a two-state transition, one can estimate the enthalpy of denaturation by using Eq. (2), which is known as the van't Hoff equation .

The heat effect of denaturation can be observed directly by scanning microcalorimetry as an intensive heat absorption peak (Fig. 6). For small single-domain proteins, the calorimetrically measured enthalpy of denaturation agrees well the van't Hoff enthalpy, showing that thermal denaturation of these proteins is well approximated by a two-state transition (13). Deviation between the calorimetric and van't Hoff enthalpies does not usually exceed 5%, which means that the population of intermediate states is rather low. This was confirmed by detailed analysis of the excess heat absorption profile (63) and has been shown subsequently for many small single-domain proteins.

**Figure 6.** Temperature-dependence of the partial molar heat capacity of pancreatic ribonuclease A (RNase) and of sperm whale met-myoglobin (Mb) in solution at pH 4.4. The dashed lines show the heat capacities of the unfolded states, extrapolated to lower temperatures according to Eq. 9 of the text. The other line is a linear extrapolation of the heat capacity of the native state of myoglobin. Reproduced from Ref. 62.



For a two-state transition at the mid-transition temperature  $T_t$ , where  $K = 1$ ,

$$-RT_t \ln K = \Delta G(T_t) = \Delta H(T_t) - T_t \Delta S(T_t) = 0 \quad (7)$$

For the entropy of transition,

$$\Delta S(T_t) = \frac{\Delta H(T_t)}{T_t} \quad (8)$$

The most remarkable feature of the thermal denaturation of proteins is that the partial heat capacity of the denatured protein is always significantly higher by 25 to 50% than that of the native state (13, 64). The value of the denaturation heat capacity increment is specific for a given protein, depends on the temperature, and decreases at higher temperatures (65-67). Several explanations for the denaturation heat capacity increment can be suggested: (1) gradual melting of residual structure in the protein with continued increase in temperature; (2) increase of configurational freedom of the polypeptide chain upon disrupting of the rigid compact native structure of the protein molecule; or (3) [hydration](#) of the groups that are exposed to water by protein unfolding. The contribution of residual structure to the denaturation heat capacity increment is not significant because the heat capacity increment observed upon thermal denaturation is almost the same as in denaturation by denaturants, if the effects of denaturant solvation are taken into account properly (46), whereas it appears that proteins denatured by denaturants (urea or GdmCl) do not generally have residual structure. According to theoretical estimates (68), the contribution of the increased configurational freedom to the observed denaturation increment of the protein heat capacity cannot be large. Thus, the main reason for the higher partial heat capacity of the protein denatured state is that in this state more groups, particularly the nonpolar groups, are exposed to solvent, that is water. It is known that transfer of nonpolar groups into water increases the heat capacity (69-72). Exposure of the polar groups to water has the opposite effect. It decreases the heat capacity (66, 72-76). Detailed calorimetric study of the transfer of model compounds into water has also revealed that aromatic groups are not identical in hydration properties to nonpolar groups but are somewhere in between the nonpolar and polar groups (75). Thus, there are three main contributors to the heat capacity increment upon protein unfolding: exposure of the nonpolar, aromatic, and polar groups that were buried in the folded conformation. Each of them contributes in proportion to their water-[accessible surface](#) areas exposed upon unfolding ( $DA_{np}$ ,  $DA_{arom}$ , and  $DA_{pol}$ , respectively). At 25°C the heat capacity increment of protein unfolding can be calculated as follows:

$$C_p^{hyd}(25^\circ\text{C}) = -(2.14\Delta A_{np} + 1.55\Delta A_{arom} - 1.27\Delta A_{pol}) \times JK^{-1} \text{ mol}^{-1} \quad (9)$$

The coefficients in this equation depend on temperature [for details, see (77)]. Ignoring this dependence and the usual identification of the aromatic groups with nonpolar groups is not justified in quantitatively analyzing thermodynamic properties of proteins (78). Upon protein unfolding, the surface area of the nonpolar plus aromatic groups exposed is greater than that of the polar groups, so unfolding of a protein results in an overall increase of the partial heat capacity (Fig. 6).

The important consequence of the denaturation heat capacity increment is that the enthalpy and entropy of denaturation are temperature-dependent functions. Indeed, because

$\frac{\partial \Delta H}{\partial T} = \Delta C_p$  and  $\frac{\partial \Delta S}{\partial T} = \frac{\Delta C_p}{T}$  for the enthalpy and entropy,

$$\Delta H(T)_{pH} = \Delta H(T_t)_{pH} + \int_{T_t}^T \Delta C_p(T) dT \quad (10)$$

and

$$\Delta S(T)_{\text{pH}} = \frac{\Delta H(T_t)_{\text{pH}}}{T_t} + \int_{T_t}^T + \frac{\Delta C_p(T)}{T} dT \quad (11)$$

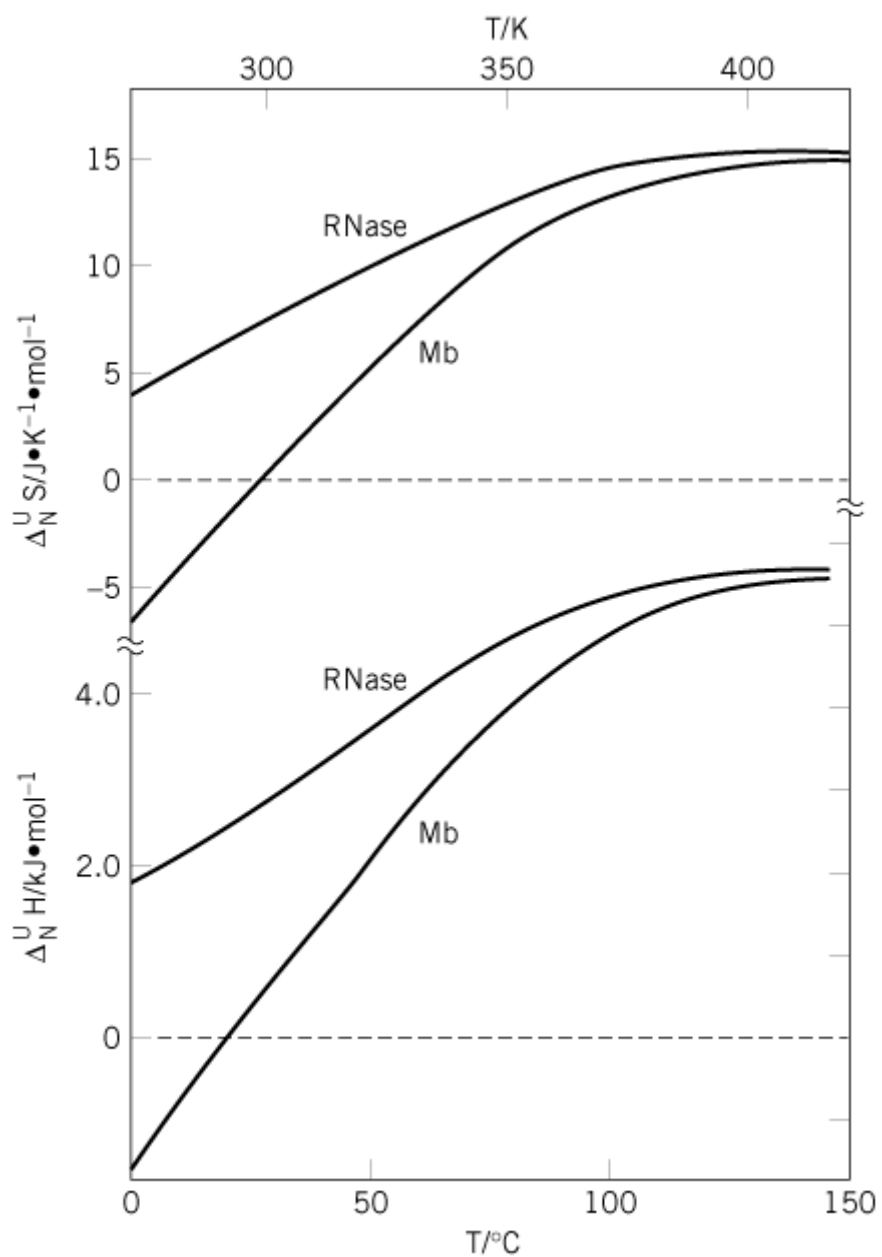
Because  $\Delta C_p$  is positive and decreases with increasing temperature, the enthalpy and entropy of protein denaturation increase asymptotically to some constant level (78, 83).

#### 2.4. Cold Denaturation

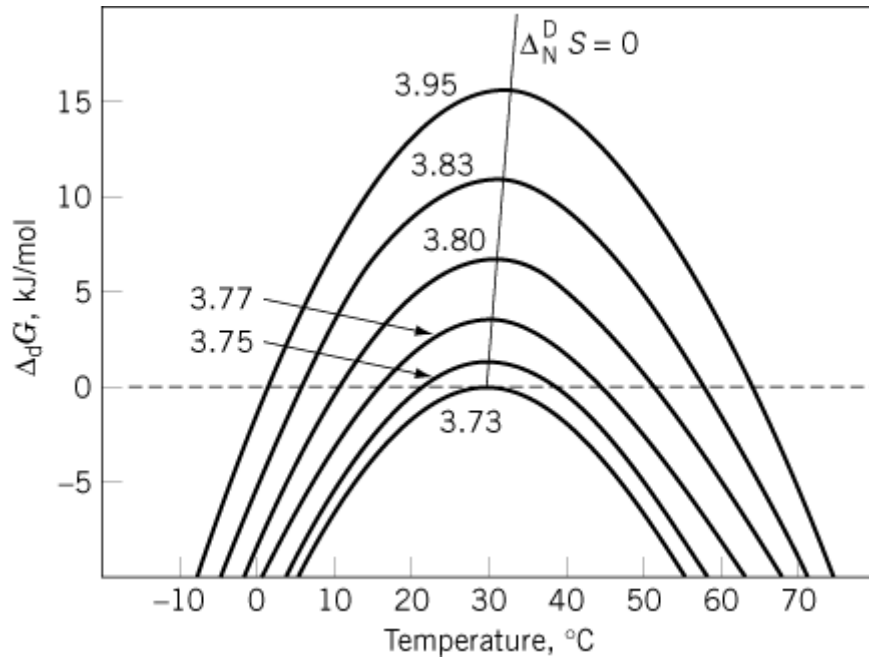
Interest in the influence of low temperatures on the stability of native protein structure began with a paper published by Hopkins in 1930 (79). Studying the precipitation of [ovalbumin](#) upon denaturation by concentrated urea at various temperatures, he noticed that the rate of denaturation was greater at 0°C than at 23°C. Similar effects were observed later by others with other globular proteins, showing that in many cases the stability of a protein to urea or GdmCl denaturation is maximal around 30°C and decreases at lower temperatures. Meanwhile, evidence accumulated showing that many enzymes in aqueous solution, without urea or GdmCl, are inactivated at low temperature (about 0°C) but their activity was restored at room temperature (for a review, see Refs. 80 and 81). Most of these proteins are **oligomeric**, and it is likely that their cold inactivation resulted from their dissociation into subunits, as was found for several supermolecular structures, such as [microtubules](#) (82, 83). The [quaternary structure](#) of proteins is less stable than the tertiary, and its stability decreases with decreasing temperature.

The dual sensitivity of proteins to both high and low temperatures was formally explained when it was realized that protein unfolding results in a significant heat capacity increase and consequently that the enthalpy and entropy of unfolding are temperature-dependent and decrease with a decrease in temperature (Fig. 7). If so, according to the thermodynamic formalism, the Gibbs free energy difference between the unfolded and folded states of the protein should be a monotonic function of temperature but should have an extremum at the temperature at which the entropy difference between the two states decreases to zero (see [Protein Stability](#)). If the entropy and enthalpy functions continue to decrease below that point, with a change of sign in entropy, one might expect that the Gibbs free energy difference should become zero at some low temperature (Fig. 8). At this temperature, the protein should denature, that is unfold, but, in contrast to thermal denaturation at high temperature, it should proceed with a release of heat and a decrease of the enthalpy and entropy, that is, an increase in order. This paradoxical conclusion was confirmed by microcalorimetric measurements of the heat effects resulting from cooling aqueous solutions of proteins. It was shown that proteins do indeed denature upon cooling, release heat (Fig. 9), and change other properties (viscosity, [circular dichroism](#) NMR spectra) indicating that the polypeptide chain unfolds (84-86). Studies of various proteins under different conditions showed that cold denaturation is a very general property of globular proteins that depends on the environmental conditions and on their specific structure (81, 87). The temperature of cold denaturation is highest for proteins that change most in heat capacity upon unfolding and can be raised significantly by appropriate choice of solvent conditions (pH, ionic strength, buffer, presence of denaturants), which facilitate observation of this process.

**Figure 7.** Temperature-dependence of the enthalpy and entropy differences between the native and denatured states of ribonuclease A (RNase) and of sperm whale met-myoglobin (Mb) per mole of amino acid residues under solvent conditions that provide the maximum stabilities of these proteins. Reproduced from Ref. 62.



**Figure 8.** The Gibbs free energy difference between the denatured and native states of metmyoglobin in acetate buffer solution at the pH values indicated on the curves. Under these conditions, the protein cold denatures within the temperature range 0° to 30°C. Each curve is at a maximum when the entropy of unfolding is zero. Reproduced from Ref. [84](#).



From the formal point of view, the cold denaturation of proteins is caused by the fact that the heat capacity of the unfolded protein is greater than that of the folded protein. The increased heat capacity upon protein unfolding is usually explained by the hydration of nonpolar groups (see [Protein Stability](#)), which also results in a decrease of entropy ( $DS^{\text{hyd}}$ ). The latter is considered to be the main cause of the [hydrophobic effect](#), which as an entropic force is proportional to absolute temperature ( $TDS^{\text{hyd}}$ ) and should decrease upon lowering the temperature, thus destabilizing the protein. However, the real situation is not so simple. Indeed, because of the hydration heat capacity effect, the negative entropy of hydration of nonpolar groups is intrinsically temperature-dependent and increases in magnitude with decreasing temperature, thus balancing to some extent the entropic contribution of nonpolar groups to the stabilization of the folded protein's compact state. On the other hand, the unfolding of a protein results in exposing and hydrating the nonpolar groups and also polar and aromatic groups. The effects of hydrating of these groups contributes negatively to stabilizing protein structure, and the magnitude increases with decreasing temperature ([77](#)).

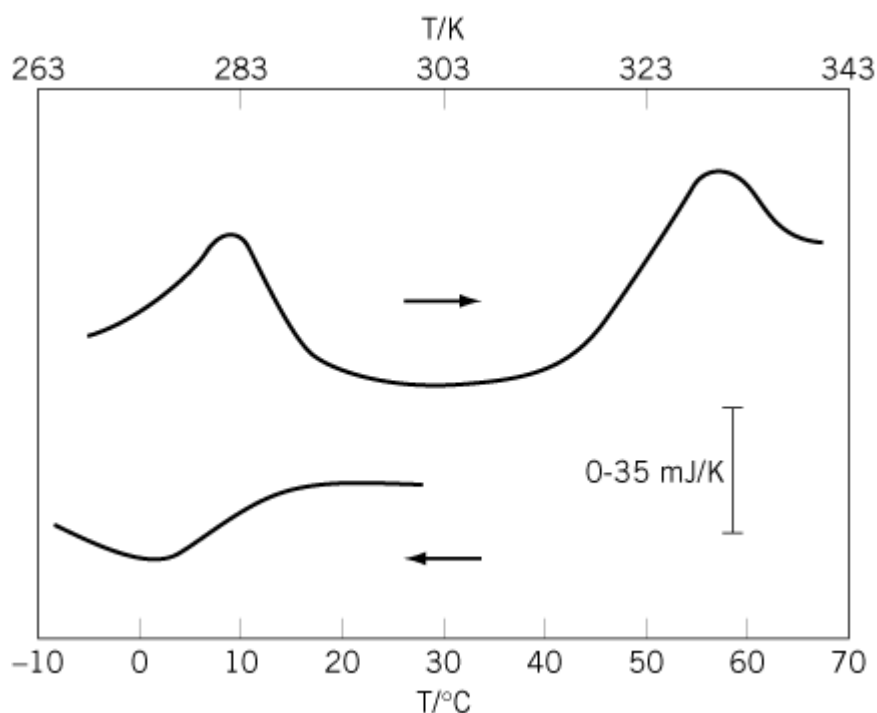
## 2.5. Pressure Denaturation

Pressure is one of the fundamental physical parameters that determine the state of any system, in particular, biological systems. Living organisms habituate at very different pressure conditions, from less than 1 bar in high mountains to thousands of bars in deep water. Therefore, it is natural that studies of pressure effects on proteins began a long time ago. In 1914, the first experiments done by Bridgmen showed that high pressure induces protein denaturation. But quantitative investigations of the effect of pressure on proteins began much later, in the middle 1960s, when it became clear that this is the only practical way to acquire information about the **volume** effects of protein folding (for reviews see Refs. [88-92](#)).

Even early studies of the influence of pressure on proteins revealed that proteins are not very sensitive to pressure, and only at extremely high pressures do they exhibit changes like those observed in temperature and pH denaturation. However, according to more recent studies, particularly by high-pressure NMR spectroscopy, pressure-denatured proteins retain more residual structure than when denatured by heat, pH, or denaturants ([86](#), [90](#), [93](#)). The pressure-induced denaturation of proteins takes place over a relatively narrow pressure interval,

depending on the temperature and solvent conditions, particularly the pH (Fig. 10). At the same time, the temperature and pH of denaturation also depend on pressure. Thus, these three parameters are interdependent, and variation of any one of them at fixed values of the others can lead to disruption of the native protein structure. Because denaturation of small, single domain proteins by temperature and pH is a highly cooperative process that is well approximated by a two-state transition, it was concluded that pressure-induced denaturation also represents a transition between two states, the native and denatured (94). This conclusion was confirmed by using an **electrophoretic** technique at high pressure (95), but it was challenged by later NMR studies (86, 93). In particular cases of multidomain proteins, the different domains undergo independent structural changes at specific well-defined limiting values of pressure (96). The immediate effect of pressure on oligomeric proteins is dissociation into subunits, followed by unfolding (80, 89, 90).

**Figure 9.** Microcalorimetric recordings of the heat effects upon cooling metmyoglobin under conditions where it cold denatures (bottom curve) and subsequent heating (top curve) that demonstrates first the refolding of the cold-denatured state and then its thermal unfolding at high temperatures. Reproduced from Ref. 84.



Assuming that pressure-induced denaturation represents a two-state transition, one can calculate the effective change of molecular volume upon unfolding, bearing in mind that according to the Le Chatelier principle, a pressure-induced shift in equilibrium should result in a decrease in volume that can be determined by Eq. (3). The great advantage of this approach to studying volume effects is that it does not require supersensitive volumetric techniques but can be realized by using any physical parameter (eg, **spectrophotometry**) that is sensitive to the conformational state of the protein and can indicate the pressure-induced change in equilibrium between the native and denatured states. This is important because proteins denature only at very high pressures, usually above 2 kbar.

It follows from the inefficiency of pressure in inducing protein denaturation that the volume



change  $ΔV_d$  upon denaturation is very small, usually less than 0.5% of the total protein volume. For example, the volume decrease in denaturation of ribonuclease A is  $-(30 ± 10)$  ml/mol at an average pressure of about 1 kbar, and it decreases by  $-4$  to  $-5$  ml/mol when the temperature increases from 25° to 50°C (97-99). Zipp and Kauzmann (100) discussed the parallelism between the thermodynamic factors governing the process of reversible denaturation of proteins and the dissolution of hydrocarbon molecules in water. If the interior of a protein resembles a liquid or solid hydrocarbon and if significant numbers of nonpolar groups are exposed to water on unfolding, one would expect a significant volume decrease upon unfolding at low pressure. With increasing pressure, the volume reduction should become less negative and should change sign at 2kbar. However, this does not happen in protein denaturation. Therefore, the simple hydrophobic model that considers globular proteins as an oil drop or a solid hydrocarbon (101) does not work. In contrast to liquid or solid hydrocarbons, the native protein interior is densely packed, but it has some cavities and loose parts. Then, increasing hydrostatic pressure forces water molecules into the protein interior, gradually filling cavities and eventually breaking up the protein structure (90). This would explain why pressure-denatured proteins retain elements of structural organization to a much greater extent than proteins denatured by other means. But this also questions how closely pressure-induced unfolding is approximated by a two-state transition. In turn, this questions the reliability of the volume effects of protein folding that are determined by pressure denaturation.

#### Bibliography

1. M. L. Anson and A. F. Mirsky (1934) *J. Gen. Physiol.* **17**, 393–398.
2. M. L. Anson (1945) *Adv. Protein Chem.* **2**, 361–386.
3. M. A. Eisenberg and G. W. Schwert (1951) *J. Gen. Physiol.* **34**, 583–606.
4. C. B. Anfinsen (1956) *J. Biol. Chem.* **221**, 405–412.
5. C. B. Anfinsen (1973) *Science* **181**, 223–230.
6. T. J. Ahern and A. M. Klibanov (1988) *Methods Biochem. Anal.* **33**, 91–127.
7. J. F. Brandts et al. (1975) *Biochemistry* **14**, 4953–4963.
8. R. L. Baldwin (1989) *Trends Biochem. Sci.* **14**, 291–294.
9. F. X. Schmid (1992) In *Protein Folding* (T.E. Creighton, ed.), Freeman, New York, pp. 197–241.
10. A. Ginsburg and W. R. Carroll (1965) *Biochemistry* **4**, 2159.
11. J. A. Thomson, B. A. Shirley, G. R. Grimsley, and C. N. Pace (1989) *J. Biol. Chem.* **264**, 11614–11620.
12. R. Lumry, R. Biltonen, and J. F. Brandts (1966) *Biopolymers* **4**, 917–944.
13. P. L. Privalov and N. N. Khechinashvili (1974) *J. Mol. Biol.* **86**, 665–684.
14. M. E. Goldberg (1969) *J. Mol. Biol.* **46**, 441–446.
15. P. L. Privalov and L. V. Medved (1982) *J. Mol. Biol.* **159**, 665–683.
16. V. V. Novokhatny, S. A. Kudinov, and P. L. Privalov (1984) *J. Mol. Biol.* **179**, 215–232.
17. D. B. Wetlaufer (1981) *Adv. Protein Chem.* **34**, 61–92.
18. M. G. Rossman and P. Argos (1981) *Annu. Rev. Biochem.* **50**, 497–532.
19. J.-R. Garel (1992) In *Protein Folding* (T. E. Creighton, ed.), Freeman New York, pp. 405–454.
20. P. L. Privalov (1982) *Adv. Protein Chem.* **35**, 1–104.

21. P. L. Privalov and S. A. Potekhin (1986) *Methods Enzymol.* **131**, 4–51.
22. E. Freire (1995) *Methods Enzymol.* **259**, 144–168.
23. K. P. Murphy and E. Freire (1992) *Adv. Protein Chem.* **41**, 313–361.
24. Y. V. Griko, E. Freire, G. P. Privalov, H. Van Dael, and P. L. Privalov (1995) *J. Mol. Biol.* **252**, 447–459.
25. L. V. Tatunashvili et al. (1990) *J. Mol. Biol.* **211**, 161–169.
26. V. V. Novokhatny, K. C. Ingham, and L. V. Medved (1991) *J. Biol. Chem.* **266**, 12994–13002.
27. I. V. Kurochkin et al. (1995) *J. Mol. Biol.* **248**, 414–430.
28. L. V. Medved (1995) *J. Biol. Chem.* **270**, 13652–13658.
29. C. C. McDonald, W. D. Phillips, and J. D. Glickson (1971) *J. Am. Chem. Soc.* **93**, 235–246.
30. F. M. Hughson, P. E. Wright, and R. L. Baldwin (1991) *Science* **249**, 1544–1548.
31. B. A. Schulman et al. (1995) *J. Mol. Biol.* **253**, 651–657.
32. S. E. Radford and C. M. Dobson (1995) *Phil. Trans. R. Soc. Lond.* **B 348**, 17–25.
33. K. A. Dill and D. Stigter (1995) *Adv. Protein Chem.* **46**, 59–104.
34. E. Freire and K. P. Murphy (1991) *J. Mol. Biol.* **222**, 687–698.
35. V. I. Abkevich, A. M. Gutin, and E. I. Shakhnovich (1995) *J. Mol. Biol.* **252**, 460–471.
36. C. N. Pace (1990) *Trends Biotech.* **8**, 93–98.
37. J. K. Myers, C. N. Pace, and J. M. Scholtz (1995) *Protein Sci.* **4**, 2138–2148.
38. L. M. Mayer and F. X. Schmid (1993) *Biochemistry* **32**, 7994–7998.
39. Y. Hagihara, S. Aimoto, A. Fink, and Y. Goto (1993) *J. Mol. Biol.* **231**, 180–184.
40. M. M. Santoro and D. W. Bolen (1992) *Biochemistry* **31**, 4901–4907.
41. D. B. Wetlauff et al. (1981) *Adv. Protein Chem.* **34**, 61–92.
42. M. Roseman and W. P. Jencks (1975) *J. Am. Chem. Soc.* **97**, 631–640.
43. R. Breslow and T. Gui (1990) *Proc. Natl. Acad. Sci. USA* **87**, 167–169.
44. C. Tanford (1970) *Adv. Protein Chem.* **24**, 1–95.
45. D. R. Robinson and W. P. Jencks (1965) *J. Am. Chem. Soc.* **87**, 2462–2470.
46. W. Pfeil and P. L. Privalov (1976) *Biophys. Chem.* **4**, 33–40.
47. G. I. Makhatadze and P. L. Privalov (1992) *J. Mol. Biol.* **226**, 491–505.
48. F. Ahmad and A. Salahuddin (1974) *Biochemistry* **13**, 245–249.
49. S. N. Timasheff (1992) *Biochemistry* **31**, 9857–9864.
50. J. A. Schellman and N. C. Gassner (1996) *Biophys. Chem.* **39**, 259–275.
51. M. M. Thayer et al. (1992) *Biophys. Chem.* **46**, 165–169.
52. A. C. W. Pike and K. R. Acharya (1994) *Protein Sci.* **3**, 706–710.
53. J. Dunbar et al. (1997) *Protein Sci.* **6**, 1727–1733.
54. J. A. Gordon and W. P. Jencks (1963) *Biochemistry* **2**, 47–57.
55. W. Pfeil and P. L. Privalov (1976a) *Biophys. Chem.* **4**, 23–32.
56. J. T. Edsall and J. Wyman (1958) *Biophysical Chemistry*, Academic Press, New York.
57. J. Wyman (1964) *Adv. Protein Chem.* **19**, 223–286.
58. W. Pfeil and P. L. Privalov (1976c) *Biophys. Chem.* **4**, 41–50.

59. D. O. V. Alonso and K. A. Dill (1991) *Biochemistry* **30**, 5974–6985.
60. A. Ikegami (1977) *Biophys. Chem.* **6**, 117–130.
61. M. Nakanishi, M. Tsuboi, and A. Ikegami (1973) *J. Mol. Biol.* **75**, 673–682.
62. P. L. Privalov (1992) In *Protein Folding* (T. E. Creighton, ed.), Freeman and New York, pp. 83–126.
63. E. Freire and R. L. Biltonen (1978) *Biopolymers* **17**, 463–479.
64. J. F. Brandts (1964) *J. Am. Chem. Soc.* **86**, 4291–4301.
65. P. L. Privalov et al. (1989) *J. Mol. Biol.* **205**, 737–750.
66. P. L. Privalov and G. I. Makhatadze (1990) *J. Mol. Biol.* **213**, 385–391.
67. G. I. Makhatadze and P. L. Privalov (1996) In *Physical Properties of Polymers* (J. E. Mark, ed.), AIP Press, Woodbury, N.Y., pp. 91–100.
68. J. M. Sturtevant (1977) *Proc. Natl. Acad. Sci. USA* **74**, 2236–2240.
69. S. J. Gill, S. F. Dec, G. Olafsson, and I. Wadso (1985) *J. Phys. Chem.* **89**, 3758–3761.
70. O. L. Privalov and S. J. Gill (1988) *Adv. Protein Chem.* **39**, 191–234.
71. J. R. Livingston, R. S. Spolar, and M. T. Record (1991) *Biochemistry* **30**, 4237–4244.
72. G. I. Makhatadze and P. L. Privalov (1990) *J. Mol. Biol.* **213**, 375–384.
73. K. P. Murphy and S. J. Gill (1991) *J. Mol. Biol.* **222**, 699–709.
74. R. S. Spolar, J. R. Livingston, and M. T. Record (1992) *Biochemistry* **31**, 3947–3955.
75. P. L. Privalov and G. I. Makhatadze (1992) *J. Mol. Biol.* **224**, 715–723.
76. H. Gomez, V. J. Hisler, D. Xie, and E. Freire (1995) *Proteins: Struct. Function Gen.* **22**, 404–412.
77. G. I. Makhatadze and P. L. Privalov (1995) *Adv. Protein Chem.* **47**, 307–425.
78. D. R. Robertson and K. P. Murphy (1997) *Chem. Rev.* **97**, 1251–1267.
79. P. G. Hopkins (1930) *Nature* **126**, 383–386.
80. R. Jaenicke (1981) *Ann. Rev. Biophys. Bioeng.* **10**, 1–67.
81. P. L. Privalov (1990) *CRC Crit. Rev. Biochem. Mol. Biol.* **25**, 281–305.
82. S. N. Timasheff (1979) *Trends Biochem. Sci.* **4**, 61–65.
83. P. Dustin (1984) *Microtubules*, Springer-Verlag, Berlin.
84. P. L. Privalov, Y. V. Griko, S. Y. Venyaminov, and V. P. Kutysenko (1986) *J. Mol. Biol.* **190**, 487–498.
85. Y. V. Griko, P. L. Privalov, S. Y. Venyaminov, and V. P. Kutysenko (1988) *J. Mol. Biol.* **202**, 127–138.
86. J. Zhang, X. Peng, A. Jonas, and J. Jonas (1995) *Biochemistry* **34**, 8631–8641.
87. F. Franks (1995) *Adv. Protein Chem.* **46**, 105–139.
88. K. Heremans (1982) *Annu. Rev. Biophys. Bioeng.* **11**, 1–21.
89. G. Weber and H. G. Drickamer (1983) *Q. Rev. Biophys.* **16**, 89–112.
90. J. L. Silva and G. Weber (1993) *Annu. Rev. Phys. Chem.* **44**, 89–113.
91. C. Royer (1995) *Methods Enzymol.* **259**, 357–377.
92. V. V. Mozhaev et al. (1996) *Proteins: Struct. Function Genet.* **24**, 81–91.
93. J. Jonas and A. Jonas (1994) *Annu. Rev. Biophys. Biomol. Struct.* **23**, 287–318.
94. J. Hermans Jr. and G. Acampora (1967) *J. Am. Chem. Soc.* **89**, 1547–1552.
95. S. A. Hawley and R. M. Mitchell (1975) *Biochemistry* **4**, 1277–1281.

96. H. G. Drickamer (1977) *High Temperatures High Pressures* **9**, 505–506.
97. S. J. Gill and R. L. Glogowsky (1965) *J. Phys. Chem.* **69**, 1515–1519.
98. D. J. F. Brandts et al. (1970) *J. Am. Chem. Soc.* **89**, 4826–4838.
99. S. A. Hawley (1971) *Biochemistry* **10**, 2436–2441.
100. A. Zipp and W. Kauzmann (1973) *Biochemistry* **12**, 4217–4228.
101. M. H. Klapper (1973) *Prog. Bioorg. Chem.* **2**, 55–132.

### Suggestions for Further Reading

102. R. L. Baldwin (1991) Molten globules: Specific or nonspecific folding intermediate? *Chemtracts-Biochem. Mol. Biol.* **2**, 379–389.
103. T. E. Creighton, ed. (1992) *Protein Folding*, Freeman, New York.
104. T. E. Creighton (1995) An unfolding story. *Curr. Biol.* **5**, 353–356.
105. R. Jaenicke (1991) Protein folding: Local structures, domains, subunits, and assemblies. *Biochemistry* **30**, 3147–3161.
106. M. N. Jones (1979) *Biochemical Thermodynamics*, Elsevier, Amsterdam-Oxford-New York.
107. C. N. Pace (1990) Conformational stability of globular proteins. *Trends Biochem. Sci* **15**, 14–17.
108. P. L. Privalov (1979) Stability of proteins. Small globular proteins. *Adv. Protein. Chem.* **33**, 167–241.
109. P. L. Privalov (1982) Stability of proteins. Proteins which do not present a single cooperative system. *Adv. Protein Chem.* **35**, 1–104.
110. J. L. Silva and G. Weber (1993) Pressure stability of proteins. *Annu. Rev. Phys. Chem.* **44**, 89–113.
111. C. Tanford (1968) Protein denaturation. *Adv. Protein Chem.* **23**, 121–275; (1970) **24**, 1–95.

## Dendrotoxins

The venom of the *Dendroaspsis mamba* snakes contains several [toxins](#), termed dendrotoxins, which bind with nanomolar affinities to voltage-dependent **potassium channels** and facilitate neurotransmitter release because they prolong membrane depolarization. Their polypeptide chains consist of 57 to 60 amino acid residues, with three disulfide bonds and folding closely similar to that of **BPTI**-type inhibitors of serine proteinases. A short amino-terminal  $3_{10}$ -helix is followed by a two-stranded  $\beta$ -sheet and a short COOH-terminal  $\alpha$ -helix ([1](#)). Dendrotoxins provide yet another example of a stable protein module that has been readapted during the course of **evolution** to perform a different biological function.

### Bibliography

1. J.-M. Lancelin et al. (1994) *Struc. Biol.* **1**, 246–250.

## Density Gradient Centrifugation

Many [centrifugation](#) separations simply centrifuge an initially homogeneous sample and physically remove any heavy and large particles or molecules from the solution, leaving them as a pellet on the bottom of the centrifugation tube. In contrast, density gradient separations employ a gradient of varying, high concentrations of a small molecule, such as sucrose or cesium chloride, distributed along the axis of the centrifugally generated force. The lowest concentration is at the solution meniscus and the highest at the base. With sufficiently high concentrations of the small molecule, the density and viscosity of the solution also vary along the gradient. The density gradient is produced directly before centrifugation, or it is generated by the centrifugal force itself. A wide variety of materials is used to form gradients (Table 1), depending on which components are to be separated.

**Table 1. Gradient Materials Used for Density Gradient Centrifugation<sup>a</sup>**

| <b>Material and Typical Use</b>   | <b>Maximum Aqueous Density<br/>(g/cm<sup>3</sup>)</b>   |
|---|---|
| Cesium salts  |   |
| Acetate: Isopycnic separation of proteoglycans, DNA, RNA, and viruses             | 2.00 at 20°C  |
| Bromide: Separation and fractionation of glycoproteins                            | 1.72 at 20°C  |
| Chloride: Most widely used density gradient material for isopycnic separations    | 1.91 at 20°C  |
| Formate: DNA hydration  | 2.10 at 20°C  |
| Oxalate: Isopycnic separation of proteoglycans, DNA, RNA                          | 2.01 at 20°C  |
| Ficoll: Isopycnic and rate-zonal separation of cells and subcellular fractions    | 1.17 at 20°C sucrose polymer  |
| Glycerol: Rate-zonal separation of RNA, proteins and enzymes sensitive to sucrose | 1.26 at 20°C  |
| Percoll: Isopycnic separation of cells and large subcellular components           | aqueous solution of polyvinylpyrrolidone-coated colloidal silica particles ~20 nm in diameter |
| Potassium salts   |   |
| Bromide: Isopycnic centrifugation of lipoproteins and density-labeled proteins    | 1.37 at 20°C  |
| Iodide: Isopycnic centrifugation of DNA   | 1.72 at 24°C  |

|  |              |
|--|--------------|
| and RNA  |              |
| Tartrate: Isopycnic centrifugation of viruses  | 1.49 at 20°C |
| Rubidium salts   |              |
| Chloride: Isopycnic banding of proteins  | 1.49 at 20°C |
| Trichloroacetate: Isopycnic DNA banding  | 1.90 at 20°C |
| Sodium salts   |              |
| Bromide: Rate-zonal and isopycnic separation of lipoproteins   | 1.53 at 20°C |
| Chloride: Rate zonal separation of DNA and lipoprotein fractionation   | 1.20 at 20°C |
| Iodide: Isopycnic centrifugation of DNA and RNA  | 1.90 at 20°C |
| Sucrose: Rate-zonal separation of DNA, RNA, subcellular fragments, proteins, and viruses; isopycnic banding of subcellular particles and viruses | 1.3 at 20°C  |

<sup>a</sup> Adapted from Spinpro Software from Beckman Instruments, a useful software tool. Unfortunately it still runs only under MS-DOS on Intel processor-based computer.

When a sample of molecules is subjected to centrifugation in a density gradient, they sediment toward the bottom if they have a greater [buoyant density](#) than the solution, or they will float toward the top if they have a lower density. The rates of sedimentation through the gradient of the various molecules can be measured, as in [sedimentation velocity centrifugation](#), and depends on their buoyant densities and [sedimentation coefficients \(s-values\)](#). Alternatively, the centrifugation can be continued until each molecule reaches that point in the gradient when the solution density is the same as its buoyant density. At this point the molecule neither sediments nor floats. Then different molecules are separated on the basis of differences in their buoyant densities, which can be measured in this way. The physical process of banding significantly improves the resolution of this technique.

## 1. Rate-Zonal Centrifugation

Separations of macromolecules or subcellular fractions with different sedimentation rates and/or sizes are accomplished by layering the sample of material as a narrow zone onto the top of a preformed gradient of the appropriate material (Table 1) and subjecting the sample to centrifugation. During centrifugation, the constituent components of the sample proceed through the gradient as individual zones, each at its specific sedimentation rate. The sedimentation process results in separating the sample components into discrete bands or zones that can be recovered from the centrifuge tube, typically by simple fraction collection (below). The shape of the gradient, that is, the change in solution density as a function of the physical length of the centrifuge tube or distance from the center of rotation, is a function of the gradient material and its concentration, the rate of centrifugation, etc. Five principal types of gradient systems used in rate-zonal separations are summarized in Table 2

**Table 2. Principle Types of Gradients in Rate-Zonal Centrifugation<sup>a</sup>**

| Type                  | Gradient Characteristics as a Function of Distance Along Centrifuge Tube   |
|-----------------------|--|
| Linear                | Solution density or gradient-forming-solute concentration versus distance is linear  |
| Isokinetic            | Solution density or gradient-forming-solute concentration versus distance is usually convex but sometimes linear; all particles of the same density sediment at a constant rate at all distances from the center of rotation |
| Step or discontinuous | Solution densities or gradient-forming-solute concentrations versus distance form discrete zones with (relatively) large differences in density (concentration) from one to another  |
| Shallow               | Solution density or gradient-forming-solute concentration changes gradually versus distance  |
| Sharp or steep        | Solution density or gradient-forming-solute concentration changes rapidly versus distance  |

<sup>a</sup> Adapted from Spinpro Software from Beckman Instruments.

Approximate values of the sedimentation coefficients (*s*-values) of macromolecules and subcellular components can be obtained from rate-zonal gradient centrifugation without purifying the protein. This is accomplished in isokinetic gradients, where particles sediment at a constant rate, as in sucrose gradient centrifugation. To determine the *s*-value, it is necessary to include sedimentation “markers” of macromolecules with known sedimentation coefficients. It is necessary only to be able to differentiate between the various molecules, perhaps by their **spectroscopic** properties or by their biological activities.

## 2. Flotation in Gradient Centrifugation

This technique is a subclass or offshoot of gradient techniques. It occurs in centrifugation when the buoyant densities of the centrifuged macromolecule or subcellular components are less than that of the solution in which they are centrifuged. The rates of flotation of the individual solution components to the top of the solution are a function of their sizes, shapes and densities, just as in sedimentation to the bottom of the tube.

## 3. Isopycnic Gradient Centrifugation

Isopycnic procedures are used to separate macromolecules on the basis of their buoyant densities. The density gradient is chosen so as to include the densities of the species to be separated from each other. The sample to be sedimented may be prepared either by homogeneously distributing it throughout the gradient or by layering it on top of the gradient. In either case, when centrifugal force is applied to the sample, the sedimenting components move along the density gradient to where they have the same buoyant density. If they are at a lower solution density initially, they tend to sediment, whereas they tend to float when at a higher solution density. When the density gradient is maintained by the centrifugal force, and perhaps even generated by it, such as with cesium chloride (Table 1), the gradient is “self-generating.” When all the molecules present are at their equilibrium positions, the centrifugal force condenses the molecules into sharp zones and maintains them. Unlike rate-zonal centrifugation, where the zones of the various molecules are constantly broadened by diffusion, but like [isoelectric focusing](#), isopycnic density gradient centrifugation condenses the sample into a narrow zone and keeps it stable and sharp in time. This results from two opposing force fields:

diffusion, which tends to dissipate the zone, and the centrifugal field, which forces any “escaping” molecule back into its appropriate zone.

If the components to be separated have small differences in buoyant densities, the gradient used should bracket the desired density range and should also be shallow, so as to afford significant separation between components. Alternatively, if the components to be separated have large differences in buoyant densities, the gradient can be adjusted so that any undesired component either floats to the top of the gradient or sediments to the bottom, whereas the desired components are within the middle of the gradient.

#### 4. Gradients of Viscous Materials

When viscous gradient materials, such as sucrose, Ficoll, Percoll, and glycerol, are used, the gradients are often prepared by layering steps of decreasing density in the centrifuge tube. Then the various steps are allowed to diffuse into each other (usually for 1 to 3 hours at room temperature) to form an approximately linear gradient. Alternatively, the sample is layered and centrifuged immediately through the series of discontinuous step zones; the sample often accumulates at an interface between two steps, so the steps are chosen so that different materials accumulate at different interfaces. Often a cushioning layer of high density gradient material is incorporated into the bottom of the tube to prevent the desired sample from pelleting onto the tube.

#### 5. Salt Gradients

Cesium chloride is sufficiently heavy that it forms a substantial density gradient under the sole influence of rapid centrifugation. In this case, the starting solution can be uniform in both the cesium chloride and the sample. During centrifugation, the density gradient of cesium chloride is generated, and the molecules of the sample migrate to their isopycnic positions. The appropriate initial cesium chloride concentration must be chosen to give the desired density gradient with the centrifugal force used. It is possible to shorten the time required for attaining equilibrium by performing the gradient, which then requires centrifugation runs of only 2 to 8 hours. Often a cushioning layer of high density gradient material is incorporated to prevent complete pelleting of the sample on the bottom of the tube.

#### 6. Fractionating Gradients After Centrifugation

The density gradients employed aid the process of recovering the samples, minimizing the disturbance of banded zones in a centrifuge tube. Methods for recovering the banded samples depend on whether the tube in which the experiment has been performed is reusable or disposable. In the former case, the tube is held rigidly while a flat-end syringe needle of appropriate length is carefully inserted at the center of the tube down through the length of the column of gradient solution to the tube base. Then the contents are slowly pumped out. In the latter case, the bottom of the tube is punctured with a needle, and the tube contents are permitted to flow out slowly. Then they are collected in a fraction collector, perhaps after passing through a UV **absorbance**, **fluorescence**, or other detector. An alternative fractionation device for this process, a centrifuge tube slicer, is available from Beckman Instruments.

Once separated into individual fractions, the gradient samples can be analyzed for the molecules of interest, such as their enzymatic activity, spectral properties, or a **radiolabel**. The exact shape of the density gradient at the end of the centrifugation is determined by measuring the refractive index of the solution.

#### 7. Centrifuge Rotors for Gradient Density Separations

Three types of rotor are used for density gradient centrifugation, varying in the orientation of the sample during centrifugation. In the swinging bucket rotor, the sample swings outward away from



the center of rotation during centrifugation and extends at a right angle to the axis of rotation. Consequently, the force applied to the column of liquid containing the sample is always along the long axis of the liquid column, and the sedimenting sample moves along the linear axis of the applied force. Swinging bucket rotors are often used with sucrose density centrifugation. In fixed-angle rotors, the liquid column is tilted at a fixed angle relative to the axis of rotation. Thus, the centrifugal force is applied at an angle to the long axis of the tube and differentially at the cell meniscus and cell base. Fixed-angle rotors are routinely used for pelleting samples and for heavier isopycnic salt gradients, where the gradient can reorient after the centrifugation run. In vertical tube rotors, the liquid column is held parallel to the axis of rotation, and the centrifugal force is applied perpendicularly to the long axis of the tube. Because the sedimentation path is short, simply the width of the tube, these rotors permit very short run times, but with good resolution, and are useful for isopycnic separations. The advantages and disadvantages of each rotor type vary significantly with the particular application.

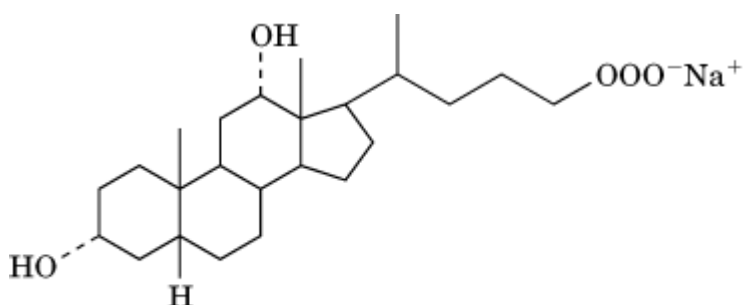
#### Suggestion for Further Reading

R. Hinton and M. Dobrota (1980) *Laboratory Techniques in Biochemistry and Molecular Biology: Density Gradient Centrifugation* (T. S. Work and E. Work, eds.), North-Holland, New York. A somewhat dated but useful compendium of techniques for preparative density gradient centrifugation.

## Deoxycholate

Sodium deoxycholate, a bile acid salt, is a strongly **denaturing**, anionic **detergent** that is structurally close to sodium cholate. The 5-OH group of cholate is replaced with a hydrogen atom in deoxycholate (Fig. 1). This detergent has been used to solubilize many **membrane-anchored** proteins, such as those that are membrane-linked via a phosphatidylinositol-containing glycolipid (termed **GPI-anchored** proteins) (1, 2). Although it has rarely been shown that sodium deoxycholate solubilizes integral **membrane proteins**, certain insoluble, extracellular proteins, such as human lens membrane **proteinase**, have been efficiently solubilized by using sodium deoxycholate (3). Sodium deoxycholate has a relatively high **critical micelle concentration** (cmc) of ~10 mM, which makes it easy to be removed by **dialysis** from the solubilized proteins and lipids. Although this detergent is strongly denaturing toward integral membrane proteins, it solubilizes membrane lipids efficiently (4) and also stimulates **G-protein-coupling** to an effector, such as **phospholipase C** (5).

**Figure 1.** The structure of sodium deoxycholate.



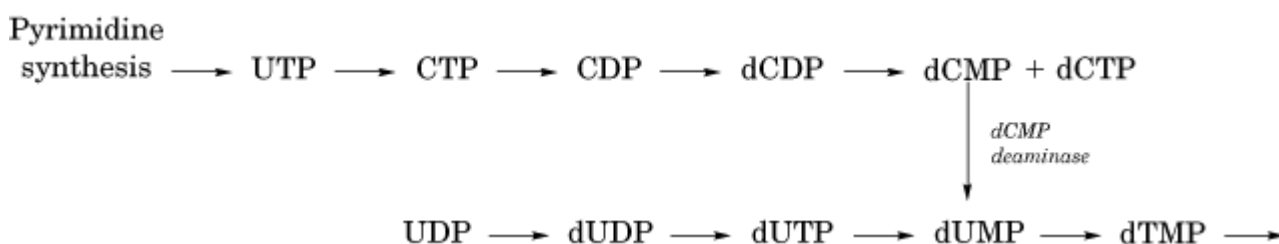
Sodium deoxycholate

## Bibliography

1. D. W. Loe, J. R. Glover, S. Head, and F. J. Sharom (1989) *Biochem. Cell Biol.* **67**, 214–223.
2. N. M. Hooper and A. J. Turner (1988) *Biochem. J.* **250**, 865–869.
3. O. P. Srivatsava and K. Srivatsava (1989) *Exp. Eye Res.* **48**, 161–175.
4. P. Banerjee J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
5. H. Nakanishi, Y. Takeyama, H. Ohyanagi, Y. Saitoh, and Y. Takai (1990) *Biochem. Biophys. Res. Commun.* **170**, 111–118.

## Deoxycytidylate Deaminase

Deoxycytidylate deaminase (dCMP deaminase) is one of two **allosteric** enzymes in the *de novo* biosynthetic pathways to **deoxyribonucleoside triphosphates** (dNTPs); the other allosteric enzyme is [ribonucleotide reductase](#). Inasmuch as dCMP deaminase sits at a branch point in the synthesis of pyrimidine dNTPs, the importance of this control in maintaining balanced dNTP pools can be readily appreciated.



It is assumed, although not yet established with certainty, that dCMP, dCDP, and dCTP are in ready equilibrium with one another via the actions of dCMP kinase and nucleoside diphosphate kinase. Clearly, dCMP is a precursor of dCTP, and the accumulation of dCTP can be a signal to the cell to divert dCMP, *via* dCMP deaminase, to pools of thymidine nucleotides. At the same time, dTTP accumulation can signal the metabolic machinery to shut off its own synthesis, and dCMP deaminase makes an appropriate regulatory target.

Scheme 1 shows two distinct pathways to dUMP, one involving dCMP deaminase and the other involving the reduction of UDP by ribonucleotide reductase with subsequent action of nucleoside diphosphate kinase and [deoxyuridine triphosphatase](#) (dUTPase). Radioisotope labeling experiments in mammalian cells show that the deaminase-based pathway is far more significant than the second pathway as a source of dTMP residues in DNA (1).

### 1. Structure and Regulation of dCMP Deaminase

Deoxycytidylate deaminase was first described in the late 1950s and early 1960s in papers by the Maleys in the United States (2, 3) and by Scarano and colleagues in Italy (4). As shown above, the enzyme catalyzes the hydrolytic deamination of deoxycytidine monophosphate to deoxyuridine

monophosphate, the immediate precursor to thymidine nucleotides via [thymidylate synthase](#). As an enzyme involved in the synthesis of a specific DNA precursor, the activity of dCMP deaminase is closely correlated with the proliferative activity of the cell (2). In a series of papers in the 1960s, both the Maley and Scarano laboratories showed that vertebrate dCMP deaminase is inhibited by dTTP and activated by dCTP (3, 4). dCTP was found to be almost completely required for activity, whereas dTTP is a strong inhibitor, with inhibition constant ( $K_i$ ) values in the micromolar range.

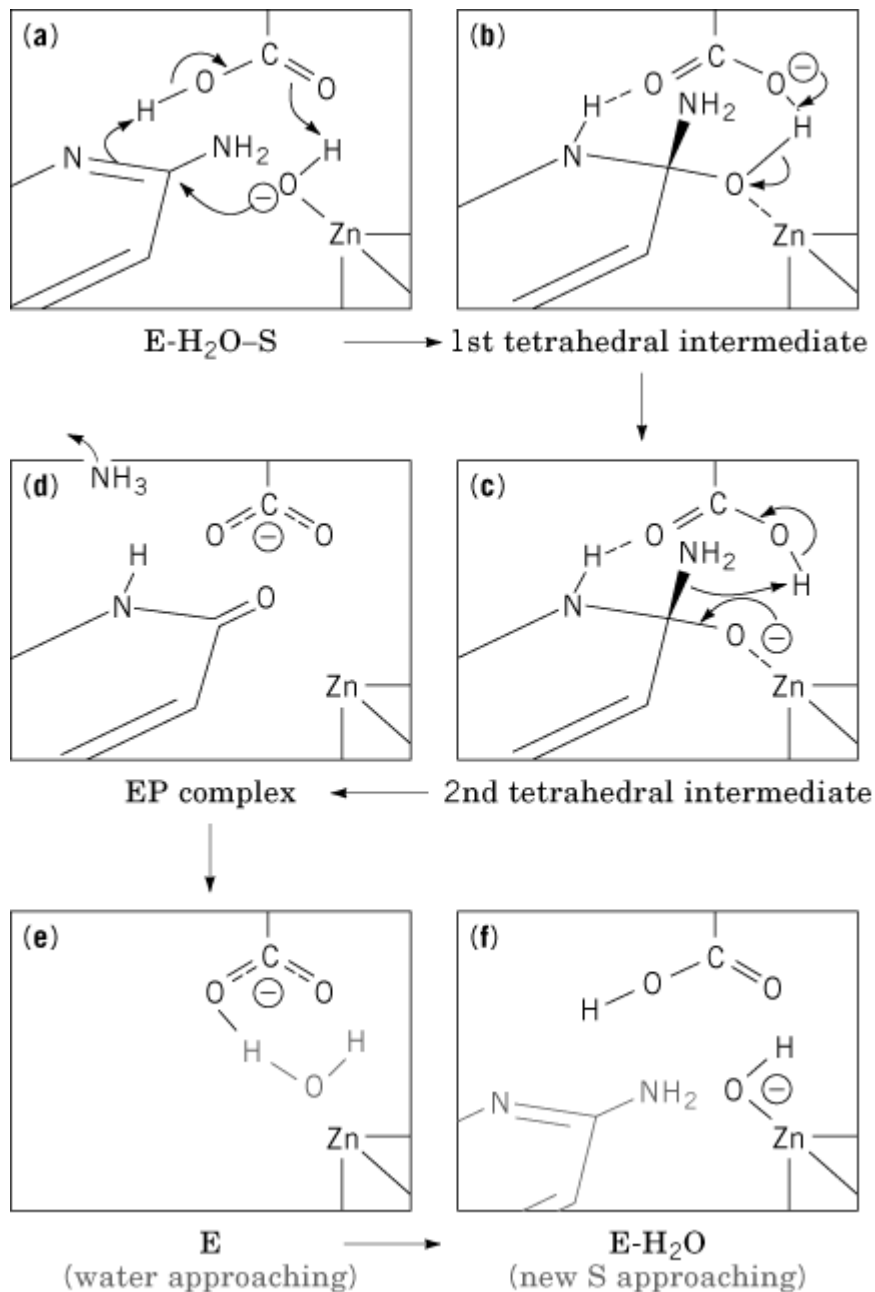
Reversal of dTTP inhibition by dCTP was shown to vary with the dCTP concentration in a highly sigmoidal manner as did the variation of the enzyme velocity with the concentration of dCMP as substrate (see [Allostery](#)). Activation involved an apparent change in subunit structure of the enzyme, for the dCTP-activated enzyme from chick embryo was shown to have a [sedimentation coefficient](#) ( $s_{20,w}$ ) of 6.78 S whereas the dTTP-inhibited enzyme had  $s_{20,w}$  of only 3.63 S (3). Conformational changes were also involved, as shown in part by large [circular dichroism](#) changes resulting from dCTP binding.

In early work, the subunit structure of dCMP deaminase was more readily defined in studies with T-even bacteriophage-infected bacteria, a much more abundant source than vertebrate tissues (5, 6). The phage T2 enzyme was shown to be a homo-hexamer of about 120,000 molecular weight, and the phage T4 and T6 enzymes are similar. The effects of dCTP and dTTP were shown to be similar to those observed with the vertebrate enzyme. However, dCTP is not a physiological regulator in these systems. T-even phage DNA contains 5-hydroxymethyl-cytosine substituted for cytosine, and the dCTP pool in infected cells is completely replaced by 5-hydroxymethyl-dCTP. This latter nucleotide was shown to be equivalent to dCTP in its effect on phage dCMP deaminases, and it is undoubtedly the true *in vivo* positive regulator of the phage enzymes. Like other phage-coded enzymes of deoxyribonucleotide synthesis, T-even phage-coded dCMP deaminase functions as part of a multienzyme complex that evidently facilitates the biosynthesis of dNTPs and possibly their incorporation into DNA (see [Thymidylate Synthase](#)). Incidentally, the host organism for T-even phages, *Escherichia coli*, does not contain dCMP deaminase. Instead, the conversion of deoxycytidine to deoxyuridine nucleotides occurs at the triphosphate level (7) through the action of a structurally unrelated enzyme, as follows:



With [recombinant DNA](#) approaches, large-scale isolation of dCMP deaminase has become possible, and structural studies are under way with the T4, yeast, and human enzymes (8, 9). The three dCMP deaminases show significant sequence [homology](#) at the C-terminus, but the yeast subunit is far larger than the T4 or human polypeptide chains. Both the T4 and human enzymes contain catalytically essential zinc (see [Zinc-Binding Proteins](#)). This fact, plus an apparent [active-site](#) homology with adenosine deaminase and cytidine deaminase, suggest a mechanism comparable to that of cytidine deaminase, whose structure has been determined (10; see Fig. 1). This mechanism is consistent with the potent inhibition of dCMP deaminase by tetrahydrodeoxyuridine monophosphate, a presumed [transition state analogue](#).

**Figure 1.** A mechanism for the cytidine deaminase reaction, based on structural analysis of the *E. coli* enzyme. The mechanism for dCMP deaminase is probably similar or identical. Reprinted with permission from L. Betts, S. Xiang, S. A. Short, R. Wolfenden, and C. W. Carter Jr. (1994) *J. Mol. Biol.* **235**, 635–656.



## 2. dCMP Deaminase is the Product of a Mutator Gene

A [mutator gene](#) is one in which the action of a mutant allele contributes toward an increase in the spontaneous [mutation](#) rate at all loci; conversely, the action of the wild-type allele contributes toward genetic stability. A classic example is the *E. coli mutH*, *mutL*, and *mutS* genes, in any of which a mutation raises the mutation rate by compromising DNA [mismatch repair](#).

In the same sense, the gene encoding dCMP deaminase is a mutator gene, as shown first in mouse S49 cells ([11](#)), where a subline evidently deficient in dCMP deaminase activity caused accumulation of dCTP and a deficiency of dTTP consistent with the loss of a major pathway to dTTP synthesis, and an increase in spontaneous mutation rates of at least an order of magnitude. Similar results were observed in yeast ([12](#)) and in *E. coli* infected with a dCMP deaminase deletion mutant of phage T4 ([13](#)). In these latter studies, additional information was gained about the mutational events occurring as a result of the dNTP pool imbalance. Infection by the T4 deaminase deletion led to dramatic accumulations of 5-hydroxymethyl-dCTP and specific increases in mutations that occur along an AT

→ GC [transition mutation](#) pathway, as expected if mutation results from misinsertion of the abundant 5-hydroxymethyl-dCMP opposite template A. On the other hand, most of the mutations in yeast that were stimulated by dCMP deaminase deficiency were GC → CG [transversion mutations](#), an event that does not follow in any logical fashion from the dNTP pool imbalance demonstrated in these cells.

Clearly, understanding mutagenesis induced by DNA precursor pool imbalance involves factors other than concentration-dependent competition between nucleotides that are correctly and incorrectly base-paired to a template base. Adding to the complexity is the existence of a dCMP deaminase-defective hamster V79 cell line in which a large increase in the [dCTP]/[dTTP] pool ratio has only a negligible effect on the spontaneous mutation rate (14). Based on analysis of synchronized Chinese hamster ovary cells (15), it is possible that most of the expanded dCTP pool in these cells is located outside the nucleus, where it cannot directly influence replication fidelity. In any event, it is apparent from the several systems that have been investigated that the biochemical consequences of dCMP deaminase deficiency are clear—cells or phage survive because an additional metabolic route to dTTP is present, but the [dCTP]/[dTTP] or [hydroxymethyl-dCTP]/[dTTP] ratio increases dramatically, as expected from the loss of a quantitatively major pathway. However, the genetic consequences of these events underscore the difficulty in understanding all the biochemical factors contributing to spontaneous mutagenesis.

### Bibliography

1. B. R. de Saint Vincent, M. Déchamps, and G. Buttin (1980) *J. Biol. Chem.* **255**, 162–167.
2. F. Maley and G. F. Maley (1960) *J. Biol. Chem.* **235**, 2968–2970.
3. G. F. Maley and F. Maley (1968) *J. Biol. Chem.* **243**, 4506–4516.
4. E. Scarano, G. Geraci, and M. Rossi (1967) *Biochemistry* **6**, 192–201.
5. W. H. Fleming and M. J. Bessman (1967) *J. Biol. Chem.* **242**, 363–371.
6. G. F. Maley and F. Maley (1982) *Biochemistry* **21**, 3780–3785.
7. G. A. O'Donovan, G. Edlin, J. A. Fuchs, J. Neuhard, and E. Thomassen (1971) *J. Bacteriol.* **105**, 666–672.
8. J. T. Moore, R. E. Silversmith, G. F. Maley, and F. Maley (1993) *J. Biol. Chem.* **268**, 2288–2291.
9. K. X. B. Weiner, R. S. Weiner, F. Maley, and G. F. Maley (1993) *J. Biol. Chem.* **268**, 12983–12989.
10. L. Betts, S. Xiang, S. A. Short, R. Wolfenden, and C. W. Carter Jr. (1994) *J. Mol. Biol.* **235**, 635–656.
11. G. Weinberg, B. Ullman, and D. W. Martin Jr. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2447–2451.
12. S. E. Kohalmi, M. Glatke, E. M. McIntosh, and B. A. Kunz (1991) *J. Mol. Biol.* **220**, 933–946.
13. R. G. Sargent and C. K. Mathews (1987) *J. Biol. Chem.* **262**, 5546–5553.
14. E. Darè, L-H, Zhang, D. Janssen, and V. Bianchi (1995) *J. Mol. Biol.* **252**, 514–521.
15. J. M. Leeds, M. B. Slabaugh, and C. K. Mathews (1985) *Mol. Cell Biol.* **5**, 3443–3450.

### Suggestions for Further Reading

16. P. Reichard (1988) Interactions between deoxyribonucleotide and DNA synthesis. *Ann. Rev. Biochem.* **57**, 349–374. This review describes the functions of dNTP biosynthetic enzymes and their roles in coordinating DNA precursor synthesis and DNA replication.
17. B. A. Kunz, S. E. Kohalmi, T. A. Kunkel, C. K. Mathews, E. M. McIntosh, and J. A. Reidy (1990) Deoxyribonucleoside triphosphate levels: a critical factor in the maintenance of genetic stability. *Mutation Res.* **318**, 1–64. A descriptive review of mutagenic processes resulting from DNA precursor pool imbalances.

## Deoxyribonucleotide Biosynthesis And Degradation

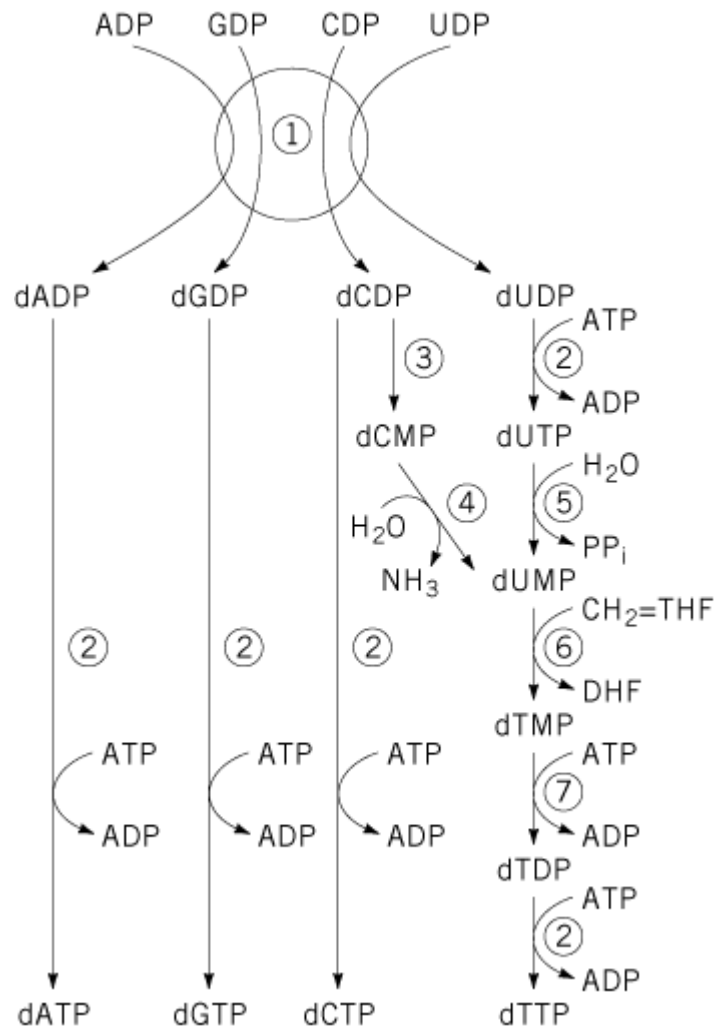
Most cells contain five to ten times as much RNA as DNA. The deoxyribonucleotide precursors to DNA use the same pathways for synthesis of purine and pyrimidine rings as do the ribonucleotides that are used for RNA synthesis and for a host of other metabolic functions (see [Purine Ribonucleotide Metabolism](#) and [Pyrimidine Ribonucleotide Metabolism](#)). Therefore, the pathways that lead from the *de novo* purine and pyrimidine synthetic pathways toward deoxyribonucleoside triphosphates (dNTPs) are quantitatively minor. They are of obvious significance, however. Not only must they be sufficiently active to supply the cell's needs for **DNA replication, recombination, and DNA repair** but they must be activated and deactivated in coordination with the timing of DNA replication during the [cell cycle](#) because dNTPs have no known metabolic roles other than their function as DNA precursors. Moreover, the pathways must be regulated so as to provide the four DNA precursors at rates corresponding to the base composition of the [genome](#) being served. dNTP pool imbalances that develop when regulatory mechanisms are perturbed have numerous biological consequences (1)—including [mutagenesis](#), recombination, induction of latent **viruses**, and cell death—caused by stimulation of [apoptosis](#) (2).

### 1. Pathways to dNTPs

#### 1.1. Chemical Differences between RNA and DNA

It is useful to begin our survey of deoxyribonucleotide biosynthesis by considering the processes through which the chemical differences between DNA nucleotides and RNA nucleotides arise. The methyl group of thymine, which distinguishes it from uracil, arises through the transfer of a single-carbon functional group to a uracil nucleotide, catalyzed by [thymidylate synthase](#). The presence of 2-deoxyribose as the sugar in DNA nucleotides rather than the ribose found in RNA comes about through reduction of the ribose sugar on a ribonucleotide substrate (see [Ribonucleotide Reductases](#)). Some aerobic bacteria and all anaerobic microorganisms studied carry out this reduction at the ribonucleoside triphosphate (rNTP) level. In all other organisms studied, however, the substrates for ribonucleotide reductase are the ribonucleoside 5'-diphosphates (rNDPs). Whether a particular reductase acts on rNDPs or rNTPs, a single enzyme reduces all four ribonucleotide substrates. Accordingly, the enzyme interacts with **allosteric** modifiers to ensure that the four DNA precursors are produced at rates commensurate with the base composition of the organism's genome. In this article, we will describe the predominant pathways to dNTPs that begin with reduction of rNDPs to deoxyribonucleoside diphosphates (dNDPs). These pathways are summarized in [Figure 1](#).

**Figure 1.** Pathways of dNTP biosynthesis *de novo*. Although these pathways are widespread, they are not universal, as indicated in the text. Enzyme 1, ribonucleoside diphosphate reductase; 2, nucleoside diphosphate kinase; 3, dCMP kinase (probably); 4, dCMP deaminase; 5, dUTPase; 6, thymidylate synthase; 7, thymidylate kinase. CH<sub>2</sub> = THF is 5,10-methylenetetrahydrofolate, and DHF is dihydrofolate.



### 1.2. The Role of Nucleoside Diphosphate Kinase

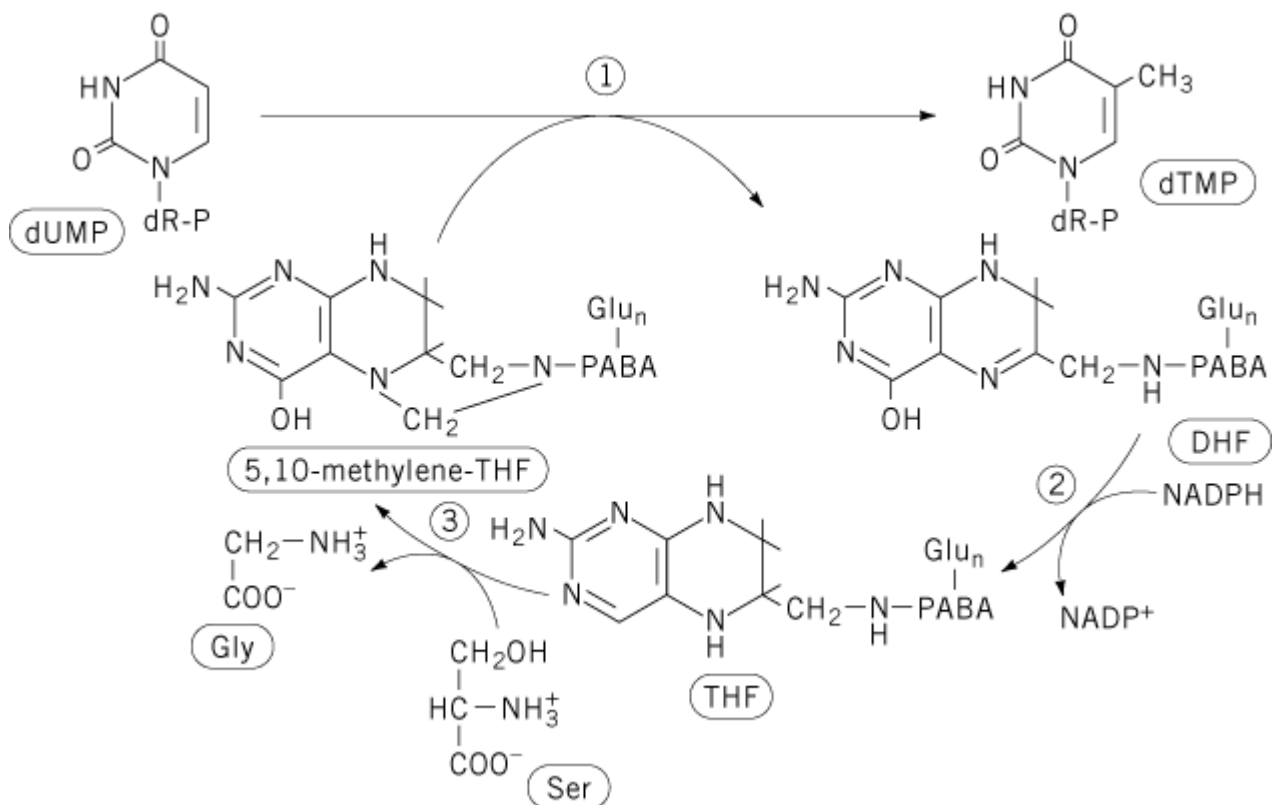
For three of the four nucleotides—ADP, CDP, and GDP—the conversion to the dNTP involves simply a phosphorylation of the rNDP reductase product, catalyzed by [nucleoside diphosphate kinase](#). In most cells, this enzyme is very active, and all known forms of the enzyme have very low specificity. Thus, the enzyme catalyzes the reversible transfer of the *g*-phosphate of any common rNTP or dNTP to the phosphate at the *b*-position of any common rNDP or dNDP. The equilibrium constant for each reaction catalyzed is close to unity. Thus, the direction in which a nucleoside diphosphate kinase-catalyzed reaction occurs *in vivo* depends on the concentrations of substrates and products. Because ATP is almost always the most abundant intracellular nucleoside triphosphate, most such reactions involve the ATP-dependent conversion of a ribo- or deoxyribonucleoside diphosphate to the corresponding triphosphate.

### 1.3. Thymine Nucleotide Biosynthesis

The biosynthesis of thymidine triphosphate is obviously more complex than that of any other dNTP given that thymine nucleotides are derived from uracil nucleotides and a methyl group must be created. In most cells, the dUDP produced by rNDP reductase action on UDP is phosphorylated to dUTP by nucleoside diphosphate kinase. dUTP is cleaved by an active pyrophosphatase, [deoxyuridine triphosphatase](#), or dUTPase. This enzyme plays a dual role: (1) Because dUTP is a good substrate for DNA polymerase, dUTPase minimizes the steady-state pool size of dUTP and, hence, helps to exclude dUMP from DNA. (2) The dUTPase reaction is a significant biosynthetic route to dUMP, the substrate for thymidylate synthase.

Thymidylate synthase catalyzes the transfer of a single-carbon functional group from 5,10-methylenetetrahydrofolate to position 5 of the pyrimidine ring in deoxyuridine monophosphate, yielding thymidine monophosphate (Fig. 2). dTMP is phosphorylated to dTDP by , and conversion to dTTP involves nucleoside diphosphate kinase.

**Figure 2.** The thymidylate synthesis cycle. Enzyme 1, thymidylate synthase; 2, dihydrofolate reductase; 3, serine transhydroxymethylase. dR-P is deoxyribose 5'-phosphate; PABA is *p*-aminobenzoate; Glu<sub>n</sub> refers to the multiple glutamate residues on naturally occurring folate coenzymes.



In the thymidylate synthase reaction, the transferred methylene group must be reduced to the methyl level, and the electron pair that brings this reduction about comes from the reduced pteridine ring of 5,10-methylenetetrahydrofolate. The coenzyme, therefore, loses both its methylene group and an electron pair, leading to dihydrofolate. Transformation of the coenzyme for reuse involves, first, its reduction to tetrahydrofolate by [dihydrofolate reductase](#) and, next, transfer of a single-carbon group to the pteridine ring, usually catalyzed by serine transhydroxymethylase. The stoichiometric requirement for the folate cofactor in the thymidylate synthase reaction probably explains the selective toxicity of dihydrofolate reductase inhibitors toward proliferating cells (3). Such inhibitors include Methotrexate, widely used in cancer chemotherapy, and Trimethoprim, an antibacterial agent that specifically inhibits dihydrofolate reductases of prokaryotic origin (see **Aminopterin**). Proliferating cells have a continuous requirement for dTTP synthesis, to sustain DNA replication. The greater the flux rate through thymidylate synthase *in vivo*, the more rapidly tetrahydrofolate pools will be depleted after administration of a dihydrofolate reductase inhibitor and, hence, the greater will be the sensitivity of those cells toward the growth-inhibiting or lethal effects of blockage of dihydrofolate reductase.

#### 1.4. Biosynthetic Routes to dUMP

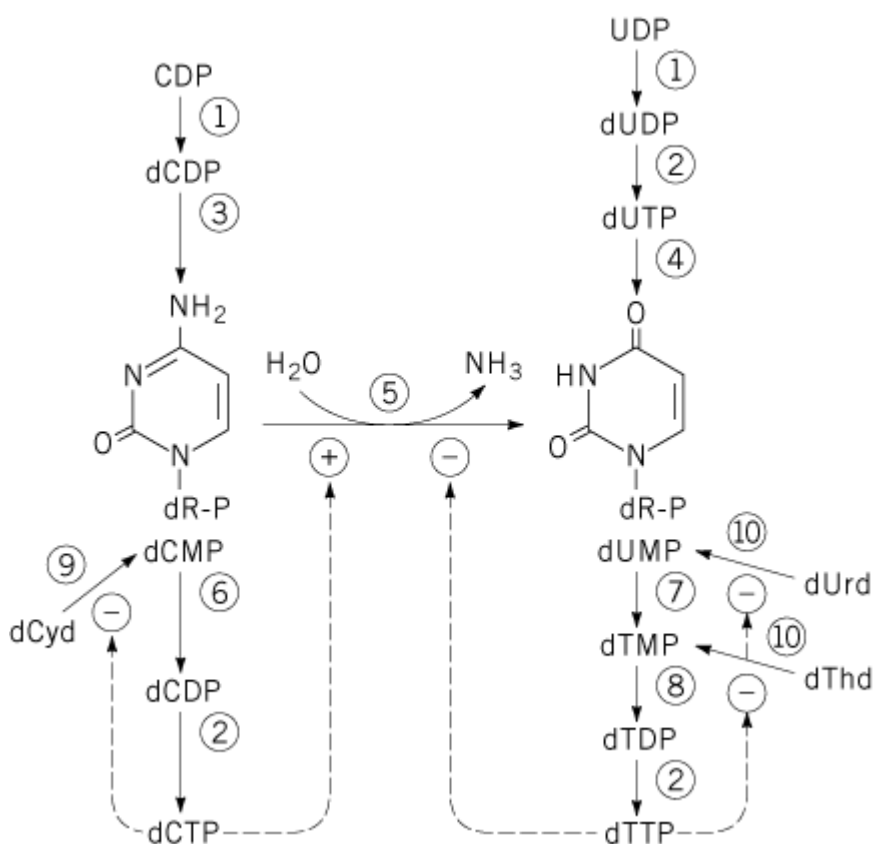
As noted earlier, the cleavage of deoxyuridine triphosphate to the corresponding monophosphate,



catalyzed by dUTPase, provides a significant source of dUMP for the synthesis of thymine nucleotides. In most cells, however, an additional route, involving deamination of deoxycytidine nucleotides, is more important in a quantitative sense. As shown in Figure 3, the hydrolytic deamination in most organisms involves conversion of dCMP to dUMP, catalyzed by **dCMP deaminase**. In a few bacteria, including *Escherichia coli*, the deamination occurs at the triphosphate level, with conversion of dCTP to dUTP, catalyzed by dCTP deaminase. In such organisms, the conversion of CDP to dUMP occurs by the following pathway:



**Figure 3.** dCMP deaminase as a branch point between dCTP and dTTP synthesis. Allosteric modifiers of dCMP deaminase are shown, as is the regulation of salvage synthetic pathways. Enzyme 1, ribonucleoside diphosphate reductase; 2, nucleoside diphosphate kinase; 3, dCMP kinase (possibly); 4, dUTPase; 5, dCMP deaminase; 6, dCMP kinase; 7, thymidylate synthase; 8, thymidylate kinase; 9, deoxycytidine kinase; 10, thymidine kinase.



The most widespread pathway, shown in Figure 3, involves deamination at the monophosphate level. In many vertebrate [cell lines](#), this is the predominant, often exclusive, pathway leading to dUMP (4). It is not clear why the *in vivo* flux rate from UDP to dUDP is so low, particularly when the activities of ribonucleotide reductase on CDP and UDP are regulated virtually identically as shown by assays of the purified enzyme *in vitro* (5). Whatever the reason, dCMP deaminase is an important metabolic branch point between routes to dTTP and dCTP. Allosteric regulation of this enzyme—activation by dCTP and inhibition by dTTP—ensures that these two dNTPs are produced at relative rates commensurate with their need for DNA synthesis. dCMP deaminase is not essential for cell viability, at least as determined in cell culture systems. However, mutant cells lacking dCMP deaminase have abnormally high dCTP pools (1, 6), and the resultant increase in the [dCTP]/[dTTP] pool ratio often brings about a **mutator** phenotype, in which the dCTP pool expansion stimulates its incorporation

opposite template nucleotides other than dGMP (1).

## 2. Deoxyuridine Nucleotide Metabolism

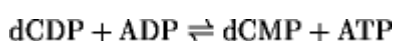
As noted earlier, dUTPase plays a role in providing dUMP for the thymidylate synthase reaction. Of equal or greater importance is its role in minimizing the incorporation of dUMP into DNA in place of dTMP (see [Deoxyuridine Triphosphatase](#)). Cells possess two independent mechanisms to minimize the amount of dUMP in DNA: (1) dUTPase acts to deplete the dUTP pool, and hence to minimize dUMP incorporation by DNA polymerase, usually opposite dAMP in template DNA. (2) The second mechanism is a base [excision repair](#) process, starting with . This enzyme scans DNA, and when it contacts a dUMP residue, it cleaves the glycosidic bond, probably by a base-flipping mechanism (7). This reaction creates an abasic site, which initiates a repair process that results in insertion of the correct nucleotide at that site (8).

Because the base-pairing properties of uracil and thymine are essentially identical, it seems likely that the mechanisms for uracil exclusion from DNA are not designed primarily to replace dUMP with dTMP, when that dUMP is base-paired with dAMP in the template DNA. Instead, the probable target for DNA uracil exclusion is dUMP residues incorrectly paired with dGMP, but DNA-uracil *N*-glycosylase simply recognizes every dUMP residue in DNA, no matter what its partner. G · U DNA base pairs would arise most frequently by deamination of a dCMP residue paired in duplex DNA with dGMP. Cytosine undergoes spontaneous deamination at appreciable rates whether as a free nucleic acid precursor or as a residue in DNA or RNA. If C in a G · C base pair underwent deamination, the resultant G · C base pair would generate a G · C and an A · U base pair in the next round of replication and, ultimately, the deamination would lead to a G · C → A · T transition mutation. Thus, the DNA-uracil exclusion systems probably act to maintain stability of the genome.

## 3. Salvage Routes to dNTPs

As noted in [Salvage pathways to nucleotide biosynthesis](#), deoxyribonucleotide salvage pathways involving uptake of extracellular precursors primarily use deoxyribonucleoside kinases. Human cells contain four such enzymes of varying specificities, two located in the cytosol and two in [mitochondria](#), whereas other organisms, such as *E. coli*, contain thymidine kinase as the only deoxyribonucleoside kinase. Thymidine kinase has received particularly intensive study, partly because of the mechanism of its cell cycle regulation (9) but largely because the enzyme is so useful as a means for incorporating radiolabel into DNA. For reasons still not clear, thymidine competes extremely effectively with the *de novo* synthetic pathway to dTTP such that, in many animal cell systems, radiolabeled thymidine is incorporated into DNA at full specific activity, often bypassing substantial endogenous pools generated by *de novo* synthesis (10). One popular experimental organism for which this does not work is **yeast**; fungi lack thymidine kinase. Investigators have circumvented this difficulty, however, by designing yeast strains that are permeable to dTMP, strains for which exogenous dTMP can be used as a labeled DNA precursor.

For salvage of deoxyribonucleoside monophosphates released by intracellular DNA degradation, the deoxyribonucleoside monophosphate kinases play the key roles. Animal cells contain four such enzymes, each specific for one deoxyribonucleotide, ie, dAMP kinase, dCMP kinase, dGMP kinase, and dTMP kinase. The enzyme phosphorylating dAMP acts also on AMP and is the well-known adenylate kinase, or myokinase. dCMP kinase acts also on UMP, and dTMP kinase acts also on dUMP. dTMP kinase is involved also in *de novo* dTTP synthesis, as shown in Figure 1. dCMP kinase may also play a role in *de novo* dNTP synthesis. The enzyme converting dCDP (produced by ribonucleotide reductase) to dCMP, en route to dUMP and dTMP, has still not been identified. Since nucleotide kinases all have equilibrium constants close to 1, it is quite possible that the role of dCMP kinase is to carry out the synthesis of dCMP:



By contrast, dAMP and dGMP kinases play roles only in nucleotide salvage reactions because the *de novo* pathways lead directly from ribonucleoside diphosphate to deoxyribonucleoside diphosphate to deoxyribonucleoside triphosphate.

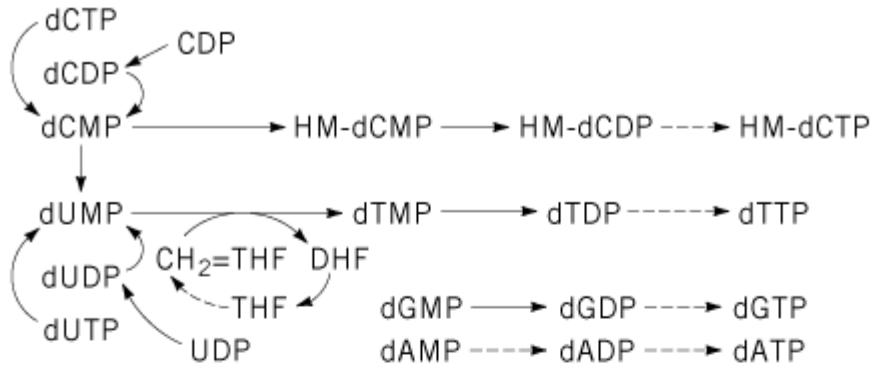
The action of nucleoside phosphorylases provides another route to salvage dNTP synthesis, eg, deoxyadenosine + P<sub>i</sub> → adenine + deoxyribose-1-phosphate. Now the deoxyribose-1-phosphate can react with another base in a nucleoside phosphorylase-catalyzed reaction proceeding in the reverse direction, eg, dR-1-P + Thy → dThd + P<sub>i</sub>. Thymidine can then become a nucleotide via the thymidine kinase reaction leading to dTMP. Although this process can significantly change relative dNTP pool sizes, it does not involve net deoxyribonucleotide synthesis; rather, it involves redistribution of the deoxyribosyl units linked to purine and pyrimidine bases.

#### 4. Biosynthesis of Novel Nucleotides in Bacteriophage Infection

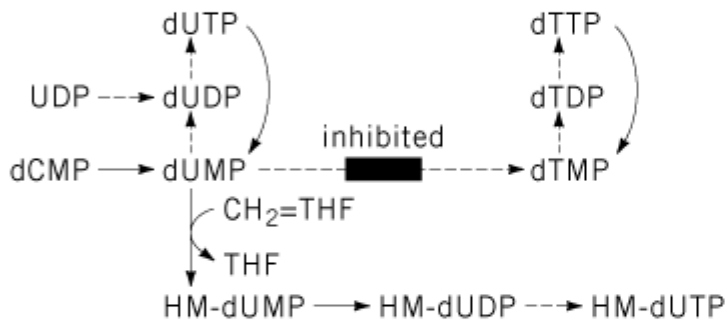
Much of what is known about pathways of dNTP synthesis in bacterial and animal cells came from investigations of the biosynthesis of unusual modified nucleotides in infection by certain bacteriophage (11). T-even bacteriophages of *E. coli* contain 5-hydroxymethylcytosine completely substituted for cytosine; many hydroxymethyldeoxycytidylate (HM-dCMP) residues in phage DNA are further modified by glycosidic links through the hydroxymethyl group to one or two glucose residues. These modifications occur through the action of phage-encoded enzymes that catalyze reactions comparable to those in cellular metabolism (Fig. 4). For example, the hydroxymethylation reaction is carried out by an enzyme, dCMP hydroxymethylase, that transfers a single-carbon group from methylenetetrahydrofolate to C-5 of dCMP, much as thymidylate synthase modifies C-5 of the pyrimidine dUMP (12); in fact, T-even phages encode a thymidylate synthase that displays significant amino acid sequence [homology](#) with dCMP hydroxymethylase. Not shown in Figure 4 is the involvement of glucosylation in the transfer of glucose to hydroxymethyl groups of HM-dCMP residues after their incorporation into DNA.

**Figure 4.** Metabolic pathways induced in bacteriophage infection. Solid arrows identify reactions known to be catalyzed by phage-coded enzymes; dashed arrows identify reactions catalyzed only by enzymes encoded by the host cell genome.

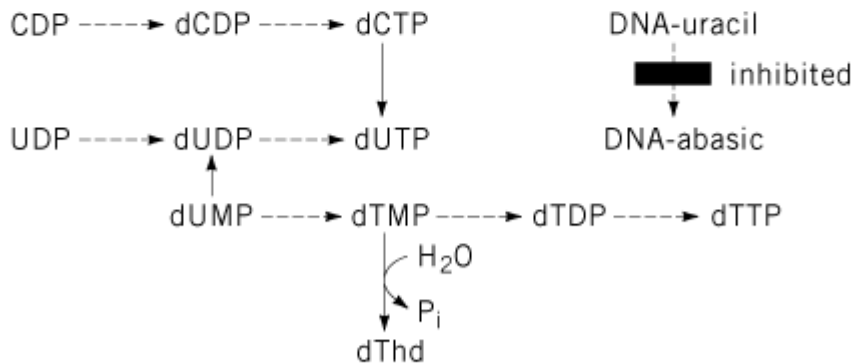
T2, T4, T6



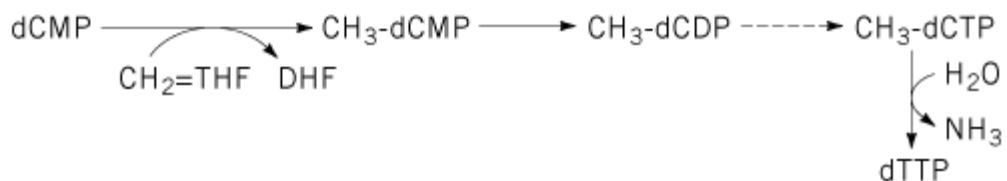
SP01, SP82



PBS2



XP12



Other modifications are carried out by bacteriophage SP01 and SP82 of *Bacillus subtilis*, in which 5-hydroxymethyluracil substitutes for thymine and two other *B. subtilis* phages, PBS1 and PBS2, which replace thymine with uracil. The relevant pathways are also shown in Figure 4. Finally, the figure shows pathways induced in *Xanthomonas oryzae* by XP12, a bacteriophage in which cytosine is completely replaced by 5-methylcytosine. These base modifications help the virus to parasitize the host bacterium, overrunning pathways used by the host cell to produce its own DNA. Perhaps it should not be surprising that these are all virulent phages.

## 5. Deoxyribonucleoside Catabolism

### 5.1. Substrate Cycles and dNTP Pool Regulation

Nucleotide kinase-catalyzed reactions are readily reversible. Whereas these enzymes have long been considered to participate primarily in dNTP synthesis, there is now good evidence that nucleotide kinases function in both directions *in vivo* and that they can participate in dNTP degradation (13). For example, when DNA replication was blocked in mammalian cells with aphidicolin, intracellular dNTPs were degraded, and deoxyribonucleosides were excreted into the medium in significant amounts (13, 14), with little change in the steady-state pool size of the four dNTPs. The implication of these findings is that dNTP pool sizes are controlled not only at the level of synthesis but are also regulated by catabolism.

Under conditions that do not interfere with DNA synthesis in cultured cells, [pulse-chase experiments](#) show that the rate of utilization of a particular dNTP for DNA synthesis is equivalent to the rate at which the pool of that dNTP turns over (15, 16). These observations argue against the compartmentation of dNTP pools in animal cells; with one known exception, all the dNTP molecules in these cells are accessible to the replication machinery. The one exception is dCTP labeled by deoxycytidine (13); much of this pool is used for synthesis of the deoxycytidine counterparts to cytidine diphosphate-containing **lipid precursors**, such as CDP-choline. It is not yet clear whether these dCDP-containing “liponucleotides” play a metabolic role distinct from those of the more abundant cytidine-containing lipids.

For cells that don't excrete deoxyribonucleosides but instead degrade them further, that further degradation usually involves attack by nucleoside phosphorylases to yield deoxyribose-1-phosphate plus the free base. Further degradation of the purine and pyrimidine bases proceeds by pathways outlined in [Purine ribonucleotide metabolism](#) and [Pyrimidine ribonucleotide metabolism](#).

### 5.2. Enzymatic Degradation of dGTP

A widely distributed enzyme is a dGTP pyrophosphatase, which catalyzes the hydrolytic cleavage of dGTP to dGMP and pyrophosphate. The metabolic role of this enzyme was obscure until it was found in *E. coli* that the enzyme is encoded by the *mutT* gene (17). The *mutT* gene product acts to prevent [transversion mutations](#) that proceed by  $G \cdot C \rightarrow T \cdot A$  and  $A \cdot T \rightarrow C \cdot G$  pathways. Mutant alleles of *mutT* cause spontaneous rates of these transversion mutations at other loci to increase by 100- to 10,000-fold. Mutations of these types are caused by oxidizing agents, such as hydrogen peroxide, which act in large part by oxidizing guanine residues in DNA to 8-oxoguanine. This oxidized purine readily base-pairs with adenine, setting in motion the events leading to mutagenesis. These realizations led to the discovery that the MutT protein is far more active in cleaving 8-oxo-dGTP than dGTP. Thus, it seems clear that the mutagenic effects of oxidizing agents include oxidation of guanine, not as a DNA residue but as a DNA precursor that is incorporated into DNA and then stimulates a transversion pathway by pairing with adenine. The MutT nucleotidase prevents this by degrading the mutagenic DNA substrate, 8-oxo-dGTP, before it has a chance to be incorporated into DNA.

*E. coli* contains a second dGTPase, which is distinctive in that it hydrolyzes dGTP to deoxyguanosine plus triphosphate. Unlike the MutT dGTPase, which is widely distributed in organisms, this other dGTPase is found only in enteric bacteria (18). Its metabolic role is still obscure.

## Bibliography

1. B. A. Kunz, S. E. Kohalmi, T. A. Kunkel, C. K. Mathews, E. M. McIntosh, and J. A. Reidy (1994) *Mutation Res.* **318**, 1–64.
2. F. J. Oliver, M. K. L. Collins, and A. López-Rivas (1996) *Biochem. J.* **316**, 421–425.
3. J. J. McCormack (1981) *Med. Res. Revs.* **1**, 303–331.
4. V. Bianchi, E. Pontis, and P. Reichard (1985) *Mol. Cell. Biol.* **7**, 4218–4224.

5. L. Thelander and P. Reichard (1979) *Ann. Rev. Biochem.* **48**, 133–158.
6. G. L. Weinberg, B. Ullman, and D. W. Martin Jr. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2447–2451.
7. R. J. Roberts (1995) *Cell* **82**, 9–12.
8. E. C. Friedberg, G. C. Walker, and W. Siede (1995) *DNA Repair and Mutagenesis*, ASM Press, Washington DC, pp. 147–153.
9. J. L. Sherley and T. J. Kelly (1988) *J. Biol. Chem.* **263**, 8350–8358.
10. C. K. Mathews (1975) *Exptl. Cell Res.* **92**, 47–56.
11. C. K. Mathews (1977) In *Comprehensive Virology*, Vol. 7, (H. Fraenkel-Conrat and R. R. Wagner, eds.), Plenum Press, New York, pp. 179–294.
12. S. S. Cohen (1968) *Virus-Induced Enzymes*. Columbia University Press, New York.
13. P. Reichard (1988) *Ann. Rev. Biochem.* **57**, 349–374.
14. B. Nicander and P. Reichard (1985) *J. Biol. Chem.* **260**, 9216–9222.
15. B. Nicander and P. Reichard (1983) *Proc. Natl. Acad. Sci. USA* **80**, 1347–1351.
16. V. Bianchi, E. Pontis, and P. Reichard (1992) *Exptl. Cell Res.* **199**, 120–128.
17. H. Maki and M. Sekiguchi (1992) *Nature* **355**, 273–275.
18. S. Quirk and M. J. Bessman (1991) *J. Bacteriol.* **173**, 6665–6669.

### Suggestion for Further Reading

19. The best single reference in this area is Chapter "2", "Biosynthesis of DNA Precursors," in A. Kornberg and T. Baker (1992) *DNA Replication*, 2nd Ed., W. H. Freeman and Co., New York, pp. 53–100.

## Deoxyuridine Triphosphatase

Living cells contain two [enzyme](#) systems whose primary role is evidently to assure the exclusion of uracil from DNA where thymine occurs instead. Deoxyuridine triphosphatase (dUTPase) catalyzes the hydrolysis of dUTP to deoxyuridine monophosphate and pyrophosphate, thereby removing from the nucleotide pool a precursor, dUTP, that is recognized by **DNA polymerases** as readily as is thymidine triphosphate for incorporation into DNA opposite template adenine residues. The other system involves a DNA [base excision repair](#) process that begins with the action of DNA-uracil *N*-glycosylase. This enzyme cleaves the glycosidic bond linking uracil to deoxyribose, and an oligonucleotide patch that contains the abasic site is removed and replaced.

Current interest in dUTPase focuses on its very high substrate specificity, its significance as a possible chemotherapeutic target, the evolution and functions of virus-coded forms of dUTPase, and the generation and regulation of organelle-specific isoforms in eukaryotic cells.

### 1. Biological Functions of dUTPase

Although dUTPase was discovered in the early 1960s ([1](#), [2](#)), the enzyme was not characterized as a highly purified protein until 1978 ([3](#)). Much of the interest in the 1970s related to the role of uracil metabolism in creating short DNA fragments that had been identified as intermediates in DNA replication. After Okazaki proposed that DNA replication involves the discontinuous synthesis of short DNA fragments—which became known as [Okazaki Fragments](#)—several laboratories reported

that nascent DNA chains are complementary to both template DNA strands, suggesting that both daughter DNA strands grow discontinuously. However, Tye et al. (4) showed that many of the short fragments arose not as intermediates in replication but through the action of the uracil exclusion systems. Incorporation of dUMP into DNA, followed by the action of uracil *N*-glycosylase and base [excision repair](#), caused DNA strand breakage that converted nascent high-molecular-weight DNA into small fragments. When corrections were made for this factor, it was found that **leading-strand** DNA replication occurs continuously, with short fragments of DNA arising as replication intermediates only in synthesis of the lagging strand.

Why uracil should be excluded from DNA was not immediately apparent because uracil is identical to thymine in its base pairing with adenine, and so its substitution for thymine was not expected to affect DNA function or genetic specificity significantly. Indeed, the existence of bacteriophages, such as the *Bacillus subtilis* phage PBS2, whose DNA contains uracil but no thymine (5), indicates that substituting U for T is *biologically* acceptable. However, it was recognized that uracil can arise in DNA either through incorporation of dUMP opposite dAMP or through deamination of dCMP in a G · C base pair. This latter event is potentially **mutagenic** because the resultant G · U base pair would be converted to an A · T base pair in subsequent rounds of replication. Because dUTPase has no direct involvement in these events, it was thought that dUTPase should be a dispensable function in most organisms. However, El-Hajj et al. (6) were unable to isolate a null mutant of dUTPase in *Escherichia coli*, leading them to speculate that the protein plays an additional, essential metabolic role; similar results were reported for yeast (7). In a later study, however, El-Hajj et al. (8) analyzed a **conditional lethal** dUTPase (*dut*) [mutation](#) of *E. coli* and found that it could be suppressed by a mutation in *dcd*, the structural gene for dCTP deaminase, the enzyme that synthesizes dUTP in enteric bacteria (see [Deoxyribonucleotide Biosynthesis And Degradation](#)). This finding, plus the construction of a multiple mutant bacterium that survived extensive substitution of uracil for thymine in its DNA (8), indicated that the primary role of dUTPase is indeed to keep the dUTP pool small and facilitate the exclusion of uracil from DNA. The lethality of dUTPase null mutants is ascribed, then, to excessive accumulation of dUMP in DNA, leading to double-strand breaks that result when repair of single-strand breaks occurs at nearby sites on opposite DNA strands.

## 2. Structure of Deoxyuridine Triphosphatase

The dUTPase molecule must be constructed to exacting design specifications in order to cleave efficiently one deoxyribonucleoside triphosphate while having little or no effect on structurally similar dNTPs needed for DNA synthesis, in particular, dUTP and dTTP. Accordingly, kinetic analysis of *E. coli* dUTPase (9) showed the ratio  $k_{\text{cat}}/K_m$  (see [Enzymes](#)) to be, remarkably, more than  $10^5$ -fold lower for dCTP than for dUTP. The nucleotides dTTP and UTP were even poorer substrates as no action of dUTPase on these nucleotides was detectable. The effects of pH on the  $K_m$  for dUTP suggested important contributions of **hydrogen bonding** involving N3 and the carbonyl oxygen on C4 in establishing the specificity of dUTP binding to dUTPase, which is consistent with structural studies (10, 11). Crystallographic analysis of both the *E. coli* and human enzymes reveals a homotrimeric structure, with each of three substrate-binding sites made up from residues in two adjacent [polypeptide chains](#). The bacterial enzyme was crystallized in the presence of a **competitive inhibitor**, dUDP, which identified the [active site](#) and provided a structural explanation of the kinetic observations.

## 3. dUTPase as a Chemotherapeutic Target

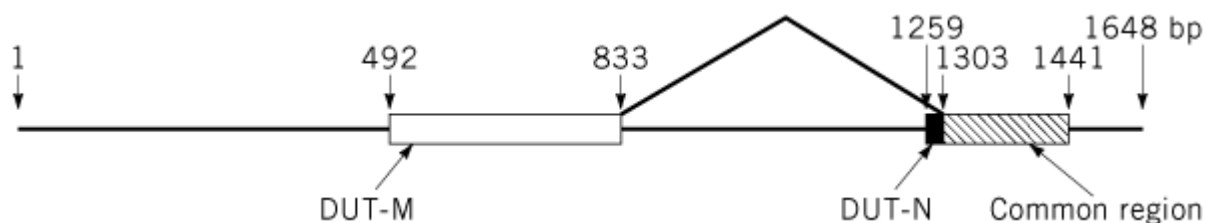
Thymidylate metabolism has long been recognized as a target for drugs that would act to inhibit DNA replication in either cancer or infectious diseases (see [Thymidylate Synthase](#); [Fluorodeoxyuridine](#); [Dihydrofolate Reductase \(DHFR\)](#)). Beginning in the early 1980s, it was recognized that a significant contributor to the lethality of inhibitors of dihydrofolate reductase or thymidylate synthase is the accumulation of dUTP and its resultant incorporation into DNA (12). Inhibition of thymidylate synthase flux either directly, by 5-fluorodeoxyuridylate, or indirectly, by

inhibition of dihydrofolate reductase, was seen to cause accumulation of dUMP which, in turn, would be converted to dUTP faster than it could be cleaved by dUTPase. Hence, the resultant accumulation of dUMP in DNA was seen to cause double-strand DNA breaks by the process discussed above. The existence of a class of fluorodeoxyuridine-resistant cell mutants with elevated levels of dUTPase (13) is consistent with this interpretation. Now that [X-ray crystallography](#) structures for dUTPase are available, this enzyme has been recognized as a suitable target for structure-assisted drug design. In principle, dUTPase inhibitors could be used, either alone or in conjunction with inhibitors of dihydrofolate reductase or thymidylate synthase, to kill proliferating cells selectively by stimulating the incorporation of dUMP into their DNA.

#### 4. dUTPase Isoforms in Eukaryotic Cells

One of the most intriguing recent developments concerning dUTPase (14) relates to the discovery of two different **isoforms** in mammalian cells, one located predominantly in the [nucleus](#) and the other localized to [mitochondria](#) (although an earlier study (15) had localized human dUTPase to the cytosol). Presumably, the mitochondrial isoform plays the same role in mitochondrial DNA synthesis as the nuclear enzyme plays for nuclear DNA replication. Of great interest is the recent finding of Ladner and Caradonna (16) that both isoforms are encoded by the same nuclear gene, with [alternative splicing](#) at the 5' end of the dUTPase gene (Fig. 1). The use of different 5' exons for the two isoforms (see [Introns, Exons](#)) gives a subunit size of 22 kDa for the cytosolic enzyme and 23 kDa for the mitochondrial isoform. This may be the first reported instance of alternative splicing of the same transcript to generate isoforms of the same enzyme localized in different cell compartments.

**Figure 1.** Alternative splicing at the 5' end of the human dUTPase gene as a means to generate nuclear and mitochondrial isoforms of dUTPase. Exons are boxed, and open reading frames of the mitochondrial, nuclear, and common sequences are shown by different shading. Intron/exon boundaries are indicated and numbered with respect to an arbitrary start point for the 5' end of the gene. The amino-terminal splice junction is indicated by joined lines. Reprinted with permission from R. D. Ladner and S. Caradonna (1997) *J. Biol. Chem.* **272**, 19072–19080.



[Cell-cycle](#) analysis of protein and [messenger RNA](#) levels shows that the mitochondrial enzyme is expressed constitutively throughout the cell cycle (as is mitochondrial DNA replication), whereas the expression of the nuclear form is closely linked to nuclear DNA replication status. The nuclear isoform undergoes **protein degradation** as cells exit the cycle. Whether this differential regulation is related to **phosphorylation** of the enzyme is an intriguing unanswered question. Ladner et al. (14) showed that the cytosolic dUTPase is phosphorylated at Ser11, in a **cyclin**-dependent protein [kinase](#) phosphorylation site, whereas the same site remains unphosphorylated in the mitochondrial isoform. The functional significance of this difference has not yet been defined.

#### 5. Virus-Encoded dUTPases

Numerous **virus**-specific forms of dUTPase have been described, including those encoded by the RNA genomes of [retroviruses](#). The first viral dUTPase to be described (17) was shown to be an activity of dCTPase, an enzyme induced after T-even bacteriophage infection, which is responsible



for exclusion of cytosine from phage DNA. This activity is essential to the observed complete replacement of cytosine by 5-hydroxymethylcytosine in the phage genome. Purification of dCTPase revealed an associated dUTPase activity (18), and both activities were shown to be intrinsic to the same protein. Also, unlike the unrelated bacterial and eukaryotic dUTPases, the phage enzyme is equally active on the respective deoxyribonucleoside diphosphates, cleaving dCDP and dUDP to dCMP and dUMP, respectively, plus inorganic phosphate. In T4, this multifunctional enzyme is encoded by gene 56. Mutants defective in this gene are lethal because they incorporate dCMP into phage DNA, and the viral genome encodes **nucleases** that specifically recognize and degrade cytosine-containing DNA.

dUTPase is encoded by the genomes of [herpes viruses](#), pox viruses, several nonprimate lentiviruses, and type B and D retroviruses. A sequence in the *gag-pro* gene region of retrovirus genomes, originally identified as the evolutionary result of the duplication of a [proteinase](#) gene, was found, as these sequences began to accumulate in the data banks, to be homologous to dUTPase genes (19). Cloning and expression of these “pseudoproteinase” genes showed that they did indeed express functional dUTPases (20, 21). It is of particular interest that the human immunodeficiency virus (HIV) genome does not encode a dUTPase, although such a gene is expressed by the closely related feline immunodeficiency virus (FIV), which is being studied as a model for HIV infection. Lerner et al (22) infected cats with wild-type FIV and with a dUTPase-negative FIV mutant. After nine months of infection, sequence analysis of viral genomes isolated from both groups of infected animals showed a fivefold increase in mutations in the dUTPase-negative infections, most of them involving G → A transitions. These could have arisen from the formation of rG · dU base pairs during **reverse transcription** in the presence of expanded dUTP pools. Of great interest is the possibility that the extreme genetic variability of HIV is related to the lack of a dUTPase gene in the wild-type HIV genome, giving HIV the same mutator **phenotype** as observed in dUTPase-deficient FIV mutants.

## Bibliography

1. L. E. Bertani, A. Häggmark, and P. Reichard (1961) *J. Biol. Chem.* **236**, PC67–PC68.
2. G. R. Greenberg and R. L. Somerville (1962) *Proc. Natl. Acad. Sci. USA* **48**, 247–257.
3. J. Shlomai and A. Kornberg (1978) *J. Biol. Chem.* **253**, 3305–3312.
4. B-K. Tye, P-O. Nyman, I. R. Lehman, S. Hochhauser, and B. Weiss (1977) *Proc. Natl. Acad. Sci. USA* **74**, 154–157.
5. A. R. Price and H. R. Warner (1975) *J. Biol. Chem.* **250**, 8804–8811.
6. H. H. El-Hajj, H. Zhang, and B. Weiss (1988) *J. Bacteriol.* **170**, 1069–1075.
7. M. H. Gadsden, E. M. McIntosh, J. C. Game, P. J. Wilson, and R. H. Haynes (1993) *EMBO J.* **12**, 4425–4431.
8. H. H. El-Hajj, L. Wang, and B. Weiss (1992) *J. Bacteriol.* **174**, 4450–4456.
9. G. Larsson, P. O. Nyman, and J-O. Kvassman (1996) *J. Biol. Chem.* **271**, 24010–24016.
10. G. Larsson, L. A. Svensson, and P. O. Nyman (1996) *Nature Struct. Biol.* **3**, 532–538.
11. C. D. Mol, J. M. Harris, E. M. McIntosh, and J. A. Tainer (1996) *Structure* **4**, 1077–1092.
12. M. Goulian, B. Bleile, and B. Y. Tseng (1980) *J. Biol. Chem.* **255**, 10630–10637.
13. C. E. Canman, T. S. Lawrence, D. S. Shewach, H. Y. Tang, and J. Maybaum (1993) *Cancer Res.* **53**, 5219–5224.
14. R. D. Ladner, D. E. McNulty, S. A. Carr, G. D. Roberts, and S. J. Caradonna (1996) *J. Biol. Chem.* **271**, 7745–7751.
15. J. A. Vilpo and H. Autio-Harminen (1983) *Scand. J. Clin. Lab. Invest.* **43**, 583–590.
16. R. D. Ladner and S. Caradonna (1997) *J. Biol. Chem.* **272**, 19072–19080.
17. G. R. Greenberg (1966) *Proc. Natl. Acad. Sci. USA* **56**, 1226–1232.
18. A. R. Price and H. R. Warner (1969) *Virology* **39**, 882–892.

19. D. J. McGeoch (1990) *Nucleic Acid Res.* **18**, 4105–4110.
20. A-C. Bergman, O. Björnberg, J. Nord, P. O. Nyman, and A. M. Rosengren (1994) *Virology* **204**, 420–424.
21. N. A. Roseman, R. K. Evans, E. L. Mayer, M. A. Rossi, and M. B. Slabaugh (1996) *J. Biol. Chem.* **271**, 23506–23511.
22. D. L. Lerner, P. C. Wagaman, T. R. Phillips, O. Prospero-Garcia, S. J. Henriksen, H. S. Fox, F. E. Bloom, and J. H. Elder (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7480–7484.

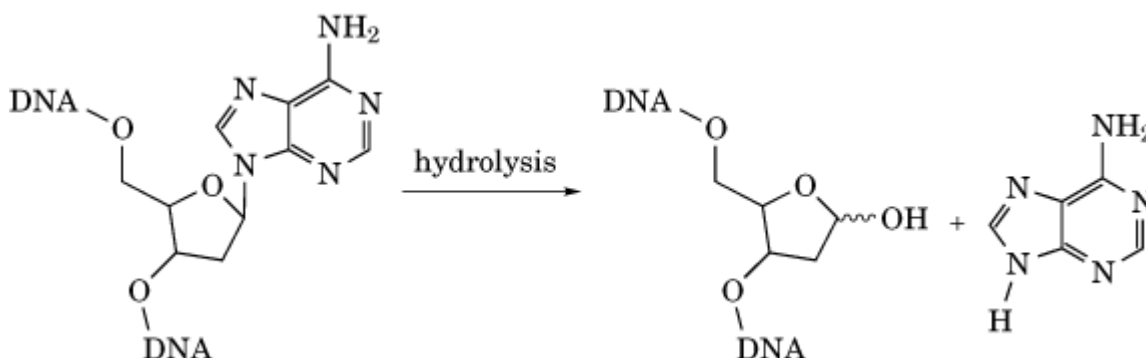
### Suggestions for Further Reading

23. A. Kornberg and T. A. Baker (1991) *DNA Replication*, 2nd ed. W. H. Freeman & Co., San Francisco. This book contains a good historical account of the discovery and elucidation of the metabolic roles of dUTPase.
24. E. M. McIntosh and R. H. Haynes (1997) dUTP pyrophosphatase as a potential target for chemotherapeutic drug development. *Acta Biochim. Polon.* **44**, 159–171. A recent review of the prospects for developing dUTPase inhibitors as useful drugs.

### Depurination

An abasic site one lacking a nucleotide base, is created at a site in DNA when the glycosidic bond connecting a **purine** base to the deoxyribose sugar is cleaved, leaving the phosphodiester backbone of the DNA intact (Fig. 1). Depurination occurs at an estimated rate of  $10^4$  depurinations per mammalian cell per day (1), often through spontaneous hydrolysis. Other processes in the cell that lead to depurination are enzymatic removal of altered bases by specific glycosylases or through chemical modifications of purine bases (primarily alkylation at N7) that labilize the glycosidic bond. There is good evidence that depurination plays a major role in spontaneous **mutagenic** events (2). Although apurinic sites inhibit chain elongation by **DNA polymerases**, bypass must occur because there is enhanced incorporation of nucleotides that are noncomplementary to the template bases. The rate of misincorporation varies with different polynucleotide complexes, but the frequency of misincorporation is directly proportional to the extent of depurination (2).

**Figure 1.** Formation of an abasic site in DNA depurination.



Weymouth and Loeb (3) and Kunkel et al. (4) studied the fidelity of DNA polymerases in copying natural DNA templates and the mutagenic implications of the misincorporation events. They measured revertants generated by misincorporation in *am3* bacteriophage fX174, a mutant containing a single [base-pair substitution](#) mutation. These studies confirmed that the copying of depurinated *am3* fX174 DNA templates by purified DNA polymerases is mutagenic, and that the frequency of revertants correlates positively with the ability of different polymerases to copy past apurinic sites. Eukaryotic DNA polymerases are error-prone and read past apurinic sites at a high frequency.

Although direct mutagenicity is readily envisaged following error-prone read through of an apurinic site, it also appears that many mutations resulting from depurination in prokaryotes are **SOS-response-dependent** (5). Miller and Low (6) examined more than 600 independent spontaneous mutations in the *lacI* gene of SOS-induced *E. coli*. Most were specifically G.C → T.A, and to a lesser extent A.T → T.A, transversions and occurred primarily at certain sites. These transversion mutations are thought to have arisen from the insertion of dAMP opposite purines or apurinic sites, possibly through the SOS-induction process altering the specificity of the [DNA replication](#) complex to favor misincorporation of dAMP opposite purines. Miller and Low suggested that many spontaneous mutations may result from replication past cryptic apurinic sites. There are also several lines of evidence that apurinic sites are intermediates in mutagenesis by chemicals that form bulky DNA adducts (2).

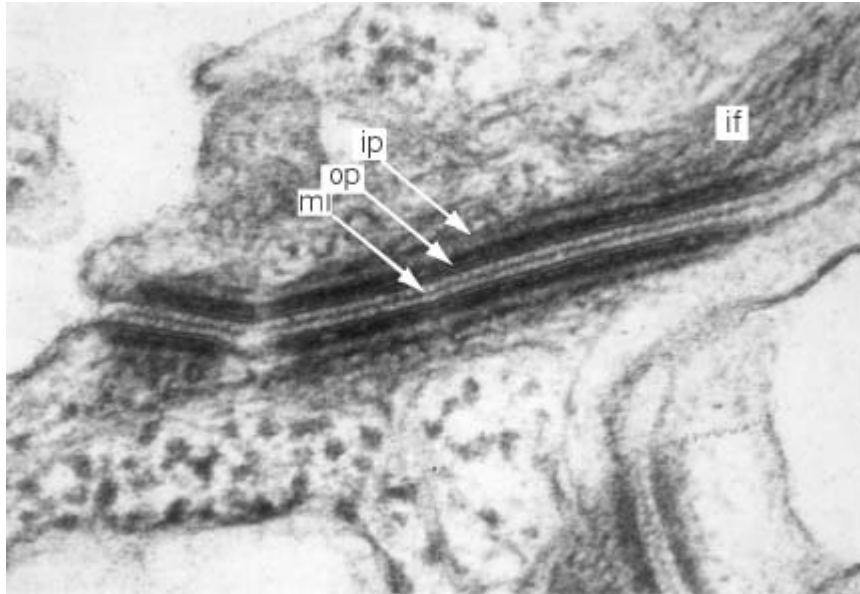
#### Bibliography

1. T. Lindahl and B. Nyberg (1972) *Biochemistry* **11**, 3610–3618.
2. L. A. Loeb (1985) *Cell* **40**, 483–484.
3. L. A. Weymouth and L. A. Loeb (1978) *Proc. Natl. Acad. Sci. USA* **75**, 1924–1928.
4. T. A. Kunkel, R. M. Schaaper, and L. A. Loeb (1983) *Biochemistry* **22**, 2378–2384.
5. T. A. Kunkel (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1494–1498.
6. J. H. Miller and K. B. Low (1984) *Cell* **37**, 675–682.

#### Desmosomes, Desmocollin, Desmoglein, and Desmoplakin

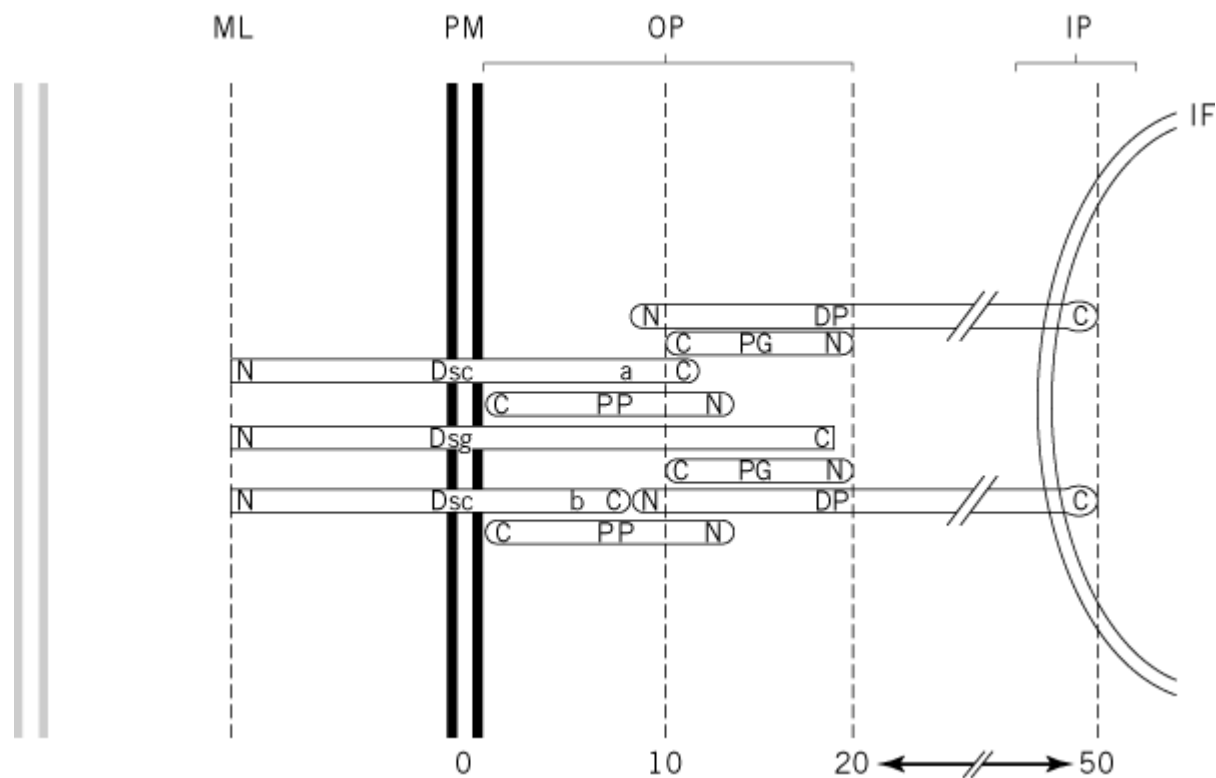
Desmosomes are punctate, adhesive intercellular [cell junctions](#) that bind cells together and provide membrane anchoring points for the intermediate-filament [cytoskeleton](#) (Fig. 1). They are circular membrane domains of up to 0.5 μm in diameter. The intercellular material, or desmoglea, has a highly organized ultrastructure consisting of an electron-dense midline that is bridged to the plasma membranes of the adhering cells across a distance of ~30 nm. Close to the cytoplasmic faces of the plasma membranes are dense outer plaques, the most consistent and easily recognizable feature of these junctions, approximately 20 nm in thickness. At about 50 nm from the plasma membranes are less dense inner plaques that appear to associate with the [intermediate filaments](#). Thus the total thickness of a desmosome from one inner plaque to the other is approximately 130 nm. The desmosomes may be thought of as the “scaffold couplings” that link the intermediate filament cytoskeleton throughout a tissue. Widely distributed, desmosomes are especially abundant in stratified epithelia, where strong intercellular adhesion is required to resist external friction. They are present in almost all epithelia, but also in cardiac muscle, the arachnoid and pia of the meninges, and the follicular dendritic cells of the lymphoid system.

**Figure 1.** An electron micrograph of a desmosome from mouse tongue epithelium, showing the midline (ml), the outer and inner plaques (op and ip), and the intermediate filament (if) or tonofilament.



Desmosomes are macromolecular complexes consisting of at least six interacting proteins (Fig. 2). Their adhesion molecules are the desmocollins and desmogleins, members of the [cadherin](#) family of calcium-dependent cell **adhesion molecules** (1, 2). These transmembrane proteins have their extracellular domains in the desmoglea and their cytoplasmic domains in the dense outer plaques. The outer plaque contains two proteins that are members of the armadillo family named after the *Drosophila* segment polarity signaling gene. These are plakoglobin (3), also known as g-catenin and present in other intercellular junctions, and plakophilin (4). The outer and inner plaques are bridged by the plakin family member desmoplakin (5), which mediates interaction between the adhesive domain and the intermediate filaments. The most recently discovered desmosomal components are the two plakin family members, envoplakin (6) and periplakin (7), which although not exclusively epidermal, appear to be involved in cornified envelope formation.

**Figure 2.** Diagram showing the major principle molecular components of the desmosome and their location within its structure. The diagram shows one-half of a desmosome from the midline (ml) in the intercellular space, through the plasma membrane (pm) and the outer plaque (op), to the inner plaque (ip) and the intermediate filaments. The numbers at the bottom represent distances in nanometers from the plasma membrane. The amino-terminus (n) and the carboxy-terminus (c) of each component is indicated. The positions of these have been mapped by immuno-gold labeling with domain-specific antibodies and electron microscopy (unpublished observations). The components are: DP, desmoplakin; Dsc a, desmocollin “a” form; Dsc b, desmocollin “b” form; Dsg, desmoglein; PG, plakoglobin; PP, plakophilin.



## 1. Desmosomal Glycoproteins

There is direct evidence that the desmocollins (Dsc) and desmogleins (Dsg) mediate desmosomal adhesion, probably by a heterophilic interaction (8). Adhesive interaction probably involves tripeptide sequences known as “cell adhesion recognition sites” near the amino termini of the glycoproteins (9). These molecules represent two subfamilies of the cadherin superfamily, each occurring as three different isoforms and the products of different genes. The molecules resemble classical cadherins in their extracellular domains, but show unique features cytoplasmically. Desmocollins have **alternatively spliced cytoplasmic domains** [the longer “a” form and a shorter “b” form (10)], and the desmogleins have long cytoplasmic domains containing unique 29-amino-acid-residue repeats (11). The extracellular domains of these glycoproteins form an ordered array, giving rise to the highly structured desmosomal midline region or desmoglea. Their cytoplasmic domains contribute to the structure of the cytoplasmic plaque. The precise roles of the different cytoplasmic domains are unclear. That of the desmocollin “a” form can support plaque assembly and intermediate filament attachment (12), but the functions of desmoglein and desmocollin “b” in plaque assembly have not been resolved. Both desmocollin and desmoglein have been shown to bind plakoglobin and plakophilin, and desmocollin “a” can also bind desmoplakin. Three isoforms of desmocollin and desmoglein show tissue-specific expression (13, 14). Dsc2 and Dsg2 are ubiquitous in all desmosome-bearing tissues, but Dsc1, Dsg1, Dsc3, and Dsg3 are largely restricted to stratified epithelia. Desmosomal glycoprotein genes are closely linked on chromosome 18q12, and this is of possible significance for regulation of the expression.

## 2. Plakoglobin

Plakoglobin is a member of the *armadillo* family of proteins. It is not exclusively desmosomal but also occurs in adherens or [intermediate junctions](#), even in non-desmosome-bearing tissues such as endothelia. It is, like b-catenin and armadillo itself, both a junctional protein and a member of the **wingless** Wnt signaling pathway (15). In its junctional role, it combines both with desmocollin and desmoglein, or with E-cadherin, where it is mutually exclusive with that of b-catenin. Null mutations

of desmoplakin in mice result in severe disruption of desmosomes in cardiac muscle and epidermis. In its signaling role, it has been shown to produce axis duplication when overexpressed in *Xenopus* embryos and combined with the adenomatous polyposis coli (APC) protein that regulates its cytoplasmic level. The cytoplasmic plakoglobin can bind to the LEF-1 [transcription factor](#) to enter the nucleus and regulate gene activity (16).

### 3. Plakophilin

This is an armadillo family protein that has important junctional properties and possibly a signaling function. It exists as two isoforms that show differential tissue specific expression, PP1 being predominantly expressed in stratified tissues (17) and PP2 ubiquitously (18). The first human genetic disease resulting from a desmosomal mutation was an epidermal dysplasia/skin fragility syndrome, effectively a PP1 null mutation (19). This defect resulted in loss of adhesion between epidermal keratinocytes and detachment of intermediate filaments from the cell periphery, indicating an important role for plakophilin in linking between the desmosomal adhesion molecules and the cytoskeleton. Plakophilin has been shown to bind to both desmosomal glycoproteins, desmoplakin and cytokeratin (20). Its potential signaling role is inferred from its armadillo-type structure and because it has been detected in cell nuclei.

### 4. Desmoplakin

This is a plakin family member, together with plectin, bullous pemphigoid antigen 1 (BAPG1), envoplakin, and periplakin (6, 7). Except for BAPG1 all have been reported to be associated with desmosomes. Desmoplakin (DP) is a ubiquitous desmosomal component. DP1 is a [coiled-coil](#) dimer with globular end **domains**. The isolated molecule has an overall length of 130 nm. DP2 lacks the central coiled rod domain and is unable to dimerize. Its length is 43 nm. Expression studies, supported by immunogold labeling and [electron microscopy](#) (see [Immunoelectron Microscopy](#)), show that the amino-terminal globular domain associates with the desmosomal plaque and the carboxy-terminal domain with the intermediate filaments (21, 22). Thus desmoplakin plays an important role in linking between the desmosomal plaque and the cytoskeleton. Although not restricted to epidermis, envoplakin and periplakin were only identified as components of the cornified envelopes of suprabasal keratinocytes (6, 7). They are essentially similar in structure to desmoplakin, and their precise role has yet to be defined. The role of plectin is more problematic, because [null mutations](#) have no apparent effect on desmosome structure and function.

Desmosomal glycoproteins are target [antigens](#) in the **autoimmune** blistering disease of pemphigus (23). Desmosomes appear at the 32- to 64-cell stage of embryonic development, where they are trophoderm specific (24).

### Bibliography

1. J. L. Holton et al. (1990) *J. Cell Sci.* **97**, 239–246.
2. P. J. Koch et al. (1990) *Eur. J. Cell Biol.* **53**, 1–12.
3. P. Cowin et al. (1986) *Cell* 1063–1073.
4. M. Hatzfeld et al. (1994) *J. Cell Sci.* **107**, 2259–2270.
5. H. Mueller and W. W. Franke (1983) *J. Mol. Biol.* **163**, 647–671.
6. C. Ruhrberg et al. (1996) *J. Cell Biol.* **134**, 715–729.
7. C. Ruhrberg et al. (1997) *J. Cell Biol.* **139**, 1835–1849.
8. N. A. Chitaev and S. M. Troyanovsky (1997) *J. Cell Biol.* **138**, 193–201.
9. C. Tselepis et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 8064–8069.
10. J. E. Collins et al. (1991) *J. Cell Biol.* **113**, 381–391.
11. G. N. Wheeler et al. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 4796–4800.
12. S. Troyanovsky et al. (1993) *Cell* **72**, 561–574.

13. P. K. Legan et al. (1994) *J. Cell Biol.* **126**, 507–518.
14. S. Schäfer, P. J. Koch, and W. W. Franke (1994) *Exp. Cell Res.* **211**, 391–399.
15. P. Cowin and B. Burke (1996) *Curr. Opin. Cell Biol.* **8**, 56–65.
16. O. Huber et al. (1996) *Mech. Dev.* **59**, 3–10.
17. H. W. Heid et al. (1994) *Differentiation* **58**, 113–131.
18. C. Mertens et al. (1996) *J. Cell Biol.* **135**, 1009–1025.
19. J. McGrath et al. (1997) *Nature Genet.* **17**, 240–244.
20. E. A. Smith and E. Fuchs (1998) *J. Cell Biol.* **141**, 1229–1241.
21. E. A. Bornslaeger et al. (1996) *J. Cell Biol.* **134**, 985–1001.
22. A. P. Kowalczyk et al. (1997) *J. Cell Biol.* **139**, 773–784.
23. J. R. Stanley (1995) In “Cell adhesion and human disease.” *Ciba Found. Symp.* **189**, 107–120.
24. T. P. Fleming et al. (1991) *Development* **112**, 527–529.

### Suggestions for Further Reading

25. D. R. Garrod et al. (1998) "Desmosomes". In *Adhesive Interactions of Cells* (D. R. Garrod, M. A. J. Chidgey, and A. J. North, eds.), JAI Press, Greenwich, CT, pp. 165–201. (A review of the latest desmosome literature and a good source of reference.)
26. K. J. Green and J. C. R. Jones (1996) *FASEB. J.* **10**, 871–881. (Another excellent review.)

## Detergents

### 1. Detergents as Amphiphilic Molecules

**Gene expression** studies in both **eukaryotes** and **prokaryotes** eventually require detecting and enriching the expression product, usually a [protein](#). When this gene product is a water-soluble, **cytosolic** protein, detecting and enriching it can begin by simply homogenizing the cells, followed by **centrifugally** separating the supernatant from the [membrane](#)-bound proteins that remain in the pellet. When the **gene** product is a plasma [membrane protein](#), however, the task is complicated because the protein has to be detached from the membrane at some stage. Although some loosely attached proteins are released by treatment with salt solutions, integral membrane proteins are released only upon treatment with amphiphilic reagents, also termed detergents .

By definition, amphiphiles harbor both [hydrophilic](#) and **hydrophobic** moieties. In a water-oil or water-lipid mixture, detergent molecules accumulate at the interface. The hydrophobic arm is inserted into oil or lipid, and the hydrophilic ([polar](#)) end is introduced into water, thus justifying the adjective “*amphiphilic*” used to describe this class of compounds. Depending on the degree of polarity, the hydrophilic end (the polar head group) of a detergent is either ionic or nonionic. Regardless of the difference in polar head groups, all detergents dissolve in water. The physicochemical properties of detergents dissolved in water have been the focus of vigorous biochemical research.

#### 1.1. Physiological Importance of Detergents

Amphiphilic substances or detergents are central to our existence as living organisms. About 80% of the [lipids](#) in the alveolar lining is dipalmitoyl phosphatidylcholine that acts as a detergent to confer the crucial nonatelectic property to the lungs. This unusual proportion of a particular phospholipid is important because lower levels of this lipid in the alveolus lead to atelectasis (collapsed alveolus) in

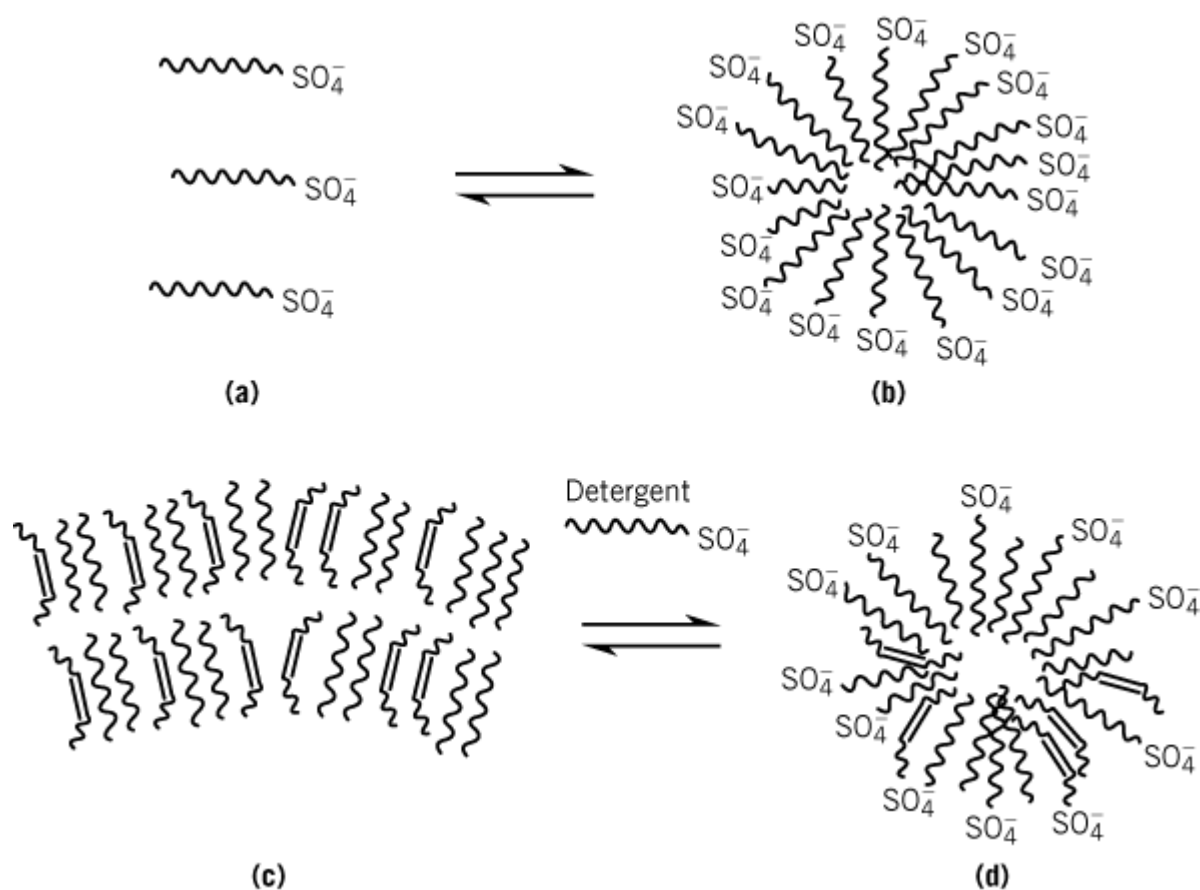
infants. Similarly, bile acids are crucial detergent molecules in the intestines. Bile acids emulsify water-insoluble fat molecules and thus help the intestinal lining to absorb them. The powerful detergents, lysophospholipids, produced in the intestine by phospholipase A<sub>2</sub>-mediated cleavage of the C<sub>2</sub> [fatty acid](#) chains of phospholipids, also aid in lipid digestion. In fact both bile acids and lysophospholipids aid in lipid absorption and also help in absorbing and removing the metabolic products of fats and lipids.

## 2. General Strategy for Obtaining Solubilized and Reconstituted Membrane Proteins

Solubilization of a membrane protein is a procedure in which the proteins and lipids, originally cradled in the membrane, are dissociated appropriately in a buffered detergent solution. This controlled dissociation of the membrane results in the forming small protein and lipid clusters that remain dissolved in the aqueous solution. Next the solubilized material is separated from any insoluble substance (aggregated proteins and lipids that are dissociated by the detergent) by [centrifugation](#). Then, the detergent is removed from the solubilized material because only after removing the detergent do the solubilized proteins and lipids reaggregate appropriately so that the activity of a functional protein is restored. Many of the functional membrane proteins retain their activity only when they are allowed to associate with other specific membrane proteins and lipids. Such reaggregation is inhibited in the presence of the dissociation-promoting detergent. Typical methods to remove detergents include (1) extensive [dialysis](#) of the detergent-solubilized fractions, whereby the detergent monomers are released from the dialysis bag; (2) **gel filtration**, which helps to retain the detergent monomers on the column while the high molecular weight, solubilized proteins flow through; or (3) [hydrophobic chromatography](#), whereby the hydrophobic detergent molecules are removed by absorption into a hydrophobic matrix. The success of each procedure depends on the physicochemical properties of the detergent, which determine the formation of detergent aggregates, termed micelles (Fig. 1). Therefore, it is important to identify and define some of the factors that regulate aggregation of detergent molecules. The [critical micelle concentration](#) (cmc) is the maximum concentration at which the detergent molecules remain as monomers. Removal of detergents by dialysis or gel filtration occurs through the monomer and is difficult when the cmc is very low.

**Figure 1.** Structures of detergents and micelles. At the critical micelle concentration (cmc), monomeric molecules (**a**), which are stable below cmc and removable by dialysis, aggregate to form a spherical or cylindrical macromolecular micelle (**b**), which is stable above the cmc and cannot be removed by dialysis. In a micelle, the hydrophobic ends of the detergent molecules are pointed inward and away from water whereas the polar ends (head groups) are pointed outward into the aqueous phase. In the presence of a plasma membrane bilayer containing saturated and unsaturated fatty acid containing phospholipids (**c**), mixed micelles (**d**) that contain detergent molecules and plasma membrane lipids and proteins are formed.





## 2.1. Different Classes of Detergents and Their Applicability in Solubilizing Membrane Proteins

As shown in Fig. 2, the various detergents can be structurally classified as follows:

1. Negatively charged, strongly ionic detergents: (a) a detergent with a long, flexible aliphatic chain: sodium dodecylsulfate (**SDS**) (cmc  $\sim$ 0.23% (w/v) or 8 mM); (b) a detergent with a rigid structure: sodium cholate (cmc  $\sim$ 0.43% or 10 mM). Sodium cholate (1%) has been used to solubilize the [adenylate cyclase](#)-G-protein complex (1), the dopamine D1 **receptor** (2) and many other receptors.
2. Detergents that form [hydrogen bonds](#): (a) detergents with flexible long-chain structures: octanoyl- *N*-methylglucamide (MEGA-8) (cmc  $\sim$ 1.9% or 60 mM), and a long chain polyoxyether, Thesit (cmc  $\sim$ 0.005% or 0.09 mM), which is structurally similar to Lubrol. Although 0.5% MEGA-8 has been used to solubilize intrinsic membrane proteins (3), a Thesit concentration of 1.2% has been used with membrane proteins, such as adenylate cyclase (4). (b) more polar but non-ionic detergents with sugar rings, *n*-dodecyl- $\beta$ -maltoside (cmc  $\sim$ 0.009% or 0.18 mM) and [n-octyl- \$\beta\$ -glucoside](#) (cmc  $\sim$ 0.7% or 23.2 mM), which have been used to solubilize the nicotinic [acetylcholine receptor](#) [1% *n*-octyl- $\beta$ -D-glucoside, (5)], opioid receptors [0.3% *n*-octyl- $\beta$ -D-glucoside, (6)], **GABA<sub>A</sub>** receptors [1% *n*-octyl- $\beta$ -D-glucoside, (7)], photoreceptor [guanylate cyclase](#) [1% *n*-dodecyl- $\beta$ -D-maltoside, (8)], and others.
3. Hydrophobic detergents with long polyoxyether chains and aromatic rings: **Triton X-100** (cmc  $\sim$ 0.013% or 0.2 mM) and Triton X-114 (cmc  $\sim$ 0.011% or 0.21 mM), which are used to solubilize the receptors for GABA [2% Tr X-100, (9)], prostacyclin (10), prolactin [1% Tr X-100, (11)], [transferrin](#) [1% Tr X-100, (12)], [insulin](#) [2% Tr X-100, (13)] and a neuronal growth cone protein, GAP-43 [1% Tr X-114, (14)].
4. **Zwitterionic** detergents: (a) a detergent with a flexible, long-chain structure: *N*-dodecyl-*N,N*-dimethyl-3-ammonio-1-propanesulfonate (in short, propanesulfonate) (cmc  $\sim$ 0.12% or 3.6 mM); (b) detergents with a rigid structure: CHAPS (cmc  $\sim$ 0.46% or 7.4 mM) and CHAPSO (cmc

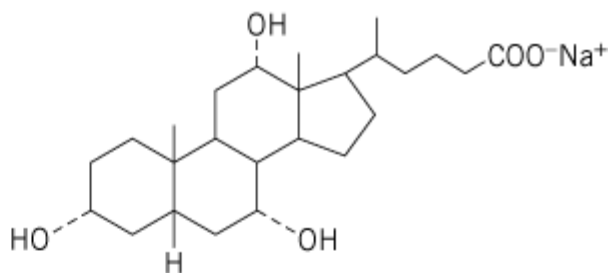
~0.5%, 8 mM), commonly used to solubilize receptors, such as those for [colony-stimulating factor](#) [0.5% CHAPS, (15)], opiates [0.6% CHAPS, (16)], neurotensin [0.6% CHAPS, (17)], **somatostatin** [0.6% CHAPS, (18)], adenosine A<sub>1</sub> [2.5% CHAPS, (19)] and human prostatic sex-hormone-binding globulin receptor [6.0% CHAPS, (20)].

5. Positively charged, strongly ionic detergents, for example, cetylpyridinium chloride.
6. Nonionic surfactants that form inclusion complexes, for example, [cyclodextrin](#).
7. Specific detergents, such as digitonin, used to solubilize opioid receptors [2% digitonin, (21)], neurotensin receptors [2% digitonin, (22)],  $\beta$ -adrenoreceptors (23), dihydropyridine Ca<sup>2+</sup> channels [1% digitonin, (24)], and many others.

**Figure 2.** Classification of detergents based on structural features.



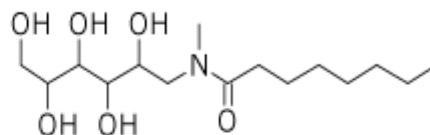
Sodium dodecyl sulfate (SDS)



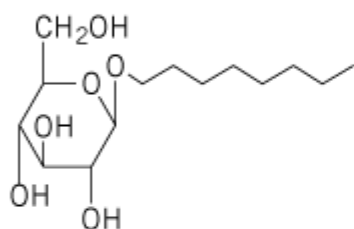
Sodium cholate (Sod. cholate)



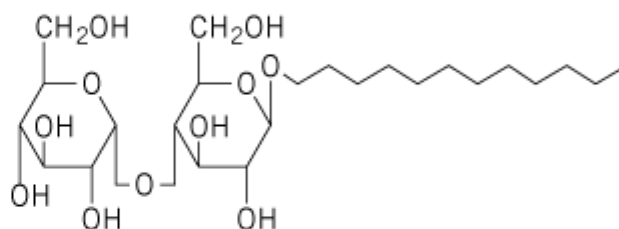
Thesit



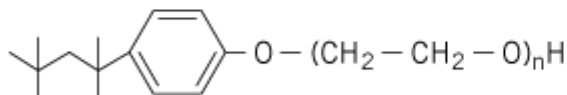
MEGA-8



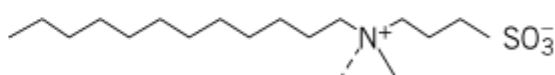
Octylglucoside (Octylglc)



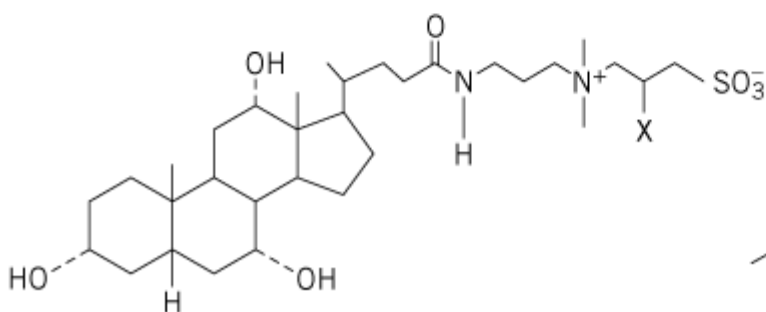
n-dodecyl- $\beta$ -D-maltoside (Dodmalt)



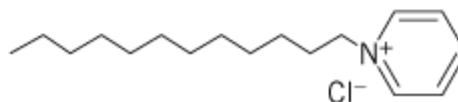
Triton X-114 (n=7) (TrX-114)  
Triton X-100 (n=10) (TrX-100)



N-Dodecyl-N, N-dimethyl-3-ammonio-1-propanesulfonate (Propansulfon)



3-[(3-Cholamidopropyl)dimethylammonio]-1-propanesulfonate(CHAPS) (X=H)  
3-[(3-Cholamidopropyl)dimethylammonio]-2-hydroxy-1-propanesulfonate(CHAPSO) (X=OH)



Cetylpyridiniumchloride (CPC)

Thus, the detergent concentrations used for solubilizing most biologically active receptor proteins are typically in the range of 0.5 to 2.0% (w/v), suggesting that for practical purposes the cmc is not

the deciding factor for choosing detergent concentration. Instead, it is obtained empirically by optimizing protein-solubilization conditions.

The other key factor in solubilizing biologically active membrane proteins is the ratio of detergent to protein maintained during solubilization. The commonly used ratios are 1:1 (16, 19) and 2 to 3:1 (4, 7, 8). Higher ratios, such as 10:1 (20) are used less often. Optimization experiments observed that the coextracted lipids are essential for reconstituting the solubilized, sheep brain serotonin 5-HT<sub>1A</sub> receptor (25), and the highest yield of the active 5-HT<sub>1A</sub> sites was obtained with 2% CHAPS at a detergent to protein ratio of 2:1 (26). This topic is discussed in more detail later.

## 2.2. Other Uses of Detergents

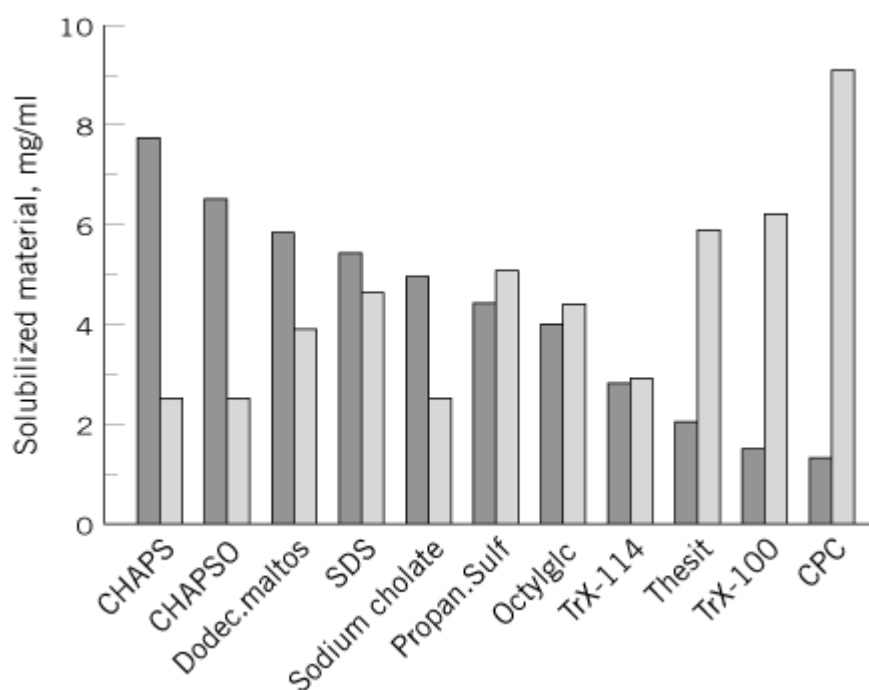
Detergents are also used in purifying DNA and RNA molecules. For example, SDS is used to denature polynucleotide molecules during the preparation of DNA, and a structurally related detergent, sodium salt of *N*-laurylsarcosine (Sarkosyl), is used in preparing RNA (27, 28). In either case, the detergent is used to denature the biopolymers and also to inhibit certain enzymes, such as **deoxyribonucleases** and **ribonucleases**, which degrade these molecules. Thus detergents are inseparable from scientific activity whether it is to solubilize lipids and proteins or to purify DNA and RNA.

## 2.3. Detergent Solubilization of Membrane Lipids

A survey of the literature (1-26) indicates that there is no general agreement on the choice of detergent for a single membrane protein, and the use of five or more different classes of detergents for a given receptor solubilization, usually at a concentration 1 to 2%, irrespective of the cmc, is quite common. This could account for the variable recovery of biological activity in different preparations, but it does not explain the differential efficacy of the different detergents in solubilizing active receptor proteins.

To resolve this question, the protein- and lipid-extracting abilities of each detergent in Fig. 2 were determined. Sufficient variability was observed to explain much of the bewildering variability in the recovery of biological activity reported in the literature (26). The receptor molecule used as a marker in this study was the sheep brain serotonin 5-HT<sub>1A</sub> receptor, which is a member of the receptor family with seven transmembrane regions, and may couple either to [adenylate cyclase](#) or [phospholipase C](#), depending on the type of cell in which it is expressed (26). Solubilized, high-affinity serotonin 5-HT<sub>1A</sub> receptor sites are reconstituted only by coextracted lipids, and adding separately isolated, pure or mixed lipids during or after solubilization is completely ineffective in reactivating an inactive preparation (25). Instead, adding lipids or lipid vesicles to already solubilized, reconstituted (in coextracted lipids) and active 5-HT<sub>1A</sub>-R preparations inhibits the **ligand-binding** activity of the receptor. Such findings indicated that coextraction of plasma membrane lipids is important in order to maintain the functionally active conformation of lipid-sensitive membrane proteins, such as the serotonin 1A receptor. So the efficiency of many detergents in solubilizing membrane lipids was studied. Considering the total amount of lipid solubilized by each detergent, detergents with long aliphatic chains that terminate in a negative charge (SDS), a zwitterionic sequence (CHAPS or CHAPSO), or a multihydroxylic-group-containing disaccharide (*n*-dodecyl- $\beta$ -D-maltoside) (DDM) clearly has greater ability to solubilize lipids, irrespective of the cmc of the detergent (Fig. 3). In contrast to detergents, such as sodium cholate or sodium [deoxycholate](#), which are also efficient in solubilizing membrane lipids, CHAPS, CHAPSO and DDM solubilized active 5-HT<sub>1A</sub> receptors. These results indicate that CHAPS, CHAPSO, and DDM, which are relatively milder and nondenaturing, allow the 5-HT<sub>1A</sub>-R-membrane lipid complex to remain undisturbed during solubilization, thus allowing recovery of the active 5-HT<sub>1A</sub>-R following removal of these detergents (25, 26).

**Figure 3.** Concentration of proteins and lipids (mg/mL) in the detergent-solubilized and reconstituted vesicles, compared to the solubilization of high-affinity 5-HT<sub>1A</sub> binding sites.



Sodium cholate has a high cmc of approximately 8 mM (or 0.5%) and is easily dialyzed like CHAPS or CHAPSO (both of which have a cmc of ~0.5%). However, although it forms vesicles (because it coextracts sufficient lipid to give a similar lipid to protein ratio of ~2:1) (Table 1 and Fig. 3), the 5-HT<sub>1A</sub> receptors obtained in the reconstituted preparations are inactive. Possible explanations include nonextraction or denaturation of the protein because of the strongly ionic nature of sodium cholate. Octylglucoside presents a different picture. As a mild, nonionic, dialyzable (cmc ~0.7%) and very mildly denaturing detergent, it is often used to extract biologically active proteins (5-7). However, its inability to coextract sufficient lipid (a lipid to protein ratio of 1:1.1) probably accounts for the low amount of solubilized 5-HT<sub>1A</sub> receptor activity, which contains less than 15% of the lipids in comparable DDM supernatants.

**Table 1. Distribution of Lipids, Protein, and Serotonin 5HT<sub>1A</sub> Receptor Binding Activity in the Vesicles Obtained after Detergent Removal Followed by Dialysis and Ultracentrifugation**

| Detergents       | Lipid, mg/mL | Protein, mg/mL | PE, % | PC, % | PS, % | PI, % | SM, % | GalCers, % | C, % | 5HT <sub>1A</sub> |
|------------------|--------------|----------------|-------|-------|-------|-------|-------|------------|------|-------------------|
|                  |              |                |       |       |       |       |       |            |      | Receptor Binding  |
| CHAPS            | 7.7          | 2.5            | 40.0  | 17.6  | 21.0  | 10.3  | 0.0   | 2.8        | 7.7  | 5500.0            |
| CHAPSO           | 6.5          | 2.5            | 27.9  | 16.2  | 34.0  | 10.0  | 0.7   | 2.7        | 8.5  | 5300.0            |
| Dodecylmaltoside | 5.8          | 3.9            | 18.3  | 14.8  | 31.3  | 5.3   | 5.9   | 6.5        | 17.7 | 3200.0            |
| SDS              | 5.4          | 4.6            | 22.0  | 28.7  | 13.0  | 6.8   | 5.0   | 11.2       | 14.0 | 0.0               |

|                  |     |     |      |      |      |      |     |      |      |       |
|------------------|-----|-----|------|------|------|------|-----|------|------|-------|
| Sodium cholate   | 4.9 | 2.5 | 18.2 | 17.0 | 27.9 | 6.0  | 8.5 | 13.4 | 8.5  | 360.0 |
| Propanesulfonate | 4.4 | 5.1 | 22.8 | 29.8 | 19.2 | 17.2 | 5.1 | 8.2  | 10.6 | 60.0  |
| Octylglucoside   | 4.0 | 4.4 | 25.0 | 20.0 | 24.4 | 6.5  | 1.6 | 8.7  | 13.6 | 600.0 |
| Triton X-114     | 2.8 | 2.9 | 32.5 | 24.2 | 14.3 | 7.6  | 1.5 | 3.8  | 15.9 | 90.0  |
| Thesit           | 2.0 | 5.9 | 14.8 | 37.0 | 13.3 | 5.2  | 3.0 | 12.6 | 14.1 | 155.0 |
| Triton X-100     | 1.5 | 6.2 | 20.5 | 27.7 | 11.6 | 9.8  | 2.7 | 9.8  | 17.9 | 80.0  |
| CPC              | 1.3 | 9.1 | 25.0 | 21.0 | 22.6 | 5.7  | 5.7 | 9.6  | 10.7 | 0.0   |

The detergents that extract large amounts of protein but are at the lower end of the lipid extraction scale (such as Thesit, Triton X-114, Triton X-100, and CPC) cannot solubilize enough lipid to produce vesicles and do not solubilize active 5-HT<sub>1A</sub> sites. Therefore, such G-protein-associated receptors require an environment containing particular lipids (enriched in PE, PS and PI) to express their biological properties fully.

The relative amounts of individual groups of lipids solubilized by the different classes of detergents are surprisingly diverse (Table 1). The inverse profiles of PC and PS are intriguing in view of earlier work showing that PC [and GalCer, sulfated GalCer (Su) and SM] are enriched in the outer membrane domain, whereas PS (and PE and PI) are enriched in the inner membrane domain (29-31). From the profile for SM it also appears that this lipid is probably localized in the PC-containing domains. This agrees well with the idea that choline phospholipids are concentrated in the outer leaflet of the membrane bilayer. It is hard to explain this asymmetrical extraction of lipids in terms other than specific detergent binding by the protein. Thus for eight of the twelve detergents tested, it was found that the solubilization-enrichment profile of PI is similar to that of PE, indicating a similar, inner membrane localization of PE, PI, and PS. PS binds to protein kinase C on the cytosolic face of the plasma membrane and results in activating this enzyme (31). This is consistent with the idea that the PI family is the source of internally released [second messengers](#) (32). The serotonin 5-HT<sub>1A</sub> receptor polypeptide has seven transmembrane **alpha-helices** and three cytosolic loops (26), which could be more strongly associated with inner membrane lipids. Differences in protein folding could explain why a strongly denaturing detergent, such as sodium cholate, is useful in solubilizing and reconstituting the dopamine D1 receptor [which can be reconstituted by adding separately isolated lipids (2)] but is ineffective in solubilizing and reconstituting the 5-HT<sub>1A</sub> receptors.

### 3. Summary

The importance of detergents in the normal functions of a living body is irrefutable. In scientific research, detergents provide a powerful means of undertaking subcellular fractionation of macromolecules, such as proteins, RNA, and DNA. Although gene expression studies do not directly begin with protein purification, the expressed protein is finally purified or detected through biochemical methods, such as column [chromatography](#), **Western blotting** analysis, or [enzyme](#) assays. If the expressed protein is **lysosomal** or **cytosolic**, simply homogenizing the cells in a hypotonic buffer, followed by separation of the soluble proteins by centrifugation, yields a supernatant that contains the required protein. However, if the expressed protein is, for example, a **mitochondrial** membrane protein, detergent solubilization of this protein is required before its assay. Although in earlier studies such choices of detergents were made empirically, the data presented here provide a stronger foundation for judicious analysis of the efficacy and applicability of many of the commonly used detergents. Thus, for immunoblot analysis, one might use Triton X-100 or Nonidet P-40. To solubilize proteins that do not require membrane lipids for acquiring functional activity, one might use either the Triton family of detergents, the bile acid detergents (eg, sodium cholate), or octylglucoside. On the other hand, if the protein is heavily lipid-dependent, one should test CHAPS,

CHAPSO, or DDM before using other detergents. Therefore, this practical overview of the structural classification of detergents with appropriate reference to their applicability to protein and lipid solubilization and to reconstitution of functional activity could be helpful in scientific research involving expression of specific genes.

#### 4. Acknowledgment

The author wishes to thank Dr. Michael Hogan, Neuroscience Program, for carefully reading the manuscript and for his helpful suggestions.

#### Bibliography

1. P. C. Sternweis and A. G. Gilman (1979) *J. Biol. Chem.* **254**, 3333–3340.
2. A. Sidhu (1990) *J. Biol. Chem.* **265**, 10065–10072.
3. J. E. K. Hildreth (1982) *Biochem. J.*, **207**, 363–366.
4. A. K. Keenan, A. Gal, and A. Levitski (1982) *Biochem. Biophys. Res. Commun.* **105**, 615–623.
5. J. M. Gonzalez-Ros, A. Paraschos, M. C. Farach, and M. Martinez-Carrion (1981) *Biochim. Biophys. Acta.* **643**, 407–420.
6. T. Fujioka, F. Inoue, S. Sumita, and M. Kuriyama (1988) *Biochem. Biophys. Res. Commun.* **156**, 54–60.
7. S. M. J. Dunn, R. A. Shelman, and M. W. Agey (1989) *Biochemistry* **28**, 2551–2557.
8. K. W. Koch (1991) *J. Biol. Chem.* **266**, 8634–8637.
9. T. N. Sato and J. H. Neal (1989) *J. Neurochem.* **52**, 1114–1122.
10. A. K. Dutta-Roy and A. K. Sinha (1987) *J. Biol. Chem.* **262**, 12685–12691.
11. H. Okamura, S. Raguette, A. Bell, J. Gagnon, and P. A. Kelly (1989) *J. Biol. Chem.* **264**, 5904–5911.
12. A. P. Turkewitz, J. F. Amatruda, D. Borkani, S. C. Harris, and A. L. Schwartz (1987) *J. Biol. Chem.* **263**, 8318–8325.
13. Y. F. Yamaguchi and J. T. Harmon (1988) *Biochemistry* **27**, 3252–3260.
14. J. H. P. Skene and I. Virag (1989) *J. Cell Biol.* **108**, 613–624.
15. R. Fukunaga, E. Ishizaka-Ikeda, and S. Nagata (1990) *J. Biol. Chem.* **265**, 14008–14015.
16. E. A. Frey, M. E. Gosse, and T. E. Cote (1989) *Eur. J. Pharmacol.* **172**, 347–356.
17. J. Mazella, J. Chabry, and J. P. Vincent (1989) *J. Biol. Chem.* **264**, 5559–5563.
18. H. T. He, K. Johnson, K. Thermos, and T. Reisine (1989) *Proc. Natl. Acad. Sci. USA* **86**, 1480–1484.
19. R. Munshi and J. Linden (1989) *J. Biol. Chem.* **264**, 14853–14859.
20. D. J. Hryb, M. S. Khan, N. A. Romas, and W. Rosner (1989) *J. Biol. Chem.* **264**, 5378–5383.
21. Y. H. Wong, C. D. Demoliou-Mason, and E. A. Barnard (1989) *J. Neurochem.* **52**, 999–1009.
22. A. Mills, C. D. Demoliou-Mason, and E. A. Barnard (1988) *J. Biol. Chem.* **263**, 13–16.
23. R. A. Cerione et al. (1983) *Proc. Natl. Acad. Sci. USA* **80**, 4899–4903.
24. U. Kanngiesseer, P. Nalik, and O. Pongs (1988) *Proc. Natl. Acad. Sci. USA* **85**, 2969–2973.
25. P. Banerjee, J. T. Buse, and G. Dawson (1990) *Biochim. Biophys. Acta.* **1044**, 305–314.
26. P. Banerjee, J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
27. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning*, 2nd ed. Cold Spring Harbor Laboratory Press, Cold Spring, NY, section 1.21.
28. P. Chomczynski and N. Sacchi (1987) *Anal. Biochem.* **162**, 156–159.
29. L. Freysz et al. (1982) In *Phospholipids in the Nervous System*, Vol. 1 (Metabolism) (L. Horrocks et al., eds.), Raven Press, New York pp. 37–47.
30. R. N. Fontaine, R. A. Harris, and F. Schroeder (1980) *J. Neurochem.* **34**, 269–277.

31. R. N. Fontaine, R. A. Harris, and F. Schroeder (1979) *Life Sci.* **24**, 295–400.
32. M. J. Berridge (1984) *Biochem. J.* **220**, 345–360.

### Suggestions for Further Reading

33. J. Lasch (1995) *Biochim. Biophys. Acta.* **1241**, 269–292. Interaction of detergent with lipid vesicles.
34. R. J. Vonk, M. Kalivianakis, D. M. Minich, C. M. Bijleveld, and H. J. Verkade (1997) *Scandinavian J. Gastroenterol.-supplement*, **222**, 65–67. The metabolic importance of unabsorbed dietary lipids in the colon.
35. G. D. Henry and B. D. Sykes (1994) *Methods Enzymol.* **239**, 515–535. Methods to study membrane protein structure in solution.
36. A. G. DeOliveira and H. Chaimovich (1993) *J. Pharm. Pharmacol.* **45**, 850–861. Effect of detergents and other amphiphiles on the stability of pharmaceutical drugs.
37. A. C. Newby, A. Chrambach, and E. M. Bailyes (1982) *Tech. Lipid Membrane Biochem.* **B409**, 1–22. The choice of detergents for molecular characterization and purification of native intrinsic membrane proteins.

## Development

Development is the process of change in an organism with time. Such change can be either morphological (such as changes in shape or structure), biochemical (such as the synthesis of specialized proteins), or both. Development usually includes all of the changes that occur during the life of an organism. These include [fertilization](#), embryogenesis, maturation, and gametogenesis.

Embryogenesis begins with fertilization of the **gametes** to form the diploid [zygote](#). The zygote undergoes cleavage divisions to form the blastoderm [embryo](#). Changes in cell shape and cell movements at **gastrulation** lead to the formation of the three germ layers: the ectoderm, endoderm, and mesoderm. The process of organogenesis forms the embryonic organs.

### Suggestions for Further Reading

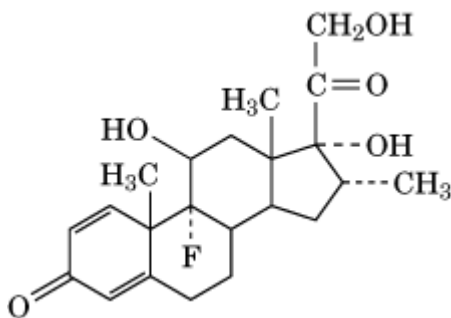
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert, R. Beddington, J. Brockes, T. Jessell, P. Lawrence, and E. Meyerowitz (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
- J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, U.K.
- A. S. Wilkins (1993) *Genetic Analysis of Animal Development*, 2nd edition, Wiley-Liss, New York.

## Dexamethasone



Dexamethasone (9a-fluoro-16a-methylprednisolone) (Fig. 1) is a semisynthetic [glucocorticoid](#) hormone demonstrating a greater binding affinity for the glucocorticoid **receptors**, a greater potency of action, and a longer half-life, than the natural glucocorticoids. The hormone is used as a drug and in research. Its molecular mode of action is the same as that of the natural glucocorticoids—that is, by binding to the glucocorticoid receptor and activating gene [transcription](#).

**Figure 1.** Structure of dexamethasone



## Dg-Dc Tailing

This procedure is used for [cloning](#) cDNA ([complementary DNA](#)) prepared from [messenger RNA](#) (1). Single-stranded cDNA prepared by **reverse transcriptase** is incubated with the enzyme terminal deoxyribonucleotide transferase and dCTP, which adds 5 to 10 C residues to the 3' ends of the DNA molecules (C-tail). The second strands are synthesized with reverse transcriptase and oligo-dT as the primer, using the [poly A](#) tails of mRNA, and the double-stranded cDNA molecules are tailed with dC as before. The resulting molecules are cloned into **vector** molecules that have been tailed with poly G (about 10 bases per end) to create complementary ends to which the dC tails on the inserts hybridize. A recent modification is to use CTP instead of dCTP, which gives more uniform tails (2).

## Bibliography

1. H. Land, M. Grez, H. Hauser, W. Linden Maier, and G. Scholtz (1981) *Nucleic Acids Res.* **9**, 2251–2266.
2. W.M. Schmidt and M.W. Mueller (1996) *Nucleic Acids Res.* **24**, 1789–1791.

## Diafiltration

**Ultrafiltration** and microfiltration are able to prepare reasonably concentrated samples of [macromolecules](#), but they permit only a limited separation of the retained solutes from the smaller or

more permeable components present in the solution needing to be filtered. In order to remove such smaller components more efficiently from the retained species, the feed solution can be washed in a process known as *diafiltration*. In this process, water or an appropriate buffer solution is added to the retentate during filtration, with the membrane-permeating species being removed as the excess fluid or buffer is filtered. Such a process may be undertaken in two different modes: continuous, in which the wash fluid is added throughout filtration, or discontinuous, in which the feed solution is first concentrated by a predetermined amount using standard ultrafiltration techniques, the wash fluid is then added, and the process repeated any number of times until the final solute concentration is achieved. When using continuous diafiltration, the amount of solute remaining in the retentate may be determined from a mass balance and expressed as

$$\frac{C_R}{C_0} \exp(-N_D S) \quad (1)$$

where  $C_R$  is the concentration of the retentate,  $C_0$  the initial concentration,  $N_D$  is the number of diavolumes, defined as the total volume of the wash fluid divided by the initial volume of the solution, and  $S$  is the solute sieving coefficient.

For discontinuous diafiltration, it is given by

$$\frac{C_R}{C_0} = \left(\frac{V}{V_0}\right)^{-(n+1)} S \quad (2)$$

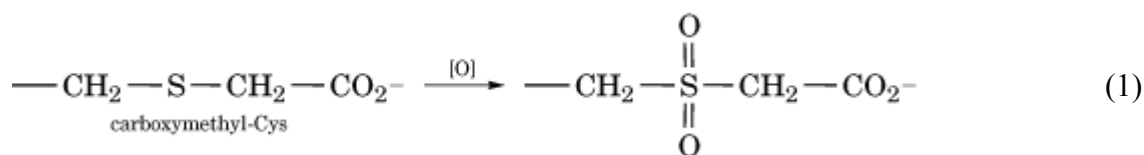
where  $n$  corresponds to the number of discontinuous diafiltration procedures.

## Diagonal Methods

It is often desirable to purify selectively those [peptides](#) from a [protein](#) that contain a particular [amino acid](#), such as all the peptides containing **cysteine residues**. This is accomplished most elegantly by so-called “diagonal” techniques (1), which rely on a change in the properties of peptides caused by the selective modification of the amino acid residues of interest. If a mixture of peptides, generated by **proteolysis** of a protein for [peptide mapping](#), is separated in one dimension and then subjected to the same procedure a second time, but at right angles to the first, all the peptides will lie on a diagonal because they had the same mobilities in both directions. The separation can be by [electrophoresis](#) or [chromatography](#), but it should be carried out on a two-dimensional medium, such as paper or thin-layer plates. To detect certain peptides, and to cause them to have a different mobility in the second dimension, the entire mixture of peptides is treated after the first separation so as to modify all the residues of a particular type and to change their separation properties. After the second dimension separation, the peptides containing these residues consequently lie off the diagonal of all the other peptides and are readily identified and isolated.

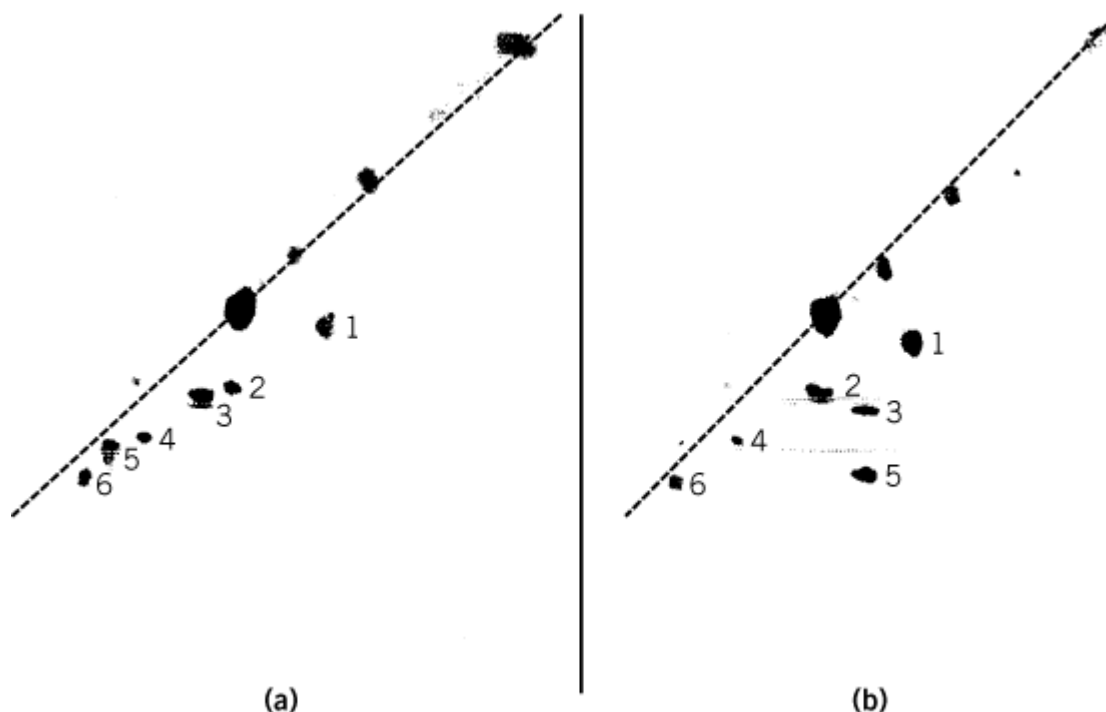
### 1. Isolating Peptides Containing Certain Amino Acids

For example, peptides containing Cys residues can be purified by blocking their [thiol groups](#) by reaction with **iodoacetic acid** (Fig. 1a) (2). A two-dimensional diagonal peptide map is prepared by electrophoresis at pH 3.5, exposing the peptides to performic acid between the two separations. This oxidizes the sulfur atom of the modified Cys residues to the sulfone:



which lowers the  $pK_a$  value of the carboxyl group, so that peptides containing Cys residues are somewhat more acidic in the second dimension; they then lie to one side of the diagonal. Alternatively, the thiol groups can be reacted initially with *N*-ethylmaleimide (see [Thiol Groups](#)); the *N*-ethylsuccinimide group introduced can then be hydrolyzed to the succinamic acid by treatment with ammonia after the first separation (3).

**Figure 1.** Isolation and identification of cysteine residues linked by disulfide bonds or modified by reaction with iodoacetate (2). The proteins were digested with trypsin, followed by chymotrypsin, and the resulting peptides were separated by paper electrophoresis at pH 3.5 in the horizontal direction (anode at the left). After exposure to performic acid vapor, the electrophoresis was repeated in the vertical direction (anode at bottom). The peptides were visualized by staining with [ninhydrin](#). The diagonal indicated by the dashed line is defined by peptides that were not altered by the performic acid and consequently had the same mobility in both dimensions; many of these peptides have migrated off the paper. (a) Reduced [BPTI](#) (bovine pancreatic trypsin inhibitor) in which the six Cys residues have been blocked by reaction with iodoacetate. The six peptides containing these residues define a second diagonal, below the first, because they were slightly more acidic in the second dimension due to their oxidation to the sulfone. These peptides are identified by the numbers of the Cys residues they contain. The amino acid residues they contain are, starting from lower left to upper right, 54 to 58 (Cys55), 27 to 33 (Cys30), 5 to 15 (Cys5 and 14), 47 to 53 (Cys51), 1 to 15 (Cys5 and 14), and 36 to 39 (Cys38). (b) The major one-disulfide intermediate in the BPTI disulfide folding pathway (see [Protein Folding In Vitro](#)) that had been trapped by reacting all the four free Cys thiol groups with iodoacetate. The two peptides containing Cys30 and Cys51 are absent from the second diagonal of modified Cys residues and had the same mobility in the first dimension, indicating that they were originally linked by an intramolecular disulfide bond between these two Cys residues.



Analogous techniques have been developed for a few other residues, using appropriate chemical modifications specific for these residues. The [amino groups](#) of **Lys**-containing peptides are initially blocked with trifluoroacetyl or maleyl groups, by treatment with the corresponding [anhydride](#); these

groups alter the charge and are easily removed after the first electrophoresis by acidification. **Met** residues are first alkylated with **iodoacetamide** to produce the charged sulfonium derivative. After the first electrophoresis, heating causes this derivative to cleave the polypeptide chain by a reaction analogous to that occurring in cleavage by cyanogen bromide. **His** residues are first blocked by dinitrophenylation of the side chain and then regenerated after the first electrophoresis by exposure of the peptides to a thiol-containing reagent. **Arg** residues are first blocked by reaction with cyclohexanedione and then regenerated at alkaline pH. **Trp** residues are modified by *o*-nitrophenylsulfenyl chloride after the first separation (4).

## 2. Characterizing Disulfide Bonds

Diagonal techniques were initially introduced by Brown and Hartley (1) for the determination of which pairs of Cys residues in a protein were linked by **disulfide bonds**. The protein is prepared for peptide mapping and the peptides are separated in the first dimension, under conditions where the disulfide bonds are stable and cannot rearrange or be reduced; the pairs of peptides linked by disulfide bonds consequently migrate together (Fig. 1b). The disulfide bonds are then cleaved by exposure to performic acid, which converts each such Cys residue to cysteic acid:



The peptides originally linked by the disulfide bond are now independent and more acidic. During separation in the second dimension, they migrate differently and lie off the diagonal, but they can be related by their common mobility in the second dimension.

Diagonal techniques are best suited to two-dimensional separations, such as paper electrophoresis and chromatography, which were used extensively in the past but have now been supplanted by more modern techniques that require less material, but operate in a single dimension, such as **HPLC**. The same basic technique can be used, but the fractions from the first separation must be separated again individually.

### Bibliography

1. J. R. Brown and B. S. Hartley (1966) *Biochem. J.* **101**, 214–228.
2. T. E. Creighton (1974) *J. Mol. Biol.* **87**, 603–624.
3. H. Gehring and P. Christen (1980) *Anal. Biochem.* **107**, 35–361.
4. T. Sasagawa et al. (1983) *Anal. Biochem.* **134**, 224–229.

## Dialysis

*Dialysis* is the transport of a solute through a membrane under the influence of a difference in the concentrations (or activities) of the solute in the solutions separated by the membrane. The separation of solutes is induced by differences in their transport by **diffusion** through the membrane matrix. A diffusible low-molecular-weight solute can thus be separated from a solution containing nondiffusible large **macromolecules**. The rate of dialysis is directly proportional to the membrane area and inversely proportional to the membrane thickness. Since the 1960s, the principal clinical use of this technique has been the treatment of patients with kidney disease.

Nonclinical applications of dialysis techniques for the separation of solutes with molecular weights

ranging from 100 to 100,00 Da may be on a small laboratory scale or with larger-scale commercial units for the recovery of alkali or in food processing applications. Dialysis is also important in biotechnology, as it provides a technique whereby products can be separated from fragile shear-sensitive or heat-sensitive solutions. A number of variants of the technique exist. These may be divided into batch and continuous processes. The simplest batch process is one in which a semipermeable membrane tubing made from cellulose and sealed at one or both ends contains the solution of interest. The tubing is suspended in a buffer solution or water. Small molecules contained within the solution bounded by the membrane diffuse into the surrounding buffer solution. The transfer of the substance across the membrane obeys Fick's law, if we assume that chemical equilibrium is established between the membrane and solution. If no water transfer occurs, and the mass transfer coefficient is constant along the diffusion path, the amount of substance transferred through the membrane  $W$  may be expressed as

$$W = DA\Delta C \quad (1)$$

where  $D$  is the overall dialysis coefficient of the substance,  $A$  the membrane area, and  $\Delta C$  the logarithmic mean bulk concentration difference across the membrane based on the initial and final values.

A variation of equilibrium dialysis is *reverse dialysis*, in which the fluid contained in the semipermeable membrane tube is placed into a water-soluble polymer, such as polyethylene glycol. This causes the water to diffuse out of the membrane tube and concentrates the macromolecule within the tube.

For batch dialysis, where mass transfer is not at steady state, the transfer rate is given by

$$kt = -m \log_e \left( \frac{C_0 - C}{C_0} \right) \quad (2)$$

where  $k$  is the dialysis rate constant,  $t$  the time,  $m$  the volume of the feed ratio to diffusate,  $C_0$  the concentration of the feed, and  $C$  the concentration of the dialysate.

Equilibration of the low-molecular-weight solutes or electrolytes by this technique can be speeded up by ensuring that the area of membrane contact is large relative to the volume of solution being dialysed. Since any molecule from the solution held within the membrane has to pass through the stagnant fluid film or boundary layer on either side of the membrane, as well as across the membrane itself, the overall dialysis rate may be enhanced by the minimization of the stagnant layers. This can be achieved by the introduction of a stirrer into the system, or by frequent changes of the dialysis fluid.

The transfer rate is also dependent on temperature. Biological materials such as proteins usually need to be dialyzed in the cold to minimize bacterial growth or protein **denaturation**. Under such conditions, the transfer of molecules will be reduced compared with that at room temperature.

If the solution requiring to be dialyzed is available in a sufficiently large quantity, continuous methods that use a *countercurrent flow* configuration to maintain a large difference in the solute concentration across the membrane may be used. For such systems, which are similar to those used in the treatment of renal failure, the amount of solute transferred across the module, if we assume that the membrane is the primary barrier to diffusion, may be expressed as

$$Q(C_{in} - C_{out}) = kA\Delta C \quad (3)$$

where  $Q$  is the flow rate of the solute undergoing dialysis,  $C_{in}$  and  $C_{out}$  are the concentrations at the

inlet and outlet,  $k$  and  $A$  are as previously defined, and  $DC$  is the log mean concentration difference between the dialysis fluid and the fluid being dialyzed.

Many substances **bind** to proteins or plasma. The solute transfer referred to above relates to an unbound solute. However, a bound solute will be in equilibrium with the unbound solute, the extent of the binding being governed by the solute **association constant**. The effect of this will be that there will be an overestimate of the concentration driving force present. As free solute is removed, there will be unbinding of the solute with the effect of understating the concentration driving force.

In a single-pass situation, ie, where the dialyzing fluid flows to waste and the solute undergoing dialysis is recirculated, the solute exchange rate within the module used for dialysis is expressed in terms of dialysance  $D$ , defined as

$$D = \frac{\text{Change in concentration of the incoming solute}}{\text{Concentration driving force}} \quad (4)$$

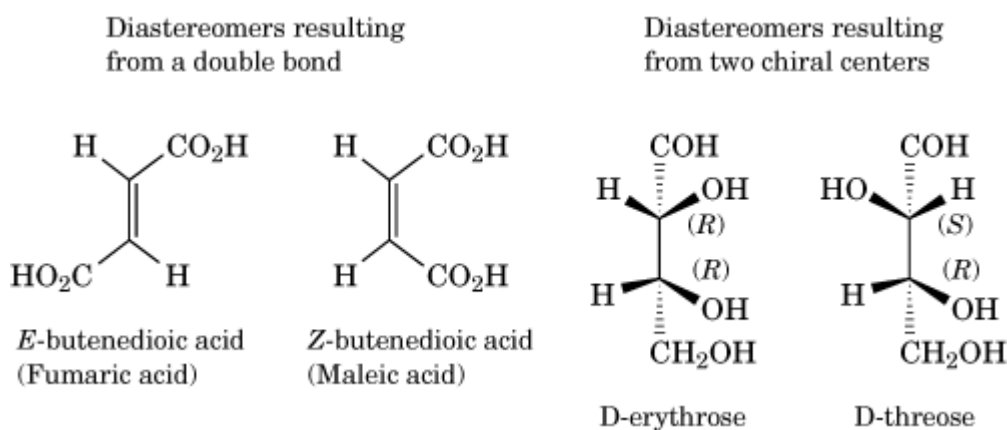
The generalized concentration driving force will be determined by the Donnan factor  $a$ , defined as the ratio of the ionic concentrations in the dialysis fluid and in the solute undergoing dialysis at equilibrium. If the solute is noncharged,  $a = 1$ , and the measurement of the transport rates may occur by the use of a single pass flow configuration on the dialyzate side, while the fluid being dialyzed is recirculated. In such a flow configuration, provided that the volume being dialyzed  $V$  is unchanged, the concentration at time  $t$ ,  $C_t$ , may be expressed in terms of the removal rate  $k$  and the initial concentration  $C_0$  by

$$C_t = C_0 \exp - \left( \frac{kt}{V} \right) \quad (5)$$

## Diastereomer

Diastereomers are [stereoisomers](#), molecules that differ only in their spatial arrangement of atoms, that are not [enantiomers](#) (1). Diastereomers have different chemical properties. There are several common types of diastereomers. Molecules containing carbon-carbon and carbon-nitrogen double bonds will exist as two different geometric diastereomers, if both of the double-bonded atoms have two different substituents. These diastereomers are differentiated with the designation **cis** and **trans** and more formally as *Z* and *E* (from the German *zusammen* for together and *entgegen* for opposite). The *Z* isomer describes the orientation when the two substituents with the higher priority are on the same side of the double bond. In potentially ambiguous cases, the priority of the substituents of each double-bonded atom is determined by the Cahn, Ingold, and Prelog rules (2) (see [Configuration](#)). The other common form of diastereomers occurs when a molecule contains more than one **chiral** center. If two stereoisomers contain at least one corresponding chiral center with the same configuration and one corresponding chiral center with the opposite configuration, the molecules will be diastereomers. These two forms of diastereomers are shown in [Figure 1](#).

**Figure 1.** The two common types of diastereomers are depicted, *cis/trans* geometric isomers and molecules containing more than a single chiral center.



Two atoms are referred to as *diastereotopic* if, on substitution with a different isotope, they would generate different diastereomers ([1](#), [3](#)) (see [Prochiral](#)). Diastereotopic atoms are chemically different and, importantly, can be distinguished by [NMR](#).

#### Bibliography

1. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), Vol. **3**, Elsevier, Amsterdam, pp. 1–48.
2. R. S. Cahn, C. K. Ingold, and V. Prelog (1966) *Angew. Chem. Int. Ed.* **5**, 385–415.
3. D. Arigoni and E. L. Eliel (1969) In *Stereochemistry* (E. L. Eliel and N. L. Allinger, eds.), Vol. **4**, Wiley-Interscience, New York, pp. 127–243.

#### Suggestions for Further Reading

4. J. March (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, pp. 104–119.
5. K. Mislow (1966) *Introduction to Stereochemistry*, W. A. Benjamin, New York, Chap. "2".
6. K. Mislow and M. Raban (1967) In *Topics in Stereochemistry* (N. L. Allinger and E. L. Eliel, eds.), Vol. **1**, Wiley-Interscience, New York, pp. 1–38.

## Dicentric Chromosome

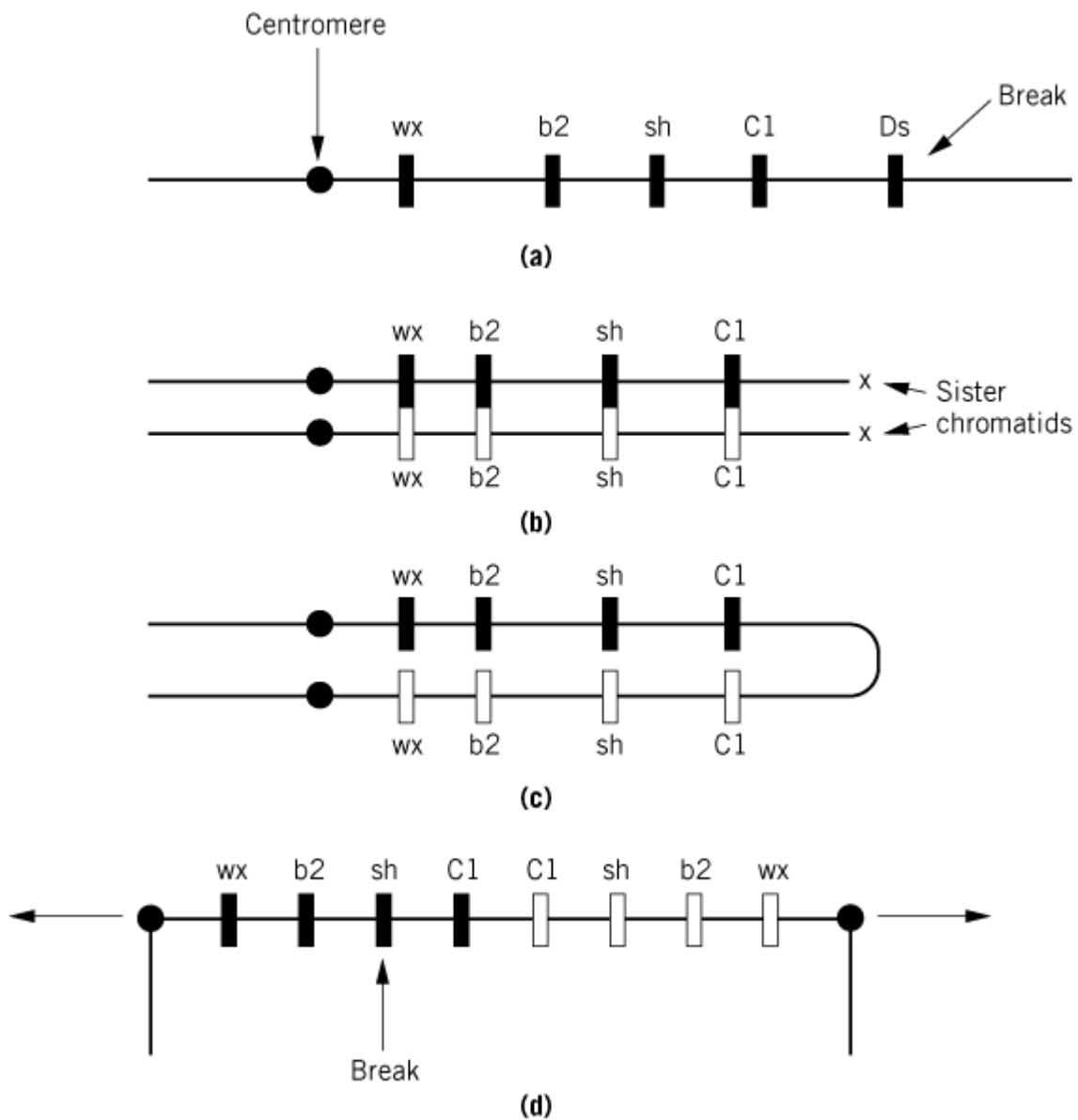
Dicentric chromosomes have two [centromeres](#). Such chromosomes are products of chromosomal damage, breakage, and fusion. Dicentric chromosomes normally break at **anaphase**, when the centromeres have separate attachments to the spindle and are subject to forces pulling them apart.

The definition of a chromosome breakage-fusion-bridge cycle in **maize** by McClintock in 1951 was an important step in defining chromosomal dynamics ([1](#)). The origin of chromosomal breakage in this case in maize is the [transposable element](#) Dissociation (*Ds*), which causes chromosomal instability in the presence of another element called Activator (*As*). *Ds* is present at specific sites in the chromosome and transposes to a different site only in the presence of *As*. Events such as these led to the discovery of transposable elements. Two forms of this cycle have been described resulting in the appearance of either dicentric chromosomes or [chromatids](#).

One type of cycle is shown in Figure [1](#). Several convenient markers are present on chromosome 9 of

maize, which include colored (C1), shrunken (sh1), bronze (bz1), and waxy (wx1). Each of the two chromosomes has a *Ds* element located at the site indicated (*Ds*). When the *Ds* element is activated, chromosome breakage may occur. The [acentric fragment](#) without a centromere is lost. The ends of the two sister chromatids fuse to form a bridge, so the resulting single DNA molecule contains two centromeres. At the anaphase, the dicentric fused chromosomal arms break, so that segments of the chromosome are either duplicated or deficient in daughter cells. This cycle occurs either at **meiosis** or **mitosis**.

**Figure 1.** Chromosome breakage-fusion-bridge cycle in maize. (a) A break in maize chromosome 9 at the site of the transposable element *Ds* is indicated. (b) At the end of the S-phase and in G2, two sister chromatids are formed that have a break at this site (sticky ends). (c) The sticky ends of the two chromatids fuse to form a chromosome containing two centromeres—a dicentric chromosome. (d) At the anaphase, the two centromeres may be forced to move in opposite directions, leading to chromosomal breakage and forming two chromosomes, one containing a short duplication and the other a deficiency.



## Bibliography

1. B. McClintock (1951) Cold Spring Harbor Symp. Quant. Biol. **16**, 13–47.



## Dictyostelium

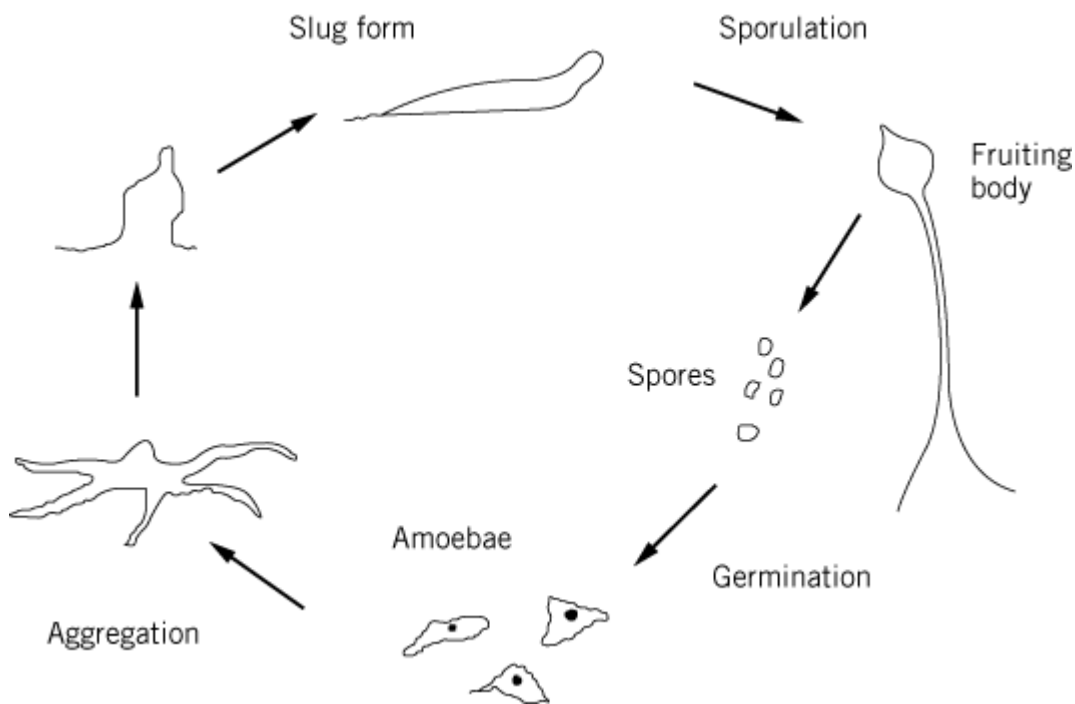
The cellular slime mold *Dictyostelium discoideum* has been studied in the laboratory extensively as a model eukaryotic organism. This soil amoeba has provided extraordinary opportunities for molecular biologists and biochemists to examine a wide variety of fundamental cellular functions.

Developmental biologists have characterized at the molecular level *Dictyostelium's* unique **developmental** pathway. Single free-living cells can respond to a number of chemical signals, aggregate, and then differentiate into a complex multicellular organism.

A powerful library of techniques has been developed for manipulating this organism in the laboratory. Large numbers of cells can be grown, and proteins can be purified for biochemical characterization and, in some cases, structural characterization. The recent utilization of peptide [proteinase inhibitors](#) have made this increasingly easy. Genetic techniques have been developed for making [mutations](#), and sophisticated methods have been applied for analyzing them. Genetic selection can be carried out with a variety of metabolic and [drug resistance](#) markers. **Plasmids** that either integrate into the [chromosome](#) or replicate extrachromosomally have been constructed from fragments of naturally occurring ones, and a variety of **promoters** are available for obtaining heterologous gene expression. Cell biological techniques are also well developed, yielding information on expression and localization in fixed and live cells.

A simplified schematic of the *Dictyostelium discoideum* developmental pathway is shown in Figure [1](#). The free-living amoebae feed on bacteria, and they can be found outside the laboratory in forest soils and detritus. In the laboratory, they can be grown in liquid media in shaking flasks or in petri plates, either in liquid media or on agar surfaces. Strains of *Dictyostelium* have been developed that can grow in synthetic and fully characterized media. On surfaces, the amoebae are motile and can use a variety of [chemotaxis](#) signals to find food and to optimize growth conditions ([1](#)).

**Figure 1.** A schematic diagram of the *Dictyostelium* developmental cycle.



Under certain conditions, such as the onset of starvation, a population of amoebae can suspend their cell cycle and begin a highly regulated set of gene inductions, which initiates the developmental pathway. In the laboratory, this can be done synchronously. Early in this process, a series of **new enzymes** are synthesized that are responsible for the synthesis, detection, and degradation of **cyclic AMP** (cAMP). These enzymes are used to regulate the chemotactic response of the cells in such a way that the cells aggregate. A founder cell initiates the aggregation process by secreting pulses of cAMP. Surrounding cells, if they are competent to respond, do so in two ways. First, they migrate toward the cAMP source. Second, they secrete a cAMP pulse themselves, which serves to amplify the signal and to increase the area over which it can travel. An extracellular phosphodiesterase degrades the cAMP between the waves.

As the cells migrate, they form cell–cell contacts that serve to establish polarity signals, which are used as the aggregate forms a mound and the cells differentiate into different types and sort into different locations. The tip of the mound becomes an organizing center which, based on environmental factors, decides whether to proceed in development to a fruiting body or to search for more favorable conditions by becoming a slug. If the organism becomes a slug, the mound tips over on the surface, and the outside cells secrete an extracellular matrix, commonly known as slime, which is left behind on the surface after the organism has passed. The tip of the slug controls its migration direction chemotactically, as well as by detecting gradients of temperature and light. Migration rates can exceed 1 cm/h.

When the mound or slug is ready to proceed to form a fruiting body, the cells have differentiated, through a pattern of differential gene expression, into either of two distinct cell types: pre-stalk or pre-spore. As the pre-stalk cells migrate to the bottom of the forming structure, they differentiate terminally into a plant-like cell with a rigid cell wall. A chlorinated alkyl phenone, DIF-1, is released to induce stalk cell differentiation. Pre-spore cells differentiate into small environmentally resistant spores, which are carried up by the forming stalk. The spore sac is quite sticky, which facilitates its use of insect legs as transport in the wild. When conditions are optimal, the spores can germinate and begin dividing once again as amoebae.

The haploid **genome** of *Dictyostelium* is about  $54 \times 10^6$  base pairs in six **chromosomes**, and a

detailed map has been assembled (2). The DNA has an unusually high content of A and T nucleotides. Complete characterization of the genome is underway. A wide variety of genetic tools have been developed for manipulation of *Dictyostelium*. Mutations can be created, and the resulting **phenotypes** can be analyzed in many ways. Naturally occurring extrachromosomal plasmids have been isolated and reconstructed into a family of vectors useful in stable or transient expression of heterologous genes. Multiple selectable markers are available, allowing the construction of strains in which [complementation](#) can be used to analyze [protein–protein interactions](#). Homologous [recombination](#) has been shown to occur, allowing gene knockouts and gene replacements to be made.

*Dictyostelium* has been particularly useful in the characterization of the fundamental aspects of transmembrane signaling and cell–cell communication in the regulation of motility (3). The organism has the ability to detect and respond to a whole series of environmental factors, such as light, temperature, and moisture. It also can detect a wide variety of chemical signals, from sugars to hormones. *Dictyostelium* uses a variety of mechanisms of intracellular signaling mechanisms, including proton gradients and electrochemical gradients, as well as  $\text{Ca}^{2+}$  waves (see [Calcium Signaling](#)). Analysis of a variety of protein kinases and phosphatases and their targets have illuminated a number of regulatory pathways. The dynamic cytoskeleton of *Dictyostelium* has been extensively studied, particularly actin, actin-binding proteins, and the myosins, as well as tubulin, tubulin-binding proteins, and microtubule-based motor proteins. Major rearrangements of the cytoskeleton are required during cytokinesis and cell motility. Analysis of genetic knockouts of individual genes and combinations of those for cytoskeletal proteins have been particularly useful in attempting to determine their cellular roles. Targeted gene disruptions have also been used to characterize the pathways of endocytosis and exocytosis. *Dictyostelium* feeds on bacteria using phagocytosis that is highly regulated and specific (4).

The highly regulated developmental pathway of *Dictyostelium* has been extremely valuable in terms of understanding gene induction and regulation. The appearance and disappearance of specific messenger RNAs have been mapped to specific time points during development, and many of the gene products have been characterized and their function determined. The appearance of different cell types in the aggregate, along with their sorting to yield the final fruiting body, has given insights into positional information and the cell cycle in morphogenesis.

Study of the cellular slime mold *Dictyostelium discoideum*, as a model eukaryotic organism, has produced many insights into basic cellular function. It will continue to make contributions as its many features are probed in more detail.

#### Bibliography

1. S. J. McRobbie (1986) Chemotaxis and cell motility in the cellular slime molds. *Crit. Rev. Microbiol.* **13**, 335–375.
2. W. F. Loomis, D. Welker, J. Huges, D. Meghakian, and A. Kuspa (1995) *Genetics* **141**, 147–157.
3. P. Devroetes (1989) *Dictyostelium discoideum*: a model system for cell–cell interactions in development. *Science* **245**, 1054–1058.
4. G. Vogel (1983) *Dictyostelium discoideum* as a model system to study recognition mechanisms in phagocytosis. *Methods Enzymol.* **98**, 431–430.

#### Suggestions for Further Reading

5. J. T. Bonner (1967) *The Cellular Slime Molds*, 2nd ed., Princeton University Press, Princeton, NJ.
6. K. B. Raper (1984) *The Dictyostelids*, Princeton University Press, Princeton, NJ.
7. F. Loomis, ed. (1982) *The Development of Dictyostelium discoideum*, Academic Press, New York.

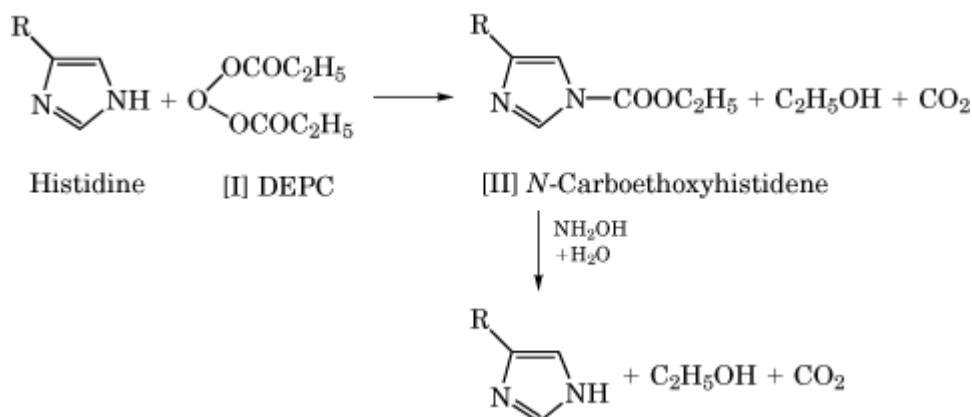
8. J. A. Spudich, ed. (1987) *Dictyostelium discoideum: Molecular Approaches to Cell Biology. Methods in Cell Biology*, Vol. **28** (L. Wilson series, ed.), Academic Press, Orlando, FL.
9. R. Mutzel (1995) *Experientia* **51**, 1103–1196. A fine collection of multi-author reviews.

## Diethylpyrocarbonate

Diethylpyrocarbonate (DEPC) [I] is an [anhydride](#) of ethoxyformic acid and is also called ethoxyformic anhydride (EFA), diethyl dicarbonate, or diethyl oxydiformate. DEPC is a liquid (density = 1.29 g/ml) with molecular weight of 162.14. It is not stable in water and hydrolyzes to two molecules of ethanol and carbon dioxide. DEPC is a good reagent for [chemical modification](#) of [histidine](#) residues. Because DEPC inhibits **ribonuclease A**, which has histidine residues at its [active site](#), DEPC is widely employed to protect **RNA** from enzymatic attack. DEPC is a suspected carcinogen, and therefore, it should be handled carefully.

Within the pH range 6 to 8, DEPC reacts with histidine residues to give *N*-carboethoxyhistidine [II], ethanol, and carbon dioxide ([Scheme 1](#)):

Diethylpyrocarbonate



The product, *N*-carboethoxyhistidine, absorbs UV at 240 nm, and the reaction can be followed spectroscopically ( $\text{De}_{240} = 3,200 \text{ M}^{-1} \text{ cm}^{-1}$ ) (1). The product is not very stable. The **half-life** of its decay is 55 h at pH 7 and much shorter at acidic or alkaline pH values. The product is regenerated to histidine by NH<sub>2</sub>OH (**hydroxylamine**). DEPC sometimes also reacts with [amino groups](#) and with the side chains of [cysteine](#), [tyrosine](#), [arginine](#), and [tryptophan](#) residues.

To prevent the action of ribonuclease, 0.2 ml DEPC is added to 100 ml of the solution that is being treated for use with RNA. To inactivate the remaining DEPC, the solution is autoclaved (2).

Bibliography

1. E.W. Miles (1977) *Methods Enzymol.* **47**, 431–442.
2. M. Gilman (1995) In *Current Protocols of Molecular Biology*, Wiley, New York, pp. 4.1.4–6.

## Difference Fourier

The difference Fourier is an extremely useful electron density map in protein [X-ray crystallography](#). The location of reagents such as inhibitors or water molecules, attached to or removed from the native structure (see **Ligand binding**) is easily determined if the native crystal structure is known. It is not necessary to perform a completely new structure determination for the derivative, so long as the new crystal is isomorphous with the old one, as in [isomorphous replacement](#). Only the intensities of the reflections in the diffraction pattern of the derivative need be collected. A normal electron density map is calculated by the equation:

$$\rho(xyz) = \frac{1}{V} \sum_h \sum_k \sum_\ell |F(hk\ell)| \times \exp[-2\pi i(hx + ky + \ell z) + i\alpha(hk\ell)]$$

If the  $|F(hk\ell)|$ 's in this equation are replaced by  $\Delta |F(hk\ell)| = |F(hk\ell)|_{\text{PH}} - |F(hk\ell)|_{\text{P}}$ , where P stands for the native protein and PH for the derivative, then the resulting difference electron density map shows the difference in electron density between native and the derivative. The density is positive for attachment and negative for removal of atoms ([1](#)):

$$\rho(xyz) = \frac{1}{V} \sum_h \sum_k \sum_\ell \Delta |F(hk\ell)| \times \exp[-2\pi i(hx + ky + \ell z) + i\alpha_{\text{P}}(hk\ell)]$$

In calculating this difference Fourier, the phase angles are those of the native protein  $\alpha_{\text{P}}(hk\ell)$ .

Because of approximations, the heights of the peaks in a difference Fourier are only half of what they would be in a normal Fourier. Nevertheless, it is a powerful method, and the attachment or removal of a few electrons can easily be detected. It also plays a useful role in locating of heavy atoms in the isomorphous replacement method. For example, after the main heavy atom site has been found from a difference [Patterson map](#), additional weakly occupied sites can be detected with a difference Fourier.

## Bibliography

1. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Difference Spectroscopy

The absorbance of a chromophore depends on its molecular environment (see [Absorption Spectroscopy](#)). It is sensitive to changes in the solvent around an exposed chromophore (see [Solvent Perturbation Spectroscopy](#)) or in the exposure of a chromophore as a consequence of conformational transitions, such as the unfolding of a protein or the melting of double-helical DNA (see **DNA melting**). These changes in absorbance are sometimes small, so it is of advantage to measure directly the difference in absorbance between two samples in which the only difference is that the macromolecule is in different environments. This can be done in a double-beam spectrophotometer, where the two samples are placed in the sample and the reference cuvette. This is called difference spectroscopy.

For the measurement of difference spectra, it is of advantage to use an instrument with a high sensitivity, so that small differences between the two samples in the different environments can be monitored reliably and with good accuracy. The development of microprocessor-controlled instruments has greatly facilitated the measurement of difference spectra. The measured data can be stored and manipulated later, eg, by changing the wavelength and absorbance scales, smoothing the data, addition and subtraction of spectra, and the calculation of derivative spectra. Also, difference spectra can be compared with other data already in the computer memory. Changes in spectra during a reaction (kinetic difference spectra) can be measured very well by using a diode array spectrophotometer, because spectra can be recorded at a very short interval.

### Suggestions for Further Reading

- S. B. Brown (1980) "Ultraviolet and visible spectroscopy". In *An Introduction to Spectroscopy for Biochemists* (S. B. Brown, ed.), Academic Press, London, pp. 14–69.
- J. W. Donovan (1973) Ultraviolet difference spectroscopy—new techniques and applications. *Methods Enzymol.* **27**, 497–525.
- R. K. Poole and C. L. Bashford (1987) "Spectra". In *Spectrophotometry and Spectrofluorimetry: A Practical Approach* (D. A. Harris and C. L. Bashford, eds.), IRL Press, Oxford, UK, pp. 23–48.
- F. X. Schmid (1997) "Optical spectroscopy to characterize protein conformation and conformational changes". In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 261–297.

## Differentiation

Differentiation is the expression of specialized functions by a cell ([1](#)). Differentiation may be reversible. For example, bacteria express sets of [enzymes](#) to facilitate the use of the nutrients that they sense in their surroundings. When the bacteria sense a change in their surroundings, such as a change in the nutrients in their culture medium, they inactivate the set of enzymes appropriate for the first medium and express a different set of enzymes to facilitate use of the new nutrients. As their surroundings change, the bacteria can change back and forth between different sets of enzymes indefinitely.

In most cases, however, differentiation is not reversible. Once a cell begins expressing specialized functions, its future choices become limited. For example, human cells that begin expressing heart-

specific proteins no longer have the potential to form skin. As cells become more specialized, they find their potential fates more restricted. The restriction of fates is called *determination*, and it was often shown to occur long before differentiation of any specialized functions. As our molecular tools have become more sophisticated, the distinction between determination and differentiation has become more blurred. Looking for the first detectable signs of differentiation, more sensitive techniques allow us to find molecular differences at earlier stages. When pushed to the limits, molecular techniques have increasingly shown that the basis of determination is the expression of a specialized gene product or function, that is, merely the first step in differentiation.

### Bibliography

1. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 468.

### Suggestions for Further Reading

2. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 468.
3. S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
4. L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
5. A. S. Wilkins (1993) *Genetic Analysis of Animal Development*, 2nd ed., Wiley-Liss, New York.

## Diffusion

Diffusion is the movement of molecules within a liquid driven by thermal fluctuations. It is important for most biological processes; indeed, the diffusion together of two reactants can be the limiting factor for many processes, when the reaction is said to be **diffusion-controlled**.

There are two types of diffusional motion: (1) *translational diffusion*, where a molecule moves its relative position within a three-dimensional system; and (2) *rotational diffusion*, where a molecule spins or rotates about one or more of its axes. The rates at which molecules undergo both types of diffusion are measured by their respective diffusion coefficients,  $D_t$  and  $D_r$ . The SI unit for  $D_t$  is the “Fick,” or  $m^2/s$ , but for historical and other reasons, biologists tend to use the cgs (Centigrade–gram–second) system unit of  $cm^2/s$ . These measurements can give important information about the sizes, structures, and physical properties of the molecules. At room temperature and in dilute solution, a small protein of molecular weight approximately 20,000 will have  $D_t$  of about  $10^{-6} cm^2/s$ ; a large virus between  $10^{-7}$  to  $10^{-8} cm^2/s$ ; a bacterial spore about  $10^{-9} cm^2/s$ . Some representative values are given in Table 1.

**Table 1. Translational Diffusion Coefficients and Derived Parameters of Some Molecules, Macromolecules and Biomolecular Assemblies**

| Substance | Molecular Weight | $10^7 \times D_{20,w}^{j0}$ , $cm^2/s$ | $r_H$ , $\text{Å}$ | $10^8 \times f$ , $g/s$ |
|-----------|------------------|--|--------------------|-------------------------|
|-----------|------------------|--|--------------------|-------------------------|

|                            |                   |      |      |      |
|----------------------------|-------------------|------|------|------|
| Water                      | 18                | 230  | —    | —    |
| Sucrose                    | 342               | 46.0 | 4.7  | 0.88 |
| Ribonuclease               | 13,700            | 11.1 | 19.3 | 3.64 |
| Ovalbumin                  | 45,000            | 7.8  | 27.5 | 5.18 |
| Fibrinogen                 | 330,000           | 2.0  | 107  | 20.2 |
| Dynein <sup>a</sup>        | $2.5 \times 10^6$ | 1.1* | 195  | 36.7 |
| Turnip yellow mosaic virus | $5.7 \times 10^6$ | 1.4  | 152  | 28.9 |

<sup>a</sup> Although the molecular weight of dynein is smaller than that of turnip yellow mosaic virus, its diffusion coefficient is smaller because it is more asymmetric. Values of  $D_{20,w}^{j0}$  and molecular weight  $M_r$  for dynein are also strongly dependent on salt concentration.

## 1. Translational Diffusion

The translational diffusion coefficient,  $D_t$ , describes the tendency of a molecule to move (translational motion) under the influence of either (1) a concentration gradient or (2) Brownian motion.

The movement of molecules in a gradient in which their concentration varies,  $dc/dx$ , where  $c$  is the concentration (in grams per milliliter) at each point  $x$ , is given by Fick's first law:

$$J = -D_t(dc/dx) \quad (1)$$

where  $J$  is the mass of particles crossing a 1-cm<sup>2</sup> cross section per second.

The same  $D_t$  characterizes the Brownian diffusion of the molecule:

$$\langle x^2 \rangle = 2D_t t \quad (2)$$

where  $t$  is the time and  $\langle x^2 \rangle$  is the average of the square of the distance the particle has moved.

The value of  $D_t$  depends not only on the intrinsic size and shape of the molecule but also on the viscosity and temperature of the medium in which it is suspended. The value of  $D_t$  must therefore be normalized to standard conditions; the standard conditions normally used are those of water at 20.0° C. The  $D_t$  corresponding to these conditions is normally designated  $D_{20,w}$ . It can be calculated from the value actually measured at absolute temperature  $T$  and in buffer "b,"  $D_{T,b}$ , with the equation

$$D_{20,w} = \left( \frac{293.15}{T} \right) \left( \frac{\eta_{T,b}}{0.01} \right) D_{T,b} \quad (3)$$

where  $\eta_{T,b}$  is the viscosity (in cgs units) of the buffer at temperature  $T$  and 0.01 is the viscosity of water at 20.0°C.

The value of  $D_{20,w}$  can also depend on the concentration of the molecule in the solution, due to



thermodynamic nonideality resulting primarily from the finite size of the molecule and its electric charge. This nonideality is represented (1) by the parameter  $k_d$  (in units of milliliters per gram) in the equation

$$D_{20,w} = D_{20,w}^0(1 + k_d c) \quad (4)$$

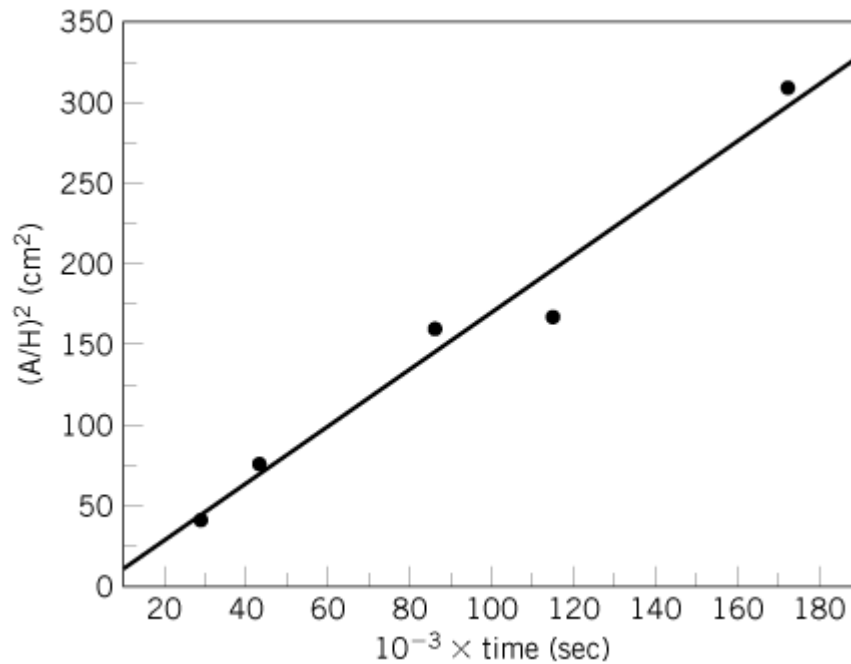
$D_{20,w}^0$  is the value of  $D_{20,w}$  at zero concentration of the molecule; its value can be estimated with equation 4 by measuring  $D_{20,w}$  at several different concentrations. For diffusion measurements, the non-ideality term is often negligible and can be neglected.

The value of  $D_t$  can be measured either by following the rate of disappearance of a concentration gradient or by dynamic light scattering. The former is the traditional method of measuring  $D_t$  (2, 3). The principle is to determine how the concentration changes with time at the boundary between two solutions, one containing the molecule of interest and the other, just the buffer or dialyzate. A special diffusion cell can be used with an appropriate optical system for recording either the concentration,  $c$ , or the concentration gradient,  $dc/dx$ , as a function of distance,  $x$ , from the center of the boundary. Alternatively, the optical system of an analytical **ultracentrifuge** can be used to measure the boundary spreading in a synthetic boundary cell. A speed is chosen to be sufficiently small that sedimentation of the molecule is negligible. The Schlieren optical system is very useful, as it gives the value of  $dc/dx$  as a function of  $x$ . The value of  $D_t$  can be determined by measuring the variation in the height of the boundary, specifically, the maximum value of  $dc/dx$ , with time,  $t$ :

$$\left(\frac{A}{H}\right)^2 = 4\pi D_t t \quad (5)$$

where  $A$  is the area under the curve of  $dc/dx$  versus  $x$ . A plot of  $(A/H)^2$  versus time will yield  $D_t$  from the slope; an example is given in Figure 1.

**Figure 1.** Measurement of the translational diffusion coefficient,  $D_p$ , for [ovalbumin](#) by the boundary spreading technique (4). As in Equation (5) of the text, the square of the ratio of area ( $A$ ) to the height ( $H$ ) of the boundary is plotted as a function of time.



Equation 5 assumes that there is no loss of material from the boundary, that is, that the area  $A$  remains constant. This assumption is reasonable for homogeneous protein preparations but may not be valid for polydisperse materials such as mucus **glycoproteins** and polysaccharides.

These methods have also generally been superseded by [dynamic light scattering](#). Nevertheless, the classical measurements are preferable with nonglobular macromolecules, with asymmetric shapes.

## 2. Interpretation of the Translational Diffusion Coefficient

The value of  $D_{20,w}$  or  $D_{20,w}^{\circ}$  obtained by either dynamic light scattering or boundary spreading can be used to provide a number of useful characteristics about a macromolecule.

### 2.1. Hydrodynamic Radius

The simplest deduction that can be made from the  $D_{20,w}^{\circ}$  is the size of the molecule as represented by its equivalent hydrodynamic radius,  $r_H$ , which is also known as the *Stokes radius*. The  $r_H$  is the radius of the equivalent sphere that would have the same  $D_{20,w}^{\circ}$ . The two parameters are related by the *Stokes–Einstein equation*:

$$r_H = \frac{k_B T}{6\pi\eta_{20,w} D_{20,w}^{\circ}} \quad (6)$$

where  $k_B$  is the Boltzmann constant ( $1.379 \times 10^{-16}$  erg/K),  $T$  is the temperature (293.15 K at 20.0°C), and  $\eta_{20,w}$  is the viscosity of water (0.01 P at 20.0°C). If the value of  $D_i$  is not corrected to standard conditions, the appropriate values for  $T$  and  $\eta$  must be used. Table 1 gives the values of  $r_H$  for several macromolecules.

### 2.2. Frictional Coefficient, $f$

The **frictional coefficient** is inversely related to the diffusion coefficient, giving a measure of the resistance to movement due to both its size and shape. It can be calculated directly from  $D_{20,w}^{\circ}$ :

$$f = \frac{RT}{N_A D_{20,w}^{\circ}} = \frac{k_B T}{D_{20,w}^{\circ}} \quad (7)$$

where  $R$  is the gas constant ( $8.314 \times 10^{-7}$  erg mol<sup>-1</sup> K<sup>-1</sup>) and  $N_A$  is Avogadro's number ( $6.02 \times 10^{23}$ ). Values of  $f$  calculated from  $D_{20,w}^{\circ}$  are given in Table 1.

It is often more informative to use the **frictional ratio**,  $f/f_0$ , which is the dimensionless ratio of the observed frictional coefficient to that of an equivalent spherical molecule of the same anhydrous mass and density.

### 2.3. Molecular Weight ( $M_r$ )

The molecular weight of a molecule can be calculated from the combination of its diffusion coefficient and its **sedimentation coefficient**,  $S_{20,w}^{\circ}$  when corrected to standard conditions. The equation analogous to (7) for the sedimentation coefficient is

$$f = \frac{M_r(1 - \bar{v}\rho_{20,w})}{N_A s_{20,w}^{\circ}} \quad (8)$$

where  $\bar{v}$  is the **partial specific volume** of the macromolecule and  $\rho_{20,w}$  is the density of the standard solvent, water at 20°C. Elimination of  $f$  between equations 7 and 8 yields the well-known Svedberg equation:

$$M_r = \frac{s_{20,w}^{\circ}}{D_{20,w}^{\circ}} \frac{RT}{(1 - \bar{v}\rho_{20,w})} \quad (9)$$

Equation 9, of course, makes it possible to calculate  $D_{20,w}^{\circ}$  if both  $S_{20,w}^{\circ}$  and  $M_r$  are known.

It is possible, in principle, to measure both  $D_{20,w}^{\circ}$  and  $S_{20,w}^{\circ}$  simultaneously from analysis of the shape of the boundary in **sedimentation velocity** analytical ultracentrifugation, although in practice this requires data of high quality and a totally homogeneous sample.

If the general shape of the macromolecule is known to a first approximation, it is possible to estimate its molecular weight from only its diffusion coefficient, just as in the case of the **sedimentation coefficient**. The power-law relation between  $M_r$  and  $D_t$  is also known as a *Mark-Houwink-Kuhn-Sakurada relation*:

$$D_t = K M_r^{-e} \quad (10)$$

where  $e = 0.333$  for a sphere, 0.85 for a rod, and 0.5–0.6 for [random coil](#) polymers. The appropriate value of the constant  $K$  is obtained from a collection of standard molecules of known  $D_t$  and  $M_r$  (5). Of course, the shape of the macromolecule will normally be known only approximately, so any molecular weight calculated in this way must be considered only an estimate.

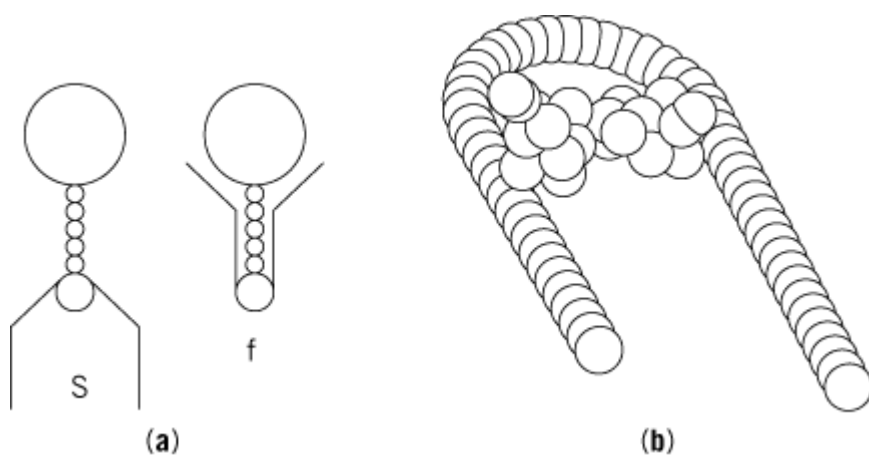
### 2.4. Shape

The translational frictional coefficient,  $f$ , or the frictional ratio,  $f/f_0$ , can be used directly to provide information about the shape of the molecule. The function defining the shape and flexibility of a macromolecule is the *Perrin translational frictional function*,  $P$ :

$$P = \frac{f}{f_0} \left( \frac{\bar{v}}{\bar{v} + (\delta/\rho_0)} \right)^{1/3} \quad (11)$$

where  $d$  is the amount of bound solvent, relative to the mass of the molecule, and  $r_0$  is the density of the bound solvent. For a molecule that is fairly rigid over a time-averaged period, the gross shape can be specified using  $P$  in terms of the axial ratio of the equivalent hydrodynamic ellipsoid or in terms of an arrangement of spheres, in hydrodynamic bead models (Figure 2). Computer programs are available for both types of modeling (6, 7). Especially with the latter approach, the diffusion coefficient should be used in conjunction with other hydrodynamic measurements, such as  $s_{20,w}$ , to obtain a unique solution to the modeling.

**Figure 2.** Use of diffusion measurements to represent the shapes of biomolecular assemblies using the hydrodynamic bead model approach: (a) from translational diffusion measurements: bead models of the slow “s” (tails out) and fast “f” (tails retracted) forms of a T-even bacteriophage (19); (b) from rotational (electric dichroism) relaxation measurements: bead models of the complex formed from a DNA fragment of 203 bp with two specific sites and two catabolite activator protein dimers (12).



### 3. Rotational Diffusion

Rotational diffusion coefficients, which result from the tumbling motion of a macromolecule about an axis or axes, are very sensitive functions of the shape of the molecule and permit the inference of some parameters of the overall shape. Unfortunately, this sensitivity comes at a severe price: measurement of rotational diffusion is usually considerably more difficult than translational diffusion.

There are two principal methods for measuring rotational diffusion. One is based on fluorescence measurements and is called *fluorescence anisotropy* decay. The other is based on electrooptical measurements and is termed *electric birefringence* (or the related *electric dichroism*) decay.

#### 3.1. Fluorescence Anisotropy Decay

Tryptophan residues provide intrinsic fluorescent chromophores in many proteins (ie, the ability to reradiate electromagnetic energy at a wavelength longer than that absorbed). If the incident monochromatic light is plane-polarized, the reradiated radiation will not only be at a longer wavelength but will also be wholly or partially depolarized depending (among other factors) on the extent of rotational Brownian motion of the macromolecule. If the macromolecule is not a protein and does not have an intrinsic fluorescent chromophore, or is a protein-containing insufficient tryptophan residues, a chromophore can be attached synthetically.

The emitted polarization is measured by two detectors, one normal to, and the other perpendicular to,

the plane of polarization; the extent or “anisotropy” of polarization,  $r(t)$ , is recorded as a function of time,  $t$ :

$$r(t) = \frac{I_{\parallel} - I_{\perp}}{I_{\parallel} + 2I_{\perp}} \quad (12)$$

where  $I_{\parallel}$  is the emitted intensity in a plane parallel to the polarization of the incident light and  $I_{\perp}$  is in a plane perpendicular. In time-resolved fluorescence depolarization anisotropy, this rapid decay in anisotropy of polarization in response to a pulse of incident radiation is recorded and averaged over many pulses (in some ways, the situation is analogous to the measurement of translational diffusion using [dynamic light scattering](#)). Allowance has to be made for the decay in intensity of the chromophore itself, specifically, the decay of the intrinsic fluorescence intensity has to be deconvoluted from the anisotropy decay function. The decay in  $r(t)$  with time can then be analyzed in terms of the rotational relaxation times of the molecule. There will be one relaxation time for a spherical particle, three for a particle with an axis of symmetry. For a general asymmetric molecule, there will be five relaxation times that need resolving:

$$r(t) = c_1 \exp(-t/\tau_1) + c_2 \exp(-t/\tau_2) + c_3 \exp(-t/\tau_3) + c_4 \exp(-t/\tau_4) + c_5 \exp(-t/\tau_5) \quad (13)$$

or, more simply  $r(t) = \sum_{i=1}^5 c_i \exp(-t/\tau_i)$  where  $i = 1-5$ . In practice, at least two pairs of relaxation times are similar; hence the problem is one of resolving three decay constants (this will be particularly true for macromolecules with an axis of symmetry). Once resolved, these can be related to macromolecular shape and hydration ([6](#), [7](#)) using relations similar to Equation ([11](#)). However, extraction of decay constants from a “multiexponential decay”—of which equation ([13](#)) is an example—is what the mathematicians call “an ill-conditioned problem” and is not easy, especially if the relaxation times are relatively similar. A further problem is that the chromophore most not relative to the rest of the macromolecule.

A much simpler procedure is to measure fluorescence depolarization or anisotropy decay in the steady state, where the light source is continuous rather than pulsed ([8](#)). It can be used to obtain the harmonic mean rotational relaxation time,  $t_h$ .  $t_h$  and  $1/t_{1-5}$  can all be related to the shape and hydration of the macromolecule ([7](#), [9](#)). A study of [fibrinogen](#) provides a good example of the application of both time-resolved and steady-state fluorescence measurements ([10](#), [11](#)).

### 3.2. Electric Birefringence Decay

Solutions of macromolecules oriented in an electric field will be *birefringent*, having different refractive indices for light polarized parallel to and perpendicular to the electric field. This is known as *electric birefringence*. A related phenomenon, for macromolecules with absorbing chromophore, is *electric dichroism*, where a solution of macromolecules oriented in an electric field exhibits different extinction coefficients parallel to and perpendicular to the electric field.

When the electric field is switched off, the birefringence (or difference in refractive indices)  $\Delta n$  will decay because of rotational motions of the macromolecule:

$$\Delta n(t) = \sum_{i=1}^5 c_i \exp(-t/\tau_i) \quad (14)$$

where  $i = 1-5$ . However, there will be just two relaxation times for molecules that can be approximated by homogeneous ellipsoidal shapes, and just one for a homogeneous ellipsoid with an axis of symmetry. An electrically homogenous spherical particle exhibits no birefringence.

Like fluorescence anisotropy decay measurements, the relaxation times  $t_1$  can be related to molecular shape and hydration (Fig. 2) (12), but there are also practical problems to be overcome. The main problem has been that of local overheating in the solution caused by the large orienting electric fields. This has meant in the past that experiments have been limited to solutions of low ionic strength. With significant advances in charge shielding in modern instrumentation, however, physiological ionic strengths are now a reality (13).

### 3.3. Diffusion of Small Molecules through Biomolecular Systems

Although most attention is given to diffusion phenomena in macromolecules, the importance of the diffusion through biopolymer matrices of small molecules and ions, even water molecules themselves, cannot be ignored; indeed, many physiological processes involve passive or active transport of water and other low-molecular-weight species through cellular and other matrices. As a direct comparison with macromolecular diffusion, Table 1 also gives the self-diffusion coefficient at 20.0°C of water ( $2.3 \times 10^{-5}$  cm<sup>2</sup>/s) and a small sugar molecule, sucrose ( $\sim 4.6 \times 10^{-6}$  cm<sup>2</sup>/s); the diffusion rates are orders of magnitude greater than those for more bulky macromolecules.

Arguably the best method for measuring diffusion of water and other small molecules is *pulsed field gradient spin echo NMR* (14). Spinning charged particles, such as atomic nuclei, will have an associated magnetic dipole moment, which generates a magnetic energy when placed within a magnetic field. Quantum-mechanical considerations restrict the energies to a limited number of discrete values. Transitions between levels—whose spacing depends on the external magnetic field—correspond to radiofrequency (RF) radiation, so the nucleus will interact or “resonate” with certain radiofrequencies if the external magnetic field is varied; this is known as nuclear magnetic resonance. These frequencies, and the strength and breadth of the resonances, depend on the particular atomic nuclei that are being examined and on the environment in which they find themselves. It is therefore possible to “home in” on a particular nuclear species; for example, such nuclei could be the hydrogen atoms in a water molecule.

With the pulsed field gradient technique, an excitation pulse of RF radiation is applied across a sample—in this case a biomolecular matrix—causing alignment of spins and the generation of an NMR signal; this signal subsequently decays as a result of diffusion. If a second RF pulse is applied after an interval  $t$ , the decay processes are reversed and a refocused “echo” signal is obtained at a time  $2t$  after application of the first pulse. This recovery process can be interrupted as follows. A pair of matched magnetic field gradient pulses that vary linearly across the sample are also applied; the first is applied between the RF pulses and the second, after the second RF pulse. The echo only has no net effect if there has been no diffusive movement of the molecules; any movement will result in an attenuation of the echo. Analysis of this attenuation can be used to determine the diffusion coefficient.

## 4. Acknowledgment

The author is grateful to Professor W. Derbyshire and Dr. I Farhat for their helpful comments.

## Bibliography

1. S. E. Harding and P. Johnson (1985) *Biochem. J.* **47**, 247–250.
2. O. Lamm (1928) *Z. Phys. Chem.* **A138**, 313.
3. L. J. Gosting (1956) *Adv. Prot. Chem.* **11**, 429–554.
4. O. Lamm and A. Polson (1936) *Biochem. J.* **30**, 528–541.
5. P. Claes, M. Dunford, A. Kenney, and P. Vardy (1992) in *Laser Light Scattering in Biochemistry* (S. E. Harding, D. B. Sattelle, and V. A. Bloomfield, eds.), Royal Society of Chemistry, Cambridge, U.K., pp. 66–76.
6. S. E. Harding, J. C. Horton, and H. Cölfen (1996) *Eur. Biophys. J.* **25**, 347–359.
7. J. Garcia de la Torre, B. Carrasco, and S. E. Harding (1997) *Eur. Biophys. J.* **25**, 361–372.

8. G. Weber (1952) *Biochem. J.* **51**, 145–155.
9. J. Garcia de la Torre and V. A. Bloomfield (1977) *Biopolymers* **16**, 1747–1763.
10. A. U. Acuna, J. Gonzalez-Rodriguez, M. P. Lillo, and K. R. Naqvi (1987) *Biophys. Chem.* **26**, 55–61.
11. A. U. Acuna, J. Gonzalez-Rodriguez, M. P. Lillo, and K. R. Naqvi (1987) *Biophys. Chem.* **26**, 63–70
12. J. Antonsiewicz and D. Porschke (1988) *J. Biomol. Struct. Dyn.* **5**, 819–837.
13. D. Porschke and A. Obst (1991) *Rev. Sci. Instrum.* **62**, 818–820.
14. E. O. Stekskal and J. E. Tanner (1965) *J. Chem. Phys.* **42**, 288–292.

### Suggestions for Further Reading

15. E. Fredericq and C. Houssier (1973) *Electric Dichroism and Electric Birefringence*, Oxford Univ. Press, Oxford (classic text, although somewhat dated).
16. K. E. Van Holde (1985) *Physical Biochemistry*, Prentice-Hall, Englewood Cliffs, N. J. [Chapters 4 (general diffusion), 6 (electro-optics) and 8 (fluorescence depolarisation) give an excellent introduction.]

## Diffusion-Controlled Reactions

Any catalytic or chemical process can be dissected, at least theoretically, into both physical and chemical steps. For a chemical reaction to occur, the reactants must come into physical contact in the correct proximity and orientation. Then whether they simply dissociate or undergo a reaction depends on a number of factors, especially the probability that the reaction will occur. With very facile reactions, every encounter results in a reaction, and the rate-limiting step for the overall reaction is the physical encounter of the two reactants, which is governed primarily by diffusion.

The second-order rate constant for a diffusional encounter of two molecules is  $10^{10}\text{s}^{-1}\text{M}^{-1}$  for macromolecules and  $10^8$  to  $10^9\text{s}^{-1}\text{M}^{-1}$  for small molecules. These numbers can, however, be altered by factors of up to  $10^2$  by attractive or repulsive interactions, especially electrostatic, between the reactants.

### Suggestion for Further Reading

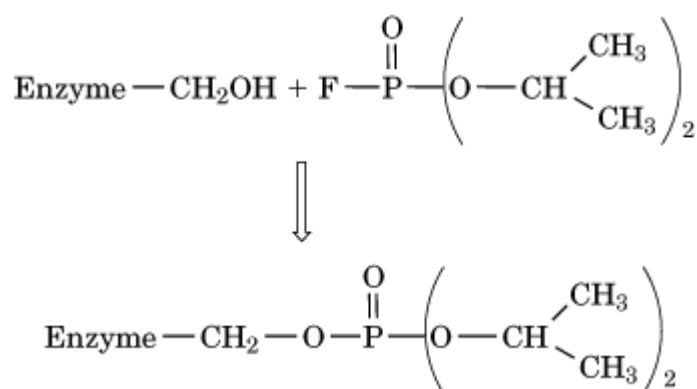
G. G. Hammes and P. S. Schimmel (1970) In *The Enzymes*, 3rd ed. (P. D. Boyer, ed.), Academic Press, New York, Vol. **2**, pp. 67–114.

## DIFP (Diisopropylfluorophosphate)

Diisopropylfluorophosphate, also known as diisopropylphosphofluoridate, DIPF, and sometimes abbreviated DFP, was developed as a nerve gas during the World War II. It is an irreversible

inhibitor of **acetylcholinesterase** and, more generally, an inhibitor of [serine proteinases](#) (1). In **chymotrypsin**, a typical serine proteinase, the [active site](#) Ser195 is phosphorylated by DIFP in a reaction that mimics the catalytic mechanism (Fig. 1). None of the other 27 serine residues of chymotrypsin is affected. The resultant diisopropylphosphoryl serine is stable to subsequent hydrolysis, and the phosphorylated enzyme is totally inactive. The reagent is a liquid at room temperature and is usually diluted into isopropanol before use (2). It is volatile and highly toxic: LD<sub>50</sub> in mice is 3.7 mg/kg when administered orally. A more convenient reagent for inactivating serine proteinases is phenylmethane sulfonyl fluoride ([PMSF](#)) (3).

**Figure 1.** Reaction of the hydroxyl group of the active site serine residue of a serine proteinase with diisopropyl fluorophosphate. The product is catalytically inactive.



## Bibliography

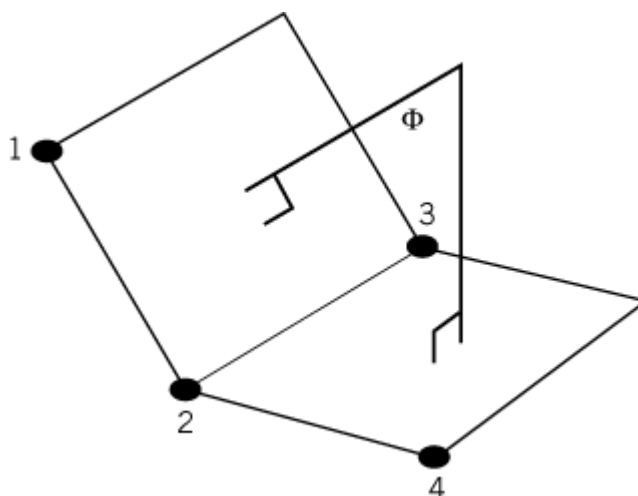
1. J. A. Cohen, R. A. Oosterbaan, and F. Berends (1967) *Methods Enzymol.* **11**, 686–702.
2. R. J. Beynon and J. S. Bond (1993) In *Proteolytic Enzymes: A Practical Approach* (R. J. Beynon and J. S. Bond, eds.), IRL Press, Oxford, U.K., p. 243.
3. B. M. Dunn (1993) In *Proteolytic Enzymes: A Practical Approach* (R. J. Beynon and J. S. Bond, eds.), IRL Press, Oxford, U.K., pp. 57–81.

## Dihedral Angle

A dihedral angle is the angle between two planes (1). This is determined as the angle formed by normal lines to each plane, as shown in Figure 1. When the rotation about a single bond in a molecule is being measured, **eg**, the rotation about the [peptide](#) N<sub>Cα</sub> bond, it is preferred nomenclature to use the conventions described for [torsion angles](#). Unlike torsion angles, it is possible to determine the dihedral angle between the planes defined by any two groups of three atoms, whether they are bonded together or not. If any four atoms are specified in order: atom-1, atom-2, atom-3, atom-4, the dihedral angle is measured between the plane containing atom-1, atom-2, and atom-3 and the plane containing atom-2, atom-3, and atom-4. If the four atoms are not sequentially bonded, this dihedral angle may be referred to as an “improper dihedral” or “improper torsion” (2).



**Figure 1.** The dihedral angle  $\Phi$  is the angle between the two planes defined by atoms 1, 2, and 3 and atoms 2, 3, and 4, respectively. The angle is determined as the angle between lines perpendicular to each plane.



#### Bibliography

1. *Oxford English Dictionary*, 2nd ed. (1989) Clarendon Press, Oxford.
2. B. R. Brooks et al. (1983) *J. Comp. Chem.* **4**, 187–217.

#### Suggestion for Further Reading

3. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer-Verlag, New York, p. 15.

### Dihydrofolate Reductase (DHFR)

Dihydrofolate reductase (DHFR) catalyzes the **NADPH**-dependent reduction of 7,8-dihydrofolate to 5,6,7,8-[tetrahydrofolate](#) (E.C. 1.5.1.3). Tetrahydrofolate, in its variously modified forms, serves as a **cofactor** in the transfer of one-carbon fragments in the biosynthesis of thymidylate and **purine** nucleotides (see **Nucleosides, nucleotides**), as well as other metabolic functions. DHFR is universally present in every kind of living cell—**prokaryotes, eukaryotes and archaeobacteria**—where it is essential for the biosynthesis of **DNA**. Thus if DHFR activity is inhibited, cellular replication cannot take place, a phenomenon that forms the basis for the action of a number of anticancer and antimicrobial pharmaceuticals (eg, **methotrexate** and **trimethoprim**) collectively known as *antifolates*. These compounds are simply inhibitors of DHFR with varying degrees of species selectivity. Though required for continuous [DNA replication](#), the role of DHFR is a rather indirect one—that is, maintenance of the intracellular pool of tetrahydrofolate in its fully reduced state. This is necessary because during the final step of thymidylate synthesis, wherein dUMP is methylated at the 5-position to yield dTMP in a reaction catalyzed by [thymidylate synthase](#), the reduced pteridine ring of the one-carbon donating cofactor  $N^5, N^{10}$ -methylenetetrahydrofolate is concomitantly oxidized to produce 7,8-dihydrofolate.

This article focuses solely on the molecular properties of DHFR, which, because of its biomedical importance, its conveniently small size, and its ready availability through [cloning](#) and [protein engineering](#), has been intensively studied by essentially every enzymological and physicochemical technique available. Omitted, however, is any discussion of the several studies of **protein folding** that have been carried out with DHFR ([1](#), [2](#)). Nor does this article discuss what is known about how the intracellular level of DHFR is regulated during [transcription](#) ([3](#)) and [translation](#) ([4](#)) and by [gene amplification](#) ([5](#)). Also passed over without further comment is the fact that the DHFR gene has been widely employed as a selectable genetic marker (see [Reporter Genes](#)).

## 1. The Molecule

The polypeptide chain composing DHFR from most species is relatively short, in the range of 159 amino acids for *Escherichia coli* to 229 for *Cryptococcus neoformans*. Exceptions, however, occur in both protozoa and at least some **plants**, where the genes for DHFR and thymidylate synthase are fused and thus encode a bifunctional enzyme molecule. The degree of sequence identity among all DHFR across the biological spectrum is about 25%, with the pattern of residue conservation leaving little doubt that all are **homologous** and folded into quite similar three-dimensional structures (see [Protein Structure](#)). That structure is well established; it has been determined by [X-ray crystallography](#) for the proteins from *E. coli*, *Lactobacillus casei*, chicken, mouse, human, the protozoon *Leishmania major*, and the fungus *Pneumocystis carinii*. In all cases, these structures contain one or more bound ligands, such as the cofactor, folate, or a variety of inhibitors, so the mode of substrate binding is known in detail. A hypothetical “movie” of the *E. coli* DHFR molecule as it proceeds through its catalytic cycle has been assembled on the basis of 47 distinct crystal structures in 12 different crystal packings ([6](#)).

The molecular architecture of DHFR is typical of a basic folding plan that is seen in the case of many other enzymes, notably **NAD(P)** or **flavin** nucleotide-dependent [oxidoreductases](#), as well as numerous enzymes that bind AMP or ATP. This kind of fold is often called an open a/b-sheet structure, because it is composed of a central twisted [b-sheet](#) flanked on both sides by a number of **a-helices**. However, this structure is quite variable in the numbers of a-helices and b-strands, their direction relative to one another, the order in which they are connected (ie, their topology), and the size and geometry of the connecting loops.

The central b-sheet of all DHFR contains eight b-strands, with the first seven strands, A–G, all parallel; the lone antiparallel strand, H, comprises the C-terminal portion of the polypeptide chain. Four a-helices flank the central b-sheet, two on each side. In vertebrate DHFR, which are slightly larger than bacterial DHFR, an extra seven-residue loop is inserted in the middle of the outermost strand G of the central b-sheet; vertebrate DHFR contain, in addition, an extra two-turn a-helix and two turns of a left-handed **polyproline**-like helix.

Both the substrate and cofactor are bound in more or less extended conformations, at approximately right angles to one another, with the pteridine ring of the substrate and the nicotinamide ring of the cofactor occupying a pocket between the carboxyl ends of strands bA and bE. They are in close contact in such a way that hydride transfer can take place between the A side (*pro-R* or *re*) of the nicotinamide ring at C4 and C6 or C7 of the pteridine ring.

## 2. Enzyme Kinetics

The kinetic pathway of DHFR is complicated by the participation of nine possible ligation states, each distinguished by whether the cofactor binding site and/or the substrate binding site is (i) unoccupied, (ii) occupied by the oxidized ligand, or (iii) occupied by the reduced ligand. Rate constants for the interconversion between these ligation states have been measured and assembled into complete kinetic schemes for DHFR from *E. coli* ([7](#)), *L. casei* ([8](#)), human ([9](#)), mouse ([10](#)), and *Pneumocystis carinii* ([11](#)). These schemes accurately predict **steady-state** kinetic behavior in computer simulations. Although the predominant pathways vary somewhat, in neither bacterial nor

vertebrate DHFR is the chemical step (ie, proton donation to N5 and hydride transfer to C6 of dihydrofolate) **rate-limiting**. Instead, the rate-limiting steps are either substrate binding or product release. These DHFR also show evidence of a conformational isomerization prior to the chemical step, possibly involved in proton donation to N5 of dihydrofolate (12). Among the few kinetic differences between bacterial and vertebrate DHFR, most prominent is the much greater efficiency of vertebrate DHFR at catalyzing the reduction of folate.

### 3. Mutagenesis Studies

A plethora of mutants of DHFR (primarily from *E. coli* and human) have been engineered by [site-directed mutagenesis](#) to probe the catalytic role of evolutionarily conserved side chains. Only a handful of these mutants have been characterized both structurally and kinetically, but interpretation of the results is not straightforward for even the most thoroughly studied. For example, mutation of conserved **Asp27** → **Asn** of *E. coli* DHFR, or the equivalent **Glu30** → **Gln** in human DHFR, results in a strong pH-dependence of  $k_{cat}$ , suggesting a role for the carboxylic acid in protonation of N5 of the substrate. However, the equivalent mutation, **Asp26** → **Asn**, in *L. casei* DHFR yields only a modest pH effect (13). Structural perturbations are minimal in these and most other such site-directed mutants, but the rate of the chemical step decreases significantly, suggesting that the entire enzyme molecule influences hydride transfer and proton donation. Further mutagenesis studies involve modification of whole  $\alpha$ -helices (14) and loops (15) or circular permutations of the polypeptide chain (16) in an effort to understand the functional role of larger structural elements in catalysis.

### 4. The Transition State

At the center of the effort to understand any enzyme lies a fundamental question: How does it catalyze the reaction? In the most general terms, elementary [catalysis](#) theory dictates that the enzyme molecule must bind the [transition state](#) of the reaction many orders of magnitude more strongly than the reactants. But to progress beyond this level of understanding for any particular enzymic reaction requires that one begin with a picture of the transition state; and this is no easy matter, because the transition state is by definition an unstable, transitory molecular species. Thus one must rely on theory, and even speculation, to suggest what the transition state might look like.

In the case of DHFR, there is at least one very clear indication from the crystal structures that the enzyme molecule is indeed “designed” to bind the transition state. Certainly the transition state must involve a close approach of the two carbon atoms between which hydride ion transfer takes place, C4 of the nicotinamide and C6 of the pteridine ring. Theoretical calculations (17) predict a C4–C6 distance of 2.6 Å in the transition state, about 1 Å shorter than **van der Waals** contact distance. Happily this is just what is seen when the separate structures of the binary complexes, DHFR · folate and DHFR · NADPH, are carefully superposed on one another in the case of both the vertebrate (18) and the *E. coli* (19) enzymes. Evidence that DHFR is trying to push the C4 and C6 atoms together is seen even in ternary complexes in which both the cofactor and the pteridine are bound simultaneously; the approach distance is 0.3 to 0.5 Å shorter than the expected van der Waals distance of 3.6 Å.

### Bibliography

1. B. E. Jones and C. R. Matthews (1995) *Protein Sci.* **4**, 167–177.
2. S. D. Hoeltzli and C. Frieden (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9318–9322.
3. J. E. Slansky and P. J. Farnham (1996) *Bioessays* **18**, 55–62.
4. E. Ercikan et al. (1993) In *Chemistry and Biology of Pteridines and Folates* (J. E. Ayling et al., eds.), Plenum Press, New York, pp. 537–540.
5. J. L. Hamlin (1992) *Mutat. Res.* **276**, 179–187.
6. M. R. Sawaya and J. Kraut (1997) *Biochemistry* **36**, 586–603.
7. C. A. Fierke, K. A. Johnson, and S. J. Benkovic (1987) *Biochemistry* **26**, 4085–4092.

8. J. Andrews et al. (1989) *Biochemistry* **28**, 5743–5750.
9. J. R. Appleman et al. (1990) *J. Biol. Chem.* **265**, 2740–2748.
10. J. Thillet, J. A. Adams, and S. J. Benkovic (1990) *Biochemistry* **29**, 5195–5202.
11. S. A. Margosiak, J. R. Appleman, D. V. Santi, and R. L. Blakley (1993) *Arch. Biochem. Biophys.* **305**, 499–508.
12. W. A. Beard et al. (1989) *J. Biol. Chem.* **264**, 9391–9399.
13. J. Basran, M. G. Casarotto, I. L. Barsukov, and G. C. K. Roberts (1995) *Biochemistry* **34**, 2872–2882.
14. L. Li and S. J. Benkovic (1991) *Biochemistry* **30**, 1470–1478.
15. L. Li, C. J. Falzone, P. E. Wright, and S. J. Benkovic (1992) *Biochemistry* **30**, 7826–7833.
16. A. Buchwalder, H. Szadkowski, and K. Kirschner (1992) *Biochemistry* **31**, 1621–1630.
17. Y.-D. Wu and K. N. Houk (1987) *J. Am. Chem. Soc.* **109**, 2226–2227.
18. J. F. Davies et al. (1990) *Biochemistry* **29**, 9467–9479.
19. V. M. Reyes, M. R. Sawaya, K. A. Brown, and J. Kraut (1995) *Biochemistry* **34**, 2710–2723.

### Suggestion for Further Reading

20. R. L. Blakley (1995) Eukaryotic dihydrofolate reductase. *Adv. Enzymol.* **70**, 23–102. (A thoroughly detailed review by a foremost authority.)

## Dimethyl Sulfate (DMS)

This is a [mutagen](#) that is a monofunctional alkylating agent of structure  $\text{CH}_3\text{-O-SO}_2\text{-O-CH}_3$ , which reacts with nucleophilic sites primarily by an  $\text{S}_{\text{N}}2$  mechanism. This mechanism involves a concerted process with the substrate, resulting in selective methylation of only the most nucleophilic sites. Thus DMS predominantly methylates nitrogen atoms on **nucleic acids**, only rarely oxygen or phosphorus. By far the most abundant lesion produced is N7-methylguanine, with minor amounts of  $\text{N}_1$ -,  $\text{N}_3$ -, and  $\text{N}_7$ -methyladenine and  $\text{N}_3$ -methyl cytosine, and trace amounts of  $\text{O}_6$ -methylguanine (1). 7-Methylguanine has base-pairing properties very similar to guanine itself, so this modification is probably not the most important cause of mutagenesis by DMS. Both 3-methyladenine and 3-methylcytosine can cause base mispairing, which may lead directly to mutagenesis by DMS (1, 2). Although probably not primarily involved in mutagenesis, methylation at the  $\text{N}_7$  of guanine significantly weakens the  $\text{N}_9$ -glycosidic bond, leading to enhanced formation of apurinic sites (3-5). [Depurination](#) may also occur at 3-methylated adenines (6, 7), and this itself has mutagenic consequences.

DMS is both **carcinogenic** and **mutagenic** in a wide variety of organisms, producing point mutations, chromosomal aberrations, and recombinational events (reviewed by Hoffman (8)). In many mutagenesis studies, DMS is reported to be a [base-pair substitution](#) mutagen, but it also causes [frameshift mutations](#) and deletions. It has been suggested that much of the data can be explained in terms of effects of DMS on DNA **repair** processes (8, 9). Lawley and Warren (2) found that *Escherichia coli* repair processes effectively remove 3-methylguanine and 3-methyladenine, but not 7-methylguanine, residues from DNA. These authors suggested that N-3 purine alkylations may block [DNA replication](#) and stimulate enzymatic repair processes. Inaccuracies in these may lead to many of the observed mutations in bacterial systems. 7-Methylguanine is removed from the DNA of

mammalian cells, although it does not interfere with DNA synthesis (10).

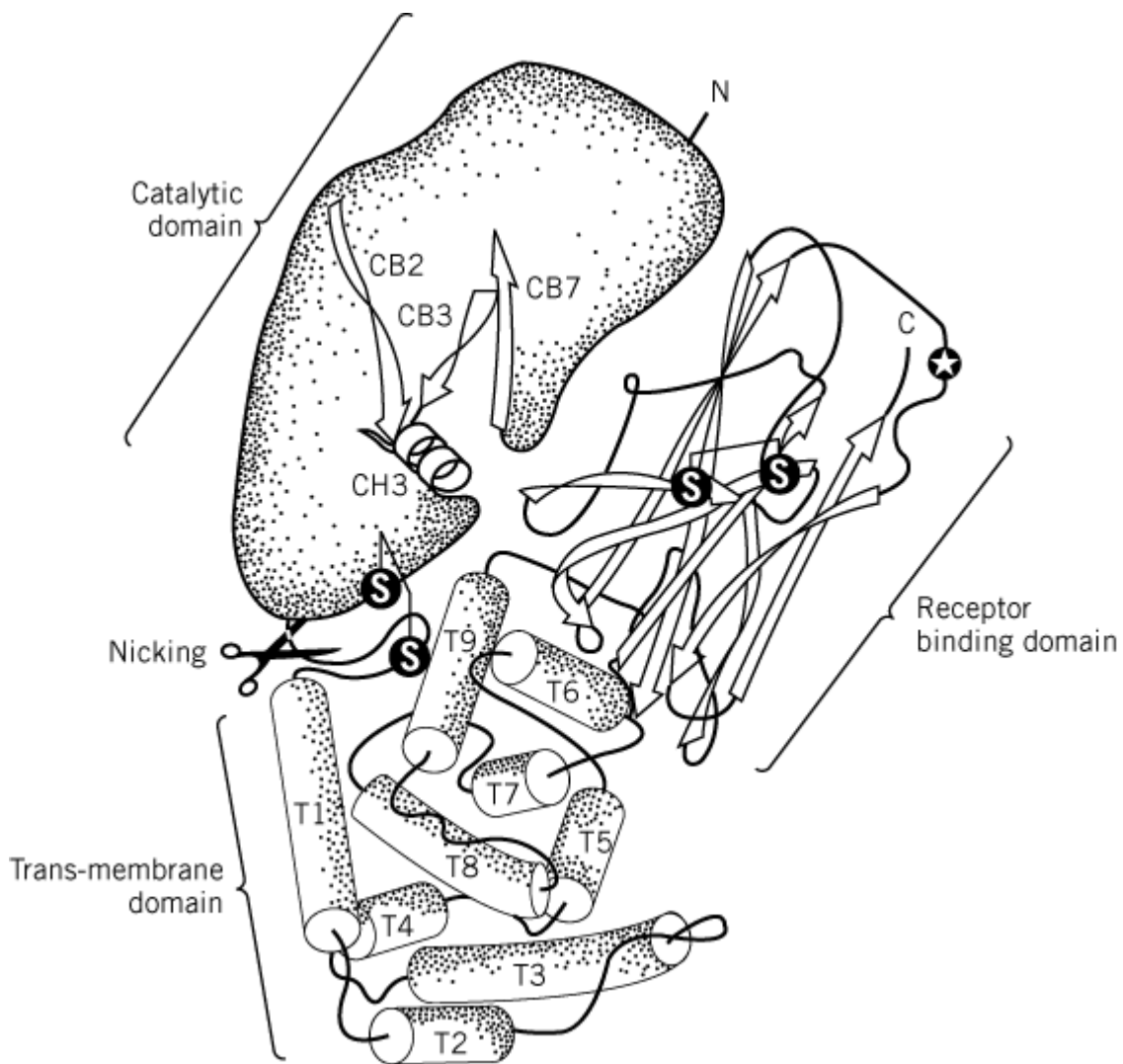
## Bibliography

1. P. D. Lawley, D. J. Orr, and S. A. Shah (1972) *Chem. Biol. Interact.* **5**, 286–288.
2. P. D. Lawley and W. Warren (1976) *Chem. Biol. Interact.* **12**, 211–220.
3. P. D. Lawley (1966) *Prog. Nucleic Acid. Res. Mol. Biol.* **5**, 89–131.
4. P. D. Lawley (1974a) In *Molecular and Environmental Aspects of Mutagenesis* (L. Prakash, ed.), Thomas, Springfield, IL, pp. 17–33.
5. P. D. Lawley (1974b) *Mutat. Res.* **23**, 283–295.
6. P. D. Lawley and P. Brookes (1963) *Biochem. J.* **89**, 127–138.
7. R. H. C. San and H. F. Stich (1975) *Int. J. Cancer* **16**, 284–291.
8. G. R. Hoffman (1980) *Mutat. Res.* **75**, 63–129.
9. J. W. Drake and R. H. Baltz (1976) *Annu. Rev. Biochem.* **45**, 11–37.
10. D. A. Scicchitano and P. C. Hanawalt (1990) *Mutat. Res.* **233**, 31–38.

## Diphtheria Toxin

Diphtheria was the first infectious disease shown to be caused solely by a [toxin \(1\)](#). Diphtheria toxin (DT) is encoded by a gene carried by temperate **bacteriophages**, such as the *o*-corynephage, that lysogenize *Corynebacterium diphtheriae* and make it toxigenic. The toxin gene is activated by the release of an iron-binding [repressor](#) protein, which binds to DNA in its holo form and in iron-depleted media is released as the apo form (1). The DT gene includes a 25-residue long amino-terminal **leader sequence** responsible for the export of DT from the cell (2, 3). After proteolytic removal of the leader sequence, DT is released by *C. diphtheriae* as a single polypeptide chain of 535 amino acid residues. This chain is rapidly cleaved *in vivo* by furin (4) and *in vitro* by different proteinases, within a loop containing three Arg residues that connects the first and second domains (Fig. 1). After such proteolytic nicking, DT is composed of two-chains: A (193 residues), responsible for the intracellular catalytic activity, and B (342 residues), responsible for binding and membrane translocation. A and B are linked by noncovalent forces and a single interchain [disulfide bond](#). Figure 1 shows the three-domain structural organization of DT (5). The amino terminal domain (residues 1 to 193) corresponds to the A subunit, which is an **ADP-ribosyltransferase** with a cleft-shaped **NAD-binding** cavity harboring the Glu148, His21, Trp50, and Tyr65 residues that are essential for activity (5). NAD binds in such a way that its nicotinamide moiety overlaps the phenolic ring of Tyr65, and its adenine group lies over the indole ring of Trp50.

**Figure 1. Structure of diphtheria toxin.** Diphtheria toxin consists of three domains, characterized by different secondary structure elements and different biological properties. The cleft in the catalytic domain is the NAD binding site. The circled star indicates Ser-508 involved in receptor binding. The toxin is activated by selective proteolysis of a furin-sensitive loop between the A domain and T domain, which is responsible for membrane translocation.



In the nicked toxin, the A domain is disulfide-linked to the T domain (residues 194 to 385), which is involved in the membrane insertion and translocation of DT. The T domain is almost entirely composed of **alpha-helices**, a type of **secondary structure** found in the membrane-embedded sector of many integral [membrane proteins](#). The third domain R (residues 386 to 535) is a flattened  $\beta$ -barrel and responsible for DT binding to its specific **receptor**, which is a protein that binds heparinlike [growth factors](#) (6). The role of this receptor in cell physiology is not known, but it was shown that CD9 enhances the sensitivity of cells to DT (7). The number of DT receptors in different cell lines correlates with their DT sensitivity. Binding is followed by internalization inside [clathrin](#)-coated **vesicles**, which uncoat and merge into early [endosomes](#), a process that requires less than 5 min in Vero cells at 37°C (8). The acidic pH of the endosomal lumen induces a conformational change of DT, with insertion of helices of the T domain into the lipid bilayer (9, 10), to create a transmembrane channel, open laterally to lipids, that mediates the translocation of the A chain into the cytosol (11). It is not known if domain A has to unfold completely in order to pass through the B channel, but it is noteworthy that A can be boiled in **SDS** and, upon detergent removal, refolds into its active conformation. The translocation and release of A are linked to the reduction of the interchain disulfide bond, which is carried out by yet unidentified reductases. Such reduction is the rate-limiting step of the entire intoxication process (8). Only about a third of the internalized DT molecules are reduced and release their A chain into the cytosol. The remaining, unreduced two-thirds, as well as the B subunits left over on the membrane after A release, are conveyed to the **lysosomes** and degraded (8).

The A chain of DT **ADP-ribosylates**, specifically [elongation factor 2](#) (EF-2), which becomes unable to transfer the tRNA-bound amino acid to the growing polypeptide chain, and cell **protein biosynthesis** are blocked ([12](#)). DT-A is active only on the EF-2 of **eukaryotes** and **archaeobacteria** because they contain diphthamide, a **post-translationally modified** His residue, which is the unique substrate of DT and of the related exotoxin A released by *Pseudomonas aeruginosa* ([12](#)). The physiological role of diphthamide *in vivo* is unknown, particularly because cell lines mutated at this His residue, or in one of the enzymes involved in its transformation into diphthamide, appear to grow normally ([12](#)). One single molecule of DT is able to kill a cell rather rapidly ([13](#)) because it can catalyze the ADP-ribosylation of about 2000 EF-2 per hour, whereas cells in culture produce only about 1500 new EF-2 per hour, and EF-2 synthesis decreases as toxin action progresses. As a result, the lethal dose of DT in sensitive animal species (those possessing DT receptors) is below 0.1 µg DT/kg.

## Bibliography

1. A. M. Pappenheimer (1982) Harvey Lect. **76**, 45–73.
2. L. Greenfield et al. (1983) Proc. Natl. Acad. Sci. USA **80**, 6853–6857.
3. G. Ratti, R. Rappuoli, and G. Giannini (1983) Nucleic Acids Res. **11**, 6589–6595.
4. J. R. Murphy (1997) In *Guidebook to Protein Toxins and Their Use in Cell Biology* (R. Rappuoli and C. Montecucco, eds.), Sambrook and Tooze, Oxford University Press, Oxford, UK.
5. S. Choe et al. (1992) Nature **367**, 216–222.
6. J. G. Naglich, J. E. Metherall, D. W. Russell, and L. Eidels (1992) Cell **69**, 1051–1061.
7. T. Mitamura et al. (1992) J. Cell Biol. **118**, 1389–1399.
8. E. Papini, R. Rappuoli, M. Murgia, and C. Montecucco (1993) J. Biol. Chem. **268**, 1567–1574.
9. K. Sandvig and S. Olsnes (1981) J. Biol. Chem. **256**, 9068–9076.
10. K. J. Oh et al. (1996) Science **273**, 810–812.
11. G. Menestrina, G. Schiavo, and C. Montecucco (1994) Mol. Aspects Med. **15**, 81–193.
12. R. J. Collier (1990) In *ADP-ribosylating Toxins and G Proteins* (J. Moss and M. Vaughan, eds.), American Society for Microbiology, Washington, DC, pp. 3–19.
13. M. Yamaizumi, E. Mekada, T. Uchida, and Y. Okada (1978) Cell **15**, 245–250.

## Direct Methods

Before an electron density map can be calculated in a crystal structure determination using [X-ray crystallography](#), the phase angles of the reflections must be found (see [Phase Problem](#)). In principle, they can be derived from the values of the amplitudes  $|F(hk\ell)|$ , and these are (apart from correction factors) proportional to the square root of the observed reflection intensities  $I(hk\ell)$ . The relationship between phases and structure factor amplitudes is due to two properties of the electron density map:

1. It is nowhere negative.
2. It is a summation of atomic electron densities.

Several types of direct methods have been derived to determine the phase angles from the structure factor amplitudes alone, and they are routinely used in structurally determining small compounds. With an increasing number of atoms, however, the power of the method becomes weaker. The maximum number of non-hydrogen atoms in the [asymmetric unit](#) is around 300. Classical direct methods are not practical for proteins unless the protein molecule is small and gives an extremely

high-resolution X-ray diffraction pattern. It also helps if the protein contains one or more heavy atoms, for example, sulfur or iron (1). Although generally not practical for a complete protein structure determination, classical direct methods are suitable for locating heavy atoms from [isomorphous replacement](#) data.

A nonclassical approach to direct phase angle determination for macromolecules has been pioneered by Bricogne (2) and is still under development. It applies the maximum entropy principle. Although not yet applicable for a *de novo* structure determination, it can already be applied if combined with available information.

#### Bibliography

1. G. M. Sheldrick et al. (1993) *Acta Crystallogr.* **D49**, 18–23.
2. G. Bricogne (1993) *Acta Crystallogr.* **D49**, 37–60.

#### Suggestion for Further Reading

3. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Direct Repeats

A repeat, or repeating, **nucleotide sequence** is any of two or more identical segments of a longer sequence or identical segments of a double-stranded **DNA**. Any DNA fragment can be inserted into a long, double-stranded molecule in either of different polarities: “head on” (taking one of the ends of the duplex as a head), or “tail on.” In the long sequence, it then appears as either one strand of the fragment or the other, complementary strand. If two identical fragments are both oriented “head on,” the repeat is called a direct repeat—that is, all identical sequences encountered in a longer sequence are direct repeats. They are normally noncontiguous and separate. Contiguous direct repeats are called [tandem repeats](#). Various types of repeats are outlined schematically in the Figure for [Tandem Repeats](#).

One important example are the direct repeats at the ends of some [transposable elements](#) (1). Integrated retroviral [genomes](#) and [retrotransposons](#) are flanked by the [long terminal repeats](#) (LTR) of 200–500 bp. For example, the yeast transposon Ty has 330-bp LTRs at its ends. The repeats are involved in the [recombination](#) events during integration or excision of the [mobile elements](#). The presence of the direct repeats in a sequence may actually serve as an indication of possible transposition events. When integrated in the same polarity, the transposons themselves are an example of direct repeats.

#### Bibliography

1. D. E. Berg and M. M. Howe, eds. (1989) *Mobile DNA* ASM, Washington, D.C.

#### Suggestion for Further Reading

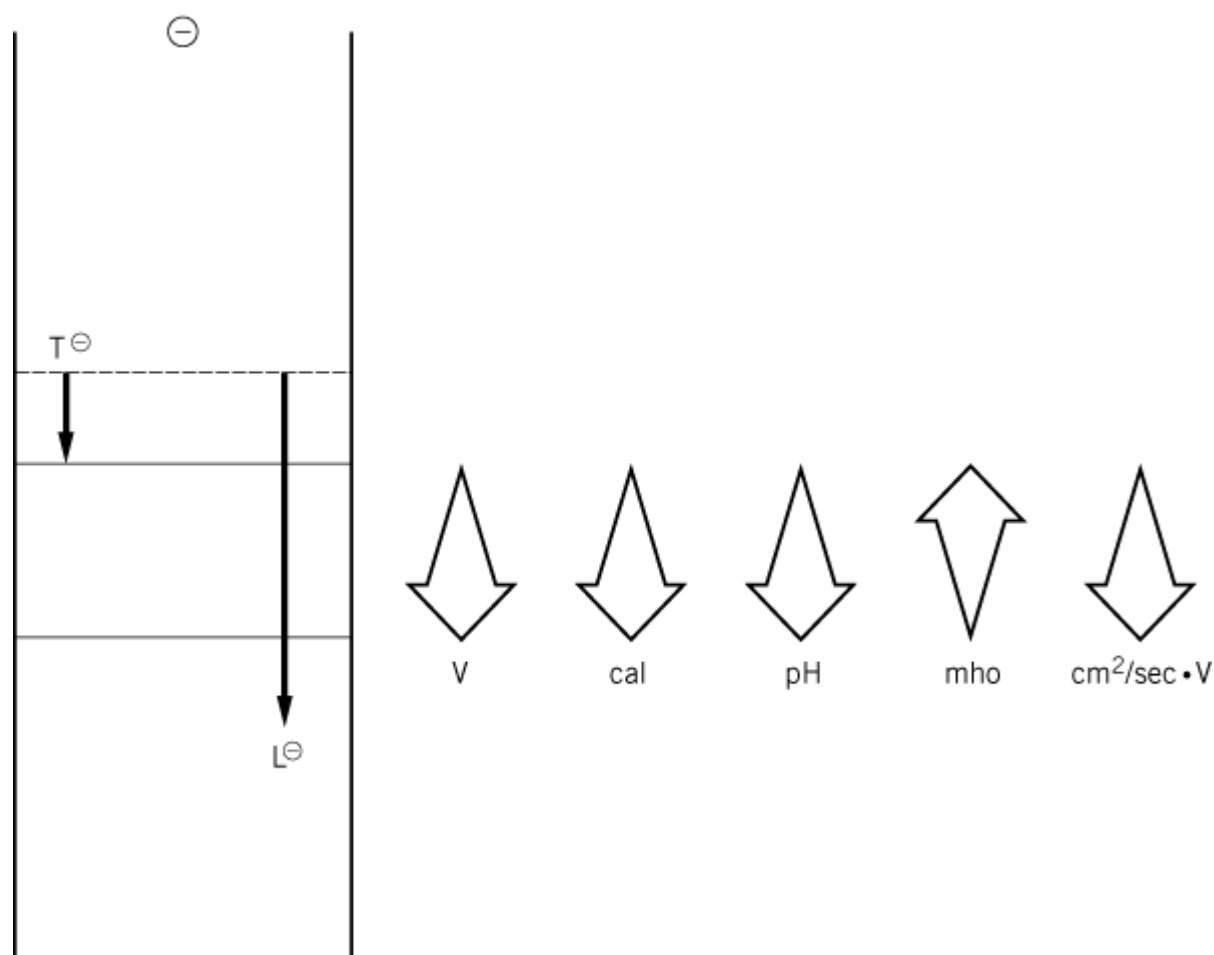
2. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.



## Disc Electrophoresis

[Gel electrophoresis](#) in *discontinuous buffer* systems is known as “disc electrophoresis” (1). Its main advantage is the provision of a highly concentrated starting zone even when the sample is dilute. It employs at least two electrophoresis buffers, i.e., a buffer discontinuity. They differ in either a cation or anion, whereas the counterion is the same. One of the differing ions, known as the “trailing ion,” has an electrophoretic mobility less than the species of interest to be separated; the other, the “leading ion,” has a greater mobility. Embedded within such a system, the species of interest will, at steady-state, migrate within a moving boundary between the leading and trailing ions, at a field strength (conductance) and concentration that is fixed (“regulated”) by the mobilities and concentrations of the leading, trailing, and common buffer ions constituting the boundary. This “regulation” imposes a “regulated” mobility on all trailing species with net mobilities equal to, or greater than, the trailing ion; hence, the system is designated as “equal-mobility-electrophoretic” or isotachophoretic (see [Isotachophoresis](#)). A number of species migrating isotachophoretically in order of their net mobilities is called a “stack.” The regulation of mobilities and concentrations behind the leading ion, resulting from the requirement for electroneutrality and the conservation of mass, gives rise to an increased field strength behind the fastest migrating (leading) ion in the electric field (Fig. 1). The gradient in field strength developing behind the leading ion causes equally oriented Joule heat and pH gradients, a gradient in mobility accelerating the trailing ions until at equilibrium they have attained the net mobility of the leading ion, and a gradient of specific conductance oriented inversely to the field strength gradient (Fig. 1). The key benefit of regulation of the trailing buffer phase by the leading ion is that a macromolecular ion present in that phase will migrate within the stack at a very high concentration, estimated to be on the order of 100 mg/mL for a 100-kDa protein, independent of the degree of dilution of the macromolecule in the original sample. Therefore, very dilute samples can produce very condensed zones of sample within the gel.

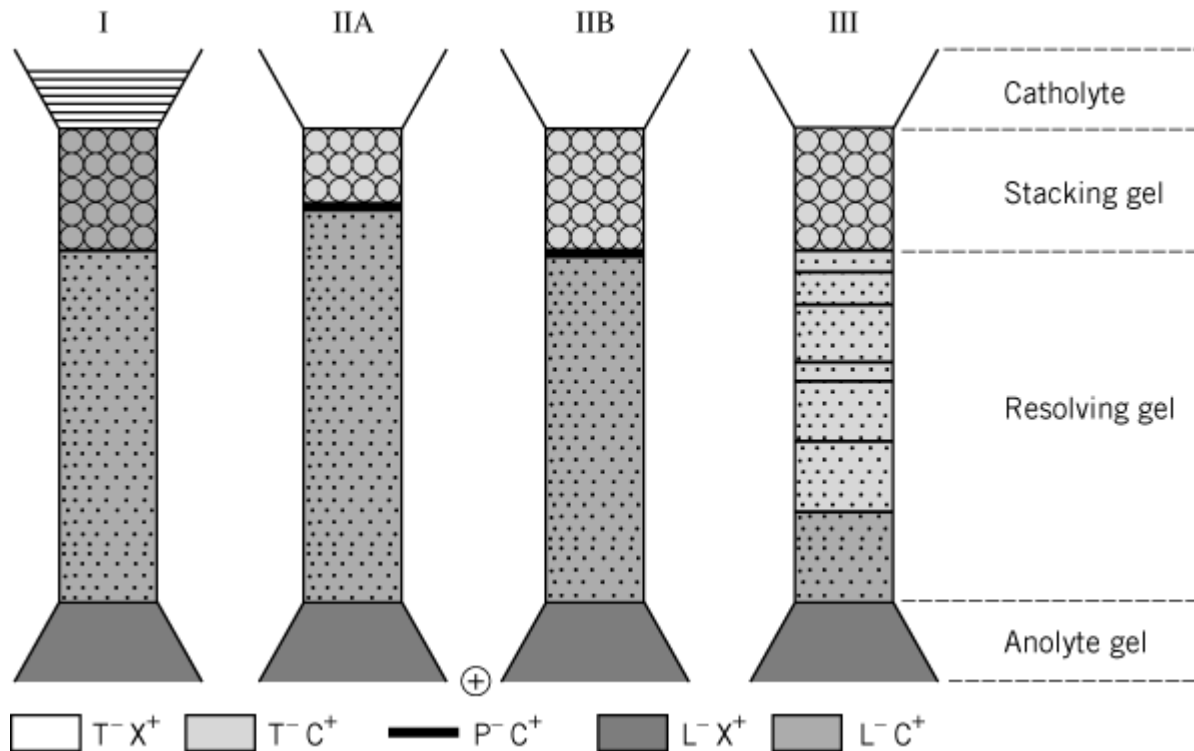
**Figure 1.** Operation of discontinuous buffer systems. Prior to electrophoresis, a stationary boundary (dotted line) separates the trailing ion T in the cathode buffer from the leading ion L in the gel buffer. As T and L electrophorese out of the starting zone, the greater instantaneous velocity of L would tend to create behind it an “ion vacuum.” This, however, would be in violation of the laws of conservation of mass and of the maintenance of electroneutrality. Operation of these two laws gives rise to a “regulation” in the zone behind L, which leads to an increase in field strength, Joule heat, and pH, and a decrease in conductance. These, in turn, produce an increased mobility of all species migrating between T and L until they attain the net mobility of L.



In disc electrophoresis, the concentrating (stacking) buffer phase is embedded in a gel that does not sieve but is present only to prevent convection; this gel is contiguous with a sieving gel, the “resolving gel.” As the stack of moving boundaries containing the concentrated sample enters the resolving gel, the buffer boundaries are unaffected by the elevated gel concentration and changed pH, whereas macromolecules are altered in their migration by both. The retardation due to the increased gel concentration causes the macromolecules to migrate at a mobility less than that of the trailing ion. Therefore, they “unstack” and separate according to their mobilities; in contrast, the stack of buffer ions remains (Fig. 2), forming a moving boundary “front” that can be marked by an appropriate [tracking dye](#). The ratio of macromolecular mobilities relative to the stacked tracking dye (in terms of relative mobility,  $R_f$ ) characterizes an electrophoretic band and can be used, when measured at several different gel concentrations, to specify the size and net charge of the macromolecule (see [Ferguson Plot](#)).

**Figure 2.** Operation of gel electrophoresis in a discontinuous buffer system. (I) A gel electrophoresis system of four phases is set up: A cathode buffer containing the trailing ion T as a salt with any counterion X plus a tracking dye; a stacking gel of large pore size, “nonrestrictive” to the migration of macromolecular particles P and containing the leading ion L as a salt with counterion C (which is common to both stacking and resolving gels) at a pH lower than that of the resolving gel; a resolving gel at a gel concentration that lowers (“restricts”) the mobility of particle P and at a pH higher than that of the stacking gel; and an anode buffer containing the common counterion C as a salt of any other ion. P is loaded on top of the stacking gel. (II) Upon initiating electrophoresis, a stack of L, P, and T forms by regulation (Fig. 1) in the stacking gel. (III) Upon arrival at the resolving gel, L and T and the tracking dye proceed as a stack at a velocity undiminished by the presence of the gel and augmented by the increased pH. In contrast, the velocity of P is lowered in the gel through “molecular sieving,” so that the mobility of P now is less than that of T. That decrease in mobility is termed “unstacking” and results in the migration of all species of P in decreasing order of

their net mobility, ie, in their separation.



Most important, disc electrophoresis is an analytical or preparative separation method that is largely independent of the volume of the original sample; this precludes the necessity for a concentration step prior to electrophoresis. This fact is important for those applications—and they are in the majority—in which the sample is available only in dilute solution. Furthermore, the stacking gel can also assist separation by excluding from the concentrated band (the stack) species of low or high mobility relative to the displacement rate of the moving boundaries associated in the stack (“selective stacking”).

Discontinuous buffer systems are available from a computer program (2) for electrophoresis across the entire pH scale, at either 0 or 25°C, and for both negatively and positively charged analytes.

#### Bibliography

1. L. Ornstein (1964) *Ann. NY Acad. Sci.* **121**, 321–349.
2. A. Chrambach (1985) In *The Practice of Quantitative Gel Electrophoresis* (V. Neuhoff and A. Maelicke, eds.), VCH, Weinheim, pp. 9–18, 85–110.

#### Suggestions for Further Reading

3. L. M. Hjelmeland and Chrambach (1982) The impact of L. G. Longworth on the theory of electrophoresis. *Electrophoresis* **3**, 9–17.
4. T. M. Jovin (1973) Multiphasic zone electrophoresis. *Biochemistry* **12**, 871–890.

#### Discontinuous DNA Replication

Elongation of the DNA chain on the lagging strand during [DNA replication](#) is discontinuous: Short segments of DNA, called [Okazaki Fragments](#), are repeatedly synthesized in the reverse direction of movement of the [replication fork](#) (1). This occurs because the two chains of double-helical DNA are antiparallel, and **DNA polymerase** can extend a DNA chain only in the 5' → 3' direction. On the leading strand, which runs 5' → 3' in the reverse direction to fork movement, the replicative enzyme carries out chain elongation continuously in a highly processive manner. On the other parent strand, the lagging strand, which runs 3' → 5' in the direction of fork movement, DNA polymerase catalyzes chain elongation only in the reverse direction of fork movement. Thus, as the replication fork proceeds, the unreplicated segment is expanded on the lagging strand. When one act of chain elongation on the lagging strand is accomplished, the next round of chain elongation must be started from a newly expanded segment on the lagging strand. To achieve completion of DNA replication on an entire region of the lagging strand, numerous [enzymes](#) cooperate at the replication fork. They are participating in [primer](#) synthesis (the initiation of Okazaki fragments), chain elongation (the extension of Okazaki fragments), and a process that connects Okazaki fragments.

In *Escherichia coli*, more than 20 different proteins participate in DNA replication (2). These were identified by screening mutants defective in DNA replication and by purifying enzymes required for *in vitro* DNA synthesis (see [dna genes](#)). From their biochemical roles at different stages of chromosomal DNA replication, it appears that at least eight proteins are involved in the discontinuous replication in *E. coli*. Those are (1) primosome proteins, including [DNA helicases](#) and [primase](#): PriA, PriB, PriC, DnaT, DnaB, DnaC, and DnaG (primase); (2) proteins required for chain elongation: DNA polymerase III (Pol III) holoenzyme and [single-stranded DNA binding protein](#) (SSB); and (3) proteins required for connecting Okazaki fragments: DNA polymerase I (Pol I), RNaseH, and [DNA Ligase](#). Among these proteins, the DnaB helicase, primase (DnaG), and Pol III holoenzyme are the basic components acting on the discontinuous DNA synthesis at the replication fork, probably forming a multiprotein complex called a *replisome*.

DNA polymerase cannot replicate duplex DNA without assistance. This enzyme requires single-stranded DNA as a [template](#) and an RNA or DNA primer annealed to the template. Two other enzymes enable polymerase to work on duplex DNA. One is a DNA helicase that opens up the duplex at the replication fork to provide a single-stranded template. The other is a primase that synthesizes a short RNA to prime DNA chain elongation. Several DNA helicases have been identified from *E. coli* and its phage, including DnaB, T7 gp4, and T4 gp41 (3-5). The biochemical characterization of these activities in *in vitro* DNA replication systems suggests that the primary replicative helicase binds to and moves on the lagging-strand template in the 5' → 3' direction, unwinding the DNA double helix as it goes. Another common property of the replicative helicases is an intimate association with a primase. The bacteriophage T7 gp4 has both primase and helicase activity within the same [polypeptide chain](#) (4). Bacteriophage T4 gp41 greatly enhances the primase activity of T4 gp61 (6). A similar functional interaction has been observed between DnaB protein and *E. coli* primase, DnaG protein (7). On most templates, DnaG exhibits a very feeble priming activity that can be greatly enhanced if DnaB first binds that DNA. This stimulation of primase activity is further increased when the DnaB helicase is activated to its processive form at the replication fork.

In *E. coli*, there are two pathways by which the DnaB helicase is loaded onto the lagging-strand template DNA: One is primosome formation directed by the PriA protein, and the other is DnaA protein-directed DnaB loading (8, 9). The former process was discovered first in the replication of bacteriophage fX174 DNA and **plasmid** ColE1 DNA, and it later appeared to be involved in the resumption of chromosomal DNA replication after replication of the *E. coli* genome has been interrupted or halted. On the other hand, the latter process was found in *oriC* plasmid DNA replication *in vitro* and is thought to form a priming complex with DnaG primase at the replication fork in *E. coli* chromosomal DNA replication.

The chain elongation of Okazaki fragments in *E. coli* is catalyzed by DNA polymerase III holoenzyme (10). This enzyme possesses a capacity to synthesize DNA with a very high processivity, sufficient for completion of about 2 kb of Okazaki fragment. In addition, the Pol III holoenzyme dissociates from the nascent Okazaki fragment and restarts the next round of Okazaki fragment synthesis from an RNA primer newly settled near the replication fork (11). Enzymes to remove primer RNA and fill the gap, such as [ribonuclease H](#) and DNA polymerase I of *E. coli*, are essential for the sealing of Okazaki fragments by DNA ligase (12). Mutants defective in either DNA polymerase I or DNA ligase show a massive accumulation of short Okazaki fragments under restrictive conditions.

Although the basic biochemical processes that occur at eukaryotic and prokaryotic replication forks are similar, there are many differences in detail (13). For example, primer synthesis in eukaryotic cells is catalyzed by DNA polymerase  $\alpha$ , which synthesizes 2 to 12 nucleotides of RNA (initiator RNA) and further adds about 20 nucleotides of DNA to the initiator RNA. The size of Okazaki fragments (40 to 300 nucleotides) in eukaryotes is significantly shorter than those observed in prokaryotes.

### Bibliography

1. R. Okazaki, T. Okazaki, S. Hirose, A. Sugino, and T. Ogawa (1976) In *DNA Synthesis and Its Regulation* (M. Goulian and P. Hanawalt, eds.), Benjamin Cummings, Menlo Park, CA, pp. 832–862.
2. K. J. Marians (1992) *Ann. Rev. Biochem.* **61** 673–720.
3. J. H. LeBowitz and R. McMacken (1986) *J. Biol. Chem.* **261**, 4738–4748.
4. J. A. Bernstein and C. C. Richardson (1989) *J. Biol. Chem.* **264**, 13066–13073.
5. M. Venkatesan, L. L. Silver, and N. G. Nossal (1982) *J. Biol. Chem.* **257**, 12426–12434.
6. D. M. Hinton and N. G. Nossal (1987) *J. Biol. Chem.* **263**, 10873–10878.
7. K. Arai and A. Kornberg (1979) *Proc. Natl. Acad. Sci. USA* **76**, 4308–4312.
8. K. Arai, R. Low, J. Kobori, J. Shlomai, and A. Kornberg (1981) *J. Biol. Chem.* **256**, 5273–5280.
9. T. A. Baker, B. E. Funnell, and A. Kornberg (1987) *J. Biol. Chem.* **262**, 6877–6885.
10. C. McHenry and A. Kornberg (1977) *J. Biol. Chem.* **252**, 6478–6484.
11. H. Maki, S. Maki, and A. Kornberg (1988) *J. Biol. Chem.* **263**, 6570–6578.
12. B. E. Funnell, T. A. Baker, and A. Kornberg (1986) *J. Biol. Chem.* **261**, 5616–5624.
13. J. J. Blow ed. (1996) *Eukaryotic DNA Replication*, IRL Press, Oxford.

### Suggestions for Further Reading

14. M. L. DePamphilis, ed. (1996) *DNA Replication in Eukaryotic Cells*, Cold Spring Harbor Laboratory Press, New York.
15. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.

### Disease Resistance Genes, Plant

Plants are the ultimate food source for many nonphotosynthetic organisms. In natural ecosystems, plant populations are usually too fragmented for major epidemics of bacteria, fungi, **viruses**, [nematodes](#), or other animals to threaten the survival of a plant species, although such pathogens are

still likely to select for variants with enhanced resistance. In the monocultures imposed by human cultivation, disease can be severe.

Plant pathogens are classified into necrotrophs and biotrophs (1). Necrotrophs overwhelm plant defenses with toxins and hydrolytic [enzymes](#) and then feed on the dead plant material. Biotrophs need a living plant to complete their life cycle and are better thought of as parasites; they can reduce plant yields substantially, without killing their hosts. Hemibiotrophs, such as *Phytophthora infestans* which causes potato late blight, have an initial biotrophic phase that rapidly switches to necrotrophy.

The interaction between plant and pathogen is said to be compatible if disease occurs, and incompatible if it does not. Most plants are resistant to most plant pathogens; the flax rust fungus does not cause disease on wheat, and vice versa. In this example of “nonhost” resistance, incompatibility can be due to a lack of “basic compatibility”; biotrophic parasites must specialize to grow on their specific hosts.

One of the most interesting and widely studied forms of disease resistance appears to involve recognition of pathogens by plants, followed by activation of defense responses that usually culminate in the death of plant cells near where the pathogen attempts ingress. This localized cell death is often referred to as the *hypersensitive response*, or HR. It should be noted that plants, unlike mammals, lack circulating defender cells, and every plant cell must be able to mount a defense response. The HR is correlated with the production of antimicrobial proteins and reactive oxygen species (2, 3).

Soon after the rediscovery of Mendel's work, breeders noted that resistance in wheat to rusts was inherited as a single dominant Mendelian trait. In the 1940s, H. H. Flor, studying the flax/flax rust interaction, carried out genetic analysis with both the rust fungus and its plant host. **Genes** in flax for rust resistance were dominant, as in wheat. Flor also showed that virulence in the pathogen segregated as a recessive trait. From this he inferred that each pathogen must carry dominant **alleles** of genes which permitted the plant to detect that pathogen race, and he termed these *avirulence* genes. Genetic interactions such as these are often referred to as “gene-for-gene.”

Flor's work underpins the current view that resistance (*R*) genes encode receptors which recognize pathogen-derived, avirulence (*Avr*) gene-encoded ligands. Pathogens evade recognition by mutations in *Avr* genes. Why would pathogens make *Avr* gene products that permit plants which carry the corresponding *R* gene to detect them? It is believed that pathogens generate a battery of pathogenicity factors (also known as compatibility factors) which are individually dispensable. Pathogens (and plant breeders) select for variant *R* genes that can detect one of these factors. This interpretation appears to be valid for interactions of plants with bacteria, viruses, fungi, and nematodes (4).

#### 1. Four Classes of “Gene-for-Gene” *R* Gene Products

The first gene-for-gene *R* gene to be isolated was *Pto* from tomato (5). *Pto* confers resistance to races of *Pseudomonas syringae* that carry the corresponding *AvrPto* gene, and encodes a **serine-threonine protein kinase** of the Raf family, with strongest homology to the *Drosophila* gene *Pelle*, which is involved in [signal transduction](#) between the *Toll* receptor and **Dorsal** transcription factor, and it may also be involved in *Drosophila* innate immunity (6, 7). *Pto* also carries a putative N terminal [myristoylation](#) site. It was surprising to discover that a gene putatively involved in recognition showed [homology](#) to genes involved in signal transduction rather than to known receptors. Remarkably, recent work shows that there is probably a direct interaction between the *Pto* and *AvrPto* gene products (8, 9). How could *AvrPto* physically interact with *Pto*? Many phytopathogenic *Pseudomonas* and *Xanthomonas* species encode a type III secretion pathway that actively delivers compatibility factors to the plant cell (10, 11). The same secretion pathway is involved in delivering proteins into mammalian cells attacked by *Shigella*, *Yersinia*, *Salmonella*, or pathogenic *Escherichia coli*. Mutations in the so-called *hrp* genes specifying this apparatus in *Pseudomonas* and

*Xanthomonas* leave the bacterium defective in both Hypersensitive Response elicitation and Pathogenicity. The *hrp* system is believed to deliver several pathogenicity factors from bacteria to plant cells, some of which (such as AvrPto) the plant evolves the capacity to recognize. The *Pto* gene is a member of a linked **gene family** of five to seven members spread over probably less than 100 kbp, one member of which, *Fen*, confers an apparent HR when treated with Fenthion insecticide. Experiments swapping **domains** between *Fen* and *Pto* have precisely localized the region involved in *AvrPto* recognition (8, 9).

The other three classes of *R* genes carry the [leucine-rich repeat](#) (LRR) motif. LRRs have been implicated in various examples of [protein–protein interactions](#) and protein–peptide interactions (12). This repeating motif creates a parallel [beta-sheet](#) that provides a recognition surface in the interaction of porcine **ribonuclease inhibitor** (PRI) with ribonuclease. Thus, the presence of LRRs is consistent with a role for *R* gene products as receptors for pathogen-specified ligands. Based on the PRI structure, a model has been proposed for the structure of various LRR proteins, including plant *R* gene products (13).

The *Cf-2*, *Cf-4*, *Cf-5* and *Cf-9* genes of tomato confer resistance to specific races of the fungal pathogen *Cladosporium fulvum* (14-16). This biotrophic fungus proliferates in the intercellular spaces of infected leaves prior to conidiation. It secretes small cysteine-rich peptides into these spaces, presumed to be compatibility factors, some of which are *Avr* gene products (17). Avr9, recognized by *Cf-9*, is a 28 amino-acid residue peptide; after secretion, Avr4 has 105 residues. The *Cf-* genes all encode **glycoproteins** that are **membrane-anchored** at the C-terminus (probably to the plasma membrane) with 25 to 38 LRRs. Binding studies showed that Avr9 binds a protein that is present even in stocks that lack *Cf-9* (18). Conceivably, this protein undergoes a conformational change upon binding Avr9, and this change is recognized by *Cf-9*. No evidence for direct interaction between purified *Cf-9* and Avr9 has yet been obtained.

Analysis of the *Cf*-gene family has revealed interesting insights into *R* gene [evolution](#) (19). At the *Cf-9* locus, there are five homologues of *Cf-9* (including *Cf-9*). At the homologous *Cf-4* locus, there are also five members, including *Cf-4*. Comparison of these sequences has revealed a pattern of evolution by sequence exchange, presumably by either **unequal crossing over** or **gene conversion**. This is consistent with the analysis of the maize *Rp1* locus, which provides genetic evidence that novel alleles arise correlated with meiotic [recombination](#) events (20). Analysis of the ratio of synonymous to nonsynonymous codon substitutions ( $K_s/K_a$ ) suggests diversifying changes in the solvent-exposed amino-acid residues of the LRR parallel beta-sheet. Consistent with this, several other members of the *Cf*-gene family confer resistance through recognition of as-yet-uncharacterized avirulence genes. A complementary study provides evidence that several of the secreted peptides of *C. fulvum* have the potential to act as *Avr* gene products (21). The picture emerging from this system reveals a battery of secreted peptides that presumably confer a pathogenicity function for *C. fulvum*, with the plant evolving the capacity to recognize and react to them through creating a recognition specificity in a *Cf*-gene homologue. How the recognition actually occurs, and what happens next, is an area of active research.

No other *R* genes have been shown to be members of this C-terminally anchored LRR class, although several [Arabidopsis](#) homologues have been identified, where in some instances they exist as linked multigene families. The nematode resistance gene *Hc1<sup>pro</sup>* is reported to be of this class (22), although the LRR sequences vary from the canonical LRR motif.

The rice *Xa21* gene confers resistance to *Xanthomonas campestris pv oryzae*, the causal agent of bacterial blight. Interestingly, it encodes a transmembrane LRR protein [kinase](#), combining features of *Pto* and *Cf-9* in one protein (23). This would suggest that the *AvrXa21* gene encodes a secreted bacterial protein, but thus far the *AvrXa21* gene has not been cloned, because little genetic variation has been detected within *Xanthomonas* at this locus. Again, *Xa21* is a member of a linked multigene family (24).

The largest class of *R* genes encode cytoplasmic (but probably membrane-associated) proteins that carry a canonical **nucleotide-binding** (NB) site and C-terminal stretch of LRRs. The NB carries the *kin1a* ([P loop](#)), *kin2*, and *kin3a* sites of protein kinases ([25](#)). In addition, recent comparisons have revealed a more extended homology with a domain of the *Caenorhabditis elegans Ced4* and the human *Apaf1* gene products, which led some authors ([26](#)) to christen this extended region the NB-ARC domain (standing for Nucleotide Binding, *Apaf*, *R* gene, and *Ced4*). *Apaf1* and *Ced4* encode proteins that are effectors of [apoptosis](#). The parallel between apoptosis and the HR has been the subject of considerable interest ([27, 28](#)), and the homologies are intriguing, but no biochemical evidence is yet available to substantiate it.

These NB-LRR genes fall into two main classes:

1. *Rps2* and *Rpm1*, *Arabidopsis* genes for resistance to *Pseudomonas* strains carrying *AvrRpt2* and *AvrRpm1*, respectively, encode a putative [leucine zipper](#) (LZ) N-terminal to the NB region ([29, 30](#)). Unusually, *Rpm1* recognizes both *AvrB* and *AvrRpm1*.
2. In contrast, the tobacco *N* gene for [tobacco mosaic virus](#) (TMV) resistance, the flax *L6* gene for rust resistance, and the *Arabidopsis RPP5* gene for downy mildew resistance carry a region at the N-terminus with homology to the cytoplasmic domain of the *Toll* receptor and *interleukin 1* receptor ([31, 33](#)). This region has recently been proposed to adopt a structure resembling the *E. coli CheY chemotaxis*-response element ([34](#)), and in plants it has been dubbed the TIR domain (standing for *Toll*, the *interleukin 1* receptor, and *resistance* gene) ([4](#)). It is striking that *R* genes from different plants conferring resistance to such distinct pathogens are so similar. The *L6* gene product carries an N-terminal putative signal anchor, which could attach it to the plasma membrane.

Additional NB-LRR subclasses may exist. The tomato *I2* gene that confers *Fusarium* wilt resistance carries a putative leucine zipper between the LRRs and the NB <sup>35</sup>. The tomato *Prf* gene was identified by mutagenesis as required for both *Pto* and *Fen* function; it encodes an extremely long region N-terminal to the NB but is clearly of the NB-LRR class. ([36](#)).

Paradoxically, even though LRR domains are more correlated with protein-protein interactions than protein kinase domains, no direct interaction has been detected between any of these NB-LRR *R* gene products and their genetically defined candidate ligands using yeast [two-hybrid systems](#). Nevertheless, transient expression experiments showed that *Rps2*([37](#)) and *Rpm1* ([38](#)) must recognize *AvrRpt2* and *AvrB* within the plant cell. *AvrBs3*, recognized by the as-yet-uncharacterized *Bs3* *R* gene product of pepper, is also recognized within the plant cell ([39](#)).

It is widely believed that the *R* gene products form part of a complex with other proteins to create a functional recognition and signal transduction machine. The nature of this complex is currently under vigorous investigation. Many other questions also remain unanswered. What are the bacterial *Avr* products doing in the compatible interaction? Are fungal resistance genes of the NB-LRR type recognizing fungal proteins that are delivered into the plant cell? A combination of genetic and biochemical approaches are being used to understand *R* gene product-dependent recognition and signal transduction events. *NDR1*, identified by mutagenesis as required for *Rps2* function, encodes a novel protein with two membrane-spanning domains ([40](#)). *EDS1*, required for *RPP5* function but not *Rps2* function ([41](#)), shows homology to [lipases](#) (J Parker, personal communication). Yeast two-hybrid analysis has led to the discovery of genes whose products interact with *Pto*; some are protein kinases and some are [transcription factors](#) ([42, 43](#)).

## 2. Resistance to Pathogen Toxins

Some pathogens produce toxins, and resistance is conferred by a dominant gene that encodes a detoxifying enzyme. The maize *Hm1* gene encodes an NADPH oxidoreductase that inactivates the



*Cochliobolus carbonum* HC toxin (44). *Hm1* was the first *R* gene to be isolated, but it does not conform to gene-for-gene genetics in that recessive mutations in the pathogen to overcome *Hm1* have never been recovered. In the 1970s, many maize varieties carried the T cytoplasm, whose mitochondria are sensitive to T toxin from *Cochliobolus heterostrophus*, the causative agent of Southern Corn leaf blight. Resistant varieties that use a different cytoplasmic male sterility genotype were then deployed, which controlled the disease.

### 3. *Mlo* Recessive Resistance Genes

There are several recessive resistance genes, of which the best characterized is barley *mlo*. *Mlo* mutants have a disease lesion-mimic phenotype and enhanced resistance to all races of the powdery mildew pathogen, *Erysiphe graminis*. The *Mlo* gene encodes a protein with seven transmembrane domains that presumably regulate negatively the defense response (45). It is also the pioneer of a new family of plant genes; analysis of the *Arabidopsis* genome sequence suggests that it will contain at least 60 genes in this class, with as yet unknown biochemical function.

### 4. Summary

Different classes of *R* genes in plants encode either recognition and response systems, or toxin detoxification or tolerance mechanisms, or they derepress plant defense mechanisms. Many *R* genes have now been isolated; the challenge is now to find out how their gene products work.

### Bibliography

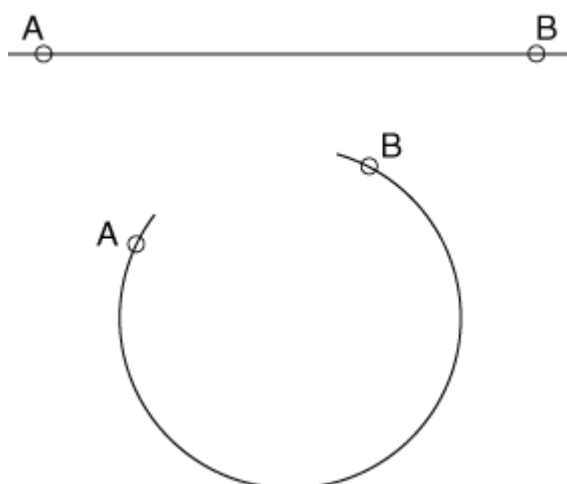
1. G. N. Agrios (1988) *Plant Pathology*, Academic Press, London.
2. K. E. Hammond-Kosack and J. D. G. Jones (1996) *Plant Cell* **8**, 1773–1791.
3. C. Lamb and R. A. Dixon (1997) *Ann. Rev. of Plant Physiol. Plant Mol. Biol.* **48**, 251–275.
4. K. E. Hammond-Kosack and J. D. G. Jones (1997) *Ann. Rev. Plant Physiol. Plant Mol. Biol.* **48**, 575–607.
5. G. B. Martin et al. (1993) *Science* **262**, 1432–1436.
6. M. P. Belvin and K. V. Anderson (1996) *Ann. Rev. Cell Develop. Biol.* **12**, 393–416.
7. J. A. Hoffmann and J. M. Reichhart (1997) *Trends Cell Biol.* **7**, 309–316.
8. S. R. Scofield et al. *Science* **274**, 2063–2065.
9. X. Tang et al. *Science* **274**, 2060–2063.
10. J. R. Alfano and A. Collmer (1996) *Plant Cell* **8**, 1683–1698.
11. G. VandenAckercken and U. Bonas (1997) *Trends Microbiol.* **5**, 394–398.
12. B. Kobe and J. Deisenhofer (1995) *Curr. Opin. Struct. Biol.* **5**, 409–416.
13. A. V. Kajava (1998) *J. Molec. Biol.* **277**, 519–527.
14. C. M. Thomas et al. (1997) *Plant Cell* **9**, 2209–2224.
15. D. A. Jones, C. M. Thomas, K. E. Hammond-Kosack, P. J. Balint-Kurti, and J. D. G. Jones (1994) *Science (Wash.)* **266**, 789–793.
16. M. S. Dixon et al. (1996) *Cell* **84**, 451–459.
17. P. DeWit (1997) *Trends Plant Sci.* **2**, 452–458.
18. M. Kooman-Gersmann, G. Honee, G. Bonnema, and P. J. G. M. De Wit (1996) *Plant Cell* **8**, 929–938.
19. M. Parniske et al. *Cell* **91**, 821–832.
20. S. H. Hulbert (1997) *Ann. Rev. Phytopathol.* **35**, 293–310.
21. R. Lauge, A. P. Dmitriev, M. Joosten, and P. DeWit (1998) *Molec. Plant-Microbe Inter.* **11**, 301–308.
22. D. G. Cai et al. *Science* **275**, 832–834.

23. W.-Y. Song et al. *Science* **270**, 1804–1806.
24. W. Y. Song, et al. (1997) *Plant Cell* **9**, 1279–1287.
25. T. W. Traut (1994) *Europ. J. Biochem.* **229**, 9–19.
26. E. A. vanderBiezen and J. D. G. Jones (1998) *Curr. Biol.* **8**, R226–R227.
27. J. B. Morel and J. L. Dangl (1997) *Cell Death and Differ.* **4**, 671–683.
28. A. M. Jones and J. L. Dangl (1996) *Trends Plant Sci.* **1**, 114–119.
29. A. F. Bent et al. *Science* **265**, 1856–1860.
30. M. R. Grant et al. *Science* **269**, 843–846.
31. G. J. Lawrence, E. J. Finnegan, M. A. Ayliffe, and J. G. Ellis *Plant Cell* **7**, 1195–1206.
32. J. E. Parker et al. (1997) *Plant Cell* **9**, 879–894.
33. S. Whitham et al. (1994) *Cell* **78**, 1011–1115.
34. F. L. Rock, G. Hardiman, J. C. Timans, R. A. Kastelein, and J. F. Bazan (1998). *Proc. Natl. Acad. Sci.* **95**, 588–593.
35. N. Ori et al. (1997) *Plant Cell* **9**, 521–532.
36. J. M. Salmeron et al. *Cell* **86**, 123–133.
37. R. T. Leister, F. M. Ausubel, and F. Katagiri (1996) *Proc. Natl. Acad. Sci.* **93**, 15497–15502.
38. S. Gopalan et al. (1996) *Plant Cell* **8**, 1095–1105.
39. G. VandenAckercken, E. Marois, and U. Bonas (1996) *Cell* **87**, 1307–1316.
40. K. S. Century et al. (1997) *Science* **278**, 1963–1965.
41. J. E. Parker et al. (1996) *Plant Cell* **8**, 2033–2046.
42. J. M. Zhou, X. Y. Tang, and G. B. Martin (1995) *EMBO J.* **16**, 3207–3218.
43. J. Zhou, Y.-T. Loh, R. A. Bressan, and G. B. Martin (1995) *Cell* **83**, 925–935.
44. G. S. Johal and S. P. Briggs (1992) *Science (Wash.)* **258**, 985–987.
45. R. Buschges et al. *Cell* **88**, 695–705.

## Distance Geometry

Many physical methods used to examine a biological [macromolecule](#) in solution produce information about the distance that separates two parts of the molecule. These methods include fluorescence **energy transfer**, chemical [cross-linking](#), [light scattering](#), and several kinds of nuclear magnetic resonance ([NMR](#)) observations. Any information about distances between different parts of a molecule limits the number of conformations that are possible for that molecule. Distance geometry is a mathematical approach that can be used to help define the conformations of a molecule that are consistent with experimental distance information ([1-6](#)).

With [X-ray crystallography](#) methods, the tertiary structure of the molecule of interest is defined by analysis of a set of diffraction data. A relation exists between every feature of the diffraction pattern exhibited by the crystal and the positions of the atoms of molecules within the crystal. When this connection is elucidated, the three-dimensional (3D) structure of the molecule is revealed. A different approach is used to find 3D structures using distance information. Consider a polymeric molecule in an extended conformation, represented by the line below.



If free rotation exists about some or all of the bonds along the molecular backbone, a huge number of conformations are possible for the molecule. If, however, an experimental observation indicates, for example, that the atoms at the position labeled A must be within 0.3 nm of the atoms at the position labeled B, then the number of possible conformations for the molecule is greatly reduced—the only conformations that need further consideration are those that are consistent with the constraint that atoms at A and B must be close to each other. Other experimentally defined distances would further reduce the number of possible 3D structures. If enough such constraints exist, only one or a small family of structures is likely that is consistent with all of them. Generally, the more constraints that can be established by experiment, the better known will be the [tertiary structure](#). NMR methods, particularly those that measure [nuclear Overhauser effects](#) (NOEs), can produce a large number of distance constraints on possible conformations of a macromolecule. An advantage in using NMR observations in this way is that the molecule of interest does not have to be in the solid (crystalline) state for successful structure determination.

An  $N \times N$  matrix of distances can be defined for a molecule made up of  $N$  atoms;  $N(N-1)/2$  unique interatomic distances exist in this matrix. If all of the distances are accurately known, mathematical methods are available to define the Cartesian coordinates of the atoms from the known distances (4). Many of the distances in the distance matrix are essentially independent of conformation and are reliably known *a priori*. Such data would include the distance between two hydrogen atoms attached to an aromatic ring. Other distances can be estimated experimentally. Information about [dihedral angles](#) and other bond angles can be used to define other distances. For some distances, however, little or no accurate distance information will be available. Thus, the distance matrix will typically be incomplete, and a unique solution for Cartesian coordinates cannot be obtained. A variety of algorithms have been developed to deal with the missing or incomplete distance information. Regardless of the approach used, distance geometry calculations typically do not produce a single conformation but an ensemble of conformations, all of which are consistent with the available distance constraints. These conformations are used as starting points for additional refinement by conformational energy minimization and [simulated annealing](#) procedures. (See also [Nuclear Overhauser Effect \(NOE\)](#), [Simulated Annealing](#)).

#### Bibliography

1. L. M. Blumethal (1970) *Theory and Applications of Distance Geometry*, Chelsea, New York.
2. G. M. Crippen and T. F. Havel (1978) *Acta Crystallogr.* **34**, 282–284.
3. G. M. Crippen (1981) *Distance Geometry and Conformational Calculations*, Wiley, Chichester, England.
4. W. Braun (1987) *Quart. Rev. Biophys.* **19**, 115–117.
5. I. D. Kuntz, J. F. Tomason, and C. M. Oshhiro (1989) *Methods Enzymol.* **177**, 159–204.

6. T. F. Havel (1991) *Prog. Biophys. Mol. Biol.* **56**, 43–78.

### **Suggestions for Further Reading**

7. C. M. Oshiro and I. D. Kuntz (1993) *Biopolymers* **33**, 107–115.

8. J. W. Shriver and S. Edmondson (1994) *Methods Enzymol.* **240**, 415–438.

9. A. M. J. J. Bonvin and A. T. Brunger (1996) *J. Biomol. NMR* **7**, 72–76.

10. P. Güntert (1997) In *Protein NMR Techniques* (D. G. Reid, ed.), Humana, Totowa, New Jersey, pp. 157–194.

11. D. M. LeMaster (1997) *J. Biomol. NMR* **9**, 79–93.

12. J. Cavanagh, W. J. Fairbrother, A. G. Palmer III and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.

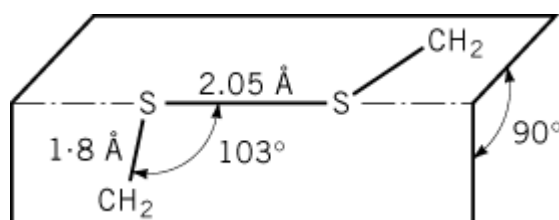
## Disulfide Bonds

The disulfide bond occurs in many biological molecules, most frequently in **secretory proteins**, in the form of a covalent linkage between the **thiol groups** of pairs of **cysteine** residues (Fig. 1). Formation of a disulfide bond from two thiol groups requires an oxidant as two electrons or hydrogen atoms are released:



The oxidant can be oxygen, but air oxidation of thiol compounds is a complex reaction involving many steps and intermediates, being catalyzed by metal ions, and releasing by-products such as peroxides.

**Figure 1.** Optimal stereochemistry of a disulfide bond between two cysteine residues. The two possible  $\text{CH}_2\text{-S-S-CH}_2$  torsion angles about the disulfide bond of  $+90$  and  $-90$  are equally favorable.

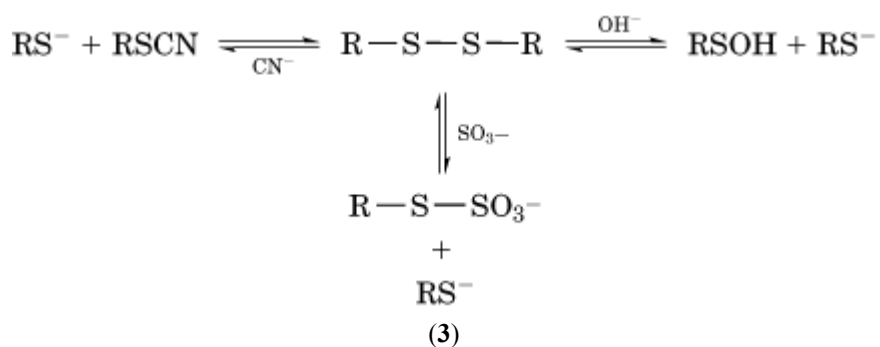


The disulfide bond is a stable covalent bond, and it would require some 60 kcal/mol (270 kJ/mol) to pull the two sulfur atoms apart. Yet this stable bond can be broken readily by a thiol group, which can react rapidly with a disulfide bond in the **thiol-disulfide exchange** reaction. With this reaction, disulfide bonds can be reduced by reacting them with an excess of a thiol reagent, RSH, such as **[β-mercaptoethanol](#)** or **[dithiothreitol](#)**:

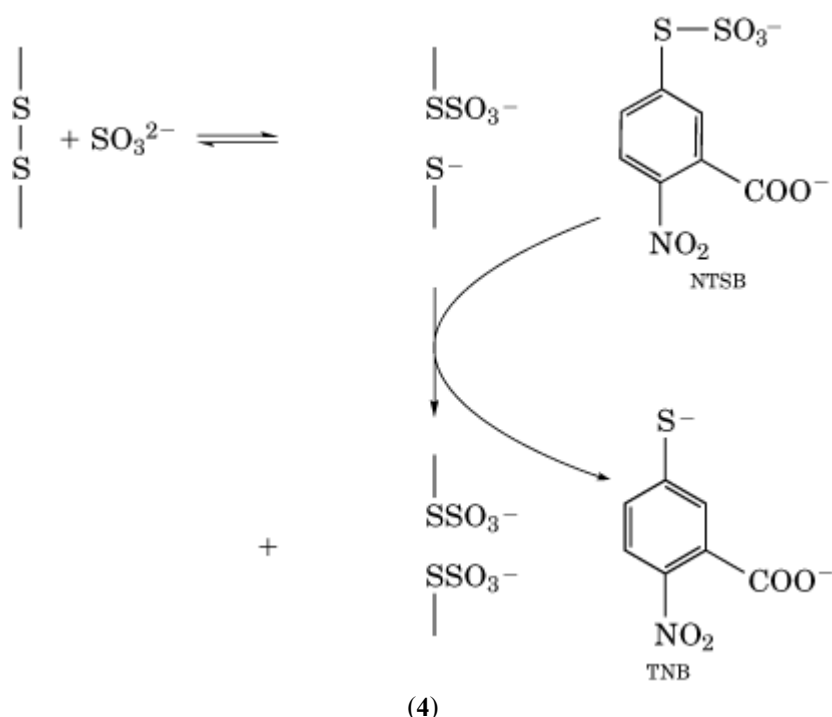


Many other reagents are available for chemically reducing disulfide bonds, such as borohydride, tributylphosphine, and tris-[2-carboxyethyl]-phosphine. The phosphines have the advantage that they function at acidic pH, where the original disulfide bond and the resulting thiol groups tend not to undergo further reactions. Furthermore, the phosphines do not react with many of the reagents that modify thiol groups, so reduction of the disulfide bond and alkylation of the thiol groups produced can be carried out simultaneously (1).

Disulfide bonds can also be cleaved by nucleophiles such as cyanide, sulfite, or hydroxide ion:



The equilibria for these reactions are such that it is difficult to drive the reactions to completion simply by adding reagent. Complete reaction can be accomplished, however, by adding a second reagent that reacts with the thiol generated. When this reagent is nitrothiosulfobenzoate (NTSB),



the reaction can go to completion, and one mole of the colored TNB product is generated for each mole of disulfide bond originally present. This is the most convenient assay for protein disulfide bonds so long as thiol groups are not also present.

The intrinsic free energy of a disulfide bond depends upon its angle of rotation (Fig. 1). It is most stable with a C-S-S-C torsion angle of  $\pm 90^\circ$ , and an angle of  $0^\circ$  or  $180^\circ$  is believed to increase the free energy by about 28 kJ/mol. This corresponds to the free energy barrier to full rotation about the disulfide bond. Cyclic 5-membered rings, such as that of lipoic acid, have a torsion angle of only  $\pm 30^\circ$  and are strained by about 15 kJ/mol. Cyclic 6-membered rings, such as that in the disulfide form of **dithiothreitol**, have torsion angles of about  $\pm 60^\circ$ .

## 1. Disulfide Bonds in Proteins

Disulfide bonds are found most frequently in molecular biology within proteins, especially in proteins that are secreted. Most of the disulfide bonds observed in the [X-ray crystallography](#) structures of proteins have favorable stereochemistry and show little evidence of conformational strain, with close to the optimal torsion angles of  $\pm 90^\circ$ . Only a few examples of hyperreactive disulfide bonds are known, such as that between the most distant cysteine residues in [alpha-lactalbumin](#) and hen [lysozyme](#) (2, 3).

As the disulfide bond is a covalent bond that is stable under appropriate conditions, which cysteine

residues in a protein are paired in disulfide bonds can be determined chemically. The classical method is to fragment the protein into peptide fragments, under conditions where the disulfide bonds are kept intact. The peptide fragments are analyzed to determine which are linked by disulfide bonds. A newer technique is to reduce the disulfide bonds individually, using phosphines at acidic pH, where the disulfide bonds do not rearrange readily, and to identify the pairs of cysteine thiol groups generated (1, 4). A more elegant approach, but one not suited to modern separation techniques, is to use [diagonal methods](#) (5).

Disulfide bonds are well known to stabilize the folded conformations of proteins in which they occur, and conversely the folded conformations stabilize the disulfide bond. Many proteins that have disulfide bonds unfold when the disulfides are reduced. In considering disulfide bonds and protein stability, it is important to distinguish between the stability of the folded conformation relative to the unfolded protein,  $U$ , (i) when the disulfide bonds are kept intact:

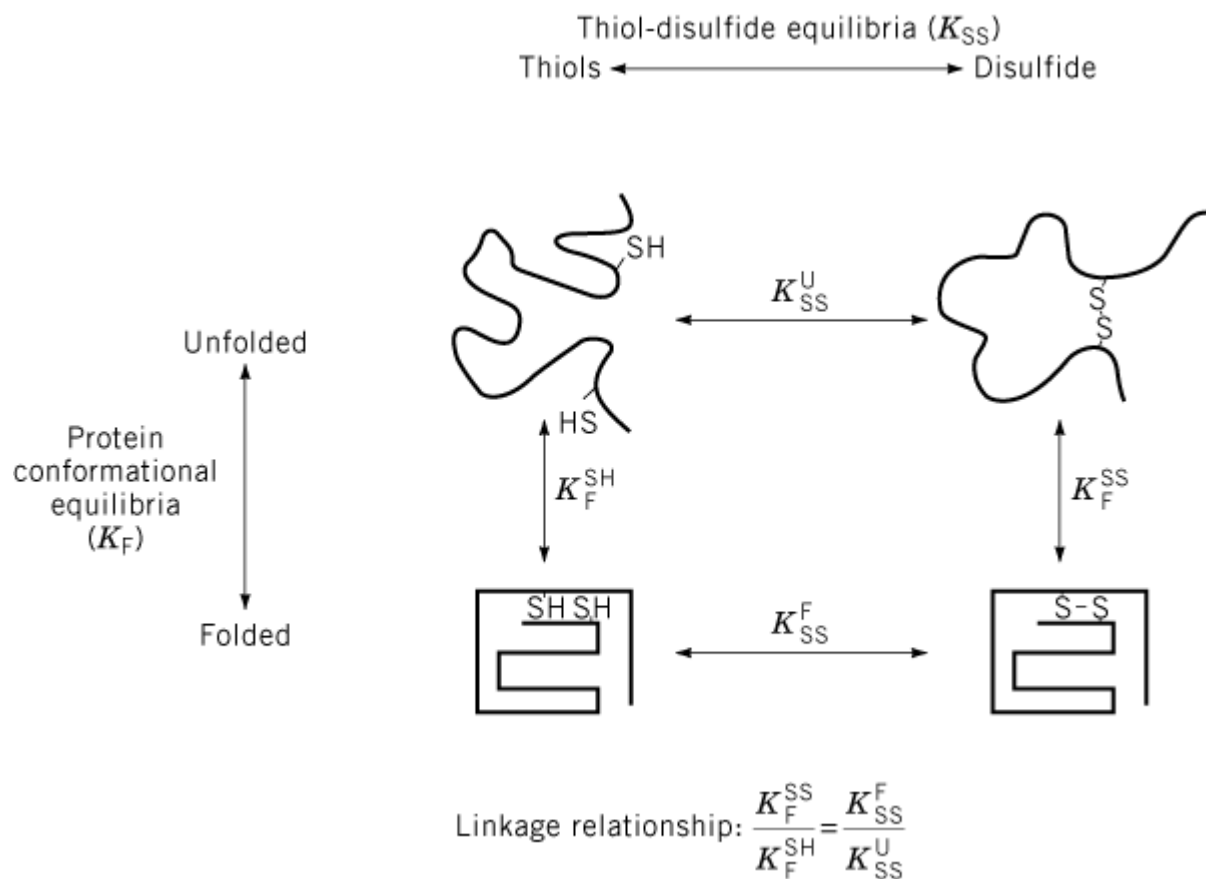


and (ii) when the disulfide bonds are reduced upon unfolding and reformed when the protein folds:



The stability of a protein disulfide bond, relative to the reduced form of the protein, with cysteine thiol groups (Eq. 6), depends upon the degree to which the protein conformation stabilizes the disulfide bond and upon the redox potential of the environment. A protein conformation that stabilizes a disulfide bond is stabilized to the same extent by the presence of that disulfide bond (Fig. 2). The greater the stability of the protein disulfide bond, the greater its contribution to stability of the folded conformation that stabilizes it.

**Figure 2.** Linkage between the stability of a particular folded conformation that favors a certain disulfide bond and the stability of that disulfide bond. A protein with two cysteine residues that can form a disulfide bond (right) is illustrated in either the unfolded (top) or folded (bottom) conformations. The indicated equilibrium constants represent the stabilities of the disulfide bonds,  $K_{SS}$ , and of the folded conformation,  $K_F$ . The equilibrium constants for disulfide stability,  $K_{SS}$ , must involve another reagent, as in Eq. 7 and 8, but this is omitted for simplicity. The linkage relationship states that whatever effect the folded conformation has on the stability of the disulfide bond, the presence of the disulfide bond must have the same quantitative effect on the stability of the folded conformation. Comparable linkage relationships pertain to all interactions within the folded conformation, not just disulfide bonds.



One way in which disulfide bonds stabilize the folded conformation, relative to the unfolded form with the disulfide bonds intact (Eq. 5), is by decreasing the conformational flexibility, or entropy, of the [unfolded protein](#). The magnitude of this effect increases with increasing distance between the cysteine residues in the [primary structure](#). It is often assumed that this is the only way that disulfide bonds stabilize the folded conformation, and equations are frequently given to predict the magnitude of the stabilization. Experimental observations indicate, however, that there are many other effects of a disulfide bond on the folded conformation, both stabilizing and destabilizing (6); the same energetically favorable disulfide bond can stabilize one folded conformation and destabilize another (7).

The redox potential determines the stability of any disulfide bond and is determined by the relative concentrations of thiol and disulfide reagents. In biological systems, this is primarily determined by the relative concentrations of the thiol and disulfide forms of [glutathione](#), GSH and GSSG, respectively:



$$\frac{[P_S^S]}{[P_{SH}^{SH}]} = K_{eq} \frac{[GSSG]}{[GSH]^2} \quad (8)$$

Within the cell cytosol, the concentration of GSH is generally 100-fold greater than that of GSSG, which is a moderately reducing environment, and a disulfide bond does not impart much stability. Within the [endoplasmic reticulum](#), however, where secreted proteins fold and form disulfide bonds, the concentrations of GSH and GSSG are almost comparable (8), so any protein disulfide bond will



be present about 100 times more frequently. These conditions are more oxidizing, but only stable protein disulfide bonds are generally formed and present under these conditions. Once the proteins are secreted from the cell, where the environment is usually much more oxidizing, the disulfide bonds contribute much more stability to the folded conformation.

Disulfide bonds are inserted into proteins naturally only when the protein conformation favors them. When the reduced protein, without disulfides, is unfolded, folding accompanies disulfide bond formation and can be used to determine the folding pathway (see [Protein Folding In Vivo](#)). Such pathways have demonstrated how only a slight tendency for an unfolded protein to stabilize a specific disulfide bond is consequently stabilized further by the presence of that disulfide bond, which can then lead to the formation of further disulfide bonds and the adoption of additional, and more stable, folded conformation.

Disulfide bonds are often stated to determine the folded conformations of the protein in which they occur. This cannot be the case, however, for such disulfide bonds would never be generated in a natural protein if the conformation did not favor them. In some cases, proteins have been shown to adopt exactly the same folded conformation in the absence of each of the disulfide bonds (9). Clearly, the disulfide bonds only stabilize the folded conformation and do not determine it.

### Bibliography

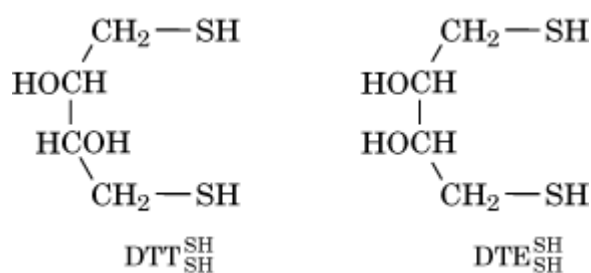
1. J. Wu and J. T. Watson (1997) *Protein Sci.* **6**, 391–398.
2. K. Kuwajima, M. Ikeguchi, T. Sugawara, Y. Hiraoka, and S. Sugai (1990) *Biochemistry* **29**, 8240–8249.
3. J. J. Ewbank and T. E. Creighton (1993) *Biochemistry* **32**, 3677–3693.
4. W. R. Gray (1997) In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, pp. 165–186.
5. T. E. Creighton (1989) In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, pp. 155–167.
6. T. E. Creighton and D. P. Goldenberg (1984) *J. Mol. Biol.* **179**, 497–526.
7. A. Zapun, J. C. A. Bardwell, and T. E. Creighton (1993) *Biochemistry* **32**, 5083–5092.
8. C. Hwang, A. J. Sinskey, and H. F. Lodish (1992) *Science* **257**, 1496–1502.
9. C. P. M. van Mierlo, N. J. Darby, D. Neuhaus, and T. E. Creighton (1991) *J. Mol. Biol.* **222**, 353–371.

### Suggestions for Further Reading

10. J. M. Thornton (1981) Disulphide bridges in globular proteins. *J. Mol. Biol.* **151**, 261–287.
11. T. E. Creighton (1984) Disulfide bond formation in proteins. *Methods Enzymol.* **107**, 305–329.
12. T. E. Creighton (1997) Protein folding coupled to disulphide bond formation. *Biol. Chem.* **378**, 731–744.

### Dithiothreitol, Dithioerythritol

Dithiothreitol (DTT, or  $\text{DTT}^{\text{SH}}$ ) and dithioerythritol (DTE) have closely related structures, each with two [thiol groups](#):



(1)

They have the very useful property of forming a stable, cyclic, intramolecular [disulfide bond](#), abbreviated here as  $\text{DTT}_{\text{S}}^{\text{S}}$  and  $\text{DTT}_{\text{S}}^{\text{S}}$ . Consequently, they are potent reductants of disulfide bonds, eg, RSSR:



For a typical model disulfide bond, like that of [glutathione](#) or [mercaptoethanol](#), the equilibrium constant for this reaction,

$$K_{\text{eq}} = \frac{[\text{DTT}_{\text{S}}^{\text{S}}][\text{RSH}]^2}{[\text{DTT}_{\text{SH}}^{\text{SH}}][\text{RSSR}]} \quad (3)$$

has the value 200 M to 2800 M at high and low pH values, respectively (1). The dependence over the pH range of about 8 to 11 arises because of the difference in  $pK_a$  values of the thiol groups involved.

The first thiol group of  $\text{DTT}_{\text{SH}}^{\text{SH}}$  to ionize has an apparent  $pK_a$  of about 9.2, while that of the second is about 10.2. Even with identical microscopic  $pK_a$  values, the first and second apparent  $pK_a$  values would be expected to differ by 0.6 pH units for statistical reasons, so there is little electrostatic interaction between the two thiol groups, ionization of one increasing the  $pK_a$  of the other by at most 0.4 pH units. The equilibrium constant has the unit of concentration because the disulfide bond of DTT is intramolecular, whereas that of RSSR can be considered to be intermolecular because its reduction generates two RSH molecules. The equilibrium constant can therefore be interpreted in terms of the [effective molarity](#) of the two thiol groups of  $\text{DTT}_{\text{SH}}^{\text{SH}}$  relative to each other. It has a large value of approximately 120 M because the two thiol groups are kept in reasonable proximity and are able to form a disulfide bond in a favorable 6-membered ring structure (2), even though the disulfide bond of  $\text{DTT}_{\text{S}}^{\text{S}}$  appears to be somewhat strained; it has a CSSC dihedral angle of only about  $60^\circ$ , rather than the most favorable  $90^\circ$  (3). Nevertheless, there is only a small change in enthalpy of 0.5 kcal/mol when  $\text{DTT}_{\text{SH}}^{\text{SH}}$  reduces model linear disulfide bonds, as in Eq. 2 (4). The thiol and disulfide forms may be distinguished readily by the UV **absorbance** of the disulfide form at 280 nm (5).

These reagents were introduced by Cleland (5) to overcome the limited potency of thiol reagents such as mercaptoethanol, and they are often called ‘‘Cleland’s reagent.’’ Although he appears to have overestimated their efficacy, measuring an equilibrium constant for Eq. 3 of  $10^4$  M, they are still some of the most effective disulfide reductants readily available.

## Bibliography

1. M.-H. Chau and J. W. Nelson (1991) FEBS Letters **291**, 296–298.
2. N. J. Darby and T. E. Creighton (1993) J. Mol. Biol. **232**, 873–896.
3. S. Capasso and Z. Zagari (1981) Acta Crystallogr. **B37**, 1437–1439.
4. S. Lapanje and J. A. Rupley (1973) Biochemistry **12**, 2370–2372.
5. W. W. Cleland (1964) Biochem. **3**, 480–482.

## Divergent Evolution

Divergence is defined as an evolutionary event in which two morphological or molecular traits arose from a common ancestor, and were initially identical, but became dissimilar during [evolution](#). Divergence is, of course, extremely common and the basis for most evolution; without divergence, all lineages would remain the same.

Divergent evolution is usually apparent at the molecular level from significant similarities between the amino acid sequences of [proteins](#) and the **nucleotide sequences** of DNA or RNA. There are so many such sequences possible that it seems most improbable that two sequences could become similar by chance. Similar sequences are said to be **homologous**, and the differences between them are assumed to be related to the time since they diverged from their last common ancestor. The molecular mechanisms of evolutionary divergence of the genetic material include nucleotide substitution and deletion/insertion; chromosomal [recombination](#), [transposition](#), and **inversion**; [gene duplication](#), and **gene conversion**; [exon shuffling](#), and **domain** shuffling; and [horizontal gene transfer](#). The number of nucleotide substitutions is a simple and useful measure of the degree of divergence between two sequences. In fact, there are a dozen methods available for estimating the number of nucleotide substitutions, using one to six parameters (1-5). Once the number of nucleotide substitutions has been estimated, a [phylogenetic tree](#) can be constructed by using those numbers, to display the pathway that the divergence has followed during evolution (6).

Divergence is contrasted with [convergent evolution](#), which is much more rare. Convergence is an evolutionary event in which two morphological or molecular traits become similar during evolution, even though the ancestor is totally different. In the case of divergence, there must be a certain degree of similarity between two traits to suggest that they had a common ancestor. For convergence, on the other hand, a certain degree of dissimilarity must exist if the two traits originated from independent ancestors. Therefore, the distinction between divergence and convergence can be difficult to judge unless additional evolutionary information is available to indicate whether the similarities or the differences are the more significant. In the case of nucleotide or amino acid sequences, any similarities greater than expected from random are usually taken to indicate divergence. For example, the divergence of nucleotide sequences between human and chicken is shown in Figure 1, where the asterisk (\*) indicates a site where nucleotide differences exist, and where 20 nucleotide sites are different and the remaining sites are the same. The probability that this similarity occurred by mere chance is extremely small. Therefore, it can be reasonably concluded that the difference in nucleotide sequences between these two species is due to divergence. Thus, it follows that divergence from the common ancestor has taken place as a result of the accumulation of nucleotide substitutions.

**Figure 1.** Divergence of nucleotide sequences between humans and chickens.

|                |     |     |     |     |     |     |     |     |     |     |     |     |     |     |     |
|----------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| <b>Human</b>   | Gly | Ile | Val | Glu | Gln | Cys | Cys | Thr | Ser | Ile | Cys | Ser | Leu | Tyr | Gln |
|                | GGC | ATT | GTG | GAA | CAA | TGC | TGT | ACC | AGC | ATC | TGC | TCC | CTC | TAC | CAG |
|                | *   |     | *   | *   |     |     | *   | *** | *   | **  | *   |     |     |     | *   |
| <b>Chicken</b> | GGG | ATT | GTT | GAG | CAA | TGC | TGC | CAT | AAC | ACG | TGT | TCC | CTC | TAC | CAA |
|                | Gly | Ile | Val | Glu | Gln | Cys | Cys | His | Asn | Thr | Cys | Ser | Leu | Tyr | Gln |
| <b>Human</b>   | Leu | Glu | Asn | Tyr | Cys | Asn |     |     |     |     |     |     |     |     |     |
|                | CTG | GAG | AAC | TAC | TGC | AAC |     |     |     |     |     |     |     |     |     |
| <b>Chicken</b> | CTG | GAG | ACC | TAC | TGC | AAC |     |     |     |     |     |     |     |     |     |
|                | Leu | Glu | Asn | Tyr | Cys | Asn |     |     |     |     |     |     |     |     |     |

During the process of molecular evolution, the number of nucleotide differences increases but probably never decreases. Therefore, in most cases, the evolutionary event can be explained by divergence only. Considering the [tertiary structures](#) of proteins, however, in some cases the structures are quite similar, even though their amino acid sequences, as well as their biological function, are very different. For these cases, convergence is a possibility, although it could also be simply a case of extreme divergence masking their original sequence and functional identity.

#### Bibliography

1. T. H. Jukes and C. R. Cantor (1969) in *Mammalian Protein Metabolism*, H. N. Munro, ed., Academic Press, New York, pp. 21–132.
2. M. Kimura and T. Ohta (1972) *J. Mol. Evol.* **2**, 87–90.
3. M. Kimura (1980) *J. Mol. Evol.* **16**, 111–120.
4. N. Takahata and M. Kimura (1981) *Genetics* **98**, 641–657.
5. T. Gojobori, K. Ishii, and M. Nei (1982) *J. Mol. Evol.* **18**, 414–423.
6. M. Nei (1987) *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.

#### DNA Chips

The start of the Human Genome Project in the late 1980s refocused scientists' attention on new, high-throughput technologies for handling and analyzing DNA. Several technologies developed by physicists were evaluated for their potential use in the study of biomolecules. For example, [mass spectrometry](#) techniques and the [scanning tunneling microscope](#) are now used for analyzing and sequencing DNA. Another technology developed by physicists and applied to biology is the microchip. The semiconductor industry manufactured silicon chips with smaller and smaller features, which allows for greater numbers of operations with a chip of a standardized size. This benefit of miniaturization was applied to the biological sciences in the form of multiparallel arrays of DNA fragments on a small chip. Arraying of DNA or RNA samples onto [nitrocellulose](#) or nylon membranes for hybridization studies is a very common analytical procedure in a molecular biology lab (1) (see [Blotting](#)). DNA chips or microarrays are, in principle, miniaturized versions of these

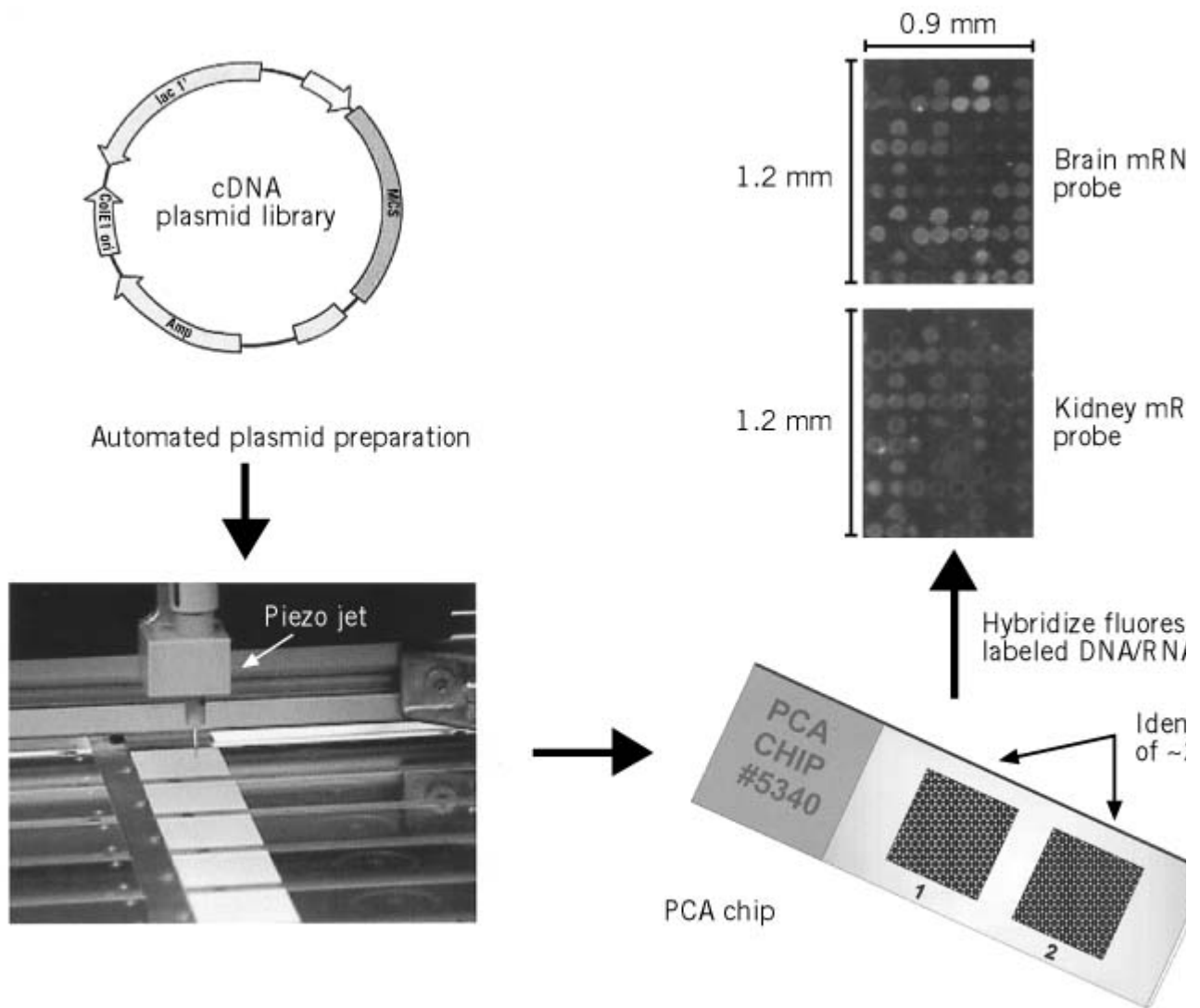
standard dot-blot techniques. With great precision, thousands of DNA fragments, often amplified by **PCR**, or DNA oligonucleotides are spotted/printed onto glass microscope slides. In the case of oligonucleotides, they can also be synthesized directly on coated silicon chips (see [DNA Synthesis](#)). These DNA chips are then used to probe fluorescent-labeled DNA or RNA samples for the presence of specific complementary sequences. DNA chip technology is currently being applied in the areas of monitoring gene expression, polymorphism analysis, gene [mutation](#) analysis, and [DNA sequencing](#).

## 1. Putting Genes on a Chip

The fabrication of high-density arrays on chips for parallel analysis of up to several thousand cDNA clones or PCR fragments on a surface area of approximately 1 cm<sup>2</sup> can be achieved by two techniques: microprinting or piezojet dispensing. The high-density arrays are deposited on the flat solid support, usually a glass slide. The key function of the printing or dispensing process is the extraction of a small volume of sample from a reservoir (eg, a microtiter plate) and the deposition of the same amount of sample fluid in each spot, at equivalent locations on each chip. The center-to-center separation of individual spots in high-density arrays is typically 150 to 500 μm. The deposited sample volume ranges from 50 to 400 pl per spot.

Microprinting devices use an array of pins with either hollow steel capillaries or flat tips. The amount of sample deposited at each spot is defined by the capillary forces acting on the droplet, which adheres to the pin after immersion into the sample reservoir. The droplets at the tips are then transferred to the spots in the array by moving the pins close to the chip surface. After each sample, the tip is cleaned by rinsing. The amount of liquid deposited depends on the size of the tip.

The piezojet is a contactless dispensing technique for subnanoliter droplets. It is based on a cylindrical piezo ceramic that compresses a glass capillary with an inner diameter of ~1 mm on application of an electric pulse. A small droplet is then expelled through the funnel-shaped orifice of the capillary. The outlet diameter of the orifice is ~50 μm. Very uniform droplets can be prepared by this method (see Fig. [1](#)).



The piezojet-dispensing process consists of three steps: (1) Aspiration of an aliquot of sample through the capillary, (2) dispensing of one subnanoliter droplet of sample to equivalent spots on each chip, and (3) discharge of remaining sample volume and cleaning of the piezojet. In order to use the sample material effectively, the aspiration volume of the piezojet should be in the microliter range.

The deposited DNA has to be immobilized on the surface in order to enable the subsequent processing steps. The most commonly used attachment method is to coat the glass slides with polylysine, which holds the negatively charged DNA backbone by [electrostatic interactions](#). Alternative methods, like **streptavidin**-coated slides to bind **biotin**-labeled DNA, and various chemical approaches to bind DNA covalently to pretreated glass, are under development.

The most elegant way to attach DNA fragments covalently onto the glass chips has been developed by Fodor and colleagues at Affymetrix (2-4). This method allows light-directed synthesis of thousands of DNA fragments (oligonucleotides) in precise locations on the microchip. To begin the process, linkers modified with a photochemically removable protecting group are attached to the solid substrate. Light is directed through a photolithographic mask, illuminating specific grid squares on the chip and causing photodeprotection, or the removal of the blocking group, in those squares. The chip is then incubated with a nucleotide harboring a photolabile protecting group at the 5' end. The cycle continues: the chip is exposed to light through the next mask, which activates new grid

sites for reaction with the subsequent photoprotected nucleotide. Using the proper set of masks and chemical steps, it is possible to construct a defined collection of oligonucleotides, generally 20 to 25 bases long, each in a predefined position on the array. A standard 1.28-cm<sup>2</sup> chip can currently be packaged with about 400,000 individual, well-defined oligonucleotides representing many thousand genes.

## 2. Monitoring Gene Expression

A bacterial [genome](#), such as that of *Escherichia coli*, encodes 4288 different genes, the yeast *Saccharomyces cerevisiae* genome about 6000 genes; *Caenorhabditis elegans* has about two times more genes, and the human genome may contain approximately 100,000 unique genes. At least partial sequence information for all these genes will be available soon. Sequence information alone, however, is insufficient for a full understanding of gene function and the control of gene expression. Only a subset of all encoded genes is expressed in any given cell; in higher eukaryotes, this subset is smaller than in bacteria or yeast cells. The levels and the timing of gene expression determine the fate of the cells, their reproduction, differentiation, function, communication, and physiology. Methods such as **Northern blots**, nuclease protection, and RT (reverse transcriptase)-**PCR** are frequently used to measure gene transcripts, but they have the inherent disadvantage of being serial, analyzing a single [messenger RNA](#) at a time. Most comparative techniques, like [subtractive hybridization](#), differential display (5) of amplified mRNA's on gels, and the SAGE (serial analysis of gene expression) method of Velculescu et al. (6), are laborious, often nonreproducible, and not particularly sensitive. Arrays of several thousand genes, the DNA chips, for the first time make it possible to obtain gene expression information on complete genomes (7) quickly, accurately, and efficiently. These gene arrays are in principle reverse Northern blots, where the DNA probes have been immobilized to identify and measure large numbers of mRNA species in parallel. Two prototype DNA microarrays for gene expression monitoring can be described:

### 2.1. cDNA Microarrays

Brown and his colleagues developed a high-capacity system to monitor the expression of many genes in parallel (8, 9). In short, the [complementary DNA](#) of expressed mRNAs from cells is collected, and individual cDNA molecules are isolated and amplified. A microsample of each cDNA from the library is then deposited by high-speed robotic printing on a polylysine-coated glass surface in an array format. Each gene has a unique location in the microarray, which may represent several thousands of genes on a few square centimeter. To compare expression of these genes between two tissues or two cell types, the [poly\(A\)<sup>+</sup>](#)-mRNA in the two samples is copied into cDNA and labeled with two differently colored fluorescent molecules. Both probe samples are applied to a single microarray and allowed to react with the DNA on the microarray. After the appropriate washing steps, each element of the microarray is scanned for fluorescence intensity from the two colors by a laser fluorescence scanner. The observed fluorescence intensity at each array element is proportional to the number of fluorescent cDNA molecules bound to it, and the ratio of the two colors is an accurate measurement of the relative expression level of the genes in the two tissue or cell samples. Twofold changes in expression can be detected (10). As an illustration, Figure 1 shows a similar DNA microarray, the plasmid chip array (PCA), made by piezojet dispensing in the labs of R. Hochstrasser and U. Certa (F. Hoffmann-La Roche Ltd., Basel, Switzerland).

### 2.2. Oligonucleotide Microarrays

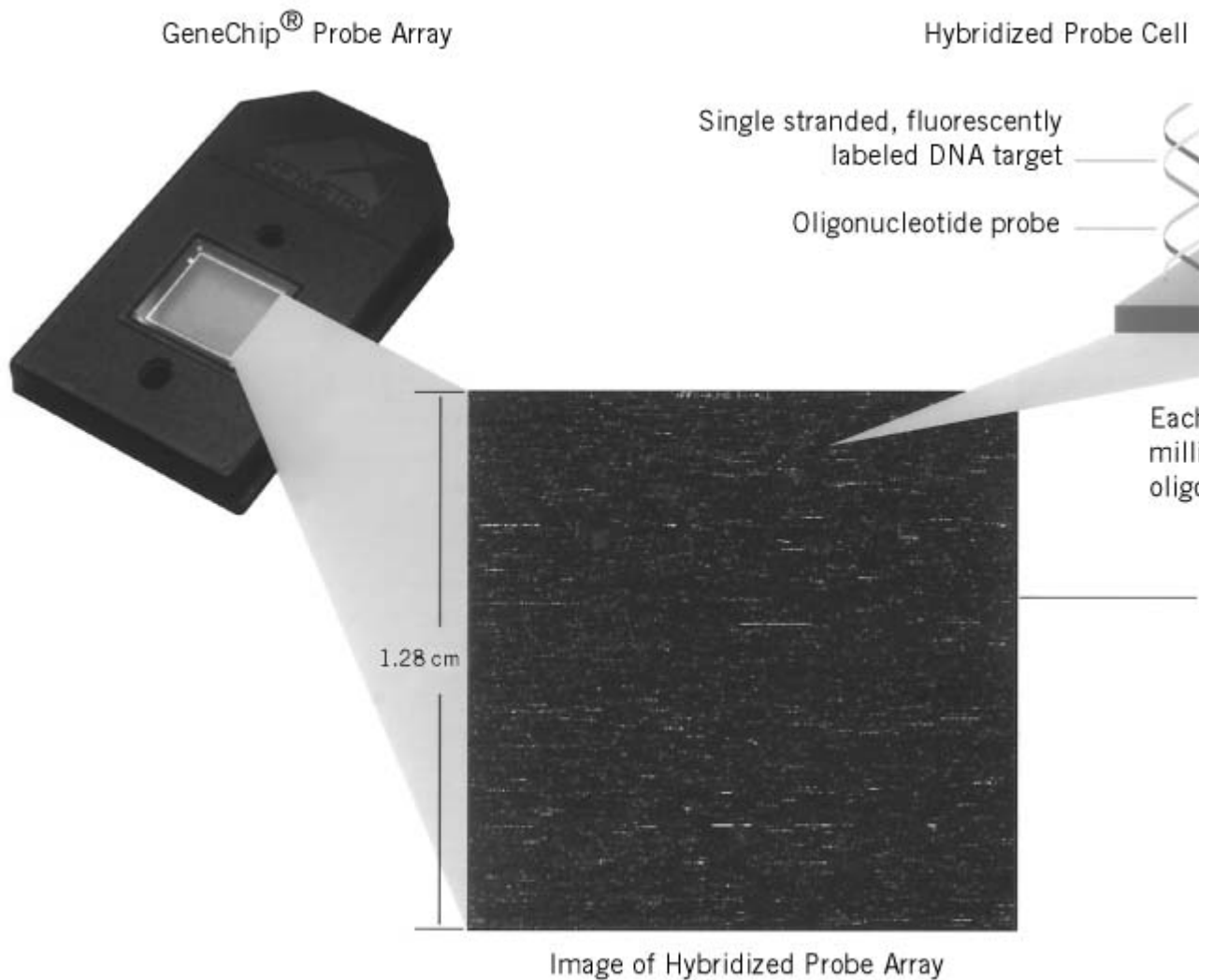
High-density arrays of oligonucleotide probes represent the other prototype method for mRNA expression monitoring. Sequence information about the expressed genes is used directly to select oligonucleotides and to design the photolithographic masks for combinatorial synthesis of the probes directly on derivatized glass, as described previously. The arrays are intentionally redundant, as they contain collections of pairs of probes for each of the RNAs being monitored. Each probe pair consists of a 20- to 25-mer oligonucleotide that is perfectly complementary to a subsequence of a particular message and a mismatch oligo that is identical except for a single base difference in a central position, which serves as an internal control for hybridization specificity. For array

hybridization experiments, target RNA is prepared from cellular mRNA by incorporating labeled ribonucleotides in an *in vitro* [transcription](#) reaction (11), or by direct chemical attachment (12). After hybridization and washing, fluorescence imaging of the arrays is accomplished with a scanning **confocal microscope**. Because oligonucleotide probes for each gene are specifically chosen and synthesized in known locations on the arrays, the hybridization patterns and intensities can be interpreted in terms of gene identity and the relative amount of each transcript. This technology has been successfully applied for monitoring gene expression in mammalian cells (11), yeast (7), and bacteria (12), measuring expression levels of less than one copy of mRNA to several hundred copies per cell in one hybridization experiment.

### 3. Accessing Genetic Diversity

High-density oligonucleotide probe arrays (Fig. 2) can also be applied to a broad range of nucleic acid [sequence analysis](#) problems, including pathogen identification, polymorphism detection and sequence checking. The significant instability of internal probe–target mismatches relative to perfect matches (4, 13) is used to design arrays of probes capable of detecting differences between nucleic acid targets. To interrogate the identity of a nucleotide as position X of a known DNA sequence, four oligonucleotides can be offered with the four different choices A, G, T, C in the middle of 15 or so flanking nucleotides that perfectly match the reference sequence. The probe with the highest intensity after hybridization would indicate the identity of the unknown base. This concept can be extended to detect polymorphisms/mutations relative to a characterized **consensus sequence**. Thus, to screen 1000 nucleotides for polymorphisms/mutations would require 4000 probes. Current applications of this technique include the survey of drug-resistance mutations in **HIV-1** reverse transcriptase and [proteinase](#) genes, a [p53](#) mutations screen, and a [cytochrome P450](#) allelic variants screen.





#### 4. Sequencing by Hybridization

Hybridization can also be used to determine the sequence of unknown DNA, ie, *de novo* [DNA sequencing](#). Sequencing by hybridization is based on the use of oligonucleotide hybridization to determine the set of constituent subsequences in a DNA fragment. This concept which was spearheaded by Crkvenjakov ([14](#), [15](#)) uses a microarray of all possible  $n$ -nucleotide oligomers, eg, 65,635 possible octamers, to identify all the  $n$ -mers present in an unknown DNA sequence. Powerful computational approaches are then used to assemble the complete sequence from the measured hybridization patterns. High-density oligonucleotide probe array technology shows significant promise in enabling *de novo* sequencing in a fast and reliable manner.

#### Bibliography

1. G. G. Lennon and H. Lehrach (1991) Trends Genet. **7**, 314–317.
2. S. P. A. Fodor, J. L. Read, M. C. Pirrung, L. Stryer, A. T. Lu, and D. Solas (1991) Science **251**, 767–773.
3. S. P. A. Fodor, R. P. Rava, X. C. Huang, A. C. Pease, C. P. Holmes, and C. L. Adams (1993) Nature **364**, 555–556.
4. A. C. Pease, D. Solas, E. J. Sullivan, M. T. Cronin, C. P. Holmes, and S. P. A. Fodor (1994) Proc. Natl. Acad. Sci. USA **91**, 5022–5026.

5. P. Liang and A. B. Pardee (1992) *Science* **257**, 967–971.
6. V. E. Velculescu, L. Zhang, B. Vogelstein, and K. W. Kinzler (1995) *Science* **270**, 484–487.
7. L. Wodicka, H. Dong, M. Mittmann, M. H. Ho, and D. J. Lockhart (1997) *Nature Biotech.* **15**, 1359–1367.
8. M. Schena, D. Shalon, R. W. Davis, and P. O. Brown (1995) *Science* **270**, 467–470.
9. D. Shalon, S. J. Smith, and P. O. Brown (1996) *Genome Res.* **6**, 639–645.
10. M. Schena, D. Shalon, R. Heller, A. Chai, P. O. Brown, and R. W. Davis (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10614–10619.
11. D. J. Lockhart, H. Dong, M. C. Byrne, M. T. Follettie, M. V. Gallo, M. S. Chee, M. Mittmann, C. Wang, M. Kobayashi, H. Horton, and E. L. Brown (1996) *Nature Biotech.* **14**, 1675–1680.
12. A. de Saizieu, U. Certa, J. Warrington, C. Gray, W. Keck, and J. Mous (1998) *Nature Biotech.* **16**, 45–48.
13. R. J. Lipshutz, D. Morris, M. Chee, E. Hubbell, N. S. Kozal, N. Shen, R. Young, and S. P. A. Fodor (1995) *BioTechniques* **19**, 442–447.
14. R. Drmanac, I. Labat, I. Brukner, and R. Crkvenjakov (1989) *Genomics* **4**, 114–128.
15. R. Drmanac, S. Drmanac, Z. Strezoska, T. Paunesku, I. Labat, M. Zeremski, J. Snoddy, W. K. Funkhouser, B. Koop, L. Hood, and R. Crkvenjakov (1993) *Science* **260**, 1649–1652.

## DNA Damage, Inducible Responses To

Agents that damage **DNA**, like other external stimuli, elicit a complex series of cellular responses in both prokaryotes and eukaryotes. The molecular mechanisms and the consequences of induction of these genotoxic responses in prokaryotes are well-understood. In eukaryotes, in contrast, the mechanisms and the effects of the multiple response reactions remain ill-defined. In *Escherichia coli* there are genotoxic response reactions that are induced mainly, but not exclusively, by bulky DNA lesions ([SOS response](#)), by alkylating agents (adaptation), and by oxidative stress (adaptive response to oxidative stress). These response reactions result in **transcriptional** activation of genes that inactivate the noxious agent, repair the damage caused by the agent (see [DNA Repair](#)), or help the cell survive the DNA damage.

### 1. SOS Response

The SOS response is a coordinated cellular response in *E. coli* and aids in the survival of the organism by affecting the expression of proteins that are involved in cellular division, **replication**, [recombination](#), and [excision repair](#) (1) (see [SOS Response](#)). About 30 genes are involved. The basic mechanism of regulation is relatively simple (2, 3). Under normal conditions, the [LexA repressor](#) binds to the SOS box, with the sequence 5'-CTG-N<sub>10</sub>-CAG-3', to maintain repression of the SOS regulon. Directly or indirectly, DNA damage or other replication blocks generate single-stranded regions that lead to the formation of single-stranded DNA, on which the RecA protein polymerizes (4, 5). The RecA-single-stranded DNA filament binds to LexA and stimulates its autoproteolysis activity (6). Cleavage of LexA inactivates its repressor activity (2-5), allowing transcription of the SOS genes, including *lexA* and *recA* themselves. Cell division stops, excision (7, 8) and recombination (2, 3) activities increase, and the capacity of replicase to bypass the damage increases (9, 10). All of these enable the cell either to eliminate the DNA damage or to survive with damaged DNA. Upon recovery, the inducing signal disappears, and LexA accumulates and represses the cognate genes. The excision repair genes *uvrA*, *uvrB*, and *uvrD* (7, 8) are under the control of the

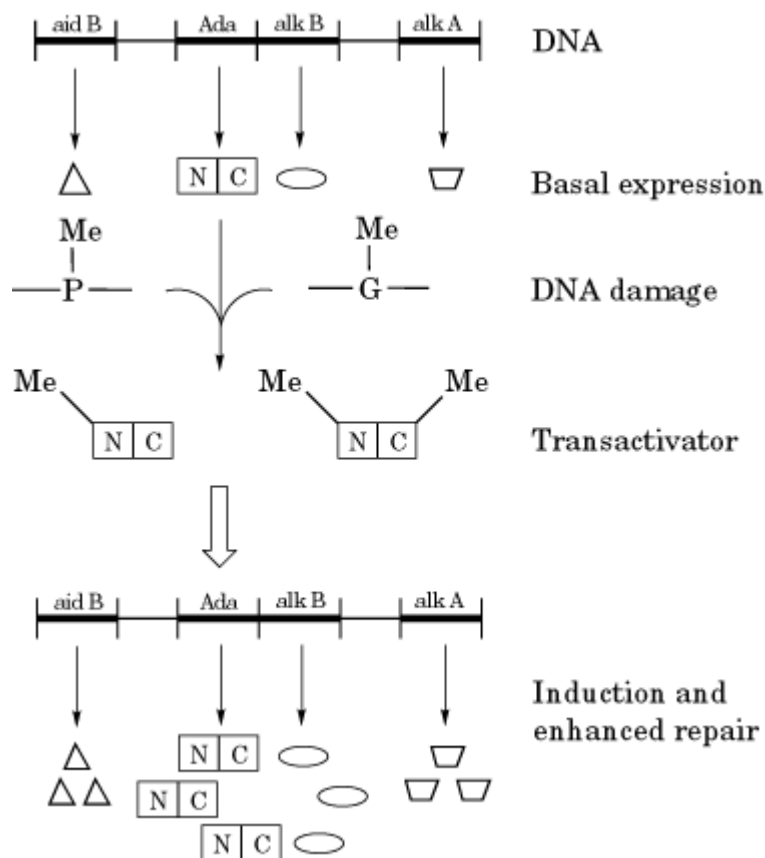
LexA protein; consequently, the cellular excision repair activity increases transiently as part of the SOS response reaction and then drops down to preinduction levels after repair.

## 2. Adaptive Response

The adaptive response is the phenomenon of increased resistance of cells to killing and mutation by alkylating agents that results from prior exposure to lower doses of the same agent (11). The phenomenon was originally discovered in *E. coli* and was later found in other prokaryotes. There is no evidence for an adaptive response in eukaryotes or archaeobacteria, although alkylating agents elicit a general stress response reaction in eukaryotes.

The molecular mechanism of the adaptive response in *E. coli* is relatively simple and known in considerable detail (12, 13) (Fig. 1). Treatment of cells with alkylating agents such as *N*-methyl-*N'*-nitro-*N*-nitrosoguanidine (MNNG) methylates multiple sites on DNA, including the backbone phosphate groups. The methylphosphotriesters that are formed by such treatments are mostly nonmutagenic and nonlethal. In contrast, the product *O*<sup>6</sup>-methylguanine is a miscoding base that induces G to A transition mutations, and 3-methyladenine is a lethal lesion. However, it is the methylphosphotriesters that constitute the signal substrate for the repair of the mutagenic and lethal alkylation products.

**Figure 1.** Adaptive response in *E. coli*. The Ada protein transfers the methyl group from methylphosphotriesters to Cys69 of the protein. This covalent modification activates Ada as a transcription factor, which induces its own gene for repair of *O*<sup>6</sup>-methylguanine, *alkA* for repair of 3-methyladenine, and *aidB* and *alkB* for protecting DNA.



The adaptive response is a positive regulatory response reaction. The central regulator in this response reaction is the 39-kDa Ada protein (*O*<sup>6</sup>-methylguanine DNA methyltransferase). The

structure of the Ada B protein in *E. coli* exhibits a **helix–turn–helix** variant **DNA-binding** motif, but this motif is separate from its active site (14). Ada acts both as a [transcription factor](#) and as a [methyltransferase](#). Upon formation of methylphosphotriesters, the Ada protein transfers the methyl group from the phosphate to a [cysteine](#) residue (Cys69) in the amino-terminal domain of the molecule. This reaction induces a conformational change in the protein, which increases its affinity for the Ada box, with the sequence A<sub>3</sub>–N<sub>3</sub>–A<sub>3</sub>–GCGCA (12, 13), upstream of cognate **promoters**.

Binding of Ada to these promoters activates transcription of target genes that are involved in preventing and repairing the damage caused by alkylating agents. The following genes are induced by Ada: the *ada* gene itself, *alkA*, *alkB* and *aidB*. The roles of these genes in cellular defense against alkylation killing and mutagenesis are summarized below.

### 2.1. Ada

The Ada protein performs three functions. It transfers the methyl group from O<sup>6</sup>-methyl-guanine (or O<sup>4</sup>-methylthymine) to Cys321 of Ada (15), it transfers the methyl group from the methylphosphotriester to Cys69, and it binds upstream of the promoters of cognate genes and acts as a positive transcriptional regulator. The order of the methyl transfer reactions is random, and they are independent of one another. In contrast, binding of Ada to promoter regions is strongly influenced by methylation of residue Cys69. Even though unmethylated Ada binds to certain promoters (such as the *ada* promoter) with low affinity, the binding is strongly stimulated by methylation of Cys69, and the binding to other promoters is completely dependent on the methylation of Cys-69. The methylation status of Cys321 has no effect on the regulatory activity of the protein. In addition to Ada, *E. coli* and some other prokaryotes contain a second alkyl transferase with no regulatory function (16).

### 2.2. AlkA

The AlkA protein removes 3-methyladenine and other alkylated purines by a [DNA glycosylase](#) activity. There are two 3-methyladenine DNA glycosylases (Tag). One is encoded by the *tag* gene (TagI), which is not regulated by Ada and has narrow substrate specificity. The other, which is regulated by Ada and encoded by *alkA* (TagII), removes 3-methyladenine and other methylated purines, as well as O<sup>4</sup>-methyl thymine, and thus has a wide substrate range.

### 2.3. AlkB

The AlkB protein is a 27-kDa monomer that enables cells to survive DNA damage induced by SN2 alkylating agents, but not by SN1. Although it is known that the protein acts at the level of DNA, the precise mechanism of the repair reaction is not known.

### 2.4. AidB

The AidB protein is a 60-kDa monomer with a high degree of homology to several mammalian [acetyl coenzyme A](#) dehydrogenases, and it exhibits isovaleryl coenzyme A dehydrogenase activity. Metabolic activation by [glutathione](#) or other nonprotein [thiol groups](#) is necessary for the toxic effect of MNNG. It is possible that AidB confers resistance by affecting the intracellular thiol pool, or by acting directly on MNNG and detoxifying it. The situation is somewhat similar to the induction of *zwf* as part of the oxidative stress response reaction. This gene encodes glucose-6-phosphate dehydrogenase; upon induction of the SoxRS regulon by superoxides, the level of the dehydrogenase increases, and the enzyme in turn increases the cellular NADPH pool, which detoxifies reactive oxygen species.

### 2.5. Mechanism of the Adaptive Response

The homeostasis of the adaptive response is maintained as follows. Upon DNA damage by MNNG, the Ada protein, which is present at about 100 copies per cell, reacts with methylphosphotriesters and transfers the methyl group to Cys69. As a result, Ada becomes a potent transcription factor, binds upstream of the *ada-alkB*, *aidB* and *alkA* operons, and turns on these genes. Upon completion of DNA repair and elimination of the alkylating agent, the methylphosphotriester inducing signal disappears. The accumulated Ada is no longer methylated; because of its high concentration, it binds

to but cannot induce the cognate genes, interfering with binding of alkylated Ada. Eventually, the alkylated Ada is degraded or diluted out by cell growth and division, and the adaptation reaction is turned off.

Human  $O^6$ -methylguanine DNA methyltransferase has been characterized in some detail (17). It has no regulatory function, but it has an unusual relation to cellular transformation. Resistance to alkylating agents has been observed in some tumor cell lines that lack  $O^6$ -methyltransferase. The resistance appears to be due to simultaneous defects in  $O^6$ -methyltransferase and in mismatch repair. An accumulation of unrepaired  $O^6$ -methylguanine lesions should lead to [cell death](#), because of a phenomenon called futile cycling, where mismatch repair systems attempt, but fail, to repair the lesion, halting replication and finally inducing [apoptosis](#). Therefore, in tumor cell lines that already lack  $O^6$ -methyltransferase, an additional defect in mismatch repair results in the loss of futile cycling and allows the cells to survive treatment with alkylating agents.

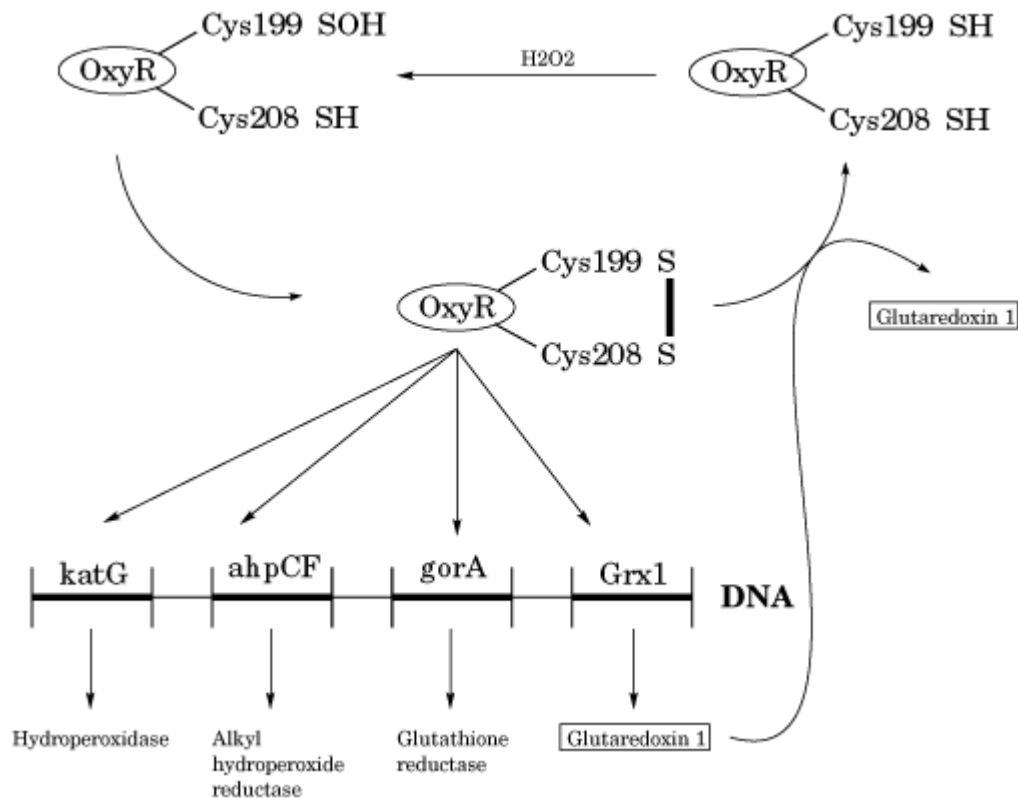
### 3. Adaptive Response to Oxidative Damage

Oxidative stress in the form of reactive oxygen species (ROS: OH,  $H_2O_2$ ,  $O_2^-$ ) induces the transcription of a variety of genes in both prokaryotes and eukaryotes. The molecular mechanisms of induction and the significance of the induction reaction for cellular survival have been elucidated in *E. coli*. In this organism there are two regulatory pathways. One is induced by  $H_2O_2$  and is referred to as the OxyR regulon. The other is induced by superoxides and is referred to as the SoxRS regulon.

#### 3.1. OxyR Regulon

The OxyR system belongs to the LysR family of single-component bacterial regulators, in which the same protein functions both as a sensor and as a signal transducer-regulator (18-20). The OxyR protein acts as both a sensor and transducer of the oxidative stress signal (Fig. 2). The activity of OxyR is regulated by the oxidation status of two cysteine residues, Cys199 and Cys208. Peroxides lead to formation of a disulfide bridge between the two cysteine residues (21). Both disulfide and dithiol forms of OxyR bind to the *oxyR* promoter, but only the disulfide form binds to the *katG* (hydroperoxidase) and *ahpC* (hydroperoxide reductase) promoters. Binding of reduced OxyR to its own promoter represses its gene, whereas the disulfide form turns on this gene, as it does to the *katG*, *ahpC*, *gorA*, and *grxI* genes. The members of this regulon are involved solely in antioxidant defense and not DNA repair.

**Figure 2.** Adaptive response to peroxides. Peroxides convert two cysteine thiol groups on OxyR to a disulfide bond. This activates OxyR as a transcription factor, which turns on genes whose products combat peroxides. In the absence of peroxides, [glutaredoxin](#) reduces OxyR and converts it to the inactive form.



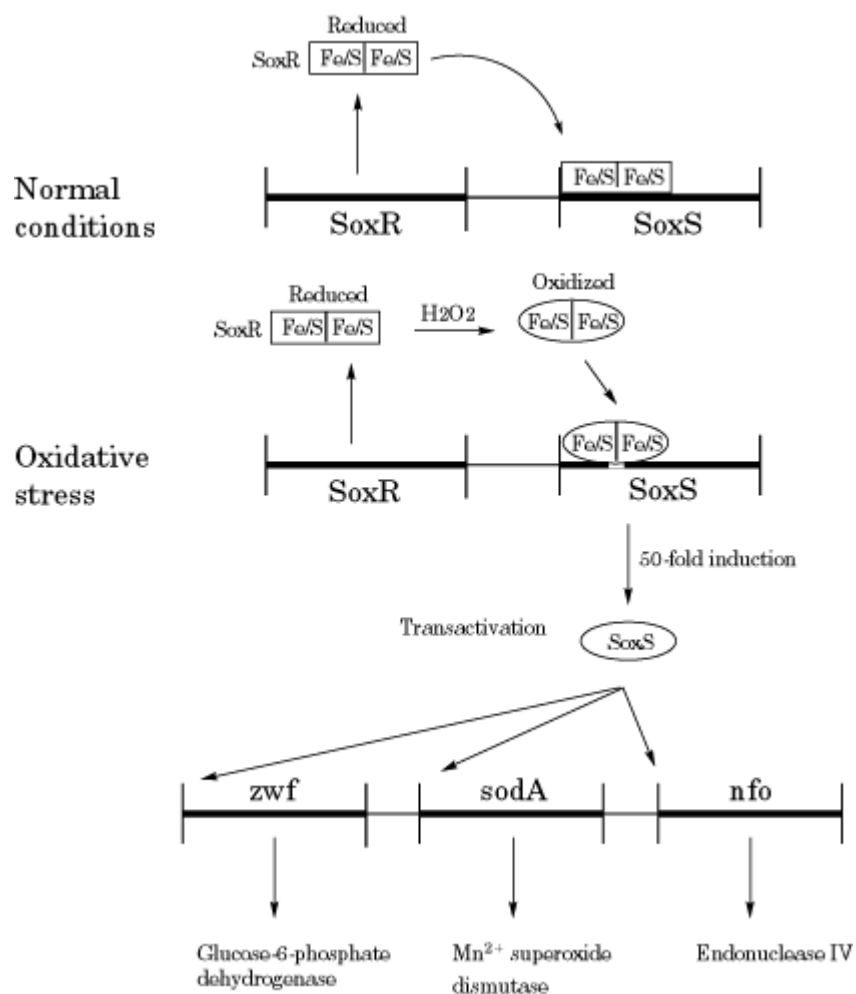
An extremely important fact regarding OxyR-mediated antioxidant response is the finding in *Mycobacterium tuberculosis* that the *oxyR* gene contains deletions and point mutations, resulting in frameshifts that make OxyR nonfunctional (22). This explains the unusual sensitivity of *M. tuberculosis* to the antituberculosis drug isoniazide. As a result of this mutation, alkylhydroperoxide reductase (encoded by *ahpC*) is not induced upon oxidative stress. Isoniazide is activated by KatG catalase-peroxidase; the reactive oxygen species that are side products of this reaction attack DNA and membrane lipids, leading to cell death. These reactive oxygen species are inactivated by alkylhydroperoxide reductase—hence the inability of *M. tuberculosis* to turn on *ahpC* in response to the oxidative burst caused by metabolizing isoniazide leads to its death by oxidative stress (22).

### 3.2. SoxRS Regulon

The SoxRS regulon is a two-component regulatory system that responds to superoxide stress (Fig. 3). SoxR is an [iron-sulfur protein](#) that is sensitive to superoxides, and when oxidized it becomes a transcriptional activator for SoxS. SoxS is a positive regulator of genes encoding superoxide-detoxifying and DNA-repair enzymes. The ferredoxin-like protein SoxR is a 34-kDa homodimer containing one essential [2Fe-2 S] cluster per polypeptide chain (23, 24). Both the reduced and oxidized forms of SoxR bind to the SoxS promoter, with comparable affinities, although binding of the reduced form does not induce transcription. Increased superoxides in the cell convert OxyR to the oxidized form and activate it as a transcription factor. Thus, the activation occurs not by increasing the affinity of SoxR for DNA, but by causing a conformational change in the protein that facilitates the binding of RNA polymerase to the *soxS* promoter and initiation of *soxS* transcription. While the *soxS* promoter is bound by the reduced form of SoxR, the distance between the two signature sequences is shorter than optimal for **RNA polymerase** binding. The conformational change induced in SoxR by oxidation results in overwinding of the promoter, which aligns the -35 and -10 sequences in register for optimal binding by RNA polymerase. The net result is about a 50-fold increase of *soxS* transcription. The SoxS protein then activates transcription of other genes that detoxify superoxides or repair DNA damage caused by superoxides. The following enzymes/proteins are up-regulated by SoxS: SodA (Mn<sup>2+</sup>-superoxide dismutase), Nfo (endonuclease IV, which is a

type-II [AP endonuclease](#)), MicI (which inhibits the synthesis of a [porin](#), thereby blocking the uptake of  $O_2^-$  generators), aconitase,  $O_2$ -insensitive fumarase (which replaces  $O_2^-$  damaged enzymes), flavodoxin, NADP<sup>+</sup> oxidoreductase (Fpr, which reactivates damaged Fe-S centers), and Zwf (glucose-6-phosphate dehydrogenase) which generates the NADPH that is used by Fpr and also perhaps directly reacts with superoxides and eliminates the offending agents).

**Figure 3.** Adaptive response to superoxides. The SoxR protein has a [2Fe-2 S] active center. The protein binds to the *soxS* promoter in reduced form, but does not turn it on. Upon oxidation by superoxides, it overwinds the promoter of *soxS*, making it a high-affinity site for RNA polymerase. Increased transcription of *soxS* leads to increased SoxS protein, which is a transcriptional activator that turns on genes involved in inactivating superoxides (*zwf*, *sodA*) or repairing DNA (*nfo*).



Induction of the members of this regulon enables the cell to reduce the uptake of superoxide generators, inactivates the superoxides, and repairs the DNA damage caused by oxidative stress. Upon disappearance of the inducing signal, the reduced state of the intracellular milieu is restored, which leads to conversion of the Sox from the oxidized to the reduced forms and the turning off of the oxidant stress response reaction.

In mammalian cells, oxidative stress induces many genes and regulates several transcription factors, including AP1 (25) and p53 (26, 27). However, the physiological relevance of these effects is not clear at present.

## Bibliography

1. C. J. Kenyon and G. C. Walker (1980) *Proc. Natl. Acad. Sci. USA* **77**, 2819–2823.
2. J. W. Little, S. H. Edmiston, Z. Pacelli, and D. W. Mount (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3225–3229.
3. R. Brent and M. Ptashne (1981) *Proc. Natl. Acad. Sci. USA* **78**, 4204–4208.
4. A. J. Clark and A. D. Margulies (1965) *Proc. Natl. Acad. Sci. USA* **53**, 451–459.
5. N. L. Craig and J. W. Roberts (1980) *Nature* **283**, 26–30.
6. J. W. Little (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1375–1379.
7. A. Sancar, G. B. Sancar, W. D. Rupp, J. W. Little, and D. W. Mount (1982) *Nature* **298**, 96–98.
8. G. B. Sancar, A. Sancar, J. W. Little, and W. D. Rupp (1982) *Cell* **28**, 523–530.
9. S. J. Elledge and G. C. Walker (1983) *J. Mol. Biol.* **164**, 175–192.
10. M. Rajagopalan, C. Lu, R. Woodgate, M. O'Donnell, M. F. Goodman, and H. Echols (1992) *Proc. Natl. Acad. Sci. USA* **89**, 10777–10781.
11. L. Samson and J. Cairns (1977) *Nature* **267**, 281–283.
12. Y. Nakabeppu and M. Sekiguchi (1986) *Proc. Natl. Acad. Sci. USA* **83**, 6297–6301.
13. I. Teo, B. Sedgwick, M. W. Kilpatrick, T. V. McCarthy, and T. Lindahl (1986) *Cell* **45**, 315–324.
14. M. H. Moore, J. M. Gulbis, E. J. Dodson, B. Demple, and P. C. E. Moody (1994) *EMBO J.* **13**, 1495–1501.
15. P. Karran, T. Lindahl, and B. Griffin (1979) *Nature* **280**, 76–77.
16. P. M. Potter, M. C. Wilkinson, J. Fitton, F. J. Carr, J. Brennand, and D. P. Cooper (1987) *Nucleic Acids Res.* **15**, 9177–9193.
17. K. Tano, S. Shiota, J. Collier, R. S. Foote, and S. Mitra (1990) *Proc. Natl. Acad. Sci. USA* **87**, 686–690.
18. M. F. Christman, R. W. Morgan, F. S. Jacobson, and B. N. Ames (1985) *Cell* **41**, 753–762.
19. G. Storz, L. A. Tartaglia, and B. N. Ames (1990) *Science* **248**, 189–194.
20. M. B. Toledano, I. Kullik, F. Trinh, P. T. Baird, T. D. Schneider, and G. Storz (1994) *Cell* **78**, 897–909.
21. M. Zheng, F. Aslund, and G. Storz (1998) *Science* **279**, 1718–1721.
22. V. Deretic, E. Pagan-Ramos, Y. Zhang, S. Dhandayuthapani, and L. E. Via (1996) *Nature Biotech.* **14**, 1557–1561.
23. P. Gaudu and B. Weiss (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10094–10098.
24. H. Ding, E. Hidalgo, and B. Demple (1996) *J. Biol. Chem.* **271**, 33173–33175.
25. S. Xanthoudakis, G. Miao, F. Wang, Y. C. E. Pan, and T. Curran (1992) *EMBO J.* **11**, 3323–3335.
26. D. Parks, R. Bolinger, and K. Mann (1997) *Nucleic Acids Res.* **25**, 1289–1295.
27. L. Jayaraman, K. G. K. Murthy, C. Zhu, T. Curran, S. Xanthoudakis, and C. Prives (1997) *Genes and Dev.* **11**, 558–570.

## Suggestions for Further Reading

28. E. M. Witkin (1976) Ultraviolet mutagenesis and inducible DNA repair in *Escherichia coli*. *Bacteriol. Rev.* **40**, 869–907.
29. J. W. Little and D. W. Mount (1982) The SOS regulatory system of *Escherichia coli*. *Cell* **29**, 11–22.
30. G. C. Walker (1984) Mutagenesis and inducible responses to DNA damage in *Escherichia coli*. *Microbiol. Rev.* **48**, 60–93.
31. B. Demple (1991) Regulation of bacterial oxidative stress genes. *Annu. Rev. Genet.* **25**, 315–



32. T. Lindahl, B. Sedgwick, M. Sekiguchi, and Y. Nakabeppu (1988) Regulation and expression of the adaptive response to alkylating agents. *Annu. Rev. Biochem.* **57**, 133–157.
33. S. Mitra and B. Kaina (1993) Regulation of repair of alkylation damage in mammalian genomes. *Progr. Nucleic Acids Res. Mol. Biol.* **44**, 109–142.

## DNA Degradation *In Vivo*

Low levels of DNA degradation occur continuously in all cells, from bacteria to humans, in conjunction with various [DNA repair](#) processes. Probably 90% of this degradation can be attributed to [excision repair](#), during which only a few **nucleotides** are released from each damaged site by the actions of repair **endo** and **exonucleases** and recycled in the overall process. Within a population of cells, high levels of DNA degradation occur in individual cells that need elimination from the population. This need occurs either because their DNA has become too heavily damaged to be successfully repaired, because the cells have been infected with bacteriophages or **viruses** or, during development and tissue remodeling in multicellular organisms, because these cells have received signals instructing them to commit suicide. In higher **eukaryotes**, this cell suicide is known as [apoptosis](#). In this case, the [chromatin](#) DNA is fragmented and packaged with other cellular materials into apoptotic bodies that are engulfed by adjacent healthy cells in the tissue and digested in their **lysosomes**. The digestion products are recycled in these cells.

Most chromatin DNA fragmentation in eukaryotic cells follows an ordered path that yields first 300-kbp double-strand (ds) fragments, then 50-kbp ds fragments and lastly a range of small ds fragments that are multiples of [nucleosome](#)-sized DNA, 180–200 bp (1). This fragmentation pattern reflects the different levels of packing of the DNA in the chromatin. Six 50-kbp loops form 300-kbp rosettes that are stacked together. Cleavage between rosettes releases the 300-kbp fragments, and cleavage at the bases of the loops releases 50-kbp linear ds DNA that is readily cleaved between nucleosomes to small fragments that have 3'-OH and 5'-P termini. When separated by [gel electrophoresis](#) in [agarose](#), the latter appear as a characteristic “ladder” of DNA. All stages of the fragmentation are  $Mg^{2+}$ -dependent, but the subsequent cleavage to small fragments is also activated by  $Ca^{2+}$ . This may indicate that there is more than one **nuclease** involved in chromatin DNA cleavage. In a few cases, the earlier stages of fragmentation are sufficient for apoptosis to proceed to completion. Random DNA degradation occurs during cell [necrosis](#) as a result of very heavy damage to the cell or tissue injury and is accompanied by lysosome disruption and cytolysis.

Except for a small body of work on *Escherichia coli*, little is known about cell suicide in bacteria. When DNA bacteriophages that have not adapted to growth in *E. coli* are cleaved by host [restriction enzyme](#) nucleases, the resulting ds DNA fragments are degraded to small oligonucleotides by a nuclease with  $Mg^{2+}$ - and ATP-dependent endo- and exonuclease activities encoded by the *recB*, *recC*, and *recD* genes, the *recBCD* nuclease. This nuclease normally acts in [recombination](#) and in recombinational ds break repair, a minor DNA repair pathway. Some bacteriophages that have adapted to *E. coli* contain genes that encode specific inhibitors of the *recBCD* nuclease. This nuclease is also responsible for degrading genomic DNA damaged beyond repair by UV light or ionizing radiation (2), but the details of the process are not known. The enzyme has **homologues** in many other species of bacteria.

The identity of the nucleases that are responsible for chromatin DNA degradation during apoptosis in

higher eukaryotes remains uncertain. Two difficulties encountered are that **nuclei** contain very low levels of active nucleases and that they are readily contaminated during isolation by nucleases from other **organelles**. The potent **mitochondrial** nuclease has been shown to contaminate nuclei isolated from calf thymus (3), a tissue that exhibits a high rate of apoptosis, although this nuclease is not known to act in chromatin DNA degradation. DNase II, an acid **deoxyribonuclease** (DNase) normally found in the lysosomes, generates a ladder of ds DNA fragments when incubated with isolated nuclei at acidic pH. A small decrease in intracellular pH (0.3 pH units) occurs during apoptosis (4), but this would not be enough to result in appreciable activation of DNase II, even if it had a bona fide nuclear location. DNase II is not metal ion-dependent and makes ds breaks with 3'-P and 5'-OH groups, termini not found on the apoptotic ds DNA fragments. However, the corresponding acid DNase of the flatworm, *Caenorhabditis elegans*, does play a role in apoptosis, namely, in digesting the apoptotic bodies when engulfed by healthy cells. A loss of function of the nuclease caused by mutation of the *nuc1* gene led to accumulation of the apoptotic bodies in the lysosomes (5).

Considerable attention has been paid to **DNase I** as a possible candidate for chromatin DNA fragmentation (6). It is  $Mg^{2+}$ -dependent and  $Ca^{2+}$ -activated and makes the appropriate ds breaks in DNA. When expressed, DNase I occurs in an inactive complex with **actin** that is activated *in vitro* by proteolysis. However, the activity is not expressed in many cell types that undergo apoptosis and, when DNase I is added to isolated nuclei, a smear of randomly cleaved DNA results, rather than a ladder of DNA fragments. Overexpression of the bovine DNase I gene in monkey COS cells induced chromatin DNA ladder formation. However, expression of any nuclease activity in the nuclei may "damage" (nick or cleave) DNA sufficiently to trigger activation of the endogenous apoptotic program.

A bewildering array of other  $Ca^{2+}$ ,  $Mg^{2+}$ -endonucleases, which vary in size from 18 to 97 kDa and generate ladders of DNA in isolated nuclei, have also been proposed to act in chromatin DNA degradation (7). These have been isolated from nuclei of apoptotic cells, but none have been fully characterized, no inactive forms have been reported, and little is known about their possible relationships. The smallest of these is NUC-18 (18 kDa), which has been identified as a **cyclophilin** (8). Cyclophilins have **peptidyl prolyl cis-trans isomerase** activity and play a role in **protein folding**. Since the specific nuclease activity associated with purified cyclophilins is very low, this activity may be associated with a contaminant. The largest (97 kDa) species cross-reacts with **antibody** raised against purified endo-exonuclease from *Neurospora crassa* (7), a  $Mg^{2+}$ -dependent enzyme with homologues in other **fungi** and **yeast** that has been shown to act in recombination and in recombinational ds break repair (3). It may be the eukaryotic counterpart of the bacterial recBCD nuclease. It occurs in both active and inactive forms and has endonuclease activity with both DNA and RNA and exonuclease activity with DNA. A mammalian endo-exonuclease, synergistically activated by  $Ca^{2+}$ , has been isolated from monkey CV-1 cells (9). It is the major degradative nuclease in nuclei of human leukemia cells, present entirely in inactive form. Proteolysis of endo-exonuclease has been detected by immunoblotting in response to different apoptotic agents and yielded polypeptides identical in sizes to the various  $Ca^{2+}$ ,  $Mg^{2+}$ -endonucleases isolated from apoptotic cells by others. Proteolysis, an essential feature of apoptosis, may result in the activation and turnover of endo-exonuclease. Direct activation of endo-exonuclease, circumventing the apoptotic signaling pathways, could lead to new therapies in eliminating unwanted cells.

## Bibliography

1. P. R. Walker, S. Pandey, and M. Sikorska (1995) Cell Death Differ. **2**, 93–100.
2. G. R. Smith (1988) Microbiol. Rev. **52**, 1–28.
3. M. J. Fraser and R. L. Low (1993) "Fungal and mitochondrial nucleases", in *Nucleases*, 2nd ed. R. J. Roberts, S. M. Linn, and S. Lloyd, eds., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 171–207.

4. M. A. Barry and A. Eastman (1993) Arch. Biochem. Biophys. **300**, 440–450.
5. J. Hevelone and P. S. Hartman (1988) Biochem. Genet. **26**, 447–460.
6. B. Polzar et al. (1993) Eur. J. Cell Biol. **62**, 397–405.
7. M. J. Fraser et al. (1996) J. Cell Sci. **109**, 2343–2360.
8. J. W. Montague, M. L. Gaido, C. Frye, and J. A. Cidlowski (1994) J. Biol. Chem. **269**, 18877–18880.
9. C. Couture and T. Y.-K. Chow (1992) Nucl. Acids Res. **20**, 1379–1385.

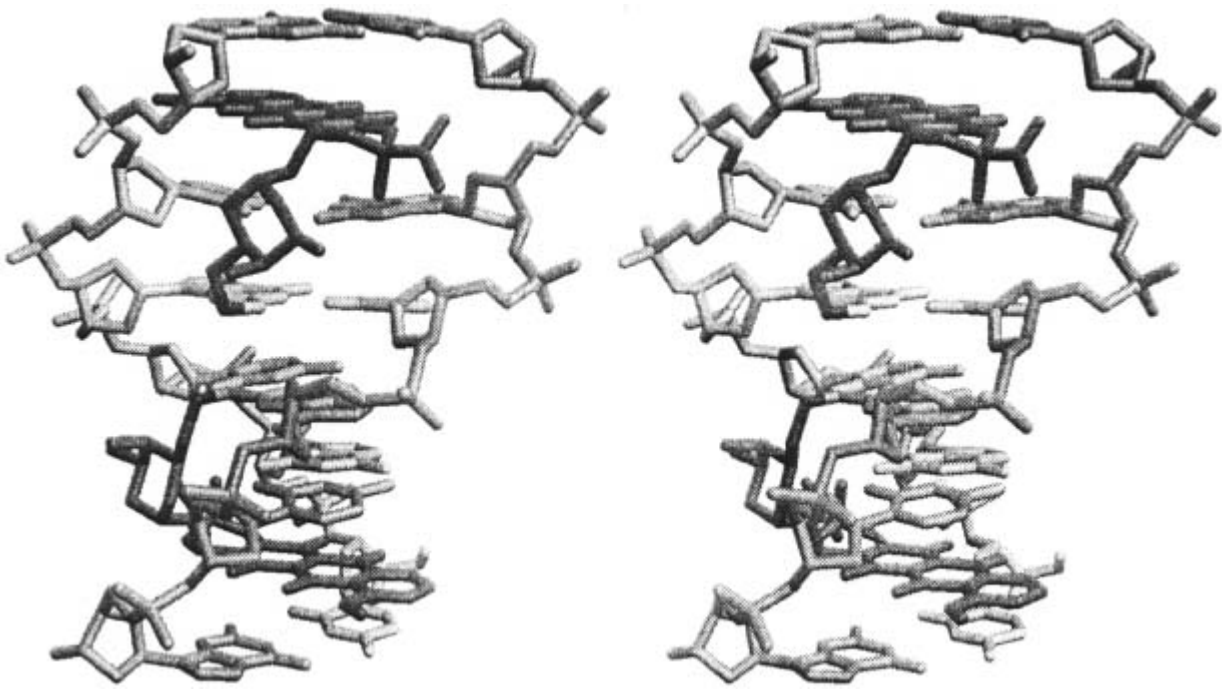
### Suggestion for Further Reading

10. M. J. Fraser (1996) *Endo-exonucleases: Actions in the Life and Death of Cells*, Bioscience Report, R. G. Landes Co., Austin, Texas, pp. 1–73.

## DNA Dynamics

DNA possesses many properties associated with a linear [polymer](#), including various mechanical (rigidity) and motional properties. The normal [B-DNA](#) is considerably more flexible than both [A-DNA](#) and [Z-DNA](#). [DNA structure](#) in solution has two classes of motional properties: large-scale and local motions. The large scale tumbling motion of a DNA duplex of the size in the range of 10–100 base pairs has a correlation time of nanoseconds to microseconds. It is related to rigidity of the duplex (of varying composition and sequence) as measured by the persistence length. This global flexibility of DNA can be used to understand the behavior of DNA in solution and in gels. There are also local motions that are important in the biological functions of DNA. The dynamic properties of DNA allow it to stretch, bend, twist, and so on. The binding of many intercalators to DNA requires that adjacent DNA [base pairs](#) be separated to 6.8 Å in order to accommodate the aromatic intercalator. This is a large conformational change (including sugar pucker rearrangements), involving many [torsion angles](#). Since the binding of intercalators is a reversible process, DNA must be in constant motions, so that the on-off binding process can happen at any time. Figure 1 shows the high resolution structure of the complex between the anticancer drug daunorubicin and CGCGCG **crosslinked** by formaldehyde (1). Finally, the motion of DNA or drug–DNA complexes may be analyzed by high resolution [X-ray crystallography](#) structure refinement that incorporates refinement of anisotropic [temperature factors](#). This demonstrates that various parts of the complex have different motional properties, as is evident in their thermal ellipsoids.

**Figure 1.** The three-dimensional structure of the daunorubicin–CGCGCG complex (Protein Database DDF023).



Some dinucleotide steps, such as CpG, TpA, or TpG (= CpA), appear to be more flexible (see [Base Pairs](#)). In fact, those steps are the preferred binding sites for intercalator binding. They are also the preferred kinked sites induced by binding of [proteins](#). The local dynamic property associated with the CpG site (and other sites) has been probed by high resolution and solid-state nuclear magnetic resonance ([NMR](#)) spectroscopy (2). The line shape analysis of the  $^2\text{H}$ -NMR spectra of the dodecamer CGCGAATTCGCG, with the  $^2\text{H}$  probe incorporated at different sites of the DNA, suggested that the deoxyribose of the C9 residue has a greater motional amplitude than do the other nucleotides. The atypical motional property of the CpG step has been attributed to be correlated with the cleavage at this site by the [restriction enzyme](#) EcoRI. Therefore, the local dynamics of DNA clearly play important roles in DNA bending induced by proteins, in binding of intercalators, in restriction enzyme action, and undoubtedly in other biological functions.

#### Bibliography

1. A. H.-J. Wang, Y.-G. Gao, Y.-C. Liaw, and Y.-K. Li (1991) *Biochemistry* **30**, 3812–3815.
2. B. H. Robinson, C. Mailer, and G. Drobny (1997) *Annu. Rev. Biophys. Biomol. Struct.* **26**, 629–638.

#### *Dna* Genes

[Mutations](#) in certain **genes** can adversely affect [DNA replication](#) and, as such, these genes have been designated as “*dna* genes.” Each *dna* gene falls within one of two groups distinguished by how the [kinetics](#) of DNA synthesis are influenced when cells harboring a conditional mutation in that gene are shifted from permissive to nonpermissive conditions: (1) Mutations of genes whose products are involved in initiation of replication exhibit a slow-stop **phenotype**; DNA synthesis is sustained until ongoing rounds of replication are completed. [Proteins](#) encoded by these genes include initiator

proteins, **heat-shock** proteins, and a protein involved in **primosome** assembly. (2) Mutations of genes whose products participate in DNA elongation exhibit an immediate or fast-stop phenotype. Genes encoding **helicases**, **primases**, and subunits of **DNA polymerases** fall within this category.

## 1. Initiation Genes

Among the first temperature-sensitive DNA replication mutants to be isolated from *Escherichia coli* were those mutated in the *dnaA* gene. These mutants exhibited a slow-stop phenotype, implicating the gene product in initiation of replication (1). It is now known that the 52.5-kDa DnaA protein initiates replication by recognizing the **origin of replication**, melting the double-stranded DNA, and recruiting replication machinery to this site. The *dnaA* gene, found at 83.6 minutes on the *E. coli* chromosomal map (2), is one of the mostly highly conserved of the *dna* genes in **prokaryotes**. **Homologues** have been identified in at least 16 species of evolutionarily distinct **gram-negative** and gram-positive bacteria (3). Two regions within the gene are conserved: a short amino terminal and a long C-terminal encoding segment are highly homologous among species (3).

In *E. coli*, the *dnaA* gene is the first gene in an **operon** also containing *dnaN* (a gene encoding the b subunit of DNA polymerase III holoenzyme) and *recF* (involved in **recombination**). The gene *gyrB* follows these three genes but probably is not part of the **transcription** unit (4). Regulation of the *dnaA* gene occurs via two **promoters** found 240 and 160 nucleotides upstream of the start site, denoted dnaA1p and dnaA2p, and a DnaA binding site, called a “DnaA box,” that is found between these two promoters. In *E. coli* and *Bacillus subtilis*, DnaA protein regulates its own expression by binding to the DnaA box and inhibiting transcription of the gene (4). Other mechanisms of regulation include sequestration of the hemimethylated form of the gene and ppGpp signaling (4).

The arrangement of the operon in which *dnaA* is found is highly conserved, with the *dnaA* gene linked to the *dnaN* gene in most bacterial species. The gene arrangement, *dnaA-dnaN-recF-gyrB*, has been identified in numerous bacterial species, including *E. coli* (3), *Caulobacter crescentus* (5), *Proteus mirabilis* (6), *Staphylococcus aureus* (7), and *Mycobacterium tuberculosis* (8). *Borrelia burgdorferi* has a unique gene arrangement, *rpmH-dnaN-dnaA-gyrB*, which is most likely the result of a genetic rearrangement during evolution (9). In addition to its proximity to the *dnaN* gene, *dnaA* is found closely linked to *oriC* in most bacterial species. An exception is *E. coli*, where the *dnaA* gene is found almost 42 kbp counterclockwise of *oriC* (4).

Numerous extragenic **suppressors** of *dnaA* mutations in *E. coli* suggest that there is a complicated interaction of DnaA with numerous other proteins. For example, specific mutations in the genes encoding the b subunit of RNA polymerase (*rpoB*), topoisomerase I (*topA*), **thioredoxin** (*trx*), and the **chaperonins** GroES (*groES*) and GroEL (*groEL*) are able to suppress temperature-sensitive *dnaA* mutants (3).

In addition to mutations of the *dnaA* gene, mutant forms of the *E. coli dnaC* gene (99 minutes) (2) hinder initiation of replication. The 27.9-kDa DnaC protein delivers DnaB helicase to the DnaA protein–*oriC* complex. *E. coli dnaC* is found within an operon that also contains *dnaT* (involved in primosome assembly) and an unknown gene, *pl8*. These genes are cotranscribed from a promoter located approximately 104 base pairs upstream of *dnaT* (10). Conservation of the *dnaC* gene among bacterial species is not well documented, but DnaC from *Salmonella typhimurium* can substitute for *E. coli* DnaC in an *in vitro* replication assay (11). While most temperature-sensitive *dnaC* mutants exhibit a slow-stop phenotype, two temperature-sensitive *dnaC* mutants, *dnaC1* and *dnaC7*, cause a fast stop in DNA synthesis, suggesting that DnaC protein also plays a role in elongation (12). In *B. subtilis*, mutations in the *dnaB*, *dnaC*, and *dnaD* genes render cells defective for initiation. The *dnaB* gene is found at 250 degrees on the chromosomal map (13), encodes a 472-residue protein, and bears no homology to any known *dna* genes (14). The *dnaB* gene product is required for the binding of the chromosomal origin to the membrane and is necessary for the initiation of replication (15). *DnaB* is cotranscribed with the *dnaI* gene, whose product is required for an unknown step of chromosomal replication (16). The *dnaC* gene product of *B. subtilis* has been suggested to be a DNA helicase and

has been shown to bear 42% homology to the *E. coli* helicase, DnaB (17). *B. subtilis dnaD* mutants are also defective in initiation, as a result of an unknown mechanism. This gene has been mapped to 200 degrees on the *B. subtilis* chromosome and encodes a 232-residue protein with no known homologue in *E. coli* (18).

Temperature-sensitive mutants of the *dnaR* gene of *E. coli* are also defective in initiation of chromosomal replication. This gene is found at 26.3 minutes on the chromosomal map (19) and encodes phosphoribosyl pyrophosphate (PRPP) synthetase; PRPP is a precursor for the synthesis of nucleotides and the amino acids histidine and tryptophan (20). Mutants are capable of producing enough PRPP for replication, implicating another role for PRPP in replication besides as a nucleotide precursor. Interestingly, wild-type *dnaR* is required for DnaA-dependent, but not for DnaA-independent, chromosomal replication (21). In addition, functional DnaA protein is required for mutations in the gene encoding ribulosephosphate epimerase (*rpe*) to suppress the temperature sensitivity of *dnaR* mutants. Some *dnaR* mutations are also suppressed by mutations in the gene for **RNA polymerase**, *rpoB*. It is hypothesized that the *dnaR* gene product interacts with the products of the *dnaA* and *rpe* genes during initiation of DNA replication, possibly in a transcriptional event (22).

## 2. Heat Shock Genes

The products of the *dnaJ* and *dnaK* genes are involved primarily in the heat shock response, where their expression is induced by increased temperatures or other environmental stresses. These proteins act as **molecular chaperones** that use nucleoside triphosphate hydrolysis to disrupt **hydrophobic** interactions and thus facilitate proper folding of proteins (12, 23) (see [Protein Folding In Vitro](#)). In addition to their role in the heat shock response, the *E. coli* proteins encoded by *dnaJ* and *dnaK* are required for initiation of replication at the [lambda phage](#) origin (24), which led to their designation as *dna* genes. For the most part, mutations in these genes do not affect *E. coli* chromosomal replication, but instead affect processes such as chromosome segregation, **plasmid** maintenance, and normal cell growth, division, and survival (12, 25). An exception is the *dnaK111* mutant, which is defective for initiating DNA replication at nonpermissive temperatures (26).

The *dnaJ* and *dnaK* genes are the most highly conserved of all the *dna* genes. Homologues of *E. coli dnaJ*, a member of the Hsp40 protein family, have been identified in bacterial species such as *Mycobacterium tuberculosis* (27), *B. subtilis* (27), *Clostridium acetobutylicum* (27), *B. burgdorferi* (28), *C. crescentus* (29), *S. aureus* (30), and *Lactococcus lactis* (27); in the **archaeon** *Methanosarcina mazei* (31); and in **eukaryotes** such as *Saccharomyces cerevisiae* (32), *Schizosaccharomyces pombe* (33), and *Homo sapien* (34). Homologues of *E. coli dnaK*, a member of the Hsp70 protein family, have been identified in bacterial species such as *B. subtilis* (28), *B. burgdorferi* (28), *C. crescentus* (35), *Zymomonas mobilis* (36), *Clostridium acetobutylicum* (31), *Chlamydia trachomatis* (37); in the archaeon *M. mazei* (31); as well as in eukaryotes such as *S. cerevisiae* (38), *Drosophila melanogaster* (38), and *H. sapien* (12). In *E. coli*, a single operon (39) of *dnaK-dnaJ* is found at 0.3 minutes (2). The gene arrangement *grpE-dnaK-dnaJ* has been identified in *E. coli* (28), *B. burgdorferi* (28), *B. subtilis* (40), *Chlamydia acetobutylicum* (28), *Clostridium acetobutylicum* (41), *Bacillus stearothermophilus* (42), and *S. aureus* (30). In contrast, the *dnaK* operons of *C. trachomatis* and *Lactococcus lactis subsp. lactis* are unique in that the *dnaJ* gene is absent from this region (43). In *E. coli*, *dnaJ* and *dnaK* are transcribed from a promoter located at the head of the *dnaK* cistron (44). In *B. subtilis*, two primary transcripts are synthesized by the *dnaK* operon. One is the entire *dnaK* operon, resulting in both *dnaK* and *dnaJ* transcripts, and the other is the result of transcription starting from an internal vegetative promoter immediately upstream of *dnaJ*. This promoter sequence is found in numerous gram-positive organisms (27, 42). Another common feature of organisms containing the *dnaK* operon is a nine-base-pair **inverted repeat** upstream of the promoter region. This region, called CIRCE (controlling inverted repeat of chaperone expression), acts as a negative *cis* element of the *dnaK* operon (45) and has been found in at least 27 eubacterial species (37). The *dnaK* operon in *C. crescentus* is regulated by two promoters; P1 is heat shock inducible, while P2, a sigma 70-like promoter, is under temporal control (29).

### 3. Primosome Assembly Genes

The *dnaT* gene of *E. coli* has been mapped to 99 minutes on the chromosome (2). The 19.5-kDa DnaT protein (formerly protein i) is required for bacteriophage FX174-type primosome assembly (12) and functions in transferring DnaB from the DnaB–DnaC complex to the PriA–PriB–DNA complex (46). This gene is found immediately upstream of *dnaC*, where it is transcribed along with *dnaC* and an unknown gene, *p18*, from a promoter located approximately 100 base pairs upstream of *dnaT* (10).

### 4. Helicase Genes

The *dnaB* gene of *E. coli* encodes a DNA helicase, DnaB, required to unwind double-stranded DNA at the replication fork (47). The *dnaB* gene consists of 1413 base pairs (48) and is found at 92 minutes on the chromosomal map (2). This gene lacks a typical [Shine-Dalgarno sequence](#), while a potential promoter sequence at nucleotides –32 to –37 has been identified (48). The *dnaB* genes of *S. typhimurium* and *E. coli* are identical in length (49), are 93% identical in their corresponding amino acid residues (49), and are functionally interchangeable (50). DnaB-like proteins include gene 4 helicase of coliphage T7 (49), gene 12 protein of *Salmonella* phage P22 (49), and ban (b analog) of P1 phage (12). The *dnaC* gene of *B. subtilis* is 42% identical to *E. coli dnaB* (17), but the *dnaC* gene product's potential role as a helicase has not been demonstrated.

### 5. Primase Genes

The 1740 base pair (51) *dnaG* gene of *E. coli* has been mapped to 67 minutes on the chromosome (2) and encodes DNA primase, a 65.6-kDa enzyme responsible for the synthesis of small RNA that serve as primers for initiation of replication (46). Homologues of *dnaG* have been found in *S. typhimurium* (52), *B. subtilis* (formerly called *dnaE*) (12), and *Rickettsia prowazekii* (53). The *dnaG* gene is part of a macromolecular synthesis operon, which contains genes for multiple processes such as translation, replication, and transcription (54). In *E. coli*, *dnaG* is the second gene in this operon, which also contains *rpsU* (encodes S21 protein) and *rpoD* (encodes the sigma 70 subunit of RNA polymerase). The key regulatory features of this operon include tandem promoters found upstream of *rpsU*, a termination codon between *rpsU* and *dnaG*, and an RNA processing site separating *dnaG* and *rpoD* (55). In addition, a poor ribosome binding site and the use of rare codons keeps primase levels relatively low compared to those of S21 and the sigma subunit (55). The gene order *rpsU-dnaG-rpoD* is found in at least 11 gram-negative species (54), including *S. typhimurium* (56) and *Pseudomonas putida* (57). The regulatory sequences within this operon are also somewhat conserved (54).

The macromolecular synthesis operons of *B. subtilis* and *Listeria monocytogenes* also contain the *dnaG* gene upstream of *rpoD*. However, unlike the *E. coli* operon, upstream of *dnaG* these operons contain **open reading frames**, *orf23* and *orf17*, respectively, whose products are homologous with each other, but not with *E. coli* S21 (58, 59). Like the *E. coli* operon, that of *B. subtilis* is regulated by multiple promoters. The use of rare codons in the *orf23* and primase genes keeps their expression at somewhat lower levels than that of *rpoD* (58).

### 6. DNA Polymerase Subunit Genes

In *E. coli*, DNA polymerase III is the replicative polymerase, and many of its subunits are encoded by *dna* genes. The core polymerase  $\alpha$ -subunit is encoded by the *dnaE* gene (60), which is 4600 base pairs long (61), is found at 4 minutes on the chromosomal map (2), and is in a proposed operon containing at least seven genes (62). The *dnaE* gene of *S. typhimurium* also encodes the  $\alpha$ -subunit. The product of this gene is similar in size and 97% identical to the corresponding *E. coli* protein (62). In *B. subtilis*, the  $\alpha$ -subunit is encoded by the *dnaF* gene. Mutations of this gene can result in either a slow- or fast-stop in DNA synthesis (12). Inherent in the *dnaF* gene product is a 3' to 5' exonuclease activity that is encoded by a separate gene in *E. coli*, *dnaQ*. As such, the amino-terminal

domain of the *B. subtilis dnaF* gene product has 26% homology to the  $\epsilon$ -subunit encoded by *E. coli dnaQ* (63). Given that the exonuclease activity is inherent in the core polymerase product of the *dnaF* gene in *B. subtilis*, it has been suggested that the *dnaQ* gene of gram-negative bacteria may have evolved from the *dnaF* gene of gram-positive bacteria (64).

The  $\beta$ -subunit (**sliding clamp**) of *E. coli* DNA polymerase III is encoded by the *dnaN* gene (65), which maps to 83 minutes on the chromosome (2). This subunit is responsible for the great processivity and high rates of association and dissociation of the holoenzyme from one template to another. Homologues of this gene have been identified in *B. subtilis* (formerly called *dnaG*) (12), *C. crescentus* (5), *M. luteus* (66), *M. capricolum* (67), and *B. burgdorferi* (9). In most organisms, the *dnaN* gene is found in an operon between *dnaA* and *recF* (68). In *E. coli*, the expression of the *dnaN* gene is dependent on three promoters within the translated region of *dnaA*. Although *dnaA* and *dnaN* are part of the same operon, transcription from these promoters can occur independently of transcription of *dnaA*, thus uncoupling the regulated expression of *dnaN* and *dnaA*. In addition, *dnaN* can be cotranscribed with *dnaA* from the two *dnaA* promoters, *dnaA1p* and *dnaA2p* (69).

The *dnaQ* (also called *mutD*) gene of *E. coli* encodes the 27.5-kDa  $\epsilon$ -subunit that functions as a 3' to 5' exonuclease (70). The *dnaQ* gene is the only *dna* gene that does not exhibit either a slow- or a fast-stop phenotype. Instead, mutants of this gene display a **mutator** phenotype in which there is an increased rate of spontaneous mutations (46). The *dnaQ* gene is adjacent to the gene encoding [ribonuclease H](#), *rnhA*. The two genes are transcribed in opposite directions from promoters located in an intervening 64 base pair region (71). In *S. typhimurium*, *dnaQ* also encodes the  $\epsilon$ -subunit (72).

The  $\tau$ - and  $\gamma$ -subunits of *E. coli* polymerase III are both encoded by the *dnaX* gene (12) located at 11 minutes on the chromosome (2). The 71-kDa  $\tau$ -subunit, the full-length product of the *dnaX* gene, functions in dimerization of the core polymerase to coordinate leading and lagging strand synthesis (73). The 47-kDa  $\gamma$ -subunit is encoded by only a portion of the *dnaX* gene. The  $\gamma$ -subunit is a component of the  $\gamma$ -complex, which functions in the loading of the  $\beta$  sliding clamp onto primed DNA (74). The  $\gamma$ -subunit was originally identified as a product of the *dnaZ* gene, which was later mapped within *dnaX*. Synthesis of the smaller  $\gamma$ -subunit occurs via an interesting mechanism during translation of the *dnaX* message; a -1 translational **frameshift** at a stretch of six adenines preceding a stable hairpin results in a UGA [stop codon](#) that terminates translation and thus produces the smaller  $\gamma$ -subunit (12). *DnaX* homologues encoding  $\tau$  and  $\gamma$ -subunits, as well as the frameshift signal, have been identified in *Thermus thermophilus* (75), *S. typhimurium* (76), and several other organisms. *C. crescentus* also contains a *dnaX* homologue whose expression, as for many of the replication genes of this organism, is under [cell cycle](#) control (77).

## 7. Other Bacterial *Dna* Genes

Mutations in the *dnaC* gene of *C. crescentus* (78) and the *dnaH* gene of *B. subtilis* (12) result in an immediate stop in replication. However, the exact roles of the gene products are unknown, and these genes do not appear to be homologous to any *E. coli* genes. In *E. coli* the *dnaY* gene encodes an arginine [transfer RNA](#) for a rare **codon**, and mutations in *dnaY* induce a fast-stop phenotype (12). Mutants of *B. subtilis dnaI* are inhibited for chromosomal replication at an unknown step, possibly primosome assembly. This gene, which has no known *E. coli* homologue, maps to 250 degrees on the chromosome and is cotranscribed with the *dnaB* gene (16).

## 8. Yeast *Dna* Genes

Yeast mutants defective in **mitotic** DNA replication have led to the identification and designation of numerous *dna* genes (79). Many of these have since been shown to be identical to genes known to be involved in progression through the cell cycle. For example, *DNA1*, *DNA6*, *DNA19*, *DNA33*, *DNA39*, and *DNA43* are identical to *CDC22*, *CDC34*, *CDC36*, *CDC65*, *DBF3*, and *MCM10*, respectively (80-82).



The essential *DNA2* gene of *S. cerevisiae* encodes a structurally unique 172-kDa 3' to 5' helicase. Biochemical and genetic evidence suggest that the *DNA2* helicase is involved in the resolution of [Okazaki Fragments \(83-85\)](#).

#### Bibliography

1. M. Kohiyama (1968) Cold Spring Harb. Symp. Quant. Biol. **33**, 317–324.
2. M. K. B. Berlyn, K. B. Low, and K. E. Rudd (1996) In *Escherichia Coli and Salmonella: Cellular Molecular Biology* (F. C. Neidhart et al., eds.), American Society for Molecular Biology, Washington, DC, pp. 1719–1902.
3. K. Skarstad and E. Boye (1994) Biochim. Biophys. Acta **1217**, 111–130.
4. W. Messer and C. Weigel (1996) In *Escherichia Coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt et al., eds.), American Society for Molecular Biology, Washington, DC, pp. 1579–1601.
5. R. C. Roberts and L. Shapiro (1997) J. Bacteriol. **179**, 2319–2330.
6. O. Skovgaard (1990) Gene **93**, 27–34.
7. J. C. Alonso and L. M. Fisher (1995) Mol. Gen. Genet. **246**, 680–686.
8. L. Salazar, H. Fsihi, E. de Rossi, G. Riccardi, C. Rios, S. T. Cole, and H. E. Takiff (1996) Mol. Microbiol. **20**, 283–293.
9. I. G. Old, D. Margarita, and I. Saint Girons (1993) Nucleic Acids Res. **21**, 3323.
10. H. Masai and K. Arai (1988) J. Biol. Chem. **263**, 15083–15093.
11. J. A. Kabori and A. Kornberg (1982) J. Biol. Chem. **257**, 13757–13762.
12. A. Kornberg and T. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.
13. N. Ogasawara, S. Moriya, P. G. Mazza, and H. Yoshikawa (1986) Nucleic Acids Res. **14**, 9989–9999.
14. T. Hoshino, T. McKenzie, S. Schmidt, T. Tanaka, and N. Sueoka (1987) Proc. Natl. Acad. Sci. USA **84**, 653–657.
15. S. Winston and N. Sueoka (1980) Proc. Natl. Acad. Sci. USA **77**, 2834–2838.
16. C. Bruand and S. D. Ehrlich (1995) Microbiology **141**, 1199–1200.
17. Y. Sakamoto, S. Nakai, S. Moriya, H. Yoshikawa, and N. Ogasawara (1995) Microbiology **141**, 641–644.
18. C. Bruand, A. Sorokin, P. Serror, and S. D. Ehrlich (1995) Microbiology **141**, 321–322.
19. Y. Sakakibara (1992) J. Mol. Biol. **226**, 979–987.
20. Y. Sakakibara (1992) J. Mol. Biol. **226**, 989–996.
21. Y. Sakakibara (1993) J. Bacteriol. **175**, 5559–5565.
22. Y. Sakakibara (1997) Mol. Microbiol. **24**, 793–801.
23. C. Gross (1996) In *Escherichia Coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt et al., eds.), American Society for Molecular Biology, Washington, DC, pp. 1382–1399.
24. K. Liberek, C. Georgopoulos, and M. Zylicz (1988) Proc. Natl. Acad. Sci. USA **85**, 6632–6636.
25. B. Bukau and G. C. Walker (1989) J. Bacteriol. **171**, 2337–2346.
26. Y. Sakakibara (1988) J. Bacteriol. **170**, 972–979.
27. M. van Asseldonk, A. Simons, H. Visser, W. M. de Vos, and G. Simons (1993) J. Bacteriol. **175**, 1637–1644.
28. K. Tilly, R. Hauser, J. Campbell, and G. J. Ostheimer (1993) Mol. Microbiol. **7**, 359–369.
29. M. Avedissian, D. Lessing, J. W. Gober, L. Shapiro, and S. L. Gomes (1995) J. Bacteriol. **177**, 3479–3484.
30. T. Ohta, K. Saito, M. Kuroda, K. Honda, H. Hirata, and H. Hayashi (1994) J. Bacteriol. **176**,

4779–4783.

31. E. Conway de Macario, C. B. Dugan, and A. J. Macario (1994) *J. Mol. Biol.* **240**, 95–101.
32. A. J. Caplan and M. G. Douglas (1991) *J. Cell Biol.* **114**, 609–621.
33. S. K. Park, S. K. Chon, and H. S. Yoo (1995) *Biochim. Biophys. Acta* **1262**, 87–90.
34. T. Raabe and J. L. Manley (1991) *Nucleic Acids Res.* **19**, 6645.
35. S. L. Gomes, J. W. Gober, and L. Shapiro (1990) *J. Bacteriol.* **172**, 3051–3059.
36. G. P. Michel (1993) *J. Bacteriol.* **175**, 3228–3231.
37. M. Tan, B. Wong, and J. N. Engel (1996) *J. Bacteriol.* **178**, 6983–6990.
38. J. C. Bardwell and E. A. Craig (1984) *Proc. Natl. Acad. Sci. USA* **81**, 848–852.
39. M. Ohki, F. Tamura, S. Nishimura, and H. Uchida (1986) *J. Biol. Chem.* **261**, 1778–1781.
40. M. Wetzstein, U. Volker, J. Dedio, S. Lobau, U. Zuber, M. Schiesswohl, C. Herget, M. Hecker, and W. Schumann (1992) *J. Bacteriol.* **174**, 3300–3310.
41. F. Narberhaus, K. Giebler, and H. Bahl (1992) *J. Bacteriol.* **174**, 3290–3299.
42. G. Homuth, S. Masuda, A. Mogk, Y. Kobayashi, and W. Schumann (1997) *J. Bacteriol.* **179**, 1153–1164.
43. T. Eaton, C. Shearman, and M. Gasson (1993) *J. Gen. Microbiol.* **139**, 3253–3264.
44. H. Saito and H. Uchida (1978) *Mol. Gen. Genet.* **164**, 1–8.
45. U. Zuber and W. Schumann (1994) *J. Bacteriol.* **176**, 1359–1363.
46. K. J. Mariani (1996) In *Escherichia Coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt et al., eds.), American Society for Molecular Biology, Washington, DC, pp. 749–763.
47. J. H. LeBowitz and R. McMacken (1986) *J. Biol. Chem.* **261**, 4738–4748.
48. N. Nakayama, N. Arai, M. W. Bond, Y. Kaziro, and K. Arai (1984) *J. Biol. Chem.* **259**, 97–101.
49. A. Wong, L. Kean, and R. Maurer (1988) *J. Bacteriol.* **170**, 2668–2675.
50. R. Maurer and A. Wong (1988) *J. Bacteriol.* **170**, 3682–3688.
51. B. L. Smiley, J. R. Lupski, P. S. Svec, R. McMacken, and G. N. Godson (1982) *Proc. Natl. Acad. Sci. USA* **79**, 4550–4554.
52. L. Rowen, J. A. Kabori, and S. Scherer (1982) *Mol. Gen. Genet.* **187**, 501–509.
53. G. L. Marks and D. O. Wood (1993) *Gene* **123**, 121–125.
54. J. Versalovic, T. Koeth, R. Britton, K. Geszvain, and J. Lupski (1993) *Mol. Microbiol.* **8**, 343–355.
55. Z. F. Burton, C. A. Gross, K. K. Watanabe, and R. R. Burgess (1983) *Cell* **32**, 335–349.
56. J. G. Scaife, J. S. Heilig, L. Rowen, and R. Calendar (1979) *Proc. Natl. Acad. Sci. USA* **76**, 6510–6514.
57. P. Szafranski, C. L. Smith, and C. R. Cantor (1997) *Biochim. Biophys. Acta* **1352**, 243–248.
58. L. F. Wang and R. H. Doi (1986) *Nucleic Acids Res.* **14**, 4293–4307.
59. R. Metzger, D. P. Brown, P. Grealish, M. J. Staver, J. Versalovic, J. R. Lupski, and L. Katz (1994) *Gene* **151**, 161–166.
60. M. M. Welch and C. S. McHenry (1982) *J. Bacteriol.* **152**, 351–356.
61. D. Shepard, R. W. Oberfelder, M. M. Welch, and C. S. McHenry (1984) *J. Bacteriol.* **158**, 455–459.
62. E. D. Lancy, M. R. Lifshits, P. Munson, and R. Maurer (1989) *J. Bacteriol.* **171**, 5581–5586.
63. B. Sanjanwala and A. T. Ganesan (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4421–4424.
64. Y. Huang, D. K. Braithwaite, and J. Ito (1997) *FEBS Lett.* **400**, 94–98.
65. P. M. Burgers, A. Kornberg, and Y. Sakakibara (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5391–5395.

66. M. Q. Fujita, H. Yoshikawa, and N. Ogasawara (1990) *Gene* **93**, 73–78.
67. M. Q. Fujita, H. Yoshikawa, and N. Ogasawara (1992) *Gene* **110**, 17–23.
68. M. E. Armengod, M. Garcia-Sogo, and E. Lambies (1988) *J. Biol. Chem.* **263**, 12109–12114.
69. A. Quinones and W. Messer (1988) *Mol. Gen. Genet.* **213**, 118–124.
70. R. Scheuermann, S. Tam, P. M. Burgers, C. Lu, and H. Echols (1983) *Proc. Natl. Acad. Sci. USA* **80**, 7085–7089.
71. H. Maki, T. Horiuchi, and M. Sekiguchi (1983) *Proc. Natl. Acad. Sci. USA* **80**, 7137–7141.
72. M. R. Lifshits, E. D. Lancy Jr., and R. Maurer (1992) *J. Bacteriol.* **174**, 6965–6973.
73. S. Kim, H. G. Dallmann, C. S. McHenry, and K. J. Marians (1996) *J. Biol. Chem.* **271**, 21406–21412.
74. H. Xiao, V. Naktinis, and M. O'Donnell (1995) *J. Biol. Chem.* **270**, 13378–13383.
75. O. Yurieva, M. Skangalis, J. Kuriyan, and M. O'Donnell (1997) *J. Biol. Chem.* **272**, 27131–27139.
76. A. Blinkova, M. F. Burkart, T. D. Owens, and J. R. Walker (1997) *J. Bacteriol.* **179**, 4438–4442.
77. E. Winzeler and L. Shapiro (1996) *J. Mol. Biol.* **264**, 412–425.
78. N. Ohta, M. Masurekar, and A. Newton (1990) *J. Bacteriol.* **172**, 7027–7034.
79. L. B. Dumas, J. P. Lussky, E. J. McFarland, and J. Shampay (1982) *Mol. Gen. Genet.* **187**, 42–46.
80. J. A. Prendergast, L. E. Murray, A. Rowley, D. R. Carruthers, R. A. Singer, and G. C. Johnston (1990) *Genetics* **124**, 81–90.
81. D. R. Evans, R. A. Singer, G. C. Johnston, and A. E. Wheals (1994) *FEMS Microbiol. Lett* **116**, 147–153.
82. A. M. Merchant, Y. Kawasaki, Y. Chen, M. Lei, and B. K. Tye (1997) *Mol. Cell. Biol.* **17**, 3261–3271.
83. M. E. Budd, W. C. Choe, and J. L. Campbell (1995) *J. Biol. Chem.* **270**, 26766–26769.
84. M. E. Budd and J. L. Campbell (1997) *Mol. Cell. Biol.* **17**, 2136–2142.
85. D. F. Fiorentino and G. R. Crabtree (1997) *Mol. Biol. Cell.* **8**, 2519–2537.

### **Suggestion for Further Reading**

86. A. Kornberg and T. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.

## **DNA Glycosylases**

[Base excision repair](#) is the system for repairing abnormal bases in DNA. The abnormal bases and simple base lesions are removed from DNA by the combined actions of DNA glycosylases ([1-3](#)) and AP (apurinic/apyrimidinic) endonucleases. The basic reaction carried out by glycosylases is the cleavage of the glycosylic bond joining the base to deoxyribose. Some of the glycosylases cleave both the glycosylic bond and the phosphodiester bond 3' to the resulting AP site in a more or less concerted reaction. Hence, these DNA glycosylases have been classified as either “pure” glycosylases or as glycosylase/AP lyases.

### 1. Uracil DNA Glycosylases

Uracil in DNA arises from two sources: (i) direct misincorporation into DNA in place of thymine by **DNA polymerases** and (ii) deamination of cytosine. Most DNA polymerases cannot discriminate between dTTP and dUTP; hence the rates of incorporation of thymine and uracil into DNA are proportional to the concentrations of the corresponding dNTPs in the nucleotide pool. The dTTP:dUTP ratio is about 1000:1, so about one uracil is incorporated into each [Okazaki fragment](#) during [DNA replication](#). The second source of uracil in DNA is that produced by deamination of cytosine. Cytosine deaminates at normal pH and temperature at a rate of  $10^{-16}\text{s}^{-1}$  for double-stranded DNA and at a rate of  $10^{-10}\text{s}^{-1}$  for single-stranded DNA. The deamination rate is greatly accelerated by low pH, high temperature, and oxidants. Uracil in DNA is removed by uracil glycosylases. Two have been identified in *Escherichia coli* and in humans (4): (i) uracil glycosylase and (ii) the double-strand-specific thymine glycosylase, which has also been shown to act on uracil residues in double-stranded DNA.

### 1.1. Uracil DNA Glycosylase (UDG)

The classic UDG is a 25- to 30-kDa monomeric protein with no cofactor or metal ion requirement. The enzyme is highly conserved between prokaryotes and eukaryotes, and it is the major uracil glycosylase in most organisms examined. The enzyme acts on single- and double-stranded DNA with equal efficiency. UDG is one of the several known DNA modification/repair enzymes that during catalysis flip out the base from inside the helix into a cavity within the enzyme (5-7). After cleaving off the base, the enzyme dissociates from DNA. This UDG has no detectable AP lyase activity.

### 1.2. Thymine–DNA Glycosylase (TDG)

In humans, this enzyme specifically removes thymines from T-G mismatches and uracils from U-G mismatches within double-stranded DNA (4-6). It is thought that the major source of mismatched thymines in mammalian cells is deamination of [5-methylcytosine](#) bases, which creates a T-G mismatch. Because about 20% of cytosine bases are methylated in mammalian cells, even low-frequency deamination of MeCyt generates high levels of T-G mismatches. Hence it is likely that TDG plays an important role in preventing mutations and maintaining genomic stability.

Even though DNA glycosylases in general, and UDG in particular, are considered to be highly specific enzymes, recent work has shown that UDG removes from double-stranded DNA many oxidized pyrimidines generated by ionizing radiation, and even normal pyrimidines (7). The  $K_m$  and  $k_{cat}$  catalytic parameters for removal of these lesions were found to be within the physiologically relevant range, so it is possible that UDG and TDG also perform auxiliary functions in repairing DNA damage caused by oxidative stress, in addition to their primary roles.

## 2. Thymine Glycol DNA Glycosylase

The prototype of thymine glycol DNA glycosylases is *E. coli* endonuclease III (8). It is a 40-kDa polypeptide chain with a [4Fe–4S] cluster (9) (see [Iron–Sulfur Proteins](#)). This enzyme has been detected in various organisms by one of its multiple activities and, accordingly, was given various names: endonuclease III, UV endonuclease, redoxendonuclease, and others. The enzyme has been highly conserved during [evolution](#). Enzymes with sequence and functional homology to *E. coli* endonuclease III exist in yeast and humans (10). These enzymes have a relatively wide substrate spectrum and act on ring-saturated, ring-contracted, and ring-rearranged pyrimidines. The crystal structure of *E. coli* endonuclease III has been determined (11). The enzyme flips out of damaged base into the active site and cleaves the glycosylic bond and the 3' phosphodiester bond of the resulting AP site, before dissociating from DNA. The iron–sulfur center within the protein has no catalytic role. Catalysis does not involve **redox** chemistry, but it occurs by a simple lyase reaction involving the  $\epsilon$ -[amino group](#) of a [lysine](#) residue.

## 3. Methylpurine DNA Glycosylase (MPG)

The MPG enzyme was originally identified as “3-methyladenine DNA glycosylase”; hence it is also referred to by that name (12, 13). There are two enzymes in *E. coli* with MPG activity, one that has a narrow substrate spectrum called TagI, and a second, TagII, with a much greater substrate range (14). The eukaryotic enzymes, including that in humans, are more similar to TagII with respect to their substrate preference than to TagI of *E. coli*. MPGs act on 3-methyladenine, 7-methyladenine, and *O*<sup>4</sup>-methylthymine. They clearly also repair alkylated bases other than methyl purines, and hence “alkylated purine DNA glycosylase” would perhaps be a more appropriate name for this group of enzymes. They have no overt AP lyase activity.

#### 4. 8-Oxoguanine Glycosylase

This enzyme was first identified as “formamidinopyrimidine DNA glycosylase” for the name of the substrate (FAPy) that is generated from guanine by ionizing radiation (15). Subsequently, it was found that the enzyme is very active on 8-oxoguanine, which is generated in vast quantities in the DNA by ionizing radiation and by oxidative stress. The enzyme is a 20-kDa monomer with no requirement for cofactors or divalent cations. The enzyme is widespread in the biological world, having been found in *E. coli*, yeast, and humans. It is a glycosylase and an AP lyase. In addition to cleaving the 3' phosphodiester bond by b-elimination, it also cleaves the 5' phosphodiester bond by d-elimination. Thus, this is a unique glycosylase that performs glycosylase/b-d-elimination by a concerted mechanism; in doing so, it generates a one-nucleotide gap.

#### 5. Pyrimidine Dimer Glycosylase

This enzyme has been found in two sources thus far: **T4 phage** and *Micrococcus luteus*. The enzymes from the two sources share sequence homology and are presumed to act by the same mechanism. Both are 18-kDa monomers, with no cofactor and no requirement for divalent cations. The enzyme cleaves the glycosylic bond of the 5' base of cyclobutane dimer and the intradimer phosphodiester bond by b-elimination. Thus, the cleavage generates a 3'-OH end and a 5' terminus with a dangling pyrimidine dimer (16). The 5' end must be further processed by a 5'- to 3'-exonuclease before DNA polymerase can produce a ligatable product. The structure of T4 endonuclease V complexed with substrate indicates that the enzyme flips out one of the adenine bases opposite the pyrimidine dimer (17). This binding mechanism explains the high specificity for pyrimidine dimers in double-stranded DNA.

#### 6. A-G Mismatch DNA Glycosylase

The 8-oxoguanine lesion is a frequent product of ionizing radiation and oxidative damage to DNA. Most often, 8-oxoguanine mispairs with adenine during replication. The resulting adenine-8-oxoguanine mispair is the substrate for A-G mismatch DNA glycosylase (18). This 30-kDa enzyme has limited sequence similarity to endonuclease III (thymine glycol endonuclease) and also has a [4Fe-4S] center. The protein in *E. coli* is encoded by the *mutY* gene; hence the enzyme is also referred to as MutY glycosylase. The enzyme also cleaves A residues in A-G mismatches, albeit at a lower efficiency than A-8-oxoG mismatches. In addition to glycosylase activity, it also possesses AP lyase activity and thus cleaves the 3' glycosylic bond 3' to an AP site, either as part of glycosylase/AP lyase concerted reaction or as an independent AP lyase acting on an isolated AP site.

#### Bibliography

1. T. Lindahl (1974) Proc. Natl. Acad. Sci. USA **71**, 3649–3653.
2. T. Lindahl (1976) Nature **259**, 64–66.
3. J. Laval (1977) Nature **269**, 828–832.
4. P. Gallinari and J. Jiricny (1996) Nature **383**, 735–738.
5. S. Klimasauskas, S. Kumar, R. J. Roberts, and X. Cheng (1994) Cell **76**, 357–369.

6. T. E. Barrett, R. Savva, G. Panayotou, T. Barlow, T. Brown, J. Jiricny, and L. H. Pearl (1998) *Cell* **92**, 117–129.
7. K. G. Berdal, R. F. Johansen, and E. Seeberg (1998) *EMBO J.* **17**, 363–367.
8. B. Demple and S. Linn (1980) *Nature* **287**, 203–208.
9. R. P. Cunningham, H. Asahara, J. F. Bank, C. P. Scholes, J. C. Salerno, K. Surerus, E. Munck, J. McCracken, J. Peisach, and M. H. Emptage (1988) *Biochemistry* **28**, 4450–4455.
10. T. P. Hilbert, W. Chaung, R. J. Boorstein, R. P. Cunningham, and G. W. Teebor (1997) *J. Biol. Chem.* **272**, 6733–6740.
11. C. Kuo, D. E. McRee, C. L. Fisher, S. F. O'Handley, R. P. Cunningham, and J. A. Tainer (1992) *Science* **258**, 434–440.
12. Y. Nakabeppu, H. Kondo, and M. Sekiguchi (1984) *J. Biol. Chem.* **259**, 13723–13729.
13. G. Evensen and E. Seeberg (1982) *Nature* **296**, 773–775.
14. J. Labahn, O. D. Scharer, A. Long, K. Ezaz-Nikapy, G. L. Verdine, and T. E. Ellenberger (1996) *Cell* **86**, 321–329.
15. S. Boiteux, T. R. O'Connor, and J. Laval (1987) *EMBO J.* **6**, 3177–3183.
16. W. A. Haseltine, L. K. Gordon, C. P. Lindan, R. H. Grafstrom, N. L. Shaper, and L. Grossman (1980) *Nature* **285**, 634–641.
17. D. G. Vassylyev, T. Kashiwagi, Y. Mikami, M. Ariyoshi, S. Iwai, E. Ohtsuka, and K. Morikawa (1995) *Cell* **83**, 773–782.
18. K. G. Au, S. Clark, J. H. Miller, and P. Modrich (1989) *Proc. Natl. Acad. Sci. USA* **86**, 8871–8881.

### Suggestions for Further Reading

19. K. Sakumi and M. Sekiguchi (1990) Structure and function of DNA glycosylases. *Mutat. Res.* **236**, 161–172.
20. E. Seeberg, L. Eide, and M. Bjoras (1995) The base excision repair pathway. *Trends Biochem. Sci.* **20**, 391–397.
21. R. J. Roberts (1995) On base flipping. *Cell* **82**, 9–12.

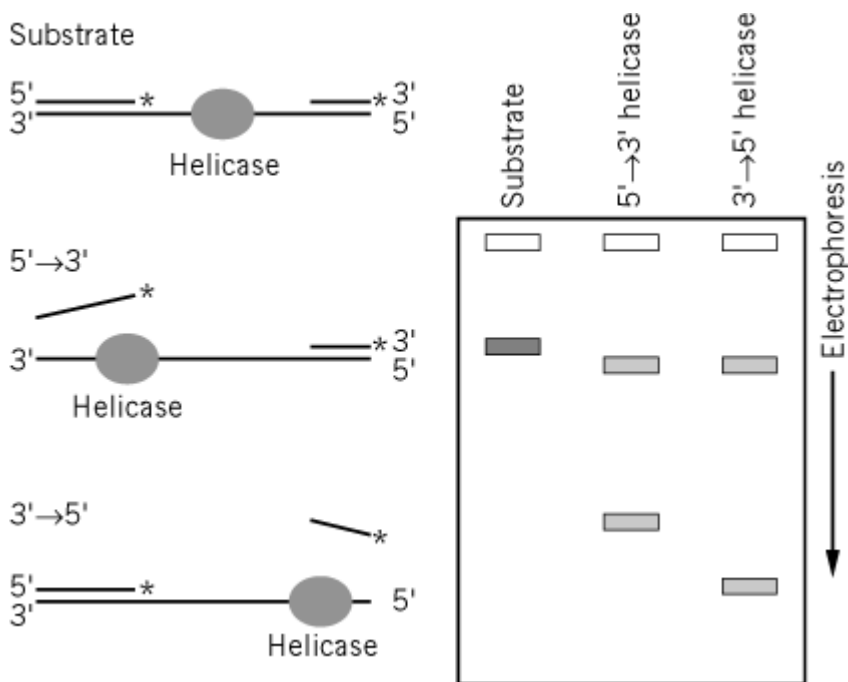
## DNA Helicase

Transient production of single-stranded **DNA** from the double helix is necessary prior to several DNA metabolic reactions, such as [DNA replication](#), [DNA repair](#), [recombination](#), and [transcription](#). DNA helicases unwind duplex DNA into two single DNA strands, ie, they disrupt the [hydrogen bonds](#) of double helical DNA by successively translocating along DNA, using the energy from hydrolysis of nucleoside triphosphates (NTP, usually ATP). Since the NTP hydrolysis is coupled to DNA binding, all helicases have DNA-dependent NTPase activity (see [ATPase](#)). Due to the involvement of helicases in various cell functions, most organisms have more than 10 helicases, and indeed, at least 12 helicases have been identified in *Escherichia coli* since the discovery of helicase I in 1976 ([1](#)). During DNA replication, a specific class of helicase unwinds duplex DNA ahead of the **replication fork**, which is the so-called “replicative helicase” and indispensable for cell proliferation. In addition to multiple helicases in cells, **bacteriophages** and eukaryotic **viruses** have their own helicases, which function in their specific propagation processes.

Unwinding of a DNA duplex by helicases is assayed independently of their particular DNA

metabolic function, using partial duplex DNA that carries a short DNA annealed to a longer single-stranded DNA. The unwinding activity is manifested by the release of the short DNA from the duplex, which is detected as a faster-migrating DNA band in [gel electrophoresis](#). The unwinding reaction proceeds in a unique direction defined as either a 5'-3' or 3'-5' polarity with respect to the strand of DNA on which the helicase is bound. The polarity used to classify the enzyme (Table 1) is determined by similar unwinding assays using two duplex DNA carrying either 5' or 3' extensions of single-stranded DNA (Fig. 1).

**Figure 1.** Helicase assay and determination of the polarity by electrophoresis.



**Table 1. DNA Helicases and Their Properties**

| Function    | Helicase       | Assembly State <sup>a</sup> | Polarity <sup>a</sup> | Source         |
|-------------|----------------|-----------------------------|-----------------------|----------------|
| Replication | DnaB           | Hexamer                     | 5' → 3'               | <i>E. coli</i> |
|             | T4 gene 41     | Hexamer                     | 5' → 3'               | T4 phage       |
|             | T7 gene 41     | Hexamer                     | 5' → 3'               | T7 phage       |
|             | SV40 T antigen | Hexamer                     | 3' → 5'               | SV40 virus     |
|             | DNA2           | Nd                          | 3' → 5'               | Yeast          |
|             | Rep            | Dimer                       | 3' → 5'               | <i>E. coli</i> |
| Repair      | RuvB           | Hexamer                     | 5' → 3'               | <i>E. coli</i> |
|             | RecBCD         | Trimer<br>(hexamer)         | Nd                    | <i>E. coli</i> |
|             | RAD26          | Nd                          | Nd                    | Yeast          |

|                           |                  |          |         |                |
|---------------------------|------------------|----------|---------|----------------|
| Repair and transcription  | (CSB/ERCC6)      |          |         | (human)        |
|                           | RAD25            | Nd       | 3' → 5' | Yeast          |
|                           | (XPB/ERCC3)      |          |         | (human)        |
|                           | RAD3 (XPD)       | Nd       | 5' → 3' | Yeast (human)  |
| Transcription termination | Rho              | Hexamer  | 5' → 3' | <i>E. coli</i> |
| Conjugation               | Helicase I/Tra I | Oligomer | 5' → 3' | F plasmid      |
| Unknown                   | Helicase III     | Oligomer | 5' → 3' | <i>E. coli</i> |
|                           | ATPase B         | Nd       | 5' → 3' | Mouse          |

<sup>a</sup> Nd = not determined.

Comparison of the [primary structures](#) of helicases revealed the presence of conserved motifs, which define several classes of the helicase family (2). A subfamily of helicases carry seven conserved motifs, but other subfamilies, which include most hexameric helicases (see below), retain only subsets of the motifs. Among the motifs, two called “A” and “B” are commonly detected in all helicases and function as a **nucleotide-binding** site. These two motifs are necessary but not sufficient for helicase functions, since the B motif corresponds to the so-called **DEAD** or DExH box, which is also observed in several [RNA-binding proteins](#). Therefore, many putative DNA or RNA helicases are proposed on the basis of just their predicted primary structures, without the identification of actual biochemical activities.

Although the precise mechanisms of their unwinding action are not understood, their characteristic oligomeric states, mainly as dimer or hexamer, propose simple explanations for their biochemical reactions. The dimeric form, such as the *E. coli* Rep helicase, provides an unwinding mechanism promoted by a rolling or inchworm movement, using two DNA-protein contacts. In this case, two DNA-binding sites alternate binding of single-stranded and duplex DNA, catalyzing unwinding. In the case of hexameric helicases, such as DnaB, T7 phage gene4 protein, and [SV40 virus large T Antigen](#), they have a common ring structure, with a central hole through which DNA passes, as analyzed by [electron microscopy](#). It is still unclear how DNA passes through the hexamer and how the unwinding reaction occurs within the molecule. Interestingly, these hexameric helicases are often members of large functional protein complexes, as in the case of typical replicative helicases, and are expected to have roles as motor molecules that may simultaneously unwind the DNA and transfer the protein complexes along DNA.

Replicative helicases have been identified only from bacteria, bacteriophage, and viruses, not from eukaryotic cells. This is mainly because there has not been any reliable *in vitro* replication system in eukaryotes to confirm whether a helicase of interest is functional in replication or not. A yeast essential gene, *dna2*, has been identified as a DNA replication mutant by screening using a permeabilized cell replication assay (3). Since this gene product has five conserved helicase motifs and actually exhibits helicase activity, it is a candidate for the yeast replicative helicase, although there is insufficient direct evidence. The known replicative helicases often have close connections with other replication activities and sometimes possess additional characteristics; eg, primase activity in T7 gene4 helicase and replication origin binding in SV40 T antigen. T antigen also interacts physically with several replication proteins and cooperatively promotes the initiation of viral DNA replication. The *E. coli* replicative helicase, DnaB, exhibits the closer physical and functional link with the replication fork components, and the rate of helicase movement is coordinately stimulated by interaction with the DNA polymerase (4).



Recent improvements in identifying the causative genes of many human heritable diseases revealed that the defects of certain helicase functions cause several human diseases, including Xeroderma pigmentosum, Cockayne's syndrome, Bloom's syndrome, Werner's syndrome, and ATR-X syndrome (5). These diseases are apparently related to the defects of genes whose functions are involved in nucleotide [excision repair](#), recombination, replication fidelity, and chromosome transmission, and they exhibit cancer-prone **phenotypes** because of the difficulty in maintaining [genome](#) integrity. Such human disease-related helicases are conserved among eukaryotes, and their counterparts are often identified as RAD genes in yeast (Table 1).

### Bibliography

1. M. Abdel-Monem et al. (1976) *Eur. J. Biochem.* **65**, 441–449.
2. A. E. Gorbalenya and E. V. Koonin (1993) *Curr. Opin. Struct. Biol.* **3**, 419–429.
3. M. E. Budd and J. L. Campbell (1995) *Proc. Natl. Acad. Sci USA* **92**, 7642–7646.
4. S. Kim et al. (1996) *Cell* **84** 643–650.
5. T. M. Lohman and K. P. Bjornson (1996) *Ann. Rev. Biochem.* **65**, 169–214.

### DNA Libraries

[Recombinant DNA](#) technologies developed in the early 1970s (1-6) form the core of all DNA-based molecular [combinatorial libraries](#). Genetic sequences of interest are recombined with a replication-competent DNA vector, such as a plasmid or bacteriophage. Each individual genetic sequence fragment is enzymatically ligated into a receptive vector, to form a unique recombinant clone. The population of individual recombinants is then reintroduced into an appropriate host cell in order to constitute a functional library. Modern cloning technologies allow the preparation of up to  $10^9$  recombinants in a single library, with greater numbers requiring a substantial scale-up that is beyond the resources of most researchers.

A different type of recombinant design is seen in some oligonucleotide-based libraries that require no vector for their survival (7-9). The replication elements in these oligonucleotides are built into their sequences in the form of primer binding sites. These libraries can be amplified *in vitro* by the polymerase chain reaction (PCR) using additional oligonucleotides capable of annealing to the termini of the library molecules (10). Other amplification technologies can be used to replicate oligonucleotide libraries, including the ligase chain reaction (11), self-sustained sequence replication (12), and strand displacement amplification (13). Because oligonucleotide-based libraries are prepared chemically and require no cloning steps, they can achieve very high sequence complexities on the order of  $1 \times 10^{15}$  different species or more.

Recent technological advances have introduced a novel type of oligonucleotide-based DNA library in which defined sets of sequences are synthesized chemically and displayed in a two-dimensional array. A highly publicized approach developed by Affymetrix employs photolithographic masking procedures to construct very high-density oligonucleotide arrays on silicon chips. One such product developed by Affymetrix is the p53 chip, which displays a large number of common disease-related [p53](#) gene mutations represented as oligonucleotide sequences. Any clinical specimen can be rapidly assayed over this chip by hybridization of nucleic acids in order to determine the genotype of p53 genes expressed by the specimen tissue. The development of further disease-related gene chips has enormous potential to carry clinical diagnosis to the molecular level where appropriate highly individualized therapies may be ascertained (eg, pharmacogenomics (14)) (see [DNA Chips](#)).

The intended use of recombinant libraries determines numerous aspects of library design and construction. The vector provides essential functions, such as an **origin of replication**, that enable replication and survival of the recombinant clone within the host cell. A selectable marker such as an [antibiotic-resistance](#) gene provided by the vector provides a means to eliminate any cells that have not taken up a recombinant clone. Unique [restriction enzyme](#) sites in the vector allow efficient insertion of foreign DNA fragments into defined contexts. Various expression control elements engineered into the vector permit selective expression of RNA or protein encoded by the cloned DNA fragment. Vectors that contain certain viral signal sequences can be packaged into viral particles in infected cells, thereby enabling highly efficient viral-mediated gene transfer into noninfected cells. In addition, [cloning](#) vectors often provide genetic markers that confer desirable properties upon the expression products. For example, DNA encoding an open reading frame may be fused in frame to the C-terminus of [glutathione-S--transferase](#) (GST) (15). The resulting GST-fusion protein can then be easily purified by [affinity chromatography](#) over [glutathione](#)-agarose resin. Finally, “shuttle” vectors contain multiple origins of replication, supporting growth of the recombinant clone in bacteria as well as in relevant experimental systems, such as yeast and mammalian cells.

Plasmid vectors are by far the most commonly used cloning vectors for relatively short DNA fragments ranging from several hundred to several thousand nucleotides. This size range is sufficient to carry the coding sequence portions of most genes. However, because genes from many organisms contain varying amounts of nonexpressed intervening sequences (**introns**), cloning of complete genes or gene clusters often requires the use of vectors that can accommodate DNA inserts exceeding 10 kilobase pairs (kbp) or more. The [lambda phage](#) particle can package approximately 50 kbp of DNA, but the lambda phage itself requires at least 30 kbp of essential phage genes for viability, limiting its insert carrying capacity to no more than 20 kbp. Some improvement is seen with “cosmid” vectors, which are composite vectors containing both plasmid replication elements and the lambda phage cos packaging signal. One type of cosmid vector can package up to 40 kbp of foreign DNA into lambda phage capsids (16). Gene mapping studies and investigations into [chromosome](#) structure require cloning vectors with even larger capacity, called *artificial chromosomes* (ACs). Thus, YACs are yeast artificial chromosomes that contain yeast [centromere](#), [telomere](#), and **autonomous replication sequence** (ARS) elements, as well as appropriate yeast selection markers (17, 18). YACs allow the cloning and maintenance of fragments of genomic DNA up to several hundred kilobase pairs in length. Similarly, ACs have been constructed that allow introduction of large foreign sequences into bacteria (called BACs) and mammalian cells (called MACs) (19).

DNA libraries may be categorized according to the source and the context of foreign DNA within the library. [Genomic libraries](#) often contain relatively large DNA inserts in phage or AC vectors. These inserts are typically not expressed, although there are exceptions. **cDNA library** inserts are derived from DNA copies of [messenger RNA](#) that are generated by the enzyme **reverse transcriptase**. mRNA is a useful source for protein-coding DNA sequences, as the nonexpressed intervening sequences have been removed through the process of pre-mRNA processing. Thus, protein-coding sequences that are interspersed throughout a gene are arranged contiguously within a cDNA, whereupon it becomes trivial to infer the encoded protein sequence using the [genetic code](#). cDNA libraries usually provide some means to express open reading frames encoded within the DNA insert, both as RNA transcripts and protein translation products. [Expression libraries](#) may, however, also be based on genomic DNA and oligonucleotides, either cloned into a vector or as linear molecules that can be amplified *in vitro*.

“Virtual libraries” are a relatively new concept in DNA libraries. Sequence information emerging from large-scale sequencing projects is incorporated into ever-expanding computer [databases](#), where cross-references can be established between any sequence and its associated functional information. Thus, a database can serve as a virtual library that can greatly expedite the identification and recovery of molecules with desired properties. Examples of public [sequence databases](#) include

GenBank and the Cancer Genome Anatomy Project (CGAP). Such virtual libraries enable tasks that once required working with actual molecules *in vitro* to be accomplished much more rapidly *in silico* using only the information contained within those molecules. For example, a scientist that identifies a sequence fragment of a protein associated with a particular disease can, through the genome databases, identify a full-length DNA sequence associated with that protein. The DNA sequence in turn provides a great deal of useful information, including the likely sequence and structure of the protein, as well as a set of primer sequences that can be used to clone the full-length gene representing that protein. Many modern database resources also maintain banks of cosmid or AC clones that are (or will be) correlated with every sequence in the database. Such resources make it possible to move from an unknown protein band obtained through [gel electrophoresis](#) to a cloned and expressed gene for that protein in a matter of days—a task that often required several years to accomplish only a decade ago.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

### Bibliography

1. R. M. Lawn, E. F. Fritsch, R. C. Parker, G. Blake, and T. Maniatis (1978) *Cell* **15**, 1157–1174.
2. S. N. Cohen, A. C. Chang, H. W. Boyer, and R. B. Helling (1973) *Proc. Natl. Acad. Sci. USA* **70**, 3240–3244.
3. N. E. Murray and K. Murray (1974) *Nature* **251**, 476–481.
4. M. Thomas, J. R. Cameron, and R. W. Davis (1974) *Proc. Natl. Acad. Sci. USA* **71**, 4579–4583.
5. T. Maniatis, R. C. Hardison, E. Lacy, J. Lauer, C. O'Connell, D. Quon, G. K. Sim, and A. Efstratiadis (1978) *Cell* **15**, 687–701.
6. G. K. Sim, F. C. Kafatos, C. W. Jones, M. D. Koehler, A. Efstratiadis, and T. Maniatis (1979) *Cell* **18**, 1303–1316.
7. C. Tuerk and L. Gold (1990) *Science* **249**, 505–510.
8. A. D. Ellington and J. W. Szostak (1990) *Nature* **346**, 818–822.
9. D. E. Tsai, D. S. Harper, and J. D. Keene (1991) *Nucleic Acids Res.* **19**, 4931–4936.
10. K. Mullis, F. Faloona, S. Scharf, R. Saiki, G. Horn, and H. Erlich (1986) *Cold Spring Harb. Symp. Quant. Biol.* **51**, 263–273.
11. F. Barany (1991) *Proc. Natl. Acad. Sci. USA* **88**, 189–193.
12. R. R. Breaker and G. F. Joyce (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6093–6097.
13. G. T. Walker, M. S. Fraiser, J. L. Schram, M. C. Little, J. G. Nadeau, and D. P. Malinowski (1992) *Nucleic Acids Res.* **20**, 1691–1696.
14. A. Persidis (1998) *Nat. Biotechnol.* **16**, 209–10.
15. D. B. Smith and K. S. Johnson (1988) *Gene* **67**, 31–40.
16. B. Hohn and J. Collins (1980) *Gene* **11**, 291–298.
17. P. Guzman and J. R. Ecker (1988) *Nucleic Acids Res.* **16**, 11091–11105.
18. A. Coulson, R. Waterston, J. Kiff, J. Sulston, and Y. Kohara (1988) *Nature* **335**, 184–186.
19. M. Ikeno, B. Grimes, T. Okazaki, M. Nakano, K. Saitoh, H. Hoshino, N. I. McGill, H. Cooke, and H. Masumoto (1998) *Nat. Biotechnol.* **16**, 431–439.

### DNA Ligase

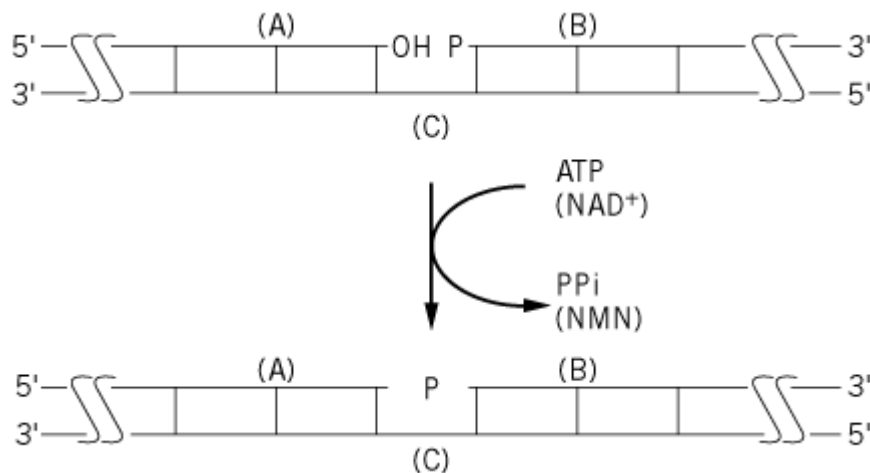
## 1. Physiological Functions

Several fundamental biological processes involve one or more steps in which an [enzyme](#), DNA ligase, catalyzes the covalent joining of one piece of **DNA** to another (1-3). For example, DNA altered spontaneously or by environmental causes is restored to its original base sequence by [DNA repair](#). DNA is duplicated by [DNA replication](#) and rearranged by homologous [recombination](#). All these processes require DNA ligase activity to link preformed pieces of DNA to one another.

During one form of DNA repair, damaged nucleotides are removed by enzymes to leave ultimately a gap in one strand of the duplex DNA, which is then filled by a **DNA polymerase** to form a product with a discontinuity (nick) between two adjacent nucleotides in the phosphodiester backbone. Likewise, semi[discontinuous DNA replication](#) leads to the formation of [Okazaki Fragments](#) that abut one another on the complementary parental DNA template strand so that a nick exists between their termini. Finally, pieces of DNA from two different chromosomes become associated during recombination so that similar nicked duplexes are formed in each chromosome. DNA ligase catalyzes the covalent joining of the DNA strands to seal the nicks and form intact, recombined duplexes. The enzyme is involved in other functions as well; whenever two pieces of DNA must be covalently joined or a DNA molecule circularized, DNA ligase probably catalyzes the reaction.

DNA fragments with appropriate termini that are aligned on a complementary strand so that their ends are juxtaposed can be joined by DNA ligase to form a phosphodiester bond (Fig. 1). For effective joining, the two strands (A and B) must be directly abutted without a gap (no missing base pairs) and must also have a 3'-OH on strand A and a 5'-PO<sub>4</sub> on strand B. The 3'-hydroxyl (-OH) of fragment A is ligated to the 5'-PO<sub>4</sub> of B to form a 3' → 5' phosphodiester bond. These two strands are abutted by base pairing with a complementary, intact strand (C). Since the formation of a phosphodiester bond is endergonic, ATP or NAD<sup>+</sup>, depending on the source of the enzyme, acts as a cofactor to supply the energy.

**Figure 1.** The overall DNA ligase reaction. A nick with a 3'-hydroxyl (—OH) on strand (A) adjacent to a 5'-phosphate (P—) terminated strand (B) on a complementary strand (C) is converted to an intact phosphodiester bond (—P—). ATP or NAD<sup>+</sup> supply the energy for bond formation (see text), and PP<sub>i</sub> or NMN, respectively, is the other product of the reaction. Long horizontal lines represent DNA strands with the polarity indicated; vertical lines represent base pairs.



## 2. Distribution

All free-living organisms are believed to produce one or more DNA ligases. Until recently, **bacteria** were thought to contain a single DNA ligase that performed all their DNA joining reactions (see text below). **Yeast** appears to have two DNA ligases. Some **bacteriophage** and animal **viruses** also encode DNA ligases. Higher **eukaryotes** probably contain four different DNA ligases that are specialized to perform particular tasks. The DNA ligases from animal and **plant** cells, **archaeobacteria**, bacteriophage, and animal viruses use ATP as the energy-supplying cofactor for the ligation reaction. **Eubacteria**, with a single reported exception, use  $\text{NAD}^+$  for this purpose. Cheng and Shuman (4) have shown that *Haemophilus influenzae* contains an ATP-dependent ligase, in addition to the expected  $\text{NAD}^+$ -dependent enzyme. Whether this finding represents an anomaly in the general cofactor-utilization classification scheme awaits the demonstration of additional ATP-dependent activities in other eubacteria. Searches of the sequences of the [genomes](#) of several other eubacteria reveal potential open reading frames that share [homologies](#) with both the expected  $\text{NAD}^+$ -dependent enzymes and with ATP-dependent DNA ligases. It will be important to determine whether any of these organisms also produce ATP-dependent DNA ligases and, if so, how they relate functionally to the  $\text{NAD}^+$ -dependent enzymes.

DNA ligase is essential because of its role in the crucial DNA transactions of the cell. Cells containing **conditional** mutants of the enzyme in *Escherichia coli* are lethal when the ligase activity is severely inactivated. They display abnormal replication, repair, and recombination **phenotypes**. Because *E. coli* cells normally contain many copies of the enzyme, small amounts of residual activity can maintain viability. The enzymes in bacteria seem to act independently in their multiple roles, that is, not in specific association with other proteins. For example, the  $\text{NAD}^+$ -dependent enzyme of *E. coli* can substitute for the inactivated bacteriophage ATP-dependent DNA ligase of an infecting bacteriophage. In eukaryotes, specific proteins may associate with DNA ligases to form complexes involved in particular functions.

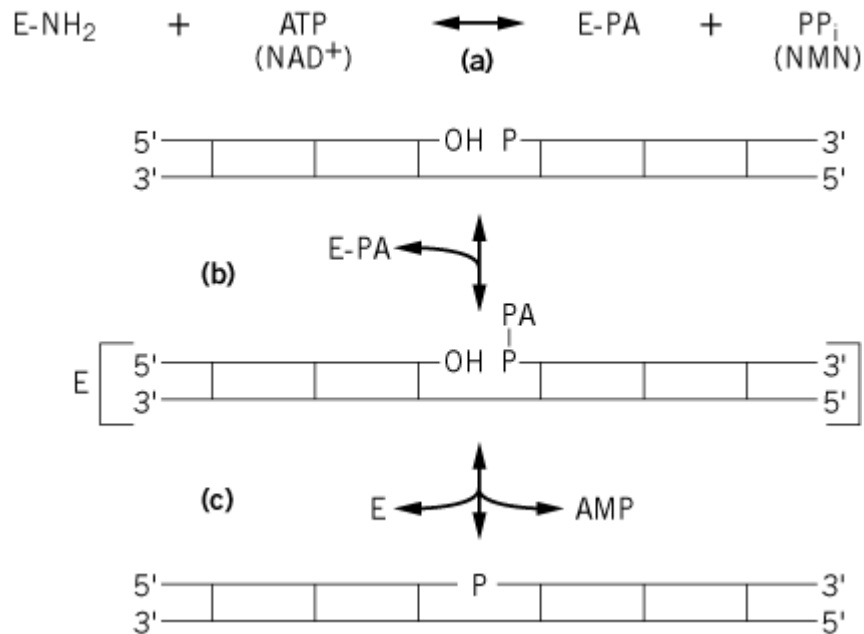
### 3. Properties

The DNA ligases comprise single [polypeptide chains](#) with molecular weights ranging from approximately 40 kDa to greater than 100 kDa, and they probably act as monomers. The **genes** of many DNA ligases have been sequenced, and several have been overexpressed and the enzymes isolated, purified, and characterized in detail. Several bacterial and bacteriophage DNA ligases are commercially available. Analyses of the sequences of genes encoding DNA ligases reveals conserved amino-acid sequence motifs among the ATP-dependent enzymes and other conserved regions among the  $\text{NAD}^+$ -dependent group. Sequence comparisons also reveal that the ATP-dependent DNA ligases are part of a superfamily that includes [RNA ligases](#) and **RNA capping** enzymes (5). The structure of only the bacteriophage T7 DNA ligase, a small, ATP-dependent enzyme, has been determined (6). It forms a bilobal structure that is reminiscent of the structures of the [methyltransferases](#) used for DNA modification.

### 4. Mechanism

Phosphodiester bond formation between template-aligned, abutted 3'-OH and 5- $\text{PO}_4$  termini proceeds in three steps, with the formation of two covalent intermediates (Fig. 2). One of these intermediates involves the enzyme, the other the DNA. In the first step of the reaction, which can occur in the absence of DNA, the enzyme reacts with ATP (or  $\text{NAD}^+$ ) to form a covalent enzyme-AMP complex. The adenylyl group is joined to the  $\epsilon$ -amino group of a [lysine](#) residue at the [active site](#) by a phosphorus-nitrogen (phosphoamide) linkage that has a large **free energy** of hydrolysis (Fig. 2a). If ATP is the cofactor,  $\text{PP}_i$  is released; if  $\text{NAD}^+$  serves as the cofactor, nicotinamide mononucleotide (NMN) is released in this step. The free energy of hydrolysis of a phosphoanhydride bond in the ATP or  $\text{NAD}^+$  is used to form the enzyme-bound, activated adenylyl group.

**Figure 2.** The DNA ligase reaction mechanism. (a) The enzyme [E-NH<sub>2</sub>] is adenylylated on a specific lysine residue [E-PA] by reaction with ATP (or NAD<sup>+</sup>) with the release of PP<sub>i</sub> (or NMN). (b) The DNA is adenylylated [AP—P—] by transfer of the adenylyl group from the enzyme to the 5'-PO<sub>4</sub> at the nick. (c) The enzyme catalyzes an attack of the 3'-OH at the nick on the activated phosphate, to form the 3' → 5' phosphodiester bond and release the enzyme. Long horizontal lines represent DNA strands with the polarity indicated; vertical lines represent base pairs.



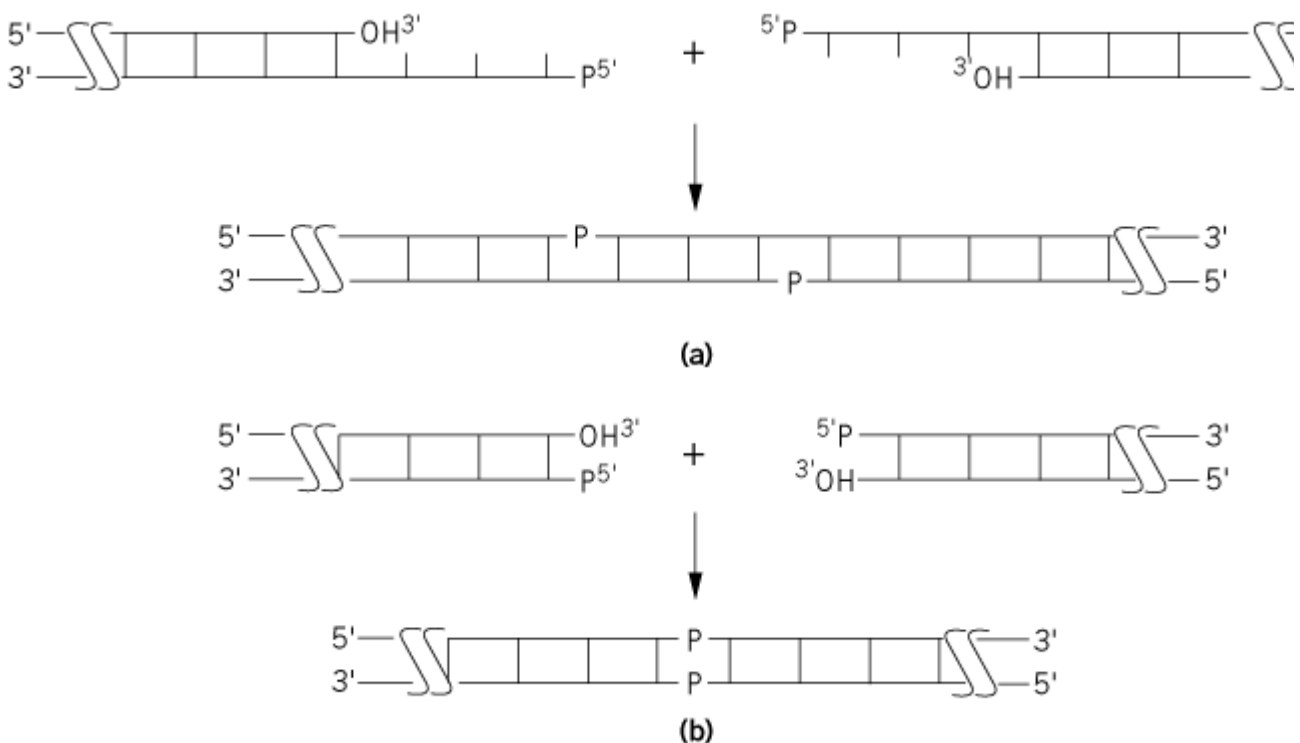
In the next step, the second covalent intermediate is formed when the adenylyl group is transferred from the enzyme to the terminal 5'-PO<sub>4</sub> at the DNA nick, to form an AMP (5' → 5') DNA phosphoanhydride linkage (Fig. 2b). In the final step of the reaction, which requires the free (unadenylylated) form of the enzyme, the AMP is displaced from the DNA by a nucleophilic attack of the 3'-OH of the adjacent DNA fragment to form the phosphodiester bond and to release the AMP from the DNA Fig. 2c).

Both covalent intermediates in the reaction have been isolated and are kinetically competent in the reaction. The reaction requires a divalent metal cation, most likely Mg<sup>2+</sup> in the cell. Steady-state kinetic analysis of the *E. coli* enzyme shows that it follows ping-pong kinetics—consistent with the described mechanism (see [Kinetic Mechanisms, Enzyme](#)). The overall reaction is reversible. The first step is readily reversed, as indicated by the facile ATP-PP<sub>i</sub> or NAD<sup>+</sup> NMN **isotope exchange reactions** catalyzed. The enzyme also catalyzes an AMP-dependent nicking-closing (topoisomerase) activity that will relax positively or negatively supercoiled DNA (see [DNA Topology](#)), indicating that the second step is also reversible.

The nature of the ends of DNA fragments affects the ability of DNA ligase to join them together. For example, duplex substrates with so called [cohesive, “sticky” ends](#) are readily joined by all DNA ligases (Fig. 3a). These substrates are formed from DNA fragments that contain short, complementary single-stranded extensions at their termini; two or more such nucleotides are required for efficient joining. These extensions can be of either strand, so long as they allow the two duplexes to anneal to one another and provide the requisite abutted 3'-OH and 5'-PO<sub>4</sub> termini. The ends of these duplex fragments hybridize to form a substrate in which there is a ligatable

discontinuity in each strand. In essence, the complementary strands serve to convert the intermolecular joining of the two DNA pieces into an intramolecular reaction by bringing them together. Another class of duplex substrates, specifically, those with **blunt ends** (Fig. 3b), which have no single-stranded extensions, are much less efficiently ligated, because the fragments to be joined are not brought into proximity by a template strand. The efficiency of this intermolecular blunt-end joining reaction depends greatly on the particular ligase. Bacteriophage T4 ligase performs the reaction best and has been extensively used for this purpose, whereas the enzyme from *E. coli* will join blunt-end duplexes only in the presence of **excluded-volume** agents, such as **polyethylene glycol**. Other DNA structures, such as those with noncomplementary, protruding ends or nonabutted fragments creating a gap on a template strand, can be ligated, but inefficiently. T4 DNA ligase is again more tolerant of these unusual structures than is the bacterial enzyme.

**Figure 3.** Common DNA ends joined by DNA ligase. (a) Fragments with a 5' extension (shown here) or with a 3' extension and complementary bases (sticky ends) form a ligatable structure. (b) Ligation of blunt ends. Long horizontal lines represent DNA strands with the polarity indicated; vertical lines represent base pairs.



## 5. Polynucleotide Substrates

The physiological polynucleotide substrates for DNA ligase are predominantly DNA chains brought into proximity by a third, complementary DNA template strand as shown in Figure 1. Studies of various combinations of RNA and DNA substrates reveal that DNA ligases will also sometimes use RNA *in vitro*. RNA can be joined to RNA on either a DNA or RNA template. RNA with a 3'-OH can also be joined to DNA with a 5'-PO<sub>4</sub> on a DNA template. Likewise, DNA with a 3'-OH can be joined to RNA with a 5'-PO<sub>4</sub> on a DNA template. The exact combinations of oligonucleotide types and the efficiencies of the joinings again depend on the particular enzyme used and the reaction conditions employed. T4 DNA ligase is again often more permissive in these “unnatural” ligations than is the *E. coli* enzyme. On binding to the DNA, the ATP-dependent enzymes may specifically recognize a nick in the duplex that contains the requisite 3'-OH and 5'-PO<sub>4</sub> termini (7). [Vaccinia](#)

[virus](#) DNA ligase appears to require a [B-DNA](#) helix to the 5' side of the nick, suggesting that substrate conformation is another parameter that can influence recognition (8). It should be noted that some studies of substrate requirements were carried out with homopolymeric substrates and that the results might be quantitatively different for substrates containing heterogeneous sequences.

## 6. Assays

Several assays have been used to quantify DNA ligase activity. [<sup>3</sup>H]poly d(A-T) folds on itself to form base pairs that create a dumbbell-like structure with a nick that can be sealed by DNA ligase to form a covalently closed, single-stranded circle. This product of the ligase reaction is resistant to **exonuclease** III, whereas the unsealed substrate is hydrolyzed to deoxyribonucleotides by the exonuclease. The rejoining, as determined by [gel electrophoresis](#), of a mixture of fragments formed by a [restriction enzyme](#) digestion of a **plasmid** is a common, qualitative assay. Synthetic oligodeoxyribonucleotides allow the scheme outlined in Figure 1 to serve as an assay. If strand A is labeled at its 5' end with <sup>32</sup>P, DNA ligase will convert the radiolabeled 5'-phosphate into a product with a phosphomonoester-resistant bond that can be quantified directly; alternatively, the longer product can be detected by gel electrophoresis. Analogously, [<sup>5'</sup>-<sup>32</sup>P]d(pT)<sub>n</sub> oligonucleotides aligned on a long poly d(A) template serve as substrates for a similar assay in which the products can be quantified. The restoration of transforming activity to DNA that has been partially inactivated by **deoxyribonuclease** treatment has also served as an assay. Finally, the adenylylated enzyme formed in the first step of the reaction can be assayed by using an appropriately **radiolabeled** ATP or NAD<sup>+</sup> and determining the acid-precipitable, protein-bound label, or by observing the difference in migration of the free and adenylylated enzyme during [polyacrylamide](#) electrophoresis. The first step of the reaction can also be assayed by measuring an ATP-PP<sub>i</sub> or NAD<sup>+</sup>-NMN exchange reaction.

## 7. Applications

DNA ligase is extensively used in [recombinant DNA](#) technology. The most frequently used enzymes are produced from the cloned genes of *E. coli*, *Thermus aquaticus*, and bacteriophage T4. Whenever one wishes to join two DNA fragments together, DNA ligase can be employed. For example, it is widely used in **vector** constructions to link DNA fragments to form chimeric molecules. DNA ligase is also used to attach duplex oligodeoxyribonucleotide linkers or adapters to a longer DNA to facilitate its further manipulation. A number of [site-directed mutagenesis](#) protocols depend on DNA ligase, in conjunction with a DNA polymerase, to form an intact, sequence-altered strand of DNA on an unaltered, circular [template](#) strand. Upon [transfection](#) and replication, progeny duplex molecules containing the mutation are produced. DNA ligase has often been used to join chemically synthesized, appropriately hybridized oligodeoxyribonucleotides for *de novo* gene synthesis. Because of the strict requirements for the structure and nature of its DNA substrate, the enzyme can also be used to determine whether a 3'-OH or 5'-PO<sub>4</sub> exists at a DNA nick or if a gap is present in one strand of a DNA duplex. It can be used to quantify the number of nicks in a DNA molecule by using a radiolabeled adenylylated enzyme and measuring the amount of radioactive AMP released. When used to circularize short fragments, the enzyme has been used to study the stiffness, twisting, bending, and helical repeat of DNA. For example, because the ends of a fragment must be brought into proper alignment for effective ligation, one can learn about the number of base pairs per turn of the DNA. Similarly, DNA ligase can be used to study the ability of another protein to bend DNA on binding, because the bending can facilitate the ligation by bringing the ends of the duplex fragment into proximity. The enzyme has also proved useful for the synthesis of defined-sequence RNA by efficiently joining RNA strands on a DNA template (9). The [ligase chain reaction](#), which depends on two pairs of oligodeoxyribonucleotides annealing to the top and bottom strands of a denatured target duplex at a given locus to form ligatable complexes, serves as an efficient DNA amplification method (10). The amplification is based on the principle that the products of the ligations on each strand can themselves serve as templates for further ligations of the excess starting oligonucleotides,



leading to an exponential increase in product on sequential heat denaturations and ligations. Use of a thermostable DNA ligase allows repeated denaturations of the nucleic acids without the need for addition of further enzyme. Because conditions can be chosen so that ligase strictly requires fully base-paired nucleotides at the nick, the ligase chain reaction can be used to detect mutations through the use of appropriately designed oligonucleotides (11). Oligonucleotides that hybridize to the target so that the nick is fully base-paired are amplified, whereas those that do not are unreactive. DNA ligase has also been used to enhance [DNA sequencing](#) by hybridization on microchips by joining fully base-paired oligonucleotides to form longer, detectable products (12). Although alternative schemes that do not require DNA ligase have been devised for accomplishing some of these functions, the enzyme remains a staple for nucleic acid manipulations.

### Bibliography

1. I. R. Lehman (1974) *Science* **186**, 790–797.
2. A. E. Tomkinson and D. S. Levin (1997) *BioEssays* **19**, 893–901.
3. M. J. Engler and C. C. Richardson (1992) in *The Enzymes*, Vol. **XV**, pp. 3–29.
4. C. Cheng and S. Shuman (1997) *Nucleic Acids Res.* **25**, 1369–1374.
5. S. Shuman and B. Schwer (1995) *Mol. Microbiol.* **17**, 405–410.
6. H. S. Subramanya, A. J. Doherty, S. R. Ashford, and D. Wigley (1996) *Cell* **85**, 607–615.
7. V. Sriskanda and S. Shuman (1998) *Nucleic Acids Res.* **26**, 525–531.
8. J. Sekiguchi and S. Shuman (1997) *Biochemistry* **36**, 9073–9079.
9. M. J. Moore and P. A. Sharp (1992) *Science* **256**, 992–997.
10. M. Wiedmann, W. J. Wilson, J. Czajka, J. Luo, F. Barany, and C. A. Batt (1994) *PCR Applications* **3**, S51–S64.
11. K. Abravaya, J. J. Carrino, S. Muldoon, and H. H. Lee (1995) *Nucleic Acids Res.* **23**, 675–682.
12. S. Dubiley, E. Kirillov, Y. Lysov, and A. Mirzabekov (1997) *Nucleic Acids Res.* **25**, 2259–2265.

### Suggestions for Further Reading

13. S. Shuman (1996) Closing the gap on DNA ligase, *Structure* **4**, 653–656. (A minireview that correlates the T7 DNA ligase structure with the sequences of the superfamily of nucleotidyl transferases to examine the mechanisms of the reactions.)
14. H.-M. Eun (1996) *Enzymology Primer for Recombinant DNA Technology*, Academic Press, San Diego, pp. 109–132. (A review of the properties, assay conditions, and applications of *E. coli* and bacteriophage T4 DNA ligases.)

### DNA Polymerase I and Klenow Fragment

*Escherichia coli* DNA polymerase I (pol I) has served for several decades as the prototype **DNA-dependent DNA polymerase**. Pol I, one of three such [polymerases](#) found in *E. coli*, has important roles during [DNA replication](#), genetic [recombination](#), and [DNA repair](#). Pol I, the first polymerase discovered in bacteria, is required for rapid, efficient growth in rich media. This enzyme contains three activities: 5′-3′ polymerase, 3′-5′ exonuclease (for proofreading), and 5′-3′ exonuclease (for **nick translation**, [excision repair](#), and hydrolysis of the RNA primers during DNA replication). Mild proteolytic digestion of *E. coli* DNA polymerase I (103 kDa) results in forming of two products: (1) the large Klenow fragment (68 kDa) that contains both the DNA polymerase and 3′-5′ exonuclease

(proofreading) activities and (2) a smaller fragment (35 kDa) that contains the 5'-3' exonuclease activity. Pol I is the most abundant polymerase in *E. coli* (400 molecules per cell), and it functions primarily to fill DNA gaps during repair and replication. The early [cloning](#) and overexpression of the Klenow fragment (1) and the determination of its three-dimensional structure (2) have allowed better understanding of polymerization mechanisms, molecular mechanisms of 3'-5' exonuclease activity, and structure-function relationships of polymerases.

## 1. Subunits of Pol I and Their Properties

DNA polymerase I, a product of the *PolA* gene which maps at minute 86 of the *E. coli* K12 chromosome (3), is a single-subunit enzyme that functions autonomously of other replicative factors. Pol I is cleaved into two subunits by mild **proteolysis** with **trypsin**. The N-terminal 5'-3' exonuclease **domain** contains 323 amino acid residues, and the C-terminal fragment contains 605 residues (4, 5), which is frequently called the Klenow fragment.

Much of the biochemical understanding of the polymerase and 3'-5' exonuclease functions stems from studies of the Klenow fragment. The Klenow fragment incorporates individual deoxyribonucleotides at a relatively modest rate of 50 nucleotides/s (6), and is moderately processive (on average, it synthesizes 50 nucleotides after binding and before dissociation). The Klenow fragment also acts as a **reverse transcriptase** by copying RNA templates (7), albeit distributively (i.e., incorporating one to two nucleotides per binding event), but the biological relevance of this activity is not known. The fidelity of Klenow polymerase during DNA-templated DNA polymerization is approximately 1/50,000 (i.e., one misincorporation per 50,000 nucleotides), and this fidelity is derived in part from its 3'-5' editing activity. Mutant Klenow fragments that have nonfunctional proofreading activity have an approximate fidelity of  $10^{-3}$  (8). Both the fidelity and processivity of the Klenow fragment are influenced by reaction conditions. Both parameters are enhanced when the pH is decreased from 9.8 to 6.2 (9).

The smaller 35-kDa domain that has 5'-3' exonuclease activities functions in pol I to degrade DNA and RNA to monomers and small oligomers. The polymerase and 5'-3' exonuclease domains act in concert during nick translation to remove [Okazaki Fragments](#) by using the 5'-3' exonuclease activity and to fill in the resulting gap. The entire *polA* gene can be deleted without effecting the viability of *E. coli* during growth in a minimal medium.

## 2. *In Vivo* Functions

Much of the understanding about the *in vivo* roles of Pol I has resulted from studying *PolA* mutants. These studies show that the entire *polA* gene can be deleted without affecting cell growth and viability during slow growth in minimal media. However, either the large 68-kDa fragment or the small 35-kDa fragment is necessary for growth in rich media (10). In the absence of sufficient pol I activity, either Pol III or Pol II (which lacks 5'-3' exonuclease activity) and [ribonuclease H](#) (the enzyme that removes RNA from RNA/DNA hybrids) potentially substitute for Pol I. The mutant strains lacking the *polA* gene are susceptible to UV light and DNA damage resulting from the alkylating agent methylmethane sulfonate. The *polA*-deficient strains also remove Okazaki fragments 10-fold slower than the wild type. A *polA* mutant that contains an [amber mutation](#) (i.e., a [stop codon](#) that produces a truncated protein), resulting in diminished ( $10^{-2}$ ) polymerase activity but nearly normal exonuclease activity, is sensitive to UV and DNA alkylation damage, suggesting pol I has an essential role during DNA repair (11). Further genetic studies with *polA-lacZ* gene fusion showed that exposure to DNA-damaging agents, including 4-nitroquinoline-*N*-oxide, UV light, mitomycin C, and methyl methanesulfonate, augments Pol I expression (12).

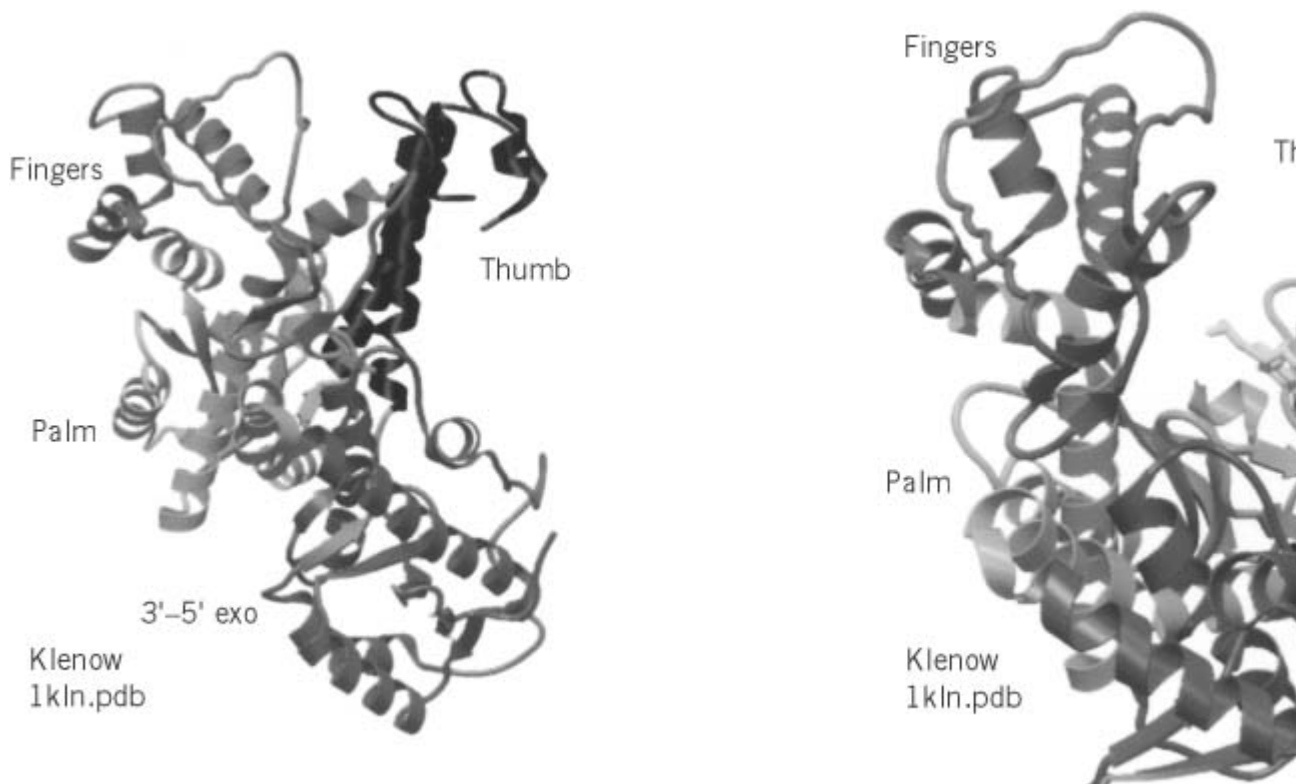
The importance of Pol I's role in nucleotide excision repair is well accepted. Nucleotide excision repair functions to remove bulky DNA adducts and pyrimidine dimers that occur after UV exposure.

During removal of pyrimidine dimers, the UvrABC protein complex binds to the lesion and forms nicks about 11 nucleotides apart that flanking the damaged residue. Next, the **helicase** UvrD displaces the DNA between the nicks, thus removing the damaged portion. Lastly, Pol I fills the resulting 11-nucleotide gap. Pol I uses its essential 5'-3' exonuclease activity to cleave the 5' nucleotides of the gap. The balance of polymerase activity and simultaneous 5-3' exonuclease activity ensures proper closure of the gap, leaving behind a nick. The last step in this repair process is ligation of the nick by a [DNA Ligase](#). DNA pol I also uses this nick translation activity to polymerize DNA and remove Okazaki fragments simultaneously during replication.

### 3. Structure of Klenow Fragment

The Klenow fragment was the first polymerase whose three-dimensional structure was solved (2). The **X-ray crystallographic** structure of the Klenow fragment morphologically resembles the human right hand, complete with finger, palm, and thumb domains (Fig. 1). The structure illustrates that the N-terminal 3'-5' exonuclease domain and the C-terminal polymerase domain are separated by about 30 Å. The structure of the Klenow fragment complexed with double-stranded DNA in an "editing mode" shows that about 5 to 8 bp of DNA are housed by the Klenow fragment in a large positively charged cleft, and about 4 nucleotides of the primer from the 3' end need to be melted from the template and then shifted to the exonuclease site during exonuclease degradation of the primer strand (13). If the DNA were housed in the polymerase **active site**, it is predicted that there would be a 90° bend, similar to that in the HIV-1 **reverse transcriptase** /DNA structure (14, 15).

**Figure 1.** Ribbon diagram of the Klenow fragment in the absence and presence of DNA. The Klenow fragment morphology contains fingers, thumb, palm, and 3'-5' exonuclease subdomains. The primer strand of the DNA is bound in the 3'-5' exonuclease site. A conformational change to the polymerizing mode allows continuation of DNA synthesis. This conformational change involves transferring the primer terminus across 30 Å, complete hybridization of the primer-template subdomain.



The polymerase active site, where dNTP binding and the chemical catalytic steps occur is located in the palm subdomain. It is thought from numerous biochemical and sequence alignment studies that the triphosphate of the incoming dNTP/divalent metal complex binds at the conserved Asp882, Glu883, and Asp705, whereas the base of the incoming nucleotide and its base pairing with the template is sensed by Tyr766 located in the conserved **a-helix O**. There is currently some confusion about the roles of Arg754 and Lys758, the two other conserved residues of helix O. Some believe that these residues bind the phosphates of the incoming dNTP. However, it is more likely that these residues are involved in template binding or perhaps pyrophosphate exchange (16).

The Klenow/double stranded DNA complex locked in an “editing mode” has given us numerous insights into the mechanism of the 3'-5' exonuclease activity. In the structure, the exonuclease domain binds four bases of partially melted single-stranded DNA. The catalytic triad of Glu357, Asp355, and Asp501 bind two divalent cations, which facilitate a nucleophilic attack on the 3' primer residue. The current model for the nucleophilic attack mechanism is that one of the metals deprotonates a [water](#) molecule, thus creating a nucleophile that, in turn, attacks the phosphate group of the 3' primer residue. The two metals are in position to stabilize the resulting pentavalent phosphorous intermediate created by the nucleophilic attack, thus reducing the [activation energy](#) of the reaction (17). A similar two-metal phosphoryl-transfer mechanism has been proposed at the polymerase active site for the Klenow fragment (18), HIV-1 reverse transcriptase (19), and DNA polymerase b (20).

#### 4. Kinetic Scheme of DNA Polymerization

A detailed kinetic mechanism has been presented using the Klenow fragment (21). The general mechanism involves these sequential steps: (1) binding of the polymerase to its template-primer, (2) binding of the appropriate dNTP to the pol-DNA complex; (3) a protein conformational change that positions the a phosphate of the incoming nucleotide near the 3'-OH group of the primer, (4) nucleophilic attack, resulting in phosphodiester bond formation (5) release of the pyrophosphate; and (6) translocation toward the new 3'-OH of the primer. Of these six steps, the protein conformational change is the slowest or rate-limiting step. The conformational change and nucleophilic attack step function together to sense proper **Watson–Crick base pairing** and thus contribute to the fidelity of this enzyme. If an improperly paired base is inserted, incorporation of the next base is very slow. This allows time for partitioning of the primer to the exonuclease site and the subsequent cleavage of the incorrect nucleoside. The presence of this proofreading 3'-5' exonuclease activity aids in promoting the fidelity of the enzyme four- to seven-fold.

#### 5. Uses in Molecular Biology

The Klenow fragment has been cloned, overexpressed, and purified, and it is widely used in laboratories. Among its uses today include incorporating radiolabeled nucleotides onto 3' DNA ends, conducting fill-in reactions of gaps in double-stranded DNA, and synthesizing blunt DNA ends after restriction digestion.

#### Bibliography

1. C. M. Joyce and N. D. Grindley (1983) *Proc. Natl. Acad. Sci. USA* **80**, 1830–1834.
2. D. L. Ollis, P. Brick, R. Hamlin, N. G. Xuong, and T. A. Steitz (1985) *Nature* **313**, 762–766.
3. B. J. Backmann (1983) *Microbiol. Rev.* **47**, 180–230.
4. D. Brutlag, M. R. Atkinson, P. Setlow, and A. Kornberg (1969) *Biochem. Biophys. Res. Commun.* **37**, 982–989.
5. H. Klenow and I. Henningsen (1970) *Proc. Natl. Acad. Sci. USA* **65**, 168–175.
6. R. D. Kuchta, P. Benkovic, and S. J. Benkovic (1988) *Biochemistry* **27**, 6716–6725.
7. L. A. Loeb, K. D. Tartof, and E. C. Travaglini (1972) *Nat. New Biol.* **242**, 66–69.

8. K. Bebenek, C. M. Joyce, M. P. Fitzgerald, and T. A. Kunkel (1990) *J. Biol. Chem.* **265**, 13878–13887.
9. K. A. Eckert and T. A. Kunkel (1993) *J. Biol. Chem.* **268**, 13462–13471.
10. C. M. Joyce and N. D. Grindley (1984) *J. Bacteriol.* **158**, 636–643.
11. W. S. Kelley and C. M. Joyce (1983) *J. Mol. Biol.* **164**, 529–560.
12. G. Wandt, S. Kubis, and A. Quinones (1997) *Mol. Gen. Genet.* **254**, 98–103.
13. L. S. Beese, V. Derbyshire, and T. A. Steitz (1993) *Science* **260**, 352–355.
14. C. M. Joyce and T. A. Steitz (1994) *Annu. Rev. Biochem.* **63**, 777–822.
15. A. Jacobo-Molina, J. Ding, R. G. Nanni, A. D. Clark Jr., X. Lu, C. Tantillo, R. L. Williams, G. Kamer, A. L. Ferris, P. Clark, A. Hizi, S. H. Hughes, and E. Arnold (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6320–6324.
16. M. Suzuki, A. K. Avicola, L. Hood, and L. A. Loeb (1997) *J. Biol. Chem.* **272**, 11228–11235.
17. L. S. Beese and T. A. Steitz (1991) *EMBO J.* **10**, 25–33.
18. T. A. Steitz and J. A. Steitz (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6498–6502.
19. P. H. Patel, A. Jacobo-Molina, J. Ding, C. Tantillo, A. D. Clark Jr., R. Raag, R. G. Nanni, S. H. Hughes, and E. Arnold (1995) *Biochemistry* **34**, 5351–5363.
20. H. Pelletier, M. R. Sawaya, A. Kumar, S. H. Wilson, and J. Kraut (1994) *Science* **264**, 1891–1903.
21. R. D. Kuchta, V. Mizrahi, P. A. Benkovic, K. A. Johnson, and S. J. Benkovic (1987) *Biochemistry* **26**, 8410–8417.

### Suggestions for Further Reading

22. A. Kornberg and T. A. Baker (1992) *DNA Replication*, Freeman, New York.
23. E. Arnold, J. Ding, S. H. Hughes, and Z. Hostomsky (1995) Structures of DNA and RNA polymerases and their interactions with nucleic acid substrates, *Curr. Opin. Struct. Biol.* **5**, 27–38

## DNA Repair

DNA repair is the set of enzymatic reactions that correct two types of structural anomalies in **DNA** (1). The first class of anomalies involves normal bases in an abnormal sequence context and includes mismatches, loops (extra nucleotides in one strand relative to the complementary strand), and bulges (a stretch of >1 noncomplementary base within the duplex). The second class, which is referred to as **DNA damage** or DNA lesions, encompasses abnormal nucleotides (modified, fragmented, cross-linked) within a normal or noncomplementary sequence context.

### 1. Mismatch Repair

Mismatches and related anomalous structures arise from errors in [DNA replication](#) and from [recombination](#) between partially **homologous** (homeologous) sequences. In addition, deamination of methylcytosine spontaneously or by acidic pH is a major cause of converting a G-C base pair to a G · T mismatch (2). The other source of simple mismatches is incorporation errors by **DNA polymerases** and the presence of a small fraction of bases in the template or in the nucleotide pool in **tautomeric** forms that lead to abnormal, yet complementary, **hydrogen bonding** and misincorporation. Loops often result from slippage/dislocation of template or nascent strand at

monotonous sequences, such as oligo dA · dT tracts or microsatellite DNA consisting of di- or trinucleotide repeats. Bulges are often caused by recombination between two related sequences with a noncomplementary sequence surrounded by homologous sequences.

Mismatches are repaired by two general mechanisms: **base-excision repair** and the general [mismatch repair](#) system (3). In the base-excision repair pathway, the mismatched base (eg, the T in the G · T mismatch) is released by a glycosylase enzyme, and the resulting abasic sugar is removed by the combined actions of *AP lyase*, which cuts the polynucleotide chain 3' to the apurinic/apyrimidinic (AP) site, and *AP endonuclease*, which cuts 5' to the AP site. The resulting one-nucleotide gap is filled in by a **DNA polymerase** and sealed by a ligase (see [Base Excision Repair](#)). In the general mismatch repair system, an endonuclease incises the newly synthesized strand anywhere from 100 to 1000 nucleotides away and either 5' or 3' to the mismatch. The intervening region is removed past the mismatch by exonucleases, and the resulting gap is filled in by DNA polymerase and is ligated to restore the normal duplex (see [Mismatch Repair](#)).

Defects in the general mismatch repair system (nucleotide excision/long-patch repair) cause hereditary nonpolyposis colorectal cancer (HNPCC) in humans. In addition, mismatch repair defects have been found in several sporadic cancers (3).

## 2. Damage Repair

DNA damage in the form of base modification, base adduction, base fragmentation, and phosphodiester bond cleavage is caused by many agents including radiation in the form of X-rays (4), ultraviolet (5, 6), and microwave; and chemicals (7) ranging in reactivity from the relatively inert water to highly reactive oxygen radicals and activated polyaromatic hydrocarbons (Table 1). The most common lesions induced in DNA by ionizing radiation and active oxygen species are oxidized and fragmented purines, pyrimidine hydrates and glycols, abasic sites, and single- and double-strand breaks (8, 9). The major lesions induced by ultraviolet light are cyclobutane pyrimidine dimers and [6–4] dipyrimidine photoproducts. There is a virtually infinite variety of lesions produced in DNA by chemical agents. Among those, benzo[*a*]pyrene-, aflatoxin-, and acetylaminofluorene–guanine adducts are carcinogenic DNA lesions produced by natural and synthetic compounds. Similarly, some anticancer drugs cause DNA damage. Mitomycin C, cisplatin, and cyclophosphamide form monoadducts with guanine bases or cause interstrand cross-links between guanines in each of the two strands of the duplex. [Psoralen](#) makes intrastrand thymine monoadducts or interstrand thymine–psoralen–thymine cross-links. Alkylating agents such as methyl methane sulfonate, nitrosoguanidine, and nitrogen– and sulfur–mustard attack essentially all nucleophilic groups in DNA, producing at least 12 different base alkylation adducts and phosphorothioesters (7).

**Table 1. Causes and Types of DNA Lesions and Mechanisms of DNA Repair**

| Type of DNA lesion                              | Cause                   | Repair Mechanism           |
|---|-------------------------|----------------------------|
| A. Base mismatches, loops and bubble structures | Replication errors      | 1. Base excision           |
|   | Recombination           | 2. General mismatch repair |
| B. DNA lesion/base lesions                      | Spontaneous deamination |                            |
|   |                         |                            |

|                                 |                                      |                                      |
|---------------------------------|--------------------------------------|--------------------------------------|
| 1. Deamination and depurination | Heat                                 | 1. Base/nucleotide excision          |
| 2. Alkylation/aryl adducts      | Natural and synthetic alkylators     | 2. Base/nucleotide excision          |
| 3. Intrastrand cross-links      | Ultraviolet light, cisplatin         | 3. Nucleotide excision               |
| 4. Interstrand cross-links      | Psoralen, melphalen, cisplatin       | 4. Nucleotide excision/recombination |
| C. Strand discontinuities       |                                      |                                      |
| 1. Single-strand break          | Ionizing radiation, oxidative stress | 1. Ligation                          |
| 2. Double-strand break          | Ionizing radiation                   | 2. Ligation/recombination            |

---

DNA lesions in the form of modified or adducted bases, interstrand cross-links, or discontinuity in one (nick) or both (break) strands interfere with replication, [transcription](#) and recombination of DNA; they also cause [mutations](#) or [cell death](#), due to the inability to replicate or to a mutation in an essential gene. Furthermore, because many of the agents that cause damage are present in the cell's normal habitat, DNA damage is a frequent occurrence that the cell must deal with on a continuous basis in order to maintain its integrity and survival (9). Hence, cells possess many biochemical reaction pathways (DNA repair mechanisms) to eliminate damage from their [chromosomes](#). These pathways fall into three main categories (1, 10): (i) base excision, (ii) nucleotide excision (or simply excision), and (iii) direct repair mechanisms. As should be apparent from this enumeration, there are extensive functional similarities between the mismatch and damage repair systems. This is not surprising, because the double-helical nature of DNA dictates the types of operations that can be applied to it to remove a mismatched or a damaged base: The mismatch or damage may be removed in the form of a base (11, 12) or a (oligo)nucleotide (10), which creates a single-stranded gap that can be filled using the intact (correct) strand as a template. However, damage repair includes additional mechanisms in its repertoire. One of these is direct repair (1). In this type of repair the chemical bonds that comprise the lesion are broken to restore the normal bases. An example of this type of repair includes photoreactivation by **photolyase**, which breaks the two sigma bonds between adjacent thymines to restore the dimer to a thymine dinucleotide (13) (see [Photolyase/Photoreactivation](#)), and  $O^6$ -methylguanine DNA methyltransferase, which breaks the single bond between the  $O_6$  of guanine and the alkyl group and thus restores the normal base (14) [see [O6-Methylguanine-DNA Methyltransferase \(MGMT\)](#)]. Ligation of a single-strand or double-strand break that is not accompanied by base loss may also be considered an example of direct repair. Finally, [recombinational repair](#) is also a repair mechanism unique to damaged DNA. This repair mechanism is employed by the cell when both strands of the duplex are damaged so that the duplex cannot be fixed by direct repair or by base and nucleotide excision repair systems, because these systems rely on an intact strand in the duplex (1). In this rather elaborate repair mechanism, the damaged region of the duplex is replaced by using a homologous duplex to retrieve genetic material and/or information. This repair system uses the enzymes of both nucleotide excision repair and genetic recombination systems (see [Recombinational Repair](#)).

## Bibliography

1. A. Sancar and G. B. Sancar (1988) DNA repair enzymes. *Annu. Rev. Biochem.* **57**, 29–67.
2. T. A. Kunkel (1992) DNA replication fidelity. *J. Biol. Chem.* **267**, 18251–18254.
3. P. Modrich and R. Lahue (1996) Mismatch repair in replication fidelity, genetic recombination, and cancer biology. *Annu. Rev. Biochem.* **65**, 101–133.

4. M. Dizdaroglu (1992) Measurement of radiation induced damage to DNA at the molecular level. *Int. J. Radiat. Biol.* **61**, 175–183.
5. D. E. Brash (1988) UV mutagenic photoproducts in *E. coli* and human cells: a molecular genetics perspective on human skin cancer. *Photochem. Photobiol.* **48**, 59–66.
6. S. Tornaletti and G. P. Pfeiffer (1996) UV damage and repair mechanisms in mammalian cells. *Bioessays* **18**, 221–228.
7. B. Singer and J. T. Kusmierek (1982) Chemical mutagenesis. *Annu. Rev. Biochem.* **51**, 655–693.
8. E. S. Henle and S. Linn (1997) Formation, prevention and repair of DNA damage by iron/hydrogen peroxide. *J. Biol. Chem.* **272**, 19095–19098.
9. K. B. Beckman and B. A. Ames (1997) Oxidative decay of DNA. *J. Biol. Chem.* **272**, 19633–19636.
10. A. Sancar (1994) Mechanisms of DNA excision repair. *Science* **266**, 1954–1956.
11. E. Seeberg, L. Eide, and M. Bjoras (1995) The base excision repair pathway. *Trends Biochem. Sci.* **20**, 391–397.
12. M. L. Dodson, M. L. Michaels, and R. S. Lloyd (1994) Unified catalytic mechanism for DNA glycosylases. *J. Biol. Chem.* **269**, 32709–32712.
13. A. Sancar (1996) No “End of History” for photolyases. *Science* **272**, 48–49.
14. L. Samson (1992) The suicidal DNA repair methyltransferases of microbes. *Mol. Microbiol.* **6**, 825–831.

## DNA Replication

DNA replication in a chemical sense is the process by which an exact copy of a **DNA** molecule having a specific base sequence is synthesized. Exact copies of a linear DNA molecule can be replicated *in vitro* using purified **DNA polymerases** and proper [primers](#). A recent advance in the long **PCR** method has made it possible to amplify DNA molecules as long as several tens of kilobases. In biology, however, DNA replication is defined as the duplication of the entire [genome](#) DNA in the cell. Basic mechanisms of replication of **plasmids**, **bacteriophages**, animal **viruses**, and bacterial [chromosomes](#) have been elucidated at the molecular level. On the other hand, knowledge of the replication of **eukaryotic** genomes is still limited, although increasing rapidly. A variety of structures and types of replication are known among plasmids, bacteriophages, plant and animal viruses. [Mitochondria](#) genomes are known to have a unique mode of replication. Here, replication of the genomes (conventionally called chromosomes) of **prokaryotes** (bacteria) and eukaryotes will be described, since the replication of genome DNA is a process essential for basic cellular functions, cell cycle, cell division, and cell differentiation.

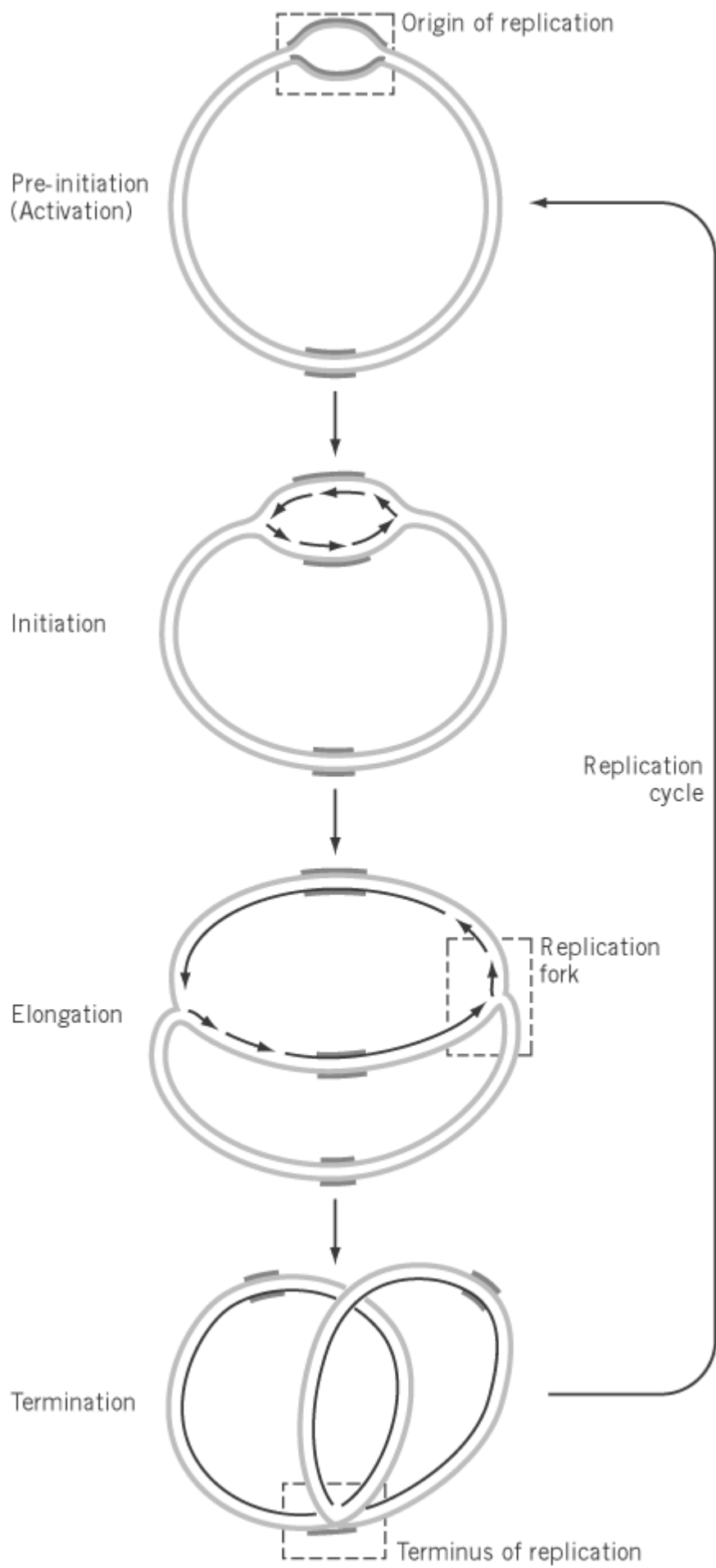
### 1. Prokaryotic Genomes (Chromosomes)

Bacterial genomes usually consist of circular DNA of one to eight megabases and comprise a single replication unit, a [replicon](#). The general aspects of genome replication have been elucidated in *Bacillus subtilis* and *Escherichia coli*, representatives of **Gram positive** and Gram negative bacteria, respectively. In both genomes, replication is initiated from a genetically defined site, oriC (see [Replication Origin](#)), proceeds bidirectionally and symmetrically at about the same elongation rate, and terminates at the defined region, terC, of the genome (Fig. 1) (see [Termination Of DNA Replication](#)). Isolation of *dna* temperature-sensitive mutants has revealed that the three processes,



initiation, elongation, and termination of replication, are regulated independently by multiple gene products. Genetic and biochemical studies subsequently revealed multiple protein complexes, the *primosome* and *replisome*, for the initiation and replication machinery, respectively (see [DNA Replication Proteins](#)). In contrast, termination is achieved by a single protein, Rtp.

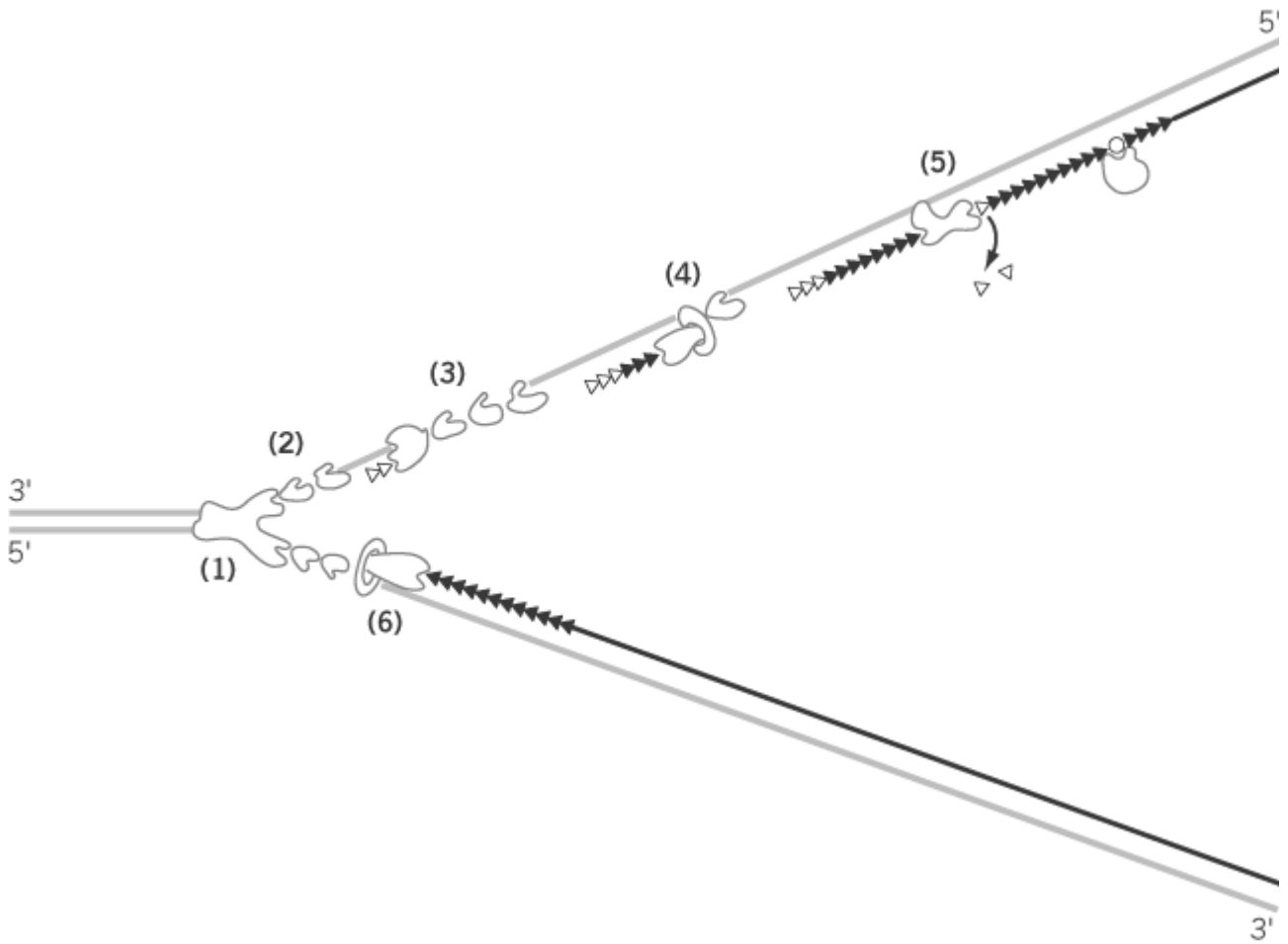
**Figure 1.** Replication of circular prokaryotic chromosome. Replication of a circular chromosome is schematically shown. Replication is initiated by unwinding of a specific site of the chromosome, the origin of replication (*oriC*), followed by the synthesis of two leading strands and lagging strands in opposite directions. Elongation proceeds bidirectionally, with the same rate of synthesis at the two replication forks. Elongation terminates at a specific site, the terminus of replication.



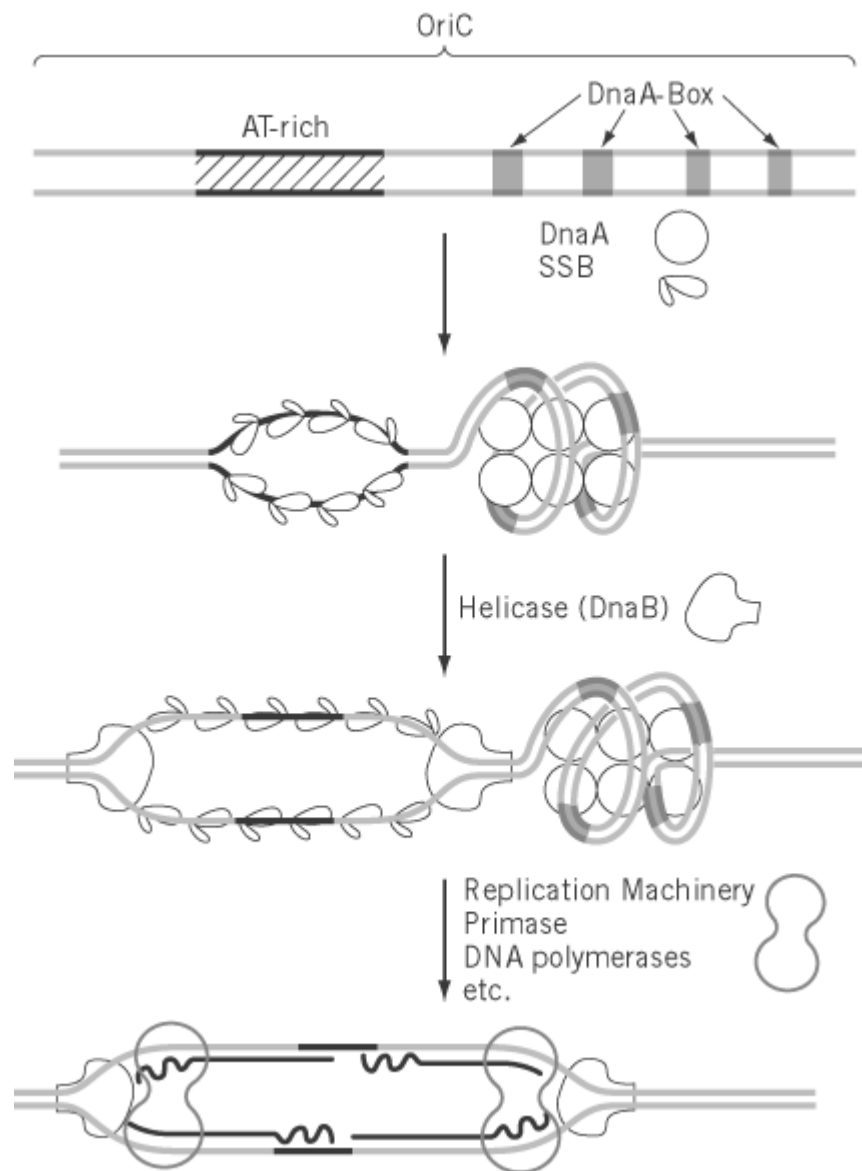
Regulatory DNA elements required for initiation, *oriC*, were identified by cloning of sequences that confer upon circular DNA the ability of autonomous replication as extra chromosomal elements in the cell (1). The *oriC* was first identified in *E. coli*, and it was then found that one of the *E. coli* Dna proteins, DnaA, functioned as an initiator of replication by binding to a sequence-specific element, the DnaA-box, to activate *oriC* and to unwind the AT-rich region within *oriC* (2). The combination of the DnaA-box and the DnaA protein functioning as *cis* and *trans* regulatory elements in the initiation of replication was subsequently found to exist commonly in many **eubacteria** (3). A replicon mechanism consisting of two regulatory elements, the replicator in *cis* and initiator in *trans*, that was proposed as early as 1963 (4) has been confirmed in prokaryotes (see [Replicon](#)).

The elongation of replication is much more complicated than it was thought in the late 1950s when [semi-conservative DNA replication](#) was demonstrated and Kornberg discovered DNA polymerase I, which apparently copied double-stranded DNA *in vitro*. This complexity is due to the fact that the two strands of DNA are oriented in opposite directions in terms of the phosphodiester bonds and DNA polymerases can function in only one, the 5' to 3' direction. Extensive analysis of growing [replication forks](#) in *E. coli* led to the discovery of asymmetric [discontinuous DNA replication](#) in which synthesis of the 5' to 3' strand proceeds sequentially (the **leading strand**), whereas the opposite strand (lagging strand) is synthesized in short pieces, using an RNA primer, and subsequently joined together. The small RNA–DNA fragments of some 500 to 1000 bases synthesized during the discontinuous replication were named [Okazaki Fragments](#) after the researcher who discovered this mechanism (5). At least five enzymes, [DNA helicase](#), [primase](#), DNA polymerase I and III, and [DNA Ligase](#), are involved in this process (Fig. 2) (see [DNA Replication Proteins](#)). These proteins are assumed to form a supramolecular complex, the replisome, and may be attached to the cell membrane. The initiation of replication is the synthesis of the first primer RNA–DNA molecule on both strands of DNA at the unwound region of *oriC*, which eventually become leading strands extended in both directions. The formation of the replisome at *oriC* requires the activation of *oriC* by DnaA and subsequent formation of the primosome complex, which includes the DNA helicase, DnaB in *E. coli*. The synthesis of lagging strands may be initiated after the formation of a considerable size of single-stranded region by elongation of the leading strand, but the detailed mechanism of synthesis of neither the first primer RNA nor the first lagging strands is known (Fig. 3).

**Figure 2.** Semiconservative and RNA-primed discontinuous replication. The mode of synthesis of leading strand and lagging strand at the replication fork is shown schematically. (1) Unwinding of double strand by helicase, (2) stabilization of single-stranded regions by [single-strand DNA binding proteins](#) (SSB) in prokaryotes (= p) or RPA in eukaryotes (= e) (3) primer RNA synthesis by primase (in p) or polymerase  $\alpha$ -primase complex (in e), (4) lagging strand synthesis by DNA polymerase III (in p) or polymerase  $\delta/\epsilon$  (in e), (5) degradation of primer RNA by polymerase I (in p) or [ribonuclease H](#) (in e), and (6) leading strand synthesis by polymerase III (in p) or polymerase  $\delta/\epsilon$  (in e).



**Figure 3.** Mechanism of initiation of replication of prokaryotic chromosomes. The OriC of *E. coli* and initiation from the oriC by the function of DnaA protein followed by DNA helicase and finally by assembly of the replication machinery is shown schematically. The DnaA box is a 9-mer sequence (consensus sequence is TTATCCACA) and DnaA protein is conserved in many eubacteria. The DnaB helicase is also conserved in at least *E. coli* and *B. subtilis*.



Elongation at the macroscopic scale proceeds bidirectionally along the circular genome at about the same rate, 50 kilobases/min, and completes at the fixed region, the terminus where two forks meet (Termination Of DNA Replication). Two sets of three (a total of six) termination signals, oriented in opposite directions, are found in both *E. coli* (6) and *B. subtilis* (3) chromosomes at about 180° from oriC, to which a protein known as replication termination protein, Rtp, binds and inhibits the elongation in one direction. The distance between the nearest two signal sets is 270 kb in *E. coli* and 59 kb in *B. subtilis*. Although the mechanism of termination is basically the same in *E. coli* and *B. subtilis*, there is no homology in these terC sequences nor in the primary structures of Rtp. This is in sharp contrast to the extensive conservation of the DnaA protein and the DnaA-box sequence among eubacteria. The tertiary structure of *B. subtilis* Rtp protein reveals the mechanism of inhibition of DNA replication at the molecular level (7). Termination of replication results in the two daughter chromosomes entwined about one another. Resolution of such a structure into two separate chromosomes can be achieved by the action of DNA gyrase (8) (see [DNA Topology](#)). In addition to the termination site, pausing sites where DNA replication slows down significantly are found near the termination sites of the *E. coli* chromosome, and their possible involvement in [recombination](#) has been discussed (9).

## 2. Eukaryotic Chromosomes

Studies on the replication of a viral genome, [SV40](#), have provided the basic knowledge about DNA

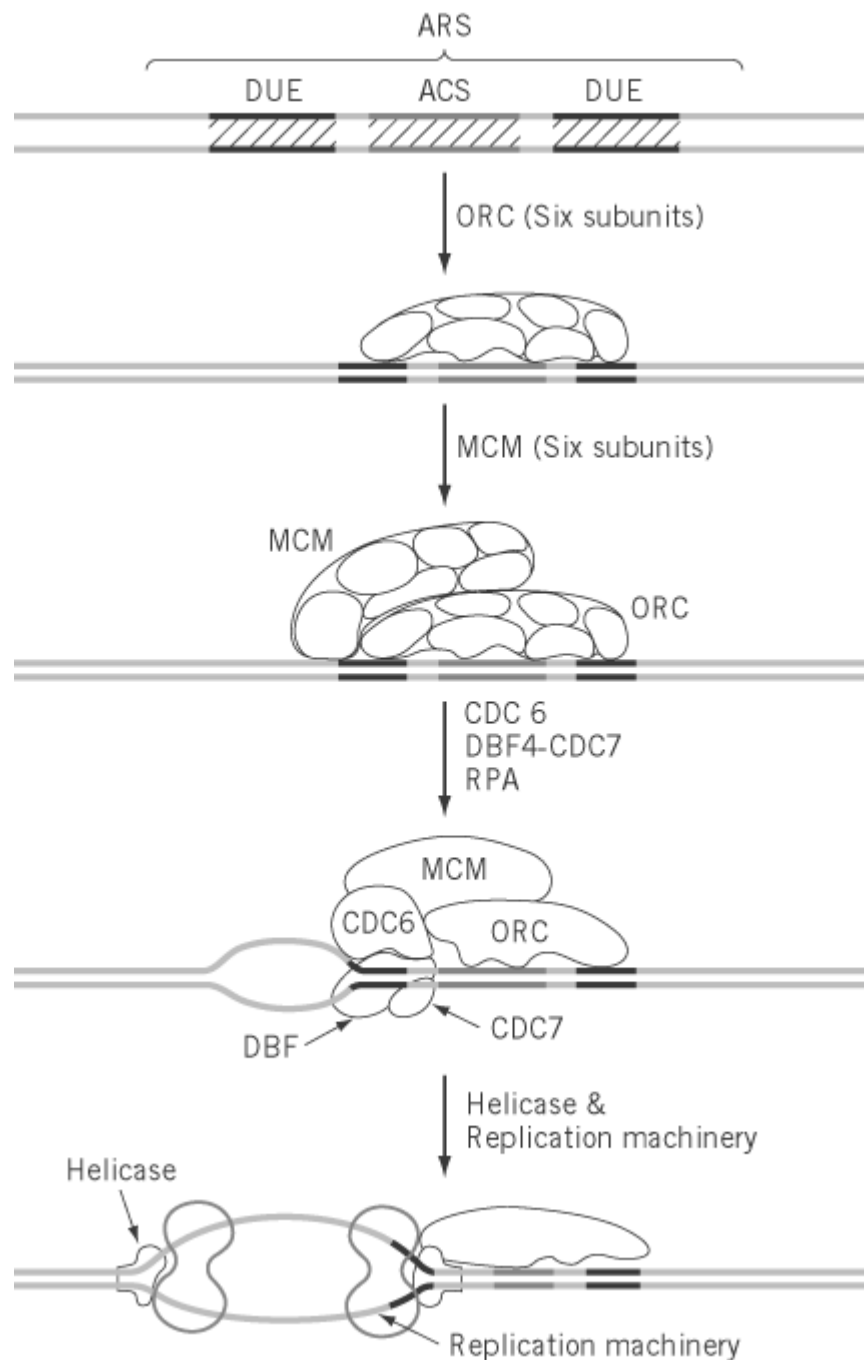
replication in eukaryotic cells, because it replicates using host proteins, except for a single viral protein, [T Antigen](#). T antigen was first identified as a protein responsible for tumor transformation of mammalian cells by the tumor virus and hence was named as Tumor antigen. Subsequently, it was found that the T antigen is an initiator protein and a helicase that acts on the origin of the 3-kb circular SV40 genome (10). Success in reconstitution of an *in vitro* replication of the entire SV40 genome led to the identification and purification of enzymes and protein factors involved in the initiation and elongation of replication and separation of replicated molecules (see [DNA Replication Proteins](#)). The basic mechanism of synthesis of leading strand by DNA polymerase  $\alpha/\epsilon$  and of lagging strand by the combination of the DNA polymerase  $\alpha$ -primase complex and polymerase  $\delta/\epsilon$  has been elucidated (11) (Fig. 2). In addition, various factors, such as RPA, RFC, and [proliferating cell nuclear antigen \(PCNA\)](#), were identified to facilitate the efficient and progressive replication of the viral genome (12). A topoisomerase, topoisomerase II, was found essential for resolving the concatenated structure formed at the end of replication of the circular genome. Since no *in vitro* replication system is available using cellular genome DNA as [template](#), the molecular mechanism of replication established by the study of SV40 genome is still the sole model of eukaryotic DNA replication and serves as a unique system to identify factors involved in DNA replication and its regulation (see [Replication Fork \(Y-Fork Intermediate\)](#)) (13).

The mechanism of initiation of replication of the SV40 genome is similar to, and rather simpler than, that of bacterial genome, because T antigen acts as helicase as well as the initiator. However, it did not provide a model for the initiation of replication of cellular chromosomes composed of multiple replicons, which was first demonstrated by the [autoradiography](#) of replicating chromosomes in 1968 by Huberman and Riggs (14). In general, the genomes of eukaryotic cells are estimated to contain about one origin every 10 to 330 kbp (15). Extensive searches for chromosomal replication origins through the cloning of [autonomously replicating sequences \(ARS\)](#) have mostly failed, except for yeast chromosomes of *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*. A number of ARS have been isolated from the 16 chromosomes of *S. cerevisiae*. In particular, 14 ARS in a 200-kbp portion of chromosome III (16) and 10 ARS in the 300-kbp entire chromosome VI (17) were identified. All ARS contain a common sequence of 11 bp ARS consensus sequence (ACS) that is essential for autonomous replication, plus additional nonspecific AT-rich sequences whose deletion causes a significant reduction in ARS activity. In some cases, a binding site for a **transcriptional factor** is located near the ACS and enhances its ARS activity. The ARS from *S. cerevisiae* is unique, and no **homologous** sequence of the approximately 100-bp ACS was found in other eukaryotes. ARS cloned from the other yeast *S. pombe* require a region of about 1 kbp containing several AT-rich clusters, with no single specific sequences like ACS essential for ARS activity. The [two-dimensional gel electrophoresis](#) of chromosomal fragments developed by Brewer and Fangman (18) has successfully detected the [eye-form intermediates](#) of replicating intermediates clearly separated from Y-form intermediates (see [Replication Fork \(Y-Fork Intermediate\)](#)). Using this method, some ARS are found to function as efficient origins of chromosomal replication in *S. cerevisiae*. However, some of the active ARS in plasmids were found to be very inefficient or silent on the chromosome, suggesting a role of **chromatin** structure in origin activity. Not all origins are fired at the same time in **S phase**, but they are initiated sequentially in a fixed order (19). A regulatory factor determining the efficiency and timing of initiation of late replicating origins have been identified (24).

Origins from genomes of higher organisms are ambiguous. Attempts to isolate ARS from mammalian chromosomes are difficult to reproduce and therefore controversial. The best-studied origin, near the [dihydrofolate reductase \(DHFR\)](#) gene in the Chinese hamster genome, showed variation in size from 0.5 to 55 kbp, depending on the methods used to detect origin activity (20). Systematic studies on large chromosomal segments have revealed that most fragments longer than 10 kbp can provide some ARS activity in mammalian cells, suggesting that DNA length is more critical than DNA sequence. In many cases, initiation occurs randomly within several kilobases, which is called the **initiation zone** rather than the origin. The concept of a replicator that is proposed in the original replicon hypothesis should be reexamined in these complex origins, in terms of the recognition site of initiators and the actual site of initiation of replication (see [Replicon](#)).

As for *trans* factors for initiation, no single initiator protein like bacterial DnaA or SV40 T antigen has been found in eukaryotic cells. Instead, a complex consisting of six proteins was found to recognize yeast ARS through binding to the ACS and named accordingly as the [origin recognition complex](#) (ORC) (21). The ORC was subsequently found to be well conserved from yeast to humans. However, the mechanism of how ORC recognizes seemingly nonspecific origin sequences on chromosomes other than *S. cerevisiae* is not clear. An *in vitro* DNA synthesis system using a [Xenopus](#) egg extract and [sperm](#) chromosomes provided a unique system for studying the biochemistry of initiation of chromosomal replication, and it led to the discovery of the [licensing factor](#) that permits replication of multiple replicons of the chromosome once in the cell cycle (22, 23). Genes homologous to the *Xenopus* licensing factor were identified among MCM genes of *S. cerevisiae*, and a MCM complex composed of six proteins that binds to ARS through interaction with ORC was subsequently discovered. Genetic studies with *S. cerevisiae* show that the formation of the MCM–ORC–ARS complex is not sufficient to activate the origin of the chromosome, and additional factors, protein [kinases](#) that interact with MCM or ORC directly or indirectly, were being discovered. Extensive studies on the regulatory network of the replication complex are expected to elucidate the molecular mechanism of [signal transduction](#) from the cell cycle engine, cdk/[cyclin](#) kinase, to the chromosomal origin of replication. The MCM complex, as well as the factors interacting with MCM and ORC, are conserved widely in eukaryotes, suggesting that the molecular mechanism elucidated in the yeast and frog provides principles that guide the studies on more complicated chromosomes including human chromosome (Fig. 4).

**Figure 4.** Mechanism of initiation of replication of a *Saccharomyces cerevisiae* chromosome. Minimal structure of the ARS contains a 11-bp sequence-specific region (ACS) that is essential for the origin, and surrounding stimulative regions that are rich in AT and called DNA unwinding elements (DUE). The ORC binds to ACS throughout cell cycle. MCM (licensing factor) complex binds to ORC at the G2-M phase of the cell cycle, and the double complex may be activated by several protein kinases, including CDC6 and DBF4/CDC7, to produce a single-stranded region that serves as the entry site for helicase and eventually for the replication machinery. The mechanism is hypothetical, as the last two steps have not been proven experimentally.



Termination of replication of adjacent replicons occurs when two forks moving in opposite direction meet. Although specific termination sites between the two adjacent replicons have not been identified, the structure of replicated replicons at the termination sites must be similar to the ends of the circular chromosomes of bacterial and viral genomes and require topoisomerases to resolve the structure to separate two sister strands. Termination of the ends of linear chromosomes requires a more specific structure, the [telomere](#), to prevent shortening of the chromosome due to the discontinuous replication mechanism.

#### Bibliography

1. Y. Hirota, S. Yasuda, M. Yamada, A. Nishimura, K. Sugimoto, H. Sugisaki, A. Oka, and M. Takanami (1978) Cold Spring Harbor Symp. Quant. Biol. **4**, 129–138.
2. D. Bramhill and A. Kornberg (1988) Cell **54**, 915–918.



3. H. Yoshikawa and R. G. Wake (1993) In *Bacillus subtilis and other Gram Positive Bacteria* (A. Sonenshein, J. A. Hoch, and R. Losick, eds.), American Society of Microbiologists, Washington DC, pp. 507–528.
4. F. Jacob, S. Brenner, and F. Cuzin (1963) Cold Spring Harbor Symp. Quant. Biol. **28**, 329–348.
5. R. Okazaki et al. (1975) In *DNA Synthesis and Its Regulation* (M. Goulian, P. Hanawalt, and F. O. Fox, eds.), Benjamin Cummings, Menlo Park, CA, pp. 832–862.
6. M. Hidaka, M. Akiyama, and T. Horiuchi (1988) Cell **55**, 467–475.
7. D. E. Bussiere, D. Bastia, and S. White (1995) Cell **80**, 651–660.
8. T. Horiuchi and Y. Fujiyama (1995) J. Bacteriol. **177**, 783–791.
9. T. R. Steek and K. Drlica (1984) Cell **36**, 1081–1088.
10. H. Stahl, P. Droge, and R. Kippers (1986) EMBO J. **5**, 1939–1944.
11. T. Tsurimoto, T. Melendy, and B. Stillman (1990) Nature **346**, 534–539.
12. T. Tsurimoto and B. Stillman (1990) Proc. Natl. Acad. Sci. USA **87**, 1023–1027.
13. B. Stillman (1994) Cell **78**, 725–728.
14. J. A. Huberman and A. D. Riggs (1968) J. Mol. Biol. **32**, 327–341.
15. R. Hand (1978) Cell **15**, 317–325.
16. S. A. Greenfeder and C. S. Newlon (1992) Mol. Biol. Cell **3**, 999–1013.
17. K. Shirahige, T. Iwasaki, M. B. Rashid, N. Ogasawara, and H. Yoshikawa (1993) Mol. Cell. Biol. **13**, 5043–5056.
18. B. J. Brewer and W. L. Fangman (1987) Cell **51**, 463–471.
19. W. L. Fangman, R. H. Hice, and E. Chlebowicz-Sledziowska (1983) Cell **32**, 831–838.
20. D. M. Gilbert, H. Miyazawa, F. S. Nallaseth, J. M. Ortega, J. J. Blow, and M. L. DePamphilis (1993) Cold Spring Harbor Symp. Quant. Biol. **58**, 475–485.
21. S. P. Bell and B. Stillman (1992) Nature **357**, 128–134.
22. J. J. Blow and R. A. Laskey (1988) Nature **332**, 546–548.
23. H. Kubota and H. Takizawa (1993) J. Cell Biol. **123**, 1321–1331.
24. K. Shirahige, Y. Hori, K. Shiraishi, M. Yamashita, K. Takahashi, C. Oluse, T. Tsurimoto, and H. Yoshikawa (1998) Nature **395**, 618–621.

### Suggestions for Further Reading

25. M. L. DePamphilis, ed. (1966) *DNA Replication in Eukaryotic Cells*, Cold Spring Harbor Laboratory Press, New York.
26. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.

## DNA Replication Proteins

Various [DNA replication](#) systems from bacteria, **plasmids**, **bacteriophage**, **viruses**, and **eukaryotic** cells have been studied to elucidate the mechanism of DNA replication. It emerged that most of them can complete the reaction with a limited number of [enzymes](#) called *replication proteins* (or factors) (Table [1](#)). The primary experiments were identification, as **conditional mutants** for bacterial growth and DNA synthesis, of the *dna genes* from *Escherichia coli* that encode replication proteins ([1](#), [2](#)).

About a dozen *dna* genes were identified and grouped as quick- or slow-stop mutants, which indicate whether the gene products are involved in the elongation or initiation stages, respectively. Development of *in vitro* DNA replication systems with *E. coli* crude lysates for single-stranded DNA phage and the *E. coli* chromosomal origin made it possible to identify and purify replication proteins by [complementation](#) of the missing activity caused by these mutations. In addition, some replication proteins were identified as proteins harboring known enzymatic activities required for replication, such as [DNA Ligase](#), [single-stranded DNA binding protein](#), [ribonuclease H](#), topoisomerase I, and DNA gyrase (see [DNA Topology](#)). DNA gyrase subunits were also identified as drug-sensitivity gene products.

**Table 1. Essential Replication Proteins from Prokaryotes and Eukaryotes**

| Stage          | Function              | <i>E. coli</i>     | I phage        | T4 phage                 | Mammalian (SV40)                      | <i>S. cerevisiae</i>         |
|----------------|-----------------------|--------------------|----------------|--------------------------|---------------------------------------|------------------------------|
| Pre-initiation |                       | DnaA               | O protein      |                          |                                       |                              |
|                | Initiator             |                    |                | ?                        | ORC? (Tag)                            | ORC (1 ~ 6)                  |
|                |                       | DnaC               | P protein      |                          | MCM?                                  | MCM (2 ~ 7)                  |
|                | Activator             |                    | RNA polymerase |                          | CDK/cyclin?                           | CDK/cyclin?                  |
| Initiation     | Replication           | DnaB               |                | Gene 41                  | CDC6?<br>CDC7/DBF4?<br>MCM? (Tag)     | DCD6?<br>CDC7/DBF4?<br>DNA2? |
|                | DNA helicase          |                    |                |                          |                                       |                              |
|                | ssDNA binding protein |                    | ssb            | Gene 32                  | RPA (replication protein A)           | yRPA (RFA1, 2, 3)            |
|                | Primase               |                    | DnaG           | Gene 61                  | DNA polymerase a/primase              | primase (PRI1, PRI2)         |
| Elongation     | DNA polymerase        | DNA polymerase III | Gene 43        | DNA polymerase a/primase | DNA polymerase I (pola: POL1 = CDC17) |                              |

|                         |           |            |   |   |
|-------------------------|-----------|------------|---|---|
|                         |           |            | DNA polymerase $\alpha$                   | DNA polymerase II (pol $\epsilon$ : POL2)       |
|                         |           |            | DNA polymerase $\delta$                   | DNA polymerase III (pol $\delta$ : POL3 = CDC2) |
| DNA polymerase          | subunit   | Gene 45    | PCNA (proliferating cell nuclear antigen) | yPCNA (POL30)                                   |
| Accessory proteins      |           |            |   |   |
|                         | g complex | Gene 44/62 | RFC (replication factor C)                | yRFC (CDC44 = RFC1, RFC2, RFC3, RFC4, RFC5)     |
| Swivel                  | Gyrase    | T4 Topo    | Topoisomerase I                           | TOP 1   |
|                         | TopA      |            | Topoisomerase II                          | TOP 2   |
| Maturation Decatenation | Gyrase    |            | Topoisomerase II                          | TOP 2   |

---

Some proteins that have roles in other cellular processes function as replication proteins at specific stages in DNA replication; eg, **RNA polymerase** is necessary to activate several replication origins or to synthesize [primer](#) RNA; a **heat-shock** protein, *grpE*, is required to initiate [lambda phage](#) DNA replication, and [thioredoxin](#) is involved in T7 phage DNA polymerase.

Replication of T4 phage DNA is dependent primarily on its own encoded proteins, and a search for its replication proteins was carried out by the **conditional mutant** technique. As a result, seven replication proteins essential for the DNA synthesis, and many proteins with related functions, were identified.

Reconstitution of these replication reactions revealed that the replication factors do not react with a [template](#) DNA individually, but as a multimeric protein complex. Furthermore, the configurations and functions of the complex alter successively, from pre-initiation to elongation stages. For example, an assembly of [DNA helicase](#) and [primase](#) (and the assembly factors in some cases), called the *primosome*, is formed at the replication origins (or the assembly sites); subsequently, the addition of two DNA polymerase subunits and their accessory proteins in this complex generates a *replisome*, a super-protein complex that synthesizes **leading and lagging** DNA strands coordinately at a [replication fork](#) (2).

Due to the difficulty of applying **genetics** to most eukaryotes, the identification of eukaryotic replication proteins was started by a search of activities known to be involved in replication, such as DNA polymerases, topoisomerases, and DNA ligases. In contrast, a unicellular eukaryote, the **yeast** *Saccharomyces cerevisiae*, is exceptionally adaptable for genetic analysis, and many genes necessary for the cell cycle to progress have been identified. However, it was difficult to identify their biochemical activities and distinguish their requirement in DNA replication, due to the absence of any *in vitro* replication systems. The first breakthrough in studies of eukaryotic replication proteins occurred with the development the [SV40 virus](#) *in vitro* DNA replication system using human cell lysates (for a review, see (3)). Replication of this virus largely relies on the replication functions of the cell, so fractionation could identify the required proteins; this has isolated several replication proteins also involved in cellular chromosomal DNA replication. Furthermore, because these replication proteins are highly conserved throughout eukaryotes, several yeast replication genes were identified by [homology](#) searches (Table 1). Using previously and newly identified replication proteins, the process of SV40 DNA replication was totally reconstituted, and the functions of the proteins in DNA replication have been well elucidated (4). The replication reaction can be roughly divided into four stages: pre-initiation, initiation, elongation, and segregation (maturation), and specific sets of replication proteins are required to process these stages. The SV40 replication system primarily reproduced the cellular replication reaction from elongation to maturation but, in contrast, insufficient information about cellular initiation proteins was obtained from this viral system. One characteristic feature of eukaryotic DNA elongation obtained from the analysis is the involvement of multiple DNA polymerases in one replication fork.

Another breakthrough was the discovery of the [origin recognition complex](#) (ORC) from yeast (5) and [licensing factor](#) (MCM) from *Xenopus* (6). ORC specifically binds to the yeast replication origin sequence in an ATP-dependent manner and is a strong candidate for the initiator protein. MCM is a key player in the licensing reaction, and yeast MCM forms a pre-replicative complex at the origin, together with ORC. Through studies on their functions, the understanding of yeast chromosomal DNA replication has progressed drastically. Furthermore, these proteins are also highly conserved among eukaryotes, implying the existence of a common mechanism to initiate eukaryotic chromosomal DNA replication.

Another specific feature of eukaryotic replication proteins is that the replication process is tightly linked with several eukaryote-specific processes, such as [cell-cycle](#) control and [chromatin](#) assembly, so the boundary between replication and these processes is ambiguous. Therefore, [cyclin](#)-dependent [kinases](#), licensing factors, and chromatin assembly factors are part of replication factors in eukaryotes and, indeed, some of them form assemblies with major replication proteins at [replication foci](#) during **S phase**.

#### Bibliography

1. T. A. Baker and A. Kornberg (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York, pp. 478–483.
2. K. J. Marians (1992) *Ann. Rev. Biochem.* **61**, 673–719.
3. B. Stillman (1989) *Ann. Rev. Cell Biol.* **5**, 197–245.
4. S. Waga and B. Stillman (1994) *Nature* **359**, 207–212.
5. S. P. Bell and B. Stillman (1992) *Nature* **357**, 128–134.
6. Y. Kubota et al. (1995) *Cell* **81**, 601–609.

#### DNA Sequencing

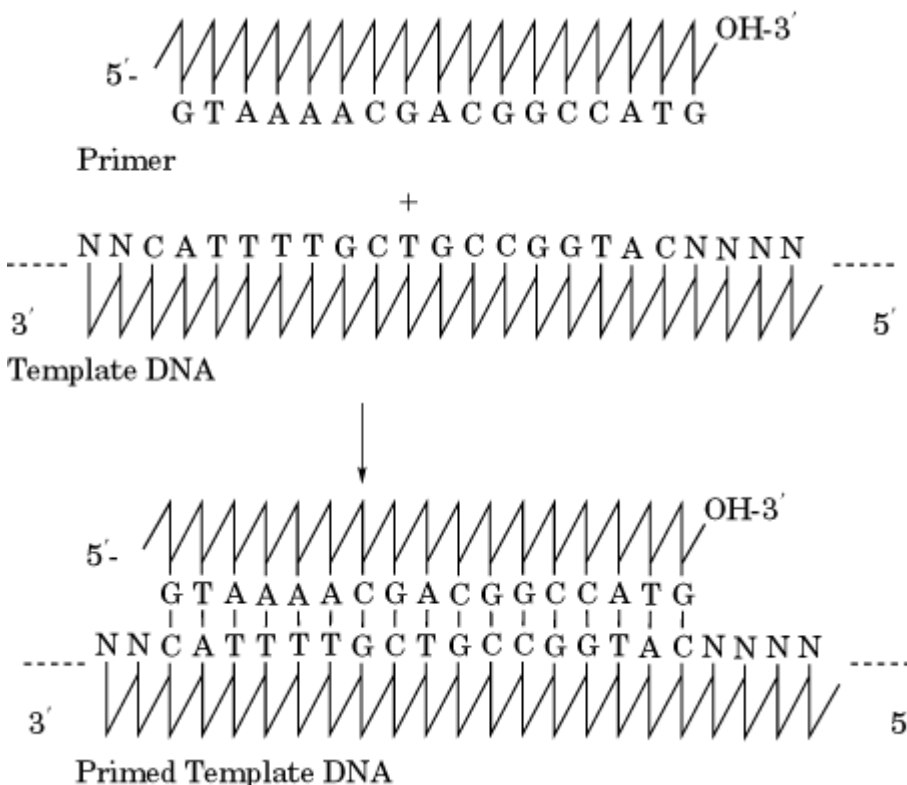
One of the fundamental methods of molecular biology is determining the sequences of bases in specific segments of DNA. The sequence information can be used to help deduce the function of the DNA segment, map its chromosomal location, determine how it might be involved in regulating gene expression or replication, or elucidate how it interacts with proteins. For example, the sequence of an unknown [complementary DNA](#) (cDNA) can be used to deduce the sequence of the protein encoded by its parent [messenger RNA](#). In turn, this protein sequence can be compared with the sequences of all known proteins, giving clues to function and evolutionary origin (see [Sequence Analysis](#)).

### 1. Chain Termination, Sanger Method

The [chain-termination](#) method of DNA sequencing was first described by Sanger in 1977. This method involves synthesizing a DNA strand by a **DNA polymerase** *in vitro*. Synthesis is initiated at only one site, where a primer anneals to the template. The growing chain is terminated by incorporating a 2',3'-dideoxynucleoside triphosphate (ddNTP) that does not support continued DNA synthesis (hence the name chain termination).

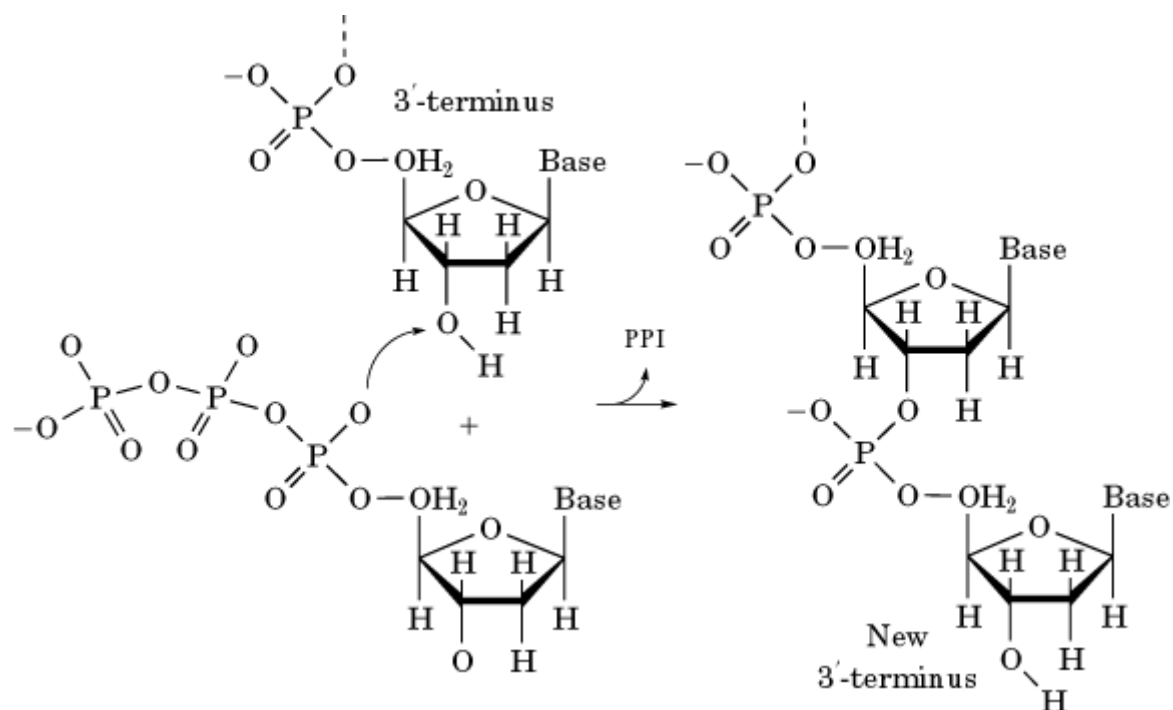
DNA polymerases initiate synthesis only at the 3'-end of a primer annealed to a DNA template. For most sequencing applications, the primer is a short synthetic oligonucleotide (18 to 35 nucleotides long) that is complementary in sequence to the template at a unique position adjacent to the region to be sequenced. The primer is hybridized to the template at the appropriate temperature (Fig. 1). Once this duplex is formed, the primer is extended by the DNA polymerase in the presence of the four deoxynucleoside triphosphates (dGTP, dATP, dTTP, and dCTP; dNTP corresponds to any one of the four).

**Figure 1.** Annealing of primer to template DNA.

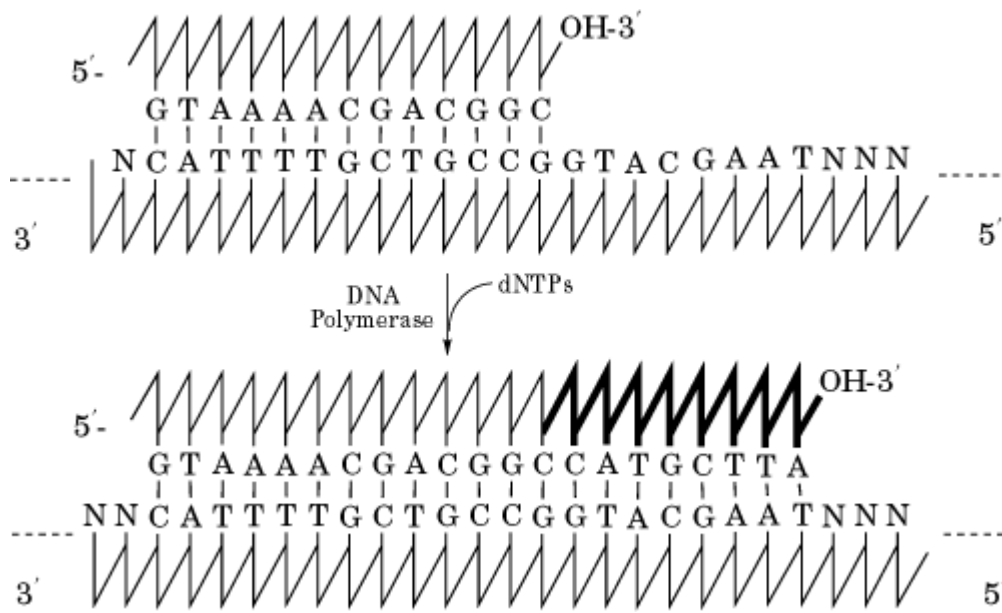


Polymerization of the new strand requires a free 3'-hydroxyl group (Fig. 2). As long as dNTPs are incorporated, there is a 3' hydroxyl group available for continued polymerization of the growing chain (Fig. 3). However, dideoxynucleoside triphosphates (ddNTPs) lack a 3'-hydroxyl group and terminate chain elongation when incorporated into the DNA. Consequently, synthesis, directed by a sample of template DNA that has a unique primer, in the presence of all four deoxy- and one dideoxynucleoside triphosphate yields a population of molecules that have a common 5'-end, plus 3'-ends all of which have the same terminal dideoxynucleotide base, but have a distribution of sizes depending on the site at which the ddNTP is incorporated (Fig. 4). Thus, the size of each fragment is determined by the sequence of the template. Typically, four separate reactions are performed, each with a different ddNTP. The products of the four reactions are analyzed by [electrophoresis](#) using a denaturing polyacrylamide gel, which accurately separates the products by size. Because the size of each fragment is determined by the template sequence, this sequence can be determined from the order of the bands on the gel.

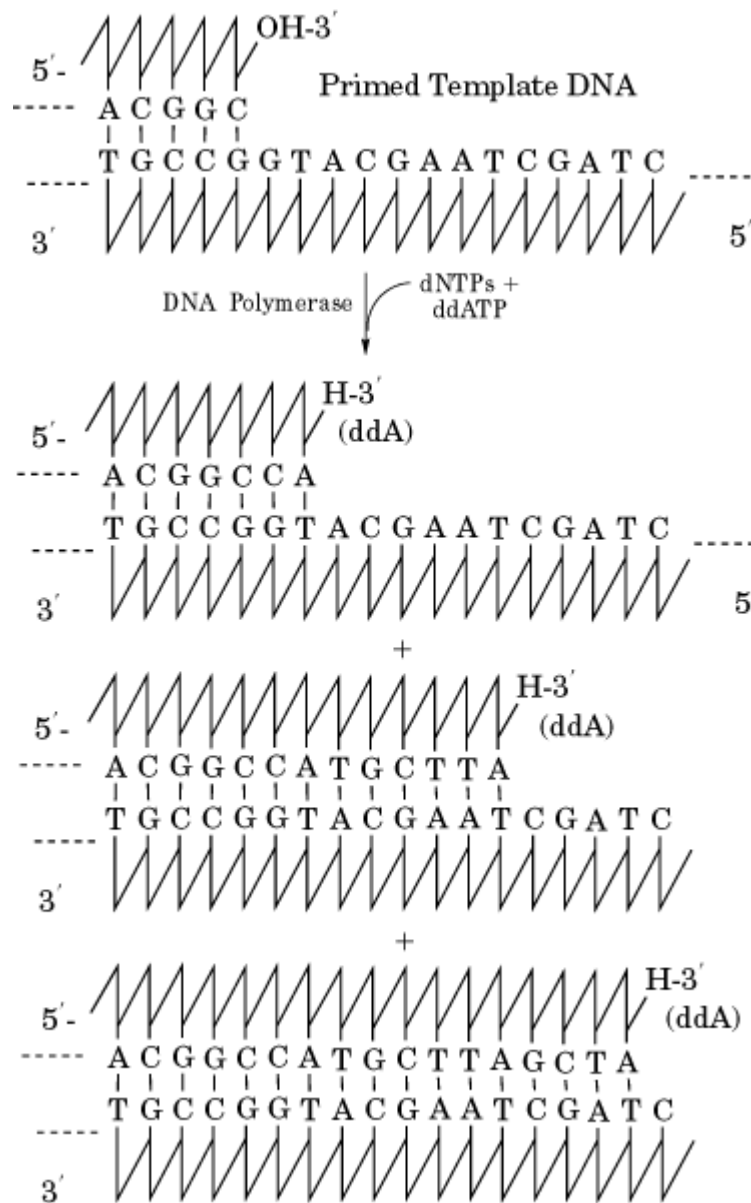
**Figure 2.** Biochemistry of chain elongation. In the presence of dNTPs, DNA polymerase catalyzes the condensation of deoxynucleoside triphosphate at the 3'-end of a primed template, releasing pyrophosphate.



**Figure 3.** Chain elongation by DNA polymerase. As long as dNTPs are incorporated, there is a free 3'-hydroxyl group available for continued polymerization of the growing chain.



**Figure 4.** Chain termination. Synthesis from a unique primer in the presence of all four deoxy- and one dideoxynucleoside triphosphate yields a population of molecules with common 5' ends, but different 3'-ends, depending on the site at which a ddNTP is incorporated.



## 2. Chain Cleavage, Maxam–Gilbert Method

Another general sequencing method, known as the chain-cleavage or Maxam–Gilbert method, has also been used extensively. This method works by a similar method of mapping DNA sequence to DNA size, but does so by degrading existing DNA chains, rather than synthesizing new ones. A sample of purified DNA to be sequenced is first labeled at one end. Then the DNA is subjected to a chemical treatment that breaks each DNA molecule at random, but only at places where one (or a defined subset) of the bases occurs. The result is a population of labeled molecules whose sizes are determined by the sequence. Determination of the sizes on samples cleaved with several sequence-specific cleavage treatments yields complete sequence information. Unfortunately, it is difficult to find chemical cleavage conditions that consistently give unambiguous, clean cleavage products, and few labeling methods can be used that provide labels stable enough for these treatments. For these and other technical reasons, chain-cleavage methods are rarely used today.

### Suggestions for Further Reading

F. Sanger, S. Nicklen, and A. R. Coulson (1977) DNA sequencing with chain-terminating inhibitors, *Proc. Natl. Acad. Sci. USA.* **74**(12), 5463–5467.



A. M. Maxam and W. Gilbert (1977) A new method for sequencing DNA, Proc. Natl. Acad. Sci. USA **74**(2), 560–564.

A. M. Maxam and W. Gilbert (1980) Sequencing end-labeled DNA routinely with base-specific chemical cleavages, Methods Enzymol. **65**, 499–560.

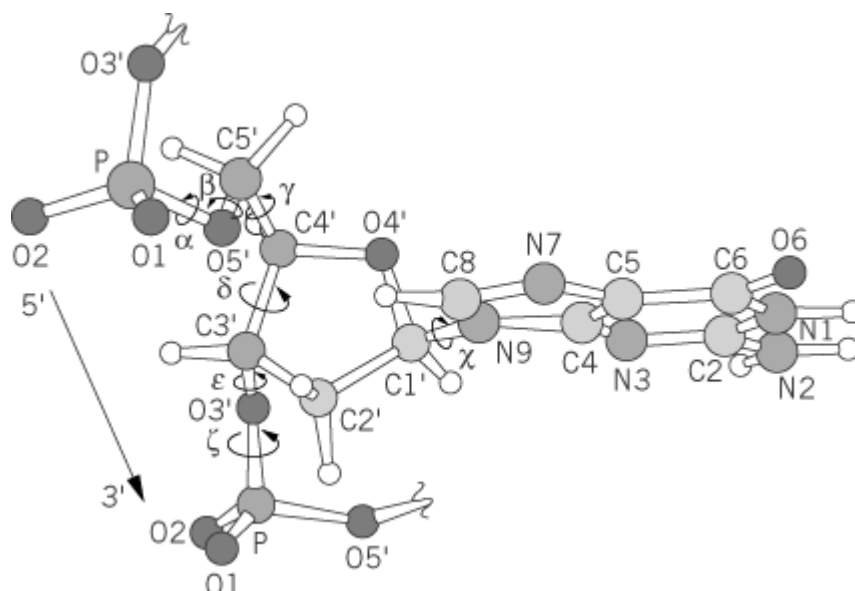
## DNA Structure

The central roles of deoxyribonucleic acid (DNA) in life, especially those associated with the storage and transfer of genetic information, are now well established. The discovery of the structure of DNA double helix by Watson and Crick nearly half a century ago ushered the modern biology into a new era and forever changed the landscape of biology. It became clear that the three-dimensional structure of DNA is intimately associated with its function. Recently, the rapid advances in the ability to determine the structures of biological [macromolecules](#) has produced a great wealth of information. In the field related to DNA structure, we begin not only to understand the subtle, yet important, sequence-dependent conformation associated with the canonical Watson–Crick DNA double helix, but also to discover new structural forms of DNA and to visualize how DNA interacts with important biomolecules, such as [proteins](#) and drugs. A better understanding of the structural basis for the function of DNA will facilitate the rational design of compounds that contribute to improve the quality of life. The out-pouring of the DNA sequence information through various [genome](#) projects has made the effort of getting structural information associated with DNA an urgent one.

### 1. Chemical Structure

DNA is a biological polymer made of the deoxynucleotide building blocks (Fig. 1). [Nucleotides, nucleosides, and nucleobases](#) are discussed in that entry. The polynucleotide chain has a directionality due to the specific internucleotide phosphodiester linkage between the O3' and O5' atoms from two neighboring deoxyriboses. By convention, a polynucleotide chain is described as going from the 5' to the 3' direction. DNA oligonucleotides can now be synthesized routinely for a wide range of applications, including structural studies (see [DNA Synthesis](#)). For convenience, DNA oligonucleotides are denoted with their sequence in a form such as d(CGCGAATTCGCG).

**Figure 1.** The torsion angles ( $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$ ,  $\zeta$ ) of a polynucleotide sugar-phosphate backbone and the glycosyl torsion angle  $\chi$  are shown. The diagram shows a fragment of B-DNA; therefore the deoxyribose conformation is C2'-*endo*, the glycosyl angle  $\chi$  is *anti*, and the  $\alpha/\zeta$  combination is *gauche*<sup>-</sup>/*gauche*<sup>-</sup>.



DNA is a polyelectrolyte because of the negative charges associated with the phosphate groups. The negatively charged DNA is neutralized by the positive charges of metal ions, polyamines, or proteins. Metal ions, such as sodium, potassium, or magnesium ions, are used in the screening of the DNA negative charges by interactions with the phosphate oxygen atoms, in enzyme reactions (eg,  $\text{Mg}^{2+}$  in **DNA polymerase's** function), or in the folding of more complex structures, such the [guanine quartet](#) structure.

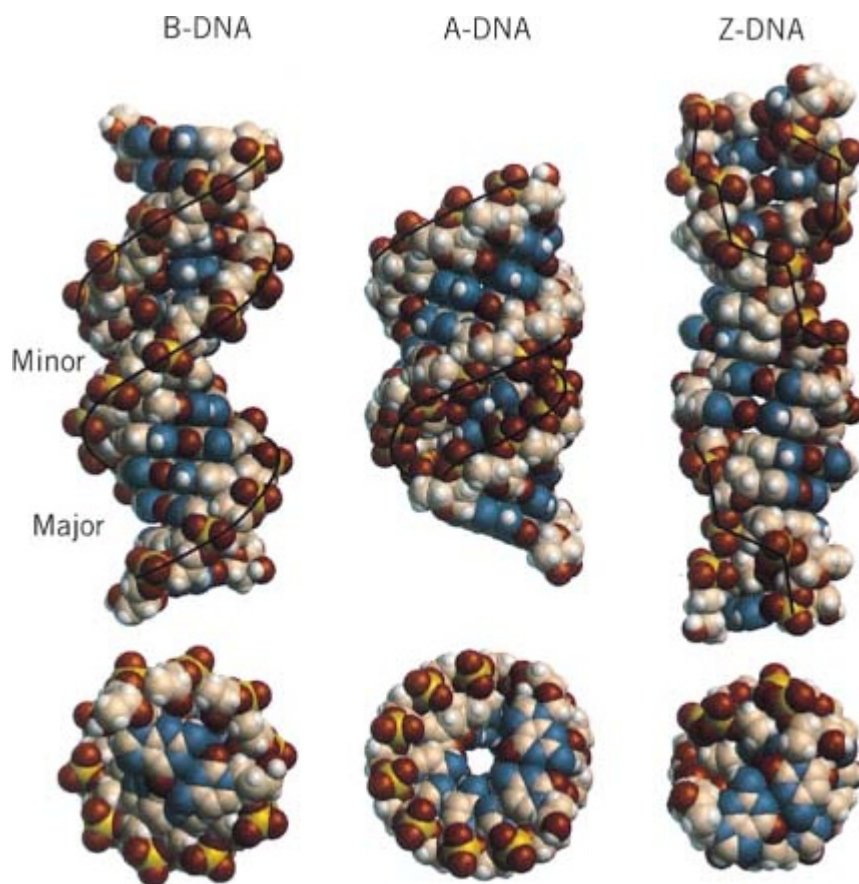
The sugar–phosphate polynucleotide chain consists of many single bonds around which the attached atoms may rotate. The definition of the [torsion angles](#) of the DNA backbone is shown in Figure 1. The deoxyribose ring is nonplanar, and adopts two preferred “puckers,” namely, *C2'-endo* and *C3'-endo* conformations. Two parameters, the sugar pseudorotation phase angle and its maximum torsion angle (ie, pseudorotation amplitude), are normally used to describe the puckering modes available in terms of the sugar torsion angles. Two preferred orientations (*anti* or *syn*) of the base with respect to the sugar (as defined by the torsion angle about the  $\text{C1}'\text{—N}$  glycosyl bond  $\epsilon$ ) are found.

DNA predominantly exists as an antiparallel double-stranded helix. The two strands are helically coiled, which maximizes the exposure of the negatively charged sugar–phosphate backbone to water and shields the **hydrophobic** aromatic bases in the middle from [water](#). In the Watson–Crick base pairs of guanine with cytosine (G:C) and adenine with thymine (A:T), the specificity of base pairing is provided by [hydrogen bonds](#). It should be noted that other types of base pairs are also known to exist, and they often are involved in unusual DNA structures.

## 2. DNA Double Helices

The polymorphism of DNA associated with different helical forms, under the influence of different environmental conditions, has begun to emerge in recent years. Much of the detailed structural information of DNA has been obtained through the high resolution [X-ray crystallography](#) analysis and nuclear magnetic resonance ([NMR](#)) spectroscopic analysis of DNA oligonucleotides with defined sequences. A survey of the Nucleic Acids Database (site currently unavailable) and their complexes with proteins and drugs indicates that hundreds of crystal structures are available. Additional structures solved by NMR can also be found in the Brookhaven Protein Database (<http://www.pdb.bnl.gov/>). The structural features of the three common types of DNA helices—namely, the [B-DNA](#), [A-DNA](#), and [Z-DNA](#)—are illustrated in Figure 2, detailed in Table 1, and described in more detail in the individual entries.

**Figure 2.** End view and side view of the B-DNA, A-DNA, and Z-DNA helices. Note the decrease in diameter, as well as the relative positions of the base pairs and backbone to the helix axis. See color insert.



**Table 1. Structural parameters of B-DNA, A-DNA, and Z-DNA**

| Structural Parameter         | B-DNA Right-handed | A-DNA Right-handed | Z-DNA Left-handed                              |
|------------------------------|--------------------|--------------------|--|
| Repeat unit (bp)             | 1                  | 1                  | 2  |
| Base pair/turn (degrees)     | 10.4               | 11                 | 12   |
| Tilt (degrees)               | ~ 0                | 19                 | -9   |
| Rise per base pair (Å)       | 3.3                | 2.3                | 3.7  |
| Helical pitch (Å)            | 34                 | 25.4               | 45   |
| Glycosyl angle               | <i>anti</i>        | <i>anti</i>        | <i>anti</i> at C/ <i>syn</i> at G              |
| Sugar pucker                 | C2'- <i>endo</i>   | C3'- <i>endo</i>   | C2'- <i>endo</i> at C<br>C3'- <i>endo</i> at G |
| Phosphate conformation (a/z) | -40/-98°           | -88° /-44°         | -146° /80° at C                                |

|                                   |                 |                  |                 |
|-----------------------------------|-----------------|------------------|-----------------|
|                                   |                 |                  | 60° / -58° at G |
| Helical diameter (Å) <sup>a</sup> | ~ 20            | ~ 25             | ~ 18            |
| Major groove                      | Wide and deep   | Narrow and deep  | Flattened       |
| Minor groove                      | Narrow and deep | Wide and shallow | Narrow and deep |

---

<sup>a</sup> Conversion factor for angstroms to meters is  $1.0 \times 10^{-10}$ .

### 3. Hydration

The [hydration](#) environment around a DNA double helix plays an important role in determining the type of conformation adopted and in determining other properties. Recent structural work suggested that some proteins, such as the Trp repressor (see [TRP Operon](#)) and the EcoRI [restriction enzyme](#), recognize DNA sequences through direct hydrogen bonds, [nonpolar](#) contacts, indirect structural effects, and, surprisingly, water-mediated interactions. Thus specific water molecules play critical roles in the sequence-specific recognition by proteins. More recently, it was found that water molecules play a different role in that they modulate the binding of sequence-nonspecific DNA-binding proteins (eg, Sac7d) to DNA of random sequence ([1](#)). Thus it is now generally accepted that the hydration shell surrounding the DNA molecule plays an important role in DNA recognition by proteins and other ligands, such as DNA-binding anticancer drugs.

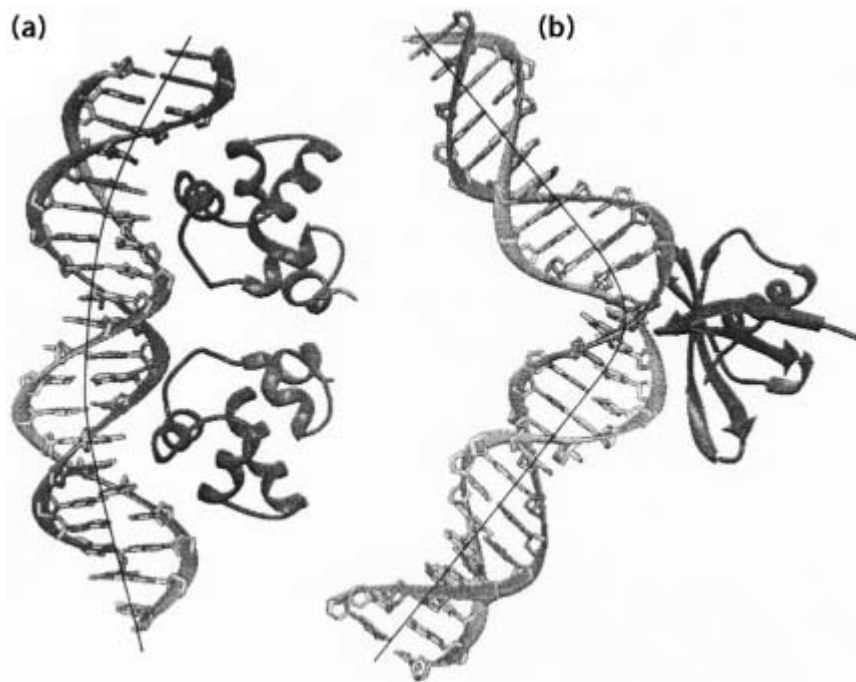
### 4. Novel DNA Structures

Certain sequences of DNA form three-dimensional structures that are not the common A-, B- or Z-DNA double helices. Those sequences form higher order structures such as a hairpin loop, **triple-stranded** structures, tetrastranded structures, and [cruciform](#). They often involve nonstandard base pairs. For example, the self-pairing of guanine bases are found in the tetrastranded [guanine-quartet](#) structure. Specific functions have already been identified for some of these higher order structures. The known multistranded structures include the [Guanine Quartet](#), the [I-Motif](#), the triple helix, and the [Cruciform \(Holliday Junction\)](#), which are described in the individual entries. Some of the novel double-stranded structures are discussed here.

#### 4.1. Bend DNA

Certain DNA sequences have abnormal mobilities in [gel electrophoresis](#). For example, DNA fragments having repeats of  $(A)_n$  nucleotides (with  $n \geq 4$ ), separated by another four to five nucleotides and phased with the helical repeat, migrate in the gel significantly more slowly than do those having mixed sequences. It was discovered that the 5'-AAAAA sequence has an intrinsic bending property that can be demonstrated by the increased efficiency of cyclization of those bend DNA fragments. The molecular basis of the intrinsic bendability of the  $(A)_n$  sequence has been investigated. The high propeller twist associated with the A-T base pair in the  $(A)_n:(T)_n$  sequence may play an important role in this property. The bend in those sequences is relatively smooth, resulting in a curved DNA structure. The opening (roll) of the bend is toward the minor groove. Many proteins induce such a smooth bending, by having many small single-step bends, which has been observed in the crystal structures of a number of protein-DNA complexes, exemplified by the structure of the 434 repressor-DNA complex (Protein Database accession number PDR015) shown in Figure [3a](#).

**Figure 3.** Two types of DNA bending modes. (a) Smooth bending found in the 434 repressor–DNA complex (Protein Database PDR015). (b) Sharp kink found in the Sac7d–DNA complex (Protein Database 1AZP).



Another type of bend DNA is found to have a sharp kink at a localized site, usually caused and stabilized by [DNA-binding proteins](#) bound to DNA. A relevant example is found in a recent structure of the complex between a chromosomal 7-kDa protein Sac7d from the hyperthermophilic archaeobacterium *Sulfolobus acidocaldarius* and DNA oligonucleotides (1). The DNA is kinked sharply at the C<sub>2</sub>pG<sub>3</sub> step in the Sac7d-GCGATCGC complex (Figure 3 b). The sharp kink is caused by the intercalation of the side chains of amino acid residues Val26 and Met29 of the Sac7d protein into DNA base pairs from the minor groove direction, widening the minor groove at this step.

This type of sharp DNA kink has been observed in the complexes of [TATA-box](#) binding protein (TBP), and two **HMG**-box containing proteins, LEF-1 and SRY, with their cognate-specific DNA sequences. Remarkably, both Sac7d and TBP use amino acids on a [beta-sheet](#) for the intercalation, whereas in LEF-1 and SRY use the amino acids located at the corner of the **alpha-helical** L-shaped HMG box for the intercalation. A more thorough discussion of the structural basis of the various types of DNA bending has appeared recently (2).

#### 4.2. Triplet-Repeat Sequences

Recently a number of human genetic diseases have been correlated with expansions of triplet repeats of the DNA sequence (CNG)<sub>n</sub>. The (CGG)<sub>n</sub> repeat in the [X-chromosome](#) is responsible for fragile-X syndrome, the (CAG)<sub>n</sub> repeat is associated with Huntington's disease and spinobulbar muscular atrophy, and finally the (CTG)<sub>n</sub> repeat is associated with myotonic dystrophy. How these unusual repetitive sequences correlate with the etiology of these diseases, and the mechanism by which those repeats are expanded during [DNA replication](#), are under intense scrutiny. Some have proposed that a “slippage” process occurs because of the ease of the formation of hairpin structures for these repeating sequences. Certain triplet repeats, such as (CAG)<sub>n</sub> and (CTG)<sub>n</sub>, but not (CGA)<sub>n</sub>, have a strong propensity to form hairpin structures. Therefore, a DNA duplex encoding the (CAG)<sub>n</sub>:(CTG)<sub>n</sub> repeats may easily exchange between duplex and [cruciform](#), especially under negative **supercoiling** strain. If there are proteins or other ligands (eg, drugs) that can stabilize the stem of the

cruciform, this process would be inhibited. The three-dimensional structures of several DNA oligonucleotides associated with those novel triplet repeats have been studied, primarily by NMR. The  $(CAG)_n$  repeat can form a stable duplex structure incorporating “sheared” G-A mismatched base pairs. The duplex structure associated with the  $(CCG)_n$  repeat appears to have the C nucleotides extruded from the helix. The structural and functional studies of those unusual repeats remain very active (3). (See also [Trinucleotide Repeats](#).)

#### 4.3. Parallel-Stranded DNA Duplex

A new addition to non-canonical DNA structures is the parallel-stranded (PS) DNA structures. The question of whether a stable DNA duplex can be parallel has been addressed previously (4). A series of A,T-containing DNA sequences was designed to form parallel duplexes using reversed Watson–Crick base pairs. The stability of those PS duplexes is modest; for example, the  $T_m$  (**melting temperature**) of a 21-mer PS-duplex is 15°C lower than that of the corresponding antiparallel duplex. A different motif was the non-Watson–Crick homo base-paired parallel-stranded DNA, called P-DNA. It was demonstrated that the d(CGA) sequence has a strong propensity to form the P-DNA structure (5).

An important requirement for a hetero base-paired parallel duplex is that the two glycosyl bonds within a base pair have to come from opposite directions, because of the identical chain polarity. For the normal nucleic acid bases, this can be accomplished using the reverse Watson–Crick base-pair conformation. However, A-T and G-C base pairs in a reverse Watson–Crick conformation are not isostructural, due to their hydrogen-bonding restrictions. Therefore, it has not been easy to design a stable PS duplex in which all four bases can be incorporated in random order.

This difficulty has been overcome by using alternative nucleosides, 2'-deoxyisoguanosine (iG) and 2'-deoxy-5-methyl-isocytosine (iC), which can form stable reverse Watson–Crick base pairs with the normal 2'-deoxycytosine (C) and 2'-deoxyguanosine (G), respectively. Indeed, oligodeoxynucleotides containing iG and iC can form remarkably stable parallel-stranded duplexes with the complementary (G,C)-containing DNA or RNA strands (6). The ability of (iG,iC)-containing DNA oligomers to form specific stable parallel-stranded duplexes may offer new opportunities for designing useful probes for applications such as antisense or aptamer molecules.

## 5. Summary

In conclusion, the importance of DNA structure and dynamics in biology is very clear. The rapid advancement in genomics, molecular biology and structural biology (including synchrotron) offers an exciting future in the investigation of the protein and nucleic acid structures associated with new and significant biological functions. One can now attack problems that were unthinkable just a few years ago. The structure of the [nucleosome](#) has been determined at 2.8 Å resolution and the detailed DNA conformation has been presented (7). The structure of the [ribosome](#) particle is on its way to be elucidated at a resolution high enough to visualize individual proteins and RNA. It can be certain that many structures of novel DNA sequences and important protein–DNA complexes will be forthcoming at a rapid pace in the next few years.

## Bibliography

1. H. Robinson, Y.-G. Gao, B. S. McCrary, S. P. Edmondson, J. W. Shriver, and A. H.-J. Wang (1998) *Nature* **392**, 202–205.
2. R. E. Dickerson (1998) *Nucleic Acids Res.* **26**, 1906–1926, and references cited therein.
3. A. M. Gacy and C. T. McMurray (1998) *Biochemistry* **37**, 9426–9434, and references cited therein.
4. J. H. van de Sande et al. (1988) *Science* **241**, 551–557.
5. H. Robinson and A. H.-J. Wang (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5224–5228.

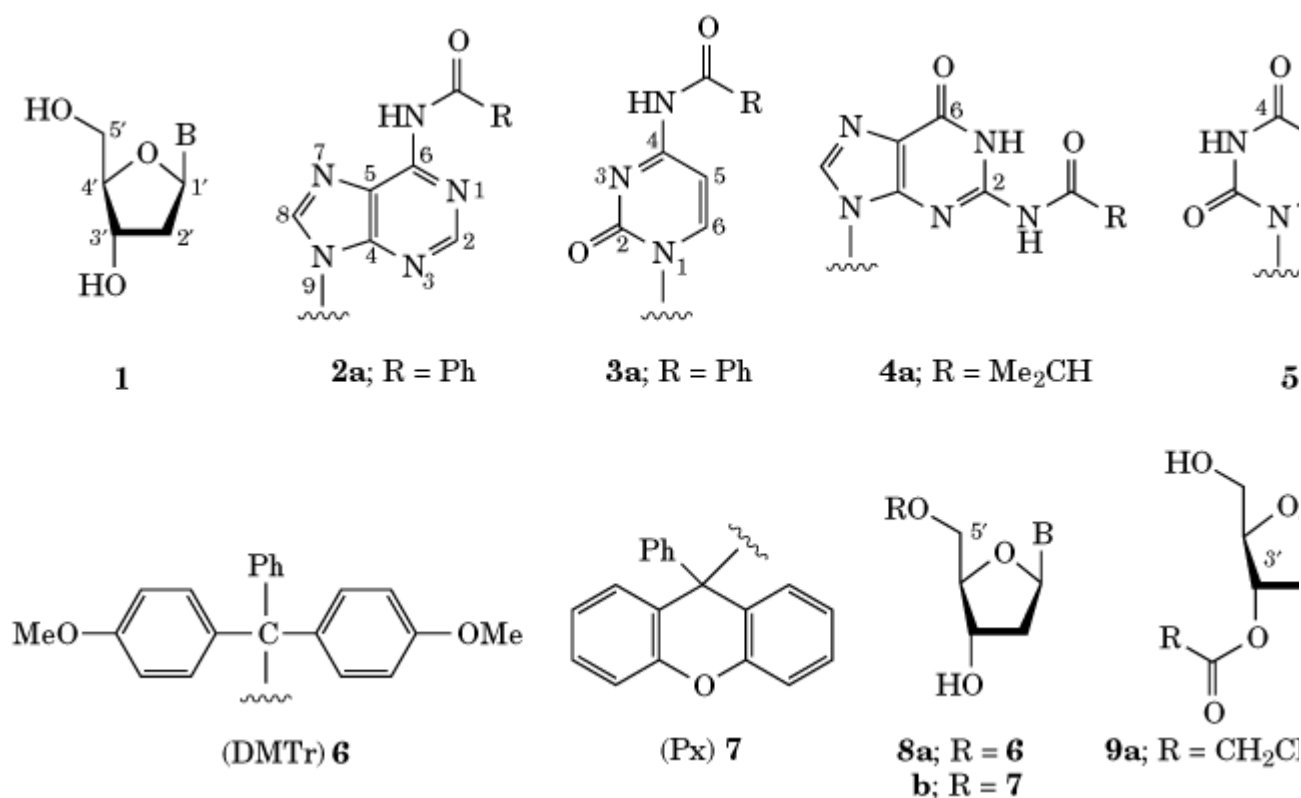
6. X.-L. Yang, H. Sugiyama, S. Ikeda, I. Saito, and A. H.-J. Wang (in press) *Biophys. J.*
7. K. Luger, A. W. Mader, R. K. Richmond, D. F. Sargent, and T. J. Richmond (1997) *Nature* **389**, 251–260.

### **Suggestions for Further Reading**

8. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer-Verlag, New York.
9. C. R. Calladine and H. R. Drew (1992) *Understanding DNA*. Academic Press, San Diego.
10. R. R. Sinden (1994) *DNA Structure and Function*, Academic Press, San Diego.
11. G. M. Blackburn and M. J. Gait, eds. (1996) *Nucleic Acids in Chemistry and Biology*, Oxford Univ. Press, New York.
12. S. M. Hecht, ed. (1996) *Bioorganic Chemistry: Nucleic Acids*, Oxford Univ. Press, New York.

### **DNA Synthesis**

The introduction of methods for the chemical synthesis of oligo- and poly- deoxyribonucleotides (DNA sequences) has had a very considerable effect on the development of molecular biology. This is clearly apparent from other sections of the Encyclopedia. The three most important factors to be taken into account in the chemical synthesis of DNA sequences are: (1) the choice of suitable protecting groups for the 2'-deoxyribonucleoside building blocks [1], (2) the development of phosphorylation procedures that are suitable for the introduction of the internucleotide linkages, and (3) the purification of the synthetic DNA sequences themselves. The choice of protecting groups is of crucial importance. The protecting groups selected should be easy to introduce; they should also remain completely intact throughout the assembly of the DNA sequences and be easily removable under conditions where the synthetic DNA is completely stable.



Structure 1. DNA synthesis.

## 1. Protecting Groups for 2'-Deoxynucleosides

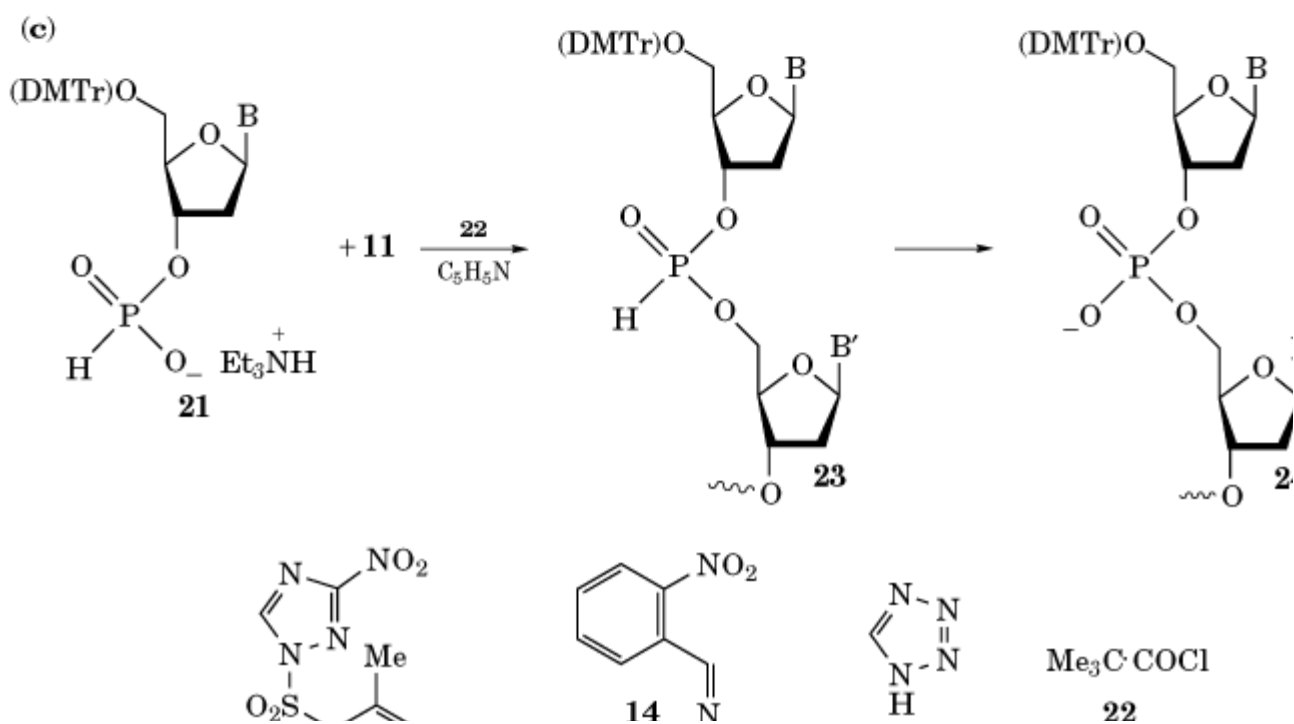
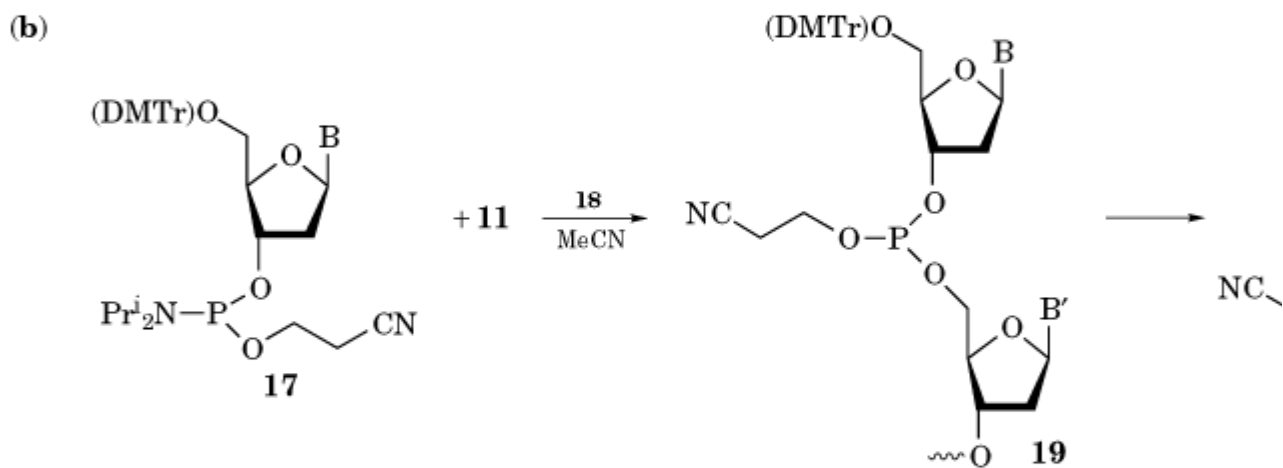
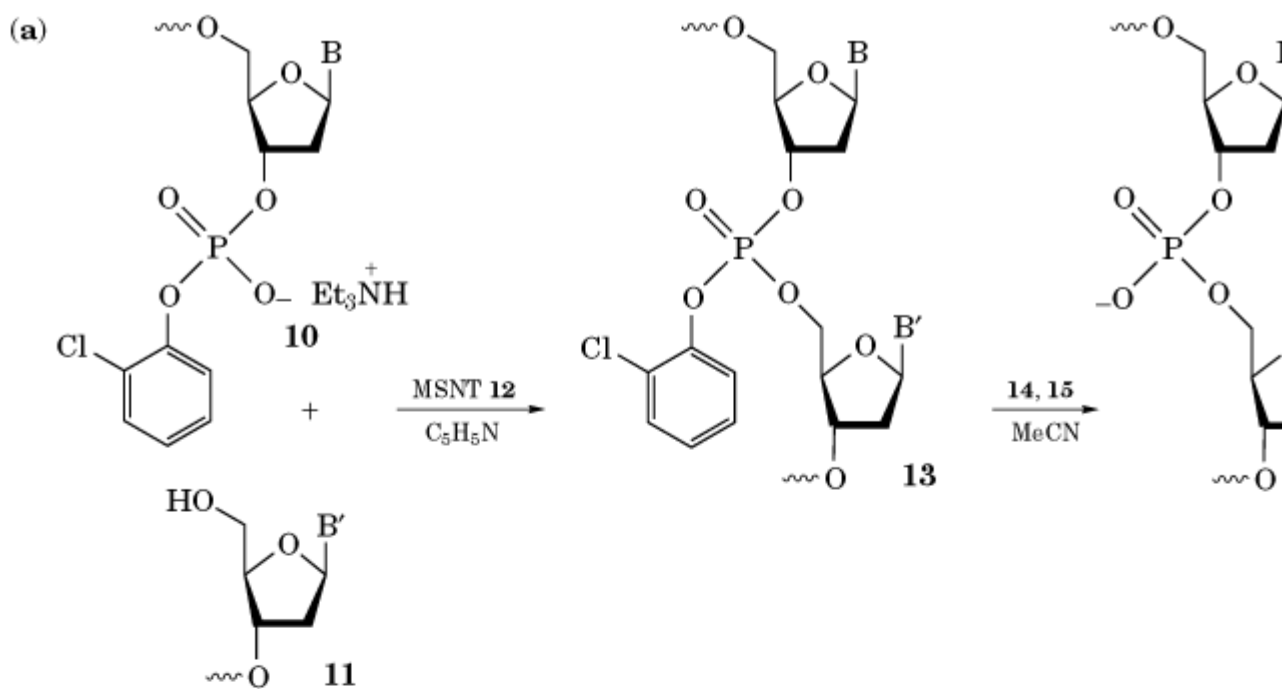
In order to avoid side reactions during phosphorylation, the adenine, cytosine, and guanine base residues are generally protected with acyl groups [e.g. as in **2a**, **3a**, and **4a**, respectively], but thymine residues [**5**] are usually left unprotected (**1**). It may sometimes be desirable (**2**) to protect thymine residues [**5**] on *O*-4 and *N*-2-acylguanine residues [**4**] also on *O*-6, but this depends on the phosphorylation procedure used. An important feature of DNA chemistry is that the glycosidic linkages, and especially those of purine deoxyribosides, are very sensitive to acid. Thus, the (di-*p*-anisyl)phenylmethyl (DMTr, [**6**]) and the 9-phenylxanthhen-9-yl (Px, [**7**]) groups, which are removable under relatively mild acidic conditions, are particularly useful for the protection of the 5'-hydroxy functions of 2'-deoxynucleosides [as in **8a** and **8b**, respectively]. 5'-O-Acyl-protected deoxynucleosides (**8**) R = acyl) have also been used as starting materials, particularly in the solution phase synthesis (see below) of DNA sequences. Acyl groups are also used to protect the 3'-hydroxy functions of deoxynucleosides [e.g. as in **9a**]. Both *N*-acyl and *O*-acyl protecting groups are usually removed at the end of the synthesis by ammonolysis under conditions where DNA is completely stable.

## 2. Phosphorylation Methods

There are essentially three phosphorylation methods (Fig. 1) that have been used successfully in the chemical synthesis of relatively high molecular weight DNA sequences.

Figure 1. Phosphorylation methods.





### 2.1. Phosphotriester Approach

The first method, which involves coupling a protected nucleoside or oligonucleotide 3'-(2-chlorophenyl) phosphate [10] with a protected nucleoside or oligonucleotide component with a free 5'-hydroxy function [11] (Fig. 1a), is usually referred to as the *phosphotriester approach* (3) as all of the internucleotide linkages are protected as phosphotriesters throughout the assembly of the target DNA sequences. In a now obsolete method known as the *phosphodiester approach*, the internucleotide linkages were left unprotected throughout the assembly of the DNA sequences. Coupling between the two components [10] and [11] is usually effected by a condensing agent such as 1-(mesitylene-2-sulfonyl)-3-nitro-1,2,4-(1*H*)-triazole (MSNT) [12] in anhydrous pyridine solution (4). The 3'-(2-chlorophenyl) phosphate component [10] is usually used in slight (*ca.* 20% to 25%) excess, and the condensing agent (e.g., MSNT), which also acts as a dehydrating agent, is used in greater excess. A very important feature of the phosphotriester approach is that the first unblocking step of the fully-assembled DNA sequence is effected by treatment (4) with an excess each of an appropriate oxime (such as (*E*)-2-nitrobenzaldoxime [14]) and  $N^1,N^1,N^3,N^3$ -tetramethylguanidine [15]. Such oximate treatment leads to complete removal of the 2-chlorophenyl protecting groups without any concomitant cleavage of the internucleotide linkages. *N*- and *O*-Acyl protecting groups [as in 2, 3, 4 and 9] are subsequently removed by treatment with concentrated aqueous ammonia, and the 5'-terminal acid-labile protecting group [as in 8a and 8b] is removed under mild conditions to leave the fully unprotected DNA sequence. The phosphotriester approach is an extremely versatile method for the synthesis of DNA sequences; it can be used both in solution phase (by the stepwise addition of mononucleotides or by the addition of oligonucleotide blocks) (4) and in **solid phase synthesis** (5).

### 2.2. H-Phosphoramidite Approach

The second phosphorylation method is generally based on protected nucleoside 3'-(2-cyanoethyl) *N,N*-di-isopropylphosphoramidites [17] and may be referred to as the *phosphoramidite approach* (6). In the presence of relatively weak acids (e.g., 1*H*-tetrazole [18]), such phosphoramidites [17] become powerful phosphitylating agents and react with protected nucleoside or oligonucleotide components with free 5'-hydroxy functions [11] (Fig. 1b) to give protected phosphite triesters [19]. These phosphite triesters are sensitive intermediates and are usually immediately oxidized (e.g., by treatment with iodine and 2,6-lutidine in aqueous tetrahydrofuran) to give the corresponding phosphotriesters [20]. The 2-cyanoethyl protecting groups are removed from the fully-assembled DNA sequences, together with the *N*- and *O*-acyl protecting groups, by treatment with concentrated aqueous ammonia. A crucial aspect of the phosphoramidite approach is that activated derivatives of *phosphorous* acid are much more reactive than corresponding activated derivatives of *phosphoric* acid. Therefore, in the phosphoramidite approach, coupling yields tend to be somewhat higher, and the coupling reactions are faster than they are in the phosphotriester approach. Largely for these reasons, the phosphoramidite approach has become the method of choice for the automated solid phase synthesis of DNA sequences (see text below); it has, therefore, been very widely used in the chemical synthesis of the relatively small quantities of the DNA sequences that are required in most molecular biological studies.

### 2.3. H-Phosphonate Approach

The third phosphorylation method is based on protected nucleoside 3'-(*H*-phosphonate) monoesters [21] as building blocks and is generally referred to as the *H-phosphonate approach* (78). Both *H*-phosphonate building blocks [21] and 2-chlorophenyl phosphates [e.g., 10] are very easy to prepare; they are also more stable and easier to handle than phosphoramidite building blocks [17]. However, partly perhaps because they are less bulky, *H*-phosphonates [21] undergo coupling reactions more rapidly than the corresponding 2-chlorophenyl phosphates. Thus, coupling between a nucleoside 3'-(*H*-phosphonate) building block [21] and a protected nucleoside or oligonucleotide with a free 5'-hydroxy function [11] (Fig. 1c) occurs very rapidly in the presence of a condensing agent such as


pivaloyl chloride [22] (8) in pyridine solution to give the corresponding H-phosphonate diester [23]. Such H-phosphonate diesters [23] are particularly sensitive to base-catalyzed hydrolysis, and for this reason almost all synthetic studies involving the H-phosphonate approach have been carried out in the solid phase (see text below). When a particular DNA sequence has been fully assembled, the H-phosphonate diester groups [as in 23] are oxidized (e.g., with iodine in the presence of a base in aqueous tetrahydrofuran) to give unprotected phosphodiester internucleotide linkages [as in 24] before the *N*- and *O*-acyl protecting groups are removed by treatment with concentrated aqueous ammonia. Despite its merits and its potential as an important synthetic method, the H-phosphonate has not been as widely used as the phosphoramidite approach in solid-phase synthesis. However, recent studies suggest that a modified H-phosphonate approach (see text below) may well become the method of choice in the solution-phase synthesis of DNA sequences. Indeed, there is no obvious reason why such a modified H-phosphonate approach could not also be applied successfully to solid-phase synthesis.

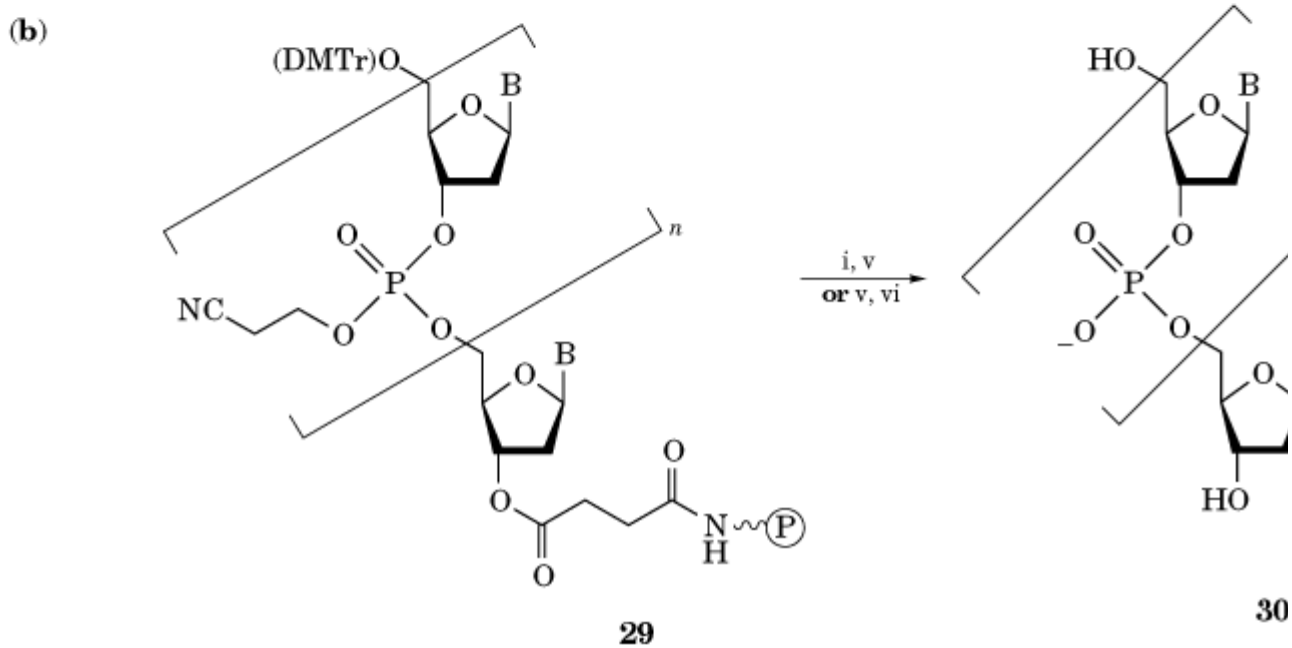
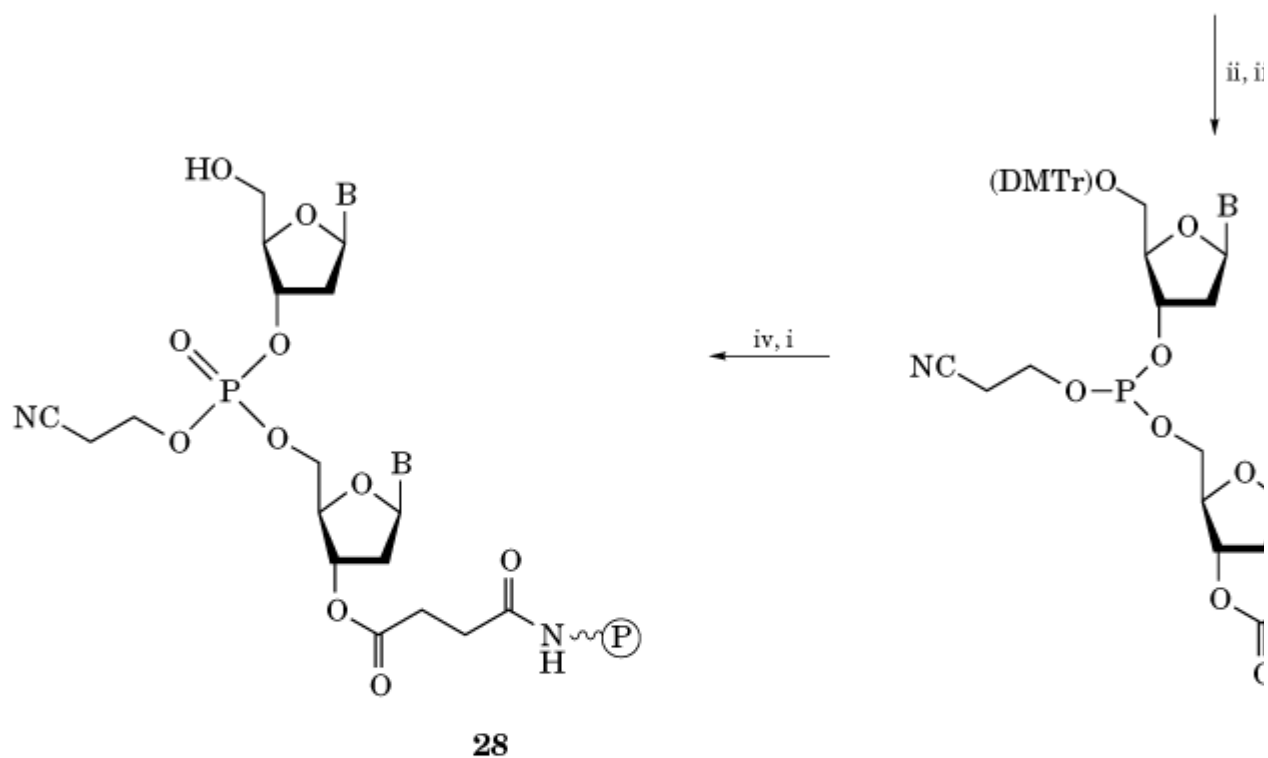
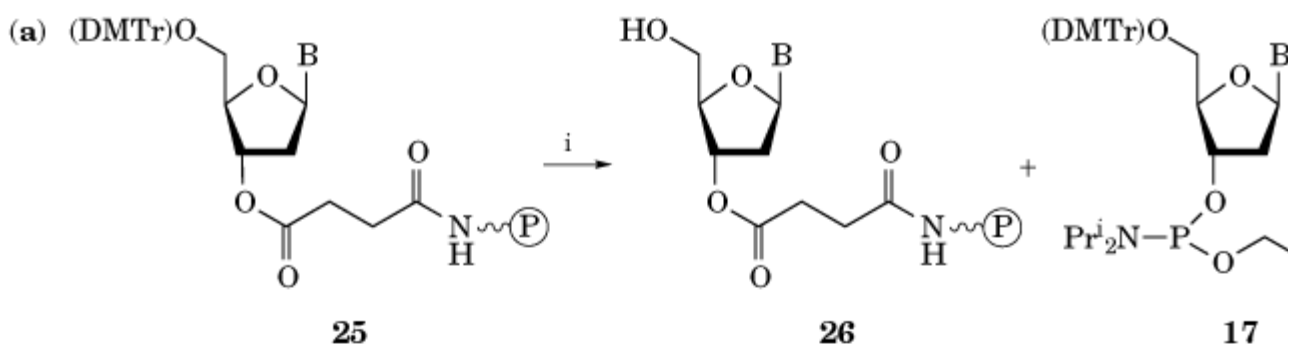
### 3. Solid-Phase DNA Synthesis

Relatively high molecular weight DNA sequences have been prepared successfully by the phosphotriester approach in solution by following essentially the procedure indicated in outline in Figure 1a. However, solution-phase synthesis is relatively laborious in that **chromatographic** purification steps are usually necessary after each coupling step. Nevertheless, if a very large quantity of a specific sequence is required (see text below), solution-phase synthesis may very well prove to be the method of choice. If, on the other hand, relatively small (i.e., milligram to gram) quantities of material are required for biological or biophysical studies, there is little doubt that solid-phase synthesis is to be preferred. While all three of the above phosphorylation methods (Fig. 1) have been used in solid-phase synthesis, the phosphoramidite approach (9) has emerged as the method of choice. This is mainly because its use leads to high coupling efficiencies and no significant side reactions. Furthermore, most commercial automatic synthesizers have been designed specifically to accommodate phosphoramidite chemistry. The main advantages of solid-phase synthesis, particularly by the phosphoramidite approach, are: (1) that it is very rapid and a DNA sequence containing, say, 50 nucleotide residues can easily be assembled and unblocked within one day; (2) only one purification step is required at the end of a synthesis as the growing DNA sequence is attached to a solid support (such as controlled pore glass [CPG] or polystyrene), and the excesses of all reagents are washed away; (3) all chemical reactions can be made to proceed in very high yield by using large excesses of reagents; and (4) the whole process may be fully automated in a DNA synthesizer. Solid-phase DNA synthesis has been developed to such an extent that the whole process can be carried out by a competent technician with no specialist knowledge of nucleotide chemistry. Automatic synthesizers, some of which are capable of assembling several different specific DNA sequences simultaneously, are readily available, and all the necessary building blocks [particularly phosphoramidites 17] and other reagents and solvents may be purchased in containers that are designed to be attached directly to the synthesizer.

The main steps in a synthetic cycle (i.e., the addition of one nucleotide residue) in solid-phase DNA synthesis by the phosphoramidite approach (10) are illustrated in outline in Figure 2a. Bottles containing solutions of each of the four main phosphoramidites [e.g. 17, B = 2a,3a,4a and 5] in acetonitrile and other reagents and solvents are attached to the appropriate ports on the synthesizer. The desired 3'-terminal nucleoside residue is generally attached to the solid support (P; usually CPG or polystyrene) through a succinoyl linker [as in 25], and this material is contained in a small column that is inserted in the synthesizer. The first step (Fig. 2a, step i) in the synthetic cycle involves treatment with acid to remove the 5'-*O*-DMTr protecting group [6] from the solid-supported 3'-terminal nucleoside residue [25]. The next step (step ii, the coupling step) is the *IH*-tetrazole-catalyzed reaction between a phosphoramidite [17] and the released 5'-hydroxy function [as in 26] to give a phosphite triester [27]. A large excess each of the appropriate phosphoramidite [17] and *IH*-tetrazole [18] are delivered automatically, leading to an anticipated coupling efficiency of *ca.* 98.5%. It is then important to "cap" any remaining 5'-hydroxy functions. This is effected by reaction with a very large excess of acetic anhydride in the presence of 1-methylimidazole and 2,6-lutidine in

tetrahydrofuran solution (step iii). “Capping” in each synthetic cycle ensures that, at the end of the synthesis, most of the truncated sequences will be of significantly lower molecular weight than the target DNA sequence, thereby facilitating its purification. The phosphite triester group [as in **27**] is then oxidized (step iv) to a phosphotriester group. The first synthetic cycle is now complete. A number of washing steps are involved in addition to steps i to iv, and the whole cycle usually requires less than 8 min. The removal of the 5'-*O*-DMTr protecting group [**6**] [to give **28**] is the first step of the second synthetic cycle, and this is followed by another coupling reaction involving the appropriate phosphoramidite building block [**17**]. It is essential that the acetonitrile solutions containing the phosphoramidites [**17**] and 1*H*-tetrazole [**18**] should be as dry as possible. The efficiency of the coupling steps may be estimated spectrophotometrically by measuring the quantities of the colored (di-*p*-anisyl)phenylmethyl (DMTr<sup>+</sup>) cations released at the beginning of each synthetic cycle. It is essential that the “deprotection” and coupling steps (steps i and ii, respectively) should proceed as nearly quantitatively as possible, and that “capping” (step iii) should be as efficient as possible. Clearly, oxidation (step iv) of the phosphite triester to the phosphotriester group should also go to completion in each cycle. If the average coupling efficiency is 98.5%, it should be possible to prepare DNA sequences containing 50 and 100 nucleotide residues in *ca.* 47% and 22% overall yields, respectively. However, if the average coupling efficiency is only 97.5%, then the overall yields of a 50-mer and a 100-mer would fall to *ca.* 29% and 8%, respectively. Therefore, if reasonable care is taken, it should be possible to prepare DNA sequences containing up to *ca.* 100 nucleotide residues by solid-phase phosphoramidite synthesis.

**Figure 2.** Reagents: i, 3% Cl<sub>3</sub>C-CO<sub>2</sub>H, CH<sub>2</sub>Cl<sub>2</sub>; ii, 1*H*-tetrazole **18**, MeCN; iii, Ac<sub>2</sub>O, 2,6-lutidine, 1-methylimidazole, T H<sub>2</sub>O-pyridine-THF; v, conc. aq. NH<sub>3</sub>, 55°C; vi, H<sup>+</sup>, H<sub>2</sub>O. : solid support, eg, CPG.



The target DNA sequence [29] is fully assembled after  $n$  complete synthetic cycles. The products must then be released from the solid support, and all the protecting groups need to be removed. Two alternative unblocking procedures are indicated in Fig. 2b. First, the 5'-terminal DMTr protecting group [6] is removed in what would be step i of the  $(n+1)$ th synthetic cycle. The loaded solid support is then treated with concentrated aqueous ammonia, usually at *ca.* 55°C (step v). This leads to the detachment of the synthetic DNA sequence from the solid support, the removal of all of the 2-cyanoethyl protecting groups from the internucleotide linkages and all of the *N*-acyl protecting groups from the base residues. The crude products obtained contain the target DNA sequence [30], contaminated with “capped” truncated sequences; this material is usually purified by high-performance liquid chromatography (HPLC) or by polyacrylamide gel electrophoresis (PAGE). In an alternative procedure (Fig. 2b, steps v and vi), the ammonolysis step is carried out before the removal of the 5'-terminal DMTr protecting group. The latter procedure is advantageous if purification of the crude products is to be carried out by **reversed-phase** HPLC as the difference between the retention times of the “capped” truncated sequences and the DNA target sequence, which is still protected with a 5'-*O*-DMTr group, is thereby increased. The DMTr protecting group is then removed by acidic hydrolysis under very mild conditions (step vi).

Automated solid-phase synthesis by the phosphoramidite approach has also been used successfully (11) in the preparation of DNA sequences in which the base residues, sugar residues, and internucleotide linkages are modified. DNA sequences with attached **fluorescent** and other **reporter** groups have also been prepared by solid-phase synthesis. Such modified DNA sequences have found numerous important applications in molecular biology.

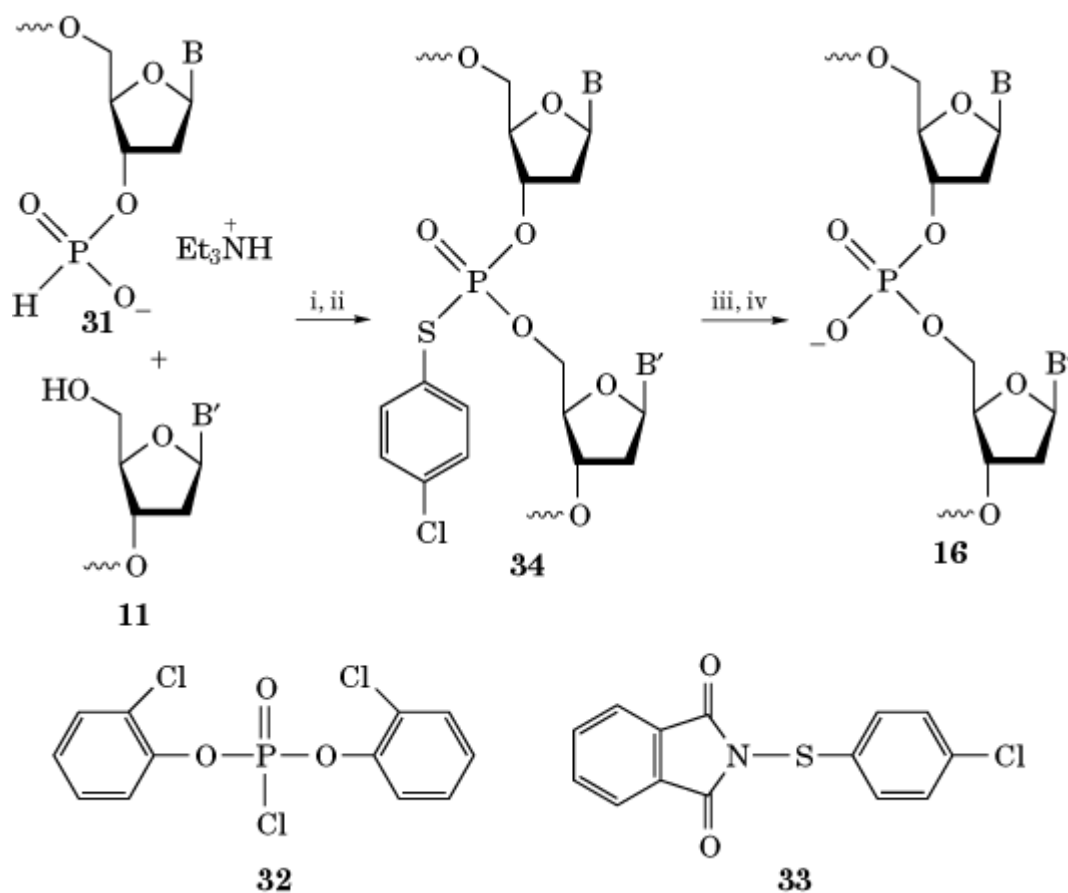
#### 4. Solution-Phase DNA Synthesis

Of the three phosphorylation methods described above, only the phosphotriester approach (Fig. 1a) is really suitable for the synthesis of DNA sequences in solution. This method, which was developed largely in the 1970s, is very versatile and is particularly suitable for the coupling of oligonucleotide blocks (i.e., the addition of two or more nucleotide residues at a time) as well as for stepwise synthesis. Phosphotriester block coupling was a key feature of the original synthesis of the human **insulin** gene (12). Although the methodology has been refined (13) since then, the development of automated solid-phase synthesis (see above) in the 1980s provided a much faster and less labor-intensive method for the preparation of the very small (usually milligram or even smaller) quantities of synthetic DNA sequences that are generally required in molecular biology. Solution-phase synthesis is much more laborious in that it is normally advisable to purify the products by chromatography after each coupling step. Although such purification processes need not necessarily amount to much more than filtration through a bed of silica gel, they are time consuming. Furthermore, solution-phase synthesis has not yet been automated. It is, nevertheless, not at all unlikely that solution-phase synthesis will become the method of choice if really large (i.e., multikilogram to tonne) quantities of moderately sized (containing *ca.* 20 nucleotide residues) DNA sequences or their analogs are required in **anti-sense** or antigene chemotherapy. Automated solid-phase synthesis has recently been scaled-up to the multigram level (14) in order to provide sufficient material for clinical trials. However, if such clinical trials are successful and very much larger quantities of pure DNA sequences and their analogs are required for drug purposes, further substantial scaling-up of solid-phase synthesis may not prove to be a practical proposition. It is quite likely that the solution-phase synthesis or perhaps a combination of solution-phase and solid-phase synthesis might lend itself much more readily to scaling-up. The phosphotriester approach has the further advantage that the fully-protected intermediates obtained are soluble in organic solvents and may, therefore, be purified by conventional chromatographic techniques, and, after all of the protecting groups have been removed, the unprotected DNA sequences obtained may, if necessary, be further purified in the same way as material that has been prepared on a solid support.

It is anticipated that success in the design of oligonucleotide drugs will stimulate further research and consequent improvements in solution-phase synthesis. Indeed, a modified phosphotriester approach involving low temperature H-phosphonate coupling (Fig. 3, step i) has recently been proposed (15).

This method appears to be more efficient than the conventional phosphotriester approach (Fig. 1a) in that virtually quantitative coupling yields are obtained, and no significant side-reactions are observed. The sensitive H-phosphonate diesters [corresponding to **23**] (Fig. 1c) are not isolated but are immediately converted at a low temperature into much more stable *S*-(4-chlorophenyl) phosphorothioates [**34**] by treatment with the phthalimide derivative [**33**] (Fig. 3, step ii). After the desired DNA sequence has been assembled, the *S*-(4-chlorophenyl) phosphorothioate groups [as in **34**] are quantitatively converted into phosphodiester [as in **16**] by standard oximate treatment (Fig. 1a).

**Figure 3.** Reagents: i, **32**, C<sub>5</sub>H<sub>5</sub>N, CH<sub>2</sub>Cl<sub>2</sub>, -40°C, ii, a, **33**, C<sub>5</sub>H<sub>5</sub>N, CH<sub>2</sub>Cl<sub>2</sub>, -40°C, b, C<sub>5</sub>H<sub>5</sub>N-H<sub>2</sub>O (1 : 1 v/v), -40°C to room temp.; iii, **14**, **15** (Fig. 1), MeCN, room temp.; iv, conc. aq. NH<sub>3</sub>, 55°C.



## Bibliography

1. H. G. Khorana (1968) *Pure Appl. Chem.* **17**, 349–381.
2. C. B. Reese and P. A. Skone (1984) *J. Chem. Soc., Perkin Trans.* **1**, 1263–1271.
3. C. B. Reese (1978) *Tetrahedron* **34**, 3143–3179.
4. J. B. Chattopadhyaya and C. B. Reese (1980) *Nucleic Acids Res.* **8**, 2039–2053.
5. C. Christodoulou (1993) In *Methods in Molecular Biology*, Vol. **20**: Protocols for Oligonucleotides and Analogs, (S. Agrawal, ed.), Humana Press, Totowa, pp. 19–31.
6. S. L. Beaucage and M. H. Caruthers (1981) *Tetrahedron Lett.* **22**, 1859–1862.
7. P. J. Garegg, T. Regberg, J. Stawinski, and R. Strömberg (1985) *Chemica Scripta* **25**, 280–282.
8. B. C. Froehler and M. D. Matteucci (1986) *Tetrahedron Lett.* **27**, 469–472.

9. S. L. Beaucage and R. P. Iyer (1992) *Tetrahedron* **48**, 2223–2311.
10. T. Brown and D. J. S. Brown (1991) in *Oligonucleotides and Analogues: A practical Approach*, (F. Eckstein, ed.), IRL Press, Oxford, pp. 1–24.
11. S. Agrawal, ed. (1994) *Methods in Molecular Biology*, Vol. **26**: Protocols for Oligonucleotide Conjugates, Humana Press, Totowa.
12. R. Crea, A. Kraszewski, T. Hirose, and K. Itakura (1978) *Proc. Natl. Acad. Sci. USA*, **75**, 5765–5769.
13. C. B. Reese and Zhang Pei-Zhuo (1993) *J. Chem. Soc., Perkin Trans. 1*, 2291–2301.
14. M. Andrade, A. S. Scozzari, D. L. Cole, and V. T. Ravikumar (1997) *Nucleosides & Nucleotides* **16**, 1617–1620.
15. C. B. Reese and Q. Song (1999) *J. Chem. Soc., Perkin Trans. 1*, 1477–1486.

### Suggestions for Further Reading

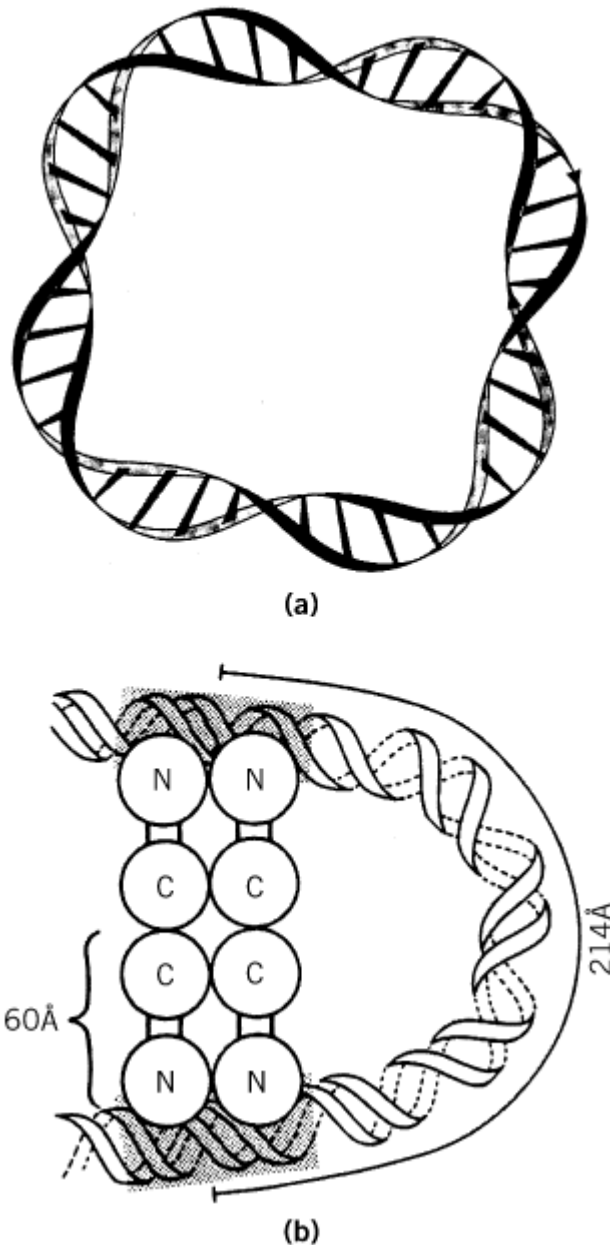
16. S. A. Narang, ed. (1987) *Synthesis and Applications of DNA and RNA*, Academic Press, Orlando.
17. F. Eckstein, ed. (1991) *Oligonucleotides and Analogues: A Practical Approach*, IRL Press, Oxford.
18. S. Agrawal, ed. (1993) *Methods in Molecular Biology*, Vol. **20**: Protocols for Oligonucleotides and Analogs, Humana Press, Totowa.
19. S. Agrawal, ed. (1993) *Methods in Molecular Biology*, Vol. **26**: Protocols for Oligonucleotide Conjugates, Humana Press, Totowa.
20. S. L. Beaucage and M. H. Caruthers (1996) In *Bioorganic Chemistry: Nucleic Acids*, (S. M. Hecht, ed.), Oxford University Press, New York, pp. 36–74.
21. S. L. Beaucage, D. E. Bergstrom, G. D. Glick, and R. A. Jones, eds. (2000) *Current Protocols in Nucleic Acid Chemistry*, Wiley, New York, pp. 3.1.1–3.4.11

## DNA Topology

Double-stranded acquire [DNA structures](#) topological properties if organized into a topological domain, which is a closed region in which the two strands are linked. The known types of elementary topological domains are closed circular duplex DNA (cdDNA), in which the two strands are separately covalently continuous; and **protein**-sealed closed DNA loops (Fig. 1) (1). The vast majority of experimental and theoretical work on DNA topological domains has been done with purified closed circular DNA. An elementary topological domain, in which the DNA is topologically monomeric, is one in which the axis of the DNA is not itself linked to that of any other DNA. DNAs of this type are widespread in nature, including bacterial plasmids and episomes, the virion and replicative forms of both eucaryotic and procaryotic viruses, mitochondrial DNA, and many others (2). The nomenclature used to describe these molecules derives in part from protein chemistry and is imprecise when applied to closed circular DNA. Among the terms commonly used are *superhelical DNA* and *supercoiled DNA*. These are not meant to indicate different structural types, but are used as synonyms to indicate that the DNA in question is more compact than either its relaxed closed circular or nicked circular counterpart. They are used interchangeably with the more generic *closed circular DNA* or, alternatively, *closed duplex DNA (cdDNA)*.



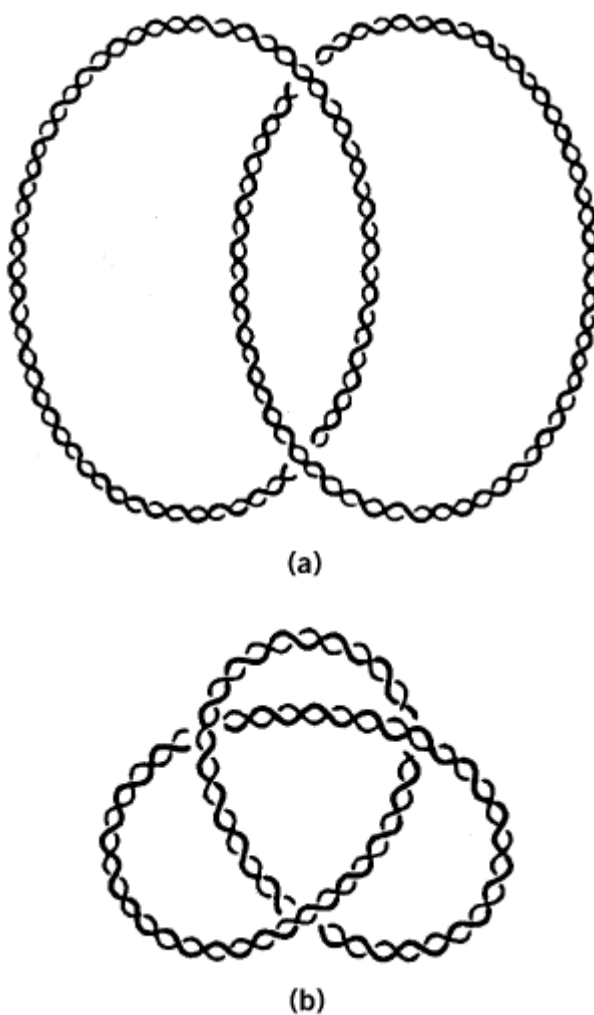
**Figure 1.** Depiction of the two kinds of elementary topological domains. (a) The structure of a closed duplex DNA (cdDNA), shown here as a relaxed duplex (rdDNA). Both strands in this DNA are covalently continuous, and the linking number is an integer. There are eight nodes in the projection shown here, and  $Lk = +4$ . (b) The structure of a DNA loop. The topological domain is maintained by tight binding to a dimeric protein. In the example shown,  $Lk \approx 6$ . The exact value of  $Lk$  depends on the precise location at which the two strands are fixed together by the protein sealant, and  $Lk$  can consequently be a fractional number. (Reproduced with permission from Ref. 1; copyright 1986, Macmillan Publishing Co.)



Protein-sealed DNA loops can occur locally in parts of DNA molecules that are not necessarily closed circular. Examples of these are the [genome](#) of *Escherichia coli*, which is organized into separate topological domains by bound proteins (3, 4), the **lambda repressor** (5, 6), several prokaryotic [transcription](#) regulation systems (7), site-specific [recombination](#) and [transposition](#) (8), and the binding of the *SfiI* [restriction enzyme](#) (9). It should be noted that protein-sealed DNA loops, in contrast to cdDNAs, are held together by noncovalent bonds. The integrity of this type of topological domain is therefore relatively less certain. Systematic determinations of the stability of DNA loops as a function of linking number have not yet been reported.

DNA can also be organized into higher order (nonelementary) topological domains, in which the duplex DNA axes of one or more circular DNAs are themselves linked. Knotted DNA, one of the possible outcomes of site-specific recombination, represents a class in which the duplex axis of a single DNA forms a knot. Catenated DNA, in which the axes of two or more DNA submolecules are linked, contains one or more higher order topological domains. The submolecules themselves may or may not contain individual elementary topological domains. Schematic depictions of DNA knots and catenanes for some simple cases are shown in Figure 2 (10).

**Figure 2.** Two kinds of higher order topological domains. (a) Two closed circular, but not supercoiled, DNAs are linked once to form a catenane. Here the higher order linking number (called the *catenation* number in this case) is 1, since the axes of the two DNAs are linked once. (b) The axis of a single DNA molecule is knotted into a trefoil. (Reproduced with permission from Ref. 10; copyright 1986, Academic Press.)

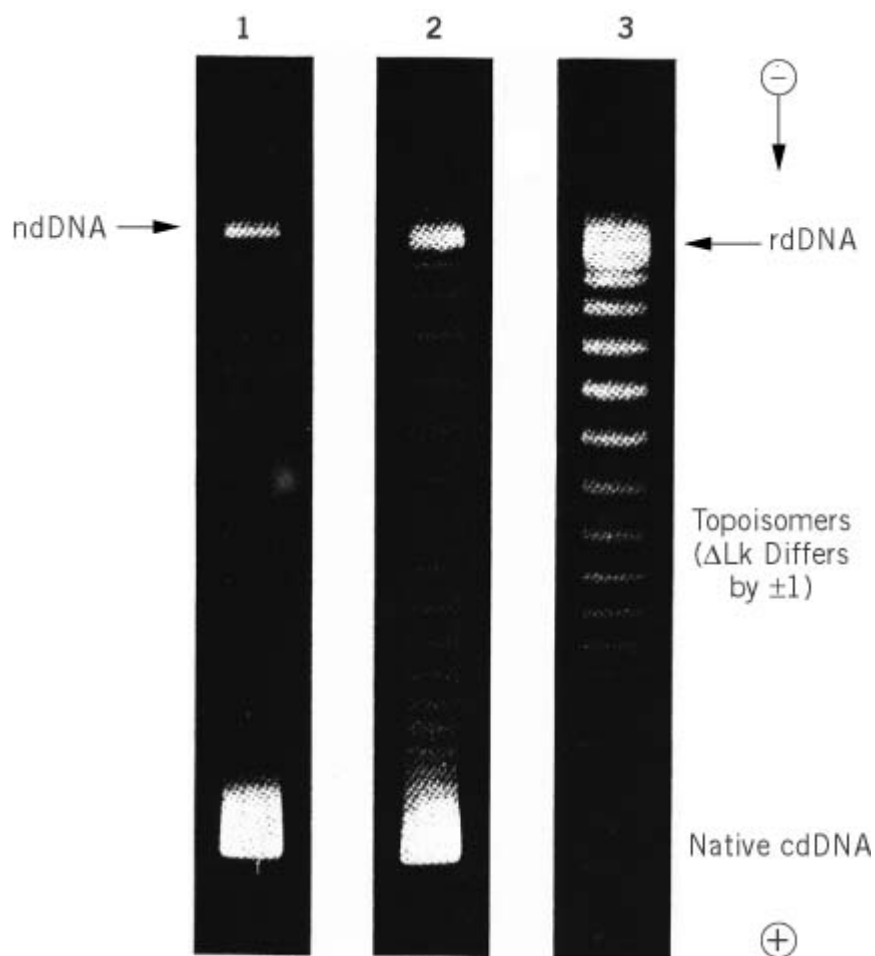


The formation of an elementary topological domain has a major influence on many DNA properties and, in addition, creates some properties that are entirely new. The properties modified may be roughly divided into two classes: structure (see [Supercoiling of DNA](#)) and energetics (see [Superhelical DNA Energetics](#)). The principal features of DNA structure affected are the duplex **twist** and the disposition in space of the DNA axis. Neat closed circular DNAs tend to form fairly regular, often branched superhelices in solution. A topological domain has a specific free energy, called either the *topological domain free energy* or the *free energy of superhelix formation*. This quantity generally increases with the extent of supercoiling of either chirality. The equilibria of all binding

reactions that involve changes in either the DNA twist or in the writhe are shifted by the inclusion of this free-energy term. Examples are the binding of intercalative drugs, such as [ethidium bromide](#), and of proteins, such as **RNA polymerase**.

The new property created by the formation of a closed circular DNA is the linking number,  $Lk$ , which is a measure of how many times the two duplex strands are interwound. It is constant for a given topological domain and can be changed only by one or more chain scissions. This makes possible the formation of DNA topoisomers, which are closed circular DNAs that differ only in  $Lk$ . Topoisomers can be reversibly interconverted with enzymes known as **topoisomerases**. The properties of individual topoisomers depend strongly on the magnitude of the associated linking number. Topoisomers can, under appropriate conditions, be physically separated by means of [gel electrophoresis](#). Figure 3 (11) shows such a fractionation for a representative closed circular DNA.

**Figure 3.** An example of the separation of DNA topoisomers by electrophoresis in an agarose gel. The three strips show the stepwise relaxation of closed circular SV40 DNA at  $0^\circ$  by incubation with a topoisomerase. Lane 1 shows the original mixture of superhelical (fastest band, labeled “Native cdDNA”) and nicked circular (slowest band, labeled “ndDNA”) species. Lanes 2 and 3 show the formation of topoisomers following 5 min and 30 min of incubation. The lowest energy product DNA is labeled “rdDNA,” and the starting superhelical DNA had  $DLk = -23$ . The individual bands contain topoisomers that differ by unity in linking number. (Because of limited gel resolution at higher  $DLk$ , not all of the higher topoisomers can be seen in this gel.) (Reproduced with permission from Ref. 11; copyright 1975, National Academy of Sciences.)



Finally, the two DNA backbone strands whose interwinding defines a topological domain contain no chemical bonds in common (except for a possible indirect connection *via* a protein sealant).

Nonetheless, the two strands can be separated only by the breakage of a covalent bond. This type of connection is sometimes termed a *topological bond*.

### Bibliography

1. M. Ptashne (1986) *Nature* **322**, 697–701.
2. W. R. Bauer (1978) *Ann. Rev. Biophys. Bioeng.* **7**, 287–313.
3. A. Worcel and E. Burgi (1972) *J. Mol. Biol.* **71**, 127–147.
4. P. Sloof, A. Maagdelyjn, and E. Boswinkel (1983) *J. Mol. Biol.* **163**, 277–297.
5. J. Griffith, A. Hochschild, and M. Ptashne (1986) *Nature* **322**, 750–752.
6. A. Hochschild and M. Ptashne (1986) *Cell* **44**, 681–687.
7. A. Hochschild (1990) In *DNA Topology and its Biological Effects* (N. R. Cozzarelli and J. C. Wang, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 107–138.
8. G. Chaconas, B. D. Lavoie, and M. A. Watson (1996) *Curr. Biol.* **6**, 817–820.
9. L. M. Wentzell and S. E. Halford (1998) *J. Mol. Biol.* **281**, 433–444.
10. A. Maxwell and M. Gellert (1986) *Adv. Protein Chem.* **38**, 69–107.
11. W. Keller and I. Wendel (1975) *Cold Spring Harb. Symp. Quant. Biol.* **39** (Pt. 1), 199–208.

### Suggestions for Further Reading

12. W. R. Bauer, F. H. C. Crick, and J. H. White (1980) Supercoiled DNA, *Sci. Am.* **243**, 118–122. (This article contains many useful illustrations of the major features of DNA topological domains.)
13. D. R. Bates and A. Maxwell (1993) *DNA Topology*, Oxford Univ. Press, London. (This is an excellent monograph that covers most aspects of DNA topology in considerable detail.)
14. N. R. Cozzarelli and J. C. Wang (1990) *DNA Topology and Its Biological Effects*, Cold Spring Harbor Press, Cold Spring Harbor, NY. (This is a more general collection of articles including the influence of DNA topology on DNA–protein interactions, superhelical structure, and DNA enzymology.)

### DNA-Binding Proteins

DNA binding proteins serve two principal functions: to organize and compact the chromosomal DNA and to regulate and effect the processes of [transcription](#), [DNA replication](#), and DNA recombination. The organization of chromosomal DNA is accomplished by abundant proteins which can bind to many sites and lack sequence specificity. In *Enterobacteria* this function is performed by FIS, H-NS, and HU proteins and in eukaryotic nuclei by the histone octamer and the linker histones, H1 and H5. By contrast, regulation of the enzymatic processes that manipulate DNA requires precise targeting to particular DNA sequences. This involves the specific recognition of the base sequence by a protein or [proteins](#). Such proteins act genetically as repressors or activators either by themselves or in combination with corepressors or coactivators.

DNA binding is specified by a large number of disparate protein motifs. Within any particular class of motif the degree of sequence selectivity is highly variable, eg, among proteins with the helix-turn-helix-motif the *lac* repressor shows a high degree of sequence specificity whereas the FIS protein possesses little. The most commonly encountered types of motif are listed in [Table 1](#).

**Table 1. Commonly Encountered Types of Motif**

| Motif                  | Target              | Sequence Specificity | Examples                           |
|------------------------|---------------------|----------------------|------------------------------------|
| Histone fold           | Backbone            | None                 | Core histones, TAFs                |
| Helix-turn-helix       | Major groove        | Highly variable      | FIS, C <sub>1</sub> repressor, CAP |
| Homeodomain            |                     |                      |                                    |
| POU domain             |                     |                      |                                    |
| Winged helix           |                     |                      | HNF3, histone H5                   |
| ETS domain             |                     |                      |                                    |
| Zinc-containing motifs |                     |                      |                                    |
| Zinc finger            | Major groove        | Generally high       | TFIIIA, Tramtrack                  |
| Receptor DBD           | Major groove        |                      | Glucocorticoid receptor            |
| Gal4 DBD               | Major groove        |                      | GAL4                               |
| GATA                   | Major+minor grooves |                      | GATA-1                             |
| bZip                   | Major groove        | High                 | GCN4, Fos-Jun                      |
| Helix-loop-helix       | Major groove        | High                 | Myc, Achaete                       |
| HMG domain             | Minor groove        | Variable, but low    | HMG1, SRY, TCF-1, UBF              |
| HU class               | Minor groove        | Variable, but low    | HU, IHF, TF1                       |
| TBP domain             | Minor groove        | High                 | TBP                                |
| GRP motif              | Minor groove        | Low, preference      | Hin recombinase, HMG14/17          |

Frequently more than one type of binding motif is found in a particular protein. In both the  $\lambda$  C<sub>1</sub> repressor and the Hin recombinase DNA binding by a helix-turn-helix motif is stabilized by additional interactions mediated by an extended strand tracking along the minor groove. Similarly binding in other proteins may be stabilized by contacts between extended peptide loops and the sugar-phosphate backbone. Bivalence, or the ability to bind to two distinct double-helical surfaces, is a characteristic of some proteins, notably the globular domain of histone H5, which can bind two adjacent duplexes, and the  $\lambda$  Int protein, which can define a DNA loop.

### DNA-Dependent DNA Polymerases

DNA-dependent DNA polymerases synthesize deoxyribonucleic acid (**DNA**), a role that is central to accurately transmitting **genetic** material from generation to generation. This family of [polymerases](#) functions in a **template**-dependent manner to insert incoming nucleotides that are encoded by the template onto a growing primer. Polymerases are highly proficient at inserting nucleotides with proper **Watson–Crick base pairing**. DNA-dependent DNA polymerase exhibits unique roles during [DNA replication](#) and [DNA repair](#). This article focuses on the characteristics of **prokaryotic** and **eukaryotic** DNA polymerases and the roles of each.

## 1. Primary Sequence, Structure, and Evolution

[Sequence analysis](#) studies (1) have divided the DNA-dependent DNA polymerase family into three subfamilies: (1) pol I type (which includes **Taq polymerase** and those from **bacteriophage T5** and **T7**); (2) pol a type (which includes phage T4 polymerase and those from **vaccinia**, [adenovirus](#), and [herpes viruses](#)); and (3) pol b type (which includes **terminal deoxynucleotidyl transferase**).

The presence of **homologous** regions in these diverse polymerases suggests that these enzymes evolved from a common ancestor (See [Polymerases](#) for a table of common motifs). All three subfamilies share two common motifs (A and C), whereas only pol I and pol a types also share a third (B). The **X-ray crystallographic** structure of the Klenow fragment of pol I suggests that motif A is involved in binding the divalent metal cofactor and deoxynucleoside triphosphates (dNTP), motif B is involved in binding the template and bases of the incoming dNTP, and motif C is involved in binding the metal cofactor. The structure of pol b shows DNA bound in an orientation 180° opposite to that observed in the complexes of DNA with **HIV-1 reverse transcriptase** and with the Klenow fragment (2). Thus, it is thought the pol b and terminal deoxynucleotidyl transferase are evolutionarily distinct from the other DNA-dependent DNA polymerases (3).

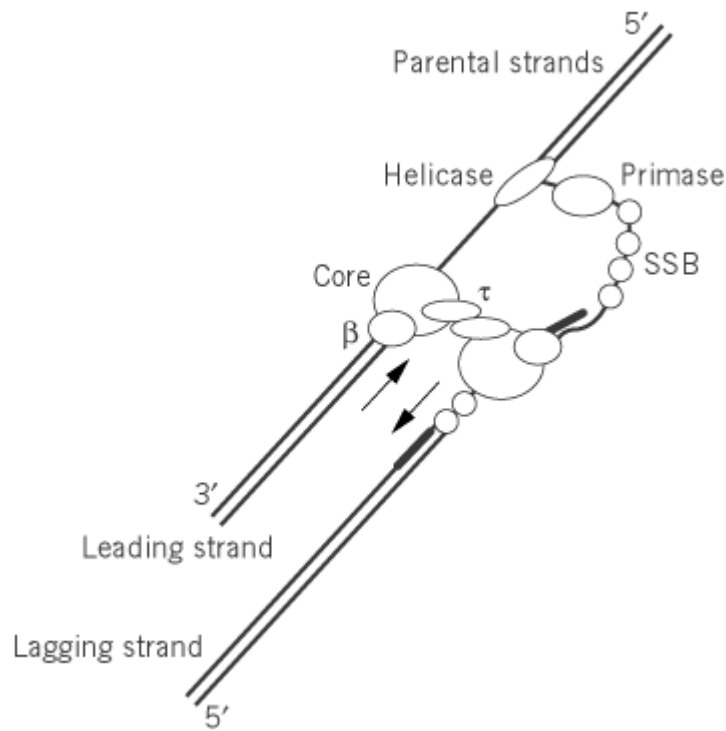
## 2. Prokaryotic DNA Polymerases

*Escherichia coli* contains three polymerases, designated DNA pol I, II, and III.

### 2.1. Pol III

The **holoenzyme** of DNA pol III (Fig. 1) is composed of 10 polypeptide subunits and is present in approximately 15 copies per bacterial cell. It is the principle polymerase responsible for replicating the [genome](#). Pol III is purified either intact or as a core enzyme separate from its accessory factors. The Pol III core consists of the polymerase subunit *a* (140 kDa; expressed from the *DnaE* gene; Ref. 4), a proofreading subunit *\** (28 kDa; from *dnaQ* or *mut D* gene; Ref. 5), and subunit *q* (10 kDa; *holE*). The processivity (i.e., the number of nucleotides synthesized per binding event) of the core enzyme is about 10 (6). This processivity is enhanced to several thousand nucleotides by the addition of subunit *b* (40 kDa), which functions as a dimer resembling a hexagonal **sliding clamp** (7). The *q* subunit (71 kDa) dimerizes readily, which, it is thought, facilitates the interaction of two polymerase core molecules. Evidence suggests that two polymerase core molecules act in concert but in opposite directions to synthesize the leading and lagging strands simultaneously (Fig. 1). The final component of the pol III holoenzyme is the *g* complex, which is composed of two subunits each of *d*, *d'*, *g*, *c*, and *j*. The *g* complex is a DNA-dependent [ATPase](#) that functions as a **clamp loader** (8). The hydrolysis of ATP facilitates loading of subunit *b* onto double-stranded DNA.

**Figure 1.** Model of the pol III replication complex during leading and lagging strand synthesis. Dimerization of the *t* domain facilitates the interaction of two core pol III molecules. This dimerization potentially allows concerted synthesis of both leading and lagging strands by one large polymerase complex. Other proteins that assist during replication in *E. coli* include DNA B helicase, primase, and single-stranded binding protein (SSB). Arrows show the 5'-3' direction of polymerization.



## 2.2. Pol I

DNA polymerase I (pol I; expressed from the *polA* gene), is a single polypeptide chain that encodes at least three separate catalytic activities. It is present at approximately 400 copies per cell, functions during [DNA repair](#) in *E. coli*, and also exhibits a minor role in DNA replication. The three catalytic activities of pol I are: (1) 5'-3' polymerization, (2) 3'-5' exonuclease (proofreading); and (3) 5'-3' exonuclease (involved in **nick translation**). Mild proteolysis of pol I yields an amino-terminal 5'-3' exonuclease domain and the COOH-terminal Klenow fragment, which contains polymerization and 3'-5' proofreading activities (9). For details on pol I functions see [DNA Polymerase I and Klenow Fragment](#).

## 2.3. Pol II

DNA polymerase II (a product of the *polB* gene) contains a 3-5' exonuclease activity, but lacks a 5'-3' exonuclease. The amount of Pol II activity is comparable to that of pol III in crude extracts from *PolA-E. coli*, and its nucleotide sequence resembles pol  $\alpha$ -type polymerases. However, the precise roles of pol II in DNA metabolism remain to be established. Mutational studies with *E. coli* show that the *pol B* gene is not required for growth or for repair after UV damage.

## 3. Eukaryotic Polymerases

The six known polymerases in eukaryotes are pol  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$ , and  $\zeta$  (see Table 1). Much of our understanding of the proteins and their roles in eukaryotic replication fork has stemmed from studying the [SV40 virus](#) replication system *in vitro* (10).

**Table 1. Eukaryotic Polymerases**

| $\alpha$ | $\beta$ | $\gamma$ | $\delta$ | $\epsilon$ | $\zeta$ |
|----------|---------|----------|----------|------------|---------|
|----------|---------|----------|----------|------------|---------|

|                              |                  |                  |              |                  |                  |           |
|------------------------------|------------------|------------------|--------------|------------------|------------------|-----------|
| Mass of core (kDa)           | 165              | 39               | 140          | 124              | 255              | 173       |
| Mass of other subunits (kDa) | 70,58,49         |                  |              | 51               | 55               | 29        |
| Chromosome map (core)        | Xp21.3-p22.1     | 8p11-p12         | 15q24        | 19q13.3-q13.4    | 12               | ?         |
| Function                     | Extends primers  | Repair           | Replication  | Replication      | Replication      | Bypass    |
| Location                     | Nucleus          | Nucleus          | Mitochondria | Nucleus          | Nucleus          | (Nucleus) |
| 3'-5' Exonucleus             | No               | No               | Yes          | Yes              | Yes              | No        |
| Fidelity                     | 10 <sup>-6</sup> | 10 <sup>-5</sup> | ?            | 10 <sup>-6</sup> | 10 <sup>-6</sup> | ?         |
| Processivity (nt)            | 10               | 5                | 100          | 10               | 1000             | ?         |
| Core subunit cloned          | Yes              | Yes              | Yes          | Yes              | Yes              | No        |

### 3.1. DNA Polymerase $\alpha$

This polymerase is composed of four subunits: (1) a 165- to 180-kDa phosphorylated **glycoprotein** that exhibits polymerase activity (its gene is located on human [chromosome Xp21.3-p22.1](#)); (2) a 70-kDa phosphoprotein with unknown function; and (3 and 4) 49- and 58-kDa subunits that have [primase](#) activity (11). An early study in which [monoclonal antibody](#) against human DNA polymerase  $\alpha$  was **microinjected** into mammalian cells showed that DNA replication is inhibited (12). Subsequent studies showed that this inhibition results from decreased extension of [Okazaki Fragments](#). It was originally thought the primary function of pol  $\alpha$  is to synthesize the lagging strand, but pol  $\alpha$ 's lack of an exonuclease domain and its poor processivity (about 10 nucleotides) cast some doubts. More recent studies of the SV40 replication system suggest that pol  $\alpha$  functions to synthesize the RNA primer for both the leading and the lagging strands *via* its primase subdomain and to extend the RNA primer by several nucleotides on both strands *via* its polymerase domain (13). The majority of leading and lagging strand synthesis can potentially be carried out by polymerases  $\delta$  and  $\epsilon$ . Mutational analysis in **yeast** of active polymerase  $\alpha$ ,  $\delta$ , and  $\epsilon$  subunits show that each is essential for growth.

### 3.2. Polymerase $\beta$

This is a 39- to 45-kDa monomeric protein in vertebrates that is constitutively expressed (its gene is located on human chromosome 8p11-12) and considered the principle enzyme involved in filling in small DNA gaps resulting from excising damaged bases (14). Although expressed as a single polypeptide chain, pol  $\beta$  contains an N-terminal 8-kDa domain, which has apurinic lyase activity responsible for cleaving the 5' phosphate-ribose residue during [base excision repair](#), and a 31-kDa C-terminal domain that contains the polymerase [active site](#). Pol  $\beta$  misincorporates noncomplementary nucleotides at a relatively high frequency, one per 5000 nucleotides polymerized (15) and lacks a 3'-5' exonuclease (proofreading) activity. Thus, DNA repair by pol  $\beta$  may be highly error-prone. Pol  $\beta$  repairs gaps of six nucleotides or less highly processively but exhibits distributive synthesis in longer gaps. Of the three types of DNA repair, [base excision repair](#) (which leaves gaps of either one nucleotide or up to 10 nucleotides *via* unique mechanisms), nucleotide [excision repair](#) (which leaves gaps of about 30 nucleotides), and [mismatch repair](#) (which leaves gaps of hundreds of nucleotides), it is thought that pol  $\beta$  is involved primarily during base excision repair. A **knockout strain** of the pol  $\beta$  gene is embryonically lethal in mice, indicating that pol  $\beta$  is essential for normal development



(16). However, the pol b-like gene of *Saccharomyces cerevisiae* (*Pol4*) may be disrupted without affecting yeast viability and only marginally affecting sensitivity to methylmethanesulfonate, a DNA alkylating agent (see [Mutagenesis](#)). Apparently, either pol d or  $\epsilon$  can substitute for pol b during base excision repair in *S. cerevisiae* (17). It has also been shown that the pol b localizes to synaptic foci during the **prophase** of **meiosis** I, suggesting that pol b activity is required to fill gaps during meiotic [recombination](#) in cells (18).

### 3.3. Polymerase g

This is a monomeric, 140-kDa enzyme that contains both polymerase and 3'-5', exonuclease activities. Although encoded in the nucleus, pol g is primarily responsible for replicating and repairing the **mitochondrial** genome. The gene that encodes human pol g (located on chromosome 15q24) has just recently been cloned. The amino acid sequence of human pol g is 40 to 50% identical to pol g from yeast and *Drosophila*, pol g is also highly homologous to *E. coli* DNA pol I (19).

### 3.4. Polymerase d

This polymerase consists of a large catalytic polypeptide chain of 124 kDa from a gene located on chromosome 19q13.3q13.4, and a small polypeptide chain of 51 kDa from a gene located on chromosome 7. Mammalian pol d was first discovered in and purified from rabbit bone marrow, the catalytic subunit of human DNA pol d has been cloned and expressed in **Baculovirus**-infected cells (20). It is one of two nuclear mammalian DNA polymerases that has 3'-5' exonuclease (proofreading) activity (21). The processivity of pol d is approximately 10 nucleotides, but this increases to well over 1000 nucleotides after the enzyme associates with a homotrimer of 36-kDa subunits of proliferating cellular nuclear antigen (PCNA), which forms a **sliding clamp** (22). In the SV40 replication system, pol d conducts most of both leading and lagging strand synthesis, complete replication in this system requires PCNA. Other proteins required for efficient pol d function include replication factor C (a phosphatase that loads PCNA onto the Pol d-DNA complex via an ATP-dependent reaction) and replication protein A (a **single-stranded binding protein** involved in reducing secondary structures of the template strand). Because of its high processivity, the Pol d-PCNA complex may also have important roles in filling DNA gaps during mismatch repair, nucleotide excision repair, and long-patch base excision repair (17).

### 3.5. Polymerase $\epsilon$

This polymerase has a large catalytic subunit of 170 kDa, whose gene is located on chromosome 12q24.3. It shares many properties with pol d, including the presence of 3'-5' exonuclease activity and very high processivity. Like pol d, pol  $\epsilon$  was also first discovered in rabbit bone marrow. Human pol  $\epsilon$  has been purified from HeLa cells and from placenta. The precise size of pol  $\epsilon$  remains controversial, as its various subunits continue to be identified. Interestingly, with multisubunit DNA polymerases like pol d and pol  $\epsilon$  the large subunit is invariably responsible for DNA polymerization. Unlike pol d, pol  $\epsilon$  is thought to be highly processive and synthesizes approximately 1000 nucleotides per binding event even when not associated with PCNA. Pol  $\epsilon$  interacts with PCNA at physiological ionic concentrations, thus further augmenting the enzyme's processivity and activity. The *in vivo* roles of pol d and  $\epsilon$  have not been adequately distinguished, but it is thought that both have roles in leading and lagging strand synthesis and in mismatch, nucleotide excision, and long-patch base excision repair.

### 3.6. Polymerase z

This sixth DNA polymerase has been described in *S. cerevisiae*. Pol z is composed of two subunits of 173 kDa (product of the yeast REV3 gene) and 29 kDa (REV7 gene). The Rev3-Rev7 complex incorporates nucleotides across damaged DNA templates (23). DNA lesions including **thymine dimers**, abasic sites, and DNA adducts are potent at arresting DNA synthesis. The translesion bypass function ascribed to pol z could be important when efficient and rapid DNA synthesis across damaged DNA is needed during replication.

## Bibliography

1. M. Delarue, O. Poch, N. Tordo, D. Moras, and P. Argos (1990) *Protein Eng.* **3**, 461–467.
2. H. Pelletier, M. R. Sawaya, A. Kumar, S. H. Wilson, and J. Kraut (1994) *Science* **264**, 1891–1903.
3. P. H. Patel, A. Jacobo-Molina, J. Ding, C. Tantillo, A. D. Clark Jr., R. Raag, R. G. Nanni, S. H. Hughes, and E. Arnold (1995) *Biochemistry* **34**, 5351–5363.
4. H. Maki and A. Kornberg (1985) *J. Biol. Chem.* **260**, 12987–12982.
5. R. H. Scheuermann and H. Echols (1984) *Proc. Natl. Acad. Sci. USA* **81**, 7747–7751.
6. P. J. Fay, K. O. Johanson, C. S. McHenry, and R. A. Bambara (1982) *J. Biol. Chem.* **257**, 5692–5699.
7. P. T. Stukenberg, P. S. Studwell-Waughan, and M. O'Donnell (1991) *J. Biol. Chem.* **266**, 11328–11334.
8. R. Onrust, P. T. Stukenberg, and M. O'Donnell (1991) *J. Biol. Chem.* **266**, 21681–21686.
9. D. Brutlag and A. Kornberg (1972) *J. Biol. Chem.* **247**, 241–248.
10. J. Li and T. J. Kelly (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6973–6977.
11. D. A. Adler, B. Y. Tseng, T. S. Wang, and C. M. Disteche (1991) *Genomics* **9**, 642–646.
12. L. Kaczmarek, M. R. Miller, R. A. Hammond, and W. E. Mercer (1986) *J. Biol. Chem.* **261**, 10802–10807.
13. R. J. Hickey and L. H. Malkas (1997) *Crit. Rev. Eukaryotic. Gene Expression* **7**, 125–157.
14. J. Abbotts, D. N. SenGupta, B. Zmudzka, S. G. Widen, V. Notario, and S. H. Wilson (1988) *Biochemistry* **27**, 901–909.
15. T. A. Kunkel and L. A. Loeb (1981) *Science* **213**, 765–767.
16. H. Gu, J. D. Marth, P. C. Orban, H. Mossmann, and K. Rajewsky (1994) *Science* **265**, 26–28.
17. A. Blank, B. Kim, and L. A. Loeb (1994) *Proc. Natl. Acad. Sci. USA* **91**, 9067–9051.
18. A. W. Plug, C. A. Clairmont, E. Sapi, T. Ashley, and J. B. Sweasy (1997) *Proc. Natl. Acad. Sci. USA* **94**, 1327–1331.
19. P. A. Ropp and W. C. Copeland (1996) *Genomics* **36**, 449–458.
20. J. Q. Zhou, C. K. Tan, A. G. So, and K. M. Downey (1996) *J. Biol. Chem.* **271**, 29740–29745.
21. J. J. Byrnes, K. M. Downey, B. G. Que, M. Y. Lee, V. L. Black, and A. G. So (1977) *Biochemistry* **16**, 3740–3746.
22. T. S. Krishna, X. P. Kong, S. Gary, P. M. Burgers, and J. Kuriyan (1994) *Cell* **79**, 1233–1243.
23. J. R. Nelson, C. W. Lawrence, and D. C. Hinkle (1996) *Science* **272**, 1646–1649.

### **Suggestion for Further Reading**

24. A. Kornberg and T. A. Baker (1992) *DNA Replication*, Freeman, New York.

### **DNA:Protein Binding Specificity**

[DNA-binding proteins](#) that bind to DNA in a sequence-specific manner have evolved to deal with a problem not normally encountered by [enzymes](#), namely that the substrate, a specific segment of DNA, is immersed in a sea of other DNA sequences on the same molecule that are chemically and structurally very similar to the specific substrate. This is a problem that is fundamentally different from the problem of the discrimination of various small molecular substrates by enzymes, which can be based at least in part on the size and shape of the substrates. More than 20 years ago, it was

pointed out that the nonspecific DNA-binding sites might be poor substrates when compared to the specific site, but they nevertheless significantly reduce the concentration of the free protein, simply because they are present in huge molar excess (1, 2). All known site-specific DNA-binding proteins also have a finite affinity for the nonspecific sites (3, 4). In this context, nonspecific DNA-binding is defined as binding that is equiprobable at any particular point along the DNA. In practice, specific and nonspecific binding constants are often determined by measuring the apparent dissociation constants of the complex of a protein with DNA containing either a specific DNA site (S) or a sequence that is completely heterologous (NS):



Table 1 compares the **dissociation constants** for specific and nonspecific binding for a number of different types of DNA-binding proteins. When referring to Table 1, it must be kept in mind that the apparent dissociation constants depend significantly on the conditions of the measurement, such as the temperature (5-15), pH (15-20), and concentration and type of cations and anions (13, 14, 16, 20, 22). A second complication stems from the fact that the stability of the specific complex is in some cases dependent on the DNA sequences flanking the specific binding site. A 17-fold difference in the dissociation constants is seen for the complexes of the [restriction enzyme EcoRV](#) and DNA sequences containing a cognate binding site with different flanking sequences (20). In addition, the apparent dissociation constants of the protein complexes with nonspecific DNA depend on the length of the DNA probe. All values for  $K_D(NS)$  in Table 1 were corrected for the length of the DNA probe by dividing the measured  $K_D(NS)$  by the twofold difference between the length of the probe and the length of the DNA-binding site of the specific complex.

**Table 1. Specificities of Prokaryotic and Eukaryotic DNA-Binding Proteins**

| Protein                    | Specific Site | $K_D$                 |                      |                   | - DDG<br>(kcal/mol) | Reference          |
|----------------------------|---------------|-----------------------|----------------------|-------------------|---------------------|--------------------|
|                            |               | $K_D(S)$<br>(M)       | $K_D(NS)$<br>(M)     | $(NS)/K_D$<br>(S) |                     |                    |
| <a href="#">Repressors</a> |               |                       |                      |                   |                     |                    |
| lac                        | $O^{sym}$     | $8.3 \times 10^{-12}$ | $2.4 \times 10^{-4}$ | $2.9 \times 10^7$ | 10.1                | <a href="#">27</a> |
| l cI                       | $O_R1$        | $1.2 \times 10^{-10}$ | $4.0 \times 10^{-4}$ | $3.3 \times 10^6$ | 8.7                 | <a href="#">28</a> |
|                            | $O_R2$        | $6.7 \times 10^{-10}$ |                      | $6.0 \times 10^5$ | 7.7                 |                    |
|                            | $O_R3$        | $6.3 \times 10^{-9}$  |                      | $6.3 \times 10^4$ | 6.4                 |                    |
| l Cro                      | $O_R1$        | $8.3 \times 10^{-12}$ | $1.5 \times 10^{-6}$ | $1.8 \times 10^5$ | 6.6                 | <a href="#">29</a> |
|                            | $O_R2$        | $1.2 \times 10^{-10}$ |                      | $1.3 \times 10^4$ | 5.1                 |                    |

|                                  |                |                       |                        |                     |         |                        |
|----------------------------------|----------------|-----------------------|------------------------|---------------------|---------|------------------------|
|                                  | $O_R^3$        | $2.0 \times 10^{-12}$ |                        | $7.5 \times 10^5$   | 7.3     |                        |
| trp                              | Wild-type op.  | $2.2 \times 10^{-9}$  | $4.2 \times 10^{-5}$   | $1.9 \times 10^4$   | 5.8     | <a href="#">8</a>      |
| P22 Arc                          | Left half-site | $2.7 \times 10^{-11}$ | $4.7 \times 10^{-8}$   | $1.7 \times 10^3$   | 4.3     | <a href="#">30</a>     |
| MetJ                             | MetBox         | $1.6 \times 10^{-6}$  | $5.0 \times 10^{-4}$   | $3.1 \times 10^2$   | 3.4     | <a href="#">9</a>      |
| <b>Restriction Endonucleases</b> |                |                       |                        |                     |         |                        |
| EcoRI                            | GAATTC         | $5.9 \times 10^{-12}$ | $2.8 \times 10^{-4}$   | $4.7 \times 10^7$   | 10.3    | <a href="#">23</a>     |
| EcoRV                            | GATATC         | $3.8 \times 10^{-10}$ | $1.0 \times 10^{-5}$   | $2.6 \times 10^4$   | 6.0     | <a href="#">20</a>     |
| <b>Transcription Factors</b>     |                |                       |                        |                     |         |                        |
| CAP                              | lac Promoter   | $1.2 \times 10^{-11}$ | $1.4 \times 10^{-6}$   | $1.2 \times 10^5$   | 6.8     | <a href="#">31</a>     |
| IHF                              | 1 attP H'      | $1.6 \times 10^{-9}$  | $2.9 \times 10^{-6}$   | $1.8 \times 10^3$   | 4.4     | <a href="#">32</a>     |
| TBP                              | TATAAAAG       | $3.7 \times 10^{-9}$  | $> 1.0 \times 10^{-5}$ | $> 2.7 \times 10^3$ | $> 4.6$ | <a href="#">33</a>     |
| MEF-2C                           | TATAAATA       | $1.1 \times 10^{-7}$  | $> 1.2 \times 10^{-5}$ | $> 1.1 \times 10^2$ | $> 2.7$ | <a href="#">34, 35</a> |
| GCN4                             | ATGACTCAT      | $3.5 \times 10^{-8}$  | $1.1 \times 10^{-6}$   | $3.1 \times 10^1$   | 2.0     | <a href="#">13</a>     |
|                                  | ATGACGTCAT     | $2.7 \times 10^{-8}$  |                        | $4.1 \times 10^1$   | 2.2     |                        |
| MASH-1                           | CAGGTG         | $5.9 \times 10^{-9}$  | $1.9 \times 10^{-8}$   | 3.2                 | 0.7     | <a href="#">15, 36</a> |
| E12                              | CAGGTG         | $1.4 \times 10^{-8}$  | $1.9 \times 10^{-7}$   | $1.4 \times 10^1$   | 1.5     | <a href="#">36</a>     |
| MyoD/E12                         | CAGGTG         | $7.5 \times 10^{-9}$  | $7.9 \times 10^{-8}$   | $1.1 \times 10^1$   | 1.4     | <a href="#">37</a>     |

Table 1 shows that the dissociation constants of the specific complexes span a range of approximately six orders of magnitude. The tightest complexes have  $K_D(S)$  values that lie in the picomolar concentration range. Interestingly, the specific DNA complexes of bacterial proteins are generally more stable than the complexes of eukaryotic DNA-binding proteins, most probably due to the longer DNA-binding sites in the prokaryotic complexes.

The dissociation constants of the nonspecific complexes listed in Table 1, on the other hand, span only approximately four orders of magnitude. For the nonspecific complexes, the eukaryotic proteins bind to DNA more tightly than the prokaryotic ones. Therefore, the DNA-binding specificity (defined as  $K_D(NS)/K_D(S)$ ) of prokaryotic DNA-binding proteins is, in most cases, significantly greater than that of eukaryotic transcription factors.

It is interesting to consider these observations in the context of the size of both the bacterial and the

mammalian [genomes](#). The *E. coli* genome consists of  $4 \times 10^6$  bp. Of the proteins listed in Table 1, the restriction enzyme EcoRI displays the highest DNA-binding specificity ( $K_D(\text{NS})/K_D(\text{S}) = 4.7 \times 10^7$ ) (23). The specific DNA-binding site of EcoRI has the sequence GAATTC. Such a hexamer sequence would occur statistically approximately 1000 times in the *E. coli* genome. As a consequence, approximately  $1.1 \times 10^4$  times more protein is bound to the specific DNA site than to the nonspecific sites:

$$[\text{PS}]/[\text{PNS}] = K_D(\text{NS}) \cdot [\text{S}]/K_D(\text{S}) \cdot [\text{NS}] \quad (2)$$

On the other hand, to bind 50% of the time to a unique binding site of the mammalian chromosome would require that the specificity of a transcription factor be  $> 3 \times 10^9$  (the size of the mammalian chromosome). Statistically, a minimal length of 16 bp is required to ensure that a given binding site is unique on the mammalian chromosome. Most transcription factors bind, however, to DNA sites that are too short to be unique on the mammalian chromosome. Proteins containing the basic **helix–loop–helix motif** (BHLH), for example, bind to the sequence CAGGTG, which occurs approximately  $7 \times 10^5$  times on a mammalian chromosome. The expression of MyoD, which recognizes DNA through a BHLH domain, can activate myogenesis in a wide variety of cell types including myoblasts and fibroblasts (24, 25), while the BHLH-protein MASH-1 promotes the differentiation of committed neuronal precursor cells (26). BHLH proteins need to bind to DNA with a specificity of approximately  $4 \times 10^3$  in order to bind with equal probability to a nonspecific site and to one of the approximately 700,000 specific sites on the mammalian chromosome. But even then, MASH-1 would still activate transcription from MyoD target promoters and vice versa. Such arguments may be part of the explanation why transcriptional regulation in higher organisms relies on multiprotein complexes with the potential for combinatorial interactions.

## Bibliography

1. P. H. von Hippel and J. D. McGhee (1972) *Annu. Rev. Biochem.* **41**, 231–300.
2. S.-y. Lin and A. D. Riggs (1975) *Cell* **4**, 107–111.
3. M. T. Record Jr. and R. S. Spolar (1990) In *The Biology of Nonspecific DNA–Protein Interactions* (A. Revzin, ed.), CRC Press, Boca Raton, FL, pp. 33–69.
4. P. H. von Hippel (1994) *Science* **263**, 769–770.
5. J.-H. Ha, R. S. Spolar, and M. T. Record (1989) *J. Mol. Biol.* **209**, 801–816.
6. Y. Takeda, P. D. Ross, and C. P. Mudd (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8180–8184.
7. L. Jin, J. Yang, and J. Carey (1993) *Biochemistry* **32**, 7302–7309.
8. J. E. Ladbury, J. G. Wright, J. M. Strutevant, and P. B. Sigler (1994) *J. Mol. Biol.* **238**, 669–681.
9. D. E. Hyre and L. D. Spicer (1995) *Biochemistry* **34**, 3212–3221.
10. P. E. Merabet and Ackers (1995) *Biochemistry* **34**, 8554–8563.
11. V. Petri, M. Hsieh, and M. Brenowitz (1995) *Biochemistry* **34**, 9977–9984.
12. T. Lundbäck, C. Cairns, J.-Å. Gustafsson, J. Carlstedt-Duke, and T. Härd (1993) *Biochemistry* **32**, 5074–5082.
13. T. Lundbäck and T. Härd (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4754–4759.
14. C. Berger, I. Jelesarov, and H. R. Bosshard (1996) *Biochemistry* **35**, 14984–14991.
15. J. H. Carra and P. L. Privalov (1997) *Biochemistry* **36**, 526–535.
16. A. G. E. Künne, M. Sieber, D. Meierhans, and R. K. Allemann (1998) *Biochemistry* **37**, 4217–4223.
17. M. T. Record Jr., J.-H. Ha, and M. A. Fisher (1991) *Methods Enzymol.* **208**, 291–343.
18. D. F. Senear and G. K. Ackers (1990) *Biochemistry* **29**, 6568–6577.

19. J. D. Taylor, I. G. Badcoe, A. R. Clarke, and S. E. Halford (1991) *Biochemistry* **30**, 8743–8753.
20. L. Li, D. von Kessler, P. A. Beachy, and K. S. Matthews (1996) *Biochemistry* **35**, 9832–9839.
21. L. E. Engler, K. K. Welch, and L. Jen-Jacobsen (1997) *J. Mol Biol.* **269**, 82–101.
22. M. T. Record Jr., T. M. Lohman, and P. de Haseth (1976) *J. Mol. Biol.* **107**, 145–158.
23. M. T. Record Jr., P. L. deHaseth, and T. M. Lohman (1977) *Biochemistry* **16**, 4791–4796.
24. L. Jen-Jacobsen (1997) *Biopolymers* **44**, 153–180.
25. C. P. Emerson (1993) *Curr. Opin. Genet. Dev.* **3**, 265–274.
26. A. B. Lassar and A. Munsterberg (1994) *Curr. Opin. Cell. Biol.* **6**, 432–442.
27. L. Sommer, N. Shah, M. Rao, and D. J. Anderson (1995) *Neuron* **15**, 1245–1258.
28. D. E. Frank, R. M. Saecker, J. P. Bond, M. W. Capp, O. V. Todikov, S. E. Melcher, M. M. Levandovski, and M. T. Record Jr. (1997) *J. Mol. Biol.* **267**, 1186–1206.
29. D. F. Seneor and R. Batey (1991) *Biochemistry* **30**, 6677–6688.
30. Y. Takeda, A. Sarai, and V. M. Rivera (1989) *Proc. Natl. Acad. Sci. USA* **86**, 439–443.
31. B. M. Brown and R. T. Sauer (1993) *Biochemistry* **32**, 1354–1363.
32. M. G. Fried and D. M. Crothers (1984) *J. Mol. Biol.* **172**, 241–262.
33. S.-W. Yang and H. A. Nash (1995) *EMBO J.* **14**, 6292–6300.
34. S. Hahn, S. Buratowski, P. A. Sharp, and L. Guarente (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5718–5722.
35. D. Meierhans, M. Sieber, and R. K. Allemann (1997) *Nucleic Acids Res.* **25**, 4537–4544.
36. D. Meierhans and R. K. Allemann (1998) *J. Biol. Chem.* **273**, 26052–26060.
37. D. Meierhans, C. el Ariss, M. Neuenschwander, M. Sieber, J. F. Stackhouse, and R. K. Allemann (1995) *Biochemistry* **34**, 11026–11036.
38. A. G. E. Künne, D. Meierhans, and R. K. Allemann (1996) *FEBS Lett.* **391**, 79–83.

## DNA:Protein Interaction Thermodynamics

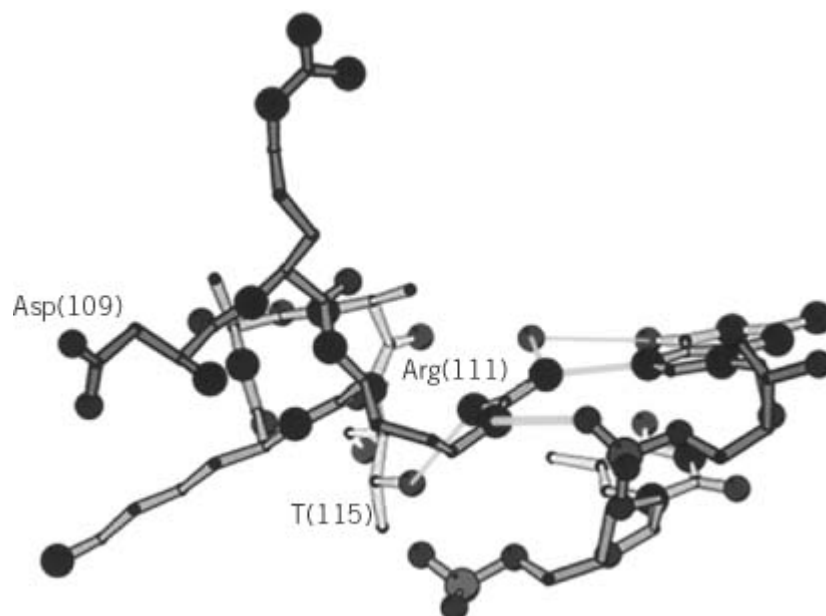
### 1. Specific Interactions

The precise recognition of a defined DNA sequence by a [DNA-binding protein](#) requires an optimal shape complementarity between the interacting species. As a result, a large number of noncovalent interactions form at the interface between the DNA and the protein. Individual interactions often contribute only a small amount to the overall stability of the complex. Nevertheless, all of these interactions are important for the preferred binding of the protein to the specific DNA site (or the discrimination of the nonspecific site).

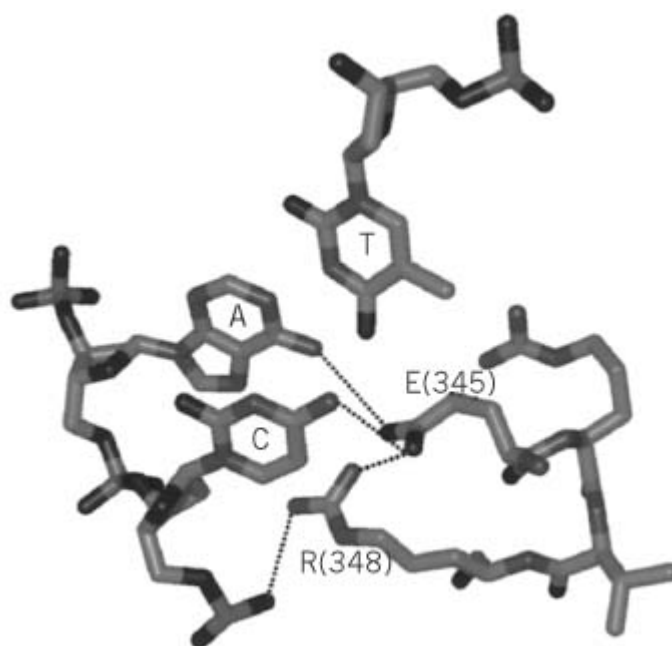
Early concepts about the determination of specificity in DNA protein complexes relied on “direct-readout,” which is based on the interaction between the functional groups of the proteins and the DNA bases ([1](#), [2](#)). Structural studies of DNA–protein complexes showed that specific interactions are mediated by [hydrogen bonds](#), [electrostatic interactions](#) between charged groups, and [van der Waals interactions](#). Constrained hydrogen bonding interactions can be made between the polypeptide backbone or short side chains and the nucleobases. Longer side chains make more flexible contacts with the DNA. They are, however, often orientated in space by other amino acid side chains or by making two contacts to the same base, to adjacent bases, or to a base and a phosphodiester group. In these cases, the **enthalpic** gain of the interaction is often reduced by a negative **entropic** contribution stemming from the reduction of conformational flexibility of the amino acid side chain. The [X-ray](#)

[crystallography](#) structures of the DNA complex of the basic **helix–loop–helix** (BHLH) proteins MyoD (3) and E47 (4) provide examples of bidentate and bridging interactions. Arg(111) of MyoD, which contacts both N(7) of guanine and an oxygen of the adjacent phosphodiester, provides an example of a bridging interaction (Fig. 1). A hydrogen bond between N(e) of Arg(111) and the hydroxyl group of a [threonine](#) residue further restricts the mobility of Arg(111). Arg(111) also forms a bidentate interaction with guanine through hydrogen bonds with N7 and O4, the latter being mediated by a [water](#) molecule. Another example is provided by the [carboxyl group](#) of a Glu(47), which makes hydrogen bonds to N7 of an adenine and to N4 of the adjacent thymine, thereby specifying the first two bases of the Glu(47) binding site (Fig. 2). A neighboring arginine residue orientates the glutamate in space through the formation of a “clamp” between the phosphate backbone and the peptide. An important consequence of the formation of such hydrogen bonding networks that severely restrict the conformational flexibility is that the specificity is achieved, at least in part, through selection of the smallest number of destabilizing interactions (5). The [restriction enzymes](#) EcoRI, EcoRV, and BamHI rely for specificity to a large extent on discrimination against noncognate sites (Ref. 6 and references therein).

**Figure 1.** Diagram showing the structure of the BHLH domain of MyoD from amino acid Asp(109) to Thr(115) and the key contacts of Arg(111) with guanine and the adjacent phosphodiester of the E-box sequence CANNTG (60). Note that the contact of Arg(111) to O6 of guanine is mediated by a water molecule and the stabilization of the arginine side chain through an hydrogen bond with the hydroxyl group of Thr(115). This display was created from the coordinates of the DNA complex of MyoD (3).



**Figure 2.** Networked hydrogen-bonding interactions between glutamate(345) and arginine(348) of E47 and the CpA dinucleotide of the E-box DNA sequence (5). E(345) is hydrogen-bonded to N4 of cytosine and to N6 of adenine. Note the “clamp” function of R(348) which connects the phosphate backbone to E(345), thereby locking the conformation of the side chain of E(345). The b and g-CH<sub>2</sub> groups of E(345) make van der Waals contacts to the methyl group of thymine. The program MacMoMo (183) was used to create this display from the coordinates of the DNA complex of E47 (4).



To the best of our knowledge, there is no site-specific DNA-binding protein for which a complete thermodynamic description of the association reaction with DNA is available, and only relatively few data are available for eukaryotic proteins. However, thermodynamic studies with the prokaryotic **Lac**, **Mnt**, **trp**, **l Cro**, **Arc**, and **l cI** repressors (7-16), the restriction endonucleases **EcoRI** (17-22), **EcoRV** (6, 23), **BamHI** (6), and **RsrI** (20), the [cyclic AMP receptor protein \(CRP\)/catabolite activating protein \(CAP\)](#) (24), the DNA-binding domain of [glucocorticoid receptor](#) (25), and the consequences of altered interactions between amino acid side chains and bases have been addressed through **site-directed mutation** of either the protein or the DNA. The results from these experiments are difficult to interpret, because of a multitude of subtle energetic changes resulting from the altered bond, perturbation of the local structure of the DNA and/or the protein, and changes in solvation. Nevertheless, a study by the group of Jen-Jacobsen of the interaction of **EcoRI** with a series of base-analogues in the absence of  $Mg^{2+}$  indicated that each hydrogen bond between the protein and the DNA stabilizes the complex by approximately 1.4 kcal/mol (19, 21). For the DNA complex of **Trp** repressor, the replacement of N7 of individual purines with carbon led to an increase in the **free energy** of binding of approximately 1 kcal/mol when the particular nitrogen atom formed a hydrogen bond to the repressor in the wild-type complex (26). Several other studies confirmed a value of 1 to 1.5 kcal/mol of stabilization for every hydrogen bond formed (23, 27). On the other hand, replacing an adenine, the N6 of which forms a hydrogen bond with a glutamate of the basic helix-loop-helix (BHLH) protein **E12** (Fig. 2), with a purine led to a reduction of only approximately 0.8 kcal/mol of the stability of the **E12** DNA complex in electrophoretic **mobility shift assay** experiments (5).

Van der Waals interactions between the protein and the methyl group of thymine that projects into the major groove in B-DNA provide another possibility for the formation of specific interactions. A careful thermodynamic study of the interaction between the methyl group of thymine -5 of the  $O_R3$  operator and a methylene group of **Lys32** of the **Cro** repressor from phage **l** indicated that the van der Waals interaction stabilizes the complex by 1.6 kcal/mol (7). When this methyl group was removed by replacing thymine with uracil, the binding enthalpy was reduced, with no observable change in the binding entropy, as would be expected for such an interaction, especially because the conformation of **Lys33** in **Cro** is stabilized by a bidentate interaction with N7 and O6 of guanine (-4) of  $O_R3$ . However, the creation of a cavity by the removal of the methyl group could potentially lead to the trapping of a water molecule. Such a case has been observed for the DNA complex of the DNA-binding domain of the glucocorticoid receptor, where a favorable enthalpic contribution



compensated for the unfavorable entropic effect of approximately 1 kcal/mol when a methyl group was removed from a thymine (28) that did not make specific contacts to the protein (29). Other studies confirmed that van der Waals interactions between a protein and a thymine methyl group stabilize the complex by between 0.5 and 2 kcal/mol (7, 19, 23, 24, 27, 30-33).

DNA-binding proteins make extensive contacts to the phosphodiester groups of the DNA backbone through charged and uncharged side chains. It is not easy to estimate the contribution from these interactions to the overall stability of the complex. One approach is to determine the contribution of the polyelectrolyte effect from the dependence of the reaction free energy on the concentration of univalent salts and to add a contribution for the interactions between the DNA phosphates and uncharged amino acid residues. It is important to keep in mind that the polyelectrolyte effect is a consequence of the fact that DNA is a polyanion of high axial charge density (reviewed in Refs. 11, 34, and 35). The association reaction is therefore driven by the release of monovalent cations (at least at low salt concentrations) (36-38). The salt-dependence of binding to various specific and nonspecific DNA sequences is often different, indicating that protein-phosphate contacts are involved in determining the specificity of DNA binding (39-44).

Anions bind to proteins with relatively small binding constants (45), and their effects are most likely not of electrostatic nature (see Hofmeister Series and Salting In, Salting Out). Anions seem to affect the dissociation constants of DNA-protein complexes, especially at high anion concentrations (46, 47). The nature of the anion is important, and replacing chloride with glutamate can increase the sequence-specific DNA-binding affinity by as much as 80-fold (43). At 25°C in a low salt buffer, the leucine zipper protein GCN4 binds to an ATF/CREB site (ATGACGTCAT) and to an AP-1 site (ATGACTCAT) with the same affinity. In a buffer containing 250 mM potassium glutamate, however, the dissociation constant of the AP-1 complex was approximately one order of magnitude smaller than for the ATF/CREB complex (48).

In order to understand fully the DNA-binding specificities displayed by proteins, it is necessary to define the energetic differences between the interactions in specific and nonspecific complexes (see DNA:Protein Binding Specificity). Model studies of the DNA-binding reactions of oligopeptides suggest that the nonspecific complex is a loose association held together by Coulombic interactions between positively charged residues of the protein and the negatively charged phosphate backbone of the DNA (49, 50). In these studies, the number of DNA phosphates contacted, as well as the number of monovalent cations thermodynamically released on binding, was found to be approximately equal to the number of positively charged side chains in the peptides (49-56). Thermodynamic studies indicate that the number of phosphate contacts is often greater in the nonspecific complexes than in the specific ones (6, 11, 57).

Only two structures of nonspecific complexes are currently known, namely those involving the restriction endonuclease EcoRV (58, 59) and the DNA-binding domain of glucocorticoid receptor (29). Glucocorticoid receptor binds to DNA sequences containing inverted repeats separated by three base pairs. In the X-ray crystallography structure analysis of the DNA complex of glucocorticoid receptor, an oligonucleotide was used where the spacing had been increased to four base pairs. Therefore, while one protein subunit bound specifically to the one inverted repeat, the other faced a nonspecific binding site. Both subunits introduced an  $\alpha$ -helix into the major groove of the oligonucleotide, which was opened up by approximately 2 Å; on the specific side. This distortion is most likely achieved through the higher number of interactions in the specific complex. Comparison of the specific and the nonspecific complexes of EcoRV revealed that a loop that penetrates the major groove is partially disordered in the nonspecific case and not well-buried in the groove. The buried surface area (*vide infra*) is more than 1570 Å<sup>2</sup>; larger in the specific complex than in the nonspecific one (58).

An interesting case is provided by the BHLH protein MASH-1. The DNA-binding specificity of MASH-1 was found to be very low for all conditions studied (39, 60). While no three-dimensional

structure for a “nonspecific” complex of MASH-1 or any other BHLH protein has been determined, the available data suggest that the specific and “nonspecific” complexes are very similar (5, 39, 60, 61), and it appears that BHLH proteins might have only one binding mode. This is most probably a consequence of the exposed recognition [a-helix](#) that adopts its helical conformation only on DNA binding (39, 60). In many other proteins that rely on an a-helix for DNA recognition, the conformation of the recognition helix is stabilized through interactions with other parts of the DNA-binding domain.

## 2. Water-Mediated Contacts

Specific interactions between a protein and DNA are often mediated by bound [water](#) molecules. General thermodynamic arguments show that the entropic cost of transferring a water molecule from bulk water to a specific site within a protein or at the interface between the protein and DNA is in the range of 0 to 7 cal mol<sup>-1</sup>K<sup>-1</sup>, corresponding to a free energy change of between 0 and 2 kcal mol<sup>-1</sup> at ambient temperature (62). The free energy cost varies thereby with the polarity of the environment of the bound water.

An example of a specific contact that is mediated by a water molecule is shown in Figure 1. The interaction between N(h) of Arg(111) of MyoD and the O6 of guanine is mediated by a well-ordered water molecule (3). Other examples of specific interactions that are mediated by well-resolved water molecules have been found in the DNA complexes of cI repressor (63), Hin recombinase (64), the restriction endonuclease BamHI (65), the DNA-binding domains of [estrogen receptor](#) (66), papillomavirus-1 E2 (67), and the [antennapedia](#)(C39 S) homeodomain (68). The DNA complex of the [transcription factor](#) GATA-1 is devoid of water molecules that mediate contacts between the protein and the nucleobases, but several water molecules were found between the protein and the phosphate backbone (69).

The 1.9 Å resolution crystal structure of the complex between the Trp repressor and its operator revealed that water-mediated contacts are the principal means for site-specific recognition (70, 71). Of the 90 water molecules that are intrinsic to the chemistry of the complex, 26 are located within the protein–DNA interface (13 in each half-site). Three of these 13 mediate four contacts between the nucleobases and the protein. Interestingly, these three hydration sites are already fully occupied in the free DNA (72). The water molecules should therefore be considered as intrinsic parts of the DNA structure and as noncovalent extensions of the DNA molecule that can be used for the stereospecific recognition of the *trp* operator (see [TRP Operon](#)). Similarly, the water molecule that mediates the contact between the backbone NH of Asn(79) and the N7 of guanine is present in all of the refined crystal structures of the uncomplexed repressor, and the water between Asn(80) and the DNA is seen in most of the refined structures of the free protein (73). These water molecules can therefore be considered functional extensions of the protein surface.

While water molecules are clearly important for both affinity and specificity in many DNA complexes of transcription factors, it must be stressed that some DNA–protein interfaces are completely devoid of water. For example, the 3100 Å resolution interface between TBP and TATA-box containing DNA sequences is characterized by a perfect complementarity and complete exclusion of water (74–77).

## 3. Dehydration Effects

Thermodynamic studies demonstrated that site-specific DNA-binding reactions are characterized by negative and relatively large changes in the **heat capacity**. As a consequence, the enthalpic ( $\Delta H$ ) and the entropic ( $\Delta S$ ) contributions to the free energy change ( $\Delta G$ ) of the binding reactions vary with temperature in an almost parallel manner, making  $\Delta G$  nearly independent of temperature. Again, there are much more data available for prokaryotic DNA-binding proteins. Only recently have good thermodynamic data for the interaction between eukaryotic transcription factors and DNA

become available. The reactions studied include the association reactions between [histones](#) and DNA (78), RNA polymerase  $E_s^{70}$  and the  $P_R$  promoter of [lambda phage](#) (79), [Lac repressor](#) and the  $lac O^+$  operator (80), the  $mnt$  repressor and the  $mnt$  operator site  $O_{mnt}$  (81), the headpiece of [Lac repressor](#) and the  $lac$  operator (82), the restriction endonuclease [EcoRI](#) and its cognate binding site (82), [Cro repressor](#) and the  $O_{R3}$  operator site (7), the DNA-binding domain of [glucocorticoid repressor](#) and a [glucocorticoid response element](#) (25, 28), [Trp repressor](#) and operator (8, 9), [1 cI repressor](#) with various combinations of the three operator sites  $O_R(1)$ ,  $O_R(2)$ , and  $O_R(3)$  (10), the transcriptional activator **cyclic AMP response protein** (CRP) complexed with two molecules of [cAMP](#) and its consensus binding site (83, 84), the transcription factor [GCN4](#) and both an [ATF/CREB](#) and an [AP-1](#) binding site (48), and the DNA-binding domain of the basic **helix–loop–helix** (BHLH) protein [MASH-1](#) and [E-box](#) containing and heterologous DNA sequences (39).

It is now generally accepted that large negative values of the heat capacity change,  $Dc_p$ , are the hallmark of biological reactions that form large highly complementary interfaces, irrespective of the overall stability of the complexes formed (9, 82). The areas buried in specific protein DNA complexes range from  $\sim 1000$  to  $\sim 5500$   $\text{Aring}^2$  (29, 39, 48, 75, 84, 85), and the heat capacity changes are caused by the removal of large amounts of nonpolar surface area from water on complex formation, accompanied by release of water. Although the change in water-accessible nonpolar surface area makes the dominant contribution to  $Dc_p$ , the changes in water accessibility of the polar surface areas (mainly due to the burial of the peptide backbone) also make a smaller contribution of opposite sign (86) (see [Hydration](#)).

Because the observed values of  $Dc_p$  for the DNA-binding reactions of some proteins are too large to be accounted for solely by the amount of buried surface area in a “rigid body” association, conformational changes of both the protein and the DNA appear to occur (see below) (84, 87, 88).

Binding of proteins to DNA in a nonspecific fashion involves almost no changes in heat capacity (7, 9, 48), indicating that the formation of nonspecific complexes does not involve major dehydration. An interesting exception is provided by the DNA-binding domain of the transcription factor [MASH-1](#). The association reaction between [MASH-1](#) and an [E-box](#)-containing oligonucleotide, the natural target of [MASH-1](#), is characterized by a heat capacity change of  $-733(\pm 99)\text{cal mol}^{-1} \text{K}^{-1}$ , while the formation of a complex with heterologous DNA results in a  $Dc_p$  of  $-575\text{cal mol}^{-1} \text{K}^{-1}$  (39). X-ray crystallography studies of the specific complexes of BHLH proteins showed that the DNA is contacted by an  $\alpha$ -helix that fits snugly into the major groove (3, 4). [Circular dichroism](#) spectroscopy suggested that the protein conformation of [MASH-1](#) was rather similar in the “specific” and the “nonspecific” complexes (39, 60). Unlike other DNA-binding proteins, all DNA complexes of [MASH-1](#) appear to show the thermodynamic characteristics of specific complexes, irrespective of the particular DNA sequence.

The number of water molecules released on formation of a complex between a protein and DNA can be estimated by measuring the dissociation constant as a function of the osmotic strength altered through the addition of neutral salts (89-91). Such osmotic-stress methods have so far been applied to only a few DNA-binding reactions, namely the restriction endonucleases [EcoRI](#) (92, 93), [EcoRV](#) and [PvuII](#) (94), the [gal repressor](#) (95), [Hin recombinase](#) (96), the cyclic AMP receptor from *E. coli* (97), and the homeodomain containing the transcriptional activators **ultrabithorax** and deformed (98).

Formation of the complex between CRP and the C1 site in the  $lac$  promoter is accompanied by the release of  $79 (\pm 11)$  water molecules, while  $56 (\pm 10)$  water molecules are taken up when CRP is transferred from the C1 site to a nonspecific site (97). Depending on the neutral salt used, between 100 and 180 water molecules are released when [gal repressor](#) binds to the  $O_I$  operator site (95). An

additional 6 ( $\pm 3$ ) waters are released when the repressor is transferred from  $O_I$  to  $O_E$ , to which it binds with enhanced affinity. Interestingly, despite the close sequence similarity in their homeodomains, water activity affects differentially the DNA binding of Ultrabithorax and Deformed (98). Between 22 and 27 water molecules were released for DNA binding by Ultrabithorax, while only 5 water molecules were released when Deformed bound to its optimal sequence. On the other hand, the DNA sequence did not exert a strong effect on the magnitude of the water release associated with DNA binding by Ultrabithorax.

The osmotic-stress methods yield a value for the number of water molecules released, which is the difference between the water molecules released and those taken up, for example as a consequence of the exposure of additional surface area due to unfolding.

#### 4. Conformational Changes of the Protein on DNA Binding

Even though in some specific complexes the tightly packed interfaces between the protein and the DNA result from the docking of well-ordered, preexisting surfaces, in an increasing number of cases the conformations of both the protein and the DNA are found to change markedly in the complex. The DNA-binding reaction of the BHLH protein MASH-1, for instance, is characterized by a transition of the peptide from a largely unfolded to a mainly  $\alpha$ -helical conformation (60). Even at concentrations well above the dimerization constant for the MASH-BHLH domain, where the HLH domain is stably folded, the basic region adopts an ordered conformation only upon binding DNA (39). A similar transition was observed with the DNA-binding domains of GCN4, Fos, and Jun, where the basic region undergoes a transition to an  $\alpha$ -helical structure upon binding to DNA (99-101). In these cases, the association does not result from a simple alignment of rigid, complementary surfaces, but rather follows what is generally known as an “[induced fit](#)” mechanism (102).

Because of a reduction in water-[accessible surface](#) area, folding transitions that occur on DNA binding result in a negative heat capacity change  $Dc_p$  (see text above). Because experimental values for  $Dc_p$  are often too large to be accounted for simply by the reduction in water-accessible surface area in a rigid body association, Spolar and Record (84) have suggested how to dissect the various contributions to the entropy change. Consider as an example the DNA-binding reaction of the BHLH domain of MASH-1. This reaction shows a strong temperature dependence for both the measured  $DH$  and  $TD S$ , which compensate to make  $DG$  almost insensitive to temperature. A notable consequence is the existence of a temperature  $T_S$ , for which  $TDS$  changes sign. Therefore, the following equation holds at  $T_S$ :

$$-\Delta S_{\text{other}} = \Delta S_{\text{HE}}(T_S) + \Delta S_{\text{RT}} + \Delta S_{\text{PE}}$$

The total change in entropy consists of a contribution  $DS_{\text{HE}}$  from the [hydrophobic effect](#), the unfavorable entropic term  $DS_{\text{RT}}$ , due to the reduction in rotational and translational degrees of freedom on association, the contribution from the polyelectrolyte effect  $DS_{\text{PE}}$ , and  $DS_{\text{other}}$ , which results primarily from conformational changes in the protein and/or the DNA (84).  $DS_{\text{HE}}(T_S)$  can be calculated from measured thermodynamic data according to the equation

$$\Delta S_{\text{HE}}(T_S) = 1.35\Delta c_p \ln(T_S/386)$$

For MASH-1,  $T_S$  was determined as 271 K and a value of  $-357 \text{ cal mol}^{-1} \text{ K}^{-1}$  was calculated for  $DS_{\text{HE}}(T_S)$  (39). The polyelectrolyte effect could be estimated from the salt dependence as  $50 \text{ cal mol}^{-1} \text{ K}^{-1}$  ((39); Meierhans, unpublished results), while  $DS_{\text{RT}}$  for a bimolecular reaction was taken to be  $-50 \text{ cal mol}^{-1} \text{ K}^{-1}$  (103, 104). The change in entropy resulting from local folding transitions

coupled to DNA could therefore be calculated as  $357 \text{ cal mol}^{-1} \text{ K}^{-1}$ . This was interpreted to indicate that approximately 54 amino acid residues of the basic region are involved in the folding reaction, or 27 residues per BHLH subunit; this interpretation is supported by CD studies of the DNA-binding reaction of MASH-BHLH and nuclear magnetic resonance (NMR) studies of the BHLH protein E47 (39, 105).

Comparison of the experimentally determined and the calculated  $Dc_p$  for the association reaction of the transcription factor GCN4 and DNA suggested that approximately 7 amino acid residues of the basic region of this basic-zipper protein underwent a transition from a random to an  $\alpha$ -helical conformation (48). The high-resolution structures of uncomplexed Trp repressor and its specific DNA complex indicate that 16 residues of helix D are disordered in the free repressor and  $\alpha$ -helical in the DNA complex (106-109) [it should be pointed out that this interpretation of the structural data has been challenged (35)]. The same value for the number of residues that change conformation upon DNA binding was obtained from thermodynamic data, where  $D S_{\text{other}}$  was determined as  $-94 \text{ cal mol}^{-1} \text{ K}^{-1}$  (8).

Structural data suggest that conformational changes also occur for the DNA-binding reaction of the DNA-binding domain of glucocorticoid receptor (110-113). These results are in good agreement with thermodynamic measurements and indicate a folding transition in approximately 18 residues per protein subunit (25). Another example of a conformational change is provided by the Antennapedia homeodomain. An *N*-terminal extension that is flexible in the free protein becomes ordered on DNA binding and contacts the minor groove of the DNA (114-116). Unfortunately, there are no thermodynamic data available for this system.

Thermodynamic studies suggest conformational changes for the DNA-binding reactions of the following proteins: the lac (117), Gal (118), and Mnt (84) repressors binding to their operator sequences, and **RNA polymerase** binding to the  $l P_R$  promoter.

For a number of DNA-binding proteins, complex formation is accompanied by changes in the **tertiary** or **quaternary structure**. For example, thermodynamic analysis of the DNA-binding reaction of the  $l$  Cro protein, under conditions where it exists as a stably folded dimer in solution, indicated a relatively small  $D S_{\text{other}}$  of  $18 \text{ cal mol}^{-1} \text{ K}^{-1}$  (7). Information from the crystal structures of both free and complexed Cro suggests that this entropy change may reflect changes in quaternary structure of the Cro dimer on DNA binding (119, 120).

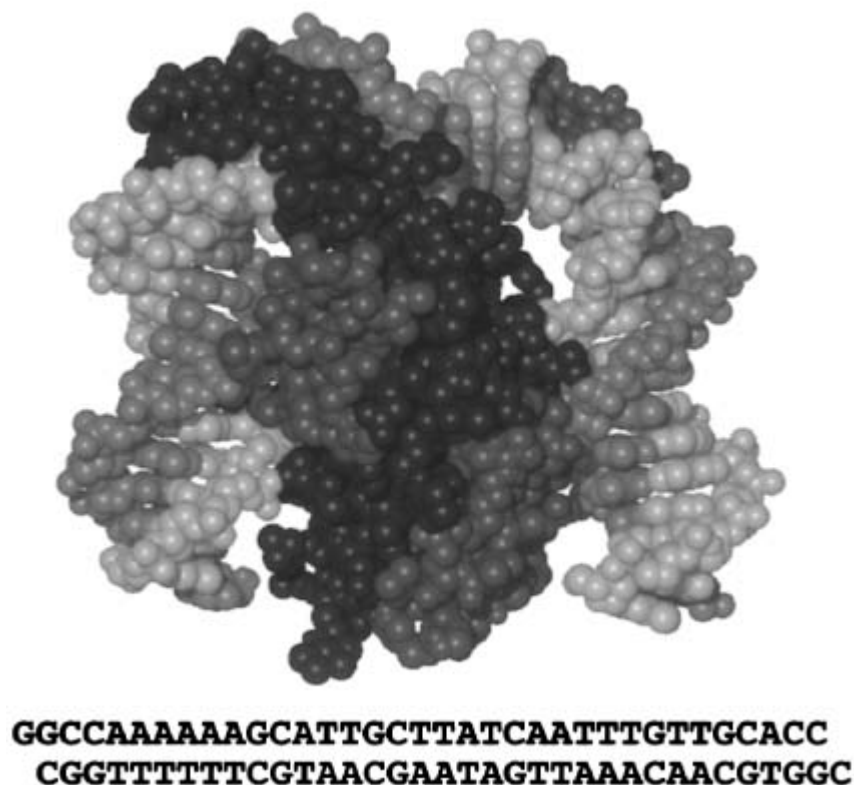
Other proteins for which structural evidence exists for coupled structural changes on binding to DNA include the endonuclease EcoRV and the transcriptional activator CRP. The DNA-binding pocket of EcoRV is not accessible in the free protein, and the cleft between the two protein subunits is too narrow to accommodate the DNA. Consequently, the pocket must be opened up through a combination of tertiary and quaternary structural transitions (58). For CRP, the relative orientations of the amino- and the carboxy-terminal domains change significantly on DNA binding. However, the structures of the DNA complexes of CRP and EcoRV revealed that optimal shape complementarity between these proteins and DNA is achieved not only through conformational adaptations of the protein, but also through changes in the structure of the DNA.

## 5. Conformational Changes of the DNA upon Protein Binding

While many proteins recognize regular B-DNA through the formation of a number of hydrogen bonds and van der Waals interactions between amino acid side chains and functional groups of the bases, it has long been recognized in many cases that DNA can adopt a bent conformation when bound to a protein (121). A dramatic example of protein-induced DNA bending is provided by the structure of integration host factor (IHF) bound to the *H'*-site of phage  $\lambda$ , in which the 34-bp piece of

DNA is literally wrapped around the protein, creating a buried protein–DNA interface of  $4600 \text{ \AA}^2$  (Fig. 3) (122, 123). The DNA is bent by more than  $160^\circ$ , thereby almost completely reversing the direction of the DNA within a short distance.

**Figure 3.** Van der Waals representation of the IHF–DNA complex (122). The drawing illustrates both the massive DNA bend induced by the protein and the large contact surface of  $\sim 4600 \text{ \AA}^2$  between protein and DNA. The DNA sequence is indicated.



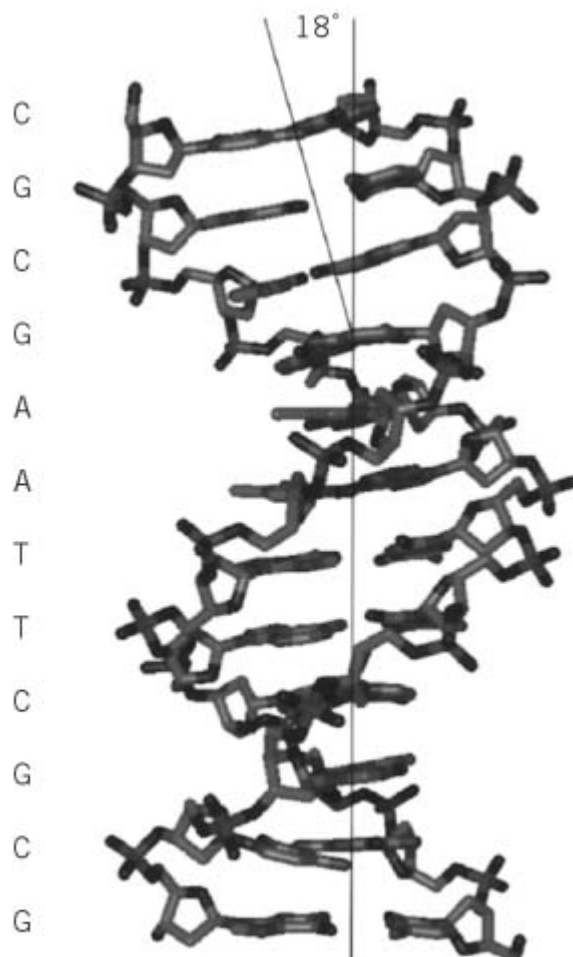
Other examples of strong protein-induced bends are provided by the structures of the DNA complexes of the purine repressor PurR (124), the sex-determining protein SRY (125), the pre-B- and T-lymphocyte-specific factor LEF-1 (126), and the human **oncogene** product ETS1 (127). In none of the above cases would it have been possible to predict the dramatic bends observed in the complexes from simple inspection of the DNA sequence.

Conversely, studies with 434 repressor indicated that its DNA-binding affinity was modulated by the flexibility of the noncontacted nucleobases at the center of the binding site (128, 129). A linear relationship was found between the free energy of the association reaction between 434 repressor and operator mutants and the flexibility of their central sequences (130). Similar results had been obtained for the Cro repressor from bacteriophage  $\lambda$  (131). These studies indicated for the first time that intrinsic sequence-dependent properties are important for the formation of DNA–protein complexes.

Sequence-dependent bending of double-stranded DNA often occurs at the junctions between regions of G–C and A–T base pairs. DNA bending of  $10^\circ$  to  $20^\circ$  has been observed in the crystal structures of oligonucleotides containing an AT-core (132). While such sequences can be bent, however, they are not necessarily bent. The transition from G–C to A–T base pairs renders this region of the DNA flexible and capable of potentially undergoing a bend. Such facultative bending is illustrated by the

X-ray structure of the dodecamer d(CGCGAATTCGCG) (Fig. 4), in which a bend of  $18^\circ$  occurs at the GC/AT junction at one end of the helix, while the helix remains unbent at the chemically equivalent junction at the other end of the helix (133, 134). The molecular mechanism of such facultative bending has been reviewed in detail elsewhere (135).

**Figure 4.** Facultative bending into the major groove at one end of the dodecamer duplex [d(CGCGAATTCGCG)]<sub>2</sub> (134). Note the positive roll at the CpG and GpA steps and the negative propeller twisting in the central four base pairs. Overall and local helical axes are shown.



The recognition sequences of **serum response factor** (SRF) (136) (see [MADS-Box Proteins](#)), TATA-binding protein (TBP) (74-77), the restriction endonucleases EcoRI and EcoRV (58, 137), and EcoRI DNA methyltransferase (138) all have AT-rich core sequences, characterized by increased bendability. The crystal structure analyses of DNA complexes of these proteins revealed that they take advantage of the inherent bendability of their DNA targets. The DNAs adopt significantly bent conformations in all of these complexes.

Myocyte enhancer factor-2C (MEF-2C) interacts with DNA in a similar fashion to SRF (139, 140), and its consensus DNA-binding site was determined in polymerase chain reaction (PCR)-mediated binding site selection assays as CTA(A/T)<sub>4</sub>TAG (141). A detailed analysis of the affinities of MEF-2C for DNA of varying sequence revealed that mutations within the central four base pairs are tolerated so long as adenine is replaced with thymine or vice versa (139, 142). However, the

replacement of the central bases with guanine and cytosine significantly diminished the affinity for MEF-2C (Table 1).

**Table 1. DNA Binding Parameters for MEF-2C (2–117) and GG-MEF-2C(1–117)**

| DNA Sequence                | [P <sub>1/2</sub> ] (nM) <sup>a</sup> |                               |
|-----------------------------|---------------------------------------|-------------------------------|
|                             | MEF-2C(2–117) <sup>b</sup>            | GG-MEF-2C(1–117) <sup>c</sup> |
| TGCTGC TATAAATA GAGTGA      | 110 (±20)                             | 103 (±9)                      |
| TGCTGC TATATATA GAGTGA      | 130 (±25)                             | 82 (±5)                       |
| TGCTGC TTTAAATA GAGTGA      | 118 (±11)                             | 112 (±28)                     |
| TGCTGC TATTAATA GAGTGA      | 117 (±23)                             | 203 (±14)                     |
| TGCTGC TAATAATA GAGTGA      | 108 (±20)                             | 227 (±22)                     |
| TGCTGC AAAAAAAAAA GAGTGA    | 742 (±163)                            | 1004 (±195)                   |
| TGCTGC TATGCATA GAGTGA      | 1072 (±200)                           | 1780 (±370)                   |
| CTGCTGC TATA-ATA GAGTGA     | 129 (±26)                             | 207 (±32)                     |
| CTGCTGC TATAAAT- GAGTGA     | 109 (±11)                             | 135 (±12)                     |
| CTGCTGC -ATAAAT-<br>GAGTGAC | 339 (±73)                             | 1896 (±417)                   |
| CTGCTGC TAT—<br>ATAGAGTGAC  | 126 (±21)                             | 115 (±12)                     |

<sup>a</sup> Protein concentration for which 50% of the DNA binding sites are filled.

<sup>b</sup> Ref. 139.

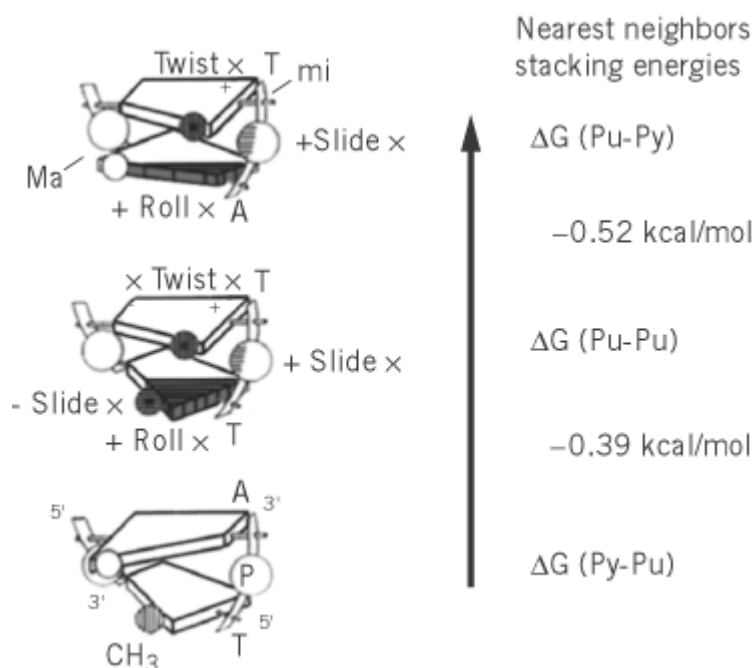
<sup>c</sup> Ref. 142.

A correlation was observed between the affinity of MEF-2C for a mutant site and the intrinsic bendability of the site. MEF-2C binds to a DNA site with an alternating run of eight thymines and adenines with maximal affinity, while the affinity for an inherently rigid A-tract sequence is reduced by more than one order of magnitude (Table 1). Because of its poor stacking (Fig. 5), the deformability of a TA step is inherently much greater than that of an AT or an AA step (74, 121, 135, 143-152). A simple mechanical model to explain this observation has been provided by Finch and co-workers (135, 153). The proximity of the methyl groups of two successive thymines and the methyl groups and the phosphate backbone make the helix rather rigid (Fig. 5). In an AT step, the stacking of the methyl group with the adjacent adenine and the intervening sugar phosphate backbone again prevents bending by a roll mechanism. However, in a TA step the methyl group projects into the major groove without any significant stacking interactions with either the adjacent adenine or the phosphate backbone, and this step therefore displays a higher deformability than the AA and AT steps.

**Figure 5.** Schematic representations of the ApT-, TpT-, and TpA-dinucleotide steps indicating the sequence specific conformational differences between these steps. The DNA is viewed from the major groove indicated by Ma. The thymine methyl groups are indicated by small spheres. The methyl groups belonging to the right strand are shaded. The twist, slide, and roll movements that could create sterically unfavorable interactions are marked with ×. The regions of



the DNA which clash due to these movements are indicated. This part of the figure was adapted from Suzuki et al. (153). The differences in the nearest-neighbor stacking energies are given on the right, indicating that the most favorable stacking interaction is observed for Pu–Py dinucleotide steps (143, 147).



Reducing the length of the AT run from eight to six nucleosides by removing the first and the eighth base of the run (which generates a consensus site for SRF) reduced the stability of the MEF-2C complex almost 20-fold (Table 1). When the six AT bases consisted of alternating thymines and adenines, however, MEF-2C bound to the corresponding oligonucleotide again with almost the same affinity as to a MEF-2C site of length eight. This observation strongly supports the proposal that the inherent bendability of the DNA-binding site is a principal determinant of the DNA-binding specificity of MEF-2C (139, 140, 142). The affinities of the various AT-rich sequences for MEF-2C might be a measure of the relative bendability of the unbound DNA.

Similarly, a comparison of the DNA structure of the free uncomplexed Trp operator with that observed in the trp repressor–operator complex provided evidence that particular DNA sequences might be predisposed to adopt a non-B-form conformation in protein complexes (72).

Circular dichroism spectroscopy and bending analysis by circular permutation assays revealed that MEF-2C potentiates the natural tendency of its DNA target to adopt a bent conformation (142). DNA binding by MEF-2C is accompanied by DNA bending of approximately 70°, irrespective of the particular DNA sequence. This observation depends on the presence of the *N*-terminal methionine residue. In its absence, DNA containing a high-affinity binding site is bent by only 49°, while heterologous sequences remain unbent (139, 142). The differences in DNA-binding affinity are much less pronounced in the absence of the *N*-terminal methionine (Table 1). The *N*-terminal methionine appears to anchor MEF-2C to the ends of the AT run, thereby orientating the protein properly on the DNA. In the minor groove, A–T base pairs can be distinguished from G–C base pairs by the lack of a heterocyclic atom at C2 of adenine. Based on the SRF structure (see [MADS-Box Proteins](#)), it had been suggested that the *N*-terminal methionine of MEF-2C is located over A4 and could specify the adenine by means of hydrophobic interaction with C2 of adenine (136). However, the observation that MEF-2C can bind with high affinity to a run of three alternating TA steps indicated that the *N*-terminal methionine does not interact with C2 of adenine (154) of the MEF site. With the short TA run, the steric clash between the side chain of methionine and the amino group on C2 of guanine (154) would reduce the stability of the complex significantly (142). It is therefore

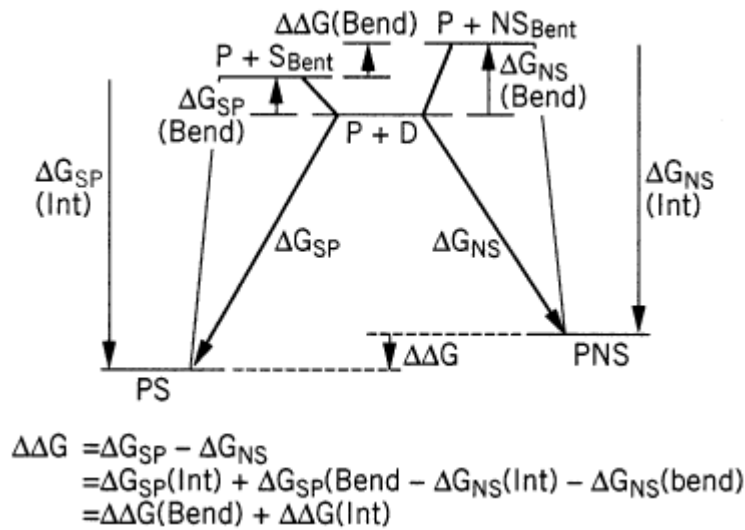
more likely that the side chain of the *N*-terminal methionine packs against the sugar ring of the nucleotide (154).

When a protein binds to a more rigid DNA sequence, it often does not bend DNA, or a substantial amount of the binding free energy must be used to bend the DNA target. Even in these cases, however, proteins often target the most deformable base step within their target sequence. The purine repressor, PurR, which is involved in the biosynthesis of purines and pyrimidines, binds to runs of four consecutive adenines and thymines, interrupted by the dinucleotide CG (155). PurR bends the DNA into the major groove by 45° through the insertion of a hinge into the minor groove and intercalation of the side chains of two leucine residues into the CG step (124). **Calorimetric** analysis suggested that purine–pyrimidine steps stack more stably than do purine–purine or pyrimidine–pyrimidine steps by approximately 0.52 kcal/mol, while the stacking energy of pyrimidine–purine steps is reduced by another 0.39 kcal/mol (Fig. 5) (143, 147). Therefore, CpG is the step that can be unstacked most easily within the pur operator. Similarly, the bending of 22° of the PRDI site within the **interferon- $\beta$**  promoter in the complex with interferon regulatory factor (IRF) occurs at the CpA step which is the most easily unstacked dinucleotide in the PRDI sequence, ACTTTCATTCTC (156). The cocrystal structure of the DNA complex of the chromosomal protein Sac7d from the hyperthermophile *Sulfolobus acidocaldarius* revealed that the DNA is kinked by 61° through the intercalation of two amino acid side chains, namely Val26 and Met29 between the bases of the dinucleotide C(2)G(3) (157). This is the base step that is most easily unstacked in the DNA sequence GCGATCGC. All these examples demonstrate that proteins recognize DNA not only through specific contacts with the nucleobases, but also through inherent, sequence-dependent properties of the DNA. Therefore, considering the properties of both the protein and the DNA makes possible an understanding of the sequence specificity and the observed DNA bend in DNA–protein complexes (135, 158).

The binding specificity of a protein is determined by the difference in binding affinities of the protein for the specific and the nonspecific sites. DNA bending is an energy costly process, even for sequences of enhanced bendability (Fig. 6). Some of the binding free energy must be used to bend a DNA molecule which, in the absence of protein, adopts either an unbent or an only slightly bent conformation. Therefore, protein-induced DNA bending can contribute to the overall specificity of a DNA-binding protein (142, 159–161). MEF-2C, for example, must use some of the binding free energy from specific interactions to induce bending at the specific site. In the case of a MEF-2C mutant lacking the *N*-terminal methionine, however, no binding free energy appears to be used for the bending of the nonspecific complex (139, 142). Consequently, the difference between the free energies of the specific and the nonspecific complexes is reduced because of the unfavorable contribution from bending the specific DNA. On the other hand, MEF-2C(1–117), which contains the *N*-terminal methionine, bends both specific and nonspecific DNA, thereby increasing the difference in the free energy between the specific and the nonspecific complexes. The increased DNA-binding specificity of MEF-2C(1–117) over that of MEF-2C(2–117) relies therefore in part on the fact that MEF-2C(1–117) bends DNA irrespective of the particular sequence. MEF-2C(1–117) appears to select DNA-binding sites that are characterized by increased bendability, because for these DNA sequences less binding free energy must be expended to force them into the bent conformation necessary for the optimal shape complementarity observed in the complexes. Interestingly, the CD spectra of the DNA complexes of MEF-2C(1–117) indicated that the conformation of the DNA was independent of its sequence, at least within the resolution of CD spectroscopy (142), while the CD spectra of the DNA in the MEF-2C(2–117) complexes varied strongly with the DNA sequence (139, 142). **Scanning tunneling microscopy** has indicated that the specific and nonspecific DNA complexes of  $\lambda$  Cro are characterized by similar bend angles. The bend angle induced by binding of a Cro dimer to DNA was determined as 69° ± 11° and 62° ± 23° for specific and non-specific DNA, respectively (159).

**Figure 6.** Free energy diagram illustrating the effect of DNA bending on the stabilities of the complexes of protein P

with DNA D. The DNA-binding specificity  $\Delta\Delta G$  is defined as the difference of the binding free energies for the formation of the specific and nonspecific complexes:  $\Delta\Delta G = \Delta G_{SP} - \Delta G_{NS}$ . The total free energy for the formation of the specific and the nonspecific complexes can be partitioned into the contributions from the (favorable) interactions between the DNA and the protein,  $\Delta G(Int)$ , and the (unfavorable) free energy change required to bend the DNA,  $\Delta G(Bend)$ .



In order to bend DNA, proteins use primarily two mechanisms. In the first, kinks are introduced through the intercalation of amino acid side chains between adjacent base pairs. In the IHF–DNA complex, two intercalating [proline](#) residues projecting from the tip of each b-hairpin arm introduce large kinks at symmetrically displaced ApA steps. SRY ([125](#)) and LEF-1 ([126](#)), which are structurally similar proteins, use nonaromatic, hydrophobic amino acids to intercalate into base steps from the minor groove side of the double helix, creating a widened minor groove and a dramatically bent DNA. TATA box-binding protein induces strong bends of approximately 45° at either end of the [TATA box](#) through the intercalation of two phenylalanine rings between two adjacent base pairs ([74-77](#)), while the hyperthermophile chromosomal proteins Sac7d bends DNA by 72°, primarily through the intercalation of Val(26) and Met(29) at a single CpA base step ([157](#)).

The second mechanism for the introduction of bends into DNA was originally proposed by Rich ([162](#)). Asymmetric neutralization of the negative charges of the phosphodiester by cationic amino acid residues could result in unbalanced Coulombic repulsions between the negative charges on DNA, causing DNA to collapse toward the bound protein ([162-164](#)). For example, the X-ray structure of **serum response factor** bound to its target DNA sequence indicated that asymmetric charge neutralization is at least partly responsible for the DNA bend angle of 72° ([136](#)). Several positively charged amino acid residues of SRF bind the phosphate groups on only one side of the SRF-site (Fig. 7). The three positively charged amino acids that interact with the distal ends of the DNA provide the hands through which SRF pulls up the far ends of the DNA. In addition, Arg143 lies in an extended conformation along the floor of the minor groove and stabilizes the relatively large propeller twists at the central base steps, which are characteristic of AT-rich regions. These interactions also facilitate the bending of DNA. The members of the MEF-2 family of proteins most probably follow a similar mechanism to induce bends into DNA ([142](#)). However, the bending of the AT-rich target sequences of SRF and of the MEF-2 proteins is helped by the intrinsically high bendability of these DNA sequences ([135](#)). The X-ray structure of the complex between the bZ protein GCN4 and an ATF-site-containing DNA revealed that the DNA was unbent ([165](#)). On the other hand, the heterodimer of the bZ proteins Jun and Fos appears to induce DNA bending when bound to an AP-1 site ([166-168](#)), although some controversy exists on this point ([169-173](#)). The identity of the amino acids just N-terminal to the basic region appears to be important for DNA bending, in that cationic amino acids in these positions contact the negatively charged phosphate diesters on only one face of the DNA, thereby inducing the DNA to bend away from the leucine zipper ([174-178](#)). The corresponding amino acids in GCN4 are the uncharged Pro–Ala–

Ala, consistent with the observation that GCN4 does not bend AP1-site-containing DNA. Analysis of the bending properties of GCN4 mutants revealed that cationic amino acids in these positions induce DNA bending toward the neutralized surface, while anionic amino acids induce bending in the opposite direction ([179](#)).

**Figure 7.** Interactions between positively charged amino acid side chains and DNA phosphates in the SRF–DNA complex ([136](#)). Note that these contacts occur only on one side of the duplex, causing the DNA to bend towards the protein. The histidines and lysines on either end of the duplex serve as handles to further bend the DNA (the DNA used is too short to allow formation of the His-phosphate contact on the left-hand side).



While DNA bending is the most dramatic change observed upon DNA binding, many other more subtle conformational rearrangements of the DNA have been observed. While a slight bending of the DNA of approximately  $25^\circ$  has been observed in the crystal structure of the 434 repressor–DNA complex ([180](#), [181](#)), the major characteristic is the significantly overwound DNA in the region of the central four base pairs of the binding site. Relative to canonical B-DNA, the net overtwisting is approximately  $20^\circ$ . Although these four base pairs are not in direct contact with the repressor, operators with A–T or T–A base pairs at these positions are bound more strongly than those bearing C–G or G–C ([182](#)). A relationship between the intrinsic twist of an operator, as determined by the sequence of its central bases, and its affinity was observed: Operators with lower affinity are undertwisted relative to operators with higher affinity ([128](#)).

#### Bibliography

1. N. C. Seeman, J. M. Rosenberg, and A. Rich (1976) *Proc. Natl. Acad. Sci. USA* **73**, 804–808.
2. J. M. Rosenberg and P. Greene (1982) *DNA* **1**, 117–124.
3. P. C. Ma, M. A. Rould, H. Weintraub, and C. O. Pabo (1994) *Cell* **77**, 451–459.
4. T. Ellenberger, D. Fass, M. Arnaud, and S. C. Harrison (1994) *Genes Dev.* **8**, 970–980.
5. M. Sieber and R. K. Allemann (1998) *Biol. Chem.* **379**, 731–735.

6. L. E. Engler, K. K. Welch, and L. Jen-Jacobsen (1997) *J. Mol Biol.* **269**, 82–101.
7. Y. Takeda, P. D. Ross, and C. P. Mudd (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8180–8184.
8. L. Jin, J. Yang, and J. Carey (1993) *Biochemistry* **32**, 7302–7309.
9. J. E. Ladbury, J. G. Wright, J. M. Strutevant, and P. B. Sigler (1994) *J. Mol. Biol.* **238**, 669–681.
10. P. E. Merabet and Ackers (1995) *Biochemistry* **34**, 8554–8563.
11. M. T. Record Jr., J.-H. Ha, and M. A. Fisher (1991) *Methods Enzymol.* **208**, 291–343.
12. D. F. Senear and G. K. Ackers (1990) *Biochemistry* **29**, 6568–6577.
13. C. D. Waldburger and R. T. Sauer (1995) *Biochemistry* **34**, 13109–13116.
14. D. E. Frank, R. M. Saecker, J. P. Bond, M. W. Capp, O. V. Todikov, S. E. Melcher, M. M. Levandovski, and M. T. Record Jr. (1997) *J. Mol. Biol.* **267**, 1186–1206.
15. B. M. Brown and R. T. Sauer (1993) *Biochemistry* **32**, 1354–1363.
16. T. L. Smith and R. T. Sauer (1995) *J. Mol. Biol.* **249**, 729–742.
17. L. Jen-Jacobsen, M. Kurpiewski, D. Lesser, J. Grable, H. W. Boyer, J. M. Rosenberg, and P. J. Greene (1983) *J. Biol. Chem.* **258**, 14638–14646.
18. L. Jen-Jacobsen, D. Lesser, and M. Kurpiewski (1986) *Cell* **45**, 619–629.
19. D. R. Lesser, M. R. Kurpiewski, and L. Jen-Jacobsen (1990) *Science* **250**, 776–786.
20. C. R. Aiken, L. W. McLaughlin, and R. I. Gumport (1991) *J. Biol. Chem.* **266**, 19070–19078.
21. D. R. Lesser, M. R. Kurpiewski, T. Waters, B. A. Connolly, and L. Jen-Jacobsen (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7548–7552.
22. L. Jen-Jacobsen (1995) *Methods Enzymol.* **259**, 305–344.
23. P. C. Newman, V. U. Nwosu, D. M. Williams, R. Cosstick, F. Seela, and B. A. Connolly (1990) *Biochemistry* **29**, 9891–9901.
24. A. Gunasekera, Y. W. Ebright, and R. H. Ebright (1992) *J. Biol. Chem.* **267**, 14713–14720.
25. T. Lundbäck, C. Cairns, J.-Å. Gustafsson, J. Carlstedt-Duke, and T. Härd (1993) *Biochemistry* **32**, 5074–5082.
26. S. A. Smith, S. B. Rajur, and L. W. McLaughlin (1994) *Nat. Struct. Biol.* (1994) **1**, 18–22.
27. J. M. Mazarelli, S. B. Rajur, P. L. Iadarola, and L. W. McLaughlin (1992) *Biochemistry* **31**, 5925–5936.
28. T. Lundbäck and T. Härd (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4754–4759.
29. B. F. Luisi, W. Xu, Z. Otwinowski, L. P. Freedman, and K. R. Yamamoto (1991) *Nature* **352**, 497–505.
30. D. V. Goeddel, D. G. Yansura, M. H. Caruthers (1977) *Nucleic Acids Res.* **4**, 3039–3052.
31. D. V. Goeddel, D. G. Yansura, and M. H. Caruthers (1978) *Proc. Natl. Acad. Sci. USA* **75**, 3578–3582.
32. E. F. Fisher and M. H. Caruthers (1979) *Nucleic Acids Res.* **7**, 401–416.
33. Y. Takeda, A. Sarai, and V. M. Rivera (1989) *Proc. Natl. Acad. Sci. USA* **86**, 439–443.
34. C. F. Anderson and M. T. Record Jr. (1995) *Annu. Rev. Phys. Chem.* **46**, 657–700.
35. T. M. Härd and T. Lundbäck (1996) *Biophys. Chem.* **62**, 121–139.
36. T. M. Lohman, P. L. deHaset, and M. T. Record Jr. (1980) *Biochemistry* **19**, 3522–3530.
37. G. S. Manning (1978) *Q. Rev. Biophys.* **11**, 179–246.
38. D. P. Mascotti and T. M. Lohman (1990) *Proc. Natl. Acad. Sci. USA* **87**, 3142–3146.
39. A. G. E. Künne, M. Sieber, D. Meierhans, and R. K. Allemann (1998) *Biochemistry* **37**, 4217–4223.

40. M. T. Record Jr., P. L. deHaseth, and T. M. Lohman (1977) *Biochemistry* **16**, 4791–4796.
41. T. M. Lohman, C. G. Wensley, J. Cina, R. R. Burgess, M. T. Record Jr. (1980) *Biochemistry* **19**, 3516–3522.
42. M. C. Mossing and M. T. Record Jr. (1985) *J. Mol. Biol.* **186**, 295–305.
43. J.-H. Ha, M. W. Capp, M. D. Hohenwalter, M. Baskerville, and M. T. Record Jr. (1992) *J. Mol. Biol.* **228**, 252–264.
44. D. F. Seneor and R. Batey (1991) *Biochemistry* **30**, 6677–6688.
45. M. T. Record Jr., C. F. Anderson, and T. M. Lohman (1978) *Q. Rev. Biophys.* **11**, 103–178.
46. S. Leirmo and M. T. Record Jr. (1990) In *Nucleic Acids and Molecular Biology*, Vol. 4 (D. M. J. Lilley and F. Eckstein, eds.), Springer, Berlin, p. 123.
47. P. H. von Hippel and T. Schleich (1969) *Acc. Chem. Res.* **2**, 257–265.
48. C. Berger, I. Jelesarov, and H. R. Bosshard (1996) *Biochemistry* **35**, 14984–14991.
49. W. Zhang, J. P. Bond, C. F. Anderson, T. M. Lohman, and M. T. Record Jr. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2511–2516.
50. S. Padmanabhan, W. Zhang, M. W. Capp, C. F. Anderson, and M. T. Record Jr. (1997) *Biochemistry* **36**, 5193–5206.
51. M. T. Record Jr., T. M. Lohman, and P. de Haseth (1976) *J. Mol. Biol.* **107**, 145–158.
52. J. D. McGhee and P. H. von Hippel (1974) *J. Mol. Biol.* **86**, 469–489.
53. W. H. Braunlin, C. F. Anderson, and M. T. Record Jr. (1987) *Biochemistry* **26**, 7724–7731.
54. G. E. Plum and V. A. Bloomfield (1988) *Biopolymers* **27**, 1045–1051.
55. D. P. Mascotti and T. M. Lohman (1992) *Biochemistry* **31**, 8932–8946.
56. D. P. Mascotti and T. M. Lohman (1993) *Biochemistry* **32**, 10568–10579.
57. M. Takahashi, B. Blazy, A. Baudras, and W. Hillen (1983) *J. Mol. Biol.* **167**, 895–899.
58. F. K. Winkler, D. W. Banner, C. Oefner, D. Tsernoglou, R. S. Brown, S. P. Heathman, R. K. Bryan, P. D. Martin, K. Petratos, and K. S. Wilson (1993) *EMBO J.* **12**, 1781–1795.
59. D. Kostrewa and F. K. Winkler (1995) *Biochemistry* **34**, 683–696.
60. D. Meierhans, C. el Ariss, M. Neuenschwander, M. Sieber, J. F. Stackhouse, and R. K. Allemann (1995) *Biochemistry* **34**, 11026–11036.
61. A. G. E. Künne, D. Meierhans, and R. K. Allemann (1996) *FEBS Lett.* **391**, 79–83.
62. J. D. Dunitz (1994) *Science* **264**, 670.
63. L. J. Beamer and C. O. Pabo (1992) *J. Mol. Biol.* **227**, 177–196.
64. J.-A. Feng, R. C. Johnson, and R. E. Dickerson (1994) *Science* **263**, 348–355.
65. M. Newman, T. Strzelecka, L. F. Dorner, I. Schildkraut, and A. K. Aggarwal (1995) *Science* **269**, 656–663.
66. J. W. R. Schwabe, L. Chapman, J. T. Finch, and D. Rhodes (1993) *Cell* **75**, 567–578.
67. R. S. Hedge, S. R. Grossman, L. A. Laimins, and P. B. Sigler (1992) *Nature* **359**, 505–512.
68. M. Billeter, P. Güntert, P. Luginbühl, and K. Wüthrich (1996) *Cell* **85**, 1057–1065.
69. G. M. Clore, A. Bax, J. G. Omichinski, and A. M. Gronenborn (1994) *Structure* **2**, 89–94.
70. P. B. Sigler (1992) In *Transcriptional Regulation*, Vol. 1 (S. L. McKnight and K. R. Yamamoto, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 475–499.
71. P. B. Sigler (1993) *Proceedings of the Robert A. Welch Foundation, 37th Conference*, pp. 63–76.
72. Z. Shakked, G. Guzikovich-Guerstein, F. Frolow, D. Rabiovich, A. Joachimiak, and P. B. Sigler (1994) *Nature* **368**, 469–473.

73. C. L. Lawson, R.-G. Zhang, R. W. Schevitz, Z. Otwinowski, A. Joachimiak, and P. B. Sigler (1988) *Proteins* **3**, 18–31.
74. Z. S. Juo, T. K. Chiu, P. M. Leiberman, I. Baikalov, A. J. Berk, and R. E. Dickerson (1996) *J. Mol. Biol.* **261**, 239–254.
75. Y. Kim, J. H. Geiger, S. Hahn, and P. B. Sigler (1993) *Nature* **365**, 512–529.
76. J. L. Kim, D. B. Nikolov, and S. K. Burley (1993) *Nature* **365**, 520–527.
77. J. L. Kim and S. K. Burley (1994) *Nat. Struct. Biol.* **1**, 638–653.
78. M. C. Polyakow, M. H. Champagne, and M. P. Daune (1972) *Eur. J. Biochem.* **26**, 212–219.
79. J. H. Roe, R. R. Burgess, and M. T. Record (1985) *J. Mol. Biol.* **184**, 441–453.
80. P. A. Whitson, J. S. Olson, and K. S. Matthews (1986) *Biochemistry* **25**, 3852–3858.
81. A. K. Vershon, S.-M. Liao, W. R. McClure, and R. T. Sauer (1987) *J. Mol. Biol.* **195**, 311–322.
82. J.-H. Ha, R. S. Spolar, and M. T. Record (1989) *J. Mol. Biol.* **209**, 801–816.
83. R. H. Ebright, Y. W. Ebright, and A. Gunasekera (1989) *Nucleic Acids Res.* **17**, 10295–10304.
84. R. S. Spolar and M. T. Record Jr. (1994) *Science* **263**, 777–784.
85. J. Janin and C. Chothia (1990) *J. Biol. Chem.* **265**, 16027–16030.
86. R. S. Spolar, J. R. Livingstone, and M. T. Record Jr. (1992) *Biochemistry* **31**, 3947–3955.
87. P. H. von Hippel (1994) *Science* **263**, 769–770.
88. S. K. Burley (1994) *Nat. Struct. Biol.* **1**, 207–208.
89. V. A. Parsegian, R. P. Rand, N. L. Fuller, and D. C. Rau (1986) *Methods Enzymol.* **127**, 400–416.
90. V. A. Parsegian, R. P. Rand, and D. C. Rau (1995) *Methods Enzymol.* **259**, 43–94.
91. R. P. Rand (1992) *Science* **256**, 618.
92. C. R. Robinson and S. G. Sligar (1993) *J. Mol. Biol.* **234**, 302–306.
93. C. R. Robinson and S. G. Sligar (1994) *Biochemistry* **33**, 3787–3793.
94. C. R. Robinson and S. G. Sligar (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3444–3448.
95. M. M. Garner and D. C. Rau (1995) *EMBO J.* **14**, 1257–1263.
96. C. R. Robinson and S. G. Sligar (1996) *Protein Sci.* **5**, 2119–2124.
97. K. M. Vossen, W. R. Daugherty, and M. G. Fried (1997) *Biochemistry* **36**, 11640–11647.
98. L. Li and K. S. Matthews (1997) *Biochemistry* **36**, 7003–7011.
99. L. Patel, C. Abate, and T. Curran (1990) *Nature* **347**, 572–575.
100. M. A. Weiss, T. Ellenberger, C. R. Wobbe, J. P. Lee, S. C. Harrison, and K. Struhl (1990) *Nature* **347**, 575–578.
101. K. T. O'Neil, R. K. Hoess, and W. F. DeGrado (1990) *Science* **249**, 774–778.
102. D. E. Koshland Jr. (1958) *Proc. Natl. Acad. Sci. USA* **44**, 98–104.
103. A. V. Finkelstein and J. Janin (1989) *Protein Eng.* **3**, 1–3.
104. J. Janin and C. Chothia (1978) *Biochemistry* **17**, 2943–2948.
105. R. Fairman, R. K. Beran-Steed, and T. M. Handel (1997) *Protein Sci.* **6**, 175–184.
106. Z. Otwinowski, R. W. Schevitz, R.-G. Zhang, C. L. Lawson, A. Joachimiak, R. Q. Marmorstein, B. F. Luisi, and P. B. Sigler (1988) *Nature* **335**, 321–329.
107. M. R. Gryk, M. D. Finucane, Z. Zheng, and O. Jardetzky (1995) *J. Mol. Biol.* **246**, 618–627.
108. D. Zhao, C. H. Arrowsmith, X. Jia, and O. Jardetzky (1993) *J. Mol. Biol.* **229**, 735–746.

109. C. H. Arrowsmith, J. Czaplicki, S. B. Iyer, and O. Jardetzky (1991) *J. Am. Chem. Soc.* **113**, 4020–4022.
110. H. Baumann, K. Paulsen, H. Kovács, H. Berglund, A. P. H. Wright, J.-A. Gustafsson, and T. Härd (1993) *Biochemistry* **32**, 13463–13471.
111. T. Härd, E. Kellenbach, R. Boelens, R. Kaptein, K. Dahlman, J. Carlstedt-Duke, L. P. Freedman, B. A. Maler, E. I. Hyde, J. A. Gustafsson, and K. R. Yamamoto (1990) *Biochemistry* **29**, 9015–9023.
112. H. Berglund, H. Kovács, K. Dahlman-Wright, J.-Å Gustafsson, and T. Härd (1992) *Biochemistry* **31**, 12001–12011.
113. M. A. L. Eriksson, H. Berglund, T. Härd, and L. Nilsson (1993) *Proteins* **17**, 375–390.
114. Y. Q. Qian, M. Billeter, G. Otting, M. Müller, W. J. Gehring, K. Wüthrich (1989) *Cell* **59**, 573–580.
115. G. Otting, Y. Q. Qian, M. Billeter, M. Müller, M. Affolter, W. J. Gehring, and K. Wüthrich (1990) *EMBO J.* **9**, 3085–3092.
116. M. Billeter, Y. Q. Qian, G. Otting, M. Müller, W. Gehring, and K. Wüthrich (1993) *J. Mol. Biol.* **234**, 1084–1097.
117. M. Lewis, G. Chang, N. C. Horton, M. A. Kercher, H. C. Pace, M. A. Schumacher, R. G. Brennan, and P. Lu (1996) *Science* **271**, 1247–1254.
118. M. Brenowitz, E. Jamison A. Majumdar, and S. Adhya (1990) *Biochemistry* **29**, 3374–3383.
119. R. G. Brennan, S. L. Roderick, Y. Takeda, and B. W. Matthews (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8165–8169.
120. R. G. Brennan (1991) *Curr. Opin. Struct. Biol.* **1**, 80–88.
121. A. Klug, A. Jack, M. A. Viswamitra, O. Kennard, Z. Shakked, and T. A. Steitz (1979) *J. Mol. Biol.* **131**, 669–680.
122. P. A. Rice, S.-W. Yang, K. Mizuuchi, and H. A. Nash (1996) *Cell* **87**, 1295–1306.
123. S.-W. Yang, and H. A. Nash (1995) *EMBO J.* **14**, 6292–6300.
124. M. A. Schumacher, K. Y. Choi, H. Zalkin, and R. G. Brennan (1994) *Science* **266**, 763–770.
125. M. H. Werner, J. R. Huth, A. M. Gronenborn, and G. M. Clore (1995) *Cell* **81**, 705–714.
126. J. J. Love, X. Li, D. A. Case, K. Giese, R. Grosschedel, and P. E. Wright (1995) *Nature* **376**, 791–795.
127. M. H. Werner, G. M. Clore, C. L. Fisher, R. J. Fisher, L. Trinh, J. Shiloach, and A. M. Gronenborn (1995) *Cell* **83**, 761–771.
128. G. B. Koudelka and P. Carlson (1992) *Nature* **355**, 89–91.
129. D. W. Rodgers and S. C. Harrison (1993) *Structure* **1**, 227–240.
130. H. R. Drew, M. J. McCall, and C. R. Calladine (1988) *Annu. Rev. Cell Biol.* **4**, 1–20.
131. A. Mondragon and S. C. Harrison (1991) *J. Mol. Biol.* **219**, 321–334.
132. R. E. Dickerson, D. Goodsell, and M. L. Kopka (1996) *J. Mol. Biol.* **256**, 108–125.
133. R. Wing, H. Drew, T. Takano, C. Broka, S. Tanaka, K. Itakura, and R. E. Dickerson (1980) *Nature* **287**, 755–758.
134. H. R. Drew, R. M. Wing, T. Takano, C. Broka, S. Tanaka, K. Itakura, and R. E. Dickerson (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2179–2183.
135. R. K. Allemann and M. Egli (1997) *Chem. Biol.* **4**, 643–650.
136. L. Pellegrini, S. Tan, and T. J. Richmond (1995) *Nature* **376**, 490–498.
137. Y. Kim, J. C. Grable, R. Love, P. J. Greene, and J. M. Rosenberg (1990) *Science* **249**, 1307–1309.



138. B. W. Allan and N. O. Reich (1996) *Biochemistry* **35**, 14757–14762.
139. D. Meierhans, M. Sieber, and R. K. Allemann (1997) *Nucleic Acids Res.* **25**, 4537–4544.
140. D. Meierhans and R. K. Allemann (1997) *Protein Expr. and Purif.* **11**, 297–303.
141. L. A. Gosset, D. J. Kelvin, E. A. Sternberg, and E. N. Olson (1989) *Mol. Cell. Biol.* **9**, 5022–5033.
142. D. Meierhans and R. K. Allemann (1998) *J. Biol. Chem.* **273**, 26052–26060.
143. O. Gotoh and Y. Tagashira (1981) *Biopolymers* **20**, 1033–1042.
144. C. R. Calladine (1982) *J. Mol. Biol.* **161**, 343–352.
145. R. E. Dickerson (1983) *J. Mol. Biol.* **166**, 419–441.
146. A. A. Travers and A. Klug (1990) In *DNA Topology and its Biological Effects* (N. R. Cozzarelli and J. C. Wang, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 57–106.
147. S. G. Delcourt and R. D. Blake (1991) *J. Biol. Chem.* **266**, 15160–15169.
148. J. R. Quintana, K. Grzeskowiak, K. Yanagi, and R. E. Dickerson (1992) *J. Mol. Biol.* **225**, 379–395.
149. A. Klug (1993) *Nature* **365**, 486–487.
150. A. A. Travers and J. W. R. Schwabe (1993) *Curr. Biol.* **3**, 898–900.
151. D. M. Goodsell, M. L. Kopka, D. Cascio, and R. E. Dickerson (1993) *Proc. Natl. Acad. Sci. USA* **90**, 2930–2934.
152. M. Suzuki and N. Yagi (1995) *Nucleic Acids Res.* **23**, 2083–2091.
153. M. Suzuki, N. Yagi, and J. T. Finch (1996) *FEBS Lett.* **379**, 148–152.
154. A. Ishihama (1988) *Trends Genet.* **4**, 282–286.
155. H. Zalkin and J. E. Dixon (1992) *Prog. Nucleic Acids Res. Mol. Biol.* **42**, 259–287.
156. C. R. Escalante, J. Yie, D. Thanos, and A. K. Aggarwal (1998) *Nature* **391**, 103–106.
157. H. Robinson, Y.-G. Gao, B. S. McCrary, S. P. Edmondson, J. W. Shriver, and A. H.-J. Wang (1998) *Nature* **392**, 202–205.
158. R. E. Dickerson (1998) *Nucleic Acids Res.* **26**, 1906–1926.
159. D. A. Erie, G. Yang, H. C. Schultz, and C. Bustamante (1994) *Science* **266**, 1562–1565.
160. A. Shepartz (1995) *Science* **269**, 989.
161. D. A. Erie and C. Bustamante (1995) *Science* **269**, 989–990.
162. A. Rich (1978) *Fed. Eur. Biochem. Soc.* **51**, 71–81.
163. J. K. Strauss and L. J. Maher III (1994) *Science* **266**, 1829–1834.
164. A. Sivolob and S. N. Khrapunov (1997) *Biophys. Chem.* **67**, 85–96.
165. W. Keller, P. König, and T. J. Richmond (1995) *J. Mol. Biol.* **254**, 657–667.
166. T. K. Kerppola and T. Curran (1991) *Cell* **66**, 317–326.
167. T. K. Kerppola and T. Curran (1991) *Science* **254**, 1210–1214.
168. T. K. Kerppola (1997) *Biochemistry* **36**, 10872–10884.
169. P. J. Hagerman (1996) *Proc. Natl. Acad. Sci. USA* **93**, 9993–9996.
170. T. K. Kerppola (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10117–10122.
171. R. J. McCormick, T. Badalian, and D. E. Fisher (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14434–14439.
172. A. Sitlani and D. M. Crothers (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3248–3252.
173. G. McGill and D. E. Fisher (1998) *Chem. Biol.* **5**, R29–R38.
174. T. K. Kerppola and T. Curran (1993) *Mol. Cell. Biol.* **13**, 5479–5489.
175. D. N. Paolella, C. R. Palmer, and A. Schepartz (1994) *Science* **264**, 1130–1133.

176. D. A. Leonard, N. Rajaram, and T. K. Kerppola (1997) *Proc. Natl. Acad. Sci. USA* **94**, 4913–4918.
177. D. N. Paolella, Y. Liu, M. A. Fabian, and A. Schepartz (1997) *Biochemistry* **36**, 10033–10038.
178. J. K. Strauss-Soukup and L. J. Maher III (1997) *Biochemistry* **36**, 10026–10032.
179. J. K. Strauss-Soukup and L. J. Maher III (1998) *Biochemistry* **37**, 1060–1066.
180. J. E. Anderson, M. Ptashne, and S. C. Harrison (1987) *Nature* **326**, 846–852.
181. A. K. Aggarwal, D. W. Rodgers, M. Drottar, M. Ptashne, and S. C. Harrison (1988) *Science* **242**, 899–907.
182. G. P. Koudelka, S. C. Harrison, and M. Ptashne (1987) *Nature* **326**, 886–888.
183. M. Dobler (1992) MoMo (color version 1.4) ETH-Zurich.

## DnaK/DnaJ Proteins

These proteins were originally discovered because mutations in their encoding genes block the replication of [lambda phage](#) in *Escherichia coli* (1). Mutations in these genes were later shown to exert pleiotropic effects on bacterial metabolism, including defects in DNA and RNA synthesis, **proteolysis**, **cell division**, **temperature-sensitive** growth, and the overproduction of **heat shock** proteins. All these effects are believed to result from changes in protein–protein interactions mediated by complexes of DnaK protein with DnaJ protein and a third protein encoded by the *grpE* gene that acts as a nucleotide-exchange factor (2-5). The DnaK and DnaJ proteins act as **molecular chaperones** in the folding (6), export (7), and degradation (8, 9) of both newly synthesized and stress-denatured polypeptide chains, and in the dissociation of oligomeric complexes essential for the initiation of phage and **plasmid DNA** replication (10). This multifaceted property stems from the ability of both DnaK and DnaJ proteins to bind and release **hydrophobic** segments of an unfolded polypeptide chain in an ATP-driven reaction cycle.

Unlike the [chaperonins](#), the DnaK and DnaJ proteins release their polypeptide substrates in unfolded states, and thus their binding serves only to shield such unfolded polypeptides transiently from premature folding during [translation](#) and from aggregation under conditions producing a **stress response**. Both the DnaK and the DnaJ proteins also have the unusual property of altering the conformations of proteins that are native, and thus seemingly fully folded; such proteins may expose chaperone recognition elements that are shielded in other proteins. The ability of both DnaK and DnaJ proteins to bind to the **heat shock** transcription factor **sigma** 32 is believed to be part of the autoregulation of stress gene expression in *E. coli* (3, 11).

The DnaK protein is a weak [ATPase](#) and is about 50% identical in primary structure to the eukaryotic family of hsp70 proteins (see [BiP \(Hsp70\)](#)), while the DnaJ protein has homologues called hsp40 proteins in the cytosol, [mitochondria](#), [chloroplasts](#), and [endoplasmic reticulum](#) of eukaryotic cells (6). Genes for homologues of both DnaK and DnaJ proteins have been identified in *E. coli* (12, 13); the expression of these genes is induced by exposure to low, rather than high, temperatures (14).

Like all the hsp70 proteins, DnaK possesses an N-terminal ATPase **domain** and a C-terminal

peptide-binding domain; the crystal structure of the latter is known (15), as is that of the ATPase domain of a **homologous** mammalian hsc70 protein (16). Polypeptides bind in extended conformations to a compact b-sandwich containing two four-stranded antiparallel b-strands. Optimal binding is produced by runs of seven hydrophobic residues; the most important determinant of peptide binding is the central binding pocket for a **leucine** residue. The ATPase domain transmits ATP-dependent conformational changes to the peptide-binding domain. All the DnaJ proteins possess a conserved *J* domain of 70 residues, which interacts with hsp70 proteins. The NMR structure of this *J* domain reveals a scaffolding of four  $\alpha$ -helices bearing at its exposed end a conserved tripeptide in a loop region (17); amino acid substitution in this tripeptide prevents binding to DnaK (18). The *J* domain is followed by a glycine-phenylalanine-rich region, and then a cysteine-rich domain resembling a **zinc finger** that is involved in binding to unfolded polypeptides (19). Thus the binding of DnaJ to DnaK combines the functions of two chaperones that have rather different specificities for binding to hydrophobic amino acid residues.

## Bibliography

1. D. E. Friedman, E. R. Olson, C. Georgopoulos, K. Tilly, I. Herskowitz, and F. Banuett (1984) *Microbiol. Rev.* **48**, 299–325.
2. C. Georgopoulos, K. Liberek, M. Zylicz, and D. Ang (1994) In *The Biology of Heat Shock Proteins and Molecular Chaperones* (R. I. Morimoto, A. Tissieres, and C. Georgopoulos, eds.), Cold Spring Harbor Press, New York, pp. 209–249.
3. B. Bukau (1993) *Mol. Microbiol.* **9**, 671–680.
4. J. Rassow, O. von Ahsen, U. Borner, and N. Pfanner (1997) *Trends Cell Biol.* **7**, 129–133.
5. D. M. Cyr, T. Langer, and M. G. Douglas (1994) *Trends Biochem. Sci.* **19**, 176–181.
6. F. U. Hartl (1996) *Nature* **381**, 571–580.
7. J. Wild, E. Altman, T. Yura, and C. A. Gross (1992) *Genes Develop.* 1165–1172.
8. S. A. Hayes and J. F. Dice (1996) *J. Cell Biol.* **132**, 255–258.
9. D. B. Straus, W. A. Walter, and C. A. Gross (1988) *Genes Develop.* **2**, 1851–1858.
10. M. Zylicz (1993) *Phil. Trans. Roy. Soc. B* **339**, 255–373.
11. J. Gamer, H. Bujard, and B. Bukau (1992) *Cell* **69**, 833–842.
12. T. H. Kawula and M. J. Levivelt (1994) *J. Bacteriol.* **176**, 610–619.
13. B. L. Seaton and L. E. Vickery (1994) *Proc. Natl. Acad. Sci. USA* **91**, 2066–2070.
14. M. J. Lelivelt and T. H. Kawula (1995) *J. Bacteriol.* **177**, 4900–4907.
15. X. Zhu, X. Zhao, W. F. Burkholder, A. Gragerov, C. M. Ogata, M. E. Gottesman, and W. A. Hendrickson (1996) *Science* **272**, 1606–1614.
16. K. M. Flaherty, C. Deluca-Flaherty, and D. B. McKay (1990) *Nature* **346**, 623–628.
17. R. B. Hill, J. M. Flanagan, and J. H. Prestegard (1995) *Biochemistry* **34**, 5587–5596.
18. D. Wall, M. Zylicz, and C. Georgopoulos (1994) *J. Biol. Chem.* **269**, 5446–5451.
19. A. Szabo, R. Korszun, F. U. Hartl, and J. Flanagan (1996) *EMBO J.* **15**, 408–417.

## DNase 1 Sensitivity

**DNase 1** (deoxyribonuclease 1) is an **enzyme** that cuts **DNA** with relatively little specificity. The enzyme is reasonably large (>30,000 Da), so that assembling DNA into **nucleosomes** and **chromatin** impedes its enzymatic activity by limiting access to the DNA double helix. Sensitivity to DNase I is

commonly believed to reflect the degree to which DNA is compacted in the chromosome. However other variables, such as proteins that like to bind DNase I and the actual structure of the DNA, might also influence enzymatic activity. DNase I likes to bind to DNA across the **minor groove** of the **double helix**.

Early experiments demonstrated the selective association of nonspecific DNA-binding proteins (prokaryotic **RNA polymerases**) with **transcriptionally** active chromatin. Following these advances, Weintraub, Felsenfeld, and colleagues showed that a comparable general accessibility to **nucleases** is associated with transcriptional activity (1, 2). This general sensitivity to nucleases includes the coding region of a gene and may extend several kilobases to either side of it, potentially defining a chromosomal **domain**. DNase I normally introduces double-strand breaks into transcriptionally active chromatin more than 10 times more frequently than into inactive chromatin. However, the exact structural basis of this generalized sensitivity is unknown. Careful analysis reveals that the untranscribed regions are just as sensitive to DNase I digestion as the transcribed regions, provided the last-cut approach to the measurement of DNase I sensitivity is used. This is defined as digestion by DNase I to fragments so small (smaller than 50 bp) that the DNA no longer hybridizes efficiently to complementary strands after **denaturation** (3). A certain length of DNA is necessary to allow specific recognition (**hybridization**) of two separated single-stranded regions. An important question not yet completely resolved is whether transcription is required to generate generalized nuclease sensitivity in certain instances or whether sensitivity always precedes transcription.

A useful example of the regulated appearance of DNase I sensitivity of a gene is provided by experiments on the action of **mitogens** on quiescent cells. In response to mitogens, a subset of genes, called the immediate-early genes, is rapidly induced (4). The most studied examples of such genes are *c-myc* and *c-fos*. These two **proto-oncogenes** are transcriptionally activated within minutes. Coincident with transcription, the **chromatin** structure of the proto-oncogenes becomes more accessible to nucleases. Once proto-oncogene transcription ceases, preferential nuclease accessibility is lost (5). Possible conformational changes in **nucleosome** structure might account for such effects. It has been proposed that **histone** H3 **cysteine** residues might become accessible in nuclease-sensitive chromatin, reflecting conformational changes within the nucleosome. However, the recruitment of **RNA polymerase** or other components of the transcriptional machinery may also provide the **thiol groups** that bind to the columns retaining active chromatin (6). Nevertheless, the rapidity of the changes in nuclease sensitivity (<90 s) and their propagation in both directions 5' and 3' to the **promoter** means that transcription and hence RNA polymerase or **HMG proteins** cannot account for all of the observed changes. This is consistent with some retention due to histone H3 in higher eukaryotes. In fact, the speed of the response suggests that changes in nuclease sensitivity precede transcription, and they may play a role in regulating *c-fos* expression.

It is interesting that one of the earliest mitogen-induced nuclear signaling events coincident with proto-oncogene induction is the rapid **phosphorylation** of histone H3 on **serine** residues within its highly charged, basic, amino-terminal **domain**. **Acetylation** of the amino-terminal domains of the core histones is also likely to be a component of transcriptional activation. Whether these changes are localized to chromatin regions containing either *c-fos*, *c-myc* or the other immediate-early genes has not yet been determined (7). An additional component contributing to the prior sensitization of the proto-oncogenes to nucleases may come from the existence of **trans-acting** factors already associated with the promoter (8). Such interactions are responsible for the second landmark in chromatin: DNase I **hypersensitive sites**. These sites are the first places where DNase I introduces a double-strand break in chromatin. They usually involve small segments of DNA sequences (100 to 200 bp) and are two or more orders of magnitude more accessible to cleavage than in inactive chromatin. As the most accessible regions of chromatin to non-histone DNA-binding proteins, DNase I hypersensitive sites generally denote DNA sequences with important functions in the nucleus. Cleavage at these sites might, however, preferentially solubilize chromatin, leading to an increase in the general accessibility of a chromatin domain to DNase I. The results emphasize the utility of DNase I as a probe for both actively transcribed genes and for detecting of regulatory DNA

in the chromosome.

### Bibliography

1. H. Weintraub and M. Groudine (1976) *Science* **193**, 848–856.
2. W. I. Wood and G. Felsenfeld (1982) *J. Biol. Chem.* **257**, 7730–7736.
3. K. Jantzen, H. P. Fritton, and T. Igo-Kemenes (1986) *Nucleic Acids Res.* **14**, 6085–6099.
4. L. F. Lau and D. Nathans (1987) *Proc. Natl. Acad. Sci. USA* **84**, 1182–1189.
5. J. Feng and B. Villeponteau (1990) *Mol. Cell. Biol.* **10**, 1126–1132.
6. J. Walker et al. (1990) *J. Biol. Chem.* **265**, 5736–5746.
7. L. C. Mahadevan, A. C. Willis, and M. J. Barrah (1991) *Cell* **65**, 775–783.
8. R. E. Herrera, P. E. Shaw, and A. Nordheim (1989) *Nature* **304**, 68–70.

### Domain, Chromosomal

A chromosomal domain can be defined at either a structural level or a functional level. At this time, there is frustratingly little information concerning the relationship between these two organizational levels. Our most thorough understanding of chromosomal organization is for the most condensed, and hence most visible, of [chromosomes](#), those at **metaphase**. Although folding of **DNA** into [nucleosomes](#) leads to a sevenfold compaction in length, and the subsequent folding of arrays of nucleosomes into the [chromatin](#) fiber to a further sevenfold compaction, a massive 250-fold compaction of DNA follows the organization of the chromatin fiber into a metaphase chromosome. One model proposed to account for this compaction suggests an organization of the fiber into loops that are radially arranged along the axis of the chromosome (1). It has also been suggested that such loops represent functional chromosomal domains.

The evidence for the organization of the chromatin fiber into loops attached to a central axis in normal cells comes from several experimental approaches. Long-standing observations on the morphology of [lampbrush chromosomes](#) in amphibian oocytes show a succession of loops emerging from a single chromosomal axis. Worcel developed a model for the *E. coli* chromosome that predicted its organization into independent domains or loops. Then these studies were extended to **interphase** chromosomes from *Drosophila* cells. Intact chromosomes were subjected to very mild digestion with **DNase I** to produce single-strand nicks. Then the resulting chromosomal fragments were examined in [sedimentation velocity centrifugation](#) experiments to determine their size. It was found that the length of the fragments decreases gradually until a plateau value of approximately 85,000 bp of DNA is reached (2). This suggests that the chromatin fiber is organized into fairly uniform domains containing about 100 kbp of DNA complexed with protein. A second approach to determining the size of structural domains in the chromosome involves **electron microscopic** examination of nuclei extracted from histones by exposure to high-salt solution. It is possible to measure the length of DNA directly on the microscope grid from where it exits a residual nuclear structure (the nuclear matrix or scaffold) to where it reenters this structure. Estimated loop sizes between 40 and 90 kbp were obtained using this technique, which were consistent with biochemical measurements (3).

The development of [pulsed-field gel electrophoresis](#) allowed more systematic analysis of the separation of cleavage sites following mild nuclease digestion of nuclei (4). This technique allows resolving very large DNA molecules by [agarose](#) electrophoresis. Nuclei of **eukaryotic** cells contain an **endonuclease** (which was first responsible for the discovery of the nucleosomal repeat) that can

be activated under controlled conditions (by the addition of exogenous  $\text{Ca}^{2+}/\text{Mg}^{2+}$  to nuclei) and which is believed to have little sequence specificity. It might be expected that cleavage of DNA by this enzyme would be inhibited by folding of DNA into the chromatin fiber, but that structural discontinuities, perhaps where the loop is attached to the chromosomal axis, might allow cleavage if the enzyme were activated. Activation of this enzyme followed by resolving the resultant DNA fragments on pulsed-field gels would potentially allow determining loop sizes. A related approach to this problem uses the observation that topoisomerase II is a major component of chromosomes and, potentially, of the nuclear matrix or scaffold. The enzymatic action of topoisomerase II introduces a double-strand break into DNA, which is then resealed. These double-strand breaks can be stabilized by the use of specific drugs (e.g., epipodophyllotoxins) that inhibit the enzyme (5). Preferential cleavage sites in nuclei spaced 50 to 300 kbp apart were detected using both the endogenous nuclease or topoisomerase II in the presence of the specific inhibitors of the rejoining of the double helix. Closer analysis revealed a hierarchy of digestion, where the 300-kbp cleavage products appeared before those of 50 kbp. It has been suggested that the 50-kbp intermediate in digestion represents the first level of organization of the chromatin fiber into loops, whereas the 300-kbp kinetic intermediate represents the next level of organization. Taken together, these observations establish a strong case that large independent loops (50 to 100 kb) of the chromatin fiber represent a unit of chromosomal structure.

Many studies have focused on the non-histone proteins in the nucleus that might represent sites of attachment of chromatin at the base of loops and the DNA sequences associated with them. The biochemical nature of the nuclear skeleton, the nuclear scaffold, and the nuclear matrix that might organize DNA in the nucleus has been the subject of much debate. Initially, the metaphase scaffold of a chromosome had a morphological definition as the complex structure at the axis of a mitotic chromosome, visualized after swelling and extraction of the histones. Biochemical extraction with high-salt (2 M NaCl) or with the detergentlike lithium diiodosalicylate (LIS) was used to define the residual nucleoprotein complex at which DNA was attached to the chromosome during interphase. Now this nuclear “matrix” (after high-salt extraction) or “scaffold” (after LIS extraction) is known to contain a substantially more complex group of proteins than the metaphase scaffold itself (6). The use of these nonphysiological extraction procedures was criticized because they cause rearrangements of protein–DNA interactions and nonspecific aggregation. An alternate strategy to study nuclear infrastructure is to encapsulate cells in agarose, to extract most of the chromatin, and to leave the nucleoprotein complexes essential for nuclear integrity under physiological conditions. This last methodology generates a nuclear “skeleton” that can transcribe and replicate DNA, which indicates that functionally relevant enzymatic complexes are retained. Recent comparative experiments confirm that the nuclear matrix or scaffold interacts with gene-poor regions of the [genome](#), whereas the nuclear skeleton interacts with gene-rich regions (7). It is now clear that the function and composition of the nuclear skeleton is very different from that of the nuclear matrix or scaffold. All define chromosomal domains, yet the reasons for generating attachments to the chromatin fiber differ.

The nuclear “scaffold” is the best-defined entity with respect to both biochemistry and genetics. A major set of proteins found in the scaffold fraction includes Sc (scaffold proteins) I, II, and III (8). The function of Sc III is unknown. Sc I (170 kDa), however, is now known to be topoisomerase II, and Sc II is known to be a heterodimeric [coiled-coil](#) protein. The enzymatic activity of topoisomerase II passes DNA strands through one another. It can cause a double strand break in DNA and rejoin it. Antibodies to topoisomerase II allowed demonstrating that the protein is an integral component of mitotic chromosomes (9). Moreover, the efficiency of recovery of total cellular topoisomerase II in the scaffold fraction (>70%) makes it unlikely that the association with the scaffold fraction is accidental. Consistent with the current view that the chromatin fiber folds further into the chromosome, immunolocalization data show that topoisomerase II is found in a large number of discrete foci scattered throughout the axial region of chromosomes. These foci are very uniform in size, suggesting that they represent discrete structural complexes. Each is believed to be an anchoring complex to which chromatin loops are attached. The presence of topoisomerase II in

these complexes can be rationalized by the necessity of unraveling DNA knots and tangles that are inevitably generated during processive enzymatic processes, such as [DNA replication](#) and [transcription](#). In fact, if topoisomerase II is inactivated *in vivo* by mutation, the mutant cells die because they cannot separate their chromosomes at the end of mitosis (10).

The Sc II protein has a much more active role in directing mitotic chromosomal condensation than topo II. Now ScII is recognized as a member of the stability and maintenance of chromosomes (SMC) family of proteins. The function of this family of proteins was first characterized in the **yeast**, *Saccharomyces cerevisiae* (11). SMC proteins are conserved in structure from **fungi** to vertebrates. Each consists of five major regions: a **nucleotide-binding region**, a region of the [alpha-helix](#) with the potential to form a coiled-coil [protein-protein interaction](#) domain, a hinge region, a second coiled-coil domain, and a carboxy-terminal region. Mutational analysis indicates that all of these domains are required for SMC protein function. SMC proteins assemble into **oligomeric** structures. There are at least four different SMC proteins in *S. cerevisiae*, and all are essential for viability, which suggests that they have essential nonoverlapping functions. The phenotypes of SMC mutant cells in yeast fail to undergo **mitosis**, and the undivided nucleus splits partially. This **phenotype** is very similar to that of topoisomerase II mutants in yeast. Closer investigation reveals that SMC mutant cells fail to condense and segregate their chromosomes (11). In spite of this clear genetic evidence that topoisomerase II and the Sc II proteins have essential roles in chromosomal organization, evidence that they have a specific association with DNA is lacking thus far. Thus the nucleic acid contribution to the assembly and function of chromatin loops or structural domain remains unclear.

It has been suggested that the nuclear matrix or scaffold contains specific DNA sequences known as scaffold or [matrix attachment regions](#) (SARs and MARs). Evidence that the sequences function in chromosomal dynamics comes from the synthesis of an artificial protein that preferentially binds to these AT-rich sequences (12). This protein interfered with the chromosomal dynamics normally observed during nuclear decondensation or chromosomal condensation in *Xenopus* [egg](#) extracts. In the absence of the identification of *bona fide* scaffold attachment sequence-binding proteins *in vivo* and examination of their function, however, whether scaffold or matrix attachments have an active role in chromosomal function remains speculative. Such attachments might promote transcriptional activity *in vivo* (13). In **plants** there is compelling evidence that AT-rich segments of the genome that have been biochemically defined as scaffold attachment regions promote the activity of genes in *cis*. Histone H1 binds preferentially to these AT-rich DNA sequences. It has been suggested that proteins, such as high mobility group (**HMG I/Y**) might displace histone H1 selectively from scaffold attachment regions that contribute to the local control of transcriptional activity (14). It must be proven that this type of selective association of proteins with scaffold attachment regions occurs *in vivo*. It has long been known that lysine-rich proteins like histone H1 interact with AT-rich DNA. Hence the significance of *in vitro* binding experiments that enrich lysine-rich proteins bound to AT-rich scaffold attachment regions remains questionable. A subfamily of sites determined by such procedures as matrix attachment sites (A elements) can be found at the boundaries of a 24-kbp region of **DNase I-sensitive** chromatin containing the chicken [lysozyme](#) gene. Sippel and colleagues demonstrated that the A elements insulate a gene from chromosomal effects in stable transformants, but they are not required in transient assays for high levels of gene activity. These A elements also give high-level position-independent, copy-number dependent expression of a *trans* gene in **transgenic** mice (15).

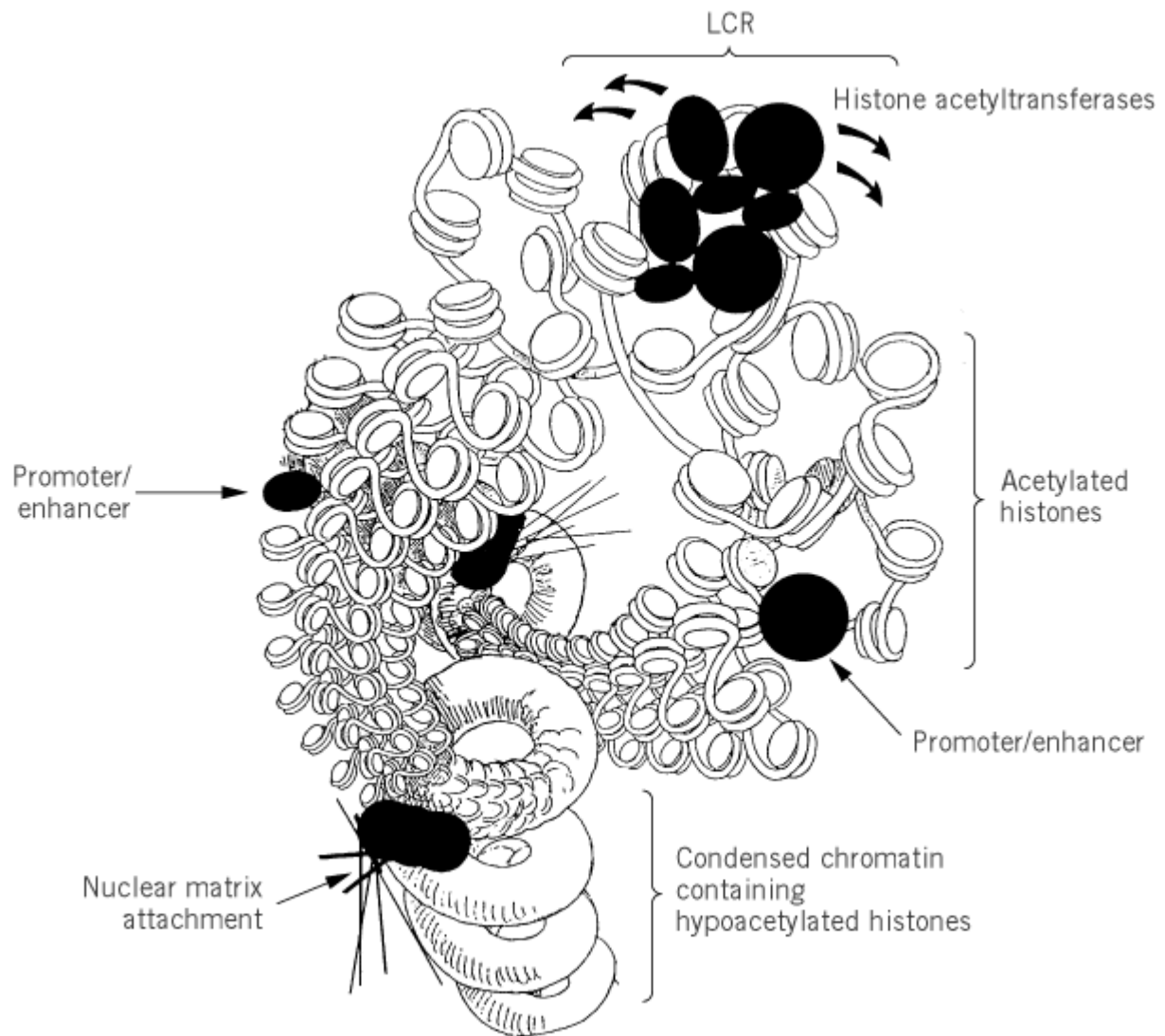
Evidence for the functional specialization of chromosomal domains came initially from cytological analysis of [polytene chromosomes](#). Immunologic staining has revealed the specific distribution of several non-histone proteins, including RNA polymerase II and proteins apparently responsible for generating inactive heterochromatin (HP1, see [Chromocenter](#)). Grossbach and colleagues used antibodies against variants of histone H1 to demonstrate that the localization of particular linker histones is highly specific for individual chromosomal domains within polytene chromosomes (16). In contrast, a similar approach with antibodies raised against HMG14-related proteins led to the general **immunofluorescent** staining of transcriptionally active domains. Turner found that

individual chromosomal domains within polytene chromosomes can be significantly enriched with forms of histone H4 that have particular states of [posttranslational modification](#) (17). This is where a covalent modification such as acetylation of lysine residues is added by an enzyme after the protein has been synthesized. These experiments indicate that specific proteins and modifications can be targeted to particular chromosomal regions.

Additional support for the division of the chromosome into individual functional domains came from the discovery of **locus control regions (LCR)**. DNase I [hypersensitive sites](#) are useful in defining important DNA sequences. DNase I cleavage of the human b-globin locus revealed the existence of four strongly nuclease-sensitive sites located 10 to 20 kbp upstream of the cluster of human b-globin genes (18) (see [Contiguous Genes](#)). These sites, known as locus control regions (LCRs), represent [cis-acting](#) elements that allow genes integrated into a chromosome to be expressed independently of chromosomal position, i.e., [position effects](#) are abolished. Consequently, LCRs allow each copy of a gene integrated in multiple copies to be expressed equivalently within a particular chromosomal domain, so that gene expression is copy-number dependent. When all four DNase I hypersensitive sites comprising the LCR are placed adjacent to [reporter genes](#), they function like [enhancers](#). However, three of the four sites do not function as enhancers in transient expression sites but only do so after incorporation into the chromosome. This suggests that the LCRs may play a special role in stabilizing an accessible chromatin structure distinct from the function of a normal enhancer element (19) (Fig. 1). Several models have been proposed for LCR function, including unraveling of the chromatin fiber, functions similar to normal enhancers, and stabilization of specific nucleoprotein complexes that are functional for transcription only in the presence of the LCR. Evidence consistent with the propagation of an altered chromatin structure comes from transgenic experiments in which LCRs confer hypersensitive sites and general DNase I sensitivity, whereas promoter elements do not. This suggests that LCRs may be able to overcome the inhibitory influences of localized regions of [heterochromatin](#). It has been suggested that LCRs or regions adjacent to them might define the boundaries of an active chromatin domain through attachment to a nuclear matrix or scaffold. At this time there is no positive evidence to suggest that this is the case *in vivo*.

**Figure 1.** A model for locus control region (LCR) function where the LCR defines a functional chromosomal domain. We propose that the LCR recruits histone acetyltransferases that acetylate histones in the functional domain and lead to an expanded and destabilized chromatin fiber. Within this domain are promoter and enhancer sequences that themselves function more effectively when chromatin higher order structure is destabilized. In this model nuclear matrix attachment sites are boundaries to separate the LCR-controlled chromosome domain from condensed chromatin containing hypoacetylated histones.





Schedl and colleagues have used [transposable elements \(P-elements\)](#) that can introduce stable transformants into the germ line of *Drosophila melanogaster* to define a distinct “insulator” element that might help delimit functional chromosomal domains. Using the *hsp70* gene, these investigators defined DNA sequence elements that contain DNase I hypersensitive sites to either side of the *hsp70* genes. These specialized chromatin structural (scs) elements conferred position-independent, **copy-number**-dependent transcription from the white **promoter**, with the exception of one insertion into heterochromatin within the *Drosophila X-chromosome*. Importantly, the scs elements do not behave as scaffold attachment regions that establish a functional separation between the two sequences (20). It has been recognized that the *Drosophila* **suppressor** of hairy-wing protein and the complex that it assembles with the gypsy transposable element have comparable “insulating” properties.

It is important to observe that scs, gypsy, or A elements block enhancer function if the element is between the enhancer and the promoter. This capacity to block activation or silencing effects has led to the current definition of these elements as insulators. Enhancers normally act over very long distances (many kilobases). It is hard to imagine how a short DNA sequence (100 to 200 bp) could inhibit DNA looping, which suggests that enhancers in this case function by an alternate mechanism, perhaps related to chromatin folding or nuclear compartmentalization (see Figure 1). Additional insight comes from experiments suggesting that gene regulation occurs by progressively altering the structure of chromatin domains that are continuous in the chromosome. In the [bithorax complex](#),

three genes ( *Ubx*, *abd-A* and *Abd-B* ), are aligned 5' to 3' in the order in which they are expressed during development and, remarkably, in the structures whose origin they control in an anterior to posterior direction in the fly. It has been suggested that each of these three large (> 20 kbp) segments of DNA are activated in succession. Several mutations in the bithorax complex are consistent with this hypothesis because by removing a boundary between two domains of chromatin, the anterior-posterior boundary between different structures is also removed. In this model, boundary elements prevent the propagation of a particular chromatin modification by functioning as a barrier. Once again, this type of genetic analysis provides strong evidence for the existence of chromosomal domain boundaries.

### Bibliography

1. J. R. Paulson and U. K. Laemmli (1977) *Cell* **12**, 817–828.
2. C. Benyajati and A. Worcel (1976) *Cell* **9**, 393–407.
3. P. R. Cook and I. A. Brazell (1975) *J. Cell Sci.* **19**, 261–279.
4. J. Filipski et al. (1990) *EMBO J.* **9**, 1319–1327.
5. G. L. Chen et al. (1984) *J. Biol. Chem.* **259**, 13560–13566.
6. S. M. Gasser, B. B. Amati, M. E. Cardenas, and J. F. Hofmann (1989) *Int. Rev. Cytol.* **119**, 57–96.
7. J. M. Craig, S. Boyle, P. Perry, and W. A. Bickmore (1977) *J. Cell Sci.*
8. C. D. Lewis and U. K. Laemmli (1982) *Cell* **29**, 171–181.
9. W. C. Earnshaw and M. M. S. Heck (1985) *J. Cell Biol.* **100**, 1716–1725.
10. S. DiNardo, K. Voelkel, and R. Sternglanz (1984) *Proc. Natl. Acad. Sci. USA* **81**, 2616–2620.
11. D. Koshland and A. Strunnikov (1996) *Ann. Rev. Cell Dev. Biol.* **12**, 305–333.
12. R. Strick and U. K. Laemmli (1995) *Cell* **83**, 1137–1148.
13. L. G. Poljak, C. Seum, T. Mattioni, and U. K. Laemmli (1994) *Nucleic Acids Res.* **22**, 4386–4394.
14. K. Zhao, E. Kas, E. Gonzalez, and U. K. Laemmli (1993) *EMBO J.* **12**, 3237–3247.
15. A. Stief, D. M. Winter, W. E. H. Stratling, and A. E. Sippel (1989) *Nature* **341**, 343–345.
16. E. Mohr, L. Trieschmann, and U. Grossbach (1989) *Proc. Natl. Acad. Sci. USA* **86**, 9308–9312.
17. B. M. Turner, L. Franchi, and H. Wallace (1990) *J. Cell Sci.* **96**, 335–346.
18. F. Grosveld, G. B. von Assendelft, D. R. Greaves, and G. Kollias (1987) *Cell* **51**, 975–985.
19. T. Jenuwein, W. C. Forrester, R. G. Qui, and R. Grosschedl (1993) *Genes Dev.* **7**, 2016–2032.
20. R. Kellum and P. Schedl (1991) *Cell* **64**, 941–950.
21. V. Orlando and R. Paro (1993) *Cell* **75**, 1187–1198.

### Suggestion for Further Reading

22. A. P. Wolffe (1998) *Chromatin: Structure and Function*, 3rd ed., Academic Press, London.

### Domain, Protein

Protein domains, sometimes referred to as protein folds or modules, are thought to be the building blocks of [protein structure](#), function, and [evolution](#). A domain is generally defined as that part of a protein's structure that can fold independently of the remainder of the protein. The [tertiary structure](#)

of a protein's [polypeptide chain](#) is assembled into one or more protein domains, each domain usually incorporating between 50 and 150 [amino acid](#) residues. Proteins having fewer than 100 residues in their polypeptide chains are usually single-domain proteins. Multidomain proteins incorporate more than one domain in their structure and have long polypeptide chains. Generally, the individual domains of a multidomain protein are formed from consecutive stretches of sequence that are connected by a single short length of linking polypeptide chain. These linking regions may be the source of [conformational](#) flexibility in the tertiary structure of the protein, and the links can be susceptible to cleavage by [proteinases](#) (thereby producing stable protein domain fragments).

Individual domains within a protein structure often have distinct roles, including catalytic, regulatory, binding, recognition, or oligomerization functions. Protein domains are categorized into different structural classes ([1](#), [2](#)) based on their **secondary structure**, including (a) those containing mostly [a-helix](#) (all-a), (b) those containing mostly [b-strands](#) arranged in antiparallel [beta-sheet](#) (all-b), (c) those containing mostly alternating a-helix and b-strand elements (a/b), and (d) those containing a-helix and b-strand elements that are separated rather than alternating in the polypeptide chain (a + b). Some also classify small proteins (<10 kDa) that are stabilized by [disulfide bonds](#) or metal ion interactions as a separate class. Recently a new class of protein domain, called the parallel [b-helix](#), that contains only parallel b-strands has been identified ([3](#)). Classification of domain structure beyond this level is complex, but [databases](#) that define different structural types of protein domains have been developed ([4-7](#)) (see [Structure Databases](#)).

[See also [Protein Structure](#).]

### Bibliography

1. J. S. Richardson (1981) *Adv. Protein Chem.* **34**, 167–339.
2. M. Levitt and C. Chothia (1976) *Nature* **261**, 552–558.
3. M. D. Yoder, N. T. Kneen, and F. Journak (1993) *Science* **260**, 1503–1507.
4. A. G. Murzin, S. E. Brenner, T. Hubbard, and C. Chothia (1995) *J. Mol. Biol.* **247**, 536–540 and <http://scop.mrc-lmb.cam.ac.uk/scop>
5. C. A. Orengo et al. (1997) *Structure* **5**, 1093–1108 and <http://www.biochem.ucl.ac.uk/bsm/cath>
6. A. S. Siddiqui and G. J. Barton (1995) *Protein Sci.* **4**, 872–884 and (site currently unavailable)
7. R. Sowdhamini, S. D. Rufino, and T. L. Blundell (1996) *Folding Des.* **1**, 209–220 and <http://www-cryst.bioc.cam.ac.uk/~ddbbase>

### Suggestions for Further Reading

8. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
9. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, New York.
10. D. G. Hardie and J. R. Coggins (1986) *Multidomain Proteins—Structure and Evolution*, Elsevier, Amsterdam.

### Domain Shuffling

[Proteins](#) are often composed of one or more functional **domains**. In particular, when a certain protein possesses more than one domain, it is called a [mosaic protein](#). Domain shuffling is a phenomenon

where the **gene** segments coding for functional domains are shuffled between different genes during evolution. The boundaries between different domains may or may not correspond to the [intron–exon](#) boundaries of the encoding gene. In fact, in certain cases a single domain is encoded by more than one exon, whereas an exon contains a number of functional domains. Therefore, the “[exon shuffling](#)” theory is not necessarily the same as the “domain shuffling” theory. Although there are many cases where domain shuffling must have taken place, it is not known what molecular event causes this, and this is important to determine. For more details, see [Exon Shuffling](#).

## Dorsal/Cactus Complex

In the *Drosophila melanogaster* embryo, dorsoventral polarity is determined by a concentration gradient of the [transcription factor](#) Dorsal in the [nucleus](#) (1-3). Prior to axis formation, Dorsal is held in the cytoplasm by an inhibitor, Cactus (4-6). During early embryonic development, **proteolysis** of Cactus in response to a signal transduced on the ventral side of the syncytial embryo releases Dorsal (7-9). Free Dorsal translocates into ventral nuclei, where it directs expression of ventral-specific **genes** and represses expression of dorsal-specific genes (10). In this manner, the spatially graded disruption of the Dorsal/Cactus complex establishes the dorsoventral body axis (11).

Dorsal and Cactus are **translated** from [messenger RNAs](#) synthesized in the ovary and form a complex distributed uniformly in the unfertilized [egg](#). Following fertilization, 14 rounds of replication and nuclear division occur prior to cellularization, so that axis formation occurs in the context of a [syncytium](#) containing hundreds, and eventually thousands, of nuclei. Within this syncytium, spatially regulated **protein degradation** of Cactus leads to a gradient of Dorsal nuclear localization, with the highest levels of nuclear Dorsal along the ventral midline and lowest levels in the dorsal half of the embryo.

Dorsal and Cactus were identified on the basis of **maternal-effect** mutations that disrupt embryonic pattern formation (12). Embryos generated by females lacking Dorsal have a dorsalized phenotype; there is no nuclear Dorsal and hence only dorsal-specific tissues are formed. Females homozygous for a mutation that inactivates Cactus generate embryos with an opposite, ventralized phenotype, in which Dorsal enters nuclei on both the ventral and dorsal sides of the embryo (4). Ten additional genes mutate to a dorsalized phenotype; these loci encode components of the signaling pathway that regulates release of Dorsal from the Dorsal/Cactus complex.

The [signal transduction](#) pathway that directs the asymmetric nuclear translocation of Dorsal is initiated by Pipe, Nudel, and Windbeutel, the products of three genes expressed in the somatic cells of the ovary (13). A ventrally localized signal is then amplified by an extraembryonic [serine proteinase](#) cascade involving the Gastrulation Defective, Snake, and Easter proteins (14-16). The result is a proteolytically activated form of Spätzle, the activating ligand for the transmembrane receptor **Toll** (17). Transduction of a signal from activated Toll to the Dorsal/Cactus complex requires a scaffolding protein, Tube, and a protein kinase, Pelle (18, 19). Dorsal binds specifically to both Tube and Pelle, which also bind to one another (20, 21). Both Dorsal and Cactus appear to undergo **phosphorylation** in response to Toll activation (8, 22, 23). Phosphorylation of Cactus is critical for signal transduction; the role of Dorsal phosphorylation in signaling is not yet clear.

Once within ventral nuclei, Dorsal serves to activate expression of ventral-specific genes (eg, *twist*) and repress expression of dorsal-specific genes, e.g., *zen* (24). By themselves, Dorsal dimers activate [transcription](#); repression requires binding to a [corepressor](#), the Groucho protein (25). Dorsal function is not uniform across the ventral half of the embryo. Instead, the gradient of nuclear localization is translated as a series of threshold responses into stripes of gene activation and cell fate determination

(26).

Whereas maternally expressed Dorsal, Cactus, Toll, Tube, and Pelle function in embryonic patterning, zygotically expressed forms of these proteins mediate an [immune response](#) in the fat body, the *Drosophila* organ most similar to the mammalian liver (27). The immune response, which directs expression of **peptide antibiotics** upon infection, also involves at least two other proteins closely related in sequence to Dorsal. These proteins, DIF (Drosophila immunity factor) and Relish, appear to be regulated in a manner similar to that observed for Dorsal (28, 29).

The immune function of the Dorsal/Cactus complex has been broadly conserved during evolution; NF- $\kappa$ B and I $\kappa$ B, homologues of Dorsal and Cactus, play a critical role in the mammalian acute phase response, as do counterparts of Toll and Pelle (30-32). Activation of NF- $\kappa$ B in response to stimulation of cells with **cytokines**, such as **tumor necrosis factor- $\alpha$**  and **interleukin-1**, leads to activation of genes that regulate both the inflammatory and immune responses. Recent data suggest that the NF- $\kappa$ B/I $\kappa$ B complex, like the Dorsal/Cactus complex, may also function in pattern formation (33, 34).

The central feature of the Dorsal protein is a multifunctional region termed the Rel homology **domain**. This region mediates the binding of Dorsal to itself, to Cactus, to other specific protein pairing partners, and to DNA (6, 20, 35, 36). The Rel homology domain also contains a nuclear localization signal (NLS) that drives translocation of free Dorsal into the nuclei of the syncytial blastoderm (see [Nuclear Import, Export](#)). It is thought that Cactus retains Dorsal outside nuclei, not by actively tethering Dorsal in the cytoplasm, but rather by occluding this NLS and thereby blocking nuclear translocation (36).

The Dorsal/Cactus complex consists of a dimer of the *Drosophila* transcription factor Dorsal bound to a monomer of its cytoplasmic inhibitor, Cactus (37). Complex formation requires an interaction between the Dorsal Rel domain and a set of six tandemly arrayed motifs in Cactus. These motifs, termed [ankyrin repeats](#), are found among proteins of widely varying function and have been shown to serve as sites of [protein-protein interaction](#). Interaction with Dorsal requires all of the ankyrin repeats in Cactus (5, 6). Carboxy-terminal to the ankyrin repeats is a **PEST domain**. PEST domains, rich in the amino acid residues proline, glutamine, serine, and threonine, are found in many proteins whose activity is regulated by proteolysis.

Two regions of Cactus govern its susceptibility to degradation (7-9). The PEST domain influences the overall stability of Cactus, whereas a motif amino-terminal to the ankyrin repeats is specifically required for degradation in response to Toll-mediated signaling. The activities of both sites appear to be regulated by phosphorylation (8, 9, 38). Similar sites govern the proteolysis of the I $\kappa$ B proteins, for which it has been directly demonstrated that phosphorylation of serine residues in the amino-terminal motif triggers a rapid and efficient degradation mediated by the [proteasome](#) (39-42).

## Bibliography

1. S. Roth, D. Stein, and C. Nüsslein-Volhard (1989) *Cell* **59**, 1189–1202.
2. C. A. Rushlow, K. Han, J. L. Manley, and M. Levine (1989) *Cell* **59**, 1165–1177.
3. R. Steward (1989) *Cell* **59**, 1179–1188.
4. S. Roth, Y. Hiromi, D. Godt, and C. Nüsslein-Volhard (1991) *Development* **112**, 371–388.
5. R. Geisler, A. Bergmann, Y. Hiromi, and C. Nüsslein-Volhard (1992) *Cell* **71**, 613–621.
6. S. Kidd (1992) *Cell* **71**, 623–635.
7. M. P. Belvin, Y. Jin, and K. V. Anderson (1995) *Genes & Development* **9**, 783–793.
8. M. Reach, R. L. Galindo, P. Towb, J. L. Allen, M. Karin, and S. A. Wasserman (1996) *Dev. Biol.* **180**, 353–364.
9. A. Bergmann, D. Stein, R. Geisler, S. Hagenmaier, B. Schmid, N. Fernandez, B. Schnell, and C. Nüsslein-Volhard (1996) *Mech. Dev.* **60**, 109–123.

10. J. Rusch and M. Levine (1996) *Curr. Opin. Genet. Dev.* **6**, 416–423.
11. D. Morisato and K. V. Anderson (1995) *Annu. Rev. Genet.* **29**, 371–399.
12. K. V. Anderson and V. C. Nusslein (1984) *Nature* **311**, 223–227.
13. D. Stein, S. Roth, E. Vogelsang, and C. Nüsslein-Volhard (1991) *Cell* **65**, 725–735.
14. K. D. Konrad, T. J. Goralski, and A. P. Mahowald (1988) *Dev. Biol.* **127**, 133–142.
15. Y. S. Jin and K. V. Anderson (1990) *Cell* **60**, 873–881.
16. C. L. Smith and R. DeLotto (1994) *Nature* **368**, 548–551.
17. D. S. Schneider, Y. Jin, D. Morisato, and K. V. Anderson (1994) *Development* **120**, 1243–1250.
18. A. Letsou, S. Alexander, K. Orth, and S. A. Wasserman (1991) *Proc. Natl. Acad. Sci. USA* **88**, 810–814.
19. C. A. Shelton and S. A. Wasserman (1993) *Cell* **72**, 515–525.
20. D. N. Edwards, P. Towb, and S. A. Wasserman (1997) *Development* **124**, 3855–3864.
21. E. A. Drier and R. Steward (1997) *Semin. Cancer Biol.* **8**, 83–92.
22. A. M. Whalen and R. Steward (1993) *J. Cell Biol.* **123**, 523–534.
23. S. K. Gillespie and S. A. Wasserman (1994) *Mol. Cell. Biol.* **14**, 3559–3568.
24. J. Jiang, C. A. Rushlow, Q. Zhou, S. Small, and M. Levine (1992) *EMBO J.* **11**, 3147–3154.
25. T. Dubnicoff, S. Valentine, G. Chen, T. Shi, J. A. Lengyel, Z. Paroush, and A. J. Courey (1997) *Genes Dev.* **22**, 2952–2957.
26. J. Jiang and M. Levine (1993) *Cell* **72**, 741–752.
27. B. Lemaitre, N. Emmanuelle, L. Michaut, J.-M. Reichhart, and J. A. Hoffmann (1996) *Cell* **86**, 973–983.
28. Y. T. Ip, M. Reach, Y. Engstrom, L. Kadalayil, H. Cai, S. Gonzalez-Crespo, K. Tatei, and M. Levine (1993) *Cell* **75**, 753–763.
29. M. S. Dushay, B. Asling, and D. Hultmark (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10343–10347.
30. I. M. Verma, J. K. Stevenson, E. M. Schwarz, D. Van Antwerp, and S. Miyamoto (1995) *Genes Dev.* **9**, 2723–2735.
31. Z. Cao, W. J. Henzel, and X. Gao (1996) *Science* **271**, 1128–1131.
32. R. Medzhitov, P. Preston-Hurlburt, and C. A. Janeway Jr. (1997) *Nature* **388**, 394–397.
33. P. B. Bushdid, D. M. Brantley, F. E. Yull, G. L. Blaeuer, L. H. Hoffman, L. Niswander, and L. D. Kerr (1998) *Nature* **392**, 615–618.
34. Y. Kanegae, A. T. Tavares, J. C. Izpisua Belmonte, and I. M. Verma (1998) *Nature* **392**, 611–614.
35. K. Tatei and M. Levine (1995) *Mol. Cell. Biol.* **15**, 3627–3634.
36. S. Govind, E. Drier, L. H. Huang, and R. Steward (1996) *Molecular & Cellular Biology* **16**, 1103–14.
37. K. Isoda and C. Nüsslein-Volhard (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5350–5354.
38. Z.-P. Liu, R. L. Galindo, and S. A. Wasserman (1997) *Genes Dev.* **11**, 3413–3422.
39. K. Brown, S. Gerstberger, L. Carlson, G. Franzoso, and U. Siebenlist (1995) *Science* **267**, 1485–1488.
40. J. DiDonato, F. Mercurio, C. Rosette, J. Wu-Li, H. Suyang, S. Ghosh, and M. Karin (1996) *Mol. Cell. Biol.* **16**, 1295–1304.
41. Z. Chen, J. Hagler, V. J. Palombella, F. Melandri, D. Scherer, D. Ballard, and T. Maniatis (1995) *Genes Dev.* **9**, 1586–1597.
42. E. B. Traenckner, H. L. Pahl, T. Henkel, K. N. Schmidt, S. Wilk, and P. A. Baeuerle (1995) *EMBO J.* **14**, 2876–2883.

### Suggestions for Further Reading

43. D. Stein, S. Roth, E. Vogelsang, and C. Nüsslein-Volhard (1991) *Cell* **65**, 725–735.
44. M. P. Belvin, Y. Jin, and K. V. Anderson (1995) *Genes Dev.* **9**, 783–793.
45. B. Lemaitre, N. Emmanuelle, L. Michaut, J.-M. Reichhart, and J. A. Hoffmann (1996) *Cell* **86**, 973–983.
46. J. Rusch and M. Levine (1996) *Curr. Opin. Genet. Dev.* **6**, 416–423.
47. D. N. Edwards, P. Towb, and S. A. Wasserman (1997) *Development* **124**, 3855–3864.

### Dorsoventral Polarity

The establishment of dorsoventral polarity in the *Drosophila* [embryo](#) has been the subject of extensive genetic, developmental, and biochemical analyses. At least three different signaling pathways are required for proper dorsoventral polarity, two that are under **maternal control** and a third pathway that requires expression of the genes in the embryo itself. All of these signaling pathways must operate correctly for normal dorsoventral patterning of the embryo.

During oogenesis, the oocyte is surrounded by specialized somatic cells called follicle cells. The follicle cells secrete the various layers of the chorion, or egg shell, during oogenesis. The chorion differs in structure along the dorsoventral axis, showing that the follicle cells have dorsoventral positional information. In addition, the oocyte itself, as well as the resulting embryo, has a dorsoventral polarity. The coordination of dorsoventral polarity between the follicle cells and the growing oocyte involves a signaling pathway of at least 13 different genes and appears to involve signaling from the oocyte to the follicle cells, and then signaling back from the follicle cells to the oocyte (1). That the oocyte sends a dorsalizing signal to the follicle cells is shown by mutations in several genes, including the [gurken](#) gene. The *gurken* gene encodes a [transforming growth factor](#) a (TGF $\alpha$ )-like protein that must be expressed in the oocyte in order for the follicle cells to make dorsal chorion structures (2). Loss of *gurken* function in the oocyte leads to ventralization of the follicle cells. That the follicle cells in turn send a dorsalizing signal is shown by mutations in the *torpedo* gene. Loss of *torpedo* function in the follicle cells causes the oocyte (and embryo) to be ventralized (3). The *torpedo* gene encodes an **EGF receptor**. Human TGF $\alpha$  has been shown to bind and activate an ectodermal growth factor (EGF) receptor (4), which suggests that *torpedo* proteins in the follicle cell membranes are receptors for *gurken* proteins secreted by the oocyte.

Signaling back from the follicle cells to the oocyte involves another pathway that is conserved in the human immune system. This pathway is the **dorsal** pathway in *Drosophila*, and it includes at least 12 maternally acting genes (5). Mutations in 11 of these genes cause the oocyte and resulting embryo to be dorsalized. These are the dorsal group of mutations, because they all seem to function to activate the *dorsal* protein in the embryo (6). The *dorsal* protein is homologous to the NF- $\kappa$ B/rel proteins required for signaling in the human immune system (7). Mutations in the *cactus* gene cause the oocyte and embryo to be ventralized. *Cactus* encodes a homologue of the I $\kappa$ B protein (8), which binds and inhibits NF- $\kappa$ B/rel. *Cactus* appears to bind and inhibit *dorsal* in *Drosophila* (9). The nuclear localization of *dorsal* protein in the embryo initiates the third signaling pathway in dorsoventral patterning in the *Drosophila* embryo.

The third signaling pathway required for dorsoventral patterning in the *Drosophila* embryo is encoded by zygotically active genes, in contrast to the maternal requirements for the first two signaling pathways. This zygotic signaling pathway appears to use the TGF $\beta$  homologue

[decapentaplegic](#) (*dpp*) as the dorsal signal. Loss of *dpp* activity in the embryo causes ventralization (10). Another TGF $\beta$  homologue, *screw*, also appears to be required zygotically for dorsalization of the embryo (11). *Dpp* and *screw* may act synergistically in dorsalizing the embryo, because injection of *dpp* [messenger RNA](#) can partially suppress the phenotype of *screw* mutants, and antimorphic *screw* mutations enhance loss of *dpp* activity (11). Three other genes required for dorsalizing the embryo encode homologues of subunits of vertebrate TGF $\beta$  receptors. These are the genes *punt*, *thick veins*, and *saxophone* (1). Activation of these receptors by the *dpp* and *screw* proteins appears to activate [transcription](#) of the *zerknüllt* (*zen*) gene. *Zen* is a **homeobox**-containing [transcription factor](#) that is required for dorsalization of the embryo (12). Thus, dorsoventral polarity in the *Drosophila* embryo begins with a signal from the developing oocyte to the follicle cells, a return signal from the follicle cells to the oocyte confirming dorsoventral polarity, and the final activation of at least one transcription factor in dorsal cells of the embryo.

## Bibliography

1. D. Morisato and K. V. Anderson (1995) *Annu. Rev. Genet.* **29**, 371–399.
2. F. S. Neuman-Silberberg and T. Schüpbach (1993) *Cell* **75**, 165–174.
3. J. V. Price, R. J. Clifford, and T. Schüpbach (1989) *Cell* **56**, 1085–1092.
4. R. Derynck, A. B. Roberts, M. E. Winkler, E. Y. Chen, and D. V. Goeddel (1984) *Cell* **38**, 287–297.
5. Reference 1, p. 379.
6. S. Roth, D. Stein, and C. Nüsslein-Volhard (1989) *Cell* **59**, 1189–1202.
7. R. Steward (1987) *Science* **238**, 692–694.
8. R. Geisler, A. Bergmann, Y. Hiromi, and C. Nüsslein-Volhard (1992) *Cell* **71**, 613–621.
9. S. Kidd (1992) *Cell* **71**, 623–635.
10. V. Irish and W. M. Gelbart (1987) *Genes Dev.* **1**, 868–879.
11. K. Arora, M. S. Levine, and M. B. O'Connor (1994) *Genes Dev.* **8**, 2588–2601.
12. C. Rushlow, H. Doyle, T. Hoey, and M. Levine (1987) *Genes Dev.* **1**, 1268–1279.

## Suggestions for Further Reading

13. D. Morisato and K. V. Anderson (1995) Signaling pathways that establish the dorsal–ventral pattern of the *Drosophila* embryo. *Annu. Rev. Genet.* **29**, 371–399.
14. A. González-Reyes, H. Elliott, and D. St. Johnston (1995) Polarization of both major body axes in *Drosophila* by *gurken*–*torpedo* signaling. *Nature* **375**, 654–658.
15. R. P. Ray and T. Schüpbach (1996) Intercellular signaling and the polarization of body axes during *Drosophila* oogenesis. *Genes Dev.* **10**, 1711–1723.
16. S. A. Holley and E. L. Ferguson (1997) Fish are like flies are like frogs: conservation of dorsal–ventral patterning mechanisms. *Bioessays* **19**, 281–284.

## Dosage Compensation Effect

Dosage compensation is a mechanism by which the activity of X-linked **genes** is made equal in the two sexes of the type XX and XY, in which one [sex](#) has two [X-chromosomes](#) and the other only one. In placental mammals, marsupials, and some monotremes, compensation is achieved by the inactivation of one X-chromosome in **somatic cells** of females (1). During embryonic development of the [mouse](#), both X-chromosomes are active, and differentiation between the active form X<sub>a</sub> and



the inactive form Xi occurs in the primitive **ectoderm**. Once inactivation has been initiated at this stage, the same X-chromosome remains inactive in the descendants of each cell after **mitosis** throughout the life of the mouse. Reactivation of chromosome Xi in females occurs at the time of **meiosis**, so that both X-chromosomes are active in oocytes. In male **germ cells**, the single X-chromosome becomes inactive at the late spermatogonial stage.

The **chromatin** of chromosome Xi is in the condensed state, and its **DNA replication** begins later than in Xa or in **autosomes**. The cytosine bases of **CpG Islands** near the 5' **promoter** regions of genes are heavily methylated on chromosome Xi, but not on Xa (2) (see **Methylation, DNA**). Differential hypersensitivity to digestion by **deoxyribonuclease** indicates altered binding of proteins to the DNA and different packaging of **nucleosomes** in the two X-chromosomes (3).

During studies of translocations between the X- and autosome chromosomes in mice, it was observed that only one of the two X-chromosomes undergoes inactivation which spreads from the inactive segment into the attached autosomal material. It appears as if an inactivation center existed on the X-chromosome, which becomes blocked on one X-chromosome (which becomes Xa) only at the time inactivation is initiated. The location of the inactivation center has been accurately mapped in human and mouse from studies of chromosomal translocations and deletions. A gene, called *XIST* in humans and *Xist* in the mouse, maps in the region of the inactivation center and is expressed by Xi but not by Xa (4). In the male, *Xist* is expressed only in the testis, suggesting a role in inactivation of the single X-chromosome in male germ cells. No protein product of *Xist* or *XIST* has yet been identified.

In *Drosophila melanogaster*, where X-chromosome is not inactivated, both X-chromosomes are transcribed in females. Compensation occurs by up-regulation and a doubled rate of **transcription** of the single X-chromosome in males. As a result, the level of transcription of the genes carried by the X-chromosome is the same in both sexes. The increase in the rate of transcription in males is caused by binding of **transcription factors**, called *m<sup>sl</sup>-1*, *m<sup>sl</sup>-2*, and *m<sup>sl</sup>-3*, to **enhancer** elements on the X-chromosome. These transcription factors appear to act at the level of the chromatin structure and have been extensively studied in the polytene salivary gland chromosomes, where they bind codependently to the same set of sites along the male X-chromosome. Molecular characterization of the protein coded by *m<sup>sl</sup>-2* has to a great extent solved the question of how *m<sup>sl</sup>*-mediated dosage compensation is restricted to males. Cloning and molecular analyses of the *m<sup>sl</sup>* genes have substantiated the proposal that the MSL proteins function as a multimeric complex to mediate dosage compensation (5), (6).

In the **nematode** *Caenorhabditis elegans*, equalization of transcription seems to occur by down-regulation in XX animals.

#### Bibliography

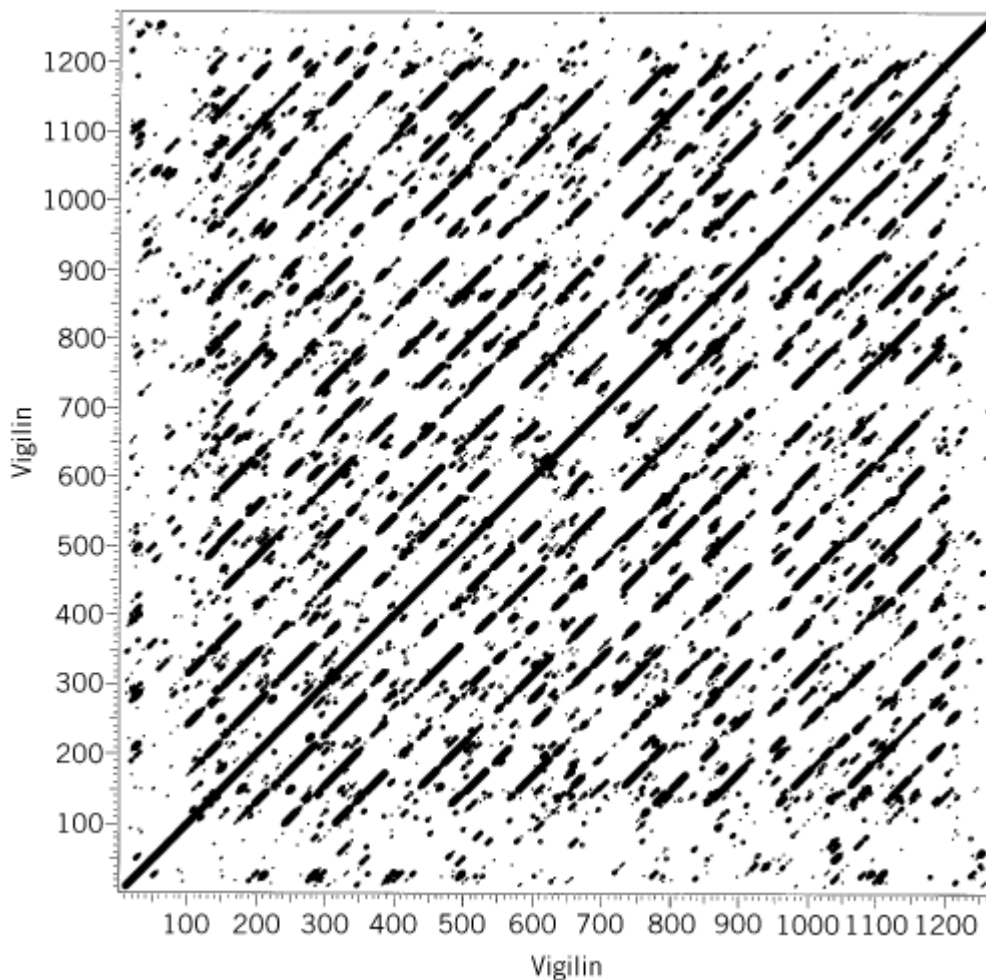
1. S. M. Gartler and A. D. Riggs (1983) *Ann. Rev. Genet.* **17**, 155–190.
2. A. D. Riggs (1990) *Phil. Trans. R. Soc. London* **B326**, 285–297.
3. G. P. Pfeifer and A. D. Riggs (1991) *Genes Dev.* **5**, 1102–1111.
4. M. F. Lyon (1991) *Trends Genet.* **7**, 69–70.
5. S. Zhou et al. (1995) *EMBO. J.* **14**, 2884–2895.
6. G. J. Bashaw and B. S. Baker (1996) *Curr. Opin. Genet. Dev.* **6**, 496–501.

#### Dot Plot

The dot plot, introduced by Gibbs and McIntyre (1), is a two-dimensional visual comparison method, useful for revealing regions of similarity between pairs of sequences or structures (2). Typically, the *X* axis represents one sequence and the *Y* axis represents the other. Every point in the plot is then scored for the similarity of the residue in sequence *X* with the similarity of the residue in sequence *Y*, and the point is plotted if the residues are scored as similar in a residue substitution matrix. Diagonal lines in the plots then indicate regions where consecutive residues are scored as similar. Insertions and deletions ([Indels](#)) are revealed by termination of a diagonal segment, followed by initiation of a new but offset diagonal. To improve the signal-to-noise ratio, it is usual to sum the scores in short windows spanning each residue: overall, fewer points are plotted, but divergent matches with rather few identities can be revealed. Cutoffs determining which points are plotted may be set using a variety of heuristic or probabilistic measures. The “double matching probability,” introduced by McLachlan (3), is an estimate of the likelihood that the given score would arise by chance in two infinitely long sequences and is widely used to provide the scale for plotting points. The user can then choose the significance level cutoff to determine which points to plot.

Dot plots can reveal recurring similarities within a sequence by self-comparison. In this case, repeated sequences appear as a series of off-center partial diagonals. [Tandem Repeats](#) result in a regular set of diagonals, consecutively decreasing in size by one repeat, moving away from the central diagonal (Fig. 1). Dispersed repeats also provide dispersed diagonals, one between each repeat element.

**Figure 1.** Dot plot showing 15 tandemly repeated KH domains in the chicken vigilin sequence. The larger the dots, the higher the matching segment score. A break in the diagonals spanning position 940 indicates a ~40-residue insertion in one of the domains.



Dot-plot self-comparisons can be used to reveal elements of double-stranded secondary structure in folded RNAs. In the simplest case, points are plotted whenever a pair of residues are able to **base pair**. More sensitive dot-plot algorithms use energy rules based on base pairing, base stacking, bulging, and looping, allied to dynamic programming algorithms (4). Therefore diagonals indicate runs of complementary residues able to form double helices. Looped out *bases*, which often occur in *RNA secondary structure*, are indicated by short offsets in consecutive diagonals.

Another variety of dot plot (5) may be used to compare two protein [tertiary structures](#) for regions of similarity, although interpretation is more complex. In this case, the plotted points typically represent distances between Ca carbons, so the plots are also known as distance plots or [contact maps](#). Characteristic recurring patterns indicate **secondary structure** elements such as [a-helices](#) and [b-strands](#). Long-range contacts are distributed irregularly in the plot. A useful automatic method to detect structural similarities by analyzing and superposing sets of self-contact plots has been implemented in the program Dali available on the WWW (6).

Dot plots are sometimes used in other ways for presenting sequence or structure information. For example, contact plots are a convenient way to summarize the residue contacts revealed in a structural investigation by nuclear magnetic resonance ([NMR](#)).

#### Bibliography

1. A. J. Gibbs and G. A. McIntyre (1970) *Eur. J. Biochem.* **16**, 1–11.
2. R. Staden (1982) *Nucleic Acids Res.* **10**, 2951–2961.
3. A. D. McLachlan (1971) *J. Mol. Biol.* **185**, 39–49.

4. J. A. Jaeger, D. H. Turner, and M. Zucker (1990) *Methods Enzymol.* **183**, 281–306.
5. D. C. Phillips (1970) *Biochem. Soc. Symp.* **31**, 11–28.
6. L. Holm and C. Sander (1993) *J. Mol. Biol.* **233**, 123–138.

### Suggestion for Further Reading

7. D. Sankoff and J. B. Kruskal (1984) *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Addison-Wesley, Reading, MA.

## Double Minute Chromosome

Double minute chromosomes (DMs) are generated as a consequence of DNA **amplification**, which is a large increase in the numbers of a specific gene and its adjacent DNA on either side, because of their selective replication. Some types of amplification result from specific developmental requirements, such as the synthesis of multiple copies of the large ribosomal RNA genes in *Xenopus laevis* **oocytes** (1). This amplification of ribosomal RNA genes meets the enormous needs for ribosomal synthesis during *Xenopus* oogenesis. More typically, DNA amplification of the type that leads to the appearance of DMs occurs as part of an environmental **stress response** in cells in culture.

[Tissue culture](#) cells exposed to a toxic substance, such as those used in cancer chemotherapy, often eventually become resistant to the substance. One way of achieving a resistant state is to amplify a gene whose product metabolizes the toxic substance into a safer nontoxic form. A typical example is the amplification of the gene encoding [dihydrofolate reductase](#) (DHFR), an enzyme involved in thymidine synthesis, in response to methotrexate, which is an inhibitor of DHFR synthesis (see [Aminopterin, Methotrexate, Trimethoprim, and Folic Acid](#)). Amplification of the DHFR gene can also be stimulated in response to other environmental conditions, such as UV light and hydroxyurea (2).

Two visible manifestations of the amplification process occur at the chromosomal level. One is the appearance of extended regions in the chromosome known as homogeneously staining regions (HSRs), whose staining characteristics differ from those of normal chromosomal. These HSRs are the sites of DNA amplification (3). The second visible change in chromosomal material is the appearance of DMs that contain the same amplified DNA, but have become extrachromosomal. The relationship between DMs and HSRs is not completely clear, but it is likely that DMs are unstable precursors of HSRs. Human primary tumors contain DMs 90% of the time, 7% contain HSRs, and 3% contain both DMs and HSRs. When cell lines are grown for extended periods of time under drug-resistant conditions, greater numbers of HSRs appear and fewer DMs are found.

DMs are basically **acentric** extrachromosomal elements that are small fragments of chromosomes. Each DM contains between 1 to 2 megabase pairs of DNA and can replicate itself. Therefore DMs function as [minichromosomes](#). The physical organization of the DNA is that of a **supercoiled** circle (4). Because DMs are acentric they are progressively lost during **cell division**, unless they replicate multiple times in any cell cycle. In general they segregate randomly between daughter cells at **mitosis**, leading to unequal numbers in daughter cells. As many as 20 DMs containing DHFR DNA are found in some methotrexate-resistant **cell lines**. **Oncogenes** might also be present in the amplified DNA of DMs and HSRs. For example, the *c-Ki-ras* oncogene is amplified more than 30-fold in DMs and HSRs in certain adrenocortical tumor cell lines (5). Double minute chromosomes

are important indicators of abnormal cellular behavior, and they may well reflect an advanced stage of malignancy.

### Bibliography

1. D. D. Brown and I. Dawid (1968) *Science* **160**, 272–280.
2. J. K. Cowell (1982) *Ann. Rev. Genet.* **16**, 21–45.
3. J. E. Looney and J. L. Hamlin (1987) *Mol. Cell Biol.* **7**, 569–579.
4. E. P. Garvey and D. V. Santi (1986) *Science* **233**, 535–538.
5. M. Schwab et al. (1983) *Nature* **303**, 497–501.

### Suggestion for Further Reading

6. M. S. Clark and W. J. Wall (1996) *Chromosomes. The Complex Code*, Chapman and Hall, London.

## Downstream

**Nucleotide sequence** elements located downstream from protein-coding sequences are defined by the direction of [transcription](#). They are primarily involved in termination of transcription and in [messenger RNA](#) (mRNA) processing, in particular, [polyadenylation](#). Prokaryotic terminators typically have runs of T at the termination point, preceded by complementary symmetrical sequences, which provide the potential for formation of hairpin structures at the ends of the mRNA (1). The hairpins presumably cause transcription to pause, a phenomenon known as **attenuation**, before the complete arrest.

The mechanism of transcription termination in eukaryotes is not fully understood (2, 3). There are usually many alternative locations for the termination points for the same gene, within a downstream region as large as several thousand bases. It appears that, in a manner similar to that in prokaryotic terminators, transcription elongation in eukaryotes first slows down, due to formation of secondary structure in the mRNA, and then terminates within a nearby U-rich sequence (3). At the noncoding 3'-end of newly formed eukaryotic RNA transcripts several polyadenylation signals usually exist with the **consensus sequence** AAUAAA. One of these sites (but not always the same in different rounds of transcription of the same gene) is somehow selected for the RNA cleavage, which occurs about 20 bases downstream from the site, and for subsequent polyadenylation of 3'-end of the mRNA. Polyadenylation of mRNA is also known in prokaryotes, although it is not as site specific or extensive as in eukaryotes (4).

The eukaryotic downstream sequences—in particular those located within 3'-ends of the RNA transcripts—are also involved in many external functions. For example, they can act in *trans* and control the efficiency of the, stability, and compartmentalization of an mRNA (5).

### Bibliography

1. K. S. Wilson and P. von Hippel (1995) *Proc. Natl. Acad. Sci. USA* **92**, 8793–8797.
2. N. J. Proudfoot (1988) *Trends Biochem. Genet.* **14**, 105–110.
3. O. Resnekov et al. (1988) *Gene* **72**, 91–104.
4. N. Sarkar (1997) *Ann. Rev. Biochem.* **66**, 173–197.

5. C. J. Decker and R. Parker (1995) *Curr. Opin. Cell Biol.* **7**, 386–392.

### Suggestion for Further Reading

6. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.

## *Drosophila* EGF Receptor

Processes of cell–cell interaction are the key to [development](#) of distinct cell and tissue fates in embryonic and postembryonic stages. Focusing on signaling pathways triggered by transmembrane receptors provides molecular clues to these events. In addition, genetic approaches allow the identification of signaling elements in the pathway and provide the hierarchy. For every pathway triggering [morphogenesis](#), it is necessary to understand how signals triggering the pathway are initiated in a defined point in space and time, and how these signals spread to activate the receptor in neighboring cells. It is also crucial to know how different levels of signaling by the receptor are translated to distinct patterns of gene expression and, finally, which mechanisms are employed to restrict or terminate signaling.

The *Drosophila* [epidermal growth factor](#) (EGF) receptor is a case in point. Only a single receptor of this family has been identified in the *Drosophila* [genome](#). This receptor carries out an extremely diverse array of functions throughout development. In all cases, a conserved signaling “cassette” is utilized. The signaling elements identified thus far constitute a robust signaling system that leads to reproducible patterning, in spite of possible variations in the levels of individual components.

### 1. Components

#### 1.1. Receptor

The EGF receptor of *Drosophila* shows a high degree of sequence [homology](#) to its vertebrate counterparts, in both its intracellular **tyrosine-kinase** domain and its extracellular ligand-binding domain ([1](#), [2](#)). The receptor is expressed ubiquitously throughout most stages of development. Thus, local activation of the pathway must be regulated by other means.

#### 1.2. Ligands

The singularity of the receptor in *Drosophila* is compensated by several ligands that confer different aspects of regulation. The primary activating ligand is called Spitz. It is a protein containing a **signal sequence**, a single [EGF motif](#), a putative dibasic cleavage signal, and a **transmembrane**, and a cytoplasmic **domain** ([3](#)). This structure is reminiscent of [transforming growth factor](#) a (TGFa), an activating ligand of the vertebrate EGF receptor. Like the receptor, the precursor form of Spitz is ubiquitously expressed. However, because the precursor form is inactive, this does not lead to receptor activation. The key to receptor activation is the **proteolytic** processing of Spitz, to release the secreted form containing the EGF motif ([4](#)). Although the precise regulation of this cleavage is still unknown, components regulating it will be described below. Secreted Spitz is a very active ligand, and its misexpression leads to dramatic changes in cell fates through high levels of activated [MAP kinase](#).

A second activating ligand is Vein, which contains a signal peptide, a single EGF motif, and an [immunoglobulin](#) domain ([5](#)). Vein is secreted constitutively and is a weaker ligand; its absence or misexpression do not induce dramatic changes in cell fates. While in some contexts Vein functions

as a single ligand (6), in other cases it appears to cooperate and complement Spitz in the same tissues.

A third activating ligand is [Gurken](#), encoding a protein with an overall structure similar to Spitz (7). There is no evidence that Gurken must be cleaved in order to be active. Gurken expression and activity are restricted to the female ovary, where it is responsible for inducing initially the posterior cell fates of the follicle cells covering the oocyte, and subsequently the dorsal cell fates. The expression of Gurken, and especially the localization of its [messenger RNA](#) in the proximity of the oocyte nucleus, is highly regulated.

Finally, a fourth ligand is Argos, encoding a secreted protein with a signal peptide and a single EGF motif (8). In contrast to the other ligands, Argos is an inhibitory ligand. Its expression abolishes the capacity of the receptor to respond to activation by Spitz (9). Transcription of Argos is induced by the activated EGF receptor pathway, thus constituting an inhibitory feedback loop (10).

### 1.3. Accessory Signaling Molecules

Among the mutants that displayed a phenotype similar to *spitz* were the genes *rhomboid* and *Star*, encoding novel transmembrane proteins with multiple and single transmembrane domains, respectively (11-13). A variety of observations suggest that these two proteins regulate the processing of Spitz. Ectopic expression of secreted Spitz can overcome the *rhomboid* and *Star* mutant phenotypes (4). In addition, Rhomboid and Star can function nonautonomously; namely, their expression in one set of cells will give rise to activation of the EGF receptor pathway in neighboring cells, consistent with their capacity to induce a diffusible signal (14). How Rhomboid and Star contribute to Spitz processing is an open question, as they have no **proteinase**-homology domains. Rhomboid homologues have been identified in [Caenorhabditis elegans](#) and humans, raising the possibility that the same molecules may be regulating ligand processing in other species.

The clue to the dynamic activation pattern of the EGF receptor pathway seems to lie in the expression pattern of Rhomboid. It is possible to follow the activation of the EGF receptor pathway, as well as that of other [tyrosine kinase receptors](#), by an [antibody](#) directed against the activated (double **phosphorylated**) form of MAP kinase. The pattern induced by EGF receptor activation can be identified, and it disappears specifically in mutants for the EGF receptor, *spitz*, *rhomboid*, or *Star*. This pattern follows precisely the expression of *rhomboid* (15).

## 2. Biological Example of Single Signaling Burst: The Embryonic Ventral Ectoderm

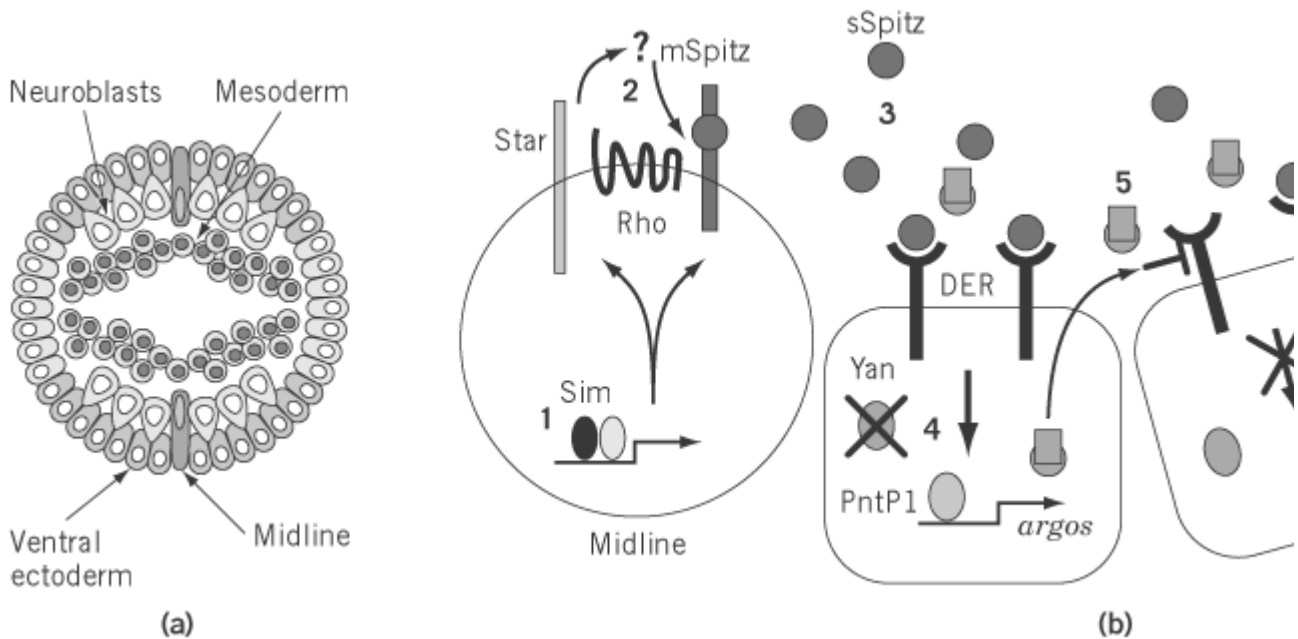
The earliest function of the EGF receptor pathway in the embryo entails the induction of ventral ectodermal cell fates. Following gastrulation, a group of ~ 8 cells on each side of the ventral midline are designated as neuroectodermal cells. Within this region, different rows of cells assume different fates, through activation of the EGF receptor pathway. The capacity of the pathway to induce graded fates stems from the fact that the activating signal, secreted Spitz, emanates from a single row of cells, the midline glial cells, positioned at the ventral midline (14). Diffusion of the ligand to neighboring cells generates graded activation of the pathway, which can be visualized by the distribution of activated MAP kinase (15).

The ectodermal cells respond differently to the varying levels of EGF receptor activation. The cells closest to the source of the ligand, induce a dual response. On the one hand, MAP kinase phosphorylates and inactivates the **transcriptional** repressor Yan, an ETS domain protein. On the other hand, MAP kinase induces the transcription of Pointed P1, an ETS domain transcriptional activator. This leads to the transcription of Pointed P1-target genes in the cell row closest to the midline (16). The transcriptional responses of the cells lying further from the midline are not known.

A central target gene for Pointed P1 is *argos*. Thus, activation of the EGF receptor leads to the production of Argos protein, which diffuses to neighboring cells and reduces the level of signaling in these cells, maintaining the EGF receptor activation gradient. This system is very robust.

Reproducible patterning of the ventral ectoderm is obtained, regardless of the absolute levels of secreted Spitz. As long as Spitz is emanating from the midline and induces the inhibitor Argos in the cells receiving maximal signaling, the pattern will be generated correctly. EGF receptor signaling in the ventral ectoderm is illustrated schematically in Figure 1.

**Figure 1.** Signalling by the EGF receptor pathway in the embryonic ventral ectoderm. **(a)** Scheme of embryonic cross section midline cells, positioned in the center of the ventral ectoderm. **(b)** (1) The fate of the midline cells is determined by the *Star* gene. The midline cells express the components necessary for Spitz processing, namely, the Spitz precursor, Rhomboid, and *Sim* surrounding ventral ectoderm in a graded manner. (4) In the cells closest to the source, maximal EGF receptor activation occurs, leading to the transcription of *PntP1*. (5) A central target gene is *argos*. Its transcription leads to the diffusion of Argos protein, which causes reduction or termination of receptor activation in the lateral cells.



### 3. Multiple Bursts of Signaling: The Ovary and Eye

The above example represents a simple case, in which the pathway is only activated once. More complicated scenarios involve multiple activation cycles of the pathway in the same tissue, each time leading to different consequences. In the ovary, signaling from the oocyte to the surrounding follicle cells determines the polarity of these cells. Initially, when the egg is small and the oocyte nucleus is positioned at the posterior end, localized Gurken, a ligand for the EGF receptor, will trigger posterior cell fates (17, 18). Subsequently, the oocyte nucleus migrates to the future dorsal anterior corner of the oocyte. Activation of the receptor by Gurken at this stage leads to the induction of dorsal follicle cell fates (19). Finally, induction of Rhomboid expression by the EGF receptor pathway leads to a third cycle of activation by allowing processing of a ligand expressed in the follicle cells (possibly Spitz), this time patterning specific follicular structures called dorsal appendages (20). It should be noted that the final cycle of activation of the receptor was obtained by inducing Rhomboid.

This aspect becomes even more dramatic in the induction of photoreceptor cell fates in the developing eye. The first photoreceptor in each ommatidium R8 is induced independent of the EGF receptor pathway. This cell serves as a source for secreted Spitz, which will activate the receptor in neighboring cells and induce the first cohort of photoreceptors (21, 22). These newly induced cells will now also produce Argos, which prevents the signal from progressing further. It seems that the differentiating cells also express the processing machinery for Spitz, and thus induce a new round of



activation. This continues for several cycles, inducing first distinct sets of photoreceptor cells and subsequently cone or pigment cells.

In conclusion, an amazingly diverse set of roles is fulfilled by the EGF receptor pathway, only a sample of which were mentioned above. They include the induction of cell fates in different contexts, as well as cell proliferation (eg, in the wing [imaginal disc](#)). In most cases, a common signaling pathway is used, both upstream and downstream of the receptor. This raises the question of specificity. It is clear that the induction of MAP kinase signaling cannot determine the final outcome. Rather, in different cell types, the same pathway induces different responses, depending on the repertoire of [transcription factors](#) that are responsive in that tissue. Differential regulation of the receptor itself appears to be directed by the multiplicity of ligands, being regulated in different ways in terms of their expression or localization, or conferring varying levels of activation.

## Bibliography

1. E. Livneh, L. Glazer, D. Segal, J. Schlessinger, and B.-Z. Shilo (1985) *Cell* **40**, 599–607.
2. E. D. Schejter, D. Segal, L. Glazer, and B.-Z. Shilo (1986) *Cell* **46**, 1091–1101.
3. B. J. Rutledge, K. Zhang, E. Bier, Y. N. Jan, and N. Perrimon (1992) *Genes Dev* **6**, 1503–1517.
4. R. Schweitzer, M. Shaharabany, R. Seger, and B.-Z. Shilo (1995) *Genes Dev* **9**, 1518–1529.
5. B. Schnepf, G. Grumblin, T. Donaldson, and A. Simcox (1996) *Genes Dev* **10**, 2302–2313.
6. T. Yarnitzky, L. Min, and T. Volk (1997) *Genes Dev* **11**, 2691–2700.
7. F. S. Neuman-Silberberg and T. Schüpbach (1993) *Cell* **75**, 165–174.
8. M. Freeman, C. Klämbt, C. S. Goodman, and G. M. Rubin (1992) *Cell* **69**, 963–975.
9. R. Schweitzer, R. Howes, R. Smith, B.-Z. Shilo, and M. Freeman (1995) *Nature* **376**, 699–702.
10. M. Golembo, R. Schweitzer, M. Freeman, and B.-Z. Shilo (1996) *Development* **122**, 223–230.
11. U. Mayer and C. Nüsslein-Volhard (1988) *Genes Dev* **2**, 1496–1511.
12. E. Bier, L. Y. Jan, and Y. N. Jan (1990) *Genes Dev* **4**, 190–203.
13. A. L. Kolodkin, A. T. Pickup, D. M. Lin, C. S. Goodman, and U. Banerjee (1994) *Development* **120**, 1731–1745.
14. M. Golembo, E. Raz, and B.-Z. Shilo (1996) *Development* **122**, 3363–3370.
15. L. Gabay, R. Seger, and B.-Z. Shilo (1997) *Science* **277**, 1103–1106.
16. L. Gabay, H. Scholz, M. Golembo, A. Klaes, B.-Z. Shilo, and C. Klämbt (1996) *Development* **122**, 3355–3362.
17. A. González-Reyes, H. Elliot, and D. St. Johnston (1995) *Nature* **375**, 654–658.
18. S. Roth, F. S. Neuman-Silberberg, G. Barcelo, and T. Schüpbach (1995) *Cell* **81**, 967–978.
19. F. S. Neuman-Silberberg and T. Schüpbach (1994) *Development* **120**, 2457–2463.
20. A. Sapir, R. Schweitzer, and B.-Z. Shilo (1998) *Development* **125**, 191–200.
21. M. Freeman (1996) *Cell* **87**, 651–660.
22. M. Tio and K. Moses (1997) *Development* **124**, 343–351.

## Suggestions for Further Reading

23. R. Schweitzer and B.-Z. Shilo (1997) A thousand and one roles for the *Drosophila* EGF receptor. *Trends Genet.* **13**, 191–196.
24. N. Perrimon and L. Perkins (1997) There must be 50 ways to rule the signal: the case of the *Drosophila* EGF receptor. *Cell* **89**, 13–16.
25. J. D. Wasserman and M. Freeman (1997) Control of EGF receptor activation in *Drosophila*. *Trends Cell Biol.* **7**, 431–436.
26. L. A. Nilson and T. Schüpbach (1998) "EGF receptor signaling in" *Drosophila* oogenesis. In *Current Topics in Developmental Biology*, Vol. **44** (R. Pederson and G. P. Schatten eds.),

## Drug Resistance

The resistance of cancer cells to cytostatic agents has been a significant impediment to the effective chemotherapy of cancer. Studies with model systems, where cancer cells grown *in vitro* were selected for resistance to specific anticancer drugs, identified two major classes of drug-resistant cells: (i) cells resistant to a single class of drugs with the same mechanism of action, and (ii) cells resistant to chemically diverse drugs with multiple mechanisms of action. The latter phenomenon was called multiple drug resistance (MDR). Some metastatic cancers are intrinsically resistant to chemotherapy, whereas others respond to treatment initially but subsequently acquire resistance, not only to the chemotherapeutic agents used against them, but to other chemically unrelated drugs as well (1). Numerous factors have been associated with MDR. These include elevated levels of protein kinase C, enhanced sodium pump activity, mutations of topoisomerase II and [tubulin](#), and altered levels of calcium and [calmodulin](#). It has not been possible, however, to correlate any of these clearly with the MDR **phenotype**. A large body of evidence, on the other hand, derived from studies of cells in culture, transfection experiments, histochemical studies of a wide variety of tumors graded for resistance to chemotherapy, and studies in reconstituted systems, strongly implicates energy-dependent pump systems that either exclude or extrude chemotherapeutic agents from MDR cells.

Besides human cancer cells and their rodent models, MDR pumps have been found in a wide array of organisms: bacteria, protozoa, and fungi, and mammalian cells; even plant cells have MDR pumps. These share the unique feature of excluding a fairly broad range of chemically unrelated compounds from the cell. Only some of these belong to the ATP Binding Cassette (ABC) superfamily of transporters, and they do not ubiquitously use ATP as their energy source. The majority of the bacterial pumps, for example, use the energy of the [proton motive force](#). The ATP-dependent system that has been most extensively characterized is the *MDR1* gene product, the multidrug transporter or [P-glycoprotein](#) (P-gp).

Many cells selected for drug resistance do not show increased levels of P-gp but nonetheless are resistant to a broad range of natural product drugs. Another member of the ABC superfamily, the MDR-associated protein (MRP1), is expressed in some of these cell lines at elevated levels. MRP1 is similar to P-gp in that it is capable of decreasing intracellular levels of drugs and is ATP-dependent. While P-gp has two membrane-spanning domains, MRP has three. The presence of a third such domain is observed in several other ABC proteins, all of which are more closely related to MRP than to any other member of this superfamily of transporters. There are now at least a dozen ABC proteins that together constitute the MRP branch of the superfamily. The most important members are (i) the human MRP, which is an anionic conjugate transporter, (ii) the multispecific organic anion transporter, MOAT or MRP-2, and (iii) SUR1, the sulfonylurea receptor. There are other significant differences between MRP and P-gp: (i) MRP1 catalyzes the efflux of anionic or neutral drugs conjugated to [glutathione](#), glucuronate, or sulfate; (ii) some substrates of P-gp, eg, taxol, are not substrates for MRP; (iii) MDR effected by P-gp is reversed by verapamil and [cyclosporin A](#), but not that caused by MRP. It is thus quite likely that both P-gp and MRP are drug transporters but that their detailed mechanisms of action may differ (2).

Additionally, a 110-kDa protein, the lung resistance protein, has been identified in non P-gp-resistant lung carcinoma cells. This protein is a major constituent of an intracellular [ribonucleoprotein](#) (vault) complex that consists of three other proteins of 210 kDa, 190 kDa, and 54 kDa, plus an RNA

molecule. There is increasing evidence that the vault complex is involved in MDR (3).

### Bibliography

1. M. M. Gottesman and I. Pastan (1993) *Ann. Rev. Biochem.* **62**, 385–427.
2. D. Lautier, Y. Cantrot, R. G. Deeley, and S. P. C. Cole (1996) *Biochem. Pharmacol.* **52**, 967–977.
3. G. L. Scheffer, P. J. Wijngaard, M. J. Flens, M. A. Izquierdo, M. L. Slovak, H. M. Pinedo, C. J. L. M. Meijer, H. C. Clevers, and R. J. Scheper (1995) *Nature Medicine* **1**, 578–582.

### Suggestions for Further Reading

4. I. B. Roninson, ed. (1991) *Molecular and Cellular Biology of Multidrug Resistance in Tumor Cells*, Plenum Press, New York. A comprehensive reference in which numerous researchers cover almost all areas of multidrug resistance in cancer cells.
5. S. V. Ambudkar and M. M. Gottesman, eds. (1998) *Methods in Enzymology*, Vol. **292**, Academic Press, San Diego. This volume of *Methods in Enzymology* is devoted to the ABC transporters and is a very good source for methodological details covering biochemical, molecular biological, and cell biological aspects.

## Duplication, Chromosomal

There are numerous types of chromosomal duplications. Each [chromosome](#) itself is duplicated during the **S-phase** of the normal [cell cycle](#). There are also various duplications that arise from chromosome breakage and reunion cycles (see [Dicentric Chromosome](#)), which might be an important contributory factor in the evolution of genomes. There are also duplication events that encompass whole [genomes](#). In the frog *Xenopus laevis*, different species have evolved that contain multiple sets of related chromosomes. This is called **polyploidy**. *Xenopus tropicalis* is the progenitor **diploid** species, whereas *Xenopus laevis* is pseudotetraploid (1). Although many duplicated genes are silenced, genes that remain active often show tissue or developmental differences in activity. Thus specialized gene family members may arise by both genomic duplication and by genetic duplication (see [Contiguous Genes](#)). Other examples of chromosomal duplication result from gene **amplification** (see [Double Minute Chromosome](#)).

Chromosomal duplication occurs through a variety of mechanisms, some that operate in evolutionary time, and others that occur in specialized cells during [development](#) or in response to environmental stress. The most common form of chromosomal duplication results from unequal sister [chromatid](#) exchange. This is caused by misalignment of two sister chromatids during **prophase** followed by **reciprocal recombination**, which duplicates one chromosome, whereas the other has a deficiency. Amplification of the large 18 S and 28 S ribosomal RNA genes in *Drosophila melanogaster* might occur through many reiterations of events of this type (2). It is also possible that the presence of thousands of short interspersed nuclear elements, such as **Alu** repeats or the long interspersed nuclear element (**LINE**)-1 in mammalian chromosomes, might lead to many misalignments and chromosomal duplications.

[Gene amplification](#) is another form of local chromosomal duplication. In cultured cells, treatment with methotrexate and hydroxyurea results in multiple replicative cycles of the **dihydrofolate reductase** (DHFR) gene. The molecular mechanisms controlling this overreplication of the DHFR gene are not understood, but it is assumed that the replicative process leads to the generation of free

DNA ends that then undergo homologous recombination (3). If the over-replicated DNA circularizes, [double minute chromosomes](#) are formed.

Chromosome duplication also occurs following modification of chromosomeal behavior during the [cell cycle](#). Normally the DNA in chromosomes replicates once during **S phase**. If multiple replicative cycles occur, however, then duplicated chromosomes can coexist in the same somatic cell nucleus. If the chromosomes separate, the cell is called endopolyploid. If the chromosomes do not separate, [polytene chromosomes](#) are formed. This repeated replicative process is called endoreduplication (4). Endopolyploid cells and those containing polytene chromosomes never divide, so they represent, a somatic dead end or a state of terminal differentiation, respectively.

#### Bibliography

1. H. R. Kobel and L. Du Pasquier (1986) *Trends Genet.* **2**, 310–313.
2. S. A. Endow and K. C. Atwood (1988) *Trends Genet.* **4**, 348–351.
3. R. T. Schimke (1988) *J. Biol. Chem.* **263**, 5989–5992.
4. W. Nagl (1978) *Endopolyploidy and Polyteny in Differentiation and Evolution*, Elsevier/North Holland, Amsterdam.

## Dynactin

Dynactin is a large [protein](#) complex that is important for the function of cytoplasmic [dynein](#) in intracellular motility (see [Motor Proteins](#)). Dynactin was first identified as a protein complex that stimulates dynein movement of membranous organelles along [microtubules](#) (1). It has been proposed that dynactin functions to bind cytoplasmic dynein to its membranous organelle cargo. Work in diverse species from fungi to mammalian neurons suggests that an alternative function of dynactin is to link cytoplasmic dynein to the [actin](#) cortex for the movement of microtubules (2).

The backbone of the dynactin molecule is a short (~40nm) filament composed of 8 to 12 molecules of actin-related protein 1 (Arp1) (3). Arp1 shares many of the binding domains of actin, including the [spectrin](#), myosin, and villin binding sites. Actin-capping protein is bound to one end of the dynactin filament, and a unique capping protein is located at the other end. The p150<sup>Glued</sup> subunit projects from the filament. This protein is the product of the *Glued* gene of *Drosophila*, and it binds both microtubules and the intermediate-chain subunit of cytoplasmic dynein. Overexpression of another major subunit of dynactin, p50, obstructs assembly of the dynactin complex, and thus it is named *dynamitin* (4). Overexpression of p50 disrupts localization of both dynein and dynactin to membranous organelles or [kinetochores](#). Genetic studies indicate that dynein and dynactin act in the same pathways (5).

#### Bibliography

1. T. A. Schroer and M. P. Sheetz (1991) *J. Cell Biol.* **115**, 1309–1318.
2. J. T. Carminati and T. Stearns (1997) *J. Cell Biol.* **138**, 629–641.
3. D. A. Schafer et al. (1994) *J. Cell Biol.* **126**, 403–412.
4. C. J. Echeverri et al. (1996) *J. Cell Biol.* **132**, 617–633.
5. T. A. Schroer (1994) *J. Cell Biol.* **127**, 1–4.

## Suggestions for Further Reading

6. E. A. Holleran, S. Karki, and E. L. Holzbaur (1998) The role of the dynactin complex in intracellular motility, *Int. Rev. Cytol.* **182**, 69–109.
7. V. Allan (1994) Organelle movement. Dynactin: protrait of a dynein regulator, *Curr. Biol.* **4**, 1000–1002.

## Dynamic Light Scattering

Dynamic light scattering (also known as [photon correlation spectroscopy](#)) from particles in solution can be used to deduce information about the size, shape, and dynamics of biological [macromolecules](#) and supramolecular assemblies. Structural analysis of individual biomolecules, conformational transitions, and intermolecular interactions can also be probed. The major use of dynamic light scattering in biology is the rapid determination of the translational [diffusion](#) constant for macromolecules in solution. The interpretation of **diffusion** constant data is frequently strengthened by combination with other information from methods, such as [electron microscopy](#), static [light scattering](#), and neutron or X-ray [small-angle scattering](#). Biological systems are well-suited to study by dynamic light scattering, because they are generally large enough to be strong scatterers at low concentrations and their diffusion constants are such that they give rise to autocorrelation functions that can be readily measured. Dynamic light scattering measurements have been used for studying objects ranging in size from proteins as small as 6000 Da to intact cells as large as micrometers.

### 1. Theory of Dynamic Light Scattering

Dynamic light scattering measures [light scattering](#) intensity fluctuations in a small volume of a sample. For small proteins (~25kDa) these fluctuations are on the microsecond timescale, while for cells they are in milliseconds. The fluctuations are related to the Brownian motion of the particles giving rise to density fluctuations caused by variations in the number of molecules in the scattering volume and random agglomerations. This measurement in the time domain is related to the spectral density of the fluctuations in the frequency domain by a Fourier transformation. In practice, diffusion coefficients are determined using an autocorrelation function that is measured by accumulating the product of the number of photons arriving at the detector from successive time intervals. This operation is repeated thousands of times and averaged. The intensity autocorrelation function,  $G_2(t)$ , is obtained by storing the average products  $I_t I_{t+t'}$ , where  $t$  is an incremented time delay, in successive channels to yield

$$G_2(t) = \langle I_t \cdot I_{t+t'} \rangle \quad (1)$$

For a solution of macromolecules, assuming a Gaussian distribution of fluctuations,  $G_2(t)$  is related to the scattered electric field autocorrelation function  $G_1(t)$  by

$$G_2(t) = \langle I^2 \rangle + G_1(t)^2 \quad (2)$$

For translation diffusion of monodisperse particles that are small with respect to the incident wavelength,  $l$ :

$$G_2(t) = A + B e^{-2\Gamma t} \quad (3)$$

where  $A$  and  $B$  are constants that depend upon experiment geometry, and

$$D = \frac{\Gamma}{Q^2} \quad (4)$$

$D$  is the translational diffusion coefficient, and  $G$  is the reciprocal of the characteristic decay time,

$$Q = \frac{4\pi \sin(\frac{Q}{2})}{\lambda} \quad (5)$$

where  $O$  is the scattering angle. For a continuous polydisperse system, equation (4) is integrated over all sizes, and hence  $G$  values, to give

$$G_1(t) = \int G(\Gamma)e^{-\Gamma t} d\Gamma \quad (6)$$

where  $G(G)$  is a distribution function that can be evaluated.

Rotational diffusion and internal dynamics can also influence the autocorrelation functions measured in a dynamic light scattering experiments. The internal dynamics of a macromolecule in solution become important only if the amplitude of the motions is not very much smaller than the wavelength of the light. For large motions, different parts of the sample molecule will scatter out of phase, and the internal dynamics will contribute one or more relaxation times that will contribute to the decay of the autocorrelation function in a single or multiple-exponential fashion.

### 1.1. Examples of Dynamic Light Scattering Applications

The [transcription](#), [translation](#), and replication of genetic information contained in the sequences of DNA is controlled by interactions of proteins with DNA and with RNA. These interactions can be probed using dynamic light scattering. For example, the enzyme DNA gyrase (see **Topoisomerases**) catalyzes the ATP-dependent **supercoiling** of DNA. The diffusion constant for DNA gyrase measured by dynamic light scattering (1) is significantly smaller than expected based on its molecular weight, and binding of DNA to the gyrase produces no change in the diffusion constant, indicating that there is no overall conformational change. These observations, combined with parallel small-angle neutron scattering experiments to determine [radius of gyration](#) values for the gyrase and its complex with DNA, have led to the conclusion that the gyrase has grooves on its surface to accommodate the DNA in a very compact complex.

Dynamic light scattering has played a role in drug development. The pharmacological effects of drugs can depend upon their behavior in solution and the characteristics of aggregates they may form (2). Lipids can be used to “solubilize” drugs that may be insoluble in aqueous media, and hence enhance their transport and efficacy. Dynamic light scattering has been used to evaluate aggregation behavior of antidepressants (3), critical micelle concentrations, and sizes of micelles designed to encapsulate the drug Indomethacin (4), as well as to evaluate how the opiate drug loperamide alters the temperature-induced phase transition of phosphatidylcholine vesicles (5).

Dynamic light scattering has also been used to monitor the assembly of enveloped **viruses**. Lyles et al. (6) used a combination of dynamic and static light scattering with stopped flow to study matrix protein binding to nucleocapsids of [vesicular stomatitis virus](#). Dynamic light scattering has recently become a common tool for assessing the crystallizability of macromolecules and macromolecular assemblies (7, 8). Because modern molecular biology techniques have provided sufficient amounts of pure materials, dynamic light scattering can be used automatically to “screen” large numbers of crystallization conditions to determine diffusion constants and hence evaluate possible aggregation states that inhibit crystallization.

## Bibliography

1. S. Krueger et al. (1990) *J. Mol Biol.* **211**, 211–220.
2. D. Attwood and P. Fletcher (1986) *J. Pharm. Pharmacol.* **38**, 494–498.
3. A. D. Atherton and B. W. Barry (1985) *J. Pharm. Pharmacol.* **37**, 854–862.
4. D. Attwood, G. Ktistis, Y. McCormick, and M. J. Story (1989) *J. Pharm. Pharmacol.* **41**, 83–86.
5. E. Hantz, A. Cao, R. S. Phadke, and E. Taillandier (1989) *Chem. Phys. Lipids* **51**, 75–82.
6. D. S. Lyles, M. O. McKenzie, and R. R. Hangton (1996) *Biochemistry* **35**, 6508–6518.
7. A. D'Arcy (1994) *Acta Crystallogr.* **D50**, 469–471.
8. A. Ferre-D'Amare and S. K. Burley (1994) *Structure (London)* **2**, 357–359.

## Suggestions for Further Reading

9. B. J. Berne and R. Pecora (1976) *Dynamic Light Scattering with Applications to Chemistry, Biology and Physics*, Wiley, New York.
10. P. A. Janmey (1993) "Application of Dynamic Light Scattering to Biological Systems. In" *Dynamic Light Scattering: The Method and Some Applications* (W. Brown ed.), Oxford University Press, Oxford, U.K.
11. R. Pecora (1985) *Dynamic Light Scattering: Applications of Photo Correlation Spectroscopy*, Plenum Press, New York.

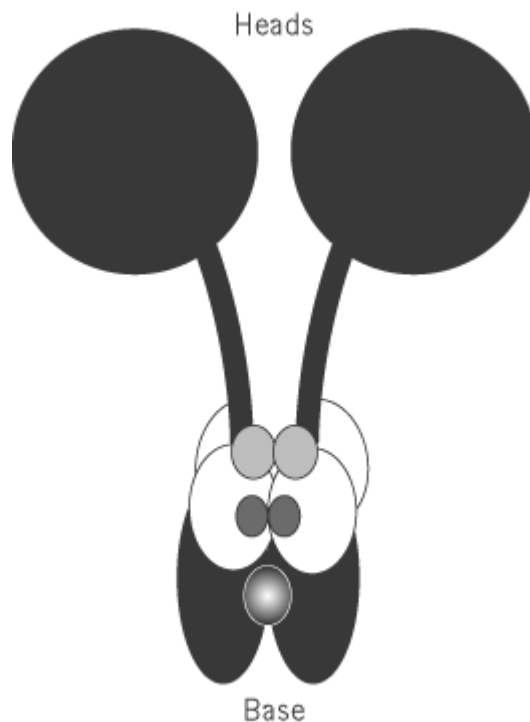
## Dynein

Dyneins are **microtubule**-based [motor proteins](#) that use the energy released by the hydrolysis of ATP to move toward the minus end of a microtubule. The dyneins are classified as either cytoplasmic or flagellar (also known as axonemal), and the flagellar dyneins are further divided into outer and inner arm dyneins. The first dyneins to be identified were the flagellar dyneins, which are the arms that project from the A microtubules of **ciliary** and flagellar axonemes. Flagellar dyneins generate the sliding force between microtubules that is the basis for the propulsive bending of flagella (and cilia) that moves cells through fluids. Cytoplasmic dyneins move membranous organelles and [kinetochores](#) along microtubules and assist in assembling the mitotic spindle. Each dynein is a large molecular complex, composed of one or more large polypeptide chains of ~500kDa, known as heavy chains, and various accessory intermediate chains (50 to 140 kDa) and light chains (8 to 45 kDa). Often the heavy chain alone is referred to as a dynein. An ~350kDa portion of each heavy chain forms a large globular head **domain**, and the rest of the heavy chain extends from the head as a thin flexible stalk. The globular head contains the motor part of the molecule, the microtubule-binding site, and the ATP hydrolysis region. Two or three dynein heavy chains are often connected by their thin stalks to a common base. The base is made up of the intermediate and light chains, which are important for binding dynein to the organelles that it transports (its cargo).

There are at least 12 different flagellar dynein heavy-chain **genes**. To date, only two uniquely cytoplasmic dynein heavy-chain genes have been identified. The major cytoplasmic dynein (Fig. [1](#)) has a native molecular weight of  $1.5 \times 10^6$ Da and contains two identical heavy chains associated with two 74-kDa intermediate chains, four 55-kDa light intermediate chains, and light chains of 22, 14, and 8 kDa. The heavy chain of this dynein is encoded by one gene, but multiple genes have been identified for the intermediate and light intermediate chains and for the 14-kDa light chain. The

heavy chain of the minor cytoplasmic dynein is encoded by a unique gene, but the intermediate or light chains, with which it is presumed to be associated, have not yet been identified. Unlike cytoplasmic dynein, flagellar outer arm dyneins have two or three different heavy chains. Inner arm dyneins are either heterodimers or monomers of heavy chains and have various accessory polypeptides.

**Figure 1.** Schematic representation of the overall structure of a cytoplasmic dynein.



Structural analysis of the heavy-chain motor domain indicates that the head is made up of seven lobes arranged in a ring surrounding a central cavity (1). Each heavy chain has four consensus **nucleotide-binding sites** known as **P loops**. The first P loop, P1, is believed to be the site of ATP binding and hydrolysis for force generation (2). The role of the other P loops is unknown. In the presence of vanadate, the heavy chains are cleaved by UV photolysis into two fragments at the first P loop. The different types of dynein hydrolyze MgATP at different rates. Dyneins can be attached to glass coverslips via their bases and their heads will bind microtubules. Using computer-enhanced video microscopy, it can be observed that this bound dynein translocates microtubules across the coverslip in the presence of ATP. The velocity of cytoplasmic dynein microtubule movement *in vitro* is ~1micron/second.

Flagellar dyneins attach via their base to one of the outer doublet microtubules (the A microtubule) of the axoneme and bind directly to **tubulin** or to other axonemal microtubule-binding proteins. Their head domains bind to the B microtubule of the neighboring outer doublet in an ATP-sensitive manner and use the hydrolysis of ATP to push them toward the tip of the axoneme. Cytoplasmic dyneins bind to other cellular organelles. The mechanism by which dynein binds to these organelles is unknown. However, cytoplasmic dynein binds directly to artificial lipid membranes, **liposomes**, so dynein may bind directly to the organelle (3). There is also evidence, however, that an accessory protein complex, **dynactin**, is also important for cytoplasmic dynein binding to cargo (4). One dynactin subunit binds the dynein intermediate chain *in vitro*. Although they bind to different cargoes, cytoplasmic and flagellar outer arm dyneins have closely related cargo-binding subunits, the intermediate chains. The C-termini of the intermediate chains are conserved, and it is believed that



they interact with the heavy chains. It is thought that the distinct intermediate-chain N-termini allow binding to different cargoes (5). Deletion of the cytoplasmic dynein heavy chain is lethal in multicellular organisms. Nonlethal mutations of cytoplasmic dynein heavy and light chains lead to axonal pathfinding defects and other pleiotropic **developmental** defects in *Drosophila*. Dyneins are regulated *in vivo*. During flagellar bending, only a subset of the flagellar dyneins function at a time. The **phosphorylation** state of an inner dynein-arm intermediate chain modulates this regulation (6). Genetic methods have also identified several proteins that make up a dynein regulatory complex, although the mechanism of the complex remains unknown. Cytoplasmic dynein heavy-chain, intermediate-chain, and light-intermediate chains are phosphorylated *in vivo*, and phosphorylation of each has been implicated in regulating dynein.

### Bibliography

1. M. Samso et al. (1998) *J. Mol. Biol.* **276**, 927–937.
2. I. R. Gibbons (1995) *Cell Motility and Cytoskeleton* **32**, 136–144.
3. K. L. Ferro and C. A. Collins (1995) *J. Biol. Chem.* **270**, 4492–4496.
4. S. R. Gill et al. (1991) *J. Cell Biol.* **115**, 1639–1650.
5. C. G. Wilkerson et al. (1995) *J. Cell Biol.* **129**, 169–178.
6. G. Habermacher and W. S. Sale (1997) *J. Cell Biol.* **136**, 167–176.

### Suggestions for Further Reading

7. I. R. Gibbons (1988) Dynein ATPases as microtubule motors, *J. Biol. Chem.* **263**, 15837–15840.
8. N. Hirokawa (1998) Kinesin and dynein superfamily proteins and the mechanism of organelle transport, *Science* **279**, 519–526.
9. R. B. Vallee, H. S. Shpetner, and B. M. Paschal (1989) The role of dynein in retrograde axonal transport, *Trends Neurosci.* **12**, 66–70.
10. R. B. Vallee and M. P. Sheetz. (1996) Targeting of motor proteins, *Science* **271**, 1539–1544.
11. G. B. Witman (1992) Axonemal dyneins, *Curr. Opinion Cell Biol.* **4**, 74–79.

### Dystrophin

Dystrophin is a very large protein (molecular weight 427 kDa) located at the cytoplasmic face of the sarcolemma membrane of striated, smooth, and cardiac muscle cells and, most specifically, in the transverse tubules of the triadic structure. It represents only about 0.01% of the total skeletal muscle protein present in the tissue (1). Its absence, however, appears to be correlated with the fatal muscle-wasting disease known as *Duchenne muscular dystrophy* (incidence about 1 in 3500 of live male births). Becker muscular dystrophy, a milder form of muscular dystrophy (incidence about 1 in 35,000 of live male births), is related to defects/deletions in the dystrophin sequence. Muscular dystrophy is not related solely to the absence or mutation of dystrophin, however, because evidence has been presented to show that some membrane **glycoproteins** that normally interact with native dystrophin are **degraded** at abnormally high rates in either its absence or when it is mutated. This indicates that dystrophin may bind to and hence stabilize the structures of these dystrophin-associated proteins *in vivo*. It seems likely that dystrophin's structural role in muscle and nerve is comparable to that played by [ankyrin](#) and [spectrin](#) with membrane proteins in kidney, brain, and erythrocytes (1).

Dystrophin is a member of the spectrin [superfamily](#) of proteins and has many structural similarities with both spectrin and  $\alpha$ -actinin. Each has a central rod domain, as well as globular regions located at the *N*- and *C*-termini. Like  $\beta$ -spectrin and  $\alpha$ -actinin, the *N*-terminal domain of dystrophin contains a sequence that has high [homology](#) with the **actin**-binding site in  $\alpha$ -actinin (2). *C*-terminal to the rod domain of dystrophin there is a **cysteine**-rich region about 150 residues that is somewhat akin to that in  $\alpha$ -actinin, followed by a further 420 residues that is unique to dystrophin (1). Two putative **calcium-binding** sites occur within the *C*-terminal domain, although the sequences do differ from the consensus [EF-Hand Motif](#), possibly significantly. Most of the defects/deletions in dystrophin that cause clinical symptoms are located within the rod domain; and very few are found in the cysteine-rich or *C*-terminal domains. These portions of the molecule *in vivo* provide the site of association with a well-characterized transmembrane protein complex. The *C*-terminal domain has been modeled as a pair of **heptad repeat**-containing regions (five and six heptads in length) separated by a proline-rich region (2). The *d* position in the heptads is largely occupied by [leucine](#) residues and is thus reminiscent of the [leucine zipper](#) structure. These regions are compatible with the formation of a two-stranded [coiled-coil](#) structure composed of parallel chains. These data suggest that dystrophin molecules may form a homodimer, or that dystrophin may form a heterodimer with *utrophin* (a dystrophin-related molecule). It is also known from [electron microscopy](#) that dystrophin molecules can assemble in a staggered manner and that these in turn may assemble further to produce end-to-end aggregates.

The rod domain in dystrophin (length estimates derived from electron microscope studies vary between 100 and 175 nm) consists of about 25 repeats, each about 109 residues long. Significant variations do occur, nonetheless; and although the repeats are homologous to those in spectrin and  $\alpha$ -actinin, they are very much less regular, especially with regard to length. Ten of the 25 repeating motifs are largely complete (3), in contrast to the remaining 15 repeats, which lack some part(s) of the three constituent  $\alpha$ -helices that comprise the three- $\alpha$ -**helix** motif (see [Spectrin](#), where details of the conformation adopted are given). It has been proposed that dystrophin contains four **proline**-rich hinge regions within the rod domain that confer flexibility to the molecule and provide sites sensitive to **proteolysis**.

#### Bibliography

1. A. P. Monaco (1989) Dystrophin, the protein product of the Duchenne/Becker muscular dystrophy gene. *Trends Biochem. Sci.* **14**, 412–415.
2. D. J. Blake, J. M. Tinsley, K. E. Davies, A. E. Knight, S. J. Winder, and J. Kendrick-Jones (1995) Coiled-coil regions in the carboxy-terminal domains of dystrophin and related proteins: potentials for protein-protein interactions. *Trends Biochem. Sci.* **20**, 133–135.
3. D. A. D. Parry, T. W. Dixon, and C. Cohen (1992) Analysis of the three- $\alpha$ -helix motif in the spectrin superfamily of proteins. *Biophys. J.* **61**, 858–867.

#### Suggestion for Further Reading

4. J. M. Ervasti and K. P. Campbell (1993) Dystrophin and the Membrane Skeleton. *Curr. Opin. Cell Biol.* **5**, 82–87.

#### Eadie–Hofstee Plot

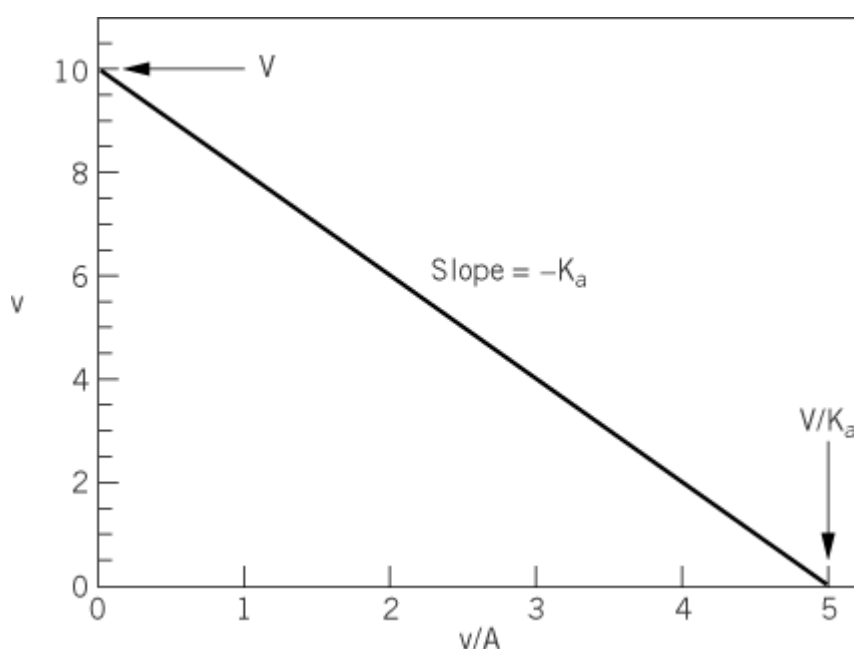
The Eadie–Hofstee plot was developed by Eadie (1) and Hofstee (2) and is one of three graphical procedures for determining values of the [kinetic](#) parameters for an [enzyme](#)-catalyzed reaction that

conforms to [Michaelis–Menten kinetics](#). It involves rearrangement of the initial velocity equation to a linear form that differs from that of the [Lineweaver–Burk plot](#). The plot is derived from the Michaelis–Menten equation describing the dependence of the enzyme-catalyzed velocity  $v$  on the concentration of the substrate  $A$  and its Michaelis constant  $K_a$ ; it is described by Equation (1):

$$v = -K_a \frac{v}{A} + V \quad (1)$$

where  $V$  is the maximum velocity of the reaction. Consequently, a plot of  $v$  against  $v/A$  yields a straight line with negative slope, equal to  $-K_a$ , a positive intercept on the abscissa that is equal to  $V/K_a$ , and a positive intercept on the vertical ordinate that is equal to  $V$  (Fig. 1).

**Figure 1.** Example of an Eadie–Hofstee plot for the determination of kinetic parameters for an enzyme-catalyzed reaction. The plot was produced by using the values of 10 and 2 (in arbitrary units) for  $V$  and  $K_a$ , respectively.



When values of  $V$  and  $K_a$  were being obtained by graphical methods, there were arguments for using this type of plot over that proposed by Lineweaver and Burk (3). But graphical procedures have now been supplanted by the use of regression analysis with computers (4).

A graphical procedure is required for the presentation of kinetic data. However, results that are presented as the variation of a dependent variable as the function of the ratio of a dependent and an independent variable, as in the Eadie–Hofstee plot, are not as readily interpreted as those presented as the variation of a dependent variable with the concentration of an independent variable. Further, the form of the kinetic equations for multisubstrate enzyme reactions becomes unnecessarily complex with the Eadie–Hofstee plot.

#### Bibliography

1. G. S. Eadie (1942) *J. Biol. Chem.* **146**, 85–93.
2. B. H. J. Hofstee (1959) *Nature* **184**, 1296–1298.
3. J. E. Dowd and D. S. Riggs (1965) *J. Biol. Chem.* **240**, 863–869.

4. W. W. Cleland (1979) Meth. Enzymol. **63**, 103–138.

## Ecdysone

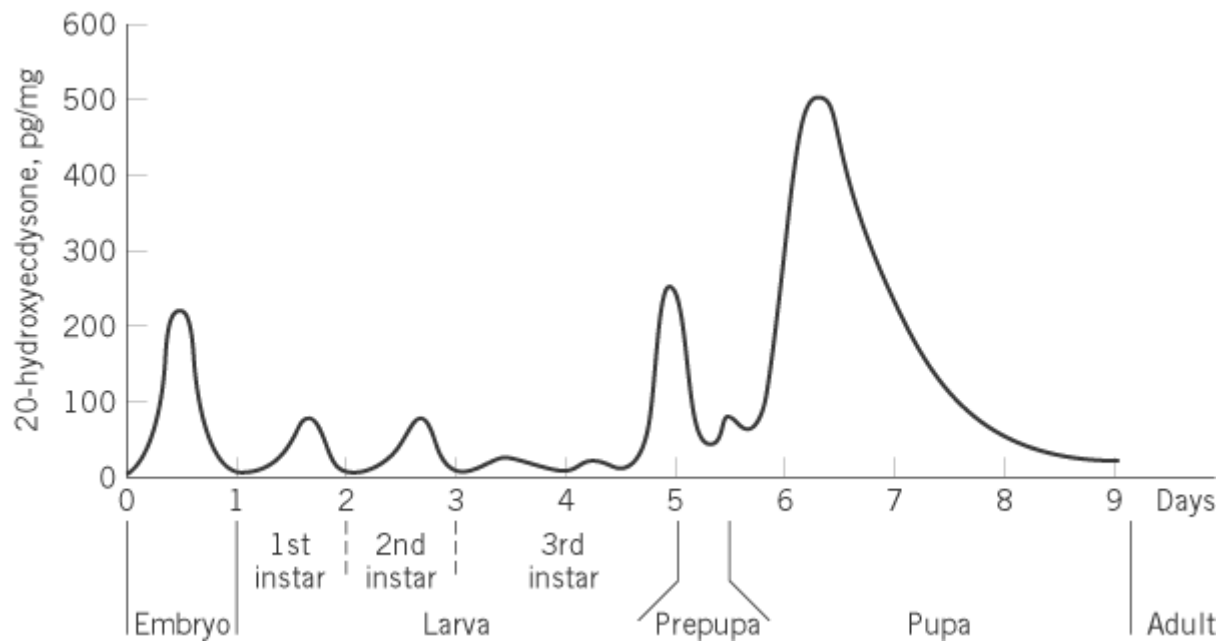
Pulses of the [steroid hormone](#) ecdysone coordinate postembryonic development in insects. In recent years, significant progress has been made in understanding the molecular mechanisms of ecdysone action in the fruit fly, *Drosophila melanogaster* (see *Drosophila*). Ecdysone, bound to its **receptor** protein, directly activates primary-response genes. A subset of these genes encodes [transcription factors](#) that induce large batteries of secondary-response target genes. Dynamic changes in ecdysone titer, combined with cross-regulatory interactions among these transcription factors, determine the timing and order of regulatory **gene expression** during the onset of metamorphosis. Furthermore, connections are now being made between these genetic regulatory hierarchies and the biological responses they direct. Studies of ecdysone action in *Drosophila* have provided a foundation for understanding the molecular mechanisms of steroid hormone action during [development](#).

### 1. Ecdysone Pulses During *Drosophila* Development

At specific times during development, and in response to environmental cues, neurons in the insect central nervous system release prothoracicotropic hormone (PTTH). This [neuropeptide](#) signals the adjacent ring gland to secrete the steroid hormone ecdysone into the hemolymph of the circulatory system ([1](#)). Some peripheral tissues, most notably the larval fat body, modify ecdysone into its biologically active form, 20-hydroxyecdysone. This hormone can then directly regulate gene expression through its interaction with the ecdysone receptor. Although insect endocrinologists use the term “ecdysone” to refer to the inactive form of the hormone, this word is synonymous with 20-hydroxyecdysone in the molecular biology literature.

Ecdysone pulses direct *Drosophila* through its life cycle, with peak hormone titers signaling the major postembryonic developmental transitions (Fig. [1](#)). Larval development progresses through three **instars**, each punctuated by ecdysone-triggered molting of the cuticle. A high-titer ecdysone pulse at the end of the third instar triggers puparium formation and the onset of prepupal development. This is followed, 10–12 hours after puparium formation, by another ecdysone pulse that initiates the 3.5 days of pupal development, followed by eclosion of the adult fly (Fig. [1](#)). Most larval tissues are destroyed during the early stages of metamorphosis and are replaced by adult tissues that develop from clusters of imaginal progenitor cells. The crawling larva thus undergoes a complete transformation, resulting in the formation of a highly motile, reproductively active adult fly. Remarkably, this transformation occurs in apparent response to a single steroid hormone, ecdysone. One focus of current research is understanding how the systemic hormonal signal is refined into the appropriate stage- and tissue-specific developmental pathways that direct this transformation.

**Figure 1.** The ecdysteroid titer profile during *Drosophila* development. The composite ecdysteroid titer is depicted, in 20-hydroxyecdysone equivalents from whole-body homogenates ([44](#)). Each of the developmental stages of the *Drosophila* life cycle are also shown, below a time scale in days. Reprinted with permission from Thummel ([45](#)). Copyright held by Cell Press.



## 2. Ecdysone Regulation of Primary-Response Genes

Ecdysone manifests its effects on development by activating its receptor, a heterodimer of two members of the nuclear receptor superfamily: EcR and USP (see [Hormone Receptors](#)) (2, 3). Interestingly, at least two of the three protein **isoforms** encoded by the *EcR* gene are expressed in a tissue-restricted manner that can contribute to the specificity of developmental responses to ecdysone (4). EcR-B1 is expressed primarily in larval tissues that are fated to die, whereas EcR-A is expressed in developing adult structures and tissues. Both high-affinity DNA binding and ecdysone binding require EcR heterodimerization with USP, a homologue of vertebrate retinoid X receptor (see [Retinoic Acids](#)) (5, 6). The ecdysone/EcR/USP complex can then bind DNA and activate the transcription of target genes (see [Hormone Response Elements](#)). It is interesting to note that a *Drosophila* **orphan receptor**, DHR38, a homologue of the vertebrate NGFI-B receptor, can also heterodimerize with USP (7). This raises the possibility that ecdysone signaling may be modified through distinct heterodimer combinations of receptors, similar to the interactions that have been demonstrated in vertebrate organisms.

Many genes are induced directly by the ecdysone-receptor complex, presumably reflecting their immediate requirement following ecdysone pulses (8). Some of these genes, like *IMP-E1* (9) and *Fbp-1* (10), are induced by ecdysone in a restricted tissue-specific manner, whereas others, like *hsp23* (11) and *Eip28 / 29* (12), are expressed in multiple ecdysone target tissues.

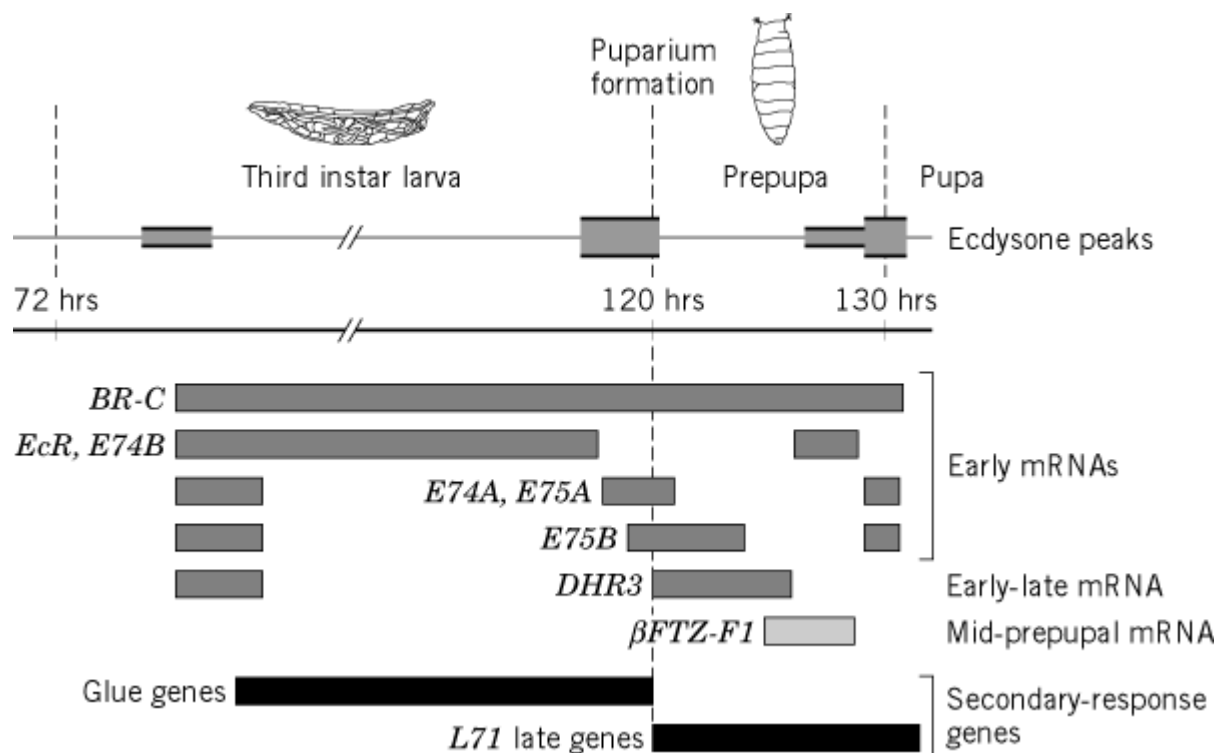
A subset of primary-response genes encode transcription factors that function at the top of ecdysone-triggered genetic regulatory hierarchies. Most of these genes were defined as **puffs** that are induced directly by ecdysone in the larval salivary gland [polytene chromosomes](#). Detailed studies by Michael Ashburner and co-workers demonstrated that these so-called early puffs appeared to encode regulatory proteins that both repress their own activity and induce large sets of secondary-response late puffs (13). The protein products of the late puffs were, in turn, thought to play a more direct role in specifying appropriate biological responses to the hormone.

Four early puff genes have been described at the molecular level. One of these, *E63-1*, encodes a calcium binding protein related to [calmodulin](#) (see [Calcium Signaling](#)), providing the possibility of cross-regulation between hormone and calcium signaling pathways (14). The remaining three genes each encode multiple transcription factor isoforms. The *Broad-Complex* (*BR-C*) encodes more than a

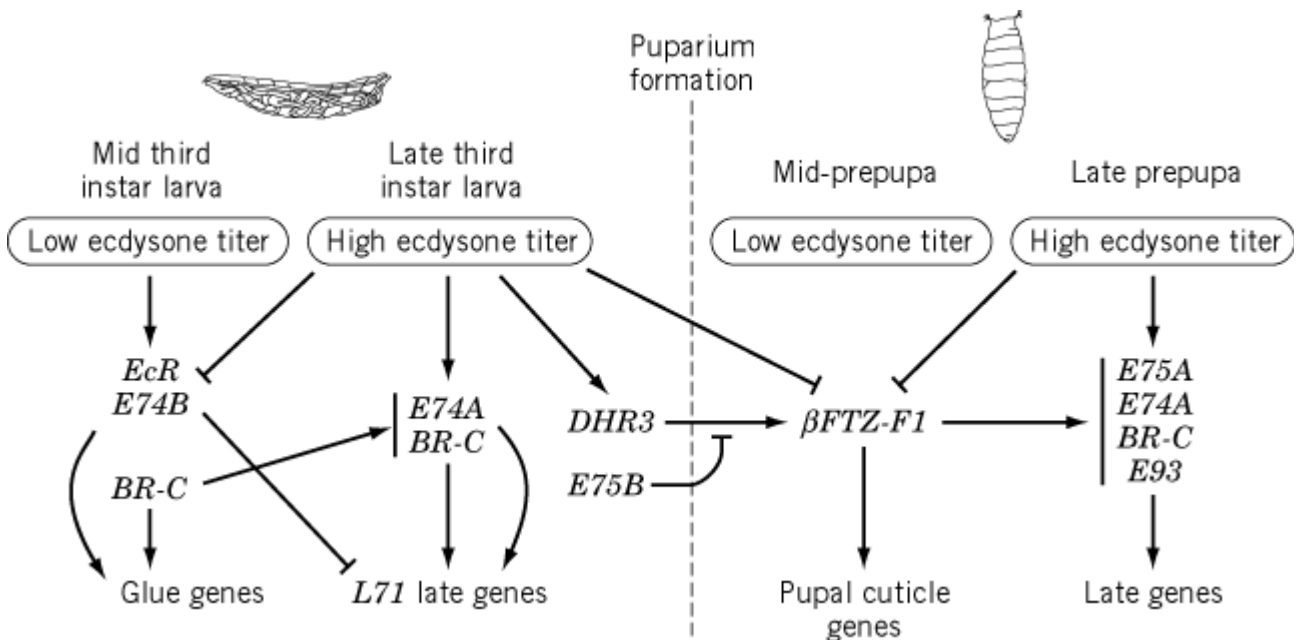
dozen proteins, each of which carries one of four possible pairs of **zinc fingers** (15, 16). The *E74* early gene contains nested **promoters** that direct the synthesis of two related proteins, *E74A* and *E74B*, that share an identical **C-terminal ETS domain** that binds DNA (17). Similarly, the *E75* early gene contains nested promoters that direct the synthesis of three orphan members of the nuclear receptor superfamily, designated *E75A*, *E75B*, and *E75C* (18). Curiously, three other orphan receptor genes, *DHR3*, *E78*, and *bFTZ-F1*, correspond to ecdysone-regulated puffs in the polytene chromosomes, suggesting that this family of transcription factors may play an important role in ecdysone signaling during development (19-21).

Ecdysone dose-response studies have revealed that each early promoter is activated at a specific critical threshold hormone concentration. Combined with ecdysone titer measurements, these dose-response studies have led to a model for the timing of early gene expression during the onset of metamorphosis (22). The *EcR*, *BR-C*, and *E74B* promoters are most sensitive to ecdysone, and are induced in mid-third instar larvae in apparent response to a low-titer hormone pulse (Figs. 2 and 3). The high-titer late larval ecdysone pulse then represses some of these **messenger RNAs**, such as *EcR* and *E74B*, as it induces the less-sensitive promoters, including *E74A*, *E75A*, *E75B*, and *DHR3* (Figs. 2 and 3) (22, 23). The ecdysone titer drops during mid-prepupal development, allowing *bFTZ-F1* to be expressed. This gene is repressed by ecdysone and thus depends on the decrease in hormone titer for its induction (24) (Fig. 3). The rapid rise in ecdysone titer in late prepupae then recapitulates the larval pattern of early gene expression, with a burst of *EcR* and *E74B* expression preceding that of *E74A* and *E75A* (Fig. 2). In addition to these dose-dependent responses to ecdysone pulses, many early genes also cross-regulate their activity, ensuring that they will be expressed at the correct time and for the appropriate duration (see text below).

**Figure 2.** Temporal patterns of ecdysone-regulated gene expression during the onset of metamorphosis. A schematic representation of the ecdysone pulses is shown at the top, with the magnitude of each pulse represented by the width of the stippled bar. Developmental time proceeds from left to right, with the major ecdysone-triggered transitions marked by dotted lines. The dotted line on the left represents the second-to-third instar larval molt, and the dotted line on the right represents head eversion and the prepupal-pupal transition. Dark grey bars show the timing and duration of primary-response regulatory gene transcription, the light grey bar represents *bFTZ-F1* transcription, and the black bars represent secondary-response gene transcription. Reprinted with permission from Thummel (46).



**Figure 3.** Multiple ecdysone-triggered regulatory hierarchies direct the onset of *Drosophila* metamorphosis. This figure summarizes regulatory interactions discussed in the text. Bars represent repressive effects, and arrows represent inductive effects. Reprinted with permission from Thummel (46).



### 3. Regulation of Ecdysone Secondary-Response Genes

Molecular and genetic studies have demonstrated that these stage-specific patterns of primary-response regulatory gene expression are transduced into corresponding patterns of secondary-response gene activity. Two classes of secondary-response genes have been studied in the larval salivary glands: the glue genes (25, 26) and the *L71* genes from the 71E late puff (27). The glue genes are induced in mid-third instar salivary glands and encode a **polypeptide** glue that is used by the animal to affix itself to a solid surface for pupariation. In contrast, the *L71* genes are of unknown function and are induced at puparium formation, as the glue genes are repressed (Fig. 2).

Both *BR-C* and *E74B* are required for proper glue gene induction, defining a mid-third instar regulatory hierarchy in the salivary glands (Fig. 3) (28). *BR-C* binding sites in the glue gene enhancers are required for proper induction, indicating that this is a direct regulatory link (29). The tissue specificity of this response is directed by the Fork head **homeodomain** transcription factor that binds adjacent to the *BR-C* in glue gene enhancers (30, 31).

The *L71* late genes also depend on *BR-C* and *E74* for their proper induction in late third instar larval salivary glands (Fig. 3). *E74B* maintains the *L71* late genes in a repressed configuration during third instar larval development. The high-titer late larval ecdysone pulse then shuts off the *E74B* repressor and induces the *E74A* activator, which, together with the *BR-C*, directly induces *L71* transcription (Fig. 3) (32-34). These studies provide molecular support for Ashburner's proposal that early puff proteins direct the induction of late puff genes. In addition, they establish an important paradigm for understanding the mechanisms of hormone signaling in vertebrates, in which regulatory hierarchies are more difficult to define.

### 4. Cross-Regulation Among Ecdysone-Induced Transcription Factors

In addition to their roles in regulating ecdysone secondary-response genes, the ecdysone-induced transcription factors also regulate one another. *BR-C* is required for maximal ecdysone-induced early gene transcription in late third instar larvae (Fig. 3) (35). In this sense, *BR-C* is functioning as a competence factor that facilitates the subsequent regulatory response to ecdysone. Similarly, the early induction of *EcR* by a mid-third instar hormone pulse allows for increased levels of ecdysone receptor in preparation for the high hormone titer at the end of larval development (4).

Evidence has also been obtained for cross-regulatory interactions among the ecdysone-regulated orphan receptor genes. *DHR3*, which is induced by ecdysone in newly formed prepupae, is an inducer of *bFTZ-F1* expression in mid-prepupae (Fig. 3) (36, 37). Furthermore, this induction can be inhibited by heterodimerization of *DHR3* with *E75B* (37). *E75B* protein levels normally decay during early prepupal development, thus determining the time at which *DHR3* can activate *bFTZ-F1* (Fig. 2). *bFTZ-F1*, in turn, appears to be a critical competence factor for maximal reinduction of the early genes (Fig. 3) (24). Furthermore, *bFTZ-F1* is sufficient to induce the stage-specific *E93* early gene in the salivary glands of late prepupae (38). The *DHR3*, *E75B*, and *bFTZ-F1* orphan receptors thus provide a regulatory link between the late larval and prepupal responses to ecdysone. Their stage-specific expression confirms that a critical developmental stage has been achieved and that the next response to the hormone should be distinct from the last response. These receptors thus provide a mechanism by which a hormonal signal can be refined into stage-specific developmental pathways.

## 5. From Gene Regulation to Biological Responses

Studies are currently under way that attempt to link ecdysone regulatory hierarchies with their appropriate biological responses during metamorphosis. As mentioned in the text above, most larval tissues are destroyed during prepupal and early pupal development. Recent studies have revealed that the destruction of the larval midgut and salivary gland occurs as stage-specific [programmed cell death](#) responses that display many of the hallmark features of [apoptosis](#) (39). Furthermore, the *Drosophila* death inducer genes *reaper* and *head involution defective* are coordinately induced in these tissues immediately preceding their destruction. A current focus of research is aimed at understanding how these death inducers are regulated by ecdysone at the appropriate time and place during metamorphosis. Specific neurons in the central nervous system also undergo programmed cell death in newly eclosed adults, in response to a decrease in ecdysone titer (40). Interestingly, these neurons can be distinguished by their high expression of the *EcR-A* isoform, suggesting that this receptor may play a key role in directing the death of these neurons.

As the larval tissues are being destroyed, the adult fly is being constructed. During prepupal development, the [imaginal discs](#) evert and elongate to form rudiments of the adult appendages. *BR-C* and *E74* both play a critical role in this process (41, 42). Studies have also shown that *bFTZ-F1* is sufficient to induce the pupal cuticle genes in mid-prepupal imaginal discs, providing a mechanism for the stage specificity of pupal cuticle deposition (Fig. 3) (43). It seems likely that future studies will provide more links between the ecdysone regulatory hierarchies and the biological responses to ecdysone during the onset of metamorphosis.

## Bibliography

1. L. I. Gilbert, R. Rybczynski, and S. S. Tobe (1996) In *Metamorphosis: postembryonic reprogramming of gene expression in amphibian and insect cells* (L. I. Gilbert, J. R. Tata, and B. G. Atkinson, eds.), Academic Press, New York, pp. 59–107.
2. A. E. Oro, M. McKeown, and R. M. Evans (1990) *Nature* **347**, 298–301.
3. M. R. Koelle et al (1991) *Cell* **67**, 59–77.
4. W. S. Talbot, E. A. Swyryd, and D. S. Hogness (1993) *Cell* **73**, 1323–1337.
5. M. R. Koelle (1992) *Ph.D. thesis*, Stanford University.
6. T. Yao et al (1993) *Nature* **366**, 476–479.
7. J. D. Sutherland, T. Kozlova, G. Tzertzinis, and F. C. Kafatos (1995) *Proc. Natl. Acad. Sci.*



USA **92**, 7966–7970.

8. A. J. Andres and C. S. Thummel (1992) *Trends Genet.* **8**, 132–138.
9. J. E. Natzle, D. K. Fristrom, and J. W. Fristrom (1988) *Dev. Biol.* **129**, 428–438.
10. J. A. Lepesant et al. (1986) *Arch. Insect. Biochem. Physiol. Supp.* **1**, 133–141.
11. E. B. Dubrovsky, G. Dretzen, and E. M. Berger (1996) *Mol. Cell. Biol.* **16**, 6542–6552.
12. A. J. Andres and P. Cherbas (1992) *Development* **116**, 865–876.
13. M. Ashburner, C. Chihara, P. Meltzer and G. Richards (1974) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 655–662.
14. A. Andres and C. S. Thummel (1995) *Development* **121**, 2667–2679.
15. P. R. DiBello et al (1991) *Genetics* **129**, 385–397.
16. C. A. Bayer, B. Holley, and J. W. Fristrom (1996) *Dev. Biol.* **177**, 1–14.
17. K. C. Burtis et al. (1990) *Cell* **61**, 85–99.
18. W. A. Segraves and D. S. Hogness (1990) *Genes and Dev.* **4**, 204–219.
19. M. R. Koelle, W. A. Segraves, and D. S. Hogness (1992) *Proc. Natl. Acad. Sci. USA* **89**, 6167–6171.
20. G. Lavorgna, F. D. Karim, C. S. Thummel, and C. Wu (1993) *Proc. Natl. Acad. Sci. USA* **90**, 3004–3008.
21. B. L. Stone and C. S. Thummel (1993) *Cell* **75**, 307–320.
22. F. D. Karim and C. S. Thummel (1992) *EMBO J.* **11**, 4083–4093.
23. M. Horner, T. Chen, and C. S. Thummel (1995) *Dev. Biol.* **168**, 490–502.
24. C. T. Woodard, E. H. Baehrecke, and C. S. Thummel (1994) *Cell* **79**, 607–615.
25. M. A. T. Muskavitch and D. S. Hogness (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7362–7366.
26. E. M. Meyerowitz and D. S. Hogness (1982) *Cell* **28**, 165–176.
27. L. G. Wright, T. Chen, C. S. Thummel, and G. M. Guild (1995) *J. Mol. Biol.* **255**, 387–400.
28. M. Lehmann (1996) *BioEssays* **18**, 47–54.
29. L. von Kalm et al. (1994) *EMBO J.* **13**, 3505–3516.
30. M. Lehmann and G. Korge (1996) *EMBO J.* **15**, 4825–4834.
31. V. Mach, K. Ohno, H. Kokubo, and Y. Suzuki (1996) *Nuc. Acids Res.* **24**, 2387–2394.
32. L. D. Urness and C. S. Thummel (1995) *EMBO J.* **14**, 6239–6246.
33. K. Crossgrove, C. A. Bayer, J. W. Fristrom, and G. M. Guild (1996) *Dev. Biol.* **180**, 745–758.
34. J. C. Fletcher, P. P. D'Avino, and C. S. Thummel (1997) *Proc. Natl. Acad. Sci. USA* **94**, 4582–4586.
35. F. D. Karim, G. M. Guild, and C. S. Thummel (1993) *Development* **118**, 977–988.
36. G. T. Lam, C. Jiang, and C. S. Thummel (1997) *Development* **124**, 1757–1769.
37. K. P. White, P. Hurban, T. Watanabe, and D. S. Hogness (1997) *Science* **276**, 114–117.
38. E. H. Baehrecke and C. S. Thummel (1995) *Dev. Biol.* **171**, 85–97.
39. C. Jiang, E. H. Baehrecke, and C. S. Thummel (1997) *Development* **124**, 4673–4683.
40. S. Robinow, W. S. Talbot, D. S. Hogness, and J. W. Truman (1993) *Development* **119**, 1251–1259.
41. I. Kiss et al (1988) *Genetics* **118**, 247–259.
42. J. C. Fletcher, K. B. Burtis, D. S. Hogness, and C. S. Thummel (1995) *Development* **121**, 1455–1465.
43. T. Murata, Y. Kageyama, S. Hirose, and H. Ueda (1996) *Mol. Cell. Biol.* **16**, 6509–6515.
44. L. M. Riddiford (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez-Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 899–939.

45. C. S. Thummel (1995) *Cell* **83**, 871–877.
46. C. S. Thummel (1996) *Trends Genet.* **12**, 306–310.

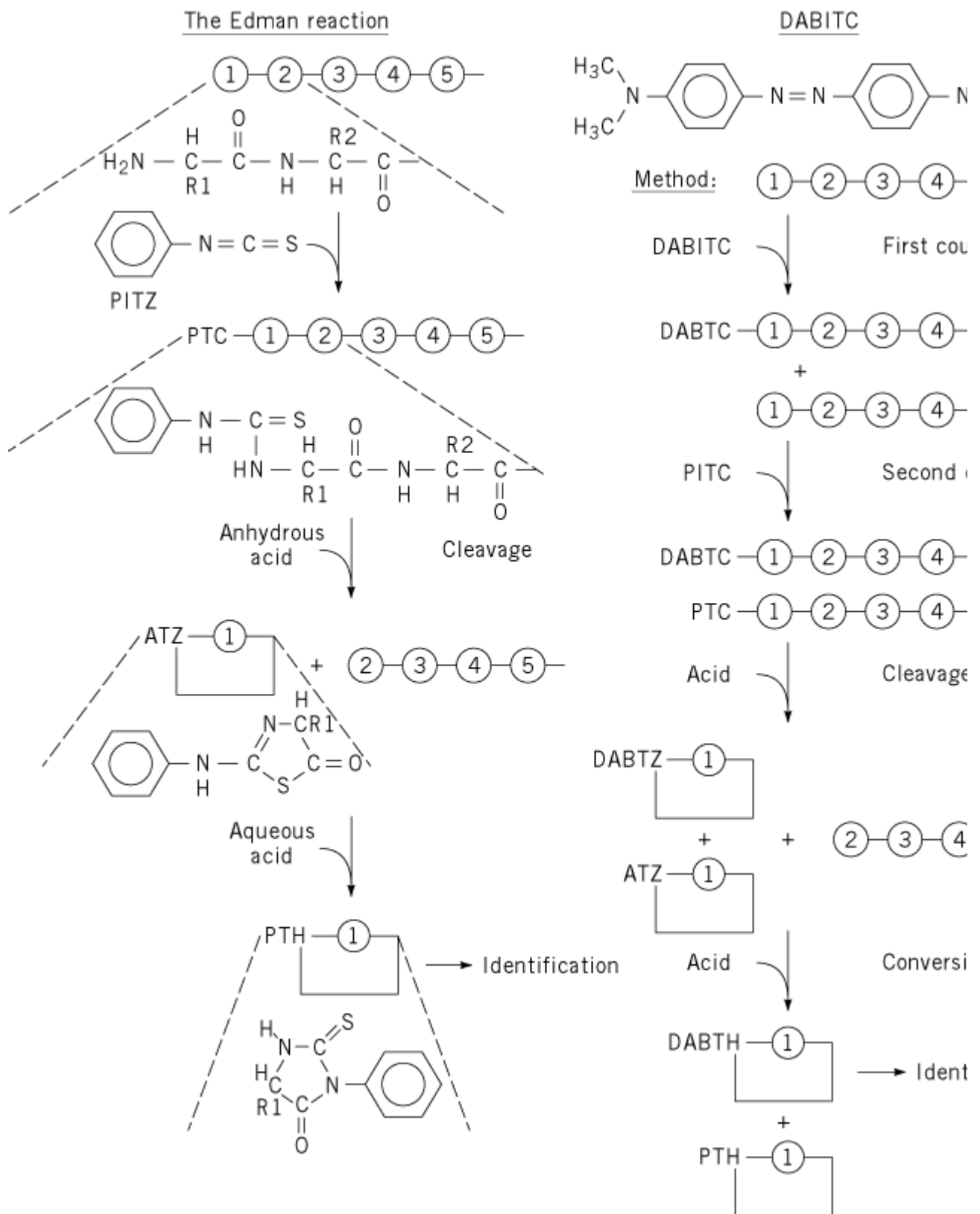
### Suggestions for Further Reading

47. L. I. Gilbert, J. R. Tata, and B. G. Atkinson (1996) *Metamorphosis: Postembryonic Reprogramming of Gene Expression in Amphibian and Insect Cells*, Academic Press, New York.
48. J. Koolman (1989) *Ecdysone, from Chemistry to Mode of Action*, Thieme Medical Publishers, Inc., New York.
49. L. M. Riddiford (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez-Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 899–939.
50. C. S. Thummel (1996) Flies on steroids—*Drosophila* metamorphosis and the mechanisms of steroid hormone action. *Trends Genet.* **12**, 306–310.

### Edman Degradation

This reaction for determining the sequences of [peptides](#) and [proteins](#) from their N-terminus was reported by Pehr Edman in the early 1950s ([1-3](#)). It now constitutes one of the oldest, if not the oldest, molecular method that is still much in use (see [Protein Sequencing](#)). It uses a three-stage reaction cycle for step-wise removal of [amino acid](#) residues from the N-terminus of a polypeptide chain. The three stages of the Edman cycle (Fig. [1](#)) are (i) coupling, (ii) cyclization, and (iii) conversion.

**Figure 1.** Principles of the Edman reaction (left). Also shown is the DABITC (right) double-coupling method. Key: PIT, phenylisothiocyanate; ATZ, anilinothiazolinone; PTH, phenylthiohydantoin; Dansyl (or DNS), 1-dimethylaminonaphtha sulphonyl; DABITC, 4-N,N-dimethylaminoazobenzene-4'-isothiocyanate.



(i) *Coupling* involves a nucleophilic attack of the N-terminal  $\alpha$ -amino group on the thiocyanate carbon of phenylisothiocyanate (PITC), to form the phenylthiocarbonyl-peptide derivative (PTC-peptide).

(ii) *Cyclization* constitutes the formation of the anilinothiazolinone derivative (ATZ-amino acid) in

anhydrous acid, thereby liberating the N-terminal amino acid residue in a cyclic form, and leaving the remaining peptide truncated by one residue (Fig. 1).

(iii) *Conversion* involves the rearrangement of the liberated ATZ derivative to the corresponding phenylthiohydantoin derivative (PTH-amino acid) by opening of the CS bond and re-closure with the CO bond.

The reaction of protein amino groups with cyanate was known before Edman's report, but the degradative yield using cyanate was low. Edman increased the yields to useful levels by introducing thiocyanates in place of cyanates, and by determining the exact reaction conditions and pointing out the need for pure chemicals.

This reaction constitutes the start of protein chemical sequencing on a routine scale; hence it really marks the start of molecular biology as a whole. The reaction is still used in most protein laboratories worldwide, it is the basis of chemical protein sequencers, and a classical method still “modern” in molecular biology at the end of the 1990s. Essentially, the only changes that have occurred during the 50 years to this reaction are:

- use of other solvents, giving alternative and more rapid protocols,
- continued automation, miniaturization, and increase in reaction speed, now making the reaction useful for routine applications at the picomole level in 30-min cycles per residue.

It was a remarkable accomplishment of Edman to design this reaction and to directly set out all the conditions, so that it has stayed unsurpassed for 50 years of worldwide protein chemistry.

## 1. Protein Sequencers

After development of the Edman reaction, Edman automated the process, developing a machine to make all additions, extractions, and lyophilisations. This constituted the birth of the protein sequencer (or sequencer); the basic concept for a complete machine was published by Edman and Begg in 1967 in a now classic paper (4) in the then newly-started *European Journal of Biochemistry*.

In the first-generation sequencers, a “spinning cup” constituted the reaction center, and the sample was kept in place by centrifugal force, making extractions and reactions possible. Extraction losses were a problem in each cycle, however, soon leading Laursen and others to the idea of solid-phase attachment of the peptide to be degraded (see **Solid-phase synthesis**). Consequently, only a few years later, the solid-phase sequencer was reported (5). Although both types of sequencers were based on the same Edman chemistry (3), their different properties, including the different attachments of the samples, made them suitable for different analyses. The liquid-phase sequencer was excellent for proteins, and the solid-phase instrument was better for short peptides, because of their different sensitivities to extractive losses and to build-up of background signals.

## 2. Manual Edman Degradations

Once the Edman reaction and the protein sequencer existed, the really limiting factor was peptide purification on the one hand, and phenylthiohydantoin (PTH)-amino acid identification on the other. For a long time, suitable chromatographic methods did not exist, hampering both purification and identification. Peptide purification, especially the separation of large peptides, was a problem. Similarly, identification of the PTHs liberated in each Edman cycle was also a problem. The PTH derivatives were initially identified by paper chromatography, which was not fully reliable (all derivatives did not separate in one analytical step) and also required the use of several solvent systems, which wasted time and material. Subsequent developments introduced gas chromatography and **thin-layer chromatography**, both increasing the speed of identification, but still not giving a reliable one-step identification of all amino acids. These limitations contributed to the development

of alternative methods, such as the dansyl and DABITC types of analysis, which did not require a major investment in equipment. One major manual method was that using dansyl (1-dimethylaminonaphthalene-5-sulphonyl) chloride to detect the new N-terminus after each Edman cycle (6) (see [Dansyl Chloride](#)). The main advantage was that it was easier to identify the residues released sequentially. Its primary importance is that it made protein sequence determination accessible to many laboratories.

A similar and later development in manual sequence analysis was the DABITC (4-N,N-dimethylaminoazobenzene-4'-isothiocyanate) method (7). In this case, protein degradation is carried out by coupling with the strongly colored DABITC, in place of PITC (Fig. 1). The DABTHs produced, which correspond to PTHs, are easily detectable by **thin-layer chromatography** (7). However, DABITC has a low coupling yield in the Edman reaction, necessitating a second coupling stage with ordinary PITC (Fig. 1, right) before the subsequent cyclization step.

At the time, both the dansyl and the DABITC methods were important, but since the 1980s these methods have gradually decreased in importance because of reliable sequencer on-line **HPLC** identification of PTH-derivatives (see below).

### 3. HPLC

The PTH-amino acid identification problem was finally solved with the introduction of **HPLC** in 1976 (8); subsequently, PTH amino acid identification became rapid, routine, and reliable. Although some identifications may still be difficult, PTH identification no longer constitutes the limiting factor in time or reliability.

At the same time, HPLC separations and the subsequent development of a whole battery of different chromatographic media soon also solved the problem of purification of peptide fragments from proteolytic digests.

### 4. Second-Generation Chemical Sequencers: Automation of All Steps

With the development of protein sequence analyzers, it became possible to determine amino acid sequences routinely, and there was an exponential increase in known protein sequences. However, analysis at this stage was still time-consuming and nonautomatic, requiring knowledge and real research. Gradually, however, a set of further inventions increased the automation and brought the actual sequence analysis more or less to the present-day automatic stage. This development essentially relied on three further inventions/improvements.

One concerned the attachment of the protein/peptide to the sequencer for analysis. This important step had several sub-steps. One was the realization that an organic cationic polymer, Polybrene (9, 10), was a suitable material capable of binding proteins and peptides to glass surfaces or other membranes. The resulting minimization of extraction losses in the washing steps of each cycle made degradations possible through to the very C-terminus of most peptides. More importantly, apart from the improvement in the lengths of degradation possible, use of Polybrene also meant that the peptide for degradation could be attached to surfaces; this therefore opened the way to the abandonment of the use of centrifugal force for attachment of the peptide for degradation. Instead, peptides could now be attached to membranes, essentially moving the degradation from the traditional "liquid phase" system to the advantageous "solid phase" (which had been started earlier by covalent attachments, cf. above) and "gas-phase" systems (11). The latter are solid phase for the attachment of the peptide to a support, and gas phase for the introduction of some reagents. These approaches are still in use and now allow rapid and sensitive analysis. Many other sub-steps in this transformation were involved. In particular, perhaps the introduction of chemical attachments of peptides to membranes (12), and the successive development of alternative, miniature column attachments for sample introductions and preparations (13), should be mentioned.

The second improvement at this stage, was the introduction of dead volume-free valve blocks, allowing use of valves with an absolute absence of cross-contamination between the reagents used in the reaction (14). In this manner, further increases in speed (because of less washing) and sensitivity (because of higher yield from lack of cross-contamination) made sequencers still more useful and rapid. Recently, this has been carried still further, and soon “chip-based blocks” may be encountered in sequencers, with all solvent delivery and removal stages in extremely small volumes of “chip-blocks” (15).

The third major advance at this stage was the introduction of automatic methods for the conversion of a thiazolinone from each cycle into the corresponding thiohydantoin. This became possible because of the introduction of a second reaction vessel (16), with separate reagents and reactions.

Once the conversion had become automated, it became possible also to link the subsequent PTH identification step to the cycle of automatic events in the sequencer, thus opening the road to on-line identification of the liberated amino acid derivative in each cycle. These on-line modes were started very early, with the introduction of HPLC, and were soon commercialized and perfected in a new set of complete sequencers, starting with the “gas-phase sequencer” that was available in the 1980s (11). Soon, the on-line approach was coupled with post-PTH-identification data treatments, allowing extensive computer interpretation at each step. All chromatograms can now be stored and compared with on-line computers and further interpreted and related to sequences in **datbanks**, analyzed in modeling programs, and submitted to further computerized adjustments.

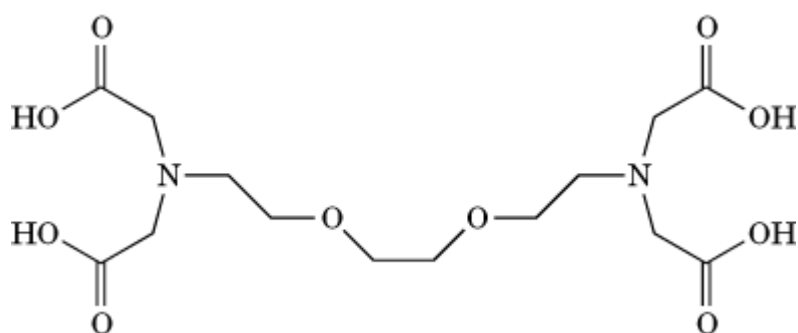
#### Bibliography

1. P. Edman (1950) *Acta Chem. Scand.* **4**, 283–293.
2. P. Edman (1953) *Acta Chem. Scand.* **7**, 700–701.
3. P. Edman (1956) *Nature* **177**, 667–668.
4. P. Edman and G. Begg (1967) *Eur. J. Biochem.* **1**, 80–91.
5. R. A. Laursen (1971) *Eur. J. Biochem.* **20**, 89–102.
6. W. R. Gray (1967) *Methods Enzymol.* **11**, 469–475 and (1972) **25**, 121–138.
7. J. Y. Chang, D. Brauer, and B. Wittmann-Liebold (1978) *FEBS Lett.* **93**, 205–214.
8. C. L. Zimmerman, E. Appella, and J. J. Pisano (1977) *Anal. Biochem.* **77**, 569–573.
9. G. E. Tarr, J. F. Beecher, M. Bell, and D. J. McKean (1978) *Anal. Biochem.* **84**, 622–627.
10. D. G. Klapper, C. E. Wilde III, and J. D. Capra (1978) *Anal. Biochem.* **85**, 126–131.
11. R. M. Hewick, M. W. Hunkapiller, L. E. Hood, and W. J. Dreyer (1981) *J. Biol. Chem.* **256**, 7990–7997.
12. J. M. Coull, D. J. C. Pappin, J. Mark, R. Aebbersold, and H. Köster (1991) *Anal. Biochem.* **194**, 110–120.
13. C. G. Miller (1994) *Methods: A companion to Methods in Enzymology* **6**, 315–333.
14. B. Wittmann-Liebold, A. W. Geissler, and E. Marzinzig (1975) *J. Supramol. Struct.* **3**, 426–447.
15. C. Wurzel and B. Wittmann-Liebold (1998) In *Microreaction Technology: Proceedings of the First International Conference on Microreaction Technology* (W. Ehrfeld, ed.), Springer Verlag, Berlin, Heidelberg, pp. 219–224.
16. W. Machleidt and H. Hofer (1980) In *Methods in Peptide and Protein Sequence Analysis* (C. Birr, ed.), Elsevier/North Holland, Amsterdam, pp. 35–47.

#### **EDTA (Ethylenediamine-*N,N,N',N'*-Tetraacetic Acid)**

EDTA is a chelating agent widely used, along with [EGTA](#), to control the concentration of divalent and trivalent cations in aqueous solutions. Its structure is given in Figure 1. It is an essential ingredient in such diverse products as laundry detergents and film developers. In the field of molecular biology, EDTA is present in many buffer solutions used for experiments where the presence of divalent and trivalent cations would interfere with the desired outcome. For example, it is used in most [buffers](#) for the initial steps of protein purification (1) from biological material, to suppress the activity of metal-dependent proteases that would otherwise degrade the desired protein.

**Figure 1.** Chemical structure of EDTA (ethylenediamine- *N,N,N',N'*-tetraacetic acid).



The flexibility of the molecule, and the variety of potential liganding sites within it, permit EDTA to accommodate a wide variety of cation types, preferred ligand geometries and bond lengths, thus contributing to its high affinity for most cations (Table 1) (2). In contrast to many other complexing agents, EDTA will form only a 1:1 complex with either divalent and trivalent cations, in which a single ion is chelated by all four carboxyl groups, which greatly contributes to the stability of the resulting complex. Hence, EDTA can maintain the concentration of free cations at very low levels. The high affinity of EDTA for cations is also used to removed metals from [metalloproteins](#) and to maintain them in their metal-free state. It is important to note that the binding of cations by EDTA is pH-dependent; protonation of the carboxylate moieties at low pH dramatically lowers their affinity for cations.

**Table 1. Affinity of EDTA for Various Cations<sup>a</sup>**

| Metal Ion        | $\log_{10} K_a$ |
|------------------|-----------------|
| Mg <sup>2+</sup> | 8.64            |
| Ca <sup>2+</sup> | 11.0            |
| Zn <sup>2+</sup> | 16.4            |
| Cd <sup>2+</sup> | 16.4            |
| Mn <sup>2+</sup> | 13.8            |
| Cu <sup>2+</sup> | 18.7            |
| Fe <sup>2+</sup> | 14.33           |

|                  |       |
|------------------|-------|
| Fe <sup>3+</sup> | 24.23 |
| Co <sup>2+</sup> | 16.3  |
| Co <sup>3+</sup> | 36    |
| Pb <sup>2+</sup> | 18.3  |

---

<sup>a</sup> The logarithm of the association constant [ $K_a$  in units of  $M^{-1}$  (reciprocal molar)] are compared at pH 7.0, 25°C, and 0.1 M ionic strength (2).

EDTA as the free acid is a white powder, has a molecular weight of 292.25 Da, decomposes at 240°C, and is relatively insoluble in water (0.34 g/100 mL at 25°C). The sodium and potassium salts of EDTA are also white powders and are much more readily soluble in water.

#### Bibliography

1. J. E. Coligan et al. (eds), *Current Protocols in Protein Science* John Wiley & Sons, Chichester, U.K. (1995).
2. A. E. Martell and R. M. Smith (1974) *Critical Stability Constants* Vol. 1, Plenum Press, New York.

#### Suggestion for Further Reading

3. R. H. Crabtree (1994) "Coordination + organometallic chemistry, principles", in *Encyclopedia of Inorganic Chemistry*, Vol. 2, R. B. King, ed., John Wiley & Sons, Chichester, U.K.

#### Effective Molarity

A reaction or interaction between two different molecules in dilute solution is relatively straightforward and is governed by their concentrations in that solution, plus the appropriate rate or equilibrium constant (see [Kinetics](#)). In more complex situations, however, the situation may differ from that expected from the bulk concentrations of two reactants. If they attract each other, for example by **electrostatic** or **hydrophobic interactions**, their *effective* molarities will be greater than their bulk concentrations. Similarly, a [nonpolar](#) reagent that partitions into the interiors of [membranes](#), or any reagent that concentrates in one organelle or compartment, will react there at a much greater rate than if it were dispersed evenly throughout the system.

A special situation occurs when two reactive or interactive groups are part of the same [macromolecule](#), such as a **protein** or **nucleic acid**. The rate constant for reaction between two such groups, or their association constant, might be known when they are individual molecules, when the observed rate or equilibrium is given by this constant times their respective concentrations, but for an intramolecular reaction or interaction, their concentration within the solution is irrelevant. Instead, it is the *effective concentration* of the two groups within the macromolecule that is important for the intramolecular interaction. The effective molarity will depend upon the structure of the macromolecule, in particular the extent to which it brings the two reactants together or keeps them apart, plus the environment in which the two groups are kept. The effective molarity for the intramolecular reaction will be essentially zero if the macromolecular structure keeps the two groups



apart (of course, groups on different molecules can still react, and their bulk concentrations still govern this *intermolecular* reaction). At the other extreme, when the macromolecular structure keeps the two groups in the correct proximity and orientation for reaction, their effective molarity can be extremely large, up to  $10^{10}$  M, concentrations that are not feasible with two independent molecules. Such values are predicted by theoretical considerations (1) and are also observed experimentally, from the ratio of the rate or equilibrium constants for the same reaction when the groups are on the same molecule and on separate molecules (2).

The large effective concentrations that are observed are believed to be due primarily to the much smaller loss of entropy that occurs when two groups on the same molecule interact. Two independent molecules that interact in solution must lose substantial translational and rotational entropy when they interact, depending upon the rigidity of the interaction. Two groups attached to the same macromolecule have already lost varying degrees of this entropy. In the ideal case, when the two reactive groups are held by the macromolecular structure in precisely the correct position for them to interact, so that there is no change in their entropy, the maximum effect and the maximum effective molarity are observed. Of course, other considerations also apply, such as if the environment of the groups is different in the macromolecule, or if they are strained and this strain is relieved upon their reaction.

Whatever the exact reason for the large effective molarities that can be measured in macromolecules, this has important consequences for understanding the stabilities of the folded structures of macromolecules. For example, the question of the role of [hydrogen bonds](#) in stabilizing protein structures has been very controversial (see [Protein Stability](#)). One might expect *a priori* that hydrogen bonds would contribute nothing to the net stability because any hydrogen bonds in the folded state would be replaced by equivalent hydrogen bonds, with the [water](#), in the [unfolded protein](#). But the two cases differ because the hydrogen bonds in the folded state are intramolecular, whereas those with the solvent are intermolecular. If the hydrogen bonding groups have effective molarities in the folded state greater than the solvent concentration, 55 M for water, the hydrogen bonds in the folded state will be more stable than those in the unfolded state, and they will stabilize the folded conformation. Another way of thinking about this is to consider the increased entropy of the water molecules that are liberated when the protein folds and forms intramolecular hydrogen bonds. Effective molarities of up to  $10^5$  M have been measured between [cysteine](#) residues in forming [disulfide bonds](#) during **protein folding** (3). Consequently, it is not surprising that folded macromolecular structures can be stabilized by many intramolecular interactions that individually are very weak in an unfolded conformation, but occur simultaneously and cooperate to generate a folded conformation (4). The large effective molarities of reactive groups in the [active sites](#) of enzymes are also important for explaining enzymatic [catalysis](#).

#### Bibliography

1. M. I. Page and W. P. Jencks (1971) *Proc. Natl. Acad. Sci. USA* **68**, 1678–1683.
2. A. J. Kirby (1980) *Adv. Phys. Org. Chem.* **17**, 183–278.
3. T. E. Creighton and D. P. Goldenberg (1984) *J. Mol. Biol.* **179**, 497–526.
4. T. E. Creighton (1983) *Biopolymers* **22**, 49–58.

#### Suggestion for Further Reading

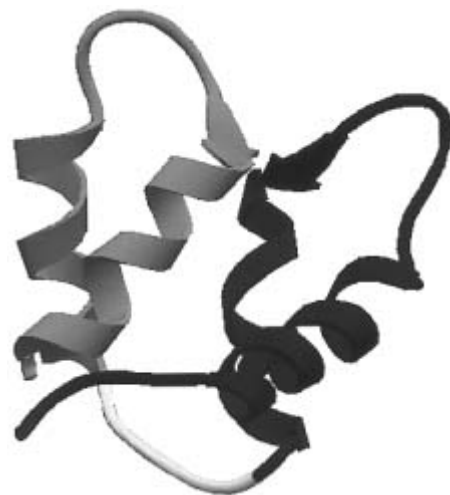
5. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York, "Chapters 4, 7", and "9".

#### EF-Hand Motif

The majority of [calcium-binding proteins](#) involved in intracellular  $\text{Ca}^{2+}$  [signal transduction](#) are characterized by a common helix–loop–helix structural [protein motif](#) in their  $\text{Ca}^{2+}$ -binding sites that is termed the EF-hand. Well over a thousand EF-hand motif sequences have been identified. An exquisite fine-tuning of the packaging of EF-hands provides these proteins with both high selectivity for  $\text{Ca}^{2+}$  and the ability to respond rapidly and efficiently to the modest (100-fold) changes in concentration associated with  $\text{Ca}^{2+}$  signals. In fact, the unique relationship between the EF-hand motif and  $\text{Ca}^{2+}$  binding extends beyond the ability to transduce  $\text{Ca}^{2+}$  signals, as other members of this protein family have critical roles in the uptake, transport and homeostasis of  $\text{Ca}^{2+}$ .

The term “EF-hand motif” was introduced by Kretsinger on examination of the three-dimensional structure of [parvalbumin](#), the first member of the EF-hand  $\text{Ca}^{2+}$ -binding protein family whose three-dimensional structure was determined (1). A great deal of structural insight into the variability of the EF-hand motif has been obtained from subsequent [X-ray crystallography](#) and [NMR](#) solution structure analysis of EF-hand calcium-binding proteins. EF-hand motifs almost always occur in pairs packed together in a face-to-face manner, forming a stable globular domain (Fig. 1). The pairing of sites is presumed to stabilize the protein conformation, increase the  $\text{Ca}^{2+}$  affinity of each site over that of isolated sites, and provide a ready means for the **cooperativity** in the binding of  $\text{Ca}^{2+}$  that is critical to their function. The conformations of EF-hands are variable and are dependent on whether ions are bound (see [Calcium-Binding Proteins](#)).

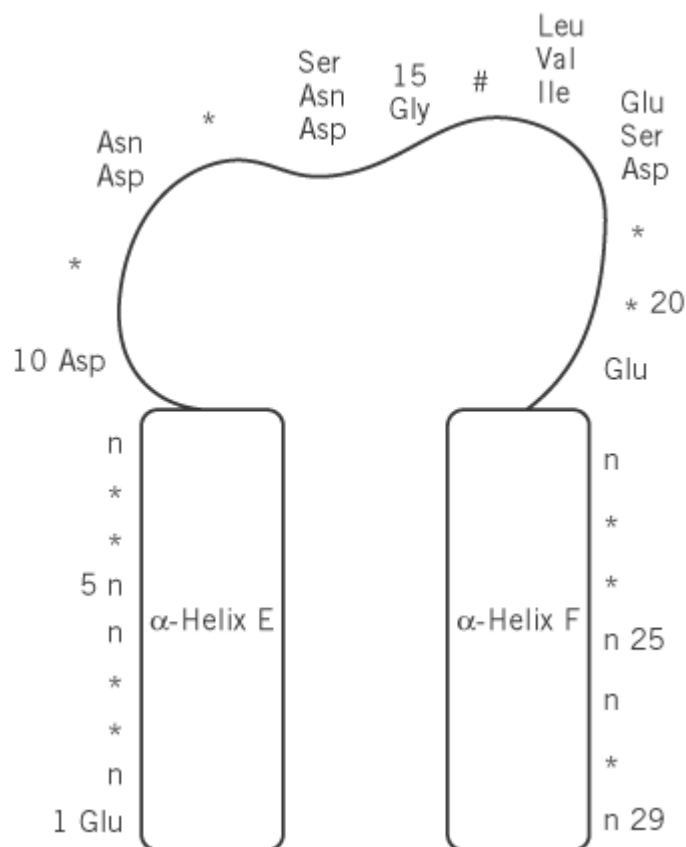
**Figure 1.** Ribbon diagram of the canonical EF-hand of the *N*-terminal domain of calmodulin using the coordinates of 1CFC (3).



The canonical EF-hand consists of a 29-residue contiguous segment of [polypeptide chain](#) containing an  $\alpha$ -helix designated as E in the original parvalbumin structure (residues 1–10), a loop around the calcium ion (residues 10–21), and a second  $\alpha$ -helix designated F (residues 19–29). The consensus sequence is shown in Figure 2. Residues 10, 12, 14, 16, 18, and 21 chelate the  $\text{Ca}^{2+}$  ion. Residues 10, 12, and 14 provide monodentate oxygen ligands, and residue 21 provides a bidentate oxygen ligand, all via side-chain carboxyl groups. Residue 16 directly coordinates  $\text{Ca}^{2+}$  via its main-chain oxygen. Residue 18 ligates  $\text{Ca}^{2+}$  indirectly, through a stably bound water molecule. Other residues have been assigned specific roles in stabilizing the EF-hand motif. For example, a central

**hydrophobic** core is formed from the side chains of residues 2, 5, 6, and 9 of helix E; 22, 25, 26, and 29 of helix F; and 17 from the binding loop.

**Figure 2.** Schematic diagram of the consensus sequence for the EF-hand  $\text{Ca}^{2+}$ -binding motif. Residues marked with \* have no constraints on which amino acids occur there. Residue 16, #, chelates the calcium by its backbone carbonyl oxygen atom. The residues of the E and F helices marked “n” must be nonpolar for packing of the two helices.



Most EF-hand proteins use this canonical motif. The S100 proteins (an important subfamily of  $\text{Ca}^{2+}$ -binding proteins), however, have a two-residue insertion and a modified coordination scheme in the first binding loop. In these pseudo-EF-hands, residues 10, 13, 15, and 18 chelate the  $\text{Ca}^{2+}$  ion via their main-chain carbonyl oxygen atoms, whereas residue 20 chelates indirectly via a water molecule, and acidic residue 23 provides bidentate ligation via its side-chain carboxylate. The role of these pseudo-EF-hands in modulating the structure or function of S100 proteins is still under investigation, although there is evidence that the packaging of a pseudo-EF-hand and a canonical EF-hand attenuates the changes in the structure and dynamics of the protein induced by the binding of calcium (2).

#### Bibliography

1. R. H. Kretsinger and C. E. Nockolds (1973) *J. Biol. Chem.* **248**, 3313–3326.
2. N. J. Skelton, J. Kördel, M. Akke, S. Forsén, and W. J. Chazin (1994) *Nature Struct. Biol.* **1**, 239–245.
3. H. Kuboniwa, N. Tjandra, S. Grzesiek, H. Ren, C. B. Klee, and A. Bax (1995) *Nature Struct. Biol.* **2**, 768–776.

#### Suggestion for Further Reading

4. N. C. J. Strynadka and M. N. G. James (1989) Crystal structures of the helix-loop-helix calcium-binding proteins, *Annu. Rev. Biochem.* **58**, 951–998.

## EGF Motif

The EGF motif is a common structural **domain** found in many [protein structures](#), especially those of proteins associated with [blood clotting](#), fibrinolysis, neural development, and **cell adhesion**. The motif is characterized as a protein module of ~ 45 residues having six conserved [cysteine](#) residues (designated CI through CVI) that form three internal [disulfide bonds](#) with the connectivity CI-III, CII-IV and CV-VI. The motif was first found in [epidermal growth factor](#) (EGF), hence the name EGF motif or EGF-like domain. Although found as a single unit in EGF and [transforming growth factor](#)  $\alpha$ , the EGF domain is often present as a repeating unit in much larger proteins. For example, the protein fibrillin has over 30 EGF domains. Furthermore, EGF motifs in proteins are frequently present in combination with repeating units of other modules such as the [Kringle domain](#).

The sequence of the EGF motif is often recognizable from just the [primary structure](#) of a protein. The consensus sequence is X-X-X-X-Cys-X(2-7)-Cys-X(1-4)-(Gly/Ala)-X-Cys-X(1-13)-t-t-a-X-Cys-X-Cys-X-X-Gly-a-X(1-6)-Gly-X-X-Cys-X, where X is any amino acid (the number in brackets defining a variable number of X residues), *a* is an aromatic residue, and *t* is a **nonhydrophobic** residue (1). The [tertiary structure](#) of the EGF motif is also highly conserved and can be described as having a two-stranded antiparallel [b-sheet](#) scaffold upon which the three disulfide bonds are positioned (Fig. 1). The disulfide bonds link the *N*-terminal and *C*-terminal loop regions to the core. Binding of **calcium** ions to some classes of EGF domain further stabilizes the *N*-terminal region and is thought to cause the formation of long helical arrangements of multiple EGF motifs (2); this property may account for their proposed function of mediating [protein-protein interactions](#).

**Figure 1.** Schematic representation of the backbone structure of an EGF domain (2). b-Strands are shown as arrows, and the three disulfide bonds are shown in light gray, with the sulfur atoms depicted as spheres. The *N*- and *C*-termini are labeled. This figure was generated using Molscript (3) and Raster3D (4, 5).



### Bibliography

1. I. D. Campbell and P. Bork (1993) *Curr. Opin. Struct. Biol.* **3**, 385–392.
2. Z. Rao, P. Handford, M. Mayhew, V. Knott, G. G. Brownlee, and D. Stuart (1995) *Cell* **82**, 131–141.
3. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
4. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
5. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestion for Further Reading

6. I. D. Campbell and P. Bork (1993) Epidermal growth factor-like modules. *Curr. Opin. Struct. Biol.* **3**, 385–392.

## Egg

An *egg* is the female **gamete**, which is originally [haploid](#) and in most cases nonmotile. After [fertilization](#) by a [sperm](#), the egg provides the basis for producing a completely new individual. Regarding this subsequent [development](#), eggs belong to the least restricted cell types of the body, as they have the power to produce every other differentiated cell type. Eggs cannot be regarded as nonspecialized cells; instead, they prove to be very specialized cells to accomplish this task.

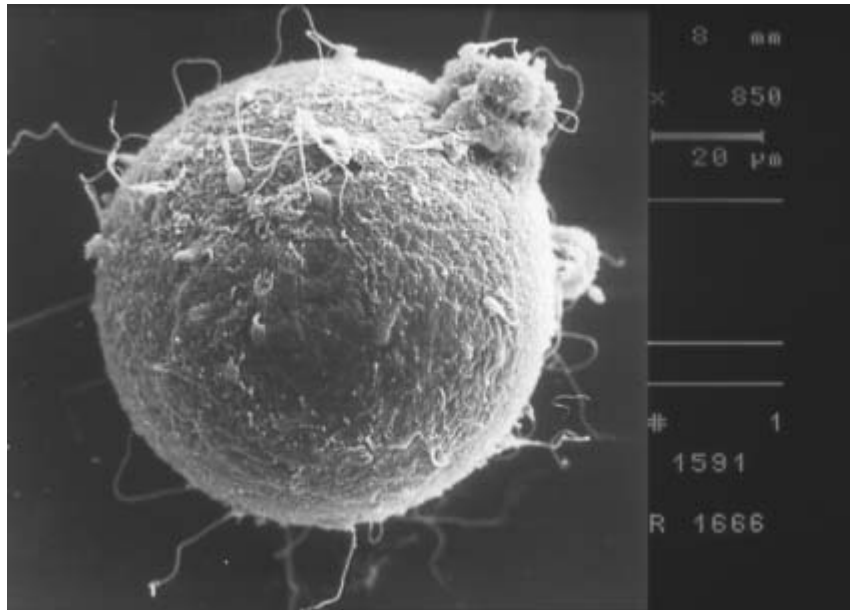
First, they must mediate sperm–egg interactions. This means that they must be able to induce metabolic pathways in the spermatozoon and biochemical signals from the spermatozoon must be understood, thus leading to intracellular reactions in the egg to promote fertilization. Second, after fertilization with a spermatozoon, the egg needs a machinery to promote biochemical actions that are necessary to fuse the male and female **pronuclei**. Third, the egg must be able to guide and support the very early cell divisions without, or with only minimum, nuclear control of the developing embryo. Thus, the egg must provide appropriate amounts and species of many different basic elements, such as [transfer RNA](#), [ribosomes](#), and many different [proteins](#), including basic elements like [transcription factors](#). These gene products have been **transcribed**, and in some cases **translated**, previously and are stored in the oocyte cytoplasm.

## 1. Morphology

Unlike other cells, eggs consist of an outer porous [extracellular matrix](#), the *zona pellucida* (1). The zona pellucida is not primarily for mechanical protection of the egg, but its foremost task is to mediate sperm–egg interactions. It is made of long filaments consisting of three types of **glycoproteins**. In all species thus far investigated, zona pellucida glycoproteins are expressed by the ZPA, ZPB, and ZPC [gene family](#) (1). [Post-translational modification](#) of the zona pellucida glycoproteins, however, varies with the species. Most of the investigations concern the situation in the mouse. Regardless of their genetic origin, but for functional reasons, the processed glycoproteins are divided into three classes. Murine zona pellucida glycoprotein 1 (ZP1) has a molecular weight of 200 kDa and is assumed to be responsible for the structural integrity of the matrix. Zona pellucida glycoproteins 2 and 3 (ZP2, ZP3) exhibit apparent molecular weights of 99 kDa and 83 kDa, respectively; both act as sperm receptors (2). In the ZP3 receptor, an [O-linked oligosaccharide](#) of 3.9 kDa molecular weight plays a central role in recognition and binding of sperm. After induction of the cortical reaction by fertilization, both receptors are converted to their inactive forms, ZP2f and ZP3f, thereby providing a secondary block against polyspermy. In species with extracorporal fertilization, such as sea urchins, the zona pellucida is surrounded by a jelly layer. In mammals, a number of corona cells adhere to the zona pellucida at the time of fertilization.

Between the oocyte plasma membrane and zona pellucida is a region that is termed the perivitelline space. On entrance of a sperm and its fusion with the oocyte plasma membrane, calcium is released from intracellular calcium stores (see [Calcium Signaling](#)). As a consequence, the cortical granules are activated and release their proteolytic contents into the perivitelline space, which subsequently diffuse into the zona pellucida. This results in the conversion of ZP2 and ZP3 into their respective inactive forms ZP2f and ZP3f. Calcium release is necessary not only to prevent fertilization from further spermatozoa, but it is also involved in (1) nuclear envelope breakdown, (2) triggering intracellular **phosphorylation** events, and (3) activation of [cell cycle](#) -control proteins leading to **mitosis** after syngamy.

**Figure 1.** Electron microscope picture of an oocyte. The oocyte is surrounded by spermatozoa that aim to penetrate the zona pellucida. (With the friendly permission of Dr. H.-W. Michelmann.)



Within the oocyte plasma membrane, the egg is filled with cytoplasm that contains basic metabolites necessary for the initial divisions of the cell. Unlike the cortical granules, not all compounds of the cell cytoplasm are located uniformly. This can be seen by light [microscopy](#) and has led to a distinction of the egg's compartments. In the frog, for example, macroscopic partitioning is visible: The nucleus is located at the **animal pole**, whereas yolk is located at the vegetal pole.

Depending on the species, eggs vary considerably in their sizes and shapes. In mammals, nutrition can be provided during the very early stages of life by diffusion and then by active transport after implantation of the embryo. Therefore, there is no need for substantial nutrition in the egg and, hence, the entire egg is no larger than 0.1 mm. In species in which extracorporeal development of the offspring is the rule, the situation is different. The sizes of eggs can range up to a couple of inches, due to large amounts of yolk, as is the case, for example, with birds.

#### Bibliography

1. J. D. Harris, D. W. Hibler, G. K. Fontenot, K. T. Hsu, E. C. Yurewicz, and A. G. Sacco (1994) DNA-Seq. **4**, 361–393.
2. P. M. Wassarman (1990) J. Reprod. Fert. Suppl., **42**, 79–87.

#### Eglin C

Eglin c is a widely studied, strong inhibitor of serine proteinases [see [Serine Proteinase Inhibitors, Protein](#)]. Its reactive site  $P_1$  Leu<sup>45</sup> residue endows it with specificity toward many, but not all, subtilisins, chymotrypsins, and elastases. It is an exceptionally strong inhibitor ( $K_1 < 10^{-13}$  M) of *Streptomyces griseus* proteinases A and B. As an efficient inhibitor of human leukocyte elastase, cathepsin G, and proteinase 3, the three human leukocyte enzymes responsible for the lung damage in emphysema and in related diseases, eglin c is of great interest to pharmaceutical companies. Eglin

c consists of a single polypeptide chain of 70 amino acid residues. In spite of its lack of disulfide bridges, it is a standard-mechanism canonical inhibitor. It is a member of the potato I family, a family of inhibitors whose members either lack disulfides or contain only one. Most members of the potato I family are isolated from plants. In contrast, eglin c was isolated from the leech *Hirudo medicinalis*, where its physiological function is not known.

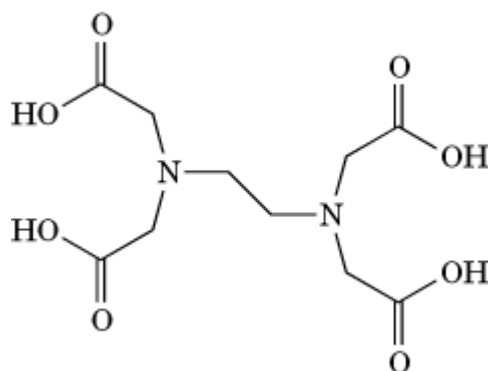
#### Suggestion for Further Reading

U. Seemüller, J. Dodt, E. Fink, and H. Fritz (1986) "Proteinase inhibitors of the leech *Hirudo medicinalis* (hirudins, bdellins, eglins)". In *Protease Inhibitors* (A. Barrett and G. Salveson, eds.), Elsevier New York, pp. 337–360.

### EGTA (ethyleneglycol bis(b-aminoethyl ether)-*N,N,N',N'*-tetraacetic acid)

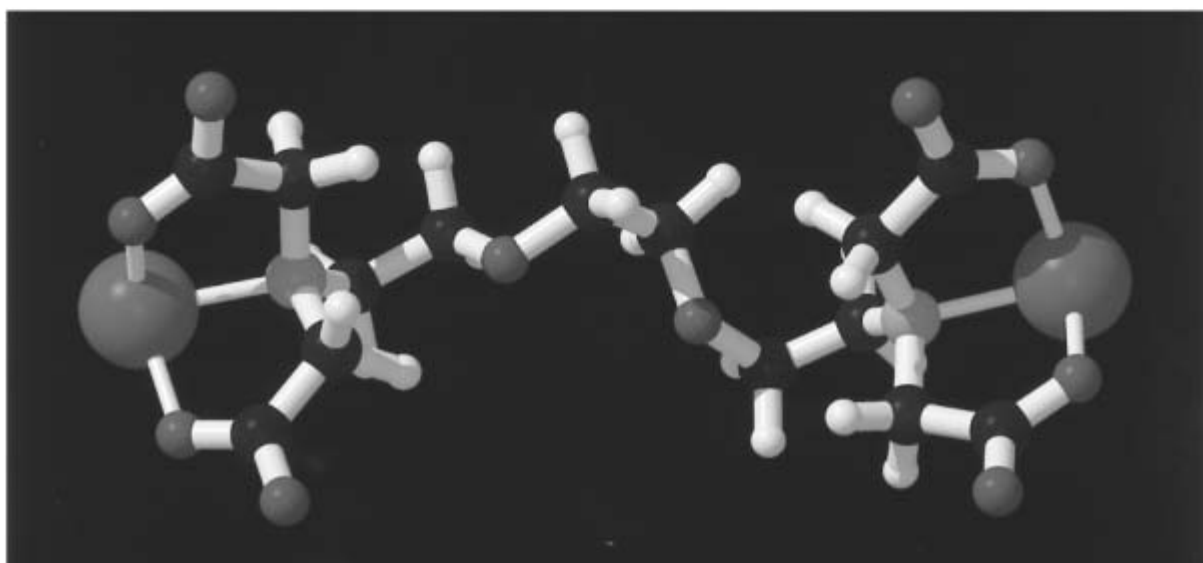
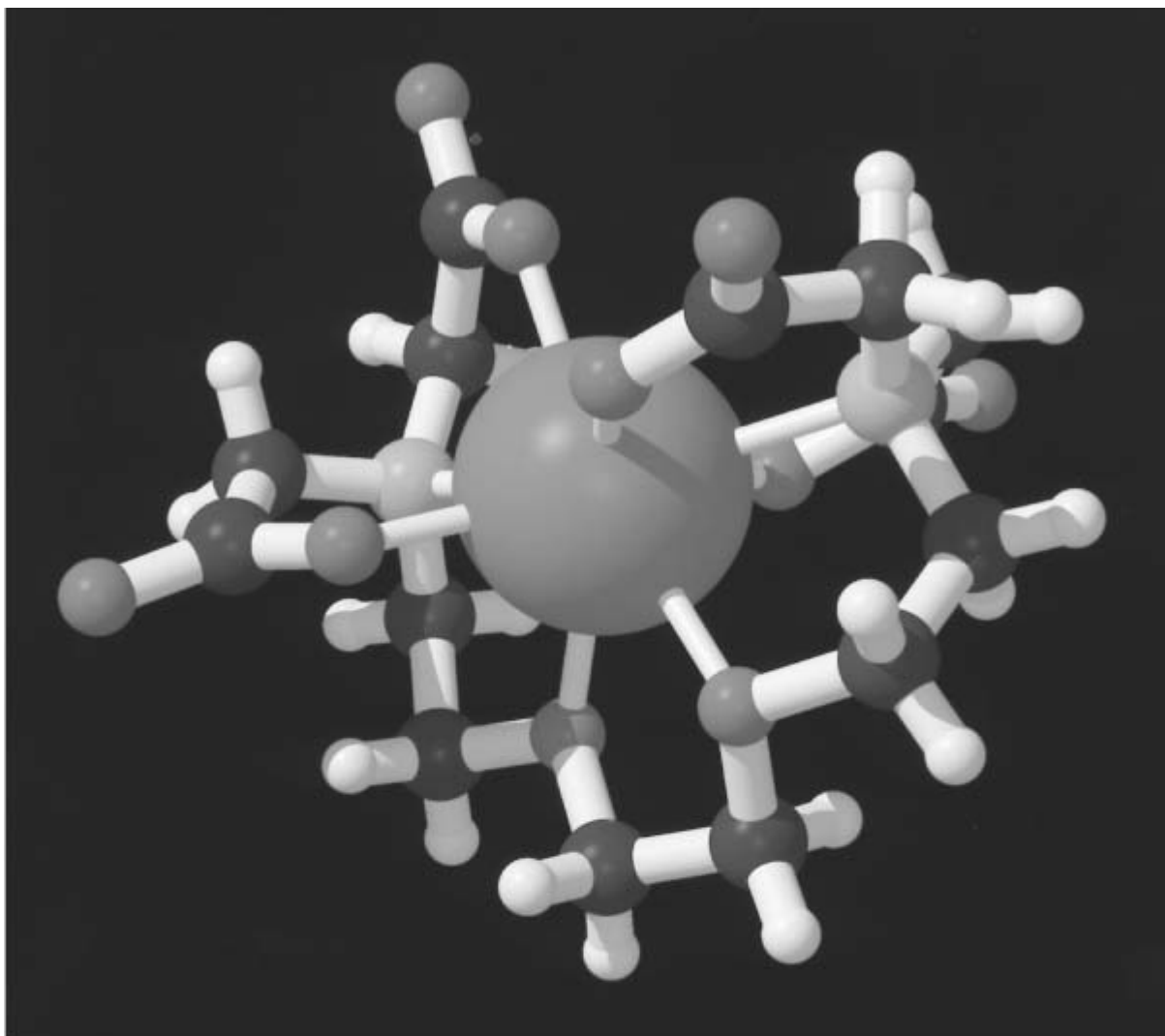
EGTA is a chelating agent widely used to control the concentration of  $\text{Ca}^{2+}$  cations in biological solutions (see also [-Tetraacetic Acid](#)). Its structure is given in Figure 1. EGTA binds  $\text{Ca}^{2+}$   $3 \times 10^5$ -fold more strongly than  $\text{Mg}^{2+}$ , which makes it useful to control  $\text{Ca}^{2+}$  levels in the presence of physiological concentrations of  $\text{Mg}^{2+}$  (1-3). The binding of cations by EGTA is pH-dependent; at low pH the carboxylate moieties will be protonated, dramatically lowering the cation affinity. The high selectivity of EGTA for  $\text{Ca}^{2+}$  over  $\text{Mg}^{2+}$  is readily explained by the structure of their respective complexes (Fig. 2, see top of next page) (4, 5). While EGTA fully wraps around  $\text{Ca}^{2+}$  in a very stable octadentate complex, it cannot similarly satisfy the more stringent binding requirements of the smaller  $\text{Mg}^{2+}$ . The preference of  $\text{Mg}^{2+}$  for oxygen ligands over nitrogen ligands is so great that it binds to EGTA via the carboxylate groups only. The remaining positions in the ligand shell are filled by water. Although these complexes can be crystallized and their structure determined, they are very dynamic in solution. NMR analysis of the  $\text{Ca}^{2+}$ -EGTA complex indicates that the carboxylate moieties interchange positions rapidly on the NMR time scale, and inversion at the nitrogen atoms is observed at elevated temperatures (5), which further indicates that the complex is very dynamic. The great stability of the  $\text{Ca}^{2+}$ -EGTA complex comes from the fact that all ligands are provided by a single molecule, specifically, the stability constants have a very high entropic contribution.

Figure 1. Molecular structure of EGTA.





**Figure 2.** Three-dimensional structures of complexes of EGTA with (a)  $\text{Ca}^{2+}$  and (b)  $\text{Mg}^{2+}$ .



EGTA as the free acid is a white crystalline powder and has a molecular weight of 340.35 Da. The affinities of EDTA and EGTA for several cations are compared in Table 1(6).

**Table 1. Comparison of Affinities of EDTA and EGTA for Various Cations<sup>a</sup>**

| Metal Ion        | $\log_{10} K_{m\text{-EDTA}}$ | $\log_{10} K_{m\text{-EGTA}}$ |
|------------------|-------------------------------|-------------------------------|
| Mg <sup>2+</sup> | 8.64                          | 5.3                           |
| Ca <sup>2+</sup> | 11.0                          | 10.9                          |
| Zn <sup>2+</sup> | 16.4                          | 12.6                          |
| Cd <sup>2+</sup> | 16.4                          | 16.5                          |

<sup>a</sup> The logarithms of their association constants ( $K_M$  in units of  $M^{-1}$ ) are compared at pH 7.0, 25°C, and 0.1 M ionic strength (6)

## Bibliography

1. A. Fabiato and F. Fabiato (1979) *J. Physiol. (Paris)* **75**, 463.
2. A. Fabiato (1991) In *Cellular Calcium: A Practical Approach*, J. G. McCormack and P. H. Cobbold, eds., IRL Press, Oxford, pp. 159–176.
3. M. Otto, P. M. May, K. Murray, and J. D. Thomas (1985) *Anal. Chem.* **57**, 1511–1517.
4. C. K. Schauer and O. P. Anderson (1987) *J. Am. Chem. Soc.* **109**, 3646–3656.
5. C. K. Schauer and O. P. Anderson (1988) *Inorg. Chem.* **27**, 3118–3130.
6. A. E. Martell and R. M. Smith (1974) *Critical Stability Constants*, Vol. **1**, Plenum Press, New York.

## Ehrlich Cells

### 1. Origin

Ehrlich cells, more properly known as Ehrlich–Lettre ascites carcinoma, were originally established as an ascites tumor of mice (1) from Ehrlich's “Strain 7” transplantable mouse carcinoma (2), which was probably of mammary origin. Several successful attempts at tissue culture have been performed, the earliest being hanging drop preparations that were passaged for 12 years (3). Subsequently, several other laboratories have cultured Ehrlich's cells (4-9), but the cell line that has been used most frequently, and which is lodged in the American Type Culture Collection (ATCC) repository as ATCC CCL-77, is the Strain E of Ehrlich–Lettre ascites (10). These cells were established as a monolayer from a 7-day-old tumor in NCTC 109 culture medium with 10% calf serum, and subcultured by scraping. The culture lodged with ATCC has been maintained in a later modification of NCTC 109, NCTC 135 supplemented with fetal calf serum.

## 2. Properties

Strain E has a short population doubling time *in vitro*: 20–24 h in Strain E and 14–16 h in Strain E (4). Strain E(4) was derived from alternate *in vitro* and *in vivo* passage, adapts rapidly to culture following *in vivo* passage, and is, presumably, the strain that was submitted to ATCC. The cells tend to have a fusiform morphology *in vitro*, suggesting an undifferentiated **phenotype**, and have a hyperdiploid **karyotype** (modal number 44,  $2C = 40$ ) with distinctive A chromosome, metacentric, and minute markers.

## 3. Usage

The Ehrlich ascites became popular because it could be expanded manyfold by *in vivo* passage, making it useful for biochemical studies involving large amounts of tissue, while it could still be maintained *in vitro* for more physiologically controlled studies. With the advent of large-scale cell culture techniques that can yield  $10^9$ – $10^{12}$  cells, ascites passage is less attractive, due to the contamination of the tumor with a variety of host inflammatory cells and the development of legislation limiting the use of ascites tumors.

## Bibliography

1. H. Lowenthal and G. Jahn (1932) *Z. Krebsforsch.* **37**, 439–447.
2. P. Ehrlich and H. Apolant (1905) *Berl. Klin. Wehr.* **42**, 871–874.
3. A. Fisher and F. Davidsohn (1939) *Nature (Lond.)* **143**, 436–437.
4. G. E. Foley, B. P. Drolet, R. E. McCarthy, K. A. Goulet, J. M. Dokos, and D. A. Filler (1960) *Cancer Res.* **20**, 930–939.
5. M. M. Guerin and J. F. Morgan (1961) *Cancer Res.* **21**, 378–382.
6. J. O. Ely and J. H. Gray (1960) *Cancer Res.* **20**, 918–922.
7. R. Cailleau and F. Costa (1961) *J. Natl. Cancer Inst.* **26**, 271–282.
8. P. W. Jackson, N. Giuffre, and D. Perlman (1960) *Can. J. Biochem. Physiol.* **38**, 1377–1378.
9. J. A. DiPaolo (1962) *Proc. Soc. Exp. Biol. Med.* **109**, 616–618.
10. C. Boone, M. Sasaki, and R. W. McKee (1965) Characterization of an *in vitro* strain of Ehrlich-Lette ascites carcinoma subjected to many periodic mouse passages. *J. Natl. Cancer Inst.* **34**, 725–730.

## Elastase

This is the designation for **enzymes** that degrade **elastin**, the elastic protein found in tissues, such as that of the lung and the aorta, that undergo repeated stretching. Elastase (E.C. 3.4.21.36) is a 25-kDa protein produced in the pancreas, in the form of an inactive precursor that is activated on secretion into the small intestine, and it helps digest dietary proteins. A different elastase found in neutrophil granulocytes (E.C. 3.4.21.37), and one of its functions may be to degrade bacterial cell walls during **phagocytosis** (1). The elastases are **serine proteinases** of the **trypsin**/chymotrypsin family, with specificity for peptide bonds involving amino acids with small side chains eg, **alanine** and **valine**).

Leukocyte elastase, a 30-kDa **glycoprotein**, degrades **extracellular matrix** proteins in addition to elastin, and in this regard it probably plays a role in neutrophil migration and perhaps tissue remodeling (2). It facilitates the activation of **plasminogen**, a protein involved in the dissolution of

**blood clots** (3), and it may play an important role in tumor metastasis (4). It also has an important pathological action in pulmonary emphysema (5). Under normal conditions, its activity is controlled by **a1-Antitrypsin**, a member of the **serpin** family of **proteinase inhibitors**. Some individuals have an inherited deficiency of this inhibitor, which leads to a proteinase/antiproteinase imbalance in lung tissue, tissue degradation, and eventually emphysema. Leukocyte elastase may also contribute to rheumatoid arthritis and inflammation.

An elastase from human alveolar macrophages has been identified as a 92-kDa enzyme that is also a *gelatinase*, that is, it degrades denatured **collagen** (6). This enzyme is not a serine proteinase but a zinc **metalloenzyme** and a member of the matrix metalloproteinase family. It is not yet clear whether this enzyme has a significant role in the pathogenesis of emphysema.

### Bibliography

1. A. Janoff et al. (1975) In *Proteases and Biological Control* (E. Reich, D. B. Rifkin, and E. Shaw, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 603–620.
2. J. G. Bieth (1988) *Pathol. Biol.* **36**, 1108–1111.
3. R. Machovich and W. G. Owen (1990) *Blood Coag. Fibrinol.* **1**, 79–90.
4. M. S. O'Reilly et al. (1994) *Cell* **79**, 315–328.
5. W. van Steenberg (1993) *Acta Clin. Belg.* **48**, 171–189.
6. S. D. Shapiro (1994) *Am. J. Respir. Crit. Care Med.* **150**(6 Pt 2), S160–S164.

### Elastin

Elastin is the major macromolecule found in the ECM of elastic tissues such as the lung, skin, and blood vessels. The molecular properties of elastin are responsible for the elasticity. Elastin is a hydrophobic protein of about 800 amino acids that assumes a random coil configuration. Networks of elastin near the cell surface are covalently cross-linked via lysine residues, and elasticity is thought to be achieved by stretching the random coils. The protein has an abundance of proline and glycine, similar to collagen, but unlike collagen, elastin has only a small amount of hydroxyproline and lacks both hydroxylysine and carbohydrate side chains. The elastic layer is strengthened by an additional layer of microfibrils made of the protein fibrillin.

Cells may interact with elastin via an unusual 67-kDa elastin receptor. The receptor protein, unlike matrix receptors in the integrin class, is a secreted rather than transmembrane protein. It appears to mediate cell adhesion to elastin via interactions with cell surface proteins. The same receptor also binds laminin, and harbors a lectin like activity for binding sugar moieties. This protein is also found inside the cell where it may act as a chaperone to facilitate tropoelastin secretion and assembly of elastic fibers.

The functional importance of elastin is demonstrated poignantly by the recent discovery that deletions in the human elastin gene give rise to a disease known as Williams syndrome. People affected with this disease are mentally retarded (although occasionally mathematically or musically gifted), have an elfin like facial appearance and are prone to blood vessel narrowing and aneurysms.

Another disease, supravalvular aortic stenosis (SVAS), results from haploinsufficiency of elastin, and is also characterized by luminal narrowing in blood vessels.

## Additional Reading

Ramirez F., Pathophysiology of the microfibril/elastic fiber system: introduction, *Matrix Biol.* **19** (6), 455–456 (2000), and other reviews in the same volume.

Debelle L and Tamburro A.M., Elastin: molecular description and function, *Int. J. Biochem Cell Biol.* **31**(2), 261–72 (1999).

## Electroelution

Molecules and particles purified by [gel electrophoresis](#) can be recovered from the gel by their [electrophoresis](#) into solution, a process termed electroelution. The procedure consists of excising the gel slices of interest, loading them onto a concentrating (stacking) gel, and electrophoresing the material contained in the slice into a recovery cup or dialysis bag attached to the bottom of the stacking gel. If the electroelution employs a single electrophoresis buffer, the sample moves directly into the recovery vessel. Electroelution of the band of interest into the stacking gel of a discontinuous buffer system of [disc electrophoresis](#) has the advantage that the material to be eluted is concentrated into a single band visualized by a [tracking dye](#), so its collection can be visually monitored. This type of electroelution is practical with microgram amounts of sample, which can be recovered with a typical yield of 70%. Polymeric impurities in the electroeluate may be removed by gel filtration (1).

More directly, electroelution can be conducted at a right angle to the direction of migration, with instant monitoring of recovery until it is complete (2). Sequential elution of bands into a buffer chamber flushed continuously or discontinuously (designated “elution-PAGE”) is available for samples of up to milligrams, but it involves dilution of the eluate and the risk of losses in the elution chamber (3).

## Bibliography

1. N. Y. Nguyen and A. Chrambach (1981) In *Gel Electrophoresis of Proteins: A Practical Approach* (B. D. Hames and D. Rickwood, eds.), IRL Press, Oxford, pp. 145–155.
2. E. Gombocz and E. Cortez (1995) *Appl. Theor. Electrophoresis* **4**, 197–209.
3. A. Chrambach and D. Rodbard (1981) In *Gel Electrophoresis of Proteins: A Practical Approach* (B. D. Hames and D. Rickwood, eds.), IRL Press, Oxford, pp. 93–143.

## Electroendosmosis

Electroendosmosis (EEO) is the flow of solvent in the electric field during [electrophoresis](#). EEO is caused by stationary charged groups in the gel (see [Agarose](#)) or along the inner wall of the vessel (see [Capillary Zone Electrophoresis](#)); negatively charged acidic groups in the gel or silicate groups embedded in glass surfaces bind mobile positively charged counterions electrostatically. The positively charged counterions migrate in the electric field, transporting their [water of hydration](#) and thereby giving rise to the water flow of EEO in that direction. Since the flow rate decreases with

increasing distance from the glass wall, EEO is responsible for the parabolic profile of the electrophoretic band when a gel is replaced by a polymer solution (see [Particle Electrophoresis](#)). EEO can be used to enhance resolution, but more frequently it is suppressed by covering the stationary charged groupings with a neutral polymer (see [Capillary Zone Electrophoresis](#)).

EEO can be measured as the displacement in the electric field of uncharged compounds, such as cyanocobalamin in the neutral pH range or mesityl oxide. EEO decreases semilogarithmically with gel concentration, in complete analogy to electrophoretic mobility (see [Ferguson Plot](#)). Quantitative measurements of electrophoretic mobilities need to be corrected for the EEO.

#### Suggestions for Further Reading

H. A. Abramson, L. S. Moyer, and M. H. Gorin (1964) *Electrophoresis of Proteins and the Chemistry of Cell Surfaces*, Hafner, New York, pp. 106–108.

S. Ghosh and D. B. Moss (1974) Electroendosmosis correction for electrophoretic mobility determined in gels. *Anal. Biochem.* **62**, 365–370.

## Electron Crystallography

Knowledge on the molecular basis of many biological processes has rapidly evolved as a result of the development of complementary techniques of molecular biology and high-resolution structural biology. Quite often the mode of action of biological [macromolecules](#) can only be understood at the atomic level. Thousands of [protein structures](#) have been solved to atomic scale resolution using X-ray crystallography, but only a handful by electron crystallography. However, [electron microscopy](#) is a technique intrinsically capable of atomic resolution. The transmission electron microscope itself can achieve resolution on the order of 0.2–0.3 nm, and the structures of small organic and inorganic molecules are routinely determined at atomic dimensions. It is now well established that the principal limitation to obtaining atomic resolution of biological specimens using electron microscopy is the specimen itself (1). Because biological molecules are sensitive to electron radiation, one must record electron micrographs of unstained specimens with very low levels of electron irradiation (this is called low-dose imaging) in order to preserve high-resolution information. Unfortunately, such low electron doses yield images with a very low signal-to-noise ratio, too low to observe any structural detail. High signal-to-noise ratio images can be obtained by averaging hundreds or thousands of low-dose images of identical structures (2). Two-dimensional crystals, thin crystals which are only one unit cell thick along one axis, currently offer the best way to average images of many molecules to build up the signal-to-noise at atomic or near-atomic resolution (3, 4). It was estimated that the average of ~5 million molecules was required to produce a 3-D density map of [bacteriorhodopsin](#) that had recognizable amino acid side chains (1, 5). Similarly, ~2.6 million molecules of light-harvesting complex II were averaged to obtain an equivalent resolution (6).

Electron crystallography encompasses the use of electron diffraction (7) and imaging of thin crystals suitable for the electron microscope. All electron images are 2-D projections of a 3-D object along the path of the electron beam. Hence, to reconstruct an object in 3-D it is necessary to collect images from different angles. The minimal number of tilted views needed to reconstruct an object at a given resolution can be calculated from geometrical considerations (8).

Once a 2-D or thin 3-D crystal is grown, the best way to evaluate its quality is to study its diffraction pattern. The quickest and easiest method is to examine electron micrographs by optical diffraction. In this technique, an electron micrograph is illuminated with laser light focused with a series of lenses

to form an optical diffraction pattern showing the distribution of Fourier transform intensities in the same manner as an X-ray diffraction pattern of a 3-D crystal. The spacing and distribution of intensities in an optical diffraction pattern offers clues to the molecular packing and secondary structure (9). Optical diffraction is also an important step in judging image quality in terms of astigmatism, defocus, and drift. More information can be obtained from the computed diffraction produced by calculating the Fourier transform of the digitized image. Unlike X-ray diffraction patterns, the computed diffraction pattern will have both structure factor amplitudes and phases in the “weak phase object approximation” (thin object and sufficiently high electron energy). Hence, a reconstruction of the object can be calculated from the Fourier transform of the image using the structure factor amplitudes and phases of the crystal; this results in a reconstructed image of the average unit cell of the crystal with noise removed. Furthermore, the phase relationship of the symmetry-related spots in the computed diffraction pattern can be used to determine the crystallographic symmetry.

Diffraction patterns, similar in principle to those obtained in X-ray diffraction, are formed in the back focal plane of the objective lens in the electron microscope. By adjusting the current in the diffraction lens, the diffraction pattern can be observed on the viewing screen. The type of pattern observed depends on the arrangement of molecules in the crystal and on how well ordered they are. How the molecules are arranged within the repeating unit of a crystal (unit cell) can be determined by the symmetry of the diffraction pattern and the intensities of the diffraction spots. In addition to the molecular arrangement, diffraction patterns can reveal the degree of specimen preservation and show the distribution of ordered regions of the specimen. The further the spots extend, the higher is the resolution of the preserved detail. Diffraction patterns can be used to assess radiation damage to the specimen with increasing electron dose, as judged by fading of the higher-order diffraction spots. The size of crystals most amenable for electron diffraction of biological macromolecules is a few microns on edge and one unit cell thick, usually described as a 2-D crystal (10). However, the electron diffraction patterns of 3-D crystals can also be examined if they are less than a few hundred Angstroms thick (11); Perkins, unpublished results). Electron diffraction can also ascertain the effects of different methods of specimen preparation. If a protein has secondary structural elements aligned in parallel, as in fibers, the spacings of the diffraction spots or streaks can indicate if there is predominantly  $\alpha$ -helix or  $\beta$ -sheet (12). One disadvantage of electron diffraction is that relatively large crystals ( $\sim 100 \times 100$  unit cells) are needed before high-resolution spots are observed. Another disadvantage is that like X-ray diffraction, electron diffraction does not produce structure factor phases; hence the structure cannot be determined from electron diffraction alone (13).

For high-resolution structural studies, advances in imaging, computer image processing, and preservation of specimens without negative stain and fixative have been critically important. The native conformation of specimens must be preserved under high vacuum conditions, and this can be achieved using one of two types of embedding media: sugars (14-16) and vitreous ice (17), which have both been shown to preserve the protein crystal structure well enough to produce images at a resolution of 0.3 nm. In addition, the development of computer-controlled spot-scan imaging has improved the efficiency of transfer of high resolution image information (18). Improvement in the rate of data acquisition has been made using cameras employing a cooled charge-coupled device (19). Because the high-resolution signal-to-noise is low, averaging over several large (20) or many more small (4) images is necessary to restore the signal above the background. Despite several noteworthy breakthroughs, the rate of progress in determining high-resolution structures of biological macromolecules by electron microscopy has been slowed by the paucity of suitable, well-ordered crystals. Recent reviews have described crystallization methods for both water-soluble and [membrane proteins](#) (21-23). Perhaps further improvements in collecting high-resolution information from single particles will obviate the need for crystals (1).

Electron crystallography offers a viable alternative to [X-ray crystallography](#) when crystals are not large enough for X-ray diffraction or for membrane proteins (24). Integral membrane proteins have the advantage that they can be routinely reconstituted in their native environment, the [lipid](#) bilayer, and often can be forced into ordered arrays through the judicious choice of lipid/protein and

detergent/protein concentrations. The first membrane protein structure revealed by electron crystallography was bacteriorhodopsin in which seven [a-helices](#) were resolved (25). Since this early work, the field of electron crystallography has advanced to the point that atomic models have been built for two membrane proteins (5, 6, 26) and good progress made toward this goal for several others (27-29). [Tubulin](#) crystals also diffract to high resolution, and electron crystallography has shown the location of the taxol-binding site (30, 31). The highest resolution structures have been generated by taking the phases from the images and the amplitudes from the electron diffraction patterns. The electron diffraction amplitudes are generally more accurate than the image amplitudes because the image amplitudes require correction for instrumental factors that affect image quality, and methods for this correction have not been fully developed (20). However, it is not always necessary for electron microscopy to be performed at high resolution to contribute to atomic resolution structures. Electron microscopy provided a critical link allowing the [actin](#) crystal structure determined by X-ray diffraction to be fit into the actin filament determined by electron microscopy (32, 33).

The structure of biological macromolecules with helical or icosahedral symmetries can be determined using the same underlying principles of electron crystallography, although the operational approaches will be different (34). Helical structures produce diffraction patterns in which the helical information is confined to regularly spaced layerlines perpendicular to the helix axis instead of discrete spots. The pattern of intensities along the layerlines usually produce X-shaped patterns. Helical symmetry provides different views of a molecule in a single image as identical molecules are arranged along a helical path, so in principle only one image is needed for a low-resolution 3-D reconstruction. As with electron crystallography, however, higher resolution requires many images in order to obtain views at many different angles and to average many images in order to improve the signal-to-noise ratio (35). The structures of helical assemblies have been solved to 0.9 nm resolution (35, 36), and higher resolution is shortly anticipated. Close to 0.7 nm resolution has been reported for a viral structure with icosahedral symmetry (37). Even with the benefits of high symmetry elements afforded by an icosahedron, more than 6000 images needed to be averaged to achieve this resolution.

## Bibliography

1. R. Henderson (1995) *Quarterly Rev. Biophys.* **28**, 171–193.
2. R. M. Glaeser and K. H. Downing (1993) *Ultramicroscopy* **52**, 478–486.
3. D. Brillinger et al. (1989) *J. Applied Stat.* **16**, 165–175.
4. G. A. Perkins, K. H. Downing, and R. M. Glaeser (1995) *Ultramicroscopy* **60**, 283–294.
5. R. Henderson et al. (1990) *J. Mol. Biol.* **213**, 899–929.
6. W. Kuhlbrandt, D. N. Wang, and Y. Fujiyoshi (1994) *Nature* **367**, 614–621.
7. D. L. Misell and E. B. Brown (1987) *Electron Diffraction: An Introduction for Biologists*, Elsevier, Amsterdam, The Netherlands.
8. A. Klug (1979) *Chem. Scr.* **14**, 245–256.
9. F. Thon (1971) In *Electron Microscopy in Material Science* (U. Valdre, ed.), Academic Press, New York, pp. 571–625.
10. T. A. Ceska and R. Henderson (1990) *J. Mol. Biol.* **213**, 539–560.
11. T. W. Jeng et al. (1984) *J. Mol. Biol.* **175**, 93–97.
12. R. D. B. Fraser and T. P. MacRae (1973) *Conformation in Fibrous Proteins*, Academic Press, New York.
13. J. M. Baldwin and R. Henderson (1984) *Ultramicroscopy* **14**, 319–335.
14. P. N. T. Unwin and R. Henderson (1975) *J. Mol. Biol.* **94**, 425–440.
15. W. Kuhlbrandt (1988) *J. Mol. Biol.* **202**, 849–864.
16. G. Perkins et al. (1993) *J. Microscopy* **169**, 61–65.
17. J. Dubochet et al. (1988) *Quarterly Rev. Biophys.* **21**, 129–228.



18. K. H. Downing (1991) *Science* **251**, 53–59.
19. A. J. Koster et al. (1992) *Ultramicroscopy* **46**, 207–228.
20. J. M. Baldwin et al. (1988) *J. Mol. Biol.* **202**, 585–591.
21. R. Kornberg and S. A. Darst (1991) *Curr. Opin. Struct. Biol.* **1**, 642–646.
22. B. K. Jap et al. (1992) *Ultramicroscopy* **46**, 45–84.
23. W. Kuhlbrandt (1992) *Quarterly Rev. Biophys.* **25**, 1–49.
24. A. Engel et al. (1992) *J. Struct. Biol.* **109**, 219–234.
25. R. Henderson and P. N. T. Unwin (1975) *Nature* **257**, 28–32.
26. N. Grigorieff et al. (1996) *J. Mol. Biol.* **259**, 393–421.
27. A. Olofsson, V. Mallouh, and A. Brisson (1994) *J. Struct. Biol.* **113**, 199–205.
28. H. Hebert, I. Schmidt-Krey, and R. Morgenstern (1995) *EMBO J.* **14**, 3864–3869.
29. T. Walz et al. (1995) *Nature Struct. Biol.* **2**, 730–732.
30. E. Nogales et al. (1995) *Nature* **375**, 424–427.
31. S. G. Wolf et al. (1996) *J. Mol. Biol.* **262**, 485–501.
32. W. Kabsch et al. (1990) *Nature* **347**, 37–44.
33. K. C. Holmes et al. (1992) In *Mechanism of Myofilament Sliding in Muscle Contraction* (H. Sugi and G. H. Pollack, eds.), Plenum Press, New York.
34. D. G. Morgan and D. De Rosier (1992) *Ultramicroscopy* **46**, 263–286.
35. T. W. Jeng et al. (1989) *J. Mol. Biol.* **205**, 251–257.
36. D. G. Morgan et al. (1995) *J. Mol. Biol.* **249**, 88–110.
37. B. Bottcher, S. A. Wynne, and R. A. Crowther (1997) *Nature* **386**, 88–91.

## Electron Imaging

Particle/wave dualism is evident in experiments with electrons. They form a beam of particles, for instance, in a synchrotron, but behave as waves in an **electron microscope**. Electrons are scattered by matter just as X-rays are (see [X-Ray Crystallography](#)). Electrons are scattered by the atomic coulombic potential distribution in the atoms. This scattering is much stronger than for X-rays and can be recorded from very tiny specimens. This is particularly valuable for obtaining information on individual regions extending over only a few [unit cells](#) of a two-dimensional crystal (see [Crystallography](#)). The imaging experiments are performed in an electron microscope at a wavelength between 0.02 and 0.04 Å, much shorter than for X-rays. Then, the Ewald sphere (see [Reciprocal Space](#)) can be regarded as a flat surface, and diffraction from a stationary crystal shows essentially the reflections in a planar section of reciprocal space. By tilting the specimen, several intersections can be obtained and combined in a more complete image of reciprocal space. The short wavelength would result in a very high image resolution if the quality of the microscope lenses were not a limiting factor. At best they allow reaching a resolution of 1.5 Å. Simple switching of lens currents changes the observation from a diffraction pattern to a real image.

Transfer of energy to the specimen and the resulting radiation damage is a limitation in applying electron scattering. This is especially serious in studying biological material. However, with very thin specimens of two-dimensional crystals and with the molecules embedded in vitreous ice or glucose, extremely interesting results have been obtained from biological specimens by short exposure times, low beam intensity, cryocooling, and a combination of diffraction data and image

analysis. The phases of the [structure factors](#) of the diffracted beams are calculated from the image (see [Phase Problem](#)). They are combined with the measured amplitudes of the diffracted beams and, exactly as in X-ray diffraction, a Fourier summation gives the atomic distribution in the specimen. At the present time, the maximum resolution is 3 Å.

### Suggestions for Further Reading

R. Henderson (1995) The potential and limitations of neutrons, electrons and X-rays for atomic resolution microscopy of unstained biological molecules, *Q. Rev. Biophys* **28**, 171–193.

International Union of Crystallography (1996) *International Tables for Crystallography*, Vol. B (U. Smueli, ed.) Kluwer Academic, Dordrecht, Boston, London, pp. 280–329.

N. Grigorieff et al. (1996) Electron-crystallographic refinement of the structure of bacteriorhodopsin, *J. Mol. Biol.* **259**, 393–421.

## Electron Paramagnetic Resonance

### 1. Introduction to Electron Paramagnetic Resonance (EPR) Spectroscopy

Electron paramagnetic resonance (EPR) refers to [spectroscopy](#) of unpaired electrons, and some aspects of the technique are termed “electron spin resonance” (ESR) or “electron magnetic resonance” (EMR). Unpaired electrons in biological systems are in much lower abundance than nuclei; thus, EPR is a technique that focuses on local sites, whereas nuclear magnetic resonance (NMR) tends to be more global. A connection exists between EPR and NMR in the variations of EPR that are termed “electron spin echo envelope modulation” (ESEEM) and “electron nuclear double resonance” (ENDOR). In these techniques, the energies of the transitions involving electrons are modulated by nearby nuclei, and this effect provides information about the type and geometry of nuclei around the site of an unpaired electron. The term “spin” in the designations of various forms of magnetic resonance spectroscopy refers to a property of electrons (or nuclei), and it is the interaction of the spin with the magnetic field that leads to separation of energy levels between which the spectroscopic transition occurs. A magnetic field is usually required for EPR, as is a source of energy to effect transitions. In addition, a few cases exist where spin transitions can be detected in the absence of a magnetic field.

**Electron transfer** systems (for instance, those in **photosynthetic** membranes) provide excellent opportunities to follow the intimate details of the transfer process with EPR spectroscopy. Other naturally occurring sources of unpaired electrons that are subjects for this form of spectroscopy include complexes of a variety of molecules with [nitric oxide](#), quinone and flavin **cofactors** in [enzymes](#), free radical enzyme intermediates, and metal ion sites in proteins (see [Metalloproteins](#)). By definition, a free radical has a homolytically broken bond with an unpaired electron. The range of applications of EPR spectroscopy is not limited to these natural sources of unpaired electrons. The EPR probe technique of [spin labeling](#) is being applied to examine dynamics of **DNA** and proteins, to study **protein folding**, and to unravel the complex motions of [membrane](#) components.

Molecular biological approaches extend the applications of spin labeling, and these applications are generically referred to as site directed spin labeling. In addition, unpaired electrons in free radical intermediates that have a short lifetime can be studied by trapping them to give a longer-lived, secondary radical in an EPR-related technique called “spin trapping.” The text below will give examples of the type of information that can be obtained by applying EPR spectroscopy to these instances of unpaired electrons in biology.

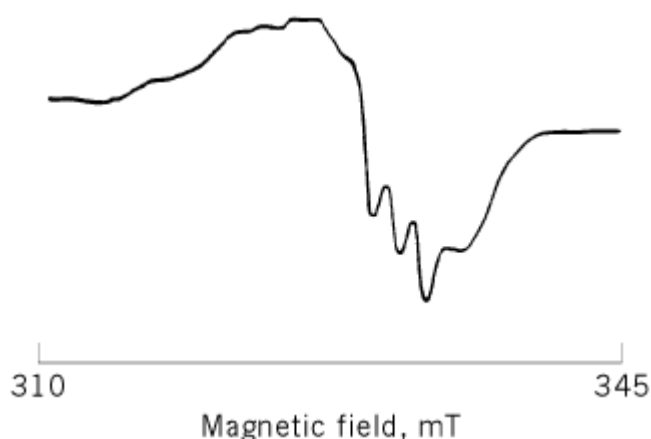
## 2. EPR in Photosynthesis

One of the long-standing puzzles of electron transfer in photosynthetic membranes is that the photoreaction center has pairs of molecular components arranged apparently symmetrically about a plane perpendicular to the membrane surface, yet electron transfer among the components follows a path involving primarily one-half of the dimer (see [Photosynthesis](#) and [Photosynthetic Reaction Center](#)) (1). EPR studies, especially of oriented crystals of reaction centers (2), provide precise information about small differences in symmetry of the two halves of the system. From EPR studies, it is possible to determine the orientation of orbitals involved in electron transfer between cofactors, and the positions of light atoms, such as hydrogen, on the cofactors can be determined by ENDOR with great precision (1). The various components of the electron transfer pathway—chlorophylls and quinones—can be selected for EPR by optical experiments, selective depletion of some components, and molecular biology. High-frequency EPR (1) helps to separate the spectra of different sites of unpaired electrons, much as higher frequency NMR gives better resolution of different nuclei. The manner in which electrons flow ultimately into the photosynthetic oxygen-evolving complex, which contains varied numbers of unpaired electrons associated with a cluster of four manganese atoms, is another challenge. ESEEM and ENDOR studies, together with isotope labeling (3) provide details of the structure of this site.

## 3. EPR of Nitric Oxide Complexes

Molecular adducts of nitric oxide give a diverse set of EPR signals that can be employed to determine the pathway of nitric oxide in tissues. Ferrous iron may or may not have unpaired spins, depending on whether it is low or high spin, but, in either case, it is, respectively, an impossible or difficult subject for EPR spectroscopy. Fortunately, the adduct of nitric oxide with ferrous iron in heme is quite suitable for EPR spectroscopy and provides a characteristic signature, shown in Figure 1. This spectrum is distinct from those of the nitric oxide adduct of either **amino-** or **thiol-** groups of protein side chains, as well as from copper adducts and adducts with other forms of iron. For further delineation of the location of nitric oxide in tissues, “spin traps” for nitric oxide are available (see section below on spin trapping). Typical spin traps for nitric oxide are ferrous ion chelated with *N*-methyl-D-glucamine dithiocarbamate or with diethylthiocarbamate.

**Figure 1.** The EPR spectrum of deoxymyoglobin to which nitric oxide (NO) gas was added. The spectrum was recorded at 77K and a frequency of 9.1 GHz. (Reproduced with permission from Ref. 4, Fig. 8.)

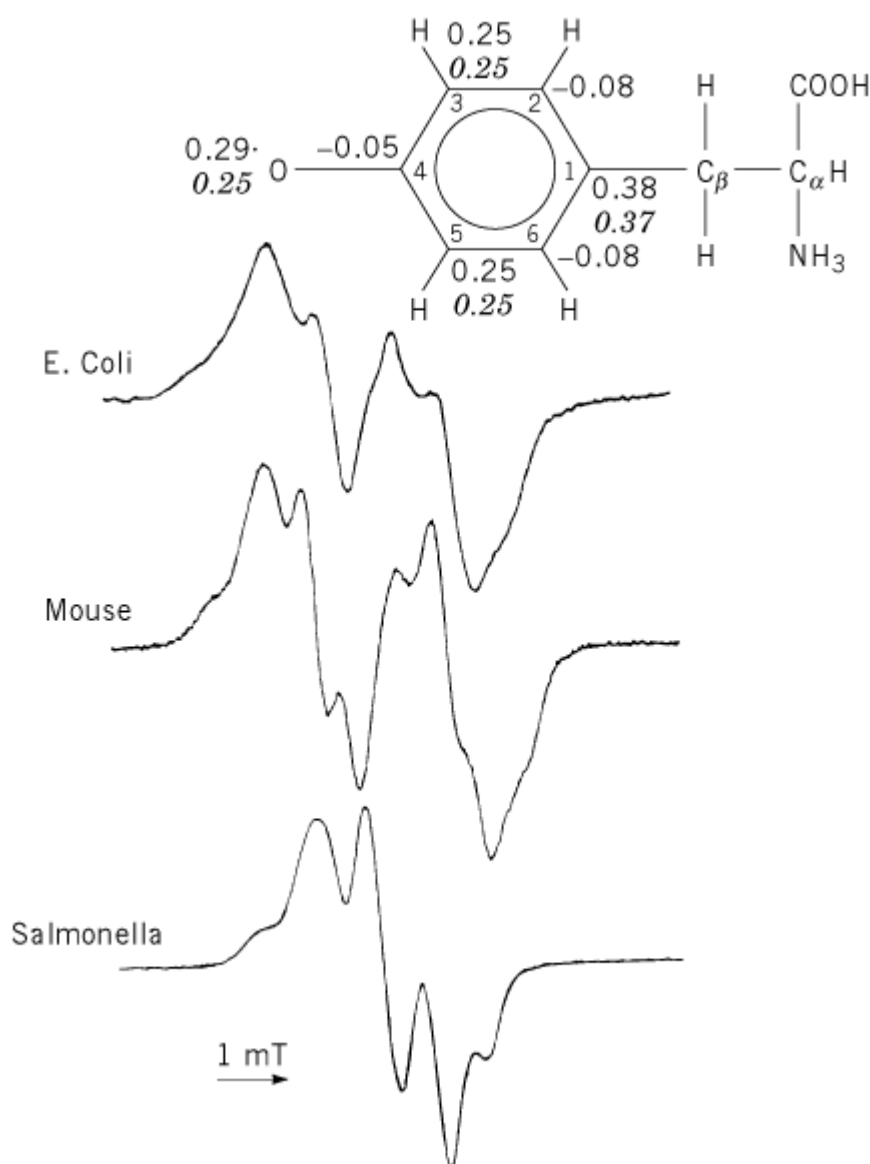


## 4. EPR of Free-Radical Protein Intermediates

[Ribonucleotide reductases](#) use metals and free radical chemistry to carry out reduction of

ribonucleotides to deoxyribonucleotides. (56) In different sources of the enzyme, the presence of free-radical intermediates has been demonstrated to involve side chains of the amino acids [tyrosine](#), [tryptophan](#), [cysteine](#) and [glycine](#). EPR spectroscopy of tyrosine radicals has been used in this system, as well as in photosynthetic membranes, to demonstrate that the unpaired electron density at the various carbon atoms of the tyrosine aromatic ring is fine tuned by [hydrogen bonds](#) to the phenolic oxygen and by the protein environment. Figure 2 gives examples of EPR of tyrosine radicals from ribonucleotide reductases of several sources. The process by which involvement of metal ions and radical side chains lead to a free radical intermediate of the ribonucleotide is thought to involve long-range electron transfer. Organic cofactors, such as flavins, also form intermediates that can be detected in EPR studies of enzyme reactions.

**Figure 2.** EPR spectra of tyrosine freeradicals of ribonucleotide reductase R2 subunits from different organisms. The EPR spectra were recorded at temperatures of 20 to 30K at ~9 GHz. The numbers given on the structure are the spin density distributions for the *E. coli* and *S. typhimurium* enzymes. (Reproduced with permission from Ref. 6, Fig. 3).

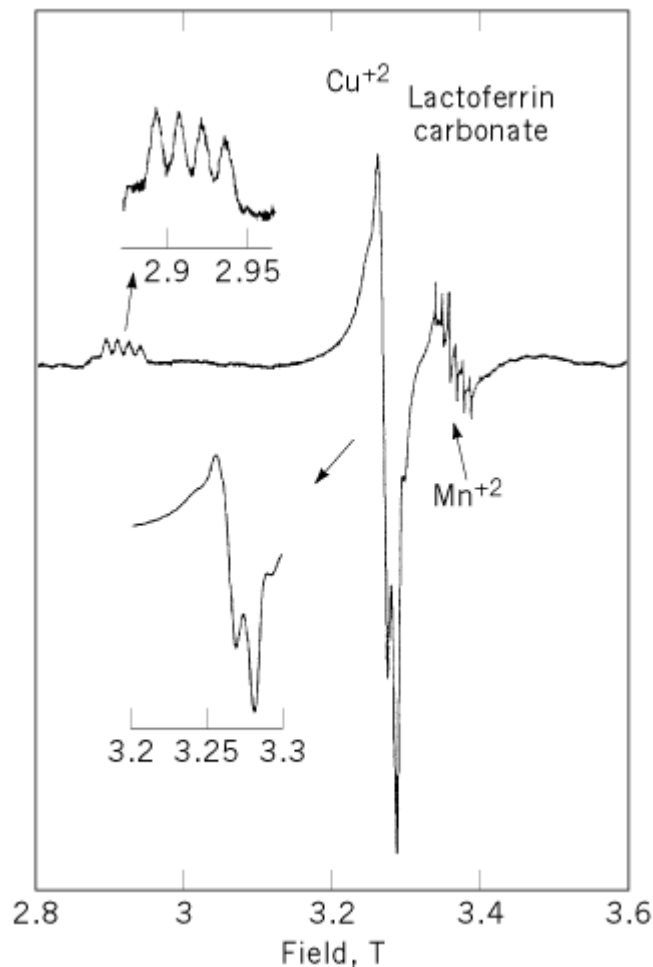


## 5. EPR of Metal Ions in Biology

All metals ions containing one or more unpaired electrons (paramagnetic ions) are in principle subjects for study by EPR, but the conditions for the experiments are quite varied (7). The paramagnetic ions in proteins that are commonly studied include manganese  $^{2+}$ , copper  $^{2+}$ , iron  $^{3+}$ , occasionally iron  $^{2+}$ , nickel $^{3+}$ , and cobalt $^{2+}$ . The conditions for the spectroscopy differ drastically, depending on which metal is the subject of study. Sometimes, manganese ions can be detected with high sensitivity in solution at room temperature, but detection of  $^{17}\text{O}$  ligands to the  $\text{Mn}^{2+}$  (for example from a [threonine](#) side chain in p21 *ras*) required both low temperature and high frequency EPR (8). ENDOR and ESEEM (1) are alternative methods for detecting nuclei, such as nitrogen or hydrogen bound to, or near, metal-ion sites. Copper sites also can be detected at room temperature, but sensitivity is improved by conducting the EPR experiments at or lower than liquid nitrogen temperature. Studies of iron are almost always done at a temperature near that of liquid helium. The two primary reasons that low temperature is used in metal ion EPR are that relaxation times are too fast for temperature studies, and the signal intensity increases inversely with temperature.

Quantitative evaluation of the number of unpaired spins in different sites within a complex biochemical electron transfer system can be made by EPR and associated studies. For experiments of this type, calculations of the theoretical spectra are often needed for interpretation. Recent advances in EPR spectroscopy, particularly high-frequency EPR, provide improved resolution in samples with multiple EPR-detectable sites. Figure 3 shows a high frequency EPR spectrum of copper ion in lactoferrin, in which the copper spectrum is well separated from the signal of a manganese impurity (8, 9). At lower EPR frequencies, the manganese signal is superimposed on the central portion of the copper signal.

**Figure 3.** High-frequency EPR spectrum of di-cupric lactoferrin. The spectrum was recorded at 40 K and 94.1 GHz. The two regions from which significant information can be obtained are amplified in insets on the left. A feature indicated on the right is the EPR signal from an impurity of manganese ion in the sample. (Reproduced with permission from Ref. 9).



## 6. Spin Trapping and EPR Imaging

The EPR spectrum shown in Figure 1 was actually from a spin trapping experiment (4). In that case, the effluent from stimulated endothelial cells was passed into a solution of deoxy [myoglobin](#), and the nitric oxide was trapped as the myoglobin heme adduct. Spin traps can also exhibit therapeutic effects in inhibition of lipoprotein oxidation (10). In a sense, EPR imaging of oxygen concentration in tissues of animals is also a spin trap experiment, but, in this case, only the magnetism of oxygen is “trapped” through its magnetic, but not chemical, interactions with a **spin label**. The degree of broadening of the EPR signal is transformed into an image (11).

### Bibliography

1. H. Levanon and K. Möbius (1997) *Ann. Rev. Biophys. Biomol. Struct.* **26**, 495–540.
2. W. Hofbauer, A. Zouni, R. Bittl, J. Kern, P. Orth, F. Lenzian, P. Fromme, H. T. Witt, and W. Lubitz (2001) *Proc. Natl. Acad. Sci. U.S.A.* **98**, 6623–2228.
3. R. D. Britt, J. M. Peloquin, and K. A. Campbell (2000) *Ann. Rev. Biophys. Biomol. Struct.* **29**, 463–495.
4. D. J. Singel and J. R. Lancaster, Jr. (1996) *Methods in Nitric Oxide Research*, M. Feelisch and J. S. Stamler, eds. John Wiley & Sons, pp. 341–356.
5. J. A. Stubbe (2001) *Trends Biochem. Sci.* **26**, 93–99.
6. A. Gräslund and M. Sahlin (1996) *Ann. Rev. Biophys. Biomol. Struct.* **27**, 259–286.
7. L. J. Berliner and J. Reuben, eds. (1993) In *Biological Magnetic Resonance*, Vol. **13**, EPR of paramagnetic molecules. Plenum, New York.

8. B. F. Bellew, C. J. Halkides, G. G. Gerfen, R. G. Griffin, and D. J. Singel (1996) *Biochem* **35**, 12186–12193.
9. B. J. Gaffney, B. C. Maguire, R. T. Weber, and G. G. Maresch (1999) *Appl. Magn. Res.*, **16**, 207–222.
10. C. E. Thomas, D. F. Ohlweiler, and B. Kalyanaraman (1994) *J. Biol. Chem.* **269**, 28055–28061.
11. A. T. Yordanov, K. Yamada, M. C. Krishna, J. B. Mitchell, E. Woller, M. Cloninger, and M. W. Brechbiel (2001) *Angew. Chem. Int. Ed.* **40**, 2690–2691.

### Suggestions for Further Reading

12. J. A. Weil, J. R. Bolton, and J. E. Wertz (1994) *Electron Paramagnetic Resonance, Elementary Theory and Practical Applications*, Wiley-Interscience, New York.
13. L. J. Berliner and J. Reuben, eds. (1998) "Biological Magnetic Resonance", Vol. **14**, *Spin labeling: The next millenium*, Plenum, New York.
14. G. R. Eaton, S. S. Eaton, and K. Ohno, eds. (1991) *EPR Imaging and In Vivo EPR*, CRC Press, Boca Raton.

## Electron Microscopy

High-resolution electron microscopy presents to the molecular biologist advantages and disadvantages compared to both [X-ray crystallography](#) and [NMR](#). The advantages are that (i) the [macromolecule](#) of interest need not necessarily be highly purified, provided that it can be unambiguously recognized, (ii) only small quantities (mg rather than mg amounts) are required, (iii) even very large macromolecules or complexes with no symmetry can be imaged easily, (iv) specimens can be imaged in their natural environment, and (v) unlike X-ray diffraction patterns, images contain both structure factor amplitudes and phases. The major disadvantages are that (i) the contrast of biological preparations in the electron microscope is inherently low, (ii) the specimen must be imaged in a high vacuum, and (iii) biological material is easily damaged by electron radiation. Electron microscopy has developed to the state where it can complement X-ray crystallography and NMR; this is especially true in the case of [cryoelectron microscopy](#) of frozen specimens with the application of computer image processing techniques. An example of how all three techniques may complement each other was the structural characterization of dihydrolipoyl transacetylase ([1-3](#)). X-ray crystallography provided the rigid core structure. NMR provided information about the flexible domains of this molecule, and electron microscopy integrated both packets of information into a detailed model.

The transmission electron microscope (TEM) is the electron counterpart to the transmission light microscope. A high-velocity, homogeneous electron beam passes through a specimen and is imaged on a phosphorescent screen or a charge-coupled device (CCD). The mechanism of image formation, however, for the TEM is by scattering, which is quite different from that of the light microscope, which is by absorption. TEMs are divided into different classes, depending on the accelerating voltages. A conventional TEM (CTEM) accelerates electrons with typical voltages between 60 and 120 kV; an intermediate-high voltage TEM (IVEM) operates between 200 and 400 kV; and a high-voltage TEM (HVEM) operates above 500 kV. The principal advantages of higher voltage TEMs are (i) the capability of imaging thicker specimens, and (ii) increased resolution of biological crystals and ice-embedded particles.

### 1. Thin-Section CTEM

The vast majority of biological specimens are much too thick as they occur in nature to be penetrated by the electron beam. A CTEM allows for quality imaging of specimens generally no thicker than 200 nm. Hence, the requirement for thin sections of biological material that are able to withstand both the high vacuum in the electron microscope and electron irradiation is apparent. Electron microscopy of thin sections of cells and tissues provides in situ visualization of cellular morphology and the spatial relationships of subcellular components. In order to preserve detailed structural information, the tissue must be well-fixed, dehydrated and embedded in plastic. The resolution achieved with tissue prepared this way is usually limited to 2–3 nm. Sections from the embedded tissue should be of a known and uniform thickness and free from breaks, folds, wrinkles, and chatter to produce standardized and reproducible observations (4). Because it is difficult to differentiate among the electron opacity of various cellular components, selective staining is necessary. Commonly used stains are heavy metals such as uranyl acetate and lead salts.

## 2. Freeze Fracturing and Etching

Thin-section techniques can be complemented by examining specimen surfaces by freeze fracturing/etching. The freeze fracture method involves the following steps: (i) cryofixation of the specimen by rapid freezing, (ii) fracturing of the specimen by cleaving it with a knife edge or breaking it, (iii) replication of the freeze-fracture plane by coating with heavy metals under high vacuum, (iv) removal of the specimen from the replica, and (v) observation of the replicas with a CTEM. The most important contribution of this technique has been the elucidation of the molecular architecture of biological [membranes](#). During freeze fracturing, membranes are split through the hydrophobic moiety, thus exposing intramembraneous surfaces. The fracture plane often follows the contours of membranes and leaves either protrusions or depressions where it passes around organelles or vesicles (4). Smooth areas represent the lipid monolayer face, whereas particles represent integral [membrane proteins](#).

*Freeze etching* is a slight modification of freeze fracturing through the process of controlled sublimation of water vapor (transition from the solid state to the gas state, bypassing the fluid state) to expose membrane surfaces or macromolecules just below the fracture plane. By etching, extracellular and intracellular structures, such as the [cytoskeleton](#), can be visualized which would otherwise be masked. The etching can either be “normal” (< ~ 100 nm) or “deep” (up to ~ 3 μm).

Related techniques involve shadowing of macromolecules adsorbed or deposited onto a flat surface and dried by methods which minimized deformation during drying. Drying is best accomplished by sublimation as in freeze-etching and requires access to a freeze-fracture instrument; this is referred to as freeze drying and shadowing (5). A simpler method developed by Branton requires only a conventional vacuum evaporator (6). This technique involves (i) preparing the macromolecule in a high-concentration glycerol solution, (ii) spraying this solution onto mica, (iii) evaporating the solvent under vacuum, (iv) rotary shadowing with a metal such as platinum, (v) carbon coating the shadowed specimen, (vi) floating the carbon replica onto water, and (vii) picking up the replica on an electron microscope grid. Rotary shadowing is particularly useful for visualizing fibers or rod-shaped complexes such as actin/myosin. The double strands of [spectrin](#) (which forms rods ~200 nm long) were resolved using this technique with tantalum/tungsten shadowing, which produced small metal grains that increase the resolution (7).

## Bibliography

1. D. DeRosier (1992) *Nature* **357**, 196–199.
2. A. Mattevi et al. (1992) *Science* **255**, 1544–1547.
3. M. A. Robin et al. (1992) *Biochemistry* **31**, 3463–3471.
4. M. A. Hayat (1989) *Principles and Techniques of Electron Microscopy*, CRC Press, Boca Raton, FL.



5. J. Kistler, U. Aepli, and M. Kellenberger (1977) *J. Ultrastruct. Res.* **59**, 76–86.
6. J. M. Tyler and D. Branton (1980) *J. Ultrastruct. Res.* **71**, 95–102.
7. J. R. Glenney (1987) In *Electron Microscopy in Molecular Biology* (J. Somerville and U. Scheer, eds.), IRL Press, Oxford, pp. 167–178.

## Electron Tomography

Electron tomography is defined as the generation of a 3-D reconstruction from an electron microscope tilt series (1). This powerful technique is closely related to computerized axial tomography used by CAT scanners in radiological imaging in the computational methods used to calculate a 3-D structure from many two-dimensional images or projections recorded over a wide range of tilt angles. The use of tomography has proved valuable not only for the determination of macromolecular complexes (2-4), but also for the visualization and analysis of relatively large, complex biological structures. Because of variable size, large-scale biological structures such as organelles elude crystallographic or single-particle approaches that require multiple images of identical structures. Electron tomography is especially powerful for these complex structures, but until recently, it has not enjoyed widespread application because of obstacles that are now being removed (5). Electron tomography presently is the imaging technique that provides the highest 3-D resolution (capable of 5–10 nm) of the internal features of organelle-size structures. This is accomplished by recording a series of images of a single specimen over a wide range of tilt angles with a small interval of tilt angle. Typically such a tilt series would consist of 61 images recorded at tilt angles from  $-60^\circ$  to  $+60^\circ$  in  $2^\circ$  increments. The individual images must be digitized (unless recorded directly in digital form using a CCD camera attached to the microscope), aligned to a common origin, and processed using computer algorithms, similar to those used in computerized tomography in medical imaging, to reconstruct a 3-D volume from a series of 2-D projections. Since whole cells and organelles, such as [mitochondria](#), are relatively large, the images must be recorded from specimens embedded in semithick sections, 0.25–1  $\mu\text{m}$  thick. Such thick specimens require the use of higher voltage electron microscopes, intermediate voltage electron microscopes (IVEM's) with accelerating voltages up to 400,000 volts, or high-voltage electron microscopes with accelerating voltages up to 1,200,000 volts. By comparison, typical transmission electron microscopes have accelerating voltages of approximately 100,000 volts. The higher voltage electron microscopes are relatively expensive, and most users wishing to do electron tomography of thick specimens must use instruments made available in national facilities.

Until recently, the bulk of our knowledge of the architecture of organelles and bacteria has come from untilted images of sections. The 3-D architecture is usually inferred from these images, sometimes with the aid of stereo imaging or serial-section reconstruction. Although thin-section electron microscopy provides relatively high-resolution images, an incorrect impression of the 3-D structure may be obtained because one is looking at only a very thin slice through a complex 3-D object. On the other hand, electron tomography offers an opportunity to map topology more accurately by providing improved resolution along the  $z$ -axis, while circumventing major deficiencies in stereo imaging and serial-section reconstructions. This technique provided the resolution necessary to observe characteristics of mitochondrial structure that differ from long-held conceptions (6,7). These newly appreciated characteristics include: (i) All observed [cristae](#) connect to the inner boundary membrane via narrow, tubular openings, termed crista junctions. (ii) Tubular cristae merge, sometimes from opposite ends of the mitochondrial periphery, to form lamellar compartments. (iii) Contact sites are not clustered about crista junctions. The full power of electron tomography is now being explored to generate 3-D distributions of specifically labeled or

immunolocalized components.

## Bibliography

1. J. Frank (1992) *Electron Tomography*, Plenum Press, New York.
2. R. A. Horowitz et al. (1994) *J. Cell Biol.* **125**, 1–10.
3. H. Mehlin, B. Daneholt, and U. Skoglund (1992) *Cell* **69**, 605–613.
4. M. Moritz et al. (1995) *Nature* **378**, 638–640.
5. J. Frank (1995) *Curr. Op. Struct. Biol.* **7**, 266–272.

## Electron Transfer Proteins

Biological electron transfer (ET) encompasses reactions that are essential for life, such as those involved in **light harvesting**, **nitrogen fixation**, and **production of the high-energy metabolites**. In proteins, ET occurs between redox centers that are separated by distances much longer than a chemical bond and that in many cases are carried by different molecules. Moreover, the temperature-dependence of ET rates is often atypical (ie, non-Arrhenius type).

A general theory (1) can describe ET and thus provide a framework to guide experimentalists in their investigations on the structural control of biological ET. From these studies, it emerged, as is frequent in biology, that different solutions have evolved in different systems to solve the same general problem.

Several families of proteins involved in biological ET, with variable structures and catalytic complexity, are known in detail (see **Redox proteins**). Although a compilation is beyond the scope of this article, a brief synopsis is provided. Many of these are one-electron carrier proteins that are usually small (about 100 amino acid residues), and display no associated enzymatic function. In these proteins, the electron resides on a special cofactor, such as a heme or a metal atom, which often can donate the electron to different partners. Examples of this group are the [cytochromes](#), cupredoxins, and [ferredoxins](#). Other more complex proteins consist of one or more **domains** and may perform chemical reactions, such as oxidation or reduction of organic molecules. Some of them can convert the flow of charge from single electron to pairs (or more) using (i) cofactors like flavins and quinones, forming stable radical intermediates, or (ii) radicals of amino acid side chains, such as those of [tyrosine](#) or [tryptophan](#), together with metal centers. Among others, it is worth mentioning the various respiratory complexes located in organelle and plasma [membranes](#), the large class of detoxifying enzymes (eg, [cytochrome P-450s](#)), and [photosynthetic reaction centers](#) .

### 1. General Theory and Intramolecular ET

The probability of electron transfer depends on the overlap between the electron-containing orbital (wavefunction) of the donor with that of the acceptor. This overlap can be very small, because redox centers in proteins are buried and often located 10–15 Å apart; nevertheless, biological ET is usually extremely fast and specific.

In intramolecular ET, the electron must move through the protein matrix, which has been assigned a low dielectric constant ( $D = 2$  to 4). In a classical picture, this energy barrier would be too high to be crossed by the electron; in a quantum-mechanical view, however, the electron can tunnel through the barrier with a finite probability. To relate tunneling theory and experiments, a simple description of

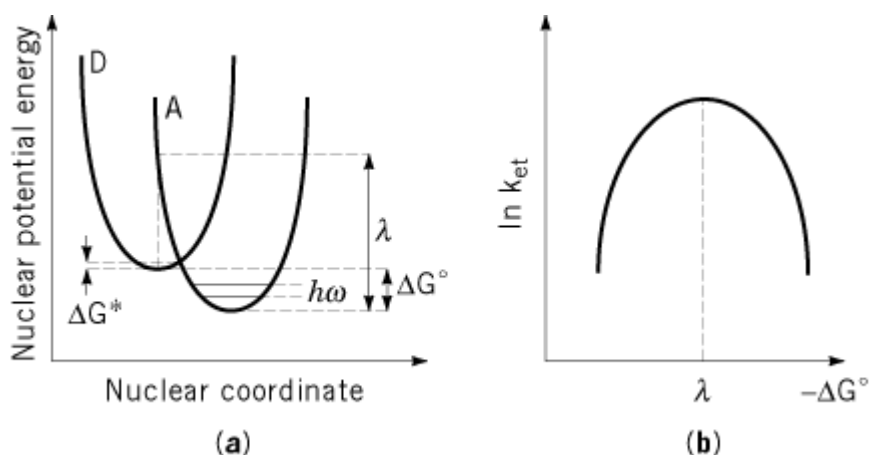
the ET rate is given by Fermi's Golden Rule, which applies to nonadiabatic ET occurring between weakly coupled redox centers:

$$k_{\text{et}} = (\text{electronic factor}) \cdot (\text{nuclear factor}) = (2\pi/\hbar)H_{\text{AB}}^2(\text{FC}) \quad (1)$$

The electronic factor ( $H_{\text{AB}}^2$ ), which describes the electronic coupling between donor and acceptor, can only be understood in quantum-mechanical terms and will be discussed below;  $\hbar$  is the Planck constant/ $2\pi$ . The nuclear factor, or Franck–Condon term (FC), states that during ET the nuclei do not have time to respond, because of their large mass relative to that of the electron.

Marcus and Sutin (1) have provided a simple description of the energy terms that comprise the FC factor using potential energy curves (Fig. 1). The energy terms that determine the rate of ET are the **free energy** change ( $\Delta G^\circ$ ), or driving force, and the reorganizational energy ( $\lambda$ ) (which is required to distort the equilibrium geometry of the reactant to that of the product, but without actually transferring the electron). ET will occur when the potential energies of reactants and products are equal.

**Figure 1.** (a) Potential energies of the ET donor (D) and acceptor (A), as a function of the reaction coordinate. The energies with their surrounding medium are approximated as harmonic oscillator potentials. The Franck–Condon factor can be described classically, provided that the nuclear vibration energy ( $\hbar\omega$ ) is not too high relative to the thermal energy ( $k_B T$ ). (b) Dependence of the logarithm of the ET rate constant on  $-\Delta G^\circ$ . The ET rate is maximal at the point where  $-\Delta G^\circ = \lambda$  (activationless regime).



In general, the activation barrier ( $\Delta G^*$ ) is given by the expression:

$$\Delta G^* = (\Delta G^\circ + \lambda)^2 / 4\lambda \quad (2)$$

Depending on the sign of  $(\Delta G^\circ + \lambda)$ , one can distinguish normal, activationless, and inverted regions, respectively with  $(\Delta G^\circ + \lambda) > 0$ ,  $= 0$ , and  $< 0$ . The ET rate goes through a maximum for  $\Delta G^\circ = -\lambda$  (activationless regime), and decreases on both sides (Fig. 1), as shown in simple systems and in several redox proteins. For example, in the photosynthetic reaction center, charge recombination is prevented because the very large driving force for the initial charge separation drives the recombination reaction into the inverted region, slowing down the corresponding rate.

Using the classical approximation for the FC factor, the Fermi equation becomes

$$k_{\text{et}} = (2\pi/\hbar)H_{\text{AB}}^2(4\pi\lambda k_{\text{B}}T)^{-1/2} \exp[-(\Delta G^\circ + \lambda)^2/4\lambda k_{\text{B}}T] \quad (3)$$

where  $k_{\text{B}}$  is the Boltzmann constant and  $T$  is the temperature. According to Eq. (3), the main temperature-dependence comes from the exponential factor, but when the [activation energy](#) is zero ( $e^0 = 1$ ), a weak temperature-dependence remains because of the preexponential factor, causing the rate to decrease as the temperature increases. Again this behavior has been observed for some reactions of the photosynthetic reaction centers.

The temperature-independence of some ET reactions can be explained, however, only by using a semiclassical description of the FC term, proposed by Hopfield (2), in which the term  $2\lambda k_{\text{B}}T$  has been replaced by  $\hbar\omega \coth(\hbar\omega/2k_{\text{B}}T)$ . If the temperature falls below  $\hbar\omega/2k_{\text{B}}$ , the rate becomes nearly temperature-independent at any value of  $\Delta G^\circ$ .

The other term influencing the ET rate, the electronic factor ( $H_{\text{AB}}$ ), represents the weak coupling of the reactant and product wavefunctions, where  $H$  stands for ‘‘Hamiltonian,’’ the quantum-mechanical description of the system.  $H_{\text{AB}}$  yields the overlap between the orbitals A and B and is often referred to as the ‘‘electronic coupling matrix element.’’ The degree of orbital overlap depends on the distance between the redox centers and the nature of the intervening medium, which in biological ET is the protein matrix and the solvent.

Because wavefunctions decay exponentially with distance, the electronic coupling will decrease with the distance ( $r-r_0$ ) according to:

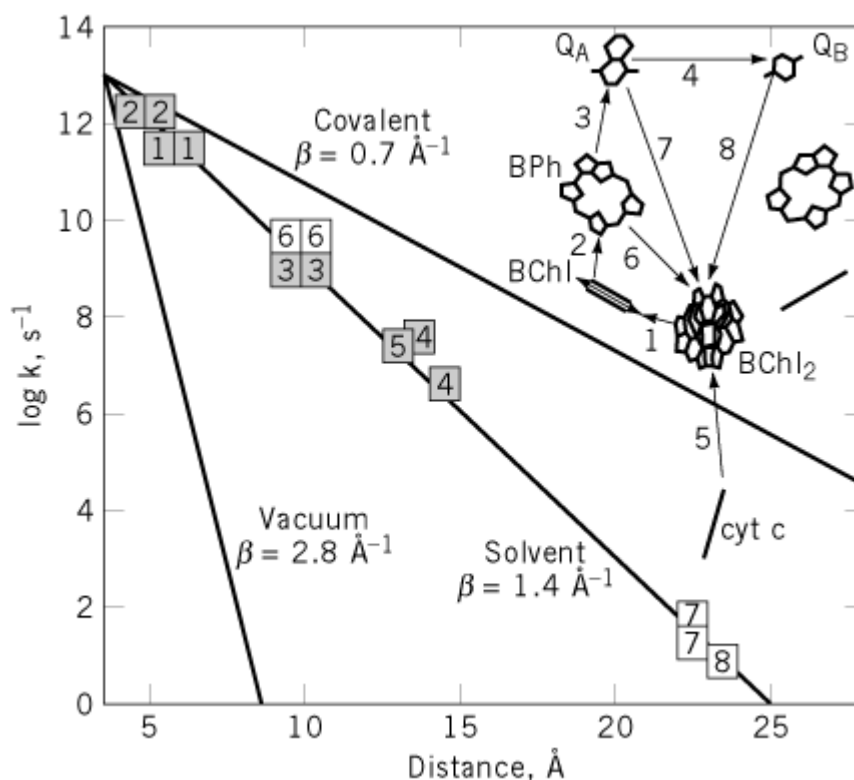
$$H_{\text{AB}}^2 = (H_{\text{AB}}^0)^2 e^{-b(r-r_0)} \quad (4)$$

where  $H_{\text{AB}}^0$  represents the electronic coupling between A and B when the redox centers are in **van der Waals** contact ( $r = r_0$ ), and its decay with distance ( $r - r_0$ ) occurs with the coefficient  $b$ . The latter therefore describes the contribution of the intervening medium in propagating the wavefunction, and its dependence on protein structure is a matter of intense research and debate.

To obtain a correct description of the role of the intervening medium and distance on  $k_{\text{et}}$ , it is necessary to work in the activationless regime ( $\Delta G^\circ = -\lambda$ ) (Fig. 1). The photosynthetic reaction center proved to be an excellent test system for several reasons: (i) The three-dimensional structures of a number of reaction centers are known at atomic resolution; (ii) it contains several redox cofactors at various fixed distances; (iii) the driving force can be varied experimentally; and (iv) ET can be studied over a wide temperature range. Using experimental data obtained with proteins where the primary quinone was substituted with other compounds of different [oxidation/reduction potential](#), Dutton and coworkers (3, 4) have shown a linear dependence of  $\ln k_{\text{et}}$  on distance (over 12 orders of magnitude), with  $b = 1.4 \text{ \AA}^{-1}$ ,  $r_0 = 3.6 \text{ \AA}$ , and a preexponential factor  $[(4\pi^2/\hbar)(H_{\text{AB}}^2)]$  of  $10^{13} \text{ s}^{-1}$  (Fig. 2). These authors concluded that the ET rate in any protein follows an exponential decay of the electronic wavefunctions with distance, implying a homogeneous intervening medium.

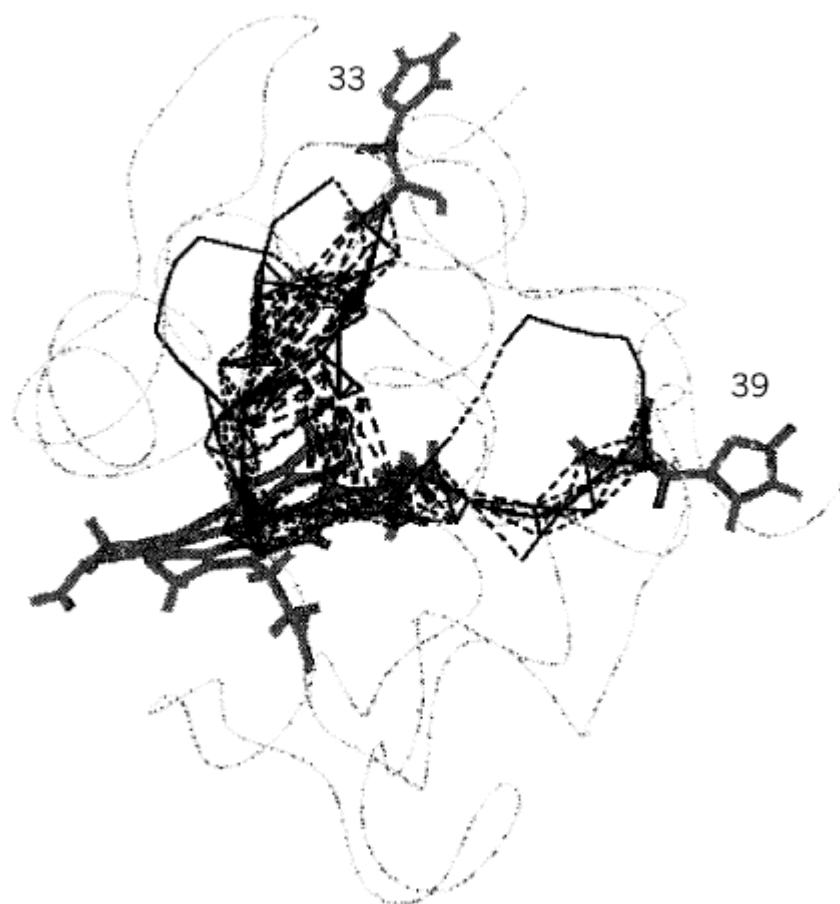
**Figure 2.** Effect of distance on the intraprotein ET rate in the photosynthetic reaction center, according to Moser and Dutton (3). Following light excitation of the bacteriochlorophyll dimer (BChl<sub>2</sub>), ET progresses sequentially through bacteriochlorophyll monomer (BChl), bacteriopheophytin (BPh), primary quinone, and secondary quinone (Q<sub>A</sub> and Q<sub>B</sub>). The resulting positively charged oxidized dimer can be re-reduced by the nearest of the four cytochromes  $c$  (cyt  $c$ ). The experimentally determined rates within the reaction center fit a straight line with  $b = 1.4 \text{ \AA}^{-1}$ , compared to  $2.8 \text{ \AA}^{-1}$

in vacuum and to  $0.7 \text{ \AA}^{-1}$  in a rigid covalent system. The line through the experimental data extrapolates to  $10^{13} \text{ s}^{-1}$  at 3 ANGSTROM, the van der Waals distance. (From Ref. 8, reproduced with permission. Copyright BIOS Scientific Publishers, 1996.)



This exceptional result (Fig. 2) does not prove, however, that the tunneling barrier between the redox centers is microscopically uniform, but rather that the observed value represents an average. Gray and coworkers (5) have used proteins modified by incorporation of ruthenium (mainly [myoglobins](#) and cytochromes) as model systems to investigate the dependence of  $b$  on the structure of the intervening medium, that is, the detailed protein structure. In this and other similar studies, a plot of  $\ln k_{\text{et}}$  against distance does not always yield a linear relationship, suggesting that the protein does not always provide a homogeneous energy barrier. Given the intrinsically heterogeneous packing inside a protein, and thus the variable dielectric constant, a model that takes into account the detailed structure of the intervening medium is needed. Onuchic and Beratan (6) have developed such a model, which views coupling between redox centers via pathways. In this model, the exponential decay of the electronic coupling depends not only on the distance between redox sites, but also on the detailed contacts along the ET pathway. Along a pathway, electronic decay factors ( $\epsilon$ ) through different connectivities, involving covalent ( $\epsilon_{\text{C}}$ ), [hydrogen-bond](#) ( $\epsilon_{\text{H}}$ ), and through-space ( $\epsilon_{\text{S}}$ ) jumps, come into play. The coupling strength decreases in the order covalent > hydrogen bond > through-space, and the total decay across the pathway is given by the sum of the individual contributions. An algorithm searching for potential ET pathways in proteins of known three-dimensional structure has been developed by these authors (Fig. 3). This model has been applied to different redox proteins in order to define the role of the intervening medium and of **secondary structure** elements on the ET rates. The model predicts that (i) hydrogen bonds provide shortcuts for the electron between otherwise unconnected protein segments, (ii) [beta-sheet](#) proteins will exhibit a greater degree of coupling than **alpha-helical** proteins, and (iii) pathways always arise in families, giving rise to the concept of “fattened” pathways and pathway tubes. The existence of multiple pathways is a challenge to [protein engineering](#) of ET proteins to enhance or inhibit ET rates by modification of residues along the most probable pathways.

**Figure 3.** Connectivity between redox centers according to Beratan, Onuchic, and coworkers. The figure depicts the structure of cytochrome *c*, with the heme group highlighted and the polypeptide chain sketched as a continuous line. The two histidine residues at position 33 and 39 are the sites of chemical modification with ruthenium. Dotted lines represent the calculated pathway for ET. The strongest coupling is provided by bonds at the center of a tube of pathways that connect the donor–acceptor couples within the protein. The decay factors were calculated to be: covalent bonds,  $\epsilon_c = 0.6$ ; hydrogen bonds,  $\epsilon_H = \epsilon_c^2 e^{-\beta(r-r_0)}$ ; through-space jumps,  $\epsilon_s = 1/2\epsilon_c e^{-\beta(r-r_0)}$ , where  $r_0$  is the reference distance for covalent bonds (1.4 Å) or hydrogen bonds (2.8 Å) and  $b$  is the decay length for through-space jumps (1.7 Å). (From Ref. 9, with permission. Copyright American Chemical Society.)



The large body of experimental data available today indicates that the two models (homogeneous and heterogeneous) are not necessarily exclusive. In some cases, the protein behaves as a homogeneous medium with “distance-dependent” exponential decay, as for photosynthetic reaction center, ruthenated [azurin](#), and ruthenated cytochrome  $b_5$ , although different  $b$  values were reported. In other cases, the coupling decay factor was shown to be anisotropic (eg, in cytochrome *c*), giving rise to a view in which “hot” and “cold” spots for ET can be identified. Of particular importance are studies on the role of dynamic fluctuations on through-space contacts.

## 2. Intermolecular ET

The principles underlying ET between two proteins offer another theme for variation and control, namely, molecular recognition and [protein–protein interactions](#) as a prerequisite for ET. Some general features that have emerged from many studies will be presented in this section.

The Marcus theory (1) described above states that electron coupling between redox centers depends primarily on distance. For efficient ET, diffusible proteins must interact in such a way to maximize proximity between the redox centers within a collisional complex (AB), given that most redox active proteins are asymmetric, with their cofactor/metal generally located on one side of the macromolecule. The kinetic model that describes interprotein ET is



Under pseudo-first order conditions (eg,  $[A] \gg [B]$ ), the overall rate constant  $k_{12}$  is given by

$$k_{12} = k_a k_{et} [A] / (k_{-a} + k_{et} + k_a [A]) \quad (6)$$

Within this model, two limiting conditions can be envisaged, where either complex formation or intracomplex ET are rate-limiting.

Marcus theory allows describing the overall rate constant  $k_{12}$  as a function of the equilibrium constant  $K_{12}$ , the self-exchange rate constants of the two partners,  $k_{11}$  and  $k_{22}$ , and the work term  $W_{12}$  involved in the configurational change of reactants and products along the reaction coordinate:

$$k_{12} = (k_{11} k_{22} K_{12} f_{12})^{1/2} W_{12} \quad (7)$$

where

$$W_{12} = \exp(w_{11} + w_{22} - w_{12} - w_{21}) / 2RT \quad (8)$$

Equation (7) can be simplified when (a) the work terms cancel one another (ie, homonuclear and heteronuclear ET involve the same energy terms) and (b)  $f_{12}$  is close to unity, which is often the case provided that  $K_{12}$  is not too large. Under these conditions, equation (7) becomes

$$k_{12} = (k_{11} k_{22} K_{12})^{1/2} \quad (9)$$

which is frequently used and is generally referred to as the “cross-relation.” For redox reactions between inorganic metal ions or complexes, Marcus theory has been experimentally verified many times, either in the general or in the “cross-relation” formulation. It has also been used to calculate reaction rates or equilibrium constants that are difficult to determine experimentally. The “cross-relation” approach has also been used to interpret protein–protein ET: Table 1 reports a set of experimentally determined rate constants for some copper and heme proteins, together with the values calculated on the basis of equation (8). When comparison with known equilibrium and self-exchange rate constants was possible, the discrepancy between calculated and measured values was found to be no greater than one order of magnitude (which is considered good). In cases in which larger deviations were observed (>100 fold), efforts have been made to rationalize this discrepancy, possibly reconsidering the applicability to the system of some of the theoretical assumptions described above. The obvious assumption that may not apply is related to the anisotropy of the redox protein and thereby to the contention that ET with any partner always occurs through one and the same contact surface.

**Table 1. Comparison of Experimentally Determined and Calculated (According to Marcus “Cross-Relation”) Second-Order Rate Constants for Different Redox Couples**

| Oxidant <sup>a</sup><br>(1) | Reductant <sup>a</sup><br>(2) | $K_{12}$ <sup>b</sup> | $R_{11}$ <sup>c</sup> (M <sup>-1</sup> s <sup>-1</sup> ) | $k_{22}$ <sup>c</sup> (M <sup>-1</sup> s <sup>-1</sup> ) | $k_{12}$ calc.<br>(M <sup>-1</sup> s <sup>-1</sup> ) | $k_{12}$ exp.<br>(M <sup>-1</sup> s <sup>-1</sup> ) |
|-----------------------------|-------------------------------|-----------------------|--|--|--|---|
| Azurin                      | Cyt <i>c</i>                  | 19                    | $8 \times 10^5$  | $2.5 \times 10^2$  | $5.5 \times 10^4$                                    | $6.4 \times 10^4$                                   |
| Azurin                      | Cyt <i>c</i> <sub>551</sub>   | 3                     | $8 \times 10^5$  | $1.2 \times 10^7$  | $4.8 \times 10^6$                                    | $6 \times 10^6$                                     |
| Cyt <i>c</i>                | Cyt <i>c</i> <sub>551</sub>   | 1                     | $2.5 \times 10^2$  | $1.2 \times 10^7$  | $5.5 \times 10^4$                                    | $5 \times 10^4$                                     |
| Plastocyanin                | Cyt <i>c</i>                  | 40.4                  | $10^3$   | $2.5 \times 10^2$  | $3.2 \times 10^3$                                    | $1.5 \times 10^6$                                   |
| Plastocyanin                | Cyt <i>c</i> <sub>551</sub>   | 40.4                  | $10^3$   | $1.2 \times 10^7$  | $7 \times 10^5$                                      | $7.5 \times 10^5$                                   |
| Stellacyanin                | Cyt <i>c</i>                  | 0.06                  | $1.2 \times 10^5$  | $2.5 \times 10^2$  | $1.3 \times 10^3$                                    | $3.5 \times 10^2$                                   |
| Stellacyanin                | Cyt <i>c</i> <sub>551</sub>   | 0.08                  | $1.2 \times 10^5$  | $1.2 \times 10^7$  | $3.3 \times 10^5$                                    | $1.6 \times 10^5$                                   |

<sup>a</sup> Proteins: *Pseudomonas aeruginosa* azurin and cytochrome *c*<sub>551</sub> (Cyt *c*<sub>551</sub>), horse heart cytochrome *c* (Cyt *c*), and parsley plastocyanin.

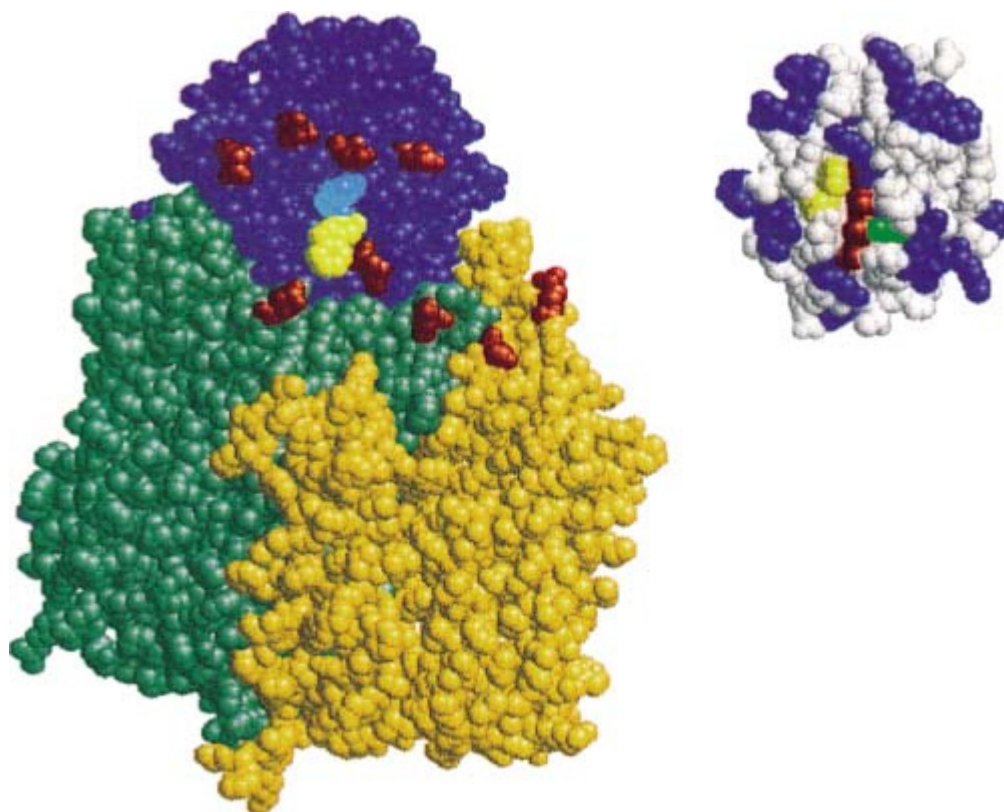
<sup>b</sup> Equilibrium constants were calculated on the basis of midpoint redox potentials ( $E_{m,7}$ ).

<sup>c</sup> Electron self-exchange rates at 20°C used in the calculation of  $k_{12}$  according to Marcus “cross-relation.”

In the **diffusion**-limited case, maximal second-order rate constants for ET proteins are in the range of  $10^8$  to  $10^9$  M<sup>-1</sup> s<sup>-1</sup>. A negative steric factor, due to the fact that only a small fraction of the total surface of a protein is near the redox active cofactor, is partially compensated for by the solvent cage effect and long-range [electrostatic interactions](#). The solvent cage effect increases the lifetime of the so-called encounter complex because, after collision, the two proteins are momentarily trapped by the solvent and thus experience a large number of mutual contacts by two-dimensional diffusion. Long-range electrostatic forces come into play, because many ET proteins display an anisotropy in surface charge distribution that results in a permanent dipole moment, favoring their association. This effect has been investigated through (i) the ionic-strength dependence of the second-order rate constant  $k_{12}$ , (ii) the effect of surface-charge modification by chemical modification or [protein engineering](#), and (iii) the influence of the computed electrostatics of the redox partners in simulated docking experiments. An interesting generality related to the electrostatic makeup of ET proteins is that the docking surfaces often show a “loose” specificity, which also explains the high degree of cross-reactivity observed both *in vitro* and *in vivo* between donors and acceptors. This is shown in [Figure 4](#) for the widely studied redox couple cytochrome *c* and cytochrome-*c*-oxidase of the respiratory chain, where the complex between the two redox partners is stabilized by multiple electrostatic interactions, and the ET pathway crucially depends on two aromatic side chains inside a patch of negative charges on the cytochrome-*c*-oxidase binding site ([7](#)).

**Figure 4.** Models of the three-dimensional structures of *Paracoccus denitrificans* cytochrome-*c* oxidase and horse heart cytochrome *c*, in a CPK representation. Color code: (i) On cytochrome *c*, *blue* represents positive residues (mainly lysines), *red* represents the heme, *yellow* represents phenylalanine 82, and *green* represents cysteine 17; (ii) on cytochrome-*c*-oxidase, *red* represents the negative residues (aspartates and glutamates) on subunit I (*green*), subunit II (*blue*) and subunit III (*brown*). *Yellow* represents residue tryptophan 121, and *light blue* represents tyrosine 122 forming the hydrophobic patch mediating ET ([7](#)). [Courtesy of Prof. F. Malatesta (University of Aquila, Italy) using the coordinates kindly provided by Prof. H. Michel (Frankfurt, Germany).] See color insert.





### Bibliography

1. R. A. Marcus and N. Sutin (1985) *Biochim. Biophys. Acta* **811**, 265–322.
2. J. J. Hopfield (1974) *Proc. Natl. Acad. Sci. USA* **71**, 3640–3644.
3. C. C. Moser and P. L. Dutton (1992) *Biochim. Biophys. Acta* **1101**, 171–176.
4. C. C. Moser, J. M. Keske, K. Warncke, R. S. Farid, and P. L. Dutton (1992) *Nature* **355**, 796–802.
5. D. N. Beratan, J. N. Onuchic, J. R. Winkler, and H. B. Gray (1992) *Science* **258**, 1740–1741.
6. J. N. Onuchic and D. N. Beratan (1990) *J. Chem. Phys.* **92**, 722–733.
7. H. Witt, F. Malatesta, F. Nicoletti, M. Brunori, and B. Ludwig (1998) *J. Biol. Chem.* **273**, 5132–5136.
8. C. C. Moser and P. L. Dutton (1996) In *Protein Electron Transfer* (D. S. Bendall, ed.), BIOS Scientific Publishers, Oxford, England, pp. 1–21.
9. J. J. Regan, S. M. Risser, D. N. Beratan, and J. N. Onuchic (1993) *J. Phys. Chem.* **97**, 13083–13088.

### Suggestions for Further Reading

10. D. S. Bendall (ed.) (1996) *Protein Electron Transfer*, BIOS Scientific Publishers, Oxford, England.
11. G. W. Canters and E. Vlijgenboom (eds.) (1998), *Biological Electron Transfer Chains: Genetics, Composition and Mode of Operation*, NATO ASI Series Vol. **512** (Series C), Kluwer Academic Publishers, Dordrecht, Boston.

## Electrophoresis

Electrophoresis is the movement of charged molecules or particles in an electric field. Electrophoresis may be conducted for the purposes of analytical identification or physical characterization of the various components of a sample, or for their preparative separation.

Electrophoretic migration of a molecule in any medium proceeds in the direction toward the electrode of opposite polarity to the net charge of the molecule (eg, negatively charged molecules migrate toward the positively charged electrode; cations move toward the cathode, anions toward the anode). The medium needs to be [buffered](#) at a pH that produces the appropriate direction and rate of migration. The mobilities of [proteins](#) generally increase the further the pH is from their [isoelectric points](#) (pI), increasing their net charge. However, resolution based on net charge differences improves as the net charge is reduced and is optimal when the net charge approaches zero at the pI (see [Isoelectric Focusing](#)). **Nucleic acids** and proteins complexed with the anionic detergent [SDS](#) (sodium dodecyl sulfate) exhibit high mobilities over a range of pH values above neutrality, due to their high surface negative charge densities. In addition to pH, the ionic strength, composition of the solvent, presence of [detergents](#) (see [Hydrophobic Electrophoresis](#)), and temperature govern the electrophoretic mobilities of all species in a sample, their separation, and the maintenance of their native conformations. For example, nucleic acid stability depends on the  $\text{Na}^+$  level. Proteins aggregate at low ionic strength; high ionic strength may exceed the Joule heat dissipation capacity of the electrophoretic apparatus and “cook” the sample during electrophoresis.

Different species can often be separated on the basis of differences in their electrophoretic mobilities, which is defined as the ratio of the migration rate (cm/s) to the electric field strength (V/cm). The mobility ( $\text{cm}^2/\text{s V}$ ) of a macromolecule is a function of the ratio of its surface net charge to its [accessible surface](#) area, plus the resistance to migration exerted by the medium in which the separation takes place. The resistance due to the viscosity of the medium affects species of all sizes equally and therefore does not contribute to their separation; it may, however, be augmented by the presence of polymer networks and permit separations based on differences in size and shape (see also [Gel Electrophoresis](#); [Agarose](#); [Polyacrylamide](#); [SDS-PAGE](#)).

For practical reasons, samples subjected to electrophoresis need to be protected from convective mixing. This can be accomplished within a liquid sample using a density gradient when the migration of the boundary of a region of sample is observed, or in tubes of less than 0.2 mm in diameter, as in [capillary zone electrophoresis](#). Much more commonly, electrophoresis is carried out with a gel when the sample usually migrates as a zone. In this case, the electrophoretic migration rate depends on the pore size of the gel. Alternatively, the sample may migrate within a stack of moving boundaries of buffer components, with a mobility regulated by that of the boundary (see [Disc Electrophoresis](#); [Isotachopheresis](#)).

Separations by electrophoresis and [chromatography](#) share many of their principles of operation, but electrophoresis excels in its ability to provide simultaneous separations of many components, to exploit charge and size differences simultaneously, to be compatible with a continuum of polymer or gel concentrations that produce a “molecular sieve,” and to be applicable to particle sizes ranging from simple organic acids to very large complexes and small cells (see [Particle Electrophoresis](#)).

### 1. Sample detection and quantification

The zone or moving boundary comprising the sample of interest can be identified and quantified by its absorbance, fluorescence, or [radioactivity](#), by specific **antibodies** (see [Protein Blots \(Western Blots\)](#); [Immuno-electrophoresis](#)), by its biological activity (see [Overlay Assay: Enzyme Zymography](#)

of [Plasminogen Activators and Inhibitors](#)), or by stains specific for the class of molecule. Staining with reagents specific for proteins or nucleic acids is the most common method. Before staining samples in electrophoresis gels, the sample normally must be fixed, ie, made insoluble, so that it does not diffuse during the staining process. Solvents such as acetic acid or trichloroacetic acid are often used to fix proteins; frequently, the fixing solution also contains the stain. The most general and sensitive staining method is [silver stain](#) (1). [Coomassie Brilliant Blue R-250](#) is the most widely used protein stain, although it requires destaining by removing the excess dye, unlike the corresponding G-250 stain (2). Moreover, the commonly used acetic acid solutions are poor fixatives compared to trichloroacetic acid solutions (2, 3). Other widely used protein stains are [Ponceau S](#) and Amido black.

Rather than stain after an electrophoretic separation, the analytes can be prestained prior to electrophoresis. Prestaining is required for capillary zone electrophoresis and automated gel electrophoresis. It must be remembered, however, that the prestaining can alter the charge and shape properties of the molecules of the sample. Prestaining of proteins commonly employs derivatization by fluorescein isothiocyanate (see [Gel Electrophoresis](#)); that of nucleic acids uses the intercalating dye [ethidium bromide](#) (which is also used for poststaining without fixation), or the more recently introduced dyes DAPI (4), TOTO, or YOYO (5).

It is often possible to estimate the quantity of a particular species present in a sample by electrophoresis. Densitometry of stained gel bands has been the classical method (6). When a sample is radioactive, it may be detected and quantified by [autoradiography](#) or by counting extracted or solubilized gel slices. There are now commercial instruments for counting radioactivity in entire gel patterns (7). Fluorescence detection is especially useful with capillary zone electrophoresis or in automated gel electrophoresis apparatus (8).

#### Bibliography

1. C. R. Merrill (1987) *Adv. Electrophoresis* **1**, 111–142.
2. B. An der Lan, C. Auzan, J. V. Sullivan, and A. Chrambach (1985) *Electrophoresis* **6**, 408–409.
3. A. Chrambach, R. A. Reisfeld, M. Wyckoff, and J. Zaccari (1967) *Anal. Biochem.* **20**, 150–154.
4. E. Buel and M. Schwartz (1993) *Appl. Theor. Electrophoresis* **3**, 253–256.
5. K. Srinivasan, S. C. Morris, J. E. Girard, M. C. Kline, and D. J. Reeder (1993) *Appl. Theor. Electrophoresis* **3**, 235–240.
6. M. M. Miller (1989) *Adv. Electrophoresis* **3**, 182–220.
7. D. M. Gersten and E. Zapolski (1991) *Adv. Electrophoresis* **4**, 49–79.
8. J. C. Sutherland (1993) *Adv. Electrophoresis* **6**, 3–43.

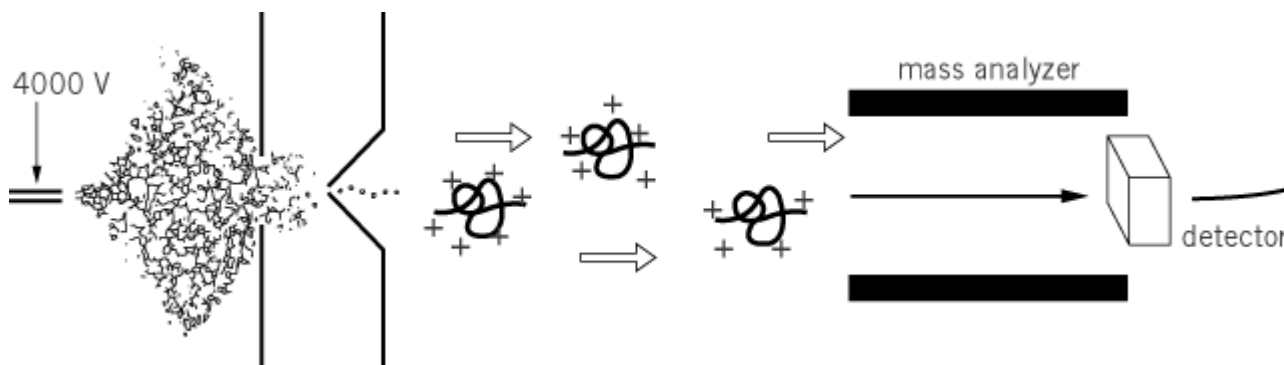
#### Suggestions for Further Reading

9. M. Bier (ed.) (1959 and 1967) *Electrophoresis: Theory, Methods and Applications*, Vols. **1** and **2**, Academic Press, New York.
10. Z. Deyl (ed.) (1979 and 1983) *Electrophoresis: A Survey of Techniques and Applications*, Elsevier, Amsterdam.
11. S. Hjerten (1978) "Analytical electrophoresis". In *Topics in Bioelectrochemistry and Bioenergetics*, Vol. **2** (G. Milazzo, ed.), Wiley, New York, pp. 89–128.
12. C. J. O. R. Morris and P. Morris (1976) *Separation Methods in Biochemistry*, 2nd ed., Wiley, New York, pp. 703–900.
13. V. Neuhoff, R. Stamm, I. Pardowitz, N. Arold, W. Ehrhardt, and D. Taube (1990) Essential problems in quantification of proteins following colloidal staining with Coomassie Brilliant Blue dyes in polyacrylamide gels, and their solution. *Electrophoresis* **11**, 101–117.
14. P. G. Righetti, C. J. Van Oss, and J. W. Vanderhoff (1979) *Electrokinetic Separation Methods*, Elsevier, Amsterdam.

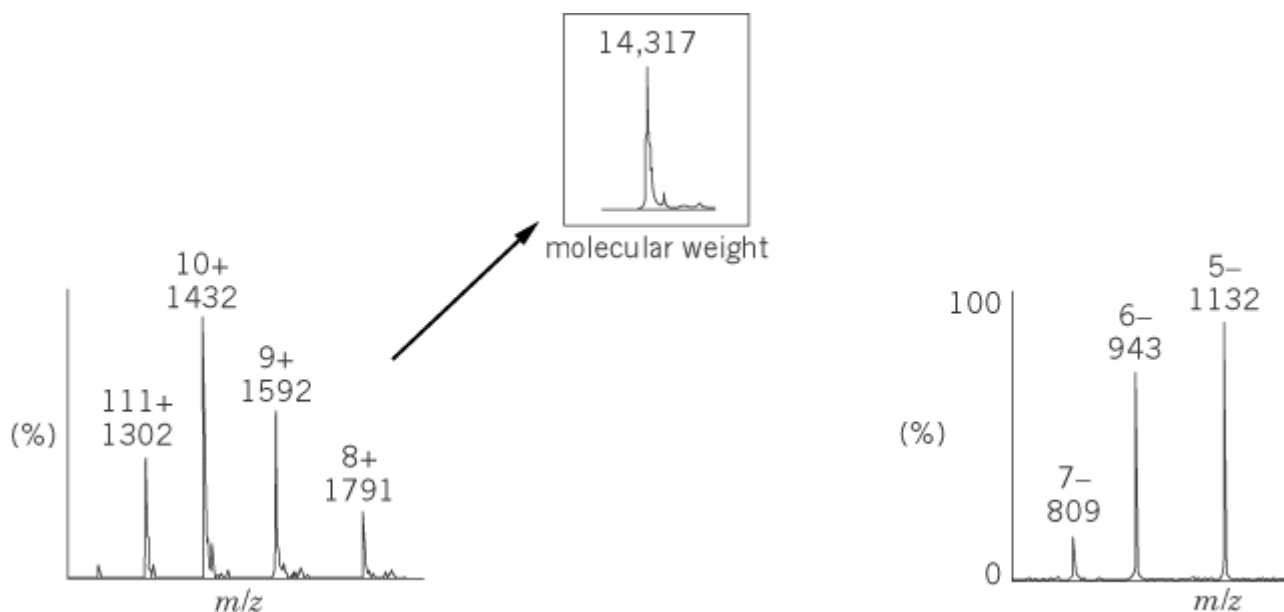
## Electrospray Ionization

*Electrospray ionization* (ESI) is a method of generating ionic forms of molecules for [mass spectrometry](#). It generates ions directly from solution (usually an aqueous or aqueous/organic solvent system) by creating a fine spray of highly charged droplets in the presence of a strong electric field (1), as shown in Fig. 1. Subsequent vaporization of these charged droplets results in the production of multiply charged gaseous ions. The number of charges retained by an analyte can depend on such factors as the composition and pH of the electro sprayed solvent, as well as the chemical nature of the sample (2-4). For large molecules, the ESI process typically gives rise to a series of multiply charged species for a given analyte. Because mass spectrometers measure the mass-to-charge ( $m/z$ ) ratio, the resultant ESI mass spectrum contains multiple peaks corresponding to the different charged states (Fig. 2).

**Figure 1.** The electrospray ionization process.



**Figure 2.** Examples of data generated on an ESI mass spectrometer. Proteins (left) typically produce positive multiply charged ions. The insets for each part show the molecular-weight spectra computer-generated



The extent of multiple charging that occurs in ESI is a unique characteristic of the technique that enables an analyte's mass to be determined with great precision, as masses can be independently calculated from several different charged states (5). The multiple charging in ESI also permits the analysis of high-molecular-weight analytes, using conventional mass analyzers that are normally limited to the detection of ions with relatively low  $m/z$  ratios. For example, a 50-kDa protein will typically retain on the order of 30 to 50 charges in ESI, yielding multiply charged species with  $m/z$  ratios between 1000 and 2000 that are easily detected with *quadrupole mass analyzers*. Another advantage of ESI-MS is its compatibility as an interface with liquid chromatography (see [Liquid Chromatography Mass Spectrometry](#)). Electrospray ionization is compared to the other common method of ionization in [Matrix-Assisted Laser Desorption/Ionization](#).

#### Bibliography

1. P. Kebarle and L. Tång (1993) *Anal. Chem.* **65**, A972–A986.
2. R. D. Smith, J. A. Loo, C. G. Edmonds, C. J. Barinaga, and H. R. Udseth (1990) *Anal. Chem.* **62**, 882–899.
3. B. T. Chait and S. B. H. Kent (1992) *Science* **257**, 1885–1894.
4. D. Arnott, J. Shabanowitz, and D. F. Hunt (1993) *Clinical Chemistry* **39**, 2005–10.
5. G. Siuzdak (1996) *Mass Spectrometry for Biotechnology*, Academic Press, San Diego.

#### Electrostatic Interactions

Electrostatic forces originate from the electric charges and dipoles of atoms and groups of atoms. They function inside macromolecules to maintain and to stabilize the molecular structure and to regulate the biological functions of catalysis and electron transfer. Electrostatic forces between molecules, facilitate specific molecular recognition and molecular assembly. The specific electrostatic field is determined by contributions from several factors, which originate from the intrinsic heterogeneity of biological macromolecules and from aqueous solvents.

## 1. Ionic Charges

In proteins, [aspartic acid](#) and [glutamic acid](#) residues have negatively ionizable side chains, and [lysine](#), [arginine](#), and [histidine](#) are positively ionizable. The amide and carboxyl termini of the polypeptide chain can also ionize, with, respectively, positive and negative charges. The intrinsic  $pK_a$  values of these groups in a normal aqueous environment are listed in Table 1. The actual  $pK_a$  value of a group in a folded protein can deviate from its intrinsic value, depending on the electrostatic field it experiences. The net charge of a protein is the sum of the positively and negatively ionized amino acid groups, plus any cofactor ions, and depends on the pH of the solvent. The pH at which a protein has no net charge is known as its [isoelectric point](#). At neutral pH, most proteins are [amphoteric](#), having both cationic and anionic groups at specific locations in the folded conformation.

**Table 1.  $pK_a$  Values for Amino Acids in Peptides<sup>a</sup> and Bases in Nucleosides<sup>b</sup>**

|                          | $pK_a$ |                            | $pK_a$ |
|--------------------------|--------|----------------------------|--------|
| Amino Acids in Peptides  |        | Bases in Nucleosides       |        |
| a-NH <sub>2</sub>        | 7.8    | Adenine (N <sub>1</sub> )  | 3.52   |
| a-COOH                   | 3.6    | Cytosine (N <sub>3</sub> ) | 4.17   |
| Asp (C <sub>g</sub> OOH) | 4.0    | Guanine (N <sub>1</sub> )  | 9.42   |
| Glu (C <sub>d</sub> OOH) | 4.5    | Guanine (N <sub>7</sub> )  | 3.3    |
| His (imidazole)          | 6.4    | Uracil (N <sub>3</sub> )   | 9.38   |
| Lys (ε-NH <sub>2</sub> ) | 10.4   | Thymine (N <sub>3</sub> )  | 9.93   |
| Arg (guanidine)          | ≈12    |                            |        |
| Tyr (OnH)                | 9.7    |                            |        |
| Cys (SgH)                | 9.1    |                            |        |

<sup>a</sup> Ref. 23.

<sup>b</sup> Ref. 24.

In **nucleic acids**, only the negative phosphate groups are charged at neutral pH, so nucleic acids do not show amphoteric properties. At very acidic or alkaline pH values, however, protonation or deprotonation of bases occurs, depending upon their  $pK_a$  values (Table 1).

## 2. Permanent Electric Dipoles

A molecule or a group of atoms has a permanent electric dipole if it has an unequal electron distribution. This dipole spontaneously produces its own electric field, and the molecule or group is then described as *polar*. One of the more important polar groups in macromolecules is the [peptide bond](#), which has a large dipole moment. The polar groups interact electrostatically with other polar

groups in the same or other molecules or with polar solvent molecules or ions. Permanent dipoles within most macromolecules are small and are distributed in different orientations, except for the [alpha-helix](#). It possesses a large permanent dipole moment due to the cooperatively aligned electric dipoles of the peptide bonds of the backbone (1)

### 3. Polarizability and Dielectric Constant

In an external electric field  $\mathbf{E}$ , an electric dipole is induced even in a neutral, nonpolar molecule or group because the outer field modifies the electrostatic balance between the electrons and the nucleus in each atom. The magnitude of this dipole is given by  $m^{\text{induced}} = a_0 \mathbf{E}$ , where the coefficient  $a_0$ , the electronic polarizability, is proportional to the volume of the molecule or group. When the molecule or group also has additional permanent dipoles, they also reorient because of the external field. Consequently, the polarizability of a molecule is the sum of its electronic and orientational polarizabilities. The contribution from the orientational polarization of a polar group is generally much larger than that of the electronic polarization (2).

When a material is physicochemically homogeneous in a macroscopic volume, the dielectric constant of the material is related to the polarizability through the *Clausius–Mossotti* equation. Thus, the dielectric constant reflects the degree of polarizability of the material. A pure, polar liquid has a much larger dielectric constant than a nonpolar liquid. For example, the dielectric constants at 20 °C of water and cyclohexane are 80 and 2, respectively. The magnitude of the electrostatic interactions between charged groups or dipoles is relatively straightforward in such a homogeneous medium. It is inversely proportional to the distance between each pair of charges and to the dielectric constant of the medium. In contrast, proteins and nucleic acids are so heterogeneous, with polar and nonpolar groups distributed nonuniformly, that it is impossible to define an dielectric constant. Roughly speaking, however, the dielectric constant should be small in a **hydrophobic** core, perhaps 2 to 4, whereas the surface regions have larger values of 10 or larger (3).

### 4. Electrostatic Shielding by Polar Solvents

When a charged, polar biological macromolecule is dissolved in a polar aqueous solvent, the solvent molecules rotate in response to their interaction with the solute macromolecule in an attempt to decrease the **free energy** of the entire system. Thus, the electrostatic force produced by the solute macromolecule is significantly shielded. This free energy is called the *self-energy*, the interaction of a solute charge or electric dipole with the dielectric environment (2). In addition, the density distribution of small ions in the solvent also changes to shield the electric field originating from the solute molecule. Small metal cations are located near the negative charges, and anions are near the positive point charges, as described by the *Debye–Hückel* approximation. Double-stranded DNA or RNA has a very high negative charge density of phosphate groups of the backbone, which is almost neutralized (about 88%) by condensed counterions, which are strongly restricted along the duplex, in addition to the Debye–Hückel type interaction. The density of the condensed counterions is relatively independent of the solvent conditions, but they are released from the polymer when the double helical conformation is deformed (4).

When the response to a given electric field is monitored directly, changes in the orientation and density distribution of the solvent molecules are the major responses to an external field. Therefore, passive measurement of the individual  $pK_a$  values of the ionizable groups (5) and of the [NMR](#) chemical shift values of the nuclei (6) are good reporters of the local electrostatic field inside macromolecules. Both measurements are possible with the aid of multidimensional NMR techniques using stable isotope labeling (7). The charges and polarities of interest in proteins can be changed by [site-directed mutagenesis](#), and the free energy difference associated with the electrostatic interaction can be evaluated (8).

### 5. Theoretical Approaches

Theoretical approaches to electrostatic interactions are divided between two types of theories: the macroscopic or continuum and the microscopic or molecular. In the continuum theory, the *Poisson* or *Poisson–Boltzmann* equation is solved for the solute–solvent system, assuming that the system is large enough to define its dielectric properties. Kirkwood first obtained an analytical solution for a simple spherical solute molecule embedded in a solvent with a high dielectric constant (9). With the aid of high-speed computational techniques, the exact shapes of biological macromolecules can be considered and the equations solved numerically (10). The computer program GRASP (11) is widely used to compute the electrostatic potential based on this method and to illustrate the electrostatic molecular surfaces of biological macromolecules and their assemblies, using red and blue colors for negative and positive charges, respectively. The drawbacks of the continuum theory (3) are (1) that the solute–solvent system is too heterogeneous to introduce the dielectric constant parameters and (2) that the calculations are based on static molecular conformations. Therefore, strict agreement between experiments and calculations is not expected. In the molecular theory, all of the atoms of the solute and solvent molecules are incorporated into the calculation, using the vacuum dielectric constant. Because the electric force is essentially long-range and the solvent molecules are significantly polar, reliable results are obtained only if the dynamics of the system are simulated for a sufficiently long period without simply cutting off the electrostatic forces.

## 6. Roles in Folded Conformations

Statistical analyses of the charge distributions on protein surfaces indicate that the charged groups are, on average, surrounded by charges of opposite sign and they tend to form [ion pairs](#), which contribute to conformational stability (12, 13). The electrostatic contribution of ion pairs to protein stability is context-dependent, and the interaction energy of a single ion pair exposed to the solvent or even buried in the protein interior, varies experimentally from 0.5 to 5 kcal/mol. The entropic cost of forming rigid [salt bridges](#) is an opposing factor that reduces the free energy of the folded protein (14). A cooperative salt-bridge network can contribute significantly to protein stability (15), as observed frequently on the surfaces of proteins of thermophilic origins. Acidic residues occur near the positive pole of the helical dipole at the N-terminus of an  $\alpha$ -helix (12) and basic residues near the negative pole. They represent the helix caps (16). These interactions are not context-dependent, but only the first one or two turns of an  $\alpha$ -helix give the essential contribution to protein stability by at most about 2 kcal/mol, independent of the length of the  $\alpha$ -helix (17).

Unfavorable electrostatic interactions destabilize the folded conformation and are important in proteins. Unpaired buried charges and isolated hydrogen bonds are energetically expensive within the hydrophobic core because of the large self-energy. Thus, **secondary structures**, such as  $\alpha$ -helices and [beta-sheets](#), can usually be depicted in protein cores as simply the backbone structures because they do not leave any backbone hydrogen bond donors or acceptors unpaired (18).

The condensation of counterions contributes to the conformation and the stability of double helices of DNA and RNA, which are highly charged due to the ionized phosphate groups of the backbone. The ordered double-stranded conformation depends markedly on the ionic strength, and the melting temperature of the conformation is increased as ionic strength increases. Divalent cations interact with these duplexes more strongly than monovalent cations.

## 7. Enzyme Catalysis

Efficient [enzyme](#) catalysis is based on the decrease in the free energy of the [transition state](#), and it is mediated by a well-designed active site structure and the electrostatic field. [Protein engineering](#) studies employing site-directed mutagenesis have highlighted the importance of ionizable or polar residues at the active site. In the catalytic reactions of proteins, the [thiol group](#) of **cysteine residues**, the hydroxyl group of [serine](#) or the imidazole group of [histidine](#) acts as a nucleophilic catalyst at the initial stage. Ionizable side chains can act as acid and base catalysts by giving and receiving a proton,



respectively. Generally, when two ionizable groups are adjacent their  $pK_a$  values deviate greatly from the intrinsic values. One value becomes higher and the other becomes lower because of the strong coupling interaction between the two ionization states. Thus, even at a neutral pH, two such ionizable side chains can form acid and/or base catalysts, which are often observed to have unusual  $pK_a$  values at the active sites of enzymes. In addition, ligated metal ions play important roles in many examples of catalysis by correctly positioning the substrate molecules, by stabilizing the transition-state conformation, or by activating water molecules to form hydroxide ions.

The **electron transfer** reaction is an essential event in energy acquisition during respiration and [photosynthesis](#). Several prosthetic groups in proteins, such as hemes, iron–sulfur clusters, and some transition metals, such as Cu, directly govern the reaction. The electron transfer reaction is finely controlled by the surrounding peptide chains of proteins through electrostatic interactions (19). In addition, because the prosthetic groups are often buried deeply in the protein, electrons are considered to transfer from the protein surface through certain, but not unique, paths in the protein molecule (20). Electrostatic interactions often trigger large-scale structural changes, especially when coupled to ATP and GTP hydrolysis. The ionic and polar interactions of proteins with the phosphate groups of ATP and GTP become reordered after their hydrolysis. This local phenomenon causes domain movement in multidomain proteins and results in large structural changes (see [ATPase](#) and [Gtpases](#)). Likewise, several **allosteric** transitions are considered to be mediated by **salt bridges**.

## 8. Molecular Recognition

Molecular recognition is a major focus of modern biology and biochemistry. Specific molecular recognition of individual biological macromolecules is involved in many biological phenomena. Protein–protein interactions are critical in biological [signal transduction](#), and they guide the assembly of multimeric proteins, generally found in cell **organelles**, **cytosol**, and cell [membranes](#). Protein molecular surfaces are involved in the recognition of other molecules. [Hydrogen bonds](#) and salt bridges often provide specificity, which is derived from the distinctive electrostatic complementarity between the molecular surfaces. At protein–protein interfaces, on average, 0.88 hydrogen bonds are found per 100 Å<sup>2</sup> of buried [accessible surface](#) area, and from 0 to 5 intersubunit salt bridges are observed between protein dimers (21).

Protein–nucleic acid interactions are the origin of complex genetic regulation. The backbone phosphate groups interact with the basic side chains of the protein, and specific base sequences are generally recognized through hydrogen bonds and hydrophobic interactions with the protein side chains (22). In the process of association, a basic protein approaches the DNA double helix by a long-range, electrostatic attractive force. Then the protein condenses, like a group of counterions, along the DNA duplex. Some of the condensed counterions are released. Then the protein slides along the double helix by one-dimensional diffusion, scanning the sequence, until it finally becomes positioned at the specific binding site. Specific basic residues in the protein molecule form ion pairs with the exposed phosphate groups of the DNA backbone. Thus, the local DNA backbone structure, which is completely neutralized, is deformed. In fact, the DNA structures bound by proteins are often deformed from the typical B-type conformation.

## Bibliography

1. A. Wada (1976) *Adv. Biophys.* **9**, 1–63.
2. J. N. Israelachvili (1985) *Intermolecular and Surface Forces*, Academic Press, London, Chaps. "4" and "5", pp. 36–64.
3. H. Nakamura, T. Sakamoto, and A. Wada (1988) *Protein Eng.* **2**, 177–183.
4. M. T. Record Jr., J. H. Ha, and M. A. Fisher (1991) *Methods Enzymol.* **208**, 291–343.
5. C. Tanford and R. Roxby (1972) *Biochemistry* **11**, 2192–2198.
6. A. C. de Dios, J. G. Pearson, and E. Oldfield (1993) *Science* **260**, 1491–1496.

7. Y. Oda et al. (1994) *Biochemistry* **33**, 5275–5284.
8. L. Serrano et al. (1990) *Biochemistry* **29**, 9434–9352.
9. J. G. Kirkwood (1934) *J. Chem. Phys.* **2**, 351–361.
10. J. Warwicker and H. C. Watson (1982) *J. Mol. Biol.* **157**, 671–679.
11. A. Nicholls, K. A. Sharp, and B. Honig (1991) *Proteins* **11**, 281–296.
12. A. Wada and H. Nakamura (1981) *Nature* **293**, 757–758.
13. D. J. Barlow and J. M. Thornton (1983) *J. Mol. Biol.* **168**, 867–885.
14. S. Dao-pin, U. Sauer, H. Nicholson, and B. W. Matthews (1991) *Biochemistry* **30**, 7142–7153.
15. A. Horovitz et al. (1990) *J. Mol. Biol.* **216**, 1031–1044.
16. J. S. Richardson and D. C. Richardson (1988) *Science* **240**, 1648–1652.
17. J. Åqvist, H. Luecke, F. A. Quioco, and A. Warshel (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2026–2030.
18. A. V. Finkelstein and O. B. Ptitsyn (1987) *Prog. Biophys. Molec. Biol.* **50**, 171–190.
19. T. Ichiye (1996) *Structure* **4**, 1009–1012.
20. A. Kuki and P. G. Wolynes (1987) *Science* **236**, 1647–1652.
21. S. Jones and J. M. Thornton (1995) *Prog. Biophys. Molec. Biol.* **63**, 31–65.
22. C. O. Pabo and R. T. Sauer (1992) *Ann. Rev. Biochem.* **61**, 1053–1095.
23. A. R. Fersht (1985) *Enzyme Structure and Mechanism*, 2nd ed., Freeman, New York, pp. 156, Table 5.1.
24. G. M. Blackburn (1996) in *Nucleic Acids in Chemistry and Biology: DNA and RNA structure*, 2nd ed. (G. M. Blackburn and M. J. Gait, eds.), Oxford University Press, Oxford, Chap. "2", p. 22, Table 2.2.

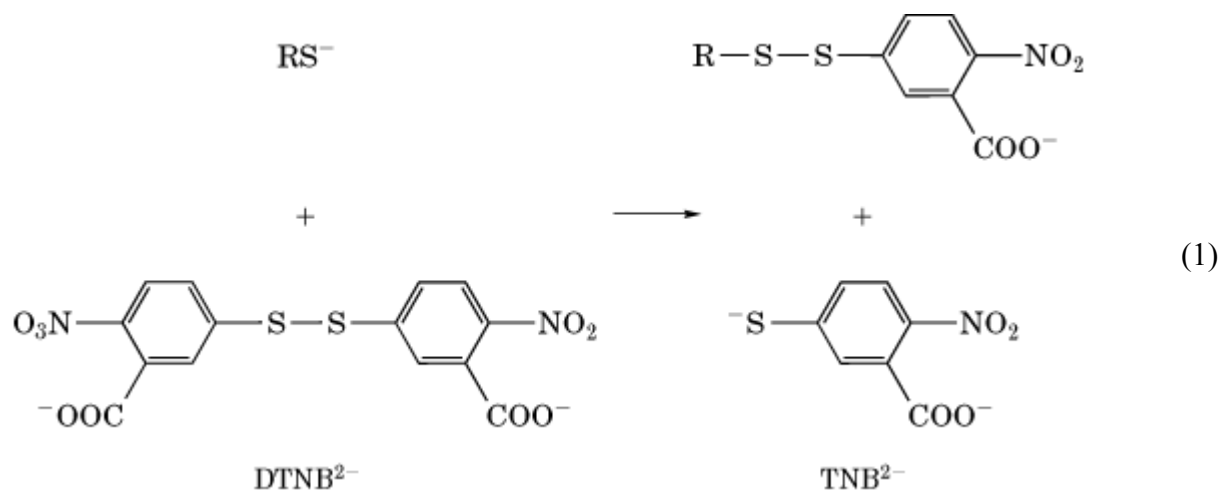
### Suggestions for Further Reading

25. B. Honig and A. Nicholls (1995) Classical electrostatics in biology and chemistry, *Science* **268**, 1144–1149. Comprehensive review of the approach of the continuum theory to electrostatic interaction in proteins and DNA.
26. H. Nakamura (1996) Roles of electrostatic interaction in proteins, *Quart. Rev. Biophys.* **29**, 1–90. Review of theoretical and experimental analyses of electrostatic interaction in proteins. Using simple models, the electrostatic energy values revealed experimentally were physically interpreted.
27. M. F. Perutz (1978) Electrostatic effects in proteins, *Science* **201**, 1187–1191. Biological importance of electrostatic interaction in proteins was reviewed in several different aspects.
28. A. Warshel and S. T. Russell (1984) Calculation of electrostatic interactions in biological systems and in solutions, *Quart. Rev. Biophys.* **17**, 283–422. Review of the molecular theory contrasted with the empirical continuum theory.
29. W. D. Wilson (1996) "Reversible interactions of nucleic acids with small molecules". In *Nucleic acids in chemistry and biology*, 2nd ed. (G. M. Blackburn and M. J. Gait, eds.), Oxford University Press, Oxford, Chap. "8", pp. 331–374. Electrostatic features of nucleic acids were well introduced.

### Ellman's Reagent, 5,5'-Dithiobis-(2-Nitrobenzoic Acid), DTNB

Barrett and Seligman (1) in 1952 used 2,2'-dihydroxy-6,6'-dinaphthyl disulfide to demonstrate [thiol groups](#) in [proteins](#) histochemically. Six years later, Ellman reported the use of bis-(*p*-nitrophenyl) disulfide for spectrophotometric measurement of thiol groups, and subsequently enhanced its water solubility by introducing carboxyl groups into the benzene rings (2). This reagent, 5,5'-dithiobis-(2-nitrobenzoic acid), DTNB, is now known as Ellman's reagent and is one of the most useful reagents for measuring thiol groups quantitatively.

The cysteine thiol groups undergo [thiol–disulfide exchange](#) with the [disulfide bond](#) of DTNB at about pH 8, to produce a mixed disulfide between them, plus 5-thio-2-nitro-benzoic acid, TNB.



The TNB is intensely colored. Its absorption maximum is normally at 409.5 nm, but shifted to 421 nm in 6 M **guanidinium chloride** (GdmCl), which is often used to unfold proteins and make their thiol groups accessible (3). For historical reasons, however, the reaction is generally monitored at 412 nm, where the molar absorbance coefficient is 14,150 M<sup>-1</sup> cm<sup>-1</sup> in standard buffers and 13,700 in 6 M GdmCl. The reaction goes essentially to completion because an excess of Ellman's reagent is used, plus the pK<sub>a</sub> of the thiol group of TNB is approximately 4.5, considerably more acidic than normal thiol groups, which have pK<sub>a</sub> values near 9. This helps to drive the reaction to completion, as thiol–disulfide exchange favors the more acidic thiol group. Irrespective of whether the original thiol groups react only with the Ellman's reagent or also with the mixed disulfide, one mole of TNB is produced per mole of original thiol group.

The ease of use of the Ellman assay makes it widely used, and many variants of Ellman's reagent have also been devised for special purposes, such as kinetic studies at acidic pH (4).

#### Bibliography

1. J. R. Barrett and A. M. Seligman (1952) *Science* **116**, 323–327.
2. G. L. Ellman (1959) *Arch. Biochem. Biophys.* **82**, 70–77.
3. P. W. Riddles, R. L. Blakeley, and B. Zerner (1983) *Methods Enzymol.* **91**, 49–60.
4. K. Brocklehurst (1979) *Int. J. Biochem.* **10**, 259–274.

#### Elongation Factors (EFs)

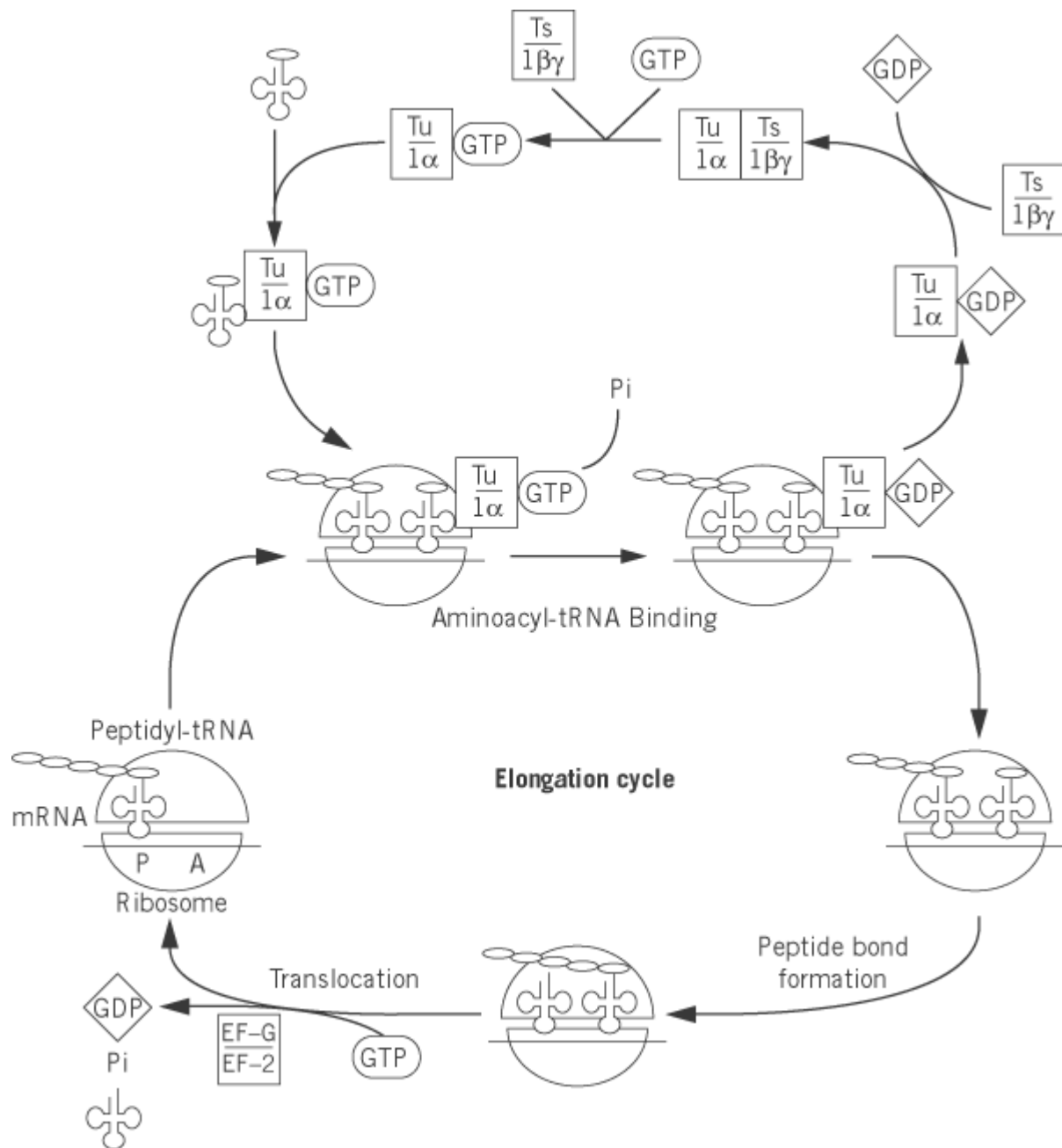
Elongation of the nascent [polypeptide chain](#) during [translation](#) of protein biosynthesis requires three elongation factors (EF-Tu, EF-Ts, and EF-G) in prokaryotes and requires four factors, (EF-1 a, EF-1 b, EF-1 g and EF-2) in eukaryotes (Table 1). The basic function of these elongation factors is to recycle the delivery of aminoacyl-[transfer RNA](#) molecules to the A site of the [ribosome](#) and to **translocate** the ribosome one **codon** forward after the **peptidyl transfer** reaction. Because of this common scenario, elongation factors share high functional and structural similarities in prokaryotes and eukaryotes: EF-Tu is equivalent to EF-1 a, EF-Ts to EF-1 bg, and EF-G to EF-2 (1).

**Table 1. *Escherichia coli* Elongation Factors**

| Factor | Molecular Mass | Number of Amino Acids | Gene        | Function                              |
|--------|----------------|-----------------------|-------------|---------------------------------------|
| EF-Tu  | 43,195         | 393                   | <i>tufA</i> | Bind aminoacyl-tRNA to A site; GTPase |
|        | 43,225         | 393                   | <i>tufB</i> | Same as above                         |
| EF-Ts  | 30,257         | 282                   | <i>tsf</i>  | GDP/GTP exchange on EF-Tu             |
| EF-G   | 77,444         | 703                   | <i>fus</i>  | Translocation; GTPase                 |

The polypeptide elongation cycle can be divided into three steps: (i) delivery of aminoacyl-tRNA to the ribosome, (ii) peptidyl transfer, and (iii) translocation of the ribosome (Fig. 1). Aminoacyl-tRNA first forms a ternary complex with EF-Tu (EF-1 a) and GTP. EF-Tu is one of the most abundant proteins. In *Escherichia coli*, it comprises about 5% of the soluble protein, and the EF-Tu/ribosome ratio is up to 7:1 in rapidly growing cells. The aminoacyl-tRNA:EF-Tu:GTP ternary complex binds to the A site of the ribosome. The proper binding of EF-Tu to the ribosome and the correct base-pairing between codon and the [anticodon](#) of the tRNA trigger hydrolysis of EF-Tu-bound GTP to GDP through activation by the [GTPase](#) active center on the ribosome. Note that both EF-Tu and EF-G carry conserved G-domains for binding guanine nucleotides, but they need the common GTPase center for stimulating the hydrolysis of GTP (ribosome-dependent GTPase). Once GTP is hydrolyzed, the GDP form of EF-Tu no longer remains bound to the ribosome because of protein conformation is altered. Common to the functioning of G-proteins is that activation occurs by binding GTP, while deactivation is achieved by GTP hydrolysis. GTP and GDP act as **allosteric** effectors on their target proteins, thus altering the conformation (see [GTP-Binding Proteins](#)).

**Figure 1.** Elongation pathway of protein synthesis. Prokaryotic (upper) and eukaryotic (lower) elongation factors are represented in boxes.



EF-Tu alone is very unstable, but is readily stabilized by GTP, GDP, or EF-Ts. EF-Ts (EF-1b) is involved in the guanine nucleotide exchange of aEF-Tu (EF-1a). EF-Ts catalyzes conversion of the EF-Tu:GDP complex to EF-Tu:GTP, promoting the recycling EF-Tu for new ternary-complex formation.

After peptidyl transfer, the elongated peptidyl-tRNA is bound to the A site of the ribosome while the deacylated tRNA is bound to the P site. EF-G (EF-2) is a translocase protein that forwards peptidyl-tRNA from the A site to the P site on the ribosome and forwards the deacylated tRNA from the P site to the E site. The GTP form of EF-G is the active state and binds to the ribosome. Translocation of the ribosome is coordinated with the hydrolysis of GTP bound to EF-G. Presumably, GTP hydrolysis induces conformational constraints on the ribosome that lead to the movement of the ribosome. Then, the GDP form of EF-G dissociates from the ribosome.

EF-3 appears to be a unique translation factor that is found only in yeast and **fungi** (2). Together

with yeast or mammalian EF-1a and EF-2, it is required for protein synthesis with ribosomes, except for mammalian ribosomes. This has prompted the suggestion that a corresponding activity in higher eukaryotes may be part of the ribosome; alternatively, EF-3 may be unique to yeast and fungi. EF-3 is known to promote the ejection of deacylated tRNA from the E site of the ribosome following translocation *in vitro*, thereby enhancing the binding of an aminoacyl-tRNA to the A site.

The three-dimensional (3D) structure of *Thermus thermophilus* EF-G comprises six subdomains, G, G', and II–V (Fig. 2, left). The C-terminal part, domains III–IV, appears to mimic the shape of tRNA (3-6). The 3D structure of the ternary complex of Phe-tRNA, *Thermus aquaticus* elongation factor EF-Tu, and the nonhydrolyzable GTP analog, GDPNP, is almost completely superimposable with EF-G:GDP [Fig. 2, right (5)], showing that domains III, IV, and V appear to mimic the shapes of the acceptor stem, anticodon helix, and T stem of tRNA in the ternary complex, respectively. Domain IV of EF-G forms a protruding “rod” conformation, which is similar to the shape of the anticodon arm of tRNA. These findings strongly suggest a common requirement for the structure and function on the ribosome. This resemblance of part of EF-G and tRNA represents a novel concept of “molecular mimicry between nucleic acid and protein.”

**Figure 2.** The structures of elongation factor EF-G:GDP (right) and EF-Tu:GDPNP:Phe-tRNA<sup>Phe</sup> (left). Domain II in EF-G and domain II in EF-Tu were placed in identical orientations, and the structures are shown side by side in a schematic representation to the same scale. (Reproduced from the diagram kindly generated by A. Ævarsson and A. Liljas.)



## Bibliography

1. W. C. Merrick and J. W. B. Hershey (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 31–69.
2. G. P. Belfield and M. F. Tuite (1993) *Mol. Microbiol.* **9**, 411–418.
3. A. Ævarsson et al. (1994) *EMBO J.* **13**, 3669–3677.
4. J. Czworkowski, J. Wang, T. A. Steitz, and P. B. Moore (1994) *EMBO J.* **13**, 3661–3668.

5. P. Nissen et al. (1995) *Science* **270**, 1464–1472.
6. K. Ito, K. Ebihara, M. Uno, and Y. Nakamura (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5443–5448.

### Suggestions for Further Reading

7. C. M. T. Spahn and K. H. Nierhaus (1998) *Biol. Chem.* **379**, 753–772.
8. C. G. Kurland, D. Hughes, and M. Ehrenberg (1996) In *Escherichia coli and Salmonella*, 2nd ed., (R. Curtiss III et al., eds.), American Society for Microbiology Press, Washington, D.C., pp. 979–1004.

## Embryo

An embryo is the juvenile stage of an animal while it is contained within the egg membranes or the mother's body. It is formed by [fertilization](#) of the egg and undergoes morphogenesis to form differentiated organs and tissues.

### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 3.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, U.K.

## Embryology

Embryology, the study of **embryos**, is the field that deals with how a single cell multiplies and changes during [development](#) to form a multicellular organism. This process by which an organism is formed is known as *embryogenesis*. Embryology can be divided into a number of smaller fields. Descriptive embryologists describe the actual events that occur during embryogenesis. Comparative embryologists study embryology in different organisms to understand how developmental mechanisms have arisen and changed during [evolution](#). Experimental embryologists manipulate embryos in the laboratory to uncover the cellular and biochemical processes, while the more recent field of molecular embryology studies the molecules that underlie early development.

A number of species have been used extensively to study the patterns and underlying mechanisms of embryogenesis. These include the [nematode](#) *Caenorhabditis elegans*, the fruit fly *Drosophila melanogaster*, several species of marine invertebrates, such as sea urchins, starfish, and marine annelids, the African clawed toad *Xenopus laevis*, the domestic chicken, and the house [mouse](#) *Mus domesticus*.

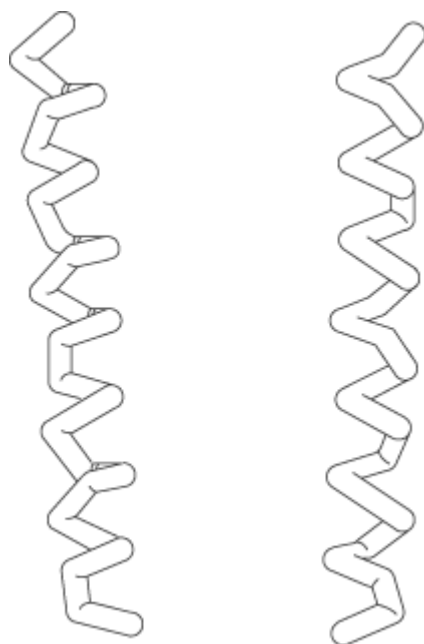
### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, U.K.

## Enantiomer

An enantiomer is an object that has a nonsuperimposable mirror image. A right-handed helix and its mirror image are enantiomers, as shown in Figure 1. Thus, a molecule and its nonsuperimposable mirror image are both enantiomers and are both necessarily **chiral** (1). Enantiomeric molecules differ in how their solutions rotate the plane of polarized light (+ and –, for dextro and levo rotations, respectively). Enantiomers interact identically with achiral compounds, but can have different interactions with other chiral objects. The enantiomers have different **configurations**. This configuration is designated by the formally accepted rules proposed by Cahn et al., (*R*) and (*S*) (2), as described under [Configuration](#). Because enantiomers differ in how their atoms are oriented in space, they are a special type of [stereoisomer](#). Most enantiomers are most readily recognized by determining if each corresponding chiral center has an opposite configuration.

**Figure 1.** Enantiomers of a helix. The left-handed helix cannot be superimposed on its mirror image, the right-handed helix. This property makes both helices chiral and enantiomers of each other.



Two identical functional groups are called enantiotopic if their substitution by an isotopic label generates enantiomers (3) (see [Prochiral](#)). Thus, the protons on the amino acid [glycine](#) are



enantiotopic because they are the same functional group, and alternate substitution of  $^2\text{H}$  for each  $^1\text{H}$  will generate the enantiomeric (*R*)- and (*S*)-[2- $^2\text{H}$ ]glycines.

### Bibliography

1. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), Vol. **3**, Elsevier, Amsterdam, pp. 1–48.
2. R. S. Cahn, C. K. Ingold, and V. Prelog (1966) *Angew. Chem. Int. Ed.* **5**, 385–415.
3. D. Arigoni, and E. L. Eliel (1969) In *Stereochemistry* (E. L. Eliel, and N. L. Allinger, eds.), Vol. **4**, Wiley-Interscience, New York, pp. 127–243.

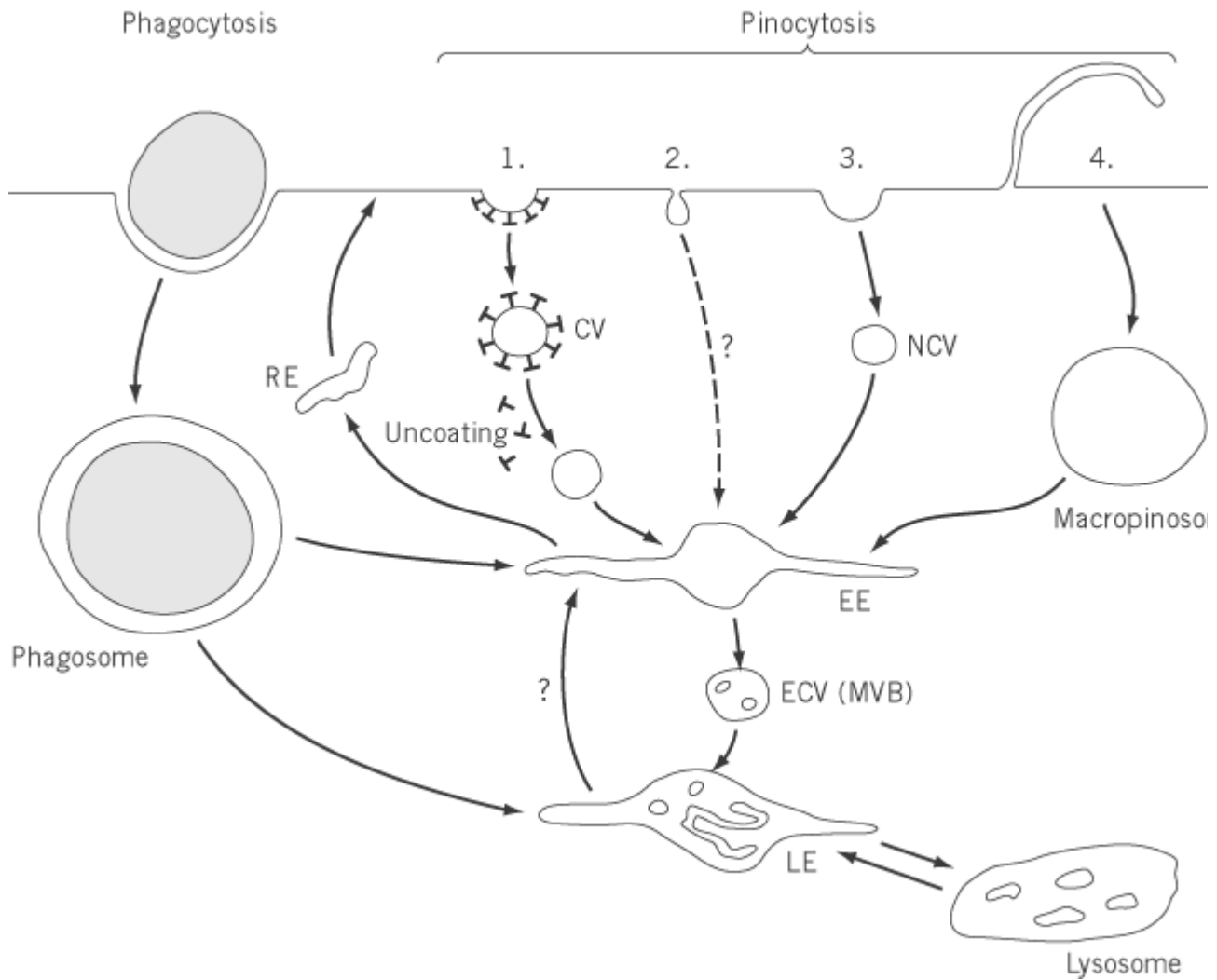
### Suggestions for Further Reading

4. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York.
5. K. Mislow (1966) *Introduction to Stereochemistry*, W. A. Benjamin, New York.
6. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), Vol. **3**, Elsevier, Amsterdam, pp. 1–48.

## Endocytosis

Endocytosis is the process whereby eukaryotic cells internalize material from their surrounding environment. Internalization is achieved by the formation of membrane-bound **vesicles** at the cell surface that arise by progressive invagination of the plasma [membrane](#), followed by pinching off (scission), and the release of free vesicles into the cytoplasm (Fig. 1). Several different types of vesicle are generated. Each carries extracellular fluid, **macromolecules**, and particles into the cell. Extracellular materials are either free within the lumen of the vesicle or bound to its luminal surface. Endocytosis is believed to be essential for cell viability and is required for numerous cellular functions, including (1) the provision of certain essential nutrients; (2) the removal of unwanted, nonfunctional or potentially harmful material from the extracellular medium; (3) the internalization, processing, and presentation of **antigens** by [major histocompatibility complex \(MHC\)](#) molecules (see [Antigen Processing, Presentation](#)); (4) the transfer of macromolecules across polarized epithelial cells; and (5) the transmission of neuronal, metabolic, and proliferative signals (see [Growth Factors](#); [Insulin](#)). In addition, endocytosis maintains cellular homeostasis by continually monitoring the plasma membrane composition and selectively removing redundant or defective **glycoproteins**, lipids, or glycolipids for repair or degradation. Endocytosis mediates the rapid recovery of membrane components inserted into the plasma membrane during secretion and is crucial for regenerating of synaptic vesicles (see [Exocytosis, Secretory Vesicles/Granules, Protein Secretion](#)). Endocytosis is also used by a range of viral [eg, [influenza virus \(1\)](#)], bacterial (eg, *Legionella pneumophila*), and protozoan (eg, *Toxoplasma gondii*) pathogens, and certain protein **toxins** produced by plants (eg, see **Ricin**) or bacteria (eg, see [Diphtheria toxin; Cholera toxin](#)), to enter target cells (2).

**Figure 1.** Schematic illustration of endocytic pathways. Four different mechanisms for pinocytosis are illustrated: (1) *via* clathrin-coated pits and vesicles (CV); (2) *via* caveolae; (3) *via* noncoated pits and vesicles (NCV); and (4) *via* macropinosomes. The abbreviations used are EE, early endosome; LE, late endosome; RE, recycling endosome; ECV, endosomal carrier vesicle; and MVB, multivesicular body.



## 1. The Endocytic Pathway

The endocytic pathway encompasses all of the organelles in the cell that are involved in endocytosis. The pathway starts at the cell surface, which provides the membrane that forms endocytic vesicles and contains the receptors responsible for binding specific ligands internalized by these vesicles. Endocytic vesicles [phagosomes, **clathrin**-coated vesicles, caveolae, **macropinosomes**, and nonclathrin vesicles (see later)] derive from the plasma membrane and are directed to fuse with and deliver their membrane and content to the intracellular organelles of the endocytic pathway, the **endosomes** and **lysosomes**.

Lysosomes are usually regarded as the graveyard of the cell, the site where ligands internalized from the cell surface and destined for degradation are delivered and hydrolyzed. Although many receptors recycle to the cell surface, others, such as growth factor receptors, remain associated with their ligands and are delivered to lysosomes for degradation. Lysosomes are usually vesicular (~0.5 to 1  $\mu\text{m}$  in diameter) and are located primarily, although not exclusively, in the perinuclear region of the cell. The pH within lysosomes can be as low as 4.8, the optimum for the many hydrolytic enzymes they contain. They derive their content from late endosomes, probably through repeated cycles of fusion with and fission from these organelles (3). In antigen-presenting cells, late endosomes and lysosomes contain large amounts of MHC II antigens and have been termed MIICs (4). In addition, lysosomes in **haematopoietic** cells are secretory organelles that undergo regulated fusion with the plasma membrane (5).

The membrane proteins, lipids, and hydrolases that make up the membrane and content of endocytic organelles are synthesized in the endoplasmic reticulum (ER) and Golgi apparatus and are subsequently transported to endosomes and lysosomes by vesicular transport. Many hydrolases carry the mannose 6-phosphate recognition marker that enables delivery *via* the [mannose 6-phosphate receptors](#). By contrast, the [membrane proteins](#) of endosomes and lysosomes carry sorting signals (see later) that facilitate their delivery.

The endocytic pathway is extremely dynamic. In some tissue culture cells, approximately 1500 coated vesicles internalize from the cell surface every minute. This is equivalent to internalizing the membrane of the entire cell surface every 30 to 60 min (6). The surface area of early endosomes is approximately one-fifth that of the cell surface, and the entire surface area of these organelles turns over approximately every 12 min (7). Receptors, such as the **low-density lipoprotein** (LDL) receptor, are internalized very rapidly at rates of 10 to 30% of the cell surface pool per min and are recycled in 10 to 15 min back to the cell surface. Each individual receptor may mediate many rounds of internalization and recycling in its life time (8). The organization of endocytic organelles within the cell is regulated by interaction with **cytoskeletal** elements. **Actin**-based systems have been implicated in the formation of endocytic vesicles from the plasma membrane (see later), and actin/myosin and the microtubule system have been implicated in moving and positioning endocytic organelles within the cell.

## 2. Endocytic Mechanisms

Classically, endocytosis has been divided into [phagocytosis](#) (“cellular eating”) and [pinocytosis](#) (“cellular drinking”) (Fig. 1). Phagocytosis describes the internalization of large particles (>0.2  $\mu\text{M}$  in diameter) following particle binding to specific plasma membrane receptors. Therefore, phagocytosis is receptor-mediated and depends on the presence of appropriate ligands. Pinocytosis is the formation of generally smaller vesicles (50 to 150 nm diameter) that transport extracellular fluid (fluid-phase endocytosis) and macromolecules specifically (receptor-mediated endocytosis) or nonspecifically (adsorptive endocytosis) bound to the plasma membrane. These vesicles are often formed constitutively from **clathrin**-coated pits independently of the presence of specific ligands, but there are alternatives to clathrin-mediated uptake (see later). The different endocytic mechanisms, more than one of which coexist and function in a single cell, are classified as [phagocytosis](#), [pinocytosis](#), **macropinocytosis**, and noncoated vesicles.

### 2.1. Noncoated vesicles

Endocytosis still occurs when clathrin-mediated uptake is blocked, as evidenced by the internalization of **ricin**, **interleukin-2**, and fluid-phase markers. Morphologically, the vesicles thought to mediate this uptake noncoated, ~80 nm in diameter, and distinct from caveolae and macropinosomes. At present, nothing is known about the molecular mechanisms responsible for this uptake process, although the continued uptake of ricin in the presence of a dominant-negative mutant form of **dynamain** suggests that it may be dynamain-independent (9).

## 3. Endocytosis Signals

Receptors and other proteins that undergo efficient endocytosis from the cell surface contain endocytosis signals. The signals associated with proteins that internalize through clathrin-coated pits are those that are best characterized. Four classes of endocytosis signal have been identified and studied in some detail. First are [tyrosine](#) residues within motifs similar to Phe-X-Asn-Pro-X-Tyr or Tyr-X-X-O (where X = any amino acid and O = a large **hydrophobic** amino acid), which are the endocytosis motifs identified in the LDL receptor and [transferrin receptor](#), respectively (10). It has been proposed that these motifs form [turn](#) structures and interact either directly with clathrin (11) or with the m2 subunit of the AP-2 complex [(10) see [Pinocytosis](#)], respectively. Generally these signals are constitutively active, and proteins that contain them undergo continuous endocytosis and recycling.

A second group of signals requires pairs of hydrophobic amino acids, frequently [leucine](#) residues, although [isoleucine](#), [methionine](#), or [valine](#) substitutes for one of the leucines in some circumstances [e.g., MHC class II invariant chain (12)]. These so-called dileucine signals are less well characterized. There is evidence, however, that they also bind AP-2 complexes through a site or sites distinct from those bound by the Tyr-based signals. There is also evidence that in some cases (e.g., CD4) the activity of these signals is regulated by [phosphorylation](#) of adjacent [serine](#) residues (13). The signal is active when the motif is phosphorylated and inactive when dephosphorylated.

The third class of internalization signal has been identified in members of the family of seven transmembrane-domain **heterotrimeric G protein**-coupled receptors that are internalized following ligand binding. For several of these receptors, in particular the  $\beta_2$ -adrenergic receptor, ligand-induced phosphorylation of serine residues in the serine-rich C-terminal domain of the molecule leads to recruitment of  $\beta$ -arrestins that uncouple associated heterotrimeric G proteins and function as adaptor complexes to recruit the receptor into clathrin-coated pits (14). Finally, in *S. cerevisiae*, a factor-induced phosphorylation of serine residues in the Ste2p receptor protein C-terminal domain leads to **ubiquitination** of Ste2p and its subsequent internalization (15). A number of mammalian cell surface receptors, including the growth hormone receptor, are also ubiquitinated following ligand binding (16). Whether this modification is required for endocytosis of the receptor and how the modification facilitates interaction with the endocytosis machinery are unclear.

Sorting signals also operate within endosomes and other sites in the endocytic pathway to specify the routing of specific proteins. The signals themselves are not well characterized. Recycling to the plasma membrane is the default pathway and may not require specific signals. However, the Tyr-X-X-O and dileucine motifs discussed previously function to target internalized proteins from early endosomes to lysosomes (12), and receptor cross-linking by multivalent ligands or antibody complexes directs sorting from early endosomes to lysosomes (17). The low pH of endocytic compartments may also influence sorting by inducing dissociation of pH-sensitive ligand-receptor interactions. Other signals, for example, a diaromatic amino acid-containing motif in the cytoplasmic domain of the cation-dependent [mannose 6-phosphate receptor](#), function as a lysosome-avoidance motif to prevent proteins from being transported to lysosomes. How these sorting signals are interpreted is unclear for the most part, although coat proteins and adaptors have been implicated in certain steps (19). A subset of the COPI coatomer complex has been found associated with early endosomes (20), as have clathrin and AP adaptor complexes.

#### 4. Medical Significance

The endocytic pathway is crucial for many cellular functions. Defects in functions associated with the endocytic pathway have been linked to a number of clinical conditions. For example, failure to internalize serum LDL is associated with atherosclerosis, and defects in producing or effectively targeting lysosomal hydrolases results in a range of inherited lysosomal storage diseases. It is believed that the problems associated with Chediak-Higashi syndrome result from defects in sorting associated with late endosomes and lysosomes. Chromosomal translocations of genes that encode proteins involved in coated vesicle formation have been found in a number of lymphomas and leukemias. In addition, many bacterial, protozoan, and viral pathogens enter cells and establish infection through the endocytic pathway. It is likely that many other conditions will be linked to endocytic functions. However, the endocytic pathway is likely to provide one of the most effective routes to target and deliver drugs, DNA, or other therapeutic agents to cells.

#### Bibliography

1. M. Marsh and A. Helenius (1989) *Adv. Virus Res.* **36**, 107–151.
2. J. M. Lord and L. M. Roberts (1998) *J. Cell Biol.* **140**, 733–736.
3. B. M. Mullock, N. A. Bright, C. W. Fearon, S. R. Gray, and J. P. Luzio (1998) *J. Cell Biol.* **140**, 591–601.
4. I. Mellman, P. Pierre, and S. Amigorena (1995) *Curr. Opin. Cell Biol.* **7**, 564–572.

5. G. M. Griffiths (1997) *Semin. Immunol.* **9**, 109–15.
6. M. Marsh and A. Helenius (1980) *J. Mol. Biol.* **142**, 439–454.
7. G. Griffiths, R. Back, and M. Marsh (1989) *J. Cell Biol.* **109**, 2703–2720.
8. J. L. Goldstein, M. S. Brown, R. G. W. Anderson, and D. W. Russell (1985) *Annu. Rev. Cell Biol.* **1**, 1–39.
9. C. Lamaze and S. L. Schmid, (1995) *Curr. Opin. Cell Biol.* **7**, 573–580.
10. T. Kirchhausen, J. S. Bonifacino, and H. Riezman (1997) *Curr. Opin. Cell Biol.* **9**, 488–495.
11. R. G. Kibbey, J. Rizo, L. M. Gierasch, and R. G. W. Anderson (1998) *J. Cell Biol.* **142**, 59–67.
12. I. Sandoval and O. Bakke (1994) *Trends Cell Biol.* **4**, 292–297.
13. M. Marsh and A. Pelchen-Matthews (1996) In *Current Topics in Microbiology and Immunology*, Vol. **205** (D. R. Littman, ed.), Springer pp. 107–135.
14. S. S. Ferguson, L. S. Barak, J. Zhang, and M. G. Caron (1996) *Can. J. Physiol. Pharmacol.* **74**, 1095–1110.
15. L. Hicke and H. Riezman (1996) *Cell* **84**, 277–287.
16. R. Govers, P. van Kerkhof, A. L. Schwartz, and G. J. Strous (1997) *EMBO J.* **16**, 4851–4858.
17. P. Ukkonen, V. Lewis, M. Marsh, A. Helenius, and I. Mellman (1986) *J. Exp. Med.* **163**, 952–971.
18. A. Schweizer, S. Kornfeld, and J. Rohrer (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14471–14476.
19. E. Diaz and S. R. Pfeffer (1998) *Cell* **93**, 433–443.
20. F. Aniento, F. Gu, R. G. Parton, and J. Gruenberg (1996) *J. Cell Biol.* **133**, 29–41.

## Endopeptidase

Peptidases are [enzymes](#) that catalyze the hydrolysis of the **peptide bonds** that link [amino acids](#) into linear **polypeptide chains**. The Enzyme Commission of the International Union of Biochemistry and Molecular Biology classifies them as E.C. 3.4, *peptide hydrolases*. The chemical reaction that occurs can be depicted in its most simplified form as  $R - CO - NH - R' + H_2O \leftrightarrow R - CO - OH + H_2N-R'$ .

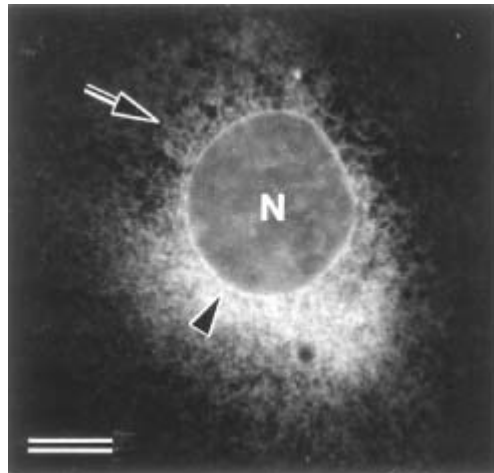
In a linear polymer of  $x$  amino acids ( $x$  can range from 2 to  $>1000$ ), there are  $x - 1$  peptide bonds. Enzymes that catalyze the hydrolysis of the first or second (or sometimes even the third) peptide bond in such a linear polymer are called *exopeptidases* and, more specifically, [aminopeptidases](#) (E.C. 3.4.11–15). This is because the amino group of the first amino acid in a linear polymer of amino acids is not part of a peptide bond, and hence this end of the chain of  $x$  amino acids is called the *amino terminus* (or N-terminus). Enzymes that catalyze the hydrolysis of the last or the penultimate peptide bond in this chain are also exopeptidases, but in this case they are [carboxypeptidases](#) (E.C. 3.4.16–19). The carboxyl group of the last amino acid in the chain is similarly not part of a peptide bond, and this end of the chain is therefore the *carboxy terminus* (or C-terminus).

Peptidases that catalyze the hydrolysis of internal peptide bonds in the polypeptide chain are called *endopeptidases* (E.C. 3.4.21–24, 99). The subclass to which they are assigned depends on the particular type of mechanism they employ to achieve catalysis (see [Proteinases](#)).

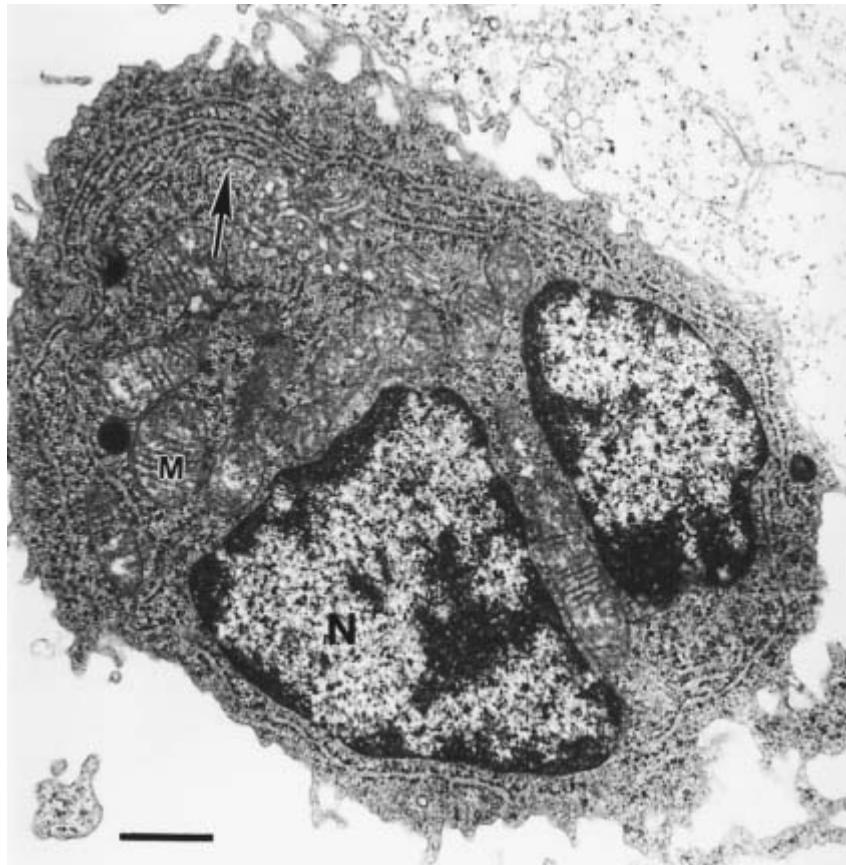
## Endoplasmic Reticulum

The endoplasmic reticulum (ER) is an extensive network of membranous cisternae that extends throughout much of the cytoplasm (Fig. 1). The three main functions of this organelle are: (i) biosynthesis of proteins that are destined to be secreted, deposited on the cell surface, or reside in other organelles of the [exocytosis](#) and [endocytosis](#) pathways (see [Protein Secretion](#)). (ii) synthesis of [lipids](#); and (iii) regulation of calcium ( $\text{Ca}^{2+}$ ) storage (see [Calcium Signaling](#)). Because the ability to perform these functions is essential for growth and differentiation of all cells, the ER is ubiquitous. Yet it undergoes dynamic changes that allow it to adapt to the changing physiological demands of various cells and tissues. The unique shape and architecture of the ER reflect its functional connections to other organelles. At the cell center, the ER is contiguous with the **nuclear** envelope. At the cell periphery it becomes a long, tubular network, often enveloping [mitochondria](#) and sometimes extending to the proximity of the plasma [membrane](#) (Figs. 1–2). This reticular architecture mediates the physical connections, as well as the chemical communication, between the ER and most other cellular compartments. A recurrent theme in the sections that follow, describing individual functional aspects of the ER, is the multitude of ways through which the ER exchanges biochemical components, be they ions, lipids, or proteins, with other organelles.

**Figure 1.** Light microscopic image of the ER. Immunofluorescence staining of a COS-1 cell expressing an ER-retained form of immunoglobulin light chain. The ER staining fills most of the cytoplasm, with the reticular appearance most evident at the periphery (arrow). The typical staining of the nuclear envelope is visible as a prominent ring (arrowhead) around the nucleus (N). Bar, 10  $\mu\text{m}$ .



**Figure 2.** Electron micrograph of a plasma cell. A murine antibody-secreting cell exhibiting a typically expanded ER. Arrow, a region of the endoplasmic reticulum with stacked flattened cisternae of rough ER, each with the typical beaded appearance due to densely bound ribosomes. N, nucleus. M, mitochondrion; bar, 1  $\mu\text{m}$ .



## 1. The Structure of the ER

The reticular architecture of the ER is maintained by numerous dynamic contacts between the ER membrane and [microtubules](#). These connections are mediated by the molecular motor complex [kinesin](#), which attaches to microtubules via its globular, *N*-terminal head **domain** and to the ER membrane via its *C*-terminal domain. The latter interaction is presumably mediated by kinectin, a 160-kDa integral [membrane protein](#) of the ER, which has been shown to be a kinesin receptor ([1](#), [2](#)). These attachments enable the ER cisternae to extend or contract along microtubules, which themselves assemble and disassemble dynamically. The dynamic nature of the ER structure is evident when cells are stressed: Drastic changes in cell shape, caused by depolymerization of microtubules, are accompanied by collapse of the reticular ER toward the cell center ([3](#)). The precise molecular regulation of the attachment to microtubules, and possibly other **cytoskeletal** elements, are yet to be elucidated. Other important unresolved issues include how the ER cisternae maintain a flattened structure, with a very high surface-to-volume ratio, and what the physiological implications are of this cisternal structure.

## 2. Subcompartments of the ER

In higher eukaryotes, that part of the ER to which [ribosomes](#) are attached is known as the rough ER, due to its distinctive beaded appearance in electron micrographs ([4](#)) (see Fig. [2](#)). In “professional” secretory cells, the rough ER is the vastly predominant form, because of the demand for high throughput of **protein biosynthesis**, folding, and trafficking (see text below). The morphology of the rough ER reflects these biosynthetic roles, and underlying it is a distinct protein composition. The rough ER membrane is the site of expression of [ribophorins](#) and many other components involved in docking ribosomes and in translocating nascent proteins across the membrane (eg, see Refs. [5](#) and [6](#)) (see [Protein Secretion](#)).

The remainder of the ER, devoid of attached ribosomes, is called the smooth ER. The smooth ER membrane is the site of residence of [cytochrome P450](#) and other components of the mixed-function oxidase system. The smooth ER is also involved in synthesis of various lipids and in the assembly of lipoprotein particles. Thus, the smooth ER predominates in cells where the expression of such proteins is high (eg, hepatocytes), while it is only a small fraction of the ER in secretory cells like pancreatic acinar cells.

The ER membrane proteins are clearly segregated between the smooth and rough subcompartments, although the molecular mechanism for this segregation is unresolved. Microscopic studies following the distribution of proteins like the binding protein [BiP](#) or [protein disulfide isomerase](#) (PDI) in a variety of cells suggest that the entire lumen of the ER is contiguous, with no obvious separation between smooth ER and rough ER lumen or between the reticular ER in the cell periphery and the nuclear envelope.

In excitable cells, the ER is specialized both morphologically and biochemically and is known as the *sarcoplasmic reticulum* (SR). Smooth cisternae of the SR make frequent contacts with the membrane of muscle cells. These terminal cisternae contain complexes of proteins involved in  $\text{Ca}^{2+}$  modulation, including membrane proteins like the ryanodine receptors and specialized SR [calcium-binding proteins](#), like calsequestrin. The segregation of proteins is clearly related to  $\text{Ca}^{2+}$  storage and release, as general luminal proteins in the SR (like BiP) do not exhibit similar subcompartmentalization (7).

### 3. Protein Translocation

The rough ER is the organelle that couples protein biosynthesis to membrane translocation, allowing segregation of secretory and membrane proteins from cytosolic or nuclear proteins. Secreted and membrane proteins are marked for translocation by a **signal sequence**, which is usually cleaved early in the biosynthesis of the proteins (8). Signal sequences are generally *N*-terminal, consist of around 20 amino acid residues, and are characterized by a central string of **hydrophobic** amino acids flanked by clusters of charged residues, and a residue with a small side chain precedes the cleavage site (9). Once a signal sequence has emerged from a ribosome, that ribosome is targeted to the ER membrane, where the nascent chain is translocated across the membrane. This targeting is accomplished primarily by a cytoplasmic ribonucleoprotein complex, the [signal recognition particle](#) (SRP) (10). SRP performs three distinct activities: (i) binding to the signal sequence that is exposed on the surface of a **translating** ribosome; (ii) retarding further elongation of the nascent polypeptide; and (iii) docking onto the SRP receptor in the ER membrane. Upon completion of this sequence of events, the SRP detaches, and the signal sequence is inserted into the membrane. Due to the intrinsic hydrophobicity of signal peptides, they themselves may interact with the lipid bilayer and thus aid the translocation process at this stage. After the initial association of the nascent chain with the membrane, the ribosomes bind more avidly to the ER by [protein-protein interactions](#) with the Sec61p complex, a seven-component transmembrane complex of proteins.

Nascent polypeptides cross the lipid bilayer through a protein channel (11), known as the *translocon*. The minimal translocon has been reconstituted *in vitro* (12, 13), and the necessary components were shown to be the SRP receptor, the Sec61p complex, and a translocating chain-associated membrane protein (TRAM). Purified mammalian and yeast Sec61p complexes in [detergent](#) solution or in proteoliposomes form cylindrical oligomers of three to four heterotrimers each, with a diameter of approximately 85 Å and a central pore of approximately 20 Å. Formation of the reconstituted translocon is stimulated by association with ribosomes or with the Sec62/63p complex (14).

At early stages of translocation, prior to signal sequence cleavage, nascent chains can be shown by chemical **cross-linking** to contact the Sec61a subunit and the signal sequence receptor (p34), via both the signal peptide and the remainder of the polypeptide chain (15, 16). At later stages of translocation, which can be trapped experimentally with the use of truncated translocation



intermediates, the proximity to Sec61p is markedly reduced. These data are in line with the interpretation that, at early stages of translocation, nascent chain-Sec61p interactions are in part mediated through interactions of the ribosome with components of the ER membrane, but that subsequent translocation can occur after dissociation of the ribosome (15).

Once the signal sequence enters the lumen, it is cleaved off most proteins. The enzyme responsible for this cleavage, signal peptidase, is a hetero-oligomeric membrane protein containing subunits of 12, 18, 21, 22 / 23, and 25 kDa. The 18- and 21-kDa subunits are mammalian homologues of Sec11p, a yeast protein that is necessary for signal peptide processing and cell viability. The topology of signal peptidase places the bulk of the protein in the ER lumen, including its putative **serine**-containing [active site](#), located in the 18- and 21-kDa subunits (17).

A second route for protein entry into the ER is post-translational and independent of SRP. In this route, which is more extensively used in yeast than in higher eukaryotes, a complete polypeptide chain is made in the cytoplasm and held in an incompletely folded state by **molecular chaperones** such as cytosolic hsp70. Fully synthesized secretory precursors then engage the ER membrane by interaction with the Sec complex, which in post-translocational translocation contains further protein subunits that are required for function, in addition to those of the co-translational Sec61 complex. It is not yet known what distinguishes a polypeptide chain that uses the co-translational route from a polypeptide that is translocated post-translationally. Productive post-translational interaction with Sec61p requires displacement of hsp70 from the precursor polypeptide. This reaction is facilitated by the chaperone Ydj1p, a homologue of the *E. coli* [DnaJ protein](#), which stimulates the inherent [ATPase](#) activity of the hsp70 chaperone. In the ER lumen, another hsp70-like chaperone, BiP, is an important part of the translocation machinery. In yeast ER, BiP is known to associate with the translocation complex via its contact with the DnaJ-homologous domain of the transmembrane protein Sec63p (18). The membrane-bound BiP can hydrolyze ATP, and the polypeptide moves through the channel formed by the Sec complex and is released at the luminal side (19). The process of post-translational translocation is made vectorial because of the cleavage of signal sequences and because the interactions with the hsp70s on either side of the membrane are unequal; cytosolic hsp70 and BiP cannot replace each other.

In the mammalian ER, as in yeast, BiP is required for the completion of polypeptide translocation. Moreover, BiP in its ATP-bound, non-substrate-binding state also functions to seal the translocon (20). In this fashion, it sustains the permeability barrier of the ER membrane, which otherwise would be compromised by the existence of thousands of translocation channels (21).

#### 4. Protein Post-translational Modifications

Membrane and secreted proteins can be modified by a number of [post-translational modifications](#), the most common being glycosylation and [phosphorylation](#). The most prevalent glycosylation is via the addition of a 14-sugar oligosaccharide onto certain [Asparagine \(Asn, N\)](#) residues, termed [N-glycosylation](#). The oligosaccharide usually consists of three glucose (Glc), nine mannose (Man) and two *N*-acetylglucosamine (GlcNAc) residues and is called *high mannose*. This oligosaccharide is assembled onto a dolichol lipid carrier by a series of [enzymes](#) that reside in the ER membrane. The first seven sugars are assembled by enzymes facing the cytosol. The Man<sub>5</sub>GlcNAc<sub>2</sub>-dolichol intermediate is then translocated across the membrane, and the final assembly steps are performed on the luminal face (22, 23).

A multisubunit enzyme, oligosaccharyltransferase (OST), is responsible for the transfer of the preassembled high-mannose oligosaccharide from the dolichol-linked donor onto Asn acceptor sites in nascent proteins on the luminal side of the ER membrane (24). These acceptor sites are within the tripeptide Asn-X-Ser/Thr (where X is any amino acid except Pro), provided that this tripeptide is accessible to OST.

Immediately upon addition to the protein, the high-mannose oligosaccharide is subjected to progressive trimming. First, the terminal Glc is removed by glucosidase I, followed by sequential removal of the next 2 Glc by glucosidase II. These steps are followed by a third trimming enzyme, a mannosidase that trims the oligosaccharide down to the Man<sub>5</sub>–Man<sub>8</sub> forms (25). Additional trimming occurs after the glycoproteins traffic to the [Golgi apparatus](#), before final remodeling to complex types of oligosaccharides.

A second type of glycosylation, [O-glycosylation](#) of serine or [Threonine \(Thr, T\)](#) residues by shorter and different carbohydrates, is a less frequent modification and is currently thought to take place mostly in the Golgi complex and not in the ER.

Another common protein modification is phosphorylation of Ser/Thr residues, and this modification is thought to play a regulatory role in the function of these proteins. A number of itinerant and resident ER proteins are known to be phosphorylated, on cytoplasmic domains as well as on luminal domains. The former phosphorylation can be performed by a number of known kinases in the cytosol. The phosphorylation of luminal domains has, in most cases, the hallmarks of a casein kinase II-like enzyme. However, no ER-specific kinase has yet been purified. The phosphorylation of Tyr residues, of major importance in signaling at the plasma membrane, has not yet been observed within the ER.

## 5. Lipids

In addition to its role in protein synthesis, the ER is a major site of lipid synthesis, and both phospholipid and sterol biosynthesis enzymes reside in this organelle. While some phospholipids are synthesized by cytosolic enzymes, much of phospholipid biosynthesis is performed by base-exchange enzymes and phospholipase D in the ER. Base-exchange enzymes alter the polar head group, allowing for interconversion of different phospholipid species. Phospholipase D hydrolyzes phospholipids to phosphatidic acid, which is used for further lipid biosynthesis.

The ER houses enzymes involved in cholesterol biosynthesis, as well as those that participate in oxidation of cholesterol to bile acids. Of them, HMG-CoA reductase is a key rate-limiting enzyme in steroidogenesis, and the intricate control of its expression is another manifestation of the communication between the ER and other organelles.

Because cholesterol and the apolipoproteins are both synthesized in the ER, much of the assembly of lipoproteins takes place in this organelle. The bulk of the very-low-density lipoprotein (VLDL) triacylglycerol and some VLDL phosphoglyceride is introduced early in the secretory pathway, soon after apolipoprotein synthesis. On the other hand, a significant fraction of VLDL phosphoglyceride associates with the resulting triacylglycerol-rich lipid–protein complexes in more distal secretory compartments, just prior to their secretion as mature VLDL.

## 6. Calcium

Changes in intracellular Ca<sup>++</sup> concentration are a common [signal transduction](#) device linking plasma membrane receptors with the initiation of specific cellular responses, including cell growth, muscle cell contractility, and exocytosis in secretory cells. The ER and its specialized form, the SR, is a major dynamic calcium storage organelle in eukaryotic cells. Its Ca<sup>2+</sup> concentration is estimated to reach 10 mM, whereas free Ca<sup>2+</sup> in the cytosol is only in the nanomolar range. The ER Ca<sup>2+</sup> storage capacity is provided by a multitude of luminal proteins. Among them are the ubiquitous and abundant chaperones, like GRP94, calreticulin (see [Calnexin/Calreticulin](#)), and the PDI-like proteins, or proteins whose expression is tissue-specific, like calsequestrin in the SR. Most of the Ca<sup>2+</sup> binding is via multiple (as many as 50) low-affinity ( $K_d = 1$  mM) binding sites per protein, but some proteins have a few high-affinity sites, including **EF-hand** domains (26). The ER membrane

regulates the flow of  $\text{Ca}^{2+}$  via three types of proteins: (i) the ER/SR-  $\text{Ca}^{2+}$ -ATPase (SERCA), which is a pump that can be inhibited by the drug thapsigargin; (ii) **calcium channels**, which are either the ubiquitous inositol trisphosphate receptors or the ryanodine receptors; (iii) and  $\text{Na}^+/\text{Ca}^{2+}$  exchangers. Together, the membrane  $\text{Ca}^{2+}$  gates and the luminal binding proteins provide a rapid response system to mobilize  $\text{Ca}^{2+}$  as physiological conditions demand. Release of ER-stored  $\text{Ca}^{2+}$  into the cytosol affects diverse sets of proteins, such as EF-hand proteins, [annexins](#), and the C2 region proteins ([26](#)).

The  $\text{Ca}^{2+}$  homeostasis function of the ER is most specialized in muscle cells, where  $\text{Ca}^{2+}$  mobilization is an important component of the contractile apparatus.  $\text{Ca}^{2+}$  binding proteins and  $\text{Ca}^{2+}$  channels are highly expressed in the SR, often employing muscle-specific versions of these proteins. For example, calsequestrin is an SR-specific  $\text{Ca}^{2+}$  buffer that has structural [homology](#) to calreticulin, which is expressed ubiquitously in higher eukaryotic cells. In nonexcitable cells, the need to store  $\text{Ca}^{2+}$  can be satisfied by many of the luminal proteins. For example, even though calreticulin can store  $\text{Ca}^{2+}$ , its ablation reveals that it is an essential modulator of [integrin](#) adhesive functions and integrin-initiated signaling, but is dispensable for the storage of luminal  $\text{Ca}^{2+}$  ([27](#)).

As discussed above, the structure and biochemical composition of the SR reflect its important role in mediating signals emanating from cell-surface receptors. The proximity to the plasma membrane of distal SR cisternae enriched in inositol trisphosphate and ryanodine receptors position the relevant receptors to interact with [second messengers](#) generated during signal transduction and to respond by modulating the storage or release of  $\text{Ca}^{2+}$ . Another example is the expression of phospholamban, a 52-residue protein that assembles into a pentamer in SR membranes. Phospholamban plays a role in regulating the resident  $\text{Ca}^{2+}$  ATPase through an inhibitory association that can be reversed by phosphorylation. The phosphorylation of phospholamban is initiated by  $\beta$ -adrenergic stimulation, identifying phospholamban as an important component in the stimulation of cardiac activity by  $\beta$ -agonists ([28](#)). These examples illustrate the functional connections that exist between the ER and the plasma membrane, connections that are important for efficient signal transduction.

## 7. Protein Folding in the ER

The lumen of the ER is the environment where all secretory polypeptide chains and those destined to reside in various cellular organelles fold. Like other **protein folding** compartments, the ER has a very high concentration of soluble proteins (estimated to be as high as 100 mg/mL), a factor that is conducive to aggregation. The ER is very different from other folding compartments in its higher [oxidation/reduction potential](#), high  $\text{Ca}^{2+}$  concentration, and lack of energy source (ATP is imported via a special permease). These conditions necessitate the presence of a number of ER enzymes and chaperones to facilitate protein folding. Moreover, a number of mechanisms have evolved to ensure that only properly configured proteins proceed along the secretory pathway and that defective molecules are disposed of ([29](#), [30](#)). Indeed, most of the abundant luminal proteins are involved in protein folding.

One of the main ER chaperones is BiP/GRP78, which is the ER member of the hsp70 family of **stress-response** proteins. *In vitro*, BiP (and other hsp70 proteins) bind hydrophobic peptides, preferring peptides of 7–10 amino acids with either alternating ([31](#)) or clustered aromatic or large aliphatic residues ([32](#)). *In vivo*, BiP probably binds similar hydrophobic sequences that are exposed in the folding polypeptide chain, but inaccessible in the native structure. One important result of this activity is to minimize aggregation of the unfolded polypeptide chain. Peptide binding stimulates the slow ATPase activity of BiP, and the substrate is released after either nucleotide exchange or hydrolysis ([33](#)), allowing it another chance to bury its hydrophobic sequences as its folding progresses. In its capacity as a peptide-binding protein, BiP is often found to associate preferentially with nonsecretable or misfolded mutants, presumably because these molecules continue to expose

binding sites that are buried more rapidly in wild-type molecules. In this fashion, BiP plays a role in the ER quality control system (34), either retarding the transport of proteins (35) or even targeting them for **protein degradation** (36).

As in the case of folding in the cytosol, an important rate-limiting step in folding in the ER can be the **cis-trans isomerization** of the **peptide bonds** of certain proline residues. This reaction is catalyzed by **peptidyl prolylcis/trans isomerases**, and it seems to be a redundant activity that several luminal proteins, including **S-cyclophilin** and **FK506-binding protein 13**, can perform.

A general property distinguishing cytosolic proteins from those that fold in the ER is the presence of **disulfide bonds** in the latter. Cysteine residues in nascent polypeptides oxidize in the lumen because of the higher redox potential, maintained by the approximately 3:1 ratio of reduced (GSH) to oxidized (GSSG) forms of **glutathione** (37). Because disulfide bonds form covalent folding intermediates, they are very important in the folding process (38). Thus, a number of luminal enzymes are capable of catalyzing the rearrangement of protein disulfide bonds. The most abundant of these enzymes are the 57-kDa protein disulfide isomerase and two structurally related proteins, ERp72 and ERp58. There is currently no evidence for the existence of disulfide isomerases in any other organelle of eukaryotes.

A second general property of proteins that fold in the ER is their frequent modification by Asn-linked glycosylation, as discussed above. The state of modification of these carbohydrates can be used to monitor the folding of **glycoproteins**. The terminal Glc residues on the initial high-mannose oligosaccharide are trimmed sequentially by glucosidase I, an ER membrane protein, and glucosidase II, a luminal protein. The enzyme that can reattach terminal Glc residues, UDP-Glc:glycoprotein glucosyltransferase, can distinguish between folded and unfolded polypeptides (39). Also involved in this monitoring system are the membrane-spanning **calnexin** and the soluble,  $\text{Ca}^{2+}$ -binding calreticulin. They are structurally related, and both can bind to monoglucosylated glycoproteins, an intermediate stage in the processing. In this way, calnexin and calreticulin bind to a partly trimmed intermediate and serve as sensors of the glycosylation state. The final Glc is then trimmed, and the glycoprotein is then either reglucosylated by the glucosyltransferase, if its state of folding is immature, or recognized as mature enough to be packaged for transport to the Golgi complex.

Although calnexin and calreticulin have similar fine specificities of binding to the trimmed oligosaccharides, they can clearly distinguish and bind different folded forms of the same protein. The basis for this difference is not clear. Nonetheless, gene ablation of calnexin and calreticulin in both *Saccharomyces cerevisiae* and mammalian cells in culture show that their roles in protein folding are redundant. An important unresolved issue is whether the **lectin-like** activity of the two chaperones is sufficient to account for their recognition of folding species, or whether the chaperones also interact with the folding polypeptides themselves.

In addition to BiP, calnexin, and calreticulin, there are several other resident ER proteins that function as molecular chaperones. Among these are GRP94 and PDI, which have been shown to associate transiently with folding intermediates. ERp72, ERp58, and the **immunophilins** are also potential chaperones or enzymes involved in folding, based on their association with other proteins and on their sequence similarity to cytosolic chaperones. GRP94 (also known as ERp99, gp96, endoplasmic, and hsp108) is a homodimeric glycoprotein that binds  $\text{Ca}^{2+}$  and ATP and makes up 5% to 10% of the total luminal protein. It is known to bind a number of peptides (40), as well as a narrow spectrum of proteins. GRP94 is a stress protein from the hsp90 family, whose synthesis is induced upon glucose starvation,  $\text{Ca}^{2+}$  depletion, or treatments with **tunicamycin** and amino acid analogues. Physiologic demands can also induce GRP94 expression: It is up-regulated during the **differentiation** of resting B lymphocytes into **antibody-secreting cells** (41) and within plasma cells in response to the production of nonsecretable immunoglobulins. The regulation of GRP94 under all these conditions is coordinated with other ER proteins, most notably BiP (41), and to a lesser extent

PDI and ERp72 (42). The role of GRP94 (and those of ERp72, ERp58 and immunophilins) in protein folding *in vivo* remains unclear, however, largely because its substrate specificity and detailed modes of action are still to be elucidated.

The growing list of ER chaperones raises several questions: Do they perform redundant functions or are they functionally distinct? How do they discriminate among the various proteins that fold in the ER? How do they participate in different aspects of the folding process? Regulation of the activity of molecular chaperones is another fundamental problem that is only beginning to be addressed. Phosphorylation has been shown to regulate the activity of BiP (43). GRP94, calnexin, and calreticulin are known to be phosphorylated *in vivo*, but the relation of the modification to their function is not yet understood. Thus, an intriguing emerging possibility is that ER chaperones cycle between active and inactive states that are marked by phosphorylation. The presence of so many folding factors clearly provides an environment that promotes efficient folding of proteins that are destined for export and in many cases are produced in large amounts.

## 8. Degradation of ER Proteins

An inevitable by-product of an active folding compartment is the presence of misfolded proteins and the need to degrade them (see [DNA Degradation In Vivo](#)). For many years it was thought that such degradation occurs within the ER itself, but this notion was challenged by the inability to demonstrate an appropriate set of ER [proteinases](#), equivalent to those that reside in the cytosol. Recent discoveries, in both mammalian systems (44, 45) and in yeast (46, 47), have now clearly demonstrated the presence of a conserved retrograde traffic pathway, whereby misfolded secretory proteins are presented to the cytosolic [proteasome](#) machinery for degradation (48). Misfolded soluble and membrane proteins are translocated out of the ER, probably through the same Sec61 channel used for incoming nascent chains, and are guided with the aid of cytosolic hsp70 to 26S proteasome particles, many of which are in proximity to the ER membrane. Often, as in degradation of cytosolic proteins, the misfolded proteins that are translocated out of the ER are marked for degradation by covalent addition of poly-[ubiquitin](#) chains. The proteasome is apparently highly processive, because no large degradation intermediates are seen, and the products are short peptides with a narrow size distribution. For many membrane and secreted proteins, proteasome-mediated degradation is now thought to be the major route of disposal if they fail to fold properly, although the existence of a proteolytic pathway within the ER itself is still a possibility.

## 9. Sorting of Itinerant and Resident ER Proteins

The transport of proteins from the ER to the Golgi complex is via **vesicles**. Cargo proteins are selectively concentrated into “exit sites” marked by 10-nm-thick electron-dense cytosolic coats that are made of two distinct biochemical complexes, called COPI and COPII (49). Under the control of a fission machinery that uses small [Gtpases](#) to regulate protein–protein interactions, COP-coated vesicles bud off the ER membrane and are targeted for fusion with the next membrane in the pathway. The mechanism of membrane budding driven by COPII shows remarkable overall similarities to that of COPI budding, the coat that is also involved in intra-Golgi transport and in retrograde transport to the ER, as well as to budding of **clathrin**-containing coated vesicles, which mediate transport from the plasma membrane and from the *trans*-Golgi network. The discrimination between cargo proteins and resident ER proteins at the level of packaging into the vesicles is thought to be mediated by a family of small 24-kDa proteins (50). Vesicle budding from the ER has been reconstituted with purified membranes and three soluble proteins from yeast: Sec13 complex, Sec23 complex, and the small GTPase Sar1p (51).

Luminal residents of the ER, like the chaperones BiP and GRP94 and the enzyme PDI, use a receptor-mediated mechanism to ensure their localization within the ER lumen. Their C-terminus consists of the tetrapeptide Lys–Asp–Glu–Leu (KDEL), which is recognized specifically by an intracellular receptor that resides mainly throughout the Golgi complex. Similar receptors with slightly altered tetrapeptide specificities are found in lower eukaryotes. KDEL proteins that “escape”

the ER are specifically retrieved by this receptor and transported back to the ER ([52](#), [53](#)).

Resident ER membrane proteins often have a different type of retrieval signal that distinguishes them from itinerant membrane proteins: a dibasic amino acid (either Lys–Lys or Arg–Arg) motif in their cytosolic domain near the C-terminus. This motif is recognized by the COPI complex, providing the basis for specifically packaging these proteins in vesicles that transport them retrogradely from the Golgi to the ER ([54](#)).

The combination of biochemical experiments reconstituting transport vesicles with the elucidation of retention mechanisms that differentiate resident ER proteins from itinerant proteins brings us closer to solving the important problem of how proteins are selectively sorted and packaged into transport vesicles, a question that is at the heart of the compartmentalization of eukaryotic cells.

## 10. Biogenesis of the ER

In fibroblasts, the ER accounts for 2% to 5% of the cytoplasmic volume, whereas in “professional” secreting cells, like pancreatic acinar cells or lymphoid plasma cells, the ER is considerably expanded, its cisternae are often stacked together ([4](#)), and it occupies up to 25% of the cytoplasm (Fig. [2](#)). For example, during the development of **B cells** into plasma cells, when synthesis and secretion of immunoglobulin increase by two orders of magnitude, the ER grows, with both the volume and membrane area increasing three- to fourfold ([41](#)), and it changes from a network surrounding the nuclear envelope to one that fills the peripheral cytoplasm. There is a coordinate increase in the synthesis of the main luminal proteins, as well as of membrane proteins involved in the translocation process.

This phenomenon is probably the normal counterpart to the induction of luminal proteins observed as a result of accumulation of unfolded proteins in the ER under a variety of metabolic stresses, the so-called [unfolded protein response](#) (UPR). The accumulation of unfolded protein in the ER causes induction of the biosynthesis of ER proteins to alleviate the situation.

A second pathway of signal transduction from the ER to the nucleus regulates sterol biosynthesis. A small family of transcription factors, the sterol [response element](#) binding proteins (SREBPs), activates sterol-regulated **promoters**, such as those of the genes for the key enzyme HMG-CoA reductase or the LDL receptor, as well as other enzymes that synthesize fatty acids ([55](#)). SREBP is synthesized as an ER membrane protein; upon sterol depletion, it is processed by two proteinases, one of them probably a thiol proteinase, releasing the active DNA-binding transcription factor and enabling its translocation to the nucleus ([56](#)). In this fashion, coordinate metabolic control over two types of lipids necessary for membrane biosynthesis is maintained and is triggered by alterations in the lipid composition of the ER membrane.

HMG-CoA-reductase is also at the center of another ER biogenesis phenomenon. When it is overexpressed as a result of compactin treatment, there is a reversible proliferation of smooth ER membranes, which pack in hexagonal arrays ([57](#)). A similar elaboration of ER-derived smooth membranes is seen in various yeast strains deficient in inositol metabolism. These may represent cases where one of the feedback mechanisms that regulate the size of the ER “overreact,” leading to excessive synthesis of some components, instead of coordinate synthesis of all components, as is required to maintain homeostasis of the ER.

## Bibliography

1. I. Toyoshima, H. Yu, E. R. Steuer, and M. P. Sheetz (1992) *J. Cell Biol.* **118**, 1121–1131.
2. H. Yu, C. V. Nicchitta, J. Kumar, M. Becker, I. Toyoshima, and M. P. Sheetz (1995) *Mol. Biol. Cell* **6**, 171–183.
3. M. Terasaki, L. B. Chen, and K. Fujiwara (1986) *J. Cell Biol.* **103**, 1557–1568.
4. G. Palade (1975) *Science* **189**, 347–358.

5. D. I. Meyer, E. Krause, and B. Dobberstein (1982) *Nature* **297**, 647–650.
6. S. Tajima, L. Lauffer, V. L. Rath, and P. Walter (1986) *J. Cell Biol.* **103**, 1167–1178.
7. R. Sitia and J. Meldolesi (1992) *Mol. Biol. Cell* **3**, 1067–1072.
8. G. Blobel and B. Dobberstein (1975) *J. Cell Biol.* **67**, 835–851.
9. G. von Heijne (1984) *J. Mol. Biol.* **173**, 243–251.
10. P. Walter, I. Ibrahimi, and G. Blobel (1981) *J. Cell Biol.* **91**, 545–550.
11. S. M. Simon and G. Blobel (1991) *Cell* **65**, 371–380.
12. C. V. Nicchitta, G. Migliaccio, and G. Blobel (1991) *Cell* **65**, 587–598.
13. D. Gorlich and T. A. Rapoport (1993) *Cell* **75**, 615–630.
14. D. Hanein, K. E. Matlack, B. Jungnickel, K. Plath, K. U. Kalies, K. R. Miller, T. A. Rapoport, and C. W. Akey (1996) *Cell* **87**, 721–732.
15. C. V. Nicchitta, E. C. R. Murphy, R. Haynes, and G. S. Shelness (1995) *J. Cell Biol.* **129**, 957–970.
16. S. High, D. Gorlich, M. Wiedmann, T. A. Rapoport, and B. Dobberstein (1991) *J. Cell Biol.* **113**, 35–44.
17. G. S. Shelness, L. Lin, and C. V. Nicchitta (1993) *J. Biol. Chem.* **268**, 5201–5208.
18. S. K. Lyman and R. Schekman (1995) *J. Cell Biol.* **131**, 1163–1171.
19. K. E. Matlack, K. Plath, B. Misselwitz, and T. A. Rapoport (1997) *Science* **277**, 938–941.
20. B. D. Hamman, L. M. Hendershot, and A. E. Johnson (1998) *Cell* **92**, 747–758.
21. T. A. Rapoport (1992) *Science* **258**, 931–936.
22. R. Kornfeld and S. Kornfeld (1985) *Annu. Rev. Biochem.* **54**, 631–664.
23. M. D. Snider and C. B. Hirschberg (1987) *Annu. Rev. Biochem.* **56**, 63–88.
24. S. Silberstein and R. Gilmore (1996) *FASEB J.* **10**, 849–858.
25. J. Bischoff and R. Kornfeld (1983) *J. Biol. Chem.* **258**, 7909.
26. J. Meldolesi and T. Pozzan, (1998) *Trends Biochem. Sci.* **23**, 10–14.
27. M. G. Coppelino, M. J. Woodside, N. Demarex, S. Grinstein, R. St. Arnaud, and S. Dedhar (1997) *Nature* **386**, 843–847.
28. I. T. Arkin, P. D. Adams, A. T. Brunger, S. O. Smith, and D. M. Engelman (1997) *Annu. Rev. Biophys. Biomol. Struct.* **26**, 157–179.
29. A. M. de Silva, W. E. Balch, and A. Helenius (1990) *J. Cell Biol.* **111**, 857–866.
30. A. Helenius, T. Marquardt, and I. Braakman (1992) *Trends Cell Biol.* **2**, 227–231.
31. M. J. Gething, S. Blond-Elguindi, J. Buchner, A. Fourie, G. Knarr, S. Modrow, L. Nanu, M. Segal, and J. Sambrook (1995) *Cold Spring Harb. Symp. Quant. Biol.* **60**, 417–428.
32. S. Rudiger, A. Buchberger, and B. Bukau (1997) *Nat. Struct. Biol.* **4**, 342–349.
33. G. C. Flynn, T. G. Chappell, and J. E. Rothman (1989) *Science* **245**, 385–390.
34. S. K. Nigam, A. L. Goldberg, S. Ho, M. F. Rohde, K. T. Bush, and M. Y. Sherman (1994) *J. Biol. Chem.* **269**, 1744–1749.
35. A. J. Dorner, L. C. Wasley, and R. J. Kaufman (1992) *EMBO J.* **11**, 1563–1571.
36. M. R. Knittler, S. Dirks, and I. G. Haas (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1764–1768.
37. C. Hwang, A. J. Sinskey, and H. F. Lodish (1992) *Science* **257**, 1496–1502.
38. T. E. Creighton (1997) *Biol. Chem.* **378**, 731–744.
39. M. C. Sousa, M. A. Ferrero-Garcia, and A. J. Parodi (1992) *Biochemistry* **31**, 97–105.
40. N. E. Blachere, Z. Li, R. Y. Chandawarkar, R. Suto, N. S. Jaikaria, S. Basu, H. Udono, and P. K. Srivastava (1997) *J. Exp. Med.* **186**, 1315–1322.
41. D. L. Wiest, J. K. Burkhardt, S. Hester, M. Hortsch, D. I. Meyer, and Y. Argon (1990) *J. Cell Biol.* **110**, 1501–1511.

42. D. P. McCauliffe, Y. S. Yang, J. Wilson, R. D. Sontheimer, and J. D. Capra (1992) *J. Biol. Chem.* **267**, 2557–2562.
43. P. J. Freiden, J. R. Gaut, and L. M. Hendershot (1992) *EMBO J.* **11**, 63–70.
44. C. L. Ward, S. Omura, and R. R. Kopito (1995) *Cell* **83**, 121–127.
45. E. J. Wiertz, D. Tortorella, M. Bogoy, J. Yu, W. Mothes, T. R. Jones, T. A. Rapoport, and H. L. Ploegh (1996) *Nature* **384**, 432–438.
46. A. A. McCracken and J. L. Brodsky (1996) *J. Cell Biol.* **132**, 291–298.
47. M. M. Hiller, A. Finger, M. Schweiger, and D. H. Wolf (1996) *Science* **273**, 1725–1728.
48. E. D. Werner, J. L. Brodsky, and A. A. McCracken (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13797–13801.
49. S. Y. Bednarek, M. Ravazzola, M. Hosobuchi, M. Amherdt, A. Perrelet, R. Schekman, and L. Orci (1995) *Cell* **83**, 1183–1196.
50. K. Fiedler, M. Veit, M. A. Stamnes, and J. E. Rothman (1996) *Science* **273**, 1396–1399.
51. M. J. Kuehn, J. M. Herrmann, and R. Schekman (1998) *Nature* **391**, 187–190.
52. S. Munro and H. R. Pelham (1987) *Cell* **48**, 899–907.
53. J. C. Semenza, K. G. Hardwick, N. Dean, and H. R. B. Pelham, (1990) *Cell* **61**, 1349–1357.
54. R. D. Teasdale and M. R. Jackson (1996) *Annu. Rev. Cell Dev. Biol.* **12**, 27–54.
55. J. B. Kim, P. Sarraf, M. Wright, K. M. Yao, E. Mueller, G. Solanes, B. B. Lowell, and B. M. Spiegelman (1998) *J. Clin. Invest.* **101**, 1–9.
56. M. S. Brown and J. L. Goldstein (1997) *Cell* **89**, 331–340.
57. R. K. Pathak, K. L. Luskey, and R. G. Anderson (1986) *J. Cell Biol.* **102**, 2158–2168.

### **Suggestions for Further Reading**

58. G. Palade (1975) Intracellular aspects of the process of protein secretion. *Science*, **189**, 347–358.
59. R. Sitia and J. Meldolesi (1992) Endoplasmic reticulum: a dynamic patchwork of specialized subregions. *Mol. Biol. Cell* **3**, 1067–1072.
60. G. Blobel and B. Dobberstein (1975) Transfer of proteins across membranes. I. Presence of proteolytically processed and unprocessed nascent immunoglobulin light chains on membrane-bound ribosomes of murine myeloma? *J. Cell Biol.* **67**, 835–851.
61. T. A. Rapoport, B. Jungnickel, and U. Kutay (1975) Protein transport across the eukaryotic endoplasmic reticulum and bacterial inner membranes. *Annu. Rev. Biochem.* **65**, 271–303.
62. J. Meldolesi and T. Pozzan (1998) The endoplasmic reticulum Ca<sup>2+</sup> store: a view from the lumen *Trends Biochem. Sci.* **23**, 10–14.
63. C. Hammond and A. Helenius (1995) Quality control in the secretory pathway. *Curr. Opin. Cell Biol.* **7**, 523–529.
64. C. Sidrauski and P. Walter (1997) The transmembrane kinase Ire 1p is a site-specific endonuclease that initiates mRNA splicing in the unfolded protein response. *Cell* **90**, 1031–1039.
65. M. S. Brown and J. L. Goldstein (1997) The SREBP pathway: regulation of cholesterol metabolism by proteolysis of a membrane-bound transcription factor. *Cell* **89**, 331–340.

### **Endosomes**



Endosomes comprise a complex series of intracellular organelles that mediate crucial functions in the [endocytosis](#) pathway (1). Ligands internalized through **clathrin**-coated vesicles are first delivered to “early endosomes.” These are tubular/vesicular organelles distributed primarily in the peripheral cytoplasm of the cell. The internal pH of these structures is mildly acidic (pH ~6.0). The low pH is established and maintained by the membrane-associated vacuolar H<sup>+</sup>bATPase and regulated by the Na<sup>+</sup>K<sup>+</sup>ATPase. The low-pH environment is essential for certain sorting functions of endosomes. For example, it facilitates dissociation of some ligands, such as **low density lipoprotein** (LDL), from their receptors, allowing the receptor to be recycled to the cell surface and reused, whereas the ligand is retained in the cell and degraded. Disruption of the endosomal pH by weak bases, ionophores, or inhibitors of the vacuolar ATPase disrupts the sorting functions of these organelles. The acidic pH within endosomes is used by certain **viruses** (eg, [influenza virus](#)) and [toxins](#) (eg, [diphtheria toxin](#)) to facilitate membrane fusion or translocation to the cytoplasm.

Early endosomes are the primary sorting organelles in the endocytic pathway, with two principal exit routes: (1) the recycling pathway to the cell surface, and (2) the pathway to late endosomes and lysosomes (1). In addition, proteins may be sorted out of early endosomes to recycling endosomes (a subset of endosomes that have nearly neutral pH and are enriched in **transferrin receptors**), the *trans* **Golgi network** (TGN), and other locations in certain differentiated cell types (for example, the GLUT 4 glucose transporter-containing compartment in adipocytes), transcytotic vesicles in polarized cells, synaptic vesicles in neurones, and MHC class II-containing vesicles in some antigen-presenting cells (2). It is believed that transport between these compartments is mediated by carrier vesicles (although these have not been characterized in detail) that fuse with their target organelles through SNARE/NSF and Rab-dependent mechanisms (3).

Late endosomes receive internalized ligands and receptors from early endosomes and provide a second sorting compartment before lysosomes. Late endosomes are again tubular/vesicular and frequently contain many membrane vesicles or lamellae. The internal pH of late endosomes is more acid (pH ~5.5) than that of early endosomes. Transport from early endosomes to late endosomes is mediated by endosomal carrier vesicles and/or multivesicular bodies that derive from early endosomes.

The different endosomal compartments are not well-characterized biochemically or morphologically. However, analyses of the distribution of Rab family of proteins have suggested that these small [Gtpases](#) might provide functional markers for endosomal subcompartments (4). Thus Rab 5 is primarily associated with coated vesicles and early endosomes, Rab 7 with late endosomes, and Rab 4 with recycling endosomes (1). Because these proteins have been implicated in the specific fusion involved in vesicular transport, their restricted and unique distributions may define specific transport steps and possibly unique SNARE-dependent fusion (4).

(See also [Phagocytosis](#); [Pinocytosis](#); [Endosome](#); [Macropinosome](#).)

## Bibliography

1. J. Gruenberg and F. R. Maxfield (1995) *Curr. Opin Cell Biol.* **7**, 552–563.
2. I. Mellman, P. Pierre, and S. Amigorena (1995) *Curr. Opin Cell Biol.* **7**, 564–72.
3. J. E. Rothman and G. Warren (1994) *Curr. Biol.* **4**, 220–233.
4. P. Novick and M. Zerial (1997) *Curr. Opin Cell Biol.* **9**, 496–504.

## Endotoxins

The majority of bacterial [toxins](#) are proteins released by **bacteria** in the medium and are termed exotoxins. However, many [gram-negative bacteria](#), upon lysis or during division, release toxic [lipid](#) components of the outer [membrane](#), which are therefore named endotoxins. These molecules are lipopolysaccharides (LPS) of a complex and variable nature that are capable of inducing a variety of local and systemic effects, including fever and toxic shock (1). LPS binds directly to the CD11/CD18 molecules and 95-kDa scavenger receptor of **macrophages**(3). After fixation to a plasma protein LBP, the complex LPS-LBP is capable of binding the CD14 molecule present on macrophages, neutrophils, and platelets (2). Such binding triggers a variety of effects. Platelets release prostaglandin E<sub>2</sub>, tromboxane A<sub>2</sub>, and PAF (platelet-activating factor), and these mediators induce platelet aggregation and coagulation and thrombotic phenomena. Activated macrophages and neutrophils release oxygen radicals and [nitric oxide](#), which causes the relaxation of smooth muscles. LPS are the most potent inducer of the release of **cytokines** from macrophages. It appears that such platelet and inflammatory cells responses have developed in order to respond rapidly to the presence of potentially pathogenic bacteria (3). However, such a defensive response may turn easily into harmful, even deadly, systemic effects when larger amount of LPS are released or bacteria multiply in the blood and release LPS in a tissue so rich in target cells.

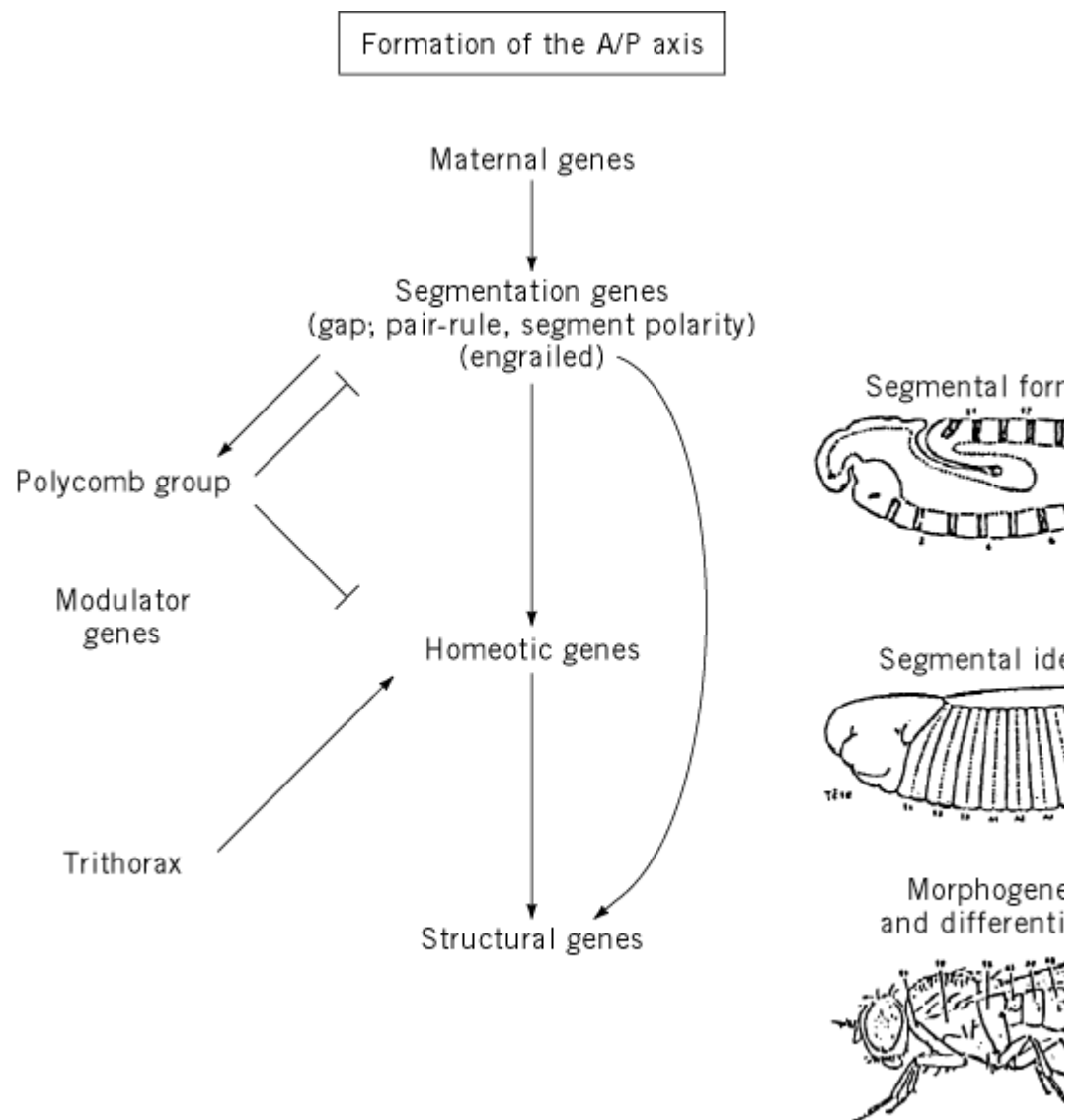
### Bibliography

1. J. A. Hewett and R. A. Roth (1993) *Pharmacol. Rev.* **45**, 381–411.
2. R. J. Ulevitch and P. S. Tobias (1995) *Ann. Rev. Immunol.* **13**, 437–457.
3. C. E. Mims (1995) *The Pathogenesis of Infectious Diseases*, Academic Press, London.

### Engrailed Gene

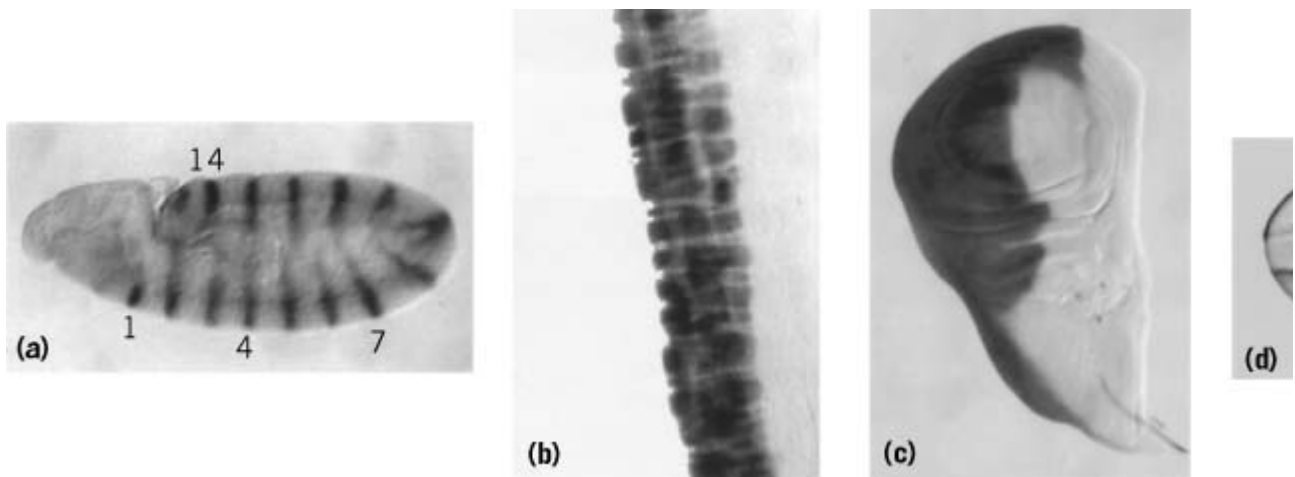
In multicellular organisms, the induction and maintenance of a **differentiated** state during [development](#) depend on the expression of a large number of genes acting at multiple levels (1). Genetic and molecular studies in *Drosophila* have largely contributed to the elaboration of this concept. During the embryonic development of *Drosophila*, the anterior-posterior (A/P) organization is generated by the activity of several maternal and zygotic genes (2, 3). Many of these genes encode transcriptional regulators that sequentially subdivide the embryo into smaller domains (segments and compartments) along the A/P axis (Fig.1). The segmentation gene *engrailed* (*en*) is one of the components of this regulatory cascade that defines the A/P axis during *Drosophila* embryogenesis. The *engrailed* function is essential to specify the posterior cell identity, as well as the A/P compartment boundaries (4-7). *en* is involved not only in the establishment of the posterior cell's fate but also in its maintenance. Indeed, *en* expression remains essential throughout the fly's life, from early embryonic development to the larval and adult appendage development (8-13) (Fig.2).

**Figure 1.** Cascade of genes that form the anterior-posterior (A/P) axis during embryogenesis. *engrailed* is a segment pol- formation of the segments. This gene directly regulates the expression of the homeotic genes involved in the identity of t expression of the *Polycomb Group* genes, in particular, *polyhomeotic* (58). It also interferes directly in morphogenesis ar structural genes, such as that for b3-tubulin (54).



**Figure 2.** Expression of *engrailed* during development. In each case, the expression of *engrailed* defines the posterior compartment (*en-lacZ*) producing an Engrailed- $\beta$ -galactosidase fusion protein has been used, and the *engrailed* gene expression for each stage is shown in (8).

- a. *engrailed* expression in a germ-band extended embryo. In this case, the transcripts were detected by *in situ* hybridization.
- b. *engrailed* expression of third instar *en-lacZ* larval hindgut.
- c. *engrailed* expression of third instar *en-lacZ* wing imaginal disk.
- d. *engrailed* expression of adult *en-lacZ* wing.



The positioning of the A/P boundaries established by the segmentation genes before gastrulation is an essential step in development because these boundaries serve as a reference for the subsequent and spatially restricted expression of the [homeotic genes](#) of the *Bithorax* and *Antennapedia* complexes (14). Each segment is characterized by the activity of a specific subset of homeotic genes, the other homeotic genes being repressed in that segment (14). The differential expression of the homeotic genes, which is essential to confer segmental identities, is established by the segmentation genes. The segmental identities are then maintained by the *Polycomb* group (*PcG*) and the *Trithorax* genes, the latter being required for the maintenance of the activation of the homeotic genes (15) and of *engrailed* (16). On the contrary, *PcG* genes are required for the maintenance of the repression of the homeotic genes (17), as well as of segmentation genes, like *engrailed* (18, 19). It has been suggested that the *PcG* proteins form multimeric complexes that maintain repression by [imprinting](#), limiting the accessibility of [transcription factors](#) to [chromatin](#) (17, 20, 21).

*engrailed* is also essential in the developing central nervous system, where it is segmentally expressed (22-24). Homologues of *engrailed* have been identified in other diptera (25) and other organisms, such as [zebrafish](#) (26), chickens (27), mice (28, 29), and humans (30). All these organisms express *engrailed* during neurogenesis. In vertebrates, the neural tube is a segmented region, where *engrailed* expression is segmentally reiterated (31). This indicates that some of the *engrailed* functions during cell determination might be conserved during evolution. For example, the genetic network regulating the development of the brain in vertebrates involves conserved genes that control segmentation in *Drosophila* (32). In particular, the *Wnt-1/engrailed* interaction was shown to be conserved from flies to mice (33).

In the *Drosophila* epidermis, *engrailed*-expressing cells also express *hedgehog* (*hh*) (34), which encodes a secreted protein (Hh) that influences the pattern of cell differentiation across the embryonic segments, indicating that these *en/hh*-expressing cells have organizer properties (35). In particular, interfaces between posterior cells, which express *engrailed*, and anterior cells, which do not, were shown to be required for pattern formation (36). In the embryos, just anterior to the *en*-expressing cells, cells express *wingless* (*wg*), which encodes a secreted protein. During early embryogenesis, interdependent regulations between *en* and *wg* occur. Later, their expression is regulated by different mechanisms and maintained by autoregulation (35-39).

Later in development, during wing morphogenesis, *en* is also involved in the specification and the maintenance of posterior patterns, as well as in creating the compartment border. *en* activity in the posterior cells directs them to express *hh*, which is secreted anteriorly (34, 40). *en* blocks the ability of the posterior cells to respond to Hh (12), whereas its absence in the anterior cells allows them to respond to Hh. Hh acts in a stripe of anterior cells abutting the A/P boundary through the Patched (*Ptc*) receptor, whose expression is also activated by the Hh signal (43). These anterior cells,

receiving the Hh signal, express *decapentaplegic* (*dpp*) (a homologue of the **transforming growth factor- $\beta$**  signaling molecule) which, in turn, exerts a long-range organizing activity on the growth and patterning of both compartments (44). Hence, *en* plays a major role in the control of **imaginal disk** growth, by regulating the expression of *dpp* (44, 45). Indeed, *en* has two separable roles: a cell autonomous role in defining the posterior compartment, by repressing genes that are expressed in the anterior compartment, such as *cubitus-interruptus* (*ci*) (46, 11), *ptc*, and *dpp* (45), and a nonautonomous role in the anterior compartment through its regulation of *hh* (34, 47). The *en*-related gene *invected* (*inv*) encodes redundant functions, even though *inv* performs only a discrete subset of the functions ascribed to *en* (40, 42). For example, *inv* does not regulate *hh* or all the genes of the *hh* signaling pathway (41).

Engrailed (En) is a homeodomain-containing protein that performs its key role during *Drosophila* development by acting as a transcription factor (48). Despite the considerable amount of research conducted on homeodomain proteins, little is known about their mechanisms of action. One issue is understanding the specificity of their action because they all bind to DNA sequences containing the same core motif 5'-TAAT-3', which is essential for their DNA-binding activity; the nucleotides flanking this core provide the specificity of binding (49, 50). The N-terminal residues of the homeodomain were found to contribute to this sequence-specific association. In particular, the identity of residue 50 of Engrailed affects the base preference, both *in vitro* and *in vivo* (51-53).

Because of this remarkable conservation of binding properties among all homeodomains, the discovery of their DNA target sequences *in vivo* has been difficult. Identification of their target genes is, however, an essential step in understanding their role during development. Such a study has been performed in order to understand how *engrailed* achieves its complex regulatory role throughout development (Fig. 2).

The identification of target genes for transcription factors has been undertaken using different approaches, based on two main ideas. The first approach is functional and is based on the existence of genetic interactions between two genes. Genetic experiments have shown that En is able to repress *ci* (46), *dpp*, *ptc* (45), and the gene for  $\beta$ 3-tubulin (54) and to activate the expression of several other genes, such as *en* itself (35), *hh* (34), *Ultrabithorax* (*Ubx*) (55, 56), and *polyhomeotic* (*ph*) (13, 54), but the direct effect of En on these genes has not always been proved.

The other approaches are based on the capacity of transcription factors to bind to specific DNA sequences. Libraries have been constructed that are enriched either in sequences that bind En *in vitro* (57) or in embryonic chromatin that binds Engrailed *in vivo* (54, 58). This latter approach has also been used successfully to identify the target genes of the homeotic protein *Ubx* (59, 60). In this process, DNA-protein complexes that exist in embryos are stabilized by UV crosslinking. The complexes formed with the endogenous En protein are then selected by immunoprecipitation with an anti-En antibody, followed by cloning.

Engrailed binding sites have also been localized directly on **polytene chromosomes** of salivary glands with an anti-Engrailed antibody. The number of En binding sites was estimated to be approximately 100 (58). Among these sites, in the 2D region, the *PcG* gene *polyhomeotic* (*ph*) was shown to be a direct target of *engrailed*, through the combined use of these different genetic and molecular approaches. In particular, the activation of *ph* by En was shown to be required for the maintenance of the posterior cell fate and the position of the anterior-posterior boundary in the wing, a role assigned to the *en* function (13).

As already mentioned, among the known target genes of En, some are repressed, whereas others are activated by *engrailed*. Transfection experiments into cultured cells have demonstrated that En can act through a consensus En binding site, referred to as the NP site ("ATTAATTGA"), as a passive repressor, by competing with specific activators for binding sites located upstream of a basal promoter (61-63). Repression by En can also occur by competition with the general transcription factor TFIID for binding to the **TATA box** (64). En has also been found to function as an active

[repressor](#) of [transcription](#) by binding to sites distinct from those bound by activators (65). The active transcriptional repression can also be enhanced in the presence of corepressors like Groucho (66, 67). A **domain** other than the homeodomain was found to be responsible for repression by En (65, 68). In particular, this repressor domain is able to confer the capacity to repress transcription in transfected cultured cells, as well as *in vivo* when transferred to a heterologous DNA-binding domain (65, 68, 69).

Analysis of the molecular mechanism of transcriptional activation by En showed that its binding fragment responsible for activation is different from that responsible for repression. For activation, the binding segment consists of a long stretch of “TAAT” motifs, where up to ten En proteins can bind cooperatively (70). A strong activation is, however, obtained only in the presence of a cofactor, Extradenticle (Exd), another homeodomain-containing protein. In the presence of En, Exd becomes able to bind to this segment (71, 70). When both proteins are bound to DNA, between En and Exd turn En into an activator (70). Interestingly, activation occurs even when such an activator-binding fragment is localized far away from the coding sequence, as long as an En consensus sequence, related to NP (“ATTAATTGA”), is located just upstream of the gene to be activated.

Interactions between distantly located En binding fragments could also be responsible for the *en* autoactivation that occurs at germ-band extension (35). The observation that the autoactivation of *en* depends on *exd* maternal expression suggests that Exd is involved in this phenomenon, as it is in the case of the activation of *ph* by En (72). The fragment responsible for the *en* autoactivation, which is able to bind Exd, has not yet been isolated.

Other peculiar observations have been described concerning the *engrailed* gene. Among them is the “homing” phenomenon, in which P elements containing *en* DNA sequences are preferentially inserted at the *en* locus (8). Moreover, preferential insertions of *Drosophila* transposable elements containing a particular fragment of the *engrailed* regulatory DNA have been described (73). Finally, this regulatory DNA fragment from the *engrailed* gene was also found to mediate transvection of the *white* gene (74, 75). This silencing activity results from the presence of a polycomb response element (PRE) within the fragment of *engrailed* regulatory DNA, which was shown to bind *pleiohomeotic* (76). The presence of a PRE in this fragment could also explain the homing and the preferential insertion of *engrailed* sequences in the genome, even though this has not yet been formally proved.

## Bibliography

1. M. Ptashne and A. A. Gann (1990) *Nature* **346**, 329–331.
2. P. W. Ingham (1998) *Nature* **335**, 25–34.
3. M. J. Pankratz and H. Jäckle (1990) *Trends Genet.* **6**, 287–292.
4. T. B. Kornberg (1981) *Proc. Natl. Acad. Sci. USA* **78**, 1095–1099.
5. T. B. Kornberg (1981) *Dev. Biol.* **86**, 363–372.
6. S. J. Poole, L. M. Kauvar, B. Drees, and T. B. Kornberg (1985) *Cell* **40**, 37–43
7. J. P. Vincent and P. H. O'Farrell (1992) *Cell* **68**, 923–931.
8. C. Hama, Z. Ali, and T. B. Kornberg (1990) *Genes Dev.* **4**, 1079–1093.
9. L. A. Raftery, M. Sanicola, R. K. Blackman, and W. M. Gelbart (1991) *Development* **113**, 27–33.
10. A. Hidalgo (1994) *Curr. Biol.* **4**, 1087–1098.
11. I. Guillen, J. L. Mullor, J. Capdevila, E. Sanchez-Herrero, G. Morata, and I. Guerrero (1995). *Development* **121**, 3447–3456.
12. M. Zecca, K. Basler, and G. Struhl (1995) *Development* **121**, 2265–2278.
13. F. Maschat, N. Serrano, N. B. Randsholt, and G. Géraud (1998) *Development* **125**, 2771–2780.
14. A. Martinez-Arias and R. A. H. White (1988) *Development* **102**, 325–338.

15. J. A. Kennison and J. W. Tamkun (1988) *Proc. Natl. Acad. Sci. USA* **85**, 8136–8140.
16. T. R. Breen, V. Chinwalla, and P. J. Harte (1995) *Mech. Dev.* **52**, 89–98.
17. R. Paro (1990) *Trends Genet.* **6**, 416–421.
18. J. M. Dura and P. H. Ingham (1988) *Development* **103**, 733–741.
19. D. Moazed, D. O'Farrell, and P. H. O'Farrell (1992) *Development* **116**, 805–810.
20. J. McKeon and H. W. Brock (1991) *Wilhelm Roux's Arch. Dev. Biol.* **199**, 387–396.
21. H. Epstein (1992) *BioEssays* **14**, 411–413.
22. N. H. Patel, E. Martin-Blanco, K. G. Coleman, S. J. Poole, M. C. Ellis, T. B. Kornberg, and C. S. Goodman (1989). *Cell* **58**, 955–968 (1989).
23. J. A. McDonald and C. Q. Doe (1997) *Development* **124**, 1079–1087.
24. C. Q. Bhat and P. Schedl (1997) *Development* **124**, 1675–1688.
25. U. Schmidt-Ott, K. Sander, and G. M. Technau (1994) *Roux's Arch. Dev. Biol.* **203**, 298–303.
26. A. Fjöse, H. G. Eiken, P. R. Njolstad, A. Molven, and I. Hordvik (1988) *FEBS Lett.* **231**, 355–360.
27. D. K. Darnell, T. B. Kornberg, and C. P. Ordahl (1986) *J. Cell Biol.* **103**, 311a.
28. A. L. Joyner and G. R. Martin (1987) *Genes Dev.* **1**, 29–38.
29. A. L. Joyner, K. Herrup, B. A. Auerbach, C. A. Davis, and J. Rossant (1991) *Science* **251**, 1239–1243.
30. S. J. Poole, M. L. Law, F. Kao, and Y. Lau (1989) *Genomics* **4**, 225–231.
31. M. Peifer and A. Bejsovec (1992) *TIG* **8**, 243–249.
32. A. L. Joyner (1996) *Trends Genet.* **12**, 15–20 (1996).
33. P. S. Danielan and A. P. McMahon (1996) *Nature* **383**, 332–334.
34. T. Tabata, S. Eaton, and T. B. Kornberg (1992) *Genes Dev.* **6**, 2635–2645.
35. J. Heemskerk, S. DiNardo, R. Kostriken, and P. H. O'Farrell (1991) *Nature* **352**, 404–352.
36. P. A. Lawrence, B. Sanson, and J. P. Vincent (1996) *Development* **122**, 4095–4103.
37. J. P. Vincent and P. A. Lawrence (1994) *Cell* **77**, 909–915.
38. J. Noordermeer, J. Klingensmith, and R. Nusse (1995) *Mech. Dev.* **51**, 145–155.
39. K. F. Yoffe, A. S. Manoukian, E. L. Wilder, A. H. Brand, and N. Perrimon (1995) *Dev. Biol.* **170**, 636–650.
40. T. Tabata, C. Schwartz, E. Gustavson, Z. Ali, and T. B. Kornberg (1995) *Development* **121**, 3359–3369.
41. A. J. Simmonds, W. J. Brook, S. M. Cohen, and J. B. Bell (1995) *Nature* **376**, 424–427.
42. E. Gustavson, A. S. Goldbourough, Z. Ali, and T. B. Kornberg (1996) *Genetics* **142**, 893–906.
43. V. Marigo, R. A. Davey, Y. Zuo, J. M. Cunningham, and C. J. Tabin (1996) *Nature* **384**, 176–179.
44. M. Sanicola, J. Sekelsky, S. Elson, and W. M. Gelbart (1995) *Genetics* **139**, 745–756.
45. J. Capdevila and I Guerrero (1994) *EMBO J.* **13**, 4459–4468.
46. S. Eaton and T. B. Kornberg (1990) *Genes Dev.* **4**, 1068–1077.
47. K. Basler and G. Struhl (1994) *Nature* **368**, 208–214.
48. C. Desplan, J. Theis, and P. H. O'Farrell (1985) *Nature* **318**, 630–635.
49. C. R. Kissinger, B. Liu, E. Martin-Blanco, T. B. Kornberg, and C. O. Pabo (1990) *Cell* **63**, 579–590.
50. A. Draganescu and T. D. Tullius (1998) *J. Mol. Biol.* **276**, 529–536.
51. J. Treisman, P. Gônczy, M. Vashita, E. Harris, and C. Desplan (1989) *Cell* **59**, 553–562.
52. S. D. Hanes and R. Brent (1991) *Science* **251**, 426–430.
53. S. E. Ades and R. T. Sauer (1994) *Biochemistry* **33**, 9187–9194.

54. N. Serrano, H. W. Brock, and F. Maschat (1997) *Development* **124**, 2527–2536.
55. A. Martinez-Arias, N. E. Baker, and P. N. Ingham (1988) *Development* **103**, 157–170.
56. R. S. Mann (1994) *Development* **120**, 3205–3212.
57. M. T. Saenz-Robles, F. Maschat, T. Tabata, M. P. Scott, and T. B. Kornberg (1995) *Mech. Dev.* **53**, 185–195.
58. N. Serrano, H. W. Brock, C. Demeret, J. M. Dura, N. B. Randsholt, T. B. Kornberg, and F. Maschat (1995) *Development* **121**, 1691–1703.
59. A. P. Gould, J. J. Brookman, D. I. Strutt, and R. A. H. White (1990) *Nature* **348**, 308–312.
60. Y. Graba, D. Aragnol, P. Laurenti, V. Garzino, D. Charmot, H. Berenger, and J. Pradel (1992) *EMBO J.* **11**, 3375–3384.
61. J. B. Jaynes and P. H. O'Farrell (1988) *Nature* **336**, 744–749.
62. K. Han, M. S. Levine, and J. L. Manley (1989) *Cell* **56**, 573–583.
63. J. Treisman, E. Harris, D. Wilson, and C. Desplan (1992) *BioEssays* **14**, 145–150.
64. Y. Ohkuma, M. Horikoshi, R. G. Roeder, and C. Desplan (1990) *Proc. Natl. Acad. Sci. USA* **87**, 2289–2293.
65. J. B. Jaynes and P. H. O'Farrell (1991) *EMBO J.* **10**, 1427–1433.
66. G. Jimenez, Z. Paroush, and D. Ish-Horowicz (1997) *Genes Dev.* **11**, 3072–3082.
67. E. N. Tolkunova, M. Fujioka, M. Kobayashi, D. Deka, and J. B. Jaynes (1998) *Mol. Cell. Biol.* **18**, 2804–2814.
68. K. Han, K. Manley, and J. L. Manley (1993) *EMBO J.* **12**, 2723–2733.
69. A. John, S. T. Smith, and J. B. Jaynes (1995) *Development* **121**, 1801–1813.
70. N. Serrano and F. Maschat (1998) *EMBO J.* **17**, 3704–3713.
71. M. A. van Dijk, L. T. C. Peltenburg, and C. Murre (1995) *Mech. Dev.* **52**, 99–108.
72. J. A. Kassis, E. Noll, E. P. VanSickle, W. F. Odenwald, and N. Perrimon (1992) *Proc. Natl. Acad. Sci. USA* **89**, 1919–1923.
73. J. A. Kassis (1994) *Genetics* **136**, 1025–1038.
74. C. Rauskolb, C. Wieschaus, and E. Wieschaus (1994) *EMBO J.* **13**, 3561–3569.
75. J. A. Kassis, E. P. Vansickle, and S. M. Sensabaugh (1991) *Genetics* **128**, 751–761.
76. J. L. Brown, D. Mucci, M. Whiteley, M. L. Dirksen, and J. A. Kassis (1998) *Mol. Cell* **1**, 1057–1064.

### **Suggestions for Further Reading:**

77. P. A. Lawrence (1992) "*Engrailed: the life history of a gene*". In *The Making of a Fly: The Genetics of Animal Design*, Blackwell Scientific, London, pp. 207–210.
78. J. Treisman, E. Harris, D. Wilson, and C. Desplan (1992) The homeodomain: a new face for the helix-turn-helix? *BioEssays* **14**, 145–150.
79. T. B. Kornberg (1993) Understanding the homeodomain. *J. Biol. Chem.* **268**, 26813–26816.
80. A. Garcia-Bellido (1998) The *engrailed* story. *Genetics* **148**, 539–544.

### **Enhancer**

#### 1. Introduction/Summary



In the regulation of gene [transcription](#), two aspects are of prime importance, namely the point of initiation and the rate of transcription. In bacteria, both of these are usually determined by short DNA sequence motifs, collectively called promoters, that are in close proximity to the initiation site. In mammals and other eukaryotes, the point of transcription initiation of protein-coding genes is determined by proximal promoter sequences, but the frequency of transcription is influenced by remote DNA sequences, termed enhancers, that can be located thousands, if not tens of thousands of basepairs from the promoter. Enhancers are typically 100–300 bp long and represent an array of binding sites for DNA-binding [transcription factors](#). Binding of such proteins to an enhancer helps, in a collaboration between enhancer and promoter, to recruit components of the transcription apparatus. These include the **RNA polymerase II** holoenzyme and [enzymes](#) that loosen the [chromatin](#) packaging of DNA. Due to its specific DNA sequence, any given enhancer binds a subset of the many hundred transcription factors of a eukaryotic organism. Some of these proteins are present only in a given cell type, others become active only upon an external stimulus, such as a [hormone](#). Accordingly, an enhancer can activate a linked gene only in the appropriate cell type, or in response to a specific stimulus. Many genes are controlled by several distinct enhancers, which ensures gene activation in response to different cues. Also, sequence rearrangements in enhancers may lead to the evolution of new patterns of gene expression in multicellular organisms. Due to their prominent role in rendering chromatin more accessible to DNA-interacting proteins, at least some transcriptional enhancers have additional overlapping functions in that they facilitate [DNA replication](#), [demethylation](#) of cytosines at CpG sites, and site-specific DNA [recombination](#), switch recombination and [somatic hypermutation](#) of [immunoglobulin](#) genes. Also, some enhancers act as a genetic switch between positive and negative gene regulation; under conditions where the linked gene has to remain silent, they serve as a platform for the binding of repressing factors/cofactors.

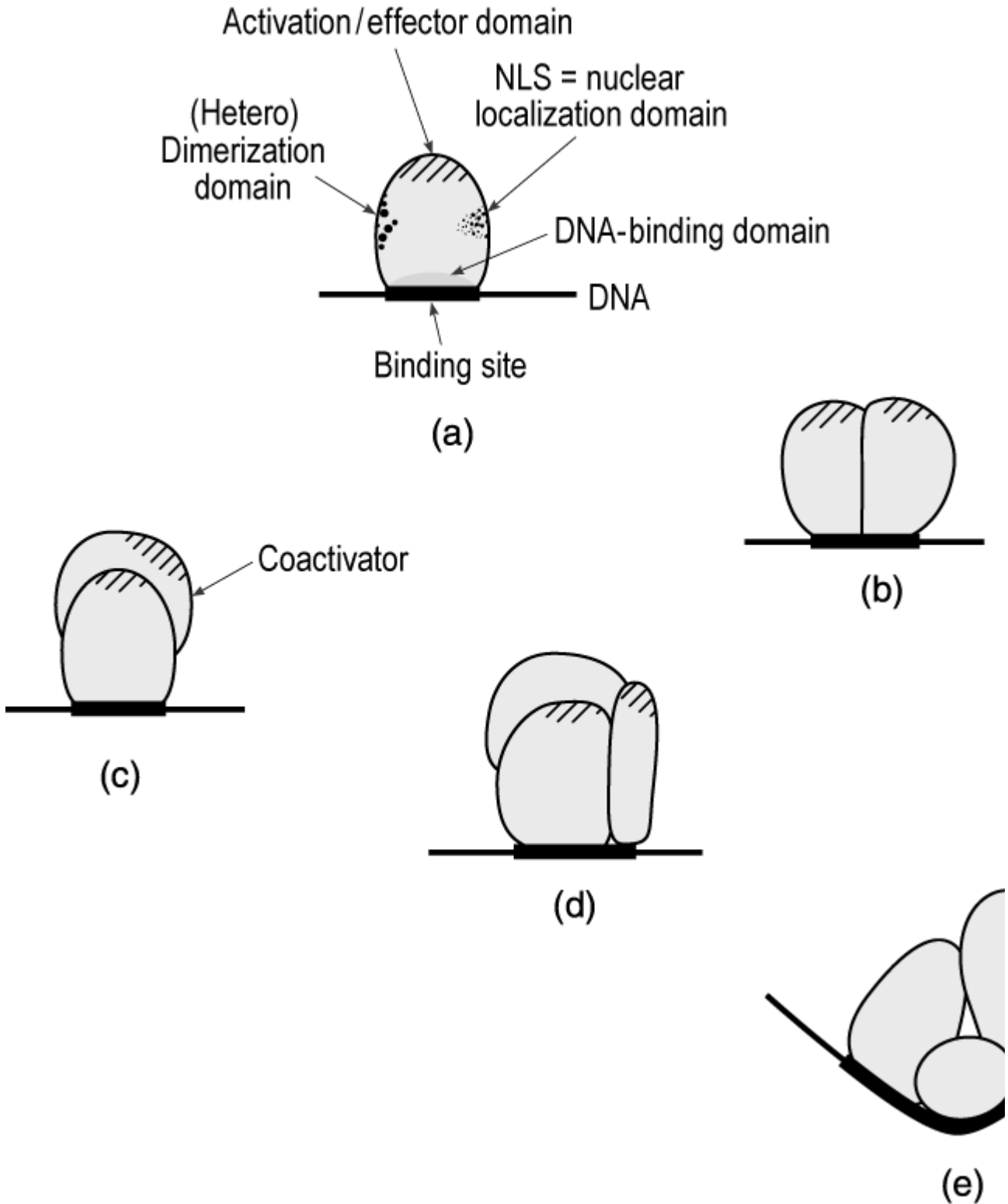
## 2. Discovery, Definition

In 1981, the first enhancer characterized was from the genome of simian virus 40 (=SV40) (1, 2). This enhancer sequence is naturally located between 300 and 100 bp upstream of the initiation site for early mRNAs, but could be experimentally placed at various distances from an unrelated [reporter gene](#), which thereby was transcriptionally activated. Characterization of enhancers from SV40 and other viruses led to the following definition of a transcription enhancer: (i) a DNA sequence, typically about 200 bp in size, that strongly activates transcription of a linked transcription unit from its correct initiation site; (ii) it activates transcription independently of its orientation; and (iii) it is able to exert its function over large distances of more than 1000 bp, from a position either upstream or downstream of the initiation site (3-5).

The first cellular enhancer was discovered in the immunoglobulin (Ig) heavy chain gene (6-8). This constituted a hallmark in the understanding of gene regulation in higher organisms, because it showed that gene control by remote enhancer sequences was not a viral peculiarity. Most importantly, the Ig enhancer stimulated transcription in a cell type-specific manner, thus being the first genetic element described to have such a property. The Ig enhancer is located downstream of the initiation ([5'-cap](#)) site, within the second [intron](#), showing that enhancers can be naturally located downstream of the site of initiation. This concept has since been widely confirmed in most eukaryotes from flies to men (Figure 1). While most of these studies were done with RNA polymerase II-transcribed genes, transcription enhancers are also associated with **ribosomal** genes that are transcribed by RNA polymerase I (9).

**Figure 1. Enhancer types in different organisms.** (a-d) Schematic representation of a mammalian gene with different enhancer locations indicated (18, 3, 4). In a gene stripped of its enhancer, transcription is low or undetectable (a). An enhancer (array of black boxes representing binding sites for transcription factors) can be located adjacent to the promoter, as in several viruses or metallothionein genes (b); or within the gene, usually as part of an intron, like in immunoglobulin (Ig) genes (c). Enhancers are also found far downstream of the transcription unit, as in T cell receptor and Ig genes (d). One and the same gene can also have several enhancers, to ensure activity in several cell types and in response to different stimuli (not shown here). In the baker's yeast, upstream activating sequences (UAS) are considered analogous to mammalian enhancer elements. They can also work in either

however, they are rarely, if at all, activating from positions downstream of the gene (e). The majority of eubacterial promoter control of RNA polymerase with sigma 70 factor and cover less than 100 bp. The promoters of genes activated upon nitrite driven by sigma 54 factor and show similarities to eukaryotic enhancers, in that an upstream binding site, in conjunction with a protein, can work from variable distances (f) (85-87).



In addition to the originally discovered intronic enhancer, the immunoglobulin heavy chain locus was also found to contain a second, extended **B cell**-specific enhancer located far downstream of the entire **locus**, as well as additional cytokine-inducible enhancers involved in Ig **class switching** (10).

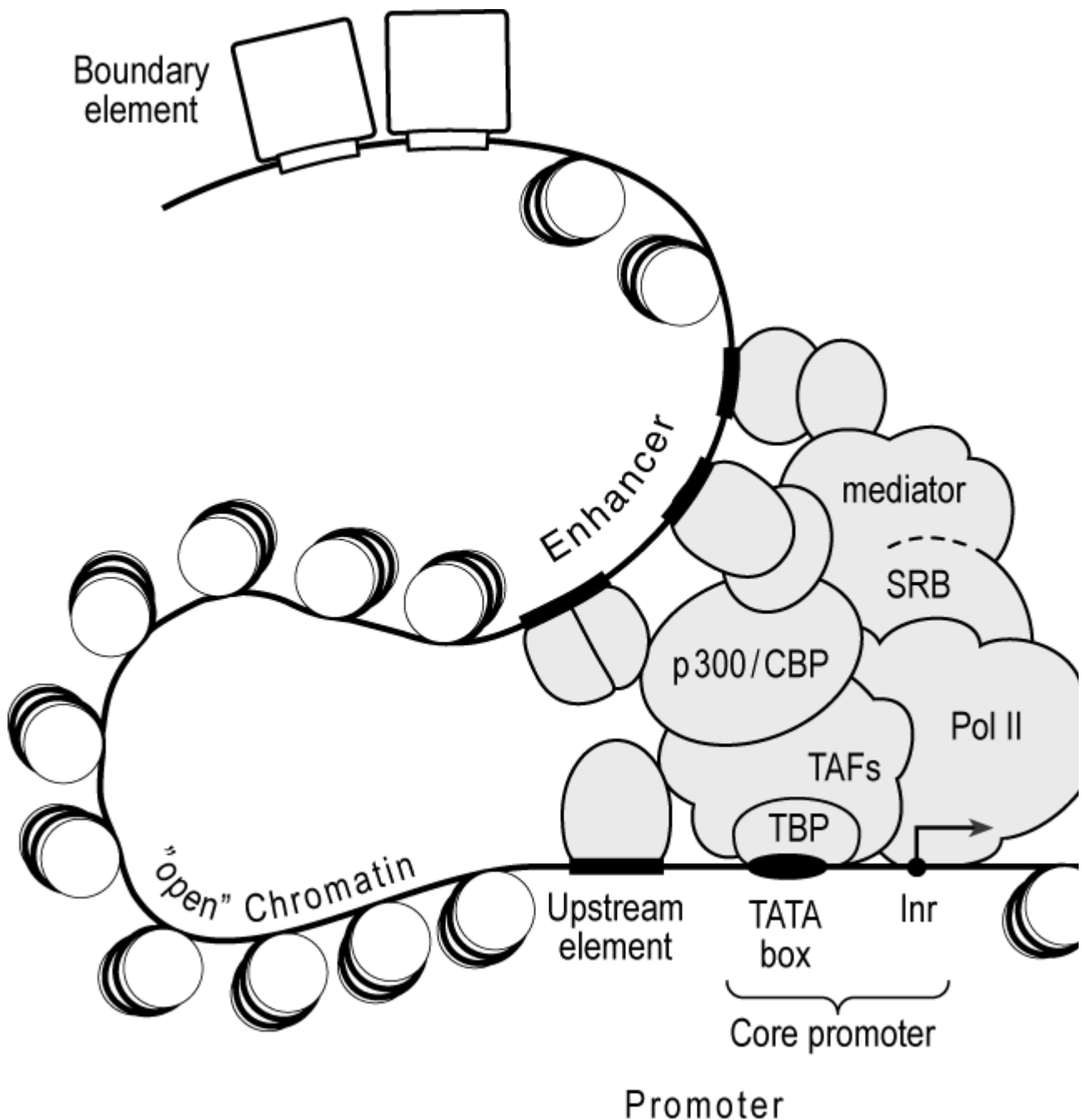
In many other instances, enhancers were established to play a crucial role in stage- and cell type-specific gene expression: In an interplay with the transcription factors that bind to them, enhancers ensure proper development of all multicellular organisms; see, for example, Refs. [11-13](#). Perhaps the most extensive studies have been done in *Drosophila*, where a gene for a developmental regulator protein is typically under the control of several separate enhancers. Depending on the stage and cell type, some enhancers can exert either positive or negative effects on transcription. For example, to establish a repetitive stripe pattern of even-skipped gene expression during early *Drosophila* embryogenesis, a single enhancer does not suffice, but one of several enhancers controls the expression of just one of these stripes ([14, 15](#); see also [16](#)).

### 3. Modular Structure of Enhancers

A typical enhancer is 100 to 300 bp long and essentially consists of an array of binding sites, each about 10 bp long, for DNA-binding transcription factors ([17, 18](#)). Every eukaryotic cell contains hundreds of different transcription factors, each with a DNA binding domain of characteristic specificity for DNA sequence recognition, plus an “activation domain”. The latter establishes contacts to other components of the transcription apparatus and thus enhances the level of gene transcription. Since under some conditions, negatively acting proteins can be recruited, an activation domain was also referred to as “effector domain”. The DNA sequence recognition specificity and affinity of many eukaryotic transcription factors does not match that of bacterial regulators. This seems to be compensated by a multitude of [protein-protein interactions](#) (see text below).

A given DNA site can, for example, be recognized by a monomer, by a homodimer reminiscent of bacterial regulators, or by heterodimers. It can also be recognized by heterodimers where only one partner binds DNA strongly, while the other makes only minor contacts to DNA and in turn may have a leading role in activation. More complex situations have also been documented (see [Figure 2](#)).

**Figure 2. Enhancer-binding transcription factors.** A typical transcription factor contains a DNA-binding domain, a signal for import into the nucleus, an activation/effector domain for contacts to coactivators/mediators and a (hetero)dimerization domain (a). Although shown here schematically, the three-dimensional structure of many transcription factors is known from X-ray or nuclear magnetic resonance (NMR) analysis. Binding to DNA can involve homo- or heterodimers (b) (e.g., see Ref. [8](#)). A transcription factor can associate with a coactivator that does not contact DNA. For example, the cell type-specific coactivator OcaB/OBF-1 binds to a ubiquitous factor (Oct-1) to bring about preferential activity in B lymphocytes (c) ([89-91](#)). A coactivator can also associate with a transcription factor to increase its DNA binding specificity and contribute a strong activation potential. For example, the herpes viral coactivator VP16 (= vFLP1) associates with the cellular Oct-1 factor and with an additional large protein termed host cell factor (HCF). By providing a strong activation potential, the virus recruits Oct-1 to a new DNA sequence motif present in viral promoters at the expense of proper host cell factor. “Architectural” proteins, that by themselves would neither activate transcription nor bind to DNA with high specificity, can also assist in the binding of *bona fide* transcription factors to DNA, as indicated here for the protein HMG I(Y) (center) that bends DNA to fit into the minor groove (e) ([20](#)); see also refs [93, 94](#). Superstructures composed of several such binding sites and DNA-bound factors are referred to as an enhanceosome ([22](#)). DNA binding transcription factors of the enhancer and the promoter, the latter including the basal transcription machinery (TBP and the general transcription factor TFIIB), recruit via multiple protein-protein contacts, further general transcription factors, adaptor- and mediator/SRB-proteins which in turn recruit RNA polymerase II to the initiation site. Individual proteins of a specific subset of factors for initiation; for example, the TBP-associated factors (TAFs) may be replaced by other multiple proteins shown here. Besides RNA polymerase, several enzymatic activities are associated with the initiation complex, including histone acetyl transferases. The latter mainly acetylate histones to ensure a loose chromatin structure permissive for transcription. Coordinated gene transcription is often flanked by boundary elements (= insulators) (f) (see text).



A typical enhancer has a modular structure and is composed of binding sites for several different transcription factors (17, 18), but synthetic enhancers can readily be obtained by multimerization of one and the same DNA sequence motif (19, 20). From these and other experiments, it is concluded that each binding site for a transcription factor (or factor dimer, see Figure 2) can be considered one enhancer unit that, in principle, is able to activate transcription by itself. However, single sites have poor, if any, activity, and multiple sites have to act synergistically to yield strong activation (19-21). More complex cases are also documented where a functional unit consists of some 60 bp of DNA, with several transcription factors stereospecifically aligned on it with the help of "architectural" proteins that bend the DNA, a superstructure that was termed the enhanceosome (22).

A DNA-binding transcription factor, whether bound to a distal enhancer or a proximal promoter, does not directly recruit RNA polymerase II, but instead acts via intermediary proteins to form large complexes (see below and Figure 2). Such multiprotein complexes seem to be quite forgiving

towards the positioning of their DNA-binding partner, which could explain the fact that an entire enhancer, subsets of it, and in many cases even individual transcription factor binding sites in an enhancer or promoter can be inverted relative to each other without loss of activity (23, 1, 17, 24). This general property of eukaryotic regulatory DNA sequences to retain activity when inverted, was first documented with an upstream segment of the [histone H2A](#) gene termed “modulator” (23).

#### 4. Mechanism, Recruitment

Transcription factors may activate transcription in more than one way, but all involve some kind of recruitment process. The activation domain of a DNA-binding transcription factor binds to additional proteins (coactivators, adaptors, mediator proteins, SRB proteins), organized in multiprotein complexes that in turn establish contacts to RNA polymerase II (25, 26). The related coactivators CBP/p300, originally found to bind the mammalian DNA-binding factor CREB that is responsive to induction by [cyclic AMP](#), contain several domains for interaction with a multitude of other transcription factors (27). Transcriptional activators are also known to contact the general transcription factor TFIIB, the [TATA-box](#) binding protein (TBP) and/or TAF proteins associated with it (28, 29).

The recruitment concept for transcriptional activation has gained credit from elegant experiments where covalent linkage of a DNA-binding transcription factor to a basal transcription factor (the TATA-binding protein TBP) rendered the *bona fide* activation domain obsolete. The same was found when a DNA-binding factor was designed to make an artificial, tight contact to the SRB complex. These findings demonstrate that linking a DNA-bound transcription factor to a component of the mediator complex is in principle sufficient to activate transcription (30-32).

The transcription factors of the enhancer are likely to bind, indirectly via cofactors, to those of the promoter to form an initiation complex, with the intervening DNA looped out. Several experiments support a looping between enhancer and promoter (eg, see Refs. 33 and 34). In spite of the compelling simplicity of the above experiments, it is unlikely that transcriptional activation is naturally brought about by a uniform mechanism. Among other evidence, this is suggested by the fact that there are different types of activation domains, when classified according to their activation properties (35, 36), and it is obvious that transcription can be controlled not only at the level of initiation but also at the level of elongation (37-39).

The coactivator and mediator protein complexes also contain several [histone acetylase](#) enzymes. Their major targets are the histones, but a number of transcription factors can also become reversibly acetylated. Acetylation of histones generally loosens chromatin structure and thus makes it more accessible to the transcription apparatus (40). Major support for this concept comes from experimental digestion of chromatin with DNaseI (41). Transcriptionally active DNA is generally DNase I sensitive, which is correlated with the presence of acetylated histones. In addition, the enhancer and promoter regions themselves are hypersensitive, due to binding of transcription factors with, at least in some cases, concomitant removal of [nucleosomes](#) (42). (Besides a genetic screen called “enhancer trap” (43, 44), the property of DNaseI hypersensitivity is the most effective method in identifying candidate enhancers from large DNA segments). Large multiprotein complexes (SRB/mediator; SWI/SNF etc.) are known to support chromatin remodelling and recruitment of the transcription apparatus (45, 25-48). Extended regions of inactive chromatin, termed heterochromatin, tend to be associated with the periphery of the nucleus (49). By fluorescence staining it was shown that a silent, peripheral gene relocalizes to a more central region upon activation of its enhancer. The general validity of this phenomenon remains to be seen, and also whether such a relocalization is the cause or the effect of transcriptional activation.

#### 5. Evolution

In higher eukaryotes, genes regulating development are often under the control of multiple enhancers, and it was argued that enhancer shuffling was the major driving force in the evolution of

new patterns of gene expression (50-53). Apart from shuffling, enhancer evolution may be facilitated by the fact that enhancer DNA sequences are more tolerant to rearrangements than protein coding sequences, which has, e.g., resulted in divergent primary enhancer sequences that nevertheless serve the same function (54). In fast-evolving viruses, a new cell type specificity is readily achieved by enhancer rearrangements, usually involving deletion-duplication of subsegments, with little alterations elsewhere in the genome (55). Duplications of enhancer subsegments, which are indicative of recent functional alterations, are also found in some cellular enhancers and in humans may even vary between ethnic groups (56). However, not every arbitrary enhancer-promoter combination can be expected to function since in any gene, promoter and enhancer(s) have co-evolved for optimal function. There are indeed documented cases of preferential, if not exclusive enhancer/promoter specificity (57-59). Nevertheless, promiscuous enhancer-promoter cooperations were found from yeast to mammals and may be the rule, rather than the exception (44, 50, 60).

## 6. Open Questions/Related Phenomena

Even though many aspects of eukaryotic transcription control have already been elucidated, thanks to a combination of genetics, reverse genetics and biochemistry, many questions remain. One of them concerns the mode of activation: does an enhancer allow for fine tuning of the activity of a gene within one individual cell, or can a gene only exist in two states, namely fully activated or completely switched off? In the latter case, the function of the enhancer would be to increase the probability of the gene being in the “on” state. While the first of these possibilities seems to make more sense, at least in the case of [hemoglobin](#) genes, there is evidence for the latter one (61).

Some extended enhancer regions have been referred to as LCR (locus control region). The most thorough studies were done with the hemoglobin beta chain locus, where the LCR consists of a group of at least five remote enhancers/DNase I hypersensitive regions contained within a DNA segment of 15 kb. Each LCR subdomain confers distinct activation properties in erythroid cells. The LCR, which is located some 50 kb from the most downstream beta-globin transcription unit, ensures the complete, stage-specific and cell type-specific expression pattern of the globin genes under its influence (62-64).

In spite of an ability to work over long distances, enhancer activity is confined to the same DNA molecule, i.e. to an activity in *cis*. In certain constellations that requires pairing of homologous chromosomes, an enhancer can also activate the promoter in the homologous locus, a phenomenon termed **transvection** (65). While well-documented in *Drosophila*, the relevance of transvection in mammals remains to be seen.

Perhaps due to its ability to loosen up chromatin structure, a transcription enhancer may be used for other activities besides transcriptional activation. For example, the activity of an enhancer *in cis* was also correlated with the initiation of DNA replication (66-68). Moreover, in immunoglobulin (Ig) genes, enhancers induce demethylation at CpG sites (70), site-specific [gene rearrangement](#) (VDJ joining) and also class switch recombination (71). In addition, the amazing process of localized Ig gene hypermutation in B memory cells is dependent on an active transcription enhancer (72). The Ig gene locus also harbors so-called [matrix attachment regions](#) (MARs), A+T-rich DNA segments which have little activity by themselves but in transgenic animals support long-range enhancer activity, apparently by counteracting the formation of transcriptionally inactive chromatin (73).

In a process opposite to activation, genes can be inactivated by recruitment of histone deacetylases and other proteins counteracting transcription to the promoter/enhancer region (74, 75). Non-acetylated histones tend to assume a chromatin structure that is refractory to the transcription apparatus, which is also reflected by a relative resistance to experimental DNase I digestion. A so-called [silencer](#) DNA segment can inactivate a linked gene over large distances (76-78). Again, for this purpose, DNA binding proteins recruit a number of cofactors which together block transcription (74, 75, 78, 79). Enhancers and silencers originally were described as non-overlapping DNA entities of opposite function. However, a number of enhancers, notably for developmental gene regulation in

*Drosophila*, also harbor binding sites for negatively acting DNA binding proteins and/or cofactors. Thus, in a cell type where the gene must not be active, an enhancer may serve a silencing purpose (80). To prevent undue spreading of an enhancer (or silencer) effect to neighboring genes, chromosomal domains of gene activity are separated from each other by so-called boundary, or insulator, sequences (81, 82)(Fig. 2).

Finally, in the process of [RNA splicing](#), a process that bears some resemblance to transcriptional enhancement has been described. So-called splice enhancers denote RNA sequences in the primary transcript which, by means of specific protein binding, facilitate RNA processing at nearby splice sites. Unlike transcription enhancers, splice enhancers are typically located in an exon (exonic splice enhancer = ESE), consist of short sequence motifs that cannot be inverted, and act over short distances only. This process is mainly used to favor one type of splice site at the expense of another, i.e., for [alternative splicing](#) (83, 84).

## Bibliography

“Enhancer” in , Vol. 2, pp. 823–828, by Walter Schaffner, Universitaet Zurich-Irchel, Zurich, Switzerland; “Enhancer” in (online), posting date: January 15, 2002, by Walter Schaffner, Universitaet Zurich-Irchel, Zurich, Switzerland.

1. J. Banerji, S. Rusconi and W. Schaffner (1981) *Cell* **27**, 299–308.
2. P. Moreau, R. Hen, B. Wasylyk, R. Everett, M. P. Gaub and P. Chambon (1981) *Nucleic Acids Research* **9**, 6047–6068.
3. M. M. Müller, T. Gerster and W. Schaffner (1988) *European Journal of Biochemistry* **176**, 485–495.
4. E. M. Blackwood and J. T. Kadonaga (1998) *Science* **281**, 60–63.
5. D. Dorsett (1999) *Current Opinion in Genetics & Development* **9**, 505–514.
6. J. Banerji, L. Olson and W. Schaffner (1983) *Cell* **33**, 729–740.
7. S. D. Gillies, S. L. Morrison, V. T. Oi and S. Tonegawa (1983) *Cell* **33**, 717–728.
8. M. S. Neuberger (1983) *The EMBO J.* **2**, 1373–1378.
9. R. H. Reeder (1990) *Trends in Genetics* **6**, 390–395.
10. V. Arulampalam, L. Eckhardt and S. Pettersson (1997) *Immunology Today* **18**, 549–554.
11. A. B. Firulli and E. N. Olson (1997) *Trends in Genetics* **13**, 364–369.
12. A. Gould, A. Morrison, G. Sproat, R. A. White and R. Krumlauf (1997) *Genes & Development* **11**, 900–913.
13. R.S. Mann (1997) *Bioessays* **19**, 661–664.
14. T. Goto, P. Macdonald and T. Maniatis (1989) *Cell* **57**, 413–422.
15. K. Harding, T. Hoey, R. Warrior and M. Levine (1989) *The EMBO Journal* **8**, 1205–1212.
16. J. Zhou and M. Levine (1999) *Cell* **99**, 567–575.
17. M. Zenke, T. Grundstrom, H. Matthes, M. Wintzerith, C. Schatz, A. Wildeman and P. Chambon (1986) *The EMBO Journal* **5**, 387–397.
18. E. Serfling, M. Jasin and W. Schaffner (1985) *Trends in Genetics* **1**, 224–230.
19. B. Ondek, A. Shepard and W. Herr (1987) *The EMBO Journal* **6**, 1017–1025.
20. M. D. Schatt, S. Rusconi and W. Schaffner (1990) *The EMBO Journal* **9**, 481–487.
21. M. Carey, Y. S. Lin, M. R. Green and M. Ptashne (1990) *Nature* **345**, 361–364.
22. T. Agalioti, S. Lomvardas, B. Parekh, J. Yie, T. Maniatis and D. Thanos (2000) *Cell* **103**, 667–678.
23. R. Grosschedl and M. L. Birnstiel (1980) *Proceedings of the National Academy of Sciences of the United States of America* **77**, 7102–7106.
24. L. Xu, M. Thali and W. Schaffner (1991) *Nucleic Acids Research* **19**, 6699–6704.
25. M. Carlson (1997) *Annual Review of Cell & Developmental Biology* **13**, 1–23.

26. L. C. Myers, C. M. Gustafsson, D. A. Bushnell, M. Lui, H. Erdjument-Bromage, P. Tempst and R. D. Kornberg (1998) *Genes & Development* **12**, 45–54.
27. R.H. Goodman and S. Smolik (2000) *Genes & Development* **14**, 1553–1577.
28. A. Hoffmann, T. Oelgeschlager and R. G. Roeder (1997) *Proceedings of the National Academy of Sciences of the United States of America* **94**, 8928–8935.
29. S.R. Albright and R. Tjian (2000) *Gene* **242**, 1–13.
30. A. Barberis, J. Pearlberg, N. Simkovich, S. Farrell, P. Reinagel, C. Bamdad, G. Sigal and M. Ptashne (1995) *Cell* **81**, 359–368.
31. S. Chatterjee and K. Struhl (1995) *Nature* **374**, 820–822.
32. S. Farrell, N. Simkovich, Y. Wu, A. Barberis and M. Ptashne (1996) *Genes & Development* **10**, 2359–2367.
33. H. P. Mueller-Sturm, J. M. Sogo and W. Schaffner (1989) *Cell* **58**, 767–777.
34. N. Dillon and F. Grosveld (1993) *Trends in Genetics* **9**, 134–137.
35. K. Seipel, O. Georgiev and W. Schaffner (1992) *The EMBO Journal* **11**, 4961–4968.
36. J. Blau, H. Xiao, S. McCracken, O. H. P, J. Greenblatt and D. Bentley (1996) *Molecular & Cellular Biology* **16**, 2044–2055.
37. D. L. Bentley (1995) *Current Opinion in Genetics & Development* **5**, 210–216.
38. S. A. Brown, C. S. Weirich, E. M. Newton and R. E. Kingston (1998) *The EMBO Journal* **17**, 3146–3154.
39. T. Wada, T. Takagi, Y. Yamaguchi, D. Watanabe and H. Handa (1998) *The EMBO Journal* **17**, 7395–7403.
40. J.C. Rice and C.D. Allis (2001) *Current Opinion in Cell Biology* **13**, 263–273.
41. H. Weintraub and M. Groudine (1976) *Science* **193**, 848–856.
42. L. Gaudreau, A. Schmid, D. Blaschke, M. Ptashne and W. Horz (1997) *Cell* **89**, 55–62.
43. F. Weber, J. de Villiers and W. Schaffner (1984) *Cell* **36**, 983–992.
44. H. J. Bellen, O. K. CJ, C. Wilson, U. Grossniklaus, R. K. Pearson and W. J. Gehring (1989) *Genes & Development* **3**, 1288–1300.
45. C. J. Wilson, D. M. Chao, A. N. Imbalzano, G. R. Schnitzler, R. E. Kingston and R. A. Young (1996) *Cell* **84**, 235–244.
46. C. Wu (1997) *Journal of Biological Chemistry* **272**, 28171–28174.
47. R.D. Kornberg and Y. Lorch (1999) *Current Opinion in Genetics & Development* **9**, 148–151.
48. C.L. Peterson and J.L. Workman (2000) *Current Opinion in Genetics & Development* **10**, 187–192.
49. M. Cockell and S.M. Gasser (1999). *Current Opinion in Genetics & Development* **9**, 199–205.
50. M. Kermekchiev, M. Pettersson, P. Matthias and W. Schaffner (1991) *Gene Expression* **1**, 71–81.
51. X. Li and M. Noll (1994) *Nature* **367**, 83–87.
52. D. Tautz (2000) *Current Opinion in Genetics & Development* **10**, 575–579.
53. B.G. Magor, D.A. Ross, L. Pilström and G.W. Warr (1999) *Immunology Today* **20**, 13–17.
54. M.Z. Ludwig, C. Bergman, N.H. Patel and M. Kreitman (2000) *Nature* **403**, 564–567.
55. S. Schirm, J. Jiricny and W. Schaffner (1987) *Genes and Development* **1**, 65–74.
56. S. Marsh, E.S. Collie-Duguid, T. Li, X. Liu, and H.L. McLeod (1999) *Genomics* **15**, 310–312.
57. X. Li and M. Noll (1994) *EMBO Journal* **13**, 400–406.
58. G. Das, C. S. Hinkley and W. Herr (1995) *Nature* **374**, 657–659.
59. S. Ohtsuki, M. Levine and H. N. Cai (1998) *Genes & Development* **12**, 547–556.
60. G. Schaffner, S. Schirm, B. Müller, F. Weber and W. Schaffner (1988) *Journal of Molecular Biology* **201**, 81–90.



61. M. Wijgerde, F. Grosveld and P. Fraser (1995) *Nature* **377**, 209–213.
62. F. Grosveld, G.B. van Assendelft, D.R. Greaves and G. Kollias (1987) *Cell* **51**, 975–985.
63. M. Bulger and M. Groudine (1999) *Genes & Development* **13**, 2465–2477.
64. J.D. Engel and K. Tanimoto (2000) *Cell* **100**, 499–502.
65. V. Pirrotta (1999) *Biochimica et Biophysica Acta* **1424**, M1–M8.
66. J. de Villiers, W. Schaffner, C. Tyndall, S. Lupton and R. Kamen (1984) *Nature* **312**, 242–246.
67. M. L. DePamphilis (1993) *Trends in Cell Biology* **3**, 161–167.
68. P. Van der Vliet (1996) Cold Spring Harbor Laboratory Press 87–118.
69. R. Li, D. S. Yu, M. Tanaka, L. Zheng, S. L. Berger and B. Stillman (1998) *Molecular and Cellular Biology* **18**, 1296–1302.
70. Y. Bergman and R. Mostoslavsky (1998) *Biological Chemistry* **379**, 401–407.
71. E. Sakai, A. Bottaro and F.W. Alt (1999) *Int. Immunology* **11**, 1709–1713.
72. A. G. Betz, C. Milstein, A. Gonzalez-Fernandez, R. Pannell, T. Larson and M. S. Neuberger (1994) *Cell* **77**, 239–248.
73. W.C. Forrester, L.A. Fernandez and R. Grosschedl (1999) *Genes & Development* **13**, 3003–3014.
74. P.A. Wade (2001) *Human Molecular Genetics* **10**, 693–698.
75. A. El-Osta and A.P. Wolffe (2000) *Gene Expression* **9**, 63–75.
76. A. H. Brand, L. Breeden, J. Abraham, R. Sternglanz and K. Nasmyth (1985) *Cell* **41**, 41–48.
77. H. N. Cai, D. N. Arnosti and M. Levine (1996) *Proceedings of the National Academy of Sciences of the United States of America* **93**, 9309–9314.
78. G. Chen and A.J. Courey (2000) *Gene* **249**, 1–16.
79. R. E. Kingston, C. A. Bunker and A. N. Imbalzano (1996) *Genes & Development* **10**, 905–920.
80. S. Gray and M. Levine (1996) *Current Opinion in Cell Biology* **8**, 358–364.
81. R. Kellum and P. Schedl (1991) *Cell* **64**, 941–950.
82. A.C. Bell and G. Felsenfeld (1999) *Current Opinion in Genetics & Development* **9**, 191–198.
83. K. J. Hertel, K. W. Lynch and T. Maniatis (1997) *Current Opinion in Cell Biology* **9**, 350–357.
84. B.J. Blencowe (2000) *Trends in Biochemical Sciences* **25**, 106–110.
85. B. Magasanik (1988) *Trends in Biochemical Sciences* **13**, 475–479.
86. J. D. Gralla (1996) *Current Opinion in Genetics & Development* **6**, 526–530.
87. B.R. Belitsky and A.L. Sonenshein (1999) *Proceedings of the National Academy of Sciences of the United States of America* **96**, 10290–10295.
88. M. Beato, S. Chavez and M. Truss (1996) *Steroids* **61**, 240–251.
89. M. Gstaiger, L. Knoepfel, O. Georgiev, W. Schaffner and C. M. Hovens (1995) *Nature* **373**, 360–362.
90. Y. Luo and R. G. Roeder (1995) *Molecular & Cellular Biology* **15**, 4115–4124.
91. M. Strubin, J. W. Newell and P. Matthias (1995) *Cell* **80**, 497–506.
92. A. C. Wilson, K. LaMarco, M. G. Peterson and W. Herr (1993) *Cell* **74**, 115–125.
93. R. Grosschedl (1995) *Current Opinion in Cell Biology* **7**, 362–370.
94. C. A. Bewley, A. M. Gronenborn and G. M. Clore (1998) *Annual Review of Biophysics & Biomolecular Structure* **27**, 105–131.

## Enkephalins

Enkephalins are pentapeptides with the sequence Tyr–Gly–Gly–Phe–X, which are the natural ligands of the opiate **receptors**. The two major enkephalins are [5-leucine] enkephalin (abbreviated [Leu] enkephalin) and [5-methionine] enkephalin (abbreviated [Met] enkephalin) (1). The peptides are derived from three precursors, pro-opiomelanocortin (POMC) and pro-enkephalins A and B (2). The prohormones are of similar size and contain multiple repeated units of enkephalins. The precursors are processed **proteolytically**, mostly in the [Golgi apparatus](#), to the hormonally active products. The two proenkephalins show considerable sequence [homology](#); they probably originated from a common ancestral gene by [gene duplication](#).

The proenkephalin peptides and the enkephalins are distributed in many types of cells, especially brain, spinal cord, and gut. POMC is present in high concentrations in the corticotropic cells of the adenohypophysis and in many endocrine-related tissues (eg, adrenal medulla). Enkephalins have analgesic properties. They are also synthesized in neurons and are abundant in various regions of the hypothalamus, acting on the secretion of hormones of the hypothalamus and hypophysis (they stimulate secretion of prolactin, thyrotropin, and [growth hormone](#) and inhibit secretion of follicle-stimulating hormone/luteinizing hormone (folitropin/lutropin; FSH/LH), adrenocorticotrophic hormone (ACTH), and vasopressin/oxytocin (3). The enkephalins are hydrolyzed by enkephalinases. Inhibitors of the degrading enzymes have been introduced as analgesics (4).

#### Bibliography

1. J. Hughes et al. (1975) *Nature* **258**, 577–579.
2. J. Douglass, O. Civelli, and E. Herbert (1984) *Annu. Rev. Biochem.* **53**, 665–715.
3. M. Ferin, D. Van Vugt, and S. Wardlaw (1984) *Recent Prog. Horm. Res.* **40**, 441–485.
4. J. M. Lecomte et al. (1986) In *Innovative Approaches in Drug Research* (A. F. Harms, ed.), Elsevier, Amsterdam, pp. 315–329.

#### Enterokinase

Enterokinase, also known as enteroprotease (E.C. 3.4.21.9), is a key [enzyme](#) in the digestion of dietary proteins in the small intestine. It catalyzes the activation of **trypsinogen** by removing six amino acids from the amino terminus of the 229-residue protein and does so with extremely high specificity. The resulting trypsin then activates all the other pancreatic proenzymes, as part of a two-step proteolytic cascade that is essential for proper digestion.

Enterokinase is a [serine proteinase](#) of the **trypsin** family and with trypsin-like specificity (1). It is synthesized in the epithelial cells of the duodenum and appears to be an intrinsic [membrane protein](#) located on the outer surface of the intestinal brush border. In fact, it is generated from a 1019-residue precursor that is proteolytically cleaved to a 235-residue light chain, which resembles trypsin and has the serine proteinase [catalytic triad](#), plus a 784-residue heavy chain (2); the two chains are linked by a [disulfide bond](#) (3). In this regard, its properties are like those of the [blood-clotting](#) enzymes (4).

It should be noted that the term [kinase](#) is normally used to denote an enzyme that catalyzes the transfer of a phosphate group from a donor, such as ATP, to a receptor, such as glucose. In this regard, enterokinase has been misnamed. It does not catalyze phosphate transfer at all and, hence, is more appropriately called *enteropeptidase*.

## Bibliography

1. A. Light and H. Janska (1989) Trends Biochem. Sci. **14**, 110–112.
2. Y. Kitamoto, R. A. Viele, H. Donis-Keller, and J. E. Sadler (1995) Biochemistry **34**, 4562–4568.
3. A. Light and P. Fonseca (1984) J. Biol. Chem. **259**, 13195–13198.
4. L. E. Anderson, K. A. Walsh, and H. Neurath (1977) Biochemistry **16**, 3354–3360.

## Enzyme Immobilization And Conjugation

Immobilization and conjugation refers to the covalent attachment of an [enzyme](#) to a solid support or to a soluble molecule through use of [bifunctional crosslinking reagents](#) (see [Crosslinking](#)). The use of immobilized enzymes in the medical, pharmaceutical, chemical, and food industries continues to increase dramatically. Future exploitation of immobilized or conjugated enzymes will undoubtedly result in the development of new biosensors and bioreactors for use in such areas as diagnostics (*in vivo* monitoring of metabolites, drugs, and proteins); environmental monitoring (detection of pathogens, toxins, and pollutants); separation sciences; the synthesis of complex, chiral organic compounds; and the detection of macromolecular interactions (reviewed in Refs. [1](#) and [2](#)).

Among the most widely used bioconjugates are reporter enzymes covalently crosslinked to [antibodies](#). These conjugates have been pivotal in the development of [enzyme-linked immunosorbent assay](#) (ELISA) systems, which allow detection of an immense variety of analytes through their specific recognition by appropriate antibodies. The two most commonly used reporter enzymes in ELISAs are [alkaline phosphatase](#) and horseradish [peroxidase](#) (HRP). The characteristics of these enzymes allow for different strategies in their conjugation to antibodies. For example, HRP contains saccharide groups that can be readily oxidized to aldehydes, which, in turn, are conjugation targets for free [amino groups](#) on the antibody. Alternatively, because HRP contains few **lysine residues**, it can be conjugated to antibodies by [glutaraldehyde](#) without significant formation of insoluble products. Several heterobifunctional crosslinking reagents, including *N*-[hydroxysuccinimide](#) ester/maleimide reagents, which selectively link **amino** and [thiol groups](#), are also widely used in enzyme–antibody conjugation ([3](#)).

Enzymes can be immobilized by conjugation to a variety of different types of supports having different properties and uses. For instance, in aqueous solution [polyethylene glycols](#) have a large exclusion volume, which encompasses those molecules conjugated to them. Thus, [proteinases](#) and antibody molecules can be excluded from enzymes conjugated to these polymers, a property that can be beneficial for applications that require exposure of the conjugates to biological fluids ([4](#)).

Reversible solubilization of an enzyme can be achieved through its conjugation to [liposomes](#), using a reagent such as a **carbodiimide**; reversible precipitation and solubilization can be brought about by manipulation of the dielectric strength of the solution. Biosensor technology utilizes proteins that are immobilized by bifunctional crosslinking reagents to diverse supports ([5](#)), including carboxymethylated dextran-coated gold film, metal-coated nylon mesh, and carboxymethyl cellulose. Immobilization of enzymes on novel solid-phase matrices, such as porous zirconium, silicates, and colloidal gold, has been used in the construction of bioreactors for the industrial generation of numerous compounds. Many supports have functional groups that must be chemically activated prior to conjugation with protein; this allows activation to be carried out under harsh conditions that would otherwise be detrimental to the protein.

## Bibliography

1. E. Katchalski-Katzir (1993) *Trends Biotechnol.* **11**, 471–478.
2. F. Svec and J. M. Frecht (1996) *Science* **273**, 205–211.
3. G. T. Hermanson (1996) *Bioconjugate Techniques*, Academic Press, San Diego, pp. 461–469.
4. H. C. Berger and S. V. Pizzo (1988) *Blood* **71**, 1641–1647.
5. G. T. Hermanson (1996) *Bioconjugate Techniques*, Academic Press, San Diego, pp. 593–629.

### Suggestions for Further Reading

6. S. S. Wong (1993) *Chemistry of Protein Conjugation and Cross-Linking*, CRC Press, Boca Raton, FL. (An outstanding book describing all aspects of the chemistry and uses of bifunctional reagents in crosslinking and conjugation, with an extensive bibliography for each chapter.)
7. G. T. Hermanson (1996) *Bioconjugate Techniques*, Academic Press, San Diego. (An excellent book covering all phases of conjugation, and not limited to proteins.)
8. *Bioconjugate Chemistry* (ISSN 1043-1802), Editor-in-Chief, C. F. Meares. (A bimonthly journal published by the American Chemical Society devoted to scientific reports describing the “joining of two molecular functions by chemical or biological means.”)

## Enzyme-Linked Immunosorbent Assay (ELISA)

Enzyme-linked immunosorbent assay (ELISA) is a technique for quantifying a substance by use of an [antibody](#) specific for that substance that has been conjugated to an [enzyme](#). The original ELISA was developed along the lines of a solid-phase [radioimmunoassay](#), substituting an enzyme for a **radioactive** label ([1](#), [2](#)). Regardless of the experimental design used for the immunological reactions, the end result is that the antibody–enzyme conjugate is immobilized. Enzymatic activity is assayed in a subsequent reaction, and the activity is related to the quantity of analyte ([antigen](#) or antibody) present in a test sample by comparison to a set of controls of known concentration.

A key reagent for all ELISAs is the enzyme conjugate. In most cases, the enzyme is [crosslinked](#) covalently to an antibody, but experimental designs with enzyme-labeled antigen are also used. Indeed, the ELISA method is adaptable to any **ligand** and **receptor** pair. The most common enzymes used are [alkaline phosphatase](#) and horseradish [peroxidase](#). These proteins are stable to the crosslinking procedures, and they catalyze reactions with highly [chromogenic substrates](#), leading to very sensitive detection limits. Optical absorbance is the most common method of measuring enzyme activity, but other detection technologies are also used, such as **fluorescence** or [chemiluminescence](#). Covalent linkage to a [chromatography](#) matrix was used to immobilize the first “immunosorbent,” but use of plastic multiwell microtiter plates is now nearly universal, as these allow simultaneous monitoring of scores or hundreds of enzymic reactions in automated instruments.

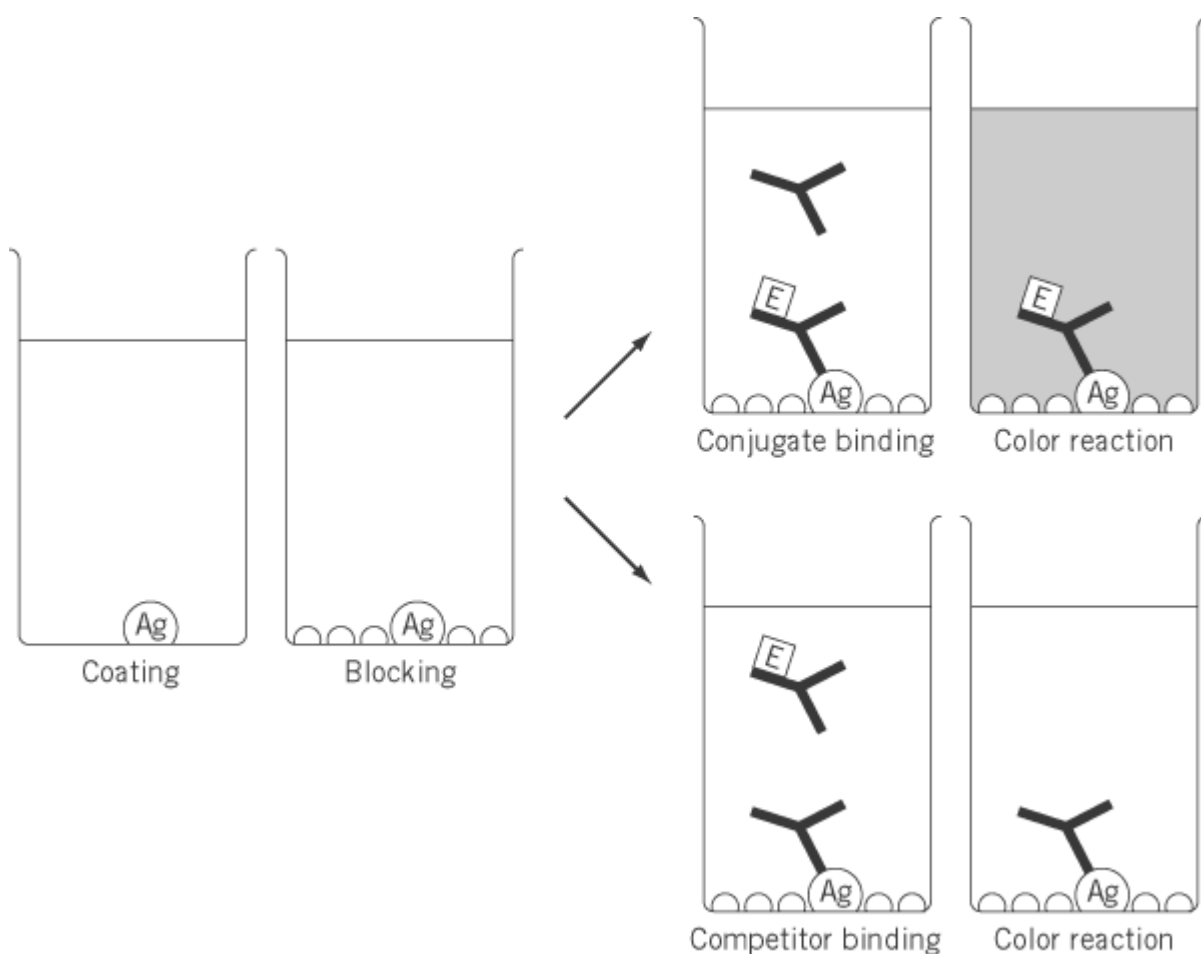
The number of experimental formats used for ELISA is immense. This discussion is limited to the simplest and most useful assay designs.

### 1. Competitive ELISA

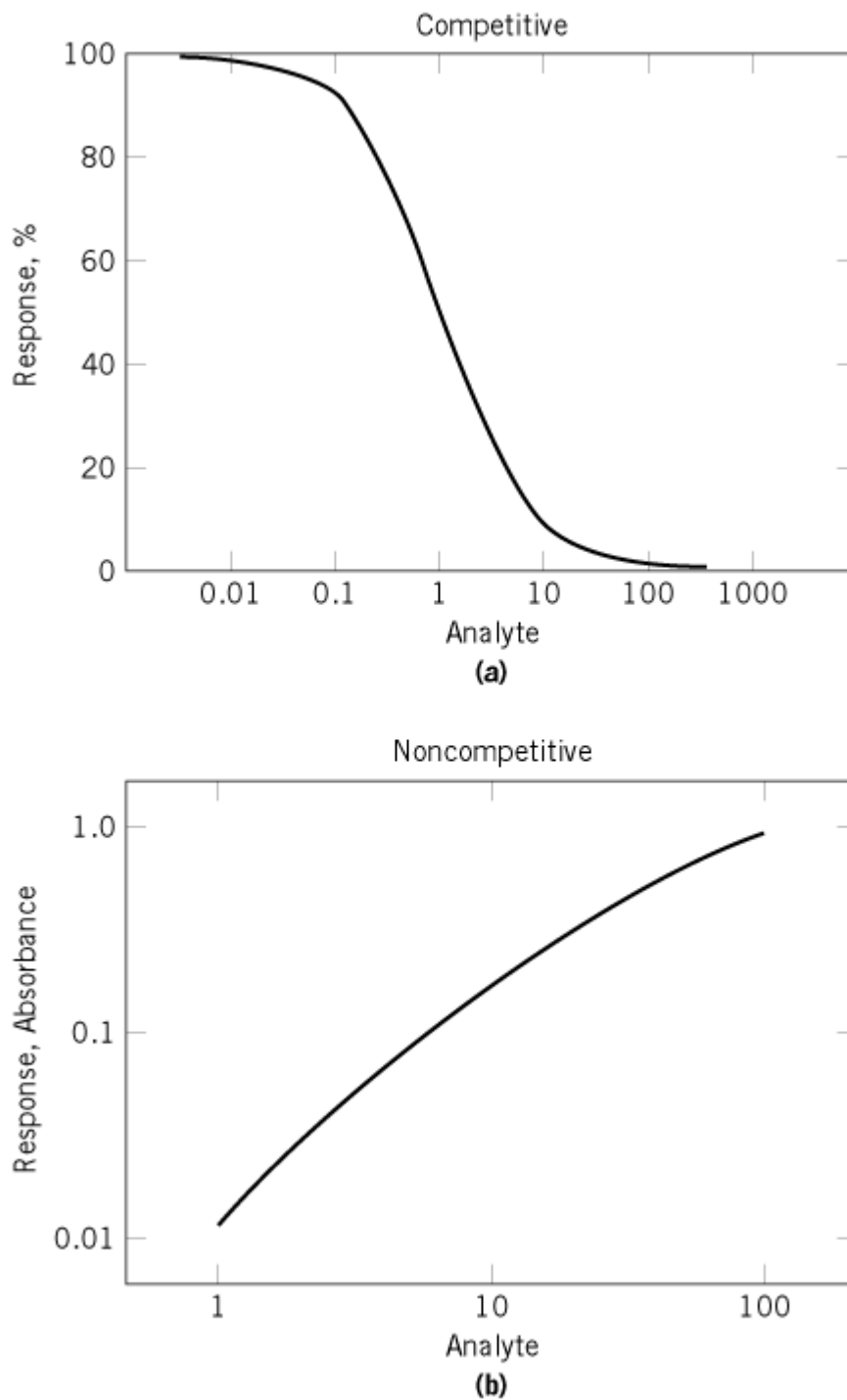
The first reported ELISA was a competition assay, in which adding an unlabeled competitor molecule reduced the amount of enzyme conjugate that could be immobilized ([1](#)). This remains an extremely common approach. The method begins with adsorption of antigen to the surface of a microtiter plate (Fig. [1](#)). This non-specific adsorption occurs spontaneously, simply by allowing the

antigen solution to stand in the wells of the plate. The linkage to the plate is noncovalent, but essentially irreversible, and the adsorbed antigen is seldom released during subsequent reactions and washing steps. Often a blocking step follows, in which a concentrated solution of a protein such as [serum albumin](#) is used to coat the remaining plastic surface, thereby preventing later adsorption of assay components. The plate is washed with a neutral buffer. The antibody solution of unknown concentration, and a series of controls, consisting of known quantities of the same antibody, are mixed with a fixed amount of the antibody–enzyme conjugate and added to different wells of the plate. The conjugate binds to the immobilized antigen, but so does the unconjugated antibody, and the binding of each is mutually exclusive. Thus, the higher the concentration of unconjugated antibody, the less the amount of enzyme immobilized. After allowing an interval for the antigen to be bound, the plate is again washed, and the addition of substrate initiates the enzymatic reaction. The rate of product release, or the accumulation of product after a fixed time, is determined for the unknown analyte and controls. Reactions to which no unlabeled competitor was added will have the highest absorbance, and are taken as a 100% response. Controls lacking conjugate, or lacking antigen, or to which a saturating amount of competitor was added, give the lowest absorbance. This value, reflecting background enzyme binding and the uncatalyzed chromogenic reaction, is taken as a 0% response. Data from the other controls are plotted to give a calibration curve, using a logarithmic concentration scale on the abscissa and an absorbance or percentage scale on the ordinate (Fig. 2a). The response from the unknown is compared to the calibration curve, and the corresponding concentration is read off the abscissa.

**Figure 1.** Competitive ELISA.



**Figure 2.** ELISA standard curves.

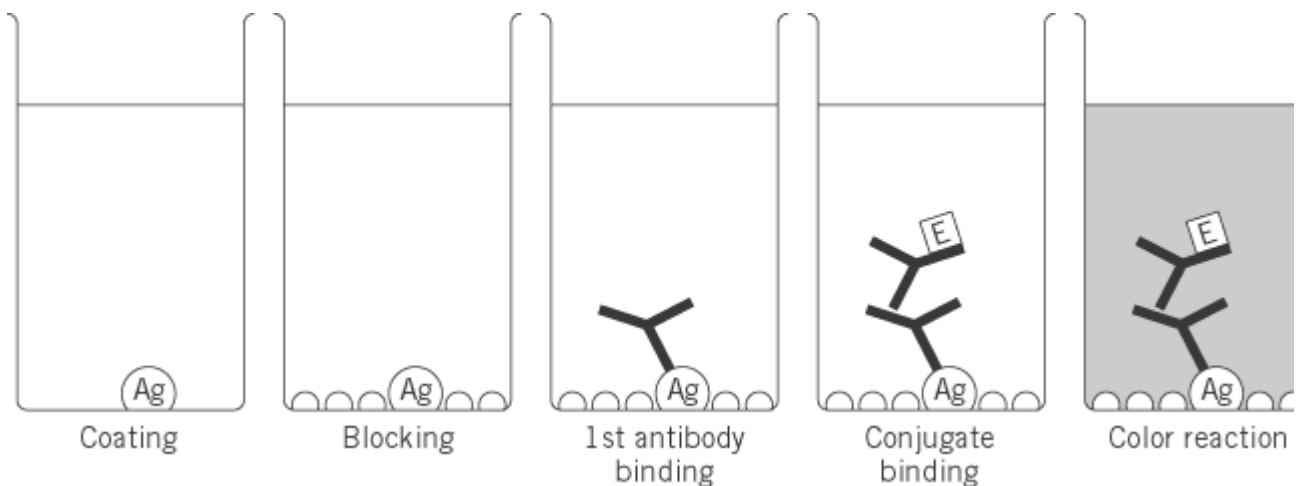


A shortcoming of this competitive design is that the response from small amounts of analyte may be almost as high as the 100% control. The change in reading due to the analyte is thus read as the difference between two large numbers. To be considered precise, the minimum readable difference must be much larger than the standard deviation of the highest point on the calibration curve, a fact that inherently limits the sensitivity of competitive ELISAs. The two noncompetitive designs that follow circumvent this limitation. The response in these assays is proportional to the analyte added; hence even a low reading may be detectable against the background level of enzymatic activity (Fig. [2b](#)).

## 2. Indirect ELISA

The enzyme conjugate in an indirect ELISA (3) is an anti-antibody; the antibody that actually recognizes that the antigen is unlabeled. Quantification of the unlabeled antibody is normally the goal sought, and the indirect ELISA is most commonly used to assess the seropositivity of subjects who have potentially been exposed to pathogens. In an indirect ELISA, antigen is immobilized and the plate blocked as above (Fig. 3). The test sample is added, along with controls consisting of antibodies of the same specificity, species, and isotype as the test sample. After an interval allowed for antigen binding, the plate is washed and the conjugate added. The antibody used to make the conjugate reagent is typically a polyclonal preparation specific for constant regions of the antibody in the test sample (see [Immunoglobulin Structure](#)). Finally, enzymic activity is assayed as above. Controls are plotted on a linear or log–log scale (Fig. 2b). In this case, the response increases with increasing amount of analyte and does not necessarily reach a plateau. Reaching a plateau at high analyte concentration may indicate that the binding capacity of the unlabeled antibody in the sample exceeds the quantity of conjugate added, or more trivially, that the optical density after the color reaction is simply too great to be measured. The concentration of analyte in the test sample is interpolated from the calibration curve, as above.

**Figure 3.** Indirect ELISA.

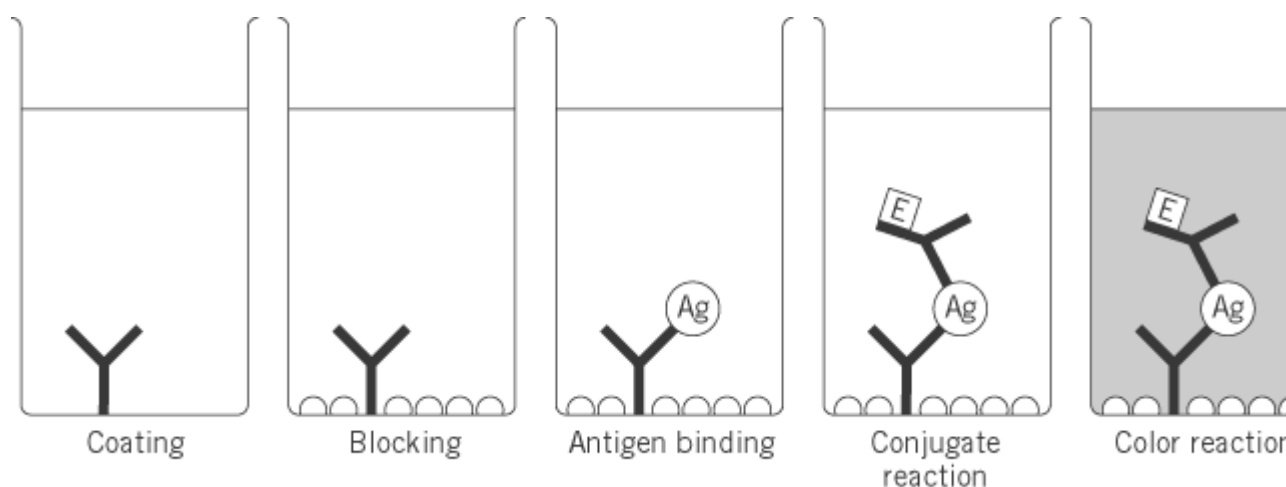


## 3. Sandwich ELISA for Antigen

This ELISA method (4) is conceptually similar to the indirect ELISA described above, but the analyte is usually a nonimmunoglobulin antigen that becomes “sandwiched” between an unlabeled immobilizing antibody and an antibody–enzyme conjugate (Fig. 4). The procedure begins with adsorption of a specific antibody to the surface of the assay plate. The plate is washed, and the test sample is added, along with controls consisting of known quantities of the same antigen. The plate is again washed, and the antibody–enzyme conjugate is added. This reagent binds to the immobilized antigen and hence becomes immobilized itself. A crucial aspect of the design of a sandwich assay is selection of an immobilizing antibody and a conjugate antibody that together will actually allow a “sandwich” to form. Obviously, two [monoclonal antibodies](#) that recognize the same [epitope](#) will exclude each other's binding. Use of a monoclonal immobilizing antibody is advantageous, however, since it occludes only a single epitope, leaving the remaining surface of the antigen available for deposition of the antibody–enzyme conjugate. After the conjugate is allowed to bind, the tray is washed a final time, and substrate is added to initiate the color reaction. A rising calibration curve

again results, from which sample concentration is interpolated as above.

**Figure 4.** Sandwich ELISA.



#### 4. Affinity Determination

Many authors have equated the analyte concentration at the 50% response point in the competitive ELISA with the equilibrium **dissociation constant** for an antibody–antigen interaction. This approximation is invalid. However, ELISA can be used to measure the equilibrium constant if the indirect ELISA method is applied after equilibrium is reached (5). The antibody and antigen in question are mixed in a series of different ratios, and the solutions are allowed to come to equilibrium. The equilibrium mixtures are then transferred to an ELISA plate that has been coated with the antigen. Standards with known amounts of free antibody are run in parallel. Free antibody will bind to the immobilized antigen. It is essential that only a small fraction of the total available amount of free antibody become bound to the plate, so that the assay does not lead to an apparent shifting of the equilibrium point (6). The plate is washed, and an antibody–enzyme conjugate that binds to the immobilized antibody is added. The plate is washed again, substrate solution added, and the amount of enzyme activity determined. The color response from the reactions with the equilibrium mixtures is compared to that from the controls, and the relative response is equated to the concentration of free antibody at equilibrium. A correction can be applied to allow for the multivalency of the antibody (7). Using the measurements of free antibody and knowledge of the total antibody and antigen present, one can infer the concentrations of free and bound antigen and antibody. These data can then be graphed on a [Scatchard Plot](#), or numerical analysis can be used to determine the association equilibrium constant for the antibody–antigen interaction.

#### Bibliography

1. E. Engvall and P. Perlman (1971) *Immunochemistry* **8**, 871–874.
2. B. K. van Weemen and A. H. W. M. Schnuurs (1971) *FEBS Lett.* **15**, 232–237.
3. E. Engvall and P. Perlmann (1972) *J. Immunol.* **109**, 129–135.
4. L. Belanger, C. Sylvestre, and D. Dufour (1973) *Clin. Chim. Acta* **48**, 15–18.
5. B. Friguet, A. F. Chaffotte, L. Djavadi-Ohanian, and M. E. Goldberg (1985) *J. Immunol. Meth.* **77**, 305–319.
6. S. Hetherington (1990) *J. Immunol. Meth.* **131**, 195–202.
7. F. S. Stevens (1987) *Mol. Immunol.* **24**, 1055–1060.



## Suggestion for Further Reading

8. D. M. Kemeny and S. J. Challacombe, eds. (1988) *ELISA and Other Solid Phase Immunoassays: Theoretical and Practical Aspects* Wiley, New York.

## Enzymes

*Enzymes* are [proteins](#) that act as biological **catalysts**. They are produced by living cells, but are independent of those cells for their catalytic activity. Like all catalysts, they speed up the rates of chemical reactions, but their special importance lies in the fact that they can do this at neutral pH and relatively low temperatures, and do so specifically, essentially for only one chemical reaction (see [Induced Fit](#)). The rate increase, which is due to the lowering of the [activation energy](#) barrier by binding and stabilizing the [transition state](#) (see [Transition State Analogue](#)), can be as much as  $10^{12}$ -fold ([1](#)). Under these circumstances, the uncatalyzed reaction would be barely detectable. As catalysts and under *in vitro* conditions, enzymes function at concentrations very much lower than those of the substrates on which they act. However, in the cell, the enzyme level often exceeds the substrate concentration. Enzymes increase the rate of a reaction in both the forward and reverse directions and, therefore, have no effect on the equilibrium constant  $K_{eq}$  of the reaction (see [Haldane Relationship](#)). In theory, enzymes as catalysts remain unchanged at the end of a reaction, but as enzymes are proteins, they are subject to **denaturation**, which can lead to loss of activity over the course of a reaction. Enzymes may show either absolute or broad specificity. For example, [catalase](#) acts only on hydrogen peroxide, whereas [alcohol dehydrogenase](#) oxidizes a number of aliphatic and aromatic alcohols.

Enzymes have the same general structure as all proteins (see [Protein Structure](#)). The basic structure consists of a [polypeptide chain](#) that is made up of 20 [amino acid](#) residues joined in a linear sequence by [peptide bonds](#). The folding of the polypeptide chain gives rise to the final enzyme structure, which contains an [active site](#) at which the substrates bind to undergo reaction. It is this folding to form the active site that is responsible for the catalytic power and specificity of enzymes. The active form of an enzyme may consist of single or multiple subunits (see [Oligomeric Proteins](#)), with molecular weights ranging from about 14,000 to 4,000,000 (Table [1](#)).

**Table 1. Molecular Weights and Subunit Composition of Selected Enzymes**

| Enzyme                                      | Molecular Weight | Number of Subunits |
|---|------------------|--------------------|
| <b>RibonucleaseA</b>                        | 13,700           | 1                  |
| <a href="#">Lysozyme</a> , hen              | 13,900           | 1                  |
| <a href="#">Carbonic anhydrase</a>          | 30,000           | 1                  |
| a-Amylase                                   | 50,000           | 2                  |
| Hexokinase                                  | 102,000          | 2                  |
| <a href="#">Alcohol dehydrogenase</a>       | 150,000          | 4                  |
| <a href="#">Aspartate transcarbamoylase</a> | 310,000          | 12                 |
| Urease                                      | 480,000          | 2                  |

|                        |           |              |
|------------------------|-----------|--------------|
| Glutamine synthetase   | 600,000   | 12           |
| Pyruvate dehydrogenase | 4,000,000 | 1 (960,000)  |
|                        |           | 24 (90,000)  |
|                        |           | 12 (112,000) |

---

Some enzymes bring about **energy transductions** that include the conversion of light energy to chemical bond energy in [photosynthesis](#), as well as the conversion of chemical bond energy to mechanical energy in muscle contraction and to pumping energy for the **membrane transport** of ions against an unfavorable gradient.

The activities of enzymes are usually subject to control, not only through changes to intracellular substrate levels, but also by activation, inhibition, and covalent modification. The activity of enzymes is influenced by pH because of the presence of ionizing amino acid residues at the active site. Some enzymes require **prosthetic groups** or metal ions (see [Coenzyme](#), [Cofactor](#)), whereas others must undergo irreversible (see [Proenzyme](#)) or reversible activation. Reversible activation is often achieved by means of **phosphorylation** or [adenylation](#). [Glycogen phosphorylase](#) is activated by the transfer of a phosphoryl group from ATP to a single serine residue and inactivated by a **phosphatase** that removes the phosphoryl group. By contrast, *glycogen synthase* is inactivated by the phosphorylation of a serine residue and activated through phosphatase action. The protein kinases responsible for the phosphorylation of phosphorylase and glycogen synthase are activated by 3', 5'-[cyclic AMP](#). The interaction of this nucleotide with the regulatory subunits on the inactive kinases causes release of the active forms of the enzymes. The reversible activation of *glutamine synthetase* involves the transfer of the AMP moiety of ATP to and from the hydroxyl group of a **tyrosine** residue. Both reactions are catalyzed by the same enzyme, *adenyl transferase*. Enzymes can be inactivated by small molecules that covalently interact with amino acid residues at the active site (see [Active Site-Directed Irreversible Inhibitors](#)) or inhibited by substrate analogues, which can behave as drugs or pesticides (see [Dead-End Inhibition](#)).

A group of enzymes that are subject to special control are the **allosteric** enzymes. The activities of these enzymes are not only more sensitive to changes in intracellular substrate concentrations, but they are also activated and/or inhibited reversibly by modifiers that are small molecules which combine at either a specific binding site on a catalytic subunit or on a special regulatory subunit. Modifiers are frequently the end-product of a metabolic pathway and have structures that differ from those of the substrates for the allosteric enzymes (see **Feedback inhibition**). The control of enzyme activity can also be achieved through their interaction with low-molecular-weight proteins, such as [calmodulin](#).

It has been demonstrated recently that it is possible to produce [monoclonal antibodies](#) which show catalytic activity (2). The generation of [catalytic antibodies](#), or abzymes, has occurred through the use of an antigen that resembles the [transition state](#) or an intermediate state of the reaction (see [Transition State Analogue](#)). The idea was to produce an antibody that would bind tightly to, and stabilize, the transition state complex so as to enhance the reaction rate. The rate enhancement and small number of antibodies produced with a particular antigen suggest that the analogues are not good mimics of the actual transition state, or that good catalytic activity requires more than simply the stabilization of a transition state complex. Catalytic activity is not restricted to proteins, for RNA molecules have been shown to function as enzymes. These relatively small molecules, which are generally found in viruses, have catalytic features in common with those of enzymes (3).

## 1. Enzyme Nomenclature

In 1956, the International Union of Biochemistry established an International Commission on Enzymes. There were several reasons why. The number of known enzymes was increasing rapidly and their naming by individual researchers was not satisfactory. Further, the names did not always convey the nature of the reaction being catalyzed and similar names were given to enzymes of different types. The first report on enzyme nomenclature was produced in 1961 and contained references to 712 enzymes. The sixth report was published in 1984 by Academic Press and contains classifications for 2,477 enzymes. The publication lists classes of enzymes as well as information about the reference number, recommended name, reaction catalyzed and other names for particular enzymes. Enzymes in this volume are categorized as [Dehydrogenases](#), [Hydrogenases](#), [Hydrolases](#), [Isomerases](#), [Synthetases/Ligases](#), [Synthases/Lyases](#), [Oxidoreductases](#), [Phosphotransferases](#), and [Transferases](#).

### Bibliography

1. L. Frick, J. P. MacNeela, and R. Wolfenden (1987) *Bioorg. Chem.* **15**, 100–108.
2. S. J. Benkovic (1992) *Ann. Rev. Biochem.* (1992) **61**, 29–54.
3. W. G. Scott and A. Klug (1996) *Trends Biochem. Sci.* **21**, 220–224.

## Epidermal Growth Factor

The epidermal growth factor (EGF) family of proteins is defined by their sequence homology and capacity to bind and activate the EGF receptor. Also described is the ErbB receptor family, which includes the EGF receptor, which is also known as ErbB-1 and HER-1, and three sequence-related receptors. Brief mention is made of the heregulin/neuregulin growth factors, which are structurally related to EGF but do not bind to the EGF receptor, although they do bind and activate other ErbB receptors.

### 1. EGF Ligands

There are six known mammalian gene products that share primary sequence homology, are recognized by the EGF receptor with high affinity (**dissociation constant** in the nanomolar range), and activate intracellular signaling pathways leading to increased cell proliferation. These ligands are known as EGF, [transforming growth factor](#) alpha (TGF $\alpha$ ), heparin-binding EGF (HB-EGF), amphiregulin, betacellulin, and epiregulin ([1](#), [2](#)). As noted below, two of these ligands, betacellulin and epiregulin, are able to bind not only to the EGF receptor, but also to ErbB-4, a separate member of the EGF receptor family. In addition to these ligands encoded by genes in mammalian cells, EGF growth factors are encoded by genes in the fruit fly *Drosophila* (*sptiz*, *gurken*, *vein*) ([3](#)), the [nematode](#) *Caenorhabditis elegans* (*lin3*) ([4](#)), and certain Pox-family viruses ([5](#)). The *Drosophila* gene *argos* apparently encodes an antagonist growth factor for the fly EGF receptor ([3](#)). To date, however, EGF receptor antagonists have not been found in other biological systems or created through synthetic chemistry.

Structurally, the EGF receptor ligands are small monomeric proteins of 6–10 kDa that are characterized by conserved sequence features—most notably six **cysteine** residues that, in the mature ligands, form three characteristically spaced disulfide bonds producing three loop structures within the molecule ([5](#)). Nuclear magnetic resonance ([NMR](#)) has been employed to determine the high-resolution [protein structures](#) of EGF and TGF $\alpha$  (see [EGF Motif](#)), while site-directed mutagenesis has

revealed which amino acid residues are essential for growth factor function (6). The structures of EGF and TGF $\alpha$  are highly superimposable, and each is dominated by [beta-sheet](#) structure. Growth factor mutagenesis (5, 6) and an NMR study of a TGF $\alpha$ :EGF receptor ectodomain complex (7) implicate residues in both the amino- and carboxy-terminal domains, which are not closely packed together in the high-resolution structure. Therefore, the growth factor may use more than one surface or domain to contact the receptor binding site. This multidomain view of growth factor binding to its receptor has also provoked models in which the ligand is bivalent and capable of associating with two receptor molecules (8).

The mature proteins, especially EGF, are chemically and biologically very stable to extreme physical conditions, such as acid and heat. All the growth factors have been produced as [recombinant proteins](#), and many are commercially available. EGF is easily derivatized with covalent probes, including  $^{125}\text{I}$ , fluorescein, [biotin](#), and various [toxins](#) (5). These modifications have been employed to quantify and/or visualize growth-factor interactions with receptors or, in the case of toxin conjugates, to kill cells having EGF receptors, particularly certain tumors that overexpress this receptor (9).

The mature EGF and EGF-like ligands are found in the extracellular space as freely diffusible proteins, but they are not products of the secretory pathway. The extracellular ligands are derived from the ectodomain of transmembrane precursor proteins present on the surface of cells (5). In these precursors, the EGF or EGF-like sequence is liberated into the extracellular environment as a diffusible growth factor following a series of **proteolytic** processing steps (5, 10). This means that the presence or concentration of extracellular mature growth factor is determined not only at the level of gene expression, but also by control of the extent of precursor proteolytic processing. Yet relatively little is known regarding the identity or biological control of the activity of the [proteinases](#) involved, although [metalloproteinase](#) activity is frequently implicated on the basis of inhibitor studies for the precursor processing of TGF $\alpha$  (11), EGF (12), HB-EGF (13), and amphiregulin (14). The existence of transmembrane precursors has raised the question of whether the precursor form of the growth factor can activate EGF receptors on adjacent cells, without the necessity for proteolytic processing. While this is not trivial to test experimentally, the available data indicate that such communication is possible (10). Whether this is biologically significant in mammals is unclear; in developmental circumstances, however, where neighboring cells obviously must communicate, such juxtacrine communication is possibly significant.

EGF and other EGF-like ligands are present in most bodily fluids (urine, milk, saliva), but they are difficult to detect in blood or serum (5). The production of these growth factors is not confined to one or two organs within the body, but it seems to occur within many tissues. Hence, the circulatory system is not essential for delivery of the growth factor to target cells. It seems likely that growth factors like the EGF ligands are examples of paracrine hormones in which the production and utilization of the ligand occurs in different cells within the same local environment of a tissue. Each EGF-like growth factor has a distinctive pattern of expression in tissues, although in certain tissues there is usually some degree of overlap with the expression of other members of this growth factor family. In the case of tumor cells, it is a frequent observation that the same cell constitutively produces and is activated by a growth factor. This is a process termed autocrine regulation and may be particularly important for TGF $\alpha$ , amphiregulin, and epiregulin, which are frequently expressed in tumor cells (15).

EGF is known to stimulate the proliferation of a great many cell types in cell culture, as its receptor is expressed in almost all nonhematopoietic cells (5). In the intact animal, however, enhanced proliferation of epithelial tissues is the most commonly observed response to the administration of exogenous EGF (5). Proliferation of skin and isolated epidermis have been the model system for an EGF-responsive tissue, but published data also show enhanced proliferation *in vivo* of internal epithelial tissues, such as trachea and gastrointestinal tissue. In the adult, it seems likely that the EGF family of growth factors participates in control of the relatively high level of homeostatic

proliferation necessary to support epithelial tissues. Administration of EGF to fetal or newborn animals often hastens certain developmental processes, such as lung maturation and eyelid opening, respectively (5). TGF $\alpha$  and amphiregulin are particularly recognized as EGF family members that are expressed during gestation (16). At least one nonmitogenic biological response has been described following EGF administration to animals: the inhibition of parietal-cell acid secretion (5). It is unclear, however, whether EGF ligands normally participate in pH control within the gastric system.

In humans, the genes for EGF, TGF $\alpha$ , amphiregulin, and HB-EGF are localized, respectively, to chromosomes 4q25  $\rightarrow$  q29, 2p11  $\rightarrow$  p13, 4q13  $\rightarrow$  q21, and 5 (17). The genes range in size from 120 kbp and 24 exons (for EGF) to 18, 14, and 10 kb and 6 exons each for TGF $\alpha$ , HB-EGF, and amphiregulin, respectively. The expression of genes in this family has been demonstrated in both nontransformed and transformed cells to involve autoinduction by the homologous growth factor, or by a different member of this growth factor family (18). The autoinduction mechanism could be a significant amplification factor in autocrine tumor growth and in Poxvirus infection.

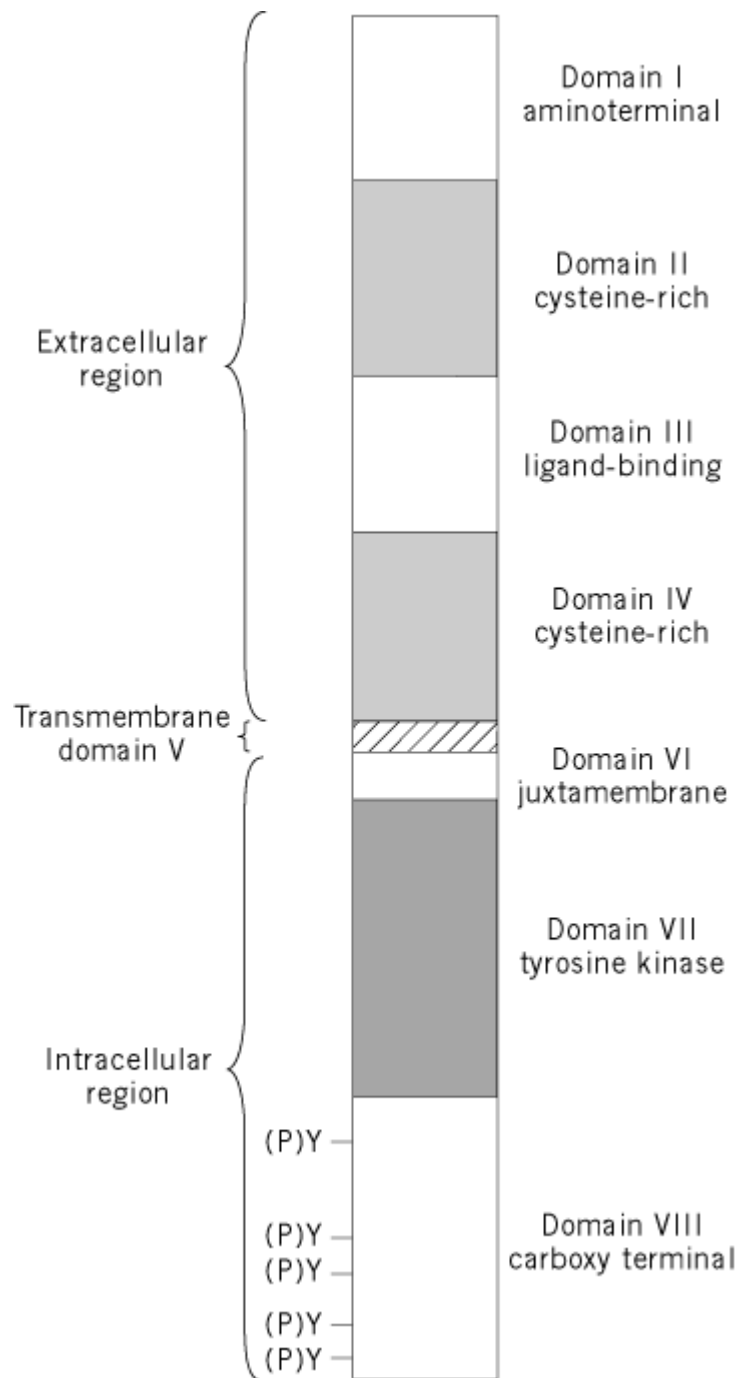
In contrast to the severe phenotype exhibited following targeted disruption of the EGF receptor gene in mice (see text below), “**knockout**” of the TGF $\alpha$  gene in mice produces only a mild disturbance in the function of the hair follicle and a consequent “wavy” hair **phenotype** (19, 20). While knockouts of other family members have yet to be reported, preliminary data for some also suggest mild changes in phenotype. Hence, while receptor function is indispensable in the animals, there does seem to be considerable functional redundancy within the EGF family of ligands.

## 2. EGF Receptors

Based primarily on sequence similarity, four members have been identified in the ErbB family of EGF receptor-related molecules (1, 2): the EGF receptor (also known as HER-1 and ErbB-1); ErbB-2 (Neu, HER-2); ErbB-3 (HER-3); and ErbB-4 (HER-4). EGF, EGF-like ligands, and heregulin (neuregulin) activate these receptors, but in distinctive and complex ways (1, 2). EGF and the EGF-like ligands (HB-EGF, TGF $\alpha$ , amphiregulin, betacellulin, epiregulin) recognize ErbB-1, while the heregulin growth factors bind to ErbB-3 and ErbB-4. However, some EGF family molecules can bind to the EGF receptor and another ErbB receptor. For example, betacellulin, HB-EGF, and epiregulin, but not EGF or TGF $\alpha$ , activate ErbB-4 as well as the EGF receptor. The heregulin isoforms bind to ErbB-3 and ErbB-4, but not to the EGF receptor. There is no known ligand for ErbB-2, which seems to function as a co-receptor with each of the other ErbB family members (21, 22).

The mature forms of these receptors contain approximately 130 kDa of protein complexed with N-linked carbohydrate (see *N*-Glycosylation) to give a molecular mass of approximately 175–185 kDa (5). Based on [cloning](#) and [sequencing](#) of [complementary DNA](#) for the EGF receptor prototype for this receptor family (23), the transmembrane orientation of the ErbB receptors is believed to consist of three regions, as shown in Figure 1. These are (i) an extracellular **ligand-binding** region, (ii) a cytoplasmic region, and (iii) a single **hydrophobic** transmembrane domain of approximately 24 amino acid residues. Among these ErbB family receptors, the cytoplasmic carboxy-terminal domains show similarity in size, but are the most heterogeneous domains in terms of sequence similarity. The EGF receptor is the most widely-studied and best-characterized molecule among the members of ErbB receptor family. The remaining portion of this article is, therefore, focused on the structure–function, activation, regulation, and genetic alterations of the EGF receptor.

**Figure 1.** Schematic representation of the EGF receptor. Based on its amino acid sequence, obtained by cDNA sequencing (23), the protein structure of the EGF receptor is subdivided into extracellular, transmembrane, and intracellular regions and into eight domains. The phosphotyrosyl residues in the carboxy-terminal domain are indicated by (P)Y.



As depicted in Figure 1, the extracellular region of the EGF receptor can be divided into four **domains** based on the presence of two cysteine-rich regions (Domains II and IV). The functions of the cysteine-rich domains are not clearly understood, although they probably provide significant structural stability, because 25 disulfide bonds are found in these two domains (24). Domain III, localized between these two cysteine-rich domains, is thought to mediate ligand binding (5). There are implications that residues in Domain I may also contribute to ligand binding (25), but there are no high-resolution data to define the receptor/ligand binding site. The ectodomain of the receptor is heavily glycosylated with 10–11 N-linked oligosaccharide chains (5). Among these carbohydrates, at least three are of the high-mannose type, while the others are complex-type oligosaccharide chains. While oligosaccharide addition is necessary for correct EGF receptor folding and processing to the cell surface, there is no evidence that carbohydrate is a determinant of ligand binding or other receptor functions at the plasma membrane.

The intracellular or cytoplasmic region of the EGF receptor is composed of three domains: the juxtamembrane domain, a **tyrosine kinase** domain, and a carboxy-terminal domain (Fig. 1). The receptor is actually an **enzyme** (tyrosine kinase), and the kinase activity is dependent on ligand binding to the receptor ectodomain (5). Clearly, the tyrosine kinase domain is critical to the receptor's capacity to elicit biological responses, namely mitogenesis, following ligand binding. Inactivation of the kinase domain by **site-directed mutagenesis** abrogates EGF responses. Five tyrosine autophosphorylation sites have been identified in the carboxy-terminal domain at Tyr992, 1068, 1086, 1148, and 1173, while modification by **phosphorylation of serine or threonine** residues by **serine–threonine kinases** has been localized to both the carboxy-terminal and juxtamembrane domains. One of these, Thr654 in the juxtamembrane region, is a substrate for protein kinase C, and its phosphorylation results in attenuation of receptor tyrosine kinase activity. Phosphorylation of the tyrosine residues of the carboxy-terminal domain results in formation of “docking sites” for many of the substrates of the receptor tyrosine kinase (26).

The EGF receptor is activated by a two-step mechanism (26). First, binding of a cognate ligand to the monomeric receptor induces the formation of noncovalent receptor dimers. Although the mechanism of dimerization is not clear, it may be due to a conformational change in the receptor external domain or to bivalent binding of EGF to the receptor. Interestingly, the EGF receptor can, in response to ligand binding, also form heterodimers with other ErbB family receptors, particularly ErbB-2 and ErbB-3 (1, 2). The second step in receptor activation occurs when dimerization of the receptors leads to juxtapositioning of the cytoplasmic domains and transphosphorylation of the two receptors present in the dimer (26). Hence, autophosphorylation is probably a mixture of inter- and intramolecular phosphorylations. At present it is unclear how autophosphorylation sites may relate to the control of tyrosine kinase activity. Unlike many other tyrosine kinase receptors, none of the EGF receptor autophosphorylation sites is within the kinase domain, where it might directly regulate the kinase active site.

Based on the studies with “dominant-negative” kinase-inactive EGF receptor mutants, it has been concluded that transphosphorylation within dimers is essential for substrate phosphorylation and concomitant EGF-mediated **signal transduction** (26). When a kinase-inactive EGF receptor mutant is coexpressed in the same cell with wild-type receptor, heterodimers of the two receptors (wild-type and mutant) are produced in response to EGF. Kinase activation and signaling are suppressed in proportion to the relative overexpression level of the mutant receptors. The reason for attenuation of the wild-type receptor function is that kinase-inactive receptors do form dimers with the wild-type receptors, but transphosphorylation does not occur.

The cytoplasmic substrates that are phosphorylated by the EGF receptor include: (i) enzymes, such as **phospholipase C-g1**; **Src**, a cytoplasmic tyrosine kinase; the **ras** GTPase-activating protein; and the phosphotyrosine phosphatase SHP-2; (ii) the transcription factor STAT; and (iii) the nonenzymatic adapter molecules GRB-2, the p85 subunit of phosphatidylinositol 3-kinase, eps8, eps15, Nck, and Shc (1). Also, heterodimerization leads to EGF-dependent phosphorylation of ErbB-2 and ErbB-3 (1, 2). Many, but not all, of these receptor-proximal proteins have in common the presence of a noncatalytic domain of approximately 100 residues termed a *src* homology, or **SH2**, domain (27). SH2 domains recognize phosphotyrosine-containing sequences and thereby mediate association of these molecules with activated receptors. This mechanism to promote receptor association serves two purposes. First, it facilitates subsequent tyrosine phosphorylation of the SH2-containing protein by the receptor tyrosine kinase domain. Second, it relocalizes cytoplasmic SH2 molecules to the receptor carboxy-terminus and hence close to the cytoplasmic face of the plasma membrane. Relocalization may be a significant factor for those molecules that function by modifying plasma membrane phospholipids (eg, phospholipase C-g1) or those (GRB-2 and Shc) that interact with membrane-localized signal transduction target proteins, such as Ras. Actually, the GRB-2 adaptor associates with the EGF receptor, but it is not detectably phosphorylated. Therefore, relocation of GRB-2 seems to be the only function of its receptor association.

Desensitization of the activated EGF receptor is achieved by two mechanisms. The first and best understood is endocytosis of the activated receptor, followed by its degradation in **lysosomes** (28). This provides a relatively slow rate of receptor inactivation and results in the down-regulation of the cellular level of receptor protein, when receptor degradation exceeds receptor biosynthesis. To initiate this process, ligand binding to the EGF receptor induces receptor clustering in **clathrin-coated pits** and rapid internalization. Following this internalization step, ligand:receptor complexes in **endosomes** are targeted to the lysosome, where both the receptor and its ligand are degraded. It should be noted that during a substantial portion of this endocytosis process, EGF remains bound to the receptor and phosphotyrosine remains detectable on the receptor. Hence, the internalized receptor remains potentially functional for some period of time. Curiously, the heregulin receptors (ErbB-3 and ErbB-4) are not subject to rapid internalization following ligand binding (29).

A second means by which activated receptors are desensitized involves the phosphorylation of serine/threonine residues within the cytoplasmic domain (30). This mechanism is not, however, well understood.

Overexpression or mutation of the EGF receptor has been reported to be associated with certain malignancies (31). For example, there is frequently detected in glioblastomas a small genomic deletion of exons encoding a small region of the EGF receptor ectodomain. This mutation inhibits ligand binding but produces a constitutively active receptor. In a sizable fraction of breast cancers, the EGF receptor is overexpressed but not mutated (31, 32). Also, the EGF co-receptor ErbB-2 is even more frequently overexpressed in breast cancer. **Antibodies** to the EGF receptor and ErbB-2 are being evaluated for therapeutic efficacy in breast cancer clinical trials (33).

In humans, the gene encoding the EGF receptor is localized to the p13 → q22 region of chromosome 7 (17). Expression of the EGF receptor gene produces two major species of **messenger RNA**: 10.5 and 5.8 kb. The 5.8-kb mRNA represents a full-length cDNA that includes 5'- and 3'-untranslated sequences (23, 34). The genesis and significance of the 10.5-kb mRNA is not clearly understood, but it probably contains a longer 3'-untranslated sequence. In certain cells, a smaller mRNA (2.0 kb) is also detected. This mRNA is formed by **RNA splicing** and encodes the EGF receptor ectodomain. In cells expressing this mRNA, the ectodomain fragment is fully glycosylated and secreted (35). The **promoter** sequence of the EGF receptor represents a typical “**housekeeping gene**” with extremely GC-rich sequences and the absence of **TATA** and **CAAT boxes** (36).

Targeted disruption of the EGF receptor gene in mice produces multiorgan failure and embryonic lethality in most, but not all, mouse strains (37-39). In some mouse strains, it has been reported that “knockouts” of the EGF receptor live up to 8 days after birth and suffer from impaired epithelial development in skin, liver, and gastrointestinal tract. However, neurodegeneration is the most common and strain-independent phenotype of postnatal mice that lack EGF receptors (40). Gene “knockouts” of ErbB-2 (41, 42) and ErbB-4 (43) both produce deficiencies in neurodevelopment and development of the embryonic heart, leading to lethality at 10 days’ gestation. Homozygous disruption of the ErbB-3 gene in mice produces severe defects in neurogenesis, specifically in Schwann cell development and the maturation of neurons, plus abnormalities in cardiac development (42, 44). In fact, defects in neural development have been noted in all ErbB knockouts, including the EGF receptor. EGF receptor-null mice, however, have not been reported to be abnormal in heart development.

### 3. Acknowledgment

Our work was supported by NIH grants CA24071 and CA75195.

### Bibliography

1. I. Alroy and Y. Yarden (1997) *FEBS Lett.* **410**, 83–86.
2. D. J. Riess II and D. F. Stern (1998) *Bioessays* **20**, 41–48.



3. J. D. Wasserman and M. Freeman (1997) *Trends Cell Biol.* **7**, 431–436.
4. M. Sundaram and M. Han (1996) *Bioessays* **18**, 473–480.
5. G. Carpenter and M. I. Wahl (1990) In *Handbook of Experimental Pharmacology*, **95/I** (M. B. Sporn and A. B. Roberts, eds.) Springer-Verlag, Berlin, pp. 69–171.
6. L. C. Groenen, E. C. Nice, and A. W. Burgess (1994) *Growth Factors* **11**, 235–257.
7. C. McInnes, D. W. Hoyt, R. N. Harkins, R. N. Pagila, M. T. Debanne, M. O'Connor-McCourt, and B. D. Sykes (1996) *J. Biol. Chem.* **271**, 32203–32211.
8. M. A. Lemmon, Z. Bu, J. E. Ladbury, M. Zhou, D. Pinchasi, I. Lax, D. E. Engelman, and J. Schlessinger (1997) *EMBO J.* **16**, 281–294.
9. I. Pastan and D. FitzGerald (1991) *Science* **254**, 1173–1177.
10. J. Massagué (1990) *J. Biol. Chem.* **265**, 21393–21396.
11. J. Arribas, L. Coodly, P. Vollmer, T. K. Kishimoto, S. Rose-John, and J. Massagué (1996) *J. Biol. Chem.* **271**, 11376–11382.
12. P. J. Dempsey, K. S. Meise, Y. Yoshitake, K. Nishikawa, and R. J. Coffey (1997) *J. Cell Biol.* **138**, 747–758.
13. M. Suzuki, G. Raab, M. A. Moses, C. A. Fernandez, and M. Klagsbrun (1997) *J. Biol. Chem.* **272**, 31730–31737.
14. C. L. Brown, K. S. Meise, G. D. Plowman, R. J. Coffey, and P. J. Dempsey (1998) *J. Biol. Chem.* **273**, 17258–17268.
15. A. W. Burgess and C. M. Thumwood (1994) *Pathology* **26**, 453–463.
16. E. D. Adamson and L. M. Wiley (1997) In *Current Topics in Developmental Biology*, Vol. **35**, Academic Press, New York, pp. 71–120.
17. C. Soler and G. Carpenter (1995) In *Guidebook to Cytokines and Their Receptors* (N. A. Nicola, ed.), Oxford University Press, Oxford, pp. 194–197.
18. J. A. Barnard, R. Graves-deal, M. R. Pittelkow, R. DuBois, P. Cook, G. W. Ramsey, P. R. Bishop, L. Damstrup, and R. J. Coffey (1994) *J. Biol. Chem.* **269**, 22817–22822.
19. G. B. Mann, K. J. Fowler, A. Gabriel, E. C. Nice, R. L. Williams, and A. R. Dunn (1993) *Cell* **73**, 249–261.
20. N. C. Leutteke, T. H. Qiu, R. L. Peiffer, P. Oliver, O. Smithies, and D. C. Lee (1993) *Cell* **73**, 263–278.
21. D. Karunakaran, E. Tzahar, R. R. Beerli, X. Chen, D. Graus-Porta, B. J. Ratzkin, R. Seger, N. E. Hynes, and Y. Yarden (1996) *EMBO J.* **15**, 254–264.
22. D. Graus-Porta, R. R. Beerli, J. M. Daly, and N. E. Hynes (1997) *EMBO J.* **16**, 1647–1655.
23. A. Ullrich, L. Coussens, J. S. Hayflick, T. J. Dull, A. Gray, A. W. Tam, J. Lee, Y. Yarden, T. A. Libermann, J. Schlessinger, J. Downward, E. L. V. Mayes, N. Whittle, M. D. Waterfield, and P. H. Seeburg (1984) *Nature* **309**, 418–424.
24. Y. Abe, M. Odaka, F. Inagaki, I. Lax, J. Schlessinger, and D. Kohda (1998) *J. Biol. Chem.* **273**, 11150–11157.
25. A. E. Summerfield, A. K. Hudnall, T. J. Lukas, C. A. Guyer, and J. V. Staros (1996) *J. Biol. Chem.* **271**, 19656–19659.
26. J. Schlessinger and A. Ullrich (1992) *Neuron* **9**, 383–391.
27. C. A. Koch, D. Anderson, M. F. Moran, C. Ellis, and T. Pawson (1991) *Science* **252**, 668–674.
28. A. Sorkin and C. M. Waters (1993) *Bioessays* **15**, 375–382.
29. J. Baulida, M. H. Kraus, M. Alimandi, P. P. Di Fiore, and G. Carpenter (1996) *J. Biol. Chem.* **271**, 5251–5257.
30. S. J. Therous, D. A. Latour, K. Stanley, D. L. Raden, and R. J. Davis (1992) *J. Biol. Chem.* **267**, 16620–16626.
31. H. Modjtahedi and C. Dean (1994) *Int. J. Oncol.* **4**, 277–296.

32. S. B. Fox and A. L. Harris (1997) *J. Mamm. Gland Biol. Neopl.* **2**, 131–141.
33. J. Baselga and J. Mendelsohn (1997) *J. Mamm. Gland Biol. Neopl.* **2**, 165–174.
34. C. R. Lin, W. S. Chen, W. Kruiger, L. S. Stolarsky, W. Weber, R. M. Evans, I. M. Verma, G. N. Gill, and M. G. Rosenfeld. (1984) *Science* **224**, 843–848.
35. W. Weber, G. N. Gill, and J. Spiess (1984) *Science* **224**, 294–297.
36. N. P. Bates and H. C. Hurst (1997) *J. Mamm. Gland Biol. Neopl.* **2**, 153–163.
37. D. W. Threadgill, A. A. Dlugosz, L. A. Hansen, T. Tennenbaum, U. Lichti, D. Yee, C. LaMantia, T. Mourton, K. Herrup, R. C. Harris, J. A. Barnard, S. H. Yuspa, R. J. Coffey, and T. Magnuson (1995) *Science* **269**, 230–234.
38. M. Sibilica and E. F. Wagner (1995) *Science* **269**, 234–238.
39. P. J. Miettinen, J. E. Berger, J. Meneses, Y. Phung, R. A. Pedersen, Z. Werb, and R. Derynck (1995) *Nature* **376**, 337–341.
40. M. Sibilica, J. P. Steinbach, L. Stingl, A. Aguzzi, and E. F. Wagner (1998) *EMBO J* **17**, 719–731.
41. K.-F. Lee, H. Simon, H. Chen, B. Bates, M.-C. Hung, and C. Hauser (1995) *Nature* **378**, 394–398.
42. S. L. Erickson, K. S. O'Shea, N. Ghaboosi, L. Loverro, G. Frantz, M. Bauer, L. H. Lu, and M. W. Moore (1997) *Development* **124**, 4999–5011.
43. M. Grassmann, F. Casagrande, D. Orioli, H. Simon, C. Lai, R. Klein, and G. Lemke (1995) *Nature* **378**, 390–386.
44. D. Riethmacher, E. Sonnenberg-Riethmacher, V. Brinkmann, T. Yamaai, G. R. Lewin, and C. Birchmeier (1997) *Nature* **389**, 725–730.

### **Suggestions for Further Reading**

45. Listed below are key review articles arranged to reflect the chronological development of EGF research.
46. G. Carpenter and S. Cohen (1979) Epidermal growth factor, *Annu. Rev. Biochem.* **48**, 193–216.
47. G. Carpenter and M. I. Wahl (1990) "The epidermal growth factor family", In *Handbook of Experimental Pharmacology*, Vol **95/I**, Springer-Verlag, Berlin, pp. 69–171.
48. J. Schlessinger and A. Ullrich (1992) Growth factor signaling by receptor tyrosine kinases, *Neuron* **9**, 383–391.
49. I. Alroy and Y. Yarden (1997) The ErbB signaling network in embryogenesis and oncogenesis: Signal diversification through combinatorial ligand–receptor interactions, *FEBS Lett.* **410**, 83–86.
50. J. B. Duffy and N. Perrimon (1996) Recent advances in understanding signal transduction pathways in worms and flies, *Curr. Opin. Cell Biol.* **8**, 231–238.

### **Epigenetics**

Epigenetics is derived from the Greek word epigenesis, which describes the Aristotelian concept that an organism develops progressively through a series of causal interactions of various parts. During the 16th and 17th centuries, the epigenetic mode of [development](#) was held in contrast to the preformation theory, which maintained that development occurred by simple enlargement of a preformed homunculus in the germ cell (for the early history of its usage, see ref. [1](#)). More modern

definitions include the concept that, during development, cells progress through microenvironments that induce a stepwise restriction of their developmental potentialities into specialized ones via changes in gene expression (2). Based on this view of development, Holliday has defined epigenetics as the study of the mechanisms “that impart temporal and spatial control on the activity of all those genes required for the development of a complex organism from the zygote to the fully formed adult” (3). More recent usage, however, has tended to emphasize the “epi” or “beyond” aspect of the term. Strohman points out that “epigenesis is the historical alternative to genetic determinism and that it establishes the basis for a level of organizational control above the genome” (4). Accordingly, molecular geneticists use the term epigenetics for a process by which the state of gene expression is modified at a given time of development and becomes heritable at the cellular or organismal level. Molecular mechanisms underlying epigenetic phenomena are known only in a few cases. The known mechanisms include modifications of DNA and changes in the [chromatin](#) structure, which are thought to be propagated through templating (5).

There exist a broad variety of epigenetic phenomena in many organisms. Perhaps one of the best-studied epigenetic phenomena is **mating-type** silencing in the yeast *Saccharomyces cerevisiae* (6, 7). Yeast cells possess one of two mating types, a or a, determined by which of the alternate mating-type **alleles** is present at the mating-type locus (MAT). In addition, cells harbor two other copies of mating-type genes at the loci called HML and HMR. Despite their identical sequences, the mating-type gene located at the MAT locus is active in transcription; those at HML and HMR are silent. The molecular basis of mating-type silencing appears to lie in the specific type of chromatin structure at the HML and HMR loci.

In the fruit fly, *Drosophila melanogaster*, several phenomena have been described as epigenetic. The first example, known as [position effect](#) variegation (PEV), is observed when a chromosomal segment of [euchromatin](#) is translocated adjacent to a heterochromatic region (8). Such a [translocation](#) results in variegated spreading of repression from the heterochromatic region into the euchromatic segment, causing a metastable state of expression of the translocated genes. Because the variable spreading of gene inactivation is inherited clonally, the PEV phenomenon is considered epigenetic. The second example includes the **silencing** of [homeotic genes](#) during development (9). Homeotic genes are expressed in specific domains of the body and silenced in others. This spatially regulated silencing is established in early embryogenesis and is maintained throughout development by a group of chromatin-associated proteins known as the [Polycomb group](#) proteins.

In mammals, genomic [imprinting](#) and [X-Chromosome Inactivation](#) in females are both examples of epigenetic effects. Genomic imprinting has been observed for some 20 genes whose expression in an early embryo depends on whether the allele is derived maternally or paternally (10). Some genes are inactive when derived from maternal gametes; in other genes those derived paternally are inactive. The molecular basis of imprinting is known to be [methylation](#) of the DNA. The state of inactivity is inherited clonally through subsequent cell divisions by templating of the methylated state. Particular CpG sequences in imprinted genes are methylated in gametes of one sex but not the other (see [CpG Islands](#)). The zygote therefore becomes heterozygous at imprinted genes, carrying a methylated allele and a nonmethylated allele. During development, methylated alleles are not expressed, while **homologous** copies that are nonmethylated are expressed and provide the necessary functions. In primordial germ cells (the precursors of eggs and sperms), imprinting is erased by demethylation of appropriate residues in the imprinted genes. These residues are remethylated *de novo* in gametes in either the maternal or paternal mode.

In female mammals, one of two **X chromosomes** is inactivated, appearing cytologically as a [Barr body](#). The X chromosome is inactivated by coating with the so-called Xist RNA (11). The Xist gene on the active chromosome is turned off by methylation of its **promoter**, while the Xist gene on the chromosome to be inactivated produces Xist RNA that coats its own chromosome. Once an X chromosome is inactivated, its inactive state is clonally inherited in its descendant cells, without the need of Xist RNA. Thus, Xist RNA appears to initiate an epigenetic state that is perpetuated, without playing a further role.

Lesser known, but equally fascinating, epigenetic phenomena are found in fungi and plants. *Neurospora crassa* and *Ascobolus immersus* seem to protect their [genomes](#) from assaults by [transposable elements](#) through mechanisms known as repeat induced-point mutation (RIP) and methylation-induced premeiotically (MIP), respectively ([12](#), [13](#)). Transposable elements, often present in multiple copies, trigger their mutual silencing in a process known as quelling by methylation or C to T [transition mutations](#) of repeated sequences. A paramutation, found in higher plants such as maize, tomato, or tobacco, is a mitotically and meiotically heritable change in gene expression that does not involve alterations in DNA sequences ([14](#)). At paramutable loci, some alleles, called paramutable alleles, are susceptible to paramutation, others, known as paramutagenic alleles, induce paramutation in the homologous copy of the same gene.

### Bibliography

1. R. J. Richards (1992) *The Meaning of Evolution*, The University of Chicago Press, Chicago, pp. 5–16.
2. C. H. Waddington (1956) *Principles of Embryology*, Allen and Unwin, London, p. 350.
3. R. Holliday (1990) *Phil. Trans. Royal Soc. Lond.* **B326**, 329–338.
4. R. C. Strohman (1997) *Nature Biotechnol.*, **15**, 194–200.
5. V. E. A. Russo, R. A. Martienssen, and A. D. Riggs, eds. (1996) *Epigenetic Mechanisms of Gene Regulation*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor.
6. S. Loo and J. Rine (1995) *Annu. Rev. Cell Dev. Biol.* **11**, 519–548.
7. S. G. Holmes, M. Braunstein, and J. R. Broach (1996) in *Epigenetic Mechanisms of Gene Regulation*, Cold Spring Harbor Laboratory Press, pp. 467–487.
8. B. T. Wakimoto (1998) *Cell* **93**, 321–324.
9. V. Pirrotta (1998) *Cell* **93**, 333–336.
10. M. A. Surani (1998) *Cell* **93**, 309–312.
11. B. Panning and R. Jaenisch (1998) *Cell* **93**, 305–308.
12. M. J. Singer and E. U. Selker (1995) In P. Meyer, ed., *Gene Silencing in Higher Plants and Related Phenomena in Other Eukaryotes*, pp. 165–177.
13. C. Goyon et al. (1994) *J. Mol. Biol.* **240**, 42–51.
14. V. L. Chandler et al. (1996) In *Epigenetic Mechanisms of Gene Regulation* (V. E. A. Russo, R. A. Martienssen, and A. D. Riggs, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, pp. 289–304.

### Suggestions for Further Reading

15. S. C. R. Elgin (1995) *Chromatin Structure and Gene Expression*, IRL Press, Oxford.
16. *Cell* **93**, 301–348 (1998) (includes several reviews on epigenetics).
17. *Trends in Genetics*, **13**, 293–341 (1997) (a special issue on epigenetics).
18. E. Jablonka and M. J. Lamb (1995) *Epigenetic Inheritance and Evolution*, Oxford University Press, Oxford.

### Epigenic Patterning

In the animal kingdom, the first developmental determination of the **germ cells** is the presence of a

*polarity*. This is necessary to determine early coordinates for further [development](#) and organogenesis. In other words, there are one or more mechanisms to determine which part of the cell is supposed to form the posterior, anterior, dorsal, and ventral poles. The morphological appearance of the [egg](#) before fertilization, and also shortly after, [fertilization](#) is not uniform but shows a more or less distinct polarity, which is termed the [animal–vegetal polarity](#). In certain species, [sperm](#) can only enter at the animal pole, indicating that there must be differences in molecular composition of the opposite poles.

There are various reasons for the development of the animal–vegetal polarity. In many cases, the feeder cells of oocytes are not uniformly located around the oocytes. Thus, they are fed from one side, which results in an asymmetrical concentration of certain molecules within the oocyte plasma. Another influence to produce polarity is probably light and gravity. However, after fertilization in space or in a fast rotating clinostate, [Xenopus](#) exhibited normal embryonic development and normal axis formation (1). In eggs containing a considerable amount of yolk, gravity forces the yolk to settle down. This event takes place in amphibia after sperm entry, which enables a free rotation of the egg **cytoplasm**. This causes the formation of the vegetal and animal polarity, in which the vegetal pole contains larger amounts of yolk than the animal pole.

The other polarity of the cell is not as easily determined in the unfertilized egg. In many cases, the posterior-anterior axis is determined by the site of sperm entry. On egg activation, a number of cytoplasmic reorganizations take place. In some species a region becomes visible, the *gray crescent*, which is located opposite to the site of the sperm entry. This is the prospective area in which **gastrulation** will take place. In amphibia, this region marks the prospective backside and distal side of the embryo. In birds, the dorsal-ventral polarity is determined by the structure of the egg, and the caudal distal polarity is assumed to be based on the influence of gravity. In mammals, there are contrary observations regarding where the inner cell mass settles. This embryonic event determines at least the dorsal side, whereas the mechanism that determines the prospective axis of the head and tail is still the subject of investigation.

Nüsslein-Volhard and others have provided more information about the situation in *Drosophila*. In this species, the presence of a genetically determined polarity of the anterior and posterior regions of the egg was demonstrated. In the anterior region, a product of the gene *bicoid* was found to determine the development of the head and thorax (2, 3). In *Drosophila* *bcd* – / – mothers, which could not produce the bicoid gene product, larvae with no heads appeared. The product of the gene *nanos* is responsible for the formation of an abdomen and germline development (4), and its highest concentration occurs in the posterior end of the egg (5). In mothers homozygous negative for the *nanos* gene, no abdomen is formed. Compensation of the missing gene products by injection of *bicoid* [messenger RNA](#) or *nanos* mRNA leads in each negative variant to the development of a normal phenotype. Nanos and bicoid form concentration gradients from the anterior to the posterior and activate in a concentration-dependent manner a number of zygotic segmentation genes (6, 7). Bicoid controls [transcription](#) of the *hunchback* gene (6). In this cascade, the gene *staufer* is important for the localization of maternal mRNA to the posterior pole and for bicoid mRNA to the anterior pole in the *Drosophila* egg (8). In contrast, the nanos protein suppresses translation of *hunchback* mRNA. *Hunchback* is a gene responsible for the formation of anterior structures; on inhibition of these anterior structures, posterior ones are formed. The hunchback gradient controls *caudal*, which is concluded to be a region-specific activator of abdominal segmentation genes (9). Furthermore, injection of *nanos* mRNA into the anterior end induced the development of larvae with two abdomen and no head. Two heads and no abdomen develop in larvae on injection of *bicoid* mRNA into the posterior end of the egg.

In *Drosophila*, axis formation involves a series of interactions between the oocyte and the surrounding somatic follicle cells. The dorsal and ventral sides is determined by the localization of the oocyte [nucleus](#) and *gurken* mRNA to the dorsal-anterior corner of the oocyte. Gurken protein is believed to act as a ligand for the *Drosophila* [epidermal growth factor](#) (EGF) **receptor** (10).

It must be admitted that the naming of the above genes is somewhat peculiar. Many of the genes are not named according to their function, but instead according to the **phenotype** produced when the gene is defective. In this respect, this means that the gene *dorsal* determines the development of ventral structures, not dorsal. The dorsal protein migrates into the nuclei of the ventral side and acts, similar to nanos or bicoid, as a **transcription factor**.

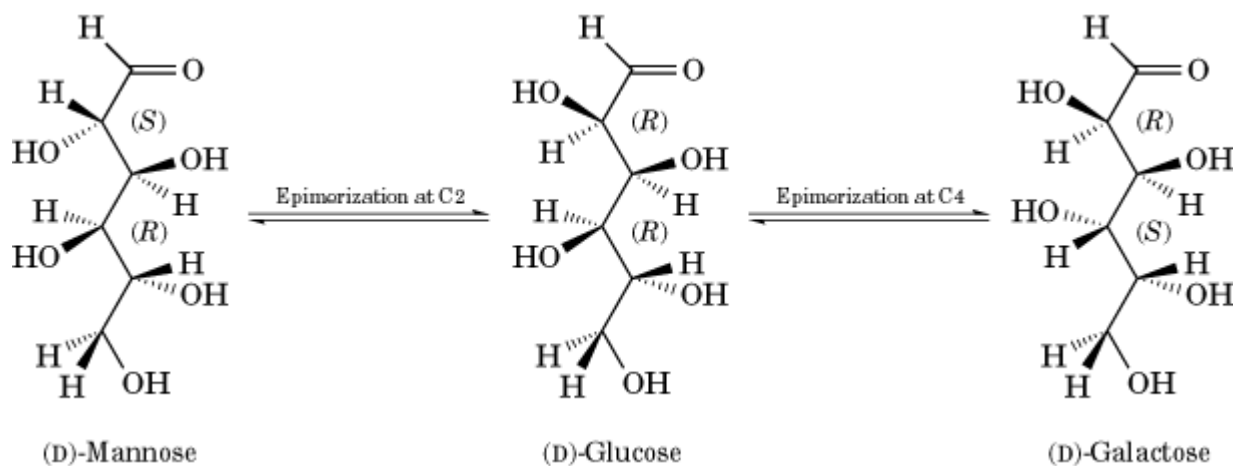
#### Bibliography

1. G. A. Ubbels (1997) Cell Mol. Life Sci. **53**, 382–409.
2. W. Driever and C. Nüsslein-Volhard (1988) Cell **54**, 83–93.
3. X. Ma, D. Yuan, K. Diepold, T. Scarborough, and J. Ma (1996) Development **122**, 1195–1206.
4. S. Kobayashi, M. Yamada, M. Asaoka, and T. Kitamura (1996) Nature **380**, 708–711.
5. E. R. Gavis, and R. Lehmann (1994) Nature **369**, 315–318.
6. W. Driever and C. Nüsslein-Volhard (1989) Nature **337**, 138–143.
7. R. Rivera-Pomar, D. Niessing, U. Schmidt-Ott, W. J. Gehring, and H. Jäckle (1996) Nature **379**, 746–749.
8. D. St. Johnston, D. Beuchle, and C. Nüsslein-Volhard (1991) Cell **66**, 51–63.
9. C. Schulz and D. Tautz (1995) Development **121**, 1023–1028.
10. S. Roth, F. S. Neuman-Silberberg, G. Barcelo, and T. Schupbach (1995) Cell **81**, 967–978.

#### Epimers and Epimerization

Epimers are [diastereomers](#) that are related by the inversion of [configuration](#) at a single **chiral** center ([1](#)). This definition extends the original meaning of epimer, which was used to identify sugars that differed in configuration at C2 ([2](#)). This definition intentionally excludes [enantiomers](#), such as D- and L-alanine, since they are not diastereomers. It also excludes diastereomers that are related by the inversion of more than a single chiral center. Thus, D-glucose and D-mannose are epimers, as are D-glucose and D-galactose. D-mannose and D-galactose are not epimers, however, because they are related by inversion at two chiral centers, C2 and C4 (Fig. [1](#))

**Figure 1.** Stereochemical drawings of glucose, mannose, and galactose, with their four chiral carbons. The configurations at C2 and C4 are labeled and distinguish these three sugars. Glucose is an epimer of both mannose and galactose because they differ by the configuration of a single chiral center. Mannose and galactose have different configurations at both C2 and C4 and are not epimers.



The chemical conversion of one epimer to another is called epimerization. If this interconversion is catalyzed by an enzyme, the enzyme is an epimerase. As an example, UDP-glucose-4-epimerase catalyzes the epimerization of the C4 carbon of glucose. In the reaction, UDP-glucose is epimerized to UDP-galactose. When the inversion of configuration occurs to interconvert enantiomers instead of diastereomers, the reaction is a **racemization**.

#### Bibliography

1. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York p. 40.
2. (1911) Chem. Abstr. V, 1606.

#### Suggestion for Further Reading

3. B. Testa (1982) "The geometry of molecules: Basic principles and nomenclatures". In *Stereochemistry* (C. Tamm, ed.), *New Comprehensive Biochemistry*, Vol. 3, Elsevier, Amsterdam, pp. 22–24.

### Epistasis

In Mendelian genetics, epistasis refers to the situation in which the phenotype caused by a [mutation](#) in one gene (A) masks the phenotype resulting from a mutation in another, nonallelic gene (B). Thus the individual carrying mutations in both genes A and B exhibits the same phenotype as the individual carrying a mutation in gene A alone. In this case, gene A is said to be epistatic over gene B, and gene B hypostatic to gene A. A classic example of epistasis is found among genes that determine mouse coat color. In the presence of albino mutant alleles, the mutant phenotypes of other coat-color determinants are not revealed, due to lack of pigment. Albino mutations, therefore, are epistatic over all other mutations that affect coat-color (1).

An epistatic relationship between two gene functions is often used to infer their relative order within a genetic pathway. Interpretation of such analysis depends on the type of genetic pathways. In biosynthetic pathways in which a series of genes modify a common intermediate substrate, epistatic genes are placed upstream of hypostatic ones. An example of such pathways is the morphogenetic pathway of **T4 phage** for head, tail, and tail-fiber assembly (2). In regulatory gene pathways, in

contrast, where a member gene regulates the expression or activity of a downstream gene or gene product, epistatic genes are generally located downstream. A good example is found in the regulatory gene pathway that determines the dorsoventral body axis of the embryo in the fruit fly, *Drosophila melanogaster* (3). Most loss-of-function mutations that affect this pathway dorsalize the embryo. In contrast, a gain-of-function mutant allele of Toll, known as Toll<sup>D</sup>, has been found that ventralizes the embryo. The presence of the opposite phenotypes allowed the investigators to combine a dorsalizing allele in a tester gene with the ventralizing allele of Toll and ask which of the two alleles is epistatic in double-mutant embryos. If the double mutant shows a dorsalized phenotype (ie, the tester allele is epistatic over Toll<sup>D</sup>), the tester gene is placed downstream of Toll. If the double mutant shows a ventralized phenotype, the tester gene is placed upstream. In this way, the genes in the pathway that affect the embryonic dorsoventral axis have been placed upstream or downstream relative to the Toll gene.

The logic of this type of analysis is straightforward if the pathway is a simple (ie, linear) series of genes or gene products. In such a pathway, the output depends solely on the activity of the terminal member of the pathway, and other members can be thought of as a series of on-off switches that determine the active or inactive state of the output gene. Regulatory gene pathways are not always simple, however, and there are situations in which epistatic genes act upstream. For excellent discussions of more complicated situations, as well as on assumptions and “rules” used in epistasis analysis, see Avery and Wasserman (4).

In population and quantitative genetics, epistasis refers to the genotype value for fitness contributed by gene interaction between two loci. The aggregate genotypic value of an individual ( $G$ ) may deviate from the sum of the genotypic values attributable to the first locus ( $G_A$ ) and the second locus ( $G_B$ ). If it does, the deviation or nonadditive component ( $I_{AB}$ ) is ascribed to epistasis, or interactions, between the two loci (5). Thus,  $G = G_A + G_B + I_{AB}$ . An example of epistasis can be seen when quantitative trait loci (QTL) from lines selected for high numbers of abdominal bristles in *Drosophila* are tested in pairwise combinations. In such tests, several pairs of QTLs for abdominal bristle number showed strong epistasis (interaction), resulting in abdominal bristle numbers that were beyond the additive effects of individual QTL (6).

#### Bibliography

1. W. K. Silvers (1979) *The coat Colors of Mice*, Springer-Verlag, New York.
2. M. Levine (1969) *Annu. Rev. Genet.* **3**, 323–342.
3. K. V. Anderson et al. (1985) *Cell* **42**, 779–789.
4. L. Avery and S. Wasserman (1992) *Trends Genet.* **8**, 312–316.
5. D. S. Falconer and T. F. C. Mackay (1996) *Introduction to Quantitative Genetics*, 4th ed, Longman, Essex, England.
6. A. D. Long et al. (1995) *Genetics* **139**, 1273–1291.

#### Suggestions for Further Reading

7. L. Avery and S. Wasserman (1992) *Trends Genet.* **8**, 312–316.
8. M. Wade (1992) In E. F. Keller and E. A. Lloyd, eds. *Keywords in Evolutionary Biology*, Harvard University Press.
9. T. F. C. Mackay (1996) *BioEssays* **18**, 113–121.



Epitope designates an area of an [antigen](#) that interacts with a specific [antibody](#). It has generally replaced the former name of antigenic determinant. By symmetry, the term *paratope* has been proposed for the antibody, instead of antibody combining site (also designated as antigen binding site), but is much less widely used. An extension of this terminology is to be found as idiotopes and allotopes that designate epitopes that define **idiotypic** and allotypic (see **Allotype, alloantigen**) specificities, respectively.

Epitopes are defined functionally by the interaction of the antigen with a particular antibody, so it does not represent a preformed antigenic structure in itself; instead, it is defined by the immune system. This means that identification of epitopes will depend upon the nature and the number of **B-cell** clones that will respond to the administration of antigen. It is therefore very likely that two individuals will have a different pattern of recognition of the same antigen, although large overlaps are the rule. A consequence of this is that the best tool to define an epitope is a [monoclonal antibody](#), which constitutes a reliable and standardized reagent. The obvious requisite for a region of the antigen to be recognized as an epitope is its accessibility—that is, that it be in contact with the solvent, which is favored by the presence of [hydrophilic](#) and/or charged amino acids.

Defining the structural basis of an epitope on natural antigens, such as [proteins](#), is an almost impossible task, because B cells, and hence antibodies, recognize epitopes directly as native structures of the antigen. This implies that most epitopes have a complex three-dimensional structural organization that is generally lost upon cleavage of the antigen. Limited [proteolysis](#) of a protein may retain some regions that still contain native epitopes, but it is quite obvious that isolation of a “pure” epitope is most often not feasible. Interesting attempts to isolate epitopes from natural proteins have nevertheless been made with some success. This is the case of [myoglobin](#) and [lysozyme](#), for which a loop of 20 amino acids clamped by a [disulfide bond](#) has long been a prototype of an “isolated” epitope. Because of their repetitive structural organization, polysaccharides are somewhat better tools to derive epitopes. This was the case for dextran, a polymer of glucose that allowed Kabat to demonstrate that various oligosaccharides were competing with the binding of dextran to its antibodies and indicating that they could be considered epitopes of this antigen. In fact, because the optimal inhibition was obtained by a hexasaccharide, it was concluded that the size of the anti-dextran combining site was that required to accommodate a molecule having the size of this inhibitor. Most of our knowledge of the antigenicity of proteins has been gained thus far with synthetic [polyamino acids](#) that constitute simplified models of proteins, because they can have a rather monotonous organization; for example, three or four amino acids might be linked covalently to a backbone of poly-L-lysine, like those extensively studied by the group of Sela in Israel. These models favored, however, sequential epitopes, as opposed to the conformational ones that are commonly encountered in natural protein antigens. The use of haptens was of course only a first approach to define an antigenic determinant, and they gave information of interest regarding the exquisite specificity of antibody recognition.

**Monoclonal antibodies** (mAbs) have proven of great value in epitope mapping the surface of an antigen. Mapping is based on inhibition of one mAb by others. When complete and reciprocal inhibition is observed between two mAbs, it is concluded that both antibodies recognize the same epitope. Although not a definitive proof, it remains an acceptable conclusion. Partial inhibition is suggestive that the two mAbs recognize overlapping epitopes. Consequently, a sufficiently large collection of mAbs permits exhaustive mapping of the epitopes, which may cover the entire surface of the antigen. One must realize, however, two important points: (1) Some regions are more frequently recognized than others; these are the immunodominant epitopes; and (2) the number of mAbs that can bind simultaneously to a given antigen molecule is necessarily limited by the steric hindrance imposed by the surface of the antibody combining site that interacts with the epitope. Having said that, it is quite obvious that the overall number of epitopes that can be detected on one antigen molecule is directly related to its surface area.

By extension, the notion of **T-cell** epitopes has been defined. From a physiological point of view, it is somewhat more complicated in that it is linked directly to the mechanism of **antigen presentation** and thus requires the interaction of the antigenic structure with two partners, the **multiple histocompatibility complex** (MHC) molecules and the **T-cell receptor (TCR)** itself. Because antigen processing and presentation have been defined for only proteins thus far, such T-cell epitopes are **peptides**. In a sense, the structural basis for a T-cell epitope is simpler and more directly accessible than in the case of B cells, because these peptides are short and are recognized as a sequential determinant. Strictly speaking, however, the term *epitope* is reserved to that part of the peptide that interacts with the TCR, whereas regions that interact with the MHC molecule, which are of decisive importance, because it conditions immunogenicity, are described sometimes as the “agretope.” Extensive analysis by [site-directed mutagenesis](#) has clearly demonstrated which key residues contributed the epitope or the agretope. Because of the differences in the structures of the cavity that bind the peptides in various MHC molecules, the size of the peptide is limited to 9 amino acid residues for those that bind to class I MHC molecules, whereas it may be somewhat larger for those that are presented by class II. A very large number of such small overlapping peptides may be derived from one protein antigen. The problem of recognition by the MHC molecules, which are in very limited number for any given individual (a few discrete molecules), is therefore very different from the huge [repertoire](#) expected from Ig and TCR molecules. What happens is that MHC molecules will “choose” only a very few peptides from one antigen, those that fit with the minimum requirements that drive the interaction with the agretope. This is why every individual will use a different way to respond to the same antigen. Using a simplified model of antigen, such as a monotonous synthetic polypeptide, the chance is high that the molecule will be **immunogenic** only in those individuals with the appropriate MHC molecules. This was the case of the strains of mice injected with synthetic polypeptides that did or did not respond, initiating the first approach to the understanding of the genetic control of the immune response and starting the elucidation of the physiological role of MHC molecules.

Both the B-cell and T-cell epitopes raise directly the problem of the repertoire of the immune system, a central notion that is directly related to the number of potential antibodies and TCR that a given organism can potentially make (see [Repertoire](#)).

See also entries [Antigen](#), [Hapten](#), [Immunogen](#), and [Repertoire](#).

#### Suggestion for Further Reading

A. C. Horsfall, F. C. Hay, A. J. Soltys, and M. G. Lones (1991) Epitope mapping. *Immunol. Today* **12**, 211–213.

## Epstein-Barr Virus

Epstein–Barr virus (EBV) is a ubiquitous human [herpesvirus](#) that is the causative agent of infectious mononucleosis and associated with several malignancies, such as Burkitt's lymphoma (BL), nasopharyngeal carcinoma (NPC), some cases of peripheral T-cell lymphoma and gastric carcinoma, and B-cell lymphoma of acquired immune deficiency syndrome (AIDS ) and transplant patients. EBV has the ability to immortalize human B-lymphocytes *in vitro*.

The EBV [genome](#) is double-stranded linear DNA with an approximate length of 175 kbp and a G+C content of 59%. Reiterated [direct repeat](#) sequences at both termini (TR) and within the genome (IR)

separate the EBV genome into five unique regions. The EBV genome circularizes after entry into cells, and it is maintained stably as a **plasmid**.

Although the entire EBV genome is maintained in immortalized lymphocytes, virus replication usually does not occur in these cells. During a latent infection, six nuclear antigens (EBNA1, EBNA2, EBNA3A, EBNA3B, EBNA3C, and EBNA-LP), three membrane proteins (LMP1, LMP2A, and LMP2B), and two small nuclear RNAs (EBER1 and EBER2) are expressed. EBNA-LP, 2, 3A, and 3C and LMP1 are essential or critical for primary B-lymphocyte growth transformation. EBNA3B, LMP2A, LMP2B, and EBERs are dispensable for immortalization.

EBNA1 from a typical EBV strain B95-8 consists of 641 amino acid residues. A long hydrophilic **domain** of residues 459 to 607 has a **DNA-binding** activity that is nucleotide sequence-specific. The specific EBNA1 binding sequence is located at oriP, which is the origin of EBV plasmid replication. Replication from oriP requires both the **cis-acting** elements (the family of repeats and the dyad symmetry elements) and the viral protein EBNA1. EBNA1 binding to oriP is also required for efficient EBNA **transcription**, where oriP serves as an **enhancer**. EBNA2 induces expression of viral genes LMP1 and LMP2 and cellular genes CD21, CD23, and c-fgr. EBNA2 responsive elements have been defined in the upstream **promoter** regions of these genes. Interaction of EBNA2 with its **response element** is mediated by a cell protein Jk, which recognizes a GTGGGAA sequence present in all known EBNA2 response elements. EBNA3A, EBNA3B, and EBNA3C are encoded by three genes that are tandemly placed in the EBV genome and are likely to have been derived a common origin. EBNA-LP is encoded by the leader of each of the EBNA **messenger RNAs** and is **translated** from these mRNAs when the first and second exons are spliced, so as to create the EBNA-LP **initiation codon**. Deletion of the EBNA-LP gene reduces greatly the efficiency of B-lymphocyte immortalization by EBV.

LMP1 consists of an amino-terminal cytoplasmic domain, six **membrane-spanning hydrophobic** domains separated by short reverse **turns**, and a carboxy-terminal cytoplasmic domain. LMP1 has the ability to transform rodent fibroblasts, and it blocks differentiation in human keratinocytes. In transgenic mice, LMP1 expression causes epithelial hyperplasia. In BL cells, LMP1 induces villous projections, growth in tight clumps, NFkB activity, and expression of activation markers (CD23 and CD40), **cell adhesion molecules** (ICAM-1, LFA-1, and LFA-3), IL-10, and the Bcl-2 proto-**oncogene**. There are two NFkB-activating regions in the carboxy-terminal cytoplasmic domain. One of these regions associates with members of the **tumor necrosis factor** receptor-associated factor (TRAF) family of proteins. CD40 is a well-known mediator of growth of B lymphocytes. LMP1, through direct interaction with TRAFs, may constitutively activate CD40 signal transduction and promote cell growth. LMP2A is a substrate for B lymphocyte **src** family **tyrosine kinases**. LMP2A expression blocks calcium mobilization in B lymphocytes. There is no known function for LMP2B.

EBERs are RNAs that are not **polyadenylated**, and they have extensive primary sequence similarity to **adenovirus** VA1 and VA2 and cellular U6 small RNAs.

In contrast to the EBV gene expression in immortalized lymphocytes, more limited numbers of EBV genes are expressed in EBV-associated tumor cells. BL and gastric carcinoma cells express EBNA1, but do not express other EBNAs and LMPs. NPC and T-cell lymphoma cells express EBNA1 and LMPs, but do not express other EBNAs. It is known that EBNAs other than EBNA1 become targets of **cytotoxic T lymphocytes**. Therefore, it is important for tumor cells not to express these antigens so as to survive under normal immunity.

EBNA2 and LMP1 are particularly important for the immortalization of primary lymphocytes. These proteins, however, are not expressed in most EBV-associated malignancies, although there is a report suggesting that malignant phenotypes are dependent on the presence of EBV genomes. These observations may be extended to the pathogenic role of EBV in other EBV-associated malignancies.

Suggestions for Further Reading

- E. Kieff (1996) "Epstein–Barr Virus and Its Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2343–2396.
- A. B. Rickinson and E. Kieff (1996) "Epstein–Barr Virus". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2397–2446.

## Equilibrium Dialysis

Equilibrium dialysis is a technique for analyzing the binding of a low molecular weight ligand to a macromolecular **receptor** (1) (see **Ligand binding**). Information that can be obtained from an equilibrium dialysis experiment includes the stoichiometry of the complex formed and the affinities of the interacting components (2), plus more subtle properties such as **cooperativity**. In the simplest form of an equilibrium dialysis experiment, receptor and ligand solutions are placed in compartments on opposite sides of a semipermeable **dialysis** membrane. The pore diameter of the membrane is chosen to allow passage of the ligand and prevent passage of the receptor. The ligand thus redistributes between the two compartments, whereas the receptor stays in its own compartment. If the receptor binds the ligand, an excess of ligand begins to accumulate in the receptor compartment. When equilibrium is reached, two conditions are met, as long as there is no nonideality in the system: (1) the free ligand concentration is equal on both sides of the membrane and (2) the ligand in the receptor compartment is partitioned between free and receptor-bound form according to the equilibrium constant for the receptor–ligand interaction. The total ligand concentration on each side of the membrane is measured after equilibration, as is the receptor concentration, if necessary. The excess ligand concentration on the receptor side is attributed to its binding to the receptor. This basic experiment is repeated for a series of receptor and ligand concentrations. Pairs of bound and free ligand concentrations may then be analyzed by a graphical method, such as a **Scatchard Plot**, or fit by least squares to a binding equation to derive the relevant parameters of the receptor–ligand interaction (3).

The primary advantage of equilibrium dialysis over other techniques for measuring ligand-binding equilibrium constants is that equilibrium dialysis does not rely on assumptions that a measured property, for example, a spectral change, correlates linearly with receptor occupancy. Development of microdialysis cells and of dialysis membranes made from purified material with well-defined molecular weight cutoffs have eliminated many artifacts once associated with equilibrium dialysis. Persisting sources of error include aggregation of the ligand, which will retard or prevent its passage through the dialysis membrane, adsorption of ligand or receptor to the membrane, and changes in volume and concentration resulting from an osmotic imbalance at the start of the experiment.

The assumption that the free ligand concentration is equal on both sides of the membrane at equilibrium is not always valid. Most **proteins** and all **nucleic acids** are charged; thus these macromolecules accumulate a set of neutralizing counterions. A ligand of opposite charge can act as a neutralizing counterion, resulting in an imbalance between free ligand concentrations in the receptor and ligand compartments of an equilibrium dialysis experiment. This imbalance, known as the Donnan effect (4), does not reflect a biologically specific ligand–receptor interaction, but occurs with whatever counter-ion is available, as a consequence of maintaining the overall electrical neutrality of the solution. The practical impact of the Donnan effect is to give spurious evidence of an association between a ligand and receptor of opposite charge. Even if a ligand is uncharged, changes in ionic strength or pH due to redistribution of buffer ions may still interfere indirectly with an experiment. The magnitude of the Donnan ratio (an ion's concentration in the receptor compartment divided by its concentration in the ligand compartment) is greatest at high receptor

concentration and low ionic strength. For example, 100  $\mu\text{M}$  receptor containing 10 charges in the presence of 1  $\text{mM}$  NaCl will yield a Donnan ratio of 1.5, meaning that a monovalent ligand will appear to be 50% enriched in the receptor compartment. If the receptor concentration is lowered to 1  $\mu\text{M}$ , or the salt concentration raised to 100  $\text{mM}$ , the Donnan ratio drops to 1.005. Although simple addition of salt essentially eliminates the Donnan effect as an interfering factor in equilibrium dialysis, the added ions may distort the experiment in a different way, by competing with the ligand for association with charged groups in the receptor binding site. An alternative to using high salt concentration is to measure the Donnan ratio directly and apply a correction factor (5).

### Bibliography

1. I. M. Klotz, F. M. Walker, and R. B. Pivan (1946) *J. Am. Chem. Soc.* **68**, 1486–1490.
2. H. N. Eisen and F. Karush (1949) *J. Am. Chem. Soc.* **71**, 363–364.
3. J. E. Fletcher and A. A. Spector (1968) *Comput. Biomed. Res.* **2**, 164–175.
4. F. G. Donnan (1911) *Z. Elektrochem.* **17**, 572.
5. P. Suter and J. P. Rosenbusch (1977) *Anal. Biochem.* **82**, 109–114.

### Suggestions for Further Reading

6. N. W. Huh, P. Berkowitz, R. G. Hiskey, and L. G. Pedersen (1991) Determination of strontium binding to macromolecules, *Anal. Biochem.* **198**, 391–393. (A modern example of equilibrium dialysis, illustrating methodology and some of the technical problems that can be encountered.)
7. I. M. Klotz (1953) "Protein interactions", in *The Proteins*, H. Neurath and K. Bailey, eds., Academic Press, New York, Chapter "8", pp. 727–806. (Early review of equilibrium dialysis.)
8. C. Tanford (1961) *Physical Chemistry of Macromolecules*, Wiley, New York. (Quantitative treatment of the magnitude of the Donnan effect is given on pp. 221–227.)

## Equilibrium Potential

Current flows resulting from ions diffusing across a phase boundary establish an electrical potential difference across that phase boundary in opposition to the ion currents. The origin of this potential lies in the difference in mobilities of anions and cations as they diffuse across the phase boundary. The electric field slows the faster ion and speeds up the slower ion, so that no net transfer of charge occurs across the phase boundary. The magnitude of the diffusion potential will depend on the difference in chemical potentials of the electrolytes in the contiguous phases, as well as differences in ionic mobility. The *Planck–Henderson equation* provides a quantitative calculation for the diffusion potential  $\Psi_D$ ; for a single electrolyte (eg, NaCl) we have

$$\Psi_D = -\frac{(u - v)RT}{(u + v)zF} \ln \frac{C_2}{C_1} \quad (1)$$

where  $u$  and  $v$  are the cationic and anionic mobilities, respectively,  $R$  is the gas constant,  $T$  is the absolute temperature,  $z$  is the valence of the electrolyte,  $F$  is the Faraday (96,500 C/mol), and  $C_i$  are the electrolyte concentrations in the two contiguous solutions. Two examples of the use of equation 1 are for (i) the *liquid junction potential* that arises when two electrolyte solutions with different concentrations form an interface and (ii) when a synthetic membrane with a selective permeability for one of the ions is inserted between the two solutions. For the ion-specific selectively permeable membrane, equation 1 becomes

$$\Psi_D = \frac{RT}{zF} \ln \frac{C_2}{C_1} \quad (2)$$

This is the equilibrium potential for the permeable ion (see [Membrane Potentials](#)).

#### Suggestions for Further Reading

D. A. MacInnes (1939) *The Principles of Electrochemistry*, Rheinhold, New York, Chapter "13". (This book presents a derivation of the Planck–Henderson diffusion equation, as well as an excellent description of the characterization of liquid junction potentials and their measurement.)

W. J. Edelman Jr.(ed.) (1971) *Biophysics and Physiology of Excitable Membranes*, Van Nostrand Reinhold, New York. (Contains a useful summary of the physical aspects of membrane potentials and their treatment by analytical methods.)

## Estrogen Receptors

Estrogen, together with the other sex steroid hormones, progesterone and testosterone, provokes the development and determination of the embryonic reproductive system, masculinizes or feminizes the brain at birth, and controls reproduction and reproductive behaviour in the adult, plus the development of secondary sexual characteristics. Furthermore, both an atheroprotective effect of estrogen and a protective effect against osteoporosis are well documented. The hormone probably enters its target cells by diffusion (see [Steroid Hormone Receptors](#)) and activates estrogen receptors (ERs) that are paradigmatic [transcription factors](#) in eukaryotes. ERs are members of the steroid hormone receptor family that, in turn, is a subfamily of the nuclear receptor superfamily of proteins (see [Steroid Hormone Receptors](#)). Effects of ERs on the [genome](#) are mediated through binding of the hormone–receptor complex to short, specific DNA sequences, termed estrogen [response elements](#) (EREs), located in the vicinity of approximately 100 estrogen-regulated genes in each cell, and subsequent modulation of gene expression (see [Hormone Response Elements](#) and [Glucocorticoid Response Element](#)). The ERE functions as an [enhancer](#) and the ERs as enhancer binding factors that are equipped with domains and surfaces responsible for hormone and DNA binding, activation of [transcription](#), and interaction with other nuclear partners.

Until recently, only one oestrogen receptor had been known, which is now called estrogen receptor a (ERa or NR3A1 according to a new nomenclature proposal (see [Steroid Hormone Receptors](#))). In 1996 a second estrogen receptor, now known as estrogen receptor b (ERb or NR3A2), was identified ([1](#)). Therefore, older literature always refers to ERa and has often to be reinterpreted in the light of at least two receptors being present in many tissues.

Apart from their involvement in the development of the reproductive system, estrogens are also implicated in the growth of breast cancers. Already in 1896, Beatson observed the regression of metastatic breast cancer after ovariectomy (*ce, after estrogen withdrawal*). Since that time, various surgical and pharmacological endocrine therapies have been employed (see below). Today, the determination of the concentration of sex steroid receptors in breast cancers is performed routinely and is very important for prognosis in general and the therapeutic regimen in particular.

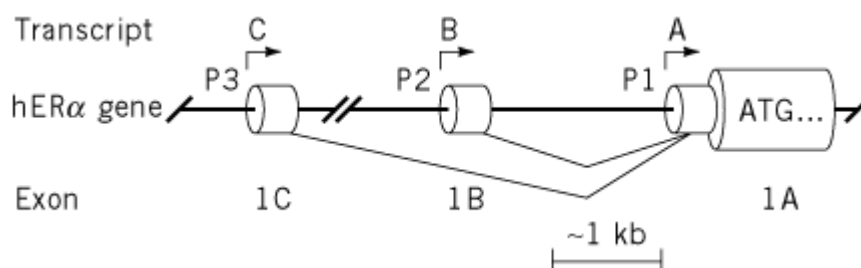
The molecular mechanisms by which ovarian hormones stimulate growth of breast cancers are unclear, but they are likely to involve an extensive crosstalk of ERs with other signalling pathways. ERa, for example, activates the **signal-transducing** Src/p21ras/Erk pathway in human breast cancer

cells via an interaction with c-Src (see [Cis Configuration](#)). In addition, ERa can integrate signals from the progesterone receptor into this signalling cascade (see [Steroid Hormone Receptors](#)). ERa is also activated when it is complexed with **cyclin D1**, which has an important role in the development of breast cancer and is required for normal breast cell proliferation and differentiation associated with pregnancy. Furthermore, such growth factors like [epidermal growth factor](#) (EGF), can activate ERa through the mitogen-activated protein, [MAP Kinases](#) pathway, which probably involves receptor phosphorylation.

## 1. Domain Structure

The human ERa was cloned from [cDNA libraries](#) prepared from the MCF-7 breast cancer cell line in 1985. This cDNA clone was later shown to contain an artifactual point mutation that resulted in the substitution of a valine (mutant receptor called HEO) for a glycine (wild-type receptor called HEGO). The hERa gene encompasses 10 exons spread over 140 kbp and is regulated by at least three different **promoters** (see Fig. 1). The ERa cDNA encodes a 595–amino acid residue protein of 66 kDa. The oestrogen receptor is a modular protein composed of distinct regions corresponding to functional and structural units called **domains** that are labeled A–F (see [Steroid Hormone Receptors](#), Fig. 1B). Transcription activation functions have been detected in regions A/B (AF-1, also called q-1 or enh-1) and E (AF-2, also called q-2 or enh-2). Region C is the DNA-binding domain (DBD), and region D is often considered as a flexible hinge region between the DNA- and ligand-binding domains. The very amino terminus of region D is also an integral part of the DBD. Region E is the ligand-(hormone)-binding domain (LBD), composed mainly of 12  $\alpha$ -helices, and contains AF-2 that forms the core of the most carboxy-terminal helix 12 (see text below). Beyond helix 12 starts region F, which is not required for hormone binding but seems to be important for the discrimination between agonistic and antagonistic ligands.

**Figure 1.** The 5' region of the hERa gene. The promoter P1 in front of exon 1A generates transcript A encompassing a 5'-untranslated region (UTR, smaller cylinder) and the start of the coding region (larger cylinder). The other two promoters, P2 and P3, generate transcripts in which the noncoding exons 1C or 1B are alternatively spliced to the same acceptor splice site within a short upstream open reading frame found in the 5'-UTR of transcript A. Depending on tissue-specific promoter utilization, different hERa transcripts can be detected in different tissues that code for the same receptor protein, but have distinct 5'-UTRs. Common to all UTRs is the presence of short upstream open reading frames that are implicated in the regulation of translation efficiency.



Regions C and E are not only responsible for DNA- and ligand-binding, respectively, but they encode other functions as well. At equilibrium the majority of ERa is in the nucleus, due to the presence of so called nuclear localization signals (NLS) that are believed to be required for **nuclear pore** recognition (see [Nuclear Import, Export](#)). A constitutive NLS is located at the border of region C to D, whereas a second NLS in the LBD is ligand-dependent.

ERb was found in rat prostate and later in human testis cDNA libraries using **PCR**-cloning strategies involving degenerate primers. With 8 exons spread over approximately 40 kbp, the human ERb gene is considerably smaller than the ERa gene. Several mRNA bands have been observed on [RNA Blots](#)

([Northern Blots](#)), possibly indicating that the ERb gene is also characterized by multiple promoters, which is a general phenomenon within the steroid hormone receptor family. The human cDNA encodes a 485-residue protein that shows a high degree of [homology](#) to ERa within the DBD (96% identity) and the LBD (58%).

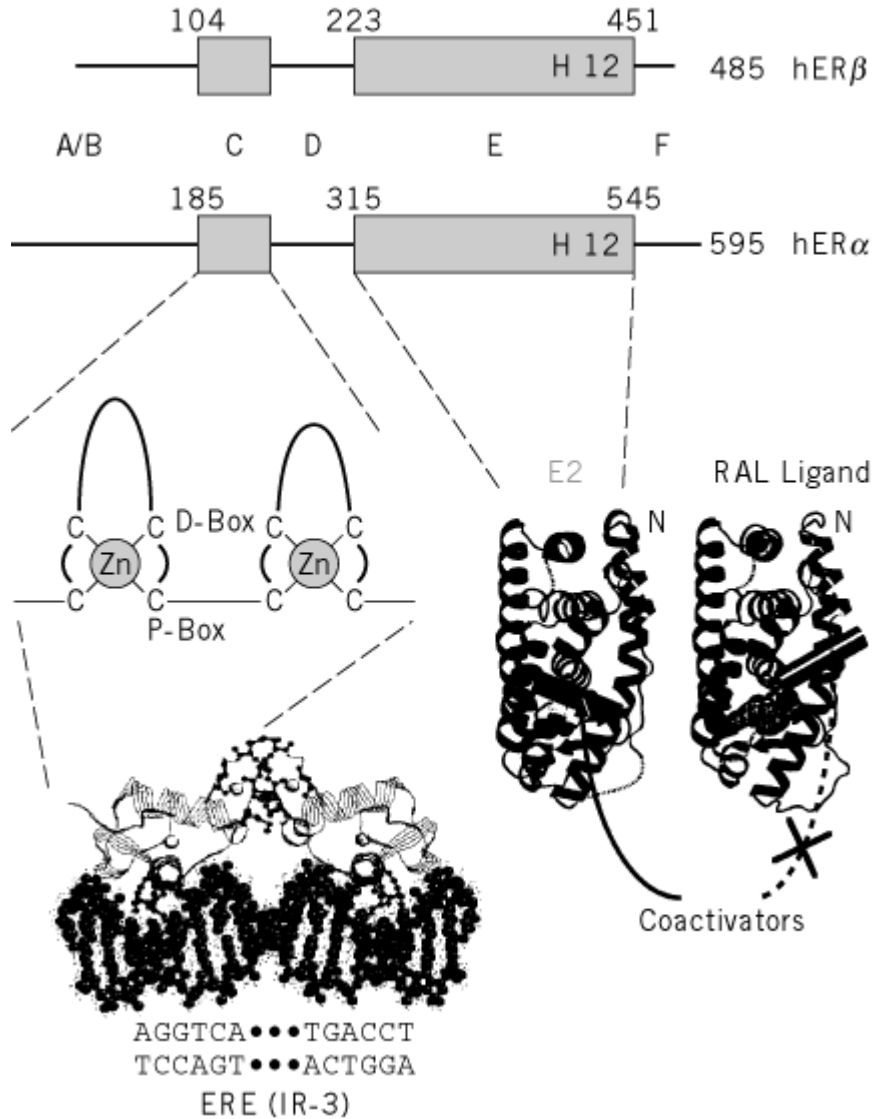
## 2. DNA Binding

All members of the steroid hormone receptor family were initially thought to bind to their cognate response elements as homodimers. Today the family also contains two oestrogen-related orphan receptors that bind to DNA as monomers (see [Steroid Hormone Receptors](#)). Moreover, the two closely related ERs, a and b, are known to bind to DNA as homodimers, but they can also form heterodimers which is an otherwise common mechanism within the nuclear receptor superfamily. Oestrogen receptor homo- or heterodimers bind to the estrogen response element (ERE) that is composed of two half-sites (AGGTCA) arranged as palindromic, [inverted repeats](#) (IRs), where the half-sites are separated by three nucleotides (IR-3). One ERE half-site is recognized by one receptor monomer. The half-site spacing of three base pairs is such that an ER homodimer binds both half sites on the same face of the DNA. The DNA-binding domain of the ERs comprises approximately 80 amino acid residues encoded by region C, plus some 14 *N*-terminal amino acids of region D. The domain contains two patterns that are reminiscent of, but clearly distinguishable from, the **zinc finger** motifs first observed in the *Xenopus laevis* transcription factor IIIA (TFIIIA; see **Zinc fingers**). The ER zinc fingers are also able to coordinate a zinc atom tetrahedrally but are of the type C<sub>2</sub>-C<sub>2</sub> (see Fig. 2). Each ER monomer contains two of these particular zinc fingers, which are folded together to form a single structural domain (2). Characteristic of this fold are two highly [amphipathic](#)  $\alpha$ -helices, one contributed by each finger, that cross at an approximately right angle at their midpoints (see Fig. 2). The **hydrophobic** surfaces of these helices form a rigid core decisive for the stability of the structure that is further strengthened by two zinc ions holding the base of a loop against the *N*-terminus of each helix. The specific interaction of an ER homodimer with DNA is determined by only a few amino acid residues located in two regions, the proximal box (P-box; see Fig. 2) and the distal box (D-box). The P-box lies at the amino terminus of the first fingers  $\alpha$ -helix (its “knuckle”). Four amino acid residues within the P-box interact with base pairs, and only three of them are required for the discrimination between an ERE and a nonestrogen steroid hormone receptor response element. When these three residues are exchanged for the three cognate residues of the glucocorticoid receptor, ERE-specific binding is converted to glucocorticoid response element-specific binding. The five D-box residues form the weak dimerization domain within the DBD and are responsible for the recognition of the correct half-site spacing (IR-3). These dimerization forces are too weak for the DBD alone to form dimers in solution, but they make DNA-binding cooperative. The complete receptor protein exists as a homodimer in solution, due to a strong dimerisation domain within the LBD.

**Figure 2.** Structure of estrogen receptors. (**Top**) The functional domains (A/B to F) of hERa and b are illustrated at the primary structure level. The numbers refer to amino acid positions. Highlighted are the DNA-binding domain (DBD) in small boxes and the ligand binding domain (LBD) in larger boxes. The A/B and F domains are drawn to scale. (**Middle left**) In an exploded view of the hERa-DBD its secondary structure, characterized by two steroid hormone receptor-specific zinc fingers, is depicted. Four cysteines each tetrahedrally coordinate two zinc ions. The proximal box (P-box) responsible for specific DNA-recognition is shown on left; the distal box (D-box) mediating DBD-dimerization is shown on right. (**Bottom left**) The secondary structure explodes into a view of the crystal structure of the hERa's DBD bound to DNA. The four amino acid side chains of the P-box interacting with DNA bases are again shown in boldface. These residues are part of an  $\alpha$ -helix responsible for specific DNA-sequence recognition that is positioned within the DNAs major groove (perpendicular to the paper plain). A second amphipathic  $\alpha$ -helix can be seen to cross the recognition helix with the red P-box residues at a right angle. The D-box residues promoting DBD dimerization are located between the two helices. Underneath the DNA the sequence of the consensus estrogen response element (ERE) is shown. Both half-sites are organised as an inverted repeat with 3-bp spacing (IR-3). (**Middle right**) The crystal structure of the ligand-binding domain complexed with the natural ligand 17 $\beta$ -estradiol (E2, left) and the antihormone raloxifene (RAL, right). Decisive for a productive interaction with coactivators is the proper positioning of helix 12 (H12) over the ligand-binding pocket. (The authors thank R. Hubbard and A. Pike,



University of York, UK, for providing this part of the figure.)



### 3. Ligand Binding

Estrogenic hormones, as well as synthetic antihormones, enter the target cell by simple diffusion and bind, within the cytoplasm, to a multiprotein complex of **molecular chaperones** and ER. Hormone binding induces a transformation of the ER complexes that is associated with an increase in affinity for DNA and leads to a tight association with the nuclear compartment. In addition to DNA-binding, the bound hormone confers transcriptional competence onto ERs that is exerted by two independent **transactivation** functions (see text above). The crystal structure (3) of the hERα ligand binding domain complexed with the endogenous hormone 17β-estradiol (PDB accession no. 1ERE; see [Structure Databases](#)) and the antihormone raloxifene (PDB accession no. 1ERR) have been determined (see Fig. 2). The LBD of hERα is folded into a three-layered antiparallel α-helical sandwich that creates a wedge-shaped molecular scaffold, in which the ligand-binding cavity is located at the narrower end of the domain. This cavity is completely partitioned from the external environment and is closed by helix 12 operating as a “lid” after 17β-estradiol has entered the binding pocket (see Fig. 2, left LBD structure). In “lid-on” position, AF-2 is able to interact with coactivators that in turn connect the receptors to the transcription machinery and activate transcription (for more details, see [Steroid Hormone Receptors](#)). When an antihormone-like raloxifene is bound instead, closure of the ligand binding pocket is not possible, and AF-2 cannot be placed in a

transcriptionally competent position (see Fig 2, right LBD structure).

Although the precise mechanism of this ligand induced switch between active and inactive AF-2 was not known until recently, a number of proteins have been rendered functionally estrogen-dependent by fusion with the hormone-binding domain of ERa (4).

#### 4. Hormones and Antihormones in Breast Cancer

Since Beatson's observation of regression of metastatic breast cancer after ovariectomy, in 1896, oncologists have been searching for an ideal antihormone that will prevent the proliferative effect of estrogen on breast cancer but preserves an atheroprotective effect and protects against osteoporosis without any side effects. Approximately 50% of all breast cancers contain clearly detectable amounts of hERa and are, thus, amenable to hormone therapy. Disease prognosis is better when receptors are present, probably because their presence reflects a greater degree of differentiation of the tumor, indicative of slower growth. Only approximately 50% of these receptor-positive cases also respond positively to an antihormone treatment. Several ER sequence variants were identified in human breast cancer cell lines and tumor specimens, but it is not known if these mutations are important for breast cancer progression or interfere with endocrine breast cancer therapies.

Currently, the most commonly used antihormone is tamoxifen, which is converted to 4-hydroxy-tamoxifen *in vivo*, the actual antihormone. Tamoxifen appears to antagonize estrogen action by competing with oestrogens for receptor binding. Receptor-bound tamoxifen, like raloxifene, is postulated to be unable to position helix 12 of the LBD in such a way that AF-2 is placed in a transcriptionally competent position (see discussion above). However, tamoxifen can still act as a partial estrogen agonist depending on the tissue and response examined. This is probably due to the ligand-independent and cell-specific activity of AF-1 within the receptor's A/B domain. Tamoxifen is clearly effective in receptor-positive breast cancer patients and has very few side effects. It still retains a protective effect against osteoporosis, but does not reduce the risk of cardiovascular disease. Its effectiveness was also proved in prevention trials, although a slightly increased risk for developing endometrial cancer and thromboembolic events must be taken into consideration. Other antiestrogens with a broader spectrum of positive effects than tamoxifen, such as raloxifene, are currently being tested in clinical trials.

With the discovery of a second estrogen receptor, all clinical data on antihormone therapies of receptor-positive breast cancer patients must be reevaluated. It is already known that some breast cancers express ERb but no ERa, and some express both receptors. Differing ERb/ERa ratios could explain cases that were typed receptor-negative but still responded positively toward antihormone therapy, as well as those receptor-positive cases that were refractory toward antihormone therapy. Therefore, the presence of both receptors in tumor specimens will probably be determined routinely in the near future.

#### 5. Results from Disruption of the Mouse Estrogen Receptor a Gene

In addition to its importance for a normal postnatal sexual development, estrogen action has always been considered as a crucial factor for prenatal development. This hypothesis is generally supported by the fact that almost no ERa mutations are known. An exceptional case, however, is a 28-year-old tall man in whom a premature stop codon in position 153 of hERa could be identified (5). The patient is homozygous for this mutation, showed no response to estrogen administration, and manifested decreased bone mineral density and increased bone turnover. This homozygous human knockout already indicates that ERa is not required for prenatal development but is very important for bone maturation and mineralization in males also. Targeted disruption of the ERa gene in mice (ERKO mice (6)) showed the same result (see [Gene Targeting](#)). Homozygous ERKO mice are viable, and the external phenotypes of both sexes are normal. Internally only the females show noticeable gross differences from normal, with hypoplastic uteri as well as hemorrhagic cystic ovaries that lack corpora lutea. In addition, the mating behavior of the females is severely

compromised, because they show no typical lordosis posture or receptiveness to wild-type males, even when treated with estrogen. Consequently, all homozygous ERKO females are infertile, as are the male mice. The latter show normal motivation to mount females, but they achieve less intromissions and virtually no ejaculations. Aggressive behaviors were dramatically reduced, and male-typical offensive attacks were rarely displayed by ERKO males. In addition to an altered mating behavior, spermatogenesis is disturbed, as manifested by low sperm numbers and defective sperm function. Therefore, expression of the ER $\alpha$  gene appears to be necessary not only for the development of the female and male reproductive systems but also for the full development of both female and male sex behavior repertoires.

With the discovery of ER $\beta$ , the question of the importance of estrogen receptors for prenatal development is open again and has to await the production of ER $\beta$  single- and ER $\alpha$ /ER $\beta$  double-knockout mice.

## 6. Functional Differences between ER $\alpha$ and ER $\beta$

In light of the high homology between ER $\alpha$  and ER $\beta$ , the question arises as to whether one receptor is a mere “backup” of the other, or whether each receptor has its own set of specific functions. ER $\beta$  was shown to have DNA-binding specificity overlapping with that of ER $\alpha$  and activates transcription of reporter gene constructs containing EREs in transient transfections in response to estradiol. The response of both receptors to a number of estrogen agonists and antagonists has been determined in such a way, and numerous differences have been observed. Furthermore, the tissue distribution of both receptors is clearly different. Whereas some tissues contain both receptors, as does the human mammary gland, others contain only a single receptor type. One prime example are the granulosa cells of the ovaries, which contain only ER $\beta$  and are devoid of ER $\alpha$ . Previously, the granulosa cells had been identified as a target for oestrogen action, although they were shown not to contain ER $\alpha$ . This puzzling finding can now be explained by the presence of ER $\beta$ . Other examples are the Leydig cells of the testis, which contain ER $\alpha$  but no ER $\beta$ , whereas the developing spermatids express ER $\beta$  and ER $\alpha$  is absent. Therefore, with the advent of a second estrogen receptor, at least two more important parameters have to be considered when analyzing the regulation of gene expression by estrogen: the ER $\alpha$ /ER $\beta$ - ratio and the presence of ER $\alpha$ /ER $\beta$ -heterodimers and their distinct functions dependent on the agonist/antagonist and promoter context investigated.

## Bibliography

1. G. G. Kuiper and J.Å. Gustafsson (1997) *FEBS Lett.* **410**, 87–90.
2. J. W. Schwabe, L. Chapman, J. T. Finch, and D. Rhodes (1993) *Cell* **75**, 567–578.
3. A. M. Brzozowski, A. C. W. Pike, Z. Dauter, R. E. Hubbard, T. Bonn et al. (1997) *Nature* **389**, 753–758.
4. T. D. Littlewood, D. C. Hancock, P. S. Danielian, M. G. Parker, and G. I. Evan (1995) *Nucleic Acids Res.* **23**, 1686–1690.
5. E. P. Smith, J. Boyd, G. R. Frank, H. Takahashi, R. M. Cohen et al. (1994) *N. Engl. J. Med.* **331**, 1056–1061.
6. D. B. Lubahn, J. S. Moyer, T. S. Golding, J. F. Couse, K. S. Korach et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 11162–11166.

## Suggestions for Further Reading

7. K. S. Korach (1994) *Science* **266**, 1524–1527.
8. M. J. Tsai and B. W. O'Malley (1994). *Annu. Rev. Biochem.* **63**, 451–486.
9. M. Beato, P. Herrlich, and G. Schütz (1995) Steroid hormone receptors: many actors in search of a plot, *Cell* **83**, 851–857.
10. H. Shibata, T. E. Spencer, S. A. Onate, G. Jenster, S. Y. Tsai et al. (1997) *Recent Prog. Horm. Res.* **52**, 141–164.

## Estrogens

Estrogens are female sex [steroid hormones](#) derived from cholesterol. The most potent estrogen is 3,17 $\beta$ -estradiol (estradiol 1,3,5-estratriene-3,17 $\beta$ -diol). The 16 $\alpha$ -hydroxylated derivative of estradiol, estriol, is a weak estrogen, whereas the 17 $\beta$ -oxidized compound, estrone, is hormonally inactive. Synthetic nonsteroid estrogens, such as diethylstilbestrol, have been introduced as therapeutic drugs in various endocrine diseases. Estrogen antagonists (eg, tamoxifen) have also found use in therapy, particularly of breast cancer.

Estradiol is synthesized in the ovary, the placenta, the adrenal cortex, and the testis. In the theca cells of the ovarian follicles, cholesterol is metabolized to progesterone and to androstendione, which in the granulosa cells is subsequently transformed to estradiol by aromatization of the A ring by the aromatase complex. This process is stimulated by FSH (folitropin).

During fetal and early postembryonic development, estrogens are responsible for inducing female characteristics and for the formation of the female genitals (oviduct, uterus, vagina). In the adult, estrogens act on the development of the secondary female sexual characteristics—for example, on the growth and differentiation of the breast glands, on the appearance of pubic and axillar hair, and on fat deposition in the skin. Estrogens affect the central nervous system and can modify sexual behavior. They also act positively on bone mass. In collaboration with progesterone, estradiol regulates the menstrual cycle. Specifically, in the first half of the cycle, estradiol stimulates proliferation of the epithelial and stromal cells of the endometrium, in preparation for ovulation.

The molecular action of estrogens, like that of the other steroid hormones, is by modulation of gene [transcription](#) (see [Response Element](#)). The hormone first binds to the [estrogen receptor](#), a member of the superfamily of nuclear receptors, which then homodimerizes and binds to specific nucleotide sequences, the estrogen [response elements](#). This triggers activation of gene transcription. Among the genes regulated by estrogens are **protooncogenes** of endometrium cells encoding [transcription factors](#) involved in cell replication (c-fos, c-jun, c-myc) (1-3), those encoding the progesterone receptor, [growth factors](#), and various structural and adhesion proteins (4). In the liver, estrogens stimulate the synthesis of steroid- and thyroxine-binding proteins. The molecular effects of estradiol on the synthesis of phosphoproteins in the chick oviduct, such as ovalbumin, have been studied in detail and were instrumental in developing concepts on the mode of action of steroid hormones (5).

### Bibliography

1. D. S. Loose-Mitchell, C. Chiappetta, and G. M. Stancel (1988) *Mol. Endocrinol.* **2**, 946–951.
2. C. Chiappetta et al. (1992) *J. Steroid Biochem. Mol. Biol.* **41**, 1134–123.
3. S. M. Hyder et al (1992) *J. Biol. Chem.* **267**, 18047–18054.
4. R. Grummer et al. (1994) *Biol. Reprod.* **51**, 1109–1116.
5. M. Kalimi et al. (1976) *J. Biol. Chem.* **251**, 516–523.

### Suggestions for Further Reading

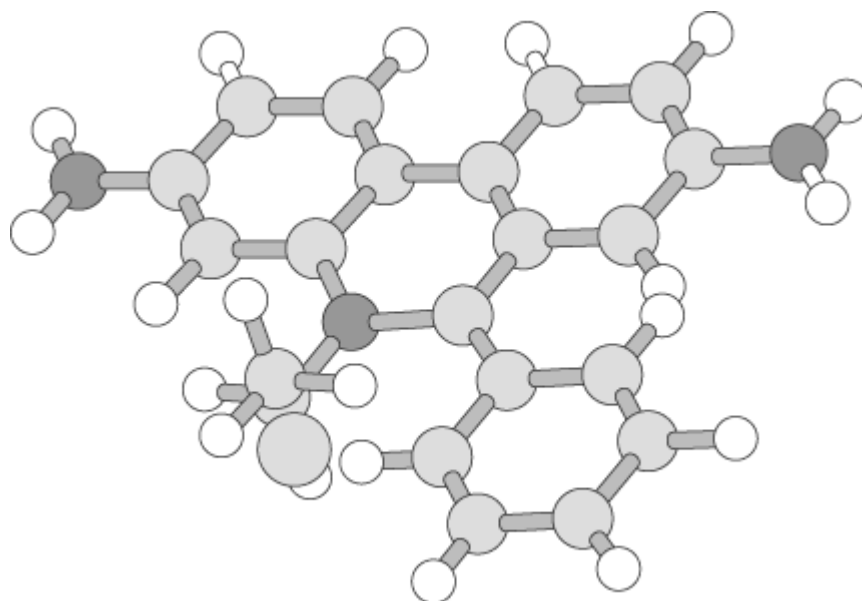
6. E. E. Baulieu and P. A. Kelly (eds) (1990) *Hormones: From Molecules to Disease*, Hermann, Paris.
7. E. Y. Adashi and P. C. K. Leung (1993) *The Ovary*, Raven Press, New York.

8. M. G. Parker (1993) *Steroid Hormone Action*, IRL Press, Oxford, U.K.

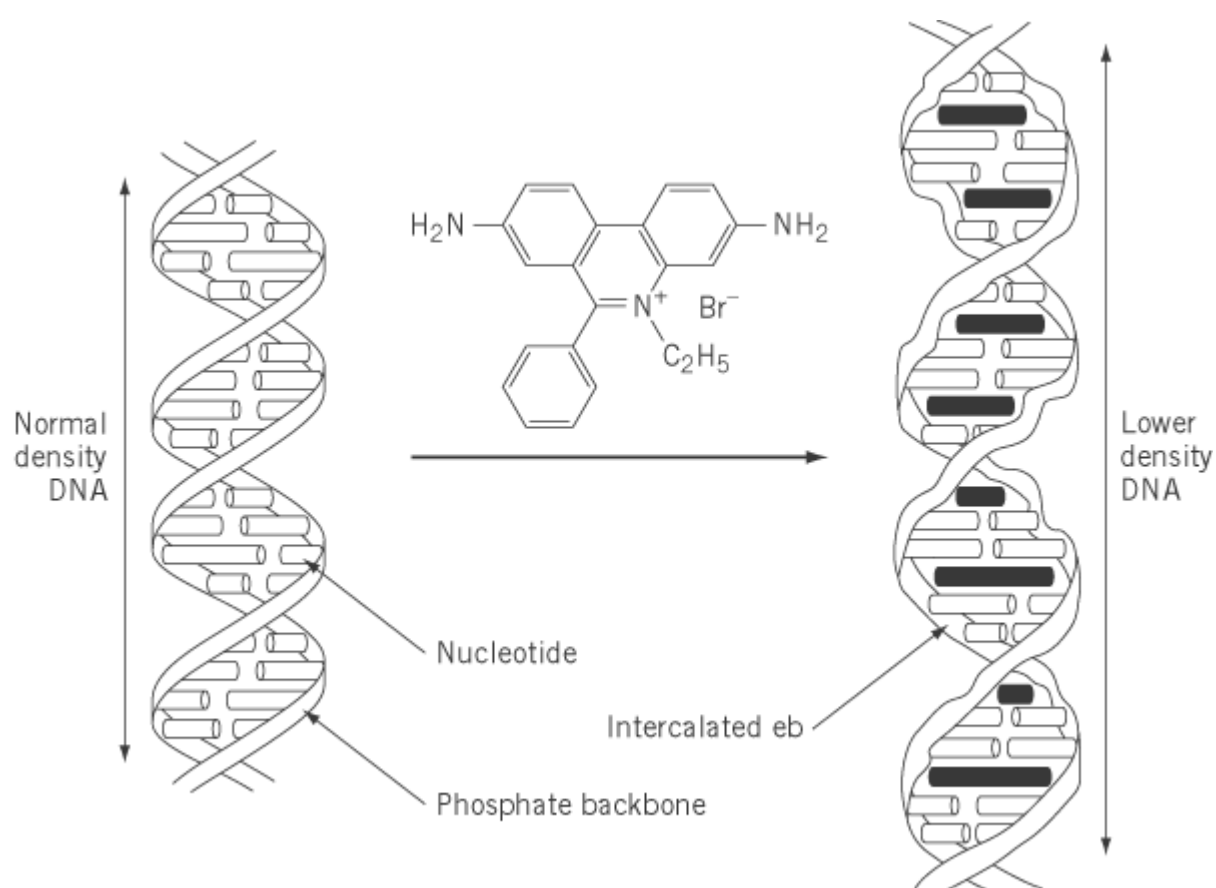
## Ethidium Bromide

Ethidium bromide (phenanthridinium, 3,8-diamino-5-ethyl-6-phenyl-, bromide) is a heteroaromatic cationic dye (an aminoacridine) with the molecular formula  $C_{21}H_{20}BrN_3$  and a relative molecular mass of 394.31 (Figs. 1, 2). It exists as a dark red crystalline or amorphous powder at standard temperatures and pressures. It has a melting point (and decomposition temperature) between 533 and 535 K, moderately high solubility in water (5 g/100 mL), and is essentially odorless. Ethidium bromide is a powerful [mutagen](#), extremely toxic by inhalation, ingestion ( $LD_{50}$  1503 mg/kg oral, rat), and skin contact, and a suspected [carcinogen](#) and reproductive toxin.

**Figure 1.** Ball-and-stick representation of the ethidium bromide molecule (phenanthridinium, 3,8-diamino-5-ethyl-6-phenyl-, bromide). RMM = 394.31.



**Figure 2.** The intercalation of ethidium bromide into a portion of a DNA double helix. The intercalated dye increases the spacing of successive base pairs, distorts the regular phosphate backbone, and reduces the pitch of the helix.



Originally, ethidium bromide (EB) found fame in the late 1940s as an antitrypanosomal, antimicrobial, antibacterial, and antiviral agent (1). It is generally agreed that these biological effects are a direct consequence of the inhibition of nucleic acid synthesis, which in turn is related to the specific binding of the drug to DNA. EB inhibits **DNA polymerase** and binds *in vitro* to both RNA and DNA. Investigation into the precise nature of the DNA-EB binding mechanism led to the discovery that the drug (and many other related molecules) binds by a mechanism termed [intercalation](#). This process has been studied extensively during the past three decades, and the photophysical changes that accompany intercalation have been successfully applied to quantitate and structurally elucidate DNA. The nature of the DNA-EB interaction and subsequent analytical applications are the focus of this article.

## 1. Photophysical Characteristics of EB

EB is dark yellow in aqueous solution and possesses a broad structureless absorption spectrum ( $\lambda_{\text{max}} \sim 480 \text{ nm}$ ). Increasing EB concentration leads to a red shift in the absorption maximum, accompanied by hypochromism. The existence of an isobestic point reflects a simple equilibrium between monomeric and dimeric species (2).

EB is weakly **fluorescent** in aqueous solution and possesses a broad structureless emission spectrum ( $\lambda_{\text{max}} \sim 617 \text{ nm}$ ). The low fluorescent quantum yield and short fluorescent lifetime (3) in aqueous solution are attributed to efficient quenching of the excited state by the transfer of a proton to an adjacent [water](#) molecule (4, 5). Moreover, the EB dimer is not fluorescent and consequently provides an efficient mechanism for fluorescent quenching when EB is complexed in the dimeric form. Both of these phenomena prove important when EB interacts with DNA.

## 2. Intercalation of Dyes with DNA

Modification of DNA by binding with dyestuffs suspected to have possible chemotherapeutic effects was noted as early as the 1940s in the pioneering work of Michaelis (6). The classic intercalation mechanism for the physical interaction between drugs and DNA is discussed elsewhere in this volume, but the characteristics of dye intercalation with DNA were established primarily with reference to aminoacridones (7) and EB (1). Consequently, the fundamental attributes of classic intercalation are outlined. Very simply, intercalation describes the insertion of planar heteroaromatic molecules between adjacent nucleotide base pairs of a DNA duplex (see [DNA Structure](#)). Intercalative binding follows the neighbor exclusion principle, where every second site along the helix remains unoccupied, and thus EB intercalation occurs through interaction of the phenantridinium ring with the duplex (8). This leads to a very elegant binding model (Fig. 2). More specifically, the EB–DNA complex involves **hydrogen bonding** between the 3,8 amino substituents of the EB molecule and the phosphate groups across the two DNA strands. These hydrogen bonds maintain the orientation of the EB cation with respect to the polymer frame (9). The simplicity of this model was naturally extended to explain the interaction of other molecules with DNA. These included antitumor antibiotics whose mode of action involves interference with nucleic acid synthesis (1). The intercalation mechanism of EB with double-stranded DNA can be thought of as consisting of three fundamental features: [1] helix extension (Fig. 2); [2] local, helical unwinding during binding; and [3] insertion of EB so that the plane of the aromatic chromophore is perpendicular to the helical axis. This binding process has been confirmed by various experimental approaches including, [X-ray crystallography](#) (10), [autoradiography](#) (11), [NMR spectroscopy](#) (12), equilibrium studies (13), **viscosity** measurements (14), and calorimetry (15). Nevertheless, optical measurements have played the dominant role in the use of the EB to DNA interaction. Consequently they are discussed in detail here.

### 3. Intercalation of EB with DNA

The formation of a metachromatic complex between EB and DNA can be simply observed by measuring the absorption spectrum. As the DNA to EB ratio increases, the absorption spectrum shifts to longer wavelengths (3, 16, 17). The formation of such complexes causes a shift due to the stabilization of the excited state upon binding (9). The shift is observed only at high DNA to EB ratios because an excess of free dye normally masks the absorbance contribution of bound EB. Importantly, a binding constant of  $2.1 \times 10^6 M^{-1}$  demonstrates the high affinity of EB of DNA (18).

More relevant to the use of EB in molecular biology are the observed variations in its fluorescent characteristics on binding to polynucleotides. Broadly, both the time-integrated fluorescent intensity and the average molecular fluorescent decaytime increase dramatically on interaction with DNA (eg, the fluorescent decay time of EB in water is about 1.8 ns, compared to 23 ns in DNA) (3, 8, 19, 20). This variation can be explained by reference to Figure 2. When EB intercalates with the double helix, it sits in the **hydrophobic** pocket of the base pair system and is shielded to a large extent from solvent molecules. Consequently, the rate of excited state quenching (via proton transfer to solvent molecules) decreases, leading to an increase in both the fluorescent decay time and the fluorescent quantum yield. Because the fluorescent intensity of EB in aqueous solution is very small (due to a low fluorescent quantum yield), measuring the fluorescent signal when bound to DNA provides the most common method for DNA detection in the molecular biology laboratory (e.g., in slab [gel electrophoresis](#) and [capillary zone electrophoresis](#)) (21-23). Furthermore, the fluorescent enhancement has resulted in extensive use of EB as a probe of the topological (24, 25) and dynamic (26-29) properties of DNA.

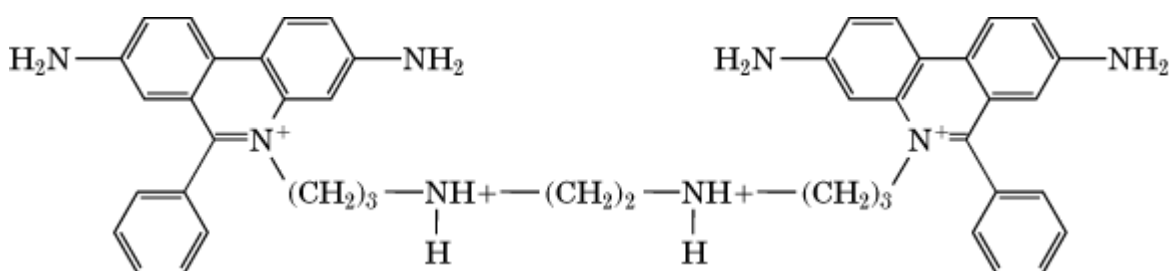
EB is widely used as a probe in clinical DNA assays. The selective binding of EB to double-stranded DNA, combined with the exquisite sensitivity of fluorescent spectroscopy, provides an obvious route to DNA quantitation on a nanogram scale (30). With excitation at 520 nm and emission at 600 nm, EB fluorescence is enhanced approximately 30-fold in upon binding to duplexes. Moreover, the use of fluorescence detection for quantitative purposes is favorable because fluorescent enhancement is

practically independent of base composition (31-33).

Nevertheless, care must be taken when using EB as a quantitative probe of DNA. Numerous studies have demonstrated a modification in the classic intercalative behavior of EB with high dye concentrations. In addition to conventional intercalation, there is also EB binds secondarily through an external site (most likely along the phosphate backbone). Binding to this secondary site occurs most readily at low ionic concentrations ( $<0.1 M$ ), after binding at primary sites has been saturated (3, 14, 19, 21). Secondary binding is most clearly evidenced by the heterogeneous nature of fluorescent decay profiles originating from high EB to DNA ratios (3, 19). The true nature of secondary binding is still poorly understood. It must be accounted for, however, when relating time-integrated fluorescent intensities to the amounts of double-stranded DNA. In addition, variations in the fluorescent quantum yield of EB intercalated into heterogeneous natural DNA (eg, calf thymus DNA) as a function of salt concentration and temperature suggest that the detailed features of the EB to DNA complex are inextricably linked to the overall structural properties of the binding sites on the duplex chain (9).

The huge interest in the clinical use of compounds that bind strongly to DNA as antitumor agents at low concentrations has led to the study of a related class of EB compounds. Appropriately designed dimers of EB have binding affinities and fluorescent enhancements much higher than those of the monomer (18, 34, 35). Glazer and co-workers synthesized an ethidium homodimer, illustrated in Figure 3, that bisintercalates at a ratio of one dimer per four to five base pairs and is stable on electrophoresis gels. The system was used successfully to detect and quantify DNA fragments with picogram sensitivity subsequent to an electrophoretic separation (36). Current applications of bisintercalators of this kind include multiplex detection of DNA restriction fragments (37), high sensitivity detection of the products from the polymerase chain reaction (PCR) (36), sizing of individual double-stranded DNA fragments by flow cytometry, and the study of DNA-protein interactions (38). Furthermore, DNA probes that have a double-stranded region (for intercalation sites) and a single-stranded region (for recognition of specific target sequences) offer an exciting application for this generation of versatile fluorescent probes.

Figure 3. Chemical structure of an ethidium homodimer (as synthesized in Ref. 36).



## Bibliography

1. M. J. Waring (1975) *Antibiotics III. Mechanism and Action of Antimicrobial and Antitumor Agents* (J. W. Corcoran and F. E. Hahn, eds.), Berlin, Heidelberg, pp. 141–165.
2. M. Guenza and C. Cuniberti (1988) *Spectrochim. Acta* **44A**(12), 1359–1364.
3. C. D. Byrne and A. J. de Mello (1998) *Biophys. Chem.* **70**, 173–184.
4. S. J. Atherton and P. C. Beaumont (1984) *Photobiochem. Photobiophys.* **8**, 103–113.
5. J. Olmsted and D. R. Kearns (1977) *Biochemistry* **16**, 3647–3654.
6. L. Michaelis (1947) *Cold Spring Harbor Symp. Quant. Biol.* **12**, 131–142.
7. R. M. Acheson (1973) *Acridines* (R. M. Acheson ed.), Wiley, London, p. 878.



8. J. B. LePecq (1976) *Biochemical Fluorescence: Concepts II* (R. Chen and H. Edelhoch eds.), Dekker, New York, pp. 711–736.
9. C. Cuniberti and M. Guenza (1990) *Biophys. Chem.* **38**, 11–22.
10. C. C. Tsai, S. C. Jain, and H. M. Sobell (1977) *J. Mol. Biol.* **114**, 301–315.
11. J. Cairns (1962) *Cold Spring Harbor Symp. Quant. Biol.* **27**, 311–318.
12. S. Chandrasakaran, R. L. Jones, and W. D. Wilson (1985) *Biopolymers* **24**, 1963–1979.
13. W. D. Wilson, C. R. Krishnamoorthy, Y. Wang, and J. C. Smith (1985) *Biopolymers* **24**, 1941–1961.
14. J. B. Le Pecq and C. Paoletti (1967) *J. Mol. Biol.* **27**, 87–106.
15. F. Quadrifoglio, V. Crescenzi, and V. Giancotti (1974) *Biophys. Chem.* **2**, 319–324.
16. M. J. Waring (1968) *Nature* **219**, 1320–1322.
17. Y. S. Babayan, G. Manzini, and F. Quadrifoglio (1988) *Mol. Biol.* **22**(4), 714–725.
18. B. Gaugain et al. (1978) *Biochemistry* **17**, 5078–5088.
19. D. P. Heller and C. L. Greenstock (1994) *Biophys. Chem.* **50**, 305–312.
20. V. W. F. Burns (1969) *Arch. Biochem. Biophys.* **133**, 420–424.
21. J. Yguerabide and A. Ceballos (1995) *Anal. Biochem.* **228**, 208–220.
22. B. Olszanska and A. Borgul (1990) *Acta Biochim. Pol.* **37**(1), 59–63.
23. S. Premaratne, S. D. Swenson, M. Mandel, and H. F. Mower (1994) *Biochim. Biophys. Acta* **1219**, 422–424.
24. T. Ide and R. Basuga (1976) *Biochemistry* **15**, 600.
25. S. A. Winkle, L. S. Rosenberg, and T. R. Krug (1982) *Nucleic Acids Res.* **10**, 8211–8223.
26. D. P. Millar, R. J. Robbins, and A. H. Zewail (1982) *J. Chem. Phys.* **76**, 2080–2094.
27. D. Magde, M. Zappala, W. J. Know, and T. M. Nordlund (1983) *J. Chem. Phys.* **87**, 3286–3288.
28. J. C. Thomas and J. M. Schurr (1983) *Biochemistry* **22**, 6194–6198.
29. T. Hard and D. R. Kearns (1986) *J. Phys. Chem.* **90**, 3437–3444.
30. A. R. Morgan et al. (1979) *Nucleic Acids Res.* **7**, 547–570.
31. M. J. Waring (1965) *J. Mol. Biol.* **13**, 269–282.
32. E. F. Gale et al. (1981). *The Molecular Basis of Antibiotic Action*, Wiley, London.
33. F. M. Pohl, T. M. Jovin, W. Baehr, and J. J. Holbrook (1972) *Proc. Natl. Acad. Sci. USA* **69**, 3805–3809.
34. B. Gaugain et al. (1978) *Biochemistry* **17**, 5071–5078.
35. J. Markovits, B. P. Roques, and J. B. L. Pecq (1979) *Anal. Biochem.* **94**, 259–264.
36. A. N. Glazer, K. Peck, and R. A. Mathies (1990) *Proc. Natl. Acad. Sci. USA* **87**, 3851–3855.
37. H. S. Rye et al. (1992) *Nucleic Acids Res.* **20**, 2803–2812.
38. H. S. Rye, B. L. Drees, H. M. N. Nelson, and A. N. Glazer (1993) *J. Biol. Chem.* **268**, 25229–25238.

### Suggestion for Further Reading

39. E. F. Gale et al. (1981). *The Molecular Basis of Antibiotic Action*, Wiley, London. This general reference book provides a complete overview of the action and use of most classes of intercalating dyes.

## Ethylene

## 1. History

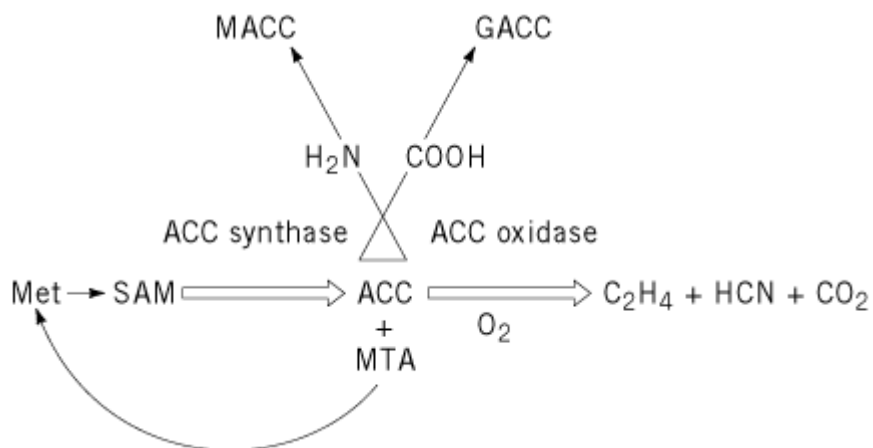
For a long time, people have recognized the ability of certain gases to stimulate fruit ripening. The ancient Chinese used to burn incense to accelerate ripening. In 1910, Cousins (1) advised the Jamaican Agricultural Department to ship oranges and bananas separately, because some emanation from the oranges (originating from fungal contamination) induced ripening of the bananas. In 1934, Gane (2) proved that ethylene is synthesized by plants and initiates fruit ripening.

The first report of effects of ethylene on vegetative growth dates from 1901, when the Russian physiologist Dimitri Neljubow (3) demonstrated that ethylene in illuminating gas causes the “triple response” in dark-grown pea seedlings: inhibited epicotyl elongation, radial swelling, and a horizontal growth habit, because of loss of geotropism (3). Deleterious effects of illuminating gas had already been observed in the previous century, when leaks in gas pipes caused precocious leaf abscission of surrounding trees.

## 2. Biosynthesis and Metabolism

Ethylene is a unique gaseous signaling molecule with a simple molecular structure ( $C_2H_4$ ). It is active at concentrations as low as 10 nL/L of air (10 ppb). The biological precursor of ethylene in higher plants is L-methionine. By the action of *S*-adenosyl-L-methionine (AdoMet) synthetase, methionine and ATP are conjugated to form AdoMet (Fig. 1). Subsequently, AdoMet is converted to 1-aminocyclopropane-1-carboxylic acid (ACC) and 5'-methylthio-adenosine (MTA) by the enzyme ACC synthase. MTA is recycled to methionine. ACC is converted to ethylene, HCN, and carbon dioxide by ACC oxidase. Alternatively, ACC can be conjugated to malonyl-ACC and *g*-L-glutamyl-ACC (4, 5). These inactive derivatives are formed to prevent precursor accumulation and conversion to ethylene. In certain cases, however, malonyl-ACC can also serve as a source of ACC (6).

**Figure 1.** The ethylene biosynthesis pathway in higher plants. ACC, 1-aminocyclopropane-1-carboxylic acid; GACC, glutamyl-ACC; MACC, malonyl-ACC; MTA, 5'-methylthio-adenosine; SAM, *S*-adenosyl-L-methionine.



Control of ethylene biosynthesis has been observed at the level of ACC synthase and ACC oxidase activities and includes an important aspect of autoregulation (7, 8). In most cases, the ethylene levels are directly correlated with the ACC synthase activity, which is known to be highly regulated during plant development and in response to a complex array of environmental stimuli (9). Data have accumulated, however, that also indicate the importance of ACC oxidase for the regulation of

ethylene synthesis (7). The synthesis of ACC does not always translate immediately into ethylene production. Interestingly, ACC is a soluble and translocatable hormone precursor, in contrast to ethylene, which, as a gas, is most probably not suited for targeted translocation processes. Translocation of ACC as a mechanism of interorgan signaling has been proposed for the regulation of plant responses to flooding stress (10, 11).

Because of its low abundance and high lability, ACC synthase proved very difficult to purify. Cloning of ACC synthase genes was achieved independently in several laboratories (12-14). The first nucleotide sequences of ACC synthase cDNA clones were reported by Van Der Straeten et al. (13) and Nakajima et al. (14). The former group obtained peptide sequences after 5000-fold purification of wound-induced pericarp tissue of ripening tomatoes, followed by oligonucleotide-directed screening of a cDNA library (13). Nakajima et al. (14) used an immunological approach to clone ACC synthase from a winter squash cDNA library. Sato and Theologis (12) isolated an ACC synthase clone by screening a cDNA expression library from zucchini with an antibody generated against purified ACC synthase, and the sequence was reported by Sato et al. (15). To date, clones encoding ACC synthase have been isolated from a large array of plant species. The enzyme is encoded by a highly divergent multigene family, and the expression of each gene is differentially regulated by a different subset of inducing or repressing conditions (16). *Arabidopsis thaliana* was shown to have at least five ACC synthase genes (17, 18). The *AT-ACSI* gene expression pattern was analyzed in detail using promoter- $\beta$ -glucuronidase reporter gene experiments (19). It was shown that *AT-ACSI* expression was high in young and immature leaf tissues, lower in stems and in senescing tissues. The transcript levels of most ACC synthase genes are inducible by cycloheximide (16, 17), which suggests either that the messenger RNAs are labile or that transcription of ACC synthase genes is under the control of a labile repressor molecule (17). In either case, it can be concluded that ACC synthase activity is regulated predominantly at the transcriptional level. In addition, the regulation of some ACC synthase genes by various inducers may involve protein kinase activity, which in turn may be controlled through inositol-1,4,5-triphosphate (IP3)-mediated  $Ca^{2+}$  mobilization (20) (see Calcium Signaling).

A better understanding of the mechanisms involved in the transcriptional regulation of ACC synthase genes will create the possibility to control ethylene formation in a highly specific manner. To date, two strategies have been successfully applied to lower the level of ACC available for ethylene synthesis. In one experiment, the overexpression of the ACC deaminase gene of *Pseudomonas* sp. strain 6G5 in tomato plants resulted in a reduction of ethylene synthesis (21). Another approach was followed by Oeller et al. (22): expression of a tomato ripening-specific ACC synthase gene in an antisense orientation repressed ethylene synthesis almost completely. These antisense tomato fruits never get ripe until they are treated with ethylene or propylene.

Numerous attempts were made to purify the enzyme downstream of ACC synthase, namely ACC oxidase, or ethylene-forming enzyme. A clone encoding this enzyme was finally isolated in a reverse genetics approach (23). As a first step, ripening-related clones from a tomato cDNA library were isolated. For one of these clones, induction of the corresponding mRNA was tightly correlated with increased ethylene synthesis (24). Transgenic tomato plants expressing the cDNA in an antisense orientation were constructed, resulting in a reduction of ethylene synthesis in a gene dosage-dependent manner. In addition, these plants were delayed in fruit ripening and leaf senescence, and ethylene synthesis was also substantially reduced in wounded leaves (23, 25). Finally, the cDNA was proven to encode ACC oxidase by heterologous expression in yeast and in *Xenopus laevis* oocytes (26, 27).

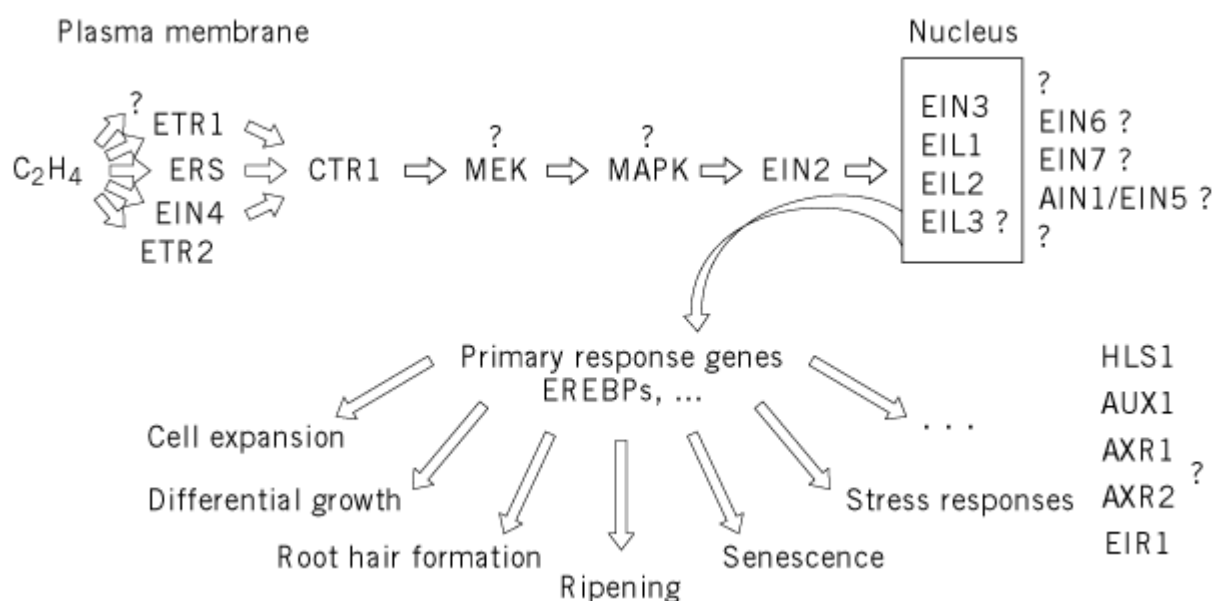
With the exception of preclimacteric fruits and flowers, it was generally believed that ACC oxidase is present constitutively in plant tissues and, therefore, would not contribute to the regulation of ethylene synthesis (28). However, Grierson et al. (7) have pointed out that the rapid induction of ACC oxidase transcripts upon wounding, when assayed in dissected plant organs, may have contributed to the false belief that this enzyme is present constitutively in all tissues. To date, we

know that the expression of ACC oxidase genes can be highly regulated throughout plant development (29). As for the expression of ACC synthase genes, the transcriptional regulation of ACC oxidases was also found to be controlled by **phosphorylation** and dephosphorylation events (30). Whereas ACC synthase expression was both **feedback-inhibited** and stimulated by ethylene (8), ACC oxidase was rapidly induced by ethylene treatments (31). Finally, it should be mentioned that ethylene can be metabolized to ethylene oxide (32).

### 3. Signal Perception and Transduction

Using the triple response for isolating ethylene mutants, and subsequent double-mutant analysis, the ethylene signal was found to be perceived through several receptors (ETR1, ERS, and maybe also EIN4 and ETR2) and was transduced via a linear pathway to the nucleus, where the expression of a set of primary ethylene-response genes is regulated (33, 34) (Fig. 2). Characterization of the *ETR1* gene suggested an explanation for the fact that different ethylene concentrations, which give different responses, can act through a common **signal transduction** pathway. The *ETR1* gene has significant homology with bacterial two-component regulators (35). This kind of receptor systems is known to be able to operate over a wide range of ligand concentrations, resulting in an output that is graded with stimulus intensity (36). The amino-terminal region of the protein contains **hydrophobic** stretches, each of which is predicted to span the membrane. ETR1 is a membrane-associated, **disulfide bond**-linked dimer in extracts of *Arabidopsis* and when expressed in yeast (37). Ethylene was shown to bind to the amino-terminal hydrophobic region (38), and all of the known *etr1* mutations were also located in this domain. In addition, the *etr1-1* mutation was shown to disrupt ethylene binding. These data are consistent with the chemical properties of ethylene, which is more soluble in lipids than in water and therefore will bind preferentially in a membrane environment (39). The dominant nature of *etr1* mutations, together with the redundancy of receptors, suggests that ETR1 inhibits the ethylene response in the absence of ethylene. This characteristic implies that *etr1* mutations are gain-of-function mutations, where the enzyme is locked in a catalytically active state. When ethylene binds to the receptor, this inhibitory activity is expected to be disrupted, allowing the ethylene response to occur. Accordingly, mutations preventing the binding of ethylene would then cause a constitutive inhibition of all ethylene responses.

**Figure 2.** Order of action of components of the ethylene response pathway in *Arabidopsis*. The pathway was established by double-mutant analysis. EIL3 is homologous to EIN3, EIL1, and EIL2; however, it has not yet been shown to be functional in ethylene signal transduction. The position of AIN1/EIN5, EIN6, and EIN7 in the pathway is unclear. AUX1, AXR1, AXR2, EIR1, and HLS1 control a subset of ethylene responses and may act in parallel to the response pathway defined by the ETR1, CTR1, and EIN2 genes.



Recently, another putative *Arabidopsis* ethylene receptor, ERS (*ethylene response sensor*), was cloned (40), supporting the finding that the *Arabidopsis* genome encodes a complex family of these signal transduction proteins (35). The *EIN4* (*ethylene-insensitive 4*) gene has not yet been isolated. However, as for *etr1* mutations, *ein4* was dominant, and *ein4 ctr1* double mutants had a **Ctrl phenotype**, indicating that EIN4 also acts upstream of CTR1 (41). Therefore, it is possible that *EIN4* belongs to the same gene family as *ETR1* and *ERS*. It should also be noted that biochemical evidence supports the existence of several ethylene-binding sites that can be classified in two types according to their rate constants for binding (42).

Downstream, the signal is transduced via CTR1, which encodes a Raf-like protein kinase (43). This observation confirms that ethylene induces a transient phosphorylation of several proteins in tobacco (44). All of the *ctr1* mutations found are predicted to disrupt the kinase activity (43). Furthermore, the homology of CTR1 to Raf-1 suggested that a MAP kinase cascade is located downstream (45). A similar combination of a “bacterial” histidine kinase and a Raf-like kinase was found for the high osmolarity glycerol pathway in yeast (46, 47). Signal transduction pathways that include Raf normally contain a MEK protein kinase, which is activated by Raf through phosphorylation (48). Subsequently, activated MEK phosphorylates a MAP kinase that, in turn, phosphorylates several downstream components. Such a MAP kinase was identified for the yeast high osmolarity glycerol pathway (49) and is also expected in the ethylene response pathway (45). To date, ethylene mutants identifying components of a putative MAP kinase cascade have not been isolated; however, screening for triple-response mutants has not yet saturated the genome (34). Given that *ctr1-1* mutations are recessive, it is expected that CTR1 is an inhibitor of the ethylene response in the absence of ethylene. If this were the case, CTR1 may be directly activated by ETR1, which in turn would be active in the absence of ethylene (34).

After CTR1, the ethylene response pathway proceeds with EIN2, followed by EIN3, EIL1, and EIL2. The order of action was based on the fact that *ein2* mutations tend to result in strong insensitivity, in contrast to *ein3* mutants, which remain more responsive to ethylene even in the case of null alleles (33, 41). This characteristic was confirmed by the cloning of *EIN3* and the related genes *EIL1*, *EIL2*, and *EIL3*. Overexpression of *EIN3* and *EIL1* in an *ein2* background resulted in a constitutive ethylene response phenotype (33), which can only be explained when EIN3 and EIL1 act downstream of EIN2, so that mutations in EIN2 do not inhibit the overexpression phenotype. *EIN3*, *EIL1*, and *EIL2* encode a novel class of nuclear proteins; and they seem to be regulating the same downstream components, because overexpression in each case complemented *ein3*. These proteins may be directly regulating gene transcription or may interact with transcription factors that have ethylene-regulated genes as targets (33).

A number of weakly insensitive mutants identified genes *AIN1/EIN5*, *EIN6*, and *EIN7*, which act downstream of CTR1, but their position relative to *EIN2*, *EIN3*, *EIL1*, and *EIL2* could not be determined (41, 50). In addition, certain mutations affect a specific part of the ethylene response. For example, the *HLS1* (*HOOKLESS1*) gene controls ethylene-regulated differential growth in the apical hook of etiolated seedlings (51), and the *EIR1* gene specifically regulates the ethylene sensitivity of roots (41).

The genetic analysis of ethylene signal transduction in *Arabidopsis* has been very successful and has resulted in the isolation of several components of the pathway. This research will now have to be complemented by biochemical studies and by detailed analysis of the various protein activities in relation to the existing physiological data on ethylene signaling. Additionally, progress is also expected from the use of new screening procedures to isolate novel ethylene mutants (52, 53). The analysis of light-grown populations of mutant seedlings can be complementary to the triple response screening, especially when considering the close interaction between light and ethylene during plant

development (39).

#### 4. Downstream Targets

Downstream from EIN3 and EIL1-3, the signal transduction pathway is expected to diverge into several branches that lead to specific ethylene responses. This view is corroborated by the existence of different **cis-acting** elements that confer ethylene inducibility in response to different inducing conditions. At least two separate mechanisms exist, one of which involves a *cis* element designated GCC box. In defense-related genes, such as chitinase (54), b-1,3-glucanase (55), and the basic pathogenesis-related (PR1) protein, a common motif was found, with the **consensus sequence** GCCGCC conserved between the *Nicotiana*, *Arabidopsis*, and *Phaseolus* GCC boxes (56). In addition, a GCC box was found in the **promoter** of the *HSL1* gene (51). Strong evidence for a key role of the GCC box in ethylene-induced transcription was given by coupling a 47-bp fragment from a tobacco b-1,3-glucanase gene containing two GCC box elements to the minimal **cauliflower mosaic virus** 35 S promoter and by demonstrating that it conveys ethylene responsiveness in transgenic tobacco plants (57). Four different cDNAs encoding GCC box-binding proteins were isolated by southwestern screening of a tobacco cDNA expression library with labeled GCC fragment as a probe (57). The predicted *trans* factors were named ethylene-**response element**-binding proteins (EREBPs) and defined a novel class of plant-specific transcriptional regulators. The DNA-binding domain was identified in a region of 59 amino acid residues that was fully conserved in the four EREBPs. Expression of the corresponding genes was ethylene-inducible. Several factors that play a role in development turned out to be members of the EREBP class of transcription factors, including APETALA2 (58), AINTEGUMENTA (59), and TINY (60). More recently, a DNA-binding protein involved in low-temperature expression was shown to be an additional member (61).

A mechanism independent of EREBPs appears to act during senescence of flowers and fruit ripening (62). A detailed analysis was made of the ethylene-regulated transcription of the **glutathione-S-transferase** gene 1 (*GST1*) during carnation petal senescence (63). A cDNA encoding a DNA-binding protein that interacts with a region of *GST1* and protects it upon **footprinting** using deoxyribonuclease 1 was obtained by southwestern screening (64). In tomato ripening, cooperative *cis* elements are required for ethylene regulation of *E4* gene transcription (65).

Finally, it should be noted that all the genes mentioned above are secondary response genes. Primary ethylene-responsive genes, activated within minutes after ethylene induction, remain to be identified.

#### 5. Effects

Ethylene has been shown to participate in a diverse array of plant developmental processes, including seed germination, cell expansion, root initiation, senescence, leaf abscission, and fruit ripening in climacteric fruits (39). In addition, it controls sex determination in some monoecious species. In mango and Bromeliads, it induces flowering. Finally, ethylene is implicated in the control of biotic as well as abiotic **stress responses** (39).

#### Bibliography

1. H. H. Cousins (1910) III. *Agricultural Experiments* Citrus. Annual Report of the Department of Agriculture, Jamaica, 7.
2. R. Gane (1934) *Nature* **134**, 1008.
3. D. Neljubow (1901) *Beih. Bot. Zentralbl.* **10**, 128–139.
4. N. Amrhein, D. Schneebeck, H. Skorupka, and S. Tophof (1981) *Naturwissenschaften* **68**, 619–620.
5. M. N. Martin, J. D. Cohen, and R. A. Saftner (1995) *Plant Physiol.* **109**, 917–926.
6. X.-Z. Jiao, S. Philosoph-Hadas, L.-Y. Su, and S. F. Yang (1986) *Plant Physiol.* **81**, 637–641.

7. D. Grierson, A. J. Hamilton, M. Bouzayen, M. Köck, G. W. Lycett, and S. Barton (1992) In *Inducible Plant Proteins* (J. L. Wray, ed.), Cambridge University Press, Cambridge, U.K., pp. 155–174.
8. R. Fluhr and A. K. Mattoo (1996) *Crit. Rev. Plant Sci.* **15**, 479–523.
9. H. Kende (1989) *Plant Physiol.* **91**, 1–4.
10. K. J. Bradford and S. F. Yang (1980) *Plant Physiol.* **65**, 322–326.
11. M. B. Jackson (1985) *Annu. Rev. Plant Physiol.* **36**, 145–174.
12. T. Sato and A. Theologis (1989) *Proc. Natl. Acad. Sci. USA* **86**, 6621–6625.
13. D. Van Der Straeten, L. Van Wiemeersch, H. M. Goodman, and M. Van Montagu (1990) *Proc. Natl. Acad. Sci. USA* **87**, 4859–4863.
14. N. Nakajima, H. Mori, K. Yamazaki, and H. Imaseki (1990) *Plant Cell Physiol.* **31**, 1021–1029.
15. T. Sato, P. W. Oeller, and A. Theologis (1991) *J. Biol. Chem.* **266**, 3752–3759.
16. T. I. Zarembinski and A. Theologis (1994) *Plant Mol. Biol.* **26**, 1579–1597.
17. X. Liang, S. Abel, J. A. Keller, N. F. Shen, and A. Theologis (1992) *Proc. Natl. Acad. Sci. USA* **89**, 11046–11050.
18. D. Van Der Straeten, R. A. Rodrigues-Pousada, R. Villarroel, S. Hanley, H. M. Goodman, and M. Van Montagu (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9969–9973.
19. R. A. Rodrigues-Pousada, R. De Rycke, A. Dedonder, W. Van Caeneghem, G. Engler, M. Van Montagu, and D. Van Der Straeten (1993) *Plant Cell* **5**, 897–911.
20. X. Liang, N. F. Shen, and A. Theologis (1996) *Plant J.* **10**, 1027–1036.
21. H. J. Klee, M. B. Hayford, K. A. Kretzmer, G. F. Barry, and G. M. Kishore (1991) *Plant Cell* **3**, 1187–1193.
22. P. W. Oeller, L. Min-Wong, L. P. Taylor, D. A. Pike, and A. Theologis (1991) *Science* **254**, 437–439.
23. A. J. Hamilton, G. W. Lycett, and D. Grierson (1990) *Nature* **346**, 284–287.
24. C. J. S. Smith, A. Slater, and D. Grierson (1986) *Planta* **168**, 94–100.
25. S. Picton, S. L. Barton, M. Bouzayen, A. J. Hamilton, and D. Grierson (1993) *Plant J.* **3**, 469–481.
26. A. J. Hamilton, M. Bouzayen, and D. Grierson (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7434–7437.
27. P. Spanu, D. Reinhardt, and T. Boller (1991) *EMBO J.* **10**, 2007–2013.
28. S. F. Yang and N. E. Hoffman (1984) *Annu. Rev. Plant Physiol.* **35**, 155–189.
29. C. S. Barry, B. Blume, M. Bouzayen, W. Cooper, A. J. Hamilton, and D. Grierson (1996) *Plant J.* **9**, 525–535.
30. J. H. Kim, W. T. Kim, B. G. Kang, and S. F. Yang (1997) *Plant J.* **11**, 399–405.
31. W. T. Kim and S. F. Yang (1994) *Planta* **194**, 223–229.
32. I. O. Sanders, A. R. Smith, and M. A. Hall (1986) *Physiol. Plant.* **66**, 723–726.
33. Q. Chao, M. Rothenberg, R. Solano, G. Roman, W. Terzaghi, and J. R. Ecker (1997) *Cell* **89**, 1133–1144.
34. J. J. Kieber (1997) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **48**, 277–296.
35. C. Chang, S. F. Kwok, A. B. Bleecker, and E. M. Meyerowitz (1993) *Science* **262**, 539–544.
36. J. S. Parkinson (1993) *Cell* **73**, 857–871.
37. G. E. Schaller, A. N. Ladd, M. B. Lanahan, J. M. Spanbauer, and A. B. Bleecker (1995) *J. Biol. Chem.* **270**, 12526–12530.
38. G. E. Schaller and A. B. Bleecker (1995) *Science* **270**, 1809–1811.
39. F. B. Abeles, P. W. Morgan, and M. E. Saltveit (1992) *Ethylene in Plant Biology*, 2nd ed., Academic Press, San Diego, CA.

40. J. Hua, C. Chang, Q. Sun, and E. M. Meyerowitz (1995) *Science* **269**, 1712–1714.
41. G. Roman, B. Lubarsky, J. J. Kieber, M. Rothenberg, and J. R. Ecker (1995) *Genetics* **139**, 1393–1409.
42. N. V. J. Harpham, A. W. Berry, M. G. Holland, I. E. Moshkov, A. R. Smith, and M. A. Hall (1996) *Plant Growth Regul.* **18**, 71–77.
43. J. J. Kieber, M. Rothenberg, G. Roman, K. A. Feldmann, and J. R. Ecker (1993) *Cell* **72**, 427–441.
44. V. Raz and R. Fluhr (1993) *Plant Cell* **5**, 523–530.
45. J. R. Ecker (1995) *Science* **268**, 667–675.
46. I. M. Ota and A. Varshavsky (1993) *Science* **262**, 566–569.
47. T. Maeda, S. M. Wurgler-Murphy, and H. Saito (1994) *Nature* **369**, 242–245.
48. D. R. Alessi, Y. Saito, D. G. Campbell, P. Cohen, G. Sithanandam, U. Rapp, A. Ashworth, C. J. Marshall, and S. Cowley (1994) *EMBO J.* **13**, 1610–1619.
49. F. Posas, S. M. Wurgler-Murphy, T. Maeda, E. A. Witten, T. C. Thai, and H. Saito (1996) *Cell* **86**, 865–875.
50. D. Van Der Straeten, A. Djudzman, W. Van Caeneghem, J. Smalle, and M. Van Montagu (1993) *Plant Physiol.* **102**, 401–408.
51. A. Lehman, R. Black, and J. R. Ecker (1996) *Cell* **85**, 183–194.
52. J. Smalle and D. Van Der Straeten (1997) *Physiol. Plant.* **100**, 593–605.
53. J. Smalle, M. Haegman, J. Kurepa, M. Van Montagu, and D. Van Der Straeten (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2756–2761.
54. K. E. Broglie, P. Biddle, R. Cressman, and R. Broglie (1989) *Plant Cell* **1**, 599–607.
55. G. Felix and F. Meins (1987) *Planta* **172**, 386–392.
56. Y. Eyal, Y. Meller, S. Lev-Yadun, and R. Fluhr (1993) *Plant J.* **4**, 225–234.
57. M. Ohme-Takagi and H. Shinshi (1995) *Plant Cell* **7**, 173–182.
58. D. Weigel (1995) *Plant Cell* **4**, 388–389.
59. K. M. Klucher, H. Chow, L. Reiser, and R. L. Fischer (1996) *Plant Cell* **8**, 137–153.
60. K. Wilson, D. Long, J. Swinburne, and G. Coupland (1996) *Plant Cell* **8**, 659–671.
61. E. J. Stockinger, S. J. Gilmour, and M. F. Thomashow (1997) *Proc. Natl. Acad. Sci. USA* **94**, 1035–1040.
62. J. Deikman (1997) *Physiol. Plant.* **100**, 561–566.
63. H. Itzhaki, J. M. Maxson, and W. R. Woodson (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8925–8929.
64. J. M. Maxon and W. R. Woodson (1996) *Plant Mol. Biol.* **31**, 751–759.
65. R. Xu, S. Goldman, S. Coupe, and J. Deikman (1996) *Plant Mol. Biol.* **31**, 1117–1127.

### **Ethyl Methane Sulfonate (EMS)**

EMS is a [mutagen](#) and a monofunctional alkylating agent of structure  $\text{CH}_3\text{-CH}_2\text{-O-SO}_2\text{-CH}_3$ , which reacts with nucleophilic sites on DNA by both  $\text{S}_{\text{N}}1$  and  $\text{S}_{\text{N}}2$  mechanisms. The  $\text{S}_{\text{N}}1$  process involves slow generation of a very reactive intermediate that alkylates DNA with only some discrimination for nitrogen atoms. Thus treatment of DNA or whole cells with EMS leads to many



DNA products. The most significant for mutagenesis is probably  $O_6$ -ethylguanine, which base-pairs readily with thymine (1). Using [bacteriophage](#) systems, Loveless (2) showed that EMS is capable of producing high mutation frequencies, with little cell killing. The mutations produced were primarily G.C → A.T transitions, and to a lesser extent transitions containing A.T base-pairs at the mutated site, as well as [frameshift mutations](#) (3, 4). These latter events may relate to effects on DNA repair processes, while the initial events may be due to the base-pair miscoding potential of the  $O_6$ -alkylguanine residues (1). EMS has often been considered the mutagen of choice for induced mutagenesis studies, being potent, easy to use, and having a well-documented mutational specificity.

### Bibliography

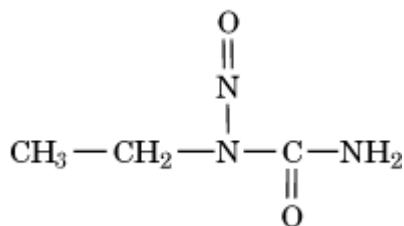
1. P. F. Swann (1990) *Mutat. Res.* **233**, 81–94.
2. A. Loveless (1958) *Nature* **181**, 1212–1213.
3. E. Bautz and E. Freese (1960) *Proc. Natl. Acad. Sci. USA* **46**, 1585–1594.
4. E. Freese (1961) *Proc. Natl. Acad. Sci. USA* **47**, 540–545.

### Ethyl-Nitrosourea

1-Ethyl-1-nitrosourea (ENU) is an ethylating and carbamoylating agent that is used as a potent [mutagen](#) in the [mouse](#). ENU is widely used in mutagenesis of male (stem cell spermatogonia; differentiating spermatogonia; post-spermatogonial stages) and female (oocytes) mouse **germ cells**, where it introduces intragenic mutations. However, it has also been used for [mutagenesis](#) in [sperm](#) and embryonic and primordial germ cells, as well as **zygotes**. Because of its germ cell potency, ENU is the most suitable chemical mutagen for the generation of desired new mutations in the mouse. A variety of ENU-induced mouse mutants have been used as animal models for human genetic diseases. To close the **phenotype** gap in the mouse, several ENU mutagenesis programs have been initiated. The phenotype gap describes the fact that, although a great number of mouse mutants exist, there is still little known about the specific phenotypes of these mice.

ENU is a yellow-pink crystalline substance, with a melting point of 103° to 104°C and a molecular mass of 117.10. The structural formula is shown in Figure 1. ENU is not stable under physiological conditions, having a **half-life** of 34.1 min at pH 7.0 and 37°C. Because ENU is sensitive to humidity and light, it is stored below –10°C.

**Figure 1.** Structural formula of 1-ethyl-1-nitrosourea (ENU).



ENU produces mainly GC/AT [transition mutations](#) and, to a smaller extent, AT/GC, AT/CG,

AT/TA, GC/CG, and GC/TA base substitutions. These effects are induced by ethylation of different positions within the four bases adenine, cytosine, guanine, and thymine, plus the phosphate groups of the DNA backbone. Additionally, ENU alkylates [transfer RNA](#) by the formation of 1,7-diethylguanosine. On removal of the ethyl group from O6-ethylguanine by the O6-alkylguanine-DNA alkyltransferase (AGT), S-ethylcysteine is formed. ENU also reacts with chromosomal proteins ([histones](#)) to an even greater extent than other mutagens, eg, 1-methyl-1-nitrosourea or 1-*n*-propyl-nitrosourea. Proteins are carbamoylated by the degradation product of ENU, isocyanic acid. It is likely that carbamoylation of the nuclear proteins influences DNA–histone interactions.

In mouse mutagenesis studies, ENU is the preferred mutagen in spermatogonia, with a linear relationship between the induction of mutations and doses greater than 100 mg ENU/kg. At concentrations lower than 100 mg/kg, mutation induction is nonlinear, as a result of [DNA repair](#) processes in stem cell spermatogonia. The major mutagenic lesion is the formation of O6-ethylguanine. However, analysis of mouse hemoglobin mutants Hbb and Hba induced by ENU has shown that base substitutions at A and T can also occur. This contrasts with findings in other systems and emphasizes the complexity of *in vivo* mammalian germ cell mutagenesis. ENU is considered to be a supermutagen in mouse stem cell spermatogonia, inducing recessive mutations with a very high frequency. This mutation rate is maximized by repeated intraperitoneal injections of at least 100-mg ENU/kg at weekly intervals to a total dose of 300- to 400-mg ENU/kg.

Chemical mutagens have also been used to induce mutations in arrested and maturing **oocytes** of female mice. However, the mutation rate is significantly lower than for treated male germ cells. In addition, a high proportion of mutants derived from treatment of oocytes are mosaics, resulting from lesions affecting only one strand of the DNA. Treatment of oocytes also generates approximately 30% of large (multilocus) lesions. Only treatment of post-spermatogonial stages results in a higher percentage (65%) of such large lesions. The differences in the mutagenic potential of ENU on oocytes and stem cell spermatogonia occur mainly because oocytes are nondividing cells until fertilization. Consequently, mutagenesis of zygotes by an exposure of 50-mg ENU/kg at the time when the female genome enters metaphase II results in much higher mutation rates than in stem cell spermatogonia.

Different approaches have to be taken to recover a mutant mouse **genotype**. Dominant mutations can be screened directly in offspring of mutagenized males, whereas recessive mutations can be analyzed using a three-generation breeding protocol. Gene mutations in mouse germ or somatic cells by mutagens have been detected by the specific locus test, dominant cataract method, mouse spot test, micronucleus test, or new test systems like the MutaMouse or Big Blue mouse [see [Mouse](#)].

#### Suggestions for Further Reading

- S. D. M. Brown and J. Peters (1996) Combining mutagenesis and genomics in the mouse—closing the phenotype gap. *Trends Genetics* **12**, 433–435.
- L. B. Russel and W. L. Russel (1992) Frequency and nature of specific-locus mutations induced in female mice by radiations and chemicals: a review. *Mutation Res.* **296**, 107–127.
- T. Shibuya and K. Morimoto (1993) A review of the genotoxicity of 1-ethyl-1-nitrosourea. *Mutation Res.* **297**, 3–38.

## Euchromatin

Euchromatin, which can be distinguished from heterochromatin contains the majority of

transcriptionally active DNA. Euchromatin, initially identified cytologically during interphase, appears less condensed and stains less intensely with dyes than heterochromatin. Features of euchromatin may be usefully summarized by comparing the active X-chromosome with the inactive X-chromosome (see [Barr Body](#), [X-Chromosome Inactivation](#)).

The euchromatin of the active X-chromosome contains hyperacetylated histones, which indicate potential for transcriptional activity. The active X-chromosome also contains reduced levels of methylated DNA and replicates early in S-phase. Euchromatin is also likely to be globally sensitive to digestion with nucleases such as DNase I (see [DNase I Sensitivity](#)).

Existence in either the euchromatic or heterochromatic states is not irreversible. This is best demonstrated by experiments in which cells containing nuclei that are predominantly heterochromatic are either transplanted into eggs or fused with other cells that are predominantly euchromatic, to generate a heterokaryon. Experimental results with heterokaryons in which two different somatic cells are fused, so that two different nuclei share a common cytoplasm, and nuclear transplantation experiments using *Xenopus* eggs have been interpreted as providing evidence for continuous regulation of a plastic differentiated state. Implicit in this model is the idea that all genes are continually regulated by *trans*-acting factors that can either activate or repress genes (1). For certain genes this is clearly true. It has also been shown, however, that considerable remodeling of chromosomal structure occurs in *Xenopus* egg and oocyte cytoplasm, during which heterochromatin is converted to euchromatin (2). A similar, albeit less impressive, remodeling of chromosomes occurs in heterokaryons. For example, the nuclei of chicken erythrocytes consist predominantly of heterochromatin containing the specialized linker histone H5. In heterokaryons formed by the fusion of chicken erythrocytes with proliferating mammalian cells, the chicken erythrocyte nuclei once again become transcriptionally active, leading to the appearance of euchromatin (3). This process is accompanied by decondensation of chromatin, enlargement of the nucleus, and the appearance of nucleoli. Transcription and replication of these nuclei are also activated.

Enlargement of the chicken erythrocyte nucleus is caused by a massive, but selective, uptake of mammalian nuclear proteins, including RNA polymerases. The specialized linker histone H5 is partially lost from the chicken erythrocyte nucleus and partially taken up by the mammalian nucleus in the heterokaryon. Histones H2A and H2B also exchange under these circumstances, but not histones H3 and H4. These results might reflect the relative affinity of the histones for DNA and their organization in the nucleosome. This reorganization is independent of DNA replication. Therefore, it is clear that chromosomal structure is quite dynamic, and some histones (H1, H2A, H2B) continually exchange with a free pool of proteins in the cytoplasm.

Several experiments suggest that at physiological ionic strength the linker histone H1 rapidly exchanges into and out of the chromatin fiber (5, 6). Histone H1 is a specific repressor for several eukaryotic genes. Presumably this dynamic property of the chromatin fiber and the nucleosome would eventually allow many *trans*-acting factors to gain access to their cognate DNA sequences. An important and unresolved question is whether this access is unlimited or whether access is restricted by chromosomal organization. It has not yet been quantitatively determined whether the level of transcriptional activity following de novo activation of a gene in a heterokaryon is identical to the transcription of the same gene in a differentiated cell. Of course, in *Xenopus* egg cytoplasm, such equivalent activation must occur for correct development to proceed through to the tadpole stage. Nuclear reprogramming, however, is more rapid here and is likely to be facilitated by DNA replication (4), although massive nucleus-wide remodeling can occur without concomitant DNA synthesis (7, 8).

## Bibliography

1. H.M. Blau and D. Baltimore, *J. Cell Biol.* **112**, 781–783 (1991).
2. J.B. Gurdon, *J. Embryol. Exp. Morph.* **36**, 523–540 (1976).
3. N.R. Ringertz, U. Nyman, and M. Bergman, *Chomosome* **91**, 391–396 (1985).

4. S. Dimitrov and A.P. Wolffe, *EMBO J.* **15**, 5897–5906 (1996).
5. M.A. Lever, J.P. Th'ng, X. Sun, and M.J. Hendzel, *Nature* **408**, 873–876 (2000).
6. A.P. Wolffe and J.C. Hansen, *Cell* **104**, 631–634 (2001).
7. N. Kikyo et al., *Science* **289**, 2360–2362 (2000).
8. N. Kikyo and A.P. Wolffe, *J. Cell Sci.* **113**, 11–20 (2000).

### **Additional Reading**

9. Wolffe A. *Chromatin: Structure and Function*, 3rd ed., Academic Press, London, U.K., 1998.

## **Euploid**

A euploid cell is a cell that has the basic [haploid](#) number of [chromosomes](#) characteristic of that species or any exact multiple of this haploid number. Contrary situations are explained in [Aneuploidy](#) and [Polyploidy](#).

## **Evolution**

Darwin defined evolution in the *Origin of Species* as “descent with modification” (1). Although there have been various interpretations, modifications, and additions, the concepts of this definition continue to be repeated as biology develops. It seems likely that Darwin's definition will never fade away. Therefore, it is reasonable to state that evolution is the process in which the structures and functions of organisms are inherited from generation to generation, along with changes at both the morphological and molecular levels.

In his theory, Darwin treated species as populations, and he recognized several features of populations that explain the process of evolution (1). In particular, he noticed that among individuals having variation in characteristics, those having advantageous characteristics, or “fitness,” can survive by producing more offspring of their own. This is the process of [natural selection](#). Darwin already had advanced the concept that organisms have intrinsically inherited material labeled “genetic elements” (and later found to be DNA) and the idea that the variation found among individuals in a species was random, but these were not widely accepted. De Vries proposed the [mutation](#) theory, which states that mutations appear suddenly in a population and that species will experience periods of rapid mutation (2). His mutation theory and Mendel's law on inheritance, as well as Darwin's natural selection, were synthesized to create the synthetic or neo-Darwin theory of evolution (3). The synthetic theory of evolution views evolution in the following way. First, mutation occurs to the genetic material, namely, DNA, in a random fashion. It is a source of genetic variation, and it is heritable to the next generation by Mendelian law. Natural selection then takes place so that individuals having advantageous characteristics can survive in a population, and the offspring have the same type of mutation, which becomes spread in a population. In this way, the genetic composition changes with the time.

The study of molecular evolution is essentially based on the synthetic theory of evolution. However,

the development of molecular biology has introduced more detailed and new information on molecular mechanisms for producing genetic variation, which have accelerated the advancement of evolutionary studies. The genome projects currently in progress will give insights to the evolutionary changes of genomes, structures, and functions.

### Bibliography

1. C. Darwin (1859) *On the Origin of Species*, Murray, London.
2. H. De Vries (1901–1903) *Die Mutationstheorie*, Von Veit, Leipzig (1909–1910) "English translation", *The Mutation Theory*, trans. J. B. Farmer and A. D. Darbishire, Open Court, Chicago.
3. E. Mayer and W. B. Provine (1980) *The Evolutionary Synthesis*, Harvard Univ. Press, Harbridge, MA.

### Evolutionary Distance

Evolutionary distance is a distance by which evolutionary closeness or remoteness can be measured quantitatively. Evolutionary distance may be separated into two distinctive levels; one morphological; the other, molecular. In the case of evolutionary distance at the molecular level, the number of nucleotide and amino acid substitutions, as well as immunological cross-reactions, are the most popular measures for evolutionary distance.

The number of nucleotide substitutions is estimated by making pairwise comparison of nucleotide sequences and correcting for multiple substitutions at the same site (see ). For correction of multiple substitutions, one needs a model of nucleotide substitution; for example, the one-parameter method invented by Jukes and Cantor (1) assumed that the rates of nucleotide substitutions between all possible pairs of different nucleotides are equal to each other. In this case, the nucleotide substitution is estimated by the equation  $K_n = -\frac{3}{4} \ln(1 - \frac{4}{3}p)$ , where  $p$  is the fraction of nucleotide differences. There are many modified and extended versions of this formula. For example, Kimura's two-parameter method was developed under the assumption that the rates of transition mutations and are different (2). In this way, four-parameter method and the six-parameter methods were also developed. The method of maximum likelihood was invented.

The number of amino acid substitutions is also used as a measure of evolutionary distance. Under the assumption that the substitution rates between any pair of amino acids are equal, the number of amino acid substitutions can be given by formula  $K_a = -\ln(1-p)$ , where  $p$  represents the fraction of amino acid differences. It is well known, however, that the substitution rate between a pair of similar amino acids is much higher than that between a pair of nonsimilar amino acids. Although Dayhoff (3) developed the algorithm to estimate the number of amino acid substitutions taking into account the similarity matrix of amino acids, the formula given was very complicated. To minimize this, Kimura (2) invented an empirical formula by adding one simple term,  $K_a = -\ln(1 - p - \frac{1}{5}p^2)$ . Even though Kimura's formula was derived empirically, the substitution number estimated by this formula is very close to Dayhoff's estimate.

Immunological distances can be also used as a measure of evolutionary distance. The measurement of immunological distance uses as a measurement the intensity of immunological crossreaction between antigens and antisera that were prepared from different species (see ). It has now become less popular, however, because the numbers of amino acid or nucleotide substitutions contain much

more quantitative information.

## Bibliography

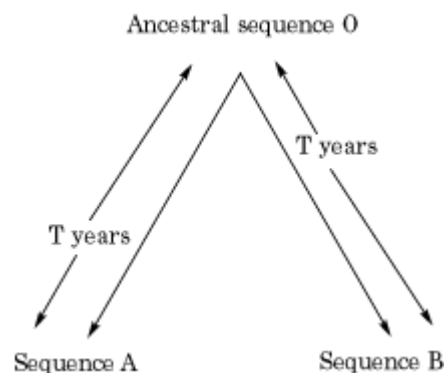
“Evolutionary Distance” in , Vol. 2, pp. 867–868, by T. Gojobori; “Evolutionary Distance” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. T. H. Jukes and C. R. Cantor (1969) *Mammalian Protein Metabolism*, Academic Press, New York.
2. M. Kimura (1983) *The Neutral Theory of Molecular Evolution*, Cambridge Univ. Press, Cambridge, UK.
3. M. O. Dayhoff (1978) *Atlas of Protein Sequence and Structure*, Vol. 5, Suppl. 3, National Biomedical Research Foundation, Washington, D.C.

## Evolutionary Rate

In the studies of molecular evolution, the rates of nucleotide or amino acid substitutions (per site per year) are usually used as the evolutionary rate (see [Sequence Analysis](#)). Compare two nucleotide sequences, A and B, as an example. Denote  $K_n$  as the number of nucleotide substitutions between sequences A and B. When the ancestral node, O, of the two sequences A and B is considered, it follows that nucleotide substitutions should have accumulated during the time of lineages OA and OB (Figure 1). If the divergence time between sequences A and B, which is measured by years, is denoted  $T$ ,  $K_n$  is the number of substitutions that have taken place during  $2T$ . This is because the divergence time of OA plus OB should be equal to  $2T$ . Thus, the rate of nucleotide substitution ( $v_n$ ) can be computed by the formula  $v_n = K_n / (2T)$ . The rate of amino acid substitution ( $v_a$ ) can be obtained in a similar way:  $v_a = K_a / (2T)$ . Note that the number of nucleotide and amino acid substitutions can be estimated by a number of different methods (see [Evolutionary Distance](#)).

**Figure 1.** Estimating the rate of evolutionary change in nucleotide or amino acid sequences from the difference between two present-day homologous sequences. Rate of nucleotide substitution:  $v_n = K_n / (2T)$ ; rate of amino acid substitution:  $v_a = K_a / (2T)$



## Excision Repair

Excision repair (1, 2) is a form of [DNA repair](#) involving the removal of a damaged nucleotide from DNA by dual incisions bracketing the damage (3, 4). The multisubunit [enzyme](#) system that makes the dual incision is referred to as *excision nuclease* or *excinuclease* (3). The damaged nucleotide is released in the form of an oligonucleotide, and the resulting single-stranded gap is filled in by **DNA polymerase** and sealed by a [DNA Ligase](#). This repair mechanism, which is also referred to as “nucleotide excision repair,” is a universal repair system. Excinuclease has been found in all free-living organisms, including [Mycoplasma genitalium](#), which is considered a model system for “minimal cell.” Although both prokaryotes and eukaryotes carry out excision repair, there are distinct differences in the mechanism in the two systems.

### 1. Substrate

In most organisms, excision repair is the sole mechanism for repairing bulky DNA lesions. These include the vast number of base adducts formed by polyaromatic hydrocarbons and by chemotherapeutic agents such as cisplatin, mitomycin C, and [psoralen](#), plus the cyclobutane pyrimidine dimers and the [6–4] photoproducts induced by ultraviolet light (see [Photolyase/Photoreactivation](#)). Nucleotide excision repair is not, however, restricted to bulky adducts. Alkylated bases including *O*<sup>6</sup>-methylguanine and 3-methyladenine, oxidative base lesions such as 8-oxoguanine and thymine glycols, and even apurinic/apyrimidinic sites ([AP sites](#)) are removed from DNA by excision repair (5-7). Mismatched bases are also removed with low efficiency by this enzyme system; unlike [mismatch repair](#), however, the excinuclease cannot discriminate between the correct base and the mismatched base, and this may result in fixation of a [mutation](#) (5). The most impressive aspect of the substrate repertoire of excinuclease is that it encompasses essentially all abnormal DNA structures. It must be noted, however, that, with the exception of bulky lesions, all other lesions can be repaired by alternative repair mechanisms, including repair by *O*<sup>6</sup>-**methylguanine DNA methyltransferase**, **DNA glycosylases**, **AP endonucleases**, and the mismatch repair systems. It is most likely that excision repair plays a critical backup role for these other repair systems, by repairing nonbulky DNA lesions when the normal repair enzymes are poorly expressed in a certain tissue or when a sudden burst of DNA damage exceeds the repair capacity of the specific repair mechanisms. It has been suggested that the neurological symptoms of xeroderma pigmentosum patients are caused by metabolism-induced nonbulky oxidative damage in brain cells, which overwhelms the main repair system for DNA damage caused by oxidative stress, [base excision repair](#) (7).

### 2. Genetics

Excision repair encompasses three basic enzymatic steps: (i) dual incisions on both sides of the damage, (ii) repair synthesis to fill in the excision gap, and (iii) ligation of the repair patch. Mutations in genes controlling any of these reactions interfere with normal functioning of excision repair. The enzymes carrying out repair synthesis and ligation are not, however, dedicated solely to repair; they also perform DNA synthesis and ligation functions during replication, recombination, and transposition. Strictly speaking, therefore, excision repair genes are the genes encoding the subunits of the excinuclease, which is the enzyme system that removes the damaged nucleotide by dual incisions bracketing the lesion. In *E. coli* and other prokaryotes, three genes—*uvrA*, *uvrB*, and *uvrC*—are necessary and sufficient for the excision reaction (3). In the yeast *Saccharomyces cerevisiae*, genes in the RAD3 epistasis group (RAD1, 2, 3, 4, 10, 14, 25) are required to reconstitute

excinuclease (8). In humans, defects in excision repair cause xeroderma pigmentosum (XP). In this hereditary recessive disease, patients are hypersensitive to sunlight and develop skin cancers, including melanomas, at an early age. In addition, the majority of XP patients exhibit neurological signs and symptoms caused by neuronal death. There are eight XP complementation groups, XP-A through XP-G, plus the XP-variant (XP-V). Patients in complementation groups XP-A to XP-G are defective in removing pyrimidine dimers and other bulky adducts from DNA; in other words, they are defective in excision repair. In contrast, XP-V-deficient individuals exhibit normal excinuclease activity, but are defective in a process called postreplication repair. The XP complementation groups define only seven of the proteins involved in the excision reaction. In addition to these proteins, the replication protein RPA, the recombination protein ERCC1 (which makes a complex with XPF), and the transcription factor TFIIH (which includes the XPB and XPD [DNA helicases](#) and four to six additional polypeptides) are also required for dual incisions (9, 10).

There is a one-to-one correspondence between the human and yeast excinuclease genes, and similar homologous genes have been identified in all other eukaryotes tested (11). Hence, it appears that the excision repair genes and mechanism have been conserved throughout the eukaryotic kingdom. In contrast, there is no homology between the excinuclease genes of eukaryotes and prokaryotes. Certain species in the third form of life, the Archaea, appear to have XPD, XPF, and XPG homologues, and certain others have *uvrA*, *B*, *C* homologues. Hence, it appears that species in this kingdom might have either the prokaryotic or eukaryotic form of excinuclease.

### 3. Biochemistry

The mechanism of the excision repair reaction is quite similar in prokaryotes and eukaryotes, even though there is no structural homology between the enzyme systems: Damage is recognized by the enzyme in an ATP-independent manner, forming an unstable DNA–protein complex; then ATP hydrolysis by the proteins within the complex promotes formation of a long-lived preincision complex. Binding of a nuclease subunit triggers dual incisions, which are concerted but nonsynchronous; that is, the 3' and 5' incisions nearly always occur in the same DNA molecule, but the 3' incision occurs first, followed within a fraction of a second by the 5' incision. Following both incisions, the excinuclease proteins are replaced by replication proteins. In *E. coli*, DNA polymerase I and DNA helicase II (UvrD protein), working together, displace the excised oligomer and the subunits of the excision nuclease that remain in the postexcision complex, and they then fill in the gap (12, 13). In yeast and humans, the replication polymerase  $\delta$  or  $\epsilon$ , working together with the polymerase **clamp loader** replication factor C (RF-C) and the polymerase clamp, the [proliferating cell nuclear antigen](#) (PCNA), fill in the excision gap (14). The last step is ligation of the repair patch to the parental DNA. A more detailed description of excision repair in *E. coli* and humans is given below.

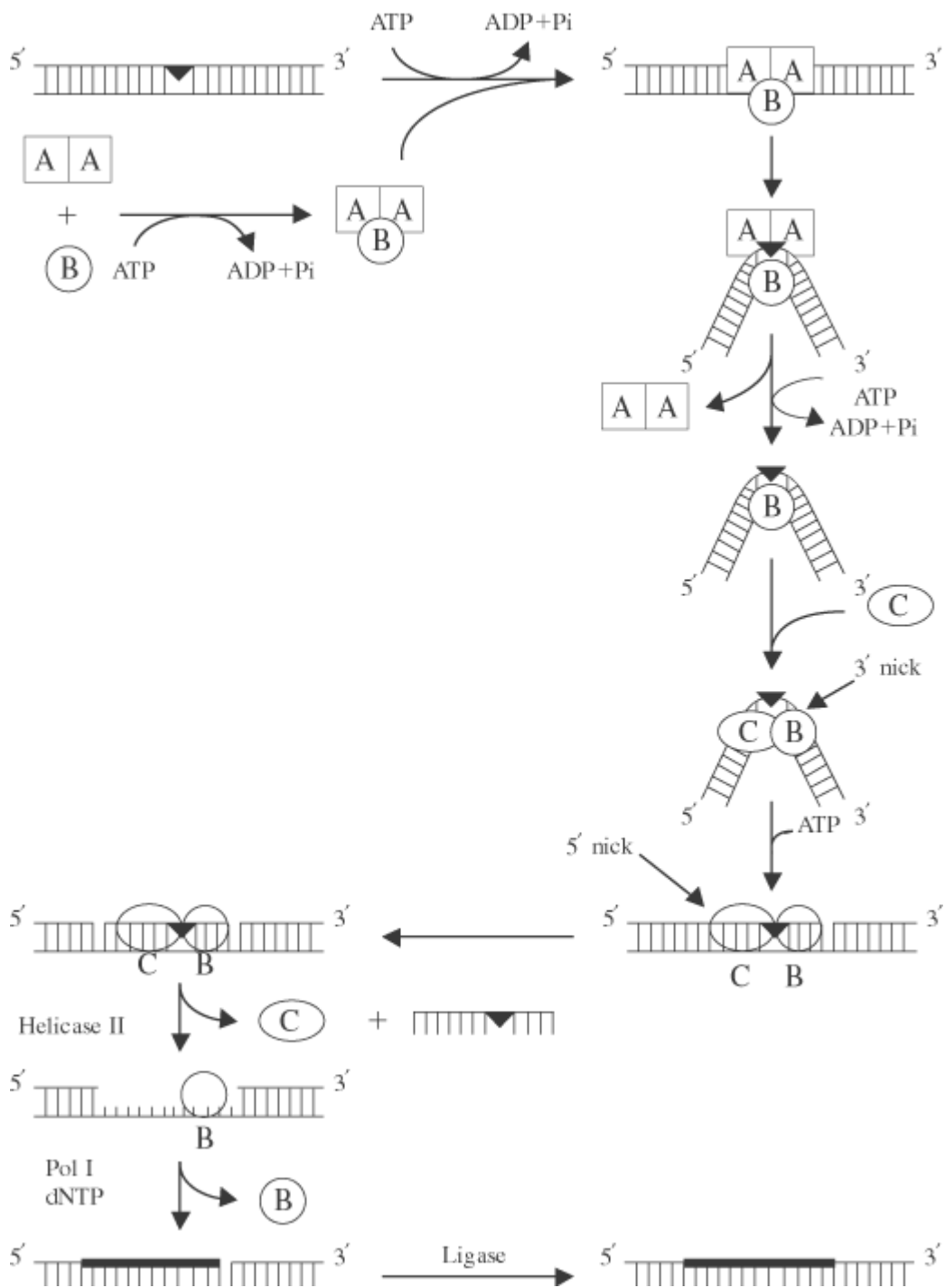
#### 3.1. Excision Repair in *E. coli* (Fig. 1)

Three subunits, UvrA, UvrB, and UvrC, are necessary and sufficient to carry out the dual incision. Their properties are summarized in Table 1. UvrA and UvrB make a heterotrimer of  $A_2B_1$  composition that, guided by the intrinsic affinity of UvrA for damaged DNA (15, 16), recognizes subtle or gross abnormalities in the duplex structure as damage. Upon binding to the damaged DNA, ATP hydrolysis by UvrA and UvrB leads to local unwinding of the DNA and formation of a stable UvrB–DNA complex, concomitant with dissociation of UvrA from the complex (16). Here UvrA plays the role of a molecular matchmaker (17): Utilizing the energy released by ATP hydrolysis, it brings together DNA and UvrB (a protein with rather weak affinity for DNA), promotes formation of a complex between the two, and then leaves so that the complex can engage in productive reactions. Following dissociation of UvrA, the UvrC subunit binds to the UvrB–DNA complex, and the 3' incision is made at the fourth or fifth phosphodiester bond from the damage by UvrB, perhaps with the contribution of [active-site](#) residues of UvrC (18). After the 3' incision, the complex undergoes another conformational change that enables UvrC, perhaps aided by residues on UvrB, to make the 5' incision at the eighth phosphodiester bond 5' to the lesion (18). After the dual incision, UvrB,



UvrC, and the “excised” dodecamer remain in a complex with the duplex. Helicase II (UvrD) binds to this complex and releases the excised oligomer and UvrC. Then, DNA polymerase I binds to the 5′ incision site and displaces UvrB, while simultaneously filling the excision gap. When the gap is filled to the 3′ end, the repaired patch is ligated to the parental DNA. There is no **nick translation** in filling in the gap and, as a consequence, the repair patch size matches precisely the size of the excision gap (13).

**Figure 1.** Excision repair in *E. coli*. UvrA and UvrB make up the A<sub>2</sub>B<sub>1</sub> heterotrimer that, guided by the affinity of UvrA to helical distortion, recognizes DNA damage. Upon recognition, the DNA is locally unwound by the helicase action of the complex and kinked by 130°. Having performed its molecular matchmaking function, UvrA dissociates, leaving behind a stable UvrB–DNA complex that is recognized by UvrC. Upon binding of UvrC, UvrB makes the 3′ incision at the fifth phosphodiester bond and UvrC makes the 5′ incision at the eighth phosphodiester bond. Helicase II (UvrD) releases the excised dodecamer and UvrC; DNA pol I displaces UvrB and fills in the gap, which is ligated.



**Table 1. Subunits of *E. coli* Excinuclease**

| Protein                           | $M_r$ | Sequence Motifs | Activity   | Role in Repair |
|-----------------------------------|-------|-----------------|------------|----------------|
| <i>Excision Nuclease Subunits</i> |       |                 |            |                |
| I. UvrA                           | (104) | (a) Walker      | ATPase (2) | (a) Damage     |

|              |    |   |                                   |                                      |
|--------------|----|---|-----------------------------------|--------------------------------------|
|              | 2  |   |                                   | recognition<br>(proximal)            |
|              |    | (b) Zinc finger (2)                           | (b) Damage-specific DNA binding   | (b) Molecular matchmaker             |
|              |    | (c) Leucine zipper                            | (c) UvrB binding                  | (c) TRC                              |
|              |    | (d) UvrA superfamily                          | (d) TRCF binding                  |                                      |
| II.<br>UvrB  | 78 | (a) Helicase motif                            | (a) Latent ATPase                 | (a) Damage recognition<br>(ultimate) |
|              |    | (b) Homology to TRCF                          | (b) Latent “helicase”             | (b) Unwinding duplex                 |
|              |    |   | (c) Damage-specific ssDNA binding |                                      |
|              |    |   | (d) Binds UvrA                    | (c) Makes 3' incision                |
|              |    |   | (e) Binds UvrC                    |                                      |
| III.<br>UvrC | 69 | (a) Limited homology to UvrB                  | (a) Nonspecific DNA binding       | (a) Induces 3' incision              |
|              |    | (b) Limited (40 amino acid) homology to ERCC1 | (b) UvrB binding                  | (b) Makes 5' incision                |

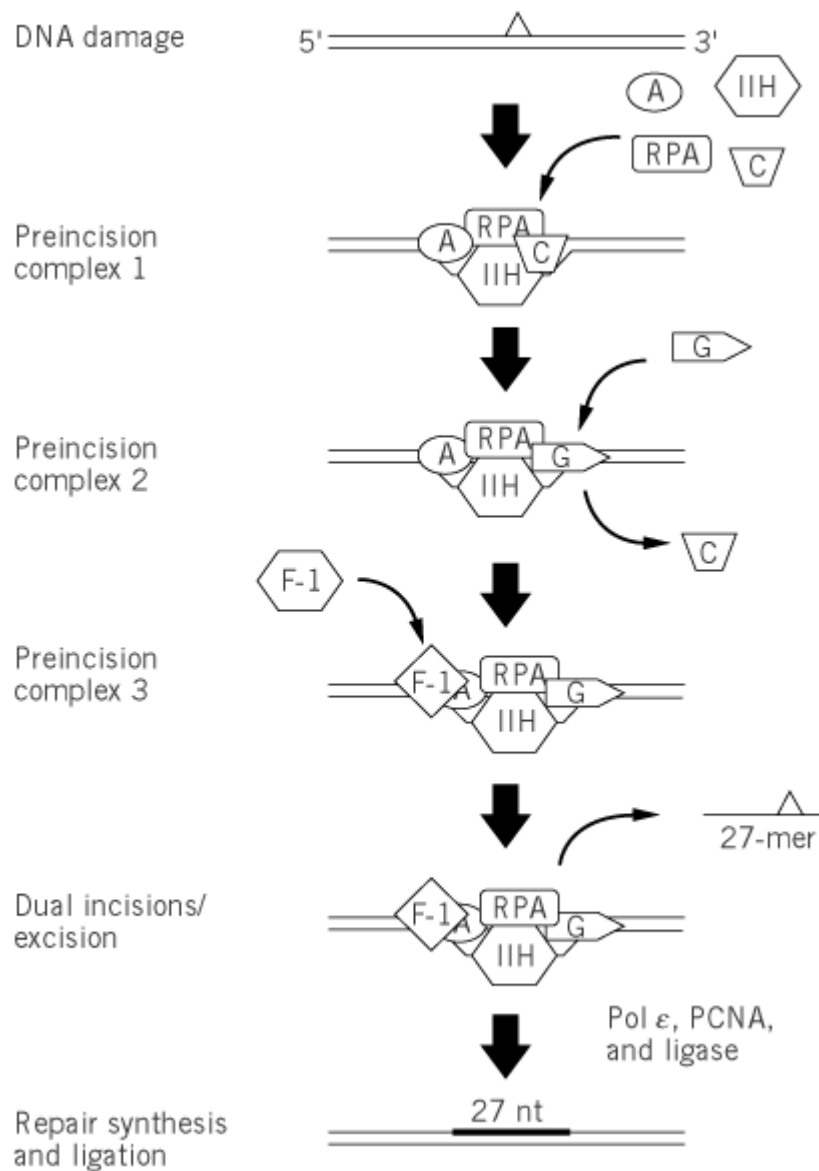
---

### 3.2. Excision Repair in Humans (Fig. 2)

In humans, the excinuclease activity of producing dual incisions results from the coordinated action of 16 polypeptide chains in six repair factors, listed in Table 2 (9, 10). The basic reaction carried out by human excinuclease is the incision of a phosphodiester bond that is  $6 \pm 3$  nucleotides 3' to the damage and  $20 \pm 5$  nucleotides 5', followed by release of the resulting 24- to 32-nucleotide oligomers carrying the damaged base (4, 11). The reaction proceeds as follows (19): XPA, RPA, and XPC · HHR23B (20) proteins form a transient complex at the site of the damage. This complex interacts with TFIIH, which is both a transcription factor and a repair factor (21), via XPA and XPC · HHR23B and recruits this multisubunit factor to the lesion site. TFIIH unwinds the duplex by about 5 nucleotides on the 3' side and 10 nucleotides on the 5' side of the damaged nucleotides (19); as a consequence, it promotes new DNA–protein interactions that lead to formation of a stable preincision complex (22). XPG has high affinity for TFIIH and enters the preincision complex, either simultaneously with TFIIH or subsequently after being recruited by TFIIH, to form a very stable DNA–protein complex in which the DNA is unwound by 15 to 20 nucleotides. Entry of XPG to the preincision complex is dependent on, and concomitant with, departure of XPC from this complex. Hence, XPC plays the role of molecular matchmaker in human excinuclease. Within this complex,

the 3' incision is made by XPG (23, 24) at the sixth  $\pm$  3 phosphodiester bond 3' to the damage. Thus, in contrast to *E. coli*, the 3' incision can be made in humans before binding of the subunit that makes the 5' incision. The XPF · ERCC1 complex is recruited to the repair assembly by the strong interaction between XPA and ERCC1, which is promoted by the presence of XPG within the complex. The XPF · ERCC1 complex makes the 5' incision at the 20th  $\pm$  5 phosphodiester bond 5' to the damage (24). Assembly of the excision nuclease is a highly concerted reaction, so that under physiological conditions the 3' incision is rapidly followed by the 5' incision; therefore, uncoupled 3' incisions do not accumulate in the cell. The dual incision generates a single-stranded DNA fragment of 24 to 32 nucleotides, which is released from the duplex (excision) (10). This is in contrast to the bacterial system, in which the “excised” oligomer remains in the postincision complex until it is released by a DNA repair auxiliary helicase (10, 12, 13). Following release of the excised oligomer, the gap is filled by DNA polymerases  $\delta$  and  $\epsilon$  in a reaction that is dependent on the polymerase clamp, the proliferating cell nuclear antigen (PCNA). The size of the repair patch precisely matches the size of the excision gap, with no gap enlargement on the 5' side or nick translation on the 3' side of the gap (25).

**Figure 2.** Excision repair in humans. XPA, RPA, XPC, and TFIIH recognize damaged DNA. TFIIH unwinds the DNA by about 20 nucleotides, and a transient DNA–protein complex forms. This structure is recognized by XPG, which enters the complex concomitant with dissociation of XPC (molecular matchmaker). XPG makes the 3' incision at the sixth phosphodiester bond and recruits the XPF · ERCC1 complex, which makes the 5' incision at the 22nd phosphodiester bond. The excised 27-mer is released along with some of the excinuclease subunits. The excision gap is filled in by DNA Pol  $\delta$  and ligated.



**Table 2. Subunits of Human Excinuclease**

| Factor  | Number | Proteins <sup>a</sup> | Sequence Motif | Activity    | Role in Repair     |
|---------|--------|-----------------------|----------------|-------------|--------------------|
| I. XPA  | 1      | XPA/p31               | Zinc finger    | DNA binding | Damage recognition |
| II. RPA | 2      | p70                   |                | DNA binding | Damage recognition |
|         | 3      | p34                   |                |             |                    |
|         | 4      | p11                   |                |             |                    |

|               |            |                |                |                          |                                      |  |
|---------------|------------|----------------|----------------|--------------------------|--------------------------------------|--|
| III.<br>TFIIH | 5          | CPB/ERCC3/p89  | Helicase       | (a) DNA-dependent ATPase | (a) Formation of preincision complex |  |
|               | 6          | XPD/ERCC2/p80  | Helicase       |                          |                                      |  |
|               | 7          | p62            |                |                          |                                      |  |
|               | 8          | p44            | Zinc finger    | (b) "Helicase"           | (b) Transcription-repair coupling    |  |
|               | 9          | Cdk7/p41       | S/T kinase     |                          |                                      |  |
|               | 10         | CycH/p38       | Cyclin         | (c) GTF                  |                                      |  |
|               | 11         | p34            | Zinc finger    | (d) CAK                  |                                      |  |
|               | IV.<br>XPC | 12             | XPC/p125       |                          | DNA binding                          | (a) Stabilization of preincision complex               |
|               |            | 13             | HHR23B/p58     | Ubiquitin                |                                      | (b) Protection of preincision complex from degradation |
|               | V.<br>XPF  | 14             | XPF/ERCC4/p112 |                          | Nuclease                             | 5' Incision  |
| 15            |            | ERCC1/p33      |                |                          |                                      |  |
| VI.<br>XPG    | 16         | XPG/ERCC5/p135 |                | Nuclease                 | 3' Incision                          |  |

<sup>a</sup> The relative molecular weights of the polypeptides are indicated, for example, by p31 for  $M_r = 31$  kDa, and so on. GTF, general transcription factor; CAK, cyclin-dependent kinase activating kinase.

#### 4. Transcription-Repair Coupling

Transcribed sequences are repaired more rapidly by excinuclease than are nontranscribed sequences, both in prokaryotes and in eukaryotes (26, 27). The open [chromatin](#) conformation of transcribed DNA in eukaryotes contributes to the accessibility of the DNA to repair enzymes, and hence it has a faster rate of repair compared to regions of condensed chromatin conformation. The main driving force of preferential repair of transcribed DNA, however, is the mechanistic coupling of transcription to repair. The best manifestation of this active process of coupling is the fact that, although both strands of the DNA in the gene are more accessible because of the open conformation, the transcribed strand is repaired faster than the nontranscribed (coding) DNA strand, which is repaired at essentially the same rate as the rest of the genome (28, 29). Hence it appears that the main cause of

repair coupled to transcription is the recruitment of the excinuclease complex to the site of damage (30). Indeed, in both prokaryotes and eukaryotes, the damage-recognition step is the slowest step of excision repair. In the act of transcription, lesions in the template strand block the progression of **RNA polymerase**—hence the stalled complex marks the site of the lesion. In other words, RNA polymerase functions as a surrogate damage-recognition protein for lesions in the template strand. Because specific binding of repair proteins to the damage itself is very slow, due to the fact that lesions removed by excinuclease have very little in common, a stalled RNA polymerase provides a fast damage-locating mechanism and a common structural element (RNA polymerase in elongation mode) for the excision repair system. It appears both in prokaryotes and eukaryotes that a specific protein called transcription-repair coupling factor (TRCF) recognizes both the stalled RNA polymerase and a subunit of the excision nuclease and recruits the repair enzyme to the site of the lesion (30). As a result of high-specificity [protein–protein interactions](#), rather than the low-specificity **DNA–protein interaction**, the rate of repair is accelerated. The transcription repair coupling factor is encoded by the *mfd* gene in *E. coli*, the RAD26 gene in yeast, and the *CSB* gene in humans. The molecular mechanism of transcription-repair coupling is relatively well-understood in *E. coli*, but not in eukaryotes.

#### 4.1. Transcription-Repair Coupling in *E. coli*

The *E. coli* TRCF is a monomer of 130 kDa, and it has helicase motifs, a UvrA binding site in the NH<sub>2</sub>-terminal half, and a RNAP-binding region in the middle (31). The protein recognizes RNA polymerase stalled at a lesion and displaces the stalled polymerase, while simultaneously recruiting the A<sub>2</sub>B<sub>1</sub> complex to the damage site. The net result is the discarding of the truncated transcript and a substantial increase in the rate of damage recognition. Furthermore, because TRCF and UvrB share the same binding site on UvrA, the TRCF also helps dissociate UvrA from the A<sub>2</sub>B<sub>1</sub>–DNA complex and facilitates rapid formation of the UvrB–DNA complex (31), which is a prerequisite for binding of UvrC. This is followed by dual incision and the remaining steps of excision repair, which are not rate-limiting. Therefore, acceleration of the damage recognition reaction and facilitation of the dissociation of UvrA by the TRCF results in overall enhancement of the excision repair rate by nearly a factor of 10 relative to excision repair unaided by RNA polymerase and TRCF. The coupling of transcription to repair results in a decline in mutation frequency of ultraviolet-irradiated cells incubated in a medium lacking amino acids prior to plating, compared to cells plated immediately on rich medium (32). In the minimal medium, the mutagenic lesions in the transcribed genes are repaired before giving rise to mutation by DNA synthesis and replication of the lesion. Hence the gene for the *E. coli* TRCF is called *mfd* (mutation frequency decline) (32).

#### 4.2. Transcription-Repair Coupling in Humans

The phenomenology of transcription-repair coupling in humans is similar to that in *E. coli*, although there are some important differences. Humans defective in transcription-repair coupling are afflicted with Cockayne's syndrome, a disease characterized by growth and developmental abnormalities and mild sensitivity to sunlight. Mutations in two genes, *CSA* and *CSB*, give rise to Cockayne's syndrome. Both genes have been **cloned**. The *CSB* gene is a 160-kDa monomer with helicase motifs and appears to be the functional homologue of the *mfd* gene of *E. coli*. The *CSA* protein is a 55-kDa protein with the WD sequence motif, a protein–protein interaction motif. The *CSB* protein is a weak [ATPase](#) and binds to RNA polymerase II either free in solution or stalled on DNA. Significantly, in contrast to the *E. coli* TRCF, the *CSB* protein, even in combination with *CSA* protein, does not disrupt the ternary complex of RNA polymerase II stalled at a lesion (33). *CSB* does bind to XPA and XPG subunits of the human excinuclease, however, and may play a role similar to that of the *E. coli* Mfd protein in recruiting the repair complex to the site of damage. The mechanistic details of coupling transcription to repair are not known at present.

Both in *E. coli* and in humans, it is the nucleotide excision repair that is tightly coupled to transcription. Recently it was found, however, that transcription also stimulates removal of thymine glycols by [base excision repair](#) in a reaction that depends on the XPG protein, but not other subunits of the human excision nuclease system. The mechanistic aspects of this coupling remain to be

elucidated.

## 5. Regulation of Excision Repair

There is a steady-state level of damage production and removal by nucleotide excision repair under physiological conditions. It would be advantageous for the cell if, upon any sudden increase of the genotoxic load, it could increase its repair capacity. Three such regulatory systems have been characterized in *E. coli*: the [SOS response](#), the adaptive response, and the adaptive response to oxidative stress. The SOS response increases the nucleotide-excision-repair capacity of the cell by increasing the level of UvrA and UvrB subunits of excinuclease. Although DNA damage induces a dramatic stress response reaction in human cells, increased repair capacity is not part of this reaction. Human excinuclease genes are not induced at the transcriptional level by DNA damage, nor is the repair activity increased as a consequence of post-transcriptional events. In particular, DNA damage causes hyperphosphorylation of the RPA p34 subunit, but this has no effect on excision repair. Similarly, certain types of damage increase the level of [p53](#) protein by [post-transcriptional](#) mechanisms. It has been reported that this increase stimulates repair in certain cell types, but not others. The mechanism of the p53 effect on repair is not known.

## Bibliography

1. R. P. Boyce and P. Howard-Flanders (1964) Proc. Natl. Acad. Sci. USA **51**, 293–300.
2. R. B. Setlow and W. L. Carrier (1964) Proc. Natl. Acad. Sci. USA **51**, 226–231.
3. A. Sancar and W. D. Rupp (1983) Cell **33**, 249–260.
4. J. C. Huang, D. L. Svoboda, J. T. Reardon, and A. Sancar (1992) Proc. Natl. Acad. Sci. USA **89**, 3664–3668.
5. J. C. Huang, D. S. Hsu, A. Kazantsev, and A. Sancar (1994) Proc. Natl. Acad. Sci. USA **91**, 12213–12217.
6. J. C. Huang, D. B. Zamble, J. T. Reardon, S. J. Lippard, and A. Sancar (1994) Proc. Natl. Acad. Sci. USA **91**, 10394–10398.
7. J. T. Reardon, T. Bessho, H. C. Kung, P. H. Bolton, and A. Sancar (1997) Proc. Natl. Acad. Sci. USA **94**, 9436–9438.
8. S. N. Guzder, Y. Habraken, P. Sung, L. Prakash, and S. Prakash (1995) J. Biol. Chem. **270**, 12973–12976.
9. D. Mu, C. H. Park, T. Matsunaga, D. S. Hsu, J. T. Reardon, and A. Sancar (1995) J. Biol. Chem. **270**, 2415–2418.
10. D. Mu, D. S. Hsu, and A. Sancar (1996) J. Biol. Chem. **271**, 8285–8294.
11. D. L. Svoboda, J. S. Taylor, J. E. Hearst, and A. Sancar (1993) J. Biol. Chem. **268**, 1931–1936.
12. P. R. Caron, S. R. Kushner, and L. Grossman (1985) Proc. Natl. Acad. Sci. USA **82**, 4925–4929.
13. I. Husain, B. Van Houten, D. C. Thomas, M. Abdel-Monem, and A. Sancar (1985) Proc. Natl. Acad. Sci. USA **82**, 6774–6778.
14. A. F. Nichols and A. Sancar (1992) Nucleic Acids Res. **20**, 2441–2446.
15. J. T. Reardon, A. F. Nichols, S. Keeney, C. A. Smith, J. S. Taylor, S. Linn, and A. Sancar (1993) J. Biol. Chem. **268**, 21301–21308.
16. D. K. Orren and A. Sancar (1989) Proc. Natl. Acad. Sci. USA **86**, 5237–5241.
17. A. Sancar and J. E. Hearst (1993) Science **259**, 1415–1420.
18. J. J. Lin and A. Sancar (1992) J. Biol. Chem. **267**, 17688–17692.
19. M. Wakasugi and A. Sancar (1997) Proc. Natl. Acad. Sci. USA **95**, 6679–6684.
20. C. Masutani, K. Sugawara, J. Yonagisawa, T. Sanoyama, M. Ui, T. Enomoto, K. Takio, K. Tanaka, P. J. van der Spek, D. Bootsma, J. H. J. Hoeijmakers, and F. Hanaoka (1994) EMBO J. **13**, 1831–1843.



21. R. Drapkin, J. T. Reardon, A. Ansari, J. C. Huang, L. Zawel, K. Ahn, A. Sancar, and D. Reinberg (1994) *Nature* **368**, 769–772.
22. M. T. Hess, U. Schwitter, M. Petretta, B. Giese, and H. Naegeli (1997) *Proc. Natl. Acad. Sci. USA* **94**, 6664–6669.
23. A. O'Donovan, A. A. Davis, J. A. Moggs, S. C. West, and R. D. Wood (1994) *Nature* **371**, 432–435.
24. T. Matsunaga, C. H. Park, T. Bessho, D. Mu, and A. Sancar (1996) *J. Biol. Chem.* **271**, 11047–11050.
25. J. T. Reardon, L. H. Thompson, and A. Sancar (1997) *Nucleic Acids Res.* **25**, 1015–1021.
26. R. C. Bockrath and J. E. Palmer (1997) *Mol. Gen. Genet.* **156**, 133–140.
27. V. A. Bohr, C. A. Smith, D. S. Okumoto, and P. C. Hanawalt (1985) *Cell* **40**, 359–369.
28. I. Mellon and P. C. Hanawalt (1989) *Nature* **342**, 95–98.
29. I. Mellon, G. Spirak, and P. C. Hanawalt (1987) *Cell* **51**, 241–249.
30. C. P. Selby and A. Sancar (1993) *Science* **260**, 53–58.
31. C. P. Selby and A. Sancar (1995) *J. Biol. Chem.* **270**, 4882–4889.
32. E. M. Witkin (1966) *Science* **152**, 1345–1353.
33. C. P. Selby and A. Sancar (1997) *J. Biol. Chem.* **272**, 1885–1890.

### Suggestions for Further Reading

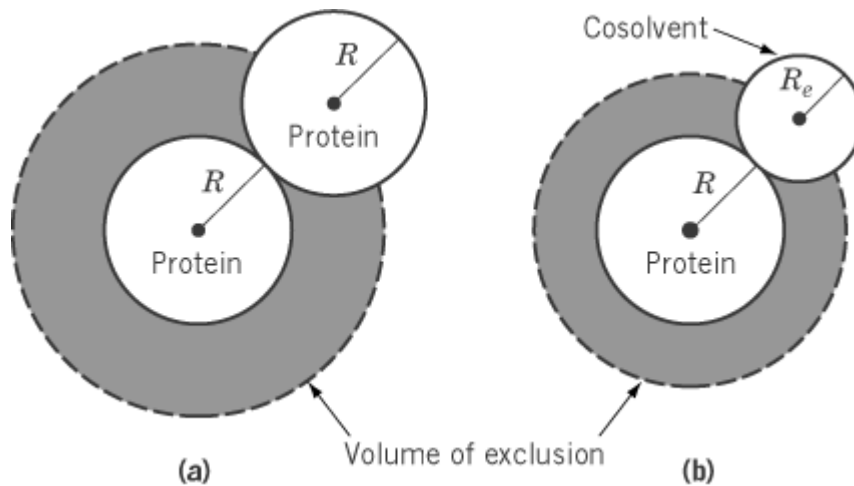
34. A. Sancar (1994) Mechanisms of DNA excision repair. *Science* **266**, 1954–1956.
35. A. Sancar (1995) Excision repair in mammalian cells. *J. Biol. Chem.* **270**, 15915–15918.
36. A. Sancar (1996) DNA excision repair. *Annu. Rev. Biochem.* **65**, 43–81.
37. S. Prakash, P. Sung, and L. Prakash (1993) DNA repair genes and proteins of *Saccharomyces cerevisiae*. *Annu. Rev. Genet.* **27**, 33–70.
38. E. C. Friedberg, G. C. Walker, and W. Siede (1995) *DNA Repair and Mutagenesis*, ASM Press.
39. R. D. Wood (1996) DNA repair in eukaryotes. *Annu. Rev. Biochem.* **65**, 135–167.
40. P. C. Hanawalt (1994) Transcription-coupled repair and human disease. *Science* **266**, 1957–1958.

### Excluded Volume

No solute molecule in a solution can be present in the same space as any other molecule in the solvent. This volume of occupation is called the excluded volume  $V_e$ . The excluded volume clearly depends on the shape and conformation of the solute molecule. In a dilute solution, the centers of two identical rigid spheres cannot be placed closer to each other than a distance twice as great as their radius (Fig. 1a). Thus, any spherical molecule excludes similar molecules from a volume eight times greater than that it physically occupies. However, the excluded volume due to the interaction of any pair of particles is counted twice, so the volume excluded by all the spherical particles in the system is four times greater than their physical volumes.

**Figure 1.** Schematic illustration of the excluded volume of spherical molecules. (a) Each solute molecule is excluded from the space occupied by other solute molecules in solution. In very dilute solutions, each molecule of radius  $R$

contributes a spherical volume of radius  $2R$ . (b) A cosolvent of radius  $R_e$  much greater than that of water statistically cannot be close to a protein molecule due to a steric exclusion principle, resulting in the preferential hydration of protein. [Taken from T. Arakawa and S. N. Timasheff (1).]



The [partial specific](#) (or molar) **volume**  $v_2$  is usually taken as the physical volume, so  $V_e = 4v_2M_2$  for a spherical solute molecule of molecular weight  $M_2$ . This should be a good approximation for globular proteins. A rigid rod-shaped solute is regarded as a cylinder, so its  $V_e$  is proportional to the length–diameter ratio of the cylinder. The excluded volume of a flexible polymer depends on its [radius of gyration](#), which is sensitive to the polymer–solvent interaction. In a good solvent, solvent molecules are preferably accessible to a polymer, which extends the chain segment and increases the excluded volume. On the contrary,  $V_e$  is small in a poor solvent, which enhances interactions between different parts of the polymer.

Many polymer solution theories have been developed for the excluded volume of flexible polymers. It is important that the excluded volume be directly related to the second virial coefficient of polymer solutions, so this is an apt criterion for the ideality of the solutions. The excluded volume effect is more significant for nucleic acids at lower salt concentration, where the stiffness or persistence length of the chain segment is influenced by polyelectrolyte effects. The excluded volume effect is also the primary basis for the ability of water-soluble polymers, such as polyethylene glycol, to **precipitate** proteins. [Preferential Hydration](#) is partly due to steric exclusion, in that a larger cosolvent is more effectively excluded from protein surface than water (Fig. 1b).

The excluded volume effect favors any chemical reaction in which the volume decreases. This includes **ligand binding** by a protein and the association of protein subunits to form [oligomers](#). The equilibrium constant for such a reaction is increased with the increasing concentration of macromolecules in the solution. At high concentrations, such as occur normally in the **cytosol** of cells, such effects can be very large (2).

#### Bibliography

1. T. Arakawa and S. N. Timasheff (1985) *Biochemistry* **24**, 6756–6762.
2. A. P. Minton (1981) *Biopolymers* **20**, 2093–2120.

## Exocytosis

Exocytosis occurs in eukaryotic cells when the [membrane](#) of a cytoplasmic **vesicle** fuses with the plasma membrane, exposing the contents of the vesicle and the luminal layer of the phospholipid bilayer to the outside world.

### 1. Examples of Exocytosis

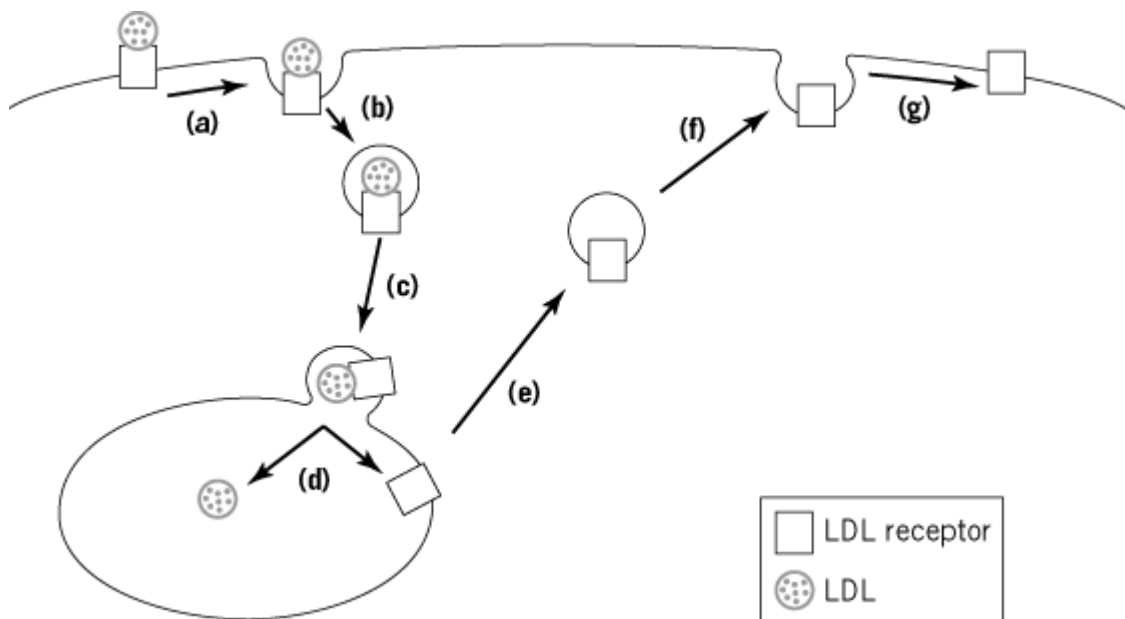
Most cells can secrete proteins (see [Protein Secretion](#)). In bacteria, the proteins can be unfolded and then threaded, or translocated, across the plasma membrane. Eukaryotes translocate newly synthesized proteins into an intracellular organelle, the [endoplasmic reticulum](#). Proteins travel from the endoplasmic reticulum by way of the **Golgi complex**, to the plasma membrane. Transport is effected by packaging the newly synthesized proteins in [secretory vesicles](#), which must fuse their membrane with that of the plasma membrane to release their contents to the extracellular space. Thus eukaryotic cells that secrete proteins must have the molecular machinery to allow exocytosis at the cell surface. Exocytosis must be highly efficient, since few secretory vesicles are seen in most cells.

Sometimes newly synthesized proteins passing through the Golgi complex are diverted selectively into cytoplasmic storage vesicles, usually called dense core secretory granules ([1](#)). To allow cytoplasmic storage, exocytosis of such secretory granules is inhibited, causing the accumulation of granules in the cytoplasm or along the plasma membrane. Accumulation of storage or secretory granules usually gives cells a characteristic morphology; examples are the granulocytes of the **hematopoietic** system, endocrine and exocrine cells, prefertilization eggs filled with cortical granules, and protozoa such as *Paramecium* and *Tetrahymena*. A signal to the outside of such protein-storage cells can remove the inhibition to exocytosis and trigger a massive efflux of protein, a process referred to as regulated secretion.

Some endocrine cells, such as chromaffin cells or neurons, actively transport small molecules from the cytoplasm into the lumen of the dense core secretory granules, giving them high internal concentrations of both protein and the small molecules ([2](#)). The energy for the active transport comes from a [proton gradient](#) across the secretory granule membrane ([3](#)). This property is also found in bone marrow-derived cells. Mast cells, for example, package histamine to high concentrations inside their secretory granules. When cells of this type undergo exocytosis, along with the stored proteins they release massive quantities of the small molecules, which usually are identical, or are related to, neurotransmitters.

In addition to having a biosynthetic pathway for protein secretion, many eukaryotic cells have an [endocytosis](#) pathway that internalizes plasma membrane components and returns most of them to the cell surface ([4](#)). This membrane recycling pathway allows a cell to move membrane to a site where it might be more needed; for example, a region of cell growth or migration ([5](#)). The recycling pathway also allows a cell to internalize nutrients bound to receptors, to remove the nutrients, and to send the receptors back to the cell surface for more. In the final step, a membrane vesicle on the recycling pathway must fuse with the plasma membrane (that is, it must undergo exocytosis). Thus, the exocytosis machinery is essential for endocytotic recycling, as well as for protein secretion (Fig. [1](#)).

**Figure 1.** Example of exocytosis in endocytosis. Exocytosis is used in other processes besides delivery of secretory proteins to the cell surface. Here, using the trafficking of low-density lipoprotein (LDL) receptor as an example, we see how exocytosis is an integral part of endocytic trafficking. (a) Cholesterol-containing LDLs bind to their receptors and (b) are internalized by endocytosis. (c) The endocytic vesicles fuse with an endosome (d) where LDL dissociates from its receptor. While LDL is eventually broken down for utilization of the cholesterol (ef), the LDL receptor is recycled back to the cell surface by exocytosis (g) so that it can be reused.



As with protein secretion, there appear to be constitutive and regulated exocytotic steps associated with membrane recycling. Membrane vesicles carrying nutrient receptors, such as the **receptors** for [transferrin](#) and low-density **lipoprotein**, fuse rapidly with the plasma membrane, independently of an external signal. In specialized cases, exocytosis is inhibited, causing an accumulation of endocytotically derived membrane vesicles (6). Two examples of such storage vesicles are synaptic vesicles and vesicles storing intracellularly a special type of glucose transporter, GLUT4 (glucose transporter 4) (7). Regulated exocytosis of the endocytotically derived vesicles can occur when the inhibition is removed, by an electrical signal in the case of the nerve cell, and by [insulin](#) in the case of GLUT4-storing fat and muscle cells.

Fortunately, similar molecular machinery seems to be used by all forms of exocytosis, regulated or constitutive, biosynthetic, or plasma-membrane recycling. Furthermore, the molecules that regulate exocytosis appear to be conserved, as are the transporters that put neurotransmitters into dense core secretory granules and into synaptic vesicles.

## 2. Measurements of Exocytosis

A common but crude measure of exocytosis is the appearance in the extracellular medium of the contents of a secretory vesicle that has undergone exocytosis. Thus exocytosis of synaptic vesicles (see [Secretory Vesicles/Granules](#)) at the nerve terminal causes release of neurotransmitter, and exocytosis of secretory granules from b-cells of the pancreatic islets of Langerhans releases insulin. Extracellular neurotransmitters and hormones can be measured biochemically, as can the constitutive secretion of enzymes such as invertase by yeast or [proteinases](#) by mammalian cells. Biochemical measurements are usually incapable of measuring the kinetics of exocytosis with accuracy.

The release of neurotransmitters, however, can be measured electrophysiologically with much greater precision. Most neurotransmitters activate the postsynaptic target cell by binding to closed **ion channels** and causing them to open. Measurement of current flow through the postsynaptic ion channels gives an almost instantaneous assay of the exocytosis rates in the presynaptic nerve terminal. The precision and sensitivity of electrophysiological measurements, pioneered by Dr. Bernhard Katz and his coworkers at University College, London, allowed detection of a quantum of neurotransmitter release, lasting a millisecond or so, generated by exocytosis of the contents of a single synaptic vesicle. Electrophysiology revealed that exocytosis is very rapid, occurring within 100  $\mu$ s of a stimulus to the nerve terminal, and that the concentration of the neurotransmitter in the

synaptic vesicle is very high, greater than 100 mM.

Postsynaptic currents generated by presynaptic release of neurotransmitter accurately measured the exocytosis because the receptor channels were 10 to 20 nm from the exocytosis site. Release of biogenic amines from cells such as chromaffin cells can likewise be measured rapidly and quantitatively by placing a carbon electrode close to the cell and detecting the biogenic amine by its **redox** potential (8). This technique was first applied to exocytosis by Wightman's group and is called amperometry.

Another technique widely used to study exocytosis uses the cell-surface capacitance. The electrical capacitance across the plasma membrane of a cell is proportional to its area. When exocytosis occurs, the capacitance increases because of the addition of membrane (9, 10). The sensitivity of capacitance measurements is such that it can measure the addition of the membrane of a single secretory granule. It cannot yet detect the addition of a single synaptic vesicle, since they are too small (50 nm diameter). When the capacitance signal is big enough, however, the capacitance changes give an instantaneous measurement of exocytosis.

To measure capacitance, a cell must be attached to a "patch electrode," which allows electrical communication between the salt solution inside the electrode and the cytoplasm of the cell.

Video techniques for quantifying fluorescent markers in cells and recording changes in **fluorescence** intensity as a function of time are becoming increasingly sophisticated. Video-microscopy techniques are used to study exocytosis by loading synaptic vesicles with a membrane-impermeable fluorescent dye (11, 12). Each synaptic vesicle can contain about 30 dye molecules. When the membrane of a synaptic vesicle fuses with that of the plasma membrane, the intensity of the fluorescence decreases by a fixed amount, which can be measured using sensitive detection devices. The dye technique for studying exocytosis was pioneered by William Betz and his colleagues at the University of Colorado.

When the secretory vesicle is big enough, it is possible to observe exocytosis in the light microscope. Some protozoa, such as *Paramecium* or *Tetrahymena*, have large cytoplasmic vesicles, packed with protein, that are docked at the plasma membrane. When the light microscope is used (13), the dense contents of the secretory granules can be seen to disappear when exocytosis is triggered. The convenience of a visual assay for exocytosis has made it possible to screen for secretion mutants in organisms of this type. Exocytosis of secretory granule content can also be detected visually in sea urchin eggs and in some classes of mast cell.

### 3. Mechanism of Exocytosis

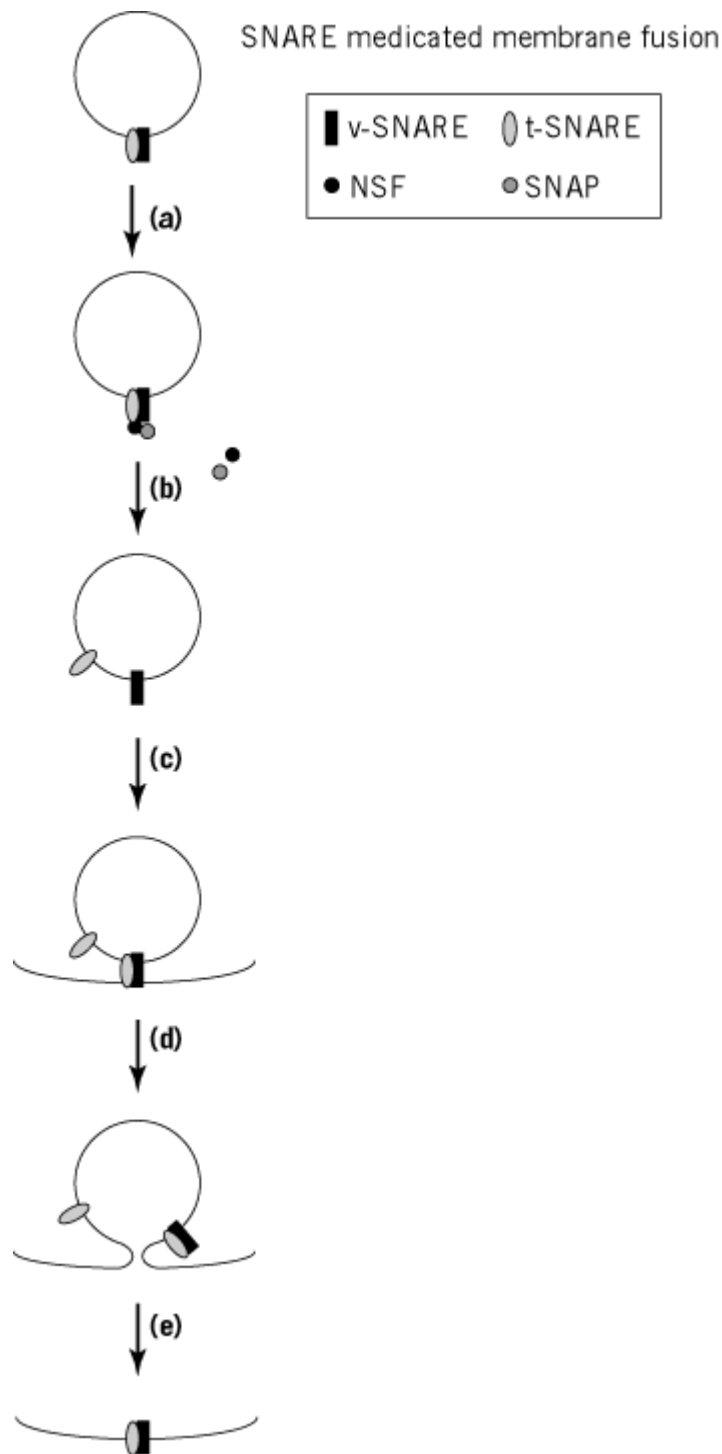
Exocytosis takes place in several stages. First the membrane of the vesicle recognizes the plasma membrane and forms a close association with it. Then the phospholipid bilayers are rearranged so that the membranes of the vesicle and those of the plasma membrane become continuous. Finally, the curved vesicle membrane flattens out, adopting the curvature of the plasma membrane.

Association of the cytoplasmic vesicle with the membrane is sometimes called docking. There are two types of vesicle docking, which sometimes get confused. Since exocytosis is not a random event but requires docking specifically between a vesicle and plasma membrane, the two membranes must recognize each other, presumably through **protein-protein interactions**. The recognition and fusion steps must be rapid processes in most cells, because few vesicles are seen adjacent to the plasma membrane or in the process of fusion. In contrast, in regulated secretory cells, such as neurons and endocrine cells, the secretory and synaptic vesicles (see **Secretory Vesicles/Granules**) are in a stable association with the plasma membrane. The long-lasting docking in neurons and neuroendocrine cells is presumably an adaptation that allows very rapid exocytosis in response to an extracellular signal. To avoid confusing the two forms of docking, we discuss first the generic, short-term docking, present in all eukaryotic cells.

Many of the molecular insights into exocytosis derive from mutations that affect protein secretion in yeast (see [Sec Mutants/Proteins](#)). Such mutants can have defects in all stages of the protein export pathway. Particularly relevant to exocytosis are the mutations that accumulate post-Golgi secretory vesicles in their cytoplasm. Of the three sets of gene products that affect exocytosis, the ones that are understood in most detail are the SNARE proteins (14). The name SNARE (SNAP REceptor) derives from one of the ways the proteins were discovered (discussed later) and is useful to describe their function. A v-SNARE protein on a vesicle targeted for exocytosis interacts with a t-SNARE on the target plasma membrane. The first identified member of the v-SNARE family was the synaptic vesicle protein initially called VAMP (Vesicle Associated Membrane Protein) (15) and later synaptobrevin (16). The characteristics of v-SNAREs on secretory vesicles include (1) a vestigial luminal or extracytoplasmic domain at the **C-terminus**, which can be as small as two amino acid residues; (2) a single conserved transmembrane domain; (3) a highly conserved juxtamembrane region that has the predicted sequence of a [coiled coil](#), and (4) a variable **N-terminal** domain. At the cell surface, the v-SNARE is believed to interact with two t-SNAREs. One is a member of the syntaxin family (17), which also has a small luminal C-terminal **domain**, a single transmembrane domain, and a coiled-coil region. The third, a member of the SNAP (synaptosomal-associated protein)-25 family (18), is a cytoplasmic peripheral protein, also with a coiled-coil domain, and is held to the plasma membrane by covalently attached **palmitoyl groups**. The three proteins are believed to form a ternary complex via their helical domains, which facilitates fusion (19). All well-studied exocytotic events, from yeast to mammalian cells, also appear to involve close **homologues** of these three SNARE proteins (20). Even synaptic-vesicle exocytosis, where the fusing vesicle is endocytotically derived, uses SNARE proteins. This has been shown using mutations in *Drosophila* that affect the neuronal SNAREs VAMP/synaptobrevin or syntaxin (21, 22). The involvement of SNAREs can also be verified by use of bacterial [neurotoxins](#). Tetanus toxin and the botulinum neurotoxins A to G are [endopeptidases](#) that specifically hydrolyze the SNARE molecules (23). Exocytosis from nerve terminals and endocrine cells is blocked by these neurotoxins, providing evidence for the involvement of the SNARE complex even in regulated exocytosis (24).

SNARE complexes form readily, even with purified components in the presence of [detergents](#). Given their high proclivity for spontaneous interaction, it is no surprise that the SNAREs appear to be sequestered in complexes that limit their reactivity. The v-SNAREs in synaptic vesicles (VAMP/synaptobrevin) is sequestered by the synaptic-vesicle protein synaptophysin, which has four transmembrane segments (25). The corresponding t-SNARE, syntaxin, is strongly bound to a cytosolic protein, which is called the *sec1* gene product in yeast and nsec1, for neuronal sec1, in the nervous system (26, 27). Since SNARE complexes are so energetically stable, energy must be expended to separate the SNAREs. Separation of the complex involves the [ATPase NSF](#) (*N*-ethyl maleimide sensitive factor) and two accessory proteins, called a- and b-SNAP (28). For NSF to bind, the SNAPs must bind the SNARE complexes. Thus the SNARE complexes received their name, SNAP receptors. The acronym SNAP stands for soluble NSF-acceptor protein in this case. For the t-SNARE SNAP-25, the acronym denotes synaptosomal-associated protein. The identity of the acronyms is coincidental but confusing. Currently, it is thought that NSF might act to dissociate the SNARE complex before fusion, but this is controversial (29, 30) (Fig. 2).

**Figure 2.** SNARE-mediated membrane fusion. (a) NSF and SNAP proteins bind to v-SNARE/t-SNARE complexes on membranes. (b) ATP hydrolysis by NSF causes dissociation of the SNARE complex so that (c) the v-SNARE can form a complex with a t-SNARE on a target membrane. (d) By an unknown mechanism complexes formed with SNAREs on opposing membranes are thought to promote membrane fusion (e), ending with the SNARE complex on the same membrane. NSF/SNAP must dissociate the complex again before either of the SNAREs can be reused.



Exocytosis in yeast is also inhibited by a mutation in a small, ras-like, [GTPase](#), the *sec4* gene product ([31](#)). The *sec4* protein is a member of a large class of small GTPases, the rab proteins. All secretory vesicles have a rab protein. In the case of synaptic vesicles, it is *rab3A*. Rab proteins appear to interact with a complex of proteins, of which the rabaptin complex is the prototype ([32](#)). The *rab* system seems to parallel the SNARE system in docking secretory vesicles at the plasma membrane. Thus, yeast mutations in the *rab* pathway can be rescued by overexpression of genes in the SNARE pathway ([24](#)). Another protein complex that appears to determine the site of exocytosis in yeast and mammalian cells is the exocyst ([33](#), [34](#)).

The final class of players in exocytosis appears to be related to [actin-binding proteins](#), because

mutations in such proteins block protein secretion (see [Sec Mutants/Proteins.](#)) The link between membrane traffic and the actin [cytoskeleton](#) is currently obscure.

#### 4. Regulated Exocytosis

Regulated exocytosis is found in cells that have a large cytoplasmic pool of secretory vesicles. In the presence of a suitable external stimulus, the probability of exocytosis per secretory vesicle increases manifold. In neurons and endocrine cells, the external stimulus commonly triggers in regions of the cytoplasm near **calcium channels** a rise of intracellular calcium ions from its normal low level of about  $0.1 \mu\text{M}$  to levels that can range up to  $100 \mu\text{M}$  and greater (35). A common conjecture is that calcium ions bind to an inhibitor of exocytosis, removing the inhibition. Good candidates for calcium sensors in neurons are members of the synaptotagmin family (36). Synaptotagmins have two **calcium-binding** domains on their cytoplasmic regions and can oligomerize. Mutations in synaptotagmin cause loss of calcium-dependent exocytosis (37). Despite the conjecture that the calcium sensor is inhibiting exocytosis, its absence in mutant organisms does not lead to high levels of exocytosis in the absence of stimulation.

Another calcium-binding protein implicated in regulated exocytosis is the CAPS protein (38). Its addition to cell-free preparations of secretory granules docked at plasma membranes restores calcium-dependent release. The CAPS protein binds to the negatively charged phospholipids, the phosphorylated forms of [phosphatidyl inositol](#). Consistent with CAPS having a role in exocytosis is the finding that secretory granules must be primed with adenine triphosphate (ATP) in cell-free systems before exocytosis can occur (39). The ATP is required to phosphorylate the 4' and 5' positions on the inositol ring, catalyzed by specific phosphatidylinositol kinases. The phosphorylation state of phosphatidyl inositol is found to determine the efficiency of several membrane trafficking steps (40), but the reason is not yet clear.

The inside surface of the plasma membrane is frequently coated with a layer of cross-linked [actin](#) filaments. To allow exocytotic vesicles access to the plasma membrane, it may be necessary to remove or rearrange the subcortical actin layer. Many reports of such rearrangements have been recorded for stimulated endocrine cells. Solubilization of the actin array may be facilitated by an actin-severing protein (41).

#### 5. The Fusion Pore

Electrophysiological measurements of exocytosis have predicted that the merging of secretory-vesicle membrane with plasma membrane is preceded by the formation of a channel, the so-called fusion pore (42). One line of evidence for such a pore comes from amperometry. When the release of a biogenic amine from a secretory granule is measured, the massive efflux of granule content is preceded by a small "step" of amine release like that predicted if the inside of the secretory granule were connected to the outside by a small channel or hole of nanometer diameter (43). A second indication of the reversible formation of a channel is "capacitance flicker." Exocytosis can be detected as a sudden increase in membrane capacitance; occasionally, however, a sudden increase is followed by a sudden decrease of exactly the same amplitude. These data are consistent with a transient and reversible fusion of secretory granule and plasma membranes. A third line of evidence supporting the idea of a fusion pore is that the current that flows through the fusion pore as the secretory-granule membrane discharges its membrane potential. From measurements of the conductance of the channel, estimates of its size can be made. When it first forms, the fusion pore seems to form a channel of about the size of a large ion channel.

The molecular nature of the fusion pore, if it exists, is unclear. Fusion between a membrane-covered virus and the plasma membrane has been studied extensively for [influenza virus](#). The viral protein triggers fusion by undergoing a conformational change that extends a 16- to 23-amino acid residue "fusion peptide," a **hydrophobic** sequence that is essential for fusion (44). Extension of the fusion peptide requires a helical hairpin to be converted to a triple-stranded [coiled-coil](#) domain (45). The



three SNARE proteins, VAMP/synaptobrevin, syntaxin, and SNAP-25, have coiled-coil domains that initiate oligomer formation. VAMP also has a small [alpha-helical](#) domain that resembles a fusion peptide. There is no experimental evidence, however, that the SNAREs form a fusion pore. Nor is it clear whether the initial hole is exclusively due to rearrangement of the lipid bilayer or to a proteinaceous channel. In influenza virus-mediated fusion, a hemifusion state is formed if the transmembrane tail of the fusion protein is replaced by a glycosylphosphatidylinositol anchor (46). This was interpreted to mean that the hemifusion state is an intermediate in virus-induced fusion.

Logically, membrane fusion must be initiated by the formation of one or several small holes through the membrane of both the exocytotic vesicle and the cell surface. Electrophysiologically, the evidence for the reversible formation of such a hole is excellent. As yet, molecular knowledge lags far behind the physiology.

## Bibliography

1. T. L. Burgess and R. B. Kelly (1987) *Annu. Rev. Cell Biol.* **3**, 243–293.
2. H. Winkler (1993) *J. Anat.* **183**, 237–252.
3. M. E. Finbow and M. A. Harrison (1997) *Biochem. J.* **324**, 697–712.
4. I. Mellman (1996) *Annu. Rev. Cell Develop. Biol.* **12**, 575–625.
5. C. R. Hopkins, A. Gibson, M. Shipman, D. K. Strickland, and I. S. Trowbridge (1994) *J. Cell Biol.* **125**, 1265–1274.
6. O. Cremona and P. De Camilli (1997) *Curr. Opin. Neurobiol.* **7**, 323–330.
7. S. Rea and D. E. James (1997) *Diabetes* **46**, 1667–1677.
8. J. M. Finnegan, K. Pihel, P. S. Cahill, L. Huang, S. E. Zerby, A. G. Ewing, R. T. Kennedy, and R. M. Wightman (1996) *J. Neurochem.* **66**, 1914–1923.
9. E. Neher and A. Marty (1982) *Proc. Natl Acad. Sci. USA* **79**, 6712–6716.
10. J. M. Fernandez, E. Neher, and B. D. Gomperts (1984) *Nature* **312**, 453–455.
11. W. J. Betz and G. S. Bewick (1993) *J. Physiol.* **460**, 287–309.
12. J. K. Angleson and W. J. Betz (1997) *Trends Neurosci.* **20**, 281–287.
13. J. C. Hutton (1997) *Proc. Natl. Acad. Sci. USA* **94**, 10490–10492.
14. J. E. Rothman (1996) *Protein Sci.* **5**, 185–194.
15. W. S. Trimble, D. M. Cowan, and R. H. Scheller (1988) *Proc. Natl. Acad. Sci. USA* **85**, 4538–4542.
16. M. Baumert, P. R. Maycox, F. Navone, P. De Camilli, and R. Jahn (1989) *EMBO J.* **8**, 379–384.
17. M. K. Bennett, J. E. Garcia-Araras, L. A. Elferink, K. Peterson, A. M. Fleming, C. D. Hazuka, and R. H. Scheller (1993) *Cell* **74**, 863–873.
18. G. A. Oyler, G. A. Higgins, R. A. Hart, E. Battenberg, M. Billingsley, F. E. Bloom, and M. C. Wilson (1989) *J. Cell Biol.* **109**, 3039–3052.
19. T. Sollner, M. K. Bennett, S. W. Whiteheart, R. H. Scheller, and J. E. Rothman (1993) *Cell* **75**, 409–418.
20. M. K. Bennett and R. H. Scheller (1993) *Proc. Natl. Acad. Sci. USA* **90**, 2559–2563.
21. K. L. Schulze, K. Broadie, M. S. Perin, and H. J. Bellen (1995) *Cell* **80**, 311–320.
22. K. Broadie, A. Prokop, H. J. Bellen, C. J. O’Kane, K. L. Schulze, and S. T. Sweeney (1995) *Neuron* **15**, 663–673.
23. G. Schiavo, O. Rossetto, and C. Montecucco (1994) *Sem. Cell Biol.* **5**, 221–229.
24. S. R. Pfeffer (1996) *Annu. Rev. Cell Develop. Biol.* **12**, 441–461.
25. L. Edelmann, P. I. Hanson, E. R. Chapman, and R. Jahn (1995) *EMBO J.* **14**, 224–231.
26. J. Pevsner, S. C. Hsu, and R. H. Scheller (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1445–1449.
27. E. P. Garcia, E. Gatti, M. Butler, J. Burton, and P. De Camilli (1994) *Proc. Natl. Acad. Sci.*

USA **91**, 2003–2007.

28. J. Hay and R. H. Scheller (1997) *Curr. Opin. Cell Biol.* **9**, 505–512.
29. A. Banerjee, V. A. Barry, B. R. DasGupta, and T. F. J. Martin (1996) *J. Biol. Chem.* **271**, 20223–20226.
30. C. Ungermann, B. J. Nichols, H. R. Pelham, and W. Wickner (1998) *J. Cell Biol.* **140**, 61–69.
31. A. Salminen and P. J. Novick (1987) *Cell* **49**, 527–538.
32. H. Stenmark, G. Vitale, O. Ullrich, and M. Zerial (1995) *Cell* **83**, 423–432.
33. D. R. TerBush, T. Maurice, D. Roth, and P. Novick (1996) *EMBO J.* **15**, 6483–6494.
34. Y. Kee, J. S. Yoo, C. D. Hazuka, K. E. Peterson, S. C. Hsu, and R. H. Scheller (1997) *Proc. Nat. Acad. Sci. USA* **94**, 14438–14443.
35. R. Llinas, M. Sugimori, and R. B. Silver (1995) *Neuropharmacology* **34**, 1443–1451.
36. S. V. Popov and M. M. Poo (1993) *Cell* **73**, 1247–1249.
37. M. Geppert, Y. Goda, R. E. Hammer, C. Li, T. W. Rosahl, C. F. Stevens, and T. C. Sudhof (1994) *Cell* **79**, 717–727.
38. K. Ann, J. A. Kowalchuk, K. M. Loyet, and T. F. Martin (1997) *J. Biol. Chem.* **272**, 19637–19640.
39. J. C. Hay, P. L. Fiset, G. H. Jenkins, K. Fukami, T. Takenawa, R. A. Anderson, and T. F. Martin (1995) *Nature* **374**, 173–177.
40. P. De Camilli, S. D. Emr, P. S. McPherson, and P. Novick (1996) *Science* **271**, 1533–1539.
41. L. D. Hernandez, L. R. Hoffman, T. G. Wolfsberg, and J. M. White (1996) *Annu. Rev. Cell Develop. Biol.* **12**, 627–661.
42. S. Muallem, K. Kwiatkowska, X. Xu, and H. L. Yin (1995) *J. Cell Biol.* **128**, 589–598.
43. M. Lindau and W. Almers (1995) *Curr. Opin. Cell Biol.* **7**, 509–517.
44. A. Albillos, G. Dernick, H. Horstmann, W. Almers, G. Alvarez de Toledo, and M. Lindau (1997) *Nature* **389**, 509–512.
45. C. M. Carr and P. S. Kim (1993) *Cell* **73**, 823–832.
46. G. W. Kemble, T. Danieli, and J. M. White (1994) *Cell* **76**, 383–391.

## Exon Shuffling

When intervening sequences were discovered to split eukaryotic **genes** into segments coding for protein, the intervening sequences were named *introns* and the coding regions, exons ([1](#)) (see **Introns/exons**). Since then, heated discussion has ensued concerning the evolutionary origin and biological significance of introns in eukaryotes, because almost all of the prokaryotic genes lack introns, and the introns have no known role.

Gilbert ([1](#)) proposed a scenario in which introns simply connect neighboring exons to each other in the **genome**; consequently, genetic **recombination** would not harm the coding portions of genes if it took place within the intron regions. Such recombination within introns would facilitate the shuffling of exons, to create new exon combinations and lead to the emergence of new genes with new functions. This is the “exon shuffling” theory (see also [Gene Splicing](#)).

Exon shuffling is related to the present controversy regarding the “early” versus “late” intron theories. The early-intron theory maintains that introns existed in the ancestor genomes of

prokaryotes and eukaryotes, but that all those in prokaryotes were deleted by some unknown mechanism. On the other hand, the late-intron theory contends that introns could have been inserted, similar to [transposable elements](#), into eukaryotic genomes quite recently, at least after the **divergence** of eukaryotes and prokaryotes. This theory is based on the observation that there are many genes in which the locations of introns are not conserved among vertebrates, invertebrates, and plants (2).

If the early-intron theory is correct, exon shuffling probably had a significant role in the creation of new genes. If the late-intron theory is right, however, exon shuffling might have been significant only after the introns were inserted. The controversy remains to be resolved.

The **domain** shuffling theory is similar to the exon shuffling theory, but in sharp contrast, maintains that the unit of shuffling during evolution is a functional protein domain, not an exon. The functional domain does not necessarily correspond to the exon. Although there are some cases where it does, most functional domains consist of more than one exon or occupy only part of a large exon. For example, the [Kringle domain](#) is known to have been shuffled as a unit during evolution. The Kringle domain is a characteristic [supersecondary structure](#) frequently found in the [serine proteinases](#) involved in the [blood clotting](#) system. Even though the Kringle domain is, in most cases, split into three exons by two introns, those found in various proteins are almost always complete forms and never found as parts corresponding to the exons. At present, however, there is no known genetic mechanism for facilitating domain shuffling, and domain shuffling seems to have taken place in prokaryotes as well (3).

### Bibliography

1. W. Gilbert (1978) Nature **271**, 501.
2. J. D. Palmer and J. M. Logsdin (1991) Curr. Opin. Gen. Devel. **1**, 565–570.
3. T. Gojobori and K. Ikeo (1994) Phil. Trans. Roy. Soc. Lond. B **344**, 411–415.

### Expansibility

The expansion of a protein molecule with increasing temperature reflects its atomic packing and flexibility, as does its [compressibility](#) with increasing pressure. [X-ray crystallography](#) studies indicate that both the crystal lattice and protein structure expand with increasing temperature. For myoglobin, a volume expansion of 5% of the crystal lattice and 3% in the protein molecule occurred as the temperature increased from 80 to 300 K, which corresponds to a thermal expansion coefficient of  $1.15 \times 10^{-4} \text{ K}^{-1}$  (1). **Ribonuclease A** exhibited a smaller expansion 0.9% for a temperature increase from 98 to 320 K (2). It is likely that a more compressible protein will also be more thermally expansible. It is believed that the observed expansions of proteins with temperature are due to motions of secondary structure and exposed surface loops and to a decrease in local atomic packing density (see [van der Waals Surface, Volume](#)). Thermal expansion of a protein is not uniform over the molecule, just as in the case of contraction by pressure.

It is known that the apparent and partial specific volumes ( $v_2$ ) of proteins in aqueous solution are linearly dependent on temperature in the range of 4 to 45°C. Positive values for  $d v_2/dT$  of  $(2.5 \text{ to } 10 \times 10^{-4} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1})$  have been reported, with the majority lying between  $(3.5 \text{ to } 5) \times 10^{-4} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1}$ . These values are greater than the expansion coefficient found for the X-ray crystal structure.

If the thermal expansion of crystal structure is caused by the cavities, the positive values of  $dv_2/dT$  for folded proteins would be dominantly attributed to the hydration effect, ie, a volume increase due to dehydration of protein molecules on elevating the temperature (see [Partial Specific \(Or Molar\) Volume](#)). The temperature effect of  $v_2$  is still pronounced for [amino acids](#) and small **peptides**, probably due to a large dehydration effect:  $dv_2/dT$  is  $1.3 \times 10^{-3} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1}$  for glycine and  $7.3 \times 10^{-4} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1}$  for triglycine. In this sense, an [unfolded protein](#) is expected to have a larger  $dv_2/dT$  value than a folded one because the accessible surface area (amount of hydration) increases on denaturation. Hawley reported that  $dv_2/dT$  increases by  $0.18 \times 10^{-4} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1}$  and  $0.51 \times 10^{-4} \text{ cm}^3 \text{ g}^{-1} \text{ K}^{-1}$  on pressure denaturation of ribonuclease A (pH 2.0) and chymotrypsinogen (pH 2.07), respectively (3). However, the temperature effects on  $v_2$  at high pressures and high concentration of denaturants, above where conformational changes or denaturation take place, are complex since the volume change of hydration  $Dv_{\text{sol}}$  may involve the contributions of changes in water structure and preferential solvent interactions.

### Bibliography

1. H. Frauenfelder et al. (1987) *Biochemistry* **26**, 254–261.
2. R. F. Tilton Jr., J. C. Dewan, and G. A. Petsko (1992) *Biochemistry* **31**, 2469–2481.
3. S. A. Hawley (1971) *Biochemistry* **10**, 2436–2442.

### Expressed Sequence Tag

An expressed sequence tag (EST) is a short stretch of **DNA** sequence that is used to identify an expressed **gene**. Although EST sequences are usually only 200 to 500 nucleotides in length, this is generally sufficient to identify the full-length **complementary DNA** (cDNA). ESTs are generated by sequencing a single segment of random **clones** from a **cDNA library**. A single sequencing reaction and automation of DNA isolation, sequencing, and analysis have allowed the rapid determination of many ESTs. Now, the majority of the sequences in sequence [databases](#) are ESTs. Although most ESTs have been isolated from humans, a large number of ESTs have been isolated from model organisms, such as *Caenorhabditis elegans* (1), *Drosophila*, rice (2), and *Arabidopsis* (3). ESTs are also being isolated from more exotic organisms, such as *Entamoeba histolytica* (4) and *Leishmania major promastigotes* (5). ESTs have numerous uses, from **genetic mapping** to analyzing gene expression, and the number of ESTs isolated from different organisms will continue to rise rapidly.

#### 1. Generation of ESTs

The most important step in generating of ESTs is producing a cDNA library. First, [messenger RNA](#) (mRNA) is extracted from the material being studied and is used as a template by **reverse transcriptase** for cDNA synthesis. Then the DNA is **cloned** into a suitable **vector** to produce a cDNA library. Random clones are isolated from the library, and one or both ends are sequenced by single-pass sequencing.

The source of the mRNA is an organism, tissue, or [cell line](#) grown under normal conditions or treated with [hormones](#), drugs, heat, etc. The source of the starting material depends on the type of genes to be identified and the reason for generating the ESTs. EST analysis is used to find previously undiscovered tissue-specific **genes** or to tag as many different expressed genes as possible. It is also

used to examine the relative abundance of expressed genes. Libraries used for the first two types of studies are often normalized, so that highly-expressed genes and rare genes are represented more equally in the library or subtracted to reduce or eliminate the number of highly abundant clones (see [Subtractive Hybridization](#)). An alternative to normalizing or subtracting libraries is to sequence a small number of ESTs from as many different tissue types and treatments as possible because different types of genes are highly expressed under different circumstances.

It is useful if the cDNAs in the library are cloned directionally, so that ESTs are isolated specifically from either the 5'- or the 3'-end of the cDNA. 3'-ESTs often represent the 3'-untranslated region of the mRNA and are used to separate members of gene families that have similar coding sequences, whereas 5'-ESTs generally represent coding sequence and give a better idea of the type of gene being expressed.

If the starting material for mRNA isolation is limiting, cDNA minilibraries suitable for EST analysis are generated by arbitrarily primed RT-PCR (6).

## 2. The Multitudinous Uses of ESTs

### 2.1. Obtaining Full-Length cDNAs from ESTs

An EST is useful because it represents an expressed gene. Once an EST has been identified, it is usually straightforward to obtain the full sequence of the “tagged” cDNA. Some sources provide EST DNA that is used in further studies. Alternatively, the EST sequence is used to design PCR primers or hybridization probes for cloning the full-length cDNA (see [Cloning](#)). Another strategy is “virtual cloning”, where computer analysis is used to array **homologous** ESTs and any other sequences from the same gene (7). Once isolated, then the full-length cDNA is used in [mutagenesis](#), **transgenic**, and expression studies to analyze the gene function.

### 2.2. Genome Analysis and Mapping with ESTs

Many organisms, for example, humans, *Arabidopsis*, and *C. elegans*, will have their entire genomic DNA sequenced in the relatively near future, and certain genomes, such those of *Saccharomyces cerevisiae* and various **bacteria**, have already been entirely sequenced. In the meantime, EST analysis provides a rapid way to identify expressed genes. Now, it is not possible to predict coding sequences in genomic DNA reliably from sequence information alone. Even when complete genomic sequences are available, ESTs are useful in analyzing the genomic sequence, for example by verifying putative coding sequences and confirming **intron** and exon boundaries. They are also helpful in distinguishing [pseudogenes](#) from real genes and in identifying **alternatively spliced** transcripts, which could never be predicted from the genomic sequence alone.

ESTs are also used to provide markers for genomic mapping by converting them into “sequence tagged sites” (STS). An STS is a set of PCR primers that identifies a single gene, and STS primers are designed from the EST sequence (8). An STS can be mapped to a specific region of the genome, and if DNA contigs are available for the genome under study, the STS can specify a genomic clone. Mapping of ESTs to genomic clones will provide a map of expressed genes for each clone and, in humans, these maps can be used to identify candidate genes for inherited human diseases that have been mapped to chromosomal regions by pedigree analysis.

### 2.3. Gene Expression Studies Using ESTs

As mentioned previously, gene expression in different organisms, tissues, or cells, before or after different treatments, can be studied by surveying ESTs isolated from libraries created from different mRNA sources. If nonnormalized and unsubtracted libraries are used to isolate ESTs, the frequency of isolation of a particular EST indicates the relative message abundance of the tagged gene (9). It is also possible to isolate novel tissue-specific genes by merely searching EST databases (10).

### 2.4. Rapid Identification of New Genes and Gene Families

As of 1998, the majority of ESTs in the databases represent genes that have not been previously

identified. Traditionally, gene sequences have been obtained one at a time, first by isolating a protein, by identifying a mutant and cloning the mutated gene by [complementation](#), or by positional cloning (see [Cloning](#)). Compared with these methods, generating ESTs is a very rapid method for identifying new gene sequences.

EST database searches also provide a rapid method for finding genes similar to a gene of interest. Methods, such as PCR, using degenerate primers, or **hybridization** at low stringency also enable the identification of similar genes, but these methods are often not as successful as database searches. Databases can be searched by using either nucleic acid sequences or protein sequences. Protein comparison searches are often more fruitful because protein sequences are better conserved. Additionally, databases can be used to search for similar genes among all the represented organisms simultaneously, thus facilitating comparative studies. Once similar genes have been identified, they can be translated into protein sequences and the database searched again to find the sequences most similar to this second set of genes. Thus it is possible to identify large families of related genes rapidly by EST database analysis.

### 2.5. Using ESTs to Clone Specific Genes

In humans, one main goal of genomic research and EST analysis is to aid in identifying and cloning human genes related to diseases. Various groups have begun to look between species to find candidate human disease genes. One group used an EST database search to identify human cDNAs homologous to previously cloned *Drosophila* genes that have an interesting developmental mutant **phenotype** (11). Then the human sequences were mapped by several methods, and the positions of the cDNAs compared with the map positions of human disease genes. Some of the cDNAs mapped to regions containing human disease genes that cause symptoms similar to the defects in the corresponding *Drosophila* mutant genes. Thus these cDNAs are candidates for the disease genes, and further studies can be undertaken to determine whether the cDNAs and the disease genes are one and the same. Another group is systematically identifying novel human ESTs related to known genes in a variety of model organisms (12). Then the ESTs are mapped to both human and [mouse](#) maps. Again, such genes can provide candidates for human disease genes, thereby potentially speeding up their cloning and aiding in their analysis by providing a model system for further studies.

ESTs are extremely useful tools for gene analysis in organisms besides humans, although to date most effort has been spent identifying human ESTs. In any organism where a mutation has been mapped to a certain genomic region or to a genomic clone, ESTs mapping to the same region identify candidate genes. In more unusual organisms, where few genes have been previously isolated, EST analysis provides a rapidly generated survey of the types of genes expressed by the organism and is used, for example, to isolate novel genes from pathogenic organisms (4, 5). In **plants**, the weed *Arabidopsis* is often used as a model system for gene isolation. EST database searches provide a rapid means of isolating similar genes from important crop plants, in much the same way that interesting genes from model systems are used to isolate corresponding genes from humans and mice. STSs from ESTs representing important plant genes could be used as markers in selective breeding programs.

### 3. Some Recent Twists on EST Generation

A technique termed *serial analysis of gene expression* (SAGE) has recently been developed to perform some of the functions of EST analysis (13). Short diagnostic sequence tags (9 to 10 bp) are randomly isolated from a tissue, concatenated, cloned, and sequenced. This method allows extremely rapid identification of thousands of genes and is used either to identify new genes or to analyze relative levels of gene expression. SAGE may prove quite useful in rapid surveys of gene expression differences between tissues or between developmental and disease states. ESTs still have many advantages over this method, however. The production and concatenation of the tags are somewhat cumbersome, and the small amount of sequence information in the tags precludes many of the interesting uses that have been found for ESTs.

An ingenious method has been developed to tag promoter-proximal sequences in mouse embryonic stem (ES) cells. Using a gene-trap retrovirus shuttle vector (14). A large number of ES cells were selected for neomycin resistance, indicating that the retrovirus containing a promoterless neomycin-resistance gene had inserted next to the promoter of an expressed gene. These neomycin-resistant ES cells were cloned. Some cloned cells were frozen and others used to extract DNA. The DNA flanking the retrovirus was isolated and sequenced to provide a promoter-proximal sequence tag (PST). PSTs are ESTs derived from genomic DNA rather than cDNA, and they are used similarly to screen sequence databases and to make STSs for mapping. However, they also have an additional advantage. Each PST represents a specific ES cell line harboring a potentially disrupted gene, and mutant ES cells are used to generate mutant mouse strains (15). PSTs thus allow rapid progression from sequence analysis of a gene to analysis of its function in a mutant mouse.

#### 4. The Future of ESTs

ESTs have already proved useful in rapidly identifying novel genes, providing markers for mapping, providing candidate disease genes, and for analyzing gene expression. Because the generation of ESTs proceeds faster than the functional identification and analysis of the genes that they represent, much future research will be directed toward determining functions for the newly identified genes and gene families. The cDNAs identified by EST analysis can be used in standard experiments designed to elucidate their cellular and developmental functions. As so many novel genes have been and are being identified, some researchers are looking for more rapid ways to analyze function or expression of large numbers of genes. A recent innovation is the use of cDNA (16) or oligonucleotide (17) microarrays that represent large numbers of genes or whole genomes. Hybridization to microarrays of labeled mRNA samples isolated from different tissues, mutant backgrounds, or developmental or growth states allows analysis of changes in gene expression for a large number of genes simultaneously. Other researchers have developed a system for monitoring the effects of gene disruption on growth in yeast for many genes simultaneously (18). The development of additional techniques for global analysis of gene expression and function will be essential for rapidly characterizing the wealth of new genes identified by EST studies.

#### Bibliography

1. W. R. McCombie et al. (1992) *Nature Genet.* **1**, 124–131.
2. K. Yamamoto and T. Sasaki (1997) *Plant Mol. Biol.* **35**, 135–144.
3. M. Delseny, R. Cooke, M. Raynal, and F. Grellet (1997) *FEBS Lett.* **405**, 129–132.
4. A. Azam, J. Paul, D. Sehgal, J. Prasad, S. Bhattacharya, and A. Bhattacharya (1996) *Gene* **181**, 113–116.
5. M. P. Levick, J. M. Blackwell, V. Conner, R. M. Coulson, A. Miles, H. E. Smith, K. L. Wan, and J. M. Ajioka (1996) *Mol. Biochem. Parasitol.* **76**, 345–348.
6. E. D. Neto, R. Harrop, R. Correa-Oliveira, R. A. Wilson, S. D. Pena, and A. J. Simpson (1997) *Gene* **186**, 135–142.
7. S. D. Rounsley, A. Glodek, G. Sutton, M. D. Adams, C. R. Somerville, J. C. Venter, and A. R. Kerlavage (1996) *Plant Physiol.* **112**, 1177–1183.
8. R. Berry et al. (1995) *Nature Genet.* **10**, 415–423.
9. K. Okubo, N. Hori, R. Matoba, T. Niiyama, A. Fukushima, Y. Kojima, and K. Matsubara (1992) *Nature Genet.* **22**, 173–179.
10. G. Vasmatazis, M. Essand, U. Brinkmann, B. Lee, and I. Pastan (1998) *Proc. Natl. Acad. Sci. USA* **95**, 300–304.
11. S. Banfi et al. (1996) *Nature Genet.* **13**, 167–174.
12. D. E. Bassett Jr., M. S. Boguski, F. Spencer, R. Reeves, S. Kim, T. Weaver, and P. Hieter (1997) *Nature Genet.* **15**, 339–334.
13. V. E. Velculescu, L. Zhang, B. Vogelstein, and K. W. Kinzler (1995) *Science* **270**, 484–487.

14. G. G. Hicks, E.-G. Shi, X.-M. Li, C.-H. Li, M. Pawlak, and H. E. Ruley (1997) *Nature Genet.* **16**, 338–343.
15. M. R. Capecchi (1989) *Science* **244**, 1288–1292.
16. M. Schena, D. Shalon, R. W. Davis, and P. O. Brown (1995) *Science* **270**, 467–470.
17. D. J. Lockhart, et al. (1996) *Nature Biotechnol.* **14**, 1675–1680.
18. V. Smith, K. N. Chou, D. Lashkari, D. Botstein, and P. O. Brown (1996) *Science* **274**, 2069–2074.

### Suggestions for Further Reading

19. D. Gerhold and C. T. Caskey (1996) It's the genes! EST access to human genome content, *BioEssays* **18**, 973–981. Excellent review of ESTs, databases and their use, and caveats.
20. E. F. Kirkness (1996) "Assessment of human gene diversity and expression patterns using expressed sequence tags. In" *Essays in Biochemistry* (D.K. Apps, ed.), Portland Press, London, U.K., pp. 1–9. Excellent clear review.
21. M. A. Marra, L. Hillier, and R.H. Waterston (1998) Expressed sequence tags-ESTablishing bridges between genomes. *Trends Genet.* **14**, 4–7. Good short review.

### Expression Libraries

An expression library is a collection of DNA molecules that can be expressed *en masse* through enzymatic methods to yield a population of DNA, RNA, or protein products that can be screened for individuals displaying desirable properties. The DNA sequences within the [library](#) can be of natural or synthetic origin. The advantage of screening expression libraries constructed from the [genome](#) or [complementary DNA](#) (see [Genomic Libraries](#) and [cDNA Libraries](#)) is that the expression products derived are usually biologically relevant. In contrast, synthetic DNA expression libraries may yield expression products that are not found in nature. However, there are two major advantages to screening synthetic DNA expression libraries. First, the expression products that are selected for a particular purpose may perform better than natural biopolymers. Second, a great deal of information is obtained by screening synthetic libraries of known composition, because the distribution of sequences selected compared with those not selected provides structure–activity relationship (SAR) data describing molecular interactions characteristic of the target being screened.

Expression libraries contain genomic or cDNA sequences cloned into expression vectors, which are plasmids or bacteriophage genomes that have been engineered to contain regulatory elements governing expression of foreign DNA inserts. For example, vectors containing a **promoter** sequence of **T7 phage** will initiate RNA [transcription](#) at a fixed position and direction. Once transcription is initiated, the **T7 RNA polymerase** will continue making an RNA copy of the DNA coding strand, including any foreign DNA that is encountered, in progressive fashion until a T7 terminator sequence is read (by convention, the coding strand is the top strand in a 5' to 3', left to right, double-stranded gene depiction). The resulting polymer that is formed is thus an RNA transcript of the DNA template. Such a DNA library can only be expressed in cells that also contain the enzyme T7 RNA polymerase. Typically, the gene encoding the T7 RNA polymerase is engineered to be inducible by the addition of a soluble sugar analogue, such as **IPTG**, so that library RNA expression can be co-induced at appropriate times in the screening process. RNA expressed from DNA libraries may be the intended chemical entity to be screened, or the RNA molecules may serve as templates for the [translation](#) of protein products. If the expressed RNA molecules contain suitable prokaryotic and/or



eukaryotic [ribosome](#) recognition sequences, ribosome-mediated **protein biosynthesis** will initiate at an AUG [start codon](#). The result will be a protein molecule encoded by the original DNA template. In general, it is easiest to screen protein expression libraries when the transcription and translation occur *in vivo* using a suitable host cell. However, in certain circumstances it may be advantageous to perform the protein expression *in vitro*, using cell-free translation extracts that contain the requisite protein synthesis machinery. In either case, the population of expressed proteins can be screened for any of a range of phenotypic properties, including [antibody](#) reactivity, affinity to a target molecule, or [catalysis](#).

The ability to engineer, at the DNA level, additional useful sequences into protein reading frames provides a powerful means of generating **fusion proteins** with desirable properties. Modern expression libraries contain many types of tags or fusion domains that facilitate **library** screening and downstream manipulation of individual clones. There are at least four types of protein fusion moieties commonly available. This list is not intended to be exhaustive but instead to illustrate the power of the technology. First, [epitope](#) tags are typically short amino acid segments fused to either the *N*- or *C*-terminus of the cloned protein. An epitope tag provides ready-made immunoreactivity to the cloned gene product through commercially available antibodies. Epitope tags are useful in both detection and purification of cloned gene products. Second, like epitope tags, affinity tags are amino acid sequences fused to either terminus of the cloned protein. The most common use of affinity tags is in affinity purification, although detection can also be performed through the tag. One commonly used affinity tag is the 6-His tag, which is simply a sequence of six **histidine** residues that bind with high affinity and specificity to nickel; therefore, 6-His fusion proteins can be purified from all other cellular proteins by [affinity chromatography](#) over nickel-agarose (1). Third, domain tags are independent units of structure and function that maintain or enhance certain activities of the proteins from which they were derived. Domain tags are useful for protein purification and detection. Common examples include **glutathione-S-transferase** (GST), which enables facile purification over glutathione-agarose resins (2), and **green fluorescent protein** (GFP), which is an autofluorescent protein that can be monitored in both quantitative and localization assays (3). Fourth, cleavage sequences are commonly introduced between the fusion tags and the cloned genes so that the tags can be removed following affinity purification. Cleavage elements are most commonly **proteinase** recognition sequences, although self-cleaving protein domains called inteins are becoming increasingly available (4) (see [Protein Splicing](#)).

A special type of cDNA expression library called the **two-hybrid system** (5) deserves special mention because it illustrates a powerful coupling of molecular engineering with combinatorial approaches to investigating protein function. The purpose of screening two hybrid libraries is to discover proteins that interact with a protein of interest. The system is termed “two-hybrid” because two separate plasmids have been engineered to express hybrid versions of transcriptional activators such as the yeast GAL4 protein. Normally, the GAL4 protein contains both a **DNA-binding** domain and a transcriptional activation domain. When intact GAL4 binds to its DNA recognition sequence, the transcriptional activation domain promotes transcription of adjacent gene sequences. In the two-hybrid system, the GAL4 DNA binding domain is expressed as a hybrid with the protein of interest, termed the “bait”. A second vector containing the GAL4 transcriptional activation domain is used to construct a library of cDNAs such that each cDNA is expressed as a hybrid protein fused to the GAL4 transcriptional activation domain. Neither of the hybrid GAL4 domains alone can promote transcription of a reporter gene, nor do they interact spontaneously. However, if one of the proteins expressed as a fusion with the GAL4 transcriptional activation domain happens to form a **protein-protein interaction** with the bait protein, the two GAL4 domains are brought into proximity where they can activate expression of a reporter gene. Many variations of this powerful method have since appeared, including two-hybrid systems for mammalian cells (6) and three-hybrid systems for detecting protein-nucleic acid interactions (7).

The other major type of expression library uses synthetic DNA to express DNA, RNA, or protein products for screening. Expression of RNA and protein can be mediated through cloning vectors in similar fashion to the cDNA expression libraries. Examples of cloned synthetic DNA expression

libraries include **phage display** and [protein engineering](#) projects in which portions of a protein-coding sequence are randomized synthetically to derive novel proteins with altered properties. Alternatively, the synthetic DNA library may be maintained as oligonucleotides that are replicated through **PCR** or some other enzymatic amplification strategy. Expression of these types of libraries always occurs *in vitro* and usually employs a synthetic phage T7 promoter to mediate expression of RNA. Thus, RNA copies of the DNA library are produced upon incubation with purified T7 RNA polymerase and ribonucleoside triphosphate (rNTP) building blocks. The resulting single-stranded RNA can be screened for affinity to specific molecular targets (RNA aptamer libraries) or for catalytic activity ([ribozyme](#) libraries). Similarly, single-stranded DNA copies can be prepared from DNA aptamer libraries by removing one of the two DNA strands, which has been previously tagged with [biotin](#), by **denaturation** and affinity chromatography over a [streptavidin](#) matrix. DNA aptamers are similar in many ways to RNA aptamers, but they are more resistant to chemical and enzymatic degradation. Finally, peptide aptamers can in principle be prepared by *in vitro* translation of RNA transcripts from oligonucleotide libraries, although this method is not commonly employed.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

### Bibliography

1. J. Schmitt, H. Hess, and H. G. Stunnenberg (1993) *Mol. Biol. Rep.* **18**, 223–230.
2. D. B. Smith and K. S. Johnson (1988) *Gene* **67**, 31–40.
3. A. Cramer, E. A. Whitehorn, E. Tate, and W. P. C. Stemmer (1996) *Nat. Biotechnol.* **14**, 315–319.
4. S. Pietrokovski (1994) *Protein Sci.* **3**, 2340–2350.
5. C. T. Chien, P. L. Bartel, R. Sternglanz, and S. Fields (1991) *Proc. Natl. Acad. Sci. USA* **88**, 9578–9582.
6. Y. Luo, A. Batalao, H. Zhou, and L. Zhu (1997) *Biotechniques* **22**, 350–352.
7. D. J. SenGupta, B. Zhang, B. Kraemer, P. Pochart, S. Fields, and M. Wickens (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8496–8501.

### Expression Systems

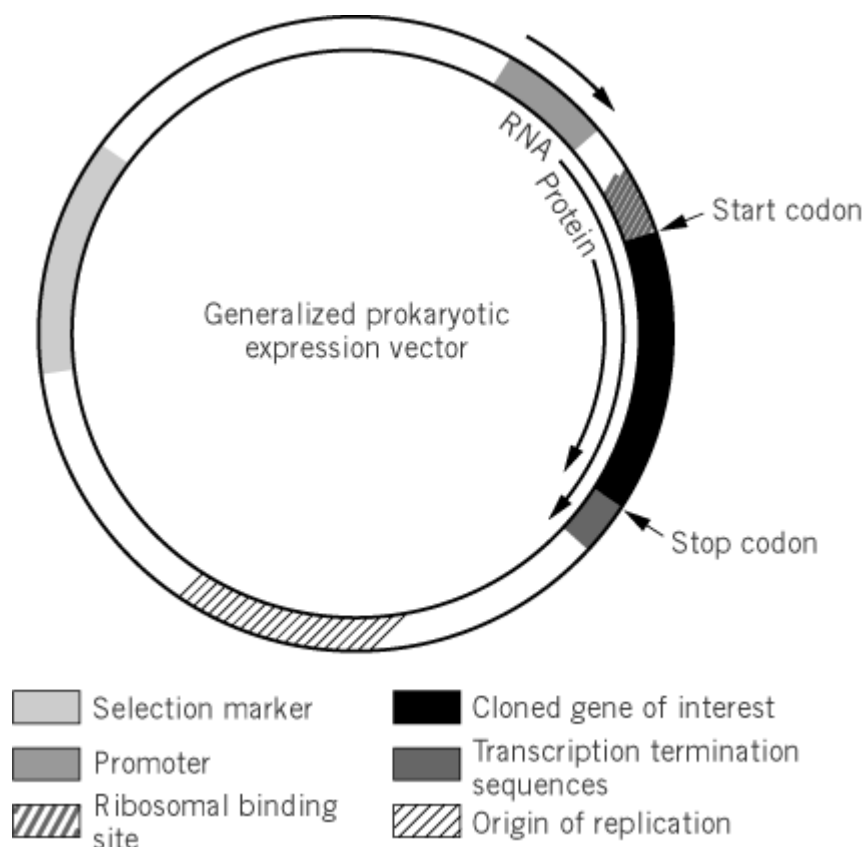
Expression systems are used for the production of [protein](#) from [recombinant DNA](#) molecules (see [Recombinant Proteins](#)). They are of widespread use in industry, health care, and scientific research because of their flexibility and ability to achieve production levels exceeding those of the native source of the protein. An expression system consists of a **vector** carrying the gene encoding the protein of interest, along with the sequences necessary for [transcription](#) of the DNA into [messenger RNA](#) and [translation](#) of the mRNA into protein, plus a host providing the **enzymatic** machinery for carrying out these processes. In a homologous expression system, the gene to be expressed derives from the same species as the host whereas, in a heterologous expression system, the gene to be expressed and the host are of different origin.

#### 1. Design of the expression vector

An expression vector not only provides the components necessary for [cloning](#), transfer, stability, and multiplication of recombinant DNA but it also delivers the elements required for correct transcription of DNA into functional mRNA and efficient translation of mRNA into the desired

protein. However, the importance of the basic elements controlling gene dosage and selection should not be ignored in the design of an expression system. A generalized **prokaryotic** expression vector is shown in Figure 1.

**Figure 1.** Generalized prokaryotic expression vector. The figure shows a typical prokaryotic expression vector, with a selection marker gene allowing stable maintenance of the plasmid, an **origin of replication** controlling the **copy number** of the vector, and a promoter and a ribosomal binding site positioned upstream of the gene of interest. One or more transcriptional termination sequences are positioned downstream. Furthermore, a eukaryotic expression vector would contain an extra origin of replication and an extra selectable marker to allow replication and maintenance in both *E. coli* and the eukaryotic host. The transcriptional termination signal would include a polyadenylation signal.



A **promoter** is a region of DNA recognized by an **RNA polymerase** and is the prerequisite for initiation of transcription. Promoters consist of characteristic sequence elements. In general, promoters from prokaryotes and **eukaryotes** differ, and a common organization pattern can be given for each. However, there can be pronounced differences between species, and more subtle variations are found within the same species, giving rise to promoters with different strengths and control mechanisms.

The strength of a promoter determines the frequency with which a gene is transcribed by controlling the rate at which the RNA polymerase and the promoter form the initiation complex. However, the strength is governed not only by the interaction with the RNA polymerase. Other proteins, so-called transcriptional activators and **repressors**, bind to the promoter region and regulate the process of transcriptional activation. Often, eukaryotic expression vectors provide transcriptional **enhancer** elements as well.

The ideal promoter for the expression of a recombinant protein not only generates the synthesis of high levels of mRNA but, even more important, it should be inducible, ie, the experimenter should

be able to control its activity. However, only a few inducible promoters are tightly controlled, so that some transcription occurs even in the uninduced state. The overproduction of a protein will most often influence the cell in a negative manner and can be detrimental in the most severe cases, where the protein of interest is toxic to the host (see [Poison Sequence](#)). If this seems to be a problem, it may prove advantageous to use a weaker promoter under tighter control. Also, when overexpression seems to outstrip the host's [post-translational modification](#), **molecular chaperone**, and proofreading systems, the use of a weaker promoter may be beneficial. For most expression systems, both constitutive and inducible promoters are available. The “on” state of the best inducible promoters can be regulated, thereby providing a means of optimization (1).

Some prokaryotic expression vectors contain antitermination elements, which stabilize the RNA polymerase on the DNA template to ensure optimal elongation of the transcript. Transcriptional terminators positioned at the 3'-end of the expressed gene restrict the size of the mRNA, minimize the sequestering of RNA polymerase, and isolate the plasmid's replication functions, thereby stabilizing the **plasmid**. Eukaryotic mRNAs are processed more extensively than prokaryotic mRNA, including such processes as **splicing** and the addition of **poly-A** tails. The signals for these steps are provided by eukaryotic expression vectors at positions downstream of the cloned gene.

The degradation of mRNA allows a biological system to adapt to changes in the environment (see [RNA Degradation In Vitro](#)). Yet the aim of the gene expression systems is to maintain a high level of mRNA for the gene of interest during its expression. The stability of mRNA can, in some cases, be increased by insertion of sequences that fold into stable **secondary structures** at the ends of the mRNA, thereby “capping” the mRNA.

The initiation of the translational process requires a **ribosomal** binding site, where the small subunit of the ribosome binds to the mRNA to form the initiation complex. The first triplet of nucleotides to be translated is the [start codon](#), which usually has the sequence AUG. In prokaryotes, a [Shine–Dalgarno sequence](#) is positioned 4 to 13 nucleotides upstream of the start codon. This purine-rich sequence base-pairs with the 3'-end of the 16 S ribosomal RNA (rRNA) and thereby regulates the specificity of ribosomal binding. So-called specialized ribosome systems can be constructed that direct ribosomes specifically to the mRNA of interest. [Site-directed mutagenesis](#) is applied to both the Shine–Dalgarno region and the rRNA to change their sequence, while maintaining the base pairing. In eukaryotes, the formation of a translational initiation complex is directed by the **5'-cap** - structure, 7MeG(5'-5')pppNp (N is the first nucleotide of the mRNA), at the 5'-end of the mRNA. However, several eukaryotic mRNAs have been found to be translated in a cap-independent manner (2).

Secondary structures in the 5'-untranslated region of the mRNA and around the start codon may hamper recognition of the start site. If potential secondary structures seem to be any problem, they can be eliminated by introduction of [silent mutations](#) at the third positions of the initial **codons** or by random [mutagenesis](#) in the region upstream of the start codon. Another approach to overcome the problem of secondary structures is to express the protein of interest as a **fusion protein** or to use a two- **cistron** system. In the first approach, a nucleotide sequence encoding a protein or peptide that is known to be well expressed is introduced in-frame between the start codon and the gene of interest. The resulting N-terminal fusion partner most often provides other advantages as well, such as increased stability or simplified purification (see [Fusion Gene, Fusion Protein](#)). In the second approach, a strong ribosomal binding site and the 5'-end of a well-translated gene, followed by a [stop codon](#), is positioned immediately upstream of the gene of interest. Translation of the protein of interest will then be reinitiated immediately after the first gene.

The [degeneracy of the genetic code](#) implies that most of the 20 [amino acids](#) are encoded by two or more codons called synonymous codons. The synonymous codons are not used with equivalent frequencies by different strains or even throughout the [genome](#) of a single organism (3, 4) (see [Codon Usage and Bias](#)). Weakly expressed genes are characterized by the occurrence of infrequently

used codons, which are typically recognized by rare tRNA species. It is advisable to avoid the use of rare codons when possible in systems for high-level expression.

The three stop codons UAG, UAA, and UGA differ in their efficiency in terminating translation. UAA seems to be favored in highly expressed genes and [should be used in] expression systems.

If a heterologous expression system is used, the gene of interest should preferably be devoid of **introns**, which may not be properly processed by a host of different origin than the gene. However, the presence of introns has, in some cases, been required for successful expression. In such cases, the intron can be provided by the expression vector downstream of the gene of interest.

## 2. Choice of host

Considerable knowledge about expression in *Escherichia coli* has accumulated over the years, and *E. coli* is often the first choice of host when a new protein has to be expressed. Furthermore, *E. coli* grows fast in inexpensive media, allowing scale-up for industrial purposes. However, proteins containing [disulfide bonds](#) are unlikely to fold correctly in the cytosol of bacteria, which is relatively reducing. Misfolding may lead to the formation of [inclusion bodies](#), which can be both an advantage and a disadvantage. The advantages are that the protein is protected from proteolytic degradation and that the purification strategy can be simplified. However, the protein needs to be denatured and refolded, a task that can cause difficulties (see [Protein Folding In Vitro](#)). The coexpression of folding factors such as molecular chaperones, [thioredoxin](#), or [protein disulfide isomerase](#), which assist the folding of protein, can be beneficial in some cases. A secretion strategy may represent another solution to the problem [see [Secretion Vector](#)]. Even though *E. coli* is able to secrete proteins, the outcome of such an approach is highly unpredictable. *Bacillus* species will often be a better choice of host in such a case. Furthermore, *Bacillus* has the advantage that it does not produce [endotoxins](#) (lipopolysaccharides) and has been classified as a GRAS (which stands for *generally regarded as safe*) organism. The range of promoters available for *Bacillus* is more limited than the range for *E. coli*. Some other prokaryotic hosts, eg, *Lactococcus lactis* (5), also are gaining popularity as expression hosts.

On the other hand, bacterial systems have this limitation: they are unable to provide many of the [post-translational modifications](#) often found in eukaryotic proteins. If these modifications are needed for obtaining proper structure or activity, a eukaryotic host should be chosen.

**Yeast**, the simplest eukaryotic expression host, offers many advantages of both prokaryotic and eukaryotic systems. Yeasts grow rapidly in inexpensive media and are easy to manipulate genetically. Furthermore, yeasts provide an environment for carrying out secretion and providing post-translational modifications that are more similar to those found in proteins from the higher eukaryotes. The yeast strain traditionally used as a host for protein expression, *Saccharomyces cerevisiae*, is regarded as safe, given its long history of use in the production of food and beverages. More recently, the methylotrophic yeasts *Pichia pastoris* (6) and *Hansenula polymorpha* have gained popularity because of their higher production levels. However, yeast shows a tendency to **hyperglycosylate** the overexpressed protein, and the resulting high-mannose polysaccharide structures may affect the activity or folding or both, and they are potentially **immunogenic**. The filamentous **fungi**, *Aspergillus* and *Trichoderma*, are becoming increasingly popular as expression hosts, not least because of the very high amounts of protein that can be obtained with secretion systems (7).

Another eukaryotic expression system of high popularity is the baculovirus/insect cell system, which is relatively easy, cheap, and fast to use. The first generation of systems utilized the strong transcription signals for expression of the polyhedron protein of the **virus**. The baculovirus system provides more, although not all, of the features characteristic of mammalian cells, and the yields are often high. However, the high yields sometimes lead to the formation of inclusion bodies. Interestingly, the baculovirus system can be used for **phage display** of large, complex disulfide-

containing proteins, thereby overcoming the limitations of the original bacterial system.

If the protein of interest is of mammalian origin and if authenticity is of utmost importance, a mammalian expression system should be chosen. It may also enable the study of an engineered protein in its natural environment; for example, a biological assay could be coupled to the expression system for the evaluation of functional effects. The basic technology is now available, but these expression systems are not very cost-efficient and are sometimes difficult to scale up. Often, it takes a long time to establish a stable system with high expression levels.

Eukaryotic expression systems can be divided into two groups: those that involve transient or stable expression of recombinant genes from transfected DNA molecules, and those that involve helper-independent viral expression vectors. Vectors used for stable expression contain a complete eukaryotic transcriptional unit inserted into a bacterial [replicon](#). The DNA integrates at a low frequency into the host genome and usually directs the expression of the desired protein at low levels. Recombinant virus systems represent powerful tools for the expression of recombinant proteins in cultured cells, animals and man. A comparison of five different eukaryotic expression systems can be found in reference [8](#).

There are several factors to take into consideration before setting up an expression system. How much protein is needed? What is its application? Is authenticity important? Can the presence of heterogeneity be accepted? Are the technologies available? The choice of an expression system will depend largely on the desired use of the expressed protein; for example, even microheterogeneity can be a problem if the protein is to be used for **X-ray crystallographic** structural analysis, as these heterogeneities may hamper the crystallization process ([9](#)).

**Proteolysis** of heterologous proteins is another problem that should be considered in the choice of a host for expression. Bacteria and lower eukaryotes use proteolysis as a primitive immune system, which attacks and eliminates “nonself.” Different proteinase-deficient strains have been constructed to overcome this problem, but these strains grow at a lower rate, and their usefulness for solving a particular problem is often unpredictable. If proteolysis is a problem, it is advisable to test a selection of these hosts. Alternatively, a secretion strategy or a fusion protein strategy can be chosen. In the first case, the gene product is removed from the cytoplasm, where most proteinases are located. In the second case, the fusion partner can have a stabilizing effect. The method of induction of expression should also be considered when proteolysis appears to be a problem. Some induction methods (eg, heat induction) activate the **heat-shock** system of the cell, which includes a whole range of proteolytic enzymes.

In particularly difficult cases, where the protein of interest is toxic or prone to form inclusion bodies, the use of a cell-free protein biosynthesis system may be beneficial. In such a system, the enzymes needed for transcription and translation are present in a cell extract instead of a live organism. One further advantage is the possibility of introducing unnatural amino acids into the protein of interest ([10](#), [11](#)).

Finally, when a suitable expression system has been established, the growth and eventual induction conditions need to be optimized to obtain the maximal yield of product. It is often advisable to decrease the growth rate below the optimal by reducing the temperature, aeration, or nutrition content of the media. The risk of overloading the protein synthesizing machinery, leading to inclusion body formation or misfolded proteins, is thereby avoided. Furthermore, the expression vector is stabilized. Also, the strategy of recovering the expressed protein will influence the product yield and thus deserves optimization.

The vast number of vectors and hosts available should be able to satisfy the demands of any protein to be expressed. Yet protein expression is a somewhat empirical process, making it difficult to foresee which system should be chosen for optimal success. Advances in the understanding of how expression systems work will undoubtedly make for a higher degree of predictability than is

currently the case.

## Bibliography

1. R. S. Donovan, C. W. Robinson, and B. R. Glick (1996) *J. Indust. Microbiol.* **16**, 145–154.
2. R. J. Kaufman (1994) *Curr. Opin. Biotech.* **5**, 550–557.
3. S. Zhang, G. Zubay, and E. Goldman (1991) *Gene* **105**, 61–72.
4. K. Wada et al. (1992) *Nucleic Acid Res.* **20**, 2111–2118.
5. O. P. Kuipers, P. G. G. A. de Ruyter, M. Kleerebezem, and W. M. de Vos (1997) *Tibtech* **15**, 135–140.
6. M. Romanos (1995) *Curr. Opin. Biotech.* **6**, 527–533.
7. R. J. Gouka, P. J. Punt, and C. A. M. J. J. van den Hondel (1997) *Appl. Microbiol. Biotechnol.* **47**, 1–11.
8. S. Geisse, H. Gram, B. Kleuser, and H. P. Kocher (1996) *Prot. Expr. Purif.* **8**, 271–282.
9. R. Giegé and V. Mikol (1989) *Trends Biotechnology* **7**, 277–282.
10. D.-M. Kim, T. Kigawa, C.-Y. Choi, and S. Yokoyama (1996) *Eur. J. Biochem.* **239**, 881–886.
11. D. Mendel, V. W. Cornish, and P. G. Schultz (1995) *Annu. Rev. Biophys. Biomol. Struct.* **24**, 435–462.

## Suggestions for Further Reading

12. P. Gacesa (1997) *Vectors: Expression Systems (Essential Techniques Series)*, John Wiley & Sons, New York. A laboratory manual that enables the reader to make the best choice of method, conditions, and equipment for working with expression vectors. A troubleshooting guide is included.
13. D. V. Goeddel (ed.) (1990) *Methods in Enzymology*, vol. **185**, Gene Expression Technology, Academic Press, San Diego, California. The reader is provided with a basic understanding of the problems that may occur while trying to express a protein in various expression systems. Furthermore, more sophisticated solutions to the problems are provided. Strongly recommended.

Many of the books suggested under the entry “**Protein engineering**” contain several chapters about different host systems.

14. In the journal *Current Opinion in Biotechnology*, one issue per volume (usually issue 5) is dedicated to reviewing the latest advances in the field of expression systems, enabling the reader to become fully updated.

## Extended X-Ray Absorption Fluorescence Spectroscopy (EXAFS)

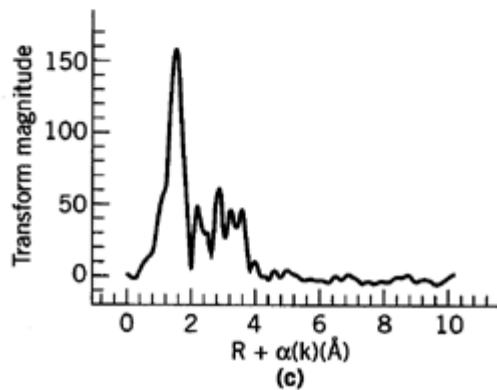
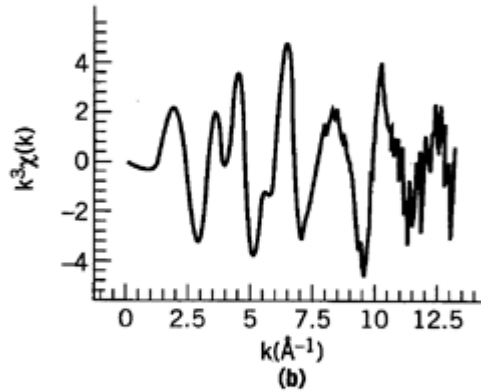
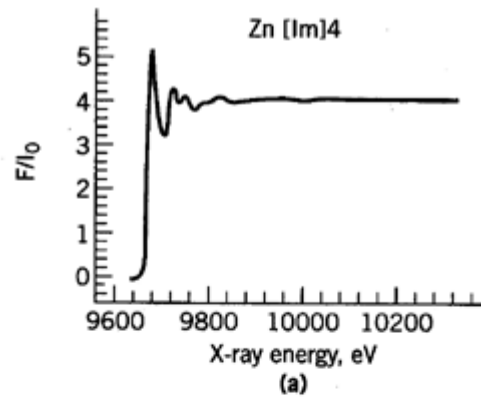
X-ray absorption spectroscopy is particularly useful for investigating the electronic structure and local environment of the metals in [metalloproteins](#) (1). It is not limited by the physical state of the sample (e.g., gel, crystal, tissue) and, unlike [X-ray crystallography](#), the entire structure of a crystalline molecule need not be determined to understand the local structure of the metal site. Obtaining data from light atoms, however, is very difficult. Atoms lighter than sulfur cannot be studied routinely because high vacuum is required. Most of the work reported to date has been with samples containing transition-metal ions. An intense, tunable X-ray source is required, and making measurements at a synchrotron facility must be done on a competitive basis. The ideal sample has a volume of 20 to 100  $\mu\text{l}$  and a concentration of at least 1 mM. It must be maintained at low

temperature while the data are being collected, so as to minimize radiation damage.

Like visible photons, X rays are absorbed by matter according to the Beer–Lambert law (see [Spectroscopy](#)) and cause changes in the energy levels of the electrons of the absorbing atom. Because the energies of X-ray photons are considerably greater than those of visible light, however, the electrons affected are those in the lower shells, principally the K- and L-shell electrons. Therefore, X rays whose energies are near those of the atomic or molecular energy levels of the metal can be used to probe the electronic structures of the metal sites. The edge region of the absorption spectrum in the top panel of Fig. 1 contains information relating to the valence state and coordination geometry of the absorbing metal atom and the chemical identities of the neighboring atoms.

**Figure 1.** X-ray absorption data for  $\text{Zn}[\text{imidazole}]_4^{2+}$  (a) Data collected as fluorescence ( $F$ ) normalized by the initial intensity  $I_0$ . (b) EXAFS data derived from the data in (a) after subtracting a background, corresponding to the metal ion without neighbors, and multiplying by the factor  $k^3$ , where  $k$  is proportional to  $E_{\text{electron}} = E_{\text{X-ray}} - E_{\text{ionizationpotential}}$ . (c) Fourier transformation of the EXAFS data in (b), showing the coordination shells around the  $\text{Zn}^{2+}$  ion. The first shell contains four imidazole N atoms at 2.40 Å, and the higher shells contain the other C and N atoms of the imidazole rings.





When the X-ray energy exceeds the ionization potential or threshold of the metal, which is the energy at which the inner shell electrons are promoted to the continuum, the absorption spectrum is characterized by an oscillatory structure (top panel of Fig. 1) at energies greater than the absorption edge, which is known as the extended X-ray absorption fine structure (EXAFS). This structure results from interference in the ejected photoelectron wave from the metal atom and the photoelectron wave backscattered from neighboring atoms. This interference phenomenon provides a very precise measurement of the distance to the neighboring atoms (typically  $\pm 0.02 \text{ \AA}$ ). In addition, the EXAFS oscillations contain information relating to the number of neighboring atoms of a limited range in atomic number (rows in the periodic table can generally be distinguished). Because this ejected photoelectron wave is spherical, the position in three dimensions is lost on average. In a simplified form, an equation describing these oscillations can be written as

$$k^3\chi(k) = \sum_j N_j A_j(k) \sin(2kr_j + \alpha_j(k))$$

where  $k$  is the wave vector, which is related to the X-ray energy,  $c$  is the oscillatory part of the X-ray

absorption spectrum,  $N$  is the number of neighbors at distance  $r$ ,  $A(k)$  is the backscattering amplitude, which depends on the chemical identity of the neighboring atoms, and  $a(k)$  is the phase shift due to the potentials of both the absorbing metal and the neighboring atoms. The summation is over all coordination shells  $j$ , but in reality data can be collected only at distances less than about 5 Å for biological samples. The previous equation is written in this interference form to facilitate describing the data analysis necessary to determine the structural parameters.

Although X-ray absorption can be measured by observing either the transmission of X rays through the sample or the fluorescence emitted as a result of absorption, the latter is used almost exclusively for biological samples, because of its enhanced sensitivity. To obtain information about the number and average distances of the neighboring atoms, the data must be analyzed to isolate these parameters in the previous equation. Figure 1 shows the results for  $\text{Zn}(\text{imidazole})_4^{2+}$  during each step of the analysis. First, a simple linear background is subtracted from the data to normalize to zero the data below the edge region (top panel). Then, a background resembling the spectrum of the metal without neighbors (ie, a free atom) is subtracted, the energy scale is converted to wave numbers ( $\text{Å}^{-1}$ ), and the data are multiplied by  $k^3$ . The results are displayed in the middle panel of Fig. 1. Now, the  $y$ -axis corresponds to the left side of the previous equation. Fourier transformation is used to isolate the contributions of the individual coordination shells, which converts the  $x$ -axis to  $r + a(k)$ . This is similar to a radial distribution function that describes the radial electron density around the absorbing metal set at the origin. The first peak of the bottom panel of Fig. 1 is the first coordination shell, and those at greater  $x$ -axis values are further from the metal. The  $x$ -axis can be converted to  $r$  if one knows the chemical identity of the scatterers in each shell. For example, the average distance of the first shell peak is  $r + a(k) = 1.52 \text{ Å}$ . It is composed of four nitrogen atoms, however, and it is known from investigating model compounds of known structure that  $a(k) = 0.88 \text{ Å}$  for Zn–N coordination. This makes the average distance of those four nitrogen ligands  $2.40 \text{ Å}$ , with an estimated error of  $\pm 0.02 \text{ Å}$ . Using model compounds, one can also determine that there are four nitrogen ligands in the first coordination shell. A more elaborate analysis, considering the higher shells of scattering atoms, would also lead to the conclusion that the ligands are imidazole. This example illustrates the need for model compound data or theoretical determination of the amplitudes and phases for each chemical type of ligand.

This example also demonstrates one limitation of EXAFS. In the absence of information about the metal-coordination environment, the data simply give an average of the number of ligands and their distances. In the best of circumstances without other biophysical or biochemical information, this method can determine only the difference between elements in the rows of the periodic table (e.g., C/N/O, S/Cl, or Se). Additional limitations occur as a result of uncertainties in the data analysis. For example, it is rarely possible to distinguish between one and two coordinated histidine residues in the presence of [thiol group](#) ligands without further information (such as knowledge of the protein [primary structure](#)). Rigorous analysis of X-ray absorption data for a complex coordination environment (e.g., metal clusters) in proteins can be impossible without methods to isolate each metal site, for example, by metal substitution.

## Bibliography

1. L. S. Powers (1982) *Bioch. Biophys. Acta* **683**, 1–38.

## Suggestions for Further Reading

2. I. Bertini, H. B. Gray, S. J. Lippard, and J. S. Valentine (eds.) (1994) *Bioinorganic Chemistry*, University Science Books, Mill Valley, CA.
3. L. Que Jr.(ed.) (1999) *Bioinorganic Spectroscopy and Magnetism*, University Science Books, Sausalito, CA.
4. G. L. Eichhorn and L. G. Marzilli (eds.) (1989) *Advances in Inorganic Biochemistry*, Elsevier, New York, Vol. **8**.

## Extracellular Matrix

Extracellular matrix (ECM) is a general term for extracellular networks of fibrous glycoproteins, proteoglycans, and carbohydrates that are secreted by cells. There are a wide variety of matrix types that fulfill many important functions in developing and mature animals and plants. ECM can be categorized as strong connective tissue ECM or loose connective tissue ECM. In strong connective tissue, such as bones, teeth, antlers and horns, a meshwork of fibrous protein is strengthened by deposition of calcium phosphate. Loose connective tissue includes specialized structures such as tendons, ligaments, fascia, cartilage, cell walls, and cuticles. Also in this category are the basal laminae (or basement membranes) that underlie epithelial cell layers and surround muscle cells, fat cells, and Schwann cells, and the collagenous interstitial matrix that is secreted by fibroblasts and lies beneath the basal laminae.

ECM contributes to the structure of organisms, acting as a scaffold for building strong connective tissue, and organizing cells into specific patterns in loose connective tissue. Tendons and ligaments contain molecules that give them their unique structural features in connecting bone and muscle. ECM in the joints provides cushioning, and elastic layers around blood vessels allow control of blood flow. The basal lamina in the kidney glomerulus acts as a molecular filter, allowing smaller molecules to pass from the blood to the urine. ECM can form important barriers in tissues to coordinate cellular movements and organization of cellular layers. For example, the basal lamina underlying endothelial cells in blood vessels prevents cells from leaving the circulation, except for certain cells that have the ability to pass through. This barrier is especially important during inflammation and metastasis.

While the important structural roles played by ECM have been recognized for a long time, recent research has uncovered crucial instructive activities that direct cellular movement and regulate cellular differentiation and gene expression. Individual molecules have now been purified, and their activities have been characterized. Many ECM molecules have been shown to support adhesion of cells, and some (such as tenascin) display anti-adhesive activities that may be important in establishing tissue boundaries or pathways of migration. Early in development, neuroblasts and myoblasts contact ECM as they migrate to their final destinations. ECM is also important in axon pathfinding and synapse formation. The basal lamina surrounding muscle cells has specialized regions in the synaptic cleft at the neuromuscular junction that contains molecules important for synapse function, such as acetylcholinesterase. In addition, synaptic basal lamina contains molecules that help guide synapse formation and regeneration, such as laminins containing the  $\beta$ -2 subunit and agrin. *beta*-2 laminins are also important in photoreceptor development and synapse formation. ECM also acts as a reservoir of positively charged growth factors, many of which bind to negative charges in ECM molecules. Growth factors and ECM molecules act in concert to direct a wide variety of cellular traffic during tissue morphogenesis. For example, fibroblast growth factor must bind to heparan sulfate chains from proteoglycans to be stable and active. In addition, the activity of some growth factors is enhanced when cells are bound to ECM, due to a synergy of their respective intracellular signaling pathways. Finally, some ECM molecules are proteolyzed to generate smaller soluble factors. Collagen XVIII is cleaved to generate endostatin, an antiangiogenic motility factor that can prevent tumor progression in mice. It has become clear that the ECM is in part the mortar that holds cells together in tissues, but it also acts as the brick mason, guiding the development and deposition of cells within the organism.

Cells interact with and respond to specific matrix molecules via cell surface receptors such as the Integrins, which bind matrix on the outside of the cell, and send signals across the plasma membrane to regulate intracellular events. Some cells make identifiable, stable junctions with ECM via

integrins. These include focal adhesions which link to the actin cytoskeleton and hemidesmosomes, which connect to intermediate filaments. Other cell interactions with ECM are more transient. Contact with the matrix can alter the polymerization dynamics of the cytoskeleton to bring about changes in cell shape important in cell adhesion, migration and polarity. Adhesion to ECM is important in cellular responses to physical/mechanical forces, such as shear stress caused by blood flow over endothelial cells. ECM can activate signaling cascades involving tyrosine kinases which control cell growth, differentiation and apoptosis (see [Integrins](#)). Other receptors may also play roles in interacting with the ECM, including proteoglycans and galactosyl transferases.

The molecular complexity of the ECM is formidable, and many components are polymerized and then interwoven by covalent and noncovalent cross-links. The specific molecular attributes of ECM components give rise to their specific functional roles. For example, collagen I, the most abundant protein in animals, forms covalently cross linked fibrils with the tensile strength of steel that are found in skin, tendon, ligament, and bone. In contrast, Collagen IV forms sheets that polymerize with heparin sulfate proteoglycans and laminin-1 to give rise to the basal lamina. The presence of negatively charged, unbranched carbohydrate polymers called glycosaminoglycans (GAGs), both as free molecules and attached to protein backbones in proteoglycans, form a cushion-like hydrated gel important in the function of joints. For example, hyaluronan, a free GAG molecule that can be as long as 20 mm if stretched from end to end, forms a viscous random coil that provides turgor pressure in spaces between cells. Aggrecan, a proteoglycan, is a key cushioning element in cartilage. The elastic properties of cross-linked elastin arrays allow lung bronchial sacs to expand and arteries to change their diameter. Many ECM proteins are composed of repeated motifs that arise from individual exons, suggesting that new proteins evolved by shuffling various exons. Electron micrographs of purified ECM molecules have revealed striking geometries, from cross-shaped structures like laminin-1, to centipede like arrays of proteoglycans in an aggrecan aggregate, to pentameric and hexameric asterisks like thrombospondin-4 and tenascin, respectively. Some ECM proteins have many different isoforms arising from both alternative splicing and multiple genes that encode homologous families of subunits. For example, there are over 20 varieties of collagens made from 38 distinct polypeptides. At least 15 types of laminin have been identified. Other ECM components such as vitronectin and fibronectin are made in one form and deposited in matrices, and a second form is an abundant component of blood plasma.

The macromolecules that make up the ECM are synthesized by cells that reside within the matrix. This includes fibroblasts in the interstitial matrix, chondroblasts in cartilage, and osteoblasts in bone. Some fibril forming ECM proteins such as collagen are made initially with N- and C-terminal extensions called propeptides that prevent polymerization until the propeptides are cleaved outside the cell. Collagen also contains posttranslationally modified hydroxyproline and hydroxy lysine, which are important for hydrogen bonds that stabilize fibrils. The enzymatic addition of the hydroxyl groups requires ascorbic acid (vitamin C), and deficiencies lead to weakened blood vessels and loose teeth characteristic of the disease Scurvy. Outside the cell, triple helical collagen monomers are cross-linked via covalent bonds to form 50 nm extracellular fibrils. Cells can assemble complex ECM arrays on their surface via the activity of integrin and dystroglycan receptors.

Once assembled, ECM is insoluble and fairly stable. For example, a collagen I molecule in bone might last 10 years before it is replaced. However, in some cases the degradation and turnover of matrix does occur, and it is carefully regulated. White blood cells must degrade the vascular basal lamina to leave blood vessels and enter inflamed tissue, and metastatic cancer cells can migrate to distant sites by violating ECM boundaries. Matrix is also degraded at certain times during development and wound healing, such as in the sprouting of new blood vessels (angiogenesis). In these cases most matrix is degraded by a family of over 20 secreted proteases called Matrix Metalloproteinases (MMPs), which require bound  $Zn^{2+}$  or  $Ca^{2+}$  for activity. MMPs are secreted as inactive precursors which are activated by proteolytic cleavage. Activity is also regulated by a family of secreted protease inhibitors called tissue inhibitors of metalloproteinases (TIMPs). The process of matrix degradation is of great importance to medical researchers interested in controlling

inflammation, tumorigenesis, and metastasis.

The importance of ECM is underscored by the consequences of mutations that ablate or alter matrix components. For example, human mutations in collagen I can result in osteogenesis imperfecta (weak bones) or Ehlers-Danlos Syndrome (defective joints), or even death. Over 1,000 mutations have been identified in 22 different collagen genes, and these lead to a variety of diseases characterized by defective connective tissue. Mutations in fibrillin results in Marfan's syndrome, a disease characterized by defects in elastic tissue such as the aorta, which is subject to aneurysms. Mutations in laminin genes can lead to congenital muscular dystrophy and junctional epidermolysis bullosa, a skin blistering disease. More recently, genetic studies in model animal systems have also demonstrated the importance of ECM. For example, loss of laminin in the worm *Caenorhabditis elegans* results in defective mesodermal cell migration and axonal pathfinding under the epidermis. Mice with null (knockout) mutations in the fibronectin gene fail to develop notochord, somites, neural tube, and heart and die by embryonic day 10.

There is still a great deal to learn about the ECM. New molecular components continue to be discovered, and their functions are not yet known. Furthermore, alternatively spliced variants of known molecules, and isoforms derived from separate genes have been discovered which differ in structure and function. Proteolytic fragments of ECM molecules, such as endostatin, angiostatin, and tumstatin, are being investigated as possible anticancer treatments. The exact molecular structure of complex ECMs such as the basal lamina are not yet clear, and how cells adhere to, migrate on, and move across matrices is still not completely understood. As more genomes are sequenced, we should get a better idea of ECM complexity and evolution. Developing better treatments for important human diseases such as arthritis, cancer, and macular degeneration may be made possible by a better understanding of the ECM.

#### Additional Reading

Kreis T. and Vale R., eds., *Guidebook to the Extracellular Matrix, Anchor, and Adhesion Proteins*, 2nd ed., Oxford University Press, Oxford, New York, 1999.

Ruoslahti E. and Engvall E., eds., Extracellular matrix components, *Methods in Enzymology*, Academic Press, San Diego, Calif., 1994, p. 245

Hay E.D., ed., *Cell Biology of Extracellular Matrix* 2nd ed., Plenum Press, New York, 1991.

## Extrachromosomal Inheritance

“Extrachromosomal inheritance” is a term to be avoided. It is used sometimes as a synonym for [cytoplasmic inheritance](#), which usually arises from the genetic information on [mitochondrial](#) or [chloroplast](#) genomes. Using the term **chromosome** in the broad sense of any DNA or RNA molecule capable of autonomous replication, the term *extrachromosomal inheritance* has a very limited content (see **Maternal effects**).

If the term [chromosome](#) is restricted to the complex structures in the **nuclei** of the eukaryotes, then *extrachromosomal inheritance* covers both cytoplasmic inheritance and the inheritance of plasmids that replicate autonomously in the nuclei of the cells. Because their distribution patterns differ from those of the major chromosomes, their inheritance diverges from Mendelian rules (see [Mendelian Inheritance](#)).

Most nuclear plasmids are phenotypically silent and can only be detected by identification of their

molecules. This is the case of the “two-micron circle,” a nuclear plasmid of *Saccharomyces cerevisiae*. Although it does not confer any phenotype to the yeast, it has been useful in the development of **vectors** for **genetic engineering**.

In bacteria, *extrachromosomal inheritance* could designate the inheritance of autonomous plasmids separate from the main chromosome (see [F Plasmid](#)).

### Suggestion for Further Reading

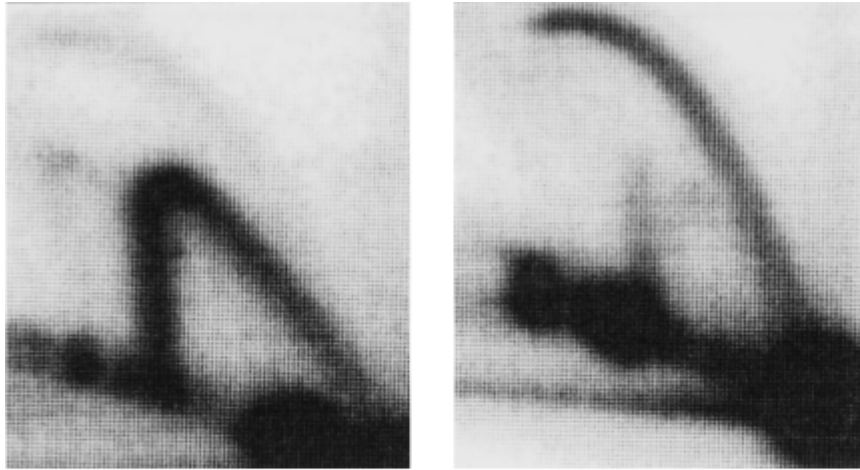
J. R. Broach and F. C. Wolkert (1991) "Circular DNA plasmids of yeasts. In" *The Molecular and Cellular Biology of the Yeast Saccharomyces* (J. R. Broach, J. R. Pringle and E. W. Jones, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 297–331

## Eye (Bubble)-Form Intermediate

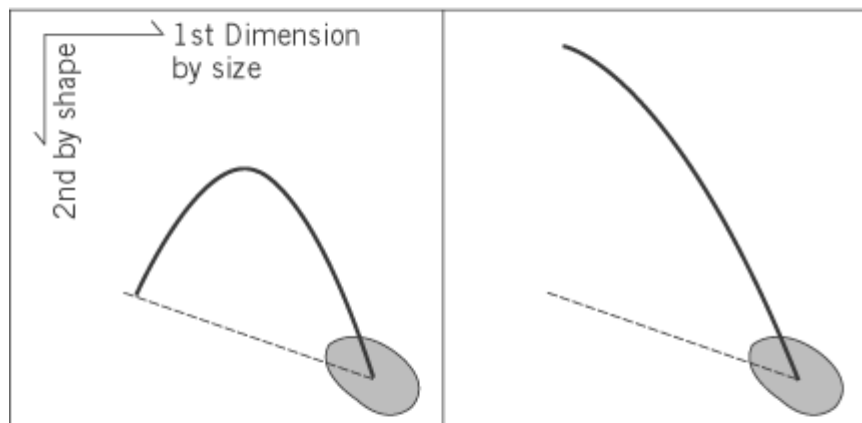
In general, [DNA replication](#) of a **chromosome** starts from a genetically fixed site and proceeds bidirectionally. At the initial stage of replication, the replicated regions form structures like an eye (or a bubble). This form is important, as it represents the origin region of the chromosome (see [Replication Origin](#)); various methods to detect eye-form intermediates have been developed. [Electron microscopy \(EM\)](#) of the replicating DNA can detect eye-forms as large as several kilobases (1). Using a combination of digestion by appropriate [restriction enzymes](#) and electron microscopy, the origin of replications of **plasmids**, **bacteriophage**, and **virus** DNA has been identified (2). Recently, the [two-dimensional gel electrophoresis](#) method was developed to detect various forms of replication intermediates (3). Chromosomal fragments produced by digestion with restriction enzyme are separated by electrophoresis by agarose gel electrophoresis, first by size and then by shape. The distribution of intermediates on the two-dimensional (2D) gel is detected by **hybridization** with a **radiolabeled** probe of a cloned DNA fragment from or near the replication origin. Since the correlation between the size and complexity of shape of eye-forms is different from that of Y-fork intermediates (see [Replication Fork \(Y-Fork Intermediate\)](#)), the former can be distinguished from the latter by the shape of arcs on the 2-dimensional gel (Fig. 1). In particular, when the origin of replication is not located at the center of the fragment, the change from eye- to Y-form occurs as replication proceeds on the fragment, which results in a discontinuous distribution of arc from eye- to Y-form. This method has been used successfully to detect the chromosomal origins of **eukaryotes**, in particular of *Saccharomyces cerevisiae* (4), *Schizosaccharomyces pombe* (5), and *Drosophila* chromosomes (6). Detection of eye-forms from mammalian chromosomes is technically difficult, due to the large amount of background DNA, and this method generally produces ambiguous results as to the site of origin of replication (7).

**Figure 1.** Two-dimensional gel electrophoresis method to detect eye-form (bubble-form) and Y-form replication intermediates. (a) Examples of patterns of Y-form (left) and Eye-form (right) intermediates of the replication of *Saccharomyces cerevisiae* chromosome are shown. Intermediates were detected by hybridization with the <sup>32</sup>P-labeled probe of cloned DNA within the fragments shown in (c). (b) The photos in (a) are drawn schematically. Dotted lines indicate the location of linear DNA fragments of corresponding size. (c) The population of the intermediates detected in (a) and (b) Left panel: Replication fork moves from the left (or right) side of the fragment and proceeds to produce Y-forks of variable size. Theoretically, a Y-fork with equal branch lengths is structurally most complex and shows least mobility in the second dimension. Right panel: Replication initiates from the center of the fragment and proceeds bidirectionally, producing various sizes of bubble intermediates. The largest bubble is like a circular DNA and shows least mobility in the second dimensions. The locations of the probe used for detection in (a) are indicated by dotted lines. Dark spots in the photos in (a), and drawn schematically in (b), indicate that the majority of DNA detected by

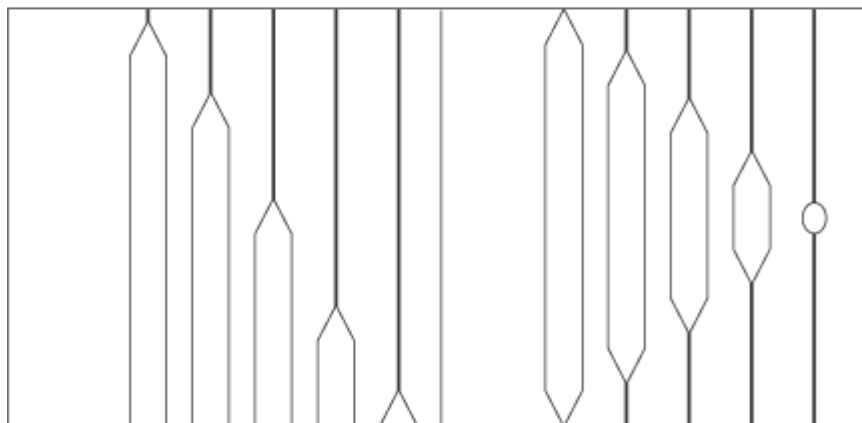
the probe is nonreplicated (or completely replicated) linear molecules.



(a)



(b)



(c)

### Bibliography

1. M. Schnos and R. B. Inman (1970) *J. Mol. Biol.* **51**, 61–73.
2. J. Wolfson, D. Dressler, and M. Magazin (1972) *Proc. Natl. Acad. Sci.* **69**, 499–504.
3. B. J. Brewer and W. L. Fangman (1987) *Cell* **51**, 463–471.
4. W. L. Fangman and B. J. Brewer (1991) *Ann. Rev. Cell Biol.* **7**, 375–402.
5. J. Zhu, D. L. Carlson, D. D. Dubey, K. Sharma, and J. A. Huberman (1994) *Chromosoma* **103**,

414–422.

6. T. Shinomiya and S. Ina (1991) *Nucl. Acids Res.* **19**, 3935–3941.

7. M. L. DePamphilis (1993) *J. Biol. Chem.* **268**, 1–4.

## F Plasmid

In 1953 a new type of genetic element was postulated in order to account for the sexual behavior of *Escherichia coli* K-12 (1, 2). This determinant spread rapidly and independently through recipient cell populations, indicating considerable autonomy, but it was completely dispensable and conferred no **phenotype** on its host other than maleness. Nevertheless, it was stably inherited by the cells it invaded. It behaved, in essence, like a small, infectious **chromosome**, and it was baptized the fertility (F) or sex factor. The nature of the proposed new element was amply confirmed: F turned out to be the prototype of the circular, double-stranded DNA **plasmid**.

F was immensely important to the development of molecular genetics. By promoting transfer of chromosomal genes, it opened the way to large-scale mapping of the bacterial **genome**; and by acquiring and transferring small portions of the host chromosome, it could create partial **diploid** conditions that enabled analyses of bacterial **gene** function and control of **gene expression**. The nature and consequences of F-specified maleness are generally understood and are described in the entry **Hfr's and F-primes**. F has also been at the forefront of investigations into how low-copy number **replicons** are maintained stably; the mechanisms that regulate F's maintenance will be the major focus here.

Early studies of F were greatly aided by the discovery of derivatives that carry chromosomal genes (F-primes). For example, isolation of F' lac<sup>+</sup> mutants (in a lac mutant strain) that remained Lac<sup>+</sup> at 30°C but segregated Lac<sup>-</sup> at 42°C demonstrated that F carries its own replication determinants (3). The selective inhibition of F maintenance by **acridine dyes** (4) also pointed to the autonomy of F replication. Measurements of β-galactosidase in F<sup>-</sup>/lac<sup>+</sup> and F' lac<sup>+</sup>/lac<sup>+</sup> strains gave the first indication that the copy number of F is about the same as that of the chromosome. The strict regulation of copy number was reflected in the phenomenon of incompatibility, the inability of two F plasmids (eg F' lac and F' gal) to be coinherited. The F replication system enables integrated F to assume control of chromosome replication in strains unable to initiate it themselves, a phenomenon known as “integrative suppression” (5, 6). Beyond the initiation of replication at the origin, however, F replication becomes strictly dependent on host functions.

The term “episome” was coined to denote plasmids like F with both autonomous and passive (integrated) modes of propagation. This property turned out not to reflect a fundamental distinction from other plasmids, and the term is no longer considered useful.

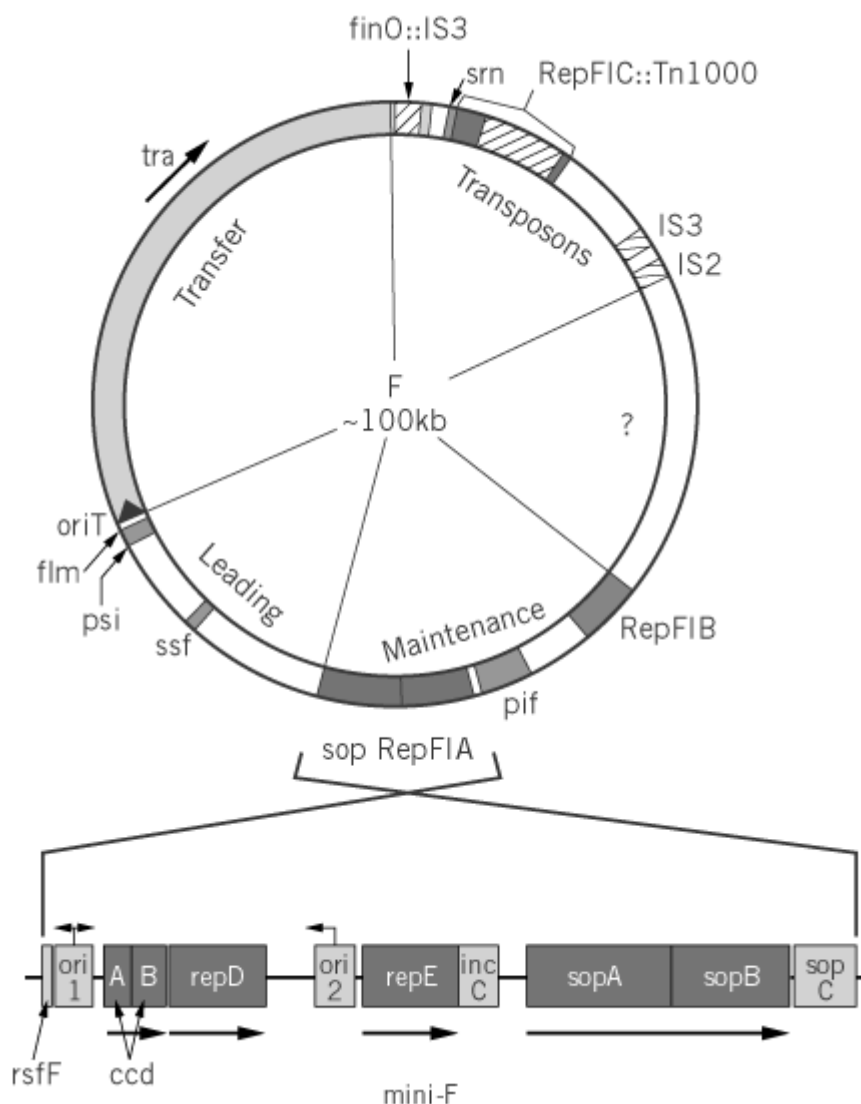
### 1. Anatomy of F, A Mosaic

F consists of four functional blocks and a region of unknown composition (Fig. 1). The 33-kbp **tra operon** encodes the essential components of donor ability, the mating apparatus and the DNA manipulation enzymes responsible for transmitting F into recipient cells, as shown by the self-transmissibility acquired by plasmid vectors into which it, along with the origin of transfer (*oriT*), has been inserted. The ~ 13-kbp leading region that follows *oriT* is the first to enter recipient cells during transfer; only a few genes within it are known, and their involvement in transfer has not been



demonstrated. The region at the other end of the *tra* operon contains a defunct replicon and appears to have acted as a trap for [transposable elements](#); the latter contribute indirectly to mating activity, because one of them is inserted in the *finO* regulator gene, causing constitutive expression of the *tra* operon and high transfer rates, and (ii) also because homologous [recombination](#) with similar elements in the chromosome is the major route of F integration to form Hfr donors. The fourth region carries all the significant determinants of F maintenance. When isolated and joined to an [antibiotic resistance](#) gene to permit selection, it behaves as a normal plasmid (mini-F) with the same replication, partition, and stability properties as its parent F (7, 8). It is, therefore, the form most often used for studies of F maintenance.

**Figure 1.** Genetic map of the F plasmid. The circular F plasmid has been divided into five sectors, roughly according to function. Gene symbols are explained in the text. The *tra* arrow indicates the direction of transcription of the *tra* operon. The *oriT* arrow indicates the direction of conjugal transfer. Unshaded blocks of the circle correspond approximately to unsequenced regions. Below is the *repFIA* mini-F; note that it is shown inverted relative to F, by convention. Arrows above *ori1* and *ori2* show bi- and unidirectional replication, respectively.



A closer view reveals F to be a mosaic of replicative determinants. Three have been identified: *repFIA*, *repFIB*, and *repFIC*. The first two can be isolated as miniplasmids, although replication of *repFIB* is partly defective, so that replication of F itself is predominantly governed by *repFIA*. The *repFIC* locus, identified as such by comparison with the *rep* sequences of other plasmids, is inactive

due to insertion of *Tn1000* in a progenitor of F. Portions of F as yet unsequenced may prove to contain other replicons. The presence of multiple replicons has proved to be a common feature of large plasmids, and it suggests that many such plasmids may have arisen by the fusion and rearrangement of a limited number of ancestral plasmids (9, 10). The fact that only *repFIA* is fully functional may account for the fact that F is stably maintained in only a narrow range of hosts: *E. coli* and similar Enterobacteria.

### 1.1. RepFIA

The *repFIA* locus is itself an ensemble of essential and accessory maintenance functions (Fig. 1). **Electron microscopic** examination of replication intermediates showed that [replication forks](#) emanated bidirectionally from a site termed *ori1* [also called *oriV* (11)]. However, miniplasmid deletion derivatives that lack this site, but show similar replication behavior, were obtained readily (12). These replicate unidirectionally from *ori2* [*oriS* (13)], both *in vivo* and *in vitro* (14, 15). It is not known why the *ori1* site is preferred in the original mini-F. In any event, the smaller mini-F, replicating from *ori2*, has been used in most investigations.

The basic *repFIA* replicon is made up of three components. One is *ori2*. The others are the single essential replication gene, *repE*, whose product regulates all aspects of initiation, and the copy control region, *incC* (13, 16, 17). A plasmid comprising only these elements replicates normally, but is nevertheless unstable. Stable inheritance requires a mechanism to ensure that plasmid copies are directed into each of the daughter cells, and this function is provided by the adjacent *sop* locus (18). *Sop* is essentially a partition module that can act independently of the *rep* locus to which it is normally linked, to stabilize plasmid vectors in which it is inserted. Together, the *rep* and *sop* loci constitute the basic F maintenance system. This arrangement is remarkably similar to that found in the P1 prophage plasmid [a detailed level, *RepFIB* is probably a more exact analogue of the mini-P1 replicon]. Studies of mini-F and mini-P1 have been reciprocally informative.

## 2. F in the Cell Cycle

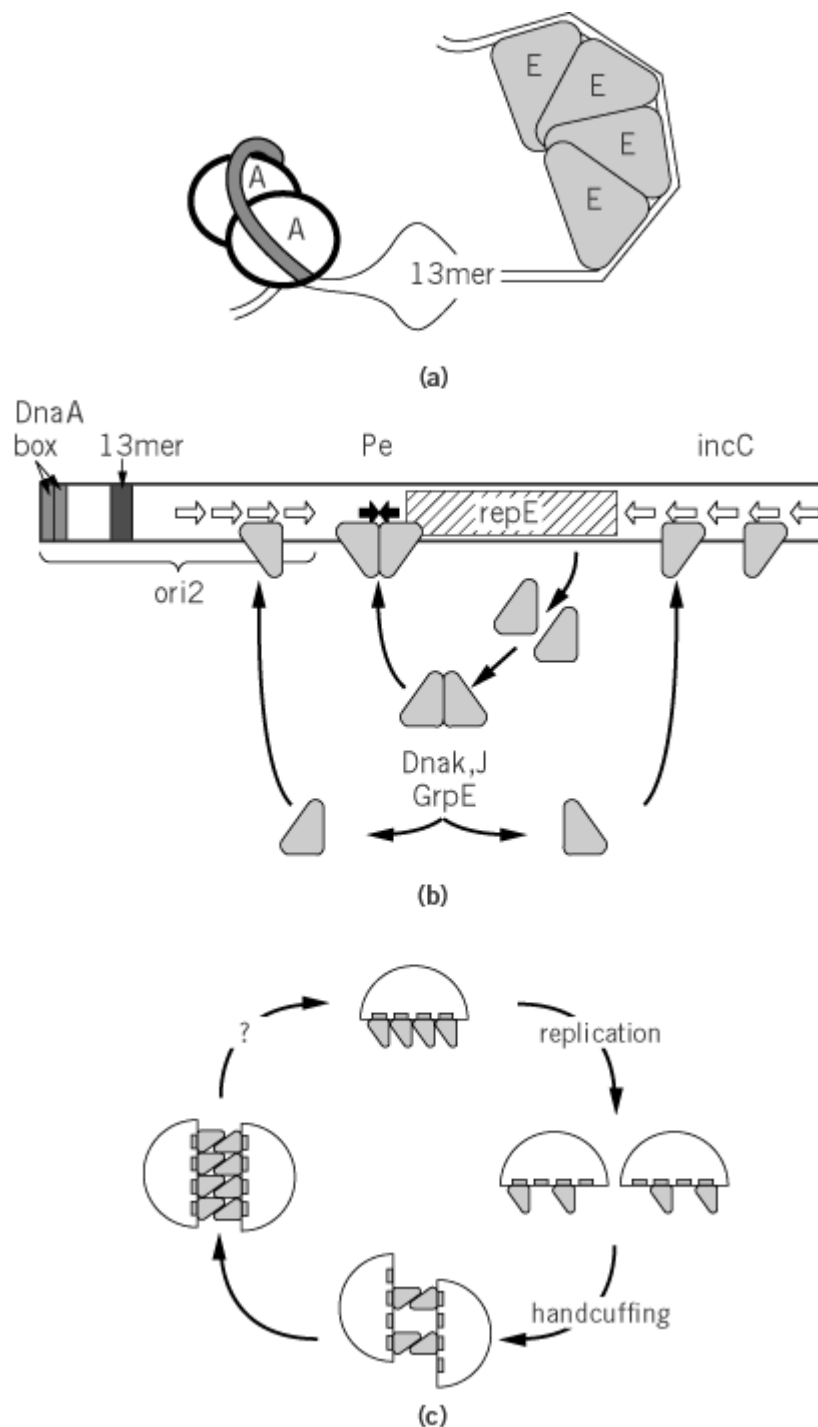
To ensure that F makes replicas available for both daughter cells at division, in spite of its low copy number, it would appear necessary for replication to be coordinated with the [cell cycle](#). Viewed simply, this would imply that F replication take place at a specific time in the cycle for given growth conditions, as does the initiation of chromosome replication. However, a number of careful studies have indicated that both F and mini-F replicate in cells of all ages, without delaying cell division (eg, Refs. 19 and 20). Despite some dissenting opinion (21), it appears reasonable at present to view replication of F as more loosely coordinated with the cell cycle than that of the chromosome. How F could achieve high hereditary stability in the absence of cell cycle coordination remains unclear.

## 3. Replication Control

F provides the regulator of its replication, while the chromosome contributes essential components to the mechanism of initiation at *ori2*. Thus, the host DnaA initiator protein binds on one side of the origin to enhance opening of an AT-rich region (Fig. 2a). Initiation itself, however, requires binding of the F RepE protein to specific sites (iterons) on the other side (22), so that the duplex becomes sufficiently deformed to allow entry of **DnaB helicase** and assembly of a replication complex (23, 24) (see [DNA Replication](#)).

**Figure 2.** Aspects of mini-F replication. (a) Opening of *ori2* at initiation; **A** is DnaA protein, bound at tandem DnaA boxes (shaded DNA), **E** is RepE monomer bound to the iterons of *ori2*; together they cause melting of the AT-rich region at a 13-mer sequence similar to those that are opened at an early stage of initiation in *oriC*. (b) Regulation of RepE synthesis and activity. Each arm of the operator (filled arrows) contains an 8-bp subsequence of the 19-mer iteron binding sites (open arrows), suggesting that initiation, negative control, and autoregulation all depend on interaction of the same RepE binding domain with DNA. Limiting RepE synthesis by autoregulated transcription, along with limiting RepE availability through binding to *incC*, could constitute a copy-control mechanism. But how

can competitive binding be effective if an autoregulatory mechanism exists to replenish the pool of free RepE? One solution is that RepE exists in two forms: It is first synthesized as a repressor, which binds to the operator but not to the origin, and part of it is then converted to an initiator form that binds to *incC* and *ori2*. In this way, titration and autoregulation become to a large extent independent of each other (30). Indeed, the complex formed by RepE with the operator is structurally distinct from that formed with *ori2* (31). The model is further supported by the finding that most RepE is present in a dimer form with operator-binding and repressor activity. The initiator is made when the dimer is converted by the DnaK–DnaJ–GrpE chaperone machine to a monomer, which binds to the *ori2* and *incC* iterons (29). These chaperones are thus needed for normal levels of mini-F replication (32). Certain mini-F mutants that have escaped chaperone-dependence make RepE that has lost the ability to dimerize, resulting in high copy numbers (29, 33). (c) Principle of handcuffing and the cycle of initiation and inhibition; the interactions shown are intermolecular, but intramolecular coupling is also possible. After replication, the RepE remaining at *ori2* pairs with RepE remaining either at *incC* or at *ori2* on a second plasmid. As more RepE is made, it also pairs, inhibiting initiation by hindering access to replication proteins or impeding the structural transition needed for initiation. Failure to pair is a feature of certain copy-mutant initiator proteins of plasmids P1, RK2, and R6K (35-37). An unknown event (question mark) ruptures the coupling to allow initiation.



The *incC* control region acts as a brake on the formation of RepE-*ori2* initiation complexes. Additional copies of *incC* in *trans* inhibit replication, while deletion of *incC* causes about a fivefold increase in copy number (16, 17, 25). Because *incC* binds RepE (Fig. 2b), it was easy to imagine that it competes with *ori2* for RepE, controlling replication by titrating the initiator. An alternative copy-control mechanism sprang from studies of plasmids P1 and RK2 (26, 27), in which it was found that initiators bound to iteron sites can also bind to each other. The iteron regions *ori2* and *incC* could thus pair, either intra- or intermolecularly, to block initiation. This model, termed “handcuffing,” would explain regulation of the frequency of initiation by a dual activity of RepE: initial reinforcement of the handcuff to prevent initiation, followed by formation of an initiation complex to start replication (Fig. 2c). A subtle modulation of RepE synthesis would seem necessary to allow the transition.

In fact, production of RepE initiator is finely regulated. First, RepE represses transcription of its own gene, by binding to the *repE* promoter (28), accounting for the limited quantities of RepE initiator. Second, newly made RepE has only this activity and is unable to initiate replication. This is because newly synthesized RepE forms dimers, which bind to the [inverted repeat](#) operator sequence in the *repE* promoter region but not to the iterons in *ori2*. To become an initiator, RepE must be converted to a monomer form (29-31). This conversion is catalyzed by the **DnaK-DnaJ-GrpE molecular chaperone** machine, a host function that is needed for normal levels of mini-F replication (29, 32, 33).

Despite the identification of these elements of F replication control, a satisfactory explanation of this process is still beyond us. The discovery of the essential role of the chaperone machine has served only to push the question of the key regulatory process one step further back: What determines the rate at which RepE dimer is converted to active monomer? If the handcuffing model is correct, how is the pairing broken to allow replication? Motor forces of the partition system have been proposed as keys (or boltcutters) for removing the handcuffs (34). But if handcuffed plasmids are pulled apart by the partition apparatus, how does unpairing occur in the case of partition-defective mini-F mutants, which replicate at a normal rate? The true status of *ori1* is unclear: If it is really the major origin in wild-type F, why is DnaA-assisted strand-opening over 2 kbp away essential for initiation?

#### 4. Partition

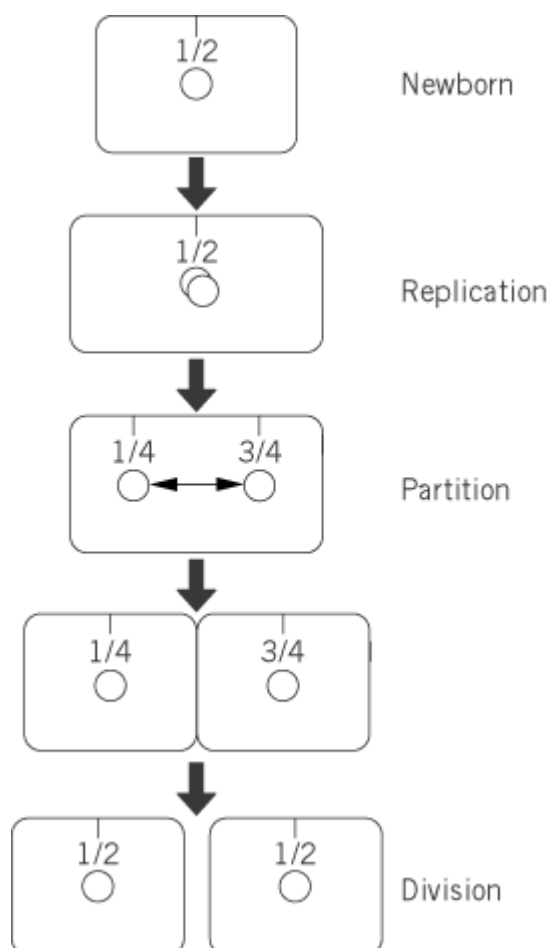
Deletion of the *sop* region causes loss of mini-F at the rate expected for random distribution of a small number of plasmid copies to daughter cells prior to division, consistent with the disruption of an active partitioning mechanism. The *sop* region consists of three elements (Fig. 1), all essential for partition (18, 38). Two partition proteins, SopA and SopB, are encoded by the *sop* operon. SopA is the main regulator of the operon; it binds to a series of short repeated sequences in the promoter to repress [transcription](#) (39). SopB may also assist this repression, because *in vitro* it enhances the affinity of SopA for its promoter (39). The major role of SopB, however, is the formation of a partition complex with *sopC*, a series of 12 [tandem repeats](#) of 43 bp that act in *cis* to ensure partition, much like a eukaryotic [centromere](#). SopB appears to bind as a dimer to an inverted repeat sequence in each 43-mer, and to wrap *sopC* DNA in positive *supercoils* to form a characteristic complex (40, 41).

How this complex acts in partition is still unknown. Current models suppose that plasmid replicas attach as paired molecules to a host structure, such as the membrane or the division septum, and that some force then splits the pairs, pulling one copy toward each pole. SopA has an [ATPase](#) activity and is presumed to have a function in partition beyond its regulatory activity. It may link the SopB-sopC complex to the host structure, where its ATPase activity could drive the separation of paired plasmid copies (42). Our ignorance of the host components with which the Sop complex interacts is presently the major barrier to understanding partition. Host mutations that influence mini-F stability have been located in the *gyrA*, *topA*, and *ugpA* genes (43-45). While the [DNA topology](#) function or

membrane localization of these proteins is suggestive, there is as yet no good explanation for their roles in partition.

Recent refinements in cytochemical methods applied to bacterial cells have made possible the localization of plasmid molecules as fluorescent foci, allowing the movement of mini-F and other plasmids to be followed in individual cells and populations (20, 46). The data thus far published indicate a preferential location at mid-cell (the site of the next division septum) for mini-F (Fig. 3). After the focus doubles, presumably reflecting replication of the plasmid, the two new foci appear to move to the 1 / 4 and 3 / 4 positions, the midpoints of the emerging daughter cells, at a speed suggestive of an active, directed mechanism. Partition-defective mini-F shows a more random distribution and is often present close to the cell poles (20). These results are consistent with old observations of the paucity of *sop*<sup>+</sup> F plasmids in anucleate cells (“minicells”) created by division near the cell poles and with the contrasting abundance of *sop* mutant plasmids (47).

**Figure 3.** Positioning and displacement of mini-F during the cell cycle, based on microscopic examination of cells labeled with fluorescent probes (20, 46). The 1 / 4 and 3 / 4 positions in older cells become the 1 / 2 positions in new cells. Note that the smallness of the interval between the appearance of two mini-Fs and their appearance at the new cell mid-poles implies active displacement (double-headed arrow).



## 5. Resolution

The stability of DNA species that replicate as q-shaped molecules is under continual threat from the formation of dimers by homologous recombination, because this decreases the number of separate molecules available for partition and reduces initiation frequency in dimer-carrying progeny. Most

such replicons have solved this problem by acquiring a site-specific recombination system that resolves dimers back to monomers. *RepFIA* is one of these. The product of the *resD* gene is necessary for the resolution, which takes place within a sequence containing an inverted repeat of 10 bp, the *rfsF* site (48, 49). ResD is presumed to be the **resolvase**, though *in vitro* evidence is lacking.

## 6. Addiction

F possesses another type of stability function that, instead of ensuring that each daughter cell inherits a plasmid copy, kills those that fail to. The host population is thus coerced into maintaining the plasmid, a situation likened to addiction (50). The *ccd* (controlled cell death) system turned out to be the first of many examples of a widespread “stable-toxin-unstable-antidote” strategy (51). The CcdB protein inhibits an essential cellular function, DNA gyrase, but in F<sup>+</sup> cells it is kept in check by complexing with the relatively unstable antidote protein, CcdA (52, 53). In an F<sup>-</sup> segregant cell, the CcdA is lost by rapid **proteolysis** and cannot be replaced, and the CcdB liberated can attack gyrase and thereby kill the cell.

Two other toxin–antidote systems, *srn* and *flm*, are carried by F. SrnB promotes the degradation of **ribosomal RNA** and **transfer RNA** by damaging the inner membrane and allowing entry of **ribonuclease I** (54), but its synthesis is normally inhibited by an unstable **anti-sense RNA**, SrnC (55). As with the Ccd system, loss of the plasmid causes a rapid decline in the antidote SrnC, and the liberated and stabilized *srnB* **messenger RNA** is then translated to produce the toxin. The locus is situated next to *repFIC::Tn1000* and was undoubtedly once part of a viable replicon, similar to the R1 plasmid whose *hok/sok* killer system it resembles. *Flm* is another *hok/sok* analogue. Its toxin (*flmA*) and antidote RNA (*flmB*) genes are located in the leading region (56). It is tempting to suggest that a burst of FlmA synthesis in the recipient modifies the membrane to facilitate conjugational transfer, but the actual relevance of *flm* is unknown.

Easy isolation of F<sup>-</sup> and mini-F<sup>-</sup> segregants suggests that these toxin–antidote systems are inefficient. This may, however, be the result of unintended laboratory selection for tolerant strains during decades of genetic experimentation. The significance of such addictive systems in the wild remains to be assessed.

## 7. Conjugal Aids

Two other loci in the leading region encode functions for which a role in conjugation can be envisaged. The *psiB* gene product is transiently expressed following transfer to a recipient cell, where it inhibits the **SOS response**, probably by preventing the activation of **RecA** protein (57). Its role might be to prevent induction of inappropriate SOS functions in the recipient. The *ssf* gene, an analogue of *E. coli* *ssb*, encodes a **single-stranded DNA binding protein** that might protect DNA entering the cell prior to its replication (58). The absence of leading region functions does not, however, impede DNA transfer or recombinant formation, at least in laboratory conditions, and their significance is thus unclear.

## 8. Phage Exclusion

Following infection of an F<sup>+</sup> cell by certain phages dubbed “female-specific”, (eg, T7 and FII), expression of phage genes begins and the cell dies, but phage development is arrested. The F locus responsible is *pif* (phage inhibition by F), an operon consisting of two genes, *pifC* and *pifA*, located between the *repFIA* and *repFIB* regions. The phenotype is not understood. It appears to be due largely to interactions between PifA and the products of phage genes 1.2 and 10, but it does not seem to be related to the normal activities of these genes. These interactions somehow interrupt entry of phage DNA and cut short expression of its late genes (59). PifC represses *pif* operon transcription by binding to a short inverted repeat sequence overlapping the *pif* promoter –10 region (60). The

participation of PifC in mini-F replication initiated at *oriI* has also been proposed, on the basis of the defective replication phenotypes of *pifC* **amber mutants** (61). It is also possible, however, that the phenotype is an indirect effect of PifC amber fragments binding to the *pif* promoter.

## 9. Perspectives

The F plasmid as a tool has largely retired from the front line of molecular genetics, to be replaced by sleeker, more powerful technologies. Nevertheless, its properties continue to be exploited, eg, its stability and low copy number make it suitable as a **vector** for [cloning](#) large fragments of eukaryotic DNA [the BAC vectors (62)]. As a biological entity, on the other hand, F has retained its mystery. We understand only dimly the control of its replication, its mechanism of partition, and the relationship of both processes to the host cell cycle, not only for F, but for all low-copy-number replicons. Finally, although whole-genome sequencing has not yet proved equal to the task of completing the sequence of F, this project is underway (R. Skurray and L. Frost, personal communications) and may well provide some surprises.

## Bibliography

1. W. Hayes (1953) *J. Gen. Microbiol.* **8**, 72–88.
2. L. L. Cavalli, J. Lederberg, and E. M. Lederberg (1953) *J. Gen. Microbiol.* **8**, 89–103.
3. F. Cuzin and F. Jacob (1967) *Ann. Inst. Pasteur* **112**, 397–418.
4. Y. Hirota (1960) *Proc. Natl. Acad. Sci. USA* **46**, 57–64.
5. Y. Nishimura, L. Caro, C. M. Berg, and Y. Hirota (1971) *J. Mol. Biol.* **55**, 441–456.
6. K. von Meyenberg and F. G. Hansen (1980) In *Mechanistic Studies of DNA Replication and Genetic Recombination*, ICN-UCLA Symposium on Molecular and Cellular Biology (B. Alberts and C. F. Fox, eds.), Academic Press, New York, pp. 137–159.
7. K. Timmis, F. Cabello, and S. N. Cohen (1975) *Proc. Natl. Acad. Sci. USA* **72**, 2242–2246.
8. M. A. Lovett and D. R. Helinski (1976) *J. Bacteriol.* **127**, 982–989.
9. D. Lane and R. C. Gardner (1979) *J. Bacteriol.* **139**, 141–151.
10. P. L. Bergquist, H. E. D. Lane, L. Malcolm, and R. A. Doward (1982) *J. Gen. Microbiol.* **128**, 223–238.
11. R. Eichenlaub, D. Figurski, and D. R. Helinski (1977) *Proc. Natl. Acad. Sci. USA* **74**, 1138–1141.
12. J. J. Manis and B. C. Kline (1977) *Mol. Gen. Genet.* **152**, 175–182.
13. D. Figurski, R. Kolter, R. Meyer, M. Kahn, R. Eichenlaub, and D. R. Helinski (1978) In *Microbiology—1978* (D. Schlessinger, ed.), ASM, Washington, D.C., pp. 105–109.
14. T. Murotsu, H. Tsutsui, and K. Matsubara (1984) *Mol. Gen. Genet.* **196**, 373–378.
15. K. Muraiso, T. Tokino, T. Murotsu, and K. Matsubara (1987) *Mol. Gen. Genet.* **206**, 519–521.
16. P. L. Bergquist, R. A. Downard, P. A. Caughey, and H. E. D. Lane (1981) *J. Bacteriol.* **147**, 888–899.
17. H. Tsutsui, A. Fujiyama, T. Murotsu, and K. Matsubara (1983) *J. Bacteriol.* **155**, 337–344.
18. T. Ogura and S. Hiraga (1983) *Cell* **32**, 351–360.
19. C. E. Helmstetter, M. Thornton, P. Zhou, J. A. Bogan, A. C. Leonard, and J. E. Grimwade (1997) *J. Bacteriol.* **179**, 1393–1399.
20. H. Niki and S. Hiraga (1997) *Cell* **90**, 951–957.
21. J. D. Keasling, B. O. Palsson, and S. Cooper (1992) *Res. Microbiol.* **143**, 541–548.
22. T. Murotsu, K. Matsubara, H. Sugisaki, and M. Takanami (1981) *Gene* **15**, 257–271.
23. D. Bramhill and A. Kornberg (1988) *Cell* **54**, 915–918.
24. Y. Kawasaki, F. Matsunaga, Y. Kano, T. Yura, and C. Wada (1996) *Mol. Gen. Genet.* **253**, 42–49.

25. H. Tsutsui and K. Matsubara (1981) *J. Bacteriol.* **147**, 509–516.
26. S. K. Pal and D. K. Chattoraj (1988) *J. Bacteriol.* **170**, 3554–3560.
27. M. J. McEachern, M. A. Bott, P. A. Tooker, and D. R. Helinski. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 7942–7946.
28. L. Sogaard-Andersen, L. A. Rokeach, and S. Molin (1984) *EMBO J* **3**, 257–262.
29. M. Ishiai, C. Wada, Y. Kawasaki, and T. Yura (1994) *Proc. Natl. Acad. Sci. USA* **91**, 3839–3843.
30. J. D. Trawick and B. C. Kline (1985) *Plasmid* **13**, 59–69.
31. L. Masson and D. S. Ray (1986) *Nucleic Acids Res.* **14**, 5693–5711.
32. Y. Kawasaki, C. Wada, and T. Yura (1990) *Mol. Gen. Genet.* **220**, 277–282.
33. M. Ishiai, C. Wada, Y. Kawasaki, and T. Yura (1992) *J. Bacteriol.* **174**, 5597–5603.
34. A. L. Abeles and S. J. Austin (1991) *Proc. Natl. Acad. Sci. USA* **88**, 9011–9015.
35. G. Mukhopadhyay, S. Sozhamannan, and D. K. Chattoraj (1994) *EMBO J* **13**, 2089–2096.
36. A. Blasina, B. L. Kittel, A. E. Toukdarian, and D. R. Helinski (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3559–3564.
37. A. Miron, I. Patel, and D. Bastia (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6438–6442.
38. H. Mori, A. Kondo, A. Ohshima, T. Ogura, and S. Hiraga (1986) *J. Mol. Biol.* **192**, 1–15.
39. H. Mori, Y. Mori, C. Ichinose, H. Niki, T. Ogura, A. Kato, and S. Hiraga (1989) *J. Biol. Chem.* **264**, 15535–15541.
40. D. P. Biek and J. Shi (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8027–8031.
41. A. S. Lynch and J. C. Wang (1994) *J. Mol. Biol.* **236**, 679–684.
42. M. Motallebi-Veshareh, D. A. Roach, and C. M. Thomas (1990) *Mol. Microbiol.* **4**, 1455–1463.
43. T. Ogura, H. Niki, H. Mori, M. Morita, M. Hasegawa, C. Ichinose, and S. Hiraga (1990) *J. Bacteriol.* **172**, 1562–1568.
44. C. A. Miller, S. L. Beaucage, and S. N. Cohen (1990) *Cell* **62**, 127–133.
45. B. Ezaki, H. Mori, T. Ogura, and S. Hiraga (1990) *Mol. Gen. Genet.* **223**, 361–368.
46. G. S. Gordon, D. Sitnikov, C. D. Webb, A. Teleman, A. Straight, R. Losick, A. W. Murray, and A. Wright (1997) *Cell* **90**, 1113–1121.
47. J. E. Hogan, B. C. Kline, and S. B. Levy (1982) *Plasmid* **8**, 36–44.
48. D. Lane, R. de Feyter, M. Kennedy, S-H. Phua, and D. Semon (1986) *Nucleic Acids Res.* **14**, 9713–9728.
49. M. B. O'Connor, J. J. Kilbane, and M. H. Malamy (1986) *J. Mol. Biol.* **189**, 85–102.
50. H. Lehnerr, E. Maguin, S. Jafri, and M. B. Yarmolinsky (1993) *J. Mol. Biol.* **233**, 414–428.
51. A. Jaffé, T. Ogura, and S. Hiraga (1985) *J. Bacteriol.* **163**, 841–849.
52. P. Bernard and M. Couturier (1992) *J. Mol. Biol.* **226**, 735–745.
53. S. Maki, S. Takiguchi, T. Miki, and T. Horiuchi (1992) *J. Biol. Chem.* **267**, 12244–12251.
54. R. Ito and Y. Ohnishi (1983) *Biochim. Biophys. Acta* **20**, 27–34.
55. A. K. Nielsen, P. Thorsted, T. Thisted, E. G. Wagner, and K. Gerdes (1991) *Mol. Microbiol.* **15**, 1961–1973.
56. S. M. Loh, D. S. Cram, and R. A. Skurray (1988) *Gene* **66**, 259–268.
57. M. Bagdasarian, A. Bailone, J. F. Angulo, P. Scholz, M. Bagdasarian, and R. Devoret (1992) *Mol. Microbiol.* **6**, 885–893.
58. A. L. Kolodkin, M. A. Capage, E. I. Golub, and K. B. Low (1983) *Proc. Natl. Acad. Sci. USA* **80**, 4422–4426.
59. C. K. Schmitt and I. J. Molineux (1991) *J. Bacteriol.* **173**, 1536–1543.
60. J. F. Miller and M. H. Malamy (1986) *Proc. Natl. Acad. Sci. USA* **83**, 1433–1437.



61. R. Eichenlaub and H. Wehlmann (1980) *Mol. Gen. Genet.* **180**, 201–204.
62. H. Shizuya, B. Birren, U. J. Kim, V. Mancino, T. Slepak, Y. Tachiiri, and M. Simon (1992) *Proc. Natl. Acad. Sci. USA* **89**, 8794–8797.

### Suggestions for Further Reading

63. N. Willets and R. Skurray (1987) "Structure and function of the F factor and mechanism of conjugation". In *Escherichia coli and Salmonella. Cellular and Molecular Biology*, 1st ed. (F. C. Neidhart et al., eds.), ASM Press, Washington D.C., pp. 1110–1133. The last comprehensive treatment of the F plasmid, providing all the essential information and references.
64. D. K. Chattoraj and T. D. Schneider (1997) Replication control of plasmid P1 and its host chromosome: the common ground. *Prog. Nucleic Acid Res. Mol. Biol.* **57**, 145–185. An up-to-date treatment of replication control that takes P1 as the reference point but generalizes the discussion to discuss fully other plasmids, such as F, with similar replicons.
65. S. Hiraga (1992) Chromosome and plasmid partition in *Escherichia coli*. *Annu. Rev. Biochem.* **61**, 283–306. A general treatment of partition in *E. coli*; now slightly dated, but it has all the basics.

## Factor Xa

Factor X is a component of the [blood clotting](#) system. It circulates in blood plasma (8 µg/mL) as an inactive precursor of a [serine proteinase](#) (1). When the blood coagulation process is initiated, factor X is converted to factor Xa (activated factor X), which forms part of a complex that converts another serine proteinase precursor, pro-thrombin, into the active enzyme [thrombin](#). Clotting occurs when thrombin converts fibrinogen to **fibrin**.

Factor X is synthesized in the liver as a 488-amino-acid residue precursor protein, pre-pro-X, the primary translation product of the factor X gene (2). Forty amino acids are cleaved from the amino terminus of pre-pro-X. The resultant 448-residue protein is modified by the attachment of sugars and, in a vitamin-K-dependent process, by the carboxylation of certain [glutamic acid](#) residues to form [g-carboxy glutamic acid](#) residues. It is cleaved by a specific proteinase that removes residues 140 to 142, so factor X circulates as a two-chain, **disulfide bond**-linked, enzymatically inactive protein. This heavy-chain–light-chain dimeric structure is typical of the blood coagulation proteins (3).

As might be expected, activation of blood coagulation is a carefully controlled process involving many clotting factors (see [Blood Clotting](#)). All the details of this process are not fully understood as yet, but it is known that factor X can be converted to an active enzyme by two distinct pathways, the *intrinsic* and the *extrinsic* (4). The former pathway involves activation by factor IXa, which together with factor VIIIa forms a **receptor** on the phospholipid surface of blood platelets. When factor X binds to this receptor complex, a 52-residue peptide is cleaved from the amino terminus of its heavy chain to form factor Xa. The latter pathway involves activation by factor VIIa which, in an analogous fashion, forms a platelet-bound complex with tissue factor. When factor X binds to this complex, the same peptide is cleaved from the heavy chain to give factor Xa. Calcium ions, which bind to the g-carboxy glutamic acid residues of the light chain (and also to the other clotting factors), are essential for factor X activation.

Once factor X is activated, it binds along with factor Va on the surface of platelets, where it interacts

with pro-thrombin to convert this zymogen into thrombin, the end product of the so-called coagulation proteinase cascade. Factor X was first recognized as an essential part of this cascade when it was shown that it could restore coagulability to the plasma of individuals who had a specific hemorrhagic disorder. It was given the name *Stuart factor*, which was subsequently changed to factor X in an effort to systematize the nomenclature of the clotting factors.

#### Bibliography

1. M. Hertzberg (1994) *Blood Rev.* **8**, 56–62.
2. C. Miao, S. P. Leytus, D. W. Chung, and E. W. Davie (1992) *J. Biol. Chem.* **267**, 7395–7401.
3. E. W. Davie, K. Fujikawa, M. E. Legaz, and H. Kato (1975) In *Proteases and Biological Control* (E. Reich, D. B. Rifkin, and E. Shaw, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 65–77.
4. M. P. McGee and L. C. Lee (1991) *J. Biol. Chem.* **266**, 8079–8085.

#### Suggestions for Further Reading

5. U. Hedner and E. W. Davie (1989) "Introduction to hemostasis and the vitamin K-dependent coagulation factors". In *The Metabolic Basis of Inherited Disease* (C. R. Scriver, A. L. Beaudet, W. S. Sly, and D. Valle, eds.), McGraw-Hill, New York, pp. 2107–2134.
6. G. J. Broze Jr. (1995) Tissue factor pathway inhibitor and the revised theory of coagulation. *Annu. Rev. Med.* **46**, 103–112.

#### Facultative Heterochromatin

Facultative heterochromatin is the term describing regions of [chromosomes](#) that appear **heterochromatic** at certain times in the [cell cycle](#) but are not always this way. This is contrasted with constitutive heterochromatin, which is always condensed and stains strongly throughout the cell cycle ([1](#)). Facultative heterochromatin represents domains of [euchromatin](#) that are condensed and inactive in a particular cell. An excellent example of the assembly of facultative heterochromatin is found in the inactive [X-chromosome](#) (see [Barr Body](#), [X-Chromosome Inactivation](#)). In the inactive X-chromosome, the [histones](#) are hypoacetylated, DNA is **methyalted**, and the chromosome replicates late in **S-phase**. All of these traits are characteristic of inactive [chromatin](#). Importantly, these are not irreversible modifications and can be reversed at certain points in development (see [Random X-Inactivation](#)).

The various characteristics of facultative heterochromatin effectively work together to stabilize a **transcriptionally** repressed state. The generation of [antibodies](#) against **acetylated** histones has allowed a number of general correlations concerning the possible functional roles of histone acetylation. There is also a strong correlation between histone acetylation and the transcriptional activity of chromatin. In *Saccharomyces cerevisiae* most of the [genome](#) is transcriptionally active and contains hyperacetylated core histones. Transcriptionally inactive **domains** of chromatin in **yeast**, such as the silent **mating type** cassettes and **telomeric** sequences, contain histone H4 that is hypoacetylated, except at one position, Lys12 ([2](#)). In higher eukaryotes, acetylation of histone H4 increases during the reactivation of transcription in the initially inactive chicken erythrocyte nucleus, following fusion of the erythrocyte with a transcriptionally active cultured cell, to form a heterokaryon (see [Euchromatin](#)). Histone acetylation is particularly prevalent over the specific **globin** genes that are actively transcribed in reticulocytes. More recent studies have demonstrated convincingly that histone hyperacetylation is actually restricted to the domain of chromatin that

contains the potentially active chicken b-globin gene locus (see [Domain, Chromosomal](#)). This result indicates very specific targeting of histone acetyltransferase activity. **Immunolabeling** of [polytene chromosomes](#) in *Chironomus* and *Drosophila* also reveals a nonrandom distribution of histone H4 acetylation that correlates with transcriptional activity. Within female mammals, the transcriptionally inactive X-chromosome is distinguished by a lack of histone H4 acetylation (3). Therefore several independent experimental approaches have shown that actively transcribed and potentially active chromatin domains are selectively enriched in hyperacetylated histones, whereas transcriptionally inactive chromatin contains hypoacetylated histones.

A role for specialized chromatin structures in mediating transcriptional silencing by methylated DNA, like that found in the facultative heterochromatin of vertebrates, has been suggested by several investigators. High levels of methyl-CpG correlate with transcriptional inactivity and **nuclease** resistance in vertebrate chromosomes (4). Methylated DNA transfected into mammalian cells is also assembled into nuclease-resistant structures within the assembled [minichromosomes](#), indicating the existence of unusual **nucleosomal** particles (5). These unusual nucleosomes migrate as large nucleoprotein complexes in [agarose](#) gel electrophoresis. These complexes are held together by higher order protein–DNA interactions, despite the presence of abundant **micrococcal nuclease** cleavage points within the DNA. Individual nucleosomes assembled on methylated DNA interact together more stably than on unmethylated templates.

The accessibility of chromatin to nucleases could also be affected directly by the stability with which the histones interact with DNA within the nucleosome. DNA methylation does not influence the association of core histones with the vast majority of DNA sequences in the [genome](#). However, for certain specific sequences, such as those found in the Fragile X mental retardation gene 1 promoter, methylation of CpG dinucleotides alters the positioning of histone–DNA contacts and the affinity with which these histones bind to DNA (see [Acentric Fragment](#)). The exact chromatin structure found *in vivo* can also result from gene activity. Linker histones, such as H1, are relatively deficient in the transcribed regions of genes. So it is not surprising that transcriptionally inactive facultative heterochromatin containing methyl-CpG should show an increase in the abundance of histone H1, whereas DNA sequences lacking methyl-CpG are deficient in H1. *In vitro* studies indicate that histone H1 interacts preferentially with methylated DNA under certain conditions.

There are also features of transcriptional repression dependent on methylated DNA that can be explained by methylation-specific [repressors](#) that operate more effectively within a facultative heterochromatin environment. Transcriptional repression is strongly related to the density of DNA methylation. A non linear relationship exists between the lack of repression observed at low densities of methyl CpG and repression at higher densities. These results led to the demonstration that local domains of high methyl-CpG density confer transcriptional repression on unmethylated **promoters** *in cis*. The importance of a nucleosomal infrastructure for transcriptional repression dependent on DNA methylation is shown by the observation that immediately after injection into *Xenopus* oocyte nuclei, methylated and unmethylated templates both have equivalent activity (6). As facultative heterochromatin is assembled, however, the methylated DNA is repressed with the loss of DNase I [hypersensitive sites](#) and the loss of engaged **RNA polymerase**. The requirement that nucleosomes exert efficient transcriptional repression dependent on DNA methylation can be explained in several ways. Methylation-specific repressors might recruit a corepressor complex that directs the modification of the chromatin template into a more stable and transcriptionally inert state. One potential candidate corepressor for MeCP2 is the SIN3-histone deacetylase complex, because inhibition of histone deacetylation reverses some of the transcriptional repression conferred by DNA methylation (7). Alternatively, like histone H1, methylation-specific repressors might bind more efficiently to nucleosomal rather than to naked DNA. Any cooperative interactions between molecules could propagate their association of methylation-specific repressors along the nucleosomal array even into unmethylated DNA segments. This latter mechanism is analogous to the nucleation of heterochromatin assembly at the yeast [telomeres](#) by the DNA-binding protein RAP1, which then recruits the repressors SIR3p and SIR4p that organize chromatin into a repressive structure (8) (see [Telomere](#)). All of these potential mechanisms could individually or together

contribute to assembling of a repressive facultative heterochromatin domain.

A final relevant issue is the significance of the timing of replicative initiation on facultative heterochromatin in **S-phase**. Facultative heterochromatin is normally replicated late in the S-phase. If replication disrupts both active and repressed chromatin structures, then the entire nucleus has to be remodeled after each replication. The accessibility of immature chromatin on newly replicated DNA provides a means to accomplish this remodeling. The reformation of nuclear structures, however, has other implications. If the [transcription factors](#) available in a cell are limiting, then a gene that is replicated early in the S-phase has more opportunity to assemble an active transcription complex than a gene that replicates late simply because the gene that replicates early is available for transcription factors to bind to it before all of the early replicating portion of the genome has sequestered these factors. Therefore a late-replicating gene in facultative heterochromatin experiences a relative deficiency in transcription factor availability. Transcriptionally active genes in euchromatin replicate early in the S-phase (see [Euchromatin](#)). The reason for this early replication is unknown, but possibilities include the local disruption of chromatin structure by transcription complexes, which makes that DNA more accessible to the replication machinery (9) and the observation that many transcription factors are in fact also replicative factors. An attractive variation of this model is that the type of chromatin assembled early in the S-phase is more accessible to transcription factors than chromatin assembled late in the S-phase. Early replicating chromatin may sequester histones that are more highly acetylated and consequently more accessible to the transcription factors that maintain continued transcriptional activity. The CENP-A protein that replaces histone H3 in the mammalian [centromere](#) within heterochromatin is synthesized at the very end of the S-phase, providing an example of how the cell cycle-dependent compartmentalization of **protein biosynthesis** might contribute to assembling a specialized chromatin structure (10). Although a general test of the significance of this model has not been made, it remains an attractive mechanism for explaining both the maintenance of specific patterns of gene expression in a proliferating cell type and the maintenance of domains of facultative heterochromatin.

#### Bibliography

1. S. W. Brown (1966) *Science* **151**, 417–425.
2. M. Braunstein et al. (1993) *Gene Dev.* **7**, 592–604.
3. P. Jeppesen and B. M. Turner (1993) *Cell* **74**, 281–291.
4. F. Antequera, D. Macleod, and A. Bird (1989) *Cell* **58**, 509–517.
5. I. Keshet, J. Lieman-Hurwitz, and H. Cedar (1986) *Cell* **44**, 535–543.
6. S. U. Kass, N. Landsberger, and A. P. Wolffe (1997) *Curr. Biol.* **7**, 157–165.
7. L. Alland et al. (1997) *Nature* **387**, 49–55.
8. M. Grunstein et al. (1995) *J. Cell Sci.* 519–536.
9. A. P. Wolffe and D. D. Brown (1988) *Science* **241**, 1626–1632.
10. R. D. Shelby, O. Vafa, and K. F. Sullivan (1997) *J. Cell Biol.* **136**, 501–513.

#### Suggestion for Further Reading

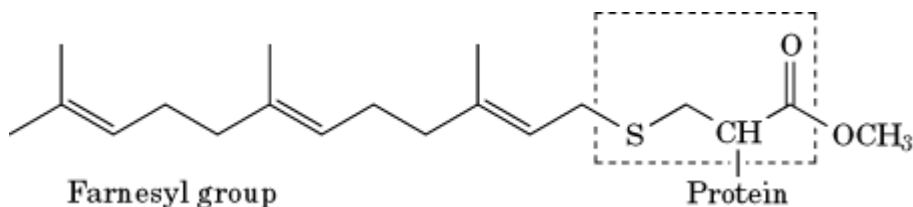
11. A. Wolffe (1998) *Chromatin: Structure and Function* 3d ed., Academic Press, London.

#### Farnesylation

Farnesylation is the process by which a [cysteine](#) residue in the C-terminal region of some eukaryotic

[proteins](#) is **posttranslationally modified** with an isoprenoid lipid (the 15-carbon farnesyl group) and the exposed carboxyl group is methylated (Fig. 1). The farnesylation and carboxylmethylation increase the affinity of the protein for the membrane and have important functional consequences. The farnesyl group is one of several lipids that act as a [membrane anchor](#) for proteins. The role of the farnesyl group and other lipids in membrane anchoring is described in more detail elsewhere [see [Prenylation](#); [Membrane Anchors](#)]

**Figure 1.** Modification of a C-terminal cysteine residue by a farnesyl group.



## Farr Assay

The Farr assay ([1](#), [2](#)) is a precipitation-based [immunoassay](#) that is especially appropriate for analysis of certain [antibody–antigen interactions](#) that do not lead to an insoluble lattice of crosslinked antibody–antigen complexes (see [Precipitin Reaction](#)). Failure to crosslink is frequently encountered with [haptens](#), small antigens, and [monoclonal antibodies](#). The Farr assay takes advantage of a common solubility property of [immunoglobulins](#). Ammonium **sulfate** at 50% saturation causes antibodies and immune complexes to precipitate, whereas free antigens and haptens often remain soluble. Measurement of the antigen-binding capacity of an uncharacterized antibody preparation can be achieved by mixing a fixed amount of [antigen](#) with serial dilutions of [antibody](#), allowing the binding reaction to reach equilibrium, precipitating free and complexed antibody with ammonium sulfate, and then measuring the antigen content of the precipitate. Antigen is precipitated quantitatively until a subsaturating dilution of antibody is reached, whereupon the antigen content of the precipitates progressively decreases. Observation of a constant amount of precipitate at saturating antibody concentrations differs from the precipitin reaction, in which excess antibody actually inhibits formation of a precipitate. The Farr assay is now obsolete, but is cited frequently in older literature.

## Bibliography

1. R. S. Farr (1958) *J. Infect. Dis.* **103**, 239–262.
2. P. Minden and R. S. Farr (1967) "The ammonium sulphate method to measure antigen binding capacity", in *Handbook of Experimental Immunology*, D. M. Weir, ed., Blackwell, Oxford, pp. 493–526.

## Fate Map

## 1. Classical Fate Maps

A fate map is a graphic representation of the cells or regions within the embryo at one point during [development](#) that will develop into specific tissues and organs at subsequent developmental stages. The concept of the fate map may be one of the earliest abstract ideas in [embryology](#), one that is implicit, for example, in the “homunculus” imagined by seventeenth century microscopists to be a miniature adult human that is “preformed” within the sperm. The concept of preformation was generally favored by scientists until the last half of the eighteenth century, when Caspar Friedrich Wolff, through observation of chick embryo development, postulated that embryos develop through a process he termed “**epigenesis**,” whereby unpatterned tissues of the early embryo reorganize and expand to form the tissues and organs of the adult. As the cellular basis for tissue organization came to be understood by the mid-nineteenth century, it eventually became clear that embryonic development could be characterized by four cellular phases common to all vertebrate and most invertebrate embryos: (i) formation of the single-celled, diploid [zygote](#) (or fertilized egg) through union of haploid sperm and egg; (ii) **cleavage**, in which the fertilized egg undergoes mitotic division without increase in size to form the multicellular **blastula** (mammalian: blastocyst); (iii) **gastrulation**, in which blastula cells rearrange and expand to form the three **primary germ layers** (ectoderm, mesoderm, and endoderm) that are organized into the basic body plan (concepts developed in the early nineteenth century by von Baer); and (iv) **organogenesis**, in which the cells within the the three primary germ layers of the embryo form the organ and tissue primordia that develop into the tissues and organs of the *fetus*. The fetal stage of development is distinguished from the embryonic phase in that the fetus is indeed a preformed miniature of the adult that contains all of the organ and tissue rudiments (in the right places and containing differentiated cells and stem cells) that need only to expand and mature to function as adult tissues and organs. The transition from embryo to fetus occurs at about the eighth week of human development.

In the late nineteenth and early twentieth centuries, embryologists began to take an experimental approach to the study of development. Fate maps of relatively simple embryos, such as ascidians, were initially established by directly observing the fate of blastula cells during gastrulation ([1](#)). The emergence of dyes that could be applied to living tissues without damaging them (“vital dyes”) permitted selective staining of parts of blastula embryos and subsequent tracing of the tissues of the embryo that developed from these stained parts. Systematic application of this procedure to trace the fate of different regions of amphibian blastulas through gastrulation led to the construction of the first detailed and comprehensive fate map of a vertebrate embryo ([2](#), [3](#)). Subsequent work by many individuals led to the construction of fate maps for organ rudiment primordia of a variety of chordate embryos (amphioxus, fish, amphibians, birds) ([4-9](#)). These maps continue today as the primary basis for our present day understanding of **morphogenetic** movements ([2](#)) that blastula cells undergo during gastrulation. Nevertheless, these maps are constantly being revised and refined by applying more modern techniques. For recent examples, see fish ([10](#), [11](#)); birds ([12-15](#)); and amphibians ([16-18](#)).

Although fate maps, such as those described previously, give a general idea how regions of the blastula give rise to the primary organ primordia of the embryo, they do not yield precise details of the way cells within those primordia give rise to the more complex tissues and organs of the fetus and adult. These events could not be adequately followed by direct observation because of the massive expansion of cells during the embryonic phase of development. Moreover, during the ensuing cell expansion and migrations, vital dyes became diluted beyond the level of detectability. The latter problem was only partially circumvented by using microscopic carbon particles to follow cell movements and rearrangements ([19](#), [20](#)).

In the late 1960s, however, Nicole Le Douarin developed a powerful method for embryonic fate mapping that allowed following the mapping of embryonic tissues indefinitely, even into full adulthood. Briefly, she discovered that quail cells could be distinguished from chicken cells by a

[nucleolus](#) in the quail cell that developed a brilliant magenta color when stained with the Feulgen reagent (21, 22). By transplanting homologous tissues from quail embryos into chick embryos, the fate of the transplanted quail cells could be established unequivocally by histological methods. Through broad application of this method, the development of the avian embryo has been mapped in detail (see Ref. 23 for many examples) Much of this newer map supports and confirms earlier work. Nevertheless, this new method brilliantly illuminated the many limitations and previous inaccuracies of earlier methods. This is particularly true for migratory mesenchyme cells within the embryo, such as the neural crest, which, as revealed by LeDouarin and co-workers, populate a previously unknown variety of tissues and organs (24-28). In addition, the embryonic development of other mobile cell populations, such as blood cell lineages (29-34) and muscle precursor cells (35-37) and many other other tissues, were extended and clarified in ways not possible previously.

## 2. Mammalian Embryo Fate Maps

Fate mapping in mammalian embryos progressed much more slowly because of the relative inaccessibility of these embryos within the mother's uterus. Therefore, the fate maps in most existing textbooks are derived from the avian embryo based on assumption that mammalian development is essentially identical. Because of the recent development of methods for *in vitro* culture of early mammalian embryos, fate mapping of early mammalian development has begun. One such fate map experiment was conducted by transplanting fragments of tissue from early mouse embryos **radiolabeled** with tritiated thymidine in unlabeled embryos, followed by localizing the labeled tissues by [autoradiography](#) (38). The recent development of highly **fluorescent** dyes that can be injected into cells has greatly facilitated fate mapping of early mouse embryo tissues (39-43) To date, the most important outcome of such studies has been to establish that existing mammalian fate maps based upon the avian model are indeed valid by and large. Therefore, it is no exaggeration to state that the methods and strategies developed by Le Douarin and the resulting detailed fate maps of the avian embryo provide the most accurate and comprehensive information garnered during the twentieth century about the cellular movements underlying vertebrate embryonic development, including mammalian and human development.

## 3. Embryonic Fate and Developmental Potential

At the onset of development, the single-celled zygote can form all of the cell and tissue types of the adult, and this potential is apportioned during embryogenesis. The relationship between developmental potential and developmental fate is the central question that developmental biologists have been trying to understand for more than a century [the reader is directed to a recent and excellent review and discussion by Viktor Hamburger of the early experiments in this area (44)]. Wilhelm Roux first addressed this problem experimentally in 1883 by ablating one blastomere of an amphibian embryo. He found that the surviving blastomere formed an incomplete embryo, leading him to conclude that fate and potential were apportioned in concert. Although this initial conclusion was flawed, Roux's novel and experimental approach laid the foundations for what came to be known as **experimental embryology**. It was not long until it was experimentally established that developmental potential is generally more widely dispersed within the early embryo than developmental fate. Thus, a specific region of the early embryo can develop into several different tissues and organs, even though its fate represents only one of those possible outcomes. As the embryo matures, the range of potentiality of these regions becomes restricted until, ultimately, potential and fate become equal, and the tissue is deemed determined (45). The mechanism(s) that control fate and potential constitute one of the central mysteries of developmental biology. The seminal experiment that established the profound depth of this mystery was reported in 1924 by Hilda Mangold, a graduate student in the laboratory of Hans Spemann. She showed that fragments of specific regions from one amphibian gastrula could induce formation of a complete second embryo when transplanted into the prospective belly region of a host gastrula (46) (an English translation of this article can be found in Ref. 47). The fact that the transplanted fragment could reprogram host tissue, fated to form abdominal structures, to form head tissues, including brain, led to the "[organizer](#)" hypothesis, in which some embryonic tissues play central roles to induce and direct the

development of neighboring tissues. That discovery provided the conceptual basis for virtually all developmental biology in the twentieth century. Despite the advances in the cellular and molecular analysis of development in the last two decades of this century, the mechanistic bases for organizer function, morphogenesis, and cell determination in vertebrates remain mysteries.

#### 4. Genetic Control of Developmental Fate

By contrast, however, research into mechanisms by which genes control embryonic development has progressed most rapidly, using the fruit fly *Drosophila*, because of the advanced state of genetic analysis in this organism. Although a fate map of the early *Drosophila* embryo was available by mid-twentieth century (48), a coherent hypothesis regarding the role of specific genes in establishing the *Drosophila* fate map was first proposed in 1980 (49-51) (see Ref. 52 for a broad discussion of the underlying embryological issues). This breakthrough in studying the relationship between embryonic fate maps and genetics has revolutionized the developmental biology of *Drosophila*. A large proportion of the **genes** that act developmentally in flies have vertebrate **homologues** that play key, but typically different, developmental roles. Thus, a definitive model for the genetic control of embryonic fate maps and pattern formation in vertebrates remains to be established.

#### 5. The Difficult Relationship Between Fate Maps and Vertebrate Genetic Expression Patterns

Molecular probes for the expressed products of genes developed during the last half of the twentieth century have allowed superimposing fate maps and gene expression maps. The first use of such a technique can be traced back to at least 1957, when Holtzer and Detwiler used a fluorescent antimyosin **antibody** to localize the earliest appearance of differentiated muscle cells within the somite myotomes of chick embryos (53). This general approach expanded enormously in recent years and now includes **in situ hybridization**, which allows detecting the **RNA** products of genes. Thus the developmental expression maps of a large number of genes of both known and unknown function are now known for development of a wide range of organisms. In instances in which the patterns derived from these two types of maps coincide, it is often assumed that the molecular expression accurately reflects that of cells whose fate is known. Although this may be generally true, the assumption carries potential pitfalls of which one should be wary, at least when unaccompanied by appropriate experimental validation. Migratory cells, for example, often intermingle transiently with and are indistinguishable from nonmigratory cells. The signal from the migratory cells in such situations can be discerned if they carry an independent lineage marker, such as the quail nucleolus referred to previously, that allows distinguishing them from their nonmigratory counterparts (54). On the other hand, some gene products (particularly RNA) can appear and be degraded extremely rapidly in patterns that reflect transcriptional waves passing through specific tissues. Pourquie and co-workers discovered such transcriptional waves and distinguished them from cell migration by demonstrating discordance over time between the pattern of injected fluorescent dyes and gene expression signals (55). Finally, cells may transiently express one gene product at one period of development and then switch to expression of another. Cepko and co-workers have developed techniques in which infective **retroviruses** are used to follow transient expression patterns of individual cells and their progeny through complex developmental pathways (56-59). These examples serve only to illustrate that gene expression patterns can only be assumed to reflect cell lineage or fate maps until experimentally verified by independent means.

#### 6. Transgenic Strategies for Fate Mapping

Molecular techniques are now being used to map developmental fate and potential and also to manipulate them experimentally. Retroviruses are widely used to deliver bioactive genes to influence development positively or negatively through the expression of exogenous gene products (for examples, see Refs. 60, 61). The emergence of methods for creating germ-line **transgenic** mice has led to a variety of new and imaginative approaches to analyzing of the relationship between cell fate and gene expression in mammalian development. Buckingham and co-workers have “knocked-in” a **beta-galactosidase** marker gene into a gene locus, which that controls the onset of myogenesis, to



locate newly “born” muscle precursor cells in the mouse embryo ([62](#), [63](#)). Recently, transgenes that must undergo recombination signals to be expressed in a cell- or tissue-specific fashion have been used by Nicolas and co-workers to track cell lineage progression in early mouse development ([64-66](#)). It is important to keep in mind that, because each of these techniques depends on cell- or tissue-specific gene regulation (which is itself only beginning to be understood), each is subject to the caveats alluded to previously for the relationship between cell fate and gene expression. Nevertheless, the power of these new techniques augers well for the future analysis of mechanisms underlying the apportionment of cell fate and potential during embryonic development.

## Bibliography

1. E. G. Conklin (1905) *J. Acad. Nat. Sci. (Philadelphia)* **2**, 13.
2. W. Vogt (1929) *Roux' Arch.* **120**, 385–706.
3. W. Vogt (1925) *Roux' Arch.* **106**, 542–610.
4. R. Wetzel (1929) *Roux' Arch.* **119**, 118–321.
5. E. G. Conklin (1932) *J. Morphol.* **54**, 69–118.
6. J. M. Oppenheimer (1936) *J. Exp. Zool.* **72**, 409–437.
7. J. Pasteels (1940) *Biol. Rev.* **15**, 59–106.
8. D. Rudnick (1943) *Q. Rev. Biol.* **19**, 187.
9. J. M. Oppenheimer (1947) *Q. Rev. Biol.* **22**, 105–118.
10. K. Woo and S. E. Fraser (1995) *Development* **121**, 2595–2609.
11. J. Shih and S. E. Fraser (1995) *Development* **121**, 2755–2765.
12. H. Eyal-Giladi and S. Kochav, (1976) *Dev. Biol.* **49**, 321–337.
13. M. A. J. Selleck and C. D. Stern (1991) *Development* **112**, 615–626.
14. Y. Hatada and C. D. Stern (1994) *Development* **120**, 2879–2890.
15. G. C. Schoenwolf, V. Garcia-Martinez, and M. S. Dias (1992) *Developmental Dynamics* **193**, 235–248.
16. S. B. Minsuk and R. E. Kelle (1996) *Developmental Biol.* **174**, 92–103.
17. T. Elul, M. A. Koehl, and R. Keller (1997) *Developmental Biol.* **191**, 243–258.
18. R. Harland and J. Gerhart (1997) *Annu. Rev. Cell Developmental Biol.* **13**, 611–667.
19. N. T. Spratt Jr. (1946) *J. Exp. Zool.* **103**, 259–304.
20. R. C. Fraser (1954) *J. Exp. Zool.* **126**, 349–399.
21. N. M. Le Douarin (1969) *Bull. Biol. Fr. Belg.* **103**, 435–452.
22. N. Le Douarin (1973) *Exp. Cell Res.* **77**, 459–468.
23. N. M. Le Douarin and M. A. Teillet, eds. (1984) *Chimeras in Developmental Biology*, Academic Press, London, San Diego.
24. N. M. Le Douarin and M. A. Teillet (1973) *J. Embryol. Exp. Morphol.* **30**(1), 31–48.
25. N. M. Le Douarin (1975) *Birth Defects Original Article Series*, **11**(7), 19–50.
26. J. Fontaine, C. Le Lievre, and N. M. Le Douarin (1977) *Gen. Comp. Endocrinol.* **33**(3), 394–404.
27. C. S. Le Lievre and N. M. Le Douarin (1975) *J. Embryol. Exp. Morphol.* **34**(1), 125–154.
28. J. M. Polak et al. (1974) *Histochemistry* **40**(3), 209–214.
29. E. Houssaint, M. Belo, and N. M. Le Douarin (1976) *Developmental Biol.* **53**(2), 250–264.
30. F. V. Jotereau, E. Houssaint, and N. M. Le Douarin (1980) *Eur. J. Immunol.* **10**(8), 620–627.
31. N. M. Le Douarin and F. V. Jotereau (1975) *J. Exp. Med.* **142**(1), 17–40.
32. N. M. Le Douarin et al. (1975) *Proc. Nat. Acad. Sci. USA* **72**(7), 2701–2705.
33. N. M. Le Douarin, E. Houssaint, and F. Jotereau (1977) *Adv. Exp. Med. Biol.* **88**, 29–37.
34. N. M. Le Douarin, F. Dieterlen-Lievre, and P. D. Oliver (1984) *Am. J. Anat.* **170**(3), 261–299.

35. B. Christ, H. J. Jacob, and M. Jacob (1974) *Experientia* **30**(12), 1446–1449.
36. B. Christ, H. J. Jacob, and M. Jacob (1977) *Anat. Embryol.* **150**(2), 171–186.
37. A. Chevallier, M. Kieny, and A. Mauger (1977) *J. Embryol. Exp. Morphol.* **41**, 245–258.
38. P. P. L. Tam and R. S. P. Beddington (1987) *Development* **99**, 109–126.
39. K. A. Lawson, J. J. Meneses, and R. A. Pederson (1991) *Development* **113**, 891–911.
40. M. Parameswaran and P. P. L. Tam (1995) *Dev. Genet.* **17**, 16–28.
41. P. P. L. Tam and S. X. Zhou (1996) *Developmental Biol.* **178**, 124–132.
42. P. P. L. Tam and R. R. Behringer (1997) *Mech. Dev.* **68**, 3–25.
43. P. P. L. Tam et al. (1997) *Development* **124**, 1631–1642.
44. V. Hamburger (1988) *The Heritage of Experimental Embryology: Hans Spemann and the Organizer*, Monographs on the History and Philosophy of Biology (R. B. R. Burian Jr., R. Lewontin, and J. M. Smith, eds.), Oxford University Press, New York and Oxford.
45. H. Spemann (1938) *Embryonic Development and Induction* Yale University Press, New Haven, CT.
46. H. Spemann and H. Mangold (1924) *Roux' Archiv.* **100**, 599–638.
47. B. H. Willier and J. Oppenheimer (1974) *Foundations of Experimental Embryology*, 2nd ed., Hafner, New York.
48. D. F. Poulson (1950) *Biology of Drosophila* (M. Demerec, ed.), Hafner, New York, pp. 168–270.
49. C. Nusslein-Volhard and E. Wieschaus (1980) *Nature* **287**, 795–801.
50. S. Roth, D. Stein, and C. Nusslein-Volhard (1989) *Cell* **59**, 1189–1202.
51. C. Nusslein-Volhard and S. Roth (1989) *Ciba Foundation Symposium* **144**, 37–55; discussion 55–64; 92–98.
52. C. Nusslein-Volhard (1979) In *Determinants of Spatial Organization* (S. S. a. I. R. Konigsberg, ed.), New York, Academic Press, pp. 185–211.
53. H. Holtzer, J. M. Marshall, and H. Finck (1957) *J. Biophys. Biochem. Cytol.* **3**, 705–729.
54. B. A. Williams and C. P. Ordahl (1994) *Development* **120**, 785–796.
55. I. Palmeirin et al. (1997) *Cell* **91**, 639–648.
56. D. L. Turner, E. Y. Snyder, and C. L. Cepko (1990) *Neuron* **4**(6), 833–845.
57. J. Price, D. Turner, and C. Cepko (1987) *Proc. Natl. Acad. Sci. USA* **84**(1), 156–160.
58. C. Cepko et al. (1995) *Methods Enzymol.* **254**, 387–419.
59. J. A. Golden, S. C. Fields-Berry, and C. L. Cepko (1995) *Proc. Natl. Acad. Sci. USA* **92**(12), 5704–5708.
60. C. P. Austin et al. (1995) *Development* **121**(11), 3637–3650.
61. D. J. Goff and C. J. Tabin (1997) *Development* **124**, 627–636.
62. S. Tajbakhsh et al. (1996) *Developmental Dynamics* **206**, 291–300.
63. S. Tajbakhsh, D. Rocancourt, and M. Buckingham (1996) *Nature* **384**, 266–270.
64. J. F. Nicolas, L. Mathis, and C. Bonnerot (1996) *Development* **122**, 2933–2946.
65. L. Mathis et al. (1997) *Development* **124**, 4089–4104.
66. D. L. Zinyk et al. (1998) *Curr. Biol.* **8**, 665–668.

## Fatty Acids

The major fatty acids of plants and animals contain even numbers (14 to 26) of aliphatic carbon atoms in straight chains, with a terminal carboxylic acid group. The acids may contain up to six nonconjugated *cis*-double bonds, generally separated by two single bonds and a single methylene group. Several shorthand nomenclatures are in use to indicate the number of carbon atoms they contain and the number and position of the double bonds; for example, linoleic acid may be written either as 18:2(*n* – 6) or 18:2w6, where the number of carbon atoms in the chain (*n* = 18) is followed by a colon and the number of double bonds (two). The position of the double bonds is set by the number of carbon atoms from the last double bond to the terminal (or w) methyl group (6); in this example, the two double bonds are between carbons 12–13 and 9–10.

Fatty acids appear as the free acids only when bound to [serum albumin](#) in blood. They occur in cell membranes primarily as the hydrophobic moiety in phospholipids and other membrane lipids as esters (see [Lipids](#)). Linoleic acid is one of a family of *essential* fatty acids (EFAs) that are required for growth and development in mammals. The EFA all contain the (*n* – 6) terminal structure, which cannot be synthesized by mammalian cells because they are unable to introduce *cis*-double bonds beyond carbon 9 in fatty acids. The EFA must therefore be obtained from plants whose enzyme systems do synthesize the (*n* – 6) terminal structure. Some of the more abundant fatty acids in cells are listed in [Table 1](#).

**Table 1. Some of the More Abundant Fatty Acids in Animals**

| Systematic Name               | Common Name          | Shorthand Notation           |
|-------------------------------|----------------------|------------------------------|
| <b>Saturated</b>              |                      |                              |
| Dodecanoic                    | Lauric               | 12:0                         |
| Tetradecanoic                 | Myristic             | 14:0                         |
| Hexadecanoic                  | Palmitic             | 16:0                         |
| Octadecanoic                  | Stearic              | 18:0                         |
| Eicosanoic                    | Arachidic            | 20:0                         |
| Docosanoic                    | Behenic              | 22:0                         |
| Tetracosanoic                 | Lignoceric           | 24:0                         |
| <b>Monoenic</b>               |                      |                              |
| <i>cis</i> -9-Dodecenoic      | Lauroleic            | 12:1( <i>n</i> – 3) 12:1w3   |
| <i>cis</i> -9-Tetradecenoic   | Myristoleic          | 14:1( <i>n</i> – 5) 14:1w5   |
| <i>cis</i> -9-Hexadecenoic    | Palmitoleic          | 16:1( <i>n</i> – 7) 16:1w7   |
| <i>cis</i> -9-Octadecenoic    | Oleic                | 18:1( <i>n</i> – 9) 18:1w9   |
| <i>cis</i> -11-Octadecenoic   | <i>cis</i> -Vaccenic | 18:1( <i>n</i> – 7) 18:1w7   |
| <i>cis</i> -9-Eicosenoic      | Gadoleic             | 20:1( <i>n</i> – 11) 20:1w11 |
| <i>cis</i> -13-Docosenoic     | Erucic               | 22:1( <i>n</i> – 9) 22:1w9   |
| <i>cis</i> -15-Tetracosanoic  | Nervonic             | 24:1( <i>n</i> – 9) 24:1w9   |
| <b>Polyunsaturated (PUFA)</b> |                      |                              |
| 9,12-Octadecadienoic          | Linoleic             | 18:2( <i>n</i> – 6) 18:2w6   |
| 6,9,12-Octadecatrienoic       | g-Linolenic          | 18:3( <i>n</i> – 6) 18:3w6   |
| 8,11,14-Eicosatrienoic        | Homo-g-linolenic     | 20:3( <i>n</i> – 6) 20:3w6   |
| 5,8,11,14-Eicosatetraenoic    | Arachidonic          | 20:4( <i>n</i> – 6) 20:4w6   |

### Suggestions for Further Reading

W. W. Christie (1982) *Lipid Analysis*, 2nd ed., Pergamon Press, Oxford, U.K. (In addition to analytical methods for analysis of lipids, this book presents a useful summary of fatty acid structures and an extensive bibliography.)

J. A. Erwin (ed.) (1973) *Lipids and Biomembranes of Eukaryotic Microorganisms*, Academic Press, New York. "Chapter 2". (A concise presentation of the comparative biochemistry of fatty acids in eukaryotic microorganisms, and the diversity of fatty acid structures in nature.)

A. I. Laskin and H. A. Lechevalier (1973) *Handbook of Microbiology*, Vol. **II**, CRC Press, Cleveland. (Contains an extensive listing of fatty acids and their properties. A sourcebook for the fatty acid composition of microorganisms.)

## Female

The female is one of the two possible **sexes** in nature, the other being the [male](#). Characteristically, mature females produce **oocytes**. In species with intracorporal [fertilization](#), the most predominant feature of a female is the limited number of **gametes** that can be produced. For example in mammals, primordial **germ cells** divide by **mitosis** in females until a certain number (200,000) of oocytes is reached. After birth, these mitotic divisions come to a standstill, and **meiosis** and further maturation of oocytes occurs after puberty at ovulation. Ovulation is not a continuous production of [eggs](#), but involves a selection process in which certain ones are able to become dominant and cause atresia in others. These events typically occur in cycles, with a duration of 21 to 28 days. However, spontaneous or inducible ovulations can also take place, as in the cat or rabbit. Although the number of oocytes produced seems to be more than sufficient, it is apparent, for example, in humans, that the fertile period in the life span of a female is, unlike in males, restricted to a certain age and, within this age, to certain days periodically. This fact is probably the most typical feature of females with regard to **reproduction**.

Considering that each individual strives to transfer its genes to the next generation, the limited number of gametes with a chance to be fertilized causes females to take greater care of each individual gamete than males do. In females, [natural selection](#) therefore proceeded in a different way than in the male. Reproductive success in females is optimized when the offspring are actually raised and attain the ability to reproduce in the next generation. As a consequence, females not only use substantial resources in producing a single gamete, by adding nutrition (yolk) (1) and all necessary substances for initial cell divisions and growth of the cells, but, especially in higher chordates, they also provide basic nutrition for the offspring born (2). This is especially the case in mammals, where nutrition and protection is provided for a prolonged time after birth during childhood by the mammary glands, which are a typical female feature. Also, as can be seen in primates, raising of the offspring is done mostly by the female, although there are exceptions (3). Natural selection does not occur in a competition in which genes must be spread, but the most successful female is the one that is best able to protect and raise the offspring. In this respect, it is evident why females in the animal kingdom are better adapted to the environment with regard to hair or feathers, because this serves as a protection against predators. The appearance of this kind of sexual dimorphism depends on the

degree of specialization among the two sexes.

In species with extracorporal fertilization, such as frogs and fishes, differences between the two sexes are smaller than in mammals. In these species, both the male and female produce a vast number of gametes to form numerous new individuals. In these species, parental care for the offspring is not the rule. As a consequence, the **phenotypically** visible differences between the sexes vanish. Nevertheless, for the sake of reproductive success, females invest more than males, because the cost of producing eggs is greater than that of producing [sperm](#).

The common biological necessity in females of a greater investment in reproduction has evolved common characteristic features among females of different species. For example, the genitalia of the female consist of **germ-cell** conducting parts (tubes) as well as of germ-cell preserving parts (uterus). Although both male and female genital cells have the same origin in the primordial urogenital tract and in the primordial germ cells, the lack of male [hormones](#) due to the absence of the testis-determining factor allows the development of mullerian ducts (see [Sex](#)). The fate of these organs is to develop finally to an uterus and a vagina during embryogenesis. Also, the hormonal situation is different in females and males. After puberty, cyclic hormonal actions take place in the female. The expression of follicle-stimulating hormone (FSH) and luteinizing hormone (LH) are directed by FSH-releasing hormone and LH-releasing hormone, respectively. FSH and LH induce ovulation and the expression of estrogen and progesterone, which occurs in the follicle.

In the cascade of sex determination and in adults, hormones are also responsible for the expression of the female phenotype. This includes psychological behavior. Hormones in the female direct fertilization events like ovulation and behavior, eg to indicate the mating time for males.

In higher organisms, the female is the homogametic sex XX, the male the heterogametic XY. Sex determination occurs by the dependence of the presence or absence of the testis-determining factor (SRY), a gene located on the [Y-Chromosome](#) in mammals. In these species, therefore, sex determination is dependent on the sperm that successfully fertilizes the oocyte. In birds, some amphibians, and insects, the situation is reversed: The homogametic sex is male, the heterogametic sex is the female. Instead of two [X-chromosomes](#), the chromosomes in the female are ZW, and sex determination occurs during ovulation (see [Sex](#)).

[Parthenogenesis](#) is possible only with eggs and occurs in some genera of lizards. In some species of fish, spermatozoa are necessary only for the induction of the egg, not for a genetic contribution.

## Bibliography

1. T. D. Williams (1994) Biol. Rev. Camb. Phil. Soc. **69**, 35–59.
2. G. N. Wade and J. E. Schneider (1992) Neurosci. Biobehav. Rev. **16**, 235–272.
3. J. W. Makin and R. H. Porter (1984) Behav. Neural. Biol. **41**, 135–151.

## Ferguson Plot

Measuring the **electrophoretic** mobility of a macromolecule at several gel (polymer) concentrations (see [Gel Electrophoresis](#)) provides information about its size and surface net charge density. In the ideal case, a plot of the logarithm of the mobility versus the gel concentration is linear and known as a Ferguson plot. The slope of such a plot gives the retardation coefficient, which is a measure of the size, shape, and conformation of the macromolecule. The intercept at zero gel concentration provides

a measure of the free electrophoretic mobility, the mobility it would have in the absence of gel, which is related to the net charge on the molecular surface. In contrast, the mere inspection of one electrophoretic pattern or measurement of the mobility of one band gives no information whatsoever concerning the size or net charge of the migrating species; even a comparison with standard macromolecules is usually unreliable.

To generate a Ferguson plot, several gel (or polymer) concentrations are used simultaneously in a multitube or multichannel electrophoresis apparatus. Alternatively, a **transverse gradient** of gel concentration can be used in a single-slab gel. The slopes and intercepts of the Ferguson plot are computed and translated to parameters descriptive of the size and net charge, including their statistical limits, using relevant computer programs (1). Even without such translation, the Ferguson plot can be used to recognize qualitative relationships between the size or net charge of related bands of an electropherogram, or to test the possibility that two species are identical. Recognition of the nature of the relationships between bands helps to decide whether to use separation methods based primarily on size (see [Gel Electrophoresis](#); [Capillary Zone Electrophoresis](#)) or charge (see [Isoelectric Focusing](#); [Isotachopheresis](#)).

A linear Ferguson plot is expected from a simple mathematical model of a rigid ball passing through an inert random fiber network of the gel matrix (the Ogston model); it is observed experimentally with relatively rigid and spherical proteins passing through gels electrophoretically. In contrast, Ferguson plots derived with [agarose](#) gels are convex, presumably due to the progressive supercoiling of double-helical agarose with decreasing concentration. They are concave when used with DNA fragments, presumably due to the progressive stretching of the DNA molecule with increasing gel concentration.

#### Bibliography

1. D. Tietz (1988) *Adv. Electrophoresis* **2**, 109–170.

#### Suggestions for Further Reading

2. A. Chrambach and D. Rodbard (1971) Polyacrylamide gel electrophoresis. *Science* **172**, 440–451.
3. J. L. Hedrick and A. J. Smith (1968) Size and charge isomer separation and estimation of molecular weights of proteins by disc gel electrophoresis. *Arch. Biochem. Biophys.* **126**, 155–164.

## Ferredoxins

Ferredoxins (Fds) are simple [iron–sulfur proteins](#); that is, they do not contain other **prosthetic groups** beside one or two iron-sulfur (FeS) clusters. The first Fd to be discovered was that from *Clostridium pasteurianum* in 1962, immediately followed by the one from plant chloroplasts (1). Fds are ubiquitous, small, generally very acidic proteins, diverse in structure and function, that are involved mainly as electron carriers of low [oxidation/reduction potential](#) in fundamental metabolic processes like [photosynthesis](#), [nitrogen fixation](#), and assimilation of hydrogen, nitrogen, and sulfur. Recently, functions other than **electron transfer** have been discovered for some members of the very large class of FeS proteins. Indeed, it has become clear that most bacteria and plants contain not just one Fd, but a variety, distinguished by amino acid sequence, FeS cluster type, redox potential, and function. The Archaeobacteria are particularly rich in Fds: *Methanococcus jannaschii* has 8 Fds and 6 polyFds, and *Archeoglobus fulgidus* has 8 Fds; as are the diazotroph bacteria (there are 8 to 9

Fds in *Azotobacter vinelandii*), followed by higher plants, from which at least six soluble Fds have been isolated. In bacteria, the number of Fds varies widely: just one Fd in *Bacillus subtilis*, but four in both *Mycobacterium tuberculosis* and *Escherichia coli*. In the first completed eukaryotic [genome](#), that of the yeast *Saccharomyces cerevisiae*, only one Fd is present, which is similar to Fds found in mammalian [mitochondria](#) (adrenodoxin/renodoxin) that function in the hydroxylations catalyzed by [cytochrome P450](#) for formation of steroid hormones, vitamin D metabolites, and bile acids.

Fds contain iron and sulfur atoms, organized in three different types of iron–sulfur clusters that are defined as 2Fe–2 S, 4Fe–4 S, and 3Fe–4 S (see [Iron–Sulfur Proteins](#)). These proteins absorb light in the 300- to 500-nm wavelength region: their solutions are brownish in color and become paler upon reduction. Typical **electron para-magnetic resonance** signals are shown by Fds in either the oxidized or reduced states. These proteins are sensitive to acids, oxidants, mercurials, and alkylating agents, but they show fair stability in the pH range 6 to 9 and to heat **denaturation** in the absence of oxygen. Fds are best divided into two classes: 2Fe–2 S Fds and 4Fe–4 S Fds, that may include the 3Fe–4 S clusters.

### 1. 2Fe–2 S Ferredoxins

These Fds are ubiquitous [hydrophilic](#) proteins of 11 to 15 kDa containing one 2Fe–2 S cluster that function as one-electron carriers. The iron-binding amino-acid sequence motif is Cys–X<sub>4</sub>–Cys–X<sub>2</sub>–Cys, which provides three of the four Cys ligands to the cluster; the fourth Cys ligand is distant in the [primary structure](#). In some bacterial Fds of the hydroxylating dioxygenases (2) and in adrenodoxin, the spacing between the two first Cys residues is X<sub>3</sub> and X<sub>5</sub>, respectively. The reduction potential varies in the range –400 to –100 mV (relative to the hydrogen electrode). The prototype of this class is the plant-type Fd (3), which has the role of distributing electrons received from the photoreduced photosystem I to several ferredoxin-dependent enzymes. The [protein structure](#) motif characteristic of these proteins is called “b-grasp” (a five-stranded **b-sheet** with an **a-helix** lying on top of the sheet). Adrenodoxin and putidaredoxin have additional small interaction **domains**.

### 2. 4Fe–4 S and 3Fe–4 S Ferredoxins

The prototype bacterial ferredoxin (4, 5) contains FeS clusters of the so-called “cubane” type, in which the four iron and four sulfide atoms are placed alternatively on the corners of a somewhat distorted cube. Coordination to the protein typically occurs through four cysteine thiolate groups to the four iron atoms. The most common bacterial ferredoxins contain two such clusters in a [polypeptide chain](#) of 50 to 60 residues. These proteins are typically found in anaerobes, where they are involved in one-electron transfer in low-potential systems, collecting reducing equivalents from a wide variety of substrates and metabolic cycles. A structural distinction has been made between the above classes of ferredoxins, which have been divided into “clostridial-type” and “photosynthetic bacterial and nif-related” ferredoxins. The former contain two Cys–X<sub>2</sub>–Cys–X<sub>2</sub>–Cys–X<sub>3</sub>–Cys–Pro motifs, while the latter contain one motif of that type and the more unusual Cys–X<sub>2</sub>–Cys–X<sub>7–9</sub>–Cys–X<sub>3</sub>–Cys–Pro motif. All these proteins are low-potential electron carriers, with typical reduction potential in the –390 to –450 mV range.

Some of the clostridial-type ferredoxins found in *Desulfovibrio spp* and in extreme **thermophiles** such as *Pyrococcus furiosus* contain a Cys–X<sub>2</sub>–Asp–X<sub>2</sub>–Cys binding motif, with a remote Cys residue as the fourth ligand and the Asp available to take the place of the missing central Cys residue of this sequence. These proteins contain a 3Fe–4 S cluster, which can be converted reversibly to a 4Fe–4 S cluster upon addition of iron and which may bind to exogenous ligands, such as [thiol groups](#) and cyanide. Cornerless cubanes containing the 3Fe–4 S clusters may be formed upon oxidation (by oxygen or ferricyanide) of regular 4Fe–4 S clusters, in a possibly reversible process.

## Bibliography

1. D. I. Arnon (1988) Trends Biochem. Sci. **13**, 30–33.
2. J. R. Mason and R. Cammack (1992) Annu. Rev. Microbiol. **46**, 277–305.
3. D. B. Knaff (1996) In *Oxygenic Photosynthesis: The Light Reactions* (D. R. Ort and C. F. Yocum, eds.), Kluwer, The Netherlands, pp. 333–361.
4. M. Bruschi and F. Guerlesquin (1988) FEMS Microbiol. Rev. **54**, 155–176.
5. J. J. G. Moura, A. L. Macedo, and P. N. Palma (1994) Methods Enzymol. **243**, 165–188.

## Suggestions for Further Reading

6. P. J. Stephens, D. R. Jollie, and A. Warshel (1996) Protein control of redox potentials of iron-sulfur proteins. Chem. Rev. **96**, 2491–2513.
7. H. Matsubara and K. Saeki (1992) Structural and functional diversity of ferredoxins and related proteins. Adv. Inorg. Chem. **38**, 223–280.
8. H. M. Holden et al. (1994) Structure–function studies of [2Fe–2 S] ferredoxins. Bioenerg. Biomembr. **26**, 67–87.

## Ferritins

Ferritins and [transferrins](#) are nonheme iron-transport proteins that provide the crucial link between dietary iron and organismal heme-iron pools ([hemoglobin](#), [myoglobin](#), and hemosiderin) in maintaining iron homeostasis in mammals. Ferritin is a large protein (ca. 500 kDa) that consists of an aggregate of so-called heavy and light chains, which form a hollow shell in the apo form. Each apoferritin molecule can take up 4500 iron atoms, effectively shielding them from harmful reactions with peroxide and superoxide.

Transferrins are much smaller proteins that take up iron by binding to the transferrin receptor (TfR) at the cell surface. The complex is subsequently internalized by [endocytosis](#), and then the iron dissociates into the cytoplasm, where it is free to be taken up by ferritin and other cellular proteins.

The [messenger RNAs](#) for transferrin receptor and for both ferritin subunits contain so-called [iron-response elements](#) (IREs), which are RNA structures containing a stem and a loop that has the consensus sequence CAGUGX ([1](#)). This regulates the [translation](#) of the mRNAs.

## Bibliography

1. M. W. Hentze and L. C. Kühn (1996) Proc. Natl. Acad. Sci. USA **93**, 8175–8182.

## Suggestions for Further Reading

2. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New York, Vols. **1–35**.
3. R. J. Cousins (1994) Metal elements and gene expression, Annu. Rev. Nutr. **14**, 449–469.



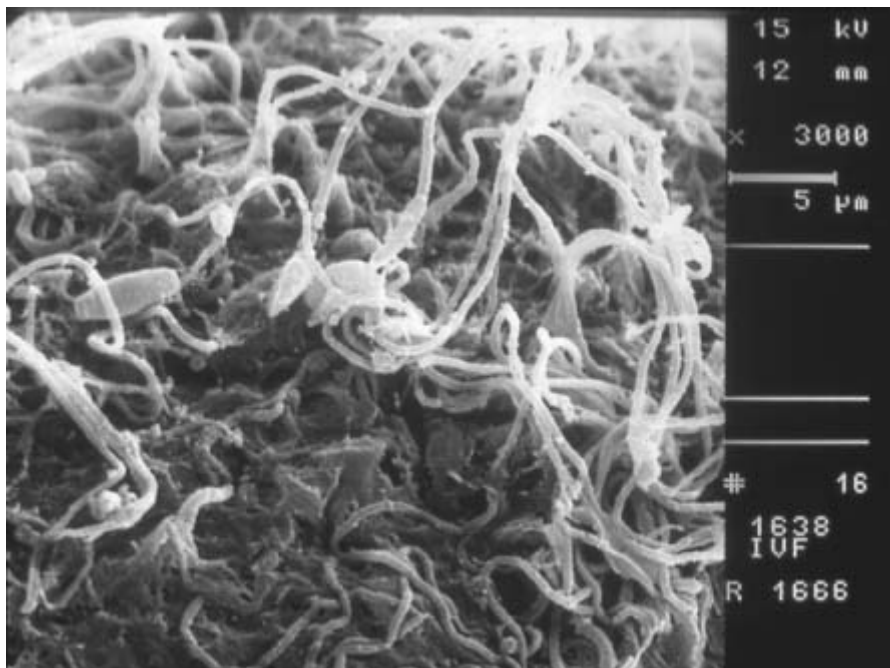
## Fertilization

In the animal kingdom, fertilization is a unique event in which a [male](#) and [female](#) haploid **gamete** join and follow a cascade of highly regulated events that results in the production of a new individual (Fig. 1). During the early stages of fertilization, the two morphologically disparate gametes, the **oocyte** and the [sperm](#), follow a very specific process of recognition and fusion. The interactions between the two gametes and with their tubal environment are mediated by highly specialized molecules located on the outer surfaces of membranes of the gametes and on the surrounding epithelia. The molecular interactions between the oocyte and spermatozoon are also unique with regard to gene function. Both gametes are [haploid](#), and furthermore the [nucleus](#) is greatly condensed in the spermatozoa. Thus, no gene [transcription](#) takes place; consequently, all biochemical and molecular biological interactions must be done by proteins and, to some extent, by pretranscribed RNA of the oocyte. This means that regulation of the nucleus does not take place; instead, proteins and RNA involved in the fertilization process act autonomously. On the spermatozoon, the adhesion molecules are divided into primary and secondary ligands (see [Acrosome](#)). Primary ligands are involved in initial gamete binding, whereas secondary ligands, such as proacrosin, come into action after the acrosome reaction has taken place. In the oocyte, interactions with the other gamete are mediated by an extracellular matrix, the *zona pellucida* (see [Egg](#)). This organelle consists of three **glycoproteins**, ZP1, ZP2, and ZP3, of which ZP1 is responsible for structural integrity, whereas ZP2 and ZP3 have been shown to have sperm receptor activity (1). In the mouse, primary ligands interact with ZP3, whereas secondary ligands like proacrosin also bind strongly to ZP2. Only after completion of all complementary signals and responses is penetration of the spermatozoon through the zona pellucida possible, which directs the transfer of the male genetic information to the oocyte. The fertilization process has been divided into individual steps by Wassarman (1), in a widely accepted scenario:

1. **Capacitation.** In the female genital tract, ejaculated spermatozoa undergo a process called *capacitation*. This maturing process of spermatozoa is necessary to gain full fertilization capacity, but the basic mechanisms are not fully understood. [Lectins](#) from the seminal cytoplasm adhering to the surface of the acrosomal cap undergo a change in conformation and adopt structures that enable coupling of primary receptors.
2. **Attachment.** *In vitro* fertilization experiments demonstrated that the oocyte and spermatozoa initially are attached only loosely. This loose attachment is not species-specific, and the spermatozoa can be removed easily from the egg by small turbulences.
3. **Primary binding.** After the initial contact, the sperm bind tightly to the zona pellucida. This species-specific interaction is mediated by primary receptors and ligands. Now the spermatozoa cannot be removed easily from the egg by physical means.
4. **Acrosome reaction.** Triggered by the contact with the zona pellucida, spermatozoa undergo [exocytosis](#), known as the *acrosome reaction*. This involves hyperactivation of the motile spermatozoa, multiple point fusion of sperm cytoplasm and the outer acrosomal membrane, formation of hybrid vesicles, and exposure of the inner acrosomal membrane. [Enzymes](#) (mainly proacrosin) and inner acrosomal components that are necessary for binding, local lysis, and penetration of the zona pellucida are released and come into action.
5. **Secondary binding.** After the acrosome reaction, secondary ligands and receptors maintain the affinity between the two gametes. The forces are strong enough to keep the spermatozoa tightly attached to the zona pellucida. Adhesive proteins like proacrosin play an important role in high affinity binding.
6. **Penetration.** After exposure of the inner acrosomal membrane, the predominant molecule of the acrosome, *acrosin*, has been generated from the precursor proacrosin. It acts as a binding molecule and also a [proteinase](#); these two phenomena together enable the sperm to bind and lyse the zona pellucida locally. The spermatozoon is driven through the zona pellucida, forced by the motile sperm tail.

7. **Fusion.** The spermatozoon that first reaches the perivitelline space then fuses with the egg plasma membrane to form a [zygote](#). On entry of the spermatozoon, the egg is activated, which leads finally to the unification of the male and female **pronuclei**.
8. **Primary block to polyspermy.** In order to prevent several spermatozoa fusing with the oocyte, a rapid depolarization of the egg plasma membrane occurs within seconds after gamete fusion.
9. **Cortical reaction and zona hardening.** The cortical granules of the egg, which can be regarded as **lysosome-like organelles**, release a variety of lytic enzymes on the fusion of a spermatozoon with the egg plasma membrane. The cortical granule contents are then deposited into the perivitelline space; when they enter the zona pellucida, the ZP2 and ZP3 glycoproteins are converted into ZP2f and ZP3f, the respective inactive forms of the receptor. This modification of the zona pellucida glycoproteins constitutes a permanent secondary block to polyspermy.

**Figure 1.** Sperm–egg interaction. Capacitated spermatozoa bind to the zona pellucida via primary ligands and receptors. This event triggers the acrosome reaction of the spermatozoa and activates secondary ligands on the sperm. (With the friendly permission of Dr. H.-W. Michelmann.)



## Bibliography

1. P. M. Wassarman (1990) *J. Reprod. Fert. Suppl.* **42**, 79–87.

## Fiber Diffraction

Fiber diffraction is a method of using X-ray diffraction to obtain structural information from a fibrous preparation of the material, rather than a three-dimensional crystal (see [X-Ray Crystallography](#)). Fibers consist of elongated molecules that are aligned parallel to each other along

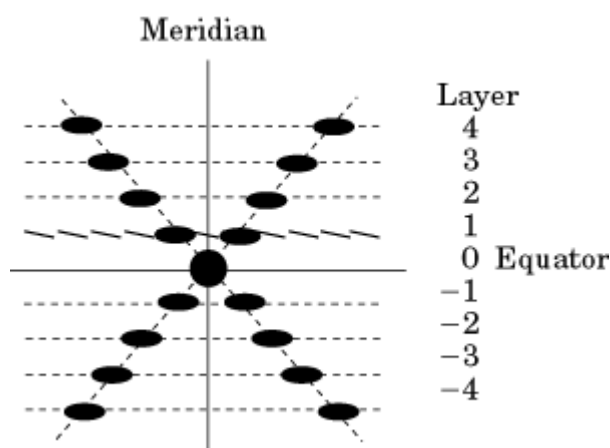
the fiber axis. In biological fibers, such as hair, [collagen](#) and **muscle**, the molecules are naturally aligned. In others the molecules must be forced to align by drawing them from a concentrated solution or by flow through a capillary. Magnetic and electric fields also have an effect.

The relatively weak X-ray scattering from fibrous specimen means that long exposures are necessary to record the diffraction pattern, and special cameras have been developed for the purpose. A fiber X-ray diffraction pattern shows series of lines of intensity perpendicular to the fiber axis, called layer lines. If the fiber axis is vertical, the layer lines are horizontal. Their distance is related to the periodicity along the fiber axis. The middle line is called the zero layer line or equator. The meridian is an imaginary line running vertically through the center of the diffraction pattern. The quality of the pattern depends on the degree of alignment of the molecules, but an X-ray fiber diffraction pattern is rarely of sufficiently high quality to determine a structure at atomic resolution. The best results are obtained with synchrotron radiation of cryocooled specimens.

In the 1930s, Astbury used fiber diffraction to group protein fibers into three types: a-fibers, b-fibers, and collagen. For a structural determination from fiber data, chemical information about the fiber is needed to postulate a model and to refine this model, obeying any symmetry constraints, until it predicts the observed intensities reasonably well. In this way, Watson and Crick arrived at the structure of **DNA** (1) and Pauling at the model for the [a-helix](#) (2). For larger structures, this is not so easy, and then [isomorphous replacement](#) is applied, as in protein crystal structure determination.

Fibers usually consist of helical arrangements of the building units. Cochran, Crick, and Vand found that the intensity diffracted by a continuous helix on the  $n$ th layer line is proportional to  $J_n^2(2pRr)$  (3). The intensity distribution is cylindrically symmetrical around the reciprocal axis parallel to the fiber axis.  $R$  is the distance in reciprocal space to this axis, and  $r$  is the radius of the helix.  $J_n$  is a Bessel function of order  $n$ . These Bessel functions are oscillating functions that fade away with the distance from the origin. Therefore, the highest values are found close to the origin, which here is the meridian in the diffraction picture. But with increasing order  $n$ , these maximum values are found at a greater distance from the origin and become smaller. The combination of these properties creates the X-shaped diffraction picture characteristic of a helical structure (Fig. 1).

**Figure 1.** Schematic depiction of the typical diffraction pattern of a continuous helix.



Actual helices are not continuous, and the effect is that the intensity on a layer line depends on the structural repeat. If  $P$  is the pitch of the helix and  $p$  the axial rise from building unit to building unit along the helix axis, layer line numbers  $\ell$  obey the equation  $\frac{\ell}{c} = \frac{n}{P} + \frac{m}{p}$ , where  $c$  is the structural

repeat distance along the helix axis,  $n$  is the order of a Bessel function, and  $m$  is (like  $n$ ) a whole number. This is illustrated by the following example. For the  $\alpha$ -helix,  $P = 5.4 \text{ \AA}$ , and  $p = 1.5 \text{ \AA}$ . The structure repeats itself along the helix axis over a distance  $c = 27 \text{ \AA}$ . Therefore,

$$P = \frac{c}{5} \quad p = \frac{c}{18} \quad \frac{\ell}{c} = \frac{n}{P} + \frac{m}{p} = \frac{5n}{c} + \frac{18m}{c}$$

For layer line  $\ell = 0$ , this can be reached with the following combinations:

---


$$\begin{aligned} n &= -36 \ -18 \ 0 \ +18 \ +36 \\ m &= +10 \ +5 \ 0 \ -5 \ -10 \end{aligned}$$


---

and for layer line  $\ell = 1$  with

---


$$\begin{aligned} n &= -25 \ -7 \ +11 \ +29 \ +47 \\ m &= +7 \ +2 \ -3 \ -8 \ -13 \end{aligned}$$


---

The intensity on every layer line is determined by a large number of Bessel functions  $J_n$ , but they have a nonnegligible value only for small  $n$ . Therefore, strong reflections occur on the zero layer line ( $\ell = 0$ ) in the  $\alpha$ -helix diffraction pattern, but not on the first one ( $\ell = 1$ ). Layer line 18 also has a strong reflection.

### Bibliography

1. J. D. Watson and F. H. C. Crick (1953) *Nature* **171**, 737–738.
2. L. Pauling, R. B. Corey, and H. R. Branson (1951) *Proc. Natl. Acad. Sci. USA* **37**, 205–211.
3. W. Cochran, F. H. C. Crick, and V. Vand (1952) *Acta Crystallogr.* **5**, 581–586.
4. A. Klug, F. H. C. Crick, and H. W. Wyckoff (1958) *Acta Crystallogr.* **11**, 199–213. A complete treatment of the theory of fiber diffraction.
5. D. A. Marvin and C. Nave (1982) In *Structural Molecular Biology, Methods and Applications* (D. B. Davies, W. Saenger, and S. S. Danyluk, eds.), Plenum, New York, London, pp. 3–44. A review.

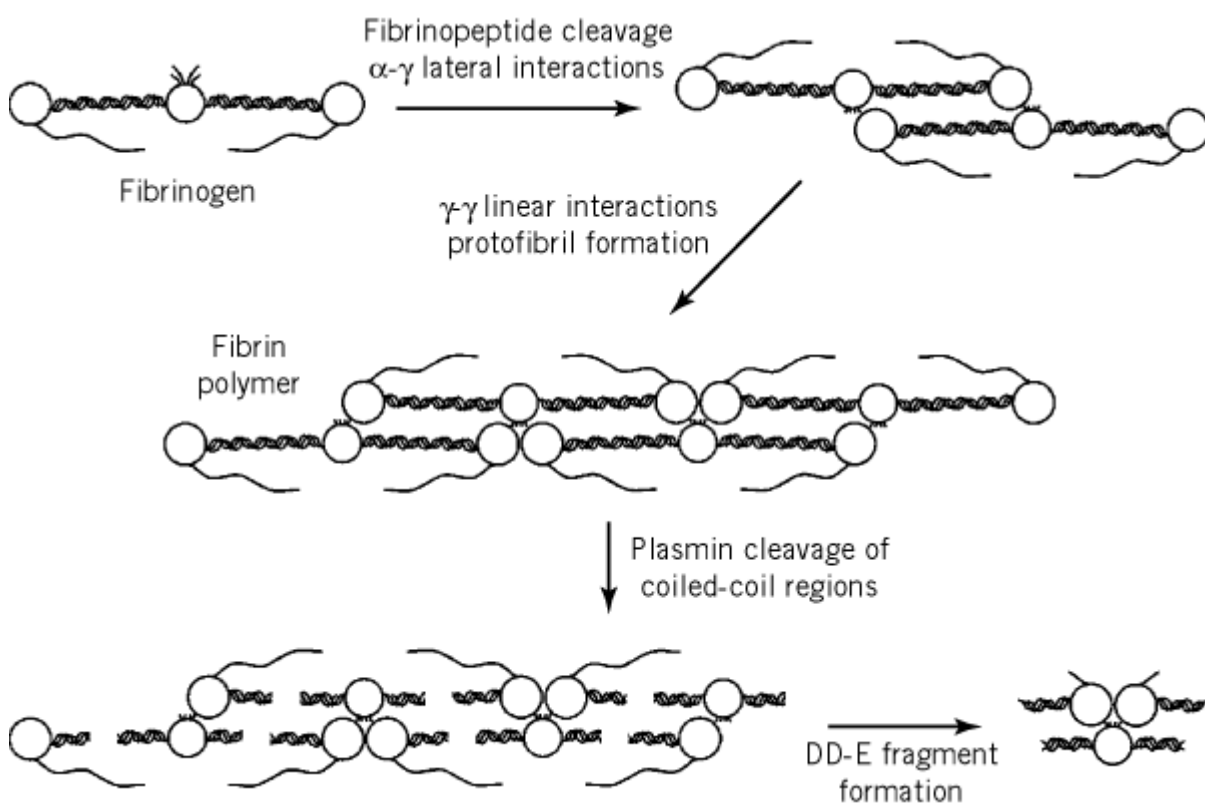
### Fibrinogen

Blood clots when soluble fibrinogen is cleaved by thrombin to generate fibrin monomers which polymerize spontaneously, non-covalently and rapidly to form the insoluble fibrin clot, which is further stabilized by covalent cross-linking catalyzed by the transglutaminase, factor XIII<sub>a</sub> [see [Blood Clotting](#)]. This remarkable property is dependent on the symmetrical multi-chain structure of fibrinogen, the proteolytic exposure of cryptic binding sites, and the high concentration of fibrinogen in blood plasma (2–4 mg/mL). The cessation of bleeding is further enhanced by fibrin acting as a cofactor in blood platelet aggregation. The fibrin clot further acts as a provisional matrix for the subsequent repair of the wound.

Fibrinogen is a disulfide-bonded hexameric protein, arranged as a dimer of Aa, Bb and g chains.

These are the products of three distinct, but related, genes located in close association on human chromosome 4. The mature chains are 625, 461 and 411 residues, respectively, resulting in a molecular weight of 350 kDa for the hexamer. The gross topological structures of fibrinogen and fibrin have been deduced by both biochemical and biophysical methods. The overall appearance of fibrinogen is a linear trinodular structure 475 Å in length, with each nodule separated by an intervening rod-like region (1) (Figure 1). The N-terminal regions of all six chains are contained in the highly disulfide-bonded central nodule, the N-terminal disulfide knot. The N-termini of the Aa and Bb chains protrude from the nodule, while the g chains have shorter N-terminal sequences. Triple-chained coiled coils connect the central to the distal nodules, each of which is composed of two homologous globular domains comprising the C-terminal regions of the Bb and g chains. The longer Aa chain folds back along the axis of the molecule, to end in a globular domain that interacts with the N-terminal disulfide knot.

**Figure 1.** Schematic model of fibrin polymerization and its subsequent dissolution.



Structures have been determined by X-ray crystallography for various C-terminal proteolytic fragments of human fibrinogen (2, 3) and, although intact human fibrinogen has yet to be crystallized, the structure of native chicken fibrinogen has been resolved at a resolution of 2.7 Å (4). The triple-chained connecting regions are susceptible to proteolysis, generating two major fragments: fragment E corresponding to the central N-terminal nodule (approximately 15% of the molecule) and fragment D containing the distal globular domains (the two identical fragments accounting for approximately 50% of the molecule). Both of these fragments contain substantial portions of the cleaved coiled coil regions with the remainder of the protein accounted for by the C-terminal regions of the Aa chains.

Fibrin polymerization essentially occurs in three sequential steps. Firstly, thrombin cleaves the N-terminal regions of the Aa and Bb chains at Arg16-Gly17 and Arg14-Gly15, respectively, releasing

fibrinopeptides A and B. The N-terminal sequence of the fibrin a chain (Gly17-Pro-Arg-Val) is the major determinant of fibrin polymerization which proceeds efficiently without further modification to the fibrin(ogen) structure, and synthetic peptides corresponding to this sequence can inhibit the polymerization of fibrin monomers (5). By contrast the b-chain sequence (Gly15-His-Arg-Pro), which is generated more slowly, although binding to fibrinogen and fragment D, cannot inhibit polymerization. Secondly, the new N-termini of fibrin monomer spontaneously interact with binding sites in the two homologous distal globular domains, with the g-chain site particularly implicated (3). Both sites are large cavities that bind a-chain Gly-Pro-Arg ligands, but the g-chain cavity is the more specific and does not bind b-chain Gly-His-Arg ligands. It is also the more distal of the two C-terminal domains and is able to intermolecularly associate in a specific “end-to-end” manner (2). Thus adjacent g-chain sites are closely situated for optimal binding of the dimeric a-chain N-terminal residues and the efficiency of polymerization derives from both a-g lateral and g-g linear interactions. The symmetry of fibrin(ogen) allows oligomerization to proceed in both directions, generating a two-molecule thick protofibril with a half-staggered overlap (Figure 1). Finally, the fibrin fibrils are reinforced by the transglutaminase activity of factor XIII<sub>a</sub>. The g-g linear interactions are the major sites for the introduction of intermolecular g-glutamyl-ε-amino-lysine isopeptide crosslinks, subsequently the globular C-terminal domains of the a chains of adjacent laterally aligned protofibrils are also crosslinked, resulting in the mature, covalently stabilized fibrin mesh. Under physiological conditions these fibers are of the order of 1 μM in diameter and thus several hundred fibrin units wide.

The assembly of fibrin also reveals epitopes that are involved in its subsequent dissolution by plasmin. Conformational changes in the vicinity of the C-terminal globular domains are thought to reveal the binding sites for tPA that are responsible for the fibrin-specific potentiation of its activity [see [Plasminogen Activators](#)]. These sites are g312-324 in the C-terminal globular domain and Aa148-160, a region interacting with the Bb-chain C-terminal globular domain. The initial binding sites for plasminogen also reside in fragment D, but subsequent plasmin degradation of fibrin reveals C-terminal Lys residues that act as supplementary sites for plasminogen binding, thereby propagating plasminogen activation. Once generated plasmin can attack 50-60 of the 358 Lys and Arg residues present in each fibrin molecule, although only 10% are cleaved rapidly. The initial cleavages release the C- and N-termini of the a and b chains, respectively, followed by the cleavages in the extended coiled-coil regions that lead to dissolution of the organized fibrin structure (see Figure 1).

## 1. Fibrin Polymerization

The generation of fibrin from fibrinogen and its polymerization essentially occur in three sequential steps.

1. Thrombin cleaves the N-terminal regions of the Aa and Bb chains of fibrinogen at the [peptide bonds](#) Arg16–Gly17 and Arg14–Gly15, respectively, releasing fibrinopeptides A and B. No further covalent modification to the fibrin(ogen) structure is necessary for polymerization to proceed efficiently. The new N-terminal sequence of the fibrin a chain (Gly17–Pro–Arg–Val) is the major determinant of fibrin polymerization; synthetic peptides corresponding to this sequence can inhibit the polymerization of fibrin monomers (4). By contrast, the new b-chain N-terminal sequence (Gly15–His–Arg–Pro), which is generated more slowly, cannot inhibit polymerization, even though it binds to fibrinogen and to fragment D.

2. The new N-termini of the fibrin monomers spontaneously interact with binding sites in the two homologous distal globular domains, with the g-chain site particularly implicated (3). Both sites are large cavities that bind the a-chain-like peptide Gly–Pro–Arg, but the g-chain cavity is the more specific and does not bind g-chain Gly–His–Arg peptides. It is also the more distal of the two C-terminal domains and is able to associate intermolecularly in a specific end-to-end manner (2). Thus adjacent g-chain sites are closely positioned for optimal binding of the dimeric a-chain N-terminal

residues. The efficiency of polymerization derives from both a–g lateral and g–g linear interactions. The symmetry of fibrin(ogen) allows oligomerization to proceed in both directions, generating a two-molecule thick protofibril with a half-staggered overlap (Fig. 1).

3. The fibrin fibrils are reinforced by the transglutaminase activity of factor XIII<sub>a</sub>. The g–g linear interactions are the major sites for the introduction of intermolecular isopeptide crosslinks between g–glutamyl carboxyl and lysine-ε-amino groups. Subsequently, the globular C-terminal domains of the a chains of adjacent laterally aligned protofibrils are also crosslinked, resulting in the mature, covalently stabilized fibrin mesh. Under physiological conditions, these fibers are of the order of 1 μm in diameter, and thus several hundred fibrin units wide.

## 2. Dissolution of the Fibrin Clot

The assembly of fibrin also makes accessible [epitopes](#) that are involved in its subsequent dissolution by plasmin, after its activation from [plasminogen](#) by tissue-type [plasminogen activator](#) (tPA). The activity of tPA is potentiated specifically by fibrin, and conformational changes in the vicinity of the C-terminal globular domains of fibrin are thought to reveal the binding sites that are responsible. These sites are residues 312–324 in the C-terminal globular domain of the g chain and 148–160 of Aa, a region that interacts with the Bb-chain C-terminal globular domain. The initial binding sites for plasminogen also reside in fragment D, but subsequent plasmin degradation of fibrin uncovers C-terminal [lysine](#) residues that act as supplementary sites for plasminogen binding, which thereby propagate plasminogen activation. Once generated, plasmin can attack 50–60 of the 358 Lys and Arg residues present in each fibrin molecule, although only 5 or 6 are cleaved rapidly. The initial cleavages release the C and N termini of the a and b chains, respectively, followed by the cleavages in the extended coiled-coil regions that lead to dissolution of the organized fibrin structure (see Fig. 1).

## Bibliography

“Fibrinogen” in , Vol. 2, pp. 907–909, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ; “Fibrinogen” in (online), posting date: January 15, 2002, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ.

1. W.R. Fowler and H.P. Erickson (1979) *J. Mol. Biol.* **134**, 241–249.
2. G. Spraggon, S.J. Everse and R.F. Doolittle (1997) *Nature*, **389**, 455–462.
3. K.P. Pratt, H.C. Côté, D.W. Chung, R.E. Stenkamp and E.W. Davie (1997) *Proc. Natl. Acad. Sci. USA*, **94**, 7176–7181.
4. Z. Yang, J.M. Kollman, L. Pandi and R.F. Doolittle (2001) *Biochemistry* (In Press)
5. A.P. Laudano and R.F. Doolittle (1980) *Biochemistry*, **19**, 1013–1019.

## Suggestions for Further Reading

6. M.W. Mosesson, K.R. Siebenlist and D.A. Meh (2001) The structure and biological features of fibrinogen and fibrin. *Ann. N Y Acad. Sci.* **936**, 11–30.

## Fibroblast Growth Factors

The fibroblast growth factors (FGFs) are a family of polypeptide [growth factors](#) that control many different biological processes *in vivo*. Their biological effects are mediated by association with specific cell-surface receptors of the intrinsic **tyrosine kinase** class (see [Tyrosine Kinase Receptors](#)). Their activity is modulated by association with cell-surface heparan sulfate containing glycosaminoglycans and other cofactors.

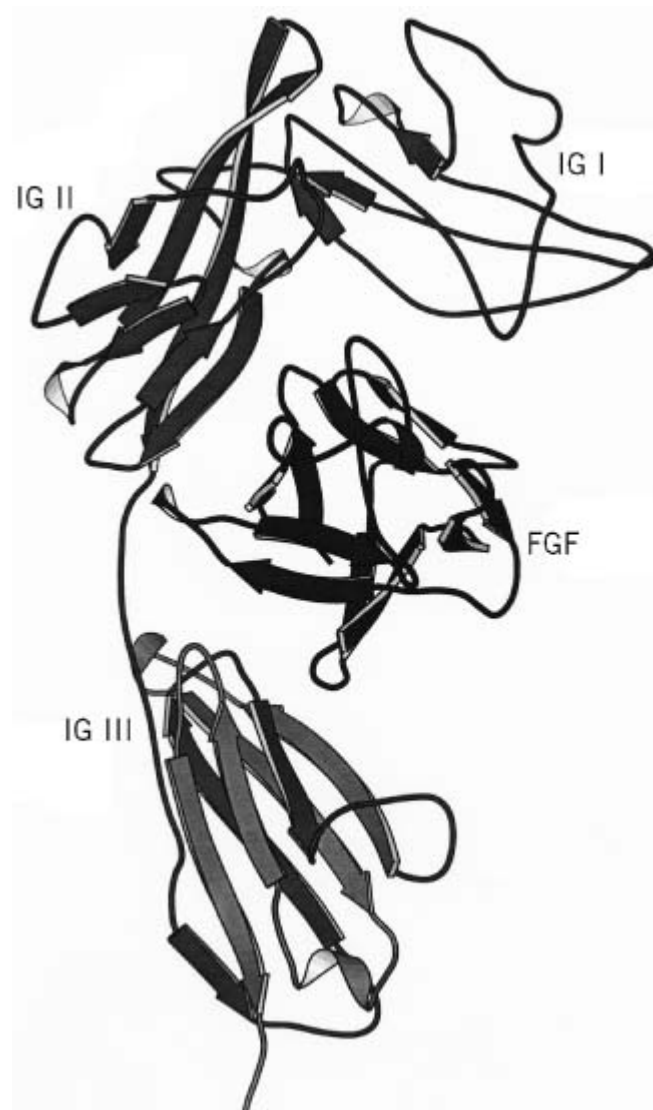
## 1. Fibroblast Growth Factor Ligands

The FGF family in mammals currently comprises 14 different [polypeptide chains](#) that act as ligands, FGF1–FGF10 and FGF15–FGF18, and share an overall sequence identity of about 40%. The FGF family is highly conserved in animals: **homologues** have been isolated from a wide variety of species, including birds, amphibians, fish, insects, and [nematodes](#) (reviewed in Ref. [1](#)). In addition to the defined family of FGF ligands, four related genes, FHF1–FGF4 (also known as FGF11–FGF14), have been isolated from [expressed sequence tag](#) (EST) projects that exhibit sequence similarity to canonical FGF ligands, but whose biological function is currently unknown ([2](#)). It is very possible that further FGF family members remained to be identified.

The [protein structures](#) of both FGF1 and FGF2 have been determined by [X-ray crystallography](#) ([3](#)). The two structures are remarkably similar and exhibit a novel 12-stranded b-trefoil topology (Fig. [1](#)). This characteristic structure is shared with a biologically unrelated growth factor, **interleukin -1** (IL1).

**Figure 1.** A molecular model of the interaction between FGF and the extracellular domain of FGF-R. The model was created by docking the crystal structure of FGF-2 ([3](#)) in place of the homologous IL1 in the crystal structure of the IL-1/IL1-R complex ([9](#)). The model shows only one half of the active complex, which is composed of two molecules of FGF and two receptors.





FGF family ligands exhibit distinct and characteristic patterns of tissue-specific gene expression during development and adulthood, suggesting that, despite their diversity, each may exhibit distinct biological functions *in vivo*. Although in most cases FGFs are expressed as a secreted molecules and elicit their biological effects in the extracellular milieu, an intriguing feature of two FGF family members, FGF1 and FGF2, is that they lack classic secretory [signal peptides](#) and, in most cases, appear to remain in an intracellular location (see [Protein Secretion](#)). This suggests that these FGFs may either require cell lysis to act or that they are released under certain conditions by a nonstandard protein export mechanism.

## 2. FGF Receptors

The biological effects of the FGF family of ligands are mediated by association with a family of related cell-surface receptors of the intrinsic tyrosine kinase class. The FGF family of receptors is encoded by four related genes, FGF-R1–FGF-R4. FGF-R genes have been isolated from mammals, birds, fish, amphibians, *Drosophila*. and nematodes (1). Individual FGF-R genes exhibit significant conservation between species and conserved features shared between the entire family. Individual FGF ligands exhibit defined patterns of receptor specificity (4), whereby some ligands interact with all FGF-R family members and others exhibit significantly more restricted affinities.

The extracellular region of the FGF-Rs comprises three domains, IgI, IgII and IgIII, which contain

[immunoglobulin](#) motifs of the immunoglobulin [superfamily](#) (IgSF). Each Ig domain is encoded by two exons. In FGF-R1–FGF-R3, the IgIII domain is generated by [alternative splicing](#) of a common IgIIIa exon onto a variant IgIIIb or IgIIIc exon. The expression of FGF-R IgIII splice variants is subject to developmentally regulated gene-specific alternative splicing mechanisms (5). The selection of IgIII exons by alternative splicing mechanisms is of major functional significance because it dictates ligand specificity of FGF-Rs (6). The analysis of chimeric FGF-Rs containing variable regions from each splice variant has revealed that the interaction of an FGF ligand with its cognate receptor involves regions of both IgII and IgIII (7). Molecular models of the FGF-R ectodomain (Fig. 1) have been generated, on the basis of relationships to other IgSF family members and the crystal structure of the structurally related ligand IL1 (8, 9). These suggest that the active complex is composed of two molecules of ligand and two molecules of receptor. The ligand interacts with [epitopes](#) in both the second and third Ig domains of the receptor. These studies also reveal the existence of an extended linker domain between the IgII and IgIII domains of IL1-R. This leads to a general scheme for FGF/receptor engagement in which two independent epitopes in the ligand interact with cognate sites in IgII and IgIII. This interaction may involve interdomain reorientation on ligand binding, facilitated by the linker domain.

### 3. FGF Cofactors

The biological actions of FGF ligands may involve a variety of cofactors, apart from the transmembrane tyrosine kinase receptors. The most prominent of these is the sulfated oligosaccharide heparan sulfate (HS), which has been shown in many studies to interact with FGF ligands with high affinity. The crystal structure of FGF2 complexed with a linear hexasaccharide (10) reveals that this interaction occurs via binding of sulfated groups of HS to a cluster of basic residues exposed on one face of the FGF ligand. The biological role of HS in FGF action is still the subject of controversy. It has been shown that HS is, in some cases, required for FGF action (5). This has been argued to result from either HS-mediated ligand dimerization, or interaction of HS with both ligand and receptor. Others have proposed that the role of HS is simply to concentrate FGF ligand in the vicinity of the receptor (reviewed in Ref. 11).

A second FGF-binding protein is the cysteine-rich FGF receptor (12) that has been identified in many tissue types. This is a transmembrane protein, unrelated to FGF-Rs, which binds FGF ligands with high affinity. The cysteine-rich FGF receptor has been identified recently as a cell adhesion ligand of E-selectin, an inducible [cell adhesion molecule](#) that mediates the binding of neutrophils to endothelial cells and functions as a  $\text{Ca}^{2+}$ -dependent [lectin](#). The signaling potential of this receptor, and the relationship between the interaction with E-selectin and FGF ligands, remain to be determined.

Several agents exist that interact directly or indirectly with FGF-Rs. Two proteins have been isolated by interaction screens in yeast that are unrelated to FGF ligands in sequence but appear to potentiate FGF-R signaling when overexpressed in [Xenopus](#) embryos (13). Genetic analysis of FGF signaling in *Drosophila* has revealed the existence of a conserved family of proteins, Sproutys, which, on genetic grounds, act to antagonize FGF receptor-mediated signaling pathways (14). Finally a variety of lines of evidence suggest that some effects of the homotypic cell adhesion molecules L1 and N-Cam may be mediated by direct, or indirect, activation of FGF-R signaling (15). Together these findings indicate that FGF signaling mediated by receptor–ligand interactions may be tightly constrained by the presence of specific cofactors *in vivo*.

### 4. Biological Actions of FGF Ligands

Genetic evidence, from gain-of-function or loss-of-function experiments, shows that the FGF family of ligands in higher vertebrates exhibits a diversity of biological functions at different times and places during embryonic and postnatal development. These include mesoderm induction, patterning of the nervous system, limb induction and outgrowth, and the development of various tissues and

organs, such as the hair, ear, lungs, and teeth (reviewed in Ref. 16). The analysis of mice harboring null mutations of FGF receptors reveals that each FGF-R has distinct biological roles *in vivo* and, in particular, that FGF signaling is required at early stages of embryonic development. Thus genetic inactivation of FGF-R1 leads to early postimplantation lethality resulting from aberrant mesodermal patterning (17). Genetic inactivation of FGF R2 leads to preimplantation lethality (18), with a phenotype similar to that of homozygous null FGF-4 mice. Homozygous null FGF-R3 mice exhibit severe and progressive bone dysplasia with enhanced and prolonged endochondrial bone growth, indicating that signaling mediated by FGF-R3 acts as a negative regulator of endochondrial ossification (19).

The FGF/FGF-R axis also has considerable significance for clinical disorders of craniofacial and limb development. One of the commonest forms of inherited dwarfism, achondroplasia, arises from specific mutations in the transmembrane domain of FGF-R3 (20). This results in activation of FGF-R signaling and premature termination of long-bone extension. Moreover, mutations in the FGF-R1–FGF-R3 genes have been identified in dominantly inherited human craniofacial malformation syndromes (reviewed in Ref. 21) such as those of Crouzon, Pfeiffer, and Apert. These syndromes are all characterized by craniosynostosis (abnormal development and/or premature fusion of the cranial sutures) and additional phenotypic features such as syndactyly (digit fusion) in the Apert syndrome. The majority of these mutations described are clustered within a specific region of the extracellular domain of the receptor. In the case of Crouzon–Pfeiffer syndromes, mutations are predicted to disrupt the pattern of **disulfide bonding** in the extracellular domain (22), leading to activation of signaling via intermolecular disulfide bonding of mutant receptors. In Apert syndrome, all cases characterized to date involve specific amino acid substitutions of two adjacent amino acid residues, Ser252 and Pro253, predicted to lie in the linker region joining Ig domains II and III. The consequence of these mutations appears to be enhanced affinity for specific FGF ligands (23), which results in activation of FGF signaling in the developing suture.

## 5. Conclusions

The FGFs exemplify many generic features of growth factors. There are multiple ligands and receptors, each of which executes specific functions in development and adult physiology. Individual FGFs can also exert several different functions, which are not limited to the control of cell multiplication. In particular, FGFs seem to play important roles in patterning tissues of the developing embryo. FGF signaling is strictly local in action and modified by the concurrent activity of various cofactors. Finally, the dynamics of signaling elicited by FGF-R activation seems to play a critical role in the development of specific tissues.

## Bibliography

1. F. Coulier, P. Pontarotti, R. Roubin, H. Hartung, M. Goldfarb, and D. Birnbasum (1997) Of worms and men: an evolutionary perspective on the fibroblast growth factor (FGF) and FGF receptor families, *J. Mol. Evol.* **44**, 43–56.
2. P. M. Smallwood, Munoz-San Juan, P. Tong, J. P. Macke, S. H. Hendry, D. J. Gilbert, N. G. Copeland, N. A. Jenkins, and J. Nathans (1996) Fibroblast growth factor (FGF) homologous factors: new members of the FGF family implicated in nervous system development, *Proc. Natl. Acad. Sci. USA* **93**, 9850–9857.
3. X. Zhu, H. Komiya, A. Chirino, S. Faham, G. M. Fox, T. Arakawa, B. T. Hsu, and D. C. Rees (1991) Three-dimensional structures of acidic and basic fibroblast growth factors, *Science* **251**, 90–93.
4. D. M. Ornitz, X. J. , J. S. Colvin, D. G. McEwen, C. A. MacArthur, F. Coulier, G. Gao, and M. Goldfarb (1996) Receptor specificity of the fibroblast growth factor family, *J. Biol. Chem.* **271**, 15292–15297.
5. D. E. Johnson, P. L. Lee, J. Lu, and L. T. Williams (1990) Diverse forms of a receptor for acidic and basic fibroblast growth factors, *Mol. Cell. Biol.* **10**, 4728–4736.

6. D. Givol and A. Yayon (1992) Complexity of FGF receptors: genetic basis for structural diversity and functional specificity, *FASEB J.* **6**, 3362–3369.
7. T. E. Gray, M. Eisenstein, T. Shimon, D. Givol, and A. Yayon (1995) Molecular modeling based mutagenesis defines ligand binding and specificity determining regions of fibroblast growth factor receptors, *Biochemistry* **34**, 10325–10333.
8. A. Bateman and C. Chothia (1995) Outline structures for the extracellular domains of the fibroblast growth factor receptors, *Nat. Struct. Biol.* **2**, 1068–1074.
9. G. Vigers, L. J. Anderson, P. Caffes, and B. J. Brandhuber (1997) Crystal structure of the type-1 interleukin-1 receptor complexed with interleukin-1 beta, *Nature* **386**, 190–193.
10. S. Faham, R. E. Hileman, J. R. Fromm, R. J. Linhardt, and D. C. Rees (1996) Heparin structure and interactions with basic fibroblast growth factor, *Science* **271**, 1116–1120.
11. D. E. Johnson, J. Lu, H. Chen, S. Werner, and L. T. Williams (1991) The human fibroblast growth factor receptor genes: a common structural arrangement underlies the mechanisms for generating receptor forms that differ in their third immunoglobulin domain, *Mol. Cell. Biol.* **11**, 4627–34.
12. A. Yayon, M. Klagsbrun, J. D. Esko, P. Leder, and D. M. Ornitz (1991) Cell surface, heparin-like molecules are required for binding of basic fibroblast growth factor to its high affinity receptor, *Cell* **64**, 814–818.
13. J. Schlessinger, I. Lax, and M. Lemmon (1995) Regulation of growth factor activation by proteoglycans: What is the role of the low affinity receptors? *Cell* **83**, 357–360.
14. L. W. Burrus, M. E. Zuber, B. A. Lueddecke, and B. B. Olwin (1992) Identification of a cysteine-rich receptor for fibroblast growth factors. *Mol. Cell. Biol.* **12**, 5600–5609.
15. N. Kinoshita, J. Minshull, and M. W. Kirschner (1995) The identification of two novel ligands of the FGF receptor by a yeast screening method and their activity in *Xenopus* development, *Cell* **83**, 621–630.
16. N. Hacohen, S. Kramer, D. Sutherland, Y. Hiromi, and M. Krasnow (1998) Sprouty encodes a novel antagonist of FGF signaling that patterns apical branching of the *Drosophila* airways, *Cell* **92**, 253–263.
17. E. J. Williams, J. Furness, F. S. Walsh, and P. Doherty (1994) activation of the FGF receptor underlies neurite outgrowth stimulated by L1, N-CAM, and N-cadherin, *Neuron* **13**, 583–594.
18. A. O. Wilkie, G. M. Morriss-Kay, E. Y. Jones, and J. K. Heath (1995) Functions of fibroblast growth factors and their receptors, *Curr. Biol.* **5**, 500–507.
19. T. P. Yamaguchi, K. Harpal, M. Henkemeyer, and J. Rossant (1994) *fgfr-1* is required for embryonic growth and mesodermal patterning during mouse gastrulation, *Genes Dev.* **8**, 3032–3044.
20. E. Arman, R. Haffner-Krausz, Y. Chen, J. K. Heath, and P. Lonai (1998) Fibroblast growth factor receptor 2 (FGFR2) is required for formation of the egg cylinder, *Proc. Natl. Acad. Sci. USA* **95**, 5082–5087.
21. C. Deng, A. Wynshaw-Boris, F. Zhou, A. Kuo, and P. Leder (1996) Fibroblast growth factor receptor 3 is a negative regulator of bone growth, *Cell* **84**, 911–21.
22. G. A. Bellus, I. McIntosh, E. A. Smith, A. S. Aylsworth, I. Kaitila, W. A. Horton, G. A. Greenhaw, J. T. Hecht, and C. A. Francomano (1995) A recurrent mutation in the tyrosine kinase domain of fibroblast growth factor receptor 3 causes hypochondroplasia, *Nat. Genet.* **10**, 357–359.
23. A. O. M. Wilkie (1997) Craniosynostosis: genes and Mechanism. *Human Mol. Genet.* **6**, 1647–1656.

### **Suggestions for Further Reading**

24. I. J. Mason (1994) The ins and outs of fibroblast growth factors, *Cell* **78**, 547–552.
25. A. O. Wilkie, G. M. Morriss-Kay, E. Y. Jones, and J. K. Heath (1995) Functions of fibroblast

growth factors and their receptors, *Curr. Biol.* **5**, 500–507.

## Fibronectin

Fibronectin is a fibrillar [extracellular matrix](#) (ECM) **glycoprotein** that binds to cells and other matrix molecules, thereby contributing to the organization of the matrix and providing important adhesive sites for cells. Fibronectin is found throughout the vertebrate body, including the brain, which lacks a recognizable ECM but expresses ECM molecules like fibronectin transiently in early [development](#). Fibronectin is a dimer of 240-kDa monomers linked at the C-terminus by a pair of [disulfide bonds](#). The protein is made of similar but nonidentical repeating domains (types I, II, III), most of which are derived from a single exon (see [Introns, Exons](#)). The type III repeat, about 90 amino acid residues in length, occurs 15 times in each monomer and is a structural element found in other proteins as well. While the protein is encoded by a single gene, multiple splicing variants give rise to diverse **isoforms** that differ in their cell binding sites and functions. One isoform is found as an abundant protein in blood and is thought to participate in wound healing, phagocytosis by white blood cells, and [blood clotting](#). Other forms are incorporated into fibrils in the ECM. The regions of the protein that bind to matrix molecules and cells have been uncovered by using mild [proteinase](#) digestion and column [chromatography](#) to isolate and characterize the fragments. The fibronectin monomer has two binding sites each for heparin and **fibrin**, and a single binding site for [collagen](#). Further analysis of small synthetic peptides has led to the identification of the sequence -Arg-Gly-Asp- (RGD), which is bound by many cells via [integrin](#) receptors alpha-5/beta-1, alpha-3/beta-b1, alpha-v/beta-1, and alpha-v/beta-3. The RGD site occurs in a type III repeat on a loop extended above the surface of the domain, which is thought to provide easy access to cells. The RGD sequence is also found in other matrix proteins recognized by integrins, such as vitronectin. An **alternatively spliced** region (type III connecting segment or V) contains a second binding motif -Leu-Asp-Val- (LDV), which is the main part of the binding site for another integrin (alpha-4/beta-1) (see [Integrins](#)).

Fibronectin is one of the best-characterized ECM glycoproteins. *In vitro*, it supports adhesion of many cell types, and it can also stimulate cell proliferation, differentiation, and migration. Cells in culture often flatten and develop prominent [actin](#) arrays in response to fibronectin. The response depends on the particular cell type and its array of fibronectin receptors. The importance of fibronectin in migration of cells has been demonstrated for neural crest cells in birds and gastrulating embryonic cells of frogs, where injected [antibodies](#) and RGD-containing [peptides](#) will block cell migration. In genetically engineered mice lacking fibronectin, death occurs early in embryogenesis in part because of a failure of mesodermal cells to develop, resulting in defects in the development of the notochord, somites, blood vessels, and heart.

### Suggestions for Further Reading

1. R. O. Hynes (1989) *Fibronectins*, Springer-Verlag, New York.
2. C. Chothia and E. Y. Jones (1997) The molecular structure of cell adhesion molecules. *Ann. Rev. Biochem.* **66**, 823–862.

## Fibrous Proteins

Although the fibrous proteins represent a diverse collection of structures, they can nonetheless be grouped conveniently into four classes: (i) the  $\alpha$ -fibrous proteins, (ii) the  $\beta$ -fibrous structures, (iii) the [collagen](#) proteins, and (iv) those proteins that assemble into filamentous arrays but are individually globular in form. In general, all fibrous proteins contain strong repeating elements in their structures, often in the form of tandem sequence motifs. These motifs not only specify the **secondary structure**, but also play a recognizable and identifiable role in the assembly to higher levels of order. Fibrous proteins have often been used as model systems to provide insights into the structures of the more complex globular proteins.

### 1. $\alpha$ -Fibrous Structures

$\alpha$ -Fibrous proteins can exist *in vivo* either as individual elongated molecules (such as [laminin](#) in basement membranes and [fibrinogen](#) in blood plasma) or as filamentous assemblies (such as **myosin** in **muscle** thick filaments, **desmin** in Type III [intermediate filaments](#), and fibrin molecules in **blood clots**). The protein sequences of all members of this family have a high propensity to form  $\alpha$ -**helices**, which are stabilized by [hydrogen bonds](#) that lie approximately parallel to the helix axis. In addition, the  $\alpha$ -fibrous proteins contain a characteristic [heptad repeat](#) (see also [Coiled-Coils](#) and [Cytokeratins](#)) in which [nonpolar](#) residues alternate three and four residues apart. The apolar stripe on the surface of the right-handed  $\alpha$ -helix becomes internalized when several such  $\alpha$ -helices aggregate and assemble into a left-handed ropelike structure (1). Although **hydrophobic** forces are very important in driving the assembly, interchain [electrostatic interactions](#) contribute significantly to determining the relative chain orientations and axial stagger, as well as the stability (1-3). In these filament-forming structures, the amino acid sequence contains a highly regular linear disposition of acidic and basic residues. The periods are usually  $180^\circ$  out of phase, which generates a simple rod structure with alternating bands of positive and negative charge. Consequently, assembly into the filamentous form is specified in large part by maximization of intermolecular ionic interactions, rather akin to an ionic zipper. Those  $\alpha$ -fibrous proteins lacking periodicities in charged residues do not form regular filamentous assemblies, as judged by data currently available. In these cases, the molecules exist as separate entities, or assembly occurs in a less regular manner via non-helical domains elsewhere in the molecule, most typically at the *N*- and *C*-terminal ends.

### 2. $\beta$ -Fibrous Structures

The  $\beta$ -fibrous structures were first observed in silk proteins. The fundamental structure is that of an extended array of  [\$\beta\$ -strands](#) held together in a  [\$\beta\$ -sheet](#) by a regular disposition of hydrogen bonds perpendicular to the axes of the chains (see [Beta-Sheet](#)). Two or more such sheets then aggregate into a  $\beta$ -crystallite. Relatively few sequence data on silks are available, but it is likely that many of them have repeating amino acid sequence motifs. This is based on the observation that the amino acid compositions are generally very simple, with perhaps only a small number of amino acids, such as glycine, alanine, serine, and glutamine, represented in any significant number. On this basis it is thought that many silks contain a dipeptide repeat (or possibly a multiple), although there are non- $\beta$ -forming examples, such as *Nematus ribesii* and *Apis mellifera*, that almost certainly have three- and seven-residue repeats, because they give rise to a collagen and  $\alpha$ -fibrous X-ray [fiber diffraction](#) pattern, respectively. Those silks with a dipeptide repeat, however, automatically give rise to a  $\beta$ -sheet with the two residues in the repeat spatially separated on opposite sides. This can give a sheet with an apolar face, for example, which can then readily assemble with the same face in a second sheet. It has been assumed that the  $\beta$ -sheets are planar in the silk structures, but this may not be the case *in vivo*. In fact, it is known that the  $\beta$ -sheets are twisted, probably in a right-handed manner, for both feather and scale [keratins](#). In all fibrous proteins, the  $\beta$ -sheets are composed of antiparallel chains: in globular proteins both parallel and antiparallel  $\beta$ -sheets are found.

### 3. Collagen Class

The collagen class of fibrous proteins is characterized by a triplet repeat of the form (Gly–X–Y)<sub>n</sub>, where X and Y are often proline and hydroxyproline (see [Collagen](#)). Type I collagen molecules, each of length 300 nm, aggregate with an axial stagger of distance D (67 nm) or a multiple to generate fibrils. Collagen types II, III, V, and XI also form D-periodic fibrils, but type VI collagen forms fibrils with a different periodicity. Other collagen types (IV, VII, VIII, IX, X, XII, XIII, and XIV) are non-fibril-forming but can be classified further as either basement membrane collagens (IV, VII), short-chain collagens (VIII, X), or fibril-associated collagens (IX, XII, XIV).

#### 4. Filamentous Arrays

The last group of fibrous proteins comprise proteins that are globular in shape but assemble either helically or as a long string to form, respectively, a regular or semiregular filamentous structure. A good example of the former system is [actin](#), in which G-actin (globular) assembles into F-actin filaments (fibrous). These, together with other proteins such as **tropomyosin** and **troponin**, form the [thin filaments](#) of muscle. **Titin**, a protein from muscle, consists of a series of globular domains strung together and is thus an example of the second type of filament-forming fibrous protein. The short arms in [laminin](#) are likewise composed of globular regions strung together in a near linear manner.

#### Bibliography

1. C. Cohen and D. A. D. Parry (1990)  $\alpha$ -helical coiled-coils and bundles: how to design an  $\alpha$ -helical bundle. *Proteins Struct. Funct. Genet.* **7**, 1–15.
2. D. Krylov, I. Mikhailenko, and C. Vinson (1994) A thermodynamic scale for leucine zipper stability and dimerization specificity: e and g interhelical interactions. *EMBO J.* **13**, 2849–2861.
3. O. D. Monera, C. M. Kay, and R. S. Hodges (1994) Electrostatic interactions control parallel and antiparallel orientation of  $\alpha$ -helical chains in two-stranded  $\alpha$ -helical coiled-coils. *Biochemistry* **33**, 3862–3871.
4. D. J. S. Hulmes, A. Miller, D. A. D. Parry, K. A. Piez, and J. Woodhead-Galloway (1973) Analysis of the primary structure of collagen for the origins of molecular packing. *J. Mol. Biol.* **79**, 137–148.

#### Suggestions for Further Reading

5. R. D. B. Fraser and T. P. MacRae (1973) *Conformation in Fibrous Proteins and Related Synthetic Polypeptides*, Academic Press, London.
6. D. J. S. Hulmes (1992) The collagen superfamily—diverse structures and assemblies. *Essays Biochem.* **27**, 49–67.

#### Fill-In Reaction

DNA fragments with 5'-overhanging ends are converted to **blunt-ended** molecules by reacting with the **Klenow fragment** of *Escherichia coli* DNA polymerase. This makes the overhanging ends double-stranded and is called the fill-in reaction. The fill-in reaction is often used to join fragments with different **cohesive ends** by converting the cohesive ends to blunt ends. This reaction is also used to introduce a [frameshift mutation](#) into a **gene**. A **plasmid** containing the gene is cut within the gene using a [restriction enzyme](#) to produce ends with a 5'-overhang. The overhang is filled in and the

ends are **religated**. As long as the length of the overhang is not divisible by three, this procedure changes the [reading frame](#) of the gene [downstream](#) of the restriction site (see **Protein biosynthesis**).

## Filter Binding Assays

The term *filter binding assay* is used to describe a variety of different techniques in biochemistry, immunology, virology and molecular biology. In molecular biology a filter binding assay is used to characterize **DNA–protein interactions**. This particular technique is also referred to as a *filter assay* or a *protein binding assay*, and it can be used to identify and characterize both DNA-binding proteins and the DNA sequences that interact with these proteins. The filter binding assay is based on the observation that proteins, but not double-stranded DNA molecules, bind to the surface of [nitrocellulose](#) membrane filters ([1](#), [2](#)). Therefore, if a DNA-binding protein is incubated with a specific DNA sequence before analysis in the filter binding assay, the resulting DNA–protein complex and any excess protein will bind to the membrane filter, while uncomplexed DNA will pass through. The amount of the DNA–protein complex retained on the membrane filter can be determined by using **radiolabeled** DNA to form the DNA–protein complexes and to quantify the amount of radioactivity retained on the nitrocellulose membrane by scintillation counting. The filter binding assay can be very quick if a multifilter vacuum filtration unit is used.

The filter binding assay offers a simple, rapid, sensitive, and versatile assay capable of providing information about the interactions between a DNA-binding protein and a specific DNA sequence. The use of radioactive DNA makes low concentrations of DNA and protein feasible. A typical experiment involves the titration of a constant amount of radiolabeled DNA with the DNA-binding protein of interest. The amount of radioactivity retained on the membrane increases with increasing protein concentration until all the radioactive DNA that can bind to the specific DNA-binding protein is depleted from the incubation mixture. Binding affinities can be determined by observing the level of specific complexes formed in the presence of a nonradioactive competitor DNA. The kinetics of association and dissociation of a specific DNA–protein complex can also be studied in the filter binding assay. However, filter binding assays cannot differentiate between dissimilar DNA–protein complexes, so this assay is effective only when one DNA–protein complex is formed. In addition, the DNA-binding protein must not be denatured during the purification process, because the filter binding assay is dependent on the interaction between the DNA-binding protein and a specific DNA sequence.

Because various membrane filters differ in binding specificities, the ability of nitrocellulose membranes to bind proteins and not native DNA is critical to the success of the filter binding assay. However, the molecular basis for this discrimination is not well understood. Both **electrostatic** and **hydrophobic** interactions have been suggested to be involved in the binding of macromolecules to nitrocellulose, but hydrophobic interactions are assumed to play the dominant role in the binding process (see [Nitrocellulose](#)). Although native DNA does not bind to nitrocellulose, heat-denatured DNA will ([3](#)) and should not be used in this assay. In addition, the native DNA tested in the assay must be free of contaminating proteins that could bind the DNA to the nitrocellulose membrane in the absence of the DNA-binding protein being tested.

Most proteins bind to nitrocellulose, but individual binding affinities are dependent on the surface characteristics of the particular protein. Therefore, the retention of a particular DNA–protein complex on the nitrocellulose filter will depend on the surface characteristics of the protein, the binding capacity of the nitrocellulose membrane, the time that the DNA–protein complex has to interact with the filter, and the regimen used to wash the filter. A rapid flow rate through the filter may not permit the binding of some DNA–protein complexes to the nitrocellulose membrane, and



extensive washing of the filter may remove DNA–protein complexes with low affinities for nitrocellulose. The optimum conditions needed to bind a DNA–protein complex to the nitrocellulose membrane will be different for each protein and should be investigated for each filter binding assay.

### Bibliography

1. O. W. Jones and P. Berg (1966) *J. Mol. Biol.* **22**, 199–209.
2. A. D. Riggs, H. Suzuki, and S. Bourgeois (1970) *J. Mol. Biol.* **48**, 67–83.
3. A. P. Nygaard and B. D. Hall (1963) *Biochem. Biophys. Res. Commun.* **12**, 98–104.

### Suggestions for Further Reading

4. O. Papoulas (1987) "Rapid separation of Protein-Bound DNA from Free DNA Using Nitrocellulose Filters". In *Current Protocols in Molecular Biology*, Vol. **2** (F. M. Ausubel et al., eds.), Wiley, New York, pp. 12.8.1–12.8.9.
5. P. G. Stockley (1994) "Filter-Binding Assays". In *Methods in Molecular Biology*, Vol. **30** (G. G. Kneate, ed.), Humana Press, Totowa, NJ, pp. 251–262.
6. A. Revzin (1996) "Nucleic Acid–Protein Complexes". In *and Molecular Medicine*, Vol. **4** (R. A. Meyers, ed.), VCH, Weinheim, Germany, pp. 243–253.

## Fingerprint, Protein

Two-dimensional [peptide mapping](#), or protein fingerprinting, was introduced by Ingram in 1956 (1). In this approach, the protein is cleaved by **trypsin** or another [proteinase](#) to obtain complete digestion into peptides. Then, the peptides are separated on a solid support by [electrophoresis](#) in the first dimension and [chromatography](#) in the second, to produce a two-dimensional pattern or map. The position of each peptide is determined by its mobilities in the two dimensions and is sensitive to its covalent structure. The advantage of a two-dimensional separation is that many peptides are separated simultaneously, and two different mobilities of each can be examined.

Fingerprints are most useful in comparing closely related proteins because the one or a few peptides that differ are often easy to detect. They end up at altered positions, whereas the majority of the peptides are unaltered and remain at the same positions in each map. The latter serve as references in detecting spots corresponding to variant peptides. The deviating peptides are eluted from the support and characterized by [amino acid analysis](#), by sequence determination, as with the [Edman Degradation](#), or by [mass spectrometry](#). This approach requires that a significant number of the peptides generated have identical sequences. Thus, two-dimensional peptide mapping is used to evaluate the extent of identity among closely homologous proteins.

A classical example of protein fingerprinting is the early work by Ingram on the **sickle cell** form of [hemoglobin](#) (1). Tryptic mapping revealed an altered position for only one pair of peptides detected after two-dimensional separation of the digests of normal and sickle cell hemoglobins. Glu6 of the normal b-chain is replaced by Val in the sickle form because of an A to T mutation in the corresponding **codon** of the b-chain gene (2).

The protein substrates for digestion are purified in solution by conventional chromatographic techniques or isolated via [SDS-PAGE](#), followed by [electroelution](#) (3). The protein is **precipitated** with **trichloroacetic acid** (30% (w/v)) at 4°C for 4 h to desalt the preparation and to concentrate the protein in a small volume before digestion. After centrifugation, removal of the supernatant, and

gently washing with acetone the “pellet” at  $-20^{\circ}\text{C}$  (the pellet is not normally visible), the appropriate digestion buffer is added, followed by the proteolytic enzyme. This is efficient if the amount of protein is greater than about 1 nmol. The protein substrate may also be digested directly in the SDS/polyacrylamide gel without previous electroelution (see [Cleveland Map](#)). It is detected by brief staining with [Coomassie Brilliant Blue](#) or by visualization using 1 M KCl, followed by excision and treatment of the gel slice by a combination of drying and rehydration in digestion buffer that contains the proteinase ([4](#), [5](#)). Small, proteolytically cleaved peptides will diffuse out, while undigested or only partially digested proteins remain in the gel piece. The outer liquid that contains the peptides is concentrated by **lyophilization** to a small volume suitable for application to the first dimension of the mapping. Frequently, the amount of peptide material is insufficient for visualization of the map via standard staining (using [ninhydrin](#)). In this case **radiolabeling** of the protein or detection by **fluorescence** techniques is necessary.

Trypsin is often used for proteolysis because of its high specificity for cleavage at [lysine](#) and [arginine](#) residues, which can be taken to completion. An enzyme to substrate ratio of about 1:50 (w:w) is often used in a volatile [buffer](#) of ammonium bicarbonate, pH 8.1, at  $37^{\circ}\text{C}$  for 2 to 24 h. Even large polypeptides are completely degraded into small peptides that are separated by the combination of electrophoresis and chromatography. However, for proteins that generate a large number of peptides, radiolabeling and [autoradiography](#) are often necessary to obtain a distinct pattern without too many spots that overlap or cluster. For this purpose, one frequently employs labeling of [cysteine](#) residues by carboxymethylation with [ $^{14}\text{C}$ ]-**iodoacetic acid**, of [tyrosine](#) residues by **iodination** with  $^{125}\text{I}$ , or of proteins **phosphorylated** with  $^{32}\text{P}$ .

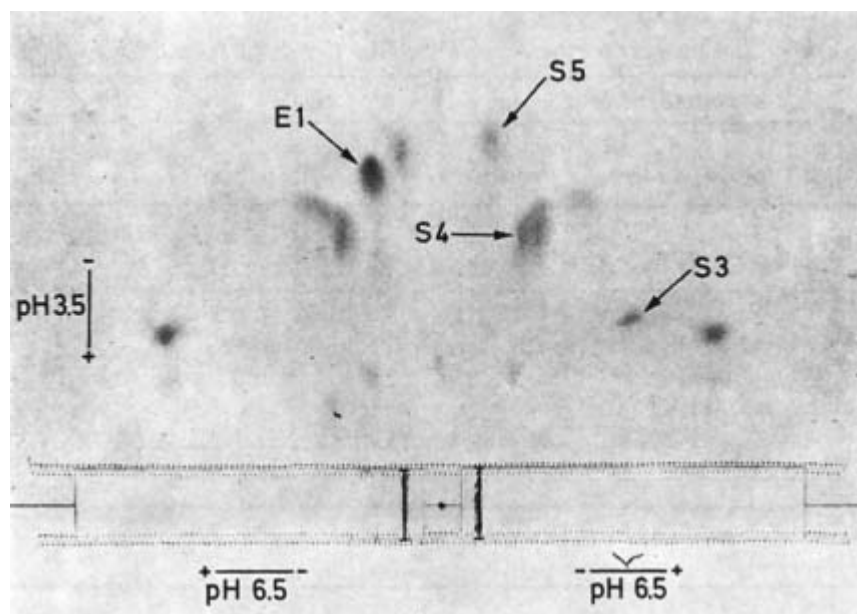
Electrophoresis is the normal separation in the first dimension before ascending **thin-layer chromatography** in the second. The glass or plastic support is covered by a thin silica or cellulose layer. The electrophoresis buffers should contain volatile components. Many solvent systems have been suggested for the second dimension ([6](#)), but they have an organic solvent in common, often an alcohol, to promote the chromatographic mobility of **hydrophobic** peptides. A classical solvent is composed of butanol, acetic acid, water, and pyridine in volume ratios of 15:3:12:10 at pH 4 to 5. After drying, the fingerprint is visualized by ninhydrin, fluorescent techniques, stains for specific residues, or autoradiography.

Fingerprinting gives information on the subunits and the homogeneity of a protein. The number of total spots, or of those containing Cys, Tyr or other residues that are detected specifically, can be compared with the number of peptides expected from [amino acid analysis](#) and the molecular weight of the protein. If the sample is homogeneous and composed of a single polypeptide chain, the number of peptides generated should fit the frequency expected from the composition. If the number of spots detected is substantially lower, it probably indicates that the protein is composed of multiple, identical, polypeptide chains. If it is significantly higher, the sample has more than one individual polypeptide.

Two-dimensional peptide mapping is also efficient in analyzing structural variants of the same polypeptide, that is, protein **isozymes** (Fig. [1](#)), or in identifying truncated forms and other types of [posttranslational modification](#). Then, peptide maps reveal whether smaller or deviating polypeptides are related to other forms of the protein (identical peptide maps except for missing or deviating fragments) or are unrelated (no peptide similarities). Finally, two identical dimensions of peptide mapping, carried through with an intermediate step of chemical modification to alter a specific amino acid type, are used to detect all peptides containing that type of residue. This was the principle for the previously so popular “**diagonal**” methods ([7](#)) based on separations on paper or thin-layer plates. The peptides that contain the modified residue are differentiated by modifying them covalently after the first dimension. The second dimension is identical to the first, so that all of the unmodified peptides have the same mobilities in both dimensions. Only the modified peptides lie running off the diagonal formed by the unaltered peptides. Now, this approach like all original mappings, can be extended to modern situations by two identical runs of HPLC, capillary

electrophoresis, or other microseparation. Again, the principle is always to detect any forms that deviate between the two separations.

**Figure 1.** Two-dimensional peptide map of acidic peptides of two similar isozymes (with E and S subunits, left and right, respectively) of an [alcohol dehydrogenase](#). Resolution was by paper electrophoresis at pH 6.5 (horizontal) and at pH 3.5 (vertical). Strips of paper from the first dimension were stitched in opposite polarity so that the maps form mirror images through a vertical center line. The spots that differ in the E and S maps are marked (E1 and S3 are clear extra spots; S4 and S5 have only slightly shifted positions). From Ref. 8 with permission.



## Bibliography

1. V. M. Ingram (1956) *Nature* **178**, 792–794.
2. V. M. Ingram (1957) *Nature* **180**, 326–328.
3. M. W. Hunkapiller et al. (1983) *Methods in Enzymology* **91**, 227–236.
4. J. Rosenfeld et al. (1992) *Anal. Biochem.* **203**, 173–179.
5. U. Hellman et al. (1995) *Anal. Biochem.* **224**, 451–455.
6. W. J. Gullick (1986) In *Practical Protein Chemistry: A Handbook* (A. Darbre, ed.), Wiley, Chichester, U.K., pp. 207–225.
7. J. R. Brown and B. S. Hartley (1966) *Biochem. J.* **101**, 214–218.
8. H. Jörnvall (1970) *Eur. J. Biochem.* **16**, 41–49.

## Fingerprinting DNA

“DNA fingerprinting” refers to methods of detecting, in **eukaryotes**, unique **DNA** patterns, which allow the identification of individuals with a probability of error similar to (or lower than) that obtained by comparing fingerprints in humans. These unique, individual patterns of DNA are the

result of Mendelian inheritance of polymorphic, **hypervariable loci** of [repetitive DNA](#). The most useful loci are those consisting of tandem repeats of short (15 to 60 bp) or very short (3 to 5 bp) specific base sequences. Different **alleles** are produced by variation in the number of repeats per locus, an event that arises by **unequal crossing over**, “slippage” during [DNA replication](#), or other means (see text below). These phenomena are favored by tandem repetition, and they are responsible for a [mutation](#) frequency several times greater than that of “traditional” point or chromosomal mutations. They can give rise to arrays of up to several hundred such repeated units.

The terms “**minisatellites**” and “**microsatellites**” are commonly used to define these loci. The names are derived from the general definition of “satellite” for repetitive DNA, because its buoyant density frequently differs from that of the bulk DNA, due to a shift in its average base composition (G and C are usually present in increased amounts). This alters the buoyant density of such DNA and causes it to separate from the bulk of fragmented DNA as a separate satellite band in CsCl buoyant [density gradient centrifugation](#) .

Historically, DNA fingerprinting was introduced in 1985 in the article “Hypervariable minisatellite regions in human DNA” by A. Jeffreys et al. (1). These authors discovered a 33-bp region repeated four times in an **intron** of the human myoglobin **gene**. The same region was found scattered randomly in other parts of the human genome, as well as in the genomes of most vertebrates. The number of repeats of this unit varies widely at different loci, but the basic 33-bp element can always be recognized, mainly from the presence of a stable 16-bp “core.”

Many other loci sharing the properties of Jeffreys’ first minisatellite were found to be scattered in the genomes of most higher organisms (2-4). Arrays belonging to the same “family” of minisatellites can be identified by the specific, stable core sequence present in every one of its tandemly repeated units, the polycore . The very great variability in the length of minisatellites in different locations of a single genome, and among genomes, represents a highly reliable means of individual identification. Its best use is in comparing individuals of single families or very small populations (see text below).

Technically, individual DNA patterns are produced by digesting genomic DNA with one or more [restriction enzymes](#) whose recognition sites (preferably not exceeding four bases) are not present in the repeated unit. Minisatellites will thus survive digestion, and the [restriction fragments](#) containing them can be resolved by agar [gel electrophoresis](#) according to their length (ie, molecular weight). They can then be transferred to a nylon membrane by [blotting](#) and can be identified by **hybridization** with an appropriate labeled oligonucleotide probe (see **Southern blotting**). A sensitive probe is a tandem repeat of sequences complementary to the polycore of the minisatellite under study. The labeling is either by [radioactivity](#) or by one of several enzymatic reactions leading to a visible product (see **ELISA**). The result is a ladder of DNA restriction fragment bands containing the polycore sequence, representing DNA fragments of different molecular weights, which is referred to as a “DNA profile.” Different individuals normally have different DNA profiles, unless they are related genetically.

Comparison of different profiles can only be made between closely similar electrophoretic lanes; it is thus essential that individual profiles to be compared are run on the same gel. The index of similarity commonly used in comparing profiles is the “band sharing coefficient ,” which gives the probability that a band of a given molecular weight will be shared within a fingerprint (5). Shared bands are those that have migrated in the gel the same distance and are therefore aligned in adjacent lanes.

Although the stretch of DNA that can be amplified by the polymerase chain reaction (**PCR**) is much shorter than the longest minisatellite variants, fingerprints have recently been produced by PCR using arbitrary primers (6, 7). This technique has proved simpler and has provided more information on types and frequencies of mutations occurring within individual minisatellites.

## 1. Microsatellites

Microsatellites represent a category of short tandem repeats whose length does not exceed 2 to 5 bp; they share the variability in repeat number of minisatellites, although not quite so great, and always remain within the limits of length that can be comfortably amplified by PCR. Technically, this is a great advantage; it only requires that the flanking sequences be known, to provide the primers required for PCR. Microsatellites are abundant in higher vertebrates, scattered in the genome at individual loci. In contrast with Jeffreys' multilocus minisatellites, microsatellites tend to be species- or taxon-specific; their abundance, however, frequently allows one to find some that are present in most species of a rather high taxon group. They have been mentioned here because of their structural similarity with minisatellites, but their use is different and will not be described here.

## 2. Multilocus and Unilocus DNA Fingerprints

The above discussion has been of multilocus fingerprinting, where the minisatellite is present at multiple loci within the genome. Some minisatellites are, however, present at only a single locus. Variability in their repeat number is still very great, however, causing heterozygosity to be the rule and a very high polymorphism in repeat number among individuals. The distinction between the two types of multilocus and unilocus minisatellites is not clear cut, however. If a family of minisatellites is hybridized with a probe at high stringency (high temperature of annealing or high salt concentration), the number of hybridizing bands in the gel will be decreased, in some cases, to a single one. This occurs when the repeats are similar, but not identical; random point mutations may cause single repeats to diverge slightly and to hybridize with the probe only under conditions of low stringency.

Hypervariable single-locus minisatellites have a number of advantages over multilocus fingerprints. The main advantage is probably their straightforward identification in different individuals and their unambiguous recognition as alleles of the same locus of any variable band identified by a specific probe. This makes single-locus fingerprinting more suitable than multilocus for population analysis.

The isolation of single-locus minisatellites is achieved by separating on the basis of size DNA fragments from a genome expected to be rich in minisatellites (G-C rich), cloning them in **cosmid vectors** made to accept inserts within a given range of sizes (charomids), and propagating them in rec-minus bacteria, to overcome the instability typical of any tandem-repeated sequence (8). Specific polycore probes will identify individual minisatellite loci in the vectors plated at low density. Those that hybridize with the probe at high stringency are potential single-locus minisatellites, which can be isolated and sequenced to construct probes that will be very specific.

## 3. Applications

The best use of multilocus fingerprinting is in the identification of individuals. The probability that two unrelated individuals share the same electrophoretic pattern is exceedingly low. The proportion of electrophoretic bands shared by two unrelated members of a population,  $a$  and  $b$ , is  $x = ((N_{ab}/N_a) + (N_{ab}/N_b))/2$ , where  $N_{ab}$  is the number of bands common to the two individuals and  $N_a$  and  $N_b$  are the total number of bands observed in the two individuals (usually up to 30). The band sharing coefficient,  $x$ , is related to the allele frequency  $q$  by the equation  $x = 2q - q^2$ , assuming that shared bands are identical alleles from the same locus and that all bands have the same population frequencies. If the bands are indeed independent markers, the mean probability that all  $N$  bands in an individual's profile are present in a second unrelated individual is  $x^N$ . For an individual with 30 scorable bands in a population with an average  $x = 0.2$ , the probability of another unrelated individual sharing the same pattern is less than  $10^{20}$ .

These characteristics prompted the forensic use of DNA fingerprinting in criminal and civil areas (9, 10). Identification of individuals from traces of biological material is possible, given the small amount (nanograms) of DNA required for a profile and the resistance to degradation of DNA, even

under the most unfavorable environmental conditions.

Multilocus fingerprinting is also the method of choice for determining close genetic relatedness ( $r = 0.5$  or, at most,  $0.25$ ). The Mendelian segregation of bands ensures that a progeny inherits, on average, a random 50% of its bands from each parent. Extra bands (ie, those not present in either parent) are found with a frequency ranging from  $10^2$  to  $10^4$ , as a result of the high but variable mutation rate of minisatellites of different sizes. In practice, both in the forensic field and in fields dealing with animal behavior, the most common problem easily solved by multilocus fingerprinting is the assessment of paternity. When the profiles of all partners (mother, progeny, and presumed father) are available, paternity can be affirmed or rejected with probabilities that depend on the total number of bands, on the proportion of bands shared, on the presence of new bands (mutations), and on the assumption of independence among loci. In individual cases, neither the probability of mutation nor that of loci not being independent severely affect the diagnosis. For example, the probability of not detecting incorrect paternity (ie, that the putative father will possess by chance say six of the paternal specific bands) is of the order of  $5 \times 10^5$  (so long as the putative father is not related to the true one).

In the animal field, this method of paternity testing has shown how widespread is the phenomenon of “extra-pair copulation” in species classified as monogamous ([11](#)).

In population studies, multilocus fingerprinting is not the best method of analysis. The average band sharing coefficient between pairs of random, unrelated members of even a large population depends upon both inbreeding and polymorphism. Specifically, high band sharing could indicate that the population is partly inbred (ie, is not a random-mating one) or that it was subjected to a recent population bottleneck that reduced drastically the number of individuals and the range of minisatellite's sizes still represented.

Population-level comparisons may be carried out only among populations effectively so small that individual variability is drastically reduced and population-specific genotypes may emerge. One of the main problems in population analysis is the impossibility of assigning bands to a specific locus: Bands shared (co-migrating) by different individuals are not necessarily identical alleles of a single locus. This problem is overcome by using single-locus probes: If more than one is available for the same study, the total amount of information (ie, number of usable bands) will approach that of multilocus fingerprinting. Presently, microsatellites have supplanted single-locus minisatellites in population genetic analysis: The type and amount of information obtained is similar, but the method of revealing microsatellite patterns is much simpler and faster.

Single-locus minisatellites found to segregate in close association with genetic traits have been very useful in the mapping of these traits ([12](#)) and, sometimes, in the isolation of the linked gene. In human medical genetics, single-locus minisatellites have been used as markers of potentially hereditary diseases, when linkage could be shown between the two in sufficiently extended pedigrees. Also, synthetic oligonucleotides, especially if rich in G and C, have been used as probes for detecting minisatellites to be used for mapping human genetic diseases ([13](#), [14](#)). Again, microsatellites have now entirely supplanted minisatellites for this purpose.

## Bibliography

1. A. J. Jeffreys et al. (1985) *Nature* **314**, 67–73.
2. T. Burke and M. W. Bruford (1987) *Nature* **327**, 149–152.
3. A. J. Jeffreys (1987) *Anim. Genet.* **18**, 1–15.
4. M. Georges et al. (1988) *Cytogenet. Cell Genet.* **47**, 127–131.
5. M. Lynch (1990) *Mol. Biol. Evol.* **7**, 478–484.
6. J. Welsh and M. McClelland (1990) *Nucleic Acids Res.* **18**, 7213–7218.
7. J. Welsh, N. Rampino, M. McClelland, and M. Perucho (1995) *Mutat. Res.* **338**, 215–229.

8. J. A. L. Armour, S. Povey, S. Jeremiah, and A. J. Jeffreys (1990) *Genomics* **8**, 501–512.
9. B. E. Dodd (1985) *Nature* **318**, 506–507.
10. A. M. Ross and H. W. J. Harding (1989) *Forensic Sci. Int.* **41**, 197–203.
11. T. Burke (1989) *TREE* **4**(5), 139–144.
12. T. G. Krontiris (1995) *Science* **269**, 1682–1683.
13. Y. Nakamura et al. (1987) *Science* **235**, 1616–1622.
14. R. Schafer et al. (1988) *Electrophoresis* **9**, 369–374.

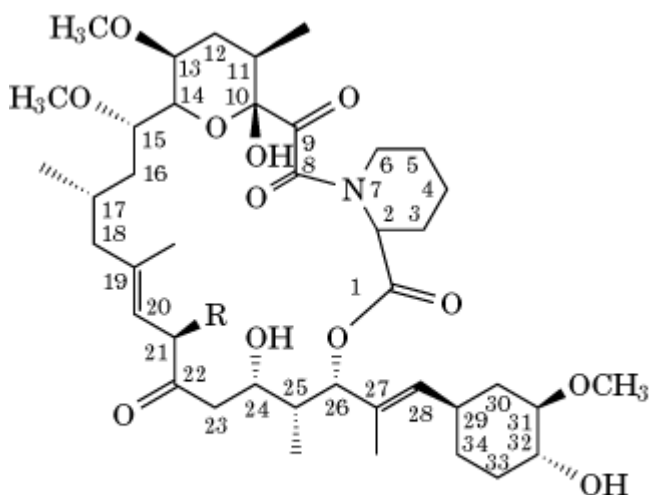
### Suggestions for Further Reading

15. L. T. Kirby (1990) *DNA Fingerprinting. An Introduction*, Stockton Press, New York.
16. A. J. Jeffreys et al. (1991) "Principles and recent advances in DNA fingerprinting". In *DNA Technology and Its Forensic Applications* (G. Berghaus, B. Brinkmann, G. Rittner, and M. Staak, eds.), Springer-Verlag, Berlin, pp. 3–19.
17. M. W. Bruford, O. Hanotte, J. F. Y. Brookfield, and T. Burke (1992) "Single locus and multilocus DNA fingerprinting". In *Molecular Genetic Analysis of Populations. A Practical Approach* (A. R. Hoelzel ed.), IRL Press, Oxford, U.K., pp. 225–269.
18. K. Rassmann, A. Zischler, and D. Tautz (1996) "DNA multilocus fingerprinting using simple repeat motif oligonucleotides". In *Molecular Genetic Approaches in Conservation* (T. B. Smith and R. K. Wayne, eds.), Oxford University Press, New York, pp. 238–250.

## FK506

The immunosuppressive drug FK506 (also known as tacrolimus) is a 23-membered macrocyclic lactam (Fig. 1) that was first isolated from a strain of *Streptomyces tsukubaensis* in Japan in 1987 (1). In addition to blocking the T cell immune response, FK506 has neuroprotective properties (2, 3) and antiparasitic effects (4). It was first used clinically for immunosuppression in 1989. The target of FK506 for immunosuppression is thought to be the cytosolic FK506-binding protein, FKBP12, which is also a **peptidyl prolyl *cis*/ *trans* isomerase** (PPIase). Upon binding to FKBP12 within a deep **hydrophobic** pocket, the **imide** bond of FK506 adopts a *trans* conformation, concomitant with reorganizing the conformation of the pipercolinic ring and of the pyranose ring. Recently, a cell-permeable, covalently linked dimer of FK506 (FK1012) was constructed as a tool to achieve controlled cross-linking of intracellular receptor domains for gene activation (5). It utilizes the high binding affinity of FK506 for the receptor FKBP12, which is effective at less than nanomolar concentrations. A number of derivatives that bind to and inhibit FKBP have been reported, some immunosuppressive and others not (6).

**Figure 1.** The structure of immunosuppressant FK506.



FK506      R: CH<sub>2</sub>CH=CH<sub>2</sub>  
 ascomycin    R: CH<sub>2</sub>CH<sub>3</sub>

Despite structural similarities to [rapamycin](#) in a major part of the macrocycle, many biological effects of FK506, including the arrest of clonal T cell expansion, are reminiscent of another class of immunosuppressants, [cyclosporin A](#) (CsA), rather than rapamycin. This finding was the first evidence of a common upstream target for the two different classes of immunosuppressants, FK506 and cyclosporin A. As monitored by [affinity chromatography](#), the FKBP12/FK506 complex binds the same protein from a calf thymus extract, the protein **phosphatase** calcineurin, as has been obtained with the complex of CsA and its target [cyclophilin](#) Cyp18. In both cases noncompetitive inhibition by the complex is observed in the nanomolar concentration range. Considering that the structures of Cyp18/CsA and FKBP12/FK506 are unrelated in amino acid sequence and spatial organization, the similar inhibition of calcineurin is striking. FK506 acts as a pro-drug, in that only the complex of the drug and its receptor is active. Two segments of FKBP12 comprising the loops in the region of residues 40 and 80 plus an effector region of FK506, the carbon chain C17 to C32, may be involved in drug activation. In the **X-ray crystallographic** structure of the human FKBP12/FK506/calcineurin complex, the FKBP12/FK506 portion contacts two distinct areas on calcineurin that do not directly obstruct both the active site and the substrate-binding cleft of the protein phosphatase. On the contrary, blocking a rather distant subsite of the phosphatase/protein phosphate recognition region is responsible for phosphatase inhibition (8).

Both CsA and FK506 block the induction of clonal T cell expansion at the early G<sub>0</sub> to G<sub>1</sub> transition of the cell cycle. They exert down-regulation of [transcription](#) of a number of similar [lymphokine](#) genes.

#### Bibliography

1. T. Kino et al. (1987) *J. Antibiot.* **40**, 1249–
2. W.E. Lyons et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 3191–3195.
3. J. Sharkey and S.P. Butcher (1994) *Nature* **371**, 336–339.
4. A. Moro et al. (1995) *EMBO J.* **14**, 2483–2490.
5. D.M. Spencer, T.J. Wandless, S.L. Schreiber, and G.R. Crabtree (1993) *Science* **262**, 1019–1024.
6. D.M. Armistead and M.W. Harding (1993) *Ann. Rep. Med. Chem.* **28**, 207–215.
7. J. Liu et al. (1991) *Cell* **66**, 807–817.
8. C.R. Kissinger et al. (1995) *Nature* **378**, 641–644.



## Suggestions for Further Reading

9. J.R. Luly (1991) Mechanism-based immunosuppressants, *Ann. Rep. Med. Chem.* **26**, 211–219.
10. T.R. Brazelton and R.E. Morris (1996) Molecular mechanisms of action of new xenobiotic immunosuppressive drugs - tacrolimus (FK506), sirolimus (rapamycin), mycophenolate (mofetil) and leflunomide, *Curr. Opin Immunol.* **8**, 710–720.
11. S.H. Snyder and D.M. Sabatini (1995) Immunophilins and the nervous system, *Nature Med.* **1**, 32–37; summarizes immunophilins in the nervous system.
12. M.E. Cardenas, D. Zhu, and J. Heitman (1995) Molecular mechanism of immunosuppression by cyclosporine, FK506 and rapamycin, *Curr. Opin Nephrol. Hypertens.* **4**, 472–477.

## Flagella—Prokaryotes

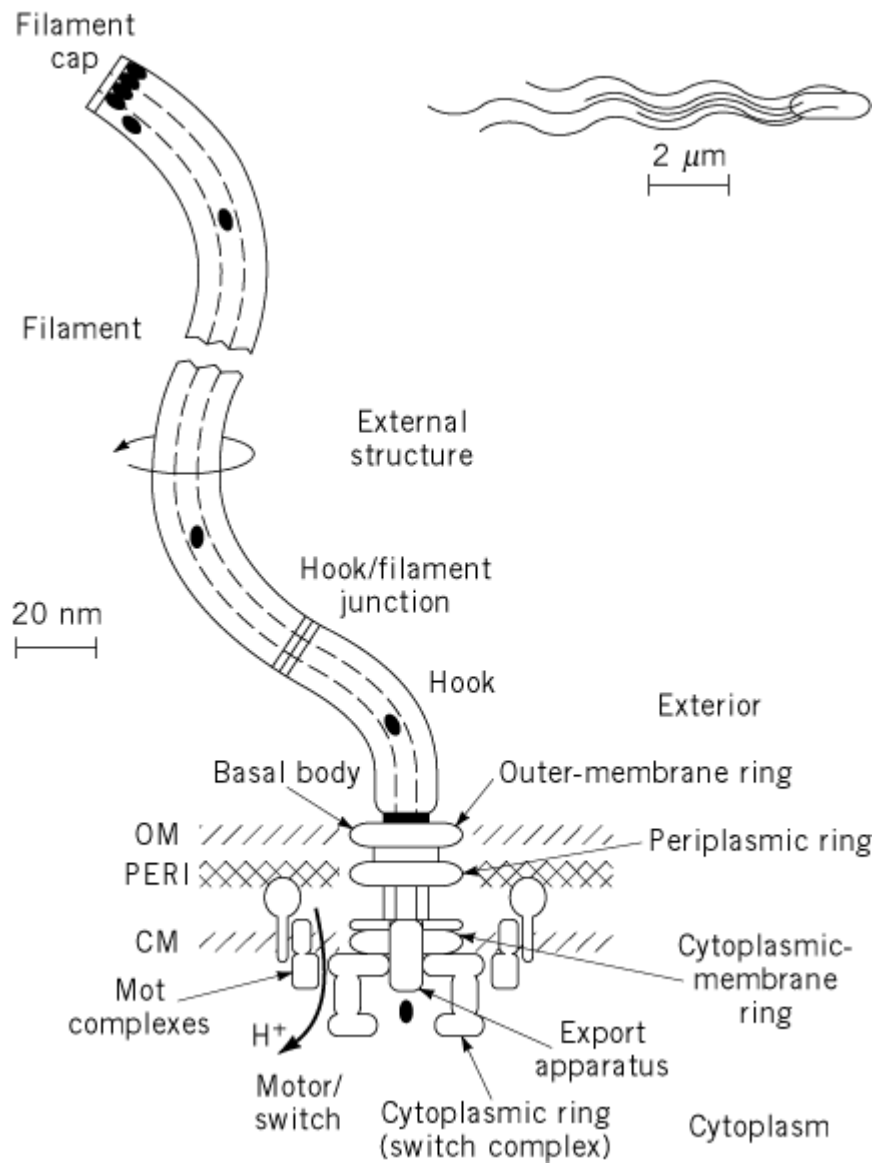
Flagella are the propulsive **organelles** of a wide variety of unicellular organisms, both **prokaryotes** and **eukaryotes**. However, prokaryotic flagella differ in almost every respect from their eukaryotic counterparts (see **Flagella—eukaryotes**).

### 1. Motility and Behavior

Prokaryotic flagella are largely external appendages that project far from the cell surface (Fig. 1, top right); for example, a typical rod-shaped **bacterium** is about 1 or 2  $\mu\text{m}$  long, but its flagellar filaments can be 10  $\mu\text{m}$  long. The motor resides at the cytoplasmic membrane and is a rotary device that transmits the rotation to the external flagellar filament. The motor is not driven by ATP hydrolysis like most mechanoenzymes (see **Energy transduction**), but by a [proton-motive force](#) (or in some marine and alkaliphilic bacteria, by a sodium-motive force); protons that have been actively pumped out of the cell by the electron transport chain are allowed to flow down their electrochemical potential gradient through the flagellar motor and surrender their potential energy in return for hydrodynamic work against the resistance of the aqueous medium (see also [Chemiosmotic Coupling](#)). The filament is very thin (about 20 nm in diameter) compared to its length, with a length/diameter ratio of about 500, and has a completely regular helical or corkscrew shape. When rotated by the motor, it therefore functions as a propeller or Archimedes screw. Depending on the handedness of the helix (left-handed or right-handed) and the direction of rotation of the motor (counterclockwise or clockwise), the cell will experience a pushing or a pulling force. The flagellar motors of most bacteria are reversible, that is, they possess a switch or gear that will alternate the direction of rotation between counterclockwise and clockwise rotation. In the simplest cases, such as *Pseudomonas* spp., which have a single flagellum at one pole of the cell, this results in a shuttling motion. In the case of species such as *Escherichia coli* or *Salmonella*, which have several flagella originating more or less randomly around the cell body, the rotation direction that corresponds to pushing results in the individual flagellar filaments coalescing into a propulsive bundle to produce the motion called *running* or *swimming*, while the opposite direction causes the bundle to fly apart and produce a chaotic motion called *tumbling*. This tumbling randomly orients the cell, so that the next episode of swimming is largely uncorrelated with the previous one. The resulting trajectory is a zigzag or random walk in three-dimensional space.

**Figure 1.** The bacterial flagellum, showing its four major structural elements: external structure (white), basal body

(light gray), motor/switch (medium gray), and export apparatus (dark gray), along with their various substructures. The outer membrane (OM), periplasmic space (PERI), and cytoplasmic membrane (CM) of the bacterium, in which the flagellum machinery is embedded, are also shown. To drive motor rotation, protons flow through the motor down their potential gradient (toward lower  $H^+$  concentration). External proteins, such as flagellin (black ovals), are selectively passed through the export apparatus, travel down a central channel in the nascent structure, and assemble at its distal end. Most of the flagellum is tiny, so that only the filaments are evident at the scale of the bacterial cell (**top right**).

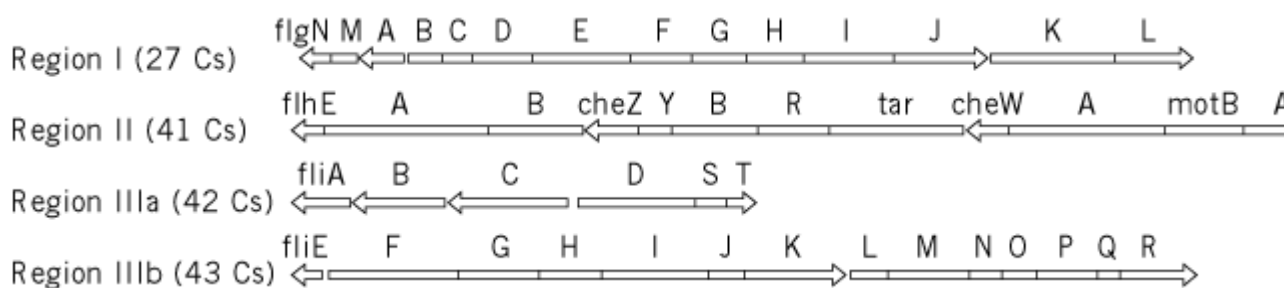


The above description applies when the cell is moving in an isotropic environment. However, if there are gradients of various chemicals (beneficial ones such as the amino acid aspartate or noxious ones such as phenol) and if the cell has **receptors** to detect these chemicals (or light, pH, temperature, etc.), it carries out an ongoing evaluation of whether its environment is improving or deteriorating. If the former is correct, it keeps going; if the latter is correct, it tumbles and then reevaluates whether its new direction is better than its old one. The net effect is a statistical drift to a better life; this behavior is called [chemotaxis](#) (or, more generally, *taxis*). The evaluation is in time, even though it is achieved by travel through space—a stationary bacterium could not make the measurement. Thus bacteria possess a rudimentary memory.

## 2. The Flagellar Gene System

Flagellar, motility, and chemotaxis **genes** (about 60 in all, depending on definition) are highly clustered into just four main regions of contiguous genes (Fig. 2). The organization of the genes shown in Figure 2 is for *Salmonella* (and is almost identical for *E. coli*); in other species, such as the developmental bacterium *Caulobacter crescentus* or the **Gram-positive** bacterium *Bacillus subtilis*, the genetic organization is somewhat different. The four gene clusters contain **operons** ranging in size from a single gene to more than 20 in some species. These operons fall into a hierarchy of expression by **transcription**, with gene products at one level playing a role in transcription of lower levels. In most cases, the controlling elements are of a more or less conventional sort, such as transcriptional activators and dedicated **sigma factors** for transcription initiation. However, there are at least two unconventional mechanisms. The first is a negative regulator, FlgM, which operationally seems like a **repressor** but is in fact an anti-sigma factor rather than a DNA-binding protein. The second is the highly unusual mechanism by which FlgM is inactivated in the last stage of the expression hierarchy; when flagellar assembly has proceeded to a fairly advanced degree, the flagellar structure itself becomes competent to export FlgM from the cell (thus inactivating it), using the same export pathway that is used for the process of flagellar export/assembly itself (see text below).

**Figure 2.** The four clusters of flagellar genes on the chromosome of *Salmonella*, along with their positions (Cs, centisomons), the operons (open arrows) and the genes they contain. Different gene symbols (*flg*, *flh*, and *fli*) are used for the four regions. Region II, which contains flagellar genes, also contains motility (*mot*) genes (see text and Fig. 1) and chemotaxis (*che*) genes.



### 3. Flagellar Structure and Function

Flagellar structure can be divided into several categories (Fig. 1):

1. *External structure.* This is the structure that is most evident by light or **electron microscopy** and includes the helical filament and a hook-shaped structure that acts as a universal joint for the filament.
2. *Basal body.* This structure is embedded in the cell surface (outer membrane, **periplasmic** space, and cytoplasmic membrane) and consists of (a) a rod whose distal end is joined to the hook, (b) an outer-membrane ring and a periplasmic ring, which together act as a bushing for the rod, and (c) a cytoplasmic-membrane ring to which the proximal end of the rod is firmly attached. The external structure and the basal body are passive components; that is, they are not part of the energy transducing machinery.
3. *Motor/switch.* The motor/switch can be subdivided into two categories. The first is a set of so-called Mot complexes that surround the cytoplasmic-membrane ring; these complexes, which are stationary, are involved in the generation of torque but not in the switching of rotation direction. The second is a cylindrical switch complex or cytoplasmic ring, mounted onto the cytoplasmic-membrane ring and projecting into the cytoplasm. The switch complex is involved both in the generation of torque (together with the Mot complexes) and in the switching of rotation direction. It is the rotating element of the motor and causes the cytoplasmic-membrane ring, the rod, the hook, and the filament to rotate as a single unit.

4. *The export apparatus*. This is the least understood structure, but available evidence suggests that it is located at the center of the cytoplasmic-membrane ring and projects out into the cytoplasm. There may also be components that shuttle between the cytoplasm and the components associated with the membrane.

The protein composition and approximate subunit stoichiometry of most of the substructures of the flagellum are known:

1. filament cap (FliD), ~10 subunits;
2. filament (FliC or flagellin), ~20,000 subunits;
3. hook/filament junction (FlgK/FlgL), ~10 subunits each.
4. hook (FlgE), ~130 subunits;
5. rod (FlgB, FlgC, FlgF, FlgG, FliE), ~5 to 25 subunits each;
6. outer-membrane ring (FlgH), ~25 subunits;
7. periplasmic ring (FlgI), ~25 subunits;
8. cytoplasmic-membrane ring (FliF), ~25 subunits;
9. cytoplasmic ring or switch complex (FliG, FliM, FliN), ~30 to 40 subunits each;
10. Mot complexes (MotA, MotB), ~12 subunits each;
11. export apparatus (composition and stoichiometry not well-characterized).

Thus there is a wide range of stoichiometries, with flagellin being by far the most abundant.

#### 4. Flagellar Morphogenesis and the Flagellar Protein Export Process

The assembly of the flagellum proceeds for the most part in a linear fashion—that is, subunit by subunit of the first substructure, then subunit by subunit of the second substructure onto the first, and so on. The first substructures to be assembled are those associated with the cytoplasmic membrane, and the last are the ones furthest away from the cell. Thus a description (oversimplified) of the order would be: cytoplasmic-membrane ring and export apparatus → switch complex and Mot complexes → basal-body rod → periplasmic and outer-membrane rings → hook → filament.

The subunits of the basal-body rod, the hook, the filament, and some other minor but essential substructures are added to the distal end of the growing structure and get there by first passing through the export apparatus (in a process that probably requires energy and certainly requires specificity of recognition) and then traveling down a central core that exists in the structure (Fig. 1).

##### 4.1. The Flagellar Export Pathway is a Type III Secretory Pathway

It has recently become evident that the pathway for export of flagellar proteins is but one member, albeit a rather special one, of a large family of pathways that are used by pathogenic bacteria for the secretion of virulence factors directed against the host; these pathways are called type III secretory pathways (see [Protein Secretion](#)). The similarities within this family extend beyond operational ones, such as lack of [signal peptide](#) cleavage, to biochemical ones. There are at least six flagellar components (FlhA, FlhB, FliI, FliP, FliQ, FliR) that have homologues in the virulence factor secretory apparatus. FliI is an [ATPase](#), which could be driving **active transport** of the protein substrates. Remarkably, it is homologous to a fundamentally important protein for bacteria (and for [mitochondria](#) and [chloroplasts](#)): the catalytic subunit of the proton-translocating  $F_0F_1$ -ATPase (see [ATP Synthase](#)).

##### Suggestions for Further Reading

1. R. M. Macnab (1996) "Flagella and motility". In *Escherichia coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low Jr, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), 2nd ed., ASM

- Press, Washington, DC, pp. 123–145. [A comprehensive review of flagella and motility in *E. coli* and *Salmonella*, with 214 references.]
2. G. W. Ordal, L. Márquez-Lagaña, and M. J. Chamberlin (1993) "Motility and chemotaxis". In *Bacillus subtilis and Other Gram-Positive Bacteria: Biochemistry, Physiology and Molecular Genetics* (A. L. Sonenshein, J. A. Hoch, and R. Losick, eds.) ASM Press, Washington, DC, pp. 765–784. [Comprehensive coverage of motility and chemotaxis in Gram-positive bacteria, with 153 references.]
  3. J. Wu and A. Newton (1997) Regulation of the *Caulobacter* flagellar gene hierarchy; not just for motility. *Mol. Microbiol.* **24**, 233–239. [A microreview of the inter-related regulation of the cell cycle and regulation of expression of the flagellar gene system in *Caulobacter crescentus*, with 47 references.]
  4. R. Losick and L. Shapiro (1993) Checkpoints that couple gene expression to morphogenesis. *Science* **262**, 1227–1228. [A concise summary of the remarkable mechanism of genetic regulation by export of the regulatory protein FlgM.]
  5. F. Van Gijsegem, C. Gough, C. Zischek, E. Niqueux, M. Arlat, S. Genin, P. Barberis, S. German, P. Castello, and C. Boucher (1995) The *hrp* gene locus of *Pseudomonas solanacearum*, which controls the production of a type III secretion system, encodes eight proteins related to components of the bacterial flagellar biogenesis complex. *Mol. Microbiol.* **15**, 1095–1114. [A primary publication describing a typical type III virulence factor secretory system and emphasizing its relationship to flagellar protein export.]

## Flow Cytometry

Flow cytometric analysis of **eukaryotic** cells, also known as fluorescence-activated cell sorting (FACS), was first developed in the late 1960s. FACS has been a tool primarily of immunologists because of its power in dissecting the complex subsets of cells derived from the blood. It has found clinical application in the classification of leukemias and lymphomas, as well as a tool for prognosis in [HIV](#) disease. Because of the relative expense of the equipment, it has been slow to progress into many other areas of biological research. However, the past decade has seen significant expansion in the techniques that are designed for FACS analysis. In particular, several developments came about that make FACS an extremely powerful tool for molecular biology analysis. These include (a) the development of FACS-based assays for [reporter genes](#) (initially, **b-galactosidase**; now including several other hydrolases) and (b) the introduction of the inherently fluorescent proteins, such as green fluorescent protein (GFP; see [Luciferases And Luciferins](#)). As the technology becomes cheaper and more widespread, these applications (and others) will become a common tool in the arsenal of molecular biology.

### 1. Molecular Biology by FACS: A Basic Overview

Flow cytometry is a technology that characterizes populations of cells, on a cell-by-cell basis, in terms of their fluorescence properties. Several different fluorescent signals can be monitored simultaneously and independently. On the basis of the amount of each fluorescence exhibited by each cell, the instrument can decide whether that cell should be sorted (separated from all others) or simply discarded. The sorted cells can be recovered in a viable state for further biochemical processing or culture. Most flow cytometric sorters can process as many as 5000 cells per second, with excellent purity and specific recovery of cells occurring as infrequently as one in a million.

Fluorescence is the fundamental parameter evaluated by a flow cytometer, so the most common molecular biology technique performed by FACS is quantification of reporter gene activity. Thus, either intrinsic or enzymatically generated fluorescence serves as the surrogate assay for the number of protein molecules synthesized from a particular gene construct introduced into the cell. The intensity of the fluorescence is, in general, linearly related to the expression from the gene construct, so the FACS can quantify directly, on a cell-by-cell basis, the **translation/transcription** activity of the reporter gene.

It is important to recognize that the measurement is made on each cell independently, so the analysis of a cell population can be displayed as a distribution of activities. Thus, heterogeneous gene expression within a population can be resolved; for example, the presence of a few highly expressing cells along with a majority of nonexpressing cells is easily detected. Furthermore, the sorting capabilities of the instrument can then be used to isolate either the positive or negative cell fraction.

Cells are passed through the interrogating laser beam one-by-one in a stream, so it is necessary that the sample be prepared as a suspension of individual cells. It is possible to monitor adherent cell populations after removing the cells from the substrate—for example, by treatment with a [proteinase](#) such as **trypsin**. It is not possible, however, to use the FACS on solid tissue samples without dispersing the cells. Virtually any cell type can be used; successful assays have been performed on **bacteria**, **yeast**, and **Drosophila** cells, as well as on a wide variety of mammalian cells.

## 2. Reporter Genes assays for FACS

The most common application of FACS in molecular biology is the quantitative measurement of reporter gene expression. A variety of reporter genes are amenable to FACS analysis, including cell surface markers, hydrolases, and GFP. Surface markers are genes coding for membrane-bound proteins that are not normally expressed by the cells of interest. These can be detected using fluorescently conjugated [monoclonal antibodies](#) and standard immunophenotyping methods. These assays are not necessarily quantitative; that is, the amount of surface expression may not always correlate with transcriptional activity. The introduction of foreign proteins on the surface of cells also prevents the introduction of such cells into a host, because of immunological rejection.

The other types of assays rely on intracellular expression of proteins that can be detected enzymatically or are inherently fluorescent, such as GFP. In 1988, Nolan et al. developed the first FACS-based assay for an enzymatic reporter gene (1-3). This assay uses the fluorogenic substrate fluorescein digalactoside (FDG). FDG is nonfluorescent until hydrolyzed; it can be easily introduced into the cytoplasm of cells. Cells that have high b-galactosidase activity (eg, that are transfected with the bacterial *lacZ* gene controlled by a promoter active in the host cells) will hydrolyze the FDG to yield fluorescent **fluorescein**. The amount of fluorescence generated is quantitatively related to the amount of cellular b-galactosidase. This assay has extreme sensitivity (it can detect cells with as few as 5 molecules of enzyme) and a large dynamic range (at least 5 orders of magnitude). More recently, assays for b-glucuronidase and b-glucosidase have also been introduced that are similar to the one for b-galactosidase (4).

The last few years has seen an explosion in the use and development of GFP as a reporter gene. A variety of mutations have been introduced into the gene to produce proteins that are more highly fluorescent and with altered fluorescence spectra. The primary advantage of the GFP assay is ease of use: the expressing cells become spontaneously fluorescent. In addition, the genes encoding the varieties of GFP are relatively small (~1.5 kb) compared to *lacZ* (~3.5 kb), giving greater flexibility in designing vectors.

## 3. Why Is FACS Such a Powerful Tool?

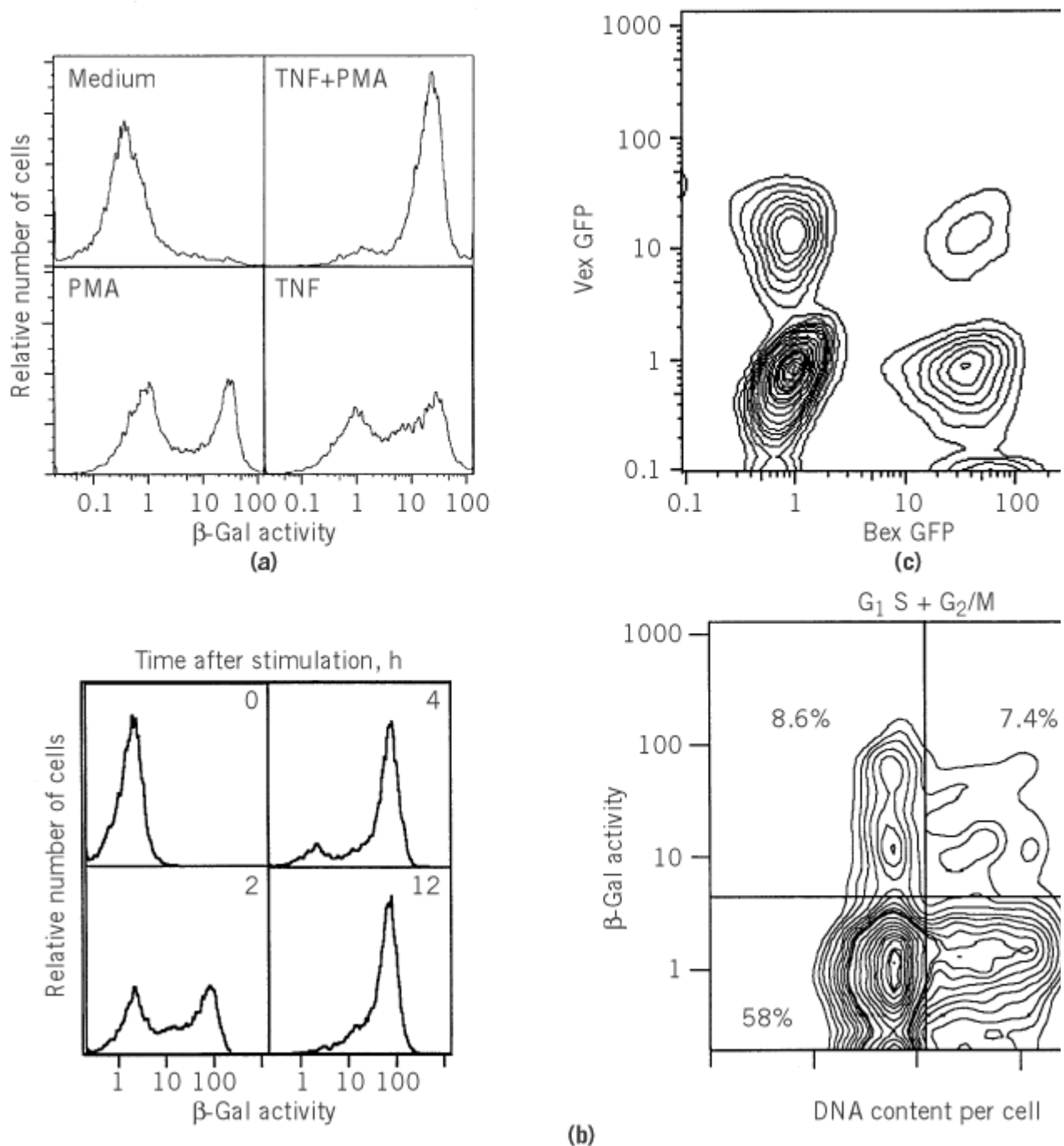
There are three principal reasons:

1. FACS analyses collect data on each individual cell, presenting the data as a distribution of activity.
2. FACS analysis is multiparametric (the ability to measure several parameters simultaneously on each cell).
3. The FACS can viably sort cells based on these measurements.

#### 4. Analysis of a Distribution

Biological systems are rarely, if ever, homogeneous. Indeed, even the responsiveness of a **clone** is often complex. This has become even more evident since the advent of FACS analysis. Figure 1 shows two examples of such heterogeneity. These examples demonstrate how a cloned population of cells (which has a uniform, single integration of a reporter gene construct) responds in a heterogeneous way to stimulations. It is immediately apparent that such complexity is completely missed by using standard, bulk assays for reporter gene activity. Bulk assays measure the average amount of activity in a population, with no information about the distribution within that population. The figure demonstrates that the signal transduction mediated by effector molecules does not affect all cells identically; some respond rapidly, whereas others remain resting for at least a period of time.

**Figure 1.** Three examples of flow cytometric analysis of gene expression. **(a)** A reporter gene construct consisting of a transcriptional fusion of the *lacZ* gene to the HIV [long-terminal repeat](#) (LTR), the promoter/enhancer region for HIV type 1, was stably transfected into 293 S cells. Basal expression is shown in the *top left* panel. After stimulation with [tumor necrosis factor](#) (TNF) or **phorbol ester** (phorbol myristate acetate, PMA) for 6 h, most cells have significantly elevated the expression of the reporter gene. Thus, the response to these factors is not homogeneous, even within this cloned cell line. The synergistic effect of TNF+PMA results in activation of the construct in virtually all cells. **(b)** In this example, the *lacZ* gene was fused to a promoter containing three adjacent binding sites for NF-AT, the nuclear factor of activated T cells, and transfected into the Jurkat T cell line. The b-galactosidase activity corresponds to the degree of activation of the nuclear factor. After stimulation of the cells with PMA+ionomycin, there is a graded response over time. At the early time points, some cells responded and others haven't; note the dichotomy in the response. Those cells that have responded have fully turned on the promoter (ie, the position of the positive peak at 2 h is essentially the same as at 12 h). What increases over time is the fraction of responding cells. *(Right)* Cells stimulated with PMA+ionomycin for 2 h were stained for b-galactosidase activity and DNA content. The multiparametric capability of the FACS allows the simultaneous measurement of these values. This analysis shows that cells progressing through the [cell cycle](#) are more likely to respond early after stimulation (nearly half of the responding cells are in S+G<sub>2</sub>M; less than one-third of nonresponding cells are in cycle). **(c)** The generation of spectral mutants of the inbred fluorescent GFP allows for simultaneous measurement of two different genetic elements in the same cell. "Vex" is violet-excited GFP, and "Bex" is blue-excited GFP. In this example, cells were infected with [retroviruses](#) expressing one or the other GFP. Uninfected cells, singly infected cells, and cells infected with both viruses can be easily distinguished by measuring the fluorescence signals by FACS.



## 5. The Power of Multiparametric Analysis

A typical FACS experiment can measure at least three different fluorescent species, in addition to measurements related to cell size and shape. Thus, one can combine a typical reporter gene expression assay (using one fluorescent channel) together with two other assays. These might include assays to measure DNA or RNA content, cell surface **phenotype**, viability, and so on. In addition, one could choose a second, independent reporter gene assay using a second color, in order to measure simultaneously two independent genetic elements (ie, [enhancers](#), promoters, [repressors](#)). Figure 1 c shows such an analysis, using two different versions of the inherently fluorescent protein GFP (5). Another analysis that can benefit from such simultaneous measurement is the dissection of



enhancer elements. By transfecting the same cells with different reporter genes fused to different enhancer elements, one can now measure the relative activity of each enhancer element in the same cell in response to stimulations.

## 6. Sorting (Selection)

Perhaps one of the most underutilized capabilities of the FACS is the ability to sort cells viably based on any of the measurement parameters. Effectively, this means that any reporter gene used in the context of the FACS becomes (in addition) a selectable marker. Selection by FACS has several advantages over standard drug-based cytotoxic selectable markers. Generally, selectable marker systems rely on a combination of toxic compounds and genes that encode resistance or sensitivity to the toxic compound. These systems work by killing or halting the growth of all cells that do not (or do) produce some threshold amount of the protein encoded by the selectable marker gene.

Survival-dependent selection systems have several disadvantages. For example:

1. The selection system is inherently toxic, often even to cells expressing the resistance gene.
2. The selection generally takes many cell divisions, or a period of weeks for mammalian cells.
3. The cells below (or above) the threshold amount of the selecting protein are dead and unavailable for further study.
4. It is difficult or impossible to select for different expression levels of the marker.
5. The cells must generally be able to reproduce in culture.

On the other hand, FACS analysis of reporter genes allows the sorting capabilities of the FACS to isolate cells based on virtually any level of expression of the nontoxic reporter genes. The FACS can analyze and sort from as many as 5000 cells per second, so 1 cell in 10 million could be selected in an hour of sorting, making FACS as efficient as drug-based selectable markers.

The ability to select cells based on a particular level of expression is a big advantage over standard selection methods, for several reasons. First of all, cells expressing very low levels of the  $\beta$ -galactosidase enzyme (eg, 10 molecules) can be detected and sorted; such low levels might not be sufficient to confer drug resistance. Second, cells expressing virtually any level of enzyme can be selected (for instance, only those expressing more than 100,000 molecules per cell). Third, cells that rapidly change expression of the reporter gene in response to a stimulus can be selectively sorted. These capabilities give FACS assays the ability to select cells that would be impossible to select using standard drug selection methods.

Finally, the ability to sort cells based on gene expression levels allows the concomitant determination of other parameters that can only be determined by, for example, biochemical assays on lysates. For instance, by sorting cells with different expression levels and then doing [messenger RNA](#) quantification on the sorted cells, a comparison of message levels within defined subsets of a population can be done (eg, see Ref. [2](#)).

## 7. Summary

Besides the analysis of reporter genes, other technologies are becoming available for FACS analysis. For instance, fluorescent *in situ* **hybridization** (FISH) for both RNA and DNA, even after *in situ* polymerase chain reaction (PCR), has been successfully performed by FACS. The FACS-based enzyme assays have also been extended to nonmammalian systems. Reporter gene analysis has been performed in **plant** cells ([6](#)), **Drosophila** cells ([7](#)), **yeast** ([8](#)), and **bacteria** ([8](#), [9](#)), complete with viable sorting of the cells for further analysis.

Flow cytometric analysis is an aspect of the majority of immunologically related research, as can be evidenced by perusing any immunology journal. Since the relatively recent introduction of this technology to molecular biology, it has been growing quickly as the applications become more

advanced and as the technology becomes more accessible. Flow cytometry has already proven to be a unique and powerful adjunct for molecular analyses; in the years to come, it will undoubtedly become a common tool.

### Bibliography

1. G. P. Nolan, S. Fiering, J. F. Nicolas, and L. A. Herzenberg (1988) Proc. Natl. Acad. Sci. USA **85**, 2603–2607.
2. S. Fiering, M. Roederer, G. P. Nolan, D. R. Micklem, D. R. Parks, and L. A. Herzenberg (1991) Cytometry **12**, 291–301.
3. M. Roederer, S. N. Fiering, and L. A. Herzenberg (1991) Methods **2**, 248–260.
4. M. Lorincz, M. Roederer, L. A. Herzenberg, and G. P. Nolan (1996) Cytometry **24**, 321–329.
5. M. T. Anderson, I. M. Tjioe, M. C. Lorincz, D. R. Parks, L. A. Herzenberg, G. P. Nolan, and L. A. Herzenberg (1996) Proc. Natl. Acad. Sci. USA **93**, 8508–8511.
6. D. W. Galbraith, G. M. Lambert, R. J. Grebenok, and J. Sheen (1995) Methods Cell Biol. **50**, 3–14.
7. M. A. Krasnow, S. Cumberledge, G. Manning, L. A. Herzenberg, and G. P. Nolan (1991) Science **251**, 81–85.
8. R. Nir, Y. Yisraeli, R. Lamed, and E. Sahar (1990) Appl. Environ. Microbiol. **56**, 3861–3866.
9. F. Russo-Marie, M. Roederer, B. Sager, L. A. Herzenberg, and D. Kaiser (1993) Proc. Natl. Acad. Sci. USA **90**, 8194–8198.

### Fluctuation Test

The fluctuation test was invented by Luria and Delbrück in 1943 to determine whether bacterial **mutants** arise before any **selection** for their presence or only after the **bacteria** are subjected to selection (1). It was enormously influential because by showing that **mutations** arise randomly before selection, it provided strong evidence that bacteria and **viruses** are normal organisms possessing heritable genetic determinants. Thus it opened the way for the genetic analysis of bacteria and their viruses and led eventually to modern molecular biology. Luria and Delbrück measured the number of mutants resistant to **bacteriophage** T1 in a large number of replicate cultures of *Escherichia coli*. They reasoned as follows. If mutants occur after the culture is exposed to the phage (ie, in response to selection), then little variation should occur among cultures in the number of mutants. In fact, the distribution would be Poisson with the variance equal to the mean. However, if mutants arise at random during nonselective growth of cells, the probability of a mutation would be constant per generation per cell. But the consequence of that mutation would depend on when during the growth of the population the mutation occurred. In a binary dividing population, a cell that sustains a mutation gives rise to a **clone** of identical descendants, each of which is a mutant. Thus a mutation during early generations gives rise to a large clone of mutant cells, whereas a late mutation gives rise to few mutant cells. Among a large set of identical cultures of dividing cells, the few cultures in which the mutation happened in the early generations (the jackpots) have a large number of mutants, whereas the majority of the cultures have none or a few mutants. The predicted distribution, called the Luria–Delbrück distribution, has a variance much larger than the Poisson distribution. This is what Luria and Delbrück observed. Their test is known as “the fluctuation test” because it measures the degree of fluctuation in the number of mutants found in replicate cultures.

Interest in the fluctuation test has been revived recently by the controversy surrounding the

phenomenon known as “directed” or “adaptive” mutation (2). By obtaining the highly variant distribution, Luria and Delbrück proved that mutations occur before selection. Because the selection they used was lethal, however, they did not prove the converse that mutations could not also arise after selection. Indeed, when the selection is not lethal, it has been found that mutations do arise in response to selection. Part, but not all, of the evidence for postselection mutation is that in a fluctuation test the distribution of the number of mutants per culture deviates from the Luria–Delbrück toward the Poisson distribution. This has inspired new models describing the distribution of mutant numbers when one or another of the underlying assumptions of the Luria–Delbrück analysis are not met (2, 3).

The fluctuation test is also useful in determining mutation rates during nonselective growth. First applied by Luria and Delbrück themselves, the computation of mutation rates has been refined since then. But all commonly applied techniques still involve assumptions and compromises. Nonetheless, the power of the fluctuation test is that a true mutation rate, ie, the number of mutational events per cell division, is obtained, not a mutant frequency that can be dramatically skewed by jackpots.

The fluctuation test consists of a large number of identical cultures inoculated with a number of cells few enough to ensure that no preexisting mutants are present. These cultures are allowed to grow, achieving at least a thousandfold increase in cell number. Then mutants are selected, usually by plating each entire culture onto a selective agar medium. After a suitable time to allow the mutant clones to grow, the resulting colonies are counted. The parameter used for calculating the mutation rate is  $m$ , the mean number of mutations (not mutants) per culture. This is, perhaps, a difficult concept to understand because it is a function of the mutation rate and also of the number of cells at risk for mutation (and thus one requirement for a proper fluctuation test is that every culture has the same number of cells). To date, however, the only solved methods to determine the mutation rate from fluctuation test results are based on the distribution of  $m$ . Then the value of  $m$  is usually divided by  $2N$ , twice the final number of cells in the culture, to obtain the mutation rate as mutations per cell per generation (because a culture of  $N$  cells contains a total of  $2N$  cells during its entire history). Some researchers attempt to correct for asynchronous divisions by dividing  $m$  by  $N/\ln 2$  instead. These calculations assume that the initial number of cells is trivial compared to the final number, that the proportion of mutants is always small, that mutants grow at the same rate as wild-type, and that reverse mutations are negligible.

The data obtained from a fluctuation test are the number of mutants  $x_i$  in each culture. From the number of mutants, there are several commonly used methods to obtain estimates of  $m$  (often called estimators), of which five are described here.

1. The  $P_0$  method (1). If the probability of mutation is constant per cell per generation, then the number of mutations per culture has a Poisson distribution (although the number of mutants has a Luria–Delbrück distribution). The probability of no mutations (and thus no mutants)  $P_0$ , is  $e^{-m}$ , the first term of the Poisson distribution. Thus, the  $P_0$  estimate of  $m$  is obtained from the proportion of cultures that have no mutants as  $m = -\ln P_0$ .
2. The method of the mean (1). Although the extreme variance of the Luria–Delbrück distribution is caused by mutations early in the growth of cultures, these are, in fact, rare events. After the critical point where the total population size (all of the cells in all of the cultures) is large enough so that the probability of a mutation approaches unity, every succeeding generation makes an equal contribution to the number of mutants (from new mutations plus the growth of preexisting mutants). Assuming that no mutations occur before the critical point, the value of  $m$  is derived as  $x_{\text{mean}} = m \ln(mC)$  where  $x_{\text{mean}}$  is the mean number of mutants per culture and  $C$  is the number of cultures.
3. The graphical method (4, 5). Because each generation produces an equal number of mutants after the critical point, the accumulated distribution of the number of mutants per culture,  $Y_x$ , is

estimated as  $m/x$ , where  $Y_x$  is the number of cultures with  $x$  or more mutants. Thus a plot of  $\log Y_x$  vs.  $\log x$  gives a straight line with a slope of  $-1$ . The intercept at  $\log x = 0$  gives  $\log m$ . The advantage of this method is that the straight line can be fit by eye to the intermediate values, which are not caused by jackpots.

4. The method of the median. Although the mean number of mutants fluctuates widely depending on the number of jackpots present, the median is not so influenced and is thus a more stable statistic. Lea and Coulson (6) found empirically that  $(x_{\text{median}}/m) - \ln(m) = 1.24$ . This equation is easily solved for  $m$  by iteration.
5. The maximum likelihood method. This method relies on solving, at least approximately, the Luria–Delbrück distribution and then generating the most likely value of  $m$  for the experimentally determined  $x_i$  values. Although this is the most accurate method of estimating  $m$ , it has been shunned in the past because of the difficulties in calculation. However, a reasonably adequate way of obtaining  $m$  is to solve Eq. (50) of Lea and Coulson (6) by iteration. In addition, Koch (7) provided exact values of  $m$  obtainable from the median and upper and lower quartiles of the distribution, which was extended to the upper quartiles by Cairns et al. (2). Simpler algorithms and sophisticated computer programs, however, make the maximum-likelihood method now computationally feasible (8).

Of these methods, the mean (2) is the least accurate because it is skewed by jackpots. The utility of the other methods depends on the value of  $m$ . If  $m \leq 1$ , the  $P_0$  method (1) is reliable, so long as  $P_0$  is determined with sufficient accuracy. When  $m = 1$  to 4, the median method (4) is convenient and adequate. At values of  $m$  larger than 4, the median can still be used, but the maximum likelihood method (5) is preferable. Another advantage of the maximum likelihood method is that the standard deviation of  $m$ , and thus of the mutation rate, can be approximated (9). Further discussions and refinements can be found in Refs. 9-11.

The phenomenon of [phenotypic lag](#), in which the expression of a new mutant phenotype is delayed, complicates determining mutation rates from fluctuation tests (3, 7, 12).

### Bibliography

1. S. E. Luria and M. Delbrück (1943) *Genetics* **28**, 491–511.
2. J. Cairns, J. Overbaugh, and S. Miller (1988) *Nature (London)* **335**, 142–145.
3. F. M. Stewart, D. M. Gordon, and B. R. Levin (1990) *Genetics* **124**, 175–185.
4. S. E. Luria (1951) *Cold Spring Harbor Symp. Quant. Biol.* **16**, 463–470.
5. J. Cairns (1980) *Nature (London)* **286**, 176–178.
6. D. E. Lea and C. A. Coulson (1949) *J. Genet.* **49**, 264–285.
7. A. L. Koch (1982) *Mutat. Res.* **95**, 129–143.
8. S. Sarkar, W. T. Ma, and G. v. H. Sandri (1992) *Genetica* **85**, 173–179.
9. F. M. Stewart (1994) *Genetics* **137**, 1139–1146.
10. M. E. Jones, S. M. Thomas, and A. Rogers (1994) *Genetics* **136**, 1209–1216.
11. G. Asteris and S. Sarkar (1996) *Genetics* **142**, 313–326.
12. P. Armitage (1952) *J. R. Stat. Soc. B* **14**, 1–40.

### Fluid Mosaic Model

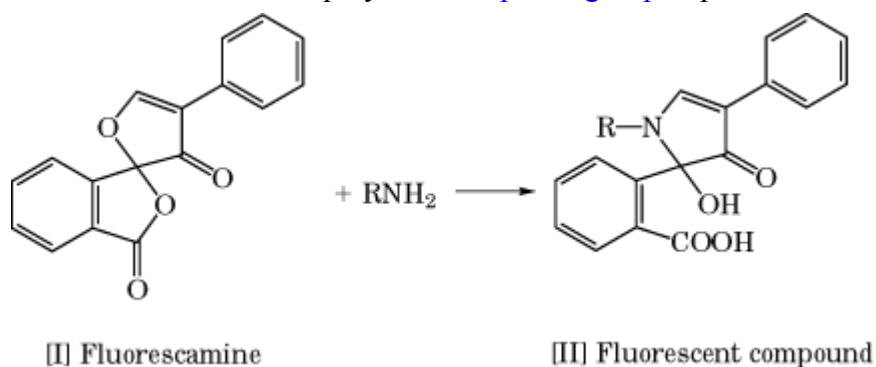
The liquid-like properties of [lipid](#) liquid crystals are assimilated into this model for the structure of cell [membranes](#) and for the organization of its [protein](#) and lipid components. The concepts that underlie this model assume that integral membrane proteins are “a heterogeneous set of globular molecules, each arranged in an [amphipathic](#) structure, i.e. with the ionic and highly [polar](#) groups protruding from the membrane into the aqueous phase, and the [nonpolar](#) groups largely buried in the **hydrophobic** interior of the membrane” (1). In this model the membrane is perceived as a mosaic of a continuous phospholipid bilayer that is periodically interrupted by proteins randomly distributed throughout the bilayer. The phospholipid bilayer is considered a viscous solvent that allows the dissolved proteins to diffuse laterally within the plane of the membrane. (See also [Membranes](#), and [Amphipathic](#).)

### Bibliography

1. S. J. Singer and G. L. Nicholson. (1972) *Science* **172**, 720–731.

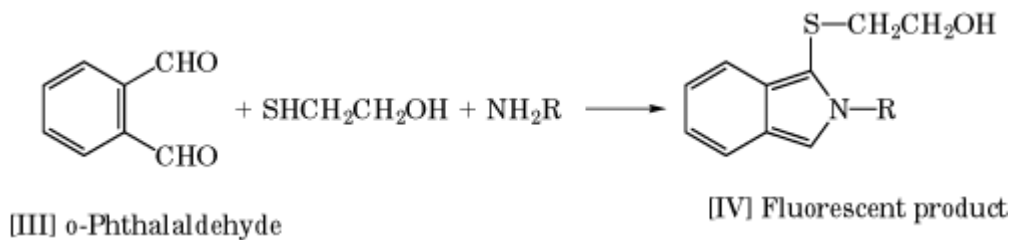
### Fluorescamine

Fluorescamine [I] reacts readily with primary [amino groups](#) to form highly **fluorescent** compounds [II] [Scheme 1](#), even though fluorescamine itself is nonfluorescent. The fluorescent products have an excitation maximum at 390 nm and an emission maximum at 475 nm. These properties make fluorescamine ideal for detecting amino groups, especially in [proteins](#), [peptides](#), and [amino acids](#). It is 10 to 100 times more sensitive in detecting primary amino groups than the [ninhydrin](#) reaction. Fluorescamine is also employed as a [reporter group](#) in protein chemistry and immunology.



The formal nomenclature for fluorescamine is 4-phenylspiro[furan-2(3H),1''-(3'H)-isobenzofuran]-3,3'-dione or 4-phenylspiro[furan-2(3H),1'-phthalan]-3,3'-dione. It has a molecular weight of 278.27 and a melting point at 154 to 155°C.

Because fluorescamine has low solubility and stability in water, *o*-phthalaldehyde [III], which reacts similarly with primary amino groups and gives highly fluorescent products [IV] [Scheme 2](#), is sometimes substituted for fluorescamine (1, 2):



To quantify proteins with fluorescamine (3), a sample solution (10 to 250  $\mu$ L, containing 500 ng to 50  $\mu$ g of protein) is made up to 1.5 mL with 0.05 M [phosphate buffer](#) (pH 8). Then 0.5 mL of fluorescamine solution, made by dissolving 30 mg of fluorescamine in 100 mL of dioxane, is added to it with vigorous stirring. The emission at 475 nm, upon excitation at 390 nm, is quickly determined. This method is based on the reaction with primary amino groups, so the extent of emission varies depending on the amino acid composition of the proteins, any contaminating compounds having primary amines interfere with the assay.

Fluorescamine reacts with amino acids promptly at pH 9 and yields fluorescent products. The reaction product with ammonia has low fluorescent intensity (less than 1% of the normal products with primary amines). Thus, this reaction is employed for automatic [amino acid analysis](#) with a fluorescence detector.

Fluorescamine is used to detect **electroblotted** proteins (4) by visualizing transferred proteins on the transfer membrane with UV light. This method is very sensitive, it is also applied to detect peptides on filter paper or in the eluate from column [chromatography](#). Fluorescamine is also used to test the completeness of other reactions of primary amines (5).

#### Bibliography

1. G.L. Peterson (1983) *Methods Enzymol.* **91**, 110–119.
3. P. Bohlen et al. (1973) *Arch. Biochem. Biophys.* **155**, 213–220.
4. S. Best and D.W. Speicher (1995) In *Current Protocols in Protein Science* (J.E. Coligan et al., eds.), Wiley, pp. 10.8.5–7.
5. A.M. Felix and M.H. Jimenez (1973) *Anal. Biochem.* **52**, 377–381.

### Fluorescence Energy Transfer

In principle, the fluorescent light emitted by a fluorophore (the energy donor) can be reabsorbed by another chromophore (the energy acceptor) in the same molecule or somewhere else in the solution, provided that the absorption spectrum of the acceptor overlaps with the fluorescence emission spectrum of the donor (see [Absorption Spectroscopy](#) and [Fluorescence Spectroscopy](#)). Such an energy transfer by emission and reabsorption of light is extremely inefficient. A much more efficient energy transfer between donor and acceptor can occur in a nonradiative process, which involves the coupling of their transition dipoles. This process is called resonance energy transfer. It depends on:

1. the extent of overlap between the fluorescence band of the donor and the absorption band of the acceptor,
2. the distance between the two chromophores, and
3. the relative orientation of their transition dipoles.

Resonance energy transfer is strongly sensitive to distance and varies with the inverse sixth power of the distance  $R$  between donor and acceptor. This is very useful to measure distances in macromolecules and changes in distances as a result of conformational changes in proteins and nucleic acids. An important number in energy transfer is the characteristic transfer distance  $R_0$ .  $R_0$  is the distance between donor and acceptor at which fluorescence emission from the donor and energy transfer to the acceptor are equally probable. Förster developed a theory to correlate the efficiency of energy transfer with  $R$  and  $R_0$ . The characteristic distance  $R_0$  can be calculated (primarily from the spectral overlap between donor and acceptor), and consequently the distance  $R$  between donor and acceptor can be determined from the measured efficiency of fluorescence energy transfer. Values for  $R_0$  often lie in the range between 1 and 5 nm. Thus, distances up to about 8 nm can be measured by energy transfer.

In proteins, energy is transferred primarily from **tyrosine** to **tryptophan** residues. Tryptophan residues can also transfer their energy to natural acceptors, such as **NADH**, riboflavin, or heme. For measuring distances in proteins or nucleic acids, often both donor and acceptor are introduced by site-directed labeling (see [Reporter Groups](#)). A naphthyl group is often used as the donor and a dansyl group as the acceptor in such experiments. A detailed list of useful donor/acceptor pairs is found in (1).

#### Bibliography

1. M. R. Eftink (1991) "Fluorescence techniques for studying protein structure". In *Methods of Biochemical Analysis* (C. H. Suelter, ed.), Vol. **35**, Wiley, New York, pp. 127–205.

#### Suggestions for Further Reading

2. C. S. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Vol. **II.**, Chap. "8", W. H. Freeman and Co., San Francisco, CA.
3. G. R. Penzer (1980) "Molecular emission spectroscopy". In *An Introduction to Spectroscopy for Biochemists* (S. B. Brown, ed.), Academic Press, London, pp. 70–114.

## Fluorescence Microscopy

Fluorescence microscopy uses the contrast created by fluorescent molecules, fluorophores, that emit light of a specific wavelength when excited by incident light of a different (shorter) wavelength. By using filters, one can select an appropriate wavelength for excitation of the fluorophore and a second filter can eliminate light of wavelengths other than that emitted by the fluorophore. The second filter is very important as it removes the much more intense excitation light, so that the emitted signal is detected in the absence of the excitation light and other stray light, giving fluorescence microscopy an excellent signal-to-noise ratio. Fluorescence microscopy universally uses epi-illumination in which light from above is reflected off the surface of the specimen by light emerging from the objective lens. Epi-illumination has the advantages that (1) unabsorbed light is directed away from the observer and (2) the alignment is simplified since the objective lens acts as its own condenser. It is important to understand the capabilities and limitations of the fluorophores. Because a specimen can be stained with a finite quantity of fluorescent dye, which decays or photobleaches when illuminated, the microscope must be configured to efficiently excite and detect the limited number of emitted photons before complete bleaching has occurred. Although a large number of fluorescent

dyes have been discovered, relatively few are suitable for fluorescence microscopy (1).

Immunofluorescence light microscopy is the most common technique of fluorescence microscopy and involves the use of antibodies labeled with fluorophores to visualize distributions of proteins and nucleic acids within a specimen. This technique can be performed by either *direct* or *indirect* antibody labeling. Direct immunofluorescence requires the conjugation of a primary antibody with a fluorophore, commonly rhodamine and fluorescein. Indirect immunofluorescence, on the other hand, uses a fluorophore conjugated secondary antibody which was raised against the [immunoglobulin](#) type of the primary antibody. For example, if the primary antibody is raised in a rabbit, the secondary antibody might be a goat anti-rabbit antibody. Indirect immunofluorescence is the more common of the two techniques because making a fluorophore conjugate of each primary antibody is time-consuming and can decrease the affinity or specificity for the [antigen](#). Furthermore, the indirect technique may be more sensitive because more than one molecule of the secondary antibody may bind to a molecule of the primary antibody (2), creating an amplification of the signal. Indirect immunofluorescence has motivated the development of a large number of commercially available secondary antibodies conjugated to various fluorophores, so in order to do immunofluorescence microscopy, one need only raise and characterize the primary antibody and then purchase the secondary antibody. A key consideration in immunofluorescence microscopy is confirmation of the specificity of both fluorophore and antibodies. In order for immunofluorescence results to be validly interpreted, one must demonstrate that (i) the secondary antibody recognizes only the primary antibody, and (ii) the primary antibody recognizes only the antigen. This means that the secondary antibodies must not cross-react with intrinsic proteins and do not recognize each other. It is also necessary to establish whether autofluorescence from compounds intrinsic to the tissue can mimic the appearance of fluorophores. This test is easily done by examining an unstained specimen using the same excitation and emission filters intended for use with the fluorophore.

Fluorescence microscopy can also be used to monitor intracellular compartments through the use of fluorescent indicator dyes, which report on their local environments. Some of the fluorescent dyes useful for microscopy are membrane potential indicators such as JC-1 (3). Using this probe, it was shown that heterogeneity in mitochondrial [membrane potential](#) exists between different [mitochondria](#) in the same cell (4). Other fluorescent dyes are ion indicators, which can report the pH or concentrations of calcium, chloride, magnesium, potassium, and sodium ions either qualitatively or quantitatively (5).

With the appropriate choice of fluorescent labels and proper equipment, one can do multilabeling [sometimes called multicolor (6)] to localize two or more molecular species using fluorophores with different excitation and/or emission wavelengths. By using multicolor immunofluorescence microscopy, the neurotransmitters, serotonin and substance P, were observed in separate populations of spinal cord fibers (6) and the motor neurons and the substance P fibers, which wrap around them, could be separately visualized. Another example is the use of DiOC6 to label mitochondria and [ethidium bromide](#) to label [mitochondrial DNA](#) in *Euglena gracilis* cells (7).

## Bibliography

1. R. Y. Tsien and A. Waggoner (1990) In *Handbook of Biological Confocal Microscopy*, 2nd ed. (J. Pawley, ed.), Plenum Press, New York, pp. 169–178.
2. L. A. Sternberger (1986) *Immunocytochemistry*, 3rd ed., Wiley, New York.
3. M. Reers, T. W. Smith, and L. B. Chen (1991) *Biochem.* **30**, 4480–4486.
4. S. T. Smiley et al. (1991) *Proc. Natl. Acad. Sci.* **88**, 3671–3675.
5. B. Herman and J. J. Lemasters (1993) *Optical Microscopy*, Academic Press, San Diego.
6. T. C. Brelje, M. W. Wessendorf, and R. L. Sorenson (1993) "Multicolor laser scanning confocal immunofluorescence microscopy: Practical application and limitations", In *Methods in Cell Biology*, Vol. **38**: Cell Biological Applications of Confocal Microscopy (B. Matsumoto, ed.), Academic Press, San Diego, pp. 98–182.



7. D. L. Spector, R. D. Goldman, and L. A. Leihwand (1998) *Cells: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Plainview, New York.

## Fluorescence Quenching

After excitation by absorption of light (see **Absorbance spectroscopy**), a chromophore can lose the absorbed energy, and thus return from the excited state to the ground state, either by emitting a photon (observed as fluorescence) or alternatively by nonradiative processes, such as exchange of heat with the solvent. For a particular chromophore, the quantum yield of the observed fluorescence depends on the relative rates of fluorescence emission and the competing nonradiative processes. Under conditions where such nonradiative processes are very slow, one photon is emitted for every photon absorbed, and the quantum yield (the number of photons emitted per number of photons absorbed) approaches one. In practice, the observed quantum yields are often much smaller than one because nonradiative processes compete efficiently with the emission of light. Such a decrease in quantum yield is called fluorescence quenching. The collision of a chromophore with certain molecules in the solution is an effective means to quench fluorescence; such a process is called collisional quenching.

Radiationless transitions can occur by two major routes: internal conversion and intersystem crossing. In internal conversion, the first excited electronic singlet state is converted to the ground singlet state by exchange of vibrational energy (ie, heat) with the solvent. In intersystem crossing, the first excited singlet state is converted to the first excited triplet state. This triplet state is more stable than the corresponding singlet state and is very long-lived, because transitions between triplet and singlet states are nominally forbidden. They are accompanied either by heat exchange or phosphorescence (ie, emission of light with low energy).

Intersystem crossing is triggered by magnetic fluctuations in the system. It is thus greatly facilitated by collisions of the excited molecules with molecules or ions that contain unpaired electrons or loosely held electron clouds (such as acrylamide or iodide anions). The efficiency of collisional quenching depends on the frequency of the collisions, ie, on the concentration of the added quencher and exposure of the fluorophore to the solvent. The susceptibility to collisional quenching can be used to measure the degree of exposure to solvent of a fluorescing group (such as a [tryptophan](#) residue in a protein). Experimentally, the ratio of the fluorescence in the absence and in the presence of a quencher is measured as a function of the concentration of the quencher and plotted in a Stern–Volmer diagram. The Stern–Volmer quenching constant is derived from this plot. It is a measure of the exposure of the fluorophore to the quencher. Static quenchers do not obey the Stern–Volmer relationship. They form nonfluorescent complexes with the fluorophore, and the binding constant can be derived from the quenching experiments.

### Suggestion for Further Reading

M. R. Eftink and C. A. Ghiron (1981) Fluorescence quenching studies with proteins. *Anal. Biochem.* **114**, 199–227.

## Fluorescence Spectroscopy

## 1. Fluorescence

Fluorescence emission of light is observed when, after excitation by the absorption of a photon (see [Absorption Spectroscopy](#)), an electron returns from the first excited state back to the ground state. Absorption and emission are both virtually instantaneous processes and occur in about  $10^{-15}$ s. The average lifetime of the excited state is about  $10^{-8}$ s, which is a short time, but seven orders of magnitude longer than the time required for light absorption and emission. In the excited state, some energy is always lost by nonradiative processes (such as transitions between vibrational states). Therefore, the energy of the emitted light is always less than that of the absorbed light, and the fluorescence of a chromophore thus occurs always at greater wavelengths than its absorption. Fluorescence emission is much more sensitive to changes in the environment of the chromophore than light absorption. As the lifetime of the excited state is long, a broad range of interactions or perturbations can influence this state and thereby the emission spectrum. Fluorescence is thus an excellent and extremely sensitive probe to investigate the structure and function of proteins.

## 2. Fluorescence of Proteins

The aromatic amino acid residues [phenylalanine](#) (Phe), [tyrosine](#) (Tyr), and [tryptophan](#) (Trp) exhibit fluorescence emission when excited in the wavelength range of their absorption spectra. In proteins that contain all three aromatic amino acids, the emitted fluorescence is usually dominated by the contribution of the Trp residues, because both their absorbance at the wavelength of excitation and their quantum yield of emission are considerably greater than the corresponding values for Tyr and Phe. This is expressed by the “sensitivity” parameter (Table 1), which is 730 for Trp and 200 for Tyr. Phenylalanine fluorescence is not observed in native proteins, because its sensitivity of 4 is very low. In addition, Phe and Tyr fluorescence is also decreased by energy transfer (see [Fluorescence Energy Transfer](#)) to Trp residues.

**Table 1. Absorbance and Fluorescence Properties of the Aromatic Amino Acids<sup>a</sup>**

| Compound      | Absorbance            |   | Fluorescence          |         | Sensitivity  |
|---------------|-----------------------|---|-----------------------|---------|--|
|               | $\lambda_{\max}$ (nm) | $\epsilon_{\max}$ (M <sup>-1</sup> cm <sup>-1</sup> ) | $\lambda_{\max}$ (nm) | $f_F^b$ | $\epsilon_{\max} \times f_F^b$ (M <sup>-1</sup> cm <sup>-1</sup> ) |
| Tryptophan    | 280                   | 5600  | 355                   | 0.13    | 730  |
| Tyrosine      | 275                   | 1400  | 304                   | 0.14    | 200  |
| Phenylalanine | 258                   | 200   | 282                   | 0.02    | 4  |

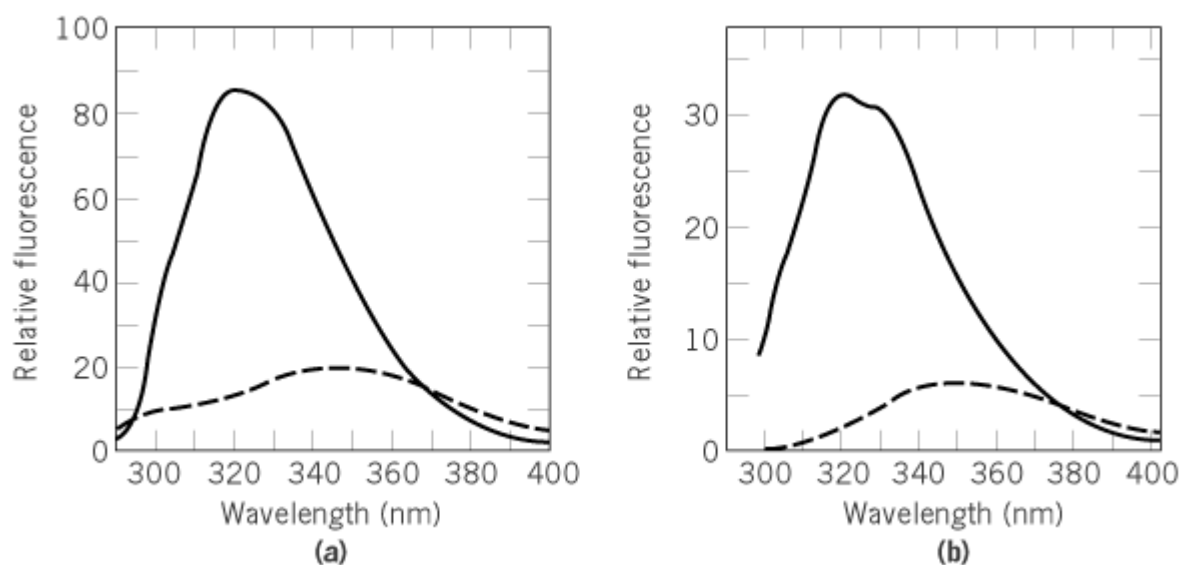
<sup>a</sup> In water at neutral pH; data are from (2).

<sup>b</sup>  $f_F$ , fluorescence quantum yield.

Changes in protein conformation, such as unfolding, very often lead to large changes in the wavelength and intensity of fluorescence emission. The fluorescence emission of Trp residues in native proteins can be either greater or smaller than the emission of tryptophan in aqueous solution. Consequently, the fluorescence intensity can either increase or decrease upon protein unfolding. The emission maximum is usually shifted from shorter wavelengths to about 350 nm, which corresponds

to the fluorescence maximum of tryptophan in aqueous solution. The exact location of this maximum depends to some extent on the nature and concentration of the [buffer](#). In a **hydrophobic** environment, such as in the interior of a folded protein, tryptophan emission occurs at shorter wavelengths (indole shows an emission maximum of 320 nm in hexane). As an example, the emission spectra of native and of unfolded **ribonuclease T<sub>1</sub>** are shown in Figure 1. RNase T<sub>1</sub> contains 9 Tyr residues and only one Trp (Trp59), which is inaccessible to solvent in the native protein.

**Figure 1.** Fluorescence emission spectra of native (—) and of unfolded (---) RNase T<sub>1</sub>. Native RNase T<sub>1</sub> (1.4 μM) was in 0.1 M sodium acetate pH 5.0; the sample of unfolded protein contained 6.0 M GdmCl in addition. Fluorescence was excited at (a) 278 nm and (b) 295 nm. The bandwidths were 3 nm for excitation and 5 nm for emission. Spectra were recorded at 25°C in 1 × 1 cm cells in a Hitachi F-4010 fluorimeter. From Ref. 1.

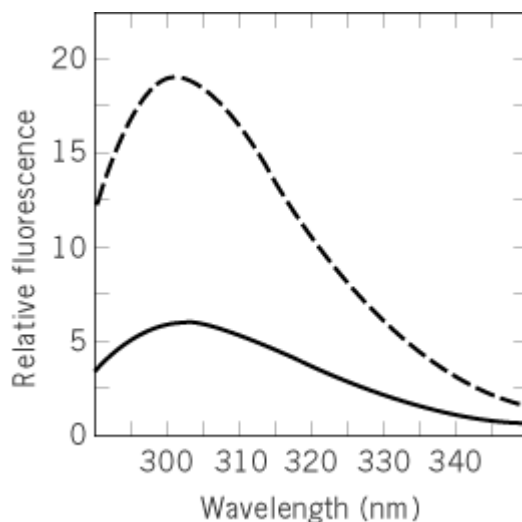


The fluorescence of the Trp residues of a protein can be investigated selectively by excitation at wavelengths greater than 295 nm. Because of the red shift and increased intensity of the absorbance spectrum of tryptophan when compared with tyrosine protein absorbance above 295 nm originates almost exclusively from Trp residues. A comparison of the emission observed after excitation at 280 and 295 nm gives information about the contribution of the Trp and Tyr residues to the observed fluorescence spectra. The data for RNase T<sub>1</sub> in Figure 1 show that the shapes of the fluorescence spectra observed after excitation at 278 and 295 nm are virtually identical. The measured emission originates almost completely from the single Trp59, which is inaccessible to solvent and hence displays a strongly blue-shifted emission maximum near 320 nm. Tyrosine emission is barely detectable in the spectrum of the native protein, because energy transfer to Trp residues occurs. Unfolding of RNase T<sub>1</sub> by **guanidinium chloride** results in a strong decrease in tryptophan fluorescence and a concomitant red shift of the maximum to about 350 nm. The distances between the Tyr residues and Trp59 increase upon unfolding, so energy transfer becomes less efficient. As a consequence, the Tyr fluorescence near 303 nm becomes visible in the spectrum of the unfolded protein when excited at 278 nm (Fig. 1a), but not when excited at 295 nm (Fig. 1b). The examples in Figure 1 indicate that, unlike absorbance, the changes in fluorescence upon folding can be very large, and the contribution of the Trp residues can be studied selectively by changing the excitation wavelength. Thus, fluorescence measurements are extremely useful to monitor conformational changes in proteins.

Multiple bands in an emission spectrum do not necessarily originate from different Trp residues of a folded protein. Figure 1 shows that even single Trp residues, such as Trp59 of RNase T<sub>1</sub>, can give rise to several emission bands. It indicates that the protein can adopt different conformations that equilibrate more slowly than the lifetime of the excited state (about 10<sup>-8</sup>s).

The fluorescence emission maximum of tyrosine remains around 303 nm, irrespective of its molecular environment. Therefore, the unfolding of proteins that contain no Trp residues is usually accompanied by changes in the intensity, but not in the wavelength of emission. As an example, the fluorescence emission spectra of the Trp59Tyr variant of RNase T<sub>1</sub> (which is devoid of Trp residues) in the native and unfolded states are shown in Figure 2. A decreased Tyr fluorescence in the native state (as in Fig. 2) is frequently observed. It is thought to originate in the folded state from **hydrogen bonding** of the tyrosyl hydroxyl group and/or the proximity of quenchers (see [Fluorescence Quenching](#)), such as [disulfide bonds](#).

**Figure 2.** Fluorescence emission spectra of 1.5 μM of native (—) and of unfolded (---) Trp59Tyr-RNase T<sub>1</sub>. This protein contains 10 Tyr residues, but no Trp residue. The native protein was in 0.1 M sodium acetate, pH 5.0; the unfolded sample contained 6.0 M GdmCl in addition. Fluorescence was excited at 278 nm, the spectra were recorded as in Figure 1. From Ref. 1.



### 2.1. Environmental Effects on Tyrosine and Tryptophan Emission

The fluorescence intensity generally decreases with increasing temperature. This decrease is substantial, and to a first approximation, Tyr and Trp emission decrease by more than 1% per degree increase in temperature. Denaturants such as [urea](#) and guanidinium chloride also influence the fluorescence of Tyr and Trp (1).

## 3. Practical considerations

### 3.1. Fluorescence Spectrophotometers

Fluorescence spectrophotometers (“fluorimeters”) are single-beam instruments. The light from an intense source (usually a Xenon lamp) passes first the excitation monochromator to select the excitation wavelength, and is then focused into the center of the cuvette. The fluorescence emission is usually monitored at a right angle relative to the exciting beam; it passes through the emission monochromator and is detected by the sample photomultiplier. Most spectrofluorimeters operate in a split-beam mode, where a small portion of the incident light is directed to a reference photomultiplier to correct the fluorescence signal for inherent instabilities of the light source. To avoid baseline drift,

some instruments interrupt the light beam periodically or employ pulsed lamps. The dark periods are used to adjust the baseline during the measurement. This increases the long-term constancy of the measured fluorescence.

### 3.2. Cuvettes for Fluorescence Measurements

Cuvettes for measuring fluorescence have polished surfaces on all four sides. Quartz cells are required for work in the UV region. For routine work, rectangular cuvettes ( $0.4 \times 1$  cm) are convenient. Since the exciting beam illuminates only the central part of the cell, the observed fluorescence intensity does not decrease necessarily when cuvettes with reduced cross sections are used. Cuvettes and glassware are easily contaminated with fluorescing substances that may leak out of plastic containers or are present in laboratory detergents. Before measurements, the cleanliness of the cell (and the distilled water) should be checked routinely by filling the cuvette with water and recording a blank in the wavelength range of interest. There should be no emission except the Raman peak of water.

### 3.3. Raman Peak of Water

In water, a Raman scattering peak is observed that is separated from the incident radiation by a fixed energy difference. After excitation at 280 nm (as in fluorescence experiments with proteins) the Raman peak of water is at 315 nm, ie, it overlaps with the emission spectrum of proteins. The Raman peak must therefore be subtracted from the measured fluorescence spectra. The performance of a fluorescence spectrometer can be checked easily by recording the signal-to-noise ratio at the Raman peak of water.

## Bibliography

1. F. X. Schmid (1997) In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 261–297.
2. M. R. Eftink (1991) In *Methods of Biochemical Analysis* (C. H. Suelter, ed.), Vol. **35**, Wiley, New York, pp. 127–205.

## Suggestions for Further Reading

3. C. L. Bashford (1987) "An introduction to spectrophotometry and fluorescence spectrometry". In *Spectrophotometry and Spectrofluorimetry: A Practical Approach* (D. A. Harris and C. L. Bashford, eds.), IRL Press, Oxford, UK, pp. 1–22.
4. C. S. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Vol. **II.**, Chap. "8", W. H. Freeman and Co., San Francisco, CA.
5. C. R. Cantor and S. N. Timasheff (1982) "Optical spectroscopy of proteins. In" *The Proteins* (H. Neurath and R. L. Hill, eds.), Academic Press, New York, pp. 145–301.
6. G. R. Penzer (1980) "Molecular emission spectroscopy". In *An Introduction to Spectroscopy for Biochemists* (S. B. Brown, ed.), Academic Press, London, pp. 70–114.

## Fluorodeoxyuridine

Fluorodeoxyuridine (FUdR) is a [mutagen](#) and a potent inhibitor of [thymidylate synthase](#) (1), which leads to depletion of cellular pools of precursors to thymidine. The accuracy of both [DNA replication](#) and [DNA repair](#) is dependent upon a balanced supply of deoxyribonucleoside triphosphate (dNTP) precursors of DNA (2, 3), and perturbation of their relative levels has far-reaching effects. These include point [mutation](#), **chromosomal** breakage, exchange or loss, as well as a stimulation of mitotic

and meiotic [recombination](#). Thus a large number of experimental systems and assays have responded positively to FUdR or its precursor 5-fluorouracil (summarized in Kunz et al. (3)), although they (unlike 5-bromo-, 5-iodo-, and 5-chlorodeoxyuridine) were not mutagenic for T4 bacteriophage in early studies (4). Unlike other nucleoside base analogues, such as [5-bromouracil](#) or [2-aminopurine](#), FUdR does not directly affect base-pairing properties, but instead causes nucleotide pool imbalances through depleting dTTP and increasing dATP and dCTP pools (3).

### Bibliography

1. S. S. Cohen, J. G. Flaks, H. D. Barner, M. R. Loeb, and J. Lichtenstein (1958) Proc. Natl. Acad. Sci. USA **44**, 1004–1012.
2. R. H. Haynes (1985) Basic Life Sci. **31**, 1–23.
3. B. A. Kunz, S. E. Kohalmi, T. A. Kunkel, C. K. Matthews, E. M. Mcintosh, and J. A. Reidy (1994) Mutat. Res. **318**, 1–64.
4. R. M. Litman and A. B. Pardee (1956) Nature **178**, 529–531.

### Fluorography

Fluorography (or photofluorography) is the technique of photographing an image produced by light emitted from a fluorescent screen or material. Light is produced by the excitation of a fluorescent material by ionizing radiation, which is produced when charged particles (electrons or beta particles) emitted from radionuclides such as tritium ( $^3\text{H}$ ), carbon-14 ( $^{14}\text{C}$ ), phosphorus-32 ( $^{32}\text{P}$ ), sulfur-35 ( $^{35}\text{S}$ ), or iodine-125 ( $^{125}\text{I}$ ) interact with the material. The emitted light exposes the film, and an image is recorded.

Fluorography is a common technique used in molecular biology [blotting](#) experiments to increase the sensitivity of medium-to-low energy beta-particle-emitting radioisotopes embedded in gel slices. A fluorescent screen is sometimes used to amplify the amount of light produced. A fluorescent screen is a sheet of material that is coated with fluorescent reagents (fluorophores, luminophores, or fluorochromes).

Fluorescence enhances the intensity of the image recorded in the film and makes it possible to obtain images in shorter time periods and with less radioactivity than conventional [autoradiography](#). Fluorographic detection of blotted samples can reduce exposure times by factors of 5 to 10 and can increase sensitivity by factors of 3 to 5 over the use of dried gels (1). Fluorography works best with flashed film and exposures of film at  $-70^\circ\text{C}$ . Film flashing helps to detect weak bands or spots. Flashing also helps to reduce the threshold of detection. Before fluorography, gels should be fixed with 30% isopropyl alcohol and 10% acetic acid to immobilize the separated proteins and to remove non-protein components that might interfere with subsequent staining.

Radiation-free fluorographic techniques for Southern blots have been developed using

chemical luminescence. Samples are placed on a chemiluminescent substrate sheet, and the blot is exposed to X-ray film for about 70 min at  $37^\circ\text{C}$ .

In other industrial and medical applications, fluorescent materials respond also to radiation from gamma-emitting [radioisotopes](#), X-ray sources, electron beams, and charged particles. A fluorograph is the photograph produced by the light from fluorescent materials, and a fluoroscope is a device for

viewing the light from fluorescing materials. A fluorescence digital imaging microscope is an instrument for viewing microscopic images of fluorescing substances, which are stimulated by ultraviolet light and filtered so that the observer sees only the fluorescence emission and not the stimulating light.

Fluorography is sometimes called abreuography in honor of Manuel de Abreau, a Brazilian physician who discovered the technique.

#### Bibliography

1. E. Quémener and F. Simonnet (1995) *Biotechniques* **18**, 100–103.

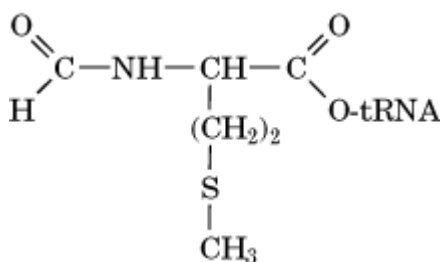
#### Suggestions for Further Reading

2. A. W. Rogers (1983) *Practical autoradiography*. Elmsford Rev. **23**, 5–39.
3. N. H. H. Heegard (1988) "Radioimmuno-detection-immunovisualization by means of <sup>14</sup>C-labeled anti-antibodies and sodium salicylate fluorography on nitrocellulose". In *CRC Handbook of Immunoblotting of Proteins*, Vol. **1**, Technical Descriptions (O. J. Bjerrum and N. H. H. Heegard, eds.), CRC Press, Boca Raton, FL, pp. 189–194.
4. Amersham International, *Amersham Review* **23** (1984) Radioisotope detection by fluorography and intensifying screens. *Amersham Res. News*, pp. 1–39.

#### fMET (*N*-Formyl Methionine)

Protein biosynthesis usually starts with *N*-formyl methionine (Fig. 1), which is directed by the [initiation codon](#) AUG, or rarely in bacteria by GUG or UUG. The same triplet AUG is decoded as a normal methionine during elongation. These distinct types of recognition are carried out by distinct [transfer RNAs](#), initiator tRNA<sup>fMet</sup> and elongator tRNA<sup>Met</sup>. The initiator tRNA carries a methionine residue that has been formylated on its amino group (*N*-formyl-methionyl-tRNA<sup>fMet</sup>, or fMet-tRNA<sup>fMet</sup>). Synthesis of fMet-tRNA<sup>fMet</sup> proceeds in two steps: (i) aminoacylation of the initiator tRNA with methionine and (ii) formyl transfer from formyltetrahydrofolate to the methionine amino group. The latter reaction is catalyzed by a formyltransferase.

Figure 1. fMet (*N*-formyl methionine).



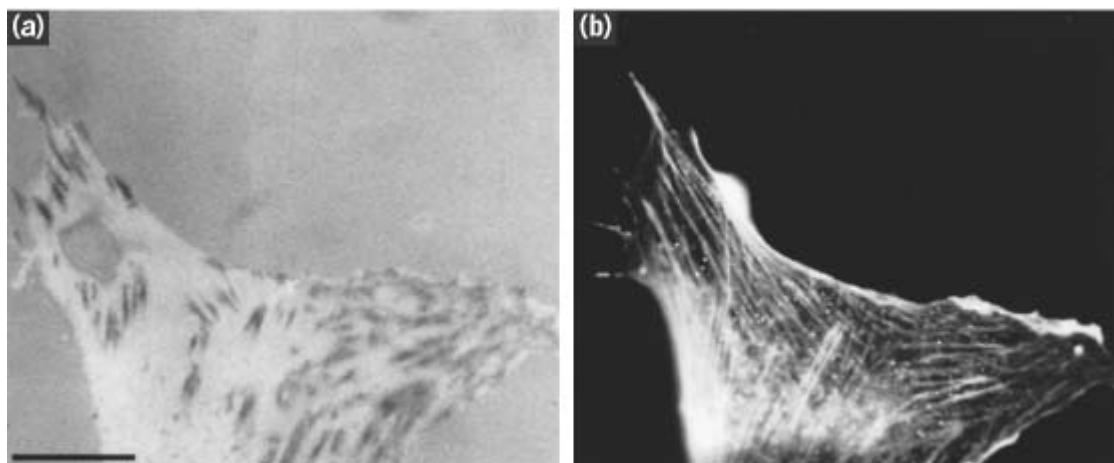
Only fMet-tRNA<sup>fMet</sup> is used for initiation of protein synthesis, because it enters part of the P site on

the prokaryotic 30S or eukaryotic 40S [ribosome](#) subunit (see [Initiation Complex](#) and [Initiation Factor \(IF\)](#)). Often the *N*-terminal formyl residue is removed by a specific deformylase enzyme. Sometimes the methionine residue at the terminus is removed by an aminopeptidase. These reactions occur rather rapidly, probably when the nascent polypeptide chain reaches the length of 15 to 30 amino acids.

## Focal Contact

Focal contacts, focal adhesions, or adhesion plaques are [cell “junctions”](#) by which many cell types in culture adhere to their substratum, being regions of closest approach between the lower cell surface and the substratum, a separation of approximately 10 nm. Elongated structures that may be 3 to 4  $\mu$ m in length, they may be recognized as black areas by interference reflection [microscopy](#) (Fig. 1), by [electron microscopy](#), or by staining with [antibodies](#) against one of their numerous components (1, 2) (see text below). These contacts are points at which the actin [cytoskeleton](#) associates with the cell surface, and are located at the cell-surface termini of actin **microfilament bundles**, called stress fibers. Focal contacts and stress fibers do not have clear equivalents *in vivo*, although their molecular composition resembles that of the dense bodies of smooth muscle. Nevertheless, the molecules, molecular interactions, and cellular functions that have been discovered by studying focal contacts are likely to be of great importance in regulating cell behavior *in vivo*. In tumor cells, disruption of focal contacts is associated with the transformed phenotype (3). The components of focal contacts are involved both in [extracellular matrix](#) adhesion and [signal transduction](#).

**Figure 1.** (a, b) Interference reflection image (a) and fluorescent image showing phalloidin staining for F-actin (b) of a cultured fibroblast, illustrating focal contacts and stress fibers. Focal contacts are the elongated black areas in a. Note how the actin stress fibres terminate at focal contacts at their peripheral ends. (Photographs provided by Dr. G. Ireland.) Bar=10  $\mu$  m.

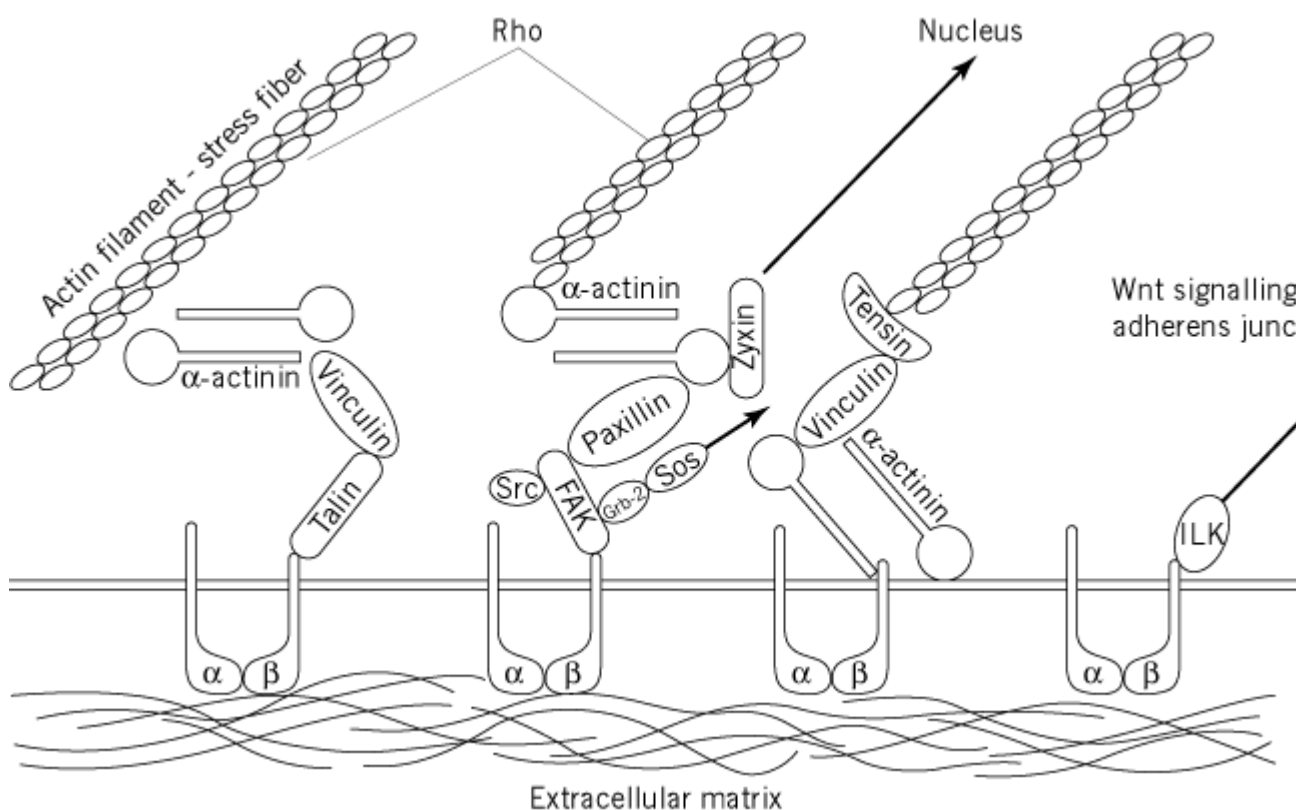


The principal adhesive components of focal contacts are [integrins](#) (4) (Fig. 2) although they also contain membrane heparan sulfate proteoglycan. These adhesion receptors bind the cell to components of the extracellular matrix, such as [fibronectin](#) and vitronectin. Several [actin-binding proteins](#) associate with the cytoplasmic side of the plasma membrane at focal contacts. These include  $\alpha$ -actinin (2), vinculin (5), talin (6), and tensin (7), the first two of which are also found in



[intermediate junctions](#). Talin and  $\alpha$ -actinin have been shown to bind to the cytoplasmic domains of integrins, particularly the  $\beta 1$  subunit (6, 8). Tensin has [actin](#) filament-capping activity. The location of these proteins to focal contacts and their binding interactions are based on a variety of types of evidence, including expression of mutant proteins and *in vitro* binding between purified, recombinant, or synthetic proteins or peptide fragments. The precise manner in which these proteins interact in focal contacts is not understood. Collectively, they provide a link between the integrin adhesion receptors and the actin cytoskeleton.

**Figure 2.** Diagram showing some of the major components of focal contacts and their possible molecular interactions. Grb 2 = growth factor receptor bound protein 2, SOS = son-of-sevenless, ILK = integrin-linked kinase,  $\alpha$  = integrin a subunit,  $\beta$  = integrin b subunit. Arrows indicate signaling pathways, details of which are not shown. For further details, see text.



Cells must be able to break and reform focal contacts in order to move. They have been shown to contain the [protease](#) calpain, which is able to cleave talin and may be involved in modulating junctional adhesion (9). Alternatively, **phosphorylation** signals generated from within the cytoplasm may modulate the adhesiveness of focal adhesions through regulation of integrin function (10). The oncogenic **tyrosine kinase** pp60<sup>src</sup> and protein kinase C, both present in focal contacts, appears to be able to modulate integrin-mediated adhesion, reducing cell-substratum adhesion (11). Such an adhesive change is likely to be important in tumorigenicity, increasing the invasive properties of cells.

An important function of focal contact components is the transduction of signals from outside the cell. Such signals may be involved in regulating a variety of cellular processes, including motility, cell division, gene expression, and [apoptosis](#). An example of a focal adhesion signaling molecule is the focal adhesion kinase (FAK) (12). This is a tyrosine kinase that is activated by ligand binding to integrins (13). FAK has been shown to bind to the cytoplasmic domain of the  $\beta 1$  integrin subunit and also to the vinculin-binding protein paxillin (14). Phosphorylation of FAK in response to integrin-

ligand binding increases its tyrosine kinase activity (15). It is also able to associate with a number of other signaling proteins via its various phosphotyrosine residues. Thus, it binds pp60<sup>src</sup> via the latter's src homology (SH2) domain (16) and similarly binds growth factor receptor-bound protein, or Grb2 (17). The latter interacts with the [guanine nucleotide exchange factor](#) mSOS1, which activates the **Ras** signaling pathway. FAK may also be involved in regulating focal contact turnover, because fibroblasts from FAK <sup>-/-</sup> mice show increased numbers of more persistent focal contacts (18).

Another focal contact-associated signaling molecule is the integrin-linked kinase (ILK), a **serine-threonine kinase** that can bind directly to the cytoplasmic domains of b1 and b3 integrins (19). Its kinase activity modulates cell–matrix interactions. Overexpression of ILK in epithelial cells results in reduction of tumorigenicity and loss of cell–cell adhesion and anchorage-independent growth. Furthermore, it causes down-regulation of expression of the zonula adherens adhesion molecule E-cadherin, translocation of b-catenin to the nucleus, and formation of a transcriptional regulating complex between b-catenin and the [transcription factor](#) LEF-1. Thus, ILK mediates cross-talk between cell–matrix and cell–cell adhesion, as well as regulation of the wnt signaling pathway (20).

Both paxillin and another focal-contact associated protein called zyxin possess so-called LIM domains, **zinc-finger** motifs that are also found in a number of proteins involved in the regulation of gene expression (21, 22). Paxillin becomes tyrosine-phosphorylated on matrix ligand binding, suggesting that it is involved in signal transduction (21). Zyxin binds to a-actinin, binds to other proteins that have potential signaling functions, such as the vasodilator-stimulated phosphoprotein (VASP), and has a nuclear export sequence. It has also been shown to shuttle between the focal contact and the nucleus, suggesting that it may provide a direct signaling link between matrix receptors at the cell periphery and gene regulation (23).

Regulation of focal contact formation is also regulated by the small [Rho GTPase](#). Overexpression of Rho in quiescent cells causes focal contact formation and stress fiber assembly, whereas inhibition of Rho leads to focal contact disassembly (24). This aspect of Rho signaling is one of several integrated signaling pathways involving small [Gtpases](#) that participate in regulating cell adhesion, cell motility, and actin cytoskeleton.

## Bibliography

1. J. Heath and G. A. Dunn (1978) *J. Cell Sci.* **29**, 197–212.
2. J. Wehland et al. (1979) *J. Cell Sci.* **37**, 257–273.
3. K. Burridge (1986) *Cancer Rev.* **8**, 18–78.
4. W. T. Chen et al. (1985) *J. Cell Biol.* **100**, 1103–1114.
5. B. Geiger et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4127–4131.
6. A. F. Horwitz et al. (1986) *Nature* **320**, 531–533.
7. A. J. Wilkins et al. (1986) *J. Cell Biol.* 102–103.
8. C. A. Otey et al. (1990) *J. Cell Biol.* **111**, 721–729.
9. M. C. Beckerle et al. (1987) *Cell* **51**, 569–577.
10. M. Cappelino et al. (1995) *J. Biol. Chem.* **270**, 23132–23138.
11. S. Kellie (1988) *Bioessays* **8**, 25–30.
12. M. D. Schaller et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5192–5196.
13. L. Lipfert et al. (1992) *J. Cell Biol.* **119**, 905–912.
14. J. D. Hildebrand et al. (1993) *J. Cell Biol.* **123**, 993–1005.
15. Burridge et al. (1992) *J. Cell Biol.* **119**, 893–903.
16. B. S. Cobb et al. (1994) *Mol. Cell Biol.* **14**, 147–155.
17. D. D. Schlaepfer et al. (1994) *Nature* **327**, 786–791.

18. D. Ilic et al. (1995) *Nature* **377**, 539–544.
19. G. E. Hannigan et al. (1996) *Nature* **379**, 91–96.
20. A. Novak et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 4374–4379.
21. C. E. Turner and J. T. Miller (1994) *J. Cell Sci.* **107**, 1583–1591.
22. I. Sadler et al. (1992) *J. Cell Biol.* **119**, 1573–1588.
23. D. A. Nix and M. C. Beckerle (1997) *J. Cell Biol.* **138**, 1139–1147.
24. A. J. Ridley and A. Hall (1992) *Cell* **70**, 389–399.

### Suggestions for Further Reading

25. R. O. Hynes (1992). Integrins: versatility, modulation and signalling in cell adhesion. *Cell* **69**, 11–25. (An excellent review of the field up to 1992.)
26. S. Dedhar and G. E. Hannigan (1996) Integrin cytoplasmic interactions and bidirectional transmembrane signalling. *Curr. Opin. Cell Biol.* **8**, 657–669. (An exciting review of many aspects of integrin signaling.)
27. M. C. Beckerle (1997) Zyxin: zinc fingers at sites of cell adhesion. *Bioessays* **19**, 949–957. (A good general consideration of focal contacts as well as an excellent zyxin review.)

### Footprinting

**Protein** binding sites on DNA can be detected *in vivo* and *in vitro* by a technique termed “footprinting”. This depends on the principle that a bound protein molecule can protect DNA from cleavage or modification by a nuclease or a chemical reagent. The binding site is then defined by a region of reduced reactivity, which is determined by analyzing the products of the cleavage reaction on denaturing polyacrylamide gel. The footprint of the protein appears as a gap in the normal pattern of cleavage.

*In vitro* DNase I is commonly used for enzymatic footprinting and dimethylsulfate and hydroxyl radicals are employed for chemical footprinting. DNase I targets the minor groove of DNA and is also sensitive to the local conformation of the bound DNA such that cleavage is enhanced where the minor groove is on the outside of DNA wrapped around the protein. Dimethylsulfate under footprinting conditions methylates the N7 position of guanine in the major groove and the N2 position of adenine in the minor groove. This methylation sensitizes the sugar-phosphate backbone to mild alkaline hydrolysis, enabling the sites of reaction to be determined. For hydroxyl radical footprinting the radicals are generated by the reaction of hydrogen peroxide with an Fe<sup>2+</sup>–EDTA complex (the Fenton reaction) and they react with a backbone sugar resulting in the cleavage of the backbone. This reagent targets the minor groove and lacks sequence specificity, although it is sensitive to groove width, cleaving preferentially where the minor groove is widened. A similar type of reagent, Mn-porphyrin, has the opposite conformational selectivity, cleaving DNA where the minor groove is narrow.

### Footprinting Proteins

Protein footprinting is a methodology to gain information about protein conformation and interactions by probing the structure of [polypeptide chain](#) that is labeled at one end. The protein of interest is end-labeled and reacted with a probe in such a way that susceptibility of a site on the protein to the probe leads to breakage of the polypeptide backbone nearby. The sites of backbone scission are deduced from the lengths of peptides containing the end label, which can be estimated readily by [SDS-PAGE](#). Alteration of the backbone cleavage pattern indicates that the end-labeled protein has undergone a conformational change and/or interaction with other macromolecules. The principle is the same as that for **footprinting nucleic acids**, which has been used widely. Protein footprinting has only recently become feasible with the advent of techniques to end-label proteins.

Information on scission sites or modification sites can be obtained by [peptide mapping](#) as well, but footprinting has three advantages:

1. *Facility* Peptide mapping requires purification of individual cleavage products and [amino acid analysis](#) and/or the *N*-terminal sequence of each (see [Protein Sequencing](#)). Footprinting determines the lengths of many peptides, and hence many cleavage sites, in a single operation without isolating individual fragments.
2. *Sensitivity* Determination of the *N*-terminal amino acid sequence of a peptide typically requires 10 pmol of peptide. Footprinting is generally more sensitive; for instance, **radiolabeling** has virtually no detection limit.
3. *Selectivity* In large multisubunit complexes, peptide mapping will be complicated by the plethora of peptides to be purified and analyzed. In footprinting, information is obtained only on the polypeptide chain that was end-labeled.

## 1. Methods to End-Label Proteins

Currently there are three ways to end-label proteins:

1. Chemical modification of the *N*-terminus ([1](#), [2](#)). The protein is first subjected to one round of [Edman Degradation](#). All [amino groups](#) become phenylthiocarbamylated, except for the one at the very *N*-terminus, which was originally the imino group of the second amino acid residue before the Edman reaction. This unique amino group can be reacted with either a radioactive or a **fluorescent** reagent, resulting in labeling of the *N*-terminus. This method inevitably destroys the folded structure of the labeled protein prior to probing, and information can be obtained only at the level of the [primary structure](#).
2. Specific [antibody](#) binding to either terminus ([3](#), [4](#)). In this method, an antibody raised against a peptide comprising the amino acid sequence of either terminus is used to detect peptides containing the terminus. A more convenient alternative is to place an [epitope](#), against which an antibody is already available, at one end of a protein and to express the epitope-tagged protein through recombinant DNA techniques (see [Protein Engineering](#)).
3. Phosphorylation at either terminus ([5](#), [6](#)). This is based on the sequence specificity of protein **kinases**, such as heart muscle kinase, which is less stringent than antigen-antibody recognition. An amino acid sequence that can be phosphorylated by a kinase is attached to a protein at the gene level. The site is labeled using the kinase with [ $\text{g-}^{32}\text{P}$ ] or [ $\text{g-}^{35}\text{S}$ ] ATP.

## 2. Methods to Probe Protein Structure

As is apparent from the principle of nucleic acid footprinting, structure probing has to result in scission of the polypeptide backbone, either directly or indirectly. There are common advantages and disadvantages for each type of method. In direct probing methods, probing immediately leads to backbone scission. An advantage of direct methods is their simplicity in comparison with indirect ones, which involve additional reactions. One major drawback is that scission of the peptide

backbone often renders the fragments generated more susceptible to ensuing scissions, particularly when the scission has occurred within a folded structural unit. This causes a problem, especially in quantitative analyses, when the first cleavage sites have to be identified. In indirect methods, chemical modification of a side chain is followed by scission of an adjacent [peptide bond](#) (or prevention thereof) after one or more subsequent steps. Because modification of a solvent-accessible side chain in one entire protein molecule is unlikely to perturb the folded structure extensively, or to enhance the reactivity of other residues of the same molecule, indirect methods should not suffer from the problems that direct methods do.

Hitherto, five probing methods, two direct and the other three indirect, have been devised and employed:

1. One direct method is **proteolysis** (4-6). Endoproteinases have been extensively used to delineate the **domain** structures of proteins, because well-folded structural units are more resistant to proteolysis than regions linking them. Interaction with other macromolecules may block some cleavage sites. In addition, conformational changes are often reflected in changes in susceptibility of sites cleavable by [proteinases](#). Although proteolysis has been useful, the biophysical basis of the interaction between the probing proteinase and the probed protein molecule is not entirely understood.
2. The other direct method uses oxygen radical, typically generated by the Fenton reaction, to cleave the backbone (7). It is likely that the reactivity of the oxygen radical with the peptide backbone is more directly related to solvent accessibility than is that of proteinases. The extent of side-chain damage by the oxygen radical per backbone breakage is not known.
3. One indirect method utilizes a combination of reversible and irreversible modifications of [lysine](#) residues (8). The protein is first subjected to limited citraconylation (see [Anhydrides](#)). Removal of the citraconyl groups after irreversible **acetylation** of the remaining lysines recovers unmodified lysyl residues only at the sites where the first modification was introduced. Complete digestion of these polypeptide chains by a lysine-specific endoproteinase generates fragments that end at a first modification site.
4. A second indirect method utilizes oxidation of [methionine](#) residues, which blocks cleavage at such residues by **cyanogen bromide** (9). Oxidation of individual methionines is observed as a decrease in the amount of peptide fragments generated by cyanogen bromide. In such cases where modification prevents cleavage, the extent of cleavage must be carefully controlled to ensure that all observable peptides will not be reduced to the shortest peptide containing the end label.
5. The last method utilizes cyanylation of cysteine residues by 2-nitro-5-thiocyanobenzoic acid (10). Raising the pH induces slow cleavage of the peptide bond at the S-cyanocysteine, thus turning the modified side chain into a cleaver. Because the cyano group can transfer to a free sulfhydryl, ie, an unmodified cysteine within the same peptide, during the slow cleavage, unreacted cysteine residues have to be blocked beforehand, eg, with N-ethylmaleimide.

The latter three methods were devised on the basis of side-chain modification schemes developed in the 1960s and 1970s (see [Chemical Modification](#)). More of them may be revived to create new footprinting methods.

### 3. Variations and Future Prospects

Various tagging technologies have been used widely in molecular and cellular biology. This indicates that use of a label is not limited to detection of peptides. For instance, some end labels can also serve to isolate terminal fragments. DNA-binding sites of a protein were delimited by immunoprecipitating peptide fragments with a terminal epitope tag, following transfer of radioactivity from DNA and partial proteolysis of the radiolabeled protein (11). Although SDS-PAGE has a fairly good resolution, other methods of determining molecular weights of proteins,

such as [mass spectrometry](#), can also be employed after peptides are separated with the use of a tag. In another example of use of a tag other than labeling, a protein kinase site was used to identify protein–protein interactions that blocked phosphorylation of the site ([12](#)).

A complex of iron and [EDTA](#), Fe-EDTA, that generates oxygen radical can be [crosslinked](#) to specific sites on macromolecules ([13](#)). This attachment of a cleavage center onto a macromolecule rather than remaining free in solution allows mapping in the vicinity of the cleavage center ([13](#)), when combined with end labeling of its interaction partner, be it a protein molecule or a DNA molecule.

Because the principle of footprinting is the same for proteins and for nucleic acids, much parallel can be drawn from the longer experience of nucleic acids footprinting. The interference experiment scheme ([14](#)), in which macromolecules are first subjected to modification and later separated according to their remaining functionality, has yet to be applied to proteins.

### Bibliography

1. D. G. Jay (1984) *J. Biol. Chem.* **259**, 15572–15578.
2. R. A. Jue and R. F. Doolittle (1985) *Biochemistry* **24**, 162–170.
3. P. Matsudaira, R. Jakes, L. Cameron, and E. Atherton (1985) *Proc. Natl. Acad. Sci. USA* **82**, 6788–6792.
4. J. E. Lindsley and J. C. Wang (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10485–10489.
5. R. Hori, S. Pyo, and M. Carey (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6047–6051.
6. V. Nktinis, J. Turner, and M. O'Donnell (1996) *Cell* **84**, 137–145.
7. T. Heyduk and N. Baichoo, and F. Heyduk (2001) *Metal Ions Biol. Syst.* **38**, 255–287.
8. R. Hanai and J. C. Wang (1994) *Proc. Natl. Acad. Sci. USA* **91**, 11904–11908.
9. M. V. de Arruda, H. Bazari, M. Wallet, and P. Matsudaira (1992) *J. Biol. Chem.* **267**, 13079–13085.
10. B. P. Tu and J. C. Wang (1999) *Proc. Natl. Acad. Sci. USA* **96**, 4862–4867.
11. N. F. Lue, A. Sharma, A. Mondragon, and J. C. Wang (1995) *Structure* **3**, 1315–1322.
12. P. T. Stukenberg, J. Turner, and M. O'Donnell (1994) *Cell* **78**, 877–887.
13. S. A. Datwyler and C. F. Meares, (2001) *Metal Ions Biol. Syst.* **38**, 213–254.
14. U. Seibenlist and R. B. Simpson, and W. Gilbert (1980) *Cell* **20**, 269–281.
15. T. M. Rana and C. F. Meares (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10578–10582.
16. U. Siebenlist, R. B. Simpson, and W. Gilbert (1980) *Cell* **20**, 269–281.

### Founder Cell

*Founder cell* is a term that has been used to describe an aspect in the initiation of aggregation in [slime molds](#), especially [Dictyostelium](#). These organisms are normally single cells, but when environmental conditions become unfavorable, they aggregate to form a composite. One cell in a population, the founder cell, will begin signaling to the others by secreting pulses of a chemoattractant. *Dictyostelium discoideum* uses the normally intracellular [second-messenger](#) molecule [cyclic AMP](#) (cAMP) as its chemoattractant. Other species use a variety of chemoattractants, such as pterin derivatives and small [peptides](#). Cells in the immediate area of the founding cell respond in two ways. They first migrate toward the chemoattractant source, and then

they secrete additional cAMP signal. An extracellular **phosphodiesterase** degrades the cAMP between the waves. The resulting temporal and spatial waves of the chemoattractant flow through the environment. As the cells respond, they move together, then pause and move again, until almost all the cells in an area are in a central aggregate. The timing of the pulses and the rate of motility, combined with the enzymatic rates of signal degradation and secretion, determine the size and shape of the aggregate (1).

Exactly what causes a single cell in a population of apparent equals to become a founder cell is unknown, but it must be tightly regulated for the aggregates to end up with an optimum number of cells. Allowing cells to divide to increase the size of the aggregate is not an option, because [cell-cycle](#) arrest occurs very early in the developmental program. Once a cell adjacent to a founder responds, that cell cannot itself become a founder. There are some data that suggest that the differentiation events that occur in the founder cell do so before neighboring cells become competent for [chemotaxis](#). Perhaps a diffusible factor secreted from the founder is involved in that process (2).

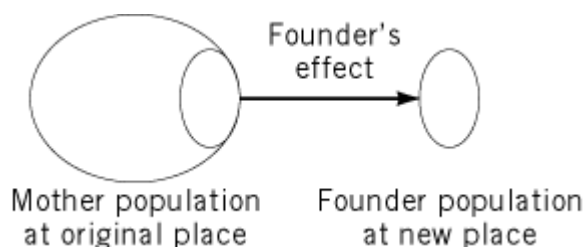
### Bibliography

1. E. Palsson et al. (1977) Proc. Natl. Acad. Sci. USA **94**, 13719–13723.
2. R. K. Raman (1976) J. Cell Sci. **20**, 497–512.

### Founder Effect

When a small number of individuals from an original population immigrate to a certain place, the new population founded by those people might have much less genetic diversity compared with those of the original population (Fig. 1). This is generally the case, because a few founder people cannot have sufficient genetic variation to reflect that of the mother population. Such a drastic decrease of genetic diversity within a population is known as the founder effect. This effect most often occurs when a segment of a larger population becomes isolated as the result of colonization of a geographically isolated area, a catastrophic event, artificial isolation, and so on. These founding populations contain only a fraction of the genetic diversity of the population from which they originated.

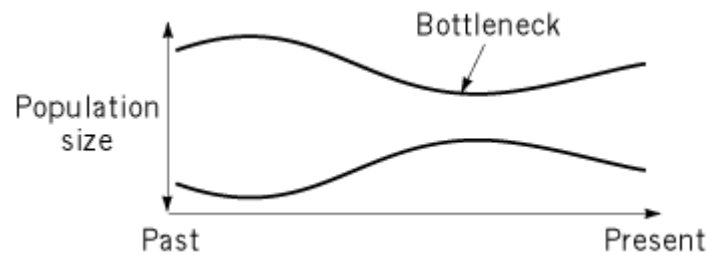
**Figure 1.** Founder effect.



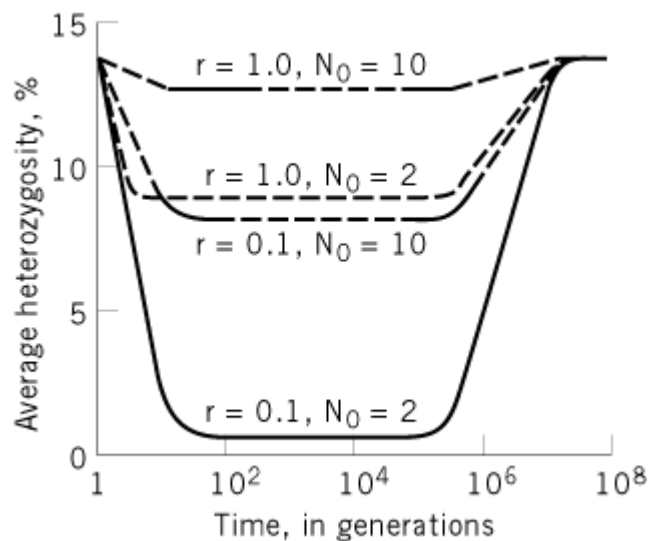
When founders are individuals sampled randomly from the mother population, the founder effect has the same influence as [genetic drift](#), which is the stochastic change of gene frequencies due to random mating. The effect of genetic drift becomes stronger when the population size is smaller (1);

therefore, the founder effect is similar to the effect of genetic drift when the size of the initial population in a given geographical area suddenly becomes extremely small. In particular, the drastic reduction in size of a population has the same effect, which is now called the “bottleneck effect” (Figs. 2 and 3). Strictly speaking, the founder effect can be considered to be a particular form of the bottleneck effect. The difference between the founder effect and the bottleneck effect is that the former occurs in a different place, whereas the latter takes place in the same location.

**Figure 2.** Bottleneck effect due to a transient decrease in the population size.



**Figure 3.** Decrease in the heterozygosity of individuals of a population resulting from the bottleneck effect. The bottleneck size is given by  $N_0$ . The intrinsic growth rate is given by  $r$ . Following the bottleneck, it is assumed that the population grows uniformly (2).



#### Bibliography

1. M. Nei, T. Maruyama, and R. Chakraborty (1975) *Evolution* **29**, 1–10
2. M. Nei (1987) *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.

#### Four-Helix Bundle Motif



The four-helix bundle is a common [protein motif](#) that is observed in [protein structures](#), in which four **a-helices** pack together lengthwise, in an antiparallel manner (Fig. 1). The four a-helices that form the bundle can derive from four independent [polypeptide chains](#) or from a single polypeptide chain. This structural motif has been observed both as an isolated three-dimensional fold and as a **domain** within much larger and more complex protein structures. The helices of four-helix bundles are longer than average (~20 residues, as opposed to the usual 10) and are [amphipathic](#), with **hydrophobic** residues buried in the core. The antiparallel packing of the helices in the bundle may be favored by interaction between the helix dipoles.

**Figure 1.** Schematic representation of the backbone of a four-helix bundle protein, taken from the structure of bullfrog ferritin (1). This structure has the up–down–down–up topology; the two “up” helices are shown as purple coils and the two “down” helices are shown as green coils. Note the long connecting loop between the second and third helix. Connecting loops and an additional C-terminal helix are shown in yellow. The N- and C-termini are labeled. This figure was generated using Molscrip (2) and Raster3D (3, 4). See color insert.



The four-helix bundle is one of the simplest protein folds and is used by a wide range of proteins with diverse functions. Proteins with a four-helix bundle structure include [growth hormones](#), **cytokines**, and [ferritins](#). The four helices may splay apart to generate a binding pocket for **cofactors** and metal ions, or they may coil around each other to form a supercoil. In the latter case, the twist of the helices in the [coiled coils](#) alters the usual helical structure so that there are 3.5 residues per turn, rather than 3.6. This means that every seventh residue in the helix is identically placed with respect to the axis of the helix. The seven positions of this [heptad repeat](#) are denoted *a* to *g*, and residues *a* and *d* point into the core of the bundle and are usually hydrophobic.

Several classes of four-helix bundle motif have been defined. These include (a) the simple up–down–up–down topology, where consecutive helices have short connections and alternating directions, (b) the up–down–down–up topology, where there is a long loop or crossover connection between helices 2 and 3, and (c) the up–up–down–down topology that is typical of the cytokine four-

helix bundles, with two long crossover connections between helices 1 and 2 and helices 3 and 4. In most cases, the helices are aligned so that each pairwise interaction is antiparallel. There are also sub-classes within each of these classifications. The simplicity of the four-helix bundle motif has made it a common target for protein design.

[See also [Alpha-Helix \(310-Helix and Pi-Helix\)](#) and [Protein Motif](#)]

#### Bibliography

1. J. Trikha, E. C. Theil, and N. M. Allewell (1995) *J. Mol. Biol.* **248**, 949–967.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

5. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
6. S. Kamtekar and M. H. Hecht (1995) The four-helix bundle: what determines a fold? *FASEB J.* **9**, 1013–1022.
7. N. L. Harris, S. R. Presnell, and F. E. Cohen (1994) Four-helix bundle diversity in globular proteins. *J. Mol. Biol.* **236**, 1356–1368. (Analyzes the diversity of four-helix bundles and gives a suggested structural taxonomy.)

## Frameshift Mutation

A frameshift mutation is a type of [mutation](#) that results from the addition or deletion of  $(3n \pm 1)$  **nucleotides** (where  $n$  is an integer). The addition or deletion of nucleotides that are not multiples of three disrupts the triplet [reading frame](#) of the gene (see [Genetic Code](#)). Technically defined, a frameshift mutation is within the coding portion of the gene, but the term is commonly used for any insertion or deletion of a few bases. Frameshift mutations were first described by Crick *et al.* ([1](#)) who used them to decipher the triplet nature of the genetic code. Although point mutations only alter one (or no) amino acids (see [Missense Mutation](#)), frameshifts alter nearly all of the amino acids downstream of the mutation until a **nonsense codon** is encountered during [translation](#) and the [polypeptide chain](#) is terminated. Frameshift mutations are induced by a wide variety of intercalating agents, such as [ethidium bromide](#).

#### Bibliography

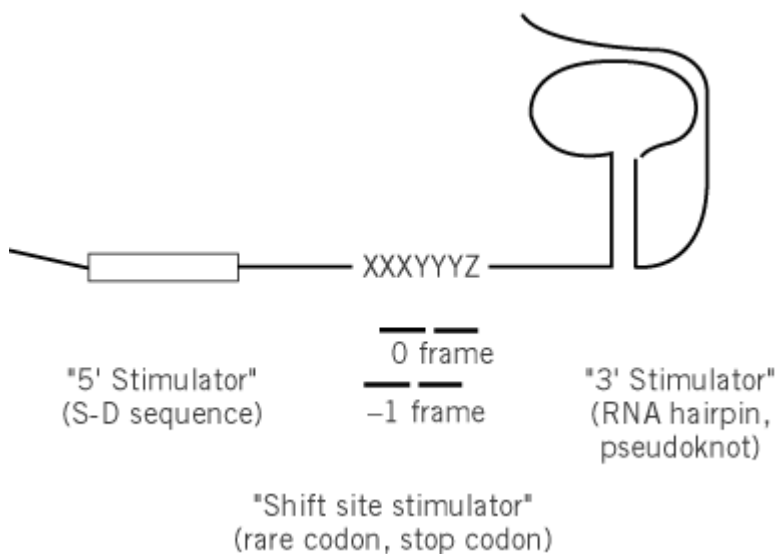
1. F. H. C. Crick, L. Barnett, S. Brenner, and R. J. Watts-Tobin (1961) *Nature (London)* **192**, 1227–1232.

## Frameshifting

Frameshifting is a change in [reading frame](#) during **gene** expression. [Transcription](#) errors, or a posttranscriptional process such as [RNA Editing](#), can cause frameshifting. For example, a fraction of the b- [amyloid precursor protein](#) and ubiquitin B [messenger RNAs](#) in the lesions of Alzheimer's disease are subject to transcriptional frameshifting (1). Most of the frameshifting cases, however, occur during translation (2).

The reading frame is accurately maintained during translation. The average occurrence of frameshift errors is estimated to be  $5 \times 10^{-5}$  per codon or lower. Aminoacyl tRNA imbalance, due to amino acid starvation or mutations in translation factors or [aminoacyl-tRNA synthetases](#), elevates the error rate. A class of [suppressor tRNAs](#) increases frameshifting at particular codons (3). In these cases, the effects are nonspecific. On the other hand, a minority of mRNAs have specific sites for frameshifting that is essential to synthesize the correct gene product. Such natural frameshifting is genetically programmed by specific signals on the mRNAs and, in some cases, is highly efficient (up to 50%). Almost all the known programmed frameshifting is either +1 (reading frame shifts to the 5' direction by one nucleotide) or -1 frameshifting (to the 3' direction by one nucleotide). An extreme exception is bypassing of 50 nucleotides within the coding region of **T4 phage** gene 60. The most crucial signal is the sequence of the frameshift site itself. Frameshift sites are often accompanied by a [cis-acting](#) sequence element (or stimulator) located either 3' or 5' to the frameshift site. Stem-loop or pseudoknot structures of RNA and the **Shine-Dalgarno** sequence are known to function as 3' and 5' stimulators, respectively (Fig. 1). In some cases, a frameshift site itself contains a rare codon or a stop codon as a shift-site stimulator. The basic function of these stimulators is to pause the elongating ribosome to trigger frameshifting. Natural frameshifting stands as a class of reprogrammed genetic decoding (or recoding).

**Figure 1.** -1 frameshift site motif and stimulators.



Most frameshifting has been found in *virus* genomes and mobile genetic elements. The former includes several [retroviruses](#) such as human immunodeficiency virus 1 (**HIV-1**) and plant RNA viruses, as well as bacteriophages such as MS2, T7, and I, while the latter includes yeast [retrotransposons](#) Ty1 and Ty3, and bacterial **insertion sequences** (Table 1). In many cases, frameshifting is required for the expression of [polymerases](#) (including **reverse transcriptase** and [RNA-dependent RNA polymerase](#)). Some retroviruses, such as human T-cell leukemia virus 1 (HTLV-1) or mouse mammary tumor virus (MMTV), require two frameshifting events to access the reverse transcriptase reading frame. Retroviral frameshifting appears to be a mechanism to keep the

synthesis rate of the functional proteins small relative to that of the structural proteins. Although the efficiencies of these frameshifting events are different from one case to another, proper efficiencies are important for optimal viral propagation or [transposition](#) in the cases investigated. In viruses and mobile elements,  $-1$  frameshifting is predominant, whereas some viruses and the yeast retrotransposons have  $+1$  frameshift sites.

**Table 1. Natural Frameshifting in Wild-Type Genes (2)**

| Designation               | Sign | Double Shift | Stimulator and Features  |
|---------------------------|------|--------------|--|
| Mammalian antizyme        | +1   | No           | 3' pseudoknot; regulated by pathway product, polyamines (autoregulation)         |
| Retrovirus <i>gag-pol</i> |      |              |  |
| RSV                       | -1   | Yes          | 3' pseudoknot  |
| HTLV-1                    | -1   | Yes          | 3' stem loop   |
| HIV-1                     | -1   | Yes          | 3' stem loop   |
| Coronavirus IBV           | -1   | Yes          | 3' pseudoknot  |
| Yeast retrotransposon Ty1 | +1   | No           | Shift site rare arginine codon   |
| Yeast virus L-A           | -1   | Yes          | 3' pseudoknot  |
| Bacterial genes           |      |              |  |
| RF-2                      | +1   | No           | 5' Shine–Dalgarno sequence; autoregulation                                       |
| <i>dnaX</i>               | -1   | Yes          | 5' Shine–Dalgarno sequence; 2 DNA polymerase subunits                            |
| IS insertion elements     | -1   | Yes          | 3' stem loop or pseudoknot   |
| Phage T7 gene 60          | +50  | No           | Matched take off and landing sites; 5' polypeptide, 3' stem loop, and stop codon |

There are only a few instances of frameshifting in bona fide cellular genes. In bacteria, two chromosomal genes, *prfB* encoding release factor 2 (RF-2; see [Release Factor](#)) and *Escherichia coli dnaX* encoding the  $\beta$  subunit of **DNA polymerase III**, are expressed by frameshifting. The reading frame of RF-2 switches to the  $+1$  frame, while DnaX shifts to the  $-1$  frame. In higher eukaryotes, there is only one cellular gene known thus far to have frameshifting, the gene for antizyme that is a regulatory protein of the cellular polyamines. Both isoforms of antizyme (AZ1 and AZ2) are expressed by  $+1$  frameshifting that is stimulated by polyamines. RF-2 and antizyme frameshifting is involved in the specific autoregulation ([4](#), [5](#)).

A majority of  $-1$  frameshift sites are characteristic heptamers with a **consensus sequence** of X XXW WWY (presented as a preshift reading frame; X and Y are one of any nucleotides and W is either A or U). In addition, most of these  $-1$  shift sites are accompanied by either pseudoknot or stem-loop

structures 3' to the frameshift site. A simultaneous shift model has been proposed for the shift mechanism, in which the ribosome and the two tRNAs in the P and A sites slip back on the mRNA by one nucleotide, forming new base pairing between the tRNAs and the mRNA. Downstream RNA structures stimulate the process by, at least in part, causing ribosomal pausing at the frameshift site (Fig. 1).

Analyses of +1 frameshifting have revealed that their mechanisms are not as uniform as -1 frameshifting. The codon 3' adjacent to the frameshift site is either a termination codon or a "hungry" codon [a codon decoded by low-abundant tRNA and thus having low efficiency (6)]. In at least two cases, tRNA slipping and pairing at the new position are not necessary.

Frameshifting can be hinted from open reading frame analysis of [complementary DNA](#) sequences or by a disagreement between the sequences of cDNA and its protein product. It is, however, sometimes difficult to prove before knowing the amino acid sequence of the protein. Frameshifting can be induced by expression of exotopic genes or artificial sequences. It is noteworthy that translational frameshifting, even with a low efficiency, can partially rescue [frameshift mutations](#) and modifies the phenotype in both experimental and naturally occurring systems, demonstrating the biological importance of frameshifting.

### Bibliography

1. F. W. van Leeuwen et al. (1998) *Science* **279**, 242–247.
2. J. F. Atkins and R. F. Gesteland (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society of Microbiology Press, Washington, D.C., pp. 471–490.
3. E. J. Murgola (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society of Microbiology Press, Washington, D.C., pp. 491–509.
4. W. J. Craigen, R. G. Cook, W. P. Tate, and C. T. Caskey (1985) *Proc. Natl. Acad. Sci. USA* **82**, 3616–3620.
5. S. Matsufuji et al. (1995) *Cell* **80**, 51–60.
6. K. Kawakami et al. (1993) *Genetics* **135**, 309–320.

### Suggestion for Further Reading

7. R. F. Gesteland and J. F. Atkins (1996) *Annu. Rev. Biochem.* **65**, 741–768.

### Free Energy Calculations

#### 1. Statistical Mechanical Averaging

With a given set of [potential functions](#), one can evaluate various average properties of the system by **computer simulation**. In particular, it is useful to simulate experimentally observed macroscopic properties using microscopic models. This can be done using the theory of statistical mechanic, which tells us that the average of a given property,  $A$  (which is independent of the momentum of the system), is given by

$$\langle \mathbf{A} \rangle = \int \mathbf{A}(\mathbf{r}) P(\mathbf{r}) d\mathbf{r} = \int \mathbf{A}(\mathbf{r}) \exp\{-U(\mathbf{r})\beta\} d\mathbf{r} / z(U)$$

$$= \int \exp\{-U(\mathbf{r})\beta\} d\mathbf{r}$$

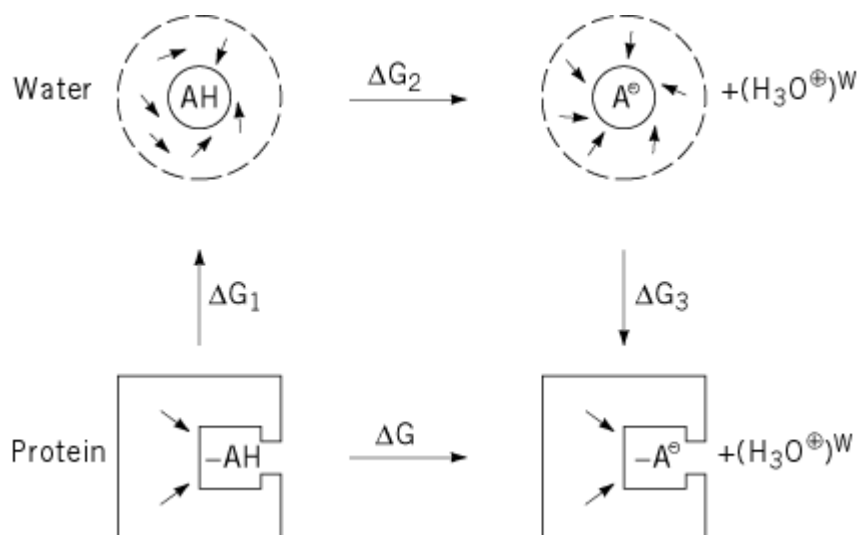
where  $b = (1/k_B T)$  ( $k_B$  is the Boltzmann constant and  $T$  is the absolute temperature), and  $d\mathbf{r}$  designates the volume element of the complete space spanned by the  $3n$  vector  $\mathbf{r}$  associated with the  $n$  atoms of the system (1). The evaluation of equation (1) requires us to explore all points in the entire conformation space of the given system. Such a study of a solvated macromolecule is clearly impossible with any of the available computers. However, one can hope that the average over a limited number of configurations will give similar results to those obtained from an average over the entire space. With this working hypothesis, we can try to look for an efficient way of spanning phase space. Present-day computers can be used to span a significant number of conformations by **Monte Carlo (MC) calculations** or by **molecular dynamics (MD) methods**.

Calculations of average properties may converge extremely slowly. This is the case, for example, for thermodynamic properties such as **enthalpy** and **entropy**. Fortunately, **free energies** converge faster than enthalpies or entropies, due to entropy-enthalpy compensation. This is quite useful because most biological processes depend on the relevant free energies and not on their individual components. We describe below the general methods for using simulations to obtain free energies.

## 2. Thermodynamic Cycles

Many biochemical processes can be described by *thermodynamic cycles* of the same type as used in chemistry (eg, the Born–Haber cycle (2, 3)). For example, the energetics associated with the deprotonation of an ionizable group in a protein can be described by the cycle as shown in Figure 1 (4). Using such a cycle allows one to replace calculation of the physical process of moving the proton from the protein to water by the much simpler comparative calculation of “mutating” the given acid from its unionized form in the protein and in water. The use of thermodynamic cycles in microscopic calculations of free energies of biological processes has been introduced in the early 1980s (4, 5), and it is now a major part of many approaches of free energy calculations of biological processes.

**Figure 1.** The thermodynamic cycle used to estimate the energetics of dissociation of an acidic group of a protein. The  $\Delta G_i$  are given by  $\Delta G_1 = G_{sol}^w(\text{AH}) - G_{sol}^p(\text{AH})$ ,  $\Delta G_2 = 2.3RT(\text{p}K_a^w - \text{pH})$ , and  $\Delta G_3 = G_{sol}^p(\text{A}^-) - G_{sol}^w(\text{A}^-)$ , where  $p$  and  $w$  designate protein and water, respectively.



### 3. Free Energy Perturbation and Related Approaches

Evaluation of free energies by statistical mechanical approaches is extremely time-consuming, due to the need to sample all the relevant configurational space. Fortunately, it is possible in some cases to obtain meaningful results using perturbation approaches. Such approaches exploit the fact that many important properties depend on local changes in the macromolecules so that the effect of the overall macromolecular potential cancels out. Such calculations are usually done by the so-called **free-energy perturbation (FEP) method** (6, 7) [also related to the umbrella sampling method (7)]. This method evaluates the free energy associated with the change of the potential surface from  $U_1$  to  $U_2$  by gradually changing the potential surface using the relationship

$$U_m(\lambda_m) = U_1(1 - \lambda_m) + U_2\lambda_m \quad (2)$$

The free-energy increment  $\delta G(1 - \lambda_m \rightarrow \lambda'_m)$  associated with the change of  $U_m$  to  $U_{m'}$ , can be obtained (7) by

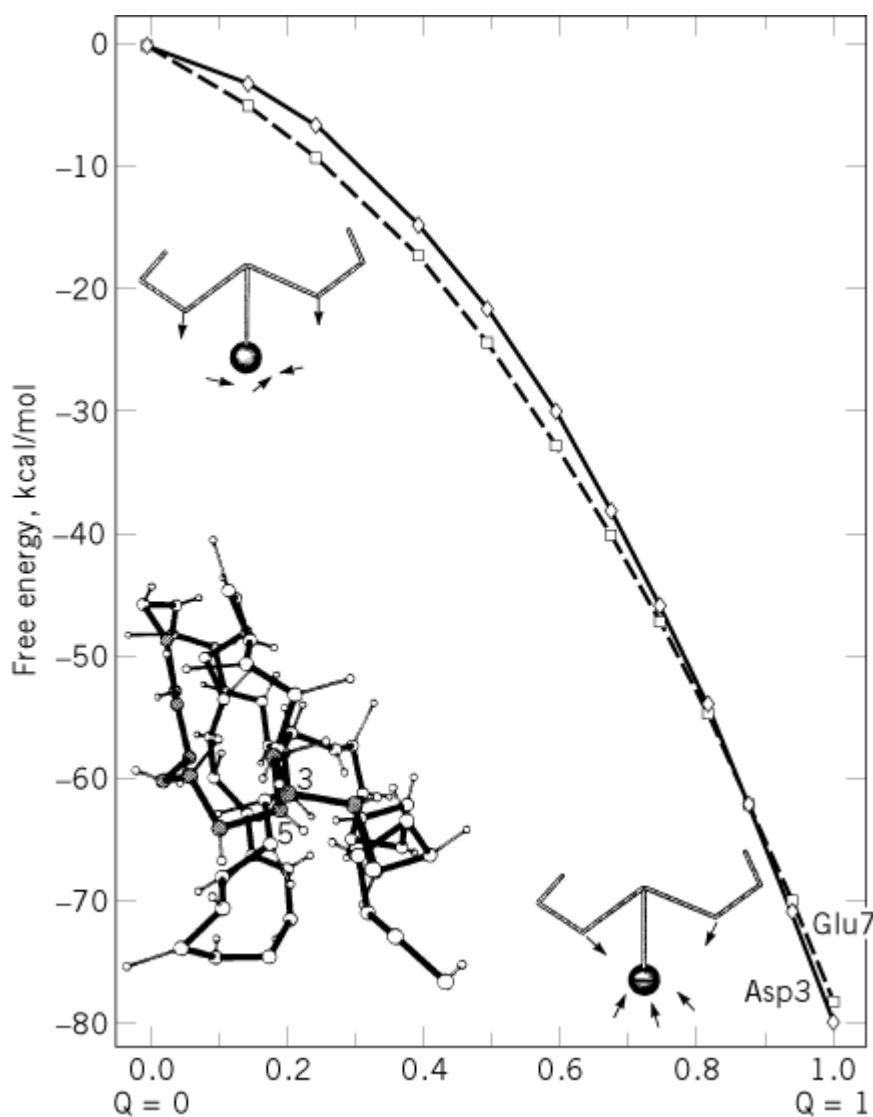
$$\exp\{-\delta G(\lambda_m \rightarrow \lambda_{m'})\beta\} = \langle \exp\{-(U_{m'} - U_m)\}\beta \rangle_m \quad (3)$$

where  $\langle \rangle_m$  indicates that the given average is evaluated by propagating trajectories over  $U_m$ . The overall free energy change is now obtained by changing the  $\lambda_m$  in  $n$  equal increments and evaluating the sum of the corresponding  $dG$ :

$$\Delta G(U_1 \rightarrow U_2) = \sum_{m=0}^{n-1} \delta G(\lambda_m \rightarrow \lambda_{m+1}) \quad (4)$$

An example of the use of Eq. (4) in calculating free energy of a biological system is given in Figure 2, which considers the solvation free energy of ionizable residues in BPTI.

**Figure 2.** Calculations of the free energy of charging Asp 3 and Glu 7 in BPTI. The figure describes the calculated “solvation” free energy obtained with the use of equation (4) as a function of the charge of the corresponding acid.



The FEP approach has been used extensively in studies of free energies of biological systems (eg, Refs. [8](#) and [9](#)) and more details are given in **Computer simulation**. It must be emphasized here that the convergence of FEP approaches is quite slow and that obtaining meaningful results requires proper treatment of long-range effects.

Many times it is very useful to estimate the free energy of biological processes by the so-called *linear response approximation* (LRA) ([10](#)). This approximation assumes that the protein and solvent environment respond linearly to the forces associated with the given process. This leads ([11](#)) to the equation

$$\Delta G_{a \rightarrow b} = \frac{1}{2}(\langle U_b - U_a \rangle_a + \langle U_b - U_a \rangle_b) \quad (5)$$

Such an approximation is, for example, the basis of macroscopic electrostatic theory, where the free energy of charging a positive ion, with a charge  $Q = 1$ , is given by the well-known result  $\Delta G = \frac{1}{2}U(Q = 1)$ . Although it is hard to accept that the LRA provides a reliable way of describing the energetics of macromolecules or of realistic molecular systems, it was found by simulation studies that this is a reasonable approximation in particular for processes that depend on electrostatic effects ([12-14](#)).



## Bibliography

1. D. A. McQuarrie (1976) *Statistical Mechanics*, Harper and Row, New York.
2. M. Born (1919) Verh. Dtsch. phys. Ges. **21**, 13.
3. F. Haber (1919) Verh. Dtsch. phys. Ges. **21**, 750.
4. A. Warshel (1981) *Biochemistry* **20**, 3167.
5. C. F. Wong and J. A. McCammon (1986) *J. Am. Chem. Soc.* **108**, 3830.
6. R. W. Zwanzig (1954) *J. Chem. Phys.* **22**, 1420.
7. J. P. Valleau and G. M. Torrie (1977) *Modern Theoretical Chemistry*, Vol. **5**, Plenum Press, New York.
8. P. Kollman (1993) *Chem. Rev.* **93**, 2395.
9. A. Warshel and J. Åqvist (1991) *Annu. Rev. Biophys. Biophys. Chem.* **20**, 267.
10. R. Kubo, M. Toda, and N. Hashitsume (1985) *Statistical Physics II: Nonequilibrium Statistical Mechanics*, Springer-Verlag, Berlin.
11. F. S. Lee, Z. T. Chu, M. B. Bolger, and A. Warshel (1992) *Protein Eng.* **5**, 215.
12. J.-K. Hwang and A. Warshel (1987) *J. Am. Chem. Soc.* **109**, 715.
13. R. A. Kuharski, J. S. Bader, D. Chandler, M. Sprik, M. L. Klein, and R. W. Impey (1988) *J. Chem. Phys.* **89**, 3248.
14. J. Åqvist and T. Hansson (1996) *J. Phys. Chem.* **100**, 9512.

## Free Energy Relationships

Understanding any chemical reaction requires knowledge about the [transition state](#) for the reaction, because the free energy of the transition state determines the rate at which the reaction occurs (see [Kinetics](#)). The transition state cannot be observed directly, but its properties can be inferred from the effects of changes in various parameters on its free energy, the rate of the reaction. One method of doing so is to study the chemical reaction in a series of closely related reactants and to correlate the variation in the rates of reaction with the physical parameters of the reactants. Such observations are referred to as free energy relationships.

In hydrolytic reactions, the reacting groups are generally classified as nucleophiles or electrophiles. Nucleophiles donate an electron pair, and are generally bases, whereas electrophiles have opposite properties and generally are acids. Protonation converts a nucleophile to an electrophile. For example, an  $\text{-NH}^2$  group is a nucleophile, whereas its protonated  $\text{NH}_3^+$  form is an electrophile. The strengths of their nucleophilicity and electrophilicity can be indicated by their respective tendencies to acquire or donate protons and can be measured by the  $\text{p}K_a$  value (see [Buffers](#)). The reaction of a nucleophile with a proton is similar to the reaction of the nucleophile with any other group, so that a relationship would be expected between the  $\text{p}K_a$  of a nucleophile and its reactivity in a chemical reaction. The  $\text{p}K_a$  is a measure of the relative free energy of the nucleophile, whereas the logarithm of the rate constant for a reaction is a measure of the free energy of the [transition state](#) for the reaction. Such relationships are generally evident from plots of the logarithm of the rate constant  $k$ , for a chemical reaction, versus the  $\text{p}K_a$  values of a series of related nucleophiles, known as Brønsted plots. Such plots can be made for a series of nucleophiles attacking a particular ester or amide bond, for example, or for one nucleophile attacking a series of molecules that have leaving groups with

different  $pK_a$  values. Linear correlations are found frequently for a given chemical reaction and a series of related reactants, so long as the reactants are sufficiently related and no special factors, such as steric hindrance, affect their rate of reaction. The general form of the linear relationship is expressed as

$$\log k = \beta pK_a + \text{constant} \quad (1)$$

When the reactant is an electrophile, the equation is the same, except that the parameter  $\beta$  is changed to  $\alpha$ . The values of  $\alpha$  are generally negative, whereas those of  $\beta$  are generally positive.

The slope of a linear Brønsted plot, the value of  $\alpha$  or  $\beta$ , provides important information about the transition state of the reaction. Their sign and magnitude are generally interpreted as an indication of the amount of charge developed in the transition state and whether the transition state is closer to the reactants or the products in its physical properties. The values of  $\alpha$  and  $\beta$  should be between 0 and 1, because complete transfer of a proton in the transition state gives a value of 1.0 and no transfer gives a value of 0. For ester hydrolysis, the values are usually 0.3 to 0.6. For example, the attack of tertiary amines on esters with poor leaving groups (e.g., alkoxides) gives values of +1.5 for the series of attacking nucleophiles and -1.5 for variation in the alcohol leaving group (1). These very large values indicate that the transition state for this reaction is close to the products. In contrast, the reaction of very basic nucleophiles with activated esters (e.g., phenolates) has relatively small values, +0.1 to 0.2 for variation of the nucleophile and -0.1 to -0.2 for variation of the leaving groups; this indicates that the transition state in this case is closer to the ester substrate.

Although free energy relationships were developed with reactions between small molecules, they are also useful in molecular biology in phenomena as diverse as **protein folding** (2).

#### Bibliography

1. A. Fersht (1985) *Enzyme Structure and Mechanism*, 2nd ed., W. H. Freeman, New York, pp. 147–154.
2. J. M. Matthews and A. R. Fersht (1995) *Biochemistry* **34**, 6805–6814.

#### Suggestion for Further Reading

3. S. Bernhard (1968) *The Structure and Function of Enzymes*, W. A. Benjamin, New York, pp. 175–183.

### Freeze Fracture

*Freeze fracture* is a preparative procedure for visualizing the distribution of **proteins** within the **lipid** bilayer of **membranes** by transmission **electron microscopy**. The membrane samples are rapidly frozen to a very low temperature and then fractured with a knife blade. Fracturing occurs along planes where the intramolecular forces are weakest. In membranes, these planes occur along the center of the phospholipid bilayer, where the apposing fatty acid chains of the monolayer leaflets interact. Freeze fracturing thus exposes the **nonpolar** interior of the membrane. Integral membrane proteins that cross the bilayer interrupt the fracturing process and are carried along with either monolayer leaflet, leaving pits in the apposed monolayer. The fracture faces are shadowed with metal, and the proteins appear as small particles. The technique may be combined with **immunolabeling** to visualize the location of specific proteins (1).

*Freeze etching* utilizes freeze-drying of the fractured frozen samples to remove water, thereby exposing the components of the surfaces of the fractured samples; the surfaces are then shadowed for viewing in the electron microscope.

#### Bibliography

1. P. Pinto da Silva (1984) In *Immunolabeling for Electron Microscopy* (J. M. Polack and I. M. Varndell, eds.), Elsevier, Amsterdam, pp. 179–188.

#### Suggestions for Further Reading

2. N. J. Severs, and D. M. Shotton (1995) *Rapid Freezing, Freeze Fracture, and Deep Etching*, Wiley-Liss, New York.
3. R. L. Roberts, R. G. Kessel, and H. Tung (1991) *Freeze Fracture Images of Cells and Tissues*, Oxford University Press, New York.

#### Frictional Coefficient, Ratio

The translational frictional coefficient,  $f$ , is a measure of the resistance to movement of a molecule; this resistance is a function of both the size and the shape of the molecule. It can be measured experimentally either by its rate of [diffusion](#) or its rate of sedimentation.

As would be expected, the translational frictional coefficient is inversely proportional to the diffusion coefficient,  $D_{20,w}^{\circ}$ , under standard conditions of water at 20.0°C and extrapolated to zero concentration:

$$f = \frac{RT}{N_A D_{20,w}^{\circ}} = \frac{k_B T}{D_{20,w}^{\circ}} \quad (1)$$

where  $R$  is the gas constant ( $8.314 \times 10^7 \text{ erg mol}^{-1} \text{ K}^{-1}$ ) and  $N_A$  is Avogadro's number ( $6.022 \times 10^{23} \text{ mol}^{-1}$ ).

Likewise, the frictional coefficient is inversely proportional to the [sedimentation coefficient](#),  $s_{20,w}^{\circ}$ , under the same standard conditions:

$$f = \frac{M_r(1 - \bar{v}\rho_{20w})}{N_A s_{20,w}^{\circ}} \quad (2)$$

where  $\bar{v}$  is the [partial specific volume](#) of the macromolecule and  $\rho_{20,w}$  is the density of the standard solvent, water at 20°C.

It is usually more informative to use the **frictional ratio**,  $f/f_0$ , which is the dimensionless ratio of the observed translational frictional coefficient to that of an equivalent spherical molecule of the same anhydrous mass and density. This corrects for the size of the molecule and gives an indication of its shape. The greater the frictional ratio, the more asymmetric it is likely to be. The frictional ratio can be calculated from the diffusion coefficient by the following equation:

$$\frac{f}{f_0} = \frac{k_B T}{6\pi\eta_{20,w}} \cdot \left( \frac{4\pi N_A}{3\bar{v}M} \right)^{1/3} \frac{1}{D_{20,w}^0} \quad (3)$$

The frictional ratio depends intrinsically on the conformation, flexibility, and degree of solvation (by water, salt ions, and any other solvent molecules) of the macromolecule. This degree of water association is termed the [hydration](#) of the macromolecule,  $d$ , and is defined as the mass in grams of associated solvent per gram of anhydrous biomolecule. This associated solvent includes both chemically bound solvent and also solvent physically entrained in the interstices in the molecule. The value of  $d$  is typically between 0.2 and 0.5 g/g for proteins, although it is a notoriously difficult parameter to pin down with any accuracy.

The function defining the shape and flexibility of the biomolecule is the *Perrin translational frictional function*,  $P$ :

$$P = \frac{f}{f_0} \cdot \left( \frac{\bar{v}}{\bar{v} + \delta/\rho_0} \right)^{1/3} \quad (4)$$

where  $\rho_0$  is the density (in grams per milliliter) of the bound solvent. For a molecule that is fairly rigid on a time-averaged basis, the gross conformation can be specified using  $P$  in terms of the axial ratio of the equivalent hydrodynamic ellipsoid or in terms of sophisticated arrangements of spheres called *hydrodynamic bead models* (see [Diffusion](#)).

#### Suggestions for Further Reading

S. E. Harding (1995) On the hydrodynamic characterisation of macromolecular conformation. *Biophys. Chem.* **55**, 69–93.

K. E. Van Holde (1985) *Physical Biochemistry*, Prentice-Hall, Englewood Cliffs, N. J. (Chapter 4 gives an excellent introduction).

## Fungi, Filamentous

The fungi represent one of the most diverse groups of organisms encountered in the natural world; approximately 70,000 different species have been described although estimates as high as 1.5 million fungal species are likely to exist. Even though a large fraction of fungi are believed to be unculturable, their importance to the natural environment of the planet cannot be underestimated. Thus, a clear understanding of their physiology at a molecular level seems prudent if not essential. The past 10 years have seen an explosion of information in this field, due primarily to the ability to transform a wide variety of these previously intractable organisms. Yet, even with the exponential advances in recombinant techniques (especially sequencing technologies) gained over the past several years, filamentous fungi appear to be the forgotten model system; perhaps it is because the completion of the first eukaryotic genome, *Saccharomyces cerevisiae*, was a fungus, albeit a yeast, or perhaps it is because of the success of the human genome project that the perception is that we need not look at additional fungi as model organisms. However, while *Saccharomyces* provides an excellent framework for studying eukaryotic systems, it clearly lacks some critical parameters that filamentous fungi can elucidate, in particular, multicellular differentiation and pathogenicity. In addition, it should be kept in mind that a large percentage of the expressed sequence tags (ESTs) identified from filamentous fungi do not have detectable homology to sequences already present in

any of the available databases, including the *Saccharomyces* genome). Therefore, the filamentous fungi as a group provide a complex and somewhat unique perspective of the microbial world that should not be overlooked.

Two of the most well known filamentous fungi, *Aspergillus nidulans* and *Neurospora crassa*, are exceptional model organisms. They have a long genetic history and are continuing to be pursued as excellent subjects in which to study development, gene regulation, circadian rhythms, and genome structure. As the functions and interactions of proteins and intracellular molecules is elucidated from a variety of organisms, including the fungi, humans appear to have more in common with their distant eukaryotic ancestors than previously thought. It is important to keep in mind that the filamentous fungi are part of the group of organisms that include the yeasts, *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*, two of the best understood eukaryotic organisms from a genetic and molecular standpoint. In general, the molecular biology of the entire kingdom may be viewed as more similar than dissimilar; not only are there are two-component signal transduction systems, transcriptional inducers and repressors, but the conservation of protein sequence is remarkable in many areas. Wherever commonalities have been searched for they have been found, with a few exceptions that may yet prove to be due to insufficient data.

## 1. Why Study Fungi?

Fungi are a biochemically diverse set of organisms, producing a wide variety of acids and degradative enzymes that support their absorptive lifestyle and play an essential role in the degradative processes of our entire ecosystem. They decompose numerous substrates by secreting a wide variety of very efficient degradative enzymes, thus providing an efficient means to recycle much of the earth's biomass. In addition, their fermentative capabilities provide a means to bake bread and ferment a wide variety of fruits and vegetables. The fungi have also been exploited commercially as efficient producers of enzymes such as glucoamylase, lipase, cellulase, and pectinase. In addition, fungal secondary products have been used in the synthesis of numerous pharmaceuticals, including [penicillin](#), cyclosporin, and mevalonin.

On the other hand, some fungi can have a devastating effect: their uncontrolled growth can irrevocably damage entire forests, crops, or agricultural feed products. Some of the most intensely studied fungi are phytopathogenic; *Magnaporthe grisea* (pathogen of rice), *Ustilago maydis* (corn), *Cryphonectria parasitica* (chestnut blight), etc. Aflatoxin, a single secondary metabolite produced by *Aspergillus flavus*, causes millions of dollar of damage to the corn industry annually. In terms of human disease, medically important species of fungi are increasing in prevalence due to a dramatic rise in the incidence of immunocompromised patients. The most significant human pathogen is *C. albicans*, a dimorphic organism that switches between a yeast-like growth phase not unlike that of its distant cousin, *S. cerevisiae*, to a hyphal growth phase that has been associated with its virulence *in vivo*. Other well-known fungal pathogens include *Aspergillus fumigatus*, *Blastomyces dermatitidis*, *Coccidioides immitis*, *Cryptococcus neoformans*, *Histoplasma capsulatum*, and *Pneumocystis carinii*. Many are associated with high mortality rates. All of these fungi are either filamentous or dimorphic in nature. A wide variety of classical genetic and modern molecular techniques are presently available to analyze these organisms, as well as the lesser well-known, but certainly not less virulent pathogens. What is important to keep in mind is the similarities between the fungi, the “big picture.” Each of the individual organisms’ details can be gleaned from further readings.

## 2. Physiology and Taxonomy

The taxonomy of the fungi is complex, and students interested in further classification or taxonomy of fungi are referred to Ref. [1](#). Over the past several decades, their special place in the taxonomy of biological organisms has been reclassified from a subdivision within the plant kingdom to part of the eukaryotic world. Within the new scheme [bacteria, archea and eukarya ([2](#))], note that the division between plants and fungi has disappeared. Even with this new assignment, there have been many proposals for their classification based upon nutritional requirements, cytology, or morphogenesis ([1](#),

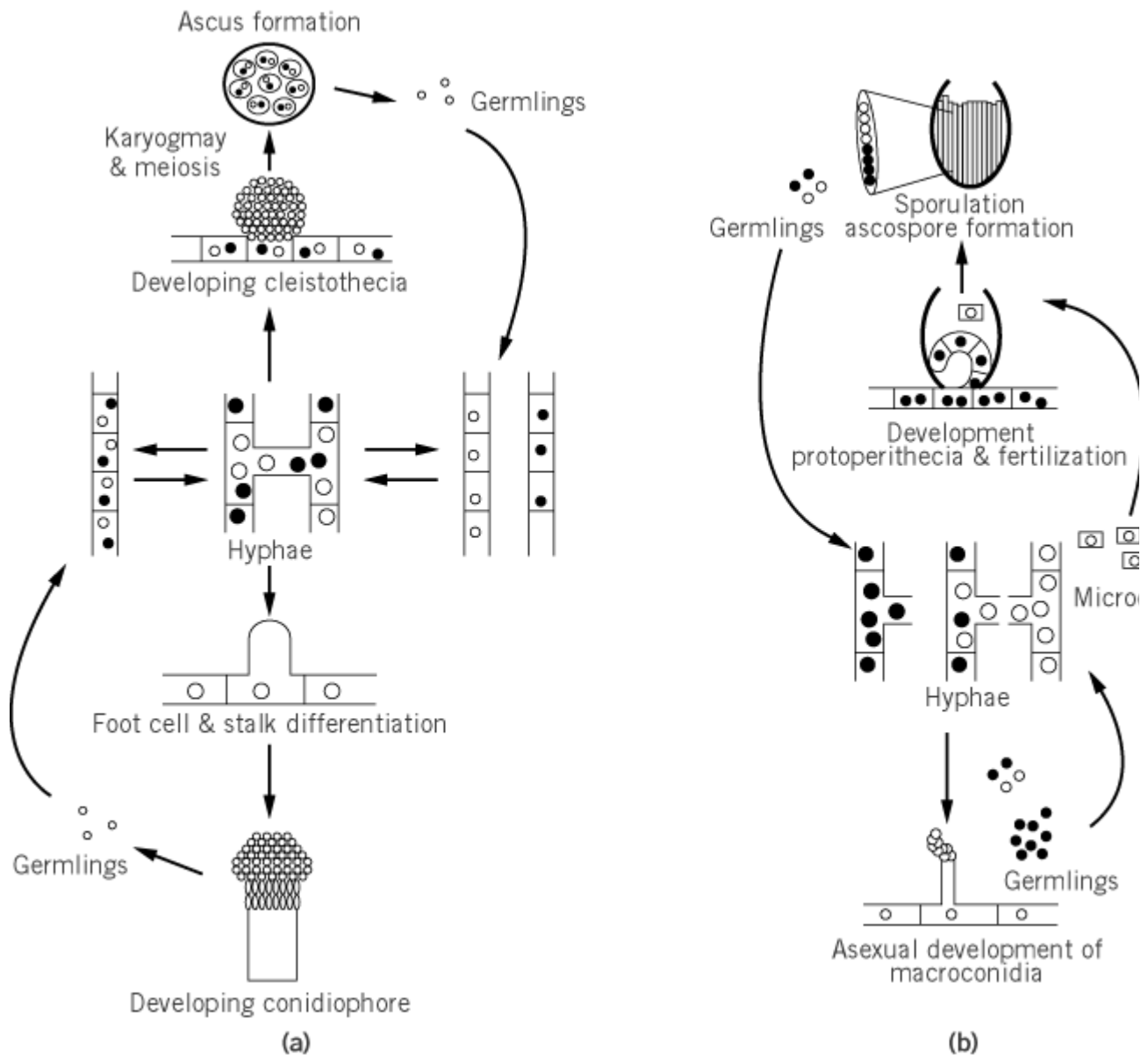
79). Perhaps the most enlightening description of the taxonomical state of the field is described by Ainsworth, “the question of the number of kinds of fungi can be more usefully approached by first considering some of the main factors which determine their recognition. ¼ Their [fungi] recorded number and distribution is closely correlated with the number and distribution of mycologists. ¼ Some taxonomists are perhaps unduly impressed by small differences ¼ others, are lumpers ¼ recognize far fewer kinds”. By and large the fungi are regarded as having five phyla: oomycetes (water molds), zygomycetes (common molds), ascomycetes (sac fungi), basidiomycetes (club fungi) and deuteromycetes (imperfect fungi having no sexual stage). With the advent of [recombinant DNA](#) technology and new [DNA sequencing](#) paradigms, it will be possible to re-evaluate the relationships between certain fungi by their sequence **homologies**, rather than by their morphology or nutritional biochemistry. Most of the research has been centered upon a few species of the ascomycetes and basidiomycetes although *C. albicans*, a member of the fungi imperfecti, has been the source of a great deal of research due to its critical role in medical mycology.

Like most of their taxonomy, fungal life cycles may vary from simple (as in the yeasts) to complex (as in the basidiomycetes). Fungi are typically non-photosynthetic heterotrophic organisms that absorb nutrients from their surrounding environment. As multicellular organisms, they display a wide range of morphological types from the microscopic yeasts and water molds to the multicellular mushrooms. They may exist as a multinucleate undifferentiated mass of cells during all or part of their life cycle, or form complex highly developed reproduction structures. In their vegetative state, fungal cells may be unicellular (yeast-like) or filamentous (mycelial or hyphal-like) in shape. Many are capable of switching between the two forms, which has been the source of a great deal of interest, especially with regard to the regulatory circuits involved and the effect on pathogenicity. All fungi are eukaryotic, but the cell may be uninucleate or multinucleate, haploid, diploid or dikaryotic during part of their life cycle. They typically contain a cell wall, which may explain their original classification within the kingdom of plants; however, the cell wall contains polymers of chitin and glucans rather than the more usual cellulose of plants.

### 3. Fungal Life Cycles and Reproduction

Fungi display a variety of reproductive choices. They may reproduce either asexually or sexually, or both. The fungal spore is specialized for reproduction, survival and, in some cases, dispersal. Not only are there heterothallic species such as *Neurospora crassa* that require different mating types in order to undergo sexual reproduction, but homothallic species (self-maters) such as *Aspergillus nidulans* and pseudohomothallic species that switch mating types (such as *Saccharomyces cerevisiae*) as well. Further, there is an assortment of fungi that do not appear to undergo sexual reproduction at all (e.g. -fungi imperfecti, including *Candida albicans*). The life cycles of *Aspergillus nidulans* and *Neurospora crassa* (Fig. 1) serve to illustrate many of the salient features found in this highly divergent group of organisms.

**Figure 1.** Schematic representation of the life cycles of *Aspergillus nidulans* and *Neurospora crassa*. **(a)** *Aspergillus nidulans*. Asexual reproduction involves formation and maturation of conidiophores (described in the text) and uninucleate spores called conidia. Vegetative haploids grow as monokaryotic or dikaryotic hyphae. These cells may then undergo sexual reproduction (upper portion of the cycle), initiated with the formation of specialized fruiting body called the cleistothecia. Eight binucleate ascospores are formed within an unordered ascus. **(b)** *Neurospora crassa*. Asexual reproduction involves the formation of conidiophores, producing millions of multinucleate spores called macroconidia. Vegetative haploids grow as monokaryotic or dikaryotic hyphae. Cells of the opposite mating type are required for sexual reproduction. Sporulation is initiated by the formation of the protoperithecia, which is fertilized by microconidia of the opposite mating type (see text). Cellular fusion is followed by a complex developmental pathway that involves synchronous nuclear division, karyogamy, and meiosis to yield hundreds of asci within the perithecium (see text for details).



### 3.1. Life Cycle of *Neurospora crassa* (3, 4)

*N. crassa* is a haploid ascomycete found in two non-switching mating types, *A* and *a* (referred to as “big A” and “little a”). Strains of either type grow as branching, thread-like cells called hyphae. These hyphal cells are typically interconnected and multinucleated. The mycelium is homokaryotic if all the nuclei present have the same genotype, and it is heterokaryotic if they have a different genetic composition. The fusion of haploids to produce diploid nuclei is rare; thus **heterokaryons** are naturally stable.

Certain stresses (a dry environment or carbon starvation) activates the asexual developmental pathway, which is also called conidiation. Aerial hyphae grow away from the mycelial substrate and develop a spore-producing structure that contains conidiophores. Conidiophores are chains of cells that grow by budding and develop into multinucleated spores called macroconidia. These bright orange spores are an ideal means of dispersal, as they are small (5–10  $\mu\text{m}$ ) and abundant.

In contrast, sexual sporulation has rather exacting nutritional requirements. It begins with the formation of protoperithecia—the female reproductive apparatus of the organism. Poorly characterized metabolic and environmental signals induce protoperithecial development in the

absence of strains of the opposite mating type. In fact, more is known about the conditions that inhibit development than those that activate it. For example, protoperithecia are sparsely formed on nutritionally complex media, on submerged liquid cultures, and at temperatures above 27°C. The absence of ammonium ions seems to be critical to induce their formation, as does a marginal degree of nitrogen starvation or intense carbon deprivation. In addition, blue light is known to stimulate their formation.

Development of microconidia occurs in the absence of strains of the opposite mating type. These asexual spores are small and uninucleate, grayish-brown in color. They are formed within the vegetative hyphae and rupture through the cell wall as they mature. Microconidia are thought to function primarily as male fertilizing agents for the developing protoperithecia.

Being heterothallic, *Neurospora* must mate with strains of the opposite mating type to form sexual spores or ascospores. For this, trychogynes, specialized hyphal elements that emanate from protoperithecia, are fertilized by male elements (macroconidia, microconidia, or vegetative cells) of the opposite mating type. Fusion, or plasmogamy, initiates the development of the perithecium, the multicellular sexual apparatus. After plasmogamy, the male- and the female-derived nuclei coexist in a heterokaryotic tissue and divide mitotically until they are sorted into dikaryotic tissue, in which each cell compartment contains only one nucleus of each mating type. The nuclei then pair and undergo a series of synchronous mitoses until the tip of the hyphal cell in which they reside bends to form a hook-shaped cell called a *crozier*. There, the two nuclei undergo a coordinate mitosis yielding, after septum formation, a uninucleate basal cell, a uninucleated lateral cell, and a penultimate ascus mother cell that contains one nucleus of each mating type. It is in this cell that karyogamy, meiosis, and postmeiotic mitoses take place resulting in an ascus containing eight haploid spores in an order that reflects their lineage. About 200 such asci are generated within a typical perithecium. Once mature, the ascospores are ejected from the ascus through the ostiole in the perighecial beak. Meiotic segregation and recombination can be studied in *Neurospora* by analyzing individual asci or random ejected spores.

### 3.2. Life Cycle of *Aspergillus nidulans* (5-7)

*A. nidulans* is a homothallic ascomycete that does not require different mating types for sexual reproduction. As a vegetative mycelium *A. nidulans* grows via apical extension in which genetically distinct nuclei can coexist in a single thallus as heterokaryons or separately, maintaining the homokaryotic state. Although the hyphae have what appear to be individual segments, the cross walls contain pores so that the mycelium is in fact a multinucleate structure. Within the growing mycelia mat, nuclei are also capable of fusion, producing diploid thalli. This is also a relatively rare event in naturally occurring populations, but may be accomplished in the laboratory with the use of complementing auxotrophies present in each of the donor hyphae (8).

Sexual reproduction involves the formation of specialized fruiting bodies called cleistothecia containing fertile hyphae. Since *A. nidulans* is homothallic, mating occurs within a colony of uninucleate cells, much like wild type strains of *S. cerevisiae*. Nuclear fusion (karyogamy) and meiosis occur within the cleistothecia followed by the formation of ascospores within individual asci. Further differentiation of the postmeiotic nuclei involves two mitotic divisions producing binucleate ascospore progeny.

Asexual reproduction in *A. nidulans* involves the formation of multicellular conidiophores and uninucleate conidia. In response to signals that are as yet poorly understood, a single hyphal segment from the vegetative mycelium develops into a foot cell, which differentiates a specialized stalk and metulae. Repeated divisions of the metulae give rise to uninucleate phialides cells which undergo additional rounds of division to produce conidia. Thus, conidia are uninucleate and haploid, whereas ascospores are binucleate. A great deal of research has been applied to understanding the genetic and biochemical mechanisms that govern conidial development in *A. nidulans*. See Ref. 9 for review.

Although both *Aspergillus nidulans* and *N. crassa* have long genetic systems that make them ideal



for study of a wide variety of eukaryotic problems, the use of recombinant DNA techniques and transformation systems in all of the fungi has enlarged the wealth of information that can be gleaned from pursuing some of the more genetically intractable organisms.

#### 4. Genome Structure

The **genomes** of fungi range in size from around 7-8 Mb for *Pneumocystis carinii* and *Ashbya gossypii* to near 40 Mbp for some of the basidiomycetes (see Table 1). The presence of introns may be rare, as for *C. albicans* and *A. gossypii*, or they may be widely distributed, as for *A. nidulans*. It is estimated that the filamentous fungi will contain anywhere from 4500 genes in the smallest fungi to 9,000–12,000 genes in the largest genomes (Phillipsen, personal communication) (10). Because most fungal chromosomes do not condense during mitosis, their karyotypes cannot be determined by cytological methods. A size estimate may be determined utilizing [pulsed-field gel electrophoresis](#) techniques in which whole chromosomes are resolved using alternating electrical fields. Some examples are given in Table 1.

**Table 1. Genome Size and Karyotype**

| Organism                           | Size (Mb) | Number of chromosomes    | References   |
|------------------------------------|-----------|--------------------------|--|
| <i>Ashbya gossypii</i>             | 8.85      | 7                        | <a href="#">22</a> , <a href="#">54</a>  |
| <i>Aspergillus nidulans</i>        | 31        | 8                        | <a href="#">55</a> , <a href="#">56</a>  |
| <i>Candida albicans</i>            | 16–17     | 8–9                      | <a href="#">57</a> , <a href="#">58</a>  |
| <i>Coprinus cinereus</i>           | 37.5      | 13                       | <a href="#">59</a>   |
| <i>Cochliobolus heterostrophus</i> | 35        | 15–16                    | <a href="#">13</a>   |
| <i>Histoplasma capsulatum</i>      | ~31       | 7                        | <a href="#">60</a>   |
| <i>Magnaporthe grisea</i>          | 40        | 7 + 1–4 mini chromosomes | <a href="#">61-63</a>  |
| <i>Neurospora crassa</i>           | 42.9      | 7                        | <a href="#">64</a> , <a href="#">65</a>  |
| <i>Pneumocystis carinii</i>        | 7–8       | 14–16                    | <a href="#">66</a>   |
| <i>Saccharomyces cerevisiae</i>    | 13.5      | 16                       | SGD: <a href="http://genome-www.stanford.edu/Saccharomyces/">http://genome-www.stanford.edu/Saccharomyces/</a> |
| <i>Schizosaccharomyces pombe</i>   | 14        | 3                        | <a href="#">67</a>   |
| <i>Ustilago hordei</i>             | 18–25     | 16–21                    | <a href="#">68</a>   |
| <i>Ustilago maydis</i>             | 22        | ~20                      | <a href="#">69</a>   |

The genomes of fungi are known to be somewhat plastic, as evidenced by the chromosome length polymorphisms (CLPs) observed for many different species using pulsed-field electrophoretic techniques to separate whole chromosomes by size. This polymorphism can be accounted for by both mitotic and meiotic events and is observed in both sexual and asexual fungi [reviewed in (11)].

The mechanisms by which these changes are initiated are largely unknown, but the presence of repetitive elements within the genome has been suggested to play a critical role. Reciprocal translocations have been demonstrated in two isolates of *C. albicans* (12) and in *C. heterostrophus* (13), and the expansion/contraction of tandem ribosomal DNA repeats has been documented as the cause of some CLPs in *C. albicans* (14), *Coprinus cinereus* (15) and *Cladosporium fulvum* (16). Subtelomeric regions of fungi are known to contain repetitive DNA sequences that have proven useful in fingerprinting individual isolates. Until recently, fungi were thought to lack [transposon](#) elements, which have been shown to influence genome fluidity in other systems. They have been identified in a number of species; a representative few are presented in Table 2. It is interesting to note that, while transposon-mediated inactivation of essential genes would, *a priori*, have severe consequences on single-celled systems such as *S. cerevisiae*, it would be of little consequence to the coenocytic lifestyles of the small and filamentous. In fact, transposons may have had significant effects on fungal genome evolution.

**Table 2. Insertional Elements in Filamentous Fungi**

| Organism                     | Element | Size     | Length of ITR (bp) | Copy no.     | Reference          |
|------------------------------|---------|----------|--------------------|--------------|--------------------|
| <i>Nectria haematococca</i>  | Nht1    | 2.2 kbp  | –                  | 0–100 copies | <a href="#">70</a> |
| <i>Fusarium oxysporum</i>    | Fot1    | 1.9 kbp  | 44                 | 4–100        | <a href="#">71</a> |
| <i>Botrytis cinerea</i>      | Flipper | 1.8 kbp  | 48                 | 0–20         | <a href="#">72</a> |
| Magnaporthe grisea           | Pot2    | 1.9 kbp  | 43                 | ~100         | <a href="#">73</a> |
|                              | MGR586  | 1.8 kbp  | 42                 | 0–50         | <a href="#">74</a> |
|                              | MAGGY   | 5.6 kbp  | –                  | 0–100        | <a href="#">75</a> |
| <i>Aspergillus niger</i>     | Tan1    | 2.3 kbp  | 44                 | 1            | <a href="#">76</a> |
|                              | Vader   | 0.44 kbp | 44                 | ~15          | <a href="#">76</a> |
| <i>Cochliobolus carbonum</i> | Fcc1    | 1.8 kbp  | 64                 | >10          | <a href="#">77</a> |
| <i>Neurospora crassa</i>     | Tad1-1  | 6.9 kbp  | –                  | ~40          | <a href="#">78</a> |

Another source of CLPs may be dispensable or “B” chromosomes. These chromosomes are not present in all members of a species, and their absence is generally thought to offer no advantage or disadvantage to the organism. Fungi, as well as other plants and animals, have small dispensable chromosomes or portions of chromosomes that may be dispensable. Their existence has been conclusively demonstrated in *C. heterostrophus* (13), and *N. haematococca* (17).

There are a number of sequencing efforts in the fungi that will shortly answer many of the questions concerning comparative genome structure (18): *Aspergillus nidulans* (19), *Candida albicans* (20), *Neurospora crassa* (21), and *Ashbya gossypii* (22) are well on their way to completion.

## 5. Conclusion

The past decade has been a particularly exciting time for those working on the filamentous fungi.

With the ability to transform these organisms with homologous and heterologous fragments of DNA, it has been possible to make great strides in a number of areas. Not only has there been a much greater understanding of developmental pathways, molecular clocks, and mating type in these organisms, but topics such as incompatibility, pathogenesis and evolution have been addressed at a molecular level. Furthermore, as more genomes are sequenced, bioinformatics and functional genomics will play an even greater role as scientists begin to uncover the function of each of the proteins expressed in a given organism. What was once the gold of biochemists and geneticists can now be mined by molecular biologists studying human disease since fungi allow us to peek into the world of mammalian biology. Like all reviews of this nature, this one is far from complete, space limitations prevented the inclusion of many interesting topics, however excellent reviews in the field are available (a partial list is given below): transformation systems ([23](#), [24](#)); mating ([25-28](#)); control of gene expression ([29](#)); protein production and secretion ([30-32](#)); development and sporulation ([4](#), [7](#), [8](#), [25](#), [27](#), [33](#)); DNA modification strategies ([34-36](#)); fungal chromosomes and plasmids ([11](#), [37-41](#)); incompatibility ([42](#)); signaling ([43](#)); virulence and pathogenesis ([44-50](#)); circadian rhythms ([51](#), [52](#)); evolution ([53](#)).

Precise sequence information on these genomes can be accessed at the following sites:

1. University of Georgia (*A. nidulans*):  
[http://fungus.genetics.uga.edu:5080/Physical\\_Maps.html](http://fungus.genetics.uga.edu:5080/Physical_Maps.html) (as of 8/31/98)
2. Oklahoma State University cosmid and cDNA sequencing (*N. crassa* and *A. nidulans*):  
<http://www.genome.ou.edu/fungal.html> (as of 8/31/98)
3. University of New Mexico (*N. crassa*): (site currently unavailable) (as of 8/31/98)
4. University of Minnesota (*C. albicans*): <http://alces.med.umn.edu/Candida.html> (as of 8/31/98).

Note: There are excellent course outlines, descriptions, and lecture notes online: For plant pathogens:

1. SUNY College of Forestry and Environmental Science:  
<http://www.esf.edu/course/jworrall/fungi.html> (as of 8/31/98)
2. University of New Mexico: <http://taipan.nmsu.edu/EPWS310/exam-1-lectures.html> (as of 8/31/98)
3. On soil fungi by Dr. D. Jones at: (site currently unavailable)
4. And a more general outline by Dr. Bell at: (site currently unavailable)

## Bibliography

1. D. L. Hawksworth, B. C. Sutton, and G. C. Ainsworth (1983) *Dictionary of the Fungi*, Commonwealth Mycological Institute, Surrey, England.
2. N. Pace (1996) *ASM News* **62**, 463.
3. M. A. Nelson (1996) *Trends Genet.* **12**, 69–74.
4. M. L. Springer (1993) *BioEssays* **15**, 365–374.
5. A. Clutterbuck (1974) *Aspergillus nidulans*. In *Handbook of Genetics: Bacteria, Bacteriophages and Fungi* (R. C. King, eds.), Plenum Press, New York.
6. B. W. Bainbridge (1987) In *Genetics of Microbes* (B. W. Bainbridge, ed.), Blackie, Glasgow, London.
7. W. E. Timberlake (1990) *Ann. Rev. Genet.* **24**, 5–36.
8. T. H. Adams, J. K. Wieser, and J.-H. Yu (1998) *Microbiol. Mol. Biol. Rev.* **62**, 35–54.
9. W. E. Timberlake and A. J. Clutterbuck (1994) *Prog. Indust. Microbiol.* **29**, 383–427.
10. D. M. Kupfer et al. (1997) *Fung. Genet. Biol.* **21**, 364–372.
11. D. D. Perkins (1997) *Adv. Genet.* **36**, 239–398.

12. S. Iwaguchi, M. Homma, H. Chibana, and K. Tanaka (1992). *J. Gen. Microbiol.* **138**, 1893–1900.
13. T. Tzeng, L. K. Lyngholm, C. F. Ford, and C. R. Bronson (1992) *Genetics*. **130**, 81–96.
14. E. P. Rustchenko, T. M. Curran, and F. Sherman (1993) *J. Bacteriol.* **175**, 7189–7199.
15. P. J. Pukkila and C. Skrzynia (1993) *Genetics* **133**, 203–211.
16. N. J. Talbot, P. Oliver, and Coddington (1991) *Mol. Gen. Genet.* **229**, 267–272.
17. V. Miao, S. F. Covert, and H. D. VanEtten (1991) *Science* **254**, 1773–1776.
18. L. Hamer (1997) *Fung. Genet. Biol.* **21**, 8–10.
19. R. A. Prade et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14564–14569.
20. E. Tait et al. (1997) *Fung. Genet. Biol.* **21**, 308–314.
21. M. A. Nelson and D. O. Natvig (1997) *Fung. Genet. Biol.* **21**, 348–363.
22. S. Stenier, J. Wendland, M. C. Wright, and P. Philippsen (1995) *Genetics* **140**, 973–987.
23. J. Rambosek and J. Leach (1987) *CRC Crit. Rev. Biotechnol.* **6**, 357–393.
24. M. J. Hynes (1997) *J. Genet.* **75**, 297–311.
25. R. Kahmann, T. Romeis, M. Bolker, and J. Kamper (1995) *Curr. Opin. Genet. Dev.* **5**, 559–564.
26. J. W. Kronstad and C. W. Staben (1997) *Ann Rev. Genet.* **31**, 245–276.
27. E. Coppin, R. Debuchy, S. Arnaise, and M. Picard (1997) *Microbiol. Mol. Biol. Rev.* **61**, 411–428.
28. L. A. Casselton and N. S. Olesnick (1998) *Microbiol. Mol. Biol. Rev.* **62**, 55–70.
29. M. S. Sachs (1998) *Fung. Genet. Biol.* **23**, 117–124.
30. D. B. Archer and P. J. F. (1997) *Crit. Rev. Biotechnol.* **17**, 273–306.
31. R. J. Gouka, P. J. Punt, and C. A. van den Hondel (1997) *Appl. Microbiol. Biotechnol.* **47**, 1–11.
32. J. P. T. W. VanDen Hombergh, L. VanDe Vondervoort, F. Tachet, and J. Visser (1997) *Trends Biotechnol.* **15**, 256–263.
33. D. J. Ebbole (1996) *J. Genet.* **75**, 361–374.
34. J. T. Irelan and E. U. Selker (1996) *J. Genet.* **75**, 313–324.
35. M. J. Singer and E. U. Selker (1995) *Curr. Top. Microbiol. Immunol.* **197**, 165–177.
36. E. U. Selker (1997) *Trend. Genet.* **13**, 296–301.
37. E. G. Barry and 1996 (1996) *J. Genet.* **75**, 255–263.
38. A. J. F. Griffiths (1995) *Microbiol. Rev.* **59**, 673–685.
39. D. R. Stadler (1996) *J. Genet.* **75**, 265–280.
40. M. Plamann (1996) *J. Genet.* **75**, 351–360.
41. F. Kempken and U. Kuck (1998) *BioEssays* **20**, 652–659.
42. J. F. Leslie and K. A. Zeller (1996) *J. Genet.* **75**, 415–424.
43. F. Banuett (1998) *Microbiol. Mol. Biol. Rev.* **62**, 249–274.
44. D. G. Panaccione (1993) *Trends Microbiol* **1**, 14–20.
45. R. Oliver and A. Osbourn (1995) *Microbiology* **141**, 1–9.
46. L. H. Hogan, B. S. Klein, and S. M. Levitz (1996) *Clin. Microbiol. Rev.* **9**, 469–488.
47. M. Hensel and D. W. Holden (1996) *Microbiol.* **142**, 1049–1058.
48. J. W. Kronstad (1997) *Trends Plant Sci.* **2**, 193–199.
49. D. J. Ebbole (1997) *Trends Microbiol.* **5**, 405–408.
50. B. E. Corner and P. T. Magee (1997) *Curr. Biol.* **7**, 691–694.
51. J. C. Dunlap (1996) *Ann. Rev. Genet.* **30**, 579–601.
52. P. Ballario and G. Macino (1997) *Trends Microbiol.* **5**, 458–462.

53. J. W. Taylor (1995) Arch. Med. Res. **26**, 307–314.
54. R. Altmann-Johl and P. P. (1996) Mol. Gen. Genet. **250**, 69–80.
55. H. Brody and J. Carbon (1989) Proc. Natl. Acad. Sci. USA **86**, 6260–6263.
56. D. M. Geiser and W. E. Timberlake (1996) Curr. Genet. **29**, 293–300.
57. B. B. Magee and P. T. Magee (1987) J. Gen. Microbiol. **133**, 425–430.
58. W. Chu, B. B. Magee, and P. T. Magee (1993) J. Bacteriol. **175**, 6637–6651.
59. P. J. Pukkila and L. A. Casselton, (1991) In *More Gene Manipulations in Fungi* (Bennett and Lasure, eds.) Academic Press, New York, pp. 127–150.
60. P. E. Steele, G. F. Carle, S. Kobayashi, and G. Medoff (1989) Mol. Cell. Biol. **9**, 983–987.
61. M. J. Orbach (1989) Fung. Genet. Newslett. **36**, 14.
62. D. Z. Skinner, H. L. Leung, and S. A. Leong (1990) *Genetic Map of the Blast Fungus Magnaporthe grisea (n = 6)*. In *Genetic Maps* (S. J. O'Brien, ed.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
63. D. Z. Skinner, A. D. Budde, and S. A. Leong (1991) Theor. Appl. Genet. **87**, 545–557.
64. M. J. Orbach, D. Vollrath, R. W. Davis, and C. Yanofsky (1988) Mol. Cell. Biol. **8**, 1469–1473.
65. A. Radford and J. H. Parish (1997) Fung. Genet. Biol. **21**, 258–266.
66. M. T. Cushion, M. Kaselis, S. L. Stringer, and J. R. Stringer (1993) Infect. Immun. **61**, 4801–4813.
67. J. B. Fan et al. (1988) Nucleic Acids Res. **17**, 2801–2818.
68. K. McCluskey and D. Mills (1990) Mol. Plant-Microb. Interact. **3**, 366–373.
69. T. G. Kinscherf and S. A. Leong (1988) Chromosoma **96**, 427–433.
70. J. Enkerli, G. Bhatt, and S. F. Covert (1997) Mol. Plant-Microb. Interact. **10**, 742–749.
71. M. J. Daboussi, T. Langin, and Y. Brygoo (1992) Mol. Gen. Genet. **232**, 12–16.
72. C. Levis, D. Fortini, and Y. Brygoo (1997) Mol. Gen. Genet. **254**, 674–680.
73. P. Kachroo, S. A. Leong, and B. B. Chatoo (1994) Mol. Gen. Genet. **245**, 339–348.
74. P. Kachroo, M. Ahuja, S. A. Leong, and B. B. Chatoo (1997) Curr. Genet. **31**, 361–369.
75. M. L. Farman, Y. Tosa, N. Nitta, and S. A. Leong (1996) Mol. Gen. Genet. **251**, 665–674.
76. E. Nyssonen et al. (1996) Mol. Gen. Genet. **253**, 50–56.
77. D. G. Panaccione, J. W. Pitkin, J. D. Walton, and S. L. Annis (1996) Gene **176**, 103–109.
78. E. B. Cambareri, J. Helber, and J. A. Kinsey (1994) Mol. Gen. Genet. **242**, 658–665.
79. J. W. Taylor (1995) Arch. Med. Res. **26**, 304–307.

## Fusion Gene, Fusion Protein

Fusion genes consist of unrelated genes, or gene fragments, that have been joined in-frame to produce a coding sequence that is **transcribed** and **translated** as a single unit. The resulting gene product is called a *fusion protein* (as well as a chimeric or hybrid protein). The fused gene elements most often encode functionally and structurally distinct protein **domains**.

In nature, several examples of fusion genes exist. Circumstantial evidence supports the hypothesis that the assembly of relatively small domains into functional proteins is an important factor in [evolution](#) (1). These natural gene fusion events result from processes such as [gene duplication](#) and

[exon shuffling](#). The possibility of deleting and adding such domains during evolution suggests that proteins are often tolerant to changes involving whole domains.

[Ubiquitin](#), a small, highly conserved **eukaryotic** protein, is another example of a naturally occurring gene fusion, as it is expressed from gene fusions to either itself (generating polyubiquitin) or to one of two **ribosomal** proteins (2). A third example is the **genomes** of many **viruses** that consist of several genes fused together. On expression, a [polyprotein](#) is produced, which is further processed by a virus-encoded [protease](#) to generate single proteins (3).

Gene fusions generated by use of [recombinant DNA](#) technology have been used experimentally for a number of purposes, such as protein expression and purification, construction of reporter systems, and medical applications, as well as in systems for directed evolution of proteins. The first-mentioned application is by far the most common and, therefore, this aspect of gene fusion technology will be treated in greatest detail (see [Expression Systems](#)).

## 1. Gene expression and protein purification

There can be several reasons for applying the fusion gene approach in a gene expression strategy. Fusion to the desired product of an [affinity chromatography](#) or purification “handle” or “tag,” with unique **ligand-binding** characteristics, allows the use of the properties of the fused handle for purification. In many cases, the handle is useful for detection purposes as well. The handle can be removed later, if necessary, to obtain the native product. The most commonly used handles seem to be [protein A](#), which binds to [immunoglobulin G](#); [glutathione S-transferase](#), which binds to its substrate [glutathione](#); [maltose-binding protein](#), which binds to maltose, and poly**histidine**, which binds to metal ions (4). A combined affinity handle with more tags integrated can allow a simplified screening of various purification and detection methods (5). These general fusion systems for expression and purification ensure a reasonable expression level (when the fusion partner gene is positioned upstream of the gene of interest) and allow rapid recovery of the gene product without the need for time-consuming optimization. Other advantages of gene fusion technology with relation to protein expression and purification are summarized in Table 1.

**Table 1. Alternative Purposes of Using a Fusion Strategy in Protein Expression and Purification**

| <b>Problem</b>                             | <b>Solution</b>   |
|--|---|
| Low solubility, caused by misfolding       | Fusion to a protein with solubilizing or chaperonelike properties, eg, protein A, thioredoxin, or ubiquitin (7, 8)  |
| Proteolytic degradation                    | A dual-affinity approach in which the protein of interest is fused between two different affinity chromatography handles (9)  |
| Proteolytic degradation; recovery problems | Fusion to trpE, cII, or an amphiphilic extension often results in the formation of <a href="#">inclusion bodies</a> , offering protection and simplified purification; <b>denaturation</b> followed by refolding is required (10, 11) |
| N-terminal heterogeneity (in bacteria)     | Addition of fusion partner, whose removal generates a homogeneous N-terminus  |
| Detection                                  | Fusion to a partner that can be readily detected by, eg, immunological methods or its enzymatic activity  |

|   |   |
|---|---|
| Toxicity; formation of <a href="#">disulfide bonds</a> (in bacteria); recovery; proteolytic degradation | The protein of interest has to be removed from the cytoplasm; fusion with a <a href="#">signal peptide</a> directs the protein of interest to a specific compartment (eg, periplasm, cell wall, culture medium) [see <a href="#">Secretion Vector</a> ] |
|---|---|

---

The removal of the affinity handle is most often the bottleneck of the purification procedure. The usefulness of chemical cleavage is limited, as these agents recognize only one amino acid and harsh conditions are often applied during the cleavage. Enzymatic methods can represent a problem if the [proteinase](#) used is not sufficiently specific. Furthermore, the cleavage efficiency can be low because of the inaccessibility of the recognition sequence by the proteinase. Finally, one or more purification steps are required to remove the proteinase again after cleavage. Now, the use of recombinant fusion proteinases equipped with a purification handle may be able to solve some of these problems. If the handle is of the same kind as that of the fusion protein to be processed, the cleavage can, in principle, be carried out directly on the column (6). Another solution to the problem of removing the fusion partner is presented in the IMPACT system (New England Biolabs), where the protein of interest is C-terminally fused to an affinity tag that binds chitin. These two fusion partners are separated by a modified [protein splicing](#) element, the intein domain. The splicing element undergoes self-cleavage at its N-terminus under reducing conditions when the fusion protein is bound on a chitin affinity column. The cleavage causes the release of the target protein while leaving the affinity tag bound to the column.

## 2. Potential medical applications

Many pharmaceutical protein drugs are rapidly cleared from the circulation after administration, making frequent supply of large doses of the therapeutic drug necessary. It has been demonstrated that the *in vivo* half-life of a drug can be increased by fusion to a carrier that is more stable *in vivo* (12).

Drug targeting is another interesting application of gene fusions. One portion of the fusion protein targets the drug to a specific cell type, for example, by interacting with a **receptor** molecule, whereas the other part, the drug itself, has the specific activity of interest.

**Immunogenic** peptides or proteins for the production of vaccines are often small and difficult to produce because high purity is necessary. A means of overcoming these difficulties is to produce the protein as a fusion protein. Normally, the native protein does not have to be released by cleavage. A particularly effective fusion partner is the serum [albumin](#) -binding region of streptococcal protein G, which seems to have inherent immunopotentiating properties, probably by providing T-cell help for [antibody](#) production (13).

The gene fusion technology has also proved valuable in gene therapy. Transcriptionally inducible systems have several critical shortcomings, such as the delay between induction and the appearance of optimal protein activity, plus the problem of activating specific genes without affecting others. **Hormone-binding domains (HBD)** of **steroid receptors** can be used to regulate heterologous proteins post-translationally. The protein to be regulated is fused to an HBD, so that it becomes inactive in the absence of hormone. The inactive state can be rapidly reversed by the addition of the cognate steroid hormone (14). Chimeric versions of the steroid hormone receptors have also been applied for the regulation of **gene expression** at the transcriptional level. The chimeric receptor is composed of an HBD, a specific DNA-binding domain, and a transcriptional activation domain, each with unique specificity. The chimeric receptor is kept in the inactive state, until an activating ligand is bound. The active chimeric receptor is now capable of binding to specific response elements

associated with the target gene, thereby inducing its expression (15).

### 3. Research Applications

Gene fusion technology has also been used to generate powerful research tools of general interest (16). The **phage display** system is based on the fusion of the gene encoding a product of interest to one of the genes encoding a **bacteriophage** coat protein. The fusion protein will thereby be displayed on the surface of the phage, and the corresponding gene will be packaged into the phage particle, allowing its later identification by [DNA sequencing](#). The system is used for the selection of a mutational variant of the protein of interest with specific characteristics from a pool of randomized versions of the protein.

The yeast [two-hybrid system](#) is a method developed for the detection of [protein–protein interactions](#). The system takes advantage of the fact that many eukaryotic transcriptional activators, including the yeast GAL4, contain two discrete domains, the DNA-binding domain and the transcriptional activation domain. If the two domains are physically separated, GAL4 becomes inactive. Two different [cloning](#) vectors are used to generate separate fusions of these GAL4 domains to genes encoding proteins that potentially might interact with each other. The recombinant hybrid proteins are coexpressed in yeast. If the non-GAL4 parts of the two hybrid proteins interact with each other, the DNA-binding domain will be tethered to the transcriptional activation domain. As a result of a two-hybrid interaction, the GAL4 transcriptional activator will be functionally reconstituted and will activate transcription of [reporter genes](#) having upstream GAL4 sites. In this way, the protein–protein interaction is made **phenotypically** detectable.

Scientists have found inspiration to develop gene fusion technology by studying nature's way of evolving proteins with new functions by fusing genes encoding distinct protein domains. Gene fusion technology has now found widespread use in biotechnology and medicine and promises to be a key technology for the future.

### Bibliography

1. W. Gilbert (1978) *Nature* **271**, 501.
2. D. Finley, B. Bartel, and A. Varshavsky (1989) *Nature* **338**, 394–401.
3. I. G. Maia, K. Séron, A.-L. Haenni, and F. Bernardi (1996) *Plant Molec. Biol.* **32**, 367–391.
4. H. M. Sassenfeld (1990) *Trends Biotechnol.* **8**, 88–93.
5. J. Nilsson et al. (1997) *Prot. Expr. Purif.* **11**, 1–16.
6. P. A. Walker et al. (1994) *Biotechnology* **12**, 601–605.
7. E. Samuelsson et al. (1991) *Biotechnology* **9**, 363–366.
8. P.-Å. Nygren, S. Ståhl, and M. Uhlén (1994) *Trends Biotechnol.* **12**, 184–188.
9. M. Murby et al. (1991) *Biotechnol. Appl. Biochem.* **14**, 336–346.
10. M. Uhlén and T. Moks (1990) *Meth. Enzymol.* **185**, 129–143.
11. J. G. Thomas, A. Ayling, and F. Baneyx (1997) *Appl. Biochem. Biotech.* **66**, 197–238.
12. S. C. Makrides et al. (1996) *J. Pharmacol. Exp. Therapeut.* **277**, 534–542.
13. A. Sjölander et al. (1997) *J. Immunol. Methods* **201**, 115–123.
14. D. Picard (1994) *Curr. Opin. Biotech.* **5**, 511–515.
15. V. E. Allgood and E. M. Eastman (1997) *Curr. Opin. Biotech.* **8**, 474–479.
16. J. B. Allen, M. W. Walberg, N. C. Edwards, and S. J. Elledge (1995) *Trends Biochem. Sci.* **20**, 511–516.

### Suggestions for Further Reading

17. E. R. LaVallie and J. M. McCoy (1995) Gene fusion expression systems in *Escherichia coli*. *Curr. Opin. Biotechnol.* **6**, 501–506. An overview of the possibilities in the *E. coli* expression



system provided by gene fusion techniques. A good reference source.

18. J. Nilsson, S. Ståhl, J. Lundeberg, M. Uhlén, and P.-Å. Nygren (1997) Affinity fusion strategies for detection, purification, and immobilization of recombinant proteins. *Prot. Expr. Purif.* **11**, 1–16. Reviews how to benefit from a fusion strategy when working with recombinant proteins. Contains an extensive reference list.
19. M. Uhlén, G. Forsberg, T. Moks, M. Hartmanis, and B. Nilsson (1992) Fusion proteins in biotechnology. *Curr. Opin. Biotechnol.* **3**, 569–575. Reviews the use of fusion proteins in biotechnology. Provides many good references.

## G Banding

G-banding is a simple technique that is usually the method of choice for identifying and studying [chromosomes](#). The first step is to fix the chromosomes with methanol and acetic acid (Carnoy's reagent), which removes the vast majority of the [histone](#) proteins. Non-histone proteins remain behind and obviously contribute most to banding patterns. The structural basis of banding is unknown, although the G-bands are presumably where DNA is accessible to thiazine dyes. The second step in G-banding is usually to pretreat the chromosomes by enzymatic digestion using a [proteinase](#), such as **trypsin** or a mixture of proteinases, such as [pronase](#). Alternatively, any other treatment that damages the protein content of chromosomes, such as high pH, can be used. As enzyme digestion proceeds, the removal of structural proteins relaxes the chromosomal compaction, leading to the selective local expansion of chromosomal volume (1). The last step is to stain the chromosomal preparations with Giemsa–Leishman's stain or Wright stain. The dark G-bands alternate with light bands in **prophase** and **metaphase** chromosomes.

General features of G-bands in terms of chromatin structure include the presence of AT-rich DNA sequences, the absence of potentially active **promoter** elements ([CpG Islands](#)), and indications that chromatin is late replicating, condenses early in mitotic and meiotic prophase, and is not **DNase I sensitive**.

G-banding is most useful in the analyzing vertebrate chromosomes, where the large number of bands and hence information content has allowed accurate **phylogenetic** comparisons between species. In medical [cytogenetics](#), banding patterns of human chromosomes achieve a resolution of up to 2000 G-bands. Chromosomes can be divided into segments dependent on the most prominent bands, which can then be further subdivided by finer analysis of visible bands. G-banding allows defining chromosomal rearrangements and breakages. For example, a chromosome break at the point Xq 27.3 (long arm, band 2, primary subband 7, subdivision 3) is associated with the form of mental retardation known as Fragile X syndrome (see [Acentric Fragment](#)).

### Bibliography

1. G. D. Burkholder and L. Duczek (1982) *Chromosoma* **87**, 425–435.

### Suggestion for Further Reading

2. R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A Synthesis*, Wiley-Liss, New York.

## G-Protein-Coupled Receptors

Almost two thousand **genes** have been identified, on the basis of their predicted [protein structure](#), that encode for [proteins](#) characterized as G-protein-coupled receptors (GPCRs). It is likely that there are actually several thousand. GPCRs are currently classified into 100 or so subfamilies, based on their structure, ligand specificity, and biological effects. These molecules have a characteristic sequence that predicts seven **transmembrane**-spanning segments, with a cytoplasmic tail of varying lengths. GPCRs are known to interact with a variety of [hormones](#), [neuropeptides](#), [chemokines](#), biogenic amines, nucleosides, eicosanoids, phospholipids, [growth factors](#), and aromatic compounds that act as olfactants or odorants. Some examples include thyrotropin, interleukin 8, epinephrine, and bombesin. Moreover, scores of “orphan” receptors have been identified by cDNA [cloning](#) and are still in search of ligands.

### 1. Receptor Structure

The wide variety of ligands for GPCRs dictate that several structural variations might allow for the variety of their regulatory features. One common feature of all of these receptors is seven stretches of 20 to 25 **hydrophobic** amino acid residues, which are thought to form [a-helices](#) that span the membrane (see [a-Helix](#)). Thus, the overall structures of these receptors resemble that of [bacteriorhodopsin](#) (1). The amino-terminal sequences lie outside the cell, and they vary in size among the family from 7 to 600 amino acid residues. In the case of some receptors, especially those for larger peptide hormones, such as thyrotropin or thrombin, the amino-terminal sequences are critically involved in ligand binding. In most cases, however, ligand binding is thought to occur primarily within key residues in the membrane-spanning regions. A series of mutational studies with the b-adrenergic receptor and other GPCRs (2) have suggested that the ligand may associate with the a-helix near the outer surface, binding in a plane parallel with the membrane. Another conserved structural feature is the loop between the fifth and sixth a-helices, the domain most critical for interactions with G proteins (see [GTP-Binding Proteins](#)). It is likely that agonist binding induces a conformational change in the receptor that permits interactions of this domain with G proteins (3). The carboxyl tail of the receptor is cytoplasmically oriented and varies in length among receptors. Mounting evidence suggests that this domain may play a crucial role in receptor cross-talk and desensitization.

Identification of the regions in GPCRs involved in ligand binding has emerged mainly from studies employing point mutations and receptor chimeras. Advances in this area have been led by studies on b-adrenergic receptors (4). Mutant forms of the receptor have been constructed and expressed in cells, to evaluate the binding of both agonists and antagonists. Studies such as these have led to models in which the ligand is thought to bind at the outer surface of the receptor, among several of the transmembrane domains. In the case of the catecholamine agonist isoproterenol, two [serine](#) residues in the fifth transmembrane helix are thought to form [hydrogen bonds](#) with the hydroxyl groups on the catechol ring (2). In addition, a phenylalanine residue in a-helix 6 participates in binding via hydrophobic interactions with the catechol ring itself. Similar types of interactions have been elucidated for nucleoside and nucleotide receptors.

Agonist binding to GPCRs is thought to induce movements among the transmembrane domains, a conformational change that is translated into G-protein interactions. Although the molecular dynamics of this process remain poorly understood, a number of models have been proposed to explain signal generation within the receptor (3). One aspect of receptor activation that may play an important role is receptor dimerization, a process that is well characterized in activation of other receptor subtypes. Recently, it has been proposed that the protonation of key transmembrane residues are required for activation. Additionally, it has been reported that intramolecular

interactions within certain intracellular loops maintain the receptor in the inactive state prior to ligand binding, which is subsequently disrupted by the interaction with activating ligand.

The specificities of receptors, both for ligands and for G proteins, have been explored by construction of chimeras. Although the overall structural features of these receptors are generally conserved, there is little sequence similarity among the superfamily. Indeed, even among the three known subtypes of the b-receptor, there is only 50% identity. Thus, chimeric receptors have revealed domains of the proteins that are necessary for binding specificity and that dictate which G proteins might be coupled to each receptor (4).

## 2. G Proteins

The proteins that serve as transducers for the GPCRs are known as [GTP-binding proteins](#), or G proteins. These molecular linkers exhibit a heterotrimeric structure, consisting of  $\alpha$ ,  $\beta$ , and  $\gamma$  subunits. The [heterotrimeric G proteins](#) belong to a superfamily of [Gtpases](#), which share a common structural core. In general, the functions of G proteins are dictated by the receptors with which they interact, an interaction mediated by the  $\alpha$  subunit. Thus far, over 30 different  $\alpha$  subunits have been identified, along with six  $\beta$  and  $\gamma$  subunits. The  $\alpha$  subunits fall into five different classes, each with a characteristic impact on its effector.

The  $\alpha$  subunit of G proteins interacts with guanyl nucleotides and in general controls the activity of the complex (5, 6). When no ligand is bound to the receptor,  $\alpha$  is maintained in the GDP-bound state. In this form,  $\alpha$  is complexed to the  $\beta$  and  $\gamma$  subunits, which prevents it from interacting with effectors.  $\beta$  and  $\gamma$  subunits together form a tight complex. Upon ligand binding, the receptor interacts with the  $\alpha$  subunit, leading to the displacement of GDP and the subsequent binding of GTP. In its GTP-bound state,  $\alpha$  then dissociates from the  $\beta\gamma$  complex and is activated, interacting with effectors such as [adenylate cyclase](#). This activation is temporal and is terminated by the intrinsic GTPase activity of  $\alpha$ , which hydrolyzes the GTP, leaving GDP bound to the protein.

The GTPase activity of G family members is subject to different types of regulation. The proteins can be frozen in their activated state by adding nonhydrolyzable analogues of GTP, such as GTP $\gamma$ S, effectively blocking the GTPase activity of the protein. Additionally, certain bacterial [toxins](#) irreversibly modify G-protein  $\alpha$  subunits, catalyzing the covalent addition of ADP-ribose from intracellular NAD. ADP-ribosylated  $\alpha$  cannot hydrolyze GTP, thus leaving the protein in a persistently active state. In the case of **cholera toxin**, this activation results in a modification of  $\alpha$ , causing the persistent activation of adenylyl cyclase, and leading to massive diarrhea (5, 6).  $\alpha$  proteins are also regulated by endogenous proteins called for regulators of G proteins (RGs), (7). This family of regulatory proteins now comprises over 20 members, and they seem to function as long-term regulators of G protein signaling.

Although they were originally thought to function only as inhibitors of  $\alpha$ ,  $\beta\gamma$  subunits in some cases can interact directly with effectors upon their release from  $\alpha$  subunits, causing either stimulation or inhibition of activity. For example, in certain neuronal systems, cells may contain different forms of adenylyl cyclase that are differentially regulated by G proteins.  $\beta\gamma$  subunits can stimulate the activity of the type III form of cyclase, but have no effect, or can even inhibit, the type II form of the enzyme. Additionally,  $\beta\gamma$  subunits can directly stimulate [phospholipase C-b](#) in some cells (8).

## 3. Effectors

In most cases, G-protein  $\alpha$  subunits interact directly with one or more effector proteins in the plasma membrane, to induce the release of soluble [second messengers](#) that in turn amplify the hormonal signal into the cell. While the number of effectors continues to grow, four distinct types have been studied in detail: [adenylate cyclases](#), [phospholipases](#), **ion channels**, and **phosphodiesterases**.

## Bibliography

1. J. Nathans and D. S. Hogness (1993) *Cell* **34**, 807–814.
2. R. J. Lefkowitz (1998) *J. Biol. Chem.* **273**, 18677–18680.
3. U. Gether and B. K. Kobilka (1998) *J. Biol. Chem.* **273**, 17979–17982.
4. T. H. Ji, M. Grossman, and I. Ji (1998) *J. Biol. Chem.* **273**, 17291–17302.
5. H. R. Bourne, D. A. Sanders, and F. McCormick (1990) *Nature* **348**, 125–132.
6. J. R. Hepler and A. G. Gilman (1992) *Trends Biochem. Sci.* **17**, 383–387.
7. V. Y. Arshavsky and E. N. Pugh (1998) *Neuron* **20**, 11–14.
8. D. E. Clapham and E. J. Neer (1993) *Nature* **365**, 403–406.

## *gal* Operon

The *gal* operon of **bacteria** encodes [enzymes](#) of galactose metabolism that constitute an amphibolic pathway and is transcribed by two **promoters**. Several [cis-acting](#) DNA sequences (**control elements**) and a variety of regulatory proteins, including a **histone**-like protein, modulate the two promoters in a multitude of ways usually found in animals and not in bacteria (Table 1). The *gal* operon has revealed several new features of gene-specific transcriptional regulation that were previously unrecognized. Perhaps such multivalent control mechanisms are required to regulate the synthesis of amphibolic enzymes.

**Table 1. Repression and Activation of Transcription of the *gal* Operon**

| Promoters | Regulatory protein and DNA elements <sup>a</sup> |              |            |
|-----------|--|--------------|------------|
|           | GalR, HU   | GalR or GalS | cAMP•CRP   |
|           | $O_E$ , $O_I$                                    | $O_E$        | <i>AS</i>  |
| <i>P1</i> | Repression                                       | Repression   | Activation |
| <i>P2</i> | Repression                                       | Activation   | Repression |

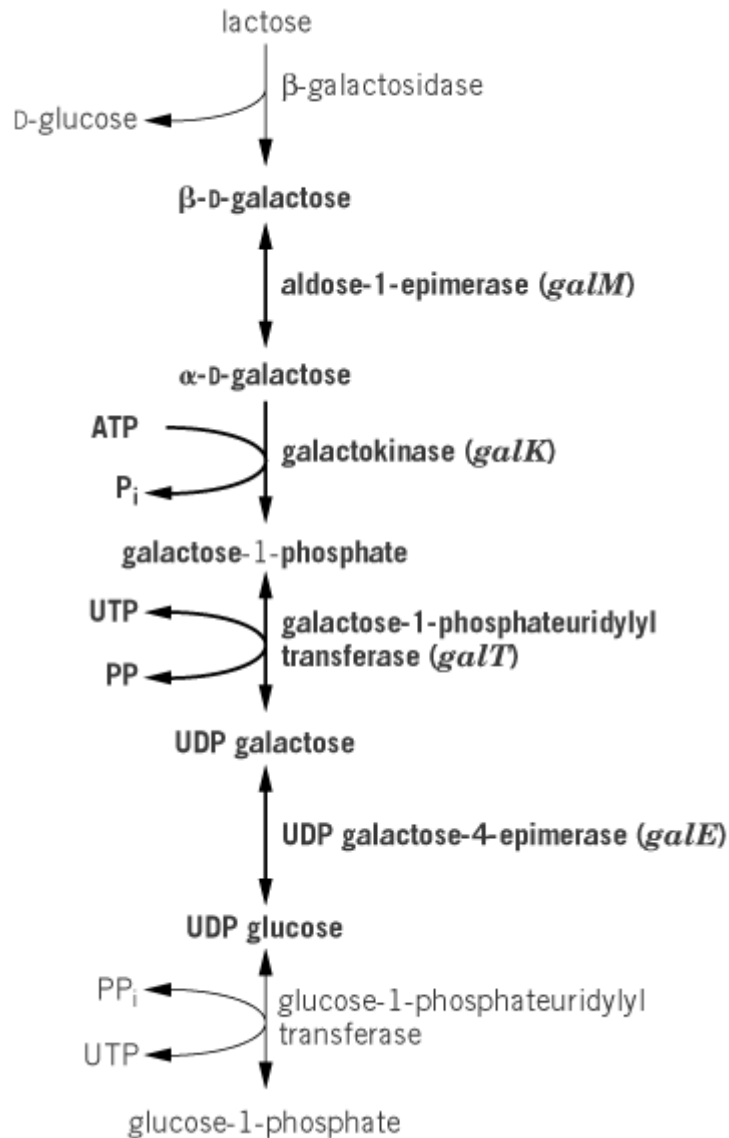
<sup>a</sup> The regulatory effects of GalR and GalS are epistatic to the effects of cAMP•CRP.

### 1. Enzymes of the *gal* Operon

The galactose metabolizing enzymes (Fig. 1) are involved in (i) the catabolism of D-galactose that was either imported into the cell by **permeases** or generated intracellularly by hydrolysis of disaccharides and (ii) the synthesis of precursors (UDP galactose and UDP glucose) of complex carbohydrates (1, 2). For catabolism, only α-D-Galactose is converted to galactose-1-phosphate by galactokinase (3). β-D-galactose, generated, for example, by hydrolysis of lactose by **beta-galactosidase**, must change to the α-anomer before it can be phosphorylated. Although β-D-galactose

can mutarotate spontaneously to the  $\alpha$ -anomer at a slow rate, the enzyme aldose-1-epimerase is largely responsible for the mutarotation *in vivo* (2). Thus, aldose-1-epimerase links the enzymes of lactose and galactose metabolism into a common pathway (Fig. 1).

**Figure 1.** The Leloir pathway of D-galactose metabolism. As shown, D-galactose is generated intracellularly by hydrolysis of the disaccharide lactose. The parts of the pathway catalyzed by enzymes of the *gal* operon are shown in bold. They are encoded by the genes shown within the parenthesis in italics.

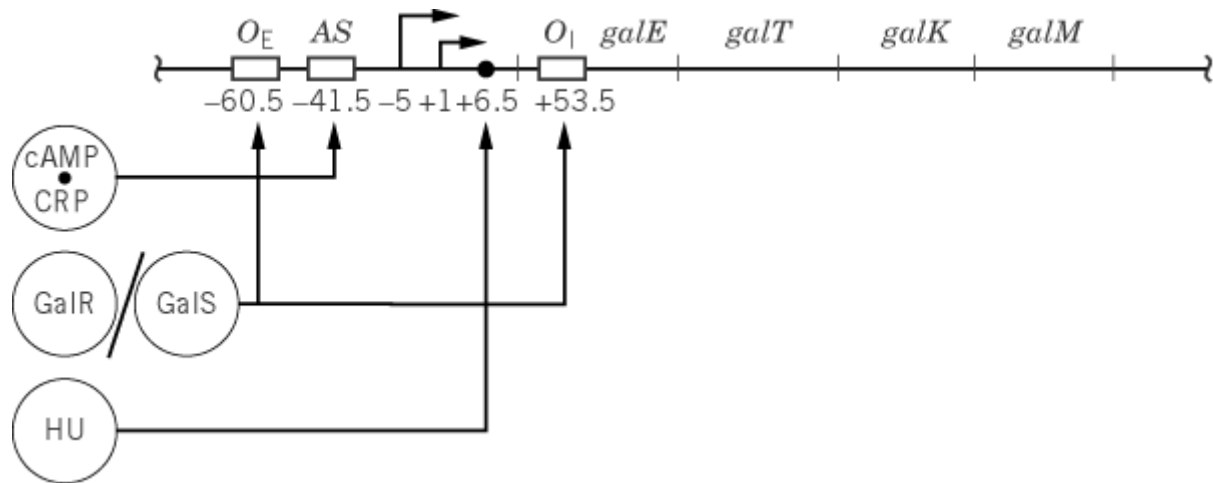


## 2. Regulation of Transcription

The structural genes and the associated regulatory elements of the *gal* operon are shown in Figure 2. The *gal* operon is transcribed from two promoters, *P*<sub>1</sub> and *P*<sub>2</sub>, which are separated by 5 base pairs (bp) (4, 5). The dominant regulator of the *gal* operon is the Gal repressor (GalR) (6). GalR, together with the histone-like protein, HU, acting as a corepressor, keeps the expression from the *gal* promoters low (7, 8). The operon is induced 15-fold in the presence of D-galactose or some of its nonmetabolizable analogs, such as D-fucose (9). The inducer lifts the repression by an **allosteric** effect on GalR. It is not known whether the true inducer is the  $\alpha$ - or the  $\beta$ -anomer of the sugar, or both. Besides the major negative control by GalR and HU, the two promoters are also regulated by

GalR alone (10), by a Gal isorepressor, GalS (11-13), and by a complex of cyclic AMP (cAMP) bound to cyclic AMP receptor protein (CRP), cAMP•CRP (4, 5, 14-16). As described below, *P1* and *P2* are regulated by these regulators in opposite directions. The *gal* promoters are also modulated by *cis*-acting DNA sequences without the participation of any regulatory proteins.

**Figure 2.** The structure of the *gal* operon. The regulators are shown as circles and their cognate DNA control elements as open bars. Their mode of actions are explained in the text. The adenine tracks present upstream of the promoter are not shown.



### 3. Regulation without Regulatory Proteins

#### 3.1. Control of *P1* by Adenine Tracks

The intrinsic strengths of the two promoters are comparable, and they are moderately active in the absence of any regulatory proteins, both *in vivo* and *in vitro*. *In vitro*, the intrinsic strength of *P1* is twofold enhanced because of the periodic presence of four to six adenine residues centered at positions -84.5, -74 and -63 on the DNA (17). The adenine tracks bend the DNA toward the face of *P1* to which **RNA polymerase** binds. The DNA curvature induced by the adenine tracks may help formation of a RNA polymerase-promoter complex (caging) that is more optimal for transcription initiation at *P1* (18).

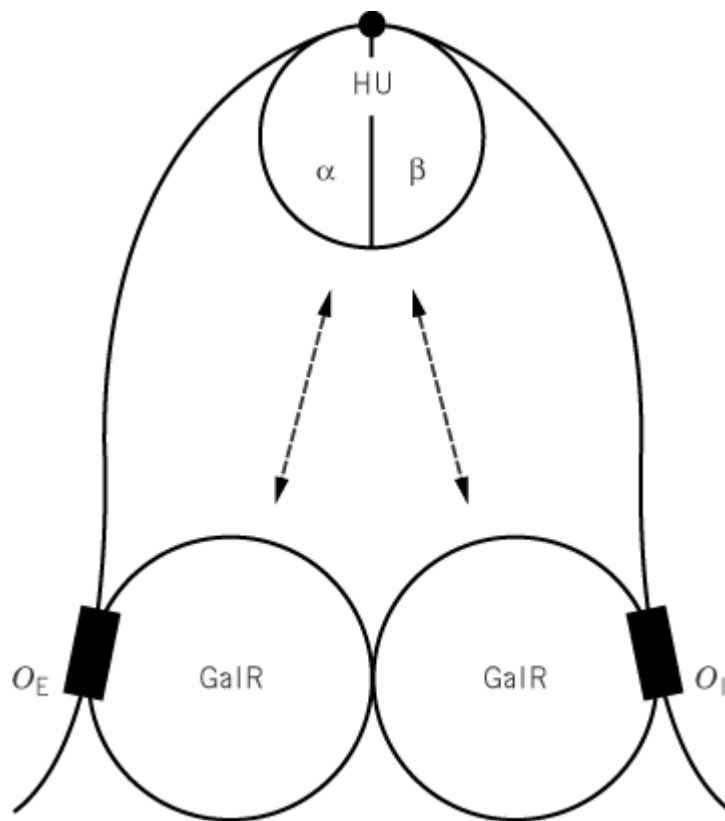
#### 3.2. Control of *P2* by UTP

Transcription of the *gal* operon from the *P2* promoter is very low when the concentration of UTP is high and *vice versa*, both *in vivo* and *in vitro* (19). UTP regulates the step of promoter clearance by RNA polymerase at the *P2* promoter in an intriguing way. *In vitro*, RNA polymerase clears at *P2* poorly and makes a large amount of aborted RNA oligomers, unlike at *P1*. At high concentrations of UTP, it also synthesizes pseudo-templated RNA oligomers of the composition pppAU<sub>n</sub> ( $n = 2$  to  $\geq 20$ ) because the enzyme “stutters” while adding uridine residues present at positions 3–5 of the *P2* RNA and not the *P1* RNA. RNA polymerase clears the promoter more efficiently and makes template-encoded normal *gal* RNA at low UTP concentrations. The involvement of UTP in the synthesis of UDP sugars by the galactose pathway may be the reason for the control by UTP (Fig. 1) analogous to the way UTP controls its own synthesis in the pyrimidine operons (20). At high UTP concentrations, the levels of UDP galactose and UDP glucose are high and inhibit the synthesis of enzymes, specifically, galactose-1-phosphate uridylyltransferase and uridine diphosphogalactose epimerase, that make them.

### 4. Repression by GalR and HU: DNA Looping

Transcription from the two *gal* promoters is coordinately repressed (negative control) by GalR and HU acting together (7, 8). The repression requires binding of GalR to two operators,  $O_E$  and  $O_I$ , which are similar 16-bp sequences with a dyad symmetry (Fig. 2) (21-23).  $O_E$  is located upstream of the promoters at position  $-60.5$ , whereas  $O_I$  is located within the structural gene *galE* at position  $+53.5$ . The two operator-bound GalR molecules associate with each other in the presence of the cofactor, HU (a heterodimer of two subunits, HUa and HUb), resulting in the formation of a DNA loop encompassing the promoters (8, 24) (Fig. 3). The DNA looping presumably changes the structure of the promoters, making them refractory to RNA polymerase caging (18, 25). The nucleoprotein complex that represses the *gal* promoters is called the Gal repressosome, because its formation is by binding of a histone-like component of the bacterial nucleoid as well as a repressor (26). Although GalR binds to  $O_E$  and  $O_I$  in both the absence and presence of HU, GalR binding alone does not bring about DNA looping and the associated simultaneous repression of the promoters (27). HU is absolutely necessary for the effect; it cannot be replaced by other histone-like proteins. Although HU is not a sequence-specific DNA binding protein, one molecule of HU heterodimer binds to and bends the *gal* DNA centered at an architecturally critical position (8). Genetic analysis and modeling defined the GalR surfaces interacting to form a stacked, V-shaped tetrameric structure (28, 29). Evaluation of the DNA elastic energies gave unambiguous preference to a DNA loop in which  $O_E$  and  $O_I$  adopt an antiparallel orientation causing undertwisting of DNA (29). Since HU binding depends on GalR binding to both  $O_E$  and  $O_I$  and the binding of GalR to the operators is enhanced by HU, GalR and HU bind **cooperatively** (8). GalR piggybacks HU to the critical position on the DNA through a specific GalR-HU interaction (31). The entire process facilitates the GalR-GalR interaction resulting in cooperativity. The GalR-HU contact may be transient and was not in the final repressosome structure. The dependence of HU binding on GalR renders the HU-containing nucleoprotein complex that brings about repression of  $P1$  and  $P2$  sensitive to inducer D-galactose for transcription derepression. Such a mechanism is an example of how DNA that is “condensed” by binding proteins and made refractory to RNA polymerase action can become available for transcription in response to specific signals.

**Figure 3.** The DNA looping of the *gal* promoter by binding of the GalR regulator and HU corepressor. The details of binding of the two proteins are discussed in the text.



## 5. Regulation in the Absence of DNA Looping: Interaction between GalR and RNA Polymerase

Without DNA looping (ie, in the absence of HU), occupation of  $O_E$  alone by GalR represses  $P1$  (by about four- to fivefold) and activates  $P2$  (twofold) (10, 27). Binding to  $O_I$  does not affect this dual control. The activation of  $P2$  and repression of  $P1$  are independent of each other. GalR exerts its specific regulatory effect on one promoter even when the other is mutated. The activation of  $P2$  or repression of  $P1$  is not an intrinsic property of the promoter; the regulation can be reversed by switching the angular orientation of the promoters relative to  $O_E$  by inserting a 5-bp segment, that is, half of a DNA helix, between  $O_E$  and the promoters. Both activation of  $P2$  and repression of  $P1$  require the formation of a specific GalR-RNA polymerase-DNA complex at each promoter. RNA polymerases containing a subunits that carry specific amino acid alterations in their carboxy terminal domain (aCTD), or that are missing the aCTD, abolish the regulatory effect of GalR without affecting intrinsic transcription, suggesting that GalR activates and represses by a direct contact with the promoter-bound RNA polymerases.

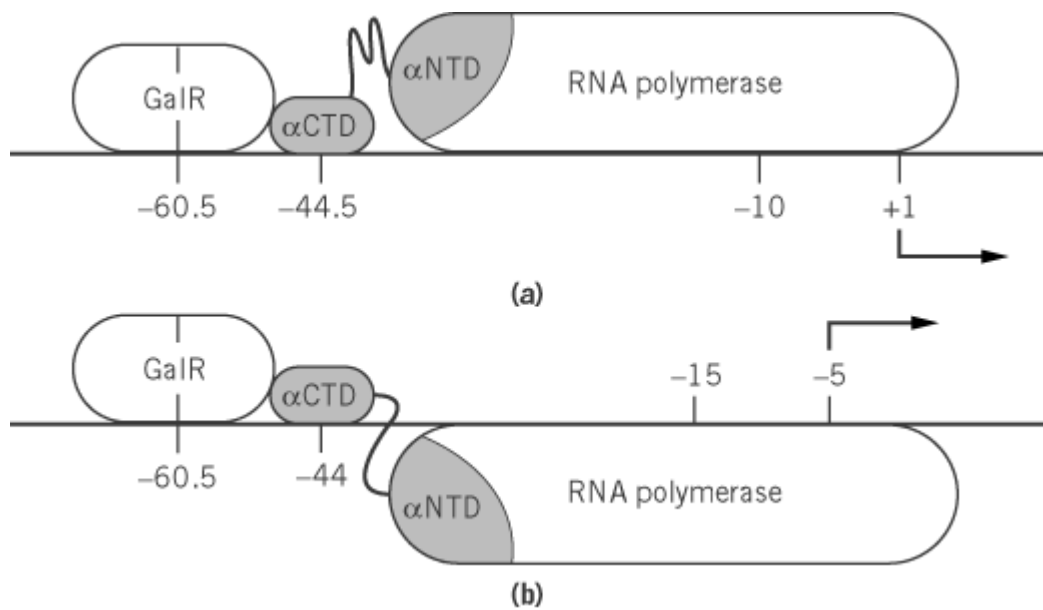
### 5.1. Activation of $P2$

Activation of  $P2$  by GalR bound to DNA at position  $-60.5$  ( $-55.5$  with respect to  $P2$ ) parallels the activation of transcription of several promoters by activators that also bind around position  $-60$  on the DNA and act by enhancing open complex formation (32). An exposed segment (not necessarily the same segment) of aCTD is contacted by the DNA-bound activator to activate transcription (33, 34). aCTD is connected to the rest of the RNA polymerase molecule by a flexible hinge (35, 36). In  $P2$ , as in these systems, aCTD binds at the region 40 bp upstream of  $P2$  (37) (Fig. 4).

**Figure 4.** Repression of  $P1$  promoter and activation of  $P2$  promoter by a contact between DNA-bound GalR and the aCTD of RNA polymerase. The two domains aNTD and aCTD of the RNA polymerase a subunit are shown shaded.



The  $\alpha$ CTD is connected to the rest of RNA polymerase by a flexible hinge. Note the differences in the topography of the two cases.



## 5.2. Repression of $P_1$

GalR occupying DNA at position  $-60.5$  represses  $P_1$ , not by hindering RNA polymerase binding, but by contacting  $\alpha$ CTD, which binds to  $P_1$  at position 45 bp upstream of  $P_1$  (10, 27, 37) (Fig. 4). GalR inhibits isomerization of RNA polymerase complex at  $P_1$ . How contacts between the same two proteins bring about opposite effects at the two promoters remains to be determined. It is not known why and under what conditions the dual behavior of GalR toward  $P_1$  and  $P_2$  in the absence of DNA looping is triggered in cells.

## 6. Gal Isorepressor

The *gal* operon is also regulated by an isorepressor (GalS). Although GalS does not seem to repress the *gal* promoters by DNA looping, the isorepressor does stimulate  $P_2$  and repress  $P_1$ , the same way that GalR does by binding to  $O_E$ , except that the effects are weaker. GalR and GalS modulate a few other operons, including those encoding both high- and low-affinity galactose **active transport** systems (12, 13, 38). The degree of regulation by GalR and GalS varies from operon to operon, perhaps to coordinate galactose metabolism and to transport efficiently under a wide range of galactose availability.

## 7. Properties of GalR and GalS

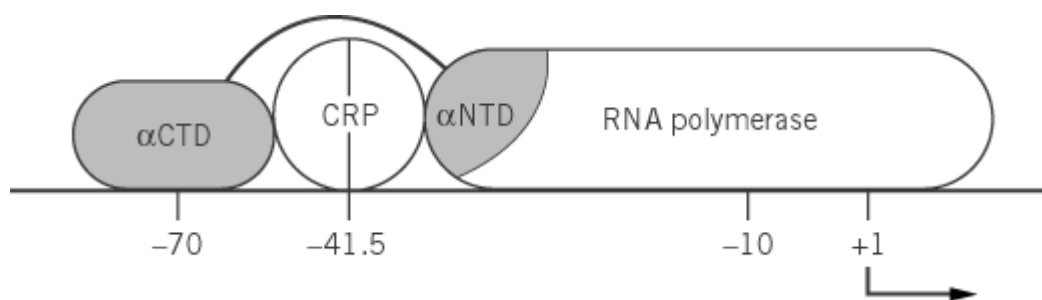
GalR and GalS are 85% similar in their amino acid sequences (12). By their homology with proteins (GalR-LacI family) whose structures are known, GalR and GalS are expected to have two-domains connected by a flexible hinge (39). A **helix–turn–helix motif** in the amino domain of each subunit of GalR and GalS dimers recognizes half of a dyad symmetry in  $O_E$  and  $O_I$ . The carboxy domains contain the inducer binding sites, which have been defined by isolation and characterization of inducer nonbinding (noninducible) repressor mutants, such as *galR<sup>S</sup>* and by modeling (40, 41). The mechanism by which the inducer derepresses the  $P_1$  promoter under conditions of nonlooping by binding to the  $O_E$ -bound GalR has been studied in detail. Despite the demonstration that inducer binding can dissociate GalR from the operator,  $P_1$  can be derepressed under conditions when  $O_E$  is occupied by a GalR-inducer complex (42). These results show that dissociation of a repressor from

an operator is not obligatory for transcription. Since the  $O_E$ -bound GalR interacts with RNA polymerase to inhibit open complex formation at  $P1$  at a postbinding step, the inducer acts by allosterically neutralizing the inhibitory contact between the proteins without dissociating the repressor from DNA (10, 17).

## 8. Regulation by cAMP•CRP

Like GalR and GalS, the global regulator cAMP•CRP complex also has differential effects on the *gal* promoters; but unlike GalR and GalS, the complex activates transcription from  $P1$  (three- to fourfold) and represses the same from  $P2$  (10-fold) (14, 16). Unlike a large group of promoters, such as the *lac* promoter, in which cAMP•CRP dimer activates transcription by binding at position  $-61.5$  on DNA and contacting RNA polymerase through the  $\alpha$ CTD (see [Lac Operon](#) and [Cyclic Amp Receptor Protein \(CRP\)/Catabolite Gene Activator Protein \(CAP\)](#)), the regulatory complex brings about dual control in the *gal* operon by binding to DNA at position  $-41.5$  (37, 43, 44) (Fig. 5). Whereas it is believed that cAMP•CRP represses  $P2$  by sterically hindering RNA polymerase binding to the overlapping  $-35$  element of the promoter, the molecular mechanism of activation of  $P1$  is different from that in the *lac* promoter. cAMP•CRP stimulates transcription initiation at  $P1$  by stimulating both RNA polymerase binding (closed-complex formation) and isomerization (45). A patch of amino acid residues (called *region 1*) of the promoter distal subunit of the cAMP•CRP dimer interacts with an  [\$\alpha\$ -helix](#) (helix 1) of the  $\alpha$ CTD and helps the latter to stretch away from the  $N$ -terminal domain ( $\alpha$ NTD) to bind to the region upstream of cAMP•CRP. A different amino acid residue patch (called region 2) in the promoter proximal subunit of cAMP•CRP interacts with a different segment of the  $\alpha$ NTD. The contact with the  $\alpha$ CTD is responsible for increasing the RNA polymerase binding, whereas the interaction with the  $\alpha$ NTD stimulates the isomerization step. This shows how the same regulatory protein can activate transcription initiation at two different biochemical steps by making entirely different contacts with RNA polymerase. Nevertheless, the dual roles of cAMP•CRP enable the *gal* operon to be expressed primarily from  $P1$  in cAMP-proficient cells and from  $P2$  in cAMP-deficient (eg, glucose grown) cells.

**Figure 5.** Activation of  $P1$  by cAMP•CRP by two different contacts with RNA polymerase. The two domains  $\alpha$ NTD and  $\alpha$ CTD of the RNA polymerase a subunit are shown shaded. Details are discussed in the text.



## 9. Summary

The study of the *gal* operon in *E. coli* has shown new ways and means of gene regulation at the level of transcription initiation: (i) how proteins bound to spatially separated sites on DNA communicate by DNA looping; (ii) how specific and nonspecific DNA binding proteins cooperate to “condense” DNA, making the latter unavailable for transcription, while remaining sensitive to an inducing signal; (iii) how the same regulator brings about opposite effects on promoters — activation and repression—by making direct contact with RNA polymerase; and (iv) how an inducer allosterically changes a repressor to neutralize the inhibitory effect of repressor still bound to DNA. Besides studying the detailed biochemical mechanisms of such controls, the two interesting questions

remain: What are the physiological reasons for the multitude of controls, and how are these diverse controls coordinated in the cell?

## Bibliography

“*gal* Operon” in , Vol. 2, pp. 952–956, by Sankar Adhya, National Cancer Institute, Developmental Genetics Section, Laboratory of Molecular Biology, 37 Convent Dr MSC 4264, Rm 5138, Bethesda MD 20892-4264, USA, (301)496-5138, [sadhya@helix.nih.gov](mailto:sadhya@helix.nih.gov); “*gal* Operon” in (online), posting date: January 15, 2002, by Sankar Adhya, National Cancer Institute, Developmental Genetics Section, Laboratory of Molecular Biology, 37 Convent Dr MSC 4264, Rm 5138, Bethesda, MD 20892-4264, USA, (301)496-5138, [sadhya@helix.nih.gov](mailto:sadhya@helix.nih.gov).

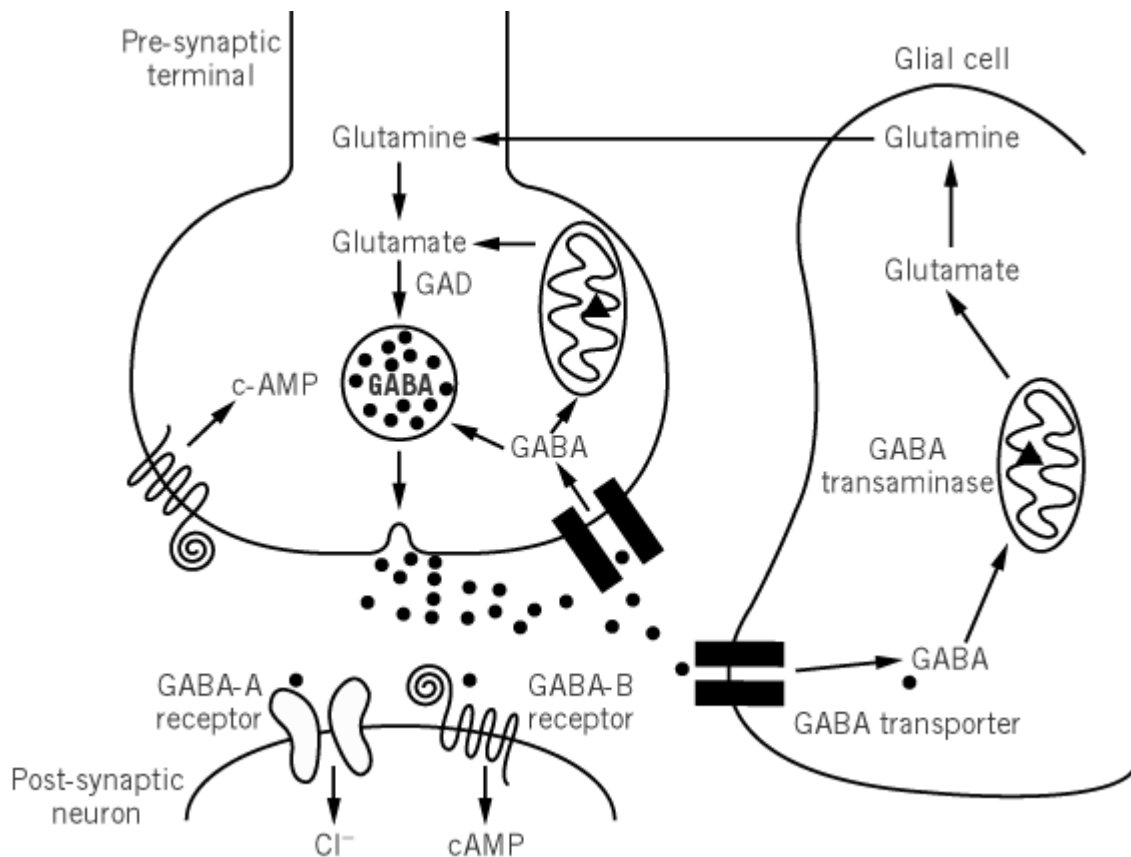
1. H. M. Kalckar (1958) *Adv. Enzymol.* **20**, 111–134.
2. G. Bouffard, K. Rudd, and S. Adhya (1994) *J. Mol. Biol.* **244**, 269–278.
3. J. R. Sherman and J. Adler (1963) *J. Biol. Chem.* **238**, 873–878.
4. S. Adhya and W. Miller (1979) *Nature* **279**, 492–494.
5. R. Musso et al. (1977) *Cell* **12**, 847–854.
6. G. Buttin (1963) *J. Mol. Biol.* **7**, 183–205.
7. T. Aki, H. E. Choy, and S. Adhya (1996) *Genes to Cells* **1**, 179–188.
8. T. Aki and S. Adhya (1997) *EMBO J.* **16**, 3666–3674.
9. G. Buttin (1963) *J. Mol. Biol.* **7**, 164–182.
10. H. E. Choy et al. (1995) *EMBO J.* **14**, 4523–4529.
11. J. P. E. Tokeson, S. Garges, and S. Adhya (1991) *J. Bacteriol.* **173**, 2319–2327.
12. M. J. Weickert and S. Adhya (1992) *J. Mol. Biol.* **226**, 69–83.
13. M. Geanacopoulos and S. Adhya (1997) *J. Bacteriol.* **179**, 228–234.
14. H. Aiba, S. Adhya, and B. deCrombrughe (1981) *J. Biol. Chem.* **256**, 11905–11910.
15. S. Busby, H. Aiba, and B. deCrombrughe (1982) *J. Mol. Biol.* **154**, 211–227.
16. H. E. Choy and S. Adhya (1993) *Proc. Natl. Acad. Sci. USA* **96**, 472–476.
17. M. Lavigne, H. Herbert, A. Kolb, and H. Buc (1992) *J. Mol. Biol.* **224**, 293–306.
18. S. Adhya et al. (1993) *Gene* **132**, 1–6.
19. D. J. Jin (1994) *J. Biol. Chem.* **269**, 17221–17227.
20. D. J. Jin, C. L. Turnbough Jr. (1994) *J. Mol. Biol.* **236**, 72–80.
21. R. diLauro et al. (1979) *Nature (London)* **279**, 494–500.
22. M. Irani, L. Orosz, and S. Adhya (1983) *Cell* **32**, 783–788.
23. H. J. Fritz et al. (1983) *EMBO J.* **2**, 2129–2135.
24. Y. Lyubchenko et al. (1997) *Nucleic Acids Res.* **25**, 873–876.
25. H. E. Choy et al. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7327–7331.
26. S. Adhya, M. Geanacopoulos, D. E. A. Lewis, S. Roy, and T. Aki (1998) *Cold Spring Harbor Symp. Quant. Biol.* **62**, 1–10.
27. H. E. Choy and S. Adhya (1992) *Proc. Natl. Acad. Sci. USA* **89**, 11264–11268.
28. M. Geanacopoulos, G. Vasmatazis, D. E. A. Lewis, S. Roy, B.-K. Lee, and S. Adhya (1999) *Genes Dev.* **13**, 1251–1262.
29. M. Geanacopoulos, G. Vasmatazis, V. B. Zhurkin, and S. Adhya (2001) *Nature Struct. Biol.* **8**, 432–436.
30. D. E. A. Lewis, M. Geanacopoulos, and S. Adhya (1999) *Mol. Microbiol.* **31**, 451–461.
31. S. Kar and S. Adhya (2001) *Genes Dev.* **15**, 2273–2281.
32. R. Ebright (1993) *Mol. Microbiol.* **8**, 797–802.
33. A. Ishihama (1993) *J. Bacteriol.* **175**, 2483–2489.

34. H. Tang et al. (1994) *Genes Develop.* **8**, 3058–3067.
35. E. Blatter et al. (1994) *Cell* **78**, 889–896.
36. Y. H. Jeon et al. (1995) *Science* **270**, 1495–1497.
37. T. A. Belyaeva et al. (1996) *Nucleic Acids Res.* **24**, 2243–2251.
38. M. J. Weickert and S. Adhya (1993) *Mol. Microbiol.* **10**, 245–251.
39. M. J. Weickert and S. Adhya (1992) *J. Biol. Chem.* **267**, 15869–15874.
40. H. Saedler et al. (1968) *Mol. Gen. Genet.* **102**, 79–88.
41. Y.-N. Zhou et al. (1995) *J. Mol. Biol.* **253**, 414–425.
42. S. Chatterjee et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2957–2062.
43. Y. Zhou, T. J. Merkel, and R. H. Ebright (1994) *J. Mol. Biol.* **243**, 603–610.
44. W. Niu et al. (1996) *Cell* **87**, 1123–1134.
45. M. Herbert, A. Kolb, and H. Buc (1986) *Proc. Natl. Acad. Sci. USA* **83**, 2807–2811.

## Gamma-Aminobutyric Acid (GABA)

g-Aminobutyric acid (GABA) is widely found in both vertebrates and invertebrates and in both neuronal and non-neuronal tissues. By far, its most understood and arguably most important role is as a neurotransmitter (Fig. 1). GABA is synthesised from glutamate by the enzyme glutamic acid decarboxylase (GAD). It is stored in vesicles within the synaptic terminal of neurons. When an action potential comes down the axon and reaches the synaptic terminal, calcium channels are activated, resulting in calcium influx. This leads to fusion of the neurotransmitter-containing vesicles with the cell membrane and release of the neurotransmitter, in this case GABA, into the synaptic cleft. The GABA diffuses across the synaptic gap and binds to its receptor. GABA receptors are of two types, ionotropic (GABA<sub>A</sub>) or metabotropic (GABA<sub>B</sub>), and their activation leads to hyperpolarization of the postsynaptic membrane and a decrease in the excitability of the cell. The response is terminated by desensitization of the receptors and by removal of the GABA from the synapse into the presynaptic terminal or surrounding glial cells. The latter is achieved through specific transporter proteins referred to as GABA transporters. Once taken up into a cell, two potential fates await it. The GABA can be taken up into mitochondria and metabolized to succinic semialdehyde (by GABA transaminase), and then to succinic acid, and enter the tricarboxylic acid pathway. This is referred to as the *GABA shunt*. Alternatively, in the presynaptic terminal, it can be recycled into synaptic vesicles and, thereby, made available for subsequent release.

**Figure 1.** A schematic of the vertebrate GABAergic synapse. GABA (filled circles) is generated from glutamate by glutamic acid decarboxylase (GAD). When released from presynaptic vesicles into the synaptic cleft, it diffuses across and binds to postsynaptic GABA<sub>A</sub> and GABA<sub>B</sub> receptors. It may also bind to presynaptic GABA<sub>B</sub> receptors. GABA is removed from the synaptic cleft into surrounding glial cells or the presynaptic terminal by GABA transporters. It is directly recycled into synaptic vesicles or taken up by mitochondria, converted by GABA transaminase (filled triangle) to succinic semialdehyde, and enters the tricarboxylic acid pathway.



## 1. Glutamic Acid Decarboxylase

GAD converts glutamic acid into GABA. In mammals there are two forms, GAD65 and GAD67 (named according to their apparent molecular weights), which are the products of two different genes (1). These two forms of GAD are expressed throughout the nervous system and also in several non-neuronal tissues including the pancreas, oviduct, and the testis. Most GABAergic neurons in the brain express both forms of GAD, but their relative abundance varies (2). GAD67 is distributed widely throughout the neuron, while GAD65 tends to be concentrated at the axon terminals (3). The role of GAD in non-neuronal tissues is, in general, poorly understood. Relatively high levels of GABA and GAD are found in the islets of Langerhans beta cells of the pancreas (4). These cells are destroyed in the autoimmune disease insulin-dependent diabetes mellitus, and one of the major antigens is GAD. GAD has some amino acid sequence identity with a coxsackievirus polypeptide, and it has been proposed that IDDM may result from molecular mimicry between the virus and GAD (5).

## 2. GABA Receptors

### 2.1. GABA<sub>A</sub> Receptors

GABA<sub>A</sub> receptors are ligand-gated ion channels. When GABA binds to specific recognition sites on the receptor, this triggers a conformational change leading to opening of the intrinsic anion channel allowing chloride ions to flow through into the cell. This opening of the ion channel is transient (10s of milliseconds). Two molecules of GABA must bind to open the ion channel. GABA<sub>A</sub> receptors are members of a larger ligand gated ion channel family, which contains nicotinic acetylcholine receptors, strychnine-sensitive glycine receptors, and 5-HT<sub>3</sub> receptors. To date, the mammalian GABA<sub>A</sub> receptor gene family consists of 15 subunits (α1–α6, β1–β6, γ1–γ3, δ, ε, ρ)(6, 7). β4 and γ4 have been identified in avian brain (8). Three other subunits, ρ1–ρ3, have also been identified. They

are expressed primarily in the retina and have a pharmacology that differs from GABA<sub>A</sub> receptors. As such, there are suggestions that they be classified as GABA<sub>C</sub> receptors (9). GABA<sub>A</sub> receptor subunits are 450 to 550 amino acids in length and have the structural features of the ligand-gated ion channel gene family: a putative signal peptide, a putative large extracellular domain, four putative transmembrane domains, the second of which is thought to line the pore of the channel. The subunits are classified by their relative deduced amino acid sequence homologies; within a group (eg, the a subunits), the subunits are approximately 70% homologous, whereas between groups (eg, a versus b) they are approximately 40% homologous. All the subunits have unique distributions of expression in the brain, and this expression is developmentally controlled (10). The native receptor is formed from the coassembly of the subunits in various combinations, resulting in a family of receptor subtypes. The precise subunit composition of receptor subtypes is not yet known. However, it is likely that most receptors are formed from the coassembly of a and b subunits with either a g, d, or e (or both g and q) subunit (6, 7). The receptor is a pentamer, the subunit stoichiometry of which is probably (a)<sub>2</sub>(b)<sub>2</sub>(g) or (a)<sub>2</sub>(b)(g)<sub>2</sub>(11, 12)).

GABA<sub>A</sub> receptors are the site of action of a number of clinically important drugs including benzodiazepines (prescribed as anxiolytics, anticonvulsants, and sedatives), barbiturates, and general anaesthetics, all of which act through unique allosteric modulatory sites on the receptor, potentiating the action of GABA in opening the channel. The use of transgenic, so called “knock-in” mice has been used to delineate which GABA<sub>A</sub> receptor subtypes are responsible for the anxiolytic properties or the sedative properties of benzodiazepines (13, 14).

GABA<sub>A</sub> receptors are also found in insects, at the neuromuscular junction and the nervous system. The subunits that have been identified (eg, *Rdl*) are homologous to vertebrate GABA<sub>A</sub> receptor subunits (15). They are the sites of action of a number of insecticides, including dieldrin. Sequencing of the *C. elegans* genome has indicated that the family of putative GABA<sub>A</sub> receptor subunits may be quite large (16).

## 2.2. GABA<sub>B</sub> Receptors

These are metabotropic receptors, interacting with a guanine nucleotide binding protein (G-protein) to produce its effect, which are slower than those of the ionotropic GABA<sub>A</sub> receptor (17). Activation of the GABA<sub>B</sub> receptor can lead to inhibition (or activation) of adenylyl cyclase, inhibition of voltage-gated calcium channels or activation of potassium channels (18). There is evidence that GABA<sub>B</sub> receptors can be both postsynaptic and presynaptic. The postsynaptic GABA<sub>B</sub> receptors, as compared with GABA<sub>A</sub> receptors, produce a slow and long lasting inhibition. Presynaptic receptors may act as autoreceptors, acting through a feedback mechanism to reduce the release of neurotransmitter (19). The first cDNA encoding the GABA<sub>B</sub> receptor was reported in 1997 (20) and is referred to as GABA<sub>B</sub>1. Subsequently, a homologous cDNA was identified (GABA<sub>B</sub>2), which heterodimerizes with GABA<sub>B</sub>1, and, indeed, the current consensus of opinion is that native GABA<sub>B</sub> receptors exist as such heterodimers (21). The receptor polypeptides have seven putative transmembrane domains and have significant amino acid sequence homology with metabotropic glutamate receptor gene family. There is no significant sequence homology with GABA<sub>A</sub> receptors or with other G-protein coupled receptors.

Baclofen, a GABA<sub>B</sub> receptor agonist, is used in the clinic to treat the spasticity that can result from spinal injury and also multiple sclerosis (22).

## 3. GABA Transporters

GABA transporters are located primarily on presynaptic terminals and surrounding glia, where they have three main functions: to regulate the concentration and duration of GABA in the synaptic cleft; to prevent the diffusion of the GABA to surrounding synapses; and to take up the GABA, thereby allowing for recycling or metabolizing. They are transmembrane proteins that use sodium and chloride ion cotransport to allow uptake of the GABA against the electrochemical gradient (23). To date, four subtypes of GABA transporter have been identified, GAT-1 (24), GAT-2 (25), GAT-3 (25), and BGT-1 (26). They are members of a larger gene family of neurotransmitter transporters including those for 5-HT, epinephrine, glycine, taurine, proline, and dopamine. The GABA transporters are approximately 600 amino acids long with 12 putative transmembrane domains and exhibit 50% to 70% amino acid sequence identity with each other. They have unique distributions in the brain, with GAT-1 being the most abundant and widespread and GAT-3 being the least abundant, with expression apparently restricted to the leptomeninges (27).

A number of compounds are known that inhibit GABA uptake by acting on GABA transporters. These include Tiagabin, which is an anti-epileptic (28).

### Bibliography

1. M. G. Erlander, N. J. K. Tillakaratne, S. Feldblum, N. Patel, and A. J. Tobin (1991) *Neuron* **7**, 91–100.
2. S. Feldblum, M. G. Erlander, and A. J. Tobin (1993) *J. Neurosci. Res.* **34**, 689–706.
3. M. Esclapez, N. J. K. Tillakaratne, D. L. Kaufman and A. J. Tobin (1994) *J. Neurosci.* **14**, 1834–1855.
4. R. L. Sorenson, D. G. Garry, and T.C. Brelji (1991) *Diabetes* **40**, 1365–1374.
5. D. L. Kaufman, M. G. Erlander, M. Clare-Satzler, M. A. Atkinson, et al (1992) *J. Clin. Invest.* **89**, 283–292.
6. E. A. Barnard, P. Skolnick, R. W. Olsen, H. Mohler, et al (1998) *Pharmacol. Rev.* **50**, 291–313.
7. T. P. Bonnert, R. M. McKernan, S. Farrar, B. le Bourdelles, et al (1999) *Proc. Natl. Acad. Sci. U.S.A.* **96**, 9891–9896.
8. M. G. Darlison and B. E. Albrecht (1995) *Semin. Neurosci.* **7**, 15–126.
9. G. A. R. Johnston (1996) *Trend Pharmacol. Sci.* **17**, 319–323.
10. W. Wisden, D. J. Laurie, H. M. Monyer, and P. H. Seeburg (1992) *J. Neurosci.* **12**, 1040–1062.
11. K. H. Backus, M. Arigoni, U. Drescher, L. Scheurer, et al (1993) *NeuroReport* **5**, 285–288.
12. Y. Chang, R. Wang, S. Barot, and D.S. Weiss (1996) *J. Neurosci.* **16**, 5415–5424.
13. U. Rudolph, F. Crestani, D. Benke, B. Brunig, et al (1999) *Nature* **401**, 796–800.
14. R. M. McKernan, T. W. Rosahl, D. S. Reynolds, C. Sur, et al (2000) *Natl. Neurosci.* **3**, 587–592.
15. R. H. French-Constant, D. P. Mortlock, C. D. Shafer, R. J. MacIntyre, et al (1991) *Proc. Natl. Acad. Sci. U.S.A.* **88**, 7209–7213.
16. The *C. elegans* Sequencing Consortium (1998) *Science* **282**, 2012–2018.
17. N. G. Bowery (1993) *Ann. Rev. Pharmacol.* **33**, 109–147.
18. P. Dutar and R. A. Nicoll (1988) *Nature* **332**, 156–158.
19. S. Thompson, M. Capogna, and M. Scanziani (1993) *Trends Neurosci.* **16**, 222–226.
20. K. Kaupmann, K. Huggel, J. Heid, P. J. Flor, et al (1997) *Nature* **386**, 239–246.
21. F. Marshall, K. Jones, K. Kaupmann, and B. Bettler (1999) *Trends Pharmacol. Sci.* **20**, 396–399.
22. H. Bittiger, W. Froestl, S. Mickel, and H. R. Olpe (1993) *Trends Pharmacol. Sci.* **1**, 391–394.
23. S. Mager, J. Naeve, M. Quick, C. Labarca, et al (1993) *Neuron* **10**, 177–188.
24. J. Guastella, N. Nelson, J. J. Nelson, L. Czyzyk, et al (1990) *Science* **249**, 1303–1306.
25. L. A. Borden, K. E. Smith, P. R. Hartig, T. A. Branchek, and R. I. Weinshank (1992) *J. Biol.*

Chem. **267**, 21098–21104.

26. A. Yamauchi, S. Uchida, H. M. Kwon, A. S. Preston, et al (1992) *J. Biol. Chem.* **267**, 649–652.
27. M. M. Durkin, E. L. Gustafson, K. E. Smith, L. A. Borden, et al (1995) *Mol. Brain Res.* **3**, 7–21.
28. P. M. Crawford, M. Engelsman, S. W. Brown, T. W. Rentmeister, et al (1993) *Epilepsia* **34**, Suppl. 2, 182.

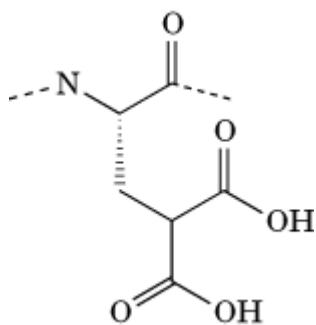
### Suggestions for Further Reading

29. N. J. K. Tillakaratne, L. Medina-Kauwe, and K.M. Gibson (1995) Gamma-aminobutyric acid metabolism in mammalian neural and nonneural tissues. *Comp. Biochem. Physiol.* **112A**, 247–263.
30. L. A. Borden (1996) GABA transporter heterogeneity: pharmacology and cellular localisation. *Neurochem. Int.* **29**, 335–356.
31. D. L. Martin and R. W. Olsen (2001) *GABA in the Central Nervous System*. Lippincott Williams and Wilkins, Philadelphia.
32. H. Mohler (2001) *Pharmacology of GABA and Glycine Neurotransmission*. Springer, Berlin.

## Gamma-Carboxyglutamic Acid

g-Carboxyglutamic acid residues are found within the highly homologous *N*-terminal domains of various vitamin K–dependent proteins (see [Calcium-Binding Proteins](#)). Such residues are generally abbreviated as “Gla”; their structure is depicted in Figure 1. Gla residues are the result of a post-translational modification of glutamic acid residues by a vitamin K-dependent carboxylase, also known as g-glutamyl carboxylase, which introduces a second carboxyl group into the side chain. The two carboxyl groups give this amino acid considerable affinity for cations, especially  $\text{Ca}^{2+}$ , as in the chelators [EDTA](#) and [EGTA](#). This modification is required for the calcium-mediated interaction of the Gla proteins with the negatively charged phospholipid [membrane](#) and is also essential for the formation of the native protein conformation.

**Figure 1.** Molecular structure of a g-carboxyglutamic acid residue.



### Suggestions for Further Reading

- J. W. Suttie (1993) *FASEB J.* **7**, 445–452 (a concise review of the post-translational modification

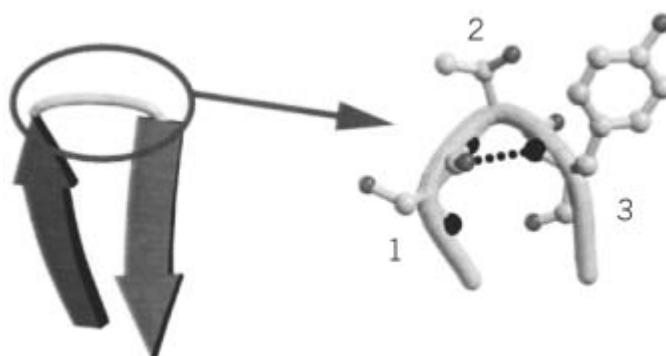


of Glu residues to Gla in vitamin K–dependent proteins).

## Gamma-Turns

A gamma-turn is a type of nonregular **secondary structure** in [protein structures](#) that causes a change in direction of the [polypeptide chain](#). In the g-turn, a [hydrogen bond](#) is formed between the [backbone](#) carbonyl oxygen of one residue ( $i$ ) and the backbone amide NH of the residue two positions further along the chain ( $i + 2$ ) (Fig. 1). This ( $i$ ) to ( $i + 2$ ) interaction distinguishes g-turns from [b-turns](#), which have an ( $i$ ) to ( $i + 3$ ) hydrogen bond. There are two types of g-turns, the classic and the inverse, which have different backbone dihedral angles for the ( $i + 1$ ) intervening residue. The two g-turn types are, however, related by an inversion of the signs of the  $\phi$  and  $\psi$  angles (classic:  $70^\circ$  to  $85^\circ$  and  $-60^\circ$  to  $-70^\circ$ , respectively; inverse:  $-70^\circ$  to  $-85^\circ$  and  $60^\circ$  to  $70^\circ$ , respectively) (see [Ramachandran Plot](#)), so their backbone conformations are mirror images of one another. This is analogous to the difference between the type I and I' and type II and II' b-turns. The classic g-turn is less common than the inverse type, because of the unfavorable  $\phi$  and  $\psi$  angles, and is found almost exclusively at the ends of **hairpins** (that is, the classic g-turn usually connects two adjacent antiparallel b-strands).

**Figure 1.** Schematic representation of a g-turn in a protein structure. **(Left)** The g-turn is shown connecting two b-strands in a hairpin motif. **(Right)** The detailed atomic structure of the g-turn is shown. Residues are numbered 1 to 3 for ( $i$ ) to ( $i + 2$ ). The hydrogen bond between the backbone carbonyl oxygen of residue 1 and the backbone amide nitrogen of residue 3 is shown as a dotted line. This example is a classic g-turn, because the backbone angles for the ( $i + 1$ ) residue are  $\phi+67^\circ$  and  $\psi-48^\circ$ . Oxygen atoms and nitrogen atoms are shown as dark spheres. This figure was generated by Molscript (1) and Raster3D (2, 3).



[See also [Secondary Structure, Protein, Turns, Beta-Turns](#) and [Omega Loop](#).]

### Bibliography

1. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
2. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
3. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

4. G. D. Rose, L. M. Gierasch, and J. A. Smith (1985) Turns in peptides and proteins. *Adv. Protein Chem.* **37**, 1–109.
5. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Gap Junction

Gap junctions are specialized [cell junctions](#), [membrane](#) plaques that provide [hydrophilic](#) channels for direct intercellular communication. They permit the passive [diffusion](#) of molecules of 1000 daltons or less between the cytoplasm of adjacent cells. They do not facilitate the transfer of large polymers, such as polypeptides, RNA, and DNA. Molecules that may pass through gap junctions are such as [cyclic AMP](#), [inositol phosphates](#), and [calcium](#) ions. They thus permit the intercellular transfer of cytoplasmic signals. Gap junctions are widespread in body tissues in both excitable and nonexcitable cells, and they are important in embryonic development.

Ultrastructurally, gap junctions appear as a pair of parallel unit plasma membranes separated by a uniform gap of 2 to 3 nm. Lanthanum infiltration reveals the presence of regularly arranged cross bridges connecting the plasma membranes (1). The junctions are composed of a regular hexagonal array of subunits, the connexons (see text below). In rat hepatocyte gap junctions, the lattice has a unit-cell dimension of 8.5 nm. The intercellular cross bridges are created by pairing of connexons between the plasma membranes of adjacent cells (2).

The connexon is the functional unit of gap junctions. This is a hexameric assembly of proteins called connexins (3) (Fig. 1). Thirteen different connexins are known from rodent tissues. The connexins range in molecular weight from 26,500 to 49,600 daltons (4). They are referred to either by a Greek nomenclature or on the basis of their molecular mass. Thus, connexin a1 is also known as Cx43, having a molecular weight of 43,000. Different connexins show greater or more specific expression in certain tissues and one, Cx38, is predominantly or exclusively expressed in the embryo. Each connexin has four helical transmembrane domains, both the amino and carboxy termini being located in the cytoplasm, and two cytoplasmic and one extracellular loops (5, 6). The third transmembrane domain is amphipathic in character, and it is believed that six of these domains from the subunits of the hexamer form the hydrophilic channel in an assembled connexon (6, 7). Two connexons in the plasma membranes of adjacent cells bind to each other by means of their extracellular loops, thus creating a hydrophilic channel between the cytoplasm of the two cells and allowing the passage of molecules of 1.5 nm or less in diameter. Formation of the gap junction intercellular channel requires a 30° rotation to interdigitate the two apposing connexons (8). Channels formed by a different connexins have differing permeability properties (9), and single-channel conductances vary from 20 to several hundred picosiemens (10). The permeability properties of gap junction channels may be modulated by **phosphorylation** and cytoplasmic calcium concentration and are voltage-dependent (10, 11). Change in the configuration of the subunits is believed to regulate the channel opening, either by blockage or by rotation, so that the channels may be either open or closed (12).

**Figure 1.** Gap junctions, connexons, and connexin. (a) The gap junction provides hydrophilic channels for intercellular communication. Connexons pair between adjacent cells, giving an intercellular gap of 2 nm with a ladder-like appearance to the intercellular space. (b) In cross section, each connexon appears as a hollow cylinder, the walls of which are made up of six connexin subunits. (c) Connexin is a molecule with four transmembrane domains: its NH<sub>2</sub>

and COOH termini inside are the cell, and there are two extracellular and one intracellular loop. (d) In the junction plaque, the connexons are arranged in a hexagonal array.

The expression of several connexons in the same individual cells enables the formation of different types of channels, either homotypic channels with identical connexins paired between cells, or heterotypic channels with connexons composed of different connexins paired between cells. Alternatively, individual connexons could be heteromeric; that is, they could be composed of mixtures of connexin subunits. There is evidence that this occurs in some situations, for example, with connexins Cx32 and Cx26 (13). On the other hand, some connexins appear to be incompatible, being unable to link via their extracellular domains to form a channel—for example, connexins Cx43 and Cx40, which are expressed in the heart Purkinje fibers and the ventricular myocardium, respectively (14). This expression pattern is believed to be important in regulating the electrical conducting properties within cardiac muscle. The full functional significance of the different connexin subunits is an important area for future investigation. Connexin mutations are involved in some human diseases. Cx32 mutations are involved in the X-linked form of Charcot–Marie–Tooth syndrome (15), Cx43 in viscerotrial heterotaxia syndrome (16) and Cx26 in hereditary nonsyndromic sensorineural deafness (16). Viscerotrial heterotaxia affects left- or right-sided development of the heart. The point mutations involved are believed to involve sites for connexin phosphorylation, implying that these are important for the regulation of connexin function. Null mutations of the connexin Cx43 in mice results in death shortly after birth due to swelling in and blockage of the right ventricular outflow of the heart (17). Other connexin null mutations cause equally specific defects. Thus, mice lacking Cx40 have cardiac conduction abnormalities resembling atrioventricular and bundle branch block (18), absence of Cx37 causes female infertility because failure of oocyte development (19), and mice lacking Cx46 develop cataracts because of proteolysis of lens crystallins (20). Signals transmitted through gap junctions have been implicated in cellular growth control and vice versa in tumorigenic cell lines, consistent with a tumor suppressor function for gap junctions (21). The importance of gap junctional communication in developmental signaling is demonstrated by the severe defects that result from blocking it (22).

## Bibliography

1. J.-P. Revel and M. Karnovsky (1967) *J. Cell Biol.* **33**, 7–12.
2. D. A. Goodenough and N. B. Gilula (1974) *J. Cell Biol.* **61**, 575–590.
3. E. L. Hertzberg and N. B. Gilula (1979) *J. Biol. Chem.* **254**, 2138–2147.
4. N. M. Kumar and N. B. Gilula (1996) *Cell* **84**, 381–388.
5. D. L. Paul (1986) *J. Cell Biol.* **103**, 123–134.
6. L. C. Milks et al. (1986) *EMBO J.* **7**, 2967–2975.
7. D. A. Goodenough et al. (1988) *J. Cell Biol.* **107**, 1817–1824.
8. G. A. Perkins et al. (1998) *J. Mol. Biol.* **877**, 171–177.
9. J. L. Brissette et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6453–6457.
10. D. C. Spray (1994). In *Molecular Mechanisms of Epithelial Cell Junctions: From Development to Disease* (S. Citi, ed.), R. G. Landes Company, Austin, TX, pp. 195–215.
11. L. S. Musil et al. (1990) *J. Cell Biol.* **111**, 2077–2088.
12. P. N. T. Unwin (1987) In “Junctional complexes of epithelial cells.” *Ciba Found. Symp.* **125**, 78–91.
13. K. A. Stauffer (1995) *J. Biol. Chem.* **270**, 6768–6772.
14. R. Bruzzone et al. (1993) *Mol. Biol. Cell* **4**, 7–20.
15. J. Bergoffen et al. (1993) *Science* **262**, 2039–2042.
16. S. H. Britz-Cunningham et al. (1995) *N. Engl. J. Med.* **332**, 1323–1329.

17. D. P. Kelsell et al. (1997) *Nature* **387**, 80–83.
18. A. M. Simon et al. (1998) *Curr. Biol.* **8**, 295–298.
19. A. M. Simon et al. (1997) *Nature* **385**, 525–529.
20. X. Gong et al. (1997) *Cell* **91**, 833–843.
21. S. W. Lee et al. (1992) *J. Cell Biol.* **8**, 295–298.
22. A. E. Warner et al. (1984) *Nature* **311**, 127–133.

### Suggestions for Further Reading

23. D. A. Goodenough, J. A. Coliger, and D. L. Paul (1996) Connexins, connexons and intercellular communication. *Annu. Rev. Biochem.* **65**, 475–502. (An excellent review.)
24. R. Bruzzone, T. W. White, and D. A. Goodenough (1996) *Bioessays* **18**, 709–718. (Shorter, but also an excellent review.)

### Gap Penalty

Aligned **homologous** protein and **nucleotide sequences** may include gaps inserted in one or more sequences to maximize the sequence similarity. These gaps are called [indels](#) and correspond to positions where insertion or deletion mutations have occurred. Superposition of homologous protein structures may reveal quite accurately the sites of indels. However, algorithms that **align sequences** must place the gaps at positions that maximize the score of the sequence alignment ([1](#)). Gap penalties are used to limit the number and length of gaps that can be opened in the sequences. The usual cost for inserting a gap is defined as  $\text{Cost} = \text{GOP} + \text{Gaplength} * \text{GEP}$ . Gaps scored in this way are called “affine gaps.” The two parameters in the equation are GOP, the gap-opening penalty, and GEP, the gap-extension penalty. The actual values of the penalties used must be in balance with the scores for aligned sequence segments. The values depend on the sequence nature and length, the number and length of genuine gaps, and the residue exchange values in the mutation matrices that are used to score substitutions. The penalties are usually calibrated by trial and error, although it should be possible to infer a likely value from the residue exchange values (while a formal hidden Markov model (HMM) should use actual gap insertion probabilities). In certain situations, more is known about where gaps are likely to occur—for example, if a [tertiary structure](#) is known ([2](#)), or if an existing alignment is being used that already has gaps in it. In this case, gap penalties can be varied in a position-specific manner, ie, lowered at sites where gap insertion is more likely. Position-specific gap penalties are applied in alignment-based profile or HMM [database](#) searches ([3](#), [4](#)) and in multiple alignments—for example, by the program Clustal W ([5](#)).

A variant form of gap penalty, called a “frame penalty,” is applied in some algorithms that compare a protein sequence to the three translation [reading frames](#) of a nucleic acid strand ([6](#)). As well as gap penalties for in-frame gaps, shifts in the reading frame are given an analogous cost composed of a frame-opening penalty (FOP) and a frame-extension penalty (FEP). Customization of the frame penalties can optimize alignment where errors in sequence determination result in [frameshifting](#), or where frameshifts affecting coding sequences arise naturally due to **introns**, [RNA Editing](#), or translational frameshifting. Introns can be crossed if the FEP is set very low. The [expressed sequence tag](#) (EST) database entries ([7](#)) have high error, and searches at the highest sensitivity require a relatively low FOP.

Optimization of gap penalties is essential for the good performance of sequence alignment

algorithms, and users should take the time to conduct their own trials. Default settings cannot be optimal for all situations.

### Bibliography

1. S. B. Needleman and C. D. Wunsch (1970) *J. Mol. Biol.* **48**, 443–453.
2. A. M. Lesk, M. Levitt, and C. Chothia *Protein Engng.* **1**, 77–78.
3. M. Gribskov, A. D. McLachlan, and D. Eisenberg (1987) *Proc. Natl. Acad. Sci. USA*, **84**, 4355–4358.
4. A. Krogh, M. Brown, S. Mian, K. Sjölander, and D. Haussler (1994) *J. Mol. Biol.* **235**, 1501–1531.
5. J. D. Thompson, D. G. Higgins, and T. J. Gibson (1994) *Nucleic Acids Res.* **22**, 4673–4680.
6. E. Birney, J. D. Thompson, and Toby J. Gibson (1996) *Nucleic Acids Res.* **24**, 2730–2739.
7. M. D. Adams et al. (1991) *Science* **252**, 1651–1656.

## Gas-Liquid Chromatography

Following the suggestion by Martin and Synge (1) that the mobile phase for partition [chromatography](#) could just as easily be a gas, James and Martin (2) introduced gas-liquid chromatography in 1952. In gas-liquid chromatography, the main mechanism of retention is the partition of an analyte, which is introduced into the mobile gas phase between the gas phase and a stationary liquid phase. The latter is deposited upon the surface of small solid particles or on the walls of a capillary column (capillary gas-liquid chromatography). The theory and application of capillary columns were first expounded by Golay (3). It was found that these columns had the advantages of low resistance to gas flow and high efficiencies equivalent to thousands of plates. Analysis of [lipids](#) is the best application of (capillary) gas-liquid chromatography in molecular biology. A detailed practical description is beyond the scope of this article, so interested readers are directed to the monograph *Gas Chromatography and Lipids: A Practical Guide* by Christie (4).

### Bibliography

1. A. J. P. Martin and R. L. M. Synge (1941) *Biochem. J.* **35**, 1358–1368.
2. A. T. James and A. J. P. Martin (1952) *Biochem. J.* **50**, 679–690.
3. M. J. E. Golay (1958) in *Gas Chromatography* (V. J. Coates, H. J. Noebels, and I. S. Fagerson, eds.) Academic Press, New York, pp. 1–13.
4. W. W. Christie (1989) *Gas Chromatography and Lipids: A Practical Guide*, The Oily Press, Ayr, Scotland.

## Gastrula, Gastrulation

The gastrula is a stage in early embryogenesis during which the process of gastrulation occurs. After [fertilization](#) and the cleavage divisions, a single layer of cells called the **blastula** forms. The blastula

may be hollow or filled with yolk. It can also consist of a large disc of cells on top of the embryo derived from large, yolky eggs [see [Meroblastic Cleavage](#)]. Regardless of the structure of the blastula, the process whereby the single cell layer gives rise to the three germ layers (endoderm, mesoderm, and ectoderm) is called gastrulation. Gastrulation usually involves changes in cell shape and cell migration. This cell migration can occur as single cells, as in sea urchin gastrulation, or the movement of sheets of cells, as in the involution of cells in fish embryos (a process called epiboly).

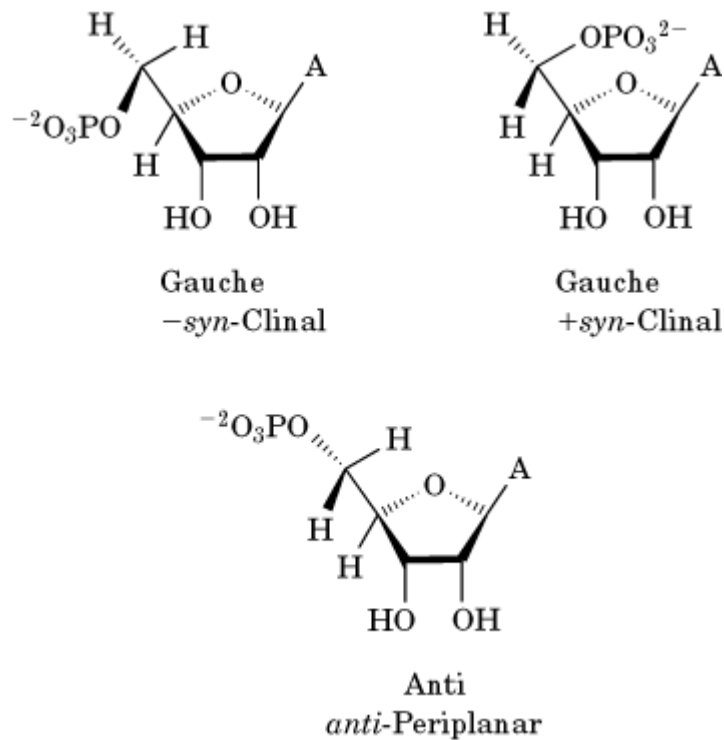
#### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 175–178.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
- J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.

### Gauche Conformation

The mathematical term gauche means not lying in a plane. In stereochemistry, specific [conformations](#) that are related by rotation about a single covalent bond between two tetrahedral atoms have been called gauche when none of the substituent atoms lie in the same plane (see [Torsion Angle](#)). Conformers where the substituents are staggered, and the torsion angle between the largest substituents on two neighboring carbons is about  $60^\circ$ , are termed gauche ([1](#)). These conformers are typically *syn* or *anti* clinal in the more rigorous designation of torsional conformers (see [Conformation](#)) ([2](#)). The alternative conformations are *anti* and eclipsed. The two gauche conformers and one *anti* conformer about the C4'—C5' ribose bond of AMP are shown in Figure [1](#).

**Figure 1.** The conformations of the three rotameric conformers of 5'—AMP. The first two are gauche because the torsion angle defined by O5'—C5'—C4'—C3' is between 0 and  $\pm 90^\circ$ . The third conformer is *anti* because this torsion angle is  $180^\circ$ .



### Bibliography

1. E. L. Eliel et al. (1967) *Conformational Analysis*, Wiley-Interscience, New York, p. 10.
2. W. Klyne and V. Prelog (1960) *Experientia* **16**, 521.

### Suggestion for Further Reading

3. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer Advanced Texts in Chemistry (C. R. Cantor, ed.), Springer-Verlag, New York, pp. 14–17.

### GC Box

A GC box is a *cis*-acting transcriptional regulatory element (see [Transcription Factors](#), Eukaryotic promoters) that contains the sequence 5'-GGGCGG-3' at its core and is bound by the [transcription factor](#) Sp1. GC boxes are found upstream of a number of genes transcribed by RNA polymerase II and increase the transcriptional rate of the gene above the level that could be achieved by the transcriptional core machinery in the absence of a GC box. Because Sp1 is a ubiquitous factor in eukaryotic cells, the presence of a GC box in a transcriptional regulatory region increases the level of expression of the gene. GC boxes may exist in one or multiple copies, depending on the promoter, and may be found in either orientation relative to the direction of transcription from the downstream initiation site. They are generally found upstream of the transcription initiation site and of the **TATA box**, if a TATA box is present. In some cases GC boxes, acting through the binding of Sp1, have an important role in assembling transcription complexes on promoters without an apparent TATA box.

## Gel Electrophoresis

The capacity of [electrophoresis](#) to resolve macromolecules differing in size, shape, or conformation (“size separation,” **molecular sieving**) is dramatically augmented when it is carried out in gels rather than simply in a buffer. This improvement in resolution arises from two causes: (1) The rate of change of electrophoretic mobility with gel concentration (the “retardation coefficient”) is directly related to the size and shape of the molecule (see [Ferguson Plot](#)); and (2) the convection and dispersion of a band of macromolecules are much reduced in gels relative to solutions. Since macromolecules (proteins, nucleic acids, and carbohydrates) are the building materials of biology, gel electrophoresis has become a ubiquitous separation method in the molecular biology laboratory. Two types of polymers capable of gelation are universally used: (1) [polyacrylamide](#) and (2) [agarose](#). The concentrations of each are selected to achieve effective molecular sieving by tailoring the pore architecture to the dimensions of the charged species to be separated.

One of the key advantages of gel electrophoresis over [chromatography](#) is the possibility of varying the gel concentration continuously so as to achieve any effective pore size within a continuum. The average pore size in a gel can range from a fraction of a nanometer (corresponding to the diameter of organic acids) to 50 to 100 nm (corresponding to the diameter of plant **viruses**) and up to nearly a micron. Another benefit of gel electrophoresis is its adaptability to gradients of pore sizes (see [Pore Gradient Electrophoresis](#)) and to the simultaneous analysis of samples containing large numbers of components: up to 30 to 50 in gels of a single dimension and up to several thousand in [two-dimensional gel electrophoresis](#).

Gel electrophoresis of multiple samples can be carried out simultaneously on multiple channels of vertically or horizontally oriented gel slabs, using either a single gel concentration or a concentration gradient. This application of gel electrophoresis allows one to compare the mobilities of bands in the various channels of the gel. Such a comparison can give information concerning the identity or nonidentity of two macromolecules, so long as the migration distances are unaffected by differences between the channels in conductance (sample ionic strength, Joule heat dissipation, or wall adherence). When the charge densities of species are equalized, as in the case of nucleic acids or with proteins coated with an ionic detergent (see [SDS-PAGE](#)), the rate of migration may depend on primarily the particle size, due to sieving by the gel. Additional and more reliable information concerning the properties of each species in a sample can be derived from gel electrophoresis at multiple gel concentrations (see [Ferguson Plot](#)).

To be effective, gel electrophoresis needs to be conducted at an optimized gel concentration and under optimized [buffer](#) conditions (see [Disc Electrophoresis](#)). The optimal resolving gel concentration is mathematically defined for any particular pair of charged species and can be derived from the Ferguson plots of those species. For large numbers of electrophoretic separations, their automation becomes important. Two kinds of apparatus automatically record and quantify the band pattern: (1) a gel electrophoresis apparatus with an **absorbance spectroscopy** detector at the end of the migration path that records the bands as they pass sequentially, and (2) a horizontal gel slab apparatus equipped with a mobile fluorescence detector that scans the migration path at desired intervals. The latter apparatus type has preparative capacity, which allows for sequential [electroelution](#) steps, with simultaneous computation of the recovery, to ensure that it is complete. It is applicable to polymer solutions (see [Particle Electrophoresis](#)), as well as gels.

Typically, gel electrophoresis zones are detected by fixation of the band in the gel and staining with a specific dye or an enzyme activity (see [Electrophoresis](#) and [Overlay Assay: Enzyme Zymography of Plasminogen Activators and Inhibitors](#)). Fixation can be replaced by [blotting](#), which permits additional methods of identification such as hybridization of DNA bands or the immunodetection of



proteins bands by **Western blots** .

#### Suggestions for Further Reading

A. Chrambach (1996) Quantitative and automated electrophoresis in sieving media: Past, present and future. *Electrophoresis* **17**, 454–464.

D. P. Goldenberg (1997) "Analysis of protein conformation by gel electrophoresis". In *Protein Structure: A Practical Approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, pp. 187–218.

D. Rodbard and A. Chrambach (1971) Estimation of molecular radius, free mobility and valence using polyacrylamide gel electrophoresis. *Anal. Biochem.* **40**, 95–134.

D. Tietz and A. Chrambach (1987) Computer simulation of the variable agarose fiber dimensions on the basis of mobility data derived from gel electrophoresis and using the Ogston theory. *Anal. Biochem.* **161**, 395–411.

D. Tietz and A. Chrambach (1992) Concave Ferguson plots of DNA fragments and convex Ferguson plots of bacteriophages: Evaluation of molecular and fiber properties, using desktop computers. *Electrophoresis* **13**, 286–294.

#### Gel Retardation Assay

The gel retardation assay, also known as the bandshift assay or electrophoretic mobility shift assay (EMSA), is used to detect [DNA-binding proteins](#) in crude cell extracts and to study the binding activity of purified DNA-binding proteins. It is also applicable to the study of RNA-protein interactions. The basis of the assay is that protein–DNA complexes remain intact when gently fractionated by gel electrophoresis and migrate as distinct bands more slowly than the free DNA fragment. The assay is simple and quick, and the use of radioactive binding-site DNA makes it highly sensitive. It can be used for both highly sequence-specific and nonspecific [proteins](#), and nonspecific [proteins](#), such as histones. Gel retardation assays can be used quantitatively to estimate dissociation constants for protein-DNA complexes. Additionally gel retardation experiments can be used to visualize protein–protein interactions between a DNA-binding protein and other non-DNA-binding proteins. The binding of a second protein to a protein-DNA complex to form a triple complex is visualized by a further retardation of mobility that is called a supershift. Supershift experiments can also be used to assay the binding of a second DNA-binding protein to the DNA or the binding of a second molecule of DNA to the protein. The assay is useful for proteins with a DNA association constant  $>10^7\text{M}^{-1}$ .

A further use of the gel retardation assay is to detect DNA bending induced by a DNA binding protein. This depends on the principle that when a protein binds to a DNA fragment of moderate length (in the range of 100–250 bp) the average distance between the ends of the DNA fragment is decreased and there is a greater resistance to the passage of the complex through the pores of the gel. The extent of the resulting retardation depends on the position of the protein binding to the ends of the fragment. The greatest retardation is observed when the binding site is positioned centrally. By comparison with proteins that induce a bend of known magnitude this method can be used to obtain an estimate of the DNA bend induced by a protein.

#### Suggestions for Further Reading

B. F. Luisi (1995) articles on "structure of DNA-binding motifs and protein-DNA complexes", A. A. Travers, articles on "DNA bending by sequence and proteins" and M. Bianchi (1995) articles on "HMG domain proteins" in D. M. J. Lilley, ed., *DNA-Protein: Structural Interactions*, Oxford

University Press, pp. 49–75.

P. A. Rice, S. Yang, K. Mizuuchi, and H. A. Nash (1996) Crystal Structure of an IHF-DNA Complex: a Protein-Induced U-Turn, *Cell*, **87**, 1295–1306.

Z. Otwinowski et al. (1988) “Crystal Structure of trp Repressor/Operator Complex at Atomic Resolution,” *Nature* **335**, 321–329.

Z. Shakked, et al. (1994) “Determinants of Repressor/Operator Recognition from the Structure of the trp Operator Binding Site,” *Nature*, **368**, 469–473.

J. G. Omichinski et al. (1997) “The Solution Structure of a Specific GAGA-factor-DNA Complex Reveals a Modular Binding Mode,” *Nat. Struct. Biol.* **4**, 122–132.

M. Buckle and A. Travers, eds., *Methods: Protein-DNA Interactions* Oxford University Press, in press.

## Gelsolin

Gelsolin is an abundant and ubiquitously expressed [actin-binding protein](#), which confers calcium sensitivity to the regulation of the dynamics of the [microfilament](#) system (1). There are also high concentrations of a gelsolin isoform in plasma. The latter is a 755-residue protein with an N-terminal extension of 25 [amino acid](#) residues. In the presence of micromolar concentrations of  $\text{Ca}^{2+}$ , gelsolin binds alongside an [actin](#) filament, severs it at that site, and then remains bound to the (+)-end on one of the resulting fragments. Gelsolin has two actin-binding sites, and in the presence of  $\text{Ca}^{2+}$  ions it nucleates filament formation from a solution of G-actin with subsequent (–)-end growth of the [polymer](#). From the minimal 2:1 complex of actin and gelsolin, one of the actin monomers can be released, if  $\text{Ca}^{2+}$  is removed by EGTA. The second, more tenaciously bound actin monomer can only be released by the binding of PtdIns 4,5-bisphosphate to the gelsolin. Therefore, both polyphosphoinositides and  $\text{Ca}^{2+}$  ions might be involved in the control of the effects of gelsolin on the turnover in the microfilament system. As shown by gene knockout experiments, gelsolin is not essential for survival. Only limited defects in motility processes have been detected, but overexpression of gelsolin in fibroblasts results in increased motile activity. Gelsolin has recently been shown to be a downstream effector of rac for fibroblast motility (2).

The structure of gelsolin has been studied extensively. It belongs to a family of actin-binding proteins that contain 1–6 repeats of 120–130 amino acid residues. Gelsolin has 6 such structurally related repeats (S1-S6), which may have arisen through gene triplication of a prototype gene, which subsequently went through [gene duplication](#). First, the three-dimensional structure of a complex between actin and subfragment 1 of the cytoplasmic form of gelsolin was solved by crystallography (3), and later the 3D-structure of the entire plasma gelsolin molecule (4). Recently, electron cryomicroscopy and helical reconstruction has been used to identify gelsolin binding sites on filamentous actin (5).

The different parts of the gelsolin molecule appear to have discrete functions. Subfragment 1 binds an actin monomer in the absence of  $\text{Ca}^{2+}$  and can cap actin filaments. Subfragment S1-S6 also binds actin monomers and competes for the same binding site on actin as fragment S1, but now the [binding](#) is  $\text{Ca}^{2+}$ -sensitive. The S2 domain by itself binds to the side of actin filaments in a  $\text{Ca}^{2+}$ -independent way. Subfragment S1-S3, can sever actin filaments. Nucleation of actin polymerization can be achieved *in vitro* with S2-S6. Thus, although gelsolin can bind to the side of the actin filament in the absence of  $\text{Ca}^{2+}$ , further action to sever the actin filament requires the involvement of the  $\text{Ca}^{2+}$ -

dependent actin binding site. Polyphosphoinositides bind to the S2 domain and dissociate the actin:gelsolin complexes. Despite efforts to determine the gelsolin binding site on F-actin by [cryoelectron microscopy](#) and helical reconstructions a convincing mechanism for the severing action of gelsolin has not yet been obtained.

### Bibliography

1. H. L. Yin (1987) *Bioessays* **7**, 176–179.
2. T. Azuma, W. Witke, T. P. Stossel, J. H. Hartwig, and D. J. Kwiatkowski (1998) *EMBO J* **17**, 1362–1370.
3. P. J. McLaughlin, J. T. Gooch, H.-G. Mannherz, and A. G. Weeds (1993) *Nature* **364**, 685–692.
4. L. D. Burtnick, E. K. Koepf, J. Grimes, E. Y. Jones, D. I. Stuart, P. J. McLaughlin, and R. C. Robinson (1997) *Cell* **90**, 661–670.
5. A. McGough, W. Chiu, and M. Way (1998) *Biophys. J.* **74**, 764–772.

## Gene Amplification

A fundamental property of living cells is their orderly transmission of genetic information from generation to generation. One aspect of this property involves a mechanism that controls **genomic** replication and ensures one complete doubling of each replicon during each cell generation. Another aspect of this property is the placement of genes so that they are expressed properly and so that each daughter cell receives the genes in the appropriate configuration. Now it is appreciated that deviations from this principle occur commonly and that amplification of DNA sequences and rearrangement of sequences occur often.

Amplification of DNA sequences, the differential increase in a specific portion of the genome compared with the remainder, occurs during development and during the vegetative growth of cells. The processes of **polyploidization** and endoreduplication, where the entire chromosome complement is multiplied in one nucleus, is not discussed here. Additionally, [aneuploidy](#) (or [trisomy](#)) can also result in a differential increase in a portion of the genome, but this is distinct from DNA amplification and is not discussed further.

Developmental amplification has been documented extensively in both germ-line cells and somatic cells of many organisms. Clearly these changes in the DNA content of the nucleus of cells are carefully regulated and lead to the appearance of the extra DNA at a predetermined time and, in many cases, the dissipation of this DNA on cue.

Another type of DNA amplification (sporadic) is detected when cells are overcoming adverse environmental conditions. Amplification of this type, usually observed by selection for a desired phenotype, has been found in bacteria, yeast, insects, and vertebrates. Salient features of this amplification process are being studied in several laboratories, and several general characteristics have emerged. These include (1) multistep process leading to the generation of highly amplified sequences; (2) a [karyotype](#) characterized by chromosomal abnormalities; (3) a genetic instability of the resistant (amplified) phenotype which may extend into a marked clonal variation among cells; (4) a spontaneous rate of detection which varies but (5) may be increased by manipulations of the cellular growth conditions (among these, treatment with [carcinogens](#)).

Carcinogenesis has been studied since Potts' initial observation that chimney sweeps are prone to

scrotal cancer. Several general characteristics of the carcinogenic process have emerged and must be accounted for in any hypothesis addressing the molecular events involved. These include (1) the seemingly multistep nature of the process; (2) the chromosomal abnormalities (in many cases non-random) that accompany most oncogenic states; (3) the clonal variation encountered in tumor cell populations; and (4) the ability of the rare (in some cases) process to be markedly increased by physical and chemical agents (known as carcinogens).

These two processes, DNA amplification and carcinogenesis, share several characteristics. Understanding of both are fundamental to our ability to understand how living organisms, single cells or complex organisms, respond to adverse environmental conditions. Several excellent reviews of both processes have appeared in the literature (1-8). This article begins by briefly reviewing programmed and sporadic amplification. The remainder addresses the measurement of sporadic DNA amplification in cells that have begun the carcinogenic process and demonstrate differing tumorigenic potential.

## 1. Gene Amplification and Development

### 1.1. In Germ-Line Cells

The earliest suggestion of developmental DNA amplification was made by King in 1908 when she described the extra "chromatin" (associated with forming nucleoli) that arose in the oocytes of the toad *Bufo* during pachytene (9). Later studies showed that these "masses" contain DNA, and **hybridization** studies demonstrated a great excess of sequences that code for ribosomal RNA (10). This differential synthesis of rRNA genes is correlated with the appearance of hundreds of **nucleolar organizers** during the pachytene stage of meiosis. Similar examples of this phenomenon are seen in other amphibians, *Xenopus*, *Rana*, *Eleutherodactylus*, *Triturus* (10), and in an echinoid worm, the surf clam (11) and many insects (12). Documentation of rDNA amplification is particularly amenable because the rDNA forms **nucleoli**, which are distinctive in appearance, and because most rDNA sequences have a relatively high guanine-cytosine content, which allows separating them from bulk DNA by CsCl **density gradient centrifugation**. These two aspects of rDNA have been important in recognizing this and other phenomena (magnification and compensation) in which the rDNA copy number changes under various genetic conditions.

The most complete study of oogenic rRNA gene amplification has been made in *Xenopus laevis*. In this animal, the early primordial **germ cells** do not contain amplified rDNA. Amplification is initiated in both oogonia and spermatogonia of the tadpole during sexual differentiation and germ cell mitosis and results in a 10 to 40-fold increase in rDNA genes. The premeiotic amplification is lost at the onset of the meiotic prophase. This loss is permanent in male germ cells but is temporary in female germ cells. Early during the meiotic prophase, the oocyte nucleus undergoes a second burst of rDNA amplification that results in a 1000-fold increase in ribosomal genes (13). The mechanism by which the first extrachromosomal rDNA copies are produced in the premeiotic stage is unknown, although circular structures have been observed (14). Because the number and placement of ribosomal **cistrons** do not change, a mechanism involving disproportionate replication is favored. The second burst of rDNA amplification probably involves a **rolling circle** intermediate. Such structures have been observed by Hourcade et al. (15) and could account for the increase of rDNA in the oocyte during the given period (16). When the oocyte nucleoli are hypotonically disrupted, DNA rings in the shape of beaded circles are visible. Although the rDNA content is constant, the sizes and number of the rings are variable, suggesting fission and fusion of nucleoli during oogenesis (17). The sizes of the circular molecules are integral multiples of a basic unit. Using molecular techniques, it was found that this basic unit is the DNA segment that codes for one precursor rRNA molecule, plus the accompanying nontranscribed spacer region. It is interesting to note that extrachromosomal, circular molecules that contain rDNA are also found at a low frequency (0.05 to 0.15% of the total number of molecules) in *Xenopus* tissue culture cells and in *Xenopus* blood cells (18). The circles are oligomeric lengths corresponding to between one and four rDNA repeat units. The presence of interlocking circles of equal size indicates that the circles are not *in vitro* cyclization artifacts but probably represent free rRNA genes. One question this raises is whether or not free DNA circles

carry other genomic sequences or if they are limited to rDNA sequences. It has been suggested that the difference in the state of rDNA between a somatic cell and an amplified germ cell may only be one of degree (14).

The cytological literature contains many other references to extrachromosomal DNA in oocytes. One of the most striking examples is found in dytiscid water beetles, such as *Rhynchosciara*. The nucleus of each oocyte contains a large [chromatin](#) mass, termed Giardina's body, in addition to the chromosomal complement. In older oocytes, this DNA is associated with multitudes of nucleoli and contains an increase in rDNA sequences (12). Hence, amplification of rDNA occurs during maturation of the oocyte. However, hybridization studies indicate that only a fraction of the extrachromosomal DNA is rDNA sequences, indicating that other DNA sequences (of unknown function) are also amplified (12).

Although the phenomenon of gene amplification in *Tetrahymena*, also results in a considerable increase in rDNA sequences, it is somewhat different from that in oocytes described previously. *Tetrahymena* contain two types of nuclei in each cell, a transcriptionally quiescent **micronucleus** that is responsible for genetic continuity and a macronucleus that is derived from the micronucleus in a developmentally regulated process whereby micronuclear sequences are eliminated, rearranged, and amplified (19-21). Considerable amplification of many micronuclear sequences occurs that results in the macronucleus containing 45 times the [haploid](#) amount of DNA. Some of this increase is accounted for by amplification of rDNA sequences, which are present as a single integrated copy in the micronucleus but in 200 copies in the macronucleus (20-22). These amplified rDNA sequences are present as linear, extrachromosomal, **palindromic** molecules (21, 23). The present data favor a model whereby excision of the single rDNA copy is followed by amplification (21).

The amplification, rearrangement, and elimination of sequences is also developmentally regulated to varying extents in the [slime molds](#), *Physarum* and *Dictostelium* (Zellweger et al., 1972) and in another ciliated protozoan, *Stylonychia* (Prescott and Murti, 1973). Studies of the molecular structures of these amplified DNA sequences have been undertaken and have revealed interesting structures at the termini.

## 1.2. In Somatic Cells

Gene amplification during differentiation of somatic cells also occurs and again was first detected by morphological criteria. Several regions of the [polytene chromosomes](#) in the larval salivary glands of the fly *Rhynchosciara americana* showed a **puffing** response to hormone treatment that was carefully defined both temporally and spatially. These puffs contain greater amounts of DNA than surrounding regions (24). Two of these puffs, characterized at the molecular level, contain a 16-fold increase of DNA sequences that code for peptides used in the synthesis of the cocoon. Amplification and puff formation depend on the developmental stage of the cell and on the cell's position within the gland (25). Puff formation is not peculiar to the salivary gland chromosomes of *Rhynchosciara*, because it has also been observed in the cells of Malpighian tubules and intestinal cells of the same insect. In addition, puff formation has also been observed in different tissues of *Chironomidae*, *Drosophila*, *Sciara*, *Hybosciara* and other diptera (24). In some cases it is thought that it is under hormonal control (26, 27). In *Rhynchosciara* and *Sciara*, the extra DNA remains associated with the original chromosomal site until lysis at pupation (27, 28). In *Hybosciara* and *Sarcophaga*, in contrast, the extra DNA is released as micronuclei from several sites, similar to that released in the oocytes of amphibians (29). It is interesting that several of the bands that puff during the larval stage in *Rhynchosciara* are identical to those bands that are overrepresented in the differentiated chromosomes of *Sciara ocellaris*. This suggests that certain germ-line events observed in *Rhynchosciara* occur in somatic cells of other species (27).

In many cases of differentiation, somatic cells acquire adequate amounts of [messenger RNA](#) to produce abundant proteins by accumulating stable mRNA molecules over a period of days. Examples of such are the silk fibroin genes, the [ovalbumin](#) genes, and the **beta-globin** genes. It is thought that this is controlled at the transcriptional or posttranscriptional level. In other cases,

however, such as during the synthesis of the insect eggshell by the ovarian follicle cells of *Drosophila*, little time is allotted to produce the specific mRNAs that are needed in large quantities. In these latter cases the rates of transcription and [translation](#) are not high enough to produce adequate amounts of protein. Spradling and Mahowald (30) found that this need is met by differential amplification of the chorion gene sequences in the ovarian follicle cells. To study the regulation of expression of the chorion gene proteins, they produced [complementary DNA](#) clones complementary to egg chamber RNA. In this way they isolated sequences that code for the major endochorion and exochorion proteins, s15, s18, s36, and s38. Using the s36 clone as a probe, they found that these sequences are more abundant in follicle cells during the terminal stages of egg chamber development than in nonovarian tissue. The sequences that code for s15, s18, and s38 were also being specifically amplified during oogenesis. Spradling (31) found that the chorion genes are located on the X-chromosome are in two clusters, s36 and s38, and are amplified 15-fold. The s15 and s18 loci are on the third chromosome and are amplified 60-fold. The genes in both clusters are amplified at the same time and both homologues are amplified equally. Sequences that flank these genes are also disproportionately replicated, but not to as great an extent. This results in a gradient of amplification that spans 90 kbp of DNA, is maximal in the center, and does not display any discrete termination sites. Amplification of the two linked chorion genes s36 and s38 (but not s15 and s18) is reduced in ocellar mutants. The ocelliless phenotype is a pleiotropic mutation characterized by homozygous females and males that totally lack ocelli (simple eyes or photosensitive organs). They also exhibit an abnormal pattern of bristles in the ocellar region of the head. Ocelliless females lack parovaria, an accessory gland of the reproducing tract, and produce abnormal chorions (which surround the egg), which leads to sterility. The ocelliless mutation has been identified as a small chromosomal inversion and has one of its break points lying near the s36 gene (32). In this inversion, the site for initiating [DNA replication](#) responsible for amplifying the chorion genes was moved near the other break point. Using cloned probes for the nonchorion genes now present near the end of the inversion, it was found that the gene copy numbers of these sequences were increased. This result suggests that clustering of the chorion genes has biological significance. Developmentally controlled activation of a single **origin of replication** can result in the amplification of several genes.

Changes in the DNA content of the nucleus of germ-line cells and somatic cells during development is common. Amplification of rRNA sequences and others that code for proteins that needed in large amounts during that particular phase of development have been extensively documented. In still other cases, the nature of the amplified DNA cannot be totally accounted for by known sequences, and it probably contains amplified sequences of unknown function. Such sequences have been implicated in the differentiation of the orchid *Cynidium* (33), the differentiation of peas (34), and in the flowering of the tobacco plant, *Nicotiana* (35). Further studies are needed to relate these instances of DNA amplification to the overall development of higher organisms.

## 2. Gene Amplification and Acquisition of a Selected Phenotype

Duplication and amplification of genetic material in cells has long been documented as a means for overcoming deleterious growth conditions. Unlike amplification that is specifically regulated during development, amplification as a means of survival is a more sporadic event whose frequency is often detected at levels much lower than those in developmental amplification.

### 2.1. Bacteria

The first well-documented example of gene duplications in bacteria that provides an adaptive advantage was seen by Novick and co-workers (36-38). Bacterial strains were grown for long periods of time in limiting concentrations of lactose in the chemostat. The bacterial strains that emerged synthesized four times the maximal normal amount of **beta-galactosidase**. The ability to produce large amounts of enzyme was unstable and could be transferred by conjugation at a time when the [lac operon](#) genes were expected to be transferred. It was concluded that the ability to overproduce beta-galactosidase was due to extra copies of the *lac* genes in these strains. The spontaneous rate of duplication was estimated at  $10^{-3}$  and  $10^{-4}$  in a similar system (39). Overproducers that have similar characteristics were subsequently reported for other enzymes, such

as ribitol dehydrogenase, [b-lactamase \(40\)](#), and several others (for an excellent review, see Ref. [41](#)). Roth and co-workers demonstrated that duplications up to a quarter of the Salmonella chromosome occur at large homologous segments, such as the rRNA genes ([41](#), [42](#)). Such large duplications depend on the *recA* system. On the other hand, duplications in the range of 10 to 30 kbp are independent of *recA* and do not involve very large homologous segments of DNA ([42-44](#)). Edlund and Normark ([45](#)) detected tandem duplications of 10 to 20 kbp at the *E. coli* chromosomal *ampC* locus that codes for b-lactamase and confers resistance to [ampicillin](#). After a stepwise selection in increasing concentrations of ampicillin, 30 to 50 copies of the duplications analyzed at the [restriction fragment](#) level had different end points. The junction points of the amplified unit in one case was sequenced, and the original duplication occurred at a sequence of 12 bp that was repeated on each side of the *ampC* locus, 10 kbp apart.

Ilstyt et al. ([46](#)) detected amplification as revertants of certain leaky Lac- mutants. The DNA sequences, which were amplified, contained the *lac* operon and anywhere from 7 to 32 kbp of flanking sequences. These regions were amplified 100-fold. More than 160 independent events were analyzed at the restriction fragment level, and the majority had different end points. These amplifications and as those studied by Normark are not mediated by large homologous sequences ([47](#)). Virtually all of the duplications detected and subsequently studied in bacteria have been tandem. Genetic techniques have permitted localizing the end points of these tandem duplications to two different classes of sequences, those composed of large homologous sequences of DNA (such as ribosomal genes or [transfer RNA](#) genes in *E. coli*) and those not. The former type may be visualized as occurring by a recombination system for which we have already identified the enzymes. The latter class of duplications raises the question of the lower limits of homology that may be used by such a system, or it raises the possibility that it is a different system entirely that mediates duplications whose end points lack long homologous stretches.

Several bacteriophages also amplify genetic markers. [Lambda phage \(45\)](#), [T4 phage \(48\)](#), and [P1 phage \(49\)](#) may use mechanisms that parallel those used by viral sequences in mammalian cells. The mechanism is yet unknown.

## 2.2. Yeast

Resistance to the toxic effects of copper in *Saccharomyces cerevisiae* is mediated by tandem gene amplification of the CUP1 locus ([50](#)) but differs from sporadic amplification events characterized in other organisms. The CUP1 locus of yeast codes for a small copper-binding protein. Copper-sensitive strains contain one copy of this locus and, when grown in elevated concentrations of copper, fail to produce resistant derivatives that have a higher gene copy number of CUP1 genes on one chromosome. Infrequently, copper-resistant strains isolated in the laboratory and to carry up to 10 tandem duplications of the CUP1 region, which has 2kbp. Further amplification could be achieved by growing the copper-resistant strains in elevated copper concentrations. The authors postulate that the mechanism of amplification in this instance proceeds by forming a disome for chromosome VIII (which carries the CUP1 locus). Copper-resistant mutants of the sensitive strain are disomic for chromosome VIII. Then, the amplification or tandem iteration could result primarily from subsequent unequal chromosome or sister chromatid exchanges. The formation of a disome may constitute an initial event in the process. Now the system is being used to estimate the frequency of sister chromatid exchanges, unequal crossovers, and gene conversions in this region. The break points occur in the nontranscribed spacer region ([51](#)).

A second example of gene amplification in yeast has molecular structures more similar to those described in other organisms. A yeast strain resistant to antimycin A, an [alcohol dehydrogenase](#) inhibitor, contains multiple copies of a nuclear gene, ADH4, an isoenzyme of alcohol dehydrogenase. The amplified copies are 42 kbp long, display a linear, extrachromosomal, palindromic structure, and contain **telomeric** sequences. Their structure resembles that of the amplified rDNA genes in the macronucleus of Tetrahymena and related ciliated protozoa, except that the nuclear copy remains within the chromosome. In contrast to that often observed in mammalian amplification (see later), the extrachromosomal copies of this gene are stable during

mitotic growth. Amplification of the ADH4 gene is relatively rare ( $\sim 10^{-10}$  mutations/cell/generation), and alternative mutations of antimycin A-resistance are in the majority.

### 2.3. Protozoans

Drug resistance in protozoan parasites is common and presents a serious problem for the chemotherapy of diseases caused by such pathogens as Trypanosoma, Leishmania, and Plasmodium (52-54). Recently, Leishmania strains resistant to the well-known chemotherapeutic agent methotrexate have been isolated and analyzed for their mechanism of resistance. Organisms that are resistant to high concentrations of methotrexate (1mM) have a 40-fold increase in [dihydrofolate reductase](#) (DHFR) which in this organism is associated with **thymidylate synthetase (TS)**. The resistant organisms contain amplified regions of DNA that can be observed directly on **ethidium bromide**-stained gels of [restriction enzyme](#) digests. The amplified region is 56 kbp long, has a copy number about 80-fold higher than that of wild-type organisms, and codes for the DHFR-TS genes and other sequences of unknown function (56). Direct visualization of the amplified DNA fragments in this system and the amplified sequences described in the *lac* region of *E. coli* (see previous), provide a valuable advantage in identifying sequences involved in generating a desired phenotype. These sequences may be observed because of the low complexity of the genomic DNA in these organisms. The targets of most drugs used for treating pathogenic protozoans and the targets of many selection protocols in bacteria are unknown. The direct visualization of amplified DNA in these organisms provides a means for obtaining these target sequences and gaining insight into their mechanisms of action.

Recent studies have shown that *Plasmodium falciparum* contains genes analogous to the **multidrug-resistance** (*mdr*) genes in mammalian cells. Parasites that resist chloroquine are also resistant to other antimalarial drugs, similar to the phenomenon of multidrug resistance seen in human tumors. Wilson et al. engineered a probe using sequences that are conserved in the mammalian [P-glycoprotein](#) to identify *mdr*-like sequences in *P. falciparum*. Their studies showed that drug-resistant parasites contain amplified copies of these specific DNA sequences when compared to their drug-sensitive siblings.

### 2.4. Invertebrates

Amplification of rRNA genes in *Drosophila* during oogenesis does not occur as has been described in *Xenopus* (see previous section). However, sporadic amplification of the rRNA genes during one generation has been observed under specific genetic conditions. Amplification of the rDNA genes in this situation results in reversion to a wild-type phenotype. Each sex chromosome carries approximately 130 to 150 rRNA genes (56). The phenotype is wild type if the diploid cell carries at least one normal locus ( $\sim 130$  genes), whereas the phenotype is altered (bobbed) if the genome carries less than 130 genes (56, 57). The intensity of the bobbed phenotype (slow development, thin chitinous cuticle, reduced body traits, and short bristles) is inversely proportional to the number of genes for rRNA. An increase in rDNA copy number is observed in the progeny of phenotypically bobbed males. It involved rapid accumulation of rDNA at either nucleolus organizer by unknown mechanisms. The rDNA that accumulates during the first generation does not have a noticeable effect on the phenotype of the fly. The phenotypic effects of the magnified rDNA become evident in the F2 progeny if the rDNA is transmitted by males and if the genotype of the F2 generation is characterized again by rDNA deficiencies. In other words, the extra copies of rDNA are eliminated if the amplified fly is crossed with a normal (bb+) female. The phenotypic inheritance of amplified rDNA requires integration and inheritance through the male germ line. Two hypotheses have been proposed to explain rDNA magnification, disproportionate replication of rDNA (58, 59) or unequal sister chromatid exchange (60). Recent experiments demonstrating the decrease of magnification frequency in organisms that carry the rDNA on a ring chromosome strongly suggest unequal crossing over as the mechanism.

Another type of amplification in *Drosophila melanogaster* differs from rDNA magnification in several characteristics. This amplification, called compensation, occurs when one nucleolus organizer of the two homologues is completely deleted (X/O or X/X-no females). In such mutants,



the remaining organizer “compensates” for the deletion of the rDNA sequences by a disproportionate replication of the remaining sequences on the intact homologue. Compensation may occur only on the X chromosomal nucleolus organizer, and the extra r-DNA is not inherited in subsequent generations. DNA magnification and compensation differ in several aspects (61): (1) In contrast to rDNA magnification, which restores only wild-type levels of rDNA sequences, compensation reaches values that exceed the rDNA copy number of the nucleolus organizer in the wild-type fly. (2) Gene compensation occurs only on the X chromosomal nucleolus organizer. The Y nucleolus organizer is refractory to compensation. As discussed earlier, amplification may occur on either. (3) rDNA accumulated during compensation is not inherited in subsequent generations as is the rDNA resulting from amplification. Compensation occurs in somatic cells and is not passed along. (4) When a genotype carries only one nucleolus organizer in the diploid chromosome set and the same genotype is partially deleted for rDNA in the remaining nucleolus organizer, the conditions for both compensation and amplification are in place. Genetic studies that compare  $Xbb^+/O$  flies with  $Xbb/O$  flies show a much greater increase in rDNA in the  $Xbb/O$  flies, which suggests that magnification and compensation can occur simultaneously. An interaction between homologous chromosomes is a prerequisite for initiating compensation. If two complete nucleolus organizers are on one X chromosome, compensation still occurs.

Using various deficiencies for the X chromosomal [heterochromatin](#), evidence has been presented for the existence of a genetic locus that regulates rDNA compensation (62). This locus, called the compensatory response (cr), is located outside the ribosomal cluster and in the X chromosomal heterochromatic region. The locus acts in *trans* to sense the presence or absence of its partner locus on the opposite homologue. If only one cr locus is present, it acts in *cis* by driving compensation (disproportionate replication) of adjacent rRNA genes. Not all embryos that have the proper genotype undergo compensatory amplification to emerge with an increased number of rRNA genes. Only a small fraction undergoes the putative compensatory amplification. In this respect, the amplification behaves like a mutagenic reversion to restore the functional phenotype.

Resistance to environmental agents, such as pesticides and toxic chemical waste, has been documented in laboratory stocks and natural populations of invertebrates. Selection of *Drosophila* larvae in increasing concentrations of cadmium yields strains that contain duplications of the [metallothionein](#) gene (63). The duplication is stably inherited in the absence of selection pressure and produces a corresponding increase in metallothionein messenger RNA. A survey of natural populations found that this is common (64) and may signal the early stages of the evolution of a gene family. The mosquito, *Culex quinquefasciatus*, develops resistance to various organophosphorus insecticides by overproducing the enzyme esterase B1. Molecular studies have demonstrated that the overproduction of the enzyme is the result of amplification of the esterase B1 gene some 250-fold (65). The resistant mosquito was normally developed and could reproduce. This observation raises questions of evolutionary significance for the duplication and amplification at least in invertebrates.

## 2.5. Plants

DNA changes in plants during response to environmental stress have been reported for flax (see Refs. 66-68). The suggestion that these changes in DNA content are induced by the environment awaits verification by further studies.

## 2.6. Vertebrates

DNA amplification in mammalian cells was first detected when murine tumor cell populations became resistant to chemotherapeutic drugs. Methotrexate, an often-used chemotherapeutic drug, inhibits the action of dihydrofolate reductase (DHFR), which is required for biosynthesizing thymidylate, glycine, and purines. Stepwise selection of cells in increasing concentrations of methotrexate generates highly resistant cells. Studies based on acquisition of resistance to methotrexate showed several alternative mechanisms of resistance. Cells resistant to the action of the drug had (1) altered the structure of the DHFR enzyme so that the folate analog was not efficiently bound (69-71); (2) altered their transport process so that the drug could no longer enter the cell (72); or (3) overproduced the DHFR enzyme (73, 74). Beidler and Spengler (75) detected chromosomal

abnormalities in the cells that overproduce DHFR and suggested that they reflect an increase in gene dosage. Schimke and co-workers obtained a cDNA for the DHFR sequence and showed that the overproduction of DHFR results from amplification of the DHFR DNA sequence (76). Now, it is known that amplification of the DNA sequence coding for the target enzyme of a metabolic inhibitor is a common mechanism for overcoming the growth restriction (2, 77-80).

Among other examples of this phenomenon subsequently found, the best studied was amplification of the CAD gene. The CAD gene codes for a multifunctional protein that catalyzes the first three steps in the synthesis of pyrimidines. The [aspartate transcarbamoylase](#) activity is inhibited by the **bisubstrate analog**, *N*-phosphoacetyl-L-aspartate (PALA). PALA-resistant cells overproduce the aspartate transcarbamoylase and also the other two enzymes, carbamoyl synthetase and dihydroorotase (81). Wahl et al. (82) showed that overproduction of these enzymes results directly from amplification of DNA coding for these proteins.

Two other enzyme activities, orotate phosphoribosyl transferase and orotidine-5'-phosphate decarboxylase, are responsible for the final steps of pyrimidine biosynthesis and are encoded in a single polypeptide chain called UMP synthetase. Cells resistant to 6-azauridine or pyrazofurin (toxic pyrimidine analogues that inhibit the decarboxylase) demonstrate a coordinate increase in both of these enzyme activities (83). Overproduction of the two activities is due to amplification of the DNA sequences that code for the UMP synthetase (84, 85). Numerous other instances of DNA amplification have been described when the growth of cells is inhibited by metabolic inhibitors, toxic agents, or altered enzymes that have reduced efficiency. Thus a leaky mutant of hypoxanthine phosphoribosyl transferase amplifies its gene, and growth in cadmium results in amplifying metallothionein (86). Mutant mouse lymphoma cells that overproduce ornithine decarboxylase (up to 15% of total cellular protein) have been generated by selection for resistance to difluoromethylornithine (DFO), an inhibitor of the enzyme. The pseudodiploid wild-type cells become **tetraploid** when resistant to high levels of DFO (87) and amplify the ornithine decarboxylase genes. Compactin-resistant Chinese hamster ovary cells contain a 100 to 1000-fold increase in the amount of 3-hydroxy-3-methylglutaryl coenzyme A reductase, a pivotal enzyme in synthesizing cholesterol (88). Molecular analysis demonstrated that the sequences that code for this reductase are amplified 15-fold in these cells. Similarly, methionine sulfoxine-resistant cells that overproduce glutamine synthetase amplify these sequences. The discovery that multi-drug resistance in cancer chemotherapy is, in some cases, mediated by amplifying the *mdr* locus has been clinically important (89, 90).

As seen previously, the overproduction of an enzyme activity often indicates amplification of specific sequences. Such overproduction has been described in the resistance of cells to such diverse agents as hydroxyurea, aphidicolin, beta aspartyl hydroxymate, mycophenolic acid, tunicamycin, 5-fluorodeoxyuridine (93), and several other drugs (for a comprehensive list and references, see Ref. 2). The demonstration that amplification is the basis of resistance in these cells waits cloning of the target genes in question. Although overproduction of specific enzyme activities often indicates gene amplification, exceptions exist, and each instance of overproduction must be properly analyzed before gene amplification is concluded. An example of this is seen in the studies of arginosuccinate synthetase (AS) overproduction. Canavanine is a toxic analog of arginine but is not an inhibitor of arginosuccinate synthetase. Canavanine-resistant cells overproduce arginosuccinate synthetase enzyme that is no longer subject to metabolic control by exogenous arginine (92, 93). Molecular studies show that this instance of overproduction is not due to amplification of the AS gene (92, 93) but is a regulatory mutation. A compactin-resistant regulatory mutant that overproduces HMG Co A reductase has also been found.

Selection pressure also leads to the amplifying sequences that have initially unknown functions. Baskin and co-workers (94) isolated cell mutants that are reciprocally cross-resistant to four diverse drugs; maytansine, adriamycin, Baker's antifol, and vincristine. The resistance is genetically unstable and is accompanied by the presence of [double minute chromosomes](#) (DMs) in the karyotype. These are small chromosomes that lack a centromeric sequence and therefore segregate randomly during

mitosis. Genetic instability and the presence of DMs are two manifestations of gene amplification (94). Other cells that are cross-resistant to other structurally diverse and structurally unrelated drugs amplify the cell surface glycoprotein called P-glycoprotein. Resistance to agents, such as [colchicine](#), [actinomycin D](#), [puromycin](#), and in some cases [vinblastine](#), generate overproduction of a 170 kDa glycoprotein that reduces membrane permeability (89, 90).

The acquisition of a specific phenotype does not necessarily have to occur by the classic selection protocols that involve growth in cytotoxic substances. The fluorescence-activated cell sorter (FACS) allows detecting single cells that overproduce products, providing one can monitor the product (see [Flow Cytometry](#)). This instrument monitors the fluorescence of individual cells as they pass through a fluorescence detector and sorts out the rare deviant with ease. Johnston et al. (95) used this technique to detect amplification of the DHFR gene without growth in methotrexate inhibitor. Chinese hamster ovary cells were stained under nonselective conditions with fluoresceinated methotrexate, which binds quantitatively to DHFR. The cells exhibit heterogeneity in DHFR content. The cells that contain the highest DHFR content were sorted, grown in nonselective conditions, and reanalyzed. After ten successive rounds of growth and sorting, the population showed a 50-fold increase in fluorescent intensity and amplified the DHFR gene 40-fold. The population was highly resistant to methotrexate.

Using the same concept, Kavathus and Herzenberg (96) isolated mouse L-cells that amplified the human T-cell differentiation antigen, Leu-2. Mouse cells were **transfected** with the herpes simplex thymidine kinase (TK) gene and total human DNA. Transfectants that expressed Leu-2 on their surface (as detected with a fluorescent [antibody](#)) were sorted on the FACS. Although most Leu-2 transfectants had a narrow range of antigen density per cell, one showed considerable variability. The brightest cells were sorted again. Six rounds of FACS “selection” for the brightest cells produced cells with a mean fluorescence 40-fold greater than the original transfectant. The karyotype of these cells contained a multitude of DM's. A [cDNA library](#) was constructed using mRNA from the amplified Leu-2 transfectant, and the Leu-2 cDNA was identified by differential screening. The mouse line that overproduced Leu-2 antigen was confirmed as having amplified the human Leu-2 DNA sequence (96).

Thus, the FACS instrument dramatically facilitates detecting cells that contain amplified sequences, which might otherwise be overlooked in classical cytotoxic selection experiments. These cells may be detected soon after they are generated, without requiring growth into a visible colony. The FACS detects as little as a twofold increase in fluorescence. This allows identifying and isolating cells that have low or increased transient or unstable levels which be. Any fluorescent probe that quantitatively correlates with a specific molecule or metabolic process allows detecting and sorting cell populations that have successive increases (or decreases) of that tag. Using this concept, one can obtain cell variants that amplify (or delete) genes for which there is no available selective agent. The FACS provides a powerful tool for studying heterogeneity in cellular populations and specific amplification of DNA sequences.

The acquisition of a selected phenotype may often result from selection pressures that are unknown at that time. In these cases a certain phenotype may be accompanied by the manifestations of gene amplification for unknown sequences. Such an instance has been described in studying the sequences that are carried on the DMs and found in the homogeneously staining regions (HSRs) of neuroblastoma cells where these structures were first described. The sequences that are amplified in these lines are cellular “onc” genes, the N-[myc oncogene](#) (97). Now, amplification of **oncogenes** has been found in several tumor types containing DMs and HSRs (97) and is discussed later.

As a class of rearrangements, tandem duplications or amplifications cause very little loss of genetic information and are therefore probably not significantly deleterious to the cell. On the other hand, the extra copies are easily dispensed with when they are no longer necessary. This would circumvent any growth disadvantages that may be encountered under normal nutritional or environmental conditions. The only limiting factor in obtaining an amplified DNA sequence is the design of a

selection system that allows detecting the sequence in question. This principle can be manipulated as a tool, that is, amplification of unknown gene sequences can be obtained to overcome defined selection conditions. For example, selection of cells that may grow in the continued presence of **DNA damaging** agents may result in some cases in amplification of DNA coding for [DNA repair](#) enzymes.

### 3. General Properties of Mammalian Gene Amplification

There are several characteristics of amplification that many of the systems discussed share in common. To illustrate these characteristics, I describe the general properties of methotrexate-resistant cells that result from amplifying the DHFR gene.

#### 3.1. Multistep Selection to Obtain Highly Amplified Sequences

Classically, mammalian cells that contain amplified DHFR genes were obtained by a stepwise selection for cells that are highly resistant to methotrexate (MTX). Before the advent of cloning, the amplified phenotype was distinguished by overproduction of DHFR enzyme and the karyotypic abnormalities that accompany the overproduction. High methotrexate resistance (by virtue of amplification of the DHFR gene) cannot be obtained by a single-step selection protocol. It is a multistep process. The initial step is rate-limiting because cells that have an increased, but low copy number, are rapidly stepped up to a high level of resistance and a high copy number. When the initial increase in gene copy was examined more closely, Brown et al. ([98](#)) and Tlsty et al. ([99](#)) found that the stringency of selection is critical in obtaining cells that have amplified DHFR. Cells that amplify DHFR genes as a mechanism of resistance to methotrexate are most readily detected within a defined window of methotrexate concentration for the first step of selection. Growth of the cells above or below this range of methotrexate concentration reduces the frequency with which amplified clones were detected. Not only do incremental increases in drug concentration promote the rapid emergence of resistance, but they also specifically promote the rapid amplification of the DHFR gene ([100](#)). A shallow increment of stepwise selection allows for emergence of methotrexate-resistant amplified cells more rapidly than in a relatively large single-step selection protocol.

#### 3.2. Karyotype Abnormalities

The second property of methotrexate resistant cells that have amplified their DHFR gene is the frequent karyotypic abnormalities in the cells. As indicated previously, abnormal chromosomal structures were associated with overproduction of the DHFR in the early studies of Biedler and Spengler ([75](#)). They described a marker chromosome in overproducing cells that have an elongated chromosomal arm. The term homogeneously staining region (HSR) was coined to describe a region of this chromosome that bands abnormally when stained with giemsa and, as was subsequently shown, is the site of the amplified DHFR sequences ([76](#)). This structure was associated with stable resistance to methotrexate. The resistant phenotype is retained even after subsequent growth in the absence of selection pressure. This is in contrast to the karyotype of cells that are unstably resistant to methotrexate. After extended growth in a nonselective medium, the resistant phenotype (amplification) diminishes rapidly and disappears. HSR structures were not found in unstably resistant cells. Close examination of the karyotype of unstably resistant cells, however, brought to light the presence of double minute chromosomes (DMs). These structures (and HSRs) were described by Balaban-Malenbaum and Gilbert ([101](#)) in cell lines obtained from human neuroblastoma. Subsequent work demonstrated that unstably resistant cells contain the amplified copies of DHFR on the double-minute chromosomes (DMs). The lack of centromeric structure in these fragments leads to their random (unequal) segregation at mitosis and a diminished number if selection pressure is no longer exerted on the cells ([102](#)).

HSRs are often found at the site corresponding to that of the amplified gene but may also be found at many other chromosomal site or sites. In several instances, multiple HSRs, each of which contain some copies of the amplified DNA sequences, are in one cell. The amplified copies of the target gene may be at the site of the resident gene or in another place or various places in the karyotype. Translocation of the amplified sequences has been observed often ([103](#)), and may provide clues

about the mechanism of amplification. These translocations may be nonrandom in (see Refs. [104](#) and [105](#) for extensive review).

The mechanism leading to the generation of multiple HSRs is not known now. Typically a cell must contain a reasonable increase in gene copy number before an HSR is observed. The degree of amplification and the length of the amplified unit both play an important part in the ability to visualize an HSR. This principle is nicely illustrated by a study where only one of two cell lines, both of which have approximately 50 copies of an amplified sequence, contained a visible HSR. It is postulated that the amplified unit in the other cell line is so short that even a 50-fold reiteration does not make it visible. Similar consideration may be applied to the visualization of DMs. The chromosomal fragments are heterogeneous in size but routinely small and difficult to see. If a cell population has been selected for resistance to a high concentration of an inhibitor, amplification is extensive and DMs are more readily detected. During the initial steps of amplification, when resistance to relatively low concentrations of inhibitors is obtained, only a few DMs may be in any one cell, and visualizing them may be difficult.

The terms stable and unstable resistance are relative designations. Although cells that contain HSR's can lose the amplified gene copies when removed from selective pressure, this occurs at a much slower rate (months) than when amplified sequences reside on DMs (weeks).

What are HSRs and DMs? What is their molecular structure? Are the amplified sequences arranged in tandem configuration, in direct or inverted orientation? The first obstacle in characterizing the amplified unit derives from its large size. The DHFR gene, which is amplified to confer methotrexate resistance, is large: 31 kbp, including introns. The size of the amplified region is greater still. Gross estimates vary from 120 to 1000 kbp as the unit of DNA that is amplified. If one could analyze the end point or end points of the amplified units, information on the structure would emerge. Although the sequence of DNA that needs to be characterized is long, cloning of neighboring fragments (chromosomal walking) has been accomplished by several laboratories in both mouse and hamster models ([106-108](#)). The information derived from these endeavors has not provided the desired portrait of the amplified unit because of another obstacle. The amplified structure is continually changing at the molecular level in the chromosomal walking studies described previously. On each amplified cell studied, the amplified sequences correlate with the cloned map only up to a certain point and then diverge. Rearrangements of DNA accompany the amplification of genes, and the basic molecular DNA amplification is obscured by the dynamic aspect of the process. Karyotypic abnormalities are detected only in cells that are highly resistant to a given metabolic inhibitor (i.e., cells that have already progressed through much of the multistep process), although Hamlin and Montoya-Zavala found that the amplified DHFR gene in Chinese hamster ovary cells is uniform in size and exists in head-to-head and head-to-tail tandem repeats.

Recent studies in several laboratories have begun addressing the karyotypic abnormalities that are evident during the initial periods of selection for resistance. In contrast to classical studies on methotrexate drug resistance, where cells were grown for long periods of time in increasing concentrations of drug (stepwise selections), the recent studies are aimed at elucidating the initial molecular events of DNA amplification and have use a single-step selection with the drug. Cells resistant to these relatively low amounts of drug amplify the target sequence to a low extent, 5- to 10-fold ([98, 99](#)). DMs and HSRs that have this low copy number are difficult or impossible to detect. Small pieces of extrachromosomal DNA that are not visible with the light microscope have been reported ([108-112](#)). Several agents increase the incidence of DNA amplification. As the cells undergo the initial steps of DNA amplification within this period, other manifestations of chromosomal instability are observed, such as increases in the frequency of sister chromatid exchange, endoreduplication, and polyploidization. The production of extrachromosomal DNA has also been detected, but it is transient if selection pressure is not exerted on the cells. Pretreatment with hydroxyurea, ultraviolet light, or methotrexate itself increases the incidence of the initial amplification of the DHFR sequences ([98-100](#)). Similar observations were made using an SV40-transformed cell system to detect amplification of SV40 sequences ([113](#)) (see next section).

### 3.3. Genetic Instability

A third characteristic of methotrexate-resistant cells that have amplified the DHFR gene is the initial genetic instability of the resistant (amplified) phenotype, which is accompanied by a marked heterogeneity in the population. For example, Chinese hamster ovary cells were grown in 0.02  $\mu\text{M}$  MTX and then analyzed for DHFR content by the FACS (114). The population of cells that emerged after the initial single-step selection showed a heterogeneous spectrum of DHFR content. Cloned cells from various portions of the spectrum were obtained and analyzed for homogeneity of DHFR content and stability of the elevated DHFR levels. CHO cells newly selected for methotrexate resistance are unstable with respect to DHFR levels, and the loss of the elevated DHFR levels was variable in the progeny of different cloned cells. This property of initial instability and heterogeneity is not to be confused with that described previously, in which stable and unstable resistance was associated with HSRs and DMs, respectively. The initial instability is observed in the newly selected population. When these cells are selected for a prolonged period in methotrexate, they have stable, elevated DHFR levels. During subsequent increases in the methotrexate concentration in the medium, the population becomes progressively more heterogeneous with respect to gene copy number. Once the initial increase in gene copy number has occurred, highly methotrexate-resistant populations are subsequently generated more rapidly. When the cells are grown for longer periods under fixed selection pressure (ie, at a constant methotrexate concentration) the population that emerges may contain either stable amplified genes or unstably amplified genes. The initial instability of the amplified DHFR genes in emerging, resistant CHO cells is consistent with the hypothesis that they are present as extrachromosomal pieces of DNA. Stabilization of the resistant phenotype could result from integration of these sequences into the chromosome at the site of amplification or elsewhere in the genome or could result from processes unknown now.

## 4. Frequency of Sporadic Amplification in Mammalian Cells

In the last few years it has become obvious that the frequency of gene amplification in different cells varies dramatically. Initially, gene amplification was measured in the models that were used to study the phenomenon, established rodent cell lines, such as S180, BHK cells, CHO cells, and 3T6 cells. The earliest type of measurement was colony formation in the presence of a given amount of cytotoxic drug, the clonogenic assay. Resistant colonies that emerged were analyzed for gene copy number. Although a critique of the various methods to measure gene amplification is beyond the scope of this review, it is important to note that the degree of stringency (concentration of the drug) is directly related to the extent of gene amplification. Therefore, incidence calculations that result from selections performed at different stringencies cannot be compared. Reported values for the rodent model systems were  $10^{-6}$  to  $10^{-4}$ . In 1983, Johnston et al. (95), used the flow cytometer and fluoresceinated methotrexate to analyze the emergence of cells that contain elevated DHFR enzyme content from a population that was not under selective pressure. They termed this the “spontaneous” rate of DHFR gene amplification. The emergence of individual cells that have elevated DHFR content was surprisingly high and varied considerably between clones (95). This frequency ( $10^{-3}$ ) is high, especially when compared with the incidence at which methotrexate resistant colonies emerge ( $10^{-5}$ ) when measured by the clonogenic assay. In a third method, (119) the entire population of cells was exposed to a stepwise selection and the populations were monitored for the appearance of subpopulations resistant to given concentrations of drug within a given time. All of these methods used rodent cell lines and concluded that amplification could occur at incidences of  $10^{-6}$  or  $10^{-4}$ , or rates that approached  $10^{-3}$  events/cell/generation.

In the 1980s, studies of the frequency of gene amplification took three major directions. One approach was to study the effect of genomic position on amplification, the second was to measure amplification in different cell types and physiological states, and the third sought to identify agents that could increase the frequency of amplification.

### 4.1. Positional Effects

In an elegant set of experiments, Wahl et al. (116) developed a system to measure the increase in

gene copy number of a specified marker placed in different locations in the genome. A CAD minigene was constructed and transfected into a Urd-hamster cell line. They found that the incidence of CAD minigene amplification varied from  $10^{-3}$  to  $10^{-3}$  and that the frequency was heritable. A parallel system (117) produced similar results by transfecting a DHFR minigene into a DHFR<sup>-</sup> cell line. Both studies concluded that the incidence of gene amplification is influenced by sequences or structures within the genome. Subsequent studies in the Wahl laboratory centered on identifying the *cis*-acting sequences that influence the frequency of gene amplification.

#### 4.2. Nontransformed Versus Transformed Cells

In a second approach, several laboratories have begun examining the incidence of gene amplification in different cell populations. Early results, using the stepwise selection of nontumorigenic and tumorigenic populations of CHEF cells, suggested that tumorigenic cells could amplify more rapidly than nontumorigenic cells (115). More recently, using the clonogenic assay at defined stringencies of selection, Otto et al. (118) found a difference in the incidence of amplification between nontumorigenic and tumorigenic cell lines. In both studies, sample size was small (three and six cell lines, respectively), but extending these studies, using the clonogenic assay, to greater than 60 cell lines confirmed that generally highly tumorigenic cells amplify at a greater frequency than nontumorigenic cells. The correlation between assayed tumorigenicity and gene amplification is not strict however, and identification of the variables that help define these frequencies awaits further experimentation.

To date, most studies of gene amplification have used immortalized cell lines and biopsied tumor samples. In two studies, however, the amplification potentials of primary diploid cells, both human and rodent, were examined and compared quantitatively to the amplification potentials of their transformed counterparts. Strikingly, the difference in amplification incidence between “normal” cells and their transformed counterparts (in some cases tumorigenic) is immense (119, 120). The amplification potential was measured at two loci, the CAD gene and the DHFR gene. Quantitative data for both normal ( $<2 \times 10^{-8}$ ) and transformed cell lines ( $10^{-4}$ ) indicated a difference in frequency greater than four orders of magnitude (120). In one study, the generation of ouabain resistance in both diploid and tumorigenic human cell lines was also measured to control for the ability to detect mutants. Although the rates of generating point mutations were identical, the generation of gene amplification at the CAD locus differed by at least four orders of magnitude. These studies suggest that there is some fundamental difference between normal cells and transformed cells that affects their ability to amplify.

In summary, these studies indicate that diploid cells lack a detectable frequency of gene amplification, whereas tumorigenic cells readily amplify DNA sequences. Present studies are designed to determine when the ability to amplify is acquired in converting a normal cell to a **neoplastic** cell and to isolate those molecular factors that regulate the process.

#### 4.3. Enhancement of Amplification

When studies were initiated to identify agents that could increase the frequency of gene amplification, it was estimated that amplification in the rodent model systems occurs in one out of every  $10^5$  cells. This is a very infrequent event if one wants to use biochemical or molecular techniques. The literature suggested that treatment of mammalian cells with carcinogenesis leads, among other things, to the production of HSRs and DMs (121). Because these two chromosomal abnormalities are linked to amplified genes, it was logical to test the possibility that treating mammalian cells with carcinogens could cause gene amplification. Studies with prototypical carcinogenic agents, such as UV light and chemical carcinogens, demonstrated that the incidence of gene amplification is increased by pretreatment with these agents (99, 122, 123). Other treatments, metabolic inhibitors, transient inhibitors of DNA synthesis, and hypoxia, could also increase the frequency of gene amplification (see Ref. 4 for a review). Most of these studies used the endogenous DHFR gene as a marker, but some also used other gene loci. In one study, MTX resistance, PALA resistance, and 5-fluorouridine resistance were cell be enhanced by treatment with ara-C, indicating

that enhancement of gene amplification is not confined to the DHFR locus (124). Viral sequences also undergo amplification that is enhanced by pretreatment with carcinogenic agents. Lavi and co-workers (113, 125) observed dramatic increases in viral sequences after the cells were treated with agents, such as benzopyrene, aflatoxin, methylnmethane sulfonate or a host of other carcinogens. The extent of enhanced of amplification ranged from a few fold to thousandfold. The basis for the enhanced gene amplification by carcinogen pretreatment is not yet known.

## 5. Considerations for Carcinogenesis

Carcinogenesis is believed to be a disruption of the homeostatic processes that control cell growth. The oncogenes identified to date contribute to this interpretation. Abnormalities in quantity or quality of [growth factors](#), growth factor receptors, signal transducers, and proteins that control entry into the [cell cycle](#) have each been implicated in the genesis of [neoplastic transformation](#). In addition, inactivation of those gene products that limit growth are also involved. These **tumor suppressor** activities contribute to the delicate balance that allows cells to grow when needed (tissue replacement, wound healing, etc.) and prevent cell growth when it is unneeded (terminal differentiation, [senescence](#)) (8). Mutations in each (or a combination) of these genes accumulates and results in progressive tumorigenicity. Continued accumulation of mutations allow the cells to acquire metastatic properties, such as invasiveness, angiogenesis, and drug resistance. Independent occurrences of these mutations allow tumor formation through alternate pathways (126). The process is believed to be multistep.

Gene amplification can be implicated in carcinogenesis in several different ways. Perhaps the most straightforward connection is the resulting drug resistance that occurs when detoxifying enzymes are overproduced, which is how gene amplification was discovered. The development of drug resistance remains one of the primary problems in fighting neoplasia. Obviously, elimination of drug resistance would do a great deal toward managing the disease. Experiments with tissue culture cells have shown us that a wide variety of loci may be amplified in mammalian cells. The amplification is usually manifested as an overproduction of the protein product that is targeted by the chemotherapeutic agent. Luria–Delbruck fluctuation analysis has demonstrated that amplification occurs spontaneously at a constant rate. It is the selective environment that allows them to be visualized. A study has compared the amplification rate in nontumorigenic and tumorigenic cells and found that the tumorigenic cells amplified the endogenous locus 100 times more than the nontumorigenic cell line (127). Restrictions on the loci that spontaneously amplify have not been encountered. Studies have also shown that more than one locus can be amplified at the same time (108).

The second connection between gene amplification and carcinogenesis involves the HSRs and DMs that have been observed in tumor cells for decades (121). Manifestations of gene amplification, HSRs and DMs have been widely reported in human tumor biopsy samples. Molecular studies have show that in many instances the genes amplified are oncogenes (97). It is assumed that the amplified oncogenes confer a growth advantage on the tumor cells that contain them. In some cases, amplification of specific genes correlates with the progression of some neoplasias. Seeger and Brodeur reported an increased extent and frequency of N-myc amplification in late stage neuroblastomas compared to early stage neuroblastomas. In this disease, the amplification of N-myc is used as a prognostic indicator for the severity of the disease. A similar story is emerging for the neu oncogene and its amplification in breast and ovarian cancer. A major unanswered question in the field of gene amplification is whether a nontransformed cell can ever undergo gene amplification or whether this type of genetic alteration is reserved for transformed cells alone.

The third implication of gene amplification in carcinogenesis relates to its nature. Gene amplification is a type of genetic rearrangement that (in mammals) is found only in tumor cells. Nowell (128) proposed a cellular scenario that accounts for the majority of observations made concerning neoplasia. He hypothesized that (at least in some cases) an early step in this process could be the acquisition of “genetic instability” (128). As the population of cells expands and generates a plethora



of heterogeneous cells, subpopulations would emerge from the selective environment in the host and become the substrate for further change. This process could lead to the evolution of a cell population that is well suited to growth in the host tissue and can acquire the molecular characteristics that would lead to malignancy.

Nowell based his hypothesis on several pointed observations: (1) Oftimes, the karyotype of tumorigenic cells is abnormal. Rearrangements, such as deletions, aneuploidy, translocations, and homogeneously staining regions (HSRs) and double minute (DMs) chromosomes (known manifestations of gene amplification) are common. It is often easy to detect changes from the original karyotype. (2) Once a tumorigenic subpopulation emerges, further karyotypic changes are common. If Nowell's hypothesis is correct, elucidation of the molecular mechanism that underlies the continual change could provide a pivotal piece of information in our studies to understand the neoplastic process. To approach this question one could measure one or several of the types of genomic rearrangements that are detected in tumor cells. The hypothesis predicts that normal or nontransformed cells should display this abnormality at a lower frequency (or rate) than tumorigenic cells. This hypothesis has been addressed using gene amplification as the molecular marker for "genetic instability" ([115](#), [118](#), [120](#)) and, as indicated before, a vast difference in amplification ability exists. These results are consistent with the Nowell hypothesis. In two cases, deletion formation and recombination with an exogenous piece of DNA (Finn et al., 1989), the genetic rearrangement is more common in tumor cells.

## 6. Summary

The literature suggests that when gene amplification occurs in "normal" tissues, it is developmentally regulated. This evidence is compiled mostly from studies on *Xenopus* and *Drosophila* (see previous section). In higher organisms, the documentation of gene amplification as a developmental event is lacking. At present, we do not know if gene amplification can be developmentally programmed in mammalian cells.

Sporadic amplification occurs in unicellular organisms, such as bacteria and yeast, but it is lacking in the normal somatic tissues of higher eukaryotes. Several reports of sporadic amplification in the germ line cells of several organisms have been reported and shown to be heritable. In all of these cases, the phenotype demonstrated increased resistance to an environmental toxin ([64](#), [65](#), [129](#)). The extensive documentation of sporadic amplification in neoplastic tissues raises questions about when the neoplastic cell acquires the ability to amplify. Are the changes that initiate neoplasia the same that initiate the ability to amplify? In this regard, the stimulation of gene amplification by agents that cause cancer becomes an important observation. An insult, such as UV light, can initiate cells into the process that leads to tumor formation. Likewise, the same event can initiate and increase amplification frequency in transformed cells. Can these same agents cause gene amplification in "normal" diploid cells? These and similar questions must be answered before we can determine the extent of contributions that gene amplification holds for carcinogenesis.

## Bibliography

1. J. L. Hamlin, J. D. Milbrandt, N. H. Heintz, and J. C. Azizkhan (1984) *Int. Rev. Cytol.* **90**, 31–82.
2. G. R. Stark and G. M. Wahl (1984) *Annu. Rev. Biochem.* **53**, 447–491.
3. G. R. Stark (1986) *Cancer Surv.* **5**, 1–23.
4. R. T. Schimke (1988) *J. Biol. Chem.* **263**, 5989–5992.
5. E. Farber (1984) *Cancer Res.* **44**, 4217–4223.
6. J. M. Bishop (1987) *Science* **235**, 305–311.
7. E. R. Fearon and B. Vogelstein (1990) *Cell* **61**, 567–579.
8. J. A. Boyd and J. C. Barrett (1990) *Pharmacol. Ther.* **46**, 469–486.
9. H. D. King (1908) *J. Morphol.* **19**, 369–438.

10. J. G. Gall (1968) *Proc. Natl. Acad. Sci. USA* **60**, 553–560.
11. D. D. Brown and I. B. Dawid (1968) *Science* **160**, 272–280.
12. J. G. Gall, H. C. MacGregor, and M. E. Kidston (1969) *Chromosoma* **26**, 169–187.
13. M. R. Kalt and J. G. Gall (1974) *J. Cell Biol.* **62**, 460–472.
14. A. P. Bird (1978) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 1179–1183.
15. D. Hourcade, D. Dressler, and J. Wolfson (1973) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 537–550.
16. J. D. Rochaix, A. P. Bird, and A. Bakken (1974) *J. Mol. Biol.* **87**, 473–487.
17. C. H. Thiebaud (1979) *Chromosoma* **73**, 37–44.
18. J. D. Rochaix and A. P. Bird (1975) *Chromosoma* **52**, 317–327.
19. M. C. Yao and M. A. Gorovsky (1974) *Chromosoma* **48**: 1–18.
20. M. C. Yao and J. G. Gall (1974) *Cell* **12**, 121–132.
21. M. C. Yao, E. Blackburn, and J. G. Gall (1978) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 1293–1296.
22. J. G. Gall and J. D. Rochaix (1974) *Proc. Natl. Acad. Sci. USA* **71**, 1819–1823.
23. M. C. Yao (1981) *Cell* **24**, 765–774.
24. M. E. Breuer and C. Pavan (1955) *Chromosoma* **7**, 371–386.
25. D. M. Glover et al. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 2947–2951.
26. C. J. Bostock and A. T. Sumner (1978) *The Eukaryotic Chromosome*, North Holland, Amsterdam, pp. 256–259.
27. C. Pavan and B. da Cunha (1969) *Genetics (suppl)* **61**, 289–304.
28. C. Pavan (1965) *Natl. Cancer Inst. Monogr.* **18**, 309–323.
29. B. Da Cunha, C. Pavan, C. Morgante, and C. Garrido (1969) *Genetics (suppl)* **61**, 304–333.
30. A. C. Spradling and A. P. Mahowald (1980) *Proc. Natl. Acad. Sci. USA* **77**, 1096–1100.
31. A. C. Spradling (1981) *Cell* **27**, 193–201.
32. A. C. Spradling and A. P. Mahowald (1981) *Cell* **27**, 203–209.
33. W. Nagl, J. Hendon, and W. Rucker (1972) *Cell Differentiation* **1**, 229–237.
34. J. Van't Hof and C. A. Bjerknes (1982) *Mol. Cell. Biol.* **2**, 339–345.
35. W. L. Wardell (1977) *Plant Physiol.* **60**, 885–891.
36. A. Novick and T. Horiuchi (1961) *Cold Spring Harbor Symp. Quant. Biol.* **21**, 239–245.
37. T. Horiuchi, S. Horiuchi, and A. Novick (1963) *Genetics* **48**, 157–169.
38. T. Horiuchi, J. I. Tomizawa, and A. Novick (1962) *Biochem. Biophys.* **55**, 152–163.
39. J. Langridge (1969) *Mol. Gen. Genetics* **105**, 74–83.
40. S. Normark et al. (1977) *J. Bacteriol.* **132**, 912–922.
41. P. Anderson and J. Roth (1981) *Proc. Natl. Acad. Sci. USA* **78**, 3113–3117.
42. P. Anderson and J. R. Roth (1977) *Annu. Rev. Microbiol.* **31**, 473–505.
43. S. W. Emmons and J. O. Thomas (1981) *J. Mol. Biol.* **91**, 147–152.
44. S. W. Emmons, V. MacCosham, and R. L. Baldwin (1975) *J. Mol. Biol.* **91**, 133–146.
45. T. Edlund and S. Normark (1981) *Nature* **292**, 269–271.
46. T. D. Tlsty, A. M. Albertini, and J. H. Miller (1984a) *Cell* **37**, 217–224.
47. S. K. Whoriskey et al. (1987) *Genes Dev.* **1**, 227–237.
48. A. Kozinski et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5064–5068.
49. J. Meyer and S. Lida (1979) *Mol. Gen. Genet.* **176**, 209–219.
50. S. Fogel, J. W. Welch, G. Cathala, and M. Karin (1983) *Curr. Genet.* **7**, 347–355.
51. M. Karin et al. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 337–341.

52. C. H. Browning (1954) The chemotherapy of trypanosomic infections. *Ann. NY Acad. Sci.* **59**, 198–213.
53. W. Peters (1974) *Ciba Found. Symp.* **20**, 309–334.
54. E. M. Rollo (1980) *Drugs Used in the Chemotherapy of Malaria in the Pharmacological Basis of Therapeutics* (A. G. Gilman, L. S. Goodman, and A. Gilman, eds.) Macmillan, New York, 6th ed, pp. 1038–1069.
55. S. M. Beverley, J. A. Coderre, D. V. Santi, and R. T. Schimke (1984) *Cell* **38**, 431–443.
56. F. M. Ritossa, K. C. Atwood, and S. Spiegelman (1966) *Genetics* **54**, 819–834.
57. F. M. Ritossa and G. Scala (1969) *Genetics* **61**, 305–37.
58. F. M. Ritossa (1968) *Proc. Natl. Acad. Sci. USA* **60**, 509–516.
59. F. M. Ritossa, C. Boncinelli, F. Graziani, and L. Polita (1971) *Proc. Natl. Acad. Sci. USA* **68**, 1580–1584.
60. K. D. Tartof (1974) *Science* **171**, 294–297.
61. W. Kunz, C. Grimm, and G. Franz (1982) *The Cell Nucleus*, Academic Press, **Vol. XII**, pp. 155–184.
62. J. D. Procunier and K. D. Tartof (1978) *Genetics* **88**, 67–79.
63. E. Otto, S. McCord, and T. D. Tlsty (1989) *J. Biol. Chem.* **264**, 3390–3396.
64. G. Maroni, J. Wise, J. E. Young, and E. Otto (1987) *Genetics* **117**, 739–744.
65. C. Mouches et al. (1986) *Science* **233**, 778–780.
66. C. A. Cullis (1983) *CRC Crit. Rev. Plant Sci.* **1**, 117–131.
67. C. A. Cullis (1977) *Heredity* **38**, 129–154.
68. C. A. Cullis (1979) *Heredity* **42**, 237–246.
69. W. F. Flintoff, S. V. Davidson, and L. Siminovich (1976a) *Somat. Cell. Genet.* **2**, 245–261.
70. W. F. Flintoff, S. M. Spindler, and L. Siminovich (1976b) *In Vitro* **12**, 749–757.
71. D. A. Haber, S. M. Beverly, M. L. Kiely, and R. T. Schimke (1981) *J. Biol. Chem.* **256**, 9501–9510.
72. F. M. Sirotnak, D. M. Moccio, L. E. Kelleher, and L. J. Goutas (1981) *Cancer Res.* **41**, 4447–4452.
73. F. W. Alt, R. E. Kellems and R. T. Schimke (1976) *J. Biol. Chem.* **251**, 3063–3074.
74. M. T. Hakala, S. F. Zakrzewski, and C. A. Nichol (1961) *J. Biol. Chem.* **236**, 952–958.
75. J. L. Biedler and B. A. Spengler (1976) *Science* **191**, 185–187.
76. F. W. Alt, R. E. Kellems, J. R. Bertino, and R. T. Schimke (1978) *J. Biol. Chem.* **253**, 1357–1361.
77. C. Tyler-Smith and C. J. Bostock (1981) *J. Mol. Biol.* **153**, 219–236.
78. B. J. Dolnick et al. (1979) *J. Cell. Biol.* **83**, 394–402.
79. W. F. Flintoff et al. (1983) *Mol. Cell. Biol.* **2**, 275–285.
80. P. W. Melera, J. A. Lewis, J. L. Biedler, and C. Hession (1980) *J. Biol. Chem.* **255**, 7024–7082.
81. T. D. Kempe, E. A. Swyryd, M. Bruist, and G. R. Stark (1976) *Cell* **9**, 541–550.
82. G. M. Wahl, R. A. Padgett, and G. R. Stark (1979) *J. Biol. Chem.* **254**, 8679–8689.
83. D. P. Suttle and G. R. Stark (1979) *J. Biol. Chem.* **254**, 4602–4607.
84. D. P. Suttle (1983) *J. Biol. Chem.* **258**, 7707–7713.
85. J. J. Kanalas and D. P. Suttle (1984) *J. Biol. Chem.* **259**, 1848–1853.
86. L. R. Beach and R. D. Palmiter (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2110–2114.
87. L. McConlogue and P. Coffino (1983) *J. Biol. Chem.* **258**, 12083–12086.
88. D. J. Chin et al. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 1185–1189.

89. I. B. Roninson et al. (1984) *Nature* **309**, 626–628.
90. I. B. Roninson et al. (1984) *Proc. Natl. Acad. Sci. USA* **83**, 4538–4542.
91. C. Rossana, L. G. Roa, and L. F. Johnson (1982) *Mol. Cell. Biol.* **2**, 1118–1125.
92. T. S. Su, H. G. O. Bock, W. E. O'Brien, and A. L. Beaudet (1981b) *J. Biol. Chem.* **256**, 11826–11831.
93. T. S. Su, A. L. Beaudet, and W. E. O'Brien (1981a) *Biochemistry* **20**, 2956–2960.
94. F. Baskin, R. N. Rosenberg, and V. Dev (1981) *Proc. Natl. Acad. Sci. USA* **78**, 3654–3658.
95. R. N. Johnston, S. M. Beverley, and R. T. Schimke (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3711–3715.
96. P. Kavathus and L. Herzenberg (1983) *Nature* **306**, 385–387.
97. M. Schwab et al. (1983) *Nature* **305**, 245–248.
98. P. C. Brown, R. N. Johnson, and R. T. Schimke (1983a) *Gene Struct. Regul. Dev.* 197–212.
99. T. D. Tlsty, P. C. Brown, and R. T. Schimke (1984b) *Mol. Cell. Biol.* **4**, 1050–1056.
100. H. Rath, T. D. Tlsty, and R. T. Schimke (1984) *Cancer Res.* **44**, 3303–3306.
101. G. Balaban-Malenbaum and F. Gilbert (1980) *Cancer Genet. Cytogenet.* **2**, 339–348.
102. R. J. Kaufman and R. T. Schimke (1981) *Mol. Cell. Biol.* **1**, 1069–1076.
103. B. Trask and J. Hamlin (1989) *Genes Dev.* **3**, 1913–1925.
104. J. K. Cowell (1982) *Annu. Rev. Genet.* **16**, 21–59.
105. J. L. Biedler, P. W. Melera, and B. A. Spengler (1980) *Cancer Genet. Cytogenet.* **2**, 47–60.
106. N. A. Federspiel, S. M. Beverley, J. W. Schilling, and R. T. Schimke (1984) *J. Biol. Chem.* **259**, 9127–9140.
107. J. Zeig et al. (1983) *Mol. Cell. Biol.* **3**, 2089–2098.
108. E. Giulotto, I. Saito, and G. R. Stark (1989) *EMBO J.* **5**, 2115–2951.
109. S. M. Carroll et al. (1987) *Mol. Cell. Biol.* **7**, 1740–1750.
110. B. J. Mauer et al. (1987) *Nature* **327**, 434–437.
111. G. M. Wahl (1989) *Cancer Res.* **49**, 1333–1340.
112. D. D. Von Hoff et al. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 4804–4808.
113. S. Lavi (1981) *Proc. Natl. Acad. Sci. USA* **78**, 6144–6148.
114. R. J. Kaufman, P. C. Brown, and R. T. Schimke (1981) *Mol. Cell. Biol.* **1**, 1084–1093.
115. R. Sager, I. Gadi, L. Stephens, and C. Grabowy (1985) *Proc. Natl. Acad. Sci. USA* **82**, 7015–7019.
116. G. M. Wahl, R. de Saint Vincent, and M. L. DeRose (1984) *Nature* **307**, 516–520.
117. C. S. Gasser, C. C. Simonsen, and R. T. Schimke (1982) *Proc. Natl. Acad. Sci. USA* **79**, 6522–6526.
118. E. Otto, J. E. Young, and G. Maroni (1989) *Proc. Natl. Acad. Sci. USA* **83**, 6025–6029.
119. J. Wright et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 1791–1795.
120. T. D. Tlsty (1990) *Proc. Natl. Acad. Sci. USA* **87**, 3132–3136.
121. P. E. Barker (1982) *Cancer Genet. Cytogenet.* **5**, 81–94.
122. T. D. Tlsty, P. C. Brown, R. Johnston, and R. T. Schimke (1982) In *Gene Amplification* (R. T. Schimke, ed.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York pp. 231–238.
123. G. Rice, C. Hoy, and R. T. Schimke (1986) *Proc. Natl. Acad. Sci. USA*.
124. B. Goz, P. Carl, and T. D. Tlsty (1989) *Mol. Pharmacol.* **36**, 360–365.
125. S. Lavi and S. Etkin (1981) *Carcinogenesis* **2**, 417–423.
126. L. Foulds (1975) *Neoplastic Development*, Vol. **1**, Academic Press, New York.
127. T. D. Tlsty, B. Margolin, and K. Lum (1989) *Proc. Natl. Acad. Sci. USA* **86**, 9441–9445.

128. P. Nowell (1976) *Science* **194**, 23–28.

129. C. Prody et al. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 690–694.

## Gene Cluster

In many cases, functionally related **genes**, each of which codes for a separate **protein**, are clustered physically on a **chromosome** in a *gene cluster*. In **bacteria**, the **gene expression** of all members of the cluster can be under the control of a single **promoter**, and a single **operator** controls their expression. This operator and the accompanying cluster of structural genes constitute an integrated topological and functional unit, known as an **operon**.

In **eukaryotes**, related genes produced by **gene duplication** are often found clustered on a chromosome. The clusters can be rather complex, ranging from two or a few adjacent related genes, as in the case of the **apolipoprotein A-1-III** and **A-IV** genes on human chromosome 11, to a tandem array of several hundred identical genes. Well-studied cases of clustering include the genes for **globins**, the **major histocompatibility complex**, and **homeotic** proteins. The members of the cluster may acquire slightly different functions or are used in different stages of **development**, such as fetal and adult **hemoglobins**.

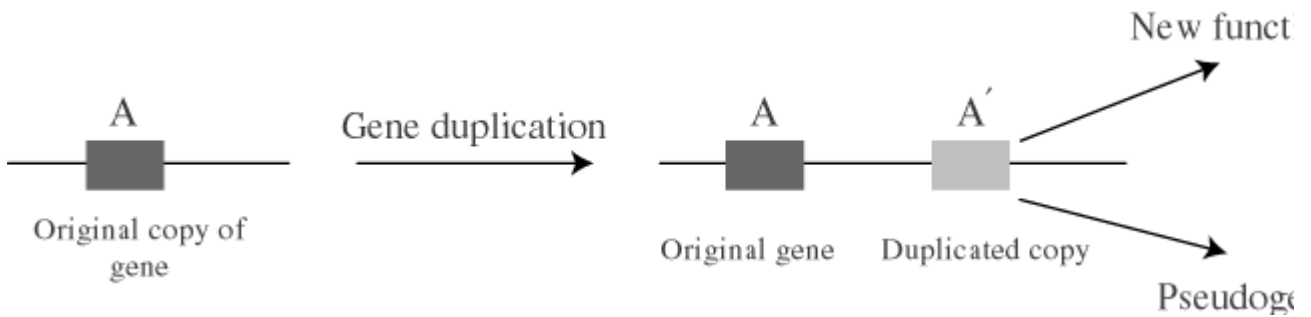
## Gene Dosage

Gene dosage refers to the number of copies of a given **gene** in a cell or in a **nucleus**. An increase in the gene dosage may result in the production of a higher level of the product of that gene, unless the gene is subject to autonomous regulation. Where sexes differ in the number of sex chromosomes, a **dosage compensation effect** may keep the levels of expression similar.

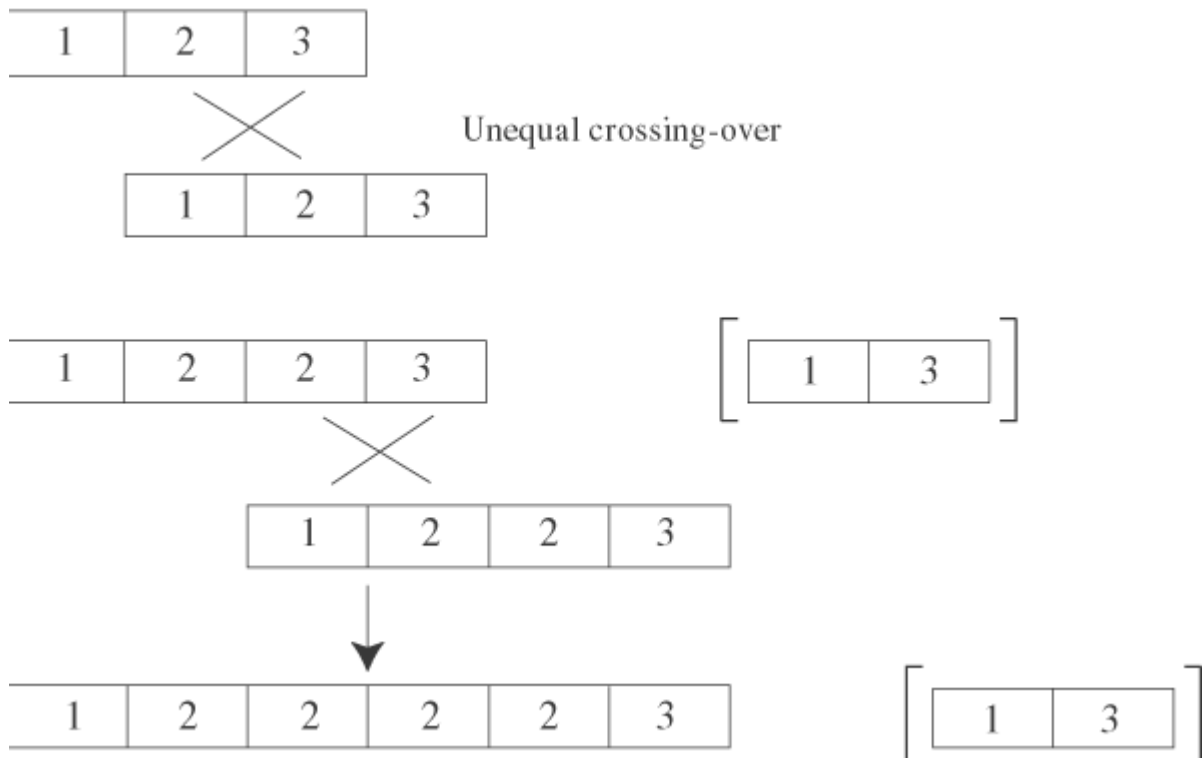
## Gene Duplication

Gene duplication is the creation of a copy of a **gene** within the (Fig. 1). Genes can be duplicated through several types of mechanisms. Among these are random **breakage and reunion**, **unequal crossing over** (Fig. 2), **retrotransposition** and **overreplication**. Gene duplication is a major source of functional genetic diversity. This diversity is achieved through the divergence and the acquisition of new functions of the copies over time (Fig. 1).

**Figure 1.** Gene duplication. A duplicated copy can acquire a new function or become a pseudogene.

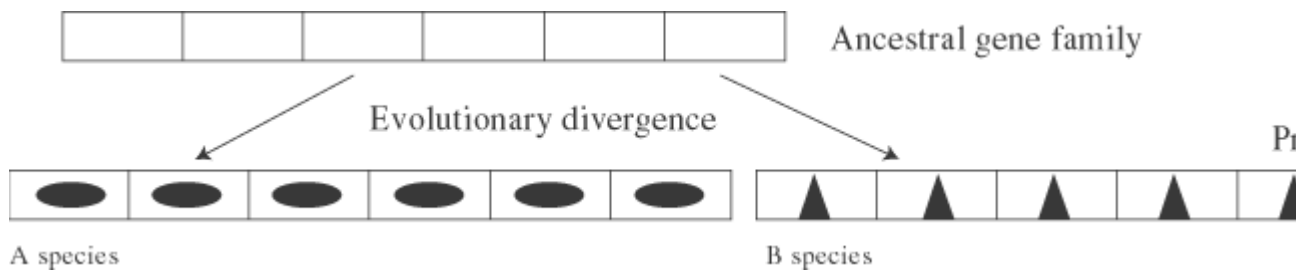


**Figure 2.** Creation of a gene family through gene duplication by unequal crossing-over.



Duplication can involve either genes or , as well as entire , and be either complete or partial. Both genome duplication and gene duplication have been important in evolution. Such duplication generates a related set of genes, or , enabling an increase of genomic diversity by creating new gene functions. It is well known that the genome contains a variety of gene families in a repeated or scattered fashion. These gene families are considered to have emerged evolutionarily from their ancestral form through gene duplication (Fig. 3). Once gene duplication has taken place by unequal crossing over of one chromosome, the other chromosome may lack certain sets of genes. Therefore, it follows that the number of gene members within a particular gene family may have changed during evolution, both increasing and decreasing.

**Figure 3.** Concerted evolution in gene families of species A and B.



Once a particular nucleotide substitution takes place in a given gene member, the gene member can propagate within the gene family through gene duplication during evolution. When gene duplication takes place through unequal crossing over, propagation of a particular gene member within the gene family in a species is faster than the evolutionary divergence of the gene families between different species. In this case, the similarity between different gene members within a given gene family in the same species becomes greater than that between corresponding gene families in different species (Fig. 3). This type of evolution is called “concerted”

**Pseudogenes** are often created after gene duplication (2). Pseudogenes are dead genes in which a copy of a gene becomes non-functional due to a mutation-causing malfunction in the DNA sequence (Fig. 1). A gene family contains a number of pseudogenes. This is also the case for the **major histocompatibility complex** gene family. Ohno (3) proposed that one copy of a gene duplication can enjoy the freedom to acquire mutations that differs from the original gene. The original gene can keep its function, even if the counterpart becomes a pseudogene. This theory implies that gene duplication can be an important source of providing new genes to the genome.

As the projects for sequencing entire genomes advance, the evolution of genome structure is starting to receive particular attention. It is known that genome size has frequently increased as organisms have evolved. To explain the increase in genome size, two major alternative hypotheses have been proposed: the genome gradual theory, which asserts that the genome size increased due to the addition of new genes created by gene duplication, and the “genome duplication theory,” which contends that the entire genome has been duplicated, resulting in almost doubling in size. Although deletions occur in a part of the genome, the change of gene number has not been taken place over an extended period of time.

## Bibliography

“Gene Duplication” in *Genetics*, Vol. 2, pp. 977–978, by T. Gojobori; “Gene Duplication” in *Genetics* (online), posted March 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. T. Ohta, (1980) *Evolution and Variation of Multigene Families. Biomathematics*, Springer-Verlag, Berlin.
2. W. H. Li, T. Gojobori, and M. Nei, (1981) *Nature* **292**, 237–239.
3. S. Ohno (1970) *Evolution by Gene Duplication*, Springer-Verlag, Berlin.

## Gene Families

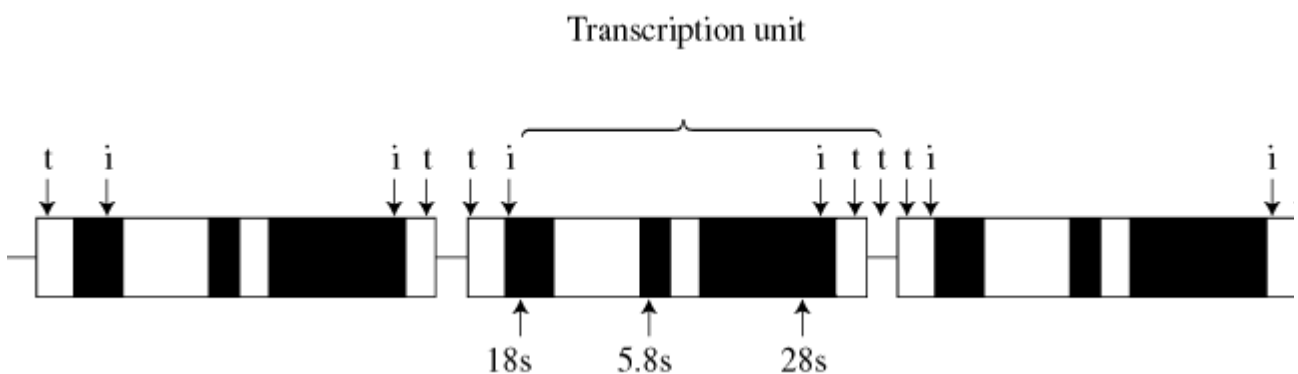
Gene families are groups of **genes** with similar sequences that have been derived from a common evolutionary ancestor. Many of the structural genes in eukaryotes are members of gene families. The number of members in a gene family can range from two to hundreds of thousands per genome. For example, the **globin** gene family has more than 500,000 copies, displaying only minor variation in their nucleotide sequences. A genome can contain hundreds of gene families, and members of a family can be either scattered over several chromosomes or clustered together on a single chromosome.

The multiple genes of a family would be expected to evolve independently. In this case, copies within a species should have the same amount of divergence as those among different species. Those within a species are usually more homogeneous, however, implying that gene family members within a species are not evolving independently and that the family members within a particular species have descended from a common ancestor after species divergence. This would suggest that there are mechanisms for maintaining homogeneity among gene family members. These mechanisms appear to be , **unequal crossing over**, and **gene conversion**. The evolution of gene families through these mechanisms is called concerted evolution. (See ).

Research in the area of gene families focuses on questions such as how the copy number of members changes during evolution and how these copies move to different locations. Further research interest in this area includes investigating the effect of copy number changes on gene diversity and identifying what mechanisms preserve similarity among family members (1).

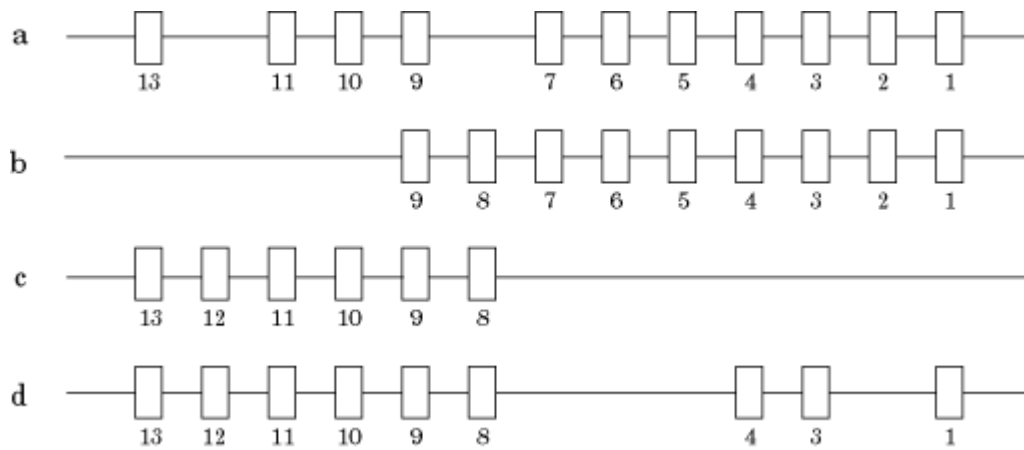
The evolutionary importance of gene families depends heavily on their specific function. In the case of **ribosomal** RNA (rRNA) gene families, for example, it is known that the gene members are highly homogeneous (Fig. 1) (2). This can be explained by the fact that of genes requires a functional ribosome. Thus, all the many gene members in this family are preserved to be highly homogeneous during evolution (3). On the other hand, the form a gene family but are known not to be homogeneous, but relatively heterogeneous. The homeotic genes are directly involved in the control of genes related to the development of organisms. In the case of mammalian homeotic gene families, a single Hox gene cluster generally contains 13 genes, and this gene cluster is repeated four times (3). The nucleotide sequences of all gene members are not homogeneous except a small portion (about 60 codons) of the gene called a **homeobox** (Fig. 2) (5). The heterogeneity among the gene members in this family is due to each gene having a specialized role in their cascade to control the development of organisms.

**Figure 1.** Animal genomes contain multiple tandem copies of the pre-ribosomal RNA transcription units (2). The three solid boxes of each transcription unit encode the 18 S rRNA, 5.8 S rRNA, and 28 S rRNA genes, respectively.



**Figure 2.** Genomic organization of *Antennapedia* class of **homeobox** genes of humans (5).





## Bibliography

“Gene Families” in , Vol. 2, pp. 978–979, by T. Gojobori; “Gene Families” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

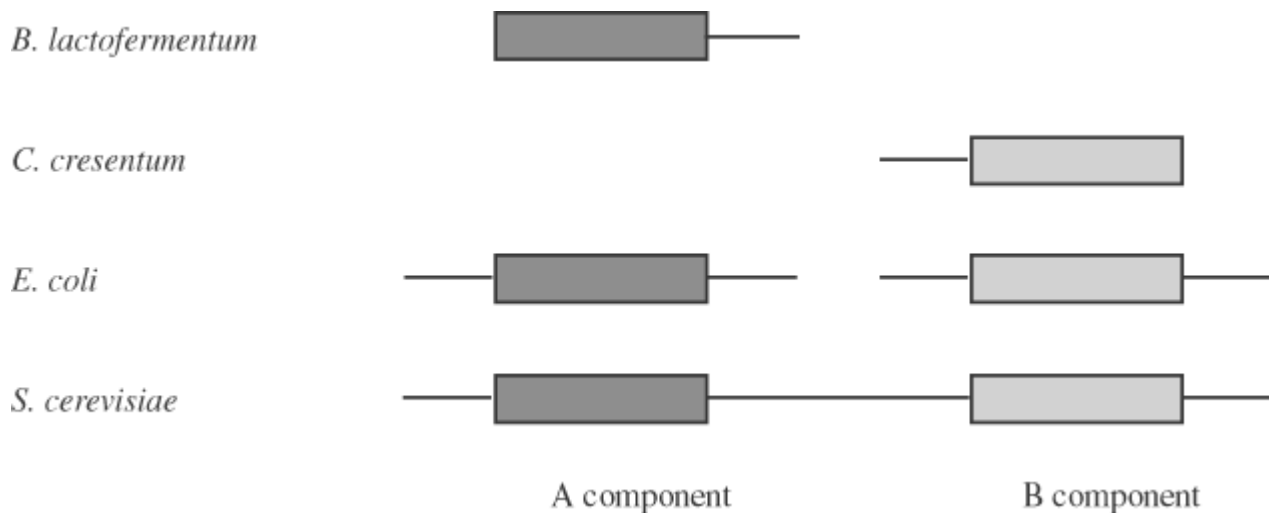
1. T. Ohta, (1980). *Evolution and Variation of Multigene Families*, Springer-Verlag, Berlin.
2. H. Lodish, D. Baltimore, A. Berk, S. L. Zipursky, P. Matsudaira, and J. Darnell, (1995). *Molecular Cell Biology*, 3rd. ed., New York.
3. N. Arnheim, M. Krystal, R. Schmickel, G. Wilson, O. Ryder, and E. Zimmer, (1980). *Proc. Natl. Acad. Sci. USA* **77**; 7323–7327.
4. C. Kappan and F. H. Ruddle, (1993) *Curr. Opin. Gen. Dev.* **3**; 931–938.
5. W. H. Li. (1997). *Molecular Evolution*, Sinauer Associates, Sunderland, MA.

## Gene Fusion

Gene fusion is a phenomenon in which two different **genes**, or parts of the genes, are fused to each other during [evolution](#), resulting in the creation of a new gene. When two particular genes are compared between two species, they can be fused in one species and separate in the other. To say with certainty that a gene is a fused gene, the separate genes should be identified as being the ancestral form by careful examination.

**Unequal crossing over** may be one of the molecular mechanisms for causing gene fusion by eliminating chromosomal material between two formerly separate genes. For example, the union in fungi between A and B components of the tryptophan synthetase enzyme (see [Operons](#)), which are normally separate in bacteria, is believed to have occurred through such gene fusion events (Fig. [1](#)).

**Figure 1.** The gene fusion event of A and B components on tryptophan synthetase gene ([1](#)).



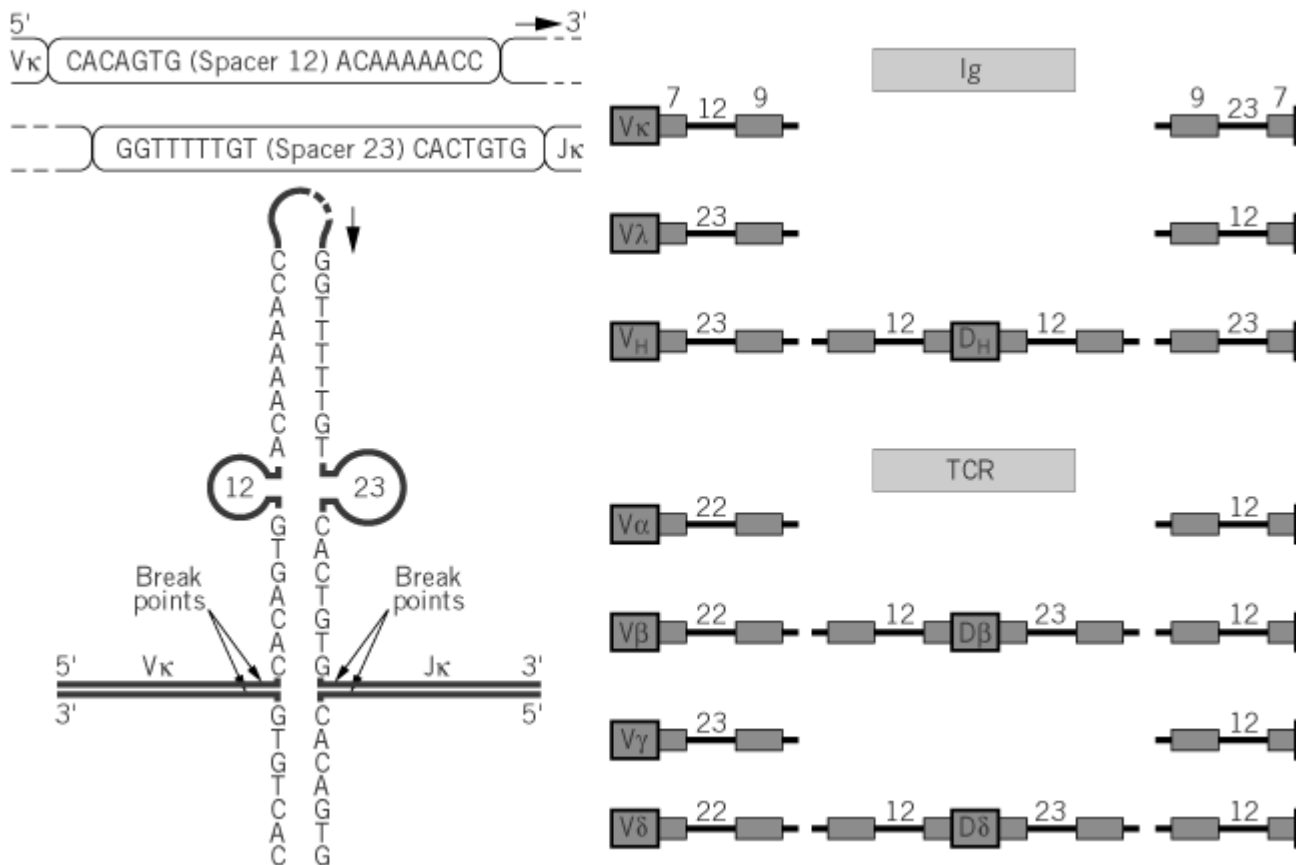
## Bibliography

1. D. M. Burns, V. Horn, J. Paluh, and C. Yanofsky (1990) *J. Biol. Chem.* **5**, 2060–2069.

## Gene Rearrangement

The unique feature of [immunoglobulin](#) gene organization is that these genes need to be rearranged before becoming functional. These rearrangements allow the immune system to generate a very large [repertoire](#) by [recombination](#) and random association of a limited number of gene segments (ie, between 200 and 300 in humans), leading to at least 10 million discrete immunoglobulins having different specificities. In the light chains, the **variable region** results from the random association of VL (Vk or Vl) and JL (Jk or Jl) genes. The heavy chains are more complex, because their diversity results from the random association of VH, D, and JH genes. Immunoglobulin genes are rearranged exclusively in the **B-cell** lineage, and these events constitute the hallmark of the B-cell [differentiation](#) in the bone marrow. Similarly, genes for the [T-cell receptor](#) (TCR) rearrange exclusively in [T cells](#) by a very similar process. Ig gene rearrangements require precise signals on the DNA, known as *recombination signal sequences* or (RSSs), that flank each of the V, D, and J gene segments and are recognized by highly specific [recombinases](#) encoded by two genes, termed RAG1 and RAG2 (for recombina-se activating genes). RSSs consist of two **palindromic** sequences that are highly conserved, although not identical, between the different loci; one is a heptamer and one is a nonamer, separated by a spacer region 12 or 23 nucleotides long (Fig. [1](#)). The two RSSs that participate to a given joint have always spacers of different lengths.

**Figure 1.** Recombination signal sequences of Ig and TCR genes. (a) The 3' end of a Vk gene is followed by an RSS that is complementary to the RSS 5' of the Jk gene onto which Vk will be connected by the recombinase. (b) Schematic organization of the RSS in Ig and TCR genes.

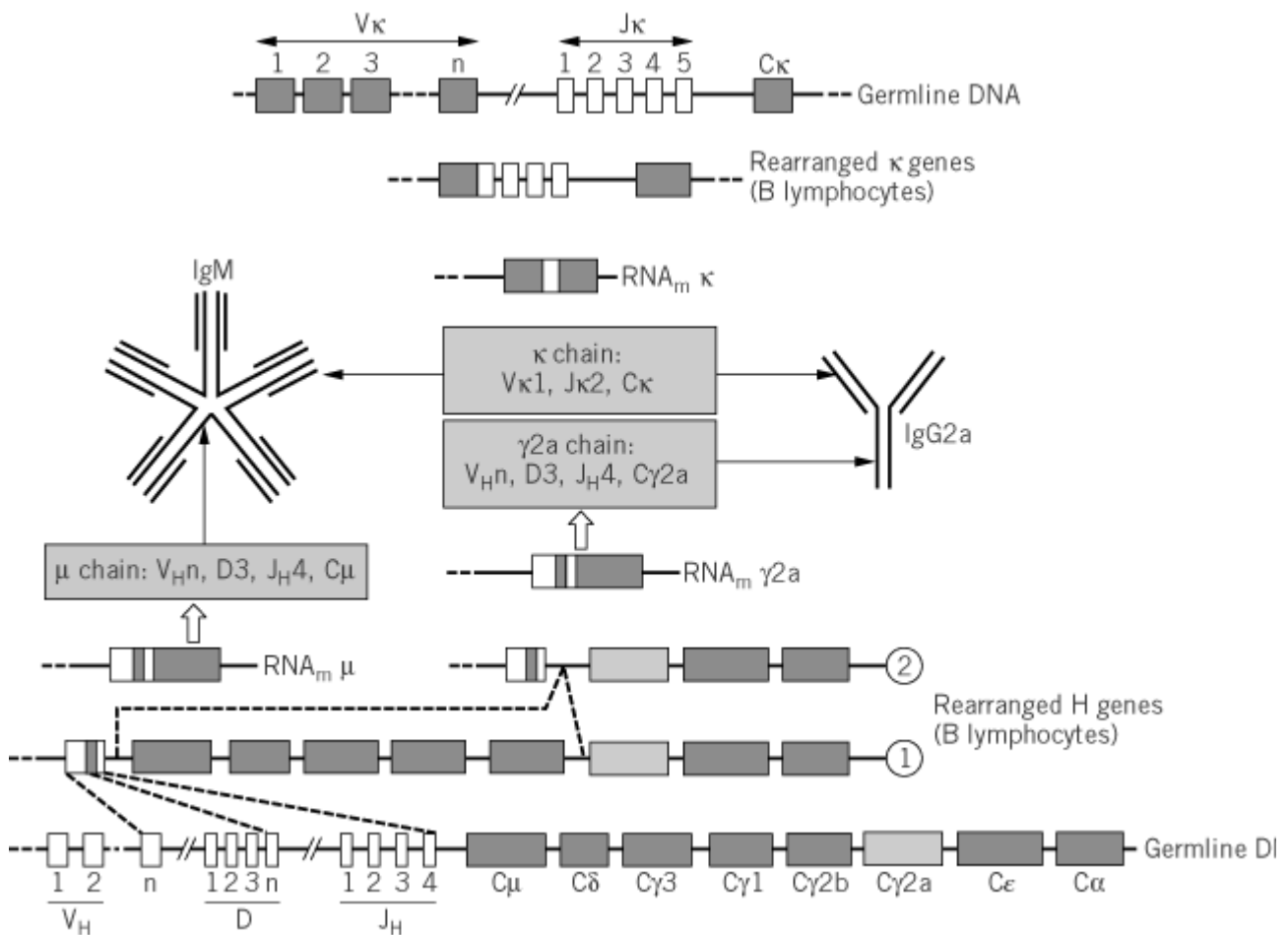


The first recombination event, DH to JH, takes place in the early proB cells and is rapidly followed by the rearrangement of one of the VH segments to D-J. Another enzyme, terminal deoxynucleotidyl-transferase, is also active in proB cells and adds nucleotides in a random fashion during the joining of D to J and of V to D-J rearrangements. Addition of these non-germline-encoded nucleotides is known as N-diversity and considerably amplifies the repertoire, because both the sequence and length of VH and VL complementarity-determining region 3 (CDR3) are modified by this mechanism. Rearrangements are strictly regulated by their end product—that is, the chain they encode once they become functional. Thus, once a m chain is synthesized, it will block the rearrangement of the second IGHV allele. This requires that the m chain be expressed as a membrane protein (mm). In a mouse rendered transgenic for a membrane VDJ-Cm chain gene, endogenous rearrangements at the IGHV locus are blocked, so that B cells will make only the transgenic m. When the same experiment is performed with a VDJ-Cm gene without the M exons (thereby encoding only the soluble form of the m chain), endogenous rearrangements are no longer blocked. The fact that m regulates rearrangement only when expressed as mm strongly suggests that it operates at the cell surface. In fact, mm is coexpressed with the so-called surrogate light chain, composed of the product of two genes, I5 and VpreB, and expressed specifically in proB and preB cells (see [B Cell](#)). I5 is highly **homologous** to I and is **disulfide-bonded** to the m chain by its penultimate **cysteine** residue, whereas VpreB resembles a V Ig domain. The surrogate light chain (YL) has been shown to contribute efficiently to the regulation of [allelic exclusion](#). The negative feedback that prevents the second allele from rearranging is probably a consequence of the turning off of the terminal deoxynucleotidyl-transferase and RAG genes. The RAG genes, but not the transferase, are reactivated as the preB cells cycle and induce light-chain gene rearrangement, with the same hierarchy of control on the second allele whenever the first has rearranged. The IgK locus is activated first. If both k alleles fail to produce a functional rearrangement, the IGL enters the game, with the same type of control. Finally it should be observed that the overall process of rearrangement is not very efficient, with only about 10% final success. This is obviously the price to

pay for the complexity and partial randomness of the system.

Once both H and L genes have successfully rearranged, the resulting immunoglobulin is expressed at the B-cell surface. When stimulated by an antigen, specific clones proliferate and colonize the secondary lymphoid organs (lymph nodes, tonsils, spleen, etc.), where they become organized in germinal centers, in close contact with other partners of the immune response (ie, dendritic cells and T cells). Contact with T helper (Th) cells will promote the last step of gene rearrangement, called *isotype switching*, by which the VDJ gene region will be associated with a new constant heavy-chain gene. The classical most frequent switch that occurs in immune responses is [IgM](#) → [IgG](#). A schematic example of all the rearrangements that lead to the production of an IgG molecule is given in Figure 2.

**Figure 2.** Rearrangement events that lead from Ig germline genes to synthesis of IgM and IgG molecules in the mouse model. Recombination first occurs at the IGHV locus, to generate a V(D)J combination (here chosen as V<sub>Hn</sub>-D3-JH4) that is transcribed with C<sub>μ</sub> and processed as a m mRNA, resulting in the expression of a m chain that will associate with k chain (here V<sub>κ1</sub>-J<sub>κ2</sub>-C<sub>κ</sub>) that results from the second wave of rearrangement targeted at the IGK locus. IgM is first expressed as a monomer at the cell surface of immature B cell (not shown). After antigenic stimulation, the pentameric form of IgM antibody is secreted by plasma cells and is rapidly replaced by an IgG (here IgG2a), which occurs by class switching (second line from bottom). Note that the switch will not change the VDJ region, so the antibody specificity is maintained.



See also entries [Recombinase](#) and [T-Cell Receptor \(TCR\)](#).

Suggestions for Further Reading

D. G. Schatz, M. A. Oettinger, and M. S. Schlissel (1993) V(D)J recombination: molecular biology and regulation. *Annu. Rev. Immunol.* **10**, 359–430.

J. Chen and F. W. Alt (1993) Gene rearrangement and B cell development. *Curr. Opin. Immunol.* **5**, 194–200.

## Gene Recruitment

Gene recruitment is a phenomenon in which a particular **gene** becomes used during [evolution](#) as a gene with a totally different function. The term “gene recruitment” was coined because a gene has evolved as if it had been recruited to exhibit a different or another function. The most remarkable example of gene recruitment occurs with the crystallin genes (Table 1). Most of the crystallin genes that are expressed in the eye were found to be expressed in other tissues as [enzymes](#). The exact same gene is expressed and used as a structural protein in one tissue, whereas it is expressed as an enzyme in other tissues. It is now thought that the enzymatic form of this protein should be considered to be the ancestral form and that it became a crystallin by gene recruitment. It seems that change of tissue specificity of gene expression is a prerequisite for gene recruitment.

**Table 1. Eye crystallins and Their Relationships to Enzymes and Stress Proteins (1)**

| Crystallin                       | Distribution                              | [Related] or Identical                       |
|----------------------------------|---|--|
| Ubiquitous Stress Crystallins    |   |  |
| a                                | All vertebrates                           | Small heat shock proteins (aB)               |
| β }<br>γ }                       | All vertebrates (embryonic gnot in birds) | [ <i>Schistosoma mansoni</i> antigen]        |
|                                  |   | [ <i>Myxococcus xanthus</i> protein S]       |
|                                  |   | [ <i>Physarum polycephalum</i> spherulin 3a] |
| Taxon-Specific EnzymeCrystallins |   |  |
| d                                | Most birds, reptiles                      | Argininosuccinate lyase (d2)                 |
| ε                                | Crocodiles, some birds                    | Lactate dehydrogenase B                      |
| z                                | Guinea pig, degu rock cavy, camel, llam   | NADPH:quinone oxidoreductase                 |
| h                                | Elephant shrews                           | Aldehyde dehydrogenase I                     |
| l                                | Rabbits, hares                            | [Hydroxyacyl CoA dehydrogenase]              |
| m                                | Kangaroos, quoll                          | [Ornithine cyclodeaminase]                   |

|   |  |                                      |
|---|--|--------------------------------------|
| r | Frogs ( <i>Rana</i> )                                    | [NAPDH-dependent reductases]         |
| t | Lamprey, turtle; moderately abundant in most vertebrates | $\alpha$ -Enolase                    |
| S | Cephalopods  | [Glutathione <i>S</i> -transferases] |
| W | Octopus  | [ALDH]                               |
| J | Cubomedusan jellyfish                                    | ?                                    |

---

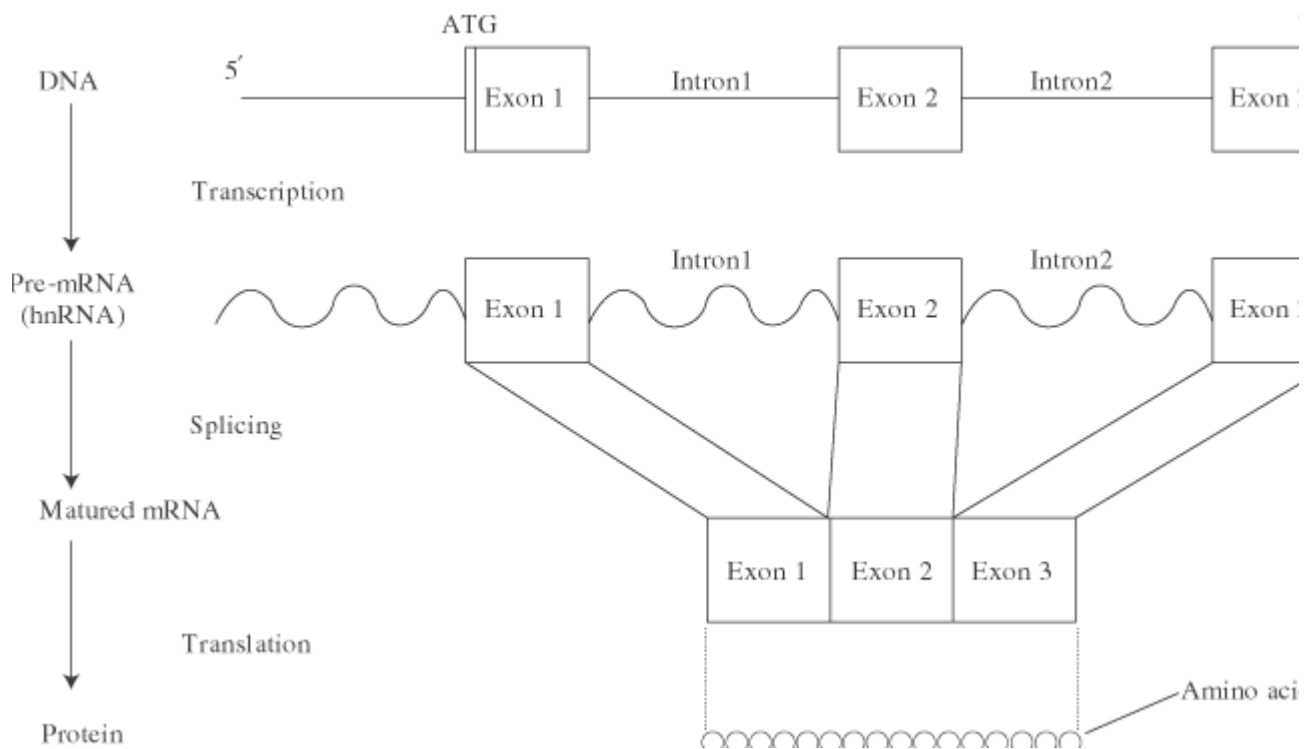
## Bibliography

1. G. Wistow (1993) Trends in Biochem. Sci. **18**, 301–306.

## Gene Splicing

Many eukaryotic **genes** are split into a number of exons by the existence of **introns**. Gilbert coined the terms “exons” and “introns” in 1978 ([1](#)). As shown in Figure [1](#), (see top of next page) the human  $\alpha$  globin gene is split into three exons by two introns. To produce the appropriate [messenger RNA](#), introns must be removed to generate a proper collection of linked exons through gene splicing; as it occurs in the mRNA, it is commonly called [RNA splicing](#). Therefore, in the case of eukaryotic genes, the essential information flow will be as follows; from DNA to mRNA precursor by [transcription](#), from precursor mRNA to mature mRNA by [RNA splicing](#), and from mature mRNA to protein by [translation](#).

**Figure 1.** Gene splicing steps in the production of mature eukaryotic mRNA as illustrated for the human  $\alpha$ -hemoglobin gene ([Gene Splicing](#)).



As for the evolutionary significance of the exon–intron structures of genes, the “exon shuffling” theory was proposed by Gilbert (1). He claimed that a variety of genes on the genome have been created by the shuffling and exchange of exons among different genes (1) (see [Exon Shuffling](#)). As introns do not possess any apparent function, it is considered that they are used as locations for genetic [recombination](#). Thus, frequent shuffling of exons may have been major forces of gene evolution in eukaryotic organisms. There has, however, been some opposition to the “exon shuffling” theory. For example, the “domain shuffling” theory contends that the unit of shuffling is not an exon but a functional **domain** of a protein that may not correspond exactly to an exon (2). In fact, it is known that the [Kringle domain](#), which is a characteristic [supersecondary structure](#) in some [serine proteinases](#), was shuffled many times, and always as a unit, even though the Kringle domain consists of three exons (3).

With only a few exceptions, on the other hand, it is well known that prokaryotic [genomes](#) do not contain introns. If an intron is considered to have existed before the **divergence** between eukaryotes and prokaryotes, the splicing mechanism should also have existed for a long time since the early stages of evolution. This is the “early-intron” theory. On the contrary, the “late-intron” theory claims that introns have been inserted as **transposon**-like elements into various genes after the divergence of eukaryotes and prokaryotes. The evolutionary significance of gene splicing depends on which theory is correct.

#### Bibliography

1. W. Gilbert (1978) *Nature* **271**, 501.
2. R. Dolittle (1993) *Sci. Am.* (Oct.) 50–56.
3. T. Gojobori and K. Ikeo (1994) *Phil. Trans. Roy. Soc. Lond. B* **344**, 411–415.

## Gene Structure

A *gene* is widely understood as the fundamental unit of genetic information, but a detailed description of a gene is not straightforward.

### 1. Intragenic Recombination and Colinearity with Polypeptide Chain

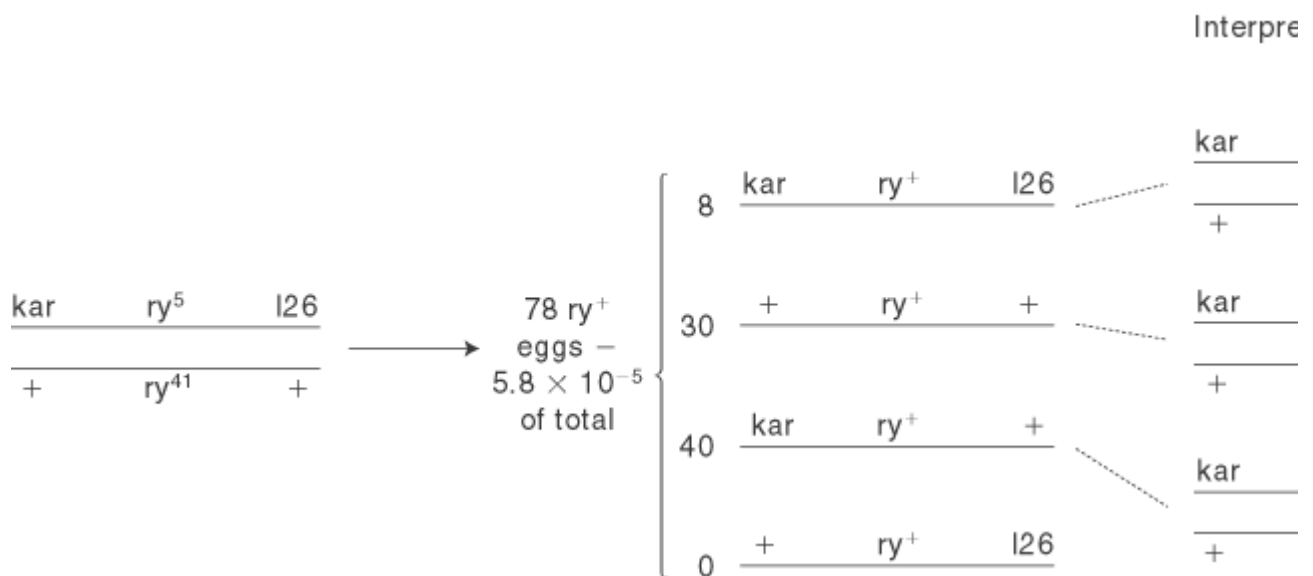
The first step in elucidating the structure of the *gene* was recognizing that it is subdivisible into linearly arranged, individually mutable sites. Until the 1950s it was assumed that genes are both units of function and indivisible units of genetic transmission. The functional criterion for different [mutants](#) affected in the same gene, that is **allelic**, was that, when brought together in a **heterozygote** or **heterokaryon**, they could not make good each other's deficiencies to produce a nonmutant **phenotype**. This criterion is still a good general guide, although, as detailed under [Interallelic Complementation](#), there are circumstances in which it fails. The indivisibility criterion meant that two different mutant alleles of the same gene should not be able to **recombine** to yield a **wild-type** allele. This was sustainable only as long as only moderate numbers of meiotic segregants from heteroallelic **diploids** were screened.

In the fruit fly *Drosophila melanogaster*, M. M Green, working with the *lozenge* gene, and E. B. Lewis, with different genes, including *white*, demonstrated rare crossing over between mutants previously thought to be allelic and dubbed these recombining mutants *pseudoallelic* (1). Then Pontecorvo and his colleagues found that the sexually reproducing fungus *Aspergillus nidulans* produces wild-type recombinants from crosses between several pairs of alleles, with a frequency on the order of 1 in 10,000, and challenged the idea of pseudoallelism, thinking it likely that mutual recombining is a typical feature of mutants that are allelic in the functional sense. This view was strengthened by Benzer's analysis of hundreds of mutations within the *rII* gene of **bacteriophage** T4, in which he showed that very nearly all pairwise combinations of noncomplementing mutants would yield wild-type virus at some frequency from mixed infection of bacterial cells (2). It was rare for different mutations within the same functional gene to fall at exactly the same site and, when at different sites, they would recombine with one another. It was shown that the same principle holds for the bacteria, *Escherichia coli* and *Salmonella typhimurium*, where the method of analysis was mainly **transduction**, and in fungi and *Drosophila*, where recombination within genes occurred during **meiosis**. In the **yeast** *Saccharomyces cerevisiae* the frequency of intragenic meiotic recombination is unusually high, up to several per cent of meiotic products.

When recombination between allelic mutants was discovered, it became possible to map sites of mutation within a gene in a linear sequence. There are two general ways of doing this. The first uses flanking markers, gene differences that have visible effects (markers) closely placed on each side of the gene under analysis. If intragene recombination is due to classical **crossing over**, the flanking markers occur mainly in one or other of the two "crossed-over" combinations among wild-type interallelic recombination products, depending on which way the mutational sites are placed with respect to the flanking markers. Figure 1 shows an example from the analysis of the *Drosophila rosy* gene by Chovnick *et al.* (3).

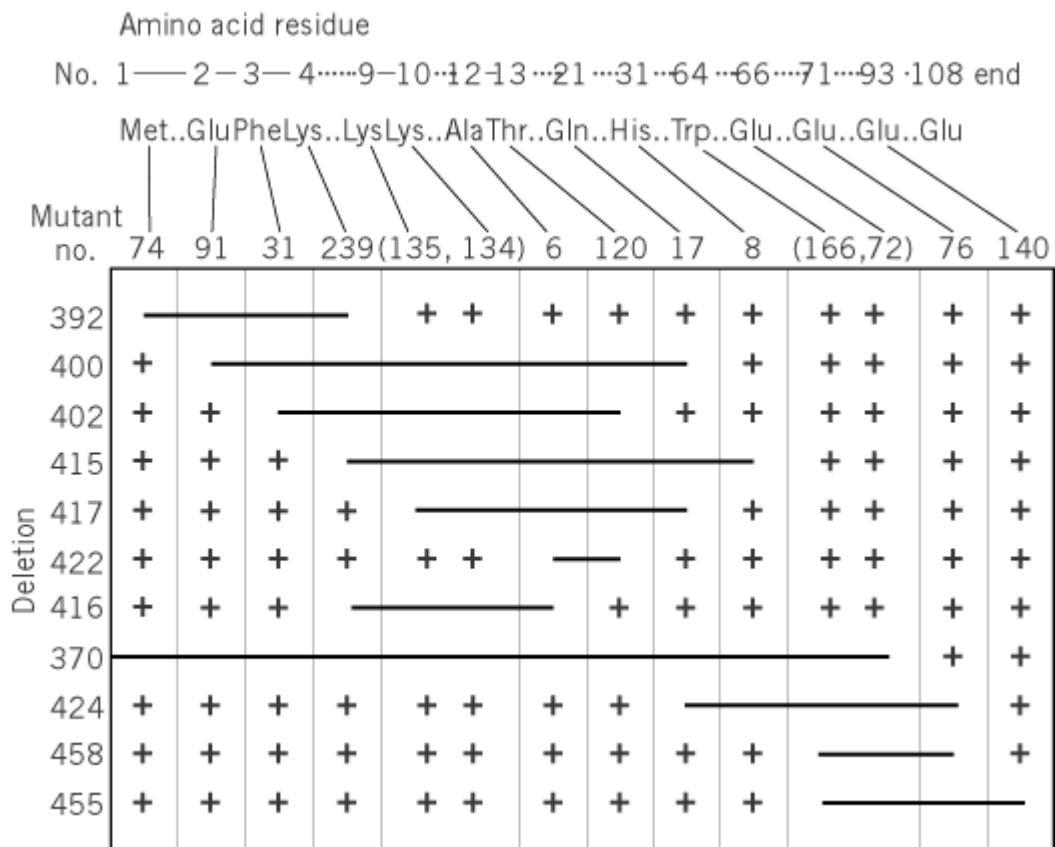
**Figure 1.** Recombination within the *rosy* (*ry*) gene of *Drosophila*, which encodes xanthine dehydrogenase (XDH). Larvae recombinant eggs from *ry5/ry41* (XDH<sup>-</sup>) females were automatically selected on purine-containing food, which kills XI parents contributed only mutant *ry*. The constitutions of the eggs with respect to mutations in closely placed flanking genes further test-crosses. The *ry*<sup>+</sup> eggs not recombined for the flanking markers are presumed to be due to conversion of one or without crossing over. The high degree of association of 5–41 recombination with one flanking marker crossover class is 126. Data from Ref. 3, whose analysis includes numerous *ry* mutant sites not shown here.





The second method, which gives clear results even when the recombination within the gene is due mainly to **gene conversion** without crossing over, relies on the principle of overlapping deletions. In any very extensive collection of allelic mutations, it is likely that some are due to deletion of segments of the gene, each of which overlaps several mutational sites and partly overlaps other deletions. No point mutation can recombine to form recombinants with a deletion that removes the corresponding nonmutant site, and neither do different deletions yield wild-type recombinants if they overlap. By crossing a suitable set of deletion mutants with each other and scoring the progeny for the presence versus absence of wild-type recombinants, one establishes the overlaps between the deletions, and then, by crossing the whole set to point mutations, localizes the mutational sites to one or other of the segments defined by the deletion overlaps. The method was first used by Benzer in his analysis of the T4 *rII* genes. An example from yeast, the *cyc-1* gene which encodes the protein of the respiratory pigment [cytochrome c](#), is explained in Fig. 2 (4).

**Figure 2.** Establishment of the sequence of mutant sites within the *Saccharomyces cerevisiae* *CYC1* gene that encodes cytochrome c. A set of overlapping deletion mutants were each tested for recombination with point mutants: + or – indicates whether or not wild-type recombinants were formed. The absence of recombinants means that the site of the point mutation falls within the segment deleted in the other parent. The deletion overlaps define the sequence of the mutant sites. The amino acid codon affected by each point mutation was determined biochemically. The sequence of codons corresponds to the sequence of mutant sites. From Ref. 4 by permission.



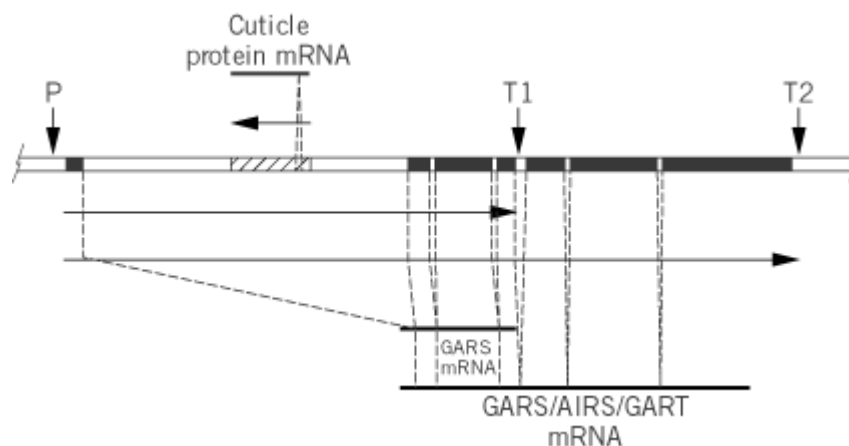
From this and many other examples, the gene emerged as a very finely subdivisible unit. Recognizing that the gene material is **DNA**, the concept, which still holds, is that every base pair of the DNA sequence is a potential site for mutation, separable from its neighboring nucleotides at some low frequency by recombination. Taken together with the idea that each gene is responsible for a single enzymic protein, which had emerged from the early *Neurospora* work, it was an obvious guess that the linear sequence of mutable and recombinable sites within the gene corresponds to the linear amino acid sequence of the polypeptide chain of the protein product.

*Drosophila rosy (ry)* and *Saccharomyces cyc-1* are particularly appropriate examples because their mutations affect well-characterized proteins, namely, the enzyme xanthine dehydrogenase in the first case and the protein of the respiratory pigment cytochrome *c* in the second. In both cases point mutations cause mostly single amino acid replacements in the polypeptide chain, and the sequence of the replacements along the polypeptide chain corresponds to the sequence of mutational sites along the gene map. The comparison is particularly extensive for yeast *cyc1* (Fig. 2). Numerous other examples could be given, and the principle, called *colinearity*, is so well established now in molecular biology that further substantiation through genetic mapping is hardly necessary. Genes, or at least those genes responsible for protein structure, are linear codes for the amino acid sequences of [polypeptide chains](#) with exceptions that are mentioned later. Parallel with examples of colinearity came the molecular identification of the gene as DNA, the demonstration of [transcription](#) of genic DNA into RNA, and the working out of the genetic **code** through which [messenger RNA](#) (mRNA) is **translated** into polypeptide sequence.

The original “one gene - one enzyme” concept was incomplete in several ways. The immediate products of genes are RNA molecules, not all of which are messengers for protein synthesis (e.g., **ribosomal**, **small nuclear RNA**, and [transfer RNA](#)), and not all proteins are enzymes. Among the enzymes, many are multifunctional, and one function can often be lost by mutation without eliminating other functions of the same polypeptide chain. As a result, analysis by mutation,

[complementation](#), and recombination may subdivide a gene into several functionally distinct domains, which may at first be confused with separate genes. A multifunctional, multidomain gene is sometimes called a *cluster gene*, to distinguish it from a [gene cluster](#). Most of the known examples are in fungi, but they occur in *Drosophila* (eg, Fig. 3) and also in mammals. This sort of gene complexity is more extensively discussed under [Interallelic Complementation](#).

**Figure 3.** The *Gart* gene of *Drosophila*, which exhibits several kinds of gene complexity: (1) an [intron-exon](#) structure (common to most *Drosophila* genes; introns are shown as open and exons as black segments); (2) the full-length mRNA encodes a trifunctional polypeptide that has three different enzymatic activities, abbreviated to GARS, AIRS, and GART, that catalyze successive steps in purine biosynthesis; (3) it has two polyadenylation/transcript termination sites, T1 and T2, one within the fourth intron resulting in a truncated mRNA that encodes only the functionally independent GARS domain; (4) it has a “nested” gene within its longest intron that encodes a cuticle protein and itself with an intron. From Ref. 11 by permission of the publisher, after Ref. 5.



## 2. Transcribed and Nontranscribed DNA Strands and the Minimal Gene

In writing about gene structure, it is useful to use the terms “*upstream*” and “*downstream*” to indicate orientation with respect to the direction of transcription of DNA into RNA. Only one strand of the DNA duplex is usually transcribed, and it is important to remember that the transcribed strand of a protein-encoding gene is not the *coding strand*, in the sense that the codon sequence is as seen in the messenger RNA, but rather the opposite-polarity complement of that sequence. The transcribed strand runs upstream-to-downstream in the chemical direction 3' to 5', whereas the mRNA and its amino acid codons run 5' to 3'. To deduce the amino acid sequence from the transcribed strand using the standard codon tables (see [Genetic Code](#)), one must mentally convert it into its complement: A to U, T to A, G to C, and C to G. The amino acid sequence can be read directly from the nontranscribed strand, which is therefore often called the [coding strand](#). It is best to distinguish the strands as transcribed or nontranscribed, bearing in mind that the sequence of the latter bears the immediately recognizable codons.

Once defined as a unit of transcription into RNA, the gene may be a relatively compact structure, at least when measured against more than  $10^9$  DNA base pairs in a typical mammalian [genome](#) or the even greater DNA contents ([C-values](#)) of some amphibia and higher plants. We can enumerate the essential components of a protein-encoding gene as follows. There must be a **promoter** segment to which **RNA polymerase** binds to initiate transcription a little way upstream of the transcription startpoint (with the exception that RNA polymerase III works with downstream promoters), and a transcription termination and [polyadenylation](#) signal. Within the transcribed region there are sequences corresponding to the different sections of the messenger transcript: an untranslated **leader sequence** (of various lengths in different genes and sometimes important for controlling translation),

a [ribosome](#) binding site and [initiation codon](#), an **open reading frame**, and a **termination codon**. The yeast *cyc1* gene (Fig. 2) is a good example of a simple gene.

### 3. Expansion of the Gene Beyond the Coding Sequence

To be transcribed into mRNA for polypeptide chains or other kinds of RNA molecule, genes need be no more than a few thousand base pairs (kilobase pairs or kbp) long. How, then, is one to account for the apparently huge surplus of DNA in the genomes of many higher organisms? Humans, for example, have about 30 kb of DNA per gene if the current estimate of about 100,000 genes is correct. Part of the answer lies in the host of repetitive and arguably functionless sequences that have become established. But probably equally important quantitatively are the intervening sequences (**introns**) that subdivide the coding sequences into fragments (exons) that are **spliced** together, and the introns spliced out, after the gene has been transcribed. The sizes and numbers of introns vary widely among different organisms. They are hardly present at all in bacteria in only a few genes and then seldom more than a hundred bases long in yeast, more numerous but still short in filamentous fungi, sometimes much longer in *Drosophila*, and often both very numerous and extremely long in mammals. The typical mammalian gene consists of rather short exons, each on the order of only hundreds or even tens of base pairs, separated by introns extending up to tens of kilobases long. Thus a gene, whose function is to encode a polypeptide chain of just a few hundred amino acid residues, may be spread over a tract of some 100,000 base pairs.

In addition to introns, other sequences, which may fall some distance outside the units of transcription, have some claim to be considered parts of the gene and to expand its domain. Although [enhancers](#) may occur within introns, they are probably more often outside the transcribed gene sequence, often upstream but sometimes downstream. The same applies to **silencers**. Insofar as enhancer sequences are [cis-acting](#), which they usually are, they may be considered as falling within the gene boundaries. If so, this may somewhat complicate the idea of the gene as a discrete functional unit if the same enhancer services two or more different transcriptional units. Thus, a single **locus control region**, an approximately 20-kbp sequence, not usually called an enhancer but at least some elements of which function as such, acts on all of the genes of the human [b-globin](#) gene cluster. Genes are also functionally linked through [chromatin](#) structure, changes in which repress or release the transcription of blocks of genes (see also [Epigenetics](#), [Position Effect](#)).

### 4. Programmed Gene Restructuring

Although most genes in most organisms have constant DNA sequences, the same in the cells where they are transcribed as in the germ cells through which they are transmitted, there are examples of genes in a wide range of organisms that are restructured in the process of cellular differentiation. In unicellular organisms, bacteria, yeasts, and protozoa, where any cell can found a new population, gene restructuring is always reversible within the cell. Thus the switching of **flagellar** antigen in the bacterium *Salmonella typhimurium* is brought about by the inversion of a DNA segment with the effect of “switching off” one **antigen**-encoding gene by separating it from its promoter and thereby switching on another antigen gene of which the first antigen is a [repressor](#) (7). The inversion is reversible. In *Saccharomyces* yeast, the well-known **mating type** switch is due to replacing a segment of DNA at a transcriptionally activating site with a segment copied from another locus, where it had been transcriptionally silent. The potential for both mating types is held at silent “cassette” loci at all times, but only one mating type is expressed from the segment present at the activating (or, more correctly, nonsilencing) locus. Similar mating-type switching occurs in the fission yeast, *Schizosaccharomyces pombe*. A system not dissimilar in principle but with more options that again involves the transfer of gene sequences from silent to expressed loci, operates to switch the major surface antigen of the pathogenic protozoan *Trypanosoma brucei* (8).

In all the cases just mentioned, the potential for switching back to the status quo ante is retained within the cell [nucleus](#). The situation is very different in ciliated protozoa, such as *Paramecium*, *Tetrahymena*, and *Oxytricha*, in which two kinds of nuclei are in the binucleate cells: a virtually

“silent” **micronucleus**, which contains the basic genetic material and has the potential for differentiation in different directions (e.g., to different mating types or different surface antigens) and an active **macronucleus**, in which a selection of genes is amplified and restructured to support a particular cell type. At meiosis, which generates haploid nuclei for sexual cross-fertilization or for a kind of self-fertilization called autogamy, the macronuclei degenerate and disappear and are replaced by division and restructuring of the micronuclei, with the opportunity for switching cell type. Some of the gene rearrangement seems quite bizarre. For example, in *Oxytricha trifallax*, a gene that encodes an [actin](#) protein whose exons are labeled 1 to 10, upstream-to-downstream in the macronucleus, was reshuffled from the order 3, 4, 6, 5, 7, 9, 10, 2, 1, 8 in the micronucleus ([9](#)).

Gene restructuring in mammals is the exception rather than the rule, but is centrally important in the immune system. The functional genes for the virtually infinitely large number of different [immunoglobulins](#) and [T-cell receptor](#) proteins are pieced together in a large number of possible permutations from DNA segments that are separated by hundreds of kilobases in undifferentiated [stem cells](#). Note that the nucleic acid splicing involved here occurs at the level of DNA, not RNA as in splicing-out of introns. This topic is dealt with in more detail under [Complex Loci](#), [Immunoglobulins](#), and [T-cell receptor](#).

A different category of incomplete genes are those made functional not by rearrangement at the DNA level but by splicing their RNA transcripts to leader sequences provided by what might be called supplementary genes elsewhere in the genome. Many or most messenger RNAs in the **Trypanosomes** acquire leader sequences and cap sites in this way ([10](#)). Something similar occurs in [nematodes](#), and such systems may be widespread in the less well explored “lower” eukaryotes.

#### Bibliography

1. E. B. Lewis (1952) Proc. Natl. Acad. Sci. USA **38**, 953–961.
2. S. Benzer (1959) Proc. Natl. Acad. Sci. USA **45**, 1607–1620.
3. A. Chovnick, A. Ballantyne, and D. G. Holm (1971) Genetics **69**, 179–209.
4. F. Sherman et al. (1975) Genetics **81**, 51–73.
5. S. Henikoff, M. A. Keene, K. Fechtel, and J. W. Fristrom (1986) Cell **44**, 33–42.
6. S. N. Krishnan, E. Frei, A. P. Schalet, and R. J. Wyman (1995) Proc. Natl. Acad. Sci. USA **92**, 2021–2025.
7. J. Zieg et al. (1980) Science **196**, 170–172.
8. J. H. J. Hoeijmakers et al. (1980) Nature **284**, 78–80.
9. M. Dubois and D. M. Prescott (1995) Proc. Natl. Acad. Sci. USA **92**, 3888–3892.
10. S. Lucke et al. (1996) EMBO J. **15**, 4380–4391.

#### Suggestion for Further Reading

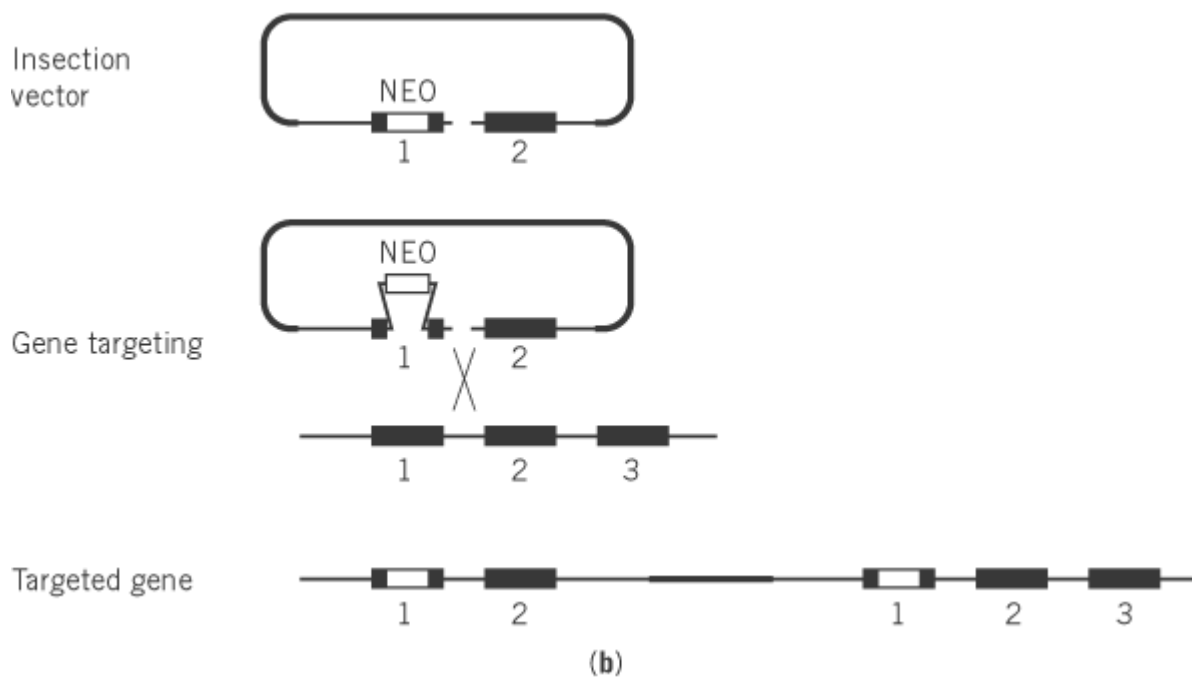
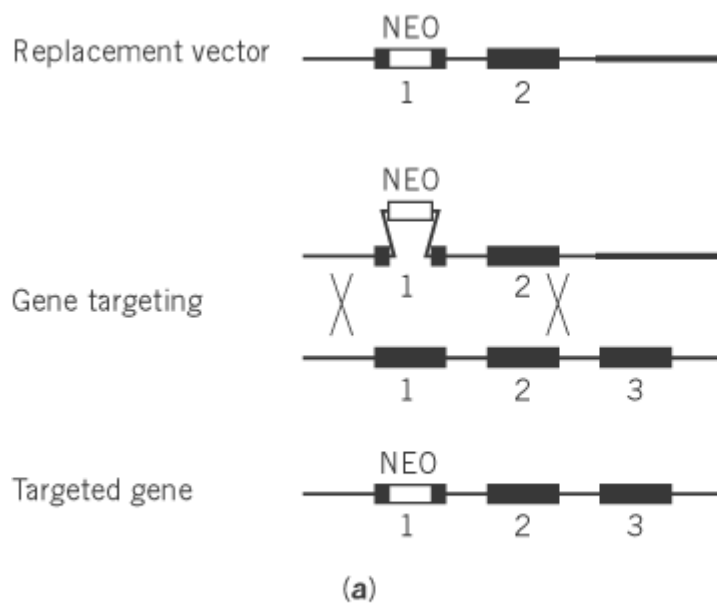
11. J. R. S. Fincham (1994) *Genetic Analysis*, Blackwell Science, Oxford.

#### Gene Targeting

Gene targeting refers to the introduction of genetic alterations into predetermined sites within the [genome](#) via homologous [recombination](#). A DNA vector containing genomic sequences surrounding the required genetic alteration is constructed by [recombinant DNA](#) techniques and introduced into the cell. Homologous recombination between the vector and genome results in insertion of the

genetic alteration at the homologous endogenous location and production of a genetically modified cell (Fig. 1).

**Figure 1.** Diagrammatic representation of the generation of genetically altered mice by gene targeting. ES cell lines derived from the mouse blastocyst are cultured *in vitro* and transfected with a gene-targeting vector. Transfected cells are selected in the appropriate growth medium and homologous recombinants are identified by molecular screening methods. Correctly targeted clones are expanded and returned to recipient mouse embryos which develop to term in pseudopregnant recipient females. Chimaeric offspring which transmit the ES cell genotype through the germline are mated to generate mice heterozygous for the targeted gene. Male and female heterozygous mice are mated to generate mice homozygous for the targeted gene.



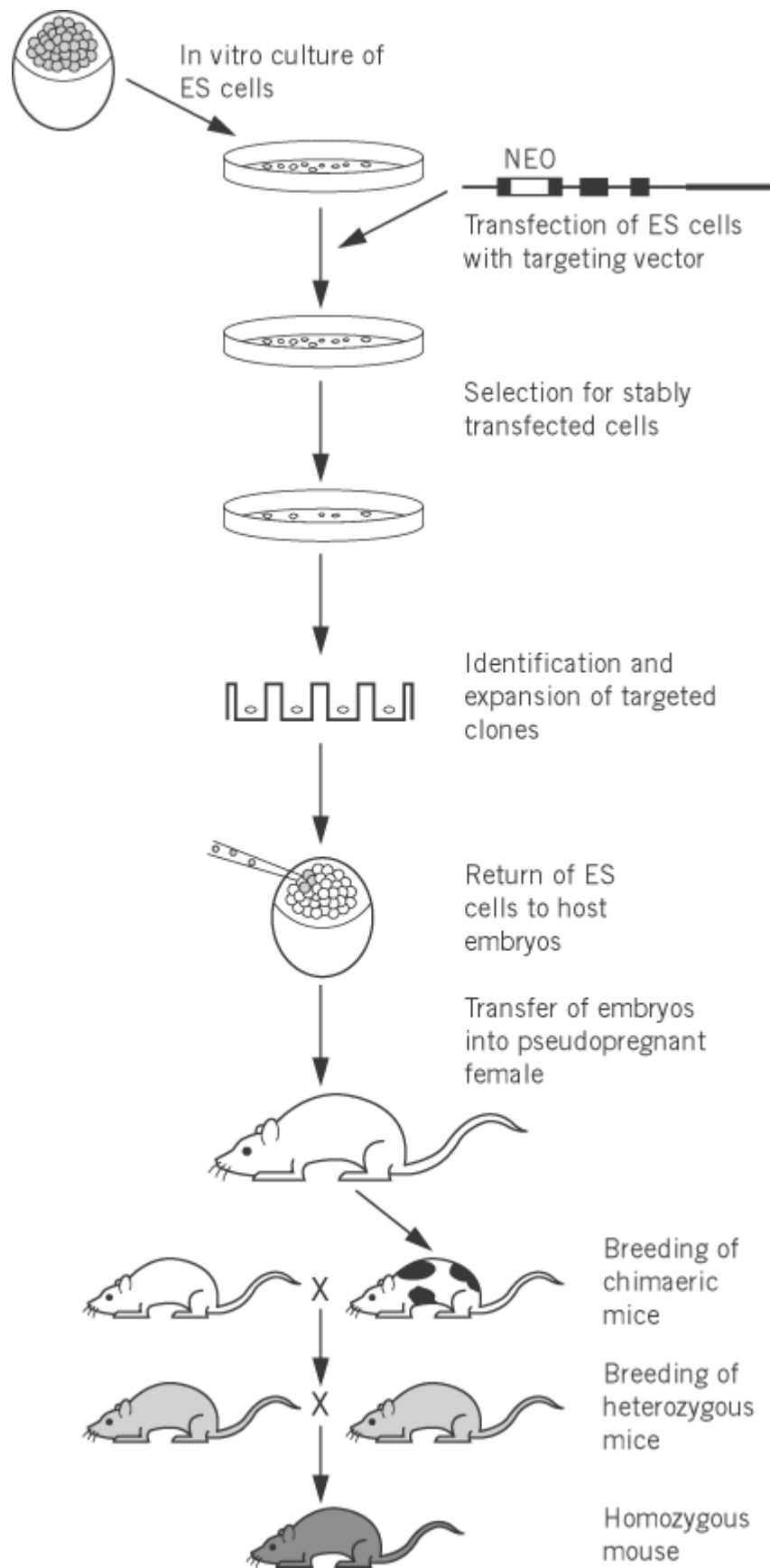
Gene targeting has been demonstrated with variable efficiencies in organisms as diverse as bacteria, yeast, and mammals, as well as in both [haploid](#) and multiploid cells. Strategies have been developed

for the introduction of diverse [mutations](#) ranging from alteration of single nucleotides to incorporation of large DNA segments. Most recent interest reflects the realization that, by coupling gene targeting and reproductive technologies, genetically altered cells can be used as vectors for the creation of genetically modified multicellular organisms. This technology has been developed furthest in the [mouse](#), and the term *gene targeting* is now most commonly applied to the generation of mutant mice. Application of gene targeting to other species has great potential for the analysis of gene function, treatment of human disease, and generation of novel agricultural products.

### 1. Generation of Null Mutant Mice via Gene Targeting in Embryonic Stem Cells

The early mouse embryo contains a pool of pluripotent cells that will give rise to all cell types of the embryo and adult. These cells can be isolated from preimplantation embryos and maintained in the undifferentiated state *in vitro* as embryonic stem (ES) cells. ES cells retain pluripotent differentiation capability even after long-term culture and can contribute differentiated progeny to all tissues, including the germline, following reintroduction into the mouse embryo (Fig. 2). Genetically altered ES cells can be used to transmit mutations constructed *in vitro* into the mouse germline, thereby enabling the creation of mice with predetermined mutations.

**Figure 2.** Diagrammatic representation of the generation of null mutations with replacement and insertion targeting vectors. **(a)** Replacement gene-targeting vector. **(b)** Insertion gene-targeting vector. Thin line, genomic DNA; thick line, plasmid DNA; numbered black boxes, exons; NEO, neomycin resistance selection cassette; X, region of homologous recombination.



The steps involved in making a mouse that carries an inactivated gene (null mutation) are shown diagrammatically in Figure 2. Standard recombinant DNA techniques are used to isolate a genomic



clone corresponding to the region to be altered and to construct a gene-targeting vector. A positive selection cassette, such as the bacterial **neomycin**-resistance gene, which provides resistance to the aminoglycoside G418, is usually positioned within the genomic clone so that it is flanked by genomic DNA sequences. The targeting vector DNA is introduced into undifferentiated ES cells grown in the presence of leukemia inhibitory factor (LIF) or a feeder cell layer. A range of standard [transfection](#) techniques allow introduction of the vector DNA into ES cells, although electroporation is most commonly used (see [Transfection](#)). In most cases the introduced DNA integrates into the genome at random sites by nonhomologous recombination, but in a few cases (in the order of 0.5% to 10%) the introduced DNA is integrated into the genome at the homologous site. Stably transfected cells are selected for expression of the selection cassette. Individual clones are expanded and screened for homologous recombination by [Southern blot](#) or [polymerase chain reaction \(PCR\)](#). Correctly targeted clones are expanded *in vitro* and reintroduced into host embryos by injection into the blastocoelic cavity of 3.5 days post coitum (d.p.c.) embryos, or by co-culture with 2.5 d.p.c. morulae. Embryos are returned to pseudopregnant recipient female mice and develop into **chimaeric** mice. Usually the coat colors of the host embryo and the mice from which the ES cells were obtained differ, so that chimaeric mice can be identified by their chimaeric coat coloration. Chimaeric mice in which targeted ES cells have contributed to the germline give rise to progeny that are heterozygous for the targeted gene in all cells of the body. To produce mice that are homozygous for the targeted gene, heterozygous male and female mice are mated.

The most common application of gene targeting technology in mice has been the creation of null mutants, usually by positioning the selectable marker so that it disrupts or replaces the gene of interest. Two types of vectors are used for the creation of null mutations: replacement and insertion vectors ([1](#)). Replacement vectors are linearized at the extremities or outside the homologous sequence, and the vectors remain co-linear with the target sequence (Fig. [2a](#)). To produce null mutations, a critical section of the coding region is replaced or disrupted with the selection cassette. Insertion vectors are linearised within the region of homology, and homologous recombination of these vectors leads to duplication of the genomic sequence (Fig. [2b](#)). To produce null mutations with insertion vectors, the coding sequence is disrupted with vector DNA that also carries a selection cassette disrupting or replacing an essential region of the coding sequence. Mice carrying specific mutations have proved extremely valuable for basic biological investigations ([2-4](#)).

## 2. Improvements to Gene Targeting Efficiency

### 2.1. Parameters

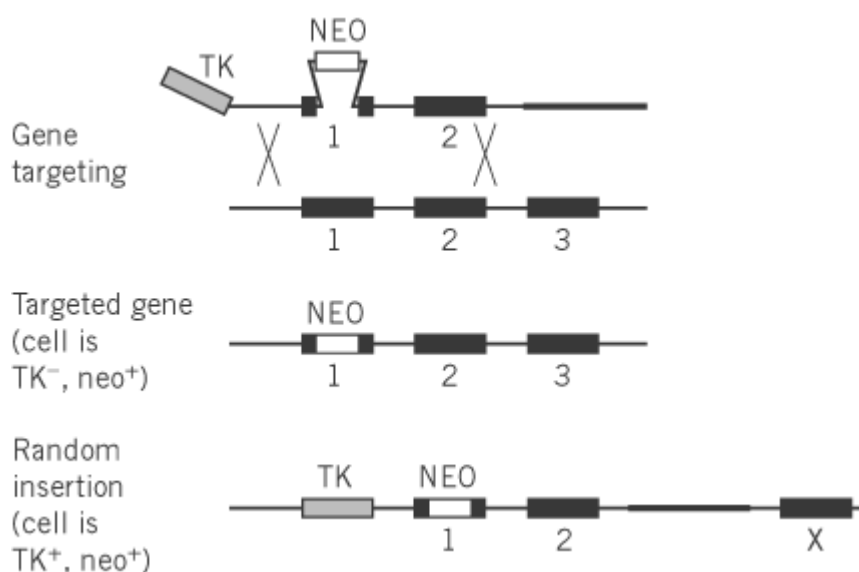
The mechanisms by which homologous recombination of replacement and insertion vectors occurs in mammalian cells are poorly understood. Genes can be targeted efficiently in ES cells regardless of whether they are expressed ([5](#)). The efficiency of homologous recombination can be increased by increasing the total amount of homology contained in a targeting vector ([1, 6, 7](#)). The relationship between the length of homology and the gene targeting efficiency is exponential between 2 and 10 kb, whereas increases in homology beyond 14 kb have little effect. The efficiency of homologous recombination can also be increased significantly by using genomic DNA for vector construction that is isogenic with the strain of mice from which the ES cells were derived ([8, 9](#)). The frequency of homologous recombination appears to be unaffected by the length of nonhomologous DNA introduced into the homologous site, at least up to 12 kb ([10](#)), and separate regions of extensive nonhomology can be introduced efficiently into the genome on a single replacement vector ([11](#)).

### 2.2. Vector Design

For loci that are targeted at low frequencies, identification of occasional homologous recombinants among a high background of stably transfected cells can be labor-intensive. Vector design can be modified to enrich for homologous recombinants. The positive-negative selection system ([12](#)) is based on the finding that nonhomologous regions at the extremities of replacement vectors are integrated during nonhomologous, but not homologous, recombination. A negative selection cassette, such as herpes simplex virus thymidine kinase (HSV-*tk*), which confers sensitivity to gangcyclovir or FIAU, is positioned outside the homologous sequence on a replacement vector.

Selection for vector integration, but against HSV-*tk* expression, enriches for homologous integrants (Fig. 3). When isogenic DNA vectors are used, enrichments in the range of 5- to 10-fold can be achieved. However, the findings that DNA ends can stimulate homologous recombination (13, 14), and that vector ends might initiate homologous recombination of replacement targeting vectors (11), suggest that the use of vectors containing homologous sequences at the extremities could also be advantageous. Selection cassettes that lack a **promoter** or a **polyadenylation** sequence (15) have also been used to improve the efficiency of gene targeting, since nonhomologous recombination of such vectors generally results in poor expression of the selection cassette. When recombined at the homologous site, the selection cassette is expressed efficiently due to appropriate positioning within the endogenous gene.

**Figure 3.** Diagrammatic representation of positive-negative selection gene targeting. Thin line, genomic DNA; thick line, plasmid DNA; black boxes, exons; NEO, neomycin resistance selection cassette; TK, thymidine kinase selection cassette; X, region of homologous recombination.



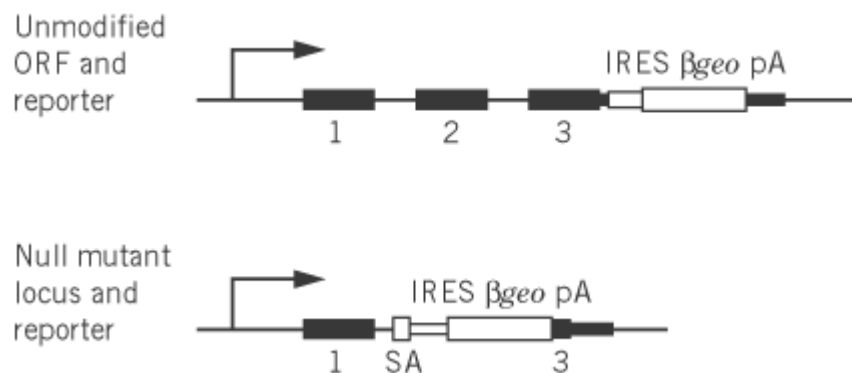
### 3. Applications of Gene Targeting

#### 3.1. Mapping Gene Expression

Integration of reporter genes into the genome so that they are expressed under the transcriptional control of the endogenous gene provides a means to visualise the normal cellular sites of gene [transcription](#) *in vivo*. Expression of reporter genes such as **b-galactosidase** (*lacZ*) can be resolved at the cellular level by histochemical techniques. Chimaeric reporter/selection cassettes, such as the b-galactosidase-neomycin fusion *bgeo* (16) that has the enzymatic activity of b-galactosidase and provides resistance to G418, can be used simultaneously as a reporter for gene expression and as a selection cassette for the identification of stably transfected ES cells. This eliminates potential transcriptional interference from promoters directing expression of the selection cassette in the targeting vector. Viral [internal ribosome-entry site](#) (IRES) sequences, which enable multiple coding regions to be translated from a single transcript in mammalian cells, can be used to couple expression of a reporter gene integrated into the 3' untranslated region to expression of an unmutated open reading frame ((17), Fig. 4).

**Figure 4.** Diagrammatic representation of gene targeting with IRES reporter vectors to map the cellular sites of gene

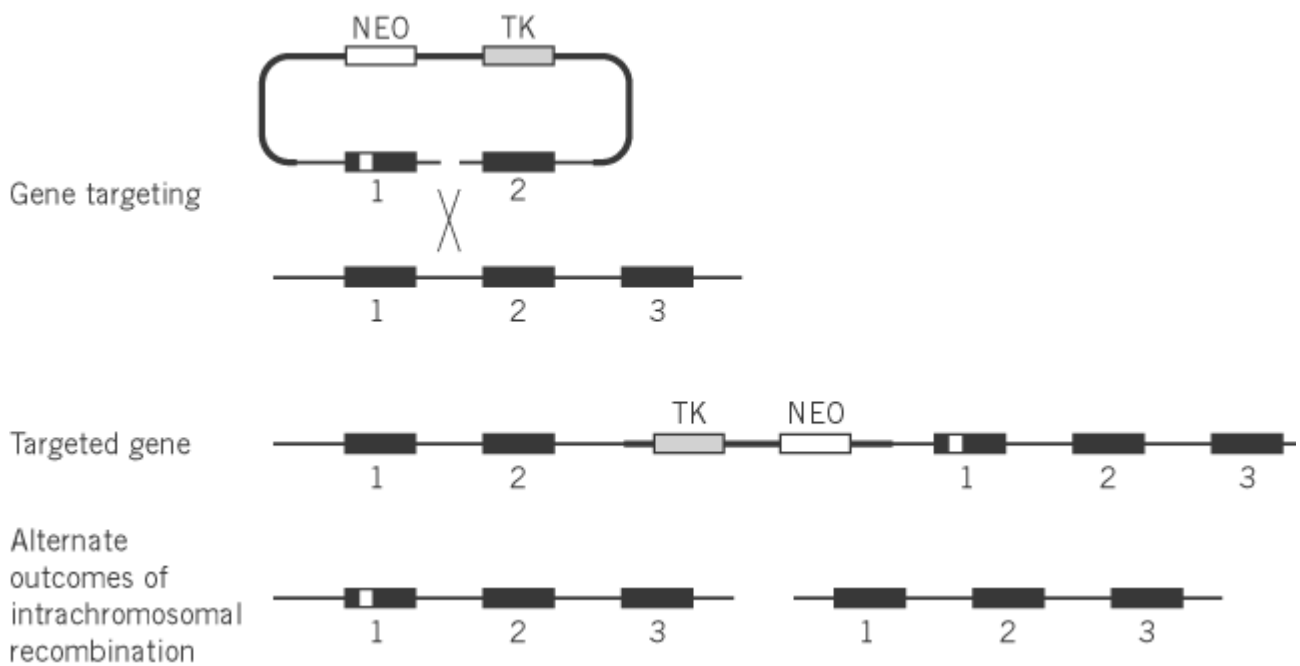
expression. Thin line, genomic DNA; thick line, untranslated region, numbered black boxes, exons; IRES, internal ribosome entry site; *bgeo* pA, b-galactosidase/neomycin resistance cassette with polyadenylation sequence; SA, splice acceptor; arrow, endogenous expression signals. Entry of ribosomes at the IRES in dicistronic mRNA results in translation of *bgeo* protein controlled by expression signals in the endogenous gene; X, region of homologous recombination.



### 3.2. Sophisticated Genetic Alterations

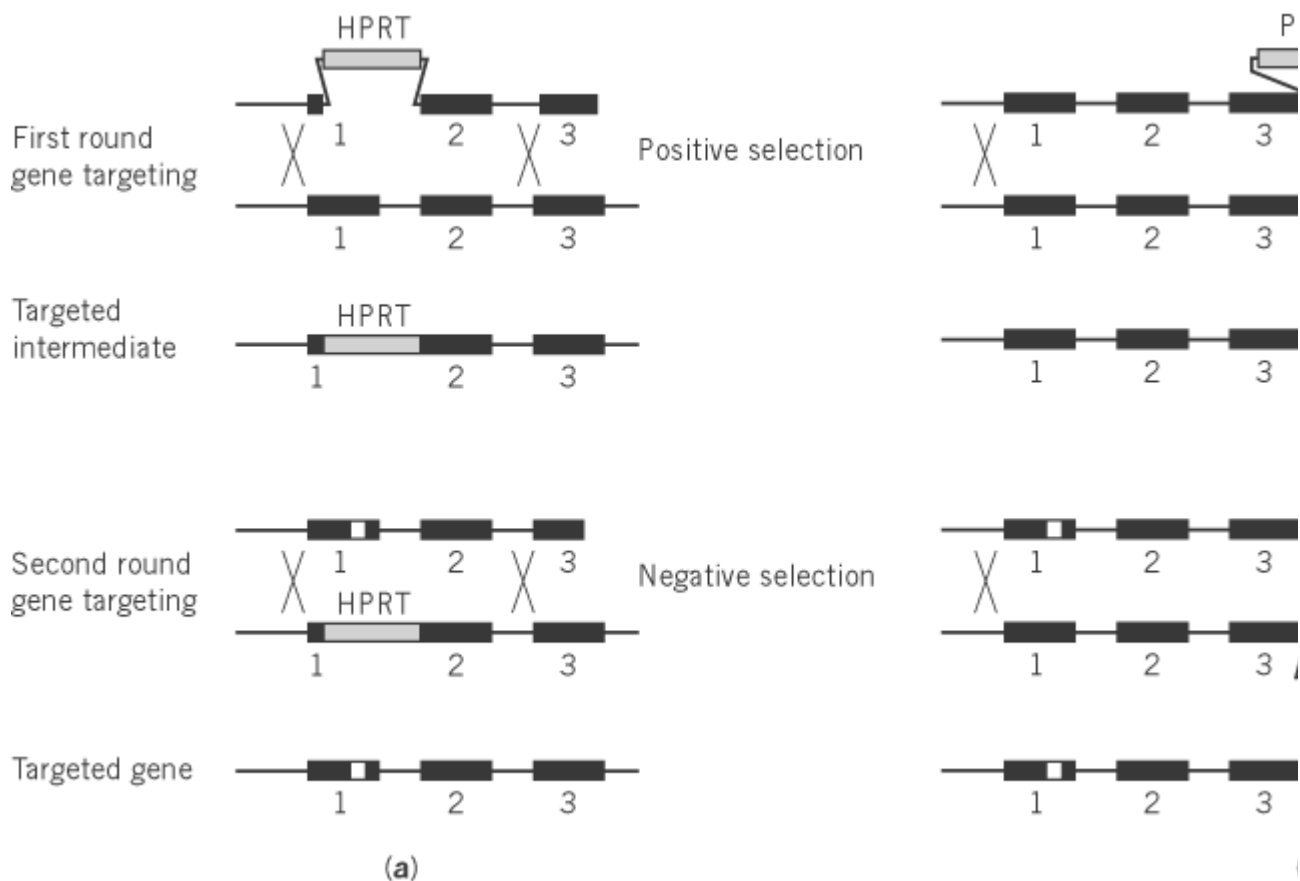
Refined gene targeting systems can be used for the introduction of subtle alterations, including single nucleotide changes, into the mammalian genome. Many of these are based on positive–negative selection cassettes such as HPRT which can be selected in hypoxanthine phosphoribosyltransferase-deficient ES cells positively in HAT medium and negatively in 6-thioguanine. In the hit-and-run system, also known as the in–out system ( (18, 19); Fig. 5), the endogenous locus is modified by homologous recombination with an insertion vector that carries the desired alteration and a positive–negative selection cassette. This leads to duplication of the homologous site and is positively selected. Intrachromosomal recombination results in cells that lose the duplication and selection cassettes and therefore survive negative selection. A proportion of these cells maintain the desired alteration.

**Figure 5.** Diagrammatic representation of the hit-and-run gene targeting approach used to generate subtle genetic modifications. Two alternate outcomes of intrachromosomal recombination are illustrated. Thin line, genomic DNA; thick line, plasmid DNA; numbered black boxes, exons; NEO, neomycin resistance selection cassette; TK, thymidine kinase selection cassette; open box, subtle mutation.



In the targeting system first suggested by Reid et al. (20), and known as double replacement (21, 22) or tag-and-exchange (23), two successive gene targeting steps are used to introduce subtle alterations into the genome (Fig. 6a). A positive–negative selection (tag) cassette is introduced into the genomic locus by gene targeting. The tag gene is replaced by the alteration of interest in a second round of gene targeting, in which direct selection for loss of the tag gene prevents survival of nonhomologous recombinants. The efficiency of second-round targeting can be compromised by a high background of resistant lines that result from physical loss of the tag cassette in the absence of homologous recombination. This most commonly results from **gene conversion** in which the tagged gene is restored to wild type by nonreciprocal transfer of genetic information from the wild-type sequence (11, 21, 22). In a modified version of tag-and-exchange targeting, termed stable tag-exchange (11), an additional positive selection cassette positioned outside the coding region is introduced on the first-round targeting vector. This cassette is not replaced during the second round of targeting and is used to select against reversion of the tagged gene to wild type. This provides a significant improvement in the efficiency of two-round gene targeting approaches. Stable tag-exchange gene targeting is therefore advantageous when investigations require the introduction of multiple, independent subtle alterations into a single locus. In cases where the presence of an exogenous selection cassette within the altered sequence may be deleterious, it can be removed by site-specific recombination using systems developed for conditional gene targeting (see text below).

**Figure 6.** Diagrammatic representation of two-step gene targeting approaches used to generate subtle genetic modifications. (a) Plug-and-socket gene targeting. Thin line, genomic DNA; numbered black boxes, HPRT, HPRT selection cassette; o selection cassette; DSocket or SocketD, alternative nonfunctional selection cassettes; Socket, functional socket selection recombination between DSocket and SocketD; X, region of homologous recombination.



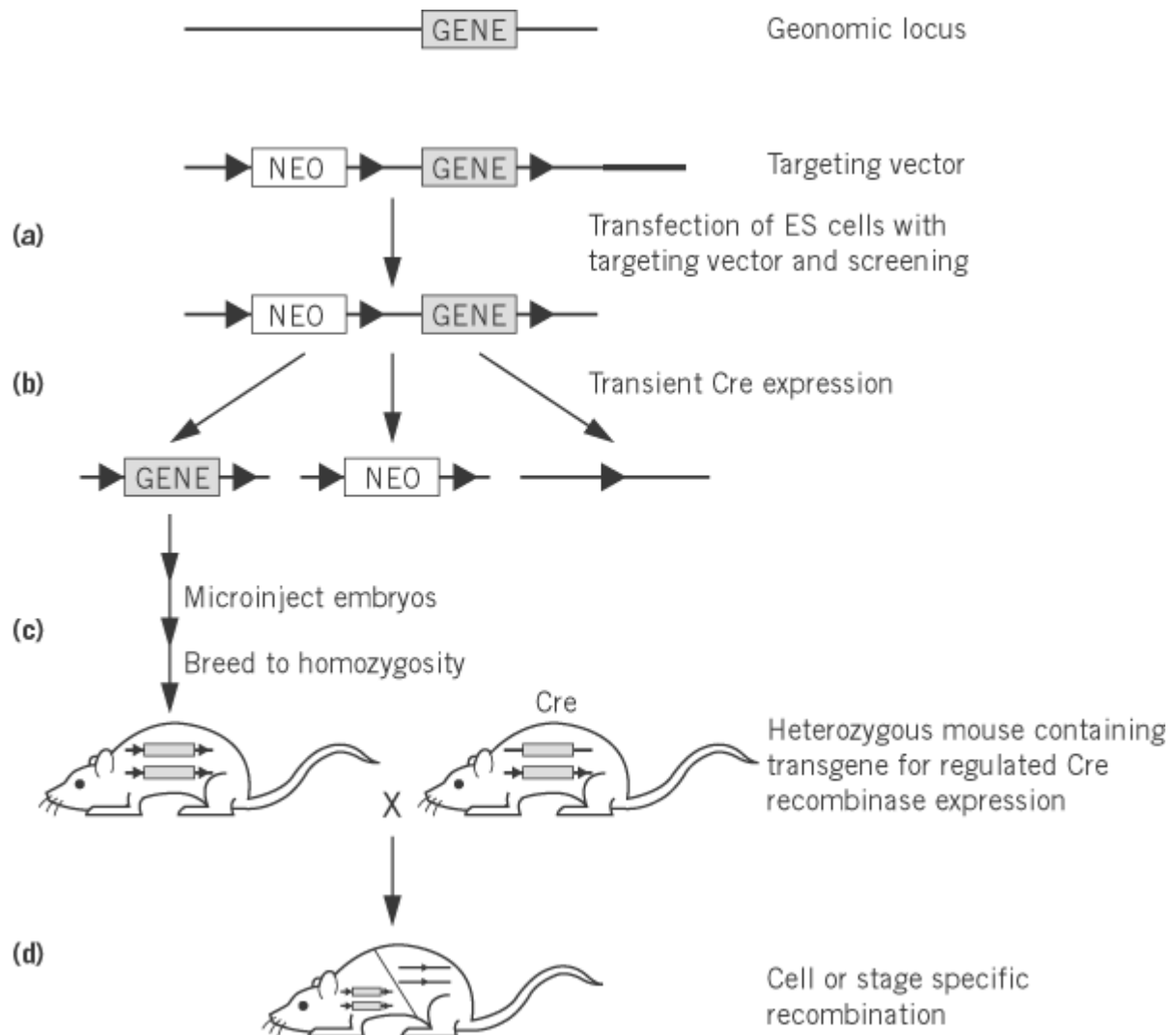
The two-step “plug and socket” system ( (24); Fig. 6b) allows use of positive selection in both rounds of gene targeting. This potentially overcomes the background of cells surviving negative second-round selection due to loss of the negative-selection cassette in the absence of homologous recombination. The locus is initially targeted with a vector that contains a functional positive “plug” selection cassette and a second nonfunctional positive “socket” selection cassette positioned outside, but close to, the locus. The vector used for second-round targeting contains the desired alteration, along with an alternative nonfunctional version of the socket selection cassette. Homologous recombination creates a functional “socket” selection cassette that is selected, and it results in replacement of the plug selection cassette with the desired alteration.

### 3.3. Conditional Gene Targeting

Null mutations produced by conventional homologous recombination have provided valuable information about the biological significance of a large number of molecules. In cases where null mutations are lethal to the embryo, however, the role of the gene product cannot be studied at subsequent stages of development. Methods for the creation of **conditional mutations**, which allow control over the site and/or time of mutation *in vivo*, are based on site-specific recombination catalyzed by a yeast or bacteriophage DNA [recombinase](#) (Fig. 7). The bacterial Cre recombinase catalyses recombination between two specific 34-bp sequences called *loxP* sites. The locus of interest is modified by gene targeting such that *loxP* sites in the same orientation are placed on either side of the selection cassette and at the end of the gene to be mutated. The ES cells are then transfected transiently with a Cre expression cassette, resulting in recombination between *loxP* sites and deletion of the intervening sequences (see [Introns](#), [Exons](#)). Cells in which the selection cassette alone has been deleted are identified and used to generate mice that carry a functional gene flanked by *loxP* sites. Homozygous offspring of these mice are crossed with mice that are heterozygous for the mutated gene and that also carry the Cre recombinase gene driven by a tissue-, cell-, or stage-specific promoter. Recombination and deletion/mutation of the *loxP* flanked gene is restricted to the

tissues or stages of development where/when the Cre recombinase is expressed and allows direct investigation of the function of the gene at these positions or stages.

**Figure 7.** Diagrammatic representation of conditional gene targeting by Cre-*lox* recombination. **(a)** A targeting construct containing the gene of interest and a neomycin resistance (NEO) gene flanked by 3 *loxP* sites (black triangles) is introduced into the ES cell genome by homologous recombination. **(b)** Transient expression of Cre recombinase in the ES cells results in loss of the gene of interest, the NEO gene or both. **(c)** Cells containing the gene of interest flanked by *loxP* sites are used to generate chimaeric mice which are bred to homozygosity. **(d)** Homozygous offspring are crossed with a heterozygous mutant which carries a Cre gene driven by a cell- or stage-specific promoter, resulting in offspring in which the *loxP* flanked gene is deleted in the cells or at the stage where the Cre gene is expressed.



In addition to conditional mutations, Cre-*lox* technology can be applied to a wide variety of targeted mutation strategies, including replacement of a gene from one species with the same gene from another species (25) and insertion of a different gene into a particular genomic location (gene “knock-in”; (26)). Cre-*lox* recombination has also been used for the creation of subtle mutations (27), to enable controlled activation of gene expression by insertion of a *loxP*-flanked stop sequence between a gene and its expression signals (28), to generate complex genomic models for chromosomal rearrangements, translocations, and deletions (29), and to model diseases caused by mutations in somatic cells, such as cancers and aging.

### 3.4. Null Mutations *in Vitro*

Cell lines in which both **alleles** of a gene are mutated provide unique opportunities for biochemical analysis of gene function. ES cells with homozygous mutations can be isolated from homozygous mutant embryos, or constructed by genetic manipulation of heterozygous ES cell lines generated by gene targeting. This can be achieved by sequential targeting using different selection markers, or by increasing selection for expression of the neomycin-resistance gene by elevating the concentration of G418 (30), thereby selecting for rare cells in which gene conversion has replaced the wild-type allele with a second copy of the targeted allele. The developmental potential of mutated ES cells can be assessed by differentiation of the cells *in vitro*, or by examining the contribution of the cells to tissues in chimaeric mice following return to the blastocyst. This can permit analysis of cell decisions that are otherwise obscured by embryonic lethality, and in some cases it can be coupled with detailed analysis of cell differentiation by measurement of marker gene expression.

Differentiated cell lines with homozygous mutations can be isolated from mutant embryos or via *in vitro* differentiation of ES cells containing two mutated alleles. These provide a valuable resource, because the effects of mutation can be examined in established *in vitro* assays specific for differentiated cell types. The availability of a homogeneous cell population allows application of biochemical techniques, such as direct measurement of downstream effector molecules.

### 3.5. Disease Models

Human diseases can have as their origin a wide variety of genetic alterations, including point mutations, deletions/truncations, large chromosomal deletions, and chromosomal translocations. Animals with genomes modified to mimic these alterations can be generated by gene targeting and used to study the cellular and biochemical basis of the disease and to test potential therapies. For example, familial hypercholesterolemia caused by loss of the low-density lipoprotein receptor has been modeled by gene knockout (31), cystic fibrosis has been modeled by the creation of point mutations within the cystic fibrosis transmembrane conductance regulator gene (32), and a familial form of Alzheimer's disease has been partially modeled by replacing the mouse Ab domain in the **amyloid precursor protein** with a mutated human sequence (33). A limitation of this approach can be the failure of mice to display the mutant phenotype seen in humans. In some cases, this can be overcome by refinement of the gene targeting approach. For example, familial adenomatous polyposis coli (FAP) results in the formation of colorectal adenomas in humans due to mutation of the APC gene. Mice homozygous for an APC deletion die *in utero*, but conditional gene targeting of the mutation in the colorectal epithelium using *Cre-lox* recombination resulted in mice that develop colorectal adenomas (34).

### 3.6. Gene Targeting in Non-mouse Species

Theoretically, gene targeting could be used in any cell that supports homologous recombination and can undergo clonal proliferation *in vitro* following selection for recombinants. However, creation of organisms carrying the modified allele can be achieved only if an intact organism can be regenerated from the genetically modified cell. The creation of gene-targeted organisms other than mice has been restricted by the failure to isolate pluripotent ES cells from these species. A potential alternative to the use of ES cell technology is foreshadowed by the demonstration that nuclei of differentiated cells can be reprogrammed or reverted to pluripotency by insertion into enucleated oocytes. The application of this nuclear transfer technology to produce offspring that are genetically identical to the donor cells has been validated for sheep, mice, and cattle, for donor nuclei from several cell types, and for genetically modified cells (35-38). Using this technique, gene targeting would be carried out in the differentiated cells, and genetically altered organisms would be created by nuclear transfer.

Gene targeting in nonmammalian metazoans is not yet well advanced. While homologous recombination has been demonstrated in moss (39), there are not yet convincing reports of gene targeting in lower vertebrates, invertebrates, or higher plants. The latter are of particular interest, given the potential for agricultural application.

### 3.7. Therapeutic Gene Targeting in Humans

An important potential application of gene targeting is the precise genetic alteration of **somatic cells** for correction of disease-causing mutations by gene therapy, delivery of therapeutic compounds to the individual, or tissue reconstitution *in vitro*. Direct genetic manipulation of, for example, somatic **stem cells** is limited by the failure of existing technologies to support isolation, transfection, and proliferation of these cell types *in vitro*. Widespread use of genetically altered cell types for human therapy may emerge from the combination of nuclear transfer technologies, which enable reprogramming of somatic cell nuclei, and gene targeting. Somatic cell nuclei isolated from an individual could be reprogrammed to a pluripotent state by nuclear transfer, modified genetically, and differentiated *in vitro* into the cell type required for transplantation. The use of gene targeting approaches could broaden the use of genetically modified cells for human therapy. This is currently restricted by difficulties associated with controlling the location, timing, and levels of gene expression, the copy number and integration site of **transgenes**, and the host **immune response** to virally based vectors.

### Bibliography

1. K. R. Thomas and M. R. Capecchi (1987) *Cell* **51**, 503–512.
2. E. P. Brandon, R. L. Idzera, and G. S. McKnight (1995) *Curr. Biol.* **5**, 625–634.
3. E. P. Brandon, R. L. Idzera, and G. S. McKnight (1995) *Curr. Biol.* **5**, 758–765.
4. E. P. Brandon, R. L. Idzera, and G. S. McKnight (1995) *Curr. Biol.* **5**, 873–881.
5. R. S. Johnson et al. (1989) *Science* **245**, 1234–1236.
6. P. Hasty, J. Rivera-Perez, and A. Bradley (1991) *Mol. Cell. Biol.* **11**, 5586–5591.
7. C. Deng and M. R. Capecchi (1992) *Mol. Cell. Biol.* **12**, 3365–3371.
8. H. te Riele, E. R. Maandag, and A. Berns (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5128–5132.
9. J. van Deursen and B. Wieringa (1992) *Nucleic Acids Res.* **20**, 3815–3820.
10. S. L. Mansour, K. R. Thomas, C. Deng, and M. R. Capecchi (1990) *Proc. Natl. Acad. Sci. USA* **87**, 7688–7692.
11. L. M. Whyatt and P. D. Rathjen (1997) *Nucl. Acids Res.* **25**, 2381–2388.
12. S. L. Mansour, K. R. Thomas, and M. R. Capecchi (1988) *Nature* **336**, 348–352.
13. K. R. Folger, E. A. Wong, G. Wahl, and M. R. Capecchi (1982) *Mol. Cell. Biol.* **2**, 1372–1387.
14. R. S. Kucherlapati et al. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3153–3157.
15. A. Bradley, P. Hasty, A. Davis, and R. Ramirez-Solis (1992) *Biotechnology* **10**, 534–539.
16. G. Friedrich and P. Soriano (1991) *Genes Dev.* **5**, 1513–1523.
17. P. Mountford and A. Smith (1995) *Trends Genet.* **11**, 179–184.
18. P. Hasty, R. Ramirez-Solis, R. Krumlauf, and A. Bradley (1991) *Nature* **350**, 243–246.
19. V. Valancius and O. Smithies (1991) *Mol. Cell. Biol.* **11**, 1402–1408.
20. L. H. Reid, R. G. Gregg, O. Smithies, and B. H. Koller (1990) *Proc. Natl. Acad. Sci. USA* **87**, 4299–4303.
21. A. Stacey et al. (1994) *Mol. Cell. Biol.* **14**, 1009–1016.
22. H. Wu, X. Liu, and R. Jaenisch (1994) *Proc. Natl. Acad. Sci. USA* **91**, 2819–2823.
23. G. R. Askew, T. Doetschman, and J. B. Lingrel (1993) *Mol. Cell. Biol.* **13**, 4115–4124.
24. P. J. Detloff et al. (1994) *Mol. Cell. Biol.* **14**, 6936–6943.
25. Y. R. Zou, W. Müller, H. Gu, and K. Rajewsky (1994) *Curr. Biol.* **4**, 1099–1103.
26. Y. Wang, P. N. Schnegelsberg, J. Dausman, and R. Jaenisch (1996) *Nature* **379**, 823–825.
27. R. M. Torres, H. Glaswinkel, M. Reth, and K. Rajewsky (1996) *Science* **272**, 1804–1808.
28. M. Lasko et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 6232–6236.
29. R. Ramirez-Solis, P. Liu, and A. Bradley (1995) *Nature* **378**, 720–724.
30. R. M. Mortensen et al. (1992) *Mol. Cell. Biol.* **12**, 2391–2395.



31. S. Ishibashi et al. (1993) *J Clin. Invest.* **92**, 883–893.
32. S. J. Delaney et al. (1996) *EMBO J.* **15**, 955–963.
33. A. G. Reaume et al. (1996) *J. Biol. Chem.* **271**, 23380–23388.
34. H. Shibata et al. (1997) *Science* **278**, 120–123.
35. K. H. S. Campbell, J. McWhir, W. A. Ritchie, and I. Wilmut (1996) *Nature* **380**, 64–66.
36. I. Wilmut et al. (1997) *Nature* **385**, 810–813.
37. J. B. Cibelli et al. (1998) *Science* **280**, 1256–1258.
38. T. Wakayama et al. (1998) *Nature* **394**, 369–374.
39. D. G. Schaefer and J. P. Zryd (1997) *Plant J.* **11**, 1195–1206.

### Suggestions for Further Reading

40. P. M. Wassarman and M. L. DePamphilis (eds.) (1993) In *Methods in Enzymology*, Vol. **225: Guide to Techniques in Mouse Development, Academic Press, San Diego.**
41. S. K. Bronson and O. Smithies (1994) *J. Biol. Chem.* **269**(44), 27155–27158.
42. M. R. Capecchi (1989) *Trends Genet.* **5**(3), 70–76.
43. P. Soriano (1995) *Annu. Rev. Neurosci.* **18**, 1–18.
44. L. A. Galli-Taliadoros, J. D. Sedgwick, S. A. Wood, and H. Korner (1995) *J. Immunol. Methods* **181**, 1–15.
45. B. Sauer (1998) *Methods (a companion to Methods Enzymol.)* **14**, 381–392.
46. A. Nagy and J. Rossant (1996) *J. Clin. Invest.* **97**, 1360–1365.

## Genetic Code

### 1. The Universal Progenitor and Standard Code

[Transfer RNA](#), tRNA, is the essential adapter in decoding the genetic information. The aminoacylation of specific tRNAs with particular [amino acids](#), and productive pairing of the [anticodons](#) of these tRNAs with **cognate** codons of the [messenger RNA](#) (mRNA) are the central processes in decoding (see [Aminoacyl tRNA Synthetases](#) and [Translation](#)). Within the past 10 years, great strides have been made in discovering the intricate relationships between tRNAs and aminoacyl-tRNA synthetases that are essential for correct aminoacylation. In contrast, the identification of mRNA codons has long been established ([1](#)), and the resulting table of specifications of amino acids by codons is what is commonly called the *genetic code* (Table [1](#)). However, the code has meaning only in terms of its readout and, since this is a dynamic process, it is necessary to consider the code and its readout together.

**Table 1. The Genetic Code**

---

| Second Position |   |   |   |   |                |
|-----------------|---|---|---|---|----------------|
| First Position  | U | C | A | G | Third Position |
| <hr/>           |   |   |   |   |                |

|   |                       |   |
|---|-----------------------|---|
|   | Phe Ser Tyr Cys       | U |
| U | Phe Ser Tyr Cys       | C |
|   | Leu Ser Stop Stop/Sec | A |
|   | Leu Ser Stop Trp      | G |
|   | Leu Pro His Arg       | U |
| C | Leu Pro His Arg       | C |
|   | Leu Pro Gln Arg       | A |
|   | Leu Pro Gln Arg       | G |
|   | Ile Thr Asn Ser       | U |
| A | Ile Thr Asn Ser       | C |
|   | Ile Thr Lys Arg       | A |
|   | Met Thr Lys Arg       | G |
|   | Val Ala Asp Gly       | U |
| G | Val Ala Asp Gly       | C |
|   | Val Ala Glu Gly       | A |
|   | Val Ala Glu Gly       | G |

---

All organisms living on earth are derived from a common ancestor, and [divergent evolution](#) from this ancestor is the central unifying theme of biology. One of the most striking facets of derivation from a common ancestor is the presence of the universal progenitor, or standard code, in nearly all organisms. Slight deviations from this standard code have appeared in certain organisms and coding niches, such as [mitochondria](#); these deviations are important and interesting, but they do not detract from the truth that *Escherichia coli*, elephants, and peas use the same genetic code for synthesizing nearly all their proteins.

Deciphering this standard code was triumphantly accomplished between 1961 and 1966 (1). Elegant genetic experiments showed that the code was read from a defined starting point in a nonoverlapping manner, with three bases (or formally a multiple thereof) specifying one amino acid, the next three bases the next amino acid, and so on (2). Confirmation of the triplet nature of the code came from cell-free **protein biosynthesis** experiments that showed which codons specified each amino acid. The template used in Nirenberg and Matthaei's pioneering experiment in 1961 was polyuridylylate [poly (U)], which led to the synthesis of polyphenylalanine. Consequently, the first codon to be deciphered was UUU for phenylalanine; very shortly afterward, AAA and CCC were found in the same way to encode lysine and proline, respectively. The use of mixed-polymer templates revealed the composition, but not the sequence, of codons specifying certain other amino acids. However, Khorana and colleagues synthesized polynucleotides with defined repeating sequences, and this provided valuable assignments (3). The other key development was the discovery by Nirenberg and Leder in 1964 (4) that trinucleotides promote the specific binding of tRNA to [ribosomes](#). Synthesis of all 64 trinucleotides led to assignments for about 50 codons by this method. Protein sequencing studies from mutants with combinations of **frameshift mutants** supported the conclusions. Furthermore, genetic studies by Brenner and Garen and their colleagues provided valuable insights into the termination codons.

By 1966, the various approaches together resulted in the complete deciphering of the code (1). It does not detract from this monumental achievement to point out that, if the code had not been established by these means, direct nucleic acid **sequencing** coupled with protein sequencing would have provided the answer shortly thereafter. Complete sequencing of the first gene was

accomplished in 1972 by “classical” RNA sequencing of the coat protein gene of the single-stranded RNA **bacteriophage** MS2 (5). Rapid DNA sequencing was developed in 1978 and led to a massive increase in sequence information. Nucleic acid sequencing, in conjunction with protein sequencing, has repeatedly confirmed the assignments of the standard code, demonstrated its utilization in decoding the main genomes of virtually all organisms, and revealed the interesting deviations.

The standard code in which the 64 possible base triplets encode 20 amino acids is presented in Table 1. The code is redundant, in that most amino acids are encoded by more than one codon. This is due in part to [wobble pairing](#) between the tRNA and the codon. Redundancy leads to variations in **codon usage** in different organisms. Three of the 64 codons do not usually code for amino acids but are used as [stop codons](#), to terminate translation. Certain codons, usually AUG, play a special role as [start codons](#), initiating translation of the mRNA. The 21st amino acid, [selenocysteine](#) (Sec), is also directly encoded, but only in special cases (see “Recoding,” below).

## 2. Specialized Code Variants: the Mitochondrial Code

It was long considered that the genetic code could not be altered during evolution. It was thought that changes in codon meaning would have such a deleterious effect on protein products that the code was “frozen.” However, quite a number of cases of genetic code variation are now known in mitochondria, as well as in the nuclear systems, of some organisms (6-9). The known variants of the genetic code in mitochondria are shown in Table 2. The small number of products encoded by animal mitochondrial genomes has undoubtedly been important in permitting code alterations in these organelles. How these alterations might have occurred is considered in the section “Candida code and code reassignment theories.”

**Table 2. Variations in Mitochondrial Genetic Code<sup>a</sup>**

| <b>Organisms</b> | <b>UGA<br/>Stop</b> | <b>AUA<br/>Ile</b> | <b>AAA<br/>Lys</b> | <b>AGR<br/>Arg</b> | <b>CUN<br/>Leu</b> | <b>UAA<br/>Stop</b> | <b>Example</b>   |
|------------------|---------------------|--------------------|--------------------|--------------------|--------------------|---------------------|--|
| Vertebrates      | Trp                 | Met                | —                  | Stop               | —                  | —                   | Human, frog  |
| Tunicates        | Trp                 | Met                | —                  | Gly                | —                  | —                   | <i>Halocynthia roretzi</i>                                   |
| Echinoderms      | Trp                 | —                  | Asn                | Ser                | —                  | —                   | Starfish, sea urchin   |
| Arthropods       | Trp                 | Met                | —                  | Ser <sup>b</sup>   | —                  | —                   | <i>Drosophila</i> spp.,<br>mosquito, honeybee                |
| Molluscs         | Trp                 | Met                | —                  | Ser                | —                  | —                   | Squid, <i>Mytilus edulis</i>                                 |
| Nematodes        | Trp                 | Met                | —                  | Ser                | —                  | —                   | <i>C. elegans</i> , <i>Ascaris suum</i>                      |
| Platyhelminthes  | Trp                 | —                  | Asn                | Ser                | —                  | Tyr                 | <i>Fasciola hepatica</i> ,<br><i>planaria</i>                |
| Cnidarians       | Trp                 | —                  | —                  | —                  | —                  | —                   | <i>Hydra</i> , <i>Metridium senile</i>                       |
| Yeasts Trp       | Met                 | —                  | —                  | Thr                | —                  | —                   | <i>Saccharomyces cerevisiae</i> , <i>Torulopsis glabrata</i> |
| Eucomycetes      | Trp                 | —                  | —                  | —                  | —                  | —                   | <i>Aspergillus nidulans</i> ,<br><i>Neurospora crassa</i>    |
| Protozoa         | Trp                 | —                  | —                  | —                  | —                  | —                   | <i>Trypanosoma brucei</i> ,<br><i>Paramecium</i> spp.        |

---

<sup>a</sup> From reference [7](#) with modifications.

<sup>b</sup> N = A, U, C, or G; R = A or G. \* Only AGA (specifying Ser) is used in *Drosophila* and mosquitoes, while both AGA and AGG (each specifies Ser) are used in honeybee, locust, and brine shrimp.

An important feature of decoding in animal mitochondria is that, in many cases, one tRNA recognizes all four members of a codon family. This simplification has resulted in a greatly reduced number of tRNAs compared to cytoplasmic protein synthesis. Key to this relaxed specificity is an unmodified U at the first position of the anticodon of many tRNAs ([10](#)). The general point is that almost all the diversity in the genetic code can be explained by variations in the anticodons of tRNAs that bind to various codons.

Generally, animal mitochondrial mRNAs have a simpler structure than those of **prokaryotes** and other **eukaryotic mRNAs**. There is no 5'-leader sequence ([11](#)) and usually only a short **poly A** sequence in the 3'-trailer sequence ([12](#)). Metazoan mitochondrial mRNA usually start translation with an AUG or AUA codon and end with a UAA or UAG codon although, in mammalian mitochondria, AGA and AGG codons are used for termination ([12](#), [13](#)). Often, only the 5' nucleotide of the UAA is directly encoded, and it is completed by addition of the poly A sequence at the 3'-end of the primary transcript of mRNAs ([12](#)).

The change of UGA from being a stop codon to specifying tryptophan (Trp) in all mitochondria, except those in plants, is due to change of the anticodon of tRNA<sup>Trp</sup> from CCA or CmCA (mC = 5-methyl C) to U\*UA, where U\* in most cases is cmnm<sup>5</sup>U (5-carboxymethyl-aminomethyl U), so that the tRNA<sup>Trp</sup> can decode both UGG and UGA codons by wobble pairing.

The tRNA<sup>Met</sup> capable of decoding both AUG and AUA codons in most mitochondria, except those of echinoderms and platyhelminthes, possesses at the first position of its anticodon 5-formyl C in vertebrates ([14](#)) and nematodes ([15](#)), cmnm<sup>5</sup>C in ascidia ([16](#)), and a modified C in addition to normal C in *Drosophila* ([17](#)). These modified nucleosides make it possible to decode A as well as G at the third codon position.

AAA codons are read as Asn, instead of the usual Lys, in starfish mitochondria; the first position of the corresponding tRNA anticodon is occupied by G ([17](#)), which may decode, at the first position, not only C and U but also A. The same is observed in Ser codons of *Drosophila* mitochondria, where AGU, AGC, and AGA codons (AGG is an unassigned codon) specify Ser. Ser is also specified by UCN codons (N = A, U, C or G). The single tRNA<sup>Ser</sup> responsible for decoding these codons possesses a GCU anticodon ([17](#)).

All AGN codons are read as Ser in most invertebrate mitochondria. In these cases, a modified G [7-methyl G in starfish and squid ([17](#), [18](#))] or an unmodified U [in nematodes ([15](#))] occupies the first anticodon position, so that these bases can recognize all four bases at the third codon position.

In ascidian mitochondria, the AGR codon (where R is A or G) is assigned to Gly. The anticodon of the corresponding tRNA<sup>Gly</sup> is cmnm<sup>5</sup>UCU ([19](#)).

### 3. Candida Code and Codon Reassignment Theories

In nine species of the asporogenic yeast, *Candida*, CUG encodes serine and not leucine as it does in the standard, nearly universal code ([20-22](#)) (Table [1](#)). The tRNA translating this nonuniversal genetic code (tRNA<sup>Ser</sup> CAG) has been characterized ([21](#)) and shown to be charged with serine ([21-26](#)).

At present, there are two theories to explain the codon reassignment: the “codon capture” (26) and the “ambiguous intermediate” (27) theories. The former postulates that a codon changes to another codon via an unassigned codon and that the codon to be changed must disappear once from coding sequences to avoid insertion of an incorrect amino acid (neutral change of codon). The latter postulates that changes in the genetic code are relatively fast processes driven by [natural selection](#) and that reassignment of codons is facilitated by a translationally ambiguous intermediate, where the transitional codon is read simultaneously as either of two amino acids by two different tRNAs.

The genetic code change found in *Candida* species may be explained by either of these two theories. Santos et al (28) favor the ambiguous intermediate theory. They provided experimental evidence that expression in yeast cells of a tRNA<sup>Ser</sup> CAG with a mutation at position 33 (5'-adjacent to the anticodon) induces the **stress response** that allows cells to acquire thermotolerance. This may provide positive selection for the genetic code change by allowing yeast to adapt to sudden changes in environmental conditions.

On the other hand, Suzuki et al. (29) found what they term a “polysemous codon” in *Candida zeylanoides*, which used the CUG codon to encode both serine and, to a slight extent (3%), leucine. This is mediated by a single tRNA<sup>Ser</sup> CAG capable of being charged with either amino acid. A genetic method demonstrated that the CUG codon is actually translated as either amino acid *in vivo*. Suzuki and coworkers emphasize that the dual assignment of the CUG codon is fulfilled by dual specificity of a single tRNA<sup>Ser</sup> CAG, not by ambiguous translation using two different tRNA. Thus, they prefer the codon capture theory to the ambiguous intermediate theory. The existence of unassigned codons is also important evidence for the codon capture theory. Further experiments are necessary, however, to clarify unambiguously the evolutionary process of codon reassignment in *Candida* species (30).

#### 4. Theories of the Origin of the Genetic Code

Two assumptions are commonly made in considering the origin of the genetic code.

1. First is the principle of continuity; namely, for any possible productive early scheme, there has to be some conceptual way that it could have evolved gradually to the modern genetic code without destroying all the fruits of prior evolution.
2. A second assumption is that the origin of the code was based on RNA. Since proteins are decoded from RNA templates and some RNA are catalytic, it is highly attractive to invoke a transitory “[RNA world](#)”. Perhaps no decoding was necessary if the genetic material (RNA) could mediate the catalytic functions directly.

Since discussion of the origin of the code which began in earnest in the late 1960s (31), one of the key imponderables has been whether there was originally a direct stereochemical relationship between primordial tRNAs and their cognate amino acids or whether the genetic code assignments were random, as in “the frozen accident theory” (32). Recent work (33) has shown selective binding of some amino acids by RNA, and the question of aminoacylation in the absence of proteins is being investigated, but the overall issue remains unresolved. Intertwined with this issue is the question of the relationship between the acceptor and anticodon branches of tRNA. Noller has suggested that both parts were initially separate, with overlapping functions, and that they were later combined (34). Another suggestion, which is difficult to evaluate, is that they are directly related (35).

Another key imponderable in decoding is the origin of ribosomes. Several different approaches, including extensive but incomplete deproteinization experiments, have pointed to the largest rRNA (23S rRNA in *E. coli*) as the key component in peptidyl transferase (reviewed in (34)). Of course, even imagining the origin of 23S rRNA is a serious problem, although an attractive suggestion is that a protodomain originated as a spin-off of a ribozyme-catalyzed function (34). A different issue is the origin of [translocation](#) of the ribosome one codon at a time along the mRNA. Modern ribosomes

have the ability of low-level synthesis of polypeptides on poly(U), (UC), or (A) in the absence of the [elongation factors](#) EF-Tu, EF-G, and GTP, suggesting that they contain the mechanism for catalyzing translocation (36) and putting the spotlight on the small subunit rRNA (16S rRNA in *E. coli*). An interesting speculation is that this originated as an RNA replicase (37). Regardless of origin, it is tempting to think that the translocation step size is determined in part by the tRNA anticodon delineating the codon size. The stability of a triplet interaction between codon and anticodon is insufficient on its own to mediate discriminatory decoding. If, for stability reasons, however, one envisages primordial tRNAs that read quintuplet codons in a primitive mRNA, it is difficult to imagine how those translation components could evolve to triplet decoding without destroying the benefits of prior evolution. Consideration of this difficulty has led to an ingenious scheme for triplet decoding involving overlapping quintuplet pairing (38, 39), but whether such a scheme is a progenitor of modern triplet decoding is speculative. It has been proposed that in present-day decoding, the triplet codon interaction is stabilized by stacking ribosomal RNA (40). The particular regions of rRNA proposed have not been implicated in genetic studies to address this issue, but it is unclear whether this is because they are not involved or because they are essential (41).

Knowledge of the structure of the ribosome will be invaluable for efforts to understand present-day decoding which, in turn, will guide thoughts about plausible schemes for its evolutionary origin.

## 5. Recoding: Reprogrammed Decoding

### 5.1. Transient Reprogramming of Decoding

Excluding some specialized niches, the genetic information in the great majority of mRNA is readout according to the standard rules with the nearly universal code specifications (Table 1). However, a minority of mRNA that use the same translational components, or essentially so, are decoded in a different manner at one or more positions in their coding sequence. At these positions, readout is reprogrammed, and the phenomenon is termed *recoding*. In the first category, the meaning of a code word is redefined so that at least some of the protein product bears the consequence of the new meaning. In the second category, a proportion of ribosomes shift [reading frame](#) at a specific place in an mRNA and continue in the new frame, following the standard rules, to give a very different product from the ribosomes that did not shift frame (see [Frameshifting](#)). In the final category, a block of nucleotides is bypassed by the translocating ribosomes.

Recoding has been most frequently encountered in **viruses**, especially plant viruses and retroviruses, and in [retrotransposons](#), such as yeast [Ty elements](#) and bacterial **insertion sequence (IS) elements**, but it is also known to be utilized in decoding several cellular genes. It probably occurs in all organisms.

### 5.2. Redefinition

The meaning of a sense codon can be redefined; for example, the codon GUG at internal positions specifies valine, but in the context in which it functions as an initiator, it specifies methionine or formylmethionine (see [Start Codons](#)). No cases are known in which the meaning of a sense codon at an internal position in a coding sequence is redefined. (Though mechanistically distinct from recoding, there are instances in which a sense codon gives the superficial appearance of acting as a stop codon. Consecutive rare codons, or some RNA structures, can cause the ribosome to pause and, on occasion, to drop off with a terminated polypeptide chain; see [Stop Codons](#).)

On the other hand, stop codons in the standard code can be redefined to specify an amino acid. The most dramatic case involves the specification of the 21st encoded amino acid, [selenocysteine](#). Until 1986 it was thought that only 20 amino acids were directly encoded (excluding the transient case of formyl methionine; see [Start Codons](#)). Then it was found—in **eubacteria**, **archae**, and mammals—that selenocysteine was directly encoded and that its codon is UGA, which in the standard code is a stop codon (42, 43). UGA specifies selenocysteine in only a very small number of species of mRNA in each cell. In nearly all other mRNA species, UGA specifies translation termination. In each

mRNA in which UGA encodes selenocysteine, there is a particular mRNA structure that acts as a recoding signal to specify that UGA means “selenocysteine” in decoding this particular mRNA. In **Gram-negative** eubacteria, this structure is a particular stem-loop structure immediately following the UGA codon (44). In mammals, the equivalent specifying element, SECIS, is in the 3' untranslated region (45). In archae, the available evidence indicates that the enabling structure can be in either the 3' or 5' untranslated region of the mRNA (46). Elegant studies in *E. coli* have shown that a special elongation factor, which does not form a complex with any tRNA other than selenocysteine tRNA, binds to the loop of the mRNA stem loop structure and effectively tethers the aminoacylated selenocysteine for delivery to the ribosome poised at the UGA (44). How the SECIS element in the 3' untranslated region of mammalian mRNA that encode selenoproteins functions is a fascinating mystery. The problem is acute in human selenoprotein P mRNA, where 10 UGA codons occur within the coding region, and all probably specify selenocysteine.

The redefinition of stop to other sense codons is also well known, and there is considerable variation in the extent of recoding signals involved in the process. With retroviruses other than Spumaretroviruses, there is no independent ribosome entry to the [pol gene](#) that encodes **reverse transcriptase** and an **endonuclease**. Instead, ribosomes that initiate translation of the upstream *gag* gene synthesize some Gag-Pol fusion protein, in addition to the separate Gag product. In murine leukemia and related viruses, the *pol* gene is in the same frame as the *gag* gene and is separated from it only by a UAG “stop” codon. Approximately 10% of the ribosomes that read the *gag* gene insert glutamine at the UAG codon and synthesize a Gag-Pol fusion. This level of “readthrough” is dependent on a recoding signal in the mRNA and is 10<sup>3</sup>-fold higher than the background error level in the absence of such signals. The recoding signal is a particular pseudoknot eight bases after the UAG (47). How this pseudoknot influences the oncoming ribosomes is unexplained. Redefinition of “stop” codons is common in plant virus decoding. In barley yellow dwarf virus, both local and distant signals are important for the redefinition (48) whereas, in [tobacco mosaic virus](#), only the sequence of the six bases following the stop codon is relevant.

### 5.3. Programmed Ribosomal Frameshifting

At specific positions in the decoding of some genes, a high proportion of ribosomes shift frame and then resume triplet decoding to give an essential product (49) (see [Frameshifting](#)). This probably happens in all organisms, but notable examples are synthesis of the Gag-Pol polyprotein of many retroviruses, including HIV-1, human antizyme, barley yellow dwarf virus polymerase, *E. coli* release factor 2, and [lambda phage](#) G-T product. Some cases of programmed frameshifting serve an autoregulatory function (eg, *E. coli* release factor 2 and human antizyme); others produce different products from a single gene in a set ratio. The recoding signals range from distant sequences to a flanking 3' pseudoknot, stem-loop structure, or rare codons, to 5' sequences, or to a combination. The frameshifting is +1 in some cases (eg, human antizyme and *E. coli* release factor 2) and -1 in others (eg, HIV, barley yellow dwarf virus and mouse mammary tumor virus). Many cases of programmed -1 frameshifting involve the simultaneous slippage of the pairing of two tRNA anticodons with the mRNA. In contrast, in programmed +1 frameshifting, pairing of a single tRNA anticodon is often disrupted and the codon that enters the A site is either a stop codon with no cognate tRNA or a rare codon whose cognate tRNA is in short supply. Repairing of the tRNA anticodon 1 nucleotide along the mRNA has greatly reduced, or no, competition from an incoming tRNA specified by the zero frame codon.

### 5.4. Bypassing

Between codons 46 and 47 of phage T4 gene 60, there are 50 nucleotides that are translationally bypassed with high efficiency (50). This remarkable instance of decoding poses many mechanistic questions, and it is not clear how widespread the phenomenon of bypassing is. However, low-level error bypassing has been encountered when overexpressing, for biotechnological purposes, a mammalian protein in *E. coli* (50).

## 6. Expanding the Code Experimentally

Redesigning the genetic code to allow the incorporation of amino acids that are not normally encoded by natural genes is a challenging task. The potential advantages of having unnatural amino acids at defined places in proteins are that they could be molecular beacons for structural studies and could open the potential for new structural and functional capabilities. Short proteins can be produced by solid-phase synthesis (see [Peptide Synthesis](#)), but longer proteins with unnatural amino acids are also desirable. The most ambitious approach is to create a 65th codon-anticodon pair from unnatural nucleoside bases having nonstandard hydrogen-bonding patterns ([52](#)). The more approachable, but still very difficult, task is to use the standard nucleotides to generate *in vivo* a tRNA that is competitive for reading a particular codon and aminoacylated with an unnatural amino acid, such as a keto-containing amino acid ([53](#)). Some success has been achieved *in vitro* ([54-58](#)) but, because of severe yield limitations, it is desirable to achieve the same goal *in vivo*.

A very promising start *in vivo* has recently been made ([53](#)). An essential requirement for a tRNA that delivers an unnatural amino acid is that it is not aminoacylated by the endogenous aminoacyl-tRNA synthetase. Instead, it needs to be recognized by a mutant synthetase and acylated with the unnatural amino acid. As a first step, Schultz and colleagues ([53](#)) isolated mutants of *E. coli* glutamyl-tRNA synthetase that aminoacylate an engineered mutant of tRNA<sup>Gln</sup> that is not aminoacylated by the wild-type synthetase. This is a considerable feat, but major hurdles remain.

## Bibliography

1. (1966) Cold Spring Harbor Symp. Quant. Biol. **31**, 1–762.
2. F. H. C. Crick, L. Barnett, S. Brenner, and R. J. Watts-Tobin (1961) *Nature* **192**, 1227–1232.
3. H. G. Khorana et al. (1966) Cold Spring Harbor Symp. Quant. Biol. **31**, 39–49.
4. M. Nirenberg and P. Leder (1964) *Science* **145**, 1399–1407.
5. W. Min Jou, G. Haegerman, M. Ysebaert, and W. Fiers (1972) *Nature* **237**, 82.
6. S. Osawa (1995) *Evolution of the Genetic Code*, Oxford University Press, Oxford, U.K., pp. 1–205.
7. S. Osawa, T. H. Jukes, K. Watanabe, and A. Muto (1992) *Microbiol. Rev.* **56**, 229–264.
8. K. Watanabe and S. Osawa (1995) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology, Washington, D.C., pp. 225–250.
9. D. R. Wolstenholme and C. M. Fauron (1995) *The Molecular Biology of Plant Mitochondria* (C. S. Levings III and I. K. Vasil, eds.), Kluwer Academic Publishers, the Netherlands, pp. 1–59.
10. J. E. Heckman et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3159–3163.
11. J. Montoya, D. Ojala, and G. Attardi (1981) *Nature* **290**, 465–470.
12. S. Anderson et al. (1981) *Nature* **290**, 457–465.
13. D. R. Wolstenholme (1992) *Intern. Rev. Cytology* **141**, 173–216.
14. J. Moriya et al. (1994) *Biochemistry* **33**, 2234–2239.
15. Y. Watanabe et al. (1994) *J. Biol. Chem.* **269**, 22902–22906.
16. A. Kondow, S. Yokobori, T. Ueda, and K. Watanabe (1998) *Nucleosides and Nucleotides*, **17**, 531–539.
17. K. Watanabe et al. (1997) in *17th International tRNA Workshop*, Abstract, Chiba, Japan, p. 12–8.
18. S. Matsuyama et al. (1998) *J. Biol. Chem.*, **273**, 3363–3368.
19. A. Kondo, S. Yokobori, T. Ueda, and K. Watanabe (1996) *Nucleic Acids Symp. Ser.* **35**, 279–280.
20. Y. Kawaguchi, H. Honda, J. Taniguchi-Morimura, and S. Iwasaki (1989) *Nature* **341**, 164–166.
21. T. Yokogawa et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 7408–7411.
22. T. Ohama et al. (1993) *Nucleic Acids Res.* **21**, 4039–4045.



23. T. Ueda et al. (1994) *Biochimie* **76**, 1217–1222.
24. M. A. S. Santos and M. F. Tuite (1995) *Nucleic Acids Res.* **23**, 1481–1486.
25. H. Sugiyama et al. (1995) *Yeast* **11**, 43–52.
26. S. Osawa and T. H. Jukes (1989) *J. Mol. Evol.* **4**, 191–198.
27. D. W. Schultz and M. Yarus (1994) *J. Mol. Biol.* **235**, 1377–1380.
28. M. A. S. Santos, V. M. Perreau, and M. F. Tuite (1995) *EMBO J.* **15**, 5060–5068.
29. T. Suzuki, T. Ueda, and K. Watanabe (1997) *EMBO J.* **16**, 1122–1134.
30. T. H. Jukes, S. Osawa, M. Yarus, and D. W. Schultz (1997) *J. Mol. Evol.* **45**, 1–8.
31. L. E. Orgel and F. H. C. Crick (1993) *FASEB J.* **7**, 238–239.
32. F. H. C. Crick (1968) *J. Mol. Biol.* **38**, 367–379.
33. M. Yarus (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 205–217.
34. H. F. Noller (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 137–156.
35. S. N. Rodin and S. Ohno (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5183–5188.
36. N. M. Rutkevitch and L. P. Gavrilova (1982) *FEBS Lett.* **143**, 115–118.
37. R. Weiss and J. Cherry (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 71–89.
38. C. R. Woese (1970) *Nature* **226**, 817–820.
39. F. H. C. Crick, S. Brenner, A. Klug, and G. Piecznik (1976) *Origins Life* **7**, 389–397.
40. H. F. Noller et al. (1986) In *Structure, Function and Genetics of Ribosomes* (B. Hardesty and G. Kramer, eds.), Springer-Verlag, New York, pp. 143–163.
41. M. O'Connor and A. E. Dahlberg (1996) *Nucl. Acids Res.* **24**, 2701–2705.
42. F. Zinoni, A. Birkmann, T. C. Stadtman, and A. Böck (1986) *Proc. Natl. Acad. Sci. USA* **83**, 4560–4564.
43. I. Chambers et al. (1986) *EMBO J.* **5**, 1221–1227.
44. A. Hüttenhofer, E. Westhof, and A. Böck (1996) *RNA* **2**, 345–366.
45. M. J. Berry, L. Banu, J. W. Harney, and P. R. Larsen (1993) *EMBO J.* **12**, 3315–3322.
46. R. Wilting, S. Schorling, B. C. Persson, and A. Böck (1997) *J. Mol. Biol.* **266**, 637–641.
47. N. M. Wills, R. F. Gesteland, and J. F. Atkins (1994) *EMBO J.* **13**, 4137–4144.
48. C. M. Brown, S. P. Dinesh-Kumar, and W. A. Miller (1996) *J. Virol.* **70**, 5884–5892.
49. J. F. Atkins, R. B. Weiss, and R. F. Gesteland (1990) *Cell* **62**, 413–423.
50. R. B. Weiss, W. M. Huang, and D. M. Dunn (1990) *Cell* **62**, 117–126.
51. J. F. Kane et al. (1992) *Nucl. Acids Res.* **20**, 6707–6712.
52. J. D. Bain, C. Switzer, A. R. Chamberlin, and S. A. Benner (1992) *Nature* **356**, 537–539.
53. D. R. Liu, T. J. Magliery, M. Pastnak, and P. G. Schultz (1997) *Proc. Natl. Acad. Sci. USA* **94**, 10092–10097.
54. C. J. Noren, S. J. Anthony-Cahill, M. C. Griffith, and P. G. Schultz (1989) *Science* **244**, 182–188.
55. J. D. Bain et al. (1989) *J. Am. Chem. Soc.* **111**, 8013–8014.
56. H.-H. Chung, D. R. Benson, and P. G. Schultz (1993) *Science* **259**, 806–809.
57. V. W. Cornish et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 2910–2914.
58. P. E. Dawson, M. C. Fitzgerald, T. W. Muir, and S. B. H. Kent (1997) *J. Amer. Chem. Soc.* **119**, 7917–7927.

### Suggestions for Further Reading

59. S. Osawa, T. H. Jukes, K. Watanabe, and A. Muto (1992) Recent evidence for evolution of the genetic code. *Microbiol. Rev.* **56**, 229–264.
60. P. Schimmel (1996) Origin of the genetic code: A needle in the haystack of tRNA sequences. *Proc. Natl. Acad. Sci. USA* **93**, 4521–4522.
61. A. Hüttenhofer and A. Böck (1997) "RNA structures involved in selenoprotein synthesis". In *RNA, Structure and Function* (R. W. Simon and M. Grunberg-Manago, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 603–639.
62. S. C. Low and M. J. Berry (1996) Knowing when not to stop: selenocysteine incorporation in eukaryotes. *Trends Biochem. Sci.* **21**, 203–208.
63. R. F. Gesteland and J. F. Atkins (1996) Recoding: Dynamic reprogramming of translation. *Ann. Rev. Biochem.* **65**, 741–768.
64. J. Castresana, G. Feldmaier-Fuchs, and S. Pääbo (1998) Codon reassignment and amino acid composition in hemichordate mitochondria. *Proc. Natl. Acad. Sci. USA* **95**, 3703–3707.

## Genetic Disease

All of the 30,000 or more human genes can mutate, and many of the mutant alleles cause or contribute to the risk of genetic disease. In some instances, pathology occurs if only one of the two alleles at a particular locus is mutant. In other instances, both alleles must be defective. There are several sources of genetic heterogeneity in disease. Different alleles at a single locus may produce different degrees of pathology or even quite different disorders. In contrast, mutations at different loci can produce the same phenotypic changes. Mutations that arise in somatic cells can also produce disease, but they cannot be transmitted to offspring.

### 1. Disease, Genes, and Mutations

A discussion of genetic disease should begin with a definition of disease. Most persons use the term to indicate a departure from normal function or health. This suggests that a clear distinction exists between the healthy state and disease. Although this dichotomy may be useful to insurance companies, it is not very descriptive of biological reality. Each of us is uniquely programmed genetically to cope with greater or less success with what the environment presents to us. And ultimately we are programmed not to survive. This perception was expressed by the great British physician A.E. Garrod in his 1931 monograph *The Inborn Factors in Disease*, "It is permissible to suppose that, as in the production of those wider departures from the standard of health which receive the name of diseases, both internal and external factors are at work, so also in the causation of those lesser departures, which it is difficult to name or classify, and which may be spoken of as trifling ailments,  $\frac{1}{4}$ "

Diseases, then, may be thought of as deviations from normality that are sufficient to interfere with function of the organism. Many such deviations are primarily associated with specific genotypes. Others are combinations of genotypic and environmental variation. Although all are of interest, our knowledge is largely limited to those diseases that are well outside the range of normal function. In the case of inherited diseases, this means that we are largely limited to considering mutations at single loci that have a large impact on the phenotype. Results from the Human Genome Project, coupled with more sophisticated statistical techniques, are making possible the detection of smaller and more complex genetic contributions to disease risk.

As of May 2001, some 9,200 loci had been identified in the human genome. Over 7,000 loci had been assigned to specific chromosomes. Many are associated with a specific inherited variation or disease. Others are defined by the existence of their protein product or messenger RNA transcript. In some instances, a gene first observed in a nonhuman species is found to have a homolog in humans. The total number of loci is estimated to be approximately 30,000. This number is derived from the complete sequence of human DNA and is much lower than earlier estimates.

All genes have the potential to mutate. It is not known, however, what portion of the total number have the potential to mutate to a form that leads to recognizable disease. This is because many genes may cause early embryonic death if mutations occur in them. For example, any mutant form of a gene that is critical in early development may be eliminated without being observed.

## 2. Genotypes and Disease

The ultimate role of genes is the production of proteins of the right kind and amount in the right places and times. The vast majority of genes code for protein structures, and associated with each gene are regulatory elements—stretches of DNA that interact with other proteins to determine the amount of transcription of the gene. Most of the understanding of inherited diseases is in terms of function or malfunction of specific proteins.

Whether a mutant gene causes disease depends on several factors. Insertion or deletion of a single nucleotide in the coding region causes a shift in the reading frame during translation, leading to complete absence of the gene product. Substitution of one nucleotide for another may be equally detrimental if it causes substitution of the wrong amino acid at a functional site. Or it may have no effect at all if the amino acid substitution does not modify the function of the protein. Many such variations exist in the normal human population. Some cause small differences in function that we consider normal variation. Mutations in regulatory regions of genes may change the amount of gene product—too little or too much—but not the structure of the protein.

The phenotypic effect of these changes in gene structure depend on the role of the gene. The amounts of many enzymes normally produced are in excess of the amount needed to maintain normal metabolism. These loci are described as haplosufficient, meaning that a heterozygote that has one normal and one nonfunctional allele is phenotypically normal. Many enzyme deficiencies are inherited as recessive traits because of the need for both alleles to be *nonfunctional* for disease to occur. Haplosufficiency also characterizes other types of proteins, such as cell cycle regulators and many receptors.

Haploinsufficiency is also observed. For example, one normal copy of a  $\beta$ -globin gene in hemoglobin is insufficient for manufacture of adequate amounts of  $\beta$ -globin, and anemia results. Again, in the case of low-density lipoprotein (LDL) receptors, one normal allele is insufficient to produce the number of functional receptors required to maintain normal levels of blood cholesterol. The transmission of these traits produces pedigrees characteristic of dominant inheritance.

One group of mutations is designated dominant negative. The name derives from the fact that the mutant allele interferes with functions of the normal allele. In a typical example, the proteins form complexes either involving one type of subunit (i.e., the product of one locus) or involving multiple subunits from different loci. The product of the mutant allele may enter into the complex but interfere with function of the complex. In the simplest case of a homodimer formed from two identical subunits, only the dimer composed of two normal subunits is functional. Both the mixed dimer and the dimer composed of two mutant subunits are nonfunctional. Because heterozygotes have reduced function, the transmission pattern is that of a dominant trait.

A variation on the previous example occurs when the mutant allele produces no protein product at all. In this case, the only subunits are produced from the normal allele, and all dimers or other multimers would be functional. This is observed in some collagen disorders, where deletion of an

allele causes no phenotypic effect, the normal allele being haplosufficient. Only homozygotes for the mutant allele would generate a mutant phenotype. However, if the mutation is a minor variation in the amino acid sequence, this may interfere with production of normal collagen. The trait would appear as dominant in pedigrees.

Not all mutations or inherited syndromes involve single genes. In contiguous gene syndromes, small chromosomal deletions remove multiple closely-linked genes. The person therefore has only one copy of the genes in that segment rather than the normal two. Any genes that are haploinsufficient will affect the phenotype. The WAGR syndrome involves deletion of a small region of chromosome 11 that includes genes for *Wilms' tumor*, *Aniridia*, *Genital anomalies*, and *Retardation of growth*. The WAGR syndrome is therefore the sum of the single gene defects.

### 3. Genetic Heterogeneity

If one considers all the possible nucleotide substitutions, insertions, and deletions in a gene that is many kilobases in size, the number of potential alleles at any one locus is enormous. In practice, one would never expect to see all possible mutant alleles in a population of finite size. However, in the case of other loci that have been extensively investigated, such as the globins, phenylalanine hydroxylase, cystic fibrosis, and Duchenne muscular dystrophy (DMD), hundreds of different mutant alleles have been identified. For the most part, these are individually very rare. In the case of the *DMD* locus, mutations are quickly eliminated by natural selection, and most mutations are found on analysis to be recent and to differ from other known mutations. In the case of other loci, a particular one or two mutations will be common, reflecting an increase in the frequency through genetic drift, founder effect, or heterozygote advantage. In the case of cystic fibrosis, some 70% of the mutant alleles in persons of European ancestry involve the identical three-nucleotide deletion, one that must have occurred many thousands of years ago to be so widely distributed.

Such allelic heterogeneity often translates into phenotypic heterogeneity. This is well illustrated in mutations of the *DMD* gene. Many of the mutations are associated with a protein product (dystrophin) that has some activity. This results in a milder form of the disease (Becker muscular dystrophy). Similar variations in severity of cystic fibrosis have been associated with different allelic combinations.

Because of the very large number of alleles at some loci, the term homozygous is often a misnomer. For rare recessive traits, it means that the two alleles are defective in their function but not that they are identical. A person with two different alleles that are functionally similar is described as a compound heterozygote. Unless the genes have been analyzed at the molecular level, one cannot be certain of true homozygosity.

Locus heterogeneity also occurs. This is the situation in which mutations at any of several loci produce the same phenotype. For example, early onset familial Alzheimer disease (FAD) is transmitted as a dominant trait and can be caused by mutations of any of three loci on three different chromosomes. The primary defect is presumably different at the molecular level, but the phenotypic results are the same.

An uncommon form of genetic heterogeneity is the production of different pathological conditions depending on where mutations occur in a gene. An example is the androgen receptor gene (*AR*), located on the X chromosome. The androgen receptor, when bound to testosterone and dihydrotestosterone, acts as a transcription factor in the nucleus. A number of mutations are known that interfere with binding and prevent the androgens from acting on target organs. The resulting disorder is known as androgen resistance. Mutations at other sites in the *AR* gene cause spinobulbar muscular atrophy.

Different populations also vary in their complements of mutant alleles. The mutations that cause cystic fibrosis occur in some 2% of persons of European origin but occur at much lower frequencies

in other populations. Tay-Sachs occurs at much higher frequencies among Ashkenazi Jews. These high frequencies of detrimental genes probably arose through genetic drift and founder effect, although selective advantage may account for the high frequencies in some instances, as in sickle cell anemia in Africans. Even when the overall frequencies of mutant alleles are similar in populations, the specific alleles may differ.

#### 4. Genetic Diseases of Somatic Cells

When inherited diseases are discussed, the traditional reference is to diseases that are transmitted from one generation to the next. Many show simple dominant or recessive inheritance according to Mendelian rules. Others are more complex, being influenced by variation at multiple loci and by environment also. The variations that can be followed in pedigrees occur in the germ lines and are therefore transmissible, even though the phenotype that we observe is based on somatic cell function.

Somatic cells, of course, have the same array of genes that are in the germ line, and they are subject to mutation as well as epigenetic alterations. The great importance of somatic cell genetics became apparent from the demonstration that mutations in somatic cells are essential parts of carcinogenesis. In the simplest case of retinoblastoma, a malignant tumor occurs when both alleles of the *RB* locus become inactive through mutation in a retinoblast. The normal restraints on cell growth no longer occur, and that cell forms a clone of cells that constitutes the tumor. The genotype of the original germline has changed in the clone of somatic cells but not in the germ cells. Retinoblastoma per se cannot be transmitted to offspring.

It has long been established that the *risk* of retinoblastoma is often transmitted as a simple dominant trait. In this case, one of the two required mutant alleles is transmitted in the germ line, and only one additional mutation is required in a retinoblast in order for the tumor to arise. This pattern has been established for a number of cancers for which high risk is found in certain families.

Somatic mutation must occur frequently at many loci. In most cases, there is no way to recognize an individual cell that has mutated. It is possible in the case of cancer because of the expansion of a clonal population from the original mutant cell. There are other instances in which somatic mutations also appear to be an essential part of the disease. For example, polycystic kidney disease is a dominantly inherited condition in which many renal cysts form. Comparison of the genotypes of the cysts with those of normal cells from the same person indicates that genetic alterations have occurred in the cysts and that each is a clone. As in cancers, clonal expansion makes it possible to detect the mutation that has occurred in a single somatic cell. Apparently two mutations are necessary, and, as in retinoblastoma, one can be transmitted through the germ line to all cells of the developing embryo. The study of somatic cell genetics is likely to be an area of great future interest as analytical techniques become ever more sensitive.

#### Additional Reading

Scriver C.R., Beaudet A.L., Sly W.S., Valle D., eds., *The Metabolic and Molecular Bases of Inherited Disease*, 7th ed., McGraw-Hill, New York, 1995.

Strachan T. and Read A.P., *Human Molecular Genetics*, 2nd ed., Wiley-Liss, New York, 1999.

On-Line Mendelian Inheritance in Man: <http://www3.ncbi.nlm.nih.gov/Omim/>

## Genetic Diversity

Genetic diversity is the biological diversity that is present in the genetic material or information. Genetic diversity may be represented at the DNA level, at the level, at the **chromosomal** level, at the level, and at other higher levels, such as gene interaction.

Heterozygosity is the general measure of the genetic variation per **locus** in a given population. Suppose that a population contains  $n$  **alleles** at a locus, designated 1, 2, 3, ...,  $n$  and with respective frequencies  $p_1, p_2, p_3, \dots, p_n$ . Designating the frequency of the  $i$ -th allele by  $p_i$ , the heterozygosity is then defined as  $H = 1 - \sum p_i^2$ . When we are interested in many loci, the average heterozygosity ( $H$ ) can be computed as an unweighted average of heterozygosities over various loci. Thus, the average heterozygosity is a simple and good measure for the genetic diversity of a population (1).

For DNA sequences, the average number of nucleotide substitutions between all possible pairs of alleles sampled in the population can be measured to evaluate the genetic diversity at the DNA level. This quantity is called the "nucleotide diversity" (1). These quantities are useful for measuring the genetic diversity within a population.

### Bibliography

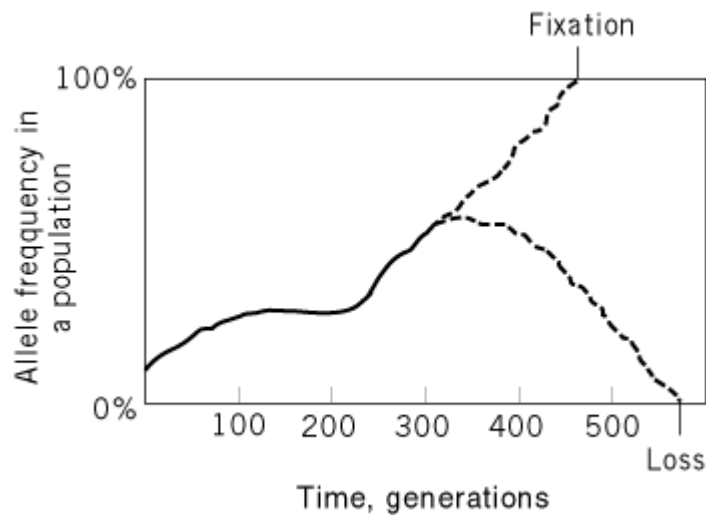
"Genetic Diversity" in , Vol. 2, p. 1000, by T. Gojobori; "Genetic Diversity" in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. M. Nei (1987). "Molecular Evolutionary Genetics". Columbia University Press, New York.

### Genetic Drift

For a given gene **locus**, the frequency of a particular **allele** in a **polymorphic** population provides useful information, in that it describes the genetic composition of a given population at a particular point of time during **evolution** (Fig. 1). In principle, the allele frequency will never remain the same; it will always change with time because of mutation, selection, and nonrandom mating. Even when random mating occurs in a population, the allele frequency will change in a random manner over time. The random change of an allele frequency over time is called "genetic drift". In general, the effect of genetic drift is greater the smaller the population size, and the effect of selection on allele frequencies dominates the effect of genetic drift when the population size is larger.

**Figure 1.** Scheme demonstrating genetic drift of an allele frequency in a population by chance over time.



Genetic drift is an important mechanism for gene evolution because it can change the genetic composition drastically over a very short period of time, particularly when the population size is extremely small. This was initially recognized by Sewall Wright (1931) and therefore known as the Wright effect (1). The Wright effect can be seen when a small number of individuals from the mother population immigrate to a new place (the [founder effect](#)) and when the population size suddenly reduces, for reasons such as climatic change or volcanic eruption (the “bottleneck effect”).

For identifying genes responsible for disease at the DNA level, it is necessary to assess [linkage disequilibrium](#) or genetic association with certain marker genes. For this practice, the effect of genetic drift should be evaluated correctly in order to discriminate between the effects of selection and of genetic drift.

#### Bibliography

1. S. Wright (1931). *Genetics* **16**, 97–159.

#### Genetic Marker

A genetic marker is a phenotypic variant due to a variation in the nucleotide sequence of a [genome](#), or the phenotypic manifestation of an **allele**. The term is used frequently as a synonym of allele.

Genetic markers may affect the morphology, behavior, physiology, or chemistry of the organism, and they are defined in relation to the phenotype of the standard **wild type**. A *marker gene* is a gene whose phenotypic variants are used as genetic markers.

From a practical viewpoint we can distinguish four groups of genetic markers:

1. *Selectable* markers which enable the individuals that carry them to grow under experimental conditions that hinder similar individuals lacking the marker. Selectable markers are nearly indispensable in the genetic analysis of very large laboratory populations of **viruses**, bacteria, unicellular eukaryotes, and cultured cells. Particularly useful are the selectable markers that are also *counterselectable*; that is, certain experimental conditions allow the growth of individuals

that lack them and hinder individuals that carry them.

2. Genetic markers that are easily recognized by visual inspection of the shape, the color, and sometimes the movements of cells and organisms. They have been particularly useful in the analysis of small populations of plants and animals.
3. Genetic markers whose detection requires specific tests. Conditional lethal markers impede growth under normal conditions, but are rescued by special circumstances of the environment or the genetic background. Chemical analyses are often indispensable—for example, to detect the presence or absence of a metabolite or to distinguish the **electrophoretic** variants of a [protein](#). Behavioral markers usually require their own specific tests. Onerous analyses can be avoided in a few cases. For example, some of the adenine auxotrophs of *Saccharomyces cerevisiae* develop a pink color under certain growth conditions; the presence of certain chemicals or certain [enzymes](#) in a cell may be revealed when they produce a colored product with certain reagents.
4. Variations in the DNA sequence that have no other phenotypic expression than the results of direct tests on the DNA (eg, [restriction maps](#), **PCR**). Their main advantage is that they occur all over the genome, even in regions in which no other genetic markers have been located.

## Genetic Suppression

Genetic suppression is the relief of a mutant phenotype by a second mutation. The second mutation gives the appearance of “suppressing” the first mutation. The suppressor mutation identifies a new site affecting the phenotype of interest, and suppression is strong evidence that the two genetic sites functionally interact. Suppressor mutations have been extremely useful for understanding genetic and biochemical phenomena. Biological systems are a tangle of interacting molecules and intersecting pathways. It should come as no surprise that mutations in one molecule are often modulated by secondary mutations in others.

Suppressor mutations can occur in the same gene as the initial mutation (intragenic suppression), or the two mutations may be in separate genes (intergenic suppression). Intragenic suppression may occur if the amino acid residues specified by the two sites must interact to give the normal function of the gene product. Intragenic suppression may also occur by the relief of a polar effect on gene expression. Intergenic suppression has been extensively used to characterize interactions between gene products. In certain cases, the suppressor mutation may establish a new pathway that supplants the pathway disrupted by the initial mutation. A special kind of suppression occurs when mutants are “rescued” by normal genes expressed from high-copy plasmids. The high concentration may facilitate an interaction between proteins that is weakened by the first mutation, or it may provide an alternative pathway to a normal phenotype. Another special kind of intergenic suppression occurs when mutations occur in transfer RNAs so that they translate stop (nonsense codons) codons as sense codons. These tRNAs (nonsense suppressors) can allow expression of genes that have suffered mutations that create stop codons within coding sequences. See “Nonsense suppression” and “tRNA suppressor” for more complete descriptions of these phenomena.

Suppressors were initially detected during studies of mutations. It was found that certain mutations express altered phenotypes only in certain genetic backgrounds. Apparently, in the non-expressing strains other genes were able to suppress these mutations. The molecular basis of these early observations has remained obscure (1), but they have been instructive for defining the principle of suppression.

In a now classic study, intragenic suppression was used to study the fundamental nature of the



genetic code. Francis Crick et al. (2) showed that genes are expressed as triplets that do not overlap one another. These workers first showed that mutations that were caused by the insertion (+) or deletion (–) of nucleotides could be suppressed by a second mutation of the opposite sign, if the two mutations occurred near one another within the gene. It was concluded that genes must be read in a specific “reading frame”. Either an insertion or a deletion would shift the reading frame such that the remainder of message sequence is decoded to specify the wrong polypeptide sequence. These workers further concluded that the second mutation of the opposite sign suppresses the first mutation by restoring the initial reading frame. In between the two mutations, the gene would be translated in a wrong reading frame, but as long as the amino acid sequence of that region was not critical for protein function then the double mutant would have a wild type phenotype. Crick also used suppressor mutations to determine the size of the “codon”. It was found that although one or two insertions invariably gave a mutant phenotype, three insertions or three deletions could give a wild type phenotype. It was concluded that because the insertion or deletion of three bases did not perturb a reading frame, the genetic code is a triplet code.

Intragenic suppression can also be used to show that two amino acid within a protein molecule interact. A mutation that eliminates a critical interaction between two amino acids can sometimes be suppressed by a change in the second amino acid that restores complementarity. The first example of such an analysis is presented by Helinsky and Yanofsky (3), in which an intragenic second-site mutation restored partial function to an *Escherichia coli trpA* mutant. Such information can complement structural studies by providing confirmation of hypothetical interactions between amino acids at specific positions. Such work is greatly facilitated by site-directed mutagenesis, which allows the creation of specific mutations to test structural predictions.

Intergenic suppression is extensively used to identify and study interactions between molecules. The principle is that mutations that destroy complementary interactions between proteins may be compensated by secondary mutations that provide suitable alternative interactions. This approach has been used to identify proteins interact within biochemical, gene expression, and signal transduction pathways. This approach is also widely used to study interactions between molecules already known or suspected to interact. Examples include studies of receptor-ligand interactions and of interactions between transcription factors and promoter elements.

Another commonly used method is the identification of cloned genes that can suppress mutations when expressed at high levels. Cells are transformed with expression libraries, and suppressed isolates are selected. Suppression may occur when the overexpressed protein compensates for the defect of the initial mutation. Several molecular mechanisms are possible, and further work is usually required to elucidate the mechanism of suppression in each instance. Possible mechanisms include the restoration of a direct interaction between proteins, increasing the expression of a partially-active mutant protein, and providing a bypass of the cellular or biochemical defect of the primary mutation. A related approach, though not usually thought of as suppression, is the use of the various “two-hybrid” systems to identify interacting proteins. These are actually highly-derived examples of intergenic suppression.

Another special but widely employed kind of intergenic suppression occurs when tRNAs are mutated such that they insert a suitable amino acid at missense, nonsense and frameshift mutations. These “informational” or “translational” suppressors have been extremely useful for studying the decoding mechanisms (4). Currently, certain translational suppressors are in development as tools for protein engineering. Systems for labeling proteins with novel amino acids at specific positions would be extremely useful. Toward this end, mutant tRNAs that decode the UAG stop (classically called the “amber” codon) are being aminoacylated with nonstandard amino acids. Then these tRNAs direct the incorporation of the novel amino acid into a nascent protein if its gene has been mutated to contain an amber codon at a specific site within its coding sequence. At this time, only a few nonstandard amino acids may be incorporated in this way, but it is anticipated that it will soon be possible to label proteins with a large variety of amino acids containing affinity tags, fluors, or reactive groups (4).

Discussion of related topics can be found in “Nonsense suppression,” “Suppressor mutation,” and “Suppressor tRNA.”.

### Bibliography

1. E. Kubli (1986) *Trends Genet.* **2**, 204–209.
2. F. H. C. Crick, L. Barnett, S. Brenner, and R. J. Watts-Tobin (1961) *Nature* **192**, 1227–1232.
3. D. R. Helinsky and C. Yanofsky (1963) *J. Biol. Chem.* **238**, 1043–1048.
4. E. J. Murgola (1994) in *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, DC, pp. 491–509.

### Genome

The *genome* is the totality of all of the **DNA** of an individual or a species. It is now taken to include all of the DNA, not just the **genes**. This definition is simple for the majority of genomes, but all of the minor DNA components must also be considered. For simple **viruses**, with a single **nucleic acid** molecule, the genome is obvious, although of course for RNA viruses it is **RNA** rather than DNA. For **haploid prokaryotes**, it is also straightforward, except for the **plasmids**, one copy of which is counted in the genome. Since prokaryotic strains and lineages differ by DNA sequence **polymorphism** and in the plasmids present, only genetic lineage genomes can be described with total accuracy. The genome of a species will be a compromise best representing all the known lineages, to make up a typical genome for the individuals of that species.

For **eukaryotes**, one haploid copy of the DNA of each of the diploid pairs of **chromosomes** (the **autosomes**) is included, plus one copy of the DNA of **sex** chromosomes. Thus the female and male genomes will differ if the sex chromosomes are different. One copy of the DNA from any **organelles** other than the **nucleus**, such as the **mitochondria** and **chloroplasts**, should be included; a complete list of such organelles is probably not yet available for any species, but many are normally unimportant. There is, of course, polymorphism of DNA sequences in a population of individuals, and the species' genome is again a compromise to obtain a typical genome. The genome can be considered to include all of the information required to specify the species, but only for certain viruses has it actually been demonstrated that the lineage can be propagated from purified DNA. When the phrase “genome size” was first used to describe the total haploid DNA content, there were strong objections, based on the concept that there is additional information in the chromosomal proteins, and even in body structures, that is necessary to propagate the lineage. While that point of view is not common at present, the proof has not been made with eukaryotes that the DNA contains all of the information required for maintenance of the lineage or species.

The majority of the DNA of a genome is not in the genes themselves and their known associated regulatory sequence. While the phenomenon of **gene regulation** is beginning to be understood, little is known of the significance of the majority of the DNA, whether it has any functions other than acting as spacer between genes. In most species, a large fraction of the DNA is repeated sequences (see **Repetitive DNA** and **Repeated DNA Sequence Interspersion**) that cause genetic **recombination** and **unequal crossing over**, resulting in genomic rearrangements, but their overall significance is not understood.

Suggestions for Further Reading

T. Cavalier-Smith, ed. (1985) *The Evolution of Genome Size*, John Wiley & Sons, Chichester U.K.,

B. John and G. L. G. Miklos (1988) *The Eukaryote Genome in Development and Evolution*, Allen & Unwin, London.

## Genomic Libraries

Early genomic libraries were viewed as a repository from which any genetic elements could be isolated for a given organism. Modern genomic libraries serve this function, but they also straddle the interface between genomic [bioinformatics](#) and functional biology. Thus, genomic sequences are increasingly understood in terms of their information content, comprising their RNA- and protein-coding sequences, as well as in terms of the overall organization of their embedded genetic program. Central to this information-based approach are new, highly sophisticated [libraries](#) of genomic DNA that maintain the structural context of genetic elements over several hundred kilobase pairs or more. Such libraries are essential for gene mapping projects, which will eventually produce an atlas of the genetic organization of many organisms. Detailed gene maps are required to interpret sequencing data that result from genome sequencing efforts. The genome sequencing projects themselves depend on high-quality genomic DNA libraries as a source of genetic material. The study of the genetic basis of disease makes use of sequence differences between individuals that can be ascertained through genomic libraries. These libraries also serve as a source of genetic probes that can be used for diagnosis and characterization of the relevant genotype of any individual.

Construction of genomic libraries requires particular attention during all stages of processing in order to minimize mechanical shearing of the cellular DNA. Nucleic acids become increasingly sensitive to shearing with increasing size, making it especially difficult to maintain the integrity of chromosome-sized species. Fortunately, specialized methods have been developed that enable the cloning of DNA fragments that are several hundred kilobase pairs or larger. The size of the genomic insert determines the type of vector into which the DNA is inserted. Inserts smaller than 10 kbp can be **cloned** into either plasmid or bacteriophage vectors. Cosmid vectors are required for fragments in the range of 10 to 40 kbp, while artificial chromosomes (ACs) are needed for even larger inserts.

Most genomic DNA library projects seek to represent a set of partially overlapping DNA fragments that collectively represent the entire [genome](#) of the subject organism. Such overlapping cloned inserts are ideal for chromosome “walking” approaches, in which the termini of one set of probes can be used to identify overlapping clones, which in turn can identify further overlapping clones downstream. In this manner, a set of clones can be derived that cover virtually any span of genetic sequences. Once prepared, such a repository facilitates the identification of genetic elements associated with any phenotype or disease. Three methods are employed for generating genomic DNA fragments of appropriate sizes for cloning. First, mechanical shearing can be used to produce relatively small inserts, although ligation into the vector with such fragments is typically inefficient. Because the shearing is random, this method produces a series of partially overlapping fragments. Second, and by far the most common method, is the use of [restriction enzymes](#) to prepare fragments. Digestion of DNA with restriction enzymes provides a great deal of control over the size of the resulting fragments. For example, endonucleases with a four-nucleotide recognition sequence cut on average every 256 nucleotides, while certain enzymes with large recognition sequences may cut approximately every  $10^6$  nucleotides. Unfortunately, complete digestion of genomic DNA with restriction enzymes eliminates any overlapping fragments; therefore, most genomic libraries employ partial digestion of genomic DNA. The degree of digestion thus influences the average fragment

size, as well as the degree of fragment overlap. Another advantage of using restriction endonucleases is that the fragment termini are defined and readily clonable. The third method of preparing genomic DNA is to perform **PCR** using primers containing a region of random sequence. In the first PCR cycle, the random primer segments anneal to both strands throughout the genome and are extended in a 5' to 3' direction. In subsequent PCR cycles, any two oppositely oriented primers that are within sufficient proximity can generate a specific PCR product. The resulting PCR products can be prepared for cloning into an appropriate vector. Alternatively, the collection of genomic DNA fragments can be maintained and propagated as PCR products using fixed sequences, introduced during the PCR, as primer binding sites for subsequent rounds of amplification.

In general, genomic libraries are used for DNA-based studies, such as gene mapping and marker analysis; however, in certain cases they may also be useful as [expression libraries](#). Many simpler organisms have few or no **introns**, so their genomic DNA can be directly expressed into [messenger RNA](#) and **translated** into protein. Furthermore, structural RNAs, such as **ribosomal RNA** and [transfer RNA](#) are not represented in most [cDNA libraries](#) and thus require genomic libraries to isolate appropriate DNA coding sequences. Finally, random genomic expression can be a productive means of generating RNA and [peptide libraries](#) for **epitope-** and **domain-**mapping exercises.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

## Germ Cell, Line

Germ cells are the cells in a multicellular organism that retain totipotency and are responsible for the propagation of succeeding generations of new individuals. Collectively, the germ cells comprise the germ line. Mature, fully differentiated germ cells in females are called [eggs](#), or oocytes, and in males [sperm](#), or spermatids. Undifferentiated germ cells present in animal embryos are called primordial germ cells, or PGCs. In dipterans such as *Drosophila*, PGCs are called pole cells. Flowering plants produce two types of germ cells, called microspores and megaspores. Microspores, or pollen grains, form within the stamen of the flower, and megaspores develop within the ovules of the carpels. A fundamental distinction between germ cells and other cells (somatic cells, or soma) is that, immediately prior to their final differentiation, germ cells undergo a specialized set of two cell divisions termed *meiosis*. Unlike *mitosis*, which is always coupled with [DNA replication](#), the first meiotic division is not in organisms with sexual reproduction. Thus, meiosis reduces the ploidy of germ cells by one-half, most commonly from **diploid** to [haploid](#). A new generation comes into being at [fertilization](#), when the two different types of mature germ cells (for instance, egg and sperm) fuse to form a [zygote](#). The egg and sperm nuclei also fuse to reconstitute the normal chromosome number. The zygote nucleus then divides mitotically, and it ultimately gives rise to all the cells of the organism.

### 1. Germ Line Specification

A common feature in animal [development](#) is that the germ line is segregated from the soma at a very early stage of embryogenesis. In many animals, germ cells develop in a region of the embryo whose cytoplasm contains specialized organelles, variously termed *P granules*, *polar granules*, or *germinal granules* (1). These organelles are assembled during oogenesis from maternally encoded gene products, and they localize asymmetrically within the egg either during oogenesis or immediately after fertilization. They are rich in RNA and protein, are not surrounded by membrane, and are associated with [ribosomes](#) and [mitochondria](#).

Transplantation experiments in *Rana pipiens* (2), *Drosophila melanogaster* (3, 4) and *Xenopus laevis* (5), and genetic experiments in *Drosophila* (6), have demonstrated that the necessary information for germ cell determination is present in the germinal cytoplasm. In *Caenorhabditis elegans*, P granules are distributed throughout the cytoplasm of the unfertilized egg. After fertilization, they become confined to the posterior end, such that they are found only in the posterior daughter of the first cleavage, called P<sub>1</sub> (7). In three subsequent mitoses, the P granules segregate to the more posterior daughter cell, called P<sub>2</sub>, P<sub>3</sub>, and P<sub>4</sub>. The P<sub>4</sub> cell is the [founder cell](#) of the entire germ line. In *Drosophila*, pole cells arise from the pole plasm at the posterior end of the preblastoderm embryo and are the first cells to form. Pole plasm is distinguished by the presence of polar granules.

Many protein and RNA components of *Drosophila* polar granules have now been identified (reviewed in refs. [8-10](#)), and their characterization broadly supports an early suggestion (11) that polar granules function by storing and regulating translation of maternal [messenger RNAs](#) required for germ cell determination. A protein component of polar granules in many organisms is the [RNA helicase](#) Vasa ([12-17](#)), and *vasa* activity is required in *Drosophila* for efficient [translation](#) of *oskar* and *nanos*, two pole plasm-localized mRNAs ([18, 19](#)). Conversely, translation of these RNAs is repressed in *Drosophila* outside of the pole plasm ([19-24](#)). Polar granules in many species are frequently associated with mitochondria, and in *Drosophila*, mitochondrial ribosomal RNAs are associated with polar granules and may be essential for pole cell formation ([25](#)). Germ cells also possess other densely staining cytoplasmic organelles, usually termed nuage ([26](#)), which are structurally related to polar granules, and which in *Drosophila* share at least one molecular component with polar granules ([12, 13, 27](#)).

In the two-cell *Xenopus* embryo, the germ plasm is present in small aggregates near the vegetal pole ([28](#)). By the time of the second cleavage division, these aggregates have moved toward the cleavage furrows, so that each of the cells of the four-cell embryo receives a portion of the germ plasm. In subsequent cleavage divisions, however, the germ plasm is associated with one of the spindle poles and therefore is segregated to only one of the two daughter cells. At gastrulation, the germ plasm becomes perinuclear and is divided symmetrically at mitosis, so the number of PGCs increases from four to about 14 by the stage at which they exit the endoderm. In [zebrafish](#), germ cells are not morphologically distinguishable until the 10- to 12-somite stage of development, when they are clustered on either side of the midline neighboring the third to fifth somite. However, analysis of the localization of zebrafish *vas* mRNA ([17](#)) suggests that, as in *Xenopus*, germ plasm is associated with the first two cleavage planes and condenses into four clumps, which do not further subdivide until after the 1000-cell stage. Microdissection studies in chick embryos suggest the germ line arises from the most central part of the blastodisc during the first 18 to 20 hours of development (the uterine stage), which includes cleavage and formation of the area pellucida ([29](#)).

Germ line specification in mammals may be quite different from that in other species, as mammalian oocytes lack a visible germ plasm ([30](#)). In mice, PGCs are first recognizable during gastrulation as a cluster of about 100 cells located within the midline extra-embryonic mesoderm posterior to the primitive streak ([31](#)). A key biological marker for totipotential embryonic cells is the [transcription factor](#) Oct-4 ([32](#)). Two tissue-specific enhancers have been mapped for the *Oct-4* gene, and [reporter gene](#) analysis indicates that the more distal one activates *Oct-4* expression in oocytes, morulae, and the inner cell mass of blastocyst-stage embryos, and then again in migratory and postmigratory germ cells. [Trans-acting](#) factors responsible for *Oct-4* activation have not yet been identified.

## 2. Germ Cell Migration

Once the germ cells have been specified, they follow an elaborate migration pathway within the interior of the embryo, culminating in an association with the gonadal mesoderm (genital ridge in vertebrates) to form the two embryonic gonads. At gastrulation, *Drosophila* pole cells are transferred

within the embryo in the posterior midgut invagination. Subsequently, they actively migrate through the midgut epithelium, move dorsally, separate into two bilateral populations, and contact the overlying mesoderm (33, 34). They are then surrounded by the gonadal mesoderm and, together with that mesoderm, form the two embryonic gonads, each containing approximately 14 pole cells. Pole cells do not divide throughout the period of migration. Pole cells lacking maternal *nanos* activity migrate into the interior of the embryo, but fail to be incorporated into gonads (35). The Wunen protein, a probable phosphatidic acid phosphatase expressed on the surface of posterior midgut cells, repels germ cells and thereby limits the range of their migration on the gut (36). Later steps in germ cell migration are influenced by a large number of genes expressed in the soma, some of which, including *tinman* and *zfh-1*, have also been implicated in gonad formation (37-39).

In fish, PGCs migrate dorsally during epiboly and early gastrulation and form two groups of cells on either side of the notochord, where they remain throughout somitogenesis. The precise position of the PGCs along the anterior-posterior axis varies among different species of fish. Later in larval development, the PGCs resume cell division and coalesce into a gonad. In *Xenopus*, the PGCs become incorporated into the larval hindgut after gastrulation, then migrate through the hindgut mesentery to reach the genital ridges on either side of the dorsal aorta (40). In the chick embryo, PGCs move to the endoderm to a site where blood vessels form, and they are taken up in the bloodstream (32). They then leave the circulatory system near the position where the hindgut develops and associate with the hindgut mesentery. In the mouse embryo, PGCs disperse from extra-embryonic tissues at 8.0 days past coitum (dpc) and migrate into the hindgut epithelium (41). From 9.5 to 11.5 dpc they move out of the hindgut, along the dorsal mesentery, and ultimately reach the gonadal primordia.

### 3. When Does the Germ Cell Lineage become fully distinct from the Somatic Lineage?

In the *C. elegans* embryo, where the complete cell lineage is known, the germ-line or P lineage becomes distinct when the P<sub>4</sub> cell is formed at the 24- to 28-cell stage of embryogenesis (42). In *Drosophila*, the germ cell lineage is distinct from the soma from the time of pole cell formation; when transplanted into host embryos, labeled pole cells do not contribute to somatic tissues (43). Conversely, in vertebrate embryos, migrating PGCs are not irreversibly determined as germ line cells. Transplantation of fluorescently labeled migrating PGCs from *Xenopus* embryos into unlabeled hosts indicates that they retain the capacity to differentiate into a large variety of somatic tissues (44). In mice, 8.5-dpc and 12.5-dpc (but not 15.5-dpc) PGCs can be cultured and give rise to undifferentiated cells called EG cells (45, 46). When EG cells are injected into host blastocysts, they frequently contribute to germ line and somatic tissues. Mammalian germ cells change the **methylation** pattern of parentally **imprinted** genes upon differentiation; for at least one such gene, *Igf2*, the methylation pattern in EG cells is similar to that in germ cells and not somatic cells (46).

### 4. How do Germ Cells Retain Totipotency?

Few details are presently understood as to how germ cells retain the capacity to generate all the different tissues of the organism, while somatic cells become increasingly specialized. It has been observed that early PGCs in *C. elegans* and *Drosophila* transcribe essentially no mRNA (47, 48). This transcriptional silencing is correlated in both organisms with a lack of a specific phosphorylated form of RNA polymerase II in premigratory and migratory germ cells (49). Furthermore, *Drosophila* pole cells that lack *nanos* activity, and fail to colonize the gonad, activate transcription of numerous reporter genes prematurely (35). It has been suggested that the lack of gene expression in embryonic germ line cells is of fundamental importance to maintaining their totipotency.

### Bibliography

1. H. W. Beams and R. G. Kessel (1974) *Int. Rev. Cytol.* **39**, 413–479.
2. L. D. Smith (1966) *Dev. Biol.* **14**, 330–347.
3. K. Illmensee and A. P. Mahowald (1974) *Proc. Natl. Acad. Sci. USA* **71**, 1016–1020.

4. M. Okada, I. A. Kleinman, and H. A. Schneiderman (1974) *Dev. Biol.* **37**, 43–54.
5. M. Wakahara (1978) *J. Exp. Zool.* **203**, 159–164.
6. A. Ephrussi and R. Lehmann (1992) *Nature* **358**, 387–392.
7. S. Strome and W. B. Wood (1983) *Cell* **35**, 15–25.
8. D. St Johnston (1993) in M. Bate and A. Martinez-Arias, eds. *The Development of Drosophila Melanogaster*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 325–363.
9. P. F. Lasko (1994) *Molecular genetics of Drosophila oogenesis*. R. G. Landes, Austin.
10. C. Rongo and R. Lehmann (1996) *Trends Genet.* **12**, 102–109.
11. A. P. Mahowald (1968) *J. Exp. Zool.* **167**, 237–262.
12. B. Hay, L. Ackerman, S. Barbel, L. Y. Jan, and Y. N. Jan (1988) *Development* **103**, 625–640.
13. L. Liang, W. Diehl-Jones, and P. F. Lasko (1994) *Development* **120**, 1201–1211.
14. M. E. Gruidl, P. A. Smith, K. A. Kuznicki, J. S. McCrone, J. Kirchner, D. L. Roussel, S. Strome, and K. L. Bennett (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13837–13842.
15. T. Komiya, K. Itoh, K. Ikenishi, and M. Furusawa (1994) *Dev. Biol.* **162**, 354–363.
16. Y. Fujiwara, T. Komiya, H. Kawabata, M. Sato, H. Fujimoto, M. Furusawa, and T. Noce (1994) *Proc. Natl. Acad. Sci. USA* **91**, 12258–12262.
17. C. Yoon, K. Kawakami, and N. Hopkins (1997) *Development* **124**, 3157–3165.
18. F.-H. Markussen, A. M. Michon, W. Breitwieser, and A. Ephrussi (1995) *Development* **121**, 3723–3732.
19. E. R. Gavis, L. Lunsford, S. E. Bergsten, and R. Lehmann (1996) *Development* **122**, 2791–2800.
20. J. Kim-Ha, K. Kerr, and P. M. Macdonald (1995) *Cell* **81**, 403–412.
21. P. J. Webster, L. Liang, C. A. Berg, P. Lasko, and P. M. Macdonald (1997) *Genes Dev.* **11**, 2510–2521.
22. A. Dahanukar and R. P. Wharton (1996) *Genes Dev.* **10**, 2610–2620.
23. C. A. Smibert, J. E. Wilson, K. Kerr, and P. M. Macdonald (1996) *Genes Dev.* **10**, 2600–2609.
24. R. P. Wharton, J. Sonoda, T. Lee, M. Patterson, and Y. Murata (1998) *Mol. Cell.* **1**, 863–872.
25. S. Kobayashi, R. Amikura, and M. Okada (1993) *Science* **260**, 1521–1523.
26. E. M. Eddy (1975) *Int. Rev. Cytol.* **43**, 229–280.
27. M. Wilsch-Bräuninger, H. Schwarz, and C. Nüsslein-Volhard (1997) *J. Cell Biol.* **139**, 817–829.
28. P. M. Whittington and K. E. Dixon (1975) *J. Embryol. Exp. Morph.* **33**, 57–74.
29. M. Ginsburg (1994) in J. Marsh and J. Goode, eds., *Germline Development*, Ciba Foundation Symposium 182, Wiley, Chichester: pp. 61–83.
30. E. M. Eddy, J. M. Clark, D. Gong, and B. A. Fenderson (1981) *Gamete Res.* **4**, 333–362.
31. M. Ginsburg, M. H. L. Snow, and A. McLaren (1990) *Development* **110**, 521–528.
32. Y. I. Yeom, G. Fuhrmann, C. E. Ovitt, A. Brehm, K. Ohbo, M. Gross, K. Hübner, and H. R. Schöler (1996) *Development* **122**, 881–894.
33. G. Callaini, M. G. Riparbelli, and R. Dallai (1995) *Dev. Biol.* **170**, 365–375.
34. M. K. Jaglarz and K. R. Howard (1995) *Development* **121**, 3495–3504.
35. S. Kobayashi, M. Yamada, M. Asaoka, and T. Kitamura (1996) *Nature* **380**, 708–711.
36. N. Zhang, J. Zhang, K. J. Purcell, Y. Cheng, and K. Howard (1997) *Nature* **385**, 64–67.
37. M. Boyle, N. Bonini, and S. DiNardo (1997) *Development* **124**, 971–982.
38. H. T. Broihier, L. A. Moore, M. Van Doren, S. Newman, and R. Lehmann (1998) *Development* **125**, 655–666.
39. L. A. Moore, H. T. Broihier, M. Van Doren, L. B. Lunsford, and R. Lehmann (1998)

Development **125**, 667–678.

40. J. Heasman and C. C. Wylie (1981) *Proc. R. Soc. Lond. B Biol. Sci.* **213**, 41–58.
41. M. Gomperts, C. Wylie, and J. Heasman (1994) *Germline Development*, Ciba Foundation Symposium 182, Wiley Chichester, pp. 121–139.
42. J. E. Sulston, E. Scheirenborg, J. G. White, and J. N. Thomson (1983) *Dev. Biol.* **100**, 64–119.
43. E. M. Underwood, J. H. Caulton, C. D. Allis, and A. P. Mahowald (1980) *Dev. Biol.* **77**, 303–314.
44. C. C. Wylie, J. Heasman, A. Snape, M. O'Driscoll, and S. Holwill (1985) *Dev. Biol.* **112**, 66–72.
45. J. L. Resnick, L. S. Bixler, L. Cheng, and P. J. Donovan (1992) *Nature* **359**, 550–551.
46. P. A. Labosky, D. P. Barlow, and B. L. M. Hogan (1994) *Development* **120**, 3197–3204.
47. M. Zalokar (1976) *Dev. Biol.* **49**, 425–437.
48. G. Seydoux, C. C. Mello, J. Pettitt, W. B. Wood, J. R. Priess, and A. Fire (1996) *Nature* **382**, 713–716.
49. G. Seydoux and M. A. Dunn (1997) *Development* **124**, 2191–2201.

### **Suggestions for Further Reading**

50. H. W. Beams and R. G. Kessel (1974) “The Problem of Germ Cell Determinants,” *Int. Rev. Cytol.* **39**, 413–479. The most complete description of the morphology of germ cell specification in a large variety of organisms.
51. J. Marsh and J. Goode, eds. (1994) *Germline Development*, Ciba Foundation Symposium, vol. 182, Wiley, Chichester, England. An excellent compendium of reviews and discussion concerning germ cell development, in many different organisms.
52. D. St Johnston (1995) “The Intracellular Localization of Messenger RNAs,” *Cell* **81**, 161–170.
53. G. Wei and A. P. Mahowald (1994) “The Germline: Familiar and Newly Uncovered Properties,” *Annu. Rev. Genet.* **28**, 309–324. This review includes a good discussion of genetic imprinting and culturing of mammalian germ cells. Very inclusive reference list.
54. A. Williamson and R. Lehmann (1996) “Germ Cell Development in *Drosophila*,” *Annu. Rev. Cell Dev. Biol.* **12**, 365–391. A detailed review of the determination, formation, migration, and differentiation of germ cells in *Drosophila*.



# Gibberellins

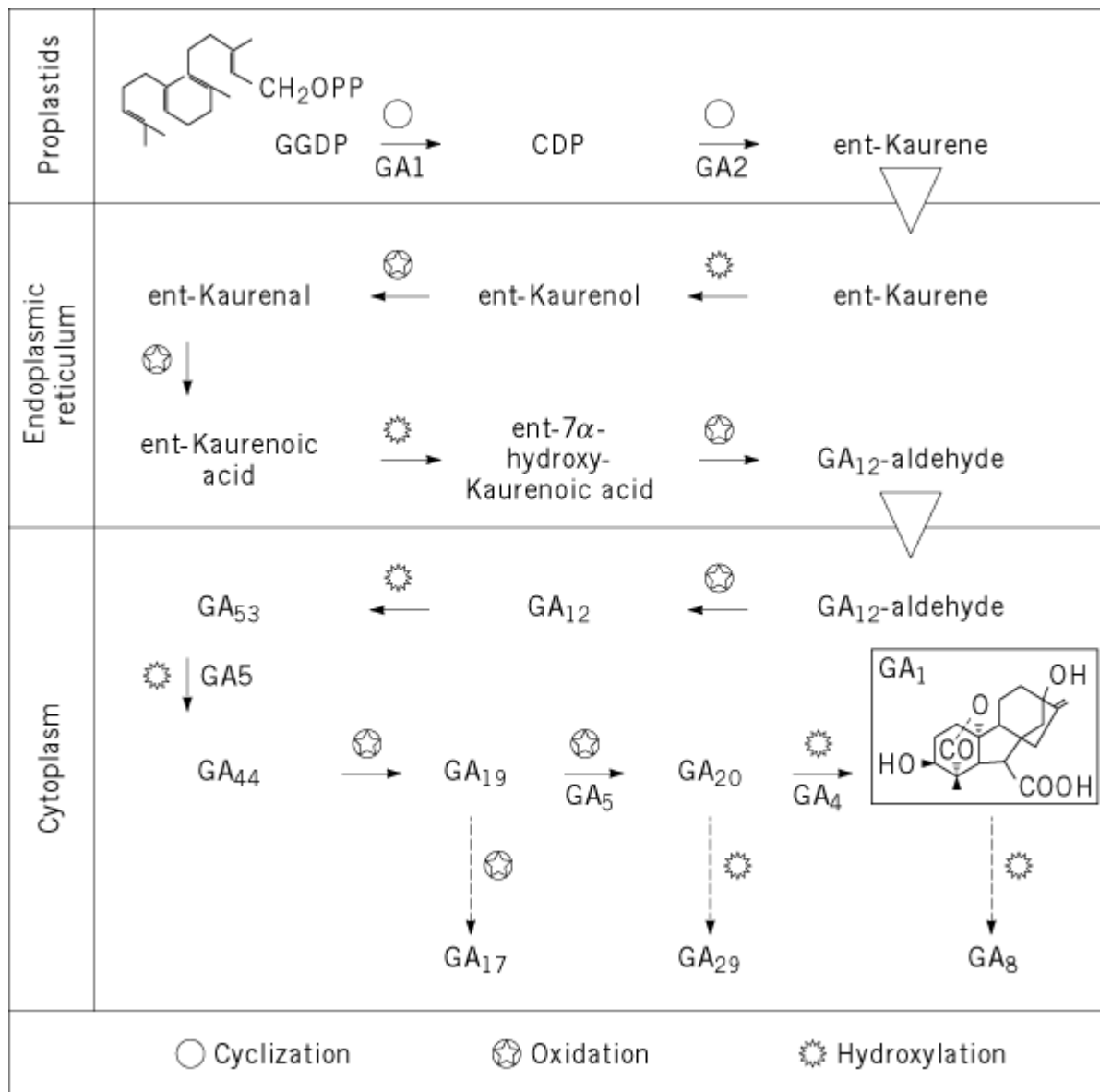
## 1. History

Gibberellins are [plant hormones](#), tetracyclic diterpene carboxylic acids that were discovered in 1926 as a phytotoxin produced by the fungus *Gibberella fujikuroi*. The latter caused a pathological longitudinal growth in rice called “foolish seedling disease” ([1](#)). The active compound was isolated from the fungus in 1930 by Yabuta and Sumiki ([2](#)) and was called *gibberellin*. In 1958, a first gibberellin (GA<sub>1</sub>) was purified from runner bean seeds (*Phaseolus coccineus*) by MacMillan and coworkers ([3](#)). Since this discovery, 112 gibberellins have been identified ([4](#)).

## 2. Biosynthesis and Metabolism

Gibberellins are isoprenoid compounds synthesized from C<sub>2</sub> units by the mevalonic acid pathway. Geranyl geranyl diphosphate, a C<sub>20</sub> molecule, serves as a donor for the entire carbon skeleton of gibberellins. Starting from geranyl geranyl diphosphate, biosynthesis of gibberellins can be divided into three stages according to the nature of the enzymes involved and their subcellular localization (Fig. [1](#)). The pathway is well-characterized, owing to the existence of gibberellin-deficient [mutants](#) that are easily identified by their dwarfed stature ([5](#), [6](#)). Gibberellin-deficient mutants are blocked at a particular step in the pathway, and their **phenotype** can be reverted by external addition of active gibberellins.

**Figure 1.** Gibberellin biosynthetic pathway in higher plants. Steps involving GA<sub>1</sub>, GA<sub>2</sub>, GA<sub>4</sub>, and GA<sub>5</sub> in *Arabidopsis* are indicated. (Adapted from Ref. [16](#))



Several **genes** encoding gibberellin-biosynthetic [enzymes](#) have been characterized, and their products seem to catalyze multiple steps. First, geranyl geranyl diphosphate is converted into *ent*-kaurene by two cyclization steps. The enzymes catalyzing these reactions are *ent*-kaurene synthase A and B, also termed *ent*-copalyl diphosphate synthase (CPS) and *ent*-kaurene synthase (KS) (7). The gene encoding the CPS has been cloned from *Arabidopsis* (*GAI*) (8), maize (*Anther ear-1*, *An1*) (9), pea (*LS*) (10), tomato (11), and pumpkin (12), whereas that for KS was cloned from pumpkin endosperm (13) and from *Arabidopsis* (*GA2*) (14).

The KS activity is mostly associated with stroma of proplastids in rapidly dividing plant tissues (15); however, *GAI* promoter activity is also present in fully expanded leaves (16), suggesting that gibberellins are transported to responsive tissues. In *Arabidopsis Arabidopsis*, *GAI* gene expression appears to be highly regulated during growth and development. In addition, *GAI* promoter activity was restricted to specific cell types, such as shoot apices, root tips, developing seeds and flowers, and vascular tissue of mature leaves (17).

The second stage converts *ent*-kaurene to GA<sub>12</sub>-aldehyde. *Ent*-kaurene is subjected to five sequential oxidation steps, resulting in GA<sub>12</sub>-aldehyde. These reactions take place in the endoplasmic reticulum

and seem to be catalyzed by a single membrane-associated [cytochrome P450](#) monooxygenase. The gene encoding this enzyme was recently cloned from maize (*Dwarf-3, D3*) ([18](#)). Its expression was observed in seedlings, developing leaves, vegetative meristems, and roots. The predicted protein carries the Fe-binding domain that is characteristic of cytochrome P450 proteins ([18](#)).

The third stage converts GA<sub>12</sub>-aldehyde to various gibberellins and takes place in the cytoplasm. The first step in this process consists of the oxidation of GA<sub>12</sub>-aldehyde to GA<sub>12</sub>. Subsequently, GA<sub>12</sub> is subjected to various reactions, depending on the organism. In the past, genes encoding one of the major enzymes involved in these reactions, GA 20-oxidase, were characterized in different species such as pumpkin ([19](#)), *Arabidopsis* (*GA5*) ([20](#)), and spinach ([21](#)). GA 20-oxidases belong to the family of 2-oxoglutarate-dependent dioxygenases. In *Arabidopsis*, GA 20-oxidases are encoded by a small [multigene family](#). Expression of the *GA5* gene in leaves was enhanced when plants were shifted from short-day to long-day conditions, indicating photoperiod regulation. In addition, *GA5* was down-regulated by GA<sub>4</sub> treatment, suggesting feedback repression of the pathway ([20](#)). Finally, a 3b-hydroxylase catalyzes the formation of biologically active gibberellins. The *GA4* locus encoding a 3b-hydroxylase was recently cloned from *Arabidopsis* ([22](#)). The gene is ubiquitously expressed, but its [messenger RNA](#) is most abundant in siliques. As with *GA5*, *GA4* is also subjected to a feedback regulatory mechanism ([23](#)).

Relatively few gibberellins appear to be biologically active, such as GA<sub>1</sub>, GA<sub>3</sub>, and GA<sub>4</sub>. The remaining ones are the precursors or the deactivated forms of active gibberellins. The latter are generally deactivated by 2b-hydroxylation. The only deactivation mutant described thus far is the slender (*sln*) genotype of *Pisum sativum*. A mutation at this locus is responsible for a large accumulation of GA<sub>20</sub>, the precursor of GA<sub>1</sub>, and gives rise to the elongated phenotype. It was suggested that *sln* encodes a regulatory protein controlling two genes involved in the degradation pathway ([24](#)). Another process in gibberellin metabolism is their conjugation to inactive derivatives, including glucosyl ethers and glucose esters, that are found mainly in seeds ([25](#)).

### 3. Signal Perception and Transduction

The gibberellin biosynthesis pathway is well-characterized, but much less is known about how gibberellins are perceived by the plant and how the signal is transduced to control the expression of gibberellin-regulated genes ([26, 27](#)). Biochemical studies on barley aleurone protoplasts suggested that gibberellin binds to a receptor located at the external face of the plasma membrane ([28](#)). This receptor has not yet been isolated, however, nor has any gibberellin receptor been cloned to date.

Our current knowledge of gibberellin signaling is based mainly on the characterization of four mutants of *Arabidopsis*: *spy* (*spindly*), *gai* (*gibberellin-insensitive*), *rga* (*repressor of gal-3*), and *pkl* (*pickle*). None of these mutants can be entirely reverted by exogenous application of gibberellin, supporting their role in gibberellin signaling.

The *spy* mutants were screened for their ability to germinate in the presence of the gibberellin biosynthesis inhibitor paclobutrazol, which blocks the cytochrome P450 monooxygenase involved in the oxidation of *ent*-kaurene to *ent*-kaurenoic acid ([29](#)). *Spy* mutants phenocopy wild-type plants that were repeatedly treated with GA<sub>3</sub>, presenting a constitutive gibberellin-response phenotype that is characterized by long hypocotyls, pale green foliage, increased stem elongation, early flowering, partial male sterility, and parthenocarpic fruit development. Moreover, *spy* plants remain gibberellin-responsive. *Spy* is partially epistatic to *gal-2*. Collectively, these data indicate that SPY is a negative regulator of at least one branch in gibberellin signaling.

The *SPY* gene was recently cloned using a T-DNA tagging approach ([30](#)). The predicted SPY protein contains a tetratricopeptide repeat of 34 amino acid residues in its *N*-terminal region. This motif is

also found in some other eukaryotic and prokaryotic proteins and is proposed to form an amphipathic [a-helix](#) that mediates [protein-protein interactions](#) (31). The tetratricopeptide repeat appears to be important for the normal function of SPY, because some of the mutant phenotypes result from a deletion in this motif. The C-terminal region of SPY shows sequence similarity to the mammalian Ser/Thr *O*-linked *N*-acetylglucosamine (*O*-GlcNAc) transferase, which plays an important role in the regulation of the activity of various nuclear and cytosolic proteins, either directly or by inhibiting their phosphorylation (32, 33). Because certain *spy* alleles are affected in the region bearing similarity to *O*-GlcNAc transferase, this activity seems to be involved in proper functioning of the *SPY* gene product.

A second player is GAI. *Gai* mutant plants present a dwarf phenotype (reduced plant height, decrease in apical dominance, and limited seed germination), comparable to gibberellin-deficient plants. The lack of response of *gai* mutants to exogenous application of gibberellin, however, suggests that the GAI protein acts in signal transduction (34, 35). *Gai* heterozygotes show an intermediate mutant phenotype, indicating that the mutation is semidominant.

The *GAI* gene was recently isolated using *Ds* [transposon](#) tagging (36). In addition, a closely related gene named *GRS* (*GAI*-related sequence) was cloned (36). Comparison of the predicted amino acid sequences of the *GRS* and *GAI* proteins with the [sequence databases](#) showed that they are members of the VHIID family of plant [transcription factors](#) defined by Di Laurenzio et al. (37). A first member of this novel class of transcription factors was called *SCARECROW* (*SCR*) and controls cell fate in *Arabidopsis* roots. The VHIID proteins are putative transcription factors that contain three main motifs: a central –Val–His–Ile–Ile–Asp–sequence, with a Leu–X–X–Leu–Leu–sequence (where X represents any amino acid residue) in its immediate vicinity, and an –Arg–Val–Glu–Arg–sequence in the C-terminal region (38). The Leu–X–X–Leu–Leu–motif is a signature motif in transcriptional coactivators that mediates binding to nuclear receptors (39). The *GAI* protein possesses a fourth sequence in its N-terminal region (–Asp–Glu–Leu–Leu–Ala–) that appears important for its normal function. A nearby Ser/Thr-rich stretch is a potential target of SPY action. *GAI* was proposed to act directly as a transcriptional repressor of target genes that promote gibberellin-mediated processes, or indirectly as a transcriptional activator that induces the expression of such a repressor (36). This hypothesis implies that gibberellin would modulate the pathway by derepression and, therefore, by inactivation of *GAI*. An extragenic suppressor mutant of *gai*, called *gar2* (*gai* suppressor 2), was genetically characterized (35).

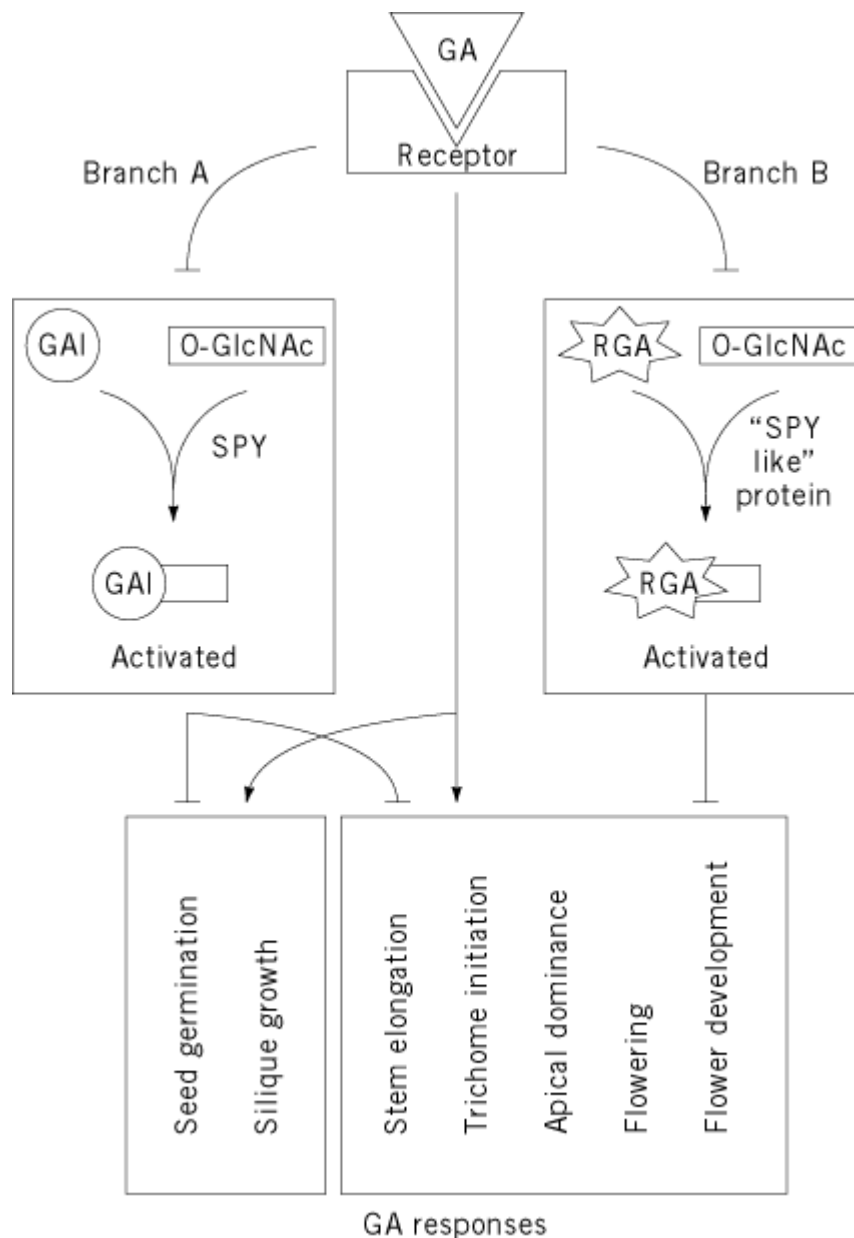
A third class of gibberellin-signaling mutants is represented by *rga* (*repressor of gal-3*) (40). *Rga* mutants were found by screening for the ability to suppress the phenotype of the *gal-3* mutant that is impaired in the first step of the gibberellin biosynthesis pathway (41). It is important to note, however, that the *rga* mutation does not suppress the gibberellin deficiency entirely, because it has no effect on seed germination or on silique growth. All *rga* alleles are recessive, indicating that *RGA* most probably encodes a negative regulator of gibberellin signaling.

The *RGA* gene was cloned using genomic subtraction and shown to be identical to *GRS* (38). The predicted protein is a member of the VHIID family of plant transcription factors described above. Although primarily repressing the gibberellin response, *RGA* may also play a role in the control of gibberellin biosynthesis (38). The *RGA* protein contains the –Asp–Glu–Leu–Leu–Ala–sequence and the Ser/Thr-rich region, both specific to proteins controlling the gibberellin response.

A final regulator, though difficult to place in the pathway at this point, is *pickle* (*pkl*). The *pkl* mutants present a phenotype reminiscent of gibberellin-deficient or gibberellin response mutants and exhibit abnormal root development (42). The primary root meristem of these mutants retained characteristics of embryonic tissues. This phenotype was partially suppressed by addition of gibberellin. Because the *gai pkl* double mutant exhibited a strong synergistic gibberellin-deficient phenotype, it was concluded that *pkl* is defective in the gibberellin-response pathway. The gene product *PKL* is proposed to be a positive regulator of gibberellin signaling (43).

Figure 2 presents a model for gibberellin signal transduction based on currently available genetic and biochemical data (38, 40, 43). Gibberellin interacts with its receptor(s), thereby activating gibberellin responses as a result of inhibition of two negative regulators (branches A and B). In branch A, SPY is acting downstream of (or modulates) GAI, because the *spy* mutant is completely epistatic to the *gai* mutation. SPY presumably activates the GAI repressor function via glycosylation. The RGA protein exerts the same function as GAI and represses the transcription of genes involved in gibberellin-regulated growth and development, or it promotes the expression of such a repressor. The repressor function of RGA is activated by SPY or a putative SPY-like protein, unknown to date. This assumption is based on the fact that the *rga* and *spy* mutations have an additive effect for the common processes controlled by gibberellins, thus indicating that SPY probably does not regulate RGA. Both O-glycosyltransferases would repress the common gibberellin response (stem elongation, trichome initiation, apical dominance, and flower development). In addition, the “GAI/SPY” branch blocks seed germination and silique growth.

**Figure 2.** Gibberellin signal transduction. The active gibberellin interacts with a membrane-associated or cytosolic receptor and activates the normal gibberellin response by blocking two repressors, named GAI (branch A) and RGA (branch B). The SPY and/or SPY-like proteins regulate the GAI and RGA proteins via glycosylation. Arrows and bars denote activation and repression, respectively. (Modified from Refs. 38, 40, and 41.)



Due to their nature, SPY/RGA/GAI would act fairly late in the pathway to induce gibberellin-regulated responses. It should be mentioned that recent evidence also suggests the involvement of heterotrimeric **G proteins** and [cyclic GMP](#) in gibberellin signaling in barley aleurone cells ([44](#), [45](#)). The former intermediates would play a role in the early phase of signal transfer.

#### 4. Downstream Targets

Gibberellin control of plant growth and development is achieved by modulating the expression of specific target genes. One of the best-characterized roles of gibberellins consists in the stimulation of production of hydrolytic enzymes by the aleurone cells. These hydrolytic enzymes, such as  $\alpha$ -amylase,  $\alpha$ -glucosidase, and a [thiol proteinase](#), are responsible for the breakdown of starch and proteins in the endosperm ([46-48](#)). The products of these hydrolyses are absorbed by the scutellum and used by the growing embryo. *De novo* synthesis of hydrolytic enzymes induced by gibberellin is often due to the presence of a gibberellin response element (GARE) in the promoter region of the corresponding genes ([49](#)). The identification of the GARE element has led to the isolation of a myb-like *trans*-acting factor that can bind to this region, which is closely related to c-Myb and v-Myb consensus sequences ([50](#)). Furthermore, gibberellins can stimulate stem elongation by increasing the xyloglucan endotransglycosylase activity, an enzyme that cleaves the rigid structure of the cell wall and is associated with growing tissues ([51](#), [52](#)). In addition, gibberellins cause microtubule reorientation favoring axial elongation ([53](#)). The tonoplast-intrinsic protein g-TIP, which maintains turgor pressure, thereby favoring wall extensibility, is also induced by gibberellin application ([54](#)). In *Petunia* flower development, gibberellin induces the expression of genes such as those of chalcone synthase, chalcone isomerase, anthocyanidin synthase, and dihydroflavonol 4-reductase, which are responsible collectively for corolla pigmentation ([55-59](#)).

#### 5. Effects

Biologically active gibberellins are involved in key processes of plant growth and development. The best-characterized roles of gibberellins consist in the control of reserve mobilization in cereals following germination and in promotion of shoot elongation through effects on both cell division and expansion. Gibberellins are also involved in the regulation of flower and fruit development and are required for germination and seedling growth in several species ([60](#)).

#### Bibliography

1. N. Takahashi, B. O. Phinney, and J. MacMillan (1991) *Gibberellins*, Springer-Verlag, New York.
2. T. Yabuta and Y. Sumiki (1938) *J. Agric. Chem. Soc. Japan* **14**, 1526.
3. J. MacMillan and P. J. Suter (1958) *Naturwissenschaften* **45**, 46.
4. T. Hisamatsu, M. Koshioka, S. Kubota, T. Nishijima, R. W. King, L. N. Mander, and D. J. Owen (1998) *Phytochemistry* **47**, 3–6.
5. R. Hooley (1994) *Plant Mol. Biol.* **26**, 1529–1555.
6. P. Hedden and Y. Kamiya (1997) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **48**, 431–460.
7. J. MacMillan (1997) *Nat. Prod. Rep.* **14**, 221–243.
8. T.-p. Sun, H. M. Goodman, and F. M. Ausubel (1992) *Plant Cell* **4**, 119–128.
9. R. J. Bensen, G. S. Johal, V. C. Crane, J. T. Tossberg, P. S. Schnable, R. B. Meeley, and S. P. Briggs (1995) *Plant Cell* **7**, 75–84.
10. T. Ait-Ali, S. M. Swain, J. B. Reid, T.-p. Sun, and Y. Kamiya (1997) *Plant J.* **11**, 443–454.
11. R. Imai, Y. Y. Yang, T. Ait-Ali, H. Kawaide, and Y. Kamaya (1996) *Plant Cell Physiol.* **37**

(Suppl.), S143 (#539).

12. M. W. Smith, S. Yamaguchi, T. Ait-Ali, and Y. Kamiya (1997) *Plant Cell Physiol.* **38** (Suppl.), S115 (#413).
13. S. Yamaguchi, T. Saito, H. Abe, H. Yamane, N. Murofushi, and Y. Kamiya (1996) *Plant J.* **10**, 203–213.
14. S. Yamaguchi, T.-p. Sun, H. Kawaide, and Y. Kamiya (1998) *Plant Physiol.* **116**, 1271–1278.
15. H. Aach, H. Bode, D. G. Robinson, and J. E. Graebe (1997) *Planta* **202**, 211–219.
16. T.-p. Sun and Y. Kamiya (1997) *Physiol. Plant.* **101**, 701–708.
17. A. L. Silverstone, C.-w. Chang, E. Krol, and T.-p. Sun (1997) *Plant J.* **12**, 9–19.
18. R. G. Winkler and T. Helentjaris (1995) *Plant Cell* **7**, 1307–1317.
19. T. Lange, P. Hedden, and J. E. Graebe (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8552–8556.
20. Y.-L. Xu, L. Li, K. Wu, A. J. M. Peeters, D. A. Gage, and J. A. D. Zeevaart (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6640–6644.
21. K. Wu, L. Li, D. A. Gage, and J. A. D. Zeevaart (1996) *Plant Physiol.* **110**, 547–554.
22. H.-H. Chiang, I. Hwang, and H. M. Goodman (1995) *Plant Cell* **7**, 195–201.
23. J. B. Reid, J. J. Ross, and S. M. Swain (1992) *Planta* **188**, 462–467.
24. J. J. Ross, J. B. Reid, S. M. Swain, O. Hasan, A. T. Poole, P. Hedden, and C. L. Willis (1995) *Plant J.* **7**, 513–523.
25. F. B. Salisbury and C. W. Ross (1992) In *Plant Physiology*, 4th ed. (F. B. Salisbury and C. W. Ross, eds.), Wadsworth, Belmont, CA, pp. 357–381.
26. S. M. Swain and N. E. Olszewski (1996) *Plant Physiol.* **112**, 11–17.
27. J. J. Ross, I. C. Murfet, and J. B. Reid (1997) *Physiol. Plant.* **100**, 550–560.
28. S. Gilroy and R. L. Jones (1994) *Plant Physiol.* **104**, 1185–1192.
29. S. E. Jacobsen and N. E. Olszewski (1993) *Plant Cell* **5**, 887–896.
30. S. E. Jacobsen, K. A. Binkowski, and N. E. Olszewski (1996) *Proc. Natl. Acad. Sci. USA* **93**, 9292–9296.
31. J. R. Lamb, S. Tugendreich, and P. Hieter (1995) *Trends Biochem. Sci.* **20**, 257–259.
32. L. K. Kreppel, M. A. Blomberg, and G. W. Hart (1997) *J. Biol. Chem.* **272**, 9308–9315.
33. W. A. Lubas, D. W. Frank, M. Krause, and J. A. Hanover (1997) *J. Biol. Chem.* **272**, 9316–9324.
34. M. Koornneef, A. Elgersma, C. J. Hanhart, E. P. van Loenen-Martinet, L. van Rijn, and J. A. D. Zeevaart (1985) *Physiol. Plant.* **65**, 33–39.
35. R. N. Wilson and C. R. Somerville (1995) *Plant Physiol.* **108**, 495–502.
36. J. Peng, P. Carol, D. E. Richards, K. E. King, R. J. Cowling, G. P. Murphy, and N. P. Harberd (1997) *Genes Dev.* **11**, 3194–3205.
37. L. Di Lorenzo, J. Wysocka-Diller, J. E. Malamy, L. Pysh, Y. Helariutta, G. Freshour, M. G. Hahn, K. A. Feldmann, and P. N. Benfey (1996) *Cell* **86**, 423–433.
38. A. L. Silverstone, C. N. Ciampaglio, and T.-p. Sun (1998) *Plant Cell* **10**, 155–169.
39. D. M. Heery, E. Kalkhoven, S. Hoare, and M. G. Parker (1997) *Nature* **387**, 733–736.
40. A. L. Silverstone, P. Y. A. Mak, E. C. Martínez, and T.-p. Sun (1997) *Genetics* **146**, 1087–1099.
41. T.-p. Sun and Y. Kamiya (1994) *Plant Cell* **6**, 1509–1518.
42. J. Ogas, J.-C. Cheng, Z. R. Sung, and C. Somerville (1997) *Science* **277**, 91–94.
43. J. Ogas (1998) *Curr. Biol.* **8**, R165–R167.
44. S. P. Penson, R. C. Schuurink, A. Fath, F. Gubler, J. V. Jacobsen, and R. L. Jones (1996) *Plant Cell* **8**, 2325–2333.
45. H. D. Jones, S. J. Smith, R. Desikan, S. Plakidou-Dymock, A. Lovegrove, and R. Hooley (1998) *Plant Cell* **10**, 245–253.

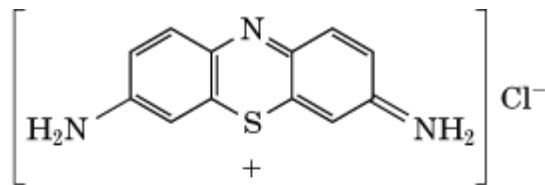
46. S. M. Koehler and T.-H. D. Ho (1990) *Plant Cell* **2**, 769–783.
47. A. K. Huttly and A. L. Phillips (1995) *Physiol. Plant.* **95**, 310–317.
48. B. K. Tibbot and R. W. Skadsen (1996) *Plant Mol. Biol.* **30**, 229–241.
49. A. K. Huttly, A. L. Phillips, and J. W. Tregear (1992) *Plant Mol. Biol.* **19**, 903–911.
50. F. Gubler, R. Kalla, J. K. Roberts, and J. V. Jacobsen (1995) *Plant Cell* **7**, 1879–1891.
51. M. C. Smith, P. R. Matthews, P. H. D. Schünmann, and P. M. Chandler (1996) *J. Exp. Bot.* **47**, 1395–1404.
52. P. H. D. Schünmann, R. C. Smith, V. Lång, P.R. Matthews, and P. M. Chandler (1997) *Plant Cell Environ.* **20**, 1439–1450.
53. H. Shibaoka (1994) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **45**, 527–544.
54. A. L. Phillips and A. K. Huttly (1994) *Plant Mol. Biol.* **24**, 603–615.
55. D. Weiss and A. H. Halevy (1989) *Planta* **179**, 89–96.
56. D. Weiss, A. J. van Tunen, A. H. Halevy, J. N. M. Mol, and A. G. M. Gerats (1990) *Plant Physiol.* **94**, 511–515.
57. D. Weiss, R. van Blokland, J. M. Kooter, J. N. M. Mol, and A. J. van Tunen (1992) *Plant Physiol.* **98**, 191–197.
58. D. Weiss, A. H. van der Luit, J. T. M. Kroon, J. N. M. Mol, and J. M. Kooter (1993) *Plant Mol. Biol.* **22**, 893–897.
59. G. Ben-Nissan and D. Weiss (1996) *Plant Mol. Biol.* **32**, 1067–1074.
60. S. M. Swain, J. B. Reid, and Y. Kamiya (1997) *Plant J.* **12**, 1329–1338.
61. H. Kende and J. A. D. Zeevaart (1997) *Plant Cell* **9**, 1197–1210.

## Giemsa Binding

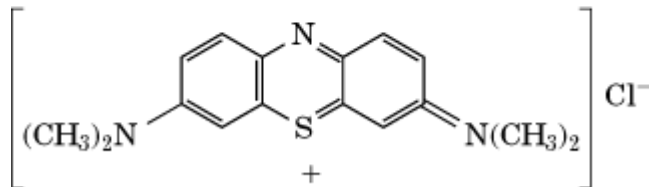
Giemsa is one of the most useful dyes for staining [chromosomes](#) in the banding patterns characteristic of different types of [heterochromatin](#) and [euchromatin](#) (see [C-Banding](#) and [G Banding](#)). It is a mixture of organic dyes including thionine, or Lauth's violet, and methylated derivatives of thionine (see [Figure 1](#)) including azure B, azure A, azure C, and methylene blue. Thionine itself stains chromosomes uniformly. The introduction of methyl groups, however, changes the properties of the molecule so that it differentially stains [chromatin](#). The thiazine ring in the Giemsa dyes interacts with the phosphodiester backbone of **DNA** and side-stacks along the **double helix**. How this might occur in a chromatin context is not yet known. Variations in chromosomal preparations, in which the histones are removed and other chromosomal proteins digested to different extents, enhance the banding patterns obtained with Giemsa dye (see [G Banding](#)).

**Figure 1.** The chemical structures of thionine and methylene blue.





**Thionine**



**Methylene blue**

### Suggestion for Further Reading

1. R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A Synthesis*, Wiley-Liss, New York.

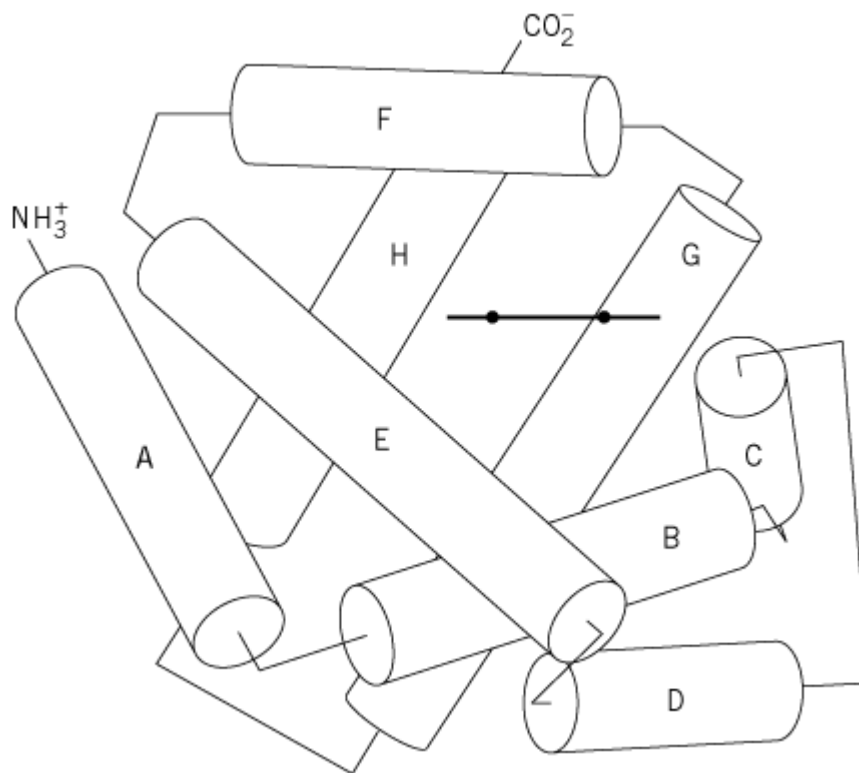
## Globins

### 1. Definition, structure, and properties

Globins are members of a family of [oxygen-binding proteins](#), characterized by **homologous** amino acid sequences and a characteristic [protein structure](#) known as the “globin fold.” As illustrated in Figure 1, the globin fold is comprised of eight **alpha-helical** segments, designated A through H, which are connected by loop segments AB, BC,  $\frac{1}{4}$  and GH. The nonhelical segments at the amino and carboxyl termini are designated as NA and HC, respectively. The globin fold is maintained even when the sequence identity between pairs of globins is as low as 16% (2). In most vertebrates, individual globin **polypeptide chains** consist of 141 to 153 amino acid residues. Another property of the globin fold is its ability to bind a heme group, protoporphyrin IX with a ferrous iron atom at its center. The heme group is inserted in a cleft in the protein structure between the E and F helices, known as the “heme pocket.” The iron atom is linked covalently to the imidazole N<sub>ε</sub>; of the “proximal” [histidine](#) residue at position F8 (the eighth residue of the F helix). There are also many [van der Waals interactions](#) between the heme group and the amino acid side chains that line the heme pocket. Under physiological conditions, the iron atom is in the ferrous state, when it is capable of binding one molecule of O<sub>2</sub> on the “distal” side of the heme group. The side chain of the distal His-E7 forms a hydrogen bond with one of the oxygen atoms, stabilizing the bound oxygen (see [Myoglobin](#)). In this state, the ferrous iron is hexacoordinated with one of the oxygen atoms, N<sub>ε</sub>; of His-F8, and the four pyrrole nitrogen atoms of protoporphyrin IX. When oxygen is not bound, the sixth coordination site for oxygen is left vacant, and the iron atom is only pentacoordinated.

**Figure 1.** The globin fold. Cylinders represent the eight  $\alpha$ -helices. The lines correspond to the interhelical regions are

not exact and signify only the connectivity of the polypeptide chain. The bold short line indicates the side view of the heme group. From Ref. 2.



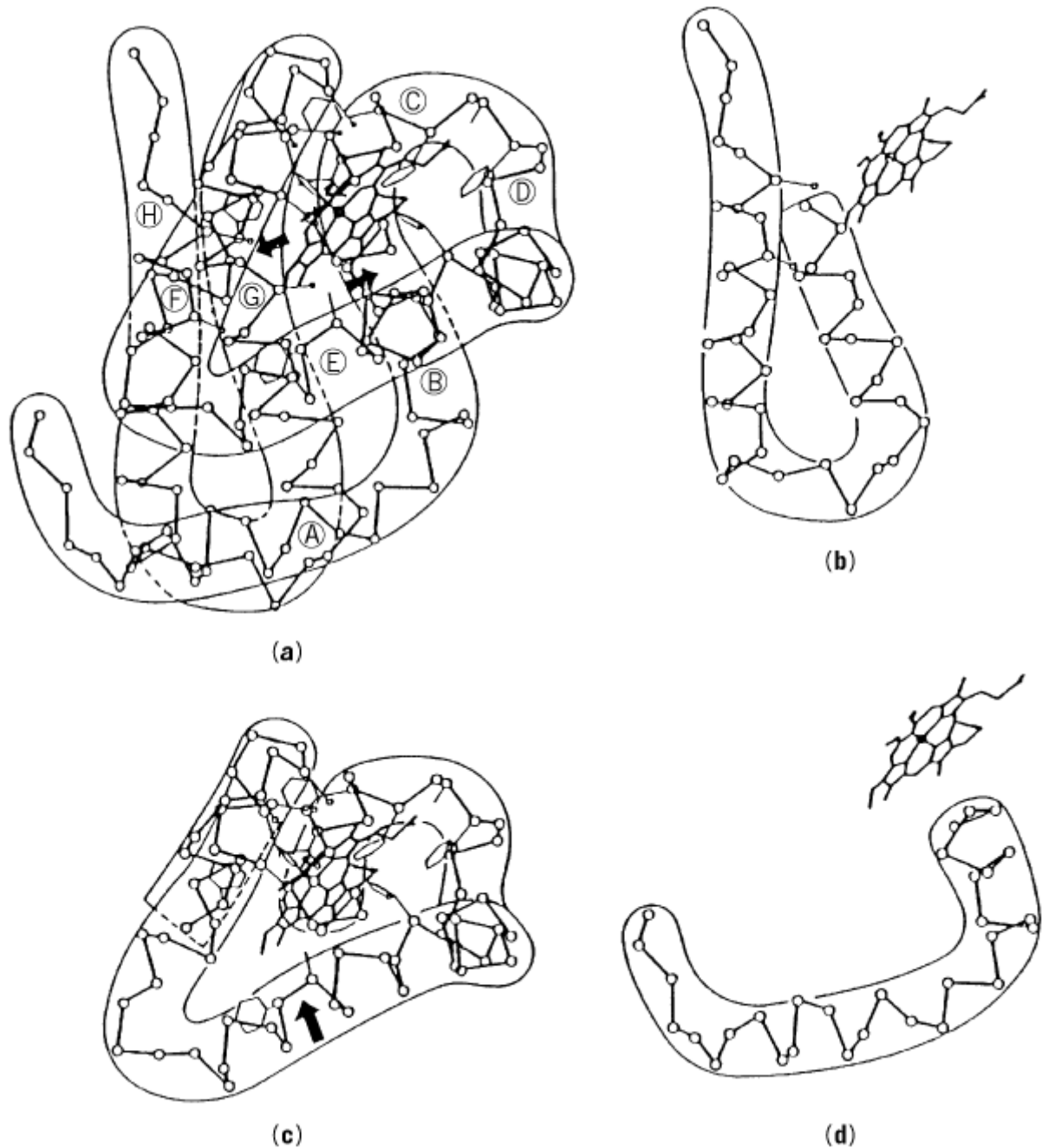
## 2. Distribution and Classification

Globins are found in organisms as [hemoglobin](#) (Hb) or [myoglobin](#) (Mb). Their common property is the reversible binding of oxygen. Mb is present in the muscles and stores oxygen, whereas Hb is present in the blood and transports oxygen from the lungs or gills to the peripheral tissues. Mb is usually a monomer, although sometimes a dimer, comprised of a single polypeptide chain and one heme group, whereas Hb is an oligomer or polymer composed two or more globin polypeptide chains, each with a heme group. Globin oxygen-binding proteins are even found in organisms that have no muscle or blood, including insects, **plants**, and single-celled organisms such as protozoa, **yeast**, and **bacteria**. These proteins are conventionally called hemoglobins, although some call them myoglobins because of their monomeric structures. Globin **genes** are widespread in modern plants, and plant hemoglobins are believed to have evolved from an **evolutionary** ancestor common to modern plants and animals (3).

## 3. Modular Structure and Evolution

The globin fold is composed of four structural units, called modules M1, M2, M3, and M4 (4) (Fig. 2). In most globins, modules M1 and M4 are encoded by two corresponding exons of the **gene** (see [Introns, Exons](#)). In contrast, modules M2 and M3, which are the structural units that hold the heme group, are encoded by a single exon. In the Hb genes of leghemoglobin (from the legume root nodule) and nematode Hb, however, this single exon is divided into two exons by the insertion of one intron. It is believed that this intron was present initially in the ancestral globin gene but lost during the molecular evolution of the other globin genes (6). Consequently, the globin fold is considered to be encoded by four exons that correspond to the four structural modules. This is one of the best indications of [protein evolution](#) by [exon shuffling](#).

**Figure 2.** Schematic representation of the three-dimensional structure of the human Hb-b chain and its constituent four structural modules. (a) The structure of the b-chain; the heme group is bound in a pocket between helices E and F. Arrows indicate the positions of the modular joints that correspond to the positions of the introns in the b-globin gene. (b) The first module, at the N-terminus, encoded by Exon 1. (c) The joined second and third modules, encoded by Exon 2; the heme group is bound by this pair of modules. (d) The fourth module, encoded by Exon 3. Reproduced from Ref. 5.



### Bibliography

1. J.C. Kendrew, R.E. Dickerson, B.E. Strandberg, R.G. Hart, D.R. Davies, D.C. Phillips, and V.C. Shore (1960) *Nature* **185**, 422–427.
2. A.M. Lesk and C. Chothia (1980) *J. Mol. Biol.* **136**, 225–270.
3. J. Landsmann, E.S. Dennis, T.J.V. Higgins, C.A. Appleby, A.A. Kortt, and W.J. Peacock (1986) *Nature* **324**, 166–168.

4. M. Go (1981) *Nature* **291**, 90–92.
5. M. Go (1983) *Shizen* (in Japanese) **7**, 26–35.
6. M. Go (1991) In *Evolution of Life: Fossils, Molecules, and Culture* (S. Osawa and T. Honjo, eds.), Springer-Verlag, Tokyo, pp. 109–122.

## Globulins

Globulins are proteins that are precipitated in both distilled water and 50% saturated ammonium sulfate. This behavior contrasts with that of [albumins](#), which show the opposite characteristics. This operational classification was very important in protein research initially, as these two characteristics were used in purification of various proteins. It is much less important today, but the name globulin is still encountered frequently.

### 1. Serum $\alpha_1$ -Globulins

In the **electrophoretic** fractionation of serum proteins, the  $\alpha$ -globulin fraction runs closest to [serum albumin](#). When electrophoresis is carried out in a Veronal buffer of pH 8.6, it is subfractionated into  $\alpha_1$  and  $\alpha_2$  fractions. The former contains  [\$\alpha\_1\$ -antitrypsin](#) (antiprotease),  $\alpha_1$ -acid glycoprotein,  $\alpha_1$ -microglobulin,  $\alpha$ -fetoprotein, and  $\alpha$ -lipoprotein (HDL). The  $\alpha_2$  fraction contains, among others,  $\alpha_2$ -microglobulin (retinol binding protein),  $\alpha_2$ -macroglobulin,  $\alpha_2$ -plasmin inhibitor, haptoglobin, and ceruloplasmin.

#### 1.1. $\alpha_1$ -Microglobulin

This glycoprotein (21.76% carbohydrate by weight) is present in plasma at a concentration of 1 mg/100 ml, with a molecular weight of 33,000 and unknown biological functions. It is heterogeneous in electrophoresis, but there is no proven heterogeneity in the amino acid sequence.

#### 1.2. $\alpha_1$ -Antitrypsin, Inter- $\alpha$ -trypsin Inhibitor

These are two [proteinase inhibitors](#) in the  $\alpha$ -globulin fraction. They inhibit many [proteinases](#), including **trypsin** and **chymotrypsin**, but their roles in homeostasis are not clearly established (see  [\$\alpha\_1\$ -Antitrypsin](#)).

#### 1.3. $\alpha$ -Fetoprotein

$\alpha$ -Fetoprotein is a protein of unknown function, but structurally related to [serum albumin](#). Its concentration in the serum is high in the fetus but very low in the adult.

#### 1.4. Orosomuroid

This is also called  $\alpha_1$ -acid glycoprotein and is a glycoprotein with an unusually high content of carbohydrate, 42% by weight. It may have some evolutionary relationship to immunoglobulins.

### 2. Serum $\alpha_2$ -Globulins

#### 2.1. $\alpha_2$ -Microglobulin

$\alpha_2$ -Microglobulin is also known as retinol binding protein (RBP) and has a molecular weight of 21,000, consisting of 181 amino acid residues, a [sedimentation coefficient](#) of 2.3 S, and no

carbohydrates; its structure is known (1). It binds one molecule of retinol (vitamin A) and transports it through the circulation in association with transthyretin (formerly called prealbumin, a homotetramer of 16-kDa subunits) as a 1:1 complex of 85 kDa. Its structure consists of a two-layered **up and down b-barrel** that encloses a primarily **hydrophobic** cavity designed to bind retinol. This protein is a member of a large family of ligand transport proteins of similar structures that includes [beta-lactoglobulin](#) and a bilin-binding protein.

## 2.2. $\alpha_2$ -Plasmin Inhibitor

$\alpha_2$ -Plasmin inhibitor is also known as  $\alpha_2$ -antiplasmin and has a molecular weight of 67,000 and carbohydrate content of 12% by weight. It inhibits [serine proteinases](#) by forming a covalent bond with the [serine](#) residue at their [active sites](#). This protein inhibits the degradation of fibrin clots by plasmin; in individuals who lack this protein, fibrinolysis occurs at an early stage of blood clotting, leading to a premature degradation of fibrin clot and causing serious bleeding.  $\alpha_2$ -plasmin inhibitor obstructs the action of plasmin toward fibrin by binding to specific regions of plasmin and its precursor [plasminogen](#).

## 2.3. $\alpha_2$ -Macroglobulin

This is a homo-tetrameric glycoprotein (7% carbohydrates) of 725,000 that inhibits almost any kind of proteinase, including those belonging to the four classes of [serine](#), [thiol](#), [metallo](#), and [carboxyl proteinases](#) (2). The mechanism of inhibition is entirely different from that of other proteinase inhibitors in that, instead of binding to the active sites of the target enzymes,  $\alpha_2$ -macroglobulin entraps medium-sized proteinases of between approximately 10,000–100,000 molecular weight within its cavity. The protein has a region in its amino acid sequence that is highly vulnerable to proteolytic attack; when one or more peptide bonds are hydrolyzed in this region, a large conformational change takes place to engulf the proteinase inside the molecular cavity and to react it with a highly reactive thiol ester bond of the macroglobulin. The active sites of the entrapped proteinases are not inhibited, and they remain active toward low-molecular-weight substrates, but are inactive toward large substrates, simply as a result of steric hindrance. The altered form of  $\alpha_2$ -macroglobulin is rapidly cleared from the blood serum, due to its internalization after binding to the receptor molecule called LRP (standing for low-density-lipoprotein receptor-related protein) on the membrane surface of peripheral cells such as fibroblasts, macrophages, hepatocytes, and so forth (see also [Macroglobulins](#)).

## 2.4. Haptoglobin

This is a glycoprotein (carbohydrate content of about 19% by weight) in the serum with the capacity to bind two half-molecules of free [hemoglobin](#) in the serum, as in the case of hemolytic anemia. It is composed of  $\alpha_1$  and  $\alpha_2$  chains of 9.1- and 16,000 molecular weight, plus two identical b chains of 40 kDa. Clearance of the haptoglobin-hemoglobin complex from the blood serum has a half-life of 10 to 20 min.

## 2.5. Ceruloplasmin

This is a blue-colored copper-containing glycoprotein, with eight copper atoms per molecule, a carbohydrate content of about 8% by weight, and a high molecular weight of 132,000. It migrates in the  $\alpha_2$  region, its sedimentation coefficient is 7 S, and it is made of four heterogeneous subunits.

## 2.6. Antithrombins

Antithrombin III is a glycoprotein (15% carbohydrate by weight) with a molecular weight of 65,000. Its concentration in plasma is 0.2 g/l. It inhibits [thrombin](#), activated forms of [blood clotting](#) factors IX, X, and XI, plasmin, and [trypsin](#). Upon association with heparin, it inhibits blood clotting.

## 2.7. Cobalamin-binding Proteins

(1) IF (standing for intrinsic factor) is a 114-kDa dimeric glycoprotein. It binds two molecules of vitamin [B12](#) (cobalamin) and helps its absorption in the intestine by binding to a specific receptor.

(2) Transcobalamin I, II, and III are vitamin B<sub>12</sub>-binding and transporting proteins in the bloodstream. I is a glycoprotein (carbohydrate content of 33% by weight) and has a molecular weight of 56,000. II is a simple protein of 60 kDa. (3) A protein known as R-binder (cobalophilin) also binds cobalamin.

### 2.8. Thyroid Hormone-binding Proteins

The [thyroid hormones](#), thyroxin (T4) and triiodothyronine (T3), are transported by thyroxin-binding pre-albumin, thyroxin-binding globulin, and serum albumin. Thyroxin and triiodothyronine in the blood are distributed among the three proteins in the ratios 70:10:20 and 40:30:30, respectively.

1. Thyroxin-binding pre-albumin migrates on electrophoresis to the anodic side of serum albumin, has a molecular weight of 54,980, contains 0.4% sugar, and has four subunits. The serum concentration of this protein becomes low in patients with liver disorders.
2. Thyroxin-binding globulin is an acidic glycoprotein with a molecular weight of 60,000; it is present in serum at a concentration of 15 µg/ml and strongly binds thyroxin and triiodothyronine.

### 2.9. Steroid Hormone-binding Proteins

These include corticosteroid-binding protein (transcortin), sex hormone-binding protein, and progesterone-binding protein. Transcortin is a glycoprotein, with a carbohydrate content of 16 to 26% by weight and molecular weight of 53,000.

## 3. Serum b-Globulins

### 3.1. Transferrin

This is an iron-binding glycoprotein in the serum with a molecular weight of 77,000. When ferric ion is bound to transferrin, it will bind to the cell-membrane transferrin receptor of erythroblasts and reticulocytes and becomes internalized as a result of receptor-linked [endocytosis](#); apotransferrin has low affinity for the receptor. The internalized transferrin releases the bound iron as the pH inside the endosome becomes lower, and then the receptor with bound apo-transferrin reassociates with the membrane, to recycle apo-transferrin into the serum. Transferrin has two iron-binding sites per molecule, but a similar protein found in hag-fish has only one iron binding site per molecule of 44 kDa. There are several genetic variants. A protein called conalbumin in egg white and lactoferrin in milk are actually apo-transferrin, with the same amino acid sequence as the transferrin in blood serum, but with different [post-translational modifications](#). More than 60% of the iron used by blood-producing cells for the hemoglobin biosynthesis is supplied by transferrin.

### 3.2. b-Lipoproteins

These include low-density and very-low-density **lipoproteins**. Serum lipoproteins are noncovalent complexes of specific proteins and lipids. They constitute a lipid transport system in the bloodstream and are classified in as chylomicrons, very-low-density lipoproteins (VLDL), low-density lipoproteins (LDL), and high-density lipoproteins (HDL). They appear in the a-globulin fraction upon electrophoresis. A further subfraction, very-high-density lipoprotein (VHDL), is often identified. Each class of lipoproteins is characterized by specific apo-proteins and immunogenicities. The density classification of lipoproteins is based on lipid content; the higher the lipid content, the lower the density. The lipid content in each class of lipoproteins is not constant but varies within a certain range of values. (More details can be found under **Lipoproteins**).

### 3.3. b<sub>2</sub>-Microglobulin

This is the smaller of the two polypeptide chains that constitute the [histocompatibility](#) complex antigen, but this protein does not have specific antigenicity. The amino acid sequence has a similarity to that of the constant domains of the [immunoglobulin](#) heavy chains.

### 3.4. Hemopexin

This is a heme-binding protein with a molecular weight of 5700 and a carbohydrate content of 20%.

#### 4. Serum g-Globulins

Serum g-globulins constitute the secretory immune system of vertebrates and are classified as IgA, IgD, IgE, IgG, and IgM, with Ig standing for [immunoglobulin](#). The basic structural unit of immunoglobulins may be seen in the  $H_2L_2$  structure of IgG, where H represents the heavy chain and L the light chain. Each light chain is disulfide-bonded to one of the heavy chains, and the two heavy chains are disulfide-bonded to form a covalently associated hetero-tetrameric subunit structure. Two [antigen](#)-binding sites are formed from the two pairs of H and L chains. The structural designs of other immunoglobulins are based on this IgG structure. More details can be found under [Immunoglobulins](#).

#### 5. Egg-White Globulins

##### 5.1. Ovoglobulins

This fraction of egg white contains [lysozyme](#) and ovomacroglobulin as the major proteins in the G1 and G2 fractions, respectively. Lysozyme is a small basic protein with a molecular weight of 14,300 (129 amino acid residues), sedimentation coefficient of 1.9 S, and an isoelectric point of 11. It is an enzyme that cleaves cell-wall mucopolysaccharides at b(1–4) linkages between N-acetylmuramic acid and N-acetyl galactosamine, and thus has bactericidal activity. Chicken egg-white lysozyme is now called c-type lysozyme with a distinction from g-type lysozymes, such as that with 185 residues purified from goose eggs. Hen lysozyme constitutes about 3% of egg-white proteins and is readily crystallized. See also [Lysozymes](#).

##### 5.2. Ovomacroglobulin

This is a large, tetrameric glycoprotein with a molecular weight of 700,000 ( $4 \times 170,000$ ). It has a **proteinase**-trapping activity similar to that of serum  $\alpha_2$ -macroglobulin. The two proteins are **homologous** in amino acid sequence and to the C3, C4, and C5 components of the [blood clotting](#) proteins. Unlike  $\alpha_2$ -macroglobulin, ovomacroglobulin lacks the internal thiol ester bond. The presence of homologous proteins has been confirmed in the chicken, goose, crocodile, and turtle egg white (3).

#### 6. Thyroid Globulins

Thyroglobulin is a glycoprotein (10% sugar) of molecular weight 670,000, two identical subunits, sedimentation coefficient of 19 S, and pI of 4.5. It is synthesized as an iodine-free protein in the epithelium of the thyroid gland; later it is iodinated on some of its [tyrosine](#) residues by thyroid peroxidase. Two iodotyrosine residues dimerize to form the [thyroid hormones](#) 3,3',5,5'-tetraiodothyronine (thyroxine) and 3,3',5-triiodothyronine. Iodination also promotes the generation of covalent dimers of thyroglobulin through disulfide bond formation. Iodination of the protein is specific to the thyroid gland of vertebrates.

#### 7. Milk Globulins

Cow's whey contains as the principal components (50 to 60%) [beta-lactoglobulin](#), which is peculiar to ruminant milk, and [alpha-lactalbumin](#). b-Lactoglobulin constitutes about 50% of the whey protein; it is fractionated as an albumin fraction that does not precipitate with 50% saturation ammonium sulfate, but the purified protein shows a globulinlike solubility in water. Whey protein is fractionated into three major fractions by centrifugation, and b-lactoglobulin is the major protein in the b fraction. It is a dimeric protein of 35 kDa, sedimentation coefficient of 2.8 S, and pI of 5.2. The lactoglobulin fraction includes b-lactoglobulin and immunoglobulin derived from the mother. The three-dimensional structure of this protein is similar to that of retinol-binding proteins (4).

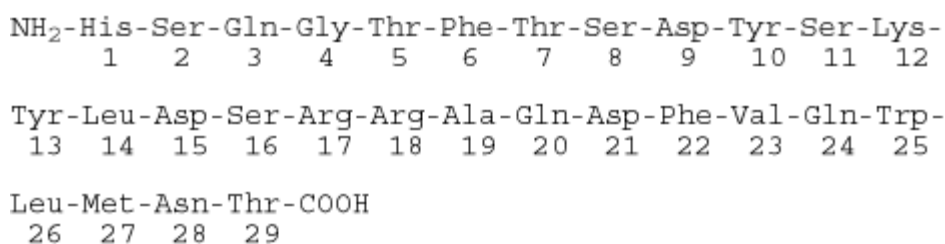
#### Bibliography

1. S. Cowan, M. E. Newcommer, and T. A. Jones (1990) Crystallographic refinement of serum retinol binding protein at 2 Å resolution, *Proteins* **8**, 44–61.
2. W. Borth, R. D. Feinman, S. L. Gonias, J. P. Quigley, and D. K. Strickland, eds. (1994). *Biology of Macroglobulin. Its Receptor, and Related Proteins*, *Annals N. Y. Acad. Sci.* **737**.
3. A. Ikai, M. Kikuchi, and M. Nishigai (1990) Internal structure of ovomacroglobulin, *J. Biol. Chem.* **256**, 8280–8284.
4. H. L. Monaco, G. Zonatti, P. Spadon, M. Bolognesi, L. Sawyer, and E. E. Eliopoulos (1987) Crystal structure of trigonal form of bovine  $\gamma$ -lactoglobulin and its complex with retinol at 2.5 Å resolution, *J. Mol. Biol.* **197**, 695–706.

## Glucagon

Glucagon is a 29-residue, single-chain [peptide hormone](#) (MW 3485) (1). Its gene is located on human chromosome 2, and its [primary structure](#) is given in Figure 1. Glucagon has a regulatory function in glucose metabolism; in general, its effects are opposite those of [insulin](#). It is synthesized by  $\alpha$ -cells in the islets of Langerhans of the pancreas and in L-cells of the intestinal mucosa as a very much larger pre-pro-glucagon molecule of 180 amino acid residues. The [signal peptide](#) is removed proteolytically upon translocation of the nascent polypeptide chain into the [endoplasmic reticulum](#), to generate proglucagon (160 residues); this contains within it other peptides. It is processed by the pro-hormone convertases 1/3 (PC-1/3) and PC-2 (see [Peptide Hormones](#)) differently in the two tissues (2). In the pancreas, proglucagon (PG) yields glucagon and glicentin-related peptide (GRPP). In the intestinal mucosa, glucagon-like peptides-1 (GLP-1) and -2 (GLP-2) and glicentin (a 69-residue peptide containing glucagon and GRPP) are produced (Fig. 2).

**Figure 1.** The primary structure of glucagon.



**Figure 2.** Proteolytic processing of proglucagon (PG) in the pancreas and in the intestine. GRPP: glicentin-related peptide. GLP: glucagon-like peptide.



|           |                 |               |    |                 |       |       |
|-----------|-----------------|---------------|----|-----------------|-------|-------|
|           | PG <sup>1</sup> |               |    |                 |       | 160   |
| Pancreas  | 1               | 30            | 33 | 61              | 72    | 158   |
|           | GRPP            | Glucagon      |    | Other Fragments |       |       |
| Intestine | 1               |               |    |                 | 69    | 78    |
|           | Glicentin       |               |    |                 | 78    | 107   |
|           |                 |               |    |                 | GLP-1 | GLP-2 |
|           | 1               | 30            | 33 |                 |       | 69    |
|           | GRPP            | Oxyntomodulin |    |                 |       |       |

The [secretory granules](#) of the  $\alpha$ -cells contain glucagon and discharge it by [exocytosis](#). Active GLP-1, a potent hormone that also antagonizes insulin action, is secreted from the intestine (3-6). Consequently, differential processing of pro-glucagon gives rise to two different active peptides in the regulatory pathways. Both glucagon and GLP-1 have their own specific receptor; their genes have been **cloned** and the proteins characterized in structure-function and [signal transduction](#) studies (7, 8).

In general, glucagon and GLP-1 antagonize many actions of insulin and regulate the levels of blood glucose. Like insulin, glucagon lacks a plasma carrier, and its circulation half-life is short, approximately 3 to 6 minutes. Glucagon is removed from the circulatory system primarily by the liver and kidney. The predominant effect of glucagon is on the liver, because it is the first organ perfused by blood with the pancreatic secretions. Glucagon binds to receptors on the plasma membrane and couples through [heterotrimeric G-proteins](#) to [GTP-binding proteins](#) and [adenylate cyclase](#). The resultant increases in [cyclic AMP](#) (cAMP) and cAMP-dependent protein kinase A (PKA) levels reverse the effects of insulin on the liver. These increases elevate the circulating glucose levels by stimulating the output of glucose from the liver by glycogenesis and promoting glycogenolysis. Glucagon secretion is inhibited by glucose. In contrast to insulin, which promotes energy storage, glucagon releases energy. It also has some effect on lipid metabolism, by reducing fatty acid synthesis and increasing fatty acid oxidation and ketogenesis. Glucagon secretion does not markedly fluctuate throughout the day, but the ratio of insulin to glucagon mediates **phosphorylation** or dephosphorylation of enzymes that regulate metabolism.

In type I diabetes, the stimulating effect of hypoglycemia on glucagon secretion is reduced. In these diabetic patients, glucose fails to inhibit glucagon secretion, but insulin injections restore the effect. On the other hand, somatostatin is an inhibitor of both glucagon and insulin secretion. A potent glucagon receptor antagonist may have therapeutic use in hypoglycemia. Consequently, the roles of glucagon and GLP-1 in glucose metabolism are presently an area of diabetes research.

#### Bibliography

1. P. J. Lefebvre, editor (1996) *Glucagon III—Handbook of Experimental Pharmacology*, vol. **123**, Springer, New York.
2. S. Dhanvantari, N. G. Seidah, and P. L. Brubaker (1996) *Mol. Endocrinol.* **10**, 342–355.
3. H. C. Fehmann, editor (1997) *The Insulinotropic Gut Hormone Glucagon-Like Peptide-1* ("Frontiers of Diabetes, vol. 13"), Karger, New York.
4. J. J. Holst (1996) *Acta Physiol. Scand.* **157**, 309–315.
5. M. E. Rothenberg et al. (1996) *J. Biol. Chem.* **270**, 10136–10146.
6. D. J. Drucker (1990) *Pancreas* **5**, 484–488.
7. J. Christophe (1995) *Biochemica et Biophysica Acta* **1241**, 45–57.

## Glucocorticoid Response Element

The [steroid hormones](#) are a group of substances derived from cholesterol that exert a wide a range of effects on biological processes (1) by stimulating the [transcription](#) of specific **genes** (2, 3). These effects on gene expression are mediated by specific [response elements](#) within the target genes. The relationships between the different response elements that mediate the response to each of these steroid hormones are discussed in the entry [Hormone response elements](#). This article discusses the sequence that mediates the response to one of the classes of steroid hormones, the glucocorticoids.

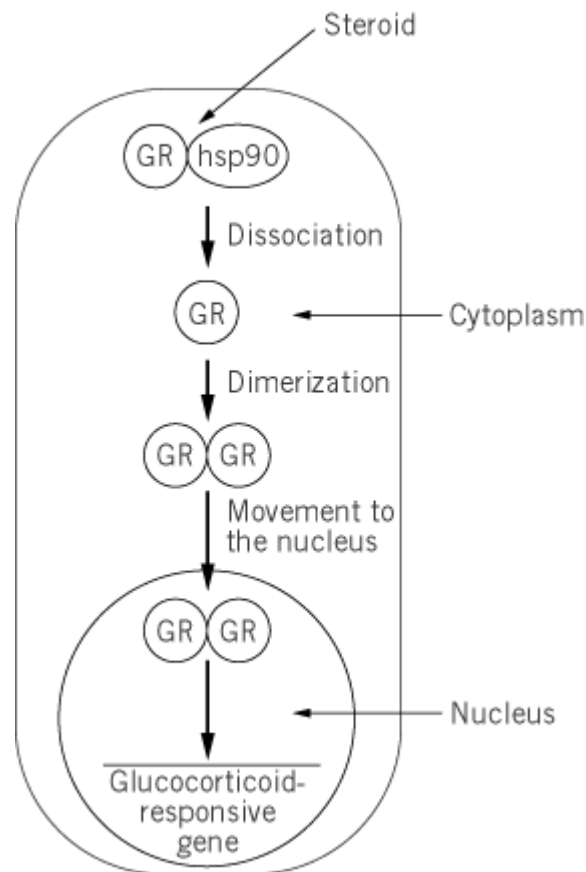
The glucocorticoid response element (GRE) is found in a number of different target genes that respond to glucocorticoid. It is a short element with the sequence

5'RGRACANNNTGTYCY3'

3'YCYTGTNNNACARGR5'

where R = purine, (ie, A or G), Y = pyrimidine (ie, C or T), and N = any base. The GRE thus consists of a 6-bp [palindrome](#) in which a 6-bp sequence (RGRACA) is repeated in the opposite orientation, with a spacing between the two halves of the palindrome of 3 bp whose sequence does not matter. This palindromic sequence is bound by two molecules of the glucocorticoid **receptor** in a dimer, which mediates the response to glucocorticoid.

In cells that have not been exposed to steroid, this receptor exists in the cytoplasm complexed to a 90-kDa heat-inducible protein, hsp90. The association with hsp90 anchors the glucocorticoid receptor in the cytoplasm and prevents it activating any genes. Upon exposure to glucocorticoid, the receptor binds to the steroid, which results in the dissociation of hsp90. This leaves the receptor free to dimerize and to move to the [nucleus](#), where the dimer binds to the palindromic GRE (Fig. 1). This binding takes place via specific regions of the receptor known as the DNA-binding **domain** and the activation domain. The bound receptor then interacts with the basal transcriptional complex of [RNA polymerase](#) and associated factors to stimulate transcription of the steroid-responsive gene. In this manner, binding the activated glucocorticoid receptor to the GRE causes the stimulation of glucocorticoid-responsive genes following exposure of the cell to steroid (4).



Interestingly, the glucocorticoid receptor can also inhibit the expression of specific genes, such as those encoding prolactin and pro-opiomelanocortin. This effect is achieved by the identical receptor–hormone complex that activates glucocorticoid-inducible genes. However, the DNA sequence element to which the complex binds when mediating its negative effect (nGRE) is distinct from the glucocorticoid response element (GRE) to which it binds when inducing gene expression, but the two sequences are related to one another:

**GRE RGRACANNNTGTYCY**

**nGRE ATYACNNNTGATCW**

where N = any base, R = purine, Y = pyrimidine, and W = A or T. The similarity between the GRE and the nGRE has led to the idea that the sequence difference causes the receptor–hormone complex to bind to the nGRE in a configuration in which its activation domain cannot interact with the basal transcription complex to activate transcription in the way that occurs following binding to the positive element (5). Indeed, it has been shown that the nGRE in the pro-opiomelanocortin gene binds the glucocorticoid receptor as a trimer, rather than the dimer form that binds to the GRE and stimulates transcription (6). When bound in this manner, the GRE is likely to act simply by preventing binding of a positively acting [transcription factor](#) to this or an adjacent site, thereby preventing gene induction. In agreement with this, the nGRE in the glycoprotein hormone  $\alpha$  subunit gene overlaps a [cyclic AMP](#) response element (CRE) and is able to inhibit gene expression only when the activating CRE is left intact (7).

Hence, the GRE and the related nGRE can bind the glucocorticoid receptor–hormone complex in distinct configurations and thereby allow glucocorticoid hormone to exert either positive or negative effects on the expression of specific target genes.

## Bibliography

1. R. J. B. King and W. I. P. Mainwaring (1974) *Steroid Cell Interactions*, Butterworths, London.
2. D. S. Latchman (1998) *Eukaryotic Transcription Factors*, 3rd ed., Academic Press, London.
3. M. Beato (1989) *Cell* **56**, 335–344.
4. W. B. Pratt, D. J. Jolly, D. V. Pratt, S. M. Hollenberg, V. Giguere, F. M. Cadepond, G. Schweizer-Groyer, M-G. Catelli, R. M. Evans, and E-E. Bahileu (1988) *J. Biol. Chem.* **263**, 267–273.
5. D. D. Sakki, S. Helms, J. Carlstedt-Duke, J. A. Gustafsson, F. M. Rottman, and K. R. Yamamoto, (1988) *Genes Develop.* **2**, 1144–1154.
6. J. Drouin, Y. L. Sun, M. Chamberland, Y. Ganthier, A. DeLea, A. M. Nemer, and T. J. Schmidt (1993) *EMBO J.* **12**, 145–156.
7. I. W. Akerblum, E. P. Slater, M. Beato, J. D. Baxter, and P. L. Mellon (1988) *Science* **241**, 350–353.

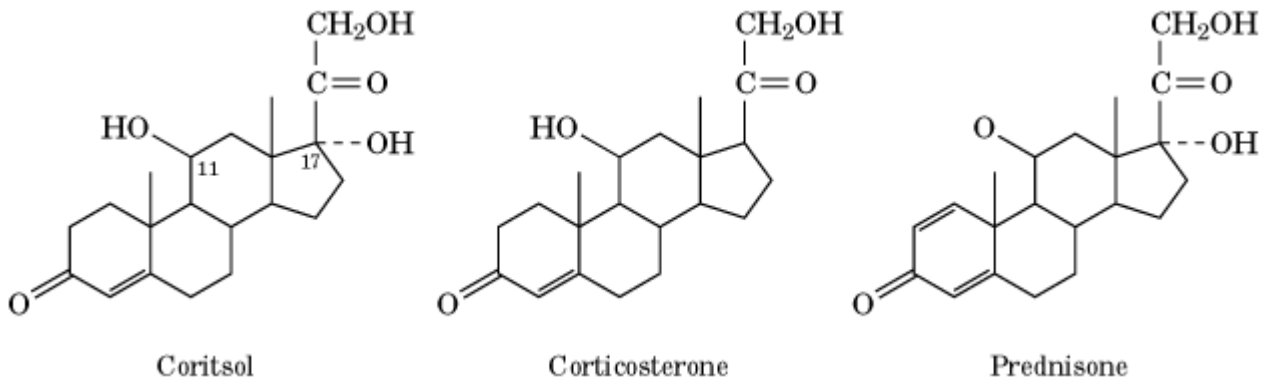
## Suggestions for Further Reading

8. H. Gronemeyer and V. Laudet (1995). Nuclear receptors. *Protein Profile* **2**, 1173–1308.
9. D. J. Mangelsdorf, C. Thummel, M. Beato, P. Herrlich, G. Schutz, K. Umesono, B. Blumberg, P. Castner, M. Mark, P. Shambon, and R. M. Evans (1995) The nuclear receptor super family; the second decade. *Cell* **83**, 835–839.
10. M. Truss and M. Beato (1993) Steroid hormone receptors: interaction with the deoxyribonucleic acid and transcription factors. *Endocr. Rev.* **14**, 459–479.

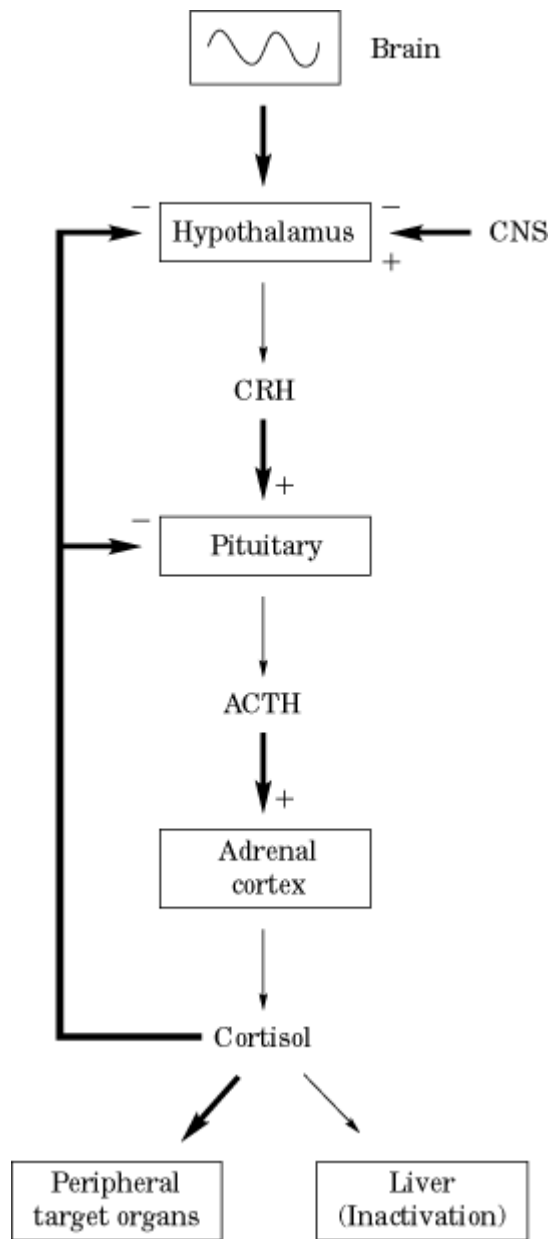
## Glucocorticoids

Glucocorticoids (or glucocorticosteroids) are natural [steroid hormones](#), or synthetic analogues, that regulate metabolic, developmental and immunological processes. The major natural glucocorticoids are cortisol (11 $\beta$ ,17,21, trihydroxy-4-pregnene-3,20-dione) (in humans) and corticosterone (11 $\beta$ ,21,d,hydroxypregn-4-ene-3,20-dione) (in rodents) (Fig. 1). Semisynthetic derivatives with greater hormonal activity than that of the natural steroids are [dexamethasone](#), prednisone, and triamcinolone (Fig. 1). The synthesis of glucocorticoids in the zona fasciculata of the adrenal cortex and their release in the blood circulation is regulated by adrenocorticotrophic hormone (ACTH), which is produced in the anterior lobe of the hypophysis, in response to stimuli coming either from the neuroendocrine cells of the paraventricular nuclei of the hypothalamus (corticotropin releasing hormone) or from the periphery (blood cortisol concentration). High blood cortisol levels reduce, by feedback inhibition, the production of both corticotropin-releasing hormone and ACTH. This regulatory circuit is called the hypothalamus–pituitary–adrenal axis (see Fig. 2). Glucocorticoid secretion is enhanced during stress conditions, such as anxiety, hunger, and trauma. Most effects of glucocorticoids (see [Hormone Receptors](#)) are mediated by the glucocorticoid receptors, which belong to the superfamily of nuclear receptors. There are two forms of glucocorticoid receptor, A and B, differing in their C-terminal regions. Upon binding the hormone, the receptor dimerizes, binds to [glucocorticoid response elements](#) (the **consensus sequence** is GGACAnnnTGTTCT), which are receptor-activated **enhancers**, and modulate gene [transcription](#).

**Figure 1.** Structure of glucocorticoids.



**Figure 2.** Regulation of cortisol production by the hypothalamus–pituitary–adrenal axis.



The glucocorticoids regulate the metabolism of sugars, lipids, and proteins. They stimulate glucose output from the liver, inhibit glucose utilization by tissues, and mobilize [fatty acids](#) from adipose tissue. The hormones act on muscle, bone, and lymph organs, inhibiting **protein biosynthesis** and enhancing [protein degradation](#). Some of the [amino acids](#) released in the circulation are used for gluconeogenesis, primarily in the liver but also in kidneys. The action of glucocorticoids on gluconeogenesis represents one of the most exploited and rewarding systems for delineating the molecular mechanisms of steroid hormone action (1, 2). In gluconeogenesis, amino acids are metabolized to two-carbon precursors to be used for glucose synthesis. The biosynthesis of enzymes like tyrosine aminotransferase, alanine aminotransferase, tryptophan dioxygenase, pyruvate carboxylase, and phosphoenolpyruvate-carboxykinase, which catalyze 2C-precursor production or utilization for glucose synthesis, are induced by glucocorticoids. All steps—from binding of the glucocorticoid ligand to its receptor, to increased transcription of the genes, to [translation](#) of the [messenger RNA](#) to the respective proteins—have been studied in detail and have supported the steroid hormone gene-activation hypothesis (3). The inhibition of protein synthesis by glucocorticoids in the lymphatic system results in decreased antibody production (immunosuppression) and increased susceptibility to infection. The suppressive effects of glucocorticoids on the inflammatory response have led to new insights and a novel mechanism of action of glucocorticoids. Namely, the hormones block binding of the [transcription factor](#) NF- $\kappa$ B on **promoters** having NF- $\kappa$ B binding sites (eg, the promoter of the [interleukin](#) gene) by direct interaction between the receptor and the transcription factor, thereby inhibiting expression of the respective (eg, interleukin) gene (4).

The role of glucocorticoids in embryonic development has been investigated in transgenic mice deficient in the glucocorticoid receptor (5). These mice die within a few hours of birth because of respiratory failure due to severely atelectic lungs, probably resulting from lowered production of surfactants and from a deficiency of a glucocorticoid-inducible sodium pump. The livers of newborns have reduced capacity to activate genes for key gluconeogenic enzymes, whereas the adrenals lack a central medulla and cannot synthesize adrenaline. The hypothalamus–pituitary–adrenal axis is impaired, resulting in increased expression of corticotropin-releasing hormone and in hypertrophy of the adrenal cortex, signifying that the hypothalamus–pituitary–adrenal axis is established during fetal development, and underlining the central role of glucocorticoids in this process.

#### Bibliography

1. G. Schutz (1988) *Biol. Chem. Hoppe Seyler* **369**, 77–89.
2. S. J. Pilkis and D. K. Granner (1992) *Annu. Rev. Physiol* **54**, 885–909.
3. P. Karlson (1963) *Perspect. Biol. Med.* **6**, 203–214.
4. R. Scheinman et al (1995) *Mol. Cell Biol.* **15**, 943–953.
5. T. J. Cole et al. (1995) *Genes Dev.*, **9**, 1608–1621.

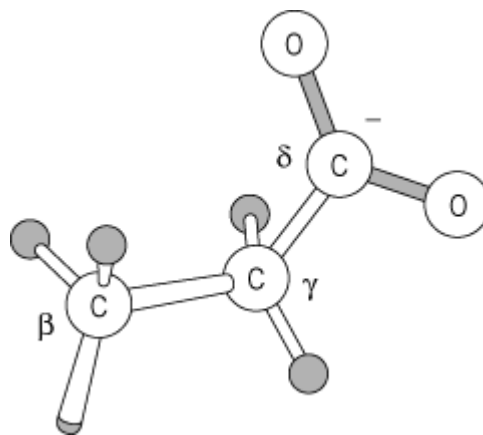
#### Suggestions for Further Reading

6. D. N. Orth, W. J. Kovacs, and C. R. DeBold (1992) In *Williams Textbook of Endocrinology*, Vol. **9**, 8th ed. (J. D. Wilson and D. W. Foster, eds.), Saunders, London, pp. 489–620.
7. G. P. Chrousos (1992) *Endocrinol. Metab. Clin. North Am.* **21**, 833–858.
8. M. Truss and M. Beato (1993) *Endocr. Rev.* **14**, 459–479.

## Glutamic Acid (Glu, E)

The **amino acid** glutamic acid is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—GAA and GAG—and represents approximately 6.2% of the residues of the proteins that have been characterized. The glutamyl residue incorporated has a mass of 129.12 Da, a **van der Waals volume** of  $109 \text{ \AA}^3$ , and an [accessible surface](#) area of  $183 \text{ \AA}^2$ . Glu residues are frequently changed during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with **aspartic acid**, **glutamine**, **alanine**, and **histidine** residues.

The side chain of Glu residues is dominated by its **carboxyl group**:



This carboxyl group is normally no more reactive than are those of corresponding organic molecules, such as acetic acid. Its intrinsic  $pK_a$  value is approximately 4.3, so Glu residues are ionized and very polar under physiological conditions; consequently, very few Glu residues are buried in folded **protein structures**, and nearly all have at least the carboxyl group on the surface. The  $pK_a$  can be shifted in folded proteins, however, and either the ionized or nonionized form can be used in the protein's function. For example, **carboxyl proteinases** have one **active-site** carboxyl group function in the ionized form, another nonionized. Asp carboxyl groups have a weak intrinsic affinity for  $\text{Ca}^{2+}$  ions and are used in many [calcium-binding proteins](#). Certain Glu residues, particularly in proteins involved in blood clotting and bone structure, are carboxylated to yield the unusual residue **g-carboxylglutamic acid**; such residues have two adjacent carboxyl groups and bind  $\text{Ca}^{2+}$  more avidly.

Glu residues differ from Asp only in having two methylene groups, rather than one, so it might be thought that they would be very similar chemically and functionally in proteins, but this is not so. The slight difference in length of the side chains causes them to have different tendencies in their chemical interactions with the peptide backbone, so they have markedly different effects on the conformation and chemical reactivity of the peptide backbone. For example, Glu residues favor the  **$\alpha$ -helical** conformation much more than do Asp. In folded proteins, Glu residues are most frequently found in  $\alpha$ -helices, whereas Asp residues occur most frequently in reverse [turns](#) .

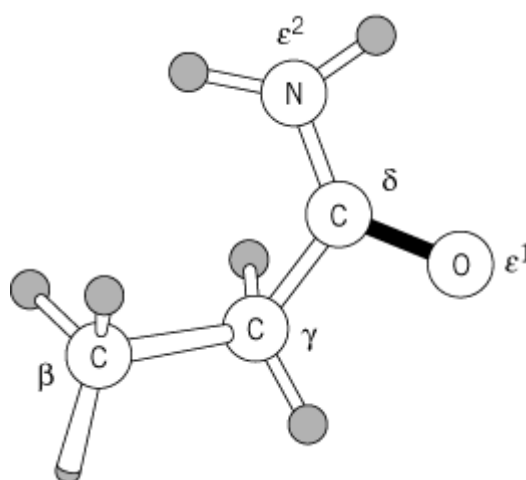
### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Glutamine (Gln, Q)

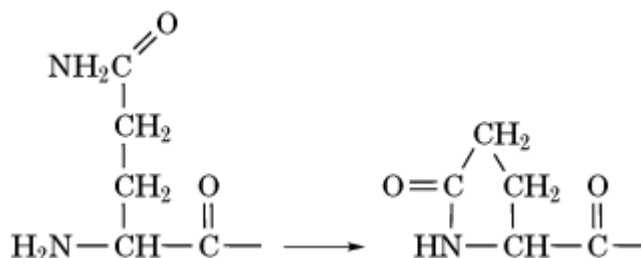
The [amino acid](#) glutamine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—CAA and CAG—and represents approximately 4.0% of the residues of the proteins that have been characterized. The glutamyl residue incorporated has a mass of 128.14 Da, a [van der Waals volume](#) of 114 Å<sup>3</sup>, and an [accessible surface](#) area of 189 Å<sup>2</sup>. Gln residues have close to average conservation during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [Glutamic Acid \(Glu, E\)](#) and [Histidine \(His, H\)](#) residues.

The side chain of Gln residues is the same as that of glutamic acid, except that the [carboxyl group](#) has been converted to the amide:



Both amino acids occur naturally and are incorporated directly into proteins during their biosynthesis. Gln residues do not arise from amidation of Glu in proteins. The amide side chain does not ionize and is not very reactive chemically. It is polar, however, as it is both a [hydrogen bond](#) donor and acceptor. The amide group is labile at extremes of pH and at high temperatures, and Gln residues can **deamidate** to Glu, although usually less rapidly than Asn residues do.

When Gln residues are at the *N*-terminus of a peptide chain, they spontaneously cyclize:



The resulting residue of pyrrolidone carboxylic acid renders the *N*-terminus unreactive in most procedures for **sequencing proteins**, such as the [Edman Degradation](#). This residue can be removed, however, by the enzyme pyroglutamyl amino peptidase

### Suggestions for Further Reading

F. Wold (1985) Reactions of the amide side-chains of glutamine and asparagine *in vivo*, Trends Biochem. Sci. **10**, 4–6.

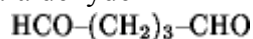
T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman,



New York.

## Glutaraldehyde

The bis-homobifunctional aldehyde glutaraldehyde



is perhaps the most frequently used chemical [crosslinking](#) agent. Because of its effectiveness as a crosslinker, glutaraldehyde has found use as a stabilizer or fixative for samples subjected to [electron microscopy](#), for bioprosthetic materials, and in drug delivery systems (reviewed in Ref. [1](#)). It is also routinely used to produce bioconjugates and to immobilize proteins (see [Enzyme Immobilization And Conjugation](#)).

Crosslinking by glutaraldehyde can occur through at least two distinct mechanisms ([2](#)). The first is attack by primary [amino groups](#) of the protein on the aldehydic carbons of glutaraldehyde to form [Schiff bases](#), which can then be rendered irreversible by suitable reductants. The second mechanism is more complex and less understood, and it results in irreversible crosslinking in the absence of reductants. This alternative mechanism involves polymerization of glutaraldehyde to form heterogeneous unsaturated polymers that differ in their length and degree of unsaturation. These polymers then undergo addition at their double bonds by other nucleophilic groups of the protein. A disadvantage of glutaraldehyde is that it often brings about extensive crosslinking and formation of insoluble complexes; however, this can be minimized through use of the two-step reaction, in which one protein is reacted with glutaraldehyde, the excess reagent removed, and the second protein then added to react with the free aldehyde groups of the first ([2](#)).

### Bibliography

1. A. Jayakrishnan and S. R. Jameela (1996) *Biomaterials* **17**, 471–484.
2. G. T. Hermanson (1996) *Bioconjugate Techniques*, Academic Press, San Diego, pp. 218–220.

## Glutaredoxin

Glutaredoxin (Grx) was discovered as a small protein that contains two [thiol groups](#) and acts as a [glutathione](#) (GSH)-dependent hydrogen donor for synthesizing deoxyribonucleotides by *Escherichia coli* [ribonucleotide reductase](#) ([1](#)). Glutaredoxin in the oxidized form has a redox-active [disulfide bond](#) (Grx-S<sub>2</sub>) and is reduced to the dithiol form (Grx-SH<sub>2</sub>) by two molecules of GSH in a **thiol-disulfide exchange** reaction via a GSH mixed-disulfide intermediate ([2](#), [3](#)). Glutaredoxin is also a general GSH-disulfide oxidoreductase that reduces either disulfide bonds by a dithiol mechanism or GSH-mixed disulfides by a monothiol mechanism ([2](#), [3](#)). The glutaredoxin system comprises NADPH, GSH, glutathione reductase, and glutaredoxin. It catalyzes the NADPH- and GSH-dependent reduction of disulfides and keeps the inside of the cell reduced.

Glutaredoxins exist in all prokaryotic and eukaryotic cells that contain glutathione, and they are also

encoded by the genes of the larger DNA **viruses**. Glutaredoxins generally have a mass of 9 to 12 kDa, and the [active site](#) has the **consensus sequence** -Cys-Pro-Tyr-Cys- (4). Glutaredoxin in mammalian cells is identical to GSH-homocystine transhydrogenase or thioltransferase. The primary structures of glutaredoxin show only limited sequence similarity to those of **thioredoxins**, but the three-dimensional structure of glutaredoxin shows that the active site dithiol/disulfide is located at the end of a [beta-strand](#) and at the beginning of an **alpha-helix**, characteristic of the thioredoxin fold. One hallmark of glutaredoxin is a binding site for GSH that is used in both the reduction of Grx-S<sub>2</sub> by GSH and for substrate recognition and reduction of GSH-mixed disulfides by Grx-(SH)<sub>2</sub>.

Glutaredoxins have numerous functions, such as acting as hydrogen donors for reductive enzymes and transmitting changes in the [oxidation/reduction potential](#) of the GSH/GSSG (glutathione disulfide) system for protein activity during thiol redox regulation and [signal transduction](#), or in restoring oxidatively damaged proteins after oxidative stress.

## 1. Glutaredoxin Functions

Glutaredoxins have been isolated from *E. coli* (1, 5), yeast (6), plants (7), and mammals (8, 9) and are encoded by **T4 bacteriophage** and [vaccinia virus](#) (10-12). *E. coli* contains three different glutaredoxins: Grx1, Grx2, and Grx3 (5). Grx1 is a highly efficient hydrogen donor for ribonucleotide reductase, either in the presence of 1 mM DTT or with the more physiological system of 4 mM GSH, NADPH, and excess glutathione reductase, when the apparent [Km \(Michaelis constant\)](#) for Grx1 is 0.13 μM (2, 3). The apparent turnover number of glutaredoxin as a dithiol hydrogen donor for ribonucleotide reductase is tenfold higher than that of thioredoxin. Glutaredoxin was discovered in an *E. coli* mutant that lacks thioredoxin (1), and double mutants that lack thioredoxin and glutaredoxin are viable (13). This includes *E. coli* cells that lack glutathione (gshA<sup>-</sup>). In combination with a thioredoxin mutation (gshA<sup>-</sup>, trxA<sup>-</sup>), such cells very strongly induce (up to 55-fold) Grx1 (14). Grx3 in *E. coli* has the same active site structure as Grx1, but only 6% of the V<sub>max</sub> with ribonucleotide reductase and a 35 mV higher redox potential (E'<sub>O</sub>) (15, 16). The function of at least Grx1 strongly depends on the GSH-concentration and the GSH:GSSG ratio, because the E'<sub>O</sub> of the redox-active disulfide bond in glutaredoxin determines whether it is reduced to the active dithiol (16).

Glutaredoxin functions reducing both sulfate and methionine sulfoxide in *E. coli* and yeast (17). Synthesis of Grx1 in *E. coli* is induced via the oxyR [transcription factor](#), and Grx1 controls its activity by reducing a disulfide bond that can be formed in the protein (18). The two glutaredoxins in yeast control the defense against superoxide and hydrogen peroxide stress, respectively, in addition to their roles in deoxyribonucleotide synthesis by ribonucleotide reductase (6).

Phage T4 induces a ribonucleotide reductase upon infection of *E. coli*, and a T4 glutaredoxin (originally called thioredoxin) is also a specific hydrogen donor for the enzyme that works with GSH as a reductant (19), and a substrate for thioredoxin reductase. T4 ribonucleotide reductase uses *E. coli* Grx1 as a hydrogen donor, but not thioredoxin (19). However, T4 glutaredoxin is reduced by thioredoxin via thioredoxin reductase, favoring phage-specific DNA synthesis (20).

The transcription factor activity of nuclear factor I is regulated by glutaredoxin *in vitro*, involving formation of mixed disulfides (21). Glutaredoxin is a direct target of **oncogenic** Jun protein (22). Other roles for glutaredoxin involve reduction of dehydroascorbic acid (23) and involvement in repairing proteins that have with mixed disulfides after oxidative stress (24) or in arsenite reduction (25).

## 2. Amino Acid Sequences

With few exceptions [eg, see Grx2 (26)], the glutaredoxins from bacterial viruses to humans are 9- to 12-kDa proteins that have both sequence and structural [homology](#) to the members of the thioredoxin

[superfamily](#) of proteins. Alignment of the amino acid sequences from diverse organisms (Fig. 1) gives a good impression of the sequence variability of these proteins. Unlike the thioredoxins, whose sequences align without gaps of more than one residue, numerous gaps of up to 5 residues are needed, in addition to the relatively long (up to 20 amino acid residues) N- and/or C-terminal extensions in some members. Several residues are strictly conserved among the glutaredoxins. The active site sequence –Cys11–Pro12–Tyr13–Cys14 (*E. coli* numbering) is the hallmark of the glutaredoxins. The tripeptide fragment that contains the proline residue with the *cis* [peptide bond](#) (Thr58–Val59–Pro60) and the residues Ile69–Gly70–Gly71–Tyr72–Thr73–Asp74 are also strictly conserved. Mammalian glutaredoxins are highly homologous to one another and, in addition to containing both N- and C-terminal extensions, contain additional non-structural cysteine residues that, analogous to the mammalian thioredoxins, might function in regulating glutaredoxin activity (9).

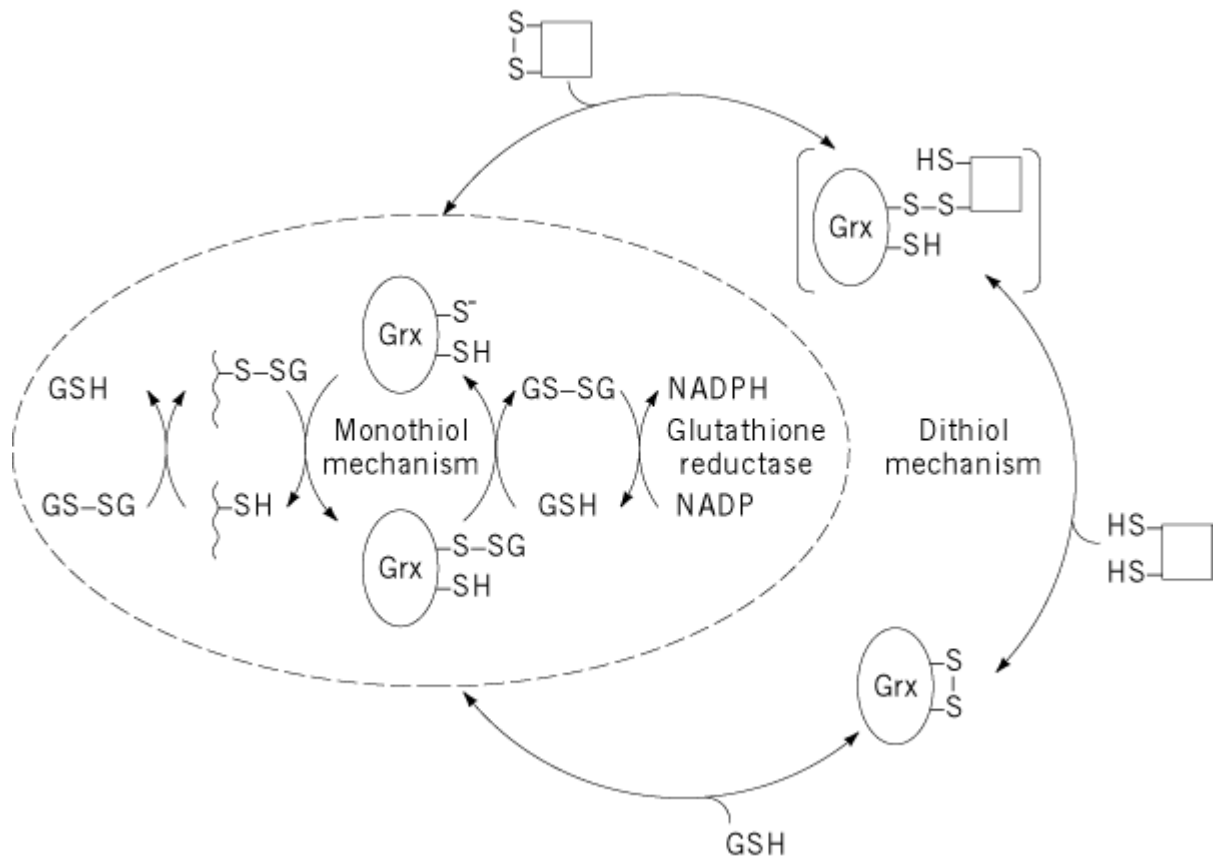
**Figure 1.** Amino acid sequence alignment of selected glutaredoxins. Sequences were taken from <sup>a</sup>(9, 32); <sup>b</sup>(33); <sup>c</sup>(34); <sup>d</sup>(acid sequences were aligned based on the three-dimensional structures of pig (29), T4 (28) and *E. coli* (Grx1) glutaredox

|                                    |                         |                         |                              |                             |         |             |
|------------------------------------|-------------------------|-------------------------|------------------------------|-----------------------------|---------|-------------|
| <b>human<sup>a</sup></b>           | ---                     | AQEFVNCKIQPGKVVVFIKP    | ---                          | TCPYCRRAQEILSQLPIKQG        | --      | LLEFV       |
| <b>pig<sup>b</sup></b>             | ---                     | AQAFVNSKIQPGKVVVFIKP    | ---                          | TCPFCRKTQELLSQLPFKEG        | --      | LLEFV       |
| <b>yeast II<sup>c</sup></b>        | MVSQETIKHVKDLIAENEIFVAS | --                      | KTYCPYCHAALNTLFEKLVPRSKVLVLQ |                             |         |             |
| <b>T4<sup>d</sup></b>              | -----                   | MFKVYGYDSNIHKCVYCDNAKRL | LLTVKKQ                      | -----                       |         | PFEFV       |
| <b><i>E.coli</i> 1<sup>e</sup></b> | -----                   | MQTVIFGRS               | ----                         | GCPYCVRAKDLAEKLSNERDDFQYQYV |         |             |
| <b>structure</b>                   |                         | $\alpha$                |                              | $\beta$                     | 10      | $\alpha$ 20 |
|                                    |                         |                         |                              |                             |         | 30          |
|                                    |                         |                         |                              |                             |         | $\beta$     |
| <b>human</b>                       | TNEIQDYLQQLTGAR         | ----                    | TVPRVFIG-KDCIGGCS            | DLVSLQ-QSGELLTRLI           |         |             |
| <b>pig</b>                         | TNEIQDYLQQLTGAR         | ----                    | TVPRVFIG-KECIGGCTD           | LESMSH-KRGELLTRLI           |         |             |
| <b>yeast II</b>                    | GADIQAALYEINGQR         | ----                    | TVPNIYI-NGKHIGGND            | LQELRETGELEELLEI            |         |             |
| <b>T4</b>                          | DDEKIAELLTKLGRDTQIGL    | TM                      | QVFPD                        | GSHIGGFDQLREYFK             |         |             |
| <b><i>E.coli</i> 1</b>             | EGITKEDLQKAGKPVE        | ---                     | TVPQIFV-DQQHIGGYTD           | FAAWVKENLDA                 |         |             |
| <b>structure</b>                   |                         | $\alpha$ 50             |                              | 60                          | $\beta$ | $\beta$ 70  |
|                                    |                         |                         |                              |                             |         | $\alpha$ 80 |
|                                    |                         |                         |                              |                             |         | $\alpha$    |

### 3. Mechanisms of Disulfide Reduction

As outlined in Fig. 2, glutaredoxin participates in reducing disulfide bonds in substrates by a dithiol mechanism that involves both thiol groups and a disulfide intermediate. Glutaredoxin also catalyzes monothiol reductions of GSH-mixed disulfides.

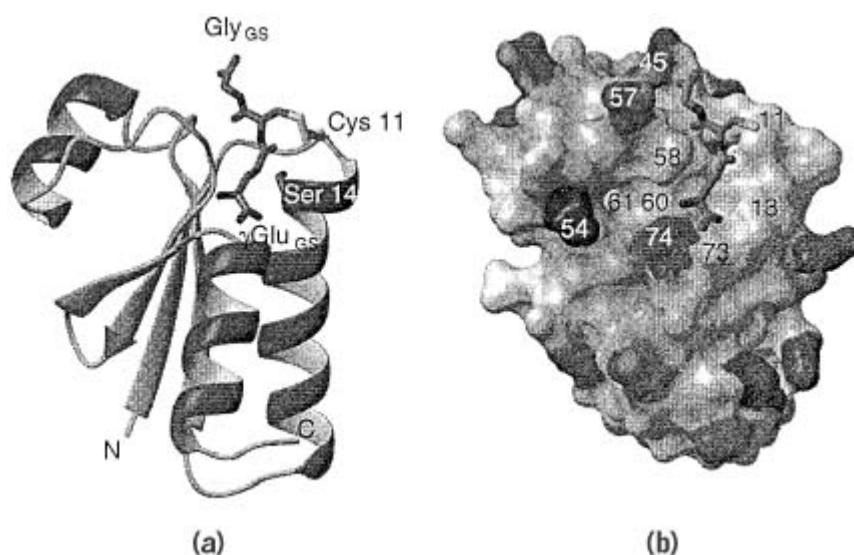
**Figure 2.** Mechanisms of glutaredoxin in GSH-disulfide oxidoreductions. The monothiol mechanism requires only the accessible N-terminal active-site Cys residue. The dithiol mechanism of disulfide reduction (like that in ribonucleotide reductase) requires both Cys residues (27, 35).



#### 4. *E. coli* Grx Structure

Three-dimensional [protein structures](#) of oxidized and reduced *E. coli* Grx-1 ([27](#)), oxidized T4 Grx ([28](#)), and oxidized pig Grx ([29](#)) have been determined. In addition, the **secondary structures** and fold topology of reduced *E. coli* Grx-3 ([15](#)) and reduced human Grx ([30](#)) have also been reported. These studies confirm the presence of the thioredoxin fold that consists of a central four-stranded mixed  $\beta$ -sheet flanked by one  $\alpha$ -helix on each face (Fig. [3](#)). In addition to this core fold, the N- and C-terminal extensions observed in sequence comparisons contain two additional helical segments.

**Figure 3.** Representative NMR solution structure of *E. coli* Grx1 (C14 S mutant) in a mixed disulfide with GSH via Cys11 ([27](#), [31](#)). (a), Cartoon of the polypeptide backbone with the side chains of Cys11 and Ser14 displayed as sticks. The glutathione ligand is also displayed as a stick model whose bonds connect non-hydrogen atoms. (b) The molecular surface colored according to the electrostatic potential (dark gray: positive, medium gray: negative, and light gray: uncharged). Residues that interact with GSH are labeled at their approximate positions.



## 5. Glutaredoxin Mixed-Disulfide Complex with GSH

The intermolecular mixed-disulfide complex between *E. coli* Grx1 and GSH has been determined by [NMR](#) techniques (27). Mutation of the C-terminal active site Cys14 residue in *E. coli* Grx1 abolishes its activity as a dithiol reductant for ribonucleotide reductase (31), but retains its activity in the monothiol mechanism (Fig. 2). Overall, the polypeptide backbone in the complex is similar to that observed in the structures of the oxidized and reduced forms, also determined by NMR. The covalently attached glutathione is found in a largely extended conformation in a cleft formed by residues from three discontinuous regions of the polypeptide chain (Fig. 3). The intimate relationship between the Grx and glutathione is also confirmed by the fact that the  $^{15}\text{N}$   $T_{1\rho}$  relaxation rates of the amide nitrogen atoms of Cys and Gly of the glutathione are similar to those observed for the Grx backbone. The floor of the cleft is made up largely of contributions from Val 59, Pro 60, and Gly 71. The sides of the cleft are formed from Thr 58 on one side and Tyr 13 and Thr 73 on the other. Interestingly, many of these residues are in positions homologous to those residues that form the hydrophobic interaction surface in the thioredoxins. Two [hydrogen bonds](#) are present in a majority of NMR conformers: HN Val 59–CO Cys<sub>GS</sub> and the HN Thr 73–aCO<sub>2</sub> gGlu<sub>GS</sub>. Additional hydrogen bonds and [salt bridges](#) are observed in several structures, including  $\epsilon\text{NH}_3$  Lys 45–aCO<sub>2</sub> Gly<sub>GS</sub>, aNH<sub>3</sub> gGlu<sub>GS</sub>–gCO<sub>2</sub> Glu 74, aCO<sub>2</sub> gGlu<sub>GS</sub>–HN Tyr 72, aCO<sub>2</sub> gGlu<sub>GS</sub>–HN Thr 73, and aCO<sub>2</sub> gGlu<sub>GS</sub>–HO Thr 73. In addition, the aCO<sub>2</sub> gGlu<sub>GS</sub> in a number of structures can be positioned directly over the N-terminus of the C-terminal  $\alpha$ -helix, allowing for a favorable interaction with the dipole of the helix. Several of these interactions (salt bridge to C-terminus of GS, H-bond from CO Cys<sub>GS</sub> to residue before the strictly conserved *cis* Pro, NH and OH of Ser to aCO<sub>2</sub> gGlu<sub>GS</sub>) have homologous counterparts in the recognition of glutathione by **glutathione S-transferases**. These observations suggest that the details of the structural basis of noncovalent glutathione binding to Grx may be representative of a wider range of substrate interactions between members of the thioredoxin superfamily.

### Bibliography

1. A. Holmgren (1976) Proc. Natl. Acad. Sci. USA **73**, 2275–2279.
2. A. Holmgren (1979) J. Biol. Chem. **254**, 3664–3671.
3. A. Holmgren (1979) J. Biol. Chem. **254**, 3672–3678.
4. J.-O. Höög, H. Jörnvall, A. Holmgren, M. Carlquist, and M. Persson (1983) Eur. J. Biochem.

136, 223–232.

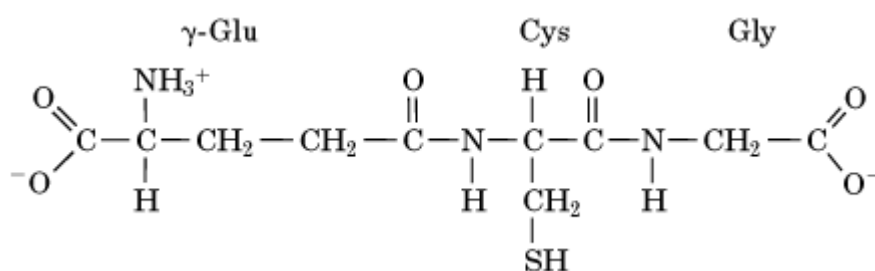
5. F. Åslund, B. Ehn, A. Miranda-Vizuete, C. Pueyo, and A. Holmgren (1994) *Proc. Natl. Acad. Sci. USA*, **91**, 9813–9817.
6. S. Luikenhuis, G. Perrone, J. W. Dawies, and C. M. Grant (1998) *Mol. Cell. Biol.* **9**, 1081–1091.
7. K. Minakucchi, T. Yabushita, T. Masumura, K. Ichihara, and K. Tanaka (1994) *FEBS Lett.* **337**, 157–160.
8. M. Luthman, S. Eriksson, A. Holmgren, and L. Thelander (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2158–2162.
9. C. A. Padilla, E. Martinez-Galisteo, J. A. Bàrcena, G. Spyrou, and A. Holmgren (1995) *Eur. J. Biochem.* **227**, 27–34.
10. B.-M. Sjöberg and A. Holmgren (1972) *J. Biol. Chem.* **247**, 8063–8068.
11. B. Y. Ahn and B. Moss (1992) *Proc. Natl. Acad. Sci. USA* **89**, 7060–7064.
12. B. O. Gvakharia, E. Hanson, EK. Koonin, and C. K. Mathews (1996) *J. Biol. Chem.* **271**, 15307–15310.
13. M. Russel and A. Holmgren (1988) *Proc. Natl. Acad. Sci. USA* **85**, 990–994.
14. A. Miranda-Vizuete, A. Rodrigues-Ariza, A. Holmgren, J. Lopez-Barea, and C. Pueyo (1996) *J. Biol. Chem.* **271**, 19099–19103.
15. F. Åslund, K. Nordstrand, K. D. Berndt, M. Nikkola, T. Bergman, H. Ponstingl, H. Jörnvall, G. Otting, and A. Holmgren (1996) *J. Biol. Chem.* **271**, 6736–6745.
16. F. Åslund, K. D. Berndt, and A. Holmgren (1997) *J. Biol. Chem.* **272**, 30780–30786.
17. M. Russel, P. Model, and A. Holmgren (1990) *J. Bacteriol.* **172**, 1923–1929.
18. M. Zheng, F. Åslund, and G. Storz (1998) *Science*, **279**, 1718–1721.
19. A. Holmgren (1978) *J. Biol. Chem.* **253**, 7424–7430.
20. O. Berglund and A. Holmgren (1975) *J. Biol. Chem.* **250**, 2778–2782.
21. S. Bandyopadhyay, D. W. Starke, J. J. Mieyal, and R. M. Gianostajski (1998) *J. Biol. Chem.* **273**, 392–397.
22. M. E. Goller, J. S. Jacovoni, P. K. Vogt, and U. Krause (1998) *Oncogene* **16**, 2945–2948.
23. W. W. Wells, D. P. Xu, Y. Yang, and P. A. Rocque (1990) *J. Biol. Chem.* **265**, 15361–15364.
24. C. M. Jung and J. A. Thomas (1996) *Arch. Biochem. Biophys.* **335**, 61–72.
25. J. Liu and B. P. Rosen (1997) *J. Biol. Chem.* **272**, 21084–21089.
26. A. Vlamis-Gardikas, F. Åslund, G. Spyrou, T. Bergman, and A. Holmgren (1997) *J. Biol. Chem.* **272**, 11236–11243.
27. J. H. Bushweller, M. Billeter, A. Holmgren, and K. Wüthrich (1994) *J. Mol. Biol.* **235**, 1585–1597.
28. B. O. Söderberg, B.-M. Sjöberg, U. Sonnerstam, and C.-I. Brändén (1978) *Proc. Natl. Acad. Sci. USA* **75**, 5827–5830.
29. S. K. Katti, A. R. Robbins, Y. Yang, and W. W. Wells (1995) *Protein Sci.* **4**, 1998–2005.
30. C. Sun, A. Holmgren, and J. H. Bushweller (1997) *Protein Sci.* **6**, 383–390.
31. J. H. Bushweller, F. Åslund, K. Wüthrich, and A. Holmgren (1992) *Biochemistry* **31**, 9288–9293.
32. M. R. Fernando, H. Sumimoto, H. Nanri, S. Kawabata, S. Iwanaga, S. Minakami, Y. Fukumaki, and K. Takeshige (1994). *Biochim. Biophys. Acta* **1218**, 229–231.
33. Y. Yang, Z. R. Gan, and W. W. Wells (1989) *Gene* **83**, 339–346.
34. C. P. Hollenberg, U. Kleinans, K. Luetzenkirchen, M. R. Rad, and G. Xu (1992) *EMBL/gen bank DDBJ databanks*.
35. S. A. Gravina and J. J. Mieyal (1993) *Biochemistry*, **32**, 3368–3376.

## Suggestions for Further Reading

36. A. Holmgren and F. Åslund (1995) Glutaredoxin, *Methods Enzymol.* **252**, 283–292.
37. *Methods in Enzymology*, Vol. **252** (1995) *Biothiols. Part B*. Glutathione and thioredoxin: Thiols in signal transduction and gene regulation.
38. J. A. Fuchs (1989) "Glutaredoxin". In *Glutathione, Chemical, Biochemical and Medical Aspects* (D. Dolphin, O. Avramovic, and R. Poulson, eds.), Part B. Wiley, New York, pp. 551–570.
39. A. Holmgren, C.-I. Brändén, H. Jörnvall and B.-M. Sjöberg (eds.) (1986) *Thioredoxin and Glutaredoxin Systems. Structure and Functions*, Raven Press, New York.

## Glutathione

Glutathione was first discovered by J. de Rey-Pailhade over 100 years ago, and its structure (ie, L-g-glutamyl-L-cysteinyl-glycine, or g-Glu-Cys-Gly),



was deduced in the 1930s. It has several important functions: (i) most importantly, as a reductant; (ii) conjugation to foreign chemicals to make them more **water-soluble**, less toxic, and excretable; (iii) transport of **amino acids** across cell **membranes**; (iv) incorporation into some leukotriene structures; (v) as a **cofactor** for some **enzymes**; (vi) as a nontoxic storage form of **cysteine**; and (vii) in forming and maintaining protein **disulfide bonds** in the **endoplasmic reticulum** (ER).

Glutathione is the predominant **thiol** compound in very many cells, both prokaryotes and eukaryotes, where its total intracellular concentration is usually in the range 0.5 to 12 mM. It is present in most eukaryotes, except for those that do not have mitochondria. It is not present in many Archaeobacteria, but in halobacteria it is replaced by g-Glu-Cys. Likewise, some eubacteria do not have glutathione, and it is found primarily in the bacteria of the cyanobacteria and the purple bacteria classes. It is believed that glutathione arose initially in the prokaryotic ancestors of mitochondria and chloroplasts and was acquired by eukaryotes along with these organelles.

Glutathione is synthesized from its three constituent amino acids. First, the side-chain **g-carboxyl group** of **glutamic acid** is joined in a **peptide bond** to the **α-amino group** of **cysteine**, by the enzyme g-glutamylcysteine synthetase. The **α-carboxyl group** of the resulting g-Glu-Cys is then joined in a normal peptide bond with the **α-amino group** of **glycine**, catalyzed by glutathione synthetase. Both steps are coupled to the hydrolysis of ATP. Some species of **plants** have a variant of glutathione, homoglutathione, in which the glycine residue is replaced by β-alanine, due to an **evolutionary** change in specificity of the enzyme glutathione synthetase. In **trypanosomes**, glutathione is largely replaced by the related trypanothione in which two g-Glu-Cys-Gly moieties are linked by peptide bonds through their carboxyl groups to a molecule of **spermidine**.

Glutathione is also degraded to its constituent amino acids *in vivo*, which is believed to give it an important role as a nontoxic storage form of cysteine. High levels of cysteine are toxic in some systems, in part because its amino, carboxyl, and thiol groups are well situated stereochemically to chelate metal ions, which catalyze air oxidation of its thiol group. A by-product of thiol oxidation is peroxides, which damage the cell. Glutathione is much more resistant to air oxidation and much less toxic, and it can occur safely at 10 to 100 times greater levels.

Glutathione is an essential cofactor for a number of enzymes, including formaldehyde dehydrogenase, glyoxylase, maleylacetoacetate isomerase, dehydrochlorinase, and prostaglandin endoperoxidase isomerase. The glutathione is involved transiently in the reactions they catalyze. For example, formaldehyde dehydrogenase uses  $\text{NAD}^+$  to catalyze the oxidation of an adduct of formaldehyde and glutathione to produce *S*-formylglutathione, which is subsequently hydrolyzed by a specific lyase to yield formate and regenerate the glutathione. Glyoxylase is similar, also generating a transient intermediate adduct between its substrate, methyl glyoxal, and glutathione.

Naturally, glutathione is found in a mixture of the thiol form, GSH, and with a [disulfide bond](#) linking two glutathione moieties, GSSG. The GSH thiol group has a  $\text{p}K_a$  value of about 8.8, while the amino and carboxyl groups have normal  $\text{p}K_a$  values that are virtually the same in the disulfide form (1).

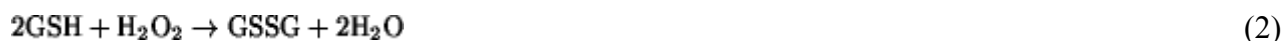
There apparently are no substantial interactions between the two glutathione moieties in GSSG, and the stability of this disulfide bond is comparable to that between two other cysteine residues on unfolded peptides or proteins; the rate and equilibrium constants for [thiol–disulfide exchange](#) are close to their expected values (2).

Within the cytosol of cells, glutathione is kept primarily in the reduced GSH form, largely by reduction of GSSG by the enzyme *glutathione reductase*, at the expense of NADPH:



The mechanism of this reaction is complex, as the enzyme contains a pair of cysteine residues that reversibly form a disulfide bond, plus one FAD molecule (3). As GSH and GSSG will react spontaneously with each other and with other thiol and disulfide compounds of the cell, in the [thiol–disulfide exchange](#) reaction, the ratio of GSH and GSSG reflects the general redox properties of the cell. Being present at such high concentrations, GSH is an important reductant within cells. For example, it is involved in the reduction of [glutaredoxin](#), and it and NADPH reduce the oxidized form of vitamin C, dehydroascorbate, to the active form, ascorbate.

GSH is used for the destruction of oxidants, such as organic peroxides and  $\text{H}_2\text{O}_2$ , which would cause irreversible damage to membranes, DNA, and numerous other cellular components; this reaction is catalyzed by the enzyme *glutathione peroxidase*:



This enzyme is unusual in that it has a [selenocysteine](#) residue at its active site (4).

Normally, GSH predominates over GSSG by a factor of 100 in the cytosol. This is a somewhat reducing environment that tends to keep protein cysteine residues in the thiol form ( $\text{P}_{\text{SH}}^{\text{SH}}$ ) and destabilizes any disulfide bonds they may make ( $\text{P}_{\text{S}}^{\text{S}}$ ),





although very stable disulfide bonds found in some proteins could be present under these conditions. The stability of a protein disulfide bond is determined by its intrinsic stability,  $K_{eq}$ , plus the ratio of the concentrations of GSH and GSSG:

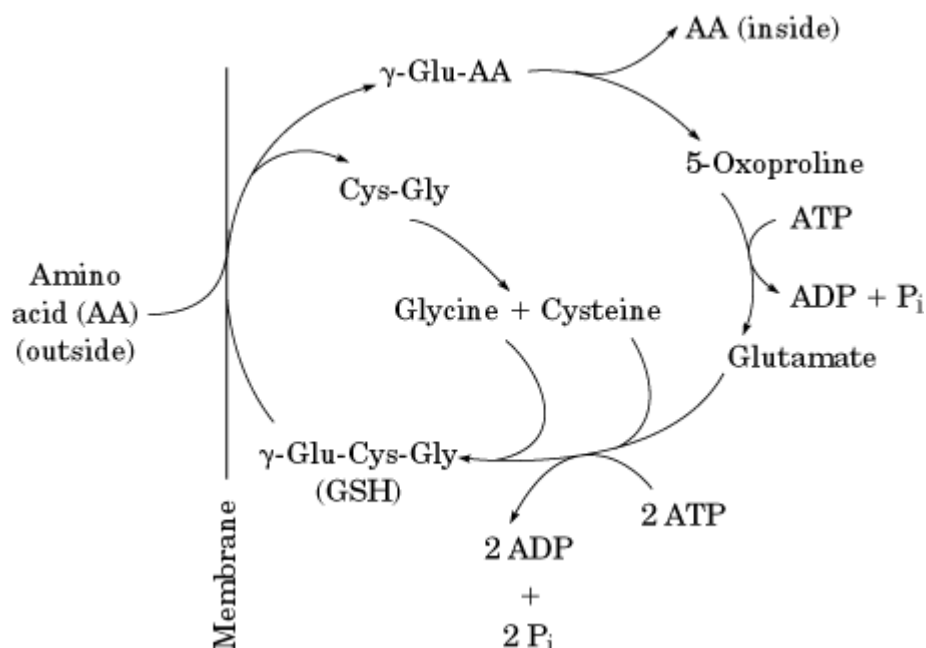
$$\frac{[P_S^S]}{[P_{SH}^{SH}]} = K_{eq} \frac{[GSSG]}{[GSH]^2} \quad (4)$$

Only very stable protein disulfide bonds, with large values of  $K_{eq}$ , could be populated in the cytosol.

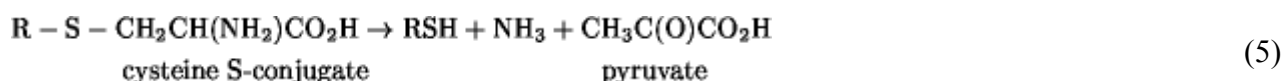
In the endoplasmic reticulum (ER) of eukaryotes, in contrast, GSH and GSSG are present in nearly equal concentrations (5). With 1 mM each of GSH and GSSG, the same disulfide bond would be present in the ER to an extent 100 times greater than in the cytosol. This is sufficient to permit proteins in the ER to form disulfide bonds if their conformations stabilize them (see [Disulfide Bonds](#)); if not, the protein cysteine residues will tend to remain in the thiol form. The formation of such stable disulfide bonds in proteins in the ER is catalyzed by [protein disulfide isomerase](#) (PDI). PDI can have an unstable disulfide bond at each of its two active sites, which is readily transferred to other proteins. The PDI active-site disulfide bonds are believed to be regenerated by chemical reaction of the protein with GSSG.

Glutathione is believed to be used in some cells, such as those of the mammalian liver and kidney, for the transport of amino acids (except for [proline](#)) across the plasma membrane (Fig. 1). The extracellular amino acid is reacted with intracellular glutathione that has been shuttled to the cell surface by the integral [membrane protein](#) g-glutamyltranspeptidase. This replaces the Cys-Gly dipeptide of GSH by the amino acid and produces g-Glu-amino acid. The amino acid is liberated upon degradation of the g-Glu-amino acid by the enzyme g-glutamyl cyclotransferase, which generates 5-oxoproline instead of glutamate. The driving force for uptake of amino acids by this “group translocation” type of mechanism is the destruction of the glutathione, for the Cys-Gly moiety is also degraded, by cysteinylglycine dipeptidase; regeneration of glutamate from 5-oxoproline by 5-oxoprolinase requires hydrolysis of ATP, and resynthesis of glutathione requires another two molecules of ATP. The disulfide form of cysteine, cystine, is an especially reactive amino acid substrate of this cycle. The g-Glu-CysSSCys generated by g-glutamyltranspeptidase is rapidly reduced in the cell cytosol to g-Glu-Cys and cysteine. The g-Glu-Cys can then be converted to GSH by addition of glycine, thereby bypassing the production of 5-oxoproline. On the other hand, humans lacking g-glutamyltranspeptidase have no difficulty in transporting amino acids, and there are many other types of amino acid **transporters**. The cycle is now believed to play only a minor role in transport of amino acids, except for cysteine.

**Figure 1.** The g-glutamyl cycle for transporting amino acids, coupled to the degradation of glutathione (GSH; g-Glu-Cys-Gly). The amino acid (AA) on the outside of the cell (left) is coupled to glutamic acid of GSH by the membrane protein g-glutamyl transpeptidase, by replacing the dipeptide cysteinyl-glycine (Cys-Gly). The amino acid is released from the g-glutamyl residue, converting it to 5-oxoproline. The concomitant degradation of GSH is the driving force for transporting the amino acid. Resynthesizing 1 mol of GSH requires the hydrolysis of 3 moles of ATP.



The g-glutamyl cycle is now believed to have a major role in higher organisms in the metabolism of leukotrienes, estrogens, prostaglandins, and many drugs and xenobiotics. All these molecules, often after their reaction with [cytochrome P-450](#), are linked covalently to the thiol group of GSH by the large enzyme family of **glutathione S-transferases**. The resulting glutathione S-conjugate is then degraded to the Cys-Gly S-conjugate and then to the cysteine S-conjugate by g-glutamyltranspeptidase and cysteinylglycine dipeptidase, respectively. The cysteine S-conjugate is converted, by a specific *N*-acetyltransferase to the corresponding mercapturate, *N*-acetylcysteine S-conjugate, which is then excreted. On the other hand, this pathway can also sometimes lead to the generation of a toxic substance. For example, halogenated cysteine S-conjugates can generate highly reactive thiol-containing fragments through the action of cysteine S-conjugate b-lyases:



where RSH is the thiol form of the xenobiotic, which can be toxic.

The generally protective, antioxidative effects of GSH are reflected in the observations that decreased levels of GSH lead to clinical symptoms in humans and decreased resistance to viruses and other diseases.

#### Bibliography

1. G. Jung, E. Beitmeier, and W. Voelter (1972) *Eur. J. Biochem.* **24**, 438–445.
2. N. J. Darby and T. E. Creighton (1995) *Biochemistry* **34**, 3576–3587.
3. P. A. Karplus and G. E. Schulz (1989) *J. Mol. Biol.* **210**, 163–180.
4. O. Epp et al. (1983) *Eur. J. Biochem.* **133**, 51–69.
5. C. Hwang, A. J. Sinskey, and H. F. Lodish (1992) *Science* **257**, 1496–1502.

#### Suggestions for Further Reading

6. A. Meister and M. E. Anderson (1983) *Glutathione*. *Ann. Rev. Biochem.* **52**, 711–760.
7. N. Taniguchi, T. Higashi, Y. Sakamoto, and A. Meister, eds. (1989) *Glutathione Centennial: Molecular Perspectives and Clinical Implications*. Academic Press, San Diego.



|                             |       |      |     |     |        |        |       |       |   |   |    |
|-----------------------------|-------|------|-----|-----|--------|--------|-------|-------|---|---|----|
| Aminochrome                 | 0.04  | 0.03 | —   | 0.8 | 150    | 0.14   |       | 0.08  | — |   |    |
| 1-Chloro-2,4-dinitrobenzene | 80    | 80   | 7.5 | 180 | 220    | 7      | 1.4   | 105   | — | — | 24 |
| Cumene hydroperoxide        | 10    | 10   | 1   | 0.6 |        |        |       | 0.03  | 3 | 7 | 3  |
| 4-Hydroxynon-2-enal         | 6     |      | 190 | 3   | 4      | 2      |       | 2     |   |   |    |
| <i>trans</i> -Stilbeneoxide | 0.002 |      |     | 5.2 | 0.0003 | 0.0004 | 0.003 | 0.002 |   |   |    |

<sup>a</sup> Specific activities (mol/min per mg pure protein) are given for representative substrates. Dashes represent no measurable activity; blank spaces signify missing data. Data largely compiled from the review articles by Mannervik and Widersten (1995) and Hayes and Pulford (1995).

## 1. Reactions Catalyzed

GSTs catalyze reactions between [glutathione](#) (GSH) and compounds with electrophilic chemical functionalities. GSH is a tripeptide, g-L-glutamyl-L-cysteinyl-glycine, that occurs in millimolar concentrations in all types of cells and in micromolar concentrations in mammalian extracellular fluids. The reactive [thiol group](#), -SH, of the GSH [cysteine](#) residue may serve as either a nucleophile or a reductant. In both cases, the reactions are dependent on the electron-donating properties of the sulfur atom. Its reactivity is further accentuated by ionization of the thiol group to the thiolate form. One of the catalytic functions of GSTs is to decrease the  $pK_a$  value of the GSH thiol group (9.2 in aqueous solution) to make GSH ionized under physiological conditions. The  $pK_a$  value of enzyme-bound GSH may be lowered by several pH units. Judging from the known affinities of GSTs for GSH, it would appear that most enzymes are fully loaded with GSH in its thiolate form and ready to react with the incoming second substrate.

The common denominator for essentially all GST-catalyzed reactions is nucleophilic attack by the GSH thiolate group on an electrophilic site in the second substrate. The electrophilic center of the second molecule is usually a carbon atom, but it may also be electronegative atoms, such as oxygen, nitrogen, or sulfur. The GST-catalyzed reactions involving O, N, and S as electrophilic sites often lead to unstable products, which react with a second molecule of GSH:



where ROOH and ROH are, respectively, a hydroperoxide and the corresponding alcohol, and GSSG is the **disulfide** form of glutathione. In effect, such reactions accomplish reduction of organic hydroperoxides, nitrate esters, and disulfides.

The carbon-centered reactions are basically of two types: additions and substitutions. Addition reactions involve molecules such as naturally occurring  $\alpha,\beta$ -unsaturated carbonyl compounds and organic isothiocyanates. The substitutions reactions involve the replacement of a chemical substituent by a glutathione group. The leaving group may be a halogenide, sulfate, or other chemical substituent with sufficient electron-withdrawing potential.

Reactions involving electrophilic carbon centers usually lead to GSH conjugates that may be excreted from the cell via ATP-dependent membrane-associated transporters, such as the multidrug resistance-associated protein (MRP) (see [Drug Resistance](#)), for subsequent disposition through the gut or the urinary tract. For urinary excretion, GSH conjugates are metabolically transformed by cleavage of the two peptide bonds of the glutathione moiety, followed by acetylation into the corresponding *N*-acetyl-cysteine derivatives, mercapturic acids. These reactions are accomplished through the action of  $\gamma$ -glutamyl transpeptidase, a dipeptidase, and an acetyltransferase (see [Glutathione](#)).

## 2. Discovery

GSTs were discovered in mammalian liver as enzymes catalyzing aromatic nucleophilic substitution reactions in which GSH replaces a halogen atom in an aromatic compound, such as a chloronitrobenzene derivative (1, 2). When other classes of organic compounds, with epoxide, alkene, and other electrophilic centers, were investigated, additional enzyme-catalyzed GSH conjugations were identified (3). The introduction of improved separation techniques led to the identification of multiple forms of GST catalyzing the same chemical reaction, i.e. **isoenzymes** in the original sense of the term (4, 5). Subsequent work with purified enzymes demonstrated that although each enzyme form has its characteristic substrate specificity profile, most GSTs catalyze reactions with a wide range of electrophiles, and the majority of the active compounds are at least to some degree substrates for several GSTs.

## 3. Substrates

GSTs as a family of enzymes are capable of catalyzing reactions involving literally thousands of electrophilic compounds, including many of xenobiotic origin. In this respect, the GSTs display similarities to **antibodies**, which may evolve to interact with and provide protection against an almost unlimited number of molecular structures.

The majority of the known GSTs are regarded as detoxication enzymes and are catalytically active with a large number of xenobiotics, including epoxides of potent **carcinogens**, such as benzo[a]pyrene and aflatoxin. Naturally occurring substrates include a wide variety of genotoxic compounds, such as epoxides, activated alkenes, hydroperoxides, and quinones, all products of oxidative metabolism. Organic isothiocyanates are abundant in edible plants, from which they are released in high concentrations by injuries caused by insects and microbial infections.

## 4. Purification

[Affinity chromatography](#) methods have been developed for the purification of GSTs. Most purification procedures involve affinity matrices based on glutathione derivatives as ligands. Immobilized *S*-hexylglutathione, first used for the purification of glyoxalase I, is linked via the  $\alpha$ -amino group of the  $\gamma$ -glutamyl residue of glutathione to a suitable matrix (6). This adsorbent potentially combines the affinities of the H subsite of the GST [active site](#) (see below) for the **hydrophobic** hexyl substituent and of the G subsite for the GSH moiety. Another commonly used affinity matrix makes use of GSH immobilized via its thiol group (7). Elution of bound GSTs may be effected by competition using GSH derivatives with affinity for the active site of the enzyme. Alternative elution procedures involve changing the pH of the eluent to extreme values, such as pH 10 or pH 2. When applied to crude tissue fractions, affinity chromatography based on these matrices will yield mixtures of GSTs, which can be resolved by other methods. For separation of the dimeric proteins in a functional form, [chromatofocusing](#) or [ion-exchange chromatography](#) may be used. For analytical purposes, the subunits may be resolved by **reversed-phase** high-performance liquid chromatography (HPLC).

Less specific affinity methods used for the purification of GSTs include Orange A dye

chromatography and **metal chelate chromatography**. The latter methods have proved useful for the isolation of the several GSTs that do not bind to affinity matrices based on glutathione derivatives.

## 5. Nomenclature

The International Enzyme Commission has recommended the name RX:glutathione R-transferase (E.C. 2.5.1.18) or the trivial form glutathione transferase, which has been adopted here. Another designation commonly used is glutathione *S*-transferase, but this is inconsistent with the rational name, since the group transferred is not the sulfur of glutathione. The generally accepted abbreviation GST has probably contributed to the persistence of the *S*-prefix.

In the time period when the first GSTs were discovered, the enzymes were named in accord with the reactions used for their identification: glutathione *S*-aryltransferase, glutathione *S*-alkyltransferase, and so on. With the exception of the membrane-bound transferase catalyzing the glutathione conjugation of leukotriene A<sub>4</sub>, named leukotriene C<sub>4</sub> synthase, the principle of designating GSTs by their substrates or products has been abandoned, because most enzymes have very broad substrate specificities so that accurate distinctions cannot be made.

Another nomenclature system arose from the finding that the cytosolic “Y-fraction” of rodent liver could be separated by [SDS-PAGE](#) into components that corresponded to GST subunits. These subunits were distinguished by lower indices, namely Y<sub>a</sub>, Y<sub>b</sub> and so on, and were further subdivided into Y<sub>a1</sub>, Y<sub>a2</sub>. Subunits with the same mobility but different structures were identified.

A rational nomenclature for the cytosolic or soluble GSTs is based on two principles. First, the native enzymes occur as binary combinations of GST protein subunits, so that the functional properties of the dimeric protein reflect those of its constituents subunits (8). Second, the GSTs can be divided into different classes primarily based on similarities in their amino acid sequences (9). In mammals, seven classes of soluble GSTs have been identified thus far: Alpha, Mu, Pi, Theta, Sigma, Kappa, and Zeta. Sequence identities within a class are normally >50%, but they may exceed 95%. Classification of GSTs from other biological species requires further structural studies, but designations for insect enzymes (Delta) and bacteria (Beta) have been used.

In the current nomenclature (Table 1), the soluble mammalian GSTs are denoted by a capital Roman letter indicating the class (A, M, P, T, S, K, and Z) as well as two hyphenated Arabic numerals showing their subunit composition (10). Thus, GST A4-4 is a homodimer of subunit 4 within the Alpha class, and GST A1-2 is a heterodimer composed of subunits 1 and 2 from the same class. There may also be reason to distinguish allelic variants: Human GST M1a-1b is a heterodimeric member of the Mu class composed of variants of subunit 1 in the Mu class encoded by the *GST M\*A* and *GST M\*B* alleles. A prefix may be used to identify the biological species from which the enzyme derives, such as “h” for the major human enzyme hGST A1-1.

The mammalian membrane-bound GST first isolated from the microsome fraction is a trimeric enzyme composed of identical subunits and cannot readily be accommodated within the nomenclature system used for the soluble GSTs. It is simply referred to as microsomal GST (MGST1). The membrane-bound GST specifically catalyzing the conjugation of leukotriene A<sub>4</sub> is referred to as leukotriene C<sub>4</sub> synthase.

## 6. Structure

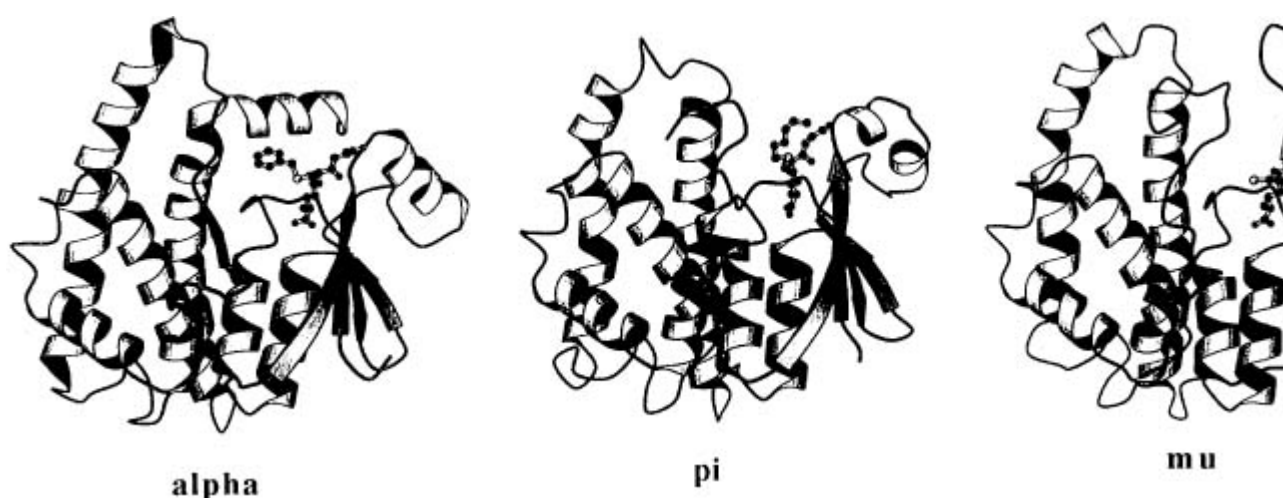
Three distinct categories of GST [protein structures](#) are known.

### 6.1. Soluble GSTs

The first category includes the soluble GSTs, which are generally dimers of equal-sized subunits,

each containing an active site composed of a pocket for glutathione (the G site) and a pocket for the largely **hydrophobic** electrophilic substrate (the H site) (11). Each subunit consists of two structural **domains**; (a) an *N*-terminal domain including the first one-third of the primary structure and (b) a domain formed essentially by the remaining two-thirds of the amino acid sequence (Fig. 1). The first domain is folded into a mixed **b-sheet** flanked by **a-helices** and provides the structural basis for the G site. The second one is formed by helical segments and provides the major contributions to the H site. The fold of the *N*-terminal domain is essentially the same as the fold of thioltransferase ([glutaredoxin](#)) and selenium-dependent glutathione peroxidase. Another salient feature of the dimer is a deep cleft between the subunits that may serve as a binding site in addition to the active site cavities. The fully formed dimer appears to be a structural requirement for a catalytically active GST.

**Figure 1.** Schematic drawings of alpha (GST A1-1), pi (GST P1-1) and mu (GST 3-3) structures with their ligands. The been chosen to show the 2 domains and C-terminal region. From Ref. (11), with permission.



All known structures contain an *N*-capping box sequence, (Ser/Thr)–X–X–Asp, as well as a hydrophobic staple motif formed by flanking amino acid residues in the core of the folded structure. As evidenced by mutational analyses of protein folding (12), these structural signatures appear to play an important role in the nucleation and orientation of a centrally located  $\alpha$ -helix.

## 6.2. Manganese-Dependent Plasmid-Encoded GST

The second category is represented by the manganese-dependent plasmid-encoded GST that is active with the antibiotic fosfomycin (13). It is a homodimer structurally related to mammalian glyoxalase I, a GSH-linked zinc protein whose structure has been determined (14). This GST is structurally unrelated to other known GSTs.

## 6.3. Membrane-Bound GSTs

The third category contains the microsomal MGST1 and, presumably, the additional membrane-bound GSTs. The detailed structure of MGST1 is unknown, but clearly it is not similar in structure to the soluble GSTs. The protein is a homotrimer, and the identical monomers appear to have a high content of  $\alpha$ -helical structure.

## 7. Catalytic Mechanism

GSTs bind the peptide moiety of GSH in a well-defined manner so that they are specific for this thiol substrate. For all GSTs investigated, the  $\alpha$ -carboxylate group of the  $\gamma$ -glutamyl residue of GSH is

required for the thiol group to serve as a substrate; when bound in the G subsite, the carboxylate appears to have a negative charge that is not balanced by a positively charged protein group. These observations suggest a contribution of the substrate carboxylate group to catalysis. At least for some GSTs, there is also evidence that binding of GSH induces a conformational change in the protein that promotes binding of the second substrate.

Binding of GSH to the G subsite affects the ionization of its sulfhydryl group, thereby making it more reactive. In addition, interactions with a **hydrogen-bond** donor in the active site may stabilize and orient the thiolate group for optimal interaction with the electrophilic substrate. In this manner the enzyme promotes catalysis by selecting a suitable ground-state conformation of GSH.

Binding to the H subsite also results in desolvation of the electrophilic substrate and selection of an orientation favorable for reaction with the sulfur of GSH. Some GSTs have [tyrosine](#) or other residues that afford binding specificity, as well as activation of substrates by hydrogen bonding and other [electrostatic interactions](#).

In some, but not all, GST-catalyzed reactions, stabilization of the [transition state](#) contributes to the catalytic efficiency. Thus, the common mechanistic feature of GSTs is the activation of GSH, whereas other components of the catalytic process may differ from substrate to substrate and from enzyme to enzyme.

## 8. Binding Function

A ligand-binding function of GSTs of possible physiological significance was originally found in the protein fraction Y, of rat liver cytosol, which binds organic compounds such as certain carcinogens, dyes, and corticosteroids. The abundance of the protein (several percent of the total cytosolic protein) and its broad specificity for binding of a variety of ligands suggested a function similar to that of [serum albumin](#) in blood plasma and led to the name of “ligandin” (intracellular albumin). In the current nomenclature, ligandin corresponds to the Alpha class GST A1-1, with some contribution of the heterodimer GST A1-2. Corticosteroid binding has also been associated with GST isoenzymes of the Mu class.

The binding function of GSTs has often been ascribed to the H subsite, at which the enzyme binds the various electrophilic substrates. Some nonsubstrate compounds, especially small molecules, can bind at the H subsite with a stoichiometry of one molecule per subunit. However, [X-ray crystallography](#) studies have also shown that the deep cleft between the two GST subunits serves as an additional binding site of the dimeric structure, explaining the stoichiometry of one ligand per two subunits observed for certain ligands.

As binding proteins, GSTs have been suggested to facilitate the intracellular transport of endogenous molecules, such as heme and bilirubin, as well as other compounds of low solubility in an aqueous environment. The transport through cellular membranes of dyes, such as indocyanine green and bromosulphophthalein, which have been used in clinical liver function tests, may also be facilitated by GSTs. An auxiliary role in **active-transport** processes involving MRP and other membrane transport systems may also be suggested. However, the full significance of GSTs as binding proteins needs to be further clarified.

## 9. Molecular Evolution

The presence of GSH is a signature of aerobic organisms ([15](#)), as is the occurrence of GSTs. Elucidation of the details of the [phylogeny](#) of GSTs requires additional structural information. The soluble (cytosolic) GSTs are found ubiquitously, however, and all their three-dimensional structures determined, from bacteria to higher plants and mammals, have the same two-domain fold, even though sequence similarities may be undetectable. A common ancestry followed by [divergent evolution](#) is therefore indicated. The most recently discovered Kappa, Theta, and Zeta classes of



GST appear to be evolutionary precursors of the more abundant mammalian enzymes from the Alpha, Mu, and Pi classes. It has been proposed, in accordance with the symbiont theory of organelle evolution, that the Kappa class GST found in [mitochondria](#) may be a descendant from a bacterial GST. The structural similarities between the *N*-terminal domain of the soluble GSTs and other GSH-linked proteins (see text above) lends support to the hypothesis that a GSH-binding protein is a functional unit that has been combined with other functional protein units and evolved to serve diverse biological functions. The GST structure itself has evolved to fulfill alternative functions, such as that of crystallins in the eye lens of cephalopods. The various GST isoenzymes may have arisen by [gene duplications](#), [recombination](#) of DNA segments, and random [mutations](#) .

The membrane-bound GSTs and the plasmid-encoded Mn-dependent GST seem to have undergone [convergent evolution](#) to acquire the property of catalyzing the nucleophilic attack of GSH on electrophiles. Thus, three evolutionarily unrelated superfamilies of GSTs have been recognized.

## 10. Differential Tissue Distribution

GSTs are found in all tissues investigated in multicellular organisms, but the expression of the multiple genes differs from tissue to tissue. Also, the enzyme distribution changes with time during the development of embryonic to adult organs. Some GSTs have a reasonably general tissue distribution, whereas others have a more restricted occurrence. Also intracellular spatiotemporal differences in GST distribution have been noted. The soluble GSTs are found mainly in the cytoplasm, but may also occur sometimes in the nucleus. The mammalian Kappa class GST is found in the mitochondrial matrix.

The differential distribution of GSTs may have toxicological consequences, since the amount and nature of a given GST will determine the capacity of a cell to resist insults from genotoxic and carcinogenic electrophiles. The presence or absence of a particular enzyme will influence the resistance **phenotype** of a tissue. In cancer cells, GSTs may contribute to drug resistance, particularly against alkylating cytostatic agents.

## 11. GST Genes

Mammalian cells generally appear to have more than 20 GST genes, each encoding a distinct GST subunit. In humans, as well as in rodents, the genes have been shown to be distributed on different [chromosomes](#), and the loci for related GSTs from the same class are generally clustered on the same chromosome. In addition, [pseudogenes](#) may occur close to the active genes. This distribution supports the classification adopted for the GSTs and is in agreement with the proposal that they have diversified in evolution by gene duplication and DNA recombination.

Genetic polymorphisms are known, and the human genes for GST M1-1 and GST T1-1 are deleted in approximately 50% and 20%, respectively, of the population. The resulting enzyme deficiencies lead to increased sensitivities to gene modifications and chromosome aberrations by certain chemical agents, and they probably also lead to increased risk of contracting certain forms of cancer. On the other hand, the activity of GST T1-1 promotes the formation of mutagens from ethylene dihalogenides, and the absence of the enzyme may be advantageous under some circumstances.

## 12. Regulation of Gene Expression and Induction of GST Activity

The differential expression of the GSTs is governed by [hormones](#) and other factors that regulate [transcription](#) of the genes. A range of **enhancers** and other regulatory elements, shown to mediate responses in other genes, have also been shown to be functional in GST gene regulation. Particular attention has been given to the **antioxidant response element**, ARE (also called the *electrophile response element*), which has proven effective in the induction of the activity of several detoxication enzymes in mammalian systems. An extremely broad range of chemical compounds containing an electrophilic group appear to serve as inducers via ARE, suggesting that enzyme induction is a

natural response to exposure to electrophiles.

### 13. Biological functions

GSTs serve in cellular detoxication systems and constitute the major line of defense against chemical electrophiles, many of which are genotoxic and carcinogenic. The substrates include numerous xenobiotics or their metabolically activated products—for example, epoxides of carcinogenic polyaromatic hydrocarbons. Plants produce a variety of toxic compounds as a defense against attack by insects and microorganisms, and some of them can be inactivated by GSH conjugation. For example, many edible plants produce toxic organic isothiocyanates that are substrates for human GSTs. Like their mammalian counterparts, plant GSTs are inducible enzymes that serve protective functions, one of which is to provide cellular resistance to electrophilic herbicides.

Among the biologically most important substrates are numerous oxidation products of normal cell constituents, such as lipids, nucleic acids, catechols, and other aromatic or unsaturated chemical compounds. Free-radical reactions, accompanied by chemical transformations involving reactive oxygen species, lead to a variety of electrophilic products that may cause damage to DNA and proteins. Lipid peroxidation gives rise to aldehydes and activated alkenes, and cellular oxidation of catecholamines produces ortho-quinones. These naturally occurring oxidation products, considered to be etiological factors in the development of atherosclerosis, cataract, cancer, Parkinson's, and other degenerative diseases, are all substrates of GSTs. Human GST A4-4 has particularly high activity with the lipid peroxidation product 4-hydroxynon-2-enal (Table 1), and human GST M2-2 is the most efficient enzyme with the ortho-quinone aminochrome derived from dopamine.

In addition to their catalytic functions, GSTs may serve as binding proteins (ligandin) that facilitate cellular transport of organic molecules. Some of the most abundant enzymes, such as GST A1-1, appear particularly well-suited for this purpose.

### 14. Biotechnology Applications

The soluble GSTs are generally stable proteins that can be produced in large quantities by heterologous expression in *Escherichia coli*. Thus, the enzymes have potential for development into useful [recombinant proteins](#) of value for biotechnical, agricultural, and medical applications. Both the catalytic and binding properties can be explored. Catalysts and binding proteins with novel specificities can be designed by a combination of mutagenesis and selection methods.

### 15. Glutathione Transferase as a Fusion Protein for Expression of Other Proteins

In the expression of proteins by **recombinant-DNA** methods, it is frequently found that the protein product is obtained in poor yield owing to low solubility, improper **codon usage**, **proteolytic** degradation, and so on. In many cases, these problems may be overcome by producing the desired protein product in fusion with a GST molecule. DNA encoding GST with a cleavable C-terminal linker to the protein to be expressed is ligated to the target DNA. GSTs are normally stable and soluble proteins that can be expressed at high levels in many cell types, and the recombinant fusion proteins usually retain these properties. The fusion protein can be purified by affinity chromatography on immobilized glutathione derivatives and, after liberation of the desired protein by proteolytic cleavage, the GST moiety can be removed by a second affinity chromatography. By this two-step purification, the desired protein may be obtained in pure form.

### Bibliography

1. J. Booth, E. Boyland, and P. Sims (1961) *Biochem. J.* **79**, 516–524.
2. B. Combes and G. S. Stakelum (1961) *J. Clin. Invest.* **40**, 981–988.
3. L. F. Chasseaud (1979) *Adv. Cancer Res.* **29**, 175–274.
4. W. B. Jakoby (1978) *Adv. Enzymol.* **46**, 383–414.

5. B. Mannervik (1985) *Adv. Enzymol. Relat. Areas Mol. Biol.* **57**, 357–417.
6. B. Mannervik and C. Guthenberg (1981) *Methods Enzymol.* **77**, 231–235.
7. P. C. Simons and D. L. Vander Jagt (1981) *Methods Enzymol.* **77**, 235–237.
8. B. Mannervik and H. Jensson (1982) *J. Biol. Chem.* **257**, 9909–9912.
9. B. Mannervik, P. Ålin, C. Guthenberg, H. Jensson, M. K. Tahir, M. Warholm, and H. Jörnvall (1985) *Proc. Natl. Acad. Sci. USA* **82**, 7202–7206.
10. B. Mannervik, Y. C. Awasthi, P. G. Board, J. D. Hayes, C. Di Ilio, B. Ketterer, I. Listowsky, R. Morgenstern, M. Muramatsu, W. R. Pearson, C. B. Pickett, K. Sato, and M. Widersten (1992) *Biochem. J.* **282**, 305–306.
11. I. Sinning, G. J. Kleywegt, S. W. Cowan, P. Reinemer, H. W. Dirr, R. Huber, G. L. Gilliland, R. N. Armstrong, X. Ji, P. G. Board, B. Olin, B. Mannervik, and T. A. Jones (1993) *J. Mol. Biol.* **232**, 192–212.
12. B. Dragani, G. Stenberg, S. Melino, R. Petruzzelli, and B. Mannervik (1997) *J. Biol. Chem.* **272**, 25518–25523.
13. B. A. Bernat, L. T. Laughlin, and R. N. Armstrong (1997) *Biochemistry* **36**, 3050–3055.
14. A. D. Cameron, B. Olin, M. Ridderström, B. Mannervik, and T. A. Jones (1997) *EMBO J.* **16**, 3386–3395.
15. R. C. Fahey and A. R. Sundquist (1991) *Adv. Enzymol. Relat. Areas Mol. Biol.* **64**, 1–53.

### Suggestions for Further Reading

16. P. D. Josephy, B. Mannervik, and P. Ortiz de Montellano (1997) *Molecular Toxicology*, Oxford University Press, New York.
17. R. N. Armstrong (1997) Structure, catalytic mechanism, and evolution of the glutathione transferases. *Chem. Res. Toxicol.* **10**, 2–18.
18. K. A. Marrs (1996) The functions and regulation of glutathione *S*-transferases in plants. *Annu. Rev. Plant Mol. Biol.* **47**, 127–158.
19. B. Ketterer (1988) Protective role of glutathione and glutathione transferases in mutagenesis and carcinogenesis. *Mutat. Res.* **202**, 343–361.
20. B. Mannervik and M. Widersten (1995) "Human glutathione transferases: classification, tissue distribution, structure and functional properties". In *Advances in Drug Metabolism in Man* (G. M. Pacifici and G. N. Fracchia. eds.), European Commission, Luxembourg, pp. 407–459.
21. J. D. Hayes and D. J. Pulford (1995) The glutathione *S*-transferase supergene family: regulation of GST and the contribution of the isoenzymes to cancer chemoprotection and drug resistance. *CRC Crit Rev. Biochem. Mol. Biol.* **30**, 445–600.

### Glycine (Gly, G)

The [amino acid](#) glycine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to four **codons**—GGU, GGC, GGA, and GGG—and represents approximately 7.2% of the residues of the proteins that have been characterized. The glycy residue incorporated has a mass of 57.05 Da, a **van der Waals volume** of 48 Å<sup>3</sup>, and an [accessible surface area](#) of 85 Å<sup>2</sup>. Gly residues are changed during [divergent evolution](#) less frequently than average; when they are, they are interchanged in **homologous** proteins most frequently with [alanine](#), [serine](#), [aspartic acid](#), and [asparagine](#) residues.

Gly is the simplest amino acid residue, with only a hydrogen atom for a side chain. Note that the  $\alpha$ -carbon atom of Gly is not asymmetric, in contrast to the other amino acids incorporated into proteins, because it is bonded to two H atoms. Consequently, this amino acid does not occur as D or L isomers.

The absence of a larger side chain gives the polypeptide backbone at Gly residues much greater conformational flexibility than at other residues. For example, about 61% of the possible values of the torsion angles  $\phi$  and  $\psi$  of a [Ramachandran Plot](#) are permitted with a Gly residue, in contrast to the 30% with other residues; moreover, the Ramachandran plot of allowed torsion angles is symmetric. This extra flexibility tends to decrease the overall average dimensions of unfolded polypeptide chains, as the Gly residues permit the chain to reverse directions more readily. This is also used in native [protein structures](#), where Gly residues occur most frequently in reverse [turns](#). The extra flexibility of the Gly residue in an unfolded polypeptide chain is believed to destabilize all fixed conformations, and Gly residues decrease the tendency to adopt an **alpha-helical** conformation.

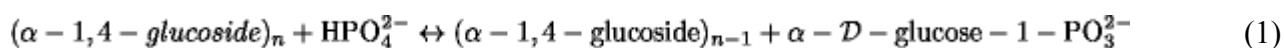
The precursors of some peptide [hormones](#), such as melanocyte stimulating hormone, have a Gly residue at the C-terminus that becomes converted to an amide group enzymatically).

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Glycogen Phosphorylase

Glycogen phosphorylase (E.C. 2.4.1.1) catalyzes the first step in the intracellular degradation of glycogen to form individual molecules of glucose that are used for energy needs. Glycogen is a polymeric storage form of sugar consisting of glucose joined into long chains by the elimination of water between the first carbon of one glucose and the fourth of another, forming  $\alpha$ -1,4 links. Every four to five glucose units, a branch is introduced by an  $\alpha$ -1,6 link, resulting in a compact, highly branched molecule, spherical in shape, and containing as many as a million glucose units. Phosphorylase carries out a phosphorylytic cleavage of the 1,4 links, as shown below:



The 1,6 branchpoint links are removed by a separate enzyme, the glycogen debranching enzyme, so that phosphorylase can then degrade virtually the entire molecule. Phosphorylase exists in two forms, *a* and *b*, differing covalently only in that Ser14 is phosphorylated in form *a*. In muscle, the nervous signal for contraction, calcium release, also activates the inactive phosphorylase *b* of resting muscle through the action of [phosphorylase kinase](#). In liver, the inactive phosphorylase of a well-fed animal is activated by the signal of the pancreatic hormone [glucagon](#) that blood sugar levels are low, again via the mediation of phosphorylase kinase. Phosphorylase was one of the first enzymes recognized to be under hormonal control by protein **phosphorylation**, as well as subject to **allosteric** controls. It can be prepared in large amounts in pure crystalline form from muscle, and this made it the subject of intense study by a number of laboratories, culminating in the determination of its three-dimensional structure at atomic resolution by [X-ray crystallography](#). The structures of several of its forms, plus associated biochemical experiments, have resulted in a better understanding of the function of this large and complicated enzyme.

## 1. Discovery and Early Developments

Carl and Gerty Cori discovered phosphorylase in 1936 when they isolated a new sugar phosphate from minced muscle mixed with water and inorganic phosphate. By characterization and synthesis, they proved that the new compound was  $\alpha$ -D-glucose-1-phosphate and showed that it was formed from glycogen by the chemical reaction depicted in equation (1). The purified enzyme required AMP for activity, and when all attempts failed to demonstrate its participation in the catalytic mechanism, as other known coenzymes did, this became recognized as perhaps the first example of an allosteric activator. In 1941, a second form of phosphorylase was isolated and crystallized that did not require AMP for activity, although the latter decreased the  $K_m$  (**Michaelis constant**) values for the substrates; this was termed the *a* form, while the first was termed the *b* form.

During the 1950s, it was discovered that phosphorylase *a* was twice the molecular weight of phosphorylase *b* and that a new enzyme, now known as a **protein phosphatase**, converted form *a* to *b*. Chemical modification of the [thiol groups](#) inactivated the enzyme, accompanied by dissociation of the *a* form into four identical subunits and dissociation of the *b* form into two identical subunits. The number of moles of AMP binding to the enzyme equaled the subunit composition of the two forms. However, under *in vivo* conditions and bound to glycogen, phosphorylase *a* is dimeric. This was the first demonstration that enzymes contain subunits (see [Quaternary Structure](#)); the only proteins known earlier to have subunits were [hemoglobin](#) and [insulin](#). The conversion of phosphorylase *b* to *a* was shown to involve the phosphorylation of residue Ser14 of each subunit by the specific phosphorylase kinase. A second organic phosphate moiety present was discovered to be [pyridoxal phosphate](#) (PLP), a derivative of vitamin B6. The **kinetic mechanism** was proven to be rapid equilibrium random Bi–Bi, meaning that either substrate can bind first, that kinetic constants approximate binding constants, and that the rate-limiting steps are to be found in the interconversion of the ternary complexes. Graves and Wang (1) have published a comprehensive review of the literature up to 1972, a convenient date because it occurs just before the crystallographic studies provided detailed structures and insights into function and regulation.

## 2. Structure–Function Relationships

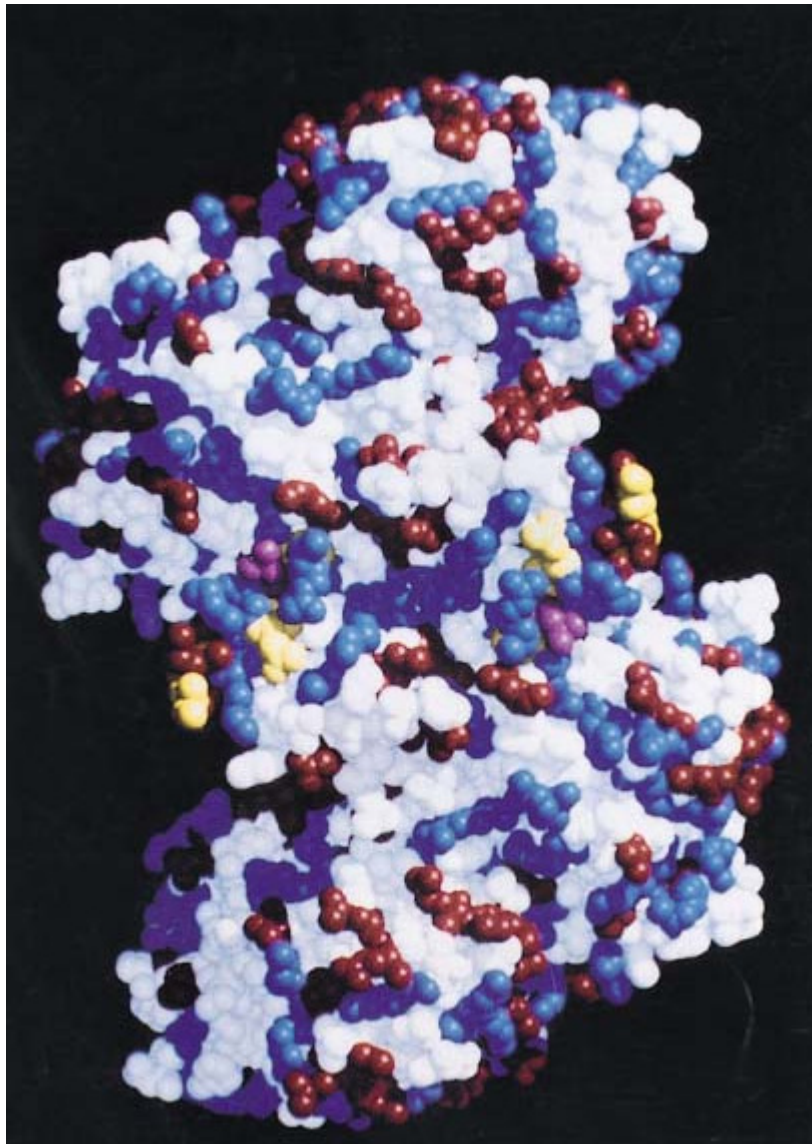
In 1970, a new tetragonal crystal form of phosphorylase *b* was discovered in the Oxford laboratory of L. N. Johnson and her colleagues. These macrocrystals belonged to [space group](#)  $P4_32_12$ ; and the [asymmetric unit](#) was one monomer, with the protein dimer formed by the crystallographic twofold axis, making the crystals suitable for X-ray crystallography. The crystals were grown at low salt concentrations in the presence of IMP. Isomorphous crystals of phosphorylase *a* were later grown in Edmonton, Canada, by using 50 mM glucose to force dimers. The structure of phosphorylase *b* was determined in Oxford by Johnson et al., while that of the *a* form was determined by Fletterick et al. in Edmonton, initially using phases from Oxford to locate the heavy atoms used in [isomorphous replacement](#). The structural studies to 1979 have been summarized in a review (2).

The complete amino acid sequence of phosphorylase was determined in the Seattle laboratories of Fischer, Neurath, Titani, and Walsh by adroit use of all conventional **sequencing** technologies, just in time for use in the three-dimensional structure determination. Each subunit of phosphorylase *a* consists of a single [polypeptide chain](#) of 842 amino acids. The *N*-terminal serine is acetylated and Ser14 is phosphorylated, while pyridoxal phosphate is bound to Lys680 via a [Schiff Base](#). Counting the three nonprotein moieties, the molecular weight of phosphorylase *a* is 97,512, while that for the *b* form subunit, often used as a molecular weight standard, is 97,434.

The tetragonal crystals of phosphorylases *a* and *b* have unit cell dimensions of  $a = b = 128.4 \text{ \AA}$  and  $c = 116.4 \text{ \AA}$ . The topology and packing are typical of a/b-type proteins, with 50% content of **a-helix** content and 30% **b-sheet** (see [Protein Structure](#)). The dimeric structures of the two forms *a* and *b* are very similar, except that in the absence of the phosphate group on Ser14, the *N*-terminal 16 residues are disordered in phosphorylase *b*. When Ser14 is phosphorylated, these residues fold into a distorted

helix that is bound to the rest of the protein. Two large **domains** compose the monomer, with the *N*-terminal domain consisting of 490 residues. The latter has a core of  $\beta$ -sheet consisting of nine  $\beta$ -strands describing a  $180^\circ$  right-handed twist flanked by helical segments, loops, and some additional  $\beta$ -strands. This domain contains the AMP binding site near Tyr75 and only 10 Å from the Ser14-phosphate group. The allosteric inhibitors ATP and glucose-6-phosphate also bind here. In a helical part of the domain, some 30 Å from either the activator site that binds AMP or the catalytic site, is found the glycogen storage site, which binds oligosaccharides in a helical form slightly distorted from that of amylose. Because binding here is some 20-fold tighter than at the [active site](#), it is suggested that phosphorylase is anchored by these sites to a glycogen particle, whether or not it is acting on other glycogen end-chains at the active site. Components of the active site and the inhibitor site for caffeine (Phe285) are found at the interface with the C-terminal domain. The latter domain contains a core of six parallel  $\beta$ -strands, with a twist and connectivity to surrounding  $\alpha$ -helices that is topologically similar to the [nucleotide-binding motif](#) of the [dehydrogenases](#), although more buried by additional helices and loops. The deep cleft between the two domains contains (a) the active site, which is buried deeply, next to the phosphate group of PLP, and (b) the entrance between Phe285 and Tyr613, which is “locked” when caffeine binds in between. The two subunits undertake many interactions when they associate to form the dimer, including shared binding sites for AMP and for the phosphate of Ser14 and long paired loops/helices (the “towers”) that extend reciprocally toward the symmetry-related active sites. The dimer possesses a concave catalytic face containing the active sites and glycogen storage sites, by which it may bind to glycogen while exposing the convex control face to the action of phosphorylase kinase or phosphatase. A computer drawn space-filling model of the control face is shown in Fig. [1](#).

**Figure 1.** Space-filling model of the control face of the phosphorylase *a* dimer viewed down the twofold axis. Positively charged residues are blue, negatively charged residues are red, and the phosphate groups of the two Ser14 residues are magenta, while noncharged residues of the *N*-terminal  $\alpha$ -helices are yellow. The *N*-terminal segments, identified by the phosphates and yellow residues, lie across the subunit interface. This T-state structure, with the phosphates fully exposed and surrounded by positive charges, represents the favored substrate for protein phosphatase. (Courtesy of Drs. R. Read and A. Muir, Department of Biochemistry, University of Alberta.) See color insert.



Studies on the catalytic mechanism of phosphorylase have been reviewed up to 1985 (3) and centered around the role of pyridoxal 5'-phosphate. The coenzyme, sandwiched between the two domains, plays an unusual structural role because its removal leads to an inactive monomer. Extensive analog studies have eliminated all parts of the coenzyme from participating in catalysis except the phosphate, which must be capable of forming a dianion. In 1977, Graves et al. found that phosphorylase *b* reconstituted with pyridoxal, and hence inactive, regained appreciable activity in the presence of inorganic phosphate ( $P_i$ ), phosphite, or fluorophosphate. Furthermore, pyrophosphate showed competitive inhibition with both glucose-1-P and the activating phosphite, while binding only one mole per monomer. That same year, crystallographic studies placed only the phosphate of the coenzyme near the substrate, 7.0 Å from the latter's phosphate in these crystals, which are in the allosteric T state. These results led to the interacting phosphates hypothesis that was strengthened by a number of experiments, including the reconstitution of the enzyme with pyridoxal pyrophosphate glucose. This compound contains both substrate and coenzyme covalently linked, and the enzyme could transfer the glucose to an added oligosaccharide. This suggested that the coenzyme phosphate might act as an electrophile (Lewis acid). However, the alternative hypothesis that the coenzyme phosphate acts as a proton donor (Brønsted acid) now appears more likely, especially considering the crystallographic studies of time-resolved catalysis in the crystalline state (4). Phosphorylase catalyzes the phosphorylation of heptenitol to heptulose 2-phosphate, and the latter, a potent

inhibitor, is found in the crystal structure with a hydrogen bond between its phosphate and that of the coenzyme. This suggests that in the normal phosphorylation of glycogen, the proton of  $P_i$  attacks the glycosidic bond, cleaving it, while the substrate  $P_i$  immediately gains a proton from the PLP phosphate. The phosphate anion stabilizes the glucosyl carbonium ion, necessary to maintain the alpha anomeric configuration. The reaction is completed by a nucleophilic attack of the phosphate on the carbonium ion to form glucose 1-phosphate.

### 3. Regulation

The two forms of the enzyme differ markedly in their allosteric regulatory properties. The *b* form is activated by AMP or IMP and is inhibited by ATP, ADP, glucose, and glucose-6-phosphate; the *a* form is not subject to these controls, but is still inhibited by glucose. Glucose binds at the active site of phosphorylase and is a competitive inhibitor of glucose-1-P, a primitive end-product **feedback inhibition** control. The crystals of phosphorylase *a* grown with glucose are in the allosterically inhibited T state, while those of phosphorylase *b* grown with IMP are also in the T state, permitting a direct comparison of the effect of phosphorylating Ser14 on the structure (5, 6). The *b* form requires AMP for activity, but full conversion to the active R state requires substrate, while ATP and glucose-6-phosphate compete for the AMP site and represent energy and end-product allosteric controls, respectively. While the allosteric ratio of T to R for unliganded *b* is of the order of 3,000, that for the *a* form is only 10, allowing full activity of the latter with substrates alone, while only the slight activation of the *a* form by AMP is inhibited by ATP. Phosphorylase *a* may be said to have “escaped allosteric control,” the latter form of regulation having been superseded by covalent modification, under extracellular hormonal and neural regulation. Because AMP binds 200 times more tightly to the *a* than the *b* form, it is obvious that a large portion of its binding energy to the *b* form is utilized in the allosteric conformational change and that the interactions of the Ser14 phosphate group must be responsible. Comparison of the two T-state tetragonal crystal structures reveal that the disordered *N*-terminal segment, residues 5–16, of the *b* form becomes an ordered helix which lies across the dimer interface in the *a* form, with eight new intersubunit [hydrogen bonds](#) between polar groups, plus additional **hydrophobic** interactions. In addition, the AMP binding site is more fully formed and capable of greater interaction with the ligand. The tighter subunit associations of the dimer in the *a* form account for it obeying the allosteric **concerted model** of Monod et al., whereas phosphorylase *b* follows the **sequential model** of Koshland et al. (with intermediate forms occurring). Thus the structural studies of phosphorylases *a* and *b* have resolved an old controversy within a single enzyme.

### Bibliography

1. D. J. Graves and J. H. Wang (1972) " -Glucan Phosphorylases", In *The Enzymes*, Vol. 7, 3rd. ed., Academic Press, New York, pp. 435–482.
2. R. J. Fletterick and N. B. Madsen (1980) *Annu. Rev. Biochem.* **49**, 31–61.
3. N. B. Madsen and S. G. Withers (1986) In *Vitamin B<sub>6</sub> Pyridoxal Phosphate*, Part B, John Wiley & Sons, New York, pp. 355–389.
4. L. N. Johnson (1989) *Carlsberg Res. Commun.* **54**, 203–229.
5. S. R. Sprang et al. (1988) *Nature* **336**, 215–221.
6. M. F. Perutz (1988) *Nature* **336**, 202–203.

### Suggestions for Further Reading

7. N. B. Madsen (1997) "Glycogen". In *Encyclopedia of Human Biology*, Vol. 4, 2nd ed., Academic Press, San Diego, pp. 341–351. [A summary of glycogen metabolism.]
8. N. B. Madsen (1986) "Glycogen Phosphorylase", In *The Enzymes*, Vol. 17, 3rd ed., pp. 365–394. [Updates the previous review, with emphasis on phosphorylation of Ser14 and its consequences.]



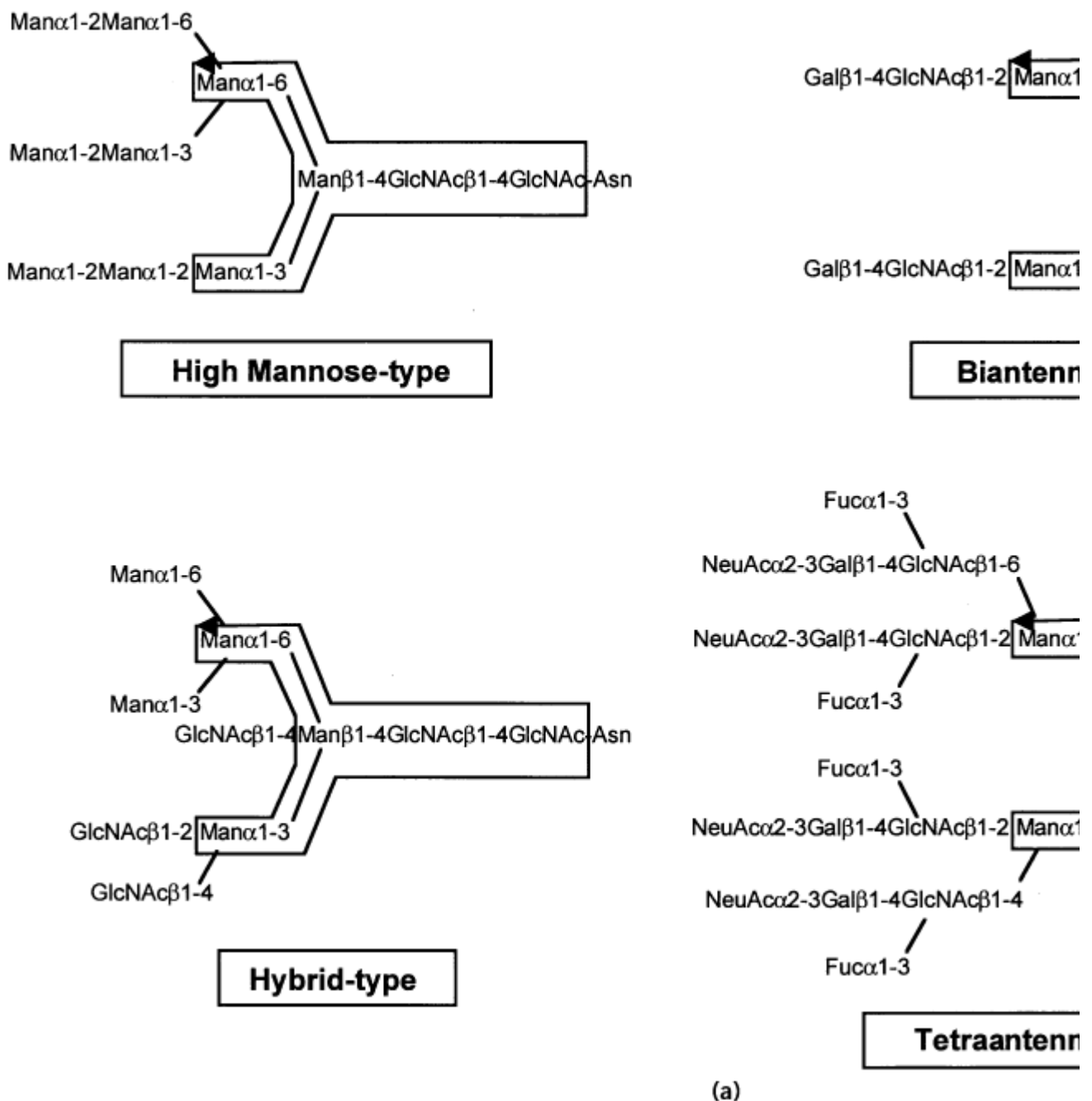
## ***N*-Glycosylation**

Glycoproteins are characterized by the presence of carbohydrates covalently linked to amino acid residues. More than a dozen different protein–carbohydrate linkages have been identified in glycoproteins, and they are often categorized as either *O*-glycosidic (*O*-glycans or *O*-linked oligosaccharides), in which carbohydrate is covalently attached to protein via a hydroxyl group in the side chains of several amino acids, or *N*-glycosidic (*N*-glycans or *N*-linked oligosaccharides), in which sugar is covalently attached to protein through the amide nitrogen of asparagine residues. Protein glycosylation is the most diverse type of post-translational modification of proteins. The biochemical processes through which these linkages are generated are commonly called [O-glycosylation](#) and *N*-glycosylation pathways, and each pathway involves dozens of different enzymes. There are dramatic differences between these pathways in terms of their subcellular localization, the enzymes involved, the structures of the oligosaccharides synthesized, and the biological functions the glycans. *N*-Glycosylation is best characterized in animal and plant cells, but *N*-linked sugars have also been observed in bacterial glycoproteins. As we shall see, *N*-glycosylation of proteins is an exceedingly complex pathway, involving the building, trimming, and elongation of relatively large oligosaccharides as they pass through the various secretory organelles.

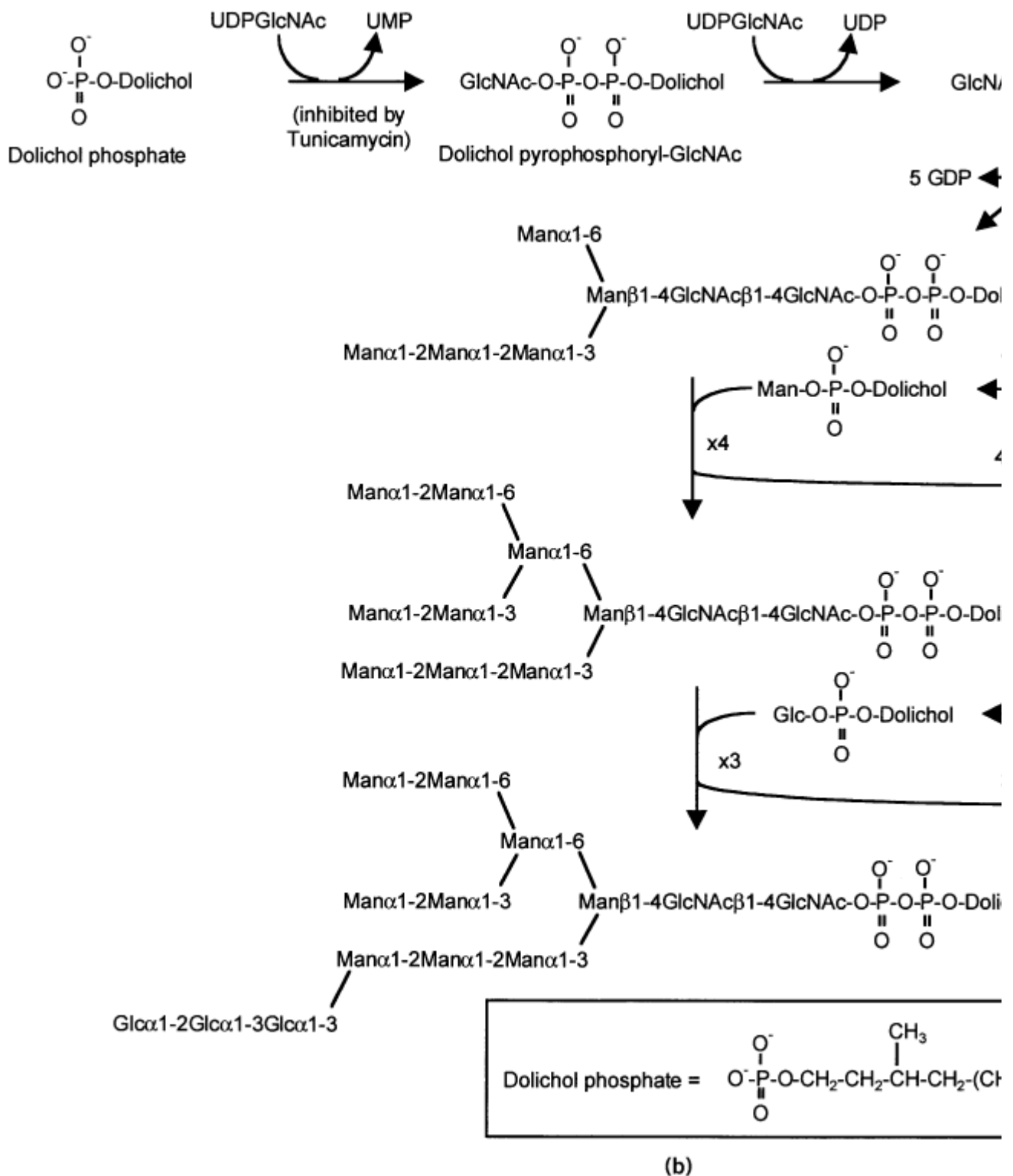
### 1. Different Types of *N*-Glycans

Examples of different types of *N*-glycan in animal cells are shown in Figure 1a. *N*-glycans contain a common pentasaccharide core composed of two *N*-acetylglucosamine (GlcNAc) residues and three mannose (Man) residues. *N*-Glycans are classified in three broad groups on the basis of their overall composition. *N*-Glycans that contain primarily Man residues in terminal, nonreducing positions are termed high mannose- or polymannose-type. In contrast, *N*-glycans that contain sugars other than Man in terminal nonreducing positions, such as sialic acid (NcuAc) or galactose (Gal), are classified as complex-type. A relatively simple example of a biantennary complex-type *N*-glycan (Fig. 1a) is compared to an example of a more highly branched version, known as a tetraantennary complex-type *N*-glycan. Those *N*-glycans that have features shared by both high mannose and complex-type are classified as hybrid-type. Thousands of structurally different *N*-glycans have been identified, and they may contain many additional sugars to those shown in Figure 1a. It should be noted that glycoconjugates are highly unusual in their structural motifs when compared to other classes of [macromolecules](#). For example, most glycoconjugates, in contrast to nucleic acids and proteins, are generally branched rather than linear. The branched nature allows the generation of highly compacted and diverse structures. In contrast to nucleic acids and proteins, which are synthesized on templates, glycoconjugates are synthesized without a template, and their assembly requires the stepwise addition of monosaccharides from activated sugar precursors by individual enzymes termed glycosyltransferases. The diversity of *N*-glycan structures is generated by the inherent branching potential, the types and activity levels of the glycosyltransferases expressed in cells, the glycoprotein species, and a host of other factors (1).

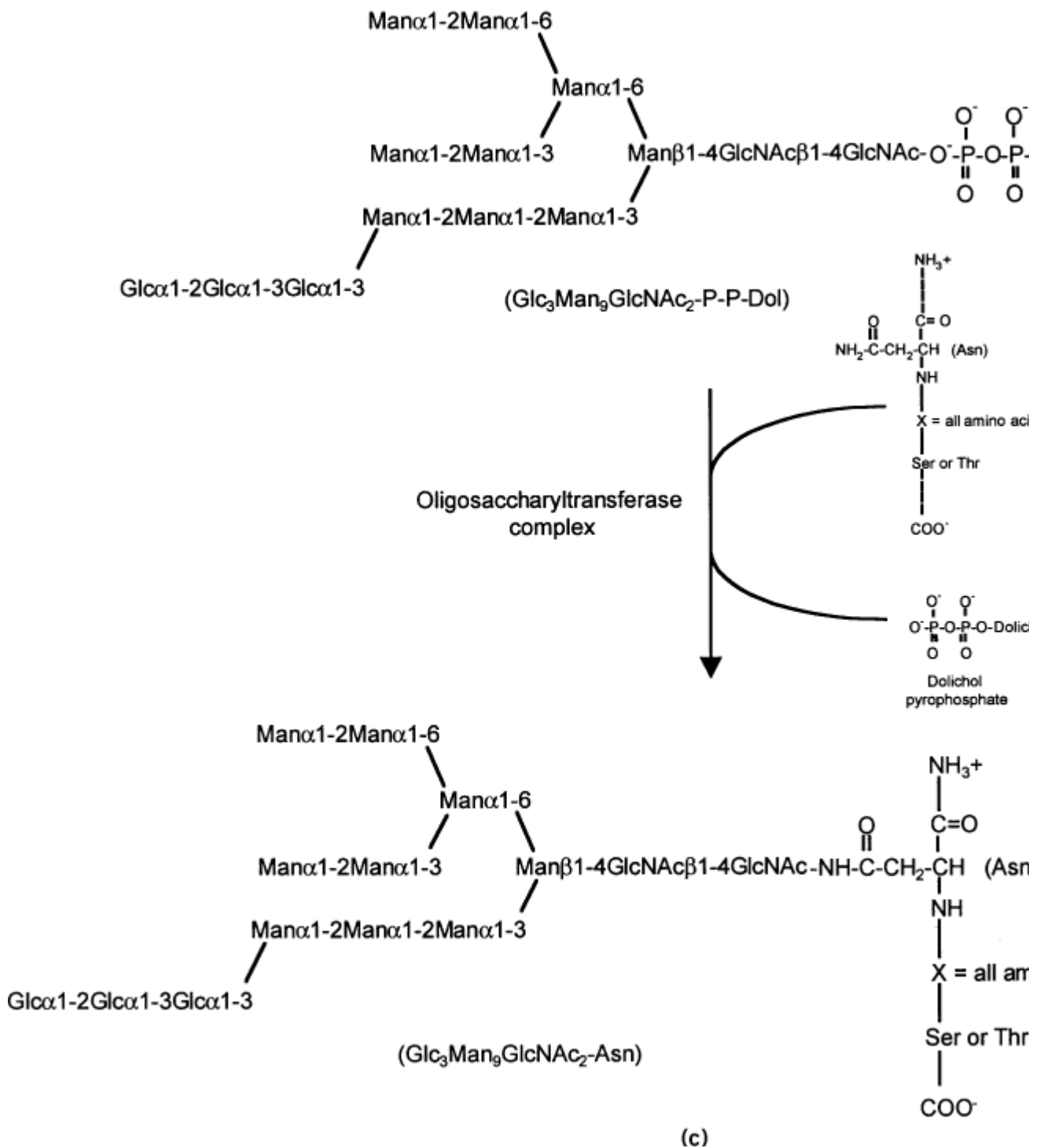
**Figure 1.** Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of Glc<sub>3</sub>Man<sub>9</sub>GlcNAc<sub>2</sub>-P-P-dolichol in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



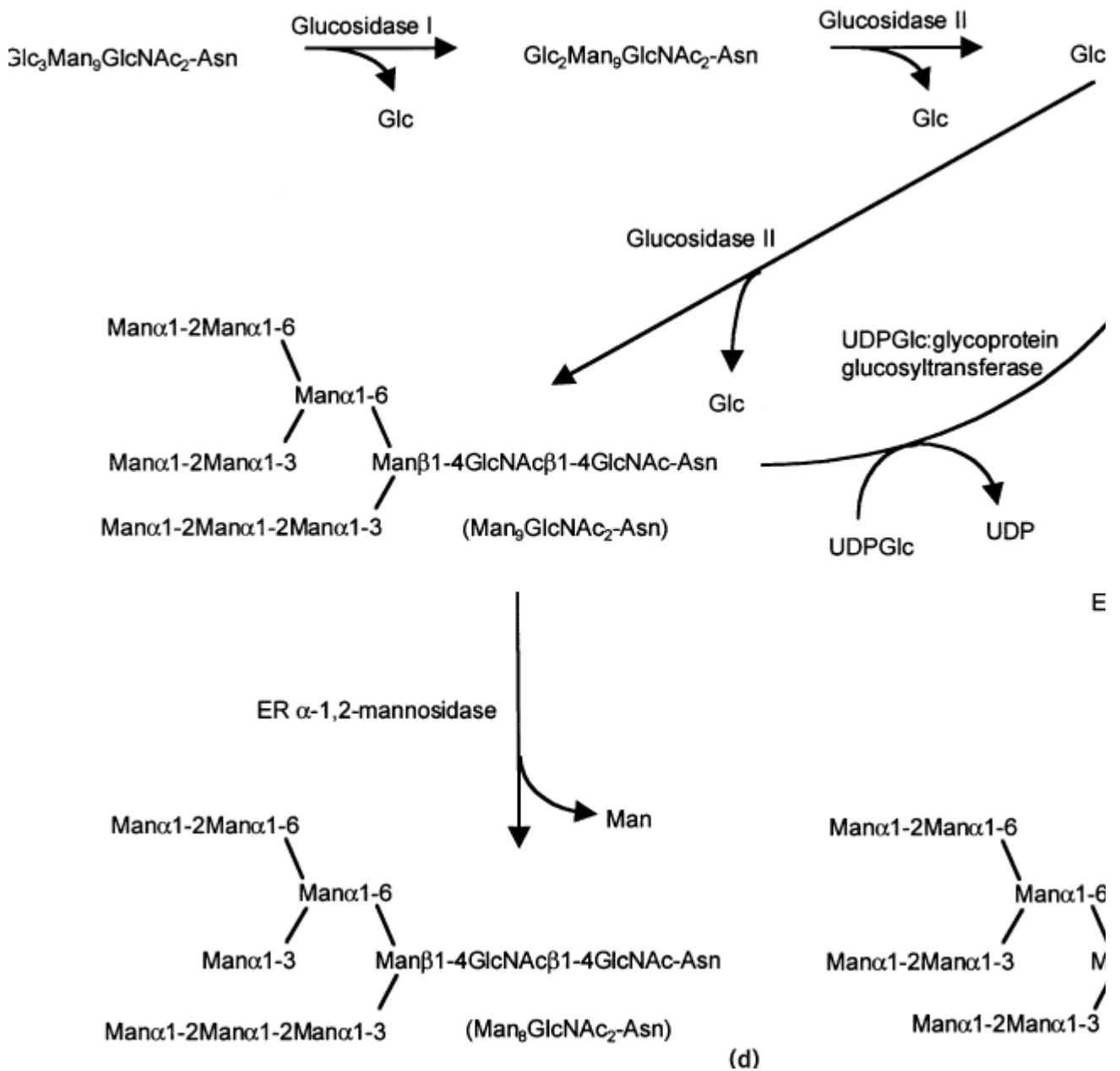
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of Glc<sub>3</sub>Man<sub>9</sub>GlcNAc<sub>2</sub>-P-P-dolichol in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



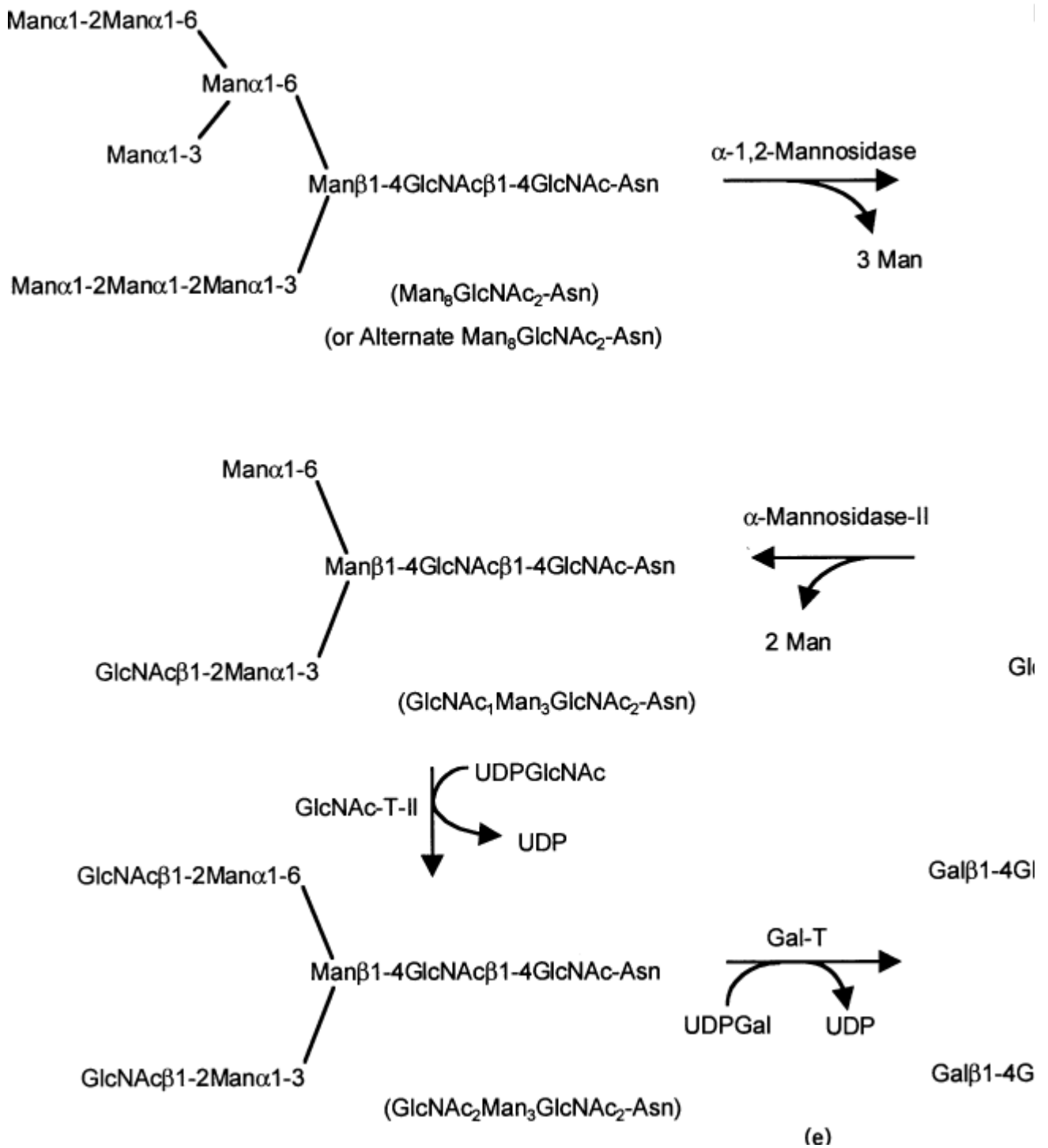
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



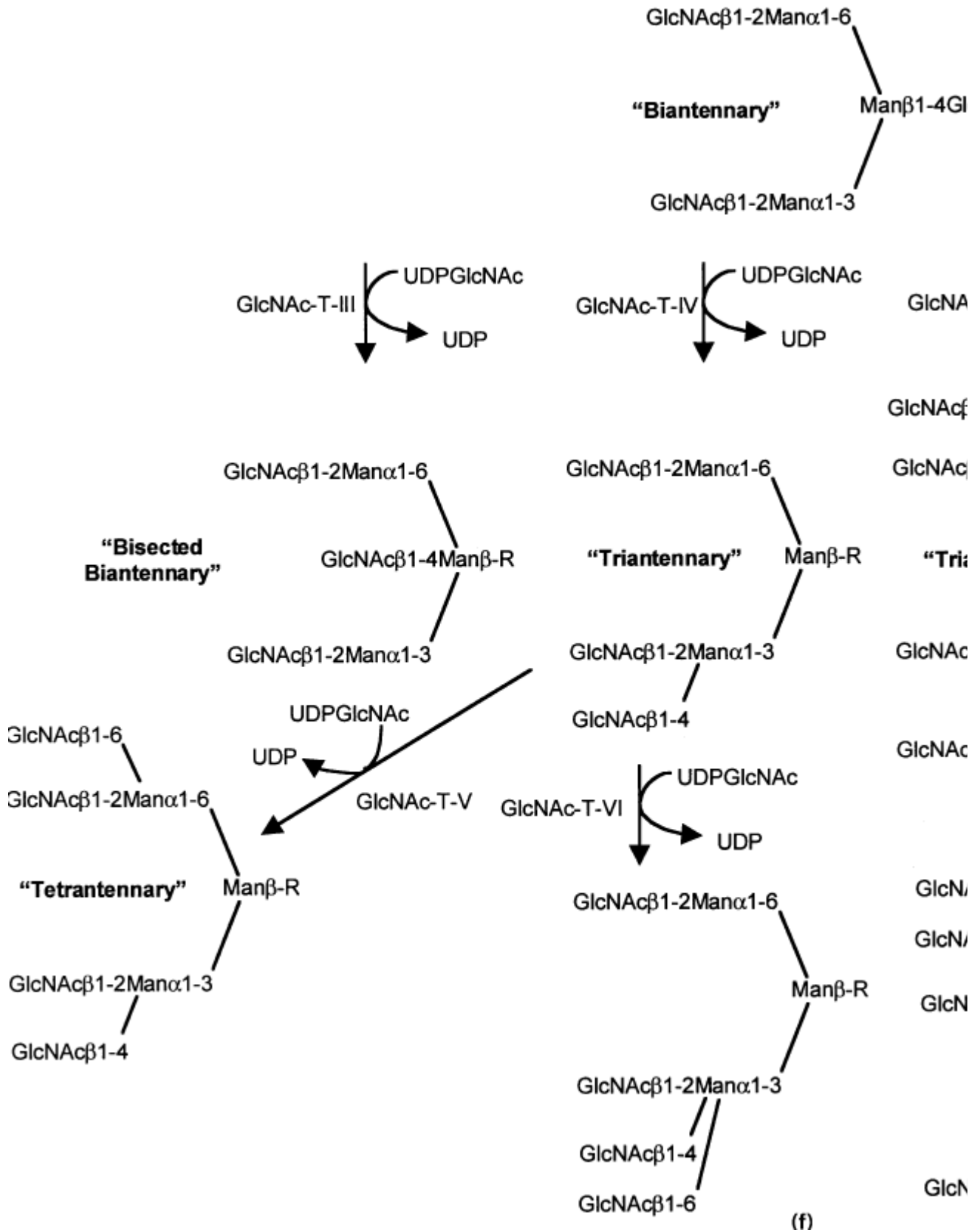
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation enzymes observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



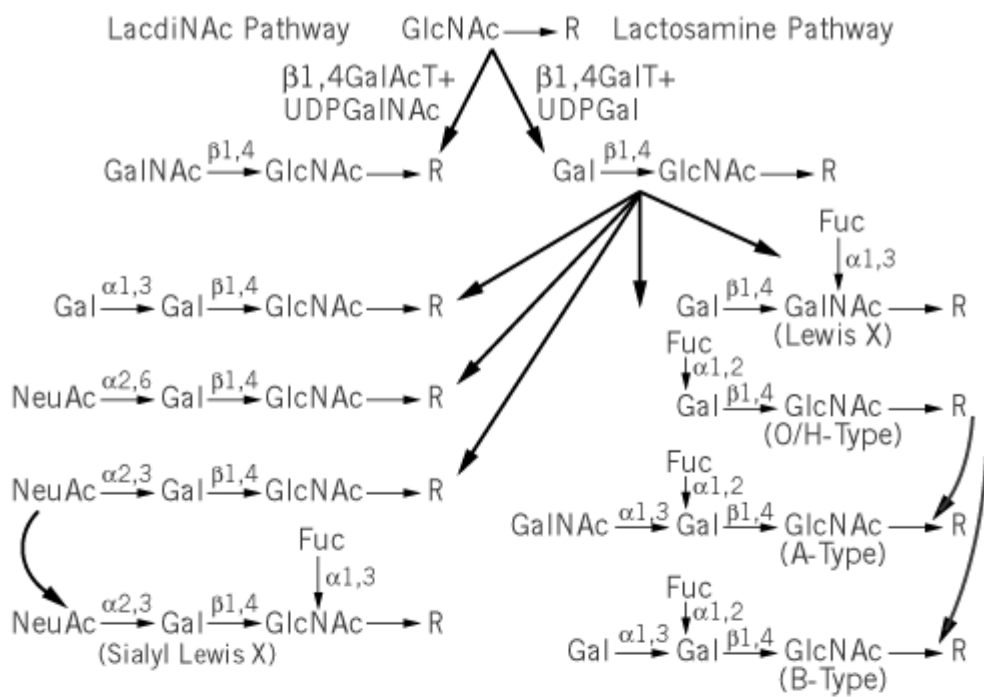
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



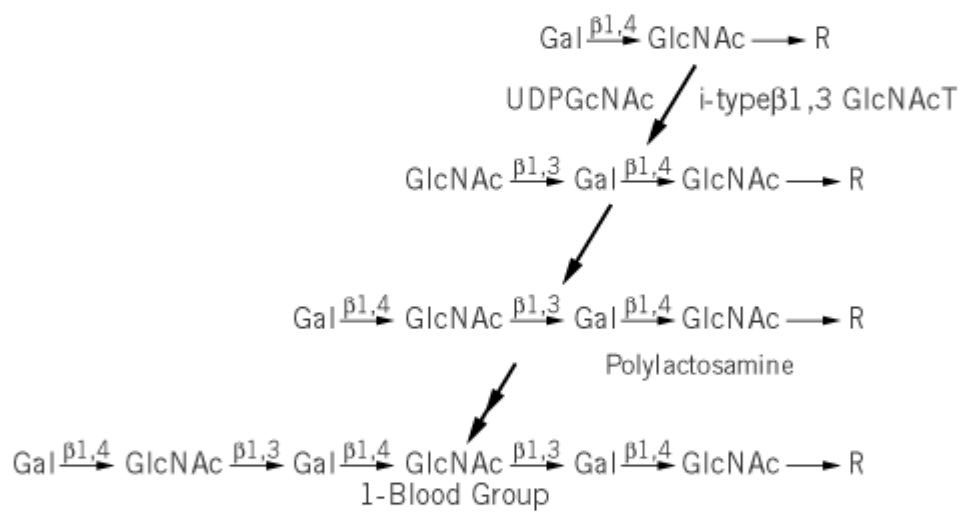
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation enzymes observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



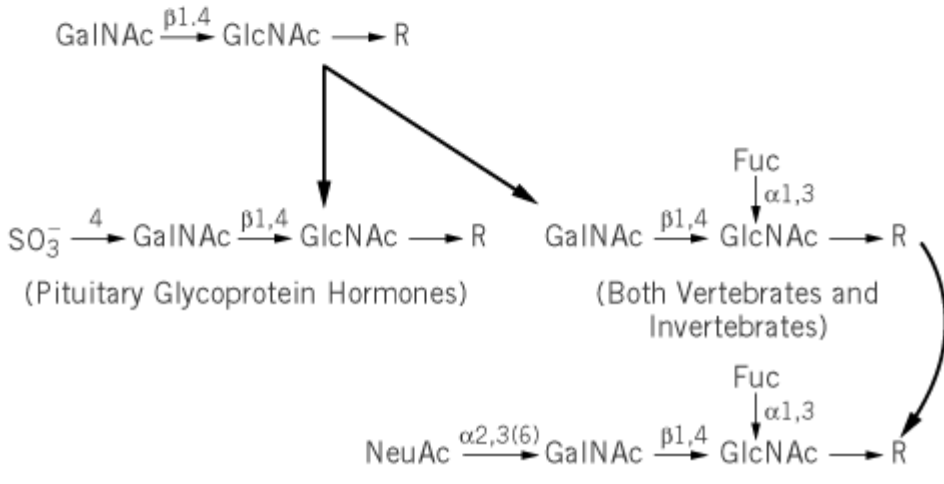
**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



Lactosamine Pathway



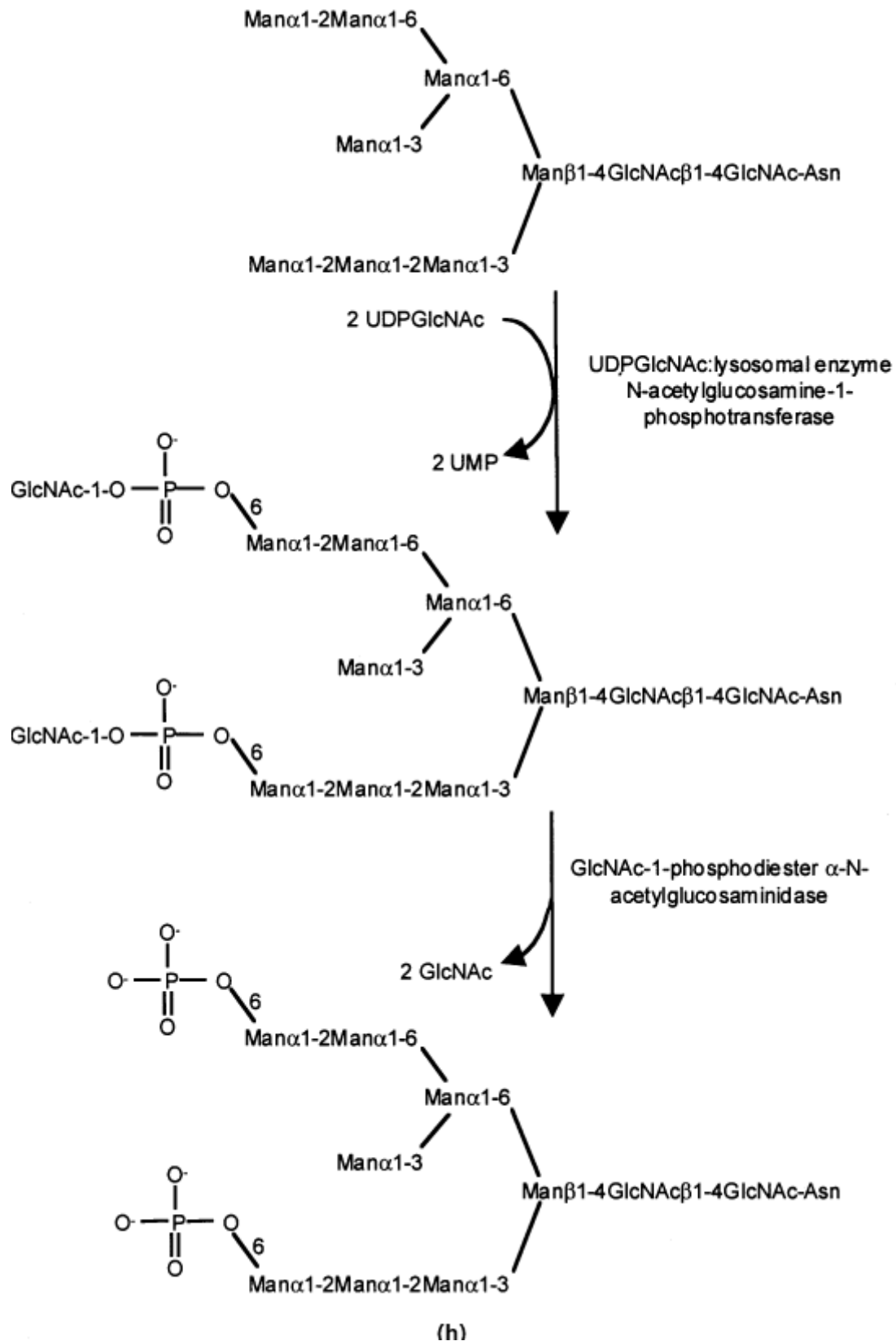
LacdiNAc Pathway



(g)



**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.



**Figure 1.** (continued) Structure and biosynthesis of *N*-glycans: (a) examples of different types of *N*-glycans; (b) steps in assembly of Glc<sub>3</sub>Man<sub>9</sub>GlcNAc<sub>2</sub>-P-P-dolichol in RER membrane; (c) reaction catalyzed by the oligosaccharyltransferase (OST); (d) early-stage trimming reactions of high mannose-type *N*-glycans in the RER and Golgi apparatus; (e) late-stage processing of high mannose-type *N*-glycans and their conversion to hybrid- and complex-type chains; (f) branching of complex-type *N*-glycans by *N*-acetylglucosaminyltransferases in the Golgi apparatus; (g) different types of terminal glycosylation observed in complex-type *N*-glycans; (h) biosynthesis of phosphorylated high mannose-type *N*-glycans in lysosomal enzymes.

| Inhibitor  | Action  |
|--|---|
| Tunicamycin  | Inhibits formation of dolichol pyrophosphoryl-GlcNAc from dolichol phosphate and UDPGlcNAc  |
| Australine<br>Castanospermine<br>1-Deoxynojirimycin<br>N-methyl-deoxynojirimycin<br>DMDP [2R,5R-bis(Hydroxymethyl)-3R,4R-dihydroxypyrrolidine] | Inhibits Glucosidase I  |
| Bromoconduritol<br>Castanospermine<br>1-Deoxynojirimycin (DNJ)<br>N-methyl-deoxynojirimycin  | Inhibits Glucosidase II   |
| 1-Deoxymannojirimycin (DMJ)<br>N-methyl-deoxynojirimycin<br>DMDP<br>Kifunensine  | Inhibits $\alpha$ 1,2-Mannosidases  |
| Glc $\alpha$ 1-3DMJ (GDMJ)   | Inhibits Endomannosidase  |
| Swainsonine<br>Mannostatin   | Inhibits $\alpha$ -Mannosidase II   |
| Brefeldin A  | Blocks translocation of proteins from ER to Golgi by blocking binding of GDP-ribosylation factor to Golgi and inhibiting GTP-GDP exchange |

## 2. Consensus Sequence for *N*-Glycosylation

The most common structure of *N*-glycans is the presence of GlcNAc-Asn linkages. The presence of these linkages in animal-derived glycoproteins was proven through the pioneering work of R. D. Marshall and A. Neuberger (2). Subsequent studies of many different glycoproteins and the enzymes involved in glycoprotein biosynthesis indicated that the primary **consensus sequence** for *N*-glycosylation of proteins in animals, plants, and yeast is Asn-X-Ser/Thr, where X can be any amino acid except Pro. This consensus sequence is known as the *N*-glycosylation sequon. However, Asn residues within the sequence Asn-X-Cys have also been found to be *N*-glycosylated in some animal glycoproteins (3), suggesting that this may be an alternative sequon in some special cases. For most glycoproteins, most of the Asn-X-Ser/Thr *N*-glycosylation sequons in the extracytoplasmic regions of the protein are *N*-glycosylated. In those cases where the Asn-X-Ser/Thr *N*-glycosylation sequons are not glycosylated, it has been assumed that the site is probably inaccessible to the enzyme

machinery that adds the sugar residues to the Asn. In some cases, however, it is clear that amino acid residues flanking the Asn-X-Ser/Thr *N*-glycosylation sequon can modulate the efficiency of *N*-glycosylation (4, 5). The numbers and positions of *N*-glycans in proteins demonstrate tremendous diversity that is glycoprotein-specific. Some glycoproteins, such as hen ovalbumin, contain a single *N*-glycan, whereas other glycoproteins may contain dozens of *N*-glycans. The envelope glycoproteins of HIV-1, the AIDS virus, contains approximately two dozen *N*-glycans. The *N*-glycans in some glycoproteins at particular sites may be high mannose-type, whereas at other sites within the same glycoprotein the *N*-glycans may be multiantennary complex type. This phenomenon is referred to as “site-specific *N*-glycosylation.”

Although the *N*-glycosylation sequon is universal among eukaryotes and typically is composed of GlcNAc-Asn linkages, certain types of bacteria also synthesize glycoproteins with other carbohydrates attached to Asn residues within the *N*-glycosylation sequon. For example, Glc-Asn and GalNAc-Asn linkages are found in the halobacterial cell surface glycoproteins and flagellins (6, 7).

### 3. Synthesis of the Lipid-Linked Oligosaccharide

Many of the early steps in *N*-glycan biosynthesis in yeast, plants, and animals are similar, resulting in the synthesis of a high mannose-type *N*-glycan precursor (8). However, after generation of the high mannose-type *N*-glycan, there is tremendous diversity between phyla in the types of modifications they make to this precursor, as discussed below. *N*-Glycosylation in eukaryotes is initiated in the rough endoplasmic reticulum (RER), and final steps in the pathway are carried out in the Golgi apparatus. Thus, *N*-glycosylation is highly compartmentalized and is associated with the secretory organelles in eukaryotes. *N*-Glycans undergo a series of trimming and processing reactions within these organelles in a highly orchestrated pathway involving the removal and addition of carbohydrate residues. All *N*-glycans arise from a common high mannose-type *N*-glycan precursor on Asn residues. This high mannose-type *N*-glycan precursor arises from the transfer of a preformed precursor oligosaccharide linked to a lipid that has the general formula  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$ . Dolichol is a long-chain polyisoprenoid (Fig. 1b), and it serves as a membrane anchor for assembly of the lipid-linked oligosaccharide. Although glucose (Glc) residues are not commonly found on *N*-glycans in mature glycoproteins, Glc residues are essential components of the lipid-linked oligosaccharide for optimal donor activity.

This lipid-linked oligosaccharide  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  is partly assembled on the cytoplasmic face of the RER; as the oligosaccharide grows in size by the stepwise addition of sugar residues (8), further elongation occurs in the lumen of the RER (9). The assembly is catalyzed by intrinsic membrane proteins in the RER (10). The first step in the assembly involves the acceptor dolichol monophosphate and the sugar nucleotide donor uridine diphospho-*N*-acetylglucosamine (UDPGlcNAc). The main product of the first reaction is dolichol pyrophosphoryl-GlcNAc and the side-product UMP (Fig. 1b). An important discovery in the elucidation of the *N*-glycosylation pathway was the finding that this first step is inhibited by tunicamycin, a toxic glucosamine derivative antibiotic produced by *Streptomyces lysosuperficus* (11). Because tunicamycin blocks the first step in the assembly process, cells and tissues treated with the drug are unable to synthesize glycoproteins containing *N*-glycans.

The dolichol phosphoryl-GlcNAc serves as a precursor for the assembly of the lipid-linked oligosaccharide and serves as the acceptor for a second reaction in which another GlcNAc residue is added from the donor UDPGlcNAc (Fig. 1b). The resulting main product is dolichol pyrophosphoryl-GlcNAc<sub>2</sub>, and the side product is UDP. The dolichol pyrophosphoryl-GlcNAc<sub>2</sub> is the acceptor for a group of mannosyltransferases in the cytoplasmic face of the RER that add five Man residues in stepwise addition from the donor guanosine diphospho-mannose (GDPMan). The product of this reaction is a lipid-linked oligosaccharide containing 5 Man residues in a structure

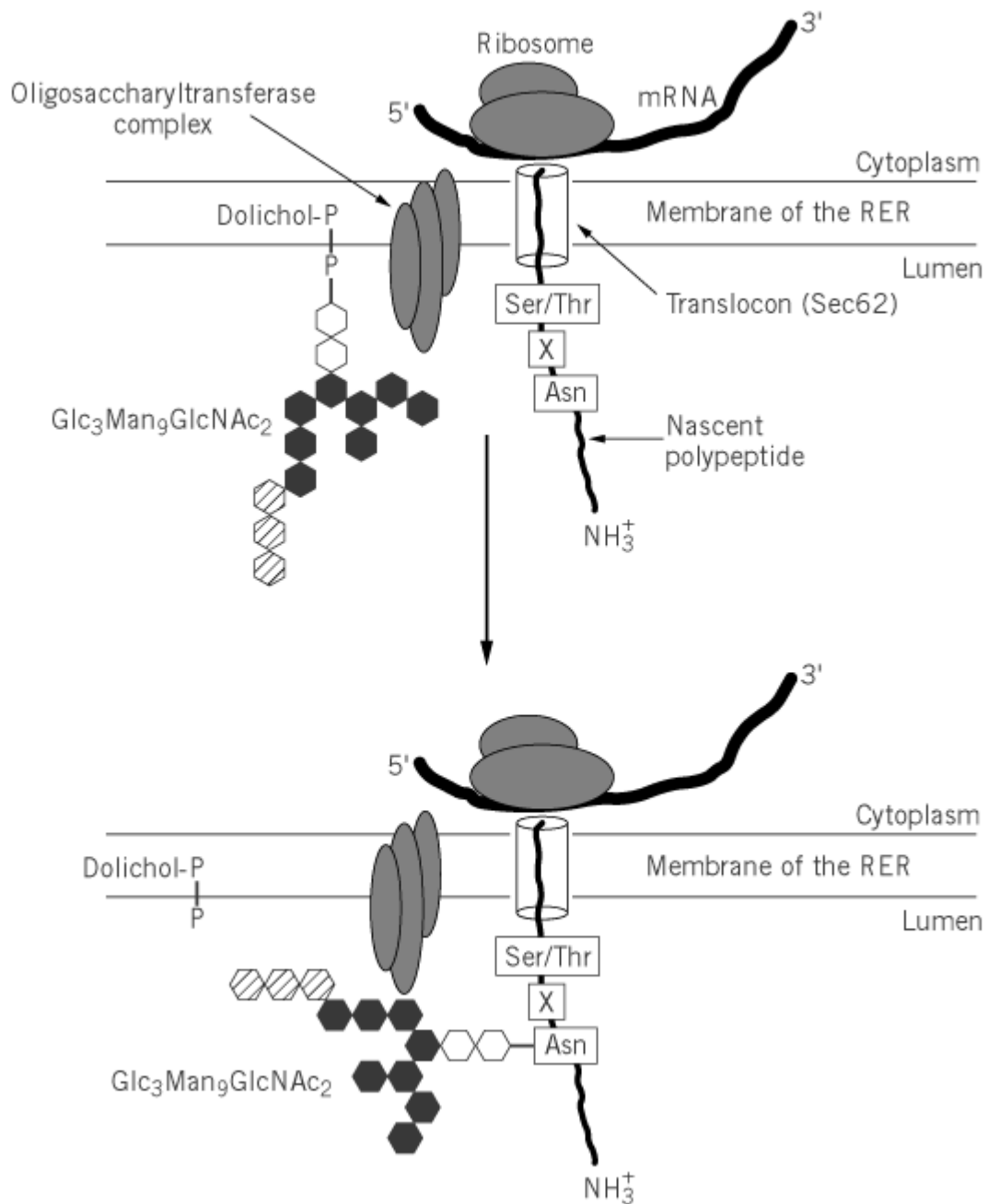
having the composition  $\text{Man}_5\text{GlcNAc}_2\text{-P-P-dolichol}$ . This intermediate is translocated into the lumen of the RER, where it is further elongated by Man addition through the donor Dol-P-Man (9). Dol-P-Man is assembled on the cytoplasmic face of the RER from dolichol phosphate and the donor GDPMan. After the addition of four Man residues from Dol-P-Man, the lipid-linked oligosaccharide has the composition  $\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$ .

This intermediate  $\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  is further elongated through addition of three Glc residues to  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  (Fig. 1b). These Glc residues are added from the donor Dol-P-Glc, which itself is assembled on the cytoplasmic face of the RER from the acceptor Dol-P and the donor uridine diphospho-glucose (UDPGlc). The final lipid-linked oligosaccharide has a total of 14 sugar residues and represents the pivotal precursor for the biosynthesis of all *N*-glycans in eukaryotes through the en bloc transfer to Asn residues in glycoproteins by the oligosaccharyltransferase complex in the membrane of the RER. It is interesting to note that 11 of the 14 residues, which are in external positions in the lipid-linked oligosaccharide, are in  $\alpha$ -linkage, whereas the 3 core residues are in  $\beta$ -linkage. Since  $\beta$ -linked sugars are usually linear in structure and  $\alpha$ -linkages promote turns in the structure, it is likely that the external residues of the lipid-linked oligosaccharide are tightly folded into a compact structure extending from the dolichol pyrophosphoryl core.

#### 4. The Oligosaccharyltransferase

The oligosaccharyltransferase (OST) is a membrane-associated protein assembly that catalyzes the transfer of  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  to Asn residues within the *N*-glycosylation sequon on nascent polypeptides in the lumen of the RER. Except for trypanosomatid protozoa, which utilize  $\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$ , the OST in all other organisms utilizes  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  as a primary donor. The OST requires  $\text{Mn}^{2+}$  for activity and catalyzes the formation of a C—N bond between the C-1 position of the reducing GlcNAc residue within  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  and the amide nitrogen of asparagine with the release of P-P-dolichol (Fig. 1c and Fig. 2). It is likely that the hydroxyl groups of Ser or Thr residues within the *N*-glycosylation sequon participate in the catalysis by activating the amide nitrogen (12, 13). It has been proposed that amino acid residues within the OST force the deprotonation of hydroxyl side chains of Ser or Thr residues, which, in turn, causes them to act as a strong oxyanion for abstraction of proton from the carboxamide nitrogen of asparagine. This possible mechanism explains the conservation of Ser and Thr in the third position of the *N*-glycosylation sequon. The activated asparagine is proposed to attack the  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  nucleophilically, resulting in displacement of P-P-dolichol from the C-1 of GlcNAc. The fact that Pro residues are not known to occur within functional *N*-glycosylation sequons suggests that close proximity between Ser and Thr residues and the Asn residue is required for optimal OST activity. Pro residues would be expected to introduce a strong bend in the polypeptide backbone, positioning Ser and Thr residues away from Asn.

**Figure 2.** Schematic of cotranslational *N*-glycosylation of nascent polypeptides in the RER by OST.



The OST has been purified from many sources (14) and is usually recovered as a protein complex containing from three to six polypeptides. However, other accessory proteins may also be important in efficient activity of the OST. Genetic mutations of genes encoding OST polypeptides have been identified in yeast, and the genes have been **cloned** by [complementation](#) (14-16). In yeast the genes identified to date are designated OST1, WBP1, OST3, SWP1, and OST2. Protein homologues of several of these genes have been found in higher animals. All the proteins encoded by these genes in yeast appear to be essential for optimal activity *in vivo*.

The OST is likely to be in close proximity to the translocon or Sec62 pore complex in the RER through which proteins cotranslationally enter the lumen of the RER. The catalytic domain of the OST is functionally located approximately 3–4 nm or  $34 \times 10^{-9}$  m from the luminal face of the RER and acts on Asn residues within the *N*-glycosylation sequon of nascent polypeptides prior to complete translation of the proteins. Thus, *N*-glycosylation generally occurs as a cotranslational process.

## 5. Early-Stage Processing Pathway for *N*-Glycans in the RER

Following transfer of the  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2$  to Asn residues, a series of reactions occurs in the lumen of the RER resulting in removal of the 3 Glc residues by glucosidases I and II. These trimming reactions are initiated on the nascent glycopeptide prior to complete translation. The newly synthesized glycoprotein is also subjected to an ER-localized  $\alpha$ 1,2-mannosidase, which usually removes one of the Man residues (Fig. 1d) (17). As discussed above, the presence of glucosyl residues on the  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2$ -P-P-dolichol is essential for optimal donor activity with the OST. But the glucosyl residues on the transferred oligosaccharide  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{Asn}$  also serve an important role in **protein folding** and retention in the RER.

## 6. Role of *N*-Glycosylation in Protein Folding

Following removal of glucosyl residues from  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{Asn}$  in newly synthesized glycoproteins, poorly folded glycoproteins are substrates for the enzyme UDPGlc:glycoprotein glucosyltransferase. This enzyme is capable of recognizing the high mannose-type *N*-glycan in poorly folded glycoprotein; it then catalyzes glucosylation of  $\text{Man}_9\text{GlcNAc}_2\text{Asn}$  sites, using UDPGlc as a donor to generate the monoglucosylated species  $\text{Glc}_1\text{Man}_9\text{GlcNAc}_2\text{Asn}$  (18).

**Molecular chaperones**, such as [calnexin and calreticulin](#), bind via carbohydrate interactions to monoglucosylated  $\text{Glc}_1\text{Man}_9\text{GlcNAc}_2\text{Asn}$  (19, 20). The interactions with these chaperones facilitates other molecular interactions, leading to correct protein folding and retarding glycoprotein exit from the endoplasmic reticulum. Interaction with chaperones is an equilibrium binding, however, and the monoglucosylated  $\text{Glc}_1\text{Man}_9\text{GlcNAc}_2\text{Asn}$  is available as a substrate for glucosidase II, which can remove the Glc residue to regenerate  $\text{Man}_9\text{GlcNAc}_2\text{Asn}$ . As indicated in Figure 1d, a cycle may develop between UDPGlc:glycoprotein glucosyltransferase and glucosidase II; this cycle continues until the glycoprotein folds in a manner that blocks its interaction with the UDPGlc:glycoprotein glucosyltransferase. The nonglucosylated glycoprotein no longer interacts with calnexin/calreticulin and exits the ER more efficiently. As discussed below, a portion of glycoprotein species with glucosylated high mannose-type *N*-glycans manage to exit the RER, and they are processed by other enzymes within the Golgi apparatus. Many glycoproteins in the RER fail to fold completely and efficiently and are degraded within the RER by [proteinases](#). Thus, the glucosylation and deglucosylation pathway may serve a gate-keeper and quality control function.

## 7. Processing of *N*-glycans in the Golgi Apparatus

Following early-stage processing in the RER, glycoproteins exit the RER by vesicular transport to the Golgi apparatus. Not all glycoproteins containing glucosylated high mannose-type *N*-glycans are retained in the RER. For those that escape the action of glucosidase II in the RER, the Golgi contains an alternate pathway to facilitate their entry into the processing pathway. The Golgi enzyme, endomannosidase, acts on glucosylated oligosaccharides, such as  $\text{Glc}_1\text{Man}_9\text{GlcNAc}_2\text{Asn}$  and  $\text{Glc}_1\text{Man}_8\text{GlcNAc}_2\text{Asn}$  to release the disaccharide  $\text{Glc}\alpha$ 1-3Man and to generate  $\text{Man}_8\text{GlcNAc}_2\text{Asn}$  and  $\text{Man}_7\text{GlcNAc}_2\text{Asn}$ , respectively (21). The enzyme can act on a variety of glucosylated high mannose-type chains. The presence of the endomannosidase in the Golgi probably accounts for the inability of chemical inhibitors of glucosidases I and II to block completely the biosynthesis of complex-type *N*-glycans in all cells.

The resulting high mannose-type *N*-glycans, such as  $\text{Man}_8\text{GlcNAc}_2\text{Asn}$  and  $\text{Man}_7\text{GlcNAc}_2\text{Asn}$ , are acted on by a number of  $\alpha$ 1,2-mannosidases in the cis-Golgi (17), resulting in the release of two to four Man residues and the formation of the processing intermediate  $\text{Man}_5\text{GlcNAc}_2\text{Asn}$  (Fig. 1e). Although earlier observations had suggested that the Golgi contained a single  $\alpha$ 1,2-mannosidase

termed mannosidase-I, recent studies have shown that there are a variety of  $\alpha$ 1,2-mannosidases. The functions of these many different  $\alpha$ 1,2-mannosidases and their specific roles in the processing pathway are not yet clear.

The resulting  $\text{Man}_5\text{GlcNAc}_2\text{Asn}$  is a substrate for a key glycosyltransferase in the biosynthesis of hybrid- and complex-type *N*-glycans. *N*-Acetylglucosaminyltransferase I (GlcNAc-T-1) catalyzes the first committed step in the biosynthetic addition of sugars to the *N*-glycan by generating the GlcNAc $\beta$ 1-2Man linkage on a specific terminal Man residue (Fig. 1d). The action of GlcNAc-T-1 commits *N*-glycans to become either hybrid- or complex-type structures. A critical development in the discovery of a critical step in the biosynthetic process was the identification of mutated Chinese hamster ovary cell lines, termed Lec1 and Clone 15B, that are deficient in GlcNAc-T-1 (22, 23). As a consequence, the biosynthesis of *N*-glycans is halted at the level of  $\text{Man}_5\text{GlcNAc}_2\text{Asn}$ , and newly synthesized glycoproteins in these cell lines lack any hybrid- or complex-type *N*-glycans. Such mutational studies of Chinese hamster ovary cells have been exploited to generate many derivative cell lines containing mutations in specific aspects of the biosynthetic pathway for *N*- and *O*-glycans (24).

Following the action of GlcNAc-T-1, the product  $\text{GlcNAc}_1\text{Man}_5\text{GlcNAc}_2\text{Asn}$  is the specific substrate of the medial Golgi enzyme  $\alpha$ -mannosidase II, which removes two Man residues resulting in the formation of  $\text{GlcNAc}_1\text{Man}_3\text{GlcNAc}_2\text{Asn}$  (Fig. 1e). Without the action of  $\alpha$ -mannosidase II, *N*-glycans are committed to remain hybrid-type structures. Although early studies had suggested that there was a single  $\alpha$ -mannosidase II, a second gene encoding a second  $\alpha$ -mannosidase II has been identified (25). Experiments in mice in which the gene for the first  $\alpha$ -mannosidase II was deleted have confirmed the presence of this functional alternative  $\alpha$ -mannosidase II enzyme (26). The overall relationship between these two  $\alpha$ -mannosidases, and whether there may be other related enzymes, is currently unclear. The product of  $\alpha$ -mannosidase II action,  $\text{GlcNAc}_1\text{Man}_3\text{GlcNAc}_2\text{Asn}$ , is the substrate for a second *N*-acetylglucosaminyltransferase, GlcNAc-T-II, which adds GlcNAc in  $\beta$ 1,2-linkage to the remaining terminal  $\alpha$ -mannosyl residue to generate a biantennary complex-type *N*-glycan (Fig. 1e). This biantennary structure is a pivotal precursor for generation of all complex-type *N*-glycans.

## 8. Terminal Steps of *N*-Glycosylation

The biosynthesis of highly branched complex-type *N*-glycans from the biantennary precursor involves potential reactions with a series of other *N*-acetylglucosaminyltransferases within the medial Golgi apparatus (27). These enzymes, termed GlcNAc-T-III through GlcNAc-T-VII, add GlcNAc residues from donor UDPGlcNAc to the C-2, C-4, and C-6 positions of mannosyl residues. The resulting multiantennary products are indicated in Figure 1f. Not all *N*-glycans are subject to these reactions, since cells differ in their expression of the particular GlcNAc-Ts, and there are site-specific features of branching in glycoproteins.

Following completion of the GlcNAc-T activity within the medial-Golgi apparatus, glycoproteins exit to the trans-Golgi apparatus and trans-Golgi network (TGN), where additional modifications may occur through the actions of terminal glycosyltransferases and sulfotransferases. Some of the potential reactions are indicated in Figure 1g. Again, each cell type differs in the expression and level of these terminal glycosyltransferases, and there are site-specific features of terminal glycosylation in mature glycoproteins. In addition, many of the terminal glycosyltransferases may compete with each other for similar substrates.

In animal cells, it is typical for complex-type *N*-glycans to contain outer galactose (Gal) and sialic acid residues (28). A common sialic acid in humans is *N*-acetylneuraminic acid (NeuAc), but dozens of different sialic acids are known, and many are expressed by other animals. Interestingly, sialic acid is not generally present in plant glycoproteins. Terminal glycosylation of *N*-glycans in plant

glycoproteins is different in several ways from those seen in animal glycoproteins (29). Many complex-type *N*-glycans in animal glycoproteins contain a primary terminal motif Galb1-4GlcNAc, known as the lactosamine or LN sequence (Fig. 1e). A second motif, however, known as the lacdiNAc or LDN pathway, has now been identified (30). LDN structures contain the alternative disaccharide GalNAcb1-4GlcNAc-R (Fig. 1e). The outer sugar residues may also be sulfated, as seen, for example, in pituitary hormones, which contain GalNAc-4-sulfate residues (31) (Fig. 1e), and in keratan sulfates, which contain GlcNAc-6-sulfate and Gal-6-sulfate residues. Sulfation of glycoproteins occurs in the trans-Golgi apparatus through the donor phosphoadenosine phosphosulfate (PAPS). PAPS, like the sugar nucleotide donors for glycosyltransferase reactions in the RER and Golgi apparatus, must be transported into those organelles by membrane transporters. These transporters, which are specific for each donor molecule, have been identified by molecular approaches (32).

## 9. Phosphorylation of *N*-Glycans in Lysosomal Enzymes

In higher eucaryotes, lysosomal hydrolases are directed to the lysosome by a mechanism that depends on specific modification of *N*-linked glycans, which allows sorting of these glycoproteins from others in the secretory pathway. The *N*-linked glycans of lysosomal hydrolases are specifically and uniquely modified in two steps to contain one or more terminal mannose 6-phosphates (M6P) (Fig. 1h). Lysosomal hydrolases bearing the M6P modification then bind to [Mannose-6-P-receptors](#) located in the trans-Golgi and are transferred to the [endosome](#) by vesicular transport. Once in the endosome, because of the lower pH, the lysosomal hydrolase disassociates from the receptor and is transferred to the lysosome. Once in the lysosome, lysosomal phosphatases remove the phosphate, rendering the process irreversible (33).

The first enzyme in the biosynthesis of the M6P modification is UDP-*N*-acetylglucosamine:lysosomal-enzyme *N*-acetylglucosamine-1-phosphotransferase (abbreviated GlcNAc-phosphotransferase). GlcNAc-phosphotransferase catalyzes the transfer of GlcNAc 1-phosphate from UDP-GlcNAc to specific  $\alpha$ 1,2-linked mannoses on high mannose oligosaccharides of lysosomal hydrolases. GlcNAc-phosphotransferase is a complex enzyme composed of three polypeptide chains with the subunit structure  $\alpha_2\beta_2\gamma_2$  (34, 35). The specificity of the M6P-dependent lysosomal targeting pathway is a reflection of the substrate specificity of GlcNAc-phosphotransferase, which can only effectively utilize native lysosomal hydrolases as substrates.

The second enzyme in the pathway, *N*-acetylglucosamine-1-phosphodiester  $\alpha$ -*N*-acetylglucosaminidase (abbreviated phosphodiester  $\alpha$ -GlcNAcase), removes the covering GlcNAc, generating a terminal M6P. Lysosomal enzymes bearing the M6P modification then bind to one of two M6P-receptors in the trans-Golgi network and are transferred to the lysosome. Phosphodiester  $\alpha$ -GlcNAcase is a tetramer of identical subunits (36). Unlike GlcNAc-phosphotransferase, phosphodiester  $\alpha$ -GlcNAcase does not exhibit any specificity for lysosomal hydrolases as substrates.

Two distinct M6P-receptors have been isolated and characterized (37) (see [Mannose 6-P Receptors](#)). A high molecular weight cation-independent M6P-receptor (CI-MPR) was first described. The CI-MPR functions as a monomer, and each monomer binds two mols of M6P. In addition to binding M6P, the CI-MPR is also a high affinity receptor for [insulin-like growth factor II](#) (IGF-II) and the urokinase-type plasminogen activator receptor (uPAR). A second receptor has been described that requires divalent cations for binding of M6P (CD-MPR). The CD-MPR functions as a dimer, with each dimer binding two mols of M6P.

Several human diseases of lysosomal hydrolase trafficking result from genetic defects in the enzymes responsible for the biosynthesis of the M6P modification. Mucopolidosis II (MLII) or I-cell disease and mucopolidosis IIIA (classic pseudo-Hurler polydystrophy) result from mutations in the GlcNAc-phosphotransferase  $\alpha$ / $\beta$ -subunit gene. Mucopolidosis IIIC (MLIIIC) or variant pseudo-Hurler polydystrophy results from mutations in the GlcNAc-phosphotransferase  $\gamma$ -subunit gene. In



the absence of GlcNAc-phosphotransferase, lysosomal enzymes fail to acquire M6P; they enter the secretory pathway, and are secreted from the cell. This mistargeting of lysosomal hydrolases results in a lysosomal deficiency of many lysosomal enzymes and a severe storage phenotype. Patients with partial deficiency of phosphodiester  $\alpha$ -GlcNAcase have been described. Like patients with GlcNAc-phosphotransferase deficiency, these patients have greatly elevated serum lysosomal enzyme levels, but do not seem to have a disease as a result of the mistargeting of lysosomal hydrolases.

The presence of the M6P or GlcNAc-1-phosphate modification blocks mannose trimming by Golgi  $\alpha$ 1,2-mannosidases. This prevents the synthesis of complex-type oligosaccharides and results in the high mannose structures typically found on lysosomal hydrolases.

## 10. Inhibitors of *N*-Glycosylation

Since the discovery of tunicamycin and its potent inhibition of the *N*-glycosylation pathway, a number of other antibiotics and sugar analogs have been identified that demonstrate specific inhibition of various enzymatic steps in the biosynthetic pathway (38). Most of the inhibitors commonly used today block glycosidase reactions in glycoprotein processing (Table 1). These include the inhibitors of glucosidases I and II,  $\alpha$ 1,2-mannosidases, the endomannosidase, and  $\alpha$ -mannosidase II. Some of these inhibitors demonstrate a high degree of specificity, whereas others are more broadly active toward several enzymes. Another drug that has proven useful in studying *N*-glycosylation pathways is [Brefeldin A](#) (BFA). BFA is a fungal metabolite that interferes with protein trafficking between the Golgi apparatus and the trans-Golgi network (TGN), resulting in a blockage in glycoprotein secretion. BFA does not directly inhibit steps in the assembly or processing of *N*-glycans. In general, BFA treatment results in decreased sialylation and fucosylation of glycoproteins in those cells in which these terminal glycosylation steps occur in the TGN.

**Table 1. Inhibitors of *N*-Glycosylation and Their Action Sites**

| <b>Inhibitor</b>  | <b>Action</b>  |
|---|--|
| Tunicamycin   | Inhibits formation of dolichol pyrophosphoryl-GlcNAc from dolichol phosphate and UDPGlcNAc |
| Australine  |  |
| Castanospermine   |  |
| 1-Deoxynojirimycin  |  |
| <i>N</i> -Methyldeoxynojirimycin  |  |
| DMDP (2 <i>R</i> ,5 <i>R</i> -bis (Hydroxymethyl)-3 <i>R</i> ,4 <i>R</i> -dihydropyrrolidine) | Inhibits glucosidase I   |
| Bromoconduritol   |  |
| Castanospermine   |  |
| 1-Deoxynojirimycin (DNJ)  |  |
| <i>N</i> -Methyldeoxynojirimycin  | Inhibits glucosidase II  |
| 1-Deoxymannojirimycin (DMJ)   |  |
| <i>N</i> -Methyl-deoxynojirimycin   |  |
| DMDP  |  |
| Kifunensine   | Inhibits $\alpha$ 1,2-mannosidases   |
| Glc $\alpha$ 1-3DMJ (GDMJ)  | Inhibits endomannosidase   |

Swainsonine  
Mannostatin  
Brefeldin A

Inhibits  $\alpha$ -mannosidase II  
Blocks translocation of proteins from ER to Golgi by blocking binding of GDP-ribosylation factor to Golgi and inhibiting GTP-GDP exchange

---

## 11. Proofreading of N-glycans

*N*-glycans may have other functions yet to be discovered in early steps of glycoprotein biosynthesis. The presence of specific endoglycosidases in animal cells capable of removing *N*-glycans from glycoproteins (39) suggested that some *N*-glycans may be removed after their addition to protein. These enzymes may be involved in general degradation of glycoproteins. It is also possible, however, that some sites are *N*-glycosylated and the *N*-glycan is subsequently removed in a “quality control” pathway. An enzyme has been identified in animal cells that is capable of catalyzing the de-*N*-glycosylation of proteins, known as a peptide:*N*-glycanase or PNGase; it may be involved in quality control of newly synthesized proteins by removing *N*-glycans from specific sites in a glycoprotein precursor, thereby reducing the number of *N*-glycans in the mature form of the protein. It is presently unclear whether this quality-control pathway operates on all glycoproteins in the ER or whether it specifically recognizes certain *N*-glycans in a restricted set of glycoproteins.

## 12. Genetic Defects in *N*-glycosylation

In recent years, a number of specific defects in the *N*-glycosylation pathway have been identified in humans. As mentioned above, the loss of the UDPGlcNAc:lysosomal enzyme *N*-acetylglucosamine-1-phosphotransferase in patients with I-cell disease results in a deficiency of lysosomal hydrolases and a generalized lysosomal storage disease with profound clinical consequences. In recent years, a number of individuals have been identified with decreased *N*-glycosylation of circulating glycoproteins. This deficiency, carbohydrate-deficient glycoprotein syndrome (CDGS), is due to blockage in several steps in the *N*-glycosylation pathway. CDGS is a genetic disease usually associated with major involvement of the central nervous system and mental retardation (40). However, some patients present with gastrointestinal disorder characterized by protein-losing enteropathy. Diagnosis of the disease usually involves identification of decreased glycosylation of serum [transferrin](#), a circulating glycoprotein that normally contains two complex-type *N*-glycans. In CDGS type II patients, there is a deficiency in the Golgi processing enzyme *N*-acetylglucosaminyltransferase II (GlcNAc-T-II). Consequently, these patients synthesize glycoproteins containing the normal number of *N*-glycans, but the glycans have a hybrid-type structure. In contrast, CDGS type I is a more complex disease. Some patients with CDGS type I are unable to synthesize glucosylated dolichol-linked oligosaccharide, leading to accumulation of  $_9\text{GlcNAc}_2\text{-P-P-dolichol}$  (41). Studies in many systems have shown that the oligosaccharyltransferase is unable to utilize efficiently the  $\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  as a donor in comparison with  $\text{Glc}_3\text{Man}_9\text{GlcNAc}_2\text{-P-P-dolichol}$  (42, 43). Other patients with CDGS type I lack phosphomannomutase, which is required for conversion of mannose-6-phosphate to mannose-1-phosphate in the biosynthesis of GDPMan (44). Still other patients with CDGS type I are deficient in phosphomannose isomerase, the enzyme required for conversion of fructose-6-phosphate to mannose-6-phosphate, eventually leading to GDPMan (45).

Hereditary erythroblastic multinuclearity with positive acidified serum (HEMPAS) lysis test (also known as congenital dyserythropoietic anemia type II) is a disease characterized by decreased *N*-

glycosylation of glycoproteins on erythrocytes, which may be due to a deficiency  $\alpha$ -mannosidase II (46). This decreased glycosylation is associated with crenation of the erythrocytes and subsequent hemolysis. Consistent with the possibility that human HEMPAS is associated with deficiency in  $\alpha$ -mannosidase II is the finding that mice lacking a functional  $\alpha$ -mannosidase II gene also develop a dyserythropoietic anemia associated with decreased *N*-glycosylation of erythrocyte glycoproteins (26).

Targeted genetic mutations in mice involving steps in the *N*-glycosylation pathway have been of great interest in recent years as scientists in the field attempt to unveil the biological functions of *N*-glycans. For example, [null mutations](#) in GlcNAc-T-I result in embryonic lethality (47, 48). It is hoped that such studies will begin to uncover many other biological functions of *N*-glycans and other potential enzymes involved in their biosynthesis.

## Bibliography

1. R. D. Cummings (1992) in *Glycoconjugates: Composition, Structure and Function* (H. J. Allen and E. C. Kisailus, eds.), Marcel Dekker, New York, pp. 333–360.
2. R. D. Marshall and A. Neuberger (1972) *Annu. Rev. Biochem.* **41**, 673–702.
3. B. A. Vance, W. Wu, R. K. Ribaud, D. M. Segal, and K. P. Kears (1997) *J. Biol. Chem.* **272**, 23117–23122.
4. E. Bause (1983) *Biochem. J.* **209**, 331–336.
5. J. L. Mellquist, L. Kasturi, S. L. Spitalnik, and S. H. Shakin-Eshleman (1998) *Biochemistry* **37**, 6833–6837.
6. P. Messner (1997) *Glycoconj. J.* **14**, 3–11.
7. J. Lechner and F. Wieland (1989) *Annu. Rev. Biochem.* **58**, 173–194.
8. R. Kornfeld and S. Kornfeld (1985) *Annu. Rev. Biochem.* **54**, 631–664.
9. C. B. Hirschberg and M. D. Snider (1987) *Annu. Rev. Biochem.* **56**, 63–87.
10. J. S. Schutzback (1997) *Glycoconj. J.* **14**, 175–182.
11. J. S. Tkacz and O. Lampen (1975) *Biochem. Biophys. Res. Commun.* **65**, 248–257.
12. E. Bause and G. Legler (1981) *Biochem. J.* **195**, 639–644.
13. B. Imperiali and T. L. Hendrickson (1995) *Bioorg. Med. Chem.* **3**, 1565–1578.
14. S. Silberstein and R. Gilmore (1996) *FASEB J.* **10**, 849–858.
15. S. te Heesen, B. Janetzky, L. Lehle, and M. Aebi (1992) *EMBO J.* **11**, 2071–2075.
16. J. H. Chi, J. Roos, and N. Dean (1996) *J. Biol. Chem.* **271**, 3132–3140.
17. K. W. Moremen, R. B. Trimble, and A. Herscovics (1994) *Glycobiology* **4**, 113–125.
18. A. J. Parodi (1998) *Braz. J. Med. Biol. Res.* **31**, 601–614.
19. D. N. Hebert, B. Foellmer, and A. Helenius (1995) *Cell* **81**, 425–433.
20. F. E. Ware, A. Vassilakos, P. A. Peterson, M. R. Jackson, M. A. Lehrman, and D. B. Williams (1995) *J. Biol. Chem.* **270**, 4697–4704.
21. M. J. Spiro, V. D. Bhoyroo, and R. G. Spiro (1997) *J. Biol. Chem.* **272**, 29356–29363.
22. C. Gottlieb, J. Baenziger, and S. Kornfeld (1975) *J. Biol. Chem.* **250**, 3303–3309.
23. P. Stanley (1985) *Mol. Cell Biol.* **5**, 923–929.
24. P. Stanley (1983) *Meth. Enzymol.* **96**, 157–184.
25. M. Misago, Y. F. Liao, S. Eto, M. G. Mattei, K. W. Moremen, and M. N. Fukuda (1995) *Proc. Natl. Acad. Sci. USA* **92**, 11766–11770.
26. D. Chui, M. Oh-Eda, Y. F. Liao, K. Panneerselvam, A. Lal, K. W. Marek, H. H. Freeze, K. W. Moremen, M. N. Fukuda, and J. D. Marth (1997) *Cell* **90**, 157–167.
27. I. Brockhausen and H. Schachter (1997) In *Glycosciences: Status and Perspectives* (H. J. Gabius and S. Gabius, eds.), Chapman and Hall, Weinheim, Germany, pp. 79–113.

28. A. Varki (1998) Trends Cell Biol. **8**, 34–40.
29. P. Lerouge, M. Cabanes-Macheteau, C. Rayon, A. C. Fishette-Laine, V. Gomord, and L. Faye (1998) Plant Mol. Biol. **38**, 31–48.
30. D. H. Van den Eijnden, A. P. Neeleman, H. Bakker, and I. Van Die (1998) Adv. Exp. Med. Biol. **435**, 3–7.
31. L. V. Hooper, S. M. Manzella, and J. U. Baenziger (1996) FASEB J. **10**, 1137–1146.
32. E. Guillen, C. Abeijon, and C. B. Hirschberg (1998) Proc. Natl. Acad. Sci. USA **95**, 7888–7892.
33. S. Kornfeld (1987) FASEB J. **1**, 462–468.
34. M. Bao, J. L. Booth, B. J. Elmendorf, and W. M. Canfield (1996) J. Biol. Chem. **271**, 31437–31445.
35. M. Bao, B. J. Elmendorf, J. L. Booth, R. R. Drake, and W. M. Canfield (1996) J. Biol. Chem. **271**, 31446–31451.
36. R. Kornfeld, M. Bao, K. Brewer, C. Noll, and W. M. Canfield (1998) J. Biol. Chem. **273**, 23203–23210.
37. N. M. Dahms, P. Lobel, and S. Kornfeld (1989) J. Biol. Chem. **264**, 12115–12118.
38. G. P. Kaushal and A. D. Elbein (1994) Meth. Enzymol. **230**, 316–329.
39. T. Suzuki, K. Kitajima, Y. Emori, Y. Inoue, and S. Inoue (1997) Proc. Natl. Acad. Sci. USA **94**, 6244–6249.
40. J. Jaeken (1991) Int. Pediatr. **6**, 56–58.
41. P. Burda, L. Borsig, J. de Rijk-van Andel, R. Wevers, J. Jaeken, H. Carchon, E. G. Berger, and M. Aebi (1998) J. Clin. Invest. **4**, 647–652.
42. L. A. Murphy and R. G. Spiro (1981) J. Biol. Chem. **256**, 7487–7494.
43. J. Stoll, R. Cacan, A. Verbert, and S. S. Krag (1992) Arch. Biochem. Biophys. **299**, 225–231.
44. E. Van Schaftingen and J. Jaeken (1995) FEBS Lett. **377**, 318–320.
45. R. Niehues, M. Haslik, G. Alton, C. Korner, M. Schiebe-Sukumar, H. G. Koch, K. P. Zimmer, R. Wu, E. Harms, K. Reiter, K. von Figura, H. H. Freeze, H. K. Harms, and T. Marquardt (1998) J. Clin. Invest. **101**, 1414–1420.
46. M. N. Fukuda, G. F. Gaetani, P. Izzo, P. Scartezzini, and A. Dell (1992) Br. J. Haematol. **82**, 745–752.
47. M. Metzler, A. Gertz, M. Sarkar, H. Schachter, J. W. Schrader, and J. D. Marth (1994) EMBO J. **13**, 2056–2065.
48. E. Ioffe and P. Stanley (1994) Proc. Natl. Acad. Sci. USA **91**, 728–732.

## O-Glycosylation

In **eukaryotes**, carbohydrate side-chains are added to certain **amino acid** residues in some **protein** sequences to form **glycoproteins**. In **N-glycosylation**, the sugars are attached in N-linkage to **asparagine residues**, and this type of glycosylation, which is initiated in the **endoplasmic reticulum** (ER), has been well studied. Glycans may also be added to **serine** and **threonine residues** in a protein, and because the hydroxyl group of these amino acids forms the link with the sugar, this type of glycosylation has been termed O-glycosylation. The **enzymes** that initiate glycosylation and add the individual sugars to extend or terminate the chain are glycosyltransferases, which catalyze the basic reaction:

nucleotide – sugar donor + R-OH receptor → nucleotide + R-O-sugar

where R is a serine or threonine residue or a sugar already attached to the acceptor protein. In most cases, the nucleotide may be UDP, as for galactose and the hexosamines, or GDP, as for fucose; these donor substances are made in the **cytoplasm**. An exception is sialic acid, which is added from cytidine monophosphate-sialic acid that is made in the **nucleus**. There are three main types of O-glycosylation in eukaryotes, namely: (1) Mucin-type O-glycosylation, in which the proteins normally carry multiple O-linked oligosaccharide chains; (2) O-Glycosylation leading to the formation of proteoglycans; and (3) O-GlcNAc glycosylation of cytoplasmic and nuclear proteins. This article will focus on O-glycosylation in mammalian cells, especially relating to mucin-type O-glycosylation. Proteins are O-glycosylated, however, even in fungal cells, and these pathways will also be discussed.

## 1. Mucin-Type O-Glycosylation

This type of glycosylation refers to the covalent attachment of O-glycans to serines and threonines in mucin core proteins, whereby the sugars are added individually and sequentially in the **Golgi apparatus**. The nucleotide sugar donors are transported into the lumen of the Golgi pathway where the glycosyltransferases are positioned. The genes coding for many of the relevant glycosyltransferases have recently been isolated, **cloned**, and expressed, and it is becoming clear from their specificities that several different enzymes may catalyze the same reaction. Moreover, the existence of enzymes that can add different sugars to the same substrate means that competition can occur for the substrate, providing the locations of the different enzymes overlap in the Golgi pathway. Thus, the locations of the enzymes involved in mucin-type O-glycosylation are as important as their level of activity in determining the final composition of O-glycans added.

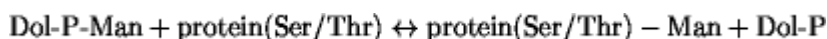
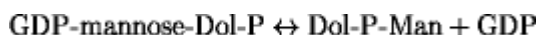
### 1.1. Initiation of Glycosylation

In mammalian cells and as far down the evolutionary scale as *Caenorhabditis elegans*, the first sugar to be added to serine or threonine in alpha linkage is N-acetylgalactosamine (GalNAc). Identifying a specific sequence that defines a glycosylation site has been difficult. Recent developments suggest that this may be largely because not one enzyme but a family of enzymes (UDP-N-acetyl-a-D-galactosamine:polypeptide-N-Acetylgalactosaminyltransferases) are involved in this reaction. These enzymes, referred to here as ppGalNAcTs, catalyze the transfer of GalNAc from UDP-GalNAc to serines and threonine residues (1). The different ppGalNAcTs have distinct but overlapping specificities for the peptide sequence, and the same enzyme can add GalNAc to both serine and threonine residues (2). Thus, the sequences flanking the serine or threonine residue influence whether or not it is glycosylated, but this cannot be inferred merely from a **database** analysis of in vivo substrates (3).

**Knockout strains** of the ppGalNAcTs are only just being prepared (4), and some redundancy is to be expected. However, the ppGalNAcTs show well-defined tissue and cell specificity of expression (5), suggesting a differential dependence of each cell type on specific enzymes. The tissue and cell specificity of the ppGalNAcTs also suggests that differences may exist in the sites of glycosylation of the same protein expressed in different cells. Identification of sites that are glycosylated in vivo has been difficult, but some data are now beginning to emerge.

The location in the ER/Golgi pathway where mucin type O-glycosylation is initiated has been controversial. There has been general agreement that GalNAc is not added in the ER, but different methods of analysis have indicated different locations (6, 7). Recent studies using **immunoelectron microscopy** to localize tagged enzymes in transfected cells suggest that three of the ppGalNAcTs (T1, T2, and T3) are found throughout the Golgi, with the profile of expression being different for each enzyme (8). With mucin proteins, in which hundreds of sugar side chains may be added, it may be that the addition of GalNAc does indeed occur throughout the Golgi. If some chains are initiated later in the pathway, this could contribute to the heterogeneity of composition of the O-glycans that are found in mucin preparations.

In **fungi**, serine and threonine residues are glycosylated via dolichyl phosphate-D mannose (Dol-P-Man) as an intermediate, a reaction that occurs in the ER and has not been observed thus far in higher eukaryotes. The following reactions have been established for *Saccharomyces cerevisiae*:



The protein molecules that are glycosylated may be secreted or are in the cell wall. Although the classification of the yeast O-glycosylation may not fall into any of the categories defined for higher eukaryotes, multiple genes code for enzymes catalyzing the transfer of mannose from Dol-P-Man to serine or threonine residues (Dol-P-Man:protein mannosyl transferases, or PMTs). In this system, the possible redundancy of the enzymes has been investigated by producing multiple mutants, and only a triple disruption is lethal, although some of the double mutants are unable to grow without osmotic stabilization (8).

### 1.2. Chain Extension in Mucin-Type O-Glycan Synthesis

After the addition of GalNAc to the mucin protein, various core structures are formed by the addition of different sugars (see Fig. 1 of [O-Linked Oligosaccharides](#)). Chains are extended from these cores by the addition of N-acetylglucosamine (GlcNAc) and galactose alternately, to give polylactosamine side chains that may be straight or branched (see Figs. 2 and 3 in [O-Linked Oligosaccharides](#)). Formation of the different cores varies with the tissue, but a pathway through core 1 and core 2 is commonly used. Two genes that can add GlcNAc in  $\beta$ 1,6 linkage to core 1 (to form core 2) have been identified (9). One of these enzymes can also catalyze the addition of GlcNAc to form the internal branch (shown in Fig. 2 of [O-linked oligosaccharides](#)), and the other (the core 2 enzyme) can only catalyze the formation of core 2 from core 1. Where biosynthesis is through core 2, the addition of the GlcNAc is crucial for chain extension.

### 1.3. Chain Termination

The O-glycans carried on mucins are usually terminated by the addition of sialic acid or fucose, and this may be subject to blood group A or B dependent transferase activity (see Fig. 3A and 3B of [O-Linked Oligosaccharides](#)). The number of enzymes that catalyze the addition of sialic acid is steadily increasing and presently is at least ten (10). Not all of these enzymes act on O-glycans, but those acting at the early stages of O-glycosylation are specific for this type of glycosylation. Important enzymes in this category are those adding sialic acid in  $\alpha$ 2,6 linkage to the first sugar GalNAc, or in  $\alpha$ 2,3 linkage to galactose in core 1 (see Fig. 2 of [O-Linked Oligosaccharides](#)). The sialyltransferase that adds sialic acid in  $\alpha$ 2,6 linkage to galactose at the end of N-glycans does not add the sugar to O-glycan chains that terminate with sialic acid in  $\alpha$ 2,3 linkage. Two candidate enzymes exist (ST3 Gal III and ST3 Gal IV) that could add sialic acid to galactose at the end of O-glycans, which have been reported to be selective for type I and type 2 chains, respectively. (See Tsuji et al (1996) for further discussion on nomenclature.)

### 1.4. Localization of Glycosyltransferases Involved in the Biosynthesis of O-Glycans in the Golgi Pathway

Organization of the Golgi apparatus can be separated into three major components: (1) the *cis*-Golgi network, (2) the Golgi stacks (*cis* medial and *trans*), and (3) the *trans* Golgi network, or TGN. The sizes of the various components of the pathway vary with cell type (see [Golgi Apparatus](#)). Movement of proteins through the pathway is vectorial, progressing from *cis* to *trans* via nonselective budding transport. Glycosylated proteins leave the *trans* face of the stack to be sorted to their various destinations in the TGN, where further glycosylation may also occur (11).

The glycosyltransferases are themselves type-II membrane glycoproteins, with a short amino-terminal cytoplasmic **domain**, a single transmembrane domain (typically 17 residues), a loosely

folded putative stem cell region (50–100 residues), and a tightly folded globular catalytic domain of more than 325 residues, which extends into the lumen of the Golgi apparatus. Early work involving the localization of these enzymes using subcellular fractionation or detection of specific oligosaccharide structures with [lectins](#) suggested an orderly compartmentalization of enzymes corresponding to the sequential addition of sugars. More recent studies show that glycosyltransferase organization is more complex, however, and that overlap can occur even between enzymes such as the core 2 enzyme responsible for chain extension and the  $\alpha$ 2,3 sialyl-transferase ST3Gal I that effects termination ([12](#)). These enzymes both use Gal b1,3 -GalNAc-R as a substrate, and changes in their levels of activity affect the final composition of the O-glycans synthesized. Definitive mapping has only been done for a few of the many enzymes involved in O-glycan initiation and biosynthesis. However, with the production of the recombinant enzymes and, in particular, specific [monoclonal antibodies](#) directed against them, the relative positions of these enzymes should be established.

The mechanisms involved in sorting the glycosyltransferases, whether these be involved in N- or O-glycosylation, are not fully understood. Two models have been suggested to explain the distribution of these enzymes in the Golgi pathway ([13](#), [14](#)) (see [Golgi Apparatus](#)).

### 1.5. Molecules Carrying Mucin-Type O-Glycans

Although some molecules (eg, erythropoietin) may carry only one O-glycan, most of the molecules carrying carbohydrate side chains attached in this way carry multiple O-glycans and are classified as mucins; the final composition of the molecule is more than 50% carbohydrate. These molecules have a common feature in containing a tandem repeat domain, rich in serine and threonine residues, to which the O-glycans are attached. Of the eight human epithelial mucin genes that have been identified by gene cloning, all but the first (MUC1) are extracellular mucins. MUC1 is a type I transmembrane protein and resembles the membrane-associated selectin ligands ([15](#)). It is much larger, however, and the sequence in the repeats is much more conserved. The addition of sugars, even only GalNAc, results in the molecule becoming highly extended, reaching above the glycocalyx, allowing the membrane-associated mucins to influence cell-cell interactions.

The multiplicity of enzymes that can initiate mucin-type O-glycosylation and the sequential nature of the biosynthesis of the O-glycans allow for an almost infinite number of possible glycoforms to be produced from the same core protein, the final structure depending on the profile of expression of the glycosyltransferases. Studies with specific [antibodies](#) to the ppGalNAcTs are showing a marked tissue and cell specificity of expression ([1](#), [5](#)), and this also holds true for many of the enzymes involved in biosynthesis of the oligosaccharides, for which studies thus far have depended on looking at [messenger RNA](#) expression. Studies on the promoters governing expression of the transferases are only now beginning, but strong indications exist that tissue-specific expression may in some cases depend on the use of alternative **promoters** and [alternative splicing](#), resulting in the production of various mRNAs with different 5'-untranslated regions ([16](#)).

### 1.6. Functions of Mucin-Type Glycoproteins

An obvious function of the extracellular mucins, which form an important component of the mucous layer covering some epithelial cells (eg, lining the gastrointestinal and respiratory tracts), is protection of the underlying epithelium from insult. The O-glycans play a major role in this protective function, not only in the formation of large oligomers, which makes the mucous layer viscous, but also in serving as receptors for invading microorganisms. The membrane-associated MUC1 mucin can also have a protective function but, like the selectin ligands, can also affect cell-cell adhesion.

The O-glycans on the membrane-associated mucins are generally heavily sialylated, and the resulting negative charge results in inhibition of cell-cell interactions unless a specific [epitope](#) in the oligosaccharide component can act as a ligand for a **receptor** on another cell, when cell-cell adhesion is enhanced. Study of the interaction of selectins with their glycoprotein ligands has led to detailed definition of the interacting oligosaccharide structure (see [O-Linked Oligosaccharides](#)). However, the specific oligosaccharide must be carried on a particular core protein, suggesting an

involvement of the protein sequence or structure, either directly or indirectly by allowing clustering of the O-glycans.

### 1.7. Glycosylation Patterns, Differentiation, and Malignancy

The composition of the side chains added to a mucin core protein can vary not only in different tissues, but also with the differentiation state of a particular cell **phenotype** and with the change to malignancy (see [Neoplastic Transformation](#)). Although the changes in malignancy may or may not relate to disease progression, those observed in differentiation are likely to be functional and relate to cell adhesion.

Changes in the expression of the core 2 enzyme have been characterized both in differentiation of [T cells](#) and in cancer. Leukosialin is a major glycoprotein expressed on leukocytes, and changes in its composition of added O-glycans occur on activation of human T cells. In resting cells, the dominant O-glycan is a tetrasaccharide (disialylated core 1), whereas in activated T cells the expression of the core 2 enzyme is increased and core 2-based structures are added ([17](#)). A similar increase in the activity of the core 2 enzyme is seen in leukocytes from patients with immunodeficiency, as such that seen in Wiscott-Aldrich syndrome ([9](#)). Furthermore, the core 2 enzyme is differentially expressed in the thymus, where expression is high in the subcapsular and cortical thymocytes and low in the medullary thymocytes. In this case, the differential expression of core 2 structures correlates with the ability of the cells to interact with a lectin (galectin) synthesized by thymic epithelial cells.

Emphasizing the importance of cell phenotype in defining profiles of glycosylation, the changes in levels of enzymes in malignancy differ in various cancers. In leukemias, the core 2 enzyme can be increased, whereas it may be decreased in breast cancers ([18](#)). In breast cancer, there is also an increase in the  $\alpha$ 2,3 sialyl-transferase that adds sialic acid to core 1, resulting in an increase in the major epitope for sialoadhesin expressed on macrophages. In the colon, where the O-glycans added to the mucins are very large and complex, changes occur in carcinomas that result in increased expression of the sialylated Le<sup>a</sup> and Le<sup>x</sup> structures—the ligands for P and E selectins, respectively.

## 2. O-glycosylation leading to the formation of proteoglycans

Proteoglycans are composed of glycosaminoglycan chains (GAGs) bound to serine residues in a protein core, via a xylose-galactose-galactose bridge. Unlike mucin-type O-glycosylation, threonine residues do not seem to act as acceptors. All GAGs (chondroitin sulfate, dermatan sulfate, heparin and heparan sulfate) except hyaluronic acid are secreted as components of proteoglycans. The post-translational processing of the core protein occurs in the Golgi apparatus ([19](#)), where, after chain initiation, monosaccharides are added stepwise from the appropriate UDP-sugars. To attain their final shape, chains are then carried through a series of modifications, including [sulfation](#) (see [O-Linked Oligosaccharides](#)). There has been more progress in isolating the genes coding for and defining the core proteins of proteoglycans than in isolating the genes coding for the enzymes involved in the biosynthesis of the carbohydrate moieties added to these proteins. The GAGs are large, and where many chains are added there may be only 10% protein in the proteoglycan.

### 2.1. Initiation of Glycosylation

The initial, rate-limiting step in the synthesis of GAGs is the transfer of xylose from UDP-xylose to a serine residue. This step is catalyzed by UDP-D-xylose: proteoglycan core protein b-D-xylosyl transferase (xyloseT). Attempts to define a sequence in the core protein that will determine whether xylose will be added suggest that, although some limitations can be placed on the sequence flanking the serine (glycine must follow toward the carboxyl end), protein **secondary structure** may play an important role ([20](#)). The linking galactose moieties are then added, and the type of GAG to be synthesized is determined by the first hexosamine to be added, which begins the biosynthesis of the main chain.

### 2.2. Chain Extension



The GAGs consist of hexosamines and either hexuronic acid or L-iduronic acid or galactose units added alternately in unbranched sequence (see [O-Linked Oligosaccharides](#)). It seems likely that the enzymes that add the first hexosamine (GalN or GlcN) are different from those acting on the more peripheral regions of the chain. Which hexosamine is added determines which GAG is subsequently synthesized, and sequences in the core protein may direct the choice, as different GAGs can be added to different core proteins in the same cell.

### 2.3. Functions of Proteoglycans

Proteoglycans may occur intracellularly in **secretory granules**, at the cell surface, and in the [extracellular matrix](#). The functions of proteoglycans are highly diverse, ranging from mechanical functions essential for maintaining the structural integrity of connective tissue, to effects on cell adhesion, motility, proliferation, differentiation, and morphogenesis. Many of these effects depend on binding of proteins to the GAG chains. As with the mucin molecules, these interactions can depend on the charge and are then relatively nonspecific and of low affinity, whereas others, involving a particular oligosaccharide with defined structure, are highly specific. The core protein, in addition to serving as a scaffold for the GAGs, may be involved in anchoring the molecule to the membrane.

### 3. O-GlcNAc glycosylation of cytoplasmic and nuclear proteins

O-linked N-acetylglucosamine linked to serine or threonine residues was discovered by Torres and Hart ([21](#)) and has since been found to be ubiquitous and abundant on nuclear and **cytoskeletal** proteins in virtually all eukaryotes, including fungi ([22](#)). This type of glycosylation occurs in the cytoplasm, not in the Golgi apparatus (the site of post-translation modification of the core proteins of mucins and proteoglycans), and it may be extremely important in modifying the activity of intracellular proteins with a wide diversity of functions. Table [1](#) lists some of the classes of proteins that have been shown to be glycosylated in this way.

**Table 1. Some Proteins Shown to be Subject to O-GlcNAc Acylation**

---

Proteins involved in transcription (Pol II and transcription factors)

---

Cytoskeletal proteins (intermediate filaments, bridging proteins)

Tumor suppressors, oncogenes (p53)

Nuclear pore proteins

---

#### 3.1. O-GlcNAc and Phosphorylation

O-GlcNAc glycosylation is dynamic; the sugar turns over much more rapidly than does the protein backbone. Most of the proteins undergoing this type of glycosylation are **phosphorylated**, in some instances on the same amino acid residue, so it is suggested that the addition of O-GlcNAc is a regulatory modification analogous to phosphorylation. Many of the sites are similar to those used by some [kinases](#), namely the glycogen synthase kinase and the MAP kinases, and in the *myc* oncogene the site has been mapped to Thr58 in the transactivation domain ([23](#)), which is also a major site of phosphorylation and a hot spot for [mutagenesis](#) in lymphomas. The glycosylation of **RNA polymerase II** at the mucin-like sequence at the C-terminal domain may indicate a role in transcriptional initiation.

#### 3.2. O-GlcNAc and Protein Interactions

Another function of O-GlcNAc appears to be in mediating cytoskeletal assembly and organization,

and a defect in the function has been implicated in Alzheimer's disease (24). Where the GlcNAc is thought to be an alternative to phosphorylation, it is accessible to modification (by a glycosyltransferase that transfers galactose to the hexosamine). When O-GlcNAc is functioning in protein-protein interactions, the sugars appear to be buried in the native molecules and are only accessible after **denaturation** or **proteolysis**.

Whereas about 50% of the sites that are glycosylated are at or near a PVS (Pro-Val-Ser) type of sequence, the other half have no apparent **consensus sequence**. It is not yet clear how many enzymes are involved in this type of glycosylation. However, a gene coding for an O-GlcNAc transferase has now been isolated and characterized. (See Hart et al (1996) for a further discussion).

### 3.3. Concluding Comments

Addition of sugars in O-linkage to proteins represents a [post-translational modification](#) that has far-reaching implications in affecting protein structure and function. The use of [recombinant DNA](#) technology to identify and catalog the families of enzymes involved in effecting the modifications shows that the number of genes devoted to this activity in higher eukaryotes is very large indeed. The functions affected by these modifications are also very wide ranging. The relation between structure and function of O-glycosylated proteins is therefore complex, but important tools are being developed, and the challenge is becoming less daunting. One might predict that there will be more entries dealing with the subject in any future issues of this work.

### Bibliography

1. T. Sorensen, T. White, H. Wandall, A. K. Kristensen, P. Roepstorff and H. Clausen (1995) *J. Biol. Chem.* **270**, 24166–24173.
2. H. H. Wandall, H. Hassan, K. Mirgorodskaya, A. K. Kristensen, P. Roepstorff, E. P. Bennett, P. A. Nielsen, M. A. Hollingsworth, J. Burchell, J. Taylor-Papadimitriou, and H. Clausen (1997) *J. Biol. Chem.* **272**, 23503–23514.
3. A. P. Elhammer, R. A. Poorman, F. Brown, L. L. Maggiora, J. C. Hoogerheide, and F. I. Kezdy (1993) *J. Biol. Chem.* **268**, 10029–10038.
4. T. Hennet, F. K. Hagen, L. A. Tabak, and J. D. Marth (1995) *Proc. Natl. Acad. Sci. USA* **92**, 12070–12074.
5. S. Röttger, J. White, H. H. Wandall, J.-C. Olivo, A. Stark, E. P. Bennett, C. Whitehouse, E. C. Berger, H. Clausen, and T. Nilsson (1998) *J. Cell Sci.* **111**, 45–60.
6. S. A. Tooze, J. Tooze, and G. Warren (1988) *J. Cell Biol.* **106**, 1475–1487.
7. E. G. Berger, T. Mandel, and U. Schilt (1981) *J. Histochem. Cytochem.* **29**, 364–370.
8. M. Gentzsch and W. Tanner (1996) *EMBO J.* **15**, 5752–5759.
9. M. F. A. Bierhuizen, K. Maemura, S. Kudo, and M. Fukuda (1995) *Glycobiology* **5**, 417–425.
10. S. Tsuji (1996) *J. Biochem.* **120**, 1–13.
11. C. Rabouille, N. Hui, F. Hunte, R. Kieckbusch, E. G. Berger, G. Warren, and T. Nilsson (1995) *J. Cell Sci.* **108**, 1617–1627.
12. C. Whitehouse, J. Burchell, S. Gschmeissner, I. Brockhausen, K. O. Lloyd, and J. Taylor-Papadimitriou (1997) *J. Cell Biol.* **137**, 1229–1241.
13. M. S. Bretscher and S. Munro (1993) *Science* **261**, 1280–1281.
14. T. Nilsson, P. Slusarewicz, M. H. Hoe, and G. Warren (1993) *FEBS Lett* **330**, 1–4.
15. Y. Shimizu and S. Shaw (1993) *Nature* **366**, 630–634.
16. M. Sekine, K. Nara, and A. Suzuki (1997) *J. Biol. Chem.* **272**, 27246–27252.
17. F. Piller, V. Piller, R. I. Fox, and M. Fukuda (1988) *J. Biol. Chem.* **263**, 15146–15150.
18. I. Brockhausen, J. Yang, J. Burchell, C. Whitehouse, and J. Taylor-Papadimitriou (1995) *Eur. J. Biochem.* **233**, 607–617.
19. N. Nuwayhid, J. H. Glaser, J. C. Johnson, H. E. Conrad, S. C. Hauser, and C. B. Hirschberg

- (1986) *J. Biol. Chem.* **261**, 12936–12941.
20. T. Brinkmann, C. Weilke, and K. Kleesiek (1997) *J. Biol. Chem.* **272**, 11171–11175.
21. C. R. Torres and G. W. Hart (1984) *J. Biol. Chem.* **259**, 3308–3317.
22. M. Machida and Y. Jigami (1994) *Biosci. Biotechnol. Biochem.* **58**, 344–348.
23. T. Y. Chou, G. W. Hart, and C. V. Dang (1995) *J. Biol. Chem.* **270**, 18961–18965.
24. L. S. Griffith and B. Schmitz (1995) *Biochem. Biophys. Res. Commun.* **213**, 424–431.

### Suggestions for Further Reading

25. H. Clausen and E. P. Bennett (1996) A family of UDP-GalNAc: polypeptide N-acetylgalactosaminyl-transferases control the initiation of mucin-type O-linked glycosylation. *Glycobiology* **6**, 635–646.
26. S. J. Gendler and A. P. Spicer (1995) Epithelial mucin genes. *Ann. Rev. Physiol.* **75**, 607–634.
27. G. W. Hart, L. K. Kreppel, F. I. Comer, C. S. Arnold, D. M. Snow, Z. Ye, X. Cheng, D. DellaManna, D. S. Caine, B. J. Earles, Y. Akimoto, R. N. Cole, and B. K. Hayes (1996) O-GlcNAcylation of key nuclear and cytoskeletal proteins: reciprocity with O-phosphorylation and putative roles in protein multimerization. *Glycobiology* **6**, 711–716.
28. R. V. Iozzo and A. D. Murdoch (1996) Proteoglycans of the extracellular environment: clues from the gene and protein side offer novel perspectives in molecular diversity and function. *FASEB J.* **10**, 598–614.
29. M. Fukuda and O. Hindsgaul (1994) *Molecular Glycobiology*, IRL Press, Oxford.
30. S. Tsuji, A. K. Datta, and J. C. Pauslon (1996) Systematic nomenclature for sialyl transferases. *Glycobiology* **6**, v–vii.

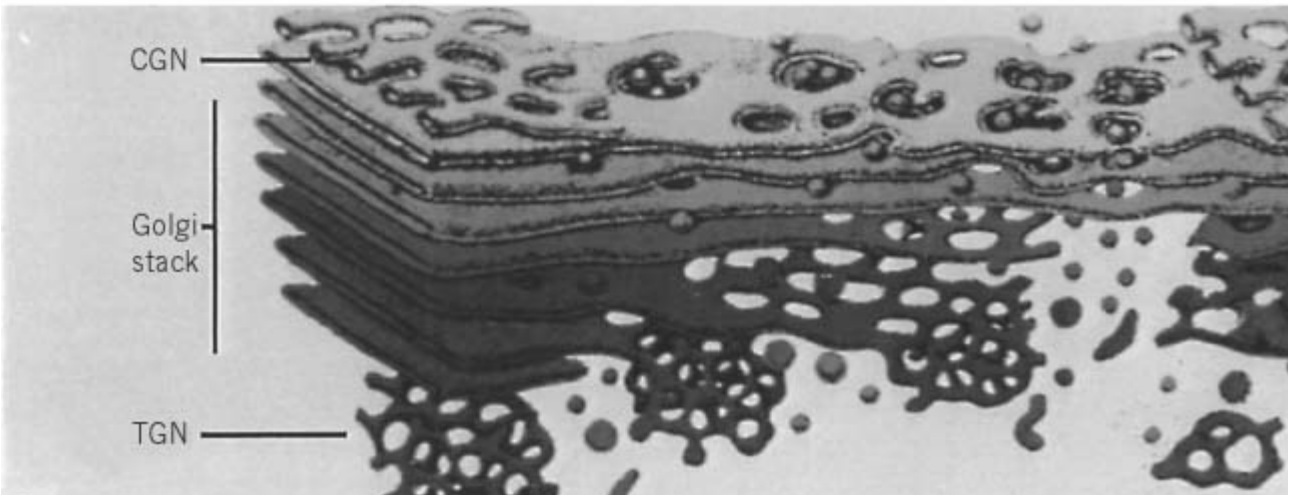
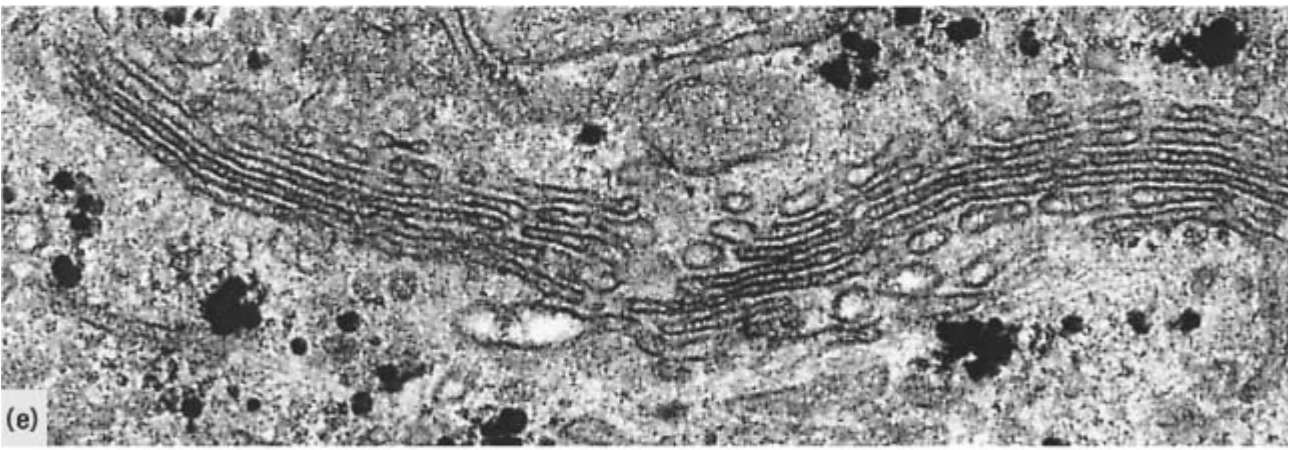
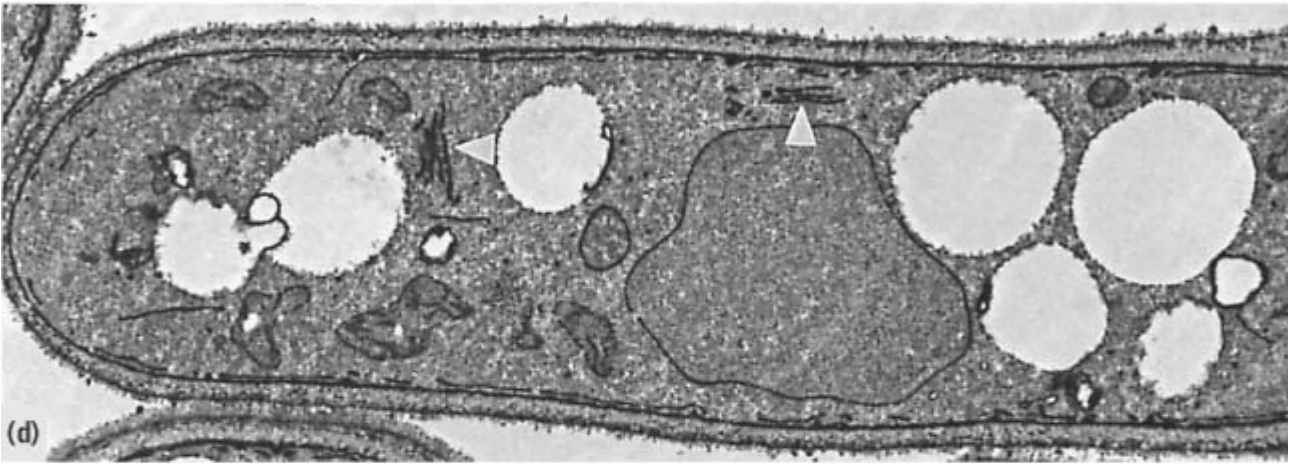
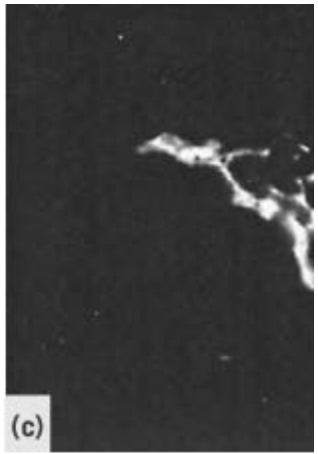
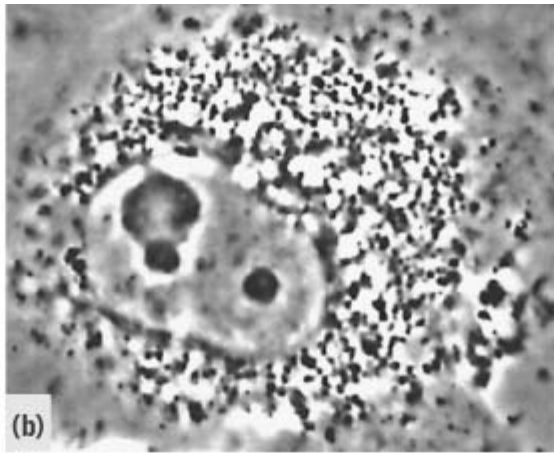
## Golgi Apparatus

The Golgi apparatus receives a varied mixture of newly synthesized [proteins](#) and [lipids](#) from the [endoplasmic reticulum](#) (ER) and sends them to their correct destination in the cell. These molecules share a common pathway through the Golgi apparatus, passing, in sequence, through an ordered array of compartments, each capable of carrying out a specific set of covalent [post-translational modifications](#). In the last compartment, they are separated from each other in preparation for final delivery.

### 1. Morphology

First described in 1898 by the Italian anatomist Camillo Golgi (1843–1926) using a silver impregnation technique (Fig. [1a](#)) ([1](#)), the Golgi apparatus has the same basic structure in all **eukaryotic** cells, comprising a stack of closely apposed and flattened cisternae (Figs. [1d](#) and [1e](#)). Cisternae are often cup-shaped, especially toward the *trans* side ([2](#), [3](#)), and there are typically three to six cisternae in the stack, although up to 40 have been reported ([4](#)). Each cisterna is about 1  $\mu\text{m}$  in diameter and comprises a core region, involved in stacking and containing resident [enzymes](#), and a fenestrated rim from which transport vesicles bud and with which they fuse ([5](#)). In **plants** and **fungi**, there are multiple copies of the Golgi dispersed throughout the cytoplasm (Fig. [1d](#)) ([3](#)). In animal cells these stacks are linked laterally, forming a bifurcating, ribbon-like structure (Figs. [1e](#) and [1f](#)) ([6](#), [7](#)), which forms a compact reticulum most often found in the pericentriolar region of the cell, adjacent to the cell [nucleus](#) (Fig. [1a](#) and [1c](#)) ([8](#)).

**Figure 1.** Different views of the Golgi apparatus. **(a)** The original “appareil réticulaire” in a Purkinje cell, described by Camille Golgi in 1898 (1) and revealed by his silver impregnation technique. **(b, c)** The Golgi reticulum in a living HeLa cell revealed by Golgi enzyme with the green fluorescent protein (138). **(b)** Phase image. **(c)** Fluorescence image. **(d)** Longitudinal section of a yeast cell (*Schizosaccharomyces pombe*) fixed with permanganate. Note the dispersed Golgi stacks (arrowheads). **(e)** Epifluorescence image of the Golgi region of a HeLa cell showing two stacks in the Golgi ribbon. **(f)** Three-dimensional reconstruction of part of the Golgi apparatus in animal cells. Sections of two Golgi stacks (left and right) are shown linked by tubules and networks connecting equivalent adjacent stacks. (Adapted from Ref. 6, with permission.) Magnification bars: **(b, c)** 10 μm, **(d)** 1 μm, **(e)** 0.2 μm.



Each face of the stack is apposed to a complex and extensive tubular network, best characterized in animal cells (and shown schematically in Fig. 1f). At the *cis* or entry face is the *cis*-Golgi network (CGN) (9), and at the *trans* or exit face is the *trans*-Golgi network (TGN) (10, 11). An intermediate compartment is interposed between the ER and the CGN, which is thought by many to ferry newly synthesized cargo molecules from the exit sites on the ER (the transitional element region) (12) to the CGN (13).

The Golgi apparatus is embedded in a matrix called the “zone of exclusion” because other cytoplasmic structures down to the size of [ribosomes](#) are excluded (14). Candidate components of this matrix have been identified, many of which have sequences predicting rod-like, [fibrous proteins](#) (15).

## 2. Biogenesis

The growth and division of the Golgi apparatus has been little studied, with the exception of the partitioning of Golgi membranes during mitosis in animal cells (16). The Golgi ribbon is converted during the early stages of **mitosis** into clusters of small **vesicles** and tubules, as well as free vesicles, a process that is thought to aid the partitioning process (17, 18). Cell-free assays (19, 20) have provided a molecular explanation of at least part of this process (21, 22).

## 3. Compartmentation

The Golgi apparatus comprises an ordered array of compartments through which newly synthesized proteins pass in sequence. The CGN is the entry point and is thought to be the last quality-control step on the pathway (23, 24). Properly folded cargo proteins proceed on to the stack. Misfolded proteins are returned to the ER for further rounds of folding (25) or degradation (26).

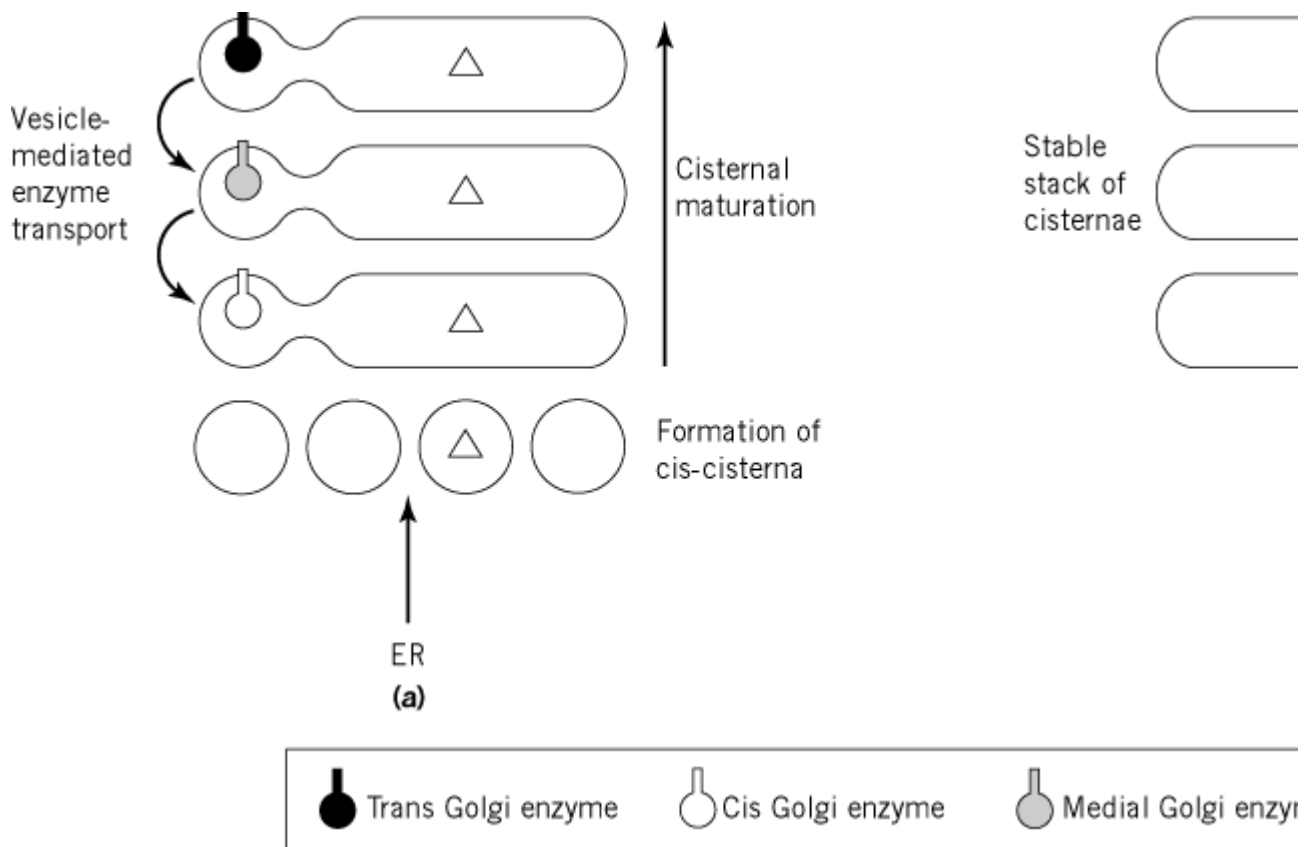
A large number of covalent modifications are carried out by the Golgi apparatus, including glycosylation, acylation, **phosphorylation**, [sulfation](#), and **proteolytic** cleavage (27-30). **N-Glycosylation** is the best characterized and comprises a sequence of processing reactions that mark the passage of proteins through the Golgi stack (27, 31). The trimming of the high-mannose oligosaccharides, which are attached co-translationally in the ER, is completed in the early part of the Golgi apparatus. Construction of complex oligosaccharides is initiated in the middle part of the Golgi stack and completed in the late part. The enzymes that carry out these steps are present in two or more adjacent cisternae, generating an overlapping distribution that might help to ensure more complete glycosylation (32). Construction of complex oligosaccharides is completed in the TGN, which also functions to sort proteins and lipids. Those destined for **lysosomes** and **secretory granules** are packaged separately from those destined for the cell surface (see text below).

**O-linked oligosaccharides** (see [O-Glycosylation](#)) and glycolipids (33) are also assembled in a stepwise fashion as they pass, in sequence, through the stacked cisternae. The precise locations at which these steps occur must await [cloning](#) and localization of the enzymes involved.

## 4. Transport Through the Golgi Apparatus

Several models have been proposed to explain the ordered transport of proteins (and some lipids) through the Golgi apparatus. The two extremes are *cisternal maturation* (34) and vesicle-mediated transport (12), illustrated in Figure 2.

**Figure 2.** Models for transport through the Golgi stack. **(a)** Cisternal maturation envisages a dynamic stack of cisternae in which cargo molecules, maturing into the next one through the retrograde transport of Golgi enzymes. **(b)** Vesicle-mediated transport envisages a stable stack of cisternae through which a mixture of cargo molecules passes, in sequence, transported from cisterna to cisterna by vesicles.



1. Cisternal maturation postulates a dynamic stack of cisternae. The *cis*-most cisterna is assembled by the fusion of transport vesicles from the ER, carrying newly synthesized cargo. Golgi processing enzymes are then delivered back from the next cisterna in the stack, which in turn receives enzymes from the next, and so on throughout the stack. In this way, each cisterna matures into the next one.

2. Vesicle-mediated transport postulates a stable stack of cisternae containing a fixed and ordered array of processing enzymes. The cargo is delivered to each cisterna in turn by **vesicles** that bud from one cisterna and fuse with the next in the stack.

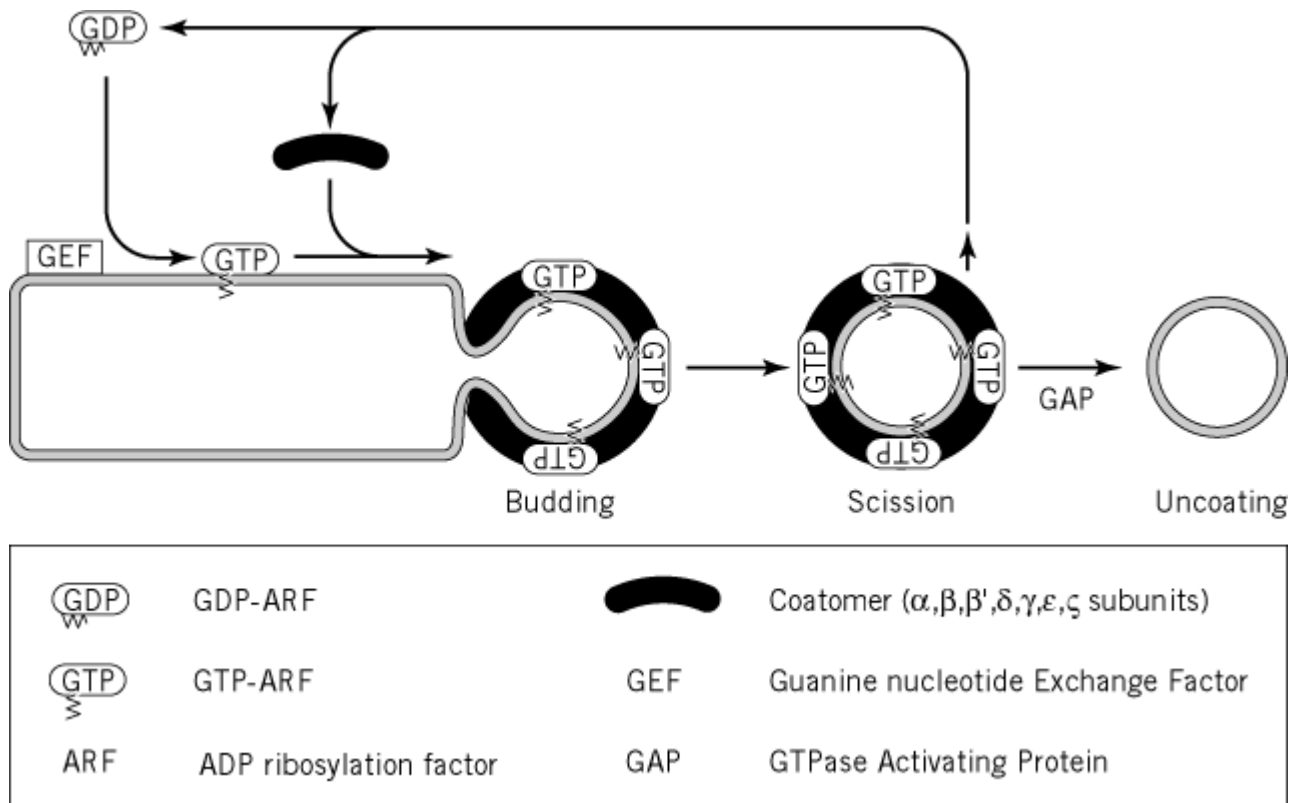
There is still controversy as to whether either model, or a hybrid one, is correct. There is, however, general agreement that the vesicles mediating these pathways are COPI (Coat protomer complex I)-coated vesicles (35).

#### 4.1. COPI-Mediated Budding

Assembly of the coat is thought to be initiated by the GTPase ARF (ADP-ribosylation factor) (36, 37) (Fig. 3). ARF exists in the cytoplasm bound to GDP, which is exchanged for GTP upon binding to Golgi membranes. This is catalyzed by a GEF (guanine nucleotide exchange factor) (38) that is sensitive to *brefeldin A* (39, 40). Exchange exposes a myristoyl group at the *N*-terminus of ARF, which aids Golgi binding (41). ARF then recruits coatamer, a complex of seven polypeptides (42, 43). This is either direct, through protein-protein interactions (44), or indirect, through the activation of phospholipase D by ARF, which generates phosphatidic acid to which the coatamer can bind (45). Stepwise binding of coatamer complexes is thought to cause incremental deformation of the membrane, generating a coated bud like that thought to occur for **clathrin**-coated vesicles (46). Scission generates a completed COPI vesicle and requires acyl-CoA, but it is otherwise

uncharacterized (47). A GAP (GTPase activating protein) catalyzes the hydrolysis of the GTP bound to ARF (48), leading to disassembly of the coat (49) and recycling of the complexes for further rounds of budding.

**Figure 3.** COPI-mediated budding of transport vesicles. Stepwise assembly of coat subunits (ARF and coatomer) generate a COPI vesicle. GTP hydrolysis releases the coat subunits for further rounds of budding, leaving an uncoated vesicle that can then fuse with the target membrane. Note that proteins such as v-SNAREs, which have to be incorporated into vesicles during budding, have been omitted for clarity.



COPI vesicles are also involved in retrograde transport from the Golgi to the ER (see text below) and perhaps also from the ER to the Golgi apparatus (50). COPI coats are found on the TGN but appear not to be involved in transport from this compartment.

There are other types of coated vesicle involved in transport from the ER to the *cis*-Golgi and from the *trans*-Golgi. Two have been characterized at the molecular level. At the *cis*-side of the Golgi apparatus, COPII coats mediate the budding of cargo-carrying vesicles from the ER (50, 51). On the *trans* side, clathrin coats package lysosomal and other proteins during transport to lysosomes or upon recycling from immature granules (see text below). Although the protein constituents of both these coats differ from COPI coats, many of the underlying principles of operation appear to be the same (52). Lipids other than phosphatidic acid, and lipid exchange proteins (53), have also been implicated in budding processes. **Diacylglycerol**, in particular, appears to play a central role (54), although it is unclear whether this role is structural or regulatory.

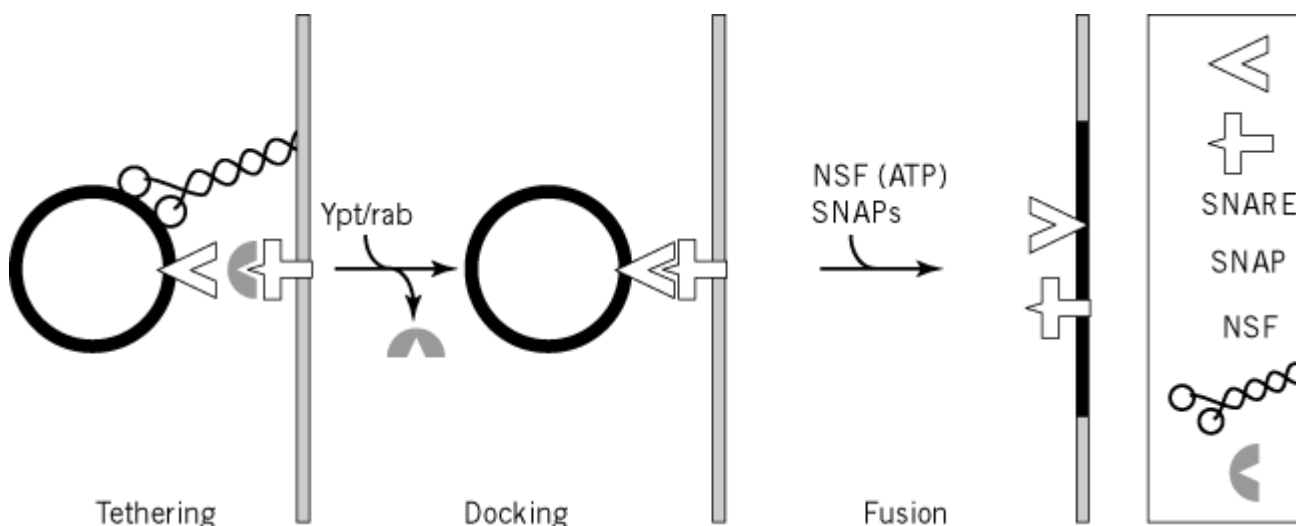
#### 4.2. Vesicle Docking and Fusion

After budding, Golgi transport vesicles must dock and fuse with the correct compartment. The SNARE hypothesis was put forward to explain the specificity of these transport steps (55) (Fig. 4). This hypothesis postulates a v- (or vesicle) SNARE (SNAP receptor) on each type of transport vesicle that interacts specifically with the t- (or target) SNARE in the recipient membrane



compartment. SNAREs were first identified at the presynaptic membrane with which synaptic vesicles fuse during neurotransmission (56-60). Using genetic and biochemical techniques, they have now been identified at many other vesicle-mediated steps in the cell (61).

**Figure 4.** Targeting and fusion of transport vesicles. Transport vesicles are tethered and then docked specifically using v process. In one model (illustrated), SNAPs mediate the binding of the NSF ATPase to the v-t SNARE pair, and hydrolyze another model, SNAPs and NSF prime the SNAREs before they can interact with each other (139).



All but one of the SNAREs studied to date are type II [membrane proteins](#), with most of their mass protruding into the cytoplasm (62); the exception has a Cys–Ala–Ala–X box at the C-terminus, so it is probably anchored by a long, hydrophobic prenyl group (63) (see [Prenylation](#)). Their sequences predict a [coiled-coil](#) structure suitable for forming paired complexes (64). Many, perhaps all, t-SNAREs are hetero-oligomers containing, in addition, a member of the SNAP25/sec9 family of proteins (65, 66). Transport vesicles can contain more than one type of v-SNARE, suggesting that multiple SNARE pairs might be used to improve the fidelity of vesicle targeting (63). Assembly of the SNARE pair is a highly regulated process. In some cases, the two membranes are initially brought together by the vesicle docking protein, p115, a process that might facilitate the efficiency with which cognate SNAREs find each other (67).

Some, perhaps all, t-SNAREs exist in an inactive form, complexed to a member of the sec1 family of proteins (68, 69), which is removed during the activation process (70). Activation of SNAREs is triggered by a member of the Ypt/rab family of small GTPases (70, 71), one or more of which is present at each vesicle-mediated transport step (72). They exist in two forms: (a) an inactive, GDP form, complexed with GDI (GDP dissociation inhibitor) in the cytoplasm and (b) an active, GTP form bound to the membrane by a prenyl group (73). Delivery of the GDP form to the correct membrane is mediated by a GDF (GDI displacement factor) (74), followed by exchange of GDP for GTP catalyzed by a GEF (75, 76). Hydrolysis of GTP by a GAP (77, 78) completes the cycle, permitting the ypt/rab to be extracted by GDI and recycled (79-81).

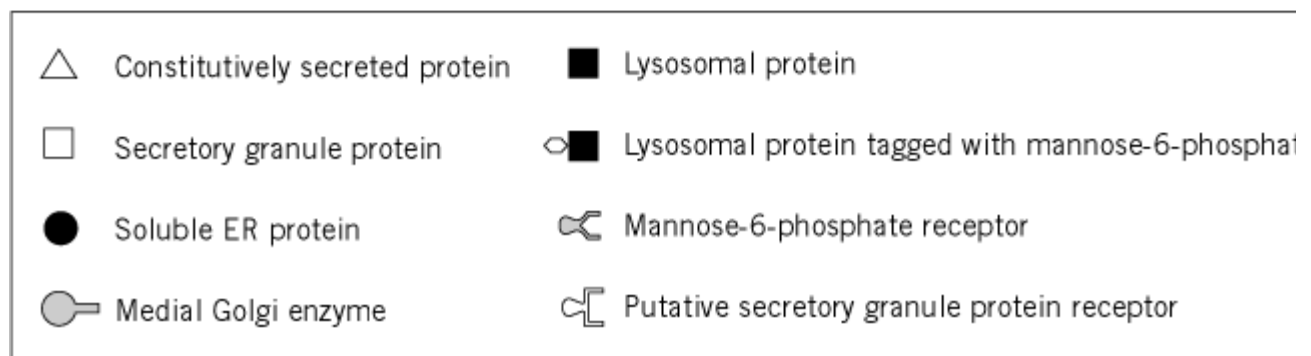
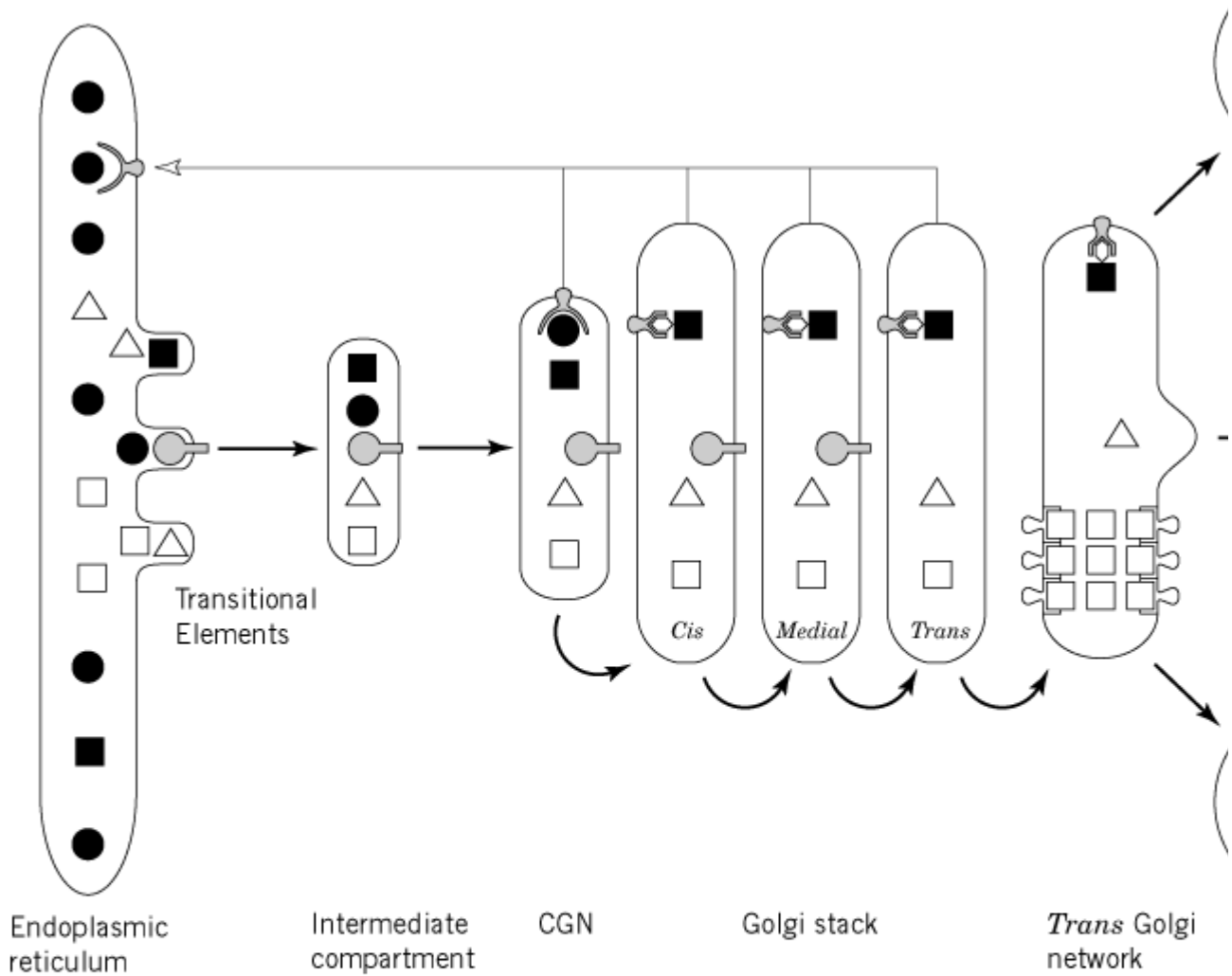
Formation of the cognate SNARE pair is thought to initiate the fusion process catalyzed by the ATPase NSF (*N*-ethylmaleimide-sensitive factor) (82). Binding of NSF to the SNARE pair is mediated by SNAPs (soluble NSF attachment proteins) (83) and leads to ATP hydrolysis, membrane fusion, and break-up of the SNARE pair (84, 85). The precise order of events is still disputed, especially the step at which NSF acts (86, 87). It is also unclear how these events lead to the physical merging of the two lipid bilayers.

More recently, the NSF-like ATPase p97 ([88](#)) has been implicated (together with NSF) in the rebuilding of the Golgi apparatus from fragments generated either by mitotic conditions ([89](#)) or by drug treatment ([90](#)). The fusion role played by p97 is distinct from that played by NSF, but the molecular mechanism is unknown.

## 5. Protein Sorting

Proteins with very different final locations within the cell are assembled in the ER and sorted in the Golgi apparatus (Fig. [5](#)). Cargo moves to the TGN, where sorting mechanisms separate lysosomal and **secretory granule** proteins from each other and from proteins destined for the cell surface. ER and Golgi proteins have retention mechanisms that inhibit movement beyond the point in the [exocytosis](#) pathway at which they usually act. Those molecules that do escape are salvaged by retrieval mechanisms. Retrieval mechanisms also recycle membrane components of the transport machinery. All these mechanisms interpret signals on the proteins undergoing transport. These signals are most often short stretches of sequence that are either read directly by the sorting mechanism or are used to construct the sorting signal.

**Figure 5.** Protein sorting during transport through the Golgi apparatus. A varied mixture of proteins is assembled in the ER. Some proteins are salvaged, mostly, but not exclusively, from the early part of the Golgi and are returned to the ER. Golgi enzymes and cargo proteins move to the TGN, where they are separately packaged and sent to the plasma membrane, late [endosomes](#)



### 5.1. The Default Pathway

The pathway from the ER to the cell surface was originally designated the default pathway because signals did not appear to be required to transport cargo through the correct sequence of Golgi compartments to the cell surface (91). This is still true for many proteins and many compartmental steps, but there is increasing evidence that some surface cargo can be selectively concentrated, especially during export from the ER (92, 93).

Two types of protein that could be involved in concentrating cargo have been described. The first are **lectin**-like proteins—for example, ERGIC53 (94) and Vip36 (95), which are thought to bind the oligosaccharides attached to many cargo proteins. In the case of ERGIC53, the lectin–cargo complex

would be concentrated in transport vesicles because the cytoplasmic tail binds to coat proteins (96). The second are the p24 family of spanning proteins (97, 98), the large luminal domains of which are thought to bind to cargo molecules by unknown means (97), aided by ancillary proteins (99). This family of proteins is selectively incorporated into transport vesicles, because their cytoplasmic tails contain a dibasic motif that binds coat subunits (100, 101).

## 5.2. Sorting in the TGN

All cargo molecules share the same pathway up to and including the TGN. It is here that lysosomal (102) and secretory granule proteins (103) are separated from surface cargo. Many soluble lysosomal enzymes are tagged in the early part of the Golgi apparatus with a mannose-6-phosphate signal that is constructed from some of the terminal mannose residues on the high-mannose oligosaccharide by two enzymes that recognize protein features unique to lysosomal enzymes (31, 104). The tagged enzyme moves through the Golgi stack to the TGN, where it binds the large **mannose-6-phosphate receptor** (105). The complex is specifically incorporated into clathrin-coated vesicles, because multiple signals in the cytoplasmic tail of the receptor (106) are recognized by an adaptor (107, 108). The complex is then ferried to a prelysosomal compartment (109) where, eventually, the low pH separates the enzyme from the receptor. The enzyme is delivered to lysosomes, whereas the receptor is recycled for further rounds of transport. Some lysosomal membrane proteins are also ferried in clathrin-coated vesicles. The same adaptor recognizes a signal in the cytoplasmic tail (110).

As well as a direct route to lysosomes from the TGN, there is also an indirect route via the cell surface. For soluble enzymes, this is mediated by the small mannose-6-phosphate receptor (106).

Secretory granule proteins begin to condense in the TGN, probably because of the lower pH of this part of the Golgi apparatus and the presence of calcium ions (111). Condensation is a sorting mechanism, because during this process all other proteins are largely excluded, including, in some cells, the granule proteins destined for other types of secretory granule. The outside of the condensing granule core is thought to bind to sorting **receptors**, thereby enveloping the core with a membrane that buds to form an immature **vacuole**. Maturation proceeds by a process involving vacuole fusion and retrieval of excess membrane, until a mature secretory granule is formed (112).

In polarized cells, there is more than one surface domain, and signals exist to direct surface cargo to one or other of these domains. The signals used include the bound *N*-linked oligosaccharides (113) or sequences within the cytoplasmic tail (114, 115). **GPI-anchored** proteins are thought to form “rafts” that are selectively incorporated into vesicles bound for the apical membrane domain (116).

## 5.3. Retention in the Golgi Apparatus

Proteins destined to remain in the Golgi apparatus are also synthesized and assembled in the ER. They are transported together with other cargo to the Golgi apparatus, but must then stop in the correct cisterna(e). The retention signal that specifies location within the Golgi apparatus is contained within the membrane-spanning domain and flanking sequences (117) and is thought to inhibit further forward transport in one of two ways. The kin recognition model (118) postulates interaction with other Golgi proteins sharing the same cisterna, generating an oligomer too large to enter the forward-moving transport vesicles. The bilayer thickness model (119) exploits the increasing thickness of Golgi membrane across the stack caused by a gradient of cholesterol. The Golgi protein moves forward until the length of the spanning domain matches that of the bilayer. Different lengths would specify different locations within the stack. There is evidence both for and against each model, but they are not mutually exclusive, and a combination of the two could serve to locate the proteins more precisely in the stack (120).

## 5.4. Retrieval by Retrograde Transport

Two types of retrieval operate from the Golgi apparatus to the ER. The first is a salvage process, recovering proteins that stray beyond the point in the exocytic pathway at which they operate (121). The need for salvage gave rise to the distillation hypothesis, which could help explain the need for a Golgi stack (122). The second type of retrieval is recycling, mostly of components of the vesicle

fusion machinery.

Salvage of ER proteins is mediated by short signals at the *N*- or *C*-terminus. Soluble proteins in the ER lumen use a *C*-terminal tetrapeptide (123), [KDEL in mammals, HDEL in budding yeast], whereas membrane proteins use either a double-lysine (type I proteins) (124) or a double-arginine (type II proteins) (125) motif (but see Ref. 126). Escaped luminal proteins bind to the KDEL receptor (127), mostly located in the CGN (128), and this triggers return of the complex (129) in COPI vesicles (130) to the ER, where a change in pH is thought to effect release (131). Escaped membrane proteins with a double-lysine motif are also returned by COPI vesicles, because this motif binds directly to the coatamer component of the COPI coats (132, 133). The mechanism for returning ER proteins with a double-arginine motif is not known.

Salvage of Golgi proteins is less well understood, although there is now good evidence that it occurs (134, 135). Salvage is also an important part of secretory granule biogenesis. During maturation of immature granules, clathrin-coated vesicles recover and return proteins destined for other pathways (136).

Recycling is exemplified by the v-SNAREs. These are present in the target compartment after vesicle fusion and must be recycled to the preceding, donor compartment for further rounds of vesicle budding. The mechanism is not known, but some evidence suggests that they are not passively recycled. Instead, they direct the vesicle to the preceding compartment, in the same way that they directed the earlier vesicles to the target compartment. A t-SNARE would be needed in each of the two compartments linked by the vesicles, but only one v-SNARE. This parsimonious solution to recycling requires a mechanism to switch the specificity of the v-SNARE, depending on which of the two compartments is being targeted (137).

## Bibliography

1. C. Golgi (1898) *Arch. Ital. Biol.* **30**, 60–71.
2. W. G. Whaley (1975) *The Golgi Apparatus*. Cell Biology Monographs, Continuation of *Protoplasmatologia*, **2**, Springer-Verlag, New York.
3. M. G. Farquhar and G. E. Palade (1981) *J. Cell Biol.* **91**, 77s–103s.
4. D. W. Fawcett (1981) *The Cell*, 2nd ed., W. B. Saunders, Philadelphia.
5. P. Weidman, R. Roth, and J. Heuser (1993) *Cell* **75**, 123–133.
6. A. Rambourg and Y. Clermont (1990) *Eur. J. Cell. Biol.* **51**, 189–200.
7. J. M. Lucocq and G. Warren (1987) *EMBO J.* **6**, 3239–3246.
8. W. C. Ho, V. J. Allan, M. G. van, E. G. Berger, and T. E. Kreis (1989) *Eur. J. Cell. Biol.* **48**, 250–263.
9. H. P. Hauri and A. Schweizer (1992) *Curr. Opin. Cell. Biol.* **4**, 600–608.
10. G. Griffiths and K. Simons (1986) *Science* **234**, 438–443.
11. M. S. Ladinsky, J. R. Kremer, P. S. Furcinitti, J. R. McIntosh, and K. E. Howell (1994) *J. Cell Biol.* **127**, 29–38.
12. G. Palade (1975) *Science* **189**, 347–358.
13. J. Saraste and K. Svensson, (1991) *J. Cell. Sci.* **100**, 415–430.
14. D. J. Morré and L. Ovtracht (1977) *Int. Rev. Cytol.* **5**, 61–188.
15. F. A. Barr and G. Warren, (1996) *Semin. Cell. Dev. Biol.* **7**, 505–510.
16. G. Warren (1993) *Annu. Rev. Biochem.* **62**, 323–348.
17. G. Warren, T. Levine, and T. Misteli, (1995) *Trends Cell Biol.* **5**, 413–416.
18. G. Warren and W. Wickner (1996) *Cell* **84**, 395–400.
19. C. Rabouille, T. Misteli, R. Watson, and G. Warren (1995) *J. Cell Biol.* **129**, 605–618.
20. T. Misteli and G. Warren (1995) *J. Cell Biol.* **130**, 1027–1039.

21. T. P. Levine, C. Rabouille, R. H. Kieckbusch, and G. Warren (1996) *J. Biol. Chem.* **271**, 17304–17311.
22. N. Nakamura, M. Lowe, T. P. Levine, C. Rabouille, and G. Warren (1997) *Cell* **89**, 445–455.
23. S. M. Hurlley and A. Helenius (1989) *Annu. Rev. Cell. Biol.* **5**, 277–307.
24. C. Hammond and A. Helenius (1995) *Curr. Opin. Cell. Biol.* **7**, 523–529.
25. D. N. Hebert, B. Foellmer, and A. Helenius (1995) *Cell* **81**, 425–433.
26. E. J. Wiertz et al. (1996) *Nature* **384**, 432–438.
27. J. Roth (1987) *Biochim. Biophys. Acta* **906**, 405–436.
28. C. Niehrs, R. Beisswanger, and W. B. Huttner (1994) *Chem. Biol. Interact.* **92**, 257–271.
29. M. F. Schmidt (1989) *Biochim. Biophys. Acta.* **988**, 411–426.
30. D. F. Steiner, S. P. Smeekens, S. Ohagi, and S. J. Chan (1992) *J. Biol. Chem.* **267**, 23435–23438.
31. R. Kornfeld and S. Kornfeld (1985) *Annu. Rev. Biochem.* **54**, 631–664.
32. C. Rabouille et al. (1995) *J. Cell. Sci.* **108**, 1617–1627.
33. K. Sandhoff and G. van Echten (1994) *Prog. Brain. Res.* **101**, 17–29.
34. B. Becker, B. Bolinger, and M. Melkonian (1995) *Trends Cell Biol.* **5**, 305–307.
35. V. Malhotra, T. Serafini, L. Orci, J. C. Shepherd, and J. E. Rothman (1989) *Cell* **58**, 329–336.
36. T. Serafini, L. Orci, M. Amherdt, M. Brunner, R. A. Kahn, and J. E. Rothman (1991) *Cell* **67**, 239–253.
37. J. Moss and M. Vaughan (1995) *J. Biol. Chem.* **270**, 12327–12330.
38. A. Peyroche, S. Paris, and C. L. Jackson (1996) *Nature* **384**, 479–481.
39. J. B. Helms and J. E. Rothman (1992) *Nature* **360**, 352–354.
40. J. G. Donaldson, D. Finazzi, and R. D. Klausner (1992) *Nature* **360**, 350–352.
41. P. A. Randazzo, T. Terui, S. Sturch, H. M. Fales, A. G. Ferrige, and R. A. Kahn (1995) *J. Biol. Chem.* **270**, 14809–4815.
42. M. G. Waters, T. Serafini, and J. E. Rothman, (1991) *Nature* **349**, 248–251.
43. T. E. Kreis, M. Lowe, and R. Pepperkok (1995) *Annu. Rev. Cell Dev. Biol.* **11**, 677–706.
44. L. Y. Zhao et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 4418–4423.
45. N. T. Ktistakis, H. A. Brown, M. G. Waters, P. C. Sternweis, and M. G. Roth (1996) *J. Cell. Biol.* **134**, 295–306.
46. J. Heuser and L. Evans, (1980) *J. Cell. Biol.* **84**, 560–583.
47. J. Ostermann et al. (1993) *Cell* **75**, 1015–1025.
48. E. Cukierman, I. Huber, M. Rotman and D. Cassel (1995) *Science* **270**, 1999–2002.
49. G. Tanigawa, L. Orci, M. Amherdt, M. Ravazzola, J. B. Helms, and J. E. Rothman (1993) *J. Cell. Biol.* **123**, 1365–1371.
50. S. Y. Bednarek, L. Orci, and R. Schekman (1996) *Trends Cell Biol.* **6**, 468–473.
51. C. Barlowe et al. (1994) *Cell* **77**, 895–907.
52. R. Schekman and L. Orci (1996) *Science* **271**, 1526–1533.
53. V. A. Bankaitis, J. R. Aitken, A. E. Cleves, and W. Dowhan (1990) *Nature* **347**, 561–562.
54. B. G. Kearns et al. (1997) *Nature* **387**, 101–105.
55. T. Sollner et al. (1993) *Nature* **362**, 318–324.
56. W. S. Trimble, D. M. Cowan, and R. H. Scheller (1988) *Proc. Natl. Acad. Sci. USA* **85**, 4538–4542.
57. M. Baumert, P. R. Maycox, F. Navone, P. De Camilli, and R. Jahn (1989) *EMBO J.* **8**, 379–384.
58. M. K. Bennett, N. Calakos, and R. H. Scheller (1992) *Science* **257**, 255–259.

59. Y. Goda (1997) *Proc. Natl. Acad. Sci. USA* **94**, 769–772.
60. G. Schiavo, O. Rossetto, and C. Montecucco (1994) *Semin. Cell Biol.* **5**, 221–229.
61. S. Ferro-Novick and R. Jahn, (1994) *Nature* **370**, 191–193.
62. J. B. Bock and R. H. Scheller (1997) *Nature* **387**, 133–135.
63. M. Sogaard et al. (1994) *Cell* **78**, 937–948.
64. A. Lupas (1997) *Trends Biochem. Sci.* **21**, 375–382.
65. G. A. Oyler et al. (1989) *J. Cell Biol.* **109**, 3039–3052.
66. P. Brennwald, B. Kearns, K. Champion, S. Keränen, V. Bankaitis, and P. Novick (1994) *Cell* **79**, 245–258.
67. S. K. Sapperstein, V. V. Lupashin, H. D. Schmitt, and M. G. Waters (1996) *J. Cell Biol.* **132**, 755–767.
68. Y. Hata, C. A. Slaughter, and T. C. Sudhof (1993) *Nature* **366**, 347–351.
69. M. K. Aalto, L. Ruohonen, K. Hosono, and S. Keranen (1991) *Yeast* **7**, 643–650.
70. V. V. Lupashin and M. G. Waters (1997) *Science* **276**, 1255–1258.
71. J. P. Lian, S. Stone, Y. Jiang, P. Lyons, and S. Ferro-Novick (1994) *Nature* **372**, 698–701.
72. K. Simons and M. Zerial (1993) *Neuron* **11**, 789–799.
73. C. Nuoffer and W. E. Balch (1994) *Annu. Rev. Biochem.* **63**, 949–990.
74. A. B. Dirac Svejstrup, T. Sumizawa, and S. R. Pfeffer (1997) *EMBO J.* **16**, 465–472.
75. O. Ullrich, H. Horiuchi, C. Bucci, and M. Zerial (1994) *Nature* **368**, 157–160.
76. T. Soldati, A. D. Shapiro, A. B. Svejstrup, and S. R. Pfeffer (1994) *Nature* **369**, 76–78.
77. E. S. Burstein and I. G. Macara (1992) *Proc. Natl. Acad. Sci. USA* **89**, 1154–1158.
78. T. J. Tan, P. Vollmer, and D. Gallwitz (1991) *FEBS Lett.* **291**, 322–326.
79. K. Fukui, T. Sasaki, K. Imazumi, Y. Matsuura, H. Nakanishi, and Y. Takai (1997) *J. Biol. Chem.* **272**, 4655–4658.
80. O. Ullrich, et al. (1993) *J. Biol. Chem.* **268**, 18143–18150.
81. S. Araki, A. Kikuchi, Y. Hata, M. Isomura, and Y. Takai (1990) *J. Biol. Chem.* **265**, 13007–13015.
82. M. R. Block and J. E. Rothman (1992) *Methods Enzymol.* **219**, 300–309.
83. D. O. Clary, I. C. Griff, and J. E. Rothman (1990) *Cell* **61**, 709–721.
84. D. W. Wilson, S. W. Whiteheart, M. Wiedmann, M. Brunner, and J. E. Rothman (1992) *J. Cell Biol.* **117**, 531–538.
85. T. Sollner, M. K. Bennett, S. W. Whiteheart, R. H. Scheller, and J. E. Rothman (1993) *Cell* **75**, 409–418.
86. A. Mayer, W. Wickner, and A. Haas (1996) *Cell* **85**, 83–94.
87. A. Morgan and R. D. Burgoyne (1995) *EMBO J.* **14**, 232–239.
88. J. M. Peters, M. J. Walsh, and W. W. Franke (1990) *EMBO J.* **9**, 1757–1767.
89. C. Rabouille, T. P. Levine, J. M. Peters, and G. Warren (1995) *Cell* **82**, 905–914.
90. U. Acharya, R. Jacobs, J. M. Peters, N. Watson, M. G. Farquhar, and V. Malhotra (1995) *Cell* **82**, 895–904.
91. F. T. Wieland, M. L. Gleason, T. A. Serafini, and J. E. Rothman (1987) *Cell* **50**, 289–300.
92. P. Quinn, G. Griffiths, and G. Warren (1984) *J Cell Biol* **98**, 2142–2147.
93. W. E. Balch, J. M. McCaffery, H. Plutner, and M. G. Farquhar (1994) *Cell* **76**, 841–852.
94. C. Arar, V. Carpentier, J. P. Lecaer, M. Monsigny, A. Legrand, and A. C. Roche (1995) *J. Biol. Chem.* **270**, 3551–3553.
95. K. Fiedler and K. Simons (1996) *J. Cell Sci.* **109**, 271–276.
96. E. J. Tisdale, H. Plutner, J. Matteson, and W. E. Balch (1997) *J. Cell Biol.* **137**, 581–593.

97. F. Schimmoller, B. Singer Kruger, S. Schroder, U. Kruger, C. Barlowe, and H. Riezman (1995) *EMBO J.* **14**, 1329–1339.
98. M. A. Stamnes et al. (1995) *Proc. Natl. Acad. Sci.* **92**, 8011–8015.
99. M. J. Kuehn, R. Schekman, and P. O. Ljungdahl (1996) *J. Cell Biol.* **135**, 585–595.
100. K. Fiedler, M. Veit, M. A. Stamnes, and J. E. Rothman (1996) *Science* **273**, 1396–1399.
101. K. Sohn et al. (1996) *J. Cell Biol.* **135**, 1239–1248.
102. S. Kornfeld and I. Mellman (1989) *Annu. Rev. Cell Biol.* **5**, 483–525.
103. T. L. Burgess and R. B. Kelly (1987) *Annu. Rev. Cell Biol.* **3**, 243–293.
104. S. R. Pfeffer and J. E. Rothman (1987) *Annu. Rev. Biochem.* **56**, 829–852.
105. R. Kornfeld (1992) *Annu. Rev. Biochem.* **61**, 307–330.
106. A. Hille Rehfeld (1995) *Biochim. Biophys. Acta* **1241**, 177–194.
107. M. S. Robinson (1992) *Trends Cell Biol.* **2**, 293–297.
108. M. S. Robinson, (1994) *Curr. Opin. Cell Biol.* **6**, 538–544.
109. G. Griffiths, B. Hoflack, K. Simons, I. Mellman, and S. Kornfeld (1988) *Cell* **52**, 329–341.
110. S. Honing and W. Hunziker (1995) *J. Cell Biol.* **128**, 321–332.
111. W. B. Huttner et al. (1995) *Cold Spring Harbor Symp. Quantit. Biol.* **LX**, 315–327.
112. S. A. Tooze (1991) *FEBS Lett* **285**, 220–224.
113. P. Scheiffele, J. Peranen, and K. Simons (1995) *Nature* **378**, 96–98.
114. K. E. Mostov and M. H. Cardone (1995) *Bioessays* **17**, 129–138.
115. K. Matter and I. Mellman (1994) *Curr. Opin. Cell Biol.* **6**, 545–554.
116. K. Simons and E. Ikonen (1997) *Nature* **387**, 569–572.
117. C. E. Machamer (1991) *Trends Cell Biol.* **1**, 141–144.
118. T. Nilsson et al. (1994) *EMBO J.* **13**, 562–574.
119. M. S. Bretscher and S. Munro (1993) *Science* **261**, 1280–1281.
120. H. R. Pelham and S. Munro (1993) *Cell* **75**, 603–605.
121. G. Warren (1987) *Nature* **327**, 17–18.
122. J. E. Rothman (1981) *Science* **213**, 1212–1219.
123. S. Munro and H. R. Pelham (1987) *Cell* **48**, 899–907.
124. M. R. Jackson, T. Nilsson, and P. A. Peterson (1990) *EMBO J.* **9**, 3153–3162.
125. M. P. Schutze, P. A. Peterson, and M. R. Jackson (1994) *EMBO J.* **13**, 1696–1705.
126. D. J. Sweet and H. R. B. Pelham (1992) *EMBO J.* **11**, 423–432.
127. J. C. Semenza, K. G. Hardwick, N. Dean, and H. R. B. Pelham (1990) *Cell* **61**, 1349–1357.
128. B. L. Tang, S. H. Wong, X. L. Qi, S. H. Low, and W. Hong (1993) *J. Cell Biol.* **120**, 325–328.
129. M. J. Lewis and H. R. Pelham (1992) *Cell* **68**, 353–364.
130. B. Sönnichsen, R. Watson, H. Clausen, T. Misteli, and G. Warren, (1996) *J. Cell Biol.* **134**, 1411–1425.
131. D. W. Wilson, M. J. Lewis, and H. R. B. Pelham (1993) *J. Biol. Chem.* **268**, 7465–7468.
132. P. Cosson and F. Letourneur (1994) *Science* **263**, 1629–1631.
133. F. Letourneur et al. (1994) *Cell* **79**, 1199–1207.
134. M. H. Hoe, P. Slusarewicz, T. Misteli, R. Watson, and G. Warren (1995) *J. Biol. Chem.* **270**, 25057–25063.
135. S. L. Harris and M. G. Waters (1996) *J. Cell Biol.* **132**, 985–998.
136. A. S. Dittie, N. Hajibagheri, and S. A. Tooze (1996) *J. Cell Biol.* **132**, 532–536.
137. M. J. Lewis and H. R. Pelham (1996) *Cell* **85**, 205–215.
138. D. T. Shima, H. Haldar, R. Pepperkok, R. Watson, and G. Warren (1997) *J. Cell Biol.* **137**,



1211–1228.

139. B. J. Nichols, C. Ungermann, H. R. B. Pelham, W. T. Wickner, and A. Haas (1997) *Nature* **387**, 199–202.

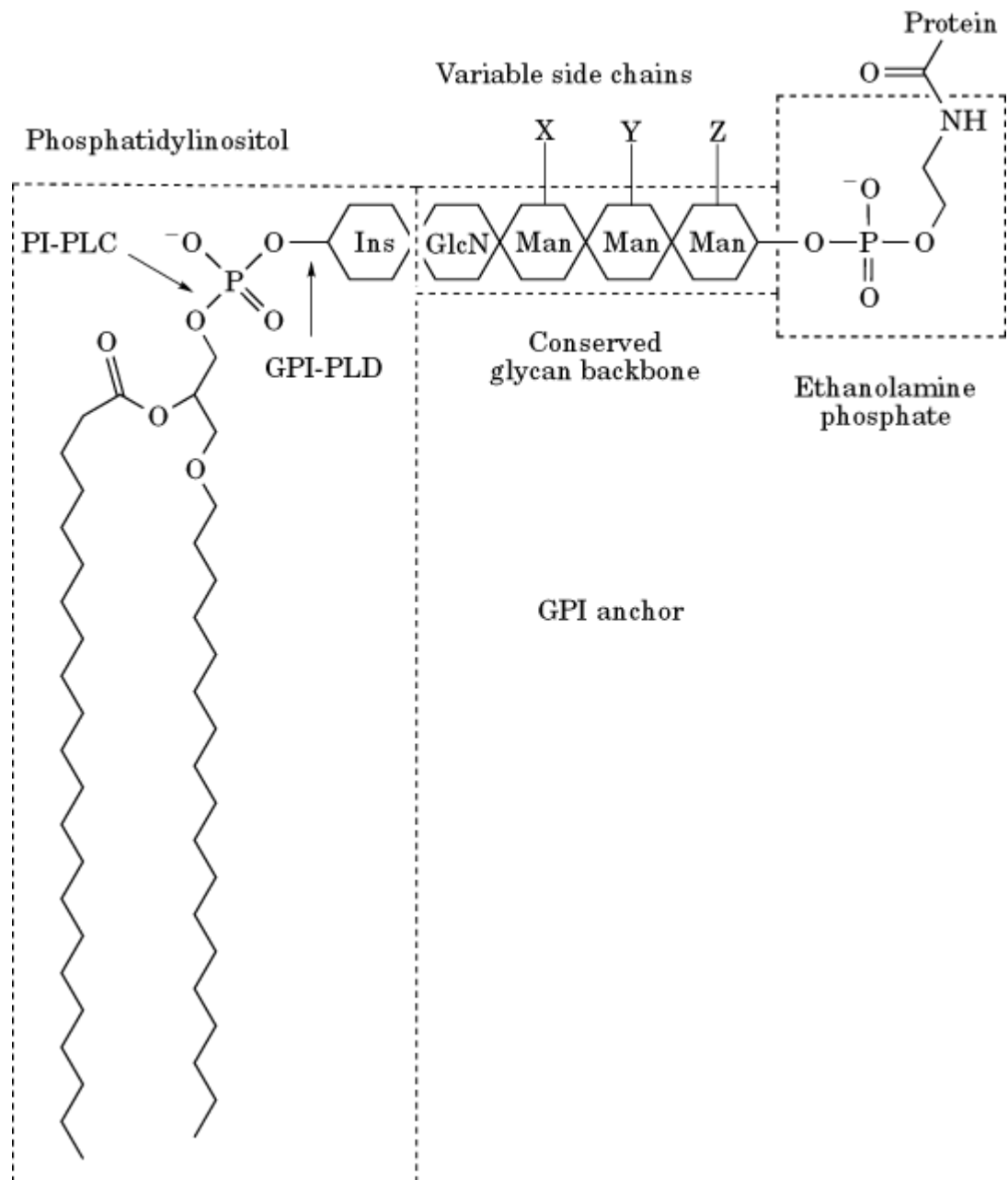
### Suggestions for Further Reading

140. H. R. B. Pelham (1989) Control of protein exit from the endoplasmic reticulum. *Annu. Rev. Cell Biol.* **5**, 1–23.
141. N. K. Pryer, L. J. Wuestehube, and R. Schekman (1992) Vesicle-mediated protein sorting. *Annu. Rev. Biochem.* **61**, 471–516.
142. J. E. Rothman and G. Warren (1994) Implications of the SNARE hypothesis for intracellular membrane topology and dynamics. *Curr. Biol.* **4**, 220–233.
143. S. R. Pfeffer (1996) Transport vesicle docking: SNAREs and associates. *Annu. Rev. Cell Dev. Biol.* **12**, 441–461.
144. J. E. Rothman and F. T. Wieland (1996) Protein sorting by transport vesicles. *Science* **272**, 227–234.

## GPI Anchor

A large number of [proteins](#) on the surfaces of eukaryotic cells are covalently linked to a glycosylphosphatidylinositol (GPI) molecule (Table [1](#)). GPI molecules consist of a glycan chain that contains ethanolamine, glucosamine, and mannose residues linked to the inositol ring of [phosphatidylinositol](#) (Fig. [1](#)). The ethanolamine in the GPI molecule is amide-linked to the carboxyl-terminal amino acid residue of the protein after it is translocated across the membrane of the endoplasmic reticulum. The GPI molecule is embedded in the lipid bilayer and is the [membrane anchor](#) responsible for retaining the protein at the cell surface. GPI-anchored proteins can be released from the membrane by highly specific phospholipases that cleave the GPI molecule. The complex structure of the GPI molecule (Fig. [1](#)), compared to the other lipids used as membrane anchors, suggests that it has additional functions.

**Figure 1.** Structure of the GPI anchor. The C-terminal residue of a protein is attached via an amide linkage to the amino group of an ethanolamine phosphate in the GPI molecule ([1](#)). A conserved glycan backbone connects the ethanolamine phosphate to the phosphatidylinositol moiety. For simplicity, the sugar residues in the glycan are shown as hexagons. The exact linkages between the sugars and with the phosphatidylinositol and ethanolamine phosphate moieties (i.e., EtNP-6Man $\alpha$ 1-2Man $\alpha$ 1-6Man $\alpha$ 1-4GlcN $\alpha$ 1-6*myo*inositol), are highly conserved between protozoa and mammals. In contrast, the side chains (shown as X, Y, and Z) that are attached to the mannose residues in the glycan backbone are highly variable and can contain a variety of other sugars and additional ethanolamine phosphate groups. The hydrophobic part of the phosphatidylinositol (which will be embedded in the lipid bilayer) is shown as a 1-alkyl, 2-acyl structure, although 1,2-diacylglycerol, 1-alkylglycerol, or ceramide are also found. The hydrocarbon chain length and degree of unsaturation are also highly variable among proteins (the example shown here is for human folate-binding protein). The inositol ring can also become acylated (not shown), which makes the GPI anchor insensitive to degradation by PI-phospholipase C (PI-PLC). The sites of action of two different phospholipases that cleave GPI anchors are also shown. PI-PLC is a phospholipase C that can also cleave phosphatidylinositol and is produced by several different pathogenic bacteria [see [Phospholipases C](#)]. Although PI-PLC is not involved in physiological release processes, it is widely used in identifying and experimentally manipulating GPI-anchored proteins on cell surfaces. GPI-PLD is a GPI-specific phospholipase D and cannot hydrolyze phosphatidylinositol. GPI-PLD is abundant in mammalian plasma and may be involved in physiologically releasing GPI-anchored proteins from cell surfaces. Ins = *myo*inositol, GlcN = glucosamine, Man = mannose, EtNP = ethanolamine phosphate.



**Table 1. Examples of GPI-Anchored Proteins**

---

Alkaline phosphatase and many other ectoenzymes  
 Variant surface protein of *Trypanosoma brucei*  
 Heparan sulfate proteoglycan  
 Scrapie prion protein  
 Thy-1  
 TAG-1 and other neural adhesion molecules  
 Carcinoembryonic antigen  
 Decay accelerating factor

Folate binding protein  
Numerous CD antigens (e.g., CD14, CD16, CD58, CD59, etc.)

---

Biosynthesis of GPI-anchored proteins takes place in two main stages: (1) The GPI molecule is synthesized by sequential addition of the glycan chain components to a phosphatidylinositol molecule at the cytoplasmic surface of the endoplasmic reticulum. Then, the GPI molecule is translocated across the membrane of the endoplasmic reticulum to the luminal surface (1-3). (2) The protein is translocated, cotranslationally, into the endoplasmic reticulum. The C-terminal region of the protein contains a special type of “GPI signal peptide” that is recognized and removed by an enzyme (believed to be a transamidase) located on the luminal surface of the endoplasmic reticulum (4). The same enzyme catalyzes coupling of the terminal ethanolamine of the GPI molecule to the newly exposed carboxyl group in an amide linkage. The attachment process is completed within a few minutes of translation/translocation, and then the GPI-anchored protein is transported through the [Golgi apparatus](#) to the cell surface by the conventional **vesicle** transport mechanisms used by other plasma [membrane proteins](#). It should be emphasized that relatively little direct biochemical information is available about the attachment process. Efficient GPI attachment may require precise spatial and temporal coordination of the transamidase with the special GPI signal peptide of the protein and translocation of the GPI precursor. Probably for that reason, reconstitution of GPI-anchoring *in vitro* from purified components presents a formidable technical challenge.

Unlike the other types of lipid anchor, the residue to which the GPI becomes attached (designated site w; see later) is quite variable and relatively difficult to identify from just the amino acid sequence of the protein. Nevertheless, rules for predicting the existence and the location of GPI anchors in a novel protein have been devised on the basis of two sources of information from about twenty naturally occurring, GPI-anchored proteins: (1) comparison of the protein sequence predicted from its gene sequence with the identity of the C-terminal residue in the mature protein, and (2) extensive mutational analysis of the C-terminal region of a limited number of model proteins (4). This region can be divided into three parts:

1. *Hydrophobic C-terminus*. A **hydrophobic** C-terminus is essential for GPI-anchoring, although the minimum number of hydrophobic residues required for GPI-anchoring is quite variable among proteins. In most known GPI-anchored proteins, the hydrophobic region is about 15 to 20 residues long, but it can be as short as eight residues. At first sight this region looks like a hydrophobic polypeptide anchor (see [Membrane Anchors](#)), and it may in fact serve this temporary function for the short interval between completion of translation/translocation and attachment of the GPI molecule. Polypeptide anchors are generally a few residues longer, however, and usually terminate in relatively hydrophilic residues. Furthermore, mutations that make this region more like a polypeptide anchor usually result in poor GPI anchoring.
2. *Hydrophilic region*: A [hydrophilic](#) region of approximately five to seven residues is located between the hydrophobic region and the attachment site. This region exhibits no strong preference for particular amino acid residues.
3. *Cleavage-GPI attachment site*: There is no strong **consensus sequence**, but there are definite amino acid preferences in the vicinity of the cleavage site: Site w (which becomes the C-terminal residue in the mature protein): Gly, Ala, Ser, Cys, Asp, and Asn preferred. Site w + 1: no preference. Site w + 2: Gly, Ala and Ser preferred. A statistical application of the w/w + 2 preferences predict the attachment site in naturally occurring proteins with approximately 80% accuracy. Apart from its predictive value, this information can be used for engineering novel GPI-anchored proteins. Thus, constructing chimeric polypeptides that contain the extracellular domain of a polypeptide-anchored protein and a GPI signal peptide has permitted expressing many proteins in a GPI-anchored form on the cell surface.

The proteins that use a GPI anchor are quite varied. More than 100 different GPI-anchored proteins have been identified (see Table 1 for examples). There is, however, no strong correlation of GPI-anchoring with particular types of protein function. For example, cell surface enzymes and adhesion molecules are found in both GPI-anchored and polypeptide anchored forms. Furthermore, although the potential for reversible membrane binding provided by a lipid anchor is useful for regulating protein distribution inside a cell, it would permit continuous loss of proteins from the cell surface (for a general discussion of factors that can affect membrane affinity of lipid-anchored proteins see [Membrane Anchors](#)). It is not known why cell surface proteins use the complex GPI molecule instead of the relatively simple lipid anchors used by cytosolic proteins. Two biophysical properties of the GPI anchor, however, may help to reduce dissociation from the cell surface: (1) The GPI anchor usually contains two relatively long (i.e., >16-carbon) unbranched acyl or alkyl chains, which in combination would have relatively high affinity for the membrane (see Fig. 1). (2) Unlike other types of lipid anchors, there is a long and relatively flexible linking molecule, the glycan backbone (see Fig. 1) between the lipid in the bilayer and the protein. This type of linkage reduces the entropy cost of membrane binding because there are minimal restrictions on translational and rotational motion of the protein (see [Membrane Anchors](#)).

Although their functional advantages are uncertain, GPI anchors confer some unusual and well-studied properties on cell-surface proteins: (1) GPI anchors are hydrolyzed by specific phospholipases that remove the hydrophobic lipid group, resulting in a protein that is no longer anchored to the membrane (see Fig. 1). GPI-phospholipase D is a secreted enzyme present in mammalian tissues and plasma and could release proteins from the cell surface. The mechanisms by which this enzyme is regulated are unknown, however (5). (2) Under some conditions, GPI-anchored proteins associate preferentially with particular types of lipid molecule in the membrane lipid bilayer. It has been suggested that GPI-anchored proteins are clustered in particular regions at the cell surface (e.g., microdomains and caveolae), but the size, composition, and functional significance of the clusters are controversial (6). (3) Cross-linking GPI-anchored proteins on the cell surface produces an activating signal in some cell types (7). Because the lipid anchor does not cross the membrane, the mechanism by which the signal is transmitted, and its possible connection to microdomains is currently an area of considerable research interest. (4) In spite of the high membrane affinity of the anchor, GPI-anchored proteins can transfer from one cell to another both *in vitro* and *in vivo* (8, 9). The precise mechanism of this process and its physiological significance are unknown. Intercellular transfer is of considerable therapeutic interest, however, because it offers a relatively efficient method for inserting novel GPI-anchored proteins (natural or engineered) into the surface of a patient's cells *ex vivo* (10).

#### Bibliography

1. M. J. McConville and M. A. J. Ferguson (1993) *Biochem. J.* **294**, 305–324.
2. V. L. Stevens (1995) *Biochem. J.* **310**, 361–370.
3. J. Vidugiriene and A. K. Menon (1994) *J. Cell Biol.* **127**, 333–341.
4. S. Udenfriend and K. Kodukula (1995) *Annu. Rev. Biochem.* **64**, 563–591.
5. J.-Y. Li, K. Hollfelder, K.-S. Huang, and M. G. Low (1994) *J. Biol. Chem.* **269**, 28963–28971.
6. S. Mayor, K. G. Rothberg, and F. R. Maxfield (1994) *Science* **264**, 1948–1951.
7. D. Brown (1993) *Curr. Opin. Immunol.* **5**, 349–354.
8. J. L. Miller, M. Giattina, E. J. Blanchette Mackie, and N. K. Dwyer (1998) *J. Lab. Clin. Med.*, **131**, 215–221.
9. S. Ilangumaran, P. J. Robinson, and D. C. Hoessli (1996) *Trends Cell Biol.* **6**, 163–167.
10. M. E. Medof, S. Nagarajan, and M. L. Tykocinski (1996) *FASEB J.* **10**, 574–586.

#### Suggestions for Further Reading

11. M. G. Low (1989) The glycosyl-phosphatidylinositol anchor of membrane proteins, *Biochim.*

Biophys. Acta **988**, 427–454.

12. M. C. Field and A. K. Menon (1993) In *Lipid Modifications of Proteins* (M. J. Schlesinger, ed.), CRC Press, Boca Raton, pp. 83–134.
13. T. Kinoshita, N. Inoue, and J. Takeda (1995) Defective glycosyl phosphatidylinositol anchor synthesis and paroxysmal nocturnal hemoglobinuria, *Adv. Immunol.* **60**, 57–103.
14. P. J. Casey and J. E. Buss (eds.) (1995) "Lipid modifications of proteins", In *Methods in Enzymology* **250**, Academic Press, New York.
15. S. Ilangumaran and D. C. Hoessli (1998) *Glycosylphosphatidylinositol-Anchored Biomolecules*, R. G. Landes, Austin, TX.

## Gram-Negative Bacteria

Two different groups of bacteria can be distinguished on the basis of the differential Gram-staining. [Gram-Positive Bacteria](#) appear bluish-purple, while Gram-negative bacteria are red. The different reactions to Gram-staining are caused by differences in the cell wall structure and composition.

### 1. The staining procedure and its history

The Danish physician Hans Christian Gram (1853–1938), by profession a pharmacologist and internist, worked temporarily with Friedländer in Berlin and described in 1884 (1) the staining of micrococci in lung tissues with Ehrlich's anilin-water solution of crystal methyl (gentian) violet, a cationic dye. After staining, the preparations were treated with Lugol's solution of iodine in potassium iodide and finally flushed with alcohol briefly. The bluish-purple color was retained by the Gram-positive cocci but discharged from the Gram-negative bacteria. Gram used Bismarckbraun or vesuvin for the counterstain. This method has been modified by many workers, but its essence remains unaltered, and it is still one of the widely used methods of differential staining in bacteriology. At present, the destained Gram-negative bacteria are counterstained with a red dye, e.g. safranin or fuchsin, to make a contrast between the Gram-positive and Gram-negative bacteria (2).

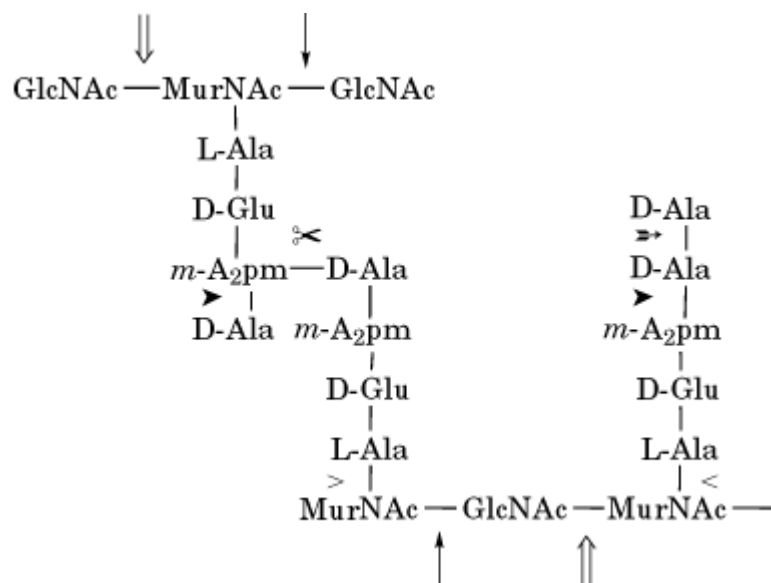
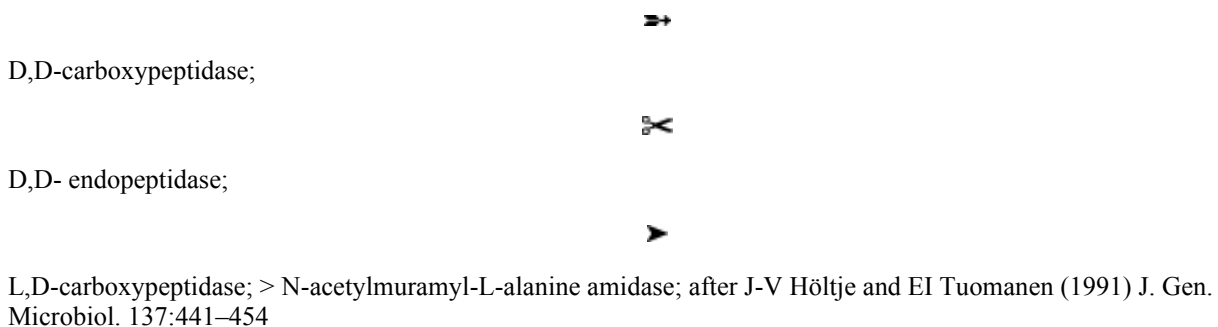
The crystal violet-iodine complex formed during Gram-staining is retained by the cell wall of Gram-positive bacteria, most particularly by the thick layer of peptidoglycan (murein). It is extracted from Gram-negative cells, and from Gram-positive cells of which the peptidoglycan has been removed or loosened by enzymatic (murein hydrolases or peptidases), mechanical, or antibiotic treatment. If the film of Gram-negative cells is too thick, or the bacteria are locally clumped together, the decolorization with alcohol will not be complete. Textbooks mention several exceptions to the rule, and some bacteria are described as being variable in the Gram reaction. Reliable results are dependent on standardization of the preparation of the bacterial films on slides and of performance of the staining procedure. In staining for transmission [electron microscopy](#) and energy-dispersive X-ray spectroscopy, the treatment with Gram's iodine solution has been replaced by treatment with trichloro(h<sup>2</sup>-ethylene) platinate (3).

### 2. The Gram-negative cell envelope

The cell envelope of Gram-negative bacteria, eg, the proteobacteria *Escherichia coli* and *Pseudomonas fluorescens*, consists of an outer membrane, which surrounds the periplasm, and the cytoplasmic or inner membrane, which delimits the protoplast. In electron microscopic pictures of thin sections, the two membranes have similar appearances. The space between the outer and the cytoplasmic membranes, the periplasm, contains soluble proteins (enzymes and components of the

transport machinery) and the murein sacculus. The latter is a single-layered peptidoglycan network. Peptidoglycan is a characteristic cell-wall component of nearly all eubacteria (bacteria) but not of Archaeobacteria (Archaea), and is shown schematically in Fig. 1. Its main function is to preserve the shape and integrity of the cell and to withstand the osmotic pressure of the protoplast. Peptidoglycan consists of long linear glycan chains interlinked with short peptide bridges. The chemical structure and the sequences of aminosugars and of amino acids of the peptide portion are conserved in most Gram-negative bacteria, but can vary in specific taxa, especially in Gram-positive bacteria.

**Figure 1.** Structure of the peptidoglycan of *Escherichia coli*. N-Acetylglucosamine (GlcNAc) and N-acetylmuramic acid (MurNAc) are interlinked in an alternating sequence by  $\beta$ -1,4 glycosidic bonds forming polysaccharide chains that are crosslinked by peptides consisting of L- and D-alanine, D-glutamic acid and *m*-diaminopimelic acid (*m*-A<sub>2</sub>pm). Murein hydrolases cleave specific bonds: muraminidase;  $\beta$ -N-acetylglucosaminidase;

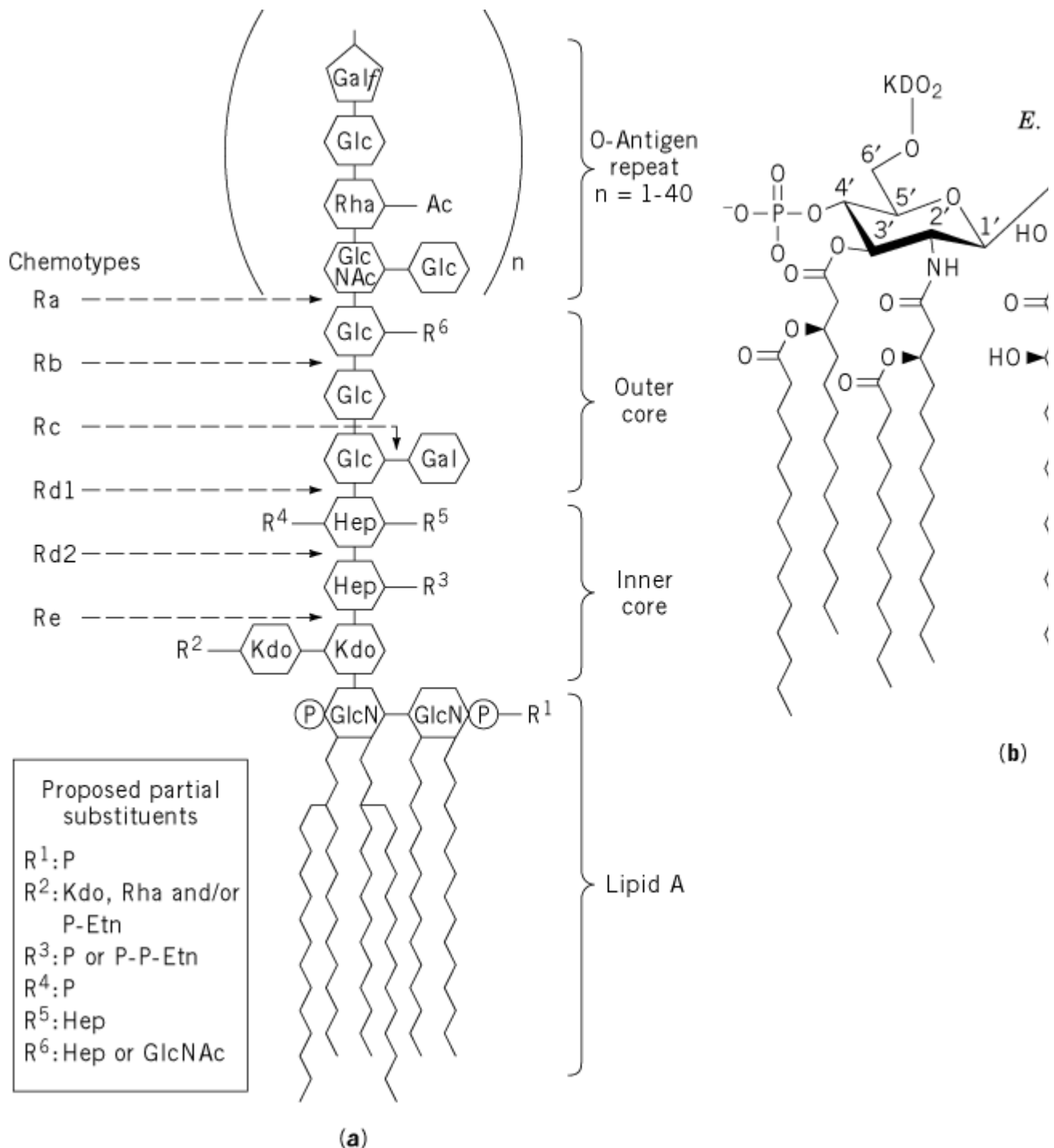


The disaccharide-peptide monomer is synthesized in the cytoplasm as the UDP-activated form and is translocated as N-acetylglucosamine (GlcNAc)  $\beta$ -1,4 N-acetylmuramic acid-(MurNAc)-pentapeptide-pyrophosphoryl-undecaprenol across the cytoplasmic membrane. In a highly controlled process, the precursor is inserted into the murein sacculus by transglycosylation in the glycan chain and subsequent crosslinking by transpeptidation between the penultimate D-alanine residue of one subunit and the *meso*-L-diaminopimelyl residue of another subunit, with release of the terminal alanine (see Fig. 1 and [Penicillin-Binding Proteins](#)). Extension of the murein sacculus is the result of a well-regulated balance between lytic and synthetic enzymatic reactions (4). Termination of the glycan strands by formation of a 1,6-anhydro-N-acetylmuramic acid residue is coupled to dephosphorylation of the undecaprenyl-P-P to undecaprenyl-P, which cycles back into the cytoplasm.

The outer membrane is organized asymmetrically. The inner leaflet consists of phospholipids, as in

the cytoplasmic membrane. The major component of the outer leaflet is the lipopolysaccharide (LPS), a glycolipid that is anchored in the outer membrane by lipid A. Lipid A is the conservative part of LPS at least in Enterobacteriaceae, eg *Salmonella*, *Shigella*, *Escherichia*, *Enterobacter*, *Klebsiella* etc. and is responsible for the physiological effects of LPS in the infected host cell ([endotoxin](#)). The structure of lipid A is shown in Fig. 2(b). Lipid A is joined to the core of LPS, which is a spacer to the outermost and variable part of LPS, the O-specific polysaccharide. The latter is the major antigenic determinant of Gram-negative bacteria; it consists of long chains of repeating units of oligosaccharides (Fig. 2a). The lipophilic character of the outer membrane is presumably responsible for the exclusion of cell-damaging compounds like the bile salts of the intestine. In several bacteria, the LPS functions as a receptor for bacteriophages, eg **T<sub>4</sub> phage**. The core-lipid A and O-antigen subunits are synthesized as activated precursors in the cytoplasm at the cytoplasmic membrane and exported, attached to bactoprenol pyrophosphate, to the periplasm and the outer membrane (5).

**Figure 2.** Schematic structure of *E. coli* K-12 lipopolysaccharide (LPS). (a) The three parts of LPS and the sugar compositions of the core and O-specific chains. Gal, D-galactose; GlcN, D-glucosamine; Kdo, 3-deoxy-D-manno-octulosonic acid; Hep, L-glycero-D-manno-heptose; GlcNAc, N-acetyl-D-glucosamine; Rha, L-rhamnose; Galf, D-galactofuranose; P, phosphate; P-Etn, phosphoethanolamine; Ac, acetate. (b) Structure of lipid A from *E. coli*. The b,1'-6-linked disaccharide of glucosamine is acylated, phosphorylated, and glycosylated with Kdo; the fatty acid residues bound to the disaccharide are laurate and myristate; after C.R.H. Raetz, Figs. 3 and 18A, in F.C. Neidhardt et al., eds., (1996) *Escherichia coli* and *Salmonella*, pp. 1035–1063. ASM Press.



[Porins](#) are transmembrane proteins formed of three identical subunits of [b-barrels](#) with short alpha-helical domains (see Alpha-helix). They are channels for the entrance and exit of small hydrophilic substances. Some porins are highly substrate-specific. Other outer membrane proteins function as phage receptors or mediate export and import. Covalently bound to peptidoglycan is the lipoprotein. Lipoproteins contain lipids covalently linked to an N-terminal [cysteine](#) residue.

The biosynthesis of all the macromolecular components of peptidoglycan and outer membrane starts in the cytoplasm. The pre-proteins and the precursors or building blocks of peptidoglycan or polysaccharides are exported across the cytoplasmic membrane by specific transport systems and are assembled in the periplasm or outer membrane. The energy for assembly cannot be provided in the periplasm by sources of the cytoplasm, so the precursors have to be in an activated state. The



precursors of LPS, peptidoglycan, and capsular polysaccharides are transported using undecaprenol phosphate as a cofactor (5).

### 3. Taxonomy

Although all Gram-negative bacteria have basically the same cell wall structure, they do not form a homogenous taxonomical group. The systematics, and especially the phylogenetic relationships (see [Phylogeny](#)) between the bacteria, are based on different principles and features (6, 7). The Gram-reaction, ie the cell-wall type, is only one, but an important, feature of a genus or species, and there are exceptions to the rule. The cyanobacteria, which perform an oxygenic type of [photosynthesis](#), contain a typical outer membrane with LPS. But the murein sacculus is thicker than that of proteobacteria, and in some species a polysaccharide is bound covalently to the murein. This may be one reason why several cyanobacteria stain Gram-positive.

The Gram-negative archaeobacteria have cell walls of different composition; these are methanochondroitin ( $[4 \rightarrow] - \beta; -D - \text{GluUA} - [1 \rightarrow 3] - D - \text{GalNAc} - [1 \rightarrow 3 \text{ or } 4] - D - \text{GalNAc}$  in *Methanosarcina*), glycoproteins, or heteropolysaccharides composed of uronic acids containing sulfate, glucose, galactose, N-acetylgalactosamine, glycine, mannose, and N-acetyl-gulosaminuronic acid groups (3, 7).

#### Bibliography

1. H. C. Gram (1884) *Fortschritte der Medizin* **2**, 185–189
2. P. Gerhardt, R. G. E. Murray, W. A. Wood, and N. R. Krieg (1994) *Methods for General and Molecular Bacteriology*, Am. Soc. Microbiol., Washington, D.C., pp. 31–32.
3. T. J. Beveridge and S. Schultze-Lam (1996) *Microbiology* **142**, 2887–2895.
4. J.-V. Holtje (1998) *Microbiol. Mol. Rev.* **62**, 181–203.
5. C. R. H. Raetz (1996) *Escherichia coli and Salmonella* (F. C. Neidhardt et al., eds.), ASM Press, Washington, D.C., pp. 1035–1063.
6. J. G. Holt, N. R. Krieg, P. H. A. Sneath, J. T. Staley, and S. T. Williams (1994) *Bergey's Manual of Determinative Bacteriology*, 9<sup>th</sup> ed., Williams & Wilkins, Baltimore.
7. E. Stackebrandt (1999) *Biology of the Prokaryotes* (J. Lengeler, G. Drews, and H. G. Schlegel, eds.), Thieme Verlag, Stuttgart, pp. 674–720.

#### Bibliography

1. H. C. Gram (1884) *Fortschritte der Medizin* **2**, 185–189
2. P. Gerhardt, R. G. E. Murray, W. A. Wood, and N. R. Krieg (1994) *Methods for General and Molecular Bacteriology*, Am. Soc. Microbiol., Washington, D.C., pp. 31–32.
3. T. J. Beveridge and S. Schultze-Lam (1996) *Microbiology* **142**, 2887–2895.
4. J.-V. Holtje (1998) *Microbiol. Mol. Rev.* **62**, 181–203.
5. C. R. H. Raetz (1996) *Escherichia coli and Salmonella* (F. C. Neidhardt et al., eds.), ASM Press, Washington, D.C., pp. 1035–1063.
6. J. G. Holt, N. R. Krieg, P. H. A. Sneath, J. T. Staley, and S. T. Williams (1994) *Bergey's Manual of Determinative Bacteriology*, 9<sup>th</sup> ed., Williams & Wilkins, Baltimore.
7. E. Stackebrandt (1999) *Biology of the Prokaryotes* (J. Lengeler, G. Drews, and H. G. Schlegel, eds.), Thieme Verlag, Stuttgart, pp. 674–720.

## Gram-Positive Bacteria

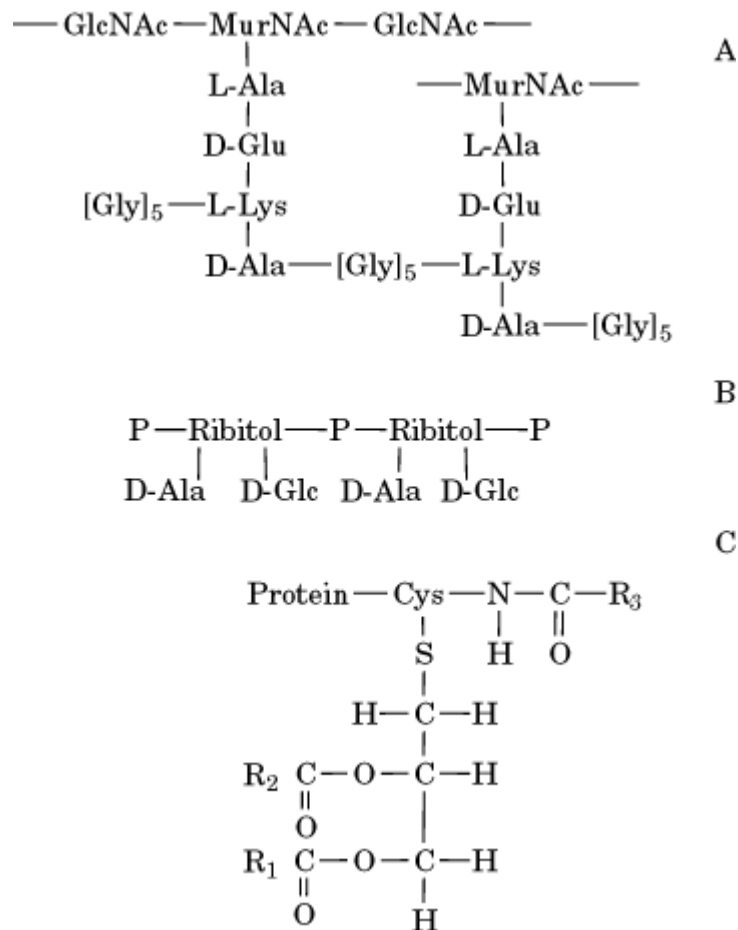
Gram-positive bacteria appear bluish-purple in the Gram stain, while [Gram-Negative Bacteria](#) are red. The crystal violet-iodine complex formed during Gram-staining is retained by the cell wall of Gram-positive bacteria.

Gram-positive eubacteria (bacteria) are found in several taxonomic groups. Typical representatives are both endospore-forming and nonsporing rods of the genera *Bacillus*, *Clostridium*, *Arthrobacter*, *Lactobacillus*, *Mycobacterium*, *Nocardia* and *Streptomyces*. A few archaeobacteria (Archaea), such as *Methanobacterium formicicum*, also stain Gram-positive. The Gram reaction depends on both the structural formats and the chemical composition of the cell envelope (see [Gram-Negative Bacteria](#)). Gram-positive bacteria have a broader range of cell-wall types than do Gram-negative bacteria (1).

### 1. Type I Gram-positive cell wall

This type of cell wall consists of a thick peptidoglycan layer mixed with teichoic acids. Peptidoglycan (murein) is a polymer of glycan strands of N-acetyl-D-glucosamine and N-acetylmuramic acid linked in alternating sequence by  $\beta$ -1,4 linkages. The glycan strands are cross-linked by peptide chains of varying complexity, which are bound in amide links to the carboxyl group of the lactic acid residue of muramic acid (Fig. 1a). In Gram-positive bacteria, the amino acid composition of the peptides cross-linking the glycan strands (interpeptide bridges; Fig. 1a) vary considerably. In these interpeptide bridges, D-amino acids dominate, while aromatic, branched-chain, and sulfur-containing amino acids, plus [histidine](#), [arginine](#), and [proline](#), are absent. In the tetrapeptide of Gram-positive bacteria, diaminopimelic acid can be replaced by [lysine](#).

**Figure 1.** Schematic structure of polymers in cell walls of Gram-positive bacteria. A. Peptidoglycan with an polyglycine interpeptide bridge, [Gly]<sub>5</sub> in *Staphylococcus aureus*. B. Teichoic acid of the ribitol-phosphate type. C. Lipoprotein with fatty acids (R<sub>1</sub> to R<sub>3</sub>) linked directly or via glycerol to the terminal cysteinyl residue.



*Teichoic acids* are polymers consisting of glycerol phosphate or ribitol phosphate that are substituted with sugars and/or amino acids (Fig. 1b). They are covalently bound to the peptidoglycan and are exposed on the cell surface. Several Gram-positive bacteria produce, in addition to teichoic acids, *lipoteichoic acids*, which are not covalently bound to peptidoglycan but are associated with the cytoplasmic membrane through **hydrophobic** interactions. They consist of poly C1-3-glycerol phosphate linked to diglucosyl-1-3-diacylglycerol. Lipoteichoic acid penetrates the peptidoglycan and reaches the cell surface. *Teichuronic acids* of Gram-positive bacteria contain hexuronic acids (glucuronic acid or N-acetylmannosamine uronic acid) instead of polyol phosphate. Teichoic, teichuronic, and lipoteichoic acids are highly-charged polymers. Lipoproteins are present in both Gram-negative and Gram-positive cell walls. The terminal [cysteine](#) residue of the protein moiety is bound to [fatty acids](#) directly or via glycerol residues (Fig. 1c). [Proteins](#) are also constituents of the Gram-positive cell wall, eg **porins** and surface layer-forming proteins.

Activated precursors of the cell wall polymers are transferred across the cytoplasmic membrane (CM) bound to undecaprenol phosphate (undec-P), which was synthesized from mevalonic acid. The undec-P is recycled via undec-P-P and undec.

## 2. Type II Gram-positive cell wall

This type of cell wall is represented by the *Corynebacterium-Mycobacterium-Nocardia* group of bacteria. These bacteria have a thick peptidoglycan layer that is, however, penetrated by glycolipids. The muramic acid may be N-glycosylated. Long chains of *lipoarabinomannan* are anchored in the cytoplasmic membrane and penetrate the peptidoglycan layer. The external layer of the cell wall consists of lipooligosaccharides, phenolic glycolipids, glycopeptidolipids, and mycolic acids esterified to arabinogalactan, which is attached to the peptidoglycan. Mycolic acids are long-chain

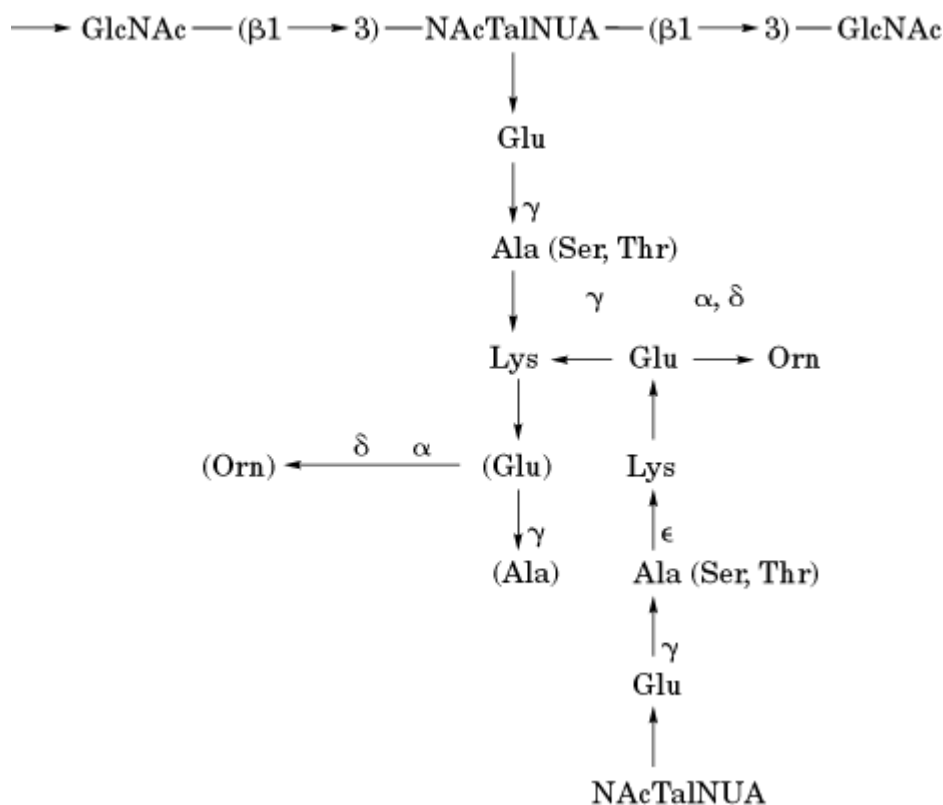
(40–90 C-atoms) alcohols and fatty acids partially unsaturated and substituted with keto-, methoxy-, epoxy- and OH-groups (2).

The type II of Gram-positive cell wall stains Ziehl-Neelsen positive (acid-alcohol fastness). A mixture of the dye basic fuchsin and phenol is heated on the microscopic slides and penetrates the lipid layer. The slides are subsequently washed in distilled water and decolorized with acid-alcohol. By this procedure the fuchsin is removed from other bacteria and cells but is retained by the acid-fast bacteria. When the cells of these bacteria are treated with alkaline ethanol the lipid content of the cell wall is reduced and the cells became non-acid-alcohol-fast but remained Gram-positive. This result shows that the Gram-positive staining depends on a thick peptidoglycan layer.

### 3. Type III Gram-Positive Cell Wall

This type of cell wall is present in few Archaea, eg *Methanobacterium formicicum*, containing pseudomurein (1). The glycan strands of pseudomurein are composed of alternating b-N-acetyl-L-talosamine-uronic acid b-1 ← 3 N-acetyl-D-glucosamine and the tetrapeptide contains L-amino acids (Fig. 2).

**Figure 2. Pseudomurein in some methanogenic Archaea.** GlcNAc, N-acetyl-D-glucosamine; NAcTalNUA, b-N-acetyl-L-talosamine uronic acid; Orn, ornithine; Ser, serine; Thr, threonine; Ala, alanine; Lys, lysine; Glu, glutamic acid.



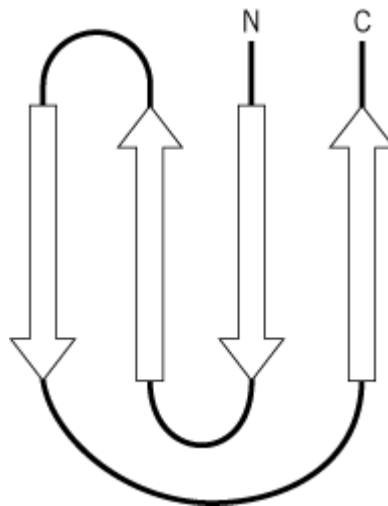
### Bibliography

1. E. Stackebrandt (1999) In *Biology of the Prokaryotes* (J. Lengeler, G. Drews, and H. G. Schlegel, eds.), Thieme Verlag, Stuttgart, pp. 674–720.
2. P. J. Brennan and H. Nikaido (1995) *Annu. Rev. Biochem.* **64**, 29–63.

## Greek Key Motif

The Greek key motif describes a particular topology for arranging four **b-strands** into an antiparallel **b-sheet** in [protein structures](#) (Fig. 1). The name comes from the similarity between this b-strand topology and a decorative pattern used in ancient Greece (also called the Greek key). A common **domain** structure in proteins is the Greek key b-barrel, a type of **antiparallel b-barrel**, where two Greek key b-sheets fold together to form an eight-stranded antiparallel b-barrel.

**Figure 1.** Schematic representation of the topology of the Greek key motif in proteins, with individual b-strands of the b-sheet depicted as arrows. The *N*- and *C*-termini of the motif are labeled.



[See also [Beta-Sheet](#) and [Antiparallel Beta-Barrel Motifs](#) and compare with [Beta-meander](#) and [Jelly roll motif](#).]

### Suggestion for Further Reading

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

## Green Fluorescent Protein (GFP)

GFP is a bioluminescent protein derived from the jellyfish *Aequorea victoria*. The **gene** was **cloned** in 1991, and the protein was found to have the unique property of emitting green light upon illumination with long-wave ultraviolet light without any requirement for additional cofactors or other substrates ([1](#), [2](#)). This has made its gene a most useful [reporter gene](#). The 27-kDa GFP protein has a novel barrel-like [protein structure](#) comprised of 11 **b-strands** that surrounds and protects a

cyclic tripeptide chromophore at its center (3). When excited by light at 470 nm, this chromophore **fluoresces** with an emission maxima of 509 nm that can be readily observed using conventional fluorescein isothiocyanate fluorescence optics.

The unique and robust properties of GFP have made it valuable for a wide range of applications. Most importantly, detection of light emission from GFP does not require any cell fixation or cell permeabilization, allowing for the first time visualization of a reporter gene in living cells. This feature has made GFP an exceptionally powerful tool in cell and transgenic biology, because expressing cells simply glow green under conventional [fluorescence microscopy](#) and [confocal microscopy](#). Furthermore, GFP is observed in living cells, so dynamic events may be followed by making real-time or time-lapse observations. GFP has now been used as a reporter gene in cells from bacteria to human. Importantly, high levels of GFP expression have not been associated with any cellular toxicity; indeed, transgenic mice with ubiquitous expression of GFP thrive without problems. However, prolonged exposure of living cells with long-wave ultraviolet light might cause cell damage, and there may be nuclear damage caused by the generation of free radicals upon GFP excitation.

The b-barrel-like structure of GFP acts to protect its [active site](#) and has allowed numerous *N*-terminal and *C*-terminal fusion genes to be made and expressed, without loss of fluorescent activity. Consequently, GFP may be used to study the behavior of tagged proteins within living cells. Proteins have been tagged for a wide range of applications in different organisms. A few of the many examples include tagging [ribonucleoprotein](#) particles in *Drosophila* to study their trafficking into the developing egg chamber (4), tagging the [histone](#) H2B protein to study the dynamics of [chromosome](#) behavior in living human tumor cells (5), and studying [cytoskeletal](#) dynamics by binding GFP to [F-actin](#) in the slime mold [Dictyostelium](#) (6).

The second major use of GFP has been for cell-lineage analysis in developmental biology. GFP has not been found to be toxic in any cell type, and the fluorescence does not bleach when excited with 480-nm light. Furthermore, the protein is highly stable; for example, protein **translated** from [messenger RNA](#) microinjected into a [Xenopus](#) oocyte may still be detected at the feeding tadpole stage of development 5 days later. Thus the fate of the progeny of labeled cells can be followed over significant periods of time during development (7).

Following initial studies with the wild-type protein, several variant GFP isoforms have been generated by mutation analysis (2). Some of these mutations have improved the absorption and emission spectra of the protein, which increase the sensitivity and the versatility of the reporter gene. The mutation Ser65Thr results in an increased brightness of emission, faster fluorophore formation, and reduced photobleaching. Other novel isoforms are particularly valuable, because they emit light at new wavelengths; currently, blue-, red-, and cyan-shifted variants are available (Clontech; [www.clontech.com](http://www.clontech.com)), and these can be used in dual or multiple labeling studies with appropriate filters. While certain mutations affect GFP function, other ‘silent’ mutations have been generated that customize the gene for expression in mammalian cells or in bacteria by optimizing **codon-usage** preferences. Advances in GFP technology can be followed at the “Fluorescent protein Newsgroup” on the Internet (<http://www.bio.net/hypermail/FLUORESCENT-PROTEINS>).

The applications of GFP as a reporter gene are largely determined by its stability and sensitivity. The b-barrel-like structure that protects the chromophore is very stable, resisting temperatures up to 65°C and ionic [detergents](#), and gives high resistance to [proteinases](#) (3). The long half-life of GFP makes it ideal for long-term cell-labeling or protein-tagging studies, but it is less suitable for monitoring dynamic changes in gene expression. To circumvent this problem, novel GFP fusion genes have been made with the [PEST region](#) of the mouse [ornithine decarboxylase](#) gene at the *C*-terminus, which destabilizes the GFP to **protein degradation**; it has a half-life of 1 to 2 hours (8).

Although GFP can be engineered to generate proteins of different half-lives and with different fluorogenic properties, less can be done to improve its sensitivity, because it lacks the ability to

amplify its signal through enzymatic turnover of added substrate. In all successful applications thus far, high levels of GFP expression have had to be obtained, achieving of the order of  $10^6$  molecules per cell or more. The sensitivity of detection is therefore orders of magnitude less than that of other enzyme-based reporter systems, and the detected expression is not tightly coupled to changes in gene transcription.

### Bibliography

1. T. Misteli and D. L. Spector (1997) *Nature Biotech.* **15**, 961–964.
2. A. B. Cubitt et al (1995) *Trend Biochem. Sci.* **20**, 448–455.
3. G. N. Phillips Jr. (1997) *Curr. Opin. Struct. Biol.* **7**, 821–827.
4. S. Wang and T. Hazelrigg (1994) *Nature* **369**, 400–403.
5. T. Kanda, K. F. Sullivan, and G. M. Wahl (1998) *Curr. Biol.* **8**, 377–385.
6. K. M. Pang, E. Lee, and D. A. Knecht (1998) *Curr. Biol.* **8**, 405–408.
7. M. Zernicka-Goetz et al. (1996) *Development* **122**, 3719–3724.
8. Clontech (1998) **XIII** (2), 16–17, Clontech Laboratories, Inc., Palo Alto, CA. (HYPERLINK <http://www.clontech.com>).

### Growth Factors

Peptide growth factors provide a critical foundation for intercellular communication in multicellular organisms and are key elements in a complex biological process. But while they promote cell growth, just as often they can have opposite, inhibitory effects on cell proliferation, or instead modulate differentiated cell function. They regulate numerous alternative cell functions that are related to, but sometimes distinct from, cell proliferation. Their activities are often distinct, one from each other, but they can also be significantly redundant; alone, they may exert one activity, but in combination have different, even opposite, effects on cells. To understand their activities, their actions must be defined individually, and then evaluated in the context of complex biological responses.

There are a number of specific characteristics that define the specific properties of a protein growth factor. First and foremost, they are [proteins](#), usually of molecular weights ranging from 5000 to 80,000 daltons. Second, they are ligands that modulate cell function through cell-surface high affinity receptors. On binding the ligand, the receptor activates a **signal-transduction** cascade that, when evaluated on cells in culture, stimulates a proliferative response that includes (but is not limited by) thymidine incorporation into DNA and cell division.

Different peptide growth factors belong to one or another family of structurally related proteins. Hence, numerous proteins that contain sequences homologous to [epidermal growth factor](#) (EGF) belong to the EGF family, others with homology to platelet derived growth factor (PDGF) belong to the PDGF family, and still others related to **nerve growth factor** (NGF), **transforming growth factor-b** (TGF b), hepatocyte growth factor (HGF) all belong to the NGF, PDGF, TGF b, or HGF family, respectively. A representative list of these families and some of their members is shown in [Table 1](#).

**Table 1. Growth Factor Families**

| Growth Factor Family               | Prototypic Target Cells               | Acronyms         |
|------------------------------------|---------------------------------------|------------------|
| Angiopoietins                      | Endothelial cells                     | AP1, AP2         |
| Bone morphogenic proteins          | Osteoblasts, osteoclasts, fibroblasts | BMP-1, BMP-2     |
| Brain-derived growth factors       | Motor and sensory neurons             | BDNF             |
| Ciliary neurotrophic factors       | Motor and sensory neurons             | CNTF             |
| Colony stimulating factors         | Macrophages, granulocytes             | M-CSF, GM-CSF    |
| Epidermal growth factors           | Epidermal cells, fibroblasts          | EGF, TGFA        |
| Erythropoietin                     | Hematopoietic stem cells              | EPO, TPO         |
| Fibroblast growth factors          | Mesenchymal neuroectoderm cells       | FGFs and KGFs    |
| Glial-derived neurotrophic factor  | Motor and sensory neurons             | GDNF             |
| Hepatocyte growth factors          | Hepatocytes, endothelial cells        | HGF-1, HGF-2     |
| Insulin-like growth factors        | Mesenchymal neuroectoderm cells       | IGF-1, IGF-2     |
| Interferons                        | Immunoregulatory cells, fibroblasts   | g-IF, a-IF       |
| Interleukins                       | Immunoregulatory cells, fibroblasts   | IL-1, IL-2, IL-3 |
| Neurotrophins                      | Motor, sensory, autonomic neurons     | NGF, NT-2, NT-3  |
| Platelet-derived growth factor     | Fibroblasts, mesenchymal cells        | PDGF-A, PDGF-B   |
| Tumor necrosis factors             | Fibroblasts, inflammatory cells       | TNF-A, TNF-B     |
| Vascular endothelial growth factor | Vascular, lymphatic endothelial cells | VEGF-A, B, C, D  |

When peptide growth factors were first being isolated and identified, they were described either by their biological activity or by the source from which they were purified. Hence, the first [fibroblast growth factors](#) (FGFs) were isolated on the basis of their ability to stimulate the proliferation of fibroblasts, NGF by its ability to stimulate nerve cells, TGFs for their ability to transform cells, HGF for its ability to stimulate hepatocytes, and PDGF because it was first purified from platelets. However, the original discovery of a peptide growth factor has seldom reflected the extent of the protein's activities or, in fact, its true physiological function. Few recall that one of the most extensively studied growth factors, EGF, was first named *urogastrone* and was originally characterized for its ability to suppress gastric acid secretion (1). A powerful peptide growth factor in its own right, EGF is a potent growth factor for epithelial cells and yet, in 1962, was isolated and purified on the basis of its ability to cause premature eruption of incisor teeth in mice and not cell proliferation per se (2). Similarly, before the advent of modern cell culture, the detection, isolation, and identification of NGF in 1953 (3), 1960 (4) and 1971 (5), respectively, were based on *in vivo*



studies and tissue explants that showed increased neurite outgrowth from nerve ganglia. But is NGF a true nerve growth factor when it stimulates neurite outgrowth, but not nerve cell proliferation? Even more importantly, NGF can stimulate hematopoietic cell colony growth and differentiation (6) and has receptors on the cell surface of fibroblasts, lymphocytes, and numerous mesenchymal cell types that are unrelated to cells of neural lineage (7, 8). The detection of both ligand and receptor in reproductive tissues points to possible functions in reproduction.

In the late 1970s, and throughout the 1980s, the advent of modern methods of cell biology gave investigators the ability to grow numerous cell types in culture, and there was a flurry of activity to define the conditions in which each cell type could grow (9). The first efforts were aimed at replacing serum with tissue extracts so as to identify the active components. Hence, the growth factor literature became saturated with descriptions of “activities” known as fibroblast growth factors, endothelial cell growth factors, cartilage growth factors, milk-derived growth factors, adrenal angiogenic factor, corpus luteum growth factor, brain growth factor, pituitary growth factor, to name only a few. On applying the then developing techniques of molecular cloning, it soon became apparent that all of these seemingly disparate activities were in fact a handful of molecules: FGFs, EGFs, HGFs and PDGFs.

With the realization that, historically, peptide growth factors are a loosely defined collection of molecules that promote cell growth, there have been numerous efforts to rename the field, reorganize its nomenclature, and refine the description of peptide growth factors. While cell biologists have named these molecules “growth factors,” immunologists have called them “**interleukins**, **lymphokines**, and cytokines”. Hematologists, in turn, have coined the phrase [colony stimulating factors](#). All descriptions are found throughout the literature and vary from field to field. In the end, the most common usage has remained the historical names given to molecules as they have been identified. When molecules are “rediscovered” because of a new heretofore unknown activity, they are classified according to their original, and ironically sometimes their least important, activity. Hence, when the adrenal angiogenic factor was sequenced and found to be identical to FGF2, it was renamed as such (10). When keratinocyte growth factor (KGF) was sequenced and found to be a member of the FGF family, it was renamed FGF7 (11). As the field has matured, and more structural and biological information has been generated, it has been possible to sort through the pleiotropic activities of growth factors and their sometimes ubiquitous distribution, to rely more heavily on their structural, rather than biological, characteristics for classification (12).

The difficulty in classifying growth factors and the attempts to distinguish them from interleukins, lymphokines, and cytokines is futile and best illustrated by an examination of their range of target cells. It has already been noted that any given growth factor can have numerous activities on any given cell type and that their activities can expand well beyond the context of their original discovery. Even more remarkable, however, is their range of target cells, which can be much greater than originally anticipated. It is here where attempts to delineate functional differences between interleukins, cytokines, and growth factors have failed. Interleukins and cytokines are generally regarded as signaling molecules with activities in the immune system. They were originally thought to be produced by the immune system, for the immune system. Yet many, if not all, interleukins are broad-spectrum proteins that, when studied in detail, regulate numerous processes outside of the immune system. Interleukin 1 (IL1), for example, acts on chondrocytes, astrocytes, fibroblasts, keratinocytes, and some neurons. Best known as a [B-cell](#) growth factor, interleukin-6 also regulates mesenchymal cell function. Perhaps even more remarkable is the observation that fibroblasts and other nonimmune cells can be sources of interleukins and not just targets, thus diminishing any differences between interleukins and protein growth factors further. As pleiotropic, ubiquitous biological effectors that control cell proliferation, the terms protein growth factor, cytokine, and interleukin describe near synonymous molecules.

Whatever classification might be conferred to growth factors, cytokines, lymphokines, or interleukins, their activities challenge the foundations of classical sciences like endocrinology, because of their peculiar, and yet necessary, mode of action. Unlike the endocrine system that relies

on tissue-to-tissue communication, the growth factor system acts at the cellular level. Thus, while the endocrine pathways rely on blood-borne transportation of signaling molecules to effect change, growth factors mediate alterations in the local biological milieu. This fundamental difference in action has led investigators to rethink how cells communicate and even how they respond to endocrine stimuli. Peptide growth factors, released into the local milieu by one cell, alter the activity of an adjacent cell that is in close proximity (paracrine target) or alternatively the cell from which it derives (autocrine target). The actions of major regulatory peptides are mediated by autocrine and paracrine mechanisms that control growth factors and normal cellular homeostasis. Thus the main difference between the growth factor and cytokine pathways and the endocrine system is in the local action of growth factors. Even lymphokines produced by circulating immune cells act locally, but only after the cell producing it (eg, lymphocytes) has traveled and delivered it to its desired site of action. All growth factors are produced locally to act locally.

With this difference between endocrine and growth factor pathways in mind, it is possible to understand that when growth factors are locally released into biological fluids, they can then act to modify the actions of classic hormones arriving on site from distal sources. Hormone action has long been known to require local cellular factors to permit, mediate, and sometimes modulate responsiveness. For example, [insulin-like growth factors](#) and their binding proteins modulate hormone action (13) and TGF- $\beta$  can modulate adrenal function (14). Furthermore, endocrine hormones can trigger gene expression of growth factors in target cells, which in turn modulate the local milieu and change the cell's responsiveness to further stimulation by the hormone (15).

In view of the fact that peptide growth factors act locally to modulate cell function, it comes as no surprise that there are numerous processes present in the local milieu to control their activities and ensure that their actions do not extend to unwanted targets. Of these processes, four are perhaps the most studied:

- Restricted production and release of the growth factor to limit its action
- Highly regulated receptor gene expression to limit target cell responsiveness
- Circulating inhibitors in blood to sequester and inactivate the ligand
- The binding and sequestration of growth factors at the cell surface and extracellular matrix to keep them in the local milieu

The notion that there is restricted production and release of growth factors and limited receptor expression is best illustrated by the difference between gene expression in embryonic development, fetal growth, and in adult tissues (16). At times as early as before gastrulation, the gene expression of growth factors becomes controlled in both a temporal and spatial fashion, and the role played by growth factors in development is essential in all aspects of patterning, tissue development, and differentiation. As exemplified by studies where the genes for different growth factors have been genetically knocked out, they have been shown to play a critical role during development, but are often suppressed in adult, quiescent tissues. In both cancer and wound healing, however, the genes for these same molecules are activated. In wound healing, they remain regulated until the injury has been resolved. In cancer, activation of these genes leads to uncontrolled cell growth, transformation, abnormal blood vessel growth, and ultimately tumor formation.

One particularly interesting mechanism of controlling growth factor activity appears unique to some of the paracrine factors involved in the injury and inflammatory response. Recent data suggest that molecules like FGF1, FGF2, IL1 $\alpha$ , and IL-1 $\beta$  can exit the cell independently of the cell's secretory [endoplasmic reticulum](#) and [Golgi apparatus](#) (17). This protein export process is mechanistically different from [protein secretion](#) and is highly specific and regulated. While its physiological significance remains unknown, the selective translocation of certain biological effectors across the plasma membrane presumably restricts their activity and allows the cell to tightly control access to the cell surface and their target cells.

The evidence for natural inhibitors of growth factors is equally compelling. For example, numerous investigators have established that genes encoding high affinity receptors for growth factors encode soluble forms of their high affinity receptors that are generated by [alternative splicing](#) of the gene's mRNA or by truncation at the cell surface (18). This form of receptor has no transmembrane domain and no signaling domain. It is secreted by cells and can bind ligand in biological fluids, thus preventing it from interacting with its cell surface (and functional) high affinity receptor. Alternatively, some cells have the ability to produce receptor antagonists that bind, but do not activate high affinity receptors. By virtue of occupying the binding site, molecules like interleukin receptor antagonist (ILRA) and angiopoietin-2 prevent a functional interaction between the ligand and its native receptor.

Perhaps the most remarkable process that appears to control the activity of growth factors stems from the observation that they often appear designed to remain in the local cellular milieu and physically sequestered outside target cells in what is presumed to be a biologically inactive form (19). This observation underlies the need for their action to remain local and predicts the existence of specific signals (eg, enzymes) to render them available to target cells.

The first type of local adhesion is characterized by a noncovalent association of growth factors with the [extracellular matrix](#), the cell surface, or pericellular structures. This binding is often salt- and/or pH-sensitive. Representative examples of growth factors with this type-1 adhesion are shown in Table 2. In contrast, the second type of adhesion is characterized by a covalent association with the extracellular surface, through a transmembrane domain, in the growth factor's precursor. Representative examples of growth factors having this type 2 adhesion are also shown in Table 2.

**Table 2. Growth Factors Associated with Cell Surface through Ionic (Type 1) or Transmembrane Domains (Type 2) in Their Precursors**

| Type 1 Association                 | Type 2 Association            |
|------------------------------------|-------------------------------|
| Fibroblast growth factors          | Epidermal growth factors      |
| Transforming growth factor b       | Tumor necrosis factors        |
| Colony stimulating factors         | G-colony stimulating factors  |
| Insulin-like growth factors        | M-colony stimulating factor   |
| Colony stimulating factor          | GM-colony stimulating factors |
| Interleukins 1 and 3               | Stem cell growth factor       |
| Hepatocyte growth factors          |                               |
| g-Interferon                       |                               |
| Vascular endothelial growth factor |                               |
| Platelet-derived growth factor     |                               |
| Bone morphogenic proteins          |                               |

Fibroblast growth factors (FGFs) are prototypic examples of type 1 adhesion. They have a high affinity for immobilized heparin, and immunohistochemical techniques localize FGF1 and FGF2 to the extracellular matrix, basement membrane, and cell surface of numerous tissues (20). Indeed, recent studies have identified the binding heparan sulfate proteoglycans that serve to sequester (and possibly deliver) FGFs to its high affinity receptors on target cells. TGFb has also been localized in extracellular and pericellular structures, but it is in what is thought to be a latent, biologically inert form (21). Recent studies have identified three kinds of receptors for this protein, one of which is a

high molecular weight proteoglycan anchored onto the plasma membrane through phosphoinositol. Other growth factors with this type 1 adhesion include granulocyte-macrophage colony stimulating factor (GM-CSF), which binds to extracellular matrix and to heparan sulfate-related glycosaminoglycans (GAG), interleukin-3 (IL-3), hepatocyte growth factor (HGF), osteogenin (OTG), and the bone morphogenic proteins (BMP).

Type 2 adhesion is characterized by growth factors that are covalently associated with the plasma membrane because their precursor encodes a transmembrane domain that anchors it to the cell surface. In order to be active, these molecules must either be processed from the precursor to generate an active ligand or must act directly on an adjacent high affinity receptor without being further processed or released. Accordingly, these molecules may have activities that are restricted by availability. The prototypic family (and the most studied) of the growth factors with type 2 adhesion to cells is EGF. It and its structural homologues, TGF $\alpha$  and the heparin-binding EGF, are contained in large precursor proteins that have transmembrane domains that lock them onto the cell surface. Most of the biological studies performed with EGFs have been with the processed unbound forms, so it is less clear if their growth factor activities are maintained when they are unprocessed, at the cell surface. But recent studies using [site-directed mutagenesis](#) support the notion that the membrane-bound precursor has growth factor activity and can interact with high affinity receptors on target cells. The ability of growth factors with this type 2 adhesion to elicit a biological response is illustrated by the recent experiment showing that PDGF anchored to the plasma membrane through mutagenesis is still capable of stimulating a mitogenic response ([22](#)). Other examples of type 2 adhesion include [tumor necrosis factor](#) (TNF) and the hematopoietic factor, CSF-1 and stem cell growth factor, the ligand for the ckit protein. In all instances, each is generated from an integral membrane protein locked onto the cell surface.

It is presumed that a growth factor locked to the cell surface through a transmembrane domain can stimulate a biological response only if it is liberated from its sites of sequestration and delivered to the target cell. The molecular mechanisms that might mediate this regulation are unknown. The extracellular matrix and the various constituents that bind the growth factor are obvious targets for the combined actions of **proteolytic** and glycolytic degradation. These include **plasmins**, **cathepsins**, collagenases, and a wide array of glycanases, such as heparinases and heparitinases. These enzymes are normally inactive due to the presence of a [macroglobulin](#), plasminogen activator inhibitors, [proteinases](#) like nexin, and collagenase inhibitors. In the case of the FGFs, however, it is particularly worthy to note that they are proteinase-resistant when bound to glycosaminoglycans ([23](#)). Thus the generation of an FGF-GAG complex by limited proteolysis provides a mechanism to generate a potent biologically active ligand from a pool of sequestered, biounavailable growth factor. Furthermore, glycosaminoglycan binding to some FGFs has been reported to enhance its binding to their high affinity receptor and in some instances is required for delivery ([24](#)). The mechanisms that regulate the type 2 adhesion of growth are not known, although there exist proteolytic enzymes capable of releasing the ligand from the cell ([25](#), [26](#)).

In the end, as the peptide growth factors that control cell growth and development have been identified and the mechanisms that regulate their activities better understood, the challenge becomes to characterize how they all act together to control normal cell function. As the use of modern molecular techniques to eliminate combinations of growth factors from the genome genetically becomes more widespread, it should be possible to provide important insight into how these factors act to control normal cell function and, ultimately, their role in whole-body homeostasis.

## Bibliography

1. H. Gregory and I. R. Wilshire (1975) *Physiol. Chem.* **356**, 1765–1774.
2. S. Cohen (1962) *J. Biol. Chem.* **237**, 1555–1562.
3. R. Levi-Montalcini and V. Hamburger (1952) *J. Exp. Zool.* **123**, 233–288.
4. S. Cohen (1960) *Proc. Natl. Acad. Sci. USA* **46**, 302–311.

5. R. H. Angeletti and R. A. Bradshaw (1971) *Proc. Natl. Acad. Sci. USA* **68**, 2417–2420.
6. H. Matsuda, M. D. Coughlin, J. Bienstock, and J. A. Denburg (1988) *Proc. Natl. Acad. Sci. USA* **85**, 6508–6512.
7. P. Ernfors, F. Hallbook, T. Ebendal, E. M. Shooter, M. J. Radeke, T. P. Misko, and H. Persson (1988) *Neuron* **1**, 983–996.
8. M. Bothwell, S. L. Patterson, G. C. Schatteman, S. Thomprun et al. (1989) *J. Neurosci.* **22**, 354–362.
9. A. Baird (1993) *Endocrinology* **132**, 487–488.
10. D. Gospodarowicz, A. Baird, J. Cheng, G. M. Lui, F. Esch, and P. Böhlen (1986) *Endocrinology* **118**, 82–90.
11. P. W. Finch, J. S. Rubin, T. Miki, D. Ron, and S. A. Aaronson (1989) *Science* **245**, 752–755.
12. A. Baird and M. Klagsbrun (1991) *Cancer Cells* **3**, 239–243.
13. G. Lamson, L. C. Giudice, and R. C. Rosenfeld (1991) *Growth Factors* **5**, 19–28.
14. J. J. Feige and A. Baird (1992) *Prog. Growth Factor Res.* **3**, 103–113.
15. J. Massague (1990) *Annu. Rev. Cell Biol.* **6**, 597–640.
16. M. Goldfarb (1996) *Cytokine Growth Factor Rev.* **7**, 311–325.
17. A. E. Cleves (1997) *Curr. Biol.* **7**, 318–320.
18. D. E. Johnson, J. Lu, H. Chen, S. Werner, and L. T. Williams (1991) *Mol. Cell Biol.* **11**, 4627–4634.
19. J.-J. Feige and A. Baird (1992) *Med. Sci.* **8**, 805–810.
20. I. Vlodavsky, J. Folkman, R. Sullivan et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 2292–2296.
21. N. L. Thompson, K. C. Flanders, J. M. Smith, L. R. Ellingsworth, A. B. Roberts, and M. B. Sporn (1989) *J. Cell Biol.* **108**, 661–669.
22. B. A. Lee and D. J. Donoghue (1991) *J. Cell Biol.* **113**, 361–370.
23. D. Gospodarowicz and J. Cheng (1986) *J. Cell Physiol.* **128**, 475–484.
24. A. Yayon, M. Klagsbrun, J. D. Esko, P. Leder, and D. M. Ornitz (1991) *Cell* **64**, 841–848.
25. J. Teixidó, R. Gilmore, D. C. Lee, and J. Massagué (1987) *Nature* **326**, 883–885.
26. A. Pandiella and J. Massagué (1991) *Proc. Natl. Acad. Sci. USA* **88**, 1726–1730.

### **Suggestions for Further Reading**

27. A. Baird (1994) Fibroblast growth factors: activities and significance of non-neurotrophin neurotrophic growth factors, *Curr. Opin. Neurobiol.* **4**, 78–86.
28. M. J. Bissel, H. G. Hall, and G. Parry (1982) How does the extracellular matrix direct gene expression? *J. Theor. Biol.* **99**, 31–68.
29. M. Y. Gordon (1991) Hemopoietic growth factors and receptors: bound and free. *Cancer Cells* **3**, 127–133.
30. M. Y. Gordon, G. P. Riley, S. M. Watt, and M. F. Greaves (1987) Compartmentalization of a haematopoietic growth factor (GM-CSF) by glycosaminoglycans in the bone marrow microenvironment, *Nature* **326**, 403–405.
31. D. C. Lee (1990) TGF- $\alpha$ : expression and biological activities of the integral membrane precursor. *Mol. Reprod. Dev.* **27**, 37–45.
32. J. Massague (1990) Transforming growth factor . A model for membrane anchored growth factors, *J. Biol. Chem.* **265**, 21393–21396.
33. B. Mroczkowski, M. Reich, K. Chen, G. I. Bell, and S. Cohen (1989) Recombinant human epidermal growth factor precursor is a glycosylated membrane protein with biological activity, *Mol. Cell Biol.* **9**, 2771–2778.
34. R. Roberts, J. Gallagher, E. Spooncer, T. D. F. Bloomfield, and T. M. Dexter (1988) Heparan sulphate bound growth factors: a mechanism for stromal cell mediated haemopoiesis, *Nature*

332, 376–378.

35. A. B. Roberts and M. B. Sporn (1990) *The Multifunctional Nature of Growth Factors in Peptide Growth Factors and Their Receptors*, Springer-Verlag, Berlin.
36. K. Tryggvason, M. Höyhty, and T. Salo (1987) Proteolytic degradation of extracellular matrix in tumor invasion, *Biochim. Biophys. Acta* **907**, 191–217.
37. S. T. Wong, L. F. Winchell, B. K. McCune et al. (1989) *Cell* **56**, 495–506.

## Growth Hormone

Growth hormone (GH), also called somatotropin or somatotrophic hormone, is a protein produced in the anterior pituitary gland and is required for normal growth in vertebrate animals. Its growth promoting action is mediated by insulin-like growth factor-1 (IGF-1) whose synthesis GH stimulates in liver and other target tissues. Growth hormone also has potent metabolic actions in mammals and promotes protein synthesis and lipid degradation while restraining carbohydrate oxidation. GH is comprised of a single chain of ~200 amino acid residues with two internal disulfide bonds and a tertiary structure that includes four alpha helices bundled in an antiparallel fashion. Growth hormone is a member of a family of proteins that are produced in the pituitary glands of all vertebrates and placentae of some mammals. Family members include pituitary prolactin (PRL), placental lactogens (PL), which are also called chorionic somatomammotropins in humans, prolactin-related proteins in ruminants and rodents, proliferins in mice, and somatolactin in fish. Genes that encode most members of the GH family are comprised of five exons and four introns and appear to have arisen by duplication from a single ancestral gene before the appearance of the vertebrates. Splice junctions have been preserved throughout vertebrate evolution. Splicing and processing variants have been described for several members of the family. Growth hormone, PRL, and, presumably, their relatives signal through single-chain membrane receptors that lack intrinsic enzymatic activity, but associate with cytosolic protein tyrosine kinases. Binding of GH to a single-receptor molecule is followed rapidly by the binding of a second receptor molecule to the opposite side of GH to form a trimeric complex that brings the intracellular domains of the two receptor molecules into close apposition and activates associated protein tyrosine kinases of the JAK (Janus kinase) family, particularly JAK2. Similar events occur with PRL signalling and probably the placental lactogens as well. A large number of proteins including the GH receptor and JAK2 become phosphorylated on tyrosine residues to form docking sites for other proteins and initiate signalling cascades. Signalling to the nucleus occurs through the agency of phosphorylated Stat (signal transduction and activation of transcription) proteins and by the MAP kinase and protein kinase C pathways as well.

### 1. History

GH was discovered in 1921 by Evans and Long (1) as the substance present in extracts of bovine pituitary glands that increased growth in rats when given daily over a period of 3 months. A decade later, Riddle's group prepared by isoelectric precipitation a pituitary extract that was rich in PRL, a hormone that promotes milk production (2). They subsequently found that the same partially purified extract promoted growth in pigeons and mice and questioned the existence of a separate and unique growth promoting hormone (3). The ensuing controversy foreshadowed later discoveries of the close relationship of GH and PRL in terms of their chemical structures as well as their overlapping biological activities in various assay systems. Available methods of protein purification and hormone assay were quite limited at that time, and, despite continued intense efforts by the Evans group, nearly a quarter of a century elapsed before it was established with virtual certainty that the growth promoting activity of the pituitary gland resided in a single molecular entity. In 1944 Li and Evans

(4) announced the successful isolation of a highly purified preparation of bovine GH that was apparently free of PRL and the activities of other anterior pituitary hormones. Despite its readily demonstrable efficacy in promoting growth in experimental animals, however, no therapeutic benefit to human pituitary dwarfs treated with bovine GH could be demonstrated in the subsequent decade of clinical trials. Knobil and Greep (5) found that GH of bovine or porcine origin was totally ineffective in rhesus monkeys and adopted the then unorthodox view that there might be species differences in the activities of GH. Fortuitously, their work coincided with large scale production of the newly developed Salk polio vaccine that required monkey kidney tissue for growth of the viruses and, hence, the sacrifice of large numbers of rhesus monkeys. They demonstrated that GH isolated from monkey pituitary glands promoted growth and nitrogen retention in hypophysectomized rhesus monkeys. Subsequent findings that GH prepared from monkey or human pituitaries is effective in humans provided the impetus for creation in 1960 of the National Pituitary Agency for large scale collection of human pituitaries at autopsy and centralization of the production of human GH in the United States. Human GH made available through this agency was used for basic research and for the treatment of more than 3,000 patients over the next quarter century (6).

The complete amino acid sequences of GHs derived from human (hGH), monkey (mGH), and bovine (bGH) pituitary glands were known by the early 1970s, and all efforts to modify animal GHs to forms that might be effective in humans or to synthesize human GH by chemical means were unsuccessful. Meanwhile, with the therapeutic need for the scarce hGH and the accumulating wealth of structural information on GH and its related proteins (7), GH was recognized as an attractive model protein for application of the emerging recombinant DNA technology. It was soon shown that GH could be synthesized in bacteria from mammalian DNA fused to a bacterial gene (8). Human GH was the first mammalian protein to be so produced from purely mammalian cDNA (9) and was targeted by the newly formed biotechnology company, Genentech, as the first company to develop and market recombinant protein commercially (10). Approval for clinical use of recombinant human GH was granted by the U.S. Food and Drug Administration in 1985 after several deaths from Creutzfeldt-Jakob disease attributed to pituitary hGH contaminated with the slow virus resulted in a ban on its further administration to human patients. In the decade between 1985 and 1995, ten times as many children were treated with the recombinant hormone as had been treated with the pituitary-derived hormone in the preceding quarter century, and, collectively, they are said to have grown to an additional height of about 4 miles (10). Ready availability of recombinant bovine and porcine GH has also fostered commercial application in agriculture where GH can be used to accelerate growth of domestic livestock, increase feed efficiency, decrease carcass fat, and increase milk yield. Political opposition, however, has limited its use in this regard. Overexpression of GH in transgenic cattle, sheep, and pigs has produced disappointing results, as these animals often become diabetic.

## 2. Biological Actions

Postnatal growth and the attainment of normal adult size is absolutely dependent upon GH in most species, but guinea pigs and some strains of chickens appear to have evolved alternative mechanisms of growth regulation that appear independent of GH and are not understood. Embryonic growth and growth in the early postnatal period in mammals is also independent of GH. Absence of GH in juvenile animals results in “dwarfism” that varies in severity from less than 25% to 65% of expected adult size. Excessive GH production during the juvenile period leads to a state of giantism in animals and humans, with some afflicted men attaining heights of over 8 feet. Transgenic mice expressing metallothionein-hGH fusion genes attained body weights that were twice more those of littermate controls (11). Cessation of growth in adult animals results from fusion and, hence, is unresponsiveness of the growth plates in the long bones rather than a lack of GH, which continues to be secreted throughout life, albeit it in decreasing amounts with increasing age. Continued overproduction of GH in human adults results in overgrowth of skin and soft tissues and deformities that result from enlargement of those bones that can continue to respond; a condition called acromegaly.

The growth-promoting effect of GH is at least partially, and perhaps entirely, attributable to the

actions of a growth factor, once known as somatomedin, and now called insulin-like growth factor-1 (IGF-1). GH enhances the synthesis and secretion of IGF-1 in a wide variety of tissues. IGF-1 was so named because of the similarity of both its structure and its biological effects to those of insulin (see **IGF-1**, this volume). Most of the IGF-1 in the blood originates in the liver and circulates tightly bound in a complex that contains two other GH-dependent proteins, the IGF binding protein 3 and an acid-labile protein. IGF-1 is also produced locally by cartilage progenitor cells in the growth plates of the long bones and by other cells that are targets for GH and acts locally in a paracrine or autocrine manner to stimulate cell division and differentiation. This local production and action of IGF-1 appears to be sufficient to account for the stimulation of growth by GH as selective disruption of IGF-1 expression by the liver, which severely reduces the plasma content of IGF-1, does not interfere with normal growth of mice (12). The synthesis of most of the eight different IGF binding proteins that have been described is either positively or negatively regulated by GH. These proteins appear to act as positive or negative modulators of IGF action.

GH is also an important regulator of energy metabolism in most homeothermic animals. It decreases the sensitivity of liver, muscle, and adipose tissue to the actions of insulin (see [Insulin](#), this volume) and promotes the utilization of fat in preference to carbohydrates to meet the demands of muscle for energy. Consequently, prolonged administration of GH leads to a reduction of carcass fat while promoting the accumulation of body protein through stimulation of protein synthesis. At the same time that insulin sensitivity is decreased, GH paradoxically increases the responsiveness of the insulin-producing cells of the pancreas to physiological signals for insulin secretion. Excess GH often results in temporary or permanent diabetes mellitus. In many species GH secretion is increased during fasting and, perhaps, facilitates the metabolic adaptations to decreased availability of dietary carbohydrates.

### 3. Secretion and Metabolism

GH is synthesized and secreted primarily by cells in the anterior pituitary gland called “somatotropes,” which appear to arise from the same pool of precursor cells that also give rise to the PRL secreting cells (lactotropes). Some cells appear able to synthesize and secrete both GH and PRL and are called “somatomammotropes” or “lactosomatotropes.” Immortalized cells derived from tumors of these cells (eg, GH3 cells) have been widely used as models for studying both synthesis and secretion of peptides. GH is the most abundant hormone produced and stored in the pituitary. Its expression level is about 800 times greater than that of any other pituitary hormone, and it has been estimated to comprise up to 1% of the dry weight of the human pituitary gland. In the living animal, GH is secreted intermittently in pulses spaced about 3 to 4 h apart (13). Positive input for secretion is provided by a 40 to 45 amino acid peptide, the GH releasing hormone (GHRH), which is secreted by neurons in the nearby hypothalamus and reaches the somatotropes through a unique capillary network called the hypophyseal portal system. Negative input for secretion is provided by somatostatin, a 14 amino acid peptide that is synthesized and secreted by other hypothalamic neurons. Pulsatile secretion of GH is thought to result from the periodic coordinated release of GHRH and cessation of somatostatin secretion. Various environmental cues (eg, “stress”) may evoke additional secretory episodes or abolish one or more secretory pulses depending upon circumstance and the species. Products of GH action, such as IGF-1, in the blood diminish pulse amplitudes by stimulating somatostatin secretion in a typical negative feedback arrangement. In addition, GH may also modulate its own secretion by exerting inhibitory actions on GHRH-secreting neurons. Recently, an additional GH secretagogue, named ghrelin, has been identified. This novel 28 amino-acid peptide is esterified with octanoate at serine 3 and appears to act on both the GHRH-secreting cells of the hypothalamus and the somatotropes.

In mammals, which have been studied more extensively than other vertebrate classes, about half of the GH found in the blood is bound to the GH-binding protein (GHBP) whose structure corresponds to the extracellular portion of the GH receptor (GHR). The GHBP is a product of the GHR gene (see below) and is thought to arise in most species by proteolytic cleavage of the GHR at the extracellular surface. The enzyme responsible appears to be the same metalloprotease that liberates the cytokine,



tumor necrosis factor- $\alpha$ , from its transmembrane precursor (14). In rodents, the GHBP is synthesized from an alternatively spliced RNA transcript of the GHR gene. The GHBP increases the apparent biopotency of GH in the intact organism by prolonging its survival time in the blood. When added to isolated tissues or cells, the GHBP reduces bioactivity of GH by competing for hormone with cellular receptors. These effects are not contradictory, as receptor binding leads to complete degradation of the hormone/receptor complex and is likely to serve as the major route of hormone degradation. Some degradation of hormone that reaches the glomerular filtrate also occurs in the kidney. Less than 1% of the GH secreted by the pituitary is excreted intact in the urine.

#### 4. The GH Family of Proteins

GH is a single-chain peptide comprised of 190 to 196 amino acid residues and has a molecular mass of approximately 22,000. Complete amino acid sequence data for hormones from more than 40 different species representing all classes of vertebrates have been determined either by direct analysis of the purified proteins or by deduction from the nucleotide sequences of their cDNAs. Characteristic of all the known GH molecules are the two disulfide bridges that form a long loop by coupling cysteines corresponding to Cys-53 and Cys-165 in the human sequence and a short loop involving Cys-182 and Cys-189. When these cysteine residues are used as reference points and gaps are allowed to maximize alignment of homologous regions, it becomes apparent that more than a third of the amino acid loci in all of the known GH molecules contain identical residues (Fig. 1). Homologies are almost twice as great when conservative substitutions are allowed. Amino acid sequences of GHs from animals as distantly related as humans and sharks have more than 50% identity, although there is considerably more divergence in the GHs within the bony fish. Nevertheless, more than 30 residues in addition to the four cysteines are identical throughout the species, and these appear in clusters with highest conservation in positions between residues 164 and 187 (15). The crystal structure of porcine GH was solved in 1987 (16). In three dimensions, the peptide chain of GH forms into four alpha helices encompassing residues 7 to 34 (helix I), 75 to 87 (helix II), 106 to 127 (helix III), and 152 to 183 (helix IV) (Fig. 2). The alpha helices are tightly packed and arranged in an antiparallel up-up-down-down orientation. The invariant and highly conserved amino acid residues are predominantly found within the alpha helices and presumably contribute to the integrity of the tertiary structure of the molecule. A long loop of 40 amino acids connects helices I and II, and shorter 18 and 24 residue loops connect helices II and III and III and IV, respectively. Short peptide chains are found at the amino and carboxyl termini.

**Figure 1.** The amino acid sequence of hGH. Shaded boxes indicate conserved amino acids in mammalian GH. The shaded lines above or below the peptide chain indicate the amino acid residues present in helices I, II, III, and IV. (From: H. M. Goodman, G. P. Frick, and S. Souza (1996) Species specificity of the primate growth hormone receptor. *News Physiol.* 5 11, 157–160).

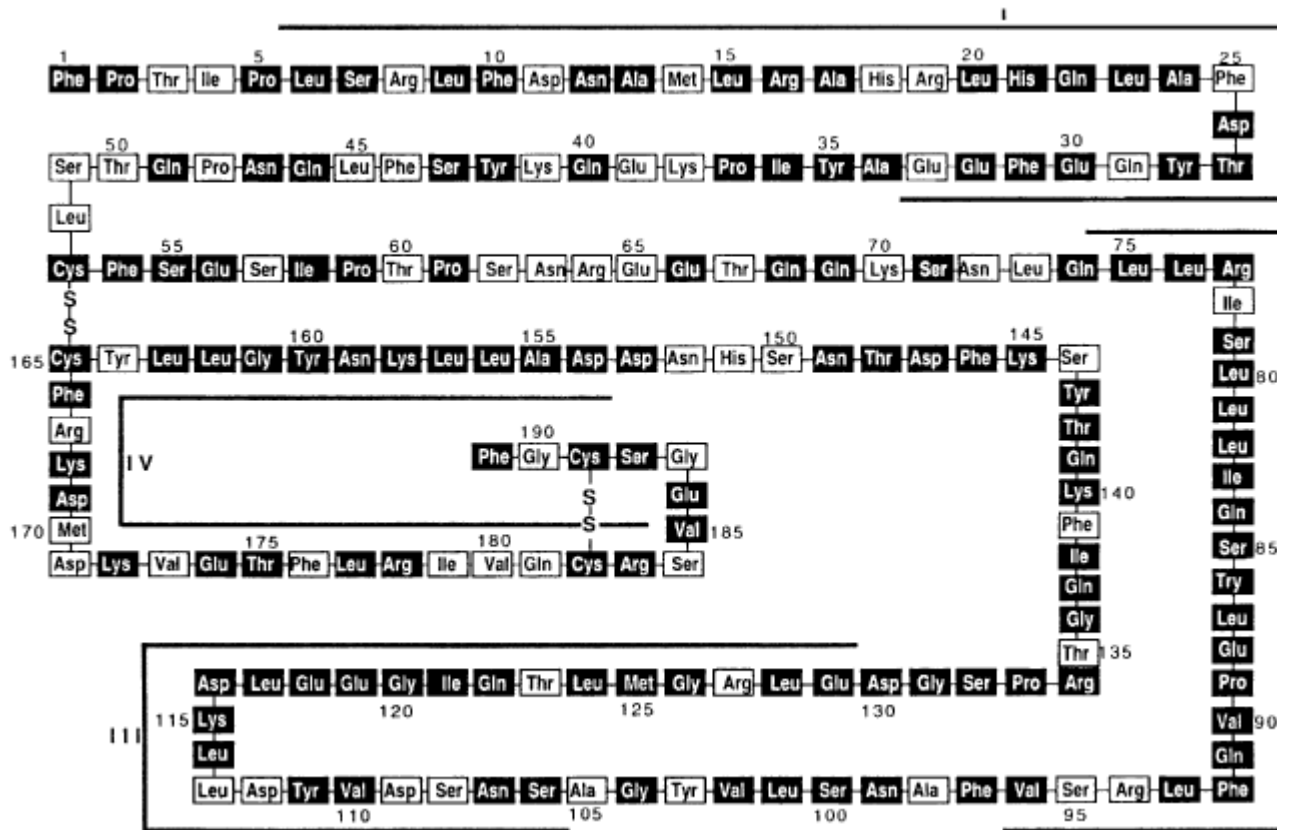
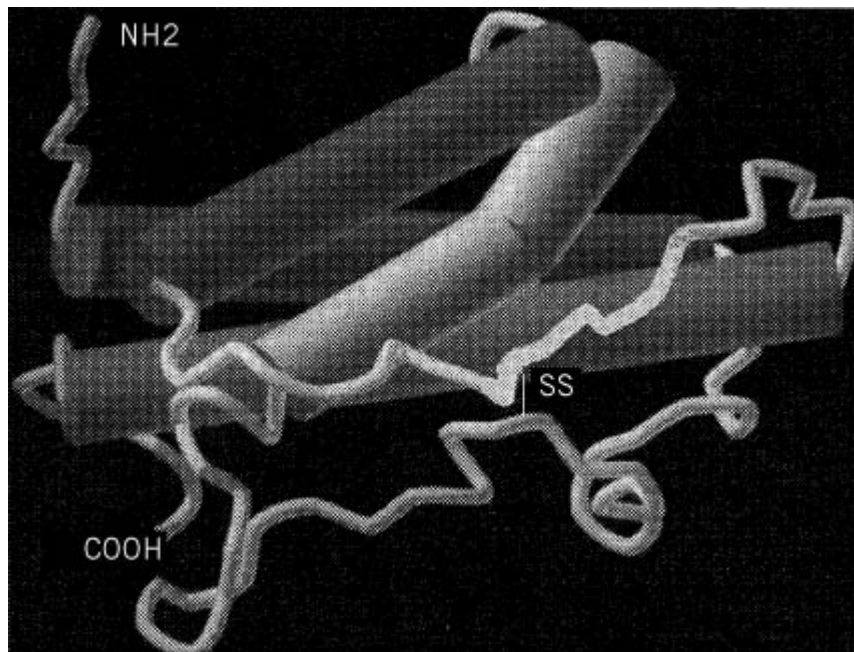


Figure 2. Tertiary structure of GH (modified from ref 14).



Variant forms of GH have been reported in several species. Of these, the most thoroughly studied is the so-called 20 k hGH. This isoform of hGH lacks 15 amino acids (residues 32 to 46) and, consequently, has a molecular mass of 20,000, instead of the 22,000 of the predominant form

discussed above (17). The 20k isoform of hGH constitutes only about 10% of the daily GH secretion and may be somewhat less potent than its 22k counterpart in some assay systems. It arises from alternative splicing of the same RNA transcript that codes for the 22k isoform (see below). Another splice variant, a 17.5 kDa form (18), has also been described, but little is known of its physiology. Processing variants include phosphorylated, glycosylated, amidated, and proteolytically cleaved forms (19). Some of these variants are found in the circulation, but their biological importance has not been established. Some modified forms of GH may be artifacts of pituitary storage or purification procedures.

As already indicated, another anterior pituitary hormone, prolactin (PRL), is closely related to GH in both its chemical and biological properties. PRL is found in all vertebrates and has a wide range of biological actions among the various species, including regulation of salt and water balance, growth, and various behaviors and functions associated with reproduction and nurturing of the young (20). PRL is named for its critical role in supporting lactation in mammals. Like GH, PRL consists of nearly 200 amino acid residues and has the same four antiparallel helix bundle tertiary structure as GH. PRL has three intrachain disulfide bridges, two of which correspond to those found in GH, and the third is located near the N terminus. PRL is about 20% identical in amino acid sequence to GH of the same species, a finding that suggested the hypothesis that the two hormones originated from a common ancestral gene (21). *A propos* of this idea, a glycoprotein with a mass of about 26 kDa that appears to be equally related to both hormones has been isolated from pituitary glands of some bony fish (22). It consists of 206 amino acids with about 24% sequence identity with both GH and PRL, and conservation of the two disulfide bridges corresponding to those in GH. Its function has not been established, but it has been named “somatolactin” because of its similarity to both GH and PRL. Somatolactin may be descended from the protein that was ancestral to both GH and PRL and, as of this writing, has not been found in tetrapod or avian pituitary glands.

Additional members of the GH family are expressed in the placenta and, hence, are confined to the mammals. The human placenta secretes a variant form of GH (GH-V), which appears to be the major form of GH present in maternal blood plasma during late pregnancy (23). GH-V differs from pituitary GH (hGH-N, for normal) in 13 of the 191 amino acids in positions that are widely distributed between residues 18 and 149 and, unlike the pituitary hormone, is glycosidated at Asn-140, which is unique to GH-V (24). It has a molecular mass of about 23 kDa. Blood plasma of pregnant women also contains a less abundant 25 kDa GH-like peptide of placental origin, which appears to arise from alternative splicing of the RNA transcript of the hGH-V gene (25) (see below). The human placenta also produces placental lactogen (hPL, more commonly called human chorionic somatomammotropin, hCS), which has both GH-like and PRL-like activities in many assay systems. It is the most abundant peptide produced by the placenta, which, in late pregnancy, may secrete as much as 1g per day. Human CS differs from the pituitary form of hGH in only 28 of the 191 amino acid residues and from hGH-V in only 40 residues, with the carboxyl terminus again showing the greatest degree of identity (24). PLs are also produced by the placentae of other species, notably ruminants and rodents, but apparently are not universally found among the mammals. Four PL and four PRL-related proteins are produced in rodent placentae, and the bovine placenta produces one PL and six PRL-related proteins. The PLs bind to GH or PRL receptors and, at least in cattle, to PL receptors, but, currently, no receptors for PRL-related proteins have been characterized. Bovine and rodent PLs and PRL-related proteins are more similar to PRL than to GH, while the placental lactogen of primates is more closely related to GH.

## 5. The GH Gene Family

Determination of the structures of the genes that encode the GH family of proteins reveals an even greater similarity and relationship among these proteins than the amino acid sequences. The structure of the GH gene has been determined in mammalian, avian, and piscine genomic libraries. The related PRL, PL, and somatolactin genes have also been sequenced in a variety of species. The mammalian GH genes span about 2 kb and are the smallest in the family, while the somatolactin gene is the largest and extends for about 16 kb. the mammalian PRL genes encompass about 10 kb, and the PL

genes resemble those of either GH or PRL, depending upon the species. The genes that encode all of these proteins contain five exons and four introns, and, despite wide variations in the sizes of the introns, the splice junction patterns and the positions of introns with respect to the coding sequences of the protein products have been remarkably conserved. In some suborders of bony fish, eg, the salmonids, however, the fifth exon has been split by insertion of a fifth intron. With the conservation of exon/intron boundaries, the sizes of the first four exons are quite similar, ~0.1, 0.2, 0.1, and 0.2 kb respectively, but exon V ranges from ~3 kb in somatotactin to ~0.3 kb in hGH, giving rise to wide variation in the length of the untranslated 3' end (15).

The human haploid genome contains five closely-linked GH-related genes at band segment q22–24 on the long arm of chromosome 17 (26, 27) and a single PRL gene on chromosome 6 (28). The hGH-N gene is expressed only in the pituitary and, to a much smaller extent, in some white blood cells (monocytes), while the three hCS genes and hGH-V are expressed only in the placenta. About 67 kb of the hGH gene locus, including intergenic sequences and up- and downstream flanking regions, have been completely sequenced (29). The five gene cluster itself spans about 48 kb (30) and contains the genes for hGH, designated as hGH-N (for normal); two genes for hPL, designated hCS-L (hCS-like) and hCS-A; the gene for hGH-V (for variant); and a third PL gene, hCS-B, arranged in order in the 5' to 3' direction. All five genes are transcriptionally oriented in the same direction and are separated by intergenic regions of 6 kb to 13 kb in length. The intergenic regions contain 48 middle repetitive Alu sequences that comprise more than 20% of the nucleotides in the hGH locus (29). A repeated P element consisting of about 1 kb is also found about 2 kb upstream from the promoters of the four placentally expressed genes. The hCS-A and hCS-B genes encode prehormones that differ by one amino acid in the leader sequence have five silent nucleotide differences, but the mature proteins are identical (31). The hCS-B gene also contains a small insertion in the 3' untranslated region that is not found in hCS-A. The nucleotide sequence of the hCS-L gene shows 93% to 94% homology with the genes for hCS-A and hCS-B, but, because a G to A transition in the splice donor site of the second intron interferes with RNA processing at this site, a complex array of alternatively spliced mRNAs is formed, the majority of which are nonfunctional. One mRNA, however, encodes a 20 kDa protein that may be secreted (32). Expression of the hCS-A and hCS-B genes in the placenta is about 1,000 times higher than that of hGH-V or hCS-L genes (29).

Unlike the primate genome, which contains a single prolactin gene and multiple GH-related genes, the genomes of rodents and ruminants appear to contain a single GH gene and multiple PRL-related genes as judged from their nucleic acid and protein structures. PLs appear to be lacking in other mammalian groups (33). The GH locus and the PRL locus in rodents and ruminants are on separate chromosomes that appear to be homologous with their human counterparts as judged from the colocalization with other known genes. Genes for the numerous PRL-related proteins that are expressed in the placentae of cattle, sheep, rats, and mice map to the same chromosomes as PRL, and, although the entire PRL chromosomal locus has not been sequenced in any of these species, it is likely that they cluster in a fashion analogous to that seen in the human GH locus (34-37).

## 6. Regulation of GH Gene Expression

There are two aspects to regulation of GH gene expression; both are complex and incompletely understood. One aspect concerns the factors that determine tissue specificity and restricts GH expression to the somatotropes. The other aspect concerns the regulation of the levels of expression by the somatotropes. Expression of hGH-N, but not hGH-V, in the somatotropes is thought to be governed by remote regions of 5' flanking DNA tens of thousands of bases upstream from the hGH-N gene. These elements, along with somatotrope-specific regulatory proteins, may reorganize the chromatin such that the transactivation factors discussed below have access to the promoter region of the downstream gene (38). Somewhat different control mechanisms may pertain to nonprimates that express only one gene at the GH locus, but, throughout the vertebrates, expression of proteins of the GH/PRL family depends upon the indispensable regulatory protein, Pit-1, which is also sometimes

called GHF-1. Pit-1 is a ~290 amino acid protein that is expressed only in certain pituitary cells, the placenta and monocytes that produce GH, and, thus provides cell-specific regulation (39, 40). Pit-1 is also crucial for differentiation and growth of somatotropes during development and for the induction of the GHRH receptor and, hence, sensitivity to hypothalamic control (41). As mentioned earlier, GHRH is the hypothalamic hormone that stimulates the somatotrope to synthesize and secrete GH. In addition, Pit-1 also activates its own transcription and, thereby, maintains the differentiated function of somatotropes, lactotropes, and somatomammotropes (42).

Physiologically, GH synthesis and secretion are driven by GHRH, which acts through receptors on the somatotrope surface to increase the production of cyclic AMP, which binds to and activates protein kinase A. Protein kinase A catalyzes the phosphorylation of a variety of proteins on serine residues and, thereby, increases or decreases their activities. Among the substrates for protein kinase A in the somatotrope are nuclear regulatory factors of the cyclic nucleotide response element binding (CREB) protein family. The hypothalamic peptide somatostatin inhibits GH synthesis and secretion by blocking the actions of GHRH on cyclic AMP production and ion permeability. In addition, peripheral hormones, triiodothyronine from the thyroid gland and glucocorticoids from the adrenal cortex, enhance GH synthesis and secretion. Other paracrine (eg, activin) and circulating factors (insulin, retinoic acid, vitamin D) also act on the somatotrope to modulate GH production either positively or negatively, while yet other regulatory influences, (eg, IGF-I, free fatty acids and glucose) regulate events in the somatotrope through their actions on GHRH and somatostatin secretion.

The proximal promoter includes the TATA box sequence, which for all of the vertebrate GH genes is TATAAAA and for the PRL and somatolactin genes is TATAAAG (15). Regulation is largely conferred by critical nucleotide sequences that lie within a few hundred bases upstream of the cap site in mammalian GH genes and are the elements that bind regulatory proteins. The hGH-N gene has two binding elements for Pit-1, at 80 and 122 base pairs 5' to the transcription initiation site. These sites are conserved at similar positions in rodent and other GH genes. The ubiquitous zinc finger protein Zn-15 binds to the so-called Z box, a highly conserved sequence located between the Pit-1 elements, and synergizes with Pit-1 (43). Other ubiquitous regulatory factors that bind to the hGH-N promoter include Sp1, AP2 (activator protein 2), NF-1 (nuclear factor-1), USF (upstream stimulatory factor), and the glucocorticoid receptor (42). There is an additional glucocorticoid binding site in the first intron of the hGH-N gene. Glucocorticoid hormones enhance GH expression by augmenting transcription and may also stabilize GH mRNA. The rat GH gene also contains a silencer element whose binding protein is specifically absent in somatotropes (44). Much less is understood concerning control of the the expression of hGH-V. Repression of hGH-V expression in the pituitary may be achieved through the agency of the P element located upstream from each of the placentally expressed genes of the GH family (29).

Despite the high degree of conservation seen in many aspects of the structure and regulation of the GH gene, some different mechanisms appear to apply for regulation of the human and rat genes. The hGH-N gene has two CREB binding sites (45), but these are lacking in the rat GH gene (46). GHRH may, therefore, directly activate the GH gene in humans, but, in rodents, its influence may be exerted indirectly via cyclic AMP response elements (CREs) in the rat Pit-1 promoter and increased synthesis of Pit-1 (47). The human Pit-1 promoter lacks recognizable CREs. The rat GH promoter, but not the hGH-N promoter contains a general hormone response element that binds to thyroid hormone receptors (48) as well as retinoic acid and vitamin D receptors (49). Accordingly, thyroid hormone enhances transcription of the rat GH gene, but not hGH-N.

## 7. Splice Variants

The preRNA transcripts are processed and spliced to produce mRNA that codes for a prehGH that contains a leader sequence of 26 amino acids that is cotranslationally removed. The first exon/intron boundary is of the class I type, falling between the first and second nucleotides of the glycine codon in hGH-N. The remaining exon/intron boundaries are of the class 0 type and fall between codons.

The translated portion of exon I includes 10 nucleotides: the first three codons and the first nucleotide of codon 4. Exon 2 begins with the second and third nucleotides of codon 4 and contains the codons for the remaining 22 amino acids of the leader sequence and the first 31 residues of the mature protein. Exon 3 codes for residues 32 to 71 of the 22k form of hGH. Exon 4 consists of the codons for residues 72 to 126, and exon 5 codes for the remaining 65 amino acids and the 3' untranslated region. Codon 46 in exon 3 contains the dinucleotide AG, which can act as an alternate splice acceptor site. In about 10% of the transcripts, the 45 nucleotides at the 5' end of exon 3 are removed along with the second intron to produce an mRNA that codes for a shorter peptide chain lacking 15 amino acids (residues 32 to 46) and with a mass of 20 kDa (17). A second, and considerably less abundant, splice variant of the hGH-N gene arises from the complete deletion of exon 3 to give rise to a 151 amino acid form with a mass of 17.5 kDa (18). Other mRNA variants lacking exons 3 and 4, and 2, 3, and 4 have been found by reverse transcription-polymerase chain reaction experiments, but it is not known if they are translated or secreted (50).

Although the hGH-V transcript is identical to the hGH-N transcript in the vicinity of the alternate splice site, no alternative splicing of hGH-V occurs at this site. This appears entirely attributable to three nucleotides present only in the hGH-N transcripts and located between the normal and the alternative splice acceptor sites at positions 17, 18, and 24 bases upstream from the alternative site (51). The transcript of hGH-V, however, is subject to alternative splicing in which the 253 nucleotides of intron 4 are retained. Because the reading frame remains open through intron 4 and into exon 5, but with a +1 frame shift, the resulting hGH-V2 transcript codes for a protein in which the 65 carboxyl terminal amino acids are replaced by a different 104 residue carboxyl tail (24). This alternatively spliced transcript represents about one-third of the hGH-V mRNA, but no hGH-V2 protein is secreted, perhaps because a hydrophobic sequence in the new carboxyl tail causes it to be retained within the synctiotrophoblast. Intron 4 may be similarly retained in alternatively spliced transcripts of hCS-A and hCS-B (52) and in bovine GH (53). Alternative splice variants of other species of GH appear not to have been studied extensively.

## 8. Mechanism of Action

### 8.1. The GH Receptor

The GH Receptor (GHR) was purified from the richest known source, the rabbit liver, by Waters and Friesen (54). It was cloned from a rabbit cDNA library using an oligonucleotide probe whose construction was based on the amino acid sequence of a fragment of the rabbit GHR (55) and identified using monoclonal antibodies raised against the purified protein in the Waters laboratory (56). Human, rodent, ruminant, porcine, and chicken GHRs were cloned soon thereafter, and the mammalian receptors were found to share 70% or greater sequence identity. The ~4 kb cDNA codes for a 620 amino acid mature protein with an 18 residue leader sequence. Hydropathy analysis indicated a single membrane-spanning region of 24 amino acid residues separating the N-terminal 246 residue extracellular domain from the 350 residue cytosolic domain. The extracellular domain contains seven cysteine residues, of which six are paired in S-S linkage, and five potential N-glycosidation sites. The intracellular domain contains no known catalytic activity but associates with a cytosolic protein tyrosine kinase, Janus kinase-2 (JAK2), which appears to bind to a proline-rich region of the receptor adjacent to the membrane-spanning domain.

The human GHR receptor is a single copy gene on chromosome 5. Its coding and nontranslated 3' tail are encoded by nine exons, designated 2 to 10, extending over 87 kb. (57). Exon 2 codes for the signal peptide, 3 through 7 for the extracellular domain, exon 8 for the transmembrane domain, exon 9 the beginning of the cytosolic domain, and the long exon 10 (3,400 bp) codes for the bulk of the intracellular domain and the long 3' untranslated region. At least eight variants of the 5' untranslated region are found in the human GHR gene (58), and multiple 5' regions are present in other GHRs indicative of multiple, perhaps tissue specific, promoters. Rodent GHR genes also contain several variants at the 5' end as well as an additional exon (exon 8A) between exons 7 and 8 (59). Alternative splicing that replaces exons 8, 9, and 10 with exon 8A produces an mRNA that codes for

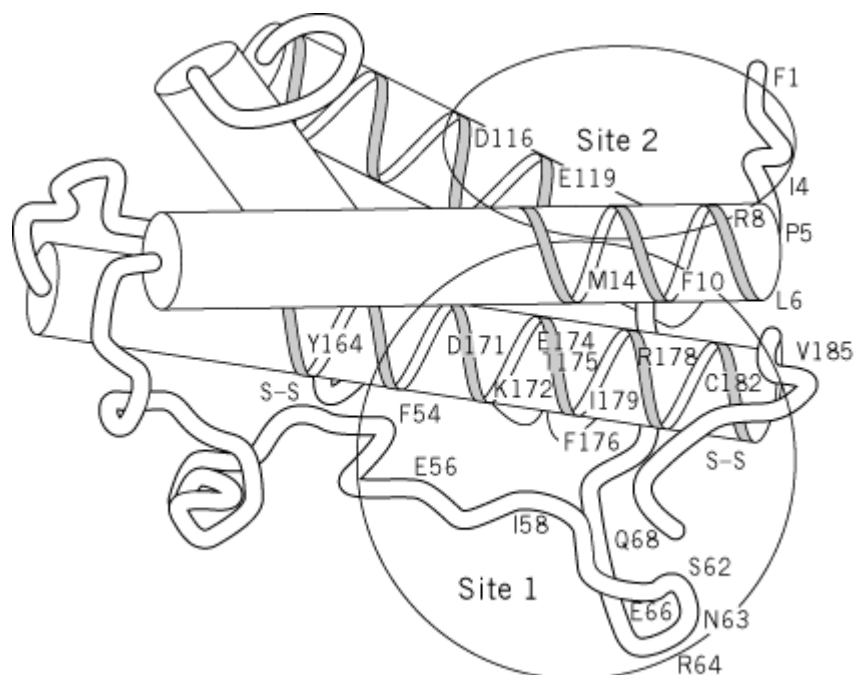
the circulating GH binding protein (GHBP) in rodents (60). The GHBP has a hydrophilic carboxyl tail instead of the transmembrane and cytosolic domains of the GHR.

As with their cognate hormones, the GH and PRL receptors share many similarities in structure (61). The human GH and PRL receptor genes map to the same locus on chromosome 5 (62). The GH and PRL receptors belong to the superfamily of Cytokine/Hematopoietin receptors that includes receptors for erythropoietin, interferon, interleukins 3 and 5, the common subunit of interleukins 2, 4, 7, 9 and 15, the common subunit (gp130) of interleukins 6 and 11, oncostatin, leukemia inhibitory factor, and several others. The relationship is based on both primary and tertiary structural features of the extracellular and cytosolic domains and is also reflected in similarities in the modes of signalling of these receptors. Ligands for these receptors are four helix bundle structures as described for GH.

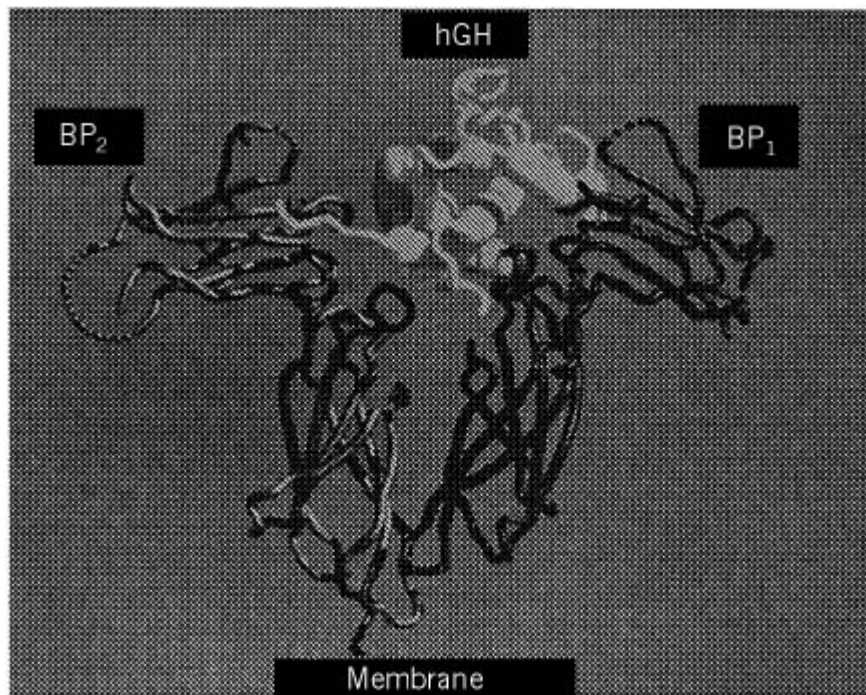
## 8.2. GH-Receptor Binding

In a brilliant series of experiments summarized in Wells and colleagues (63) at Genentech mapped the surfaces of the hGH and receptor molecules that come in contact upon binding and determined the stoichiometry of binding to be two receptors per hGH molecule (Fig. 3). High resolution X-ray crystallography of the complex of hGH bound to the extracellular domain of the GHR (64) revealed that virtually identical binding sites on the two receptor molecules bind to nonidentical sites on the surface of the four helix bundle GH molecule and came in contact with each other in their C-terminal regions in a manner that would bring the transmembrane and cytosolic domains in close apposition (Fig. 4). Their data indicated that GH sequentially binds to one receptor molecule at the so called site 1 on GH and then to a second receptor at site 2 on the hormone, the consequences of which produce a receptor dimer and the initial step in signalling. Binding is ordered due to stronger interaction between hGH and the GHR at site 1 than at site 2. Dimerization of the GH receptor upon ligand binding is prototypical of the behavior of the large superfamily of cytokine receptors and provides a basis for understanding how hormone binding can be transduced into a signal by other tyrosine kinase-dependent receptors as well.

**Figure 3.** Patches on the surface of hGH that bind the hGHR. Key amino acid residues whose mutation to alanine decreases binding are indicated. GH binds to its receptor first at **Site 1** and then at **Site 2**. Note that the orientation of the hGH molecule has been rotated 180° horizontally from that shown in Figure 2 (modified from ref 62).



**Figure 4.** Three-dimensional structure of the GH:GHR complex. Two molecules of the extracellular domain of the GHR, indicated here as **BP1** and **BP2** bind to opposite faces of **hGH** and interact with other in the juxtamembrane region. (from ref. [62](#)).



These studies defined the amino acid residues in the GH molecule that contribute to receptor binding either because they come in direct contact with the surface of the receptor or because they contribute to the tertiary structure in a way that presents the interacting amino acids in the proper orientation. “Designer” GH molecules with greater or lesser affinity for lactogenic receptors ([65](#)) have been prepared as well as GH molecules with greater affinity for the receptor or an antagonist (G120R) with a mutation in the third helix that interferes with binding at site 2 ([66](#)). The region of highest conservation in GH is near the carboxyl terminus, which encompasses amino acids that contribute heavily to site 1. One of the nonconserved amino acids in this region is located at position 171, which in all nonprimate GHs is histidine, but in all primate GHs is aspartic acid. The coordinating residue on the hGHR is arginine in position 43 in place of leucine all other receptors. The incompatibility of His171 in nonprimate hormones with Arg 43 in the human receptor decreases binding affinity 200-fold and is largely responsible for “species specificity” ([67](#)).

### 8.3. Signal Transduction

The primary event in GH signal transduction is the GH-induced dimerization of the GHR, which is thought to bring the cytosolic domains in close apposition resulting in activation and perhaps recruitment of the cytoplasmic tyrosine kinase, Janus kinase-2 (JAK2) and other members of the Janus kinase family ([68](#)). Activation of JAK2 results in phosphorylation of tyrosine residues in the GHR and of JAK2 itself. JAK2 appears to associate with the GHR at the N-terminal third of the cytosolic domain, within which is a critical proline-rich region adjacent to the membrane-spanning domain (box 1) that is conserved within the receptor superfamily. Because a relatively large number of cytosolic and nuclear proteins become tyrosine phosphorylated, it is likely that JAK2 has additional substrates or that other tyrosine kinases are recruited and become activated. Phosphorylated tyrosine residues serve as docking sites for other proteins that contain Src Homology 2 (SH2) domains (see **Src Homology 2**, this volume), thereby recruiting other specific proteins that may be substrates for JAK2 or adapter proteins (e.g., SHC or GRB2) that recruit or activate other



proteins to the complex and initiate signaling pathways. As with other tyrosine kinase-dependent signalling, multiple pathways appear to be activated (69). One of the most prominent of these pathways involves the Stat proteins (for signal transduction and activation of transcription) which, upon activation by tyrosine phosphorylation, form homo- or hetero-dimers with other Stat proteins, migrate to the nucleus, bind DNA, and activate transcription of particular GH-sensitive genes including the early response gene *c-fos* (70) and the hepatic serine protease inhibitor-1 (*Spi-1*) gene (71). Activation of JAK2 also directly or indirectly leads to tyrosine phosphorylation of the insulin receptor substrates, IRS-1 and 2, and to activation of phosphatidylinositol-3-kinase and can, thereby, set in motion various metabolic responses to GH (68). Another prominent signaling cascade, the MAP kinase (for *mitogen activated protein*) pathway is also activated. MAP kinases catalyze the phosphorylation on serine residues of diverse proteins including many other protein kinases, Stat proteins, cytoskeletal proteins, and transcription factors such as Elk1, which, when phosphorylated, bind to the *c-fos* promoter. Stimulation of cells with GH also results in the activation of phospholipase(s), which cleaves membrane phospholipids to release diacylglycerol, which may activate various isoforms of protein kinase C. Protein kinase C catalyzes serine phosphorylation of diverse cellular proteins including DNA regulatory proteins. Finally, GH also causes a rapid increase in intracellular-free calcium (72) by a mechanism that is likely independent of JAK2 as this effect is evident in cells that only express a mutant form of the GHR that cannot bind JAK2 (73). Most of these signalling cascades participate in signal transduction by a variety of hormones, cytokines, and other extracellular signalling molecules, raising the issue of how specific and unique responses to GH are achieved in cells that are also responsive to other agonists. It is likely that the various signalling pathways activated by GH intersect, reinforce, or obtund each other in patterns that are unique to GH and to the particular target cells to produce specific responses to GH.

#### Bibliography

1. H. M. Evans and J. A. Long (1921) *Anat. Rec.* **21**, 62–63.
2. O. Riddle, R. W. Bates, and S. W. Dykeshorn (1932) *Proc. Soc. Exp. Biol. Med.* **29**, 1211–1215.
3. R. W. Bates, T. Laanes, and O. Riddle (1935) *Proc. Soc. Exp. Biol. Med.* **33**, 446–450.
4. C. H. Li and H. M. Evans (1944) *Science* **99**, 183–184.
5. E. Knobil and R. O. Greep (1959) *Rec. Prog. Horm. Res.* **15**, 1–58.
6. S. D. Frasier (1997) *J. Pediatr* **131**, S1–S4.
7. H. D. Niall, M. L. Hogan, G. W. Tregear, G. V. Segre, P. Hwang, and H. G. Friesen (1973) *Rec. Prog. Horm. Res.* **29**, 387–416.
8. P. H. Seeburg, J. Shine, J. A. Martial, R. D. Ivarie, J. A. Morris, A. Ullrich, J. D. Baxter, and H. M. Goodman (1978) *Nature* **276**, 795–798.
9. J. A. Martial, R. A. Hallelwell, J. D. Baxter, and H. M. Goodman (1979) *Science* **205**, 602–607.
10. M. J. Cronin (1997) *J. Pediatr.* **131**, S5–7.
11. R. D. Palmiter, G. Norstedt, R. E. Gelinis, R. E. Hammer, and R. L. Brinster (1983) *Science* **222**, 809–814.
12. J. L. Liu, S. Yakar D. LeRoith (2000) *Endocrinology* **141**, 4436–4441.
13. L. A. Frohman and J. O. Jansson (1986) *Endocr. Rev.* **7**, 223–253.
14. M. Rand-Weaver, H. Kawauchi, and M. Ono (1992) "Evolution of the structure of the growth hormone prolactin family". In *The Endocrinology of Growth, Development, and Metabolism in Vertebrates*, (M. P. Schreibman, C. G. Scanes, and P. K. T. Pang, eds.) Academic Press, San Diego, pp. 13–42.
15. Y. Zhang, J. Jiang, R. A. Black, G. Baumann, and S. J. Frank (2000) *Endocrinology* **141**, 4342–4348.
16. S. S. Abdel-Meguid, H. S. Shieh, W. W. Smith, H. E. Dayringer, B. N. Violand, and L. A. Bente (1987) *Proc. Natl. Acad. Sci. USA* **84**, 6434–6437.
17. F. M. DeNoto, D. D. Moore, and H. M. Goodman (1981) *Nucleic Acids Res.* **9**, 3719–3730.

18. C. M. Lecomte, A. Renard, and J. A. Martial (1987) *Nucleic Acid Res.* **15**, 6331–6348.
19. U. J. Lewis (1984) *Annu. Rev. Physiol.* **46**, 33–42.
20. C. S. Nicoll (1974) "Physiological actions of prolactin". In *Handbook of Physiology, Endocrinology*, Vol. **4**, The Pituitary Gland and Its Neuroendocrine Control (E. Knobil and W. H. Sawyer, eds.) American Physiological Society, Washington, pp. 253–292.
21. H. D. Niall, M. L. Hogan, R. Sayer, I. Y. Rosenblum, and F. C. Greenwood (1971) *Proc. Natl. Acad. Sci. U.S.A.* **68**, 866–870.
22. M. Rand-Weaver, T. Noso, K. Muramoto, and H. Kawauchi (1991) *J. Biochem.* **30**, 1509–1515.
23. G. Hennen, F. Frankenne, J. Closset, F. Gomez, G. Pirens, and N. el Khayat (1985) *Int. J. Fertil.* **30**, 27–33.
24. J. Ray, B. K. Jones, and S. A. Liebhaber (1989) *Endocrinology* **125**, 566–568.
25. N. E. Cooke, J. Ray, J. G. Emery, and S. A. Liebhaber (1988) Two distinct species of human growth hormone-variant mRNA in the human placenta predict the expression of novel growth hormone proteins. *J. Biol. Chem.* **263**, 9001–9006.
26. D. L. George, J. A. Phillips III, U. Francke, and P. H. Seeburg (1981) *Hum. Genet.* **57**, 138–141.
27. D. Owerbach, W. J. Rutter, J. A. Martial, and J. D. Baxter (1980) *Science* **209**, 289–291.
28. D. Owerbach, W. J. Rutter, N. E. Cooke, J. A. Martial, and T. B. Shows (1981) *Science* **212**, 815–816.
29. E. Y. Chen, Y-C. Liao, D. H. Smith, Barrera-Saldaña, R. E. Galinas, and P. H. Seeburg (1989) *Genomics* **4**, 479–497.
30. G. S. Barsh, P. H. Seeburg, and R. E. Galinas (1983) *Nucleic Acids Res.* **11**, 3939–3958.
31. H. A. Barrera-Saldaña, P. H. Seeburg, and G. F. Saunders (1983) *J. Biol. Chem.* **258**, 3787–3793.
32. A. Misra-Press, N. E. Cooke, S. A. Liebhaber (1994) Complex alternative splicing partially inactivate the human chorionic somatomammotropin-like (hCS-L) gene. *J. Biol. Chem.* **269**, 23220–23229.
33. F. Talamantes, L. Ogren, E. Markoff, S. Woodard, and J. Madrid (1980) *Fed. Proc.* **39**, 2582–2587.
34. R. Hediger, S. E. Johnson, W. Barendse, R. D. Drinkwater, S. S. Moore, and J. Hetzel (1990) *Genomics* **8**, 171–174.
35. L. L. Jacson-Grusby, D. Pravtcheva, F. H. Ruddle, D. I. H. Linzer (1988) *Endocrinology* **122**, 2462–2466.
36. A. B. Dietz, M. Georges, D. W. Threadgill, J. E. Womack, and L. A. Schuler (1992) *Genomics* **14**, 137–143.
37. N. E. Cooke, C. Szpirer, and G. Levan (1986) *Endocrinology* **119**, 2451–2454.
38. B. K. Jones, B. R. Monks, S. A. Liebhaber, and N. E. Cooke (1995) *Mol. Cell Biol.* **15**, 7010–7021.
39. C. Lefevre, M. Imagawa, S. Dana, J. Grindlay, M. Bodner, and M. Karin (1987) *EMBO J.* **6**, 971–981.
40. C. Nelson, V. R. Albert, H. P. Elsholtz, I.-W. Lu, and M. Rosenfeld (1988) *Science* **239**, 1400–1405.
41. C. Lin, S.-C. Lin, C.-P. Chang, and M. G. Rosenfeld (1992) *Nature* **390**, 765–768.
42. L. E. Theill and M. Karin (1993) *Endocr. Rev.* **14**, 670–689.
43. S. N. Lipkin, A. M. Näär, K. A. Kalla, R. A. Sack, and M. G. Rosenfeld (1993) *Genes Devel.* **7**, 1674–1687.
44. R. J. Roy, P. Gosselin, M. J. Anzivino, D. D. Moore, and S. L. Guerin, (1992) *Nucleic Acids Res.* **20**, 401–408.

45. A. R. Shepard, W. Zhang, and Eberhard N. L. (1994) *J. Biol. Chem.* **269**, 18094–1814.
46. F. Argenton, S. Bernardini, S. Puttini, L. Colombo, and M. A. T. Bortousi (1996) *Eur. J. Biochem.* **238**, 591–598.
47. A. McCormick, H. Rady, J. Fukushima, and M. Karin, (1991) *Genes Devel.* **4** 1490–1503.
48. A. Sugawara, P. M. Yen, and W. W. Chin, (1994) *Endocrinology* **135**, 1956–1962.
49. P. Garcia-Villalba, Jimenez-Lara, and A. Aranda (1996) *Mol. Cell Biol.* **16**, 318–327.
50. A. Palmethofer, D. Zechner, T. A. Luger, and A. Barta, (1995) *Mol. Cell Endocrinol.* **22**, 225–234.
51. P. Estes, N. E. Cooke, and S. A. Liebhaber, (1990) *J. Biol. Chem.* **265**, 19863–1987.
52. J. N. Macleod, A. K. Lee, S. A. Liebhaber, and N. E. Cooke (1992) *J. Biol. Chem.* **267**, 14219–14226.
53. R. K. Hampson and F. M. Rottman (1987) *Proc. Natl. Acad. Sci. U.S.A.* **84**, 2673–2677.
54. M. J. Waters and H. G. Friesen (1979) *J. Biol. Chem.* **254**, 6815–6825.
55. R. Barnard, P. Bundesen, D. Rylatt, and M. J. Waters (1984) *Endocrinology* **115**, 1805–1813.
56. D. W. Leung, S. A. Spencer, G. Cachianes, R. G. Hammonds, C. Collins, W. J. Henzel, R. Barnard, M. J. Waters, and W. I. Wood (1987) *Nature* **330**, 537–543.
57. P. J. Godowski, D. W. Leung, L. R. Meacham, J. P. Galgani, R. Hellmiss, R. Keret, P. S. Rotwein, J. S. Parks, Z. Laron, W. I. Wood (1989) *Proc. Natl. Acad. Sci. U.S.A.* **86**, 8083–8087.
58. R. I. Pekhletsky, B. K. Chernov, and P. M. Rubtsov (1992) *Mol. Cell Endocrinol.* **90** 103–109.
59. A. Edens, J. N. Southard, and F. Talamantes (1994) *Endocrinology* **135**, 2802–2805.
60. W. C. Smith, J. Kuniyoshi, and F. Talamantes (1989) *Mol. Endocrinol.* **3**, 984–990.
61. P. A. Kelly, J. Djiane, M. C. Postel-Vinay, and M. Ederly (1991) *Endocr. Rev.* **12**, 235–251.
62. K. C. Arden, J. M. Boutin, P. A. Kelly, J. Dijane, and W. K. Cavenee (1990) *Cytogenet. Cell Genet.* **53**, 161–165.
63. J. A. Wells, B. C. Cunningham, G. Fuh, H. B. Lowman, S. H. Bass, M. G. Mulkerrin, M. Ultsch, and A. M. de Vos (1993) *Rec. Prog. Horm. Res.* **48**, 253–275.
64. A. M. de Vos, M. Ultsch, A. A. Kossiakoff (1992) *Science* **255**, 306–312.
65. B. C. Cunningham and J. A. Wells (1991) *Proc Natl Acad Sci U.S.A.* **88**, 3407–3411.
66. G. Fuh, B. C. Cunningham, R. Fukunaga, S. Nagata, D. V. Goeddel, and J. A. Wells (1992) *Science* **256**, 1677–1680.
67. S. C. Souza, G. P. Frick, X. Wang, J. J. Kopchick, R. B. Lobo, and H. M. Goodman (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 959–963.
68. L. S. Argetsinger, G. S. Campbell, X. Yang, A. Witthuhn, O. Silvennoinen, J. N. Ihle and C. Carter-Su (1993) *Cell* **74** 237–244.
69. L. S. Argetsinger and C. Carter-Su (1996) *Physiol. Rev.* **76**, 1089–1107.
70. D. J. Meyer, G. S. Campbell, B. H. Cochran, L. S. Argetsinger, A. C. Lerner, D. S. Finbloom, C. Carter-Su, J. Schwartz, (1994) *J. Biol. Chem.* **269**, 4701–4704.
71. P. L. Bergad, H.-M. Shih, H. C. Towle, S. J. Schwartzberg, and S. A. Berry (1995) *J. Biol. Chem.* **270**, 24903–24910.
72. Y. Schwartz, H. M. Goodman, and H. Yamaguchi (1991) *Proc. Natl. Acad. Sci. U.S.A.* **88**, 6790–6794.
73. N. Billestrup, P. Bouchelouche, G. Allevato, M. Ilondo, and J. H. Nielsen (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 2725–2729.

### **Suggestions for Further Reading**

74. N. E. Cooke and S. A. Liebhaber (1996) Molecular biology of the growth hormone-prolactin gene system. *Vitamins Hormones* **50**, 385–459.

75. V. Goffin, K. T. Shiverick, P. A. Kelly, and J. A. Martial (1996) Sequence-function relationships within the expanding family of prolactin, growth hormone, placental lactogen, and related proteins in mammals. *Endocrinol. Rev.* **17**, 385–410.
76. J. A. Wells and A. M. de Vos (1993) Structure and function of human growth hormone: Implications for the hematopoietins. *Annu. Rev. Biophys. Biomol. Struct.* **22**, 329–351.
77. R. V. Anthony, R. Liang, E. P. Kayl, and S. L. Pratt (1995) The growth hormone/prolactin gene family in ruminant placentae. *J. Reproduct. Fertil. Suppl.* **49**, 83–95.
78. J. L. Kostyo ed. (1999). "The endocrine system: Hormonal control of growth", In *Handbook of Physiology*, Oxford University Press, New York.

## Gtpases

Regulatory GTPases form a large, evolutionarily related [superfamily](#) of [enzymes](#) that catalyze the hydrolysis of GTP to GDP and inorganic orthophosphate. The importance of these [GTP-binding proteins](#) lies, however, in the diverse regulatory activities that they acquire upon GTP binding and that are terminated by GTP hydrolysis. The GTPases thus act as molecular switches, active in the GTP-bound form and inactive in the GDP-bound form. GTP binding and GTP hydrolysis are both tightly regulated by other proteins (see [Guanine Nucleotide Exchange Factors](#) and [RGS Proteins](#)). These proteins thus control a GTPase catalytic cycle and determine the fraction of time that the GTPase will spend in the active states. Alternatively, the regulatory proteins can independently turn the switch on or off in response to appropriate signals.

All the GTPases share a conserved structure around the GTP-binding site, with three identifiable regions of conserved amino-acid sequence. More closely related families share further sequence similarities. The individual families of GTPases are defined by the activities that they regulate and by their conserved structures. Each family is discussed under separate entries.

1. Multiple [initiation factors](#) and [elongation factors](#) required for translation on [ribosomes](#) are regulatory GTPases. Their GTP-triggered functions include the recruitment of aminoacyl-tRNAs to the ribosome, participation in [peptide bond](#) formation and the associated movement of the nascent protein on the ribosome (see [Translation](#)). This structurally and functionally diverse group also includes proteins involved in translocation of nascent proteins across membranes, such as the [signal recognition particle](#). Elongation factor Tu (EF-Tu) from *Escherichia coli* is the historic prototype of the GTPase superfamily. Several proteins involved with co-translational, **signal sequence**-mediated translocation of nascent proteins through membranes are GTP-binding proteins more distantly related to this group.
2. Small, monomeric GTP-binding proteins (SMGs) constitute a large family of structurally related GTPases, most in the 20- to 25-kDa size range. They mediate numerous functions in eukaryotic cells, including activation of protein [kinase](#) signaling cascades, regulation of [cytoskeleton](#) structure, **nuclear import/export**, budding and docking of membrane **vesicles** involved in subcellular organellar trafficking, and so on.
3. **Heterotrimeric GTP-binding** regulatory proteins, often referred to as G proteins, organize and convey signals from cell surface receptors to cellular effector proteins, such as **adenylyl cyclase**, cyclic GMP phosphodiesterase, **ion channels**, [phospholipase C](#), protein kinases, and so on. Individual G proteins (about 20 in total) display distinct but nonunique selectivity among receptors and effectors. Both the GTP-binding subunits (~40 kDa), which include a domain similar to the SMGs, and the tightly associated b and g subunits convey signals to distinct effectors.

4. **Dynamins** are GTP-binding cytoskeletal proteins involved in the budding of endocytic vesicles. They are homologous to the other GTP-binding proteins only in the immediate region of the GTP-binding pocket. Dynamins also interact with the bg subunits of the heterotrimeric G proteins.
5. **Tubulins**, the structural proteins of eukaryotic [microtubules](#), are also GTPases but are only distantly related to the regulatory GTPases. Polymerization of microtubules is regulated in part by GTP binding to tubulin.

## GTP-Binding Proteins

GTP-binding proteins are members of a [superfamily](#) of regulatory proteins whose activity is controlled by binding GTP. They act as conformational molecular switches that regulate numerous and diverse processes throughout the eukaryotes. Many GTP-binding proteins are devoted to [signal transduction](#), mediating information transfer from receptors to intracellular effector proteins. Others regulate such diverse processes as the traffic of [membrane proteins](#) and lipids via cytoplasmic transport **vesicles**, the **cytoskeleton**-mediated control of cell morphology, **nuclear import/export**, and such membrane-reforming processes as [endocytosis](#), [protein secretion](#), and budding. In their GTP-bound, active forms, GTP-binding proteins can stimulate or inhibit the activity of cellular effector proteins, recruit effectors to specific intracellular sites, or sequester effectors and thus block their activity. These proteins thus lie at the convergence of multiple regulatory pathways and function as major integrators of cellular information.

GTP-binding proteins are divided into two major classes, each of which is further divided into individual families. Generally, these structurally determined groupings also indicate related functional specialization and biochemical properties. GTP-binding proteins are ubiquitous throughout the eukaryotes, however, and structurally related proteins may perform strikingly different functions in distantly related organisms.

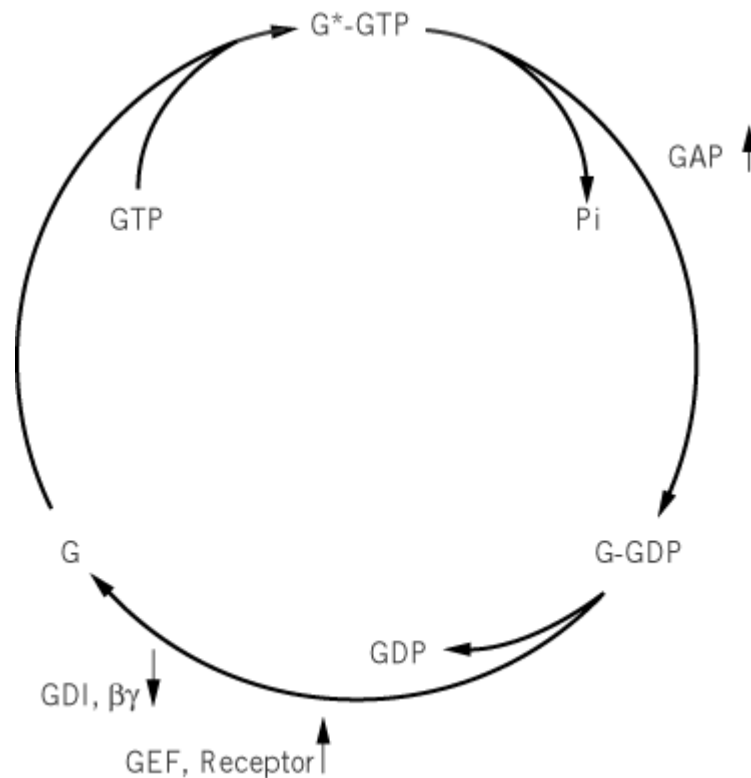
The two major groups of GTP-binding proteins are the [heterotrimeric G proteins](#) and the small, monomeric GTP-binding (SMG) proteins. SMGs are also commonly referred to as G proteins, or small G proteins, and the nomenclature is not well established. Heterotrimeric G proteins are composed of a GTP-binding a subunit (~40 to 45 kDa) and a dimer of b and g regulatory subunits (~35 kDa and 7 to 8 kDa). Heterotrimeric G proteins convey signals from receptors at the cell surface to cellular effector proteins (see [Heterotrimeric G Proteins](#)). The SMG proteins (20 to 25 kDa) are involved in more diverse functions. In addition to these two major classes, other GTP-binding proteins with related primary structures, tertiary structures, and GTPase cycles include the [elongation factors](#) of [translation](#) (with bacterial EF-Tu as prototype) and **dynamin** (see [Gtpases](#)).

Heterotrimeric G proteins are grouped according to the structures and functions of their a subunits and are named with subscripts originally linked to their functions (1, 2). Members of the  $G_s$  family were the first G proteins to be discovered as stimulators (hence “s”) of **adenylyl cyclase**. The  $G_i$ 's were initially defined as cyclase inhibitors (“i”), but they also perform many other functions. The  $G_q$ 's stimulate the phosphatidylinositol bisphosphate-specific **phospholipase C-b**'s. Functions of the  $G_{12/13}$  family include regulation of the **rho** SMG family and of  $Na^+/H^+$  exchangers. The bg subunits are also diverse: There are 5 Gb subunits and at least 11 Gg subunits in mammals. While not all of the 55 conceivable Gbg dimers occur, their potential diversity is huge.

The SMG proteins cluster in five major subfamilies according to their structure and function. Smaller families also exist, and assignment of some SMGs to families remains uncertain. The ras proteins constitute a family of about five proteins in mammals, although gem/kir proteins are also assigned to this group (3-6). Ras proteins were the first SMGs to be identified, because they are active and frequently observed **oncoproteins** when mutated to constitutively active forms or when overexpressed (7). Ras proteins are the prototypical and best-studied examples of the SMG group (see [Oncogenes](#), [Oncoproteins](#)). The rac family (rac, rho, CDC42) are involved in both [signal transduction](#) and in the regulation of cytoskeletal control of cell morphology (8, 9). Rab family members (at least 20) are involved in [protein secretion](#) and [endocytosis](#) (10, 11). ARFs (at least 6) mediate the budding and fusion of intracellular membrane vesicles (12, 13). Ran mediates both inward and outward traffic of macromolecules from the nucleus (14, 15) (see **Nuclear import/export**).

Both heterotrimeric G proteins and SMG proteins act by traversing a tightly controlled cycle of GTP binding and hydrolysis. GTP binding activates G proteins (Fig. 1) (1, 16). Activation is terminated when the bound GTP is hydrolyzed by the G protein's intrinsic [GTPase](#) activity. The GDP-bound form is inactive, as is the unliganded form. Each step in this cycle—GTP binding, GTP hydrolysis, and GDP release—is relatively slow, and each can be regulated both positively and negatively. Activation—GDP dissociation and GTP binding—is promoted by [guanine nucleotide exchange factors](#) (GEFs). GEFs are also referred to as exchange catalysts. They can be either cell-surface receptors that act in response to extracellular signals or cytosolic proteins that are regulated by **allosteric** ligands or **phosphorylation**. Activation is inhibited by GDP dissociation inhibitors (GDIs), which are regulated independently. Deactivation is also accelerated by GTPase-activating proteins (GAPs), providing a third regulatory input. The net output represents the integration of all three inputs.

**Figure 1.** The regulatory GTPase cycle for GTP-binding proteins. GTP-binding proteins, G, are converted by GTP to the active state G\*-GTP in which they can bind and activate effector proteins. Activation is terminated by hydrolysis of the bound GTP to GDP, which is usually a slow process that is accelerated by GTPase-activating proteins (GAPs). The inactive G-GDP complex is relatively stable, unless nucleotide exchange is catalyzed by a guanine nucleotide exchange factor (GEF). In the case of heterotrimeric G proteins, the GEFs are a family of cell-surface receptors that bind ligands on the extracellular face of the plasma membrane. The inactive, GDP-bound state is stabilized by GDP-dissociation inhibitors (GDI), which in the case of the heterotrimeric G proteins are the G $\beta\gamma$  subunits.



The differing rates of each step in the GTPase cycle, along with their various modes of regulation, allow GTP-binding proteins to operate in one of three conceptually different switching modes: toggles, meters, or timers. (1) A GTP-binding protein that traverses the GTPase cycle only once with respect to the process that it regulates behaves as an on/off switch. It is toggled to the active state by GEF-catalyzed activation and maintains its activity until a GAP turns it off. (2) A G protein that traverses the GTPase cycle quickly produces a graded signal in which its fractional activation, averaged over time, is a meter of the relative inputs to activation and deactivation. (3) A G protein can simply allow a process to occur for a fixed period of time, the lifetime of the bound GTP prior to its hydrolysis. While these basic patterns are formally the same, the distinctive kinetic properties of each determine how a G protein will function in the cell. The ability to work in these three different ways are vital to the functional diversity of G-protein action.

### Bibliography

All references are reviews and also serve as suggested reading.

1. E. M. Ross (1989) Signal sorting and amplification through G protein-coupled receptors. *Neuron* **3**, 141–152.
2. A. G. Gilman (1987) G proteins: transducers of receptor-generated signals, *Annu. Rev. Biochem.* **56**, 615–649.
3. H. R. Bourne, D. A. Sanders, and F. McCormick (1991) The GTPase superfamily: conserved structure and molecular mechanism. *Nature* **349**, 117–127.
4. H. R. Bourne, D. A. Sanders, and F. McCormick (1990) The GTPase superfamily I. A conserved switch for diverse cell functions. *Nature* **348**, 125–132.
5. G. Bollag and F. McCormick (1991) Regulators and effectors of *ras* proteins. *Annu. Rev. Cell Biol.* **7**, 601–632.
6. G. M. Bokoch (1996) Interplay between Ras-related and heterotrimeric GTP binding proteins: lifestyles of the big and little. *FASEB J.* **10**, 1290–1295.
7. M. Barbacid (1987) *ras* Genes. *Annu. Rev. Biochem.* **56**, 779–827.

8. L. Van Aelst and C. D'Souza-Schorey (1997) Rho GTPases and signaling networks. *Genes Dev.* **11**, 2295–2322.
9. A. Hall (1998) Rho GTPases and the actin cytoskeleton. *Science* **279**, 509–514.
10. P. Novick and M. Zerial (1997) The diversity of Rab proteins in vesicle transport. *Curr. Opin. Cell Biol.* **9**, 496–504.
11. V. M. Olkkonen and H. Stenmark (1997) Role of Rab GTPases in membrane traffic. *Intl. Rev. Cytol.* **176**, 1–85.
12. M. G. Roth and P. C. Sternweis (1997) The role of lipid signaling in constitutive membrane traffic. *Curr. Opin. Cell Biol.* **9**, 519–526.
13. A. L. Boman and R. A. Kahn (1995) Arf proteins: the membrane traffic police? *Trends Biochem. Sci.* **20**, 147–150.
14. J. M. Avis and P. R. Clarke (1996) Ran, a GTPase involved in nuclear processes: its regulators and effectors. *J. Cell Sci.* **109**, 2423–2427.
15. D. Gorlich and I. W. Mattaj (1996) Nucleocytoplasmic transport. *Science* **271**, 1513–1518.
16. M. S. Boguski and F. McCormick (1993) Proteins regulating *ras* and its relatives. *Nature* **366**, 643–654.

## Guanidination

Guanidination is a method for [chemical modification](#) of [amino groups](#) in proteins. The reaction is relatively slow and proceeds primarily with the ε-amino groups of [lysine](#) residues, converting them to homoarginine, much less with α-amino groups. The modified amino groups retain their positive charge with increased  $pK_a$ . Therefore, structural effects on the modified proteins are usually minimal. The function of the modified protein is usually affected only if the amino group modified intimately participates in the function.

The reaction is followed by [amino acid analysis](#), measuring the appearance of homoarginine and the decrease in lysine. The reaction is quite specific to amino groups, but some side reactions occur with the **thiol** and imidazole groups of [cysteine](#) and [histidine](#) residues, respectively.

### 1. Guanidination with *O*-methylisourea

Almost all ε-amino groups of a protein are guanidinated (**1**) after treatment with 0.5 to 1.0 M *O*-methylisourea [**1**] at pH 10.1 to 11.0 and 0 to 25°C for 2 to 5 days (Scheme 1):

**Scheme 1.** x



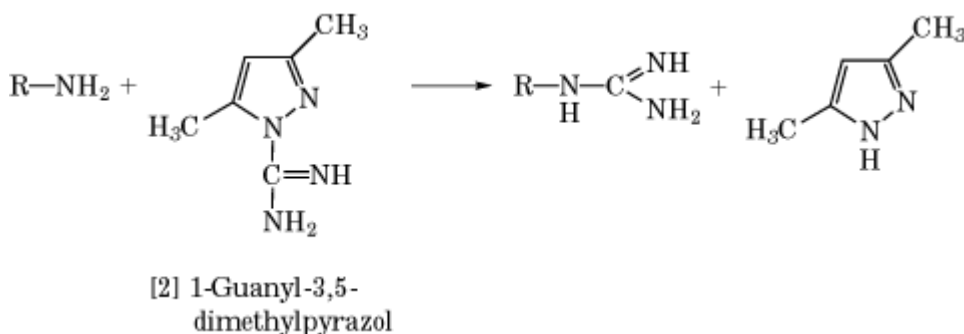
[**1**] *O*-Methylisourea



## 2. Guanidination with 1-guanyl-3,4-dimethylpyrazol

1-Guanyl-3,4-dimethylpyrazol [II] guanidinates proteins under milder conditions than *O*-methylisourea (2). Proteins need to be treated with 0.3 to 0.7 M 1-guanyl-3,4-dimethylpyrazol at pH 9.5 and 0°C for 3–7 days (Scheme 2):

Scheme 2. x



## Bibliography

1. J.R. Kimmel (1967) *Methods Enzymol.* **11**, 584–589.
2. A.F.S.A. Habeeb (1972) *Methods Enzymol.* **25**, 558–566.

## Guanidinium Salts

Guanidinium chloride (GdmCl) is one of the most commonly used **denaturants**. Classically, it is employed at 6M concentration, at which it unfolds most globular [proteins](#) into the form of a [random coil](#). The Gdm<sup>+</sup> ion acts principally by competitively breaking intraprotein [hydrogen bonds](#) and forming hydrogen bonds with protein [peptide bonds](#), although the Gdm<sup>+</sup> ion has also a **hydrophobic** character and can interact in this manner with protein nonpolar residues, principally aromatic groups (1).

In using guanidinium ion denaturation, one has to be careful about the selection of the entire salt, since the anion (X<sup>-</sup>) exerts its own effect and the net result is the summation of the actions of the Gdm<sup>+</sup> ion and X<sup>-</sup>. The effectiveness of Gdm<sup>+</sup> salts in **protein unfolding** follows strictly the [Hofmeister series](#), with the order of denaturing ability being:



GdmSCN and GdmCl are strong denaturants, Gdm acetate is much weaker, and Gdm<sub>2</sub>SO<sub>4</sub> is a [protein stabilizer](#) (1). It is clear that for the SCN<sup>-</sup> salt the effect is due to the summation of two denaturing ions; the strong denaturing capacity of Gdm<sup>+</sup> can overcome the weak stabilization by Cl<sup>-</sup> and CH<sub>3</sub>COO<sup>-</sup>. In contrast, the strongly stabilizing SO<sub>4</sub><sup>2-</sup> ion overcomes the denaturing ability of

Gdm<sup>+</sup>. This progression of the Gdm<sup>+</sup> salts along the Hofmeister series is known to reflect their preferential interactions with proteins: GdmCl exhibits [preferential binding](#); Gdm<sub>2</sub>SO<sub>4</sub> [preferential hydration](#), ie, is excluded from the [accessible surfaces](#) of proteins (2).

## Bibliography

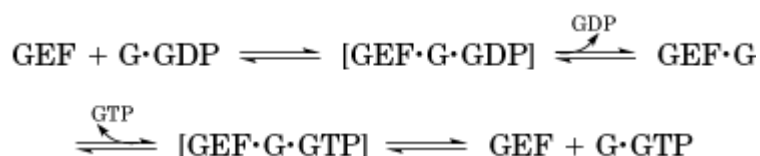
1. P. von Hippel and T. Schleich (1969) In *Structure and Stability of Biological Macromolecules* (S. N. Timasheff and G. D. Fasman, eds.), Marcel Dekker, New York, Chap. "6".
2. T. Arakawa and S. N. Timasheff (1984) *Biochemistry* **23**, 5924–5929.

## Guanine Nucleotide Exchange Factors

Guanine nucleotide exchange factors (GEFs) are a structurally diverse group of proteins that are defined by a common function: catalyzing the exchange of bound GDP for GTP on GTP-binding regulatory proteins (see [Gtpases](#) and [GTP-Binding Proteins](#)). They act on target proteins as diverse as [initiation factors](#) and [elongation factors](#) of [translation](#), the ras-like monomeric GTP-binding proteins, and the [heterotrimeric G proteins](#). Because the GTP-binding target proteins are activated when GTP binds, the GEFs are stimulatory regulators. The GEFs are also referred to as GDP dissociation stimulators (GDSs), GDP release factors (GRFs), and GDP/GTP exchange catalysts.

The prototypical GEF is the bacterial translational elongation factor EF-Ts, which is the GEF for EF-Tu. GEFs related to EF-Ts are found in all organisms. GEFs for the heterotrimeric G proteins are cell-surface **receptors** that initiate G-protein signaling. They are based on a **hydrophobic** core of seven membrane-spanning **a-helices** that binds ligand on the extracellular face and binds G proteins on the cytoplasmic face. GEFs for the small, monomeric GTP-binding regulatory proteins in eukaryotes also fall into classes specific for their targets. While there is no globally identifiable GEF **consensus sequence** or [protein structure](#), GEFs for a particular class of GTP-binding proteins usually share clear sequence [homology](#), although GEFs are frequently multi-**domain** proteins that combine a GEF domain with domains that fulfill other signaling or [protein-protein interaction](#) functions. GEFs may be constitutively active, as are the GEFs for the translation factors. Alternatively, they may be regulated by ligand binding, by **phosphorylation**, by binding of yet other regulatory proteins, or by recruitment to their sites of action when they or their anchors are phosphorylated.

Despite the diversity of their structures and targets, GEFs catalyze GDP/GTP exchange by a common mechanism, shown below (where G is the target GTP-binding protein):



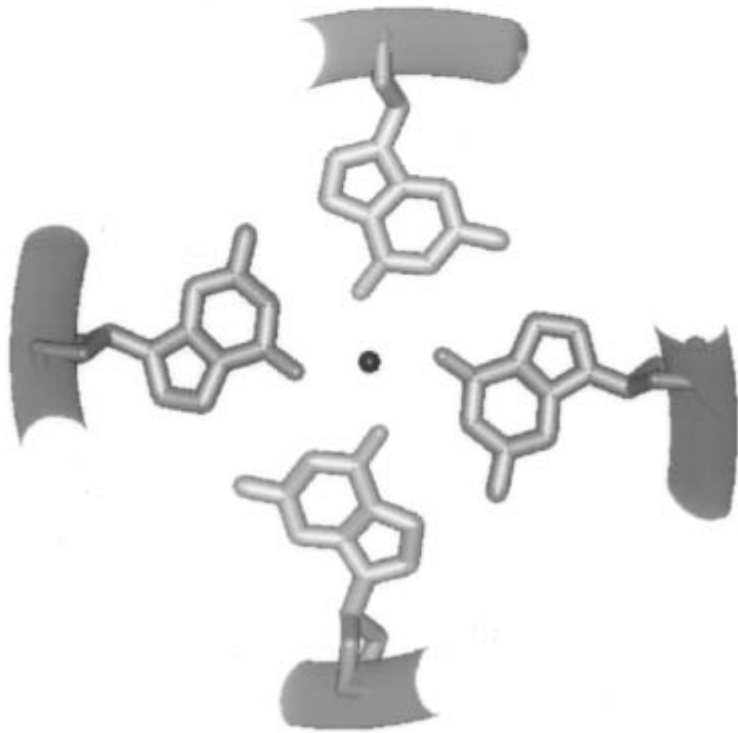
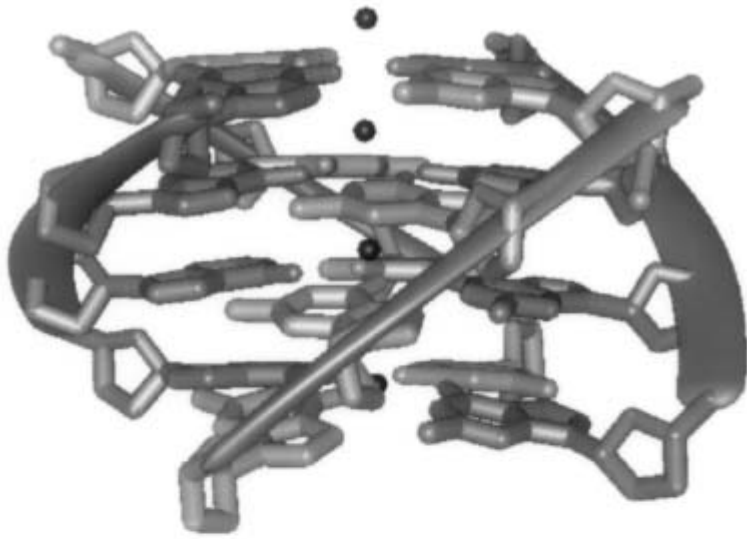
Formally, binding of a GEF to its target protein decreases the affinity of the target for guanine nucleotides: It apparently converts the nucleotide binding site to a more “open” configuration. When the GEF binds, dissociation of GDP is accelerated, leaving the GEF bound to the unliganded target. Binding of guanine nucleotide to the target protein must reciprocally decrease its affinity for the GEF, and cytoplasmic GTP is in large excess over GDP, so GTP binds and displaces GEF by the same mechanism. Release of the GEF after it has promoted GDP/GTP exchange allows a single GEF

molecule to act sequentially on multiple molecules of target protein. Thus GEFs are true **catalysts** of exchange: They accelerate GDP/GTP exchange but do not alter the equilibrium balance of bound GTP and GDP.

## Guanine Quartet

**Telomeric** DNA plays an important role in stabilizing [chromosomes](#). *Telomeres* refer to the terminal protein-DNA complexes of eukaryotic chromosomes, and the DNA sequences found in telomeres are rich in G and contain a short highly-repetitive sequence, such as (TTGGGG)<sub>n</sub>. The G-rich strand at the 3' end exists as a single-stranded overhang and is capable of forming a higher-ordered [DNA structure](#), known as the guanine- or G-quartet, in which four strands associate. There are several ways in which the G-strand can form a very stable four-stranded structure. Indeed, the structures of a number of telomeric sequences have been determined, showing the various forms of the G-quartet structures. The d(GGGGTTTTGGGG) sequence from the protozoa *Oxytricha* forms a hairpin, and two such hairpins joins to make an antiparallel hairpin dimer, thereby resulting in a four-stranded helical structure, with the thymines forming loops at either end. The G-quartet structure was formed by the guanine residues being held together by a cyclic, head-to-tail network of [hydrogen bonds](#) and a metal ion (eg, [TeXnical Error]) located at the center. The four guanine residues were found to have a glycosyl conformation that alternates between *anti* and *syn* (see [DNA Structure](#)). Another arrangement of the tetra-stranded G-quartet is a fourfold symmetric motif in which the four G-strands are in the parallel orientation (Fig. 1) (1). The polymorphism of the quadruplex structure has been reviewed elsewhere (2). More recently, DNA molecules containing a string of isoguanines have been demonstrated to form very stable parallel tetraplex structures as well.

**Figure 1.** Side and end views of a parallel G-quartet DNA structure (Nucleic acid data base No. UDD023).



### Bibliography

1. K. Phillips, Z. Dauter, A. I. H. Murchie, D. M. J. Lilley, and B. Luisi (1997) *J. Mol. Biol.* **273**, 171–182.
2. J. R. Williamson (1993) *Current Opin. Struct. Biol.* **3**, 357–362.

### Guanylate Cyclases

[Cyclic GMP](#) is synthesized from GTP by the enzyme guanylate cyclase (EC 4.6.1.2) and yields pyrophosphate as a product. Cyclic GMP (cGMP) acts as a [second messenger](#), activating cGMP-dependent protein kinases (see **Phosphorylation protein**) and/or regulating cGMP-sensitive gated **ion channels**. The role of cGMP as an intracellular messenger in vascular smooth muscle relaxation and retinal photo-transduction is well established (see [-Monophosphate, cGMP](#)). Garbers, Goeddel and co-workers (1-3) found that the catalytic centers of guanylyl cyclases are strongly related to eukaryotic class III [adenylate cyclases](#). Thus guanylate cyclases form a single class, in contrast to the four classes of adenylate cyclases, but they are of at least two very different types, linked to the function of the enzyme in the cell and either soluble or membrane-bound in eukaryotic cells. These two forms differ in their structure, regulation, biochemical, and physicochemical properties (4-7).

## 1. Membrane-bound Guanylate Cyclases

Guanylate cyclases in family I act as sensors and are often **receptors** for [hormones](#), such as atrial natriuretic peptide (ANP), which is involved in controlling of osmotic pressure and sodium excretion in mammals. These guanylate cyclases are complex, membrane-bound enzymes comprising a receptor for specific hormones coupled to a catalytic domain similar in sequence and structure to the catalytic domain of class III adenylate cyclases (1-3). The known guanylate cyclase receptors form several subfamilies (2, 6,8). Those from sea urchins recognize speract and resact, small peptides that stimulate sperm motility and metabolism ((9),(10)). The receptors for natriuretic peptides (ANF) exist in two forms, both of which synthesize cGMP: guanylate cyclase-A (also named ANP-A), which is specific for ANP, and guanylate cyclase-B (or ANP-B) which is stimulated by brain natriuretic peptide rather than by ANP. There are at least three ANP receptors, two with guanylate cyclase activity (ANP-A and ANP-B) and one (ANP-C) that is responsible for clearing ANP from the circulation without a role in signal transduction (11). Intestinal cells contain the receptor for the *Escherichia coli* heat-stable enterotoxin (guanylate cyclase C). The endogenous ligand for this intestinal receptor is a small peptide called guanylin.

Odorant information is encoded by a series of intracellular [signal transduction](#) events thought to be mediated primarily by the second messenger cAMP. But a subset of olfactory neurons expresses a cGMP-stimulated phosphodiesterase (PDE2) and a guanylate cyclase of the receptor type (guanylate cyclase D), which demonstrate that cGMP has an important regulatory function in olfactory signaling (see [-Monophosphate, cGMP](#)) (8). Finally, retinal guanylate cyclases (often named retGC) exist in at least two forms, retGC-E and retGC-F. They play a specific functional role in the rods and/or cones of photoreceptors and trigger a protein phosphorylation cascade (see [-Monophosphate, cGMP](#)) (12). They consist of an apparent extracellular domain linked by a single **transmembrane region** to an intracellular domain. They are coupled to guanylate cyclase activating protein-2, which is a  $\text{Ca}^{2+}$ -binding protein that activates retGC-1 in a  $\text{Ca}^{2+}$ -sensitive manner (see [Calcium Signaling](#)). It is not known whether retGCs act as receptors, but their structures are similar to that of the other plasma membrane-bound guanylate cyclases.

The organization of all these guanylate cyclases is similar: they have an N-terminal extracellular domain that acts as the **ligand-binding** receptor region, then a transmembrane region, followed by a large cytoplasmic C-terminal region that can be subdivided into two domains, a guanylate cyclase catalytic domain and a protein kinase-like domain that is important for controlling the protein phosphorylation cascade linked to the specific signal recognized by the cognate guanylate cyclase (13, 14).

## 2. Soluble guanylate cyclases

The second family of guanylate cyclases is cytoplasmic and soluble and often designated sGC. They have completely different regulatory functions and always form heterodimers. The two subunits, alpha and beta, are proteins that, although different in length (from 70 to 82 kDa) and sequence, are

highly related (15). Two forms of beta subunits are currently known, beta-1, which is expressed in lung and brain, and beta-2, which is more abundant in kidney and liver. The most fascinating feature of these subunits is that they bind a heme **prosthetic group**. Upon binding of **nitric oxide** (NO) to this heme group, the sGC catalytic activity is stimulated and generates the intracellular signaling molecule cGMP. NO, discovered in 1987, is a signal transduction molecule. Its importance has been emphasized by its role in blood circulation and cardiac muscle functioning (16, 17). Carbon monoxide (CO) also plays a role similar to that of NO (18).

Fifteen conserved **cysteine** residues of sGC have been mutated to serine by *in vitro* **site-directed mutagenesis**. All of the resulting recombinant enzymes synthesize cGMP, demonstrating that these residues are not directly involved in catalysis. On the other hand, mutation of two cysteine residues located in the N-terminal, putative heme-binding region of the beta subunit yields proteins that are insensitive to NO and lose their heme prosthetic group. In contrast, mutation of the corresponding cysteine residues of the alpha subunit does not alter their NO responsiveness, indicating that heme binding is probably a specific feature of the N-terminal domain of the beta subunit (19).

### 3. Similarities Between the Membrane and Soluble Guanylate Cyclases and Class III Adenylate Cyclases

In general, the genetic organizations of both enzyme types have been conserved by the localization of at least one catalytic domain in the carboxy-terminal part of the protein, coupled to a variety (in length and in sequence) of amino-terminal parts. Both forms of guanylate cyclases share a conserved domain that is fundamental for the catalytic activity of the enzyme. A similar domain is also found twice in the different forms of membrane-bound, class III adenylate cyclases from mammals, **slime mold**, or *Drosophila* (see **Adenylate Cyclases**). A polypeptide consensus pattern detects both domains of class-III adenylate kinases and guanylate cyclases: Gly-Val-[Leu/Ile/Val/Met]-X<sub>0,1</sub>-Gly-X<sub>5</sub>-[Phe/Tyr]-X-[Leu/Ile/Val/Met]-[Phe/Tyr/Trp]-[Gly/Ser]-[Asp/Asn/Thr/His/Lys/Trp]-[Asp/Asn/Thr]-[Ile/Val]-[Asp/Asn/Thr/Ala]-X<sub>5</sub>-[Asp/Glu] (1, 3).

The common evolutionary origin of the adenylate cyclases and guanylate cyclases, evidenced by the apparent facility with which it is possible to build up a purine nucleoside triphosphate cyclase of broad specificity, might be relevant to the phylogeny of the catalytic center of class III adenylate cyclases (20). This similarity suggests the existence of a common ancestral purine nucleotide triphosphate cyclase. As a consequence of this interpretation of the sequence data, one might wonder whether evolution has not permitted the existence of some overlap in their specificities because a given cyclase is triggered, under appropriate regulatory conditions, to synthesize either cAMP or cGMP alternatively. This could have been used in cyclic nucleotide-mediated controls existing in eukaryotes, and it might explain the older observations that, in some cases at least, the cAMP and cGMP concentrations vary in opposite directions.

### Bibliography

1. M. S. Chang, D. G. Lowe, M. Lewis, R. Hellmiss, E. Chen, and D. V. Goeddel, (1989) *Nature* **341**, 68–71.
2. D. L. Garbers, (1992) *Cell* **71**, 1–4.
3. O. Barzu and A. Danchin (1994) *Prog. Nucleic Acids Res. Mol. Biol.* **49**, 241–283.
4. D. L. Garbers (1990) *The New Biologist*, **2**, 499–504.
5. D. Koesling, E. Böhme, and G. Schultz, (1991) *FASEB J.* **5**, 2785–2791.
6. S. Schulz, M. Chinkers, and D. L. Garbers, (1989) *FASEB J.* **3**, 2026–2035.
7. S. Schulz, P. S. T. Yuen, and D. L. Garbers, (1991) *Trends Pharm. Sci.* **12**, 116–120.
8. S. Yu, L. Avery, E. Baude, and D. L. Garbers, (1997) *Proc. Natl Acad. Sci. USA* **94**, 3384–3387.
9. D. L. Garbers, J. K. Bentley, L. J. Dangott, C. S. Ramarao, H. Shimomura, N. Suzuki, and D.

- Thorpe, (1986) *Adv. Exp. Med. Biol.* **207**, 315–357.
10. H. Shimomura, L. J. Dangott, and D. L. Garbers, (1986) *J. Biol. Chem.*, **261**, 15778–15782.
  11. M. Chinkers and D. L. Garbers (1991) *Adv. Cyclic Nucleic Acid Res.* **60**, 553–575.
  12. R. B. Yang, D. C. Foster, D. L. Garbers, and H. J. Fulle (1995) *Proc. Natl. Acad. Sci. USA* **92**, 602–606.
  13. J. G. Aparicio and M. L. Applebury (1996) *J. Biol. Chem.* **271**, 27083–27089.
  14. D. Koesling, G. Schultz, and E. Böhme (1991) *FEBS Lett.* **280**, 301–306.
  15. B. Wedel, C. Harteneck, J. Foerster, A. Friebe, G. Schultz, and D. Koesling, (1995) *J. Biol. Chem.* **270**, 24871–24875.
  16. H. Kook, S. E. Lee, Y. H. Baik, S. S. Chung, and J. H. Rhee, (1996) *Life Sciences* **59**, 41–47.
  17. O. W. Griffith and R. G. Kilbourn (1997) *Adv. Enzyme Regul.* **37**, 171–194.
  18. A. Friebe, G. Schultz, and D. Koesling, (1996) *EMBO J.* **15**, 6863–6868.
  19. A. Friebe, B. Wedel, C. Harteneck, J. Foerster, G. Schultz, and D. Koesling, (1997) *Biochemistry* **36**, 1194–1198.
  20. A. Beuve and A. Danchin (1992) *J. Mol. Biol.* **225**, 933–938.

## Guide RNA

Guide RNAs are cofactors for RNA **editing** in kinetoplastid protozoan parasites. They are localized in the [mitochondria](#), where they function to direct insertion and/or deletion of uridylates (U) at specific sites of the pre-[messenger RNA](#) to generate the functional mRNA ([1-5](#)) (see [RNA Editing](#)). The guide RNAs have three domains: (i) the 5' region sequence is complementary to the substrate pre-mRNA and acts as an “anchor.” *In vitro*, this duplex formation is essential for editing to occur ([6](#)). (ii) The central domain of the guide RNA contains the information necessary to insert and/or delete uridylates in the pre-mRNA to make the mature edited sequence, which is normally around 30–40 nucleotides in length. (iii) The 3' end of the guide RNA is characterized by an oligo (U) tail that averages 12 nucleotides in length and is added post-transcriptionally; the function of this tail is unclear. With the development of an *in vitro* system for editing, a critical role in the specificity of editing for the cognate guide RNA has been established ([6](#)).

## Bibliography

1. R. Benne (1992) *Mol. Biol. Rep.* **16**, 217–227.
2. G. J. Arts, P. Sloof, and R. Benne (1995) *Mol. Biochem. Parasitol.* **73**, 211–222.
3. G. J. Arts, H. Vanderspek, D. Speijer, J. Vandenburg, H. Vansteeg, P. Sloof, and R. Benne (1993) *EMBO J.* **12**, 1523–1532.
4. K. Stuart, T. E. Allen, M. L. Kable, and S. Lawson (1997) *Curr. Opin. Chem. Biol.* **1**, 340–346.
5. H. Vanderspek, G. J. Arts, R. R. Zwaal, J. Vandenburg, P. Sloof, and R. Benne (1991) *EMBO J.* **10**, 1217–1224.
6. E. M. Byrne, G. J. Connell, and L. Simpson (1996) *EMBO J.* **15**, 6758–6765.

## Suggestions for Further Reading

7. R. Benne (1996) RNA editing—how a message is changed, *Curr. Opin. Genet. Devel.* **6**, 221–231.
8. M. L. Kable, S. Heidmann, and K. D. Stuart (1997) RNA editing: getting U into RNA, *Trends*

Biochem. Sci. **22**, 162–166.

9. P. Sloof and R. Benne (1997) RNA editing in kinetoplastid parasites: what to do with U, Trends Microbiol. **5**, 189–195

## Gurken

The *gurken* gene of *Drosophila melanogaster* encodes a protein with [homology](#) to the [transforming growth factor](#) alpha (TGF- $\alpha$ ) class of signaling molecules. During *Drosophila* oogenesis, the *gurken* gene is **transcribed** in the cells of the germline. The *gurken* [messenger RNA](#) accumulates in the oocyte and is **translated** throughout oogenesis. The Gurken protein is transported to the oocyte membrane, and activates the *Drosophila* **epidermal growth factor receptor** (Egfr; also abbreviated DER), which is expressed on the adjacent follicle cells that surround the oocyte. The activation of this major transmembrane receptor regulates the expression of a number of genes in the follicle cells, and thus initiates a series of events that lead to the correct patterning of both the anteroposterior and the dorsoventral axis of the egg and embryo.

The *Drosophila* Egfr is a [tyrosine kinase receptor](#) that is used in multiple tissues and at various stages in development to regulate growth, patterning, and [differentiation](#) of various cell types. Several ligands have been found to activate the receptor during development (1). Gurken acts as a ligand that is specific to oogenesis. The Gurken protein consists of an extracellular **domain**, which harbors a single [epidermal growth factor](#) (EGF)-like domain, a **transmembrane** domain, and a short cytoplasmic domain (2). It is **homologous** to the gene *spitz* of *Drosophila*, which acts as ligand of Egfr in the embryo and in [imaginal discs](#).

During *Drosophila* oogenesis, eggs are produced through the cooperation of the three different cell types that make up individual egg chambers; within each egg chamber, an individual oocyte is connected at its anterior to a group of fifteen nurse cells, and this cluster of 16 cells is surrounded by a follicle cell epithelium (3). Early in oogenesis, *gurken* mRNA is found in the developing oocyte, and the Gurken protein accumulates in the oocyte membrane (4). At this stage, the oocyte occupies only a small part of the volume of the egg chamber, and therefore only a limited number of follicle cells at the posterior of the egg chamber are in contact with the oocyte membrane. In this group of follicle cells, the Egfr is activated by Gurken protein, which results in the specification of these follicle cells as posterior cells. In flies that are mutant for severe loss-of-function **alleles** of *gurken*, Gurken protein is absent, and the follicle cells at the posterior end of the egg chamber are not induced to assume a posterior cell fate. They consequently develop as anterior follicle cells, which results in the production of an egg with two anterior ends. The posterior follicle cells normally send a signal back to the oocyte, which organizes the [cytoskeleton](#) of the oocyte and specifies the posterior end of the oocyte. In the mutant egg chambers, no such signal is sent from the follicle cells, and the cytoskeleton of the oocyte is misorganized. RNAs such as [bicoid](#) or *oskar*, which should normally be localized to one end of the egg **RNA localization** are mislocalized in these mutant oocytes, and the embryos that develop inside such eggs have an abnormal antero-posterior pattern (5, 6).

In midoogenesis, the pattern of *gurken* RNA accumulation changes dramatically. At this stage, the oocyte has grown and now occupies about one-third to one-half of the egg chamber. The oocyte nucleus has moved from its initial, symmetric position to one side of the oocyte. The *gurken* RNA accumulates in the region around the oocyte nucleus, and the Gurken protein is now found in a very restricted part of the oocyte membrane, directly overlying the oocyte nucleus (2, 4). At this stage, Gurken activates the Egfr in the lateral follicle cells that contact the oocyte on the side where the



nucleus is situated. Activation of the *Egfr* in these lateral follicle cells induces them to become dorsal follicle cells. In the absence of *gurken* signaling, the lateral follicle cells develop into ventral follicle cells (7). The ventral follicle cells normally regulate the production of a ventral signal that activates the **Toll** receptor protein on the ventral side of the egg and is responsible for inducing ventral cell fates in the developing embryo. In the strong *gurken* mutants, therefore, the ventralization of the follicle cell epithelium leads to an overproduction of the ventral signal, and consequently to a ventralized embryo (8).

The activation of *Egfr* in the follicle cells by the Gurken protein therefore has two different functions. In the early egg chamber, activation of the receptor in follicle cells situated at the termini of the egg chamber results in the induction of posterior cell fates. In midoogenesis, activation in the follicle cells on the lateral side of the egg chamber induces dorsal cell fates. The difference in response to receptor activation reflects a pre patterning of the follicle cells along the antero-posterior axis, which, in part, depends on the prior activity of the [Notch signaling](#) system (9, 10).

In the establishment of dorsal follicle cell fates, *gurken* signaling presumably leads to the formation of a broad field of dorsal cell fates, but secondary patterning mechanisms appear to operate to lead to the final, complex pattern of cell fates of the mature egg. In addition, *gurken* signaling also interacts with signaling through [decapentaplegic](#) (*dpp*), a molecule with homology to TGF- $\beta$ . In follicle cells that receive both the *gurken* signal and the *dpp* signal, formation of dorsal appendages is repressed, and formation of operculum cell fate is induced (11, 12). The regulation of the embryonic ventral signal by activation of *Egfr* seems, however, to be independent of the production of dorsal anterior follicle cell fates. In situations where ectopic activation of *Egfr* in follicle cells is induced in parts of the follicle cell epithelium, embryos result that show only regional dorsalization, corresponding to the region of the follicle cell epithelium where *Egfr* was ectopically activated (13).

## Bibliography

1. R. Schweitzer and B.-Z. Shilo (1997) *Trends Gen.* **13**, 191–196.
2. F. S. Neuman-Silberberg and T. Schüpbach (1993) *Cell* **75**, 165–174.
3. A. C. Spradling (1993) "Developmental genetics of oogenesis", In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–70.
4. F. S. Neuman-Silberberg and T. Schüpbach (1996) *Mech. Dev* **59**, 105–113.
5. A. Gonzales-Reyes, H. Elliott, and D. St Johnston (1995) *Nature* **375**, 654–658.
6. S. Roth, F. S. Neuman-Silberberg, G. Barcelo, and T. Schüpbach (1995) *Cell* **81**, 967–978.
7. T. Schüpbach (1987) *Cell* **49**, 699–707.
8. S. Roth and T. Schüpbach (1994) *Development* **120**, 2245–2257.
9. H. Ruohola, K. A. Bremer, D. Baker, J. R. Sedlow, L. Y. Jan, and Y. N. Jan (1991) *Cell* **66**, 433–449.
10. A. Gonzales-Reyes and D. St Johnston (1998) *Development* **125**, 2837–2846.
11. V. Twombly, R. K. Blackman, H. Jin, J. B. Graff, R. W. Padgett, and W. M. Gelbart (1996) *Development* **122**, 1555–1565.
12. W. M. Deng and M. Bownes (1997) *Development* **124**, 4639–4647.
13. A. M. Queenan, A. Ghabrial, and T. Schüpbach (1997) *Development* **124**, 3871–3880.

## Suggestions for Further Reading

14. D. Morisato and K. Anderson (1995) Signaling pathways that establish the dorsal-ventral pattern of the *Drosophila* embryo, *Ann. Rev. Gen.* **29**, 371–399.
15. R. P. Ray and T. Schüpbach (1996) Intercellular signaling and the polarization of body axes during *Drosophila* oogenesis, *Genes Dev.* **10**, 1711–1723.
16. L. A. Nilson and T. Schüpbach (1998) EGF receptor signaling in *Drosophila* oogenesis, *Curr.*

## Haldane Relationship

This relationship of [enzyme](#) catalysis was first recognized by J. B. S. Haldane in 1930 (1) and is generally known as a Haldane. It relates the kinetic parameters for an enzyme-catalyzed reaction to the equilibrium constant  $K_{eq}$  for that reaction. As  $K_{eq}$  is a thermodynamic measurement that cannot be altered by an enzyme, the ratio of the rates in the two directions also must not change in the presence of the enzyme; consequently, a Haldane expresses the relationship between the kinetic and thermodynamic properties of an enzymic reaction. Although there are both kinetic and thermodynamic Haldanes (2), the thermodynamic Haldanes do not apply to all kinetic mechanisms and hence this discussion will be restricted to kinetic Haldanes.

For a simple Uni–Uni reaction,  $A \rightleftharpoons P$ , that is catalyzed by a mutase or [isomerase](#) (see **Kinetic mechanisms**), at concentrations of reactants well below the values of their [Km \(Michaelis constant\)](#), the rates of the forward and reverse reactions are given by the expressions  $(V_1/K_a) A$  and  $(V_2/K_p) P$ , respectively (see [Michaelis–Menten Kinetics](#)). Here,  $A$  and  $P$  represent the concentrations of the reactants, which have  $K_m$  values of  $K_a$  and  $K_p$ , respectively.  $V_1$  and  $V_2$  represent maximum velocities for the forward and reverse reactions, respectively. At equilibrium, the rates of the forward and reverse reactions must be equal, and hence it follows that the relationships given in Equation 1 must hold:

$$\frac{P}{A} = K_{eq} = \frac{V_1 K_p}{V_2 K_a} \quad (1)$$

In general terms, a Haldane is the ratio of the apparent rate constants for the reactions in the forward and reverse directions when all the substrate concentrations are relatively low. When there is more than one substrate, the apparent rate constant is given by the  $V/K_m$  value for the last substrate to add to the enzyme, multiplied by the reciprocals of the dissociation constants for the substrates that added previously. The same definition applies for the products that act as substrates for the reaction in the reverse direction. Thus, for an ordered Bi–Bi mechanism, where substrate A binds to the enzyme before B and product P dissociates before Q:

$$K_{eq} = \frac{V_1}{K_b} \frac{1}{K_{i\alpha}} \bigg/ \frac{V_2}{K_p} \frac{1}{K_{i\beta}} = \frac{V_1 K_p K_{i\beta}}{V_2 K_{i\alpha} K_b} \quad (2)$$

For a Bi–Bi ping-pong mechanism, in which substrate A is converted to product P, which dissociates leaving a modified enzyme, before substrate B is bound and converted to product Q, there is a Haldane for each half-reaction (see **Kinetic mechanisms** and **Isotope exchange**). When multiplied together, they give the overall equilibrium constant as shown in Equation 3 (2):

(3)

$$K_{eq} = K_{eq1} * K_{eq2} = \frac{V_1/K_a}{V_2/K_p} * \frac{V_1/K_b}{V_2/K_q} = \left(\frac{V_1}{V_2}\right)^2 \frac{K_p K_q}{K_a K_b}$$

Haldane relationships provide a means of checking on the correctness of a proposed kinetic mechanism. The value for  $K_{eq}$ , as calculated by using values determined for the kinetic parameters, must be in agreement with the experimentally determined value. They also have implications with respect to the catalytic efficiency of enzymes. It appears that the wide variation in the maximum velocities for enzyme-catalyzed reactions is due to the differences in the  $K_{eq}$  values for the reactions, as well as to differences in the physiological concentrations of the substrates, which are usually in the region of their  $K_m$  values (3).

### Bibliography

1. J. B. S. Haldane (1930) *Enzymes*, Longman, Green & Co. Ltd, London, p. 80. (Reprinted in 1965 by MIT Press, Cambridge, Mass.)
2. W. W. Cleland (1982) *Meth. Enzymol.* **87**, 366–369.
3. W. W. Cleland (1975) *Acc. Chem. Res.* **8**, 145–151.

### Half-Life, Half-Time

The half-life of a [radioisotope](#) or of a chemical reaction can often be described by the time required for the original reactant to decrease to one-half of its original concentration. This parameter is most appropriate for first-order reactions of the type  $S \rightarrow P$  in which the velocity is proportional only to the concentration of  $[S]$ :

$$\text{Velocity} = \frac{-d[S]}{dt} = k_{app}[S] \quad (1)$$

The decay of radioisotopes follows this relationship, as do all unimolecular reactions. Even when there is more than one reactant, the reaction occurs with pseudo-first-order kinetics if the other reactants are present in excess and their concentrations do not alter substantially during the course of the reaction (see [Kinetics](#)). With first-order reactions, the velocity decreases exponentially with time, as the value of  $[S]$  decreases:

$$\ln_e \frac{[S]}{[S]_0 - [S]} = 2.303 \log_{10} \frac{[S]}{[S]_0 - [S]} = k_{app}t \quad (2)$$

where  $t$  is time,  $k_{app}$  is the apparent rate constant, and  $[S]_0$  is the starting concentration of S. When  $[S]$  has reached half the value of  $[S]_0$  at the half-life  $t_{1/2}$ :

$$\ln_e 2 = k_{app}t_{1/2} \quad (3)$$

$$t_{1/2} = \frac{\ln_e 2}{k_{app}} = \frac{0.693}{k_{app}} \quad (4)$$

When known, the half-time can be used to calculate the apparent rate constant for the reaction.

The half-life can also be used to determine the final extent of a reaction. Because the velocity of a first-order reaction decreases exponentially with time (Eq. (2)), it takes a very long time to reach completion or equilibrium. For example, after one to eight half-lives, the extent of a reaction has reached respectively 50, 75, 87.5, 93.75, 96.83, 98.44, 99.23, and 99.61% of completion. After following a reaction for several half-lives, the final extent of reaction can be predicted.

#### Suggestion for Further Reading

W. P. Jencks (1969) *Catalysis in Chemistry and Enzymology*, McGraw–Hill, New York, pp. 557–564.

## Haploid

A *haploid* cell has one complete set of [chromosomes](#) characteristic of that species, except for the sex chromosomes (see [X-chromosome](#) and [Y-Chromosome](#)). It has a [ploidy](#) of one. In sexual species, only the **gamete** cells are usually haploid. The normal cell is **diploid** and has complete sets of chromosomes.

## Haptens

Haptens are small molecules that can interact with specific **antibodies** but cannot induce an immune response by themselves. They can be rendered **immunogenic** by coupling to a carrier protein that will provide the necessary substrate for **antigen processing and presentation**, whereas the hapten part will contribute the [epitopes](#). The hapten model was put forward and studied extensively by Landsteiner in the 1930s. Many haptens were variations on a cyclic organic compound, especially the phenyl ring that provided a convenient tool to play with a large variety of groups that could be replaced or transposed from one position to another. Among the most popular were 2,4-dinitrophenol and phenylarsonate. Haptens are constructed so to have an active group that can form a covalent interaction with a side chain of the protein carrier, most often a [tyrosine](#) residue (for diazotation) or [lysine](#) (formation of a covalent bond to the **ε-amino group**). The essence of the elegant work of Landsteiner was to show clearly the exquisite specificity of antibodies, because simply replacing one group by another on the phenyl ring (for instance a NH<sub>2</sub> group by a carbonyl) or changing the position of one group (eg, NH<sub>2</sub> from ortho to meta or para) was sufficient to abolish or at least severely diminish recognition by the antibody. The model also provided a nice structural basis to study cross-reactions, by systematic permutations of groups on the hapten itself or by comparing the antigenicity of the same hapten conjugated to different carriers. This approach was used much later to analyze the so-called “carrier effect” which was a first step toward elucidation of the [major histocompatibility complex](#) (MHC) restriction and understanding of the mechanisms of processing

and presentation. The hapten model had, however, one major drawback on the impact of the selective theory first proposed by Ehrlich in 1901. The essence of Ehrlich's theory was that introduction of antigen selected antibodies that were supposed to preexist at the lymphocyte surface (which was demonstrated to be true some 60 years later!). As a result of the Landsteiner approach, for about 20 years, in the 1940s and 1950s, it seemed implausible that the immune system might have generated preexisting antibodies to artificial molecules that were created by chemists. This conclusion, however, simply did not take into account the enormous size of the antibody [repertoire](#) or the relative degeneracy of antibody recognition, as finally realized later and revisited by the [clonal selection theory](#).

Immunological detection of small molecules of biological interest, such as [hormones](#), neuromediators, or diverse ligands of small size, faces the problem of their immunogenicity and requires that these molecules be first conjugated onto a carrier protein, and there are many examples of this in the recent literature. Once the appropriate antibodies are available, detection of these molecules in a biological sample may be performed by a great variety of inhibition techniques, such as the [radioimmunoassay](#) (RIA) or the [enzyme-linked immunosorbent assay \(ELISA\)](#). Conversely, it may happen that small molecules behave spontaneously as immunogens, possibly inducing an undesired hypersensitization. This is the case with some reactive compounds currently used in laboratories, such as **iodoacetamide**, a commonly used alkylating agent that can induce delayed-type hypersensitivity with severe rashes. It is thought to react with cysteine [thiol groups](#) of proteins that then behave as carrier and render the small molecule immunogenic. A similar phenomenon is observed with b-lactam antibiotics, of which [penicillin](#) is the prototype, occasionally leading to the possible occurrence of anaphylactic reactions.

See also entries [Antigen](#), [Epitope](#), and [Immunogen](#).

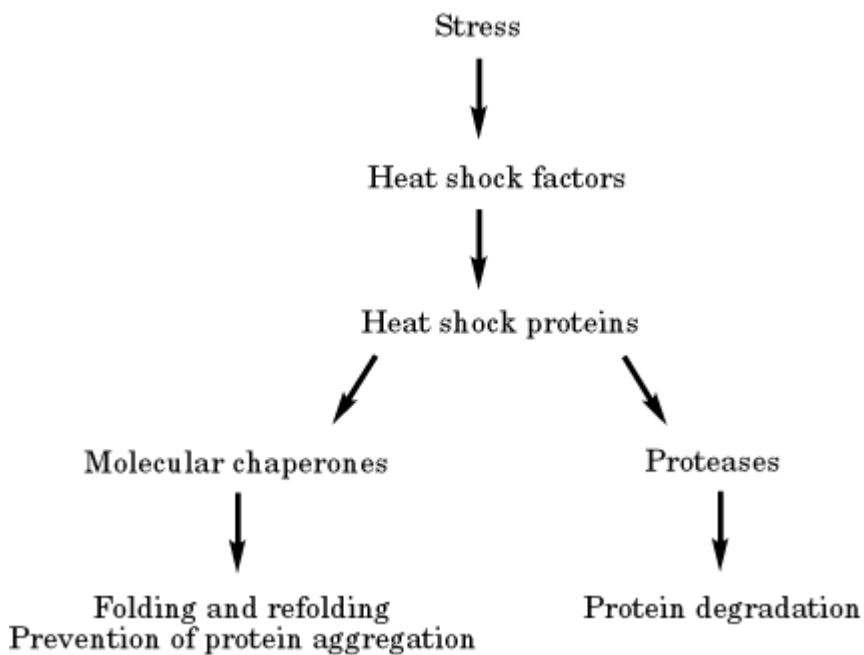
#### Suggestion for Further Reading

O. Mäkelä, I. J. T. Seppälä (1986) "Haptens and carriers". In *Handbook of Experimental Immunology*, Vol. 1: Immunochemistry (D. M. Weir, L. A. Herzenberg, C. Blackwell, L. A. Herzenberg, eds.), Blackwell, Oxford, U.K., pp. 1–13.

## Heat-Shock Response

The heat-shock response represents the orderly cellular response to diverse forms of environmental and physiological stress, which results in the **transcriptional** activation of **genes** encoding **molecular chaperones**, [proteinases](#), and other inducible genes for protection and recovery from stress (Fig. 1). Although the exposure of cells and organisms to heat shock may have its origins with the stress of the primordial environment, the heat-shock response has since adapted to include a diverse array of stresses, including infection with **viruses** and bacteria, exposure to transition heavy metals, amino acid analogues, pharmacologically active small molecules, and oxidants. Common to these stresses are the effects on protein biogenesis defined by events associated with **protein biosynthesis**, **protein folding**, translocation, and assembly into the final native protein or protein complex. Stress challenges protein homeostasis and results in an increased flux of non-native proteins, which, if left unprotected, are prone to misfolding and aggregation. Consequently, through the elevated synthesis of molecular chaperones and proteinases, the heat-shock response responds rapidly and precisely to the intensity and duration of specific environmental and physiological stress signals by repairing protein damage to reestablish protein homeostasis.

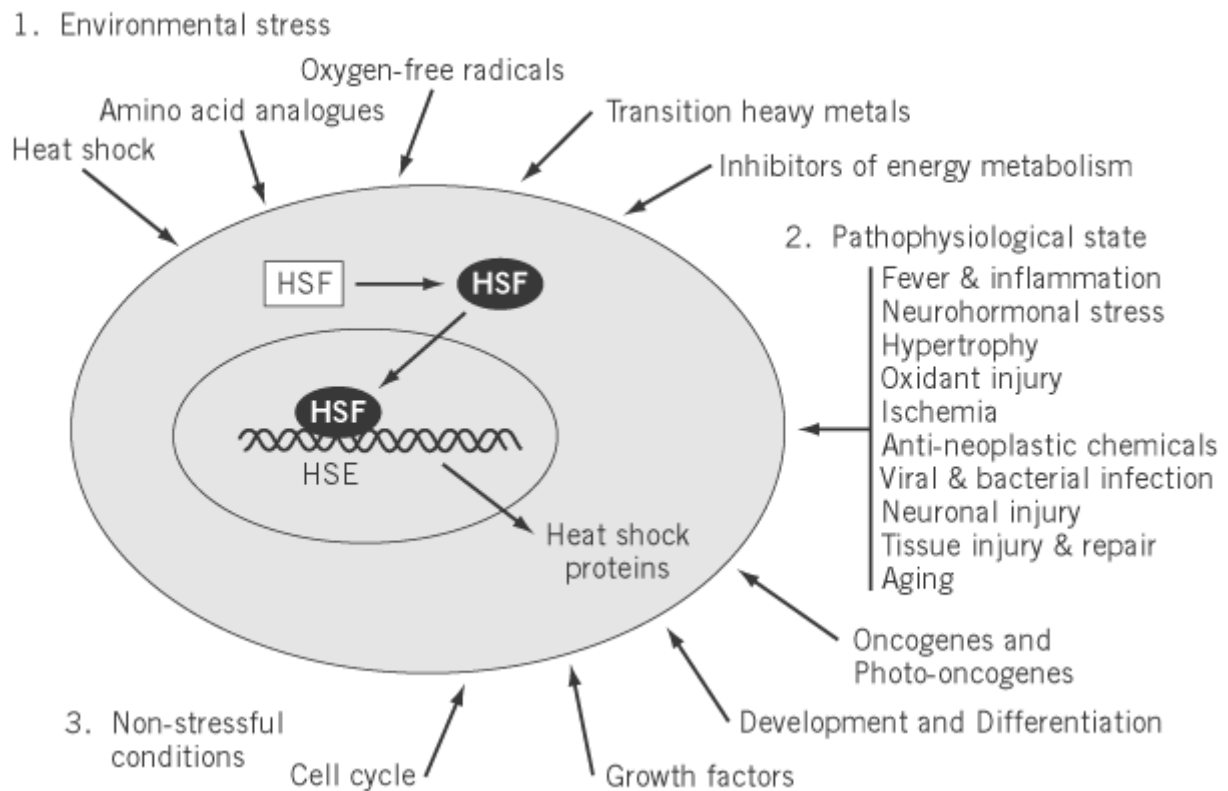
**Figure 1.** Stress activation of heat-shock proteins.



Studies on the heat-shock response have revealed how the cell senses stress and the role of heat-shock proteins in protecting against and repairing protein damage. While prolonged exposures to extreme stress are harmful and can result in cell and tissue death, activation of the heat-shock response by prior exposure to heat shock results in stress tolerance and cytoprotection against stress-induced molecular damage (1). Transient exposure to somewhat elevated temperatures, or reduced levels of chemical stress, cross-protects cells, tissues, and organisms against sustained, normally lethal, exposures to stress. This reveals a valuable survival strategy that “a little stress is good.”

The complexity that underlies the transcriptional regulation of the heat-shock genes is reflected by the diverse array of stress conditions, which can be distinguished into three major categories (Fig. 2): (i) *environmental stresses*, namely, heat shock, amino acid analogues, drugs, oxidative stress, toxic chemicals, heavy metals, and pharmacologically active small molecules; (ii) *nonstress conditions*, namely, the [cell cycle](#), [growth factors](#), serum stimulation, development, differentiation, and activation by certain **oncogenes**; and (iii) *physiologic stress, pathophysiologic and disease states*, namely, neuroendocrine hormones, tissue injury and repair, fever, inflammation, infection, ischemia and reperfusion, neural degenerative diseases, and cancer.

**Figure 2.** Conditions that induce the heat-shock response. Heat-shock gene expression is induced by environmental and physiological stress, nonstressful conditions including cell growth and development, and pathophysiological states.



This entry addresses (i) how cells utilize the stress-sensing network of the heat-shock response to regulate the transcription of genes that encode heat-shock proteins and molecular chaperones; (ii) the role of molecular chaperones in protein folding, translocation, and protection against the deleterious consequences of misfolded proteins; and (iii) the role of the heat-shock response and heat-shock proteins in cytoprotection against disease.

### 1. Stress-Induced Transcriptional Regulation of Heat-Shock Genes in *Escherichia coli*

The heat-shock response in *E. coli* leads to the elevated synthesis of more than 20 proteins (Table 1). Upon temperature elevation, transcriptional activation of the respective heat-shock genes is dependent upon the **s factor** s32, the product of the rpoH (htpR) gene, which associates with the **RNA polymerase** core enzyme (E) to form the E s32 holoenzyme that binds specifically to **promoters** of most heat-shock genes (2-4). s32 is essential for normal growth; cells lacking s32 are highly restricted in cell growth conditions and are hypersensitive to stress (5). Transcriptional activation of the principal heat-shock genes htpG, **dnaK**, dnaJ, grpE, and groEL/S requires stress-induced changes in the levels of s32. Under normal growth conditions, the concentration of s32 is very low, approximately 10 to 30 copies per cell at 30°C (6) as a result of its short half-life ( $t_{1/2} = 1\text{min}$ ) (7, 8). Following heat shock, s32 levels increase rapidly, due to its increased rate of synthesis of s32 (7, 9-11) and increased stability during heat shock. The consequence of elevated levels of s32 is a burst of heat-shock gene transcription (12-14).

**Table 1. A Brief Summary of the Nomenclature, Location, and Function of the Major Chaperone Families**

| Organism | Chaperone | Compartment | Functions |
|----------|-----------|-------------|-----------|
|----------|-----------|-------------|-----------|

|        |                      |                |                 |   |
|--------|----------------------|----------------|-----------------|---|
| hsp100 | <i>E. coli</i>       | ClpA,B,C       | Cytosol         | Roles in stress tolerance; aids in the resolubilization of heat-inactivated proteins from insoluble aggregates; regulates yeast prions  |
|        | <i>S. cerevisiae</i> | Hsp 104        | Cytosol         |   |
| hsp90  | <i>E. coli</i>       | HtpG           | Cytosol         | Functions in signal transduction (eg, steroid hormone receptor activity); modulates protein kinase activities; refolds and maintains denatured proteins <i>in vitro</i>   |
|        | <i>S. cerevisiae</i> | Hsp83          | Cytosol         |   |
|        | Mammals              | Hsp90<br>Grp94 | Cytosol<br>ER   |   |
| hsp70  | <i>E. coli</i>       | DnaK           | Cytosol         | Functions in lambda phage replication; autoregulates the heat shock response; interacts with nascent chain polypeptides; may facilitate protein degradation by acting as a cofactor in proteolytic systems; dissociates clathrin; functions in interorganellar transport; roles in signal transduction (eg, steroid hormone receptor function); refolds and maintains denatured proteins <i>in vitro and in vivo</i> together with the Hdj-1/Hsp40 co-chaperone |
|        | <i>S. cerevisiae</i> | Ssa 1-4        | Cytosol         |   |
|        |                      | Ssb 1,2        | Cytosol         |   |
|        |                      | Kar2           | ER              |   |
|        |                      | Ssc1           | Mitochondria    |   |
|        | Mammals              | Hsc70          | Cytosol/nucleus |   |
|        |                      | Hsp70          | Cytosol/nucleus |   |
|        |                      | Bip            | ER              |   |



|       |                      |                             |                         |  |
|-------|----------------------|-----------------------------|-------------------------|--|
| hsp60 | <i>E. coli</i>       | mhsp70<br>groEL             | Mitochondria<br>Cytosol | Functions in assembly of bacteriohages and Rubisco; refolds or prevents aggregation of denatured proteins <i>in vitro</i> ; may facilitate protein degradation by acting as a cofactor in proteolytic systems; functions in folding of organellar proteins |
|       | <i>S. cerevisiae</i> | Hsp60                       | Mitochondria            |  |
| hsp40 | Plants               | Cpn60                       | Chloroplasts            | Chaperone activity, essential co-chaperone activity with Hsp70 proteins to enhance ATPase rate and substrate release   |
|       | Mammals              | Hsp60                       | Mitochondria            |  |
|       | <i>E. coli</i>       | dnaJ                        | Cytosol                 |  |
|       | <i>S. cerevisiae</i> | Ydj1                        | Cytosol/nucleus         |  |
| shsp  | Mammals              | Hdj-1/Hsp40<br>Hdj-2, Hsj-1 | Cytosol                 | Suppresses aggregation and heat inactivation of proteins <i>in vitro</i> ; confers thermotolerance through stabilization of microfilaments; possible roles in cell growth and differentiation  |
|       | <i>E. coli</i>       | IbpA and B                  |                         |  |
|       | <i>S. cerevisiae</i> | hsp27                       | Cytosol                 |  |
|       | Mammals              | aA and aB-crystallin        | Cytosol                 |  |
|       |                      | hsp27                       | Cytosol                 |  |

---

A complementary stress response to misfolded proteins in the periplasm and outer membrane of *E. coli* is mediated by a second heat-shock factor, sE (s24) (15, 16). Activation of sE leads to the induction of at least 10 additional proteins, four of which have been identified: the periplasmic proteinase degP, s32, the periplasmic [peptidyl prolyl cis/trans isomerase](#) fkpA, and sE itself (17-23). Cells lacking sE are sensitive to elevated temperatures, sodium dodecyl sulfate (SDS)/ethylenediaminetetraacetic acid (EDTA), and crystal violet (16, 21, 22, 24). The temperature-

sensitive phenotype of sE mutants can be restored by activation of a second [signal transduction](#) cascade, the Cpx pathway, revealing that *E. coli* has at least two partially overlapping stress signal cascades capable of relieving extracytoplasmic stress (25).

Studies on the regulation of the heat-shock response provide evidence for autoregulation in that the activity and/or levels of heat-shock proteins influences the expression of heat-shock genes. Mutations in the heat-shock proteins dnaK, dnaJ, or GrpE, or reduced expression of groEL/S, enhance heat-shock gene expression at normal growth temperatures and result in an extended heat-shock response after shift to high temperatures (26-28). These cells are defective in the control of s32 synthesis and exhibit increased stability of s32 (8, 27, 29). The degradation of s32 involves dnaK, dnaJ, and GrpE by targeting s32 to the ATP-dependent proteinases FtsH, Lon, Clp, and HslVU (30, 31). Negative regulation of s32 is mediated by interaction of a specific region of s32 with dnaK, such that deletion or mutations within this region lead to its sustained synthesis and prolonged half-life (32, 33). The interaction of s32 with dnaK, dnaJ, and GrpE is ATP-dependent, so that dnaJ, in the presence of nucleotide, increases the formation of dnaK–dnaJ–s32 complexes in the ADP state, and GrpE stimulates nucleotide exchange and dissociation of the complex (34, 35). Association of dnaK and dnaJ with s32 prevents s32 from binding to RNA polymerase, resulting in the arrest of heat-shock gene transcription (34, 36-38). In addition to sequestration of s32 from RNA polymerase, the dnaK chaperone machinery autoregulates the *E. coli* heat-shock response by the simultaneous reactivation of heat-aggregated s70 and the preferential assembly of s70 with RNA polymerase (39). The equilibrium between free active s32 and DnaK/DnaJ-bound inactive s32, therefore, constitutes an important feature of the homeostatic control of heat-shock gene regulation in *E. coli*.

## 2. Stress-Induced Transcriptional Regulation of Heat-Shock genes in Eukaryotes

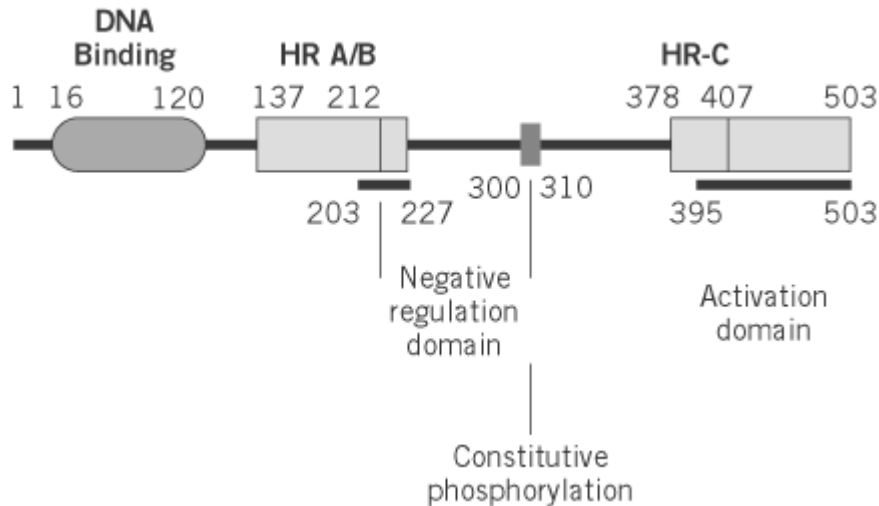
### 2.1. Family of Heat-Shock Factors

In higher eukaryotes, stress-induced regulation of the heat-shock response occurs principally by stress-induced activation of heat-shock factors (HSFs) from an inert to a transcription-competent **DNA-binding** state, binding to the heat-shock promoter element, and interaction with components of the basal transcription machinery, thus resulting in inducible transcription (40-43). Among eukaryotes, the complexity of HSF-dependent regulation varies among species, with yeast and *Drosophila* encoding a single HSF gene (44-46), whereas four HSF genes (HSF1 to HSF4) have been isolated and characterized in the [mouse](#), chicken, and human **genomes**, and a distinct but functionally related family of HSFs have been characterized in plants (47-54).

Comparison of the various **cloned** HSF genes reveals an overall sequence identity of 40%, with a high degree of structural conservation (Fig. 3) in the winged [helix–turn–helix](#) motif DNA-binding domain (43, 53, 55-57), an extended 80-residue **hydrophobic** repeat (HR-A/B) involved in trimerization (46, 58, 59), and a carboxyl-terminal localized transactivation domain (51, 60-63). With the exception of the HSF in budding yeasts and human HSF4, another hydrophobic repeat (HR-C) is located adjacent to the carboxyl-terminal transactivation domain, which has been suggested to function in suppression of trimer formation by interaction with HR-A/B (50, 64-68). Also positioned between HR-A/B and HR-C are sequences that function as negative regulatory regions that repress either or both DNA binding and transcriptional activation (61, 62, 68). Other features unique to specific HSFs are the presence of an amino-terminal transactivation domain in the *Saccharomyces cerevisiae* HSF and the lack of a functional transactivation domain in human HSF4 (51). A common feature of HSFs other than from *S. cerevisiae* is the negative regulation of DNA binding or transcriptional activity.

**Figure 3.** General features of heat-shock factors. Schematic representation of HSF1 structural motifs [DNA binding, hydrophobic repeats (HR-A/B and HR-C), and the transcriptional activation domain] and the negative regulatory domains which influence HSF1 activity. The relative positions of these domains are indicated by the amino acid

residues indicated.

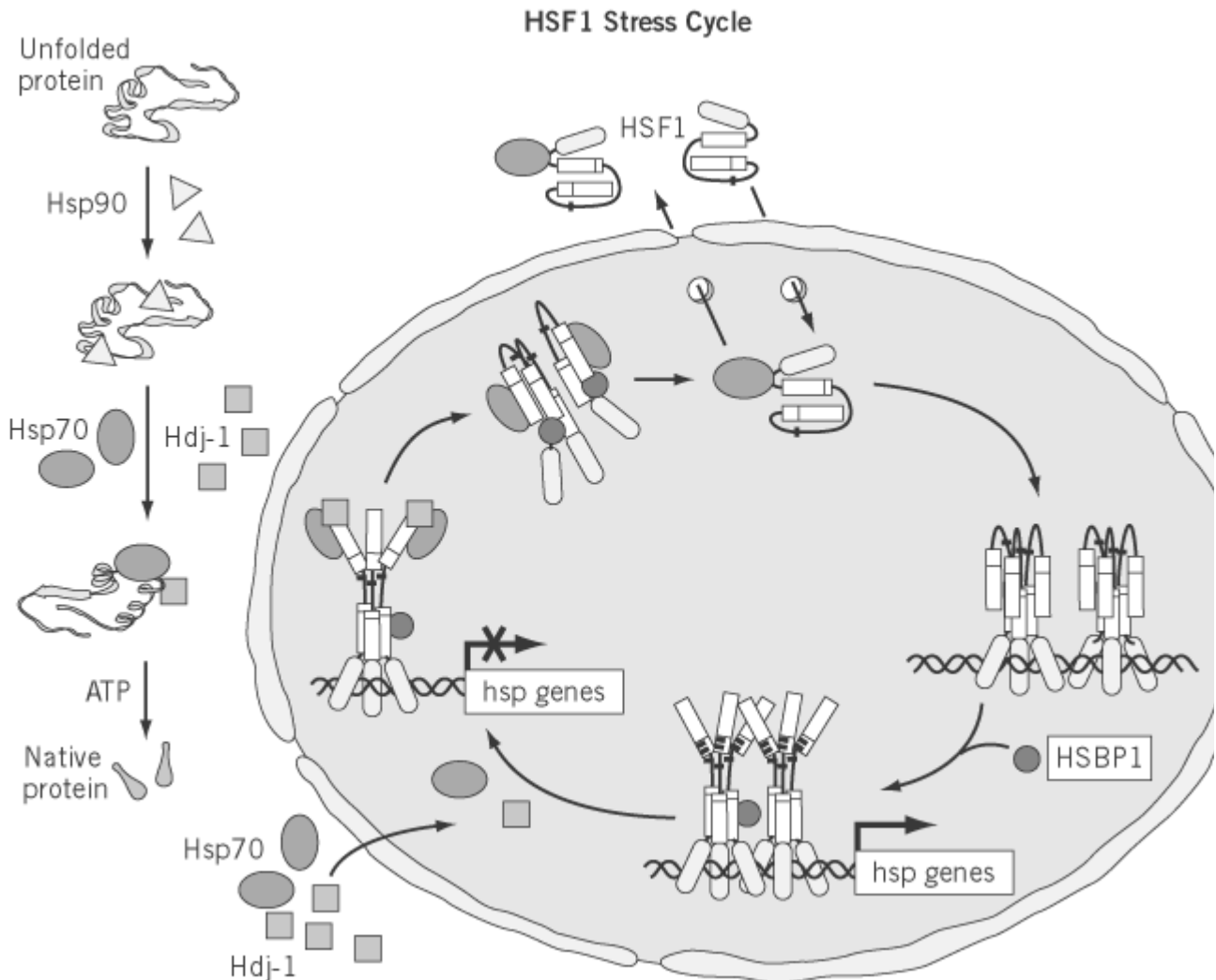


Of the HSFs coexpressed in vertebrates, HSF1 has the properties of the predominant stress-induced transcriptional activator, which acquires both DNA binding and transcriptional activity and is functionally analogous to yeast and *Drosophila* HSF (69, 70). In avian cells, HSF1 and HSF3 are coexpressed and coactivated by chemical and physiological stress, leading to the suggestion of possible redundancy, yet cells deleted for HSF3 and expressing HSF1 are severely compromised for the heat-shock transcriptional response (71-73). HSF3 also interacts with other transcription factors and can be activated, independent of stress, by the *Myb oncogene* via direct **protein-protein interaction** with the HSF3 DNA-binding domain (74). Myb is a growth-regulated **transcription factor**; consequently the interaction between Myb and HSF3 reveals a new pathway for genetic crosstalk between cell growth and the stress response. Evidence is also accumulating to indicate that HSF1 may negatively influence the transcriptional activation of other cellular genes, independent of direct binding to heat-shock elements, perhaps via protein-protein interactions with HSF1 (75). Another member of the HSF family, HSF2, in contrast to HSF1, regulates heat-shock gene expression during cell growth and differentiation. HSF2 activity was first described in human erythroleukemia (K562) cells exposed to hemin (76-78), and it was subsequently detected during murine spermatocyte differentiation (79) and embryogenesis (80, 81). The stress-sensing pathway for HSF2 activation is associated with **protein degradation**. HSF2 activity is induced, in a cell-type independent manner, in mammalian tissue culture cells by exposure to specific inhibitors of the **ubiquitin-dependent proteasome** or in cells harboring conditional mutations in components of the proteasome machinery (82). The function of the ubiquitin-dependent proteasome machinery is to degrade short-lived and misfolded proteins; therefore, activation of HSF2 by down-regulation of the protein degradative machinery reveals a requirement for heat-shock proteins, perhaps to prevent misfolding and aggregation of non-native proteins targeted for degradation.

## 2.2. Regulation and Autoregulation of HSF

Stress-induced activation of HSF1 is a multistep process (Fig. 4). In response to heat shock and other stresses, HSF1 translocates and relocalizes within the nucleus, oligomerizes to a DNA-binding and transcriptionally inert state, becomes inducibly **phosphorylated**, and activates the transcription of heat-shock genes. During prolonged heat shock, HSF1 activity attenuates, as characterized by transcriptional repression of heat-shock genes. Key events during attenuation include the association of HSF1 with the molecular chaperones Hsp70 and Hdj1, interaction with the negative regulator HSBP1, dissociation of HSF1 trimers from DNA, and refolding of the trimer to the monomer. Translocation of HSF1 from the cytoplasm to the nucleus has been described in *Drosophila* and mammalian cells (43, 57, 69, 70, 78, 83-86), yet others have also observed that HSF1 is localized constitutively to the **nucleus** (86).

**Figure 4.** Regulation of the heat-shock response. Activation of heat-shock factor (HSF1) is linked to the appearance of native proteins and the requirement for molecular chaperones (Hsp90, Hsp70, and Hdj1) to prevent the appearance of misfolded proteins. HSF1 exists in the control state as an inert monomer (shown as intramolecularly and negatively regulated for DNA binding and transcriptional activity) and undergoes stepwise activation to (i) a DNA-binding competent state that is transcriptionally inert, (ii) acquisition of inducible phosphorylation, resulting in complete activation and inducible transcription of heat-shock genes, and (iii) attenuation of HSF1 activity. HSF1 activity is negatively regulated by heat-shock factor binding protein 1 (HSBP1), which binds to the region of HSF1 corresponding to the [heptad repeat](#), and Hsp70 and Hdj1, which bind to the transcriptional *trans*-activation domain thus repressing HSF1 activity.



Maintenance of HSF1 in a repressed state is delicately balanced and easily disrupted; for example, overexpression of HSF1 by transient transfection results in constitutively active HSF1; likewise, mutation of a single critical phosphoserine residue has the same derepressing effect ([70](#), [87](#), [88](#)). Acquisition of HSF1 DNA-binding activity is insufficient to activate transcription, as demonstrated by the effects of the anti-inflammatory drugs sodium salicylate, indomethacin, or ibuprofen, which induce DNA binding-competent trimers that, nevertheless, lack inducible phosphorylation and are transcriptionally active ([89-91](#)). Rather than corresponding to a nonfunctional form of HSF1, however, the salicylate-induced HSF1 can be converted *in vivo* to the fully active HSF1 by a subsequent exposure to heat shock ([91](#)).

Exposure of human cells to 42°C heat shock results in the rapid and transient inducible transcription of heat-shock genes, with maximal rates of hsp70 and hsp90 transcription occurring within 30 to 60

minutes; thereafter, the rate of heat-shock gene transcription declines rapidly, prior to the loss of HSF1 DNA-binding activity (88, 92, 93). The increased levels of heat-shock messenger RNAs following heat shock result from both transcriptional activation of heat-shock genes and the increased stability of heat-shock messenger RNAs during stress (93, 94). This results in the elevated synthesis and accumulation of heat-shock proteins.

The initial evidence that the eukaryotic heat-shock response is autoregulated was based on the effects of amino acid analogues on *Drosophila* cells. Rather than the expected transient heat-shock response, treatment with azetidine, a proline analogue, resulted in the continuous activation of heat-shock gene expression (95). An explanation for this result was that amino acid analogues induce the heat-shock response as a result of their incorporation into nascent polypeptides that misfold. Such amino acid analogue-containing nascent polypeptides would be expected to associate with Hsp70 (96); the sequestration of Hsp70 by these permanently unfolded polypeptides in turn results in activation of heat-shock gene transcription. However, unlike the heat-shock response, where attenuation is linked to the *de novo* synthesis of heat-shock proteins, Hsp70 synthesized in the presence of amino acid analogues is also misfolded, and therefore nonfunctional as a chaperone, with the consequence that the heat-shock response does not attenuate. These observations led to the proposal that heat-shock proteins are involved in autoregulation of the heat-shock response. Genetic evidence to support the autoregulation of the heat-shock response has shown that overexpression of the yeast Ssa1p (cytosolic Hsp70) dampens the heat-shock response from two different (*ssa1* and *ssa4*) promoters for Hsp70 (97). Likewise, deletion of the *ssa1* and *ssa2* Hsp70 genes resulted in an unusually high level of expression of another Hsp70 gene (*Ssa3p*) and other heat-shock proteins in a HSF-dependent manner; these results strongly implicated HSF as a potential target for autoregulation by members of the Hsp70 family (98). The relationship between HSF and Hsp70 was further supported by a search for extragenic suppressors of the temperature-sensitive phenotype of an Hsp70 mutant (*ssa1ssa2* strain), which also uncovered HSF as an interactive component of the regulatory response (99). A spontaneous mutant EXA3 that could reverse the growth defect of the Hsp70 mutant (*ssa1ssa2*) is very closely linked with the gene encoding HSF, and another mutation identified in the genetic screen maps to the HSF gene (100).

Biochemical evidence to support a role for molecular chaperones in the regulation of the heat-shock response have demonstrated that stress-induced HSF1 trimers were associated with Hsp70 (101-105). The Hsp70-HSF complexes are sensitive to ATP, which is typical of chaperone-substrate interactions, and can be reconstituted *in vitro* (101, 105). Through the use of HSF1 deletion mutants and direct *in vitro* binding assays, a site for Hsp70 binding was mapped on the transactivation domain of HSF1 (105). Whereas overexpression of Hsp70 inhibited completely the induction of heat-shock gene transcription, there was little or no effect on the formation of HSF1 trimers or on inducible phosphorylation of HSF1 (103, 105). This demonstrates that Hsp70 can also function as a negative regulator of HSF1 activity and that the repression of heat-shock gene transcription that occurs during attenuation is due to the repression of the HSF1 transactivation domain by direct binding to Hsp70 (101-103, 105). Although Hsp70 has a demonstrated role in the autoregulation of the heat-shock response, other molecular chaperones, such as Hdj1 or Hsp90, have also been implicated in the regulation of HSF1 activity. Hdj1 interacts with HSF1 in higher eukaryotes and negatively regulates HSF1 transcriptional activity (105). A role for members of the DnaJ family in regulation of the heat-shock response is supported by the observation in *S. cerevisiae* that the DnaJ homologue, SIS1, negatively regulates its own expression (106). However, SIS1 autoregulation requires the heat-shock element (HSE) and other sequences, suggesting that additional regulatory molecules might be involved. HSF may also associate with Hsp90; such interactions have been detected for yeast, rat, and rabbit Hsp90 (107, 108) but not with human Hsp90 (102, 103, 105).

Identification of the transactivation domain of HSF1 as a chaperone-binding site offers a number of intriguing possibilities for Hsp70 as a competitor for HSF1 interaction with the basal transcriptional machinery (109) or the role of Hsp70 in the conformational change of HSF1. The consequence of Hsp70 binding may render the transactivation domain inaccessible to the transcriptional machinery, thus resulting in the transcriptional repression. A related role for chaperones as regulators of

transcriptional activators has also been described for the family of **steroid receptors**, although in this case multiple chaperones are recruited to maintain the activator in a repressed state by formation of a stable chaperone-substrate complex (110). It is tempting to consider that the interactions observed between p53 and Hsp70 (Hsc70) could also reflect a form of chaperone-dependent regulation of the folded state that affects the properties of the transcriptional activator (111).

### 3. Molecular Chaperones: Roles in Translocation, Protein Folding, and Refolding

#### 3.1. General Features of Molecular Chaperones

Although some proteins refold spontaneously *in vitro* when diluted at low concentrations from **denaturants**, larger, multidomain proteins often have a propensity to misfold and aggregate (see **Protein folding**). Consequently, the challenge within the densely packed environment of the cell is to capture and sequester efficiently non-native proteins that may lead subsequently to either refolding or degradation. Molecular chaperones of the Hsp90, Hsp70, and Hsp60 class accomplish this by capturing non-native intermediates and, together with co-chaperones and ATP, facilitate the appearance of the folded native state (112). The Hsp70 chaperones, for example, recognize stretches in polypeptides that are rich in **hydrophobic** residues and are transiently exposed in early folding intermediates, but typically confined to the hydrophobic core in the native state (113). This contrasts to the Hsp60/GroEL **chaperonin**, which creates a protected environment with the properties of a “protein-folding test-tube” in which non-native proteins can undergo rounds of binding and release to acquire the native state (114, 115). The consequence of chaperone interactions with early folding intermediates, therefore, is to shift the equilibrium of protein folding toward on-pathway events and to minimize the appearance of nonproductive intermediates that may have a propensity to aggregate as misfolded species.

The family of molecular chaperones is large; they are diverse in size and apparent structure, yet highly conserved throughout evolution and abundant under nonstress conditions (40, 112, 116). Most chaperones are abundant in growing cells and can attain concentrations of 1% to 5% of total cell protein. They are classified according to molecular size, that is, Hsp100, Hsp90, Hsp70, Hsp60, Hsp40, and small Hsps (Table 1). Biochemical studies on chaperones have shown that a common function of the chaperones Hsp104, Hsp90, Hsp70, the small Hsps, **immunophilins** (FKBP52 and CyP40), the steroid aporeceptor protein p23, and Hip (Hsp70 and Hsp90 interacting protein) is to prevent the *in vitro* aggregation of model protein substrates and to maintain the substrate in an intermediate folded state competent for subsequent refolding to the native state. A distinction among proteins with chaperone activities is that refolding to the native state requires the activity of a specific subset of chaperones, such as Hsp90, Hsp70, or Hsp60/GroEL, which are **nucleotide-binding** proteins, and the regulatory properties of chaperone-specific co-chaperones, such as p23, the immunophilins, which enhance Hsp90, dnaJ, or Hip proteins that associate with Hsp70, and Hsp10/groES, which stimulates Hsp60/groEL. The suggestion that chaperones exist as heteromeric complexes that associate with different folded intermediates allows the substrate to acquire different regulatory states.

#### 3.2. Representative Analysis of the Hsp70 Family

Chaperones are found in nature as functionally related protein families; for example, the eukaryotic Hsp70 proteins are found in the cytoplasm, **nucleus**, **mitochondria**, and **endoplasmic reticulum**, where they maintain proteins in intermediate folded states competent for translocation across the **membranes** of the mitochondria or endoplasmic reticulum and for subsequent refolding to the native state (116-119). In *E. coli*, the Hsp70 homologue, dnaK, originally identified in a genetic screen for host proteins required for replication of **lambda phage**, is also required for growth at nonpermissive temperatures (120, 121). Analysis of the structure and functional properties of Hsp70 reveals a multi-domain organization corresponding to a 45-kDa amino-terminal **ATPase** domain, the 18-kDa carboxyl-terminal peptide-binding domain, and a 7-kDa intra- and interdomain EEVD (Glu–Glu–Val–Asp) regulatory motif (122-124).

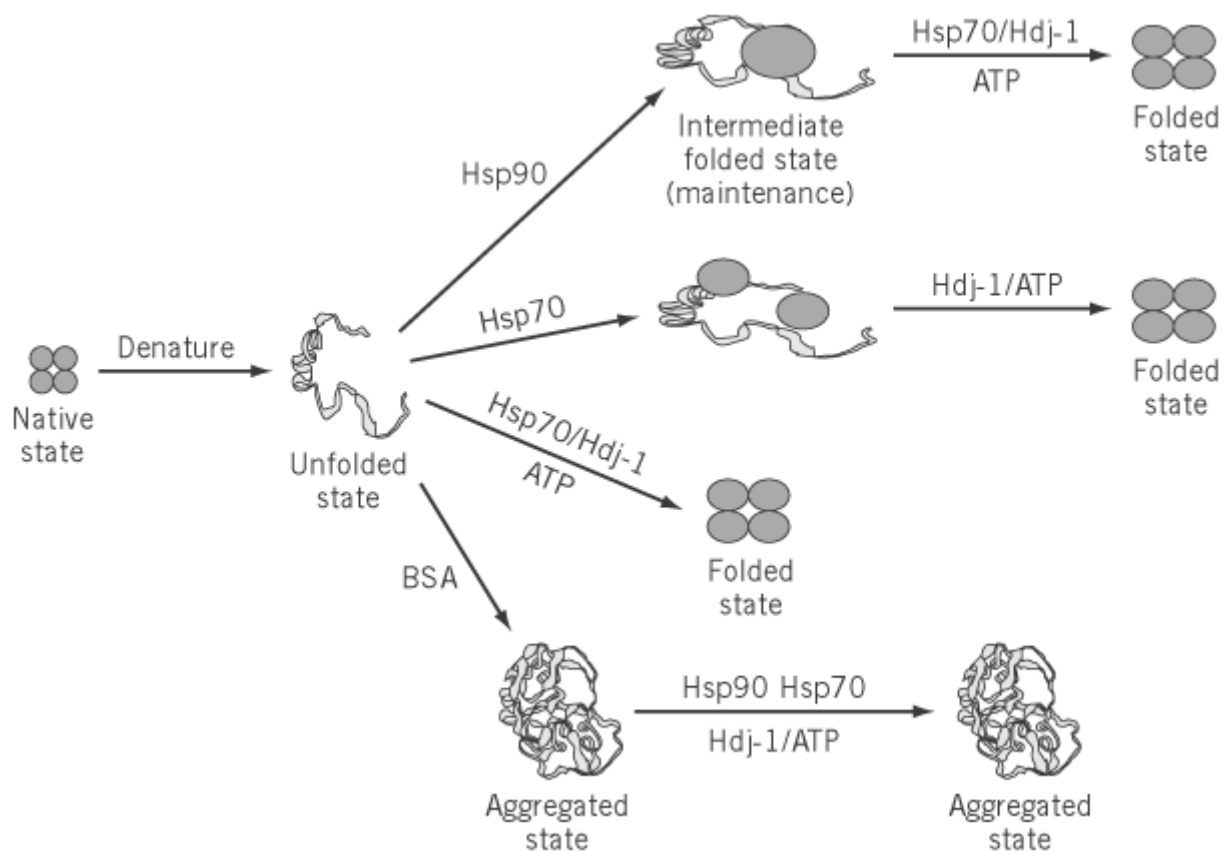
Hsp70 interacts with unfolded polypeptides and peptides by binding to arrays of hydrophobic

residues in which acidic residues are excluded and the surrounding environment enriched in basic residues (113). Typically, such sites are often buried within the hydrophobic core of folded proteins and are only transiently exposed during protein folding. Affinity of the peptide-binding domain for denatured substrates is dependent on the nucleotide state of the ATPase. Hsp70 exhibits a high affinity for substrate when bound to ADP; upon release of ADP and re-binding of ATP, the affinity for peptides decreases and the substrate is released (125, 126). ATP is hydrolyzed to ADP, and the cycle begins again. The intrinsic ATPase of Hsp70 is exceedingly slow and stimulated by increased temperature, the presence of potassium, calcium, and magnesium ions, the presence of substrates, and the co-chaperones dnaJ and GrpE (126-129). DnaJ stimulates the dnaK ATPase and is required for both the *in vivo* and *in vitro* folding activities of dnaK. dnaJ also functions independently as a chaperone with “holding” activity and has been shown to prevent the aggregation of unfolded rhodanase and **luciferase** in an ATP-independent manner (130). The chaperone activity of DnaJ may not be essential, however, because the human homologue, Hdj-1, neither binds non-native substrates nor prevents protein aggregation. Hdj-1 is necessary for Hsp70-mediated refolding (131). The second dnaK co-chaperone is GrpE, which functions in concert with dnaJ to enhance ADP release, thus allowing the unfolded substrate to dissociate from the dnaK/dnaJ/GrpE complex, also known as the Hsp70 chaperone machine.

### 3.3. Biochemical Properties in Protein Folding

Among their multitude of biochemical activities, certain chaperones interact with unfolded polypeptides to “hold” non-native proteins in a refolding-competent state, whereas others can also “refold” the polypeptide to its final native state. The relationships between these events in protein folding are indicated in Figure 5. Representatives of the former class include abundant molecular chaperones, such as Hsp90 and the small Hsps. *In vitro*, these heat-shock proteins stably maintain non-native polypeptides in intermediate refoldable states that can then be subsequently refolded by the addition of nucleotide, Hsp70, and Hdj-1 (131, 132). Chaperones of the latter class include Hsp70 and groEL (Hsp60), which can both hold and refold (112, 131). The refolding activity requires accessory proteins or co-chaperones, nucleotide binding, and hydrolysis. In the absence of these co-factors, neither Hsp70 nor groEL are effective in refolding a polypeptide to the native state (124, 131, 133).

**Figure 5.** Biochemical properties of molecular chaperones in protein folding. The biochemical fate of an unfolded protein is presented schematically. The native protein substrate is denatured (either chemically or thermally) and acquires an unfolded state. Dilution of the denatured protein into buffer containing a nonspecific protein, bovine serum albumin, leads to the aggregated state that is not resolubilized by chaperones (Hsp90, Hsp70, Hdj1) in the presence of ATP. In contrast, dilution of the unfolded protein into the molecular chaperones Hsp70, Hdj1, and ATP results in refolding to the native, enzymatically active state. The holding or intermediate folded state is acquired by dilution into either Hsp90 or Hsp70 alone; subsequent addition of Hsp70, ATP, and the co-chaperone Hdj-1 results in folding to the native state.



While many of the biochemical properties of molecular chaperones are based on *in vitro* approaches, the situation *in vivo* is more complex due to the coexistence of multiple families of chaperones and diverse substrates within the same subcellular compartments. Most chaperones are expressed constitutively and are relatively abundant in each subcellular compartment, where they function in activities including cell signaling and regulating protein [kinase](#) activities. Chaperones including Hsp70, Hsp90, p23, and immunophilins are essential components of steroid aporeceptor complexes and necessary for steroid responsiveness ([134](#), [135](#)). In the unstressed cell, chaperones also modulate the folded state of polypeptides; for example, members of the Hsp70 family have a role in **protein biosynthesis** in [translation](#) elongation ([117](#), [136](#)), they maintain mitochondrial-targeted polypeptides as unfolded intermediates in a translocation-competent state, and they facilitate the refolding of translocation-competent proteins in concert with Hsp60 ([112](#), [118](#)). During periods of stress, the expression of many, but not all, chaperones is induced to cope with an increased volume of unfolded proteins. The association of unfolded proteins with chaperones may retard protein aggregation, facilitate refolding of non-native intermediates, and provide an interface with the **protein degradation** machinery.

Molecular chaperones stimulate on-pathway folding, principally by their ability to (a) function as kinetic traps to prevent off-pathway intermediates that may lead to the formation of aggregates and (b) favor folding to the native state (Fig. 5). Chaperones exhibit transient and stable interactions with non-native protein substrates; under normal conditions of cell growth, chaperones are associated with kinases and steroid aporeceptors ([137](#)). Mutations that affect protein function often lead to stable complexes with chaperones; for example, mutant forms of [p53](#) are stably associated with Hsp70 ([111](#), [138](#)). In such examples, modulation of the activity or half-life of the chaperone-associated substrate may influence cellular events and pathogenesis. Both the absolute and relative abundance of chaperones influences protein biogenesis in the unstressed cell and during periods of stress, where the expression of many but not all chaperones is induced to cope with an increase in the flux of non-native proteins.



#### 4. Role of Heat-Shock Proteins and Stress-Induced Cytoprotection in Disease

The abnormal expression of heat-shock proteins has been well documented for diseases including ischemia and reperfusion damage, cardiac hypertrophy, fever, inflammation, metabolic diseases, bacterial and viral infection, cell and tissue injury, aging, and cancer (40, 139, 140). Serum **antibodies** to the heat-shock proteins Hsp90, Hsp70, Hsp60, and small heat-shock proteins have been detected in individuals with infectious diseases and [autoimmune diseases](#), including rheumatoid arthritis and insulin-dependent diabetes (141). These observations have led to questions whether stimulation of the [immune response](#) is the result of, or causes, cell damage, and whether the appearance of circulating antibodies to heat-shock proteins reflects the organismal response to stress and damage. Although the acute response to stress may be critical for recovery and long-term survival, the chronic expression of heat-shock proteins in diseased tissues may be deleterious to protein biogenesis and have negative effects on cell growth. Because the ability of a cell to survive stress is dependent on both the kinetics of activation and the duration of the response, changes in the environment of the cell will shift the stress-sensing capacity of the heat-shock response, with profound effects on survival, repair, and recovery.

The cytoprotective role of heat-shock proteins offers novel approaches for the treatment of diseases involving metabolic disorders, inflammation, infection, and ischemia (142, 143). Ischemia of the brain, kidney, heart, and liver results from decreased blood flow to tissues; the associated pathology that ensues reflects adaptation to decreased levels of oxygen and nutrients. Expression of heat-shock proteins are rapidly induced, during ischemia and to an even greater degree upon reperfusion associated with the appearance of oxygen-free radicals. Exposure to environmental toxins, such as xenobiotics and aromatic hydrocarbons that promote oxidative damage, could potentiate the damaging effects of ischemia and reperfusion. Cytoprotection and repair against the deleterious effects of stress and trauma can be accomplished by the overexpression of one or more heat-shock protein genes. Yeast cells engineered to overexpress Hsp70, the small heat-shock proteins, or Hsp 104 are also protected against exposure to lethal heat-shock temperatures, H<sub>2</sub>O<sub>2</sub>, heavy metals, arsenite, anoxia, and ethanol toxicity (1). In vertebrates, modulation of the heat-shock response or the expression of specific heat-shock proteins can limit or prevent the pathologies associated with certain chronic diseases. A class of diseases associated with the appearance of misfolded proteins are the [prion](#) diseases, which from yeast to humans are associated with changes in the folded state of the prion that correlates with conversion from a noninfective to an infective state. Among the fascinating features of prion diseases is the [epigenetic](#) transmission of an altered protein conformation, thus endowing the infectious prion with unique biological properties. Recent data on the yeast prions Ura3 and Psi, which have features in common with the scrapie prion protein, reveal that the molecular chaperone Hsp 104 has an essential role in solubility of the Psi protein and the regulation of the Psi+ phenotype (144). Whether these observations have direct implications for heat-shock proteins in neurodegenerative diseases remains an exciting dimension for future studies.

In the pathology of myocardial disease, the induction of Hsp70 is associated with both ischemia and reperfusion and therefore could reflect the appearance of damaged proteins during adaptation of the myocardium. The elevated synthesis of heat-shock proteins, therefore, could reflect the response of myocytes to survive the stress by repairing protein damage. For example, the induction of Hsp70 following aortic constriction or work-overload-induced cardiac hypertrophy in animal models could reflect a response to the aberrant synthesis, accumulation, or degradation of proteins or, alternatively, a response to events that occur during the partial reentry of the myocardial cell into the growth cycle (142, 143). Does activation of the heat-shock response reflect protein damage and does this reveal a potential strategy to detect or manipulate the course of events? A correlation between the levels of Hsp70 and the degree of myocardial protection has been demonstrated; moreover, transgenic mice overexpressing Hsp70 exhibited an enhanced resistance to myocardial ischemic stress, thus providing direct evidence for a role of heat-shock proteins in cytoprotection (145-147). These studies reveal that induction of the heat-shock response by the stressed myocardium is

proportional to protein damage; therefore a potential strategy would be to enhance the expression of heat-shock proteins such as Hsp70, which in turn may confer a more rapid reestablishment of normal cardiac protein synthesis and myocardial function. The activation of a heat-shock response restores normal cardiac function after injury, possibly by removal of misfolded cardiac proteins and reestablishment of normal cardiac protein synthesis (142).

As an alternative to genetic manipulation to alter heat-shock gene expression, it may be more fruitful to search for small molecules that enhance either the expression or function of heat-shock proteins. Such a pharmacological approach has been suggested, for example, based on experiments using herbimycin-A as an inducer of Hsp70 with protective effects in simulated ischemia on rat neonatal cardiomyocytes (148) and the cytoprotective activity of a hydroxylamine derivative (Bimoclomol) during ischemia and wound healing (149). Other classes of small molecules with heat-shock regulatory properties are nonsteroidal anti-inflammatory drugs (NSAIDs), cyclopentenone prostaglandins, **serine proteinase inhibitors**, and inhibitors of the ATP-dependent **ubiquitin-dependent proteasome** (89, 150, 151). Salicylates and other NSAIDs activate an intermediate form of HSF1, which potentiates the cellular response to stress (89). Pretreatment with NSAIDs such as salicylates or indomethacin (at subthreshold concentrations) both decreases the temperature threshold of the heat-shock response and confers cytoprotection (150). Exposure to aspirin or indomethacin at concentrations comparable to clinical levels results in the priming of human cells for subsequent exposure to heat shock and other stresses, the enhanced transcription of heat-shock genes, and cytoprotection from thermal injury (152). Considering that the expression of heat-shock genes occurs in response to pathological conditions, it would not be surprising if the ability of NSAIDs to activate the DNA-binding activity of HSF and Hsp70 expression contributes to their pharmacological efficacy. Another class of pharmacologically active small molecules in the pathway of arachidonate metabolism are cyclopentenone prostaglandins, which activate the heat-shock response and protect against thermal injury and viral infection (151). It is intriguing to note that this class of molecules has cardioprotective value (153). The development of pharmacologically active small molecules that influence either the regulation or function of heat-shock proteins offers a means to harness the positive cytoprotective value of the heat-shock response.

## 5. Summary

- Exposure of cells to stresses such as heat shock, oxidant injury, toxic chemicals, and heavy metals causes an imbalance in protein metabolism that challenges the cell to respond rapidly, yet precisely, to minimize the deleterious effects of environmental and physiological stress.
- The heat-shock response, by activation of heat-shock transcription factors, results in the elevated expression of heat-shock genes and the concomitant synthesis of heat-shock proteins and molecular chaperones.
- Molecular chaperones function in a variety of protein biosynthetic events to protect proteins from the deleterious effects of acute or chronic stress by stabilizing and refolding protein folding intermediates or facilitating protein degradation.
- Accumulation of misfolded proteins is of central importance to diseases of protein folding, including **sickle-cell** hemoglobin, cystic fibrosis, and prion diseases, in addition to complex multifactorial diseases, including bacterial and viral infections, myocardial ischemia, neurodegenerative diseases, and cancer.

## 6. Acknowledgments

The author was supported by a research grant from the NIH and the Fogarty International Fellowship. The assistance of K. Veraldi and S. Fox in preparation of this manuscript is appreciated.

## Bibliography

1. D. A. Parsell and S. Lindquist (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 457–494.
2. A. D. Grossman, J. W. Erickson, and C. A. Gross (1984) *Cell* **38**, 383–390.
3. D. W. Cowing et al. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 2679–2683.
4. M. Bloom et al. (1986) *J. Bacteriol.* **166**, 380–384.
5. T. Yura, H. Nagai, and H. Mori (1993) *Annu. Rev. Microbiol.* **45**, 301–325.
6. E. A. Craig and C. A. Gross (1991) *Trends Biochem. Sci.* **16**, 135–140.
7. A. D. Grossman, D. B. Straus, W. A. Walter, and C. A. Gross (1987) *Genes Dev.* **1**, 179–184.
8. K. Tilly, J. Spence, and C. Georgopoulos (1989) *J. Bacteriol.* **171**, 1585–1589.
9. D. B. Straus, W. A. Walter, and C. A. Gross (1987) *Nature* **329**, 348–351.
10. A. S. Kamath-Loeb and C. A. Gross (1991) *J. Bacteriol.* **173**, 3904–3906.
11. H. Nagai, H. Yuzawa, and T. Yura (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10515–10519.
12. S. A. Lesley, N. E. Thompson, and R. R. Burgess (1987) *J. Biol. Chem.* **262**, 5404–5407.
13. S. Skelly et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8365–8369.
14. D. B. Straus, W. A. Walter, and C. A. Gross (1987) *Nature* **329**, 348–351.
15. J. Meccas et al. (1993) *Genes Dev.* **7**, 2618–2628.
16. P. E. Rouviere et al. (1995) *EMBO J.* **14**, 1032–1042.
17. J. W. Erickson et al. (1987) *Genes Dev.* **3**, 1462–1471.
18. B. Lipinska, S. Sharma, and C. Georgopoulos (1988) *Nucleic Acids Res.* **16**, 10053–10067.
19. J. W. Erickson and C. A. Gross (1989) *Genes Dev.* **3**, 1462–1471.
20. Q. P. Wang and J. M. Kaguni (1989) *J. Bacteriol.* **171**, 4248–4253.
21. S. Raina, D. Missiakas, and C. Georgopoulos (1995) *EMBO J.* **14**, 1043–1055.
22. D. Missiakas and S. Raina (1997) *EMBO J.* **16**, 1670–1685.
23. P. N. Danese and T. J. Silhavy (1997) *Genes Dev.* **11**, 1183–1193.
24. K. Hiratsu et al. (1995) *J. Bacteriol.* **177**, 2918–2922.
25. L. Connolly A. De Las Penas, B. M. Alba, and C. A. Gross (1997) *Genes Dev.* **11**, 2012–2021.
26. K. Tilly, N. McKittrick, M. Zylicz, and C. Georgopoulos (1983) *Cell* **34**, 641–646.
27. D. B. Straus, W. A. Walter, and C. A. Gross (1990) *Genes Dev.* **4**, 2202–2209.
28. M. Kanemori, H. Mori, and T. Yura (1994) *J. Bacteriol.* **176**, 4235–4242.
29. A. D. Grossman, D. B. Straus, W. A. Walter, and C. A. Gross (1987) *Genes Dev.* **1**, 179–184.
30. T. Tomoyasu et al. (1995) *EMBO J.* **14**, 2551–2560.
31. M. Kanemori, K. Nishihara, H. Yanagi, and T. Yura (1997) *J. Bacteriol.* **179**, 7219–7225.
32. H. Nagai, H. Yuzawa, M. Kanemori, and T. Yura (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10280–10284.
33. J. S. McCarty et al. (1996) *J. Mol. Biol.* **256**, 829–837.
34. K. Liberek and C. Georgopoulos (1993) *Proc. Natl. Acad. Sci. USA* **90**, 11019–11023.
35. J. Gamer et al. (1996) *EMBO J.* **15**, 607–617.
36. S. Skelly et al. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 5497–5501.
37. K. Liberek, T. P. Galitski, M. Zylicz, and C. Georgopoulos (1992) *Proc. Natl. Acad. Sci.* **89**, 3516–3520.
38. J. Gamer, H. Bujard, and B. Bukau (1992) *Cell* **69**, 833–842.
39. A. Blaszczyk, M. Zylicz, C. Georgopoulos, and K. Liberek (1995) *EMBO J.* **14**, 5085–5093.
40. R. I. Morimoto, A. Tissieres, and C. Georgopoulos (1990) *Stress Proteins in Biology and Medicine*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–36.
41. J. T. Lis and C. Wu (1993) *Cell* **74**, 1–20.

42. R. I. Morimoto (1993) *Science* **269**, 1409–1410.
43. C. Wu (1995) *Annu. Rev. Cell Dev. Biol.* **11**, 441–469.
44. G. Wiederrecht, D. Seto, and C. S. Parker (1988) *Cell* **54**, 841–853.
45. P. K. Sorger and H. R. Pelham (1988) *Cell* **54**, 855–864.
46. J. Clos et al. (1990) *Cell* **63**, 1085–1097.
47. S. K. Rabindran, G. Giorgi, J. Clos, and C. Wu (1991) *Science* **259**, 230–234.
48. K. D. Sarge et al. (1991) *Genes Dev.* **5**, 1902–1911.
49. T. J. Scheutz et al. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 6911–6915.
50. A. Nakai and R. I. Morimoto (1993) *Mol. Cell Biol.* **13**, 1983–1997.
51. A. Nakai et al. (1997) *Mol. Cell Biol.* **17**, 469–481.
52. K. D. Scharf et al. (1990) *EMBO J.* **9**, 4495–4501.
53. L. Nover et al. (1996) *Cell Stress Chap.* **1**, 215–223.
54. L. Nover and K. D. Scharf (1997) *Cell. Mol. Life Sci.* **53**, 80–103.
55. C. J. Harrison, A. A. Bohm, and H. C. M. Nelson (1994) *Science* **263**, 224–227.
56. G. W. Vuister, S. J. Kim, C. Wu, and A. Bax (1994) *Biochem.* **33**, 10–16.
57. C. Wu, J. Clos et al. (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 417–455.
58. P. K. Sorger and H. C. M. Nelson (1989) *Cell* **59**, 807–813.
59. R. Peteranderl and H. C. M. Nelson (1992) *Biochem.* **31**, 12272–12276.
60. M. Green, T. J. Scheutz, E. K. Sullivan, and R. E. Kingston (1995) *Mol. Cell Biol.* **15**, 3354–3362.
61. Y. Shi, P. E. Kroeger, and R. I. Morimoto (1995) *Mol. Cell Biol.* **15**, 4309–4318.
62. J. Zuo, D. Rungger, and R. Voellmy (1995) *Mol. Cell Biol.* **15**, 4319–4330.
63. J. Wisniewski, A. Orosz, R. Allada, and C. Wu (1996) *Nucleic Acids Res.* **24**, 367–374.
64. S. K. Rabindran et al. (1993) *Science* **259**, 230–234.
65. J. Zuo, R. Baler, G. Dahl, and R. Voellmy (1994) *Mol. Cell Biol.* **14**, 7557–7568.
66. A. Orosz, J. Wisniewski, and C. Wu (1996) *Mol. Cell Biol.* **16**, 7018–7030.
67. E. Zandi, T. N. Tran, W. Chamberlain, and C. S. Parker (1997) *Genes Dev.* **11**, 1299–1314.
68. T. Farkas, Y. A. Kutsikova, and V. Zimarino (1998) *Mol. Cell Biol.* **18**, 906–918.
69. R. Baler, G. Dahl, and R. Voellmy (1993) *Mol. Cell Biol.* **13**, 2486–2496.
70. K. Sarge, S. P. Murphy, and R. I. Morimoto (1993) *Mol. Cell Biol.* **13**, 1392–1407.
71. A. Nakai et al. (1995) *Mol. Cell Biol.* **15**, 5268–5278.
72. M. Tanabe, A. Nakai, Y. Kawazoe, and K. Nagata (1997) *J. Biol. Chem.* **272**, 15389–15395.
73. M. Tanabe et al. (1998) *EMBO J.* **17**, 1750–1758.
74. C. Kanei-Ishii et al. (1997) *Science* **277**, 246–248.
75. C. M. Cahill et al. (1997) *Adv. Exp. Med. Biol.* **400**, 625–630.
76. N. G. Theodorakis et al. (1989) *Mol. Cell Biol.* **9**, 3166–3173.
77. L. Sistonen et al. (1992) *Mol. Cell Biol.* **12**, 4104–4111.
78. L. Sistonen, K. D. Sarge, and R. I. Morimoto (1994) *Mol. Cell Biol.* **14**, 2087–2099.
79. K. D. Sarge et al. (1994) *Biol. Reprod.* **50**, 1334–1343.
80. V. Mezger et al. (1994) *Dev. Biol.* **166**, 819–822.
81. M. Rallu et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2392–2397.
82. A. Mathew, S. K. Mathur, and R. I. Morimoto (1998) *Mol. and Cell Biol.* **18**, 5091–5098.
83. J. S. Larson, T. J. Schuetz, and R. E. Kingson (1988) *Nature* **335**, 372–375.
84. D. D. Mosser, P. T. Kotzbauer, K. D. Sarge, and R. I. Morimoto (1990) *Proc. Natl. Acad. Sci.*

USA **87**, 3748–3752.

85. V. Zimarino, S. Wilson, and C. Wu (1990) *Science* **249**, 546–549.
86. J. T. Westwood, J. Clos, and C. Wu (1991) *Nature* **353**, 822–827.
87. U. Knauf, E. M. Newton, J. Kyriakis, and R. E. Kingson (1996) *Genes Dev.* **10**, 2782–2793.
88. M. P. Kline and R. I. Morimoto (1997) *Mol. Cell Biol.* **17**, 2107–2115.
89. D. A. Jurivich, L. Sistonen, R. A. Kroes, and R. I. Morimoto (1992) *Science* **255**, 1243–1245.
90. C. Giardina and J. T. Lis (1995) *J. Biol. Chem.* **270**, 10369–10372.
91. J. J. Cotto, M. Kline, and R. I. Morimoto (1996) *J. Biol. Chem.* **271**, 3355–3358.
92. K. Abravaya, B. Phillips, and R. I. Morimoto (1991) *Genes Dev.* **5**, 2117–2127.
93. D. D. Mosser, N. G. Theodorakis, and R. I. Morimoto (1988) *Mol. Cell Biol.* **8**, 4736–4744.
94. N. G. Theodorakis and R. I. Morimoto (1987) *Mol. Cell Biol.* **7**, 4357–4368.
95. B. J. DiDomenico, G. Bugaisky, and S. Lindquist (1982) *Cell* **31**, 593–603.
96. R. P. Beckmann, L. A. Mizzen, and W. J. Welch (1990) *Science* **248**, 850–854.
97. D. E. Stone and E. A. Craig (1990) *Mol. Cell Biol.* **10**, 1622–1632.
98. W. R. Boorstein and E. A. Craig (1990) *Mol. Cell Biol.* **10**, 3262–3267.
99. R. J. Nelson, M. F. Heschl, and E. A. Craig (1992) *Genetics* **131**, 277–285.
100. J. T. Halladay and E. A. Craig (1995) *Mol. Cell Biol.* **15**, 4890–4897.
101. K. Abravaya, M. P. Myers, S. P. Murphy, and R. I. Morimoto (1992) *Genes Dev.* **6**, 1153–1164.
102. R. Baler, W. J. Welch, and R. Voellmy (1992) *J. Cell Biol.* **117**, 1151–1159.
103. S. K. Rabindran et al. (1994) *Mol. Cell Biol.* **14**, 6552–6560.
104. S. L. Nunes and S. K. Calderwood (1995) *Biochem. Biophys. Res. Commu.* **213**, 1–6.
105. Y. Shi, D. D. Mosser, and R. I. Morimoto (1998) *Genes Dev.* **12**, 654–666.
106. T. Zhong, M. M. Luke, and K. T. Arndt (1996) *J. Biol. Chem.* **271**, 1349–1356.
107. K. Nadeau, A. Das, and C. T. Walsh (1993) *J. Biol. Chem.* **268**, 1479–1487.
108. S. C. Nair et al. (1996) *Cell Stress Chap.* **1**, 237–250.
109. P. B. Mason Jr. and J. T. Lis (1997) *J. Biol. Chem.* **272**, 33227–33233.
110. S. P. Bohlen and K. R. Yamamoto (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 313–334.
111. T. R. Hupp, D. W. Meek, C. A. Midgley, and D. P. Lane (1992) *Cell* **71**, 875–886.
112. F. U. Hartl (1996) *Nature* **381**, 571–579.
113. S. Rudiger, L. Germeroth, J. Schneider-Mergener, and B. Bukau (1997) *EMBO J.* **16**, 1501–1507.
114. A. L. Horwich et al. (1993) *Cell* **74**, 909–917.
115. M. Mayhew et al. (1996) *Nature* **379**, 420–426.
116. M. J. Gething (1997) *Molecular Chaperones and Protein-Folding Catalysts*, Oxford University Press, New York, pp. 18–21.
117. E. A. Craig et al. (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 31–52.
118. T. Langer and W. Neupert (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 53–84.
119. B. Miao, J. Davis, and E. Craig (1997) *Molecular Chaperones and Protein-Folding Catalysts*, Oxford University Press, New York, pp. 3–13.
120. C. Georgopoulos, K. Liberek, M. Zylicz, and D. Ang (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 209–249.

121. A. Wawrzynow and M. Zylicz (1997) *Molecular Chaperones and Protein-Folding Catalysts*, Oxford University Press, New York, pp. 481–483.
122. D. B. McKay et al. (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 153–177.
123. X. Zhu et al. (1996) *Science* **272**, 1606–1614.
124. B. C. Freeman, M. P. Myers, R. Schumacher, and R. I. Morimoto (1995) *EMBO J.* **14**, 2281–2292.
125. K. Liberek, Skowyra, M. Zylicz, and C. Georgopoulos (1991) *J. Biol. Chem.* **266**, 14491–15596.
126. D. R. Palleros, W. J. Welch, and A. L. Fink (1991) *Proc. Natl. Acad. Sci. USA* **88**, 5719–5723.
127. K. Liberek et al. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2874–2878.
128. S. Sadis and L. E. Hightower (1992) *Biochemistry* **31**, 9406–9412.
129. M. C. O'Brian and D. B. McKay (1995) *J. Biol. Chem.* **270**, 2247–2250.
130. H. Schroder, T. Langer, F. U. Hartl, and B. Bukau (1993) *EMBO J.* **12**, 4137–4144.
131. B. C. Freeman and R. I. Morimoto (1996) *EMBO J.* **15**, 2969–2979.
132. U. Jakob, M. Gaestel, K. Engel, and J. Buchner (1993) *J. Biol. Chem.* **268**, 1517–1520.
133. P. Goloubinoff, J. T. Christeller, A. A. Gatenby, and G. H. Lorimer (1989) *Nature (London)* **342**, 884–889.
134. D. F. Smith (1997) *Molecular Chaperones and Protein-Folding Catalysts*, Oxford University Press, New York, pp. 518–521.
135. W. B. Pratt and M. J. Welsh (1994) *Semin. Cell Biol.* **5**, 83–93.
136. R. J. Nelson et al. (1992) *Cell* **71**, 97–105.
137. W. B. Pratt, U. Gehring, and D. O. Toft (1997) *Stress-Inducible Cellular Responses*, Birkhauser Verlag, Basel, pp. 79–95.
138. O. Pinhasi-Kimhi, D. Michalovitz, A. Ben-Zeev, and M. Oren (1986) *Nature* **320**, 182–184.
139. R. I. Morimoto et al. (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 417–455.
140. U. Feige, R. I. Morimoto, I. Yahara, and B. S. Polla (1996) *Stress-Inducible Cellular Responses*, Birkhauser Verlag, Basel.
141. U. Feige and W. Van Eden (1997) *Stress-Inducible Cellular Responses*, Birkhauser Verlag, Basel, pp. 359–374.
142. I. J. Benjamin and R. S. Williams (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 533–552.
143. T. S. Nowak and H. Abe (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 553–575.
144. M. M. Patino, J. J. Liu, J. R. Glover, and S. Lindquist (1996) *Science* **273**, 622–626.
145. R. Mestril et al. (1994) *J. Clin. Invest.* **93**, 759–767.
146. M. S. Marber et al. (1995) *J. Clin. Invest.* **95**(4), 1446–1456.
147. J. C. L. Plumier et al. (1995) *J. Clin. Invest.* **95**, 1854–1860.
148. S. D. Morris, D. V. Cumming, D. S. Latchman, and D. M. Yellon (1996) *J. Clin. Invest.* **97**, 706–712.
149. L. Vigh et al. (1997) *Nat. Med.* **3**(10), 1150–1154.
150. B. Lee et al. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7207–7211.
151. M. G. Santoro (1997) *Trends Microbiol.* **5**, 276–281.
152. C. Amici, A. Rossi, and M. G. Santoro (1995) *Cancer Res.* **55**(19), 4452–4457.
153. A. Rossi, G. Elia, and M. G. Santoro (1998) *J. Biol. Chem.* **273**, 16446–16452.

## Hedgehog Signaling

Nusslein-Volhard and Weischaus (1) first identified the *hedgehog* (*hh*) **gene** in their now famous genetic screen that generated many [mutations](#) affecting segment number and polarity in *Drosophila*. Since then, most of the genes discovered in that screen have been **cloned**, and the studies elucidating their functions have provided an in-depth understanding of the molecular mechanisms driving the differentiation of cells with developmentally equivalent properties. The *hh* gene encodes a secreted [morphogen](#) that induces an intracellular signaling cascade. This cascade is an integral part of the mechanisms that establish the polarity of the segments during embryogenesis and limb development. *Hh* signaling is well conserved among vertebrates as well, and genes homologous to *hh* have been cloned and studied in various groups, including human, mouse, chicken, [Xenopus](#), and [zebrafish](#) (2).

### 1. Function of *hedgehog* Signaling in *Drosophila* Development

The *hh* gene plays important roles in both embryogenesis and disc development in *Drosophila melanogaster*. Embryogenesis of *Drosophila* is a dynamic process involving intricate hierarchical interactions among different gene products whose expression is defined in a temporally and spatially regulated manner. The maternal and zygotic gap gene products lay out the domains from which the insect metamer units evolve. The [pair-rule genes](#), together with the maternal and gap gene activities, establish the repeating embryonic parasegments. The segment polarity genes then specify the anterior and posterior positional information within each evolving segment. Differential expression of the **homeotic** genes expressed within each segment establishes the unique identities of the segments. All of these interactions occurring during embryogenesis define the body plan and serve as a blueprint for pattern formation in latter development.

*Hh* and **wingless** (*wg*, a gene that encodes another signaling morphogen) are two segment polarity genes whose main function is to define the anterior/posterior polarity within each of the embryonic segments. Embryos that lack either *hh* or *wg* function are smaller than wild-type embryos, secrete cuticles that are covered with a lawn of denticles, and do not make naked cuticle (1). This phenotype characterizes many of the segment polarity genes involved in transducing either the *hh* or *wg* signals. The expression of *hh* and *wg* is initiated at the blastoderm by pair-rule gene products in the posterior and anterior compartment of the segments, respectively, and is maintained during gastrulation by interactions among the segment polarity gene products (3, 4). Thus, the loss of *wg* function does not affect initial *hh* expression but causes the gradual loss of *hh* expression at the onset of gastrulation (5-7). Similarly, the loss of *hh* function does not affect the establishment of early *wg* gene expression, but at later stages it results in the loss of *wg* expression along the posterior edge of the anterior/posterior (A/P) segmental boundary (8, 9). Although *hh* is expressed in only two rows of cells that flank the A/P boundary, *hh* mutations affect the development of the entire segment. Genetic studies have demonstrated that these long-range effects of *hh* are mediated through secondary signals such as *wg* and **decapentaplegic** (*dpp*) (10-13).

The *hh* signaling pathway is one of the best-studied cascades involved in embryogenesis and limb development, and it is well conserved between insects and vertebrates. In *Drosophila*, the components identified in this pathway include the putative *hh* receptor *patched* (*ptc*) (14, 15), the coreceptor *smoothed* (*smo*) (16, 17), the transcription factor **cubitus interruptus** (*ci*) (18), the **serine/threonine kinases** protein kinase A (PKA) (19) and *fused* (*fu*) (20), *suppressor of fused* (*su*) (*fu*) (21), *costal-2* (*cos -2*) (20), and *slimb* (*slmb*) (22), which encodes a protein involved in the **protein degradation** of CI, the protein product of the *ci* gene.

In a developing embryo, *hh* is secreted from the posterior compartment of each segment and binds to

its postulated receptor, *ptc*, which is found on the surface of the cells anterior and adjacent to the *hh*-secreting cells. The *hh* precursor protein, HH, is composed of 471 amino acid residues and consists of three functional **domains**: the [signal peptide](#) domain (residues 63–83), the N terminus, and the C terminus. The signal peptide is involved in secretion via the [endoplasmic reticulum](#) and is cleaved on secretion at the cell surface (5). The 46-kDa full-length HH precursor (F) is not detected in *Drosophila* embryos or in cell culture, indicating that the signal peptide is cleaved very efficiently on biosynthesis of HH. The resulting 39-kDa HH protein is called Uncleaved (U) and is detected in *Drosophila* embryos and discs and in cell culture. U undergoes further internal **proteolytic** cleavage, generating a 25-kDa C-terminal protein and the 19-kDa N-terminal fragment that is the active HH molecule responsible for both the short- and long-range *hh* activities. This cleavage is intramolecular and is catalyzed by sequences in the C terminus that are conserved across species and display some similarity to those of [serine proteinases](#) (12, 23). Mutation of the invariant [histidine](#) residue in this region abolishes HH autocleavage, indicating that the activity of the C terminus is crucial for HH processing (23). The autocleavage occurs between Gly257 and Cys258. In the course of the catalyzed reaction, the nucleophile that attacks the Gly257 carboxyl group is not the usual [water](#) molecule, but cholesterol, which consequently remains covalently attached to the N-terminal HH peptide after cleavage. The presence of cholesterol increases the concentration of biologically active HH at the surface of cells secreting HH (24, 25). How the cholesterol-coupled form of HH exerts its short- and long-range effects is not clear. The current hypothesis is that the short-range effect is mediated directly by HH, and that the long-range effect is mediated by the molecules whose expression is induced by HH signaling.

Both genetic and biochemical evidence suggests that the transmembrane protein PTC, the product of the *patched* gene, is the receptor for HH (14, 15, 26). However, *ptc* mutations cause the expression of *hh* target genes—including *wg*, *dpp*, and *ptc* itself—to be up-regulated and to become independent of an *hh* signal (11, 13, 27-29). This result suggests that *ptc* activity antagonizes the *hh* [signal transduction](#) cascade in the absence of an *hh* signal. The gene *smo* also encodes a transmembrane protein and, like *hh*, is a positive regulator of *wg* expression (16, 17). However, SMO does not appear to bind to HH (30). The simplest explanation for these findings is that *ptc* inhibits *smo* activity in the absence of an *hh* signal. When HH binds to PTC, the *ptc* inhibition of *smo* is relieved, and the intracellular cascade can induce *hh* target gene expression. That SMO and PTC **immunoprecipitate** together from cell culture in the presence or absence of HH suggests that PTC-SMO interactions are determined by a conformational change, rather than a dissociation event (31). The activation of *smo* does not require *ptc*, as SMO is constitutively active in a *ptc* mutant background (28, 30).

Genetic [epistasis](#) experiments have demonstrated that *Su(fu)*, *fu*, *cos-2*, and *ci* act downstream from *ptc* and are responsible for transducing the *hh* signal into the [nucleus](#), where *ci* activates the *hh* target genes (8, 20, 32). A perplexing phenomenon associated with *ci* as a [transcription factor](#) is that its full-length activated form has never been detected in the nucleus in vivo, not even in the nuclei of cells located at the A/P boundary, where *ci* exerts its transactivation function (18), or in mutants (eg, *ptc*, *cos-2*, PKA mutants) that exhibit elevated levels of CI protein (29, 33). The accepted explanation for this phenomenon is that cells require undetectable amounts of CI in the nucleus to activate the expression of target genes. CI exists in the cytoplasm as a complex consisting of COS-2, FU, and other unidentified proteins. In the absence of an *hh* signal, this complex is mainly tethered to the cytosolic microtubules (see [Microtubule-Associated Proteins \(MAPs\)](#)) (34-36). CI also exists as two forms: a 155-kDa transcriptionally active form and a 75-kDa N-terminal proteolytic fragment that contains the zinc-finger **DNA-binding** domain. Both forms are detected in the cytosol. In the absence of an *hh* signal, CI is degraded, and the N-terminal fragment of CI is detected in the nucleus, where it functions as a [repressor](#), possibly because it fails to recruit its co-activator *Drosophila* CBP (cyclic-AMP response element binding protein) (*dCBP*) (34, 37). When the cells receive an *hh* signal, CI dissociates from the microtubule-protein complex, is not degraded, and the full-length form of CI activates *hh* target gene expression (34-36).

A large body of evidence has demonstrated that protein kinase A (PKA) antagonizes *hh* signaling.



PKA mutant clones in the anterior compartment of imaginal discs express *hh* target genes and generate the pattern duplications associated with ectopic *dpp* and *wg* expression (38, 39). As expected, increasing the PKA activity by expressing a constitutively active PKA catalytic subunit antagonizes the *hh*-induced target gene expression along the A/P boundary of the imaginal disc (29, 40). In embryos, PKA has a dual effect on *hh* signaling. Whereas loss of PKA activity leads to increased *hh* target gene expression, the constitutive activity of PKA causes up-regulation of *wg* expression in a *smo*- and *ci*-dependent manner, suggesting that PKA also has a positive function in *hh* signaling during embryogenesis (41). Because an increase in PKA activity does not interfere with *hh* function in either embryogenesis or disc development, it is unlikely that *hh* regulates PKA activity (29, 40, 41). *Ci* is at least one of the direct targets of PKA regulation in the *hh* signaling pathway. In the absence of an *hh* signal, cell culture studies show that CI is **phosphorylated** by the basal PKA activity, and this CI phosphorylation of CI leads to the proteolysis of CI to the N-terminal repressor form (42). When cells receive an *hh* signal, the proteolysis of CI is inhibited, resulting in elevated levels of activated full-length CI. This event leads to the transcriptional activation of target genes in the nucleus, presumably through the interaction with its co-activator *dCBP* (37).

The *hh* gene is also important for the patterning and growth of adult structures. The adult wings and legs are derived from pouches of larval epithelial cells called [imaginal discs](#). Imaginal discs are composed of anterior and posterior compartments derived from groups of cells that originate from the anterior and posterior compartment of certain embryonic segments. Because homozygous *hh* mutants die before hatching, knowledge of *hh* function in patterning the adult appendages was obtained from flies that ectopically express *hh* or from the clonal analysis of *hh* mutant cells. The ectopic expression of *hh* in the anterior compartment of wing and leg discs leads to the duplication of anterior disc patterns or the development of supernumerary wings or legs, depending on the position of the HH-producing clones (10, 43, 44). Loss of endogenous *hh* function in wing and leg discs is associated with loss of patterning in both anterior and posterior compartments and loss of proximal-distal outgrowth that results in smaller appendages (10). The changes in compartment and proximal-distal patterning are mediated by *hh*-induced *dpp* and *wg* (10, 44, 45). In the wing disc, *dpp* is normally expressed along the A/P boundary. Its expression is induced in clones of anterior cells that ectopically express HH, and its expression is abolished in the absence of *hh* function (10, 13, 43). Cells located dorsally in the leg disc express high levels of *dpp*, whereas cells along the ventral A/P boundary express low levels of *dpp* and high levels of *wg*. Dorsal-anterior HH-producing clones induce ectopic *dpp* expression, whereas ventral-anterior HH-producing clones induce ectopic *wg* expression. Similarly, loss of function of *hh* in posterior cells results in a size reduction of the adult leg (10, 45). *Hh* functions similarly in other imaginal discs, such as genital and eye-antennal discs. Loss or gain of function of *hh* is associated with loss or gain of genital and eye-antennal patterns, which is most likely mediated through *dpp* (44, 46, 47).

In the eye disc, pattern formation occurs as a wave of differentiation, the Morphogenetic Furrow (MF), that moves across the disc from posterior to anterior. Undifferentiated cells lie in front of the MF, whereas a patterned array of differentiated photoreceptors and neurons remains behind. HH expression in the cells posterior to the MF drives the progression of the MF, but the mediator of the *hh* signal in the eye-antennal disc is not known. *Dpp* promotes ommatidial differentiation, and its expression is under the control of *hh* before MF formation, after which time, *dpp* autoregulates its own expression and suppresses *wg* expression to promote differentiation behind the MF (5, 23, 48-50).

## 2. Function of *hedgehog* in Vertebrate Development

Vertebrates carry multiple genes homologous to *hh*, including ***Sonic hedgehog*** (*Shh*), ***Indian hedgehog*** (*Ihh*), and ***Desert hedgehog*** (*Dhh*) (51). *Ihh* has been identified in mice, humans, and chicken, and its expression is detected in the developing midgut, lung, and cartilage of the bones. During bone morphogenesis, *Ihh* is expressed in prehypertrophic chondrocytes and controls the balance between ossification and proliferation through a second signal, parathyroid-related protein (PTHrP), which is secreted by perichondrial cells. Ectopic expression of *Ihh* leads to the ectopic

induction of PTHrP. Misexpression of *Ihh* or PTHrP results in the disruption of bone development: cartilage elements are broadened, and the bones fail to ossify. On the other hand, premature ossification and insufficient growth of the bones are associated with loss of function of PTHrP or PTH/PTHrP receptor. The genetic analysis of PTHrP-mutant mice indicates that PTHrP is epistatic to *Ihh* and demonstrates that the differentiating prehypertrophic chondrocytes express *Ihh* to induce PTHrP expression in the adjacent perichondrium. The induced PTHrP signal acts as a negative-feedback cue on chondrocytes expressing PTH/PTHrP receptors, to ensure the proliferative status of the cells, and suppresses their early differentiation (52, 53). The mouse *Dhh* gene is involved in spermatogenesis and is proposed to support the survival of sperm cells (54).

### 2.1. *Sonic hedgehog* (*Shh*)

*Shh* is the best-studied *hedgehog* molecule because it is expressed early in embryogenesis and patterns multiple processes during development. These include the establishment of left-right asymmetry in early embryogenesis, specification of ventral neuron fate in the central nervous system, patterning of the dorsal-ventral axis of somites, establishment of the anterior-posterior axis, outgrowth of the limb buds, and eye development. *Shh* expression is first detected in the 8.0 days postcoitum (dpc) mouse embryo, prior to the first somite formation, in the node and the node-derived cells of the head process, which underlies the neural plate. As embryogenesis progresses and the embryo axis extends caudally, *Shh* is detected in the notochord, an extension of the head process, and in the floor plate cells overlying the notochord in the ventral midline of the neural tube. The timing of *Shh* expression in these midline structures matches the inductive activities manifested by the notochord and the floor plate (51, 55). Misexpression and explantation experiments have demonstrated that *Shh* secreted from the notochord induces the differentiation of floor plate cells and that *Shh* from both the notochord and the floor plate specifies the ventral neuronal differentiation in the brain and the motor neuron fate in the spinal cord (56-58). Misexpression of *Shh* by the Wnt-1 promoter in mouse embryos results in ectopic floor plate cell differentiation in the dorsal neural tube, as demonstrated by the expression of the floor plate cell marker, *HNF-3b* (hepatocyte nuclear factor 3b) (51). Treatment of neural tube explants with *Shh* conditioned medium also leads to induction of *HNF-3b* (56). Ablation of notochord, loss of *Shh* function, or antibody blockage experimentation provides further evidence that *Shh*, synthesized in the underlying notochord, is required for the induction of floor plate cell differentiation (58-60). The LIM (*lin-11*, *Isl1*, *mec-3*) homeodomain protein, *Islet-1* (*Isl-1*), is an early ventral neuron marker. Neural tube explants, grown in contact with floor plate tissue or in *Shh*-conditioned medium, acquire ventral neuron fate by expressing *Isl-1*. However, explants from various levels of the neuraxis acquire different identities and express a number of different genes in addition to *Isl-1*, depending on their position along the neural axis. For example, when treated with *Shh*, neural plate cells caudal to the forebrain (rhombencephalon and spinal cord) express both *Isl-1* and *SCI*, which is characteristic of motor neurons, whereas explants from forebrain express *Isl-1* and *Nkx-2.1*, markers that are characteristic of ventral forebrain neurons. The ability of *Shh* to induce different types of ventral neuron fates along the rostrocaudal axis of the neural tube suggests that an early pre-pattern of positional information exists that dictates the differential response to *Shh* (56, 57).

In addition to the patterning of CNS, *Shh* also specifies the dorsal-ventral polarity of somites. Somites originate as spherical balls of epithelial cells derived from paraxial mesoderm flanking the notochord and neural tube. As development proceeds, somites undergo dorsal and ventral differentiation, giving rise to ventral-medial mesenchymal sclerotome and dorsal-lateral epithelial dermomyotome. The sclerotome differentiates into axial skeleton including the vertebral column and ribs. The dorsal dermomyotome is subdivided into medial and lateral compartments, which give rise to epaxial (deep back) and hypaxial (limb and body wall) muscles, respectively. The dermis of the back is derived from the superficial epithelial cells that comprise the dermatome. *Shh* is required for both sclerotome differentiation and dermomyotome proliferation. Grafting notochord or floor plate on the dorsal somite enhances sclerotome formation and suppresses dermomyotome formation, suggesting that a substance secreted from the notochord and floor plate specifies ventral sclerotome differentiation (61, 62). Misexpression of *Shh* in the chick dorsal somite results in the dorsal expansion of *Pax-1*, a sclerotomal marker, and down-regulation of the dermomyotomal marker *Pax-*

3 (63) (see [Pax Genes](#)). Furthermore, mouse explants co-cultured with *Shh*-expressing cells have increased levels of *Pax-1* and decreased *Pax-3* expression, supporting the hypothesis that *Shh* from axial structures induces ventral sclerotomal cell fate (64). The differentiation and survival of dorsal-medial myotome cells require the inductive functions of both *Shh* from notochord/floor plate and the *wingless* homologue *Wnt* gene products that are derived from the dorsal neural tube (65-68).

The vertebrate limb evolves from both the lateral plate and the lateral edge of the somite. Mesenchymal cells from the dorsal-lateral somite adjacent to the presumptive limb bud migrate laterally and give rise to the future limb musculature, and cells from the lateral plate proliferate and develop into the remaining limb tissues, including the bones and tendons. Limb development is a dynamic process orchestrated by multiple inductive signals from the posteriorly located zone of polarization activity (ZPA), the dorsal ectoderm, and the apical ectodermal ridge (AER) that abuts the anterior-posterior axis of the limb, at the distal tip where the dorsal and ventral surfaces meet. *Shh* is expressed in the ZPA and has been shown to establish anterior-posterior polarity in the limb and promote limb outgrowth (69). However, *Shh* alone is not sufficient to induce the expression of its target genes, *Bmp-2* and *Hox* gene cluster, which transduce the *Shh* signal to regulate limb patterning and growth (69, 70). The induction of these genes requires an additional signal, possibly **fibroblast growth factor-4** (*FGF-4*), that is secreted from the overlying AER (71). *Shh* polarizes and maintains a functional AER. In a similar fashion, AER maintains *Shh* expression in the ZPA, establishing a positive feedback loop and an interdependence between these two inductive signaling centers (71). *Shh* and *FGFs* act together to induce the expression of *Bmp-2*, which relays a signal to induce or maintain *Hox* gene expression (71). Examination of lung and eye development indicates that *Shh* is also an inductive and mitogenic signal governing the growth and differentiation in these organs (72, 73).

### 3. Comparison of *hedgehog* Signaling in Flies and Vertebrates

The *hh* signaling cascade is well conserved between vertebrates and flies. Homologues of both *ptc* and *smo* have been identified in mammals, and their functions are conserved. Mutations in the mammalian *ptch* and *smo* are associated with familial and sporadic basal cell carcinoma, suggesting that these two homologues function as **tumor suppressor** genes (31, 74-76). Three homologues of *ci* exist in vertebrates: *Gli1*, *Gli2*, and *Gli3*, which exhibit distinct and overlapping expression patterns during development (77-80). Experiments with *Gli1* antisense and *Gli1* ectopic expression in different systems have shown that *Gli1* is the major transcriptional activator that mediates the *Shh* signaling pathway, whereas *Gli3* functions as a repressor (81-85). The function of *Gli2* overlaps both *Gli1* and *Gli3* in different developmental processes (85, 86). Mutations in *Gli1* are associated with glioblastoma and skin cancer (87), and those in the *Gli3* gene cause a dominant extra-toe phenotype in mice and are associated with Greg cephalopolysyndactyly syndrome in humans (88-90), suggesting that *Gli1* and *Gli3* function as tumor suppressors. The negative regulation of *Shh* signaling pathways by PKA is also conserved in vertebrates. Ectopic expression of a dominant-negative PKA mutant in the murine dorsal neural tube results in the ectopic expression of *Shh* target genes such as *HNF-3b*, *path*, and *Gli1*, and the suppression of the dorsal neuronal cell marker *Pax-3* (91). PKA activity has also been shown to antagonize the *Shh* signaling pathway in developing somites (64). Although the *hh* signaling pathways are conserved in most of the systems examined, there are differences between the vertebrate and *Drosophila* systems. The regulation of CI proteolysis in flies makes CI both an activator and a repressor, whereas in vertebrates it appears that the activation and repression functions are separated into two distinct genes, namely *Gli1* and *Gli3*, respectively. The level of the activated form of CI is regulated post-translationally in flies, whereas the *Gli* levels are regulated transcriptionally in vertebrates. Continued analysis of these important regulatory cascades will undoubtedly increase our understanding of developmental signaling pathways, their similarities, and their evolutionary divergence.

### Bibliography

1. C. Nusslein-Volhard and E. Wieschaus (1980) *Nature* **287**, 795–801.

2. M. Hammerschmidt, A. Brook, and A. P. McMahon (1997) *Trends Genet.* **13**, 14–21.
3. S. DiNardo, E. Sher, J. Heemskerk-Jongens, J. A. Kassis, and P. H. O'Farrell (1988) *Nature* **332**, 604–609.
4. A. M. Arias, N. E. Baker, and P. W. Ingham (1988) *Development* **103**, 157–170.
5. J. J. Lee, D. P. V. Kessler, S. Parks, and P. A. Beachy (1992) *Cell* **71**, 33–50.
6. J. Mohler and K. Vani (1992) *Development* **115**, 957–971.
7. T. Tabata, S. Eaton, and T. B. Kornberg (1992) *Genes Dev.* **6**, 2635–2645.
8. P. W. Ingham, A. M. Taylor, and Y. Nakano (1991) *Nature* **353**, 184–187.
9. P. W. Ingham, (1993) *Nature* **366**, 560–562.
10. K. Basler and G. Struhl, (1994) *Nature* **368**, 208–214.
11. J. Capdevila and I. Guerrero (1994) *EMBO J.* **13**, 4459–4468.
12. J. A. Porter, D. P. v. Kessler, S. C. Ekker, K. E. Young, J. J. Lee, K. Moses, and P. A. Beachy (1995) *Nature* **374**, 363–366.
13. T. Tabata and T. B. Kornberg (1994) *Cell* **76**, 89–102.
14. J. E. Hooper and M. P. Scott (1989) *Cell* **59**, 751–765.
15. Y. Nakano, I. Guerrero, A. Hidalgo, A. Taylor, J. R. S. Whittle, and P. W. Ingham (1989) *Nature* **341**, 508–513.
16. M. v. d. Heuvel and P. W. Ingham (1996) *Nature* **382**, 547–551.
17. J. Alcedo, M. Ayzenzon, T. V. Ohlen, M. Noll, and J. E. Hooper (1996) *Cell* **86**, 221–232.
18. T. V. Orenic, D. C. Slusarski, K. L. Kroll, and R. A. Holmgren (1990) *Genes Dev.* **4**, 1053–1067.
19. D. Kalderon and G. M. Rubbin (1988) *Genes Dev.* **2**, 1539–1556.
20. A. J. Forbes, Y. Nakano, A. M. Taylor, and P. W. Ingham (1993) *Development (Supplement)*, 115–124.
21. A. Pham, P. Therond, G., Alves, F. B. Tounier, D. Busson, C. Lamour-Isnard, B. L. Bouchon, T. Preat, and H. Tricoire (1995) *Genetics* **140**, 587–598.
22. J. Jiang and G. Struhl (1998) *Nature* **391**, 493–496.
23. J. J. Lee, S. C. Ekker, D. P. v. Kessler, J. A. Porter, B. I. Sun, and P. A. Beachy (1994) *Science* **266**, 1528–1537.
24. J. A. Porter, S. C. Ekker, W.-J. Park, D. P. von-Kessler, K. E. Young, C.-H. Chen, Y. Ma, A. s. Woods, R. J. Cotter, E. V. Koonin, and P. A. Beachy (1996) *Cell* **86**, 21–34.
25. J. A. Porter, K. E. Young, and P. A. Beachy (1996) *Science* **274**, 255–259.
26. V. Marigo, R. A. Davey, Y. Zuo, J. M. Cunningham, and C. J. Tabin (1996) *Nature* **384**, 176–179.
27. A. Martinez-Arias, N. E. Baker, and P. W. Ingham (1988) *Development* **103**, 157–170.
28. J. Capdevila, M. P. Estrada, E. Sanchez-Herrero, and I. Guerrero (1994) *EMBO J.* **13**, 71–82.
29. W. Li, J. T. Ohlmeyer, M. E. Lane, and D. Kalderon (1995) *Cell* **80**, 553–562.
30. Y. Chen and G. Struhl (1996) *Cell* **87**, 553–563.
31. S. M. Stone, M. Hynes, M. Armanini, T. A. Swanson, Q. Gu, R. L. Johnson, M. P. Scott, D. Pennica, A. Goddard, H. Phillips, M. Noll, J. E. Hooper, F. d. Sauvage, and A. Rosenthal (1996) *Nature* **384**, 129–134.
32. C. K. Motzny and R. Holmgren (1995) *Mech. Dev.* **52**, 137–150.
33. J. Hepker, Q.-T. Wang, C. K. Motzny, R. Holmgren, and T. v. Orenic (1997) *Development* **124**, 549–558.
34. P. Aza-Blanc, F.-A. Ramirez-Webe, and T. B. Kornberg (1997) *Cell* **89**, 1043–1053.
35. D. J. Robbins, K. E. Nybakken, R. Kobayashi, J. C. Sisson, J. M. Bishop, and P. P. Therond (1997) *Cell* **90**, 225–234.

36. J. C. Sisson, K. S. Ho, K. Suyama, and M. P. Scott (1997) *Cell* **90**, 235–245.
37. H. Akimaru, Y. Chen, P. Dai, D.-X. Hou, M. Nonaka, S. M. Smolik, S. Armstrong, R. H. Goodman, and S. Ishii (1997) *Nature* **386**, 735–738.
38. T. Lepage, S. M. Cohen, F. J. Diaz-Benjumea, and S. M. Parkhurst (1995) *Nature* **373**, 711–715.
39. D. Pan and G. M. Rubin (1995) *Cell* **80**, 543–552.
40. J. Jiang and G. Struhl (1995) *Cell* **80**, 563–572.
41. J. T. Ohlmeyer and D. Kalderon (1997) *Genes Dev.* **11**, 2250–2258.
42. Y. Chen, N. Gallaher, R. H. Goodman, and S. M. Smolik (1998) *Proc. Natl. Acad. Sci. USA* **95**, 2349–2354.
43. T. Kojima, T. Michiue, M. Orihara, and K. Saigo (1994) *Gene* **148**, 211–217.
44. A. L. Felsenfeld and J. A. Kennison (1995) *Development* **121**, 1–10.
45. F. J. Diaz-Benjumea, B. Cohen, and S. M. Cohen (1994) *Nature* **372**, 175–179.
46. J. Mohler (1988) *Genetics* **120**, 1061–1072.
47. E. H. Chen and B. S. Baker (1997) *Development* **124**, 205–218.
48. J. Royet and R. Finkelstein (1997) *Development* **124**, 4793–4800.
49. m. Dominguez and E. Hafen (1997) *Genes Dev.* **11**, 3254–3264.
50. F. Chanut and U. Heberlein (1997) *Development* **124**, 559–567.
51. Y. Echelard, D. J. Epstein, B. St-Jacques, L. Shen, J. Mohler, J. A. McMahon, and A. P. McMahon (1993) *Cell* **75**, 1417–1430.
52. A. Vortkamp, K. Lee, B. Lanske, G. V. Segre, H. M. Kronenberg, and C. J. Tabin (1996) *Science* **273**, 613–621.
53. B. Lanske, A. C. Karaplis, K. Lee, A. Luz, A. Vortkamp, A. Pirro, M. Karperien, L. H. K. Defize, C. Ho, R. C. Mulligan, A.-B. Abou-Samra, H. Juppner, G. V. Segre, and H. M. Kronenberg (1996) *Science* **273**, 663–666.
54. M. J. Bitgood, L. Shen, and A. P. McMahon (1996) *Curr. Biol.* **6**, 298–304.
55. E. Marti, R. Takada, D. A. Bumcrot, H. Sasaki, and A. P. McMahon (1995) *Development* **121**, 2537–2547.
56. H. Roelink, J. A. Porter, C. Chiang, Y. Tanabe, D. T. Chang, P. A. Beachy, and T. M. Jessell (1995) *Cell* **81**, 445–455.
57. J. Ericson, J. Muhr, M. Piaczek, T. Lints, T. M. Jessell, and T. Edlund (1995) *Cell* **81**, 747–756.
58. J. Ericson, S. Morton, A. Kawakami, H. Roelink, and T. M. Jessell (1996) *Cell* **87**, 661–673.
59. R. i. Altaba (1992) *Development* **115**, 67–80.
60. C. Chiang, Y. Litingtung, E. Lee, K. E. Young, J. L. Corden, H. Westphal, and P. A. Beachy (1996) *Nature* **383**, 407–413.
61. O. Pourquie, M. Coltey, M.-A. Teillet, C. Ordahl, and N. M. Le-Douarin (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5242–5246.
62. M. Goulding, A. Lumsden, and A. J. Paquette (1994) *Development* **120**, 957–971.
63. R. L. Johnson, E. Laufer, R. D. Riddle, and C. Tabin (1994) *Cell* **79**, 1165–1173.
64. C.-M. Fan, J. A. Porter, C. Chiang, D. T. Chang, P. A. Beachy, and M. Tessier-Lavigne (1995) *Cell* **81**, 457–465.
65. M. E. Pownall, K. E. Strunk, and C. P. Emerson, Jr. (1996) *Development* **122**, 1475–1488.
66. A. E. Munsterberg, J. Kitajewski, D. A. Bumcrot, A. P. McMahon, and A. B. Lassar (1995) *Genes Dev.* **9**, 2911–2922.
67. E. Hirsinger, D. Duprez, C. Jouve, P. Malapert, J. Cooke, and O. Pourquie (1997) *Development* **124**, 4605–4614.
68. M.-A. Teillet, Y. Watanabe, P. Jeffs, D. Duprez, F. Lapointe, and N. M. L. Douarin (1998)

Development **125**, 2019–2030.

69. R. D. Riddle, R. L. Johnson, E. Laufer, and C. Tabin (1993) *Cell* **75**, 1401–1416.
70. Y. Yang, G. Drossopoulou, P.-T. Chuang, D. Duprez, E. Marti, D. Bumcrot, N. Vargesson, J. Clarke, L. Niswander, A. McMahon, and C. Tickle (1997) *Development* **124**, 4393–4404.
71. E. Laufer, C. E. Nelson, R. L. Johnson, B. A. Morgan, and C. Tabin (1994) *Cell* **79**, 993–1003.
72. S. Bellusci, Y. Furuta, M. G. Rush, R. Henderson, G. Winnier, and B. L. M. Hogan (1997) *Development* **124**, 53–63.
73. A. M. Jensen and V. A. Wallace (1997) *Development* **124**, 363–371.
74. L. V. Goodrich, R. L. Johnson, L. Milenkovic, J. A. McMahon, and M. P. Scott (1996) *Genes Dev.* **10**, 301–312.
75. V. Marigo, M. P. Scott, R. L. Johnson, L. V. Goodrich, and C. J. Tabin (1996) *Development* **122**, 1225–1233.
76. J. Xie, M. Murone, S.-M. Luoh, A. Ryan, Q. Gu, C. Zhang, J. M. Bonifas, C.-W. Lam, M. Hynes, A. Goddard, A. Rosenthal, E. H. Epstein Jr., and F. J. d. Sauvage (1998) *Nature* **391**, 90–92.
77. K. W. Kinzler, J. M. Ruppert, S. H. Bigner, and B. Vogelstein (1988) *Nature* **332**, 371–374.
78. D. C. Hughes, J. Allen, G. Morley, K. Sutherland, W. Ahmed, J. Prosser, I. Lettice, G. Allan, M.-G. Mattei, M. Farrall, and R. E. Hill (1997) *Genomics* **39**, 205–215.
79. J. M. Ruppert, B. Vogelstein, K. Arheden, and K. W. Kinzler (1990) *Mol. Cell. Biol.* **10**, 5408–5415.
80. C.-c. Hui, D. Slusarski, K. A. Platt, R. Holmgren, and A. L. Joyner (1994) *Dev. Biol.* **162**, 402–413.
81. J. Lee, K. A. Platt, P. Censullo, and A. R. i. Altaba (1997) *Development* **124**, 2537–2552.
82. M. Hynes, D. M. Stone, M. Dowd, S. Pitts-Meek, A. Goddard, A. Gurney, and A. Rosenthal (1997) *Neuron* **19**, 15–26.
83. H. Masuya, T. Sagai, K. Moriwaki, and T. Shiroishi (1997) *Dev. Biol.* **182**, 42–51.
84. H. Sasaki, C.-c. Hui, M. Nakafuku, and H. Kondoh (1997) *Development* **124**, 1313–1322.
85. A. R. i. Altaba (1998) *Development* **125**, 2203–2212.
86. Q. Ding, J. Motoyama, S. Gasca, R. Mo, H. Sasaki, J. Rossant, and C.-c. Hui (1998) *Development* **125**, 2533–2543.
87. N. Dahmane, J. Lee, P. Robins, P. Heller, and A. R. i. Altaba (1997) *Nature* **389**, 876–881.
88. T. M. Pohl, M.-G. Mattei, and U. Ruther (1990) *Development* **110**, 1153–1157.
89. C.-c. Hui and A. L. Joyner (1993) *Nature Genetics* **3**, 241–246.
90. A. Wild, M. Kalff-Suske, A. Vortkamp, D. Bornholdt, R. Konig, and K.-H. Grzeschik (1997) *Hum. Mol. Gen.* **6**, 1979–1984.
91. D. J. Epstein, E. Marti, M. P. Scott, and A. P. McMahon (1996) *Development* **122**, 2885–2894.

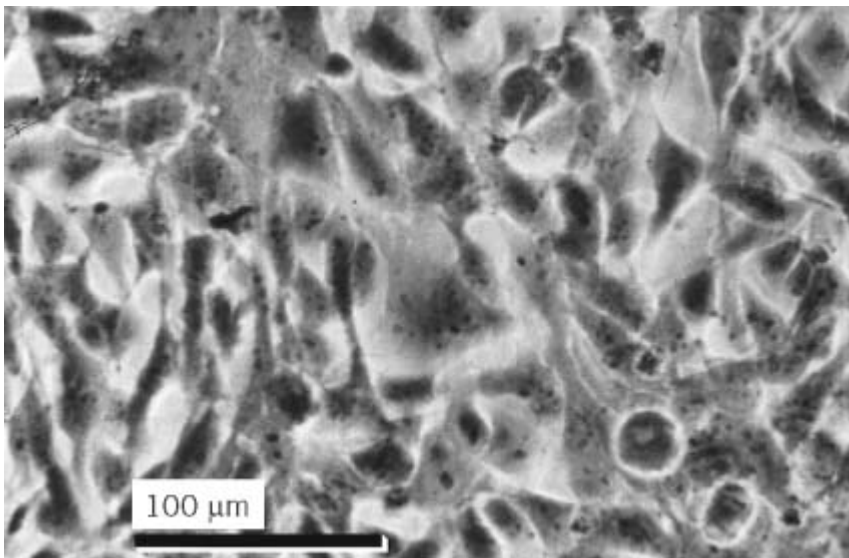
### **Suggestions for Further Reading**

92. M. Hammerschmidt, A. Brook, and A. P. McMahon (1997) The world according to *hedgehog*. *Trends Genet.* **13**, 14–21.
93. P. W. Ingham (1998) Transducing hedgehog: the story so far. *EMBO J.* **17**, 3505–3511.
94. R. L. Johnson and C. J. Tabin (1997) Molecular models for vertebrate limb development. *Cell* **90**, 979–990.

## HeLa Cells

The HeLa cell line was originated by George Gey and co-workers in 1952 (1) from a cervical adenocarcinoma from a 31-year-old Afro-American female. It has proved to be a very vigorously growing [cell line](#) and is one of the most extensively used cell lines in the world. It is lodged at the American Type Culture Collection (ATCC) as CCL-2. It is even more extensively used than many people realize, because it is also guilty of cross-contaminating more cell lines than any other. HeLa was also the first human cell line to be cloned, when HeLa-S3 was isolated by Puck in 1956 (2). HeLa-S3 (Fig. 1) is in the ATCC catalogue as CCL-2.2. HeLa cells were originally propagated in a combination of human and calf serum (see [Serum Dependence](#)), but it has been shown that human serum is not necessary, and they are usually maintained in Eagle's Minimal Essential Medium supplemented with 10% fetal bovine serum. HeLa-S3 can be propagated in suspension using low-calcium medium (3).

**Figure 1.** Sub-confluent culture of HeLa-S3 cells. Phase contrast, Olympus CK microscope, 20× objective.



### 1. Properties

HeLa is a continuous cell line with an **aneuploid** and heteroploid [karyotype](#) containing a number of marker [chromosomes](#) (4). It expresses the A subtype of glucose-6-phosphate dehydrogenase, which is present in Negroes and absent in Caucasians. The cells are [cytokeratin](#) positive, confirming their epithelial origin, contain **papilloma virus** sequences (HPV-18), and have low expression of [p53](#) and normal **retinoblastoma** proteins (5). They have a short population doubling time (18–22 h) and high plating efficiency (60–80% in monolayer).

### 2. Usage

HeLa cells have been used for propagation and assay of **viruses** and have been shown to be susceptible to **poliovirus-1,2,3**, **adenovirus-3**, encephalomyocarditis, and coxsackievirus B<sup>5</sup>. They have been used in molecular studies (6, 7), in somatic hybridization analysis of senescence (8) and invasion (9), in fermentor culture for biotechnology (10), and for the development of serum-free media (11). As they are deficient in [gap-junction](#) communication, they have been used for

[transfection](#) analysis of gap-junctional proteins (12). They have also been used for studies on molecular signaling (13).

HeLa-S3, a clone of HeLa (2), is one of the more rapidly dividing of the HeLa strains and can be grown in suspension with mechanical agitation. Because it is less well attached than the parental strain, mitotic cells are readily detached by shaking the flask, following partial cell-cycle blockade induced by short-term storage at 4°C, (14, 15). H-HeLa is a subline suitable for the passage and titration of [rhinoviruses](#) (16).

### 3. Cross-Contamination

HeLa was the first human cell line to be propagated widely in a large number of laboratories throughout the world. It is also a vigorously growing cell line with a high plating efficiency and is, therefore, ideally placed to cross-contaminate other cell lines. That this happened was first demonstrated many years ago by **isoenzyme** analysis (17) and subsequently by chromosomal analysis (18). Combined with more recent studies of DNA **fingerprinting**, it is now clear that a large number of continuous cell lines established in the 1960s and 1970s were cross-contaminated with HeLa cells, including lines such as KB (19) and Hep-2 (19, 20), which are still in regular use. The most serious aspect of this is that, although continued use may be quite acceptable in certain circumstances, few authors acknowledge that these lines are HeLa-contaminated, suggesting that a large part of the scientific community is still unaware of this problem, over 30 years after it was first revealed (21). Reports published by the ATCC (22), and the European Collection of Animal Cell Cultures (ECACC) (23), have confirmed that a significant proportion of cell lines in common use may be HeLa cross-contaminants.

### Bibliography

1. G. O. Gey, W. D. Coffman, and M. T. Kubicek (1952) *Cancer Res.* **12**, 364–365.
2. T. T. Puck and P. I. Marcus (1955) *Proc. Natl. Acad. Sci. USA* **41**, 432–437.
3. R. I. Freshney (1994) *Culture of Animal Cells, a Manual of Basic Technique*, Wiley-Liss, New York, p. 85.
4. T. R. Chen (1988) *Cytogenet. Cell Genet.* **48**, 19–24.
5. American Type Culture Collection catalogue, ATCC, P.O. Box 1549, Manassas, VA 20108 – 1549 [www.atcc.org](http://www.atcc.org), CCL-2.
6. W. K. Hansen, W. A. Deutsch, A. Yacoub, Y. Xu, D. A. Williams, and M. R. Kelley (1998) *J. Biol. Chem.* **273**, 756–762.
7. D. A. Jackson and A. Pombo, (1998) *J. Cell Biol.* **140**, 1285–1295.
8. M. D. Waterfield, M. A. Baker, M. F. Greaves, and E. J. Stanbridge (1986) *J. Biol. Chem.* **261**, 2418–2424.
9. K. Ess, H. Chen, A. Kier, and R. Brackenbury (1995) *J. Cell. Physiol.* **162**, 341–347.
10. Y. Chen, T. L. LaPorte, S. S. Wang, and J. Shevitz (1992) *Cytotechnology* **8**, 85–88.
11. G. J. Blaker, J. R. Birch, and S. J. Pirt (1971) *J. Cell. Sci.* **9**, 529–537.
12. B. Hertlein, A. Butterweck, S. Haubrich, K. Willecke, and O. Traub (1998) *J. Memb. Biol.* **162**, 247–257.
13. A. W. Gagnon, L. Kallal, and J. L. Benovic (1998) *J. Biol. Chem.* **273**, 6976–6981.
14. A. A. Newton and P. Wildy (1959) *Exp. Cell Res.* **16**, 624–635.
15. B. Lesser and T. P. Brent (1970) *Exp. Cell Res.* **62**, 470–473.
16. R. R. Rueckert et al. (1971) *Virology* **44**, 259–270.
17. S. M. Gartler (1967) *Second Decennial Review Conference on Cell, Tissue and Organ Culture*; NCI Monograph, pp. 167–195.
18. W. Nelson-Rees and R. R. Flandermeyer (1977) *Science* **195**, 1343–1344.



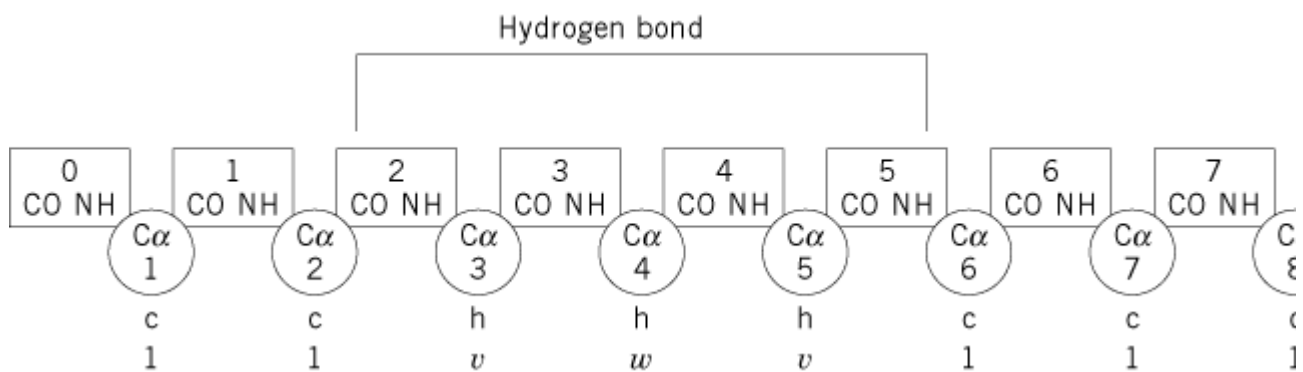
19. S. J. O'Brien, R. Olson, and J. E. Shannon (1977) *Science* **195**, 1345–1348.
20. T. R. Chen (1988) *Cytogenet. Cell Genet.* **48**, 19–24.
21. C. S. Stulberg, W. D. Peterson Jr., and W. F. Simpson (1976) *Am. J. Hematol.* **1**, 237–242.
22. K. S. Lavappa (1978) **14**, 469–475.
23. G. Stacey, B. Bolton A. Doyle, and B. Griffiths (1992) *Cytotechnology* **9**, 211–216.

## Helix–Coil Theory

Unlike many equilibrium folding–unfolding reactions in proteins that can be described adequately by a simple two-state folding mechanism, [a-helix formation](#) in polypeptides is not a simple two-state reaction from fully helical molecules to fully unstructured, or [random coil](#), peptides. This multistate behavior results from mixtures of partially helical molecules, primarily with the ends of the peptide “frayed.” From the architecture of the  $\alpha$ -helix, this is easy to understand. The ends of an  $\alpha$ -helix are different from positions in its center because at the *N*-terminus the NH groups, and at the *C*-terminus the C=O groups, are not paired in [hydrogen bonds](#), which are important for stabilizing the helical conformation. Therefore, under optimal helix-forming conditions, the population of helical molecules represent a distribution of partially helical states. Several models, based on simple statistical mechanics, have been developed to describe the macroscopic properties of the system. Collectively, these are known as *helix–coil theories*, although helix to random coil (or nonhelical structure) is a more accurate description of the models.

The two most popular models for the helix–coil transition in peptides were developed independently by Zimm and Bragg (1) and by Lifson and Roig (2) in the late 1950s and early 1960s (see [Zimm–Bragg Model](#) and [Lifson–Roig Model](#)). The two models share similar features, with a few important differences. For each model, helix formation is a function of three basic parameters: (i) the chain length or number of residues in the polypeptide, (ii) a helix propagation parameter ( $s$  or  $w$  value), and (iii) a helix nucleation parameter ( $v$  or  $v^2$ ). In each model, a single unit or residue can exist in only one of two possible conformations [helix or coil ( $h$  or  $c$ )] in the Lifson–Roig treatment, and the statistical weights ( $v$  or  $w$ ) reflect the microscopic equilibrium constant for the transition between  $h$  and  $c$ . Helical residues,  $h$ , flanked on both sides by other  $h$  residues are assigned a weight  $w$ , whereas  $h$  residues at the ends of a stretch of  $h$  residues are assigned a weight of  $v$ . The nucleation parameter ( $v^2$  in the Lifson–Roig treatment) reflects the “penalty,” considered primarily entropic, to order three consecutive residues into a stabilized helical segment. Figure 1 illustrates this for the simple case of a helical peptide with one hydrogen bond stabilizing the helical conformation. The energetic difficulting in nucleating the  $\alpha$ -helix imparts a degree of **cooperativity** to helix formation. Enumeration of all possible combinations of  $h$  and  $c$  residues for a polypeptide chain of a given length, along with the assigned statistical weights, provides the partition function for the system. The partition function can then be used to calculate average properties of the population, notably the fraction of helical residues, the average length and number of helical stretches, and other characteristics of the population.

**Figure 1.** Depiction of a peptide with one helical hydrogen bond formed. The squares represent the backbone peptide with the appropriate side chains. The  $C_{\alpha}$  groups are numbered, and the conformation of each residue is denoted as either peptide, the statistical weights are defined according to the Lifson–Roig model, using the coil state as a reference (weight  $v$  and internal helical residues  $w$ ).



The helix propagation parameter,  $w$  or  $s$ , is most closely akin to the helix propensity for a given residue. These helix propensities have recently been measured for all twenty amino acids found in proteins (see [Alpha-Helix Formation](#)). The Lifson–Roig model is also amenable to modification (3, 4) to include the energetics of other helix-stabilizing interactions, such as  $N$ - and  $C$ -capping and the interactions between side chains. A quantitative determination of these parameters, in conjunction with the Lifson–Roig model for the helix to coil transition, provides a reliable estimate of the amount of helix present in a peptide of any sequence.

#### Bibliography

1. B. H. Zimm and J. K. Bragg (1959) *J. Chem. Phys.* **31**, 526–535.
2. S. Lifson and A. Roig (1961) *J. Chem. Phys.* **34**, 1963–1974.
3. A. J. Doig, A. Charkrabarty, T. M. Klinger, and R. L. Baldwin (1994) *Biochemistry* **33**, 3396–3403.
4. B. J. Stapley, C. A. Rohl, and A. J. Doig (1995) *Protein Sci.* **4**, 2383–2391.

#### Suggestions for Further Reading

5. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, W. H. Freeman, San Francisco. Chapter 20 provides a good introduction to models for helix formation.
6. D. Poland and H. A. Scheraga (1970) *Theory of Helix–Coil Transitions in Biopolymers*, Academic Press, New York. A broad discussion of helix–coil theory. The book also reprints many of the seminal articles on helix–coil theory.

#### Helix-Turn-Helix Motif

The helix-turn-helix motif was first identified as a conserved sequence element in the repressors encoded by the lambdoid phages of *E. coli* and *Salmonella typhimurium*. Subsequently this element has been found in a large variety of DNA-binding proteins, both in prokaryotes and in eukaryotes, and shows a very high degree of structural conservation. This DNA-binding domain consists of two short  $\alpha$ -helices that usually are separated by a glycine residue. This amino acid in concert with its neighbors acts as a flexible hinge allowing the polypeptide chain to bend between the two helices so that they can make hydrophobic contacts with each other. These contacts preserve the relative orientation of the two helices and result in the formation of a compact tertiary structural domain. The

distal helix lies in the DNA major groove and the proximal helix contacts the sugar-phosphate backbone. Although the structure of this domain is highly conserved, the orientation of the recognition helix within the major groove is quite variable. This feature suggests that the motif is a structure that possesses a rigidity which is necessary for precise sequence recognition.

Related to the prokaryotic **helix-turn-helix** fold are protein folds of essentially the same fundamental structure, but which contain longer turns or loops. These eukaryotic variants include the homeodomain fold, which contains both a helix-turn-helix element and a short extended minor groove-tracking sequence, the POU-specific domain, and the helix-wing-helix fold. The latter fold contains an extended  $\beta$  sheet immediately following the recognition helix and is typified by hepatocyte nuclear factor 3 (HNF3) and the globular domain of the linker histones, H5 and H1. Whereas the prokaryotic proteins usually bind as homodimers and bind to palindromic DNA targets, the eukaryotic counterparts normally bind as monomers.

## Bibliography

“Helix-Turn-Helix Motif” in , Vol. 2, pp. 1105–1106, by Andrew Travers, MRC Laboratory of Molecular Biology, Cambridge; “Helix-Turn-Helix Motif” in (online), posting date: January 15, 2002, by Andrew Travers, MRC Laboratory of Molecular Biology, Cambridge.

## Helper Plasmid

In the context of genetic [transformation](#) of **plants**, a helper plasmid is a **plasmid** present in [Agrobacterium](#) that provides functions required by the **bacteria** for transferring foreign **DNA** to a plant cell. They have been extremely important in [plant genetic engineering](#). Generally, helper plasmids are derivatives of the [Ti plasmid](#) that contain an active virulence region, but from which the T-DNA has been removed.

Central to the concept of developing the binary vector system for *Agrobacterium*-mediated plant cell transformation was the finding that the **genes** contained on the T-DNA are not required for its transfer to the plant cell, that this function is provided by gene products of the virulence region of the Ti plasmid, and that they are [trans-acting](#) (1). In practice, this means that foreign DNA to be transferred to the plant cell can be **cloned** between the T-DNA border sequences in any plasmid that is stable in *Agrobacterium*. Then, such a plasmid, known as the binary vector, can be transferred to *Agrobacterium* that contains a helper plasmid by **conjugation**, electroporation, or transformation (see [Transfection](#)). Helper plasmids are based on deletion derivatives of Ti plasmid that lack the T-DNA. The helper plasmid provides the functions required for transferring the foreign DNA located between the T-DNA borders of the binary vector to the plant cell. The most widely used binary vector/helper plasmid system is the pBin binary vector (2) used in combination with the LBA4404 helper plasmid (1).

## Bibliography

1. A. Hoekema, P. R. Hirsch, P. J. Hooykaas, and R. A. Schilperoot (1983) *Nature* **303**, 179–181.
2. M. W. Bevan (1984) *Nucleic Acids Res.* **12**, 8711–8721.

## Suggestions for Further Reading

3. G. An (1995) In *Agrobacterium Protocols* (K. M. A. Gartland and M. R. Davey, eds.), Humana

Press Totowa, New Jersey. pp. 47–58.

4. F. F. White (1993) "Vectors for gene transfer in higher plants", In *Transgenic Plants*, Vol. 1 (D. Kung and R. Wu, eds.), Academic Press, San Diego, pp. 15–48.

## Helper Virus

Some **viruses** require the gene functions of other viruses to replicate. The viruses that rescue such replication-deficient viruses are known as helper viruses. [Adenoviruses](#), herpes simplex virus 1 and 2 (HSV-1, HSV-2) (see [Herpesvirus](#)), [cytomegalovirus](#), and pseudorabies virus are helpers for adeno-associated viruses (AAVs), dependoviruses of *Parvoviridae* family. In addition, [hepatitis B virus](#) is a helper virus for hepatitis D virus, and most standard viruses are helpers for their deletion mutants, which are called defective interfering particles. Defective viruses and their helper viruses are usually closely related in many replication processes (see [Rous Sarcoma Virus \(RSV\)](#)). It is now considered, however, that the helper functions are not simple and that they involve the interaction between virus and cells in addition to between defective virus and helper virus.

Genetic analysis of helper functions has been most extensive for adenoviruses. Five regions of the adenovirus genome are required for a permissive AAV infection. They are the E1A, E1B, E2A, E4, and VA regions. The E1A region is required for the *trans*-activation of the AAV p5 and p19 promoters. The E1B-coded 55-kDa transforming protein and the E4 gene product, a 35-kDa protein, are both required for the efficient accumulation of AAV [messenger RNAs](#). Their role in AAV replication is believed to be to stabilize the AAV mRNAs and assist their transport to the cytoplasm. The E2A gene product, an adenovirus [DNA-binding protein](#), and the adenovirus VA RNAs are required primarily for the efficient [translation](#) of AAV mRNA to produce the viral capsid protein.

Although a helper virus is necessary for fully permissive AAV infection, helper virus genes appear not to be directly involved in AAV [DNA replication](#). AAV DNA replication may use primarily cellular factors that exist in cells. The role of helper virus gene products may stimulate the synthesis of cellular genes that are required for AAV DNA replication.

### Suggestions for Further Reading

- K. I. Berns (1996) "*Parvoviridae: The Viruses and Their Replication*". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2173–2197.
- J. M. Taylor (1996) "Hepatitis Delta Virus and Its Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2809–2818.

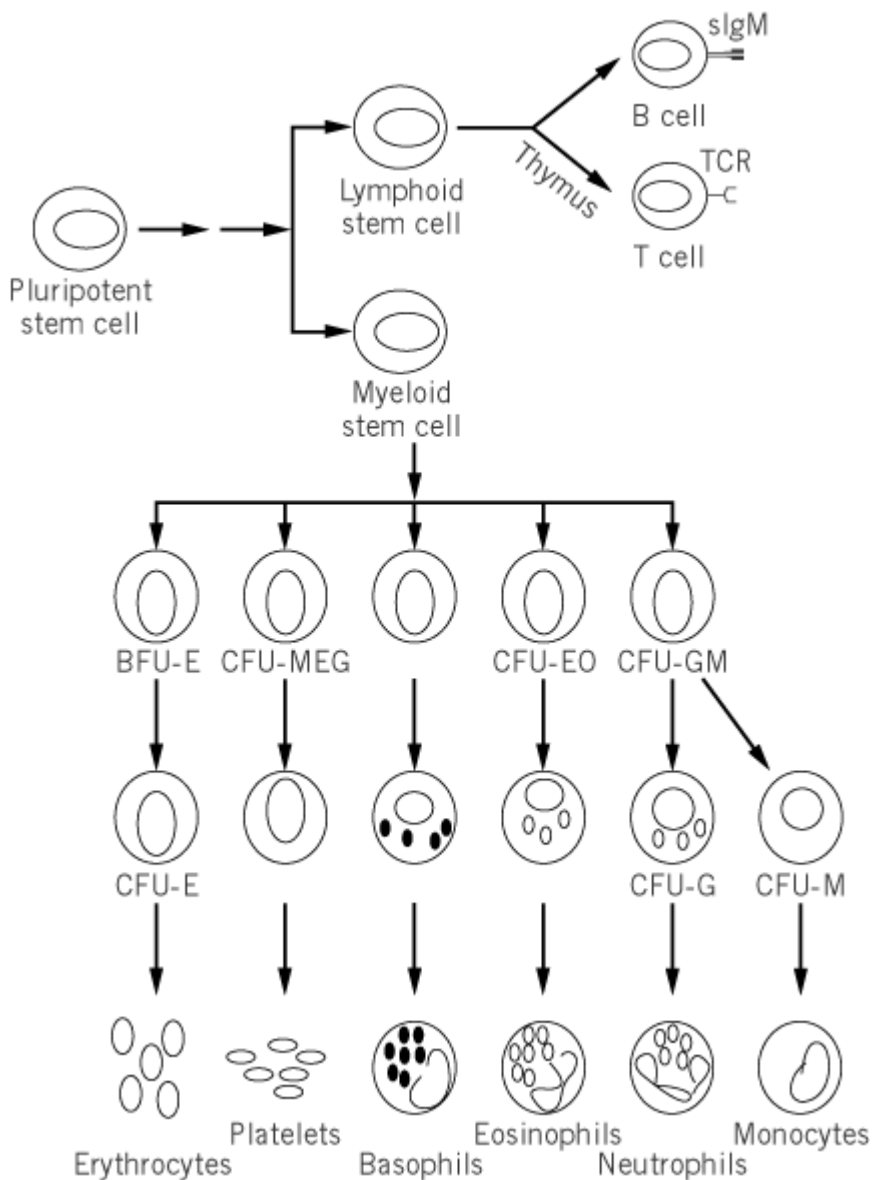
## Hematopoiesis

*Hematopoiesis* is a dynamic process in which blood cells of at least eight distinct lineages must be continually replaced from a population of pluripotent hematopoietic [stem cells](#) throughout the life span of the animal. The hematopoietic stem cell is a relatively rare cell, estimated to be

approximately 0.001% of nucleated bone marrow cells. Stem cells are normally quiescent, in either the G0 or G1 phase of the [cell cycle](#). Hematopoiesis involves the progressive restriction of the developmental potential of hematopoietic progenitors, with a concomitant increase in their proliferative capacity. Conversely, the self-renewal capacity of hematopoietic progenitor cells progressively decreases with the restriction of their developmental potential, with most mature cell types having short life spans and little or no proliferative capacity. In order to produce the correct numbers of mature blood cells, this balance between proliferation and differentiation of hematopoietic progenitors must be tightly controlled. Perturbation of this delicate balance can result in leukemogenesis or in anemias, immunodeficiencies, and other cellular deficiency disorders.

The existence of hematopoietic progenitor and stem cells can be inferred by the progeny they produce. Transplantation of lethally irradiated mice with cells from syngeneic donors has been the classical approach to demonstrate the presence of repopulating hematopoietic stem cells. The first 4 to 6 months after engraftment is characterized by clonal fluctuations in stem cell proliferation and differentiation, followed by the emergence of a stable hematopoietic system, dominated by a small number of totipotent clones. In addition, the extensive self-renewal capacity of these cells can be demonstrated by their ability to reconstitute recipient animals after serial transplantations. Competitive repopulation assays, in which two genetically distinct populations of bone marrow cells are transplanted together into the same recipient, can be used to determine the effects of a specific genetic alteration on the viability of repopulating hematopoietic stem cells.

The earliest restriction in developmental potential of hematopoietic stem cells involves the generation of putative lymphoid- and myeloid-restricted stem cells from the pluripotent myeloid-lymphoid stem cell. Myeloid-restricted stem cells can be detected by the classic spleen colony assay. In this assay, limiting numbers of bone marrow cells are injected into lethally irradiated recipient mice, giving rise to the formation of macroscopic hematopoietic colonies in the spleen within eight days, each derived from a single progenitor cell. These cells, termed CFU-S (colony-forming unit-spleen), are pluripotent and can give rise to all hematopoietic cells of the myeloid lineage. They possess only limited self-renewal capacity, as measured by their ability, on retransplantation, to give rise to secondary spleen colonies. In addition, CFU-S cannot readily confer hematopoietic repopulation and long-term survival of recipient mice. Additionally, progenitor cells that form spleen colonies at day 12 to 14 after transplantation are more potent in terms of repopulating ability than are day 8 CFU-S cells, suggesting that the spleen colony assay detects a range of cells that differ in their proliferative or repopulation potential.



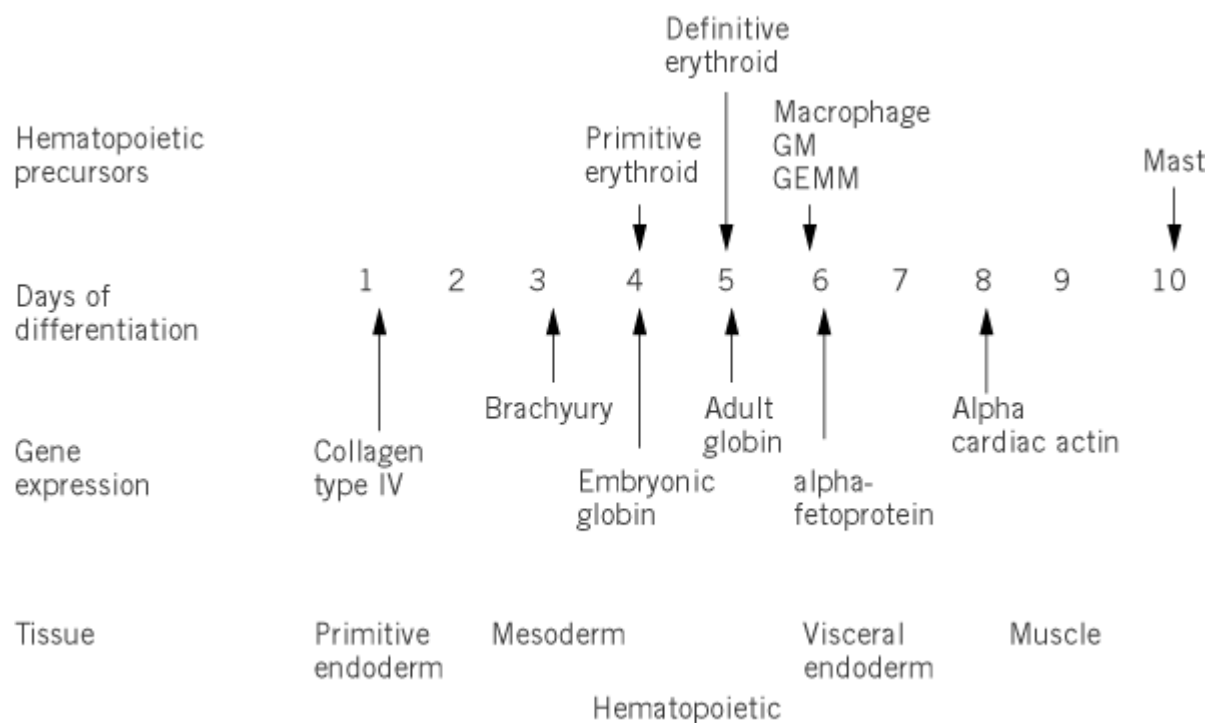
The presence of more restricted hematopoietic progenitor cells can be demonstrated by their ability to form colonies of differentiated blood cells in a semisolid medium in the presence of hematopoietic [growth factors](#). Multipotential hematopoietic [cytokines](#), including [interleukins](#) IL-3 (multipotential CF), IL-4, and granulocyte/macrophage [colony stimulating factor](#) (GM-CSF), are capable of supporting the proliferation of multilineage progenitor cells, resulting in the formation of colonies containing multiple hematopoietic lineages. Lineage-specific growth factors, such as erythropoietin, granulocyte CSF (G-CSF), or macrophage CSF (M-CSF or CSF-1), enhance proliferation and differentiation of committed progenitor cells forming colonies of nucleated red blood cells, neutrophils, and [macrophages](#), respectively. A final group of cytokines, including IL-6, IL-11, and Steel (stem cell factor, c-kit ligand), act synergistically with the multilineage and lineage-specific cytokines to enhance colony formation in culture by shortening the G<sub>0</sub> period of the hematopoietic progenitors.

During mouse embryogenesis, the first visible sign of hematopoiesis appears in the extraembryonic mesoderm of the day 7.5 yolk sac in the form of blood islands. These blood islands contain primarily primitive erythrocytes, which are large, nucleated cells expressing embryonic [hemoglobin](#), although adult globins are synthesized and accumulate in these cells at later stages of development. Cells of the monocyte-macrophage lineage also appear in the yolk sac. Although the differentiated

hematopoietic cells found in the yolk sac are restricted to cells of the erythroid and macrophage lineages, progenitor cells found in the yolk sac at later stages of development can differentiate into various hematopoietic lineages *in vitro* and have been reported to be capable of forming CFU-S *in vivo*, demonstrating the multipotency of these cells. Furthermore, repopulation of recipient animals with yolk sac-derived hematopoietic cells has been reported, suggesting the presence of repopulating hematopoietic stem cells within the yolk sac environment.

Around day 10 of embryonic development, the primary site of hematopoiesis switches from the yolk sac to the fetal liver, the main site of fetal hematopoiesis. The fetal liver produces predominantly cells of the definitive erythroid lineage, expressing adult hemoglobins, as well as some myeloid and lymphoid cells. Various progenitor and stem cells, including CFU-S and pluripotent hematopoietic stem cells capable of long-term repopulation of recipient mice, are all present at a high frequency in the fetal liver. The fetal liver remains the primary erythropoietic organ until shortly before birth, when the major site of hematopoiesis shifts, once again, to the bone marrow. Within the intersinusoidal spaces of the medullary cavity of the bone, hematopoietic cells are found closely associated with the stroma. These cells directly influence hematopoiesis through cell-cell interactions and the production of cytokines, forming a microenvironment that supports blood cell production.

**Figure 2.** Time course of embryoid body development. Gene expression and presence of hematopoietic precursors are shown.



It has long been accepted that, during embryogenesis, the hematopoietic stem cells (HSCs) that originate in the yolk sac migrate to the fetal liver to initiate the second wave of hematopoiesis. However, recent findings of hematopoietic cells in the para-aortic splanchnopleure of mouse embryos prior to colonization of the fetal liver (1-3) suggest the possibility of another source of hematopoiesis in the embryo proper. It is unclear at this time whether the cells involved in primitive hematopoiesis in the yolk sac are precursors to the cells involved in definitive hematopoiesis or whether definitive hematopoiesis develops independently in the embryo proper. Transplantation studies using the adult mouse as a recipient demonstrated HSC activity in the aorta, gonads, and

mesonephros (AGM) region of the embryo at day 10, followed by the fetal liver and yolk sac on day 11 (3, 4). These studies suggest that the first definitive HSCs arise in the AGM region and migrate to the fetal liver to establish definitive hematopoiesis. However, transplantation studies utilizing newborn or embryonic mice as recipients demonstrated the presence of HSCs contributing to definitive hematopoiesis in the yolk sac as early as day 9 (5), lending support to previous models. Although the origin of definitive hematopoiesis remains in debate, microenvironment is clearly a crucial element in the regulation of HSC potential.

The study of murine hematopoietic development, and the role of the microenvironment in this process, has been greatly facilitated by the establishment of *in vitro* model systems in which embryonic stem (ES) cells are differentiated to form cells of the hematopoietic system. Coculture of ES cells with primary bone marrow stroma or stromal cell lines results in the formation of a variety of hematopoietic cells including those of the erythroid, myeloid, B, and T cell lineages (6-8). Additionally, if the leukemia inhibitory factor is removed from the medium and the cells are grown in suspension, ES cells will form primitive embryonic structures termed *embryoid bodies* (EBs), which contain a vasculature and closely resemble the visceral yolk sac of the post-implantation embryo. These EBs contain many differentiated cell types, including those of the hematopoietic system, easily detectable by the appearance of visible blood islands containing red blood cells with hemoglobin.

ES-derived EBs express many genes that serve as markers of tissue restriction and lineage commitment during hematopoietic development (9-11), including the full complement of globin genes in the correct temporal order (12). EBs also contain precursors of the primitive erythroid, definitive erythroid, macrophage, neutrophil, and mast cell lineages, as demonstrated by colony-forming assays (10, 11, 13). Later in EB development, differentiated megakaryocytes, neutrophils, eosinophils, and mast cells also appear (14). Additionally, ES cell-derived EBs contain B and T cell precursors (15) and markers indicative of B and T lymphoid cells (16).

The *in vitro* differentiation of ES cells into cells of the hematopoietic system has been utilized to address the question of the origin of primitive versus definitive hematopoiesis. It remains unclear, however, whether primitive and definitive erythrocytes develop from common or distinct precursors in this system (17, 18). Recent evidence from *Xenopus* suggests the presence of bipotential primitive-definitive hematopoietic progenitors in the vertebrate embryo (19). However, the detectability of primitive erythrocytes in the yolk sac prior to the appearance of hematopoietic progenitor and stem cells in that tissue indicates that the traditional hierarchy of adult hematopoiesis may not apply in the context of primitive hematopoiesis.

The molecular events that regulate the development and maintenance of the hematopoietic system are beginning to come to light. A great deal of the work that forms the basis of this knowledge has been elucidated using the *mouse* as a model system. The combination of naturally occurring mutations and those generated by homologous *recombination* in ES cells has produced a battery of mutant ES and mouse lines, which have been instrumental in revealing the molecular biology of hematopoiesis. Furthermore, *in vitro* differentiation of homozygous mutant ES cells has provided additional information regarding the specific role of the targeted genes during hematopoietic development. Some of the genes that have been identified as essential for the development of hematopoietic stem cells during embryogenesis, or the functioning of these cells in the adult, are listed in Table 1.

**Table 1. Hematopoietic Mutants**

---

| <b>Disrupted Gene</b> | <b>Encoded Protein</b> | <b>Model</b> | <b>Phenotype</b> |
|-----------------------|------------------------|--------------|------------------|
|-----------------------|------------------------|--------------|------------------|



|                   |  |                                    |   |
|-------------------|--|------------------------------------|---|
| CBFa2 (AML1) (42) | Runt homology domain (RHD) TF <sup>a</sup> | KO mice                            | Block in definitive hematopoiesis<br>Stem cell defect                           |
| CBFb (43)         | Subunit of CBF                             | KO mice                            | Same as CBFa2   |
| EKLF (32)         | Zinc finger TF                             | KO mice                            | Block in definitive erythropoiesis  |
| EpoR (29, 30)     | Erythropoietin receptor                    | KO mice                            | Block in definitive erythropoiesis  |
| Flk1 (27)         | RTK  | KO mice/EBs                        | Stem cell defect  |
| GATA-1 (36, 37)   | Zinc Finger TF                             | KO mice/EBs                        | Block in primitive erythropoiesis<br>Block in definitive erythrocyte maturation |
| GATA-2 (38, 39)   | Zinc Finger TF                             | KO mice/EBs                        | Stem cell defect  |
| Kit (20)          | RTK  | <i>Dominant White Spotting (W)</i> | Stem cell defect<br>Macrocytic anemia<br>Mast cell deficiency                   |
| Myb (40)          | Myb domain TF                              | KO mice                            | Reduced definitive erythroid and myeloid precursors                             |
| PU.1 (Spi-1) (41) | Ets domain TF                              | KO mice                            | Block in definitive erythroid maturation<br>Stem cell defect                    |
| Rbtn2 (31)        | LIM domain TF                              | KO mice/EBs                        | Block in erythroid maturation   |
| Shp-1 (22, 23)    | Tyrosine phosphatase                       | <i>Motheaten (me)</i>              | Increased myeloopoiesis   |
| Steel (20)        | Ligand for <i>c-kit</i>                    | <i>Steel (Sl)</i>                  | Same as <i>W</i>  |
| Tal-1/SCL (33-35) | Basic HLH TF                               | KO mice                            | Lack of erythroid and myeloid precursors  |

<sup>a</sup> Key: TF, transcription factor; RTK, receptor tyrosine kinase; KO, knockout; EBs, embryoid bodies.

The process of hematopoiesis is controlled by a number of hematopoietic growth factors, which bind to specific cell surface **receptors** and generate a complex series of intracellular signaling events. These signaling events culminate in changes in **gene expression**, which direct the cell toward proliferation or differentiation. The major strategy by which hematopoietic cells interpret extracellular signals is through the activation of receptors that either possess intrinsic **tyrosine kinase** activity or are coupled to cytoplasmic proteins with tyrosine kinase activity. Several of these receptors have been shown to play a role in the regulation of hematopoietic stem cells and to be crucial for the proper development of the hematopoietic system. Naturally occurring mutations in the receptor tyrosine kinase, *c-kit* (*W*), and its ligand, steel factor (*Sl*), have uncovered an important role for this signaling pathway in the regulation of hematopoietic stem cells (20). In addition to macrocytic anemia and mast cell deficiency, hematopoietic cells from *W* mice show defects in spleen

colony formation and in long-term reconstitution of lethally irradiated recipient mice. In contrast, *Sl* mutant mice have a defective microenvironment that is unable to support the growth of normal stem cells. Genetic and biochemical analysis has revealed two other naturally occurring mutations that form part of the Kit-Steel signaling pathway (21). These mouse loci, *motheaten* (*me*) (22, 23) and *microphthalmia* (*mi*) (24, 25), encode a protein tyrosine phosphatase (Shp-1) and a basic helix-loop-helix [transcription factor](#) (MITF), respectively.

[Gene targeting](#) strategies have also led to the identification of several hematopoietic receptors that are essential for normal hematopoietic development. One of these receptors is Flk1, a receptor tyrosine kinase that binds *vascular endothelial growth factor* (VEGF). Mice homozygous for a mutation in Flk1 lack both mature endothelial and hematopoietic cells (26). In addition, Flk1 <sup>-/-</sup> ES cells do not contribute to primitive or definitive hematopoiesis *in vitro* or *in vivo* (27), demonstrating a close association between the molecular and cellular regulation of vasculogenesis and hematopoiesis. Gene targeting of STK, a tyrosine kinase receptor of the Met family isolated from a hematopoietic stem cell [cDNA library](#), has demonstrated a role for this receptor in the negative regulation of macrophage activation and [nitric oxide](#) production during a cell-mediated immune response (28).

A second class of cell surface receptors that are important for hematopoietic differentiation contains the cytokine receptors, which signal through the Jak/Stat pathway. An example of this class of receptors, which effects the proliferation and survival of erythroid progenitors, is the erythropoietin (epo) receptor. Mice homozygous for mutations in *epo* or the *epo* receptor die at embryonic day 13 from a failure of definitive fetal liver erythropoiesis (29, 30). Committed BFU-E (burst-forming unit-erythroid) and CFU-E (colony-forming unit-erythroid) progenitors are present in these mice, suggesting that the *epo* signaling pathway is necessary for the proliferation and survival of CFU-E progenitors but not for erythroid lineage commitment.

The signaling events generated in response to the engagement of hematopoietic receptors by their respective ligands culminate in changes in gene expression, which direct the cell toward proliferation or differentiation. Genetic analysis in mice has also revealed several transcription factors that are essential for the development of hematopoietic precursors. Collectively, these transcription factors regulate the expression of hematopoietic-specific genes that are important for the differentiation of multiple blood cell lineages from the pluripotential hematopoietic stem cell. Some of these factors are specific for the development of the erythropoietic lineage (GATA-1, EKLF, Rbtl2), suggesting a role in erythroid commitment and/or maturation. Other factors play a more global role in regulating the proliferation versus differentiation of primitive hematopoietic progenitors (PU.1, GATA-2, Myb, SCL/Tal-1). Interestingly, these studies suggest that primitive and definitive hematopoiesis are controlled by somewhat distinct molecular mechanisms. Disruption of genes encoding Rbtl2 (31), SCL/Tal-1 (32-35), GATA-1 (36, 37), or GATA-2 (38, 39) block both primitive and definitive erythropoiesis, whereas mutations in *c-myb* (40), PU.1 (41), EKLF (32), and CBF a and b (42, 43) affect only definitive erythropoiesis.

The list of mouse mutations effecting the development of the hematopoietic system is growing rapidly. Furthermore, many mutations exist that result in mild alterations in the regulation of adult hematopoiesis rather than in a complete block in hematopoiesis. This finding suggests some redundancy in the pathways that ultimately culminate in proliferation and differentiation of hematopoietic cells. The generation of mice with mutations in multiple genes involved in the regulation of hematopoiesis will be instrumental in the characterization of hematopoietic signaling pathways, as well as in the identification of points at which these pathways intersect.

## Bibliography

1. I. E. Godin et al. (1993) *Nature* **364**, 67–70.
2. A. L. Medvinsky et al. (1993) *Nature* **364**, 64–67.
3. A. M. Muller et al. (1994) *Immunity* **1**, 291–301.

4. A. Medivinsky and E. Dzierzak (1996) *Cell* **86**, 897–906.
5. M. C. Yoder et al. (1997) *Immunity* **7**, 335–344.
6. T. Nakano, H. Kodama, and T. Honjo (1994) *Science* **265**, 1098–1011.
7. R. Palacios, E. Golunski, and J. Samaridis (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7530–7534.
8. J. C. Gutierrez-Ramos and R. Palacios (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9171–9175.
9. T. McClanahan et al. (1993) *Blood* **81**, 2903–2915.
10. G. Keller et al. (1993) *Mol. Cell. Biol.* **13**, 473–486.
11. R. M. Schmitt, E. Bruyns, and H. R. Snodgrass (1991) *Genes Dev.* **5**, 728–740.
12. M. H. Lindenbaum and F. Grosveld (1990) *Genes Dev.* **4**, 2075–2085.
13. M. V. Wiles and G. Keller (1991) *Development* **111**, 259–267.
14. U. Burkert, T. von Ruden, and E. F. Wagner (1991) *New Biol.* **3**, 698–708.
15. A. J. Potocnik, P. J. Nielsen, and K. Eichmann (1994) *EMBO J.* **13**, 5274–5283.
16. U. Chen, M. Kosco, and U. Staerz (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2541–2545.
17. M. Kennedy et al. (1997) *Nature* **386**, 488–493.
18. T. Nakano, H. Kodama, and T. Honjo (1996) *Science* **272**, 722–724.
19. J. B. Turpen et al. (1997) *Immunity* **7**, 325–334.
20. A. Reith, and A. Bernstein (1991) *Genome Analysis 3: Genes and Phenotypes* (K. Davies and S. Tilghman, eds.) 105–133.
21. R. F. Paulson et al. (1996) *Nature Genetics* **13**, 309–315.
22. L. D. Shultz et al. (1993) *Cell* **73**, 1445–1454.
23. H. W. Tsui et al. (1993) *Nature Genet* **4**, 124–129.
24. C. A. Hodgkinson et al. (1993) *Cell* **74**, 395–404.
25. M. J. Hughes et al. (1993) *J Biol Chem* **28**, 20687–20690.
26. F. Shalaby et al. (1995) *Nature* **376**, 62–66.
27. F. Shalaby et al. (1997) *Cell* **89**, 981–990.
28. P. H. Correll et al. (1997) *Genes Function* **1**, 69–83.
29. H. Wu et al. (1995) *Cell* **83**, 59–67.
30. C.-S. Liu et al. (1996) *Genes Dev.* **10**, 154–164.
31. A. J. Warren et al. (1994) *Cell* **78**, 45–57.
32. B. Nuez et al. (1995) *Nature* **375**, 316–322.
33. C. Porcher et al. (1996) *Cell* **86**, 47–57.
34. L. Robb et al. (1996) *EMBO J.* **15**, 4123–4129.
35. R. A. Shivdasani, E. L. Mayer, and S. H. Orkin (1995) *Nature* **373**, 432–434.
36. L. Pevny et al. (1991) *Nature* **349**, 257–260.
37. M. J. Weiss, G. Keller, and S. H. Orkin (1994) *Genes Dev* **8**, 1184–1197.
38. F.-Y. Tsai and S. H. Orkin (1997) *Blood* **89**, 3636–3643.
39. F.-Y. Tsai et al. (1994) *nature* **371**, 221–226.
40. M. L. Mucenski et al. (1991) *Cell* **65**, 677–689.
41. E. W. Scott et al. (1994) *Science* **265**, 1573–1577.
42. T. Okuda et al. (1996) *Cell* **84**, 321–330.
43. Q. Wang et al. (1996) *Cell* **87**, 697–708.

### Suggestions for Further Reading

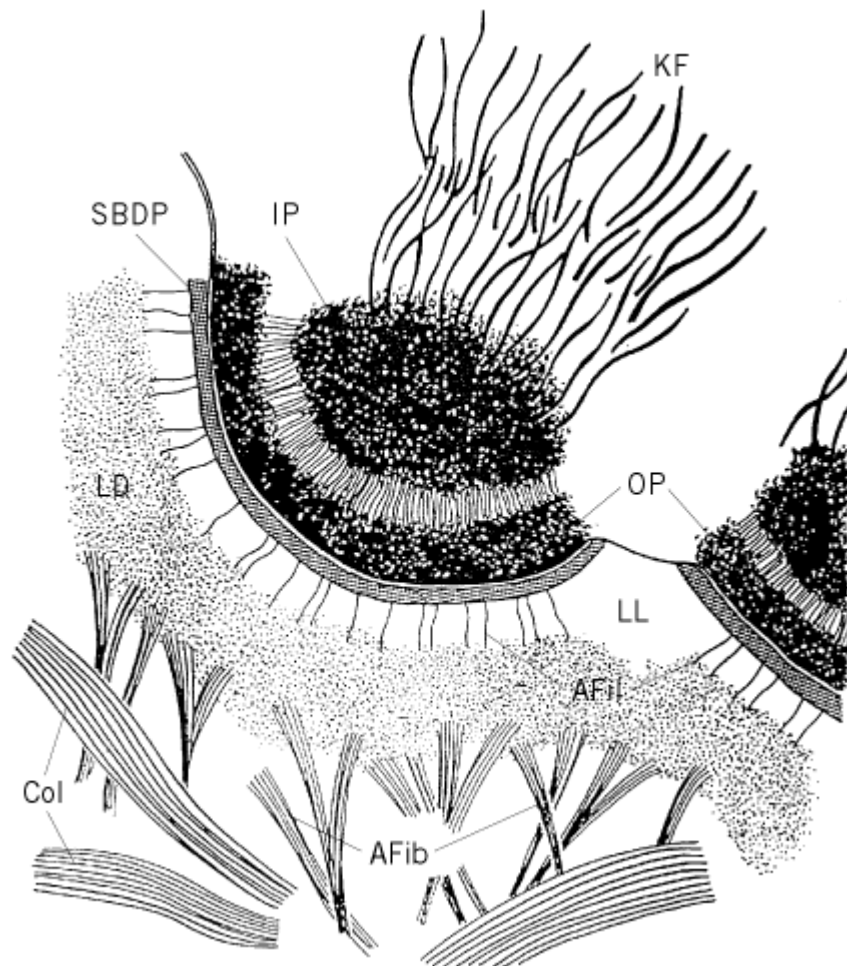
44. A. Muller and E. Dzierzak (1993) ES cells as a model of embryonic hematopoiesis? *Dev. Biol.* **4**, 341–349.

45. S. H. Orkin (1995) Transcription factors and hematopoietic development. *J. Biol. Chem.* **270**, 4955–4958.
46. R. F. Paulson and A. Bernstein (1995) Receptor tyrosine kinases and the regulation of hematopoiesis. *Sem in Immunol* **7**, 267–277.
47. L. I. Zon (1995) Developmental biology of hematopoiesis. *Blood* **86**, 2876–2891.

## Hemidesmosomes

Hemidesmosomes are [cell junctions](#) that are recognized by their ultrastructural appearance (Fig. 1). Located at the basal cell surface, their most prominent feature in cross section is a dense outer plaque closely applied to the inner surface of the plasma membrane (1). The outer plaque is separated by a less dense region from an inner dense plaque, from which [intermediate filaments](#) extend into the cytoplasm. The whole structure forms a circular membrane domain of no more than 0.5 micrometres in diameter and approximately 150 nm deep from the plasma membrane to the inner surface of the inner plaque. The extracellular face of the hemidesmosome is in contact with the lamina lucida of the basement membrane. This region is characterized by the presence of fine anchoring filaments that extend from the plasma membrane across the lamina lucida to the lamina densa. In some sections a linear density called the sub-basal dense plate is seen in the lamina lucida, closely associated with the extracellular face of the hemidesmosomal plasma membrane. Particularly in the epidermis and amnion, the matrix beneath the basement membrane contains numerous banded structures, called anchoring fibrils, which insert into the lamina densa. In favorable sections, filamentous continuity between the anchoring fibrils and the anchoring filaments can be observed. This means that there is continuity of structure in epidermis from the dermis right through to the intermediate filament [cytoskeleton](#) of the epidermal cells, in which the hemidesmosome itself provides the transmembrane connection. Such continuity appears to be essential for strong adhesive binding of the epidermis to the dermis, because genetic defects leading to absence or abnormality of anchoring fibrils, anchoring filaments, hemidesmosomes, or basal cell [keratin](#) filaments give rise to blistering diseases involving loss of adhesion between the epidermis and the dermis (see below). Although hemidesmosomes appear to resemble half **desmosomes**, they are in fact quite different structures, with a unique and distinct molecular composition (2).

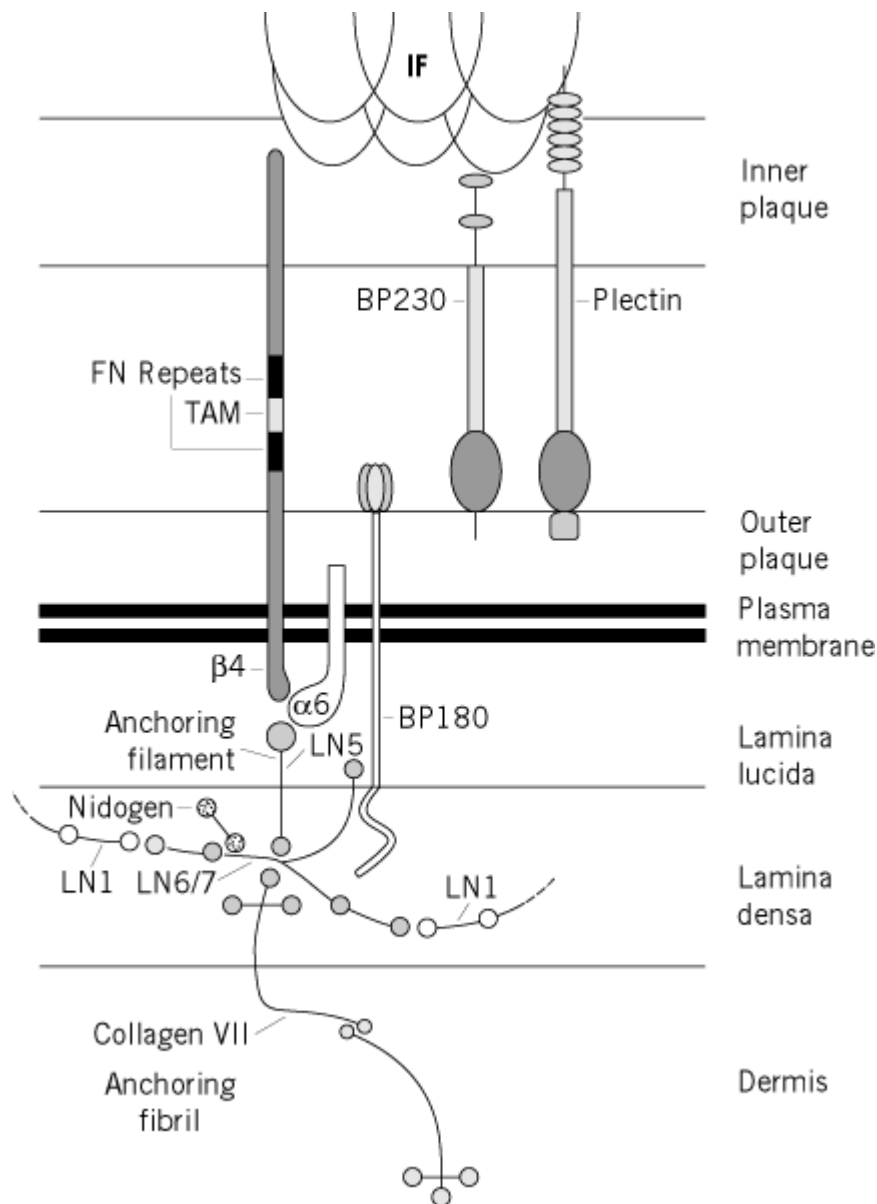
**Figure 1.** Drawing of an epidermal hemidesmosome seen by [electron microscopy](#). Afib = anchoring fibril; Afil = anchoring filament; Col=dermal collagen fiber; IP = inner plaque; KF = keratin filaments; LD = lamina densa; LL = lamina lucida; SBDP = sub-basal dense plate.



The major molecular components of hemidesmosomes are shown in Figure 2. The principal adhesion molecule is  $\alpha 6 \beta 4$  [integrin](#) (3), a heterodimer that is extremely abundant on the basal surface of epidermal cells. The  $\beta 4$  subunit is unique among integrin  $\beta$  subunits in having an extremely long cytoplasmic domain of approximately 1000 amino acid residues. The extracellular **domains** of the integrin subunits lie in the lamina lucida of the basement membrane, and their cytoplasmic domains in the hemidesmosomal plaque. The other adhesion molecule is bullous pemphigoid antigen 2 (BPAG2 or BP180), a Type 2 [membrane protein](#) whose amino-terminus lies within the desmosomal plaque and the carboxy-terminus in the basement membrane (4). BPAG2 is also referred to as [collagen](#) XVII, because it has collagenous repeats in its extracellular domain.  $\alpha 6 \beta 4$  integrin binds to the truncated **laminin** molecule, laminin 5, which is covalently complexed with laminin 6 or laminin 7 (5). Laminin 5, together with the extracellular domain of BP180, constitutes the anchoring filaments. Laminin 6 and 7 interact with other laminins in the basement membrane and also with nidogen, which links them to other basement membrane components. In addition, they bind to the specific component of anchoring fibrils, collagen VII (6). The anchoring fibrils descend from the lamina densa of the basement membrane into the underlying dermis. The principal plaque components of the hemidesmosome are bullous pemphigoid antigen 1 (BPAG1 or BP230) and plectin (7, 8). Both are members of the plakin family, which also includes desmoplakin, envoplakin, and periplakin of the **desmosome**. BP230 and plectin are believed to be involved in linking between the hemidesmosomal plaque and the intermediate filaments. Although not shown in the diagram, direct interaction between plectin and  $\alpha 6 \beta 4$  integrin has been demonstrated (9).

**Figure 2.** Diagrammatic impression of the molecular composition of a hemidesmosome employing information from a

number of sources.  $\alpha 6$  =  $\alpha 6$  integrin subunit;  $\beta 4$  =  $\beta 4$  subunit; BP180 = 180kDa bullous pemphigoid antigen (BPAG2); BP230 = 230kDa bullous pemphigoid antigen (BPAG1); FN = fibronectin; IF = intermediate filaments; LN 1, 5, 6, 7 = laminins 1, 5, 6, and 7; TAM = tyrosine activation motif.



A number of human diseases have been extremely important in elucidating the composition and function of hemidesmosomes. In the **autoimmune** blistering disease bullous pemphigoid, **antibodies** against BPAG1 and 2 are commonly found. Those against an extracellular juxta membrane epitope of BPAG2 have been shown to be pathogenic (10). A number of genetic diseases, collectively known as epidermolysis bullosa (EB), affect basal cells of epidermis and epidermal-basement membrane adhesion. EB simplex involves the basal cell keratins, K5 and K14 (see [Keratins](#)). Such mutations disrupt the assembly of the basal cell intermediate filaments and result in cell breakage at a level above the hemidesmosomes (11). The junctional form of EB (EBJ) involves mutations in the  $\alpha$ ,  $\beta$ , or  $\gamma$  chains of laminin 5 or in  $\alpha 6 \beta 4$  integrin (12). This is a lethal disease resulting in epidermal detachment at the level of the basement membrane. EB dystrophica results from mutations of the anchoring fibril component Type VII collagen (13). Anchoring fibrils fail to form, and the epidermis becomes detached below the level of the basement membrane. The importance of plectin in hemidesmosomes has been revealed by the disease EB with muscular dystrophy (14). Here attachment of intermediate filaments from the hemidesmosomal plaque occurs in basal epidermal

cells, resulting in epidermal blistering. This is followed later in life by muscular weakening, arising from absence of plectin from the submembrane cytoskeleton of skeletal muscle fibers. Null mutations of *a6b4* integrin, plectin, and BP230 produce phenotypes consistent with the defects found in human disease (15-18). This work suggests that BPAG1 forms a similar function to plectin in attaching intermediate filaments to the plaque. The BPAG1 null mutation revealed an additional function for this molecule in the nervous system (17). BPAG1<sup>-/-</sup> mice develop severe dystonia, resulting in abnormal limb positioning and loss of involvement through muscle and nerve regeneration. This is reminiscent of a hereditary neurodegenerative disease of mice called *dystonia musculorum (dt/dt)*. *dt* is allelic with BP230. It appears that [alternative splicing](#) of the amino-terminus of BPAG1 in the nervous system promotes interaction with actin, and its absence results in disorganisation of the neuronal cytoskeleton (19). It is noteworthy that the dermal epidermal junction requires extremely strong adhesion of cells to the extracellular matrix. This may be why hemidesmosomes have two adhesion molecules and two molecules to bridge between the plaque and the intermediate filaments.

### Bibliography

1. J. E. Ellison and D. R. Garrod (1984) *J. Cell Sci.* **72**, 163–172.
2. P. K. Legan et al. (1992) *Bioessays* **14**, 385–393.
3. A. Sonnenberg et al. (1991) *J. Cell Biol.* **113**, 907–917.
4. Y. Hirako et al. (1996) *J. Biol. Chem.* **271**, 13739–13745.
5. M. F. Champlaud (1996) *J. Cell Biol.* **132**, 1189–1198.
6. H. P. Bächinger et al. (1990) *J. Biol. Chem.* **265**, 10095–10101.
7. D. H. Klatte et al. (1989) *J. Cell Biol.* **109**, 3377–3390.
8. G. Wiche et al. (1983) *J. Cell Biol.* **97**, 887–901.
9. C. M. Neissen (1997) *J. Cell Sci.* **110**, 1705–1716.
10. Z. Liu et al. (1993) *J. Clin. Invest.* **92**, 2480–2488.
11. W. H. I. McLean and E. B. Lane (1995) *Curr. Opin. Cell Biol.* **7**, 118–125.
12. R. A. J. Eady and M. G. Dunhill (1994) *Arch. Dermatol. Res.* **287**, 2–9.
13. A. M. Christiano et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 3549–3553.
14. W. H. I. McLean et al. (1996) *Genes. Dev.* **10**, 1724–1735.
15. E. Georges-Labouesse et al. (1996) *Nature Genet.* **13**, 370–373.
16. R. Van der Neut et al. (1996) *Nature Genet.* **13**, 366–369.
17. K. Andrä et al. (1997) *Genes Dev.* **11**, 3143–3156.
18. L. Guo et al. (1995) *Cell* **81**, 233–243.
19. Y. Yang et al. (1996) *Cell* **86**, 655–665.

### Suggestions for Further Reading

20. D. R. Garrod (1993) *Curr. Opin. Cell Biol.* **5**, 30–40. (A detailed molecular comparison of desmosomes and hemidesmosomes.)
21. K. J. Green and J. C. R. Jones (1996) *FASEB. J.* **10**, 871–881. (An excellent review.)
22. D. R. Garrod and J. E. Collins (eds.) (1994) *The Molecular Biology of Desmosomes and Hemidesmosomes*, R. G. Landes, Austin, TX. (Articles by some of the leaders in the field.)

### Hemizygote

The notion of *hemizyosity* has historical importance in bacterial **genetics**. After the early demonstration by Lederberg of genetic [recombination](#) in *Escherichia coli* and the extension of genetic analysis to include a much wider range of unselected markers, a number of complicated anomalies of [chromosome](#) segregation began to appear showing that this bacterial system is unconventional. A particular segment of the chromosome of one of the parents, identified by a pair of linked markers, was virtually missing from the **heterozygotes** generated initially and failed to appear in their [haploid](#) progeny. These initial **diploids** were incomplete (*hemizygous*). It was first wrongly thought that the missing segment was eliminated from them after their formation. The phenomenon was even given the name of *postzygotic exclusion*. Analysis of this phenomenon led to the discovery of sexuality in bacteria and to the notion of [circular chromosome](#), which resolved the difficulties. To avoid confusion, the term *merozygote* should be used instead of hemizygote for these incomplete and transient heterozygotes.

The term hemizygote should be restricted to a diploid cell or individual in which one of the two copies of a given **gene** has been lost either by deletion or loss of a chromosome. Heterogametic species, in particular human males, are said to be hemizygous for genes located on the [X-chromosome](#) and not on the [Y-Chromosome](#) (and vice versa).

#### Suggestions for Further Reading

W. Hayes (1968) *The Genetics of Bacteria and their Viruses*, 2nd ed., Blackwell Scientific Publications, Oxford and Edinburgh, pp. 651–653.

F. Jacob and E. L. Wollman (1961) *Sexuality and the Genetics of Bacteria*, Academic Press, New York, pp. 34–36.

## Hemoglobin

Hemoglobin (Hb) is one of the heme-containing [oxygen-binding proteins](#), that generally contain a [globin](#) polypeptide chain and a protoheme IX **prosthetic group** that contains a ferrous iron atom. It is found in circulating fluids (blood or hemolymph) of animals, where its fundamental function is transport of molecular oxygen. Its **phylogenetic** distribution is wide, but somewhat capricious. All of the vertebrates, except icefish, have a Hb within their red blood cells that is a tetramer made up of two pairs of two different subunits, each a globin polypeptide chain that contains one heme group. In contrast, invertebrate Hbs are strikingly variable, in their architecture and size, and also in their function. They can be either intracellular or extracellular. Generally, an intracellular Hb is relatively small, whereas extracellular Hbs have a molecular mass ranging from 250 to 3500 kDa. The latter are polymeric and are dissolved directly in the blood plasma, hemolymph, or coelomic fluid. The invertebrate extracellular Hbs have been known traditionally as *erythrocruorins*, but this term lost its significance when it became apparent that there is no particular distinction between the subunits of the extracellular and of the vertebrate Hbs in their prosthetic group, amino acid sequence [homology](#), or three-dimensional [protein structure](#), the globin fold.

### 1. Vertebrate Hemoglobin

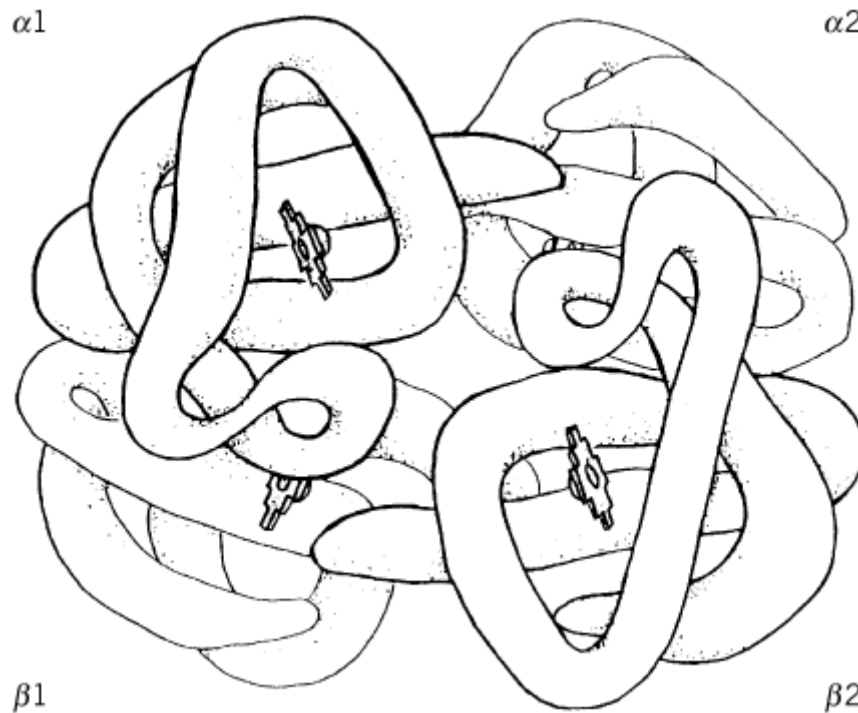
#### 1.1. Structure

The adult form of human Hb, HbA, is composed of two  $\alpha$ -subunits and two  $\beta$ -subunits, each of which has one bound heme group (Fig. 1). The  $\alpha$ - and  $\beta$ -chains consist of 141 and 146 amino acid



residues, respectively. All of the vertebrate Hb chains have similar numbers of residues and homologous amino acid sequences. All adopt very similar globin folds comprised of eight **alpha-helices**, designated A through H, although  $\alpha$ -chains lack  $\alpha$ -helix D. Early in [development](#), related polypeptide chains replace the  $\alpha$ - and  $\beta$ -chains.

**Figure 1.** The structure of the vertebrate hemoglobins. The molecule is a tetramer composed of two pairs of  $\alpha$ - and  $\beta$ -subunits. The  $\alpha$ - and  $\beta$ -chains are homologous, have similar globin folds, and each has a bound heme group. How the individual chains are distinguished is described in the text.



Within the  $\alpha_2\beta_2$  tetramer, there are relatively few contacts between the pair of  $\alpha$ -subunits and the pair of  $\beta$ -subunits, whereas contacts are extensive between the unlike subunits (Fig. 1). If the  $\alpha$ -subunits are designated  $\alpha 1$  and  $\alpha 2$ , the  $\beta$ -subunits are numbered so that the distance between the heme iron atoms of the  $\alpha 1/\beta 2$  pair is shorter than that in the  $\alpha 1/\beta 1$  pair. Because of symmetry, the  $\alpha 1/\beta 1$  pair is identical to  $\alpha 2/\beta 2$ , and  $\alpha 1/\beta 2$  is the same as  $\alpha 2/\beta 1$  (see [Quaternary Structure](#) and [Oligomeric Proteins](#)). The tetrameric assembly is stabilized by weak noncovalent bonds, by many [van der Waals interactions](#), a few [hydrogen bonds](#) between unlike subunits, and several [salt bridges](#) between the two  $\beta$ -subunits.

## 1.2. Ligand Binding

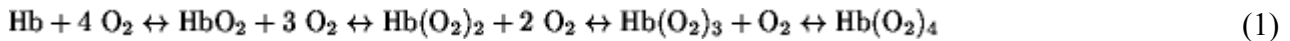
One molecule of  $O_2$  binds to the sixth coordination site of the iron atom which is on the distal side of the heme group (see [Globins](#) and [Myoglobin](#)). The heme pocket in the globin fold provides a highly **hydrophobic** environment that maintains the heme iron atom in the ferrous state. Oxidation of the iron atom to ferric to produce methemoglobin destroys its ability to bind oxygen. Upon binding oxygen, the quaternary structure undergoes an extensive change in the relative orientations of the four subunits, especially due to changes in the  $\alpha 1/\beta 2$  and  $\alpha 2/\beta 1$  interfaces. This switch between two alternative quaternary structures is believed to be involved in the **cooperativity** of oxygen binding. Hb was the primary inspiration of the **concerted model** of [allostery](#), and it remains one of the proteins that best fits that model.

Several other classes of ligand are bound at other specific sites on the deoxy quaternary structure of hemoglobin. These ligands include protons (see [Bohr Effect](#)), carbon dioxide, chloride ions, and 2,3-diphosphoglycerate (2,3-DPG). As the ligands bind more avidly to the deoxy form of Hb, they compete in effect with oxygen for binding, and the presence of one of these ligands affects the binding of the others.

### 1.3. Physiological Function

The physiological functions of vertebrate Hb are (1) transport of oxygen from the lungs to peripheral tissues; (2) transport of carbon dioxide and enhancement of its transport by the blood plasma; and (3) pH regulation (acid–base balance) of the blood.

The transport of oxygen by Hb is well described by its oxygen dissociation curve (or oxygen equilibrium curve), which expresses the dependence of oxygen saturation of Hb ( $Y$ ) upon the partial pressure of oxygen ( $P_{O_2}$ ) at chemical equilibrium (Fig. 2). The general features of this curve are (1) its sigmoid shape and (2) its shift to the right to higher oxygen partial pressures at higher temperatures or in the presence of the nonheme ligands that bind specifically to deoxy-Hb, protons,  $CO_2$ ,  $Cl^-$ , and 2,3-DPG. The sigmoid shape of the binding curve is ascribed to stepwise enhancement of the oxygen-binding affinity:



and

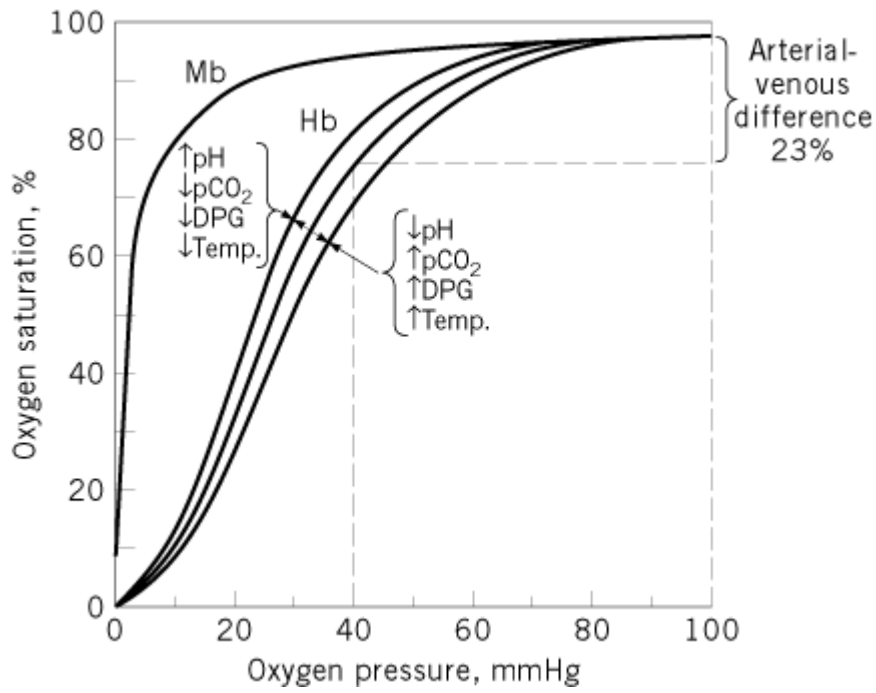
$$Y = \frac{[HbO_2] + 2[Hb(O_2)_2] + 3[Hb(O_2)_3] + 4[Hb(O_2)_4]}{4([Hb] + [HbO_2] + [Hb(O_2)_2] + [Hb(O_2)_3] + [Hb(O_2)_4])} \quad (2)$$

$$= \frac{K_1 p + 3K_1 K_2 p^2 + 3K_1 K_2 K_3 p^3 + K_1 K_2 K_3 K_4 p^4}{1 + 4K_1 p + 6K_1 K_2 p^2 + 4K_1 K_2 K_3 p^3 + K_1 K_2 K_3 K_4 p^4}$$

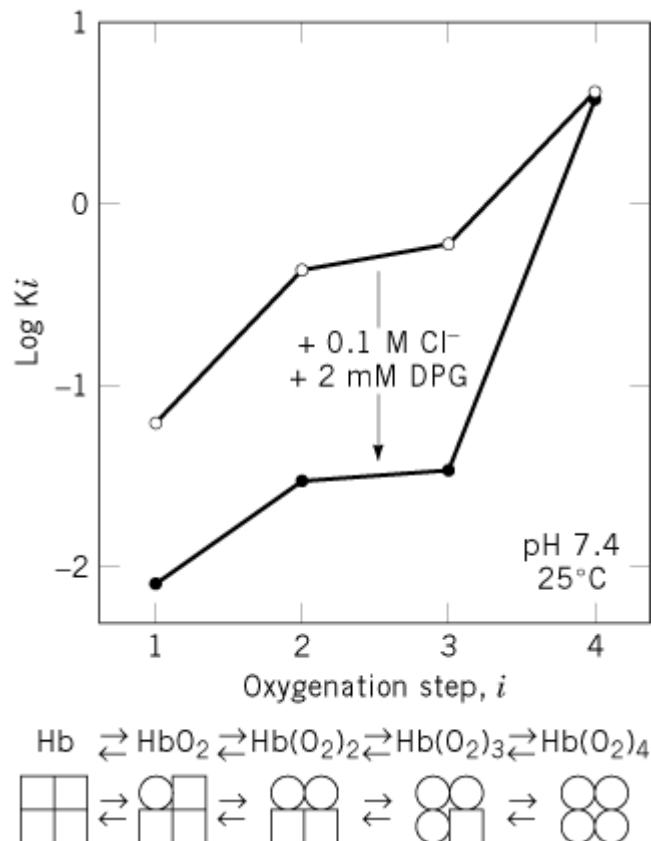
Here,  $p = P_{O_2}$  and  $K_i$  ( $i = 1, 2, 3, 4$ ) are the intrinsic oxygen **association constants** for the  $i$ th molecule to bind. If the four heme groups of  $\alpha_2\beta_2$  Hb were equivalent and independent of each other,  $K_1 = K_2 = K_3 = K_4$ , which would produce a hyperbolic oxygen dissociation curve, like that of [myoglobin](#) (Mb) (Fig. 2). Mb is a monomeric globin that contains a single heme group. In contrast, HbA gives a sigmoidal curve in which  $K_1 < K_2$ ,  $K_3 \ll K_4$  and the ratio of  $K_4$  to  $K_1$  is 500 (2) (Fig. 3).

This cooperative binding of oxygen produces the sigmoid curve, in which apparent interactions between the four heme groups of the molecule are mediated by conformational changes in the protein (3). The sigmoidal shape of the oxygen-binding curve contributes to increasing the difference in oxygen saturation between arterial and venous blood and enhances the amount of oxygen released at the tissues (Fig. 2).

**Figure 2.** Oxygen dissociation curves for hemoglobin (Hb) and myoglobin (Mb). The oxygen concentration is specified by its partial pressure in mm of Hg. The dissociation curve for the monomeric Mb has the expected hyperbolic shape, whereas that for Hb is sigmoid, indicating positive cooperativity in binding the four oxygen molecules by each Hb tetramer. The dissociation curve for Hb is shifted to the left or right upon changing the solution conditions or temperature, as indicated, whereas that for Mb is insensitive to changes in conditions, except for the temperature. The middle of the three curves for Hb is for standard conditions: 37°C, pH 7.4, 40 mm Hg  $CO_2$ , and 5 mM 2,3-DPG. The partial pressure of oxygen for arterial blood, 100 mm Hg, and that for mixed venous blood, 40 mm Hg, produce a difference in oxygen saturation of 23%, whereas there would be a much smaller release of oxygen for a hyperbolic curve like that of Mb.



**Figure 3.** Changes in oxygen affinity of human hemoglobin with oxygen binding at the four heme groups and those upon addition of the cofactors 0.1 M  $\text{Cl}^-$  and 2 mM 2,3-DPG. The abscissa indicates the  $i$ th oxygen molecule to bind, and the ordinate gives the logarithm of the intrinsic association constant for the  $i$ th step. Open circles correspond to the absence of cofactors, closed circles to their presence. The four-step oxygen binding is illustrated schematically below the graph. Squares and circles represent deoxy and oxy subunits, respectively. Data taken from K. Imai (1982) *Allosteric Effects in Haemoglobin*, Cambridge University Press, London.



Under physiological conditions, the nonheme ligands lower the oxygen affinity and increase the cooperativity of oxygen binding by HbA. The effect of these ligands on the intrinsic binding constants,  $K_p$ , is not uniform with respect to  $i$  (Fig. 3). The shift of the oxygen binding curve to higher partial oxygen pressures, that is, the lowering of the overall oxygen affinity upon acidification or upon increasing the CO<sub>2</sub> concentration or temperature has physiological consequences because this enhances the release of oxygen at the peripheral tissues, where such changes result from active metabolism. The effect of pH on the oxygen affinity is known as the [Bohr effect](#). In some cases of chronic hypoxia found in anemia, chronic heart and lung diseases, etc., there is an increase in the intracellular concentration of 2,3-DPG, a decrease in Hb's oxygen affinity, and consequently an enhancement of oxygen release at tissues (4).

Carbon dioxide produced in tissue cells diffuses into the red blood cells and reacts with water, catalyzed by [carbonic anhydrase](#):



The HCO<sub>3</sub><sup>-</sup> bicarbonate ions thus produced move to the blood plasma in exchange for chloride ions and are transported to the lungs. Reaction 3 is greatly shifted to the right when Hb absorbs the protons produced. Moreover, dissolution of the carbon dioxide is greater in the tissue capillaries than in the alveolar because of the uptake of protons upon oxygen release by Hb:



This oxygen-linked proton binding known as the Haldane effect is the reciprocal expression of the Bohr effect. The release of oxygen in the tissue capillaries enhances the uptake of protons by Hb, thereby causing a further shift of Eq. 3 to the right. Part of the CO<sub>2</sub> is transported linked directly to Hb as a carbamino moiety:



where -NH<sub>2</sub> is the α-amino groups of the α- and β-chains. This CO<sub>2</sub> binding is also oxygen-linked. Dissociation of oxygen in tissue capillaries enhances the CO<sub>2</sub> binding by Hb, increasing the CO<sub>2</sub> content of the venous blood. The remainder of the CO<sub>2</sub> is dissolved directly in the blood plasma. The major portion (85%) of the total CO<sub>2</sub> is transported as bicarbonate, 10% as the carbamino form, and the remaining 5% as free CO<sub>2</sub>. Consequently, as a result of its allosteric properties, Hb plays a key role in CO<sub>2</sub> transport and oxygen transport.

The properties of HbA described here also apply to other vertebrate Hb, except for some slight differences. In some species, 2,3 DPG is replaced by other cellular metabolites, such as ATP, GTP, or **inositol pentaphosphate**, and CO<sub>2</sub> is replaced by bicarbonate ion.

## 2. Annelid Hemoglobins

Annelida possess giant extracellular HB of molecular mass 3500 kDa. The molecular architecture is a hexagonal bilayer whose diameter is about 30 nm and whose thickness is about 20 nm. The molecule is composed of 12 identical spherical protein units, known as “submultiples,” six of which make one hexagonal layer. These submultiples are assembled to make up an entire molecule using many linker subunits. The submultiple is composed of 16-kDa globin subunits that have the same globin fold as the vertebrate Hb subunits, but the linker subunits are unrelated, heme-free

polypeptide chains. The extracellular Hb from the earthworm *Lumbricus terrestris* has been most extensively studied. Its submultiple is a dodecamer composed of three copies of each of four different kinds of globin chains, and 12 such dodecamers are assembled with 36 linker chains and approximately 57 tightly bound calcium ions. Consequently, there are  $12 \times 12 (= 144)$  heme groups in each molecule. Some of the globin and linker chains have attached carbohydrates, which are thought to contribute to stabilization of the quaternary structure through their noncovalent, **lectin**-like binding (5, 6).

Annelid Hbs demonstrate a striking diversity of oxygen-binding properties. Most demonstrate a Bohr effect and cooperativity of oxygen binding like those of vertebrate Hb, and sometimes of even greater magnitude, but others have none. The annelid Hb do not respond to cofactors, such as 2,3-DPG or other organic phosphate compounds. Instead, their oxygen affinity is increased by divalent calcium and magnesium ions, which are abundant in the annelid blood. There is some evidence that the functional allosteric unit of annelid Hb is the submultiple.

### 3. Hemoglobins of Other Species

A great variety of Hbs are found in lower organisms. They can be (1) one-domain monomers, that have a single globin chain of about 16 kDa and one heme group; (2) polymers of such single-domain monomers; (3) two-domain polymers; and (4) multiple-domain polymers. Most of the single-domain monomer Hbs are present in the cytoplasm. Some are in noncirculating cells, whereas others are in coelomic cells of annelida, such as *Glycera*. Although those in noncirculating cells do not apparently participate in oxygen transport, they are still classified as hemoglobins. Leghemoglobin is a monomeric Hb contained in the cytoplasm of root nodules of legumes. It has an extremely high oxygen affinity, which is important for the completely removing of traces of oxygen from the symbiotic bacteria, because nitrogen fixation is a strictly anaerobic process. Consequently, the function of leghemoglobin is to eliminate oxygen, rather than to supply it to tissues.

Finding the single-domain Hb from the gram-negative bacterium *Vitreoscilla* was the first indication that the evolutionary origin of globins may date back to the common ancestor with prokaryotes. Hbs isolated from **yeasts** and *Escherichia coli* have two domains, one of which is a globin chain and the other a flavoprotein. The NADH reductase activity of the second domain prevents the heme iron of the first domain from auto-oxidizing to maintain its capability of reversibly binding oxygen. *Chironomus* larvae possess Hbs that are polymorphic in various species, tissues, and developmental stages. Among the 12 Hbs of *C. thummi thummi*, seven exist as homodimers, four as monomers, and one as a monomer or dimer (see [Chironomus](#)). The [nematode](#) *Ascaris lumbricoides* Hbs consists of two-domain globin chains, one of which has lost the ability to bind heme. Hbs from nematodes, including *Ascaris*, have extremely high oxygen affinities. The brine shrimp *Artemia salina* Hbs have a molecular mass 250 kDa and are composed of two polypeptide chains, each containing eight heme-binding domains. The blood clam *Scapharca inaequivalvis* has intracellular homodimeric and heterotetrameric Hbs, which bind oxygen cooperatively. The most conspicuous feature of these Hbs is that their globin chains have an additional  $\alpha$ -helix at their N-terminus and that the E and F helices are external, exposed to solvent, whereas the G and H helices are internal, involved in subunit interactions. Thus, the clam Hb tetramer is “back to front” relative to the vertebrate Hb tetramer.

The various species of Hb demonstrate oxygen affinities that vary 26,000-fold. The lowest affinity, measured by the value of the partial pressure of oxygen at half saturation, is 26 mm of Hg, that of human Hb at 37°C, whereas that for the highest affinity is only 0.001 mm Hg, for *Ascaris* Hb at 20°C. The species-dependence of the oxygen affinity of Hb is a manifestation of its adaptation to the environments that the organisms inhabit. All of the Hbs have the same protoheme IX heme group, so this tremendous variation in oxygen affinity arises from species differences in the amino acid sequences of the globin moiety. How this is accomplished is not known.

## Bibliography

1. G. S. Adair (1925) *J. Biol. Chem.* **63**, 529–545.
2. K. Imai (1979) *J. Mol. Biol.* **133**, 233–247.
3. M. F. Perutz (1970) *Nature* **228**, 726–739.
4. F. A. Oski, A. J. Gottlieb, M. Delivoria-Papadopoulos, and W. W. Miller (1969) *N. Engl. J. Med.* **280**, 1165–1166.
5. S. Ebina, K. Matsubara, K. Nagayama, M. Yamaki, and T. Gotoh (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7367–7371.
6. P. D. Martin, A. R. Kuchumov, B. N. Green, R. W. A. Oliver, E. H. Braswell, J. S. Wall, and S. N. Vinogradov (1996) *J. Mol. Biol.* **255**, 154–169.

## Suggestions for Further Reading

7. R. E. Weber (1980) Functions of invertebrate hemoglobins with special reference to adaptations to environmental hypoxia, *Am. Zool.* **20**, 79–101.
8. K. Imai (1982) *Allosteric Effects in Haemoglobin*, Cambridge University Press, London.
9. S. N. Vinogradov (1985) The structure of invertebrate extracellular hemoglobins (erythrocruorins and chlorocruorins), *Comp. Biochem. Physiol.* **82B**, 1–15.
10. M. F. Perutz (1989) Mechanisms of cooperativity and allosteric regulation in proteins, *Q. Rev. Biophys.* **22**, 139–236.
11. T. Gotoh and T. Suzuki (1990) Molecular assembly and evolution of multi-subunit extracellular annelid hemoglobins, *Zool. Sci.* **7**, 1–16.
12. M. Bolognesi, D. Bordo, M. Rizzi, C. Tarricone, and P. Ascenzi (1997) Nonvertebrate hemoglobins: Structural basis for reactivity, *Prog. Biophys. Mol. Biol.* **68**, 29–68.
13. M. F. Perutz, A. J. Wilkinson, M. Paoli, and G. G. Dodson (1998) The stereochemical mechanism of the cooperative effects in hemoglobin revisited, *Annu. Rev. Biophys. Biomol. Struct.* **27**, 1–34.

## Hemoglobin Mutations

The study of human hemoglobin mutations opened the era of molecular genetics. In 1949, Linus Pauling and colleagues first described the alteration of the electrophoretic mobility of hemoglobin (Hb) by the **sickle cell** anemia mutation. The basis for the more positive electrical charge of sickle cell hemoglobin (Hb S) was not known, but the fact that it was a property of individual Hb molecules led Pauling to designate sickle cell anemia a molecular disease. In 1956 Vernon Ingram demonstrated that the charge difference is due to substitution of valine for glutamic acid in position six of the b-globin peptide chain. This observation was the first to prove that [mutations](#) can alter the amino acid sequence of a polypeptide chain and, by extension, it was proposed that the normal function of **genes** is to specify the [primary structures](#) of [proteins](#), a hypothesis that is well established now.

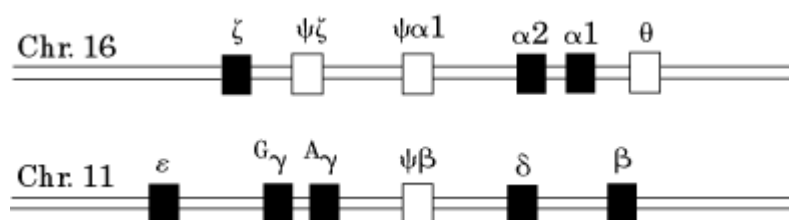
It is no accident that the initial connection between mutations and protein structure emerged from studies of Hb. In higher animals, Hb is the protein that is most readily available in relatively pure form. In the case of mammals, a lysate of washed red cells can be used for Hb classification without

further purification. Hb levels are measured as part of routine clinical examinations, and persons with diminished synthetic ability or with functionally aberrant Hb are likely to be identified and studied further. Hb is probably the protein that has been most extensively studied worldwide for genetic variations, and the array of mutations known is remarkable in number and in effect.

## 1. The Globin Loci

A [hemoglobin](#) molecule in vertebrates is tetrameric and consists of two each of two globin subunits. The folded [tertiary structure](#) of each globin chain binds one molecule of heme. The two kinds of globin subunits in Hb A, the major human adult Hb (actually present from birth onward), are labeled a-globin and b-globin. The a-globin polypeptide chain consists of 141 amino acid residues, and the b-globin chain has 146 residues. The molecular structure can be written  $\alpha_2\beta_2$ . An additional adult Hb, designated Hb A<sub>2</sub>, comprises 2.5% of the total. Its molecular structure is  $\alpha_2\delta_2$ , so Hb's A and A<sub>2</sub> differ only in having different “b-like” globins. Other Hb's present at various stages of human development are fetal Hb's,  $\alpha_2^G\beta_2$  and  $\alpha_2^A\beta_2$ , and embryonic Hb's,  $\alpha_2^\epsilon\beta_2$ ,  $\alpha_2^\zeta\beta_2$ , and  $\alpha_2^z\beta_2$ . All of the globin chains resemble each other in amino acid sequence and folding pattern and are similar to [myoglobin](#) of muscle, which occurs as a monomer. Genetically, the a and z genes form a cluster on human [chromosome](#) 16, and the b, d, G<sub>g</sub>, A<sub>g</sub>, and e genes form a cluster on chromosome 11 (Fig. 1). All of the functional genes have three **exons** and two **introns**. In addition to the functional genes, there are several [pseudogenes](#).

**Figure 1.** Diagrams of the organization of the a-globin complex on chromosome 16 and the b-globin complex on chromosome 11. Functional genes are in black. Gene sizes and relative positions are not drawn to scale.



The globin genes constitute a [gene family](#) that, along with myoglobin, evolved from a single ancestral gene through [gene duplication](#) and evolutionary **divergence**. Recent duplication has occurred in the case of the two a genes, which are identical, and the two g genes, whose products differ by a single codon. The b and d genes are also similar.

An interesting feature of the globin genes is the similarity in order of expression during development with order on the chromosomes. During the embryonic period, the z,  $\epsilon$ , g, and a genes are actively transcribed, forming  $\alpha_2^\epsilon\beta_2$ ,  $\alpha_2^z\beta_2$ , and  $\alpha_2^g\beta_2$  Hb's. The primary Hb during the fetal period is  $\alpha_2^G\beta_2$ . At about birth, but not as a result of birth, production changes to  $\alpha_2\beta_2$  and  $\alpha_2\delta_2$ . The molecular basis for these switches has not been established.

## 2. Mutations of Hb Genes

Over 500 different mutations have been observed in the Hb loci, exemplifying almost every known mutational mechanism. Those involving the b locus are especially apparent, because this locus becomes active at birth. If a mutation is limited to the b locus, embryonic and fetal development should be normal. In contrast, the a locus is active from the embryonic period. A major defect in its function is likely to interfere with fetal development and lead to abortion. The same is true of the g

loci. Defects in function that are limited to the d locus are likely to go unnoticed because of the small amount of Hb A<sub>2</sub> normally made.

The phenotypic consequences of Hb mutations can be classified roughly into four categories: (1) complete absence of a gene product, (2) normal function but diminished amount, (3) abnormal function with or without diminished amount; and (4) normal function and amount. The first group includes genetic deletions, which can extend over two or more loci. Single [base-pair substitutions](#) can produce any of the four effects, depending on the site of the substitution. The inability to produce adequate amounts of Hb is known as [thalassemia](#).

### 3. Absence of a Gene Product

The classical examples are the a- and b-thalassemias that involve deletions of all or part of the a and b genes, respectively. (See [Thalassemia](#).) Other causes are [frameshift mutations](#) that involve insertion or deletion of one or two nucleotides, which switch the [reading frame](#). Unless these changes occur very near the C-terminal end of the polypeptide, no recognizable product is produced. Several such mutations, known in the b globin gene, result in b thalassemia.

Single nucleotide substitutions also lead to complete absence of a gene product. In a number of thalassemia mutants, a nucleotide substitution in the coding region of a gene introduces a [stop codon](#), producing a truncated, nonfunctional, and unrecognizable gene product. If a nucleotide substitution occurs at a critical site in a regulatory region, [transcription](#) of the gene may not occur. This has been shown in several cases of thalassemia, in which substitutions in the [TATA box](#) and other 5' sites interfere with transcription. Substitutions in the 3' **poly-A** site leads to absence of product, presumably by interference with [translation](#).

One unusual cause for the absence of b-globin is nonfunction of the switch mechanism that causes transcription to change from the g genes to the d and b genes. In such cases, the g genes continue to function, and transcription of the d and b genes is not initiated. The major Hb is  $\alpha_2\gamma_2$  (fetal Hb or Hb F), and the condition is known as hereditary persistence of fetal hemoglobin. The Hb F is synthesized in amounts adequate for normal function, and homozygotes for this mutation are normal. The responsible mutations occur in two regions. One is a site between the g genes and the d gene. The other sites are in the **promoter** regions of the g genes.

### 4. Diminished Amount of a Normally Functional Hb

Mutations outside the coding regions of genes substantially suppress transcription without obliterating it completely. Several nucleotide substitutions in the promoter region of the b-globin gene do that, resulting in very low amounts of normal Hb A (so-called b<sup>+</sup>-thalassemia). A number of substitutions in the coding region also reduce the amount of globin. In some instances, the globin is quite functional. Hb C results from the substitution of lysine for glutamic acid in residue six of the b-chain. This does not interfere with function but the amount of Hb C produced is less than that of Hb A. Hb C is widespread in West Africa and offers some protection against falciparum malaria. Another polymorphic mutant, Hb E, occurs widely in East Asia and appears to protect against malaria also. Hb E results from a change in codon 26 of the b gene, near the boundary between exon 1 and intron 1. There are several cryptic splice sites in this region. In the case of the Hb E mutation, a cryptic splice site in codon 25 is substantially activated, and only a portion of the primary transcript is spliced correctly. Presumably this mechanism occurs in a number of other mutations that have not been studied so extensively.

### 5. Functionally Abnormal Globins

The prime example of a functionally abnormal Hb is **sickle cell Hb** (Hb S). Hb S transports oxygen normally. However, in the deoxygenated state, the Hb S molecules aggregate to form long crystal-



like structures that disrupt the red cell. The cause is substitution of a **hydrophobic** valine residue in place of glutamic acid in residue six of the b-globin polypeptide chain. This causes the N-terminal part of the b chain to fold differently, leading to aggregation.

Other examples are the mutations that cause **methemoglobinemia**. These involve amino acid substitutions in the heme-binding region. In normal hemoglobins, the Fe of the heme is stabilized in the Fe<sup>2+</sup> state, which is required for O<sub>2</sub> binding. In Hb M mutations, the Fe is predominantly in the more stable Fe<sup>3+</sup> form (methemoglobin, metHb). The amount of Hb M made is normal. Therefore a heterozygote has approximately half normal Hb A and half Hb M. Such a person is essentially normal, comparable to a person who is heterozygous for b-thalassemia. The high level of metHb gives them a gray-blue color, however. Hb M mutations are rather rare, and homozygotes have not been observed. Presumably, they would be lethal.

## 6. Mutations Normal in Function and Amount

The list of these mutations is not extensive because they normally escape detection unless populations are screened. An interesting example is Hb Lisbon, which results from the substitution of aspartic acid for glutamic acid at residue 23 of the a<sub>1</sub>-globin. This rare human variant is the standard “wild-type” form in gorillas. All other amino acids are identical in the a<sub>1</sub> globins of the two species.

Another interesting normal variation is the tandem duplication of one of the a-globin genes in populations in the southwest Pacific. Having five copies of the a-globin genes rather than four is a harmless variation.

### Suggestion for Further Reading

Web site: *On-Line Mendelian Inheritance in Man*, <http://www3.ncbi.nlm.nih.gov/Omim/>.

## Hemophilia

Hemophilia is the term used to describe uncontrolled bleeding, which may range from mild (bleeding after surgery or significant trauma) to severe (spontaneous bleeding into joints, muscle, and internal organs). Most examples are attributable to one of the two classic inherited hemophilias: the sex-linked disorders hemophilia A and hemophilia B. The most common of these is hemophilia A, which is caused by [mutations](#) in the **gene** on the human [X-chromosome](#) coding for the blood coagulation cofactor molecule, factor VIII (see [Blood Clotting](#)). Hemophilia B is caused by mutations to another X-chromosome-encoded gene, that for the [zymogen](#) of the blood coagulation [serine proteinase](#) factor IX. Less common congenital bleeding disorders are caused by mutations in the genes coding for other [proteinase](#) components of the blood coagulation system (factors VII, X, XI, and pro[thrombin](#)), the cofactor molecule factor V, [fibrinogen](#), and the [serpin](#) a<sub>2</sub>-antiplasmin.

Other causes of bleeding are due to defects in platelet aggregation caused by mutations in the genes coding for the [integrin](#) a<sub>IIb</sub> b<sub>3</sub> (glycoprotein IIb-IIIa, the platelet-fibrinogen binding site), the [cell adhesion molecule](#) von Willebrand factor, or its platelet binding-site glycoprotein Ib-IX.

The molecular basis of both hemophilia A and B has been studied extensively. Historically, bleeding defects have provided the basis for defining the reactions of the blood coagulation pathways and,

more recently, elucidating of structure–function relationships in the coagulation factors. Hemophilias A and B have also been valuable models in the study of mutational mechanisms. All possible types of gene disruption have been detected, including point mutations, deletions, insertions, duplications, and inversions, and they can lead to either absence of the protein, reduced levels of protein, protein with reduced functional activity, or nonfunctional protein.

In various populations, the frequency of hemophilia A is 1 in  $510 \times 10^3$  male births. Over 80 point mutations in the factor VIII gene have been detected. Of these, 27 cause severe hemophilia and are mainly [nonsense mutations](#), while 29 cause moderate to severe hemophilia and all are missense mutations. The most functionally interesting of the latter are those that lead to normal levels of protein but with severely reduced activity, that is, with reduced specific activity, which account for 5% of the total point mutations. These are all mutations of the Arg residue at the P1 position, immediately preceding the [peptide bond](#) to be cleaved, of the thrombin cleavage sites of factor VIII, and involve the hypermutable CpG dinucleotide (which is susceptible to [5-methylcytosine](#) deamination) (see [CpG Islands](#)).

The frequency of hemophilia B is one-sixth that of hemophilia A. More than 300 point mutations have been detected, 40% of which lead to factor IX with reduced specific activity; these mutations affect all the **domains** of factor IX. The mutations that can be most readily rationalized in structure–function terms include those leading to abnormal [post-translational modification](#) of the Gla domain (ie, generation of **g-carboxyglutamic acid** and b-hydroxyaspartic acid residues), mutation of the [arginine](#) residues at the proteolytic activation sites, and mutations affecting the [active-site](#) serine residue. Over half of these mutations are in CpG dinucleotides, with C → T or G → A [transition mutations](#).

#### Suggestion for Further Reading

E. G. D. Tuddenham and D. N. Cooper (1994) *The Molecular Genetics of Haemostasis and Its Inherited Disorders*, Oxford Monographs on Medical Genetics No. 25, Oxford University Press.

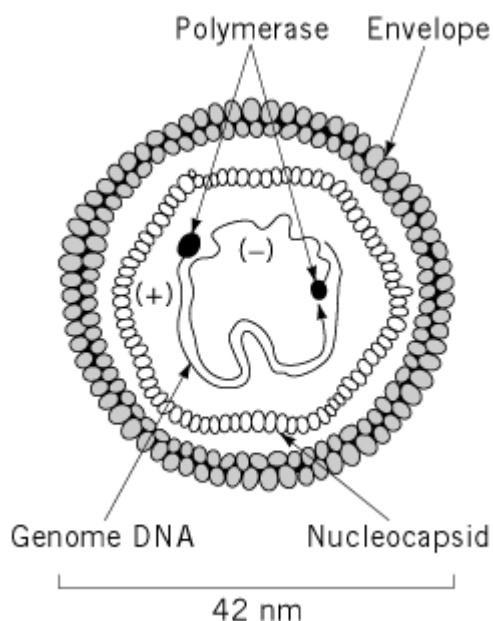
## Hepatitis B Virus

Hepatitis B virus (HBV) is one of the hepadnaviruses and causes acute, self-limited or persistent infection that can potentially be passed vertically from carrier mother to infant, or horizontally by the percutaneous and transmucosal exchange of infected blood and body fluids. Persistent infection by HBV is often associated with chronic liver disease that can lead to the development of hepatocellular carcinoma (HCC). The neonatal transmission route is thought to be an important mechanism for HBV persistence because of the immunological immaturity of the infected host. The concomitant infection is usually self-limited, with clearance of virus particles from the liver and blood and the development of lasting immunity to reinfection.

HBV is a spherical particle, 42 nm in diameter (Fig. 1). The genome of HBV is a circular DNA duplex of 3.2 kbp that replicates by **reverse transcription** (see [HIV \(Human Immunodeficiency Virus\)](#) and [Rous Sarcoma Virus \(RSV\)](#)) of a pre-genome RNA and contains a single-stranded gap region of variable length in one strand. [DNA replication](#) and [transcription](#) of viral genes and the mechanism of virion assembly within infected hepatocytes have been unraveled so some extent, but less is known about the virus–cell interactions that control virus attachment, uncoating, and entry into susceptible cells, as well as the integration of HBV DNA into the host cell [genome](#). Although viral sequences implicated in the interaction of the virion with the cellular [membrane](#) are located in a

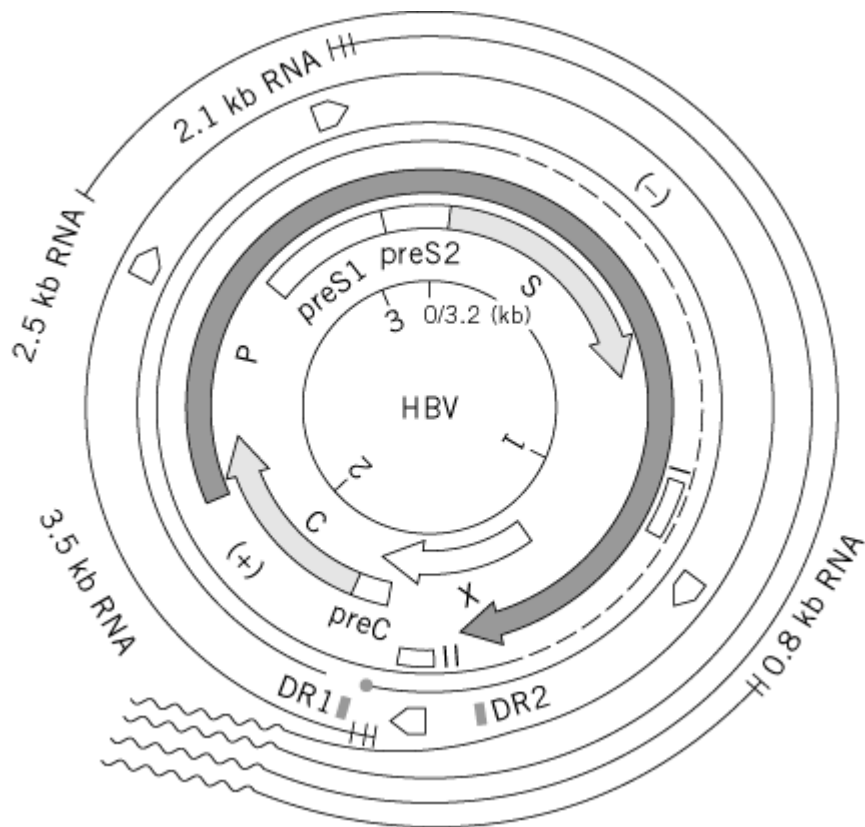
variable region of the amino-terminal part of the pre-S1 domain, the cellular receptor for HBV has not yet been fully characterized.

**Figure 1.** Schematic diagram of a hepatitis B virus particle.



**Nucleotide sequences** and those amino acid sequences deduced from HBV DNA from different virus subtypes (adw, adr, ayw, ayr) have revealed a compact, overlapping coding organization of the genome, where the entire genome is coding (Fig. 2). Four open reading frames (ORFs) are localized on one viral strand in the same transcriptional polarity. Two of them, the C and S regions, specify structural proteins: the 183-residue nucleocapsid protein HBcAg and the 226-residue envelope major protein HBsAg, respectively. The smallest ORF, X, codes for a [trans-acting](#) 154-residue regulatory protein that interacts with various cellular proteins and may play a major role in hepatic cell transformation. Overlapping these coding regions, the largest ORF, P, encodes a viral polymerase of 843 residues, which is necessary for viral replication and has activities of **DNA polymerase**, **reverse transcriptase**, plus [ribonuclease H](#). This overlapping of the P region is a characteristic feature and, with its unusual overlapping of coding and regulatory sequences (such as **promoters**, two specific [enhancers](#), and a [poly A](#) termination signal), makes this virus unique. In addition to the central homologous region, two strong [homology](#) blocks in the amino- and carboxy-terminal parts of the X protein might correspond to a conservative selective pressure for functional activity of the X as a *trans*-regulator, and of the viral ribonuclease H encoded by C-terminal sequences of P, which overlap with the 5' end of the X protein.

**Figure 2.** Genetic organization and transcription map of the hepatitis B virus genome. The HBV DNA was derived from virus of the adr subtype, with a nucleotide sequence of 3215 bp. The positions of promoters and two enhancers are indicated by *p* and *f*, respectively. The four different transcripts are indicated with a poly(A) tail (~). The terminal protein domain of polymerase (•) is attached to the 5' end of the minus strand (—). DR1 and DR2 represent 11-bp direct repeats.



Another important characteristic of HBV is a high [mutation](#) rate due to the lack of a proofreading function of the reverse transcriptase that replicates its genome. This may increase the production of mutated variants during persistent infection, which enables the virus to escape from the host immune response or enables a growth advantage of mutant viral genomes. Substantial nucleotide and amino acid sequence heterogeneity in different isolates of HBV, plus specific mutations in the HBV genome in individual patients, have been observed, and partially deleted HBV genomes have also been observed in human hepatocellular carcinoma. However, there is no experimental evidence that these defective HBV genomes might harbor oncogenic potential, nor that any particular HBV subtype might be more oncogenic than others. Comparison of the HBV genome with those of animal hepadnaviruses reveals extensive homologies among mammalian hepadnaviruses (WHV, GSHV), which share a basically identical genomic organization.

The [antibody](#) response to HBV envelope antigen (HBsAg) is well characterized to be **T-cell**-dependent and plays an important role in clearance of the circulating virus and preventing their attachment and uptake by susceptible hosts. The role of the antibody response to the nucleocapsid (HBcAg and its derivative in serum, HBeAg) and nonstructural proteins (X and P) in HBV pathogenesis is less clear yet. Persistent infection of HBV is mainly caused by a weak antiviral immune response that cannot kill or cure all of the infected cells. Alternatively, chronic hepatitis may be caused by an inappropriate T-cell response that can kill but cannot cure the infected hepatocytes without enough production of [cytokines](#) to clear the infection.

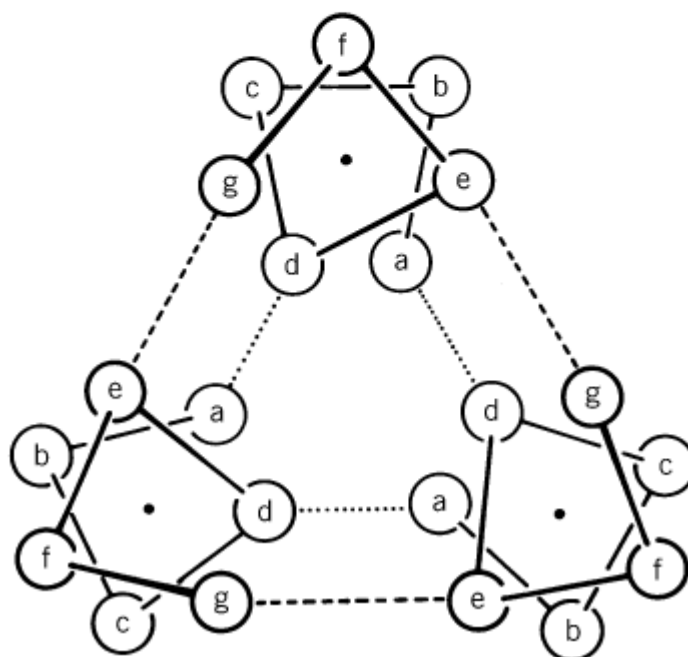
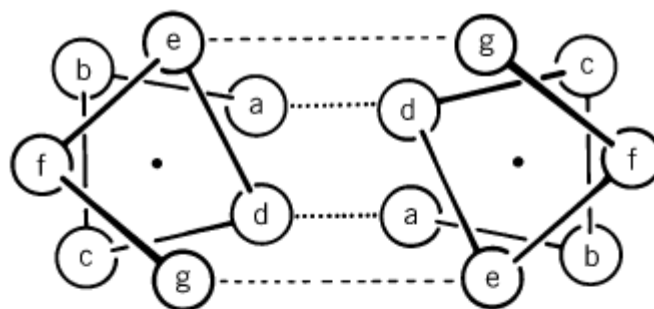
#### Suggestion for Further Reading

D. Ganem (1996) "*Hepadnaviridae* and Their Replication". In *Fields Virology*, 3rd ed. (B. N.Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2703–2737.

## Heptad Repeat

[Fibrous proteins](#) of the  $\alpha$ -type, as well as many globular proteins, contain a heptad repeat of the form  $(a-b-c-d-e-f-g)_n$  in which [nonpolar](#) amino acid residues are found with high probability in positions  $a$  and  $d$ . Regions displaying a heptad repeat invariably have the  $\alpha$ -**helix** conformation and give rise to a **coiled-coil** molecule in which the apolar residues in  $a$  and  $d$  in several chains pack in a regular manner along the axis (1) (Fig. 1) (see also [Cytokeratins](#)). [Electrostatic interactions](#) between charged residues in positions  $e$  and  $g$  of adjacent polypeptide chains help to specify the relative orientation and axial stagger of the chains. Various positions within the heptad repeat are not occupied equally by different amino acids (1-3), and this has provided the basis for distinguishing between two- and three-stranded coiled-coils in fibrous proteins (4). Consider the  $a$ ,  $d$ ,  $e$ , and  $g$  positions in turn. Leucine, lysine, asparagine, and arginine residues occur more commonly in position  $a$  in two-stranded coiled-coils than in three-stranded ones, whereas the reverse is true for alanine, glutamine, serine, and threonine. In the  $d$  positions, leucine is favored in two-stranded coiled-coils, but isoleucine, serine, threonine, and valine are found more commonly in the three-stranded conformations. In  $e$  positions, glutamic acid, isoleucine, and glutamine are more commonly observed in two-stranded ropes than three-stranded ones, but the reverse is true for phenylalanine, leucine, serine, tryptophan, and tyrosine. Finally, in position  $g$ , glutamic acid, glutamine, and arginine are favored in two-stranded ropes, whereas isoleucine, leucine, asparagine, tryptophan, and tyrosine occur significantly more frequently in three-stranded coiled-coils.

**Figure 1.** Heptad repeat seen in axial projection for two- and three-stranded  $\alpha$ -helical coiled-coils. Residues in positions  $a$  and  $d$  are generally apolar and are shielded from water as a result of systematic close packing along the axis of the rope-like structure so formed (dotted lines). Interchain ionic interactions are also possible, the most common type being those that occur between oppositely charged residues in positions  $e$  and  $g$  (dashed lines). (From Ref. 66, with permission.)



Groupings of amino acids also clearly provide a means for distinguishing one multichain coiled-coil from another (Table 1). These figures indicate, for instance, that in two-stranded structures the acidic residues occur very infrequently in position *a* but quite commonly in position *d*. In contrast, basic residues are found commonly in *a* but much less so in *d*. In three-stranded coiled-coils, the preferences for the charged residues are almost reversed in these positions. In all cases the *a* and *d* positions are occupied by apolar residues at a frequency of about 75%. In three-stranded structures the occupancy of charged residues in positions *e* and *g* is reduced, and the occupancy of apolar residues is increased. This trend continues as the number of strands in the coiled-coil rope increases.

**Table 1. Percentage of Apolar (ap), Acidic (ac), and Basic (ba) Residues in Positions *a*, *d*, *e*, and *g***

| Position | Group | Two-Stranded Structure | Three-Stranded Structure |
|----------|-------|------------------------|--------------------------|
| <i>a</i> | ap    | 72.6                   | 77.5                     |
| <i>a</i> | ac    | 1.4                    | 3.1                      |
| <i>a</i> | ba    | 13.8                   | 2.2                      |

|          |    |      |                   |
|----------|----|------|-------------------|
| <i>d</i> | ap | 76.1 | 79.3              |
| <i>d</i> | ac | 6.3  | 1.0               |
| <i>d</i> | ba | 3.8  | 3.9               |
| <i>e</i> | ap | 16.9 | 20.9              |
| <i>e</i> | ac | 27.9 | 22.9              |
| <i>e</i> | ba | 23.7 | 22.5              |
| <i>g</i> | ap | 20.0 | 25.8              |
| <i>g</i> | ac | 32.3 | 22.7              |
| <i>g</i> | ba | 19.0 | 17.3 <sup>a</sup> |

<sup>a</sup> Acidic residues (ac) are aspartic acid and glutamic acid; basic residues (ba) are arginine and lysine; apolar residues (ap) are leucine, isoleucine, valine, methionine, phenylalanine, tyrosine, and alanine.

A run of consecutive heptad repeats is frequently interrupted, sometimes by a region of sequence with different structural characteristics and sometimes by the insertion/deletion of a few residues in an otherwise perfect and continuous heptad substructure. Analysis shows, however, that the vast majority of such discontinuities fall into one of three categories (5): (i) the deletion of three residues (a “*stutter*”), (ii) the insertion of one residue (a “*skip*”) which is structurally equivalent to two stutters, and (iii) the deletion of four residues (a “*stammer*”). Stutters lead to a local unwinding of the coiled-coil structure, whereas the stammer results in a local overwinding. A skip thus leads to a longer distance axially over which the coiled-coil is partially unwound.

#### Bibliography

1. C. Cohen and D. A. D. Parry (1986) -Helical coiled-coils—a widespread motif in proteins. *Trends Biochem. Sci.* **11**, 245–248.
2. J. F. Conway and D. A. D. Parry (1990) Structural features in the heptad substructure and longer range repeats of two-stranded -fibrous proteins. *Int. J. Biol. Macromol.* **12**, 328–334.
3. J. F. Conway and D. A. D. Parry (1991) Three-stranded -fibrous proteins: the heptad repeat and its implications for structure. *Int. J. Biol. Macromol.* **13**, 14–16.
4. D. N. Woolfson and T. Alber (1995) Predicting oligomerization states of coiled coils. *Protein Sci.* **4**, 1596–1607.
5. J. H. Brown, C. Cohen and D. A. D. Parry (1996) Heptad breaks in -helical coiled coils: stutters and stammers. *Proteins Struct. Funct. Genet.* **26**, 134–145.
6. C. Cohen and D. A. D. Parry (1990) -Helical coiled-coils and bundles: how to design an -helical protein. *Proteins Struct. Funct. Genet.* **7**, 1–15.

#### Suggestions for Further Reading

7. C. Cohen and D. A. D. Parry (1990) -Helical coiled-coils and bundles: how to design an -helical protein. *Proteins Struct. Funct. Genet.* **7**, 1–15.
8. A. Lupas (1996) Coiled-coils: new structures and new functions. *Trends Biochem. Sci.* **21**, 375–382.

#### Hermaphrodite

In the animal kingdom, [natural selection](#) has promoted the development of two opposite [sexes](#), the [male](#) and [female](#) genders. In most species, individuals belong to either sex, developing either a female **phenotype** with female external genitalia: a uterus, tubes, and two ovaries, or, depending on the presence of a testis-determining factor, a male phenotype with male external genitalia: epididymis and a testis (see [Sex](#)). In a hermaphrodite, either in a true hermaphrodite or pseudohermaphrodite, elements of both genders are present in one individual. Some [nematodes](#) are physiologically true hermaphrodites (eg, *Caenorhabditis elegans*). Here, the molecular switch to determine whether a male or hermaphrodite is formed depends on an [X-chromosome](#) counting system (1), which involves primarily the xol-1 (2, 3) and tra-1 (4) genes. In mammals, however, male and female organs are not normally present in one individual.

Due to the severity of the aberration, different types of hermaphrodites can be found. The clinical appearance of a hermaphrodite is characterized by intermediate sex forms of the external genitalia and **germ line** cells. Gonads that diagnose a true hermaphrodite are ovary and testis, or a combination of the two, the ovotestis. The combination of an ovary and ovotestis is most commonly found in hermaphrodites, followed by the combination of two ovotestes. When close to an ovary, the tubes and rudiments of an uterus are most commonly detected; when close to a testis, a vas deferens and an epididymis are present. In some cases, however, with a lower frequency of appearance, the opposite gender organs, or none, are developed: Tubes and uteri and/or vas deferens and epididymis are found, together with an ovotestis. This can occur because sex differentiation is highly dependent on [hormones](#). As sexual organs are formed from primordial **germ cells**, as well as from soma cells, the gonadal sex does not fit to the genetical and/or cytological sex under certain circumstances.

Besides true hermaphrodites, there are also pseudohermaphrodites, in which all other organs are female but a testis has developed. Due to the higher temperature in the body than in the scrotum, these cells are prone to the development of cancer and will be removed.

Regarding the [karyotype](#) of hermaphrodites, those with a [Y-Chromosome](#) are more likely to develop a testis than those without. The ovotestis contains both ovarian and testicular cells. The ovarian cells are in many cases normal, and primary follicles are also found. In contrast, the testicular cells show histological aberrations in true hermaphrodites, regardless of whether these cells are found in a testis or an ovotestis.

Regarding the hormone production of hermaphrodites, they are characterized by normal levels of estrogens and gestagens. However, the levels are also dependent on the amount of normal ovarian stroma present in the hermaphrodite. Spermatogenesis does not occur, but testosterone is produced, although it does not reach the levels of normal males. The reduced amounts are responsible for the virilization. There is only one case reported in which a hermaphrodite became a father. However, there are hermaphrodites that are fertile after removal of the testicular tissue and can give birth to children. This certainly does not depend solely on the presence of an ovary, but also on a physiological development of at least one of the tubes and of the uterus. Interestingly, all children born from female hermaphrodites were male. Roughly half the hermaphrodites were identified before the third year of life, but many hermaphrodites were identified only during or after puberty.

During embryogenesis, germ cells become either spermatogonia or oogonia, depending on their genetic constitution. The **genotype** sex normally determines which type of gonads are developed. Further differentiation depends on the hormonal secretion of the gonads themselves, which consequently determines the **phenotypic** sex. The primordial gonad develops to an ovary, and the inner and external genitalia also tend to feminize, so long as the production of male hormones is not stimulated. In the case of secretion of male hormones, especially testosterone and mullerian inhibiting hormone, mullerian ducts do not develop; instead, wolffian ducts are formed that change the fate of the organs, and epididymis, vas deferens, and seminal vesicles evolve.



In cattle and goats, a special situation is found in a male and female twin. Due to the transfer of anastomoses, the development of müllerian ducts is inhibited in the female twin, which results in an infertile female. The external genitalia are female, but the uterus is rudimentary or missing. The male twin is fertile under normal circumstances. Apart from these special cases in twin births in cattle and goats, the biological reason for pathological hermaphroditism in mammals is still unclear.

Other genetically determined disorders related to gender are the Turner (X0) syndrome (5, 6) and surplus X-chromosomes, such as the XXY constitution (Klinefelter syndrome), which can now be detected using PCR techniques (7). Due to the absence of the Y-chromosome, the Turner syndrome is mostly linked with a female phenotype, although a mild masculinization has been reported (8). An additional X-chromosome leads to infertility and a reduced IQ. Additional Y-chromosomes, as found in XYY karyotypes, lead to extreme tallness, but not necessarily to crime, as was assumed in former times. In every case, the presence of a Y-chromosome leads automatically to a male phenotype, regardless of how many X-chromosomes were inherited. However, it must be kept in mind that it is not the Y-chromosome itself that triggers the development of male fate, but the SRY gene (see Sex). Mutations within this gene can disrupt its function and consequently lead to female development, even though a Y-chromosome is present.

### Bibliography

1. M. Nicoll, C. C. Akerib, and B. J. Meyer (1997) *Nature* **388**, 200–204.
2. L. M. Miller, J. D. Plenefisch, L. P. Casson, and B. J. Meyer (1988) *Cell* **55**, 167–183.
3. N. R. Rhind, L. M. Miller, J. B. Kocpczynski, and B. J. Meyer (1995) *Cell* **80**, 71–82.
4. C. P. Hunter and W. B. Wood (1990) *Cell* **63**, 1193–1204.
5. T. Ogata and N. Matsuo (1995) *Hum. Genet.* **95**, 607–629.
6. C. Geerkens, W. Just, K. R. Held, and W. Vogel (1996) *Hum. Genet.* **97**, 39–44.
7. A. Kleinheinz and W. Schulze (1994) *Andrologia* **26**, 127–129.
8. R. Medlej, J. M. Lobaccaro, P. Berta, C. Belon, B. Leheup, J. E. Toublanc, J. Weill, C. Chevalier, R. Dumas, and C. Sultan (1992) *J. Clin. Endocrinol. Metab.* **75**, 1289–1292.

### Herpesvirus

Herpesviruses are large DNA **viruses** whose **genomes** consist of a large linear double-stranded DNA molecule, 125 kbp to 229 kbp in length. Their hosts range from lower vertebrates to humans. At present, approximately 100 herpesviruses have been isolated and at least partially characterized. In humans, five viruses [herpes simplex virus 1 (HSV-1) and 2 (HSV-2), varicella zoster virus (VZV), human **cytomegalovirus** (HCMV) and **Epstein–Barr virus** (EBV)] have been identified as members of the herpesvirus family, *Herpesviridae*. In the past decade, however, three new herpesviruses, namely, human herpesvirus 6, 7, and 8 (HHV-6, HHV-7, and HHV-8), were isolated from healthy people and AIDS patients. These human herpesviruses are responsible for a wide variety of diseases.

Although herpesviruses share common features in their structure, gene organization, replication style, and so forth, the family *Herpesviridae* is further divided into three subfamilies, *Alphaherpesvirinae*, *Betaherpesvirinae*, and *Gammapherpesvirinae*. *Alphaherpesvirinae* includes HSV-1, HSV-2, VZV, pseudorabies virus, equine herpesvirus 1 (EHV-1), and so on. They have a relatively short replication cycle, induce cytotoxic damage in infected cells, and establish latency primarily in sensory ganglia. *Betaherpesvirinae* includes HCMV, HHV-6, HHV-7, and so on, and is characterized by its restricted host range, long replication cycle, and slow spread of infection from

cell to cell in culture. They can establish latency in lymphoreticular cells, secretory glands, kidneys, and other tissues. *Gammaherpesvirinae* contains EBV, HHV-8 (Kaposi's sarcoma-associated herpesvirus), herpesvirus saimiri (HVS), and so on. They replicate in lymphoid cells, and latency is frequently demonstrated in lymphoid tissue. A most important feature of this subfamily is that some are associated with naturally occurring malignant tumors.

The virions of herpesviruses consist of four morphological elements: an inner core, an icosahedral capsid with 162 capsomers, a surrounding amorphous layer, known as the tegument, and an envelope containing a number of **glycoproteins**. The tegument of the HSV virion is known to contain about 20 distinct structural proteins, some of which have important regulatory functions for the initiation of viral replication. The **genomes** of eight human herpesviruses have been entirely sequenced, and the complete DNA sequences are also available for several animal herpesviruses, including EHV-1, EHV-2, HVS, and mouse CMV. Although the number of viral genes encoded in the genomes varies, a set of approximately 40 genes is conserved in all herpesviruses. These encode capsid proteins, envelope glycoproteins, proteins involved in DNA cleavage/packaging, and **enzymes** such as **DNA polymerase**, **DNA helicase**, **DNA primase**, uracil-DNA glycosylase, dUTPase, **ribonucleotide reductase**, and protein **kinase**. The herpesvirus genomes contain a number of accessory genes that are dispensable for viral replication in cell culture. For example, 45 of 84 HSV genes have been shown to be nonessential. These accessory gene products are supposed, however, to be important for viral growth and spread in their natural hosts.

The principal strategy of herpesvirus replication is similar in all of them and can be summarized as follows: (i) Fusion of the virion envelope with the plasma membrane occurs at the cell surface in a pH-independent manner; (ii) viral capsids with a portion of the tegument are transported to a **nuclear pore**, where viral DNA is released into the **nucleus**, circularized, and then **transcribed** by cellular **RNA polymerase II**; (iii) viral gene expression is coordinately regulated and sequentially ordered in a cascade fashion; (iv) viral DNA replicates by a **rolling circle** mechanism, and unit-length viral DNA is cleaved from newly synthesized concatemers and packaged into preformed empty capsids; (v) full capsids pass through the inner and outer membranes of the **nuclear envelope** by budding and re-envelopment and associate with tegument proteins in the perinuclear region; (vi) enveloped virions accumulate in the **endoplasmic reticulum** and then mature in the **Golgi apparatus**; and finally, (vii) mature virions are released into the extracellular space by **exocytosis**.

All herpesviruses have the ability to persist in their natural hosts throughout life. During latency, the viral genomes are present in the form of an **episome** in the nucleus, and only a small subset of viral genes is transcribed. The synthesis of specific viral proteins appears to be essential for some herpesviruses to maintain latency, but not for others. Various stimuli, such as stress and immunosuppression, induce reactivation, which triggers viral replication and shedding of infectious virus. The molecular mechanism of herpesvirus reactivation remains to be elucidated.

#### Suggestions for Further Reading

B. Roizman (1996) "Herpesviridae" . In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2221–2230.

B. Roizman and A. E. Sears (1996) "Herpes Simplex Viruses and Their Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2231–2295.

Y. Becker and G. Darai (1994) *Pathogenicity of Human Herpesviruses due to Specific Pathogenicity Genes*, Springer-Verlag, Berlin.

## Heterochromatin

The term *heterochromatin* was originally used to describe [chromatin](#) that stained intensely with particular dyes. This is in contrast to the rest of chromatin, which stained less well and is called euchromatin. Heterochromatin comes in two varieties, constitutive, which is permanently heterochromatic, and facultative, which can change its state. The chromatin regions of centromeres and telomeres are examples of constitutive heterochromatin. In contrast, the inactive X-chromosome (see [Barr Body](#), [X-Chromosome Inactivation](#)) and the vast bulk of inactive chromatin in any chromosome are composed of facultative heterochromatin.

The most common observed consequence of heterochromatin formation is the repression of transcription, either in the heterochromatin itself or in regions of chromatin adjacent to the heterochromatin domain. The variability in gene expression at the border of the heterochromatin is described as “position effect variegation.” Three explanations have been offered for this phenomenon. The first is that special proteins, such as heterochromatin protein 1 (HP1, see following text), cause heterochromatin to adopt its distinct structure and that these proteins can “spill over” into regions of normal chromatin and exert repressive effects on gene expression. The second is that heterochromatin represents the sequestration of chromosomal domains in specialized nuclear compartments from which the transcriptional machinery is excluded, for example, the chromocenter. The third applies only to *Drosophila* and other insects with polytene chromosomes in that, following placement adjacent to heterochromatin, a gene undergoes fewer rounds of replication than normal. Fewer copies of the gene would cause a concomitant reduction in transcription. Most investigators accept that heterochromatin-specific proteins diffuse onto normal chromatin and thereby influence the gene expression of juxtaposed genes. Although initially defined from work on insects, there is an increasing body of evidence that heterochromatin domains exist in all eukaryotic chromosomes, including those of yeast. Moreover, it is clear that metazoans have used this type of chromosomal organization for the developmental regulation of gene expression.

An approach to the molecular basis of heterochromatin formation has been to look for mutations in *Drosophila* or more recently in *Saccharomyces cerevisiae* that enhance or suppress position effects on gene expression (modifiers of position effect variegation). Modifiers of position effect variegation include mutations in the genes encoding the chromosomal proteins involved in forming heterochromatin. Using this approach, the gene encoding a nonhistone protein, HP1, was identified (see [Chromocenter](#)). Mutations of the HP1 gene reduce position effects on gene expression (1). HP1 is preferentially associated with the heterochromatin regions of polytene chromosomes. Other proteins homologous to HP1 include Polycomb, which is also chromatin-associated and from genetic experiments it is known that it influences the expression of many genes in normal chromatin (2). Neither HP1 nor Polycomb interacts with DNA directly, but they presumably recognize some aspect of nucleosome or chromatin fiber structure. Recent experiments reveal that certain proteins that affect position effect variegation, for example, SUV39H1, are histone methyltransferases that selectively modify lysine #9 in histone H3 (3), and that the resulting modified histone is a preferential target for binding by HP1 (4). Combined with earlier data that overexpression of this enzyme leads to HP1 redistribution from heterochromatin (5), as does the histone hyperacetylation via inhibition of histone deacetylases (6), these data strongly support a model that idiosyncratically modified nucleosomes form a landing pad for HP1. Polycomb and HP1 share a common amino acid sequence known as the chromodomain (chromatin modification) (see [Chromocenter](#)). This domain is highly conserved through evolution and can be found even in yeast. For example, the *S. pombe* SW16 gene encodes a chromodomain protein. The SW16 protein is involved in the assembling the chromosomal domain containing the transcriptionally silent mating type cassettes.

Position effect variegation also depends on the presence of normal levels of the histone proteins and posttranslational modifications of the N-terminal tails of the core histones, such as acetylation. The histones are needed for forming nucleosomes and for subsequently forming higher order structures. Acetylating the N-terminal tails alters the interaction of the histones with DNA and, potentially, the folding of nucleosomal arrays. Position effect variegation also depends on proteins that recognize

nucleosomal arrays, but not naked DNA, such as HP1 and Polycomb (see [Position Effect](#)). HP1 is associated only with inactive chromosomal regions, but Polycomb interacts with at least 50 different sites within *Drosophila* chromosomes, including developmentally important loci encoding homeodomain proteins. However, mutations in a second gene family that interacts with the histone proteins allow genes associated with Polycomb to remain active. Thus, Polycomb is likely to exert its effects on gene expression via chromatin structures dependent on the histone proteins (7, 8).

Position effect variegation is recognized now as a universal phenomenon in eukaryotic chromosomes. Genes integrated into yeast chromosomes near the silent mating loci or close to the telomeres are repressed in a way that reflects their proximity to these sites in the chromosome. This silencing effect spreads over at least 5 to 10 kbp of contiguous DNA, but not as much as 20 to 30 kbp in yeast. As in *Drosophila*, the genes influencing the position effect in yeast encode structural components of chromatin or enzymes associated with the modification of chromatin (see [Position Effect](#)).

Examples of mammalian chromosomal regions that contain heterochromatin are also the centromere and the telomere (9). The telomeres of mammalian chromosomes have an unusual chromatin structure in which nucleosomes are closely packed (see [Telomere](#)). Mammalian telomeres consist of the sequence (TTAGGG)<sub>n</sub> repeated for 10 to 100 kbp. Heterochromatin at the centromere contains tandemly repeated simple sequence “satellite” DNA, for example the a-satellite DNA at the human centromere (see [Centromeres](#)). This a-satellite heterochromatin plays a structural role by mediating attachment of the kinetochore. In these instances, specialized heterochromatin structures have an important architectural role and protective function in the chromosome.

Occasionally, an entire nucleus becomes heterochromatinized. One example is the inactivation of the erythrocyte nucleus in chicken. Here the special linker histone variant H5, which represses transcription and compacts nucleosomal arrays very effectively, accumulates. Histone H5 is more arginine-rich than the normal linker histone H1 found in somatic cells. This increase in arginine content probably strengthens the interaction of H5 with DNA and stabilizes chromatin structure. In this instance, the assembly of a specialized heterochromatin structure reflects the terminal differentiation of this particular specialized cell type.

#### Bibliography

1. P.B. Singh, *J. Cell Sci.* **107**, 2653–2668 (1994).
2. S.C.R. Elgin, *Curr. Opin. Cell Biol.* **2**, 437–445 (1990).
3. S. Rea et al., *Nature* **406**, 596–599 (2000).
4. M. Lachner et al., *Nature* **410**, 116–120 (2001).
5. M. Melcher et al., *Mol. Cell. Biol.* **20**, 3728–3741 (2000).
6. A. Taddei, C. Maison, D. Roche, and G. Almouzni, *Nat. Cell Biol.* **3**, 114–120 (2001).
7. H.W. Brock and M. van Lohuizen, *Curr. Opin. Genet. Dev.* **11**, 175–181 (2001).
8. B. Zink and R. Paro, *Nature* **337**, 468–471 (1989).
9. D. Shore, *Curr. Opin. Genet. Dev.* **11**, 189–198 (2001).

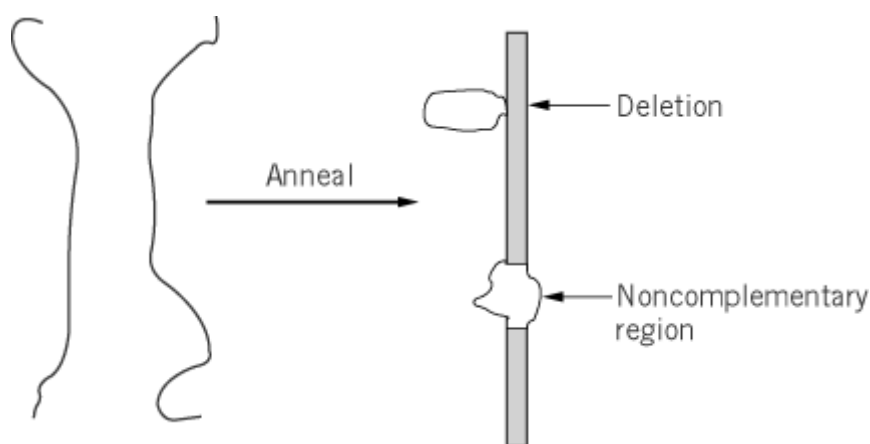
#### Additional Reading

10. Wolffe A., *Chromatin: Structure and Function*, 3rd ed., Academic Press, London, U.K., 1998.

#### Heteroduplex

Two single-stranded nucleic acid polymers with at least partial complementarity can form duplexes with stretches of double helix, corresponding to the regions of nucleotide sequence [homology](#), interleaved with unpaired stretches, corresponding to noncomplementary regions. Two types of single-stranded loops are introduced by noncomplementarity in the heteroduplex (Fig. 1). A loop in only one of the strands indicates a deletion. A loop composed of single-strand regions in both strands indicates a region of sequence incompatibility. A heteroduplex can be composed of two DNA strands, one DNA strand and one RNA strand, or two RNA strands. *In vivo*, heteroduplexes (DNA–DNA) can arise in the process of [recombination](#) or as a consequence of errors (mismatches) in [DNA replication](#) (DNA–DNA) or [transcription](#) (DNA–RNA). Hybridization of nucleic acids of different origins leads to heteroduplexes *in vitro*.

**Figure 1.** Loops in a heteroduplex indicate either a deletion or a region of noncomplementarity in the two strands.



The level of sequence homology in the heteroduplex can be evaluated by monitoring the change in ultraviolet **absorbance** on heteroduplex formation or by selective enzymatic digestion of single-stranded regions, followed by [hydroxyapatite chromatography](#). These techniques are analogous to those used for studies of reassociation kinetics (1).

[Electron microscopy](#) can be used to visualize the regions of sequence homology in the heteroduplex, even though nucleic acids are not sufficiently dense to be easily observable directly in the electron microscope. The heteroduplex is spread on the surface of the electron microscope grid. Heavy metals and proteins are applied to the nucleic acid by one of several techniques, thereby making the polymers visible in the electron microscope.

#### Bibliography

1. R. J. Britten, D. E. Graham, and B. R. Neufeld (1974) *Meth. Enzymol.* **29**, 363.

#### Suggestion for Further Reading

2. R. W. Davis, M. Simon, and N. G. Davidson (1971) Electron microscope heteroduplex methods for mapping base sequence homology in nucleic acids, *Meth. Enzymol.* **21**, 413–428.

## Heterokaryon

A *heterokaryon* is a combined cell with two separate **nuclei**. Pontecorvo (1) who worked with the **fungus** *Aspergillus nidulans*, which lacks a sexual cycle, observed that genetic **recombination** occurs during cell division in **mitosis** and employed it in the genetic analysis of fungi. If asexual **spores** from two genetically dissimilar strains of a filamentous fungus are planted very close together on the surface of a solid medium, the emerging hyphae fuse at points of contact. This fusion is followed by the transfer of nuclei from one hypha to the other, so that *heterokaryons* are formed whose hyphae contain nuclei from both strains. As the heterokaryotic hyphae grow, the nuclei multiply independently.

When asexual spores are formed, however, each initial sporogenous cell receives only one nucleus of one type or the other. Subsequent multiplication of these cells thus produces chains of spores. Each chain consists of genetically identical spores, **homokaryons**, although the spores of adjacent chains may differ in phenotype. Very occasionally, however, two of the nuclei in the heterokaryotic hyphae fuse to form a single **diploid** nucleus. This nucleus multiplies as such and is incorporated into spores in the same way as the haploid nuclei, so rare chains of diploid spores arise which, on subsequent germination, yield diploid individuals. If the haploid parental strains differ in complementary nutritional requirements, eg, one of genotype Ab (requiring B in the growth medium) and the other aB (requiring A), the diploid nucleus (Ab/aB) will possess genetic determinants to synthesize both A and B substances needed for growth of the parents. Therefore the existence of diploid spores can be recognized and cultures obtained from them by plating on media lacking the growth factors A and B (2).

Cell fusion is usually carried out by treating a suspension of cells with some inactivated **viruses** or with polyethylene glycol, to alter the plasma membrane of cells so that they are induced to fuse with each other.

Heterokaryons provide a way of mixing the components of two separate cells to study their interactions (see [Hybrid Cell](#)).

### Bibliography

1. G. Pontecorvo (1954) *Caryologia*, suppl. 6, 192–200.
2. J. A. Roper (1992) *Experientia* **8**, 14–15.

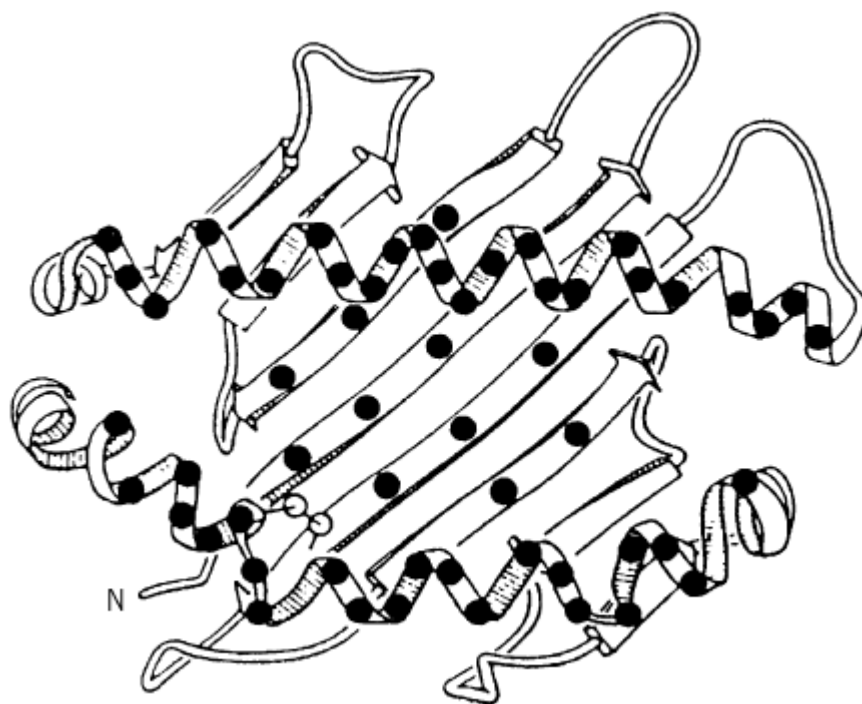
## Heterosis

Heterosis is the phenomenon in which **heterozygotes** exhibit superiority of fitness over **homozygotes**. This superiority often occurs in the areas of viability, longevity, fecundity, and resistance to disease. Therefore, heterosis is sometimes referred to as “hybrid vigor”. Heterosis occurs when two different, often **inbred** strains are crossed. The effect of homozygous deleterious alleles is reduced in the heterozygotes produced through the cross. The increase in vigor may be due to overdominance or superiority of the heterozygote for particular gene differences, or to the introduction of favorable dominant alleles to loci previously homozygous for deleterious recessive **alleles**. Therefore, the terms “heterosis” and “overdominance” are used interchangeably; both confer heterozygotes with an evolutionary advantage over homozygotes.

When survival or reproduction is lowered with recessive alleles, such alleles are eliminated faster from inbred homozygous populations than from outbred, heterozygous populations. Homozygous deleterious alleles will first increase, and then [natural selection](#) will tend to remove the deleterious alleles, resulting in less genetic variation and greater average fitness as the heterozygous alleles increase. However, deleterious recessive genes will not be totally eliminated by selection, because they will not be exposed to selection in the heterozygotes.

The [major histocompatibility complex](#) (MHC) genes in humans and [mouse](#) are known to have extremely high numbers of alleles (Figure 1) (1). Since the number of alleles is much greater than expected if the alleles were all equivalent functionally, it has been long speculated that positive selection is operating on these complex genes. By making combinatorial sets of different alleles, an organism can become resistant to viral infection, which confers superiority of fitness over the homozygote. Thus, the MHC genes are counted as one example of overdominance.

**Figure 1.** The protein structure of the MHC human protein, HLA Class IA. The residues involved in binding antigen are shown by the closed circles (1).



Hughes and Nei (2) examined the numbers of synonymous and nonsynonymous substitutions in the MHC genes. They found that the number of nonsynonymous substitutions, which change the amino acid, was significantly greater than that of synonymous substitutions in the **antigen**-recognition sites, even though the opposite occurred elsewhere. They concluded that positive selection, most probably by overdominance or heterosis, is operating only on the antigen recognition sites.

#### Bibliography

1. T. Gojobori and T. Imanishi (1991) *Transplantation Now* 4, 26.
2. A. L. Hughes and M. Nei (1988) *Nature* 335, 167–170.

## Heterotrimeric G Proteins

Heterotrimeric G proteins are **GTP-binding** regulatory proteins that amplify, integrate, and transmit information from cell-surface receptors to cellular signaling proteins, known as effectors. G-protein signaling pathways are found in all eukaryotes; but they have been best studied in animals, where their number is greatest. Animals express over a thousand G-protein-coupled receptors that bind a variety of ligands: biogenic amines and lipids, **peptides**, **protein hormones**, odorants, and other compounds. Fungi use G proteins to convey signals from receptors for pheromones and nutrients, and higher plants may use G proteins in response to pathogens or light.

G proteins convey information by traversing a tightly controlled cycle of GTP binding and hydrolysis. G proteins are activated when GTP binds to the Ga subunit; but GTP binding is extremely slow, unless it is promoted by an agonist-bound receptor. GTP-activated G protein can then bind and activate the effector. Activation is terminated when the G protein hydrolyzes bound GTP to GDP, which does not activate (see [Gtpases](#)). Signal amplitude is thus a balance between receptor-stimulated GTP binding and GTP hydrolysis.

G proteins are defined by their GTP-binding subunits, which are primarily responsible for selective recognition of receptors and effectors (Table 1). Vertebrate genomes encode about 20 Ga subunits, while *Caenorhabditis elegans* has closer to 30. Plants and fungi have only one or two. Diversity is further enhanced in some cases by [alternative splicing](#) of RNA transcripts. Ga subunits form heterotrimers by reversible binding to Gbg dimers. Vertebrates express at least 5 Gb's and 12 Gg's. Although possible combinations of Gb with Gg and of Ga with Gbg dimer are limited by both their cellular expression and intrinsic affinity, the potential number of heterotrimers still numbers over 100. Gabg trimers are bound firmly to the inner face of the plasma membrane, with the exception of [transducin](#) in photoreceptor membranes. Membrane attachment depends on Gbg, although Ga subunits are somewhat **hydrophobic** and are further modified by two lipid groups (see [Membrane Anchors](#)).

**Table 1. Heterotrimeric G Proteins in Animals<sup>a</sup>**

| Family | Ga Genes        | Splicing Products         |
|--------|-----------------|---------------------------|
| s      | s               | 2 "long," 2 "short"       |
|        | olf             |                           |
| i      | i1              |                           |
|        | i2              |                           |
|        | i3              |                           |
|        | z               |                           |
|        | o               | A,B                       |
|        | t1 (transducin) |                           |
|        | t2 (transducin) |                           |
|        | g (gustducin)   |                           |
| q      | q               | (3 in <i>Drosophila</i> ) |
|        | 11              |                           |



|    |                                      |
|----|--------------------------------------|
|    | 14                                   |
|    | 15 / 16 (orthologous)                |
| 12 | 12                                   |
|    | 13                                   |
|    | con (concertina, <i>Drosophila</i> ) |

---

<sup>a</sup> This list is most representative of vertebrates, although flies and roundworms also have members of each family. In the roundworm *C. elegans*, for example, there are also 13 “worm-specific” Ga subunits that are expressed primarily in sensory neurons (57). Ga subunits have been identified in the coelenterates, but not yet in sponges or protozoa. Ga subunits in fungi and plants do not fall neatly into any of the animal classes.

The Gabg trimer is quite stable, with **dissociation constants** ( $K_d$ ) for release of Gbg estimated at  $10^{-9}$  M or less. GTP binding decreases this affinity and drives dissociation, to form free Ga-GTP plus free Gbg. Thus, activation of a Ga by receptor-promoted GTP binding also releases an active Gbg subunit, which regulates a distinct spectrum of effectors. Activation of Ga allows Gbg to generate a second signaling output that can be subsequently quenched by Ga-GDP.

## 1. Ga Subunits

Ga subunits are globular proteins with molecular weights in the range 39 to 44 kDa (1). They are based on two discrete structural **domains**. The GTP-binding domain, formed by the *N*-terminal and *C*-terminal thirds of Ga, has a [tertiary structure](#) and amino acid sequence similar to the small monomeric G proteins (see [GTP-Binding Proteins](#)), and it is therefore referred to as the ras-like domain, in analogy to p21<sup>ras</sup>. The middle third of the Ga sequence folds into a discrete **a-helix-rich** domain that acts in part as a cover to the GTP-binding site. Ras-like domains and helical domains can be expressed separately in bacteria and will associate with each other in solution (2). The linkage between the two domains is crucial. The more *C*-terminal linker peptide is also one of three “switch” regions that undergo profound conformational change upon activation; it is referred to as switch I. It resides near switch II, which is totally in the ras-like domain. The switch regions are involved in the binding of Gbg, a family of related proteins known as regulators of G-protein signaling, or [RGS proteins](#), and effectors (1). The second major interaction site involves the *N*- and *C*-termini of Ga proteins, which are close together in the folded structure. This region is involved in recognition of receptors and effectors, Gbg, and RGS proteins. In the inactive heterotrimer, the *N*-terminus of Ga lies as an a-helix along the side of the Gb subunit. The *N*-terminus is apparently quite flexible, however, because its conformation or orientation is apparently different in the abg trimer, free Ga-GDP, or activated free Ga or when bound to an RGS protein.

The *N*-terminal helix of Ga contains two important covalent lipid modifications. The terminal [amino group](#) of Ga subunits is blocked by lipid: in the Ga<sub>i</sub> family by a **myristoyl** amide group that is added shortly after translation; in Ga<sub>s</sub> by an unidentified but hydrophobic modification; and in the Ga<sub>q</sub> family by an unidentified modification. In addition, a [cysteine](#) residue close to the *N*-terminus is linked to a **palmitoyl group** in all Ga's except Ga<sub>t</sub>. In contrast to myristoylation, palmitoylation occurs through a thioester linkage and is therefore reversible. Palmitoyl groups turn over in cells and, in several cases, palmitate turnover is promoted by G-protein activation, so that depalmitoylation follows G-protein activation and repalmitoylation ensues thereafter. Myristoylation and the unknown modification of Ga<sub>s</sub> are important for establishing high-affinity interactions with effectors, Gbg, and RGS proteins. The regulatory roles of palmitoylation are less clear. Palmitoylation contributes to the

overall hydrophobicity of Ga's, and the absence of palmitate may explain in part the solubility of Ga<sub>t</sub> (3). Palmitoylation contributes to membrane attachment, but it is unclear whether this contribution reflects interaction with a specific anchoring protein (Gbg perhaps) or merely an increase in hydrophobicity. It also enhances the affinity for Gbg and decreases the affinity for RGS proteins (4).

Ga-GTP regulates effectors that include enzymes that synthesize or degrade [second messengers](#), **ion channels**, protein **kinases**, regulators of small G proteins, and transport proteins (Table 2). The list is growing constantly. Thus far, it appears that individual Ga subunits have relatively restricted ranges of effector targets. Thus, Ga<sub>s</sub> stimulates **adenylyl cyclase** (5), but other targets have not been identified. G<sub>t</sub> stimulates cyclic GMP phosphodiesterase (6). G<sub>q</sub> family members all stimulate **phospholipase C**'s (7), and there is one report of G<sub>q</sub> stimulation of Bruton's **tyrosine kinase** (8). Ga<sub>i</sub> is relatively unique in that it can inhibit two different effectors: several of the adenylyl cyclase **isoforms** (5) and, at least in some cases, some of the same K<sup>+</sup> channels that are stimulated by Gbg (9).

**Table 2. G-Protein-Regulated Effectors<sup>a</sup>**

| Effector                             | a      | bg | GAP      |
|--------------------------------------|--------|----|----------|
| Adenylyl cyclase                     | s, i↓  | ↓  | -        |
| Cyclic GMP phosphodiesterase (PDE) t |        |    | “co-GAP” |
| Phospholipase C                      | q      |    | +        |
| GIRK K <sup>+</sup> channel          | i↓     |    |          |
| PI-3-kinase                          |        |    |          |
| Ca <sup>2+</sup> channel             |        | ↓  |          |
| Rho GEF                              | 13(12) |    | +        |
| Protein kinases                      | q?     | ↓  |          |
| Transporters                         | ?      | ?  |          |

<sup>a</sup> The arrows indicate whether the activity is increased or decreased by the G protein. Question marks indicate uncertainty of direct interactions. GAP activity of effectors has not been assayed in all cases. Cyclic GMP PDE is shown as a co-GAP because it potentiates the GAP activity of RGS9. The list of known G-protein-regulated effectors is constantly growing.

## 2. Gbg Subunits

Gbg is a stable heterodimer of a ~35kDa Gb subunit and a ~8kDa Gg subunit, which fold and associate shortly after translation. Nascent Gb will denature if g is not present, and, once formed, Gbg can only be dissociated under denaturing conditions. In the native dimer, Gb is a torus of seven pie-wedge domains known as a b-propeller, because each wedge is composed almost exclusively of **b-strands**. The sequence of Gb is a sevenfold repeat of the “WD40” motif. However, each blade of the propeller is not one individual WD repeat, but is composed of roughly half of each of two adjacent repeats. The ~8kDa Gg subunit is extended across the face of the b-propeller, and the Gg N-terminus binds the N-terminus of Gb in an extended [coiled coil](#). Gbg binds to Ga along the face of the b-propeller opposite to Gg, with the central hole in the propeller aligned with the switch II region

of Ga. The other substantial Ga–Gbg contact is between the *N*-terminal  $\alpha$ -helix of Ga and the edge of Gb (10-12). Specificity in Gbg dimer formation is still poorly understood, and most Gbg dimers form readily. In mammals, Gb<sub>1–4</sub> are strikingly **homologous** to each other, and Gb<sub>5</sub>, the most divergent, is still identical in over 50% of its residues. Gg's are much more diverse.

Although Gbg subunits were initially considered to be nonspecific inhibitors of G-protein activation, it is now clear that they regulate as many or more effectors as do the Ga's (Table 2) and fulfill many other signaling functions as well (12). Gbg can either stimulate or inhibit different isoforms of adenylyl cyclase, activate PIP<sub>2</sub>-specific phospholipase C- $\beta$ , open K<sup>+</sup> channels and close Ca<sup>2+</sup> channels, and either activate or inhibit both **tyrosine kinases** and **serine/threonine protein kinases**. New Gbg-mediated functions are being discovered regularly. Gbg regulates effectors when it is released from the Gabg trimer upon activation, reflecting the preferential binding of Gbg to the inactive, GDP-bound Ga. Full dissociation of Ga-GTP from Gb may actually be an extreme outcome of a conformational change that alters the Ga–Gbg interface to allow Gbg to regulate effectors while remaining in physical contact with Ga.

In addition to controlling effector proteins, Gbg exerts multiple regulatory effects directly on Ga's. First, Gbg is required for stable membrane attachment of Ga. Gbg binds tightly to bilayers, in part because of *C*-terminal **prenylation** of the g subunit (see text below). Sprang proposed a plausible orientation of Gabg with respect to the bilayer, so that both the prenylated *C*-terminus of Gg and the lipid-modified *N*-terminus of Ga can interact directly with lipids (10). Next, Gbg is essentially required for regulation of Ga by receptors, in part by anchoring and orienting Ga with respect to the membrane surface. Anchoring Ga at the membrane surface increases its local concentration at the receptor several thousandfold, increasing its apparent affinity for the receptor proportionately. Direct effects of Gbg on receptor–Ga interactions are also likely. Last, Gbg binds preferentially to Ga in its inactive, GDP-bound conformation, thereby slowing dissociation of GDP from Ga. Gbg is thus a GDI (GDP-dissociation inhibitor) for its Ga partner, suppressing the basal signal output when there is no stimulatory input from the receptor (see text below). Gbg generally decreases the rate of dissociation of GDP from Ga about 10-fold, and even greater effects may occur in cells. Through a combination of these mechanisms, Gbg can facilitate the appropriate activation of a Ga by receptor and suppress its basal activation, thereby decreasing the background signal while permitting a rapid and efficient response to agonist.

Gbg contributes to agonist-dependent desensitization of receptor signaling by recruiting members of a family of protein kinases that phosphorylate G-protein-coupled receptors to initiate their desensitization and down-regulation. G-protein-coupled receptor kinases (GRKs) selectively **phosphorylate** agonist-liganded receptors at one or more **serine** or **threonine** residues in the *C*-terminal region or in the third cytoplasmic loop. Phosphorylation recruits scaffolding proteins, known as *arrestins*, which both inhibit interactions with Ga and link receptors to nascent clathrin-coated pits to initiate their endocytosis (13, 14).

The selectivity of Gbg among receptors, Ga subunits, and effectors—and the structural basis thereof—remains unclear. It has not been possible to isolate abg trimers with a unique subunit composition except from specialized cells. One exception, Gb<sub>1</sub>g<sub>1</sub>, is found exclusively in photoreceptors and interacts well only with Ga<sub>s</sub>, but such selectivity is unusual. *In vitro* studies have usually failed to observe marked selectivity of Gbg dimers in any context. The results of cellular studies have been mixed. Most notably, Kleuss et al. (15, 16) found that suppression of individual Gb or Gg molecules by microinjection of **antisense** messenger RNA leads to receptor-selective inhibition of Gbg signaling functions.

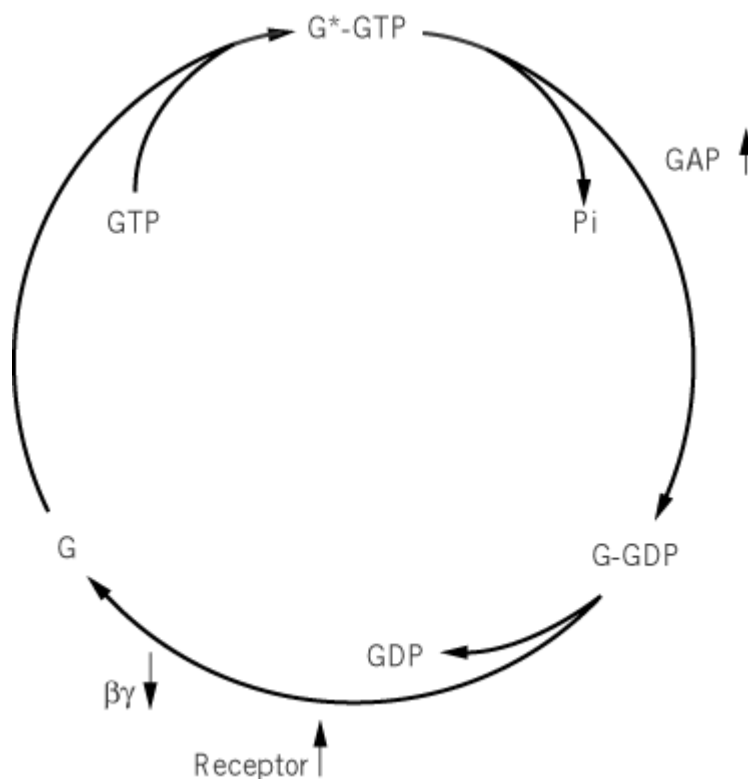
In addition to its controlled release from Ga, the availability and activity of Gbg is also controlled independently. Phosducin, a ubiquitously expressed protein that is unrelated to Ga, binds Gbg with high affinity and blocks its interactions with other proteins. Phosducin is thought to act as a Gbg

chelator, and its ability to bind Gbg is regulated by protein kinase A-catalyzed phosphorylation (17). Phosphorylation of phosphoducin decreases its affinity for Gbg, thereby regulating Gbg in response to cellular cyclic AMP. Gbg is also known to bind both actin and calmodulin-Ca<sup>2+</sup>, but it is unclear whether these interactions actually convey regulatory information to or from Gbg.

### 3. Regulatory GTPase Cycle

The regulatory GTPase cycle, proposed by Cassel and Selinger in 1977 (18, 19), provides the biochemical mechanism of G-protein signaling. Its kinetic parameters determine the rates of signal initiation and termination and the signal amplitude, both basal and agonist-stimulated (Fig. 1). Each step is tightly controlled. A G protein traverses the GTPase cycle relatively quickly with respect to the duration of the signal that it mediates, undergoing repeated activation, deactivation, and nucleotide exchange. The signal amplitude is therefore proportional to the fraction of time that the G protein spends in the GTP-bound, active state. The amplitude can be increased by accelerating GDP release and GTP binding or by inhibiting GTP hydrolysis. Similarly, the amplitude is decreased by inhibiting nucleotide exchange or by stimulating GTP hydrolysis. The rate of steady-state GTP hydrolysis and the fractional activation of the G protein at steady state are therefore independent: A G protein can cycle quickly or slowly, while maintaining either a high or low fractional regulatory activity.

**Figure 1.** The regulatory GTPase cycle. GTP-binding proteins are converted by binding GTP to the active state G\*-GTP such that they can in turn bind and activate effector proteins. Activation is terminated by hydrolysis of bound GTP to GDP, usually a slow process that is accelerated by GTPase-activating proteins, GAPs. The inactive G-GDP complex is relatively stable unless nucleotide exchange is catalyzed by a receptor. The inactive, GDP-bound state is stabilized by Gbg subunits.

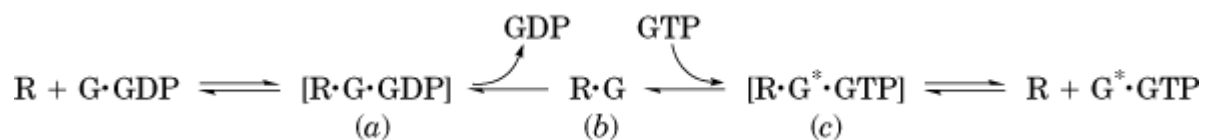


#### 3.1. GDP-GTP Exchange

The release of bound GDP and the binding of GTP is a slow process in which the initial release of GDP is rate-limiting under physiological conditions. *In vitro*, dissociation of GDP can take minutes

to hours, and a G protein retains a molecule of bound GDP after purification. Once the GDP has dissociated, GTP binding is fast under most experimental conditions. Association rate constants are estimated at  $10^6$  to  $10^8 \text{M}^{-1} \cdot \text{min}^{-1}$ , so the binding of GTP at a  $50 \mu \text{M}$  cytosolic concentration takes place in a second or less. Activation of the G protein subsequent to binding of GTP is also fast, although a GTP-bound but nonactivated  $G_\alpha$  subunit can be prepared at low temperatures or in the absence of divalent cations. Unliganded G proteins can be prepared *in vitro*, although harsh conditions may be required to strip the GDP (20).

G-protein-coupled receptors initiate signaling by accelerating GDP/GTP exchange. Agonist-liganded receptor initially binds to G-GDP to form a transient low-affinity complex of receptor, G protein, and GDP, from which the GDP dissociates [(a) in equation 1]. In the absence of added guanine nucleotide, the resulting receptor–G-protein complex [(b) in equation 1] is itself relatively stable; this maintains the nucleotide binding site in an “open” conformation. In the presence of GTP, a second transient complex, receptor–G–GTP [(c) in equation 1] forms and activates the G protein. At this point, the receptor can dissociate (to leave free, activated G protein) and is itself free to catalyze further G-protein activation. In cases where receptor-catalyzed GDP/GTP exchange occurs much more rapidly than G-protein deactivation, a single agonist-bound receptor can maintain the activation of multiple G-protein molecules and consequently generate considerable signal amplification. In specialized cells, such as photoreceptors, several thousandfold amplification can be achieved in this way (21).



Formally, receptor-catalyzed GDP/GTP exchange reflects the **negative cooperativity** of binding of receptors and guanine nucleotides to G proteins: Each decreases the G protein's affinity for the other, and the decreased affinity appears as an increased rate of dissociation (22). According to this mechanism, a receptor can accelerate the exchange of any guanine nucleotides, and it is the high cellular ratio of GTP to GDP that drives exchange in the physiological direction.

GDP/GTP exchange is also modulated by Gbg, which exerts two separate effects in the GTPase cycle. First, Gbg is a GDP-dissociation inhibitor (GDI) for  $G_\alpha$  subunits. By stabilizing the binding of GDP to  $G_\alpha$ , Gbg decreases the frequency of spontaneous GDP/GTP exchange and thus suppresses background noise in G-protein signaling pathways. Gbg can also inhibit receptor-stimulated GDP/GTP exchange and thereby inhibit G-protein signaling. Release of Gbg when one G protein is activated can inhibit the activation of another, with the extent of inhibition determined by the relative affinities of the two  $G_\alpha$ 's for Gbg. Second, and in seeming contrast to its inhibitory effects, Gbg is essentially required for the activation of G proteins by receptors. As discussed above, this effect reflects the membrane anchoring role of Gbg, but it may also have a direct regulatory component.

### 3.2. GTP Hydrolysis

$G_\alpha$  subunits catalyze hydrolysis of bound GTP by facilitating direct attack by [water](#) on the phosphoryl group, and GTP hydrolysis is followed quickly by dissociation of inorganic orthophosphate (1). The intrinsic rates of hydrolysis of bound GTP vary markedly, however. For  $G_s$ ,  $G_i$ , and  $G_o$ , GTP hydrolysis proceeds with a rate constant of  $3$  to  $4 \text{min}^{-1}$ , such that the active state has a half-life of 10 to 15 s.  $G_q$  is somewhat slower, with a half-life of 1 min.  $G_z$  is slowest, with a half-life of about 7 min. These rates are still much faster than those of the small, monomeric G proteins, evidently because the helical domain of heterotrimeric G proteins provides a catalytically important [arginine](#) residue that the small GTP binding proteins lacks.

If GTP hydrolysis is the primary deactivating event in G protein signaling, the rate of hydrolysis of bound GTP observed *in vitro* should be commensurate with the rate of signal termination upon removal of receptor agonist. This is true in some systems, such as adenylyl cyclase. In others, however, physiological termination of a signal is much faster than the rate at which G-protein-bound GTP is hydrolyzed by the appropriate G protein *in vitro* (23). This discrepancy was resolved by the discovery of GTPase-activating proteins, GAPs, which accelerate the hydrolysis of Ga-bound GTP up to several hundredfold (24).

The first GAP for a heterotrimeric G protein to be identified was also a G-protein-regulated effector, and two such effectors are now known to have GAP activity. Phospholipase C- $\beta$ 's both are regulated by  $G_q$  and act as a  $G_q$  GAP. PLC- $\beta$ 1 binds  $G_{\alpha_q}$ -GTP with a  $K_d$  of  $\sim 2nM$  and increases its rate of hydrolysis more than 50-fold (24, 25). A  $G_{13}$ -regulated effector, a GEF for the monomeric G protein rho, is also a  $G_{13}$  GAP (26). The second group of G-protein GAPs comprises the RGS proteins. About 30 mammalian RGS proteins have been identified, and RGS proteins are found throughout the eukaryotes. They are GAPs for members of the  $G_i$  and  $G_q$  families, but with a specificity that is otherwise still poorly understood. RGS proteins have not been demonstrated to have effector activity, although definable functional domains of many RGS proteins outside the conserved GAP domain suggests that they may have other functions. The binding interface between RGS proteins and the switch II region of Ga subunits is similar in overall structure to that between GAPs for the small, monomeric GTP-binding proteins and their targets. However, GAPs for the small G proteins contribute the necessary arginine residue that is provided by the helical domain in Ga subunits (1).

The physiological functions of G-protein GAPs are still largely speculative. GAPs can obviously inhibit G proteins by decreasing the lifetime of the active state, and one role of inducible RGS proteins is attenuation of G-protein signals. The yeast RGS protein Sst2p, which is induced by the **mating** factor signaling pathway, is clearly a feedback inhibitor that is induced by mating factor and serves to turn off the mating response after its induction. GAPs are also clearly required for fast responses to the removal of agonists, and they contribute indirectly to the speed of the stimulatory phase of the GTPase cycle also (24, 27, 28). Such kinetic effects are more likely functions for the GAP activities of effector proteins.

#### 4. Mechanism-Based Tools for G-Protein Research

Knowledge of the regulatory GTPase cycle has led to the development and use of highly specific tools for modulation of G-protein function in cells and *in vitro*. These tools include guanine-nucleotide analogs, bacterial [toxins](#), and a wide array of mutations that can interrupt or alter specific reactions in the GTPase cycle.

##### 4.1. Guanine Nucleotides

Nonmetabolizable analogs of GTP and GDP are used both as manipulators and as probes of G-protein signaling. Prototypes are GTP $\gamma$ S and Gpp(NH)p, GTP analogs that bind and activate G proteins but do not undergo chemical hydrolysis. Because they bind with relatively high affinity (particularly GTP $\gamma$ S), they lock G proteins in the active state. In cells and in membranes, they activate slowly in the absence of agonists, but more quickly in response to receptor-mediated nucleotide exchange. [ $^{35}$ S]GTP $\gamma$ S is also the standard radioligand used for quantifying G proteins. The Ga-[ $^{35}$ S]GTP $\gamma$ S complex is efficiently trapped on [nitrocellulose](#) filters, providing a convenient assay method. GDP $\gamma$ S, which mimics GDP and also binds tightly to Ga subunits, blocks activation and acts as a signaling antagonist. All three compounds are quite useful for intervention in G-protein pathways, but they have little, if any, selectivity among different G proteins, from the small monomeric G proteins through to the heterotrimeric G proteins. Greater specificity can be achieved by using a specific receptor agonist to accelerate nucleotide exchange, so that in a short time interval

only the G-protein target of that receptor will bind the nucleotide. An alternative to nonhydrolyzable GTP analogs, a combination of  $\text{Al}^{3+}$  and  $\text{F}^-$ , activates heterotrimeric G proteins directly by binding as the  $\text{AlF}_4^-$  ion to the Ga-GDP complex in place of the  $\gamma$ -phosphoryl group of GTP. Al/F thus converts G-GDP to the active form. While activation by Gpp(NH)p or GTP $\gamma$ S requires nucleotide exchange, Al/F activates almost immediately. Moreover, Al/F is selective for heterotrimeric G proteins and does not bind to monomeric G proteins under normal conditions. **X-ray crystallographic** studies showed that  $\text{GDP/AlF}_4^-$  is not a true analog of GTP, but instead is an analog of the octagonal bipyramidal [transition state](#) that occurs during GTP hydrolysis (1). The reason that small monomeric G proteins do not bind Al/F is that they are such poor GTPases that they have low affinity for the transition state structure, although they can bind Al/F if a GAP is present to stabilize the transition state (29).

#### 4.2. Toxins

**Cholera toxin** and [pertussis toxin](#) are widely used in G protein research because they are highly specific probes for  $G_s$  and for  $G_i$  or  $G_o$ , respectively. They **ADP-ribosylate** the  $\alpha$  subunits of these G proteins with distinct chemistry, specificity, and effect, and they can be used both *in vitro* and in cells. Cholera toxin, after entering the cell through the endocytic pathway, catalyzes the ADP-ribosylation of an arginine residue at the GTP-binding site of  $G_{\alpha_s}$ , thereby blocking the hydrolysis of bound GTP.  $G_s$  that has been ADP-ribosylated by cholera toxin is thus constitutively activated by cellular GTP. Cholera toxin is a specific activator of  $G_s$  signaling pathways because the cognate arginine residue in other Ga subunits is ADP-ribosylated much less efficiently. Pertussis toxin specifically ADP-ribosylates a conserved cysteine residue near the C-terminus of  $G_{\alpha_i}$ ,  $G_{\alpha_o}$ , and  $G_{\alpha_t}$ , thereby rendering the modified Ga insensitive to receptor. Thus, while the Ga retains its intrinsic regulatory and GTPase activities, it cannot be activated, and its function in cells is consequently blocked. Pertussis toxin is most useful for defining  $G_i$ -mediated pathways in cells, because the cysteine target residue is absent in the  $G_s$ ,  $G_q$ , and  $G_{12}$  families. Pertussis-toxin treatment of cultured cells can also be used to create a functionally  $G_i$ -free background in which to study mutant  $G_i$  family members whose target cysteine residue has been replaced.

#### 4.3. Mechanism-Based Mutants in G Proteins

Knowledge of the structure and genetics of the Ga family has yielded a group of reasonably reliable mutations that can be used to block GTP hydrolysis, to compromise GTP binding, or to block interaction with effectors or GEFs (1). Constitutively activating mutations that interfere with GTP hydrolysis are the oldest member of this class; the first is based on a naturally occurring, activating mutation of Gly12 in p21<sup>ras</sup> that makes it oncogenic. The corresponding mutation in Ga subunits also activates, but usually to a lesser extent. In Ga subunits, mutation of a conserved glutamine residue to leucine or of a conserved arginine (the cholera toxin modification site) to cysteine will produce constitutively active proteins that are expressed at normal levels and are otherwise unaltered. The leucine mutant is also insensitive to the GAP activity of RGS proteins (30). Apparent constitutive activation in cells can also be achieved by mutations that increase the rate of GDP/GTP exchange by destabilizing GDP binding (31). Strong dominant negative mutations in Ga subunits have been based on the mutation of Val17 in p21<sup>ras</sup>. This mutation, which presumably acts by creating a nonactivating Ga that chelates receptors, has been less widely used.

Mutations can alter the specificity with which a Ga recognizes specific receptors or effectors. Experience with chimeric mutations that replace either the N- or C-terminal few amino acids of a Ga with those of another have permitted the construction of G proteins that can respond to receptors of interest and regulate effectors of interest in transfected [cell lines](#). While the construction of these chimeras is not yet standard, the regions of interest for a specific interactions are reasonably well-defined, and enough mutants have been created to guide investigators in new cases.

## 5. G-Protein-Regulated Effectors

Heterotrimeric G proteins stimulate or inhibit a wide variety of effector proteins: enzymes that synthesize or degrade second messengers, protein kinases, ion channels, exchange factors for monomeric G proteins and, probably, a fairly large number of transporters (Table 2). As suggested by this list, the structures of effector proteins vary widely: They are either integral or peripheral membrane proteins that can be monomeric, homo-oligomeric, or hetero-oligomeric.

The mechanisms of effector regulation by G proteins also vary widely. Regulation can be stimulatory or inhibitory, additive, or synergistic, and both  $\alpha$  and  $\beta\gamma$  subunits independently regulate different spectra of effectors. Moreover, multiple isoforms of individual effectors often display differential sensitivity to regulation by  $G_{\alpha}$  and/or  $G_{\beta\gamma}$  subunits or by other signaling inputs. Individual isoforms are differentially expressed according to cell type, developmental stage, or acute regulatory signals. Such diversity allows a cell to respond to a single G-protein input with unique second messenger outputs, depending on the array of effector isoforms that it expresses. Last, an effector may also act as a GAP to regulate the GTPase activity of the  $G_{\alpha}$  subunit by which it itself is regulated. A receptor can thus modulate the amplitude and speed of the inputs it receives.

### 5.1. Adenylyl Cyclase

[Adenylate cyclase](#) catalyzes the formation of [cyclic AMP](#), the prototypical and ubiquitous [second messenger](#). Eukaryotic adenylyl cyclase is itself the prototypical G-protein-regulated effector. Its stimulation by the concerted action of hormones and guanine nucleotides led to the discovery of G proteins (32, 33). At least 10 isoforms of adenylyl cyclase are found in vertebrates, all of which are stimulated by the activated form of  $G_{\alpha_s}$ . Individual isoforms may be inhibited by the activated form of  $G_{\alpha_i}$ , and they may be either stimulated or inhibited by (or insensitive to)  $G_{\beta\gamma}$ . Cell-specific expression of these cyclase isoforms allows the construction of flexible signaling networks in which receptor-initiated signals to different G proteins can combine additively, synergistically, or negatively, to produce an integrated cyclic AMP signal. In addition to direct G-protein regulation, adenylyl cyclase isoforms can be further regulated by [calmodulin](#), by  $Ca^{2+}$  itself, or by phosphorylation by several protein kinases (34). In yeast, adenylyl cyclase is regulated by a monomeric, ras-related G protein in response to nutritional signals (35).

Adenylyl cyclase in animals is a monomeric, integral membrane protein whose architecture is based on a pseudodimeric structure (36, 37). Each half of the adenylyl cyclase molecule is a cluster of six membrane-spanning  $\alpha$ -helices, followed by a cytoplasmic domain. The [active site](#) and the site of all known regulatory interactions are contained within the pseudodimer of the two cytoplasmic domains, C1 and C2. Both  $G_{\alpha_i}$ ,  $G_{\alpha_s}$ , and  $G_{\beta\gamma}$  bind to C1:C2 dimers, expressed either separately or when fused. Determinants of sensitivity to  $G_{\alpha_s}$ ,  $G_{\alpha_i}$ , and  $G_{\beta\gamma}$  are found in distinct and nonoverlapping binding sites in the C1:C2 pseudodimer. The function of the membrane-spanning domain is unknown. Its apparent similarity to membrane transport proteins suggests that it may have an evolutionary history as a transporter. The possibility that an adenylyl cyclase from paramecium is an ion channel (38) is at least consistent with this idea. It is of evolutionary interest that adenylyl cyclase in yeast, which is stimulated by a ras-like monomeric G protein, is a large (~200kDa) soluble protein that is only vaguely related to animal cyclases in the catalytic domains (39).

### 5.2. Cyclic GMP Phosphodiesterase

A distinct cyclic GMP phosphodiesterase (PDE) is found in vertebrate photoreceptor cells (rods and cones), where it degrades [cyclic GMP](#) in response to light signals to cause the closure of cyclic GMP-gated cation channels (6). Retinal PDE is stimulated by the active form of  $G_t$  ([transducin](#)), which in turn is activated by the photobleached (meta II) form of [rhodopsin](#), a receptor whose ligand is the covalently bound, light-sensing chromophore retinal.  $G_t$ -sensitive PDE is a heterotetramer of an  $\alpha\beta$  catalytic dimer and two inhibitory  $\gamma$  subunits, to which activated  $G_t$  binds. In contrast to



adenylyl cyclase, which is activated by the binding of  $G_{\alpha_s}$ , the PDE is stimulated by the dissociation of the  $\beta$  subunits when they are sequestered by  $G_{\alpha_t}$ -GTP. The rhodopsin- $G_{\alpha_t}$ -phosphodiesterase pathway is also unusual in that its components are expressed at extraordinarily high levels in photoreceptor membranes, where rhodopsin represents about 70% the total membrane protein and  $G_{\alpha_t}$  is about half of the remainder (40).

### 5.3. Phospholipase C

PIP<sub>2</sub>-specific **phospholipase C-b** (PLC-b) isozymes are peripheral membrane proteins that are stimulated either by  $G_{\alpha_q}$  family members, by Gbg, or by both, depending on the PLC-b isoform (7). They hydrolyze PIP<sub>2</sub> to the second messengers **diacylglycerol**, an activator of protein kinase C, and inositol-1,4,5-trisphosphate, which releases stored  $Ca^{2+}$  into the cytoplasm (see **Inositol triphosphates**). PLC-bs, approximately 150 kDa in size, are based on a core catalytic domain that is common to all PIP<sub>2</sub>-selective PLCs. In addition, the PLC-bs also contain a long, C-terminal extension that binds  $G_{\alpha_q}$  and mediates its stimulatory interaction (41). Gbg subunits bind at sites in the N-terminal half of the molecule. Differential stimulation of PLC-b activity by  $G_{\alpha}$  or Gbg depends on the isoform. Both *in vitro* and *in vivo*, PLC-b-1 is primarily  $G_{\alpha_q}$ -sensitive; PLC-b-2 and b-3 are also markedly stimulated by Gbg; regulation of PLC-b4 has not been described in detail, but it does respond to  $G_{\alpha_q}$ . In cells, stimulation of  $G_{\alpha_i}$  also can lead to PLC-b stimulation, but the response is mediated by Gbg. Sensitivity of PLC activation to pertussis toxin is therefore indicative of Gbg signaling. PLC-bs were the first G-protein-coupled effectors to be shown to have intrinsic GAP activity, which is specific for  $G_{\alpha_q}$ , its activator (25). Such intrinsic GAP activity in an effector is presumably involved in rapid temporal control of signaling and in determination of receptor specificity, rather than in autoinhibition (24). The  $G_{\alpha_q}$  PLC-b pathway, while ubiquitous in animals, is particularly important for vision in mollusks and arthropods, where, in photoreceptor cells, rhodopsin regulates a  $G_{\alpha_q}$  that in turn stimulates a PLC-b.

### 5.4. Inward Rectifier K<sup>+</sup> Channels

G-protein-gated, inward rectifier K<sup>+</sup> channels (GIRKs), which are found in heart, neurons, and other cells, cause a cellular hyperpolarization that is particularly important in the control of neurons, cardiac muscle, and smooth muscle. GIRKs are stimulated by Gbg subunits that in most cases are released by stimulation of  $G_{\alpha_i}$ , so that GIRK activation is usually sensitive to pertussis toxin (12). A  $G_{\alpha_i}$ -regulated K<sup>+</sup> channel is either a homotetramer or heterotetramer of members of the GIRK protomer family. Free Gbg binds directly to the GIRK channel, and addition of GDP-bound (inactive)  $G_{\alpha}$  can inhibit GIRK stimulation by binding Gbg and blocking its action. It has been reported that activated  $G_{\alpha_i}$  can inhibit GIRK-mediated currents, although it is not clear how this inhibition relates to the ability of nonactivated  $G_{\alpha_i}$  to chelate Gbg and thus terminate its activating effect (9).

### 5.5. Other Ion Channels

Voltage-gated  $Ca^{2+}$  channels of the N and P/Q (dihydropyridine-insensitive) type are inhibited by Gbg. These are hetero-oligomeric channels that have a complex [quaternary structure](#) of integral and peripheral protein components. The location of the binding sites for Gbg are not known, but are presumed to lie in the  $\beta$  subunit, whose different isoforms determine the sensitivity to Gbg (42, 43). Some voltage-gated Na<sup>+</sup> channels are also stimulated by Gbg, perhaps by a direct binding mechanism (44).

### 5.6. Protein Kinases

Direct regulation of protein kinases by G proteins was first demonstrated in the yeast *Saccharomyces*

*cerevisiae*, where Gbg (Ste4p, Ste18p) released from the Ga subunit (Gpa1p) in response to a pheromone receptor activates an ERK signaling cascade to control the mating response (45). An increasing number of protein kinases are now known to be regulated by G proteins in all eukaryotes, and these signaling pathways are under intense study.

At least three groups of protein kinases are regulated by G proteins, and their activation can lead to control of cell growth, cell differentiation, or acute physiological responses. First, Gbg both stimulates and inhibits different serine/threonine kinases in several MAP/ERK kinase cascades in animals and fungi, and probably in plants too. In the case of Raf1, Gbg probably stimulates the kinase activity directly, but other protein kinases may simply be recruited to their sites of action by Gbg. Second, the G-protein-coupled receptor kinases that initiate desensitization and endocytosis of G-protein-coupled receptors also bind directly to Gbg. Binding to Gbg localizes these kinases to the plasma membrane and facilitates the phosphorylation of receptors, although it is not known whether Gbg actively increases protein kinase catalytic activity (14). In animal cells, activation of G-protein pathways also leads to increased phosphorylation of protein tyrosine residues by a combination of indirect mechanisms and, probably, direct stimulation of the kinases. Both  $G_{\alpha_q}$  and Gbg have been reported to stimulate the activity of Bruton tyrosine kinase (8, 46).

### 5.7. Phosphatidylinositide-3-kinase

Phosphatidylinositide-3-kinases (PI-3-Ks) catalyze the phosphorylation of PI-4,5-P<sub>2</sub> to PI-3,4,5-P<sub>3</sub>, which activates at least one protein kinase and may trigger other cellular responses. PI-3-Ks are heterodimers of a catalytic subunit, which is a member of a conserved family of PI-3-K isozymes and one of a diverse group of regulatory subunits (47). The G-protein-regulated PI-3-K has unique regulatory and inhibitory subunits not closely related to other members of the family. It is stimulated directly by Gbg, and both subunits are required for regulation by Gbg.

### 5.8. Rho GEF

A [guanine nucleotide exchange factor](#) (GEF) for the Rho class of small monomeric G proteins is both stimulated by  $G_{13}$  and acts as a GAP for  $G_{12}$  and  $G_{13}$  (26, 48).

### 5.9. Transporters

Transporters of  $Na^+/H^+$ , glucose,  $Mg^{2+}$ , and  $Ca^{2+}$  have all been reported to be regulated, positively or negatively, by G-protein signaling pathways (49-52). In some cases, their behaviors are consistent with direct regulation by either Ga or Gbg subunits.

## 6. G-Protein-Coupled Receptors

G-protein-coupled receptors (GPCRs) bind G proteins at the inner face of the plasma membrane and, in response to extracellular agonists, catalyze the exchange of GDP for GTP. GPCRs make up a huge [multigene family](#) of proteins whose members are found throughout the eukaryotes. Genes for GPCRs in humans number about 1000, exclusive of the multiple odorant receptors expressed in olfactory epithelium, and account for over 1% of the total [genome](#). Appropriately, GPCRs respond to a diverse range of agonists: biogenic amines, lipids, peptides, chemically diverse pheromones, nucleotides and nucleosides, protein hormones,  $Ca^{2+}$ , odorants and tastants, photons (with the help of an attached 11-*cis*-retinal chromophore), and [proteinases](#) (by exposure of endogenous agonist peptides) (53, 54).

The GPCRs share a common overall structure and mechanism of action. They are based on a common core structure, a bundle of seven hydrophobic, bilayer-spanning helices in the plasma membrane. The *N*-terminus, usually glycosylated, is extracellular, and the *C*-terminus is cytoplasmic. Both are of variable length. Interhelix loops are short, with the exception of the third cytoplasmic loop, which links  $\alpha$ -helices 5 and 6: It can contain up to about 200 amino acid residues. Thus, while the membrane-spanning domain of GPCRs has a minimal size of less than 30 kDa, receptors range in total size to more than 100 kDa. The high-resolution [protein structure](#) of a GPCR is not yet known, but a low-resolution structure is available for visual rhodopsin (55). This structure orients the  $\alpha$ -

helices with respect to each other; and it distinguishes their packing from that of [bacteriorhodopsin](#), an unrelated seven-helix **proton pump**. The rhodopsin structure has been combined with biochemical data, comparison of similar receptors, and the results of numerous mutagenesis experiments, to produce a useful model of a receptor's hydrophobic core ([56](#)).

The amino acid sequences of GPCRs in animals are most highly conserved in the membrane-spanning helices and short loops. Several signature sequences are identifiable in these regions, and receptors closely related by sequence usually bind similar extracellular ligands ([54](#)). Animal GPCRs include three subfamilies. The major subfamily accounts for receptors for most ligands, including those for most small molecules, proteins, and many peptides. The other two families are the receptors for glutamate and for a group of [peptide hormones](#). Enough GPCR sequences are now available that it is frequently possible to predict the ligand specificity of a receptor from just its sequence. GPCRs in fungi and slime molds are also based on a bundle of seven membrane-spanning  $\alpha$ -helices, but any homology of these receptors to animal receptors or to each other is not obvious. Their evolutionary origins are uncertain. Higher plants probably also use seven-span receptors to regulate heterotrimeric G proteins, but identification of candidate proteins as legitimate GPCRs is not yet certain.

The GPCR family is extraordinarily important to drug design: About 30% of prescription drugs act on GPCRs. Both because of this impetus and because of basic interest in their function and biology, extraordinary efforts have gone toward detailed analysis of these receptors and their ligand-binding sites. Ligand-binding sites have been mapped chemically, by mutagenesis and by the study of responses to structurally defined ligands. Experimental approaches have been complemented by intense molecular modeling. Two general structures have emerged. In receptors for small ligands such as biogenic amines, the ligand-binding site is a relatively deep cleft on the extracellular face of the bundle of membrane-spanning helices. Ligands make contact with residues from multiple helices, and it is thought that these contacts initiate the levering of the helices that constitutes receptor activation. Receptors for other ligands, notably the  $\alpha\beta$  dimeric glycoprotein hormones, have extensive *N*-terminal domains that themselves have ligand-binding activity. These *N*-terminal ligand-binding domains may position the ligand near the external face of the helices to promote helix–helix levering.

In contrast to the external ligand-binding site, the cytoplasmic site on GPCRs that binds and regulates G proteins has been harder to identify. Mutagenesis experiments—mostly the creation of chimeric receptors that use segments of one receptor to alter the G protein specificity of another—suggest that the second and third intracellular loops are the primary components of the G-protein-binding site. Other mutations in the third loop can create constitutively active receptors. In those cases where the third intracellular loop is large, regions near spans 5 and 6 are more important for G-protein regulation, and the central segment fulfills other functions. It has not been possible to identify sequence motifs that determine specificity among G proteins, even when the general regions important for selectivity have been established. This divergence of sequence in receptors that bind the same G protein suggests that receptor subfamilies acquired specificity for extracellular ligands prior to the evolution of the multiple G-protein families.

The process of GPCR activation by ligand is best understood as a conformational equilibrium: An unliganded receptor is usually in an inactive conformation, but the binding of agonist favors the activated state. Partial agonists promote activation incompletely. In any case, the active receptor can catalyze nucleotide exchange on the correct  $G_\alpha$  subunit. In constitutively active mutant receptors, an unliganded receptor can catalyze nucleotide exchange. As predicted, these mutants generally bind agonist ligands more tightly than do the wild-type receptor, indicating that they are constitutively in the agonist-favored conformation. The coupling of receptors and G proteins leads to a stereotyped effect of guanine nucleotides on the affinity of agonist binding. As discussed above, binding of an active (agonist-liganded) receptor to G protein drives dissociation of guanine nucleotide to form an agonist–receptor–G-protein complex that can be dissociated by the addition of GTP (or GDP) (equation 1). The receptor–G-protein complex displays higher affinity for agonist than does free

receptor. Because addition of guanine nucleotide to membranes drives dissociation of this complex, to leave free receptor and nucleotide-bound G protein, nucleotides are often observed to decrease the affinity of receptors to agonist—but not antagonist—ligands. The nucleotide-driven decrease in agonist affinity has become a useful test for whether a receptor of unknown mechanism interacts with a G protein (22).

Cells regulate their level of response to G-protein stimulation in large part by modulating the expression and activity of the appropriate GPCRs. GPCRs undergo acute desensitization upon stimulation as a result of phosphorylation of serine and threonine residues on their cytoplasmic faces, usually in the third cytoplasmic loop or in the C-terminal region. Phosphorylation is catalyzed both by second messenger-activated protein kinases and by a group of GPCR kinases (GRKs). GRKs preferentially phosphorylate receptors in their agonist-bound forms. Phosphorylation limits receptor interaction with G proteins to varying extents. More importantly, phosphorylation helps recruit inhibitory proteins, known as arrestins, which both further inhibit receptor activity and promote receptor-endocytosis and consequent degradation by linking receptors to [clathrin](#), the major coat protein of coated pits (13, 14).

### Bibliography

1. S. R. Sprang (1997) *Annu. Rev. Biochem.* **66**, 639–678.
2. D. W. Markby, R. Onrust, and H. R. Bourne (1993) *Science* **262**, 1895–1901.
3. E. M. Ross (1995) *Curr. Biol.* **5**, 107–109.
4. Y. Tu, J. Wang, and E. M. Ross (1997) *Science* **278**, 1132–1135.
5. R. K. Sunahara, C. W. Dessauer, and A. G. Gilman (1996) *Annu. Rev. Pharmacol. Toxicol.* **36**, 461–480.
6. L. Stryer (1986) *Annu. Rev. Neurosci.* **9**, 87–119.
7. W. D. Singer, H. A. Brown, and P. C. Sternweis (1997) *Annu. Rev. Biochem.* **66**, 475–509.
8. K. Bence, W. Ma, T. Kozasa, and X.-Y. Huang (1997) *Nature* **389**, 296–299.
9. W. Schreibmayer, C. W. Dessauer, D. Vorobiov, A. G. Gilman, H. A. Lester, N. Davidson, and N. Dascal (1996) *Nature* **380**, 624–627.
10. M. A. Wall, D. E. Coleman, E. Lee, J. A. Iñiguez-Lluhi, B. A. Posner, A. G. Gilman, and S. R. Sprang (1995) *Cell* **83**, 1047–1058.
11. D. G. Lambright, J. Sondek, A. Bohm, N. P. Skiba, H. E. Hamm, and P. B. Sigler (1996) *Nature* **379**, 311–319.
12. D. E. Clapham and E. J. Neer (1997) *Annu. Rev. Pharmacol. Toxicol.* **37**, 167–203.
13. O. B. Goodman Jr., J. G. Krupnick, F. Santini, V. V. Gurevich, R. B. Penn, A. W. Gagnon, J. H. Keen, and J. L. Benovic (1996) *Nature* **383**, 447–450.
14. J. G. Krupnick and J. L. Benovic (1998) *Annu. Rev. Pharmacol. Toxicol.* **38**, 289–319.
15. C. Kleuss, J. Hescheler, C. Ewel, W. Rosenthal, G. Schultz, and B. Wittig (1991) *Nature* **353**, 43–48.
16. C. Kleuss, H. Scherübl, J. Hescheler, G. Schultz, and B. Wittig (1993) *Science* **259**, 832–834.
17. P. H. Bauer, S. Müller, M. Puzicha, S. Pippig, B. Obermaier, E. J. M. Helmreich, and M. J. Lohse (1992) *Nature* **358**, 73–76.
18. D. Cassel and Z. Selinger (1976) *Biochim. Biophys. Acta* **452**, 538–551.
19. D. Cassel, H. Levkovitz, and Z. Selinger (1977) *J. Cyclic Nucleotide Res.* **3**, 393–406.
20. K. M. Ferguson and T. Higashijima (1991) *Methods Enzymol.* **195**, 188–192.
21. P. A. Liebman, K. R. Parker, and E. A. Dratz (1987) *Annu. Rev. Physiol.* **49**, 765–791.
22. E. M. Ross (1992) In *An Introduction to Molecular Neurobiology* (Z. W. Hall, ed.), Sinauer Associates, Sunderland, Mass., pp. 181–206.
23. G. E. Breitwieser and G. Szabo (1988) *J. Gen. Physiol.* **91**, 469–493.

24. E. M. Ross (1995) *Rec. Prog. Hormone Res.* **50**, 207–221 [Abstract].
25. G. Berstein, J. L. Blank, D.-Y. Jhon, J. H. Exton, S. G. Rhee, and E. M. Ross (1992) *Cell* **70**, 411–418.
26. T. Kozasa, X. Jiang, M. J. Hart, P. M. Sternweis, W. D. Singer, A. G. Gilman, G. Bollag, and P. C. Sternweis (1998) *Science* **280**, 2109–2111.
27. G. H. Biddlecome, G. Berstein, and E. M. Ross (1996) *J. Biol. Chem.* **271**, 7999–8007.
28. N. Zerangue and L. Y. Jan (1998) *Curr. Biol.* **8**, 313–316.
29. R. Mittal, M. R. Ahmadian, R. S. Goody, and A. Wittinghofer (1996) *Science* **273**, 115–117.
30. D. M. Berman, T. M. Wilkie, and A. G. Gilman (1996) *Cell* **86**, 445–452.
31. T. Iiri, P. Herzmark, J. M. Nakamoto, C. Van Dop, and H. R. Bourne (1994) *Nature* **371**, 164–168.
32. E. M. Ross and A. G. Gilman (1977) *J. Biol. Chem.* **252**, 6966–6969.
33. E. M. Ross and A. G. Gilman (1980) *Annu. Rev. Biochem.* **49**, 533–564.
34. R. Taussig, W.-J. Tang, J. R. Hepler, and A. G. Gilman (1994) *J. Biol. Chem.* **269**, 6093–6100.
35. T. Toda, I. Uno, T. Ishikawa, S. Powers, T. Kataoka, D. Brock, S. Cameron, J. Broach, K. Matsumoto, and M. Wigler (1985) *Cell* **40**, 27–36.
36. J. J. G. Tesmer, R. K. Sunahara, A. G. Gilman, and S. R. Sprang (1997) *Science* **278**, 1907–1916.
37. J. Krupinski, F. Coussen, H. A. Bakalyar, W.-J. Tang, P. G. Feinstein, K. Orth, C. Slaughter, R. R. Reed, and A. G. Gilman (1989) *Science* **244**, 1558–1564.
38. J. E. Schultz, S. Klumpp, R. Benz, W. J. Schurhoff-Goeters, and A. Schmid (1992) *Science* **255**, 600–603.
39. S. Marcus, M. Wigler, H. P. Xu, R. Ballester, M. Kawamukai, and A. Polverino (1993) *Ciba Found. Symp.* **176**, 53–61.
40. H. E. Hamm and M. D. Bownds (1986) *Biochemistry* **25**, 4512–4523.
41. S. G. Rhee and K. D. Choi (1992) In *Advances in Second Messenger and Phosphoprotein Research* (J. W. Putney, Jr. ed.), Raven Press, New York, pp. 35–61.
42. S. Herlitze, D. E. Garcia, K. Mackie, B. Hille, T. Scheuer, and W. A. Catterall (1996) *Nature* **380**, 258–262.
43. S. R. Ikeda (1996) *Nature* **380**, 255–258.
44. J. Y. Ma, W. A. Catterall, and T. Scheuer (1997) *Neuron* **19**, 443–452.
45. J. Kurjan (1992) *Annu. Rev. Biochem.* **61**, 1097–1129.
46. S. Tsukada, M. I. Simon, O. N. Witte, and A. Katz (1994) *Proc. Natl. Acad. Sci. USA* **91**, 11256–11260.
47. L. R. Stephens, A. Eguinoa, H. Erdjument-Bromage, M. Lui, F. Cooke, J. Coadwell, A. S. Smrcka, M. Thelen, K. Cadwallader, P. Tempst, and P. T. Hawkins (1997) *Cell* **89**, 105–114.
48. M. J. Hart, X. Jiang, T. Kozasa, W. Roscoe, W. D. Singer, A. G. Gilman, P. C. Sternweis, and G. Bollag (1998) *Science* **280**, 2112–2114.
49. M. E. Maguire and J. J. Erdos (1980) *J. Biol. Chem.* **255**, 1030–1035.
50. R. Hooley, C.-Y. Yu, M. Symons, and D. L. Barber (1996) *J. Biol. Chem.* **271**, 6152–6158.
51. R. A. Hall, R. T. Premont, C.-W. Chow, J. T. Blitzer, J. A. Pitcher, A. Claing, R. H. Stoffel, L. S. Barak, S. Shenolikar, E. J. Weinman, S. Grinstein, and R. J. Lefkowitz (1998) *Nature* **392**, 626–630.
52. M. Kuroda, R. C. Honnor, S. W. Cushman, C. Londos, and I. A. Simpson (1987) *J. Biol. Chem.* **262**, 245–253.
53. S. Watson and S. Arkininstall (1994) *The G-Protein Linked Receptor Facts Book*, Academic Press, San Diego.
54. C. D. Strader, T. M. Fong, M. R. Tota, D. Underwood, and R. A. F. Dixon (1994) *Annu. Rev.*

Biochem. **63**, 101–132.

55. V. M. Unger and G. F. X. Schertler (1995) *Biophys. J.* **68**, 1776–1786.

56. J. M. Baldwin, G. F. X. Schertler, and V. M. Unger (1997) *J. Mol. Biol.* **272**, 144–164.

57. C. I. Bargmann and J. M. Kaplan (1998) *Annu. Rev. Neurosci.* **21**, 279–308.

### Suggestions for Further Reading

58. D. E. Clapham and E. J. Neer (1997) G protein subunits. *Annu. Rev. Pharmacol. Toxicol.* **37**, 167–203.

59. H. G. Dohlman, J. Thorner, M. G. Caron, and R. J. Lefkowitz (1991) Model systems for the study of seven-transmembrane-segment receptors. *Annu. Rev. Biochem.* **60**, 653–688.

60. A. G. Gilman (1987) G proteins: Transducers of receptor-generated signals. *Annu. Rev. Biochem.* **56**, 615–649.

61. T. Gudermann, F. Kalkbrenner, and G. Schultz (1996) Diversity and selectivity of receptor–G protein interaction. *Annu. Rev. Pharmacol. Toxicol.* **36**, 429–459.

62. Y. Kaziro, H. Itoh, T. Kozasa, M. Nakafuku, and T. Satoh (1991) Structure and function of signal-transducing GTP-binding proteins. *Annu. Rev. Biochem.* **60**, 349–400.

63. J. G. Krupnick and J. L. Benovic (1998) The role of receptor kinases and arrestins in G protein-coupled receptor regulation. *Annu. Rev. Pharmacol. Toxicol.* **38**, 289–319.

64. E. M. Ross (1989) Signal sorting and amplification through G protein-coupled receptors. *Neuron* **3**, 141–152.

65. S. R. Sprang (1997) G protein mechanisms: insights from structural analysis. *Annu. Rev. Biochem.* **66**, 639–678.

## Hfr'S And F-Primes

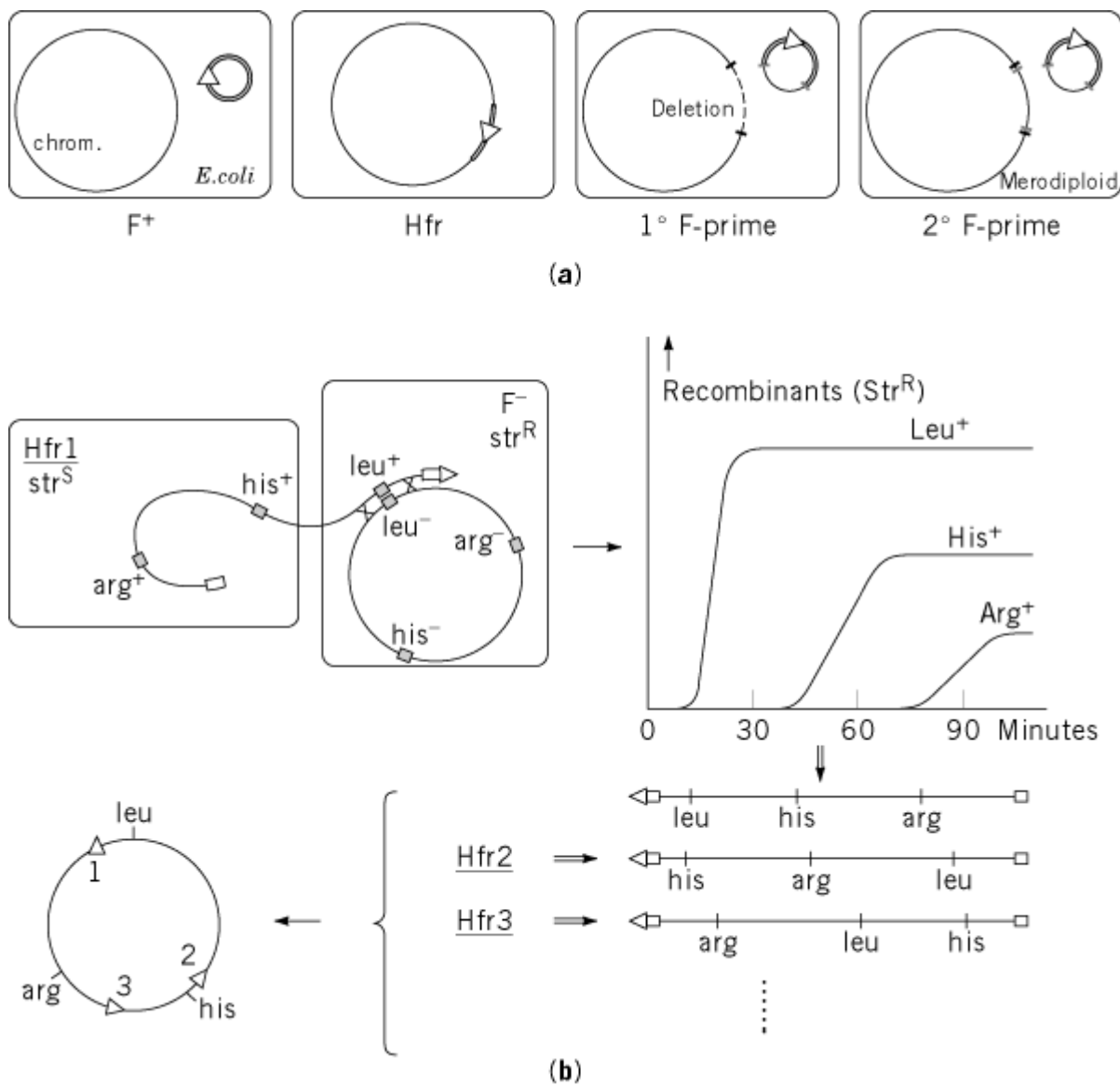
### 1. Bacterial Sex

Immediately after its discovery by Lederberg and Tatum in 1946 ([1](#)), bacterial mating was found to be unorthodox. Conventional sex mixes the entire **genomes** of both parents. But fertile strains of *Escherichia coli* K-12 donated only parts of their **chromosomes** during mating, transferring them into the other partner, the recipient, to recombine with its genome. The genetic determinant of this donor ability was also unusual: It was able to transfer itself rapidly from cell to cell as an autonomous unit without any associated chromosomal genes. This determinant was named the F (or sex) factor, later shown to be a **plasmid**, and the recombination process it enabled proved to be the key to our understanding of the genetic and physical organization of the bacterial chromosome (see [F Plasmid](#)).

### 2. Hfr's

Cultures of F<sup>+</sup> strains were able to transfer the F factor itself at high frequency to F<sup>-</sup> recipients, thereby converting them to F<sup>+</sup> donors, but chromosomal markers were transferred much less efficiently, generating recombinant progeny at frequencies of about one per million parental cells. By chance, certain isolates from an F<sup>+</sup> culture were found with the opposite character: They had virtually lost the ability to transfer donor ability but now generated recombinants at greatly elevated frequencies. They were consequently called Hfr (high-frequency recombination) strains (Fig. [1a](#)).

**Figure 1.** Bacterial genetics by Hfr's and F-primes. **(a)** Genetic donor types. Left to right: An F<sup>+</sup> cell containing the autonomous F plasmid; an Hfr, in which F has integrated into the chromosome; a primary F' strain, in which F has excised with a portion of chromosome; a secondary F' strain, diploid for the chromosomal portion carried by the F' as a result of transfer to a wild-type recipient (ie, merodiploid). The arrow indicates the origin (*oriT*) and orientation of F transfer. **(b)** Mapping by interrupted mating. Following donor-recipient contact, the Hfr chromosome breaks at *oriT* (open arrow) to initiate transfer of a single donor strand into the recipient; concomitant replication in the donor restores its chromosome (not shown). The entering strand recombines with the resident chromosome, replacing the mutant alleles (*leu*, *his*, *arg*) with wild-type counterparts, and recombinants are selected by plating on medium with [streptomycin](#) (to stop growth of donors) and without one of the [amino acids](#) (to stop growth of nonrecombinant recipients). The recombination frequency declines as distance from *oriT* increases, owing to random breakage, and recombinant formation after sampling is prevented by forced breakage, allowing measurement of distances between *oriT* and the gene. For a given Hfr, the map can be read directly from the time-of-entry data. The combined results from several matings (F<sup>-</sup>×Hfr1, Hfr2, Hfr3) showed that the *E. coli* chromosome is genetically circular.



### 3. F-Primes

Although Hfr strains were generally stable, F<sup>+</sup> “revertant” derivatives identical to the original strain were readily isolated. Some of these derivatives showed an intermediate **phenotype**: They transferred not only the F<sup>+</sup> character at high frequency, but a small set of chromosomal markers as

well, so that the normally haploid recipient became [merodiploid](#), that is, **diploid** for only a part of its genome. Such derivatives were termed F-primes (F'; Fig [1a](#)).

The physical basis of these early observations is now understood, even if some aspects of the transfer mechanism itself continue to present challenges (see **Conjugation**). The molecular basis of Hfr and plasmid-prime formation was worked out almost entirely using the F plasmid in *E. coli* K12, and we therefore focus our discussion mostly on F. Host gene mobilization involving other plasmids in other bacterial genera conform, in general, to the F pattern. Nevertheless, Hfr's and primes derived from other plasmids justify inclusion in this chapter on the grounds of their utility, and we briefly describe methods for isolating them.

#### 4. The Nature of Hfr Transfer

Wollman and Jacob ([2](#)) exploited the one-way transfer of DNA by developing the interrupted mating experiment to map genes (Fig. [1b](#)); in doing so, they revealed the nature of Hfr transfer. Mating is started by mixing cultures of Hfr and F<sup>-</sup> strains. Samples of the mixture are taken at intervals and mating pairs ruptured by violent shaking; portions are then plated on agar to select for growth of recombinants. Each Hfr strain introduced its markers into the F<sup>-</sup> recipient in a precise order and at a fixed time. Recombination frequencies were greatest for markers transferred early ( $\sim 10^{-1}$ ) and declined progressively for later entering markers (to  $\sim 10^{-5}$ ); this presumably reflects the sensitivity of conjugal pairs to spontaneous breakage. The results established the linkage relationships of various markers as a simple function of time of entry into the recipient. Moreover, collation of the results from matings with several Hfr strains indicated that the *E. coli* genome consists of a single, circular linkage group (Fig. [1b](#)) ([3](#)). This striking genetic result has been corroborated by various biophysical means and confirmed by [DNA sequencing](#) of the entire genome ([4](#)).

Hfr donors transferred donor ability at only very low frequencies and in association with late-entering markers. This observation, together with the time of entry results, suggested that Hfr strains arose by a single recombinational crossover between a circular F plasmid and a circular chromosome, such that the F was integrated into the chromosome (Fig. [1a](#)) ([5](#), [6](#)). During mating, the F is broken to provide the leading end, or origin of transfer (*oriT*), which passes into the recipient, while the genes that encode the mating functions are left at the other end, which is therefore the last, and least likely, to enter (Fig. [1b](#)). The order of entry for a given Hfr is determined simply by the site and orientation of F integration into the chromosome. The donor ability of F<sup>+</sup> cultures is then explained by the formation of many different Hfr's, each with F integrated at a distinct site, so that the F<sup>+</sup> population as a whole transfers all markers at about the same rate.

The recombination process is completed in the recipient, where a portion of the linear DNA transferred can be substituted for the equivalent segment of the recipient's circular chromosome by homologous recombination; if different alleles are present in the recombining region, the result is formation of a genetic recombinant.

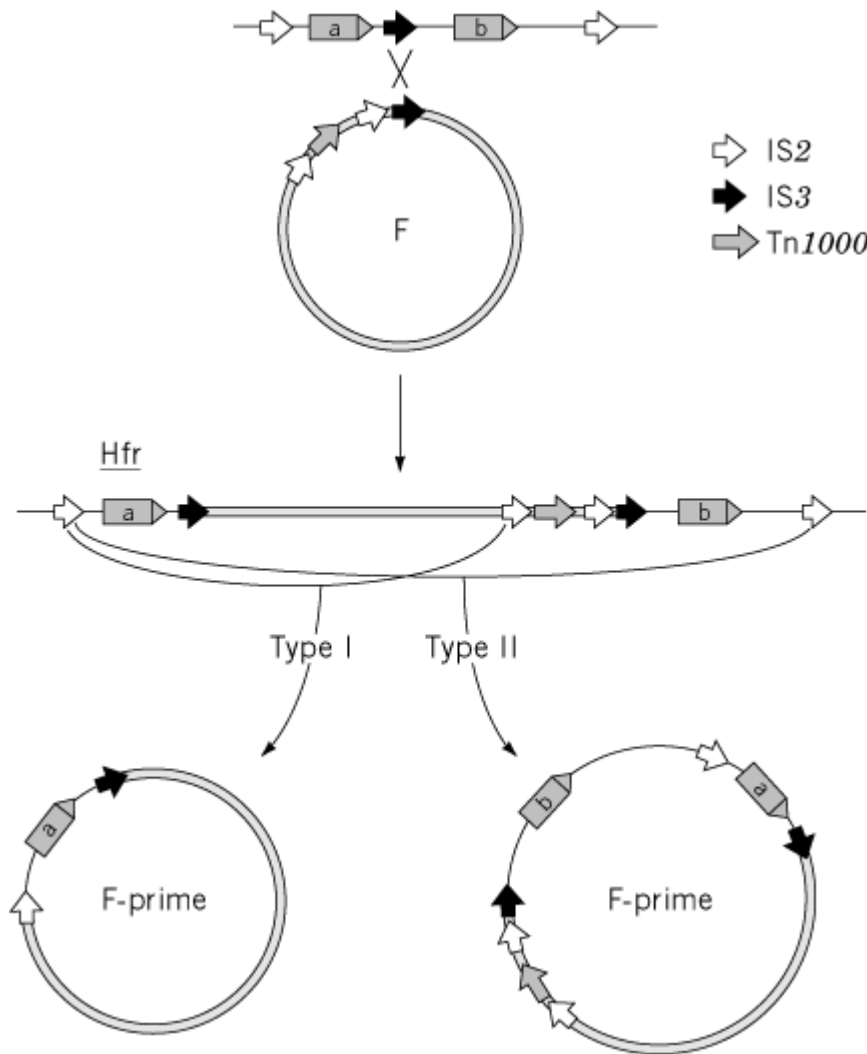
#### 5. Genesis of Hfr Strains

As more Hfr's were isolated and their points of origin determined, it became clear that certain sectors of the chromosome were favored for insertion of F ([7](#)). Moreover, the ability of F<sup>+</sup> strains to generate recombinants was severely reduced in *recA* mutants ([8](#)), implying a need for homologous recombination in Hfr formation. The reasons became evident later, when **electron microscopic** analysis of heteroduplexes formed between F and other extrachromosomal DNAs showed that F carries four known [transposable elements](#): Tn1000, IS2, and two IS3s ([9](#)). About 13 IS2s and 6 IS3s are scattered throughout the *E. coli* chromosome, and their locations correspond with the origins of transfer of most Hfr's ([10](#)). Thus, homologous recombination via a single crossover between F and the chromosome is the dominant mode of Hfr formation by F (Fig. [2](#)). Nevertheless, other



mechanisms exist, as suggested by the isolation of Hfr's in which the F is flanked by copies of Tn1000, implying integration by replicative transposition (11).

**Figure 2.** Hfr and F' formation. Homologous recombination between a transposable element on F and one on the chromosome, situated between two genes, **a** and **b**, integrates F, while the same process involving a different pair of elements excises it to form an F-prime. Transposition into targets on either the chromosome or the same plasmid creates other varieties of F-prime (not shown; see Ref. 9). Notice how these simple recombination events have changed the relationship of **a** and **b**; the Type I event has placed them on different molecules, while the Type II has reversed their order. This genome has evolved.



It is not only in their role as portable regions of [homology](#) that the transposable elements of F promote recombination in *E. coli*. At some point in the evolution of F, one of the IS3s transposed into the *finO* regulator gene (12), resulting in constitutive overexpression of the transfer system and the high transfer frequencies that enabled Lederberg and Tatum to find recombinants. F and *E. coli* K-12 thus owe their preeminence in the development of molecular genetics to these transposition events.

## 6. Genesis of F-Primes

Adelberg and Burns (13) observed that one Hfr “revertant”, which now transferred maleness at high efficiency, had undergone alterations in both its F plasmid and its chromosome, such that each now

recombined with the other at very high frequency, as if they retained “memories” of each other. The molecular basis of the memory was in fact an exchange of segments upon excision of F (14, 15).

The various forms that these exchanges take are shown in Figure 2. Once again, the patches of sequence homology provided by the transposable elements play a major role (9). Recombination between two transposable elements, one in F and one in the chromosome, excises a plasmid that now carries one chromosomal flank but has left a dispensable part of itself in the chromosome (16). Plasmids formed in this way are termed Type I F-primes. Alternatively, recombination between chromosomal sequences on both sides of the F integration site creates an undeleted F-prime, termed Type II. Any chromosomal repeated sequence is sufficient: Type II F-primes that arose by recombination between IS5 elements (17) or between **ribosomal** RNA genes (18) have been found. Type I F-primes may also arise by transposition of Tn1000 into a chromosomal target, with associated deletion of F sequences on the other side of the transposon. Type II RP4-primes might also be formed this way: An RP4::mini-Mu plasmid integrated by transposition can be excised by replicative transposition of one of the bordering mini-Mu elements into chromosomal DNA (19). An unusual event, in which F *oriT* had apparently undergone recombination with certain chromosomal loci, generated Type I F-primes lacking the *tra* operon (20).

F-prime formation has been likened to *in vivo* cloning. There is in principle no limit to the amount of chromosomal DNA that an F-prime can carry, and F-primes carrying up to 30% of the chromosome have been found. In practice however, large F-primes tend to be unstable or growth inhibitory.

Spontaneous rearrangements may be capable of generating a great variety of F-primes from a given Hfr. After excision of the initial F-prime, intramolecular transposition may delete or invert parts of the chromosomal sequence, or the F-prime may reintegrate elsewhere in the chromosome, effectively transposing its chromosomal DNA, only to re-excise as a different F-prime with a new juxtaposition of chromosomal sequences (21, 22). The potential consequences for evolution of the host genome alone are obvious (23), even without considering the possibilities raised by transfer to other species (24-26).

## 7. Utility of Hfr's and F-Primes

The major use of Hfr's has been the long-range mapping of genetic markers, and circular maps of several **gram-negative** bacteria have been established. How long their utility in this regard will continue is not certain, given the increasing rate at which entire microbial genomes are being sequenced. Hfr mating is still, however, an efficient way to assign new mutations to approximate locations on the chromosome, and it may often be the most convenient way of constructing genetically marked strains.

Whereas Hfr's find their main use in mapping, F-primes have been most important in analysis of function. Transfer of the appropriate F-prime to mutant strains allows tests of **complementation** and **dominance** at roughly equivalent copy numbers, as illustrated by the central role played by *F'**lac* plasmids in deducing the basis of regulated expression of the *lac* operon (27). F-primes have also been particularly useful in studies of the biology of plasmids themselves, as surrogates for parent plasmids that confer no readily observable phenotype. Plasmid-primes mobilized by the promiscuous conjugation system of RP4 are used for **mutagenesis** of many species, as vehicles for delivering transposons or mutant genes for allele replacement (28).

## 8. Isolating Hfr's

The relatively frequent formation of *E. coli* Hfr strains by F integration is not usually mimicked by other plasmid-host pairs. Therefore various selective methods for isolating strains with integrated plasmids have been developed.

### 8.1. Treatment with Mutagens

Nitrosoguanidine treatment enabled isolation of Hfr-ColV strains of *E. coli* (29) and Hfr-FP2 strains of *Pseudomonas aeruginosa* (30). Ultraviolet irradiation stimulates Hfr formation by F in *E. coli*. Why these treatments promote Hfr formation is unknown; stimulation of recombination by **DNA damage** might play a part.

8.2. Selection for a Plasmid-Borne Marker Under Conditions in Which the Plasmid Cannot Replicate  
Treatment of an *Erwinia* strain carrying F'*lac* with **acridine** orange, which specifically inhibits F replication, led to isolation of an Hfr strain (31). Growth at high temperatures of cells carrying temperature-sensitive (*ts*) IncP plasmids allowed isolation of Hfr's in *P. aeruginosa* (32) and *Rhodobacter capsulatus* (33). Exploiting restricted host range has served the same purpose; Hfr's have been selected after transfer of a hybrid plasmid, containing the ColE1 **replicon** linked to the RP4 *tra* region, into species of *Rhodobacter* (34) and *Methylobacillus* (35), in which ColE1 cannot replicate. Tn10-mediated transposition of the *E. coli* chromosome into a *ts* vector carrying the *tra* region generated Hfr's of any predetermined origin and orientation (36).

8.3. Rescue of Replication-Defective Chromosomes by Plasmid Replicons  
Strains of *E. coli* and *Salmonella typhimurium* with *dnaA*ts mutations are unable to grow at 42°C because initiation of chromosome replication is blocked. Temperature resistance results from the integration of plasmids able to assume control of chromosome replication, and this can be used to select Hfr strains. Several R and Col plasmids, as well as F, have given rise to Hfr's in this way (see Ref. 37 for review).

8.4. Incorporation of a Transposon into the Conjugative Plasmid  
This approach, which recalls the formation of Hfr's by Tn 1000 transposition, is more generally applicable. Two plasmids in particular have been used to transfer chromosomal DNA in a wide variety of gram-negative bacteria. RP4 carrying either the prophage Mu transposon or its mini-Mu derivative is integrated into the chromosome by replicative transposition, such that the inserted RP4 is flanked on each side by a copy of the Mu element (19). R68.45 carries a tandem duplication of IS21: Cutting between the IS21 elements and pasting into a chromosomal target site integrates the conjugative plasmid, once again flanked by copies of the transposon (38).

8.5. Conjugational Transposons  
Conjugation determinants are inserted into a transposon, and strains are selected that carry the transposon at various sites in the chromosome. Each derivative is an Hfr. Yakobsen and Guiney (39) created bipartite Hfr derivatives of *R. meliloti*: Tn5 carrying *oriT* of RP4 was transposed into the chromosome, and RP4 transfer functions were provided by a plasmid, permitting oriented transfer of chromosomal DNA. This approach essentially recapitulates the mobilization of nonconjugative plasmids, such as ColE1, by self-transmissible ones, such as F. Alternative systems in which all the transfer functions are inserted into the transposon have also been constructed (40).

8.6. Transposons as Regions of Portable Homology  
Again mimicking natural F integration, this time by homologous recombination, a plasmid carrying a transposable element can be transferred to a recipient that already has the transposable element inserted into its chromosome; growth of exconjugants in conditions restrictive for plasmid replication (temperature-sensitivity, narrow host range, etc.) allows the direct selection of Hfr strains. This approach has been used to obtain Hfr's by integration of F' *ts*::Tn10 in *S. typhimurium* (41) and R91-5::Tn5 in *Pseudomonas putida* (42). It lends itself both to the generation of banks of Hfr derivatives and to the isolation of individual Hfr's with predetermined origins of transfer.

8.7. Potential Problems  
Three potential problems beset the isolation of useful Hfr strains: (i) paucity of insertion sites, (ii) instability, and (iii) deleterious integration. Homology-based integration depends on the chance that the transmissible plasmid and the chromosome have sequences in common, and it also leaves the integrated DNA susceptible to **RecA**- or **resolvase**-mediated excision (43). Fortunately, the transposon-based strategies for Hfr isolation largely overcome these problems: Because transposon

target specificity is usually very low, the range of Hfr's obtainable is, for practical purposes, unlimited; and spontaneous transposon excision is sufficiently rare to permit formation of stable Hfr's. Most conjugative plasmids are of somewhat higher copy number than F, and their integration is often poorly tolerated (32, 43, 44). Here again, a transposon-based strategy solves the problem, because incorporation of transfer functions in a transposon uncouples them from the replicon, as described above.

## 9. Isolating F-Primes

As in the case of Hfr's, several techniques for isolating particular plasmid-primes have been developed to overcome the limitations of natural mechanisms. One of the earliest used was the selection for early transfer of a marker normally transferred late in an Hfr cross (15). F-primes carrying proximal markers can also be selected by using a *recA* mutant recipient to impair integration of the selected marker (45). This approach has been used to select an R-prime in *P. aeruginosa* (46), and a variation of it in which recombination is prevented by transfer to a recipient of a different species has led to the isolation of many R-primes (eg, Ref. 47). Some recombinants, formed by mating one Hfr strain with another that had been converted to the F<sup>-</sup> phenocopy state by prolonged incubation in stationary phase, proved to be “double males,” containing two integrated Fs: recombination between them created F<sup>-</sup> primes carrying the intervening DNA (48). A similar event appears to be responsible for the formation of certain R-primes.

The above methods all rely on genetic maneuvers to select spontaneously formed plasmid-primes. A quite distinct method has been used to produce finer gradations in the range of F-primes created. It is based on the essentially random cleavage of host DNA that occurs during **phage P1** infection and consists of simply **transducing** F (or R) DNA into F<sup>-</sup> cells (49). It has been used thus far to create shorter versions of existing F-primes, for genetic analysis (50), and to isolate an R'<sup>lac</sup> plasmid directly by infection of an Hfr strain (51).

## 10. Prospects

With the radical, and no doubt permanent, shift to physical mapping and sequencing as the preferred means of analyzing genomes, Hfr's and plasmid-primes will now serve mainly as adjuncts to other methods. Even the role of F-primes in genetic complementation has been largely taken over by convenient low copy-number cloning vectors. Nevertheless, these traditional methods should continue to be useful during the initial stages of genetic characterization of new organisms able to host promiscuous conjugative plasmids, such as RP4. Moreover, recent insights into the roles of gene mobilization in the pathogenicity of bacteria underline the importance of understanding the mechanisms of plasmid–chromosome interactions (52). It is not unreasonable to anticipate novel ones. It is also interesting to speculate on giant plasmid-primes as the origin of multichromosome bacteria (53, 54).

## Bibliography

1. J. Lederberg and E. Tatum (1946) *Nature* **158**, 558.
2. E. L. Wollman and F. Jacob (1955) *C. R. Acad. Sci.* **240**, 2449–2451.
3. F. Jacob and E. L. Wollman (1958) *Symp. Soc. Exp. Biol.* **12**, 75–92.
4. F. R. Blattner et. al. (1997) *Science* **277**, 1453–1474.
5. A. L. Taylor and E. A. Adelberg (1961) *Biochem. Biophys. Res. Commun.* **5**, 400–404.
6. A. M. Campbell (1962) *Adv. Genet.* **11**, 101–145.
7. R. Curtiss, III and D. R. Stallions (1969) *Genetics* **63**, 27–38.
8. R. C. Clowes and E. E. M. Moody (1966) *Genetics* **53**, 717–726.
9. N. Davidson, R. C. Deonier, S. Hu, and E. Ohtsubo (1975) In *Microbiology—1974* (D. Schlessinger, ed.), ASM, Washington, D.C., pp. 56–65.

10. M. Umeda and E. Ohtsubo (1989) *J. Mol. Biol.* **208**, 601–614.
11. M. S. Guyer, R. R. Reed, J. A. Steitz, and K. B. Low (1981) *Cold Spring Harbor Symp. Quant. Biol.* **45**, 135–140.
12. K.-C. Cheah and R. A. Skurray (1986) *J. Gen. Microbiol.* **132**, 3269–3275.
13. E. A. Adelberg and S. N. Burns (1960) *J. Bacteriol.* **79**, 321–330.
14. A. Richter (1957) *Genetics* **42**, 391.
15. F. Jacob and E. A. Adelberg (1959) *C. R. Acad. Sci.* **249**, 189–191.
16. P. A. Sharp, M.-T. Hsu, E. Ohtsubo, and N. Davidson (1972) *J. Mol. Biol.* **71**, 471–497.
17. M. S. Timmons, A. M. Bogardus, and R. C. Deonier (1983) *J. Bacteriol.* **153**, 395–407.
18. D. L. Blazey and R. O. Burns (1983) *J. Bacteriol.* **156**, 1344–1348.
19. F. van Gijsegem and A. Toussaint (1982) *Plasmid* **7**, 30–44.
20. R. G. Hadley and R. C. Deonier (1980) *J. Bacteriol.* **143**, 680–692.
21. J. M. Buysse and S. Palchaudhuri (1984) *Mol. Gen. Genet.* **193**, 543–553.
22. E. Ohtsubo and M.-T. Hsu (1978) *J. Bacteriol.* **134**, 795–800.
23. M. Riley and A. Anilionis (1978) *Annu. Rev. Microbiol.* **32**, 519–560.
24. K. E. Sanderson and C. A. Hall (1970) *Genetics* **64**, 215–228.
25. M. Mergeay and J. Gerits (1978) *J. Bacteriol.* **135**, 18–28.
26. L. A. Wiater, A. Marra, and H. A. Shuman (1994) *Plasmid* **32**, 280–294.
27. F. Jacob, D. Perrin, C. Sanchez, and J. Monod (1960) *C. R. Acad. Sci.* **250**, 1727–1729.
28. R. Simon, U. Priefer, and A. Pühler (1983) *Biotechnology* **1**, 784–791.
29. P. L. Kahn (1968) *J. Bacteriol.* **96**, 205–214.
30. H. Herrmann and E. Günter (1984) *Mol. Gen. Genet.* **197**, 286–291.
31. A. K. Chatterjee and M. P. Starr (1973) *J. Bacteriol.* **116**, 1100–1106.
32. C. Riemann, M. Rella, and D. Haas (1988) *J. Gen. Microbiol.* **134**, 1515–1523.
33. A. Colbeau, J.-P. Magnin, B. Cauvin, B. T. Champion, and P. M. Vignais (1990) *Mol. Gen. Genet.* **220**, 393–399.
34. S. V. Kameneva, T. P. Politseva, N. V. Belavina, and S. V. Shestakov (1986) *Genetika* **22**, 2664–2672.
35. I. G. Serebrijski, S. M. Kazakova, and Y. D. Tsygankov (1989) *FEMS Microbiol. Lett.* **59**, 203–206.
36. V. Fran, A. Conter, and J.-M. Louarn (1990) *J. Bacteriol.* **172**, 1436–1440.
37. C. Reimann and D. Haas (1993) In *Bacterial Conjugation* (D. B. Clewell, ed.), Plenum Press, New York, pp. 137–188.
38. N. S. Willets, C. Crowther, and B. W. Holloway (1981) *Plasmid* **6**, 30–52.
39. E. A. Jakobson and D. G. Guiney (1984) *J. Bacteriol.* **160**, 451–453.
40. D. A. Johnson (1988) *Plasmid* **20**, 249–258.
41. F. G. Crumley, R. Menzel, and J. R. Roth (1979) *Genetics* **91**, 639–655.
42. A. D. Strom, R. Hirst, J. Petering, and A. Morgan (1990) *Genetics* **126**, 497–503.
43. C. Riemann and D. Haas (1986) *Mol. Gen. Genet.* **203**, 511–519.
44. D. Haas and B. W. Holloway (1976) *Mol. Gen. Genet.* **144**, 243–251.
45. B. Low (1968) *Proc. Natl. Acad. Sci. USA* **60**, 160–167.
46. B. W. Holloway (1978) *J. Bacteriol.* **133**, 1078–1082.
47. R. W. Hedges, A. E. Jacob, and I. P. Crawford (1977) *Nature* **267**, 283–284.
48. A. J. Clark (1963) *Genetics* **48**, 105–120.
49. J. Pittard and E. A. Adelberg (1963) *J. Bacteriol.* **85**, 1402–1408.

50. N. J. Marsh and D. E. Duggan (1972) *J. Bacteriol.* **109**, 730–740.
51. A. Nishimura, Y. Nishimura, and L. Caro (1973) *J. Bacteriol.* **116**, 1107–1112.
52. S. Falkow (1996) In *Escherichia coli and Salmonella. Cellular and Molecular Biology*, 2nd ed. (F. C. Neidhardt et al., eds.), ASM Press, Washington, D.C., pp. 2723–2729.
53. W. Hayes (1968) *The Genetics of Bacteria and Their Viruses*, 2nd ed., Blackwell Scientific Publications, Oxford, p. 674.
54. C. Rosenberg, P. Boistard, J. Denarié, and F. Casse-Delbart (1981) *Mol. Gen. Genet.* **184**, 326–333.

### Suggestions for Further Reading

55. F. C. Neidhardt et al., eds. (1996) *Escherichia coli and Salmonella. Cellular and Molecular Biology*, 2nd ed., ASM Press, Washington, D.C. Chapters 126–129 in Vol. **2**, Section C provide a comprehensive treatment of the subject, centered on *E. coli*, as well as an exhaustive bibliography.
56. C. Reimann and D. Haas (1993) "Mobilization of chromosomes and non-conjugative plasmids by cointegrative mechanisms." In *Bacterial Conjugation*, (D. B. Clewell, ed.), Plenum Press, New York. An excellent starting point for accessing the available information on conjugal transfer in a wide range of bacteria.
57. W. Hayes (1968) *The Genetics of Bacteria and Their Viruses*, 2nd ed. Blackwell Scientific Publications, Oxford Chapters 22 ("Conjugation") and 24 ("Sex Factors and Other Plasmids") of this textbook is a lucid account of the historical development of the subject, written by one of its pioneers.
58. T. D. Brock (1990) *The Emergence of Bacterial Genetics*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York. Chapter 5 ("Mating") is another historical perspective. The science, enlivened by anecdote, is treated more lightly than in Hayes, but without sacrifice of rigor or clarity.

### High Mobility Group (HMG) Proteins

The major proteins of [chromosomes](#) and [chromatin](#) are the [histones](#), but other, nonhistone chromosomal components are also present, in much lower amounts. The most abundant and best characterized of these is the class of high mobility group (HMG) proteins. (The name is an operational definition coined when the proteins were first isolated in the 1970s as proteins that were soluble in 2% trichloroacetic acid or 5% perchloric acid and migrated with a high mobility in polyacrylamide [gel electrophoresis](#).) HMG proteins in metazoans fall into three structurally and functionally distinct classes: the HMG 1, 2 family, HMG 14 and 17, and the HMG I(Y) family. (The numbering is also a legacy from the past!) They appear to be present in all cells of higher eukaryotes. They are relatively abundant, being present in the [nucleus](#) on average at about 1 molecule per 10 to 15 [nucleosomes](#) [considerably less for HMG I(Y)]. All have intriguing properties, and their cellular roles are not fully understood, but they are likely to play important roles in **gene** expression. None of the three classes shows specificity for particular DNA sequences; instead, they bind to particular structures in DNA or to chromatin. Both HMG 1,2 and HMG I(Y) have "architectural" roles, and both bind primarily to DNA in the minor groove, although they are wholly unrelated in amino acid sequence and structure. The early work on these proteins is described in a definitive book ([1](#)). A recent review ([2](#)) gives a comprehensive account of more recent work and is an excellent source of references up to 1996, which in general will not be duplicated here. For more general background,

the reader should consult two volumes on chromatin ([3](#), [4](#)).

## 1. The HMG 1 and 2 Family

In vertebrates this class consists of HMG1 (molecular mass ~25,000Da) and two closely related forms of HMG2, which are a few residues shorter. They are the products of **homologous** genes and show a high degree of evolutionary conservation. They have a tripartite structure consisting of two tandem homologous regions of about 80 amino acid residues (the two DNA-binding “HMG boxes”) and a long acidic tail of about 30 (HMG 1) or 20 (HMG 2) consecutive **aspartic** and **glutamic acid** residues linked to the second box by a short basic region. HMG 1 and 2 bind to DNA with a “structure-preference” rather than a sequence preference. They bind to bent or bulged DNA and to DNA kinked by the antitumor drug *cis*-platin, in preference to linear DNA; they also prefer four-way junctions, **supercoiled** DNA and DNA minicircles, and they stabilize DNA loops (the juxtaposed ends probably resemble crossovers in supercoiled DNA). They bend linear DNA, as shown by the ability to cause circularization (in the presence of **DNA Ligase** to join the ends) of short DNA fragments (~80bp) that will not circularize unaided, and they constrain negative supercoils in relaxed circular DNA. The ability to distort DNA, as well as to recognize DNA distortions reflected in these properties, is likely to reflect the *in vivo* role(s) of HMG 1 and 2.

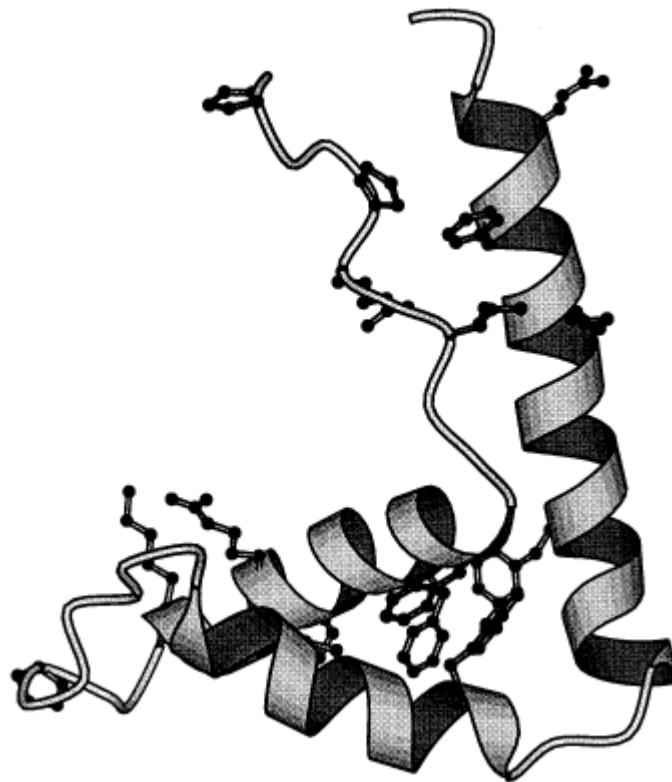
Relatively abundant HMG-box proteins that bind to DNA without sequence specificity also occur in yeast (NHP6A and B), insects (eg, HMG-D and HMG-Z in *Drosophila*; HMG1a and b in *Chironomus*), and plants (HMGa, etc.). In all these cases, the proteins contain only one HMG box and the acidic tail is often shorter than in vertebrate HMG1,2 (eg, 10 residues in HMG-D). All these proteins bind preferentially to distorted DNA substrates (although the *Chironomus* HMG 1a and 1b bind A-T-rich linear and four-way junction DNA with equal affinity). The HMG-box motif also occurs in a large and growing number of sequence-specific **transcription factors**, as a single copy, with no acidic tail, and embedded in unrelated sequence ([5](#), [6](#)). The first examples were the product of the male sex-determining gene on the mouse Y chromosome (SRY) and the lymphocyte enhancer binding protein (LEF-1). The specificity of these proteins for their DNA targets resides in the HMG box. They also bind to four-way DNA junctions, and they bend linear DNA containing their DNA binding sites ([7](#)). The estimated bend angles are ~130° for LEF-1 and ~90° for mouse and human SRY. Proteins with two or more HMG boxes bind DNA with relatively low sequence specificity. These include the **RNA polymerase I** transcription factor UBF (upstream binding factor; 4–6 HMG boxes, depending on source, and a very long acidic tail) and the mitochondrial factors mtTF1 and ABF2.

This diverse group of sequence-specific and non-sequence-specific HMG-box **DNA-binding proteins** is united by the ability to distort and bend DNA, which is probably central to their biological roles. They are regarded as “architectural elements in the assembly of nucleoprotein structures” ([7](#)). Their properties are analogous to those of the bacterial IHF (integration host factor) and **HU proteins**. Indeed HMG-D, NHP6A and HM can all functionally replace HU in *Escherichia coli*, and HMG 1 and 2 can functionally replace HU in facilitating assembly of the invertasome during Hin-mediated **recombination** .

**Protein structures** have been determined by **NMR** for the A and B HMG boxes of HMG 1, the HMG boxes of *Drosophila* HMG-D and the transcription factor Sox4, and for the HMG boxes of SRY and LEF-1 complexed with oligonucleotides. Although there are slight differences that might be significant, the protein folds are all generally very similar, showing that there is no gross change in the HMG-box structure on binding to DNA. The HMG box is L-shaped, the relative positions of the two arms being fixed by a tightly packed **hydrophobic** cluster; Figure [1](#) shows the A-box of HMG1 ([8](#)). The “long” arm contains an *N*-terminal extended **b-strand** packed against **a-helix III** in the *C*-terminal half of the molecule; the “short” arm consists of a-helices I and II. DNA binding occurs on the concave face of the protein, which interacts in the minor groove of the DNA. The SRY/DNA and LEF-1/DNA structures ([9](#), [10](#)) show that the protein partially intercalates into the minor groove and causes the DNA to bend toward the major groove, away from the protein (Fig. [2](#)). The minor groove

is also opened up (becomes wider and shallower), forming a larger surface for interaction with the protein, and the DNA is locally untwisted. The protein thus sits on the *outside* of the bend and is very reminiscent of the way in which the **TATA box**-binding protein TBP interacts with DNA. In both structures a bulky hydrophobic residue (Ile in SRY, Met in LEF-1) intercalates between base pairs near the center of the DNA to cause bending. Binding is stabilized by interaction of a basic region *C*-terminal to the box with the major groove opposite the widened minor groove. All HMG boxes are likely to bind to DNA in essentially the same way, except that the role of intercalation in the non-sequence-specific proteins is less clear. Sequence-specific and non-sequence-specific HMG boxes also show characteristic differences in residues at key positions in the *N*-terminal b-strand.

**Figure 1.** Structure of the HMG box in the A-domain of HMG1 (residues 11-83 of the protein) determined by NMR spectroscopy (8).



**Figure 2.** Structure of the LEF-1 HMG box complexed with DNA complex, determined by NMR spectroscopy (10).





The isolated A and B boxes of HMG 1 bend and distort DNA like the entire protein; there are slight differences between them that might be related to the small differences in structure between the two boxes (8). Tandem boxes have a higher affinity for DNA than single boxes and are more effective than single boxes in DNA binding and bending; the basic linker between the boxes and the acidic tail further enhances the affinity (eg, Ref. 11). The acidic tail works in the opposite direction on most DNA substrates *in vitro*, by lowering the affinity of the HMG boxes for DNA, possibly by general charge repulsion, and possibly by interacting electrostatically with the DNA-binding face of the box (es). The acidic tail of HMG-D, however, increases the affinity and selectivity of the protein for mini-circles (<100bp), which are highly constrained and effectively present a “pre-bent” substrate (12); the same is true for HMG 1. Perhaps the tail stabilizes a conformation of the protein that is suited to recognition of a highly pre-distorted substrate.

Surprisingly little is known for certain about the binding of HMG 1 and 2 to [chromatin](#). Early observations on micrococcal nuclease (**Staphylococcal nuclease**) digestion products of chromatin (3, 4) suggested that HMG 1 and 2 were associated with [linker DNA](#), and that there was a population of HMG1-containing nucleosomes lacking histone H1. This led to the suggestion that HMG 1 binds to sites normally occupied by H1, but the evidence for this is still circumstantial. However, binding of H1 and HMG 1 to nucleosomes reconstituted on to the *Xenopus* 5 S rRNA gene resulted in similar protection of chromosome length DNA against micrococcal nuclease digestion (see [Chromatin](#) and [Nucleosome](#)), which was interpreted in terms of a shared structural role (13). The finding that, like HMG 1, H1 and its globular domain would also bind to synthetic four-way junctions (14) has been taken to reinforce this (the four-way junction is assumed to mimic the crossover of the entering and exiting duplexes around the dyad of the nucleosome). However, it is not clear that the two proteins are recognizing the same features of the four-way junction. Moreover, given the nature of H1 binding to the nucleosome (see [Nucleosome](#) or [Chromatin](#)), even if H1 and HMG 1 bind in the same general vicinity, the details of their interactions in chromatin would be expected to be very different. *In vivo* evidence for a shared *functional* role in at least some circumstances is more compelling. In *Drosophila* early embryos, where cell division is rapid, H1 is absent and the condensed chromosomes contain HMG-D (which has only one HMG box). At later stages, the chromosomes contain H1 rather than HMG-D, and gene [transcription](#), which was previously absent, starts (15). This apparently reciprocal relationship suggests that H1 and HMG-D

may have a common function at different stages of embryogenesis in *Drosophila*, although not necessarily an identical binding site on the nucleosome. Rigorous reconstitution studies are needed to settle the question of how HMG 1 and 2 bind to chromatin. It is, of course, possible that they function only at **promoters** or **enhancers**, where a nucleosome is absent. For example, by bending DNA they might facilitate the interaction of two proteins, bound on either side of the kink, that is required for some functional purpose; or they could facilitate the binding of a protein that binds best to a DNA bend, by pre-bending the DNA in an essentially catalytic role; or they might also stabilize a DNA loop, such as might be formed to bring enhancer and promoter elements into proximity.

HMG 1 and 2 have been variously implicated in **DNA replication**, transcription, **DNA repair** and **recombination**, and, given their versatility in DNA distortion and recognition of distorted DNA, this is not altogether surprising. In recent *in vitro* studies of V(D)J recombination of **immunoglobulin biosynthesis** (see **Gene Rearrangement**), HMG 1 and 2 were found to stimulate V(D)J cleavage (16) and to be components of a stable post-cleavage complex between synapsed recombination signals (17). An obvious question is whether HMG 1 and 2 have roles in transcriptional activation; distinct roles for the HMG boxes and the acidic tails might be envisaged. *In vitro* studies (for references see Ref. 2) have given conflicting answers. It has been suggested that activation might result from stabilization of an activated conformation of the TFIID-TFIIA initiation complex on the promoters of activated genes; or, based on *in vitro* binding assays, that HMG 1 and 2 might act indirectly to promote transcription by promoting the binding of various transcription factors to their cognate DNA binding sites, probably by bending the DNA. *In vivo* results (**transfection** experiments) show enhancement of transcription by HMG 1 and suggest that the acidic tail acts as a transactivation domain, although this is not the case for the heterologous mammalian HMG 1 expressed in yeast (2). However, HMG 1 and 2 have also been reported to repress transcription by RNA polymerase II *in vitro*—either by interacting with TBP in the presence of a TATA-box-containing oligonucleotide and thus preventing binding of TFIIB and formation of a preinitiation complex, or by acting later, after the assembly of the TBP-TFII promoter complex. There are now several examples of facilitation of transcription factor binding to DNA by HMG 1 and 2, and this seems likely to be important. For example, binding of the progesterone receptor to an oligonucleotide containing the progesterone **response element** is enhanced ~10-fold by HMG 1 and 2, and HMG 2 stimulates the sequence-specific binding of the octamer transcription factors Oct 1 and Oct 2, by interacting with the **POU domains** (2). HMG 1 stimulates DNA binding of human HOXD9, by interaction with the **homeobox** domain (18), and of the tumor suppressor **p53** (19). Ternary complexes between the transcription factor, HMG protein, and DNA are likely, but in at least some cases the HMG must be only weakly bound, because such complexes cannot be detected by gel electrophoresis.

## 2. The HMG I(Y) Family

HMG I and Y are isoforms generated by **alternative splicing**, differing only in HMG Y (96 amino acid residues) having an 11-residue deletion compared with HMG I (107 residues). HMG C (105 residues) is the third member of this family and is related (~50% sequence identity) to HMG I(Y). Like HMG 1 and 2, they appear to function as “architectural” transcription factors, but the two classes are structurally unrelated. The proteins are subject to **post-translational modifications**, notably **cell cycle-dependent reversible phosphorylation**, probably by p34cdc2/cyclin (“cdc2 kinase”). They are expressed at higher levels (~15 to 50 times) in transformed and undifferentiated cells than in differentiated somatic cells, relatively independent of cellular growth rate; there is also a correlation with increased metastatic tumor potential. Further background and references may be found in Ref. 2.

The HMG I(Y) proteins have a distinctive primary structure containing three regions designated “A-T hooks,” separated by flexible chain. HMG C is similar in the A-T hook regions, but much less so elsewhere. The A-T hook binds to the minor grooves of A-T-rich sequences, occupying 5 or 6 bp; three tandem AT hooks would therefore occupy 15 to 18 bp of continuous A-T DNA. The **consensus sequence** of the A-T hook motif is -Pro-Arg-Gly-Arg-Pro- flanked by basic residues; the motif also occurs in a number of unrelated DNA-binding proteins in a range of species. The polypeptide

backbone structure is predicted to be similar in shape to the drugs distamycin A, netropsin, and Hoechst 33258, which bind to the minor groove of B DNA and can displace HMG I(Y) from A-T-rich DNA. The structural similarity is borne out by NMR spectroscopy (see text below). In addition to minor groove binding—although in completely different ways—HMG I(Y) also shares with HMG 1 and 2 the property of binding to A-rich non-**B-form** DNA, such as four-way junctions, in preference to the corresponding duplexes (20). HMG I(Y) also bends DNA (but in this case bending is attributed simply to asymmetric charge neutralization on one face of the DNA duplex by the basic amino acid side chains in the A-T hook DNA-binding domains) and supercoils relaxed circular DNA in the presence of **topoisomerase** I, probably through strand unwinding, the mechanism of which is unclear. Removal of the negatively charged C-terminal domain increases supercoiling by 8- to 10-fold, apparently without affecting the affinity of HMG I(Y) for A-T-rich DNA (2).

Earlier studies suggested a variety of roles for HMG I(Y) proteins (2). More recently, more direct experiments point to a role in transcriptional regulation (negative or positive) of a number of mammalian genes with A-T-rich promoter or enhancer sequences. In positive regulation, HMG I(Y) has been proposed to act as an “architectural” transcription factor that bends DNA and contacts components of a multiprotein complex, as well as promoting contacts between other members. A well-documented example (21) is the involvement of HMG I(Y) in the virus-induced expression of the human b-*interferon* gene (IFN-b). HMG I(Y) appears, through a combination of DNA bending and [protein-protein interactions](#), to mediate the assembly of a complex containing HMG I(Y) and the transcription factors NF-κB and ATF-2/c-Jun, all bound simultaneously to two separate “positive regulatory domains” (PRDII and IV) of DNA in the 5′ promoter/enhancer region of the b-*interferon* gene. NF-κB binds to the major groove and HMG I(Y) to the minor groove. The structure of a 39-residue peptide containing the second and third DNA-binding domains complexed with a 12-bp oligonucleotide containing a 5-bp A-T tract from the PRDII element of the b-*interferon* enhancer has recently been determined by NMR spectroscopy (22) and reveals extensive hydrophobic and polar contacts in the minor groove centered around an Arg-Gly-Arg motif. In contrast to the HMG box–DNA complexes (see text above), the minor groove is not widened and the DNA is not bent. However, this is a short piece of DNA and, although the protein does not directly cause bending at its binding site, as the HMG box does, bending might well be induced outside the binding site and would be apparent only with a longer DNA segment. Principles similar to those involved at the b-*interferon* promoter are likely to prove a common feature of the role of HMG I(Y) in many systems—for example, in the function of HIV-1 pre-integration complexes (23).

The situation with respect to the binding of HMG I(Y) to chromatin is not clear. The proteins will bind to nucleosome core particles *in vitro*, at up to four molecules per nucleosome, in a noncooperative manner. However, they bind to A-T-rich DNA in preference to mixed-sequence nucleosomes and, as pointed out (2), if A-T sequences were present in the linker DNA—or in a nucleosome-free region—that would probably be the preferred binding site for HMG I(Y) in chromatin. There is an intriguing possible connection between HMG I(Y) and chromatin structure through histone H1 (see [Histones](#)). *In vitro* H1 will bind preferentially to A-T-rich fragments of naked DNA, possibly through “-SPKK-” (Ser-Pro-Lys-Lys) motifs in the C-terminal tail, and can be displaced by HMG I(Y). In the nucleus, scaffold attachment regions of chromatin are A-T-rich and have been suggested to be focal points for displacement of H1 by HMG I(Y), perhaps leading to “open chromatin” domains (24); in neoplastic cells, greatly increased levels of HMG I(Y) might contribute to altered patterns of gene expression. This is still speculation, but it is highly suggestive that HMG I(Y) co-localize (as shown by immunofluorescence) with A-T-rich scaffold regions in metaphase chromosomes, which contain the interphase scaffold attachment regions brought together (25).

### 3. HMG 14 and 17

These evolutionarily related proteins are present in all tissues of most higher organisms. Their precise cellular roles are not understood, although various lines of evidence suggest that they facilitate transcription of chromatin, and it now appears that this is at the level of higher-order

structure. Early work reported a correlation between DNase I-sensitivity of chromatin, which itself correlates with transcriptional competence (see [Chromatin](#)), and the presence of HMG 14 and 17 (see Refs. [2](#) and [3](#) for background). HMG 14 and 17 contain, respectively, 98 and 89 amino acid residues and have a high content of lysine, alanine and proline and a negatively charged C-terminal region (although there is no continuous run of acidic residues, as in HMG 1 and 2). They bind to nucleosomes in preference to free DNA through a 30-residue, basic sequence that is conserved within the HMG 14 and HMG 17 classes, but differs in detail between the two ([2](#)). HMG 14 and 17, despite their common ancestry, are only 60% identical in sequence and are likely to perform different roles in the cell.

Two copies of HMG 14 and 17 bind cooperatively to 146-bp nucleosome core particles (see [Nucleosome](#)) at roughly physiological ionic strength, presumably by recognizing particular structural features. In native chromatin, nucleosomes containing only HMG 14 or HMG 17 are clustered in runs of, on average, six nucleosomes, potentially giving functionally distinct domains along the chromatin fiber, the significance of which is not yet clear ([26](#)). The free proteins appear to be unstructured in solution but could well become structured upon interaction with the nucleosome core. Binding of HMG 14 is weakened by **phosphorylation** of Ser 6, which is an early event in the induction of immediate-early genes on mitogenic stimulation ([27](#)). [Footprinting](#) and [cross-linking](#) data have led to a model for the location of the two HMG 14 or 17 molecules on the core particle, in which the HMG proteins are bound by their basic N-terminal regions to sites 20 or 30 base pairs from the end of the core particle DNA and then loop under one of the DNA ends to contact, in each case, a major groove adjacent to the dyad on the central turn of DNA ([2](#)). This could explain why HMG 14 and 17 stabilize nucleosome core particles, by preventing charge repulsion that would lead to unraveling of the ends. In contrast, HMG 14 and 17 appear to destabilize chromatin higher-order structure in some way, although they do not prevent its formation ([2](#)), consistent with the facilitating effect of HMG 14 and 17 on transcription. When tested for activator function as fusions to DNA-binding proteins in yeast, HMG 14 and 17 did not act as classical transcription factors, despite an acidic C-terminal region. However, transcription from [minichromosomes](#) assembled in *Xenopus* or *Drosophila* embryo extracts was stimulated by HMG 14 and 17, in contrast with transcription from naked DNA, consistent with a role for HMG 14 and 17 in destabilizing chromatin structure ([2](#)). Enhancement of transcription from assembled [SV40](#) minichromosomes by HMG14 was interpreted as relief of H1-mediated repression, resulting in unfolding, for which the acidic C-terminal region of HMG 14 is needed ([28](#)). Specific interactions with the N-terminal tail of H3 may be involved ([29](#))—yet another role for the N-terminal tails (see [Chromatin](#)).

## Bibliography

1. E. W. Johns (1982) *The HMG Chromosomal Proteins*, Academic Press, London.
2. M. Bustin and R. Reeves (1996) *Prog. Nucleic Acid Res. and Mol. Biol.* **54**, 35–100.
3. K. van Holde (1988) *Chromatin*, Springer-Verlag, New York.
4. A. P. Wolffe (1995) *Chromatin: Structure and Function*, 2nd ed., Academic Press, New York.
5. V. Laudet, D. Stehelin, and H. Clevers (1993) *Nucleic Acids Res.* **21**, 2493–2501.
6. A. D. Baxevanis and D. Landsman (1995) *Nucleic Acids Res.* **23**, 1604–1613.
7. R. Grosschedl, K. Giese, and I. Pagel (1994) *Trends Genet.* **10**, 94–100.
8. C. H. Hardman et al. (1995) *Biochemistry* **34**, 16596–16607.
9. M. H. Werner, J. R. Huth, A. M. Gronenborn, and G. M. Clore (1995) *Cell* **81**, 705–714.
10. J. J. Love et al. (1995) *Nature* **376**, 791–795.
11. K. D. Grasser et al. (1998) *Eur. J. Biochem.* **253**, 787–795.
12. D. Payet and A. A. Travers (1997) *J. Mol. Biol.* **266**, 66–75.
13. K. Nightingale, K. Dimitrov, R. Reeves, and A. P. Wolffe (1996) *EMBO J.* **15**, 548–561.
14. P. Varga-Weisz, K. van Holde, and J. Zlatanova (1993) *J. Biol. Chem.* **268**, 20699.
15. S. S. Ner and A. A. Travers (1994) *EMBO J.* **13**, 1817–1822.

16. D. C. van Gent, K. Hiom, T. T. Paull, and M. Gellert (1997) *EMBO J.* **16**, 2665–2670.
17. A. Agrawal and D. G. Schatz (1997) *Cell* **89**, 43–53.
18. V. Zappavigna et al. (1996) *EMBO J.* **15**, 4981–4991.
19. L. Jayaraman et al. (1998) *Genes Dev.* **12**, 462–472.
20. D. A. Hill and R. Reeves (1997) *Nucleic Acids Res.* **25**, 3523–3531.
21. W. Du, D. Thanos, and T. Maniatis (1993) *Cell* **74**, 887–898.
22. J. R. Huth et al. (1997) *Nat. Struct. Biol.* **4**, 657–665.
23. C. M. Farnet and F. D. Bushman (1997) *Cell* **88**, 483–492.
24. K. Zhao, E. Kas, E. Gonzalez, and U. K. Laemmli (1993) *EMBO J.* **12**, 3237–3247.
25. Y. Saitoh and U. K. Laemmli (1994) *Cell* **76**, 609–622.
26. Y. V. Postnikov et al. (1997) *J. Mol. Biol.* **274**, 454–465.
27. J. M. Barratt, C. A. Hazzalin, E. Cano, and L. C. Mahadevan (1994) *Proc. Natl. Acad. Sci. USA* **91**, 4781–4785.
28. H. F. Ding, M. Bustin, and U. Hansen (1997) *Mol. Cell. Biol.* **17**, 5843–5855.
29. L. Trieschmann, B. Martin, and M. Bustin (1998) *Proc. Natl. Acad. Sci. USA* **95**, 5468–5473.

### **Suggestions for Further Reading**

30. D. M. Crothers (1993) Architectural elements in nucleoprotein structures. *Curr. Biol.* **3**, 675–676.
31. K. D. Grasser (1998) HMG1 and HU proteins: architectural elements in plant chromatin. *Trends Plant Sci.* **3**, 260–265.
32. R. Grosschedl (1995) Higher-order nucleoprotein complexes in transcription: analogies with site-specific recombination. *Curr. Opin. Cell Biol.* **7**, 362–370.
33. S. S. Ner, A. A. Travers, and M. E. Churchill (1994) Harnessing the writhe: a role for DNA chaperones in nucleoprotein complex formation. *Trends Biochem. Sci.* **19**, 185–187.

## **High-Performance Liquid Chromatography (HPLC)**

HPLC is a modern high-resolution liquid **chromatographic** technique best described as high-performance liquid chromatography because the essence of the technique is the high resolution of the separations achieved by using uniform microparticulate, chromatographic supports and well-designed equipment. HPLC has also been called high-pressure or high-speed liquid chromatography but these acronyms do not reflect the essential features of the technique.

The development of HPLC has been in several areas including theory, special columns, and equipment. Though the theoretical principles of HPLC were firmly established by Martin and Synge (1) in the early 1950s, HPLC in practice did not appear until the late 1960s because of instrumental problems. Horváth et al. (2) constructed one of the first practical HPLC apparatuses for use in their research on **nucleotides**. Subsequently, dramatic developments in packing materials, particle sizes (as small as 3  $\mu\text{m}$ ), narrow-bore columns (as small as 0.5 mm inner diameter for microbore columns) and high column inlet pressures (up to 12,000 psi) have been achieved for high resolution and efficiency.

HPLC is carried out in all classical modes of column chromatography involving a liquid mobile

phase. The most important are liquid-solid adsorption, liquid-liquid and organo-bonded partition, **ion-exchange**, **size-exclusion**, and [affinity chromatography](#). HPLC has the advantages over classical liquid chromatography that the columns are reusable, the sample introduction can be automated, flow rates can be controlled precisely, and detection and quantification can be achieved by using continuous flow detectors. These features have led to improved analytical accuracy and precision.

The widespread application of HPLC in biochemical studies is evidenced by countless publications in the field. For example, HPLC is widely used by molecular biologists to isolate **nucleic acids** from synthetic oligonucleotides to natural **plasmids** (see reviews 3, 4 and Chapter 11.1 in 5), and by protein biochemists to purify **peptides** and [proteins](#) (see reviews 6–10 and Chapter 11.2 in 5). Those requiring more detailed practical information are directed to the reviews cited previously and also to a number of excellent monographs (5), (11-13).

## Bibliography

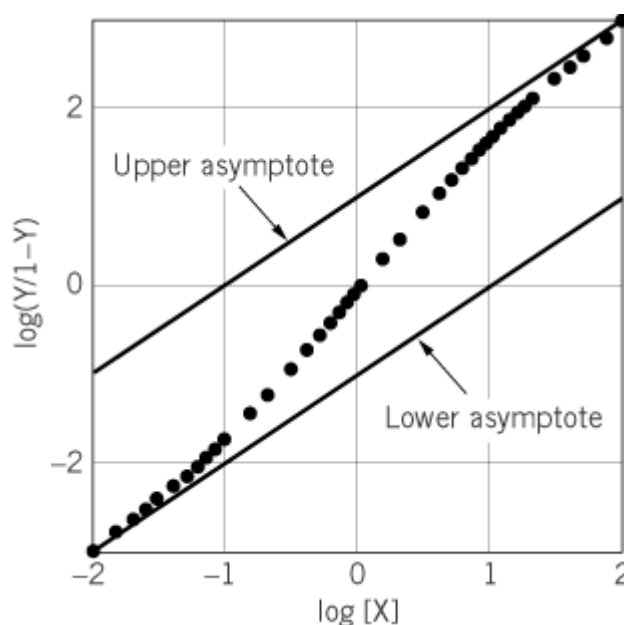
1. A. J. P. Martin and R. L. M. Synge (1941) *Biochem. J.* **35**, 1358–1368.
2. C. Horváth, B. A. Preiss, and S. R. Lipsky (1967) *Anal. Chem.* **39**, 1422–1428.
3. M. Colpan and D. Riesner (1988) in *Modern Physical Methods in Biochemistry, Part B* (A. Neuberger and L. L. M. Van Deenen, eds.), Elsevier, Amsterdam, pp. 85–105.
4. P. R. Brown (ed.) (1984) *HPLC in Nucleic Acid Research* (Chromatographic Science, Vol.28), Marcel Dekker, New York.
5. A. Fallon, R. F. G. Booth, and L. D. Bell (1987) *Applications of HPLC in Biochemistry*, Elsevier, Amsterdam.
6. W. S. Hancock (1984) *Handbook of HPLC for the Separation of Amino Acids, Peptides, and Proteins, Vols. I and II*, CRC Press, Boca Raton.
7. M. T. W. Hearn and M. I. Aguilar (1988) in *Modern Physical Methods in Biochemistry, Part B* (A. Neuberger and L. L. M. Van Deenen, eds.), Elsevier, Amsterdam, pp. 107–142.
8. A. R. Kerlavage (ed.) (1989) *The Use of HPLC in Receptor Biochemistry*, Alan R. Liss, New York.
9. C. T. Mant and R. S. Hodges (eds.) (1991) *High-Performance Liquid Chromatography of Peptides and Proteins: Separation, Analysis, and Conformation*, CRC Press, Boca Raton.
10. R. M. Chicz and F. E. Regnier (1990) in *Guide to Protein Purification* (M. P. Deutscher, ed.), *Methods in Enzymology* 182, Academic Press, New York, pp. 392–421.
11. H. Engelhardt (ed.) (1986) *Practice of High Performance Liquid Chromatography: Applications, Equipment and Quantitative Analysis*, Springer-Verlag, Berlin.
12. U. D. Neue (1997) *HPLC Columns: Theory, Technology, and Practice*, Wiley-VCH, New York.
13. A. Weston and P. R. Brown (1997) *HPLC and CE: Principles and Practice*, Academic Press, San Diego.

## Hill Coefficient, Plot

The Hill coefficient is used to provide a quantitative measure of cooperativity of ligand binding. It is usually estimated graphically from a Hill plot. A hypothetical example of a Hill plot is shown in Figure 1. To construct a Hill plot,  $y$ , the fractional saturation of the binding sites by a ligand X, is determined experimentally. The data are plotted as  $\log (y/1-y)$  versus  $\log [X]$ . The Hill coefficient,

$n_H$  is given by the slope of this plot at  $\log(y/1-y) = 0$ —that is, at  $y = 0.5$  or 50% saturation of the X binding sites.

**Figure 1.** Hill plot. A hypothetical Hill plot for positively cooperative binding of ligand X to a macromolecule with two binding sites for X is shown. The data are presented as points, and the limiting asymptotes with slope 1 at high and low extents of ligand binding are shown as lines. The intrinsic stepwise binding constants,  $K_1$  and  $K_2$ , given by the intercepts at  $\langle rf \rangle \log \langle rf \rangle [X] = 0$ , are 0.1 and 10, respectively. The Hill coefficient is 1.8.



The value of the Hill coefficient provides two important items of information. First, its value relative to 1 indicates the type of cooperativity. For noncooperative ligand binding,  $n_H = 1$ . For positively cooperative ligand binding,  $n_H > 1$ ; and for negatively cooperative ligand binding,  $n_H < 1$ . Second, for positively cooperative ligand binding, the value of  $n_H$  puts an upper limit on the numbers of interacting binding sites for X. For example, for a tetrameric protein with four binding sites for X,  $n_H \leq 4$ .

The Hill plot was initially developed by A. V. Hill in the context of proposing a mechanism for the positive cooperativity of oxygen binding to [hemoglobin](#) (1). The central concept of his proposal was recognition that cooperativity in ligand binding requires interactions between multiple binding sites. In the limit that the cooperativity between  $n$  binding sites is very strong, then only unliganded and fully  $n$ -liganded species of the protein exist. This can be expressed by the binding reaction



An apparent association constant may be written as

$$K = \frac{[(P \cdot X)_n]}{[P_n][X]^n} \quad (2)$$

The fractional saturation,  $y$ , is then given by

$$\bar{y} = \frac{[X]^n}{(1/K) + [X]^n} \quad (3)$$

Equation 3 can be converted into a form useful for graphical analysis by taking the ratio of bound sites to free sites,  $y/(1-y)$ , and taking the logarithm

$$\log(\bar{y}/[1 - y]) = n \log[X]_f + \log K \quad (4)$$

Equation 4 is the basis for the Hill plot.

Analysis of the properties of Equation 4 as a function of  $[X]$  shows that the slope approaches +1 at both small and large values of  $y$  (2-4). These are shown as the lines in Figure 1 that are labeled lower asymptote and upper asymptote, respectively. As described below, in principle these asymptotes provide information about ligand binding. However, accurate determinations of the asymptotes in practice are difficult because they require measurements at very low and very high extents of ligand binding. For this reason, most uses of the Hill plot focus on the slope at 50% saturation with ligand X, the Hill coefficient.

The Hill coefficient can usually be determined with good precision for cooperative binding systems, and it provides an accurate phenomenological description of the cooperativity. For systems with two interacting binding sites for X (eg, dimeric proteins), it may be possible to determine apparent dissociation constants for each site from the Hill plot (see text below). For larger numbers of interacting binding sites, it may be possible to determine the dissociation constants for the first and last ligands bound and the interaction energy between these sites. However, determinations of dissociation constants for other ligand bindings steps are not in general possible. (See Ref. 3 for discussion relative to oxygen binding to hemoglobin.)

To understand the relations between cooperative ligand binding and the Hill coefficient, it is instructive to consider the case of a symmetrical dimer with one binding site for ligand X on each subunit. The stepwise occupancy of the binding sites can be written



where  $K_1$  and  $K_2$  are intrinsic stepwise association constants defined by the equations

$$K_1 = \frac{[P_2X]}{[P_2][X]}, \quad K_2 = \frac{2[P_2X_2]}{[P_2X][X]} \quad (6)$$

The Hill coefficient at  $y = 0.5$  is given by

$$n_H = \frac{2}{1 + (K_1/K_2)^{1/2}} \quad (7)$$

For  $(K_1/K_2) > 1$ ,  $n_H < 1$ , indicating negative cooperativity. For  $(K_1/K_2) < 1$ ,  $n_H > 1$ , indicating positive cooperativity. The value of  $n_H$  cannot exceed 2. The intercepts of the asymptotes in Figure 1, when  $\log [X] = 0$  are  $\log K_1$  (lower) and  $\log K_2$  (upper). Thus they provide values for the intrinsic stepwise binding constants for the initial and final ligand binding steps. Although this is a general result, experimental determination of the asymptotes is in practice difficult. An analogous treatment



of cooperativity in the kinetic behavior of a symmetric dimeric enzyme using Weber's approach of coupling free energies (4) has been presented and shows the relations between the magnitude of the coupling free energies and the value of the Hill coefficient (5). A treatment of the Hill coefficient in the context of the concerted allosteric model is presented on page 135 of Ref. 3.

### Bibliography

1. A. V. Hill (1910) *J. Physiol. (London)* **40**, iv–vii.
2. J. T. Edsall and H. Gutfreund (1983), *Biothermodynamics: The Study of Biochemical Processes at Equilibrium*, Monographs in Molecular Biophysics and Biochemistry, (H. Gutfreund, ed.), Wiley, New York, p. 248.
3. J. Wyman and S. J. Gill (1990), *Binding and Linkage: Functional Chemistry of Biological Macromolecules*, University Science Books, Mill Valley, CA.
4. G. Weber (1992), *Protein Interactions*, Chapman and Hall, New York, p. 293.
5. G. D. Reinhart (1988) *Biophys. Chem.* **30**, 159–172.

### Hirudin

Hirudin is a protein inhibitor of thrombin in the saliva of the medicinal leech *Hirudo medicinalis* [see Proteinase inhibitors, protein and Serine proteinase inhibitors]. The target enzyme here seems obvious. It is a single polypeptide chain of 65 amino acid residues crosslinked by three disulfide bridges. It combines extremely tightly ( $K_1 10^{-14} \text{M}$ ) with thrombin. The crystal structure of a thrombin-hirudin complex reveals that hirudin is not a standard-mechanism canonical inhibitor. The active site of thrombin is embedded in a canyon restricting the entry of many substrates and of most trypsin inhibitors that do not inhibit thrombin. Hirudin blocks the active site, but none of its residues embeds into the  $S_1$  cavity, a hallmark of standard-mechanism canonical inhibition. The highly flexible COOH terminus of hirudin (residues 50–65) binds to the fibrinogen-binding exosite of thrombin. Many hirudin-based drugs block this fibrinogen-binding exosite rather than the active site of thrombin.

### Suggestion for Further Reading

- M. G. Grütter et al. (1990) Crystal structure of the thrombin-hirudin complex: A novel mode of serine protease inhibition. *EMBO J.* **9**, 2361–2365.

### His Operon

Energy equivalent to about 41 ATP molecules is required to synthesize one molecule of the amino acid [histidine](#) (1). The considerable metabolic cost of histidine biosynthesis presumably accounts for the evolution of multiple strategies to regulate the rate of synthesis of the amino acid in response to environmental changes. Checkpoints regulate both the flow of intermediates through the biosynthetic pathway and the amounts of histidine-biosynthetic [enzymes](#) present. Expression of the genes for

these enzymes is regulated in bacterial cells by mechanisms that are both general (metabolic regulation, elongation control) and specific (attenuation control, segmental stabilization of the distal part of the [messenger RNA](#)).

## 1. Structural organization of the operon

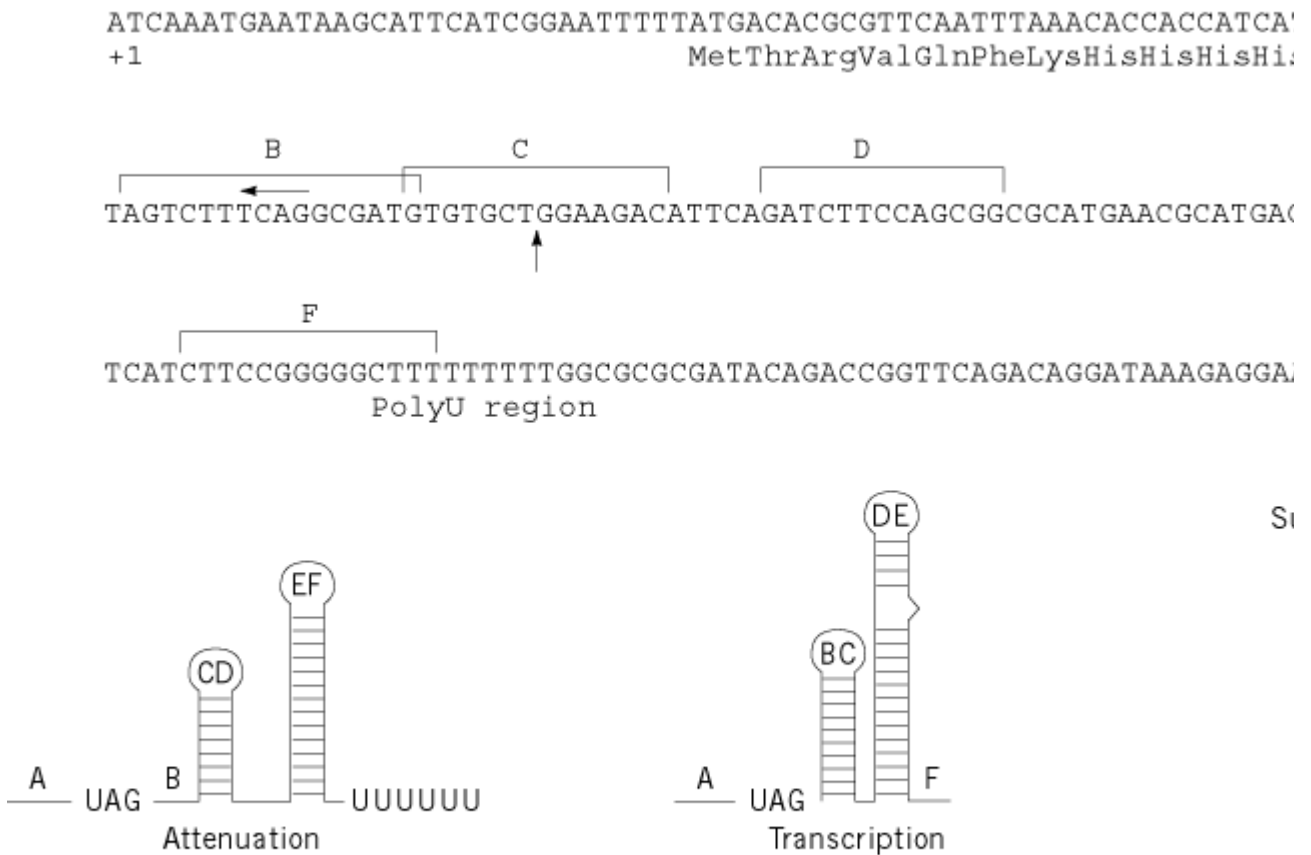
In *Escherichia coli* and *Salmonella typhimurium* the enzymes responsible for the biosynthesis of histidine are encoded by eight genes tightly clustered in a single, large [operon](#) (*his* operon). In both species, [transcription](#) produces a single polycistronic mRNA about 7300 nucleotides long, extending from a primary **promoter** (*hisP1*) to a Rho-independent terminator. Two weak internal promoters, *hisP2* and *hisP3*, are located within the *hisC* and *hisF* genes, respectively. The structural organization of the operon is essentially the same in the two species, and in both the translational [stop codon](#) of each **cistron** overlaps the translational [initiation codon](#) of the downstream cistron (2). This organization allows [ribosomes](#) to initiate the translation of a new cistron without moving away from the mRNA after terminating translation of the preceding one. Such a translational coupling mechanism probably guarantees equimolar synthesis of the corresponding gene products (3).

## 2. Control of transcription initiation and elongation

Transcription of the *his* operon is about four-fold more efficient in bacteria growing in minimal-glucose medium than when growing in rich medium. This form of control, called *metabolic regulation*, adjusts the expression of the operon to the amino acid supply in the cell. It is mediated by the “alarmone” guanosine 5'-diphosphate 3'-diphosphate (ppGpp), which is the effector of the stringent response (see [Stringency](#)). The alarmone regulates the *his* operon positively by stimulating the primary promoter *hisP1* under conditions of moderate amino acid starvation (4).

In addition to this general metabolic control, *his* operon transcription is specifically regulated by [attenuation of transcription](#), a mechanism in which a regulatory element, located upstream of the first structural gene of the cluster, modulates the level of expression of the histidine biosynthetic enzymes in response to the intracellular levels of charged histidyl-transfer RNA, His-tRNA<sup>His</sup> (see [Transfer RNA](#)) (5). The *his*-specific regulatory element is transcribed in a 180-nucleotide RNA leader, which exhibits two prominent features: (i) a 16-residue coding sequence including seven consecutive codons specifying histidine, and (ii) overlapping regions of dyad symmetry capable of folding into mutually exclusive, alternative **secondary structures** that signal either transcription termination or antitermination (6, 7). Six RNA segments are involved in base pairing (Fig. 1 A to F) and the stem-loop structure formed by the E and F RNA regions, plus the adjacent run of uridylylate residues, constitutes the attenuator, a strong Rho-independent transcription terminator (Fig 1). Translational control of *his* operon transcription is determined by ribosome occupancy of the leader RNA, which in turn depends, given the peculiar composition of the *his* leader peptide, on the availability of His-tRNA<sup>His</sup>. High levels of His-tRNA<sup>His</sup> allow rapid movement of ribosomes up to the B segment; in this case, formation of the C:D and E:F stem-loop structures will result in premature transcription termination (Fig. 1, Attenuation). In the presence of low levels of charged tRNA<sup>His</sup>, ribosomes stall at the consecutive histidine codons of the leader peptide and prevent the A:B pairing by masking the A segment. Base pairing between the B and C and between the D and E RNA regions prevents formation of the attenuator and determines the antitermination conformation (Fig. 1, Transcription). In the case of severe limitation of the intracellular pool of all charged tRNAs, translation of the leader peptide fails to initiate: under these conditions, the A:B, C:D and E:F stem-loop structures form sequentially, producing a strong transcription termination (Fig. 1, Superattenuation). **RNA polymerase** pauses after synthesis of the first RNA hairpin (A:B). This pausing is believed to synchronize transcription and translation of the leader region by halting the elongating RNA polymerase until a ribosome starts translation of the leader peptide (8). The pause hairpin (Fig. 1) is the only portion of the structure thought to form when RNA polymerase resides at the pause site.

**Figure 1.** Regulation of translation of the *his* operon messenger RNA. Top: the nucleotide sequence of the leader region *typhimurium* from the transcription initiation site (+1) to the first structural gene, *hisG*. Brackets above the nucleotide segments (A to F) capable of forming alternative, mutually exclusive secondary structures. The convergent arrows indicate structure required for transcriptional pausing at the downstream site indicated by a vertical arrow (see text). The amino peptide and of the amino-terminal region of *hisG* are shown. Bottom: Schematic representation of the different conformations (see text). The run of U's at the 3' end indicates the formation of the terminator hairpin E:F. The position of the terminator (UAG) in each RNA configuration is also indicated.



St

Because the absolute amount of charged tRNA<sup>His</sup> controls the level of *his* attenuation (5), mutants exhibiting high *his* operon expression contain defects in tRNA<sup>His</sup> biosynthesis, aminoacylation with histidine, or tRNA<sup>His</sup> modification and processing. The *hisR* gene encodes the single cellular tRNA<sup>His</sup>; and mutations in the *hisR* promoter reduce the total cellular content of tRNA<sup>His</sup> molecules by about 50% and thereby cause increased readthrough transcription of the *his* attenuator (9). The *hisS* gene encodes histidyl-[aminoacyl tRNA synthetase](#), which aminoacylates tRNA<sup>His</sup> molecules with histidine. Mutations that lower the activity of the histidyl-tRNA synthetase or decrease the enzyme's affinity for histidine, tRNA<sup>His</sup>, or ATP, affect the level of *his* attenuation by reducing the percentage of tRNA<sup>His</sup> molecules charged with histidine (10). The *hisT* gene encodes pseudouridine synthase I, which catalyzes the formation of pseudouridine residues in the [anticodon](#) region of several tRNA species, including tRNA<sup>His</sup>. Although the undermodified tRNA<sup>His</sup> molecules are charged with histidine to the same extent as in wild-type strains, transcription termination at the *his*

attenuator is greatly decreased, because the slow rate of translation of the consecutive histidine codons causes stalling of ribosomes (11).

The overall contribution of the internal promoter *hisp2* to the expression of the distal genes of the operon is negligible when transcription proceeds from *hisp1*, because *hisp2* is inhibited by transcription readthrough, a phenomenon known as *promoter occlusion* (12). *hisp2* is also subjected to metabolic regulation, although to a lesser extent than *hisp1*.

Elongation of the *his*-mRNA is modulated by a non-specific mechanism operating at the level of intracistronic *transcription termination elements* (TTEs) (13). These elements consist of cytosine-rich and guanosine-poor RNA regions and are the binding-activation sites of the transcription-termination Rho factor, which is responsible for polar effects in polycistronic operons (14). A premature arrest of translation, produced by [nonsense mutations](#), favors the binding of Rho to the TTE on the nascent transcript. The subsequent interaction of Rho with the elongating RNA polymerase causes a premature release of transcripts. Polarity results, with reduced expression of the genes located downstream from the TTE.

### 3. Decay and Segmental Stabilization of the *his*-mRNA

The primary 7300-nucleotide *his*-mRNA has a half-life of about 3 minutes in cells growing in minimal-glucose medium and is degraded in a net 5' → 3' direction. Three major processed species, 6300, 5000, and 3900 nucleotides long, encompassing the last seven, six, and five cistrons, respectively, are generated in the decay process (12). The 6300- and the 5000-nucleotide RNAs, which have half-lives of 5 and 6 minutes, respectively, have heterogeneous 5' ends generated by ribonuclease E cleavage (see [RNA Degradation In Vitro](#)). The 3900-nucleotide processed RNA species has a unique 5' end and an uncommon stability, having a half-life of about 15 minutes. This RNA species is generated by specific processing events requiring sequential cleavages by two different endonucleases. RNase E triggers the process by cleaving at a major target site located in the *hisC* cistron, 620 nucleotides upstream of the 5' end of the processed species. Subsequently, ribonuclease P cleaves the processed RNA species generated by RNase E at a site located 76 nucleotides upstream of the start codon of the *hisB* cistron, producing the mature 5' end. The observation that the RNase P-catalyzed reaction requires the presence of ribosomes suggests that translation of the *hisB* cistron might favor formation of the structure recognized by RNase P (15).

### Bibliography

1. M. Brenner and B. N. Ames (1971) In *Metabolic Pathways* (H. J. Vogel, ed.), Academic Press, Inc., N.Y., vol. 5, pp. 349–387.
2. M. S. Carlomagno, L. Chiariotti, P. Alifano, A. G. Nappo, and C. B. Bruni (1988) *J. Mol. Biol.* **203**, 585–606.
3. C. Yanofsky, T. Platt, I. P. Crawford, B. P. Nichols, G. E. Christie, H. Horowitz, M. Van Cleemput, and A. M. Wu (1981) *Nucleic Acids Res.* **9**, 6647–6668.
4. J. C. Stephens, S. W. Artz, and B. N. Ames (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4389–4393.
5. F. Blasi and C. B. Bruni (1981) *Curr. Top. Cell. Regul.* **19**, 1–45.
6. W. M. Barnes (1978) *Proc. Natl. Acad. Sci. USA* **75**, 4281–4285.
7. P. P. Di Nocera, F. Blasi, R. Di Lauro, R. Frunzio, and C. B. Bruni (1978) *Proc. Natl. Acad. Sci. USA* **75**, 4276–4280.
8. C. L. Chan and R. Landick (1989) *J. Biol. Chem.* **264**, 20796–20804.
9. H. M. Johnston, W. M. Barnes, F. G. Chumley, L. Bossi, and J. R. Roth (1980) *Proc. Natl. Acad. Sci. USA* **77**, 508–512.
10. J. A. Lewis and B. N. Ames (1972) *J. Mol. Biol.* **66**, 131–142.
11. C. F. Singer, G. R. Smith, R. Cortese, and B. N. Ames (1972) *Nature New Biol.* **238**, 72–74.

12. P. Alifano, C. Piscitelli, V. Blasi, F. Rivellini, A. G. Nappo, C. B. Bruni, and M. S. Carlomagno (1992) *Mol. Microbiol.* **6**, 787–798.
13. P. Alifano, F. Rivellini, D. Limauro, C. B. Bruni, and M. S. Carlomagno (1991) *Cell* **64**, 553–563.
14. P. Alifano, M. S. Ciampi, A. G. Nappo, C. B. Bruni, and M. S. Carlomagno (1988) *Cell* **55**, 351–360.
15. P. Alifano, F. Rivellini, C. Piscitelli, C. M. Arraiano, C. B. Bruni, and M. S. Carlomagno (1994) *Genes Dev.* **8**, 3021–3031.

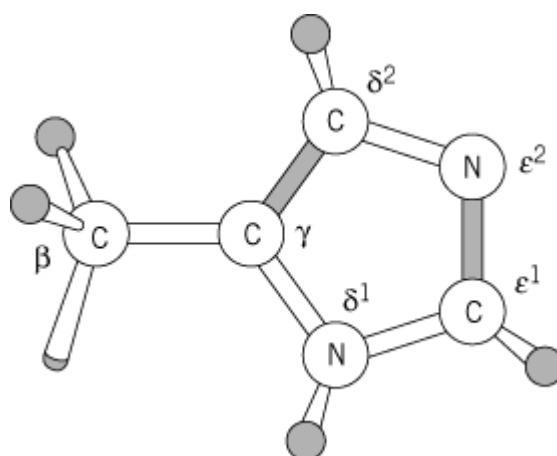
### Suggestions for Further Reading

16. M. Winkler (1996) "Biosynthesis of histidine. In" *Escherichia coli and Salmonella typhimurium*, 2nd Ed. (F. C. Neidhart, ed.), American Society for Microbiology, Washington DC, pp 485–505.
17. P. Alifano, R. Fani, P. Liò, A. Lazcano, M. Bazzicalupo, M. S. Carlomagno, and C. B. Bruni (1996) Histidine biosynthetic pathway and genes: Structure, regulation and evolution. *Microb. Rev.* **60**, 44–69.

## Histidine (His, H)

The [amino acid](#) histidine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—CAU and CAC—and represents approximately 2.2% of the residues of the proteins that have been characterized. The histidyl residue incorporated has a mass of 137.14 Da, a **van der Waals volume** of 118 Å<sup>3</sup>, and an [accessible surface](#) area of 194 Å<sup>2</sup>. His residues are changed during [divergent evolution](#) with average frequency; they are interchanged in **homologous** proteins most frequently with [asparagine](#) and [glutamine](#) residues.

The side chain of His residues consists of a b-methylene carbon and an imidazole group:



The imidazole groups possesses several special properties that make it extremely effective as a nucleophilic catalyst. It is an amine, which is much more reactive than hydroxide ion in terms of basicity. Furthermore, it is a tertiary amine, which is intrinsically more nucleophilic than primary or secondary amines. The enhanced reactivity of tertiary amines is usually canceled by their greater steric hindrance, but in imidazole the atoms bonded to the two nitrogen atoms are held back in a

five-membered ring and cause relatively little steric hindrance. Imidazole has a  $pK_a$  value near 7, so it is one of the strongest bases that can exist at neutral pH. A weaker base would have a lower nucleophilic reactivity, whereas a stronger base would be protonated to a greater extent at neutral pH and would be correspondingly less reactive.

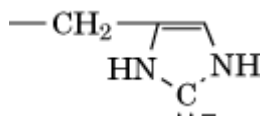
In the nonionized form of the imidazole ring, the nitrogen atom with the hydrogen atom is an electrophile and a donor for **hydrogen bonding**, while the other nitrogen atom is a nucleophile and a hydrogen bond acceptor. Consequently, this one side chain is extremely versatile, almost the chemical equivalent of being ambidextrous. Not surprisingly, His residues are often found at the [active sites](#) of enzymes and involved directly in catalysis.

The two nitrogen atoms of the His side chain are designated here as d1 and  $\epsilon$ 2, but they are also known, respectively, as p and q or as N-1 and N-3. The latter designation is often ambiguous, because biochemists usually assign the number 1 to the nitrogen atom adjacent to the  $C_\beta$ , whereas organic chemists tend to reverse the numbers. The nonionized imidazole ring can exist as two [tautomers](#), with the hydrogen atom on either the d1 and  $\epsilon$ 2 nitrogen atom. In model peptides, the hydrogen atom is usually predominantly on the  $\epsilon$ 2 nitrogen atom, which has a  $pK_a$  value about 0.6 pH unit higher than that of the d1 atom. The relative affinities of the two nitrogen atoms for protons can vary with their local environment, however, and both forms are found in folded proteins.

The nonionized His side chain is readily protonated, with a  $pK_a$  value near 7 at the second N atom, which destroys its nucleophilicity. The positive charge is shared by the two N atoms by resonance:



The  $C_\epsilon^1$  atom, between the two N atoms, is observed to exchange its hydrogen atom slowly with aqueous solvents, indicating that it has a very small probability of being deprotonated:



This [hydrogen exchange](#) reaction provides a useful probe of the environments of His residues in proteins. His residues are especially useful in  $^1\text{H-NMR}$  studies of proteins because the hydrogen atom on the  $C_\epsilon^1$  atom is usually well resolved from the multitude of resonances of the other hydrogen atoms in proteins. Its resonance is usually shifted by about 1 ppm to a lower field strength on protonation of the side chain, often making it relatively easy to determine the  $pK_a$  values of individual His residues, even in large proteins.

His is one of the residues least favoring the *alpha*-**helical** conformation in model peptides, but it occurs in  $\alpha$ -helices in native [protein structures](#) slightly more frequently than in *beta*-**sheets** or reverse [turns](#).

The imidazole group is, in principle, capable of undergoing numerous chemical modification reactions, but most of these occur much more readily with **amino** and [thiol groups](#), so very few are suitable for modifying His residues specifically. The imidazole nitrogen atoms have an intrinsic affinity for metal ions, especially zinc, iron, and copper. His residues are frequently involved in binding such ligands to proteins, as in [metal-requiring enzymes](#), **zinc fingers**, [myoglobin](#), and [hemoglobin](#).

Suggestions for Further Reading

E. A. Barnard and W. D. Stein (1959) The roles of imidazole in biological systems, *Adv. Enzymol.* **20**, 51–110.

M. Tanokura (1983) <sup>1</sup>H-NMR study on the tautomerism of the imidazole ring of histidine residues. I. Microscopic pK values and molar ratios of tautomers in histidine-containing peptides, *Biochim. Biophys. Acta* **742**, 576–585.

## Histocompatibility

The transplantation of tissue from one individual to another most often results in an adverse immune reaction between the host and the transplanted tissue, leading to graft rejection. The inclusion of lymphoid effector cells in a tissue graft can lead to graft-versus-host disease as well. These phenomena stem primarily from discordant expression of MHC cell surface molecules encoded by **polymorphic** class I and class II loci of the [major histocompatibility complex](#) (MHC). Graft acceptance or rejection is most strongly associated with identity or difference at class I loci, followed by class II and so-called minor histocompatibility antigens, as described below. The clinical evaluation of donor–host compatibility involves “tissue typing” individuals to determine the alleles present at the HLA-A, -B, and -DR loci. In practice, perfect matches at these loci are rare (requiring identity at as many as six independent loci, because **gene** products from maternal and paternal [chromosomes](#) are co-dominantly expressed); the majority of transplants rely mainly on the use of immunosuppressive agents postoperatively to increase graft survival. Host individuals are additionally screened for the presence of serum **antibodies** to MHC molecules, arising from prior exposure to non-self MHC through a blood transfusion or pregnancy.

Allograft rejection involves a variety of phenomena; the discussion here is confined to mismatched donor–host class I alleles (whose gene products have often been referred to as “classical transplantation antigens”) as the basis for tissue incompatibility. MHC class I molecules present peptide **antigens** to the clonally expressed [T-cell receptors](#) (TCRs) of [cytotoxic T lymphocytes](#) (CTLs). The composition of the peptides presented by a given class I molecule is influenced by allomorph-specific peptide binding preferences, which in turn are a function of polymorphic amino acid residues found in the peptide-binding region of the class I molecule. Under typical circumstances, CTLs are rendered tolerant to the self class I/self peptide complexes expressed in an individual. The presence of an allograft introduces a new cohort of class I/peptide complexes into the host system, and the array of peptides displayed on the cells of the graft will reflect the binding preferences of the donor's class I allomorphs. Among these allo-MHC/peptide complexes are potentially numerous **epitopes** to which the host individual is not tolerant, thereby marking cells of the transplanted tissue for destruction by circulating CTL.

The cellular immune response to non-self class I/peptide complexes is termed “alloreactivity.” The molecular basis for TCR recognition of allo-MHC molecules has long been an area of intensive study and debate in immunology. The preceding paragraph implies a scheme in which allorecognition is a peptide-specific phenomenon and fundamentally similar to the recognition of self class I+foreign peptide. In this view, conserved regions of the majority of class I molecules are themselves likely to provide similar molecular surfaces for contact with the TCR, regardless of the class I allomorph being considered. Specific interactions thus occur between variable regions of the TCR and the class I-bound peptide. There is now considerable evidence to support this model of allorecognition (1). However, allotypic residues on the surface of a class I molecule may also provide critical determinants for allorecognition by some TCR (2), and it is possible that both modes of recognition are utilized in the polyclonal T-cell response to an allograft.

The peptide-based nature of allorecognition also accounts for the involvement of “minor histocompatibility antigens” in graft rejection. Even when genetically nonidentical donor and host are HLA matched, they also manifest differences stemming from the expression of variants of several other polymorphic proteins. A well-known example is the H-Y antigen, encoded on the [Y-Chromosome](#) of males, which presents a barrier to successful male-to-female transplantation. This protein is expressed in nearly all tissues and differs in amino acid sequence from its [X-chromosome](#) homologue (present in both males and females) (3). Consequently, some H-Y antigen-derived peptides that are presented by MHC class I molecules (3, 4) will constitute novel epitopes in a female host receiving tissue from a male donor, provoking a cellular immune response. The total number of human minor histocompatibility loci and the identity of the proteins they encode remain as yet undetermined.

## Bibliography

1. L. A. Sherman and S. Chattopadhyay (1993) *Annu. Rev. Immunol.* **11**, 385–402.
2. R. Brock, K. H. Wiesmuller, G. Jung, and P. Walden (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13108–13113.
3. W. Wang et al. (1995) *Science* **269**, 1588–1590.
4. O. Rotzschke, K. Falk, H. J. Wallny, S. Faath, and H. G. Rammensee (1990) *Science* **249**, 283–287.

## Histone Acetylation

Acetylation of [chromatin](#) by histone acetyltransferases (HATs), using [acetyl coenzyme A](#) as the acetyl donor, occurs at the [e-amino groups](#) of particular [lysine](#) side chains in the *N*-terminal basic tails of the core [histones](#), with varying consequences. Histone deacetylases (HDACs) remove the acetylation. Histone acetylation has roles in chromatin assembly linked to [DNA replication](#), in gene [transcription](#), where the histones are hyperacetylated, and, apparently paradoxically, in [heterochromatin](#) (1). In the centric heterochromatin of the [polytene chromosomes](#) of the fruit fly *Drosophila* (2) and at the silent mating type loci in the yeast *Saccharomyces cerevisiae* (3), the histones are hypoacetylated, as might be expected, but an acetyl group on Lys12 in the exposed *N*-terminal tail of H4 (see [Nucleosome](#)) is essential for heterochromatin formation and may act as a recognition site for the assembly of other heterochromatin proteins. The importance of the pattern of site-specific acetylation, rather than the overall number of acetyl groups, is becoming increasingly clear, but is still poorly understood. Elevated levels of Lys16 acetylation on H4 in the male [X-chromosome](#) in *Drosophila* appears to be important, in some way, in **dosage compensation**, which results in a doubling of transcription from X-linked genes in the male (4).

Histone H4 deposited with H3 at the **replication fork** by the core histone **molecular chaperone** CAF-1 (chromatin assembly factor 1) during chromatin assembly is acetylated at Lys5 and Lys12 primarily, by a so-called B-type (deposition-linked) acetyltransferase. Deacetylation occurs soon after deposition. Transcription-related hyperacetylation occurs in chromatin at lysines 5, 8, 12, and 26 in H4, and lysines 9 and 14 in H3. Its role appears to be to “open” chromatin, primarily by destabilizing higher-order structure, thus permitting access by the transcription machinery. Suggestions of the interactions that might be disrupted have come recently from the high-resolution structure of the nucleosome core particle. Because the *N*-terminal tails are not bound within the nucleosome core, acetylation is unlikely to act at this level. The facilitation of **transcription factor** binding due to acetylation observed in *in vitro* assays with mononucleosomes (eg, Ref. 5) may be



due to disruption of nonspecific interactions of the histone tails with the DNA gyres, in the absence of any possibilities for higher-order structure formation. Direct evidence that acetylation *per se* does result in increased transcription has come from studies of the *hsp26* gene reconstituted using a *Drosophila* embryo extract into chromatin with either unacetylated or acetylated histones (6). In nuclei, large domains of chromatin-encompassing active genes appear to be acetylated; this state precedes transcription, and the chromatin is said to be in a transcriptionally competent state. In at least one well-documented case (7), the chicken  $\beta$ -globin gene cluster, the genes are embedded in a domain of greater than 30 kbp, where the histones are hyperacetylated and which shows increased sensitivity to the **nuclease** DNase I (see [Chromatin](#)), suggesting a less compact structure.

## 1. Histone Acetylases

The finding that the *Tetrahymena* HAT, HAT-A, is **homologous** to a yeast transcriptional co-activator, GCN5 (8), solved the mystery of how acetylation might be targeted to particular regions of chromatin. It has since been shown that the acetyltransferase activity of GCN5 (which functions as the GCN5-Ada2-Ada3 complex) is essential for its activator activity (9, 10) and that acetylation occurs in the **promoter** region. GCN5 also functions with Ada proteins in the yeast SAGA complex, which also contains Spt proteins (transcriptional regulators) (SAGA = Spt3-Ada-GCN5-Acetyltransferase) and is involved in the regulation of many yeast genes. It is now clear (1) that a growing number of proteins known to be transcriptional co-activators and co-activator-associated proteins have HAT activity [eg, p300/CBP, P/CAF (a human homologue of GCN5), the TFIID component TAF<sub>II</sub>250, which is likely to be associated with all promoters during transcriptional initiation, and many others]. (CBP = cAMP response element binding protein; CAF = CBP-associated factor.) Most, including GCN5, function within large complexes, where one subunit recognizes a DNA-bound transcription factor and another has acetyltransferase activity. p300/CBP interacts with transcription factors in many different signaling pathways (eg, nuclear [hormone receptors](#), c-Fos, c-Jun/v-Jun), and acetylation is therefore the end point of many [signal transduction](#) pathways. Acetylation of different histones, and different sites within a particular histone, may turn out to require different HATs; GCN5 *in vivo* acetylates preferentially Lys8 and Lys16 of H4 and Lys14 of H3 (11). Histones may not be the only targets of all of the enzymes designated as histone acetyltransferases, for example, HMG 1 and HMG I(Y) are known to be acetylated *in vivo*. The transcription factor/tumor suppressor [p53](#) is acetylated both *in vitro* and *in vivo* by p300, its coactivator. Two general transcription initiation factors, TFIIEb and TFIIF, can also be acetylated by P/CAF and p300 (12). Because the HAT activities that have recently been studied are tethered to sequence-specific DNA binding proteins at promoters (and [enhancers](#)), their effect will necessarily be local. It is not clear how this relates to whole domain acetylation over tens of kilobases. One possibility is that the loosening/disruption of the chromatin structure close to the promoter might be transmitted down the 30-nm filament (see [Chromatin](#)) through disruption of cooperative interactions, and that the temporarily loosened tails become targets for a pool of untargeted acetylases. There are other possibilities, and this area needs further study. Various models for how histone acetylases and deacetylases selectively affect gene expression have recently been discussed (13).

The first structure of a HAT catalytic subunit (Hat1 from the yeast *Saccharomyces cerevisiae*, which was the first HAT to be isolated) has recently been determined, complexed with acetyl CoA (14); *in vivo* Hat1 exists as a heterodimer with the regulatory subunit HAT2. Hat1's sequence preference for acetylation is Lys12>Lys5, and it acetylates free rather than nucleosomal histones, so it seems likely to be a deposition-related HAT. Nonetheless, Hat1 in both yeast and human cells appears to be primarily nuclear, rather than cytoplasmic as is generally assumed (see Refs. 14 and 15). Hat1 has sequence motifs in common with the GCN5-related superfamily of *N*-acetyltransferases (the GNAT superfamily) and should therefore be an excellent model for the other HATs. Indeed the Hat1 structure explains mutagenesis data on other HATs, such as Gcn5 and CBP; some of the mutations map to the acetyl CoA binding site. The structure of a HAT with a histone or histone tail bound is now needed to reveal the basis of the site specificity of acetylation, which is in the meantime suggested by modeling (14); the structure of a Hat1/Hat2 complex would also be of considerable

interest. Hat2 is similar to the p48 subunit of CAF-1 (see [Chromatin](#)) and to similar subunits in histone deacetylase HD1 in humans and a related deacetylase in *Drosophila*. These related proteins may be histone-binding subunits.

## 2. Histone Deacetylases

Histone deacetylases (HDACs) are enzymes that remove acetyl groups, which generally results in repression of transcription. They also appear to be targeted to particular sites by gene regulatory proteins, in this case [repressors](#), in organisms from yeast to mammals. (Intriguingly, the yeast deacetylase RPD3, which is also a repressor of **silencing** at [telomeres](#) and mating type loci, specifically deacetylates Lys5 and Lys12 in H4 *in vitro*, consistent with the requirement for Lys12 in an acetylated form for heterochromatin.) Like the HATs, HDACs function in the context of large complexes in which interaction of the HDAC and repressor appears to be mediated by [corepressors](#). In yeast the deacetylase RPD3 is complexed with the corepressor sin3. The mammalian homologue, mSin3, has four paired [amphipathic](#) helix (PAH) domains that can interact with several repressors and corepressors. One is the Mad/Max repressor (Mad antagonizes the transcriptional activation and transformation functions of the **Myc oncoprotein**); another is N-CoR, a repressor for the [thyroid hormone](#) receptor, which exists in a complex with mSin3A and the mammalian homologue of RPD3. Progress in this rapidly evolving field has been succinctly summarized ([1](#), [16](#)). The finding that the [5-methylcytosine](#) binding protein, MeCP2, recruits Sin3 ([17](#), [18](#)) and a deacetylase suggests a way of ensuring that chromatin in methylated (inactive) regions of the genome remains unacetylated and in a stable higher-order structure.

## Bibliography

1. M. Grunstein (1997) *Nature* **389**, 349–352.
2. B. M. Turner, A. J. Birley, and J. Lavender (1992) *Cell* **69**, 375–384.
3. M. Braunstein et al. (1996) *Mol. Cell. Biol.* **16**, 4349–4356.
4. J. R. Bone et al. (1994) *Genes Dev.* **8**, 96–104.
5. M. Vettese-Dadey et al. (1996). *EMBO J.* **15**, 2508–2518.
6. K. P. Nightingale, R. E. Wellinger, J. M. Sogo, and P. B. Becker (1998) *EMBO J.* **17**, 2865–2876.
7. T. R. Hebbes, A. W. Thorne, and C. Crane-Robinson (1988) *EMBO J.* **7**, 1395–1402.
8. J. E. Brownell et al. (1996) *Cell* **84**, 843–851.
9. M.-H. Kuo et al. (1998) *Genes Dev.* **12**, 627–639.
10. L. Wang, L. Liu, and S. Berger (1998) *Genes Dev.* **12**, 640–653.
11. M.-H. Kuo et al. (1996) *Nature* **383**, 269–272.
12. A. Imhof et al. (1997) *Curr. Biol.* **7**, 689–692.
13. K. Struhl (1998) *Genes Dev.* **12**, 599–606.
14. R. N. Dutnall, S. T. Tafrov, R. Sternglanz, and V. Ramakrishnan (1998) *Cell* **94**, 427–438.
15. A. Verreault, P. D. Kaufman, R. Kobayashi, and B. Stillman (1997) *Curr. Biol.* **8**, 96–108.
16. M. J. Pazin and J. T. Kadonaga (1997) *Cell* **89**, 325–328.
17. X. Nan et al. (1998) *Nature* **393**, 386–389.
18. P. L. Jones et al. (1998) *Nat. Genet.* **19**, 187–191.

## Suggestions for Further Reading

19. J. E. Brownell and C. D. Allis (1996) Special HATs for special occasions. Linking histone acetylation to chromatin assembly and gene activation. *Curr. Opin. Genet. Dev.* **6**, 176–184.
20. P. A. Grant et al. (1998) The SAGA unfolds: convergence of transcription regulators in chromatin-modifying complexes. *Trends Cell Biol.* **8**, 193–197.
21. G. A. Hartzog and F. Winston (1997) Nucleosomes and transcription: recent lessons from

- genetics. *Curr. Opin. Genet. Dev.* **7**, 192–198.
22. P. Kaufman (1996) Nucleosome assembly: the CAF and the HAT. *Curr. Opin. Cell Biol.* **8**, 369–373.
  23. S. Y. Roth and C. D. Allis (1996) Histone acetylation and chromatin assembly: a single escort, multiple dances? *Cell* **87**, 5–8.
  24. B. M. Turner (1998) Histone acetylation as an epigenetic determinant of long-term transcriptional competence. *Cell. Mol. Life Sci.* **54**, 21–31.
  25. P. A. Wade, D. Pruss, and A. P. Wolffe (1997) Histone acetylation: chromatin in action. *Trends Biochem. Sci.* **22**, 128–132.

## Histone Fold

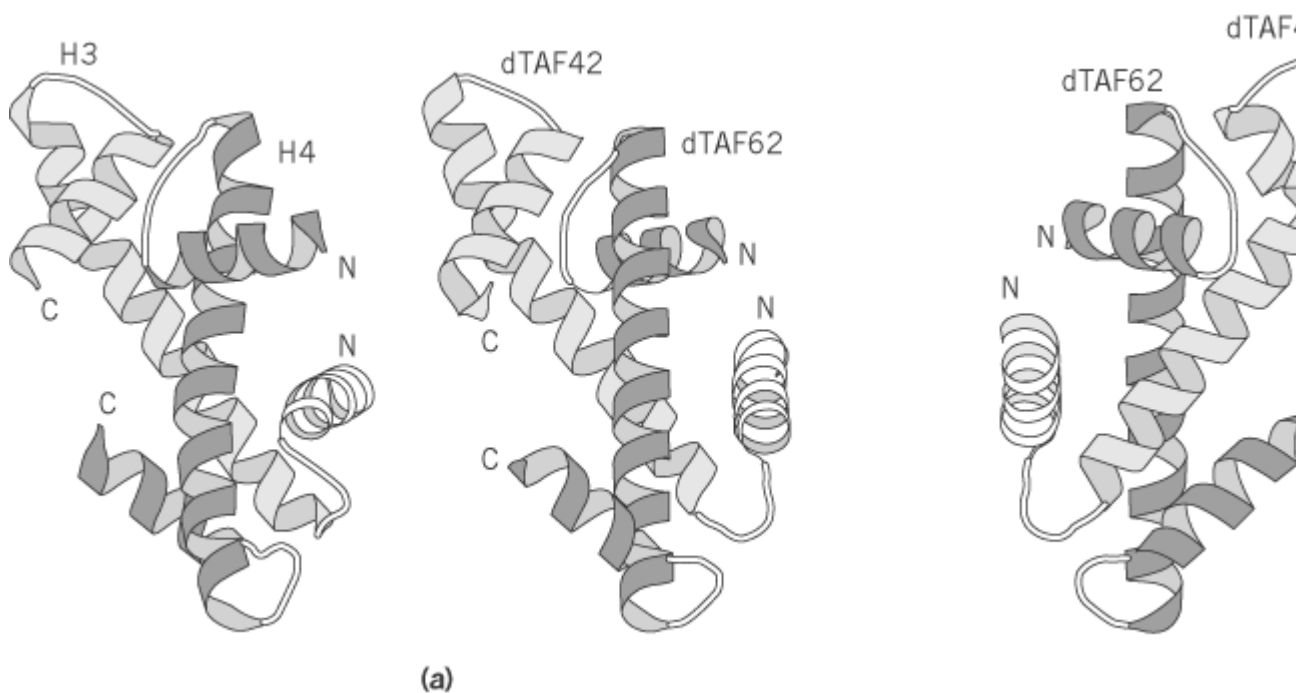
The histone-fold motif is a [protein structure](#) feature common to all four [histones](#) that was first revealed in the structure of the core histone octamer (1). The mode of interaction of this motif with DNA, suggested by the octamer structure, is shown in detail in the structure of the [nucleosome](#) core particle (2). The histone fold has since been recognized in the structures of several other proteins, spanning the evolutionary range from archaeobacteria to mammals, that are involved, like the histones, in **DNA:protein recognition** and [protein–protein interactions](#). Some of the histone-fold proteins thus far identified show little or no sequence [homology](#) with the histones (3).

The histone fold consists of a long central [a-helix](#) connected by loops to two shorter a-helices. It may have arisen by [gene duplication](#) of a **helix–strand–helix motif**, with helix fusion resulting in the long central a-helix (4). The histones in the nucleosome core particle occur as heterodimers through extensive interactions between the two histone folds in a “handshake motif,” with extensive **hydrophobic** interactions between the [nonpolar](#) faces of the two long amphipathic a-helices. The same basic fold occurs in the homodimeric archaeal histone HMfB (from the hyperthermophile *Methanothermus fervidus*) (5), which shows some sequence similarity to the histones, and in the presence of DNA forms tetramers similar to the H3<sub>2</sub>H4<sub>2</sub> tetramer (see [Histones](#)). In HMfB, however, there are no regions similar to the extensions to the basic fold in H3 that contribute to the organization of the DNA in the nucleosome core particle.

Strikingly, the histone fold has also been found to occur in other proteins. The first reports were for **TATA box**-binding protein-associated factors (TAFs), components of the [transcription](#) initiation factor TFIID. The 2 Å [X-ray crystallography](#) structure of the complex of TAF(II)42 and TAF(II)62 from *Drosophila* showed that it was a tetramer (a dimer of heterodimers), with the interacting regions (which constitute a relatively small part of the total protein) organized in histone folds similar to those of H3 and H4 (6) (Fig. 1). The finding that human TAF(II)20 (homologous to *Drosophila* TAF(II)30a) probably also contains a histone-fold motif similar to that in H2B (7) (no H2A homologues have been found) led to the suggestion that there is a histone octamer-like TAF complex within human TFIID, composed of a tetramer of TAF(II)80 and TAF(II)31 and two homodimers of TAF(II)20 (in *Drosophila* a tetramer of TAF(II)62 and TAF(II)42) and two homodimers of TAF(II)30a (6, 7). Surprisingly, three histone-like TAFs also occur amongst the 20 subunits of the P/CAF acetyltransferase complex (see [Histone Acetylation](#)), namely, TAF(II)31, TAF(II)20 / 15 and PAF65a (similar to human TAF(II)80), which resemble H3, H2B, and H4, respectively, and an octamer-like structure has again been suggested (8). The yeast 1.8-MDa SAGA (*Spt-Ada-GCN5*-acetyltransferase) acetyltransferase complex, which is similar to the human P/CAF complex, also contains three histone-fold TAFs with homology to H2B, H3, and H4 as integral

components (9), and the Spt3 protein appears to have histone folds in its *N*- and *C*-terminal regions, which have been suggested to interact intramolecularly (10). X-ray crystallography has recently shown that yet another pair of TAFs, human TAF(II)18 and TAF(II)28, which show only weak homology to histones, also form a heterodimer through a histone fold (9). The histone-fold motif thus appears to be widespread, and it may simply be a robust structural motif particularly well-suited to protein dimerization in large complexes. It remains to be seen, however, whether histone octamer-like TAF structures have a particular role to play in processes involving histones. Despite the presence of the histone fold, any interaction of TAFs with DNA is likely, *a priori*, to be different from that in the nucleosome core because the arginine side chains of the core histones that insert into the minor groove of the surrounding DNA are absent; in the case of *Drosophila* TAF(II)42 and TAF(II)62, they are replaced by serine and glutamate. The **four-helix-bundle** interface and the orientation of the two heterodimers within the TAFII tetramer are also different from those in the H<sub>3</sub>H<sub>4</sub><sub>2</sub> tetramer, although this in itself might not necessarily preclude DNA binding.

**Figure 1.** Drawings of the histone folds in H3 and H4, and *Drosophila* TAF(II)42 and TAF(II)62. (a) The TAF(II)42/TAF(II)62 heterodimer from which the H<sub>3</sub>H<sub>4</sub><sub>2</sub> tetramer in chromatin is formed; (b) the tetramer TAF(II)42<sub>2</sub>.TAF(II)62<sub>2</sub> (residues TAF62). (Adapted from Ref. 6 and reproduced with permission.)



Although no high-resolution structural data are yet available, it seems likely, on the basis of sequence homology with H2A and H2B, that the histone-fold motif also occurs in two subunits of the trimeric mammalian **transcription factor** CBF [CCAAT-binding factor (11)] and its yeast homologue HAP (CCAAT boxes are widespread promoter elements; see **CAAT Box**) and in the human protein NC2 [negative cofactor 2 (12); also known as DRAP1] and its yeast homologue, which repress transcription of class II genes by binding to TBP and to DNA and inhibiting preinitiation complex formation by preventing TFIIB binding. It is striking that the histone fold already present in the *Archaea* has been adapted not only for DNA packaging in eukaryotes but also apparently as a dimerization motif in gene regulatory proteins, which sometimes show little, if any, sequence similarity. This is reminiscent of the “winged helix” DNA-binding motif in the globular domain of histone H5 (see **Histones**), which is also found in bacterial **cyclic AMP receptor protein** and the liver transcription factor HNF3g, a member of a large family of transcription factors (the

HNF/forkhead family) from several organisms.

## Bibliography

1. G. Arents et al. (1991) Proc. Natl. Acad. Sci. USA **88**, 10148–10152.
2. K. Luger et al. (1997) Nature **389**, 251–260.
3. A. D. Baxevanis and D. Landsman (1998) Nucleic Acids Res. **26**, 372–375.
4. G. Arents and E. N. Moudrianakis (1995) Proc. Natl. Acad. Sci. USA **92**, 11170–11174.
5. M. R. Starich, K. Sandman, J. N. Reeve, and M. F. Summers (1996) J. Mol. Biol., **255**, 187–203.
6. X. Xie et al. (1996) Nature **380**, 316–322.
7. A. Hoffmann et al. (1996) Nature **380**, 356–359.
8. V. V. Ogryzko et al. (1998) Cell **94**, 35–44.
9. P. A. Grant et al. (1998) Cell **94**, 45–53.
10. C. Birck et al. (1998) Cell **94**, 239–249.
11. S. N. Maity and B. deCrombrughe (1998) Trends Biochem. Sci. **23**, 174–178.
12. A. Goppelt, G. Stelzer, F. Lottspeich, and M. Meisterernst (1996) EMBO J. **15**, 3105–3116.

## Suggestions for Further Reading

13. S. K. Burley, X. Xie, K. L. Clark, and F. Shu (1997) Histone-like transcription factors in eukaryotes. Curr. Opin. Struct. Biol. **7**, 94–102.
14. A. Hoffmann, T. Oelgeschlager, and R. G. Roeder (1997) Considerations of transcriptional control mechanisms: Do TFIID-core promoter complexes recapitulate nucleosome-like function? Proc. Natl. Acad. Sci. USA, **94**, 8928–8935.
15. K. Struhl and Z. Moqtaderi (1988) The TAFs in the HAT, Cell **94**, 1–4.

## Histones

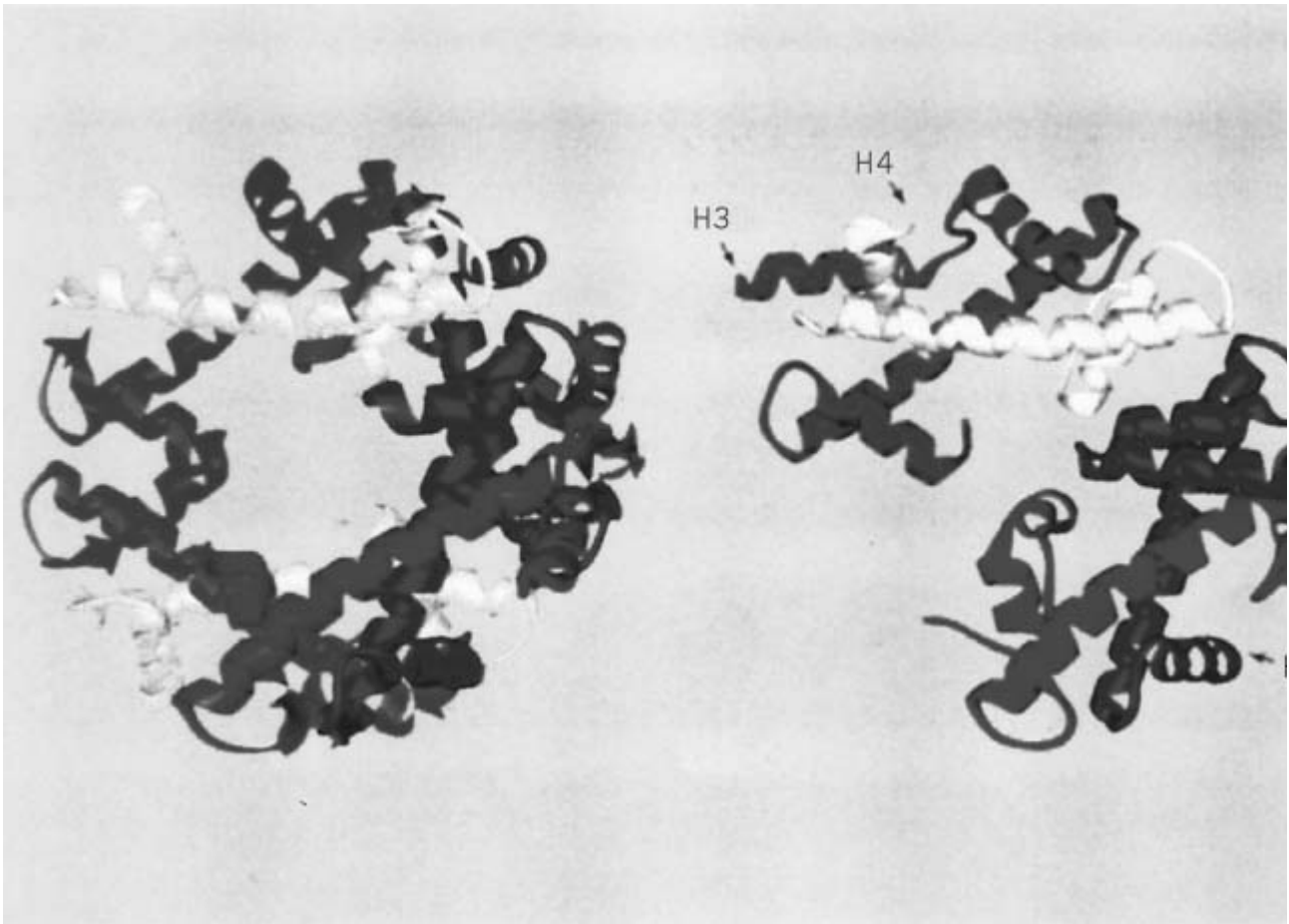
Histones are small basic (positively charged) [proteins](#) that are complexed with DNA in the [nucleus](#) of almost all eukaryotic cells (not dinoflagellates or mammalian spermatozoa) and serve to package the DNA into [nucleosomes](#) in [chromatin](#). It has recently become clear that they also play an important regulatory role in gene [transcription](#) and repression. Together, the histones occur in a mass roughly equal to that of DNA. There are five histone types, falling into two functional classes: the core histones (H3, H4, H2A, and H2B) and the larger linker histones (H1 and its variants). Two copies of each of the core histones form the histone octamer that constitutes the protein core of the nucleosome, around which are wound two complete left-handed superhelical turns of DNA. The nucleosome is stabilized by the binding of one molecule of the histone H1 class, which also plays a role in stabilization of the higher-order structure of chromatin. It is becoming clear that H1 may also play a specific role in the regulation of some genes. Both the core histones and H1 undergo dynamic reversible modification. Acetylation of the core histones plays an important role in processes as diverse as transcription, nucleosome assembly, and [heterochromatin](#) formation (see [Histone Acetylation](#)). Histone H1 is **phosphorylated** during S-phase and mitosis in the [cell cycle](#), and its extreme cell-type specific variants in avian erythrocytes and sea urchin sperm are phosphorylated up until the final stages of chromatin compaction. Two volumes contain a wealth of information about the histones ([1](#), [2](#)).

## 1. Core Histones: The Histone Octamer

The core histones [molecular masses ~11,300Da (H4) to 15,300 (H3)] show a high degree of evolutionary sequence conservation, consistent with their role as organizing structural components in a common mode of DNA packaging in eukaryotes. Histones H3 and H4 are the most conserved proteins known, with only two (conservative) amino acid differences out of 102 residues between pea and calf H4. All the core histones have high contents of the basic amino acids [lysine](#) and [arginine](#) (whose side chains are protonated, and hence positively charged, at neutral pH). H2A and H2B are relatively lysine-rich, and H3 and H4 are relatively arginine-rich. The basic residues are unevenly distributed along the amino acid sequence and are concentrated toward the *N*-terminal regions of about 25 to 40 residues. These regions (the *N*-terminal “tails”) are implicated in stabilization of chromatin higher-order structure and are essential for a number of vital functions. They are also the sites of several [post-translational modifications](#), including acetylation, which has been well-studied (see [Histone Acetylation](#)).

The four core histones exist in chromatin as an octamer that is structurally organized as a tetramer of H3 and H4 (H<sub>3</sub><sub>2</sub>H<sub>4</sub><sub>2</sub>) and two dimers of H2A and H2B (H2A.H2B). The octamer may be isolated as an entity from chromatin at high ionic strength (eg, 2 M NaCl), which disrupts the [electrostatic interactions](#) between the histones and DNA and, in addition, preserves the association of tetramers and dimers (3). At physiological ionic strength, in the absence of DNA, the isolated octamer falls apart into tetramers and dimers. The structure of the octamer, determined to 3.1 Å resolution by [X-ray crystallography](#) (4), reveals two distinctive features of the histone associations (Fig. 1). First, the four histones (which have no obvious sequence [homology](#)) share a common fold, now designated the “[histone fold](#).” This comprises a long central [α-helix](#) and two shorter flanking α-helices connected by loops. Second, the histone pairs H3, H4 and H2A, H2B share a common mode of association as heterodimers in a “handshake motif,” creating a crescent-shaped structure. Two H3.H4 heterodimers interact through α-helical regions on H3 to form the rather flat tetramer, which resembles a twisted horseshoe; H2A.H2B dimers flank this on each face, through interactions between α-helical regions on H2B and H4. The octamer crystal structure confirmed the presence of a left-handed helical ramp (for DNA binding) on the octamer surface, which had been inferred from earlier work (5), and suggested how the histone fold might form DNA contacts. The details of the interaction of the octamer with DNA have since been observed directly in the high-resolution structure of the nucleosome core particle (see [Nucleosome](#)), which also reveals the positions of the histone tails that were not visible in the structure of the octamer alone, because they are disordered. Intriguingly, the histone fold has since been found embedded in much larger proteins, notably several TAFs ([TATA-box](#) binding protein-associated factors) and [transcription factors](#) (see [Histone Fold](#)).

**Figure 1.** The structure of the histone octamer at 3.1 Å resolution, showing the histone folds (4).



The *N*-terminal histone tails of the core histones are well-conserved, despite being external to the structured core of the octamer, suggesting some structural or functional role. The tails do not appear to be required for nucleosome core integrity and may be removed by **proteolysis** without causing structural disruption. There is good evidence that they play a role, together with the basic tails of H1, in stabilization of chromatin higher-order structure. In the structure of the nucleosome core particle (6) (see [Nucleosome](#)) they form few, if any, contacts within the nucleosome core, but are extended from it and are well-placed to interact with neighboring nucleosomes or other chromosomal components. The role of the tails has been extensively studied genetically in yeast (7) (see [Nucleosome](#)). It is increasingly recognized that the tails may function as recognition sites for protein factors important for chromatin function. For example, particular sequences in the *N*-terminal tails of H4 (residues 16 to 29) and H3 (residues 4 to 20) are required, together with other proteins (Rap1, Sir2, Sir3, and Sir4) for transcriptional **silencing** (through formation of a more repressed chromatin structure) at yeast [telomeres](#) and silent mating type loci (see [Chromatin](#)). The amino-termini of H3 and H4 are also required for a1–a2 repression in yeast (a regulatory mechanism involved in distinguishing **diploid** from **haploid** yeast cells) (8) and for the action *in vitro* of at least one chromatin remodeling machine, *Drosophila* NURF (9).

The *N*-terminal tails are the sites of several types of post-translational modification, the best studied of which is acetylation (see [Histone Acetylation](#)). It is involved in transcription and chromatin assembly, the acetylation patterns being different in the two cases, and could also provide additional markers for specific recognition by various partner proteins. The specific acetylation of Lys12 only in heterochromatin in *Drosophila* and yeast could perform such a role. Transcription-linked hyperacetylation, primarily on H4 and H3, is believed to result in relaxation of the repressive higher-order structure by disruption of internucleosome contacts. Crystal contacts in the nucleosome core particle suggest contacts that might be involved (6) (see [Nucleosome](#)). Other post-translational

modifications have been less well studied and are less well understood (1, 2). Phosphorylation occurs constitutively at the *N*-terminal serine of H2A and, in a small proportion of histones, at Ser 10 and Ser 28 of H3 in response to mitogen stimulation in mammalian cells (10). Phosphorylation and acetylation appear to coexist in the same nucleosomes after mitogen stimulation and might act synergistically to disrupt internucleosomal contacts involving H3 tails and to facilitate transcription. The role of other modifications of the core histones (methylation of lysine side chains, and—in the case of H2B—**ADP-ribosylation**) is unclear; ADP-ribosylation may have a role in DNA repair, probably by disruption of chromatin structure. Also unclear is the role of ubiquitination of lysine side chains near the *C*-terminus of H2A and H2B (ie, covalent attachment of the small protein [ubiquitin](#) through an isopeptide linkage), but it appears not to be related to the well-recognized role of ubiquitination as a signal for **protein degradation**.

Core histone variants, the products of distinct genes, are found in all organisms. Their precise role is not clear, but they presumably allow fine-tuning of nucleosome structure and stability for particular purposes. A well-documented example is the developmental pattern of expression of five H2A variants and four H2B variants, which changes through the cleavage, blastula, and gastrula stages of sea urchin embryogenesis (11). Significantly, there is no variation in H3 and H4, which form the structural core of the nucleosome. During spermatogenesis in the same organism, there is a global replacement of somatic H2B with the larger and more basic sperm-specific variant spH2B, which has a long *N*-terminal extension with several “-SPKK-” (-Ser/Thr-Pro-X/Lys-Lys/Arg-) motifs (see text below) and helps in the tight compaction of chromatin in the sperm head. On fertilization, this is phosphorylated and then replaced during subsequent cell divisions with somatic variants. *Drosophila* has two evolutionarily conserved H2A variants: H2A.X, which has a *C*-terminal extension (like wheat H2A1, which has 19 extra residues and binds to [linker DNA](#)), and H2A.vD, which is essential for *Drosophila* development and has counterparts in mammals (H2A.Z), chickens (H2A.F/Z), and the ciliated protozoan *Tetrahymena* (hv1), where it is found in actively transcribed chromatin. H2A.Z is interesting because it has an *N*-terminal tail resembling that of H4 and indeed is acetylated more than the canonical H2A, thus providing obvious additional possibilities for regulation of function. For further details, see Ref. 2.

A class of more extreme so-called core histone variants are really hybrid proteins, with a region of histone sequence fused to a wholly unrelated region of unknown function. Two such proteins are the mammalian centromere-specific protein CENP-A (the yeast homologue is CSE4) and macroH2A. CENP-A (molecular mass ~17,000Da) is similar to H3 in its *C*-terminal domain, but has a highly divergent *N*-terminal region and appears to be associated, although not solely, with [a-satellite DNA](#) (12). Assuming that it replaces the normal H3, it presumably imparts special properties to centromeric chromatin. MacroH2A is a highly conserved protein in which the *N*-terminal third resembles H2A and the *C*-terminal two-thirds contains a [coiled-coil](#) protein dimerization motif. There are two macroH2A subtypes, which are (a) highly conserved and (b) identical in the histone region. It came as a surprise to find that one subtype is localized to the inactive [X-chromosome](#) (in mouse) and is distributed throughout the chromatin (13); it is tempting to speculate that the coiled-coil might have a role in protein oligomerization and condensation of chromatin. MacroH2A is thus one more distinctive component of the inactive X-chromosome, the others being heavily **methyated** DNA, hypoacetylated H4, and association with a large **cis-acting** nuclear RNA (termed Xist).

## 2. Linker Histones/Histone H1

Linker histones (H1 and its variants and subtypes) in higher eukaryotes are larger than the core histones (molecular mass ~20,000 to 25,000 Da) and are particularly lysine-rich. They have a tripartite domain structure, in which a central globular domain of about 80 amino acid residues, which is highly conserved between species, is flanked by basic *N*- and *C*-terminal domains (“tails”) that are much more divergent. In the absence of DNA, the tails are disordered and may be selectively removed by proteolysis, leaving the globular domain intact. In a typical mammalian H1, the *N*- and *C*-terminal tails are about 40 and 100 residues long, respectively; the lengths differ in some species-specific and cell-type-specific variants. Nucleosomes released from chromatin by digestion with



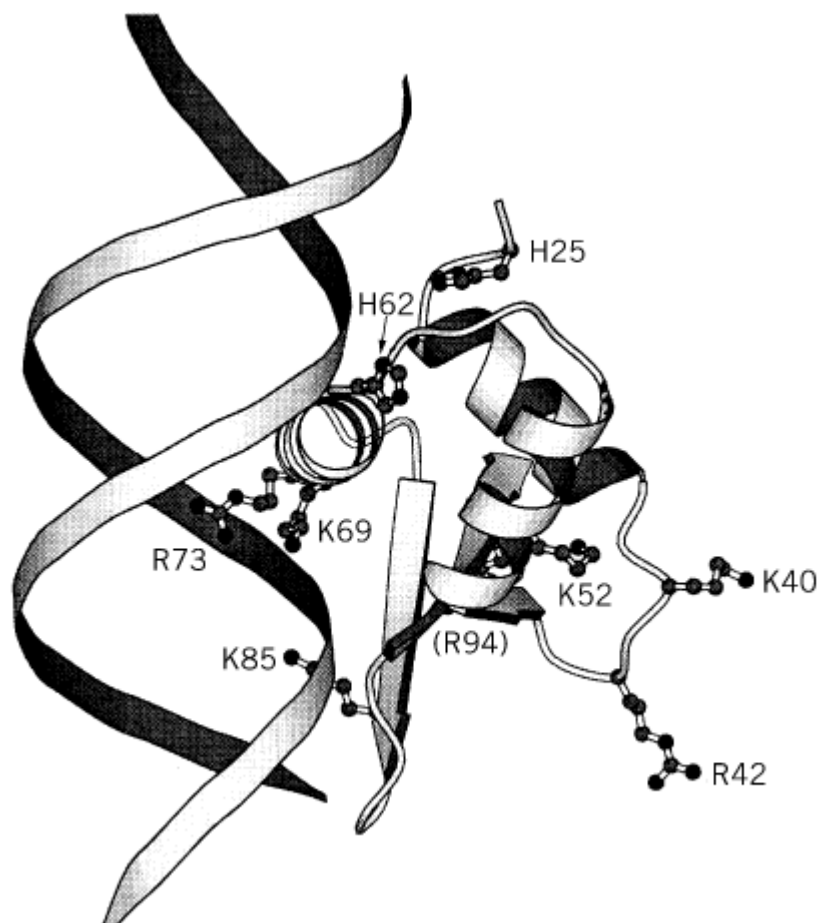
micrococcal nuclease (**Staphylococcal nuclease**) are rapidly trimmed by further digestion so that the ~200 bp of DNA is reduced to ~166 bp, giving a chromatosome. If H1 is removed before digestion, however, only 146 bp remain protected, as digestion proceeds to the limit nucleosome core particle. The globular domain of H1 alone is sufficient to protect 166 bp, so this is the region of H1 that binds close to the nucleosome core, stabilizing the nucleosome and protecting an extra 20 bp from digestion. The basic C-terminal tail (and possibly also the N-terminal tail) binds to the linker DNA between adjacent nucleosomes (hence the name *linker histones*), partially neutralizing its negative charge and promoting folding of the nucleosome filament into a higher-order structure, the 30-nm filament; H1 is located on the inside of the 30-nm filament (see [Chromatin](#)). The role of H1 is thus to stabilize both the nucleosome and an ordered higher-order chromatin structure. H1 may be phosphorylated (see text below) and poly-ADP-ribosylated; the role of the latter modification is not clear, but it is probably involved in [DNA repair](#) .

There are many subtypes and more extreme sequence variants of H1, which all share the same domain organization. These may coexist in the same cell type (eg, there are six H1 subtypes in chicken erythrocytes), and the assumption is that they might stabilize different chromatin higher-order structures. This would be likely to occur only if they were clustered along a chromatin filament, and this indeed appears to be the case for one of the seven H1 variants in the interphase polytene chromosomes in the Dipteran insect *Chironomus thummi* (14). Information is lacking for other systems or for other H1s in *Chironomus*. There are several instances of changes in linker histone variants during developmentally regulated processes (11, 15). During sea urchin embryogenesis, six distinct H1 subtypes are produced at particular stages; and *Xenopus* has an embryonic form of H1, namely B4, that is very different from the normal somatic H1. The C-terminal tail is less basic, and this might result in a less condensed chromatin structure in *Xenopus* embryos, compatible with rapid nuclear division and [DNA replication](#). Extreme variants may also be produced to shut down transcription in particular cell types, and they appear to be associated with more stable chromatin higher-order structures (see [Chromatin](#)). For example, in transcriptionally inert mature avian erythrocytes, H5 (which binds more tightly than H1, probably due to its higher arginine content, although it has shorter N- and C-terminal tails) has largely replaced H1; and in sea urchin sperm, somatic H1 is replaced by a larger and more basic sperm-specific H1, namely spH1, which, together with sperm-specific core histones, is very effective in condensation of the DNA. (In mammalian sperm, histones are replaced by [protamines](#), small arginine-rich proteins, rather than by specialized histone variants.) H1<sup>0</sup>, an H1 variant whose appearance in mammalian cells appears to correlate with terminal differentiation (eg, in neurons), is much more similar to H5 than to H1 in its globular domain.

The globular domains of histones H1 and H5 (GH1 and GH5) can bind to the body of the nucleosome (see text above) and protect an additional 20 bp beyond the core particle length from digestion. The structure of GH5, determined by X-ray crystallography, shows a “winged helix” motif (a variant on the **helix–turn–helix** DNA-binding motif) which is also found in two proteins that bind to specific DNA sequences, namely the liver transcription factor HNF<sub>3</sub>g (hepatocyte nuclear factor 3g,) and the prokaryotic protein **cyclic AMP receptor protein** (CRP) (16). It consists of a three-helix bundle with a b-hairpin (or “wing”) at its C-terminus. Based on the X-ray crystal structures of these proteins bound to DNA, the likely mode of binding of GH5 to DNA, in the major groove, has been deduced. Biochemical evidence points to a second DNA binding site on GH5, and this has been identified on the opposite face of the globular domain (17) (Fig. 2). It is likely that the two sites are occupied in the nucleosome by two of the three duplexes (the central turn and the entering and exiting DNA) that are present in the vicinity of the dyad axis of the nucleosome. Current evidence for bulk nucleosomes (18) suggests that GH5 bridges the central turn, close to the dyad, and one of the entering/exiting duplexes. The orientation proposed would place the C-terminal tail of H1 directed toward linker DNA, as expected from various lines of evidence. A different mode of binding of the globular domain to the reconstituted *Xenopus* 5 S nucleosome has been proposed, but it is not yet clear whether this is a special property of the 5 S nucleosome (see [Nucleosome](#)). The existence of two binding sites on the globular domain is probably the basis of its preference (over

linear DNA) for four-way DNA junctions (19), which mimic a pair of duplexes at the DNA crossover point near the nucleosome dyad.

**Figure 2.** The structure of the globular domain of histone H5 (GH5) showing the two clusters of basic residues at the two proposed DNA-binding sites on opposite faces of the domain (17). Binding to DNA at one site is by analogy with the structures of HNF3g and CRP (see text).



There is no high-resolution structural information for the *N*- and *C*-terminal tails of H1 or any of its variants. The role of the *N*-terminal tail (~40 amino acid residues in the canonical H1) is unclear, but there is some evidence that it may anchor the globular domain correctly in the presence of the *C*-terminal tail. The basic *C*-terminal tails (~90 to 130 amino acid residues, depending on the variant) have unusually high contents (~50%) of the basic residues lysine and arginine, as well as relatively high contents of alanine and proline. They are essential for chromatin condensation and are believed to bind to [linker DNA](#). They are disordered in the free protein, but they are likely in the presence of DNA to exist as  $\alpha$ -helical segments rich in lysine and alanine, the segments being separated by proline residues (20). The proline residues are often embedded in so-called “-SPKK-” sequence motifs (-Ser/Thr-Pro-X/Lys-Lys/Arg-), where the serine is a potential phosphorylation site. There is some phosphorylation at S-phase of the cell cycle and considerably more at mitosis, probably due to the action of the cyclin-dependent kinase,  $p34^{cdc2}$  (21). The role of phosphorylation is likely to be to loosen the interactions between the H1 tails and linker DNA, and at mitosis to permit other factors and interactions to drive chromosome condensation. Some members of the repertoire of H1s produced during different stages of sea urchin embryogenesis lack -SPKK- motifs, suggesting that the binding of different H1s may be differentially regulated. Phosphorylation also plays a role in the control of chromatin transitions unrelated to the cell cycle, namely, in the late stages of avian

erythropoiesis, or sea urchin sperm spermatogenesis, where there is no cell division and where special histone variants that bind more tightly than somatic H1 (ie, H5 and spH1, respectively) suppress transcription by promoting chromatin condensation. H5 and spH1 are phosphorylated at the -SPKK- motifs (of which there are six copies in the *N*-terminal tail of spH1 and further copies in the C-tail) until the final stage of maturation, when a final dephosphorylation step in the spermatozoon or mature erythrocyte results in the final stage of chromatin condensation and even closer chromatin packing in the sperm head or erythrocyte nucleus; removal of the phosphate increases the net positive charge on the histone and promotes histone–DNA interactions, leading to condensation. An analogous situation exists in the amitotic nucleus of *Tetrahymena* (22).

There are two very atypical H1s, or—in one case—candidate H1, in lower eukaryotes, which lack the characteristic domain organization of the canonical H1. The ciliated protozoan *Tetrahymena thermophila*, which has a transcriptionally active but amitotic **macronucleus** and a transcriptionally inert mitotic **micronucleus**, contains a small, basic, macronuclear protein of 163 residues, designated H1, that appears to share some of the functions of linker histones. (A set of distinctive polypeptides appears to substitute for H1 in the micronucleus.) The macronuclear H1 condenses chromatin in the nucleus, and phosphorylation at ‘TPVK’ sites results in an increase in nuclear volume, suggesting that phosphorylation loosens interactions with chromatin, as expected (22, 23). It lacks the distinctive globular domain of the canonical H1, however, although its *C*-terminal region has similarities to the *C*-terminal tail of H1. Exactly how the *Tetrahymena* H1 is bound to chromatin is unclear. An unusual, but quite different, domain organization also occurs in a candidate H1 identified from the recently determined complete sequence of the yeast (*Saccharomyces cerevisiae*) genome. Previous attempts to isolate H1 from yeast had been unsuccessful, and one view was that because the linker length of yeast chromatin was virtually zero (chromatin repeat length ~166 base pairs) there was no requirement for neutralization of linker DNA charge or, therefore, for a canonical linker histone. Analysis of the yeast genome revealed an open reading frame (the *HHO1* gene) encoding a protein (Hho1p) with regions of sequence homology to the globular domain of H1, which is regarded as a candidate H1 (24). It has *two* globular domains of about 80 residues, with a basic *N*-terminal extension and a basic connecting linker with some resemblance to the *C*-terminal tail of the canonical H1. Recombinant Hho1p appears to have some of the distinctive properties of H1 in a standard chromatin reconstitution, namely protection of an additional ~20bp of DNA beyond the core particle length against exonuclease digestion (25). Its unique structure, however, means that there are likely to be differences in its detailed mode of binding to chromatin compared with canonical H1, and it is by no means clear at present whether Hho1p functions as a true H1 or is instead a transcription factor with domains homologous to the globular domain of H1 [cf. some TAFs, which have sequence homology to core histones (see [Histone Fold](#)); or even HNF3<sub>g</sub> which contains the same structural fold as in the globular domains of H5 and H1 (see text above), although in that case there is no sequence homology between the two proteins].

Gene disruptions and deletions have been used to ask whether H1 is essential. Disruption (“knockout”) of the H1<sup>0</sup> gene in mouse appeared to be without consequence; it is likely, however, that one or some of the six other subtypes compensate (26); in other words, there is functional redundancy. Deletion of the entire complement of somatic H1 genes has not been achieved. Early *Xenopus* embryos contain the unusual H1 variant B4 (see text above) which is later replaced by somatic H1; elimination of B4 had little effect on nuclear assembly or on the development of the organism (27). Deletion of the single gene for the atypical macronuclear H1 of *Tetrahymena* was not lethal, but was not without effect either (23). Vegetative growth, general transcription, and general nucleosome repeat length were unaffected, but there were changes in the nuclear volume (presumably consistent with the role of the protein in chromatin condensation and packing in the nucleus), in the efficiency of meiotic division, and, significantly, in the transcriptional regulation of specific genes. However, although some genes were activated, consistent with the expected repressive role for H1, others were repressed (eg, the *CyP* gene, which encodes a [thiol proteinase](#)). One explanation would be that the binding of the *Tetrahymena* H1 is needed in some nucleosomes to position them in such a way that certain sequences necessary for the binding of activators are

exposed. In yeast, deletion of the yeast *HHO1* gene had no detectable effect on cell growth, viability, or mating, or on telomeric silencing, basal transcriptional repression, or efficient **sporulation** (25, 28). It also did not affect transcription of the SW1/SNF-dependent *SUC2* gene or the repression of the silent *a1/a2* genes, or activation at a distance of a GAL1 promoter (28). It remains to be seen whether the Hho1p protein is really a *bona fide* H1 or a transcription factor evolutionarily related to the globular domain of histone H1, despite assuming some of the functional properties of H1 in an *in vitro* assay (25). These and other studies suggest strongly that H1 and its variants probably have roles beyond that of simple repression and stabilization of higher-order structure in chromatin. Gene-specific effects of H1 are clear in one well-documented case, namely, the 5 S rRNA genes of *Xenopus laevis*. Replacement of the unusual embryonic variant B4 with the somatic H1 during embryogenesis has been shown to be causal for the selective repression of the oocyte 5 S, leaving the somatic genes active (29). Gene-specific effects of H1, in addition to a general, default, repressive role, may turn out to be more common than might have been imagined.

## Bibliography

1. K. van Holde (1988) *Chromatin*, Springer-Verlag, New York.
2. A. P. Wolffe (1995) *Chromatin: Structure and function*, 2nd ed., Academic Press, London.
3. J. O. Thomas and R. D. Kornberg (1975) Proc. Natl. Acad. Sci. USA **72**, 2626–2630.
4. G. Arents et al. (1991) Proc. Natl. Acad. Sci. USA **88**, 10148–10152.
5. A. Klug et al. (1980). Nature **287**, 509–516.
6. K. Luger et al. (1997) Nature **389**, 251–260.
7. M. Grunstein (1997) Nature **389**, 349–352.
8. L. Huang, W. Z. Zhang, and S. Y. Roth (1997) Mol. Cell. Biol. **17**, 6555–6562.
9. P. T. Georgel, T. Tsukiyama, and C. Wu (1997) EMBO J. **16**, 4717–4726.
10. M. J. Barratt, C. A. Hazzalin, E. Cano, and L. C. Mahadevan (1994) Proc. Natl. Acad. Sci. USA, **91**, 4781–4785.
11. D. Poccia (1986) Int. Rev. Cytol. **105**, 1–65.
12. K. F. Sullivan, M. Hechenberger, and K. Masri (1994) J. Cell Biol. **127**, 581–592.
13. C. Costanzi and J. R. Pehrson (1998) Nature **393**, 599–601.
14. E. Mohr, L. Trieschmann, and U. Grossbach (1989) Proc. Natl. Acad. Sci. USA **86**, 9308–9312.
15. S. Khochbin and A. P. Wolffe (1994) Eur. J. Biochem. **225**, 501–510.
16. V. Ramakrishnan (1994) Curr. Opin. Struct. Biol. **4**, 44–50.
17. F. A. Goytisolo et al. (1996) EMBO J. **15**, 3421–3429.
18. Y.-B. Zhou et al. (1998) Nature **395**, 402–405.
19. P. Varga-Weisz et al. (1994) Proc. Natl. Acad. Sci. USA **91**, 3525–3529.
20. D. J. Clark, C. S. Hill, S. R. Martin, and J. O. Thomas (1988) EMBO J. **7**, 69–75.
21. R. W. Lennox and L. H. Cohen (1988) In *Chromosomes and Chromatin*, Vol. **1** (K. W. Adolph, ed.) CRC Press, Boca Raton, FL. pp. 33–56.
22. S. Y. Roth and C. D. Allis (1992) Trends Biochem. Sci. **17**, 93–98.
23. X. T. Shen and M. A. Gorovsky (1996) Cell **86**, 475–483.
24. S. C. Ushinsky et al. (1997) Yeast **13**, 151–161.
25. H. G. Patterson et al. (1998) J. Biol. Chem. **273**, 7268–7276.
26. A. M. Sirotkin et al. (1995) Proc. Natl. Acad. Sci. USA **92**, 6434–6438.
27. M. Dasso, S. Dimitrov, and A. P. Wolffe (1994) Proc. Natl. Acad. Sci. USA **91**, 12477–12481.
28. D. Escher and W. Schaffner (1997) Mol. Gen. Genet. **256**, 456–461.
29. P. Bouvet, S. Dimitrov, and A. P. Wolffe (1994) Genes Dev. **8**, 1147–1159.

## Suggestions for Further Reading

30. V. Ramakrishnan (1997) Histone structure and the organization of the nucleosome. *Ann. Rev. Biophys. Biomol. Struct.* **26**, 83–112.
31. V. Ramakrishnan (1997) Histone H1 and chromatin higher-order structure. *Crit. Rev. Eukaryot. Gene Expr.* **7**, 215–230.
32. K. van Holde and J. Zlatanova (1996) The linker histones and chromatin structure: new twists. *Prog. Nucleic Acid Res. Mol. Biol.* **52**, 217–259.
33. A. A. Travers (1999) Towards a higher order structure for chromatin—the location of the linker histone in the nucleosome. *Trends Biochem. Sci.* (in press).
34. R. D. Cole (1987) Microheterogeneity in H1 histones and its consequences. *Int. J. Protein Res.* **39**, 433–449.

## HIV (Human Immunodeficiency Virus)

Human immunodeficiency virus (HIV), a causative virus of acquired immune deficiency syndrome (AIDS), is a member of the *Lentivirus* genus of the family *Retroviridae* (see [Rous Sarcoma Virus \(RSV\)](#)). The first isolate was reported in 1983 and was named lymphadenopathy-associated virus (LAV). Other isolates, such as human T-lymphotropic virus type III (HTLV-III) and AIDS-associated retrovirus (ARV), were reported in 1984. These HIV isolates were later categorized as HIV-1, because a serologically unique HIV isolate was discovered in 1986 and designated as HIV-2. Great sequence variation, especially in the *env* gene, is characteristic of HIV, and at least nine subtypes of HIV-1 (A to I) and/or five of HIV-2 (A to E) have been classified.

HIV virions are spherical and about 110 nm in diameter. The HIV virion has a lipid-bilayer envelope, and approximately 72 knobs (spikes) protrude from the envelope. Each knob is composed of oligomeric surface protein (SU:gp120 in HIV-1 and gp130 in HIV-2) and transmembrane protein (TM:gp41 in both HIV-1 and HIV-2). Inside the envelope, the capsid proteins (CA:p24) form the bullet-shaped capsid shell. The matrix protein (MA:p17) is present between the envelope and CA. The nucleocapsid protein (NC:p9) binds to the viral genome. SU, TM, MA, CA, and NC are the major structural proteins. Two identical copies of single-stranded RNA of positive polarity are contained in a HIV virion as the [genome](#). The enzymes [proteinase](#) (PR), **reverse transcriptase** (RT) (see [Rous Sarcoma Virus \(RSV\)](#) and [Hepatitis B Virus](#)), and [integrase](#) (IN) are also contained in the capsid shell.

The genome RNA of about 9.2 kb is, like cellular [messenger RNAs](#), capped at the 5' end (see [Cap](#)) and tailed by [poly\(A\)](#) at the 3' end. The genome has short [direct repeats](#) (R) and unique sequences at both ends (U5 at the 5' end and U3 at the 3' end). These sequences are important to form [long terminal repeats](#) (LTRs) during replication. The coding sequences of HIV have three major genes of retroviruses: *gag*, *pol*, and *env*. The *gag* gene codes for MA, CA, and NC; the *pol* gene codes for PR, RT, and IN; the *env* gene codes for SU and TM. It is characteristic of HIV to have accessory genes. Two of them, *tat* and *rev*, are essential for replication in both HIV-1 and HIV-2. HIV-1 has *vif*, *vpr*, *nef*, and *vpu*, while HIV-2 has *vif*, *vpr*, *nef*, and *vpx* instead of *vpu* as nonessential accessory genes.

The replicative cycle starts with an interaction between the virion and the cell surface receptor. The principal receptor for HIV is the CD4 molecule, which is expressed on CD4-positive T lymphocytes, macrophages, and so on. The CD4-binding sites on the virion are on SU (gp120). The next crucial step after binding is fusion of the virion envelope and the plasma membrane. It has been postulated that an *N*-terminal hydrophobic **domain** in TM is the actual fusion protein, which needs to be

exposed by a conformational change in the envelope glycoproteins caused by SU binding to CD4.

Interestingly, an additional cell-surface molecule, actually a [chemokine](#) receptor like CCR5 or CXCR4, is needed to initiate the fusion process. These chemokine receptors are referred to as “co-receptors” for HIV infection. Once the nucleoproteins get into the cytoplasm, synthesis of the minus-strand DNA starts using RT and host [transfer RNA](#)<sup>Lys</sup> attached near the 5' end of the viral genome as a primer. Synthesis of viral DNA is similar to that of other mammalian and avian retroviruses, but lentiviruses, including HIV, start the plus-strand DNA synthesis at two sites; one is a polypurine tract near U3 like many other retroviruses, and the other is a central polypurine tract at the end of pol gene. Double-stranded viral DNA with the LTR (see [Rous Sarcoma Virus \(RSV\)](#)) on both ends is carried to the [nucleus](#) as preintegration complex. Presumably, nuclear localization signals in MA and/or Vpr play an important role in nuclear translocation of the preintegration complex. While mammalian and avian retroviruses require cell division for efficient virus production, HIV can infect nondividing, terminally differentiated cells. Linear double-stranded viral DNA is the substrate for IN and the direct precursor of integrated provirus.

Integrated proviral DNA has an LTR on both ends. Full-length viral RNA is transcribed from the U3/R boundary in the 5'LTR to the R/U5 boundary in the 3'LTR, and 5'-capped and 3'-poly(A)-tailed viral RNA serves as both genome RNA and mRNA for [translation](#). The mRNA undergoes complex [RNA splicing](#). The env and accessory gene products are translated from spliced mRNAs. The env gene products (Env) are translated as a precursor protein, gp160, in the rough [endoplasmic reticulum](#), processed to gp120 (SU) and gp41 (TM), and expressed on the cell surface. The gag and pol gene products (Gag and Pol) are translated from unspliced mRNA. The pol gene is not in the same reading frame as gag gene, however, and pol gene products are translated as a Gag–Pol fusion protein by ribosomal [frameshifting](#) at the gag–pol boundary. Viral assembly occurs just below the plasma membrane. Virions bud out from the cells, covered by the plasma membrane (viral envelope) with SU and TM. Viral proteinase (PR) digests Gag and Gag–Pol fusion proteins to MA, CA, NC, PR, RT, and IN, and the virions are matured.

#### Suggestion for Further Reading

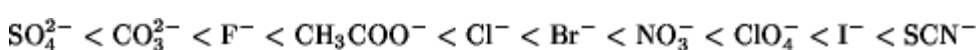
P. A. Luciw (1996) "Human Immunodeficiency Viruses and Their Replication. In" *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 1881–1952.

## Hofmeister Series

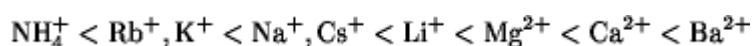
The Hofmeister (or *lyotropic*) series, first described by Hofmeister in 1888, is the ordering of ions in their effectiveness in salting out proteins (see [Precipitation](#)). Since then, it has been established that the same order is maintained in the effects of these ions on **nucleic acid** and [protein](#) stability (see **Stabilization and destabilization by co**), association–dissociation equilibria between macromolecules, [enzyme](#) activity, and various other biological and biochemical functions (1). The ordering is as follows:

← Stabilization    Unfolding →  
← Salting in        Salting out →

Anions:



Cations:



Anions and cations are essentially independent of each other in these actions, and their effects are additive. For example, it has been found that LiBr is a good salting in agent, whereas NaCl is a salting out agent (2). Reshuffling the ions gives the result that both NaBr and LiCl have no effect on the solubility, in other words, the salting in capacity of one ion ( $\text{Li}^+$  or  $\text{Br}^-$ ) is compensated by the salting out characteristic of the co-ion ( $\text{Cl}^-$  or  $\text{Na}^+$ ).

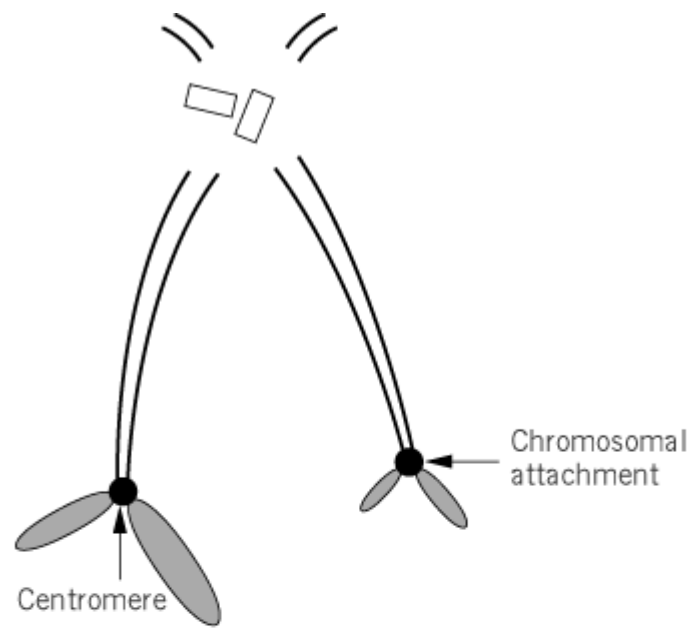
#### Bibliography

1. P. von Hippel and T. Schleich (1969) In *Structure and Stability of Biological Macromolecules* (S. N. Timasheff and G. D. Fasman, eds.), Marcel Dekker, New York, "Chap. 6".
2. D. R. Robinson and W. P. Jencks (1965) *J. Am. Chem. Soc.* **87**, 2470–2479.

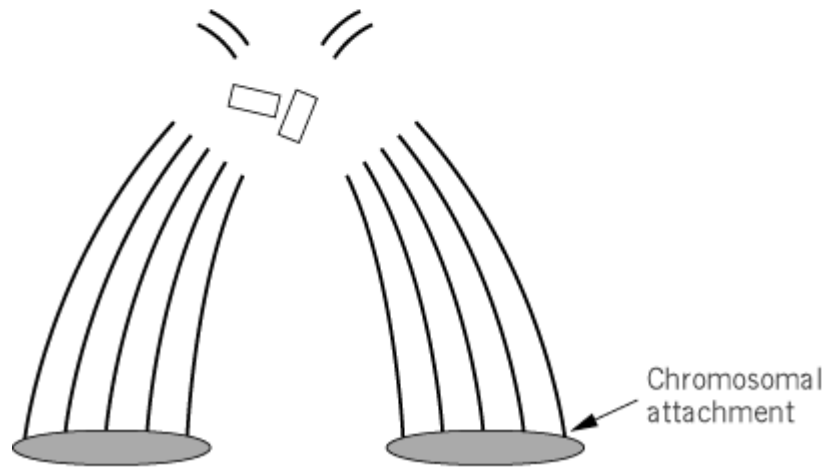
### Holocentric Chromosome

Holocentric or polycentric [chromosomes](#) have centromeric [heterochromatin](#) (see [Centromeres](#)) distributed at multiple sites along the length of the chromosome (1). Organisms with holocentric chromosomes include [nematodes](#), such as *Ascaris*, insects, such as *Lepidoptera*, and **plants** in the genus *Lazula*. Spindle fibers attach along the whole length of holocentric chromosomes (Fig. 1). Sites of attachment of the spindle to centromeric chromatin are interspersed with noncentromeric segments, as in *Ascaris*, or all the chromatin in the chromosome is competent for attachment, as in the coccid bugs (2). In this latter case, centromeric functions are appropriately described as diffuse. In some organisms, like *Lazula*, localized centromeric function might be retained in mitotic chromosomes but shift to different positions in meiotic chromosomes. Mammalian [tissue culture](#) chromosomes also persist and segregate when they have multiple centromeres, but in this case only one centromere remains functional (4).

**Figure 1.** Monocentric and holocentric chromosomes. (a) Monocentric chromosomes have a single site of chromosomal attachment to the spindle. (b) Holocentric chromosomes have multiple sites of attachment to the spindle.



(a)



(b)

### Bibliography

1. S. Dimpinelli and C. Coday (1989) *Trends Genet.* **5**, 310–314.
2. S. Hughes-Schrader (1948) *Adv. Genet.* **2**, 127–203.
3. J. P. Braselton (1981) *Chromosoma* **82**, 143–153.
4. B. K. Vig (1984) *Chromosoma* **90**, 39–45.

### Suggestion for Further Reading

5. R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A synthesis*, Wiley-Liss, New York.



## Holoenzyme, Apoenzyme

A number of [enzymes](#) rely on the presence of a **prosthetic group** at their [active site](#) in order for them to be catalytically active (see [Coenzyme, Cofactor](#)). For example, aminotransferases require [pyridoxal phosphate](#). Prosthetic groups are small molecules that are covalently or tightly bound at the active site. They provide groups that the protein component of the enzyme does not possess, but that are essential for catalysis. The enzyme with a prosthetic group bound is referred to as a holoenzyme, whereas after removal of the prosthetic group, the protein component alone is termed an apoenzyme. Thus,

$$\text{Holoenzyme} = \text{apoenzyme} + \text{prosthetic group} \quad (1)$$

The term apoenzyme is also applied to **allosteric** enzymes such as [aspartate transcarbamoylase](#) from which the regulatory subunits have been removed.

## Homeobox Genes

Homeobox genes play fundamental roles in [development](#) and [evolution](#). They are perhaps the best examples of key regulators of gene [transcription](#) that are at the heart of the genetic circuitry, regulating different pathways. Homeobox genes have been highly conserved throughout evolution and are involved in the genetic control of the body plan, the determination of cell fate, and several other basic developmental processes. Common among the hundreds of different homeobox genes known today is a highly homologous, 180-bp structural motif, the homeobox. The homeobox encodes a 60-amino acid residue [polypeptide chain](#), the **homeodomain**, that represents the **DNA-binding** domain of the respective proteins.

### 1. History of Homeobox Genes

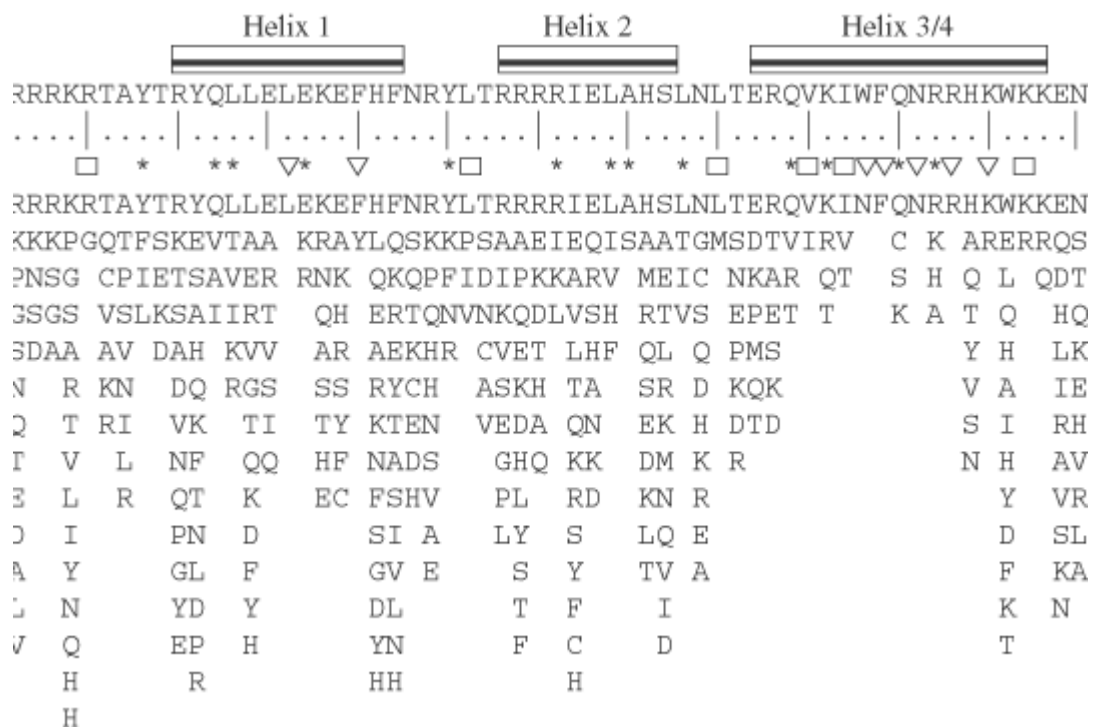
Homeotic mutations lead to segmental transformations, suggesting that they are involved in the genetic control of a body plan. The first homeotic mutant, bithorax, was discovered in 1915 (1). This *Drosophila* gene was later shown to be part of a gene cluster termed the [bithorax complex](#) (BX-C) (2). Single mutations in the bithorax complex induced major modifications in the body plan, for example, generating four-winged flies. In 1983, the first homeobox gene locus, *Antennapedia*, was isolated (3), which is part of the **Antennapedia complex** (Antp). Comparative sequence analysis of the first three cloned homeobox genes, *Antennapedia*, *fushi tarazu*, and **Ultrabithorax**, established a common DNA motif, coined the homeobox (4, 5). This motif was not confined to *Drosophila*, but was also found in vertebrates, including mice and humans. In 1984, the first homeobox gene in the mouse was **cloned** by cross-hybridization using a *Drosophila* homeobox probe (6).

### 2. Structure and Function of the Homeodomain

The highly conserved **homeodomain** is typically composed of 60 amino acid residues. Figure 1 shows the amino acid **consensus sequence** based on 346 different homeodomain sequences (7). Comparison of the [X-ray crystallography](#) structures of the homeodomains of **engrailed** and other homeobox genes of the yeast **mating-type** proteins revealed a **helix-turn-helix** binding motif (8), suggesting that all homeodomains bind to DNA in the same manner. The three **alpha-helical** regions are composed of residues 10–22, 28–38 and 42–58 of the homeodomain (see Fig. 1). Helix 3, the

recognition  $\alpha$ -helix, binds to the major groove of DNA, whereas helices 1 and 2 lie close to each other in an antiparallel orientation outside the double helix. Thus, only a small number of residues in helix 3 (in particular, that at position 50) and in the *N*-terminal arm are responsible for the specificity of contacts with the DNA. Furthermore, *in vitro* binding studies with protein extracts of the locus *Antennapedia* showed that the homeodomain binds as a monomer to its DNA-binding site, dissociated only slowly, with an estimated half-life of 90 min (9). Other examples suggest that homeobox proteins may dimerize. Transcriptional regulation of homeodomain proteins *in vivo* was first demonstrated for *bicoid*, interacting with its target gene *hunchback* (10-12).

**Figure 1.** The primary and secondary structures of the homeodomain. The homeodomain contains three well-defined  $\alpha$ -helices and a more flexible fourth helix. The schematic illustration of the structure represents a composite derived from the structures of the *Antennapedia*, *engrailed*, and *MATa2* homeodomains. The amino acid consensus sequence is based on 346 homeodomain sequences. For each position, the amino acid encountered most frequently is listed at the top, while other amino acids are listed beneath in decreasing order of frequency of occurrence. Amino acids occurring fewer than 5 times (1.5%) are not shown. The symbols on top indicate highly conserved positions:  $\nabla$  are the most highly conserved positions, with only one or two amino acids found at that position. *f* indicates highly conserved positions with three to five different amino acids found at a particular position, and \* indicates conserved positions, with not more than nine different amino acids.



### 3. HOX Gene Clusters

Homeobox genes are clustered in *Drosophila* and organized in two separate clusters, the bithorax (BC-C, 2) and *Antennapedia* (*Antp*) complexes (13, 14), and are referred to as *HOM-C complexes*. The genes of these two gene clusters have been termed *HOX genes*. HOX genes play a critical role in pattern formation. The order of HOX genes within each cluster correlates with the temporal order of expression of these genes, which are expressed along the anteroposterior axis of the *embryo*. Interestingly, homeobox genes with sequence similarity to *Drosophila HOX genes* are also clustered in vertebrates. Both humans and mice have four HOX gene clusters, which reside on different chromosomes. They are termed *Hoxa* to *Hoxd* in mice and *HOXA* to *HOXD* in humans (15). Altogether, 38 genes reside in these four gene clusters, and all are orientated in the same 5'-3' direction of transcription (16, 17). Furthermore, in mammals the order of HOX genes within a cluster

dictates the time of expression.

#### 4. Homeotic Mutations

The term *homeotic mutations* actually refers to HOX gene mutations within the homeotic clusters. For example, a homeotic mutation in the *Antennapedia* gene complex transforms an antenna into a leg. A mutation in the *bithorax* gene results in a fly with four wings. A mutation in human HOXD13 results in synpolydactyly. Some of the mouse HOX mutations have been shown to lead to only relatively mild phenotypes. This may be due to functional compensation or to different parallel regulatory pathways.

In a broader sense, homeotic mutations may encompass mutations of all homeobox genes. Homeotic mutations may represent the classic examples of loss-of-function mutations or gain-of-function mutations. It is interesting that missense mutations causing loss-of-function mutations are often concentrated in the homeodomain or in other functionally important parts of the protein. In contrast, nonsense mutations leading to gene truncations are distributed more widely across the genes. Also, a 50% reduction in the gene product (haploinsufficiency) or overexpression may lead to different clinical phenotypes, suggesting that correct dosage is crucial at certain times in development. This was shown for the first time for PAX6, which leads to aniridia when underexpressed and to other severe eye abnormalities when overexpressed ([18](#), [19](#)) (see [Pax Genes](#)).

#### 5. Homeodomain Gene Families

Using degenerate oligonucleotide probes or amplification by **PCR**, a wide variety of novel homeobox genes that represent non-HOX homeobox genes have been identified. Some of these “isolated” (“orphan”) homeodomain gene families contain additional highly conserved domains. Examples of these conserved motifs are the **paired-**, **POU-**, **LIM-**, **engrailed-**, and **zinc-finger** motifs, which specify the individual classes of homeobox genes. Each of the two independent DNA-binding domains contacts a specific short sequence on the major groove of the DNA double helix. Examples of the homeobox genes with a paired domain include several members of the PAX gene family (PAX 2, 3, 6); POU domain genes are **Pit1**, **Oct1**, and **unc 86**; genes with a LIM motif include the nematode *lin-11* and putative human **oncogene** *rhombotin1*; the engrailed class of homeobox genes includes *EH1-5*, and a gene with a zinc-finger motif is the human *ATBF1* gene.

A homeobox page on the different homeobox genes is maintained by Thomas R. Bürglin and can be accessed via the Internet (site currently unavailable).

#### 6. Homeobox Genes in Disease

In the last few years, gene disruption experiments in transgenic mice, as well as positional cloning in humans, have defined a number of homeobox genes associated with genetic diseases and congenital syndromes. The first human HOX gene mutation, HOXD13, was described in families with synpolydactyly (webbing and duplication of fingers) ([20](#)). Here the homeodomain was still intact, and the affected members of these families had an expansion of 7–10 additional alanine residues. RIEG, a bicoid-related homeobox gene, was shown to be involved in Rieger syndrome, a complex human disorder with dental, ocular, and further abnormalities ([21](#)). Just recently, the chick and *Xenopus* homologues of RIEG (*Pitx2*) were shown to determine the left–right asymmetry of internal organs in vertebrates ([22](#)). SHOX, a **paired**-related gene, was shown to play a major role in human growth ([23](#)), affecting the height in Turner and Léri-Weill syndrome patients, Langer dwarfism, and idiopathic short stature patients. The MSX1/MSX2 family of genes has been found to be mutated in some forms of craniosynostosis ([24](#)). PAX2 and PAX6 genes play a role in eye development, and PAX3 in the Waardenburg syndrome, the most common form of inherited deafness ([25](#)). Two gene families in particular, *Emx* and *Otx*, have been reported to play a role in brain development, and *EMX2* mutations, for example, have been associated with schizencephaly ([26](#)).

## 7. Evolutionary Conservation

Because of their fundamental importance in embryonic development, homeobox genes have been highly conserved throughout evolution. Homeobox genes have been found in plants, fungi, insects, and vertebrates. Both the **secondary structures** and **tertiary structures** of homeodomains of different species have been highly conserved. The two HOM-C clusters in *Drosophila*, and the four HOX clusters in mice and humans, arose by **gene duplication** and **divergence** from a common ancestral cluster (27). Such a single gene cluster was found in *Amphioxus*, the closest invertebrate relative of the vertebrates.

Functional homology between *Drosophila* and mouse HOX genes has been shown by gene knockout and transgenic mice (28, 29). Most striking are the gene transfer experiments carried out between mouse and *Drosophila*. They demonstrated, for example, that the mouse HOXB-6 gene can induce antennal legs when ectopically expressed in *Drosophila* (30). The ability to compare and extrapolate among species as diverse as *Drosophila*, mice, and humans is also highlighted in experiments with the mutants of eyeless (*ey*) in *Drosophila*, small eye (*Pax6*) in mice, and Aniridia (*PAX6*) in humans, which share a 90% sequence identity in their homeodomains. The murine *Pax6* is not only strongly homologous to the *Drosophila* *PAX6*, but both can induce ectopic eye structures on the wings, legs, and antennae by targeted expression (31).

### Bibliography

1. C. Bridges and T. H. Morgan (1923) Carnegie Inst. Wash. Publ. **327**, 1–251.
2. E. B. Lewis (1978) Nature **276**, 565–570.
3. R. L. Garber, A. Kuroiwa, and W. J. Gehring (1983) EMBO J. **2**, 2027–2036.
4. M. P. Scott and A. J. Weiner (1984) Proc. Natl. Acad. Sci. USA **81**, 4115–4119.
5. W. McGinnis, M. S. Levine, E. Hafen, A. Kuroiwa, and W. J. Gehring (1984) Nature **308**, 428–433.
6. W. McGinnis, C. P. Hart, W. J. Gehring, and F. H. Ruddle (1984) Cell **38**, 675–680.
7. T. R. Bürglin (site currently unavailable).
8. A. Laughon and M. P. Scott (1984) Nature **310**, 25–31.
9. M. Affolter, A. Percival-Smith, M. Müller, W. Leupin, and W. J. Gehring (1990) Proc. Natl. Acad. Sci. USA **87**, 4093–4097.
10. W. Driever and C. Nüsslein-Volhard (1989) Nature **337**, 138–143.
11. W. Driever, G. Thoma, and C. Nüsslein-Volhard (1989) Nature **340**, 363–367.
12. G. Struhl, K. Struhl, and P. M. Macdonald (1989) Cell **57**, 1259–1273.
13. T. C. Kaufman, R. Lewis, and B. Wakimoto (1980) Genetics **27**, 309–362.
14. T. C. Kaufman, M. Seeger, and G. Olsen (1990) Adv. Gen. **27**, 309–362.
15. D. Duboule, ed. (1994) *Guidebook to the Homeobox Genes*, IRL Press, Oxford.
16. M. Kessel and P. Gruss (1990) Science **249**, 374–379.
17. R. Krumlauf (1992) BioEssays **14**, 245–252.
18. T. Jordan, I. Hanson, D. Zaletayev, S. Hodgson, J. Prosser, A. Seawright, N. Hastie, and V. van Heyningen (1992) Nat. Gen. **1**, 328–332.
19. A. Schedl, A. Ross, M. Lee, D. Engelkamp, P. Rashbass, V. van Heyningen, and N. D. Hastie (1996) Cell **86**, 71–82.
20. Y. Muragaki, S. Mundlos, J. Upton and B. R. Olsen (1996) Science **272**, 548–551.
21. E. V. Semina, R. Reiter, N. J. Leysens, W. L. M. Alward, K. W. Small, N. A. Datson, J. Siegel-Bartelt, D. Bierke-Nelson, P. Bitoun, B. U. Zabel, J. C. Carey, and J. C. Murray (1996) Nat. Gen. **14**, 392–399.
22. A. K. Ryan, B. Blumberg, C. Rodriguez-Esteban, S. Yonei-Tamura, K. Tamura, T. Tsukui, J. de la Peña, W. Sabbagh, J. Greenwald, S. Choe, D. P. Norris, E. J. Robertson, R. M. Evans, M. G.

- Rosenfeld, and J. C. I. Belmonte (1998) *Nature* **394**, 545–551.
23. E. Rao, B. Weiss, M. Fukami, A. Rump, B. Niesler, A. Mertz, K. Muroya, G. Binder, S. Kirsch, M. Winkelmann, G. Nordsiek, U. Heinrich, M. H. Breuning, M. B. Ranke, A. Rosenthal, T. Ogata, and G. A. Rappold (1997) *Nat. Gen.* **16**, 54–63.
  24. D. Davidson (1995) *Trends Gen.* **11**, 405–411.
  25. D. Engelkamp and V. van Heyningen (1996) *Curr. Opin. Gen. Dev.* **6**, 334–342.
  26. E. Boncinelli (1997) *Curr. Opin. Gen. Dev.* **7**, 331–337.
  27. C. Kenyon (1994) *Cell* **78**, 175–180.
  28. M. Kessel and P. Gruss (1991) *Cell* **67**, 89–104.
  29. H. LeMouellic, Y. Lallemand, and P. Brûlet (1992) *Cell* **69**, 251–264.
  30. J. Malicki, K. Schughart, and W. McGinnis (1990) *Cell* **63**, 961–967.
  31. G. Halder, P. Callaerts, and W. J. Gehring (1995) *Science* **267**, 1788–1792.

### **Suggestions for Further Reading**

32. D. Duboule, ed. (1994) *Guidebook to the Homeobox Genes*, IRL Press, Oxford.
33. E. Boncinelli (1997) Homeobox genes and disease, *Curr. Opin. Gen. Dev.* **7**, 331–337.

## **Homeostasis**

Homeostasis is the tendency of an organism to maintain an equilibrium among its internal physiological functions and between the organism and its environment. Metabolism within the cell is regulated to maintain a constant inner environment. Organ systems for respiration, digestion, metabolism, and excretion all work to maintain homeostasis within the organism. Some vertebrates (birds and mammals) also have homeostatic mechanisms to maintain their internal temperatures.

### **Suggestions for Further Reading**

- G. G. Simpson and W. S. Beck (1965) *Life, an Introduction to Biology*, 2nd ed., Harcourt, Brace, and World, New York, p. 106.
- H. W. Smith (1953) *From Fish to Philosopher*, Little, Brown, Boston. (A wonderful little book on the function and evolution of excretory systems in vertebrates.)
- A. S. Romer (1962) *The Vertebrate Body*, 2nd ed., W. B. Saunders, Philadelphia.

## **Homeotic Genes**

Homeotic genes are required during the development of plants and animals to control the differentiation of repeated homologous structures, such as vertebrae or flower organs. [Mutations](#) in homeotic genes result in the transformation of one of these homologous structures into the likeness of another structure normally present in a different position. Although most homeotic genes encode [transcription factors](#) (see below), a homeotic gene should be considered as such using anatomical,

and not molecular, criteria.

The term “homeosis” was first coined to describe certain spontaneous aberrations seen in the wild in which one part of a homologous series is transformed into the likeness of another (1). In plants, it is frequent to observe these transformations between the different organs of the flower (petals transformed into stamens). In animals, these transformations can happen between appendices (antenna into leg transformations in insects), parts of a segment (a lumbar vertebra into a thoracic one in vertebrates), or result in the formation of supernumerary organs in more anterior or posterior positions (extra mammary glands). Bateson also included under the term homeosis certain bilateral transformations in which both primordia develop in an animal in which normally only the left or the right primordia gives rise to the adult structure (tusk of the Narwhal; ovary and oviduct of fowl). The first homeotic mutant in animals, *bithorax*, was described in the fruit fly *Drosophila melanogaster* by Bridges and Morgan in 1923. In bithorax mutant flies, part of the metathorax transforms into mesothorax (halteres into wing). Mutations in many genes resulting in homeotic transformations of organs or segments have been isolated in plants, insects, and vertebrates. All these mutant transformations fit within Bateson's anatomical definition of homeosis, and, therefore, the genes that mutate to give such phenotypes are termed “homeotic genes.”

In animals, the concept of a homeotic gene is associated historically with that of homeobox and Hox gene from which it should be distinguished. The association originates from the early studies in *Drosophila* of the molecular nature of the homeotic genes. Genetic analysis had showed that some homeotic genes are in different locations in the [genome](#), but many cluster in two complexes: the Antennapedia complex contains homeotic genes required for the development of the head and the anterior thorax, and the Bithorax complex contains genes required for the development of more posterior segments. Molecular analysis of the homeotic genes of the Bithorax and Antennapedia complexes showed that they encode transcription factors with a common protein domain (2, 3). Because this new protein domain was first found in homeotic genes and was common to all of them, it was called the homeodomain, and the DNA sequence encoding it was called homeobox. Despite this historical link, not all homeotic genes encode homeodomain proteins and vice versa. In animals there are homeobox genes required for segmentation or dorsal-ventral specification that, when mutated, do not result in homeotic transformations. In plants, this lack of correlation is more evident as most homeotic genes are encoded by transcription factors of the **MADS** domain class, and the mutation of homeoproteins like “knotted” does not result in homeotic transformations (4).

According to their similarity, homeodomain sequences can be classified into at least 30 classes (5). The homeodomains encoded by the genes of the Antennapedia and Bithorax complexes are most related by sequence, and homologues of these genes have been found not only in vertebrates and arthropods but also in unsegmented worms like *Caenorhabditis elegans* (6, 7). These evolutionarily related genes are required for the specification of structures along the anterior–posterior axis and have in common the property of being organized in clusters. To reflect this common function and evolutionary origin, the term “Hox genes” was coined. Hox genes are expressed in the anterior-posterior axis of the organism in a collinear order with their location in the cluster. Genes located 3' in the cluster are expressed in more anterior positions than those located more 5'. Hox genes are homeotic genes; however, certain mutations in Hox genes affect cell properties like migration or cell differentiation, without resulting in clear homeotic transformations.

In plants, the ABC model explains satisfactorily how the homeotic genes form the different organs of the flower (8). According to this model, there are three classes of homeotic genes (termed A, B, and C) acting in combination to form the four flower organs: petals, sepals, stamens, and carpels. These organs are arranged in concentric rings known as whorls. The first (outer) whorl is formed by sepals that express class A genes during development; the second whorl is formed by petals that express A and B genes; the third is formed by stamens that express B and C genes; and the fourth (innermost) is formed by carpels that express C genes. This restricted spatial expression can be represented by the following formula: 1A, 2AB, 3BC, 4C, where numbers represent whorls and letters represent the

gene class expressed in those whorls. The ABC model also proposes that the A and C genes repress each other; so, in a class A mutant, C genes are expressed in the region where A genes are normally expressed, and vice versa. In contrast, the mutation of B genes does not affect the spatial expression of A or C genes. With this premise, most of the available expression and mutant data can be explained. In an A mutation C genes will be expressed in all four whorls (1C, 2BC, 3BC, 4C), resulting in a flower expressing C in the outer whorl, which consequently develops as carpels; expressing BC in the second and third whorls, which develop as stamens; and expressing only C in the fourth whorl, which then develop as carpels. With the same logic, a mutation in C leads to a flower composed of sepals, petals, petals, sepals (1A, 2AB, 3AB, 4A). Class B mutants lead to flowers composed of sepals, sepals, carpels, carpels (1A, 2A, 3C, 4C). A triple mutant lacking one gene of each class lacks all the floral organs, and the whorls develop as leaves.

A common characteristic of all homeotic genes is their capability to organize the development of entire segments or structures. To reflect this property, the homeotic genes have been named “selector genes” (9) (and also “identity genes”) as they are high in a genetic hierarchy and can “select” what kind of organ is formed in a certain position. Homeotic genes have this property because they control groups of downstream genes, which are ultimately responsible for the shape of the organs by controlling cell behaviors like cell division, adhesion, etc. In *Drosophila*, where more are known, the downstream genes encode diverse proteins, such as signaling molecules, adhesion molecules, transcription factors, etc. (10).

One property that is essential for a homeotic gene is that it is active only in a subset of the homologous series of organs. Where a certain homeotic gene is active is controlled in both plants and animals at the transcriptional level. Only the organ primordia in which a particular homeotic gene is expressed will have the set of characteristics that this gene confers. In a loss of function mutant for this homeotic gene, the characteristics of the organ are replaced with those of another homologous organ. On the other hand, if the function of a homeotic gene is activated in the primordia of a homologous organ that normally does not express it, this organ will have a homeotic transformation.

There is a complex system involved in activating, modulating, and maintaining homeotic gene expression. The best known example is the regulation of Hox gene expression. In *Drosophila* it is the segmentation genes that activate spatially restricted patterns of Hox gene expression (11). Vertebrate Hox gene expression is thought to be initiated by the retinoic acid morphogen (12). Recent studies show that besides the Retinoic acid receptors, the vertebrate caudal homologs are activators of Hox expression in frogs and mice (13, 14). There is a group of genes have been well characterized in both vertebrates and invertebrates that are responsible for the maintenance of Hox gene expression. These are mainly positive regulators, the *Trithorax* genes, and negative regulators, the *Polycomb* genes. Mutations in the *Trithorax* and *Polycomb* genes result in homeotic transformations; therefore, they must also be considered homeotic genes. In plants, some genes required for the formation of the floral meristem act as early activators of flower homeotic genes (15, 16). Interestingly, a negative homeotic gene regulator has been isolated in plants that is homologous to a *Polycomb* gene (17).

Most homeotic genes studied to date encode transcription factors, but, as homeosis is an anatomical concept, it is possible that other classes of genes result in homeotic transformations. In fact, the transformations of wing toward notum observed in mutants for the signalling gene *wingless* have been considered a homeotic transformation (18).

## Bibliography

1. W. Bateson (1894) MacMillan & Co, London.
2. M. P. Scott and A. Weiner (1984) Proc. Natl. Acad. Sci. U.S.A. **81**, 4115–4119.
3. W. McGinnis, *et al.* (1984) Nature **308**, 428–433.
4. J. A. Long, *et al.* (1996) Nature **379**, 66–69.

5. C. Kappen, K. Schughart, and F. H. Ruddle (1993) *Genomics* **18**, 54–70.
6. C. Kenyon (1994) *Cell* **78**, 175–180.
7. R. Krumlauf (1994) *Cell* **78**, 191–201.
8. D. Weigel and E. M. Meyerowitz (1994) *Cell* **78**, 203–209.
9. A. García-Bellido (1975) CIBA Foundation Symposium ed. *Cell Patterning*, Vol. **29**, Elsevier, Amsterdam. 161–182.
10. Y. Graba, D. Aragnol, and J. Pradel (1997) *BioEssays* **19**, 379–388.
11. M. Bienz and J. Müller (1995) *BioEssays* **17**, 775–784.
12. H. Marshall, *et al.* (1996) *FASEB J.* **10**, 969–978.
13. H. V. Isaacs, M. E. Pownall, and J. M. Slack (1998) *Embo J.* **17**, 3413–27.
14. M. Houle, *et al.* (2000) *Mol. Cell Biol.* **20**, 6579–86.
15. D. Wagner, R. W. Sablowski, and E. M. Meyerowitz (1999) *Science* **285**, 582–4.
16. M. A. Busch, K. Bombliès, and D. Weigel (1999) *Science* **285**, 585–7.
17. J. Goodrich, *et al.* (1997) *Nature* **386**, 44–51.
18. G. Morata and P. A. Lawrence (1977) *Dev. Biol.* **56**, 227–240.

### Suggestions for Further Reading

19. S. B. Carroll (1995) Homeotic genes and the evolution of arthropods and chordates. *Nature*, **376**, 479–485.
20. Y. Graba, D. Aragnol, and J. Pradel (1997). *Drosophila* Hox complex downstream targets and the function of homeotic genes. *BioEssays* **19**, 379–388.
21. R. Krumlauf (1994) *Hox* genes in vertebrate development. *Cell* **78**, 191–201.
22. D. Weigel and E. M. Meyerowitz (1994) The ABCs of floral homeotic genes. *Cell* **78**, 203–209.

### Homokaryon

A *homokaryon* is a hyphal cell, mycelium, organism or **spore of fungi** in which all the **nuclei** are genetically identical. The alternative situation is a [heterokaryon](#).

### Homological Modeling

[Proteins](#) that have similar [amino acid](#) sequences, or [primary structures](#), adopt the same fold, or [conformation](#) of the **polypeptide** backbone, similar [tertiary structures](#). In other words, the relationship of the three-dimensional [protein structure](#) to the amino acid sequence (see [Protein Structure Prediction](#)) does not have a one-to-one correspondence, but one-to-many. Therefore, when the 3-D structure of one of the proteins encoded by a **gene family** is known, it can be assumed that all of the other **homologous** protein members of the family adopt essentially the same fold. This is a strictly empirical observation that has held valid for natural proteins, and it provides the basis for homological modeling. It is not easy to give a precise numerical threshold, but two proteins with



sequences that are identical at 30% or more of the amino acid residues over a span of more than 100 residues, are almost certainly homologous to each other and belong to the same family (1). Again, this is an empirical rule deduced from natural proteins, which have evolved from a common **evolutionary** ancestor by the accumulation of individual [mutations](#), and it need not be applicable to *de novo* designed proteins. In fact, if the sequence is deliberately designed in one step, the fold of a protein can be converted from totally **alpha-helical** to [beta-sheet](#) while keeping the sequences 50% identical (2). There is, however, no guarantee that *de novo*-designed proteins will fold to a unique 3-D structure.

Homological modeling begins with alignment of a query sequence against the sequence of a protein homologue of known structure (see [Aligning Sequences](#)). The sequence alignment is best performed by a mathematical technique called dynamic programming. The two sequences aligned may contain insertions or deletions ([indels](#)) here and there, shown as gaps in one of the sequences. As the sequence similarity decreases, the number of gaps increases, and the entire alignment becomes less certain. If the structure is modeled according to an incorrect alignment, the resulting model will also be incorrect. Thus, a sequence identity of 40 to 50% or more is usually required for accurate homological modeling. Given the sequence alignment, the query amino acid sequence is mounted onto the known structure, which supplies a template backbone, and the necessary amino-acid side chains are replaced according to the sequence alignment.

Once a proper protein of known structure has been found for a query sequence, the main problems of homological modeling are twofold. The first is to fill in any missing polypeptide backbone by generating an appropriate loop structure. The other is to put all of the new side chains into the correct orientation. If the query structure has residues inserted, there is no template for that part of the sequence. The procedure for generating loops should generate additional polypeptide backbone that joins its two termini smoothly to the template structure and also has an energetically favorable conformation. Indels generally occur at the protein surface, so there are few interactions or steric hindrance to guide the structural design. The orientations of the new side chains of interior residues are determined by the packing of all of the atoms within the protein interior. A simple way to incorporate the new side chains is to adjust their conformation against the fixed conformations of nonsubstituted side chains and the polypeptide backbone. The conformations can also be selected from those observed most frequently in known protein structures, collected in “rotamer libraries,” and from those calculated to have the most favorable energies. A “dead-end elimination” algorithm is more advanced (3). More automatically, the [simulated annealing](#) method (4) can be applied to the entire model structure, allowing even the backbone conformation to vary, while seeking the energetically most stable and optimum conformation as a whole. Several computer packages for homological modeling are commercially available.

## Bibliography

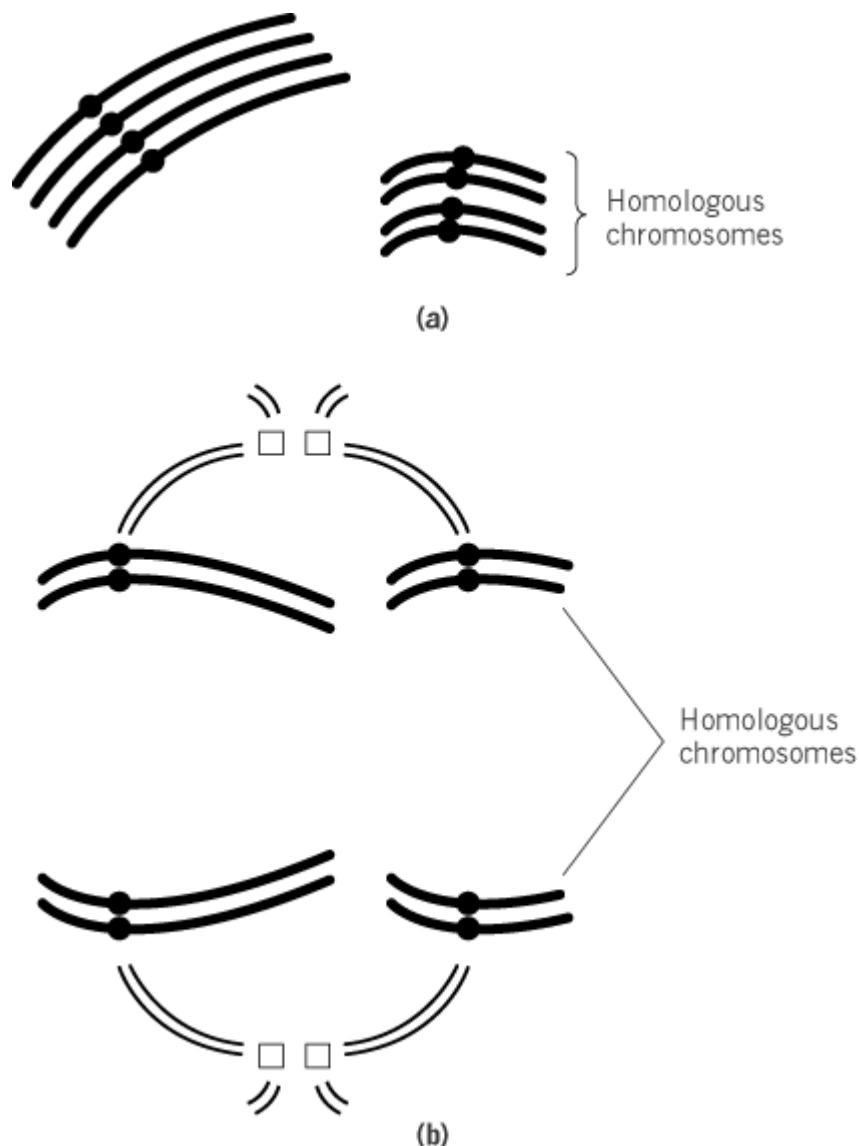
1. C. Sander and R. Schneider (1991) *Proteins: Struct Function Genet.* **9**, 56–68.
2. S. Dalal, S. Balasubramanian, and L. Regan (1997) *Nat. Struct. Biol.* **4**, 548–552.
3. A. Desmet, M. D. Maeyer, B. Hazes, and I. Lasters (1992) *Nature* **356**, 539–542.
4. A. Sali and T. L. Blundell (1990) *J. Mol. Biol.* **212**, 403–428.

## Homologous Chromosomes

Each **diploid** somatic [nucleus](#) contains two copies of each [chromosome](#), known as homologues. The existence of pairs of homologous chromosomes is important for providing at least two copies of any

particular gene. This provides the cell with two opportunities to generate a functional gene product. However, it also creates the problem of segregating the homologous chromosomes into the **gametes**. Homologous chromosomes always pair during the first meiotic cell division, when they have to be separated to help create **haploid** gametes (Fig. 1). In certain instances, homologous chromosomes also pair in the somatic **interphase**. This pairing is common in *Dipteran* insects and in some **plants** (1). The techniques of fluorescent *in situ* hybridization (FISH) and **confocal microscopy** (see **Denaturation Mapping**) allow for the definition of chromosomal territories within the nucleus. Homologous chromosomes frequently occupy adjacent territories in insects and plants (2). Nevertheless, homologues do not often pair in the interphase of mammalian cells. Homologous chromosomes are occasionally found in contact with the same nuclear structures, such as the **nucleolus** or the nuclear membrane, but the distances between homologues are quite variable (3).

**Figure 1.** Homologous chromosomes during meiosis. (a) Homologous chromosomes align during the pachytene of the first meiotic division. (b) Homologous chromosomes are separated during the anaphase of the first meiotic division.



Pairing of homologous chromosomes in **meiosis** occurs simultaneously with the beginning of chromosome condensation early in the **prophase** of the first meiotic division. In *Coprinus*, a **fungus**, the interaction between homologues, known as **synapsis** begins on the longitudinally aligned

chromosomes at points towards the ends of the chromosomes and moves towards the [centromere](#). It is probable that there are specific sites within the chromosome that have evolved to mediate chromosomal pairing (4). How homologous sequences are positioned next to each other during this chromosomal alignment process is also unknown, but recognition of similar [nucleoprotein](#) structures is likely to be involved (see [Chromocenter](#)). However [heterochromatin](#) association does not have a major role in this process, nor do the specialized chromosomal structures found at centromeres and [telomeres](#). The end result of chromosomal adjustment is stable physical linkage of the homologues by the [synaptonemal complex](#) (SC) (5). This is a ribbonlike structure consisting of two electron-dense lateral elements (25 nm wide) that are separated by a less dense central region (105 nm wide). The major components of the SC are synthesized during the meiotic prophase. Most of the chromatin is located outside of the SC in loops attached to the lateral elements. The SC has important functions in providing a structural framework for efficient [recombination](#), thereby stabilizing the crossover sites where chromatin bridges form at **chiasmata** (sites of **crossing-over** between homologous DNA sequences). SCs might also introduce constraints on the formation of chiasmata so that they cannot occur too close together. Visualization of SCs under the **electron microscope** after special staining procedures is a useful for revealing any misalignments of chromosomes due to deletions and duplications. This is a useful adjunct to conventional analysis of chromosomal bands.

### Bibliography

1. B. John (1976) *Chromosoma* **46**, 279–286.
2. L. Avivi and M. Feldman (1989) *Hum. Genet.* **55**, 281–285.
3. L. Manmelides (1984) *Proc. Natl. Acad. Sci. USA* **181**, 3123–3128.
4. R. S. Hawley (1980) *Genetics* **94**, 625–637.
5. C. B. Gillies (1984) *CRC Crit. Rev. Plant Sci.* **2**, 81–116.

### Suggestion for Further Reading

6. R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A Synthesis* Wiley-Liss, New York.

## Homologous Recombination

Homologous, or generalized, recombination requires extensive nucleotide sequence [homology](#) between the interacting [chromosomes](#) or chromosomal segments. Recombination occurs anywhere within the homologous region, although the frequency varies throughout the region. Near hot spots of recombination, such as the **Chi sites** of *Escherichia coli*, recombination occurs at higher than average frequency. Cold spots, near which recombination occurs at lower than average frequency, are also known.

Homologous recombination plays multiple biological roles. It was first recognized by its production of new combinations of **alleles** from genetically marked parental chromosomes. Later, it was recognized that it plays an important role in the segregation of replicated homologues in the first meiotic division (without further replication sister chromatids subsequently segregate in the second meiotic division). Finally, homologous recombination is the principal method by which double-strand DNA breaks are repaired with fidelity, at least in microorganisms (see [DNA Repair](#)).

Homologous recombination can be classified in numerous ways. One way distinguishes reciprocal recombination (crossing over) from nonreciprocal recombination (gene conversion). Reciprocal

recombinant types are produced in a single event in the former, but not in the latter. A strict classification requires assessing all of the parental chromosomes entering the event and all of the progeny chromosomes emerging from it. These are most readily assessed in meiosis of certain **fungi** (ascomycetes) that enclose the four products of meiosis, spores, in a sac called an ascus.

Consider a cross between marked parental chromosomes  $abc$  and  $+++$ . After [DNA replication](#), recombination may occur between a nonsister pair before the first meiotic division. After the two meiotic divisions, four meiotic products emerge.

In gene conversion, an  $a+c$  recombinant chromosome is produced but the reciprocal type,  $+b+$ , is not. In a crossover, both an  $a++$  recombinant chromosome and its reciprocal type  $+bc$  are produced. Both a gene conversion and a crossover may occur in the same event. In a population of meioses, statistical association of the two events is commonly observed; that is, among the meioses that have a gene conversion at  $b$ , a higher frequency of crossing over occurs in the  $ab$  or  $bc$  interval than in the population as a whole.

Many models have been proposed for homologous recombination. Most models attempt to account for gene conversion and crossing over as alternative outcomes of a common initial event. Common themes in most models of homologous recombination are (1) an initiating break, either single-stranded or double-stranded, in one parental chromosome; (2) the formation of a single-strand end on that chromosome; (3) invasion of the single-stranded tail into the homologous chromosome and the displacement of the identical (or nearly identical) single strand, to form a displacement or D-loop; (4) formation of a more complex joint molecule, such as a Holliday junction; and (5) resolution of this intermediate into one or two progeny chromosomes.

The joint molecule intermediates contain hybrid DNA; that is, DNA with one strand from each parent. If the parental DNA molecules differ genetically in this region, the hybrid DNA is [heteroduplex](#); that is, it contains a noncomplementary nucleotide sequence. If the parental chromosomes differ at a single base pair, the heteroduplex contains a simple mismatch. [Mismatch repair](#) converts the heteroduplex to a homoduplex of one or the other parental type and produces a gene convertant. More complex differences between the parental DNA molecules, such as deletions, insertions, or multiple base-pair changes, also produce heteroduplex DNA acted upon by repair enzymes to produce gene convertants. In the absence of such mismatch repair, the heteroduplex DNA is replicated and produces a postmeiotic segregant, that is, a meiotic product that produces two genetic types at the first mitotic division.

D-loops may arise by a single-strand invading tail with a 3'-end, which primes DNA synthesis. Such synthesis is necessary to repair a double-strand break that has been enlarged into a double-strand gap by loss of nucleotides. If the double-strand gap occurs at a site of genetic difference between the parental chromosomes, DNA synthesis templated by the intact chromosome repairs the double-strand gap and produces a gene convertant without mismatch repair.

In many models of homologous recombination, resolution of a Holliday junction in alternate ways gives crossing over (or not) with respect to markers flanking the site of hybrid DNA formation. This feature accounts for the statistical association of crossing over and gene conversion, yet allows for gene conversion without crossing over. These two types of events may, however, sometimes occur by distinct mechanisms. For example, a single-strand invasion without formation of a Holliday junction could readily give a gene convertant, but not a crossover.

Models of homologous recombination can be classified as follows: (1) In breakage-and-reunion there is no loss or resynthesis of nucleotides except for that associated with mismatch repair. This mechanism is exemplified by the Holliday model and, with the addition of limited DNA synthesis, the double-strand gap repair model. (2) In break-copy mechanisms, one parental DNA is broken, and one of its ends primes DNA synthesis (copying) using the other parental DNA as template. This mechanism occurs in recombination of bacteriophages, bacteria, and some types of mitotic double-

strand break repair. (3) In copy-choice recombination the initial stage of DNA replication is templated by one parental chromosome, but then replication switches to the other parental DNA at a homologous point. This type of recombination may occur in RNA viruses.

Another type of homologous recombination, single-strand annealing, occurs in the special cases of deletion between direct repeats in a chromosome or in other instances when DNA can be lost, for example, in phage recombination. In this case, a double-strand break may occur between the direct repeats. Resection of one strand at each end exposes complementary single strands that anneal and produce hybrid DNA with one strand from each copy of the repeat. Appropriate trimming of single-strand ends, filling of gaps, and ligations produce a recombinant. This occurs between homologous regions on separate chromosomes if two initiating breaks occur at appropriate places.

Single-strand annealing is an example of a nonconservative recombination event, that is, when two homologous segments or chromosomes interact to produce only one segment or chromosome. In conservative recombination, exemplified by breakage-and-reunion mechanisms, two parental chromosomes produce two progeny chromosomes, one or both of which is recombinant.

Homologous recombination is a complex process promoted, in the cases studied, by multiple enzymes. A central activity, exemplified by the RecA protein of *E. coli*, promotes the association of homologous nucleotide sequences. RecA protein promotes the formation of D-loops from single-strand tails and a homologous duplex. Proteins of related amino acid sequence and function are widespread in many organisms, ranging from bacteria to humans. Other enzymes prepare substrates for RecA protein. For example, the RecBCD enzyme of *E. coli* generates a single-strand intermediate from linear duplex DNA. Other enzymes resolve recombination intermediates into products. For example, the RuvC protein of *E. coli* cleaves a Holliday junction near the exchange point into two separate duplexes. Their nicks are ligated by DNA ligase to produce intact, possibly recombinant DNA molecules.

#### Suggestions for Further Reading

1. R. Kucherlapati and G. R. Smith (eds.) (1988) *Genetic Recombination*, American Society for Microbiology, Washington, D.C.
2. D. F. R. Leach (1996) *Genetic Recombination*, Blackwell Science, Oxford, England.

## Homology

Homology is a similarity due to a shared common ancestor in . Therefore, homology is the result of ; in contrast, is a consequence of . Originally, homology was used for morphological characters, such as organs. For example, homologous organs were defined as organs that are related to each other through a common descent, even though now they perhaps exhibit different functions. At the present time, however, the term is used in relation to molecular traits, in particular, **nucleotide sequences** and amino acid sequences.

Any two sequences can be compared by alignment, where maximum similarity is used as the basis of measurement (see ). Once the alignment is made, the proportion of identical nucleotides, site by site, can be computed. This proportion is the simplest measure for the “degree of sequence similarity.” When the proportion is statistically significant, it is reasonable to say that the two sequences being compared are homologous, because the probability that the similarity was derived by mere chance is, in most cases, extremely small. In general, if the proportion of identical sites is at least 25–30% in an

amino acid sequence, or over 40% in a nucleotide sequence, the two sequences may be said to be homologous. This is because there is an extremely high possibility that this similarity is due to a common ancestor and a very low probability that it occurred randomly. The degree of sequence similarity can then be called the “degree of homology.”

By using sequence similarity, one can search for homologous sequences in the nucleotide and amino acid . A given sequence is used as a query sequence and is compared with a sequence in an entry of the database by sliding one site in each comparison, until all possible comparisons have been made. This procedure is repeated for all entries of the database. Only the entries where the sequence homology is statistically significant to the query sequence are noted. This search is called a “homology search.” In practice, algorithms called “dynamic programming” are used to make the homology search more efficient.

The homology search is particularly useful for predicting the function of a newly identified gene or protein. This is based on the following logic: In general, conservation of a sequence is known to be stronger for functionally important regions of a protein. If we search a sequence and find a region where conservation is strong, it is reasonable to conclude that the region is functionally important and that the function is possibly shared by the homologous sequences. Doolittle constructed his own public database of amino acid sequences, which were deduced from nucleotide sequences (1). He conducted a primitive homology search by using a viral **oncogene** (*sis*) as a query sequence and then found homology with human (Figure 1). He found that the homology was very high, and for this reason he proposed that the *sis* gene may have a function as a . Following Doolittle's prediction, experiments confirmed the hypothesis (2). In fact, it was found that the *sis* gene was the host gene that was inserted into the viral genome. This was the first report in which a homology search was successful for predicting biological function of an unknown sequence. Since then, the homology search has become a very popular tool in molecular biology.

**Figure 1.** Sequence homology between mouse PDGF gene and viral *sis* oncogene. Identical amino acids are indicated by similar amino acids by a period.

```

PDGF (mouse) : MNRCWALFLPLCCYLRLVSAEGDPIPEELYEMLSDHSIRSFDLQRLLRD:
               . . . . . : . . . . . :
V-sis       :          MTLTWQGDPIPEELYKMLSGHSIRSFDLQRLRQGD:

PDGF (mouse) : DLNMTRAHSGVELESSSRGRRSLGSLAAAEPAVIAECKTRTEVFQISRNLII
               . . . . . : . . . . . :
V-sis:       DLNMTRSHSGGELESLARGKRSLGSLSVAEPAMIAECKTRTEVFEIS----

```

### Bibliography

“Homology” in , Vol. 2, pp. 1165–1166, by T. Gojobori; “Homology” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. R. F. Doolittle (1981) *Science* **214**, 149–159.
2. T. F. Deuel, J. S. Huang, S. S. Huang, P. Stroobant, and M. D. Waterfield (1983) *Science* **221**, 1348–1350.

## Homozygote

The **phenotype** of a **diploid** cell may depend upon the interaction of two **alleles** of a gene because a **recessive** gene can express itself only in the absence of its **dominant** allele (note, however, that some genes are neither dominant or recessive). This brings us to two very important words in the language of genetics, *homozygote* and **heterozygote** that describe the general genetic constitution of diploid cells.

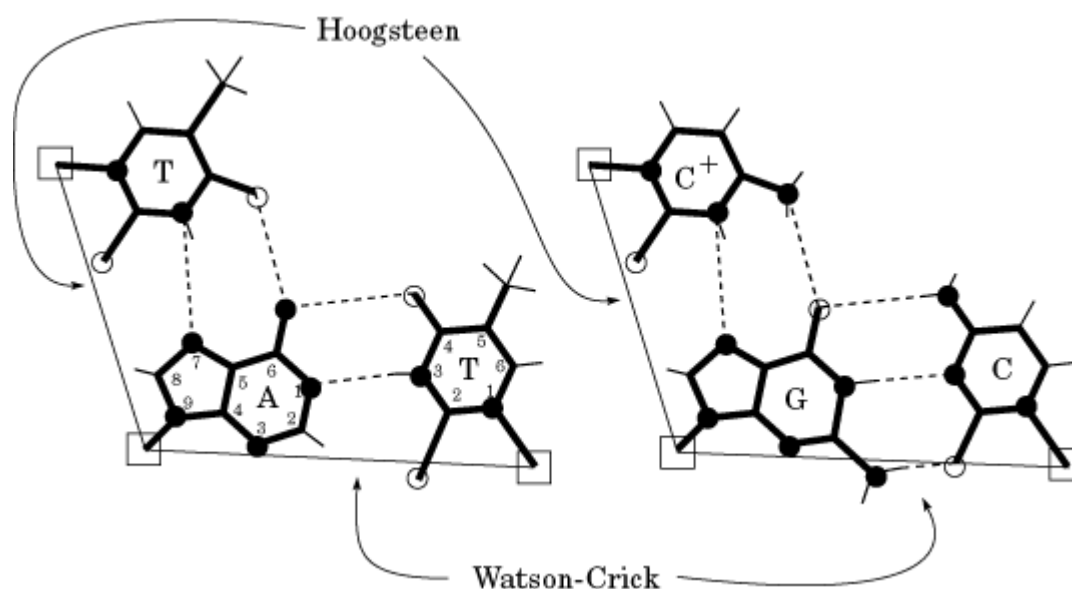
A homozygote is a [zygote](#), or other diploid cell, derived from the union of two **gametes** of identical **genotype**. In practice, however, homozygosity requires evidence from genetic crosses that a particular pair of genes is identical, and it is obviously impossible to prove that *all* the thousands of pairs of genes in what appears to be a homozygous individual are in fact identical. So the words homozygote and heterozygote have a narrower connotation and denote the relationship between specific [homologous chromosome](#) regions or genes.

We say that a diploid cell is homozygous for gene *A*, which means that it has the genetic constitution either *AA* or *aa*, or, as it is usually written, *A/A* and *a/a*, where *A* and *a* are two distinguishable alleles. A cell may be homozygous for gene *A* but heterozygous for gene *B* (*AB/Ab*).

## Hoogsteen Base Pairing

Nucleic acid bases can form a variety of base pairs. In all canonical duplex structures including B-DNA, A-DNA, and Z-DNA, the so-called Watson–Crick base pairs are found. The Watson–Crick base-pairings of guanine with cytosine and adenine with thymine are shown in Figure 1. The G to C and A to T complementarities are provided by specific hydrogen bond donors and acceptors between the bases. In a G:C base pair three hydrogen bonds are found between G-O6 and C-N4, G-N1 and C-N3, and G-N2 and C-O2 positions, and in an A:T base pair two hydrogen bonds are found between A-N6 and T-O4, and A-N1 and T-N3 positions. The two C1' atoms within a G:C base pair and an A:T base pair are equidistant (~10.5 Å). Thus the G:C and A:T base pairs in the Watson–Crick conformation are nearly iso-structural, from the point of the view of the sugar-phosphate backbone.

**Figure 1.** Schematic illustration of the A:T and G:C Watson–Crick and Hoogsteen base pairs. Hydrogen bonds are shown as dashed lines, and the distance between two C1' atoms (shown as open squares) is drawn with a thin line. The numbering system of the purine and pyrimidine rings is shown.



In earlier crystallographic analyses of nucleic acid base pairs, it was found that an A:T base pair can also adopt a different conformation. Instead of using the N1 position of adenine to base pair with N3 of thymine, the N7 position of adenine was used. This type of A:T base pair was named Hoogsteen base pair (1). A striking difference between those two types of base pairs is that the C1'–C1' distance in the Hoogsteen base pair of 8.65 Å is significantly shorter than that of the Watson–Crick base pair of 10.5 Å. Therefore the Hoogsteen base pair is not compatible structurally with the Watson–Crick base pair in B-DNA. The incorporation of a Hoogsteen base pair in B-DNA destabilizes the duplex structure. A similar G:C Hoogsteen base pair is not stable unless the C is protonated at the N3 position.

However, certain chemical modifications may enhance the stability of the Hoogsteen base pair. It was shown that the modified nucleoside 3-isodeoxyadenosine (iA) forms a stable base pair with thymine (T) using the Hoogsteen conformation, and the iA:T base pair is fully compatible with the normal Watson–Crick base pair as evident from the stable duplex of d(CG[iA]TCG)<sub>2</sub> shown by NMR analysis (2). Another modification which stabilizes the Hoogsteen base pair involves the use of 8-amino-adenine in which the amino group at the C8 position can form an additional hydrogen bond with the O2 of thymine (3). Those modified bases may find applications when the Hoogsteen base pairing is needed in nucleic acid structures.

The binding of protein or drug to DNA may induce alternative base pairs, including Hoogsteen base pair, to form. It has been shown that the quinoxaline bis-intercalator antibiotics triostine A and echinomycin bind to -XCGY- sequences with the drug bracketing the CG dinucleotide base pairs. In the structure of the triostine-CGTACG and triostine-GCGTACGCT complexes, the A:T and G:C<sup>+</sup> base pairs adjacent to the quinoxaline rings adopt Hoogsteen base pairs (4). The molecular basis for the Hoogsteen base pair in the structure is to improve the stacking interactions between the quinoxaline ring and the adjacent base pairs. It is likely that such a conformational rearrangement from the Watson–Crick to Hoogsteen base pairs may occur in other biological systems.

An important role for the A:T and G:C<sup>+</sup> Hoogsteen base pairs is found in the structure of a nucleic acid triple helix. The formation of a poly(A).poly(U).poly(U) triplex requires that the third strand of poly(U) is bound to the poly(A) strand of the poly(A).poly(U) duplex using Hoogsteen base pairing. The T.A.T. and the C.G.C<sup>+</sup> triple base pairs are shown in Figure 1. (Note that the triple C.G.C. base pair is not stable unless the second C is protonated.)



It should be pointed out that Hoogsteen base pairing is not restricted in DNA. In transfer RNA, the nucleic acid bases are often modified and involved in tertiary interactions. For example, in yeast tRNA<sup>Phe</sup> the nucleotide at position 58 is 1-methyl-adenosine, which is base paired with T54 (5). Because the N1 position of this adenine is methylated, it cannot participate in the normal Watson–Crick base pair. Instead m<sup>1</sup>A58 and T54 form a reverse Hoogsteen base pair. It is likely that those noncanonical base pairs, including Hoogsteen and reverse Hoogsteen base pairs, will be found in higher-ordered RNA structures.

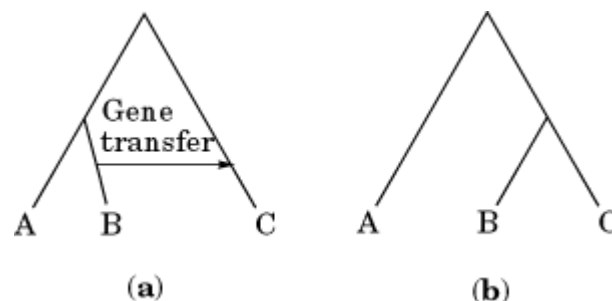
### Bibliography

1. K. Hoogsteen (1963) *Acta Crystallogr.* **16**, 907–916.
2. B. Bhat, Neelma, N. J. Leonard, H. Robinson, and A. H.-J. Wang (1996) *J. Am. Chem. Soc.* **118**, 3065–3066.
3. K. Kawai, I. Saito, and H. Sugiyama (1998) *Tetrahedron Lett.* **39**, 5221–5224.
4. A. H.-J. Wang et al. (1984) *Science* **225**, 1115–1121.
5. S. H. Kim et al. (1974) *Science* **185**, 435–440.

### Horizontal Gene Transfer

In general, **genes** are transferred from one generation to the next in a vertical fashion, from parent to progeny through the gametes donated from each parent to form the zygote from which the offspring develops (see [Vertical Gene Transfer](#)). In horizontal gene transfer, a gene is transferred from one species to another, often with a **virus** or bacterium acting as the vehicle (Fig. 1). In vertical gene transfer, the genes must originate in germline cells in order to be transferred. In horizontal gene transfer, the gene may originate in a somatic cell, but the receiving cell must be a germline cell in order for the transferred gene to be passed on to the next generation. Strictly speaking, transfers of genes between somatic cells of two species are not called horizontal gene transfer.

**Figure 1.** Illustration of horizontal gene transfer and how it affects the (a) species and (b) gene phylogenetic trees.



Horizontal gene transfer is often difficult to detect. It is often suspected, however, when there is marked discontinuity in the [phylogenetic tree](#) for a gene, such as the appearance of a gene in one species, but not in any of its most closely related species. Discrepancies between the phylogenetic trees for species and for genes can also indicate that horizontal gene transfer may have occurred.

The **nucleotide sequences** of entire [genomes](#) show that horizontal gene transfers must have occurred often. When the complete genome sequences of three representative species from the three major organismic worlds, namely, *Escherichia coli* (Eubacterium), *Methanococcus jannaschi* (Archaeobacterium), and *Saccharomyces cerevisiae* (Eukaryote) were compared with each other, approximately 60% of all the genes of the archaeobacterial genome are closer to those of the eubacterium than to those of eukaryotes, whereas about 30% of the archaeobacterial genes were closer to eukaryotes. The remaining 10% of the genes are archaeobacteria-specific (1). Therefore, the archaeobacterial genome looks like a mosaic structure, a mixture of genetic material from various genomes. At present, only horizontal gene transfer can explain the mosaic features of the genomes. It implies that a tremendous amount of horizontal gene transfer should have taken place during long-term evolution. Genomes appear to have a robustness that allows major changes in gene order and arrangement without loss of function. This provides an environment in which horizontal gene transfer can occur and provides a source for the emergence of new gene systems and new gene interactions. Thus, the unit of function appears to be the gene, and the positions of genes in the genome have little effect on function. There appear to be only rare examples in which a set of genes must work as a unit. This elasticity of the genome allows it to absorb any disturbances caused by horizontal gene transfer, and it demonstrates the evolutionary significance of the genome itself.

#### Bibliography

1. H. Watanabe, T. Gojobori, and M. Kin-Ichiro, (1997) *Gene* **205**, 7–18.

## Hormone Receptors

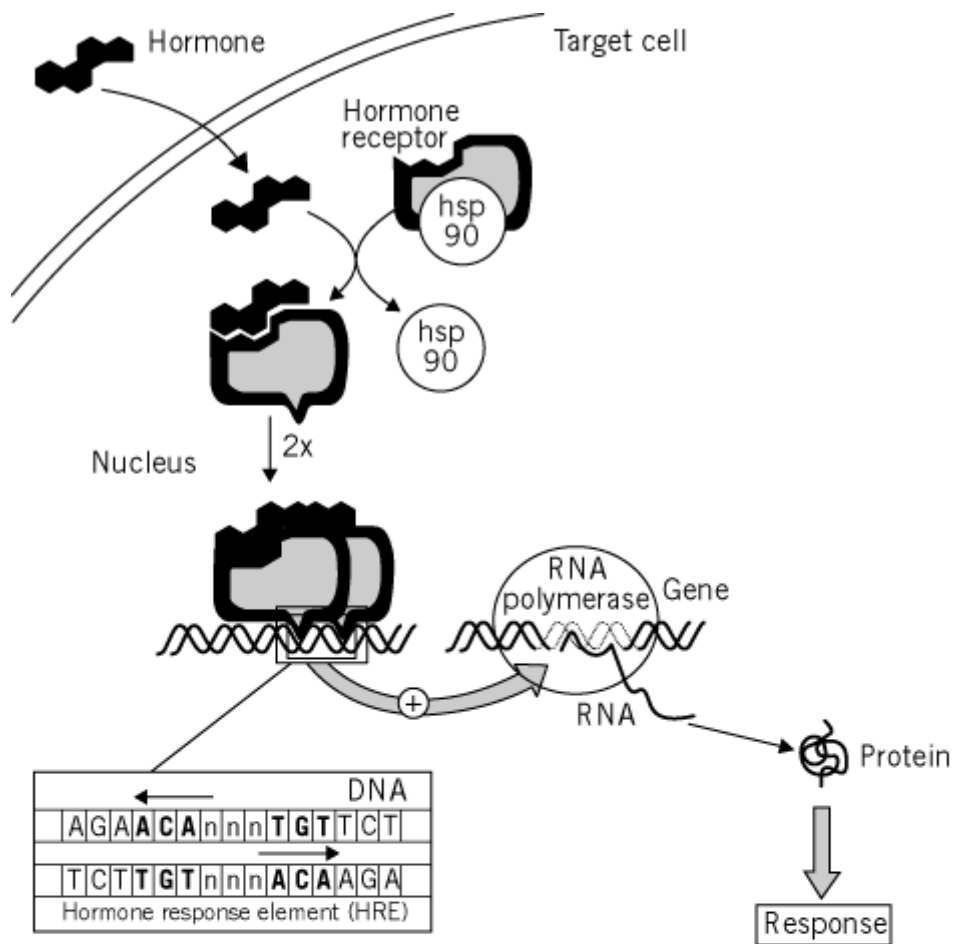
[Hormones](#) regulate a wide variety of physiological functions, encompassing intermediary metabolism, growth, and cell differentiation. In contrast to the diversity of their phenotypic effects, they have only two fundamental mechanisms of action, depending on their physical chemical characteristics. The lipophilic [steroid hormones](#) and [thyroid hormones](#) (vitamin D<sub>2</sub> and [retinoic acid](#) have similar modes of action) are **hydrophobic** and primarily act intracellularly, modulating gene [transcription](#), whereas the [peptide hormones](#), adrenaline, and melatonin are [hydrophilic](#) and act primarily at the cell [membrane](#), triggering a cascade of [signal transduction](#) events leading to intracellular regulatory effects. The first step, and a prerequisite for the elicitation of the hormonal effect of both lipophilic or hydrophilic hormones, is binding of the hormone to its appropriate **receptor** protein. The [hormone receptors](#) have the following characteristics: (a) high specificity, so that the receptors can recognize their specific ligand and discriminate between the incoming signals; (b) high affinity for their ligand ( $K_D$  values in the range of  $10^{-7}$  to  $10^{-10}$ M); and (c) limited capacity (no more than a few thousand receptor molecules per cell).

### 1. Lipophilic Hormones

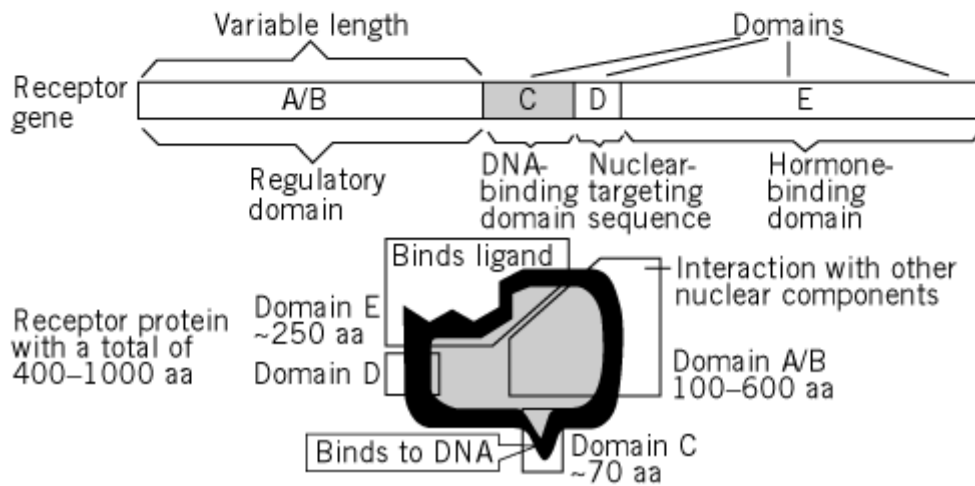
The hydrophobic hormones appear to diffuse across the membrane into the cell (although there are claims of a specific transport mechanism) and bind with high affinity and specificity to their respective receptor protein (see [Estrogen Receptors](#)). These receptors, also known as nuclear receptors, belong to a large family of proteins with similar structures, characterized by various functional **domains**, through which the biological activity of the molecule is expressed (Fig. 1). The central DNA-binding domain is responsible for recognition and specific interaction with specific **palindromic** DNA sequences, the [hormone response elements](#). Characteristic of DNA-binding

domains are those comprised of two **zinc finger** structures: One is on the amino-terminal side making specific contact with the DNA major groove, and the other is responsible for the specific homodimerization of the activated receptor. Other important domains are (a) the ligand-binding domain, which is at the carboxy terminus of the receptor, (b) [nuclear import, export](#) motifs, (c) **heat-shock** protein binding sites, and (d) **transactivation** domains, which can be at either the *N*- or *C*-terminal regions of the receptor. In the absence of hormone, the receptor is complexed to other proteins, such as the heat shock protein hsp90, [immunophilins](#), and smaller peptides. Upon hormone binding, the complex dissociates and the receptor protein is released and forms homodimers, which then bind to the hormone response elements (HREs) (see [Hormone Response Elements, Glucocorticoid Response Element](#)). The HREs are [enhancer](#) elements that are activated by the ligand–receptor complex. The result of this binding is stimulation (or, in some cases, inhibition) of [transcription](#) of the appropriate gene. The mechanism is not yet elucidated, but it also involves interaction of the hormone–receptor complex with other nuclear proteins, such as [transcription factors](#), “adaptor proteins,” “coactivators,” and transacetylases (1-4). The primary transcript is subsequently processed (see [RNA Splicing](#)), and the resulting [messenger RNA](#) is **translated** into the respective protein. The nature of the protein determines the phenotypic expression of the hormone action: For example, if the protein is an [enzyme](#) involved in gluconeogenesis, the blood glucose level will change; if it is a [growth factor](#), cell replication will be affected.

**Figure 1.** Mechanism of action of lipophilic hormones. (a) Mechanism of action of lipophilic hormones. (b) Receptors of lipophilic hormones. [From Ref. [10](#), with permission (slightly modified).]



(a)



(b)

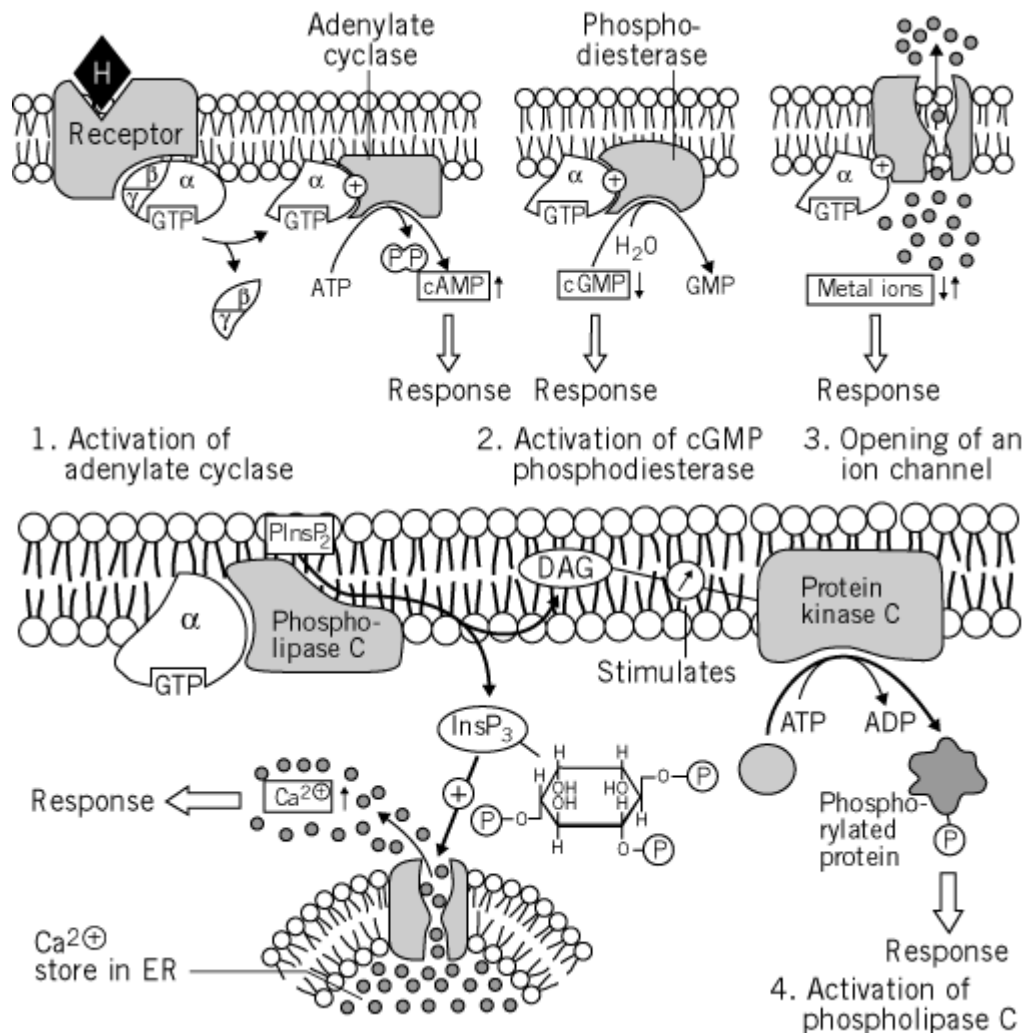
The effects of steroid hormones are usually long-term, needing hours to be exerted. It is becoming increasingly evident, however, that some steroid hormones act very rapidly, by binding to membranes and altering the intracellular concentrations of  $\text{Ca}^{2+}$ ,  $\text{Na}^+$ , and other ions (5). The nature of the membrane-binding sites (eg, whether they are parts of a pore complex) is not known.

## 2. Hydrophilic Hormones

The hydrophilic hormones are not able to penetrate the cell and instead interact with membrane

receptors (see [Hormone Receptors](#)) (Fig. 2). Two types of membrane receptors for hydrophilic hormones have been characterized, Type I and Type III. The third type, Type II, represents [ion channels](#), which are mainly regulated by neurotransmitters.

**Figure 2.** Mechanism of action of hydrophilic hormones. [From Ref. 10, with permission (slightly modified).]



### 2.1. Type III Receptors

These receptors are integral [membrane proteins](#), with seven **a-helical** domains spanning the membrane, an extracellular domain that is important for ligand recognition, and an intracellular one important for signal transduction. Binding of the cognate ligand induces a conformational change in the receptor that is further transmitted through a family of guanine nucleotide-binding proteins (**G proteins**) associated with the receptor (Fig. 2). The G proteins are trimers of α, β, and γ subunits. In the resting state, the α subunit has GDP bound. Interaction of the hormone with the receptor leads to exchange of GDP for GTP on the α subunit and dissociation of the G complex into a separate α subunit and a β/γ dimer. The free α subunit is the effector molecule that activates enzymes involved in the production of [second messengers](#), such as the [adenylate cyclases](#), which catalyze 3',5'-**cyclic AMP** synthesis from ATP, and **phospholipase C**, which catalyzes the production of inositol 1,4,5-trisphosphate (IP<sub>3</sub>) and **diacylglycerol** (DAG) from phosphatidylinositol bisphosphate. Cyclic AMP and DAG are **allosteric** modulators of protein kinases A and C, respectively, whereas the hydrophilic IP<sub>3</sub> reaches the endoplasmic reticulum and acts there by releasing Ca<sup>2+</sup> from storage

sites (6-8). Activation of the protein kinases leads to **phosphorylation** at serine/threonine residues of a large variety of proteins (see [Serine/Threonine Kinases and Phosphatases](#)), which are activated in their various functions, such as enzymes, structural proteins, and transcription factors.

The free  $\alpha$  subunit of the G protein, due to its intrinsic [GTPase](#) activity, hydrolyzes its bound GTP to GDP. When the hormonal stimulus ceases, the GDP is not replaced by GTP, and the G-protein trimer is reconstituted. Consequently, the system thus returns to its initial inactive condition.

## 2.2. Type I Receptors

These receptors are best exemplified by the [insulin](#) receptor. It is composed of two identical  $\alpha$  subunits that are extracellular and two identical  $\beta$  subunits, with extracellular, transmembrane, and intracellular domains. The intracellular domains possess intrinsic protein kinase activity. Binding of the cognate ligand induces a conformational change and dimerization of the receptor molecule. This leads to activation of the kinase activity of the receptor and phosphorylation of [tyrosine](#) residues within the kinase domain and at other sites of the receptor (see [Tyrosine Kinases and Phosphatases](#)). These last sites create docking sites for downstream signal transduction molecules containing **SH2 domains**, which are activated by a variety of mechanisms (eg, tyrosine phosphorylation, conformational changes). Subsequently, by way of intermediate second messenger molecules, they trigger metabolic processes, gene transcription, and [DNA replication](#). The importance of second messenger formation is that the original hormonal signal is amplified by two to three orders of magnitude and that multiple signal effects are integrated by using the same second messenger.

## 3. Medical Disorders

A series of endocrine diseases result from genetic defects of the receptors for both the lipophilic and the hydrophilic hormones. They are manifested as resistance syndromes, with elevated blood levels of the specific hormone, because the regulatory feedback inhibitory cycle is not functioning. Mutations have been found in various domains of the receptors, leading to partial or total loss of response to hormonal stimuli. In more rare cases, the receptors have lost their function, but in addition they inhibit the action of the wild-type receptor (dominant negative mutants). Table 1 shows an overview of diseases caused by mutations in genes coding for members of the nuclear hormone receptor family.

**Table 1. Some Diseases Caused by Mutations in Genes Coding for Members of the Nuclear Hormone Receptor Superfamily<sup>a</sup>**

| Mutant Receptors  | Disease                               | Remarks   |
|-------------------|---------------------------------------|---|
| Androgen          | Androgen insensitivity syndrome (AIS) | Complete AIS leads to a complete female phenotype.  |
| Glucocorticoid    | Glucocorticoid resistance syndrome    | Symptoms of fatigue, hirsutism, and hypertension.   |
| Mineralocorticoid | Mineralocorticoid resistance syndrome | Rare autosomal disease, with symptoms of hyponatraemia and dehydration.   |
| Estrogen          | Estrogen resistance syndrome          | Only one case with complete resistance has been reported. Symptoms of incomplete epiphyseal closure, prolonged linear growth. |
| Thyroid hormone   | Thyroid hormone                       | Caused by mutations in the $\beta$ -  |

|  |   |   |
|--|---|---|
|  | resistance syndrome                           | receptor gene. Goiter, learning disabilities, hearing defects, delayed bone maturation, mental retardation. |
| Vitamin D <sub>3</sub><br>(calcitriol) | Vitamin D <sub>3</sub><br>resistance syndrome | Severe rickets, secondary hypocalcemia, hyperparathyroidism.  |

---

<sup>a</sup> Ref. 9, with permission (slightly modified).

## Bibliography

1. U. Clever and P. Karlson (1960) *Exp. Cell Res.* **20**, 623–626.
2. M. Beato, P. Herrlich, and G. Schutz (1995) *Cell* **83**, 851–857.
3. D. M. Heery et al. (1997) *Nature* **387**, 733–736.
4. J. Iorchia et al. (1997) *Nature* **387**, 677–683.
5. M. Wehling (1994) *Annu. Rev. Physiol.* **59**, 365–393.
6. A. Hall (1990) *Science* **249**, 635–640.
7. W. J. Tang and A. G. Gilman (1992) *Cell* **70**, 869–872.
8. N. DiVecha and R. F. Irvine (1995) *Cell* **80**, 269–278.
9. A. Baniahmad, M. Eggert, and R. Renkawitz (1997) In *Transcription Factors in Eukaryotes* (A. G. Papavassiliou, ed.), Springer, Heidelberg.
10. J. Koolman and K.-H. Rohm (1996) *Color Atlas of Biochemistry*, Thieme, Stuttgart.

## Suggestions for Further Reading

11. R. E. J. Ribeiro, P. J. Kushner, and J. D. Baxter (1995) The nuclear hormone receptor gene superfamily. *Annu. Rev. Med.* **46**, 443–453.
12. M. G. Parker (1993). *Steroid Hormone Action*, IRL Press, Oxford, U.K.
13. R. White and M. G. Parker (1998) Molecular mechanisms of steroid hormone action. *Endocrine-Related Cancer*, **5**, 1–4.

## Hormone Response Elements

Lipophilic [hormones](#) are able to enter cells and exert their actions by binding to specific **receptors**. Binding of hormones of this type, such as the **steroids**, [retinoic acid](#), [thyroid hormone](#) and vitamin D<sub>3</sub>, results in the activation of their receptors, which then bind to specific [response elements](#) and activate gene [transcription](#). The detailed discussion of this process for one specific hormone and its response elements is provided in [Glucocorticoid response element](#). This article describes the relationships between the different response elements for these various hormones ([1](#)).

In each case, the particular response element confers a specific response to a particular hormone because the response element binds only the receptor for that hormone, and not other receptors ([2](#), [3](#)). Thus, the glucocorticoid response element (GRE) binds the glucocorticoid receptor–glucocorticoid hormone complex and hence mediates response to glucocorticoid hormone (Table [1](#),

part a). However, alteration of an A residue at the third position of the palindromic GRE sequence to a T residue converts the GRE into an ERE ([estrogen](#) response element), which binds the estrogen receptor–hormone complex and therefore activates gene expression in response to estrogen (Table 1, part a). Similarly, while the glucocorticoid and estrogen response elements have a 3-bp space between the two halves of the response element, the removal of this space to make the two halves of the response element contiguous with one another results in a thyroid hormone response element (TRE) that binds the thyroid hormone receptor. Hence, relatively small changes in sequence can affect the receptor that is bound by the response element and therefore result in it mediating a different pattern of gene expression.

**Table 1. Relationships Among Various Hormone Response Elements**

| <b>(a) Palindromic repeats</b>  |                              |
|---------------------------------|------------------------------|
| Glucocorticoid                  | RGRACANNNTGTYCY <sup>a</sup> |
| Oestrogen                       | RGGTCANNNTGACCY              |
| Thyroid                         | RGGTCA—TGACCY                |
| <b>(b) Direct Repeats</b>       |                              |
| 9- <i>cis</i> -Retinoic acid    | AGGTCAN <sub>1</sub> AGGTCA  |
| All- <i>trans</i> retinoic acid | AGGTCAN <sub>2</sub> AGGTCA  |
|                                 | AGGTCAN <sub>5</sub> AGGTCA  |
| Vitamin D <sub>3</sub>          | AGGTCAN <sub>3</sub> AGGTCA  |
| Thyroid hormone                 | AGGTCAN <sub>4</sub> AGGTCA  |

<sup>a</sup> N indicates that any base can be present at that position, R indicates a purine ie A or G, Y indicates a pyrimidine ie. C or T. A dash indicates that no base is present, the gap having been introduced to align the sequence with the other sequences.

Although a response to thyroid hormone can be obtained by having the palindromic element illustrated in Table 1, part a, most thyroid hormone response elements are of a different form, in which the GGTCA sequence is directly repeated, with a 4-bp gap between the two direct repeats (Table 1, part b). As with the palindromic elements, the spacing between these direct repeats can affect the nature of the receptor that is bound: Direct repeats with a spacing of 1, 2, or 5 bp produce a response to different forms of retinoic acid, a 3-bp spacing produces a response to vitamin D<sub>3</sub>, and a 4-bp space results in a response to thyroid hormone. Detailed structural studies of the different receptors involved have recently revealed the manner in which their different structures allow them to recognize these different response elements specifically (4).

Hence, a family of related response elements that bind a family of related receptor proteins allow different genes to respond to different, but related, lipophilic hormones.

#### Bibliography

1. M. Beato (1989) *Cell* **56**, 335–344.
2. D. J. Mangelsdorf, C. Thummel, M. Beato, P. Herrlich, G. Schutz, K. Umesono, B. Blumberg, P. Castner, M. Mark, P. Shambon, and R. M. Evans (1995) *Cell* **83**, 835–839.



3. M. G. Parker (1993) *Curr. Opin. Cell Biol.* **1**, 512–518.
4. F. Rastinejad, T. Pearlmann, R. M. Evans, and P. B. Sigler, (1995) *Nature* **375**, 203–211.

### Suggestions for Further Reading

5. H. Gronemeyer and V. Laudet, (1995) Nuclear receptors. *Protein Profile* **2**, 1173–1308.
6. H. Gronemeyer, and D. Moras (1995) How to finger DNA. *Nature* **375**, 190–191.
7. D. J. Mangelsdorf, and R. M. Evans (1995) The RXR heterodimers and orphan receptors (1995). *Cell* **83**, 841–850.

## Hormones

Hormones are chemical signaling molecules that regulate metabolism, growth, and differentiation in multicellular organisms ([1](#)). Hormones are produced in specialized cells of the endocrine glands and reach their target cells by way of the blood circulation. The term “hormone” is sometimes applied loosely to encompass various other regulatory compounds, including those of lower animals and plants, that show comparable effects. Several hormones (primarily [steroid hormones](#) and [thyroid hormones](#)) are transported in the circulation while bound to specific binding proteins. These proteins serve as hormone depots, releasing the hormone when needed and also protecting it from rapid inactivation. Most hormones are either steroids (see [Steroid Hormones](#)) or peptides and proteins (see [Peptide Hormones](#)) (Table [1](#)). Adrenaline, the thyroid hormones, and melatonin are hormones derived from [amino acids](#). In insects and crustaceans, peptide hormones, steroid hormones (ecdysteroids), and sesquiterpenoid derivatives (juvenile hormones) have been isolated.

**Table 1. List of Hormones and Sites of Production**

| Hormones                                   | Site of Production            |
|--|-------------------------------|
| <b>Steroids</b>                            |                               |
| Glucocorticoids (cortisol, corticosterone) | Adrenal cortex                |
| Mineralocorticoids (aldosterone)           | Adrenal cortex                |
| Androgens (testosterone)                   | Testis                        |
| Estrogens (17 $\beta$ -estradiol)          | Ovary (follicles)             |
| Progestins (progesterone)                  | Ovary (corpus luteum)         |
| Calcitriol <sup>a</sup>                    | Kidneys                       |
| <b>Peptide and Proteo Hormones</b>         |                               |
| <b>(Partial List)</b>                      |                               |
| Parathormone                               | Parathyroids                  |
| Insulin                                    | Pancreas                      |
| Glucagon                                   | Pancreas                      |
| Ocytocin                                   | Neurohypophysis               |
| Vasopressin                                | Neurohypophysis               |
| Melanotropin                               | Pars intermedia of hypophysis |

|                               |                 |
|-------------------------------|-----------------|
| Somatotropin                  | Adenohypophysis |
| Corticotropin                 | Adenohypophysis |
| Thyrotropin                   | Adenohypophysis |
| FSH (folitropin)              | Adenohypophysis |
| LH (lutropin)                 | Adenohypophysis |
| <b>Amino Acid Derivatives</b> |                 |
| Thyroxin, triiodothyronine    | Thyroid gland   |
| Adrenaline                    | Adrenal cortex  |
| Melatonin                     | Epiphysis       |

---

<sup>a</sup> Steroid derivative, not a steroid, but with molecular action like steroid hormones.

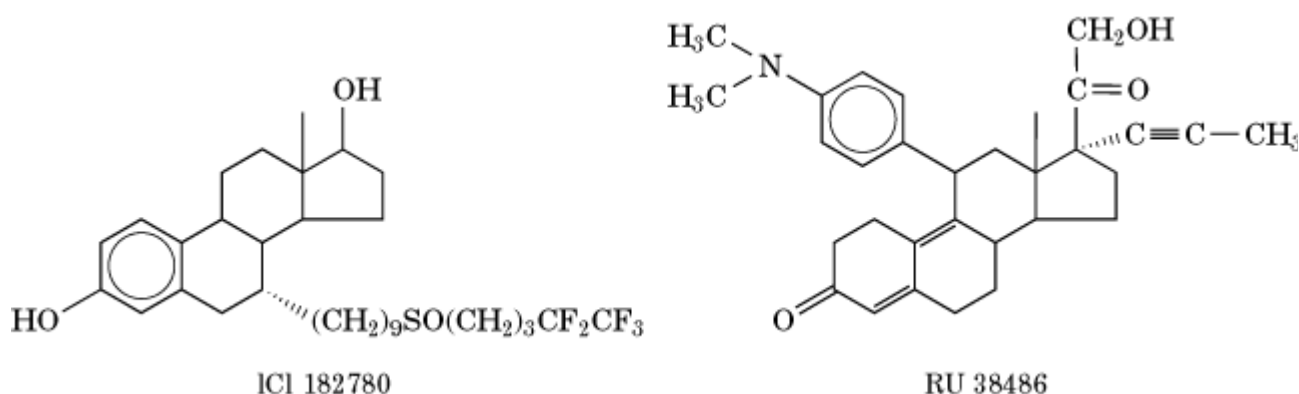
The steroid hormones are derived from cholesterol by a series of enzymatic reactions that take place in the cytosol and in [mitochondria](#) of primarily cells of the adrenal cortex, ovary, and testis. In some cases, the steroid hormone must be subjected to modification in the target tissue, either to be activated or to produce a more active derivative. The steroid hormones are inactivated mainly in the liver in various ways—for example, by reduction of keto groups or unsaturated double bonds, ring hydroxylation, oxidation of hydroxyl groups, and cleavage of the side chain. They are excreted after being conjugated to sulfuric or glucuronic acids. Most of the peptide and proteohormones are synthesized in the form of precursor proteins (prohormones) and are stored in the endocrine cell. Before being released into the circulation, the prohormones are cleaved to the active hormone. Their inactivation involves **proteolytic** cleavage by intracellular [proteinases](#). The production and release of many hormones is controlled by “tropic” hormones produced by the hypophysis, an endocrine “master” gland; this phenomenon is known as *hormone hierarchy*. The hypophysis responds to chemical signals (metabolites, hormones) from the periphery or to neuroendocrine stimuli from the hypothalamus (liberins, which act positively on hypophysial hormone release, and statins, which inhibit their release). From an evolutionary standpoint, the endocrine system developed earlier than the other two major regulatory systems of higher organisms, the nervous and the immune systems.

The molecular mechanism of hormone action depends on the physical chemical characteristics of the hormones. Steroid and thyroid hormones, which are **hydrophobic**, can enter the cell and bind intracellularly to their respective [hormone receptors](#) (2-4). These proteins belong to the superfamily of nuclear receptors and are present in the quiescent cell as complexes with other proteins, among them **heat-shock** proteins. This complex dissociates upon binding of the hormone ligand, the receptor homodimerizes, and it then binds to specific nucleotide sequences, the [hormone response elements](#) (HRE). These represent receptor-activated [enhancer](#) elements situated at varying distances from the **promoter** of the hormonally regulated genes. Interaction of the hormone–receptor homodimer complex with the HRE results, in a still undefined way, in either stimulation or inhibition of gene [transcription](#), and consequently of expression of the gene. The nature of the protein synthesized determines the **phenotypic** effect. The [hydrophilic](#) hormones (peptide and proteohormones, adrenaline and melatonin) cannot enter the cell, but bind to membrane [hormone receptors](#). This results, by way of a **G-protein** effector system, in activation of enzyme systems (eg, [adenylate cyclase](#), **phospholipase C**) generating [second messengers](#), which are **allosteric** effectors of protein kinases primarily (5-7) (see [Phosphorylation, Protein](#)). The hormonal effect depends upon the nature of the phosphorylated protein. One class of membrane receptors eg, the [insulin](#) receptor, first dimerizes upon binding the hormone and is then autophosphorylated; in this state, it can phosphorylate and activate other intracellular proteins (see [Hormone Receptors](#)).

Analogue of hormones can be either more or less potent than the natural hormones, or they can be inhibitors of hormone action; many have been synthesized and used in research or as therapeutics.

The more potent analogues of steroid hormones exhibit higher affinity for the cognate receptor than the natural hormone or greater stability, due to lower rates of inactivation and excretion. Important examples of such analogs are diethylstilbestrol, an [estrogen](#), and [dexamethasone](#) and triamcinolone, [glucocorticoids](#). Competitive inhibitors also bind with high affinity to the receptor but cannot elicit the hormonal effect; they can inhibit receptor homodimerization, binding of the receptor to hormone response elements, or interaction of the receptor with the transcriptional machinery. Some such inhibitors (Fig. 1) that have found clinical application include the antiestrogen ICI 182780 used in treatment of breast cancer, the antiandrogen cyproteron in prostate cancer and hirsutismus, and RU486, an antiprogestosterone compound used in pregnancy termination.

**Figure 1.** Synthetic hormone antagonists.



### Bibliography

1. E. T. Baulieu and P. A. Kelly (eds.) (1990) *Hormones: From Molecules to Disease*, Hermann, Paris.
2. E. M. Evans (1988) *Science* **240**, 889–895.
3. D. J. Mangelsdorf et al. (1995) *Cell* **83**, 835–839.
4. M. Beato, P. Herrlich, and G. Schutz (1995) *Cell* **83**, 851–857.
5. A. Hall (1990) *Science* **249**, 635–640.
6. W. J. Tang and A. G. Gelman (1992) *Cell* **70**, 869–872.
7. N. Divecha and R. F. Irvine (1995) *Cell* **80**, 269–278.

### Suggestions for Further Reading

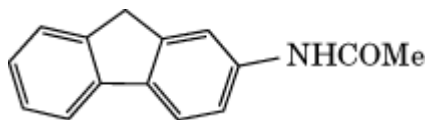
8. J. D. Wilson and D. W. Foster, eds. (1992) *Williams Textbook of Endocrinology*, 8th ed., Saunders, London.
9. M. G. Parker (1993) *Steroid Hormone Action*, IRL Press, Oxford, U.K.
10. R. L. J. Ribeiro, P. J. Kushner, and J. D. Baxter (1995) The nuclear hormone receptor superfamily. *Annu. Rev. Med.* **46**, 443–453.

### Hotspot

DNA sequence analysis can be used to reveal the distribution of [mutations](#) within a gene due to spontaneous [mutagenesis](#) or following treatment by specific chemicals. Such mutation spectra reveal that neither spontaneous nor induced mutagenesis is random, but mutations cluster in specific regions within a gene. Those DNA sites at which a large number of mutants occur are referred to as “hot spots.” For example, Halliday and Glickman (1) found that 71% of the 729 spontaneous mutants in the *lac I* gene of *E. coli* carried mutations in a small region (positions 620–632) with base sequence [5'TGGC3']<sub>3</sub>. In contrast, 98% of ICR-191-induced mutations in the *lac I* gene were +1 and –1 [frameshift mutations](#) in runs of repetitive G.C base pairs (2).

There are several possible explanations for genetic hot spots. A trivial one is that there could be a selection bias that favors the detection of mutations at some sites in preference to those at others, but this is probably not true of most observations in the literature. The most likely explanations of most hot spots are unequal distribution of premutational lesions in the gene, differences in processing of lesions in different regions, or some sort of differential [DNA repair](#) process. The frameshift mutations produced by *N*<sub>2</sub>-(acetylamino)fluorene (AAF, Fig. 1) provide an example in which the reasons for hot spots have been explored in some detail, especially by Fuchs and co-workers (3). This model carcinogen forms adducts on the C8 position of guanine bases. However, the distribution of AAF-induced mutations fails to correspond to the distribution of AAF adducts in DNA. Only 19% of the DNA adducts appear at the sites of mutational hot spots, even though these sites produce 89% of the mutations (4). Subsequent studies have suggested that these hot spots result from differences in processing of DNA damage at particular sites, with premutational events being more readily converted into a mutation at hot spots than at other sites in DNA.

**Figure 1.** Structure of AAF (2-acetylaminofluorene).



## Bibliography

1. J. A. Halliday and B. W. Glickman (1991) *Mutat. Res.* **250**, 55–71.
2. M. P. Calos and J. H. Miller (1984) *J. Mol. Biol.* **153**, 39–66.
3. G. R. Hoffmann and R. P. P. Fuchs (1997) *Chem. Res. Toxicol.* **10**, 347–359.
4. R. P. P. Fuchs (1984) *J. Mol. Biol.* **177**, 173–180.

## HU And Related Proteins

The bacterial type-II [DNA-binding proteins](#) include the abundant HU and integration host factor (IHF) proteins and the phage-encoded TF1. They are believed to condense their cognate **genomes** by binding at multiple sites and inducing coherent bends. HU and TF1 possess little, if any, binding sequence specificity, but IHF, which is also required for site-specific **recombination**, **DNA replication**, and [transcription](#), binds at specific sites characterized by a limited **consensus sequence**.

All proteins of this type bind to DNA as dimers and bend the DNA to a substantial extent. In particular, IHF induces a bend angle of at least  $160^\circ$ , and possibly in excess of  $180^\circ$ , within the  $1/2$  turns of double helix that comprise its binding site. IHF is normally a heterodimer, while HU can exist either as a homo- or as a heterodimer. In both IHF and HU, the two  $\sim 10$  kDa subunits intertwine to form a compact core, from which two long [beta-sheet](#) ribbon arms extend. These arms track along the minor groove from the inside to the outside of the wrapped DNA, where they terminate at the two substantial kinks. In addition to these interactions via the b arms, IHF also clamps the hairpin by minor-groove contacts with **alpha-helices** from both subunits in the core of the dimer. All the contacts to the DNA are either in the minor groove or are part of an extensive network of [electrostatic interactions](#) with the backbone phosphate groups.

The two kinks are induced by the partial intercalation between adjacent base pairs of an absolutely conserved [proline](#) sidechain located at the tips of the b arms. In the IHF–DNA complex, the DNA bend is maintained by two distinct mechanisms. On the outside of the bend, the hydrophobic intercalation stabilizes the opening of the minor groove. On the inside, charge neutralization counteracts the enhanced repulsion between the phosphates on opposite sides of the narrowed grooves.

The sequence dependency of IHF binding is determined by the conformation of DNA rather than by base-specific contacts. The “consensus” sequence consists of two short elements separated by approximately half a turn in only one of the two half-sites. The conserved sequence CAA at the kink site can accommodate the severe distortion induced by the protein better than other short sequences, and so is favored.

#### Suggestion for Further Reading

P. A. Rice, S. Yang, K. Mizuuchi, and H. A. Nash (1996) Crystal structure of an IHF–DNA complex: a protein-induced U-turn. *Cell* **87**, 1295–1306.

## Hybrid Cell

A [heterokaryon](#), with two separate **nuclei**, eventually proceeds to **mitosis**. This produces a *hybrid cell* in which the two nuclear envelopes have been disassembled, allowing all the [chromosomes](#) to be brought together in a single large nucleus. Although such hybrid cells can be cloned to produce hybrid [cell lines](#), the initial hybrid cells are unstable and lose chromosomes. For unknown reasons, mouse-human hybrid cells predominantly lose human chromosomes (see [Hybridomas](#)). The loss is random and gives rise to a variety of hybrid cell lines, each of which contains only one or a few human chromosomes. This phenomenon has been used in mapping the locations of genes in the human [genome](#).

## Hybrid Dysgenesis

Hybrid dysgenesis is a syndrome observed in *Drosophila* that results from the high-frequency

germline [transposition](#) of a particular *Drosophila* [transposable element](#) called the **P element** (1). The phenotypes associated with hybrid dysgenesis include a high frequency of [mutation](#) and sterility. Hybrid dysgenesis results when males from a strain that contains P elements in its genome are crossed with females from a non-P strain. The lack of high-frequency transposition when P<sup>+</sup> (or P<sup>-</sup>) strains are crossed among themselves results from the fact that the P element encodes, in addition to a [transposase](#), an inhibitor of transposition. Thus, hybrid dysgenesis occurs only when the egg is from a P<sup>-</sup> strain, lacking the P element and its inhibitor, and when the P-element-containing DNA arrives via sperm without the inhibitor.

Regulation at several levels can account for hybrid dysgenesis. Transposition of the P element is restricted to the germ line because an intact transposase protein is produced only in the germ line. Transposase is encoded by four separate exons whose [messenger RNA](#) needs to be assembled by [RNA splicing](#). This splicing event can occur only in the germ line because a somatic protein that binds to Exon 2 RNA near the 5' splice site inhibits splicing in somatic tissue (2, 3). Another level of control of P element transposition is the transposition inhibitor encoded by the P element itself (1). This inhibitor derives from the P element mRNA, so the transposition inhibitor is related to, but is not identical to, transposase. The mechanism of inhibition appears to involve repression of [transcription](#), which decreases the amount of available transposase.

The high-frequency P-element transposition in the germline that leads to hybrid dysgenesis thus results from the germline-specific synthesis of transposase from the DNA of the sperm of a P<sup>+</sup> male in the absence of a transposition inhibitor; the inhibitor is lacking in the maternal cytoplasm because this fly lacks a P element to generate the inhibitor.

## Bibliography

1. W. R. Engels (1996) *Curr. Top. Microbiol. Immunol.* **204**, 103–124.
2. F. A. Laski, D. C. Rio, and G. M. Rubin (1986) *Cell* **44**, 7–19.
3. C. W. Siebel and D. C. Rio (1990) *Science* **248**, 1200–1208.

## Hybrid Individual

A hybrid is the offspring of genetically different individuals. The term is used often to refer to the offspring of distantly related members of the same species or of closely related species. The potential production of fertile hybrids is a fundamental criterion to assign different individuals, strains, or races to a single biological species.

A *heterozygote* is an individual whose nuclei carry at least two different versions of at least part of the genetic information. The difference may consist in the nucleic acid sequence, the number of chromosomes, or the order of the genes in the chromosomes.

### 1. Mendelian Hybrid

In a restricted sense, a Mendelian hybrid is the result of crossing two true-breeding strains (homozygotes) (see [Homozygote](#), [True Breeding](#)).

A hybrid may involve two **alleles** *A* and *a* of a single gene (a monohybrid *Aa*), or two pairs of alleles at two different genes (the dihybrid *AaBb*), and so on. Hybrids produce characteristic Mendelian

segregations in crosses with each other or with their parents (see [Mendelian Inheritance, Crosses](#)).

Genes that appear in heterozygosis in a polyhybrid may be carried in the same or different [chromosomes](#). The dihybrid  $AaBb$  is said to be in the *cis-configuration*, or the *coupling* phase, when it contains two dominant alleles on the same chromosome and two recessive alleles on the homologous chromosome of the diploid; to make a graphic distinction, the genotypes of the homologous chromosomes are separated by a slash:  $AB/ab$ . The dihybrid  $AaBb$  is said to be in the *trans-configuration*, or the *repulsion* phase, when each chromosome carries a dominant and a recessive allele:  $Ab/aB$ . When two genes are sufficiently close in the same chromosome, their alleles do not form the random combinations expected from Mendelian inheritance; alleles in the same chromosome tend to remain together.

Hybrids may be larger, more vigorous, or more resistant to adverse conditions than their parents. This phenomenon, called [heterosis](#) or *hybrid vigor*, occurs frequently in plants and animals and has led to the commercial production of hybrid seed for agriculture. Depending on the actual techniques that are used, hybrid seed may be too expensive and replaced in practical application by “*double hybrid seed*,” the offspring of two hybrids, each with different parents.

## Hybridization, Nucleic Acids

Shortly after the discovery of the double helix, it was found that the two strands of DNA could be separated and that duplex DNA could be reformed. Further it was found that DNA from different sources could be combined to form hybrid duplexes (1). Hybridization techniques have become one of the workhorses of molecular biology (2). The properties that promote formation of duplexes from partially or exactly complementary single-stranded nucleic acids are the foundation of hybridization techniques. A requirement for most hybridization experiments is labeled nucleic acid, which provides a detectable signal reporting the quantity of hybrid duplex formed with high sensitivity. Many techniques in nucleic acid chemistry rely on formation of a hybrid duplex composed of single-stranded nucleic acids from different origins.

Hybridization techniques vary markedly in their details, but there are features common to all. Under appropriate solution conditions, two nucleic acid strands with at least partial sequence complementarity can form a hybrid duplex. Each strand may be either DNA or RNA. Most hybridization experiments consist of a minimum of three steps: (1) preparation of the target, (2) formation of the hybrid duplex, and (3) detection of the hybrid duplex. Each step can be technically challenging, and proper procedure requires multiple accompanying control experiments.

The nucleic acids to be hybridized should be free of interfering contaminants. If any of them are part of a duplex, a denaturation step, either by heat or alkali, is included. Frequently, a partial degradation step prior to denaturation is included in the procedure. The degradation is accomplished either by shear or by endonuclease digestion. Selection of the temperature and solution conditions for the hybridization reaction is critical in determining the outcome (see [Stringency](#)). Under high stringency conditions, only regions that are highly complementary will form hybrid duplexes. As the stringency is lowered, the fraction of bases participating in hybrid duplexes will increase because of the greater tolerance for defect in sequence alignment under such conditions.

The power of hybridization techniques comes from the dependence of the stability of the hybrid duplex on the extent of complementarity of the constituent strands. The stability of the hybrid duplex depends strongly on solution conditions and temperature, which must be designed carefully to achieve the desired level of sequence discrimination. Most hybridization solutions contain EDTA to

bind divalent cations, a polymer utilized to enhance hybrid formation (presumably via an [excluded volume](#) effect), and a **denaturant**, usually formamide, in a pH buffered salt solution. The combination of solution conditions (especially the formamide and NaCl concentrations) and temperature defines the stringency of the hybridization reaction. Optimizing the stringency of the reaction is of paramount importance in the success of any hybridization experiment. Defining the optimum depends on the objectives of the experiment. High stringency conditions reduce the stability of the hybrid duplex and select for high sequence homology. Such conditions are used when high target specificity is desired. Low stringency conditions tolerant greater sequence variability. Lower stringency conditions are useful when it is desirable to relax sequence selectivity, such as when probing for a family of related sequences.

When hybridization is performed in solution, the selective detection of hybrid duplexes is difficult. Some progress has been made on this issue by application of fluorescence **energy transfer** techniques; however, selective detection of hybrids remains problematic. This has limited the application of solution hybridization. Its primary application is for formation of [heteroduplexes](#) that are subsequently analyzed to determine the extent of sequence complementarity. Another difficulty encountered in solution hybridization is the competition between formation of the hybrid duplex and reformation of the original duplex. There are several approaches to address this problem. For example, when RNA–DNA hybrids are formed, high formamide concentrations (80%) can be used; DNA–DNA duplexes are more destabilized under such conditions. The most common way to circumvent the competitive equilibrium problem is by immobilization of the denatured nucleic acid on a surface. Immobilization prevents reassociation of the denatured duplex. Immobilization of a target nucleic acid forms the basis for a wide variety of probe hybridization techniques. It is in these techniques that the true utility of hybridization for studying nucleic acids is realized.

Traditionally, hybridization reactions have been confined to formation of RNA–DNA and DNA–DNA duplexes. The recent and rapid developments in synthesis of unnatural nucleic acid analogues with base and backbone modifications is expanding the number of possible combinations of duplexes that can be formed. These developments are likely to be exploited in new hybridization reactions. Also, the recent development of triplex-forming nucleic acids for use as probes has expanded the range of reactions that can be usefully included under the heading hybridization.

#### Bibliography

1. B. J. McCarthy and R. B. Church (1970) *Annu. Rev. Biochem.* **39**, 131–150.
2. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning: A Laboratory Manual*, 2nd ed., Cold Spring Harbor Press, New York.

#### Suggestion for Further Reading

3. P. Tijssen (1993) *Hybridization with Nucleic Acid Probes: Part 1, Theory and Nucleic Acid Preparation; Part 2, Probe Labeling and Hybridization Techniques*, Elsevier, Amsterdam.

## Hybridomas

Classical seminal studies by Köhler and Milstein in 1975 (1) led to the generation of [monoclonal antibodies](#) with defined specificity. Hybridoma methodology was based primarily on fusing murine immune lymphocytes with a malignant myeloma lymphocyte cell line that possessed certain important genetic and biochemical properties. This important methodology was necessitated by the fact that immune antibody-producing B lymphocytes were difficult to maintain as stable [cell lines](#)



under [tissue culture](#) conditions. However, fusion of antibody-producing cells with a malignant myeloma cell line resulted in proliferative cells that can be maintained indefinitely.

Spontaneous fusions between lymphocytes are rare events, so the frequency of fusion was increased to generate large numbers of fused cells by using surface-active agents such as inactivated [Sendai virus](#) or the chemical polyethylene glycol (PEG). Upon mixing immune lymphocytes with myeloma cells and the fusing reagent, a selection system was mandatory to eliminate unfused parental cells and to assay for appropriately fused cells that synthesized antibodies with the desired specificity. The tissue culture selection process was dependent on the procurement of myeloma cell lines that no longer synthesized an immunoglobulin and lacked genes for **thymidine kinase** (TK) and/or **hypoxanthine-guanine-phosphoribosyl transferase** (HGPRT). These enzymes catalyze important steps in the salvage pathway used in DNA synthesis by dividing cells. Unfused myeloma cells are unable to proliferate in the HAT selection medium (containing hypoxanthine, [aminopterin](#), and thymidine) because aminopterin inhibited the *de novo* DNA biosynthetic pathway, forcing the cells to use the nonfunctional salvage pathway. Fused myeloma cells proliferate in HAT medium because the immune lymphocytes (fusion partner) furnish the essential TK and HGPRT enzymes for the salvage pathway. Thus, cells that survived incubation with HAT medium represented a successful immune B-cell–myeloma cell fusion, usually on a 1:1 cellular basis. The second phase in hybridoma production was the selection of those fused cells that synthesized and secreted antibody of the desired specificity. This was generally achieved by diluting fused cell populations to a single cell and then generating clones by proliferation of the single cell. The extracellular fluid or media from the single cell and resulting progeny were assayed for the appropriate antigen binding. The methods and techniques used to detect monoclonal antibody activity are usually based upon some variation of the solid-phase immunoassay (eg, ELISA or [radioimmunoassay](#) formats).

Hybridoma methodology and consistent generation of relatively large amounts of monoclonal antibodies with defined specificities has proven to be a source of important immunological reagents.

1. The provision of adequate quantities of monoclonal antibody has facilitated the routine determination of primary structures of both the H and L chains at the protein level.
2. Various monoclonal antibodies can be crystallized on a regular basis, providing atomic, 2- to 3-Å resolution structures by **X-ray crystallographic** procedures. In general, **Fab** fragments (four-domain substructures of antibodies, including the variable domains of both chains) with bound ligand crystallize most efficiently. Approximately 100 resolved structures of monoclonal antibodies (Fab fragments) have now been reported.
3. This methodology provides large homogeneous hybridoma cell populations from which [messenger RNA](#) responsible for the synthesis of H and L chains can be obtained. Through development of polynucleotide primers, polymerase chain reaction (**PCR**) technology can be used to generate [complementary DNA](#) copies for sequencing at the gene level.
4. Gene cloning experimentation with antibody genes has been important in successful [site-directed mutagenesis](#) studies. Such studies have been critical in deciphering structure–function relationships governing antibody activity.
5. Hybridoma methodology has led to the development of antibody derivatives, such as single-chain antibodies (2). The latter represent ~25,000–Da structures containing only the variable domains of the H and L chains of a specific antibody molecule. Variable domains constituting single-chain antibodies are usually attached covalently through the use of flexible polylinkers (10 to 15 amino acid residues) that are encoded in the single-chain synthetic gene. Single-chain genes are subsequently incorporated into the appropriate plasmids and expressed in either **prokaryote** or **eukaryote** cell lines.
6. Antigen–antibody interactions can now be studied on a homogeneous basis, yielding important information regarding the thermodynamics of binding.
7. The provision of monoclonal antibodies by hybridoma technology has led to the development of standard diagnostic procedures as well as therapeutically useful immunochemical reagents.

Monoclonal antibodies of defined specificity generated through hybridoma methodology have been used to solve many issues related to structure–function relationships within the antibody molecule. Based on studies with monoclonal antibodies, the property of Fab (active site) segmental flexibility within the IgG class of antibodies was solved. Homogeneous binding at both active sites within a bivalent molecule ruled out **allosteric** effects within antibody molecules. Thus, binding at one active site does not influence antigen binding at the adjacent site. Conformational changes transmitted throughout the molecule subsequent to binding of antigen were dismissed as an explanation for such important phenomena as **complement** binding and fixation. In all cases, the availability of monoclonal antibodies proved important to examine definitively these important questions.

Although hybridoma technology has been the method of choice to produce monoclonal antibodies of defined specificity, new techniques have now emerged. Similar to the construction of single-chain antibodies, gene segments encoding the H- and L-chain variable domains are genetically fused to genes encoding a **bacteriophage** coat protein (3). The engineered bacteriophage infect bacteria, and the resulting phage particles express active antibody products on their surface. The resulting [phage display library](#) expresses many different antigen-binding domains. Phage that specifically bind antigen are selected and used to infect bacteria in a second cycle. Each selected phage produces a monoclonal antigen-reactive particle. Primary structures of the variable domains can be determined, and those genes can be fused to the antibody constant region genes to reconstruct a monoclonal antibody. Such genes transfected into myeloma cell lines are expressed, and the antibody products are secreted. Phage display methodology has important implications for the future.

#### Bibliography

1. G. Köhler and C. Milstein (1975). *Nature* **256**, 495–497.
2. R. E. Bird, K. D. Hardman, J. W. Jacobson, S. Johnson, B. M. Kaufman, S.-M. Lee, T. Lee, S. H. Pope, G. S. Riordan, and M. Whitlow (1988) *Science* **242**, 423–426.
3. I. Roitt, J. Brostoff, and D. Male (1996) *Immunology*, 4th ed. Mosby, London, p. 28.9.

#### Hydration

The total hydration of a protein (or nucleic acid) is the effective amount of [water](#) immobilized by (bound to) a protein molecule or other macromolecule. This consists of a summation of interactions of water molecules with individual sites on the [accessible surface](#) of a protein molecule. Such interactions vary from (1) very strong ones, such as water molecules trapped within cavities on a protein molecule and involved in the actual folded structure of the protein, (2) weak interactions, such as water molecules that hydrate charged and other **polar groups** on the surface of a protein molecule; and (3) very weak interactions that comprise water molecules whose rotation or translation is momentarily perturbed by the proximity of a protein molecule (by weak attraction or repulsion). This last type of interaction is primarily **entropic** in nature and reflected by a very small value of the **free energy**. The summation of all these interactions manifests itself as effective [binding](#) (ie, immobilization of the water by the protein) that may have a small fractional value at each protein surface site, but which sums up to some whole numbers over the entire molecule.

*Total hydration* is very difficult to measure, and its values may be a function of the techniques used. [X-ray crystallography](#) detects the more strongly interacting water molecules. The [NMR](#) technique that detects water molecules whose freezing is perturbed by the presence of the protein (1) and vapor

pressure osmometry (2) give similar values of total hydration. These values are of similar magnitude to those obtained from [equilibrium dialysis](#) in those cases in which there is total exclusion of ligand from the protein surface. Furthermore, a measurement of site occupancy by the ligand by a contact detecting technique (total binding), when combined with equilibrium dialysis, yields the total hydration (see Eq. 1 of [Binding](#)). See also [Preferential Hydration](#).

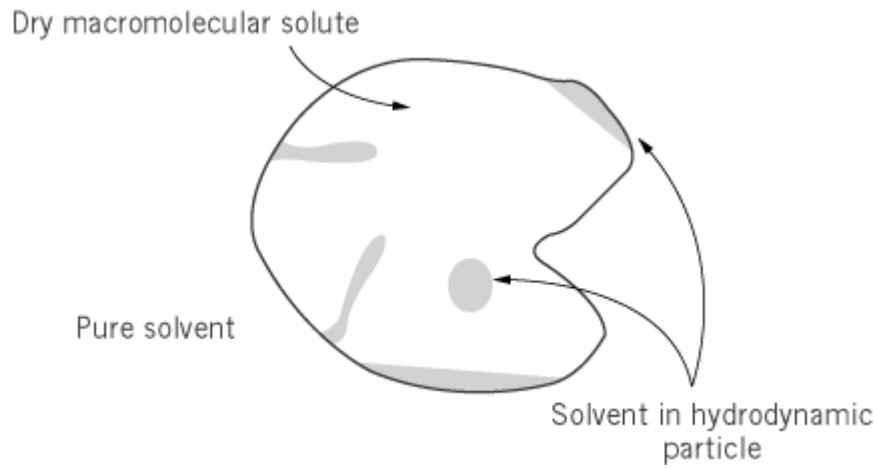
#### Bibliography

1. I. D. Kuntz and W. Kauzmann (1974) *Adv. Protein Chem.* **28**, 339–345.
2. H. B. Bull and K. Breese (1974) *Arch. Biochem. Biophys.* **161**, 665–670.

### Hydrodynamic Volume

Measurement of the hydrodynamic properties of a biological [macromolecule](#) can provide low resolution information about its size and shape. This information can be useful in modeling any interactions in which the macromolecule is involved. Among the techniques used for hydrodynamic measurements are **light scattering**, **sedimentation velocity** and [sedimentation equilibrium centrifugation](#), and [size exclusion chromatography](#). The hydrodynamic properties of a macromolecule are often expressed in terms of a macromolecular volume (see [Stokes Radius](#)). In considering the volume of a macromolecule, however, one must include not only the macromolecule itself but also the role of solvent. Solvent that is “trapped” is an inherent part of a macromolecule and will, therefore, be reflected in measurements of its transport properties. A simple model of [hydration](#) of a biological macromolecule is shown in Figure 1. As indicated, [water](#) may be bound either at the surface of the macromolecule or in interior spaces. Some solvent molecules may be bound specifically, and others may be simply trapped. In the thermodynamic treatment of this bound water, no distinction is made between the various classes of bound-water molecules. Instead, they are treated as a smooth surface or as internal pools. In any hydrodynamic measurement, these water molecules remain tightly associated with the macromolecule, so the total volume of the macromolecule derived from such studies will necessarily include the bound solvent. A mathematical treatment of hydrodynamic volume in which the bound water molecules are explicitly included is provided here, based on References [1-3](#).

**Figure 1.** A cartoon diagram of a macromolecule with bound water. The water can be associated with the surface of the macromolecule or may be accumulated in “pools” in the interior of the molecule. The thermodynamic treatment of this water presented here does not distinguish between these different types of bound-water molecules.



Assume that a given macromolecule has  $d_1$  grams of water bound per gram of macromolecule. If  $V_1^*$  is the average specific volume of this water in  $\text{cm}^3/\text{gm}$ , and if a solution consists of  $g_1$  grams of water and  $g_2$  grams of anhydrous macromolecule, the total volume of the solution can be calculated using the following equation:

$$V_{\text{tot}} = (g_1 - g_2 d_1) \bar{V}_1^\Delta + g_2 V_2 + \delta_1 V_1^* g_2 \quad (1)$$

where  $V_{D1}$  is the specific volume of pure water in  $\text{cm}^3/\text{gm}$ . The first term in the equation is the volume of unbound water; it is assumed to have the same specific volume as pure water. The second term is the volume occupied by the macromolecule; the specific volume of the macromolecule in solution,  $V_2$ , is not necessarily equivalent to the specific volume of the pure solid macromolecule. The third term is the volume of bound water. The sum of the second and third terms is the volume of the hydrated macromolecule. This volume can be calculated using the following equation:

$$V_h = \frac{M}{N_0} (V_2 + \delta_1 V_1^*) \quad (2)$$

where  $M$  is the **molecular weight** of the macromolecular solute and  $N_0$  is Avogadro's number. Unfortunately, this equation is not very useful, since the terms  $V_2$ ,  $d_1$ , and  $V_1^*$  are not known. These variables can, however, as indicated below, be replaced by those that are experimentally accessible.

The **partial specific volume** of a solute is the change in solution volume when a small increment of the solute is added, at the limit of infinite dilution. This quantity is determined experimentally by measuring the density of a solution as a function of the weight concentration of the solute. The expected result of this measurement can be derived from Equation (1)

$$\bar{V}_2 = \left[ \frac{\partial V_{\text{tot}}}{\partial g_2} \right]_{g_1, T} = -\delta_1 \bar{V}_1^\Delta + V_2 + \delta_1 V_1^* \quad (3)$$

This derivation assumes that the macromolecule is sufficiently dilute that the degree of its hydration is independent of its concentration. Using this equation, we can replace  $V_2$  and  $V_1^*$  in Equation (2) to yield

$$V_h = \frac{M}{N_0}(\bar{V}_2 + \delta_1 \bar{V}_1) \quad (4)$$

The term  $V_{D1}$  represents the specific volume of water, which is a pure substance, and can be replaced by its partial specific volume,  $V_1$ , which is simply the inverse of the density of pure water. It follows that the hydrated volume of the macromolecule is

$$V_h = \frac{M}{N_0}(\bar{V}_2 + \delta_1 \bar{V}_1) \quad (5)$$

The partial specific volume of the macromolecular solute,  $V_2$ , can be measured directly by determining the change in solution volume as a function of added macromolecular solute. The final term is the hydration, and  $d_1$  is an unknown, but it can be reasonably estimated.

Direct determination of the partial specific volumes ( $V_2$  in Eq. (5)) of macromolecules is accomplished by determining the change in solution volume as a function of added macromolecular solute. This can be a difficult measurement and may require a large amount of material. The availability of accurate microbalances can minimize the latter difficulty. Partial specific volumes can also be determined by ultracentrifugation in various mixtures of  $H_2O/D_2O$  with varying densities (3). The advantages of this technique are that it is both rapid and requires very little material. The partial specific volumes of proteins can also be reasonably estimated from their amino acid compositions, using estimates of specific volumes of individual amino acid residues (4).

Hydration of biological macromolecules ( $d_1$  in Eq. (5)) has been measured using a variety of techniques, including NMR spectroscopy of frozen samples and calorimetric studies. Results from these studies on proteins indicate that the variation in bound water is related to the amino acid composition of the protein. Based on these results, methods for calculating the amount of water bound to a protein from its amino acid composition have been developed. Results of the calculations are in good agreement with the results of experimental measurement (5).

#### Bibliography

1. C. Tanford (1961) *Physical Chemistry of Macromolecules*, Wiley, New York, pp. 339–341.
2. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part II, Freeman, San Francisco, pp. 550–555.
3. S. J. Edelstein and H. Schachman (1973) *Meth. Enzymol.* **27**, 82–98.
4. E. J. Cohn and J. T. Edsall (1943) *Proteins, Amino Acids and Peptides as Ions and Dipolar Ions*, Reinhold, New York, p. 157.
5. I. D. Kuntz Jr. and W. Kauzmann (1974) *Adv. Protein Chem.*, **28**, 239–345.

#### Hydrogen Bond

A hydrogen bond is an attractive interaction between the donor pair  $X-H$  and the acceptor atom  $Y$ , where  $H$  is a hydrogen atom that is covalently bound to an electronegative atom  $X$  but attracted to another electronegative atom,  $Y$ . The hydrogen bond is commonly represented by a dotted line, as in  $X-H \cdots Y$ . The electron density in the  $X-H$  bond is attracted more to  $X$ , so  $X$  has a partial negative charge and  $H$  a partial positive charge, which has been demonstrated by neutron and X-ray

diffraction studies of hydrogen bonded crystals. Since the hydrogen bond involves at least three and sometimes four atoms, physically it is characterized by the bond length ( $\text{—H}^{1/4}\text{Y—}$ ), the angle formed by  $\text{X—H}^{1/4}\text{Y}$ , and the enthalpy of formation,  $DH$ . Typical donors are  $\text{F—H}$ ,  $>\text{NH}$ ,  $\text{—NH}_2$ ,  $\text{—OH}$ ,  $\text{—COOH}$ , and less frequently,  $\text{—S—H}$ ,  $\rightarrow\text{C—H}$  groups. Hydrogen bonds involving  $\text{C—H}$  are among the weakest. The acceptors are commonly  $\text{—N}$ ,  $>\text{C=O}$ ,  $\text{—C—F}$ ,  $>\text{NH}$ ,  $>\text{O}$ , and p-electron systems of unsaturated or aromatic hydrocarbons. Examples of hydrogen bonded pairs commonly found in biological structures are  $[\text{P—OH}^{1/4}\text{O=P}]$ ,  $[\text{O—H}^{1/4}\text{O=C}]$ ,  $[\text{N—H}^{1/4}\text{O=C}]$ , and  $[\text{O—H}^{1/4}\text{O—H}]$  (1).

Historically, the presence of associative interactions among certain types of organic and inorganic molecules was recognized in the early 1900s. The remarkable unusual properties of liquid [water](#) were attributed to the network of hydrogen-bonded  $\text{H}_2\text{O}$  molecules (for a summary of the early history and references, see Ref. 2). The idea of a hydrogen bond provided an extra valence for the hydrogen atom, which already has a closed electron shell structure. The early observations were summarized in a unified view and presented as the “hydrogen bond” by L. C. Pauling (3).

## 1. Methods of Detection

For their ease and sensitivity, the [vibrational spectroscopy](#) techniques of infrared (IR) and, less frequently, Raman spectroscopy, plus [NMR](#), are the preferred methods for detecting the presence of hydrogen bonds. Neutron diffraction from single crystals gives the most precise elucidation of the hydrogen bond geometry, as it detects the hydrogen atom directly, in contrast to X-ray crystallography, where the presence of a hydrogen bond can be inferred only by the shorter than normal distance between the heavier atoms X and Y (1, 2, 4, 5).

### 1.1. X-H Stretching Mode in Infrared Spectroscopy

A low frequency shift in IR spectra in the frequency range of  $1500\text{--}4000\text{ cm}^{-1}$  that is highly temperature-dependent and accompanied by a dramatic increase in absorbance is usually taken as evidence of hydrogen bonds involving  $\text{O—H}$  or  $\text{N—H}$  groups as donors. Bending of  $\text{X—H}^{1/4}\text{Y}$  bonds in and out of the plane causes absorbance in the region of  $1700\text{--}1800\text{ cm}^{-1}$  and  $400\text{--}900\text{ cm}^{-1}$ , respectively. The stretching and bending vibrations of  $\text{H}^{1/4}\text{Y}$  hydrogen bonds have much lower resonance frequencies of  $50\text{--}600\text{ cm}^{-1}$  and less than  $50\text{ cm}^{-1}$ , respectively (1). Such spectral changes allow the determination of an equilibrium constant for formation of the hydrogen bond, and its temperature dependence can give the enthalpy of formation. The low frequency vibrational modes make important contributions to the thermodynamic properties of hydrogen-bonded systems.

### 1.2. Deshielded Proton Signals in NMR

NMR spectroscopy can detect hydrogen bonds because the hydrogen atoms involved are less shielded and consequently resonate at a lower magnetic field. In the case of nucleic acids, these hydrogen atoms are known to have typical chemical shifts of 9–15 ppm downfield of the internal standard (6). If the acceptor group is of aromatic p electrons, the resonance is shifted upfield relative to non-hydrogen-bonded atoms.

### 1.3. Position of the H Atom by Neutron and X-Ray Diffraction

Neutron diffraction gives the most accurate positions of hydrogen atoms, as they scatter about half as much as C, N, and O atoms. In contrast, hydrogen atoms are barely detectable by X-ray diffraction, and hydrogen bonds are evident primarily from the closeness of the X and Y atoms. On the other hand, there are much more X-ray crystallography data available than neutron diffraction data. The controversy over the existence of weak  $\text{C—H}^{1/4}\text{O}$  (or N, S, Cl) hydrogen bonds has been discussed (7).

## 2. Energetics

## 2.1. Enthalpy of Formation

The enthalpy of formation of typical hydrogen bonds is commonly determined as the temperature dependence of the association constant for hydrogen-bonded dimers in isolation. The values measured range from a few to 100 kJ/mol, and the majority are in the range 20–30 kJ/mol (Table 1). The bond energy of a hydrogen bond is therefore less than 0.1 that of typical C—C or C—N covalent bonds.

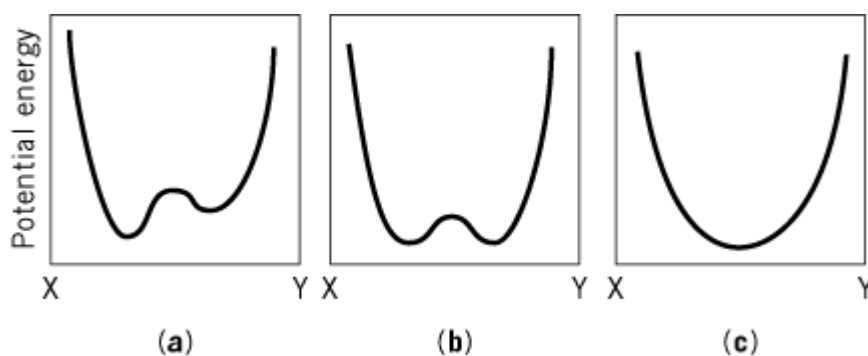
**Table 1. Thermodynamics of Hydrogen Bond Formation**

| Species                     | $-DH$ , kJ/mol | $-DS$ , e.u./mol-H bond |
|-----------------------------|----------------|-------------------------|
| Formic acid dimer (gas)     | 28             | 18                      |
| Acetic acid dimer (gas)     | 31             | 18                      |
| Propionic acid dimer (gas)  | 32             | 18                      |
| Stearic acid dimer (liquid) | 28             | 14                      |
| Water (gas)                 | 18–21          |                         |
| Water (liquid)              | 14             |                         |
| Methanol (gas)              | 13–25          |                         |
| Methanol (liquid)           | 20             |                         |
| Ethanol (gas)               | 17             |                         |
| Ethanol (liquid)            | ~17            |                         |
| Ammonia (gas)               | ~17            |                         |
| HF (gas)                    | 28–29          |                         |
| HCN (gas)                   | 14             |                         |

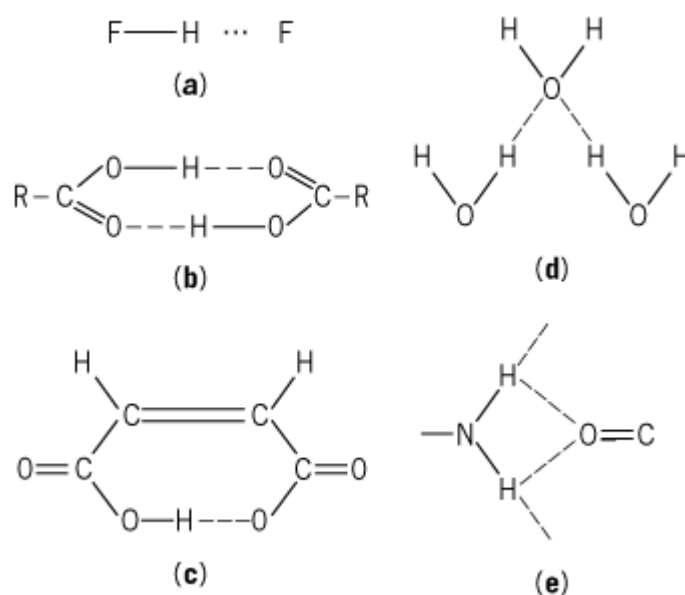
## 2.2. Hydrogen Bond Potential

Typical potential-energy curves for hydrogen bonds are presented in Figure 1. The first curve (a) demonstrates that most common unsymmetric potential, where the H atom is more strongly associated with atom X than with Y. The second curve (b) illustrates a symmetric case, where H is bonded equally to X and to Y. As the potential barrier becomes lower, the position of the hydrogen atom becomes sensitive to environmental factors. The third curve (c) illustrates the rare case where the potential-energy curve has a single minimum, and the hydrogen bond is truly centric; a good example of this may be found in intramolecular, symmetric hydrogen bonds illustrated in Figure 2.

**Figure 1.** Potential-energy curves of hydrogen-bonded systems: (a) a typical case of an unsymmetric potential, where the H atom is more strongly bonded to X; (b) a symmetric case such as Fig. 1c; and (c) a rare case of a centric potential.

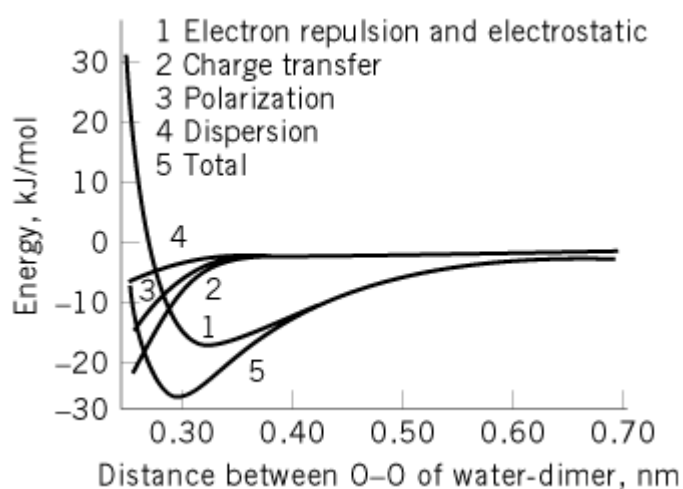


**Figure 2.** Schematic illustration of various geometries of hydrogen-bonded complexes: (a) a strong hydrogen bond between hydrogen fluoride and fluoride anion; (b) a cyclic dimer of aliphatic acids; (c) a symmetric example of intramolecular hydrogen bonding; (d) a hydrogen bond network of water molecules; (e) an example of a multi-centered hydrogen bond system. Taken from [9](#) with permission.



The Morokuma potential decomposes the total hydrogen bond energy into the following terms and evaluates the relative contribution of each term from the distance between the atoms,  $r$ : (1) electrostatic ( $r^{-1}$ ), (2) polarization ( $r^{-4}$ ), (3) exchange repulsion ( $r^{-12}$ ), (4) charge transfer  $\exp(-r)$ , and (5) dispersion interactions ( $r^{-6} + r^{-8} + r^{-10}$ ) ([1](#), [8](#)) (see [Electrostatic Interactions](#) and [van der Waals Interactions](#)). The dependence of the total hydrogen bond energy calculated for a water dimer in the gas phase is given in Figure [3](#) in terms of the individual interactions above.

**Figure 3.** Decomposition of the total hydrogen bond energy of the water dimer, demonstrating an important contribution of the electrostatic term. (Taken from Ref. [1](#) with permission.)



### 3. Hydrogen Bond Geometry



The acceptor Y atom of a hydrogen bond often has electrons in lone pairs that are directional, and consequently the hydrogen bond is also highly directional. Because of this property, the hydrogen bond is often a determining factor for the geometry of hydrogen-bonded crystals and in the conformation of polymers. In Figure 2, the geometries of a variety of hydrogen-bonded molecular complexes are presented in simplified forms. The distinction between inter- and intramolecular hydrogen bonds is important for small molecules, because the geometric constraints are more severe for the latter. For small molecules in solution, hydrogen-bonded pairs are constantly changing their counterparts, due to the low potential and activation energies of hydrogen bonds.

### 3.1. Bond Length

The distance parameters of  $X-H\cdots Y$  hydrogen bonds may be defined in several ways. When the position of the H atom is known, the distance between H and Y is defined as the hydrogen bond length. When it is not known, the distance between X and Y is a less sensitive, but in many cases still a good, guide for the presence of a hydrogen bond. For example, high resolution neutron diffraction studies indicate that nonbonded distances between two O atoms of  $<0.30$ ,  $0.34$ , and  $>0.37$  nm correspond to respective hydrogen bond probabilities of nearly 100%, about 50%, and almost 0% (1, 9). The strength of a hydrogen bond is no more than 0.1 that of a common covalent bond, so the geometries of weak to moderate hydrogen bonds are sensitive to their environment. The hydrogen bond distance  $H\cdots Y$  can vary by as much as  $\pm 20\%$  from the averages given in Table 2. For example, a compilation of more than 120 bond lengths determined by neutron diffraction studies gave a distribution of  $O-H$  distances of  $0.095-0.120$  nm, of  $H\cdots O$  from  $0.120$  to  $0.200$  nm (1).

**Table 2. Bond Lengths in Typical Hydrogen Systems**

| Hydrogen-Bonded Systems | X—H             | X—Y      | H $\cdots$ Y | Examples                                    |
|-------------------------|-----------------|----------|--------------|---|
| O—H $\cdots$ O          | $\sim 0.100$ nm | 0.275 nm | $\sim 0.17$  | Ice   |
| N—H $\cdots$ O          | $\sim 0.090$ nm | 0.299 nm | $\sim 0.21$  | Urea  |
| O—H $\cdots$ N          | $\sim 0.110$ nm | 0.278 nm | $\sim 0.17$  | Acetoxime                                   |
| N—H $\cdots$ N          | $\sim 0.090$ nm | 0.308 nm | $\sim 0.22$  | Guanine·HCl· $\frac{1}{2}$ H <sub>2</sub> O |

In the nucleosides and nucleotides, some of the H $\cdots$ Y distances are as follows (1):

|                               |                    |                    |
|-------------------------------|--------------------|--------------------|
| P—OH $\cdots$ OP              | 0.155–<br>0.169 nm |                    |
| P—OH $\cdots$ OH              | 0.165–<br>0.189 nm |                    |
| O <sub>W</sub> —H $\cdots$ OP |                    | 0.166–<br>0.188 nm |
| P—O—H $\cdots$ O <sub>W</sub> | 0.159–<br>0.168 nm |                    |
| Other O—H $\cdots$ O          | 0.174–<br>0.218 nm |                    |
| N—H $\cdots$ O=P              | 0.158–<br>0.189 nm |                    |
| N—H $\cdots$ O=C              | 0.169–<br>0.232 nm |                    |
| N—H $\cdots$ N                | 0.173–<br>0.223 nm |                    |
| NH—H $\cdots$ O=P             | 0.167–             |                    |

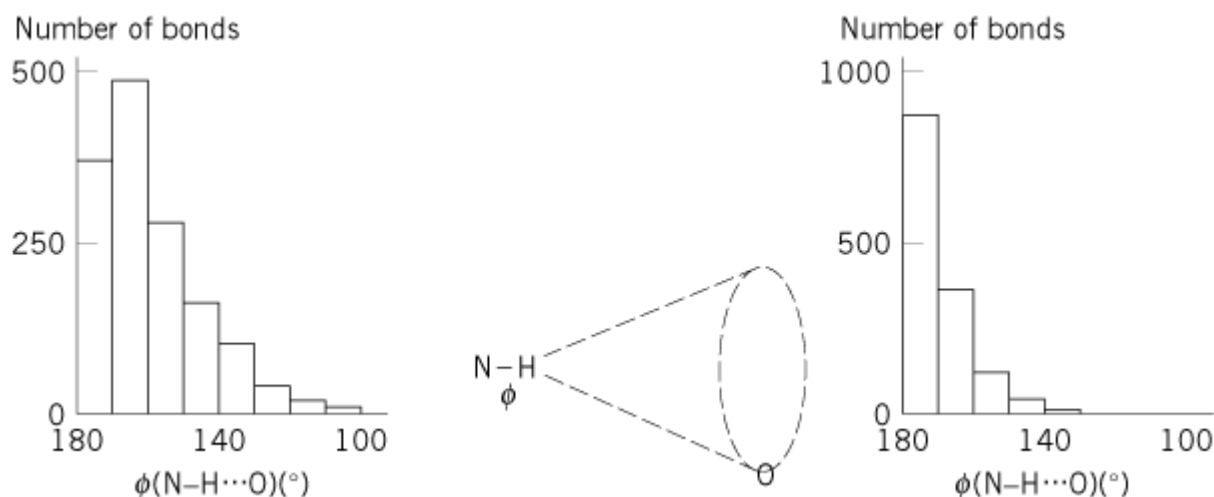
|          |          |
|----------|----------|
|          | 0.207 nm |
| NH—H¼O=C | 0.168–   |
|          | 0.273 nm |
| NH—H¼N   | 0.185–   |
|          | 0.276 nm |

---

### 3.2. Bond Angle

Although the three atoms involved in an X—H¼Y hydrogen bond tend to lie on a straight line, and unless there are two or more acceptors, the final geometry is a compromise of all the environmental factors. A hydrogen bond is therefore sometimes said to be “soft.” As a result, the angle formed by the three atoms X—H¼Y is close to 180° for a strong hydrogen bond, and the flexibility about this angle increases as the bond becomes weaker. Bifurcated hydrogen bonds cannot be linear. Figure 4 (10) illustrates the distribution of X—H¼Y angles.

**Figure 4.** The distribution of hydrogen bond angles. The three atoms involved in the formation of a hydrogen bond have preferred configuration for the hydrogen bond to be most effective. The distribution on the left is normalized with respect to the statistical factor of  $2 \sin \phi$ . (Taken from Ref. 9 with permission.)



A very good example of the geometry of a hydrogen-bonded dimer is found in acetic acid in the gas phase, which is similar to scheme **b** of Figure 2. The O—H¼O angle is 180°, and the O—H and H¼O bond lengths are 0.104 and 0.164 nm, respectively.

## 4. Hydrogen Bonds in Molecular Biology

### 4.1. Hydrogen Bonds in Water

Water is well known to differ from the liquids formed by comparable molecules, such as H<sub>2</sub>S, in its higher melting and boiling temperatures, its greater heat of fusion and of evaporation, and in its greater viscosity, properties that are important for life on earth. The [hydrophobic effect](#) that is one aspect of the interactions between [nonpolar](#) molecules and groups in water is a result of the physical properties of hydrogen-bonded networks in liquid water. When nonpolar molecules dissolve in water, they create a cavity within the hydrogen-bonded network of water molecules, which disrupts the network and forces a rearrangement of the water molecules; this reduces the entropy of the system. It is thermodynamically preferable for the nonpolar molecules to associate, thereby reducing the total area of the water–solute interface and minimizing the disruption of the hydrogen-bonded network of water.

#### 4.2. Hydrogen Bonds in Protein Structure

Hydrogen bonds in folded [protein structures](#) are formed primarily (about 90%) between the N—H and C=O groups of the polypeptide backbone, and about 85% of these are in either [alpha-helix](#) or [beta-sheet](#). In the  $\alpha$ -helix, hydrogen bonds are formed between peptide groups every four residues apart; in the case of  $\beta$ -sheets, every two residues. The average hydrogen bond length (H $\cdots$ O) is shorter for  $\beta$ -sheets (0.197 nm) than for the  $\alpha$ -helix (0.206 nm). Three-centered, bifurcated hydrogen bonds are observed in helices, but are not as common as in small molecules. Polar side chains with hydrogen atoms, such as those of [tyrosine](#), [serine](#), [threonine](#), [lysine](#), [arginine](#), [glutamic acid](#), [aspartic acid](#), [glutamine](#), and [asparagine](#), can act as hydrogen bond donors. Turn structures, such as [beta-](#) and [gamma-turns](#), are also stabilized by main-chain hydrogen bonds ([11](#), [12](#)). The concept of low-barrier hydrogen bonds has been introduced to explain the fast and efficient proton transfer within and between proteins ([13](#)), but their presence, although persuasive in the gas phase, is still controversial and not supported by an NMR study ([14](#)).

#### 4.3. Hydrogen Bond Network in the Active Site of Serine Proteinases

A hydrogen bond network was proposed in the [active sites](#) of [serine proteinases](#) on the basis of the X-ray crystallographic structure of [chymotrypsin](#), so as to explain the unusual chemical reactivity of the serine residue at its active site ([15](#), [16](#)). At neutral pH, the hydroxyl group of an ordinary serine residue is not able to carry out an electrophilic attack on the carbonyl carbon of [peptide bonds](#). That the active site Ser195 residue is in an activated state because of its participation in a hydrogen bond network with Asp102 and His57 has been well accepted. In the process of hydrolyzing a substrate peptide bond, the carbonyl group of the susceptible bond forms a covalent bond with the hydroxyl group of Ser195, to form an acyl intermediate, cleaving the peptide bond, and permitting the amino-terminus of the carboxyl-terminal portion of the peptide to form a hydrogen bond with His57. The amino group is then replaced by a water molecule, which will cleave the covalent bond to Ser195 and the amino-terminal portion of the substrate will dissociate, to complete the [proteolysis](#) reaction.

Hydrogen bond networks are also present in the active sites of [carboxyl proteinases](#) and in the functions of [DNA-binding proteins](#), such as [zinc fingers](#). The importance of hydrogen bonds in the subunit interactions of [oligomeric proteins](#) has been discussed thoroughly by Perutz ([17](#)).

#### 4.4. Hydrogen Bonds in Nucleic Acids

The importance of hydrogen bonds in **nucleic acids** is twofold. Hydrogen bonds between the bases stabilize the **double-helix** structures of DNA and RNA, but a more important aspect is that hydrogen bonds are formed only between specified pairs of adenine:thymine, A:T (uracil for RNA) and guanine:cytosine, G:C. Formation of hydrogen bonds only between such specified pairs permits the basic mechanism of transmission of genetic information in the nucleic acids. The biological significance of the Watson–Crick model of DNA was in the specific hydrogen-bonding patterns between A:T and G:C complementary nucleotide pairs, which have been amply verified. There are other possibilities of hydrogen bond patterns among the four bases, but, under physiological conditions, the **Watson–Crick base pairs** are the most stable. An important aspect is that the G:C base pair is more stable than the A:T pair, in part because of the greater number of hydrogen bonds.

#### 4.5. Hydrogen Bonds in Carbohydrates and Lipids

Carbohydrate molecules have a large number of —OH groups, all of which are candidates for forming hydrogen bonds. There is no doubt that the conformations of carbohydrates, and especially those of polymeric carbohydrates, are strongly influenced by hydrogen bond formation. Formation of hydrogen-bonded (often three-centered) structures in their fibrous and crystalline states is discussed by Jeffrey ([1](#)). The functional groups in lipid molecules that can form hydrogen bonds are limited to the hydrophilic head groups. Consequently, the conformational variation of lipid molecules through hydrogen bond formation is rather limited, compared to the preceding examples.

### 5. Prospectives

The importance of hydrogen bonds in chemistry and especially in biology is increasingly appreciated

because of the advance of our knowledge of the natural and synthetic multimolecular systems known as supra-molecular systems. These complex structures are stabilized by multitudes of weak secondary interactions such as **van der Waals**, **electrostatic**, chelating, and hydrogen bonding interactions. The importance of supramolecular structures resides in the fact, as well as in our expectation, that a structure held together weakly and therefore dynamically will exert more subtle functions than will other developed molecules and molecular systems. The hydrogen bond is based on neutral groups and strongly controls the geometry of the resulting molecular associations, which will be important in controlling the structure and activity of synthetic and hybrid structures once the nature of hydrogen bond is understood more precisely. New methods are emerging for the study of hydrogen bonded systems with completely different perspectives. For example, the strength and specificity of hydrogen bonding can now be measured using [scanning probe techniques](#) applied to a group of molecules (18) and at the single-molecule level (19).

## Bibliography

1. G. A. Jeffrey (1997) *An Introduction to Hydrogen Bonding*, Oxford Univ. Press, New York.
2. G. C. Pimentel and A. L. McClellan (1960) *The Hydrogen Bond*, Freeman, San Francisco.
3. L. C. Pauling (1960) *The Nature of Chemical Bond and the Structure of Molecules and Crystals: An Introduction to Modern Structural Chemistry*, Cornell Univ. Press, Ithaca.
4. P. Schuster, G. Zundel, and C. Sandorfy (1976) *The Hydrogen Bond: Recent Developments in Theory and Experiments*: Vol. 1 I, *Theory*; Vol. 2, *Structure and Spectroscopy*; Vol. 3, *Dynamics, Thermodynamics and Special Systems*, North-Holland, Amsterdam.
5. G. A. Jeffrey and W. Saenger (1991) *Hydrogen Bonding in Biological Structures*, Springer-Verlag, Berlin.
6. C. W. Hilbers (1979) "Hydrogen-bonded proton exchange and its effect on NMR spectra of nucleic acids", in *Biological Applications of Magnetic Resonance*, R. G. Shulman, ed., Academic Press, New York, Chapter "1", pp. 1–44.
7. R. Taylor and O. Kennard (1982) *J. Am. Chem. Soc.* **104**, 5063–5070.
8. U. C. Singh and P. A. Kollman (1985) *J. Chem. Phys.* **83**, 4033–4040.
9. C. Ceccarelli, G. A. Jeffrey, and A. Taylor (1981) *J. Mol. Struct.* **70**, 255–271.
10. R. Taylor, O. Kennard, and W. Versichel (1984) *J. Am. Chem. Soc.* **106**, 244–248.
11. S. N. Vinogradov (1980) "Structural aspects of hydrogen bonding in amino acids, peptides, proteins, and model systems", in *Molecular Interactions*, H. Ratajczak and W. J. Orville-Thomas, eds., Wiley, New York, Vol. 2, pp. 179–229.
12. G. E. Schulz and R. H. Shirmer (1979) *Principles of Protein Structure*, Springer-Verlag, New York.
13. R. A. Copeland and S. I. Chan (1989) Proton translocation in proteins, *Annu. Rev. Phys. Chem.* **40**, 671–698.
14. E. L. Ash, J. L. Sudmeier, E. C. De Fabo, and W. W. Bachovchin (1997) A low-barrier hydrogen bond in the catalytic triad of serine proteases? Theory versus experiment, *Science* **278**, 1128–1132.
15. D. M. Blow, J. J. Birktoft, and B. S. Hartley (1969) Role of a buried acid group in the mechanism of action of chymotrypsin, *Nature* **221**, 337–340.
16. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland Publishing, New York, p. 237.
17. M. Perutz (1992) *Protein Structure: New Approaches to Disease and Therapy*, Freeman, New York.
18. A. Noy, D. V. Vezenov, and C. M. Lieber (1997) Chemical force microscopy, *Annu. Rev. Mat. Sci.* **27**, 381–421.
19. A. Ikai (1996) STM and AFM of bio/organic molecules and structures, *Surf. Sci. Rep.* **26**, 261–332.

## Hydrogen Exchange

Hydrogen exchange is the interchange of hydrogen atoms between a molecule of interest, usually a [peptide](#), [protein](#), or [nucleic acid](#), and [water](#). Usually the exchange is detected because the hydrogens are different isotopes, the normal protium ( $^1\text{H}$ ), or deuterium ( $^2\text{H} = \text{D}$ ), or tritium ( $^3\text{H} = \text{T}$ ).

Designating the two isotopes as H and H', the remainder of the macromolecule as M, and O as the oxygen atom of water, the reaction can be represented as



Hydrogen exchange experiments provide information about the flexibilities of biological macromolecules. These are not static structures but exhibit a flexibility that is often important in their function. They undergo spontaneous thermal fluctuations and folding and unfolding processes, and they may display variations of their conformations. These motions vary with the binding of other molecules, with **allosteric** interactions, and with other environmental influences. All of these affect the rates of hydrogen exchange, whose measurement then helps to elucidate the dynamic structure of the molecule. For example, it is possible to measure the rate of base-pair opening in the DNA **double helix**, which is prerequisite for **replication** and [transcription](#).

Rates of exchange of different types of hydrogen atoms are quite variable. Those attached to oxygen atoms of sugars in carbohydrates or of [serine](#), [threonine](#), and [tyrosine](#) residues in proteins exchange so rapidly that they are difficult to measure. Exchange of hydrogens attached to carbon is too slow. Those of most experimental interest are attached to nitrogen, as in the CONH peptide backbone or the [asparagine](#) and [glutamine](#)  $\text{CONH}_2$  side chains of peptides and proteins and in the amino ( $\text{NH}_2$ ) and imino (NH) groups of the nucleic acid bases **adenine**, **thymidine**, **uracil**, **guanine**, and **cytosine**. These exchange in milliseconds, or even faster with appropriate catalysts, when free and accessible to the solvent, but much slower, by as much as  $10^{12}$ -fold, when buried in the interior of a macromolecule and inaccessible to solvent. These retardations reveal information about the stability of the folded macromolecular structure. The rate at which exchange actually occurs provides information about [molecular dynamics](#).

### 1. Rates and Mechanisms of Chemical Exchange in Model Small Molecules

The rate of exchange  $R_{\text{ex}}$  or the rate of disappearance of MH and the formation of MH' [Eq. (1)] is given by Eq. (2), where  $[\text{MH}]_{\text{eq}}$  and  $[\text{MH}' ]_{\text{eq}}$  are the equilibrium concentrations of MH and MH' and  $k_{\text{ex}}$  is the rate constant for the reaction (see [Molecular Dynamics](#)). For the amide groups of peptides and proteins, it is observed that  $k_{\text{ex}}$  has the form in Eq. (3), corresponding to acid catalysis by  $\text{H}^+$  and base catalysis by  $\text{OH}^-$ . The catalytic constants are  $k_{\text{H}}$  and  $k_{\text{OH}}$ , respectively, and a term that is usually too small to detect, except in solvents less polar than water, corresponds to water catalysis. For the NH group of a typical amide,  $k_{\text{H}} \sim 10^{-1} \text{M}^{-1} \text{s}^{-1}$  and  $k_{\text{OH}} \sim 10^7 \text{M}^{-1} \text{s}^{-1}$ . A graph of  $\log k_{\text{ex}}$  versus pH is shown in Fig. 1. This shows a minimum given by Eq. (4), where  $K_{\text{w}} = [\text{H}^+][\text{OH}^-] = 10^{-14} \text{M}^2$ .

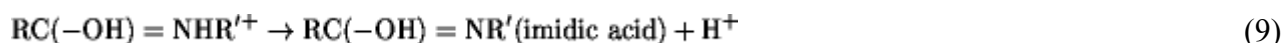
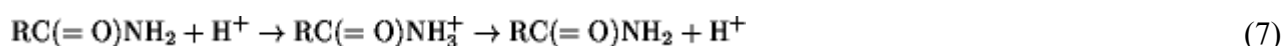
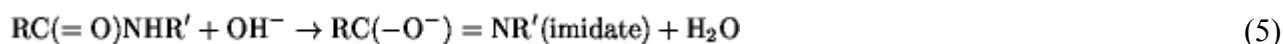
$$R_{\text{ex}} = -\frac{d([\text{MH}] - [\text{MH}]_{\text{eq}})}{dt} = \frac{d([\text{MH}'] - [\text{MH}']_{\text{eq}})}{dt} \quad (2)$$

$$= k_{\text{ex}}([\text{MH}] - [\text{MH}]_{\text{eq}})$$

$$k_{\text{ex}} = k_{\text{H}}[\text{H}^+] + k_{\text{OH}}[\text{OH}^-] + k_{\text{w}} \quad (3)$$

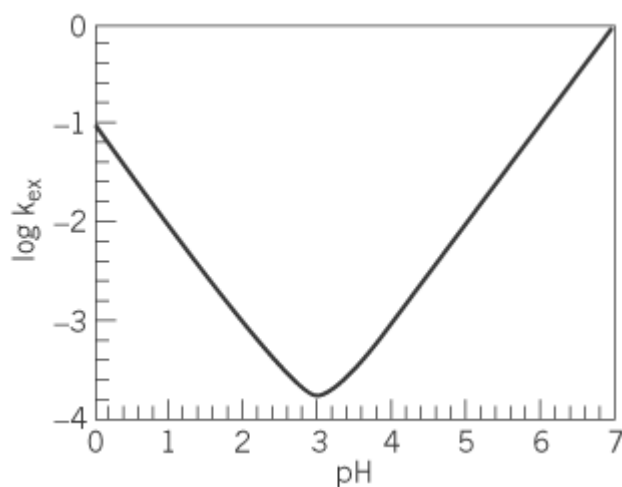
$$pH_{\text{min}} = \frac{1}{2} \log \left( \frac{k_{\text{H}}}{K_{\text{w}}k_{\text{OH}}} \right) \quad (4)$$

The mechanism of the base-catalyzed reaction involves removing the amide NH hydrogen atom [Eq. (5)], to create the imidate anion (the conjugate base of the amide), which then abstracts an  $\text{H}^+$  from water to regenerate the amide [Eq. (6)]. The mechanism of the acid-catalyzed reaction depends on the amide. For the  $\text{CONH}_2$  side chains of asparagine and glutamine residues, the exchange occurs simply by transient protonation of the nitrogen, which produces an unstable conjugate acid of the amide, followed by removal of a different hydrogen as  $\text{H}^+$  [Eq. (7)]. For backbone CONH groups, exchange occurs by attachment of  $\text{H}^+$  to the oxygen atom, which is much more basic than the nitrogen (see [Kinetics](#)), and produces another conjugate acid [Eq. (8)], followed by removal of  $\text{H}^+$  from nitrogen, to produce the imidic acid, the unstable **tautomer** of an amide [Eq. (9)]. Attachment of a solvent  $\text{H}^+$  to the nitrogen, followed by removal of  $\text{H}^+$  from oxygen, regenerates the amide, but with its H exchanged [Eq. (10)].



The rates vary with the substitution pattern in the amide (1). Table 1 lists logarithms of relative rate constants  $k_{\text{H}}$  and  $k_{\text{OH}}$  for various amino acid residues (2). Electron-withdrawing groups increase  $k_{\text{OH}}$  because they stabilize the negative charge that is created in the imidate anion. They also decrease  $k_{\text{H}}$  because they destabilize the positive charge created in either mechanism. As a result,  $pH_{\text{min}}$  [Eq. (4)] is decreased and the minimum in Fig. 1 moves to the left. Rates also increase with increasing temperature. In acid, the [activation energy](#) is 14 kcal/mol. In base, at constant pH, the apparent activation energy is 17.5 kcal/mol, but this greater value largely results from the temperature-dependence of  $K_{\text{w}}$ .

**Figure 1.** pH dependence of the value of the intrinsic exchange rate constant  $k_{\text{ex}}$  for the backbone amide of a typical peptide.



**Table 1. Rate Constants for Acid- and Base-Catalyzed Hydrogen Exchange in  $\text{CH}_3\text{C}(=\text{O})\text{NH}_a\text{CHRC}(=\text{O})\text{NH}_b\text{CH}_3$ , relative to  $\text{R}=\text{CH}_3$ <sup>a</sup>**

| R  | $\log_{10}k_{\text{H},a}$ | $\log_{10}k_{\text{H},b}$ | $\log_{10}k_{\text{OH},a}$ | $\log_{10}k_{\text{OH},b}$ |
|--|---------------------------|---------------------------|----------------------------|----------------------------|
| $\text{CH}_3$  | $\equiv 0$                | $\equiv 0$                | $\equiv 0$                 | $\equiv 0$                 |
| $((\text{CH}_2)_3\text{NHC}(\text{NH}_2)_2^+$          | -0.59                     | -0.32                     | 0.08                       | 0.22                       |
| $\text{CH}_2\text{CONH}_2$                             | -0.58                     | -0.13                     | 0.49                       | 0.32                       |
| $\text{CH}_2\text{CO}_2^-$                             | 0.9                       | 0.58                      | -0.30                      | -0.18                      |
| $\text{CH}_2\text{COOH}$                               | -0.9                      | -0.12                     | 0.69                       | 0.6                        |
| $\text{CH}_2\text{SH}$                                 | -0.54                     | -0.46                     | 0.62                       | 0.55                       |
| $\text{CH}_2\text{S}^-)_2$ (Cys <sub>2</sub> )         | -0.74                     | -0.58                     | 0.55                       | 0.46                       |
| H  | -0.22                     | 0.22                      | 0.27                       | 0.17                       |
| $(\text{CH}_2)_2\text{CONH}_2$                         | -0.47                     | -0.27                     | 0.06                       | 0.20                       |
| $(\text{CH}_2)_2\text{CO}_2^-$                         | -0.9                      | 0.31                      | -0.51                      | -0.15                      |
| $(\text{CH}_2)_2\text{COOH}$                           | -0.6                      | -0.27                     | 0.24                       | 0.39                       |
| $\text{CH}_2\text{Im}$ (His)                           | -                         | -                         | -0.10                      | 0.14                       |
| $\text{CH}_2\text{ImH}^+(\text{His} \cdot \text{H}^+)$ | -0.8                      | -0.51                     | 0.8                        | 0.83                       |
| $\text{CH}(\text{CH}_3)\text{CH}_2\text{CH}_3$         | -0.91                     | -0.59                     | -0.73                      | -0.23                      |
| $\text{CH}_2\text{CH}(\text{CH}_3)_2$                  | -0.57                     | -0.13                     | -0.58                      | -0.21                      |

|   |       |                        |       |       |
|---|-------|------------------------|-------|-------|
| $(\text{CH}_2)_4\text{NH}_3^+$  | -0.56 | -0.29                  | -0.04 | 0.12  |
| $(\text{CH}_2)_2\text{SCH}_3$   | -0.64 | -0.28                  | -0.01 | 0.11  |
| $\text{CH}_2\text{Ph}$  | -0.52 | -0.43                  | -0.24 | 0.06  |
| <i>Cis</i> - $(\text{CH}_2)_3\text{N}_a$ (Pro) ( $\text{H}_a$ absent)   | -0.19 | ( $\text{H}_a$ absent) | -0.24 |       |
| <i>Trans</i> - $(\text{CH}_2)_3\text{N}_a$ (Pro) ( $\text{H}_a$ absent) | -0.85 | ( $\text{H}_a$ absent) | 0.60  |       |
| $\text{CH}_2\text{OH}$  | -0.44 | -0.39                  | 0.37  | 0.30  |
| $\text{CH}(\text{OH})\text{CH}_3$                                       | -0.79 | -0.47                  | -0.07 | 0.20  |
| $\text{CH}_2\text{Indol}$ (Trp)   | -0.40 | -0.44                  | -0.41 | -0.11 |
| $\text{CH}_2\text{C}_6\text{H}_4\text{OH}$ (Tyr)                        | -0.41 | -0.37                  | -0.27 | 0.05  |
| $\text{CH}(\text{CH}_3)_2$  | -0.74 | -0.30                  | -0.70 | -0.14 |
| N-term <sup>b</sup>   | -     | -1.32                  | -     | 1.62  |
| C-term <sup>c</sup>   | 0.96  | -                      | -1.8  | -     |
| C-term <sup>d</sup>   | ~0.05 | -                      | -     | -     |

<sup>a</sup> Ref. 2.

<sup>b</sup> For  $^+\text{H}_3\text{NCH}_2\text{CONHCH}_3$ .

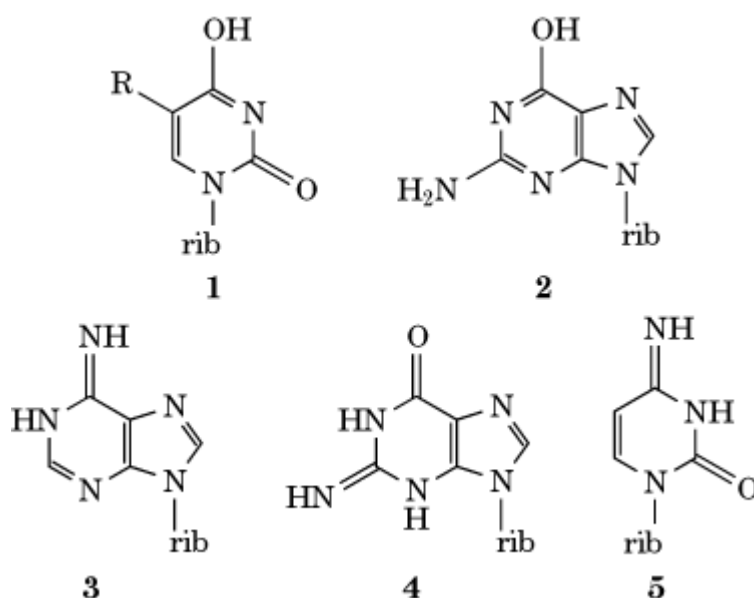
<sup>c</sup> For  $\text{CH}_3\text{C}(=\text{O})\text{NHCH}_2\text{CO}^-_2$ .

<sup>d</sup> For  $\text{CH}_3\text{C}(=\text{O})\text{NHCH}_2\text{COOH}$ .

The mechanisms of exchange in nucleic acids are analogous, except that the equivalent of Eq. (7) never contributes. The imino NH groups of uracil, thymidine, and guanine are sufficiently acidic that they can be removed by  $\text{OH}^-$  at every encounter, in a [diffusion-controlled reaction](#) that has a rate constant  $k_{\text{OH}} \sim 10^{10} \text{M}^{-1} \text{s}^{-1}$ . Weaker bases, such as components of the [buffer](#), also remove these hydrogen atoms. Acid-catalyzed exchange proceeds by protonation of oxygen, followed by deprotonation of the NH, to produce the iminol tautomer (**1**, R = H or  $\text{CH}_3$ , or **2**; Fig. 2) as a transient intermediate. Acid catalysis for guanine is masked by conversion to the conjugate acid. Exchange of an amino  $\text{NH}_2$  group of adenine, guanine, or cytosine is catalyzed by acid and base, and there is also a pH-independent rate near neutrality. The base-catalyzed reaction involves slow removal of an NH by  $\text{OH}^-$ . The acid-catalyzed reaction involves attachment of  $\text{H}^+$  to a ring nitrogen, followed by removal of one of the amino hydrogen atoms by water, to produce an abnormal tautomer (**3**, **4**, **5**; Fig. 2) as a transient intermediate. In the pH-independent region, that NH is removed by  $\text{OH}^-$  instead of water, to form the same intermediate. In each of these substrates, the acid catalysis is masked at low pH by complete conversion of substrate to its conjugate acid.

**Figure 2.** Abnormal tautomers of nucleic acid bases that are transient intermediates in hydrogen exchange.





Because charged species, either cationic or anionic, must be formed as intermediates, the local environment affects the rate of hydrogen exchange. For example, nearby charges stabilize or destabilize the charged intermediate by [electrostatic interactions](#). These charges include the phosphate backbone of nucleic acids and the side chains of acidic or basic residues in proteins. However, their effect is moderated by ions of salts that screen the interaction. As a result, the intrinsic rates are affected by the ionic strength of the solution. Moreover, the charges on nearby sites titrate with pH, leading to a variability of  $k_{\text{ex}}$  that is more complicated than in Fig. 1. Another effect on rates arises in less polar solvents, where those charged intermediates are less stable and lead to a reduction in rate. This may also contribute to a retardation of exchange in the [nonpolar, hydrophobic](#) interior of a protein. It is widely believed that the major retardation of exchange in biological macromolecules occurs when the hydrogen atom is involved in a [hydrogen bond](#) to some adjacent group. Generally, that hydrogen bond must be broken to expose the NH to solvent and render it susceptible to exchange. But simple interpretation is confounded by the fact that hydrogen bonding and inaccessibility usually occur simultaneously in folded macromolecules (see [Nonpolar](#)).

## 2. Mechanisms of Chemical Exchange in Macromolecules

An NH group at the surface of a macromolecule and hydrogen-bonded only to water has an intrinsic exchange rate characteristic of a model small molecule. Hydrogen atoms that are protected in the interior of a macromolecule and by internal hydrogen bonding exchange more slowly. The degree of this retardation provides information about the stability of the macromolecular structure and the flexibility that permits exchange.

The process by which inaccessible hydrogen atoms gain access to water so that they can undergo exchange is still a matter of debate, but there are two limiting and extreme models. According to the solvent-penetration model, water and catalytic  $\text{H}^+$  or  $\text{OH}^-$  gain access to the interior of the macromolecule through channels, perhaps generated transiently by conformational fluctuations. Then the rate of exchange is governed by the extent to which those buried hydrogen atoms come into contact with the solvent atoms. According to the other model, called the local-unfolding model, cooperative motion of a segment of the molecule breaks a set of hydrogen bonds and other interactions that hold the macromolecule in its stable conformation. Thereby, the NH groups of that segment are exposed to solvent, where they can undergo exchange. This model agrees nicely with the dynamic nature of macromolecular structure, although it does not explain why exchange is not decreased substantially by incorporating the macromolecule into a crystal lattice, and it is an

oversimplification to assume a single unfolded state from which exchange occurs. Macromolecules probably achieve a variety of unfolded states, from those involving only minor local fluctuations to those involving global unfolding.

The local-unfolding model can be analyzed kinetically, as in Eq. (11). Here  $k_{\text{unfold}}$  is the rate constant for the motion that converts the closed to the open form,  $k_{\text{fold}}$  is the rate constant for reversion back to the closed form, and  $k_{\text{ex}}$  is the rate constant for hydrogen exchange, as in Eqs. (2), (3). The general expression for the observed rate constant is given by Eq. (12). Two limiting behaviors can be distinguished:



$$k_{\text{obs}} = \frac{k_{\text{unfold}} k_{\text{ex}}}{k_{\text{fold}} + k_{\text{ex}}} \quad (12)$$

If  $k_{\text{ex}} \gg k_{\text{fold}}$ , then  $k_{\text{obs}}$  [Eq. (12)] simplifies to  $k_{\text{unfold}}$ . The rate-limiting step is the opening to produce the locally unfolded state, from which exchange is rapid. Under these conditions, called “EX1,” the rate is independent (except insofar as  $k_{\text{unfold}}$  is pH-dependent) of the concentrations of acid or base catalyst, which enter into Eq. (3). In principle, these conditions can be induced at high pH, where  $k_{\text{ex}}$  becomes sufficiently large. In practice, it is rare to achieve these conditions for a protein without **denaturation**, and exchange generally occurs from the totally unfolded conformation. In contrast, this is possible for DNA, and extrapolation to high pH or catalyst concentration provides  $k_{\text{unfold}}$ , the rate at which the double helix breaks the hydrogen bonds and allows an NH group to exchange.

If  $k_{\text{fold}} \gg k_{\text{ex}}$ , then  $k_{\text{obs}}$  [Eq. (12)] simplifies to  $K_{\text{op}} k_{\text{ex}}$ , where  $K_{\text{op}} = k_{\text{unfold}}/k_{\text{fold}}$  is the equilibrium constant, [open]/[closed], for forming the locally unfolded state. These conditions are called “EX2”, and they correspond to exchange that occurs only during that fraction of the time when the macromolecule is unfolded. Because  $k_{\text{ex}}$  can be approximated by the value (including acid or base catalysis, according to pH) for an appropriate model small molecule (1, 2), then the experimentally determined value of  $k_{\text{obs}}$  provides an estimate of  $K_{\text{op}}$ . Moreover, conversion to the corresponding **free energy**,  $DG^\circ = -RT \ln K_{\text{op}}$ , provides a measure of the energy that stabilizes the closed native state of the macromolecule or of the energy required to allow the local unfolding. Alternatively,  $k_{\text{ex}}/k_{\text{obs}}$  ( $= 1/K_{\text{op}}$ ) can be called a protection factor P, which represents the extent to which exchange is retarded in the native state.

### 3. Applications

Study of hydrogen exchange has the advantage that the measurement need not perturb the macromolecule and involves only the normal solvent, water. Separate samples can be exposed to a labeling pulse for varying time periods. Then the exchange is quenched by adjusting the pH to trap the label, and the subsequent time course of the loss of label from each site is monitored. The pulse conditions can be chosen to label selectively the sites of interest for the particular application.

Hydrogen exchange has long been used to provide information about the [tertiary structures](#) of proteins and about protein dynamics. Characteristic patterns of protection factors are associated with hydrogen atoms in **a-helices** or **b-sheets**. To the extent that hydrogen exchange also involves cooperative unfolding of an entire structural **domain**, all of the hydrogens in that domain should show the same protection factor.

Other aspects that have been investigated include the [quaternary structure](#) of proteins and their interactions with small molecules and other macromolecules. For example, selective labeling allows recognizing the segments involved in allosteric transitions in such proteins as [hemoglobin](#). Because exchange of some NH atoms requires local unfolding, the observed rates also provide information about the stability of the native structure and the way it responds to changes in the environment. Intermolecular interactions that have been probed include protein–peptide interactions, enzyme–substrate and enzyme–inhibitor binding, protein–protein interactions, such as those between [cytochrome c](#) and cytochrome c peroxidase or between a protein and an [antibody](#) or a [chaperonin](#), and protein–nucleic acid interactions.

One of the most recent applications is characterization of transient intermediates in **protein folding**. Exposed hydrogen atoms in a deuterated protein that is undergoing refolding are labeled during a rapid pulse with normal water. Then the pattern of a normal hydrogen isotope trapped in the refolded native protein can be analyzed to infer the structure of the intermediate. This method has been applied to several proteins, including **ribonuclease A**, cytochrome *c*, **barnase**, and hen egg white [lysozyme](#).

For duplex DNA, it is possible to measure  $k_{\text{unfold}}$ , the rate of base-pair opening. Most rate constants measured are between  $2 \times 10^1$  and  $2 \times 10^3 \text{s}^{-1}$ , but lower values are evidence of unusual structural features. It is also possible to study the influence of **intercalators** and of **triple helices**.

#### 4. Methods of Measurement

The simplest method for detecting hydrogen exchange is measuring the overall extent of exchange by measuring the isotopic content of the macromolecule. This is achieved by counting the [radioactivity](#) of tritium or by UV-visible or [vibrational spectroscopy](#) that is sensitive to isotopic content. These methods measure only the total content of isotope, however, and cannot distinguish among the various hydrogen atoms in the macromolecule.

Methods involving various aspects of nuclear magnetic resonance ( [NMR](#) ) distinguish the individual hydrogen atoms in a molecule, and methods for assigning each signal have been developed. Two-dimensional Fourier-transform NMR is especially powerful at spreading out the spectrum so that each site is resolved. Pulse methods have been developed to select only particular nuclei, such as those attached to  $^{15}\text{N}$ . Then the simplest method for measuring rates of exchange is transferring the sample to  $\text{D}_2\text{O}$  and watching the disappearance of each NH signal.

Three other techniques utilize the unique capability of NMR to detect exchange and measure rates under equilibrium conditions, when no net reaction is proceeding. The derivatives in Eq. (2) are zero, but an NH remains in an MH site for an average time of  $1/k_{\text{ex}}$ . (1) The signal of an NH that is exchanging with solvent  $\text{H}_2\text{O}$  displays an NMR broadening that increases with the exchange rate, until the signal becomes so broad that it disappears into the base line. (2) The signal of a CH adjacent to an NH is split into a doublet. As the NH exchange rate increases, the doublet components broaden, overlap, and coalesce into a broad single peak that then sharpens. This method provides rate constants for small molecules but is much less suitable for macromolecules, where the CH signals are not resolvable. (3) Various pulse sequences also label  $\text{H}_2\text{O}$  or individual NH signals with a characteristic magnetization (saturation or inversion). As hydrogen atoms exchange among these sites, the magnetization is transferred from one site to another. Analysis of the magnitude and time course of these transfers provides the rate constants for exchange. All of these techniques have the further advantage of not being limited by rapid exchange during mixing.

[Mass Spectrometry](#) is used to measure the overall exchange process because the different isotopes slightly alter the molecular weight of the macromolecule. It has only recently been applied to

measuring hydrogen exchange rates at individual sites because it is necessary to cleave the macromolecule into small pieces that can be analyzed individually. This must be done under conditions that preserve the label, using enzymes that function at 0°C and at pH<sub>min</sub> (Fig. 1), and then the fragments are separated and analyzed for isotopic content.

### Bibliography

1. R. S. Molday, S. W. Englander, and R. G. Kallen (1972) *Biochemistry* **11**, 150–158.
2. Y. Bai, J. S. Milne, L. Mayne, and S. W. Englander (1993) *Proteins* **17**, 75–86.

### Suggestions for Further Reading

3. S. W. Englander and N. R. Kallenbach (1984) Hydrogen exchange and structural dynamics of proteins and nucleic acids. *Q. Rev. Biophys.* **16**, 521–655 (an extensive, broad review).
4. C. L. Perrin (1989) Proton exchange in amides: Surprises from simple systems, *Acc. Chem. Res.* **22**, 268–275 (review of mechanisms of exchange in models, with NMR methods).
5. C. K. Woodward (1994) Hydrogen exchange rates and protein folding, *Curr. Opinion Struct. Biol.* **4**, 112–116.
6. M. Guéron and J.-L. Leroy (1995) Studies of base pair kinetics by NMR measurement of proton exchange, *Methods Enzymol.* **261**, 383–413. (applications to nucleic acids).
7. S. W. Englander, T. R. Sosnick, J. J. Englander, and L. Mayne (1996) Mechanisms and uses of hydrogen exchange. *Curr. Opinion Struct. Biol.* **6**, 18–23 (brief recent review including methods and applications).

## Hydrogen Isotopes

Hydrogen is element number 1 in the periodic table and the most abundant element in the universe. Hydrogen has a valence state of 1 (1); it occurs mainly in nature in combination with oxygen as [water](#) in all living organisms and is present in all organic matter.

Three isotopes of hydrogen exist in nature: <sup>1</sup>H (protium, with an abundance of 99.985%), <sup>2</sup>H (deuterium, with an abundance of 0.015%), and <sup>3</sup>H (*tritium*, which exists naturally only in trace amounts) (see [Radioactivity](#) and [Radioisotopes](#)). Stable deuterium, which was first prepared by Urey (1), is used as heavy water (<sup>2</sup>H<sub>2</sub>O or D<sub>2</sub>O) for contrast variation in **neutron scattering**, altering the density of aqueous solvents, and in studying [hydrogen exchange](#) reactions in macromolecules, as well as in moderating (slowing down) neutrons in some nuclear reactors. It is also used as a proton source in nuclear magnetic resonance spectroscopy and imaging.

Among the hydrogen isotopes, only tritium is **radioactive**, with a half-life of 12.35 years (2). It decays by beta-minus emission to <sup>3</sup>H, which is stable. Tritium yields one beta particle per decay, with an energy of 0.0186 MeV maximum, 0.00568 MeV on average. Tritium is produced in a nuclear reactor by irradiation of enriched lithium targets according to the reaction <sup>6</sup>Li(*n,a*)<sup>3</sup>H, or by irradiation of <sup>3</sup>H gas targets according to the reaction <sup>3</sup>He(*n,p*)<sup>3</sup>H. Tritium is used as an ingredient in thermonuclear weapons, in luminous lighting (airport runway lights, exit signs and luminous paints), and as a tracer in biomedical studies.

Tritium was first made available for biological research by the U.S. Atomic Energy Commission in

1948. The first tritium uptake studies in humans were conducted in 1951. [Streptomycin](#) and tetracycline were first labeled with tritium (2), and tritiated water was first used to measure total-body water in 1956 (3). Hundreds of different tritiated compounds are now commercially available in generous amounts and with high specific activity from commercial suppliers. Methods for tritium labeling typically involve exchange, substitution, or addition reactions, such as exchange by heterogeneous catalysis, exchange by homogeneous catalysis, exchange catalyzed by radiation, substitution by chemical reduction, or hydrogenation.

The short range of tritium's weak beta particle, its low toxicity, relatively low cost, and usefulness for high resolution [autoradiography](#) make it a convenient choice for tracer studies. A frequently encountered problem is pH-dependent loss of tritium by exchange with solvent water.

In molecular biology, tritium is widely used as a tracer and label for **nucleic acid** precursors. One of the most common tritiated compounds used is  $^3\text{H}$ -**thymidine**, which may be introduced into the laboratory animal by direct injection or added to excised tissues. The metabolism and uptake of  $^3\text{H}$ -thymidine is used to monitor [DNA replication](#). Tritiated thymidine has also been used to study the effectiveness of therapeutic drugs that either stimulate or block DNA synthesis, such as chemotherapeutic agents. Tritium-labeled [amino acids](#) are incorporated into proteins and are used to follow protein biosynthesis.

The presence of tritium may be determined by liquid scintillation counting or chromatography. In [autoradiography](#), specimens containing tritiated compounds may be coated with photographic emulsion and incubated under refrigeration to expose the emulsion to tritium beta particles. The emulsion is then developed and fixed, and the tissues are stained for microscopy to identify the labeled organelles (4) (see [Autoradiography](#)).

#### Bibliography

1. D. R. Lide and H. Pr. Frederikse, eds. (1995) *CRC Handbook of Chemistry and Physics*, CRC Press, Boca Raton, Fla.
2. Knolls Atomic Power Laboratory (1966) *Chart of the Nuclides*, 15th editor, available from General Electric Company, San Jose, Calif.
3. F. A. Odekwa, D. D. Kozoll, and K. A. Meyer (1963) *J. Nucl. Med.* **4**, 60.
4. L. G. Caro and R. P. Van Tubergen (1962) *J. Cell. Biol.* **15**, 173–188.

#### Suggestions for Further Reading

5. *The Radiochemical Manual*, 2nd ed. (1966) Amersham, Bucks, England.
6. National Council on Radiation Protection and Measurement (1979) *Tritium and Other Radionuclide Labeled Organic Compounds Incorporated in Genetic Material*, NCRP Report No. 63, Washington D.C.,

## Hydrogenase

Hydrogenases are [enzymes](#) that catalyze the reversible oxidation of molecular hydrogen and play an important role in the anaerobic metabolism of microorganisms. Hydrogen uptake and production occur in sulfate-reducing **bacteria** and fermentative **anaerobes**, as well as in aerobic **nitrogen-fixing** bacteria, photosynthetic bacteria, and **cyanobacteria** (1). Three distinct categories of

hydrogenases have been described:

- *Iron hydrogenases*. With two [4Fe-4 S] centers in addition to a special type of [4Fe-4 S] center.
- *Nickel–iron hydrogenases*. With a nickel center coordinated to a sulfur atom, two [4Fe-4 S] centers, and a [3Fe-xS] center.
- *Nickel–iron-selenium hydrogenases*. With a nickel center coordinated to a selenium atom as well as two [4Fe-4 S] centers.

## Bibliography

1. D. S. Patil (1994) *Methods. Enzymol.* **243**, 68–94.

## Hydrolase

Hydrolases are [enzymes](#) that catalyze the hydrolysis of a range of bonds. They frequently show broad specificity, and this characteristic has posed problems in naming them. The systematic name for an enzyme that performs a hydrolytic reaction includes the term hydrolase, whereas the trivial name is formed by the addition of the suffix *ase* to the name of the substrate or the class of substrate on which the enzyme acts. Those hydrolases that act on carboxylic esters, thioesters, phosphomonoesters, phosphodiester, phosphotriester, diphosphomonoesters, **DNA** and **RNA** have been classified as a group that acts on ester bonds. The glycosidases hydrolyze *O*-glycosyl, *N*-glycosyl, and *S*-glycosyl compounds; etherases act on ether and thioether bonds; peptidases hydrolyze aminoacylpeptide, dipeptide, and dipeptidylpeptide bonds. Other categories of hydrolase are the serine [carboxypeptidases](#) and [serine proteinases](#), which contain a [DIFP](#) (diisopropyl fluorophosphate)-sensitive a [serine](#) residue at their active sites, the [thiol proteinases](#), the [carboxyl proteinases](#), and the [metalloproteinases](#). There are also hydrolases that act on C—N, rather than peptide, bonds, and the substrates for these enzymes include amides, cyclic amides, linear amidines, cyclic amidines, and nitriles. Other hydrolases act on acid anhydrides and C—C, halide, P—N, S—N, and C—P bonds.

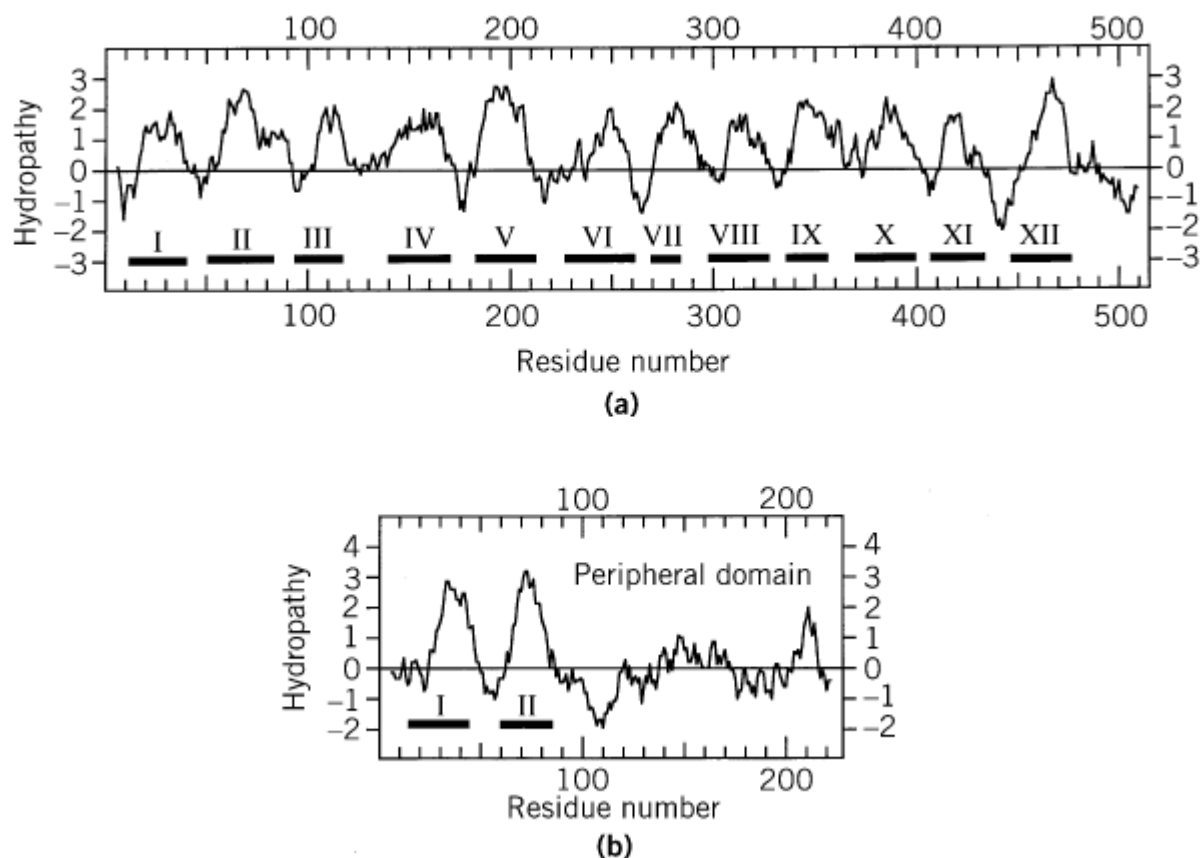
## Hydropathy

[\$\alpha\$ -helices](#) can be identified in the [primary structure](#) of a [membrane protein](#) with a relatively high confidence using hydropathy analysis, in which each residue is assigned a characteristic [hydrophobicity](#) value. Hydropathy is a term that refers to the spectrum of amino acid [side chains](#) in the [hydrophilicity](#)–hydrophobicity scale. It was originally coined by Kyte and Doolittle(1) in an article describing an algorithm that identifies helical transmembrane spans and that came to be widely used. Their scale is based on the [transfer free energies](#) from water to vapor of model compounds for the amino acid side chains, as well as on the exterior–interior distribution of amino acids in [protein structures](#). The hydropathy values of some residues were adjusted subjectively by the authors. Another frequently used hydrophobicity scale has been devised by Engelman, Steitz, and Goldman(2). This is a physicochemical scale based on the free energies of transfer from water to oil of amino acid residues that are parts of a helix. These free energy changes have been corrected for

unfavorable hydrophilic contributions due to polar atoms and protonation/deprotonation events of charged residues upon transfer to the [nonpolar](#) phase. A number of other scales exist. A useful review on hydrophobicity scales has been written by Eisenberg(3), and references to recent literature are provided by Wimley and White(4), who have measured the free energies of transfer of all 20 amino acid residues from the interface of a lipid bilayer to water.

To calculate a hydropathy profile, the hydrophobicity of a chosen number of residues within a window is summed and divided by the number of residues in the window. The resulting score is assigned to the residue in the middle of the window, after which the window is moved on by one residue and the procedure is repeated. An example of a hydropathy prediction is shown in Figure 1. In the cases of [bacteriorhodopsin](#), cytochrome oxidase, and bacterial [photosynthetic reaction center](#), this simple analysis correctly identified the approximate positions of all the transmembrane helices, with no false positives.

**Figure 1.** A Kyte–Doolittle hydropathy analysis of subunits I (a) and II (b) of the bovine cytochrome *c* oxidase, with the true locations of the transmembrane helices marked beneath the plot. A window length of 11 residues was used. The C-terminal peripheral domain of subunit II is indicated. Helix VII of subunit I is least well predicted. It contains only 17 residues, whereas hydropathy analysis suggested a length of about 26 that partially overlaps the residues seen in the actual helix in the X-ray structure.



Several more advanced methods for prediction of transmembrane helices have been developed. In addition to the hydropathy analysis, these make use of, for example, the observation that deletions and insertions occur in **homologous** sequences only outside the transmembrane sequences. Thus, use of a multiple sequence alignment, instead of a single sequence, as the input for the algorithm improves the prediction. Likewise, the modern algorithms make use of statistical analyses of amino

acid preferences within transmembrane helices and in their flanking loops.

The “positive inside” rule has been incorporated into the prediction protocols to determine the orientation of the polypeptide in the membrane. This rule is based on the fact that **arginine** and **lysine** residues occur more frequently in the loops between transmembrane helices that reside in the inner side of the membrane (and are shorter than 60 residues) than in the loops on the outer surface of the membrane.

### Bibliography

1. J. Kyte and R. F. Doolittle (1982) A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **157**, 105–132.
2. D. M. Engelman, T. A. Steitz, and A. Goldman (1986) Identifying nonpolar transbilayer helices in amino acid sequences of membrane proteins. *Annu. Rev. Biophys. Biophys. Chem.* **15**, 321–353.
3. D. Eisenberg (1984) Three-dimensional structure of membrane and surface proteins. *Annu. Rev. Biochem.* **53**, 595–623.
4. W. C. Wimley and S. H. White (1996) Experimentally determined hydrophobicity scale for proteins at membrane interphases. *Nature Struct. Biol.* **3**, 842–848.

## Hydrophilic

Most biological molecules contain electronegative atoms such as O, N, or S. These molecules or groups are **polar** and can act as **hydrogen bond** acceptors or donors, allowing them to interact more or less strongly with other such polar groups, with groups possessing charges of opposite sign, or with solvent **water**, which is equipped with partial positive and negative charges. Molecules or groups that bear full electrostatic charges, such as the  $\text{NH}_3^+$  and  $\text{COO}^-$  groups of amino acids in neutral solution, are extremely polar and especially difficult to remove from water. **Hydrophilicity**, or hydrophilic character, provides a quantitative index of a molecule's polarity, by indicating the difficulty of removing it from solvent water.

## Hydrophilicity

When biological molecules combine or react with each other, or change shape, their interactions with the solvent change. The noncovalent binding of **ligands** by **enzymes**, **antibodies**, and other **receptors**, eg, involve removal of the interacting partners from the solvent to which they were previously exposed; also, the ease of breaking and making chemical bonds depends on the strength of interaction of the reactants and products with solvent. Most biological molecules contain electronegative atoms such as O, N, or S, which can act as **hydrogen bond** acceptors or donors and allow them to interact more or less strongly with solvent **water**; they are said to be **hydrophilic**. The rates and equilibria of most biological processes, including binding interactions and metabolic transformations, are therefore strongly influenced by changing interactions with the surrounding solvent, which typically is ~80% water in intracellular fluid. Many biological compounds are so

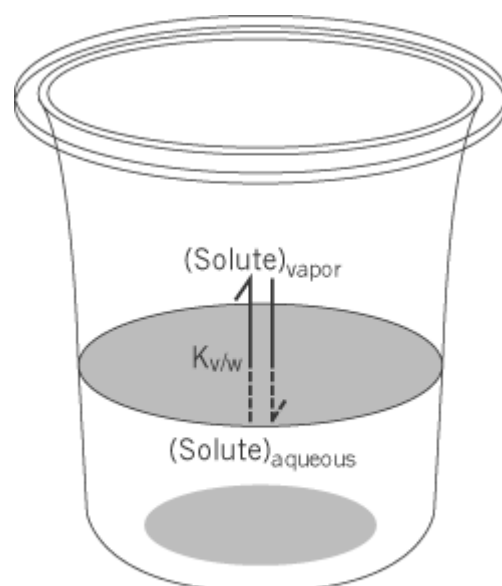


[polar](#) that their removal from solvent water, an operation essential to measuring the strength of their solvation by water, or *hydrophilicity*, is difficult. The very difficulty of measuring hydrophilicity suggests the magnitude of the solvation forces and the potential magnitude of their influence on rates and equilibria.

Hydrophilicity can be defined as the dimensionless equilibrium constant for the transfer of a molecule between the dilute vapor phase, in which intermolecular contacts are negligible, to dilute aqueous solution, in which each solute molecule is completely surrounded by solvent water. The concentration is measured in the same units, such as moles/liter, in each phase. Hydrophilic character is simply  $\log_{10}$  of the equilibrium constant for transfer from vapor to water (1). Hydrophilicity and hydrophilic character values provide the simplest measure of a molecule's affinity for watery surroundings, since they do not depend on any reference phase in which the molecule interacts with other partners. In that sense, they differ from [hydrophobicity](#), which involves a second solvent that is immiscible with water, and with which the solute molecule interacts more or less strongly (see [Partition Coefficient](#)). Hydrophilic character also differs from solubility, in which the reference phase is the pure solute in liquid or crystalline form: Such an arbitrary reference is of questionable value if different solute molecules are to be compared against each other.

Hydrophilic character can be determined in some cases by measuring the water solubility of a gas under an atmosphere of known partial pressure, or conversely, by measuring the concentration of a solute in the gas space over an aqueous solution of known concentration (Fig. 1). In more typical cases, of the kind often encountered in biological molecules, the concentration of gas in the vapor phase may be too low to detect directly. In such cases, a large volume of carrier gas (usually argon or nitrogen) is bubbled through an aqueous solution of known concentration, then through an efficient trap, which may be water itself. From the amount of solute collecting in the trap, plotted as a function of the measured volume of carrier gas that has passed through the system, it is a simple matter to calculate the concentration of solute in the vapor phase. Using [radioisotope](#) labels, [UV spectroscopy](#), or high-field [NMR](#) measurements, it is possible in this way to examine molecules as hydrophilic as **peptides** (2) or guanidino compounds (3), and to determine water-to-vapor distribution coefficients as low as  $10^{-10}$ . A scale of such equilibrium constants for simple organic compounds *in uncharged form* is shown in Figure 2. More extreme values are beyond direct measurement. However, there is much evidence to suggest that the free energy of solvation by water tends to be an additive function of a molecule's constituent groups (1), allowing values for complex molecules to be estimated by considering their substituents. From the behavior of methylamine and acetic acid, eg, it can be inferred that the concentration of glycine in the vapor phase, over a 1- *M* aqueous solution, is approximately  $10^{-14}M$ , even though that value is too low to measure directly.

**Figure 1.** The equilibrium for transfer of a molecule from dilute aqueous solution to the vapor phase. Hydrophilicity is the reciprocal of this equilibrium constant, and it measures the affinity of a molecule for watery surroundings.



$$K_{v/w} = \frac{(\text{solute})_{\text{vapor}}}{(\text{solute})_{\text{aqueous}}}$$

**Figure 2.** Scale of hydrophilicities for simple organic compounds in uncharged form.

$$K_{eq} = \frac{\text{Moles/liter (vapor)}}{\text{Moles/liter (aqueous)}}$$

|                    |                  |  |                                      |
|--------------------|------------------|--|--------------------------------------|
|                    | 10 <sup>2</sup>  |  |                                      |
| Alkanes            |                  | CH <sub>3</sub> —CH <sub>3</sub>   |                                      |
| Alkenes            | 10               | CH <sub>2</sub> =CH <sub>2</sub>   |                                      |
|                    |                  |  |                                      |
| Alkynes            | 1                | CH≡CH  |                                      |
|                    |                  |  |                                      |
| Thiols             |                  | C <sub>2</sub> H <sub>5</sub> —SH  |                                      |
| Chlorides          | 10 <sup>-1</sup> | C <sub>2</sub> H <sub>5</sub> —Cl  |                                      |
| Thioethers, ethers |                  | CH <sub>3</sub> —S—CH <sub>3</sub>   |                                      |
|                    |                  | CH <sub>3</sub> —O—CH <sub>3</sub>   |                                      |
|                    |                  |  |                                      |
| Esters             | 10 <sup>-2</sup> | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \\ \diagdown \quad \diagup \\ \quad \quad \text{O} - \text{CH}_3 \end{array}$                |                                      |
| Aldehydes          |                  |  |                                      |
| Ketones            |                  | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \begin{array}{l} \text{(H)} \\ \diagdown \end{array} \\ \quad \quad \text{CH}_3 \end{array}$ |                                      |
| Nitriles           | 10 <sup>-3</sup> | CH <sub>3</sub> —CN  |                                      |
|                    |                  |  |                                      |
| Amines             |                  | C <sub>2</sub> H <sub>5</sub> —NH <sub>2</sub>   |                                      |
| Alcohols           |                  | C <sub>2</sub> H <sub>5</sub> —OH  |                                      |
|                    | 10 <sup>-4</sup> |  |                                      |
| Water              |                  | H <sub>2</sub> O   |                                      |
| Acids              | 10 <sup>-5</sup> | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \\ \diagdown \quad \diagup \\ \quad \quad \text{OH} \end{array}$                             |                                      |
|                    |                  |  |                                      |
| Phosphotriesters   | 10 <sup>-6</sup> | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \\ \diagdown \quad \diagup \\ \quad \quad \text{N}(\text{CH}_3)_2 \end{array}$               | (CH <sub>3</sub> O) <sub>3</sub> P=O |
|                    |                  |  |                                      |
| Amides             |                  | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \\ \diagdown \quad \diagup \\ \quad \quad \text{NH}_2 \end{array}$                           | CH <sub>2</sub> OH                   |
| Diols              | 10 <sup>-7</sup> |  |                                      |
| Peptides           |                  | $\begin{array}{c} \text{O} \\ \parallel \\ \text{CH}_3 - \text{C} \\ \diagdown \quad \diagup \\ \quad \quad \text{NH}(\text{CH}_3) \end{array}$                | CH <sub>2</sub> OH                   |

These methods have revealed that the hydrophilicities of the protein amino acid side-chains, expressed in terms of their equilibria of transfer from the vapor phase to neutral aqueous solution at 25°C, span a range of approximately 10<sup>16</sup>-fold(3, 5). These values are as follows:

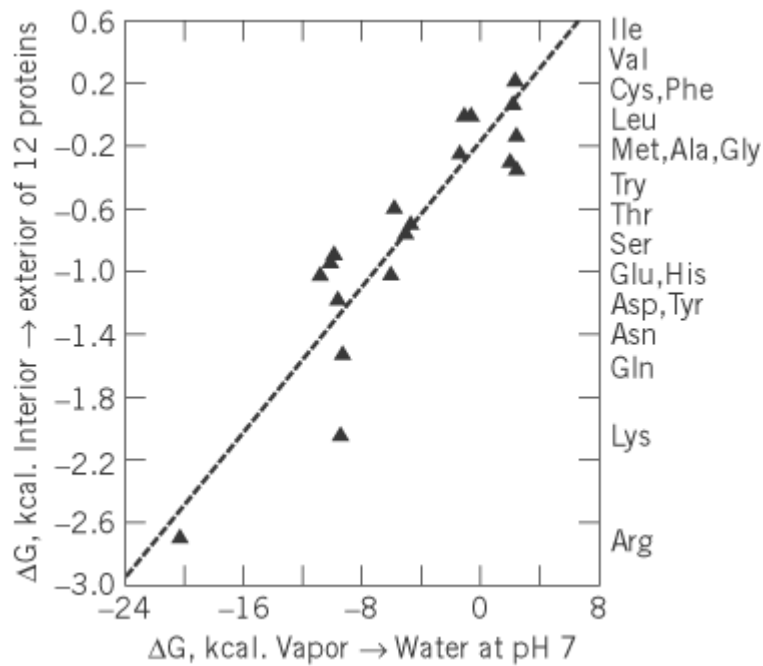
---

| $\log_{10}K_{v \rightarrow w, \text{pH7}}$ |       |
|--|-------|
| Glycine                                    | -1.75 |
| Leucine                                    | -1.68 |
| Isoleucine                                 | -1.58 |
| Valine                                     | -1.46 |
| Alanine                                    | -1.43 |
| Proline                                    | -1.10 |
| Methionine                                 | -1.09 |
| Phenylalanine                              | -0.56 |
| Cysteine                                   | 0.91  |
| Threonine                                  | 3.59  |
| Serine                                     | 3.72  |
| Tryptophan                                 | 4.32  |
| Tyrosine                                   | 4.49  |
| Glutamine                                  | 6.90  |
| Lysine                                     | 7.00  |
| Asparagine                                 | 7.12  |
| Histidine                                  | 7.55  |
| Glutamate                                  | 7.50  |
| Aspartate                                  | 8.05  |
| Arginine                                   | 14.65 |

---

These values were found to be closely related to the tendencies of the corresponding amino acids found at the surfaces of proteins, as revealed by their structures determined by [X-ray crystallography](#). Figure 3 shows the relationship between hydrophilicity and the nominal [transfer free energy](#) of each side-chain from the interior to the surface of a protein, as indicated by the relative populations of surface and buried residues for each amino acid. In addition, a remarkable bias in the [genetic code](#) was noted, such that purines (usually A) serve as the second code letter for all the more hydrophilic amino acids (3, 4). Figure 4 shows the amino acid side-chains in order of increasing hydrophilicity, alongside the code letters at the second position for each of the amino acids, but the significance of this correlation is not understood.

**Figure 3.** Tendencies of the 20 amino-acid residues to appear at the surfaces of globular proteins, compared with their free energies of transfer from the vapor phase to neutral aqueous solution (3, 4).



**Figure 4.** Code letters and anticode letters at the second position for each of the amino acids.

|                  | <u>2nd codeletter (mRNA)</u> | <u>2nd anticodeletter (DNA)</u> |
|------------------|------------------------------|---------------------------------|
| Gly              | G                            | C                               |
| Leu              | U                            | A                               |
| Ile              | U                            | A                               |
| Val              | U                            | A                               |
| Ala              | C                            | G                               |
| Phe              | U                            | A                               |
| Cys              | G                            | C                               |
| Met              | U                            | A                               |
| Thr              | C                            | G                               |
| Ser              | C(G)                         | G(C)                            |
| Trp              | G                            | C                               |
| Tyr              | A                            | T                               |
| Gln              | A                            | T                               |
| Lys              | A                            | T                               |
| Asn              | A                            | T                               |
| Glu              | A                            | T                               |
| His              | A                            | T                               |
| Asp              | A                            | T                               |
| Most hydrophilic | Arg                          | G                               |

The close relationship observed between hydrophilicity and [protein structure](#) (Fig. 3) suggested that this scale might have some value in predicting **protein folding**. It was also used to propose a scale of [hydropathy](#) values for predicting the relative tendencies of different parts of a protein sequence to be found in [membranes](#) (6). More recently, this scale of hydrophilicities has been

compared and contrasted with scales of hydrophobicity, obtained by reference to condensed phases including 1- *n*-octanol and to cyclohexane (7). Completion of a three-sided cycle allows estimation of equilibria of transfer between the vapor phase and cyclohexane, a pure measure of **van der Waals** forces. As expected, these are found to be roughly proportional to the [accessible surface](#) area. Since the interiors of protein molecules represent condensed phases, scales of hydrophobicity can be considered more directly related to folding equilibria than are scales of hydrophilicity using the vapor phase as a reference. Hydrophilicity, on the other hand, is more amenable to theoretical study because it does not entail the involvement of a second condensed phase in addition to water.

### Bibliography

1. J. Hine and P. K. Mookerjee (1975) *J. Org. Chem.* **40**, 292–298.
2. R. Wolfenden (1978) *Biochemistry* **17**, 201–204.
3. R. Wolfenden, L. Andersson, P. M. Cullis, and C. C. B. Southgate (1981) *Biochemistry* **20**, 849–855.
4. R. Wolfenden, P. M. Cullis, and C. C. B. Southgate (1979) *Science* **206**, 575–577.
5. P. R. Gibbs, A. Radzicka, and R. Wolfenden (1991) *J. Amer. Chem. Soc.* **113**, 4714–4715.
6. J. Kyte and R. F. Doolittle (1982) *J. Mol. Biol.* **157**, 105–132.
7. A. Radzicka and R. Wolfenden (1988) *Biochemistry* **27**, 1664–1670.

### Suggestion for Further Reading

8. R. Wolfenden (1983) Waterlogged molecules. *Science* **222**, 1087–1093.

## Hydrophobic Chromatography

### 1. The Basic Principle

Hydrophobic chromatography (HC) is a [chromatographic](#) method for the characterization, separation, and isolation of proteins and cells (1). It is based on **hydrophobic** (water-repelling) interactions between (a) hydrocarbon chains of different lengths and shapes that are anchored on inert beads and (b) the hydrophobic “pockets” or “patches” that are found on the [accessible surface](#) of [proteins](#). Such hydrocarbon chains interact with the protein molecules submerged in [water](#) and bind them to minimize the surface of the hydrophobic area that is exposed to the water. This binding brings about the release of “ordered” water molecules, which forms the basis for the hydrophobic interactions. It turns out that proteins can be distinguished and resolved from one another by HC, because they differ in the size, number, and water repellence of their hydrophobic “patches” or “pockets.” This basic principle in chromatography (2, 3) is also known as hydrophobic interaction chromatography and **reversed-phase** high-performance liquid chromatography ( [HPLC](#)). It has now become a very powerful method of separation and has been used for research in biochemistry, molecular biology, and biotechnology—for example, in the preparation and characterization of drugs such as [enzymes](#), [hormones](#), [antibodies](#), and receptors.

### 2. The Development of HC

HC was developed as a result of an unexpected observation made while attempting to make use of [affinity chromatography](#) (AC) (4). The basic concept in AC is based on the use of the specific affinity between the biorecognition site of a protein (an enzyme, an antibody, a receptor, etc.) and

one of its biospecific ligands (a substrate, an inhibitor, a hormone, or a regulatory metabolite) that is covalently linked to an inert matrix backbone and can be regarded as an immobilized “bait” for the protein. The column with the immobilized bait is then used to “fish out” the desired protein from a mixture. While developing the basic principles of AC, it was observed that the purification of a desired protein is often improved by interposing a hydrocarbon chain (an “arm”) between the ligand and the matrix backbone (4, 5). It was presumed that such an arm relieves the steric restrictions imposed by the backbone on the ligand, thereby increasing its flexibility as well as its availability to the protein (5, 6). It was also implied that the hydrocarbon arms do not alter the inert nature of the matrix, a condition that obviously has to be ensured to preserve a strictly biospecific (active-site mediated) adsorption of the extracted protein. This assumption seemed reasonable at the time (at least with water-soluble proteins), because it had been shown in the case of [lysozyme](#) that “most of the markedly nonpolar and hydrophobic side chains  $\frac{1}{4}$  are shielded from the surrounding liquid by more polar parts of the molecule” and that, as predicted by Sir Eric Rideal and Irving Langmuir, “lysozyme is quite well described as an oil drop with a polar coat” (7).

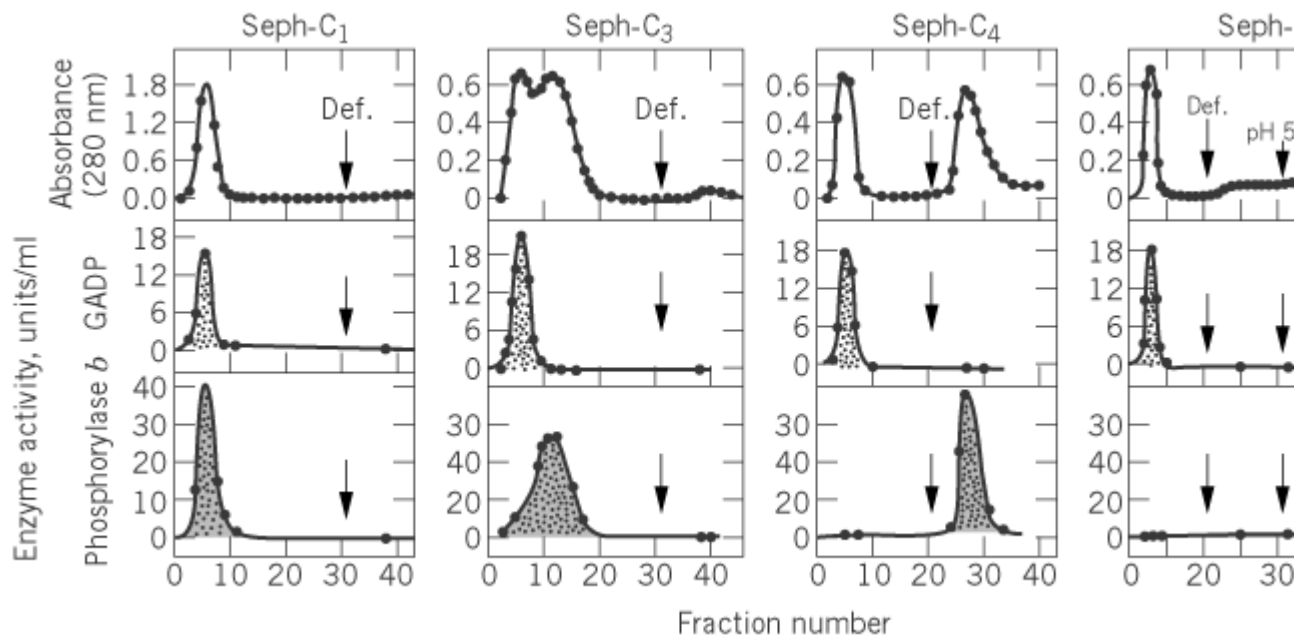
For the purification of enzymes involved in glycogen metabolism, glycogen (the “bait”) was anchored onto beaded [agarose](#). Unexpectedly, two glycogen-coated agarose preparations that differed only in the length of the hydrocarbon chains bridging the ligand and the agarose bead (Seph-C<sub>8</sub>-glycogen and Seph-C<sub>4</sub>-glycogen) exhibited different behaviors: Whereas the column material with eight-carbon-atom bridges adsorbed [glycogen phosphorylase b](#), that with four-carbon-atom bridges did not even retard it (2). Control experiments revealed that both the adsorption of the enzyme by Seph-C<sub>8</sub>-glycogen and its exclusion by Seph-C<sub>4</sub>-glycogen could be reproduced with columns that contained no glycogen whatsoever (Seph-C<sub>8</sub>-NH<sub>2</sub> and Seph-C<sub>4</sub>-NH<sub>2</sub>), indicating that free w-aminoethyl chains (which had not reacted with glycogen) were mainly responsible for the retention of phosphorylase on Seph-C<sub>8</sub>-glycogen and also raising the possibility that the very length of the hydrocarbon chain determines the retention power of these columns.

### 3. Homologous Series of Alkylagarose Columns

To challenge this working hypothesis, a homologous series of alkylagarose columns was synthesized, and their retention power was tested with two enzymes, glycogen phosphorylase *b* and glyceraldehyde-3-phosphate dehydrogenase (Fig. 1). This “homologous series of alkylagarose columns” (Seph-C<sub>*n*</sub>) refers to a set of columns that are identical in their ligand density, net charge, and ultrastructure. Such a set of columns can be considered as a series of consecutive controls in which each column differs from the preceding one in the series only in having hydrocarbon chains one carbon atom longer. In general, under the same experimental conditions, different proteins display different adsorption profiles on the Seph-C<sub>*n*</sub> columns (1-3). For example, under the conditions used to retain phosphorylase *b* ( $n \geq 4$ ), D-glyceraldehyde-3-phosphate dehydrogenase was not even retarded on any of the columns tested ( $n = 14$ ). Thus, while Seph-C<sub>1</sub> and Seph-C<sub>2</sub> would not distinguish between the two enzymes and would exclude both of them, Seph-C<sub>3</sub> would begin to resolve them, and Seph-C<sub>4</sub> would separate them efficiently, indicating that the ability to bind and to discriminate between two proteins (ie, to resolve and separate them) is a function of the number of carbon atoms in the hydrocarbon chains of the columns (Fig. 1). In addition, the above results showed that it is possible to adjust the tightness of adsorption of a given protein and avoid an overly strong retention, which would then require drastic conditions for elution that might **denature** the eluted protein. Finally, these studies illustrated the use of “deforming agents” (local, fully reversible **denaturants** that loosen the structure of the proteins) as efficient and delicate eluents in detaching adsorbed proteins from HC columns.

**Figure 1.** Preferential adsorption of glycogen phosphorylase *b* on alkylagarose columns varying in the length of their alk

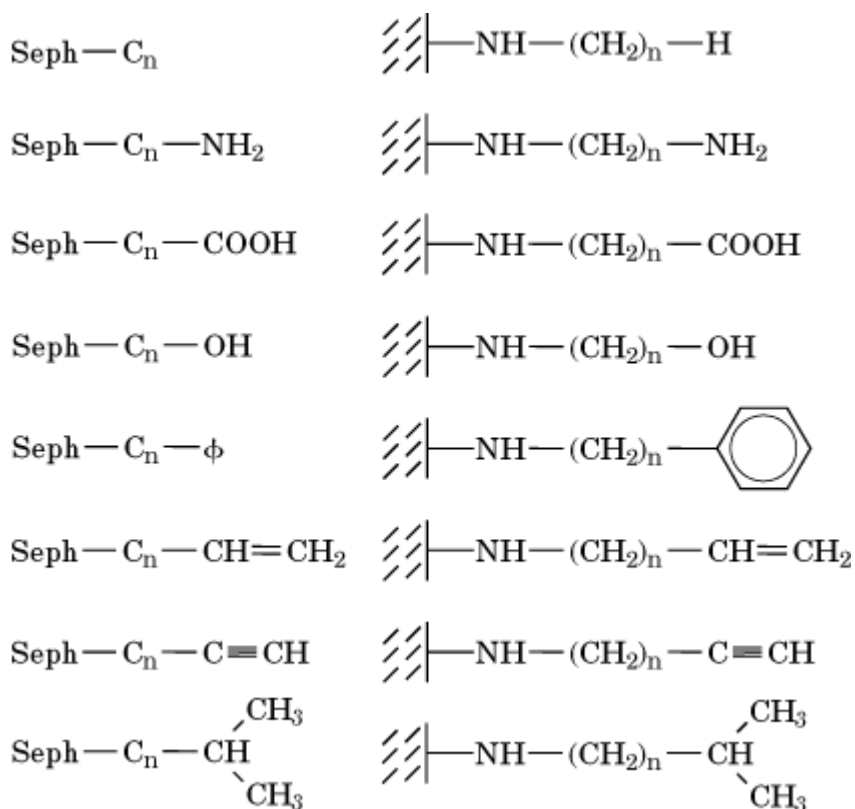
chains. Samples containing D-glyceradehyde-3-phosphate dehydrogenase (GAPD) and glycogen phosphorylase *b* were a onto the columns. Nonadsorbed proteins were washed off, then (*arrow*) a deforming buffer (Def.) was applied, which el phosphorylase from Seph-C<sub>4</sub> but not from Seph-C<sub>6</sub>. Acetic acid (0.2 M) was needed to detach glycogen phosphorylase *b* Seph-C<sub>6</sub>, which in this case came out in an inactive (denatured) form.



HC provides now a general, systematic approach to the purification of water-soluble, as well as lipophilic, proteins. This approach makes use of homologous series of alkylagaroses and their derivatives (Seph-C<sub>n</sub>-X where X = H, NH<sub>2</sub>, COOH, OH, C<sub>6</sub>H<sub>5</sub>, etc.; Fig. 2) to achieve resolution and purification of proteins and cells. Each member in these series offers flexible hydrophobic arms or “yardsticks” that interact with accessible hydrophobic “patches” or “pockets” in the various proteins, retaining only some proteins out of a mixture. Further resolution is then achieved by gradually changing the nature of the eluent. The wide applicability of this type of chromatography in protein purification and in separating cells is now well accepted. Procedures for the synthesis and the characterization of Seph-C<sub>n</sub> ( $n = 112$ ) are described in Ref. 1.

**Figure 2.** Examples of homologous series of alkylagarose derivatives (Seph-C<sub>n</sub>-X) that can be used for HC (1).

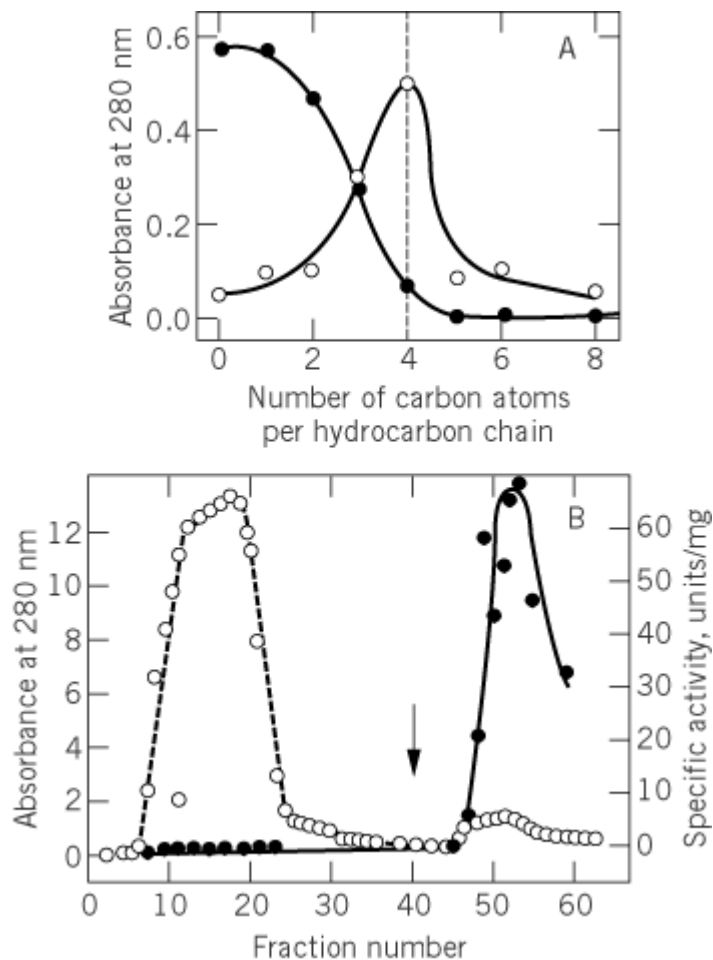




#### 4. Selecting a Column for a Given Purification—the Exploratory Kit

The selection of a suitable substituted agarose within a given homologous series (eg, Seph-C<sub>n</sub>) is achieved by means of an exploratory group of column materials designed for this purpose (1). Each column in the kit contains one member of the series. If, under certain conditions, a specific protein is adsorbed onto hydrocarbon chains X carbon atoms long, and if a change in the eluent promotes exclusion from Seph-C<sub>x</sub>, then this is the column of choice. An example of an adsorption-elution profile of glycogen phosphorylase *b*, on the Seph-C<sub>n</sub> series, is given in Figure 3a. As seen in this figure, as the length of the hydrocarbon chains increases, less protein is excluded from the column. While Seph-C<sub>0</sub> and Seph-C<sub>1</sub> exclude all the protein, Seph-C<sub>2</sub> and Seph-C<sub>3</sub> retard it, whereas Seph-C<sub>4</sub> and higher members of the series retain the protein. By passage of a 0.4 M imidazolium-citrate buffer, a “deformer” for glycogen phosphorylase *b* (8, 9), it is possible to elute the enzyme quantitatively from the Seph-C<sub>4</sub> column, but not from higher members in the series. In fact, recovery of the enzyme from Seph-C<sub>6</sub> was achieved only with a drastic eluent, 0.2 M CH<sub>3</sub>COOH, which displaces the enzyme but also denatures it. On the basis of such results, Seph-C<sub>4</sub> should be chosen for the preparative purification of this enzyme; and indeed this column allowed a 60- to 100-fold purification in one step, with >95% recovery of enzyme activity (Fig. 3b).

**Figure 3.** Hydrophobic chromatography (a) Adsorption profile of glycogen phosphorylase *b* on a Seph-C<sub>n</sub> kit. The filled circles indicate the amount of protein passing straight through the column. Empty circles indicate the amount of protein eluted by a “deforming buffer.” (b) Preparative purification of the enzyme from a muscle extract. Absorbance (empty circles) and specific activity (filled circles) were monitored. Nonadsorbed protein was removed by washing, and elution with the “deforming buffer” was initiated at the fraction indicated by the arrow (1, 2).



## 5. Means of Elution

Studies aimed at the optimization of elution from alkylagaroses have shown that proteins can be desorbed from such columns by a variety of means: polarity-reducing agents, specific deformers, mild [detergents](#), low concentration of denaturants, alteration in pH or temperature, and changes in ionic strength and ionic composition. Because the availability of hydrophobic crevices or patches on the surface of a protein appears to depend on its conformation, and because the retention of proteins by alkylagaroses depends largely on the lipophilicity, size, shape, and number of these crevices and patches, the above-mentioned means of elution may function either by directly disrupting the hydrophobic interactions between the column material and the protein or by changing the conformation of the protein. In fact, both mechanisms may often operate simultaneously. High selectivity in elution may also be achieved by using biospecific ligands, such as coenzymes, substrates, specific metal ions, and [allosteric](#) effectors, which often bring about ligand-induced conformational changes in proteins.

## 6. Useful Features of Hydrophobic Chromatography

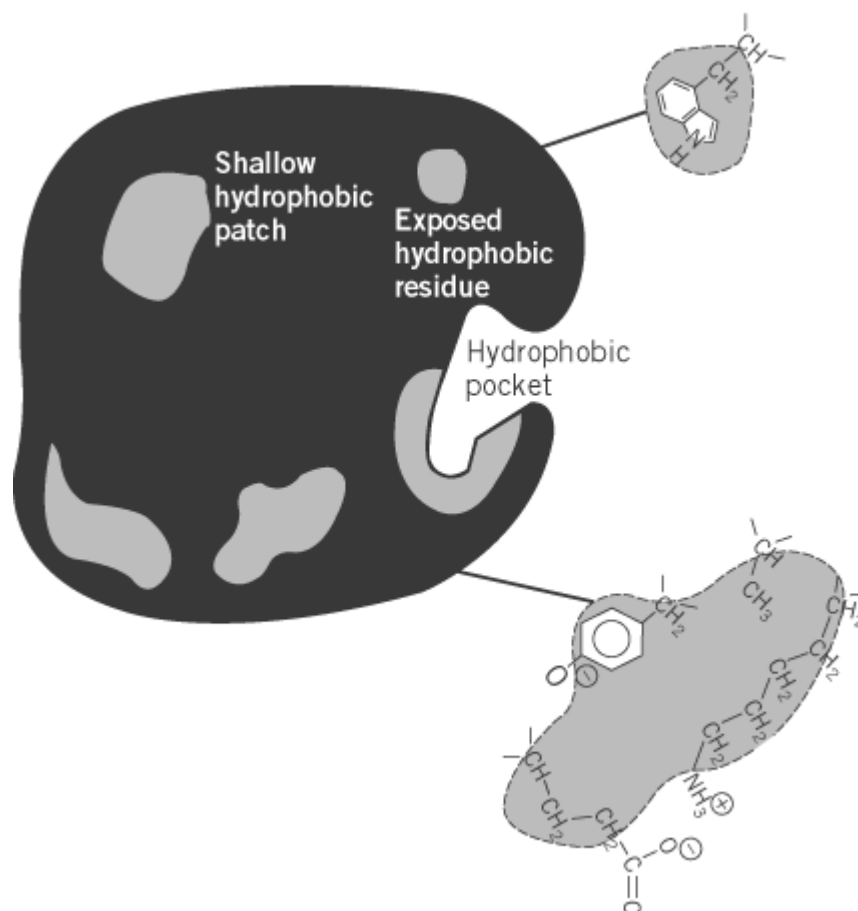
The important features of hydrophobic chromatography are as follows: (i) It provides an independent criterion for the resolution of macromolecules. (ii) It provides means for gradually adjusting the adsorption forces and thereby avoiding an overly strong retention. (iii) It can be used for the detection of conformational changes in proteins when these are reflected in the size or distribution of hydrophobic crevices or patches. A potential use of these column materials may be in the purification and study of lipophilic membrane-bound proteins that have accessible hydrophobic regions used for their localization within the membrane. (iv) It provides a basis for the development

of automatic consecutive fractionators (1) for the purification of proteins.

## 7. On the Mechanism of Action of Alkylagaroses

In their functional conformation, water-soluble proteins are folded so as to “bury” as many as possible of their hydrophobic side chains in the interior of the molecule and to expose as many as possible of their polar, charged side chains to interaction with water. It is now clear, however, that complete burying of all hydrophobic groups is generally not achieved, leaving some hydrophobic groups exposed at the surface of the protein. Together with hydrophobic components of charged amino acids, such as the (-CH<sub>2</sub>-), stretches of [lysine](#) and [glutamic acid](#) residues, or the phenyl ring of tyrosine residues, these form hydrophobic “patches” or “pockets” at the surface of the molecule (Fig. 4). A sufficiently large hydrophobic patch may constitute a binding site for the hydrocarbon chains implanted on the hydrophilic agarose matrix and form “hydrophobic bonds,” freeing “ordered” water molecules and allowing them to interact with each other (10-12). The available hydrophobic patches and pockets of different proteins vary in number, size, shape, and lipophilicity, and these variations are reflected in the relative affinities of different proteins for a specific alkylagarose. It is the properties of such patches, and perhaps their distribution on the surface of different proteins, that play a major role in the resolution of proteins on alkylagaroses.

**Figure 4.** Schematic representation of the surface of a protein molecule with an exposed hydrophobic amino acid residue (tryptophan) and hydrophobic “pockets” or “patches” (*hatched*). The scheme also illustrates how hydrophobic constituents of several nonhydrophobic amino acids bearing hydrophobic functional groups can together form a site capable of accommodating a hydrophobic hydrocarbon chain (1).



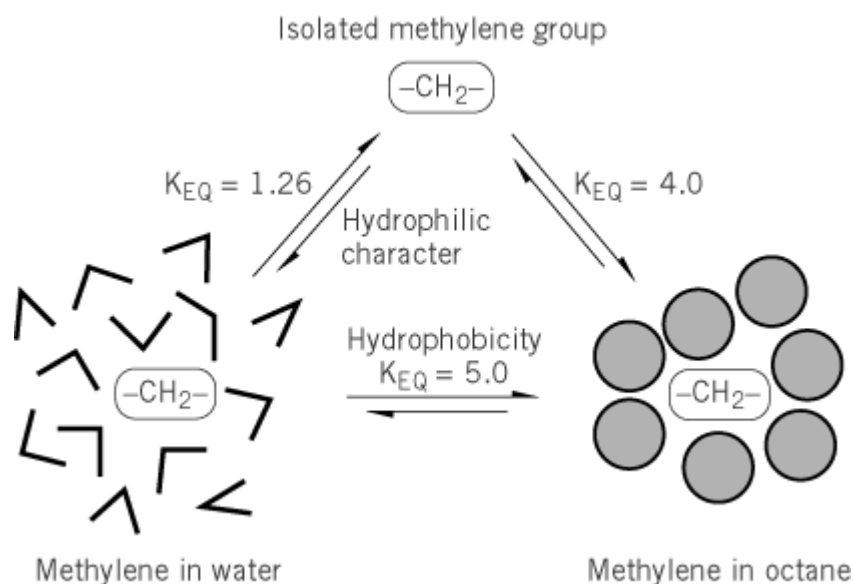
## Bibliography

1. S. Shaltiel (1984) *Methods Enzymol.* **104**, 69.
2. Z. Er-el, Y. Zaidenzaig, and S. Shaltiel (1972) *Biochem. Biophys. Res. Commun.* **49**, 383.
3. S. Shaltiel and Z. Er-el (1973) *Proc. Natl. Acad. Sci. USA* **70**, 778.
4. P. Cuatrecasas, M. Wilchek, and C. B. Anfinsen (1968) *Proc. Natl. Acad. Sci. USA* **61**, 636.
5. P. Cuatrecasas (1970) *J. Biol. Chem.* **245**, 3059.
6. P. Cuatrecasas and C. B. Anfinsen (1971) *Annu. Rev. Biochem.* **40**, 259.
7. D. C. Phillips (1966) *Sci. Am.* **November**, 78.
8. S. Shaltiel, J. L. Hedrick, and E. H. Fischer (1966) *Biochemistry* **5**, 2108.
9. J. L. Hedrick, S. Shaltiel, and E. H. Fischer (1969) *Biochemistry* **8**, 2422.
10. G. S. Hartley (1936) *Aqueous Solutions of Paraffin-Chain Salts* Hermann, Paris.
11. H. S. Frank and M. W. Evans (1945) *J. Chem. Phys.* **13**, 507.
12. W. Kauzmann (1959) *Adv. Protein Chem.* **14**, 1.

## Hydrophobic Effect

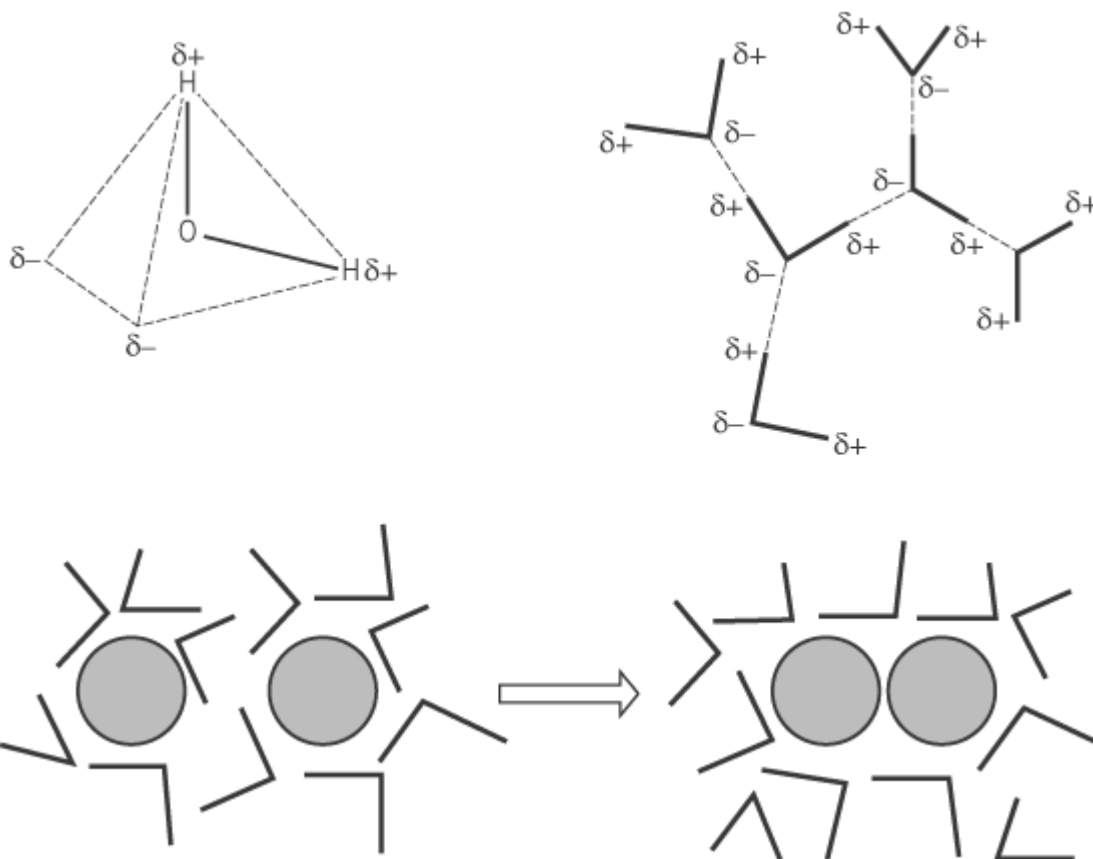
Cooks and physicists are aware of the tendency of oil to separate from [water](#), a solvent that is very self-cohesive. In biology, this effect is expressed in the tendency of uncharged (or [nonpolar](#)) molecules to escape from water by entering less polar surroundings, or by adhering to each other (Fig. 1). This phenomenon, termed the *hydrophobic effect* (1, 2), is believed to play a decisive role in maintaining the stability of biological [membranes](#), in the proper folding of protein molecules (see [Protein Stability](#)), and in determining the relative affinities of [hormones](#), [antibodies](#), and substrates for proteins that bind them. The tendency of any particular nonpolar molecule or chemical group at equilibrium to favor transfer from water to a nonpolar phase, such as a hydrocarbon solvent, is termed its [hydrophobicity](#).

**Figure 1.** A partial view of the origins of the hydrophobic effect.



The physical origins of the hydrophobic effect remain controversial, because the properties of water as a solvent are not yet fully understood. It is of interest to ask whether nonpolar molecules, such as methane or ethane with no [polar](#) groups, tend to leave water and enter less polar solvents (octane is a familiar example) mainly because they are repelled by water; or whether they do so because they are attracted to the less polar solvent (or to each other). That question can be addressed by referring to an absolute standard, the vapor phase, that neither attracts nor repels solutes, and in which “solute” molecules exist in isolation most of the time. Using the vapor phase as a reference, single molecules of methane are found to exhibit an appreciable tendency to leave water, as indicated by their equilibria of transfer from water to the vapor phase. For methane, this equilibrium constant is 27 in favor of the vapor phase at 25°C. This tendency increases gradually with the size of a normal alkane, as can be seen in the series  $\text{CH}_4$  27,  $\text{C}_2\text{H}_6$  20,  $\text{C}_3\text{H}_8$  29,  $\text{C}_4\text{H}_{10}$  38,  $\text{C}_5\text{H}_{12}$  51,  $\text{C}_6\text{H}_{14}$  74,  $\text{C}_7\text{H}_{16}$  83. For each  $\text{CH}_2$  increment, transfer to the vapor phase is enhanced by an average factor of 1.26. In contrast, each  $\text{CH}_2$  increment enhances the equilibrium of transfer from the vapor phase to a nonpolar solvent by a factor of 4.0. As a result, each  $\text{CH}_2$  increment increases the distribution coefficient for transfer from water to a nonpolar solvent by a factor of 5.0 (Fig. 2). Similar behavior has been observed in the normal series of hydrocarbons, esters of acetic acid, amines, and alcohols (3). In summary, single molecules of hydrocarbons have an appreciable tendency to leave water and enter the vapor phase, where they exist in isolation. These molecules are thus truly “hydrophobic” in the sense that they have an absolute tendency to escape from water. That is hardly surprising, since water is among the most self-cohesive molecules known, with an extremely high surface tension. However, this tendency to leave water is only slightly enhanced by increasing the size of the hydrocarbon (addition of a  $-\text{CH}_2$  increment), whereas increasing size greatly enhances the tendency of a nonpolar solute to leave water and enter a nonpolar solvent. One is led to conclude that attraction to the nonpolar solvent also plays a major role in the hydrophobic effect, and this is believed to be due to [van der Waals interactions](#), or London dispersion forces.

**Figure 2.** Analysis of the hydrophobic effect in terms of the equilibrium affinity of a solute for water and for a nonpolar liquid. (a) Electrons are more strongly attracted by the oxygen atom of the water molecule, leaving partial negative charges near the corners of a tetrahedron, and partial positive charges at the other two corners. (b) Water molecules are strongly attracted to each other by their partial charges, forming H-bonds between H and O atoms. As a result, liquid water is very self-cohesive. (c) Nonpolar molecules, like oil drops, break up the structure of water in which they are dissolved. The self-cohesiveness of water encourages these molecules or groups to coalesce, so that the part of the solvent can return to its preferred structure. This is called the hydrophilic effect.



The hydrophobic effect is unique, among the noncovalent forces involved in maintaining biological structures, in that it grows stronger with increasing temperature (1, 4), ie, **enthalpy** is released ( $DH$  is positive) when the force takes effect. If there is a net free energy of attraction (ie,  $DG$  is negative), then **entropy** must be released ( $TDS$  must be positive and greater in magnitude than  $DH$ ), since  $DG = DH - TDS$ . If a protein's stability were dependent entirely on the hydrophobic effect, then the protein would be expected to lose its structure at low temperatures, and some proteins do in fact undergo "cold denaturation" (see **Protein denaturation**). There has been continuing interest in the probable structural origins of the entropy (and enthalpy) increases that accompany the formation of hydrophobic interactions. Changes in the properties of water in the immediate neighborhood of the solute almost certainly account, at least in part, for the loss of entropy that accompanies the introduction of a nonpolar molecule or methylene increment into water from the vapor phase or a nonpolar solvent. Frank and Evans (5) suggested that this might imply formation around solutes of a kind of *clathrate* or *iceberg* structure, in which water molecules were more ordered than in the bulk solute. These authors did not suggest that water surrounding a solute resembled Ice I in any literal sense, nor does current experimental evidence seem to support such a possibility. Entropy losses associated with introducing nonpolar solutes into water might also arise, at least in part, from restrictions on the mobility of solutes when they are introduced into a structured aqueous environment. In normal aliphatic compounds of increasing size, solubility might (according to this view) be reduced by progressive restrictions on internal rotation, and compounds with internal rotations already restricted would not be affected to the same extent. In apparent accord with this view, steroids and cycloalkanes display lower activity coefficients in water than nonrigid compounds of similar size, with reference to both nonpolar solvents and to the vapor phase (6), and the alicyclic amino acid proline is more hydrophilic than would be expected for an amino acid with an aliphatic side-chain of the same size (7). These effects could alternatively be due, at least in part, to the reduced surface area of cyclic as compared with open-chain compounds, resulting in the imposition of less severe restrictions on the solvent. These uncertainties illustrate the general difficulty of deconvoluting thermodynamic functions, which describe equations of state, without making

assumptions that are themselves open to question. A clearer view of the origins of the structural origins of the hydrophobic effect awaits a more complete understanding of the properties of water and aqueous solutions.

All molecules, even those that are truly nonpolar, are attracted to each other by London dispersion forces. These forces are responsible for the fact that octane, and the [lipid](#) molecules present in biological membranes, are liquids at room temperature. They also make a major contribution to the hydrophobic effect discussed above, and to the resulting tendencies of protein molecules to adopt the specific structures needed for their activities. Anyone who has observed the wheels of a locomotive, or two cyclists riding a tandem bicycle, has a rough mental image of the London dispersion force. If we picture two atoms of argon near each other, it is easy to imagine that temporary fluctuations in the distribution of electrons around one of these atoms may, at any given instant, place a partial negative charge toward the other atom. In response to this fluctuation, electrons in the orbitals of the neighboring atom tend to adjust slightly, placing a partial positive charge in the vicinity of the first atom and creating a small electrostatic attraction between them. A moment later, the roles may be reversed, but a net attraction remains. Because of the small charges involved, and their fluctuating character, London dispersion forces are individually weak and intensely dependent on the distances separating the atoms. These forces reach a maximum at an optimal distance equivalent to the sum of the **van der Waals** radii of the two atoms (typical van der Waals radii are 0.10 nm for H, 0.16 nm for C, 0.135 nm for O, 0.145 nm for N, and 0.17 nm for S). Beyond that distance, their mutual attraction falls off with the 6th power of the distance separating the two atoms. At distances closer than optimal, repulsion between the electron clouds of the two atoms builds up rapidly (unless they are able to share their electrons in a covalent bond), with the 12th power of the distance separating the atoms. The number and aggregate strength of the attractions between nonpolar molecules can be considerable, even though they are individually weak. They account for a substantial part of the hydrophobic effect discussed above and are thus largely responsible for maintaining the integrity of biological membranes and the proper folding of protein molecules.

### Bibliography

1. W. Kauzmann (1959) *Adv. Prot. Chem.* **14**, 1–45.
2. C. Tanford (1974) *The Hydrophobic Effect*, Wiley, New York, pp. 1–11.
3. R. Wolfenden and C. A. Lewis (1976) *J. Theor. Biol.* **59**, 231–235.
4. G. Némethy and H. A. Scheraga (1962) *J. Chem. Phys.* **36**, 3382–3400.
5. H. S. Frank and M. W. Evans (1945) *J. Chem. Phys.* **13**, 507–532.
6. M. Osinga (1979) *J. Amer. Chem. Soc.* **101**, 1621–1622.
7. P. R. Gibbs, A. Radzicka, and R. Wolfenden (1991) *J. Amer. Chem. Soc.* **113**, 4714–4715.

### Suggestions for Further Reading

8. J. T. Edsall and H. A. McKenzie (1978) Water and proteins. I. The significance and structure of water, its interaction with electrolytes and non-electrolytes. II. The location and dynamics of water in protein systems and its relation to their stability and properties. *Adv. Biophys.* **10**, 137–207.
9. W. P. Jencks (1969) *Catalysis in Chemistry and Enzymology*, McGraw-Hill, New York, pp. 393–436.
10. W. Kauzmann (1959) Some factors in the interpretation of protein denaturation. *Adv. Prot. Chem.* **14**, 1–45.
11. C. Tanford (1974) *The Hydrophobic Effect*, Wiley, New York, pp. 1–11.

## Hydrophobic Electrophoresis

Although aqueous buffers are the appropriate solvents for [gel electrophoresis](#) of native [hydrophilic](#) macromolecules, separations of **hydrophobic** species require the presence of [detergents](#) or miscible organic solvents. Uncharged or [amphoteric](#) detergents are compatible with the formation of a gel and [electrophoresis](#). Alternatively, charged detergents can be added to the buffer and have a negligible effect on electrophoresis, so long as conductance is not greatly elevated. The most widely used detergent for that purpose is the strongly dissociating and denaturing detergent, **SDS**. However, the **molecular sieving** of charged detergent micelles in electrophoresis as a function of gel concentration needs to be considered, as well as the generation of moving boundaries in discontinuous buffer systems, with the detergent as a leading or trailing ion (see [Disc Electrophoresis](#)). Miscible organic solvents are particularly suitable for relatively hydrophobic gels, such as some of the “Hydrolink” gels (1) or the new spongelike copolymer gel media(2).

When it is desirable to preserve native or relatively undissociated structures and their biological activities, gel electrophoresis must be conducted in “nondenaturing detergents,” ie, under conditions that render the particular detergent relatively innocuous with regard to the disruption and inactivation of active structures and complexes. For that purpose, seven groups of amphoteric or uncharged detergents are particularly useful(3). Their effectiveness in solubilization depends sharply on their concentration, [critical micelle concentration](#) (cmc) and micellar weights, and the protein/detergent ratio.

### Bibliography

1. R. J. Molinari, M. Connors, and R. G. L. Shorr (1993) *Adv. Electrophoresis* **6**, 44–64.
2. M. G. Harrington, K. H. Lee, E. Bailey, and L. E. Hood (1994) *Electrophoresis* **15**, 187–199.
3. L. M. Hjelmeland and A. Chrambach (1984) *Methods Enzymol.* **104**, 305–318.

### Suggestion for Further Reading

4. L. M. Hjelmeland and A. Chrambach (1981) *Electrophoresis and electrofocusing in detergent containing media: A discussion of basic concepts.* *Electrophoresis* **2**, 1–11.

## Hydrophobic-Interaction Chromatography

**Hydrophobic** (literally meaning “water hating”) describes [nonpolar](#) compounds that have little affinity for [polar](#) molecules, such as [water](#), and self-associate in a polar environment because of the [hydrophobic effect](#). Hydrophobic interactions are energetically favorable because of an increase in **entropy** ( $\Delta S > 0$ ) that originates when the solvent water molecules leave the more ordered structure around the nonassociated nonpolar compounds for the more unstructured, bulk polar environment. If matrices in a [chromatography](#) column provide hydrophobic surfaces to interact with analytes and one can separate mixtures of molecules with hydrophobic moieties, we are dealing with hydrophobic-interaction chromatography.

The principle of hydrophobic-interaction chromatography was originally known (1) by the name “**saltin-out** chromatography,” but it acquired its current name only after the synthesis of specific stationary phases in the early 1970s (2-4). In hydrophobic-interaction chromatography, retention is



promoted by high concentrations of an appropriate salt, and elution is achieved with low salt concentrations. The matrices in the column are very [hydrophilic](#) and have only a very mild hydrophobicity. The increase in concentration of a salting-out salt drives the analytes out of the mobile phase and onto the stationary phase, much like the [precipitation](#) of highly water-soluble proteins or hydrophilic polymers that occurs under conditions of high ionic strength (see [Salting In](#), [Salting Out](#) and [Preferential Hydration](#)).

Although many view hydrophobic-interaction chromatography as an extension of [reversed-phase chromatography](#) to an aqueous mobile phase without organic modifiers, there are a number of basic differences between these two separation techniques. In hydrophobic-interaction chromatography, the matrix is hydrophilic and is substituted with short-chain groups, such as phenyl, butyl, or octyl (see Ref. [5](#) for a list of different matrices), and the mobile phase is usually an aqueous salt solution. In reversed-phase chromatography, the matrix is an octyl (C8) or octadecyl (C18)-derived silica, and the mobile phase is usually a mixture of water and a less polar organic modifier (eg, methanol or acetonitrile). These two techniques exploit the different sources of the hydrophobicity of proteins. Hydrophobic-interaction chromatography depends on surface hydrophobic groups, which arise from the protein's **tertiary** and [quaternary structures](#), and is carried out under conditions that maintain the integrity of the protein conformation. Reversed-phase chromatography is carried out on the [unfolded protein](#), in which nearly all the hydrophobic groups are exposed to the matrix, and depends on the protein's [primary structure](#). Proteins frequently **denature** under conditions of reversed-phase chromatography because organic solvents are needed for elution and/or because of the strength of the interaction with the stationary phase, whereas the mild elution conditions and mild hydrophobic interaction of hydrophobic-interaction chromatography matrices allow the elution and recovery of native proteins.

The most interesting aspect of hydrophobic-interaction chromatography is the frequent increase of retention with increasing temperature because this chromatographic method is an entropy-driven process. This is the opposite of any other type of retention chromatography. It should also be pointed out that the proteins are loaded onto the hydrophobic-interaction chromatography column using a buffer with a high salt concentration as the starting mobile phase and are then eluted by a gradient with decreasing ionic strength. This is the opposite of the procedure employed in [ion-exchange chromatography](#).

Hydrophobic-interaction chromatography is extremely useful for purifying proteins and peptides. The principle of this chromatographic method is orthogonal to those of ion-exchange and [size-exclusion chromatography](#) methods. Therefore it is used effectively in schemes that combine all three chromatographic methods to separate proteins from complex mixtures. However, the solubilities of proteins in high-salt buffers are limited by the salting-out effect. This constrains the applications of hydrophobic-interaction chromatography for preparative purposes. A detailed operational protocol is beyond the scope of this article, so interested readers are directed to excellent reviews [[5](#), [6](#), and chapter 13 (pp. 250–259) of Ref. [7](#)].

## Bibliography

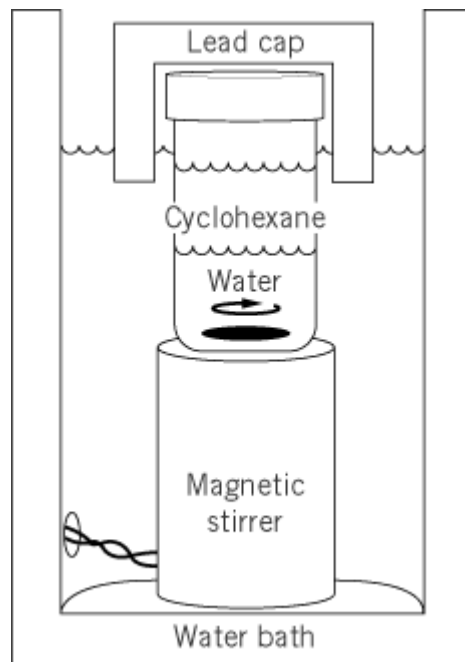
1. A. Tiselius (1948) *Ark. Kemi. Mineral. Geol.* **B26**, 1.
2. S. Shaltiel and Z. Er-el (1973) *Proc. Natl. Acad. Sci. USA* **70**, 778–781.
3. B. H. J. Hofstee (1973) *Anal. Biochem.* **52**, 430–448.
4. S. Hjertén, J. Rosengren, and S. Pahlman (1974) *J. Chromatogr.* **101**, 281–288.
5. K.-O. Eriksson (1989) in *Protein Purification: Principles, High Resolution Methods, and Applications* (J.-C. Janson, and L. Rydén, eds.), VCH, New York, pp. 207–226.
6. R. M. Kennedy (1990) in *Guide to Protein Purification* (M. P. Deutscher, ed), *Methods in Enzymology* **182**, Academic Press, New York, pp. 339–343.
7. U. D. Neue (1997) *HPLC Columns: Theory, Technology, and Practice*, Wiley-VCH, New York.

## Hydrophobicity

A molecule is said to be *hydrophobic* if it has a tendency to move from the [polar](#) solvent [water](#) to a [nonpolar](#), oily surrounding. The tendency of a molecule or chemical group to undergo transfer at equilibrium from water to a nonpolar phase, such as a hydrocarbon solvent, is termed its *hydrophobicity*. The inverse properties are [hydrophilic](#) and [hydrophilicity](#). Hydrophobicity is important in allowing drugs to cross membranes, and its analysis has therefore been important in the development of the pharmaceutical industry (1, 2). The relative hydrophobicities of the [amino acids](#) are of particular interest to biochemists and molecular biologists, because of their role in **protein folding**. Thus, the least hydrophobic amino acids are found almost exclusively at the [accessible surfaces](#) of proteins in which they occur, exposed to solvent water, whereas the interiors of proteins are almost uniformly composed of hydrophobic amino acids (3). For reasons that are not yet fully understood, but may be related to protein folding, the hydrophobicities of amino acids are also associated with a pronounced bias in the composition of the [genetic code](#) (4) (see also [Hydrophilicity](#)).

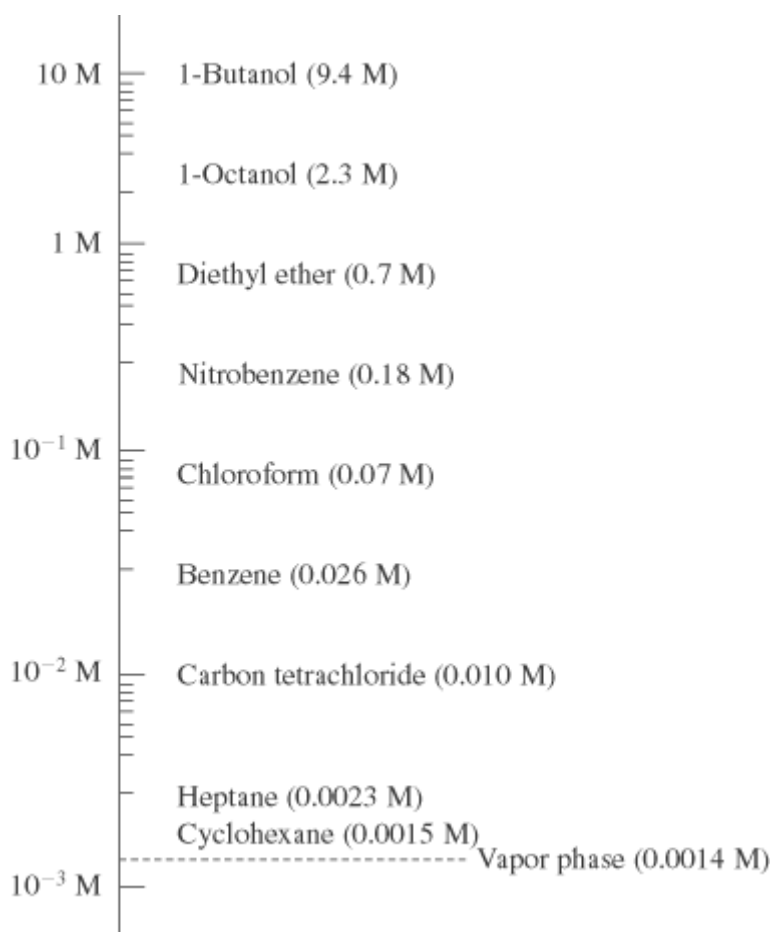
Hydrophobicity is usually expressed quantitatively as either the equilibrium constant or as the negative Gibbs **free energy** ( $-RT \ln K_{eq}$ ) for transfer of a solute from dilute aqueous solution to the nonpolar phase. This dimensionless equilibrium constant, the [partition coefficient](#), is determined experimentally by measuring the concentration of the solute that is dissolved in each phase after thorough mixing, at equilibrium. After mixing equilibrium is achieved, a separatory funnel or similar device is used to separate the phases before measuring their solute content. A simple device for performing these experiments rapidly, with small volumes of material, is illustrated in Figure 1. Some solutes show a pronounced tendency to self-associate, either in water or in the nonpolar solvent. Self-association is indicated by an increasing tendency to favor one phase, that in which aggregation occurs more readily, as the solute concentration is increased. When such behavior is encountered, partition coefficients must be extended to sufficient dilution to avoid self-association in both phases. The results can then be used as a basis for calculating the concentration-independent distribution coefficient between phases, the hydrophobicity.

**Figure 1.** A simple apparatus for measuring hydrophobicity, consisting of a submersible magnetic stirrer, 20-mL scintillation vial, and lead cap to hold the vial in place. Measurements are made using aqueous solutions buffered in a pH range in which the solute is uncharged, using the solute at high dilution to avoid self-association, and a sensitive method of detection (typically NMR or radioisotopic assay). To estimate effective  $K_{dist}$  at pH 7, the observed  $K_{dist}$  is corrected for the fraction of solute that is ionized.



Certain immiscible solvents, such as 1-*n*-butanol and 1-*n*-octanol, are much more polar than others and are miscible with water to a considerable extent. Figure 2 shows that when a distribution measurement has been performed using octanol, the nonpolar phase contains 2.3 *M* water; in contrast, water-saturated cyclohexane contains only  $1.6 \times 10^{-3}M$  water. Solutes experience relatively little change in surroundings in passing between wet octanol and water, so there is a “leveling” effect that tends to result in distribution coefficients that approach unity. To provide a more sensitive measure of differences in solute polarity, the nonpolar solvent that is used as a reference in these experiments should be as nonpolar as possible, ideally a liquid hydrocarbon such as cyclohexane. Adopting a scale of this kind, the side-chains of the amino acids that occur in proteins vary in their distribution coefficients between cyclohexane and neutral aqueous solution over a range of  $4 \times 10^{14}$ , from arginine (the most polar) to leucine and isoleucine (the least polar). For octanol-to-water distribution, the range is only  $3 \times 10^2$ . There is also evidence of specific interactions between solutes and octanol, as might be expected from the **hydrogen-bonding** potential of octanol's hydroxyl groups (2).

**Figure 2.** Solvents compared with respect to their water content at saturation, at 25°C.



For a highly [polar](#) compound, such as a carboxylic acid or an amine, the solute may partition so strongly in favor of water that the concentration in cyclohexane is difficult to detect. Sensitivity can be improved by using a very large volume (eg, 1000 mL) of the nonpolar phase to extract the solute from a small amount (1 mL) of water, separating the phases, then “back-extracting” the solute into a small volume (1 mL) of water. After appropriate correction, distribution coefficients as large as  $10^7$  can be measured in this way using [NMR](#), [UV spectroscopy](#) or [radioisotope](#) methods as a means of detection. [Ionization](#) tends to render a solute extremely polar, and its hydrophobicity is then more easily determined at pH values where the compound is uncharged. Its hydrophobicity at pH 7 can then be calculated by multiplying the value for the uncharged compound by the fraction that is actually uncharged at pH 7.

Many molecules are found to exhibit behavior consistent with the view that a molecule's hydrophobicity, expressed as free energy, is a roughly additive function of the hydrophobicities of its constituent chemical groups. A comprehensive collection of early data from the literature, with correlations and exceptions noted, has been published by Hansch and Leo (2). More recently, data have been collected for the amino acid side-chains (3), **nucleic acid** bases (5), phosphoric acid esters and amides (6), *cis*- and *trans*- [peptide bonds](#) (7) and carbohydrates (8). The hydrophobicities of the reactants and products of hydrolysis of ATP have been shown to differ to such an extent that ATP hydrolysis, which is highly exergonic in water (see [Adenylate Charge](#)), becomes endergonic in its absence, ie, changing the solvation can be considered to provide the entire driving force for ATP hydrolysis (8).

Hydrophobicities of the amino acid side chains occurring in proteins, expressed as  $\log_{10}$  of the equilibrium constant at 25°C, for transfer from neutral aqueous solution to cyclohexane (3), are

---

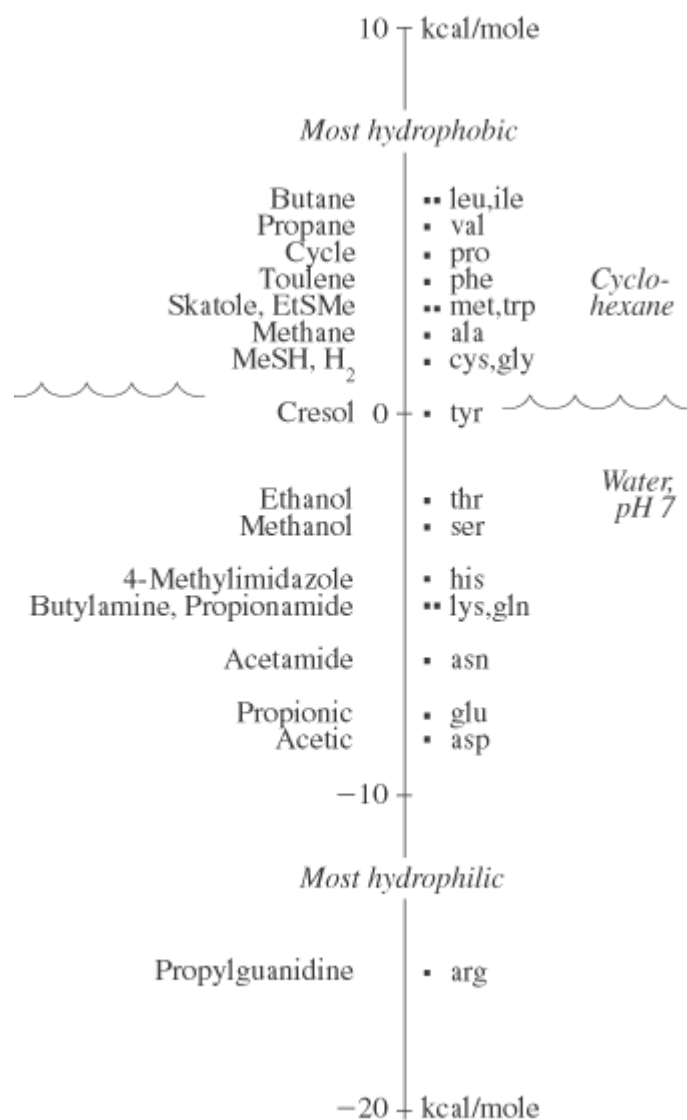
|               |        |
|---------------|--------|
| Leucine       | 3.62   |
| Isoleucine    | 3.62   |
| Valine        | 2.97   |
| Proline       | 2.51   |
| Phenylalanine | 2.19   |
| Methionine    | 1.73   |
| Tryptophan    | 1.71   |
| Alanine       | 1.33   |
| Cysteine      | 0.93   |
| Glycine       | 0.69   |
| Tyrosine      | -0.10  |
| Threonine     | -1.89  |
| Serine        | -2.49  |
| Histidine     | -3.42  |
| Glutamine     | -4.07  |
| Lysine        | -4.08  |
| Asparagine    | -4.88  |
| Glutamate     | -5.00  |
| Aspartate     | -6.41  |
| Arginine      | -10.97 |

---

Proline differs from the other amino acids in that its nitrogen atom is involved in a secondary amide linkage. From a family of related compounds, it can be inferred that proline is 2.9-fold less hydrophobic than valine (9), placing its value at 2.51 on the above scale.

It is of interest to compare the hydrophobicities of the amino acids that occur in proteins with their tendencies to appear in buried positions or at the surfaces of proteins. Such a comparison, shown in Figure 3, suggests that there is a close correlation. It should be noted that the slope of the line is very shallow, consistent with the view that the environment of buried amino acid residues is much less nonpolar than cyclohexane, especially in the case of the more polar amino acids.

**Figure 3.** Distributions of amino-acid side-chains, from surface to buried positions in proteins (10), treated as an equilibrium constant ( $K_{\text{exposed} \rightarrow \text{buried}}$ ), compared with distribution coefficients for the corresponding side-chains from neutral aqueous solution to cyclohexane (3).



## Bibliography

1. A. J. Leo, C. Hansch, and D. Elkins (1971) *Chem. Rev.* **71**, 525–680.
2. C. Hansch and A. J. Leo (1979) *Substituent Constants for Correlation Analysis in Chemistry and Biology*, Wiley, New York.
3. A. Radzicka and R. Wolfenden (1988) *Biochemistry* **27**, 1664–1670.
4. R. Wolfenden, L. Andersson, P. M. Cullis, and R. Wolfenden (1979) *Science* **206**, 575–577.
5. P. M. Cullis and R. Wolfenden, *Biochemistry* **20**, 3024–3028.
6. R. Wolfenden and R. Williams (1983) *J. Amer. Chem. Soc.* **105**, 1028–1031.
7. A. Radzicka, L. Pedersen, and R. Wolfenden (1987) *Biochemistry* **27**, 4358–4362.
8. R. Williams and R. Wolfenden (1985) *J. Amer. Chem. Soc.* **104**, 4345–4351.
9. P. R. Gibbs, A. Radzicka, and R. Wolfenden (1991) *J. Amer. Chem. Soc.* **113**, 4714–4715.
10. C. Chothia (1976) *J. Mol. Biol.* **105**, 1–18.

## Suggestion for Further Reading

11. R. Wolfenden (1983) Waterlogged molecules. *Science* **222**, 1087–1093.

## Hydroxyapatite Chromatography

Hydroxyapatite, also known as hydroxylapatite and abbreviated as HA, is a modified form of crystalline calcium phosphate,  $\text{Ca}_{10}(\text{PO}_4)_6(\text{OH})_2$ , that has proven to be a very useful inorganic matrix for purifying and fractionating diverse [macromolecules](#). Hydroxyapatite [chromatography](#) was introduced by Tiselius et al. in 1956 (1). Bernardi (2, 3) undertook extensive studies of **protein** behavior on hydroxyapatite, worked out standard procedures for the operation of HA columns, and discussed the mechanisms of adsorption and elution. Gorbunoff (4-6) extended Bernardi's work and deduced a mechanism of protein–HA interaction based on three fundamental conclusions:

1. Adsorption and elution cannot be regarded as simple reversals of a single process.
2. **Amino** and [carboxyl groups](#) act differently in the adsorption of proteins to HA.
3. Elution of basic and acidic proteins by different salts occurs by different mechanisms.

Briefly, in the adsorption of proteins to HA columns equilibrated with phosphate buffer, amino groups act through nonspecific [electrostatic interactions](#) between their positive charges and the general negative charges on the HA column, whereas carboxyl groups are repelled electrostatically from the negative charge of the column and by binding specifically to calcium sites on the column. Basic proteins are eluted as a result of normal Debye–Hückel charge screening by  $\text{F}^-$ ,  $\text{Cl}^-$ ,  $\text{ClO}_4^-$ ,  $\text{SCN}^-$  and phosphate anions, or by specific displacement by  $\text{Ca}^{2+}$  and  $\text{Mg}^{2+}$  ions. Acidic proteins are eluted by displacement of their carboxyl groups from HA calcium sites by ions, such as fluoride or phosphate, that form stronger complexes with calcium than carboxyl groups.

Hydroxyapatite chromatography provides a unique fractionation method that has properties distinct from other separation techniques, such as preparative [electrophoresis](#), [ion-exchange chromatography](#) and **size-exclusion chromatography**. Therefore HA chromatography has been used to advantage in purification schemes, frequently as a final purification step, and is used widely in biochemical research (see Ref. 7, a booklet provided by Calbiochem). The reader requiring more detailed and practical information is directed to the excellent review by Gorbunoff (8).

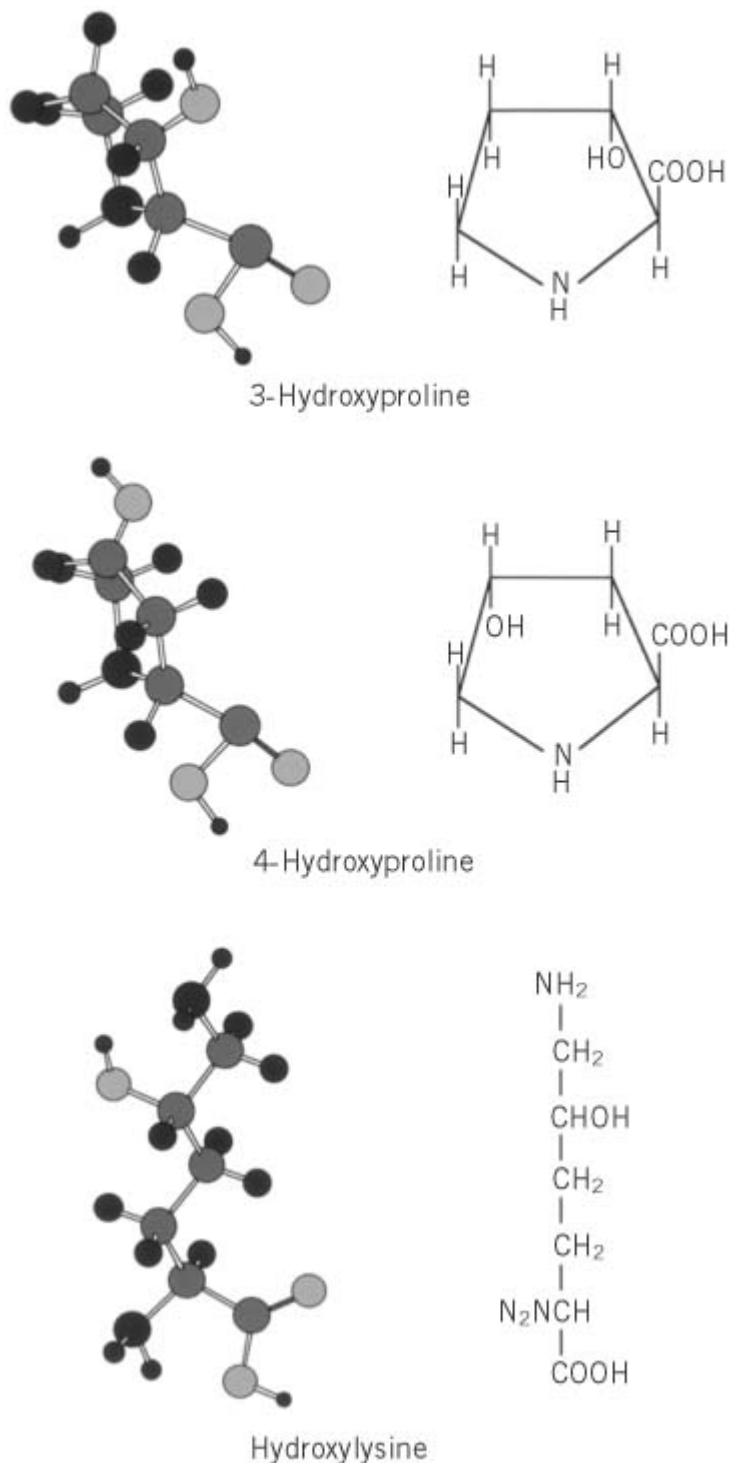
### Bibliography

1. A. Tiselius, S. Hjertén, and Ö. Levin (1956) *Arch. Biochem. Biophys.* **65**, 132–155.
2. G. Bernardi (1971) *Methods Enzymol.* **22**, 325–339.
3. G. Bernardi (1973) *Methods Enzymol.* **27**, 471–429.
4. M. J. Gorbunoff (1984) *Anal. Biochem.* **136**, 425–432.
5. M. J. Gorbunoff (1984) *Anal. Biochem.* **136**, 433–439.
6. M. J. Gorbunoff and S. N. Timasheff (1984) *Anal. Biochem.* **136**, 440–445.
7. T. L. Brooks (1985) *Hydroxylapatite: Fast Flow and High Resolution*, Calbiochem Biochemicals, San Diego.
8. M. J. Gorbunoff (1990) in *Guide to Protein Purification* (M. P. Deutscher, ed.), *Methods in Enzymology* **182**, Academic Press, New York, pp. 329–339.

## Hydroxylation (Lysine, Proline)

The principal [posttranslational modifications](#) of the polypeptides of [collagen](#) are the hydroxylation of proline and lysine residues to yield 4-hydroxyproline, 3-hydroxyproline (Hyp) and hydroxylysine (Hyl), plus glycosylation of the hydroxylysyl residues (Fig. 1). These modifications are catalyzed by three hydroxylases and two glycosyl transferases. The reactions occur until the polypeptides form the triple-helical collagen structure, which prevents any further modification reactions.

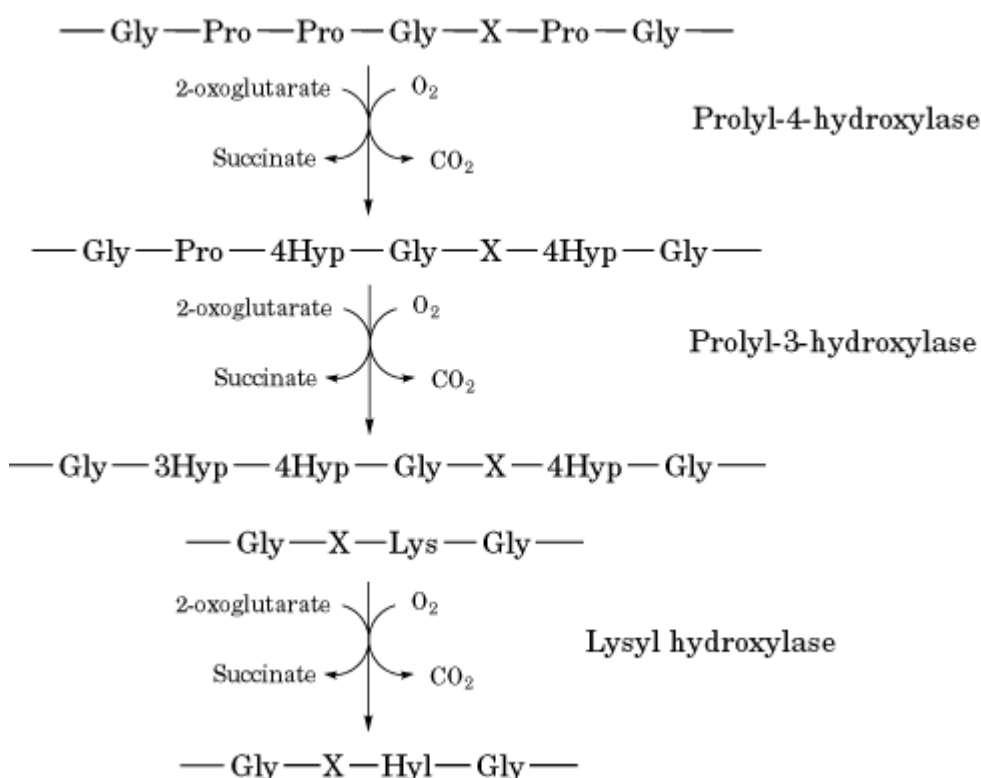
**Figure 1.** Structures of 4-hydroxyproline, 3-hydroxyproline, and hydroxylysine, as the free amino acids.





The hydroxylation of prolyl and lysyl residues in polypeptides is catalyzed by three separate glycoprotein members of the 2-oxoglutarate dioxygenase family of enzymes: prolyl 4-hydroxylase, prolyl 3-hydroxylase, and lysyl hydroxylase (Fig. 2). All three enzymes require a nonhelical polypeptide substrate. The minimum substrate structures for these enzymes are the triplets  $-X-Pro-Gly-$ ,  $-Pro-Hyp-Gly-$ , and  $-X-Lys-Gly-$ , respectively. The activities of both prolyl 4-hydroxylase and lysyl hydroxylase are influenced by the adjacent amino acids and by peptide chain length.

**Figure 2.** Hydroxylation reactions on proline and lysine residues of collagen.



The reaction mechanisms of all three hydroxylase enzymes are similar, and each requires as cosubstrates  $Fe^{2+}$ , 2-oxoglutarate, and molecular  $O_2$ . The 2-oxoglutarate is stoichiometrically decarboxylated during hydroxylation; one atom of the  $O_2$  molecule is incorporated into succinate and the other, into the hydroxyl group formed on the proline or lysine residue. The main function of ascorbate in the collagen hydroxylase reaction may well be to act as an antioxidant for the oxidation of  $Fe^{2+}$  to  $Fe^{3+}$ . Ascorbate is not consumed, and the hydroxylation reaction may proceed in its absence. The iron chelator,  $\alpha, \alpha'$ -dipyridyl, is a potent inhibitor of these hydroxylase reactions. In its presence, nonhydroxylated collagen is synthesized (in the case of type I collagen) or nonhelical chains may be secreted (in the case of type IV collagen).

### 1. Prolyl 4-Hydroxylase

The active prolyl 4-hydroxylase in vertebrates is a tetramer of molecular weight 240,000, and consists of two copies of each of two different types of subunits (a subunits, 64 kDa; b subunits, 60

kDa). It has two catalytic sites, one per pair of dissimilar subunits. The peptide- and 2-oxoglutarate-binding sites are located on the a subunit, while the ascorbate-binding sites may be composed of both a and b subunits. The principal contribution to the catalytic sites is from the a subunits, but some parts of the large catalytic sites may come from both the a and b subunits. The avian and mammalian a subunit sequences are highly conserved. The b subunit has proved to be a highly unusual multifunctional polypeptide that in isolation is [protein disulfide isomerase](#) (PDI).

## 2. Prolyl 3-Hydroxylase

Prolyl 3-hydroxylase (procollagen proline, 2-oxoglutarate 3-dioxygenase) requires a specific substrate sequence, catalyzing hydroxylation of prolyl residues only in the sequence –Pro–Hyp–Gly–. Prolyl 3-hydroxylase requires the same specific cofactors as does prolyl 4-hydroxylase. The reason for this modification is unknown, and the relative abundance of 3-hydroxyproline varies markedly in the different collagens. The individual chains of type I collagen each contain a single 3-hydroxyproline residue, while type IV collagen can contain up to 20 residues per 1000 amino acids.

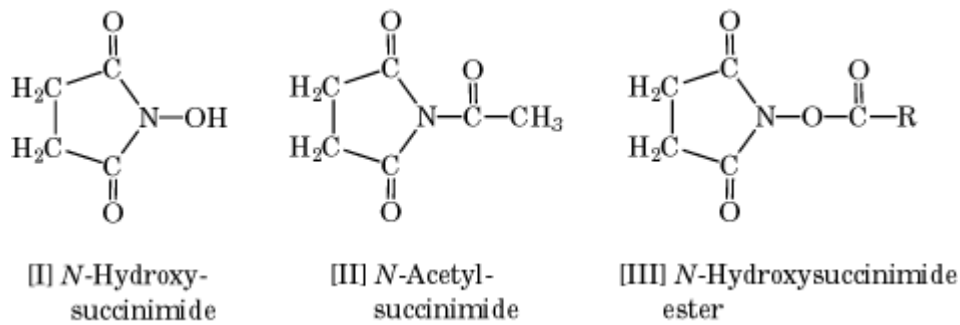
## 3. Lysyl Hydroxylase

The enzyme lysyl hydroxylase (procollagen-lysine, 2-oxoglutarate 5-dioxygenase) catalyzes the hydroxylation of lysyl residues at the d carbon atom (Fig. 1) in –X–Lys–Gly–triplet sequences by a mechanism similar to that described for prolyl 4-hydroxylase. The hydroxyl groups of the hydroxylysyl residues have two important functions—they act as attachment sites for carbohydrate residues (either galactose or glucosyl-galactose), and they play a role in the formation of intra- and intermolecular collagen crosslinks (see [Collagen](#)). The hydroxylysine content of the a chains of the different collagen types varies from 5 to 70 residues per 1000 amino acids.

The enzyme lysyl hydroxylase in vertebrates is a **glycoprotein** containing asparagine-linked carbohydrate moieties. The enzyme is a dimer of 190,000 Da, composed of two monomers with different carbohydrate contents. The polypeptide chain of chick lysyl hydroxylase consists of 710 amino acids, with a 20-residue signal sequence. It is not very similar in sequence to the homologous, catalytically important a subunit of prolyl 4-hydroxylase, even though the enzymes share the same cofactor and substrate specificities.

## ***N*-Hydroxysuccinimide**

*N*-Hydroxysuccinimide [I] is employed as an **acetylation** or acylation agent in its acetylated [II] or ester [III] forms. The activated ester is prepared by reacting *N*-hydroxysuccinimide with **carbodiimide**-activated [carboxyl groups](#). Acetylation of [amino groups](#) in [proteins](#) is conducted under milder conditions with *N*-acetylsuccinimide [II] or *N*-hydroxysuccinimide acetate than with acetic anhydride.



Various *N*-hydroxysuccinimide esters are prepared and employed to acylate proteins. These reactions are used for [chemical modification](#) and labeling of proteins (1). *N*-Hydroxysuccinimide esters are often employed as one part of the reactive groups in **bifunctional reagents** for [cross-linking](#) proteins (2). *N*-Hydroxysuccinimide esters are also employed to prepare activated column resins for [affinity chromatography](#). *N*-Hydroxysuccinimide esters of [amino acids](#) are obtained by condensing of *N*-protected amino acids and *N*-hydroxysuccinimide with carbodiimide, and they are widely employed for [peptide synthesis](#).

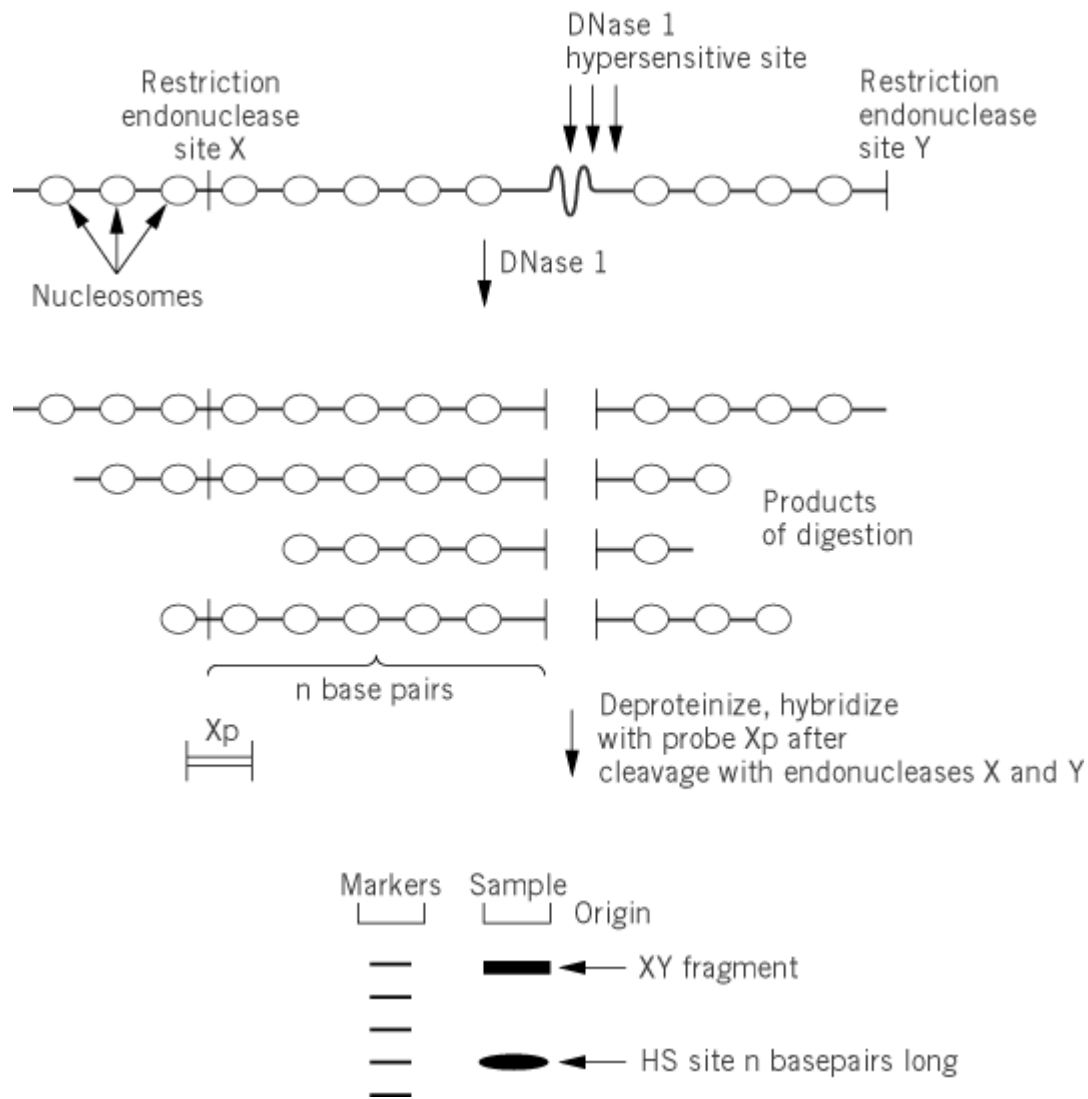
#### Bibliography

1. Y.W. Shing and A. Ruoho (1981) *Anal. Biochem.* **133**, 171–175.
2. T.H. Ji (1983) *Methods Enzymol.* **91**, 580–609.

#### Hypersensitive Site

**DNase 1** hypersensitive sites are the first place this **nuclease** introduces a double-strand break in [chromatin](#) (see [DNase 1 Sensitivity](#)). These sites usually involve small segments of DNA sequences (100 to 200 bp) and are two or more orders of magnitude more accessible to DNase 1 cleavage than inactive chromatin (1, 2). DNase 1 hypersensitive sites result from the assembly of specific [nucleoprotein](#) complexes that generally contain [transcription factors](#) flanked by positioned [nucleosomes](#). DNase I hypersensitive sites also reflect the stable association of a transcription factor on the surface of a nucleosome (3). DNase I hypersensitive sites and nucleosome arrays are usually detected by DNase 1 or **micrococcal nuclease** digestion, followed by indirect end-labeling methodologies to map the hypersensitive site relative to a site of [restriction enzyme](#) cleavage (Fig. 1). As the most accessible regions of chromatin to non-histone DNA-binding proteins, DNase I hypersensitive sites generally denote DNA sequences with important functions in the nucleus.

**Figure 1.** A scheme for mapping the position of DNase I hypersensitive sites in chromatin relative to restriction endonuclease cleavage sites.

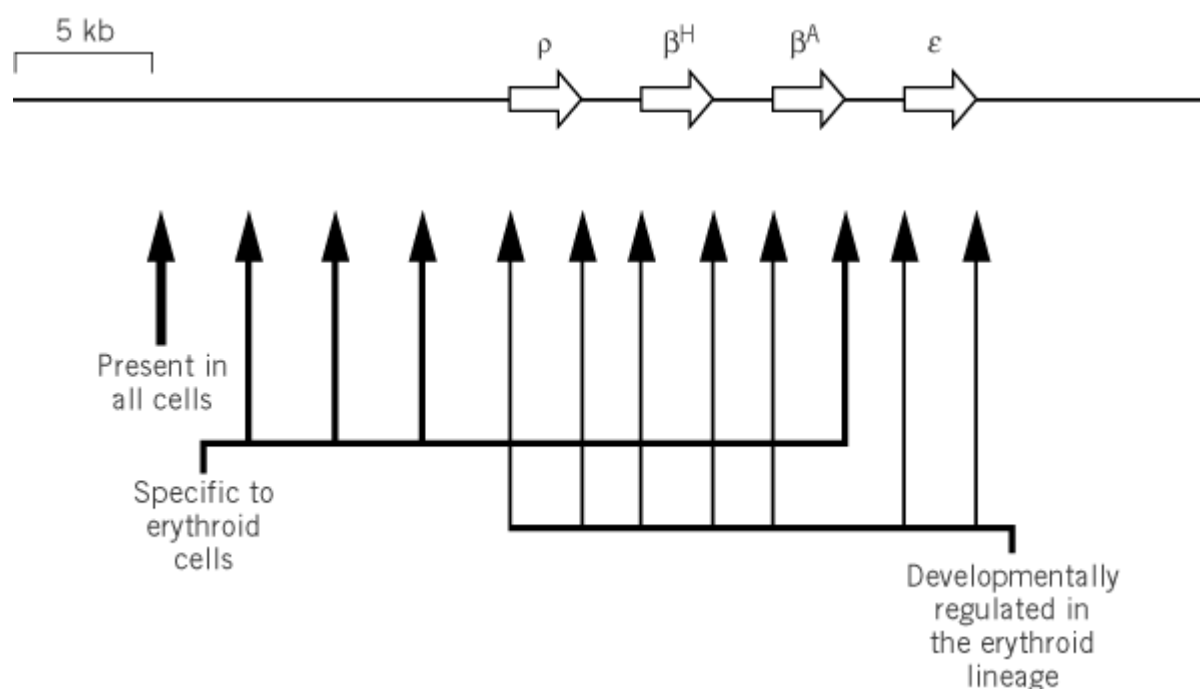


DNase I hypersensitive sites were first detected in the SV40 [minichromosome](#), at the region that functions as an **origin of replication**, and in *Drosophila* chromatin (4). In general, these sites are accessible to all [enzymes](#) or reagents that cut duplex DNA. These sites have been mapped to a large number of functional segments of DNA, including **promoters**, [enhancers](#), **locus control regions**, transcriptional **silencers**, origins of replication, [recombination](#) elements, and structural sites within or around [telomeres](#) (5).

These sites often fall into a hierarchy of patterns around regulated genes. Twelve DNase I hypersensitive sites are found in the chicken **b-globin** locus that contains four globin genes (5'-s-b<sup>H</sup>-b<sup>A</sup>-e-3') and covers more than 65 kbp of DNA (Fig. 2). One site is present in all cells, independent of whether the genes are transcriptionally active or not. Three sites, upstream of the s-globin gene, were present only in erythroid cells destined to express the globin genes. These sites were initially without clear functional significance. However, a similar site was found between the b<sup>A</sup> and e genes that corresponds to an enhancer element. Four sites were found over the promoters of each gene, depending on whether the gene was transcriptionally active, and three sites were found downstream of the genes, corresponding to transcription termination elements (the b<sup>A</sup> gene excluded). It is important to note that the formation of DNase I hypersensitive sites at the promoters of the globin genes is a relatively late step in the commitment of these genes to become transcriptionally active. However, it is clear that the formation of such sites precedes the actual initiation of transcription by

**RNA polymerase.** Indeed, the generation of these sites may account for a component of the general nuclease sensitivity of a gene (see [DNase I Sensitivity](#)) (6).

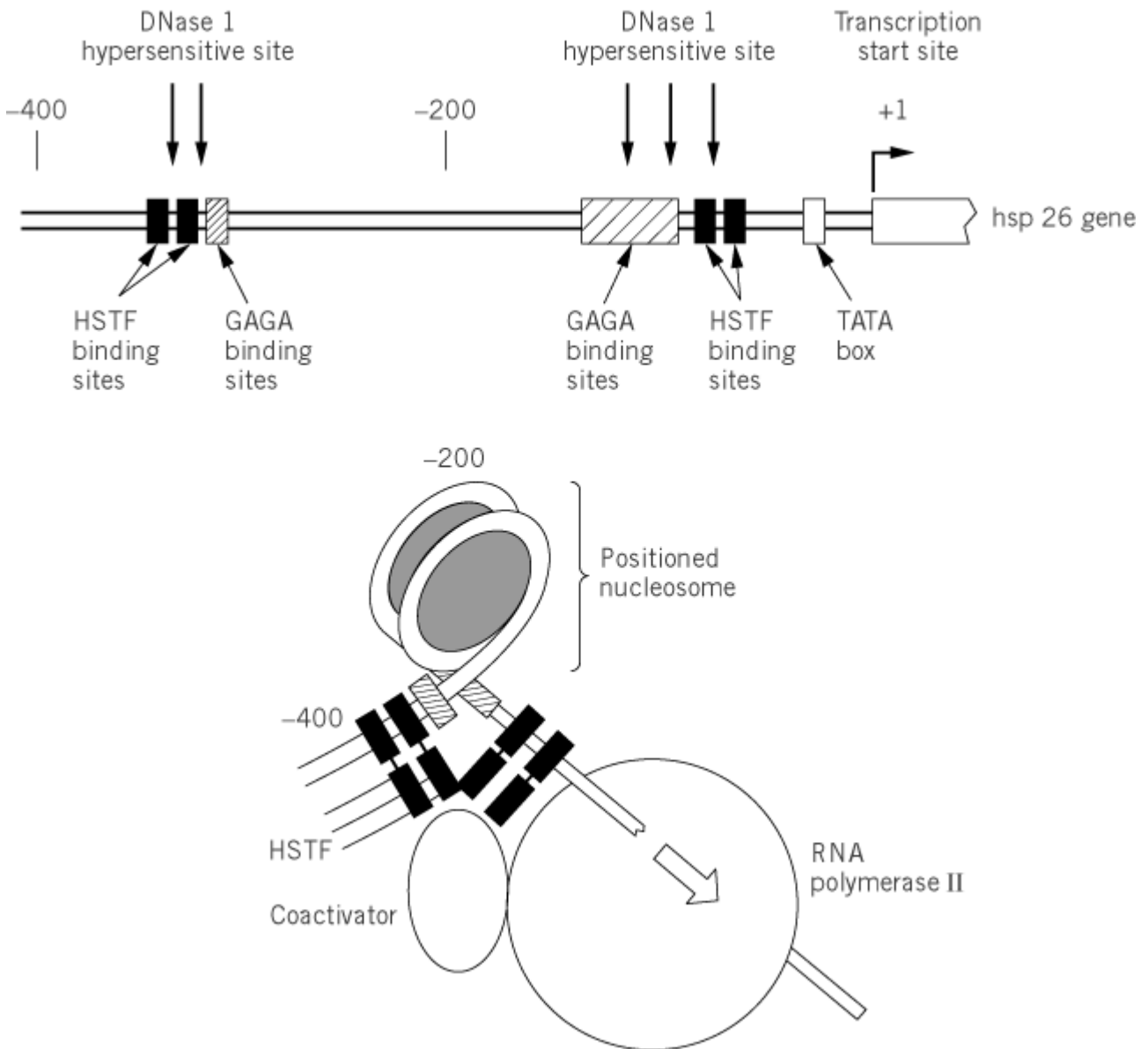
**Figure 2.** DNase I hypersensitive sites of the chicken b-globin locus. Contiguous genes are indicated as open arrows, the three distinct types of DNase I hypersensitive sites present in all cells specific to erythroid cells, and those that are developmentally regulated within the erythroid lineage are indicated.



One of the most thorough dissections of a DNase I hypersensitive site was carried out by Elgin and colleagues (4). The *Drosophila* **heat-shock** protein *hsp26* gene is very rapidly activated transcriptionally by raising the temperature of a fly to a stressful level (a heat shock of 34°C). Two DNase I hypersensitive sites exist at the promoter of the *hsp26* gene, including recognition sequences for the promoter-specific, heat-shock transcription factor (HSTF, a [leucine zipper](#) protein) (Fig. 3). Following heat shock, HSTF binds to these sites. In contrast, transcription factor TFIID is bound to the [TATA box](#) before and after heat shock. TFIID alone is insufficient to cause the *hsp26* gene to be transcribed, and the specific association of the HSTF protein is also required. High-resolution analysis revealed that a nucleosome is positioned between the proximal and distal binding sites for HSTF, i.e. between the two DNase I hypersensitive sites. In this case, the exact position of histone-DNA contacts within this nucleosome depends on the DNA sequence to which the histones bind (from -300 to -140 relative to the start site of transcription at +1), and also on adjacent DNA sequences. These are repeats of the type  $(CT)_n(GA)_n$ , which bind a specific *transacting* factor, the GAGA protein (7). The  $(CT)_n(GA)_n$  repeat regions are located to either side of the positioned nucleosome at -347 to -341 and at -135 to -85. The GAGA factor bound to these repeats may function as a “bookend” to determine exactly where the nucleosome will be positioned. Recent evidence suggests that the GAGA factor might function through ATP-dependent mechanisms to direct the assembly of a particular chromatin architecture actively. Transcription of the gene is regulated through the association of the HSTF with recognition elements at -51, -170, -269, and -340. The sites at -170 and -269 will be wrapped around the core histones in rotational frames that prevent HSTF association. However, the wrapping of DNA around the nucleosome will also bring the HSTF molecules bound to the sites at -340 and -51 into juxtaposition, and the histones may potentially facilitate transcription by causing a clustering of HSTF activation domains. Direct

evidence for transcriptional activation mediated by this nucleosome is yet to be established. It is clear, however, that the positioning of a nucleosome in this particular way allows key transcription factors to obtain access to essential regulatory elements in spite of the assembly of the gene into chromatin.

**Figure 3.** The specific chromatin organization of the *Drosophila* hsp26 promoter. Key cis-acting elements are indicated relative to the start site of transcription (hooked arrow). The organization of these sites on a specific nucleosomal scaffold is indicated together with the interactions necessary to prevent or activate transcription. A tethered RNA polymerase II molecule is released through events initiated by the binding of HSTF.



For a DNase I hypersensitive site on the vitellogenin genes, nucleosome positioning directed by the DNA sequence occurs between  $-300$  and  $-140$  relative to the start site of transcription at  $+1$ . The binding sites for the stimulatory transcription factors, the [estrogen receptor](#) and nuclear factor *1*, lie outside the region of DNA that is wrapped around the histones, at  $-300$  to  $-330$  and at  $-120$  to  $-110$ , respectively (8). When these sites are brought together either by positioning a nucleosome in between them (or artificially by deleting the intervening DNA), transcription is enhanced about 5- to 10-fold. This moderate stimulatory effect is much more significant than it might appear because

assembling a nonspecific chromatin structure would normally lead to a >20-fold repression of transcription as the binding sites for transcription factors are occluded by the histones. Thus in the vitellogenin gene example, nucleosome positioning has two roles: (1) to provide a scaffold that allows transcription factors to communicate more effectively; and (2) to prevent the formation of repressive histone-DNA interactions that may prevent any transcription factor from gaining access to a chromatin template.

DNase I hypersensitive sites have proven very useful in defining important DNA sequences. Among these are the four strongly nuclease-sensitive sites located 10 to 20 kbp upstream of the cluster of human  $\beta$ -globin genes. These sites, known as locus control regions (LCRs), represent *cis-acting* elements that allow genes integrated into a chromosome to be expressed independently of chromosomal position, i.e. [position effects](#) are abolished. A consequence of this is that LCRs allow each copy of a gene integrated in multiple copies to be expressed equivalently, so that gene expression is **copy-number** dependent. Therefore the LCR functions to control gene activity over an entire chromatin **domain**.

### Bibliography

1. C. Wu et al. (1979) *Cell* **16**, 797–806.
2. J. D. McGhee et al. (1981) *Cell* **27**, 45–55.
3. J. Wong et al. (1997) *EMBO J.* **16**, 7130–7145.
4. S. C. R. Elgin (1988) *J. Biol. Chem.* **263**, 19259–19262.
5. D. S. Gross and W. T. Garrard (1988) *Ann. Rev. Biochem.* **57**, 159–197.
6. M. Reitman and G. Felsenfeld (1990) *Mol. Cell. Biol.* **10**, 2774–2786.
7. Q. Lu, L. L. Wallrath, H. Granok, and S. C. R. Elgin (1993) *Mol. Cell. Biol.* **13**, 2802–2814.
8. C. Schild, F.-X. Claret, W. Wahli, and A. P. Wolffe (1993) *EMBO J.* **12**, 423–433.

### Hypervariable Locus

In general terms, a “hypervariable locus” is a genetic locus that shows a great degree of variability, or **polymorphism**, in the [genomes](#) of a population. High degrees of polymorphism occur at [hot spots](#) of [mutation](#) or [recombination](#), or they can be the consequence of a continuously varying strong [natural selection](#) pressure on a locus, such as that exerted on the **major histocompatibility locus**. The term “hypervariability” is normally used, however, to indicate the instability of [minisatellites](#), which is related to their sequence of 15 to 60 bp being tandemly repeated. The mutation frequency for changes in the number of tandem repeats has been estimated to vary between  $10^2$  and  $10^4$  per kbp per generation, which is one to three orders of magnitude higher than the rate of point mutation (1). **Unequal crossing over** within such a tandem repeat is likely to be an important generator of asymmetrical new alleles (2). Jeffreys noted that the core of his first hypervariable minisatellite shared some sequence characteristics with the **Chi sequence** of [lambda phage](#), a signal for homologous recombination, which is probably achieved by binding a **Rec** protein that locally unwinds and nicks DNA (3). The hypothesis that minisatellites are recombination hotspots has been partially confirmed by the observation that most highly variable human minisatellites tend to be located in the subtelomeric regions of [chromosomes](#) (see [Telomere](#)), precisely those that exhibit high recombination rates (4).

Changes in the lengths of minisatellite repeats occur both in the **germline** and in **somatic cells**,

possibly by sister [chromatid](#) exchanges. Somatic length mutations may represent markers of tumor cell progression and variation (5). With some hypervariable minisatellites, no difference has been found between the frequency of germline mutations in **oocytes** or **spermatocytes**. Given the large difference in the number of **mitosis** events that follow **meiosis** in the female and male germlines, it would seem that the cause of mutation is not always mitosis-dependent (6). Also, for some minisatellites, closely related **alleles** show the same flanking DNA sequences, whereas for other minisatellites flanking **haplotype** switching is observed. It should be concluded that recombinational events within minisatellites are not always the cause of length mutation.

Another mechanism possibly responsible for length mutation is slippage of the **DNA polymerase** during replication of the minisatellite, with the consequence of variation in one or very few units (7). A computer simulation to estimate the expected number and size range of alleles under a stepwise, replication slippage mutation model explained the available data better with the shorter [microsatellites](#) than with minisatellites (8).

Recently, methods have been devised for analyzing mutant alleles of single minisatellites by the polymerase chain reaction (**PCR**). Several features of the variability observed, such as its polarity, the addition of elements without evidence of **crossing over**, and rearrangement of elements with respect to their sequence in the homologous chromosome, suggest the involvement of complex **gene-conversion**-like events in the generation of mutant alleles (9).

### Bibliography

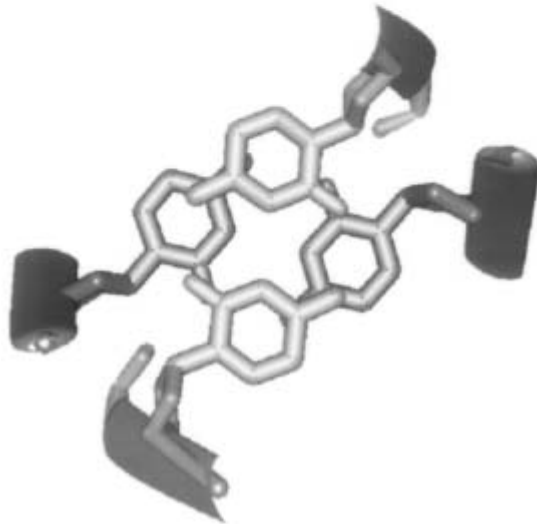
1. A. J. Jeffreys, V. Wilson, and S. L. Thein (1985) *Nature* **314**, 67–73.
2. A. J. Jeffreys, N. J. Royle, V. Wilson, and Z. Wong (1988) *Nature* **332**, 278–281.
3. G. R. Smith (1983) *Cell* **34**, 709–710.
4. A. P. Jarman and R. A. Wells (1989) *Trends Genet.* **5**, 367–371.
5. M. Perucho et al. (1995) *Cold Spring Harbor Symp. Quant. Biol.* **59**, 339–348.
6. A. J. Jeffreys et al. (1995) *Electrophoresis* **16**, 1577–1585.
7. G. Levinson and G. A. Gutman (1987) *Mol. Biol. Evol.* **4**, 203–221.
8. M. D. Shriver, L. Jin, R. Chakraborty, and E. Boerwinkle (1993) *Genetics* **134**, 983–993.
9. A. J. Jeffreys et al. (1994) *Nature Genet.* **6**, 136–145.

### I-Motif

The complementary DNA strand of the [telomere](#) TTGGGG sequence, CCCCAA, adopts an unusual [DNA structure](#) at low pH. The structure was first proposed by a [NMR](#) analysis (1) and later confirmed by [X-ray crystallography](#) (2) of cytosine-rich sequences under low pH conditions. It contains a highly unusual, mutually intercalated, four-stranded structure. This structure has been dubbed the i-motif, where two parallel duplexes intercalate with one another in an antiparallel orientation (Fig. 1). Each parallel duplex consists of hemiprotonated C:C<sup>+</sup> base pairs. It is unclear whether such a structural motif is associated with any biological function at present.

**Figure 1.** Side and end views of the structure of the i-motif (Nucleic Acids Database UDF062).





## Bibliography

1. K. Gehring, J. L. Leroy, and M. Gueron (1993) *Nature* **363**, 561–565.
2. L. Chen, L. Cai, X. Zhang, and A. Rich (1994) *Biochemistry* **33**, 13540–13546.

## Ideogram

An ideogram is a diagrammatic representation of the [karyotype](#) that shows all of the pairs of [homologous chromosomes](#) in the [nucleus](#). The pairs of chromosomes are lined up in order of size, so that the [centromeres](#) are aligned and the short arm is uppermost. An ideogram is a useful point of

reference for analyzing [mutations](#).

## Idiotypes

The concept of the idiotypic of **immunoglobulins** emerged in the early 1960s from two different approaches, one by Oudin with rabbit **antibodies**, the other by Kunkel and analysis of the immunochemical characteristics of human myeloma [proteins](#). At the time, Oudin had described the allotypic specificities as defining **antigenic** characteristics of a group of individuals within a given animal species. This was observed for rabbit immunoglobulins and was the consequence of **allelic** variants of both the heavy and light chains. It was then found that when rabbit antibodies prepared against *Salmonella typhi* were injected into another rabbit expressing the same (known) allotypes, they induced the synthesis of antibodies that specifically recognized the anti-*Salmonella* antibodies raised in the first rabbit. It was also shown that they did not react with normal rabbit serum taken before the *Salmonella* injection, thereby excluding that they identified a new allotypic specificity. Idiotypic specificities were thus defined as antigenic specificities characteristic of one antibody (the idiotypic) produced by one animal and specific for one given antigen. Antibodies produced by the second animal were termed **anti-idiotypic antibodies**. Idiotypes are also referred to as Ab1 while anti-idiotypes are referred to as Ab2, in a simplified nomenclature.

At the same time, Kunkel was comparing the immunochemical characteristics of human myeloma proteins with those of normal immunoglobulins. He found that antibodies raised against one given myeloma protein, which cross-reacted against normal immunoglobulins, became specific for the immunizing protein once extensively adsorbed on normal immunoglobulins. This was taken as indicative of determinants specific for any given monoclonal product, although it was not entirely clear whether it was also linked to the pathological “abnormal” nature of the myeloma protein.

Structural analysis confirmed, as anticipated, that idiotypic determinants, or idiotopes, were present on the Fab fragments and, more precisely, on the variable regions of H and/or L chains (see [Immunoglobulin Structure](#)), depending on the idiotypic. Extensive analysis of idiotypic structure performed by a combination of immunochemical and structural analysis of anti-**hapten** antibodies revealed several types of idiotopes. Some Ab1–Ab2 interactions could be inhibited by a hapten of Ab1, indicating that the corresponding idiotopes were part of the antibody combining site, whereas other were not. Genetic analysis also revealed that some idiotypic specificities were common to several idiotypes, whereas others were strictly specific for one given Ab1, leading to the distinction of public and private idiotypes, a notion that was directly related to the dual origin of antibody diversity, germline encoded and somatically generated. Idiotypy was studied extensively when it was discovered that the cascade  $\text{Ag(X)} \rightarrow \text{Ab1} \rightarrow \text{Ab2}$  could be continued and amplified in a large idiotypic network of interactions, providing a basis for autoregulation of the immune system. An especially interesting observation was that Ab2 could induce certain Ab3 molecules that resembled Ab1 in that they could bind the original Ag(X) antigen. This led to the definition of the “internal image” of the antigen, proposed by Jerne, containing the idea that the collection of normal immunoglobulins of an individual could represent a huge population of internal images of the outside antigen world. It also stimulated approaches of idiotypic vaccines that could have been used in place of classical antigens, an interesting idea whenever antigens were either difficult to identify or poorly immunogenic.

See also entries [Autoimmunity](#), and [Autoimmune Diseases](#).

Suggestions for Further Reading

- A. Coutinho (1995) The network theory: 21 years later. *Scand. J. Immunol.* **42**, 3–8.
- C. A. Bona (1996) Internal image concept revisited. *Proc. Soc. Exp. Biol. Med.* **213**, 32–42.

## IgA

Immunoglobulin A (IgA), accounts for about 10% to 15% of the circulating **immunoglobulins**. The main feature of IgA is its presence in high concentrations in all mucosal secretions (intestine, lung, colostrum, tears, etc.), and for this reason it is generally considered of prime importance in fighting pathogens, against which it constitutes a first line of defense. The presence of elevated amounts of IgA in secretions correlates with their special structural features, namely their ability to dimerize and to have an extra Ig secretory component added upon passage through epithelial cells.

The basic Ig structure of IgA is a classical H<sub>2</sub> L<sub>2</sub> monomer in which the heavy chains are a chains, composed of the usual V<sub>H</sub> region and three C<sub>H</sub> domains. In humans there are two **isotypes**, IgA1 and IgA2, encoded by distinct C<sub>H</sub> genes. As a result, the two isotypes differ from each other by small amino acid sequence differences and by the length of their hinge regions. In addition, the IgA2 isotype has two **allelic** variants, recognized by specific [antisera](#) as two distinct allotypes, IgA2 m(1) and IgA2 m(2). The IgA2 m(1) has a special structural feature, because its heavy chain lacks the [cysteine](#) residue that is covalently **disulfide-bonded** to the light chains in other isotypes; consequently, the H–L interactions are exclusively noncovalent, whereas the light chains are covalently dimerized. Carbohydrate moieties account for about 10% of the molecule and are mostly attached in the vicinity of the hinge region. In the serum, the IgA are essentially monomeric, whereas they polymerize in the vicinity of epithelial cells, with the formation of dimers (major), but also trimers and tetramers. When polymerized, the IgA gains an extra chain, termed the [J chain](#) [which is not the same as the **J (joining) region**], which is disulfide-bonded to part of the IgA heavy chains. Once dimerized, the IgA molecule is covalently bound to the poly-Ig receptor, expressed at the basal pole of the epithelial cell. This complex is internalized, and the polymerized IgA is released at the apical pole of the cell in the lumen of the corresponding organ upon **proteolytic** cleavage between two **domains** of the poly-Ig receptor, so that one subunit remains covalently linked to the IgA polymer, forming the so-called secretory component. Besides allowing the IgA to migrate through the epithelial cell, the secretory component is thought to protect the immunoglobulin from cleavage by [proteinases](#). The local concentration of IgA may be very high and, due to its polymerized state, this Ig has an increased avidity for [antigen](#). Because of its abundance in colostrum, it is of crucial importance for passive antibody transfer to the newborn and thus ensures a temporary protection against pathogens, because the immune system will remain immature for some time. It should be recalled, however, that the protective effect of colostrum and/or milk greatly varies from one species to the other. Whenever the placental barrier is easily passed by an active transfer of [IgG](#), the importance of colostrum IgA is minimized. In humans, protection is ensured by transplacental transfer. An example of the reverse situation is seen in cattle.

Mucosal immunity, as immunologists like to refer to it whenever IgA is involved, is presently a rapidly expanding field. It has not been studied sufficiently, considering that it appears to be of prime importance in fighting against pathogenic microorganisms.

See also entries [Class Switching](#), [Isotype](#), and [J Chain](#).

Suggestions for Further Reading

M. E. Lamm (1997) Interaction of antigens and antibodies at mucosal surfaces. *Annu. Rev. Microbiol.* **51**, 311–340.

J. Mestecky and J. R. McGhee (1987) Immunoglobulin A: molecular and cellular interactions involved in IgA biosynthesis and immune response. *Adv. Immunol.* **40**, 153–245.

B. A. Hendrickson et al. (1996) Lack of association of secretory component with IgA in J chain-deficient mice. *J. Immunol.* **157**, 750–754.

## IgD

IgD is an [immunoglobulin](#) that is present on the surface of mature [B cells](#). Its circulating concentration is in the range of 40–400  $\mu\text{g/mL}$ , which does not contribute significantly to the overall amount of secreted Ig (compare with the 10 mg/mL range for [IgG](#)). IgD exists as only one [isotype](#) and is exclusively a  $d_2k_2$  or  $d_2l_2$  monomer. The d chain in humans has three **constant domains**, whereas there are only two in murine species, in which  $Cd_2$  is replaced by an unusually long hinge region. The main characteristic of IgD is to be essentially coexpressed at the surface of the same B cell with [IgM](#). This does not contradict the [clonal selection theory](#), because the two B-cell receptors express the same variable regions and hence have the same antibody specificity. This is because the [switch region](#) that is located at the 5' position of all C gene isotypes is absent from the Cd locus. Coexpression of m and d chains results from [alternative splicing](#), thus leading to both IgM and IgD that are associated to same k or l chains within a given B cell. Occasionally, clones expressing solely IgD without IgM can be encountered. This is most apparent in rare myeloma cells that are IgD producers. In that case, the absence of IgM probably results from a deletion of the Cm genes at the genomic level. This seems also to be the case of a recent observation of [B cells](#) detected in human tonsils that express exclusively IgD and have surprisingly accumulated an unusually large number of somatic mutations. Despite their very peculiar specialization as surface Ig, the role of IgD is still very poorly understood, although some arguments point to their possible role in B-cell induction of **tolerance**.

See also entries [Isotype](#).

### Suggestions for Further Reading

F. R. Blattner and P. H. Tucker (1984) The molecular biology of IgD. *Nature* **307**, 417–422.

J. Wienands, J. Hornbach, A. Radbruch, C. Riesterer, and M. Reth (1990) Molecular components of the B cell antigen receptor complex of class IgD differ partly from those of IgM. *EMBO J.* **9**, 449–455.

C. C. Goodnow et al. (1988) Altered immunoglobulin expression and functional silencing of self-reactive B lymphocytes in transgenic mice. *Nature* **334**, 482–484.

Y.-J. Liu et al. (1996) Normal human  $\text{IgD}^+\text{IgM}^-$  germinal center B cells can express up to 80 mutations in the variable region of their IgD transcripts. *Immunity* **4**, 603–613.

## IgE

IgE is the [immunoglobulin](#) most commonly known as the main cause of anaphylaxis or, more generally speaking, of the immediate hypersensitivity states, such as asthma, hay fever, atopic dermatitis, and other well-known manifestations of atopic patients. IgE is present in normal individuals at very low concentrations, normally somewhat less than 1  $\mu\text{g/mL}$ . This concentration may be considerably elevated in patients suffering from one or another form of immediate hypersensitivity, and also in patients with internal parasite infections. This points to a major role as support of anaphylactoid states and suggests a possible role in the defense against parasites, a view that has been seriously reinforced of late.

IgE is a monomer of molecular formula  $\epsilon_2\kappa_2$  or  $\epsilon_2\lambda_2$ . It may, as all the other Ig subclasses, be expressed as cell-surface or circulating Ig, and it has an elevated carbohydrate content (>10% by weight). The  $\epsilon$  chains have four constant C<sub>D</sub> **domains**, and the molecular weight of IgE averages 180 kDa. IgE does not fix **complement** and is heat-labile. Its role in driving the immediate hypersensitivity reactions is due to the fixation of circulating IgE antibodies on **Fc receptors** of very high affinity; FcR $\epsilon$ I is located on several cell types, mostly basophils and mast cells. These large cells contain granules that contain a high concentration of pharmacologically active substances, like histamine, neutral [proteinases](#), proteoglycans, and **tumor necrosis factor  $\alpha$** ; (TNF $\alpha$ ). Contact of an [antigen](#) (here denominated an allergen) with specific IgE fixed to such cells of an atopic patient who has already been sensitized will activate the target cell and release the pharmacologically active substances, thus initiating a large cascade of secondary reactions. Histamine will bind selectively to H1-type receptors and exert various deleterious effects, such as vasodilatation, local exudation, and spasmogenic activity of intestinal, uterine, or bronchial muscles, depending on the localization of the target. TNF $\alpha$  liberated upon mast-cell activation may have a beneficial effect in antimicrobial defense. Other mediators are synthesized as a response to IgE triggering. Those include prostaglandins, PGD<sub>2</sub>, leukotrienes, and so on, and a large number of **cytokines** that considerably amplify the response and especially initiate a strong inflammatory reaction. Fixation of IgE on other cell types, like eosinophils, **macrophages**, or platelets, can also occur through the Fc $\epsilon$ RII, which is of lower affinity than Fc $\epsilon$ RI, and trigger liberation of many mediators, including free radicals that may be beneficial by helping to fight against pathogens, although most of the effects are quite aggressive for the host tissues.

Discovery of the two main populations of helper T cells, Th1 and Th2, which drive a fine regulation of the immune response through a complex network of cytokines, has opened the way to a reinvestigation of the possibility of modulation of the expression of IgE, which is of prime importance for atopic patients, but also to fight more effectively against parasites, which remain quite elusive pathogens to date. IgE production is essentially driven by cytokines of the Th2 group, particularly [interleukin](#) 4, which was shown to favor the **class switch** to this Ig class.

See also entries [Antibody](#), [Class Switching](#), [Immunoglobulin](#), and [Isotype](#).

#### Suggestions for Further Reading

B. J. Sutton and H. J. Gould (1993) The human IgE network. *Nature* **366**, 421–428.

T. Mossmann and R. L. Coffman (1989) Th1 and Th2 cells. Different patterns of lymphokine secretion lead to different functional properties. *Annu. Rev. Immunol.* **7**, 145–174.

## IgG

IgG is the major class of [immunoglobulin](#) (Ig) and has long been taken as the basic Ig model. IgG are found either as free circulating molecules in the bloodstream, where they account for about 70% of total Ig, or anchored at the [B-cell](#) surface, constituting the B-cell receptor in association with the IgαIγb heterodimer signaling module. They contain two g and two light chains, either k or l, and about 2% by weight of carbohydrate, located on the C<sub>H2</sub> **domains**. IgG is expressed as four distinct subclasses, termed IgG1, IgG2, IgG3, and IgG4 in humans, resulting from the presence of the corresponding g1 to g4 heavy chains encoded by separate germline Cg genes. The four **isotypes** have minor structural differences, both in the length of their hinge region and in [primary structure](#), which remain 95% identical; this suggests that the amplification of g isotypes is a rather recent event. IgG subclasses also exist in other species. In the mouse there are also four discrete isotypes called IgG1, IgG2a, IgG2b, and IgG3. That there is no correspondence between the various human and mouse isotypes confirms that the subclass expansion is characteristic of each animal species and occurred after they diverged. The different isotypes are endowed with somewhat different biological properties. All except for IgG4 bind **complement**. The first component of the complement cascade, C1q, has a binding site on the C<sub>H2</sub> domain. This site becomes accessible once the antibody molecule has interacted with antigen, presumably as the result of an opening of the Fab–Fc angle, which is made possible because of the flexibility of the hinge region (see [Immunoglobulin Structure](#)). The IgG subclasses bind to different **Fc receptors**. The major subclass is IgG1, followed by IgG2. Expression of the various subclasses during [development](#) follows different kinetics, with IgG2 and IgG4 being synthesized only slowly after birth. Allotypic markers termed Gm factors have long been identified on IgG and have provided an interesting basis for studying human population genetics.

IgG can be cleaved, in the native state, by a variety of **proteolytic** enzymes, due to the flexibility of the hinge region and to the somewhat accessible short regions that separate the various domains. Porter (1) reported the isolation of Fab and Fc fragments upon cleavage with [papain](#); IgG may also be cleaved with **pepsin**, which largely destroys the C<sub>H2</sub> domain, leaving a (Fab')<sub>2</sub> fragment, which has two pieces slightly larger than an Fab **disulfide-bonded** through the inter-heavy-chain bridges that have remained untouched, and the Fd' fragment, which covers roughly the C<sub>H3</sub> COOH-terminal domain. Cleavage with **plasmin** cleaves off the last C<sub>H3</sub> domain, leaving a major Fab<sub>b</sub> fragment (for antibody and complement binding). Finally, it is possible in a limited number of cases to release the variable V<sub>H</sub>–V<sub>L</sub> domains noncovalently bound as the Fv fragment, which contains one intact antibody combining site. The three-dimensional [protein structures](#) of the isolated domains are quite independent of their interactions with other parts of the IgG molecule. Extensive [protein engineering](#) has been made on this structure, by introducing the corresponding [complementary DNAs](#) into vectors that are expressed in bacteria or other systems, such as the baculovirus. Extensive [phage display libraries](#) have been constructed that provide a huge combination of V<sub>H</sub>–V<sub>L</sub> pairs having an almost infinite set of discrete recognition Ig-like molecules.

During the early phases of the primary [immunization](#), IgM is rapidly replaced by IgG upon [class switching](#), a molecular mechanism that operates at the genomic level, leading to the replacement of one isotype by the other, without changing the variable regions and thus conserving the antibody specificity. [Somatic hypermutations](#) occur primarily in IgG molecules, so they become increasingly adapted to the exquisite recognition of **epitopes** of the immunizing [antigen](#). This implies both an enzymatic mechanism for mutations and a strong selection exerted by the antigen itself on IgG-producing B cells. Some of the selected clones with high affinity will ultimately differentiate into plasma cells and produce circulating IgG, whereas others will become [memory cells](#), expressing membrane IgG of high affinity that will rapidly react in a subsequent reexposure to the antigen, initiating a secondary response.

See also entries [Antibody](#), [B Cell](#), [Class Switching](#), [Immune Response](#), and [Immunoglobulin](#).

## Bibliography

1. R. R. Porter (1959) The hydrolysis of rabbit  $\gamma$ -globulin and antibodies with crystalline papain. *Biochem. J.* **73**, 119–126.

## Suggestion for Further Reading

2. D. Burton (1990) Antibody: the flexible adaptor molecule. *Trends Biochem. Sci.* **15**, 64–69.

## IgM

Immunoglobulin M (IgM) is the first [immunoglobulin](#) to appear after [immunization](#), the first to occur in ontogeny and in [phylogeny](#), and the first to be expressed at the surface of [B cells](#) during their [differentiation](#). Most of these properties are a result of the fact that the genes encoding their  $C_H$  regions are the closest to the  $V_H$  locus. IgM is built on the classical  $H_2L_2$  model as either  $m_2k_2$  or  $m_2l_2$  monomers (see [Immunoglobulin Structure](#)). As such, these monomers are expressed at the B-cell surface, where they represent, associated with the  $I\alpha$ – $I\beta$  heterodimer signaling module, the B-cell receptor that interacts with the [immunogen](#) upon primary immunization. In the serum, most of the circulating IgM is a pentamer,  $(m_2k_2)_5$  or  $(m_2l_2)_5$ , in which monomers are covalently **disulfide-bonded**, a linkage that involves the addition of the [J chain](#), similar to that present in the [IgA](#) dimer. The constant region of the m chain does not contain a hinge region, but instead has an additional C domain, as compared with [IgG](#). IgM has a high content of carbohydrate (10% to 12% by weight) which is bound at [N-glycosylation](#) sites of the first three  $C_m$  domains.

The basic monomer has a molecular weight of 180 kDa, and the pentamer has a molecular weight of about 970 kDa. Because of their size, circulating IgM subclasses were originally designated as [macroglobulins](#). They are found in a lymphoproliferative disease, the Waldenstrom macroglobulinemia, where they are overexpressed as the result of the malignant amplification of an IgM-producing B-cell clone. In normal serum, IgM accounts for 10% to 15% of the total circulating immunoglobulins. IgM behaves as expected as decavalent antibodies for binding to **epitopes**, provided that they are located on a small antigen, such as a [hapten](#). On larger antigens, as a result of steric hindrance, IgM pentamers will interact with at most five molecules of antigen. Individual antibody-combining sites are most often of low affinity but the pentamerization potentiates the binding efficiency, demonstrated by the “avidity” of the antibody, a typical immunologist parameter not really defined in mathematical terms. IgM is very efficient in **complement** binding and thus is highly lytic for cellular targets, which may be bacteria or eukaryotic cells like red blood cells; in this case one uses the term *hemolysins*, rather than *hemagglutinins*, when the complement component is not participating in the final reaction.

The so-called natural antibodies are most often of the IgM type. This is the case of natural anti-A and anti-B hemagglutinins, described as the classical **alloantibodies** responsible for major transfusion incompatibilities. More generally speaking, natural antibodies are frequently endowed with **autoimmune** recognition properties and thus may bind a large variety of molecules, such as [tubulin](#), [actin](#), double-stranded DNA, and so on. These natural antibodies are harmless and most likely represent the basic [repertoire](#) of germline-encoded immunoglobulins, which therefore have nearly no somatic mutations. As a result, their specificity is very low, which makes them bind many antigens, although usually only weakly. Why these antibodies are not aggressive to the organism is still

subject to debate, although the combination of two parameters, low affinity and low concentration, is possibly the answer.

As already mentioned, IgM is the first group of isotypes to be produced upon primary immunization. Whenever the antigen is T-dependent, which is the most frequent situation, IgM is rapidly replaced by IgG, which will warrant a better affinity because of the acquisition and selection of [somatic hypermutations](#). Immunization with polysaccharides, which are essentially T-independent antigens, will thus be limited to the production of IgM antibodies. In phylogeny, IgM is also the first group of immunoglobulin classes that have been identified in primitive fishes. It should be mentioned that the classical immune response, with the existence of circulating antibodies, seems to have occurred with the lower vertebrates. Invertebrates have, however, other types of responses to fight very efficiently against pathogens.

See also entries [Antibody](#), [Immune Response](#), [Immunoglobulin](#), and [J Chain](#).

#### Suggestion for Further Reading

A. C. Davis and M. J. Shulman (1989) IgM-molecular requirements for its assembly and function. *Immunol. Today* **10**, 118–128.

## Illegitimate Recombination

Illegitimate recombination includes those events that do not fall into any of the more clearly defined categories of [recombination](#) ([homologous recombination](#), **site-specific recombination**, and [transposition](#)). Illegitimate recombination may itself include an assortment of types of events. Parental DNA molecules or DNA segments in illegitimate recombination often share a few homologous base pairs. Some of these events result from [DNA replication](#) errors and others from joining broken DNA ends. Illegitimate recombination generates deletions and translocations and therefore can have profound genetic consequences.

One type of illegitimate recombination widely studied in bacteria results from slippage between the template and newly replicated DNA strands. The prototype for this mechanism is that suggested by Streisenger et al. (1) for generating [frameshift mutations](#) at points of repeated nucleotides. For example, DNA containing on one strand  $A_5$  can generate  $A_4$  or  $A_6$  by loss or gain of a base pair by the following proposed mechanism. Replication pauses within the repeat. The newly synthesized strand unwinds and anneals with the template, but out of register, leaving one or more nucleotides as a single-stranded loop on either the template or the newly synthesized strand. Replication resumes, and this newly synthesized strand has fewer or more nucleotides than the parent. In the absence of repair synthesis, another round of faithful replication fixes the mutation.

A similar event can occur between direct repeats of any sequence that may or may not be separated by intervening base pairs. If the repeats do not themselves have internal repeats, interaction between them produces a deletion or duplication of one copy of the repeat plus the intervening sequence. One could imagine that these events occur by homologous recombination where the repeats are out of register (unequal crossing over), but in *Escherichia coli* these events are independent of known homologous recombination functions. The events are not reciprocal, that is, a deletion is not formed concurrently with a duplication. Deletions can occur between repeats separated by thousands of base pairs, although the frequency decreases as this separation becomes larger (2).



Linear DNA introduced into mammalian cells by [transfection](#) often undergoes illegitimate recombination by integration into the [genome](#) or by end-joining. End-joining has been extensively studied and occurs by ligation of the ends without any apparent requirement for [homology](#). If excision of nucleotides from one or both ends occurs before or after transfection, a few homologous base pairs may be exposed, and annealing may occur between them. Filling of gaps and ligation can generate a novel joint. These events may be paradigms for integration of exogenous linear DNA and for chromosomal translocations. These events typically occur between DNA molecules sharing little or no homology and may occur between broken chromosome ends that swap partners and rejoin.

The DNA ends that provoke illegitimate recombination may be generated by a variety of mechanisms. Linear DNA may be introduced into cells, as in transfection. [Restriction Enzymes](#) may be induced to cut DNA at special sites in cells. DNA **topoisomerases** make transient double-strand breaks that may be illegitimately rejoined. This process is stimulated by topoisomerase inhibitors. Pausing or termination of DNA replication generates single-strand gaps whose ends may provoke illegitimate recombination, perhaps after conversion of the gap into a double-strand break. The variety of mechanisms that generate breaks may account for the variety of types of illegitimate recombination observed in cells.

Illegitimate recombination is not entirely a haphazard process. Regulated forms of this process occur in generating functional [immunoglobulin](#) genes by [gene rearrangements](#) and in eliminating micronuclear DNA in *Tetrahymena* **macronucleus** development. Other types of recombination previously considered illegitimate are recognized now as site-specific recombination and transposition. Further studies of illegitimate recombination may reveal unappreciated mechanisms and the biological roles of these events.

#### Bibliography

1. G. S. Sreisinger et al. (1967) Cold Spring Harbor Symp. Quant. Biol. **31**, 77–84.
2. S. T. Lovett et al. (1994) Mol. Gen. Genet. **245**, 294–300.

#### Suggestions for Further Reading

3. N. D. Allgood and T. J. Silhavy (1988) "Illegitimate recombination in bacteria", In *Genetic Recombination* (R. Kucherlapati and G. R. Smith, eds.), American Society for Microbiology, Washington, D.C., pp. 309–330.
4. D. Erhlich (1989) "Illegitimate recombination in bacteria", In *Mobile DNA* (D. E. Berg and M. M. Howe, eds.), American Society for Microbiology, Washington, D.C., pp. 799–832
5. S. D. Ehrlich et al. (1993) Mechanisms of illegitimate recombination, *Gene* **135**, 161–166.
6. M. Meuth (1989) "Illegitimate recombination in mammalian cells", In *Mobile DNA* (D. E. Berg and M. M. Howe, eds.), American Society for Microbiology, Washington, D.C., pp. 833–860.
7. D. Roth and J. Wilson (1988) "Illegitimate recombination in mammalian cells", In *Genetic Recombination* (R. Kucherlapati and G. R. Smith, eds.), American Society for Microbiology, Washington, D.C., pp. 621–653.

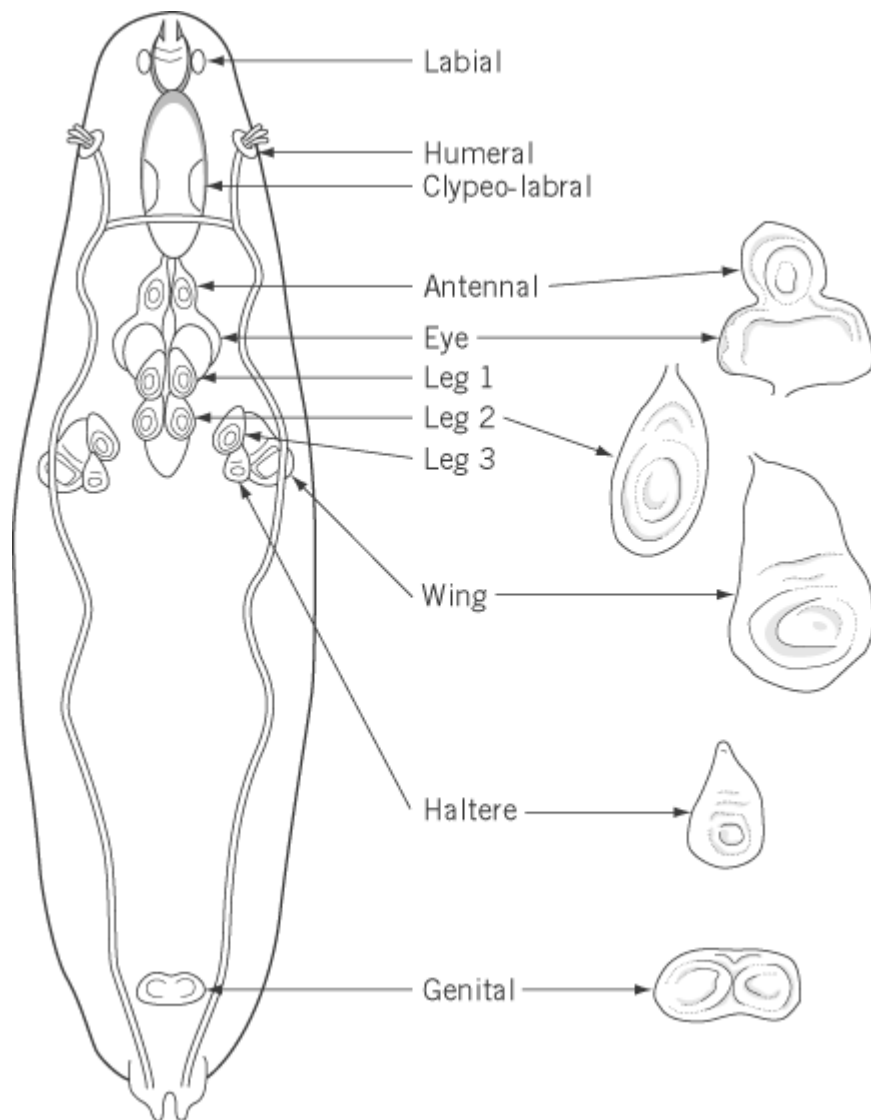
#### Imaginal Disc

The past decade has witnessed enormous progress in understanding the molecular mechanisms that control the [development](#) of imaginal discs in insects. Lessons obtained from the studies of imaginal discs are being used to unravel mechanisms underlying vertebrate development, because many **genes**

and developmental pathways are conserved throughout evolution.

The development of imaginal discs is a feature unique to the holometabolous insects (beetles, moths, flies, and butterflies). The life cycles of holometabolous insects include four stages: embryo, larva, pupa, and adult (also known as *imago*). Development through these stages does not occur simply by increasing the sizes of exiting structures; instead, it involves a dramatic transformation from larva to adult. The transformation is achieved by the replacement of larval tissues with a special population of cells, known as *imaginal discs*, which form sac-like structures residing inside the larva. Imaginal discs have been studied most extensively in the fruit fly *Drosophila melanogaster* (1-4). There are 10 major pairs of imaginal discs, which give rise to adult cuticular structures of the head and the thorax, and a single genital disc, which forms the adult genitalia. The location of the imaginal discs in the mature larva (third instar stage) is illustrated in Figure 1. Studies in *Drosophila* have shown that imaginal discs are initially set aside during embryogenesis, apart from most embryonic cells that will contribute to the larva. In the embryo, the imaginal disc precursor cells can be visualized by the expression of molecular markers or by their particular sizes, shapes, and movement. The imaginal discs appear in the newly hatched larva as local thickenings of the epidermis and, in the case of the wing disc, contain around 40 cells (5). Subsequently, the imaginal disc cells proliferate rapidly during larval stages. By the end of the larval stage, the largest imaginal disc, that of the wing, contains approximately 50,000 cells (6). During pupation, the imaginal disc cells go through morphogenetic movements, in response to a pulse of the molting hormone [ecdysone](#), to give rise to specific adult structures, accompanied by the breaking-down of the larval tissues.

**Figure 1.** Locations of imaginal discs in a third instar larva, and the morphology of several major imaginal discs. There are 10 major pairs of imaginal discs, which give rise to adult cuticular structures of the head and the thorax, and a single genital disc, which forms the adult genitalia. (After Ref. [86](#))



The molecular studies of imaginal discs in *Drosophila* have been addressing the following major questions. First, how are the imaginal disc precursor cells set aside (or established) during embryogenesis? Second, how is the growth of imaginal discs in the larva regulated? Third, how do developmental patterns form in the mature imaginal disc along its anterior–posterior, dorsal–ventral, and proximal–distal axes? And fourth, how is the identity of a given imaginal disc specified? The following is a brief summary of our current understanding on these topics.

### 1. Establishment of the Imaginal Disc Precursor Cells in the Embryo

The fate of imaginal disc precursor cells is determined by intersecting anterior–posterior and dorsal–ventral signals in the embryo. These signals are provided by two sets of patterning genes along the embryonic anterior–posterior and dorsal–ventral axes. Along the anterior–posterior axis, a regulatory hierarchy involving gap genes, [pair-rule genes](#), and segment polarity genes determines the different cell fates. In particular, the activities of several segment polarity genes, including *wingless* (*wg*), *engrailed* (*en*), *patched* (*ptc*), and *hedgehog* (*hh*), are required for the consolidation of the parasegment boundaries, which appear during early development of the embryo (7, 8). Among these genes, *wg*, encoding a secreted molecule of the Wnt family of proteins (9), is required for the establishment of the thoracic wing and leg discs, which originate from primordia spanning the parasegment boundaries (10, 11). Like the anterior–posterior axis, the embryonic dorsal–ventral axis is also patterned by a hierarchy of maternal and zygotic genes. In particular, the [decapentaplegic](#)

(dpp) protein, a member of the **transforming growth factor b** (TGF- $\beta$ ) gene family of signaling molecules (12), and other molecules in the dpp signaling pathway function to specify dorsal cell fates (13, 14). In contrast, the DER protein, the *Drosophila* **EGF receptor** homologue, and other genes in the DER signaling pathway are required to specify ventral cell fates (15, 16). The positions of the thoracic disc precursor cells along the dorsal–ventral axis are determined by the combined repressing activities of the dpp and the DER proteins, which limit the thoracic disc precursor cells to lateral positions of the embryo (17). Therefore, the thoracic disc precursor cells are established by the intersecting signals of Wg along the anterior–posterior axis, and dpp and DER along the dorsal–ventral axis.

## 2. Regulation of the Imaginal Disc Growth in the Larva

As mentioned above, imaginal discs increase their number of cells by thousands of fold during larval development. Misregulation of disc growth can result in imaginal discs with abnormal sizes and morphology. Genetic studies have identified a large number of mutants that are either defective or excessive in disc growth. The disc overgrowth mutants are of particular interest, because they may provide genetic models for [tumor suppressor genes](#) in humans. Recessive mutations in several dozens of genes have been identified that cause two types of disc overgrowth phenotypes: hyperplasia and neoplasia (18-20). Hyperplastic discs retain the normal epithelial organization and the capacity to differentiate, whereas the neoplastic discs do not. These several dozens of genes are known as tumor suppressor genes in *Drosophila*, and over 10 of them have been characterized molecularly. For example, one of these genes, *lethal(1)discs large* (*dlg*), is the first identified member of a evolutionary conserved family of proteins, called membrane-associated guanylate kinase homologues (MAGUKs) (21). In *Drosophila*, *dlg* is localized at and required for the formation of septate junctions in epithelial cells, and loss of *dlg* function causes epithelia to lose their organization and overgrow (22). Further studies of these *Drosophila* tumor suppressor genes are likely to shed light on general mechanisms in the regulation of cell proliferation.

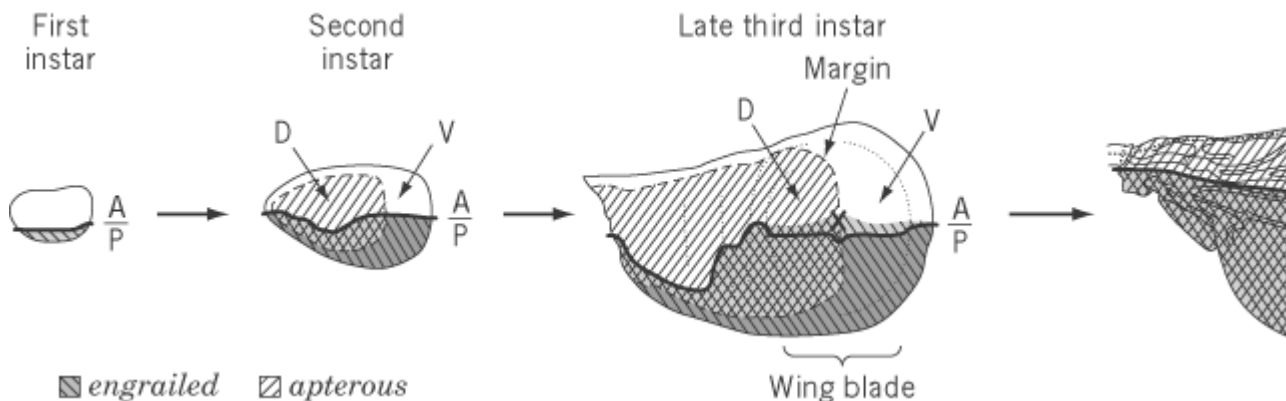
## 3. Pattern Formation in the Imaginal Discs

### 3.1. Anterior–Posterior Patterning

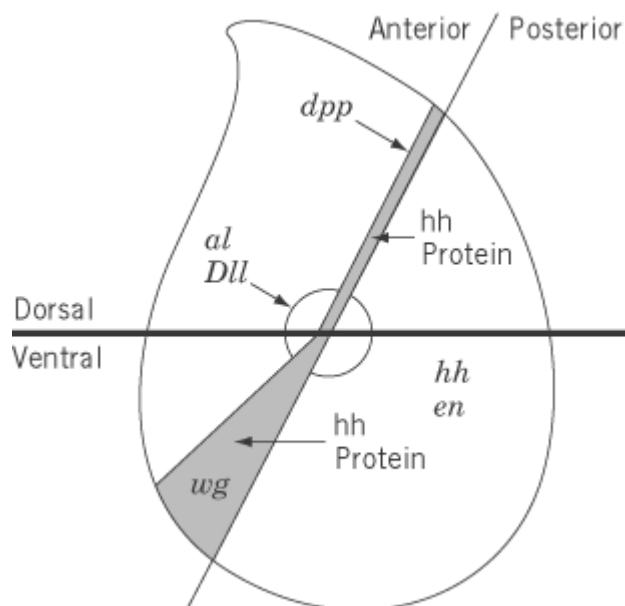
Previous studies have shown that imaginal discs are subdivided into distinct cell populations, called compartments (23, 24). The compartmental boundary serves as a line of lineage restriction, and cells in one compartment do not intermingle with cells in the other. The anterior and posterior compartments of the thoracic discs can be traced back to the embryonic stage when the disc precursor cells are specified initially along the anterior–posterior axis of the embryos (25). In these discs, the posterior compartment is defined by the expression of the [transcription factor](#) *engrailed* (*en*) (26), which activates the expression of *hedgehog* (*hh*) in the same compartment (27, 28). In the wing disc, the hedgehog protein synthesized in the posterior compartment diffuses a short distance into the anterior compartment to induce the expression of *dpp* in a stripe of cells just anterior to the anterior–posterior boundary (27, 29). The decapentaplegic protein in turn acts as a long-range [morphogen](#) to organize the growth and patterning of the whole wing (30, 31) (Figure 2). The leg disc is very similar to the wing disc, except that in the leg disc the induction of *dpp* by the hedgehog protein is limited to the adjacent dorsal anterior cells, whereas the adjacent ventral anterior cells are induced by the hedgehog protein to express *wg* (32, 53). In addition, *dpp* and *wg* mutually repress the [transcription](#) of each other (Figure 3) (34, 35). The unpaired genital disc consists of three primordia, the female genital, the male genital, and the anal primordia. Each primordium of the genital disc is divided into anterior and posterior compartments. Here, *hh*, *dpp*, and *wg* are used to pattern each primordium in a manner similar to how they function in the leg disc (36, 37). In contrast to the wing, leg, and genital discs, the eye disc is not divided into any lineage-restricted compartments. Pattern formation in the eye disc is accompanied by movement of the morphogenetic furrow (MF), a morphological distinguishable indentation, which sweeps across the eye disc from posterior to anterior. Cells anterior to the MF are undifferentiated, whereas cells posterior to the MF undergo cellular differentiation (37). The initiation of eye disc differentiation is regulated by interactions between the positive regulator *dpp* and the negative regulator *wg*. Subsequently, the progression of

differentiation requires *hh*. Despite the lack of the anterior and posterior compartments, the regulatory relationships between *hh*, *dpp*, and *wg* are retained in the eye disc. It has been shown that *dpp* and *wg* suppress the expression of each other, and the hedgehog protein produced by cells posterior to the MF induces *dpp* expression within the MF (39, 40).

**Figure 2.** Axis formation and gene expression in the developing wing. In the first instar larva, the wing disc has already posterior compartment is identified by the expression of *engrailed* (26). In the second instar larva, the dorsal–ventral axis is identified by the expression of *apterous* (33). The dorsal-specific expression of *apterous* is maintained during third larva begins to be expressed in a region slightly anterior to the anterior–posterior boundary (87). The dorsal and ventral surface of the late third instar disc: the perspective ventral surface folds under the dorsal one, while the wing blade unfolds and the wing blade region (X) is placed at the tip of the wing. (After Ref. 54.)



**Figure 3.** Spatial expression and regulatory interaction between *hedgehog* and other genes in the leg imaginal disc. The *hh* protein diffuses across the anterior–posterior compartment boundary; it induces the adjacent dorsal anterior cells to express *wingless* (32, 55). The central region of the disc, where *decapentaplegic* and *wingless* are coexpressed and *aristaless* and defines the distal tip of the leg (68-70). (After Ref. 32.)



Several other components of the hedgehog signal transduction pathway have been identified. Positively-acting components of the pathway include the seven transmembrane protein smoothed

(41), the kinase fused (42), and the transcription factor **cubitus interruptus** (43); negatively-acting components of the pathway include the transmembrane protein patched (44, 45), **protein kinase A** (46-49), and the **kinesin**-related protein costal2 (50). The current model suggests that the patched and smoothed proteins constitute the receptor for the hedgehog signal (51). The cubitus interruptus protein, which is expressed only in the anterior cells due to the repression by the engrailed protein in the posterior compartment, forms a protein complex with the fused and costal2 proteins (50, 52). In some anterior cells, the cubitus interruptus protein is cleaved to generate a truncated form that translocates to the nucleus and represses *dpp* transcription. The hedgehog protein inhibits the proteolysis of the cubitus interruptus protein, leading to the expression of *dpp* in a stripe of anterior cells adjacent to the anterior–posterior boundary (53).

### 3.2. Dorsal–Ventral Patterning

During the second instar larval stage, the wing disc is also divided into dorsal and ventral compartments (23). The dorsal–ventral compartmental boundary lies at the future wing margin, which separates the upper surface of the wing from the lower one. Wing formation depends on the interactions between the dorsal and ventral cells (54). The dorsal selector gene *apterous* (*ap*), which encodes a LIM **homeodomain** protein, is expressed and required in the dorsal cells (33, 55-57), and it activates the expression of *fringe* (*fng*) in these cells (58). The fringe protein acts through two ligands of the **Notch** receptor, *Serrate* and *Delta* (59-63), to restrict the expression of downstream genes *wg* and *vestigial* (*vg*) along the dorsal–ventral boundary, which in turn coordinate wing growth and patterning (64-66).

### 3.3. Proximal–Distal Patterning

Besides the anterior–posterior and dorsal–ventral axes, appendages (wing and leg) also possess a proximal–distal axis. It appears that interactions between the anterior–posterior and dorsal–ventral axes of the disc are critical for patterning along the proximal–distal axis (67). In the leg, the hedgehog protein secreted from the posterior compartment induces the adjacent dorsal anterior cells to express *dpp* and ventral anterior cells to express *wg*. The region where *wg* and *dpp* are coexpressed expresses *Distal-less* and *aristalless*, which defines the most distal part of the leg (68-71). Accordingly, a secondary proximal–distal axis can be induced when a new intersection point of *wg* and *dpp* expression is generated in the leg disc. Similarly in the wing disc, the intersection of *dpp* expression along the anterior–posterior boundary and *wg*–*vg* expression along the dorsal–ventral boundary corresponds to the distal tip of the wing.

## 4. Specification of the Imaginal Disc Identity

Although patterning of all imaginal discs share some common features, the identity of an individual imaginal disc is specified by the **homeotic** and other selector genes. Homeotic genes in the **Bithorax** and **Antennapedia** complexes are involved in the specification of head and thoracic discs (72). For example, *Ultrabithorax* (*Ubx*) is expressed in the haltere, but not in the wing disc, and specifies the haltere fate (73-75). Loss of *Ubx* function in the haltere results in a haltere-to-wing transformation, generating the famous four-winged fly (76). The homeotic gene *Antennapedia* (*Antp*) specifies the leg disc, and ectopic expression of *Antp* in the eye–antenna disc results in antenna-to-leg transformation (77, 78). However, homeotic genes are not the only genes that are used to specify disc identity. The *Drosophila Pax* family gene *eyeless* (*ey*), which encodes a protein containing a paired domain and a homeodomain, is required for the determination of the eye disc fate (79). Ectopic expression of *ey* can induce ectopic eye development in other imaginal discs (80). Several genes have been shown to act downstream of *ey* in this process, including *eyes absent*, *sine oculis*, and *dachshund* (81, 82).

Studies of imaginal disc development in *Drosophila* have shed new light on various aspects of vertebrate development. For example, gene expression and functional studies suggest a striking similarity between fly appendages and vertebrate limbs (32, 67, 83, 84). The dorsal–ventral boundary of the wing disc is analogous to the apical ectoderm ridge (AER), which is a major signaling center for the developing vertebrate limb. The posterior compartment of the wing disc is

analogous to the zone of polarizing activity (ZPA), the source of the morphogen, [sonic hedgehog](#), that patterns the anterior–posterior axis of the vertebrate limb bud. In addition, several key regulatory genes in patterning imaginal discs have been implicated in human diseases. For example, mutations in the human *patched* gene are responsible for basal cell carcinoma, a common form of human cancer ([85](#)).

## Bibliography

1. S. M. Cohen (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor NY, pp. 747–841.
2. P. J. Bryant (1978) In *The Genetics and Biology of Drosophila* (M. Ashburner and T. R. F. Wright, eds.), Academic Press, London, pp. 229–335.
3. H. Oberlander (1985) In *Comprehensive Insect Physiology, Biochemistry and Pharmacology*, Vol. 2 (G. A. Kerkut and L. I. Gilbert, eds.), Pergamon Press, Oxford, pp. 151–182.
4. J. R. S. Whittle (1990) In *Seminars in Cell Biology*, Vol. 1, (P. W. Ingham, ed.), Saunders, Philadelphia, pp. 241–252.
5. M. M. Madhavan and H. A. Schneiderman (1977) *Wilhem Roux's Arch.* **183**, 269–305.
6. J. W. Fristrom (1972) In *The Biology of Imaginal Discs* (H. Ursprung and R. Nothiger, eds.), Springer-Verlag, Berlin, pp. 109–154.
7. J. Hooper and M. P. Scott (1992) In *Early Embryonic Development of Animals*, Vol. 18, (W. Hennig, ed.), Springer-Verlag, Berlin, pp. 1–48.
8. A. Martinez Arias (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 517–608.
9. R. Nusse and H. E. Varmus (1992) *Cell* **69**, 1073–1087.
10. S. M. Cohen (1990) *Nature* **343**, 173–177.
11. B. Cohen, A. A. Simcox, and S. M. Cohen (1993) *Development* **117**, 597–608.
12. D. M. Kingsley (1994) *Genes Dev.* **8**, 133–146.
13. E. L. Ferguson and K. V. Anderson (1992) *Cell* **71**, 451–461.
14. D. Morisato and K. V. Anderson (1995) *Annu. Rev. Genet.* **29**, 371–399.
15. E. Raz and B. Z. Shilo (1993) *Genes Dev.* **7**, 1937–1948.
16. N. Perrimon and L. A. Perkins (1997) *Cell* **89**, 13–16.
17. S. Goto and S. Hayashi (1997) *Development* **124**, 125–132.
18. E. Gateff (1994) *Int. J. Dev. Biol.* **38**, 565–590.
19. E. Gateff et al. (1996) *Int. J. Dev. Biol.* **40**, 149–156.
20. K. L. Watson, R. W. Justice, and P. J. Bryant (1994) *J. Cell Sci. Suppl.* **18**, 19–33.
21. D. F. Woods and P. J. Bryant (1991) *Cell* **66**, 451–464.
22. D. F. Woods and P. J. Bryant (1989) *Dev. Biol.* **134**, 222–235.
23. A. Garcia-Bellido, P. Ripoll, and G. Morata (1973) *Nature New Biol.* **245**, 251–253.
24. P. A. Lawrence and G. Morata (1976) *Dev. Biol.* **50**, 321–337.
25. A. Martinez-Arias and P. A. Lawrence (1985) *Nature* **313**, 639–642.
26. T. Kornberg, I. Siden, P. O'Farrell, and M. Simon (1985) *Cell* **40**, 45–53.
27. J. J. Lee, D. P. von Kessler, S. Parks, and P. A. Beachy (1992) *Cell* **71**, 33–50.
28. T. Tabata, S. Eaton, and T. B. Kornberg (1992) *Genes Dev.* **6**, 2635–2645.
29. T. Tabata and T. B. Kornberg (1994) *Cell* **76**, 89–102.
30. J. Capdevila and I. Guerrero (1994) *EMBO J* **13**, 4459–4468.
31. D. Nellen, R. Burke, G. Struhl, and K. Basler (1996) *Cell* **85**, 357–368.
32. N. Perrimon (1995) *Cell* **80**, 517–520.

33. B. Cohen, M. E. McGuffin, C. Pfeifle, D. Segal, and S. M. Cohen (1992) *Genes Dev.* **6**, 715–729.
34. W. J. Brook and S. M. Cohen (1996) *Science* **273**, 1373–1377.
35. J. Jiang and G. Struhl (1996) *Cell* **86**, 401–409.
36. E. H. Chen and B. S. Baker (1997) *Development* **124**, 205–218.
37. F. Casares, L. Sánchez, I. Guerrero, and E. Sánchez-Herrero (1997) *Dev. Genes Evol.* 216–228.
38. T. Wolff and D. F. Ready (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1277–1326.
39. U. Heberlein and K. Moses (1995) *Cell* **81**, 987–990.
40. J. E. Treisman and U. Heberlein (1998) *Curr. Top. Dev. Biol.* **39**, 119–158.
41. J. Alcedo, M. Ayzenzon, T. Von Ohlen, M. Noll, and J. E. Hooper (1996) *Cell* **86**, 221–232.
42. T. Prémat et al. (1990) *Nature* **347**, 87–89.
43. S. Eaton and T. B. Kornberg (1990) *Genes Dev.* **4**, 1068–1077.
44. J. E. Hooper and M. P. Scott, *Cell* **59**, 751–765 (1989).
45. R. G. Phillips, I. J. Roberts, P. W. Ingham, and J. R. Whittle (1990) *Development* **110**, 105–114.
46. D. Pan and G. M. Rubin (1995) *Cell* **80**, 543–552.
47. J. Jiang and G. Struhl (1995) *Cell* **80**, 563–572.
48. W. Li, J. T. Ohlmeyer, M. E. Lane, and D. Kalderon (1995) *Cell* **80**, 553–562.
49. T. Lepage, S. M. Cohen, F. J. Diaz-Benjumea, and S. M. Parkhurst (1995) *Nature* **373**, 711–715.
50. J. C. Sisson, K. S. Ho, K. Suyama, and M. P. Scott (1997) *Cell* **90**, 235–245.
51. D. M. Stone et al. (1996) *Nature* **384**, 129–134.
52. D. J. Robbins et al. (1997) *Cell* **90**, 225–234.
53. P. Aza-Blanc, F. A. Ramirez-Weber, M. P. Laget, C. Schwartz, and T. B. Kornberg (1997) *Cell* **89**, 1043–1053.
54. S. S. Blair (1995) *Bioessays* **17**, 299–309.
55. J. A. Williams, S. W. Paddock, and S. B. Carroll (1993) *Development* **117**, 571–584.
56. F. J. Diaz-Benjumea and S. M. Cohen (1993) *Cell* **75**, 741–752.
57. S. S. Blair, D. L. Brower, J. B. Thomas, and M. Zavortink (1994) *Development* **120**, 1805–1815.
58. K. D. Irvine and E. Wieschaus (1994) *Cell* **79**, 595–606.
59. S. Artavanis-Tsakonas, K. Matsuno, and M. E. Fortini (1995) *Science* **268**, 225–232.
60. J. F. de Celis, A. Garcia-Bellido, and S. J. Bray (1996) *Development* **122**, 359–369.
61. F. J. Diaz-Benjumea and S. M. Cohen (1995) *Development* **121**, 4215–4225.
62. D. Doherty, G. Feger, S. Younger-Shepherd, L. Y. Jan, and Y. N. Jan (1996) *Genes Dev.* **10**, 421–434.
63. R. J. Fleming, Y. Gu, and N. A. Hukriede (1997) *Development* **124**, 2973–2981.
64. J. Kim, K. D. Irvine, and S. B. Carroll (1995) *Cell* **82**, 795–802.
65. V. M. Panin, V. Papayannopoulos, R. Wilson, and K. D. Irvine (1997) *Nature* **387**, 908–912.
66. J. Kim et al. (1996) *Nature* **382**, 133–138.
67. W. J. Brook, F. J. Diaz-Benjumea, and S. M. Cohen (1996) *Annu. Rev. Cell Dev. Biol.* **12**, 161–180.
68. G. Campbell, T. Weaver, and A. Tomlinson (1993) *Cell* **74**, 1113–1123.
69. K. Basler and G. Struhl (1994) *Nature* **368**, 208–214.
70. F. J. Diaz-Benjumea, B. Cohen, and S. M. Cohen (1994) *Nature* **372**, 175–179.



71. T. Lecuit and S. M. Cohen (1997) *Nature* **388**, 139–145.
72. A. Garcia-Bellido (1977) *Am. Zool.* **17**, 613–629.
73. P. A. Beachy, S. L. Helfand, and D. S. Hogness (1985) *Nature* **313**, 545–551.
74. G. Struhl (1982) *Proc. Natl. Acad. Sci. USA* **79**, 7380–7384.
75. S. D. Weatherbee, G. Halder, J. Kim, A. Hudson, and S. Carroll (1998) *Genes Dev.* **12**, 1474–1482.
76. E. B. Lewis (1978) *Nature* **276**, 565–570.
77. S. Schneuwly, R. Klemenz, and W. J. Gehring (1987) *Nature* **325**, 816–818.
78. G. Gibson and W. J. Gehring (1988) *Development* **102**, 657–675.
79. R. Quiring, U. Walldorf, U. Kloter, and W. J. Gehring (1994) *Science* **265**, 785–789.
80. G. Halder, P. Callaerts, and W. J. Gehring (1995) *Science* **267**, 1788–1792.
81. F. Pignoni et al. (1997) *Cell* **91**, 881–891.
82. R. Chen, M. Amoui, Z. Zhang, and G. Mardon (1997) *Cell* **91**, 893–903.
83. S. S. Blair (1997) *Curr. Biol.* **7**, R686–R390.
84. K. D. Irvine and T. F. Vogt (1997) *Curr. Opin. Cell Biol.* **9**, 867–876.
85. R. L. Johnson et al. (1996) *Science* **272**, 1668–71.
86. D. Fristrom and J. W. Fristrom (1993) In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 843–897.
87. S. S. Blair (1992) *Development* **115**, 21–33.

### Suggestions for Further Reading

88. S. M. Cohen (1993) "Imaginal disc development." In *The Development of Drosophila melanogaster* (M. Bate and A. Martinez Arias, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 747–841.
89. P. A. Lawrence and G. Struhl (1996) Morphogens, compartments, and pattern: lessons from *Drosophila*? *Cell* **85**, 951–961.
90. W. J. Brook, F. J. Diaz-Benjumea, and S. M. Cohen (1996) Organizing spatial pattern in limb development. *Annu. Rev. Cell Dev. Biol.* **12**, 161–180.
91. K. D. Irvine and T. F. Vogt (1997) Dorsal–ventral signaling in limb development. *Curr. Opin. Cell Biol.* **9**, 867–876.
92. J. E. Treisman and U. Heberlein (1998) Eye development in *Drosophila*: formation of the eye field and control of differentiation. *Curr. Top. Dev. Biol.* **39**, 119–158.

### Immortalization

Cell lines isolated from normal tissue continue to proliferate for a fixed number of cell population doublings, after which the culture undergoes [senescence](#); that is, it will undergo a culture age-dependent cessation of growth ([1](#)). Such cell lines are called *finite cell lines*, in contrast to *continuous cell lines* (see [Cell Line](#)). Several potential explanations have been proposed for their limited lifespan *in vitro*: nutritional deficiency, cumulative genetic damage, delayed interaction with an infectious agent, or genetically programmed aging. Although there are, undoubtedly, cell lines that die out soon after isolation, due to nutritional deficiency or infectious agents such as [mycoplasma](#), there is now

overwhelming evidence that the replicative lifespan is programmed into cells (2-4). This has been established by **somatic cell** hybridization and microcell fusion experiments, which have shown that there are at least ten dominant **genes** that will re-establish a finite lifespan in continuous cell lines after fusion with cells from a finite cell line (5, 6). The identity of these genes is not clear, but their interaction with **cell cycle** regulators has been suggested (7). As finite cells show progressive shortening of their **telomeres** with successive generations in culture, until replication is no longer possible, regulation of **telomerase** (8-10) or an alternative mechanism for maintenance of telomere length (11), is now the most likely regulator of replicative lifespan. The vast majority of continuous cell lines from tumors (11, 12) re-express telomerase, which is capable of regenerating the truncated telomere at each cycle, thereby preventing cessation of cell division. Reintroduction of the telomerase gene extends the lifespan (13, 14) and may, ultimately, prove to be the most satisfactory means of immortalizing cells, particularly if under the control of a regulatable promoter.

Agents used to immortalize cell lines, such as irradiation, mutagens, **oncogenes**, or viral genes such as **SV40 virus** large **T Antigen**, **adenovirus** E1a, or HPV E7, interact with cell-cycle regulatory proteins, such as **p53** and **retinoblastoma** proteins and the **cyclin**/Cdk complexes, either directly or via mutation (15). Mutations may be induced directly by irradiation; alternatively, interactions with p53 and other genome guardian proteins may allow spontaneous mutations to be propagated and not repaired or eliminated from the cell population. While these mutational events may lead directly to inactivation of senescence genes, it is also possible that reduced cell cycle control and DNA surveillance allows cells to proceed beyond their normal lifespan, allowing time for small minority populations of cells expressing telomerase to appear, particularly if one of the mutations was in a telomerase **repressor** (16).

Senescence was first described with human diploid fibroblasts but has been shown to occur for other cell types (10) and other species (17). There is a general relationship between the proliferative lifespan *in vitro*, the age of the donor individual, and the average lifespan of the species. Cell lines derived from embryos and young individuals will generally survive longer than cell lines from older donors; moreover, cell lines from long-lived species generally survive longer than those from short-lived species, although this correlation is less well established than that with donor age. Not all cell lines will senesce; notable exceptions are cell lines derived from mice or from tumors of many species. In these cases, evidence of senescence can still be seen, but transformed cells appear within the population and eventually overgrow, producing continuous cell lines. Such immortalized cell lines have been generated spontaneously from many species, such as Syrian (18) and Chinese hamsters (19), cows (20), and monkeys (21, 22), but rarely from normal human or chick cells. Because spontaneous immortalization is such a rare event in normal human cell lines, when it does occur, it begs the question as to whether the cell lines from which immortal lines do arise are really normal, or whether they were already genetically predisposed by one or more initiating mutations in a positively acting **oncogene** by or deletion or inactivation of a negatively acting **tumor suppressor gene**.

*Transformation* is a term with many different meanings, depending on the context in which it is used (see **Neoplastic Transformation**). Transformation of cultured cells implies a heritable **phenotypic** change resulting from a spontaneous or experimentally induced genomic alteration. It is usually associated with an increased life-span and alterations in growth control, and it may, but does not always, give rise to a tumorigenic cell line. There are at least three discernible processes: immortalization, development of aberrant growth control, and neoplastic transformation. A cell line may become immortal without becoming neoplastic or losing many aspects of normal growth control, such as **contact inhibition**, density limitation of cell proliferation, or anchorage dependence. Examples of this are NIH-3T3 primitive mouse embryo mesodermal cells and BHK21 baby hamster kidney fibroblasts. However, both of these may progress to fully transformed lines: 3T3 cells spontaneously (23) (by maintenance at high cell density) and BHK21 by infection with a transforming virus, such as **polyoma** (10). It is not clear whether immortalization is a prerequisite for neoplastic transformation or aberrant growth control, although life-threatening progressive neoplasia almost certainly requires immortalization of one or more stem lines of the tumor.

Immortalization can be defined as the acquisition of an infinite lifespan, usually taken as in excess of 100 population doublings, while, in this context, transformation implies loss of growth control mechanisms, such as growth factor or [serum dependence](#), G1 cell cycle arrest, contact inhibition of cell motility, density limitation of cell proliferation, dependence of cell proliferation on anchorage to, and spreading on, a substrate, and increased production of secreted [proteinases](#) such as [plasminogen](#) activator. Neoplastic transformation implies that the cells will grow as an invasive tumor *in vivo* and usually incorporates all, or most of, the growth control aberrations listed above. It is possible, therefore, that there are only two stages, immortalization and transformation, which need not have a fixed temporal relationship to one another.

Immortalization can also be induced by agents that interfere with growth regulatory genes, or senescence genes ([5](#), [6](#)). These may be physicochemical agents, such as high energy irradiation ([24](#)), oncogenes ([25](#)) or viral genes, such as SV40 large T antigen ([26](#)), E1a from adenovirus, E6 and E7 from papilloma viruses ([27](#)), and [Epstein–Barr virus](#) genes ([28](#)). Viral immortalization can be achieved by [transfection](#) of the appropriate viral genes or infection with whole virus. The second has biohazard safety implications and requires, for safe handling, demonstration that the immortalized cells do not produce infectious virus. For this reason, transfection of cloned genes, most commonly SV40 large T antigen, is generally preferred. Typically, transfected cultures arrest at the M1 stage of the cell cycle, but a few morphologically altered clones will continue to proliferate. These can be picked or selected as resistant if cotransfected with a selectable marker, such as resistance to G418 or hygromycin B ([29](#)). These clones are usually pooled and subsequently undergo crisis (M2), a cessation of growth, from which a small fraction of cells may grow through to form a continuous cell line. The frequency of occurrence of immortalization within a cell population may be as low as  $1 \times 10^9$ , or as high as  $1 \times 10^5$ .

Continuous cell lines have many advantages for routine culture (Table [1](#)), mostly due to their higher growth rate and saturation density, their ability to proliferate in suspension, and their high plating efficiencies. These properties are associated with transformation, however, and continuous cell lines, which are immortalized but not transformed, will not give such high yields or plating efficiencies. As some interference with normal growth control mechanisms is implied in both immortalized and transformed cells, no continuous cell line can be said to be normal. However, cell lines that are immortalized, but not neoplastically transformed, may be said to be closer to the normal phenotype. In both cases, regrettably, there is a tendency to lose lineage markers and the capacity to differentiate ([30](#), [31](#)). However, there are a number of cases where differentiation and lineage markers are retained ([24](#), [32](#)) and there are good prospects for normal **gene expression**, post-transcriptional and [post-translational modification](#), and phenotypic expression. Such cell lines may be useful for the production of proteins by [protein engineering](#), such as factor VIII ([33](#)), or as vectors for further transfection studies. The advent of telomerase-induced immortalization, particularly if the gene is introduced with a promoter that is temperature- or hormone-sensitive, may prove to be the best way forward for generating immortalized cell lines capable of both indefinite replication and expression of differentiation under appropriate conditions.

**Table 1. Properties of Transformed Cells<sup>a</sup>**

| <b>Growth Characteristics</b> | <b>Genetic Properties</b>      | <b>Structural Alterations</b> | <b>Neoplastic Properties</b> |
|-------------------------------|--------------------------------|-------------------------------|------------------------------|
| Immortal                      | High spontaneous mutation rate | Modified actin cytoskeleton   | Tumorigenic                  |

|   |                                       |  |  |
|---|---------------------------------------|--|--|
| Anchorage independent: clone in agar, may grow in stirred suspension                                      | Aneuploid                             | Loss of cell-surface-associated fibronectin                          | Angiogenic   |
| Loss of contact inhibition  | Heteroploid                           | Increased lectin agglutination.                                      | Enhanced proteinase secretion, eg, plasminogen activator |
| Growth on confluent monolayers of homologous cells "focus" formation                                      | Overexpressed or mutated oncogenes    | Modified extracellular matrix  | Invasive   |
| Reduced density limitation of growth: high saturation density, high growth fraction at saturation density | Deleted or mutated suppressor genes   | Altered expression of cell adhesion molecules (cadherins, integrins) |  |
| Low serum requirement   | Stable or elongated telomeres         | Disruption in cell polarity  |  |
| Growth-factor-independent   | Overexpressed telomerase or ALT genes |  |  |
| High plating efficiency   |                                       |  |  |
| Shorter population doubling time  |                                       |  |  |

<sup>a</sup> Modified from Freshney: *Culture of Animal Cells. A Multimedia Guide*, 1999 Wiley-Liss, New York.

## Bibliography

1. L. Hayflick and P. S. Moorhead (1961) *Exp. Cell Res.* **25**, 585–621.
2. S. Goldstein, S. Murano, H. Benes, E. J. Moerman, R. A. Jones, and R. Thweatt (1989) *Exp. Geront.* **24**, (5–6) 461–468.
3. S. E. Holt, J. W. Shay, and W. E. Wright (1996) *Nature Biotechnol.* **14**, (July) 836–839.
4. M. Sasaki, T. Honda, H. Yamada, N. Wake, and J. C. Barrett (1994) *Cancer Res.* **54**, 6090–6093.
5. O. Pereira-Smith and J. Smith (1988) *Proc. Natl. Acad. Sci. USA* **85**, 6042–6046.
6. T. Ogata, D. Ayusawa, M. Namba, E. Takahashi, M. Oshimura, and M. Oishi (1993) *Mol. Cell. Biol.* **13**, 6063–6043.
7. G. H. Stein, L. F. Drullinger, R. S. Robetorye, O. M. Pereira-Smith, and J. R. Smith (1991) *Proc. Natl. Acad. Sci. USA* **88**, 11012–11016.
8. J. W. Shay, H. Werbin, and W. E. Wright (1994) *Mol. Cell. Diff.* **2**, 1–21.
9. C. W. Greider and E. H. Blackburn (1996) *Sci. Am. (Feb.)* **274**, 92–97.
10. P. Konig and D. Rhodes (1997) *Trends Biochem. Sci. (Feb.)* **22**, 43–47.

11. T. M. Bryan and R. R. Reddel (1997) *Eur. J. Cancer* **33**, 767–773.
12. S. Bachetti and C. M. Counter (1995) *Intl. J. Oncol.* **7**, 423–432.
13. A. G. Bodnar, M. Ouellette, M. Frolkis, S. E. Holt, C.-P. Chiu, G. B. Morin, C. B. Harley, J. W. Shay, S. Lichtsteiner, and W. E. Wright (1998) *Science* **279**, 349–352.
14. H. Vaziri and S. Benchimol (1998) *Curr. Biol.* **8**, 279–282.
15. E. K. Parkinson (1996) in *Culture of Immortalized Cells*, R. I. Freshney, and M. G. Freshney eds., Wiley-Liss, New York, pp. 1–23.
16. M. Oshimura and J. C. Barrett (1997) *Eur. J. Cancer* **33**, 710–715.
17. R. J. Hay and B. L. Strehler (1967) *Exp. Gerontol.* **2**, 123.
18. I. Macpherson, and M. Stoker (1962) *Virology* **16**, 147.
19. T. T. Puck, S. J. Cieciura, and A. Robinson (1958) *J. Exp. Med.* **108**, 945–956.
20. P. Del Vecchio and J. R. Smith (1981) *J. Cell. Physiol.* **108**, 337–345.
21. H. Hopps, B. C. Bernheim, A. Nisalak, J. H. Tjio, J. E. Smadel (1963) *J. Immunol.* **91**, 416–424.
22. Y. Gluzman (1981) *Cell* **23**, 175–182.
23. D. Brouty-Boyé, R. W. Tucker, J. Folkman (1980) *Intl. J. Cancer*, **26** 501–507.
24. N. A. Punchard, D. Watson, R. Thompson, and M. Shaw (1996) in *Culture of Immortalized Cells*, R. I. Freshney, M. G. Freshney eds., Wiley-Liss, New York, pp. 203–237.
25. T. L. Morgan, D. Yang, D. G. Fry, P. J. Hurlin, S. K. Kohler, V. M. Maher, and J. J. McCormick (1991) *Exp. Cell Res.* **197**, 125–136.
26. L. V. Mayne, T. N. C. Price, K. Moorwood, J. F. Burke (1996) in *Culture of Immortalized Cells*, R. I. Freshney, M. G. Freshney eds., Wiley-Liss, New York, pp. 77–93.
27. G. A. Seigel (1996) *In Vitro Cell Dev. Biol.—Animal* **32**, 66–68.
28. B. J. Bolton and N. K. Spurr (1996) in *Culture of Immortalized Cells*, R. I. Freshney and M. G. Freshney eds., Wiley-Liss, New York, pp. 283–298.
29. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning, a Laboratory Manual*, 2nd ed., **3** vols. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
30. H. C. Isom and I. Georgoff (1984) *Proc. Natl. Acad. Sci. USA* **85**, 9783–9787.
31. D. Wynford-Thomas (1996) in *Culture of Immortalized Cells*, R. I. Freshney, M. G. Freshney eds., Wiley-Liss, New York, pp. 183–201.
32. J. F. Burke, T. N. C. Price, and L. V. Mayne (1996) in *Culture of Immortalized Cells*, R. I. Freshney, M. G. Freshney, eds. Wiley-Liss, New York, pp. 299–313.
33. F. Moldovan, H. Benanni, J. Fiet, O. Cussenot, J. Dumas, C. Darbord, and H. R. Soliman (1996) *In Vitro Cell Dev. Biol.—Animal* **32**, 16–23.

## Immune Response

The immune response is a complex sequence of events initiated by the introduction of an [immunogen](#) in an organism. In vertebrates, it is the consequence of the stimulation of B and/or T lymphocytes, leading either to the production of specific **antibodies** from [B cells](#) or to the occurrence of effector **T cells**—for example, [cytotoxic T lymphocytes](#). Immunologists classically make a clear distinction between immune responses that are characterized by antibody production and those that result primarily in the emergence of cytotoxic T cells or delayed-type hypersensitivity, which are depicted as humoral and cellular immunity, respectively. In fact, penetration of complex

**antigens**, such as most bacteria and **viruses**, trigger both types of immune responses. Furthermore, antibody responses directed against T-dependent antigens, the most common case, require cellular cooperation between B and T cells, and the humoral and cellular responses are quite intricate.

It must be realized that both the nature and the intensity of the immune response is largely dependent on the nature, dose, and mode of penetration of the antigen. For example, when given subcutaneously or intravenously, proteins stimulate a good circulating antibody response, provided that they are given in an immunogenic range of doses. If given orally, mucosal antibody production will occur preferentially. By contrast, intradermal injection will induce delayed hypersensitivity, mediated by T lymphocytes. The level of response may be substantially increased by admixing antigen with [adjuvants](#) that ensure a long-lasting liberation of antigen, in addition to providing nonspecific stimulation of the immune system.

Whatever the nature of the response, it always follows basic principles that are most simply quantified in a classical humoral response. Following antigen administration, two parameters may be followed: One is the disappearance of antigen from the body fluid, and the other is the occurrence of circulating antibodies. After a first immunization, the blood concentration of antigen decreases rapidly within a few hours, remains somewhat steady for a few days, and then rapidly decreases again, until it is no longer detectable. This occurs after 5 to 7 days, when antibodies start to be produced. As they circulate in the bloodstream, they bind to the remaining antigens, and the immune complexes are eliminated (which may be a major problem for individuals suffering renal failure). The antibodies that first appear are [IgM](#). Their yield increases for a few days, until they are replaced by [IgG](#) antibodies that result from [class switching](#) of the [isotype](#). The antibody titer slowly diminishes and becomes barely detectable. This is a typical primary response. As B cells start making IgG antibodies in germinal centers of peripheral lymph nodes, they rapidly accumulate [somatic hypermutations](#) so that antibodies with higher affinity will occur and be positively selected by antigen. A fraction of these B cells will constitute the long-lived [memory cells](#) that will react immediately upon a second administration of antigen. The secondary response then takes place, characterized by a rapid elevation of the IgG antibody titer, which will remain high for a much longer period of time. If a different antigen were used instead of the first one, a typical primary response would have been observed, and this is illustrative of the high specificity of the immune response. Antigen may be given several times, ensuring a successive boost of responses, leading to a hyperimmunized state.

The general kinetics described above are typical of T-dependent soluble antigens. With polysaccharides, there is no secondary response, and the antibodies remain essentially of the IgM type (with some IgG2a and IgA2 in humans). The immune response to particulate antigens, such as bacteria, is generally more complex, due to antigenic diversity. The respective contribution of humoral and cellular responses is largely dependent on the microorganism.

See also entries [Antibody](#), [Antigen](#), [Immunogen](#), and [Immunization](#).

#### Suggestions for Further Reading

P. F. Searle and L. S. Young (1996) Immunotherapy II: antigens, receptors and costimulation. *Cancer Metastasis Rev.* **15**, 329–349.

J. Mestecky, S. M. Michalek, Z. Moldoveanu, and M. W. Russell (1997) Routes of immunization and antigen delivery systems for optimal mucosal immune responses in humans. *Behring Inst. Mitt.* **98**, 33–43.

Immunity may be defined as a resistant state conferred on a living organism by a wide variety of cells and molecules that tend to eliminate pathogens or aggressive agents, thus ensuring the individual's integrity. Immunologists often make a distinction between natural, or innate, and adaptive immunity, although the latter is also quite “natural” when understood as a basic physiological phenomenon. It might be advisable to use instead the concept of nonspecific versus specific, but even this is not entirely satisfactory, because the limits are not so clear. This is well illustrated by the opposed views that were expressed at the end of the nineteenth century by those, like Metchnikoff, who supported the cellular basis for immunity, centered on the phagocytic properties of the [macrophage](#) and those, like Ehrlich, who gave a major role to circulating **antibodies**, thus supporting humoral immunity. In fact, both mechanisms are part of the organism defenses; and the macrophage, which was initially considered to have nonspecific phagocytic potentialities, turned out to be able to acquire some specificity by opsonisation—that is, the binding of antibodies of the various isotypes to its **Fc receptors**. Much more recently, the ability of macrophages to process and present protein antigens proved definitively that all these partners could be linked and that a clear-cut distinction was primarily a scholastic dispute. Defense mechanisms are quite diversified in the living world, although vertebrates have developed a highly complex immune system endowed with an adaptive response that can face a huge potential antigenic [repertoire](#) because of a sophisticated genetic organization that may express at any given time millions of different and specific **immunoglobulins** or **T cell receptors** (TCRs) produced by [B cells](#) and [T cells](#), respectively. Invertebrates have also developed immune defenses against pathogens that involve, in addition to phagocytic cells or soluble [enzymes](#), mechanisms like those described in insects, which have a limited repertoire of genetically encoded peptides endowed with diverse bacteriolytic, antiviral, or antiparasite activities.

One may consider that the defense mechanisms that operate in vertebrates are organized on successive levels. The first one involves cutaneous and mucosal barriers that ensure a physical protection, as well as a chemical one, including enzymes such as [lysozyme](#) that can have a bactericidal effect. The second level is used whenever pathogens have entered the bloodstream. Examples of efficient but nonspecific mechanisms that can be mobilized quickly are (a) phagocytosis by monocytes and macrophages, (b) inflammatory reaction that produces many cytolytic factors, (c) lytic activity initiated by components of the **complement** cascade, (d) natural killer (NK) cells, and (e) various **cytokines**, with special reference to [interferons](#) that block viral replication, and so on. Simultaneously, the adaptive immune response is initiated, leading to the occurrence of the first wave of [IgM](#) antibodies, which are soon replaced by more efficient [IgG](#), and/or to the emergence of [cytotoxic T lymphocytes](#), which are of prime importance in destroying virus-infected cells. Vaccination will reinforce immunity by stimulating a specific immune response to major dangerous pathogens, allowing the organism to react quickly with the appropriate B- and T-cell responses when a natural infection takes place.

See also entries [Antigen](#), [Antibody](#), [Autoimmunity](#), [Immunogen](#), [Immune Response](#).

#### Suggestions for Further Reading

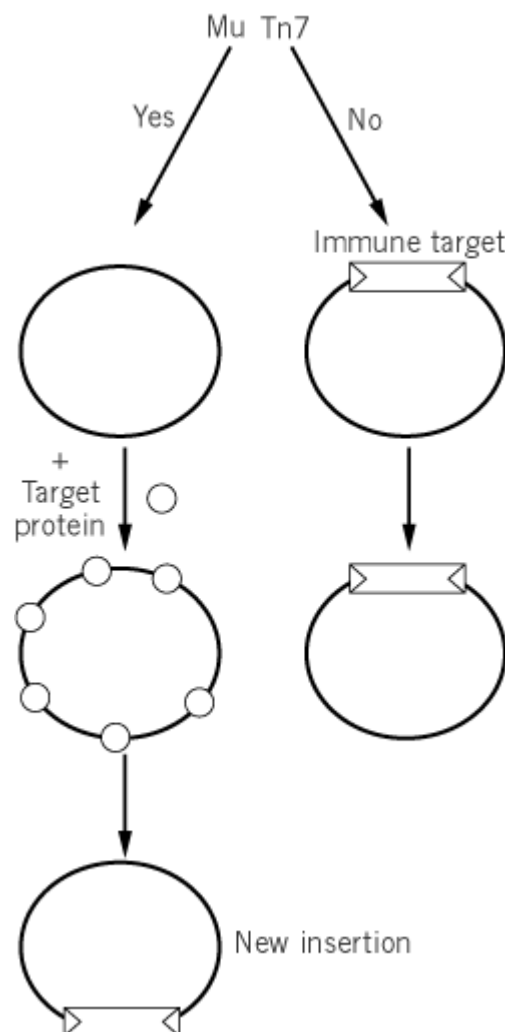
M. C. Carroll and A. P. Prodeus (1998) Linkages of innate and adaptive immunity. *Curr. Opin. Immunol.* **10**, 36–40.

H. Kohler and S. Paul (1998) Superantibody activities: new players in innate and adaptive immune responses. *Immunol. Today* **19**, 221–227.

## Immunity, Transposon

For some [transposable elements](#), the frequency at which they insert into a target DNA that already contains a copy of the element is much reduced from their frequency of insertion into that same DNA lacking a copy of the transposable element (Fig. 1). This inhibition does not represent a global inhibition of [transposition](#) within the cell: The inhibition is specifically [cis-acting](#) on the DNA that contains the transposable element. This inhibition is called *transposition immunity* or *target immunity*. Although target immunity is usually evaluated in intermolecular transposition reactions, it is likely that a principal function of this phenomenon is to discourage intramolecular transposition—that is, to prevent the insertion of a mobile element into itself or into adjacent donor backbone. Target immunity has the effect of discouraging insertions locally and promoting insertion at a distance from the donor site of the [mobile element](#).

**Figure 1.** Target immunity. For some elements such as Mu and Tn7, the frequency with which they insert into a target DNA already containing a copy of the element is much reduced. A target DNA already containing a copy of the transposon is a poor target because the target binding protein (small circle) is not effectively bound to the immune target. Binding of the target protein is discouraged by transposition proteins that bind specifically to the transposon ends in the immune target DNA.





Target immunity has been observed for a number of bacterial elements. The phenomenon was first observed with Tn3-like transposons and subsequently with bacteriophage **Mu** (1), and Tn7 (2, 3), a **drug-resistance** element. Phage Mu uses multiple rounds of transposition to set up [replication forks](#) that promote phage replication during lytic growth. Through transposition immunity, Mu ensures that new insertions made during the replication process will not be into extant copies of Mu, destroying these copies; thus transposition immunity facilitates Mu replication.

The mechanism of immunity has been dissected *in vitro* with Mu (4, 5), and Tn7 (6, 7). Common features of these elements are that the activity of the [transposase](#) is controlled by an ATP-dependent regulatory protein that interacts with the transposase and the target DNA via an ATP-dependent **DNA-binding** activity. Transposition occurs when the transposase interacts with this regulatory protein in an ATP-bound state on the target DNA. The ATP state of the regulatory protein is key to its activity, promoting transposition when it exists in an ATP-bound state on the target DNA and being inactive in a non-ATP state removed from the target DNA. As described below, the state of the ATP regulatory protein appears to be the key to immunity. A reasonable model is that the establishment of immunity involves discouraging the interaction of the ATP-dependent regulatory protein with the target DNA (Fig. 1).

In all systems examined, target immunity is imposed on a potential target DNA by the DNA sequence ends of the transposon—that is, the *cis*-acting [recombination](#) sequences at the ends of the transposon that are the transposase recognition sites. An attractive model for immunity is that when the transposase interacts with transposon end sequences on the potential target DNA, it exists on the target DNA in high local concentration and actively discourages the interaction of the ATP-dependent regulatory protein with the target DNA, thus effectively clearing the target DNA of targeting protein and discouraging transposition into that DNA. The ends of the transposon with bound transposase appear to interact with the target protein by simple [protein–protein interactions](#) facilitated by DNA looping (5, 7).

It seems likely that different forms of the transposase will interact with the target protein to execute transposition and transposition immunity. In Tn7, for example, immunity can be imposed by one of the protein components of the transposase, while actual transposition requires collaboration between both protein subunits of the transposase (7). In other cases, such as Mu, different forms of transposase, for example in different oligomeric states, may signal these different reactions.

Retroviruses also show a form of immunity in that intramolecular insertion is discouraged, although they show no inhibition for insertion into DNAs containing retroviral ends (8). The ability to block intramolecular insertion has been correlated with a requirement for a particular human protein (9), but the role of this protein and the mechanism of retroviral immunity remain to be established.

## Bibliography

1. A. Darzins, N. E. Kent, M. S. Buckwalter, and M. J. Casadaban (1988) Proc. Natl. Acad. Sci. USA **85**, 6826–6830.
2. L. K. Arciszewska, D. Drake, and N. L. Craig (1989) J. Mol. Biol. **207**, 35–52.
3. B. Hauer and J. A. Shapiro (1984) Mol Gen Genet. **194**, 149–158.
4. K. Adzuma and K. Mizuuchi (1988) Cell **53**, 257–266.
5. K. Adzuma and K. Mizuuchi (1989) Cell **57**, 41–47.
6. R. J. Bainton, K. M. Kubo, J.-N. Feng, and N. L. Craig (1993) Cell **72**, 931–943.
7. A. Stellwagen and N. L. Craig (1997) EMBO J. **16**, 6823–6834.
8. M. S. Lee and R. Craigie (1994) Proc. Natl. Acad. Sci. USA **91**, 9823–9827.
9. M. S. Lee and R. Craigie (1998) Proc. Natl. Acad. Sci. USA **95**, 1528–1533.

## Immunization

Immunization occurs as the result of the penetration of an [immunogen](#) in an organism, leading to the stimulation of the immune system and the production of **antibodies** and/or [T-cell](#) effectors that specifically recognize the initiating [antigen](#). Immunization may occur naturally, upon penetration of a pathogen, bacteria, virus, or parasite, and requires several days before becoming fully effective fighting the invader. This may be too slow to ensure efficient protection, and this led to the development of systematic vaccinations against major pathogens. The word *vaccination* originates from the first vaccine that was directed against smallpox, using the cross-reactive [vaccinia virus](#). Used on a worldwide basis, it led to the complete eradication of this severe disease, and thus is no longer in use. Major vaccines used in humans include tetanus, **diphtheria**, measles, **polio**, mumps, rubella, **influenza**, rabies, **hepatitis B**, and a few others. Vaccines are most frequently administered by injection, but sometimes orally, as in the case of the live attenuated poliovirus vaccine. A vaccine may be used as the full pathogen after it has been inactivated, or as an attenuated genetic variant. One also can use a purified molecule that has been isolated from the pathogen or genetically engineered.

For laboratory purposes, immunization generally aims at producing antibodies for immunochemical analysis. Conventional polyspecific antibodies are generally prepared in mice, guinea pigs, rabbits, or goats. There is no standard immunization procedure, but only some general guide lines that define basic conditions for good immunogenicity. Besides using an appropriate dose, the antigen might have to be conjugated to a carrier protein. Use of [adjuvant](#), most frequently the complete Freund's adjuvant, is admixed with the immunogen. Finally, repeated injections will increase both the affinity and the titer of the antibodies generated.

See also entries [Antigen](#), [Antigen Processing, Presentation](#), [Antibody](#), [Immune Response](#), and [Immunogen](#).

### Suggestions for Further Reading

W. M. McDonnell and F. K. Askari (1997) Immunization. *JAMA* **278**, 2000–2007.

E. Manickan, K. L. Karem, and B. T. Rouse (1997) DNA vaccines. A modern gimick or a boon to vaccinology? *Crit. Rev. Immunol.* **17**, 139–154.

## Immunoaffinity Chromatography

Immunoaffinity chromatography (IAC) is a widely used method for the selective extraction of **antibodies** or other biomolecules on the basis of the biospecific affinity between an antibody (Ab) and its [antigen](#) (Ag) (**1**). The antigen may be a small molecule (a [hapten](#)) or a macromolecule (M) such as an [enzyme](#), a [hormone](#), a [growth factor](#), and so on. IAC has been used for many years for the purification of antibodies by immobilized antigens and also for the purification of antigens (enzymes, receptors, hormones, growth factors, etc.) by antibodies anchored to inert carriers. The antibodies used can be polyclonal, provided that a pure antigen is injected to obtain the antibodies and that strict specificity is ascertained. It is often advisable to use anti-peptide (sequence-oriented)

polyclonal antibodies to reduce the danger of heterogeneity in the specificity of the antibody population. [Monoclonal antibodies](#) are usually a very good choice because they are homogeneous and directed against a distinct epitope. Under ideal conditions, antibody columns allow separation of specific peptides or proteins from crude mixtures in one step. The procedure for an IAC purification involves the same steps described for [affinity chromatography](#).

A major difficulty in the use of immobilized antibodies is the high affinity of some antibodies for their antigens. This sometimes makes the recovery of active enzymes difficult, because the harsh conditions that are required in the elution step bring about either an inactivation of the enzyme or the loss of an important regulatory function. Among the eluents that have (and are) being used in IAC are buffers with a low pH (2.2) or a high pH (11.5), **chaotropes** (5 M KSCN), or mild **denaturants** such as [urea](#) (3.5 to 8 M) or **guanidinium** salts (2–6 M). Therefore, attempts are sometimes made to purify by exclusion (also known as “reverse immunoadsorption”). In this procedure, the contaminating proteins are the ones to be adsorbed and removed, while the desired protein is excluded (2).

With the introduction of monoclonal antibodies in IAC, it became possible to use homogeneous immunoglobulins not only directed against a distinct antigenic determinant (3, 4) and to tailor their affinity. A lowered affinity permits milder conditions for elution and reduces the possibility of irreversible denaturation during the elution step. As expected, the use of monoclonal antibodies has had a significant impact in biochemistry, in molecular biology, and in medicine.

#### Bibliography

1. M. A. J. Godfrey (1997) In *Affinity Separations* (P. Matejtschuk, ed.), IRL Press, Oxford University Press, New York, p. 141.
2. J. A. Weare, J. T. Gafford, N. S. Ur, and E. G. Erdos (1982) *Anal. Biochem.* **123**, 310.
3. G. Kohler and C. Milstein (1975) *Nature (London)* **256**, 495.
4. J. W. Goding (1980) *J. Immunol. Methods* **39**, 285.

#### Immunoassays

Immunoassays are a collection of techniques for measuring biological or synthetic materials, using **antibodies** specific for those materials. Immunoassays are also applied in reverse, in which a purified antigen or hapten is used to quantify an uncharacterized antibody preparation. The common element in all immunoassays is the interaction of an antibody with a ligand. The assay methods differ in the accessory chemical reactions that take place prior to, or as a consequence of, ligand binding, and in the instrumentation used to detect the binding event or subsidiary reaction. [Radioimmunoassays](#) and [enzyme-linked immunosorbent assays](#) (ELISA) are especially widely used types of immunoassay. Others are [complement fixation](#), [agglutination](#), [Farr Assay](#), and [Ouchterlony Double Diffusion](#). They are listed as separate entries and not discussed here. Below, we sketch the principles underlying the most common immunoassay methods relevant to molecular biology.

##### 1. Fluorescence Quench

The emission of a photon as a molecule returns to its electronic ground state from an excited state produces **fluorescence**. “Quenching of fluorescence” refers to chemical interactions that allow an excited state to dissipate its energy through a nonradiative mechanism (see [Fluorescence Quenching](#)). Common mechanisms that can cause quenching include collisions with solvent,

protonation or deprotonation reactions, resonance energy transfer, and collisions or static contact with quenchers. Quenchers encompass a huge variety of substances, from molecular oxygen to heme proteins, and the binding of a quencher to an antibody can easily be measured. Quenchers most relevant to immunoassays are aromatic organic molecules, which are frequently used as [haptens](#). These molecules are thought to cause nonradiative deexcitation by transiently accepting electrons from the excited state of a fluorophore.

[Tryptophan](#) is the only strongly fluorescent naturally occurring amino acid, and antibody sequences usually show an abundance of exposed Trp residues in the antigen-combining site (1). Consequently, binding of a quenching antigen or hapten to an antibody nearly always causes decrease of Trp fluorescence from the antibody. Furthermore, the fluorescence decrease is generally linear with the number of combining sites filled; hence it is a direct measure of the degree of saturation of the antibody (2). Bound and free Ab as a function of added ligand can be calculated and displayed as a [Scatchard Plot](#) or used for direct fitting to a binding equation to determine the antibody valency and the association equilibrium constant for binding the ligand. Fluorescence quench on antigen binding is a very desirable technique for immunoassay, because it allows direct measurement of antibodies and antigens interacting in solution and can be performed without chemical labeling of the molecules under study with [reporter groups](#).

## 2. Fluorescence Anisotropy or Fluorescence Polarization

This method relies on the change in molecular tumbling rate that occurs when an antibody–antigen complex forms. Fluorescence anisotropy is most appropriate for antibody-antigen pairs that are of different molecular size (3), and may be applicable even if no overt change in fluorescence signal results from complex formation.

Plane-polarized light absorbed by a fluorophore will (to a first approximation) be emitted with the same plane of polarization, unless the fluorophore molecule reorients during the lifetime of its excited state. Since the tumbling rate is a function of molecular size, polarization of emitted light is very sensitive to complex formation. A small molecule will tumble rapidly, hence the orientation of an excited state fluorophore will randomize, and fluorescence emission will be isotropic. In contrast, a large molecule will tumble slowly; hence emission from a fluorophore will remain polarized. A small molecule that binds to a large molecule will take on the tumbling characteristics of the large molecule. For example, anisotropic emission from a fluorescent tag on a small antigen will sharply increase on complexation with a large antibody. To use this principle for a quantitative immunoassay, fluorescence anisotropy is measured separately for the free antigen and fully complexed antigen. Between these two limiting values, anisotropy will be linear with the extent of antigen complexation; hence the degree of complex formation in a test sample may be inferred from the observed anisotropy.

Fluorescence anisotropy is measured by exciting a fluorophore with polarized light and observing fluorescence intensity through a polarizing filter oriented parallel ( $I_{\parallel}$ ) or perpendicular ( $I_{\perp}$ ) to the excitation plane. The relative intensities from the parallel and perpendicular emissions are used to calculate anisotropy or polarization from the equations:

$$\text{Anisotropy} = \frac{I_{\parallel} - I_{\perp}}{I_{\parallel} + 2I_{\perp}} \quad (1)$$

$$\text{Polarization} = \frac{I_{\parallel} - I_{\perp}}{I_{\parallel} + I_{\perp}} \quad (2)$$

Although polarization and anisotropy are very closely related, equations relating observed intensities to ligand binding and other phenomena are generally simpler when expressed in terms of anisotropy; hence this is the preferred parameter. Fluorophores used for anisotropy measurements of proteins, such as fluorescein and **dansyl**, typically have excited-state lifetimes of 5 ns. This lifetime is on the same order as the tumbling times of protein antigens. The small, globular protein **lysozyme**, for example, also has a rotational correlation time of 5 ns (4); hence a fluorescein label on lysozyme would show considerable depolarization of fluorescence. Intact **IgG** has a rotational correlation time of 220 ns; hence its complex with the labeled lysozyme would show strong fluorescence anisotropy.

One complication to the fluorescence anisotropy technique is the occurrence of local motions affecting the fluorophore, often referred to as “segmental flexibility.” The timescale for these local motions can be sufficiently fast to reorient the fluorophore and depolarize emission, even if the fluorophore is attached to a large, slowly tumbling macromolecule (5).

### 3. Phosphorescence or Time-Resolved Fluorescence

Phosphorescence is a fluorescence phenomenon in which, for quantum-mechanical reasons, the excited state of a fluorophore cannot easily return to the ground state by emission of a photon. Emission does occur, but on a much longer timescale than is normal for fluorescence. Because of this long lifetime, the presence of even small amounts of quenching agents such as  $O_2$  is sufficient to suppress observation of any signal from naturally phosphorescent biological compounds. Phosphorescence is measured by exciting a sample with a pulsed light source, waiting an interval for short-lived fluorescence states to decay, then integrating photon counts during the remaining dark period before the next pulse. This ability to resolve short-lived from long-lived fluorescence is the chief advantage of a phosphorescent probe, as background phosphorescence of biological samples is almost nil. By contrast, the steady-state fluorescence background is often quite high, and the signal-to-background ratio in short-lived fluorescence-based assays is a more common limitation of sensitivity of an assay than the absolute magnitude of the signal.

Europium is the most common phosphorescent probe in time-resolved fluorescence immunoassays (6). This metal ion is attached to antibodies or antigens through bifunctional linker molecules containing a chelate moiety at one end and a group reactive with protein residues at the other. In a typical assay, a trapping antibody is immobilized on a solid support, antigen in a test sample is allowed to bind, and then the  $Eu^{3+}$ -labelled antibody is allowed to bind the trapped antigen. After washing away unbound reagents,  $Eu^{3+}$  is extracted from the antibody in low pH buffer containing phosphorescence enhancing agents (7). These agents act both by sequestering the  $Eu^{3+}$  ion from quenchers and transferring absorbed light to the  $Eu^{3+}$  ion. The fluorescence emission is integrated repetitively over several hundred microseconds, using a stroboscopic fluorimeter. Comparison of the response from the test sample to a standard curve constructed with known concentrations of analyte allows determination of the concentration of analyte in the test sample.

### 4. Chemiluminescence

**Chemiluminescence** is a phenomenon related to fluorescence, in which an electronic excited state of a molecule is reached through a chemical reaction. Emission of a photon, which is optically detected, returns the molecule to the ground state. A chemiluminescent immunoassay involves an antibody or antigen that has been covalently tagged with a molecule capable of participating in a chemiluminescent reaction. The tag remains inert through immunologic steps such as antigen or secondary antibody binding. Subsequent mixing with a trigger solution initiates the luminescent reaction.

A chemiluminescent immunoassay to detect an antigen typically begins with immobilization of a capture antibody on a solid support. The test sample is added, followed by the tagged detection antibody. Luminol was the first luminogenic tag to be used for an immunoassay (8), but derivatized

acridinium esters are now more common (9, 10). Unbound reagents are washed away, and luminescence is determined in a luminometer. This device combines two automated functions (11). First, trigger solutions (alkaline  $H_2O_2$  in the case of acridinium labels) are added to the sample to initiate the luminescence reaction. Mixing is usually complete in less than 1 s. Next, a detector, usually with photon-counting electronics, measures light emission from the sample over a defined time interval, typically 10 s. As with other immunoassays, the response of a test sample is compared to the responses from a set of standards to determine the analyte concentration in the test sample.

Because photon emission from luminescent tags can be induced to occur quantitatively over the span of a minute, as compared to radioactive decay of isotopic isotopes, which continues over months in the case of  $^{125}I$ , chemiluminescent assays are inherently very sensitive. ELISA assays that use a luminogenic substrate, which are sometimes referred to as “chemiluminescent immunoassays,” are perhaps the most sensitive of all assays. Combination of an enzyme tag on the antibody and a luminogenic enzymatic reaction yields a degree of amplification that can detect as little as 1000 molecules of analyte (12).

### Bibliography

1. E. A. Padlan (1990) *Proteins* **7**, 112–124.
2. H. N. Eisen (1964) *Meth. Med. Res.* **10**, 115–121.
3. E. Haber and J. C. Bennett (1962) *Proc. Natl. Acad. Sci. USA* **48**, 1935–1942.
4. B. P. Maliwal and J. R. Lakowicz (1984) *Biophys. Chem.* **19**, 337–344.
5. W. B. Dandliker and G. A. Feiger (1961) *Biochem. Biophys. Res. Commun.* **5**, 299–304.
6. K. Pettersson, H. Siitari, I. Hemmilä, E. Soini, T. Lövgren, V. Hänninen, P. Tanner, and U.-H. Stenman (1983) *Clin. Chem.* **29**, 60–64.
7. I. Hemmilä, S. Dakubu, V.-M. Mukkala, H. Siitari, and T. Lövgren (1984) *Anal. Biochem.* **137**, 335–343.
8. J. S. A. Simpson, A. K. Campbell, M. E. T. Ryall, and J. S. Woodhead (1979) *Nature* **279**, 646–647.
9. I. Weeks, I. Behesgti, F. McCapra, A. K. Campbell, and J. S. Woodhead (1983) *Clin. Chem.* **29**, 1474–1479.
10. I. Weeks, M. Sturgess, R. C. Brown, and J. S. Woodhead (1986) *Meth. Enzymol.* **133**, 366–387.
11. P. E. Stanley (1986) *Meth. Enzymol.* **133**, 587–603.
12. L. J. Kricka (1993) *Clin. Biochem.* **26**, 325–331.

### Suggestions for Further Reading

13. J. R. Lakowicz (1983) *Principles of Fluorescence Spectroscopy*, Plenum Press, New York. (Rigorous but very readable description of fluorescence and other optical techniques and their use in biological applications.)
14. C. P. Price and D. J. Newman, eds. (1997) *Principles and Practice of Immunoassay*, 2nd ed., Macmillan, London. (Comprehensive treatment covering many types of immunoassay.)
15. E. P. Diamandis and T. K. Christopoulos, eds. (1996) *Immunoassay*, Academic Press: New York.
16. R. F. Masseyeff, ed. (1993) *Methods of Immunological Analysis*, VCH, New York. (An ambitious 12-volume series that is truly a comprehensive treatment of immunoassay methods. Volume 1, W. H. Albert and N. A. Staines, eds., *Fundamentals*, contains many treatments that are of general utility regardless of instrumental method.)

## Immunoelectron Microscopy

In molecular biology, the location of specific molecules in organelles and cells at the microscopic level has gained prominence in recent years through the development of a wide range of cytochemical techniques. These techniques are based upon the property that specific markers can be chemically coupled to a wide variety of molecules. Examples are fluorescent markers for light microscopy and electron-dense markers, such as [ferritin](#) or colloidal gold, for [electron microscopy](#). The use of poly- and [monoclonal antibodies](#) has widened the possibility of localization to essentially every particle that can invoke an immunogenic response. The antibodies produced can be bound to markers used in microscopy and the distribution of the molecules of interest visualized directly by 1-, 2-, or 3-step labeling procedures ([1](#)).

Immunocytochemistry is the application of antibodies in labeling protocols and must be performed in an aqueous environment to ensure that the [antigen](#) detectability is accurate and reproducible. Hence, cryotechniques are much preferred over the classical technique of chemical fixation, dehydration, and plastic embedment, which usually results in a loss of antigenicity. Along with cryofixation and cryosectioning, the use of colloidal gold is now commonplace. Colloidal gold has numerous advantages, among which are (i) it is electron dense and easily visualized in electron microscope images; (ii) various sizes of gold particles can be obtained, allowing double-labeling experiments; (iii) it can be easily complexed with a wide variety of ligands; and (iv) the maintenance of bioactivity of complexed molecules is excellent. Typically, the primary antibody binds to the antigen of interest and is visualized by labeling with a secondary antibody to which colloidal gold is conjugated in analogy to the use of fluorescent antibodies in immunofluorescence. With immunogold labeling, certain considerations should be addressed before starting ([2](#)). (i) Where is the antigen located? Additional preparative steps may be necessary to make the antigen accessible to the gold antibody. (ii) What is the nature of the antibody? A monoclonal antibody may not bind the antigenic site if the site is modified during electron microscope preparation; this problem is generally much less severe with cryotechniques. Polyclonal antiserum, which contains a number of different antibodies directed against more than one epitope on the antigen, is less susceptible to conformational changes or modification of the antigen. (iii) From which species were the antibodies derived? Care should be exercised to ensure that antibodies are raised to pure undenatured antigens in a species compatible with the labeling system.

A recent advance in photooxidation has allowed for precise correlative immunolocalization on the same specimen using fluorescent light and immunoelectron microscopies ([3](#)). With use of this technique, proteins can first be localized to individual cells or tissue types at the light microscope level with an antibody containing a suitable fluorescent label such as eosin. Subsequently, the protein can be further localized at the electron microscope level by illuminating the fluorophore in the presence of oxygen and diaminobenzidine. Fluorophores with a relatively poor quantum yield, such as eosin, produce reactive oxygen intermediates that oxidize the diaminobenzidine, producing a brown osmophilic [polymer](#). Preparation of these specimens for electron microscopy by osmium fixation yields a highly localized osmium stain at the antigen. Fluorescence photooxidation can also be used with nucleic acid sequences by using biotinylated probes and an eosin–streptavidin conjugate (see [Streptomycin](#)).

### Bibliography

1. M. J. Dykstra (1993) *A Manual of Applied Techniques for Biological Electron Microscopy*, Plenum Press, New York.
2. J. R. Harris (1991) *Electron Microscopy in Biology*, IRL Press, Oxford.
3. T. J. Deerinck et al. (1994) *J. Cell Biol.* **126**, 901–910.

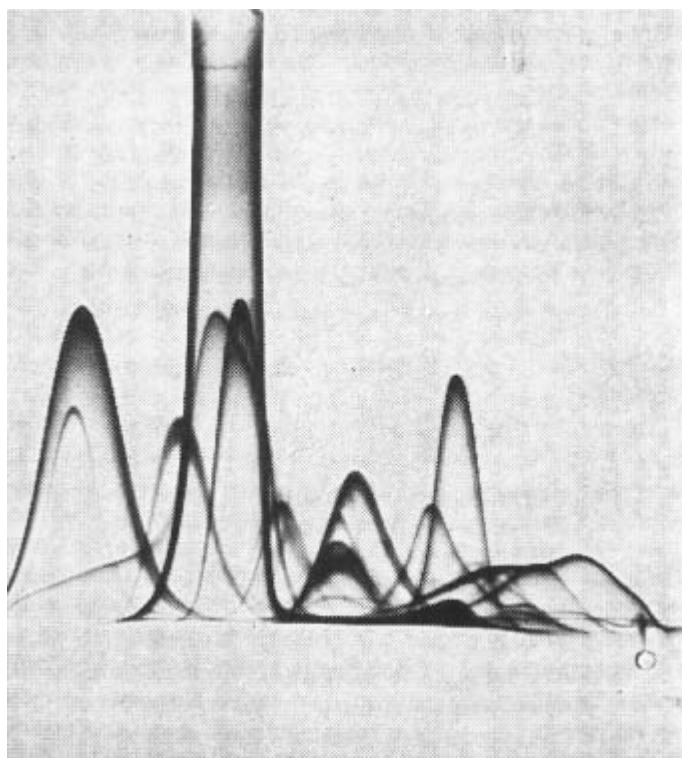
## Immuno-electrophoresis

If an [antibody](#) specific for the protein of interest is available, immunological detection and immuno-electrophoresis are the preferred methods for identification of that protein. Specifically, since the immuno-reactive site on a protein need not be affected by any [post-translational modifications](#) of a protein, immunological identification is the principal method for the recognition of a particular gene product in an electrophoretic pattern. [Blotting](#) of protein [gel electrophoresis](#) patterns to a membrane, with subsequent exposure of the blot to antibody, is the primary tool of immunological identification. More than any other factor, it is responsible for the possibility to construct protein databases from the hundreds to thousands of spots separable by [two-dimensional gel electrophoresis](#) .

Immuno-electrophoresis is the electrophoresis of one or more macromolecule in gels containing specific antibodies against them. Depending on the relative concentrations of the antibody and macromolecule, plus their affinity for each other, they interact to form complexes, which can then **immunoprecipitate** within the gel. The precipitates are visible within the gel. In this way, it is possible to quantify the amount of macromolecule present, by measurement of the size of the “rocket” produced and comparison with a known standard concentration (1). The antigenic relationships between proteins can be inferred from the precipitation patterns generated when the two are subjected to immuno-electrophoresis in adjacent wells of the gel, so that their precipitation patterns overlap (Fig. 1). If the overlapping rockets behave independently, the two molecules responsible are not immunologically related. If, however, they generate “fused rockets,” at least some of the same antibody molecules recognize both macromolecules. Crossed immuno-electrophoresis (Fig. 1) consists of an initial gel electrophoresis separation of the macromolecules of interest, followed by their electrophoresis at a right angle to the initial separation into a gel containing the antibody; this technique can yield particularly high resolution (3).

**Figure 1.** Crossed immuno-electropherogram of human serum. First dimension: Electrophoresis occurred in 1% (w/v) [agarose](#), sodium barbital buffer, pH 8.6, 0.02 M ionic strength, at 10 V/cm and 15°C, for 70 min, with the anode to the left. The circular well where the sample of human serum was applied is visible at the right. Second dimension: Electrophoresis occurred in the same gel and gel buffer, but containing 12.5 µl rabbit antibody against human serum/cm<sup>2</sup> of gel, at 2 V/cm, for 20 h, with the anode at the top. The antibody is isoelectric at pH 8.6 and therefore stationary during electrophoresis. (Fig. 3 of (1)).





### Bibliography

1. N. H. Axelsen, J. Kroll, and B. Weeke (1973) *Scand. J. Immunology*, Supplement 1, **2**, 1–169.
2. J. Clausen (1969) In *Immunochemical Techniques for the Identification and Estimation of Macromolecules* (T. S. Work and E. Work, eds.), North Holland, Amsterdam, pp. 400–572.
3. C.-B. Laurell (1972) *Scand. J. Clin. Lab. Invest.*, Supplement 124, **29**, 21–37.

### Suggestion for Further Reading

4. C.-B. Laurell (1965) Antigen-antibody crossed electrophoresis. *Anal. Biochem.* **10**, 358–361.

### Immunogen

Immunologists make the distinction between [antigen](#) and immunogen. This came from the observation that small molecules termed [haptens](#) could interact with **antibodies**, even though they were unable by themselves to stimulate an [immune response](#). Raising anti-hapten antibodies therefore necessitates preparing a conjugate of the hapten with a carrier protein. On this basis, it was proposed to distinguish antigenicity from immunogenicity and, consequently, an antigen from an immunogen. Antigenicity describes the chemical structures that contribute to interaction with an antibody, whereas immunogenicity defines the properties of a molecule for inducing an immune response. Natural immunogens, as opposed to chemically derived hapten conjugates, possess the dual properties of inducing the immune response and interacting with the resulting antibodies. Thus one must consider that a protein molecule, which is recognized by the antibodies through its [epitopes](#), also behaves as its own carrier. This simply results from the [antigen](#)

[processing/presentation](#) mechanisms that take place in the antigen-presenting cells.

Because immunogenicity is linked to the presentation phenomenon, which is genetically regulated because it involves the active participation of molecules of the [Major histocompatibility complex](#) (MHC), it is largely dependent upon the genetic constitution of the animal that receives the antigen. It has long been noticed that one strain of mice will not respond to a given antigen, whereas another one would develop an excellent response. This was strictly correlated with the presence of certain **alleles** of the MHC genes and provided a clear demonstration that the immune response was under strict (dominant) genetic control. This is why two outbred individuals, chosen at random, will respond differently to the same antigen. See [Antigen Processing, Presentation](#) for a more comprehensive explanation of immunogenicity.

Finally it should be remembered that, depending on the conditions used for [immunization](#), the same antigen may reveal itself as immunogenic or not. One important parameter is the dose of antigen given. Either too low or too high a dose may prevent immunization and even induce a **tolerant** state, so the correct dose must be carefully defined before deciding whether a molecule is immunogenic or not. Another parameter is the route of injection, which conditions the secondary lymphoid organs that will be encountered by the antigen. Also, immunogenicity may be greatly enhanced by the addition of [adjuvants](#) to the antigen solution. Immunogenicity therefore results from a large variety of causes and is far from depending solely on the nature of the immunogen.

See also entries [Antibody](#), [Antigen](#), [Antigen Processing, Presentation](#), [Epitope](#), [Haptens](#), [Immune Response](#), and [Immunity](#).

#### Suggestions for Further Reading

M. Sela (1969) Antigenicity: some molecular aspects. *Science*, **166**, 1365–1374.

H. O. McDevitt et al. (1972) Genetic control of the immune response. Mapping of the Ir-1 locus. *J. Exp. Med.* **135**, 1259–1278.

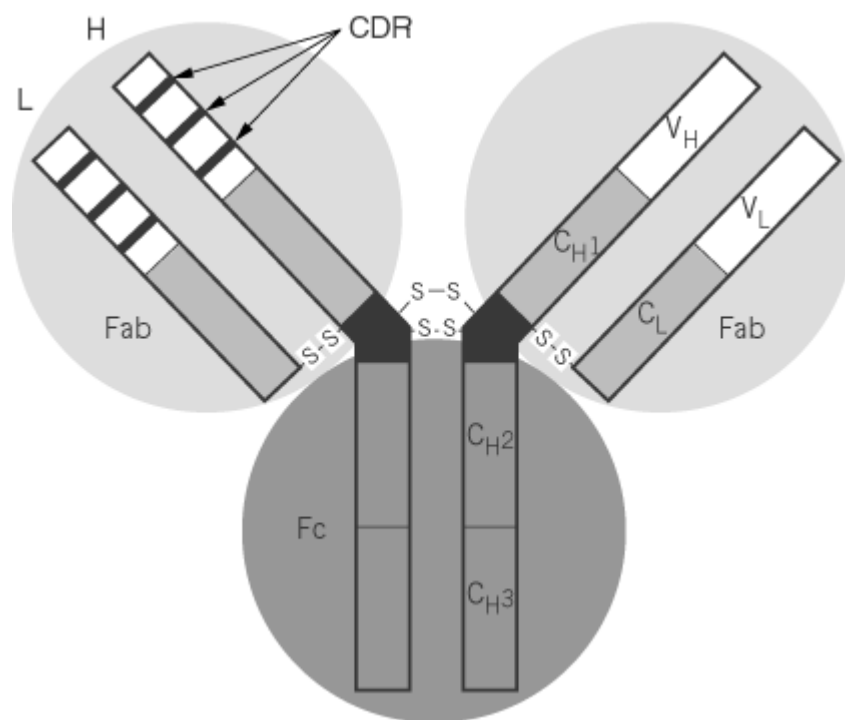
## Immunoglobulin

Immunoglobulins (Igs) are **glycoproteins** synthesized by B lymphocytes that are the agents for [B-cell](#) recognition in the immune system. They are expressed either at the surface of B lymphocytes, where they represent the B-cell receptor (BCR), or in a soluble form in the body fluids, especially in the bloodstream, where they are also designated as **antibodies**. Immunoglobulins exert a dual function: (a) They recognize native **epitopes** on the antigen surface and (b) they are endowed with biological or effector functions such as **complement** fixation, **active transport** through the placental barrier in mammals, fixation to a variety of cells that express specific receptors for the various Ig isotypes [**Fc receptors** (FcRs)], and so on. The number of potential epitopes that immunoglobulins may recognize is enormous. This is the [repertoire](#) problem. By contrast, the number of effector functions is very limited, which immediately raises a major problem for the organization of the Ig molecules, which must accommodate a high diversity of antibody combining sites at one end and rather conserved structures at the other end. This difficulty has been solved by a very special structural organization of Ig, which is the consequence of a very sophisticated organization of the **genes** that encode the Ig polypeptide chains.

### 1. The Basic IgG Topological Model

Immunoglobulins have a multichain structure characterized by the basic pattern of **IgG**, the major Ig class or **isotype**, which contains two identical heavy (H) **polypeptide chains** and two identical light (L) chains (Fig. 1). This basic pattern was defined in the 1960s as the result of many independent workers, but predominantly by Porter and Edelman, who received the Nobel Prize in 1972, and by Hilschmann, who reported the first amino acid sequences of human light chains. The H chains have a molecular weight of about 52 kDa, and the L chains have a molecular weight of 23 kDa. The IgG molecule is symmetrical and contains two identical antibody-combining (or antigen-binding) sites. It contains about 2% by weight carbohydrates, which are located on the COOH-terminal half of the heavy chains. B cells are clonally organized, so that each one makes an Ig of one particular specificity; the entire population of Ig of one individual is extremely heterogeneous. Hence, structural studies had to be performed primarily on monoclonal materials, either derived from the product of malignant cells (myeloma proteins) or from **monoclonal antibodies**. Until the beginning of the 1980s, most structural data were derived from amino acid sequencing, thereafter by **DNA sequencing** of **complementary DNA** or **genomic DNA**.

**Figure 1.** Topological model of the IgG molecule. The symmetrical molecule contains two identical H and two identical L chains, each having a variable (V<sub>H</sub> and V<sub>L</sub>) and a constant (C<sub>L</sub>, C<sub>H</sub>) region. The constant region of the heavy chain contains three homology regions (C<sub>H1</sub>, C<sub>H2</sub>, C<sub>H3</sub>) and a flexible hinge between C<sub>H1</sub> and C<sub>H2</sub>, which makes the molecule accessible to proteolytic enzymes that will release the Fab (antigen binding) and Fc (crystallizable) fragments. Homology regions provide the structural basis to the domain organization of Ig molecules.



Structural analysis revealed that Igs were built in a very peculiar way, with a general organization centered on a pseudo-unit or homology region of about 110 amino acid residues containing an intrachain **disulfide bond**. The light chain contains two such regions: one NH<sub>2</sub>-terminal, which is highly variable from one L chain to another, the so-called **variable region** or V<sub>L</sub>, and one COOH-terminal, which is highly conserved between the various chains, termed C<sub>L</sub>. There are two types of light chains, **kappa** (κ) and **lambda** (λ), built on the same model, but encoded by genes at different loci, K and L, respectively. The relative proportion of either type varies extensively from one animal species to another. In the mouse, κ chain Ig account for 95% of the total, whereas the proportion is

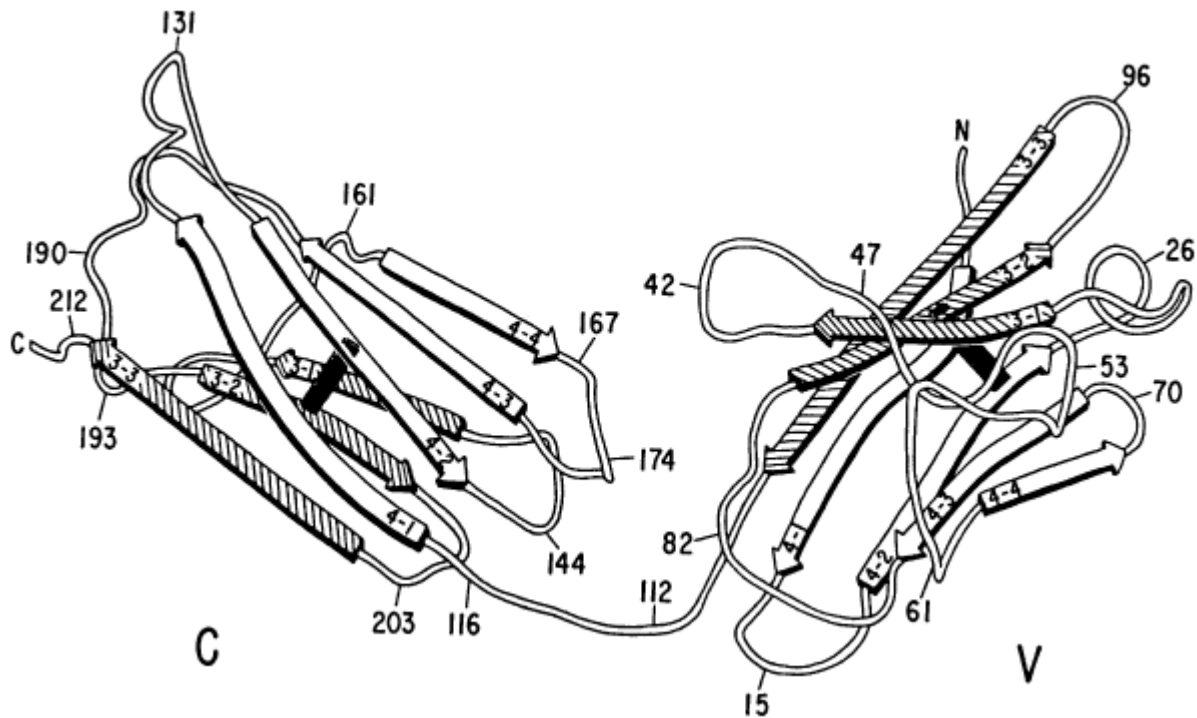
inverted in cattle or horses. In humans, the k/l ratio is approximately 60/40. Heavy chains are built on the same principle: They contain a  $V_H$  region of somewhat similar size and contain a  $C_H$  region that contains three homology regions, termed  $C_{H1}$ ,  $C_{H2}$ , and  $C_{H3}$ . The  $NH_2$ -terminal half of the heavy chain, composed of  $V_H$ - $C_{H1}$  and designated the Fd fragment, is covalently bound by a disulfide bond to the light chain. The ensemble Fd-L is the Fab (for antigen binding fragment), whereas the COOH-terminal halves of the two heavy chains form the Fc (for crystallizable fragment) because they were initially obtained by Porter upon cleavage with papain. The complete IgG therefore contains two Fab fragments, each having one identical antigen binding site, because the molecule is a symmetrical one, and one Fc; each of these fragments has roughly the same molecular weight of 150 kDa. The short section of the heavy chains located between the Fab and the Fc fragments is flexible and, for that reason, accessible to **proteolytic** cleavage, and it is called the *hinge region*. Within both the  $V_H$  and  $V_L$ , cumulative sequence data have revealed the presence of **hypervariable** regions that have been shown to interact directly with epitopes and that were for this reason called *complementarity determining regions*, or CDR (see [Immunoglobulin Structure](#)).

## 2. Domain Organization and the Basic Ig Fold

The special general architecture of Ig, characterized by the existence of the [homology](#) regions, is of particular interest for several reasons. First, it strongly suggests that this is the result of a long series of gene duplications from an ancestral gene that might have encoded a basic 110-residue polypeptide chain from which constant and variable regions would have evolved, making possible the dual functions of Ig. Second, the gene duplication phenomenon has allowed effector functions to be amplified by the emergence of other immunoglobulin **isotypes**, which fall into five classes in higher vertebrates, termed [IgG](#), [IgM](#), [IgA](#), [IgD](#), and [IgE](#). These molecules differ from each other structurally by the nature of the constant region of their heavy chains, which are g, m, a, d, and  $\epsilon$ , respectively, whereas their light chains are the same. In the course of [immunization](#), antibodies are the first of the IgM class, and there is a **class switch** later on to the production of another isotype, usually IgG, which represents 70% of circulating Ig. Because the  $V_H$  region is conserved during this switch and because the light chain remains unchanged, the antibody specificity is maintained, whereas new specialized biological functions have occurred, as described separately under specific headings.

The existence of homology regions has not only provided a basis for understanding the molecular evolution of immunoglobulins, but has also raised the point that, if these regions had diverged and been selected for the emergence of new biological functions, they had to be endowed with thermodynamic autonomy, so that each may acquire a three-dimensional (3-D) structure independently of the others; this is a prerequisite to achieve separate and specialized functions. This was the essence of the **domain** hypothesis, proposed in 1970 by Edelman, which was later fully verified when [X-ray crystallography](#) studies of antibodies and antibody-antigen complexes were reported. In fact, because of the presence of a floppy hinge region in the IgG molecule, it was not possible to get crystals of a complete molecule but, instead, of the separated Fab and Fc fragments, which turned out to have the same general multiglobular structure, precisely centered on each homology region or domain. This led to the definition of the domain “immunoglobulin fold”, depicted in Fig. 2 and well-defined from the 3-D structure of a light chain. Both the V and the C domains are similarly organized as a b-barrel, characterized by two planes containing antiparallel **b-sheets** stabilized by the intrachain disulfide bond. The various **b-strands** are connected by loops of variable length, and the V domain has extra loops that contribute the antigen binding site. There are other some minor differences between the V and the C domains, because nine b-strands are present in the  $V_H$  and  $V_L$  domains, whereas only seven are in  $C_L$  and  $C_{H1}$  and eight are in  $C_{H2}$  and  $C_{H3}$ . All CDRs are located in loops, and the 3-D analysis of an antibody-antigen complex confirms that all participate in antigen binding. The immunoglobulin fold was subsequently identified in many other proteins, including many molecules of immunological interest, providing a basis for the definition of the immunoglobulin superfamily.

**Figure 2.** Schematic organization of light-chain V and C domains, indicating the basic immunoglobulin fold. Arrows indicate the antiparallel  $\beta$ -strands as arranged in each domain as two distinct  $\beta$ -pleated sheets tightly anchored by an intrachain disulfide bond. (From Ref. 1959, with permission.)



#### Suggestions for Further Reading

- R. R. Porter (1959) The hydrolysis of rabbit gammaglobulin and antibodies by crystalline papain. *Biochem. J.* **73**, 119–126.
- G. M. Edelman and M. D. Poulik (1961) Studies on structural units of the  $\gamma$ -globulins. *J. Exp. Med.* **113**, 861–884.
- J. B. Fleischman, J. B. Pain, and R. R. Porter (1962) Reduction of gammaglobulins. *Arch. Biochem. Biophys. Suppl* **1**, 174–180.
- G. M. Edelman, B. A. Cunningham, W. E. Gall, P. D. Gottlieb, U. Rutishauser, and M. J. Waxdal (1969) The covalent structure of an entire gamma G immunoglobulin molecule. *Proc. Natl. Acad. Sci. USA.* **63**, 78–85.
- N. Hilschmann and L. Craig (1965) Amino acid sequence studies with Bence-Jones proteins. *Proc. Natl. Acad. Sci. USA* **53**, 1403–1409.
- W. J. Dreyer and J. C. Bennett (1965) The molecular basis of antibody formation: a paradox. *Proc. Natl. Acad. Sci. USA.* **54**, 864–869.
- T. T. Wu and E. A. Kabat (1970) An analysis of the sequences of the variable regions of the Bence-Jones proteins and myeloma light chains and their implications for antibody complementarity. *J. Exp. Med.* **132**, 211–250.
- L. M. Amzel and R. J. Poljak (1979) Three-dimensional structure of immunoglobulins. *Annu. Rev. Biochem.* **48**, 961.

## Immunoglobulin Biosynthesis

The biosynthesis of [immunoglobulins](#) (Igs) is a particularly complex process, because it necessitates a long sequence of events: (a) Ig [gene rearrangements](#) that take place during early [B-cell](#) differentiation in bone marrow; (b) assembly of heavy (H) and light (L) **polypeptide chains** to form the basic  $H_2L_2$  monomer; (c) transport to the cell surface, through the various intracellular compartments that permit the protein core to be sequentially glycosylated; (d) pentamerization for the [IgM](#) isotype, dimerization for [IgA](#), with binding of the [J chain](#); and (e) shedding of the final Ig as a circulating [antibody](#). In addition to these steps, one must also include **classing**, **switching** of [isotype](#), and the [alternative splicing](#) that allows the heavy chain to be synthesized in either membrane or secreted forms (see [Membrane Proteins](#) and [Protein Secretion](#)).

The phenomenon of Ig gene rearrangement is described in other entries (see [Gene Rearrangement](#), [Allelic Exclusion](#), [B Cell](#), [V Genes](#), and [Recombinase](#)). In brief, before [transcription](#), any Ig gene must first rearrange the various elements that will encode the **variable regions** of both the H and L chains. This takes place in bone marrow, where the **antigen-independent** early steps of B-cell differentiation, which are mostly devoted to Ig gene rearrangement, a highly controlled sequential succession of events ( $D \rightarrow J_H$ ,  $V_H \rightarrow D - J_H$  and  $V_L \rightarrow D - J_L$ ), occur. As a result of these rearrangements, a membrane-bound IgM is expressed at the surface of the immature B cell. Finally, as the cell leaves the marrow to enter the periphery, coexpression of IgD and IgM occurs in mature B cells. Further differentiation to plasma cells is **antigen-driven**. At this terminal stage, Igs are secreted as circulating antibodies. Expression of other isotypes implies that the switch mechanism has been operative, leading the way to all remaining steps of Ig biosynthesis mentioned above.

Each heavy and light chain is translated separately on discrete **polysomes** of the [endoplasmic reticulum](#). This is a rapid process, as for most polypeptides (about 1 min for H, and 30 s for L). The heavy chain is retained by [BiP](#), a **heat-shock** protein, Hsp70, until it binds a light chain. Two pathways have long been known, depending on the isotype and the animal species. In the first, most frequent one, half-molecules (HL) are formed that rapidly dimerize. In the other, heavy chains dimerize, and then  $H_2L$  and finally  $H_2L_2$  are formed. All these events take place in the ergastoplasm cisternae where the first carbohydrate, an *N*-acetylglucosamine, is added. Monomers are transported to the Golgi apparatus after formation of an ergastoplasmic vesicle. Subsequent carbohydrates are added, and the glycosylated monomer is included in Golgi vesicles that will fuse at the plasma membrane of the cell. After addition of a fucose residue, the molecule is either shed or anchored at the cell surface, depending upon whether this heavy chain is of the membrane or secreted type.

All surface isotypes are expressed as  $H_2L_2$  monomers. Pentamerization of secreted IgM takes place just before leaving the cell, in the immediate vicinity of the plasma membrane, simultaneously with formation of a [disulfide bond](#) to the J chain. Dimerization of the IgA isotype occurs in mucosal plasma cells, where disulfide bonding to the J chain takes place. This causes covalent binding to the poly-Ig receptor, expressed at the basal pole of the epithelial cell. After internalization, the complex is transported through the epithelial cell, and the polymerized IgA is released at the apical pole of the cell in the lumen of the corresponding organ, along with **proteolytic** cleavage between two **domains** of the poly-Ig receptor, so that one subunit remains covalently linked to the IgA polymer, forming the so-called secretory component. Besides allowing the IgA to migrate through the epithelial cell, the secretory component is thought to protect the immunoglobulin from cleavage by proteolytic enzymes.

See also entries [Antibody](#), [B Cell](#), [Gene Rearrangement](#), and [Immunoglobulin](#).

Suggestions for Further Reading

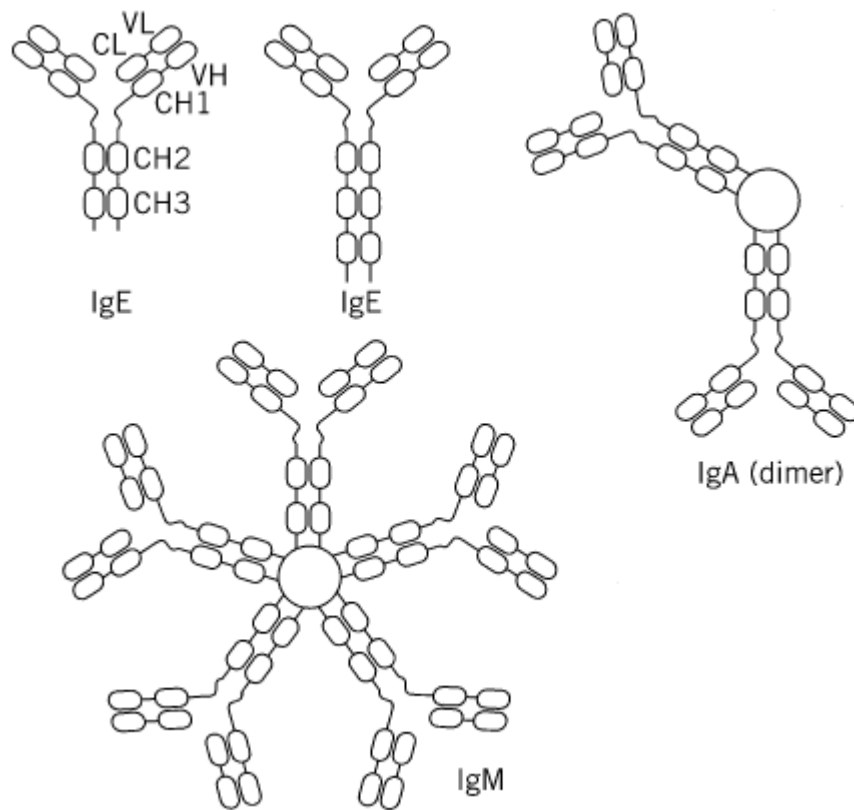
- T. K. Blackwell and F. W. Alt (1989) Mechanism and developmental program of immunoglobulin gene rearrangement in mammals. *Annu. Rev. Genet.* **23**, 605–636.
- E. S. Vitetta and J. W. Uhr (1973) Synthesis, transport, dynamics and fate of cell surface Ig and alloantigens in murine lymphocytes. *Transpl. Rev.* **14**, 50–75.
- R. Wall and M. Kuehl (1983) Biosynthesis and regulation of immunoglobulins. *Annu. Rev. Immunol.* **1**, 393–422.

## Immunoglobulin Structure

The **immune** system has evolved to recognize the entire organic world by producing complementary **immunoglobulins**, or **antibodies**, to identify and defend the vertebrate body against foreign substances (for a general reference, see (1)); an entry point to the many web sites containing information about antibodies is <http://www.antibodyresource.com>). Antibody molecules, like **enzymes** and many other **proteins**, have specific binding sites that interact with **ligands**, often in a true lock-and-key manner (see [Antibody–Antigen Interactions](#)). They also interact with other molecules involved in intercellular signaling, for example, the **Fc receptor**. Antibodies are part of a **superfamily** of proteins, some but not all of which participate in the immune response. Molecules related to antibodies that appear on cell surfaces mediate cell–cell recognition and signaling processes to trigger proliferation of particular cells in response to antigenic challenge.

The mechanism by which antibody diversity is created is now understood in considerable detail at the genetic (or sequence) level and at the protein structural level (see [Immunoglobulin Biosynthesis](#)). Common features of immunoglobulins were first recognized physicochemically and serologically, then described in terms of their amino acid sequences, and only much later revealed in atomic detail by [X-ray crystallography](#). Limited **proteolytic** digestion and cleavage of **disulfide bonds** shows that antibodies contain multiple **polypeptide chains** (Fig. 1). The chains are classified by size into light (L) chains (about 23 kDa) and heavy (H) chains (about 60 kDa). Systematic internal **homologies** in the amino acid sequences suggest that immunoglobulins are composed of multiple copies of related units, each forming an individual, quasi-independent, three-dimensional structure, or **domain**, that has a common folding pattern. Light chains contain two domains, and heavy chains contain four or five. Heavy chains were distinguished originally on the basis of their amino acid sequences into k and l classes or **isotypes**. In humans, k and l light chains are present in comparable proportions. In the [mouse](#), k light chains predominate.

**Figure 1.** Domain structures of the most common immunoglobulins. Most immunoglobulins contain both light and heavy chains, each of which comprises one variable domain and varying numbers of constant domains. The combination of two light plus two heavy chains is a higher order building block, of which immunoglobulins of different classes (IgG, IgA, IgM, IgD and IgE) show different states of oligomerization. Additional chains are linkers where necessary.



Different types of antibody molecules show variations on a general structural theme. They are multichain proteins built from domains about 100 amino acid residues long known as the immunoglobulin fold. Before the first crystal structure determinations of immunoglobulins, Kabat and co-workers identified the regions involved in antigen binding by analyzing the distribution of variability in aligned immunoglobulin sequences (2). Two types of domains were recognized on the basis of the variability of their sequences: variable (V) and constant (C) domains. Within the variable domains, they noticed regions of still greater variability called the hypervariable regions or complementarity determining regions (CDR), which they correctly and presciently predicted would be responsible for antibody specificity for different antigens. The regions of the variable domain outside the CDRs are called the **framework regions**. The framework provides a scaffolding of nearly constant structure to which the antigen-binding loops are affixed.

Indeed, when the first structures of immunoglobulin domains were determined by X-ray crystallography, it appeared that both constant and variable domains had a similar, double [beta-sheet](#) structure in the framework region, that the hypervariable regions corresponded to surface loops in the variable domains, and that these loops did indeed interact with bound antigens. In almost all immunoglobulin domains, a conserved disulfide bond links the two sheets, and a **tryptophan residue** packs against it. Chothia and Lesk subsequently carried out the three-dimensional analog of the sequence comparisons of Kabat and coworkers by using the crystal structures (3). This led to some revision of the boundaries between the framework and CDR regions.

Different classes of immunoglobulins—IgG, IgA, IgM, IgD, and IgE—differ in their assembly from chains and domains. Figure 1 shows the domain structures of the different classes of immunoglobulins. Molecules in the class best known structurally, the IgGs, usually have two heavy chains that contains four domains and two light chains that contain two domains. The antigen-combining site in most IgGs is formed from three loops from the light chain and three from the heavy chain. Each IgG has two copies of the binding site, one from each light chain–heavy chain pair, which permits multiple interactions with antigens to form aggregates.



The IgG molecule contains four polypeptide chains, two identical light chains, each containing one variable and one constant domain, denoted  $V_L$  and  $C_L$ , and two identical heavy chains, containing one variable and three constant domains, denoted  $V_H$ ,  $C_{H1}$ ,  $C_{H2}$ , and  $C_{H3}$  (see Fig. 2). The  $C_{H1}$  and  $C_{H2}$  domains are linked by a segment of polypeptide chain called the hinge region. The angle between the  $V_L$ - $V_H$  domain pair and the  $C_L$ - $C_{H1}$  domain pair is called the elbow angle. These joints are usually flexible in solution. Immunoglobulins also contain carbohydrate moieties introduced by *N*-glycosylation that are not shown.

**Figure 2.** A schematic diagram of the structure of an IgG molecule (a), showing the distribution of domains in the heavy and light chains and (b) the interchain disulfide bridges, and the definitions of the fragments - Fab, Fab' and Fc - produced by limited proteolytic cleavage.

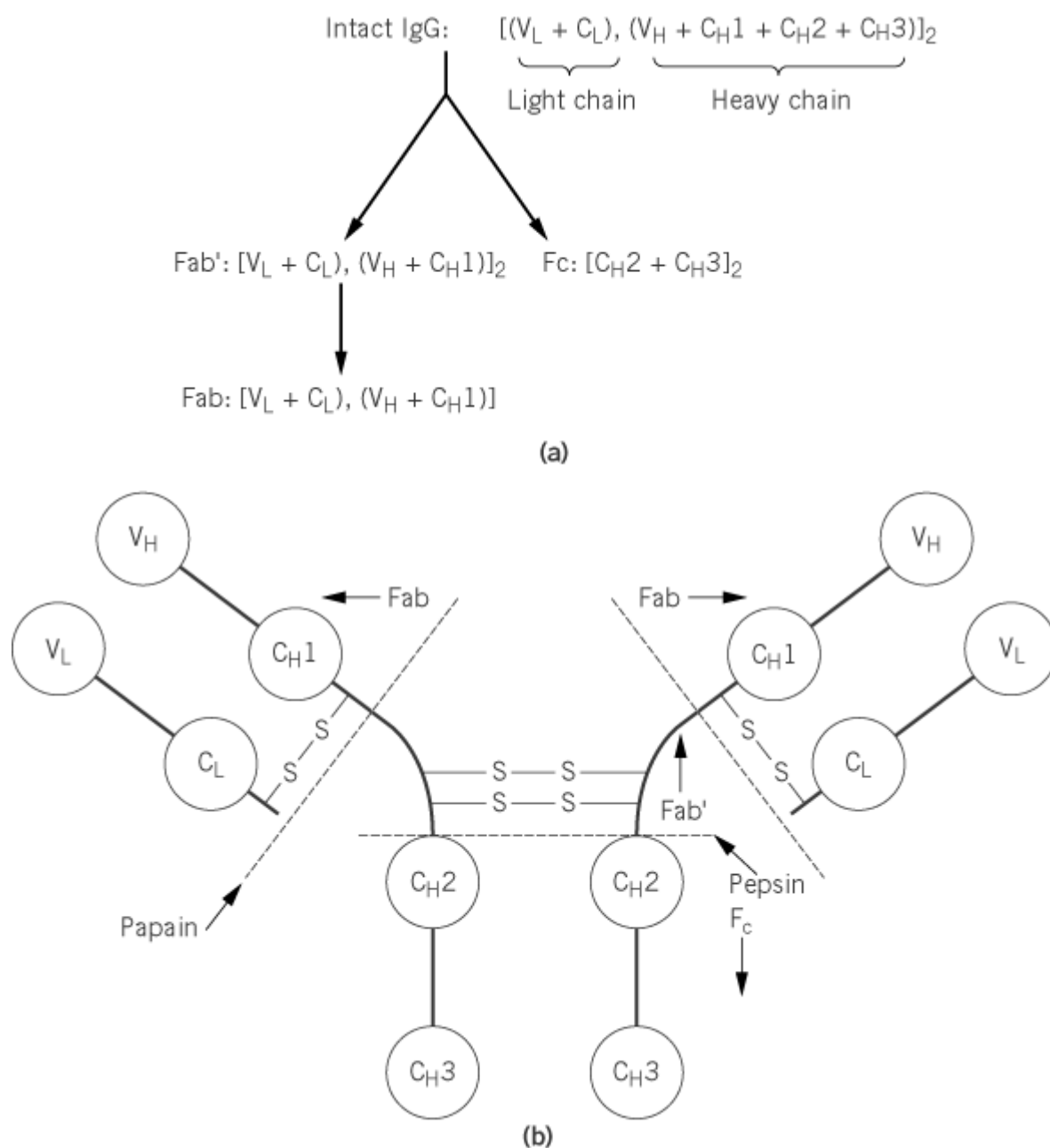
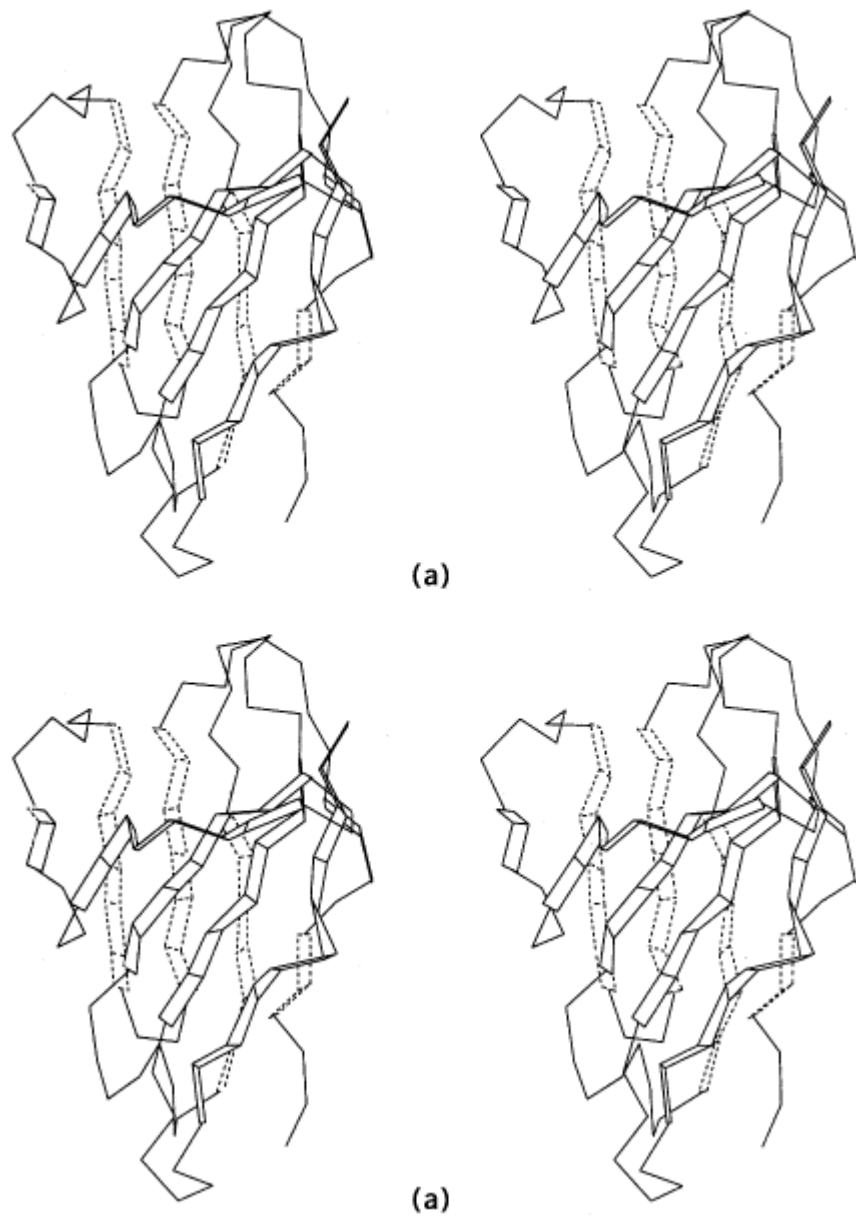


Figure 2 defines the fragments produced by limited proteolytic cleavage. Most of the structural analysis has been carried out on Fab fragments (4). Single-chain Fv fragments, in which the C-terminus of the  $V_H$  domain is linked by a flexible peptide to the N-terminus of the  $V_L$  domain, or vice versa, are important in [protein engineering](#) (5).

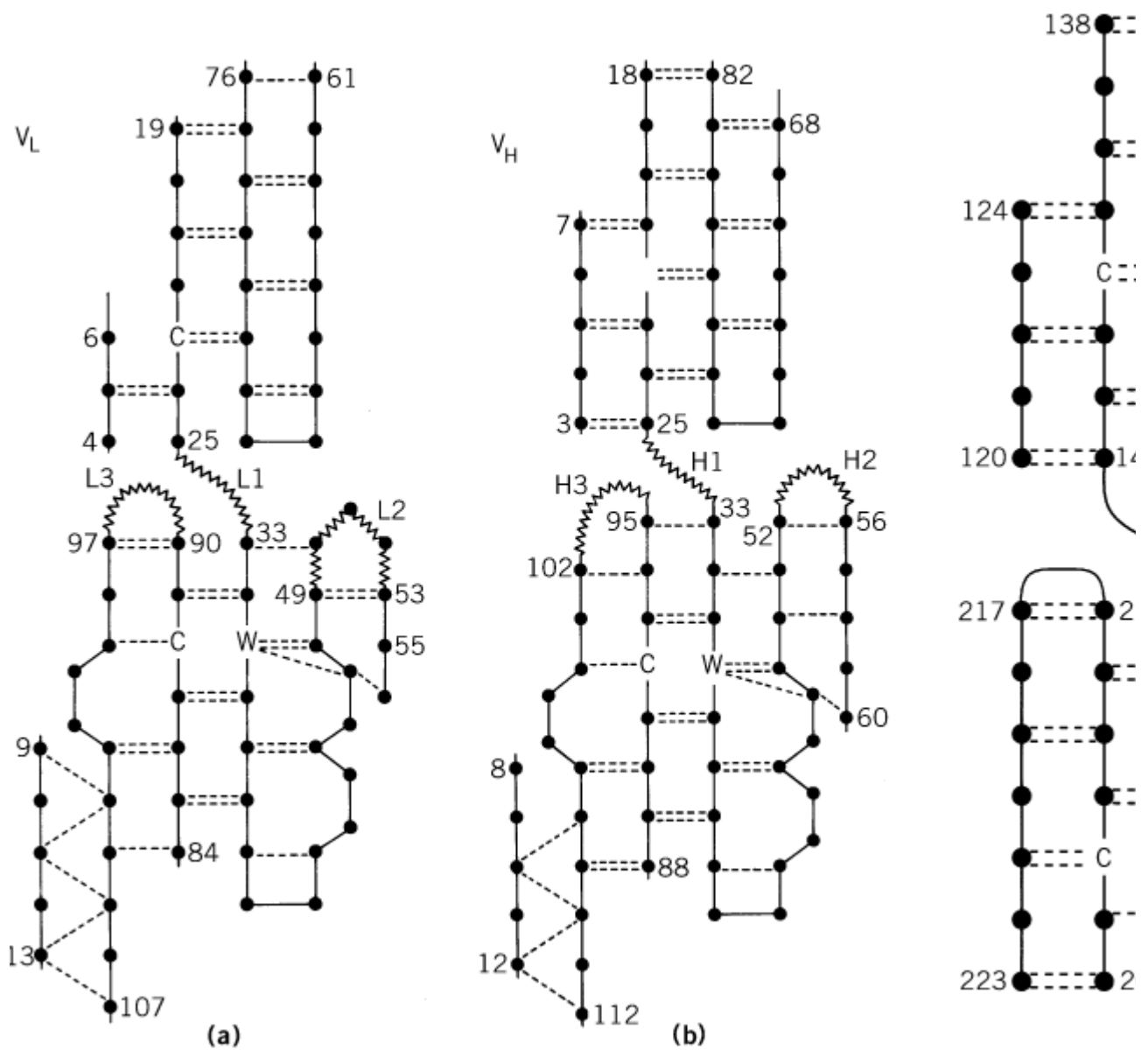
The interactions between the different polypeptide chains in an IgG include the disulfide bonds (see Fig. 2) and the interfaces between corresponding domains— $V_L$ – $V_H$ ,  $C_L$ – $C_{H1}$ , and  $C_{H3}$ – $C_{H3}$ , but not  $C_{H2}$ – $C_{H2}$ —that pack together to form extensive **van der Waals** contacts. It is characteristic of proteins that interior interfaces are formed by the packing of complementary surfaces. This complementarity of fit fixes the relative spatial disposition of the pieces that interact. A tendency to conserve the residues involved in these interfaces explains why different light and heavy chains pair fairly freely to form complete immunoglobulins. In the  $V_L$ – $V_H$  interaction, conservation of the relative geometry has the additional important consequence that, to a reasonable first approximation, the double b-sheet frameworks of  $V_L$  and  $V_H$  domains form a scaffolding of nearly constant structure on which the antigen-binding site is erected.

Figures 3 and 4 show the **secondary structures** and [tertiary structures](#), respectively, of the V and C domains. Each has the form of a double b-sheet “sandwich.” In immunoglobulin domains, the two sheets are oriented so that their strands approximately parallel. Figure 5 shows the interface between  $V_L$  and  $V_H$  domains (6).

**Figure 3.** The folding patterns of (a) variable domains and (b) constant domains. The domains shown are the  $V_L$  and  $C_H$ . Strands of b-sheet are shown as polygonal ribbons, and the loops joining them as lines that join the  $C_{\alpha}$  atoms of successive amino acids. Each domain contains two b-sheets. One domain is shown in solid lines and the other in broken lines.



**Figure 4.** Schematic representations of typical hydrogen-bonding patterns of (a)  $V_L$ , (b)  $V_H$  and (c) C domains. The two opened like a book. Strands in the upper sheets in this figure were indicated by ribbons with broken lines in Figure 3. Str in this figure correspond to the strands drawn with solid lines in Figure 3.



**Figure 5.** The interface between the  $V_L$  and  $V_H$  domains in the antilysozyme antibody D1.3 (8). The  $V_L$  domain is trace  $V_H$  domain by a broken line. In this perspective, the antigen-binding site is on the upper surface of the molecule. Side of the  $V_L$  chain that contact bound antigen are shown with filled circles, and those from the  $V_H$  chain with open circles. The strongly twisted.

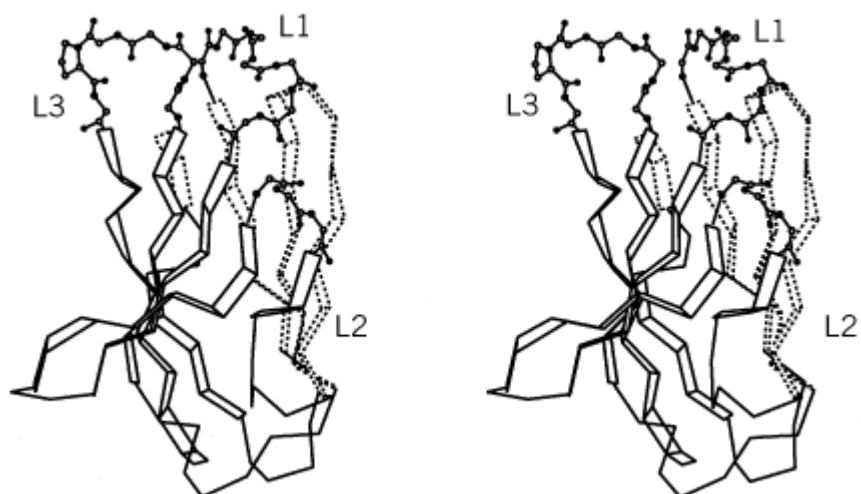


### 0.1. The Antigen-Binding Site

The study of antibody antigen-binding sites lies at the intersection of two topics of interest. One is the immune response: What is the relationship among the sequences of immunoglobulins generated by genetic combination and [somatic mutation](#) and the three-dimensional conformations of the resulting sites? The second is the question of how the recognition and binding of an antigen by antibodies relate to our understanding of protein–ligand interactions, and the extent to which we can quantitatively explain a specificity at the level of interatomic interactions.

Most antibody antigen-binding sites are formed primarily from six loops—three from the  $V_L$  domain and three from the  $V_H$  domain. The three loops from the  $V_L$  domain are called L1, L2, and L3, in order of their appearance in the amino acid sequence. Alternatively, they are called CDR1, CDR2, and CDR3.  $V_H$  domains contain three complementary loops H1, H2, and H3. In contrast, the antibodies of the camel and related animals contain only heavy chain  $V_H$  CDRs suffice to create the antigen-binding site. This observation has obvious applications to selecting artificial antibodies for therapy. Four of the loops, L2, H2, L3, and H3, are “**hairpins**,” in that they link antiparallel strands of a single  $\beta$ -sheet, connecting strands of sheet that are hydrogen-bonded to each other. L1 and H1 bridge a strand from one of the two  $\beta$ -sheets to a strand in the other (Fig. 6).

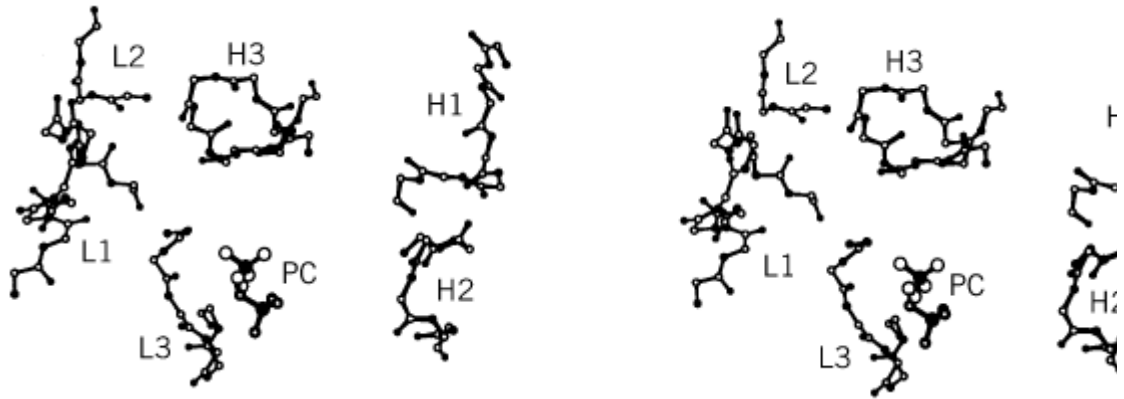
**Figure 6.** The  $V_L$  domain of REI, indicating the antigen-binding loops L1, L2, and L3.



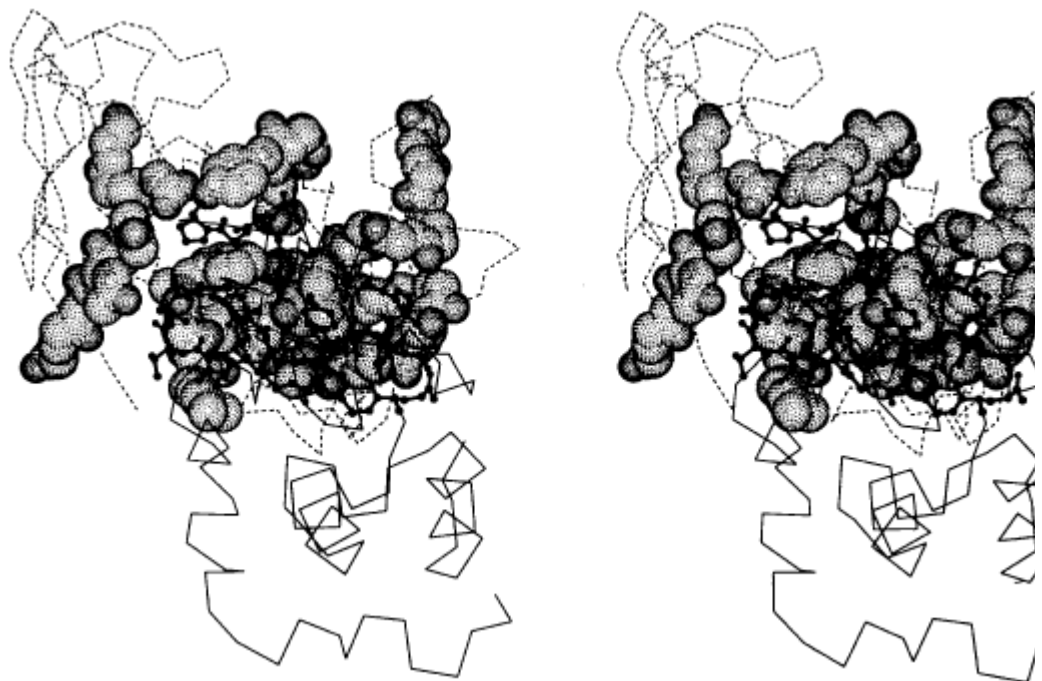
PC603 (Fig. 7) and the other the complex between Fab HyHEL-5 and hen egg white [lysozyme](#) (Fig. 7). The orientations were chosen looking down onto the antigen-binding site, that is, an “antigen’s-eye view” of the site. The light-chain loops are on the left and the heavy-chain loops are on the right. Reading clockwise starting at

loops appear in the order L1, L2, H3, H1, H2, and L3. The arrangement of the loops is roughly sym opposite H1, L2 is opposite H2, and L3 opposite H3. The central position of H3 is noteworthy.

**Figure 7.** The spatial distribution of antigen-binding loops, or CDRs, in the antigen-binding site of McPC603. PC = 1 phosphorylcholine. The orientation is chosen to give a view down onto the antigen-binding site, an “antigen's-eye view”



**Figure 8.** Interactions between HyHEL-5 and hen egg white lysozyme. Residues in lysozyme that make contact with t in bold ball-and-stick representation. Residues in the antibody that make contact with the lysozyme are shown in skele Broken lines indicate hydrogen bonds between antigen and antibody.



In many cases, antigen-antibody interactions involve the rigid association of the two structures. In mechanism analogous to “[induced fit](#)” takes place, in which the structures of the free and ligated a significant structural changes. The structures of complexes between antibodies and protein antigen antibodies specific for different **epitopes** of lysozyme and one specific for **neuraminidase**, demon general features of the interactions are similar to those seen when other proteins associate without l conformational changes (7) (see [Protein-Protein Interactions](#)). The number of residues in contact g

from 27 to 39, distributed approximately equally between antibody and antigen. The [accessible surface area](#) of the antibody–antigen interface varies between 1250 and 1940 Å<sup>2</sup>, and the interfaces are densely packed. In many cases several buried [water](#) molecules occupy sites within the interface. There are typically 10 [hydrogen bonds](#) across the interface.

In the complex between Fab D1.3 and hen egg white lysozyme, for example, the interface is a relatively flat surface with dimensions **20 Å × 30 Å**, although one glutamine side chain from the lysozyme inserts into a pocket (8). A total of 17 residues from the antibody make contact with a total of 17 residues from the lysozyme. Pairs of residues are in contact. Of the 17 residues from the antibody, three are framework residues adjacent to CDRs.

In the lysozyme–antibody complexes of known structure, there is rather little change in structure on binding relative to their unliganded states. Therefore, these complexes fit the picture of recognition by fairly rigid bodies that have preformed complementarity. This is not observed, however, in complexes of antibodies and other proteins clipped out of protein structures, because the latter are flexible in solution. In the complex between a nine-residue region from influenza virus hemagglutinin (9), the peptide in the complex has the conformation of a type I b-turn, somewhat different from its conformation in the native hemagglutinin structure. Comparison of the structure of the Fab fragment of the antibody in free and liganded forms shows that there has also been a rearrangement of the structure of loop H3. The relative orientations of V<sub>L</sub> and V<sub>H</sub> domains change on binding antigen by a rotation, typically of about 3°, but up to 11° in known structures.

Another peptide–Fab complex, B13I2 (10), contains the seven-residue epitope of a 19-residue antigen derived from the C-helix of myohemerythrin. The residues of the peptide in the complex have the conformation of a β-turn, not helical as in the native protein. What is most significant is that the antibody cross-reacts with apomyohemerythrin, which it could not do while retaining its native structure. It is interesting that apomyohemerythrin binds more strongly to the antibody. Apomyohemerythrin lacks the iron atoms, one of which binds to the C-helix in the native state. This implies that myohemerythrin changes conformation to bind the antibody. The removal of a stabilizing interaction—a barrier to conformational change—enhances binding. (The structure of the B13I2–myohemerythrin complex has not been determined.)

### 0.3. The Sequence Repertoire of the Immune System

Analysis of sequences of the loci on human [chromosomes](#) that carry the genes for antibodies and comparison of these results with the amino acid sequences of expressed antibodies reveal what genetic resources the immune system possesses and how it uses them. The following discussion is restricted to human sequences.

The gene for the variable region of the light chain is formed during B-cell development by combinatorial joining of three segments: V<sub>L</sub> (variable), D (diversity) and J (joining). Each individual locus contains approximately 51 functional V<sub>L</sub> segments (11), 25 D segments, and six functional J segments. The number of possibilities implied by these various combinations, they apparently are not used with equal frequency. Comparing a database of expressed antibodies with the set of V<sub>L</sub> germ-line sequences, of a total of 51 germ-line V<sub>L</sub> segments, only about 11 are used frequently, and others appear in only a minority of expressed sequences, 10% of the functional V<sub>L</sub> gene segments produce 50% of the expressed V<sub>L</sub> repertoire (12).

The nonuniform appearance of sequences is also reflected in a bias in the structures of the hypervariable regions of the associated primary antibodies. As discussed later, antigen-binding loops of antibodies usually adopt a few conformations chosen from a small repertoire of discrete canonical structures. The antibodies studied to date show strong preferences for particular main-chain conformations in their light-chain hypervariable regions.

corresponding analysis of the  $V_H$  segments—which code for a region containing the first two hyper variable regions of the  $V_H$  domain—revealed similar nonuniform structural preferences (14). For example, sequences of the CDR1 and CDR2 segments, which code for the part of the  $V_H$  domain that includes loops H1 and H2, suggest that there are two different canonical structures of H1 and six of H2, potentially producing 18 different structures. However, the sequences contain only seven different combinations of canonical structures of H1 and H2 (14).

#### 0.4. Somatic Mutation and the Maturation of the Antibody Response

The distribution of residues that differ from their germ-line identities in expressed antibodies has been studied by Tomlinson and co-workers (15) by comparing sequence diversity on a residue-by-residue basis between germ-line antibodies and mature ones. Whereas the diversity in the germ-line antibodies is greatest in H3 and at the center of the antigen-binding site (because of junctional diversity in gene assembly; see [Immunoglobulin Biosynthesis](#)), somatic mutations spread the diversity to its periphery. In almost all cases, somatic mutations are isolated single-residue substitutions, not insertions or deletions. The total change is less than the total change between germ-line segments, that is, if one thinks of the germ-line segments as “islands in a sequence sea,” somatic mutations extend the domains of the individual germ-line segments, rather than filling in between them. Thus, there was no ambiguity in identifying the germ-line segment from which each mature antibody arose. Somatic mutation is a perturbation of a germ-line segment rather than a major reconfiguration of the antigen-binding site.

Wedemeyer et al. (16) studied the maturation of an anti-nitrophenyl phosphonate catalytic antibody by X-ray crystallography. They solved the structures of the germ-line and somatically mutated Fab fragment without ligand. Their results reveal both the sites of mutation and the structural effects. There are nine amino acid sequence changes between germ-line and mature antibody, three in the light chain and six in the heavy chain. The light chain mutants are in the L1 loop. The other two are at positions near in the sequence to L1. The heavy-chain mutants appear in H2, three are in regions adjacent to or in contact with the antigen-binding site, and two are surface residues at the opposite ends of the domains, which presumably have little effect on the antigen-binding site. No mutation is at a position directly in contact with the antigen. No mutation appears in H3, although there are extensive contacts between H3 and the hapten. The conformations of two of the loops (H1 and H2) differ between germ-line and mature antibody.

The antigen-binding site of the mature antibody has a similar conformation with and without bound antigen, a conformation that is complementary to the hapten. This structure is shared by the ligated state of the antibody, but the unligated state of the germ-line antibody differs in conformation, that is, the germ-line antibody does not adopt the hapten-binding conformation, but this is induced by the ligand. In contrast, the mature antibody adopts the conformation even in the absence of ligand.

Wedemeyer et al. (16) suggest that conformational flexibility in the primary antibody response represents a mechanism for generating antibody diversity. Indeed, it is possible that still other conformations of the antibody might be induced by alternative, related ligands. In this way the same primary antibody could serve as a starting point for alternative maturation pathways to produce secondary antibodies to different ligands. However, there is no direct structural evidence that this happens, nor how different the induced conformation would be.

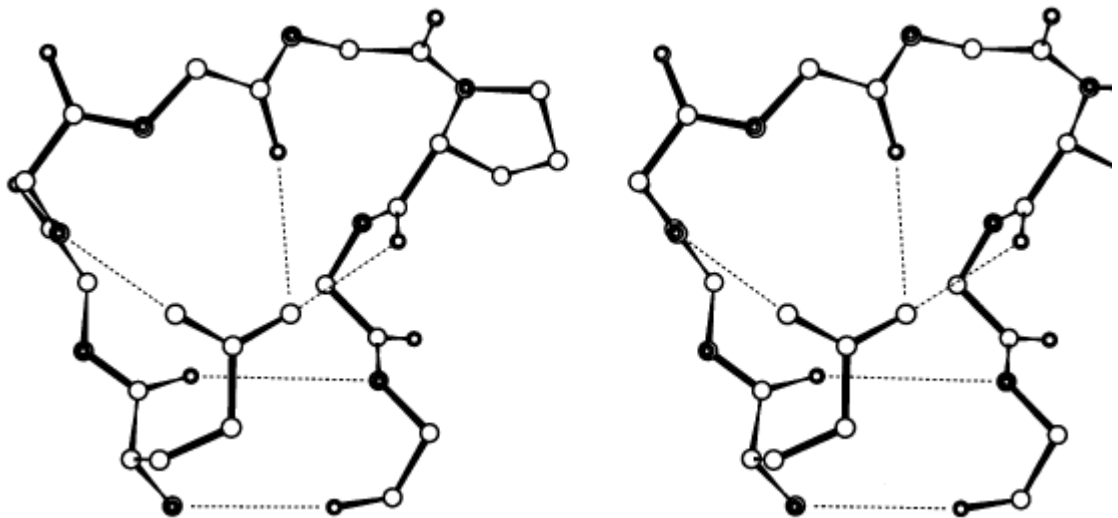
#### 0.5. Structures of Loops in Proteins

The many known protein structures demonstrate that generally loops connecting strands of a  $\beta$ -sheet adopt characteristic sets of conformations, which depend on their lengths, sequences, and interactions with their surroundings. For short hairpins in proteins, the conformation is usually determined primarily by the sequence of the loop itself (17). For a polypeptide chain to reverse direction in a “tight” turn, it usually takes advantage of a special residue, such as Gly or Pro, that allows a nonstandard main-chain conformation. Conformations of many short hairpins depend on the length and position of such a special residue. This also applies to many of the shorter CDR hairpins, L2, L3, and H3. H2 is exceptional, however, in that a single residue strongly influences its conformation (18).

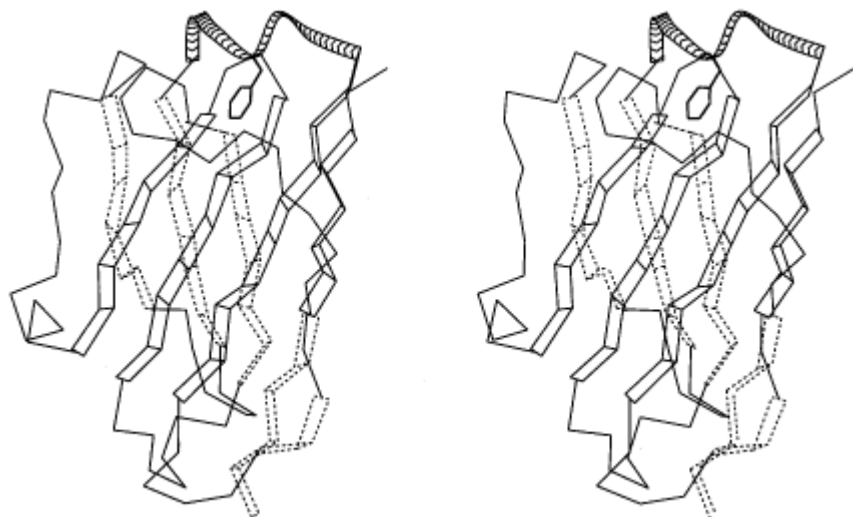


Medium-sized loops fall into two classes (19). The conformations of those that have ends close together are usually determined by hydrogen-bonding interactions of polar atoms that point into the interior of the loop. A typical  $V_k$  L3 loop, illustrated in Figure 9, exemplifies this case.) The conformations of more extended loops, such as H1, depend on the packing of bulky side chains against or into the core of the protein. In H1, for example, residue 29 of the heavy chain is usually a Phe, which is packed into the framework (see Fig. 10).

**Figure 9.** The most common canonical structure of  $V_k$  L3, that from REI. This shows the characteristic *cis*-proline at the C-terminus and the polar side-chain at the residue just before the loop that forms hydrogen bonds to stabilize the conformation of the loop.



**Figure 10.** The H1 loop of Fab J539. An important contribution to stabilizing the conformation of this loop is the packing of Phe29 into the region between the two sheets.



### The Canonical Structure of Antigen-Binding Loops of Antibodies

Studies of immunoglobulin structures have shown that their antigen-binding loops have only a small number of discrete main-chain conformations. These conformations are determined by a few particular residues outside the loop but interacting with it. Among corresponding loops of the same length, only these

be conserved that maintain the conformation of the loop (20). The conserved residues may be those special main-chain conformations or that form special hydrogen-bonding or packing interactions. The residues in the sequences of the loops are left free to vary to modulate the surface topography and distribution of the antigen-binding site.

The biological significance of this robustness of the main-chain conformation of the antigen-binding site is to play a role in tuning antibodies by somatic mutation. In the system as it has evolved, a change in the antigen-combining site in most cases produces only a conservative structural change. This is what happens in a system that already shows affinity and requires only “fine-tuning.” Consider the alternative. If one completely altered the main-chain conformation of the loops, the effect would produce a successive series of primary responses, rather than a secondary response whose structures are perturbed in only minor (but important) ways from a set of already selected primary antibodies. One should not underestimate the effect of tiny structural changes to produce very large effects on binding affinity.

The ability to isolate the determinants of loop conformation in a few particular residues in the sequence makes it possible to analyze the distribution of loop conformations in the many known immunoglobulin structures. Out of the six hypervariable regions of antibodies—L1, L2, L3, H1, and H2—have only a few main-chain conformations or “canonical structures” (H3 is discussed separately later). Most sequence variation is on the surface by altering the side chains on the same canonical main-chain structure. Sequence changes at a few positions switch the main chain to a different canonical conformation.

As an example, Figure 9 shows the L3 loop from V<sub>k</sub> REI. There is a proline residue at position 90, just N-terminal to the peptide bond. Hydrogen bonds between the side chain of the residue at position 90, just N-terminal to the main-chain atoms of residues in the loop stabilize the conformation. The side chain is a Gln in REI and in other V<sub>k</sub> chains. The combination of this loop length, one of these polar side chains at position 90, and the proline constitute the “signature” of this conformation in this loop, from which it can be recognized in a sequence of an immunoglobulin for which the structure has not been experimentally determined. The conformation of the loop (10) illustrates a nonhairpin loop, which joins strands in separate sheets. Its conformation is stabilized by inward-pointing side chains against the framework. General support for the canonical structural determinants of antigen-binding loops has come from analysis of known sequences and structures of immunoglobulin structures. The conformations observed in the relatively few known immunoglobulin structures account for the large proportion of all immunoglobulins.

When the canonical structural model was proposed, analysis of the large compilation of immunoglobulin sequences of Kabat and co-workers (21) showed that over 90% of the hypervariable regions in V<sub>k</sub> and V<sub>h</sub> have the same length and contain the “signature” residues of a known canonical structure and therefore have conformations close to those found in the currently known immunoglobulin structures. Indeed, the model was conservative in that it includes only those hypervariable regions that matched the sequence exactly. Of course one was dealing then with a limited number of known structures, and the model was overturned if subsequent structural determinations produced an unbounded proliferation of main-chain conformations. However, the conclusion that the repertoire is limited and discrete is justified, although many canonical structures have appeared. Recently, 244 hypervariable regions in 49 immunoglobulin F<sub>1</sub> structures determined at resolutions of between 1.7 and 3.1 Å were analyzed according to their structure (22). Some 85%, and even more in subsequent work, were clustered into groups identified as known structures and had very similar conformations. The structures on which the canonical structure are based are primarily human and mouse antibody fragments. However, analysis of the sequences of shark and other vertebrate antibodies suggest that the basic pattern of canonical structures goes back to the very early stages in the evolution of the immune system (23). These observations suggest that the main-chain conformations of five of the regions—L1, L2, L3, H1, and H2—are described by a small repertoire of canonical structures. Most of the structures present among the immunoglobulin structures already known, and the main-chain conformations of an immunoglobulin of unknown structure can very often be predicted from the sequence by the pr

specific residues (24).

A roster of current canonical structures is available (20). Its accuracy and completeness depend on (1) the determination of the sets of residues responsible for the observed conformations and (2) the assumption that changes in the identities of residues at other sites do not significantly affect the conformation of the binding loops. These results have been tested and refined by using them to predict the atomic structures of binding sites in immunoglobulin structures before determining them by X-ray crystallography (24).

#### 0.6.1. Precision of the Canonical Structural Model

Qualitatively, the canonical structure model states that the main-chain conformations of five of the binding loops fall into small, discrete clusters. It is important to know quantitatively how similar are different exemplars of the same canonical structure in different antibodies, for a proper understanding of the mechanism of immune recognition and for the designing or modeling antibody structures. The two questions addressed are, what is the spread in the main-chain conformations of the canonical structures and what is the spread in their positions and orientations with respect to the framework? When the canonical structure was initially proposed, this question was difficult to answer, partly because the differences in observed conformations were comparable with the experimental error. Some of the structures had been determined only at low resolution and, even more seriously, some contained qualitative errors.

Now, the precision of canonical structures has been examined by analyzing and comparing the conformations of the antigen-binding loops in 17 immunoglobulins for which accurate structures have been determined at high resolution (20). The structures used were determined to a resolution of 1.6 to 2.3 Å and refined to an R-factor greater than 21%. These contain 7 V<sub>L</sub>, 10 V<sub>K</sub>, and 15 V<sub>H</sub> chains. Although the hypervariable regions of particular immunoglobulins do not cover all of the known canonical structures, they do include all of the canonical ones. The results indicate that the main-chain conformations of different examples of the same canonical structure in different antibodies are usually conserved to within about 0.3 to 0.4 Å rms deviation, which is a high degree of similarity.

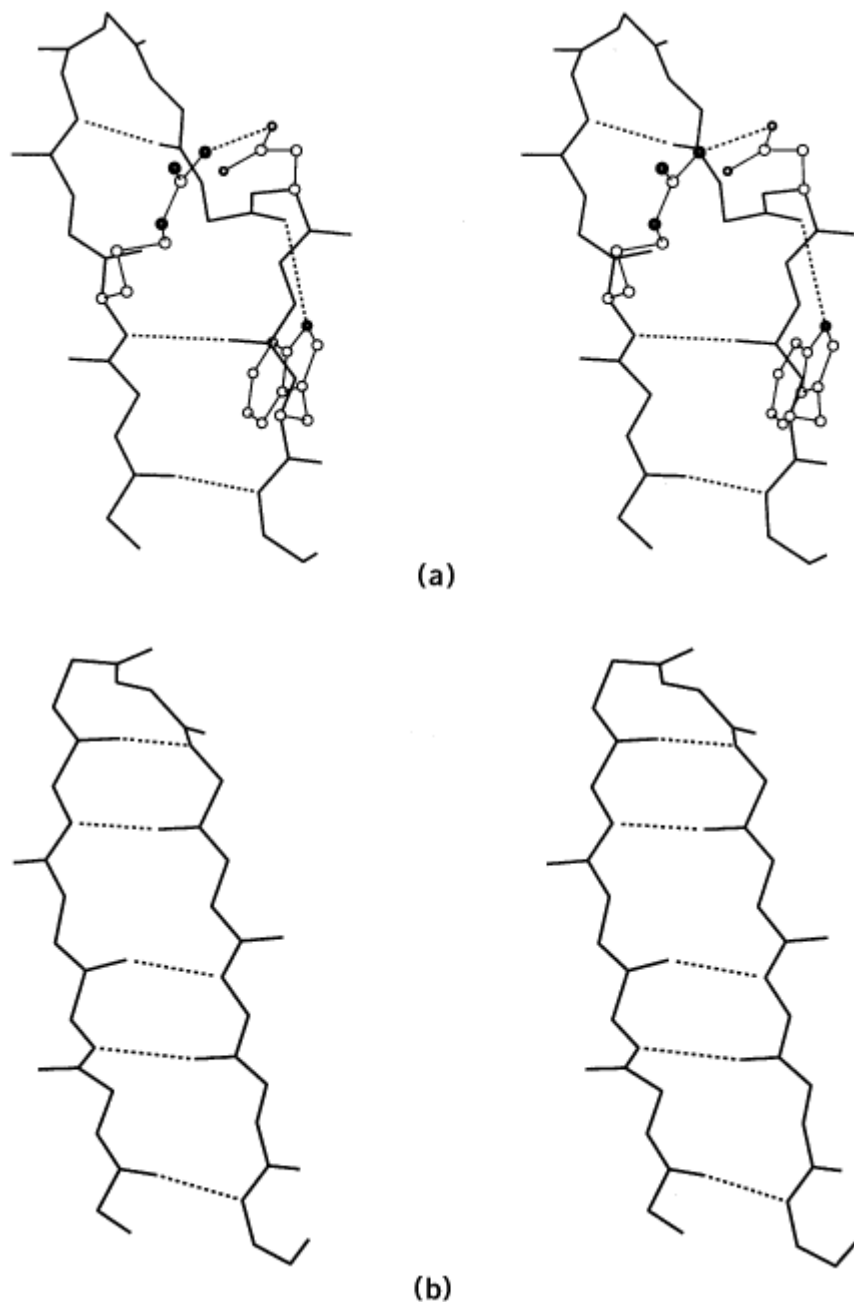
#### 0.6.2. The H3 Loop is Even More Variable

H3, the third hypervariable region of the heavy chain, is far more variable in length, sequence, and position than the other antigen-binding loops (21). Therefore, it cannot be included in the canonical structural model. The conformational repertoire that applies to the three hypervariable regions of V<sub>L</sub> chains and the first two hypervariable regions of V<sub>H</sub> chains. Because the H3 loop falls in the region of the V-D-J join in the assembly of the immunoglobulin heavy chain gene (Fig. 11), several mechanisms contribute to generating its diversity, including combinations of V<sub>H</sub>, D, and J gene segments and [alternative splicing](#) patterns at the junctions.

**Figure 11.** The roles of the V, D, and J gene segments in determining the immunoglobulin fold. Loop H3 arises from the V-D-J join in the assembly of the heavy chain gene. In the V<sub>H</sub> domain of McPC603 shown here, the region encoded by the D gene segment coincides almost exactly with the H3 loop region.



**Figure 12.** Two H3 loops with alternative torso conformations: (a) McPC603 - the bulged torso conformation, (b) 4 nonbulged torso conformation.



In expressed antibodies, H3 is prominently at the center of the antigen-binding site (see figs. 7,8) central position, H3 makes significant interactions with other loops with the framework, with the partner, and with ligands, which influence its conformation. Thus H3, in contrast to the other five binding loops, has a conformation that depends strongly on its molecular environment. Indeed, a containing the  $V_H$  domain of antibody B1-8 combined with two different  $V_L$  domains has two ve conformations of H3 (25). This important observation implies that general rules governing the cc H3 (unlike the other five antigen-binding loops) must involve interactions outside the local regio

To include all residues that contribute to determinants of the conformation of the H3 region and l conserved residues Cys92 and Gly104 provide useful landmarks to identify the H3 region in a ne H3 region is defined as the residues from Cys92 to Gly104. Analysis of the conformations of H3 divided the H3 region into a torso section comprising a head, the apex of the loop, and residues p framework—four residues starting from the conserved Cys92 and from the N-terminus and six res

C-terminus. For H3 structures that contain more than 10 residues, there are two main classes of conformations (see Fig. 11). In the major class, the conformation of the torso has a **b<sub>12</sub>-bulge** at a few H3 regions, the torso does not contain a bulge, but continues the regular hydrogen-bonding sheet.

The choice of bulged or nonbulged torso conformation is dictated primarily by the sequence. A bulge is formed whenever residues are present that permit formation of a [salt bridge](#) between the side chain of a Lys residue at position 94 and Asp101, unless a residue at position 93 forms the salt bridge. In the minor class, the torso region were chosen as the largest set of residues proximal to the framework that have a repertoire of conformations, as in the canonical structure model of other antigen-binding loops. The heads or apices of the H3 loops have a very wide variety of conformations. In shorter H3 regions containing the nonbulged torso conformation, the heads follow the rules relating sequence to structure. For longer H3 regions containing the bulged torso conformation, there are many very different conformations of the head, which can be catalogued but are difficult to classify.

Accurately predicting the conformation of the torso of the H3 region from amino acid sequence is possible in most cases. Predicting the conformation of the head is possible in some cases. However, our understanding of sequence–structure relationships has reduced the uncertainty to no more than a few residues at the head of the H3 region, but these residues appear in a crucial position within the antigen-binding site!

#### 0.7. Prediction of Antigen-Binding Sites

The possibility of accurately predicting the three-dimensional structures of immunoglobulins at the antigen-binding site is centrally important for engineering antibodies with prescribed specificity. Several methods for predicting the structures of antigen-binding sites have been proposed. Some apply information from known structures, while others are based on *a priori* methods.

The canonical structural analysis of the conformations of the hypervariable loops suggests a model for the H3 loop: Construct the frameworks by standard **homology modeling** techniques. Then identify, if possible, canonical structures of the loops and graft them onto the modeled framework.

Another procedure for predicting the structures of antigen-binding loops relies on conformational search calculations. The main-chain conformation of an antigen-binding loop attached to a given framework is constrained that the chain must connect two fixed end points with a specified number of residues. If the number of residues is fewer than about six residues, it is possible to enumerate a fairly complete set of feasible main-chain conformations, which bridge the given end points and do not make steric collisions within the loop or between the loop and the rest of the molecule. The search procedure can be fine enough to include the correct one in the many possibilities enumerated. To choose one of them as the predicted conformation, it is possible to estimate conformational energies and to evaluate the accessible surface areas of each conformation in the context of the remainder of the protein and set criteria for selecting the most favorable one. Typical conformational energy calculations include terms representing hydrogen bond, van der Waals, and [interactions](#). Accessible surface area calculations estimate the interaction between the protein and antigen. Note that these procedures are completely general, independent of the current state of the structure.

Procedures for conformation generation and evaluation have been implemented in a number of computer programs. An application to predict all six antigen-binding loops of McPC603 and HyHEL-5 has been reported (29). The resulting models of the antigen-binding loops of McPC603 and HyHEL-5 have rms deviations of 0.3 to 2.6 Å (backbone) and 1.4 to 4.1 Å (all atoms).

Jones and Thirup (30) developed an alternative approach to modeling loops, based on selecting loops from proteins in the database of known structures that span the given end points and overlap with peptide termini. For this procedure to be useful, it is required (1) that the loops in immunoglobulin antigen-binding sites

appear in other proteins in the database; (2) that the geometrical relationship between the loop and that flank it are the same in immunoglobulin antigen-binding sites and proteins in the database of the same conformation; and (3) if several loops are found in the database that properly join the flanking an antigen-binding loop, it must be possible to select the correct one. This problem is faced in methods in which loop conformations are generated *a priori*. If the selection from candidates identified from homologous regions of other immunoglobulins by database searching is based on patterns in the sequences, the procedure reduces in effect to the canonical structural method.

The main-chain conformations of most antigen-binding loops recur in other antibody structures (of course, one of the elements of the canonical structural model) and also in unrelated proteins (18), exceptions are unusually long loops, such as L1 of McPC603 or H3 of KOL. In other words, the loops whose structures commonly appear in antibodies are not unique to antibody structures. However, the geometrical relationship between the loop and its flanking peptides usually differs when comparing a binding loop with a loop of similar conformation in an unrelated protein (except for short hairpin conformation). In these cases it would not be possible to build a correct model of an unknown loop. A protocol that combines these approaches, by using a CONGEN procedure alone for short loops and supplements it with database search techniques for longer loops, has been described (32).

## 1. Conclusions: How Does the Immune System Generate and Then Refine Molecules of General

Based on the studies of sequences and structures, we can provide some fairly detailed answers to questions about what the immune system is doing at the structural level. Here we are concerned with the generation of diversity in structure and thereby in affinity and specificity but do not treat the cellular mechanisms which challenge by antigens elicits the proliferation of those cells that produce antibodies that bind ([Immunoglobulin Biosynthesis](#)).

The antigen-combining site has a main-chain conformation determined by the canonical structures L2, L3, H1, and H2, plus the more variable H3. The number of potential main-chain structures in the binding site is equal to the product of the number of possible canonical structures for loops L1, L2, L3, H1, H2, and of the number of possible structures of H3. Although putting numbers into these statements yields large numbers of potential antibody sequences, there are preferences for structure at the level of the canonical structures. Because only a relatively few residues are required to fix the canonical structure of the main chain, most of the antigen-binding loop residues can be freely decorated with side chains to vary the topography and charge distribution of the binding site.

The immune system is one of many biological examples of molecular recognition. Individual antigen-antibody complexes fit the classical ideas of lock-and-key complementarity and induced fit, but the system is more complex than, for example, enzymes, in which the lock and key are in most cases fixed and do not change within the lifetime of the organism. In the primary immune response, the immune system achieves a wide variety of keys by providing many locks, but to achieve the spectacular affinity of the secondary response, it uses a mechanism to perturb the locks for better fit. The integration of the immune response over the system involving both sequence and structural dimensions together with the clinical and biotechnological applications makes it a fascinating topic to study.

## 2. Acknowledgment

We thank Drs. I. Tomlinson and E. Gherardi for a critical reading of this manuscript.

## Bibliography

1. A. K. Abbas, A. H. Lichtman, and J. S. Pober *Cellular and Molecular Immunology*, 3rd ed., Philadelphia, 1997.

2. T. T. Wu and E. A. Kabat (1970) *J. Exp. Med.* **132**, 211–249.
3. C. Chothia and A. M. Lesk (1987) *J. Mol. Biol.* **196**, 901–917.
4. E. Padlan (1996) *Adv. Protein Chem.* **49**, 57–133.
5. R. Glockshuber, M. Malia, I. Pfitzinger, and A. Plückthun (1990) *Biochemistry* **29**, 1362–1369.
6. C. Chothia, J. Novotny, R. E. Bruccoleri, and M. Karplus (1985) *J. Mol. Biol.* **186**, 651–663.
7. J. Janin and C. Chothia (1990) *J. Biol. Chem.* **265**, 16027–16030.
8. T. N. Bhat, G. A. Bentley, G. Boulot, M. I. Green, D. Tello, W. Dall'acqua, H. Souchon, F. J. A. Mariuzza, and R. J. Poljak (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1089–1093.
9. J. M. Rini, U. Schulze-Gahmen, and I. A. Wilson (1992) *Science* **255**, 959–965.
10. R. L. Stanfield, T. M. Fieser, R. A. Lerner, and I. A. Wilson (1990) *Science* **248**, 712–719.
11. I. M. Tomlinson, G. P. Cook, N. P. Carter, R. Elaswarapu, S. Smith, G. Walter, L. Buluwela and G. Winter (1994) *Hum. Mol. Genet.* **6**, 853–860.
12. J. P. L. Cox, I. M. Tomlinson, and G. A. Winter (1994) *Eur. J. Immunol.* **24**, 827–836.
13. O. Ignatovich, I. M. Tomlinson, P. T. Jones, and G. Winter (1997) *J. Mol. Biol.* **268**, 69–77.
14. C. Chothia, A. M. Lesk, E. Gherardi, I. Tomlinson, G. Walter, J. D. Marks, M. B. Llewelyn, (1992) *J. Mol. Biol.* **227**, 799–817.
15. I. M. Tomlinson, G. Walter, P. T. Jones, P. H. Dear, E. L. L. Sonnhammer, and G. Winter (1994) *J. Mol. Biol.* **256**, 813–817.
16. G. J. Wedemayer, P. A. Patten, L. H. Wang, P. G. Schultz, and R. C. Stevens (1997) *Science* **275**, 1669.
17. B. L. Sibanda, T. L. Blundell, and J. M. Thornton (1989) *J. Mol. Biol.* **206**, 759–777.
18. A. Tramontano, C. Chothia, and A. M. Lesk (1989) *Proteins: Struct. Function Genet.* **6**, 382–391.
19. A. Tramontano, C. Chothia, and A. M. Lesk (1990) *J. Mol. Biol.* **215**, 175–182.
20. B. Al-Lazikani, A. M. Lesk, and C. Chothia (1997) *J. Mol. Biol.* **273**, 927–948.
21. E. A. Kabat, T. T. Wu, H. M. Perry, K. S. Gottesman, and C. Foeller (1991) *Sequences of Proteins of Immunological Interest*, 5th ed., Public Health Service, N.I.H, Washington, D.C..
22. A. C. R. Martin and J. M. Thornton (1996) *J. Mol. Biol.* **263**, 800–815.
23. S. Barré, A. S. Greenberg, M. F. Flajnik, and C. Chothia (1994) *Nat. Struct. Biol.* **1**, 915–920.
24. C. Chothia, A. M. Lesk, A. Tramontano, M. Levitt, S. J. Smith-Gill, G. Air, S. Sheriff, E. A. Davies, W. R. Tulip, P. M. Colman, S. Spinelli, P. M. Alzari, and R. L. Poljak (1989) *Nature* **338**, 325–329.
25. X. Y. Pei, P. Holliger, A. G. Murzin, and R. L. Williams (1997) *Proc. Natl. Acad. Sci. USA* **94**, 1175–1179.
26. H. Shirai, A. Kidera, and H. Nakamura (1996) *FEBS Lett.* **399**, 1–8.
27. V. Morea, A. Tramontano, M. Rustici, C. Chothia, and A. M. Lesk (1998) *J. Mol. Biol.* **275**, 101–110.
28. A. Sali (1995) *Mol. Med. Today* **1**, 270–277.
29. R. E. Bruccoleri, E. Haber, and J. Novotny (1988) *Nature* **335**, 564–568.
30. T. A. Jones and S. Thirup (1986) *EMBO J.* **5**, 819–822.
31. A. Tramontano and A. M. Lesk (1992) *Proteins: Struct. Function Genet.* **13**, 231–245.
32. A. C. R. Martin, J. C. Cheetham, and A. R. Rees *Proc. Natl. Acad. Sci. USA* **86**, 9268–9272.
33. G. Winter, A. D. Griffiths, R. E. Hawkins, and H. R. Hoogenboom (1994) *Ann. Rev. Immunol.* **12**, 455.
34. P. G. Schultz and R. A. Lerner (1995) *Science* **269**, 1835–1842.



## Immunophilin

Cytosolic receptor proteins for immunosuppressant drugs in mammals of the [cyclosporin](#) and peptidomacrolide type have been named immunophilins because of their putative involvement in the signaling mechanism of the cellular [immune response](#) (1). It was further pointed out that the immunophilin/immunosuppressant complex is active, not the individual free drug or receptor. A composite surface formed by the bound drug and by the surrounding immunophilin residues functions by binding other cellular proteins involved in [signal transduction](#). Initially, because of the limited knowledge about the diversity of these enzymes, the groups of immunophilic receptor proteins coincided with the enzyme classes of [peptidyl prolyl cis/trans isomerases](#) (PPIase), cyclophilins, and FK506-binding proteins (FKBP). Now, cyclophilins and FKBP have been isolated that display only modest affinity for these drugs or cannot form active enzyme/drug complexes. Furthermore, the new PPIase families of [trigger factors](#) and parvulins do not bind the immunosuppressive drugs at all. On the other hand, there are now also many other proteins that have considerable affinity for the immunosuppressants and are not PPIase (2). To avoid confusion with nomenclature, the term immunophilin should be limited to PPIases with proven involvement in immunosuppression.

### Bibliography

1. R. F. Standaert, A. Galat, G. L. Verdine, and S. L. Schreiber (1990) *Nature* **346**, 671–674.
2. G. Fischer (1994) *Angew. Chem., Int. Ed. Engl.* **33**, 1415–1436.

### Suggestions for Further Reading

3. G. Wiederrecht and F. Etzkorn (1994) *The immunophilins, Perspect. Drug Discovery Design* **2**, 57–84.
4. S. L. Schreiber and G. R. Crabtree (1992) *The mechanism of action of cyclosporin-A and FK506, Immunol. Today* **13**, 136–142.
5. S. H. Snyder and D. M. Sabatini (1995) *Immunophilins and the nervous system, Nature Med.* **1**, 32–37; summarizes immunophilins of the nervous system.

## Imprinting

### 1. Mammalian Imprinting: From Genomes to Genes

Sexual reproduction, which involves the fusion of two gametes produced by two distinct individuals, is the most common form of reproduction encountered among higher eukaryotes. Although sexual reproduction has obvious long-term benefits, such as being a powerful source of genetic variability and in triggering evolution, many short-term disadvantages, including what is called the “cost of meiosis” (1), could have acted against its maintenance during evolution. In this regard, one can view imprinting as a way to “maintain” sexual reproduction, by providing an absolute short-term

advantage. Parental imprinting, also known as genomic or gametic imprinting, is a genetic effect in which the expression (behavior) of **genomes, chromosomes**, or chromosomal regions differs according to their parental origin. As a result, individuals from species in which imprinting occurs require a contribution of both parental genomes to develop properly. This phenomenon is observed in many biological systems, from bugs, where sexual determination involves the inactivation and sometimes even the loss of the whole paternal genome (2), to mammals, which will be the main subject of this review.

In mammals, evidence that maternal and paternal genomes are imprinted to function differentially during development initially came from the study of parthenogenetic and androgenetic mouse embryos generated by pronuclear transplantations (3-5). In general, parthenogenetic embryos (two maternal genomes) fail to develop normal extraembryonic components, whereas androgenetic embryos (two paternal genomes) have a very poorly developed embryo proper. Mouse genetic studies involving the analysis of maternal and paternal disomies and/or duplications of defined chromosome regions have subsequently demonstrated that this phenomenon does not affect the whole genome *per se*, but rather only specific regions (6). Ten regions on six mouse autosomes have been shown to be subject to imprinting [chr 2, 6, 7, 11, 12, and 17, (6)]. The maternal and paternal duplications for five of the imprinted regions have distinct consequences, indicating that there are in total 15 known imprinting effects. One cannot exclude the possibility, however, that there are additional imprinted regions with phenotypic effects too subtle to be detected.

Further investigations have revealed that the genetic basis for these imprinting effects is likely to be the genes themselves, because some of the genes lying within these regions are subject to imprinting, which means that they are monoallelically expressed, either from the paternal or the maternal chromosome, depending on the gene. In the same way, monoallelic expression can occur in most tissues of the individual or can be restricted to a very specific developmental stage and/or [cell lineage](#). Now, more than 25 imprinted genes have been identified in mouse, and most of these are conserved in human (Table 1). Despite the growing list of imprinted genes and their more detailed characterization, the function of imprinting remains a mystery. The abnormalities displayed by parthenogenetic and androgenetic mouse embryos underscore, however, that imprinting plays a pivotal role in development.

**Table 1. List of Imprinted Genes in Human and Mouse and Some of Their Characteristics<sup>a</sup>**

| Gene  | Function      | Chr | Exp       | Human   |   |    |
|---|---------------|-----|-----------|---|---|----|
|   |               |     |           | Specificity   | Methylation                             | Ch |
| <b>Human Chromosome 11/Mouse Chromosome 7</b> |               |     |           |   |   |    |
| IGF2/Igf2                                     | Growth factor |     | pat       | Mostly monoallelic except in choroid plexus,leptomeninges and adult liver |   |    |
| INS2/Ins2                                     | Insulin       |     | pat + mat |   |   |    |
| H19   | Noncoding RNA |     | mat       | Mostly monoallelic except in early placenta                               | Hypermethylation of the paternal allele |    |
| HASH2/Mash2                                   | Transcription |     | mat       | Mostly monoallelic  |   |    |

|                     |                             |     |  |
|---------------------|-----------------------------|-----|--|
| KVLQT1/Kvlqt1       | factor<br>Potassium channel | mat | Mostly monoallelic except heart  |
| p57 <sup>KIP2</sup> | Cyclin-cdk inhibitor        | mat | Mostly monoallelic with weak expression of the paternal allele in most tissues; biallelic in fetal brain |
| IPL/Ipl             | Role in apoptosis?          | mat | Mostly monoallelic with some mosaicism   |
| IMPT1/Impt1         | Metabolite transporter?     | mat | Allelic bias   |

### Human Chromosome 15/Mouse Chromosome 7

|                       |                                 |     |   |  |
|-----------------------|---------------------------------|-----|---|--|
| ZNF127/Zfp127         | Zinc-finger protein             | pat | Mostly monoallelic                      | Hypermethylation of the maternal allele  |
| ZNF127A<br>S/Zfp127as |                                 |     |   | pat  |
| NECDIN/Necdin         |                                 | pat | Mostly monoallelic                      | Hypermethyl of the maternal allele   |
| SNRPN/Snrpn           | Splicing                        | pat | Mostly monoallelic                      | 5' Region methylated on the maternal allele<br>3' region methylated on the paternal allele |
| PAR-SN                | No ORF                          | pat | Monoallelic (lymphoblast)               |  |
| PAR-5                 |                                 | pat | Monoallelic (fibroblast)                |  |
| IPW                   | No ORF                          | pat | Mostly monoallelic                      |  |
| PAR-1                 |                                 | pat | Monoallelic (fibroblast)                |  |
| UBE3A/Ube3A           | Ubiquitin protein ligase        | mat | Moallelic only in brain                 | No allelic methylation identified  |
| DST                   |                                 | mat | Moallelic only in brain                 |  |
| UBE3A/DST-antisense   | Regulation of UBE3A imprinting? | pat | Expressed and monoallelic only in brain |  |

### Miscellaneous Genes

|     |       |      |     |                    |    |
|-----|-------|------|-----|--------------------|----|
| p73 | Tumor | 1p36 | mat | Mostly monoallelic | NT |
|-----|-------|------|-----|--------------------|----|

|                               |   |             |           |   |  |    |
|-------------------------------|---|-------------|-----------|---|--|----|
| IGF2R/Igf2r                   | suppressor<br>IGF2<br>receptor                                | 6q25-<br>26 | mat + pat | Mostly biallelic.                           | Region 2<br>hypermethylated<br>on the maternal<br>allele | 17 |
|                               |   |             |           | Some cases of<br>monoallelic<br>expression  |  |    |
| Igf2r <sup>-</sup> -antisense | Regulation<br>of Igf2r<br>imprinting?                         |             |           |   |  | 17 |
| MAS/Mas                       | Cell surface<br>receptor                                      | 6q25-<br>26 |           | Monoallelic only in<br>adult breast         |  | 17 |
| PEG1-MET/Peg<br>1-Mest        | a/b<br>Hydrolase<br>fold family                               | 7q32        | pat       | Monoallelic in fetus,<br>biallelic in blood | Maternal<br>methylation                                  | 6  |
| HTR2                          | Serotonin<br>receptor   | 13q14       | mat       | Monoallelic<br>(fibroblast)                 | Maternal<br>methylation                                  | 14 |
| U2af1-rs1                     | Splicing  | NC          |           |   |  | 11 |
| Peg3/Pw1                      | ?   |             |           |   |  | 7  |
| Grf1                          | Regulates<br>synthesis and<br>release of<br>growth<br>hormone |             |           |   |  | 9  |
| Peg5/neuronatin               |   |             |           |   |  | 2  |
| Impact                        |   |             |           |   |  | 18 |
| XIST/Xist                     | X-inactiva-<br>tion   | X           |           |   |  | X  |

<sup>a</sup> The first part of the table corresponds to the human chromosome 11/mouse chromosome 7 cluster; the second part of the table corresponds to the human chromosome 15/mouse chromosome 7 cluster; and the third part of the table corresponds to miscellaneous genes. Chr, chromosomal localization; Exp, allele expressing the gene; pat, paternal; mat, maternal. Specificity: specificity of the monoallelic expression. Methylation: allelic methylation. Noted as G when inherited from the gametes, S (somatic) when not. NC (in the case of U2af1-rs1 in man) indicates not conserved at the location homologous to that in mouse. Blank cells indicate not tested or not published results.

This is also true in humans, because the presence of two paternal genomes without a maternal one results in the formation of hydatidiform moles (7), which is essentially equivalent to that found with the androgenetic embryos in the mouse. In addition, the deregulation of imprinting causes disease, and the expression of several inherited diseases is strongly influenced by the sex of the transmitted parent. The best-known examples are the Beckwith–Wiedemann syndrome (BWS) and the Prader–Willi and Angelman syndromes (PWS and AS) (8). BWS is often associated with lack of a biparental inheritance of chromosome 11p15 and is characterized by pre- and postnatal overgrowth, including hemihypertrophy, visceromegaly, and macroglossia (enlarged tongue), and an increased probability of developing childhood cancers, the most frequent ones being Wilm's tumor and adenocortical carcinoma. The involvement of imprinting in this syndrome was suggested by the preferential maternal transmission of mutations in BWS (9, 10) as well as the occurrence of paternal uniparental disomies (UPDs) of chromosome 11p in several sporadic cases of BWS (11).

PWS and AS are two clinically distinct syndromes that share a chromosomal location (15q11-q13) but display opposite parental origins of transmission. PWS is characterized by hypotonia, hypogonadism, variable mental retardation and obesity (12), while the AS clinical features include hyperactivity, ataxic gait, “puppet-like” arm movements, bouts of inappropriate laughter, seizures, EEG abnormalities and severe mental retardation with absence of speech (13, 14). Most (60% to 70%) PWS and AS cases are associated with *de novo* cytogenetic deletions of the 15q11-q13 region. In PWS, the deletion occurs exclusively on paternal chromosome 15 (15), whereas AS chromosome 15q11-q13 deletions are of maternal origin (16). Maternal and paternal uniparental disomy (UPD) also occur in some PWS (17) and AS (18) cases, respectively. Roughly 25% of PWS patients display UPD with both chromosomes 15 being inherited from the mother and no chromosome 15 being derived from the father. Only about 5% of AS cases result from paternal UPD and maternal nullisomy for chromosome 15. These observations indicate that the gene(s) responsible for PWS are expressed only from the paternal chromosome, whereas the AS gene(s) are transcribed from the maternal chromosome.

The list of human diseases resulting from abnormal imprinting will no doubt expand as the catalogue of imprinted genes (see Table 1) becomes more complete. The identification of other loci subject to genomic imprinting requires some knowledge of the properties specific to such regions.

## 2. Characteristics of Imprinted Regions and Genes

The occurrence of imprinting and its consequent monoallelic expression implies that both alleles of the same gene can be recognized and differentially regulated within the same nucleus. The fact that imprinting occurs in laboratory inbred mouse strains, which are genetically identical at all loci, indicates that the basis of imprinting must be *epigenetic*. Imprinting is a two-step process that must fulfill two conditions. The first step is the apposition of the mark, the imprint, which will allow one allele of a gene (or a chromosomal region) to be distinguished from its homologue. This “labeling” must occur during gametogenesis, which is the only stage where both alleles are physically separated. The question remains as to which allele, the repressed or the expressed one, carries the imprint, although recent evidence suggests that it is the silent allele that is modified in the gamete (see text below). The second step of the process is the reading of the imprint in the appropriate developmental stage and/or tissue, which implies that the imprint has not been erased during development and is stably maintained. Stability is one of the required characteristics of imprinting. Once established, the imprint must indeed be clonally inherited and maintained through cell divisions. Reversibility is the other necessary property of imprinting, because the imprint must be erased and properly reestablished at each generation. Both stability and reversibility are features of, and can be explained by, epigenetic modification of the chromatin.

Differential **methylation** of the cytosine residue in **CpG** dinucleotides is an attractive candidate as the “label” of imprinted genes and regions, because it is an epigenetic modification that is reversible and heritable, and it can play a role in the regulation of gene (19, 20). The association of DNA

methylation with transcriptional activity and [chromatin](#) conformation is most evident in the case of the mammalian [X chromosome](#), where the **promoters** of several **housekeeping genes** on the inactive X chromosome are enriched for 5-methylcytosine. These results suggest that repression of gene activity can be mediated through methylation of 5' regulatory elements. This association is less clear for tissue-specific autosomal genes where, despite a significant body of suggestive evidence, there is no direct proof that DNA methylation modulates transcriptional activity.

A significant difference in the degree of methylation between the paternal and maternal genomes has been elegantly demonstrated by Rougier et al. (21). These authors used **antibodies** to 5-methylcytosine to label metaphase chromosomes in mouse preimplantation embryos, and they observed distinct differences in the labeling patterns at the one-cell stage. The **euchromatic** arms of the maternal chromosomes were intensely stained and demonstrated an R-like banding pattern. In contrast, the antibody staining of the paternal set of chromosomes was quite faint (21). While these data demonstrate the marked difference in global DNA methylation between the two parental chromosome complements in preimplantation development, more direct evidence of a role for DNA methylation in the process of epigenetic marking has been derived by the analysis of sites of allelic methylation within imprinted genes. Most imprinted genes indeed display parent-specific methylation (Table 1). For some of them, including *H19* and *SNRPN* (22-25), methylation is associated with the repressed allele, whereas for others, like *Igf2*, it is the expressed allele that is methylated (26). Strong evidence for the pivotal role of DNA methylation in the regulation of imprinted gene expression comes from studies in mice bearing a disruption of the DNA methyltransferase gene (*Dnmt*), the enzyme responsible for the maintenance of cytosine methylation (27). Such mice display embryonic lethality and a content of 5-methylcytosine that is reduced to about 30% of wild-type levels (27). Moreover, the expression of several imprinted genes in *Dnmt*<sup>-/-</sup> mice is abnormal with, for example, activation of the normally repressed paternal *H19* allele and silencing of the normally active paternal *Igf2* allele (28). While these results clearly indicate the functional role of allelic DNA methylation in imprinted gene expression, the imprinting of some genes, such as *Mash2*, appears normal in *Dnmt*<sup>-/-</sup> mice, suggesting the possibility of different mechanisms of imprinting in terms of its association with methylation (29).

Another property of imprinted regions is asynchronous DNA replication. Chromosome replication banding studies generally indicate that homologous regions on both autosomes of a pair replicate at the same time in the S phase of the [cell cycle](#) (30). Using this approach, Izumikawa et al. first observed that this was not the case for the imprinted chromosomal region 15q11-q13 (31). That imprinted chromosomal regions replicate asynchronously has been confirmed in other replication banding studies (32, 33) and by using a different assay of DNA replication involving the technique of fluorescence [in situ hybridization](#) (FISH) (34). The use of the FISH assay has revealed that DNA replication asynchrony is a feature of other imprinted regions (34-36). The parental origin of the earlier replicating allele was investigated for several probes from different imprinted regions, and it was generally found that [DNA replication](#) occurs earlier on the paternal homologue (34, 37-39). The patterns of allele-specific replication have also been investigated in the cells of PWS and AS patients (34, 37-40), with one important finding being that the DNA replication of various probes from the PWS/AS region is synchronous in the cells of PWS and AS UPD individuals (37-40), suggesting the requirement of a biparental contribution for the regulation of replication asynchrony (38). These findings led to the suggestion that allelic cross-talk might be involved in the normal pattern of imprinting. As several such cross-talk models of transcriptional regulation involve pairing of homologous chromosomes, the possibility that the maternal and paternal chromosome 15 homologues physically associate during somatic cell division was investigated.

Three-dimensional analysis of FISH in spatially intact nuclei was performed by laser scanning [confocal microscopy](#) to determine the distance between homologous chromosomes in cell cycle-fractionated human lymphocytes (41). The results demonstrated that in normal cells the chromosome 15 homologues pair during the late S phase. The most tightly paired region was observed to be 15q11-q13, consistent with the hypothesis that allelic cross-talk could be involved in the normal pattern of imprinting in chromosome 15q11-q13 (41). Homologous association was also observed at

the imprinted IGF2 locus on chromosome 11p15 (41), suggesting that the physical association of chromosome homologues is another manifestation of imprinting in normal human somatic cells and that regulation in *trans* could play an important role in the establishment and/or maintenance of imprinting.

The hypothesis that *trans*-sensing effects may play some role in the regulation of imprinting has been confirmed to some extent at the molecular level using mouse models, where imprinted genes have been modified or deleted by **gene-targeting** mutagenesis (42, 43). Crosses between mouse strains bearing different *H19* deletions revealed that transmission through the maternal germline affects the DNA methylation pattern of the paternal *Igf2* allele. In these crosses, where the maternal (active) *H19* allele is deleted, the paternal copy of *Igf2* is hypomethylated relative to that in wild-type animals (42). It has also been shown that imprinted expression of another gene, *Ins2*, which is tightly linked to *Igf2*, can be affected by the *trans* allele (43). In the latter experiments, the level of transcription of the maternal copy of a **reporter gene** (Lac Z) construct inserted at the *Ins2* locus depended on what was present on the paternal copy. Expression of the maternal copy of LacZ was repressed in the case where the paternal copy was an insertion of a different reporter (Neo), but not in mice where the paternal homologue is wild type. These results suggest that interaction between two alleles in *trans* can modulate both allele-specific transcription and DNA methylation and thus play a fundamental role in the establishment and/or maintenance of imprinting.

The properties of imprinted regions summarized above underscore the fundamental role of epigenetic modification of chromatin in the establishment and the maintenance of the imprint. There is, in addition, some molecular evidence for allele-specific changes in chromatin structure at imprinted loci. For the imprinted gene *U2af1-rs1* (see Table 1), there is a correlation between imprinted expression and an “open” or “active” chromatin expression, because the gene shows both expression and the presence of nuclease-hypersensitivity at sites in the promoter region of the paternal, but not the maternal, allele (44). No such allele-specific differences in nuclease sensitivity has been identified for the imprinted *Igf2* gene (45), however, suggesting that differential chromatin structure, at least as assayed by these approaches, is not a general feature of imprinted genes. Alternatively, these discrepant findings may reflect different mechanisms of imprinting for *U2af1-rs1* and *Igf2*. Further evidence for an association between imprinting and chromatin structure at the global level comes from the observation of sex-specific meiotic recombination frequencies in imprinted regions (46, 47).

### 3. A Catalog of Imprinted Genes and Their Characteristics

Although the occurrence of imprinting in mammals was demonstrated in 1984 (3, 4), one had to wait several years for the first imprinted genes to be described. As it is often the case in biological research, several papers, published at approximately the same time, reported imprinting of three genes, *Igf2r*, *Igf2*, and *H19* (48-50). These are the most studied of the imprinted genes in terms of understanding imprinting mechanisms (see text below). Since these initial reports, the number of genes identified as imprinted has continuously increased. There were eight reported by 1994 (51), and four years later more than 25 genes are identified as being monoallelically expressed in a parental-dependent way, at least at some point during development (Table 1). Due to limited space, these cannot be described in detail here, and only a few have been chosen as examples.

*Igf2r* was mapped to mouse chromosome 17, in a region known to contain the Tme locus, a **maternal-effect** mutation that results in the death of embryos at day 15 of gestation (52). **Northern blot** hybridizations of RNA from embryos having inherited a deletion of Tme from their mother or their father showed that only the maternal allele of *Igf2r* is expressed (48). Methylation analysis of a 130-kbp region spanning the *Igf2r* locus has identified two **CpG island** regions that display allelic methylation (26). Region 1, which corresponds to the 5' end of the transcription unit, is hypermethylated on the transcriptionally silenced, paternal allele. Region 2, localized in the second **intron** of the gene, shows the opposite pattern, with hypermethylation of the expressed, maternal allele. Only the region 2 hypermethylation profile is inherited from the gametes and remains constant

in all developmental stages and must therefore represent the imprint for the *Igf2r* gene. The functional role of this region in the imprinted expression of *Igf2r* was demonstrated by generating transgenic mice carrying a yeast artificial chromosome (YAC) from that region, which successfully reproduced methylation and expression of the endogenous *Igf2r* (53). Deletion of region 2 from these transgenes results in loss of imprinting and biallelic expression of *Igf2r*. Moreover, the region 2 CpG island was shown to be responsible for the expression of an **antisense** RNA (AS) whose expression profile is opposite to that of *Igf2r*. Only the paternal allele is normally transcribed and, in the case where *Igf2r* biallelic expression is associated with deletion of region 2, both alleles are repressed. It has been suggested that the antisense is the “imprintor” gene that regulates in *cis* the imprinted expression of *Igf2r* (54) (see text below).

Evidence for the imprinting of *Igf2* came initially from genetic studies of knockout mice generated to investigate the potential role of this gene during development. Heterozygotes having inherited the deletion from the father are smaller in size than normal animals and are indistinguishable in appearance from homozygotes (55). In contrast, when the mutation is transmitted through the mother, the heterozygotes are phenotypically normal. Molecular analysis of *Igf2* expression in genotyped embryos using ribonuclease protection assay subsequently confirmed that the paternal *Igf2* gene is transcriptionally active, whereas the maternal transcript is not detectable, except in choroid plexus and leptomeninges (49).

*H19* is an abundant RNA in the developing mouse embryo. Despite the absence of coding potential of *H19*, it is conserved significantly between human and mouse (56). Overexpression of *H19* in transgenic mice results in late prenatal lethality, suggesting that its dosage must be strictly controlled during development (57). In addition, *H19*, like *Igf2*, was mapped to mouse chromosome 7, in a region known to display an imprinting effect. Imprinted expression of *H19* was assessed using a ribonuclease protection assay that could distinguish between **polymorphic** alleles from different mouse subspecies. Only the maternal allele was shown to be expressed in liver and skeletal muscle of neonates and adults (50). Fine mapping subsequently localized *H19* less than 90 kbp from *Igf2*, with both genes having the same transcriptional orientation (58). Thus, both genes are physically clustered and imprinted, although in the opposite direction. In addition, *Igf2* and *H19* display a very similar, if not identical, pattern of expression during development (59). These findings led to the hypothesis that *H19* and *Igf2* operate from a common set of regulatory elements (60). Indeed, two endoderm-specific **enhancers**, located distal to *H19*, were shown to be responsible for the maternal expression of *H19* and the paternal expression of *Igf2* (61). The choice of the gene (*H19* or *Igf2*) to be activated by these enhancers is dictated by their methylation status. Both *H19* and *Igf2* have a region of hypermethylation on the paternal chromosome (23, 62). Only in the case of *H19*, however, is the methylation inherited from the gamete, suggesting a role as the imprint to distinguish the two alleles of *H19* (63). This paternal-specific *H19* methylation is thought to result in the strong preference for *Igf2* expression on the paternal chromosome. Methylation of *H19* on the paternal allele thus indirectly triggers the expression of *Igf2*. On the maternal chromosome, the preference for *H19* expression could well result from its closer proximity to the enhancers, a possibility that has been demonstrated by creating transgenic mice with duplicated and mislocated enhancers (64). *H19* can therefore be classified as a gene that is directly imprinted, whereas *Igf2* is imprinted indirectly, leading to the suggestion that the role of *H19* is to regulate the expression of *Igf2* and other flanking genes. This role for *H19* is analogous to that of an imprintor as mentioned above (54).

As shown above, initial reports of imprinted genes were based on genetic evidence: either the observation of a parental effect associated with a mutation (*Igf2*), or the mapping of a gene to a region known to display an imprinting effect (*Igf2r* and *H19*). This latter “strategy” has led to the identification of other clusters of imprinted genes, including that of mouse chr. 7/human chr. 11 *H19/Igf2* cluster and the mouse chr. 7/human chr. 15 cluster (Table 1).

The human 11p15 region, which contains *H19* and *IGF2*, is associated with the Beckwith–Wiedemann syndrome (BWS). *IGF2* has long been a strong candidate for this syndrome, since it is highly expressed in the tissues that are the most affected in BWS. Moreover, *IGF2* imprinting is lost



in 80% of BWS cases, resulting in biallelic expression of the gene (65), a finding that could explain the overgrowth characteristic associated with this syndrome. However, no chromosome rearrangement or mutation affecting the expression of *IGF2* alone has been identified in BWS patients, and the search for additional candidate gene(s) has led to a detailed transcription and imprinting map of the region. *p57<sup>KIP2</sup>*, an imprinted gene (66) that encodes a cyclin-dependant kinase (CDK) inhibitor, is the first gene shown to be mutated in BWS patients (67). *KVLQT1* was identified by positional cloning and mapped between *p57<sup>KIP2</sup>* and *IGF2*. *KVLQT1* is disrupted by chromosomal rearrangements in BWS patients (68) and is imprinted, with predominant maternal expression in most tissues. One exception is heart, where biallelic expression is observed. In several specimens, including different fetal tissues, however, monoallelic expression is not complete and, although the maternal allele is predominately expressed, residual expression from the paternal gene is detected. This “weak” or “leaky” imprinting is even more pronounced in the case of the mouse homologue *Kvlqt1*. While early expression of the gene is maternal in origin, this strong maternal bias is gradually lost in a tissue-specific manner as embryogenesis proceeds. Adult animals show complete biallelism (69). *IMPT1*, an other imprinted gene lying in the human 11p15.5/mouse distal 7 regions, also shows an allelic bias (“unequal allelic expression”) toward the maternal allele (70).

Incomplete imprinting has also been reported for *UBE3A*, the gene involved in Angelman syndrome. This gene displays highly tissue-specific imprinting, with predominant maternal expression in brain, but the paternal allele still expresses [TeXnical Error] of the total (71). In that case, this is likely to be due to the presence of different cell populations in the brain samples studied, one displaying complete monoallelic expression and a second, less abundant one, showing normal, biallelic expression; in mouse, *Ube3A* imprinting was shown to be restricted not only to brain but, more specifically, to certain types of neurons (72). Whether similar explanations are valid for most cases of incomplete imprinting remains to be determined. Most techniques used so far to study the imprinting status of genes, such as Northern blot, reverse transcriptase–polymerase chain reaction (RT-PCR), or ribonuclease protection, are performed on tissue samples or even whole individuals (embryos), which implies a mixture of several cell types. Analysis at the cellular level will require the use of more sophisticated technologies, like single-cell PCR or RNA-FISH. As suggested by the recent description of “weakly” imprinted genes, more detailed analyses of the expression pattern of genes may prompt us to reconsider our definition of imprinting. Will any gene showing monoallelic expression in a reasonable number of cells at any developmental stage or in any tissue be considered as imprinted? One important point to keep in mind, however, is that if imprinting implies monoallelic expression at some point during development or in some cell types, monoallelic expression doesn't necessarily indicate imprinting. This is evident from the observation that genes such as those encoding the olfactory receptors (73) and the lymphokine IL-2 (74) are monoallelically expressed but cannot, however, be considered as imprinted, because their expression can be from either the maternal or the paternal allele (random allelic expression).

It may also be very difficult to exclude that any gene is imprinted, because of the difficulty in analyzing the status of allelic expression in all cell populations within any tissue at any developmental stage. In other words, if the demonstration of nonrandom allelic expression is the molecular proof of imprinting, the “failure” to detect such an expression pattern is not proof of the absence of imprinting. One example that illustrates this problem is the Angelman syndrome gene, *UBE3A*. This gene was shown to be responsible for Angelman syndrome based on its chromosomal localization and the presence of point mutations in several patients. However, *UBE3A* displayed biallelic expression in fibroblasts and lymphocytes since the only cell types that had initially been tested, this gene was considered as not imprinted. It was shown only later that *UBE3A* is indeed imprinted, but that this is restricted to brain (71, 75) or, even more specifically, to certain types of neurons (72).

#### 4. Mechanisms of Imprinting

Imprinted genes appear to be preferentially organized into clusters, rather than dispersed throughout

the genome, although this may be biased by the fact that the search for imprinted genes has been mostly based on their localization within a region displaying an imprinting effect. This organization into clusters, along with additional molecular evidence, suggests that each imprinted gene is not regulated independently from the other member of the cluster and that [cis-acting](#) elements are crucial for their concerted regulation. Two main types of *cis*-elements will be described here, imprinting centers and imprintors, and their relation to each other will be discussed.

#### 4.1. Imprinting Center

As mentioned above, PWS and AS are two distinct syndromes that map to the chromosome 15q11-q13 region. Analysis of several patients carrying distinct rearrangements has allowed definition of two nonoverlapping regions responsible for these syndromes, with the AS region being smaller and lying distal to the PWS region ([76](#)). However, small deletions of less than 50 kbp within the PWS interval, in the 5' region of the paternally expressed *SNRPN* gene (Table [1](#)), were shown to be associated not only with PWS upon paternal transmission, but also with AS when transmitted maternally ([77](#)). These deletions affect the methylation profile and the gene expression pattern of the entire PWS region and, probably, although this has not yet been proven, the AS region. In such PWS patients, both chromosomes carry a maternal-specific methylation profile, and the paternally expressed genes of the region, such as *SNRPN*, *PAR-1*, and *PAR-5*, are transcriptionally silent ([78](#)). In AS patients, these genes are expressed biallelically ([79](#)) and the paternal methylation pattern is found on both chromosomes. The expression profile of the maternally expressed, AS gene *UBE3A* has not been tested, but it is expected to be silenced on both alleles. Because of these long-range effects, the deletions are thought to affect an imprinting control region, the imprinting center (IC), and are called “imprinting mutations.” Fine mapping of the IC has led to the identification of alternative *SNRPN* 5' exons, the BD exons, which are spliced to exon 2 of *SNRPN*, skipping the first exon and giving rise to another transcript, different from the *SNRPN* transcript but also paternally expressed ([80](#)).

In general, AS imprinting mutations result from the deletion of the BD exons, while PWS imprinting mutations are associated with deletion of at least *SNRPN* exon 1 ([80](#)), indicating that the IC has a bipartite structure. One model presented to explain how the IC regulates the imprinting process in PWS and AS invokes a mechanism of imprint switching ([77](#), [80](#)). In this model, the *SNRPN* exon 1 acts as an initiation site, functioning in the paternal germline to switch the chromosome epigenetic state from maternal to paternal. Mutation or deletion of this region causes a failure in the imprint switch during paternal gametogenesis, leading to PWS, because the paternally inherited chromosome has a maternal imprint. The other *cis* element of the bipartite IC would function via an interaction in *trans* with a female-specific factor in the maternal-to-paternal switch of epigenetic states. Mutation or deletion of this structure results in a failure of switching in the maternal germline and a chromosome 15 abnormality where the maternally transmitted chromosome carries a paternal imprint. An individual inheriting this chromosome would have AS due to the lack of a normal maternal contribution of the 15q11-q13 region. While the elements that control the imprint switching process remain to be identified, it has been proposed that the BD exons regulate imprint switching in the maternal germline ([80](#)).

It would appear that imprinting mutations may also be involved in the molecular pathogenesis of BWS ([36](#)). The evidence for this comes from a translocation family where *IGF2* is expressed biallelically, rather than paternally, although monoallelic DNA methylation of sites in the 3'-end of the gene is observed in most individuals and tissues. Such an uncoupling of allelic DNA methylation from transcription suggests that the disease results from an imprint mutation.

#### 4.2. Imprintors

The findings that only a subset of imprinted genes display a site(s) of allelic methylation that are inherited from the gametes has led to the concept of imprintor genes ([54](#)). Such genes inherit the allelic methylation signal from the gametes and thereby act in *cis* to control the expression of the other imprinted “target” genes in the domain by a mechanism of promoter/enhancer competition, as discussed above for the coordinate regulation of *H19* and *Igf2* ([60](#)), with *H19* acting as the imprintor

in this domain. This concept has also been invoked to explain the regulation of *Igf2r* and its antisense (53). In the latter case, the maternal allele-specific methylation of a CpG island within region 2 (see above) is inherited from the mother and persists in embryonic and adult tissues (26). This region 2 corresponds to a promoter for an antisense that is specifically transcribed from the paternal allele, but repressed on the maternal, strongly suggesting involvement of the antisense in the regulation of *Igf2r* imprinting (53). In the context of the model of Barlow (54), the antisense acts as the imprintor in the *IGF2r* domain.

There is also evidence to indicate that antisense may also play a role in imprinting at the AS locus due to the recent observation of a *UBE3A* antisense transcript (81). This antisense, which also encompasses another sense transcript that is downstream of *UBE3A*, displays the opposite pattern of allelic expression (paternal), and this only in brain. In other tissues where *UBE3A* is biallelically expressed, such as lymphocyte and fibroblast, the antisense is silenced, suggesting that expression of the two transcripts is mutually exclusive (81). These findings once again implicate a role for antisense in regulation of imprinted gene expression and suggest, in two cases, that the putative imprintor is an antisense, although this does not hold true for *H19*. One property shared by the *Igf2r* and *UBE3A* antisense transcripts and *H19* is that these do not have protein coding potential, which does not contradict their acting as imprintors, given that the sole function of an imprintor is to regulate the imprinted expression of other genes.

#### 4.3. Relation Between Imprinting Centers and Imprintors

There are two difficulties with the imprint switch model as it relates to the mechanism of action of the IC in the AS/PWS region. The first is that the BD exons display paternal allele-specific expression, making it difficult to understand how these might function to initiate switching in the maternal germline. Because of this puzzling observation and because these alternative *SNRPN* exons are not conserved in mouse, an alternative model involving competition for a common set of regulatory elements has been proposed to explain the mechanism of the IC (82).

The latter model was proposed initially to explain the opposite imprinting of *H19* and *Igf2* (60) (see above). In the case of the bipartite IC in the AS/PWS region, the element that overlaps the region spanning the BD exons elements is a *cis*-acting maternal epigenetic mark. This *cis*-acting site is extensively methylated in the female germline and represses the paternal genes via spreading of DNA methylation (82). Deletion of this *cis*-acting element would cause AS by transmission of an unmethylated maternal chromosome which thus carries a paternal epigenotype. The *SNRPN* promoter would somehow maintain the PWS region in an unmethylated state in the soma, and its deletion would result in a failure of the maintenance of this state and cause PWS. A neuronal enhancer would also be located in the PWS/AS region to explain the repression of maternally expressed genes, including *UBE3A*, on the maternal chromosome. *SNRPN* and *UBE3A* would compete for this enhancer, with access of *UBE3A* to this enhancer being regulated by the methylation status of the *SNRPN* promoter, which is normally paternally methylated (82). There is, however, evidence to suggest that the region spanning the BD exons is biallelically methylated (83) and that some sites in the PWS region are more methylated on the paternal chromosome (24, 77). Both these findings contradict the enhancer-competition model to explain imprinting in chromosome 15q11-q13 (83).

A different concept of the imprintor, which includes aspects of both models, has been presented to explain the action of the PWS/AS IC (54). On the one hand, the paternally expressed *SNRPN* transcripts encoded that contain the alternative BD exons could function as the imprintor expression competition with *UBE3A*. Alternatively, an imprintor gene encoded from within the AS critical region, but not from the IC, would compete for expression with *UBE3A* (54). The finding of a *UBE3A* antisense provides some support for this latter version of the imprintor model.

The arguments summarized above indicate that, despite much progress, the function and mechanism of action of imprinting centers is not yet understood. Neither the imprint switch, the imprintor, nor the original enhancer-competition model explain an important property of the IC—that is, how these

centers act over extensive (1–2 Mb) genomic regions. A mouse model, in which the region syntenic to the PWS-IC has been deleted, not only has some features of PWS, but also displays lack of expression of several paternally expressed genes, including *Ndn*, a gene that is about 1 Mb proximal of the IC (84). Further characterization of these mutant mice should be quite useful in further studies of the function of the IC during development. Finally, the observation that a fragment spanning the *SNRPN* exon 1 can function as a **silencer** element in transgenic *Drosophila* suggests that elements within the IC may regulate imprinting via a silencing mechanism that is evolutionarily conserved (85). It remains to be determined, however, if the silencing element identified by this approach can act over genomic distances in the Mbp size range, as would be predicted for the functional element of an IC.

## 5. Evolution and Function

Although evidence is accumulating that clearly implicates parental imprinting in developmental and pathological processes, its function in normal development raises more questions than provides answers. Several theories have been developed to explain the acquisition and, more importantly, the maintenance (fixation) of imprinting in mammals. One of these is the subtle and precise control on gene expression that imprinting allows during development. Imprinting has also been described as having evolved from a host defense mechanism (86) and, alternatively, as a surveillance mechanism for chromosome loss (87). However, these hypotheses are often dismissed, because there is no explanation for the fact that imprinting, as described here, is restricted to mammals. In contrast, the parental conflict theory is based on the hallmark manifestation of mammalian evolution—that is, the acquisition of a placenta, which mediates the interaction between the embryo and the mother during intrauterine growth (88). Following the latter theory, the evolutionary interest of the father is to promote growth of its own progeny by maximizing the amount of resource extracted from the mother. On the other hand, it is in the interest of the mother to maximize her total number of surviving offspring without showing favoritism toward the offspring of particular males. This is particularly true for offspring within litters from several fathers. This theory predicts that paternally controlled genes would favor fetal growth, whereas maternally controlled genes would tend to reduce offspring size. This is in agreement with the global observation that androgenetic embryos display abnormal extraembryonic tissues and that parthenogenetic embryos have a poorly developed embryonic portion. At the molecular level, the case of *Igf2/Igf2r* also favors this theory; *Igf2*, a paternally expressed gene, is involved in fetal growth, whereas *Igf2r*, which is maternally transcribed, is supposed to facilitate the degradation of *IL2*. A systematic analysis of the effects of UPDs on fetal growth has, however, shown several contradictions to the predictions of the conflict model (89).

As the number of imprinted genes identified will increase, theories proposed for the evolution of imprinting will probably multiply. The general picture emerging presently is that more genes than initially suspected are imprinted. Whether the imprinted status is “necessary” for their function remains, however, to be determined. Some can be imprinted as “bystanders” because of their location in an imprinted region. It also becomes more and more evident that imprinting is not an “all or nothing” phenomenon. It can affect genes with a high cellular specificity and/or only partially, with the “silent” allele still expressed, albeit at a low level. As a result, we might then consider imprinted expression to be the rule rather than an exception.

## Bibliography

1. P.-H. Gouyon, S. Maurice, X. Reboud, and I. Till-Bottraud (1993) *La Recherche* **24**, 70–76.
2. J. J. Bull (1983) *Evolution of Sex Determining Mechanisms*, Benjamin Cummings, Menlo Park, CA.
3. J. McGrath and D. Solter (1984) *Cell* **37**, 179–183.
4. M. A. H. Surani, S. C. Carton, and M. L. Norris (1984) *Nature* **308**, 548–550.
5. D. Solter (1988) *Annu. Rev. Genet.* **22**, 127–146.

6. B. M. Cattanach and J. Jones (1994) *J. Inherit. Metab. Dis.* **17**, 403–420.
7. S. D. Lawler, S. Povey, R. A. Fisher, and V. J. Pickthall (1982) *Ann. Hum. Genet.* **46**, 209–222.
8. M. Lalande (1997) *Annu. Rev. Genet.* **30**, 173–195.
9. C. Moutou, C. Junien, I. Henry, and C. Bonaiti-Pellie (1992) *J. Med. Genet.* **29**, 217–220.
10. D. Viljoen and R. Ramesar (1992) *J. Med. Genet.* **29**, 221–225.
11. I. Henry, C. Bonaiti-Pellie, V. Chehensse, C. Beldjord, C. Schwartz, G. Utermann, and C. Junien (1991) *Nature* **351**, 665–667.
12. M. G. Butler (1990) *Am. J. Med. Genet.* **35**, 319–332.
13. J. Clayton-Smith and M. E. Pembrey (1992) *J. Med. Genet.* **29**, 412–415.
14. C. Rougeulle and M. Lalande (1998) *Neurogenetics* **1**, 229–237.
15. M. G. Butler and C. G. Palmer (1983) *Lancet* **1**, 1285–1286.
16. J. H. Knoll, R. D. Nicholls, R. E. Magenis, J. M. Graham, M. Lalande, and S. A. Latt (1989) *Am. J. Med. Genet.* **32**, 285–290.
17. R. D. Nicholls, J. H. Knoll, M. G. Butler, S. Karam, and M. Lalande (1989) *Nature* **342**, 281–285.
18. S. Malcolm, J. Clayton-Smith, M. Nichols, S. Robb, T. Webb, A. L. Armour, A. J. Jeffreys, and M. E. Pembrey (1991) *Lancet* **337**, 694–697.
19. A. Razin and H. Cedar (1994) *Cell* **77**, 473–476.
20. R. Jaenisch (1997) *Trends Genet.* **13**, 323–329.
21. N. Rougier, D. Bourc'his, D. Molina Gomes, A. Niveleau, M. Plachot, A. Paldi, and E. Viegas-Péquignot (1998) *Genes Dev.* **12**, 2108–2113.
22. M. Brandeis, T. Kafri, M. Ariel, J. R. Chaillet, J. McCarrey, A. Razin, and H. Cedar (1993) *EMBO J.* **12**, 3669–3677.
23. A. C. Ferguson-Smith, H. Sasaki, B. M. Cattanach, and M. A. Surani (1993) *Nature* **362**, 751–755.
24. C. C. Glenn, S. Saitoh, M. T. Jong, M. M. Filbrandt, U. Surti, D. J. Driscoll, and R. D. Nicholls (1996) *Am. J. Hum. Genet.* **58**, 335–346.
25. R. Shemer, Y. Birger, A. D. Riggs, and A. Razin (1997) *Proc. Natl. Sci. USA* **94**, 10267–10272.
26. R. Stoger, P. Kubicka, C. G. Liu, T. Kafri, A. Razin, H. Cedar, and D. P. Barlow (1993) *Cell* **73**, 61–71.
27. E. Li, T. H. Bestor, and R. Jaenisch (1992) *Cell* **69**, 915–926.
28. E. Li, C. Beard, and R. Jaenisch (1993) *Nature* **366**, 362–365.
29. T. Caspary, M. A. Cleary, C. C. Baker, X.-J. Guan, and S. M. Tilghman (1998) *Mol. Cell. Biol.* **18**, 3466–3474.
30. R. Drouin, G. P. Holmquist, and C.-L. Richer (1994) *Adv. Hum. Genet.* **22**, 47–115.
31. Y. Izumikawa, K. Naritomi, and K. Hirayama (1991) *Hum. Genet.* **87**, 1–5.
32. M. S. Lin, A. Zhang, and A. Fujimoto (1995) *Hum. Genet.* **96**, 572–576.
33. R. Drouin, M. Boutouil, R. Fetni, G. P. Holmquist, P. Scott, C.-L. Richer, and N. Lemieux (1997) *Chromosoma* **106**, 405–411.
34. D. Kitsberg, S. Selig, M. Brandeis, I. Simon, I. Keshet, D. J. Driscoll, R. D. Nicholls, and H. Cedar (1994) *Nature* **364**, 459–463.
35. O. W. Smrzka, I. Fae, R. Stoger, R. Kurzbauer, G. F. Fischer, T. Henn, A. Weith, and D. P. Barlow (1995) *Hum. Mol. Genet.* **4**, 1945–1952.
36. K. W. Brown, A. J. Villar, W. Bickmore, J. Clayton-Smith, D. Catchpoole, E. R. Maher, and W. Reik (1996) *Hum. Mol. Genet.* **5**, 2027–2032.
37. J. H. Knoll, S. D. Cheng, and M. Lalande (1994) *Nat. Genet.* **6**, 41–46.
38. J. M. LaSalle and M. Lalande (1995) *Nat. Genet.* **9**, 386–394.

39. P. H. Gunaratne, M. Nakao, D. H. Ledbetter, J. S. Sutcliffe, and A. C. Chinault (1995) *Genes Dev.* **9**, 808–820.
40. L. M. White, P. K. Rogan, R. D. Nicholls, B.-L. Wu, B. Korf, and J. H. M. Knoll (1996) *Am. J. Hum. Genet.* **59**, 423–430.
41. J. M. LaSalle, and M. Lalande (1996) *Science* **272**, 725–728.
42. T. Forné, J. Oswald, W. Dean, J. R. Saam, B. Bailleul, L. Dandolo, S. M. Tilghman, J. Walter, and W. Reik (1997) *Proc. Natl. Acad. Sci. USA* **94**, 10243–10248.
43. B. Duvillié, D. Bucchini, T. Tang, J. Jami, and A. Paldi (1998) *Genomics* **47**, 52–57.
44. R. Feil, M. D. Boyano, N. D. Allen, and G. Kelsey (1997) *J. Biol. Chem.* **272**, 20893–20900.
45. R. Feil, M. A. Handel, N. D. Allen, and W. Reik (1995) *Dev. Genet.* **17**, 240–252.
46. W. P. Robinson, and M. Lalande (1995) *Hum. Mol. Genet.* **4**, 801–806.
47. A. Paldi, G. Gyapay, and J. Jami (1995) *Curr. Biol.* **5**, 1030–1035.
48. D. P. Barlow, R. Stoger, B. G. Herrmann, K. Saito, and N. Schweifer (1991) *Nature* **349**, 84–87.
49. T. M. DeChiara, E. J. Roberston, and A. Efstratiadis (1991) *Cell* **64**, 849–859.
50. M. S. Bartolomei, S. Zemel, and S. M. Tilghman (1991) *Nature* **351**, 153–155.
51. M. S. Bartolomei (1994) *Nature Genet.* **6**, 220–221.
52. K. Artzt, D. P. Barlow, W. Dove, K. Fisher-Lindahl, M. Lyon, J. Klein, and L. M. Silver (1991) *Mamm. Genome* **1**, S281–S301.
53. A. Wutz, O. W. Smrzka, N. Schweifer, K. Schellander, E. F. Wagner, and D. P. Barlow (1997) *Nature* **389**, 745–749.
54. D. P. Barlow (1997) *EMBO J.* **16**, 6899–6905.
55. T. M. DeChiara, A. Estratiadis, and E. J. Robertson (1990) *Nature* **345**, 78–80.
56. C. I. Brannan, E. Claire Dees, R. S. Ingram, and S. M. Tilghman (1990) *Mol. Cell. Biol.* **10**, 28–36.
57. M. E. Brunkow and S. M. Tilghman (1991) *Genes Dev.* **5**, 1092–1101.
58. S. Zemel, M. S. Bartolomei, and S. M. Tilghman (1992) *Nat. Genet.* **2**, 61–65.
59. R. Ohlsson, F. Hedborg, L. Holmgren, C. Walsh, and T. J. Ekstrom (1994) *Development* **120**, 361–368.
60. M. S. Bartolomei, A. L. Webber, M. E. Brunkow, and S. M. Tilghman (1993) *Genes Devel.* **7**, 1663–1673.
61. P. A. Leighton, J. R. Saam, R. S. Ingram, C. L. Stewart, and S. M. Tilghman (1995) *Genes Dev.* **9**, 2079–2089.
62. T. Sasaki, R. Scott Hansen, and S. M. Gartler (1992) *Mol. Cell. Biol.* **12**, 3819–3826.
63. K. D. Tremblay, J. Saam, R. S. Ingram, S. M. Tilghman, and M. S. Bartolomei (1995) *Nature Genet.* **9**, 403–413.
64. A. L. Webber, R. S. Ingram, J. M. Levorse, and S. M. Tilghman (1998) *Nature* **391**, 711–715.
65. R. Weksberg, S. D. R., F. Y. L., S. Q. L., and S. J. (1993) *Nat. Genet.* **5**, 143–150
66. I. Hatada, J. Inazawa, T. Abe, M. Nakayama, Y. Kaneko, Y. Jinno, N. Niikawa, H. Ohashi, Y. Fukushima, K. Iida, C. Yutani, S.-i. Takahashi, Y. Chiba, S. Ohishi, and T. Mukai (1996) *Hum. Mol. Genet.* **5**, 783–788.
67. I. Hatada, H. Ohashi, Y. Fukushima, Y. Kaneko, M. Inoue, Y. Komoto, A. Okada, S. Ohishi, A. Nabetani, H. Morisaki, M. Nakayama, N. Niikawa, and T. Mukai (1996) *Nature Genet.* **14**, 171–173.
68. M. P. Lee, R.-J. Hu, L. A. Johnson, and A. P. Feinberg (1997) *Nature Genet.* **15**, 181–185.
69. T. D. Gould and K. Pfeifer (1998) *Hum. Mol. Genet.* **7**, 483–487.
70. D. Dao, D. Frank, N. Qian, D. O'Keefe, R. J. Vosatka, C. P. Walsh, and B. Tycko (1998) *Hum. Mol. Genet.* **7**, 597–608.

71. C. Rougeulle, H. Glatt, and M. Lalande (1997) *Nature Genet.* **17**, 14–15.
72. U. Albrecht, J. S. Sutcliffe, B. M. Cattanach, C. V. Beechey, D. Armstrong, G. Eichele, and A. L. Beaudet (1997) *Nature Genet.* **17**, 75–78.
73. A. Chess, I. Simon, H. Cedar, and R. Axel (1994) *Cell* **78**, 823–834.
74. G. A. Holländer, S. Zuklys, C. Morel, E. Mizoguchi, K. Mobisson, S. Simpson, C. Terhorst, W. Wishart, D. E. Golan, A. K. Bhan, and S. J. Burakoff (1998) *Science* **279**, 2118–2121.
75. T. H. Vu and A. R. Hoffman (1997) *Nature Genet.* **17**, 12–13.
76. J. S. Sutcliffe, Y.-h. Jiang, R.-J. Galjaard, T. Matsuura, P. Fang, T. Kubota, S. L. Christian, J. Bressler, B. Cattanach, D. H. Ledbetter, and A. L. Beaudet (1997) *Genome Res* **7**, 368–377.
77. K. Buiting, S. Saitoh, S. Gross, B. Dittrich, S. Schwartz, R. D. Nicholls, and B. Horsthemke (1995) *Nat. Genet.* **9**, 395–400.
78. J. S. Sutcliffe, M. Nakao, S. Christian, K. H. Orstavik, N. Tommerup, D. H. Ledbetter, and A. L. Beaudet (1994) *Nat. Genet.* **8**, 52–58.
79. S. Saitoh, K. Buiting, P. K. Rogan, J. L. Buxton, D. J. Driscoll, J. Arnemann, R. König, S. Malcolm, B. Horsthemke, and R. D. Nicholls (1996) *Proc. Natl. Acad. Sci. USA* **93**, 7811–7815.
80. B. Dittrich, K. Buiting, B. Korn, S. Rickard, J. Buxton, S. Saitoh, R. D. Nicholls, A. Poustka, A. Winterpacht, B. Zabel, and B. Horsthemke (1996) *Nature Genet.* **14**, 163–170.
81. C. Rougeulle, C. Cardoso, M. Fontes, L. Colleaux, and M. Lalande (1998) *Nature Genet.* **19**, 15–16.
82. S. M. Tilghman, T. Caspary, and R. S. Ingram (1998) *Nature Genet.* **18**, 206–208.
83. A. Schumacher, K. Buiting, M. Zeschnigk, W. Doerfler, and B. Horsthemke (1998) *Nature Genet.* **19**, 324–325.
84. T. Yang, T. E. Adamsom, J. L. Resnick, S. Leff, R. Wevrick, U. Francke, N. A. Jenkins, N. G. Copeland, and C. I. Brannan (1998) *Nature Genet.* **19**, 25–31.
85. F. Lyko, K. Buiting, B. Horsthemke, and R. Paro (1998) *Proc. Natl. Acad. Sci. USA* **95**, 1698–1702.
86. D. P. Barlow (1993) *Science* **260**, 309–310.
87. J. H. Thomas (1995) *Proc. Natl. Acad. Sci. USA* **92**, 480–482.
88. T. Moore and D. Haig (1991) *Trends Genet.* **7**, 45–49.
89. L. D. Hurst and G. T. McVean (1997) *Trends Genet.* **13**, 436–443.

## ***In Situ* Hybridization**

The spatial distribution of a target nucleic acid in a tissue, cell, or [chromosome](#) can be visualized by [in situ hybridization](#) techniques. As with other hybridization techniques, there are many variations. Each procedure shares a set of common steps: (1) fixing the sample to preserve its structure, (2) limited proteinase digestion that provides access to the hybridization probe, (3) addition of the hybridizing probe, (4) washing to remove unhybridized probe, and (5) detection of hybrids.

Small probes are usually preferred over larger probes for *in situ* hybridization. The probe must **diffuse** efficiently into the sample so that the hybridization reaction can occur. The advantage imparted by faster diffusion of small probes is partially countered by their tendency to be lost in the washing step. Because the stability of the hybridized duplex depends on concentration for short

oligonucleotides, the dilution that accompanies the washing step can result in release of some of the probe from the hybrid duplexes.

Visualization is difficult when a small number of target molecules are present. This limited most of the applications of the method to hybridization of [messenger RNA](#). Low copy-number DNA and RNA are, however, physiologically important. Viral DNA also may be present in cells in low copy number. The need to increase the sensitivity of the technique to permit visualization of low copy-number targets has led to incorporation of amplification techniques into *in situ* hybridization techniques. Several amplification techniques have been employed, with polymerase chain reaction (**PCR**) being the most widely applied.

[Radioisotope](#) or **fluorescence** detection of the hybrid duplex can be employed, but the current method of choice is bioluminescence or [chemiluminescence](#). A [biotin](#) label is detected by [streptavidin](#) coupled to [alkaline phosphatase](#) or horseradish [peroxidase](#). These enzymes catalyze several reactions that can begin a cascade of reactions leading to production of light. Because biotin is a natural metabolite, high background signals may be observed in some tissues. One of the many alternatives is a digoxigenin label that is recognized and bound by a specific [antibody](#) linked to alkaline phosphatase or horseradish peroxidase. Counter staining permits identification of cellular structures and, in turn, observation of the distribution of hybrid duplexes within the cell or tissue.

#### Suggestions for Further Reading

- D. G. Wilkinson, ed. (1993) *In situ Hybridization: A Practical Approach*, IRL Press, New York.
- J. Gu, ed. (1995) *In situ Polymerase Chain Reaction and Related Technology*, Eaton Press, Natick, MA.
- L. Kricka, ed. (1992) *Nonisotopic DNA Probe Techniques*, Academic Press, San Diego.

## Inborn Errors of Metabolism

Inborn errors of metabolism is the term applied to genetic disorders caused by loss of function of an enzyme. Enzyme activity may be low or lacking for a variety of reasons, such as lack of transcription and translation of the corresponding gene, improper folding of the polypeptide, unstable tertiary structure, a defective active site, failure to bind the corresponding coenzyme under physiological conditions, or failure to enter into complexes with other proteins. A single good copy of the gene is often sufficient, and many inborn errors are therefore recessive traits. The severity of the phenotype is highly variable among inborn errors, depending on the role of the biochemical pathway in metabolism. A few inborn errors can be treated by diet. Most cannot. Inborn errors were the first genetic disorders to be considered for gene therapy.

### 1. Early Recognition of Inborn Errors

The phrase *inborn errors of metabolism* appears to have been first used by A.E. Garrod in the Croonian Lectures of 1908, published as a monograph with that title in 1909. In 1902, Garrod published an account of a patient with alkaptonuria and speculated that it was due to an inherited derangement of metabolism. Subsequently, he encountered other rare patients who suffered from apparent metabolic defects. The fact that the parents of such patients were often related to each other suggested that the metabolic defects might be inherited as recessive traits. In his 1909 publication, he suggested four candidate disorders as examples of inborn errors: alkaptonuria, albinism, pentosuria, and cystinuria. Although the sciences of genetics and biochemistry were in their infancies at the



time, Garrod was remarkably perceptive in recognizing the possibility of inherited variations in metabolism.

Garrod's prescient observations were largely ignored except as a footnote in the biochemical literature. The geneticists were not ready for biochemistry, nor were biochemists quite ready to deal with genetics. Enzymes were just being identified as the important agents in metabolism, and little was known about the enzymes themselves or other protein functions. Nothing was known about the chemical nature of genes or how they act. It was not until the late 1930s that geneticists began to employ biochemical approaches to study what genes do. Soon thereafter, G.W. Beadle, E.L. Tatum, and others, based primarily on studies of *Neurospora*, observed that each mutant form studied appeared to be defective in a specific enzymatic conversion. Their work led to the *one gene-one enzyme* hypothesis, which equates the action of each genetic locus with the activity of a single enzyme.

We know now that many proteins are not enzymes. Also, many enzymes are complexes of different proteins. (Some RNA molecules have enzyme activity—the ribozymes—but these will not be considered here and are not active in metabolic pathways as that term is usually defined.) We also know now that genes are composed of DNA and their function is to code for the primary structure of proteins. The inborn errors of metabolism envisioned by Garrod are readily explained by the role of genes in coding for specific enzymes or enzyme subunits. Errors in metabolism reflect loss of activity of a particular enzyme due to an error in a gene that codes for that enzyme. The terms *metabolic block* and *metabolic error* are used as synonyms for inborn errors of metabolism.

## 2. Some Properties of Enzymes

It is helpful to consider how enzymes function if one is to understand the many variations that lead to malfunction. First, the polypeptide chain that is coded by a gene must be folded into the correct three-dimensional structure. Protein folding is a consequence of the amino acid sequence and the interaction of the side chains with each other and with water. Many detrimental mutations act by substituting one amino acid for another, interfering with the electrostatic and hydrophobic interactions of the functional enzyme. In the case of glucose 6-phosphate dehydrogenase deficiency, hundreds of mutant forms have been identified. A common feature of many is decreased stability of the folded enzyme.

To function, enzymes must have an active site that determines substrate access and binding. They often require coenzymes, small organic molecules that bind and help carry out the chemical transformations at the active site of the enzyme. Some coenzymes are made by the host organism. Others, of which the B vitamins are examples, must be supplied in the diet in higher animals. The coenzymes must bind at specific sites on the apoenzymes to form active holoenzyme. Mutation at those sites can reduce or abolish the affinity of the enzyme for the coenzyme. In some instances when the affinity is reduced, increase in the coenzyme level can compensate. For example, homocystinuria results from a defect in cystathionine synthase, which catalyzes the conversion of homocysteine to cystathionine. The homocysteine substrate accumulates and is excreted in urine as homocysteine. Some cases are treatable with large doses of vitamin B<sub>6</sub>, which is converted to the coenzyme pyridoxal phosphate and increases the level of active enzyme. In other patients, the molecular defect is in another part of the enzyme, or the enzyme may be missing. Vitamin B<sub>6</sub> is of no value in those cases. The two types of cases are not distinguishable clinically apart from the response to therapy.

Coenzymes that are synthesized in the body may also be absent due to defects in their synthesis. This causes loss of activity of the enzymes of which they are a normal part. An example is found in phenylketonuria (PKU). PKU is a classic inborn error of metabolism usually caused by mutation of the phenylalanine hydroxylase gene. This enzyme uses tetrahydrobiopterin as a coenzyme. In a small proportion of cases, the genetic defect is in synthesis of tetrahydrobiopterin. The end result is the

same. Phenylalanine cannot be broken down by the normal pathway and accumulates, causing the classic features of PKU.

Many enzymes and other proteins function as dimers, higher multimers, or part of a multisubunit complex. This imposes a requirement that the monomeric form have a configuration that permits complex formation. Changes in the amino acid sequence can alter the ability of a monomeric enzyme to enter into complexes, either with itself to form an active dimer or with other subunits to form a complex. The mutant subunit might still have activity under some circumstances, but the conditions under which the substrate is encountered inside a cell may be very different.

This list of protein functions is not exhaustive. Many properties of an enzyme can be affected by structural changes. Any change that substantially or completely blocks the activity so that a metabolic conversion does not occur can be the cause of an inborn error of metabolism.

### 3. Metabolic Errors and Phenotypes

Although some inborn errors are relatively benign, many are quite detrimental, often fatal. The impact of a metabolic error depends greatly on the metabolic pathway and on the ability of other metabolic and physiological systems to compensate for the error. Alkaptonuria is relatively benign. The pathway that is blocked is in the catabolic breakdown of homogentisic acid, and no essential products are missing because of this block. Homogentisic acid is not reabsorbed by the kidneys and is excreted in the urine, preventing accumulation of high levels in the body. Over time, oxidation products of homogentisate are deposited in joints, causing a form of arthritis. In contrast, the block in phenylketonuria is in the conversion of phenylalanine to tyrosine. Phenylalanine is efficiently reabsorbed by the kidneys and accumulates to high levels in the blood, generating by-products that interfere with development of the central nervous system.

Several blocks in the breakdown of gangliosides are responsible for degenerative diseases of the central nervous system (e.g., Tay-Sachs disease). In these instances, breakdown occurs in the lysosomes. The large gangliosides accumulate in the lysosomes, which become very large, eventually interfering with cell function and ultimately causing death. Such disorders are examples of lysosomal storage diseases. A characteristic of these disorders is an initially normal phenotype, followed by degeneration as the lysosomes store increasing amounts of material.

Another group of metabolic blocks depends on environmental agents for expression. Galactosemia is one of the earliest known examples. Affected persons cannot metabolize galactose, which comes from lactose in milk. The galactose 1-phosphate that accumulates interferes with other metabolic processes and can lead to death unless milk is avoided. Other rare individuals are homozygous for absence of functional plasma cholinesterase. If succinylcholine is administered as a muscle relaxant during electroshock or other treatment, they have very prolonged recovery of muscle function.

### 4. Inborn Errors and Genetics

In most cases, if a locus that codes for an enzyme has one normal allele, a normal phenotype results; it is haplosufficient. Two nonfunctional alleles are required for the metabolic error to occur. This means that the condition is inherited as a recessive trait. If  $q$  is the frequency of the mutant allele, the frequency of homozygotes would be  $q^2$ , and the frequency of heterozygous "carriers" would be  $2pq$ . For  $q = 0.01$ , the frequency of homozygotes therefore is  $q^2 = 0.0001$  and  $2pq = 0.0198$ . The bulk of the recessive alleles in the population are in heterozygous combination and not subject to elimination by selection. Some genes that code for enzymes are haploinsufficient. The mutant allele is associated with a modified phenotype and is transmitted therefore in pedigrees as a dominant trait. The phenotype of the heterozygote need not be identical to the homozygous mutant phenotype, which might be more severe.

With the advent of prenatal diagnosis for nearly all metabolic errors, once the gene is identified, the frequency of affected children born can be substantially reduced if a mating is known to be at risk. Identification of such matings usually occurs because of a prior pregnancy that involved an affected offspring. Without massive population screening of normal persons, however, the frequency of matings at risk will not change.

## 5. Treatment of Inborn Errors

Reference has been made to treatment of phenylketonuria and galactosemia by elimination of components of diet. Some forms of homocystinuria respond to high levels of dietary vitamin B<sub>6</sub>. The number of inborn errors that are managed by diet is small. Because inborn errors are usually associated with loss of activity of a single enzyme, the idea of enzyme replacement has been an enticing goal for many years. How to do this technically is a major challenge. The principle has been demonstrated in a sense in a boy with phenylketonuria who received a liver transplant for reasons unrelated to the PKU. The condition was cured. Because PKU is effectively treated by diet, liver transplants are not considered an appropriate therapy.

There are a number of inborn errors of great severity that are candidates for gene therapy. The first attempts at gene therapy involved a patient with adenosine deaminase deficiency, which causes severe immunodeficiency. The gene is well characterized and can be incorporated into viral vectors. The results were very promising, but a satisfactory protocol has yet to be achieved. Another candidate disorder is Lesch-Nyhan syndrome, which is due to loss of activity of hypoxanthine phosphoribosyl transferase (hprt). Patients are mentally retarded and engage in self-mutilation. A very low level of hprt activity is sufficient to prevent the disorder. Gene therapy need not be very efficient in order to be effective. Success has been elusive, however.

Despite the disappointments in early attempts at gene therapy, it remains the most promising approach for metabolic errors as well as for certain other genetic disorders.

### Additional Reading

Harris H., *Garrod's Inborn Errors of Metabolism*, Oxford University Press, London, U.K., 1963, xi+207 pp. Included are reprints of the original 1902 and 1909 publications of Garrod.

Scriber C.R., Beaudet A.L., Sly W.S., Valle D., eds., *The Metabolic and Molecular Bases of Inherited Disease*, 7th ed., McGraw-Hill, New York, 1995.

Strachan T. and Read A.P., *Human Molecular Genetics*, 2nd ed., Wiley-Liss, New York, 1999.

*On-Line Mendelian Inheritance in Man*: <http://www3.ncbi.nlm.nih.gov/Omim/>. This database can be searched for any inborn error of metabolism. It provides an up-to-date description of research and links to related sites.

## Inclusion Bodies

Inclusion bodies are insoluble, biologically inactive aggregates of partially folded **protein**, which are often produced when a heterologous or [recombinant protein](#) is overexpressed, especially in the popular and frequently used *Escherichia coli* bacteria (see [Expression Systems](#)). Inclusion body formation is frequently seen in both **Gram-negative** and Gram-positive bacteria, but less often in yeast, insect, or mammalian cells. Inclusion bodies are also called *refractile bodies* because they

refract light and are visible in the phase-contrast **microscope**. Usually, there is only one inclusion body per cell, which can span the whole diameter of the cell, depending on the expression level of the protein. Inclusion bodies can occur in both the cytosolic and **periplasmic** location in bacteria. If the initial expression level is sufficiently high, the inclusion bodies tend to be a rather homogeneous aggregate of the overexpressed protein. However, several other proteins are also found as coprecipitants, such as [elongation factor](#) EF-Tu, outer membrane proteins, and small **heat-shock** proteins.

Systematic analysis of inclusion body formation shows no correlation with the size, relative [hydrophobicity](#), or number of [proline](#) residues (see [Cis/Trans Isomerization](#)) in the protein. However, the relative charge and the number of amino acid residues that prefer the [turn](#) conformation correlate reasonably well with the tendency for inclusion body formation. Also, there is an excellent correlation between inclusion body formation and high expression level of the protein. Even native, cytosolic proteins can form inclusion bodies at a sufficiently high expression level. A [kinetic](#) model has been suggested that takes these observations into account. It states that inclusion body formation is dependent on the relative rates of **protein biosynthesis**, **protein folding**, and aggregation.

Overexpression of a particular protein can be deleterious to cell viability (1). The protein may have a biochemical activity that is directly toxic, or the overexpression state itself is toxic. This is marked by an inability of the cell to recover after a period of overexpression. One major contributing factor is the observed destruction of [ribosomes](#) (2). Finally, a word of caution. The most abundant *E. coli* protein is elongation factor EF-Tu, reaching 5% to 10% of the total cellular protein. Any protein that is expressed at higher levels challenges the biosynthetic capacity of the cell, especially if it is made from a gene with a high portion of rare **codons**. The most negative consequences are less growth capacity and high error levels in translation.

## 1. How Can the Formation of Inclusion Bodies Be Avoided?

In most cases, the native, soluble form of a protein is desired in high quantity and purity. It is therefore important to avoid inclusion body formation, which requires that the protein be solubilized, generally in a **denatured** state, and then folded. Furthermore, the quality of protein folded from inclusion bodies can be inferior (3).

There are several ways to avoid inclusion body formation. First, fermentation parameters can be varied, such as the richness of the medium, its pH, and the level of aeration or micronutrients. A major improvement is often seen when the growth temperature is lowered. Second, the expression level of the protein can be varied. This can be done by selecting an appropriate expression **vector** with sufficiently weak **promoter** and/or [translation](#) initiation elements. Also, an expression vector that allows regulated but variable steady-state levels of gene expression can be used (4). However, most overexpression systems are inducible, so that varying the level of **inducer** may suffice. Third, coexpression of protein folding modulators, such as the **molecular chaperones** GroES/EL can help. Coexpression of the **DnaK/DnaJ** chaperones may help when GroES/EL is ineffective. Also, coexpression of *E. coli* [thioredoxin](#) can counteract inclusion body formation (3). Other folding modulators of potential use are [peptidyl-prolyl cis/trans isomerases](#) and thiol/disulfide oxidoreductases. Fourth, genetic engineering can improve soluble expression. The simplest solution is to fuse the protein of interest to other proteins that are known to improve the yield of soluble protein, such as the bacterial thioredoxin and IgG binding domains of *S. aureus* [protein A](#) (see [Fusion Gene, Fusion Protein](#)). A more sophisticated approach involves the rational design of higher solubility by the introduction of specific point [mutations](#) by [site-directed mutagenesis](#), and elimination of amino acid sequences involved in aggregation. Finally, a **signal sequence** can cause the transport of the protein of interest to the **periplasm** or to the extracellular medium, where the environment favors the folding of some proteins.

## 2. How Can Inclusion Bodies Be Used?

Even inclusion body formation can be turned into an advantage so long as the protein can be refolded. The primary advantage is the ease of purification of inclusion bodies; those with a recombinant protein fraction of more than 60% can be isolated by centrifugation of crude cell lysate at moderate rotor speeds (5000 to 12,000 g). The purity of an inclusion body protein can be improved further by resuspension and centrifugation in a buffer containing 2% sodium [deoxycholate](#) (5). Inclusion bodies are often the only way to obtain a protein that is highly toxic in its native form. Any protein can usually be isolated via inclusion bodies if it is fused to an amphiphilic extension (6) or if it is expressed at a sufficiently high level (eg, using T7 promoter vectors). Also, host strains with a defective **heat-shock** response or with a mutation in an important molecular chaperone system may be useful. Inclusion bodies are often an advantage when the yield of the native protein is low because of extensive **proteolysis**.

The key to efficient use of inclusion bodies is to know how to solubilize and refold the protein in high yield. The strong denaturants 4 to 6 M **guanidinium chloride** or 6 to 10 M [urea](#) are the most common solubilizing reagents. Another method uses [detergents](#), such as cetylmethylammonium chloride or [sarkosyl](#), or alkaline solution (7). Often, it is necessary to use **thiol** reagents along with the **chaotropes** to facilitate reduction of nonnative [disulfide bonds](#).

All methods for folding the solubilized and denatured protein are based on the removal or dilution of the denaturant and reducing agent. This can be done by one-step dilution, by [dialysis](#) or, as when sarkosyl is used, by the addition of a micelle-forming detergent (eg, **Triton X-100**) that sequesters the protein-bound sarkosyl (8).

A major drawback to most *in vitro* folding approaches is the presence of aggregation side reactions, which compete with the proper folding reaction. It is therefore often necessary to refold at low protein concentration or to use pulse renaturation techniques (9). Furthermore, low temperatures or the composition of the refolding buffer may affect the yield of native protein. A major improvement can be observed when the denatured protein is refolded on a solid support, which effectively prevents aggregation by protein-protein contacts (10, 11).

Molecular chaperones, such as GroES/EL, have also been used for *in vitro* refolding of denatured proteins. A promising new development is the advent of refolding chromatography, where immobilized chaperones or fragments thereof are used for refolding of protein. Elaborate “artificial” chaperone systems have also been constructed, applying nonprotein refolding agents or changes in the physical environment (6). A combination of detergents and [cyclodextrins](#) has been used to mimic the two-step action of GroES/EL in protein folding. Systems have also been made that are based on temperature-jump techniques or refolding on a solid support.

## Bibliography

1. B. Miroux and J. E. Walker (1996) *J. Mol. Biol.* **260**, 289–298.
2. H. Dong, L. Nilsson, and C. G. Kurland (1995) *J. Bacteriol.* **177**, 1497–1504.
3. T. Yasukawa et al. (1995) *J. Biol. Chem.* **270**, 25328–25331.
4. L. M. Guzman, D. Belin, M. J. Carson, and J. Beckwith (1995) *J. Bacteriol.* **177**, 4121–4130.
5. L. Rao et al. (1996) *Anal. Biochem.* **241**, 173–179.
6. J. G. Thomas, A. Ayling, and F. Baneyx (1997) *Appl. Biochem. Biotech.* **66**, 197–238.
7. J. Suttner et al. (1994) *J. Chromatog. B* **656**, 123–126.
8. J. V. Frangioni and B. G. Neel (1993) *Anal. Biochem.* **210**, 179–187.
9. R. Rudolph (1996) In *Protein Engineering: Principles and Practice* (J. L. Cleland and C. S. Craik, eds.), Wiley-Liss, Inc., New York, pp. 283–298.
10. T. E. Creighton (1986) In *Protein, Structure, Folding and Design*, (D. Oxender, ed.), Alan R Liss, New York, pp. 249–257.
11. A. Holzinger, K. S. Phillips, and T. E. Weaver (1996) *Biotechniques* **20**, 804–808.

### Suggestions for Further Reading

12. R. S. Donovan, C. W. Robinson, and B. R. Glick (1996) Optimizing inducer and culture conditions for expression of foreign proteins under the control of the *lac* promoter. *J. Industrial Microbiol.* **16**, 145–154.
13. G. Georgiou (1996) "Expression of proteins in biotechnology." In *Protein Engineering: Principles and Practice* (J. L. Cleland and C. S. Craik, eds.), Wiley-Liss, Inc., New York, pp. 101–127.
14. C. Kurland and J. Gallant (1996) Errors of heterologous protein expression. *Curr. Op. Biotech.* **7**, 489–493.
15. R. Rudolph (1996) "Successful protein folding on an industrial scale". In *Protein Engineering: Principles and Practice* (J. L. Cleland and C. S. Craik, eds.), Wiley-Liss, Inc., New York, pp. 283–298.
16. R. Rudolph and H. Lilie (1996) *In vitro* folding of inclusion body proteins. *FASEB J.* **10**, 49–56.
17. J. G. Wall and A. Plückthun (1995) Effects of overexpressing folding modulators on the *in vivo* folding of heterologous proteins in *Escherichia coli*. *Curr. Op. Biotech.* **6**, 507–516.

### Indel

An “indel” is a site of insertion or deletion in aligned, homologous protein or nucleotide sequences (1). Mutational events in nucleic acid sequences include insertions and deletions, as well as base substitutions. Often, it is not possible to distinguish between insertional or deletional events, because this would require either access to the (presumably dead) common ancestor of the aligned sequences or else a sufficient number of related sequences to reveal the nature of a particular event in a recently derived lineage. In globular proteins, indels are largely restricted to sites (usually loops) at the protein surface (2). A gapped region of a sequence alignment may actually have a history of multiple insertions and deletions, and it is then a fruitless exercise to try to retrace the actual events. Indel is therefore a useful term as it avoids the need to make an error-prone assignment of the actual mutational event(s). The most sensitive algorithms for [aligning sequences](#) allow for the possibility of indels being placed at any position in any sequence, using [gap penalties](#) to apply a cost to indel insertion, to limit their frequency (3, 4).

### Bibliography

1. D. Sankoff and J. B. Kruskal (1984) *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Addison-Wesley, Reading, MA, p. 11.
2. S. Pascarella and P. Argos (1992) *J. Mol. Biol.* **224**, 461–471.
3. S. B. Needleman and C. D. Wunsch (1970) *J. Mol. Biol.* **48**, 443–453.
4. T. F. Smith and M. S. Waterman (1981) *Adv. Appl. Math.* **2**, 482–489.

### Suggestion for Further Reading

5. D. Sankoff and J. B. Kruskal (1984) *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*, Addison-Wesley, Reading, MA.

## Induced Fit

The specificity of [enzymes](#) is of importance for their function and involves the fit of a substrate into the [active site](#) of the enzyme. The concept of enzyme specificity, which implies that there is a close steric relationship between an enzyme and its substrate, was proposed in 1894 by Emil Fischer, who developed the lock and key hypothesis of enzyme action. He explored the idea that the active site of an enzyme was simply a rigid template for the substrate and this would account for the ability of an enzyme to act on only one or a limited number of substrates.

The lock and key hypothesis was accepted for several decades. But it was questioned in 1958 by Koshland ([1](#)), who recognized that it did not account for the low activity of many enzymes with [water](#), which is always present at a concentration of 55 M, as a substrate or for the failure of enzymes to act on compounds whose structure was similar to that of the normal substrate. He proposed that:

The active site of an enzyme is not necessarily a template for the substrate.

The binding of the substrate at the active site produces a structural, or conformational, change.

Conformational change brings about the correct alignment of the catalytic groups of the enzyme with those of the substrate so that a reaction occurs.

The induced fit concept can explain why molecules of a structure similar to that of the normal substrate can combine at the active site, but not undergo reaction, because they lack the structural features required for catalysis.

An early, and now classical, demonstration of substrate-induced conformational change is the binding of glucose to hexokinase ([2](#)). This enzyme has two lobes, and between them is a cleft in which the active site is located. The binding of glucose induces the rotation of the two lobes and closure of the cleft; these movements have the effect of excluding water and facilitating transfer of the phosphoryl group from ATP to glucose, rather than to water. Kinetic evidence for conformational changes with hexokinase comes from studies on the [ATPase](#) activity of the enzyme that occurs in the absence of glucose. This activity is low relative to the rate of phosphorylation of glucose, but is stimulated by lyxose. This sugar does not undergo phosphorylation, but interacts at the binding site for glucose to increase the maximum velocity of the reaction and lower the [K<sub>m</sub>](#) for  $\text{MgATP}^{2-}$  ([3](#)). The ability of substrates to cause conformational changes at the active site of enzymes has now been demonstrated by a variety of techniques, including [X-ray crystallography](#). These changes are also a characteristic feature of the catalytic function of **allosteric** enzymes.

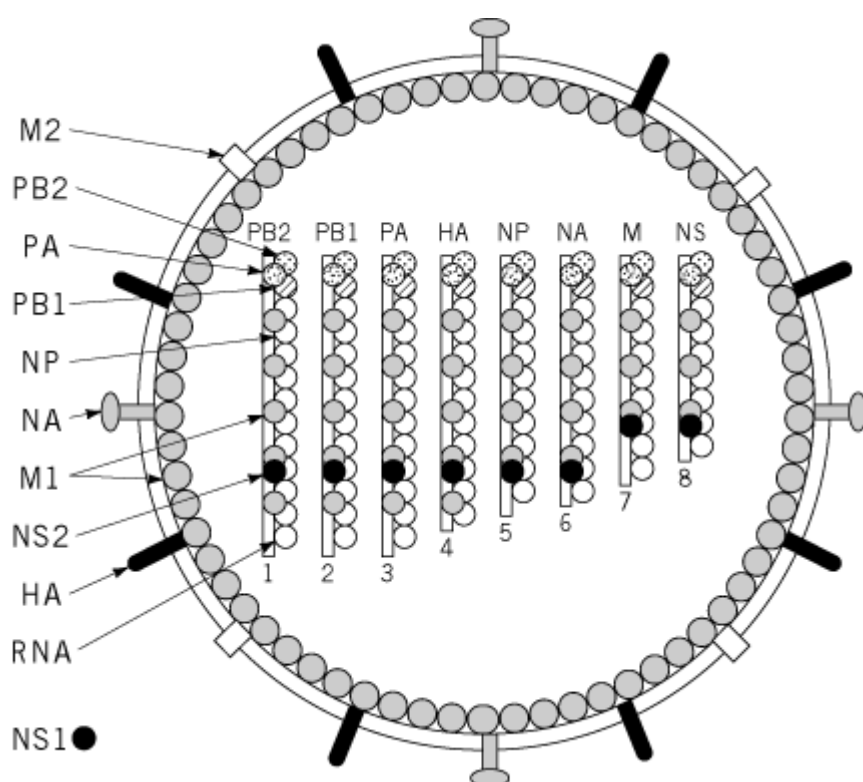
### Bibliography

1. D. E. Koshland (1960) *Adv. Enzymol.* **22**, 45–97.
2. C. M. Anderson, F. H. Zuckerman, and T. A. Steitz (1979) *Science* **204**, 375–380.
3. G. DelaFuente, R. Lagunas, and A. Sols (1970) *Eur. J. Biochem.* **16**, 226–233.

## Influenza Virus

Influenza viruses, which belong to the family of *Orthomyxoviridae*, are enveloped **viruses** with single-stranded, negative-sense RNA as the **genome**. The lipid envelope of influenza viruses is derived from the plasma membrane of the host cell in which the virus was grown. Influenza A and B viruses each contain eight species of RNA segments (Fig. 1), whereas influenza C virus contains seven segments. Influenza A, B, and C viruses are divided on the basis of antigenic differences between their nucleoprotein (NP) and matrix (M) proteins. Influenza A viruses are further subdivided into subtypes based on the antigenic differences of the hemagglutinin (HA) and neuraminidase (NA) glycoproteins (currently into 15 HA and 9 NA subtypes). Influenza A and B viruses grown in eggs or tissue-culture cells have a spherical structure, are 80 to 120 nm in diameter with some pleomorphism, and exhibit variation throughout the life cycle. Viruses freshly isolated from humans or animals exhibit greater heterogeneity and pleomorphism, including the presence of greatly elongated forms.

**Figure 1.** Schematic diagram of influenza A virus particle.



Influenza A and B virus envelopes each contain two kinds of 10- to 14-nm spikes radiating outward (the HA and the NA) (Fig. 1). The HA has a trimeric rod-shaped structure, while the NA has a tetrameric mushroom-shaped structure. The HA has the receptor (sialic acid)-binding and fusion activities. The NA has a neuraminidase activity that cleaves the virus receptor. Influenza A virus envelope also contains a tetrameric M2 **ion channel**, while influenza B virus contains a NB ion channel. The influenza C virus envelope contains only single spikes (HE), organized into orderly hexagonal arrays, that have both hemagglutinin and esterase activities.

The first isolation of the Influenza A virus from swine was reported in 1931, and that from humans was reported in 1933. The complete **nucleotide sequence** of the influenza virus was first reported for the A/PR/8/34 (PR8) strain in 1982, and it contains 13,588 nucleotides. The PR8 virus RNA segment 1 (2341 nucleotides) codes for PB2 (759 residues, 30 to 60 molecules/virion), 2 (2341 nucleotides) for PB1 (757 residues, 30 to 60 molecules/virion), 3 (2233 nucleotides) for PA (716 residues, 30 to



60 molecules/virion), 4 (1778 nucleotides) for HA (566 residues, 500 molecules/virion), 5 (1565 nucleotides) for NP (498 residues, 1000 molecules/virion), 6 (1413 nucleotides) for NA (454 residues, 100 molecules/virion), 7 (1027 nucleotides) for M1 (252 residues, 3000 molecules/virion), and M2 (97 residues, 20 to 60 molecules/virion) via splicing, 8 (890 nucleotides) for NS1 (230 residues, nonstructural) and NS2 (121 residues, 130 to 200 molecules/virion) via splicing (Fig. 1).

The viral transcription and replication takes place in the [nucleus](#) of the infected cell. The PB1, PB2, and PA proteins form a tetramer that has RNA-dependent **RNA polymerase** (replicase and transcriptase) activity. The PB2 protein has **endonuclease** activity and cuts the host mRNA to generate 10- to 13-nucleotide-long capped-RNA primer for the viral transcription (known as “cap snatching”). The viral RNA genomes and the NP protein form helical [ribonucleoprotein](#) (RNP) complexes. Twelve nucleotides at the 3' end and 13 at the 5' end of the RNA genomes are well-conserved among all influenza viruses, and they comprise the **promoter** where the virus polymerase binds and starts transcription and replication. The genomes also form a “panhandle structure,” in which the ends of an otherwise single-stranded RNA are held together by base pairing. The M1 is thought to underlie and probably integrate into the inner layer of the lipid bilayer. The M1 also interacts with the RNP, as well as with the HA and the NA. The NS1 protein is not incorporated into the virion. It interferes with the splicing and nuclear export of the cellular mRNAs. It also regulates the viral translation. The NS2 (it is currently known as NEP) has a nuclear export signal (NES, see [Nuclear Import, Export](#)).

Influenza virus is a major causative agent of the acute respiratory influenza illness, which is sometimes lethal for the elderly. In the pandemic of 1918 to 1919, called Spanish influenza, 20 to 40 million people were killed by infection with influenza A virus (H1N1). Influenza A viruses naturally infect humans, swine, horses, seals, a wide variety of birds, and some other mammalian species. Wild aquatic birds are probably the primordial reservoir of all influenza A viruses. On the other hand, influenza B viruses infect only humans. Influenza C virus infects humans and swine.

Immunity to influenza is long-lived and is subtype-specific. The HA contains the major neutralizing antigenic **epitopes**. The accumulation of point [mutations](#) (antigenic drift) alters the antigenicity of the HA. The appearance of new subtypes occurs by the reassortment of the HA or the NA gene among influenza A viruses (antigenic shift). Inactivated influenza vaccines are currently used widely worldwide, but live attenuated vaccines are being developed. Amantadine hydrochloride is a licensed antiviral agent for prophylactic and therapeutic use against all influenza A virus subtypes in humans.

Genetic manipulation of influenza A and B viruses has been successful, and some influenza A virus vector systems also have been developed.

### Suggestions for Further Reading

- R. A. Lamb and R. M. Krug (1996) "*Orthomyxoviridae: The Viruses and Their Replication*". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 1353–1395.
- B. R. Murphy and R. G. Webster (1996) "Orthomyxoviruses". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 1397–1445.

## Initiation Codon

Protein biosynthesis starts using the specific genetic code, named the initiation codon or start codon. In general, the AUG codon is used universally in both prokaryotes and eukaryotes. Occasionally,

GUG is used as an alternate initiation codon in about 10% of bacterial genes, and very rarely is UUG also used. In eukaryotes, UUG is never used, but CUG, ACG, and GUG are occasionally or rarely used as an alternate initiation codon. There is only one known case in bacteria in which AUU functions as an initiation codon. The initiation codon is recognized by the specific [transfer RNA](#) referred to as initiator tRNA, tRNA<sup>iMet</sup> (eukaryotes) or tRNA<sup>fMet</sup> (prokaryotes) and is decoded as *N*-formyl methionine (fMet) (see [fMET \(N-Formyl Methionine\)](#)).

In prokaryotes, polycistronic genes are often regulated cooperatively at the translational level. This phenomenon, referred to as translational coupling, is observed when the downstream initiation codon overlaps the upstream stop codon in the sequence AUGA, or the downstream initiation codon is sequestered by the RNA stem structure. Ribosomes can enhance translation initiation of these downstream genes by slippage of the 30S ribosomal subunit or by unwinding the inhibitory RNA stem structure.

## Initiation Complex

Initiation of protein synthesis is believed to be a rate-limiting step, because it requires multiple [initiation factors](#) and specific signals on [messenger RNA](#) that coordinately interact and function to form initiation complexes. The obvious difference between prokaryotic and eukaryotic mechanisms is based on the difference of where the [ribosome](#) enters the mRNA. Prokaryotic mRNAs are often polycistronic, and ribosomes can gain entry to the mRNA directly at the translation start site within the mRNA sequence to form initiation complexes, while eukaryotic mRNAs are monocistronic and ribosomes must gain entry to mRNA at its [Cap](#) and move downstream to the translation start site to form initiation complexes (see [Scanning Hypothesis](#)). Correct recognition of the site for initiation complex formation by prokaryotic ribosomes is mediated by complementary base-pairing between the **Shine–Dalgarno** sequence on mRNA and the anti-Shine–Dalgarno sequence on the 3' terminus of 16S rRNA. The eukaryotic initiation process is in sharp contrast to the prokaryotic process, because it involves numerous initiation factors. This reflects the fact that specific protein factors are required for 5' Cap recognition by the ribosome and for scanning of the ribosome to the downstream initiation site (1).

An initiation complex is first formed as a preinitiation complex composed of the small ribosomal particle (30S for prokaryotes, 40S for eukaryotes), [fMet-transfer RNA](#), initiation factors, and GTP at the start site of mRNA. Upon formation of proper preinitiation complexes, the large ribosomal particle (50S for prokaryotes or 60S for eukaryotes) joins the preinitiation complex, to form the initiation complex (2). During this step, GTP is hydrolyzed to GDP while the fMet-tRNA-carrier factors, IF-2 (prokaryotes) and eIF-2 (eukaryotes), dissociate from the ribosome, leaving fMet-tRNA at the P site of the ribosome.

Although the translation of virtually all eukaryotic cellular mRNAs and the majority of viral RNAs is initiated via the scanning ribosome mechanism, there is a small, yet not insignificant, number of RNAs (particularly viral RNAs) that are translated by an internal initiation mechanism [see [Internal Ribosome Entry Site \(IRES\)](#)(3)].

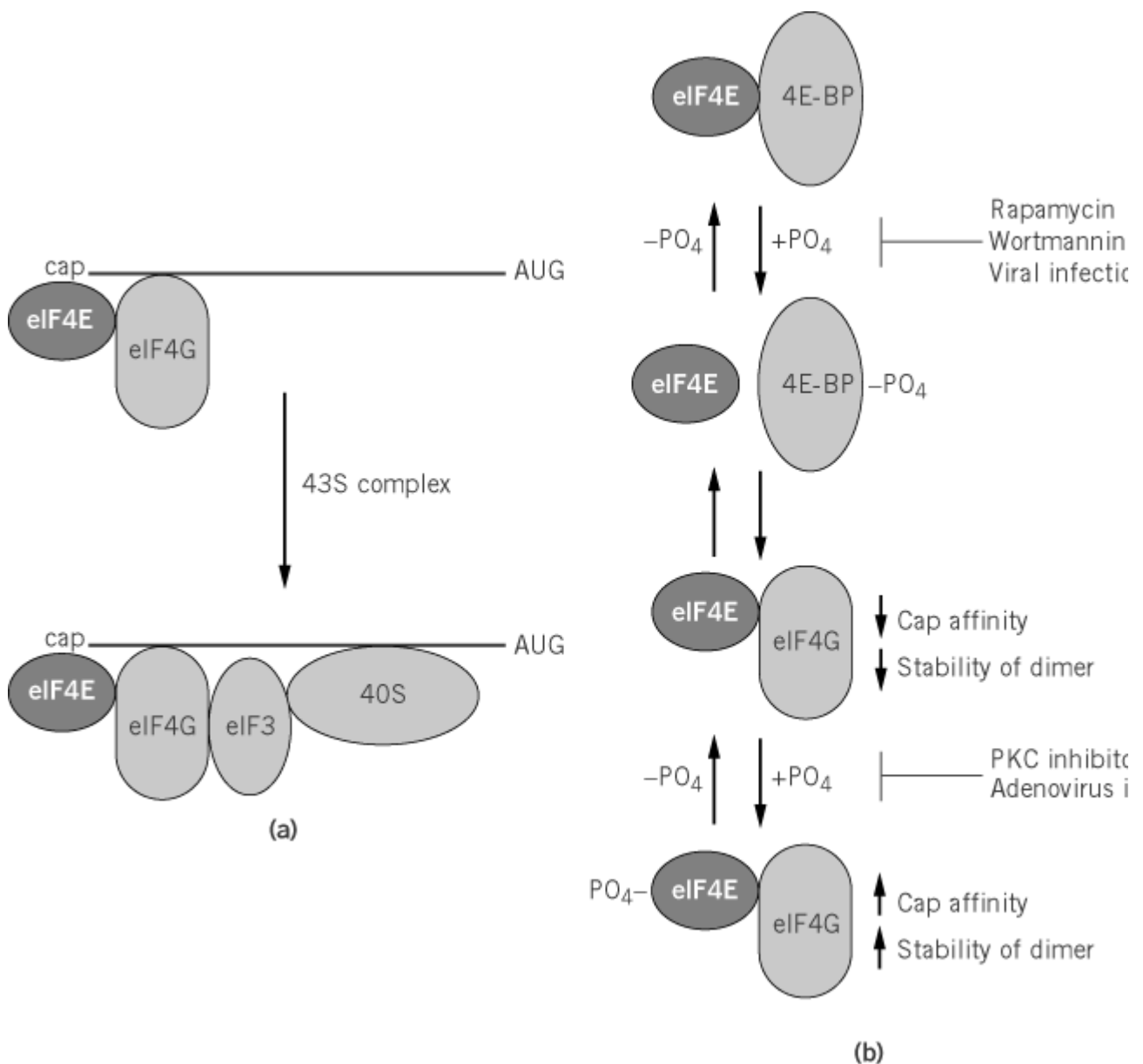
Initiation of protein synthesis is a good target for regulating gene expression. Most of the bacterial ribosomal proteins are negatively autoregulated at the initiation step of their synthesis. When present in excess over rRNA, a set of ribosomal proteins (**translational repressors**) bind to their own mRNA and stop its translation. The ribosomal protein-binding site on its mRNA (or **operator**) overlaps the ribosome binding site, thereby blocking the formation of initiation complexes. The

operator consists of a specific RNA structure that mimics the rRNA structure where the ribosomal protein (repressor) binds. Like the ribosomal protein repressor, *E. coli* threonyl-tRNA synthetase (ThrRS) is negatively autoregulated at the translation level when in excess of its substrate tRNA<sup>thr</sup>. ThrRS binds to the leader region of its own mRNA at the translational initiation site and represses its translation by preventing ribosome binding. The ThrRS operator mimics the structure tRNA<sup>thr</sup>.

Initiation of protein synthesis in eukaryotes can be regulated by **phosphorylation** of eIF-2a by mammalian protein kinase PKR (4). The PKR activity was originally detected as the enzyme responsible for double-stranded (ds) RNA-dependent protein synthesis inhibition in **poliovirus**-infected [HeLa Cells](#), and in cell-free translation systems made from rabbit reticulocytes. The target of dsRNA's mysterious ability to inhibit protein synthesis initiation was eventually shown to be the protein kinase PKR, which blocks protein synthesis by phosphorylating initiation factor eIF-2a. Higher eukaryotic cells possess an intrinsic defense system that is mediated by [interferon](#). Interferon release by virus-infected cells sensitizes neighboring cells by inducing the [transcription](#) of genes that encode anti-viral products. Some of these, notably the protein kinase PKR and 2,5A synthase, act at the level of protein synthesis. (2,5A synthase makes an oligonucleotide that activates a latent ribonuclease.) Viruses, in turn, have elaborated products (e.g., the VA RNA produced by [adenovirus](#)) that neutralize these cellular defenses, thereby enabling the production of viral proteins and progeny virions.

Recruitment of the 40S small ribosomal subunit to eukaryotic mRNA is mediated by interactions between a limited set of translation initiation factors (2). One of these factors, eIF-3, is a 40S subunit-associated factor comprised of at least eight subunits in mammalian cells that interacts with the mRNA-associated initiation factor eIF-4F (Fig. 1). eIF-4F consists of two core subunits. These are the mRNA **5'-Cap**-binding protein eIF-4E and the large subunit eIF-4G. Recent studies on eIF-4G have revealed that it binds to eIF-4E and eIF-3, as well as to the **poly(A)**-binding protein (Pab1p in yeast) (5). The multipurpose adapter nature of eIF-4G allows it to recruit the 40S ribosome to mRNA via the simultaneous association of eIF-4G with both eIF-4E and eIF-3 (Fig. 1). Knowledge of eIF-4E's interaction with the 5'-Cap permits an understanding of cellular and viral strategies to control Cap-stimulated translation. For example, cells express a small family of inhibitory proteins that regulate eIF-4E assembly. These are called the 4E-binding proteins [4E-BPs (6)]. The 4E-BPs share an amino acid motif with the N-terminal domain of eIF-4G that is known to be required for eIF-4G's interaction with eIF-4E. In their nonphosphorylated form, the 4E-BPs act as competitive inhibitors of the eIF-4G-eIF-5E interaction (Fig. 1).

**Figure 1.** Control of initiation complex formation via Cap binding in eukaryotes (7). (a) eIF-4E recruits the 40S subunit of mRNA by protein interaction with eIF-4G and eIF-3. (b) Regulation of initiation complex formation by phosphorylated eIF-4E binding proteins, 4E-BPs.



## Bibliography

1. W. C. Merrick and J. W. B. Hershey (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews and N. Sonenberg, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 31–69.
2. R. J. Jackson (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 71–112.
3. E. Ehrenfeld (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 549–573.
4. H. D. Robertson and M. B. Mathews (1996) *Biochimie* **78**, 909–914.
5. M. W. Hentze (1997) *Science* **275**, 500–501.
6. A. Pause et al. (1994) *Nature* **371**, 762–767.
7. A. B. Sachs, P. Sarnow, and M. W. Hentze (1997) *Cell* **89**, 831–838.

## Suggestions for Further Reading

8. M. Springer, C. Portier, and M. Grunberg-Manago (1997) In *RNA Structure and Function* (R. W.

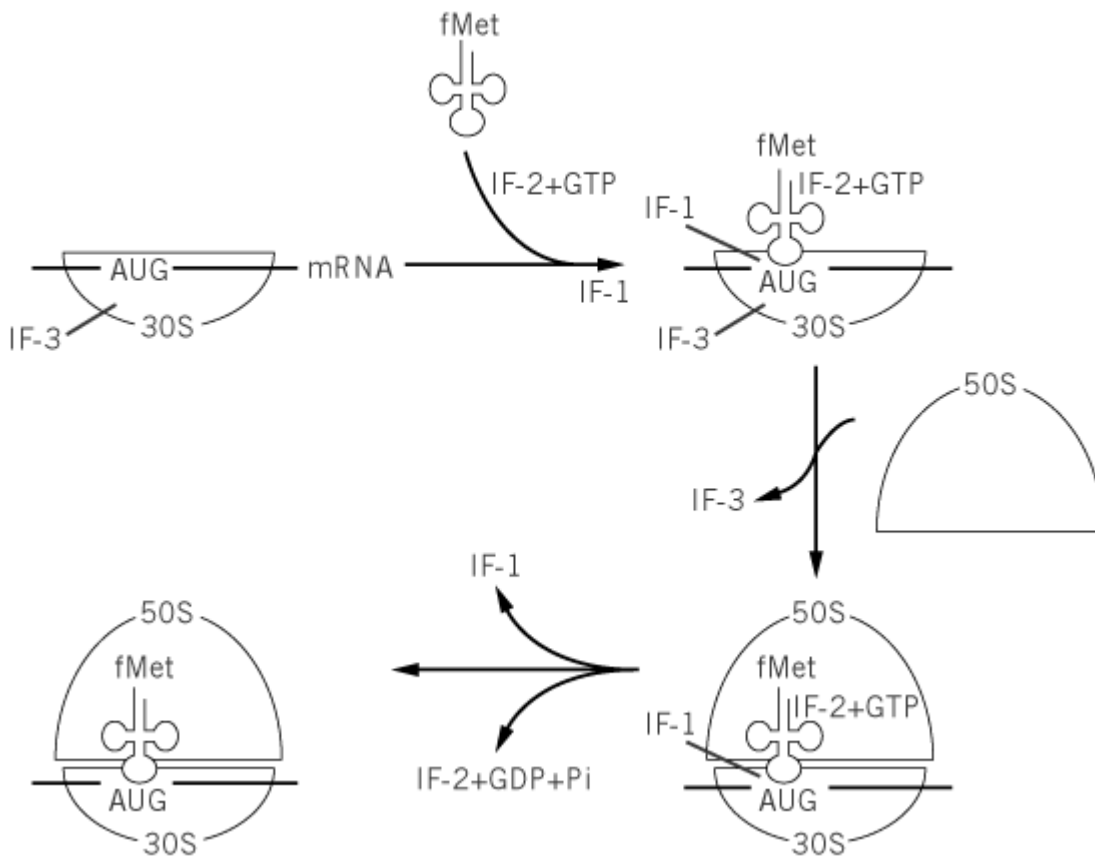
Simons and M. Grunberg-Manago, eds.), Cold Spring Harbor, Laboratory, Cold Spring Harbor, NY, pp. 377–413.

9. A. B. Sachs, P. Sarnow, and M. W. Hentze (1997) *Cell* **89**, 831–838.

## Initiation Factor (IF)

Formation of an initiation complex for [translation](#) of protein biosynthesis on 70 S [ribosomes](#) in bacteria involves, in addition to initiator fMet-tRNA<sup>fMet</sup>, [messenger RNA](#), and GTP, the cooperative interaction of three proteins, initiation factors IF-1, IF-2 and IF-3 (Table [1](#)). Until now, no convincing data have indicated the existence of additional factors. Numerous *in vitro* experiments demonstrate that IF-1 and IF-3 shift the ribosome-dissociation equilibrium from 70 S ribosomes toward dissociation into 30 S and 50 S subunits. Under physiological conditions, the ribosomal subunit equilibrium is toward 70 S formation, and IF-1 serves to increase the rate constant for both the dissociation and association of 70 S ribosome, whereas IF-3 shifts the equilibrium strongly toward dissociation by binding to the 30 S particle. IF-3 also selects the initiator tRNA in the presence of the initiation codon triplet, as well as on a natural mRNA template. IF-3 action distinguishes the initiator [transfer RNA](#) from elongator tRNAs by specific domains in the tRNA and not by its formyl methionine ([fMet](#)). Part of the information is located in the [anticodon](#) stem and loop region. IF-3 has been shown to destabilize selectively polynucleotide-primed 30 S complexes formed with elongator tRNAs. The third initiation factor, IF-2, directs fMet-tRNA<sup>fMet</sup> binding to the 30 S subunit. This initiation factor might also exclude noninitiator tRNAs from the P site, while allowing formylated charged initiator tRNA to enter. This reaction is stimulated by IF-1 and GTP. In addition, IF-2 helps the association of the two subunits and has a GTPase activity in the presence of ribosomes. The process of initiation complex formation in prokaryotes is summarized in Fig. [1](#) ([1](#)).

**Figure 1.** The initiation pathway of protein synthesis in prokaryotes.



**Table 1. *Escherichia coli* Initiation Factors**

| Factor | Molecular Mass | Number of Amino Acids | Gene        | Function  |
|--------|----------------|-----------------------|-------------|---|
| IF-1   | 8,119          | 71                    | <i>infA</i> | Promotes IF-2 and IF-3 functions  |
| IF-2a  | 97,300         | 889                   | <i>infB</i> | Binds fMet-tRNA and GTPase  |
| IF-2b  | 79,700         | 732                   | <i>infB</i> |   |
| IF-3   | 20,668         | 181                   | <i>infC</i> | Ribosome dissociation and mRNA binding; increases the specificity of the initiation complex for fMet-tRNA |

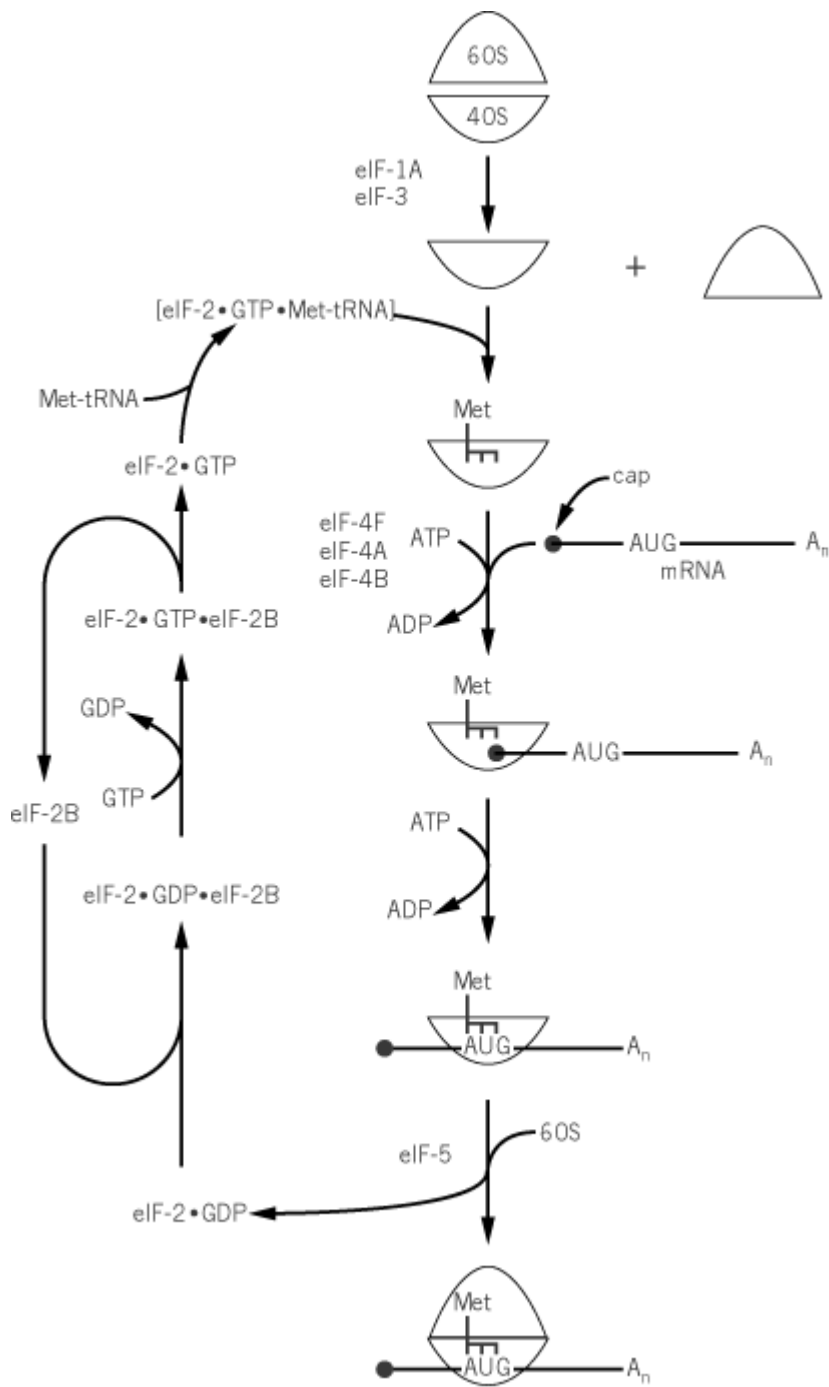
In the absence of tRNA or initiation factors, the 30 S particle by itself binds to the proper initiation region on mRNAs. The sequence of events directed by initiation factors shows that IF-1 and IF-3 are released before or during the joining of the subunits, whereas IF-2 is released after their joining and requires GTP hydrolysis. The three factors are present in cells at sufficiently high concentration to saturate the free 30 S particle. The steady-state levels of initiation factors are coordinately regulated in exponentially growing bacteria relative to one another and relative to ribosome levels. All levels

rise as a function of increasing growth rate, suggesting that the cellular levels of initiation factors are under metabolic control. Neither IF-2 and IF-3 synthesis is under [stringent control](#) *in vivo*.

All of the initiation factor genes have been cloned and mapped: *infA* for IF-1, *infB* for IF-2, and *infC* for IF-3. They are dispersed on the bacterial chromosome and, with the exception of *infA*, are associated with other components of the protein biosynthesis apparatus (threonyl and phenylalanyl-tRNA synthetase, and ribosomal protein L20 and L35 for IF-3; tRNA<sup>fMet</sup>, NusA, and ribosomal protein S15 for IF-2). The *infB* gene produces two forms of IF-2, IF-2a (97.3 kDa), and IF-2b (79.7 kDa), using two in-frame initiator codons. Both forms are implicated in fMet-tRNA binding to 30 S ribosomes and possess a ribosome-dependent GTPase activity. The meaning of the existence of two forms of IF-2 in bacteria, as well as their functional difference, is not known.

The eukaryotic initiation steps are promoted by at least 10 different initiation factors, some of which are themselves complexes of numerous protein subunits. Figure 2 summarizes where the more important of the initiation factors function in the pathway (2). eIF-1A and eIF-3 are involved in the dissociation of ribosomes; eIF-2 and eIF-2B are primarily concerned with tRNA<sup>iMet</sup> binding to 40 S subunits; eIF-4A, eIF-4B, and eIF-4F promote mRNA binding and scanning; and eIF-5A is involved in the junction step. Essentially all the initiation factor polypeptide cDNAs have been **cloned** and sequenced (Table 2).

**Figure 2.** The initiation pathway of protein synthesis in eukaryotes.



n pathway of protein synthesis in eukaryotes. [\[Full View\]](#)

**Table 2. Eukaryotic Initiation Factors (2)**

| Factor | Subunit | Molecular Mass (kDa) | Function  |
|--------|---------|----------------------|---|
| eIF-1  |         | 12.6                 | Enhances initiation complex formation                             |
| eIF-1A |         | 16.5                 | Ribosomal dissociation; promotes Met-tRNA <sup>iMet</sup> binding |



|        |       |   |
|--------|-------|---|
| eIF-2  | 125   | Binds Met-tRNA <sup>iMet</sup> and GTP                                    |
| a      | 36.1  | Phosphorylated subunit  |
| b      | 38.4  | Binds Met-tRNA <sup>iMet</sup>  |
| g      | 51.8  | Binds GTP; Met-tRNA <sup>iMet</sup>                                       |
| eIF-2B | 270   | Guanine nucleotide exchange factor for eIF-2                              |
| a      | 33.7  |   |
| b      | 39.0  | Binds GTP   |
| g      | 58    | Binds ATP   |
| d      | 57.1  | Binds ATP   |
| e      | 80.2  | Phosphorylated subunit  |
| eIF-2C | 94    | Stabilizes ternary complex in presence of RNA                             |
| eIF-3  | 650   | Dissociates ribosomes; promotes Met-tRNA <sup>iMet</sup> and mRNA binding |
| p35    | 35    | Highly phosphorylated   |
| p36    | 36.5  |   |
| p40    | 39.9  |   |
| p44    | 35.4  | RNA binding motif   |
| p47    | 47    |   |
| p48    | 52.5  | Identical to Int-6 (oncogene) homologue                                   |
| p66    | 64.0  | Binds RNA   |
| p110   | 105.3 |   |
| p116   | 98.9  | Major phosphorylated subunit  |
| p170   | 166.5 |   |
| eIF-4A | 44.4  | ATPase, helicase, binds RNA   |
| eIF-4B | 69.8  | Binds RNA, promotes helicase activity                                     |
| eIF-4F | 250   | Binds 5' Cap, helicase  |
| eIF-4E | 25.1  | Cap-binding subunit   |
| eIF-4A | 44.4  | ATPase, helicase  |
| eIF-4G | 153.4 | Binds eIF-4A, eIF-4E, eIF-3, and Pab1p                                    |
| eIF-5  | 48.9  | Promotes GTPase with eIF-2 and ejection of eIFs                           |
| eIF-6  | 25    |   |

eIF-2 forms a ternary complex with Met-tRNA<sup>iMet</sup> and GTP that binds to the 40 S ribosome, where it is established by eIF-1A and eIF-3. The resulting 43 S preinitiation complex then interacts with mRNA prepared by eIF-4F, eIF-4A, and eIF-4B. eIF-4F is a heterodimeric complex that binds to the [-Cap](#) structure of mRNA. It is composed of one molecule each of eIF-4E, eIF-4A, and eIF-4G. eIF-4E is of lower abundance in the cell than the other subunits and initiation factors, thereby making the eIF-4F complex rate-limiting under most circumstances. eIF-4A is an RNA-dependent [ATPase](#) and

binds RNA only weakly. Its sequences contain ATP-binding elements found in the **DEAD box** family of [RNA Helicases](#). eIF-4B is an RNA-binding protein that contains two **RNA-binding** motifs: the RNP consensus octamer and hexamer elements and an arginine-rich domain. eIF-4F (bound to the 5' Cap) together with eIF-4B possess RNA helicase activity dependent on ATP hydrolysis. When mRNA secondary structure is melted, the 40 S ribosome is able to bind to the 5' Cap region and begin scanning. Proper interaction of the Met-tRNA<sup>iMet</sup> anticodon with the AUG initiator codon is monitored by eIF-5, which promotes GTP hydrolysis, followed by the ejection of initiation factors and junction of the 40 S initiation complex with 60 S subunits. eIF-2 leaves as a complex with GDP and requires eIF-2B-catalyzed exchange of the GDP for GTP in order to promote another round of initiation.

### Bibliography

1. D. E. Draper (1996) In *Escherichia coli and Salmonella*, 2nd ed. (R. Curtiss III et al., eds.), American Society for Microbiology Press, Washington, D.C., pp. 902–908.
2. W. C. Merrick and J. W. B. Hershey (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 31–69.

### Suggestion for Further Reading

3. M. Springer, C. Portier, and M. Grunberg-Manago (1997) In *RNA Structure and Function* (R. W. Simons and M. Grunberg-Manago, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 377–413.

## Initiation Of DNA Replication

How replication of a [chromosome](#) is initiated has been a major subject of molecular genetics, because the regulation of chromosomal replication occurs primarily at the stage of initiation (see [DNA Replication](#)). Molecules that might be involved in the initiation of replication of the bacterial [circular chromosome](#) were proposed by Jacob et al. in 1963 (1). They assumed that each [replicon](#) consists of two elements participating in positive regulation of the initiation: (1) a *cis*-acting DNA signal (the replicator, which is now called the [replication origin](#)), and (2) a *trans*-acting protein factor (the initiator), which would interact directly with the replicator. The positive regulation for the initiation of chromosomal replication arose from the observation that **protein biosynthesis** is required for the initiation of each round of replication in *Escherichia coli*. Later, it was found that initiation occurs when each individual cell reaches a fixed amount of cell mass relative to the number of the replication origins, regardless of the growth rate (2). To explain these phenomena, the initiator accumulation theory was proposed (3). In this model, an initiator is synthesized and accumulates in the cell at the same rate as the cell mass increases, and initiation takes place every time a certain number of initiators accumulate for each replication origin. The model also proposed that the initiator should act as a [repressor](#) of its own **gene expression** to maintain the cellular content of the protein (autogenous regulation) as constant and so it is accumulated in concert with an increase in cell mass. Although Jacob and his colleagues merely stated that the model was a simplified one, their far-sighted hypothesis has now proven to be correct concerning the prokaryotic replicons, those of bacteria, **bacteriophage**, and **plasmids** (see [DNA Replication](#)).

### 1. Initiation of Bacterial Plasmid Replication

The replication of bacterial plasmids and phage has been shown to start from a unique site on the replicon (the origin, see [Replication Origin](#)), and most plasmids and phage encode a gene, often designated *rep*, for a replicon-specific protein required for the initiation of DNA replication. The replicon-specific Rep proteins have been shown to recognize specific sequences at the origin and to act positively to recruit host-encoded replication enzymes (4). The initiation frequency of the replication of plasmids is regulated so as to maintain a steady-state **copy number** characteristic of the particular plasmid in a given host bacterium (5). The majority of plasmids in **Gram-negative** bacteria contain direct repeats of the Rep binding sequence (iterons) within the origin. Many of these plasmids also contain indirect repeats of sequences related to the iterons within the **promoter** region of the *rep* gene, and the expression of the Rep protein is autoregulated, as proposed in the initiator accumulation theory. Direct repeats within the functional origin of an iteron-containing plasmid, when inserted into a plasmid normally compatible with the iteron-containing plasmid, show strong incompatibility with the plasmid. In addition, insertion of additional homologous iterons into an iteron-containing plasmid results in a reduction of the copy number of the plasmid. Based on these observations, it was proposed initially that replication begins when a sufficient amount of Rep protein is accumulated to cover the iterons at the origin, and additional Rep binding sequences reduce the effective amount of the initiator (6). However, extensive studies conducted on the regulation of initiation of the iteron-containing plasmids have revealed properties that argue against the simple titration model (5). For example, increasing the concentration of the Rep protein of several plasmids *in vivo* did not result in a proportional increase in the copy number. Several Rep proteins have been shown to exist as two forms; dimers that bind to the iterons in the origin and monomers that regulate its own expression. Although new models that modify the simple titration model to explain the regulation of the initiation of plasmid replication have been proposed, much remains to be learned to understand it completely (5).

## 2. Initiation of Replication of Bacterial Chromosome

A number of genes have been identified as being involved in the initiation of chromosomal replication in bacteria. The functions of these genes during the initiation process have now been well defined through genetic and biochemical studies in *E. coli* (4, 7). DnaA protein was found to function at the first step in the initiation process by binding to specific DNA repeats (the DnaA box) in the replication origin of the chromosome (*oriC*), and then opening the double strand for formation of the primosome, the machinery for the first primer synthesis (see [DNA Replication](#)). The combination of a DnaA box and DnaA protein functioning as *cis* and *trans* regulatory elements, respectively, in the initiation of chromosomal replication was subsequently found to be common in **eubacteria** (7, 8).

Expression of the *E. coli dnaA* gene is negatively regulated by the DnaA protein itself, by the protein binding to DnaA boxes located in the promoter region of the gene, and the concentration of the DnaA protein is kept proportional to the origin concentration. Overproduction of the DnaA protein stimulated the initiation of chromosomal replication up to two-fold, and that stimulation was proportional to the increase in the DnaA protein concentration within this limited range. DnaA overproduction also resulted in a two-fold increase in [minichromosome](#) (*oriC* plasmid) copy number. Thus, *E. coli* DnaA protein is considered to fulfill the qualifications required for an initiator that controls the timing of initiation of replication during the cell cycle, and the time of initiation is assumed to be set by the DnaA protein accumulating to a threshold level (9). However, it should be noted that other factors affecting initiation control have been identified. The characteristic of *oriC* of enterobacteria is the presence of many GATC sequences, which are sites for the Dam **restriction-modification** methylase. It was found that newly replicated *oriC* DNA strands are not methylated during up to one-third of the cell cycle after the replication. Such hemi-methylated DNA tends to bind to the cell membrane and does not act as a template for activation by DnaA protein. Furthermore, initiation of replication occurs more or less at random in relation to the cell cycle in a mutant lacking Dam methylase (7). The DnaA protein can be modified in different ways by binding ATP or ADP. Phospholipids were found to reactivate the inactive form (ADP-DnaA) *in vitro* by exchanging ADP for ATP. Initiation from *oriC* seemed to be inhibited *in vivo* in a mutant in which

the synthesis of phospholipids was limited ([10](#)).

The *dnaA* gene is also autoregulated in *Bacillus subtilis*. However, expression of the *B. subtilis dnaA* gene is coupled to initiation of replication of the chromosome. It is assumed that the autorepression by interaction of DnaA protein with many DnaA boxes near the promoter is so strong in *B. subtilis* that the gene is expressed only when the chromosome undergoes vigorous conformational changes during the replication of the gene itself ([11](#)). In *B. subtilis*, fewer than one *oriC* plasmid/cell is allowed to coexist with the chromosome. Introduction of extra copies of DnaA boxes would be incompatible with chromosomal replication. This stringent control of the initiation observed may be ascribed to the strong negative control of *dnaA* expression. Accumulation of the DnaA protein to a threshold level was shown to be prerequisite for the initiation of replication in *B. subtilis*. However, control of the timing of the initiation may require other factors, because the DnaA protein seems to be synthesized only in the early stage of the replication cycle in this cell. Although the role of DnaA as a key factor for regulation of initiation is clear in both *E. coli* and *B. subtilis*, the overall view of the regulation of the timing of initiation is not still clear. And, different devices to control the timing seem to have been evolved in two bacteria.

### 3. Initiation in Eukaryotes

A unique initiator protein ([T Antigen](#)) has been identified to regulate the initiation of replication of the eukaryotic **SV40 virus**. However, no single initiator protein like bacterial DnaA has been found for eukaryotic chromosomes. Instead, multiprotein complexes have been identified to initiate replication of the multiple replicons of eukaryotes. Furthermore, a complex network of [protein–protein interactions](#) (complex formation) and protein modifications ( **phosphorylation**) regulates the initiation (see [DNA Replication](#)).

### Bibliography

1. F. Jacob, S. Brenner, and F. Cuzin (1963) Cold Spring Harbor Symp. Quant. Biol. **28**, 329–348.
2. W. D. Donachie (1968) Nature **219**, 1077–1079.
3. L. Sompayrac and O. Maaloe (1973) Nature New Biol. **241**, 133–135.
4. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman and Company, New York.
5. D. R. Helinski, A. E. Toukdarian, and R. Novick (1996) in *Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 2295–2324.
6. H. Tsutsui and K. Matsubara (1981) J. Bacteriol. **147**, 509–516.
7. W. Messer and C. Weigel (1996) In *Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington DC, pp. 1579–1601.
8. H. Yoshikawa and N. Ogasawara (1991) Mol. Microbiol. **5**, 2589–2597.
9. J. Herrick, M. Kohiyama, T. Atlung, and F. G. Hansen (1996) Mol. Microbiol. **19**, 659–666.
10. W. Xia and W. Dowhan (1995) Proc. Natl. Acad. Sci. USA **92**, 783–787.
11. H. Yoshikawa and R. G. Wake (1993) In *Bacillus subtilis and Other Gram-positive Bacteria* (L. Sonenshein, J. A. Hoch, and R. Losick, eds.), American Society for Microbiology, Washington, DC, pp. 507–528.

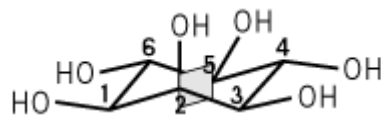
### Inositol Lipids and Phosphates

*Myo*inositol, a carbohydrate first found in the mid-1800s in muscle, was defined in the 1930s as a [growth factor](#) for yeasts and mammals and then as a component of membrane phospholipids (phosphoinositides) (1, 2). In the 1950s, Hokin and Hokin found that secretory stimuli often accelerate phosphoinositide turnover in their target cells (2, 3). This initial observation led ultimately to recognition of inositol lipids and phosphates as crucially important regulatory molecules in eukaryote cells.

## 1. *Myo*inositol

*Myo*inositol is one of nine stereoisomeric hexahydroxycyclo-hexanes (the inositols). Several inositols occur in nature, but only *myo*inositol is used widely in the headgroups of membrane phospholipids and in water-soluble inositol polyphosphates. When in the preferred “chair” configuration, *myo*inositol has five equatorial hydroxyl groups (1-6): only the 2-OH is axial (Fig. 1) (1, 2, 4). IUPAC recommendations decree that D- *myo*inositol 1-phosphate (abbreviated as Ins1*P*; the same molecule as L-*myo*inositol 3-phosphate) is the reference molecule for the numbering of biological *myo*inositol derivatives (4). Inositol is made biosynthetically by cyclization of the central metabolite glucose 6-phosphate (by way of 5-ketoglucose 6-phosphate) to Ins3*P*, which is then dephosphorylated (1).

**Figure 1.** The configuration of, and carbon atom numbering (D-configuration) in, *myo*inositol. The mirror plane () identifies the plane of symmetry through the 2- and 5-carbons of the inositol ring.



It is not clear how [evolution](#) “selected” *myo*inositol (from here on simply termed *inositol*) to take on multiple biological roles. Maybe inositol's stereochemistry gives its derivatives some biologically desirable property. Alternatively, evolution may have “needed” a conveniently synthesized and stable polyol, and selected a molecule that comes, by a very short route, from a ubiquitous cell constituent.

The stereochemical versatility of inositol has allowed biology to put it to many uses. Because it has an axis of symmetry through its 2- and 5-carbons (Fig. 1), addition of one substituent to C-1, C-3, C-4, or C-5 immediately renders every carbon in the ring stereochemically unique. For example, Ins1*P* and Ins3*P* are enantiomers, as are the inositol tetrakisphosphates Ins(1,3,4,5)*P*<sub>4</sub> and Ins(3,4,5,6)*P*<sub>4</sub> (L-Ins(1,3,4,5)*P*<sub>4</sub>) (1, 4). Thus, the possible variety of stereochemically (and biologically) unique inositol derivatives can be almost endless.

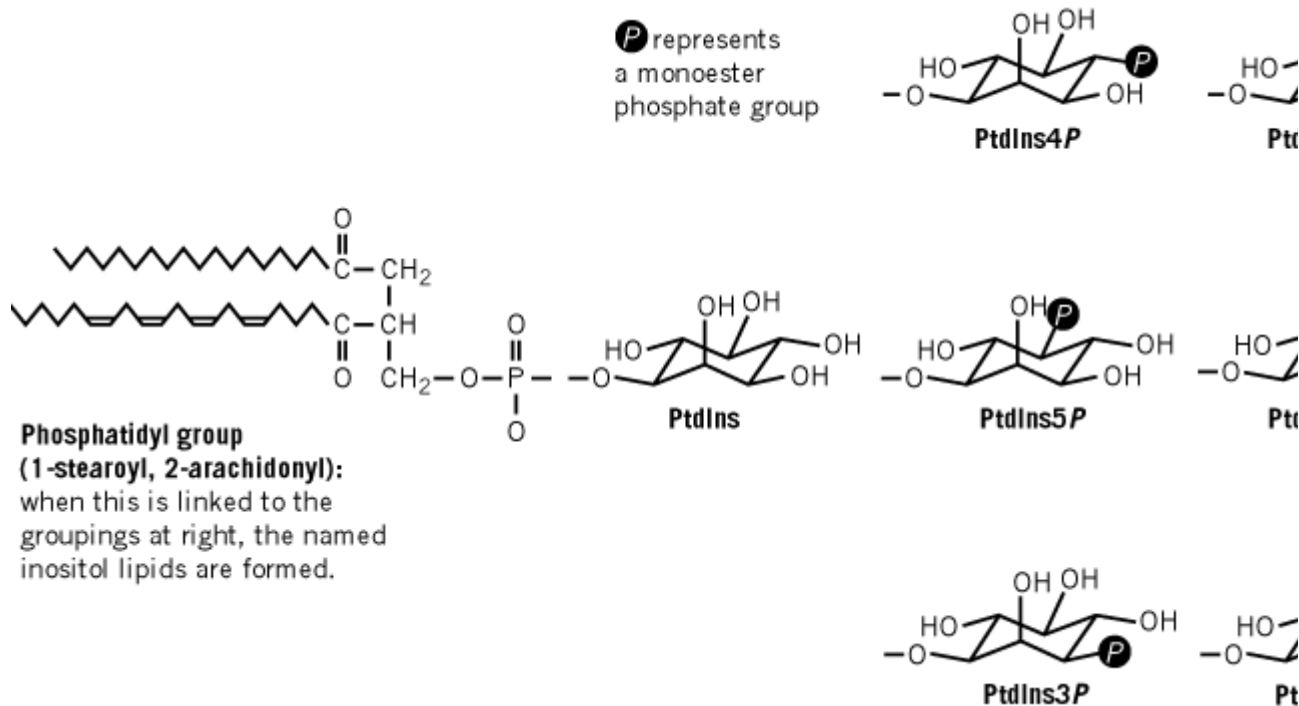
## 2. Inositol lipids and phosphates: variety and evolution

Inositol seems to be essential to all eukaryotic cells, and some members of both bacterial kingdoms (Eubacteria and Archaeobacteria) also contain inositol lipids and phosphates. Many of the functions of the inositol lipids are becoming understood in remarkable detail, but only a few functions have been definitively ascribed to particular inositol polyphosphates.

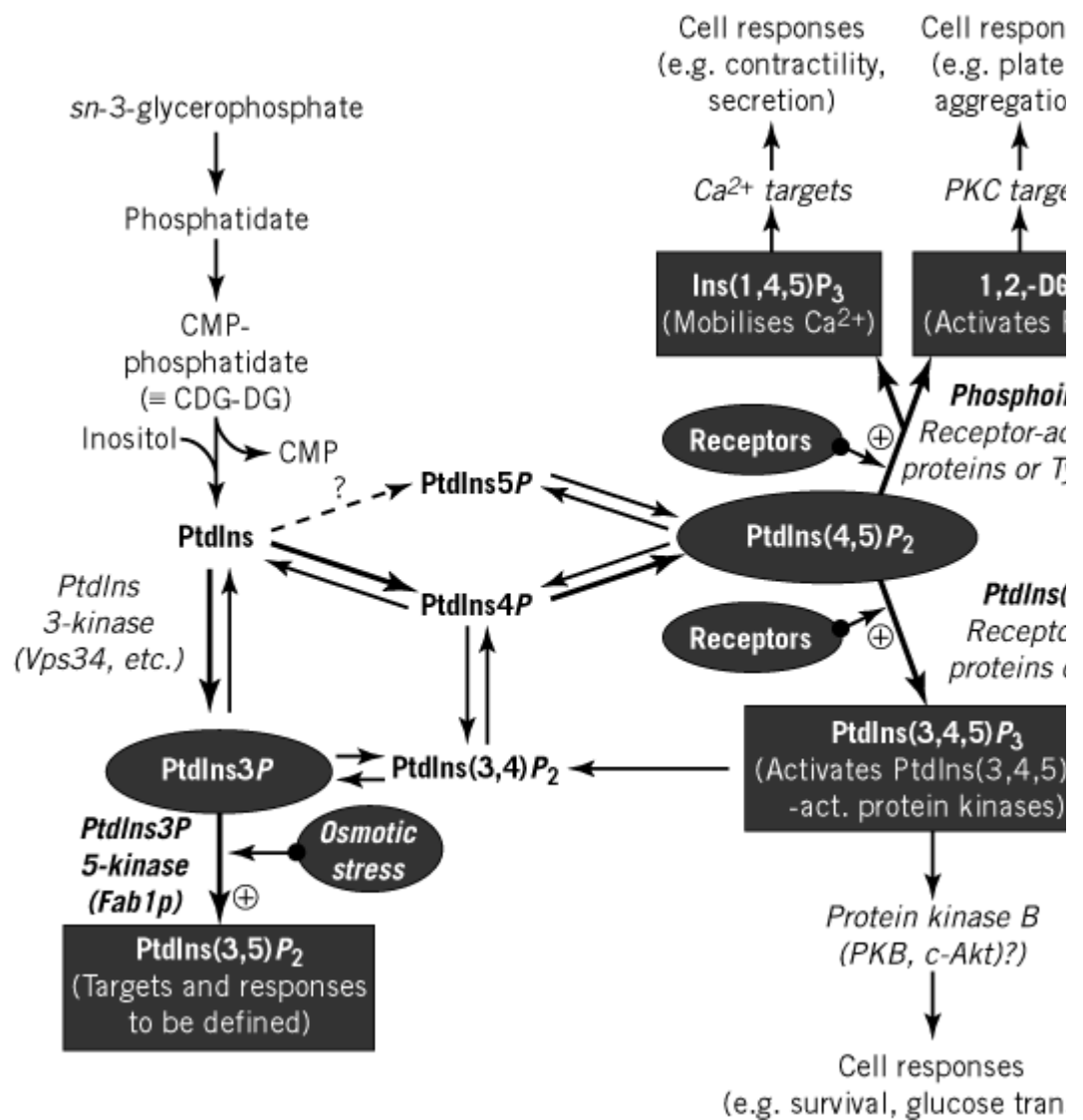
Before about 1950, cells were thought to contain a single inositol glycerophospholipid, namely phosphatidylinositol (PtdIns; *sn*-1,2-diacyl-glycero-3-phospho-D-1-*myo*inositol) (Fig. 2). This is generally the most abundant inositol glycerophospholipid, and it and its phosphorylated derivative often have a 1-arachidonyl, 2-stearoyl acyl fatty acyl pairing in mammalian cells. Seven other

glycerophosphoinositides, all of which are phosphorylated derivatives of PtdIns, have been identified since 1949 (Fig. 2). Figure 3 summarizes the likely metabolic and functional interrelationships between these glycerolipids.

**Figure 2.** The “simple” inositol glycerolipids of eukaryotic cells.



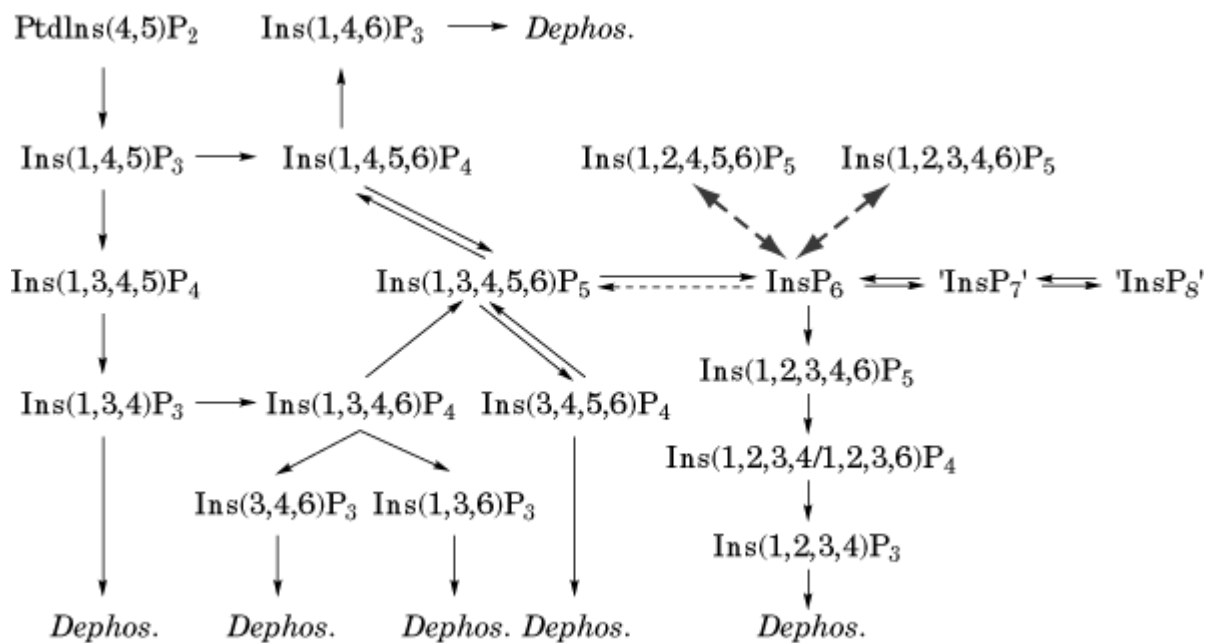
**Figure 3.** Likely pathways involved in eukaryotic glycerophosphoinositide-based signaling. Substrates of signaling path in shaded ovals, and proven or probable second messengers are in shaded boxes. Note that PtdIns serves as the common and that PtdIns(4,5) $P_2$  3-kinase signaling pathways and PtdIns(4,5) $P_2$ -directed (Type I) and PtdIns-directed (Type III) in functions.



In addition to these relatively simple membrane lipids, the membrane-penetrating lipid anchors of many cell-surface (glyco)proteins and their precursors are PtdIns glycans of varying complexity (see [GPI Anchor](#)). Some eukaryotic cells, notably plants and yeasts, also contain complex sphingolipids with inositol-containing headgroups. Although these do not yet have proven functions, their metabolism in yeast is somehow linked to vesicle trafficking (5).

Cells also contain many water-soluble inositol (poly)phosphates. For example, *Dictyostelium* and mammalian [cell lines](#), such as WRK1 (mammary tumor) and HL60 (myeloid), contain 25 or more different inositol phosphate species (eg. ref. 6). Figure 4 summarizes the metabolic links between many of these inositol phosphates, primarily as in mammalian cells. A variety of inositol phosphate interconversion pathways have been reported from plants and yeasts, including a remarkably direct conversion of  $\text{Ins}(1,4,5)\text{P}_3$  to  $\text{InsP}_6$ , apparently by a single [enzyme](#) (7). Plant seeds store much of the phosphate essential to their successful germination and development as phytate (an insoluble bivalent cation complex of  $\text{InsP}_6$ ), alleged to be the most abundant phosphate ester in nature.

**Figure 4.** Pathways of inositol polyphosphate interconversion in eukaryote cells. For simplicity, individual inositol monophosphates and bisphosphates on the dephosphorylation pathways have been omitted.



Some families of Eubacteria have additional inositol-containing phospholipids, the best characterized being the PtdIns mannosides of mycobacteria. Archaeobacteria, which inhabit some of the most hostile environments on the planet, have glycerophospholipids that are very distinctive, both because the steric configuration of the diradylglycerol backbone of these lipids is opposite to that of the eukaryotic and eubacterial lipids and because they use very different, and chemically very stable, hydrocarbon side chains. Again, however, some of these have inositol headgroups (8). Moreover, hyperthermophilic archaeobacteria use inositol in another remarkable way. When they are exposed to temperatures that are even more extreme than usual, they synthesize very high cytosolic concentrations (~0.5M) of the K<sup>+</sup> salts of water-soluble inositol phosphate derivatives, such as dimyoinositol 1,1'-(3,3')-phosphate, as cytoprotective solutes (9).

When the primeval cells from which all later organisms evolved were emerging, a fairly early step was almost certainly enclosure of the newly evolving cellular metabolic machinery within a selectively permeable membrane envelope made of amphiphilic lipids. The fact that *myo*inositol, an unusually stable cyclic polyol, is a component both of the diether lipids of Archaeobacteria and of the (mainly) diacyl-glycerolipids of Eubacteria (some) and eukaryotes suggests that it was recruited into membrane lipids very early. This probably happened before the two-billion-year-old evolutionary radiation that gave rise to the major biological kingdoms, making inositol a very ancient cell component.

Only since the early 1980s, however, has it been recognized how inositol's stereochemical versatility has let cells employ its derivatives for many different jobs. Throughout this recent period, frequent discoveries of “new” inositol-containing cell constituents have regularly been followed by the definition of “new” functions for inositol derivatives.

### 3. Inositol lipids and phosphates in signaling

One of the most essential characteristics of cells is their ability to respond with changes in their behavior to chemical changes in their environment. Many of these extracellular chemical signals cannot penetrate the plasma membrane, and so the cell surface receptors that respond to them couple to transmembrane [signal transduction](#) mechanisms that pass information to the cell interior. Inositol lipids play essential roles as the metabolic substrates and/or the output signals of at least two, and maybe three or more, widespread cell signaling systems.

#### 3.1. Signaling through Phospholipase C Activation



The best characterized of these signaling pathways uses receptor activation of PtdIns(4,5)  $P_2$  hydrolysis by phosphoinositidase C ([phospholipase C](#), PLC) to make the [second messengers](#) inositol 1,4,5-trisphosphate (Ins(1,4,5)  $P_3$ ) and 1,2-diacylglycerol, as outlined in Fig 3 (top right). Despite extensive early work by the Hokins (3), only in the 1970s was it realized that inositol lipid hydrolysis has the properties of a coupling reaction that links cell-surface receptors to the intracellular mobilization of  $Ca^{2+}$  (2) (see [Calcium Signaling](#)). PtdIns(4,5) $P_2$  hydrolysis was then identified as the primary stimulated reaction (10). Meanwhile, Nishizuka defined the polyunsaturated 1,2-diacylglycerol (1,2-DG) liberated from PtdIns(4,5) $P_2$  as a second messenger that activates many members of the protein kinase C family (11). Finally, Streb et al. showed that Ins(1,4,5)  $P_3$  released  $Ca^{2+}$  from an intracellular store (12). We now know that Ins(1,4,5) $P_3$  triggers intracellular  $Ca^{2+}$  release through the intrinsic cation **channels** of at least four subtypes of Ins(1,4,5) $P_3$  “receptor” that are expressed primarily in various elements of the [endoplasmic reticulum](#) (13).

These revelations established two new principles. First, chemical species very different from nucleotide cyclic phosphates (see [-cyclic AMP, cAMP](#); [-Monophosphate, cGMP](#)) and the  $Ca^{2+}$  ion, the prototypic intracellular messenger molecules that were recognized between the mid-1950s and the mid-1970s, can function as “second messengers.” Second, receptor-generated intracellular messenger molecules need not be water-soluble and thus mobile within the cell. Instead, protein kinase C homes to the 1,2-DAG that is formed at the plasma membrane as a consequence of PtdIns(4,5) $P_2$  hydrolysis. There it is activated by simultaneous interactions with 1,2-DAG, anionic aminophospholipids, notably phosphatidylserine, in the inner leaflet of the plasma membrane, and the ambient concentration of cytoplasmic  $Ca^{2+}$ .

There are numerous excellent reviews of the PLC signaling pathway and the extraordinary variety of receptors that harness it (eg, refs. 11-15; and see [Phospholipases C](#)): some of the important receptors that harness this signaling pathway are listed in Table 1 (see also [Acetylcholine Receptor](#), [Rhodopsin](#)). There is a substantial family of PLCs that can hydrolyze PtdIns(4,5) $P_2$ , and these are activated by two or more fundamentally different mechanisms (16). PLCs of the b-family are activated by some subfamilies of receptors in the 7-span or “serpentine” superfamily. Different serpentine receptors activate PLC-bs, either through the GTP-ligated  $\alpha$ -subunits of [heterotrimeric G proteins](#) of the  $G_Q$  subfamily (**pertussis toxin-insensitive**) or  $\beta\gamma$ -subunit dimers that are liberated following activation of G proteins in the  $G_i/G_o$  subfamily (**pertussis toxin-sensitive**). PLCs of the g-subfamily are attracted to the cytoplasmic face of the plasma membrane by certain phosphotyrosine-containing peptide motifs that are made by activated **tyrosine kinases**, and there they are tyrosine-phosphorylated: these two events cooperate to activate PtdIns(4,5) $P_2$  hydrolysis by the PLCgs (16). PLCds exist alongside PLCbs and PLCgs in animal cells, but only PLCds seem to be present in yeasts and plants. To some degree, the PLCds remain a puzzle, particularly because we still do not fully understand how they are activated. They can be activated by increases in cytosolic ( $Ca^{2+}$ ) that are fairly large but physiologically plausible, but there may be other activating factors yet to be discovered.

**Table 1. Stimuli that Transmit Signals into Target Cells at Least in Part Through Phosphoinositide-Based Signaling Pathways**

---

PLC-Catalyzed PtdIns(4,5) $P_2$  Hydrolysis

|   |   |
|---|---|
| Platelet-derived growth factor<br>(PDGF)  | Angiotensin I (AT <sub>1</sub> receptor)                |
| Adrenergic (α <sub>1</sub> receptor)  | Endothelins (3 receptors)                               |
| Histamine (H <sub>1</sub> receptor)   | Rhodopsin (invertebrates only)                          |
| Acetylcholine (muscarinic, 2 subtypes)  | Thromboxane A <sub>2</sub>                              |
| Vasopressin (V <sub>1</sub> receptor)   | Interleukin 8   |
| Oxytocin  | Glutamate (metabotropic receptors)                      |
| Substance P and other tachykinins   | Some taste receptors?                                   |
| PtdIns(3,4,5)P <sub>3</sub> Synthesis Catalyzed by Type I Phosphoinositide3-Kinases |   |
| Platelet-derived growth factor<br>(PDGF)  | The αβ/CD3 antigen receptor<br>complex of T lymphocytes |
| Insulin   |   |
| IGF-1   | βMetLeuPhe  |
| Colony-stimulating factor-1<br>(CSF-1)  |   |
| PtdIns3P 5-Kinase-Catalyzed PtdIns(3,5)P <sub>2</sub> Synthesis                     |   |
| Hyperosmotic stress (yeast)   |   |

---

A particularly clear illustration of the essential nature of the phospholipase C signaling pathway has emerged from genetic studies of the visual systems of insects (eg, *Drosophila*), which use rhodopsin-linked phospholipase C activation as their primary signaling pathway in response to the light activation of rhodopsin (17). Fly retinas show a striking set of genetic defects in this signaling pathway. They do not function properly if they lack: (1) a receptor-coupled phosphoinositidase C; (2) a diacylglycerol kinase and CDP-diacylglycerol synthase necessary for the resynthesis of PtdIns from the hydrolyzed PtdIns(4,5)P<sub>2</sub>; or (3) a PtdIns transfer protein that is assumed somehow to maintain the PtdIns supply to the photoreceptor membranes (see Table 2 and ref. 5).

**Table 2. Biological Effects for Eye Function in *Drosophila* of Mutations in, or Changes in, the Expression Level of Proteins Involved in Phosphoinositide-Based Signaling (See Ref. 3 for References to Original Observations)**

| Gene or Protein | Target Protein  | Phenotype  |
|-----------------|---|--|
| <i>NorpA</i>    | PLC-b, coupled to light-activated <a href="#">rhodopsin</a> | Eye elements fail to generate a light-triggered receptor potential |
| <i>RdgA</i>     | Diacylglycerol kinase                                       | Photoreceptors degenerate  |
| <i>RdgB</i>     | PtdIns transfer protein                                     | Light-induced degeneration of photoreceptors                       |
| <i>Cds</i>      | CDP-diacylglycerol  | Reduced signal gain  |

|                            |  |   |
|----------------------------|--|---|
|                            | synthase (a step in PtdIns biosynthesis) | during photoreception; light-dependent degeneration of photoreceptors   |
| Type-I                     | Organ-targeted Type I PI3K               | Eye size correlates positively with the concentration of expressed PI3K |
| PtdIns(4,5) $P_2$ 3-kinase |  |   |

---

### 3.2. Signaling through 3-Phosphorylation of PtdIns(4,5) $P_2$ to Phosphatidylinositol 3,4,5-Trisphosphate (PtdIns(3,4,5) $P_3$ )

In the late 1980s, eukaryote cells yielded up both PtdIns3  $P$  and phosphoinositide 3-kinases that can make PtdIns3  $P$  from PtdIns (generically abbreviated below as PI3Ks). Soon, however, it was found that some PI3Ks can use PtdIns, PtdIns4 $P$ , and PtdIns(4,5) $P_2$  as substrates and so can also make two other novel lipids: PtdIns(3,4) $P_2$  and PtdIns(3,4,5) $P_3$ . Moreover, neutrophils stimulated by the chemotactic attractant *f*MetLeuPhe (see [Chemotaxis](#)) were seen to make a lipid with the charge of a phosphatidylinositol trisphosphate (PtdIns  $P_3$ ): this proved to be PtdIns(3,4,5) $P_3$ .

Soon came recognition that many activated receptor tyrosine kinases (see [Growth Factors](#); [Platelet-Derived Growth Factor](#); [Insulin](#)) and a few G-protein-coupled serpentine receptors (especially those in hematopoietic cells) can very rapidly (within seconds) stimulate a subset of PI3Ks that phosphorylate PtdIns(4,5) $P_2$  to PtdIns(3,4,5) $P_3$  (see [Table 1](#)). These are now known as the Type I subfamily ([18](#)).

The PtdIns(3,4,5) $P_3$  generated in response to receptor stimulation, which is the most highly charged membrane phosphoinositide present in cells, is presumed to be located in the inner leaflet of the plasma membrane. This identified it as another potential receptor-generated signaling molecule. Several key observations soon combined to confirm this view ([18-20](#)). First, the PtdIns(4,5) $P_2$  3-kinase responsible for PtdIns(3,4,5)  $P_3$  synthesis in response to activated tyrosine kinases is a heterodimer of a catalytic subunit (p110) and a **SH2 domain**-containing regulatory subunit (p85). Immediately after the activation of some growth factor receptors, this p110/p85 PI3K becomes physically complexed with the activated receptor tyrosine kinase. This binding is mediated by the SH2 domain of p85, which recognizes a particular phosphotyrosine-containing motif (typically *P*Tyr-Met-Asp/Pro-Met) in the cytoplasmic “tail” of the activated receptor tyrosine kinase. Interaction with this motif localizes the kinase to the inner surface of the plasma membrane, where it becomes tyrosine-phosphorylated and is activated. Second, mutations of the “kinase insert” sequence of the [platelet-derived growth factor \(PDGF\)](#) receptor abolish this *P*Tyr-containing motif and reduce the growth-stimulating and survival-promoting activity of PDGF. Third, the tyrosine kinase activities of the activated [insulin](#) and insulin-like growth factor I receptors phosphorylate insulin receptor protein 1 (IRS-1), generating multiple PI3K-activating motifs. This activation of PI3K activity is essential both for the survival-promoting activities of these hormones and for a number of their classical actions on carbohydrate metabolism ([19](#), [20](#)). Fourth, overexpressed or constitutively active PI3Ks are growth-promoting and can be oncogenic ([18](#)). Finally, neutrophils contain a different Type I PtdIns(4,5) $P_2$ -directed PI3K, which is activated by the G protein  $\beta\gamma$ -subunit complexes that are released following activation of some G protein-coupled receptors (eg, the *f*MetLeuPhe receptor) ([21](#)).

The fungal metabolite wortmannin A was a key reagent in many of these studies. Interest in this compound was aroused when it was found to inhibit multiple, and apparently unrelated, cell responses to external stimulation. These included:

1. The rapid oxidative burst characteristic of stimulated neutrophils;
2. Insulin-stimulated glucose uptake into target tissues, notably adipose and skeletal muscle, which is achieved primarily by fusion of vesicles enriched in glucose-permeable channels with the plasma membrane;
3. Maintenance of cell survival in response to “growth/survival factors” such as IGF1; and
4. Activation of the formation of lamellipodia at the leading edge of migrating cells (see [Cytoskeleton](#)).

Once it was realized that wortmannin is a very potent inhibitor of the Type-I PI3Ks that make PtdIns(3,4,5) $P_3$ , this became an obvious candidate as the immediate activator of some or all of these responses. Although subsequent studies showed that wortmannin can be a less than scrupulously precise pharmacological reagent, other evidence has subsequently confirmed that many of the biological responses that are most sensitive to wortmannin inhibition are indeed consequences of PI3K activation. Used carefully, therefore, this compound has remained valuable for the initial identification of responses that might be consequences of receptor-stimulated PtdIns(3,4,5) $P_3$  synthesis.

Some details of how cells sense and respond to a PI3K-driven rise in PtdIns(3,4,5) $P_3$  concentration as a regulatory signal have been elucidated. Most notably, accumulation of PtdIns(3,4,5) $P_3$  very quickly leads to activation of protein kinase B (PKB, also known as c-Akt), a **serine-threonine protein kinase** and the product of a proto-oncogene (22). To become fully activated, cytosolic PKB must become firmly localized to the inner face of the plasma membrane and must be phosphorylated on two widely separated residues (Thr473 and Ser308). PKB physically associates with the plasma membrane, primarily through a direct interaction between its [pleckstrin homology \(PH\) domain](#) and the normal PtdIns(4,5) $P_2$  complement of that membrane, but this association may be strengthened by newly formed PtdIns(3,4,5) $P_3$ .

Receptor-stimulated formation of PtdIns(3,4,5) $P_3$  triggers phosphorylation of the membrane-bound PKB. In this process, one PtdIns(3,4,5) $P_3$ -activated protein kinase (phosphoinositide-dependent kinase 1; PDK1) phosphorylates PKB on Thr308, and another may phosphorylate it on Thr473 (Fig. 2; refs. 23, 24). The interaction that activates PDK1 seems to be between the newly formed PtdIns(3,4,5) $P_3$  and a PtdIns(3,4,5) $P_3$ -selective PH domain in PDK1.

A substantial number of other candidate target proteins avidly and selectively bind PtdIns(3,4,5) $P_3$ : these include the so-called centaurins (25, 26) and Bruton's tyrosine kinase, an enzyme essential for normal B lymphocyte development (27). Which, if any, of these serve as directly-controlled regulatory targets of this receptor-responsive polyphosphoinositide remains to be determined (20).

#### 4. Phosphoinositides and membrane trafficking

The earliest work on receptor-stimulated phosphoinositide turnover, in the 1950s and 1960s, focused primarily on its possible relationship to the triggering of secretory processes (3). The emphasis then moved to unraveling the central roles of phosphoinositides in signaling. Recently, however, there has been a remarkable resurgence of interest in how phosphoinositides contribute to intracellular membrane and protein trafficking (28).

#### 4.1. PtdIns(4,5)P<sub>2</sub> Synthesis and Exocytosis

It has long been known that **secretory vesicles** (eg, adrenal medullary chromaffin “granules”) support rapid PtdIns4P biosynthesis, and recently it was recognized that **exocytotic** secretory processes in yeast and in mammalian cells fail in the absence of a functional cytosolic “PtdIns transfer protein” (PITP). In yeast, *Sec14*, one of many *Sec* genes implicated in secretory protein transit early in the pathway between the [Golgi apparatus](#) and the plasma membrane, encodes a PITP and, in permeabilized PC12 cells (which mimic catecholamine secretion by adrenal medulla), a PITP is among the proteins that are needed to reconstitute ATP-driven exocytotic secretion.

PtdIns is made at the endoplasmic reticulum but is involved in signaling at the plasma membrane. Soon after the discovery of PITPs in the 1970s, it was recognized that these molecules might serve as “ferries” that move PtdIns around cells ([2](#)), an idea that was recently validated ([29](#), [30](#)). PITPs bind PtdIns and/or phosphatidylcholine, one molecule at a time, and shuttle these lipids across the cytosol from membrane to membrane. Mammals have two very similar PITPs (a and b), either of which can support PtdIns(4,5)P<sub>2</sub>-dependent signaling and exocytosis. Although the yeast protein encoded by *Sec14* is functionally quite similar, it has an unrelated amino acid sequence.

A second cytosolic protein that is essential for catecholamine secretion by PC12 cells is a PtdIns4P 5-kinase. Successful exocytosis, therefore, requires PtdIns to be delivered to the secretory vesicle membrane and phosphorylated to PtdIns(4,5)P<sub>2</sub> ([31](#)). Accumulation of a substantial amount of PtdIns(4,5)P<sub>2</sub> in the secretory vesicle membrane may somehow “prime” the vesicles, switching them to a state in which they can fuse with the plasma membrane in a subsequent Ca<sup>2+</sup>-activated exocytosis step.

#### 4.2. PtdIns3P and Membrane Trafficking

PtdIns3P, by contrast, is somehow implicated in vesicular protein sorting to the vacuole of yeast and the lysosomes of mammalian cells (which are functionally homologous organelles), and also in membrane fusion between elements of a membrane vesicle compartment (the “early” [endosomes](#)) that plays a central role in the import into eukaryotic cells of macromolecules such as the iron-carrier protein [transferrin](#). This first became apparent when the protein encoded by *vps34*, one of many genes whose disruption interferes with vesicular protein sorting in cells, was identified as a PtdIns-specific PI3K: Vps34p and other Vps34p-like PI3Ks make up the Type III PI3K subfamily. It is not clear how PtdIns3P contributes to the membrane fusion and fission events in these membrane trafficking processes. However, detailed studies of EEA1, a protein characteristic of early endosomes, have shown that a **zinc finger**like protein domain known as a “FYVE finger” serves as a highly selective PtdIns3P-binding domain; similar FYVE finger domains are present in several other mammalian and yeast proteins involved in membrane trafficking processes ([32](#), [33](#)).

The importance of the 3-phosphorylation of PtdIns by Type III PI3Ks was further emphasized when it was discovered that the PtdIns3 P that Vps34p makes is the substrate for the synthesis of phosphatidylinositol 3,5-bisphosphate (PtdIns(3,5)P<sub>2</sub>), the most recently discovered of the polyphosphoinositides ([34](#), [35](#)). Whereas PtdIns3P synthesis is a constitutive activity of unstressed cells—as it should be if it is involved in the continuous vesicle-mediated protein flows into cells from the exterior and from Golgi to lysosomes—PtdIns(3,5)P<sub>2</sub> synthesis is acutely regulated. In both fission and budding yeasts, hyperosmotic stress provokes a striking acceleration of PtdIns(3,5)P<sub>2</sub> synthesis ([35](#)). The exact function of this striking response is yet to be defined.

### 5. Inositol lipid-binding domains in proteins

One common theme that has emerged particularly strongly from many recent studies of the functions of inositol lipids and phosphates is the involvement of stereospecific, high-affinity interactions between particular polyphosphorylated inositol lipid or phosphate species and **protein**

**domains** whose primary function appears to be the selective recognition of these molecules. Many of these domains are members of a large family of so-called [pleckstrin homology \(or PH\) domains](#) (15, 36), but see also the mention above of ‘FYVE fingers.’

PH domains have a characteristic cluster of basic amino acid residues that, in each protein, ligates a particular subset of polyphosphorylated inositol derivatives [eg, (37)]. For example, the PH domains of PKB/Akt and of Bruton's tyrosine kinase preferentially bind PtdIns(4,5) $P_2$  and PtdIns(3,4,5) $P_3$ , respectively (23, 27). As the name suggests, the prototypic PH domain was found in pleckstrin: this is an abundant platelet cytosol protein that becomes heavily phosphorylated by protein kinase C downstream of receptors (eg, for [thrombin](#)) that activate phospholipase C but whose function is still not entirely clear. A major function of PH domains is to mediate noncovalent binding of proteins to the cytoplasmic surfaces of polyphosphoinositide-containing membranes—notably the plasma membrane, but also other structures such as exocytotic secretory vesicles—so as to localize them at their functional sites. Some such effects are likely to be constitutive, as when a PH domain-containing protein homes to the steady-state complement of plasma membrane PtdIns(4,5)  $P_2$ . In other circumstances, PH domains probably play a major regulatory role. For example, PtdIns(3,4,5)  $P_3$ -binding PH domains will relocate proteins to the plasma membrane, often for activation, only after cells have been exposed to a PtdIns(3,4,5) $P_3$ -generating stimulus.

## 6. The Functions of inositol phosphates

Although the biological function of Ins(1,4,5) $P_3$  is clear, why cells contain so many other inositol polyphosphates remains uncertain. That eukaryotic cells commit substantial genetic resources and metabolic energy to complex inositol polyphosphate interconversion pathways (see Fig. 4) surely indicates that at least some of these molecules have important cellular functions, but few of these are established with any certainty. Some of these are summarized in Table 3, and discussion of others can be found in recent reviews (5, 6, 15, 36, 42, 43).

**Table 3. Defined and Probable Functions of Some Inositol Polyphosphates and Glycerophosphoinositol Polyphosphates (See the Text for Further References)**

| Compound   | Function  | References |
|--|---|------------|
| <i>Ins(1,4,5)<math>P_3</math></i> : formed by hydrolysis of PtdIns(4,5) $P_2$                    | Signaling, through Ins(1,4,5) $P_3$ -stimulated $Ca^{2+}$ mobilization  | (13, 14)   |
| <i>Ins(1,4,5,6)<math>P_4</math></i> : formed from Ins(1,4,5) $P_3$ by Ins(1,4,5) $P_3$ 3-kinases | Modulation of $Ca^{2+}$ entry following depletion of intracellular $Ca^{2+}$ stores? Modulation of GTPase-activating protein                              | (15)       |
| <i>Ins(3,4,5,6)<math>P_4</math></i>  | Slowly accumulates during stimulation of PLC-coupled receptors. Inhibits a $Ca^{2+}$ /calmodulin-dependent and kinase-regulated epithelial $Cl^-$ channel | (38)       |
| <i>Ins(1,3,4,5,6)<math>P_5</math></i> (often   | Modulation of $O_2$ binding to avian  | (39)       |

mimicked experimentally and reptilian **hemoglobins** with  $\text{InsP}_6$ )

**GroPIns4 P** (and/or **GroPIns(4,5) P<sub>2</sub>**) Accumulate in Ras-transformed cells, inhibit [adenylate cyclase](#) (40)

**Several pyrophosphate derivatives of  $\text{InsP}_6$**  (shown as  $\text{InsP}_7$  and  $\text{InsP}_8$  in Fig. 4) Widespread, rapid metabolic turnover of pyrophosphate groups, high free energy of hydrolysis, exact structures vary between cell types (41)

---

## Bibliography

1. R. Parasarathy and F. Eisenberg (1986) *Biochem. J.* **235**, 313–322.
2. R. H. Michell (1975) *Biochim. Biophys. Acta* **415**, 81–147.
3. L. E. Hokin and M. R. Hokin-Neaverson (1989) *Biochim. Biophys. Acta* **1000**, 465–469.
4. IUPAC (1989) *Biochem. J.* **258**, 1–2.
5. R. H. Michell (1997) *Essays Biochem.* **32**, 31–47.
6. P. J. Hughes and R. H. Michell (1993) *Curr. Opin. Neurobiol.* **3**, 383–400.
7. P. P. Ongusaha, P. J. Hughes, J. Davey, and R. H. Michell (1998) *Biochem. J.* **335**, 671–680.
8. Y. Koga, M. Nishihara, H. Morii, and M. Akagawa-Matsushita (1993) *Microbiol. Rev.* **57**, 164–182.
9. L. O. Martins, R. Huber, H. Huber, K. O. Stetter, M. S. DaCosta, and H. Santos (1997) *Appl. Environ. Microbiol.* **63**, 896–902.
10. J. A. Creba, C. P. Downes, P. T. Hawkins, G. Brewster, R. H. Michell, and C. J. Kirk (1983) *Biochem. J.* **212**, 733.
11. Y. Nishizuka (1990) *Nature* **334**, 661–665.
12. M. J. Berridge (1987) *Ann. Rev. Biochem.* **56**, 159–194.
13. K. Mikoshiba (1997) *Curr. Opin. Neurobiol.* **7**, 339–345.
14. M. J. Berridge (1993) *Nature*, **361**, 315–325.
15. N. Divecha and R.F. Irvine (1995) *Cell* **80**, 269–278
16. S. B. Lee and S. G. Rhee (1995) *Curr. Opin. Cell Biol.* **7**, 183–189.
17. R. Ranganathan, D. M. Malicki, and C. S. Zuker (1995) *Ann. Rev. Neurosci.* **18**, 283–317.
18. B. Vanhaesebrock, S. J. Leever, G. Panayatou, and Waterfield (1997) *Trends Biochem. Sci.* **22**, 267–272.
19. L. Stephens, T. R. Jackson, and Hawkins, P. T. (1993) *Biochim. Biophys. Acta* **1179**, 27–75.
20. A. Toker and L. C. Cantley (1997) *Nature* **387**, 673–676.
21. L. R. Stephens, A. Eguinoa, H. Erdjument-Bromage, M. Lui, F. Cooke, J. Coadwell, A. S. Smrcka, M. Thelen, K. Cadwallader, P. Tempst, and P. T. Hawkins (1997) *Cell* **89**, 105–114.
22. B. Marte and J. Downward (1997) *Trends Biochem. Sci.* **22**, 355–358.
23. D. Stokoe, L. R. Stephens, T. Copeland, P. R. Gaffney, C. B. Reese, G. F. Painter, A. B. Holmes, F. McCormick, and P. T. Hawkins (1997) *Science* **277**, 567–570.
24. D. R. Alessi and P. Cohen (1998) *Curr. Opin. Genet. Devel.* **8**, 55–62.
25. L. P. Hammonds-Odie, T. R. Jackson, A. A. Profit, I. J. Blader, C. W. Turck, G. D. Prestwich, and A. B. Theibert (1996) *J. Biol. Chem.* **271**, 18859–18868.

26. K. Tanaka, S. Imajoh-Ohmi, T. Sawada, R. Shirai, Y. Hashimoto, S. Iwasaki, K. Kaibuchi, Y. Kanaho, T. Shirai, Y. Terada, K. Kimura, S. Nagata, and Y. Fukui (1997) *Eur. J. Biochem.* **245**, 512–519.
27. M. Fukuda, T. Kojima, H. Kabayama, and K. Mikoshiba, (1996) *J. Biol. Chem.* **48**, 30303–30306.
28. P. De Camilli, S. D. Emr, P. S. McPherson, and P. Novick (1996) *Science* **271**, 1533–1539.
29. S. Cockcroft (1997) *FEBS Lett.* **410**, 44–48.
30. K. W. A. Wirtz (1997) *Biochem. J.* **324**, 353–360.
31. J. C. Hay, P. L. Fiset, G. H. Jenkins, K. Fukami, T. Takenawa, R. A. Anderson, and T. F. J. Martin (1995) *Nature* **374**, 173–177.
32. J.-M. Gaullier, A. Simonsen, A. D'Arrigo, B. Bremnes, H. Stenmark, and R. Aasland (1998) *Nature* **394**, 432–433.
33. V. Patki, D. C. Lawe, S. Corvera, J. V. Virbasius, and A. Chawla (1998) *Nature* **394**, 433–434.
34. C. A. Whiteford, C. A. Brearley, and E. T. Ulug (1997) *Biochem. J.* **323**, 597–601.
35. S. K. Dove, F. Cooke, M. R. Douglas, L. G. Sayers, P. Parker, and R. H. Michell (1997) *Nature*, **390**, 187–192.
36. M. Fukuda and K. Mikoshiba (1997) *BioEssays* **19**, 593–603.
37. M. Hyvonen and M. Saraste (1997) *EMBO J.* **16**, 3396–3404.
38. I. I. Ismailov, C. M. Fuller, B. K. Berdiev, V. G. Shlyonsky, D. J. Benos, and E. Barrett (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10505–10509.
39. H. Yanase, S. Cahill, J. J. M. Dellano, L. R. Manning, K. Schneider, B. T. Chait, K. D. Vandegriff, R. M. Winslow, and J. M. Manning (1994) *Protein Sci.* **8**, 1213–1223.
40. M. Falasca, A. Carvelli, C. Iurisci, R-G. Qiu, M. H. Symons, and D. Corda (1997) *Mol. Biol. Cell.* **8**, 443–453.
41. C. Albert, S. T. Safrany, M. E. Bembenek, K. M. Raddy, K. K. Reddy, J. R. Falck, M. Brocker, S. B. Shears, and G. W. Mayr (1997) *Biochem. J.* **327**, 553–560.
42. F. S. Menniti, K. G. Oliver, J. W. Putney, and S. B. Shears (1993) *Trends Biochem. Sci.* **18**, 53–56.
43. N. Sasakawa, M. Sharif, and M. R. Hanley (1995) *Biochem. Pharmacol.* **50**, 137–146.

## Insertion Sequence Elements

Insertion sequence (IS) elements are small (0.75 to 1.5 kbp) [transposable elements](#) found in bacteria (1). They usually encode only a [transposase](#) gene and special sequences at the tips of the element on which the recombinase acts to move the element from place to place. IS elements thus structurally resemble, for example, [P elements](#) from *Drosophila*, the widespread Tc1/mariner elements, and Ac elements from maize, although IS elements are generally much smaller. Like P, Tc1/mariner, and Ac elements, IS elements alter DNA in the host genome upon insertion and can cause gene inactivation.

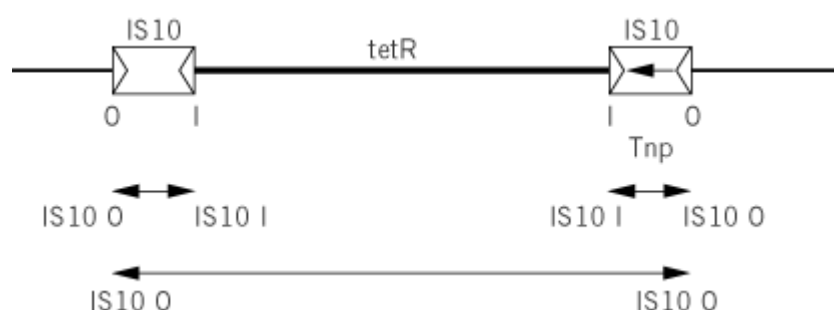
IS elements were among the first nonmaize transposable elements discovered, being identified in bacteria in the late 1960s as the cause of highly polar mutations in bacterial [operons](#) and also as components of translocatable [antibiotic-resistance](#) segments. It was possible to manipulate bacterial



DNA molecules at this time because of their relatively small size, something that was not possible for other organisms such as maize, and the direct characterization of the DNA of these polar mutations revealed that they resulted from an insertion of a discrete piece of DNA.

Some IS elements also associate to form what are called *composite elements* in which two ISs flank a region of DNA that encodes a selectable trait, for example an antibiotic-resistance determinant (Fig. 1). The action of transposase at the two external edges of the IS elements can move the entire IS–drug resistance gene–IS element from place to place (Fig. 1). Each IS element can also move independently, although this does not result in movement of the antibiotic resistance determinant. IS elements are found as individuals at various positions in bacterial chromosome.

**Figure 1.** IS elements can associate to form composite transposons. IS elements encode information to promote their translocation, a transposase (arrow) and special sequences at their ends (triangles). Composite transposons are formed by two IS elements flanking a drug resistance determinant, here tetracycline. When the inside (I) and outside (O) ends of an individual IS element are the substrates for recombination, the IS element alone translocates. When the outside (O) ends of the two ISs flanking the antibiotic resistance determinant are the substrates for recombination, the entire IS–drug resistance–IS segment translocates.



Particularly well-studied IS sequences are IS10 (2) and IS50 (3). The composite transposon Tn10 consists of two IS10 elements flanking a DNA segment that encodes a tetracycline-resistance determinant. The composite transposon Tn5 is composed of two IS50 elements flanking a **kanamycin**-resistance encoding segment. Both of these elements translocate by a cut-and-paste mechanism, in which the element is excised from the donor DNA and then inserts into the target DNA. They and their derivatives have also been invaluable reagents in the genetic analysis of bacteria (4). The strategy of using transposon mutagenesis to tag genes physically to allow **gene** isolation has been a very valuable technique.

Another class of IS elements—members of the IS3 family that includes IS3 and IS911—have an interesting strategy for expressing transposase. The element contains two open reading frames, OrfA and OrfB. OrfA encodes the DNA-binding determinant that directs this protein to the ends of the transposon, but OrfA has no transposase activity. The actual transposase OrfAB results from a **frameshifting** event near the end of OrfA that joins OrfA and OrfB (5). The amount of transposase, and thus the frequency of transposition, is determined by the amount of frameshifting. These elements also appear to transpose by an unusual mechanism (6, 7). Transposition begins by cleavage at one end of the transposon, exposing a 3'OH end of the element, followed by an intramolecular attack of that 3'OH end to just outside the 5' end of that same strand; note that this is the same chemistry as is involved in other transposition reactions, the only difference being that the joining reaction is intramolecular, rather than intermolecular. By either another such single-strand cleavage and intramolecular joining, or by replication, a double-stranded DNA circle results in which the IS ends are closely juxtaposed, with just a few nucleotides separating them. This “circle junction” version of the IS element then interacts with the target site and is inserted by breakage at 3' ends of the transposon and joining the exposed ends to the target DNA.

It should be appreciated that although some of the well-studied elements and widely used bacterial elements such as Tn10 and Tn5 contain IS elements, certainly not all bacterial transposable elements contains IS sequences. Tn3, [Mu phage](#), and Tn7 are all examples of bacterial elements that lack IS sequences, and these elements do not associate to form composite transposons.

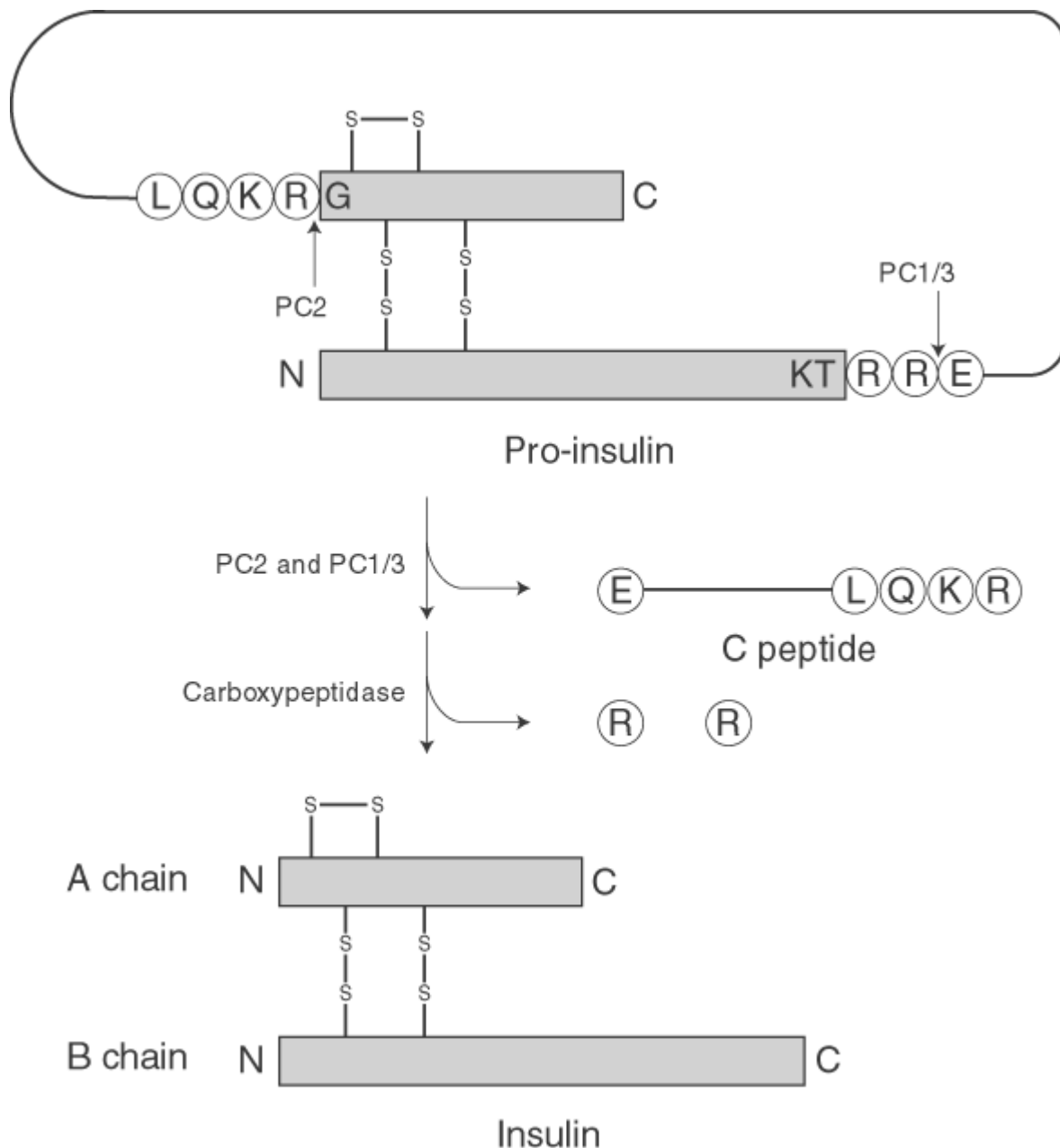
### Bibliography

1. E. Ohtsubo and Y. Sekine (1996) *Curr. Top. Microbiol. Immunol.* **204**, 1–26.
2. N. Kleckner, R. Chalmers, D. Kwon, J. Sakai, and S. Bolland (1996) *Curr. Top. Microbiol. Immunol.* **204**, 49–82.
3. W. S. Reznikoff (1993) *Annu. Rev. Microbiol.* **47**, 945–963.
4. C. M. Berg and E. Berg (1995) In *Mobile Genetic Elements* (D. J. Sherratt, ed.), IRL Press, Oxford, England, pp. 38–68.
5. M. Chandler and O. Fayet (1993) *Mol. Microbiol.* **7**, 497–503.
6. P. Polard and M. Chandler (1995) *Genes Dev.* **9**, 2846–2858.
7. B. Ton-Hoang, P. Polard, and M. Chandler (1998) *EMBO J.* **17**, 1169–1181.

### Insulin

Insulin is a [peptide hormone](#) that is secreted by the b-cells of the islets of Langerhans in the pancreas. The principal role of this hormone is the regulation of glucose metabolism and [homeostasis](#) by regulating the concentrations and activities of numerous [enzymes](#) involved with catabolism and anabolism. The molecule (MW 5808) consists of two [polypeptide chains](#), A and B, linked by two [disulfide bonds](#) (Fig. [1](#)). Reduction of these disulfide bonds results in the loss of biological activity.

**Figure 1.** The processing of pro-insulin (top) to the two polypeptide chains of insulin linked by disulfide bonds (bottom). The sequence of the pro-insulin polypeptide chain is shown, by one-letter abbreviations, only near the sites of cleavage by the convertases PC1 / 3 and PC2 and the carboxypeptidase; the remainder of the connecting C peptide that is removed proteolytically is shown as the solid line only.



In humans, the insulin **gene** is located on chromosome 11. Several **alleles** coding both normal and pathological forms of insulin have been identified. An individual's genes are co-dominantly expressed. Insulin is synthesized initially on **polyribosomes** of the rough **endoplasmic reticulum** (ER) as pre-pro-insulin (see **Polyproteins**), a single-chain precursor protein of molecular weight 11,500. The *N*-terminal **signal peptide** of 23 or 24 primarily **hydrophobic** residues is cleaved immediately by signal peptidase after penetration of the peptide into the ER lumen. The resulting pro-insulin (MW 9000) undergoes **protein folding** and disulfide bond formation. It is then transported to the **Golgi apparatus** and is sorted into **clathrin-coated secretory granules**. **Posttranslational modification** by **membrane-bound endoproteinas**, pro-hormone convertases PC-2 and PC-1 / 3, and a carboxypeptidase converts proinsulin into equimolar amounts of insulin and free C peptide (Fig. 1). Other than to facilitate proper folding, the C peptide has no known function. The insulin stored in the mature secretory granules assumes a hexameric **quaternary structure** with two bound zinc ions (1). In response to plasma glucose levels, the mature granule translocates to the cell membrane, fuses with the membrane, and discharges its contents. Other than glucose, several other substances (eg, mannose, ribose, arginine, and leucine) also stimulate release.

Upon secretion by the pancreas, insulin is directly infused into the portal vein and the liver, where it exerts its primary metabolic effects. Monomeric insulin generates intracellular effects by binding to its plasma-membrane receptor (see [Peptide Hormones](#)). The precise signal transduction mechanisms (see [Signal Transduction](#)) are still not completely known (2).

The principal target tissues of insulin are liver, muscle, and fat. In most tissues, insulin increases the number of plasma membrane glucose **transporters**, such as GLUT 4 (3, 4), but glucose uptake in the liver is also increased because of the increased activity of regulatory enzymes that play a role in glycolysis. Glycogenolysis is inhibited. Glucokinase, **phosphofructokinase** -1, pyruvate kinase, and insulin-dependent **phosphodiesterase** are all involved. In addition, specific **phosphatases** increase the activity of glycogen synthase and thus lead to increased glycolysis. Glucose-6-phosphatase is inhibited. Overall, hepatocyte glucose and its phosphorylated forms are increased, resulting in lower blood levels of glucose. The actions are complex, because all the metabolic pathways must be coordinated and maintained to effect homeostasis.

The insulin-receptor gene is located on chromosome 19. The insulin receptor (see [Hormone Receptors Peptide Hormones](#)) is a **tyrosine kinase** and comprised of two identical extracellular  $\alpha$ -subunits (MW 135 kDa) and two identical transmembrane  $\beta$ -subunits (MW 95 kDa). The  $\alpha$ -subunit contains the insulin-binding domain, and the membrane-spanning  $\beta$ -subunit transduces the signal. After the hormone-receptor complex is formed, it is internalized by [endocytosis](#) mechanisms. The receptor is recycled back to the membrane, while the insulin molecule is degraded (5).

Mutations occur in some forms of diabetes mellitus in both the genes for insulin and its receptor (6), as well as defects in ER and Golgi processing (7). The mutations may result in impaired binding, impaired transmembrane signaling, altered endocytosis and recycling, or impaired transport of the receptor to the cell membrane. Diabetes is thus a catabolic disease caused by an apparent insulin deficiency. There are genetic components to both type I (insulin dependent) and type II (insulin independent) diabetes. Moreover, type I diabetes can result from an **autoimmune** response that destroys the pancreatic  $\beta$ -cells and insulin production (8).

Therapy consists of the subcutaneous administration of insulin. [Protein engineering](#) of human insulin production in *Escherichia coli* now provides the therapeutic drug. Use of insulin-secreting cell lines as **gene therapy** is presently only in the early stages of development (9).

## Bibliography

1. H. R. Petty (1993) *Molecular Biology of Membrane: Structure and Function*, Plenum, New York, pp. 255–262.
2. P. Stralfors (1997) *BioEssays* **19**, 327–335.
3. O. Ezaki (1997) *Biochem. Biophys. Res. Commun.* **241**, 1–6.
4. J. E. Rice, C. Livingstone, and G. W. Gould (1996) *Biochem. Soc. Trans.* **24**, 540–546.
5. G. D. Holman and M. Kasuga (1997) *Diabetologica* **40**, 991–1003.
6. D. A. Brooks (1997) *FEBS Lett.* **409**, 115–120.
7. D. Accili (1995) *Diabetes-Metabolism Rev.* **11**, 47–62.
8. A. A. Alzaid (1996) *Acta Diabetologica* **33**, 87–99.
9. M. Tiedge and S. Lenzen (1995) *Experimental and Clinical Endocrinology and Diabetes* **103S**, 46–55.

## Insulin-Like Growth Factors

The insulin-like growth factors (IGFs) are among the most ubiquitous class of peptide [growth factors](#) in the body and are found in most tissues and in blood. Two types, IGF-I and -II, have been purified from serum, each with molecular masses of approximately 7.6 kDa. Liver is a major site of their expression, both in fetal and postnatal life, although almost all tissues have been shown to express these peptides in the human and animal fetus ([1](#), [2](#)); this suggests a predominantly autocrine or paracrine role, in which the IGF affects primarily the same cell or only nearby cells. In the fetus, the most abundant isomer is IGF-II, but in some species, such as rat, but not humans, IGF-II is absent from adult serum and is replaced by IGF-I. In humans, IGF-II persists throughout life, although the relative abundance of IGF-I increases postnatally. A high-affinity type 1 IGF receptor is ubiquitous in developing tissues and recognizes IGF-I with an order of magnitude greater binding affinity than it does IGF-II. Consequently, IGF-I is a more potent **mitogen** with a half-maximal response on most cell types at approximately 1 to 3 nM. The IGFs also have proven roles in the control of cell differentiation and are preventive of [apoptosis](#) in many cell types. In postnatal life, synergy between endogenous IGF-I and trophic endocrine hormones facilitates endocrine glandular function. For example, thyroid-stimulating hormone synergizes with IGF-I to support thyroxine synthesis, while a synergy of IGFs with follicle stimulating hormone allows sex steroid production in the ovary ([3](#)). The type 1 IGF receptor has an intracellular [tyrosine](#) kinase domain that is capable of **phosphorylating** the [insulin](#) receptor substrate proteins (IRS-1 and -2). IGFs are able to activate the proto- **oncogene** *ras* in many cell types, which subsequently signals gene [transcription](#) changes in cell [nuclei](#), often by a mitogen-activated protein (MAP) kinase pathway, but also by the activation of phosphoinositol-3 kinase (see [Inositol Lipids and Phosphates](#)), culminating in the activation of the [transcription factors](#) *fos*, *jun*, and *myc*. An additional high-affinity receptor that specifically binds IGF-II, the type II or cation-independent **mannose-6-phosphate receptor**, is ubiquitous but has no consistent intracellular signaling pathway or biological endpoint.

The IGFs are seldom found in free form but are complexed to one of six distinct classes of specific binding protein, termed IGFBP-1 to -6. These are found both in serum and in extracellular fluids, and they serve not only as carrier proteins to extend the biological half-life of the ligands but also modulate their biological actions by either interacting or competing with the type 1 IGF receptors. While all six IGFBPs have a conserved core structure, differences in their *N*- and *C*-termini confer individual relative binding affinities for IGF-I and -II and an ability to interact with both the [extracellular matrix](#) and the cell surface ([4](#)). Two of these IGFBPs, IGFBP-1 and -2, contain an **integrin**-binding motif that allows binding to the cell surface  $\alpha_5\beta_1$  integrin, which is the [fibronectin](#) receptor. The ability of these IGFBPs to potentiate IGF action is related to their integrin-binding activity, which may facilitate an advantageous presentation of the ligand to its high-affinity receptor. All IGFs except IGFBP-1 also contain heparin-binding **domains**, allowing binding to sulfated glycosaminoglycans in the extracellular matrix and on the cell membrane. The majority of IGF-I and -II in blood is carried on IGFBP-3. From the second half of fetal development, IGFBP-3, together with an IGF molecule, is associated with an acid-labile subunit in the circulation to generate a tertiary complex of 150 kDa. In this form, IGFs cannot leave the circulation; the fraction accessible to extracellular fluids is carried by IGFBP-1, -2, -4, and -5.

A large proportion of the IGF:IGFBP complexes in extracellular fluids and stored within the extracellular matrix are probably inaccessible to the cell-surface receptors. Their availability depends on modification of the IGFBPs by specific [proteinases](#), resulting in a reduced binding affinity for IGFs. Such proteinases have been identified for IGFBPs-2 to -5. An IGFBP-3-degrading proteinase appears in maternal serum from the second trimester of human pregnancy until term, which reduces the amount of IGF carried by IGFBP-3 and increases its transcapillary passage in association with other IGFBPs ([5](#)). A naturally occurring tissue proteinase can also remove the three *N*-terminal amino acid residues from IGF-I, resulting in a much reduced binding affinity for IGFs. Thus, while IGF-I and -II are present in the circulation, this serves predominantly as an extracellular store.

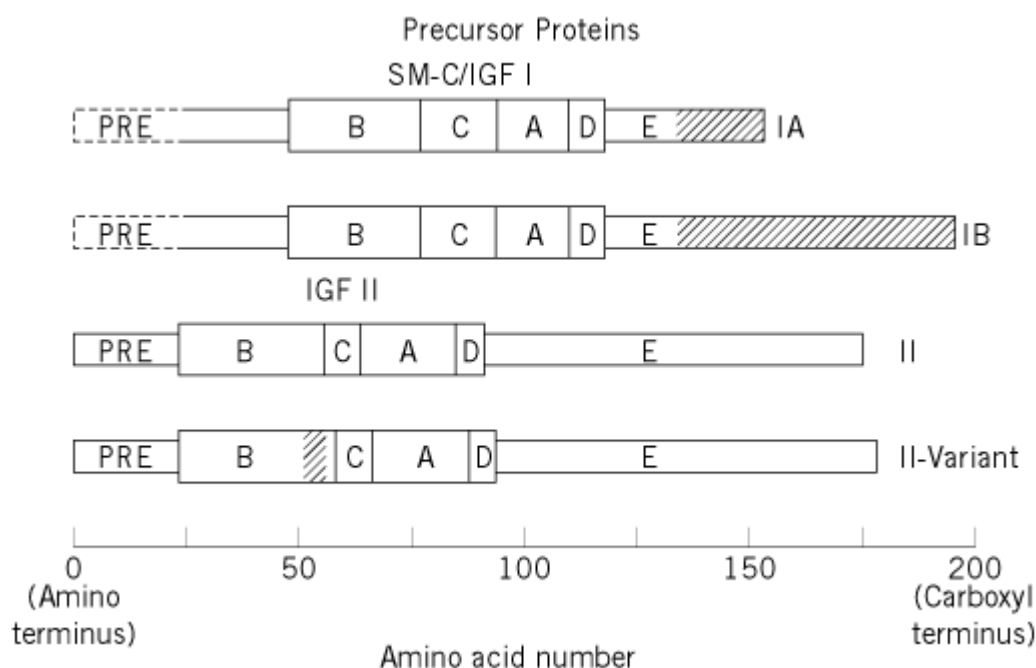
Controlled **proteolysis** of IGF-BPs and extracellular matrix molecules is likely to be the key regulatory step in the bioavailability and subsequent actions of IGF-I and -II.

## 1. IGF gene structure

### 1.1. IGF-I

The human IGF-I **gene** is located on the long arm of chromosome 12 (6). It comprises at least five exons (see [Introns, Exons](#)) spread over more than 90 kbp (7). The gene encodes two precursor forms of IGF-I, represented by transcripts with distinct 3' regions that arise by [alternative splicing](#) of the [pre-messenger RNA](#) (Fig. 1). Each transcript includes sequences of the three 5' exons (exons 1, 2, and 3) but only one of the two 3' exons (exons 4 and 5) of the gene. The resulting precursors, IGF-Ia and IGF-Ib, have identical amino-terminal segments of 48 residues that are encoded by exons 1 and 2. The identities of IGF-Ia and -Ib are maintained in their mature IGF-I regions, which comprise 70 amino acid residues encoded by exons 2 and 3, and in the first 16 residues at their C-terminus, encoded by exon 3. They are different in the remainder of their C-terminus, where IGF-Ia has an additional 19 amino acid residues encoded by exon 5 and IGF-Ib has an additional 61 residues encoded by exon 4.

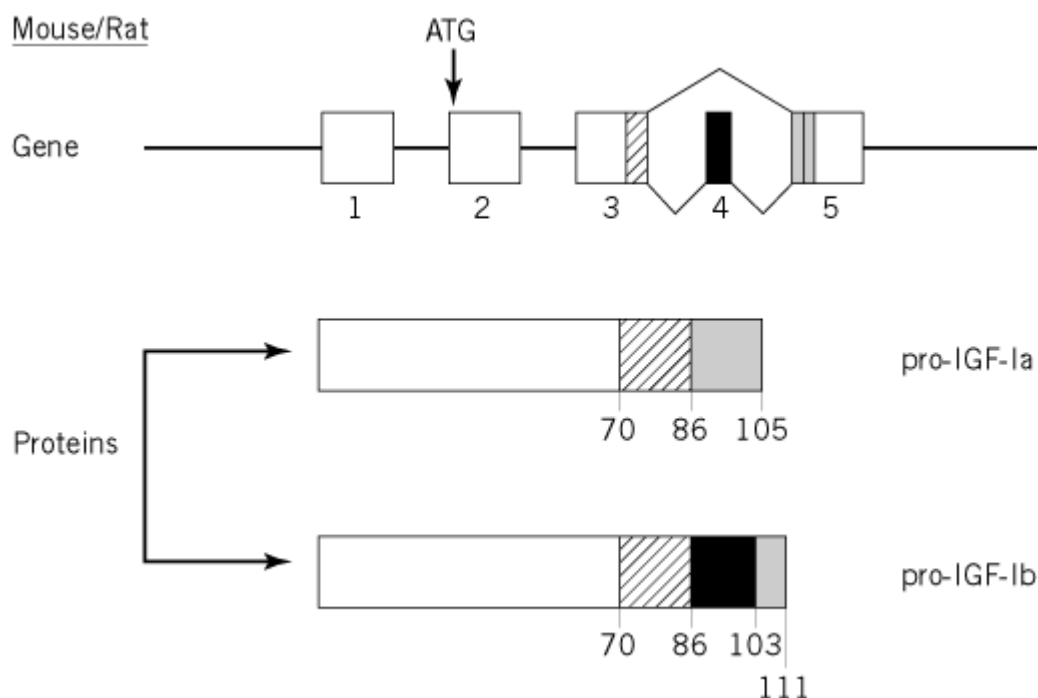
**Figure 1. Schematic representation of human IGF-I and -II precursor proteins.** The hatched areas in IGF-Ia and IGF-Ib highlight the alternative sequences arising from alternative splicing (exon 3 to exon 4 in IGF-Ia, and exon 3 to exon 5 in IGF-Ib). The hatched area in the predicted IGF-II variant highlights the area at the juncture of the B and C domains close to an RNA splice site, where differences in the variant human IGF-II molecule tend to occur.



The rat IGF-I gene has six exons, of which five are **homologous** to those of the human gene, with extensive identity of nucleotide sequences (8). The two alternative **signal sequences** (I and IA) encode variant amino-termini for the pre-pro-IGF-I signal peptide, depending on which translation [initiation codon](#) is used *in vivo*. IGF-I cDNA clones with alternative 5'-untranslated regions have also been identified in the mouse (9) and other species. Like the human gene, the rat IGF-I gene encodes two precursors, represented by distinct mRNA transcripts. The rat IGF-Ia precursor is analogous to the human product, having a C-terminus encoded by exons 3 and 5. However, the C-terminal region of rat IGF-Ib is dissimilar to that of human IGF-Ib (which is encoded by exons 3 and

4), comprising sequences encoded by exons 3, 4, and 5 (Fig. 2).

**Figure 2. Representation of alternative splicing resulting in the formation of alternative IGF-I precursor peptides in the rodent.** Exons are represented by the open boxes. Introns are represented by a thin line. The rodent (and human) IGF-I gene contains five exons. It is assumed that translation is initiated at the indicated ATG codon. Exon 2 contains the signal peptide and residues 1 to 25 of the B domain. Exon 3 contains the remainder of the B domain; the complete C, A, and D domains; and residues 70 to 86 of the E domain (hatched). Exons 4 and 5 contain alternate sequences for the remainder of the E domain and the 3'-untranslated region. Splicing of exon 3 to exon 5 gives rise to the IGF-Ia precursor (protein-coding region of exon 5 is stippled). The IGF-Ib precursor arises by splicing exon 3 to exon 4 (protein-coding region of exon 4 is solid), which is spliced to exon 5. Pro-IGF-Ib terminates earlier in exon 5 than pro-IGF-Ia (indicated by the vertical line in the stippled area of exon 5).

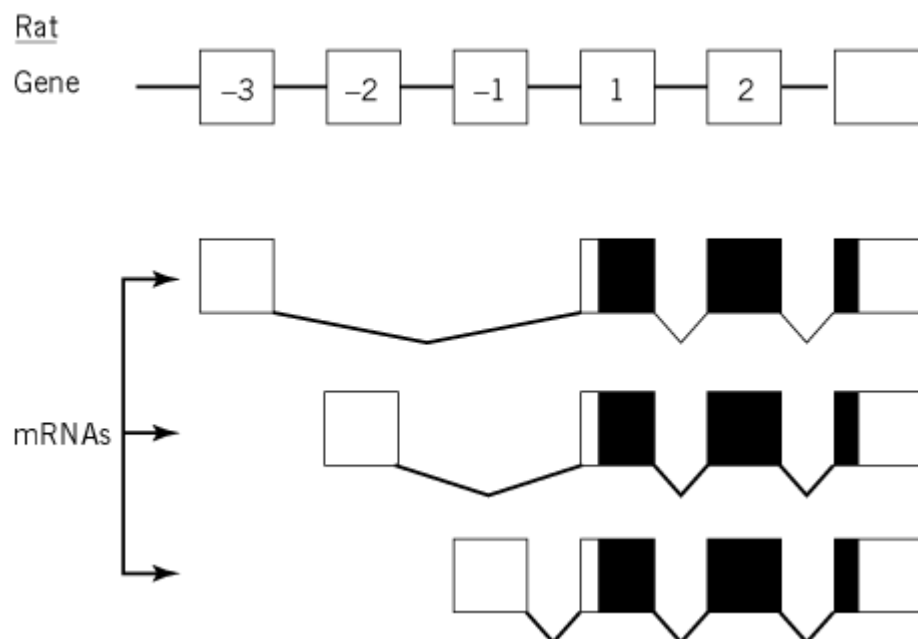


Messenger (mRNA) species encoding IGF-I of at least three size classes are detected by **Northern blot** analysis, ranging from approximately 0.8 kb to 7.8 kb. At least three factors contribute to this size heterogeneity. Alternative splicing pathways can account for two classes of transcripts encoding the precursors IGF-Ia and IGF-Ib. Additionally, transcripts have been identified with distinct 5'-untranslated sequences in mouse and rat and with variable lengths of 3'-untranslated sequences. Using probes that distinguish IGF-Ia and -Ib mRNAs in tissues of humans, it was shown that transcripts encoding both precursors are represented within each of the three size classes of IGF-I transcripts. In all rodent tissues, transcripts encoding IGF-Ia are more abundant than those encoding IGF-Ib, although the relative abundance of these mRNA species differs between hepatic and nonhepatic tissues; transcripts encoding IGF-Ib represent approximately 13% of IGF-I mRNAs in the liver of adult rats but only 2 to 5% of the total elsewhere. When [growth hormone](#) was given to hypophysectomized rats, the levels of both mRNA types were elevated coordinately in nonhepatic tissues, but in liver IGF-Ib mRNA was raised threefold more than that of IGF-Ia mRNA. The functional significance of the two pre-pro-IGF-I precursors remains unknown. Most of the heterogeneity in IGF-I transcript size is due to differing lengths of 3'-untranslated sequences. These arise from multiple [polyadenylation](#) signal sequences previously within exon 5 of the rat gene, and they may affect the stability of IGF-I mRNAs, providing a mechanism for [posttranscriptional regulation](#) of expression. Evidence for multiple polyadenylation sites in IGF-I gene also exists for the mouse, human, and sheep proteins.

## 1.2. IGF-II

The IGF-II genes in human, rat, and mouse each encode a 180-residue precursor (pre-pro-IGF-II), consisting of a 24-residue *N*-terminal signal peptide, a 67-residue mature IGF-II, and an 89-residue *C*-terminal pro-peptide. As with pre-pro-IGF-I, the pre- and pro-peptide fragments are removed during protein biosynthesis. Sequences for the cDNA diverge within the 5' and 3' noncoding regions, indicating the existence of IGF-II transcripts differing in their untranslated sequences. The human IGF-II gene is found on the short arm of chromosome 11, where it is closely linked to the [insulin](#) locus (10). The rodent IGF-II gene is similarly adjacent to the insulin gene on chromosome 1 in the rat, or chromosome 7 in the mouse (11). The genes for both human and rat IGF-II cover at least 30 kbp and 12 kbp, respectively, and contain multiple exons and **promoters**. Exons E4, E5, and E6 encode pre-pro-IGF-II protein, while exons E1, E2, and E3 are noncoding and used differentially, forming alternative, 5'-untranslated regions of different RNA molecules (Fig. 3). The human IGF-II gene uses three additional noncoding exons. Consensus eukaryote promoter sequences, such as a [TATA box](#) and [GC boxes](#) (SPI recognition sequences), are present in the regions upstream of human and rat exons E2 (promoter P2) and E3 (promoter P3). The promoter region (P1) upstream of rat, human, and mouse exon E1 lacks TATA and GC boxes, and the transcription start sites are heterogeneous. The region upstream of human exon I (P1) lies within 1.4 kbp of the insulin gene. A polyadenylation signal, conforming to the **consensus sequence** AATAAA, is located 20 nucleotides from the 3' end of exon 6, and further potential polyadenylation sites are also present in exon 6 of both the human and rat IGF-II genes.

**Figure 3. Schematic representation of the rat IGF-II gene.** Transcripts arising at each of three promoters and terminal protein-coding regions are solid, untranslated regions are open.



IGF-II mRNAs differ predominately in their 5'-untranslated sequences, with the major transcripts in humans approximately 4.8 kbp, and in rodents approximately 3.8, 4.6, and 8.6 kbp. These arise from the promoters P1, P2, or P3 of the noncoding exon E1, E2, or E3 spliced to the three coding exons, E4, E5, and E6, with transcription termination at the polyadenylation site. Transcripts derived from P2 are most abundant in human tissues, whereas those from P3 are most abundant in rodent tissues. Transcripts arising from P1 are least abundant in both species. Termination of transcription at the polyadenylation site (consensus AATAAA) in exon 6 gives rise to short transcripts of 1 to 3 kb found in rat and human tissues.

IGF-II cDNAs encoding variant proteins have been found in humans (12). These transcripts arise by alternative splicing of a weak splice-acceptor site upstream of the major splice-acceptor site of exon 5. The additional nine nucleotide sequences predict the replacement of Ser-29 with the tetrapeptide Arg-Leu-Pro-Gly, and this IGF-II

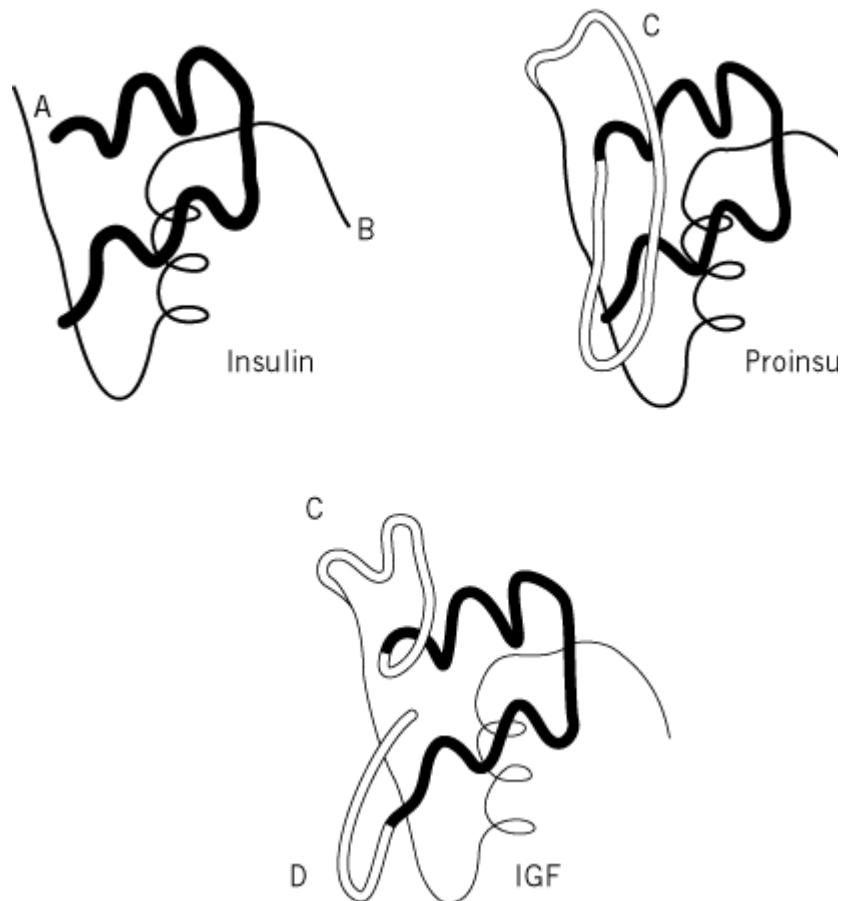


(13). The expression of IGF-II is regulated by p53 through interaction with the IGF-II P3 promoter, and tumorigenesis when variant p53 forms are present (14).

## 2. Structure of IGF-I and -II peptides

The IGFs are single-chain polypeptides that show amino acid sequence homology with pro-insulin. This allows the recognition of equivalents of the insulin B, C, and A domains within the IGFs, together with the D domain, which is not found in pro-insulin (Fig. 4). All amino acid residues common between species IGFs with the exception of residue A21, which is Asn in insulin but Ala in the IGFs. Conservation of the tertiary structure similar to that of insulin, and the hydrophobic core of insulin, consisting of A2-B24 Phe, is also conserved. The most obvious differences from pro-insulin are in the C domain. The insulin and IGF-I are identical. IGF-I consists of 70 amino acid residues, with a predicted molecule weight of 7440 detected by mass spectrometry (15). Its isoelectric point is 8.8. Human IGF-II consists of 67 amino acid residues with a molecular weight of 7469, which was reported as 7440 by mass spectrometry. IGF-II has an isoelectric point of 8.8. The difference between predicted and detected sizes may be due to covalent changes caused by the extraction procedure.

**Figure 4. Structural homology of insulin, pro-insulin, and IGF.** Schematic representation of the three-dimensional structure of insulin (with additional C chain) and the insulin-like growth factors (with C chain and D chain) demonstrating their close structural homology.



There are two different IGF-I precursors, one derived from exons 1, 2, 3, and 5 (IGF-Ia) and the other (IGF-1b). They differ in their C-terminal precursor sequence (E domain) and result in identical mature peptides. IGF-I lacking the first three N-terminal residues [des(1-3)IGF-I] was reported by Sara et al. (17) in 1977. Des(1-3)-IGF-I has a 100-fold reduction in affinity for IGF-BPs-1 and -2 as compared to IGF-I. The affinity is only two to three times lower than is intact IGF-I. Studies on receptor interaction and biological activity

and recombinant des(1–3)IGF-I. The truncated molecule has a much increased biological activity than intact IGF-I. The *in vivo* half-life of the peptide is shorter than that of intact IGF-I in the rat, due to this. In growth-deficient animals it stimulates growth in length more than does recombinant IGF-I (18).

Variation also exists in IGF-II molecules. Jansen et al. (19) isolated and characterized a cDNA derived from liver encoding an IGF-II variant in which B29 Ser is replaced by the sequence Arg–Leu–Pro–Gly. This cDNA was obtained from the same library that contained the cDNA encoding classical IGF-II and was splicing of the pre-mRNA. The peptide has been isolated from human serum and found to constitute 1% of total IGF-II (13). Its functional characteristics have shown a two- to threefold reduction in affinity for purified IGF-II receptor in placenta and reduced mitogenicity. Another IGF-II variant described by Zumstein et al. (20) is a 10-kDa IGF-II in which Ser is substituted by Cys–Gly–Asp, and it contains the first 21 residues of the C-terminal precursor. This substitution does not occur at an intron-exon junction, allelic variation must be responsible for its occurrence. It has been found in human spinal fluid, serum, and brain (21–23). Presumably these molecules represent unprocessed precursor forms.

### 3. IGF binding proteins

The amino acid sequences at the NH<sub>2</sub>- and COOH-termini of the six IGFBPs are highly conserved, with the NH<sub>2</sub>-terminus of exon I, including 12 cysteines, identical in human and rat IGFBP-1–5 (4) (Fig. 5). Human and rat IGFBP-1 share a Gly–Cys–Gly–Cys–Cys segment that is conserved in other IGFBPs. Species differences also occur in the COOH-terminal segment. In the COOH-terminal region, 18 residues, including four cysteines, out of 28 are identical in human and rat. Related molecules designated by some as IGFBPs-7 to -10 have been identified, but they lack biologically meaningful homology. The middle region of the IGFBP proteins, corresponding to the C-terminal part of exon 2, is not conserved in different IGFBPs. It includes N-glycosylation sites that are used for the attachment of additional cysteine residues in IGFBP-4. IGFBP-1, -2, and -3 proteins are encoded by four exons; IGFBP-4 is encoded by five exons, the last three comprised exclusively of 3'-untranslated sequences. There is no evidence that alternative splicing could account for the observed variation. The IGFBPs do not appear to be synthesized as larger precursor proteins.

**Figure 5. Alignment of the amino acid sequences of the six rat IGFBPs.** Amino acid residues are shown in one-letter code. Asterisks indicate maximal homology alignment. Identical residues at corresponding positions in more than three of the IGFBPs are shown. Residues are marked by asterisks, except for the additional two cysteines (underlined) present only in IGFBP-4. The RG is indicated by an arrow. The internal amino acid sequences that could not be aligned were grouped by parentheses.

```

* * * * *
Rat IGFBP-1 APQPWHCAPCTAERLELCPVP-AS-CPEISRPA GCGCCPTCALPLGAACGVAI
Rat IGFBP-2 EVLFRCPCTPERLAACGPPDPAP-CAELVREPGCGCCSVCARQEGEACGVYI
Rat IGFBP-3 GAGAVGAGPVVRCPCDARALAQCAPPPTAPACTELVREPGCGCCLTCALREGDACGVYI
Rat IGFBP-4 DEAIHCPCSEKRLARCRPPVG---CEELVREPGCGCCATCALGLGMPCGVYI
Rat IGFBP-5 LGSFVHCEPCDEKALSMC-PPSP-LGC-ELVKEPGCGCCMTCALAEGQSCGVYI
Rat IGFBP-6 ALAGCPGC----(GPGVQEEEDAGSPAD)----GCAETGGCFRREGQPCGVYI

```

```

(VLEPPPPATSSLSGSQHEEAKAAVASEDELAESPEMTEEQLLDSFHLMAPSREDQPILWNAISTYSSMRAREITDLI
(EKRRVGATPQQVADSEDDHSEGLVENHVDGTMNMLGGSSAGRKPPKSGMKELAVFREKVNEQHRQMGKGAKHLSLI
(ANASAASNLSAYLPSQPSPGNTTESEEDHNAGSVESQVVPSTHRVTD SKFHPLH SKMEVI IKGQARDSORYKVDYES
(TELSEIEAIQESLQTS DKDESEHPNNSFNPCSAHDHRC LQKHM AKVRDRSKMKVVGTPREEPRPVPQG)
(LNEKSYGEQTKIERDSREHEEPTTSEMAEETYS PKVFRPKHTRI SELKAEAVKKERRKLTQSKFVGGAEHTAHPRI
(QRARGPSEETT KESKPHGGASRPRDRDRQKNPRTSAAPIRPSVQD GEMG)

```

```

* * * * *
PCQRELYKVLERLAAAQOKA--GD-E-IYKFYLPNCNKNGFYHSKQ CETSLDGEAGLCWCVYPWSGKKIPGSLETRG
PCQQELDQVLERI STMRLPDDRGP L EHL YSLHIPNCDKHGLYNL KQCKMSLNGRGE CWCVNPNNTGKPIQGAPTIRG
PCRREMEDTLNHLKFLNVLSPRGV-----HIPNCDKKG F YKKKQCRPSKGRKRGF CWCVDKYGQPLPGYDTKGKD
SCQSELHRALERLAASQ--S-RTH-EDLFIIPNCDRNGMFHPKQCHPALDQGRGKCWCVD RKTGVKLP GGLEPKG
PCRRHMEASLQEFKASPRMVPRV-----YLPNCDRKG F YKRKQCKPSRGRKRGI CWCVDKY-GNKLPGMEYVDG
PCRRHLD SVLQQLQTE-VF--RGGANGL---YVPNCDLRG F YRKQCRSSQGNRRGPCWCVDPM-GQPLPVSPDGQG

```

The major carrier protein for IGFs in blood is IGFBP-3 (24). It is also able to bind an acid-labile non-predominant 150-kDa IGFBP complex in adult serum. Concentrations of ALS in human plasma are higher, and have a longer half-life in the circulation when it is associated with the 150-kDa complex. The human IGFBP-3 gene consists of four protein-coding exons homologous to those of the IGFBP-1 and IGFBP-2 genes, and a fifth exon is located on chromosome 7, 20 kbp from the IGFBP-1 gene, in tail-to-tail orientation. The transcript is 132 bp upstream from the adenine–thymine–guanine (ATG) translation initiation codon. IGFBP-3 is abundant in liver but was also present at high levels in kidney, stomach, heart, adrenal, ovary, and contains a TATA box 30 bp upstream from the transcription initiation site, a GC element (97 bp upstream) that overlaps the GC element.

The human IGFBP-1 gene is 5.2 kbp in length and contains four protein-coding exons. Its mRNA is abundant in decidua in late pregnancy but was not detected in placenta. The other major site where a 1.5-kb mRNA transcript is found. IGFBP-1 mRNA was increased 100-fold in rat liver two days after streptozotocin and was reversed after three days of insulin treatment. The increase in hepatic IGFBP-1 in kidney results primarily from increased transcription of the IGFBP-1 gene. Insulin treatment decreases and IGFBP-1 mRNA levels within one hour, most likely by direct effects on gene transcription. IGFBP-1 has insulin-potentiating effects *in vitro*. The presence of an Arg–Gly–Asp (RGD) sequence at the COOH-terminus allows interaction with cell-surface receptors for proteins in the extracellular matrix. Serine phosphorylation–dephosphorylation is a major biological action of IGFBP-1. Dephosphorylation decreases the affinity of IGFBP-1 for IGF-I, promoting the formation of an inactive complex.

The IGFBP-2 gene has been localized to human chromosome 2 and to the homologous mouse chromosome 2. The human genes are organized into four exons, each of which contains coding sequences. The human gene is 3.5 kbp in size, with a large size of intron 1 (27 kbp). In Northern blots of RNA from fetal rats, a 1.5-kb IGFBP-2 mRNA was found in kidney, and lung and was less abundant in intestine, muscle, heart, and skin. In midgestation rat brain, the epithelium of the choroid plexus, suggesting that IGFBP-2 might participate in the transport of IGFs. IGFBP-2 mRNA in rat liver increased 10-fold after 48 hours of fasting, which was reversed after one day of insulin treatment. mRNA was increased 10- to 20-fold in streptozotocin-treated rats, and this was reversed by insulin treatment. IGFBP-2 has an RGD sequence at its COOH-terminus. Human IGFBP-2 binds IGF-II with 10-fold greater affinity than IGF-I.

similar for rat IGFBP-2. IGFBP-2 inhibits the mitogenic effects of IGFs *in vitro*.

The IGFBP-4 gene is located on chromosome 17, and IGFBP-4 mRNA has been estimated as 2.5 kb. Hybridization is seen in liver, and a weaker hybridization signal exists with RNA from adrenal, testes, hypothalamus, and brain cortex. In addition to the 18 cysteine residues found in other IGFBPs, IGFBP-4 has the nonconserved middle region of the molecule. Glycosylated and nonglycosylated forms of IGFBP-4 are present in serum.

The IGFBP-5 gene is located on human chromosome 5. In adult rat tissues, a 6-kb IGFBP-5 mRNA is present at high levels in lung, heart, and stomach. A 1.7-kb transcript, presumably differing in the length, has also been reported. IGFBP-5 binds IGF-II with 10-fold higher affinity than IGF-I.

IGFBP-6 was purified from human cerebrospinal fluid and human and rat serum. It lacks two of the cysteine residues found in other IGFBPs. Human, but not rat, IGFBP-6 contains a potential *N*-glycosylation site. IGFBP-6 binds IGF-II with high affinity and IGF-I with greater than 70-fold lower affinity. The human IGFBP-6 gene has been cloned. IGFBP-6 mRNA is widely expressed in tissues of adult male rats, with a 1.3-kb transcript that is abundant in intestine, adrenal, kidney, stomach, spleen, heart, brain, and liver.

#### 4. Selective deletions of IGF genes

Gene deletion by homologous [recombination](#) is a powerful way of examining the morphological and physiological consequences of growth factor deficiency. Some of the most dramatic findings from gene targeting have been obtained with the use of homologous recombination. Homologous recombination has been used to disrupt either the IGF-I, the IGF-II, or the type 1 IGF receptor gene. "knockouts" have then been obtained. Deletion of the IGF-I gene by homologous recombination resulted in animals with a birthweight about 60% of normal, of which some died within six hours of birth (25). Some of the mutant females were infertile due to a failure of ovarian follicular development. Using a similar strategy, the IGF-II gene is parentally **imprinted** and is only transmitted from the male allele in the majority of tissues, including the placenta, plexus and meninges, where both alleles of the gene are active (26). IGF-II-deficient homozygotes have been obtained in animals lacking IGF-I, demonstrating that both isomers have a role in prenatal growth. However, in humans, there is a single published record of a natural deletion of the IGF-I gene, which was associated with growth retardation and short stature in the child (27).

Deletion of the type 1 IGF receptor, which is primarily responsible for both the mitogenic signaling by IGF-I and IGF-II, yielded homozygous animals that were only 45% of normal weight at delivery and died within a few days to a failure to breathe and probably resulted from a widespread muscle hypoplasia, including that of the diaphragm. There was an increase in neuronal cell density in the spinal cord and brainstem of the mutant animals, while the skull, the stratum spinosum, and bone ossification was delayed by about two fetal days. Double gene knockout of the type 1 receptor genes resulted in a similar phenotype to that found after deletion of the receptor alone. Deletion of the type 1 receptor yielded a subgroup of animals with only 30% of normal birth weight at term and gross growth retardation. This suggests that an additional receptor to the type 1 form may also contribute to IGF-II signaling. Deletion of the type 2 receptor is deleted in a naturally occurring gene deletion identified by the lack of the imprinted locus on chromosome 15 at embryonic stage (29). Whether this receptor can contribute to IGF-II signaling *in vivo* is not clear.

The aforementioned series of studies demonstrates that neither IGF-I nor IGF-II is crucial for key morphological development, but that they act as "true" growth factors, contributing to the expansion of stem cell populations and differentiation.

#### Bibliography

1. V. K. M. Han, P. K. Lund, D. C. Lee, and A. J. D'Ercole (1988) *J. Clin. Endocrinol. Metab.* **66**, 1000-1004.
2. A. L. Brown, D. E. Graham, S. P. Nissley, D. J. Hill, A. J. Strain, and M. M. Rechler (1986) *J. Biol. Chem.* **261**, 11500-11504.
3. D. J. Hill and J. Hogg (1989) In *Clinical Endocrinology and Metabolism, Perinatal Endocrinology*, Bailliere Tindall, London U.K., pp. 579-625.
4. R. H. McCusker and D. R. Clemmons (1992) In *The Insulin-Like Growth Factors, Structure and Function*, Plenum Press, New York, pp. 1-10.

Schofield, ed.), Oxford University Press, Oxford, U.K., pp. 110–150.

5. L. C. Giudice, E. M. Farrell, H. Pham, G. Lamson, and R. G. Rosenfeld (1990) *J. Clin. Endocrinol.* **71**, 103–110.
6. J. E. Brisenden, A. Ullrich, and U. Francke (1984) *Nature* **310**, 781–784.
7. P. Rotwein, K. M. Pollock, D. K. Didier, and G. G. Krivi (1986) *J. Biol. Chem.* **261**, 4828–4832.
8. A. Shimatsu and P. Rotwein (1987) *Nucleic Acids Res.* **15**, 7196–7197.
9. G. I. Bell, M. M. Stempien, N. M. Fong, and L. B. Rall (1986) *Nucleic Acids Res.* **14**, 7873–7880.
10. K. L. O'Malley and P. Rotwein (1988) *Nucleic Acids Res.* **16**, 4437–4446.
11. P. Rotwein and L. J. Hall (1990) *DNA and Cell Biol.* **9**, 725–735.
12. P. N. Schofield and V. E. Tate (1987) *Development* **101**, 793–803.
13. B. Hampton, W. H. Burgess, D. R. Marshak, K. J. Cullen, and J. F. Perdue (1989) *J. Biol. Chem.* **264**, 1030–1034.
14. L. Zhang, F. Kashanchi, Q. Zhan, S. Zhan, J. N. Brady, A. J. Fornace, P. Seth, and L. J. Helman (1990) *J. Biol. Chem.* **265**, 1030–1034.
15. J. L. Van den Brande, C. M. Hoogerbrugge, K. Beyreuther, P. Roepstorff, J. Jansen, and S. C. van den Berghe (1988) *Endocrinol.* **122**, 683–695.
16. P. Rotwein (1986) *Proc. Nat. Acad. Sci. U.S.A.* **83**, 77–81.
17. V. R. Sara et al. (1986) *Proc. Nat. Acad. Sci. U.S.A.* **83**, 4904–4907.
18. C. Gillespie, L. C. Read, C. J. Bagley, and F. J. Ballard (1990) *J. Endocrinol.* **127**, 401–405.
19. M. Jansen, F. M. A. van Schaik, H. van Tol, J. L. Van den Brande, and J. S. Sussenbach (1985) *J. Biol. Chem.* **260**, 1030–1034.
20. P. P. Zumstein, C. Luthi, and R. E. Humbel (1985) *Proc. Nat. Acad. Sci. U.S.A.* **82**, 3169–3172.
21. G. K. Haselbacher and R. Humbel (1982) *Endocrinology* **110**, 1822–1824.
22. G. K. Haselbacher, M. E. Schwab, A. Pasi, and R. E. Humbel (1985) *Proc. Nat. Acad. Sci. U.S.A.* **82**, 3169–3172.
23. L. K. Gowan, B. Hampton, D. J. Hill, R. J. Schlueter, and J. F. Perdue (1987) *Endocrinology* **121**, 1822–1824.
24. S. Rajaram, D. J. Baylink, and S. Mohan (1997) *Endocrine Rev.* **18**, 801–831.
25. J-P. Liu, J. Baker, A. S. Perkins, E. J. Robertson, and A. Efstratiadis (1993) *Cell* **75**, 59–72.
26. T. M. De Chiara, A. Efstratiadis, and E. J. Robertson (1990) *Nature* **345**, 78–80.
27. K. A. Woods, C. Camacho-Hubner, D. Barter, A. J. Clark, and M. O. Savage (1997) *Acta Paediatr.* **86**, 1030–1034.
28. J. Baker, J-P. Liu, E. J. Robertson, and A. Efstratiadis (1993) *Cell* **75**, 73–82.
29. D. P. Barlow, R. Stoger, B. G. Herrmann, K. Saito, and N. Schweifer (1991) *Nature* **349**, 84–87.

### Suggestions for Further Reading

30. W. H. Daughaday and P. Rotwein (1989) *Endocrine Rev.* **10**, 68–91.
31. J. S. Sussenbach (1989) *Prog. Growth Factor Res.* **1**, 33–48.
32. P. N. Schofield (ed.) (1992) *The Insulin-Like Growth Factors: Structure and Biological Function*. Oxford University Press, U.K.

## Integrases

Integrases represent a subset of DNA (poly)nucleotidyltransferase [enzymes](#) each of which catalyzes the insertion of a specific **extrachromosomal** segment of **DNA** into the [genome](#) of its host. They form central elements in the replicative cycles of several prokaryotic and eukaryotic **viruses** and in the intracellular [translocation](#) of [retrotransposons](#). The best studied of these are the integrases

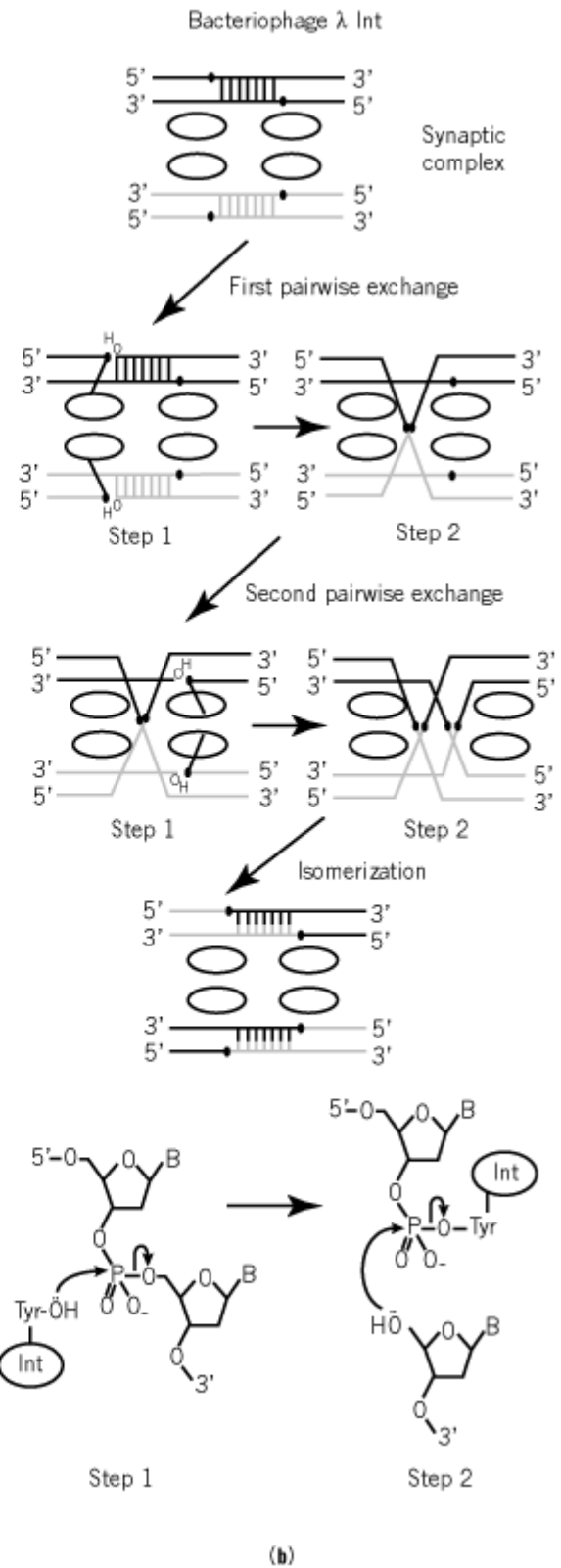
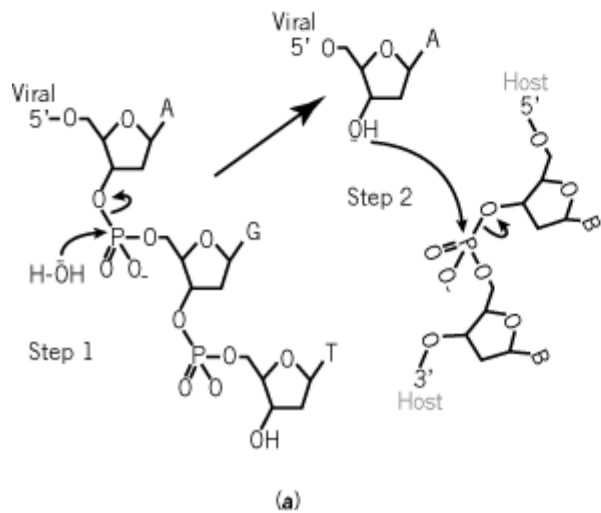
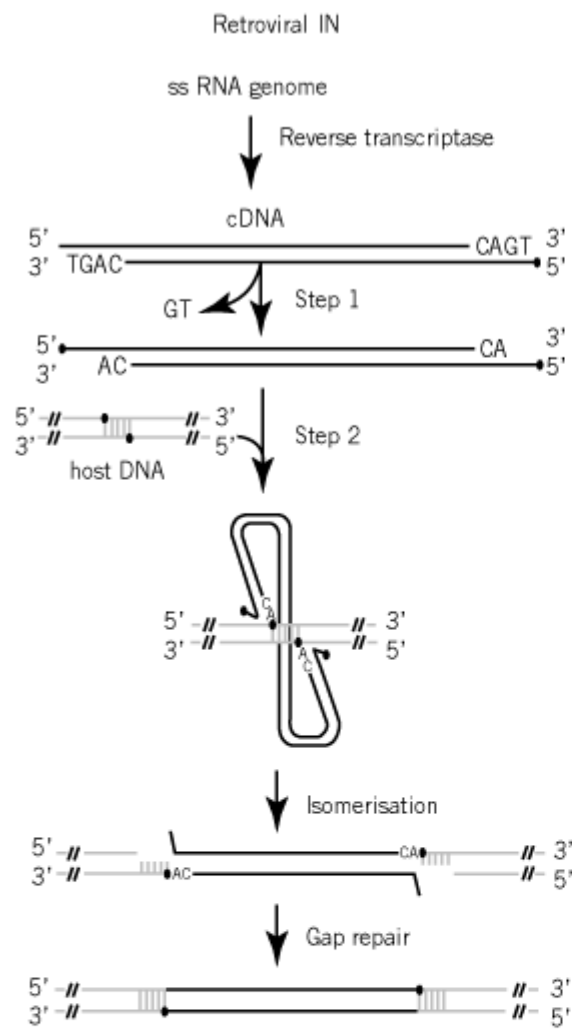
encoded by the [pol gene](#) of various retroelements, such as the retroviridae, and those of temperate bacteriophages, most notably [lambda phage](#). The different strategies employed to break and reunite the host DNA with the incoming polynucleotide distinguish two distinct families of integrases. A key difference is whether or not the reaction takes place by forming a covalent integrase-DNA intermediate. In the lambda Int family of site-specific recombinases, the integrase forms a transient covalent linkage with the 3'-ends of the bacteriophage and host DNA at the points of recombination. In contrast, retroviral integrase (IN) catalyzes DNA insertion via a noncovalent intermediate, using a process that is chemically similar to the “cut and paste” mechanism of classical transpositional [recombination](#). Unraveling DNA integration in molecular terms will necessitate understanding just how the underlying chemical steps are harnessed by each distinct integrase nucleoprotein complex to meet the different DNA site specificities and topologies demanded by their diverse genetic rearrangements.

## 1. Retroviral Integrases

### 1.1. DNA Integration Mechanism

DNA integration is a critical step in the replicative cycle of all single-stranded (ss) RNA retroviruses, including notably **HIV-1**, the causative agent of AIDS. Following **reverse transcription**, the resultant [complementary DNA](#) copy of the retroviral ssRNA genome must be inserted into a chromosome of the infected host cell to elicit a productive infection. The retroviral integrase (IN) catalyzes two consecutive nucleophilic substitution reactions in this process (Fig 1a; Steps 1 and 2) at the phosphodiester bonds that ultimately form the sites of attachment between the viral and host DNA. The first reaction, known as 3'-OH end processing, results in the nucleolytic expulsion of the two external nucleotides from both 3'-ends of the linear viral cDNA. Here, the 5'-phosphoryl group of the leaving dinucleotide is transferred directly to the electronegative group of the attacking nucleophile, generally donated by a [water](#) molecule. In the second step (polynucleotidyl transfer), the recessed 3'-OH groups are in turn coordinated as nucleophiles to attack specific positions on the phosphodiester backbone of the host DNA. In line with the displacement of a host 5'-phosphoryl group, a direct transesterification with the attacking viral 3'-OH end results in the insertion of the viral DNA at this site.

**Figure 1.** Comparison of retroviral cDNA and bacteriophage  $\lambda$  DNA integration. Both mechanisms involve two consecutive nucleophilic substitution reactions (Steps 1 and 2), shown side-by-side at bottom of figure. (a) Retroviral IN. Step 1: Endonucleolytic 3'-OH end processing. The two outermost bases of the cDNA adjacent to a subterminal CA dinucleotide are removed from both 3'-termini, yielding a staggered intermediate with 5'-protrusions. Step 2: Polynucleotidyl transfer. Each processed 3'-OH end is reattached in a staggered fashion to opposite strands of the host DNA (for HIV-1 at phosphodiester bonds positioned 5 bp apart). Gap repair, probably by host proteins, results in duplication of the host DNA (shown in gray) at the proviral junction. (b) Bacteriophage  $\lambda$  Int. Step 1: Breakage of partner strands. A pair of Int protomers (ovals) attack the scissile phosphodiester bonds (black dots) at parallel positions on the juxtapositioned recombination strands, resulting in the formation of covalent 3'-phosphotyrosine linkages and free 5'-OH ends. Step 2: Polynucleotidyl transfer (resolution). Each 5'-OH end rejoins with the 3'-phosphoryl of its facing partner strand by branch migration and resolution of the covalent protein-DNA intermediate. Recombination intermediates with the target DNA (shown in gray) are drawn to emphasize the crossover sites [after (20)] rather than its topological structure, which more accurately corresponds to a **Holliday junction** intermediate.



Two notable hallmarks of this type of DNA integration mechanism persist at the borders of the resultant proviral–host DNA junction. First, the provirus is two bases shorter than the preintegrated cDNA, where its 3'-boundaries with the host DNA are delineated at each end by a highly conserved CA-3' dinucleotide. Secondly, a region of 4 to 6 bp of host DNA is duplicated at the junction of proviral attachment that arises from the staggered integration of each viral end on opposite stands of the host DNA. The length of this duplication is specific to each retrovirus rather than its host.

Although retroviral DNA integration is clearly nonrandom, it may occur at any one of a large number of potential sites within the genome, which are dictated principally by the accessibility and structure of the [chromatin](#) (1). Yeast retrotransposons ([Ty elements](#)) integrate much more precisely, despite an essentially identical DNA integration mechanism, at or near nonessential genes or transcriptionally “silent” regions of the genome. Ty3 cDNA is drawn into the required locus by [protein–protein interactions](#) between a component of the Ty3 integration complex, possibly the Ty3 IN itself, and a polIII promoter-bound protein complex composed minimally of the [transcription factors](#) TFIIB and TFIIC, thus securing the DNA integration functions at that site (2). A similar, but as yet unsubstantiated, proposition has been put forward for HIV-1 IN, which interacts functionally with integrase interactor 1 (INI 1), a homologue of the yeast transcription factor SNF5, which could potentially direct integration into regions of open chromatin (3).

Unlike bacteriophage integrases, the IN proteins of retroelements cannot excise the integrated viral template from the host genome. Instead, viral propagation occurs via new generations of transcribed ssRNA progeny, which are packaged into virions with retroviral proteins and released to infect a new host cell. Nevertheless, IN can resolve the partial integration intermediates formed in this process and catalyzes the excision of the viral DNA from “one-ended” intermediates of a Y-type structure by a process named “disintegration” (4). During excision of the viral part, the broken target DNA strand is resealed by a similar phosphoryl transfer mechanism to that of the forward integration reaction, although the “reverse” process differs in the sense that it displays a relaxed specificity for viral sequences.

## 1.2. The IN Protein

Retroviral cDNA (>8 kb) exists in the cytoplasm of infected cells as a high molecular weight preintegration complex (PIC) and copurifies with several retroviral, and at least some host proteins. One of these is the high mobility group (**HMG**) protein HMG1(Y) (5), which among other global activities bends DNA and stimulates the concerted integration activity of purified IN (6). Increasing evidence suggests that such accessory proteins contribute important architectural features necessary to guide the integration functions of the PIC (7). Nonetheless, purified IN alone encompasses all of the minimal requirements for catalyzing both nucleophilic substitution (end processing and transfer) reactions with synthetic viral termini *in vitro* (reviewed in Ref. 8).

Although the IN-catalyzed functions exhibit some degree of specificity for the nucleotide sequence (s) of the viral termini and for the CA dinucleotide in particular, purified IN does not bind preferentially to its cognate sequence. The factors within the virion that direct IN to the two distal termini of the nascent cDNA and coordinate their concerted integration at a discrete host target site remain to be established. Neither function is reproduced efficiently *in vitro* with the purified protein.

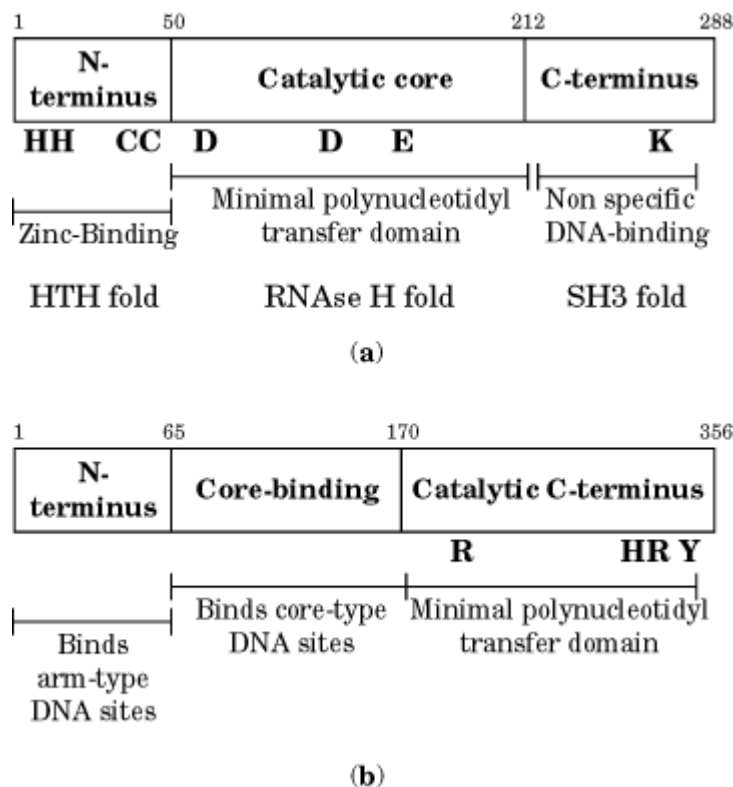
Structure-function studies on deleted derivatives of the integrase indicate that the protein is organized into three major **domains** (Fig. 2a). Although crystallizing the full length protein has proven difficult because of its inherent aggregative properties, now the structures of each of the isolated domains of the human immunodeficiency virus type-1 (HIV-1) integrase have been determined independently.

1. The N-terminal region (residues 1 to 50) bears a highly conserved HHCC motif (His-X<sub>3-7</sub>-His-X<sub>23-32</sub>-Cys-X<sub>2</sub>-Cys) that binds Zn<sup>2+</sup> stoichiometrically (9). Contrary to the organization of a classical **zinc finger**, its structure more closely resembles that of helical [DNA-binding proteins](#) that possess a [helix-turn-helix \(HTH\) motif](#) and is most akin to that of the trp repressor (10) (see [TRP Operon](#)). The helical bundle arrangement is ordered by Zn<sup>2+</sup>, which is tetrahedrally coordinated by the conserved His and Cys residues. The binding of Zn<sup>2+</sup> stimulates multimerization and the catalytic functions of the integrase (9, 11). The N-terminus is essential for 3-OH processing and integration but is dispensable for the disintegration reaction, which minimally requires the core region.



2. The central **protease**-resistant core domain (residues 50 to 212) comprises a minimal [active site](#) for the polynucleotidyl transfer reaction. This domain carries a triad (D,D,E or Asp,Asp,Glu) of essential catalytic residues (organized Asp-X<sub>39-58</sub>-Asp-X<sub>35</sub>-Glu) that coordinate a divalent cation (Mg<sup>2+</sup>, and also Mn<sup>2+</sup>) necessary for phosphoryl transfer. The core domain binds and catalyzes the resolution of Y-structure disintegration intermediates but cannot alone catalyze the forward reactions of the integration mechanism with linear DNA substrates. Several residues of the core are highly conserved in retroviral and retrotransposon integrases and in analogous prokaryotic transposases (12). Most remarkable is its structural [homology](#) to the [ribonuclease H](#) and RuvC enzymes, comprising a central 5-stranded [beta-sheet](#) flanked by **alpha-helical** regions (13). Despite the fact that these proteins do not catalyze DNA integration or transposition, nonetheless, they are involved in similar transesterification reactions using [aspartic acid](#) and [glutamic acid](#) catalytic residues. The catalytic residues superimpose well on the corresponding structures, indicating a general structural fold for a topologically related [superfamily](#) of DNA nucleotidyltransferases, all of which catalyze nucleophilic substitutions at the phosphodiester bonds of polynucleotides.
3. The C-terminal region is less well conserved both in sequence and in size and exhibits an independent nonspecific DNA-binding activity. Several integrases also possess C-terminal nuclear localization signals that, it is proposed, direct the active uptake of the large PIC structure through the undisturbed nuclear membrane [eg (14)] (see [Nuclear Import, Export](#)). The organization of its 5 [beta-strands](#) into a b-barrel consisting of two antiparallel b-sheets is very similar in structure to Src-homology 3 (SH3) domains, which are often associated with protein-protein interactions (15) (see [SH2, SH3 Domains](#)). It is likely that the basic character of this domain that has solvent-exposed [lysine](#) and [arginine](#) residues is involved in long-range [electrostatic interactions](#) with the negatively charged DNA phosphate groups. A saddle-shaped groove formed between two of the b-strands may provide a channel for this interaction (16).

**Figure 2.** Structure and functional organization of the domains of (a) retroviral (HIV-1) IN and (b) bacteriophage I Int proteins. Numbering is of the amino acid residues.



All three domains form dimers when analyzed as isolated polypeptide chains. Most likely, interactions between the domains of IN play a central role in forming the requisite multimeric configuration (17). The precise oligomeric form adopted by the protein *in vivo* is still not known, although high-order assemblies have been implicated in forming stable, catalytically active nucleoprotein complexes *in vitro* [eg (18)]. Despite progress made in this respect via protein–DNA [cross-linking](#) studies (19), the overall three dimensional view of these nucleoprotein intermediates remains sketchy.

## 2. Bacteriophage Lambda Integrase

### 2.1. DNA Integration Mechanism

The chromosomal integration of the viral genomes of temperate bacteriophages, exemplified by the reaction catalyzed by lambda integrase (Int), is a conservative site-specific recombination mechanism in which the recombination partner (*att*) sites are reciprocally exchanged by a set of isoenergetic cleavage and rejoining reactions (20). Like retroviral integration (see previous discussion), each strand-exchange reaction takes place via two consecutive nucleophilic substitutions, except, in this case, (1) a covalent Int–DNA intermediate is formed with each of the four recombination strands, and (2) the strand exchange is precise, neither adding nor deleting nucleotides at the site of insertion (Fig. 1b). Integrative recombination takes place through two sequential sets of such strand exchanges. Each is begun by the nucleophilic attack of Int at equivalent positions on the donor (*attP*) and target (*attB*) recombination strands. The energy of the primary substitution reaction is conserved via a 3'-phosphotyrosine (Fig. 1b, Step 1), which, after strand migration, is subsequently resolved by the exposed 5'-OH end of the corresponding recombination partner (Fig. 1b, Step 2). This second substitution results in covalently rejoining the two strands, and, in doing so, liberates the unmodified Int protein. Excision of the integrated viral genome (between *attL* and *attR* sites generated by the original recombination) is also catalyzed by Int. Although this occurs by the same reaction chemistry, the stimulation of excision versus integration is accomplished with the assistance of other viral and host factors. Xis and Fis promote excision, whereas IHF is required for both.

### 2.2. The Int Protein

The bacteriophage  $\lambda$  integrase (Int) is a 40-kDa (356 amino acid residue) site-specific DNA-binding protein that has type I topoisomerase activity. It recognizes *att* regions that contain two distinct types of Int binding sites distributed over ~270 bp: (1) a central, imperfectly symmetrical region equating to the crossover site (~30 bp) that consists of two inverted “core” DNA-binding sites, separated by the 7 bp that delineate the overlap between the boundaries of the points of recombination, and (2) more distant flanking motifs, termed “arm” DNA-binding sites. In conjunction with other accessory protein-binding sites, arm sites are control points through which the higher order synaptic complex may adopt the appropriate **tertiary conformation** for the correct positioning and hence, efficient recognition of the core-binding sites.

Several structure-function studies have pinpointed the primary DNA-binding properties of  $\lambda$  Int to within discrete domains of the protein (Fig. 2b). The N-terminus (residues 1 to 169) is required for site-specific binding to the core site but is devoid of catalytic activity, whereas the first 64 residues are sufficient to bind tightly to arm-type DNA substrates. The central domain (residues 65 to 170) recognizes core-type sites, albeit with a low binding affinity. The C-terminal domain (residues 170 to 356) contains the minimal catalytic site required for polynucleotidyl transfer. Sequence alignments show that several residues are highly conserved with other site-specific recombinases (21), particularly in the position of the catalytic [tyrosine](#) residue, Tyr342, and three essential active-site residues (Arg212, His308, and Arg311, or R-H-R). Its overall topology is that of a mixed a-b structure comprising a bundle of seven a-helices cradled by seven b-strands (organized into one b-sheet at the foot of the structure and two flanking [beta-turns](#)) (22). The catalytic R-H-R triad clusters at the base of a shallow groove, proposed as the DNA-binding surface, that runs along one face of

the protein. These residues are separated by ~20 Å from the catalytic tyrosine residue, which is located on an exposed and highly flexible 17-residue loop. The speculative function of the loop may be to bend backward bringing the Tyr342 nucleophile into the vicinity of the clustered R-H-R residues ([cis-acting](#)) or alternatively to extend outward and reach into the catalytic pocket of a paired Int protomer ([trans-acting](#)).

### Bibliography

1. R. Craigie (1992) *Trends Genet.* **8**, 187–190.
2. J. Kirchner, C. M. Connolly, and S. B. Sandmeyer (1995) *Science* **267**, 1488–1491.
3. G. V. Kalpana, S. Marmon, W. Wang, G. R. Crabtree, and S. P. Goff (1994) *Science* **266**, 2002–2006.
4. S. A. Chow, K. A. Vincent, V. Ellison, and P. O. Brown (1992) *Science* **255**, 723–726.
5. C. M. Farnet and F. D. Bushman (1997) *Cell* **88**, 483–492.
6. A. Aiyar, P. Hindmarsh, A. M. Skalka, and J. Leis (1996) *J. Virol.* **70**, 3571–3580.
7. S-Q. Wei, K. Mizuuchi, and R. Craigie (1997) *EMBO J.* **16**, 7511–7520.
8. C. Vink and R. H. A. Plasterk (1993) *Trends Genet.* **9**, 433–437.
9. R. Zheng, T. M. Jenkins, and R. Craigie (1996). *Proc. Natl. Acad. Sci. USA* **93**, 13659–13664.
10. M. Cai, R. Zheng, M. Caffrey, R. Craigie, G. M. Clore, and A. M. Gronenborn (1997) *Nat. Struct. Biol.* **4**, 567–577.
11. S. P. Lee and M. K. Han (1996) *Biochemistry*, **35**, 3837–3844.
12. O. Fayet, P. Ramond, P. Polard, M. F. Prère, and M. Chandler (1990) *Mol. Microbiol.* **4**, 1771–1777.
13. F. Dyda, A. B. Hickman, T. M. Jenkins, A. Engleman, R. Craigie, and D. R. Davies (1994) *Science* **266**, 1981–1986.
14. S. P. Moore, L. A. Rinckel, and D. J. Garfinkel (1998) *Mol. Cell. Biol.* **18**, 1105–1114.
15. A. P. Eijkelenboom, R. A. Lutzke, R. Boelens, R. H. Plasterk, R. Kaptein, and K. Hard (1995) *Nat. Struct. Biol.* **2**, 807–810.
16. P. J. Lodi et al. (1995) *Biochemistry* **34**, 9826–9833.
17. M. D. Andrade and A. M. Skalka (1996) *J. Biol. Chem.* **271**, 19633–19636.
18. I. K. Pemberton, M. Buckle, and H. Buc (1996) *J. Biol. Chem.* **271**, 1498–1506.
19. T. S. Heuer and P. O. Brown (1998) *Biochemistry* **37**, 6667–6678.
20. W. M. Stark, M. R. Boocock, and D. J. Sherratt (1992) *Trends Genet.* **8**, 432–439.
21. S. E. Nunes-Düby, H. Y. Kwon, R. S. Tirumalai, T. Ellenberger, and A. Landy (1998) *Nucleic Acid. Res.* **26**, 391–406.
22. H. J. Kwon, R. Tirumalai, A. Landy, and T. Ellenberger (1997) *Science* **276**, 126–131.

### Suggestions for Further Reading

23. A. Landy (1989) Dynamic, structural, and regulatory aspects of  $\lambda$  site-specific recombination. *Annu. Rev. Biochem.* **58**, 913–949.
24. K. Mizuuchi (1997) Polynucleotidyl transfer reactions in site-specific DNA recombination, *Genes to Cells* **2**, 1–12.
25. P. Rice, R. Craigie, and D. R. Davis (1996) Retroviral integrases and their cousins, *Curr. Opinion Struct. Biol.* **6**, 76–83.

### Integrative Suppression

[Initiation Of DNA Replication](#) of a [replicon](#) depends on an interaction between an initiator **protein** and a [replication origin](#) sequence specific for each replicon. In *Escherichia coli*, the DnaA protein binds to the origin and melts the duplex DNA strands to recruit the machinery for DNA synthesis, and it is the sole essential protein required specifically for the initiation of the chromosomal replication. **Conditional lethal** *dnaA* mutants defective in the initiation of replication under [nonpermissive conditions](#) can be rescued if the initiator gene and the replication origin of another replicon are integrated somewhere in the host chromosome. Such cells can survive by initiating chromosome replication from the origin of the integrated replicon, instead of the native origin. This phenomenon is called integrative suppression. In contrast, integration of the other replicon cannot relieve mutations in genes that are required for the replication of the integrated replicon.

The integrative suppression of *E. coli dnaA* mutations was observed using **plasmid** replicons, such as **F factor** and R1, and phage replicons, such as P2 (1). Furthermore, deletion mutants of the *dnaA* gene or the *oriC* sequence could be obtained. Similar integrative suppression was reported recently for *Bacillus subtilis* (2). The initiation site and frequency of chromosome replication could be changed artificially using integrative [suppression mutations](#), and the regulation of cell cycle progression in bacteria has been studied based on the **phenotypic** changes in the mutants (2, 3).

#### Bibliography

1. W. Messer and C. Weigel (1996) in *Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 1579–1601.
2. A. K. M. Hassan et al. (1997) *J. Bacteriol.* **179**, 2494–2502.
3. K. Nordstrom, R. Bernander, and S. Dasgupta (1991) *Mol. Microbiol.* **5**, 769–774.

## Integrins

Integrins are a large family of related, transmembrane, cell surface **receptor** proteins that bind [extracellular matrix](#) (ECM) and [cell adhesion molecules](#) (CAMs). Almost all cell types express some type of integrin, as these receptors are often responsible for positioning cells within their specific environment by associating them with ECMs or CAMs on adjacent cells. Integrins play crucial roles in cell–matrix and cell–cell interactions important in the [development](#) and function of organisms ranging from plants to *Drosophila* to humans. Many cells require an attachment to the extracellular matrix to survive and grow. Cells also make use of integrins during migratory events that are important in neural development, inflammation, and metastasis.

### 1. Integrin structure

Integrins are noncovalently associated heterodimers consisting of an alpha and a beta subunit, both of which span the [membrane](#) lipid bilayer. The ligand binding site outside the cell requires contributions from both subunits. Both subunits have a single membrane-spanning domain and short cytoplasmic domains. The latter interact with the [actin](#) cytoskeleton and transmit signals to the interior of the cell. When ligand is bound, integrins signal across the membrane by altering actin polymerization and activating a [tyrosine kinase](#). The name *integrin* derives from the way that these receptors “integrate” fibrils of ECM on the outside of the cell in direct register with actin filaments underneath the plasma membrane.

Alpha subunits range in size from approximately 200 to 120 kDa and are 25% to 45% identical in amino acid sequence. The alpha subunits are characterized by a series of 3 to 4 repeated [EF-hand motifs](#) in the extracellular domain that are important for binding divalent cations.  $Mn^{2+}$  and  $Mg^{2+}$  ions are required for integrin binding, and disruption of integrin activity with [EDTA](#) is often used to dislodge tissue culture cells from the substratum. These divalent cations participate in the binding of ligand. Some integrin alpha subunits contain an extra inserted **domain** in the N-terminus, called an *I domain*, which is thought to expand the **ligand binding** ability of integrins. Some alpha subunits are proteolytically cleaved in the extracellular domain near the membrane, and the resulting heavy and light chains are linked by [disulfide bonds](#). Just cytoplasmic to the transmembrane domain, all alpha subunits have a highly conserved GFFKR motif (-Gly-Phe-Phe-Lys-Arg-) whose function is unknown. Other than this sequence, the cytoplasmic domains of the alpha subunits are generally not highly conserved. Differential [RNA splicing](#) is known to occur for some alpha subunit [messenger RNA](#) (eg, alpha-6, alpha-7), and the [alternative splicing](#) variants have specific patterns of expression and may have altered binding properties as well.

The beta subunits are about 95 to 120 kDa in size, with the notable exception of the 200-kDa beta-4 subunit. The amino acid sequences, again about 40% to 50% identical between subunits, include four conserved [cysteine](#)-rich domains in the extracellular region. The short cytoplasmic domain (40 to 50 amino acid residues) is generally more conserved in the beta subunits, except for beta-4, which has an intracellular domain of 1000 residues containing fibronectin type III domains. Interestingly, the beta-4 integrin interacts with [intermediate filaments](#) rather than with the actin [cytoskeleton](#) as other integrins do. Splice variants of the beta subunit mRNA have been identified which, in the case of the beta-1 subunit, is thought to result in different signaling properties, with one such variant displaying a dominant negative activity. The beta-1 subunit cytoplasmic domain also contains NPXY motifs (-Asn-Pro-X-Tyr-), thought to be important in mediating [endocytosis](#) by [clathrin](#)-coated vesicles.

## 2. Integrin Signal Transduction

A conformational change is propagated across the lipid bilayer to alter proteins associated with the cytoplasmic domains. The beta-1 subunit, and probably also other beta subunits, interacts with the actin cytoskeleton via associations with the proteins talin and alpha-**actinin**. These proteins colocalize at structures called *focal adhesions* (adhesion plaque, [focal contact](#)), which are areas closely opposed to the substratum in cultured cells. Protein tyrosine kinases also associate with focal adhesions, and *focal adhesion kinase* associates with the cytoplasmic domain of the beta-1 subunit and is activated on integrin ligation. This kinase activity is thought to alter the dynamics of actin and to signal to the nucleus to regulate **gene expression**. Integrins also bring about cellular alkalization and can lead to increases in the cytoplasmic calcium ion concentration and activation of protein kinase C (see [Calcium Signaling](#)). A variety of [enzymes](#) are known to increase their activity after integrin stimulation, including phospholipase A2, the GTPase **Rho**, and [phosphatidyl inositol](#) phosphate kinase (PIP 5-kinase), and enzymes associated with the Ras/MAP kinase pathway.

[Growth Factors](#), such as [platelet-derived growth factor](#), are known to synergize with ECM integrin ligands to stimulate cellular replication, and this results from a convergence of signaling pathways. Other cell surface receptors may interact with integrins and activate integrin signaling.

Integrin activity is regulated post-translationally in a cell-specific manner. For example, [transfection](#) of the same alpha-2 [complementary DNA](#) into different cell types that express the beta-1 subunit can result in either an inactive **receptor**, a [collagen](#) receptor, or a dual-function collagen and laminin receptor. The activity appears to be regulated by altering the integrin conformation, but the factors controlling this process are not understood. Activating antibodies have been generated that bind to the extracellular domain of integrins and activate increased ligand binding. Intracellular changes can also activate integrins through an inside-out signaling mechanism. For example, in T lymphocytes, the alpha-1/beta-2 integrin is activated to bind I-CAM only after the [T cell receptor](#) binds [antigen](#). In

this way, the T cell becomes more adhesive to endothelial cells lining the blood vessels in the area surrounding the antigen. This type of affinity modulation is also seen in platelets, where ADP or [thrombin](#) can induce high-affinity [fibrinogen](#) binding in the alpha-IIb/beta-3 integrin, which is important for blood clotting. Tyrosine phosphorylation of the beta-1 subunit by pp60src inhibits ligand binding.

Regulation of integrin function is crucial to many processes in [development](#). For example, decreasing integrin expression in skin epithelial cells leads to their movement away from the basal lamina and their terminal differentiation. Inactivation of integrins in retinal neurons is correlated with the timing of axon pathfinding, and this may lead to a turn by the growth cone away from ECM ligands and towards CAM ligand.

### 3. The Integrin Family

At least 16 alpha subunits and 8 beta subunits, constituting at least 21 heterodimers, have been discovered to date. Different alpha-beta combinations give rise to receptors with different ligand-binding specificity, although not all possible alpha-beta pairings have been detected (Table 1). Receptors are often multifunctional, with abilities to bind more than one ligand. In addition, there is an apparent redundancy in function among the heterodimers. For example, at least seven integrin laminin receptors have been identified. It has become clear that some integrins bind to different sites on the same ligand and bind different ligand isoforms with varying avidity. For example, integrin alpha-5/beta-1 binds to the RGD site in fibronectin, but the alpha-4/beta-1 receptor binds to an LDV motif in an alternately spliced fibronectin. Some beta subunits are able to pair with multiple alpha subunits and, likewise, some alpha subunits are capable of forming heterodimers with more than one beta subunit.

**Table 1. Integrin Heterodimers and Some of Their Ligands**

|            | <b>b1</b>                                  | <b>b2</b> | <b>b3</b>     | <b>b4</b>                       | <b>b5</b> | <b>b6</b>        | <b>b7</b>                          | <b>b8</b> |
|------------|--|-----------|---------------|---------------------------------|-----------|------------------|------------------------------------|-----------|
| <b>a1</b>  | C, <sup>a</sup> L <sub>1</sub>             |           |               |                                 |           |                  |                                    |           |
| <b>a2</b>  | C, L <sub>1</sub> , a3,<br>Ch              |           |               |                                 |           |                  |                                    |           |
| <b>a3</b>  | L <sub>5</sub> , C, F, E,<br>a2, 3, T      |           |               |                                 |           |                  |                                    |           |
| <b>a4</b>  | VCAM,<br>F <sub>LDV</sub> , a4, T,<br>O, W |           |               |                                 |           |                  | VCAM, F <sub>LDV</sub> ,<br>MadCAM |           |
| <b>a5</b>  | F <sub>RGD</sub> , T                       |           |               |                                 |           |                  |                                    |           |
| <b>a6</b>  | L <sub>1</sub> , L <sub>5</sub>            |           |               | L <sub>1</sub> , L <sub>5</sub> |           |                  |                                    |           |
| <b>a7</b>  | L <sub>1</sub>                             |           |               |                                 |           |                  |                                    |           |
| <b>a8</b>  | F <sub>RGD</sub> , Tn,<br>O                |           |               |                                 |           |                  |                                    |           |
| <b>a9</b>  | C, L <sub>1</sub>                          |           |               |                                 |           |                  |                                    |           |
| <b>a10</b> | C  |           |               |                                 |           |                  |                                    |           |
| <b>aV</b>  | V, F <sub>RGD</sub>                        |           | Fb, F, V, OB, |                                 | V         | F <sub>RGD</sub> |                                    | V         |

|             |                   |                  |           |
|-------------|-------------------|------------------|-----------|
|             |                   | <b>W, T</b>      |           |
| <b>aL</b>   | <b>ICAM 1,2,3</b> |                  |           |
| <b>aM</b>   | <b>C3bi,FX,Fb</b> |                  |           |
| <b>aX</b>   | <b>FB, C3bi</b>   |                  |           |
| <b>aD</b>   | <b>ICAM3</b>      |                  |           |
| <b>aIIB</b> |                   | <b>Fb,F,V W,</b> |           |
|             |                   | <b>ICAM</b>      |           |
| <b>aE</b>   |                   |                  | <b>aE</b> |

<sup>a</sup> Key: C, collagens; F, fibronectin (RGD site or LDV site indicated); L, laminin (isoform number indicated); V, vitronectin; VCAM, vascular cell adhesion molecule 1; T, thrombospondin; Tn, tenascin/cytotactin; ICAM, intercellular adhesion molecule 1,2,3; MadCAM, mucosa addressin cell adhesion molecule. Fb, fibrinogen; W, von Willebrand factor; FX, factor X; C3bi, inactivated component of complement 3B; O, osteopontin; B, bone sialoprotein; Ch, chondroadherin.

The beta-1 family of integrins includes receptors for [fibronectin](#), [laminins](#), collagens, vitronectin, thrombospondin, and a variety of other matrix ligands. These integrins mediate many types of interactions between cells and the ECM. Beta-1 integrins are involved in important basic developmental cellular movements, such as **gastrulation**. Beta-1 laminin receptors are important in the development of neurons and muscle cells, from migration of neuroblasts and myoblasts to axon and dendrite outgrowth to synapse formation. The alpha-2/beta-1 collagen/laminin receptor mediates the contraction of collagen gels by resident cells during the remodeling of the interstitial ECM, and the alpha-5/beta-1 receptor has been implicated in organizing fibronectin networks. The alpha-4/beta-1 receptor, which binds both fibronectin and the vascular **cell adhesion molecule-1** (VCAM-1), allows white blood cells to bind to inflamed endothelial cells (which express VCAM-1) and to exit out of the blood vessel to mediate an immune response. Some beta-1 integrins are also thought to interact with each other or to associate with other integrins on adjacent cells.

The beta-2 integrins are expressed by white blood cells and mediate cell–cell interactions important to the immune response by interacting with [immunoglobulin](#) super family molecules (ICAMs) and other ligands. Once neutrophils are activated, the alpha-L/beta-2 on their surface adheres to ICAMs on the endothelial cell surface, working in concert with the alpha-4/beta-1 receptor to allow the cell to exit the circulation. Beta-3 integrins include vitronectin receptors on bone and blood vessel endothelial cells (alpha-v/beta-3), and the alpha-IIA/beta-3 fibrinogen receptor important in platelet aggregation. Angiogenesis can be blocked by inhibitors of the alpha-v/beta-3 receptor on endothelial cells, which causes them to undergo [apoptosis](#), and this same receptor is also involved in phagocytotic processes.

Other beta-subunit families, although smaller and less well characterized, may prove to be equally important. For example the alpha-6/beta-4 integrin is part of the [hemidesmosomes](#) in epithelial cells that attach cells to anchoring filaments made of laminin-5, and the beta-8 integrins are abundantly expressed on neurons.

The importance of integrins is underscored by genetic studies in a variety of species. In the fruit fly *Drosophila melanogaster*, integrin homologues were first identified as “position-specific antigens” with interesting patterns of expression during development. Loss of function of one beta subunit results in the myospheroid **phenotype**, where muscle cells cannot adhere to the matrix and assume a nonfunctional rounded shape. Loss of the alpha-PS3 gene results in numerous defects in tissue morphogenesis and cellular movements. Another alpha subunit is required for short-term memory in flies, suggesting a new role for integrins at the synapse. [Nematode](#) mutants show defects in cellular migrations, and mouse **knockout strains** in individual subunits have thus far proved lethal. For

example, knockout of the alpha-4 subunit results in lethality at embryonic day 10, in part because of defective heart formation. Cellular interactions during heart development may require interactions between alpha-4/beta-1 and VCAM-1. Beta-1 subunit mutants fail to implant in the uterus. In humans, nonfunctional beta-2 subunit results in leukocyte adhesion deficiency disease, which renders patients prone to bacterial infections. Loss of the beta-3 subunit results in a clotting disorder called Glanzmann's thrombasthenia, and mutations in the alpha-6/beta-4 heterodimer lead to a blistering skin disease known as epidermolysis bullosa.

#### Suggestions for Further Reading

R. O. Hynes and A. D. Lander (1992) Contact and adhesive specificities in the associations, migrations and targeting of cells and axons. *Cell* **68**, 303–322.

M. A. Schwartz, M. D. Schaller, and M. H. Ginsberg (1995) Integrins: emerging paradigms of signal transduction. *Ann. Rev. Cell Dev. Biol.* **11**, 549–599.

E. Ruoslahti (1996) RGD and other recognition sequences for integrins. *Ann. Rev. Cell Dev. Biol.* **12**, 697–715.

## Interallelic Complementation

### 1. What counts as allelic?

The use of complementation tests for defining the gene as a unit of function is explained under the headings [Complementation](#) and [Complementation tests](#). It was assumed as a general rule in those articles that complementation between two mutants occurs only when their mutations are in different **genes (cistrons)**. Complementation tests carried out, for example, on [auxotrophic](#) (nutritionally exacting) mutants of the budding **yeast** *Saccharomyces cerevisiae* or the filamentous **fungus** *Neurospora crassa*, indeed split the mutants cleanly into different complementation groups. Little or no complementation occurs within groups, and universal complementation occurs between them (see Fig. 3 of [Complementation Tests](#)). However, the patterns that emerge from complementation tests are not always clear-cut, and many examples can be cited of complementation between mutants that one would want to call allelic on other grounds.

An **allele** is a particular form of a gene, and so what counts as allelic complementation depends on how the gene is defined. As a working rule, most would agree to define the gene as a segment of DNA that is transcribed into a single molecule of RNA, although that definition needs qualification to accommodate various special systems of gene organization. It may also be considered reasonable to include within the definition those adjoining *cis*-acting, nontranscribed DNA sequences that regulate gene [transcription](#), though this extension of the definition may lead to difficulties in deciding how far the gene's territory extends and whether it is allowed to overlap with other genes.

In the great majority of the eukaryotic organisms studied, genes that have protein-encoding functions are virtually always units of [translation](#) and of transcription. Each gene specifies just one [polypeptide chain](#). In bacteria, however, and in the [nematode](#) *Caenorhabditis*, many genes are organized in [operons](#), which are units of transcription that contain two or several nonoverlapping sequences each of which is translated into a different polypeptide chain. Thus, the *Escherichia coli lac* operon [messenger RNA](#) encodes **b-galactosidase** (*LacZ*), b-galactoside permease (*LacY*), and galactoside **transacetylase** (*LacA*). Here the term gene refers to the single-protein units rather than the whole unit of transcription. Even though *LacZ* and *LacY* are within the same unit of transcription, they have sufficient independence of function to rank as separate genes, and the normal phenotype of a  $Z^- Y^-$  partial diploid (see [Complementation Tests](#)) poses no problem and is not counted as allelic



complementation.

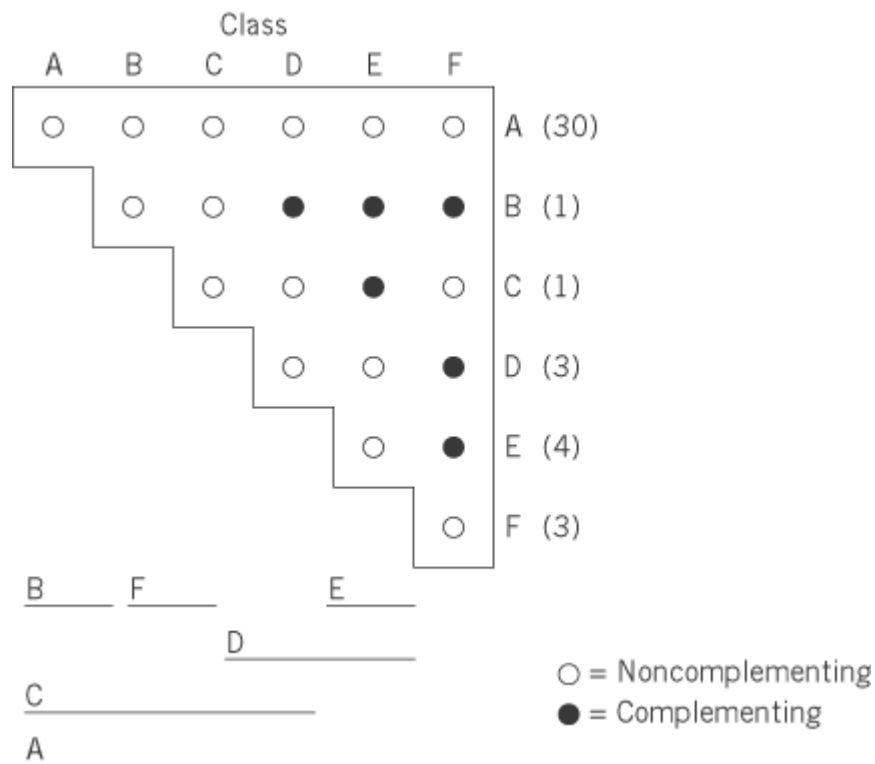
Attention was first focused on allelic complementation as a result of intensive studies of auxotrophic (nutritionally exacting) mutants of microorganisms, first in the fungus *Neurospora crassa* and then in budding yeast and bacteria. The *Neurospora* mutants gave the main early encouragement to the view that each gene is responsible for the specificity (at first rather vaguely defined, but later taken to mean amino acid sequence) of a single [enzyme](#). Each enzyme was expected to catalyze a single reaction, and indeed, nearly all auxotrophic mutants were defective in carrying out just one biosynthetic step. The finding, in *Escherichia coli*, that **tryptophan synthetase** deficiency results from mutation of either of two clearly distinct genes was clarified by the demonstration that this enzyme consists of two different polypeptide chains, one for each gene. Subsequently it was recognized that “one gene—one polypeptide” was a better working hypothesis than “one gene—one enzyme.”

## 2. Complementation Maps

Because genes were expected to encode single polypeptide chains, there was no obvious possibility of allelic complementation. That would require recombination between polypeptide chains, which did not occur. However, as more mutants were screened, it was commonly found that, although most mutants defective in some single biosynthetic step were indeed noncomplementing in all pairwise combinations, testing of larger numbers revealed a few that showed complementation with some of the other members of the series, though still noncomplementing with most of them.

The complementation relationships within a series of mutants of the same enzyme were summarized in the form of complementation maps, in which mutants were represented as linear segments, nonoverlapping or overlapping depending on whether or not they complemented one another. [Figure 1](#) shows a relatively simple example. Each group of mutants had some unity and justified their being called allelic because many or most of them were noncomplementing with all of the others. The overlaps in the complementation maps were crucial because they made it difficult to consider that the series of mutants falls in more than one gene, perhaps encoding different polypeptide components of the same enzyme. In principle, the overlaps could represent deletions of DNA that extended between genes, but deletions as a general explanation could be ruled out on two grounds. First, the mutational sites of “overlapping” mutations often recombined in **meiosis**. Secondly, they could, in many cases, undergo reverse mutation to wild type.

**Figure 1.** Complementation matrix and complementation map for a set of 42 *Neurospora crassa arg-1* mutants that are deficient in the enzyme argininosuccinate synthetase. The number of mutants in each of the complementation classes A–F is shown in parenthesis ([1](#)).



Up to a point, it was possible to place all of the segments of most complementation maps in a simple line (eg, Fig. 1), encouraging the view that perhaps the map had some relationship to the linear structure of the gene. But with the inclusion of more mutants, it usually became impossible to accommodate all of them without making the map two-dimensional with closed loops and tails. And where a complementation map and a fine-structure gene map were both available, usually a simple and consistent relationship did not exist between the two (3).

No general theory explains the many and various forms of complementation maps, but at least some special explanations have emerged from studies of particular cases. Much depends on the whether the enzyme (or other protein) gene product has one function or several.

### 3. Dimeric/Oligomeric Enzymes with Single Functions—Conformational Correction

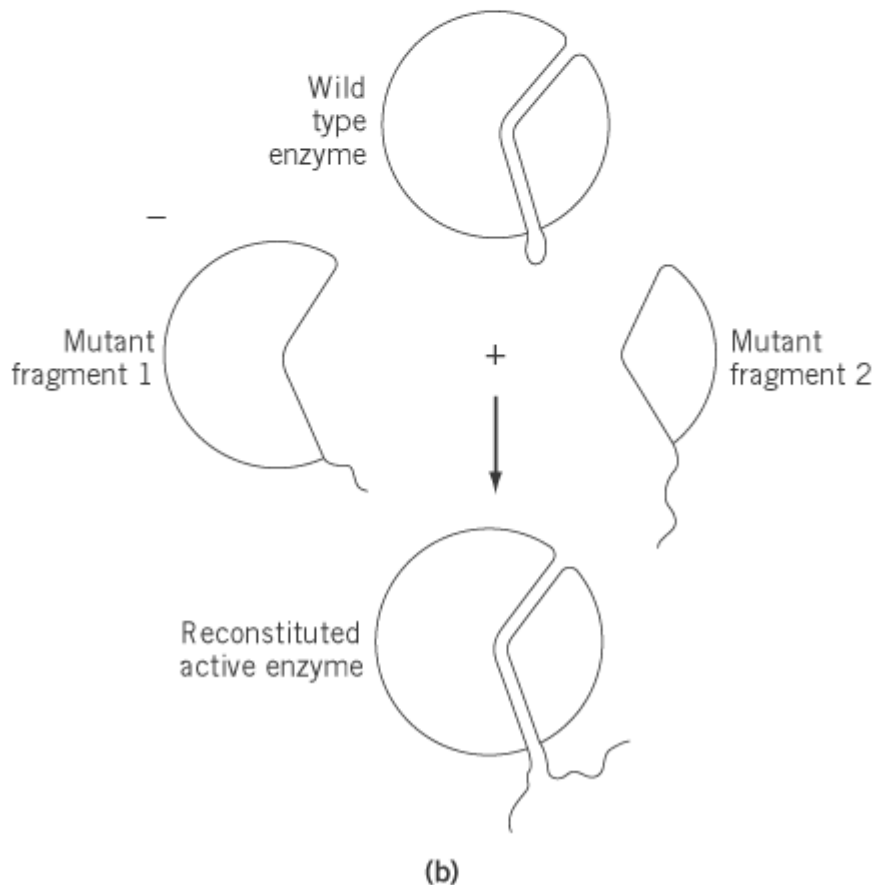
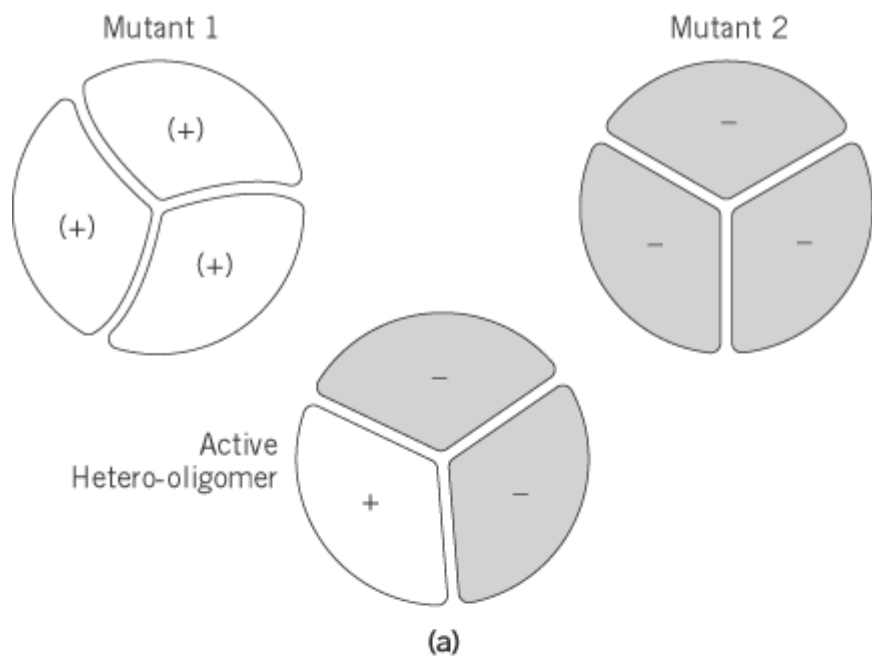
In the first *Neurospora* examples, it was clear that the mutations affect single enzymes that catalyze single reactions. Thus the *ad-4* mutants, unable to make adenine, are defective in adenylosuccinase (4), and the *am* mutants that require  $\alpha$ -amino nitrogen lack NADP-specific glutamate dehydrogenase (GDH) (5). Within each of these mutant series, some pairs show complementation, and in both cases it was shown that, under certain conditions, some degree of enzymatic activity could be obtained from mixed protein preparations of complementing mutants. The mutually complementing *am* mutants produce inactive, but physically nearly normal, forms of GDH protein, and active GDH is readily obtained by freezing and thawing mixtures of purified complementing mutant proteins in the presence of 0.1 M sodium chloride. This is a procedure that brings about dissociation and reassociation of protein subunits in other systems (6). NADP-specific GDH is a hexamer (six identical polypeptide subunits), and the enzymically active complementation products are mixed hexamers, though the complex mixture of hybrid proteins (all combinations from 5:1 to 1:5) was resolved only partly into its components (7).

The hybrid protein hypothesis for allelic complementation had already gained support from experiments by Schlesinger and Levinthal (8) on mutants of *Escherichia coli* deficient in [alkaline phosphatase](#), a dimeric enzyme. Here a mixture of two individually almost inactive mutant forms of

the enzyme dissociated to monomers by acidification and then allowed to reassociate was to form highly active mixed dimers (8).

Crick and Orgel (9) proposed that individually inactive monomers of dimeric enzymes might remedy each other's defects through mutual correction of their folded conformation. This is especially likely when the enzyme involved is [allosteric](#) and switches between active and inactive forms by ligand binding. One likely reason for mutational inactivation of an allosteric enzyme is excessive stabilization in the inactive form. Allosteric proteins are generally dimeric or [oligomeric](#) and, according to the simplest version of allosteric theory, their constituent monomers shift their conformation in a concerted, all-or-none, fashion. If so, one would expect that the conformational preference of the majority or more strongly stabilized component would prevail in a mixed oligomer. The idea is explained diagrammatically in Fig. 2a. It receives some support from studies of the *Neurospora am* mutants. *Neurospora* NADP-GDH is indeed an allosteric enzyme and, in most of the complementing pairs of mutants, one partner produces a GDH variety that is potentially active, but switches to the active conformation only under unphysiologically extreme conditions. In the best investigated examples, the other partner is a form of the protein that is unconditionally inactive because it has lost an essential [side-chain](#) in a substrate-binding site (see [Enzymes](#)), but is conformationally nearly normal. The potential activity of the first component was realized upon **hybridization** with the second (6). Another *am* mutant (10) that complements either of the first two kinds is best interpreted as producing a GDH variety that cannot form stable hexamers by itself.

**Figure 2.** Two modes of complementary interaction between mutant derivatives of the same protein. (a) Formation of mixed oligomers (shown here as trimers) that have one mutant monomer potentially active (+) but stabilized in an inactive conformation and the other always inactive but conformationally normal. In the mixed trimer (hetero-oligomer), the activity of the (+) component is realized by conformational correction. (b) Formation of active enzyme by association of two independently folded fragments. The shapes depicted in this figure are imaginary.



These examples of allelic complementation still leave some possibility of rescuing complementation as a criterion for allelism. Instead of insisting that alleles should never complement each other, one can adopt a softer criterion, namely that they should never complement to form a truly wild-type protein product. Then, the *cis/trans* criterion for allelism survives in modified form, and the crucial *cis-trans* difference is no longer between function and no function but between full function and (perhaps marginally) subnormal function. In the case of the *Neurospora am* mutants, it was fairly

easy to demonstrate that their complementation products are subwild and have relatively low activity per unit protein and reduced **thermostability**. But this amended criterion for allelism is not always easy to apply. The best indication that a series of mutants belongs to a single gene is that, notwithstanding complementation between some pairs, a substantial proportion fails to complement any of the others.

On the conformational correction hypothesis, complementation maps that relate to single enzymatic functions must be presumed to represent interactions within the three-dimensional structures of proteins. Though most are too complex to interpret, a limited number of nonoverlapping segments can be explained in terms of global properties of mutant proteins, as in the case of the *Neurospora am* mutants mentioned previously. In any case, conformational correction applies only when the protein is a dimer or oligomer.

#### 4. Fragment Complementation with a Monomeric Enzyme

At one time allelic complementation was seen as diagnostic of a dimeric or oligomeric gene product. What was overlooked, however, was the possibility of piecing together a functional monomeric protein from fragments. There is good evidence from several sources that this happens.

The best known example (11) involves *Escherichia coli* **b-galactosidase**, the *LacZ* product. This enzyme consists of a single, long (1025-residue) polypeptide chain. Mutants that have deletions of the N-terminal (“upstream”) end of the chain complement virtually any mutant that has a defect in a substantial section at the C-terminal or downstream end (the w-region) provided that the deletion does not extend into that region. The upstream deletions were called w-donors and the downstream deletions complementing mutants w-acceptors. The latter class includes chain-terminating mutants in which the w region is hardly present at all. Thus an interaction between complementary fragments of a single polypeptide chain, called in this case w-complementation, result in reconstituting active enzyme.

The general explanation for complementary interaction of polypeptide fragments is that each of the two polypeptide sequences, whether continuous as in the wild type or separated as in w-complementation, fold up independently to form three-dimensional “**domains**”—protein globules that fit closely and stably together like two pieces of a three-dimensional jigsaw puzzle. Their mutual affinity enables the initially separated domains to find one another when synthesized in the same cell or, indeed, when mixed *in vitro*. w-Complementation creates a fully active and stable b-galactosidase, and the lack of a covalent connection between the two fragments does not matter. The overlap between them increases the molecular weight of the enzyme somewhat but does not impede its function.

Now, complementation between *LacZ* sequences is used routinely in **plasmid** technology (12). An *E. coli lacZ* deletion mutant, M15, lacks just 30 amino acid residues at the N-terminus of the b-galactosidase chain, and this nearly complete polypeptide is complemented by many different short N-terminal fragments that evidently can repair the N-terminal (a) domain. A number of different cloning plasmids have been constructed so as to encode one of these a fragments, and, if introduced into M15, signal their presence by producing active a-complemented b-galactosidase and the appearance of blue colonies on a growth medium containing a substance that is cleaved by the enzyme to give the color, **X-gal**.

It is still not clear how important interactions between polypeptide fragments, as opposed to formation of mixed oligomers, will prove to be an explanation of allelic complementation. For oligomeric proteins, the two possibilities are not mutually exclusive. Their occurrence together within the same allelic series could be responsible for the complexity of some complementation maps.

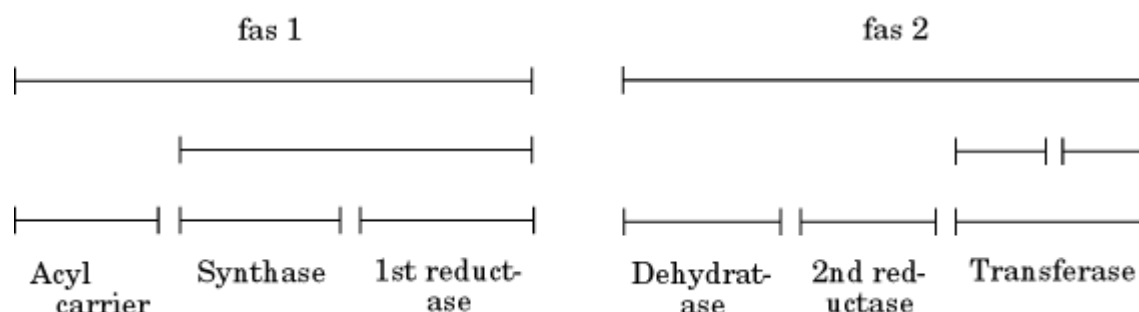
#### 5. Enzymes with Multiple Functions

Another explanation for allelic complementation, which applies more commonly than at one time thought, is that the protein product of the gene has several enzymatic functions, usually involved with different steps in the same biosynthetic pathway. An amino acid replacement that eliminates one enzymatic activity does not necessarily have any critical effect on the others. In terms of protein structure, this means presumably that the polypeptide chain is folded into a number of domains that function largely independently. In this kind of situation, different mutant alleles can in principle complement one another without mutual conformational correction or even the formation of mixed oligomers, though these interactions may also occur.

The genetic control of multifunctional enzymes has been most thoroughly investigated in filamentous fungi and budding yeast because extensive collections of relevant mutants have been available in these organisms. Here we briefly review three examples, two from the yeast *Saccharomyces cerevisiae* and one from *Neurospora crassa*.

Schweizer and his colleagues (13, 14) collected a large number of yeast mutants that had lost the ability to synthesize [fatty acids](#). The pathway of synthesis, starting with acetic and malonic acids and leading to long-chain fatty acids, could be divided biochemically into six steps, corresponding to separately assayable enzymatic activities. The whole pathway was catalyzed by a protein complex composed of just two long polypeptide chains. The mutants were assigned by genetic mapping and complementation tests to two genes, *fas1* and *fas2*, unlinked to each other. Every *fas1* mutant complemented every *fas2* mutant. Extensive complementation also occurred within *fas1* and within *fas2*. The complementation maps are shown in Fig. 3. Three complementation groups in each map are complementary in all combinations, and these correspond to losses of different single enzymatic activities. In both *fas1* and *fas2*, there are also mutants that failed to complement any of the others and some that have more restricted overlaps in the map. The comprehensively noncomplementing mutants are mostly chain terminators or, less often, frameshifts, which failed to synthesize a complete polypeptide chain (see [Frameshift Mutation](#)) and, as a result have lost all the enzyme activities normally associated with it. Although the different domains of the polypeptide chains could function independently so long as complete chains were produced, a shortened chain failed to form any of the active domains.

**Figure 3.** Complementation relationships of *Saccharomyces cerevisiae* *fas* mutants deficient in one or more activities of the fatty acid synthesizing enzyme complex. By complementation testing, they split into two distinct groups, *fas1* and *fas2*, that correspond to two unlinked genes and two different polypeptide chains. Within each group are three mutually complementing classes that correspond to losses of three different enzymatic activities of a single polypeptide chain and numerous mutants (mostly chain-terminating) that fail to complement all of the others (13, 14).

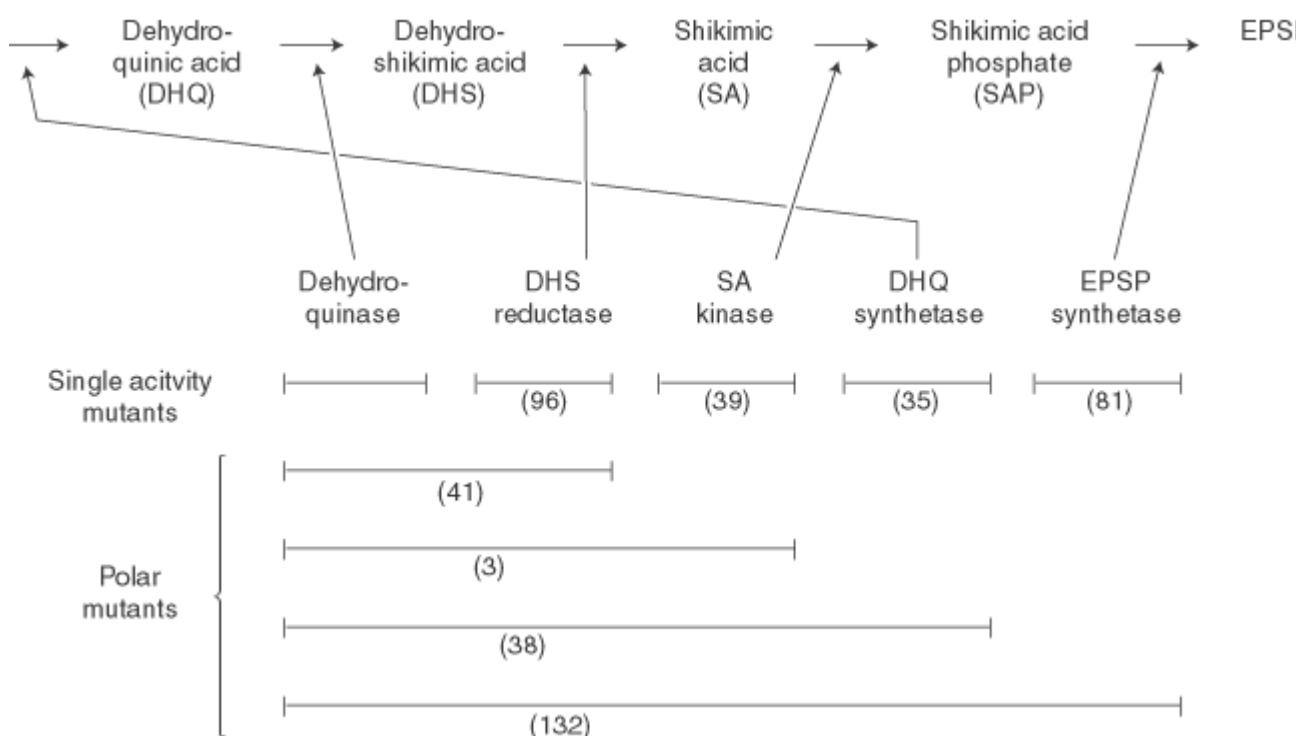


The *Neurospora* example, probably the most complex and extensively analyzed of any, involves the pathway for synthesizing aromatic compounds. Mutants blocked in this pathway have to be supplied with a complete mixture of essential aromatic metabolites ([phenylalanine](#), [tyrosine](#), [tryptophan](#), and *p*-aminobenzoic acid) because they cannot make the common precursor of the benzene ring,

chorismic acid. The pathway segment relevant here comprises the five steps leading to enolpyruvylshikimic acid phosphate, the immediate precursor of chorismate. The five distinct enzymatic activities were at first thought to be due to separate proteins encoded by a close cluster of genes because a number of closely linked and mutually complementing mutants were found that lack single enzymatic activities. However, doubt was thrown on the separate gene interpretation by the numerous mutants that failed to complement two or more of the single-activity groups.

The complementation map is shown in Fig. 4. The map shows a polarized pattern, and all the overlaps start at the left (DHS reductase) end of the map and extend for different distances to the right. The situation became much clearer following the biochemical demonstration that the entire complex of reactions is catalyzed by a dimeric multifunctional enzyme that has just one kind of polypeptide chain. As with yeast *fas1/fas2*, the mutants represented as overlaps on the complementation map are largely chain terminators. But the *aro* protein differs from the *fas1/fas2* complex in the polarized pattern of its complementation map. The obvious explanation for this is that the gene, as represented in Fig. 4, is transcribed and translated from right to left, and prematurely terminated forms of the polypeptide chain retain the enzymatic activities of the domains upstream of the termination point. A remaining anomaly is that, to give coherence to the complementation map, the EPSP synthetase segment has to be placed at the N-terminus of the chain, upstream of DHQ synthetase, whereas in the actual polypeptide chain the two segments are the other way round. The fact that the enzymatic complex is a dimer opens the door to further complications. There is the possibility of conformational interactions between monomers, and the unanswered question of how prematurely terminated monomers are accommodated in the dimeric structure.

**Figure 4.** Relationships among *Neurospora crassa* mutants that are deficient in one or more of five steps in the synthesis aromatic ring, depicted in the upper line (EPSP, the precursor of chorismic acid, is enolpyruvyl shikimic acid phosphate) different enzymatic activities (second line) are properties of different domains of a single polypeptide chain. The complementation map below shows five classes of *aro* mutants, each of which is deficient in a single enzyme activity, complementing each other in all combinations, and “polar” mutants each of which lacks two or more activities. The num mutants found in each class in the first screen are shown in parentheses. Initially no singly deficient mutants were found DHQase class because this function is duplicated by a catabolic DHQase. The missing class was found after the latter en had been eliminated by mutation (15).



The second example of polarized complementation, this time from yeast, provides a link to similar situations in insects and mammals. Among several classes of **uridine** or **cytidine**-requiring (URA) mutants, the *URA2* group, a tightly linked group, that maps to a first approximation to the same chromosome locus (*URA2*), includes two mutually complementing subgroups, A and C, and other mutants (B) that complement neither (16). When the mutational sites were placed in order in a recombinational gene map, the A and C subgroups mapped in two close but nonoverlapping segments, whereas nearly all of the B mutations, which eliminate both A and C functions, mapped within the C segment. It was shown that the B overlapping class consists mostly of chain-terminating mutations, with some frameshifting insertions or deletions, which are expected to eliminate or make nonsense of “downstream” polypeptide sequences. Essentially the same situation has been found in two other fungi, *Neurospora* and the ink-cap mushroom *Coprinus*.

The enzymological explanation of all three fungal examples is the same. The protein gene product is multifunctional, and has two enzymatic activities, carbamoyl phosphate synthase (CPSase) and **aspartate trans carbamoylase** (ATCase), that catalyze the first two steps of the pyrimidine biosynthetic pathway. The yeast *URA2*-A mutants lacked ATCase, the C mutants lacked CPSase, the B mutants lacked both, and the same applied to the analogous mutant classes in the other two fungi. The polypeptide also has a third enzymatic activity, glutamine amidotransferase (GATase), an initially unrecognized component of CPSase activity. The single polypeptide chain is synthesized from the GATase/CPSase end, so that downstream chain-terminating mutations could eliminate ATCase and not CPSase, and upstream terminators could only eliminate both.

A similar situation applies in the pyrimidine biosynthetic pathway of *Drosophila*. The difference here is that the analogous gene product has an additional enzyme function, dihydro-orotase, that catalyzes a further step in the pathway, as well as GATase, CPSase and ATCase (17). The gene is called *rudimentary* (*r*) referring to the effect on the wings and bristles of the fly, which are sensitive indicators of pyrimidine limitation. The *r* series of mutants provided the first demonstration in *Drosophila* of allelic complementation comparable to the fungal examples (18). Later work showed that generally the mutually complementing *r* mutants are deficient in different enzymatic activities of the single-gene product.

Multifunctional enzymes encoded by single genes (sometimes called “cluster genes” to point out their superficial resemblance to gene clusters) are not uncommon in higher organisms. For example, the first steps of pyrimidine biosynthesis in the mouse are catalyzed by a single protein very like the *Drosophila r* gene product, and in *Drosophila* a single protein (encoded by the *Gart* gene) is responsible for catalyzing several steps of purine synthesis, not all of them sequential. If complementation between mutant alleles has not been a feature of these other systems, it is because the necessary collections of viable mutants have not been obtained.

Although a degree of complementation between different mutant derivatives of multifunctional enzymes is to be expected, it is not clear whether it always, or even ever, results in maximally efficient enzymatic function. The catalysis of several sequential reactions by a single protein suggests that the organism's metabolism derives some benefit from the proximity of the different catalytic protein domains. In the fungal examples just reviewed, the separated enzymatic activities function together well enough to support normal growth under laboratory conditions, but this is perhaps not a very stringent test.

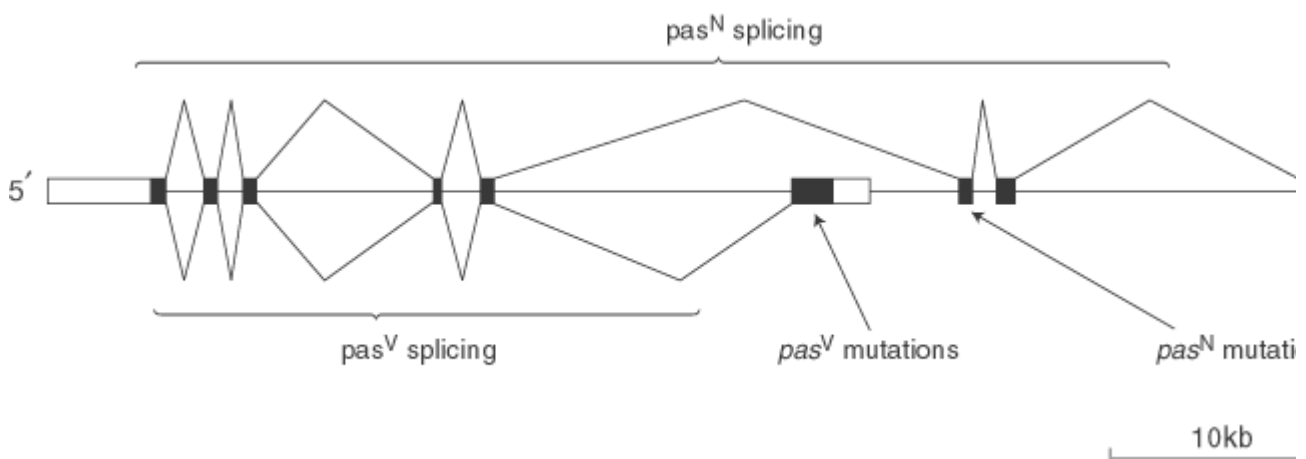
## 6. Single Genes that Encode More Than One Protein

Investigation of the exon-intron structures of genes has brought the realization that a single gene, that is, with a single primary RNA transcript, can encode two or more different proteins that have different and complementary functions as a result of different modes of splicing of the pre-mRNA (see [Alternative Splicing](#)). The different protein products always, at least in the known examples, have much of their amino acid sequence in common, but two alternative products may each have



sequences that the other lacks. Mutations in these differentiating segments affect different functions and should show complementation. Thus the *Drosophila Passover* gene (19) encodes two alternative proteins respectively required for normal synaptic connections in the neural system and for some other function necessary for viability. Certain mutations, although allelic in the sense that they affect the same primary transcript, show complementation because they damage different proteins (Fig. 5).

**Figure 5.** Complementation between alleles at the *Passover* gene of *Drosophila melanogaster* due to alternative modes of splicing introns from the pre-mRNA. The  $pas^V$  (recessive lethal) and  $pas^N$  (recessive neural dysfunction) mutations complement each other to produce a wild-type phenotype. They fall within exons that are mutually exclusive in the differently spliced mRNAs. Mutations ( $pas^{VN}$ ) that fall within the common upstream exons complement neither  $pas^V$  nor  $pas^N$ . Fill open boxes represent translated and untranslated exon sequences, respectively, and the lines between them are introns (1



## 7. Morphologically Pleiotropic Genes

**Pleiotropic** means having two or more distinct effects. If a pleiotropic gene mutates so as to lose one of its functions and not the other(s), then there are obvious possibilities for allelic complementation. This point has already been made in connection with genes that encode multifunctional enzymes, and it also applies to pleiotropic genes that have effects at the morphological level.

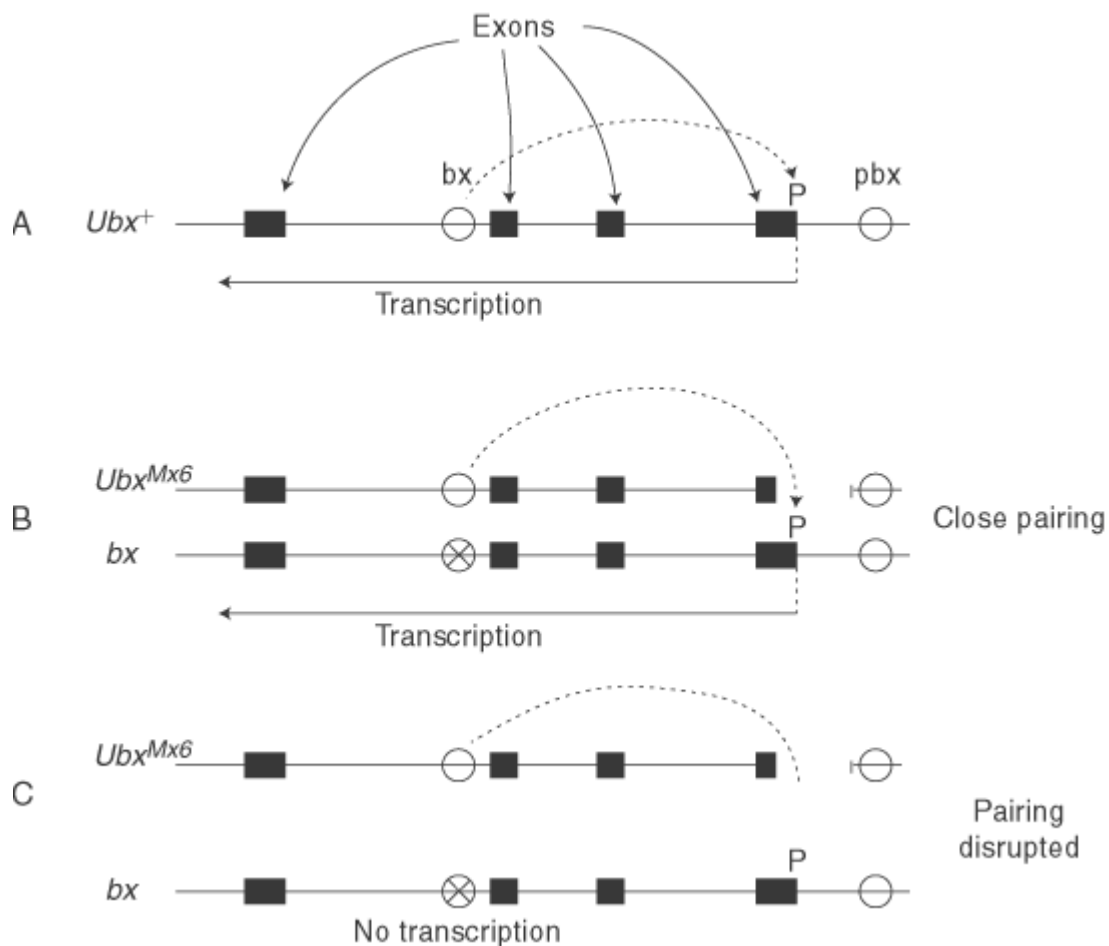
A good example is provided by the *dumpy* (*dp*) gene of *Drosophila melanogaster* (20). Different recessive mutant alleles of this gene, when homozygous, have three different effects: (1) a skewing of the thoracic bristles (called “vortex”, *v*); (2) a shortening and oblique orientation of the wings (*o*); and (3) recessive lethality (*l*). In some alleles, two or three of these functional defects occur together. *o*, *v*, and *ov* alleles complement *l* in diploids, whereas *ol*, *vl* and *ovl* do not. In combination, the effects of the alleles are additive. Each partner supplies that part of normal gene function that it possesses and the other lacks. No new function is created but each allele merely contributes what it can do by itself. A somewhat similar mutational effect but with less variety was discovered many years earlier by N. P. Dubinin in his investigation of the *Drosophila* **achaete-scute** locus (21). Here different mutations eliminated different bristles, and heterozygotes with two scute mutant alleles lacked only those bristles that were missing in both parental homozygotes.

The whole array of *dp* mutants can be attributed to a single gene on two grounds. The mutants that have lost two or three functions cannot be due to deletions overlapping more than one gene because they all form wild-type recombinants in crosses with the single-function mutants; and the sites of the different mutations, mapped by recombinational analysis, do not fall into functionally different segments but are rather intermingled.

To the extent that their multiple functions can be separately affected by mutation, pleiotropic morphology-determining genes, such as *dp* in *Drosophila*, bear some formal similarity to the “cluster genes” of fungi. But *dp*, at least, is not resolvable into single-function segments, and the reason for its pleiotropy is likely to be different. Another, better understood *Drosophila* gene, **Ultrabithorax** (*Ubx*), that encodes a regulatory protein and is differentially expressed in different tissues, may provide a better model.

The *Ultrabithorax* (*Ubx*) gene in *Drosophila* plays a vital role in establishing segment identity in the developing fly embryo (22). It is expressed in a zone extending from the end of the abdomen to the middle of the second thoracic segment. Loss of *Ubx* function by mutation causes lethal transformation of the third thoracic and first abdominal segments to likenesses of the second thoracic segment. Mutations of less drastic effect in the *Ubx* gene cause losses of *Ubx*<sup>+</sup> function in specific **parasegments**, which comprise the posterior half of one segment and the anterior half of the next segment to the rear. These mutations, it is thought, affect parasegment-specific **enhancer** sequences that are necessary in some way for the local **transcription** of *Ubx*. Fly embryos homozygous for a recessive mutation called **bithorax** (*bx*), located within the *Ubx* large intron (Fig. 6), fail to transcribe *Ubx* in parasegment 5 that covers the anterior half of the third thoracic segment (T3). T3 normally carries vestigial wings (“balancers” or halteres). The true wings are on the second thoracic segment (T2). In *bx/bx* flies, the front half of T3 becomes T2-like, and the front half of the haltere is expanded, tending toward the wing. Another recessive mutant, *postbithorax* (*pbx*), that maps just upstream of the *Ubx* transcriptional start point, has an analogous effect on *Ubx* expression in parasegment 6 that covers the posterior half of T3. When homozygous, *pbx* causes wing-like enlargement of the posterior half of the haltere. The homozygous double mutant *bx pbx/bx pbx*, develops (when it develops at all) into a four-winged fly. Both the *cis* and the *trans* double heterozygotes, *pbx bx/++* and *bx +/- pbx* are phenotypically normal.

**Figure 6.** Complementation by transvection in *Drosophila melanogaster*. (a) The normal expression of the *Ubx*<sup>+</sup> gene in the posterior thoracic segment depends on an enhancer (○) located in a *Ubx* intron. (b) In the *bithorax* (*bx*) allele (recessive, causing partial segmental transformation with enlarged halteres), this enhancer has been disrupted (Å). The *Ubx*<sup>M×6</sup> allele (recessive lethal) has a deletion that takes out the *Ubx* promoter (P). The complementation of the two alleles to give a very nearly wild-type fly (24) is due to the enhancer, still present in *Ubx*<sup>M×6</sup>, that acts in *trans* on the promoter P. (c) The effect is due to close somatic pairing of chromosomes. When this is disrupted by a structural rearrangement, the *bx* mutant phenotype reappears. The position of the upstream enhancer disrupted in *pbx* mutants is also shown. *bx* and *pbx* complement each other fully in *trans* (see text).



Based on their complementation in *trans*, one might, in the absence of further information, regard *bx* and *pbx* as mutations in different genes. This, indeed, is how they were originally interpreted. However, neither allele is complemented to more than a very limited extent by a *Ubx* null mutant allele, and so by this criterion they are functionally part of *Ubx*—complementary regulators of the same transcription unit. So *bx* and *pbx* can be regarded as alleles of *Ubx*, and the wild-type *bx/pbx* phenotype provides an example of another kind of allelic complementation, based on its tissue-specific controls of expression, not on the structure of the gene product or on the multiplicity of its effects (though it does have multiple effects).

The *Ubx* gene bears on the topic of allelic complementation in another and more surprising way. E. B. Lewis (23) discovered that *Ubx/bx* heterozygotes have a more extreme mutant phenotype when the somatic pairing of homologous chromosomes (which is characteristic of flies) is disrupted by a structural rearrangement with a breakpoint close to the *Ubx* locus. He called this phenomenon **transvection**. The idea is that the two mutant alleles somehow help each other toward normal function when they are physically close. It has been found in a few other *Drosophila* genes as well as *Ubx*.

More recently, the *Ubx/bx* phenomenon has been reinvestigated and given a new interpretation (24). One of the *Ubx* alleles used in the recent work, *Ubx*<sup>Mx6</sup>, is null (and recessive lethal) by virtue of a deletion that removes the promoter region of the gene. Compared with other null *Ubx* alleles, it shows unusually strong complementation with the standard *bx* allele, which has a *gypsy* transposable element (see [Transposon](#)) inserted into the largest *Ubx* intron, disrupting an enhancer of *Ubx* transcription in parasegment 5. The *Ubx*<sup>Mx6</sup>/*bx* heterozygote is very nearly wild type morphologically except when the somatic pairing of the homologous chromosomes is impeded by a

structural rearrangement, when the front half of the haltere becomes enlarged. The proposed explanation of the strong complementation is that the  $bx^+$  enhancer, still present in the *Ubx* allele, acts in *trans* on the enhancer-deficient *bx* allele. The idea is shown in Fig. 6. [Enhancers](#) are usually supposed to act only in *cis* and bind proteins that help to form the transcriptional initiation complex at the gene promoter on the same chromosome, but this interpretation of transvection suggests that, with close pairing of homologues (as occurs in *Drosophila* polytene chromosomes), an enhancer on one homologue can reach across and help to activate the allele on the other chromosome of the pair. In the case of this particular *Ubx/bx* interaction, the *trans* effect may have been particularly strong because, as a result of the deletion in the *Ubx* allele, there was no competition for the services of the enhancer from a promoter in *cis*.

If we grant that the enhancer is a part of the gene (difficult to deny when it is within one of the gene's introns), then transvection is yet another mechanism for allelic complementation.

### Bibliography

1. D. G. Catcheside and A. Overton (1959) Cold Spring Harbor Symp. Quant. Biol. **23**, 137–140.
2. U. Leupold and H. Gutz (1965) In Genetics Today (Proc. 11th Int. Congr. Genetics) **2**, 31–33.
3. M. E. Case and N. H. Giles (1960) Proc. Natl. Acad. Sci. USA **46**, 657–676.
4. D. O. Woodward, C. W. H. Partridge, and N. H. Giles (1958) Proc. Natl. Acad. Sci. USA **44**, 1237–1244.
5. J. R. S. Fincham (1962) J. Mol. Biol. **4**, 257–274.
6. A. Coddington, T. K. Sundaram, and J. R. S. Fincham (1966) J. Mol. Biol. **17**, 503–512.
7. D. H. Watson and J. C. Wootton (1978) Biochem. J. **175**, 1125–1133.
8. M. J. Schlesinger and C. Levinthal (1963) J. Mol. Biol. **7**, 1–12.
9. F. H. C. Crick and L. E. Orgel (1964) J. Mol. Biol. **8**, 161–165.
10. J. R. S. Fincham and A. J. Baron (1977) J. Mol. Biol. **110**, 627–642.
11. A. Ullman and D. Perrin (1970) In *The Lactose Operon* (edited by J. R. Beckwith and D. Zipser, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 143–172.
12. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning*, 2nd ed. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York pp. 4.7–4.9.
13. P. Tauro et al. (1974) Mol. Gen. Genet. **12**, 131–148.
14. A. Knobling, D. Schiffman, H. D. Sickinger, and E. Schweitzer (1975) Eur. J. Biochem. **56**, 359–367.
15. H. W. Rines, M. E. Case, and N. H. Giles (1969) Genetics **61**, 789–800.
16. M. Denis-Duphil and J. G. Kaplan (1976) Mol. Gen. Genet. **145**, 259–271.
17. J. N. Freund and B. P. Jarry (1987) J. Mol. Biol. **193**, 1–13.
18. O. G. Fahmy and M. J. Fahmy (1959) Nature **184**, 1927–1929.
19. S. H. Krishnan, E. Frei, A. P. Schalet, and R. J. Wyman (1995) Proc. Natl. Acad. Sci. USA **92**, 2021–2025.
20. E. A. Carlson (1959) Genetics **44**, 347–373.
21. N. P. Dubinin (1933) J. Genet. **26**, 37–58.
22. P. A. Lawrence (1992) *The Making of a Fly*, Blackwell Scientific, Oxford.
23. E. B. Lewis (1955) Amer. Nat. **89**, 73–89.
24. A. Martinez-Laborda, A. Gonzales-Reyes, and G. Morata (1992) EMBO J. **11**, 3543–3562.

## Interband

The dense bands of [polytene chromosomes](#) stain more intensely with DNA-sensitive dyes, indicating that the DNA concentration and hence the compaction of [chromatin](#) are greater in these regions. These bands broadly correlate with the number and positions of known genetic loci (1). Nevertheless, each band contains much more DNA than found in a typical gene. In certain well-defined cases, a band contains small clusters of genes, such as those encoding the chorion proteins. In contrast, the [interband](#) regions contain less condensed DNA and chromatin, as has been shown directly for the *Notch* gene. [In situ hybridization](#) reveals that the 5' regulatory regions of the *Notch* gene lie in the interband, and the gene itself occupies a single cytologically defined band (2). Transcriptionally active bands become differentiated structures known as **puffs**, in which the chromosomal structure is disrupted as a consequence active [transcription](#).

### Bibliography

1. B. H. Judd and M. W. Young (1973) Cold Spring Harbor Symp. Quant. Biol. **38**, 573–579.
2. M. C. Rykowski, S. J. Parmalee, D. A. Agard, and J. W. Sedat (1988) Cell **54**, 46–72.

### Suggestion for Further Reading

3. D. DePomerai (1990) *From Gene to Animal* 2nd ed., Cambridge University Press, Cambridge, UK.

## Intercalation

Virtually any molecule with sufficiently large aromatic chromophore will bind to the double-helical [DNA structure](#) by insertion of the flat chromophore between the [base pairs](#). The primary driving force for intercalative binding is stacking (charge-transfer and dipole-induced dipole) interactions between the intercalator and the adjacent base pairs, and many compounds have larger and more electron-deficient chromophores than a DNA base pair. The residence time of the intercalator at a particular site varies widely according to its structure, ranging from a few milliseconds to many hours. Intercalation disrupts the DNA structure, forcing the base pairs apart by about 3.5 Å and unwinding the helix by 10–25° (depending on the structure of the intercalator) (1). **Acridines** are classic examples of DNA intercalators, and this DNA interaction is considered integral to their activity as frameshift mutagens.

A wide range of other planar polyaromatic molecules apart from acridines have been shown to intercalate and also to cause [frameshift mutations](#) (2, 3). The most common types of frameshift mutations involve the gain or loss of one or two base pairs, typically in repetitive sequences of DNA. In these cases, every codon after the mutation site is altered within a coding region, typically resulting in the formation of truncated proteins because of the generation of [nonsense mutations](#). Nonfunctional gene products are thus a common result of frameshift mutations, and these may have major phenotypic effects, depending upon the role of that specific gene function in cell viability and metabolism.

In 1966, Streisinger et al. (4) developed a model in which he suggested that frameshift mutations are generated because of a localized pairing out of register, or slippage, during [DNA replication](#). This

slippage occurs at repetitive runs of one or more bases, leaving an unpaired base or a few bases out of register, ultimately resulting in the incorporation of an inappropriate number of bases during DNA replication. Intercalating agents may bind to the DNA bulges associated with the mispairing, thereby stabilizing these structures and enhancing the spontaneous rates of frameshift mutagenesis at this site. This model, albeit with modifications, still provides a useful explanation of most of the data available to the present day (5, 6).

## Bibliography

1. L. P. G. Wakelin (1986) *Med. Res. Rev.* **6**, 275–340.
2. W. A. Denny, P. M. Turner, G. J. Atwell, G. W. Rewcastle, and L. R. Ferguson (1990) *Mutat. Res.* **232**, 233–241.
3. L. R. Ferguson, P. M. Turner, and W. A. Denny (1990) *Mutat. Res.* **232**, 337–343.
4. G. Streisinger, Y. Okada, J. Emrich, J. Newton, et al. (1966) *Cold Spring Harbor Symp. Quant. Biol.* **31**, 77–84.
5. L. R. Ferguson and W. A. Denny (1990) *Mutagenesis* **5**, 529–540.
6. G. R. Hoffmann and R. P. P. Fuchs (1997) *Chem. Res. Toxicol.* **10**, 347–359.

## Interferons

### 1. General Introduction to the Interferon System

Interferons (IFNs) are members of the group of secretory proteins called cytokines and are defined by their ability to interfere with viral infection. However, IFNs are multifunctional proteins that also inhibit cell proliferation, regulate cell differentiation, regulate apoptosis (programmed cell death), and modulate activities of the immune system. A wide range of cell types produce IFNs in response to infection by viruses, bacteria, and mycoplasma. Noninfectious agents such as growth factors, other cytokines, and double-stranded RNA (dsRNA) molecules also induce IFN expression. In fact, dsRNA may be a physiologically relevant mediator of IFN synthesis since it is an important intermediate in viral replication. IFNs are extremely potent biological response modifiers with specific activities of at least  $10^6$  units per milligram of protein, with 1 unit defined as the amount of IFN necessary to inhibit virus replication by 50%.

IFNs were originally classified according to the cell types that generally produce them. The type I IFNs (IFN $\alpha$  and IFN $\beta$ ) are usually produced by leukocytes and fibroblasts, respectively, while type II IFN (IFN $\gamma$ ) is produced by lymphocytes. However, cloning of individual human type I IFN genes has shown that the family is rather large and includes 15 different IFN $\alpha$  subtypes (with four pseudogenes), one IFN $\beta$ , and two IFN $\omega$  (with five pseudogenes). Type I IFNs are acid-stable, have similar protein structure and biological activities, bind and transduce signals through a common multiprotein cell surface receptor, are induced in response to viral and other inducers, and share a common gene locus on human chromosome 9. All the genes for type I IFNs are devoid of introns. In contrast, the single type II IFN, IFN $\gamma$ , is acid-labile, and its expression is limited mainly to T cells and NK cells. The structural gene for IFN $\gamma$  is located on human chromosome 12 and contains three introns.

Although IFNs share similar biological actions, there are distinctive differences in the biological actions of type I and type II IFNs. For example, only IFN $\gamma$  induces class II histocompatibility antigens (all IFNs induce class I expression) and inhibits the Th2 subset of helper T cells, thereby altering the balance between the humoral and cellular immune responses. Therefore, while IFN $\gamma$  is

more active in immunomodulation than type I IFNs, type I IFNs in general are more efficient antiviral agents. Thus, although all IFNs are antiviral, the biological roles of the different IFNs in vivo are distinct, as shown in mouse models for infectious disease. Mice lacking a functional type I IFN system (by knockout of one of the type I IFN receptor genes) were highly sensitive to a wide variety of viral infections (1, 2). In contrast, mice lacking a functional type II IFN system (by knockout of the IFN $\gamma$  gene or a IFN $\gamma$  receptor gene) were more sensitive to bacterial and parasitic infections, but generally remained resistant to viral infections (1).

IFNs are not directly antiviral, but elicit their biological effects by first binding to specific receptors on the surface of target cells. The interaction of the ligand with its cognate receptor transduces signals to the nucleus that results in new gene expression: the transcription of new mRNAs and their translation into new proteins. While many of the proteins induced by type I and type II IFNs are distinct, there is a certain degree of overlap. Nonetheless, the IFN-induced proteins are ultimately responsible for biological actions of IFNs. The transcription of some IFN-induced genes does not require de novo protein synthesis, such as the family of immediate early genes called IFN-stimulated genes (ISGs). Other IFN-induced genes require protein synthesis and are involved in the delayed response to IFNs, such as class II histocompatibility antigens and IFN-response factors (IRFs).

## 2. IFN Receptors

All type I IFNs (IFN  $\alpha$ ,  $\beta$  and  $\omega$ ) bind to a ubiquitously expressed cell surface receptor, varying in number from 500 to 20,000 high-affinity ( $K_d \approx 20$ –100 pM) receptors/cell (3, 4). The cDNAs coding for two subunits of the type I IFN receptor, IFNAR1 and IFNAR2, have been cloned; both of these are located on human chromosome 21 (5-8). The human IFNAR1 chain (~110–135 kDa glycoprotein) is a 557–amino acid protein, with a 100–amino acid intracellular domain and a 21–amino acid transmembrane segment. The IFNAR1 chain is a signal transducing subunit, but does not directly bind IFN (9). The major ligand binding chain of the type I IFN receptor is IFNAR2, which exists as three different polypeptides due to alternative splicing and differential use of polyadenylation sites (7). The longest form, IFNAR2c (~100 kDa glycoprotein), is a 515–amino acid protein that acts as a functional receptor in type I IFN signaling.

Two subunits of the type II IFN $\gamma$  receptor, IFNGR1 and IFNGR2, have also been cloned (10-12). The human IFNGR1 chain (90-kDa glycoprotein) is the ligand-binding subunit of the receptor and is a 472–amino acid protein with a cytoplasmic domain of 221 amino acids and a 23–amino acid transmembrane segment. The gene for IFNGR1 is located on human chromosome 6. Although transfection of the human IFNGR1 gene into murine cells confers high-affinity ligand binding, a species-specific component (IFNGR2) encoded by human chromosome 21 is necessary for a fully functional receptor. The human IFNGR2 chain (~65-kDa glycoprotein) is a 316–amino acid polypeptide with a 66–amino acid cytoplasmic domain and a 24–amino acid transmembrane segment.

Cytokine receptors have been classified into two groups on the basis of structural similarity in their extracellular domains. The class I cytokine receptor family includes receptors for growth hormone, prolactin, erythropoietin, interleukins-2, -3, -4, -6, and -7, granulocyte-CSF, and granulocyte macrophage-CSF and are typified by conserved tryptophan, NH<sub>2</sub>-terminal and COOH-terminal cysteine pairs, and a membrane proximal WSxWS box. The IFNAR1, IFNAR2, IFNGR1, IFNGR2 and the interleukin (IL)-10 receptor constitute the type II cytokine receptor family, based on a characteristic 200–amino acid extracellular domain (13). These receptors have no inherent enzyme activity, but on ligand binding generate cytoplasmic signals in combination with JAK tyrosine kinases and STAT proteins. Members of this family in general have a ligand binding chain with either low or high affinity, and at least one chain involved in signal transduction, which has either relatively low binding affinity for the cytokine or no detectable binding activity at all. Receptor activation by this family of cytokines requires the ligand-dependent interaction between receptor subunits to form oligomers.

### 3. Signal Transduction by IFNs Through the JAK-STAT Pathway

IFNs elicit their biological effects by transducing a signal to the nucleus, which results in selective gene expression (14, 15). ISGs are a family of genes that are transcriptionally activated within minutes by IFNs. ISG transcriptional activation is regulated through the tyrosine phosphorylation of STAT proteins (for *S*ignal *T*ransducers and *A*ctivators of *T*ranscription), which is mediated by members of the Janus family of nonreceptor protein tyrosine kinases (JAKs) (16-19). JAKs are so named because they have a functional kinase domain as well as a second kinase-like domain, and hence are two-faced like the Roman god Janus. JAKs constitutively associate with membrane proximal regions of the intracellular domains of IFN receptor subunits to provide catalytic activity and transduce intracellular signals by phosphorylating a variety of substrates on tyrosine residues (20). In the case of the type I IFN receptor, the TYK2 and JAK1 Janus kinases are constitutively bound to the IFNAR1 and IFNAR2 subunits, respectively (8, 21). On ligand binding the JAKs cross-phosphorylate each other, with Tyk2 becoming phosphorylated first on type I IFN addition. In the case of the IFN $\gamma$  receptor, the JAK1 and JAK2 kinases are constitutively bound to the IFNGR2 and IFNGR1 subunits, respectively, and on ligand binding JAK1 appears to be activated first (22, 23). The activated JAKs phosphorylate specific tyrosine residues within the cytoplasmic tails of the IFN receptor chains to provide docking sites for STAT proteins (24, 25). STAT proteins are transcription factors and hence contain DNA-binding domains (with the exception of STAT2, which apparently can only function in a complex with STAT1). STAT proteins also contain src homology-2 (SH2) regions (19), which are modular noncatalytic domains of about 100 amino acids that are found in various cytoplasmic signaling proteins. SH2 domains recognize phosphotyrosine motifs, with specificity dictated by the primary sequence immediately COOH-terminal from the phospho-Tyr (26). The separate docking sites for the binding domains of SH2 domains from different signaling molecules have been mapped on many receptors and cytoplasmic signaling proteins.

Type I IFNs activate STAT1 and STAT2 by inducing their phosphorylation on discrete tyrosine residues. These STATs dimerize, complex with a 48-kDa DNA-binding protein, and then move into the nucleus, recognizing a highly conserved IFN-stimulus response element (ISRE) in the promoter of ISGs to directly activate these genes (27, 28). This conserved 15-base pair element (consensus: AGTTTCNNTTTCNC/T) is present in the promoters of nearly 20 type I IFN-inducible genes examined so far and is both necessary and sufficient for gene induction.

As an added complexity, some type I IFN-inducible genes like IRF-1 do not contain an ISRE and are activated by STAT1 dimers (29). Type I IFNs also activate STAT3, the transcription factor for acute phase response genes. STAT3 recognizes a conserved promoter element found within these genes and is activated by a wide variety of cytokines, suggesting that it may integrate diverse signals into common transcriptional responses (25, 30-32). STAT3 binds directly to the tyrosine-phosphorylated IFNAR1 chain and also undergoes IFN-dependent tyrosine phosphorylation (25). In addition, other members of the STAT family (STAT4, STAT5, and STAT6) have been reported to be activated by type I IFNs in a cell type-restricted manner (33, 34). Although much is known about the role of STAT1-STAT2 dimers in ISG transcriptional activation, the physiological role of other STAT protein complexes in IFN-induced gene expression is not known.

For type II IFNs, a variety of response elements have been described for genes regulated by IFN $\gamma$ . Most notable is the IFN $\gamma$ -activated sequence (GAS, consensus: TTNCNNNA), which is found in IFN $\gamma$ -induced genes such as the genes for guanylate binding protein, interferon response factor (IRF)-1, and the Fc $\gamma$  receptor I. IFN $\gamma$  induces the oligomerization of the IFNGR1 and IFNGR2 chains of its cognate receptor and the activation of JAK1 and JAK2 kinases (35). The cytoplasmic tail of the IFNGR1 chain undergoes tyrosine phosphorylation and hence provides a docking site for STAT1. Subsequently, STAT1 also undergoes phosphorylation, forms homodimers, translocates to the nucleus, and binds to the GAS element, resulting in the transcriptional activation of IFN $\gamma$ -induced genes. As is the case for type I IFNs, IFN $\gamma$  also activates several members of the STAT family, but the physiological relevance of this remains to be established.



Besides STATs, a number of other *trans*-acting factors (such as IFN-regulatory factors 1 and 2, IFN consensus sequence-binding protein, etc.) also appear to regulate IFN-inducible gene expression. For example, type I IFNs activate nuclear factor  $\kappa$ B (NF- $\kappa$ B) (36). Active DNA-binding forms of NF- $\kappa$ B are dimers of various members of the Rel/NF- $\kappa$ B family of polypeptides (p50, p52, c-Rel, v-Rel, RelA, and RelB). NF- $\kappa$ B normally exists as a dormant complex in the cytoplasm through the binding of inhibitor of  $\kappa$ B (I $\kappa$ B) proteins. Viruses, cytokines, and lipopolysaccharides induce inactive NF- $\kappa$ B/I $\kappa$ B complexes to dissociate, allowing NF- $\kappa$ B to enter the nucleus and bind DNA. NF- $\kappa$ B binds to sites in the promoters and enhancers of key cellular genes regulating apoptotic, immune, and inflammatory responses. NF- $\kappa$ B suppresses apoptosis (37-40), and hence activation of NF- $\kappa$ B by IFN protects cells from apoptotic death (36). Blockage of NF- $\kappa$ B nuclear translocation and activation by the introduction of super-repressor forms of I $\kappa$ B enhances apoptotic killing by IFN.

The signal transduction pathway activated by many other cytokines is similar to that of IFNs but involves different combinations of JAKs (there are four: JAK1, JAK2, JAK3, and TYK2) and STAT proteins (there are seven STATs). Thus, the IFN-activated JAK/STAT pathway serves as a paradigm for understanding cytokine signal transduction in general. The critical roles that JAK1, TYK2, STAT1, STAT2, and STAT3 play in type I IFN signaling on the one hand, and that JAK1, JAK2, and STAT1 play in type II IFN signaling on the other hand, have been demonstrated through the use of cell lines with defined individual defects in these signaling components (20, 35, 41). Although the tyrosine phosphorylation of STATs is responsible for transcriptional activation of IFN-stimulated genes, serine phosphorylation events are also critical for optimal IFN-induced transcription and for the biological response to IFN (42-45). A variety of serine kinases have been implicated in IFN signaling, including protein kinase C isoforms, the mitogen-activated protein kinase cascade, phosphatidylinositol 3-kinase, RNA-dependent protein kinase (PKR), and protein kinase B (PKB)/Akt (43, 45-50).

#### 4. IFN-Stimulated Genes and Their Functions

Since there is overlap in the signal transduction pathways activated by all IFNs, there is significant overlap in the genes induced by type I and type II IFNs. However, distinctive expression, as well as significant overlap, of gene transcripts regulated by either type I (a, b) or type II (g) IFNs has been demonstrated by the use of oligonucleotide arrays corresponding to more than 6,800 human genes (51). A short discussion of some of the well-characterized gene products induced by IFNs (there are over 30 different products identified so far) follows (52).

- a. Protein Kinase R (PKR) This type I IFN-induced protein kinase requires dsRNA for activity. Activated PKR phosphorylates several key cellular proteins, most notably the eukaryotic initiation factor-2a, resulting in the cessation of polypeptide chain initiation. The major cellular function of PKR appears to be growth control by regulating the rates of cellular and viral protein synthesis.
- b. 2',5' Oligoadenylate Synthetase (OAS) OAS is a family of type I IFN-induced enzymes that require dsRNA as a cofactor. OAS catalyzes the formation 2'-5' oligoadenylates and thereby activates RNase L, which cleaves viral and cellular RNAs.
- c. Mx Proteins Mx proteins are type I IFN-inducible proteins that are involved in the inhibition of myxovirus (such as influenza virus) replication in mice by a poorly defined pathway. Mx genes and proteins have been identified in other species, including humans, but their antiviral activity is broader than that of murine Mx.
- d. Indoleamine 2,3-Dioxygenase (IDO) IDO is a type I and type II IFN-inducible intracellular enzyme that degrades tryptophan to kynurenine and plays a role in the inhibition of the growth of intracellular parasites such as *Toxoplasma gondii*.
- e. IFN-Regulatory Factor-1 (IRF-1) IRF-1 is a type I and type II IFN-inducible transcription factor that binds to a DNA element found in the IFN $\beta$  promoter. IRF-1 appears to play a key regulatory role in growth control and tumorigenesis.

- f. Class I Histocompatibility Antigens Class I antigens are induced by both type I and type II IFNs and are broadly expressed on almost all cell types. Class I antigens play critical roles in antigen recognition by cytotoxic T cells (CD8+).
- g. Class II Histocompatibility Antigens Class II antigens are induced by IFN $\gamma$  and have a restricted tissue distribution (present on B cells, macrophages, and activated T cells). These antigens participate in antigen recognition by helper T cells (CD4+).

## 5. Conclusion

This article has focused on the present understanding of the highly complex IFN system. IFNs were the first cytokines to be discovered (over 40 years ago) and characterized. Rapid advances have been made in the last decade on the structure and regulation of IFN genes, the characterization and cloning of the components of the multisubunit IFN receptors, and the signal transduction pathways by which IFNs exert their effects. However, many gaps remain in our knowledge. Future work will undoubtedly refine our understanding of this complex cytokine system.

## 6. Literature Cited

An apology is made to colleagues whose work was not cited due to space restrictions.

### Bibliography

1. U. Muller et al., *Science* **264**, 1918–1921 (1994).
2. C.M. Owczarek et al., *J. Biol. Chem.* **272**, 23865–23870 (1997).
3. C. Vanden Broecke and L.M. Pfeffer, *J. Interferon Res.* **8**, 803–811 (1988).
4. L.C. Plataniias, L.M. Pfeffer, R. Cruciani, and O.R. Colamonici, *J. Immunol.* **150**, 3382–3388 (1993).
5. G. Uze, G. Lutfalla, and I. Gresser, *Cell* **60**, 225–234 (1990).
6. P. Domanski et al., *J. Biol. Chem.* **270**, 21606–21611 (1995).
7. G. Lutfalla et al., *EMBO J.* **14**, 5100–5108 (1995).
8. D. Novick, B. Cohen, and M. Rubinstein, *Cell* **77**, 391–400 (1994).
9. O.R. Colamonici et al., *J. Biol. Chem.* **269**, 9598–9602 (1994).
10. S. Hemmi et al., *Cell* **76**, 803–810 (1994).
11. M. Aguet, Z. Dembic, and G. Merlin, *Cell* **55**, 273–280 (1988).
12. J. Soh et al., *Cell* **76**, 793–802 (1994).
13. J.F. Bazan, *Proc. Natl. Acad. Sci. U.S.A.* **87**, 6934–6938 (1990).
14. A.C. Larner, A. Chaudhuri, and J.E. Darnell Jr., *J. Biol. Chem.* **261**, 453–459 (1986).
15. A.C. Larner et al., *Proc. Natl. Acad. Sci. U.S.A.* **81**, 6733–6737 (1984).
16. M. Muller et al., *Nature* **366**, 129–135 (1993).
17. L. Velazquez, M. Fellous, G.R. Stark, and S. Pelligrini, *Cell* **70**, 313–322 (1992).
18. C. Schindler, K. Shuai, V.R. Prezioso, and J.E. Darnell Jr., *Science* **257**, 809–813 (1992).
19. X.Y. Fu, *Cell* **70**, 323–335 (1992).
20. J.N. Ihle et al., *TIBS* **19**, 222–229 (1994).
21. O.R. Colamonici et al., *J. Biol. Chem.* **269**, 3518–3522 (1994).
22. D. Kaplan et al., *J. Biol. Chem.* **271**, 9–12 (1996).
23. M. Sakatsume et al., *J. Biol. Chem.* **270**, 17528–17534 (1995).
24. A.C. Greenlund, M.A. Farrar, B.L. Viviano, and R.D. Schreiber, *EMBO J.* **13**, 1591–1600 (1994).

25. C.H. Yang et al., *J. Biol. Chem.* **271**, 8057–8061 (1996).
26. T. Pawson and J. Schlessinger, *Curr. Biol.* **3**, 434–442 (1993).
27. R.L. Friedman and G.R. Stark, *Nature* **314**, 637–639 (1985).
28. N.C. Reich et al., *Proc. Natl. Acad. Sci. U.S.A.* **84**, 6394–6398 (1987).
29. R. Pine, A. Canova, and C. Schindler, *EMBO J.* **13**, 158–167 (1994).
30. G.S. Campbell et al., *J. Biol. Chem.* **270**, 3974–3979 (1995).
31. T. Boulton et al., *Proc. Natl. Acad. Sci. U.S.A.* **92**, 6915–6919 (1995).
32. Z. Zhong, Z. Wen, and J.E. Darnell Jr., *Science* **264**, 95–98 (1994).
33. E. Fasler-Kan et al., *Eur. J. Biochem.* **254**, 514–519 (1998).
34. J.D. Farrar et al., *J. Biol. Chem.* **275**, 2693–2697 (2000).
35. J.E. Darnell Jr., I.M. Kerr, and G.R. Stark, *Science* **264**, 1415–1421 (1994).
36. C.H. Yang et al., *Proc. Natl. Acad. Sci. U.S.A.* **97**, 13631–13636 (2000).
37. A.A. Beg et al., *Nature* **376**, 167–170 (1995).
38. A.A. Beg and D. Baltimore, *Science* **274**, 782–784 (1996).
39. D.J. Van Antwerp et al., *Science* **274**, 787–789 (1996).
40. C.Y. Wang, M.W. Mayo, and A.S. Baldwin Jr., *Science* **274**, 784–787 (1996).
41. C.H. Yang, A. Murti, and L.M. Pfeffer, *Proc. Natl. Acad. Sci. U.S.A.* **95**, 5568–5572 (1998).
42. L.M. Pfeffer et al., *Proc. Natl. Acad. Sci. U.S.A.* **88**, 7988–7992 (1991).
43. N.C. Reich and L.M. Pfeffer, *Proc. Natl. Acad. Sci. U.S.A.* **87**, 8761–8765 (1990).
44. H. Akai and A.C. Lerner, *J. Biol. Chem.* **264**, 3252–3255 (1989).
45. C.R. Faltynek et al., *J. Biol. Chem.* **264**, 14305–14311 (1989).
46. L.M. Pfeffer, B. Strulovici, and A.R. Saltiel, *Proc. Natl. Acad. Sci. U.S.A.* **87**, 6537–6541 (1990).
47. M. David et al., *Science* **269**, 1721–1723 (1995).
48. S. Uddin et al., *J. Immunol.* **158**, 2390–2397 (1997).
49. L.M. Pfeffer et al., *Science* **276**, 1418–1420 (1997).
50. C.H. Yang et al., *J. Biol. Chem.* **276**, 13756–13761 (2001).
51. S.D. Der, A. Zhou, B.R. Williams, and R.H. Silverman, *Proc. Natl. Acad. Sci. U.S.A.* **95**, 15623–15628 (1998).
52. D.H. Kalvakolanu and E.C. Borden, *Cancer Invest.* **14**, 25–53 (1996).

### **Additional Reading**

53. Pfeffer L.M., *Mechanisms of Interferon Action*, CRC Press, Boca Raton, Fla., 1987.
54. Pestka S., Langer J.A., Zoon K.C., and Samuel C.E., *Rev. Biochem.* **56**, 727–777 (1987).
55. Darnell J.E. Jr., Kerr I.M., and Stark G.R., *Science* **264**, 1415–1421 (1994).
56. Ihle J.N. et al., *TIBS* **19**, 222–229 (1994).
57. Kalvakolanu D.H. and Borden E.C., *Cancer Invest.* **14**, 25–53 (1996).
58. Leaman D.W., *Prog. Mol. Subcell. Biol.* **20**, 101–42 (1998).
59. Sen G.C., *Semin. Cancer Biol.* **10**, 93–101 (2000).

### **Interleukin-1 Motif**

The interleukin-1 motif is a common [protein motif](#) found in [protein structures](#). It is present in those of interleukin-1, [fibroblast growth factors](#), and Kunitz-type [trypsin inhibitors](#). The motif is characterized by an unusual pseudo-threefold symmetry. The secondary structure of the interleukin-1 motif has 12 [b-strands](#) that form six **hairpins**. Three of the hairpins form a six-stranded **antiparallel b-barrel** with a **hydrophobic core**; the other three form a hairpin triplet structure that caps the barrel (Fig. 1). There is little sequence identity between proteins adopting this three-dimensional structure, apart from a preference for large hydrophobic residues in those regions of the sequence that form the core.

**Figure 1.** Schematic representation of the backbone structure of human interleukin-1b (1) showing the characteristic fold. Strands are shown as arrows, and connecting loops are in yellow. The three roughly equivalent repeating b-structural units are shown in different colors (red, green, and blue) to highlight their similarity. Two orthogonal orientations are shown. (Left) The six-stranded b-barrel forms the lower half of the structure and is capped by the triangular array of b-hairpins. (Right) view is shown looking down the axis of the triangular array and b-barrel showing the pseudo-threefold symmetry of the structure. This figure was generated using Molscript (2) and Raster3D (3, 4). See color insert.



The interleukin-1 motif is sometimes called a b-trefoil, because it can be dissected into three equivalent Y-shaped (or trefoil) b-structural units. However, the b-trefoil terminology should not be confused with the “trefoil peptide domain” that represents a very different structural motif of gastrointestinal protective peptides. Furthermore, it should be noted that the known structures of all other [interleukins](#) (apart from interleukin-1) are **four-helix bundles**.

[See also [Beta-Sheet](#) and [Antiparallel Beta-Barrel Motifs](#).]

#### Bibliography

1. J. P. Priestle, H.-P. Schar, and M. G. Grutter (1989) Proc. Natl. Acad. Sci. USA **86**, 9667–9671.
2. P. J. Kraulis (1991) J. Appl. Crystall. **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) Acta Crystallogr. **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) J. Mol. Graphics **6**, 219–222.

#### Suggestions for Further Reading

5. M. B. Swindells and J. M. Thornton (1993) A study of structural determinants in the interleukin-1 fold. *Protein Eng.* **6**, 711–715.
6. A. G. Murzin, A. M. Lesk, and C. Chothia (1992) -Trefoil fold: patterns of structure and sequence in the Kunitz inhibitors, interleukins-1 and -1 $\beta$ , and fibroblast growth factors. *J. Mol. Biol.* **223**, 531–543.

## Interleukins

Multicellular organisms are composed of a few hundred cell types whose survival, growth, differentiation, and effector function in specific tissues or organs need to be exquisitely controlled. This is accomplished mostly by a myriad of secreted regulatory proteins, broadly classified as cytokines (1), that are secreted locally and modulate the cell physiology of each tissue. Cytokines embrace a wide variety of regulatory proteins that include [growth factors](#), [colony-stimulating factors](#), interleukins, [lymphokines](#), monokines, and [interferons](#). Typically, the molecular masses of cytokines range from a few thousand to a few tens of thousands. Cytokines, in contrast to [hormones](#), are synthesized locally in a transient manner, often by a variety of nonspecialized cells, and function locally at their sites of production. In general, their half-life is short. Classical hormones (eg, [insulin](#)), on the other hand, are produced by specialized cells within particular organs, are disseminated throughout the body via the circulatory system, and act at a distance from their sites of production.

The interleukins constitute a subset of cytokines that are produced by and regulate cells involved in immune defense mechanisms. The generic term interleukin (IL) and a corresponding numbering system of nomenclature were introduced in 1979 (2) to reflect that these regulatory proteins are communicating signals for leukocytes, and not just lymphocytes. However, this system has not been completely successful, because now it has been shown that some interleukins may be produced by nonleukocytes and act on other cell types. At present, new leukocyte regulatory proteins are always assigned IL numbers upon establishing a gene sequence. IL-18 is the most recent.

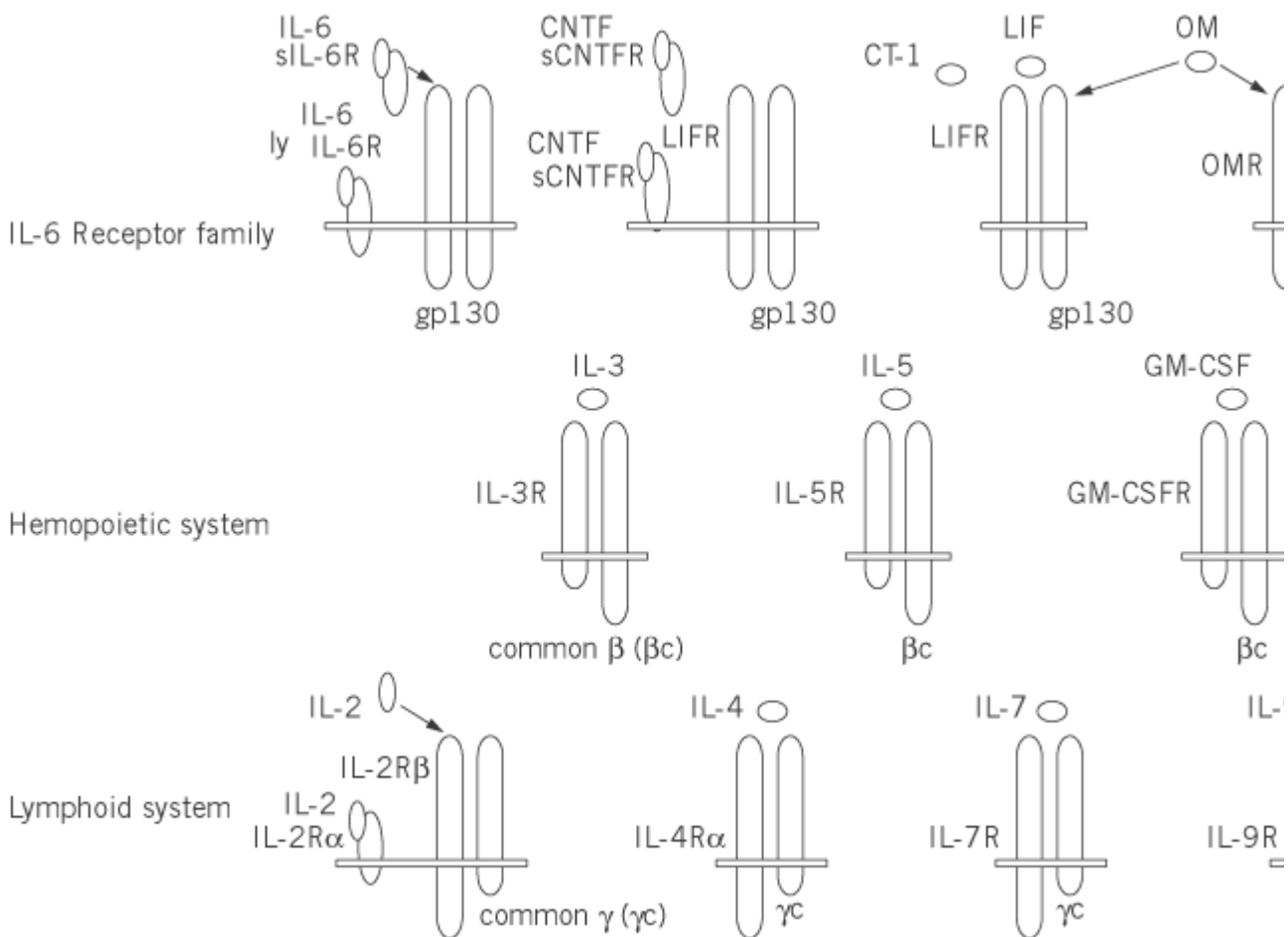
### 1. General Properties of Interleukins

Table 1 lists the major biological actions of the interleukins. These proteins, which have been characterized most thoroughly in the human and mouse systems, are noted for their pleiotropy and redundant biological activity. The receptors for many of the interleukins, such as IL-1, -2, -3, -5, and -6, have been well characterized. Typically, they consist of multiple protein subunits. Shared signalling receptor subunits are seen for a number of cytokines, for example, IL-3/IL-5/granulocyte-macrophage colony stimulating factor and also IL-2/IL-4/IL-7/IL-9/IL-15 (see Fig. 1). To illustrate this further, Fig. 1 shows the structure of the IL-6 receptor complex as a paradigm for the interleukins. The biological activity of IL-6 is exerted upon binding with low affinity to its cognate receptor (IL-6 receptor or  $\alpha$ -subunit) located on the surface of its target cells. Then, the low-affinity IL-6/IL-6 receptor complex binds with high affinity to the  $\beta$ -subunit of the receptor complex, the **signal-transducing** molecule gp130, to exert the biological activities of this interleukin. The subunit gp130 is also shared by IL-11 and even other cytokines (eg, ciliary neurotrophic factor, leukemia inhibitory factor, oncostatin M, and cardiotropin-1) which, in part, explains the overlapping activities of these cytokines at the molecular level.

**Table 1. Biological Actions of the Interleukins**

| Interleukin | Action  | References                              |
|-------------|---|---|
| IL-1a/b     | Stimulates thymocyte proliferation; acute phaseresponse   | <a href="#">5</a> , <a href="#">6</a>   |
| IL-2        | Supports proliferation and activation of T lymphocytes  | <a href="#">7</a>                       |
| IL-3        | Stimulates hematopoiesis; production of cells in all myeloid lineages   | <a href="#">8</a>                       |
| IL-4        | Stimulates proliferation of B lymphocytes   | <a href="#">9</a>                       |
| IL-5        | Activation and differentiation of eosinophils; accessory factor neutrophils   | <a href="#">10</a> , <a href="#">11</a> |
| IL-6        | Stimulates hemopoietic progenitor cells, B-lymphoid, and myeloma cells; induction of acute phase protein synthesis in the liver     | <a href="#">12</a> , <a href="#">13</a> |
| IL-7        | Stimulates proliferation of thymocytes, T lymphocytes including cytotoxic T cells, and early B-lymphocyte precursors                | <a href="#">14</a>                      |
| IL-8        | Stimulates chemotaxis by neutrophils and other cell types   | <a href="#">15</a>                      |
| IL-9        | Stimulates growth of T-lymphoid cell lines; activation of mast cells; potentiates proliferative effect of IL-2 on fetal lymphocytes | <a href="#">16</a>                      |
| IL-10       | Suppresses cytokine production by T lymphocytes   | <a href="#">17</a> , <a href="#">18</a> |
| IL-11       | Hematopoietic growth factor; induces acute phase protein synthesis  | <a href="#">19</a>                      |
| IL-12       | Stimulates cytokine production by natural killer cells and T cells  | <a href="#">20</a>                      |
| IL-13       | Regulates the function of monocytes and B cells; anti-inflammatory action   | <a href="#">21</a>                      |
| IL-14       | Stimulates proliferation of B cells   | <a href="#">22</a>                      |
| IL-15       | Stimulates proliferation of T cells   | <a href="#">23</a>                      |
| IL-16       | Lymphocyte chemoattractant factor; modulates T-cell activation  | <a href="#">(24)</a>                    |
| IL-17       | Stimulates production of proinflammatory and hemopoietic cytokines in epithelial, endothelial, fibroblastic, or stromal cells       | <a href="#">25</a>                      |
| IL-18       | Stimulates proliferation of T cells; induces IFN-g production   | <a href="#">26</a>                      |

**Figure 1.** Structure of the Intelleukin-6 receptor complex as a paradigm for interleukins. IL-, interleukin; R, receptor; CN leukemia inhibitory factor; OM, oncostatin-M, GM-CSF, granulocyte macrophage-colony stimulating factor; gc, commo



The signal-transduction mechanisms used by the interleukins is presently being unraveled. Although binding of interleukin to its cognate receptor complex stimulates protein **tyrosine kinase** activity, the receptors themselves do not possess any intrinsic tyrosine kinase activity. Rather, nonreceptor tyrosine kinases, such as members of the **src** and **JAK** families, may become associated with the interleukin receptor signaling subunits (typically, binding a ligand to its receptor induces the dimerization of signaling subunits) leading to phosphorylation of tyrosine residues of the signalling subunit. These events lead to the activation of multiple signal-transduction pathways, such as the **STAT**, **Ras-MAP kinase** and **phosphatidyl inositol kinase** pathways (3, 4).

Recent reviews for all of the interleukins that outline their role in pathogenesis and disease therapy are listed in Table 1. Other features of the interleukins are listed in Table 2. For access to current information regarding the interleukins (as well as other cytokines), the reader is referred to the following internet web sites:

**GroCyte** home page: (site currently unavailable)

**COPE** (Cytokines Online Pathfinder Encyclopedia) home page:

<http://bioinfo.weizmann.ac.il/cgi-bin/cope/cope.pl>

**Table 2. Properties of Interleukins**

---

|           |       |                            |               |        |
|-----------|-------|----------------------------|---------------|--------|
| Molecular | Amino | Glycosylation <sup>c</sup> | Cysteines and | Three- |
|-----------|-------|----------------------------|---------------|--------|

| Interleukin <sup>a</sup> | Weight<br>( $\times 10^{-3}$ ) <sup>b</sup> | Acid<br>Residues       | pI  | Disulfide |    | Disulfide<br>Connectivities <sup>d</sup>   | Structure <sup>e</sup>  | Dimensional<br>Structure          |
|--------------------------|---|------------------------|-----|-----------|----|--|---|-----------------------------------|
|                          |   |                        |     | N-        | O- |  |   |                                   |
| IL-1a (IL1A<br>Human)    | 18  | 159                    | 5.3 | 1N-       |    | 1  | IL-1R type<br>I<br>IL-1R type<br>II                                   | X ray (2II)                       |
| IL-1b (IL1B<br>Human)    | 17.4  | 153                    | 6.0 | 1N-       |    | 2  | IL-1R type<br>I<br>IL-1R type<br>II                                   | X ray (1I)<br>NMR (6I)            |
| IL-2 (IL2<br>Human)      | 15.4  | 133                    | 7.2 |           | O- | 3<br><br>(Cys58–<br>Cys105)  | IL-2Ra<br>(low<br>affinity)<br>IL-2Rb<br>and gc<br>(high<br>affinity) | X ray<br>(1ITB),<br>NMR (1I)      |
| IL-3 (IL3<br>Human)      | 15  | 133                    | 7.3 | 2N-       |    | 2<br><br>(Cys16–Cys84)   | IL-3Ra<br>IL-3Rb  | NMR (1J)                          |
| IL-4 (IL4<br>Human)      | 15  | 129                    | 9.1 | 2N-       |    | 6<br><br>(Cys2–<br>Cys126),<br>(Cys23–<br>Cys64),<br>(Cys45–Cys98)                       | IL-4R   | X ray<br>(1HIK),<br>NMR<br>(1RCB) |
| IL-5 (IL5<br>Human)      | 13.1 (26.2<br>homodimer)                    | 115 (130<br>homodimer) | 7.3 | 2N-       |    | 2 per chain (4<br>involved in<br>interchain<br>disulphide<br>homodimer)<br>(Cys44–Cys86) | IL-5Ra IL-<br>3Rb   | X ray<br>(1HUL)                   |
| IL-6 (IL6<br>Human)      | 20.8  | 183                    | 6.4 | 2N-       |    | 4<br><br>(Cys43–Cys49)<br><br>(Cys72–Cys82)  | IL-6Ra<br>(low<br>affinity)<br>gp130<br>(high<br>affinity)            | X ray,<br>(1ALU)<br>NMR (1I)      |
| IL-7 (IL7<br>Human)      | 17.4  | 152                    | 8.4 | 3N-       |    | 6  | IL-7R<br><br>IL-2Rgc<br>Possibly                                      | Theoretic<br>model<br>(1IL7)      |



|                                |       |       |         |     |    |   |                               |
|--------------------------------|-------|-------|---------|-----|----|---|-------------------------------|
| IL-8 (IL8 Human)               | 8–8.9 | 69–77 | 8.5–9.0 |     | 4  | another IL-8RA<br>IL-8RB<br>(Cys12–Cys39)<br>(Cys14–Cys55)        | NMR (1I)                      |
| IL-9 (IL9 Human)               | 14.1  | 126   | 8.6     | 4N- | 10 | IL-9R   |                               |
| IL-10 (IL10 Human)             | 18.6  | 160   | 7.8     | 1N- | 4  | IL-2Rgc<br>IL10R<br><br>Possibly others                           | X ray (1I)                    |
| IL-11 (IL11 Human)             | 19.1  | 178   | 11.3    |     | 0  | IL-11Ra   |                               |
| IL-12-p35 subunit (I12A Human) | 22.5  | 197   | 6.4     | 3N- | 7  | gp130<br>Possibly one other<br>IL12Rb <sub>2</sub>                |                               |
| IL-12-p40 subunit (I12B Human) | 34.7  | 306   | 5.4     | 4N- | 10 | IL12Rb <sub>1</sub>   |                               |
| IL-13 (IL13 Human)             | 12.3  | 112   | 8.5     | 4N- | 4  | IL13Ra <sub>1,a2</sub><br>IL4Ra<br>(Cys28–Cys56)<br>(Cys44–Cys70) | Theoretic model<br>(3ITR, 3I) |
| IL-14 (IL14 Human)             | 53    | 483   | 8.8     | 3N- | 34 |   |                               |
| IL-15 (IL15 Human)             | 12.8  | 114   | 4.4     | 2N- | 4  | IL15Ra<br><br>(Cys54–Cys104)<br>(Cys61–Cys107)                    | IL-2Rb & gc                   |
| IL-16 (IL16 Human)             | 13.3  | 130   | 4.7     |     | 1  | CD4   | NMR (1I)                      |
| IL-17 (IL17 Human)             | 15.1  | 132   | 8.3     | 1N- | 6  | IL17R   |                               |

<sup>a</sup> SwissProt database accession filename is shown in parenthesis: <http://www.expasy.ch/srs5/>.

<sup>b</sup> Molecular weight calculated from the amino acid weight of the expressed functional protein.

<sup>c</sup> Includes potential N-glycosylation sites (ie, Asn-X-Ser/Thr motif).

<sup>d</sup> Cysteine connectivities determined experimentally and, in some cases by homology.

<sup>e</sup> Pdb database filename is shown in parenthesis: <http://www.pdb.bnl.gov/>.

<sup>f</sup> Genbank database accession filename is shown in parenthesis: <http://www.expasy.ch/srs5/>.

## Bibliography

1. S. Cohen, P. E. Bigazzi, and T. Yoshida (1974) *Cell. Immunol.* **12**, 150–159.
2. L. A. Aarden et al. (1979) *J. Immunol.* **123**, 2928–2929.
3. P. C. Heinrich et al. (1998) *Biochem. J.* **334**, 297–314.
4. K. D. Liu, S. L. Gaffen, and M. A. Goldsmith (1998) *Curr. Opinion Immunol.* **10**, 271–278.
5. F. S. di Giovine and G. W. Duff (1990) *Immunol. Today* **11**, 13–20.
6. C. A. Dinarello (1997) *Cytol. Growth Factor Rev.* **8**, 253–265.
7. K. A. Smith (1988) *Science* **240**, 1169–1176.
8. J. W. Schrader (1998) In *The Cytokine Handbook*, 3rd ed. (A. Thomson, ed.), Academic Press, San Diego, pp. 105–132.
9. J.-L. Boulay and W. E. Paul (1992) *Curr. Opinion Immunol.* **4**, 294–298.
10. K. Takatsu (1992) *Curr. Opinion Immunol.* **4**, 299–306.
11. C. J. Sanderson, H. D. Campbell, and I. G. Young (1988) *Immunol. Rev.* **102**, 29–50.
12. R. J. Simpson et al. (1997) *Protein Sci.* **6**, 929–955.
13. T. Kishimoto, S. Akira, M. Narazaki, and T. Taga (1995) *Blood* **86**, 1243–1254.
14. R. Murray (1996) *Curr. Opinion Hematol.* **3**, 230–234.
15. M. Baggiolini, B. Dewald, and B. Moser (1994) *Adv. Immunol.* **55**, 97–179.
16. J.-C. Renaud et al. (1993) *Adv. Immunol.* **54**, 79–97.
17. R. de Waal Malefyt et al. (1992) *Curr. Opinion Immunol.* **4**, 314–320.
18. M. Howard and A. O'Garra (1992) *Immunol. Today* **13**, 198–200.
19. X. Du and D. A. Williams (1997) *Blood* **89**, 3897–3908.
20. G. Trinchieri (1993) *Immunol. Today* **14**, 335–338.
21. G. Zurawski (1994) *Immunol. Today* **15**, 19–26.
22. J. L. Ambrus et al. (1993) *Proc. Nat. Acad. Sci. USA* **90**, 6330–6334.
23. S. Bulfone-Paus et al. (1997) *Nat. Med.* **3**, 1124–1128.
24. T. Kishimoto (1996) *J. Immunol. Methods* **196**, 103–104.
25. F. Fossiez et al. (1996) *J. Exp. Med.* **183**, 2593–2603.
26. M. T. Gillespie and N.J. Horwood (1998) *Cytol. Growth Factor Rev.* **9**, 109–116.

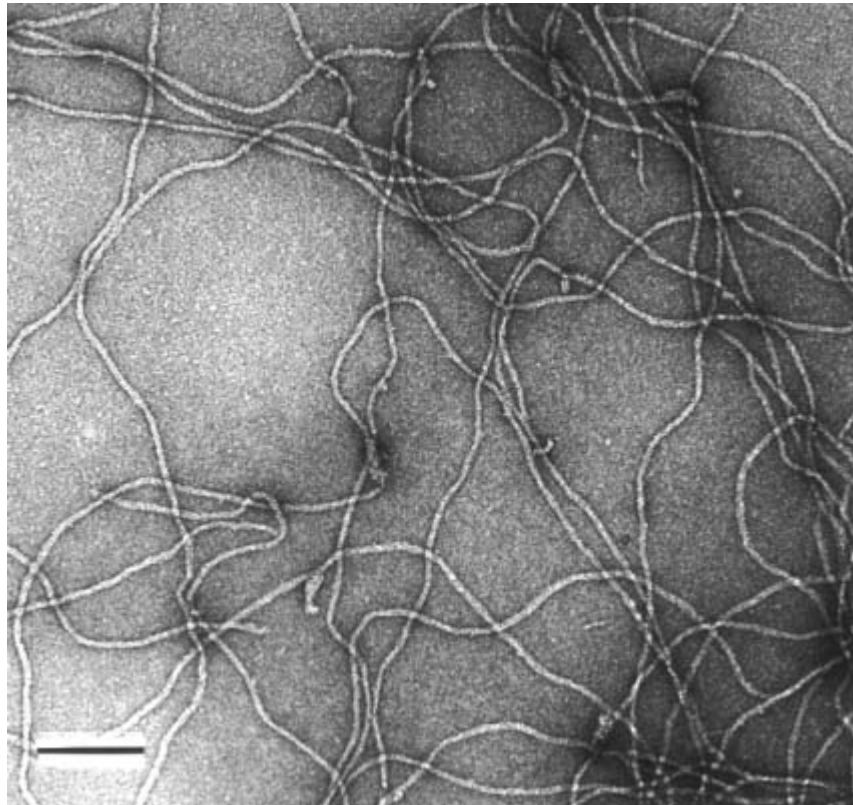
## Suggestions for Further Reading

27. A. W. Thomson (ed.) (1998) *The Cytokine Handbook* 3rd ed., Academic Press, San Diego.
28. N. A. Nicola (ed.) (1994) *Guidebook to Cytokines and Their Receptors*, Oxford University Press, Oxford.

## Intermediate Filaments

Intermediate filaments (IFs), which are about 10 nm in diameter, are so named because they were thought to be intermediate in size between the well-known [microtubules](#) and the **actin**-containing [microfilaments](#). They are ubiquitous components of virtually all eukaryotic cells (Fig. 1), being found not only in the **cytoplasm** but also within the cell [nucleus](#). The IF are dynamic structures capable of assembling and disassembling during cell division and of making and breaking interactions with microtubules, microfilaments, and other cellular components.

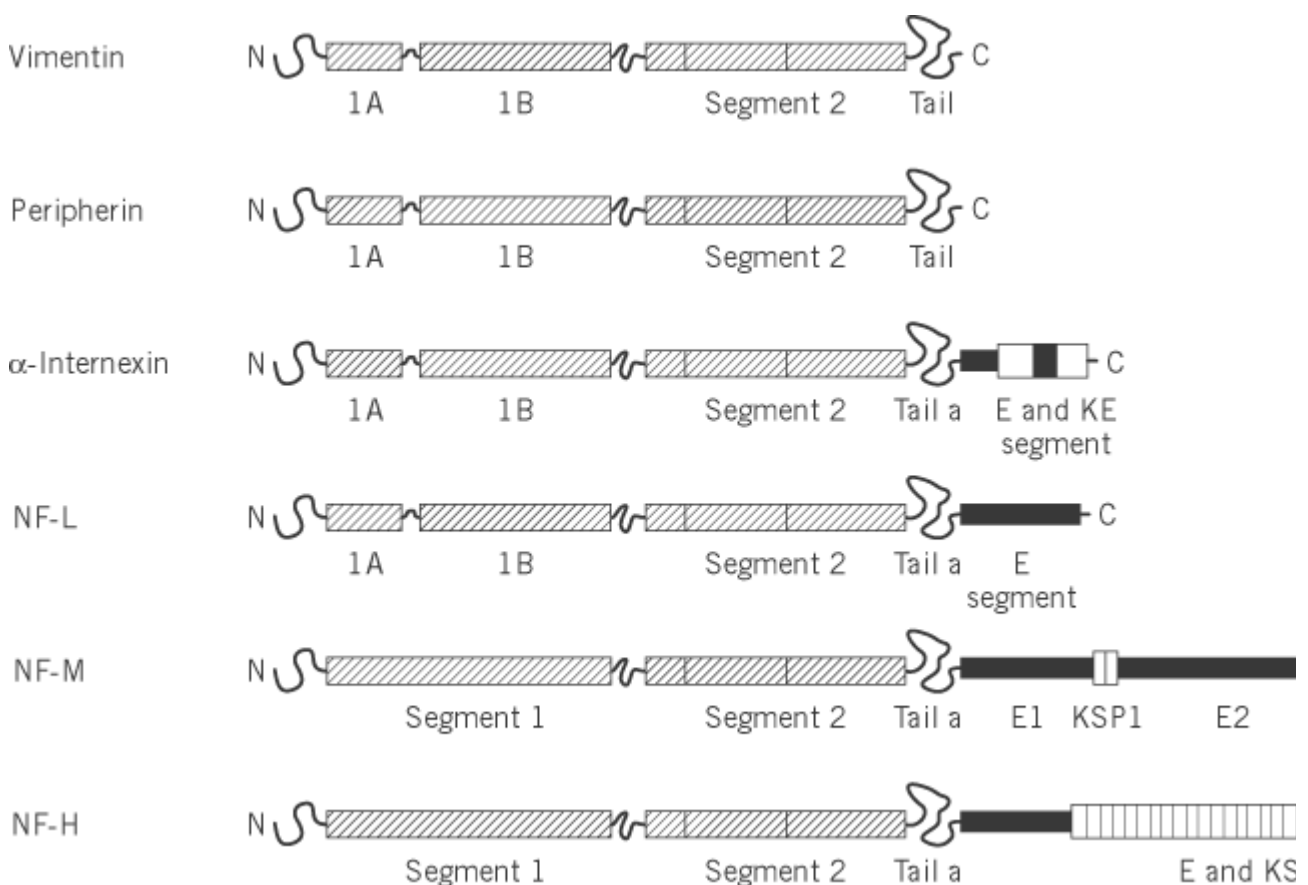
**Figure 1.** Electron micrograph of vimentin IF showing long sinuous filaments about 10 nm in diameter. Bar = 100nm. (From Ref. 4, with permission.)



[Sequence analysis](#) of IF protein chains has shown common and variable **domains**. In essence, each protein has a largely conserved central, **heptad repeat**-containing domain that separates regions at either end of the chain that are chemically and structurally diverse (Fig. 2). These similarities and differences, together with subsequent data on their **gene** structures, have allowed a simple classification system to be formulated. Type Ia and IIa IF chains occur in hard  $\alpha$ -[keratin](#), type Ib and IIb chains occur in epidermal keratin, type III chains occur in **desmin**, **vimentin**, glial fibrillary acidic protein and peripherin, type IV chains occur in  $\alpha$ -internexin and light-, medium- and heavy-neurofilament chains, type V chains occur in the nuclear **lamins**, and type VI chains occur in nestin. The lamin group is unique in their nuclear localization and in the insertion of six additional heptads within the central domain. (Note, however, that invertebrate IF also contain a six-heptad insert in the

same place.) One IF family, the [cytokeratins](#), is discussed in detail in that entry. The keratins form obligate heterodimeric structures (type I/type II) in contrast to the homodimeric conformations of type III IF (and type IV and type V IF too). Nonetheless, heterodimers of the form III/IV and III/VI are also found *in vivo*. The heptad-containing regions thus form a coiled-coil rope of predicted length 45.7 nm in hard  $\alpha$ -keratin, 46.2 nm in epidermal keratins, and 45.3 nm in vimentin, other type III molecules and in type IV molecules too. Type V IF molecules are about 52 nm in length. The values for epidermal keratin and vimentin are determined from an analysis of the disposition of the residues that are cross-linked (1-3). Electron micrographs of rotary shadowed IFs have been obtained that are in agreement with these values. [Scanning transmission electron microscopy](#) (STEM) data indicate that native IFs *in vivo* contain about 32 chains in cross section, although some aggregates with about 16, 24, and 40 chains have also been observed *in vitro*. This has been taken to indicate that the IF may contain an eight-chain protofibrillar substructure. *In vitro* two, three, four, or five such protofibrils may aggregate, but *in vivo* it seems more probable that only the 32-chain (four-protofibril) possibility occurs. The data thus show that the rod domains form the core of the IF and that the *N*- and *C*-terminal domains lie on the surface, where they may interact with other elements of the cell. As noted before, some types of IF in either the cytoplasm or the nucleus are able to assemble and reassemble during certain stages of the [cell cycle](#), most particularly **mitosis**. It now seems clear that **phosphorylation** generally plays a significant role in this process. Almost all IF chains contain phosphorylation sites in their *N*- and *C*-terminal domains. Phosphorylation resulting from the action of a range of possible **kinases** results in the disassembly of preformed IF or prevents formation of IF *in vitro*. Conversely, dephosphorylation allows reassembly to proceed. The complete story is somewhat more complicated than indicated here, but without doubt phosphorylation/dephosphorylation plays a central role in regulating the formation of IF.

**Figure 2.** Domain structures of intermediate filaments of type III (vimentin and peripherin) and type IV ( $\alpha$ -internexin and NF-L, NF-M, NF-H). While the rod domains of all IF proteins, segments 1 (or 1A and 1B) and 2, are **homologous** (though distinct), the *C*-terminal proteins display considerable variation. KSP, for example, refers to the multiple Lys-Ser-Pro motifs that are important in glutamic acid-rich regions. (From Ref. 4, with permission.)



IF function has long been a controversial topic, but real progress has been made in recent years. Consider, for example, just one of the IFs, namely, [neurofilaments](#). The long, highly phosphorylated C-terminal domains in the medium and heavy neurofilament chains (NF-M and NF-H) are critically involved in determining the caliber of the axon, as well as its organization and function. The C-terminal sequences, when phosphorylated, are believed to control slow axonal transport down the microtubule and to have enhanced flexibility that allow them to modulate the density of the neurofilaments, as well as their separations. Overexpression of NF-L and NF-H produce experimental phenotypes that resemble motor neuron disease. Furthermore, mutations in the 52 – Lys–Ser–Pro–phosphorylation sites in the C-terminal domain of NF-H are the proximal cause of some cases of sporadic amyotrophic lateral sclerosis. The C-terminal domains of the NF-M and NF-H chains are known to interact with the C-terminal ends of the  $\alpha$ - and  $\beta$ -[tubulin](#) subunits in [microtubules](#), although [microtubule-associated proteins](#) such as  $\tau$  (**tau protein**) and MAP2 are also believed to be intimately involved in this phosphorylation-regulated interaction. Terminal domains in other IF proteins are also important in specifying their functions.

### Bibliography

1. P. M. Steinert, L. N. Marekov, R. D. B. Fraser, and D. A. D. Parry (1993) Keratin intermediate filament structure: crosslinking studies yield quantitative information on molecular dimensions and mechanism of assembly. *J. Mol. Biol.* **230**, 436–452.
2. P. M. Steinert, L. N. Marekov, and D. A. D. Parry (1993) Conservation of the structure of keratin intermediate filaments: molecular mechanism by which different keratin molecules integrate into pre-existing keratin intermediate filaments during differentiation. *Biochemistry* **32**, 10046–10056.
3. P. M. Steinert, L. N. Marekov, and D. A. D. Parry (1993) Diversity of intermediate filament structure: evidence that the alignment of coiled-coil molecules in vimentin is different from that in keratin intermediate filaments. *J. Biol. Chem.* **268**, 24916–24925.
4. D. A. D. Parry and P. M. Steinert (1990) *Intermediate Filament Structure*, Springer-Verlag, Heidelberg.

### Suggestions for Further Reading

5. D. A. D. Parry and P. M. Steinert, 1995 *Intermediate Filament Structure*, Springer-Verlag, Heidelberg.
6. R. D. Goldman and P. M. Steinert (eds.) (1990) *Cellular and Molecular Biology of Intermediate Filaments*, Plenum Press, New York.
7. E. Fuchs and K. Weber (1994) Intermediate filaments: structure, dynamics, function and diseases. *Annu. Rev. Biochem.* **63**, 345–382.

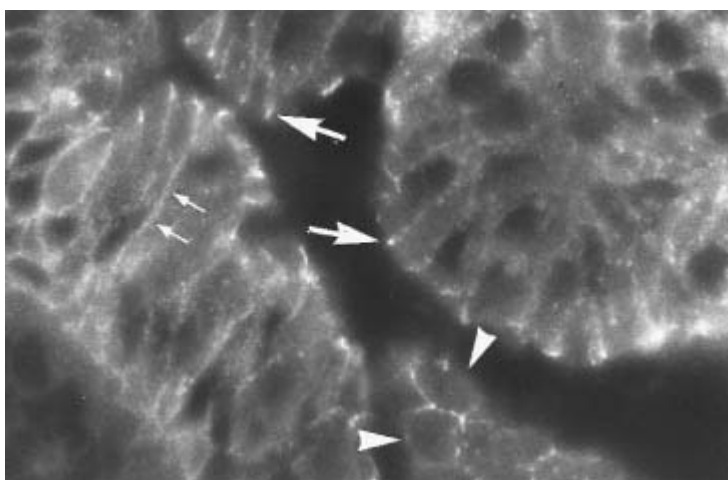
### Intermediate Junction

The intermediate or adherens junction (see [Cell Junctions](#)) is a cell–cell adhesive that occurs in a number of different forms. The category includes the zonula adherens of simple epithelial cells (1), forming a subapical ring beneath the [tight junctions](#), spot-like punctae adherentes (2) of small intestinal mucosa, and the ribbon-like fasciae adherentes of the intercalated discs of cardiac muscle

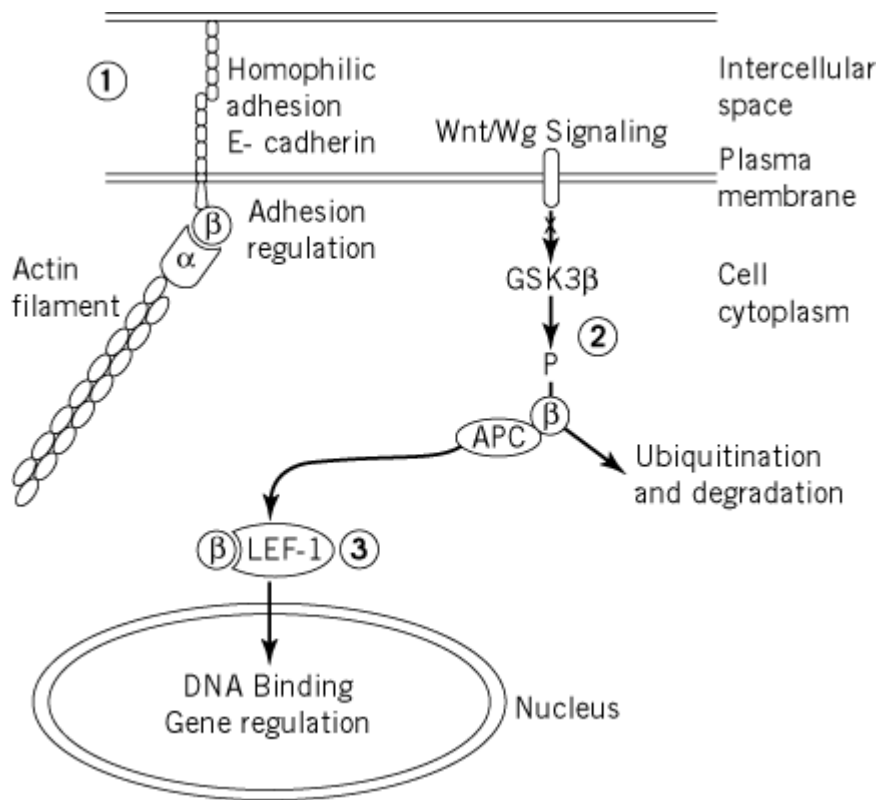
(3). Ultrastructurally, these junctions are characterized by membrane separation of 20 nm and a cytoplasmic plaque or membrane undercoat that is associated with the [actin](#) or microfilament [cytoskeleton](#). The principal functions of adherens junctions are (i) cell adhesion, (ii) organization of the basolateral domain/terminal web of epithelia, (iii) [signal transduction](#), and (iv) contractility.

The principal adhesion molecules of adherens junctions are classical [cadherins](#) such as E-cadherin and *N*-cadherin (Figs. [1](#) and [2](#)) ([4](#)). These mediate homophilic binding between the apposed plasma membranes of adjacent cells. The cytoplasmic domains of cadherins are linked to the actin cytoskeleton via molecules known as **catenins** ([5](#)).  $\beta$ -Catenin is a member of the armadillo family of junctional/signaling molecules ([6](#));  $\alpha$ -catenin is an [actin-binding protein](#) that binds to E-cadherin via  $\beta$ -catenin ([7](#)).  $\alpha$ -catenin is related to vinculin, another actin-associated protein located in intermediate junctions ([8](#)). Alternatively,  $\gamma$ -catenin or plakoglobin may participate in this linkage in place of  $\beta$ -catenin. A large number of other proteins that participate in actin cytoskeletal assembly and function have also been found in association with junctions of the adherens type (see [Table 1](#)). The actin filaments extend toward the membrane with the barbed ends outwards after heavy myosin decoration. Several candidates for actin-capping and -binding proteins are present. These include vinculin, ezrin, radixin and moesin, and  $\alpha$ -actinin ([9](#)). E-cadherin may be found in epithelial cells as a complex with fodrin, adducin, and [ankyrin](#), members of the submembrane cytoskeleton. The latter group of molecules also form complexes with the basolaterally located sodium potassium  $\text{Na}^+/\text{K}^+$  ATPase. These associations between adhesive components of the adherens junction and components of the membrane cytoskeleton of the basolateral epithelial cell membrane are believed to be important in establishing epithelial cell polarity ([10](#)). It is significant that the **tyrosine kinases** c-src and c-yes associate with the junction, since tyrosine **phosphorylation** appears to be an important regulator of the adhesive function of E-cadherin and adherens junctions. In particular, tyrosine phosphorylation of  $\beta$ -catenin had been shown to down-regulate the adhesive function of the junction, resulting in loss of cell adhesion, and to promote invasive behaviour on epithelial cells ([11](#)). Tyrosine phosphorylation of another armadillo protein, p120, also appears to regulate junctional adhesion ([12](#)). Association of the epidermal growth factor receptor (EGF-R) tyrosine kinase with  $\beta$ -catenin has been reported ([13](#)). This emphasizes the role of tyrosine phosphorylation in regulating the function of the junction.

**Figure 1.** Zonulae adherentes of mouse small intestine stained with a [monoclonal antibody](#) to E-cadherin. The zonula form of the junction is evident in transverse sections of cells (arrowheads). In vertical sections, the concentration of stain in the junctions is clearly seen (large arrows), as well as the presence of weaker staining along with the entire lateral borders of the cells (small arrows). (Reproduced from Ref. [23](#), with permission of Chapman and Hall.)



**Figure 2.** Diagram showing the various cellular roles of the adherens junction protein b-catenin. (1) b-Catenin associates with the cytoplasmic domain of the adhesion molecule, E-cadherin where its function is to provide a regulatory link, via a-catenin, to the actin cytoskeleton. This link is essential for the adhesive function of E-cadherin in the intermediate junction. (2) Cytoplasmic b-catenin, when phosphorylated by GSK3b, associates with APC protein and is degraded. (3) Wnt or wingless signaling inhibits GSK3b, allowing unphosphorylated b-catenin to associate with LEF-1 transcription factor to enter the nucleus and regulate gene activity.



**Table 1. Components of the Zonula Adherens Junction in Simple Epithelial Cells<sup>a</sup>**

|                                      | MW ( $M_r$ )<br>( $\times 10^3$ ) | Function   |
|--------------------------------------|-----------------------------------|--|
| <a href="#">E-cadherin</a>           | 120                               | Calcium-dependent cell adhesion molecule.                    |
| <b>b-Catenin</b>                     | 88                                | Mediates E-cad, binding to cytoskeleton. Signaltransduction. |
| <b>Plakoglobin (gamma-catenin)</b>   | 83                                | Mediates E-cad, binding to cytoskeleton. Signaltransduction. |
| <b>Armadillo (<i>Drosophila</i>)</b> | 116                               | Mediates cad binding to cytoskeleton. Signaltransduction.    |
| <b>p120</b>                          | 120                               | Signal transduction.   |
| <b>Vinculin</b>                      | 130                               | Actin-associated protein.                                    |
| <b>a-Catenin</b>                     | 102                               | Links E-cad, to cytoskeleton via b-cat.                      |
| <b>a-Actinin</b>                     | 100                               | Actin binding protein.                                       |

|                          |           |   |
|--------------------------|-----------|---|
| <b>Ezrin</b>             | 85        |   |
| <b>Radixin</b>           | 82        | Membrane-actin cytoskeleton interaction                             |
| <b>Moesin</b>            | 75        |   |
| <b>Fodrin (spectrin)</b> | 240 / 220 | Membrane skeletal protein: binds ankyrin, actin, band 4.1, adducin. |
| <a href="#">Ankyrin</a>  | 200–220   | Links membrane to cytoskeleton                                      |
| <b>Protein 4.1</b>       | 78        | Links spectrin to actin and membrane.                               |
| <b>Adducin</b>           | 103+93    | Mediates actin-spectin binding.                                     |
| <b>Zyxin</b>             | 82        | Actin binding (via a actin).  |
| <b>Tropomyosin</b>       | 35        | Actin binding.  |
| <b>Myosin II</b>         | 260       | Actin binding contractile.  |
| <b>c-src</b>             | 60        | Tyrosine kinase   |
| <b>c-yes</b>             | 60        | Oncogene products   |

<sup>a</sup> The proteins in this table may be found to be associated with cell–cell adhesive junctions of the zonula adherens type but are not restricted to the zonula adherens.

b-Catenin, like other members of the armadillo family, including plakoglobin or g-catenin and the *Drosophila* armadillo protein itself, have a dual function as adherens junction components and in signal transduction and the regulation of gene activity. These proteins participate in the [wingless signaling](#) pathway in *Drosophila* and the Wnt pathway in vertebrates. Thus, both b-catenin and g-catenin have been shown to cause axis duplication when overexpressed in *Xenopus* embryos (14). Furthermore, each may form a complex with the LEF-1 [transcription factor](#) and enter the cell [nucleus](#) to regulate gene expression (15). The adenomatous polyposis coli (APC) **tumor suppressor** protein, another member of the armadillo family, is important in regulating the cytoplasmic level of b-catenin. Thus, b-catenin may associate either with E-cadherin at the cell surface or with APC (16). In the latter complex, it is targeted for attachment to ubiquitin and **protein degradation**. This complex forms when b-catenin is phosphorylated by the **serine-threonine kinase** glycogen kinase synthase kinase-3 b (GSK3b). However, dephosphorylation causes dissociation of the b-catenin-APC complex and increases the cytoplasmic level of b-catenin. This promotes complex formation with T-cell-factor/lymphoid enhancer-binding factor 1 (TCF/LEF-1) and entry into the nucleus. Activation of this signaling pathway may occur either by Wnt signaling or through mutation of APC (17). Mutations of the latter are associated with familial adenomatous polyposis coli, where patients form numerous colonic polyps, an important risk factor in colorectal cancer. The equivalent pathway in *Drosophila*, the wingless signaling pathway, is important in development of segment polarity (18). Mutations of E-cadherin and b-catenin are also important in tumor promotion. Inherited mutations in E-cadherin have been shown to predispose affected individuals to carcinomas (19).

The contractile function of the actin microfilament ring associated with adherens junctions is important in aspects of morphogenesis during embryonic development. In neural tube rolling in the amphibian embryo, for example, contraction of apicolateral microfilament rings in the cuboidal epithelial cells of the neural plate results in curvature of the epithelium by narrowing of the apical surface (20, 21). The adhesive components of the adherens junctions to which the microfilament rings attach bind the cells together. It has also been shown that reagents, such as [cytochalasin](#), that affect actin [microfilament](#) assembly modulate the permeability of tight junctions (See [Tight Junction](#)) in simple epithelia. This may indicate that the zonula adherens plays a role in regulating the function of tight junction (22).



## Bibliography

1. M. G. Farquhar and G. E. Palade (1963) *J. Cell Biol.* **17**, 375–412.
2. D. Drenckhan and H. Franz (1986) *J. Cell Biol.* **102**, 1843–1852.
3. T. Volk and B. Geiger (1984) *EMBO J.* **3**, 2249–2260.
4. K. Boller et al. (1985) *J. Cell Biol.* **100**, 327–332.
5. M. Ozawa et al. (1989) *EMBO J.* **8**, 1711–1717.
6. P. D. McCrea et al. (1991) *Science* **254**, 1359–1361.
7. A. Nagafuchi et al. (1991) *Cell* **65**, 845–857.
8. B. Geiger et al. (1981) *J. Cell Biol.* **91**, 614–628.
9. S. Tsukita et al. (1993) *J. Cell Biol.* **123**, 1049–1053.
10. E. Rodriguez-Boulan and W. J. Nelson (1989) *Science* **245**, 718–725.
11. J. Behrens et al. (1993) *J. Cell Biol.* **120**, 757–766.
12. M. J. Ratcliffe et al. (1997) *J. Biol. Chem.* **272**, 31894–31901.
13. H. Hotschuetzky et al. (1996) *J. Cell Biol.* **127**, 1375–1380.
14. N. Funayama et al. (1995) *J. Cell Biol.* **128**, 959–968.
15. J. Behrens et al. (1996) *Nature* **382**, 638–642.
16. B. Rubinfeld et al. (1997) *Science* **275**, 1790–1792.
17. V. Korinek et al. (1997) *Science* **275**, 1784–1787.
18. M. Van de Wetering et al. (1997) *Cell* **88**, 789–799.
19. P. Guilford et al. (1998) *Nature* **392**, 402–405.
20. P. C. Baker and T. E. Schroeder (1967) *Dev. Biol.* **15**, 432.
21. P. Karfunkel (1971) *Dev. Biol.* **25**, 30.
22. J. L. Madara et al. (1986) *J. Cell Biol.* **102**, 2125–2136.
23. D. R. Garrod and J. E. Collins (1992) "Intercellular Junctions and Cell Adhesion in Epithelial Cells. In" *Epithelial Organization and Development* (T. P. Fleming, ed.), Chapman and Hall, London, pp. 1–52.

## Suggestions for Further Reading

24. D. R. Garrod and J. E. Collins (1992) "Intercellular Junctions and Cell Adhesion in Epithelial Cells. In" *Epithelial Organisation and Development* (T. P. Fleming, ed.), Chapman and Hall, London, pp. 1–52.
25. A. Ben-Ze 'ev (1999) "Focal Adhesion and Adherens Junctions: Their Role in Tumourgenesis". In *Adhesive Interactions of Cells* (D. R. Garrod, M. A. J. Chidgey, and A. J. North, eds.), JAI Press. Greenwich, CT, pp. 135–163.
26. A. I. M. Barth et al. (1997) Cadherins, catenins and APC protein: interplay between cytoskeletal complexes and signalling pathways. *Curr. Opin. Cell Biol.* **9**, 683–690.

## Internal Guide Sequence

In the group I [self-splicing introns](#), the 5' [splice site](#) is determined by the formation of intramolecular base pairing ([1](#)). The site is located in the 5' half of the stem in a stem-loop structure known as P1. A short, **intron**-encoded sequence in the 3' half of the P1 stem is called the internal guide sequence (or

IGS) and pairs with the end of the 5' exon (see "[Splice Sites](#)"). The name suggests that it binds to both splicing junctions to bring the two exon ends together, but this original idea has been proven to be unnecessary. The functionally similar element in group II introns is the exon-binding site (or EBS). Intron-encoded guide sequences have also been observed at the [RNA Editing](#) sites of glutamate receptors by the double-stranded RNA adenosine deaminase.

In connection with the [RNA world](#) hypothesis, a resemblance between the **codon-anticodon** interaction and the IGS-exon pairing has been pointed out (2). It has been reported that an **fMet**-oligo RNA in which the RNA sequence is complementary to an artificial IGS created in a group I intron-derived **ribozyme** can be hydrolyzed at the acylation bond, mimicking a prototype of **peptidyl transfer**.

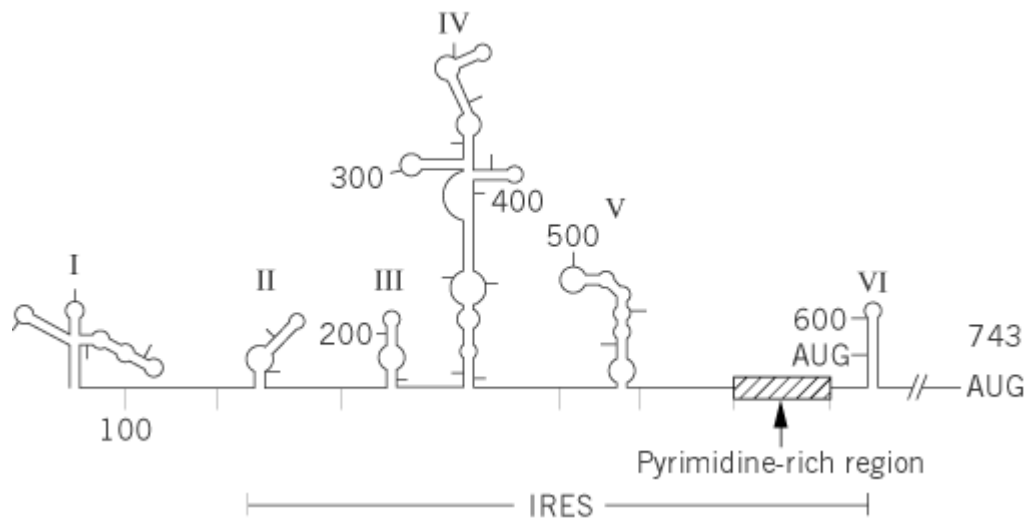
#### Bibliography

1. T. R. Cech (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 239–269.
2. H. F. Noller (1991) *Nature* **353**, 302–303.

#### Internal Ribosome Entry Site (IRES)

Most eukaryotic [messenger RNAs](#) have the Cap structure at 5' termini (see [Cap](#)), where 40 S **ribosomal** subunits enter the mRNA with association with the eIF-4F complex (cap-binding protein complex; see [Initiation Factor \(IF\)](#)). The 40S initiation complex then scans the 5' untranslated region (or 5'UTR) to find the initiation codon (see [Scanning Hypothesis](#)). The mRNAs of picornaviruses, represented by [poliovirus](#), however, have no 5'-Cap structure, but contain an internal ribosome entry site (IRES) on the 5'-untranslated region (UTR), which functions as an alternative entry site for the ribosome (1). Upon infection, poliovirus produces a proteinase that cleaves eIF-4G, the largest subunit of eIF-4F, thereby shutting down the Cap-dependent initiation of cellular mRNAs, while activating the viral translation by IRES-dependent initiation. Picornavirus mRNA has a long 5'UTR (600 to 1200 nucleotides), and the IRES corresponds to the 3' proximal 450-nucleotide region that involves multiple stem-loop structures (Fig. 1). One motif that may act as the actual ribosome binding site is a pyrimidine-rich region of about 20 nucleotides preceding an AUG triplet with a spacer of 20 to 25 nucleotides. In some types of picornavirus, including poliovirus, ribosomes do not initiate translation at this AUG triplet but at the next available AUG. On the other hand, the other picornaviruses, such as encephalomyocarditis virus, initiate [translation](#) at the first AUG. The IRES-dependent initiation requires a subset of ordinal initiation factors, as well as specific cellular factors such as La protein and polypyrimidine-tract binding protein. Viruses other than picornavirus, such as hepatitis C virus (HCV) and plant RNA viruses, also contain IRES. In addition, at least two cellular mRNAs, **antennapedia** mRNA of *Drosophila* and **BiP** mRNA of mammals, show IRES activities. No sequence resemblance has been noted among the IRES of picornavirus, HCV, and the cellular mRNAs.

**Figure 1.** Schematic representation of the predicted stem-loop structure of the poliovirus IRES (1).



### Bibliography

1. E. Ehrenfeld (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 549–573.

### Suggestion for Further Reading

2. R. J. Jackson (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews and N. Sonenberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 71–112.

## Interspersed DNA Elements

To understand the architecture of the human [genome](#), we need a complete definition of all *interspersed DNA elements*, the families of sequences repeated throughout the genome, as these make up the majority of human DNA. These repeat sequences are arbitrarily classified into [CpG Islands](#), [Sines](#) (short interspersed DNA sequence elements comprising [Alu sequences](#) in primates), [Lines](#) (long interspersed DNA sequence elements), **minisatellites** and **microsatellites**. All these entities are repeated many times in the genome, and are arbitrarily defined as about  $10^4$  copies or more per **diploid** genome.

## Introns, Exons

**Introns** are sections of primary RNA **transcripts** that are removed as part of the maturation of functional **messenger** RNA or, rarely, of ribosomal RNA. They are also known as *intervening sequences*. The sections between the introns that are retained in the mature RNA product are called *exons*. Introns occur very widely in the nuclear genes of **eukaryotes**, and nuclear introns command

the most attention. But there are two other, rather different kinds of intron, called **Group I** and **Group II**, that have been found in eukaryotic **organellar** genomes and occasionally in **bacteria** and **bacteriophage**. Introns of another kind have been found in a few species of **Archaeobacteria** (1), but these are not covered in this article.

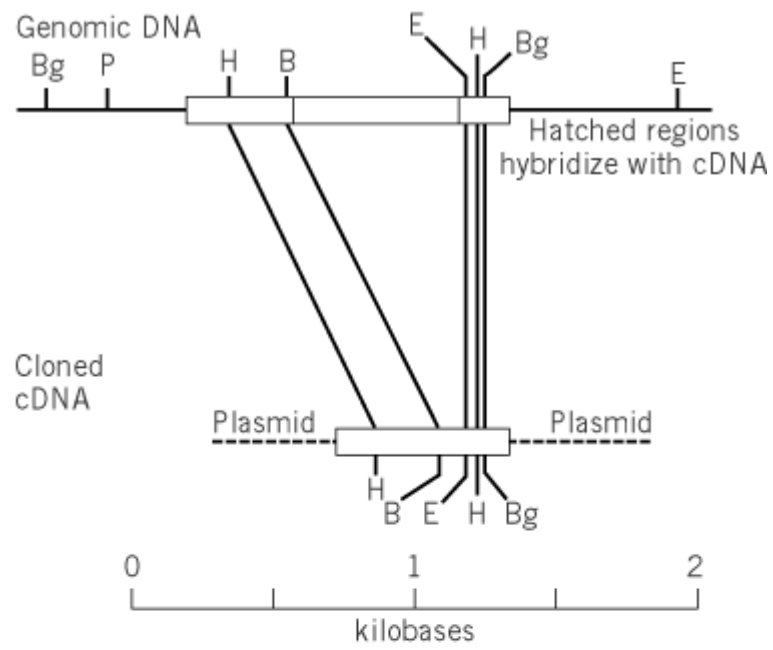
## 1. Eukaryotic Nuclear Introns

### 1.1. Occurrence and Detection

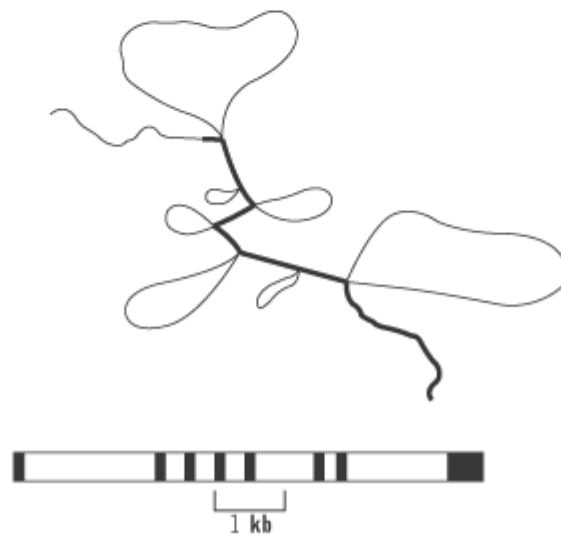
Introns are found in some **genes** of all eukaryotes thus far studied, but their frequency and average size increase enormously as **genomes** become larger and organisms more complex. In the budding **yeast** *Saccharomyces cerevisiae*, most genes do not have introns, and, where they occur, the introns are seldom more than 100 bp long. In filamentous **fungi**, such as *Neurospora* and *Aspergillus*, most genes that have been sequenced so far have introns but seldom more than two or three and usually in the size range 50 to 100 bp. The nematode worm *Caenorhabditis elegans* has some much longer introns, but the transcribed sequences still consist mostly of exons. In the fruit fly *Drosophila* many genes are interrupted by long introns, sometimes kilobases long. In mammals most genes are often expanded 100-fold or more by numerous introns that far exceed the exons in average length, though their coding sequences may not amount to more than a few kilobases. Some examples of exon/intron proportions and arrangements from a number of representative eukaryotes are presented in the article on Gene Structure.

Introns are defined by comparisons between the sequences, or **restriction maps** of chromosomal (genomic) DNA and **complementary DNA** (cDNA) obtained by **reverse transcription** of processed mRNA. cDNA is also used as a probe to search for **homology** with ordered genomic restriction fragments. Some of the first discovered introns were made visible by electron microscopy of artificially constructed genomic DNA/cDNA or genomic DNA/mRNA heteroduplexes. The introns appeared as unpaired single-stranded loops (Fig. 1). Shorter introns are identified, at least tentatively, as segments interrupting **open reading frames** with intron consensus sequences at their ends (Table 1).

**Figure 1.** Two ways of detecting introns (a) Comparison of restriction-site maps of the rabbit b-globin gene and cDNA, reverse transcribed from the b-globin messenger RNA (2). A second, smaller intron was not revealed by this analysis and is not shown in the diagram. (b) The appearance of a heteroduplex formed between chicken ovalbumin gene DNA and its mRNA, drawn from the electron micrograph of Dugaiczky *et al.* (3).



(a)



(b)

**Table 1. Intron Consensus Sequences in Various Organisms** <sup>a, b</sup>

| Organism      | 5' Splice | Branch Point | 3' Splice        |
|---------------|-----------|--------------|------------------|
| Vertebrates   | AG/GUAAGU | UNCURAC      | (Y-rich)NCAG/G   |
| Saccharomyces | AG/GUAUGU | UACU AAC     | (Y-rich)YAG/     |
| Plants        | AG/GUAAGU | UA-rich      | (U-rich)UGCAG/GU |

<sup>a</sup> From Ref. 4.

<sup>b</sup> Key: Y = pyrimidine; R = purine; N = any nucleotide. underlined residues are invariant, or nearly so. The slash (/) divides exon from intron.

Although the internal sequences of nuclear introns bear no relationship to each other (or, indeed, to anything else), they show some consensus at or near to their 5' and 3' (i.e., upstream and downstream) termini, sometimes called the donor and receptor sites (Table 1). At the upstream end, the consensus is 5'-GUAUGU for *Saccharomyces* and 5'-GUPuAGU in mammalian introns. The GU is nearly universal, and the following bases are somewhat variable (U in the messenger RNA corresponds to T in the DNA coding strand). Introns are terminated virtually universally at their downstream (3') ends by YAG (Y for pyrimidine), usually immediately preceded by a pyrimidine-rich sequence. Some 15 to 40 bp upstream of the 3' AG is a sequence called (for reasons that will emerge) the *branch point*, which in yeast has the invariant sequence 5'-UACUAACA, reduced in other fungi to CURAC, and represented in mammals by the barely recognizable YXYRAY (X standing for any base and R for either purine). The conserved A is most important. The somewhat different consensus sequences in different groups of organisms are shown in Table 1).

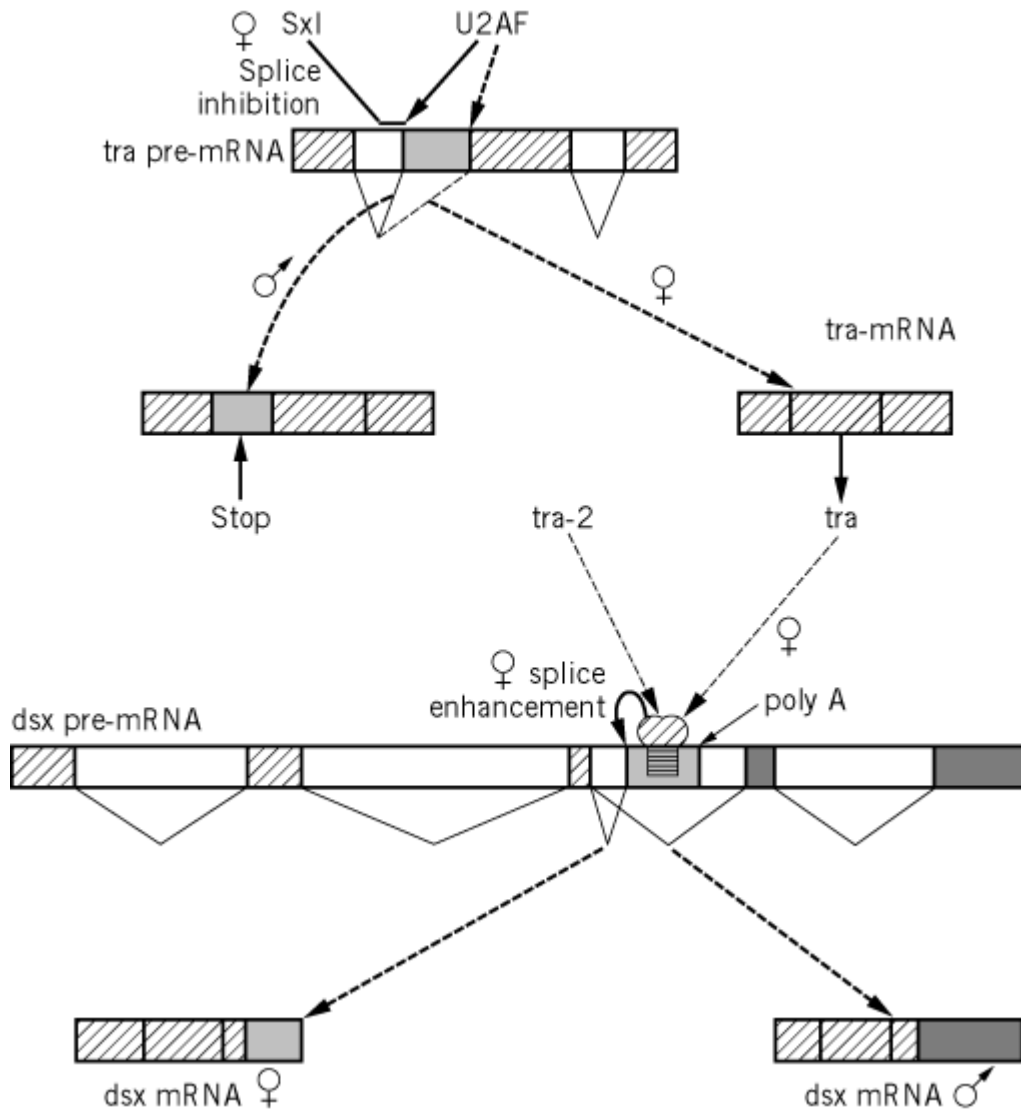
Intron excision occurs in two steps. First, the upstream guanosine 5'-phosphate is cleaved from the upstream exon and condensed with the 2'-hydroxyl group of the branch-point adenosine to form a loop, often called a “**lariat**” structure, still joined to the downstream exon (see Fig. 4 later). Then the freed 3'-hydroxyl of the upstream exon attacks the phosphate diester link between the lariat and the downstream exon to free the lariat and join the two exons. This sequence of interactions is catalyzed by an array of five **small nuclear ribonucleoprotein** complexes (snRNPs), U1, U2, U4, U5, and U6, each of which contains a short unique RNA molecule and a number of proteins, some of which, the Sm proteins, are RNA-binding. Collectively the U complexes form a supercomplex, called the **spliceosome**.

To a considerable extent, the process of splicing can be explained in terms of **Watson–Crick base pairing** between RNA sequences, shown by suppression of the effects of **mutation** in splicing sequences by compensating changes in U-RNAs or by chemical cross-linking of the paired sequences. Just as with **codon–anticodon** pairings in messenger RNA **translation**, most of the RNA–RNA pairings must be too weak to hold by themselves, and presumably they act by guiding much stronger protein–RNA and protein–protein interactions, about which much remains to be learned. The course of events has been most thoroughly worked out for *Saccharomyces*. With only minor variations, it is the same in mammals and presumably in eukaryotes generally. (see **Splicing RNA**).

## 1.2. Splicing Options

Some intron excision is regulated in all eukaryotes that have introns in tissue or cell-specific ways. Thus, in *Saccharomyces* yeast, which has relatively few introns, the splicing-out of the intron in the *MER1* gene, which functions in meiotic **recombination**, depends on the protein product of the **meiosis**-specific *MER2* (5). In *Drosophila*, one of the three introns in the **transposase** gene of the mobile **P-element** is spliced out only in the **germ line**, which explains why the element does not transpose in somatic tissues (6). Sexual development in the fly depends on a splicing cascade. The *Sexless* (*Sxl*) gene product, produced only in females, promotes both the splicing of its own pre-mRNA and that of the *transformer* (*tra*) gene, so that the latter encodes a protein which, acting together with the product of another gene *tra2*, splices the sexless (*sxl*) pre-mRNA to produce mRNA for female development (Fig. 2). The alternative (“default”) splicing option is female-determining (8). In mammals that have their abundant introns, the splicing of pre-mRNAs in different ways in different tissues is almost more the rule than the exception. Two other examples of alternative patterns of pre-mRNA splicing are illustrated under **Gene Structure** and **Interallelic Complementation**.

**Figure 2.** Two examples of splicing option control from *Drosophila*. (a) The female-specific *Sxl* (*sex-lethal*) protein binds at the 3' end of intron 1, blocking the otherwise favored male splicing option and diverting splicing action of U2AF to a splice site in exon 2, bypassing a stop codon and allowing formation of functional female-specific mRNA. See text for further explanation. (b) *tra-2* and female-specific *tra* proteins bind to a splicing-enhancer region (shown horizontally hatched) in exon 4 of the *dsx* (doublesex) gene to activate splicing at the 3' end of intron 3, thus preempting the otherwise favored alternative of splicing exon 3 directly to exon 5. The two modes of splicing result in two functional mRNAs, one for males and one for females. There is a polyadenylation signal at the 3' end of exon 4, thus effectively terminating the female mRNA at this point. Information from Ref. 7.



So far as it is understood, selective splicing of introns depends mostly on splicing enhancers, RNA motifs to which proteins bind to facilitate splicing at nearby sites. A family of RNA-binding proteins, the so-called SR proteins, have been especially implicated in enhancer function (7, 9). Characteristically, they have two kinds of domains, one RNA-binding and the other rich in arginine-serine dipeptide sequences (hence the designation SR). The latter, it is thought, are involved in protein-protein binding and recruiting of components of the splicing apparatus.

One of the most thoroughly investigated systems of splicing control by an enhancer is the final step of the *Drosophila* sex-determining splicing cascade, alluded to previously (10) (Fig. 2). The enhancer sequence that determines male-type splicing of *Drosophila* *dsx* mRNA and interacts with the protein products of *tra* and *tra2* is located in the fourth exon between 300 and 570 base pairs

downstream of the splice junction that it controls. It includes six tandemly arranged approximate tandem repeats, each 13 bp, with an 18-bp purine-rich sequence between the fifth and sixth repeats. Both of the repeats and the purine-rich element are necessary for efficient splicing of the preceding intron. *In vitro* protein–RNA binding studies have shown that the *tra2* protein binds to the repeats and the *tra* protein to the purine-rich segment (10). Both are proteins that have SR domains, and their binding to the RNA is reinforced by cooperative interaction with one or more other SR proteins, which would not bind specifically to the enhancer by themselves. It is thought that the whole SR complex so formed binds to splicing factors, perhaps U2 components, and relays them to the 3' splice site located 300 to 320 bp upstream (Fig. 2). A problem to which evolution has presumably found an answer is that of reconciling the two quite different functions of the *dsx* splice enhancer. On the one hand, it has to provide binding sites for the splice-promoting proteins, and on the other hand, it encodes a part of the amino acid sequence of the female-specific *dsx* protein that has its own special functions.

Purine-rich splice-enhancer sequences are found in other genes, including some from mammals (11). Binding to SR proteins, either site-specific or more general in function, is likely to be a general property of splice enhancers.

A different method of splice-site regulation is exemplified by the first step of the *Drosophila* sex cascade (12). The *Sxl* protein, which acts in the female on *tra* pre-mRNA to produce the female-specific messenger, is a kind of degenerate SR protein that retains its RNA-binding domain but has lost its SR domains. It binds with high affinity and specificity to the polypyrimidine tract adjacent to the 3' splice site of the third *tra* intron to the exclusion of a complete SR protein, called U2AF, which is necessary for recruiting the U2 complex, a function that SR-less *Sxl* protein is unable to perform. As a result, the third exon becomes spliced to the fifth instead of the fourth exon ((12); Fig. 2).

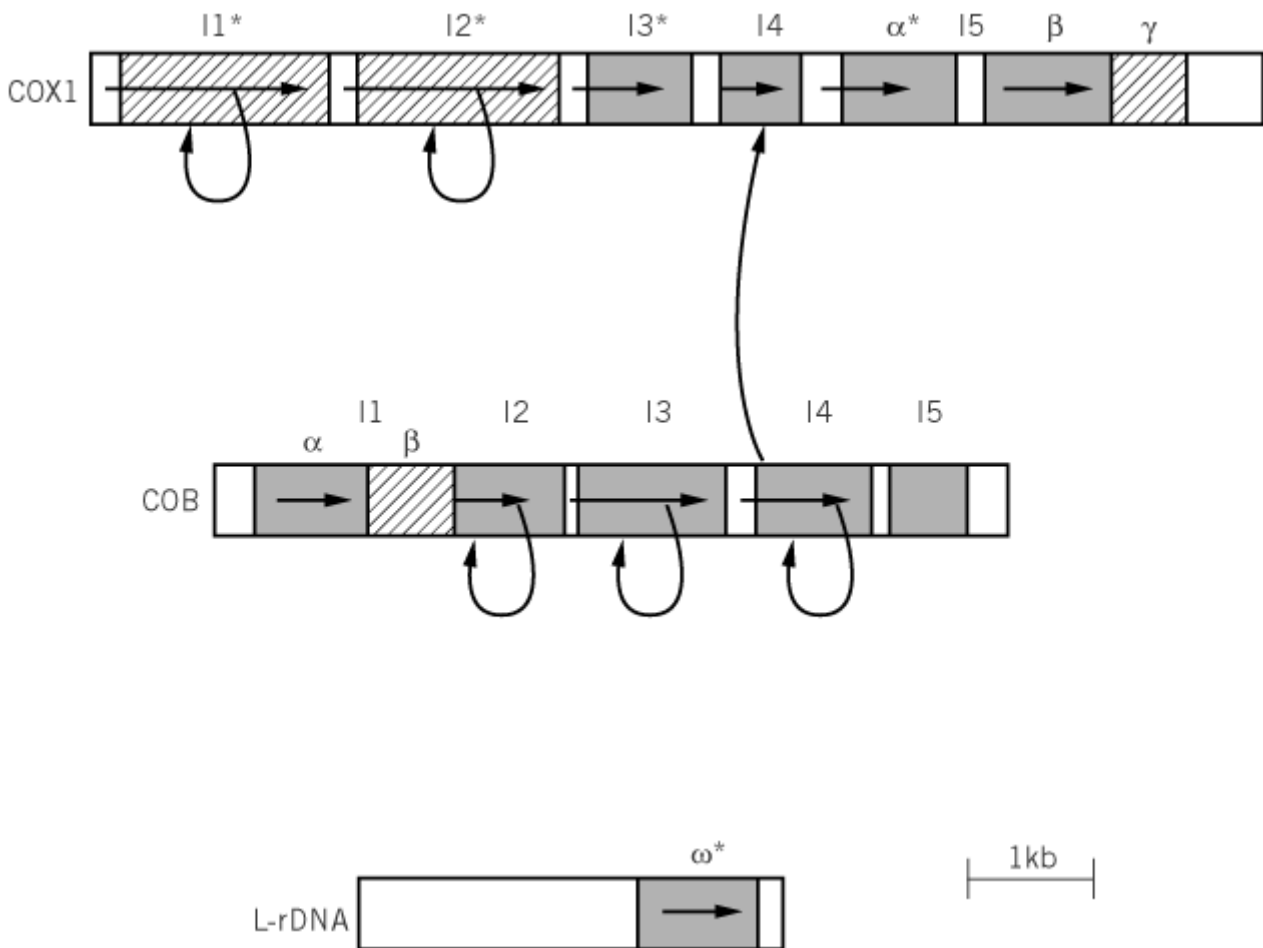
In both of the *Drosophila* examples just reviewed, the splicing control is exercised at the 3' end of the intron. Other splicing options involve choices among different 5' splice sites. It is presumed that, here too, the options are controlled by specific RNA-binding proteins, particularly of the SR class, that have either enhancing or inhibitory effects on specific splice sites. A very great deal remains to be discovered in this field.

## 2. Group I and Group II Introns

These introns are most commonly found within the genes of [mitochondria](#) and [chloroplasts](#) and more rarely in **bacteria** and **bacteriophages** (Table 2). **Group I introns** are particularly prevalent in mitochondrial DNA of fungi, also accompanied at least in *Saccharomyces* by group II (Fig. 3). Both groups are common in **plastid** DNA of many kinds of photosynthetic eukaryotes. Group I introns are found particularly in [transfer RNA](#) genes. With the exceptions shown in the Table 2, they are not found in nuclear genes.

**Figure 3.** Group I and group II introns in mitochondrial genes of *Saccharomyces cerevisiae*. *COB* and *COX1* encode cytochrome b and subunit 1 of cytochrome oxidase. The L-rDNA transcript is processed to form large ribosomal RNA. Open, stippled, and hatched boxes represent, respectively exons, group I introns, and group II introns. Straight arrows indicate intron open reading frames. Note that several of them start in upstream exons. Curved arrows indicate that the protein products facilitate the splicing of their own introns or, in one case, an intron in another gene. The asterisks indicate high-frequency transfer to intronless alleles. This diagram is a composite. These introns are not all present together in any one yeast strain. Information from Ref. 1.





**Table 2. Distribution of Group I and Group II Introns<sup>a</sup>**

|                      | <b>Group I</b>  | <b>Group II</b>                     |
|----------------------|---|-------------------------------------|
| <b>Mitochondrial</b> |   |                                     |
| Yeasts               | <i>Saccharomyces</i><br><i>Schizosaccharomyces</i><br><i>Klyveromyces</i> | <i>Saccharomyces</i>                |
| Sea anemone          | <i>Metridium</i>  |                                     |
| Filamentous fungi    | <i>Neurospora</i><br><br><i>Podospora</i><br><i>Aspergillus</i>           |                                     |
| <b>Bacteria</b>      | <i>Agrobacterium</i>  | <i>Calothrix</i> <i>Escherichia</i> |
| Blue-green algae     | <i>Anabaena</i><br><i>Cyanophora</i>                                      |                                     |
| Bacteriophage        | T4 ( <i>E. coli</i> )<br>SPO1 ( <i>Bacillus subtilis</i> )                |                                     |

## Chloroplast

|                  |  |  |
|------------------|--|--|
| Flowering plants | <i>Nicotiana</i>                         | <i>Nicotiana</i> (tobacco)               |
|                  | <i>Zea</i> (corn)                        |  |
|                  | <i>Vicia</i> (bean)                      |  |
| Liverwort        | <i>Marchantia</i>                        | <i>Marchantia</i>                        |
| Algae            | <i>Dictyota</i> (brown seaweed)          |  |
|                  | <i>Vaucheria</i>                         |  |
|                  |  | <i>Bryopsis</i><br>(green, filamentous)  |
|                  | <i>Ochromonas</i> (yellow-green unicell) |  |
|                  | <i>Chlamydomonas</i>                     | <i>Chlamydomonas</i><br>(green unicells) |
|                  | <i>Chlorella</i>                         |  |
|                  |  | <i>Euglena</i>                           |

## Nuclear

|            |                    |
|------------|--------------------|
| Slime mold | <i>Physarum</i>    |
| Ciliate    | <i>Tetrahymena</i> |

---

<sup>a</sup> [1](#), [13](#), and references therein.

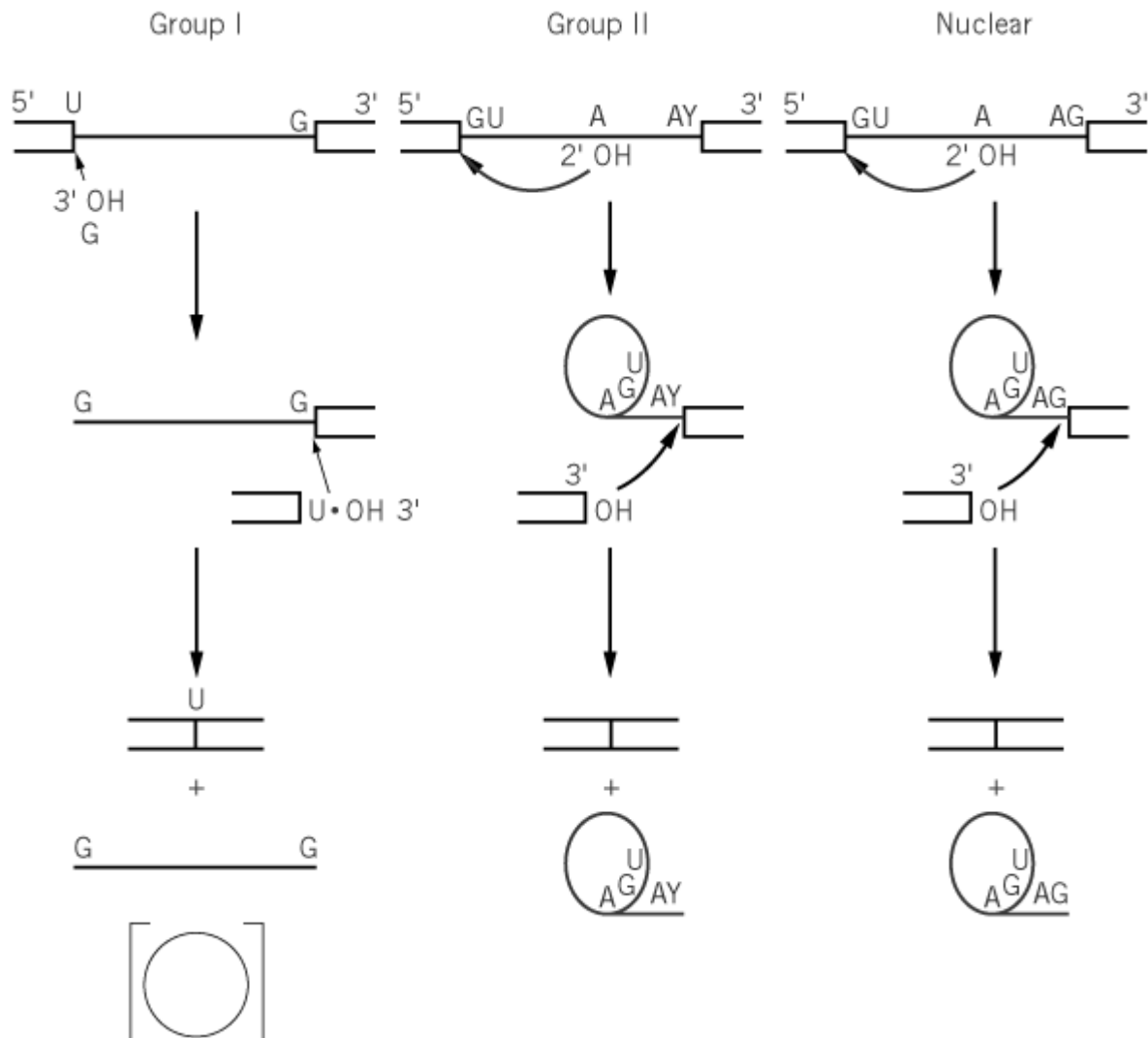
Group I and group II introns differ fundamentally from the common introns of nuclear pre-mRNAs in their requirements for splicing. They have characteristic internal secondary structures that differ between the two groups, in the form of double-stranded stems, stabilized by base-pairings, and single-stranded loops. This internal structure, it is thought, serves the same purpose in aligning the intron ends for splicing as does base pairing with the RNAs of the snRNPs in splicing nuclear introns.

At least some group I and (to a lesser extent) group II introns act as **ribozymes** to catalyze their own splicing *in vitro* without the aid of protein, though mostly under unphysiological conditions, such as high temperature and high magnesium ion and general salt concentration. *In vivo*, their splicing is facilitated by proteins, generally called *maturases*, some of which are encoded in open reading frames within the introns themselves ([1](#), [13](#), [14](#)). This has been shown genetically for some of the *Saccharomyces* introns (Fig. [3](#)), where mutations in the open reading frames result in failure to splice, a deficiency that can be remedied by the same kind of nonmutant introns acting in *trans*. For example, in the *Saccharomyces cerevisiae* *cob* and *cox1* genes, that respectively encode cytochrome b and subunit 1 of cytochrome oxidase, numerous introns, some group I and some group II, depend mostly for their splicing on intron-encoded proteins, and some require nuclear-encoded proteins, either instead or in addition.

Groups I and II are distinguished from one another in their secondary structures, in the nature of the proteins encoded by their open reading frames (see later), and in the chemistry of their splicing reactions (Fig. [4](#)). In Group I, the phosphodiester linkage at the 3' end of the upstream exon is attacked by the 3'-hydroxyl of a free guanosine molecule. A trans-esterification reaction leaves the junction cleaved with a free 3'-hydroxyl on the exon side and the guanosine attached to the free 5' end of the intron. Then, in a second transesterification step, the 3'-hydroxyl of the upstream exon reacts with the phosphodiester bond at the 3' end of the intron (always a guanylate residue), releasing the intron, sometimes but not always as a circle, and ligating the exons together (Fig. [4](#)). In one case,

a group I intron in the fungus *Podospora anserina*, the excised circles reverse-transcribed into DNA, has embarked on a damaging career of its own, replicating freely in the mitochondria and causing senescence (1).

**Figure 4.** The modes of excision of group I, group II, and nuclear introns compared. Based on Ref. 1 by permission.



The splicing of group II introns is similar to that of nuclear pre-mRNA introns (Fig. 4), and a lariat intermediate is formed. G is always at the 5' end of the intron (consensus GUGYG), and generally AY (not AG) is at the 3' end. A crucial A also serves as a branch-point residue, although the sequence surrounding it, usually composed of U and/or C residues, has little similarity to the corresponding sequences of nuclear introns. The 5' G interacts with the branch-point A to form the lariat and then, again as in nuclear introns, the free 3'-hydroxyl of the upstream exon attacks the phosphodiester linkage at the 3' end of the intron to fuse the exons and release the intron as a lariat.

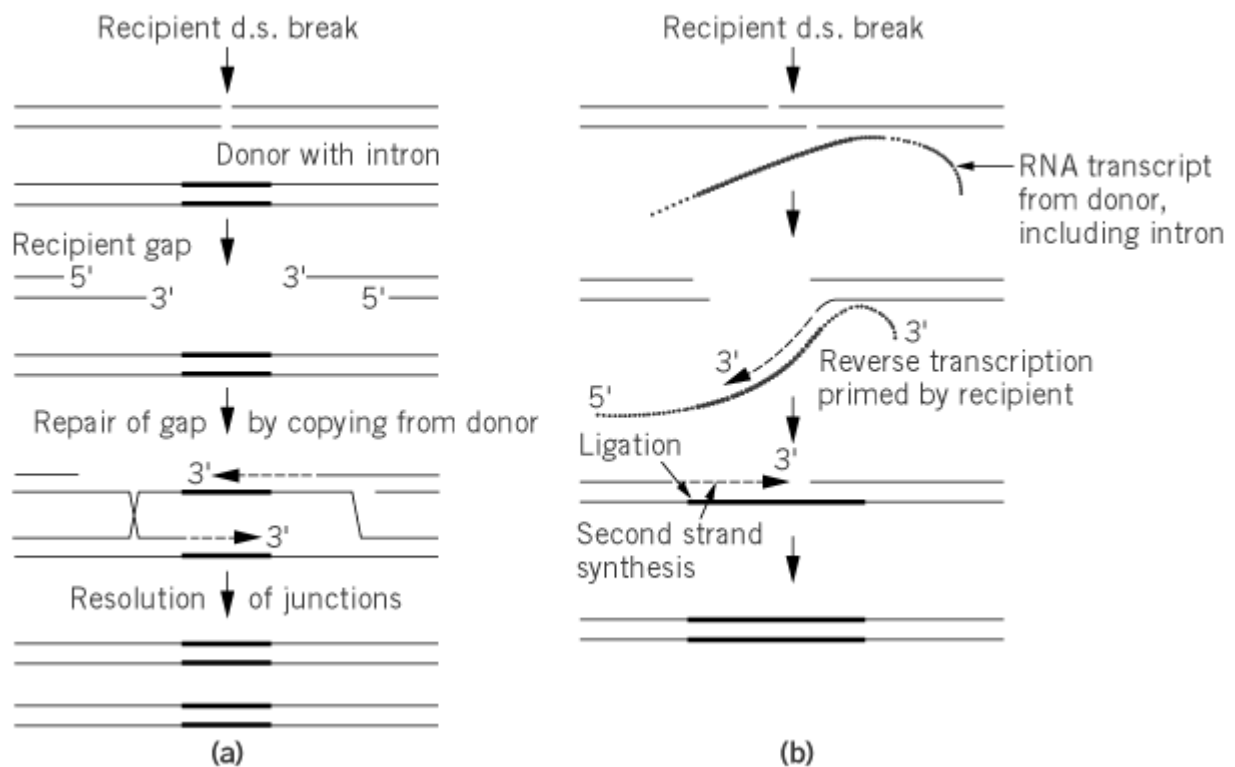
An extraordinary feature of group I introns, found also but less frequently in group II, is their ability to transpose as DNA from one allele to another. The transposition is conservative in the sense that the recipient acquires the intron without the donor losing it. One of the most thoroughly investigated examples involves the group I w-intron in the mitochondrial large rRNA gene of some strains of *Saccharomyces cerevisiae* (14). When a strain possessing the intron (w+) is crossed with one that

lacks it ( $w^-$ ) (the mitochondria of the two strains are distinguished by drug resistance or other markers), the progeny emerging from the cross have recombined as well as predominantly parental-type mitochondrial genomes, virtually all now carrying  $w$ . And although most of the mitochondrial genetic markers are inherited from either parent with equal probability, two drug-resistance markers, [chloramphenicol](#) and **erythromycin** resistance, located within the rRNA gene and hence close to  $w$ , are inherited predominantly from the  $w^+$  parent. Apparently  $w$  is unilaterally transferred from one mitochondrial genome to the other and carries flanking DNA of variable length along with it. The explanation for this conversion (and coconversion) is that  $w$  includes an open reading frame that encodes a double-strand **endonuclease** that is specific for a site within the rRNA gene. This enzyme makes a double-strand break in the mitochondrial DNA of the  $w^-$  parent, which is then enlarged to a gap that is repaired by copying from the  $w^+$  parent ([1](#), [14](#)). This explanation bears a strong resemblance to the double-strand break model currently proposed for [recombination](#) in meiosis.

The high-frequency transposition of  $w$  is strictly homologous, that is to say, it is targeted to a sequence similar to the that from which it came, a process sometimes called “homing.” The same property has been identified in a number of other group I introns (Fig. [3](#)), and the same explanation applies to these, namely, the intron encodes an endonuclease that makes a double-stranded cut at or close to the intron's home site. As indicated before, at least some of the group I-encoded endonucleases also act as maturases to facilitate splicing.

The open reading frames in the approximately 20% of group II introns that possess them encode proteins of quite a different family. Virtually all have a consensus domain, called X, a possible but unproved function of which is maturase activity ([15](#)). Many of them also have two other recognizable domains, one characteristic of **reverse transcriptase** and the other **zinc-finger**-like with double-strand endonuclease activity specific for the homing site ([16](#)). The transposition of the yeast *COXI* intron, which has been analyzed *in vitro*, occurs by reverse transcription, primed by the 3' OH-terminated transcribed DNA strand on one side of the double-strand break, followed by ligation of the 3' end of the single-strand, reverse transcript into the recipient DNA and the repair synthesis of its complementary strand (Fig. [5](#)).

**Figure 5.** Proposed mechanisms of “homing” transposition of group I (a) and group II (b) introns. Dotted lines are RNA, broken lines new synthesis, and thick segments intron sequence. After Refs. [1](#) and [14](#) by permission.



An interesting feature of many of the open reading frames in both class I and class II introns is that, although they are mainly within the intron, they start in the upstream exon (Fig. 3). This mutual adaptation of exon and intron must mean that their association is very ancient and that the intron cannot be a recent invader. This conclusion is borne out by observations of positional conservation of some group I introns over great lengths of evolutionary time. However, their patterns of distribution that have similar kinds of introns present in diverse genes and homologous genes with different sets of introns are hard to explain except on the assumption that both group I and group II introns can move, even if rarely, to new sites. In some instances, such movement has been directly demonstrated. To further complicate the picture, there is evidence that open reading frames within introns can sometimes transpose from one intron to another (17). Again, the most plausible, though unsubstantiated, explanation for such events is that transcripts can be reverse-transcribed and the cDNAs reintegrated into the genome, usually into sequences homologous to those from which they came from but occasionally at nonhomologous, that is, ectopic, sites.

### 3. Uses for Introns

From the elaborate mechanisms that organisms have evolved to cope with introns, one would suppose that they must confer some benefit. There are several ways in which they can contribute to the function of the organism.

First, introns in nuclear genes give additional flexibility to gene function. Splicing options allow a gene to encode two or more different proteins to serve different functions at different times. Some sequences are in common and others are unique. Some examples were given before and others can be found in the articles on [Gene Structure](#) and [Allelic Complementation](#). Investigating the control of pre-mRNA splicing is still at an early stage, but it promises to reveal a degree of complexity to rival that existing in control of gene [transcription](#).

Secondly, introns can be involved in controlling transcription. There are fairly numerous examples of downstream enhancers located in introns of genes. One, the *bithorax* enhancer within the *Drosophila Ubx* gene, is mentioned under [Gene Structure](#).

Thirdly, there are several examples, which may become much more numerous as more of the very long introns in mammalian genes are sequenced, of open reading frames that encode significant looking protein sequences, “nested” within the introns of other genes. An example is cited in the article on [Gene Structure](#). The *Drosophila Gart* gene, identified because of its function in **purine biosynthesis**, has an open reading frame in its largest intron that encodes cuticle protein, used at the pupation stage. Here the mRNAs of *Gart* and the nested gene are transcribed from opposite strands of the DNA, and therefore, presumably, cannot, both be produced at the same time. There are also a few examples of genes nested within mitochondrial introns. The sea anemone *Metridium senile* harbors a group I intron within its mitochondrial gene for subunit 5 of NADH dehydrogenase, and this intron includes the genes for subunits 1 and 3 of the same enzyme (18). In this case all three genes use the same DNA strand for transcription, and so the excised intron becomes a bicistronic mRNA for the nested genes.

A less well-substantiated, but very real possibility is that introns are sometimes used as timing devices for gene expression. Transcription is not instantaneous. In *Drosophila* it has been estimated that it proceeds at the rate of about 1.4 kb per minute. This implies that, for a gene expanded to 60 kb because of its intron content, as is sometimes the case in *Drosophila* and often in mammals, there is an interval of nearly an hour between the start and completion of transcription. There is one example in *Drosophila* where intron length is crucial (19). Two genes, *knirps* (*knr*) and *knirps-like* (*knr-1*) encode very similar proteins but differ greatly in intron length, about 1 kb in *kni* and 19 kb in *knr-1*. The gene *kni* is essential at the blastoderm stage of embryonic development, when nuclear division proceeds very rapidly. Mutants that have no functional *kni* can be partially rescued by transcription of *knr-1* cDNA, which lacks introns, but not by genomic *knr-1*, presumably because the latter does not have time to be transcribed between one nuclear division and the next. Although it has not been demonstrated, it may be that for some developmental genes a long lead-time between transcription-induction and expression of the protein product is appropriate. In all events, organisms have become adjusted to the various transcriptional times of their genes, and hence to some extent depend on maintaining their accustomed intron lengths.

#### 4. The Origin of Introns

Discussion about the origin of introns has been dominated by the idea, first put forward by Walter Gilbert (20), that exons were once separate genes and that introns are the present-day relics of the gaps between them when they were first brought together early in evolution. This is sometimes called the “introns early” hypothesis. The strongest evidence in favor of this view is the observation that where proteins consist of different **domains**, somewhat separate in the three-dimensional structure, and distinct even if mutually dependent functions, they are encoded by different exons (20). The heavy and light chains of [immunoglobulins](#) are good examples. And where pre-mRNA are spliced in different ways to yield functionally different protein products, it is clear that at least some of the introns mark the boundaries between rearrangeable protein domains that are functionally distinct modules. It is also true (see [Gene Structure](#)) that numerous multifunctional genes in eukaryotes are clearly homologous to separate single-function genes in bacteria, which are generally seen as more akin to the earliest organisms (though, in fact, they have been evolving for just as long as anything else).

There are difficulties, however, in imagining how introns that have an elaborate splicing-out mechanism could have evolved independently in such a similar form in all of the spaces between originally separate domains when they were first brought together. The alternative, the “introns late” hypothesis, is that genes, however complex their functions and however they were formed, came into existence without the benefit of introns and that the introns came in later, perhaps in a form akin to present-day group II introns, which are both mobile and similar to nuclear introns in their splicing mechanism. [Natural selection](#) could explain the correlation between intron positions and interdomain regions, since that these are likely to be the positions where any disruption due to insertion would do the least damage.

Comparisons between homologous genes in different organisms give contradictory messages. Some introns have retained their positions, although not their internal sequences, over vast tracts of evolutionary time, but others come and go even between quite closely related organisms, often into and out of positions that cannot be considered boundaries between functionally distinct domains. The strict introns-early hypothesis, if anyone held strictly to it, would suppose that every intron position that we see now was there from the origin of the gene and that the reason that all the positions are not occupied now is that introns can be lost, most plausibly by gene conversion, when patches of cDNA replace the intron-containing genomic sequences. But as has been pointed out (22), if introns were present in the primeval gene in all their present-day positions, the exons must have been so short as to be incapable of encoding any functionally distinct part of a protein.

An unanswered, perhaps unanswerable, question is whether the longest introns were always so long and, if not, how they grew.

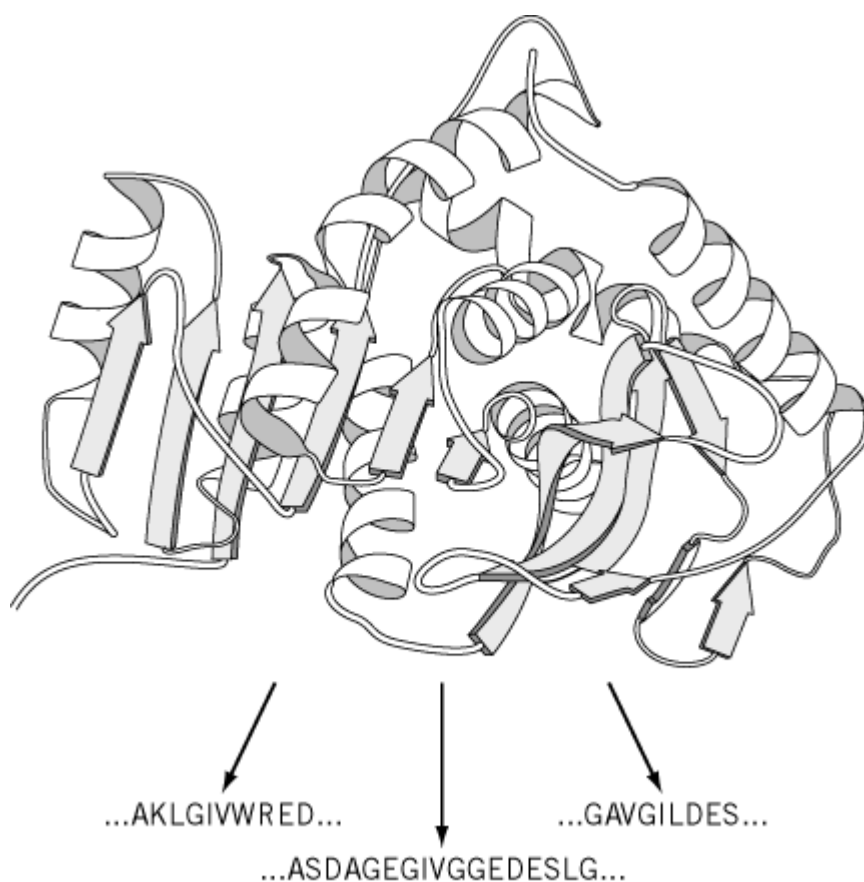
### Bibliography

1. A. M. Lambowitz and M. Belfort (1993) *Annu. Rev. Biochem.* **62**, 587–622.
2. A. J. Jeffries and R. A. Flavell (1977) *Cell* **12**, 1097–1108.
3. A. Dugaiczuk et al. (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2253–2257.
4. W. Filipowicz, in Ref. 22.
5. J. A. Engelbrecht, K. Voelker-Meiman, and G. S. Roeder (1993) *Cell* **66**, 1257–1268.
6. C. W. Siebel, L. D. Fresco, and D. C. Rio (1992) *Genes Dev.* **6**, 1386–1401.
7. J. Valcarcel, R. Singh, and M. R. Green, in Ref (22).
8. B. S. Baker (1989) *Nature* **340**, 521–524.
9. J. L. Manley and R. F. Tacke (1996) *Genes Dev.* **10**, 1569–1579.
10. K. W. Lynch and T. Maniatis (1995) *Genes Dev.* **9**, 284–293.
11. Q. Sun et al. (1993) *Genes Dev.* **7**, 2598–2608.
12. J. Valcarcel, R. Singh, and P. D. Zamore et al. (1993) *Nature* **362**, 171–176.
13. T. R. Cech (1990) *Annu. Rev. Biochem.* **59**, 543–568.
14. B. Dujon (1989) *Gene* **82**, 91–114.
15. F. Michel and J. L. Ferat (1995) *Annu Rev. Biochem.* **64**, 435–461.
16. S. Zimmerly, H. Guo, P. S. Perlman, and A. M. Lambowitz (1995) *Cell* **82**, 545–554.
17. C. H. Sellem, Y. d'Aubenton-Carafa, J. L. Rossingnol, and L. Belcour (1996) *Genetics* **143**, 777–788.
18. C. T. Beagley, N. A. Okada, and D. R. Wolstenholme (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5619–5923.
19. M. Rothe M. Pehl H., Taubert, and H. Jackle (1992) *Nature* **359**, 156–159.
20. W. Gilbert (1978) *Nature* **271**, 501.
21. C. C. F. Blake, (1979) *Nature* **277**, 598 and *Inter.Rev. Cytol.* **95**, 149–185.
22. A. Stoltzfus, J. M. Logsdon Jr., J. D. Palmer, and W. F. Doolittle (1979) *Proc.Nat. Acad. Sci. USA* **94**, 10739–19744.

### Inverse Folding Problem

This problem arose from attempts to design [proteins \*de novo\*](#). Here one begins by selecting a certain known target [protein structure](#), for example, the [globin fold](#), and freely designs amino acid sequences ([primary structures](#)) that should adopt this [conformation](#) spontaneously. There will be a number of sequences compatible for a particular target because **homologous** proteins are known to adopt the same fold (see [Homological Modeling](#)). A solution becomes unique, however, if one searches for the sequence that would fold and would give the most stable structure. Such a sequence is likely to differ even from the native sequence of the target structure. In this procedure, one starts from a 3-D structure and works toward a sequence (Fig. 1), opposite to the direction of **protein folding** and [protein structure prediction](#), where one starts with the sequence and tries to find the 3-D structure.

**Figure 1.** Inverse folding problem. Given a particular 3-D structure, search for any amino acid sequence that would preferably adopt the fold. Note that the thinking is reversed in direction compared with that of protein folding.



One of the obstacles in solving the inverse folding problem is dealing with the amino acid residue side chains. The protein backbone can be fixed just the same as in the target structure, but the side chains vary with the sequence. At each stage, the optimal packing of the atoms of the side chains within the protein interior needs to be determined. To avoid this annoying task, a simplified treatment of side chains has been introduced (1); each of the 20 side chains of the normal amino acids is represented by one point (the  $C_{\beta}$  atom for all residues except [glycine](#)), and its bulkiness and other physicochemical properties are effectively included in the interaction potential used, which is called the “mean-force energy potential” or the “Sippl potential.” The position of the  $C_{\beta}$  atom is fully determined by the dihedral angles of the backbone alone (see [Ramachandran Plot](#)), so it is not necessary to worry about the side-chain conformation. The Sippl potential is a function of the



distance between the C<sub>β</sub> atoms of two side chains and of the amino acid type. The quality of the side-chain packing is given quantitatively by the sum of the interaction energies estimated by the Sippl potential between the side chains. The applicability of the Sippl potential has been established by the “Sippl test (2).” The native amino acid sequence *a* of one structure (A) is selected from a structural library of various known structures and mounted onto another structure (B) that is larger in size than A. Mounting sequence *a* onto structure B is possible with various alignments, shifting residues one-by-one without introducing any gaps. Each time, the total energy of the structure is estimated by using the Sippl potential. Upon completion, structure B is replaced by the next structure of the structural library, including structure A itself. If the native combination of sequence and structure gives the lowest energy among all the alternatives, it is concluded that the potential function is effective and useful. The Sippl test has been performed successfully with various Sippl-type potentials (3, 4), but this is not sufficient for the inverse folding problem. The Sippl test is a recognition test for sequence recognizing structure, just opposite to the inverse folding problem of structure recognizing sequence.

Another difficulty with the inverse folding problem is whether the energy values of different systems (proteins with different sequences) can be compared to each other (5, 6). Of course, it is almost meaningless to compare the energy of two different substances, such as ethane and ethanol. Apart from the absolute energy, however, the energy of protein stability is defined as the energy difference between the folded and **unfolded** protein states, which can be compared among different proteins. Ota et al. (7) introduced a simple approximation to estimate the energy of protein stabilization in combination with the Sippl-type potential energy and applied the method to analyzing the thermal stability of mutant proteins. Similar treatments could conceivably deal with the inverse folding problem and achieve the *de novo* design of proteins.

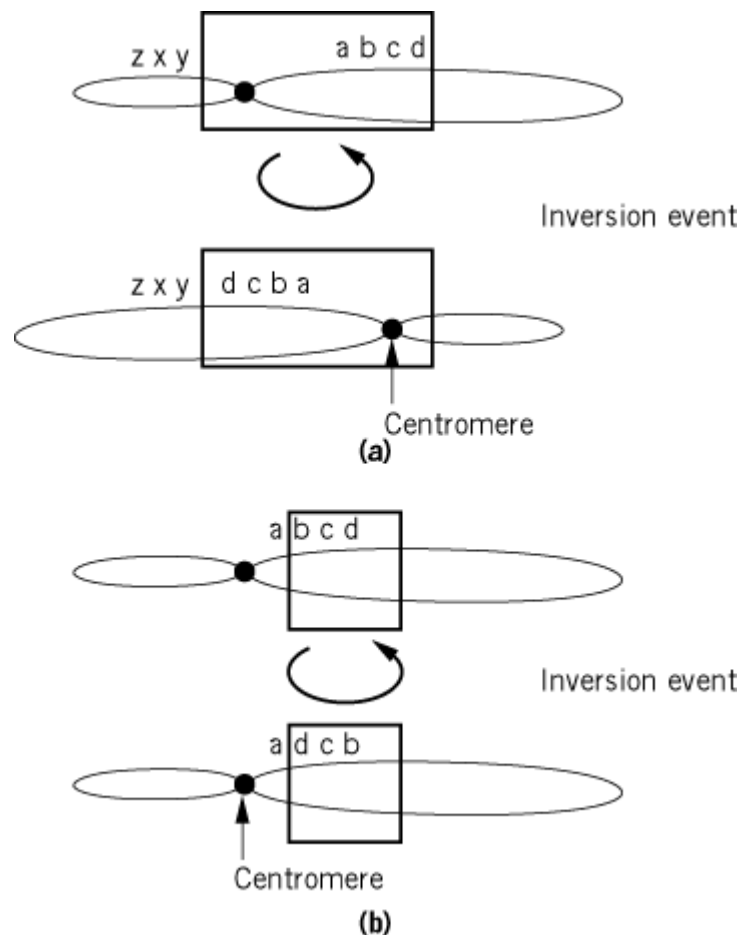
#### Bibliography

1. M. J. Sippl (1990) *J. Mol. Biol.* **213**, 859–883.
2. M. Hendlich, P. Lackner, S. Weitckus, H. Floeckner, R. Froschauer, K. Gottsbacher, G. Casari, and M. J. Sippl (1990) *J. Mol. Biol.* **216**, 167–180.
3. C. Ouzounis, C. Sander, M. Scharf, and R. Schneider (1993) *J. Mol. Biol.* **232**, 805–825.
4. Y. Matsuo, H. Nakamura, and K. Nishikawa (1995) *J. Biochem.* **118**, 137–148.
5. M. J. Rooman and S. J. Wodak (1995) *Protein Eng.* **8**, 849–858.
6. G. M. Crippen (1996) *Struct. Function Genet.* **26**, 167–171.
7. M. Ota, S. Kanaya, and K. Nishikawa (1995) *J. Mol. Biol.* **248**, 733–738.

#### Inversion, Chromosomal

The major sources of variation in the eukaryotic [genome](#) are chromosomal breakages and reunions, leading to inversions and [translocations](#). If two breaks occur in a [chromosome](#) and the excised fragment is rotated 180° before being reunited, then an inversion occurs. Three breaks in a chromosome can lead to an intrachromosomal translocation. In humans, both pericentric (the position of the [centromere](#) is changed) and paracentric (the position of the centromere is not changed) inversions occur in 1 to 2% of the population. Surprisingly often these gross chromosomal inversions do not have any **phenotypic** consequence (1). Most of these inversions are pericentric and involve nonrandom chromosomal breakpoints (Fig. 1).

**Figure 1.** Chromosomal inversions. Those that change the position of the centromere are called pericentric inversions (a), and those that maintain the position of the centromere are called paracentric (b).



The [polytene chromosomes](#) of *Drosophila* are especially suited for studying gross chromosomal rearrangements. It has been shown that the evolution of distinct species of *Drosophila* was associated with the accumulation of many paracentric and pericentric inversions, so that on average two new inversions can be identified in each species (2). Comparison of the chromosomal banding patterns of humans with that of various apes also reveals that several chromosomes have undergone inversions leading to detectable differences (3). Inversions have also been well documented in **plants** and **fungi**. Nevertheless, inversions are an unwelcome occurrence for sexually reproducing organisms because of the consequence for **meiosis**. Chromosomal inversions cause [homologous chromosomes](#) to become locally heterozygous, which can create problems during the chromosomal alignment process that happens in meiotic **prophase**. This heterozygosity is normally overcome by the formation of an inversion loop, so that chromosomal alignment is maintained where possible. Inversions also create problems for the appropriate segregation of chromosomal material if **crossovers** occur within the inversion loop.

In humans, fertility is reduced because paracentric inversion prevents appropriate chromosomal segregation when a crossover occurs within an inversion loop. The products of a paracentric inversion are an [acentric fragment](#) and a [dicentric chromosome](#), which breaks during **anaphase**. Only the two [chromatids](#) that are not involved in the crossover remain normal. The [embryos](#) generated from **gametes** without the full complement of chromosomes are generally inviable. With pericentric inversions, the risk of producing an abnormal embryo is increased, depending on the origin of the gamete. The products of a pericentric inversion have [gene duplications](#) and deficiencies.

These duplications potentially lead to defective gametes. With male gametes, there is a 4% risk of abnormality, and with female gametes there is an 8% risk. This presumably reflects the greater contribution of the female gamete and the maternal chromosomes to early embryonic [development](#). Comparable decreases in fertility are found in flowering plants. This leads to a strong evolutionary selection against chromosomal inversion ([4](#)).

#### Bibliography

1. D. S. Moorehead (1976) *Am. J. Human. Genet.* **28**, 294–301.
2. H. Carson (1983) *Genetics* **103**, 465–489.
3. J. J. Yunis and O. Drakash (1982) *Science* **215**, 1525–1529.
4. S. Wright (1977) *Evolution and Genetics of Populations*, Vol. **4**: Variability within and among natural populations, University of Chicago Press, Chicago.

#### Suggestion for Further Reading

5. M. S. Clark and W. J. Wall (1996) "Chromosomes". *The Complex Code*, Chapman and Hall, London.

## Inverted Repeats

A repeat, or repeating, **nucleotide sequence** is any of two or more identical segments of a longer sequence or identical segments in a double-stranded **DNA**. If two such DNA fragments are inserted into a long DNA duplex with opposite polarity (one “head on,” the other “tail on”), the corresponding pair are known as inverted repeats. The various types of repeats are schematically outlined in the Figure for [Tandem Repeats](#). The sequences of the pair of inverted repeats, if read from the same strand, are complementary to one another. Inverted repeats placed immediately after one another make a [palindrome](#). If the separation between the inverted repeats is comparatively short (up to a few tens of bases), the repeats may correspond to a hairpin in a RNA transcript generated from such a structure, or to a [cruciform](#) structure in DNA under torsional constraint ([1](#)).

Inverted repeats are an important element of several types of mobile DNA elements, or [transposons](#) ([2](#)). They are located at the ends of the transposons when integrated into a [genome](#). For example, bacterial **insertion sequences** (IS) are flanked by 10- to 40-bp inverted repeats. The insertion sequences themselves can make a pair of inverted repeats, as in the case of the bacterial transposon Tn10, which is flanked by 1400-bp IS10 sequences of opposite polarity. The inverted repeats are involved in the [recombination](#) events that occur during integration or excision of the transposons.

#### Bibliography

1. Y. Timsit and D. Moras (1996) *Q. Rev. Biophys.* **29**, 279–307.
2. *Mobile DNA* (1989) (D. E. Berg and M. M. Howe, eds.) ASM, Washington, D.C.

#### Suggestion for Further Reading

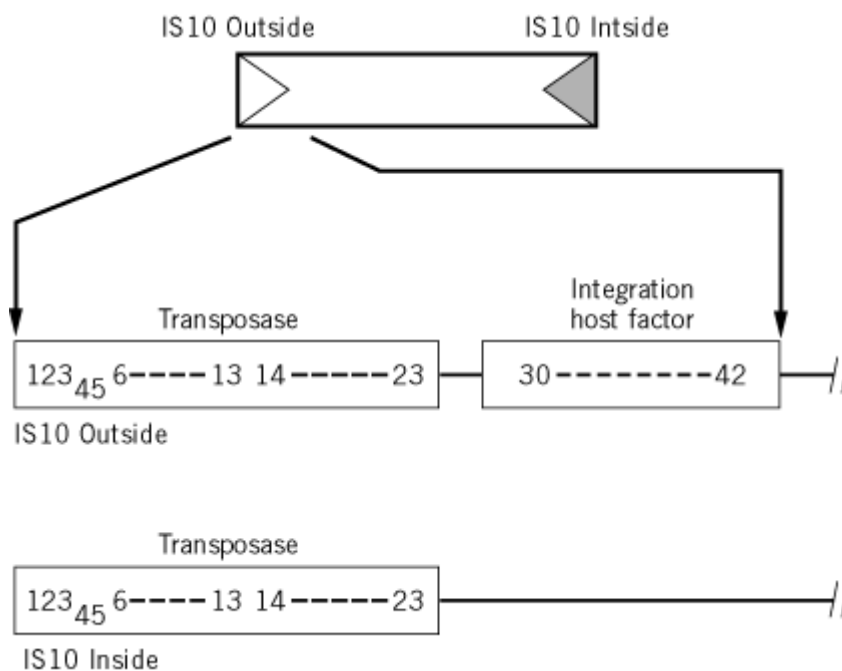
3. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.

## Inverted Terminal Repeat

The movement of [transposable elements](#) from place to place in [genomes](#) results from the actions of a special element-encoded [recombinase](#) on special DNA sequences at the tips of the element. In virtually all such elements, these special [recombination](#) sequences at the ends of the transposon include transposase binding sites arranged as inverted repeats. It is these sites that define where the breakage and joining events must occur. Because identical (or nearly identical) terminal sequences are arranged as inverted repeats, the disposition of the transposase at both ends relative to the transposon ends is identical, so that the breakage and joining will occur in the same relative positions at both ends.

The terminal inverted repeats are usually composed of special sequences at the extreme tips of the transposon that signal exactly where that breakage and joining should occur, plus more internal binding sites for the transposase (1, 2) (Fig. 1). In some elements, there are multiple binding sites for the transposase, including the one that is part of the terminal inverted repeat and others internal (3). These interior sites probably play a role in loading the transposase appropriately onto the termini and may be dispensable after the initial stages of transposition have occurred, for example, after synapsis of the transposon ends. For some elements, a binding site for a host factor important for recombination may also be part of the terminal repeat.

**Figure 1.** The ends of IS10. Special sequences at the ends of the element specify the position of transposase binding and also provide signals to prompt DNA breakage and joining. Transposase binding is specified by sequences from position 6 to position 23, and the positions from 6 to 13 are most critical. Although particular sequences at positions 1, 2, and 3 are not required for transposase binding, they are necessary to prompt the chemical steps of DNA breakage and joining. The identity of positions 4 and 5 is not critical. The outside and inside ends of IS10 have the same transposase binding motifs, but the outside end has in addition a binding site for integration host factor (positions 30–42), a bacterial protein, that is a cofactor in reactions involving the outside ends.



Although multiple transposase binding sites may be present, breakage and joining is restricted to the terminal inverted repeats because of the presence of a transposase binding site adjacent to the special sequences at the extreme tips of the element that provoke cleavage (1, 2). Alteration of these “cleavage signal” nucleotides generally blocks cleavage, but transposase binding still occurs. Often there is a spacer region of a few nucleotides between the transposon binding site and the cleavage nucleotides at the tip.

## Bibliography

1. D. Haniford and N. Kleckner (1994) EMBO J. **13**, 3401–3411.
2. J. Sakai, R. M. Chalmers, and N. Kleckner (1995) EMBO J. **14**, 4374–4383.
3. B. D. Lavoie and G. Chaconas (1996) Curr. Top. Microbiol. Immunol. **204**, 83–102.

## Iodine Isotopes

Iodine is element number 53 in the periodic table and has valence states of 1, 3, 5, or 7 (1). Iodine is an essential element in humans in small amounts, and lack of iodine causes thyroid goiter. The thyroid concentrates iodine, synthesizes and stores the thyroid hormone **thyroxine**, and releases it to maintain physiological balance. Excessive exposure to chemical iodine can cause skin, eye, and mucous membrane irritations. Iodine compounds are used as chemical reagents, as antimicrobial agents, as radiopharmaceuticals in medicine, and in photographic chemicals.

Thirty-five isotopes of iodine have been identified (2), ranging in atomic mass from  $^{108}\text{I}$  (half-life = 36 ms) to  $^{142}\text{I}$  (half-life = 0.25 s) (see [Radioactivity](#) and [Radioisotopes](#)). Only one isotope of iodine is stable ( $^{127}\text{I}$ ). The longest lived isotope ( $^{129}\text{I}$ ) has a half-life of 15.7 million years.

The four most important radioactive isotopes of iodine (or *radioiodine*) for biomedical use are:  $^{123}\text{I}$ ,  $^{124}\text{I}$ ,  $^{125}\text{I}$ , and  $^{131}\text{I}$  (Table 1). Each radioiodine is prepared uniquely and has a variety of different applications, advantages, and disadvantages (3). Iodine is a useful **radiolabel** for many studies, as it can be chemically incorporated into proteins relatively easily. All radioiodine compounds must be handled with care, as radioiodine is volatile and readily concentrated in the human thyroid gland.

**Table 1. Principal Radioisotopes of Iodine in Medicine and Research: Their Radioactive Emissions, Yields, and Energies (2)**

| Element | Mass (amu) <sup>a</sup> | Physical Half-life | Beta-Particle Yield | Average Beta Energy (MeV) <sup>b</sup> | Gamma-Ray Yield         | Gamma Energy (MeV)      |
|---------|-------------------------|--------------------|---------------------|--|-------------------------|-------------------------|
| Iodine  | 123                     | 13.22 hours        | —                   | —                                      | 0.833                   | 0.159                   |
|         | 124                     | 4.18 days          | 0.229               | 0.824                                  | 0.453<br>0.611<br>0.101 | 0.511<br>0.603<br>0.723 |

|     |            |     |       |        |        |
|-----|------------|-----|-------|--------|--------|
|     |            |     |       | 0.105  | 1.691  |
| 125 | 59.6 days  | —   | —     | 0.0667 | 0.0355 |
| 131 | 8.021 days | 1.0 | 0.182 | 0.0606 | 0.284  |
|     |            |     |       | 0.812  | 0.364  |
|     |            |     |       | 0.0727 | 0.637  |

---

<sup>a</sup> Atomic mass units.

<sup>b</sup> Million electronvolts.

High purity iodine-123 is usually produced in an accelerator by the reaction  $^{127}\text{I}(p,5n)^{123}\text{Xe} \rightarrow ^{123}\text{I} + b^-$ , or by the reaction  $^{124}\text{Xe}(p,g)^{123}\text{I}$ . Iodine-124 is an undesirable impurity in  $^{123}\text{I}$  preparations. Iodine-123 is used as a diagnostic imaging agent in nuclear medicine for thyroid imaging and uptake studies. Its relatively short half-life of 13 h. is a disadvantage for transportation from site of production to hospitals for clinical use.

Iodine-124 can be prepared in an accelerator by deuteron bombardment of tellurium according to the reaction  $^{124}\text{Te}(d,2n)^{124}\text{I}$ . Iodine-124 is a positron (beta-plus) emitter having useful applications in positron-emission tomography in diagnostic nuclear medicine. Iodine-124 has also been used as a label for 5-iodo-2'-deoxyuridine, an analog of thymidine, for DNA uptake studies, and tumor imaging. It has potential application as a therapeutic radionuclide.

Carrier-free  $^{125}\text{I}$  is prepared in a reactor using an  $^{124}\text{Xe}$  target according to the reaction  $^{124}\text{Xe}(n,g)^{125}\text{Xe} \rightarrow ^{125}\text{I} + b^-$ . It may also be prepared in an accelerator by deuteron bombardment of tellurium according to the reaction  $^{124}\text{Te}(d,n)^{125}\text{I}$ . Iodine-125 emits low energy (0.035 MeV) gamma rays and thus is *not* useful for external imaging. However,  $^{125}\text{I}$  is used extensively to **radiolabel** small specimens and [proteins](#) for a variety of applications in molecular biology, including [radioimmunoassay](#), [autoradiography](#), and [fluorography](#). Iodine-125 is a common label for 5-iodo-2'-deoxyuridine in DNA uptake and **hybridization** studies. It has possible applications as a therapeutic radionuclide because of the high linear energy transfer of Auger electrons emitted by  $^{125}\text{I}$ -labeled DNA in cell nuclei. Iodine-125 is also used in permanent seed implants (brachytherapy) for treatment of slowly growing tumors, such as carcinoma of the prostate.

Iodine-131 may be prepared in an accelerator according to the reaction  $^{130}\text{Te}(d,n)^{131}\text{I}$ . It is more commonly prepared in a reactor by tellurium capture according to the reaction  $^{130}\text{Te}(n,g)^{131}\text{Te} \rightarrow ^{131}\text{I} + b^-$ , or as a  $^{235}\text{U}$  fission product. Iodine-131 was discovered and used for thyroid function studies in 1938 (4). It was the first radioisotope to be approved (in 1951) by the U.S. Food and Drug Administration for human use as a new drug.

Today,  $^{131}\text{I}$  is used extensively in clinical nuclear medicine for many purposes (3). These include thyroid uptake studies (555 kBq, or kilobecquerels, administered), scans with metastatic thyroid cancer in athyroid patients, and treatment of benign thyroid disease (185 MBq), for thyroid carcinoma ablation (5.55 GBq), and as a labeled [monoclonal antibody](#) in targeted-cell therapy of lymphoma and leukemia (up to 30 GBq). Other applications are radioiodine labels for dyes, macromolecules, proteins, and hormones, including  $^{131}\text{I}$ -norcholesterol for diagnosis of aldosterone-producing tumors in the adrenals and  $^{131}\text{I}$ -metaiodobenzylguanidine for imaging pheochromocytoma and other neuroendocrine tumors. Iodine-131 has the advantage of a relatively long half-life (8.021

days) for production, transportation, and labeling, but the disadvantage of abundant gamma rays, which can be a radiation hazard.

### Bibliography

1. D. R. Lide and H. Pr. Frederikse, eds. (1995) *CRC Handbook of Chemistry and Physics*, CRC Press, Boca Raton, Fla.
2. Knolls Atomic Power Laboratory (1966) *Chart of the Nuclides*, 15th ed., available from General Electric Company, San Jose, Calif.
3. J. C. Harbert (1996) in *Nuclear Medicine: Diagnosis and Therapy*, J. C. Harbert, W. C. Eckelman, and R. D. Neuman, eds., Chapter "20", pp. 407–427.
4. M. Brucer (1990) *A Chronology of Nuclear Medicine*, Heritage Publications, Inc., St. Louis, Mo.

### Suggestions for further reading

5. L. K. Bustad, ed. (1964) *Biology of Radioiodine*, Pergamon Press, Oxford.
6. J. C. Harbert, W. C. Eckelman, and R. D. Neumann (1996) *Nuclear Medicine: Diagnosis and Therapy*, Thieme Medical Publishers, New York.

## Iodoacetamide, Iodoacetate

Iodoacetamide and iodoacetate are two of the most common reagents used for chemically modifying [thiol groups](#) irreversibly. They have a long history, dating back to 1874 when it was reported that bromoacetate injected subcutaneously into frogs killed them, as a result of progressive respiratory and cardiac depression; iodoacetate was subsequently found to be more reactive. These studies were important in elucidating the mechanism of glycolysis, which led to the realization that iodoacetate and other such compounds acted by reacting with the thiol groups of proteins, in this case at the [active site](#) of the [enzyme](#) glyceraldehyde-3-phosphate dehydrogenase. Iodoacetamide became popular because it can cross cell membranes much more readily than can iodoacetate and therefore is more potent with cells and organs.

Both iodoacetamide and iodoacetate react rapidly and specifically with the ionized form of thiol groups, to generate very stable adducts:



These adducts are commonly known as the carboxamidomethyl and carboxymethyl, respectively. They are stable to conditions routinely used to hydrolyze proteins to their substituent [amino acids](#), except that the carboxamidomethyl group gets hydrolyzed to the carboxymethyl (see [Amino Acid Analysis](#)).

The reaction with iodoacetamide or iodoacetate occurs by the thiolate anion attacking and displacing the iodine atom in a nucleophilic reaction. Consequently, the rate of the reaction increases with increasing pH, up to just beyond the  $\text{p}K_a$  of the thiol group, which typically is about 9 for most biological substances. The pH dependence of the rate of the reaction can be used to measure the

apparent  $pK_a$  value of the thiol group, although there are exceptional cases in catalytic proteins where the results are misleading, presumably because the protein modifies the energy of the transition state for the reaction (1). The reaction occurs rapidly with thiol groups. The second-order rate constant for reaction at 25°C of iodoacetamide with a fully ionized cysteine thiol group, with a  $pK_a$  of about 8.7, is approximately  $25 \text{ s}^{-1}\text{M}^{-1}$ . Consequently, the half-time for modification of such a thiol group upon adding 0.1 M iodoacetamide is only 0.3 sec. Iodoacetate reacts about three times less rapidly than does iodoacetamide.

It is important to realize that protons are released by the reaction. Frequently, high concentrations of a thiol reagent are used to reduce all cysteine residues, and then a large excess of iodoacetate or iodoacetamide is added to react with both the protein and the reagent. In this case, the pH will drop and the rate of the reaction will slow, and even stop before completion, unless very high concentrations of a [buffer](#) are present.

The two reagents differ only in one being an amide, and neutral, and the other an acid, and negatively charged. This charge difference is useful for [counting residues](#) of thiol groups in a protein. Mixtures of the two reagents are reacted with a protein containing  $N$  thiol groups, to produce a spectrum of molecules with net charge differences of 0 to  $N$ . If this spectrum is resolved by [electrophoresis](#), [isoelectric focusing](#), [ion-exchange chromatography](#), or other techniques dependent only upon the charge of the protein, the number of species with different charges present can be counted, giving the integer value of  $N$ .

Both iodoacetamide and iodoacetate also react with [amino groups](#), but at a lower rate, and primarily at very alkaline pH values, where the reactive nonionized amino groups are present. Both [lysine](#) and [histidine](#) residues can be modified in this way. The reaction of **ribonuclease A**, which has no free thiol groups, with the two histidine residues at its [active site](#), each reacting with different N atoms, was important in demonstrating the specificity of the interactions that could take place in enzyme-active sites (2). They can also react with the sulfur atom of **methionine residues** to generate the positively charged sulfonium salt, but at a much lower rate (3). This rate is essentially independent of pH, however, so it can become the predominant modification reaction of a protein at very low pH values. The reactions of iodoacetamide and iodoacetic acid can be used in [diagonal methods](#) to isolate specifically cysteine- and methionine-containing peptides.

Many variants of iodoacetate and iodoacetamide have been devised by adding further groups with absorbance and fluorescence properties that serve as reporters of structure in proteins.

#### Bibliography

1. J. W. Nelson and T. E. Creighton (1994) *Biochemistry* **33**, 5974–5983.
2. R. G. Fruchter and A. M. Crestfield (1967) *J. Biol. Chem.* **242**, 5807–5812.
3. H. G. Gundlach, S. Moore, and W. H. Stein (1959) *J. Biol. Chem.* **234**, 1761–1764.

#### Suggestions for Further Reading

4. J. L. Webb (1966) *Enzyme and Metabolic Inhibitors*, Vol. **III**, "Chapter 1", Academic Press, New York.
5. F. R. N. Gurd (1967) Carboxymethylation. *Methods Enzymol.* **11**, 532–541.

## Ion Channel Receptors



A diverse group of neurotransmitter receptors are ligand-gated **ion channels**. These can be subgrouped into receptors that are either excitatory (acetylcholine, glutamate, serotonin) or inhibitory [g-aminobutyric acid (GABA), glycine]; they have similar structural features, but different regulatory properties. These receptors are characterized by fast- and short-acting responses to ligand. Additionally, these receptors are subject to targeting by docking proteins that induce aggregation. The increased permeability to ions results in the influx of [second messengers](#), resulting in activation of downstream targets.

## 1. Receptor Structure

Because it was one of the first receptors studied in detail, more is known about the nicotinic [acetylcholine receptor](#) than any other (1). This receptor consists of four different **polypeptide chains**, with the composition  $\alpha_2\beta\gamma\delta$ . Each subunit has a single **transmembrane** domain. Although the structure of the channel is not known, the five subunits are thought to assemble together to form a pore, with the neurotransmitter binding site located in the extracellular domain of the subunits. Numerous mutational analyses have led to an understanding of the molecular dynamics of the regulation of the channel by ligands. Thus, acetylcholine binding induces a conformational change, opening the gate to allow influx of  $\text{Na}^+$  and  $\text{K}^+$  ions (1).

The structures of the inhibitory ligand-gated anion channels are similar to the nicotinic cholinergic receptor. Both GABA and glycine activate  $\text{Cl}^-$  channels, in a manner analogous to the regulation of  $\text{K}^+$  and  $\text{Na}^+$  channels by acetylcholine.

## 2. Regulation of Ligand-Activated Channels

Ligand-activated channels are generally found to be clustered on the postsynaptic side of a synapse. Recent studies have suggested that the clustering is mediated by interactions with a group of targeting proteins, such as PSD-95 and GRIP, that contain PDZ domains. These regions are **domains** involved in [protein-protein interactions](#) and are known to interact with specific sequences in the C-terminal regions of transmembrane proteins (2).

Ligand-gated ion channels are also known to be regulated by **phosphorylation**. The [primary structures](#) of these receptors contain numerous **consensus sequences** for various protein [kinases](#). Both *in vivo* and *in vitro* experiments have demonstrated phosphorylation of both excitatory and inhibitory receptors, by both **tyrosine kinases** and **serine/threonine kinases**. These phosphorylations are likely to play a role in modulating the activation of the receptor by ligands, and they may have consequences for understanding long-term modulation of receptor function (3).

See also [Calcium Signaling](#).

## Bibliography

1. A. Karlin (1993) *Curr. Opin. Neurobiol.* **3**, 289–309.
2. H.-C. Kornau, P. H. Seeburg, and M. B. Kennedy (1997) *Curr. Opin. Neurobiol.* **7**, 368–373.
3. T. C. Smart (1997) *Curr. Opin. Neurobiol.* **7**, 358–367.

## Ion Pair

An ion pair is two molecules or groups of opposite electric charge in close proximity due to favorable [electrostatic interactions](#) between them. The term ion pair is generally used in all fields of liquid chemistry, in contrast to a [salt bridge](#), but in biochemistry the term frequently refers to a salt bridge.

## **Ion-Exchange Chromatography**

Ion exchange is the reversible interchange of ions of the same charge between a solution and a solid, insoluble material in contact with it. It proceeds by equivalents. One chemical equivalent of one kind of ion enters the solid, and one chemical equivalent of another kind of ion must leave. Therefore, electrical neutrality is always maintained. An ion exchanger consists of an insoluble matrix to which charged groups have been covalently bound. The charged groups associate with mobile counterions. These counterions are reversibly exchanged with other ions of like charge (ie, ion exchange) without altering the matrix. Two main forms of ion exchangers are commonly used: anion exchangers and cation exchangers. Positively charged exchangers have negatively charged counterions (anions) available for exchange and so are termed anion exchangers. Negatively charged exchangers have positively charged counterions (cations) and are termed cation exchangers.

Ion-exchange chromatography is an inclusive term for all **chromatographic** separations of ionic substances carried out by using an insoluble ion exchanger as the stationary phase. As an ionic solution passes over and through the exchanger, ions from the solution exchange for those already on the solid matrix. This technique came into prominence during World War II as a separation procedure for the rare earth and transuranium elements. The technique was first used by Taylor and Urey (2) to separate lithium and potassium isotopes using zeolite resins, and Samuelson (3) demonstrated the potential of synthetic resins. In the mid 1950s, Sober and Peterson (4-6) synthesized the diethylaminoethyl (DEAE) and carboxymethyl (CM) derivatives of cellulose, still in use today for ion-exchange chromatography of proteins.

Although all ion exchanges are reversible, the equilibrium distributions vary widely, and this variation makes ion-exchange chromatography possible. In other words, the basic principle of ion-exchange chromatography involves competition between different ions of the same charge for binding to an oppositely charged chromatographic matrix, the ion exchanger. The interaction between the analytes and the ion exchanger depends on several factors: the charge properties of the analytes; the ionic strength and the nature of the particular ions in solution; the pH; and the presence of other additives in the mobile phase, such as organic solvents. Because retention in ion-exchange chromatography involves an [electrostatic interaction](#) between the fixed charges of the ion-exchange matrix and those of the analytes, it is clear that the more highly charged an analyte, the more strongly it binds to a given, oppositely charged ion exchanger. Similarly, more highly charged ion exchangers, those with a greater degree of substitution with charged groups, bind analytes more effectively than weakly charged ones. It should be mentioned that the terms strong and weak ion exchangers derive from the  $pK_a$  values of their charged groups and do not say anything about the strength with which they bind analytes. At pH values far from the  $pK_a$ , binding equally strong either to a weak or a strong ion exchanger. In addition to the ion exchange effect, other types of binding also occur in the ion-exchange process. These effects are generally small and are mainly due to **van der Waals** forces and [polar](#) interactions.

Because of its versatility, its high resolving power, its high capacity, and its straightforward basic

principle, ion-exchange chromatography is one of the chromatographic techniques used most often in biochemical studies [see section IV of Ref. [1](#) (pp. 170–270), Ref. [7](#)]. Detailed operational protocols and the lists of its applications are beyond the scope of this volume, and the interested reader is directed to several excellent reviews ([7-10](#)).

## Bibliography

1. C. T. Mant and R. S. Hodges (eds.) (1991) *High-Performance Liquid Chromatography of Peptides and Proteins: Separation, Analysis, and Conformation*, CRC Press, Boca Raton.
2. T. I. Taylor and H. C. Urey (1938) *J. Chem. Phys.* **6**, 429–438.
3. O. Samuelson (1939) *Z. Anal. Chem.* **116**, 328.
4. H. A. Sober and E. A. Peterson (1954) *J. Am. Chem. Soc.* **76**, 1711–1712.
5. E. A. Peterson and H. A. Sober (1956) *J. Am. Chem. Soc.* **78**, 751–755.
6. H. A. Sober, F. J. Gutter, M. M. Wycoff, and E. A. Peterson (1956) *J. Am. Chem. Soc.* **78**, 756–763.
7. E. Karlsson, L. Rydén, and J. Brewer (1989) in *Protein Purification: Principles, High Resolution Methods, and Applications* (J.-C. Janson, and L. Rydén, eds.), VCH, New York, pp. 107–148.
8. E. F. Rossomando (1990) in *Guide to Protein Purification* (M. P. Deutscher, ed.), *Methods in Enzymology* **182**, Academic Press, New York, pp. 309–317.
9. H. F. Walton (1992) in *Chromatography* (E. Heftmann, ed.) (*Journal of Chromatography Library*, Vol. **51A**), Elsevier, Amsterdam, pp. A227–A265.
10. J. Swadesh (1997) in *HPLC: Practical and Industrial Applications* (J. Swadesh, ed.), CRC Press, Boca Raton, pp. 171–243.

## Ionization

Ionizable groups in biological macromolecules are either protonated or deprotonated to have positive or negative ionic charges, depending on the pH condition of the solvent. At a neutral pH, the phosphate groups in nucleic acids are all ionized, and the bases are almost neutral. In contrast, in proteins at a neutral pH, the charge states of the ionizable groups are exquisitely dependent on the individual  $pK_a$  values of the ionizable residues. They directly affect the electronic structures around active sites, and the stabilities of protein structures (see [Electrostatic Interactions](#)).

The degree of ionization of an acid,  $i$ , is  $10^{-pK_{ai}} / (10^{-pK_{ai}} + 10^{-pH})$ , and the degree of ionization of a base,  $j$ , is  $10^{-pH} / (10^{-pK_{aj}} + 10^{-pH})$ . Therefore, the total net change of an [amphoteric](#) molecule, such as a protein, is given by

$$-\sum_i 10^{-pK_{ai}} / (10^{-pK_{ai}} + 10^{-pH}) + \sum_j 10^{-pH} / (10^{-pK_{aj}} + 10^{-pH})$$

The charge state of an ionizable group is governed by its intrinsic acidity and also by the electrostatic environment, which depends on the ionic strength of the solvent. When an ionizable group is surrounded by other anions, the  $pK_a$  value increases. In contrast, when it is surrounded by cations,

the  $pK_a$  value decreases. Therefore, when a [salt bridge](#) is formed between an acid and a base, the  $pK_a$  of the acid decreases, and the  $pK_a$  of the base increases. Thus, each ionizable group in a protein has a unique  $pK_a$  value, and this changes with changes in the ionization of nearby groups. The apparent  $pK_a$  value is determined by titration experiments or estimated by theoretical calculations based on the tertiary molecular structure. Even when the only structural information about a protein is its amino acid composition, the total net charge of the protein can be approximated, and the [isoionic point](#) estimated, by using only the intrinsic  $pK_a^{\text{int}}$  values, listed in Table 1 of [Electrostatic interactions](#).

## Iron-Binding Proteins

Iron is the most abundant transition metal in the earth's crust, but it is poorly available as a nutrient because it forms insoluble ferric hydroxides in the presence of  $O_2$  at neutral pH. Organisms have coped with this challenge by evolving various uptake strategies based on chelation. Once internalized, free iron is quite toxic to cells:  $Fe^{2+}$  and  $Fe^{3+}$  ions react with  $H_2O_2$  and  $O_2^-$  to generate hydroxyl radicals, which indiscriminately oxidize cellular constituents. Hence, there is a clear need for some form of iron homeostasis (1).

Microbial iron metalloregulation (2) is best understood in *Escherichia coli*. In this bacterium, the ferric uptake regulation (*fur*) gene regulates many genes, including more than a dozen associated with siderophoric biosynthesis and transport. The monomeric Fur protein is approximately 18 kDa in mass and binds a single  $Fe^{2+}$  ion. With  $Fe^{2+}$  bound, Fur targets a GATAATGATAATCATTATC **consensus sequence** near the **promoter** of iron-regulated genes and acts as a dimeric [repressor](#).

*Corynebacterium diphtheriae* (and other bacterial and fungal pathogens) experiences a state of iron starvation upon entering a mammalian host, and it responds by producing [diphtheria toxin](#) at maximal rates. The diphtheria *tox* repressor (DtxR) plays the primary role in regulating toxin biosynthesis.  $Fe^{2+}$ -activated DtxR is a functional dimer, which binds at the *tox* **operator** locus. **X-Ray crystallographic** structures of DtxR suggest that ferrous ions cause a caliper-like conformational change in the relative orientation of the monomers in the dimeric [protein structure](#), leading to a decrease in the distance between the two DNA-recognition [alpha-helices](#) to match that in the target DNA (3) (see [DNA-Binding Proteins](#) and [Repressors](#)).

Bacteria sequester iron during times of plenty and produce bacterioferritins. Mammalian analogs of these iron storage proteins, [ferritins](#), are much better understood.

### Bibliography

1. G. Winkelmann, D. van der Helm, and J.B. Neilands (eds.) (1987) *Iron Transport in Microbes, Plants and Animals*, VCH, Weinheim, Germany.
2. J. H. Crosa (1997) *Microbiol. Mol. Biol. Rev.* **61**, 319–336.
3. N. Schiering, X. Tao, H. Zeng, J. R. Murphy, G.A. Petsko, and D. Ringe (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9843–9850.

### Suggestions for Further Reading

4. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New

York, Vols. 1–35.

5. R. J. Cousins (1994) Metal elements and gene expression, *Annu. Rev. Nutr.* **14**, 449–469.

## Iron-Response Elements

The [messenger RNAs](#) for [iron-binding proteins](#), such as the [transferrin](#) receptor and both subunits of [ferritin](#), contain so-called iron-response elements (IREs). These are RNA structures that contain a stem and a loop that has the **consensus sequence** CAGUGX (1). The transducer proteins, termed iron-responsive proteins (IRPs), bind to the IRE loop when the free iron concentration is low, whereas they are inactivated by high cellular iron levels. Formation of the IRE–IRP complex has a variety of physiological effects, as indicated in Table 1. IRP-binding stabilizes the transferrin mRNA against **RNA degradation**. Regulation of the transferrin receptor expression occurs at the **transcriptional** level. On the other hand, ferritin IRE–IRP complexes block association of the small **ribosomal** subunit with the mRNA of both ferritin subunits and prevent their [translation](#). Consequently, regulation of ferritin expression occurs at the translational level.

**Table 1. Physiological Effects of Key Mammalian IREs**

| mRNA-Containing IRE        | Low [Fe]      | High [Fe]          |
|----------------------------|---------------|--------------------|
| Transferrin receptor       | Fe uptake ↑   | Fe uptake ↓        |
| Ferritin H- and L-chains   | Fe storage ↓  | Fe storage ↑       |
| 5-Aminolevulinate synthase | [porphyrin] ↓ | [porphyrin] normal |

5-Aminolevulinate synthase, the first enzyme in the heme biosynthetic pathway, is also regulated by IRPs. One of its two IRPs is strikingly homologous to aconitase, an iron-sulfur enzyme in the Krebs cycle. Interestingly, [nitric oxide](#) inhibits iron uptake by both IRPs (they remain active and capable of IRE recognition) and hence functions as a vertebrate [hormone](#).

### Bibliography

1. M. W. Hentze and L. C. Kühn (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8175–8182.

### Suggestions for Further Reading

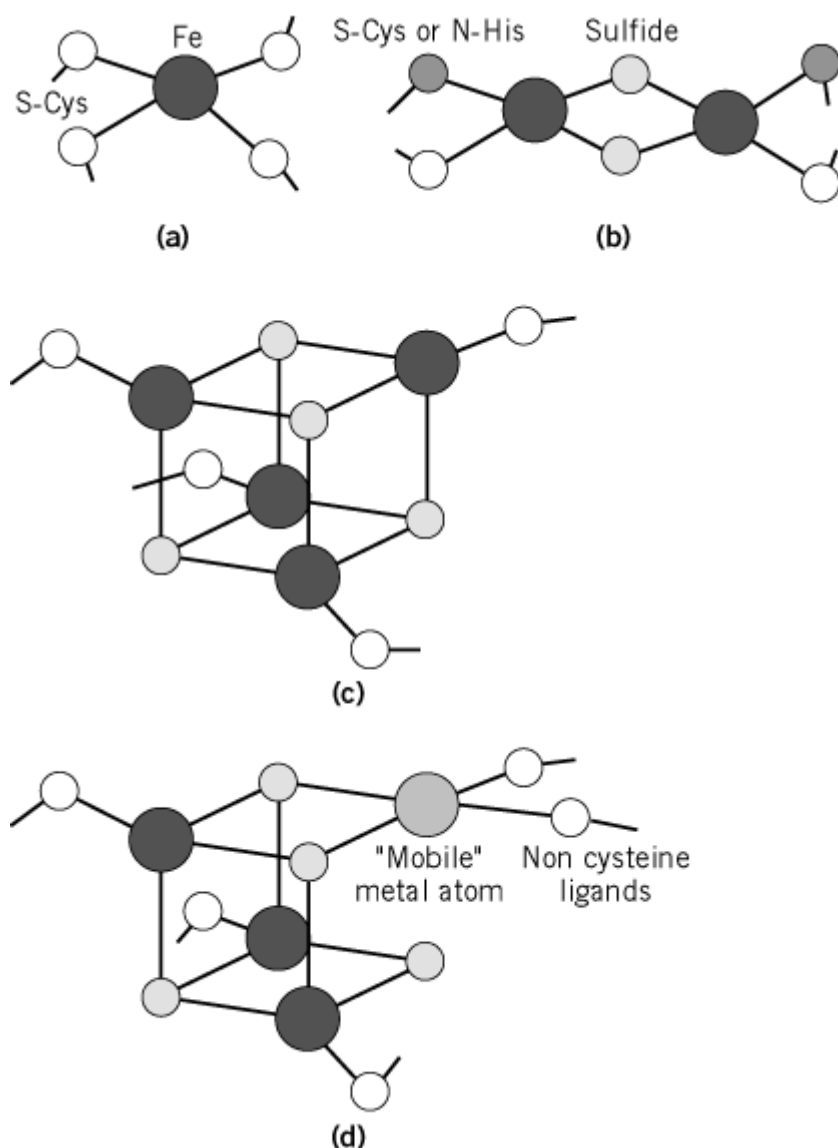
2. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New York, Vols. 1–35.

3. R. J. Cousins (1994) Metal elements and gene expression, *Annu. Rev. Nutr.* **14**, 449–469.

## Iron–Sulfur Proteins

This class of [proteins](#) comprises several types of proteins and [enzymes](#) that contain iron atoms coordinated by sulfur atoms in clusters of different types (Fig. 1). The simplest FeS proteins are the rubredoxins (1Fe, no inorganic sulfur) and the [ferredoxins](#) (2Fe, 3Fe and 4Fe). In more complex FeS proteins, other **prosthetic groups** (FAD, FMN, heme, Mo, Ni, etc.) are present, in addition to the iron-sulfur cluster(s). Thus, FeS proteins fulfill several roles: **electron transfer** and storage, catalysis, structural, and regulatory functions ([1](#), [2](#)).

**Figure 1.** Schematic structures of iron-sulfur clusters found in FeS proteins as determined by [X-ray crystallography](#) studies: (a) the Fe(S-Cys)<sub>4</sub> unit in rubredoxin; (b) the 2Fe–2S cluster in eukaryotic ferredoxins and in Rieske proteins; (c) the cubane-like 4Fe–4S structure of bacterial ferredoxins, HiPIPs, IRPs, and active aconitase; (d) the cuboidal 3Fe–4S structure found in some bacterial ferredoxins, in inactive aconitase, and in some complex redox enzymes; (e) the P cluster of dinitrogenase; (f) the FeMoCo cofactor of dinitrogenase.



### 1. High Potential Electron Carriers: Rubredoxins, HiPIPs, and Rieske Proteins

The simplest of all FeS proteins is rubredoxin, containing a single Fe atom ligated by a four-cysteine-residue unit (Fig. 1a), and no inorganic sulfide, in a polypeptide of 50 to 60 amino acid residues. Rubredoxins are found in strictly anaerobic bacteria and archaea (*Clostridia*, *Desulfovibrionaceae*, *Pyrococcales*) and have [oxidation/reduction potentials](#) (+30 to 60 mV, relative to the hydrogen electrode) that appear to be too high to be relevant for bacterial metabolism. Although assumed to function as electron carriers, the redox partners of rubredoxins in anaerobes have not been identified. Recently, a rubredoxin–oxygen oxidoreductase was found to function *in vitro* at the end of a soluble electron transfer chain that couples NADH oxidation to oxygen consumption in *Desulfovibrio gigas*. There is another protein in *D. gigas* containing a single Fe(Cys)<sub>4</sub> unit in each of its two small subunits, named desulfiredoxin. Rubredoxin-like proteins that contain two Fe(Cys)<sub>4</sub> clusters are also found in *Pseudomonaceae*. These aerobic rubredoxins transfer electrons to enzymes that catalyze the oxidation of alkanes.

Photosynthetic microorganisms such as *Chromatium spp* and *Rhodospira spp* contain the so-called HiPIPs (high-potential iron–sulfur proteins). These proteins contain a single, cysteine-coordinated 4Fe–4S cubane-like structure (Fig. 1c), but these protein structures favor the high-potential 3+/2+ transition of the metal cluster (reduction potentials +90 to +450mV), rather than the 2+/1+ transition typical of [4Fe–4S] clusters in bacterial ferredoxins.

Among the high-potential FeS proteins of the **mitochondrial** and photosynthetic electron transfer chains, a singularity is represented by the so-called “Rieske” protein in the cytochromes b-c<sub>1</sub> and b<sub>6</sub>-f complexes. These “Rieske” proteins contain a 2Fe–2S cluster with two histidine and two cysteine ligands in a unique binding motif (Fig. 1b). The same type of coordination is found in other FeS proteins, most notably in ferroxidase (found in the intermembrane space of mitochondria and required for inserting iron in tetrapyrroles) and in some bacterial dioxygenases.

## 2. Complex FeS Proteins

The redox properties of FeS clusters are utilized in multielectron transfer systems, which can be made more efficient by combining donors and acceptors into a single, often multimeric structure. This evolutionary advantage led to the appearance of complex FeS proteins, in which 4Fe–4S or 2Fe–2S clusters are combined with other inorganic or organic cofactors.

*Desulfovibrionaceae* and some *Clostridia* also contain what apparently amounts to a natural chimeric protein, in which a rubredoxin-like **domain** is linked to a domain that contains a binuclear, oxo-bridged iron center, such as those found in [ribonucleotide reductase](#) or in hemerythrin—hence the trivial names of these proteins: rubrerythrin and nigerythrin. Rubrerythrin was found to have ferroxidase activity, and an oxygen-scavenging function was hypothesized for these proteins. Yet another combination of metallic clusters has been found in *D. gigas* desulfoferredoxin, a protein in which a desulfiredoxin site is associated with an iron site having nitrogen ligands. No biological function is known for desulfiredoxin or desulfoferredoxin.

Nitrogenase components 1 and 2 represent the best-known examples of combinations of different and unusual metal–sulfur centers. In the dimeric component 2 (dinitrogenase reductase), a 4Fe–4S cluster (Fig. 1c) is bound at the dimer interface. Each of the subunits provides two of the coordinating cysteine residues, and the cluster structure and properties are modified upon [adenylation](#) of the dimer. The a<sub>2</sub>b<sub>2</sub> nitrogenase component 1 (the true dinitrogenase) contains a number of the so-called P-clusters (derived from 4Fe–4S cubanes) (Fig. 1e), again harbored jointly by adjacent subunits. In the same intersubunit fashion, dinitrogenase also binds FeMoCo, the metal cluster at which the multielectronic reduction of dinitrogen to ammonia occurs. FeMoCo is a modified double cubane structure, in which one of the corners is occupied by a Mo rather than a Fe atom, and is coordinated also by histidine and by an isocitrate molecule (Fig. 1f). Another

elementary process in bacterial physiology is hydrogen evolution, which uses an enzyme containing an Ni–FeS cluster, along with 4Fe–4S clusters.

FeS proteins containing multiple iron centers and organic cofactors (most typically flavins and pterins) abound in the electron transport chains of both prokaryotes and eukaryotes. These membrane-bound enzymes are typically made up of several subunits, with a variety of arrangements, and types of clusters and cofactors spanning a broad range of reduction potential values, often modulated through the binding of suitable effectors. Also, an active role for FeS clusters in energy transduction has been proposed. Soluble complex iron–sulfur enzymes play key roles in bacterial metabolism (glutamate synthase or anaerobic ribonucleotide reductase) and in the assimilation pathways of inorganic nutrients in plants (as in sulfite, nitrate, and nitrite reductases), often in conjunction with unique cofactors. A mammalian soluble protein of great complexity is xanthine oxidase, containing flavin, iron–sulfur, and a molybdopterin cofactor. Formylmethanofuran dehydrogenases are FeS proteins found in methyloprophs, which contain Mo or W and a pterin dinucleotide.

### 3. FeS Proteins Having No Electron-Transfer Function

There is a remarkable structural similarity between the Fe(Cys)<sub>4</sub> unit found in rubredoxin-like proteins and the Zn(Cys)<sub>4</sub> unit found in DNA-binding **zinc fingers**, and [transcription factors](#) containing a Fe(Cys)<sub>4</sub> unit have been reported for a number of organisms. No DNA-binding activity has been reported for rubredoxins themselves. Examples of FeS proteins displaying the ability of binding polynucleotides are not limited to rubredoxin-like proteins. *Escherichia coli* endonuclease III contains a single 4Fe–4S cluster, which apparently only plays a structural role.

To alleviate the toxicity of the byproducts of aerobic respiration, *E. coli* induces the synthesis of protective enzymes, and this induction is controlled by the regulatory FeS proteins SoxRS, OxyR, and ArcAB. ArcAB, Fnr, SoxRS, and OxyR function in concert, so that *E. coli* can optimize its energy production and growth rate. Fnr and SoxRS are [DNA-binding proteins](#), and they utilize 4Fe–4S and 2Fe–2S clusters as direct sensors of the redox environment (3).

Perhaps the best-known example of the regulatory function of FeS proteins is represented by the [RNA-binding proteins](#) involved in iron homeostasis in mammalian cells. When cells are depleted of iron, iron regulatory proteins IRP1 and IRP2 bind with high affinity to specific RNA stem-loop structures (**iron responsive elements**, IREs) in [messenger RNA](#) transcripts. Binding of IRPs to IRE stabilizes the mRNA coding for the transferrin receptor against degradation, while blocking the [translation](#) of the mRNA coding for the iron-storage protein, [ferritin](#). In the presence of iron, IRP1 assembles a 4Fe–4S cluster and loses its IRE-binding properties, whereas IRP2 is rapidly subject to **protein degradation**. Thus, the presence of iron allows synthesis of ferritin (as well as of d-aminolevulinate synthase, mitochondrial succinate dehydrogenase, and aconitase), while repressing the synthesis of the transferrin receptor and therefore diminishing iron import into the cell (4).

The role of iron in citrate isomerization catalyzed by aconitase was recognized many years ago, and aconitase was found to contain a 4Fe–4S cluster, which can be converted reversibly to an inactive 3Fe–4S form (Fig. 1d). This conversion is made easier by the absence of a coordinating cysteine residue to the “labile” iron atom. Also complex FeS proteins are not limited to redox reactions. A FAD and 4Fe–4S enzyme from *Clostridium aminobutyricum* catalyzes the reversible dehydration of 4-hydroxybutyryl-CoA to crotonyl-CoA, a reaction involving the cleavage of an unactivated C–H bond; similar activities have been reported for other bacteria.

### 4. Processing of FeS Proteins and Cluster Assembly

From a chemical standpoint, iron–sulfur clusters are among the easiest cofactors to assemble. Synthetic FeS complexes containing 1, 2, 4, or more iron atoms and an appropriate complement of



sulfide readily self-assemble in mixtures of iron, [thiol groups](#), and a suitable source of sulfide (1). This observation, plus the positive charge of the pyrite-like fundamental unit of FeS clusters, and evolutionary considerations, prompted some hypotheses on a role of FeS clusters in the emergence of life in the reducing and acidic environment near submarine hot springs in primordial oceans (5).

A different situation is encountered in the cell, where the toxicity of iron and sulfide must be accounted for, and where **protein folding** or assembly of multimeric proteins must occur, along with **protein targeting** of nuclear-encoded proteins to the appropriate compartments. In cells and organelles, sulfide may be derived from cysteine or from thiosulfate through appropriate enzymes, while the actual carriers of iron remain unknown. FeS proteins are imported into organelles as metal-free precursors, and cluster insertion and acquisition of a native structure are very often assisted by more or less complex scaffolding systems. These scaffolding systems are most relevant when assembly of interprotein multimetallic clusters (as in dinitrogenase) or of multiple-component, multiple-cofactor enzymes is required.

### Bibliography

1. H. Beinert, R. H. Holm, and E. Münck (1997) *Science* **277**, 653–659.
2. M. K. Johnson (1998) *Curr. Opin. Chem. Biol.* **2**, 173–181.
3. E. Hidalgo, H. Ding, and B. Dimple (1997) *Trends Biochem. Sci.* **22**, 207–210.
4. T. A. Rouault and R. D. Klausner (1996) *Trends Biochem. Sci.* **21**, 174–177.
5. G. Wächtershäuser (1992) *Prog. Biophys. Mol. Biol.* **58**, 85–201.

### Suggestions for Further Reading

6. R. H. Holm, P. Kennepohl, and E. I. Solomon (1996) Structural and functional aspects of metal sites in biology. *Chem. Rev.* **96**, 2239–2314.
7. D. M. Kurtz (1997) Structural similarity and functional diversity in diiron-oxo proteins. *J. Biol. Inorg. Chem.* **2**, 159–167.
8. H. Beinert, M. C. Kennedy, and C. D. Stout (1996) Aconitase as iron–sulfur protein, enzyme and iron-regulatory protein. *Chem. Rev.* **96**, 2335–2373.
9. R. Cammack (1992) Iron–sulfur clusters in enzymes: themes and variations. *Adv. Inorg. Chem.* **38**, 281–322.
10. D. H. Flint and R. M. Allen (1996) Iron–sulfur proteins with nonredox functions. *Chem. Rev.* **96**, 2315–2334.

### Isochores

The vertebrate [genome](#) is compartmentalized into compositionally homogeneous regions, the *isochores*, originally identified (1) as segments of DNA longer than 300 kb that are homogeneous in base composition over segments at least 3 kb long. Isochores can be arranged into several families distinguished by different base compositions, as shown initially by [density gradient centrifugation](#) of DNA molecules derived from isochores through the breaks occurring during DNA preparation. Typical of the genomes of most mammals, the human genome has two GC-poor isochore families, L1 and L2 (which have few G and C bases, primarily A and T), which together represent 62% of the genome. Two GC-rich families H1 and H2 represent, respectively, 22 and 9% of the genome, and a very GC-rich H3 family represents only 3–4%. The remaining 4% of the genome consists of satellite and **ribosomal** DNA (2).

The localization of specific genes to defined isochores indicates that genes are not distributed randomly in the human genome. The density of genes is low in GC-poor isochores, increases with increasing GC content in the H1 and H2 families, and reaches a maximum in the very GC-rich family H3, where the gene concentration is at least 20-fold higher than in the GC-poor isochores. The H3 family corresponds to **telomeric** bands of **metaphase** chromosomes, and the L1 + L2 families correspond to **Giemsa bands**, whereas reverse bands comprise both GC-poor and GC-rich isochores.

## 1. Isochores and Evolution

The compositional distributions of large genome fragments (greater than 100 kb, such as those comprising standard DNA preparations) of **exons** and their codon positions and of **introns** are correlated with each other. They represent compositional patterns and are very different between the genomes of cold- and warm-blooded vertebrates, mainly because the former are much less heterogeneous in base composition and never reach the high GC levels attained by the latter. Only small compositional differences are found among the genomes of either cold- or warm-blooded vertebrates. Compositional patterns allow defining two modes of genome evolution: a conservative mode with no compositional change and a transitional (or shifting) mode with compositional changes. The conservative mode is observed among both cold- and warm-blooded animals. The transitional mode comprises both major and minor compositional changes. In vertebrate genomes, the major changes are associated with GC-rich and very GC-rich isochores in mammalian and avian genomes. Thus, the compositional changes reflect genome **phenotypes** differing not only between warm and cold-blooded vertebrates (eg, *Xenopus* versus mammals and birds), but also between mammals and birds (human and mouse versus chicken), and even among mammals. For a comprehensive review, see refs. [3](#) and [4](#).

### Bibliography

1. G. Macaya, J-P Thiéry, and G. Bernardi (1976) *J. Mol. Biol.* **108**, 237–254.
2. D. Mouchiroud et al. (1991) *Gene* **100**, 181–87.
3. G. Bernardi (1993) *Mol. Biol. Evol.* **10**, 186–204.
4. G. Bernardi (1995) *Ann. Rev. Genetics* **29**, 445–476.

## Isochromosome

Isochromosomes are chromosomes that contain multiple identical chromosomal arms. They may arise during the [cell cycle](#) when the two sister [chromatids](#) of an **acrocentric** or **telocentric chromosome** remain fused at their [centromeres](#), leading to the creation of a **metacentric** isochromosome. Alternatively, if [homologous chromosomes](#) fail to pair during the first division of **meiosis**, the centromere of the univalent chromosome might split transversely, i.e., between the two arms, which results in two telocentric chromosomes that could be segregated during cell division. These telocentric chromosomes are unstable, and they may fuse to form an isochromosome in which one whole chromosome arm is deleted and the other duplicated. Isochromosomes in humans have been used to demonstrate that the centromere has variable activity, dependent on chromosomal context. For example, an isochromosome X containing two centromeres can be stably inherited through **cell division** if one is inactivated.

### Suggestion for Further Reading

M. S. Clark and W. J. Wall (1996) *Chromosomes. The Complex Code*, Chapman and Hall, London, pp. 43–44.

## Isocratic Elution Chromatography

Continuous elution of analytes in a **chromatographic** separation with only one mobile phase (a single solvent or a constant-composition solvent mixture) throughout the separation is known as isocratic-elution chromatography. This method is frequently very time-consuming, and the analytes are eluted as broad bands if they are strongly adsorbed on the column. Therefore, it is common to follow an initial period of isocratic separation with an elution gradient, in which the polarity of the initial mobile phase is progressively changed by the continuous programmed addition of increasing quantities of a second mobile phase to speed up the rate of elution and to sharpen eluent zones that are trailing or fronting. A description of the theory of isocratic-elution chromatography is beyond the scope of this article, and interested readers are directed to the excellent chapter written by Jandera and Churáček (1).

### Bibliography

1. P. Jandera and J. Churáček (1985) *Gradient Elution in Column Liquid Chromatography: Theory and Practice (Journal of Chromatography Library, Vol. 31)*, Elsevier, Amsterdam, pp. 1–55.

## Isoelectric Focusing

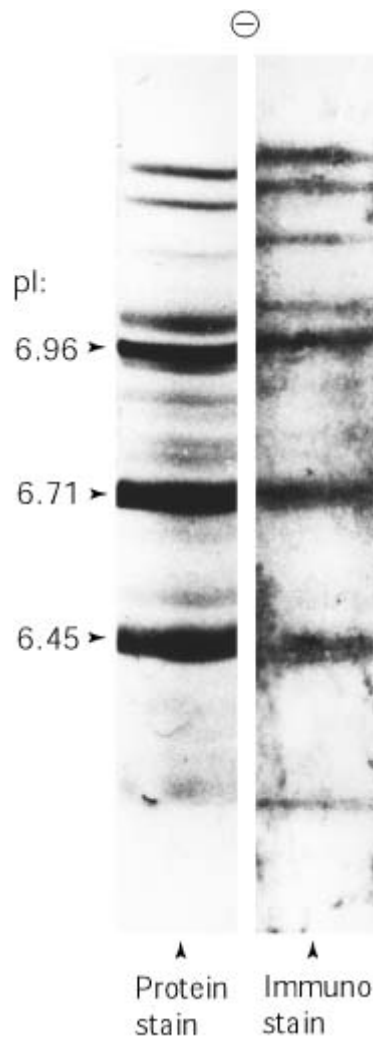
Conventional isoelectric focusing (IEF) in soluble, [amphoteric buffers](#) known as carrier ampholytes (CAs) was reported by Svensson-Rilbe in 1961 (1-3). In this technique, the macromolecules of interest are subjected to electrophoresis in a continuous pH gradient. They migrate electrophoretically toward the pH of their [isoelectric point](#) (pI), where they have zero mobility. Therefore, the system comes to equilibrium, with each species of macromolecule tightly focused at its pI; any [diffusion](#) away is reversed by the electrophoretic migration. At about the same time, Meselson et al. (4) described **isopycnic centrifugation**, a related high-resolution technique that is another member of a family called isoperichoric focusing by Kolin (5). Unlike conventional **chromatographic** and **electrophoretic** techniques, where peaks of material are constantly dissipated by diffusion, IEF and isopycnic centrifugation have built-in mechanisms opposing this. As the analyte reaches an environment (the *perichoron* in Greek) in which its physicochemical parameters are equal (*iso*) to those of its surroundings, it focuses, or condenses, in an ultrathin zone, kept stable and sharp in time by two opposing force fields: diffusion (tending to dissipate the zone) and external fields (voltage gradients in IEF or centrifugal fields in isopycnic centrifugation) forcing the “escaping” analyte back into its “focusing” zone.

Both conventional IEF using carrier ampholytes and the new version using immobilized pH gradients (IPG) (6) rely on the equation for resolving power (expressed as  $DpI$ , i.e., the difference in isoelectric points between a protein and a just resolved, nearest contaminant) (7):

$$\Delta pI = 3.17 \sqrt{(D_t [\partial(\text{pH})/\partial x] / E [\partial u/\partial(\text{pH})])}$$

where  $D_t$  is the protein translational diffusion coefficient,  $E$  is the voltage gradient (V/cm) applied,  $\partial(\text{pH})/\partial x$  is the slope of the pH gradient along the separation axis, and  $\partial u/\partial(\text{pH})$  is the titration curve of the protein in terms of its mobility as a function of pH. Experimental conditions minimizing DpI offer the greatest resolution, which are shallow pH gradients and high voltage gradients. The best DpI attainable is of the order of 0.01 pH unit in CA-IEF but 0.001 in IPG, where more shallow pH gradients and higher voltages are possible. Consequently, complex mixtures of macromolecules can be resolved into very many species, and what appears to be a single species by many separation techniques is frequently resolved into several by IPG (Fig. 1). A brief description of both methods will be given.

**Figure 1.** An example of the separations possible with isoelectric focusing. This soluble form of the **receptor** for **epidermal growth factor** (EGF) gave a single band on **SDS-PAGE**, but it was resolved by IPG into three major isoforms and several minor bands. The pH gradient was from pH 5 to pH 8, in a 10-cm-long 5% T, 4% C **polyacrylamide gel** equilibrated with 30% ethylene glycol. A total of 50  $\mu\text{g}$  of protein was loaded at the anode, and focusing was carried out with the cathode uppermost. The resulting gels were stained with **Coomassie Brilliant Blue R-250** (in ethanol, acetic acid, and copper sulfate) on the left and by Western blotting on the right. The pI values of the major components are indicated. (From Ref. 18, with permission.)

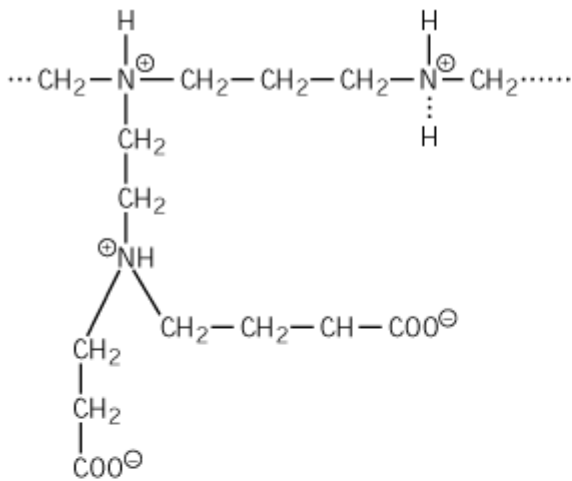


## 1. Conventional Isoelectric Focusing in Soluble, Amphoteric Buffers

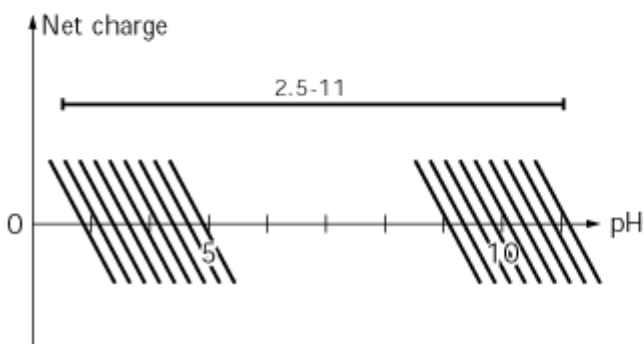
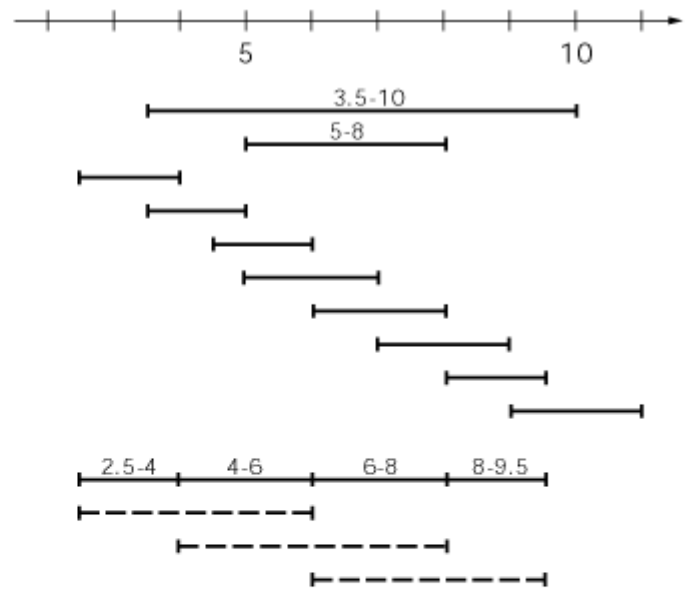
Carrier ampholytes set up the pH gradient for IEF by distributing themselves electrophoretically at their isoelectric points. For a range of pH values, an appropriate mixture of ampholytes is necessary. They must not only be amphoteric, but must also have suitable buffering power and conductivity at their pI (7). The hallmark of a carrier ampholyte is the absolute value of the difference between its pI and the  $pK_a$  of its buffering groups; the smaller this value, the greater its conductivity and buffering capacity at its pI (8). Very few natural compounds have this property.

A breakthrough came in 1964 with Versterberg's synthesis of Ampholines (trade name of the commercial product from LKB). They are "oligoprotic amino carboxylic acids, each containing at least four weak protolytic groups, at least one being a carboxyl group and at least one a basic nitrogen atom, but no peptide bonds" (9). Figure 2 reports a hypothetical general structural formula of these ampholytes. Very narrow pH ranges of ampholytes are produced, except for the very extreme pH 2.5 to 4 and pH 9 to 11 ranges, by subfractionating the "wide pH-range" mixture in multicompartiment electrolyzers. The synthesis of Ampholines is a genuine "chaotic" process, in that very heterogeneous mixtures of different oligoamines (typically pentaethylene hexamine, tetraethylene pentamine, and triethylene tetramine, including their branched isomers) were reacted with an  $\alpha$ - $\beta$  unsaturated acid (typically acrylic acid) at an appropriate ratio (usually 2 N atoms/1 carboxyl group);  $\beta$ -propionic acid residues result. By this synthetic approach, >600 ampholyte species could be generated in the pH 3 to 10 interval. Subsequently, Pharmalyte species containing >5000 chemically distinct amphoteres in the pH 2.5 to 11 interval were produced. Almgren (10) had predicted theoretically (assuming equimolar distribution of the CA species, even distribution of their DpI along the pH scale, and the same  $DpK_a$ ) that only 30 amphoteres/pH unit had to be present to generate a stepless pH course. Because the oligoamines prepared had  $pK_a$  values well distributed along the pH scale, they could provide the buffering power and conductivity required for IEF. Additionally, the extra  $pK_a$  value of the "titrant" acrylic acid grafted onto the oligoamino backbone endowed acidic carrier ampholytes with extra buffering power. This is one of the reasons why focusing with CA buffers is usually more successful in acidic pH ranges than alkaline. The other reasons are that this synthesis generates fewer alkaline ampholytes and that open-face IEF gel slabs at pH >8 absorb atmospheric  $CO_2$ , unless submerged under a thin layer of light paraffin oil.

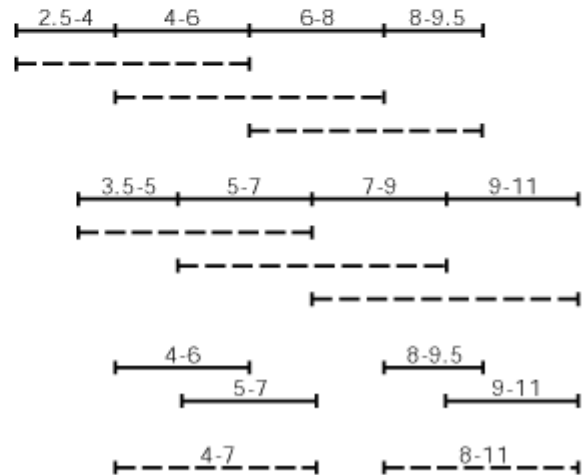
**Figure 2.** Composition of ampholines. (a) Representative chemical formula (aliphatic oligo-amino, oligo-carboxylic acid) (b) Portions of hypothetical titration curves of Ampholines. (c) Different pH cuts for wide- and narrow-range carrier ampholytes. (Courtesy of LKB Producter AB.)



(a)



(b)



(c)

CA-IEF offers unrivaled resolving power in a vast number of separations and is now a standard technique in molecular biology. In addition to its current uses, at least two particular applications are worth mentioning: (a) its use in **two-dimensional (2-D) electrophoresis** to generate 2-D maps ([11](#), [12](#)) and (b) for generating titration curves of macromolecules. In the latter case, the CA-IEF slab gel is focused to generate the desired pH gradient. The slab is then turned 90° and the sample is loaded in a long, thin trench spanning from anode to cathode. The sample is then subjected to electrophoresis perpendicular to the stationary pH gradient. As a result, each protein exhibits its own pH/mobility curve in the gel slab ([13](#)).

There are, however, a number of problems with the CA-IEF gels: They have uncertain chemical environments of low ionic strength and uneven buffering capacity and conductivity; cathodic drift results in extensive loss of proteins at the gel cathodic extremity upon prolonged runs; and the pH gradients possible are limited by the nature of the carrier ampholytes available. The low ionic strength within the gel often induces near-isoelectric precipitation and smears of proteins, even in analytical runs with small amounts of protein. For all these reasons, the technique of IPG was launched in 1982 ([6](#)).

## 2. Immobilized pH Gradients

Immobilized pH gradients (IPGs) are unique in that they are based on the insolubilization of the entire set of buffers and titrants responsible for generating and maintaining the pH gradient along the separation axis. Contrary to CA-IEF, where a multitude of soluble, amphoteric buffers is expected to create and sustain the pH gradient during the electrophoretic run, IPGs are based on only a few well-defined, nonamphoteric buffers and titrants, able to perform in a highly reproducible manner.

Advantages of the IPG technique are: (a) increased resolution (by at least one order of magnitude); (b) unlimited stability of the pH gradient; (c) flexibility in the choice of pH interval; (d) increased loading ability; (e) high reproducibility; (f) minimal distortion by salts in the sample; (g) full control of pH, buffering capacity and ionic strength; and (h) easy separation of sample from buffering ions in preparative runs.

In IPGs, the pH gradient exists prior to the IEF run and is copolymerized, and thus immobilized, within the [polyacrylamide](#) matrix. This is achieved by using, as [buffers](#), a set of six nonamphoteric weak acids and bases, having the general chemical composition  $\text{CH}_2\text{CH}-\text{CO}-\text{NH}-\text{R}$ , where R denotes either two different weak carboxyl groups, with  $\text{p}K_a$  values of 3.6 and 4.6, or four tertiary amino groups, with  $\text{p}K_a$  values of 6.2, 7.0, 8.5, and 9.3 (available under the trade name Immobiline from Pharmacia-LKB). A more extensive set, comprising 10 chemicals (with the addition of a  $\text{p}K_a$  3.1 acidic buffer, a  $\text{p}K_a$  10.3 alkaline species, and two strong titrants, a  $\text{p}K_a$  1 acid and a  $\text{p}K_a >12$  quaternary base), is available as “pI select” from Fluka AG, Buchs, Switzerland. During gel polymerization, these buffering species are efficiently incorporated into the gel (84% to 86% conversion efficiency at 50°C for 1 h). Immobiline-based pH gradients can be cast in the same way as conventional gradient polyacrylamide gels, by using a density gradient to stabilize the Immobiline concentration gradient, with the aid of a standard two-vessel gradient mixer (see [Transverse Gradient Gel Electrophoresis \(Tgge\)](#)). The buffers are not amphoteric, but are bifunctional: At one end of the molecule is located the buffering (or titrating) group, and at the other is the acrylic double bond, which will participate in the polymerization process. Acidic and alkaline Immobilines have different temperature coefficients ( $\partial\text{p}K_a/\partial T$ ), so temperature affects Immobiline pH gradients, as do the ionic strength and additives that change the water structure ([chaotropic](#) agents, such as [urea](#)) or lower its dielectric constant. The largest changes are due to the presence of urea: Acidic Immobilines increase their  $\text{p}K_a$  values in 8 M urea by as much as 0.9 pH units, and basic Immobilines increase their  $\text{p}K_a$  values by only 0.45 pH unit. [Detergents](#) in the gel (up to 2% w/v) do not alter the Immobiline  $\text{p}K_a$ , suggesting that acidic and basic groups attached to the gel are not incorporated into surfactant micelles. For generating extended pH gradients, one should use two additional strong titrants having  $\text{p}K_a$  values well outside the desired pH range: QAE (quaternary amino ethyl)-acrylamide ( $\text{p}K_a >12$ ) and AMPS (2-acrylamido-2-methyl propane sulfonic acid,  $\text{p}K_a \approx 1$ ).

With the IPG technology, ultranarrow (0.1), narrow (1.0), and extended (up to 8.0) pH gradients can be engineered with high precision and reproducibility. Recipes have been tabulated listing 58 1-pH-unit-wide gradients, separated by 0.1 pH unit increments, starting with the 3.8 to 4.8 pH interval and ending with the pH 9.5 to 10.5 range ([14](#)). If a narrower pH gradient is needed, it can be derived from any of the 58 pH intervals tabulated by a simple linear interpolation of intermediate Immobiline molarities. Recipes for gradients with pH spans from 2 pH units up to 6 pH unit spans are available ([15](#)). All the formulations are normalized to give the same average value of buffering power of  $3\text{mequiv}\cdot\text{L}^{-1}\text{pH}^{-1}$ , adequate for producing highly stable pH gradients. For pH intervals covering  $>4\text{pH}$  units, the best solution is to mix a total of 10 Immobiline species, 8 of them buffering ions and two the strong acidic and basic titrants. Nonlinear (eg, concave and convex exponential, as well as sigmoidal) IPG gradients can be generated and optimized ([16](#)).

After casting, IPG gels should be washed extensively, so as to remove nonpolymerized material, salts, and polymerization catalysts. This also produces clean matrices, devoid of toxic materials (eg, unreacted acrylic double bonds) that could modify proteins by reacting with [thiol groups](#) and terminal [amino groups](#). As IPG gels are cast onto plastic supports and are rather thin (0.5 mm) and porous, they can be stored dry and then reswollen with any desired additives just prior to use. IPG are very reproducible and effective at very alkaline pH values, up to pH 12, where CA-IEF would simply fail; histones could thus be focused to steady state ([17](#)).

## Bibliography

1. H. Svensson (1961) *Acta Chem. Scand.* **15**, 325–341.
2. H. Svensson (1962) *Acta Chem. Scand.* **16**, 456–466.
3. H. Svensson (1962) *Arch. Biochem. Biophys. Suppl.* **1**, 132–140.
4. M. Meselson, F. W. Stahl, and J. Vinogradov (1957) *Proc. Natl. Acad. Sci. USA* **43**, 581–585.
5. A. Kolin (1977) In *Electrofocusing and Isotachopheresis* (B. J. Radola and D. Graesslin, eds.), de Gruyter, Berlin, pp. 3–33.
6. B. Bjellqvist, K. Ek, P. G. Righetti, E. Gianazza, A. Görg, W. Postel, and R. Westermeier (1982) *J. Biochem. Biophys. Methods* **6**, 317–339.
7. H. Rilbe (1973) *Ann. N.Y. Acad. Sci.* **209**, 11–22.
8. P. G. Righetti, M. Fazio, and C. Tonani (1988) *J. Chromatogr.* **440**, 367–377.
9. O. Vesterberg (1973) *Ann. N.Y. Acad. Sci.* **209**, 23–33.
10. M. Almgren (1971) *Chem. Scripta* **1**, 69–75.
11. D. S. Young and N. G. Anderson (Guest Eds.) (1982) (Special Issue on 2-D Gel Electrophoresis) *Clin. Chem.* **28**, 737–1092.
12. L. Anderson and N. Anderson (Guest Eds.) (1984) (Special Issue on 2-D Gel Electrophoresis) *Clin. Chem.* **30**, 1897–2108.
13. P. G. Righetti (1996) In *Current Protocols in Protein Science* (J. E. Coligan, B. M. Dunn, H. L. Ploegh, D. W. Speicher, and P. T. Wingfield, eds.) Vol. **1**, suppl. 3, Wiley, New York, pp. 7.3.13–7.3.26.
14. P. G. Righetti, A. Bossi, and C. Gelfi (1998) In *Gel Electrophoresis of Proteins. A Practical Approach* (B. D. Hames ed.), IRL Press, Oxford, U.K., pp. 127–187.
15. P. G. Righetti (1989) In *Protein Structure, A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, U.K., pp. 23–63.
16. E. Giuffreda, C. Tonani, and P. G. Righetti (1993) *J. Chromatogr.* **630**, 313–323.
17. A. Bossi, C. Gelfi, A. Orsi, and P. G. Righetti (1994) *J. Chromatogr. A* **686**, 121–128.
18. E. Wenisch, P. G. Righetti, and W. Weber (1992) *Electrophoresis* **13**, 668–673.

## Suggestions for Further Reading

19. P. G. Righetti (1983) *Isoelectric Focusing: Theory, Methodology and Applications*, Elsevier, Amsterdam, pp. 1–386.
20. P. G. Righetti (1990) *Immobilized pH Gradients: Theory and Methodology*, Elsevier, Amsterdam, pp. 1–400.
21. H. Rilbe (1996) *pH and Buffer Theory. A New Approach*, Wiley, Chichester, pp. 1–192.

## Isoelectric Point (PI)

The isoelectric point is the pH at which the total net charge of an [amphoteric](#) molecule becomes zero. The isoelectric point depends on the ionic strength of the solvent, the counterion species, and the solute concentration, whereas the [isoionic point](#) is the isoelectric point in the absence of counterions and at low solute concentrations. The isoelectric point is readily measured as the pH at which the molecule is immobile in an electrophoresis experiment (see [Isoelectric Focusing](#)). The isoelectric point approaches the isoionic point at low ionic strength, when all bound ions are removed. At the isoelectric point, [amphoteric](#) molecules aggregate and become insoluble because there are no



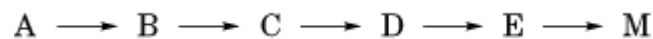
electrostatic repulsions between molecules with the same net charge.

## Isofunctional Proteins

In some instances, it is necessary for an organism to produce two or more forms of the same protein, *isofunctional proteins* that have the same function, so as to optimize their metabolic regulation. Cells generally have two primary mechanisms for regulating their metabolic reaction rates in response to changes in the levels of various metabolites: (i) **allosteric** inhibition or activation of the activities of certain [enzymes](#) and (ii) regulation of the biosynthesis of the enzymes.

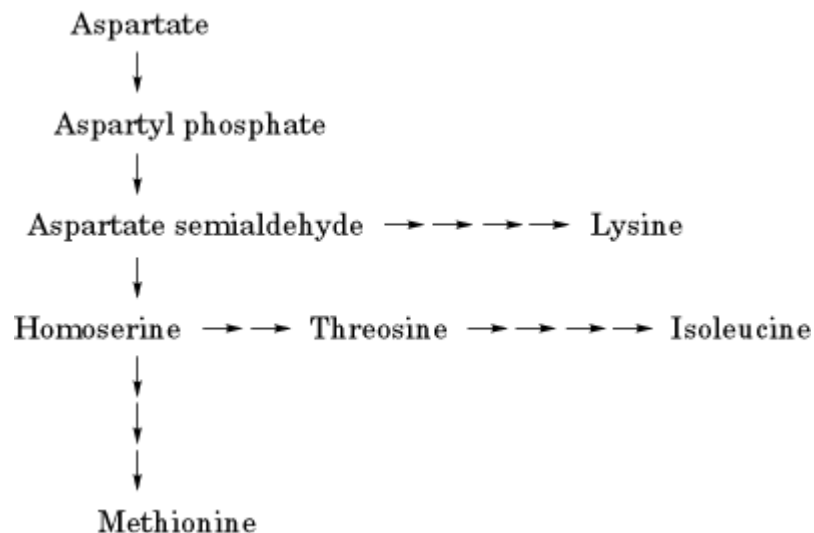
### 1. Regulation of Biosynthetic Pathways

Consider a linear sequence of reactions



leading to the synthesis of the essential metabolite M. In the presence of excess M in the intracellular pool, the enzyme carrying out the conversion  $A \rightarrow B$  will often be subject to **feedback inhibition** by M if, as is generally the case, this enzyme is an allosteric protein. M may also **repress** the synthesis of the enzymes of the entire pathway. These control phenomena are economical and efficient ways of keeping the pool size of M relatively constant and do not cause any detrimental side effects.

Now consider a branched reaction biosynthetic sequence, such as that of the aspartate family of amino acids:



Regulation of this pathway by feedback inhibition of the first enzyme, aspartokinase, or by repressing expression of the enzymes of the entire pathway by the various end products (lysine, threonine, isoleucine, and methionine) could not ensure balanced production of all end products when one or just a few are present in excess. The solution to this potential dilemma differs according to the organism studied, but that of *Escherichia coli* illustrates the use of isofunctional proteins.

### 2. Three Aspartokinases in *Escherichia Coli*

Of all the natural amino acids occurring naturally as protein constituents, only two, L-lysine and L-threonine, influence the aspartokinase activity of crude extracts of *E. coli*. Neither totally inhibit the activity, and the total inhibition, when they are present simultaneously, is the sum of that observed for each independently because one aspartokinase is sensitive to inhibition by lysine and another is sensitive to threonine (1). The synthesis of the enzyme aspartokinase is repressed by lysine, and that of the latter is subject to a multivalent repression by threonine and isoleucine ((1),(2)). A third aspartate kinase is inhibited by none of the end-products of the pathway, but its synthesis is repressed by methionine (see [Methionine Regulon](#)). The three aspartokinases are highly homologous and share a common ancestor (3-5).

The rational solution in *E. coli* to the dilemma posed by the regulation of branched biosynthetic pathways was the synthesis of multiple, *isofunctional* enzymes that catalyze the synthesis of the common precursor, aspartyl phosphate. Each of these enzymes is independently subject to feedback inhibition and to repression by different end-product metabolites. Another example is encountered in the *E. coli* biosynthetic pathway of aromatic amino acids. Three distinct isofunctional 3-deoxy-d-arabino-heptulosonate-7-phosphate synthetases are, respectively regulated by allosteric inhibition and/or repression by the three end products, phenylalanine, tyrosine, and tryptophan (6) (see [TRP Operon](#)).

#### Bibliography

1. E. R. Stadtman et al. (1961) *J. Biol. Chem.* **236**, 2033–2038.
2. G. N. Cohen and J-C. Patte (1963) *Cold Spring Harbor Symp. Quant. Biol.* **28**, 513–515.
3. M. Cassan et al. (1986) *J. Biol. Chem.* **261**, 1052–1057.
4. M. Katinka et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5730–5733.
5. M. M. Zakin et al. (1983) *J. Biol. Chem.* **258**, 3028–3031.
6. A. J. Pittard (1996) In *Escherichia coli and Salmonella; Cellular and Molecular Biology* (F. C. Neidhardt, ed.), ASM Press, Washington, DC, Vol. I, p. 460.

#### Suggestion for Further Reading

7. G. N. Cohen (1994) *Biosyntheses*, Chapman and Hall, New York and London, pp. 147–200.

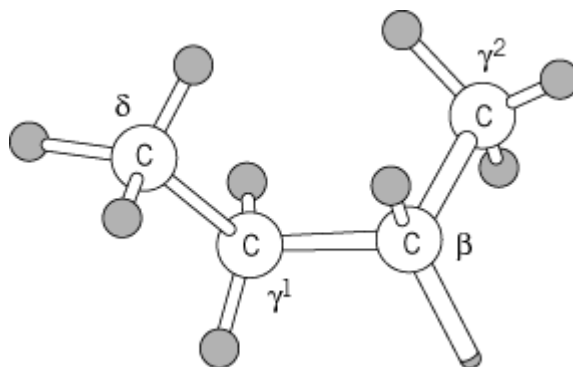
### Isoionic Point (Isoionic pH)

The isoionic point is the pH at which the total net charge of an [amphoteric](#) molecule becomes zero in pure [water](#). The ideal isoelectric point is that in the complete absence of salt and at very low solute concentrations. No counterions should be present except for the protons and the hydroxyl ions of the solvent. Any metal ions bound to the macromolecules should be removed. The isoionic point can be estimated from the  $pK_a$  values of the individual ionizable groups.

### Isoleucine (Ile, I)

The [amino acid](#) isoleucine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to three **codons**—AUU, AUC, and AUA—and represents approximately 5.2% of the residues of the proteins that have been characterized. The isoleucyl residue incorporated has a mass of 113.16 Da, a **van der Waals volume** of 124 Å<sup>3</sup>, and an [accessible surface](#) area of 182 Å<sup>2</sup>. Ile residues are frequently changed during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [valine](#), [leucine](#), and [threonine](#) residues.

The Ile side chain is [nonpolar](#) with no functional or reactive groups:



Consequently, Ile is one of the most **hydrophobic** amino acid residues, and 60% of Ile residues are completely buried in native [protein structures](#). Ile residues do not favor the **alpha-helical** conformation in model peptides, and they occur most frequently in the [beta-sheet](#) type of **secondary structure** in native proteins. Note that the Ile side chain has a center of asymmetry, at C<sup>b</sup>, and that only the one isomer occurs naturally and is incorporated into proteins.

#### Suggestions for Further Reading

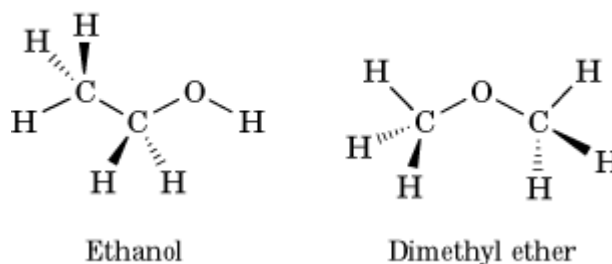
T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Isomer

The word isomer indicates that two molecules share the same elemental formula, being derived from the Greek isos, for “equal” and meros for “part.” The importance of isomers in chemistry and biochemistry derives from the historical importance of elemental analysis in determining structural identity. A simple example of isomers is ethanol and dimethyl ether, both compounds having the elemental formula C<sub>2</sub>H<sub>6</sub>O (Fig. 1).

**Figure 1.** The structural differences between the isomers ethanol and dimethyl ether. These demonstrate how two molecules having the same elemental formulas can have different covalent structures.

### Structural isomers



Since two molecules of common elemental formula can differ in many ways, many times the mode of differentiation is indicated. Ethanol and dimethyl ether are structural isomers, as are glucose and fructose, because the number and type of bonds are different. Geometric isomers differ in geometrical arrangement of bonds, eg, fumarate and maleate. [Stereoisomers](#), [diastereomers](#), and [enantiomers](#) are specialized cases of isomerism, where the differences in the molecules are determined solely by the spatial orientation of the bonds. [Tautomers](#) are another specialized class of isomers that are differentiated by their ability to equilibrate rapidly.

A common mistake is to regard conformers as isomers. Conformers are generated by rotation about a single bond and not a difference in the structure or configuration of the bonds present, which is required for isomerism ([1](#)).

### Bibliography

1. E. L. Eliel et al. (1967) *Conformational Analysis*, Wiley-Interscience, New York, p. 1.

### Suggestions for Further Reading

2. K. Mislow (1962) In *Comprehensive Biochemistry, I: Atomic and Molecular Structure* (M. Florkin and E. H. Stotz, eds.), Elsevier, Amsterdam, pp. 192–241.
3. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), New Comprehensive Biochemistry 3, Elsevier, Amsterdam, pp. 1–48.

## Isomerases

Isomerases form a broad group of [enzymes](#) that catalyze intramolecular rearrangements. The subdivisions within this class are:

- *Racemases*. Act on [amino acids](#) and hydroxy acids with **chiral** centers.
- *Epimerases*. Act on carbohydrates and their derivatives, as well as other compounds with chiral centers.
- *Isomerases*. Catalyze [cis trans isomerization](#) reactions, keto-enol transformation, the interconversion of aldoses and ketoses, shifts of double bonds, and the rearrangement of [disulfide bonds](#).
- *Mutases*. Transfer acyl, phosphoryl, amino, and other groups from one intramolecular position to another

## Isomorphous Replacement

To determine the structure of a macromolecule using [X-ray crystallography](#), an electron density map of the structure is computed as a Fourier summation:

$$\rho(xyz) = \frac{1}{V} \sum_h \sum_k \sum_\ell |F(hk\ell)| \times \exp[-2\pi i(hx + ky + \ell z) + i\alpha(hk\ell)] \quad (1)$$

where  $\rho(xyz)$  is the electron density at position  $x, y, z$  in the [unit cell](#) and  $F(hk\ell)$  is the amplitude of the [structure factor](#) for reflection  $h k \ell$ . Apart from correction factors,  $|F(hk\ell)| = \sqrt{I(hk\ell)}$ , where  $I(hk\ell)$  is the intensity measured for the reflection  $h k \ell$ . The phase angle for this reflection is  $\alpha(hk\ell)$ . The value of  $\alpha(hk\ell)$  cannot be derived directly from the experimental data for macromolecular structures (see [Phase Problem](#)). Several methods do exist, however, to find the values for the phase angles. For macromolecular structures, such as proteins, isomorphous replacement is the oldest and the most general method. It was applied for the first time by Perutz to solve the structure of [hemoglobin](#) (1).

It requires the X-ray diffraction pattern of native protein crystals and also of one or more derivatives that contain heavy atoms attached to the native protein. The attachment must be isomorphous, which means that the structure of the protein itself should not be affected by the binding of the heavy atom reagents. A first check for isomorphism is comparing the [unit cell](#) dimensions of the native and derivative crystals. An easy way to introduce heavy atoms is to soak a native crystal in an appropriate solvent that contains the heavy atom. It uses the large voids and channels filled with the solvent that are present between protein molecules in a crystal. Small reagents easily diffuse into the crystal and reach reactive sites on all of the molecular surfaces.

After processing the diffraction data, the position of the heavy atoms in the unit cell must be found. This can be done with a difference Patterson summation which has as coefficients  $(\Delta F)$

$\Delta F = (F_{PH} - F_P)$  and all phase angles equal to zero:

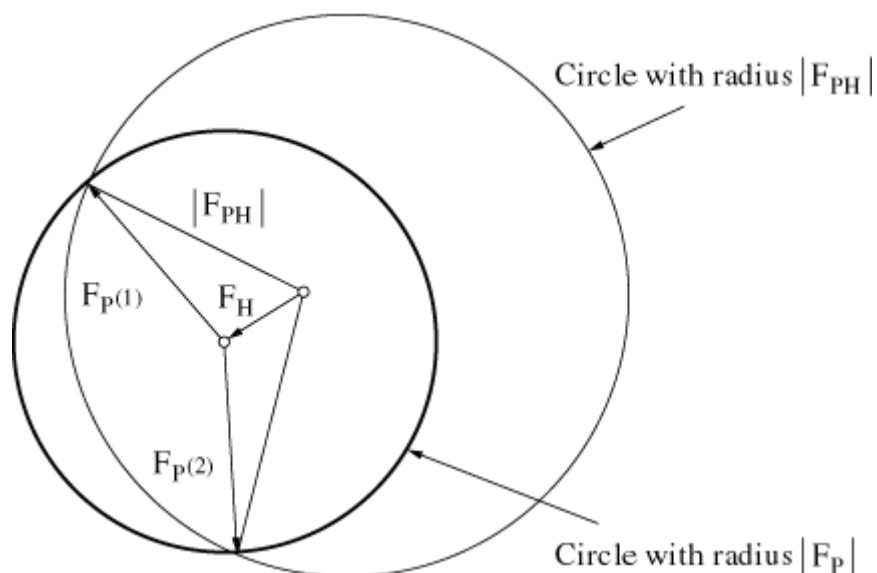
$$P(uvw) = \frac{1}{V} \sum_h \sum_k \sum_\ell \{\Delta F(hk\ell)\}^2 \cos[2\pi(hu + kv + \ell w)]$$

where  $u, v$ , and  $w$  are coordinates in the Patterson cell [see [Patterson Map](#)].  $F_{PH}$  and  $F_P$  are the amplitudes of the structure factor for one reflection of the derivative and the native structure, respectively.

After determining and refining the heavy atom position and parameters, the heavy atom contribution of the vector  $\mathbf{F}_H$  to the structure factor can be calculated. The knowledge of  $F_{PH}$ ,  $F_P$  and  $\mathbf{F}_H$  is sufficient to determine the phase angle of the reflection from the crystal that contains only the protein. In principle, this is done by a circle construction method due to Harker. In this method, the structure factors are regarded as vectors in a plane (Argand diagram). The length of such a vector is equal to the amplitude of the structure factor and its direction is given by the phase angle of the reflected beam. For each reflection, a circle is drawn whose radius is  $F_P$  (Fig. 1). Next, the vector  $\mathbf{F}_H$

is pointed to the center of this  $F_P$  circle. A second circle is drawn with radius  $F_{PH}$  whose center is at the origin of vector  $F_H$ . The points of intersection of the two circles are possible end points of the structure factor of the protein, because only at those positions is the principle of isomorphous replacement obeyed:  $F_P + F_H = F_{PH}$ .

**Figure 1.** The Harker construction for determining the protein phase angle in the isomorphous replacement method.



The choice between the two possible vectors,  $F_P$  (1) and  $F_P$  (2), can be made by repeating the procedure with a second, or even further, heavy atom derivatives. This is called multiple isomorphous replacement (MIR). Because of experimental errors and nonideal isomorphism, the correct phase angle may differ somewhat from the exact point of intersection. Therefore, weighting probabilities are introduced in most procedures for calculating the phase angle, depending on the degree of closure of the vector triangle  $F_P + F_H = F_{PH}$  (2). With only one derivative, the method is called single isomorphous replacement (SIR). With SIR, the protein electron density map is calculated from  $F_P$  (1) +  $F_P$  (2). It is expected that the correct values for the many  $F_P$ 's contribute to a correct electron density map and the incorrect values contribute to noise.

#### Bibliography

1. D. W. Green, V. M. Ingram, and M. F. Perutz (1954) Proc. Roy. Soc. A **225**, 287–307.
2. D. M. Blow and F. H. C. Crick (1959) Acta Crystallogr. **12**, 794–802.

#### Suggestion for Further Reading

3. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York, Chap. "7".

#### Isopropyl-1-Thio-β-D-Galactopyranoside

*Escherichia coli* grown in lactose (1-4-β-galactosyl-glucose) produces about a thousand times more **β-galactosidase** than when grown in glucose. What happens? The [kinetics](#) of the process could be revealing, but lactose is a substrate of β-galactosidase. Lactose is hydrolyzed into glucose and galactose by this [enzyme](#), and it is also isomerized. In fact, one of the isomerization products, *allolactose* (1-6-β-galactosyl-glucose) is the actual **inducer**. Lactose itself does not induce! Therefore, lactose is badly suited for the kinetic study of this induction process.

What would possibly be an ideal inducer that is not hydrolyzed by β-galactosidase? The German biochemist Otto Meyerhof had left a collection of glycosides in Paris when fleeing the German army in 1940. One of them was phenyl-1-thio-β-D-galactopyranoside, synthesized in 1912 by Konrad Delbrück in Emil Fischer's laboratory in Berlin. In contrast to ordinary glycosides, it is stable under extreme conditions. Jacques Monod found the collection after the war. He picked up the idea and interested two German chemists, Dietmar Türk and Burckhardt Helferich in Bonn, in synthesizing a whole set of 1-thio-β-D-galactopyranosides. The best inducer was isopropyl-1-thio-β-D-galactopyranoside (IPTG). Because IPTG is not measurably hydrolyzed by β-galactosidase, it was called a gratuitous inducer. Using IPTG, it could be shown that induced β-galactosidase synthesis starts three minutes after addition of the inducer and continues from then on at with a constant rate (1) (see [Lac Operon](#)). Later, **radiolabeled** IPTG was used to isolate Lac repressor (see [Lac Repressor](#)) (2). See [Lac Operon](#) for further details and reference 3 for its general history.

### Bibliography

1. J. Monod (1956) In *Enzymes: Units of Biological Structure and Function*, Academic Press, New York, pp. 7–28.
2. W. Gilbert and B. Müller-Hill (1966) Proc. Natl. Acad. Sci. USA **56**, 1891–1898.
3. B. Müller-Hill (1996) *The lac Operon. A Short History of a Genetic Paradigm*, de Gruyter, Berlin.

### Isoschizomer

*Isoschizomers* are distinct, different restriction endonucleases (see [Restriction Enzymes](#) and [Restriction–Modification Systems](#)) that recognize the same DNA sequence and cleave at the same position in the DNA. Isoschizomers are derived from different organisms and sometimes have different amino acid sequences. For example, R•*EcoRI* (from *Escherichia coli*) and R•*RsrI* (from *Rhodobacter sphaeroides*) are isoschizomers that both recognize the duplex DNA sequence G<sup>^</sup>AATTC and cleave (designated by the caret symbol “<sup>^</sup>”) the phosphodiester bond after the G on each strand (G<sup>^</sup>AATTC). Another isoschizomeric pair are the type II enzymes *Ksp632I* and *Ear I*, which recognize CTCTTC and cleave (1/4) to the right side (3′) of the sequence (1). The (1/4) designation indicates that the (5′–3′) strand is cleaved one nucleotide residue to the right of the canonical sequence and that the complementary strand is cleaved on the same side 4 nucleotide residues beyond the site.

Except in a few cases, like the R•*EcoRI*/R•*RsrI* pair, isoschizomers show little amino acid sequence similarity to one other. DNA methyltransferases (MTases; see [Methyltransferase, DNA](#)), the companion enzymes of endonucleases, isolated from different organisms and having the same DNA sequence and methylation specificity, are referred to as *isoprostomers* or *isomethylators* (2). For example, M•*EcoRI* and M•*RsrI* both recognize GAATTC and methylate the second adenine base to give GAmATTC, where mA represents the methylated adenine product.

A related type of isomeric restriction endonucleases is termed *neoschizomer* or *heteroschizomer* (2). These endonucleases are also derived from different organisms and have the same sequence specificity, but cleave at different positions within the sequence. For instance, R•*Asp*718I and R•*Kpn*I are neoschizomers that produce G ^ GTACC and GG TAC ^ C cleaved sequences, respectively. Another neoschizomeric pair are R•*Fo*KI and R•*Sts*I, which recognize the sequence dGGATG and cleave (9/13) and (10/14) nucleotides beyond the recognition sequence, respectively (3). MTases that have the same sequence specificity but methylate at different positions are termed *heteroprostomers* (2). For example, M•*Hpa*II and M•*Msp*I both recognize CCGG, but M•*Hpa*II methylates the first C to give mCCGG and M•*Msp*I, the second C to give CmCCG.

Depending on the particular endonuclease or MTase, one **isoenzyme** may provide a more economical or suitable alternative to another. Some isomers are less costly to produce, have less [star activity](#), can be more easily freed of nuclease contamination, or have greater stability or different sensitivities to methylation. Isomeric endonucleases and MTases are useful in several applications. For example, some isoschizomers can be used to analyze the methylation state of DNA. For instance, to determine the level of 5-methylcytosine in CG dinucleotides in mammalian [genomes](#), isoschizomers with different sensitivities to methylation may be used. R•*Msp*I, R•*Hpa*II, and R•*Hap*II all recognize and cut after the first C of C ^ CGG. R•*Msp*I cuts at methylated sequences CmCCG, whereas the isoschizomers R•*Hpa*II and R•*Hap* II do not. In a related example, MTase isomers can be used to modulate endonuclease activity. M•*Hpa*II (CmCCG) and M•*Msp*I (mCCGG) are a heteroprostomeric pair. M•*Msp*I methylation inhibits R•*Msp*I and R•*Hpa*II but not R•*Hap*II, whereas M•*Hpa*II methylation inhibits R•*Hpa*II and R•*Hap*II but not M•*Msp*I.

Neoschizomers are also useful in molecular cloning. For instance, R•*Sac*I generates a 5' overhang, [staggered cut](#) on cleavage of the canonical sequence GAGCT ^ C. Its neoschizomer, R•*Eco*ICRI, generates a blunt cut: GAG ^ CTC. *Blunt-end* DNA generated by R•*Eco*ICRI can be ligated (see [DNA Ligase](#)) to other blunt-end DNA fragments, thereby eliminating the need for compatibility of cohesive termini for ligation of DNA fragments produced by R•*Sac*I cuts. The existence of restriction–modification isoenzymes expands our capacity to manipulate DNA.

## Bibliography

1. A. Pingould, J. Alves, and R. Geiger (1993) in *Methods in Molecular Biology*, Vol. 16, *Enzymes of Molecular Biology*, M. M. Burrell, ed., Humana Press, Inc., Totowa, N.J., p. 167.
2. C. Kessler and V. Manta (1990) *Gene* 92, 1–248.
3. R. J. Roberts and S. E. Halford (1993) in *Nucleases*, 2nd ed., (S. M. Lynn, R. S. Lloyd, and R. J. Roberts, eds.) Cold Spring Harbor Laboratory Press, New York, pp. 35–88.

## Isotachopheresis

Isotachopheresis (also known as “stacking”) is an application of discontinuous [buffer](#) systems for [electrophoresis](#) in which sequential moving boundaries of all the components present (the stack) are set up and displaced in the electric field at a single velocity. In principle, such a system does not differ from the stacking phase of [disc electrophoresis](#). The sequential leading and trailing ion moving boundaries may be separated from each other by [proteins](#) or other ionic species with intermediate net mobilities (“spacers”). Both [amino acids](#) and carrier ampholyte mixtures (see [Isoelectric Focusing](#)) have been used as spacers between stacked protein zones. At equilibrium, contiguous zones of stacked proteins or spacers migrate at the same velocity; hence the name



*isotachophoresis* . The width of each zone is proportional to the amount of protein or spacer in the zone. Thus, when the sample load is increased, the distance between the leading and trailing edges of a moving boundary (the zone within the stack) increases, making it possible to isolate a homogeneous protein from the center of each zone of the “extended stack.” The concentration of protein is “regulated” at very high values by the mobilities and concentrations of the ions constituting the discontinuous buffer system that generates the stack, and so isotachophoresis is an ideal preparative technique. Theoretically, isotachophoresis is not only synonymous with steady-state stacking and disc electrophoresis but becomes also mechanistically the same as isoelectric focusing when the concentration of the common counterion approaches zero.

The spaced zones can be detected by differences in their temperature or conductance, because each moving boundary at the steady-state is characterized by a specific field strength that provides an equal (“iso”) mobility in spite of the differences in net mobility between components of the stack. Such detectors affixed to capillaries are the preferred instrumentation of analytical isotachophoresis, and they preceded by many years the use of absorbance and fluorescence detectors in capillary zone electrophoresis.

#### Suggestions for Further Reading

F. M. Everaerts, J. L. Beckers, and P. E. M. Verheggen (eds.) (1976) *Isotachophoresis*. Elsevier, Amsterdam.

G. Baumann and A. Chrambach (1976) Gram-preparative protein fractionation by isotachophoresis: Isolation of human growth hormone isohormones. *Proc. Natl. Acad. Sci. USA* **73**, 732–736.

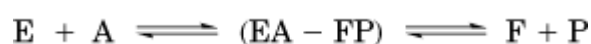
A. Chrambach and L. M. Hjelmeland (1984) "Some recent conceptual advances in moving boundary electrophoresis and their practical implications. In" *Electrophoresis '83* (H. Hirai, ed.), W. de Gruyter, Berlin, pp. 81–97.

A. Adam and C. Schots (eds.) (1980) *Biochemical and Biological Applications of Isotachophoresis* . Elsevier, Amsterdam.

## Isotope Exchange at Equilibrium

Measuring the rate of exchange of isotopes between substrates and products in the presence of an [enzyme](#) has been very useful for determining the **kinetic mechanisms** of enzyme-catalyzed reactions. It relies on the incorporation of **radioactive** atoms from a substrate (or product) into the product (or substrate) under conditions where there is no net reaction. Thus, it differs from other kinetic techniques that rely on the measurement of the rates of product formation. The procedure that is utilized depends on the mechanism under study

A single-site ping-pong mechanism consists of two partial reactions, so the overall reaction is a double displacement (see **Kinetic mechanisms**); First, substrate A is converted to P, leaving a modified form of the enzyme that then converts substrate B to product Q. For the first partial reaction



the addition of only substrate A to the enzyme will result in the release of P and the production of a modified form of the enzyme F. The maximum amount of P that can be released is equivalent to the

enzyme concentration and that will occur only when the equilibrium of the partial reaction lies far to the right or the concentration of A is high. Usually, the concentration of enzyme that is used to detect a partial exchange reaction is low and in the micromolar region, so the release of P is small. In the absence of substrate B, the system is at thermodynamic equilibrium, and the A and P reactants are shuttling back and forth. Thus, on the addition of isotopically labeled P (or A), it will be converted to A (or P) and the isotope is incorporated into A (or P). Given sufficient time, the label will be distributed in the same proportions in A and P, and the specific radioactivities of A and P will be identical. The rate at which A is converted to P must be the same as that for the conversion of P to A. It is this partial exchange reaction in the absence of the other substrate that characterizes a ping-pong mechanism. But it is important to demonstrate that the exchange reaction occurs with both halves of the reaction, ie, with both A/P and with B/Q. If only one exchange reaction occurs, it could be due to the presence of an impurity in the enzyme preparation.

The equation for the initial rate of the exchange reaction observed between A and P, starting with labeled A,  $v_{A-P}$ , is given by the expression

$$v_{A-P} = \frac{(V_1 K_{ia}/K_a)AP}{K_{ia}P + K_{ip}A + AP} \quad (1)$$

Here,  $V_1$  is the maximum rate of the enzyme-catalyzed reaction,  $A$  and  $P$  are the concentrations of the two reactants, and  $K_a$  is the  $K_m$  of A. The equation has the same form as that for the initial velocity of a ping-pong mechanism (see **Kinetic mechanisms**), except that the denominator contains **dissociation constants** for the EA ( $K_{ia}$ ) and FP ( $K_{ip}$ ) complexes (1). In the absence of the formation of a dead-end FA complex (see [Substrate Inhibition](#)), a plot of  $1/v$  against  $1/A$  at different concentrations of  $P$  yields a family of parallel straight lines (see [Lineweaver-Burk Plot](#)). Analysis of the kinetic data can give values for  $K_{ia}$  and  $K_{ip}$  as well as for the apparent maximum velocity of the isotope exchange reaction, which is equal to either  $V_1 K_{ia}/K_a$  or  $V_2 K_{ip}/K_p$ , where  $V_2$  is the maximum rate of the reverse reaction. An equation of the same form applies for measurement of the initial rate of the B-Q exchange:

$$v_{B-Q} = \frac{(V_1 K_{ib}/K_b)BQ}{K_{iq}B + K_{ib}Q + BQ} \quad (2)$$

The maximum velocity of the exchange reaction, relative to the maximum velocities of the overall chemical reactions, depends on the values of the dissociation and Michaelis constants. But there is a relationship between the maximum rates of the chemical ( $V_1$  and  $V_2$ ) and exchange ( $V_{A-P}$  and  $V_{BQ}$ ) reactions (Eq. 3):

$$\frac{1}{V_{A-P}} + \frac{1}{V_{B-Q}} = \frac{1}{V_1} + \frac{1}{V_2} \quad (3)$$

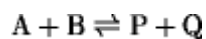
A further relationship that holds for a single-site ping-pong mechanism is given in Equation 4:

$$\frac{1}{V_{A-Q}} = \frac{1}{V_{A-P}} + \frac{1}{V_{B-Q}} \quad (4)$$

The addition of both substrates or both products to measure the rate of the overall exchange reaction would reduce the rate of exchange of the partial exchange reactions because of the formation of the FB and/or EQ complexes, which would reduce the concentration of the forms of enzyme that participate in the partial exchange reactions.

The procedure that is utilized for the study of exchange reactions with sequential reactions differs

from that for a ping-pong mechanism, because sequential mechanisms do not exhibit partial exchange reactions. The technique is applied by allowing the sequential reaction to come to equilibrium and then disturbing that equilibrium by the addition of a small chemical concentration of a highly **radiolabeled** reactant. The initial rate of the exchange of the label from substrate to product is measured as a function of the concentration of a substrate–product pair that is increased in its equilibrium ratio. For the Bi–Bi reaction



$K_{eq}$  is equal to  $(P)(Q)/(A)(B)$  and it is possible to raise the concentration of any one of four pairs of reactants,  $P/A$ ,  $P/B$ ,  $Q/A$ , or  $Q/B$ , without disturbing the equilibrium. The exchange rate can be measured for any one of three possible pairs, but for the most part, exchanges are measured between like substrate–product pairs with the absolute concentration of the same pairs being raised in constant ratio. For the creatine kinase reaction



the three possible exchanges would be creatine-phosphocreatine, ATP-phosphocreatine, and ATP-ADP. No atoms are exchanged between creatine and ADP.

The exchange patterns for ordered and rapid equilibrium, random mechanisms are given in Table 1. A straightforward ordered mechanism is characterized by the fact that the A–Q exchange is inhibited as the concentration of the B/P pair is increased (1). This result contrasts with that for a rapid equilibrium, random mechanism where there is no inhibition of the A–Q exchange rate with increasing concentrations of the B/P pair. The latter mechanism is further characterized by the equality of all three exchange reactions, as demonstrated for the creatine kinase reaction (2). This will not be the case unless catalysis is the sole rate-limiting step of the reaction sequence (3). An enzyme that catalyzes a Bi–Bi reaction which conforms to a rapid equilibrium, random mechanism has distinct binding sites for each of the two substrates involved in the reaction. Thus, it becomes possible for one substrate and the product of the other substrate to be present on the enzyme at the same time. Two such dead-end complexes could form. The one involving the smaller substrate and smaller product (EBQ) will always form, whereas the other dead-end complex may, or may not, do so. Increasing the concentration of the reactant pair that is involved in dead-end complex formation will lead to inhibition of all three exchanges.

**Table 1. Isotope Exchange Patterns**

| Mechanism                 | Varied Substrate–Product Pair | Exchange Measured | Type of Double-reciprocal Plot    |
|---------------------------|-------------------------------|-------------------|-----------------------------------|
| Ordered                   | A/Q                           | A – Q             | Linear                            |
|                           |                               | B – P             | Linear                            |
|                           | B/P                           | A – Q             | Substrate inhibition <sup>a</sup> |
|                           |                               | B – P             | Linear                            |
| Rapid equilibrium, random | A/Q                           | A – Q             | Linear                            |
|                           |                               | B – P             | Linear                            |
|                           | B/P                           | A – Q             | Linear                            |
|                           |                               | B – P             | Linear                            |

|     |       |                      |
|-----|-------|----------------------|
| B/Q | A – Q | Substrate inhibition |
|     | B – P | Substrate inhibition |
|     | A – P | Substrate inhibition |

---

<sup>a</sup> This term is used to describe the fall-off in the initial rate of isotope exchange as the concentration of a substrate–product pair is raised in constant ratio.

The equations that describe the initial rates of isotope exchange for sequential reaction mechanisms are considerably more complex than those for the partial exchange reactions with a ping-pong mechanism. The equations for ordered (2), and rapid equilibrium, random (3) and partly random (4) mechanisms have been presented, and these references also contain a more complete discussion of isotope exchange at equilibrium for sequential reaction mechanisms.

### Bibliography

1. W. W. Cleland (1986) *Invest. Rates Mech. Reactions* **6**, 791–870.
2. J. F. Morrison and W. W. Cleland (1966) *J. Biol. Chem.* **241**, 673–683.
3. P. F. Cook, G. L. Kenyon, and W. W. Cleland (1981) *Biochemistry* **20**, 1204–1210.
4. G. R. Ainslie and W. W. Cleland (1972) *J. Biol. Chem.* **247**, 946–951.

### Isotope Filtering, Editing

Selective detection of the nuclear magnetic resonance (NMR) signals of hydrogen atom (proton) spins that are **scalar coupled** to another type of nucleus produces an isotope-filtered, or isotope-edited, proton NMR spectrum. NMR experiments are now well established as a means of determining the [tertiary structures](#) of biological macromolecules and examining the complexes formed between these molecules and various ligands. Most such experiments involve observation of proton (<sup>1</sup>H) NMR signals, often with the goal of determining [nuclear Overhauser effects](#) (NOEs) between specific atoms of the system. As the size of the molecule increases, the number of hydrogen atoms in the structure increases. The proton NMR spectrum then becomes so crowded that resolution of a particular signal from nearby signals is difficult. Achieving the necessary resolution is further hindered because the widths of the NMR spectral lines generally also increase with increasing molecular size. One way to overcome these difficulties is to introduce <sup>13</sup>C and <sup>15</sup>N atoms into the molecule at nearly 100% abundance, rather than the much lower natural abundance. The <sup>1</sup>H NMR signals for protons directly attached to these isotopes can be detected selectively, producing a proton spectrum of the system that is simplified and, thus, easier to analyze for the desired information.

A wide variety of experimental methods are available to achieve isotope editing or filtering of one-, two- or three-dimensional NMR spectra (1-5). All experimental methods depend on the existence of a substantial spin coupling (J-coupled) interaction between <sup>13</sup>C or <sup>15</sup>N atoms and the hydrogen atoms that are attached to them. The one-bond coupling constant <sup>1</sup>J<sub>CH</sub> is typically 120–160 Hz, whereas <sup>1</sup>J<sub>NH</sub> for the peptide proton–peptide amide nitrogen interaction is about 90 Hz; thus, the necessary conditions for a successful isotope filtering experiment are well met in systems that have

been enriched in these isotopes. The filtering experiments can potentially also be extended to other isotopes. For example, Cd<sup>2+</sup> in [metalloproteins](#) couples strongly to nearby hydrogen atoms, and the proton NMR signals for these can be deduced by an isotope filtering experiment that relies on Cd-H coupling.

A particularly powerful application of isotope filtering is the study of complexes formed between a [receptor](#) and a **ligand**. If the receptor molecule is used in a form in which carbon and nitrogen positions are enriched to high levels with <sup>13</sup>C and <sup>15</sup>N, isotope editing would permit selective detection of the proton signals from the receptor when a complex forms between the receptor and a ligand. Alternatively, positions of the ligand could be enriched with these isotopes, in which case an isotope editing experiment with the receptor-ligand complex would display proton signals from the only ligand (6, 7). (See also [Triple Resonance](#), [Nuclear Overhauser Effect \(NOE\)](#).)

### Bibliography

1. R. H. Griffey and A. G. Redfield (1987) *Quart. Rev. Biophys.* **19**, 51–82.
2. L. P. McIntoch and F. W. Dahlquist (1990) *Quart. Rev. Biophys.* **23**, 1–38.
3. G. Otting and K. Wuthrich (1990) *Quart. Rev. Biophys.* **23**, 39–96.
4. S. W. Fesik and E. R. P. Zuiderweg (1990) *Quart. Rev. Biophys.* **23**, 97–131.
5. H. Oschkinat, T. Muller, and T. Dieckmann (1994) *Angew. Chem. Int. Ed. Engl.* **33**, 277–293.
6. S. W. Fesik (1988) *Nature* **332**, 865–866.
7. G. Wider, C. Weber, and K. Wuthrich (1991) *J. Am. Chem. Soc.* **113**, 4676–4678.

### Suggestions for Further Reading

8. J. Cavanagh, W. J. Fairbrother, A. G. Palmer III, and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.
9. *NMR of Proteins* (1993) (G. M. Clore and A. M. Gronenborn, eds.), CRC Press, Boca Raton.
10. K. Wuthrich (1976) *NMR in Biological Research: peptides and proteins*, American Elsevier, New York.

## Isotype

Isotype specificities were first described by Oudin for **immunoglobulins** as a characteristic of one molecule of one given animal species. Immunoglobulin isotypes designate the different classes or subclasses that differ from each other by the structures of their heavy chains, leading to the classical [IgG](#), [IgM](#), [IgD](#), [IgA](#), and [IgE](#), corresponding to the g, m, d, a, and \* heavy chains, respectively. IgG are furthermore divided into four subclasses that are also discrete isotypes, numbered from IgG1 to IgG4, whereas IgA has two subclasses, IgA1 and IgA2. The word also applies to light chains, with the **kappa (k) and lambda (l) light chain** isotypes. Each isotype is encoded by a separate constant-region gene that is split in as many **domains** (including the hinge regions) as are present in the corresponding chain. Kappa and lambda chains are present in all Ig classes and subclasses, whereas heavy chains are isotype-specific. The discrete isotypes have distinct effector functions, such as **complement** fixation, active transplacental transfer, and binding to discrete **Fc receptors**, thereby ensuring a selective fixation to different cell types.

During [immunization](#), IgM antibodies appear first, a few days after [immunogen](#) administration. They are rapidly replaced by IgG that will acquire [somatic hypermutations](#), thus becoming of higher

affinity. Replacement of IgM by IgG involves rearrangement at the DNA level, known as isotype [class switching](#) . The switch may involve any switch region present at the 5' position of each CH region gene (with the exception of IgD), leading to the expression of either isotype. Targeting of the isotype that will be involved in switching is under the control of **cytokines** produced by T cells and therefore only occurs in T-dependent **immune responses**.

See also entries **Antibody**, **Class switching**, **IgA**, **IgD**, **IgE**, **IgG**, and [IgM](#).

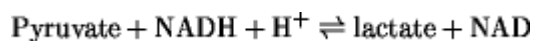
#### Suggestion for Further Reading

J. B. Natvig and H. G. Kunkel (1973) Human immunoglobulin: classes, subclasses, genetic variants and idiotypes. *Adv. Immunol.* **16**, 1–59.

## Isozyme, Isoenzyme

Isozymes, or isoenzymes, are [enzymes](#) that catalyze the same reaction, but exist in different molecular forms, possess different properties, and show different tissue distributions. They are usually recognized by the different **electrophoretic** mobilities they possess. Two classical examples of isozymes are **lactate dehydrogenase** and creatine kinase.

Lactate dehydrogenase catalyzes the reaction



It is a tetrameric enzyme with a molecular weight of 140,000. It can be composed of two different types of subunits, each of which has a molecular weight of 35,000. These subunits are the products of different genes and are designated H (heart) and M (muscle). With a tetrameric structure and two types of subunits, five different isozyme forms are possible:  $\text{H}_4$ ,  $\text{H}_3\text{M}$ ,  $\text{H}_2\text{M}_2$ ,  $\text{HM}_3$ , and  $\text{M}_4$ . All five types are found.  $\text{H}_4$  and  $\text{H}_3\text{M}$  predominate in heart muscle,  $\text{H}_2\text{M}_2$  in brain and kidney, whereas  $\text{M}_4$  is found in muscle and the liver. They differ with respect to their [substrate inhibition](#) by pyruvate, with  $\text{H}_4$  being the most sensitive (1).

Creatine kinase is a dimeric enzyme with a molecular weight of 80,000 that catalyzes the reaction



The two types of subunits are products of two different genes, have molecular weights of 40,000, and are referred to as M (muscle) and B (brain). The three isozymes of creatine kinase are MM, MB, and BB, and they are found in muscle, the heart and brain, respectively. They differ in their electrophoretic and kinetic properties (2).

The isozymes of both lactate dehydrogenase and creatine kinase play an important role in the clinical diagnosis of myocardial infarction and muscle disease. This is because particular forms of each of these enzymes occur in these organs and are present in blood only after cell necrosis.

## Bibliography

1. E. H. Braswell (1975) In *Isozymes I: Molecular Structure* (C. I. Markert, ed.), Academic Press, New York, p. 119.
2. D. C. Watts (1973) *The Enzymes* **8**, 383–455.

## J Chain

The [immunoglobulin J chain](#) is a glycosylated [polypeptide chain](#) with a molecular weight of 15 to 16 kDa, which is linked covalently to polymerized forms of [IgA](#) and [IgM](#) by [disulfide bonds](#). It is encoded by a **gene** containing four exons (and three **introns**) located on chromosome 5 in mice and 4 in humans and expressed exclusively in the [B cell](#) lineage. It is found in mucosal plasma cells of

the intestinal or respiratory tract and interstitium of mammary, salivary, or lacrimal gland and is present only in the cytoplasm. The J chain folds into an eight-stranded  $\beta$ -barrel that is presumably organized as one single Ig domain (see [Immunoglobulin Structure](#)). The positions of the [cysteine](#) residues involved in binding the J chain to either the  $\alpha$  or  $\mu$  chain have been determined by [site-directed mutagenesis](#). The J chain contains eight cysteine residues, of which two interact with the penultimate cysteines of  $\alpha$  or  $\mu$  chains, to form disulfide bridges. There is always one J chain per polymer, whatever the number of monomers. IgA dimers are generally seen as bound end-to-end through interactions with J chain. It has long been proposed that J chain was necessary for binding of the IgA dimer to the polymeric receptor involved in transepithelial transport of the immunoglobulin. Recent experiments indicated that dimerization was retained in J-chain-deficient mice, although transepithelial transport was maintained, probably by a mechanism different from that observed under physiological conditions.

See also entries [Antibody](#), [Class Switching](#), [IgA](#), [IgM](#), [Immune Response](#), and [Immunoglobulin Biosynthesis](#).

#### Suggestions for Further Reading

- M. E. Koshland (1985) The coming of age of the immunoglobulin J chain. *Annu. Rev. Immunol.* **3**, 425–453.
- E. J. Wiersma et al. (1997) Analysis of IgM structures involved in J chain incorporation. *J. Immunol.* **159**, 1719–1726.
- B. A. Hendrickson et al. (1995) Lack of association of secretory component with IgA in J chain-deficient mice. *J. Immunol.* **157**, 750–754.
- J. Mestecky et al. (1997) Immunoglobulin J chain. An early differentiation marker of human B cells. *Ann. NY Acad. Sci.* **815**, 111–113.

## J (Joining) Gene Segments

[Gene rearrangements](#) of [immunoglobulin](#) (Ig) genes take place during [B-cell](#) differentiation and operate between discrete gene segments, V-D-J for the heavy chains and V-J for the light chains. The **codons** that encompass this zone of rearrangement encode the third complementarity determining regions (CDR3) of H and L chains (see [Immunoglobulin Structure](#)). J gene segments are organized differently in the various Ig loci. In humans, there are six JH and four Jk that are clustered at the 5' position of the constant-region genes. The JI genes are organized differently, because they are individually associated with a discrete CI gene. There are seven JI-CI tandemly associated genes, although only four are functional. With the exception of these minor organizational differences, all J genes, whatever the locus, have the same general structure, with a translatable region having an average length of 13 codons, and a recombination signal sequence (RSS) involved in gene rearrangement that follows the usual organization with a heptamer, a spacer, and a nonamer sequence, and located at the 5' position of the coding region. The spacer has a length of 23 nucleotides for the JH and the JI genes, whereas it is of 12 for the Jk gene segments. The coding region may not be used in its entirety, depending on the rearrangements. Interestingly, when a 5' Jk gene has been selected to form a Vk-Jk functional gene, a second rearrangement involving a more 3' J gene and a more 5' V gene may take place, provided that the RAG (recombinase activating gene) machinery is reactivated (see [Recombinase](#)). This leads to a secondary rearrangement, also known as *editing*, that may allow a B cell to escape negative selection in the bone marrow, or even in the periphery, as part of the mechanisms used to set up **tolerance**.



See also entries [Antibody](#), [Gene Rearrangement](#), [Immunoglobulin](#), and **Tolerance**.

### Suggestions for Further Reading

H. G. Zachau (1993) The immunoglobulin kappa locus-or-what has been learned from looking closely at one-tenth of a percent of the human genome. *Gene* **135**, 167–173.

J.-P. Fripiat et al. (1995) Organization of the human immunoglobulin lambda light-chain on chromosome 22q11.2. *Hum. Mol Genet.* **4**, 983–991.

G. P. Cook et al. (1994) A map of the human immunoglobulin V(H) locus completed by analysis of the telomeric region of chromosome 14q. *Nature Genet.* **7**, 162–168.

## JAK/STAT Signaling

JAKs are ubiquitously expressed, cytoplasmic proteins that are associated with [tyrosine kinase receptors](#); they are required for the biological effects of most **cytokines**. They can couple cytokine receptors to a variety of downstream [signal transduction](#) pathways, but are specific for activation of a unique family of [transcription factors](#), known as STATs.

The JAK family consists of four members, each with a molecular weight in the range of 120 to 130 kDa. These proteins have a [kinase](#) domain in the C-terminus, which is preceded by a pseudokinase domain, the function of which is uncertain. The association of JAKs with the receptor depend upon the receptor structure (1). Single-chain receptors (such as those for GM-CSF [colony stimulating factor](#) or [growth hormone](#)) interact with JAK through a JAK-binding domain in the juxtamembrane region. Receptor aggregation is then thought to bring the kinases into proximity, causing their phosphorylation of each other and activation. Receptors for ligands such as the [interleukins](#) IL-3 and IL-6 contain two subunits. In this case, the b chain, which may be shared among different ligand-binding chains, interacts with JAK in a similar juxtamembrane region, undergoing activation upon aggregation. Some receptors (such as those for [interferons](#)) contain three independent subunits, in which b and g chains are required for signaling through JAK.

### 1. Downstream Substrates

Like [tyrosine kinase receptors](#), the JAK-associated cytokine receptors activate a variety of signaling pathways in cells, depending on the combinatorial diversity of the **SH2 domain**-containing proteins that are involved (1). A number of the signaling subunits for the cytokine receptor have consensus binding sites for proteins such as **g**, SHP2, **phosphoinositide 3'-kinase**, Shc, and others. Additionally, the JAK proteins can phosphorylate a variety of intracellular substrates, including IRS proteins, Shc, and others, that lead to activation of the [MAP kinase](#) and phosphoinositide 3'-kinase-initiated cascades. However, the one family of intracellular substrates that is more or less specific for this family of receptors are the STAT proteins.

### 2. STATs

STATs are SH2-containing substrates of the JAK kinases that act as [transcription factors](#), directly linking the receptor to the regulation of gene expression. All of these proteins contain a C-terminal SH2 domain, an SH3 domain, and several additional domains at the amino terminus. The SH2 domains of the different STAT family members recognize discrete docking sites on the receptor complex. STATs are recruited to the receptor upon cytokine binding, when they undergo

**phosphorylation** of tyrosine residues by JAK. Phosphorylation then triggers formation of another complex with other cellular proteins, or in some cases dimerization. This dimerization is mediated by intermolecular SH2 interactions. These phosphorylation-induced [protein–protein interactions](#) are accompanied by nuclear localization of the protein, where it can bind to specialized **upstream activating sequences** in the **promoters** of responsive genes.

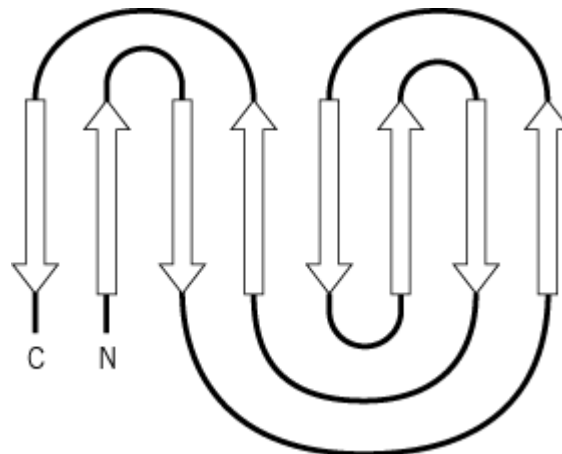
#### Bibliography

1. J. N. Ihle (1995) *Nature* **377**, 591–597.

### Jelly Roll Motif

The jelly roll motif describes a particular **topology** for arranging eight [b-strands](#) into an antiparallel [b-sheet](#) that is frequently found in [protein structures](#) (Fig. 1). The name comes from the similarity between this b-strand topology and a slice of rolled cake, called a jelly roll or Swiss roll. The jelly roll topology forms a common **domain** structure in proteins called the jelly roll b-barrel, a type of **antiparallel b-barrel**.

**Figure 1.** Schematic representation of the topology of the jelly roll motif in proteins, with individual b-strands of the b-sheet depicted as arrows. The *N*- and *C*-termini of the motif are labeled.



[See also [Beta-Sheet](#).]

#### Suggestion for Further Reading

1. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

### Jumping Gene

“Jumping gene” is another term for “[transposable element](#),” a discrete piece of DNA that can move between nonhomologous positions in a [genome](#). Such elements were first discovered by Barbara McClintock in the 1940s during her study of irregular pigment patterns in maize kernels. We now know that such transposable elements are very widespread, being found in virtually all organisms examined.

The lively term “jumping genes” conveys a sense of the dynamic nature of some segments of DNA. This dynamic view contrasts with the view of DNA as a static and unchangeable molecule, being passed from generation to generation without alteration. While it is certainly true that in most organisms DNA is not subject to large-scale change, the identification of “jumping genes” revealed that some segments of DNA could move.

The “jumping” of an element from one site to another is usually mediated by information encoded by the [transposon](#). The transposon encodes a special [recombinase](#), a [transposase](#), that mediates element translocation. This transposase binds to and acts on special sequences at the ends of the element to execute the DNA breakage reactions that separate the element from the donor site and then join it to the insertion site. When a transposable element translocates from one site in the genome to another, it can have considerable effect. Insertion into a host **gene** by such an element will generally inactivate that gene, resulting in a [mutation](#). Indeed, a considerable fraction of spontaneous mutations in some organisms result from transposable element insertion.

Some transposable elements encode only a transposase and have no “activity” other than their ability to cause mutation. Other elements, however, encode additional genes that can have considerable effect of the host organism. For example, many bacterial elements carry an [antibiotic-resistance](#) gene so that bacteria containing this element are resistant to this drug. Thus when such a bacterial element transposes from a plasmid on the chromosome, the antibiotic marker gene becomes covalently linked to the chromosome and the bacterium will continue to be **drug-resistant** in the absence of the plasmid. Conversely, when a drug-resistance transposon translocates from the chromosome onto plasmids that can move from one cell to another, this drug-resistance gene will be transmitted to many other bacteria. We now appreciate that transposable elements are the basis of the rapid spread of antibiotic resistance determinants, an increasingly widespread clinical problem.

Another important type of transposable elements are **viruses** in which transposition is used to link the viral **chromosome** to the genome of the infected cell ([1](#)). One important class of transposable viruses that affects humans is the human immunodeficiency virus ([HIV](#)) virus that can lead to acquired immune deficiency syndrome (AIDS). When the HIV virus enters a human cell, its DNA is integrated into the host human genome by transposition; thus this viral DNA is now a stable part of the genome. Because the virus has integrated, it is always present in the genome and cannot be lost. Thus strategies for treating HIV infection and AIDS rely upon affecting HIV gene expression and action, rather than trying to eliminate the virus.

#### Bibliography

1. J. D. Boeke and J. P. Stoye (1997) In *Retroviruses* (H. Varmus, S. Hughes, and J. Coffin, ed.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 343–435.

## Junctional Diversity

Junctional diversity in [immunoglobulins](#) refers originally to the multiple **codons** that can be generated at position 96 of the **kappa (k) light chain** as the consequence of various possible breakpoints when Vk and Jk genes undergo [gene rearrangement](#) in [B cells](#). This phenomenon may be schematically illustrated as follows:

| Codon Number           | 95 | 96                     | 97  |
|------------------------|----|------------------------|-----|
|                        |    | Vk gene C C T C C C 3' |     |
|                        |    | 5' T G G A C G Jk Gene |     |
| Junction 1             |    | C C T T G G A C G      |     |
| Protein sequence 1 PRO |    | TRP                    | THR |
| Junction 2             |    | C C T C G G A C G      |     |
| Protein sequence 2 PRO |    | ARG                    | THR |
| Junction 3             |    | C C T C C G A C G      |     |
| Protein sequence 3 PRO |    | PRO                    | THR |
| Junction 4             |    | C C T C C C A C G      |     |
| Protein sequence 4 PRO |    | PRO                    | THR |

The term *junctional diversity* has been extended to include events that take place during gene rearrangements when DNA strands are cleaved and before the D–J, V → D–J, or V → J coding joints are formed. This includes removal of a variable number of nucleotides by exonucleases and subsequent addition of a random sequence of some nucleotides by terminal deoxynucleotidyl transferase (TDT). The frequent result of these various events is a highly increased diversity at the joining regions of rearranged genes.

See also entries [Gene Rearrangement](#).

#### Suggestions for Further Reading

S. Tonegawa (1983) Somatic generation of antibody diversity. *Nature* **302**, 575–581.

M. Gellert (1997) Recent advances in understanding V(D)J recombination. *Adv. Immunol.* **64**, 39–64.

### $K_m$ (Michaelis Constant)

The Michaelis constant is named after Leonor Michaelis who, together with Maud Menten, was responsible for laying the foundations of steady-state [enzyme](#) kinetic theory (1). On the basis of their studies on invertase, they formulated the kinetic scheme:



where E, S, and P represent the enzyme, substrate, and product, respectively (see [Michaelis–Menten](#)

[Kinetics](#)). The initial velocity equation for the reaction under steady-state conditions was presented as:

$$v = \frac{k_3 E_t S}{K_s + S} = \frac{V S}{K_m + S} \quad (2)$$

where  $K_s$  was considered to be the **dissociation constant** of the ES complex and equal to  $k_2/k_1$ .

Subsequently, it was demonstrated that the rate of product formation is not always slow compared with the rate at which ES dissociates back into E and S. Therefore, in the more general formulation of the initial-rate equation for an enzyme-catalyzed reaction,  $K_s$  was replaced by  $K_m$ , with the  $m$  chosen in honor of Michaelis.  $K_m$  is not assumed to be a dissociation constant and becomes equal to  $(k_2 + k_3)/k_1$ .  $K_m$  can be considered a generic term.

For a single-substrate reaction,  $K_m$  always gives a measure of the concentration of substrate required to yield half-maximum velocity, which usually also gives an indication of the intracellular concentration of the substrate. An advantage associated with the use of Michaelis constants to characterize enzymes and their substrates is that the value is independent of enzyme concentration. An impure enzyme preparation can be used for such purposes, provided that it does not contain interfering enzyme activities, inhibitors, or activators.

The expressions for a Michaelis constant increase in complexity with an increase in the number of substrates and products associated with a reaction (2). In general terms, with multisubstrate enzymes, the Michaelis constant is considered to be the concentration of substrate that yields half-maximum velocity when all other substrates are saturating. For this type of reaction, it is now common practice to use A, B, etc. as the symbols for substrates and P, Q, etc. as symbols for products, with the Michaelis constants for these reactants being denoted by  $K_a$ ,  $K_b$ ,  $K_p$ ,  $K_q$ , etc. Thus, the expressions for the initial-rate equations are made as simple as possible, with retention of the significance of the Michaelis constant. Inhibition constants for a substrate or product are distinguished from Michaelis constants by having the letter  $i$  precede the substrate symbol, eg,  $K_{ia}$ ,  $K_{ip}$ .

## Bibliography

1. L. Michaelis and M. L. Menten (1913) *Biochem. Z.* **49**, 333–369. [Translated into English by T. R. Caine Boyde (1980) In *Foundation Stones of Biochemistry*, Voile et Aviron, Hong Kong, pp. 287–316.]
2. W. W. Cleland (1963) *Biochim. Biophys. Acta* **67**, 104–137.

## Kallikreins

Kallikreins are a subgroup of [serine proteinases](#) that, by definition, function physiologically to release small vasoactive peptides, called *kinins*, from circulating plasma proteins, called *kininogens* (1, 2). Proteinases that release kinins from kininogens are also known as *kininogenases*. Thus, kallikreins are also considered a subgroup of kininogenases. Kallikreins were initially discovered as having an important role in the regulation of blood pressure through kinin production, but they have subsequently been implicated in a number of other physiological processes, including inflammation, [blood clotting](#), fibrinolysis, smooth muscle contraction, and the activation of prohormones and proenzymes. The name *kallikrein* was derived from the Greek word *kallikreas* for pancreas, because

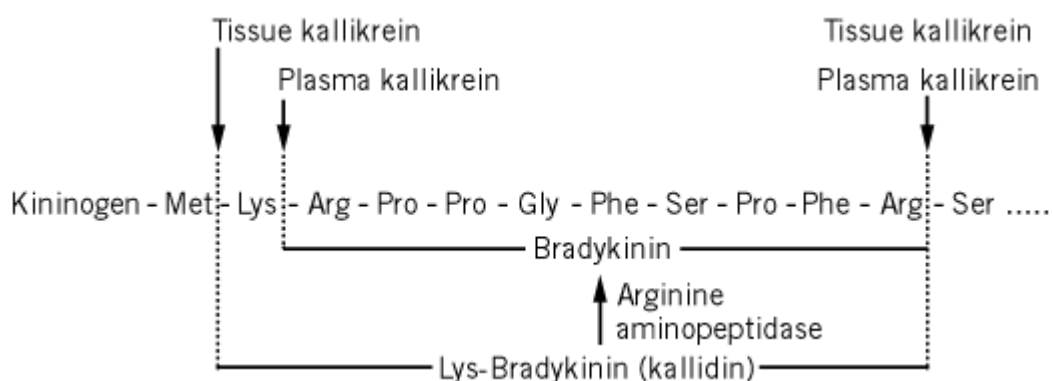
it was originally thought that the tissue kallikrein found in urine originated from the pancreas. Two types of kallikreins have been identified in mammals: plasma kallikrein and tissue kallikrein. These two types of kallikreins differ in their substrate specificity, molecular weight, type of kinin released, and physiological roles.

## 1. Basic Properties and Genomic Structure

### 1.1. Plasma Kallikrein

Plasma kallikrein (PK) is expressed solely in the liver and is coded by a single 22-kbp **gene** (2, 3). This gene contains 15 exons (see [Introns, Exons](#)), where exons 1 and 2 code for the **signal sequence**. The PK gene is very similar in organization and **chromosomal** location to that for factor XI (antecedent thromboplastin plasma), a factor involved in the blood-clotting cascade. These two genes are closely located on the same chromosome (chromosome 4q34-35 in humans and chromosome 8 in mouse) and have 15 exons, and their [complementary DNA](#) sequences are 58% identical. This suggests that PK and factor XI originated by gene duplication from a common ancestor. PK is secreted by hepatocytes as its [zymogen](#) form, pre-PK (Fletcher factor). The PK gene has been cloned, and mature PK purified to homogeneity from rat, mouse, and human. The human PK has 75% amino acid sequence identity with rat and mouse PK, indicating a substantial conservation between these species. PK hydrolyzes its natural substrate, high-molecular-weight kininogen (HK), to release the vasoactive nonapeptide, bradykinin (Fig 1).

**Figure 1.** Kinin formation from kininogen by plasma kallikrein and tissue kallikrein.



### 1.2. Tissue Kallikreins

Tissue kallikreins are encoded by a multigene family of similar genes, termed *KLK* (1, 4). Many genes have been identified in various species (Table 1). The genes for all members of the *KLK* family in various species contain five exons and four introns, and the proteins are serine proteinases. Only one member of each family functions exclusively as a true kallikrein, although other members may have limited kininogenase activity. The nomenclature among the different species has been standardized so that the true kallikrein gene in each family is designated as *KLK1* and its product as K1—for example, the *hKLK1* gene and hK1 protein in humans (4). Physiologically, the K1 proteins hydrolyze a different form of kininogen, called low-molecular-weight kininogen (LK), to release the decapeptide, Lys-bradykinin (Lys-BK) or kallidin (Fig. 1). *KLK1* is expressed predominantly in the pancreas, salivary glands, and kidney, but it is also found at low levels in other tissues. Other members of the tissue kallikrein family, although not true kallikreins, share significant amino acid and structural [homology](#) with the *KLK1*-encoded protein. However, the substrates and the roles of many of these proteins have not been elucidated. The mouse *KLK* gene family (*mKlk*) was the first tissue kallikrein family to be discovered. It contains 24 genes; 14 of these are functional and 10 are inactive [pseudogenes](#) (4). All of the 24 genes are clustered in a 310-kbp region of chromosome 7. The rat kallikrein family (*rKLK*) contains 13 genes; 10 of these are expressed and 3 are pseudogenes

(1, 5). All 13 genes are located within a 440-kbp BssHIII fragment of the rat chromosome. The chromosomal location of the rat *KLK* gene family has not been published. The *KLK* families in monkey, dog, pig, and guinea pig are much smaller (1) (Table 1).

**Table 1. Tissue Kallikreins**

| Species   | Gene Name     | Size (aa)   | Protein   |
|---|---------------|-------------|---|
| <b>Human</b>                                    | <i>hKLK1</i>  | 262         | hK1 (true kallikrein)   |
|   | <i>hKLK2</i>  | 261         | hK2   |
|   | <i>hKLK3</i>  | 261         | hK3 (prostate-specific antigen)   |
| <b>Monkey</b>                                   | <i>cmKLK1</i> | 257         | cmK1 (true kallikrein)  |
|   | <i>rhKLK1</i> | 261         | rhK1 (prostate-specific antigen)  |
| <b>Dog</b>                                      | <i>dKLK1</i>  | 261         | dK1 (true kallikrein)   |
|   | <i>dKLK2</i>  | 260         | dK2 (canine prostatic arginine esterase)  |
| <b>Pig</b>                                      | <i>pKLK1</i>  |             | pK1 (true kallikrein)   |
| <b>G. pig</b>                                   | <i>gpKLK1</i> |             | gpK1 (true Kallikrein)  |
|   | <i>gpKLK2</i> |             | gpK2 (prostate tissue kallikrein)   |
| <b>Mouse</b>                                    | <i>mklk1</i>  | 261         | mK1 (true kallikrein)   |
|   | <i>mklk3</i>  | 261         | mK3 (g-nerve growth factor)   |
|   | <i>mklk4</i>  | 261         | mK4 (a-nerve growth factor)   |
|   | <i>mklk5</i>  | 261         | mK5   |
|   | <i>mklk6</i>  | 261         | mK6   |
|   | <i>mklk8</i>  | 261         | mK8   |
|   | <i>mklk9</i>  | 261         | mK9 (epidermal growth factor-binding proteinC)                                      |
|   | <i>mklk11</i> | 261         | mK11  |
|   | <i>mklk13</i> | 261         | mK13 (epidermal growth factor-binding proteinB/prorenin converting enzyme)          |
|   | <i>mklk14</i> | 261         | mK14,   |
|   | <i>mklk16</i> | 261         | mK16 (g-renin)  |
|   | <i>mklk21</i> | 261         | mK21  |
|   | <i>mklk22</i> | 261         | mK22 (epidermal growth factor-binding proteinB/b-nerve growth factor endopeptidase) |
|   | <i>mklk24</i> | 261         | mK24  |
|   | <i>mklk26</i> | 261         | mK26 (prorenin converting enzyme 2)   |
| <i>mklk2, 7, 10, 12, 15, 17,18, 19, 23, 25,</i> |               | Pseudogenes |   |
| <b>Rat</b>                                      | <i>rKLK1</i>  | 261         | rK1 (True kallikrein)   |
|   | <i>rKLK2</i>  | 261         | rK2 (Tonin)   |
|   | <i>rKLK3</i>  | 261         | rK3   |
|   | <i>rKLK4</i>  | 261         | rK4   |

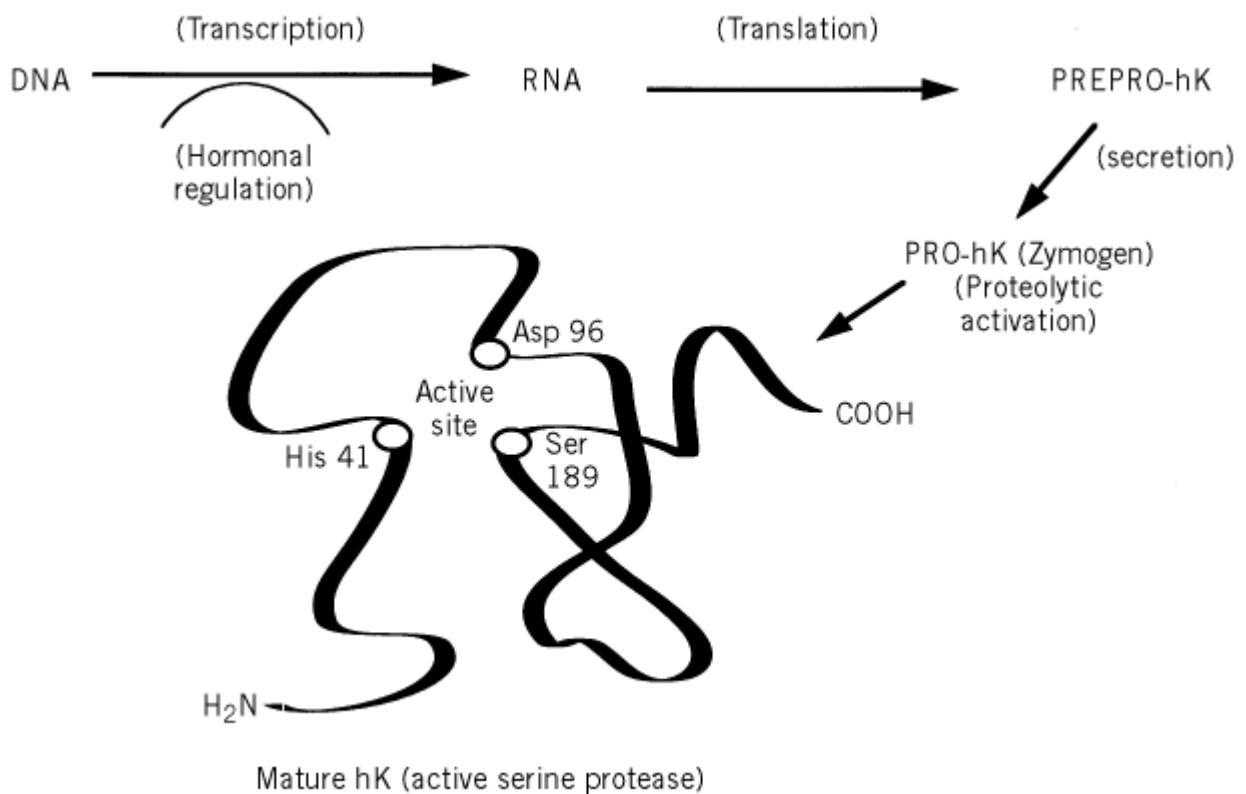
|                      |     |   |
|----------------------|-----|---|
| <i>rKLK6</i>         | 261 | rK6   |
| <i>rKLK7</i>         | 261 | rK7 (Esterase B/proteinase A/K7)                      |
| <i>rKLK8</i>         | 261 | rK8   |
| <i>rKLK9</i>         | 261 | rK9 (Submandibular enzymatic vasoconstrictor[SEV]/S3) |
| <i>rKLK10</i>        | 261 | rK10 (T-kininogenase/proteinase)                      |
| <i>rKLK12</i>        | 261 | rK12  |
| <i>rKLK5, 11, 13</i> |     | Pseudogenes   |

---

The tissue kallikrein family in humans is composed of three genes, located in a 60-kbp region on chromosome 19 q13.3-13.4 (1, 6). The *KLK1*, *KLK2* and *KLK3* genes code for hK1, hK2 (human kallikrein 2), and hK3 (prostate-specific antigen, PSA), respectively. The three human tissue kallikreins (HKS) are highly homologous: hK3 has 60% and 80% sequence identity with hK1 and hK2, respectively (7). The expression and protein processing of all three human tissue kallikreins (hK) are very similar and are represented schematically in Fig. 2. Similar to PK, all three hKs are synthesized as their enzymatically inactive prepro form. The prepro region of all three hKs encodes the signal sequence and is composed of a 17-residue pre-region and a 7-residue pro-region (7). The pro form of hK (PHK) is secreted from the cells and is converted to mature, enzymatically active hK by **trypsin**-like extracellular proteinases (2, 8-10). The pro forms of hK1, hK2, and hK3 have all been found in human body fluids (2, 11, 12). Although not true kallikreins, hK3 (PSA) and hK2 have recently received substantial attention. PSA is perhaps the best-known human kallikrein because it is a well-established serum marker for prostate cancer and benign prostatic hyperplasia (7). hK2 has recently been a focus of intense investigation and has been proposed as a potential diagnostic marker for prostate cancer (7). Because kallikreins in various species are very similar, the rest of this section will emphasize the expression, enzymatic characteristics, and physiological roles of human kallikreins, with reference to kallikreins from other species when necessary.

**Figure 2.** Expression of human tissue kallikrein (hK).





## 2. Gene Regulation and Protein Expression

### 2.1. Plasma Kallikrein

PK is expressed only in the liver. The **promoter** of the rat gene contains the consensus TATTAA box (see [TATA Box](#)) and CCAAT box (see [CAAT Box](#)) elements. The liver-specific regulatory sequences of PK are not known, but a putative estrogen [response element](#), which may be responsible for the elevated levels of PK [messenger RNA](#) in female rats, has been localized between nucleotides -343 and -330 (3).

### 2.2. Tissue Kallikreins

The expression of all three *hKLLK* genes are regulated *in vivo* by steroids. Presumably both [cis-acting](#) and [trans-acting](#) elements responsible for tissue specificity and steroid responsiveness are embedded within the promoters. Both [glucocorticoids](#) and [estrogen](#) have been shown to regulate *KLK1* gene expression. Glucocorticosteroids, which regulate a wide range of inflammatory responses, decrease biosynthesis of rat K1 (13). Furthermore, cortisol treatment of adrenalectomized rats increases rat salivary gland tissue kallikrein activity and alters *rKLK1* mRNA levels. Estrogen increases both the *rKLK1* mRNA level and rK1 enzyme activity and immunoreactivity in normal rat anterior pituitary. Similar observations have been made in *hKLLK1* expression and enzyme activity during the menstrual cycle (14). Recent studies in rats and mice have shown that *KLK1* expression is not regulated by androgens or [thyroid hormone](#) (14).

The hormonal and other regulatory mechanisms that control *KLK1* expression at the transcriptional level have not been studied in detail. The 5' flanking regions of human *KLK 1* contain several **consensus sequences** for estrogen, progesterin, glucocorticosteroids, or [cyclic AMP](#) (cAMP) response elements. However, these putative elements have not been functionally tested. [Sequence analysis](#) of mouse, rat, and human *KLK1* genes revealed a conserved putative cAMP responsive element between nucleotides -235 and -214. It contains the sequence CCCCACCC (15), a motif that has been identified as a binding site for AP-2, a [transcription factor](#) that has been implicated in the

regulation of many genes. Some of the above-mentioned hormonal responses may not be the direct result of primary hormone activation, but secondary to other cellular changes.

As opposed to *KLK1*, androgen and thyroid hormones are considered the primary regulators of gene expression and enzyme activity for the remainder of the *KLK* genes. Androgen upregulates the expression of these genes. Salivary gland mRNA levels of these genes are higher in males than in females. Furthermore, castration of mice and rats reduces *KLK* gene expression, but androgen administration to these animals restores *KLK* gene expression. Thyroid hormone administration also increases *KLK* gene expression in both rats and mice. Transcriptional activation of *hKLK2* and *hKLK3* by androgen has been extensively studied. The promoters of both *hKLK2* and *hKLK3* contain androgen response elements (AREs), which are necessary for transcriptional up-regulation by androgens (16, 17). The AREs in *hKLK3* have been mapped to the regions starting from nucleotides –170 and –394. The molecular mechanism by which thyroid hormone regulates the expression and activity of *KLK* genes is not understood at this time.

To identify the elements that may be contributing towards tissue-specific expression, at least 500 bp of the 5'-flanking region of all the mouse, rat and human *KLK* genes have been sequenced and compared. These regions are highly conserved (70–90% identity), particularly within species. There are at least two conserved elements at nucleotides –235 to –214 and –177 to –157 (18). The reverse complement of the latter sequence is similar to a *cis*-regulatory element in the gene for a rat pancreatic serine proteinase, suggesting that it may be responsible for directing pancreatic-specific expression of *hK1*. Transgenic mice studies have also suggested a pancreatic tissue-specific sequence within the 500-bp region upstream of the *KLK1* (19). The regions responsible for tissue-specific expression of *hKLK2* and *hKLK3* have not been localized.

### 3. Protein Characteristics and Physiological Roles

#### 3.1. Plasma Kallikrein

Mature PK (EC 3.4.21.34) is a 619-residue glycoprotein ( [isoelectric point](#), pI, of 8.9) that circulates in serum as pre-PK bound to its substrate, HK. Factor XIIa activates both pre-PK and factor XI (which is also bound to HK) (20). PK plays an essential role in the blood clotting cascade. The activation of pre-PK involves a single peptide bond cleavage at Arg371 to form a two-chain molecule attached by disulfide bonds. This creates a protein with a heavy chain (371 residues) and a light chain (248 residues) with independent functions, both of which are necessary for promotion of clot formation. The heavy chain is responsible for binding to factor XIIa, while the proteolytic activity resides in the light chain (21). The light chain of PK acts as a serine proteinase with some structural characteristics not unlike the tissue kallikreins: a [catalytic triad](#) composed of His, Asp, and Ser residues, and an Asp residue in the binding pocket that directs cleavage after basic residues.

#### 3.2. Tissue Kallikrein

Tissue kallikreins (EC 3.4.21.35) share sequence and structural homology with other serine proteinases, such as trypsin, [chymotrypsin](#), [plasmin](#), [thrombin](#), and [elastase](#). The high homology with these well-studied proteinases has allowed structural comparisons (see [Homological Modeling](#)). The human kallikreins share highly conserved amino acid residues that comprise the catalytic triad: His 41, Asp 96, and Ser 189—the hydrolytic amino acid (22) (see [Catalytic Triad](#)). Most kallikreins are described as having trypsin-like activity, which is characterized by peptide bond cleavage after the basic amino acids, [lysine](#) and [arginine](#). Molecular modeling studies based on the crystal structures of known serine proteinases suggest that certain kallikrein residues are important for substrate specificity. In the human kallikrein family, this is exemplified by *hK3* (PSA). *hK3* is distinct from *hK1* and *hK2* in that it possesses a Ser residue at position 183, a site that is highly conserved as an Asp among all other kallikreins. Molecular modeling studies indicate that Asp183 is responsible for the coordination of basic amino acids in the substrate binding pocket. The substitution of Ser for Asp183 in *hK3* is thought to be the reason that *hK3* has chymotrypsin and not trypsin-like specificity (22). *hK2* also has unique properties within the human tissue kallikrein family, even though it has an Asp183 and displays trypsin-like activity like *hK1*. Also found in the

prostate, like hK3, hK2 is distinct from hK1 and hK3 in both substrate specificity and apparent function. It is the only human kallikrein that can autoactivate (12) and can activate pro-hK3 (10).

#### 4. Physiological Roles

PK is associated with a number of metabolic processes but is most commonly associated with blood clotting, fibrinolysis, inflammation, and vascular tone (2). In the coagulation cascade, both coagulating factor XI and pre-PK are bound to circulating HK, which becomes attached to tissue surfaces following vascular damage. PK can activate **urokinase**-type [plasminogen](#) activator, which in turn activates plasminogen and leads to fibrinolysis. K1 is associated with hypertension and inflammation through regulation of local blood flow, Na<sup>+</sup>/water homeostasis, and vascular permeability (1, 2).

Both PK and tissue K1 play a primary role in inflammation through the release of kinins. The inflammatory response is mediated by neutrophil invasion. Both HK and LK are bound to the external surface of neutrophils. The kinins, BK and Lys-BK, may be released from HK and LK by the action of PK (also bound on the surface of neutrophils) and K1, respectively. The kinins released during this process are potent pain-producing agents, and they mediate local vasodilation and vascular permeability.

Physiologically, PK is inhibited by C1 inhibitor (which also inhibits factor XIIa) but may also be inhibited by the serum serine proteinase inhibitors  $\alpha_2$ -**macroglobulin** and [antithrombin](#) III. The activity of K1 appears to be regulated by a specific physiological inhibitor called kallistatin. Kallistatin, expressed primarily in the liver, is a member of the [serpin](#) (serine proteinase inhibitor) family (23). Kallistatin itself has been linked directly to blood pressure regulation by experiments that have shown it to reverse the hypotensive effects of tissue kallikrein expression in transgenic mice (24).

PK and K1 are largely associated with kinin production, whereas other members of the tissue kallikrein family have different functions. An example of kallikrein diversity can be seen in the case of the human tissue kallikrein family. Neither hK2 nor hK3 shows significant kininogenase activity. hK2 has trypsin-like specificity, like hK1 (25), but hK2 releases kinins from serum at levels 1000-fold lower than hK1 (26) and shows no ability to release lys-bradykinin (27). hK3 has chymotrypsin-like activity and therefore no kininogenase activity. Both hK2 and hK3 are found in the prostate but exhibit entirely different functions from one another (7). One physiological role for hK3 is the liquefaction of the seminal clot in seminal plasma after ejaculation, through the cleavage of semenogelin I and II. Because hK2 is unique among the human kallikreins in its ability to activate pro-hK3 (Pro-PSA), it may function in the regulation of hK3 (PSA) activity in the prostate.

#### Bibliography

1. J. A. Clements (1997) In *The Kinin System*, S. G. Farmer (ed.), San Diego: Academic Press, pp. 72–97.
2. K. D. Bhoola, C. D. Figueroa, and K. Worthy (1992) *Pharmacol. Rev.* **44**, 1–80.
3. G. Beaubien, I. Rosinski-Chupin, M. G. Mattei, M. Mbikay, M. Chretien, and N. G. Seidah (1991) *Biochemistry* **30**, 1628–1635.
4. O. A. Carretero, L. A. Carbini, and G. Scicli (1993) *J. Hypertens.* **11**, 693–697.
5. A. G. Scicli, L. A. Carbini, and O. A. Carretero (1993) *J. Hypertens.* **11**, 775–780.
6. L. A. Carbini, A. G. Scicli, and X. Carretero (1993) *J. Hypertens.* **11**, 893–898.
7. H. G. Rittenhouse, J. A. Finlay, S. D. Mikolajczyk, and A. W. Partin (1998) *Crit. Rev. Clin. Lab. Sci.* **35**, 275–368.
8. S. Takahashi, A. Irie, Y. Katayama, K. Ito, and X. Miyake (1986) *J. Biochem.* **99**, 989–992.
9. A. Kumar, A. Goel, T. Hill, S. Mikolajczyk, L. Millar, K. Kuus-Reichel, and M. Saedi (1996)

Cancer Res. **56**, 5397–5402.

10. A. Kumar, S. D. Mikolajczyk, A. S. Goel, L. S. Millar, and M. S. Saedi (1997) Cancer Res. **57**, 3111–3114.
11. M. S. Saedi, T. M. Hill, K. Kuus-Reichel, A. Kumar, J. Payne, S. D. Mikolajczyk, R. Wolfert, and H. G. Rittenhouse (1998) Clin. Chem. **44**, 2115–2119.
12. S. D. Mikolajczyk, L. S. Grauer, L. S. Millar, T. M. Hill, A. Kumar, H. G. Rittenhouse, R. L. Wolfert, and M. S. Saedi (1997) Urology **50**, 710–714.
13. A. A. Jaffa, D. H. Miller, R. H. Silva, H. S. Margolius, and R. K. Mayfield (1990) Kidney Int. **38**, 212–218.
14. J. A. Clements (1994) Mol. Cell Endocrinol. **99**, C1–C6.
15. D. R. Wines, J. M. Brady, D. B. Pritchett, J. L. Roberts, and R. J. MacDonald (1989) J. Biol. Chem. **264**, 7653–7662.
16. P. H. Riegman, R. J. Vlietstra, J. A. van der Korput, A. O. Brinkmann, and J. Trapman (1991) Mol. Endocrinol. **5**, 1921–1930.
17. P. Murtha, D. J. Tindall, and C. Y. Young (1993) Biochemistry **32**, 6459–6464.
18. D. R. Wines, J. M. Brady, E. M. Southard, and R. J. MacDonald (1991) J. Mol. Evol. **32**, 476–492.
19. J. A. Simson, J. Wang, J. Chao, and L. Chao (1994) Lab. Invest. **71**, 680–687.
20. N. M. Kaplan and M. Silverberg (1987) Blood **70**, 1–16.
21. F. van der Graaf, G. Tans, B. N. Bouma, and J. H. Griffen (1982) J. Biol. Chem. **263**, 4698–4703.
22. J. Clements (1989) Endocr. Rev. **10**, 393–419.
23. K. X. Chai, D. C. Ward, J. Chao, and L. Chao (1994) Genomics **23**, 370–378.
24. J. X. Ma, Z. R. Yang, J. Chao, and L. Chao (1995) J. Biol. Chem. **270**, 451–455.
25. S. D. Mikolajczyk, L. S. Millar, A. Kumar, and M. S. Saedi (1998) Prostate **34**, 44–50.
26. D. Deperthes, F. Marceau, G. Frenette, C. Lazure, R. R. Tremblay, and J. Y. Dube (1997) Biochim. Biophys. Acta **1343**, 102–106.
27. L. Bourgeois, M. Brillard-Bourdet, D. Deperthes, M. A. Juliano, L. Juliano, R. R. Tremblay, J. Y. Dube, and F. Gauthier (1997) J. Biol. Chem. **47**, 29590–29595.

### Suggestions for Further Reading

28. A. R. Khan and N. G. J Michael (1998) Molecular mechanisms for the conversion of zymogens to active proteolytic enzymes. Protein Sci. **7**, 815–836.
29. W. J. Catalona, A. W. Partin, K. M. Slawin, M. D. Brawer, et al. (1998) Use of percentage of free prostate-specific antigen to enhance differentiation of prostate cancer from benign prostate disease. JAMA **279**, 1542–1547. (Discusses the latest developments and utility of hK3 (PSA) as diagnostic marker for prostate cancer.)

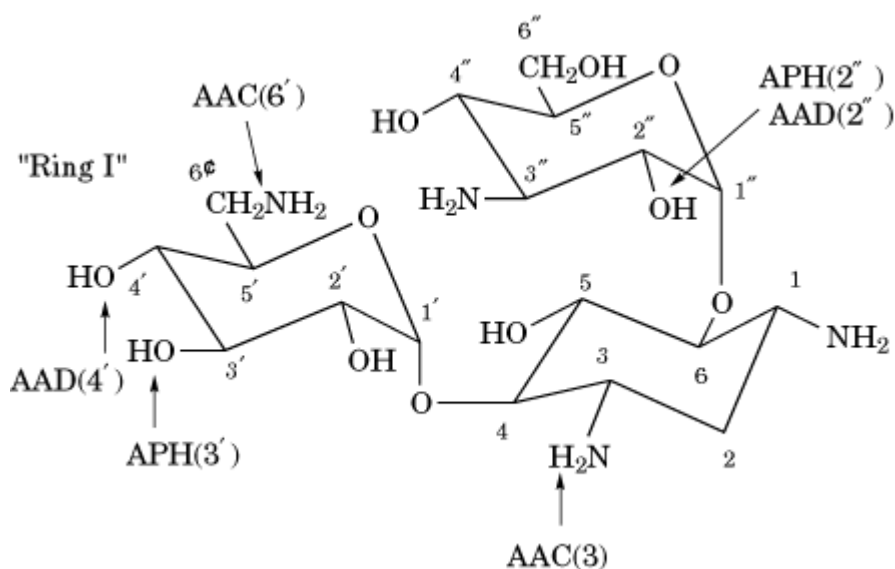
## Kanamycin

Kanamycin (KM) is an aminoglycoside antibiotic that was discovered in culture supernatants of *Streptomyces kanamyceticus* by Umezawa et al. in 1957 ([1](#)). It has been a widely used drug for the treatment of infections due to aerobic **Gram-negative** and **Gram-positive** bacteria and especially as

a second-line antibiotic for the treatment of tuberculosis. Although its usefulness had declined because of its toxicity and the appearance of KM-resistant microorganisms, in recent years, it has been revived because of the emergence of multiple **drug resistant** (MDR) *Mycobacterium tuberculosis*. In addition, derivatives of KM have proved convenient in selection systems of eukaryotic cells for molecular genetic studies.

KM is the major component (A) of three products of *S. kanamyceticus*: KM A, B, and C. It belongs to the same oligosaccharide group of water-soluble antibiotics as [streptomycin](#). It consists of two amino sugars (6-D-glucosamine and 3-D-glucosamine) linked to a centrally located 2-deoxystreptoamine moiety ( $C_{18}H_{36}N_4O_{11}$ , molecular weight 484.50) (Fig. 1) (2).

**Figure 1.** Structure of kanamycin. The sites of KM modification by various enzymes are indicated. The saccharide Ring I structure that is implicated in binding to 16S rRNA of *Escherichia coli* is indicated. Key: AAC, *N*-acetyltransferase; AAD, *O*-adenyltransferase; APH, *O*-phosphotransferase.



KM exhibits bactericidal activity by inhibiting **protein biosynthesis**. Its mechanism of transport into the cytoplasm of microorganisms is the same as that of other aminoglycoside antibiotics (see [Streptomycin](#)) (3). KM then binds to both the 30S and the 50S subunits of [ribosomes](#) in the cytoplasm (4). Recently, Fourmy et al. (5) reported that the common Ring I structure of paromomycin, KM, [neomycin](#), and gentamicin (Fig. 1) binds to the region around nucleotide 1400 of 16S ribosomal RNA (rRNA) of *Escherichia coli*. Two adenine residues at positions 1408 and 1493 of the 16S rRNA are essential for antibiotic binding and comprise a specific binding pocket; the pairing of nucleotides 1409 and 1491 of 16S rRNA provides the floor of the antibiotic-binding pocket. Prokaryotic 16S rRNA sequences have an adenine residue at position 1408, whereas this position is a guanine residue in eukaryotic sequences. Therefore, KM acts preferentially on prokaryotic organisms. Prokaryotic ribosomes are sensitive to antibiotic concentrations that are 10 to 15 times lower than that inhibiting eukaryotic ribosomes (6). KM inhibits both the initiation and the elongation of [translation](#) by fixing the peptidyl transfer RNA (tRNA) to the A site of the ribosome, and it also inhibits release of the [transfer RNA](#) (7). A complex of KM and 16S rRNA decreases the rate of dissociation of aminoacyl-tRNA and inhibits translational processivity, resulting in miscoding. KM causes more miscoding of messenger RNA template than does streptomycin (8). KM uses denatured DNA, rRNA, and tRNA as a template, resulting in abnormal protein synthesis (9, 10).

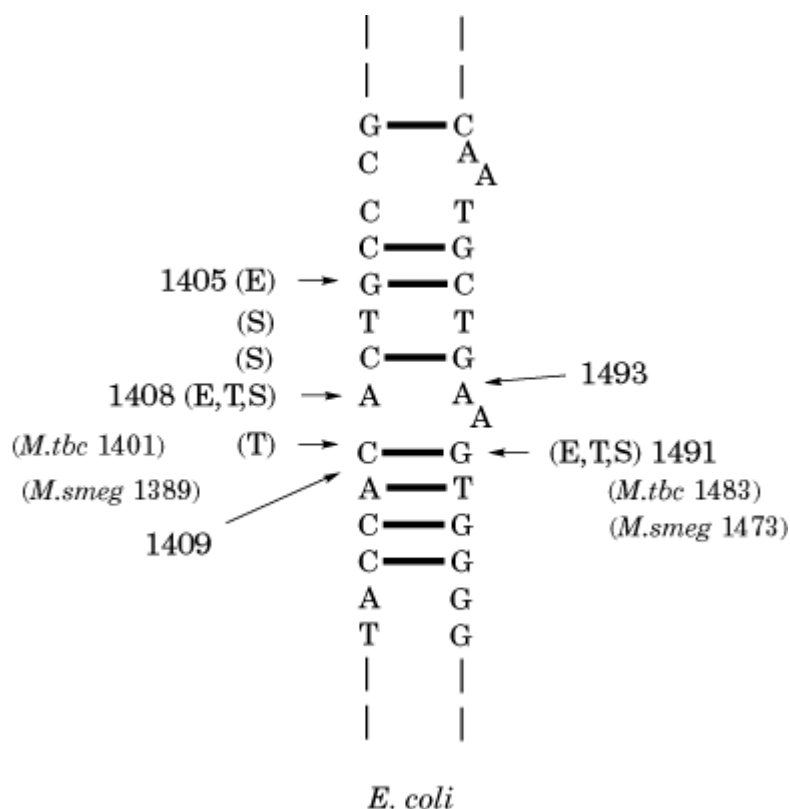
Resistance to KM has been found in many microorganisms, often concomitant with resistance to

other aminoglycosides, such as neomycin, viomycin-capreomycin, and streptomycin (11, 12). Two mechanisms of resistance to KM have been described (13, 14); one is KM modification by modifying enzymes (15-18), and the other is alteration of target sites (19). Most of the KM resistance of clinically isolated microorganisms is caused by the modification by enzymes, whereas modification enzymes have not been found in mycobacteria.

The modification enzymes are usually encoded by resistant **plasmids** and **transposons** and appear to be located in the **periplasmic** space in the cell; the modification of KM is believed to interfere with the active transport of the drugs into the cell (20). Three modification enzymes of KM have been reported. The *O*-phosphotransferase (APH) enzyme encoded by the *aph* gene (the *neo* gene) phosphorylates the 3'-OH or 2''-OH groups of KM. The *O*-adenyltransferase (AAD) enzyme encoded by the *aad* gene adenylates the 4'-OH or 2''-OH groups of KM. The *N*-acetyltransferase (AAC) enzyme encoded by the *aac* gene acetylates the 3-NH<sub>2</sub> or 6'-NH<sub>2</sub> groups of KM (Fig. 1). Recently, the *aph* gene has been widely used in selection systems of eukaryotic cells. APH enzymes that are carried on the bacterial transposons Tn5 (14), Tn1545 (15), and Tn903 (13) confer cross-resistance to KM, neomycin, and geneticin (G418), which inhibit protein biosynthesis in both prokaryotic and eukaryotic cells. Eukaryotic cells do not normally express an endogenous APH activity, but they can express the enzymes encoded by the bacterial transposon (21).

The alterations of KM target sites that produce resistance are thought to be **methylation** or substitution **mutations** of the 50S ribosomal subunit (22), the S12 protein (the *rpsL* gene) (23), and the 16S rRNA (the *rrs* gene) (24). In *E. coli*, methylation at position 1405 or 1408 of 16S rRNA causes cross-resistance to KM and gentamicin and to KM and apramycin (25), respectively, and a G→C or T substitution mutation at position 1491 of the *rrs* gene confers cross-resistance to many aminoglycoside antibiotics including KM (Fig. 2) (24). In *M. tuberculosis* (Fig. 2), an A→G mutation at position 1400 (corresponding to 1408 of *E. coli*), a C→T mutation at position 1401 (corresponding to 1409 of *E. coli*), and double mutations of C→A at position 1401 and G→T at position 1483 (corresponding to 1491 of *E. coli*) cause KM resistance (26, 27). Suzuki et al (26) reported that approximately 70% of clinically isolated KM-resistant *M. tuberculosis* strains possessed substitution mutations at these three positions. In *Mycobacterium smegmatis*, the substitution mutation A→G at position 1389 (corresponding to position 1408 of the *rrs* gene of *E. coli*) causes high-level resistance to KM (28, 29). A T→A or C mutation at position 1387 or a C→G mutation at position 1388 (corresponding to positions 1406 and 1407 of *E. coli*, respectively) causes low-level resistance. A G→A mutation at position 1473 (corresponding to 1491 of the *rrs* gene of *E. coli*) confers resistance to viomycin and capreomycin, whereas a G→T mutation at the same site shows high-level cross resistance to KM also (24). Viomycin and capreomycin belong to the water-soluble peptide group of antibiotics and have also been used as second-line drugs for the treatment of tuberculosis. KM resistance by the substitution mutation A→G at position 1389 (corresponding to position 1408 of the *rrs* gene of *E. coli*) is recessive to wild type, which was revealed by the conjugation system in *M. smegmatis* (29). Therefore, mutation of a single gene copy cannot confer the resistance phenotype to organisms possessing multiple *rrs* gene copies.

**Figure 2.** Mutations in the A site of 16S rRNA that produces KM resistance. E, T, and S indicate mutation sites in *E. coli*, *Mycobacterium tuberculosis*, and *Mycobacterium smegmatis*, respectively. The large nucleotide numbers are those of the *rrs* gene of *E. coli*; some nucleotide numbers of the other species are indicated in smaller numbers. The highlighted secondary structure at nucleotides 1408 and 1491 (and at 1409 and 1493), is important for binding KM.



The KM-producing strain, *S. kanamyceticus*, possesses the modification enzyme AAC (6') and a resistant ribosome owing to a methylated 16S rRNA (30). Hotta et al. (31) reported that a mutant of the KM-producing strain that is resistant to 100 µg/ml KM possesses a point mutation at the putative **promoter** region of a cryptic *kan* gene, and the mutation leads to the enhancement of transcription, conferring high-level resistance to KM.

KM is effective against aerobic Gram-negative bacilli, *Staphylococcus*, and mycobacteria. However, KM is not effective against *Pseudomonas aeruginosa*, *Serratia* sp., *Streptococcus pyogenes*, *Streptococcus pneumoniae*, anaerobic organisms, and fungi. Therefore, new aminoglycosides, eg, gentamicin (32) and amikacin (33), are being used now instead of KM (3). KM is still used, however, for the treatment of infections caused by aerobic Gram-negative bacteria, such as dysentery and gonorrhea, and by MDR-*Staphylococci* and *M. tuberculosis*. In addition, KM is an important agent against septicemia and renal infection due to Gram-negative bacilli. Because of the appearance of MDR-*M. tuberculosis* strains against first-line drugs such as isoniazid (INH), rifampicin, or streptomycin, KM is being revived as an important therapeutic agent.

For therapeutic applications, kanamycin sulfate is usually taken by injection or orally. Its absorption, distribution, elimination, and adverse effects are the same as those of streptomycin. Its gastrointestinal absorption is poor, and it is not metabolized but is excreted by the kidney. Its concentration in the blood depends on the patient's renal function.

KM shows nephrotoxicity and auditory or vestibular toxicity similar to that seen with other aminoglycosides. KM especially causes damage to the eighth cranial nerve and is associated with severe auditory toxicity. Because of its nephrotoxicity, its use must be limited for patients with renal impairment.

#### Bibliography

1. H. Umezawa, M. Ueda, K. Maeda, K. Yagishita, S. Konda, Y. Okami, R. Utahara, T. Osato, K.

- Nitta, and T. Takeuchi (1957) *J. Antibiot.* **10A**, 181.
2. *The Merck Index: An Encyclopedia of Chemicals, Drugs, and Biologicals* (1996). (S. Buda-vari, ed.) Merck Research Laboratories, Whitehouse Station, N.J., pp. 900–901.
  3. G. L. Mandel and W. A. Petri Jr. (1996) In *Goodman & Gilman's The Pharmacological Basis of Therapeutics* (J. G. Hardman and L. E. Limbird, eds.), McGraw-Hill, New York, pp. 1159–1161.
  4. M. Misumi, T. Nishimura, T. Komai, and N. Tanaka (1978) *Biochem. Biophys. Res. Commun.* **84**, 358–365.
  5. D. Fourmy, M. I. Recht, S. C. Blanchard, and J. D. Puglisi (1996) *Science* **274**, 1367–1371.
  6. J. M. Wilhelm, S. E. Pettitt, and J. J. Jessop (1978) *Biochemistry* **17**, 1143–1149.
  7. M. Misumi and N. Tanaka (1980) *Biochem. Biophys. Res. Commun.* **92**, 647–654.
  8. J. Davies and B. D. Davis (1968) *J. Biol. Chem.* **243**, 3312–3316.
  9. B. J. McCarthy and J. J. Holland (1965) *Proc. Natl. Acad. Sci. USA* **54**, 880–886.
  10. H. Masukawa and N. Tanaka (1966) *J. Biochem.* **62**, 202.
  11. D. Apirion and D. Schlessinger (1968) *J. Bacteriol.* **96**, 768–776.
  12. K. Suga and Y. Mizuguchi (1974) *Jap. J. Microbiol.* **18**, 139–147.
  13. H. Umezawa and S. Kondo (1983) *Clinical chemotherapy* (H. P. Kiemmerle, ed.), Thieme-Stratton Inc., New York, pp. 120–146.
  14. H. Umezawa and S. Kondo (1982) *Handbook of Experimental pharmacology* (H. Umezawa and I. R. Hooper, eds.), Springer-Verlag, Berlin, pp. 267–292.
  15. A. Oka, H. Sugisaki, and M. Takanami (1981) *J. Mol. Biol.* **147**, 217–226.
  16. P. Mazodier, P. Cossart, E. Giraud, and F. Gasser (1985) *Nucleic Acids Res.* **13**, 195–205.
  17. P. Courvalin and C. Carlier (1987) *Mol. Gen. Genet.* **206**, 259–264.
  18. B. G. Spratt (1994) *Science* **264**, 388–393.
  19. J. Davies (1994) *Science* **264**, 375–382.
  20. J. Davies and D. I. Smith (1978) *Annu. Rev. Microbiol.* **32**, 469.
  21. T. Maniatis, E. F. Frisch, and J. Sambrook (1982) *Molecular Cloning*, Cold Spring Harbor Laboratory Press, New York.
  22. E. C. Choi, T. Nishimura, and N. Tanaka (1980) *Biochem. Biophys. Res. Commun.* **94**, 755–762.
  23. H. Masukawa (1969) *J. Antibiot.* **22**, 612–623.
  24. E. A. De Stasio, D. Moazed, H. F. Noller, and A. E. Dahlberg (1989) *EMBO J.* **8**, 1213–1216.
  25. A. A. Beauclerk and E. Cundliffe (1987) *J. Mol. Biol.* **193**, 661–671.
  26. Y. Suzuki, C. Katsukawa, A. Tamaru, C. Abe, M. Makino, Y. Mizuguchi, and H. Taniguchi (1998) *J. Clin. Microbiol.* **36**, 1220–1225.
  27. G. J. Alangaden, B. N. Kreiswirth, A. Aouad, M. Khetarpal, F. R. Igno, S. L. Moghazeh, E. K. Manavathu, and S. A. Lerner (1998) *Antimicrob. Agents Chemother.* **42**, 1295–1297.
  28. P. Sander, T. Prammananan, and E. C. Bottger (1996) *Mol. Microb.* **22**, 841–848.
  29. H. Taniguchi, B. Chang, C. Abe, Y. Nikaido, Y. Mizuguchi, and S. Yoshida (1997) *J. Bacteriol.* **179**, 4795–4801.
  30. E. Cundliffe (1989) *Annu. Rev. Microbiol.* **43**, 207–233.
  31. K. Hotta, J. Ishikawa, T. Ogata, and S. Mizuno (1992) *Gene* **115**, 113–117.
  32. J. P. Roselot, J. Marquez, E. Meseck, A. Murawski, A. Hamdan, C. Joyner, R. Schmidt, D. Migliore, and H. L. Herzog (1964) In *Antimicrobial Agents and Chemotherapy-1963* (J. C. Sylvester, ed.), American Society for Microbiology, Ann Arbor, Mich., pp. 14–16.
  33. H. Kawaguchi, T. Naito, S. Nakagawa, and K. I. Fujisawa (1972) *J. Antibiot. (Tokyo)* **25**, 695–708.



## Kappa (k) and Lambda (l) Light Chains

The distinction between k and l light chains of [immunoglobulins](#) (Igs) was first based on serological arguments, because it was observed that [antibodies](#) raised against [Bence-Jones proteins](#) (BJPs) could distinguish two [serotypes](#) that did not cross-react and were later termed K and L. These [antigen](#) specificities were also present on [myeloma proteins](#), and everything became clarified when it was shown that BJPs were the Ig light chains that were also expressed in the corresponding myeloma protein, which was, of course, a complete immunoglobulin. From a historical standpoint, it is interesting to recall that such a serological approach was one of the first tools that allowed immunologists to gain evidence in favor of the clonal theory of [immunoglobulin biosynthesis](#); it also offered an elegant way to refine analysis of [immunoglobulin structure](#). It was, for example, clearly demonstrated that monoclonal myeloma proteins had either k chains or l chains and never both, bringing the first indication that Igs were symmetrical molecules, which is a requirement for the clonal theory. Because the K and L serotypes were identified on discrete Ig classes, it also clearly demonstrated the ubiquitous nature of k and l chains, as opposed to heavy chains that were class-specific. When Ig gene organization became known, it was demonstrated that both types of light chain, which constitute discrete [isotypes](#), were encoded by two separate loci, IGK and IGL, located on human chromosomes 2 and 22, respectively. Although the serologic K or L “characters” were due to [epitopes](#) located on the constant regions of either chains, the “kappaness” and the “lambdaness” are also present on the **variable regions**, which have distinct primary sequences as encoded by separate sets of [V genes](#).

The relative contribution of k and l chains to the overall pool of light chains expressed in immunoglobulin molecules is highly variable from one species to another. For example, horses and cattle possess only l chains, whereas k is present on 95% of Ig in the mouse. In humans, the k/l ratio is about 3/2. There is no obvious functional reason for this, because both k- and l-containing Igs are equally “good” antibodies; It appears simply to be related directly to the number of V genes. In humans, this number is somewhat balanced, whereas in mice there are only three V<sub>l</sub> genes for more than 50 V<sub>k</sub>. These differences also indicate that [gene duplication](#) and amplification are species characters, and there is no clear evidence of **vertical transmission** of the multigenic organization of V genes.

In mouse and humans, there is only one C<sub>k</sub> gene, so there are no isotypes of k chains. There is, however, some degree of allelic **polymorphism** in humans, due to the presence of three K<sub>m</sub> allotypes that result from point mutations at two positions of the constant region. In the rabbit, several **allotypes** have been described that correlate with multiple amino acid substitutions in the constant region. The situation is different for the l chain, as a consequence of the organization of the C<sub>l</sub> genes. In the mouse, in which four J<sub>l</sub>–C<sub>l</sub> tandem genes have been identified, only three isotypes are expressed in Ig, because no I<sub>2</sub> was found. No allelic variant has been reported. In humans, four C<sub>l</sub> functional genes have been identified, also as part of J<sub>l</sub>–C<sub>l</sub> tandems. They are expressed as four discrete l isotypes that may also be serologically distinguished.

Based on amino acid sequence homologies of the variable regions, k and l chains have long been grouped in families that contain variable numbers of genes, among which some are more frequently expressed. Selective pressure that is responsible for this preferential expression is not clear, although some arguments would suggest that light chains that have the highest binding affinity for heavy chains would be favored.

See also entries [Bence-Jones Proteins](#), [Immunoglobulin](#), and [Isotype](#).

### Suggestions for Further Reading

T. Kirschbaum, R. Jaenichen, and H. G. Zachau (1996) The mouse immunoglobulin kappa locus contains about 140 variable gene segments. *Eur. J. Immunol.* **26**, 1613–1620.

H. G. Zachau (1993) The immunoglobulin kappa locus-or-what has been learned from looking closely at one-tenth of a percent of the human genome. *Gene* **135**, 167–173.

J.-P. Fripiat et al. (1995) Organization of the human immunoglobulin lambda light-chain on chromosome 22q11.2. *Hum. Mol Genet.* **4**, 983–991.

## Karyogamy

In sexual **reproduction**, the final objective is the creation of new unique individuals. In the [fertilization](#) process, two morphologically disparate **gametes**, the [egg](#) and [sperm](#), recognize, bind, and fuse with each other. After this, the [male](#) and [female](#) haploid [genomes](#) unite, which is called karyogamy. Most of this process is directed by maternal **gene** products, as there is little or virtually no gene [transcription](#) in the [nucleus](#) of the **oocyte** at this early stage of nuclear fusion. It was shown in cows that protein biosynthesis is not required for formation of the male pronucleus, since it was not blocked in the presence of the protein synthesis inhibitor [cycloheximide](#) (1).

Unification of the [haploid](#) genomes is accompanied by lively cytoplasmic movements. However, the mechanism of how the male and female **pronuclei** are brought together is still the subject of speculation. It has been proposed that a specific  $\alpha$ -[tubulin](#) isoform is necessary for the normal cellular function of a [kinesin](#) motor protein. Embryos with the **genotype**  $\alpha$ -Tub67C/+ showed delayed **meiosis** and defective female pronucleus formation (2). In artificially activated ova of sea urchins, migration of the female pronucleus was observed in the absence of a male pronucleus, indicating that male components are not required.

With regard to karyogamy, two different types of fertilization are observed, either (1) the male and female chromosomes intermix without fusion of the pronuclei (ascaris type) or (2) the pronuclei fuse first, after which the newly formed nucleus undergoes **mitosis** (sea urchin type).

1. **Ascaris type.** In the ascaris type of fertilization, which is found in some mollusks, a number of mammals, and humans, insemination occurs prior to the completion of meiosis. The male and female pronuclei, which have spheroid shapes, become closely apposed and form a nucleoplasmic projection. **Asters** become situated to either side of the associated pronuclei and form the prospective poles of the **mitotic spindle**. [Chromosomes](#) start to condense at this time and become visible. Then the pronuclear envelopes of the male and female pronuclei break down. The paternal and maternal chromosomal groups mix with each other and proceed to the metaphase stage of mitosis. This means that the parental genomes will be associated within a single nucleus for the first time in [blastomeres](#) at the two-cell stage.
2. **Sea urchin type.** The sea urchin type of fertilization differs from the ascaris type with respect to the time the egg is inseminated. Meiotic divisions have been already completed in the oocyte when a spermatozoon enters. After swelling of the male pronucleus, both nuclei are brought together, and pronuclear association occurs. The male and female pronuclear envelopes fuse and form an internuclear bridge, by which the nuclei join together and form a zygote nucleus.

New investigations in humans using time-lapse video cinematography after intracytoplasmic sperm injection have provided a clearer view concerning the events during pronuclear formation and early development (3). Circular waves of granulation within the oocyte cytoplasm were observed that had a periodicity of 20 to 53 min. During the granulation phase, the sperm head decondenses and the second polar body is extruded, followed by the central formation of the male pronucleus. The female pronucleus is formed in the cytoplasm at the same time or slightly after the formation of the male pronucleus, and they are then drawn together. In mice, modulation of either male or female pronuclear formations is proposed to be dependent on [imprinting](#) and on the developmental consequences of early asynchrony between male and female pronuclear development (4). Organelles contract from the cortex toward the center of the oocyte (3).

One of the most important questions is how the egg is activated. Among the possible mechanisms, it is generally assumed that a protein released by the spermatozoon triggers elevation of calcium levels in the oocyte cytoplasm. This hypothesis is supported by experiments in which purified human cytosolic factors were injected into mouse oocytes and induced oscillations in calcium concentration. Moreover, human oocytes that remain unfertilized after the intracytoplasmic injection of spermatozoa are activated, and some of them contained three or more pronuclei (5).

### Bibliography

1. R. C. Chian and M. A. Sirard (1996) *Zygote* **4**, 41–48.
2. D. J. Komma and S. A. Endow (1997) *J. Cell Sci.* **110**, 229–237.
3. D. Payne, S. P. Flaherty, M. F. Barry, and C. D. Matthews (1997) *Hum. Reprod.* **12**, 532–541.
4. J. Fulka Jr., M. Horska, R. M. Moor, J. Fulka, and J. Kanka (1996) *Dev. Biol.* **178**, 1–12.
5. G. D. Palermo, O. M. Avrech, L. T. Colombero, H. Wu, Y. M. Wolny, R. A. Fissore, and Z. Rosenwaks (1997) *Mol. Hum. Reprod.* **3**, 367–374.

### Suggestions for Further Reading

6. J. Y. Nothias, S. Majumder, K. J. Kaneko, and M. L. DePamphilis (1995) Regulation of gene expression at the beginning of mammalian development. *J. Biol. Chem.* **270**, 22077–22080. Describes the regulation of gene expression at the beginning of mammalian development.
7. F. J. Longo, ed. (1997) *Fertilization*, 2nd ed., Chapman & Hall, London, Weinheim, New York, Tokyo, Melbourne, Madras, p. 181–211. Highly recommended to gain deeper insight into the process of karyogamy and plasmogamy.

## Karyoplasm

In Greek, “karyo” means nut ([nucleus](#)) and “plasm” means substance. Karyoplasm is the nuclear substance, in contrast to the **cytoplasm** on the outside of the [nuclear envelope](#). The term karyoplasm was first used by Strasburger in 1882 to differentiate the nuclear plasm from the cytoplasm (1). Although the terms are perfectly synonymous, [nucleoplasm](#) is more often used today than karyoplasm. Nevertheless, the term “karyophylic” is commonly used to refer to aspects of the nucleus.

### Bibliography

1. E. B. Wilson (1925) *The Cell in Development and Heredity*, MacMillan, New York.

## Karyotype

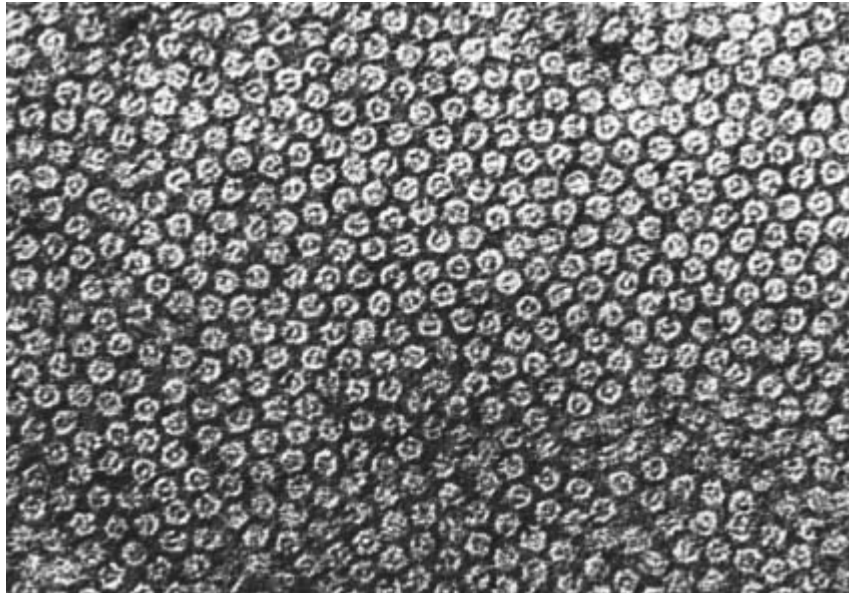
The karyotype is the complete **chromosomal** set of the [nucleus](#) of the cell. A diagrammatic representation of the karyotype with all of the pairs of chromosomes arranged in order of size is called an [ideogram](#). Every organism has a standard karyotype, which provides a frame of reference for the analysis of [mutations](#). Karyotypic analysis is the study of all of the visible traits of chromosomes in a typical cell. This is useful in studying speciation. For example, distinct species of *Drosophila* that are very similar in appearance are distinguished easily by karyotypic analysis, which reveals numerous chromosomal **inversions** and [translocations](#).

## Keratins

There are three groups of keratins: the a-keratins, the b-keratins, and the feather keratins ([1](#), [2](#)). The a-keratins can be subdivided into (i) the *hard* a-keratins of hair, nails, claws, beaks, quills, hooves, baleen, and horns and (ii) the (*soft*) epidermal or [cytokeratins](#) of the *stratum corneum*, corns, and calluses. The b-keratins are derived from the a-keratins as a result of pressure and temperature, and consequently do not represent an *in vivo* structure. Feather keratins are found in feathers and scales and in parts of claws and beaks.

Hard a-keratin and epidermal cytokeratin form filamentous structures known as [intermediate filaments](#) (IFs) with diameters of about 10 nm (Fig [1](#)). The cytokeratins are discussed separately under the entry [Cytokeratins](#). Hard a-keratin IF are embedded in a matrix of proteins that are rich in [disulfide bonds](#), which endow this form of keratin with much of its special mechanical properties ([1](#), [2](#)). The hard a-keratin chains in the IF consist of nonhelical *N*- and *C*-terminal **domains** separated by a region of **a-helical** sequence dominated by a [heptad repeat](#) that aggregates in a [coiled-coil](#). The terminal domains are rich in [cysteine](#) residues, and these too form disulfide bonds *in vivo*, either with the matrix proteins or with other intermediate filament proteins. For details of the coiled-coil structure of the rod domain, as well as the structural hierarchy, see the entries on [Cytokeratins](#) and on [Intermediate Filaments](#).

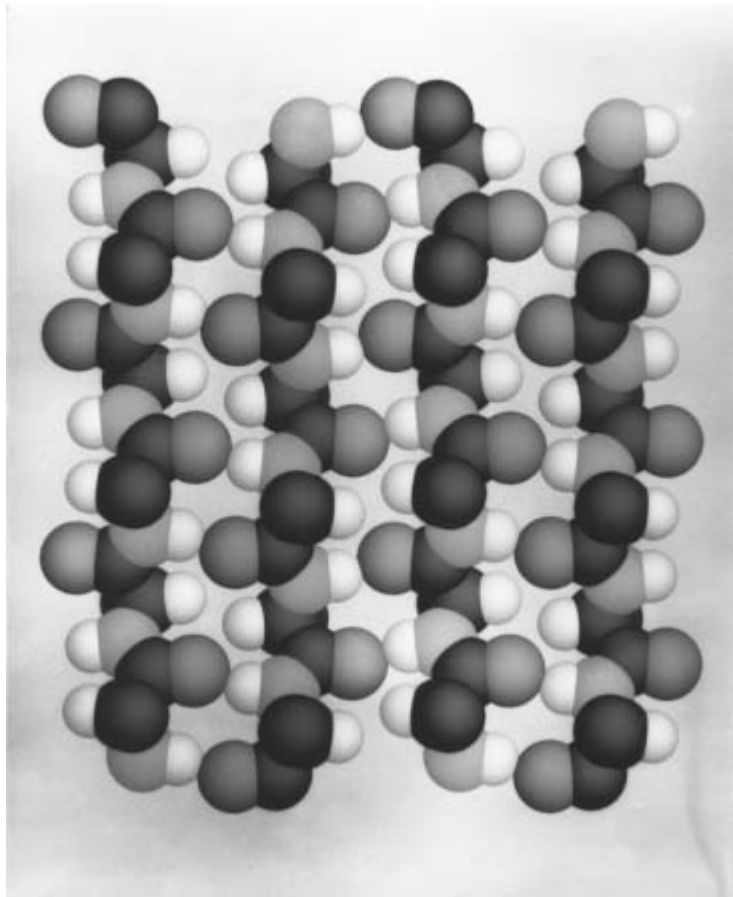
**Figure 1.** Electron micrograph of a transverse section of hard a-keratin intermediate filaments (IFs). The IFs are about 8 nm in diameter and are embedded in an osmiophilic matrix. (Courtesy of G. E. Rogers and R. D. B. Fraser.)



The matrix proteins in hard  $\alpha$ -keratin, which constitute the group of intermediate filament-associated proteins, can be categorized into three classes: the high-sulfur proteins, the ultra-high-sulfur proteins, and the high-tyrosine proteins. Each is a heterogeneous group, and all are involved in a vast network of covalent disulfide bonds (largely intramolecular) that stabilize the assembly of molecules into a rigid and mechanically tough structure (3). The high-sulfur proteins can be placed in either the B1, B2, BIIIA, or BIIIB families. The B2 family contains two 75-residue domains, each consisting of tandem pentapeptide repeats, linked by a 22-residue glycine-rich structure (4). The BIII family contains one of these 75-residue blocks and a glycine-rich sequence close to its *N*-terminus (4). The pentapeptide repeat has the idealized sequence Cys–Gln/Arg–Pro–Ser/Thr–Cys and is believed to form a  $\beta$ -bend, stabilized by a disulfide bond between the two cysteine residues four apart. This results in a structure in which the majority of disulfide bonds are intramolecular and in which there are well-defined and rigid  $\beta$ -bend “knots,” with relatively free rotation about the single peptide bond that connects one to another. The BIIIB family does not contain any repeating motifs, nor does it have any similarity with the B2 or BIIIA families of proteins. The high-tyrosine family can be divided into type I and type II subgroups. In these cases, however, little evidence of structural regularity can be recognized in the amino acid sequences. The content of all the matrix proteins varies greatly from one source of hard  $\alpha$ -keratin to another. Indeed, they can vary along the length of a single fiber, as well as with the nutritional state of the animal (3). There is no matrix in epidermal keratin, but filaggrin—an intermediate filament-associated protein—plays an important part in assembling IF into functional macroaggregates (5).

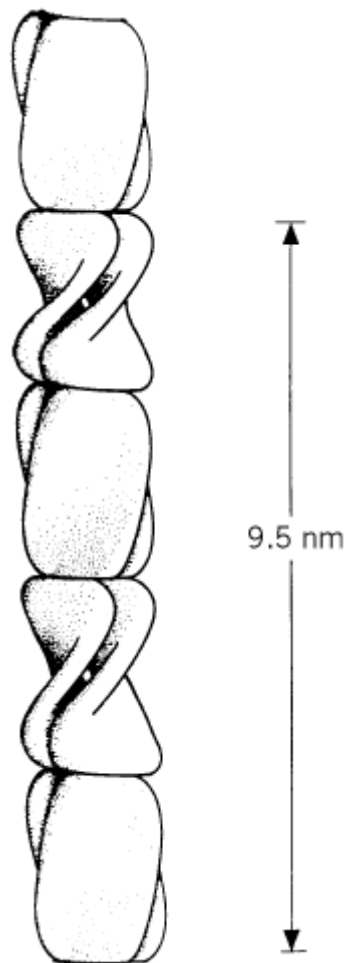
The refined structure of  $\beta$ -keratin, which is derived from the  $\alpha$ -keratins as a result of pressure and temperature, shows a regular packing of antiparallel  $\beta$ -strands within a  $\beta$ -sheet, but with the  $\beta$ -sheets displaying a stacking disorder (6) (Fig 2). There are about 10 chains per sheet and two or three sheets per  $\beta$ -crystallite, giving a maximum of about 30 chains. This value closely matches that estimated from the studies on keratin IF and strongly implies that each  $\beta$ -crystallite is derived from a single IF and that the latter contains equal numbers of up- and down-pointing chains (as is indeed now established independently from studies on the  $\alpha$ -keratin IF).

**Figure 2.** Space-filling model of  $\beta$ -keratin showing antiparallel chains forming one of the two to three  $\beta$ -pleated sheets present. The axial rise per residue is 0.334 nm, the average chain separation is 0.47 nm, and the sheets are 0.97 nm apart. (From R. D. B. Fraser and T. P. MacRae, with permission.)



Feather and scale keratins are similar at both the ultrastructural and molecular level. Each consists of filaments about 3.3 nm in diameter embedded in a “matrix.” However, each is composed of only a single species of polypeptide chain, implying that the filament and matrix arise from different parts of the same protein chain. The molecular weights of the constituent chains do differ significantly, however (about 10 and 15 kDa in feather and scale, respectively). From X-ray [fiber diffraction](#), infrared [vibrational spectroscopy](#), and modeling studies, a detailed structure has been derived for feather keratin in which successive subunits are related to one another by a left-handed fourfold screw axis (7). Thus consecutive units display a relative axial displacement of 2.36 nm and a unit twist of  $-90^\circ$ , giving rise to a helix of pitch length 9.5 nm. Each structural unit comprises two molecular chains that associate in an antiparallel manner via regions of sequence that adopt a four-stranded, antiparallel  $\beta$ -sheet structure (Fig 3). This feature is readily recognizable in the amino acid sequence of both feather and scale keratin chains as a 30-residue element containing proline residues about eight apart, and with alternating nonpolar residues in between the prolines. This gives rise to a  $\beta$ -sheet with an apolar face. Packing of this face from two  $\beta$ -sheets through apolar interactions stabilizes the unit and provides the core of the filament. The remainder of the sequence provides the matrix in which the filament appears to be embedded. Scale keratin has the same filament structure as feather keratin, but the degree of order in the packing of the filaments is greater. The sequence difference between these two proteins is essentially confined to a single region in which scale (but not feather) has a fourfold tandem repeat of 13 residues. Each repeat is predicted to form a pair of antiparallel  $\beta$ -strands that puts [tyrosine](#) residues in positions where they can interact optimally with those in similar sheets of other molecules (8). It is clear that feather and scale keratins have a common ancestor. Evolutionarily, scales preceded feathers, and the probability is that feathers evolved from scales as a consequence of deletion of the fourfold tandem repeat.

**Figure 3.** Model of the structure of feather keratin. Each molecule contains four strands of antiparallel twisted  $\beta$ -sheet that interacts **isologously** in an antiparallel manner with the same region in a second molecule. These subunit pairs aggregate **heterologously** and are related to one another by the operation of a fourfold screw axis. The  $\beta$ -sheets thus form the core of the filaments of diameter 3.3 nm, whereas the remainder of the polypeptide chains form the “matrix” in which the filaments are present. (From Ref. 8, with permission.)



### Bibliography

1. R. D. B. Fraser, T. P. MacRae, and G. E. Rogers, (1972) *Keratins: Their Composition, Structure and Biosynthesis*, Charles C Thomas, Springfield, Il.
2. R. D. B. Fraser and T. P. MacRae (1973) *Conformation in Fibrous Proteins and Related Synthetic Polypeptides*, Academic Press, London.
3. J. M. Gillespie (1990) "The Proteins of Hair and Other Hard -Keratins". In *Cellular and Molecular Biology of Intermediate Filaments* (R. D. Goldman and P. M. Steinert, eds), Plenum Press, New York, pp. 95–128.
4. D. A. D. Parry, R. D. B. Fraser, and T. P. MacRae (1979) Repeating patterns of amino acid residues in the sequences of some high-sulphur proteins from -keratin. *Int. J. Biol. Macromol.* **1**, 17–22.
5. J. W. Mack, A. C. Steven, and P. M. Steinert (1993) The mechanism of interaction of filaggrin with intermediate filaments: the ionic zipper hypothesis. *J. Mol. Biol.* **232**, 50–66.
6. R. D. B. Fraser, T. P. MacRae, D. A. D. Parry, and E. Suzuki (1969) The structure of -keratin. *Polymer* **10**, 810–826.
7. R. D. B. Fraser, T. P. MacRae, D. A. D. Parry, and E. Suzuki (1971) The structure of feather keratin. *Polymer* **12**, 35–56.

8. K. Gregg, S. D. Wilton, D. A. D. Parry, and G. E. Rogers (1984) A comparison of genomic coding sequences for feather and scale keratins: structural and evolutionary implications. *EMBO J.* **3**, 175–178.

### Suggestions for Further Reading

9. D. A. D. Parry (1996) "Keratins". In *Polymeric Materials Encyclopedia*, Vol. **5**, (J. C. Salamone, ed), CRC Press, Boca Raton, FL, pp. 3515–3523.
10. D. A. D. Parry and P. M. Steinert (1995) *Intermediate Filament Structure*, Springer-Verlag, Heidelberg, pp. 1–193.

### Kex-2 Gene

*Kex-2* is a **gene** of the yeast *Saccharomyces cerevisiae* that was detected because mutations in it cause a deficit in the secretion of killer toxin. *Kex-2* has been **cloned**, sequenced, and characterized (1). Kexin, the *Kex-2* gene product, is a calcium-dependent [serine proteinase](#) of the [subtilisin](#) family that possesses a signature [catalytic triad](#) of three particular Asp, Ser, and His residues. The “charge-relay” system that these three residues comprise is similar to that found in the [trypsin](#) serine proteinases, but the remainder of the structure is completely different, and subtilisin and kexin are not **homologous** to the trypsin family. Kexin was the first endoproteinase identified to be responsible for limited [proteolytic](#) processing at specific pairs of basic residues in [polyprotein](#) precursors to generate bioactive [peptide hormones](#). Kexstatin is a strong inhibitor of this *Kex-2* proteinase (2-4).

*Kex-2* has been a fundamental discovery, because the [protein secretion](#) pathways of yeast share many conserved structural features, functional similarities, and common mechanisms with mammalian cells. Thus *Kex-2* has led to the isolation and characterization of the corresponding mammalian gene, *FUR*, and its product, furin (5). The isolation of both of these genes has proven useful in the study of the molecular mechanisms involved in polyprotein processing. Furthermore, there is sequence homology between kexin and the mammalian pro-hormone convertases (PCs), also called subtilisin-like proprotein convertases (SPCs), or paired basic amino acid cleaving enzymes (PACES), which has led to the isolation of at least seven related proteinases by cloning of their [complementary DNA](#). These include furin/PC-1, PC-3/PACE, PC-2, Pace-4/PC-7, PC-4, PC-5/PC-6, and PC-8/LPC, which are also known as SPC-1 to SPC-7, respectively (6, 7).

Both heterologous and homologous signals can be used to direct the secretion of [recombinant proteins](#) in yeast. The expression of heterologous proteins in the fungus yeast has the advantage that this organism is capable of glycosylating proteins. Other advantages of the yeast expression system include well-characterized inducible **promoters**, [posttranslation modification](#) pathways similar to those of mammals, and high-level expression and secretion of recombinant proteins, approaching the range of 1 gram per liter of culture.

The examples of somatostatin, [insulin](#), and phospholipase A-2 (see [Lipases](#)) generated in yeast demonstrated that the yeast secretory apparatus (its [secretory granules/vesicles](#)) gives correct and efficient [transcription](#), [translation](#), processing, and secretion of these proteins (8, 9). In addition, [transfection](#) of the yeast *Kex-2* gene into mouse fibroblasts suggests that kexin is active in the [Golgi apparatus](#) (10). Thus, the yeast *Kex-2* gene is becoming a useful tool, not only for engineering the cloning and expression of [recombinant proteins](#) but also for studying the maturation of pro-hormones and other precursor-derived proteins and, in particular, for studying the molecular mechanisms of proteolytic processing in the secretory pathways of bioactive peptide hormones (see [Peptide](#)



[Hormones](#)).

## Bibliography

1. E. Stimson, Y. D. Julius, A. Blake, L. Blair, R. Kussinawa, and J. Thorner (1984) *Cell* **37**, 1075–1089.
2. R. Fuller, R. F. Sterne, and J. Thorner (1988) *Annual Rev. Physiol.* **50**, 345–362.
3. R. S. Fuller, A. J. Blake, and J. Thorner (1989) *Science* **246**, 482–486.
4. K. Mizuno, T. Nakamura, T. Ohshima, S. Tanaka, and H. Matsuo (1988) *Biochem. Biophys. Res. Commun.* **156**, 246–254.
5. R. Leduc and J. B. Denault (1966) *FEBS Lett.* **379**, 113–116.
6. Y. Rouille, S. J. Duguay, K. Lund, M. Furuta, Q. Gong, G. Lipkind, A. A. Oliva, S. J. Chan, and D. F. Steiner (1995) *Frontiers of Neuroendocrinology*, **16**, 322–361.
7. N. D. Rawlings and A. J. Barrett (1994) *Method. Enzymol.* **224**, 19–61.
8. Y. Bourbonnias, D. Bolin, and D. Shields (1988) *J. Biol. Chem.* **263**, 15342–15347.
9. L. Thim, M. T. Hansen, K. Norris, I. Hoegh, E. Boel, and J. Forstrom (1986) *Proc. Natl. Acad. Sci. USA* **83**, 6766–6770.
10. D. Germain, L. Zollinger, C. Racine, F. Gossard, D. Dignard, D. Y. Thomas, P. Crine, and G. Boileau (1990) *Molecular Endocrinology* **4**, 1572–1579.

## Kinase

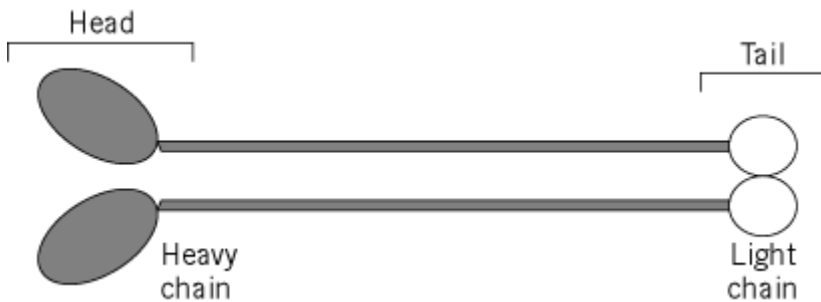
The term kinase is a trivial name that has been given to individual [enzymes](#) that belong to the group of [phosphotransferases](#). These are the enzymes that catalyze the reversible transfer of the g-phosphoryl group of a nucleoside triphosphate (NTP, usually ATP or GTP) to an acceptor molecule, which may be an alcohol, carboxylic acid, a nitrogenous compound, or a phosphorylated compound. The full description of such an enzyme is NTP:acceptor phosphotransferase. As the term kinase implies that a reaction involves phosphoryl group transfer, it is unnecessary to use the term phosphokinase to describe the reaction. The latter term has been used quite widely in connection with the reaction catalyzed by ATP:creatine phosphotransferase, or creatine kinase.

## Kinesin

Kinesin is a **microtubule**-based [motor protein](#) that transports organelles to the plus ends of microtubules. Kinesin was first identified in neural tissue as the protein responsible for anterograde fast axonal transport, the movement of membrane bounded organelles toward the synapse in neurons. The molecule has a native molecular weight of ~360kDa and is a heterotetramer composed of two identical heavy chains of ~120kDa and two accessory 64-kDa light chains. The heavy chain has an amino-terminal globular **domain** of ~40kDa, a central region that forms [coiled coils](#), and a **C-terminal** region. The central regions of two heavy chains dimerize, and the kinesin dimer is an ~100nm rod that has two globular domains at one end (the heads) and a fan-shaped tail at the other (Fig. [1](#)). The light chains are located at the tail. The head domains generate the motive force, bind

microtubules, and hydrolyze ATP. The heavy chain tail and the light chains are important for binding kinesin to cargo. In the presence of the nonhydrolyzable ATP analog AMPPNP (5' adenylylimidodiphosphate), kinesin binds tightly to microtubules. This property is used to purify kinesin from other microtubule-binding proteins. When kinesin is attached to a glass coverslip to which microtubules and ATP are added, microtubules can be observed moving along the surface of the coverslip. In this *in vitro* microtubule gliding assay, kinesin moves microtubules at a velocity of ~0.5microns/second.

**Figure 1.** Schematic illustration of the overall structure of a kinesin.



**Immunofluorescence microscopy** with [antibodies](#) to the heavy and light chains localizes kinesin to membranous organelles in cultured cells. Antibodies to the kinesin heavy chain inhibit organelle movement *in vivo* (1), and **antisense** probes to the kinesin heavy chain **gene** block anterograde axonal transport in cultured neurons (2). Mutations of the kinesin heavy chain also disrupt axonal transport (3).

Kinesin was the first member of the rapidly growing kinesin [superfamily](#) identified. The different members of the superfamily are sometimes called Kifs (kinesin family members). Each family member shares the conserved heavy chain motor (head) domain, which binds microtubules and hydrolyzes ATP. However, the family members differ in the structure of the rest of the heavy chain and the composition of the accessory proteins. The location of the motor head domain is not conserved. The head domain of different family members can be found at either the amino terminus, the carboxy terminus, or in the center of the polypeptide chain. The archetypal kinesin is involved only in movement of membrane-bounded organelles, but many of the other kinesin family members are involved in assembling and maintaining the spindle and in chromosomal movement. Several of the kinesin family members have been well characterized, and a description of three members demonstrates the variety of molecules and their functions. The BimC family of kinesins is important for chromosomal separation in mitosis and in maintaining spindles. They are composed of four ~120kDa heavy chains that form a bipolar (antiparallel) minifilament that has motor domains at both ends of the filament. This homotetramer has a “slow” motor activity that moves microtubules at 0.04 microns/second in an *in vitro* motility assay (4). Members of the kinesin II family are heterotrimers composed of two distinct motor molecules (~85 and 95 kDa) and a third noncatalytic subunit (~120kDa) in a 1:1:1 ratio. These molecules move protein complexes, “rafts,” along the inner surface of the flagellar membrane and are important for flagellar function (5). The Kar3 family of kinesins has motor domains at the carboxy terminus of the heavy chain, and they are the exception to the rule that kinesins move toward the plus ends of microtubules. They generate force toward the minus ends of microtubules. These proteins are required for forming and maintaining the meiotic and mitotic spindles. How the Kar3 head domain moves in the opposite direction along microtubules compared to the other kinesin head domains is under investigation (6).

Bibliography

1. J. Lippincott-Schwartz et al. *J. Cell Biol.* **128**, 293–306.
2. A. J. Ferreira et al. (1992) *J. Cell Biol.* **117**, 595–606.
3. D. D. Hurd and W. M. Saxton (1996) *Genetics* **144**, 1075–1085.
4. A. S. Kshina et al. (1996) *Nature* **379**, 270–272.
5. J. M. Scholey (1996) *J. Cell Biol.* **133**, 1–4.
6. U. Henningsten and M. Schliwa (1997) *Nature* **389**, 93–96.

### Suggestions for Further Reading

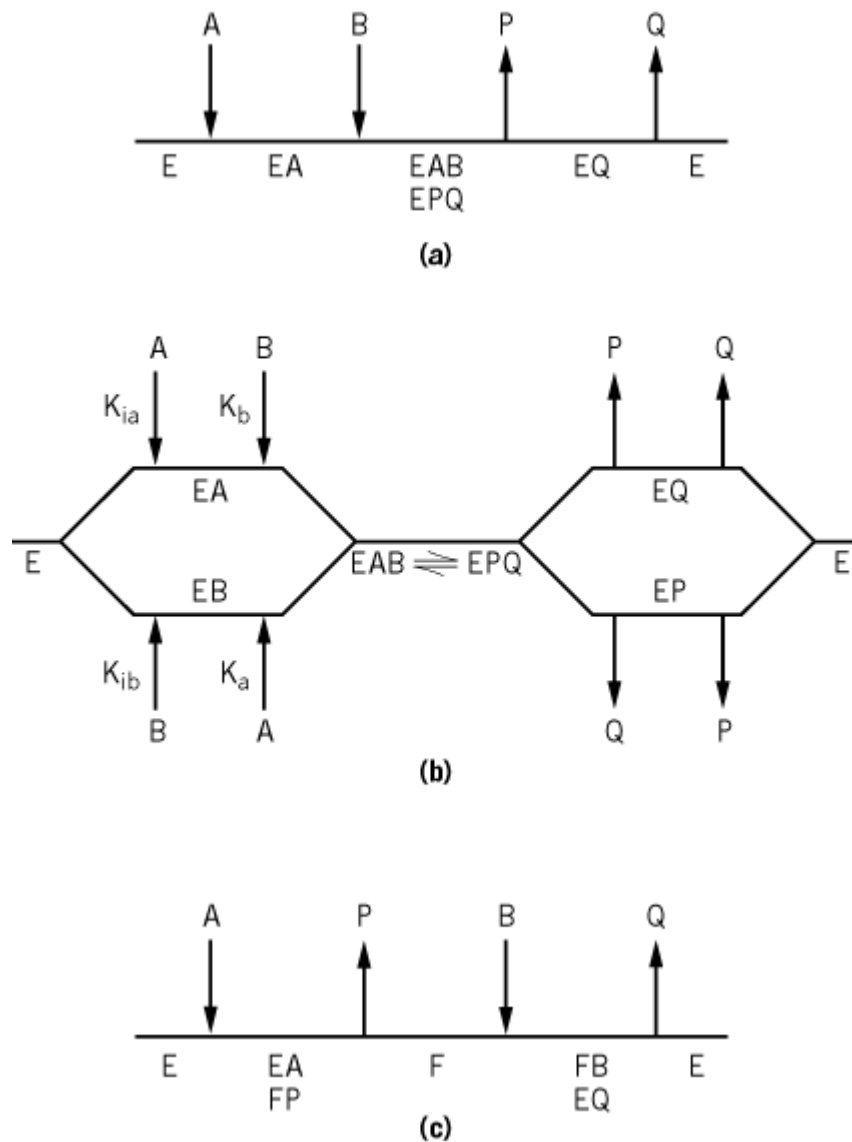
7. S. M. Block (1998) Kinesin: What gives? *Cell* **93**, 5–8.
8. G. S. Bloom and S. A. Endow (1995) Motor proteins 1: kinesins, *Protein Profiles* **2**, 1109–1111.
9. N. Hirokawa (1998) Kinesin and dynein superfamily proteins and the mechanism of organelle transport, *Science* **279**, 519–526.
10. R. B. Vallee and M. P. Sheetz (1996) Targeting of motor proteins, *Science* **271**, 1539–1544.

### Kinetic Mechanisms, Enzyme

The term kinetic mechanism used with [enzymes](#) refers to the order of substrate addition to, and product release from, the enzyme. The procedure that is now used widely for describing kinetic mechanisms was introduced by W. W. Cleland in 1963 ([1](#)). It involves the use of the letters A, B, and C for substrates; P, Q, and R for products; E and F for stable enzyme forms. For illustrative purposes, the enzyme is represented by a horizontal line, and the enzyme-reactant complexes are shown below the line. The reversible addition of substrates, and release of products, are denoted by single, vertical arrows. When required, rate constants for the forward and reverse steps of each reaction are written to the left and right of the vertical arrows, respectively.

The types of kinetic mechanisms that apply, and the shorthand procedure for describing them, can be well illustrated by reference to reactions involving two substrates and two products. These Bi–Bi reactions form a high proportion of all known enzymic reactions and include those catalyzed by [dehydrogenases](#) and phosphotransferases (see [Kinase](#)). The mechanisms fall into two broad categories: sequential and nonsequential (Fig. [1](#)).

**Figure 1.** Kinetic mechanisms for reactions involving two substrates (A and B) and two products (P and Q). (a) Ordered. (b) Random. (c) Ping-pong.



## 1. Sequential Mechanisms

Sequential mechanisms are characterized by the fact that both substrates must add to the enzyme before either product is released. This category has two basic subdivisions: ordered and random.

### 1.1. Ordered Mechanisms

In an ordered mechanism, there is a compulsory order of substrate addition and product release (Fig. 1a). The native enzyme possesses a binding site for substrate A, and formation of the EA complex results in the creation of a binding site for substrate B. This is most likely due to the binary enzyme–substrate complex undergoing a conformational change (see [Induced Fit](#)) or the highly synergistic binding of the two substrates, in which case the tightly bound one appears to bind first (2). After formation of the ternary EAB complex, the chemical reactions occur, with the formation of an EPQ complex. These two ternary, or central, complexes are treated as a single entity because steady-state kinetic studies, except for isotope effects, do not permit the detection of unimolecular isomerizations and cannot distinguish between them. Further, the inclusion of such isomerization steps does not affect the form of the steady-state rate equations. The mechanism is usually symmetrical, with the product of the first substrate to add being the last to dissociate from the enzyme.

Many dehydrogenases conform to this mechanism, with **NAD** or **NADP** being the first substrate to

add, and NADH or NADPH the last product to be released (3). When no assumptions are made about the relative rates of the various reaction steps, the mechanism is referred to as steady-state-ordered. The Theorell–Chance mechanism is a special case of an ordered mechanism where the steady-state concentrations of the ternary complexes, such as E-NAD-ethanol and E-NADPH-acetaldehyde, are kinetically insignificant (4). However, this is just a matter of degree, and the formation of these ternary complexes can be demonstrated by means of product inhibition studies, or better still by [isotope exchange at equilibrium](#).

Another special case in the ordered category is the *equilibrium ordered* mechanism, where the interaction between E and A occurs at thermodynamic equilibrium. This mechanism is observed when a metal ion is an essential activator and must interact to form an active enzyme with which the substrate then reacts (5). It can also occur when the off rate constant is considerably greater than the [turnover number](#).

## 1.2. Random Mechanisms

Enzymes conforming to a random kinetic mechanism possess distinct binding sites for each of the two substrate–product pairs. Thus, either substrate can add to the enzyme either before or after the other substrate (Fig. 1b). The same principle applies to product release. The presence of one substrate on the enzyme can enhance, hinder, or have no effect on the binding of the other. It is also possible for one substrate, and the product of the other substrate, to be present on the enzyme at the same time. The formation of such a dead-end complex will always occur with the smaller substrate–product pair, but not necessarily with the larger substrate–product pair if there is sufficient steric hindrance (see [Product Inhibition](#)). Several phosphotransferases ([kinases](#)) exhibit random kinetic mechanisms (6, 7).

## 2. Nonsequential Mechanisms

### 2.1. Ping-Pong Mechanisms

The representative of the nonsequential mechanisms is the single-site ping-pong mechanism (Fig. 1c). The first substrate to bind possesses a donor group that is transferred to the native enzyme E, to form a new stable form of enzyme F and the product P is released. The second substrate B reacts at the site that P has vacated, and the donor group on the enzyme is transferred to B to form Q. Aminotransferases have this mechanism, with the [amino group](#) of an [amino acid](#) being transferred to the cofactor [pyridoxal phosphate](#), to form a pyridoxamine-phosphate enzyme and subsequently to an  $\alpha$ -keto acid (8). Nucleoside diphosphokinase also has a ping-pong mechanism, with F being a phosphorylated enzyme (9). The F form of enzyme can be isolated and used for stoichiometric reaction with B. A characteristic feature of enzymes conforming to a ping-pong reaction mechanism is that they catalyze partial exchange reactions (see [Isotope exchange](#)) and have substrates of a similar structure. There are also enzymes, such as transcarboxylase (10), that catalyze two-site ping-pong mechanisms, and several that exhibit multisite ping-pong mechanisms (11). In these cases, an enzyme-bound carrier such as [biotin](#), lipoic acid, or 4-phosphopantetheine serves to carry the transferred moiety between the separate sites.

## 3. Initial Rate Equations for Sequential Mechanisms

### 3.1. Steady-State-Ordered and Rapid-Equilibrium Random Mechanisms

Under conditions where catalysis is not solely rate-limiting, initial velocity data for a random mechanism would yield nonlinear double-reciprocal plots (see [Lineweaver–Burk Plot](#)). But such nonlinearity cannot generally be detected with steady-state kinetic techniques, and random mechanisms behave as though they are of the rapid-equilibrium random type, with the interconversions of the ternary EAB and EPQ complexes being the slow steps of the reaction sequence in each direction. The initial velocity equation for rapid-equilibrium random and steady-state-ordered mechanisms has the same form (Eq. (1)):

$$v = \frac{VAB}{K_{ia}K_b + K_aB + K_bA + AB} \quad (1)$$

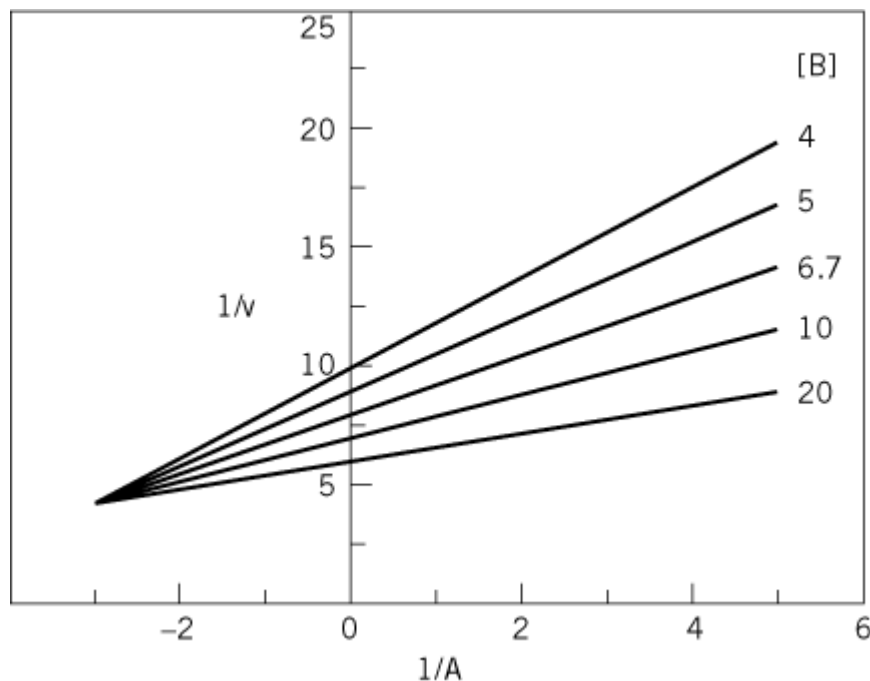
where  $V$  is the maximum velocity of the reaction;  $K_a$  and  $K_b$  are the  $K_m$  (Michaelis constants) for the substrates A and B, respectively;  $A$  and  $B$  are their respective concentrations; and  $K_{ia}$  is the **dissociation constant** for the complex EA. The numerator of Equation (1) represents the driving force of the reaction, and the denominator represents the distribution of the enzyme among its various forms. For the rapid-equilibrium random reaction, the four terms of the denominator denote the proportion of total enzyme that is present as free enzyme, EB, EA, and EAB, respectively. Because catalysis is slow relative to product release, no significant concentrations of EP and EQ will be present (Fig. 1b).

For a steady-state ordered mechanism, the  $K_{ia}K_b$  and  $K_aB$  terms represent the free enzyme,  $K_bA$  the EA complex, and AB the combined concentrations of the EAB, EPQ, and EQ complexes (Fig. 1a). The reciprocal form of Equation (1), with  $A$  as the variable substrate, is given in Equation (2):

$$\frac{1}{v} = \frac{K_a}{V} \left[ \frac{K_{ia}K_b}{K_aB} + 1 \right] \frac{1}{A} + \frac{1}{V} \left[ \frac{K_b}{B} + 1 \right] \quad (2)$$

This equation indicates that a double-reciprocal plot of  $1/v$  against  $1/A$  at different fixed concentrations of B would yield a family of straight lines that intersect to the left of the vertical ordinate, with the slopes and intercepts of the lines varying as a function of the concentration of B (Fig. 2). This is known as a primary plot. Equation (1) is symmetrical, so the double-reciprocal plot with B as the variable substrate will have the same general form as Equation (2). The coordinates of the crossover point would be  $-1/K_{ia}$  and  $(1/V) [1 - K_d/K_{ia}]$  with  $A$  as the variable substrate and  $-K_d/(K_{ia}K_b)$  and  $(1/V) [1 - K_d/K_{ia}]$  with  $B$  as the variable substrate.

**Figure 2.** Intersecting initial velocity pattern for a steady-state-ordered or rapid-equilibrium random mechanism obtained upon varying substrate A, at different fixed concentrations of substrate B. A similar pattern would be observed with B as the variable substrate. Values chosen for the kinetic parameters were  $V$ , 0.2;  $K_{ia}$ , 0.33;  $K_a$ , 0.05;  $K_b$ , 4.0, in arbitrary units.



For a rapid-equilibrium random mechanism,  $K_a/(K_{ia}K_b)$  would be equal to  $K/K_{ib}$ , as for this mechanism  $K_{ia}K_b = K_aK_{ib}$  (Fig. 1b). It should be noted that the value of  $K_a/K_{ia}$  determines whether the lines cross above, on, or below the abscissa and the position of the crossover point, in relation to the abscissa, is independent of which substrate is varied. The initial velocity patterns for the steady-state-ordered and rapid-equilibrium random mechanisms are of the symmetrical intersecting type. From replots of the slopes and intercepts of the primary plots as a function of the reciprocals of the concentrations of the changing fixed substrate concentration, it is possible to determine values for each of the kinetic parameters. However, it is preferable to obtain these values by performing an overall fit of the initial velocity data to Equation (1) (12). Initial velocity patterns do not allow a distinction to be made between a steady-state-ordered and rapid-equilibrium random mechanism. Further, if the mechanism is of the steady-state-ordered type, it is not possible to assign the  $K_{ia}$  value to a particular substrate.

### 3.2. Equilibrium Ordered Mechanism

An equilibrium ordered-mechanism is described by Equation (3):

$$v = \frac{VAB}{K_{ia}K_b + K_bA + AB} \quad (3)$$

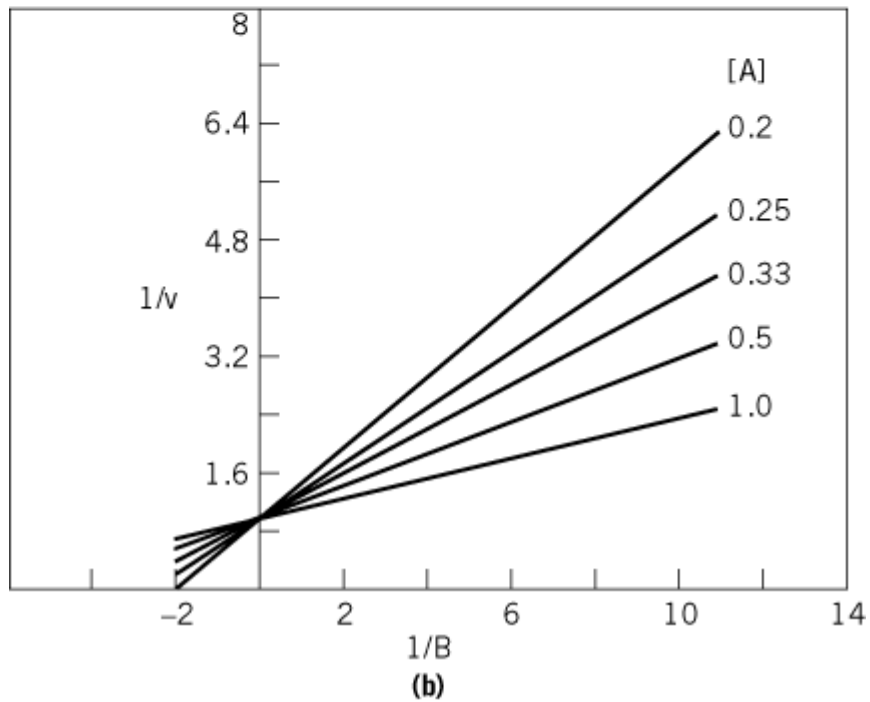
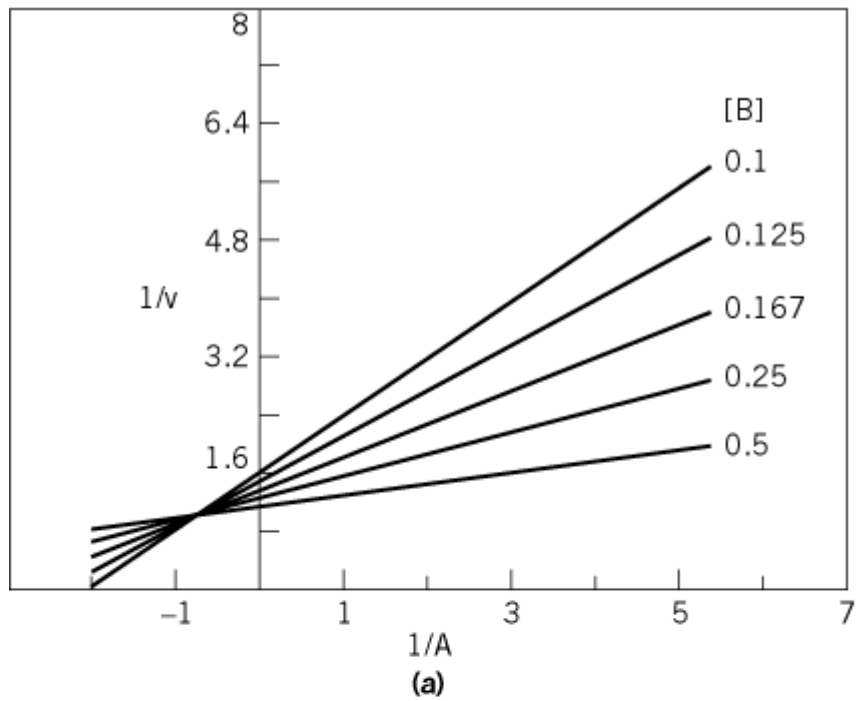
It differs from Equation (1) in not having a denominator term in  $B$ . That is, the value for  $K_a$  is essentially zero, and the free form of enzyme is represented by the single term,  $K_{ia}K_b$ . In contrast to Equation (1), Equation (3) is asymmetric. Rearrangement of this equation in double-reciprocal form yields Equations (4) and (5), with  $A$  and  $B$  as the variable substrates:

$$\frac{1}{v} = \frac{K_{ia}}{V} \left[ \frac{K_b}{B} \right] \frac{1}{A} + \frac{1}{V} \left[ \frac{K_b}{B} + 1 \right] \quad (4)$$

$$\frac{1}{v} = \frac{K_b}{V} \left[ \frac{K_{ia}}{A} + 1 \right] \frac{1}{B} + \frac{1}{V} \quad (5)$$

These equations indicate that although the initial velocity patterns are of the intersecting type, they are not symmetrical. The pattern with  $A$  as the variable substrate intersects to the left of the vertical ordinate, with coordinates of the crossover point being  $-1/K_{ia}$  and  $1/V$ . With  $B$  as the variable substrate, the lines intersect on the vertical ordinate (Fig. 3). The crossover point for this mechanism, with either  $A$  or  $B$  as the variable substrate, must lie above the abscissa. Further, a replot of the slopes of the primary plot with  $A$  as the variable substrate vs  $1/B$  yields a straight line that passes through the origin. A replot of the intercepts of the primary plot with  $A$  as the variable substrate vs  $1/B$  would give values for  $V$  and  $K_b$  (Eq. (4)), and a replot of the slopes of the primary plot with  $B$  as the variable substrate vs  $1/A$  would yield a value for  $K_{ia}$ .

**Figure 3.** Asymmetric intersecting initial velocity patterns for an equilibrium-ordered mechanism. (a) Varying substrate A, at different fixed concentrations of substrate B. (b) Varying substrate B, at different fixed concentrations of substrate A. The patterns were produced by using values in arbitrary units for  $V$ ,  $K_{ia}$ , and  $K_b$  of 1.0, 1.35, and 0.06, respectively.



### 3.3. Ping-Pong Mechanism

A reaction conforming to a single-site ping-pong mechanism is described by Equation (6):

$$v = \frac{VAB}{K_aB + K_bA + AB} \quad (6)$$

This is a symmetrical equation, without a constant term in the denominator, which contains only the Michaelis constants for the two substrates.  $K_aB$  and  $K_bA$  represent the E and F forms of enzyme, respectively, and  $AB$  represents all the binary enzyme-reactant complexes. Rearrangement of Equation (6) in double-reciprocal form with  $A$  and  $B$  as variable substrates gives Equations (7) and (8), respectively:

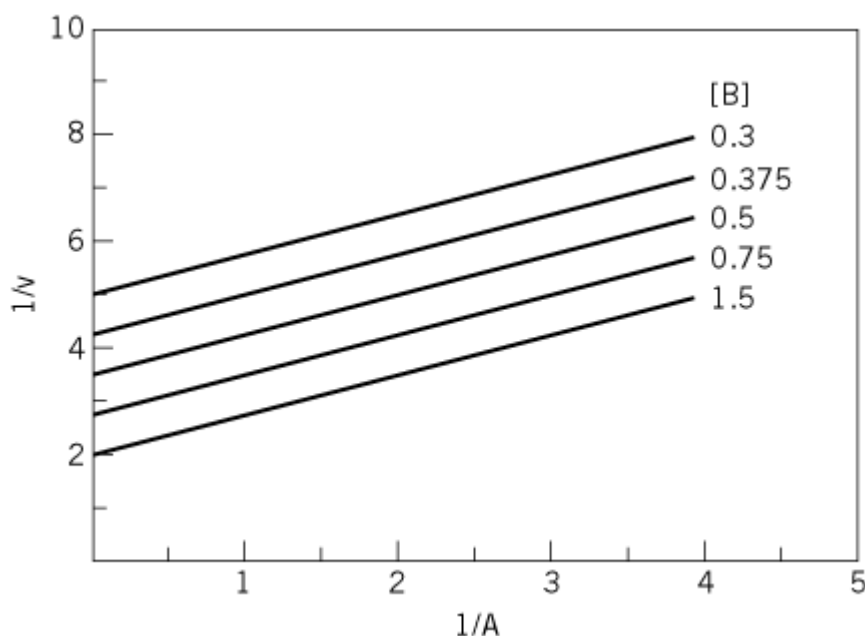


$$\frac{1}{v} = \frac{K_a}{V} \frac{1}{A} + \frac{1}{V} \left[ \frac{K_b}{B} + 1 \right] \quad (7)$$

$$\frac{1}{v} = \frac{K_b}{V} \frac{1}{B} + \frac{1}{V} \left[ \frac{K_a}{A} + 1 \right] \quad (8)$$

These equations indicate that the slopes of the lines do not vary with the concentration of the changing fixed substrate ( $B$  or  $A$ ), and thus, the initial velocity patterns consist of families of parallel straight lines (Fig. 4) The patterns are symmetrical. Values for  $V$ ,  $K_a$ , and  $K_b$  can be obtained from replots of the intercepts as a function of the reciprocals of the concentrations of the changing fixed substrate.

**Figure 4.** Parallel initial velocity pattern for a ping-pong mechanism as obtained by varying substrate A at different fixed concentrations of substrate B. The values in arbitrary units for  $V$ ,  $K_a$ , and  $K_b$  were 0.8, 0.6, and 0.9, respectively.



Some caution must be exercised when an initial velocity pattern seems to consist of parallel straight lines. It may be that the lines intersect below the abscissa and far to the left of the vertical ordinate, so that the pattern only appears to be parallel. The test is to study the reaction in both directions. A true ping-pong mechanism will give rise to parallel initial velocity patterns in each direction, whereas an apparent ping-pong mechanism will show a parallel pattern in one direction and an intersecting pattern in the other.

#### 4. Obtaining the Values of Parameters

The plotting of initial velocity data as described above is useful for determining the function to which the kinetic data should be fitted by nonlinear regression analysis (12), which should be used to obtain estimates of the values of kinetic parameters, and for illustrating the agreement between the experimental data and theoretical fit. Determination of values for kinetic parameters, together with estimates of their standard errors, should come from the fit of all the data to an assumed rate equation.

## Bibliography

1. W. W. Cleland (1963) *Biochim. Biophys. Acta* **67**, 104–137.
2. W. B. Knight and W. W. Cleland (1989) *Biochemistry* **28**, 5728–5734.
3. K. Dalziel (1975) *The Enzymes* **11**, 1–60.
4. H. Theorell and B. Chance (1951) *Acta Chem. Scand.* **5**, 1127–1144.
5. J. F. Morrison and K. E. Ebner (1971) *J. Biol. Chem.* **246**, 3977–3984.
6. J. F. Morrison and E. Heyde (1972) *Ann. Rev. Biochem.* **41**, 29–54.
7. P. A. Frey (1992) *The Enzymes* **20**, 141–186.
8. E. S. Snell and S. J. Di Mari (1970) *The Enzymes* **2**, 335–370.
9. E. Garces and W. W. Cleland (1969) *Biochemistry* **8**, 633–640.
10. D. Northrop (1969) *J. Biol. Chem.* **244**, 5808–5819.
11. W. W. Cleland (1986) *Invest. Rates Mech. Reactions* **6**, 791–870.
12. W. W. Cleland (1979) *Meth. Enzymol.* **63**, 103–138.

## Kinetics

The kinetics of a reaction deals with all of the factors that affect the rate of the reaction and with its underlying mechanism. **Thermodynamics** is concerned with how much of a substance S can be converted into a product P at equilibrium, whereas kinetics is concerned with how fast and by what pathway the conversion takes place. Understanding the difference is critically important in molecular biology. Viewed on a macroscale, the human body appears to be quite stable. On the microscale, however, the biochemicals that make up the body are in constant flux. Chemical reactions in biological systems are catalyzed by enzymes in microseconds to milliseconds. Moreover, biochemicals are continually being synthesized, converted to new forms, and degraded. Over time, a necessary language has developed to describe the time-dependence of a reaction and the factors that contribute to [catalysis](#). It is necessary to understand what is meant by the “order” of a reaction, its [half-life](#) and [relaxation time](#), its rate-determining step, [transition state](#), and [activation energy](#), [diffusion-controlled reactions](#), [free energy relationships](#), and the usefulness of [pulse-chase experiments](#).

On one level, kinetics describes the time-dependence of a reaction in mathematical terms. Chemical reactions can be described in terms of the molecularity and the “order” of the process. The molecularity is defined by the number of molecules participating in the reaction. Thus, the conversion of  $S \rightarrow P$  is a unimolecular reaction, whereas the reactions  $S + A \rightarrow P$  and  $S + A + B \rightarrow P$  are bimolecular and trimolecular reactions, respectively. The order of the reaction is defined by the power of the concentration-dependence of the rate of conversion of the starting substance S into products P. In general, the rate of a reaction is defined by the rate of disappearance of the reactants or the rate of appearance of the product:

$$\text{Rate} = \frac{-d[S]}{dt} = \frac{d[P]}{dt} \quad (1)$$

The observed rate depends on the concentration of each of the reactants (eg, [S]), raised to some power  $i$  depending on the number of molecules involved before the rate-determining step:

$$\text{Rate} = \frac{-d[\text{S}]}{dt} = k_{\text{app}}[\text{S}]^i \quad (2)$$

The proportionality constant is the apparent rate constant for the reaction. Reactions are variously said to be zero-, first-, or second- order when  $i$  is, respectively, 0, 1, or 2. Even greater values of  $i$  are possible. For a reaction with multiple reactants, the rate of the reaction depends on some order of each of them:

$$\text{rate} = \frac{-d[\text{S}]}{dt} = k_r[\text{S}]^i[\text{A}]^j[\text{B}]^k \quad (3)$$

The order of a reaction depends on which reactants are involved in the rate-determining step of a reaction, which need not be the same as the molecularity of the reaction.

### 1. First-Order Kinetics

If the rate of disappearance of S depends on the first power of the concentration of S, the reaction is said to be first order in [S].

$$\text{rate} = \frac{-d[\text{S}]}{dt} = k_{\text{uni}}[\text{S}] \quad (4)$$

The superscript one is usually not written but is implied by its absence. Examples of such a unimolecular reaction are the isomerization of a molecule or the conversion of one [radioisotope](#) into another, in neither case is another molecule necessarily involved. The apparent rate constant  $k_{\text{uni}}$  is first-order and has units of ( $\text{time}^{-1}$ ) (see Table 1), so that the rate has the required units of (concentration per time). The progress of such a reaction can be predicted by integrating Eq. (4). When the reaction is irreversible,

$$\frac{-d[\text{S}]}{[\text{S}]} = k_{\text{uni}} dt \quad (5)$$

$$\ln_e \frac{[\text{S}]_0}{[\text{S}]_0 - [\text{S}]} = 2.303 \log_{10} \frac{[\text{S}]_0}{[\text{S}]_0 - [\text{S}]} = k_{\text{uni}} t \quad (6)$$

where  $[\text{S}]_0$  is the starting concentration of S. The velocity of the reaction decreases with time, as the concentration of S diminishes.

**Table 1. Kinetic Order of a Reaction**

| Reaction Order | Reaction                              | Equation                          | Units of $k_a$ |
|----------------|---------------------------------------|-----------------------------------|----------------|
| Zero           | $\text{S} \xrightarrow{k_a} \text{P}$ | $-d[\text{S}]/dt = k_a$           | Conc $t^{-1}$  |
| First          | $\text{S} \xrightarrow{k_a} \text{P}$ | $-d[\text{S}]/dt = k_a[\text{S}]$ | $t^{-1}$       |

|        |  |  |               |
|--------|--|--|---------------|
| Second | $S + S \xrightarrow{k_3} P$                            | $-d[S]/dt = k_a[S]^2$                      | Conc $t^{-1}$ |
| Second | $A + S \xrightarrow{k_3} P$                            | $-d[S]/dt = k_a[A][S]$                     | Conc $t^{-1}$ |
| Mixed  | $A + S \xrightleftharpoons{K_c} C \xrightarrow{k_3} P$ | $-d[S]/dt = \frac{k_a[A][S]}{(K_c + [S])}$ | $t^{-1}$      |

---

Reactants are never completely converted to products because the product can reform the reactants. Thus an equilibrium is established between S and P:



When the reaction is significantly reversible so that only some of S is converted to P at equilibrium, both the forward and reverse rate constants,  $k_f$  and  $k_r$ , are involved in the kinetics:

$$-d[S]/dt = k_f[S] - k_r[P] \quad (8)$$

Because the concentration of the reactant(s) and product(s) at any time must always equal the initial reactant concentration  $S_0$ , the reaction follows first-order kinetics:

$$-d[S]/dt = (k_f + k_r)[S] + C \quad (9)$$

where  $C = -k_r[S_0]$ . The apparent first-order rate constant  $k_{app}$  for the reaction, is the sum of the forward and reverse rate constants:

$$k_{app} = k_f + k_r \quad (10)$$

At equilibrium the rate of the forward reaction is equal to the rate of the reverse reaction:

$$k_f[S] = k_r[P] \quad (11)$$

When rearranged, this expression yields

$$K_{eq} = \frac{k_f}{k_r} = \frac{[P]}{[S]} \quad (12)$$

By convention the concentration of product(s) is divided by the concentration of reactant(s). The constant  $K_{eq}$  is called the *equilibrium constant* for the reaction. The magnitude of the equilibrium constant is a function of the difference in the free energy  $DG$  of the reactant(s) and product(s). A large negative  $DG$  is reflected in a reaction that goes to completion.

## 2. Second-order kinetics

A second-order reaction can be either a reaction between two molecules of S or between S and another molecule A (Table 1). In these two cases, the rate equations are given, respectively, by

$$\text{Rate} = \frac{-d[S]}{dt} = k_{bi}[S]^2 \quad (13)$$

$$\text{Rate} = \frac{-d[\text{S}]}{dt} = \frac{-d[\text{A}]}{dt} = k_{\text{bi}}[\text{S}][\text{A}] \quad (14)$$

In the latter case, however, one reactant might be present in much greater concentration than the other, so that its concentration remains constant throughout the reaction. The disappearance of the limiting reactant, say S, is then pseudo-first-order:

$$\text{Rate} = \frac{-d[\text{S}]}{dt} = k_{\text{app}}[\text{S}] \quad (15)$$

Then the value of the pseudo-first-order rate constant is proportional to the concentration of A, and the second-order rate constant can be determined with this information:

$$k_{\text{app}} = k_{\text{bi}}[\text{A}] \quad (16)$$

### 3. Zero-order kinetics

If the rate of disappearance of S is independent of its concentration, the reaction is said to be zero-order in S. This nomenclature is used because any number raised to the zero power is one. Therefore, the term  $[\text{S}]^0$  would be absent:

$$\text{Rate} = \frac{-d[\text{S}]}{dt} = k_{\text{app}}[\text{S}]^0 = k_{\text{app}} \quad (17)$$

Such kinetics is observed with multireactant reactions if this particular reactant is not involved in the rate-determining step or possibly if it has been converted to another form before the rate-determining step, as in the case of the **Michaelis–Menten** complex of substrates with [enzymes](#).

### 4. Rate-determining step

With a number of reactants, any chemical reaction is likely to occur in several sequential steps, with just one or two of the reactants involved in each step. One of these steps is likely to be slower than the others, and its rate determines the rate of the over-all reaction. This is the rate-determining or rate-limiting step. The rate equation for the overall reaction is given by that for this step. For example, for the following reaction:



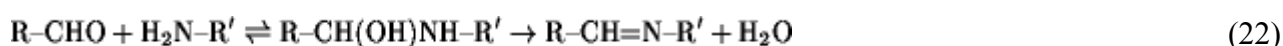
the second step might be rate-limiting, so that the rate equation is

$$\text{Rate} = \frac{d[\text{P}]}{dt} = k_{\text{app}}[\text{C}][\text{D}] \quad (20)$$

C is not, however, an initial reactant, so that the rate equation depends on how C was formed from A and B. For example, if C were in rapid equilibrium with A and B during the reaction, with equilibrium constant  $K_{\text{eq}}$ , the overall rate equation would be

$$\text{Rate} = \frac{d[\text{P}]}{dt} = k_{\text{app}} K_{\text{eq}} [\text{A}][\text{B}][\text{D}] \quad (21)$$

The identity of the rate-determining step can change, as the concentrations of the reactants change, which can lead to complex kinetics that are of mixed-order. For example (1), the conversion of an aldehyde R-CHO into an amine R-CH=N-R', occurs by the aldehyde reacting first with an amine H<sub>2</sub>N-R' to form a carbinolamine:



The first step is reversible, with equilibrium constant  $K_{\text{eq}}$ . The rate constant for the second step is  $k_{\text{uni}}$ . At low concentrations of amine and low pH, the rate-determining step is the formation of the carbinolamine. At high concentrations of amine, however, the carbinolamine is formed completely, and the rate of the reaction no longer depends on the concentration of the amine. The kinetics are zero-order with respect to the amine. A change in the rate-determining step has occurred, and the overall order is mixed (Table 1). At concentrations of amine less than 0.1 the value of  $K_{\text{eq}}$ , the reaction becomes pseudo-first-order:

$$-\frac{d[\text{R-CHO}]}{dt} = k_{\text{app}} [\text{R-CHO}] \quad (23)$$

$$k_{\text{app}} = \frac{k_{\text{uni}} [\text{H}_2\text{N-R}']}{K_{\text{eq}}} \quad (24)$$

This type of kinetics is also observed with enzymes because of the occurrence of the Michaelis-Menten complex.

## 5. Steady-state kinetics

With multi step reactions, such as



the concentrations of the intermediate species A, B and C often change relatively little after the reaction is initiated. This occurs generally when the intermediates are present at only low concentrations. They are generated from S at about the same rate as they are converted to P, and their concentrations decrease only as that of S diminishes. Such a situation is in a steady state. An example is the steady state of the level of water in a dam if the amount of water released from the dam equals the amount of water flowing into it during the same moment. Similarly, metabolites in cells are often at a steady state because their individual rates of formation are equal to their rates of breakdown.

The steady-state approximation is often very useful in analyzing the kinetics of complex reactions. Many reactions are studied in the steady state. The steady-state concentration of each species can be expressed in terms of the rate equations for each step in its formation and breakdown, where the rates of formation and breakdown are equal. Then the steady-state concentration can be used in the equations for the overall reaction. It is only an approximation, however, and a complete kinetic analysis requires measuring the individual rate constants by following the rapid approach to the steady state.

## 6. Transient state kinetics

Measuring the approach to a steady state generally requires rapid reaction techniques, because individual reactions can occur on very short time scales. A stopped-flow instrument can rapidly mix reactants within about 1 ms, so kinetic reactions with rate constants approaching  $10^3 \text{ sec}^{-1}$  can be monitored. Rapid reactions can often be slowed markedly by decreasing temperature (2). [Relaxation spectrometry](#) can be used to study reactions in which a sufficient equilibrium occurs.

### Bibliography

1. D. S. Auld and T. C. Bruice (1967) *J. Amer. Chem. Soc.* **89**, 2083–2089.
2. D. S. Auld (1993) *Methods Enzymol.* **226**, 553–564.

### Suggestions for Further Reading

3. A. A. Frost and R. G. Pearson (1962) *Kinetics and Mechanism*, Wiley, New York.
4. T. C. Bruice and S. J. Benkovic (1966) *Bioorganic Mechanisms*, Vols. **I** and **II**, Benjamin, New York.
5. W. P. Jencks (1969) *Catalysis in Chemistry and Enzymology*, McGraw–Hill, New York.
6. D. Piszkiwicz (1977) *Kinetics of Chemical and Enzyme-Catalyzed Reactions*, Oxford University Press, New York.
7. K. J. Laidler and P. S. Bunting (1973) *The Chemical Kinetics of Enzyme Action* 2nd ed., Oxford University Press, London.
8. A. Cornish–Bowden (1979) *Fundamentals of Enzyme Kinetics* Butterworths, London, Boston, pp. 1–15.

## Kinetochores

The kinetochore is the point where the [centromere](#) of a [chromosome](#) attaches to the **mitotic spindle**. This platelike structure is essential for chromosomal movement during **mitosis**. The mitotic spindle consists of numerous spindle fibers composed of [tubulin](#) and associated proteins. Tubulin exists as two closely related proteins,  $\alpha$ - and  $\beta$ -tubulin, that form dimers of 100,000 Da. The centromeric spindle fibers connect the kinetochore to the **centrosome**. The organization of tubulin in [microtubules](#) is an active function of the kinetochores and has an important role in directing chromosomal movement during **anaphase** (1). Chromosomal movement toward the poles of the dividing cell depends on removing tubulin from microtubules at the kinetochoric ends. This shortens the centromeric spindle fibers and moves them toward the centrosomes at the poles of the cell (2). In **yeast**, each kinetochore is associated with a single spindle fiber, whereas in mammalian cells the kinetochore is much larger and is attached to many fibers.

Among the mammalian centromeric proteins (CENPs) that have been biochemically defined are CENP-A, CENP-B, CENP-C, and CENP-E (see [Centromeres](#)). These proteins are found on the surface of constitutive [heterochromatin](#) and comprise components of this specialized chromatin structure and of the kinetochore itself (3). For example, CENP-C is concentrated in a narrow band immediately below the inner kinetochoric plate at the interface between heterochromatin and the kinetochore. CENP-C is required for normal kinetochoric assembly and is found only at the active centromere of a stable [dicentric chromosome](#), suggesting a direct role in centromeric function (4).

Components of the kinetochore include CENP-E, which resembles microtubule-binding motor proteins and functions as a kinetochore motor during the early part of mitosis (5).

### Bibliography

1. T. J. Mitchison (1988) *Ann. Rev. Cell Biol.* **4**, 527–545.
2. B. R. Brindley (1985) *Ann. Rev. Cell Biol.* **1**, 145–185.
3. A. F. Pluta et al. (1995) *Science* **270**, 1591–1594.
4. J. Tomkiel et al. (1994) *J. Cell Biol.* **125**, 531–545.
5. K. D. Brown, K. W. Wood, and D. W. Cleveland (1996) *J. Cell Sci.* **109**, 961–969.

### Suggestion for Further Reading

6. A. Wolffe (1998) *Chromatin: Structure and Function*, 3rd ed., Academic Press, London.

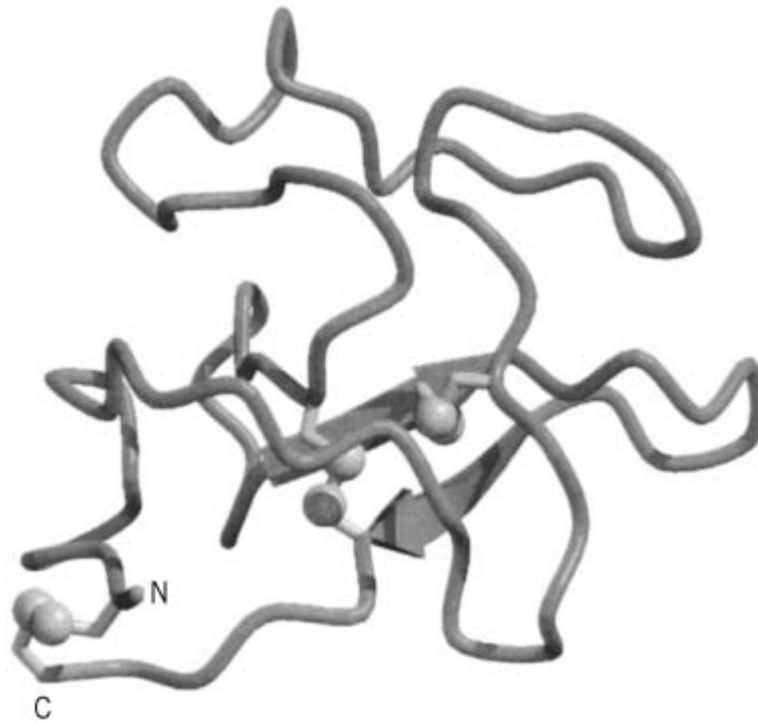
## Kringle Domain

The Kringle domain is a unit of both [protein structure](#) and function. It is usually found in [proteinases](#) associated with [blood clotting](#) and fibrinolysis such as [thrombin](#), [plasminogen](#), and plasminogen activators. The Kringle domain is characterized as an *N*-terminal module of ~80 residues having six conserved [cysteine](#) residues (called CI through to CVI) that form three internal [disulfide bonds](#), with the connectivity I–VI, II–IV and III–V. The Kringle domain is usually present as 1 to 5 repeating units in the [polypeptide chain](#) (although apolipoprotein (a) has 38 Kringle domains), in combination with repeating units of other protein modules, such as the [EGF motif](#). The function of the Kringle domain appears to be that of molecular recognition, especially that of protein fragments having an *N*-terminal [lysine](#) residue.

The amino acid sequence of the Kringle domain is highly conserved and usually can be recognized from the [primary structure](#) of a protein. The [tertiary structure](#) of Kringles is also highly conserved; a two-stranded antiparallel [b-sheet](#) at the core of the domain is formed through close packing of two of the three disulfide bonds (Fig. 1). The Kringle structure also includes many [b-turns](#) and several additional regions of short (2-residue) b-sheet. Recognition of a *C*-terminal lysine residue of another protein by Kringle domains occurs through three interactions: (1) two conserved **aspartate** residues interact with *N<sub>z</sub>* of the lysine side chain; (2) two conserved [tryptophan](#) residues form a **hydrophobic** groove that interacts with the methylene groups of the lysine side chain; and (3) an [arginine](#) or lysine residue interacts with the negatively charged *C*-terminal carboxylate of the lysine.

**Figure 1.** Schematic representation of the backbone structure of a Kringle domain (taken from the coordinates of the human plasminogen Kringle 4 domain, with Protein Data Bank accession code 1KRN). b-Strands are shown as arrows, and the three disulfide bonds are shown in ball-and-stick representation, with the sulfur atoms depicted as spheres. The *N*- and *C*-termini are labeled. This figure was generated using Molscript (1) and Raster3D (2, 3).





#### Bibliography

1. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
2. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
3. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

4. L. Patthy et al. (1984) Kringles: modules specialized for protein binding. *FEBS Lett.* **171**, 131–136 (Early description of the Kringle domain based on sequence.)
5. I. I. Matthews, P. Vanderhoff-Hanaver, F. J. Castellino, and A. Tulinsky (1996) Crystal structures of the recombinant Kringle 1 domain of human plasminogen in complexes with the ligands *e*-aminocaproic acid and *trans*-4-(aminomethyl)cyclohexane-1- carboxylic acid. *Biochemistry* **35**, 2567–2576.

#### LA (Lon) Proteinase

The first [proteinase](#) found to be ATP-dependent was proteinase La (lon). It is required in *Escherichia coli* for the **protein degradation** of most abnormal proteins and certain regulatory proteins. Some bacteria contain multiple forms of this proteinase which serve distinct physiological functions. Cells lacking this enzyme have reduced capacity to degrade abnormal proteins and have phenotypic effects (UV sensitivity, mucoidy) due to failure to degrade cell septation inhibitors or regulators of capsular polysaccharide synthesis. It is a homopolymer of four identical subunits, each containing both [ATPase](#) and proteolytic [active sites](#). The proteinase hydrolyzes unfolded proteins in an ATP-dependent reaction in which, on average, two ATP molecules are consumed for every [peptide bond](#)

cleaved. This proteinase is also a protein-activated ATPase, and the binding of a polypeptide substrate stimulates ATP hydrolysis and activates the proteolytic mechanism. In fact, when it is isolated, the enzyme is latent, but ATP binding allows the active sites to form: ATP hydrolysis then temporarily inactivates the enzyme until additional ATP is bound. Thus, the enzyme appears to go through a complex activation-inactivation cycle that is triggered by substrates and prevents inappropriate or nonspecific proteolysis *in vivo*.

### Suggestions for Further Reading

A. L. Goldberg, R. P. Moerschell, C. H. Chung, and M. R. Maurizi (1994) ATP-dependent protease La (Ion) from *Escherichia coli*. *Methods Enzymol.* **244**, 350–375.

S. Gottesman, M. R. Maurizi, and S. Wickner (1997) Regulatory subunits of energy-dependent proteases. *Cell* **91**, 435–438.

## Lac Operon

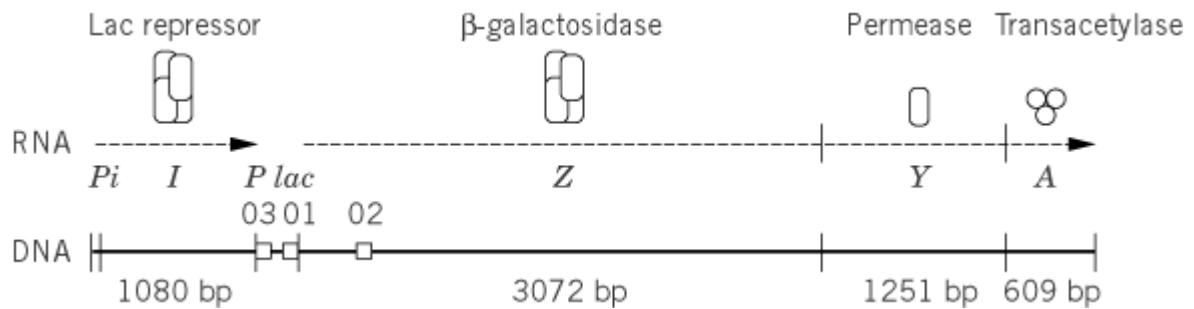
Lactose turns on the synthesis of **b-galactosidase** in *Escherichia coli* a thousandfold. How does it do this? What are the parts of the induction machinery? The first correct answers to these questions were given in 1960 by François Jacob and Jacques Monod (1). A specific [repressor](#), the product of the *lac I* **gene** (see [Lac Repressor](#)), interacts with the *lac* **operator** and thereby prohibits the production of *lac* [messenger RNA](#) (mRNA). The **inducer**, lactose, somehow specifically inactivates the Lac repressor. One year later, Jacob and Monod presented more evidence for this model of negative regulation in the *lac* and the Lambda *phage* systems (2, 3). They proposed that all systems of *E. coli*, and even **eukaryotes** (4), function in this way. Positive control of **gene expression** was generally ruled out. In retrospect, we see that this was wrong, in particular considering the situation in eukaryotes. The *lac* operon looked so perfect that a function for its two additional operator-like sequences was not considered for almost thirty years. Only recently was it demonstrated that both play a substantial role in repression (5-7).

When the [operon](#) model was proposed for the *lac* system in 1961, the experimental evidence supporting this model was overwhelming (2, 3). First, it was shown unambiguously that gratuitous inducers, like **IPTG**, which are not hydrolyzed by b-galactosidase, induce *de novo* synthesis of b-galactosidase. This rules out the instruction hypothesis, which claimed that a hypothetical pre-galactosidase folds irreversibly around the inducer, like **antibodies** were then thought to get their specificity by folding around **antigens**. Second, the **dominance** of wild type  $I^+$  over **constitutive**  $I^-$  Lac repressor [mutants](#) was easily explained by the operon model. The presence of Lac repressor, of course, is dominant over its absence (see [Lac Repressor](#)). The existence of the repressor predicted the existence of a genetic entity, called the *lac* operator (*lac O*), a stretch of DNA or RNA with which the repressor interacts. If the operator is destroyed by mutation, it should not bind repressor any longer. Repression should decrease drastically. But the effect should be seen in only *cis* with the particular genes that are directly linked with the operator (see [Cis-Acting](#)). This could be tested experimentally. Mutants were indeed found that had this quality.

This test was possible because the *lac* operon is **polycistronic**. It consists of the *Z* gene that codes for b-galactosidase, the closely-linked *Y* gene that codes for Lac **permease**, a [membrane protein](#) necessary for the uptake of lactose, and the *A* gene that codes for a transacetylase (Fig. 1). The existence of episomes allows the crucial dominance test: in a construct  $I^+O^+Z^+Y^-/F'I^+O^CZ^-Y^+$ , where  $O^+$  or  $O^C$  signify wild-type operator or operator carrying a constitutive mutation, synthesis of

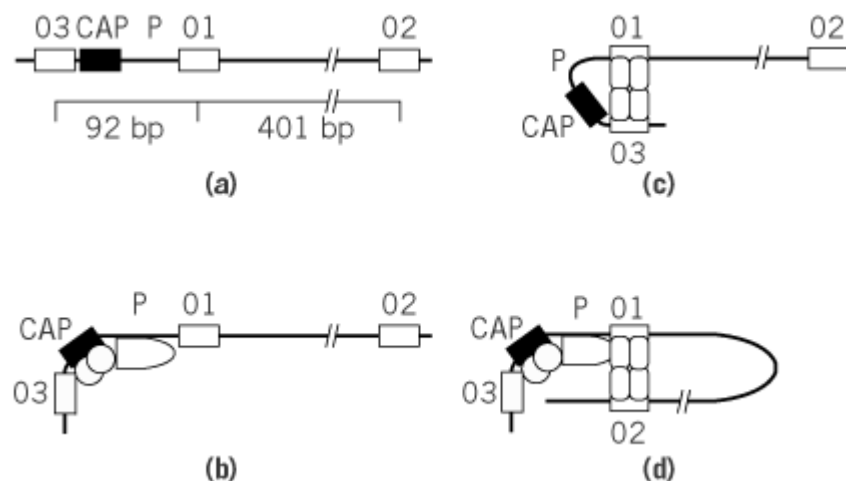
Lac permease should be constitutive, whereas synthesis of  $\beta$ -galactosidase should remain repressed. Indeed, only transcription of the mRNA in *cis* of the operator mutation is turned on (1-3). Later, it was demonstrated by Walter Gilbert that the central base pair of the **palindromic** *lac* operator is positioned 11 bp downstream of the start of [transcription](#) (8). Thus, Lac repressor most likely acts by steric hindrance of **RNA polymerase**.

**Figure 1.** The *lac* operon. At the bottom is the gene structure with the open reading frames coding for the I, Z, Y and A gene products, plus the promoters *P<sub>i</sub>* and *P<sub>lac</sub>* and the operators *O<sub>1</sub>*, *O<sub>2</sub>*, and *O<sub>3</sub>*. Two RNA products are produced by transcription, as shown. The second is polycistronic and contains the Z, Y, and A open reading frames. The four proteins that result from the translation of the messenger RNA are shown at the top.



The operon model for the *lac* system gained immensely in credibility because evidence for similar regulation of the totally unrelated Lambda system was provided (2, 3) (see **Lambda repressor**). The data looked so convincing that the authors and most readers thought that all systems in *E. coli*, and even in eukaryotes, are regulated like *lac* or Lambda (4). It took years until it became clear that this is not so. Even the *lac* promoter itself needs the presence of the **catabolite activator protein (CAP)** and of **cyclic AMP** to be fully active. The CAP binding site of the *lac* operon is situated 61.5 base pairs upstream of the transcription start site (Fig. 2). Transcription of the *lac* operon drops down to 2% in the absence of CAP protein bound there.

**Figure 2.** The *P<sub>lac</sub>* promoter region. (a) The structure of the promoter region of the *lac* operon (see Figure 1), indicating the promoter site (*P*), the site for binding catabolite activator protein (*CAP*), and the three operator regions *O<sub>1</sub>*, *O<sub>2</sub>*, and *O<sub>3</sub>*. (b) Unrepressed transcription of the operon, showing RNA polymerase binding to the promoter site and the dimeric CAP protein with bound cyclic AMP binding to its site. (c) Repression by the tetrameric Lac repressor binding simultaneously to *O<sub>1</sub>* and *O<sub>3</sub>*; the *CAP* and *P* sites are shown unoccupied, although there is no evidence for or against this. (d) Repression by the Lac repressor that binds simultaneously to *O<sub>1</sub>* and *O<sub>2</sub>*. It is likely that both the CAP protein and RNA polymerase are bound to their sites, although this has not been demonstrated explicitly.



The model of the *lac* operon and its control is present in all textbooks of genetics, molecular biology, and biochemistry. Yet if a *lac* operon were constructed as it is in the textbooks, repression would be only 1 to 2% efficient. Missing in all textbooks are the two *lac* operator-like structures, called pseudo-operators *O2* and *O3* by their discoverers (Figure 2). *O2* is located 401 bp downstream of *O1*, and *O3* is 92 bp upstream. In the absence of both *O2* and *O3*, repression drops 70-fold! But if only one of them is destroyed, repression drops only two- to threefold (5). How can one explain this behavior? It was shown before that wild-type Lac repressor forms tetramers and that two subunits bind to one operator. Tetrameric Lac repressor thus can bind to two operators if they are on different pieces of DNA or if they are on the same piece at an appropriate distance. If two operators are presented on a single linear piece of DNA, binding is optimal at distances of integral turns of the DNA double helix (assuming a turn every 10.5 bp) (9). Such loops are not obtained with dimeric Lac repressor that has its C-terminal residues deleted and thus is unable to provide the four-helical bundle that forms the tetramer. Such dimeric Lac repressor represses only weakly compared to tetrameric wild-type repressor (5). Thus, it was concluded that tetrameric Lac repressor forms stable loops between either *O1* and *O2* or between *O1* and *O3*. Thus, *O2* and *O3* are thus auxiliary operators (6, 7). The evolution of such "redundant" structures is a mystery that asks for an answer.

Important concepts and also widely used methods and techniques, have been developed out of the *lac* operon: see **Alpha complementation**, **Fusion gene**, **protein**, **Reporter gene**, **X-gal**, (chemical) **DNA sequencing**, **Gel retardation assay**, **Footprinting nucleic acids**. The main articles by Monod have been reprinted (10). One book reports on the *lac* operon and another on operons in general (11, 12). Finally, the history and present state of research on the paradigmatic *lac* operon have been written recently (13).

#### Bibliography

1. F. Jacob, D. Perrin, C. Sanchez, and J. Monod (1960) *Comptes Rendues* **250**, 1727–1729.
2. F. Jacob and J. Monod (1961) *J. Mol. Biol.* **3**, 318–351.
3. F. Jacob and J. Monod (1961) *Cold Spring Harbor Symp. Quant. Biol.* **26**, 193–211.
4. J. Monod and F. Jacob (1961) *Cold Spring Harbor Symp. Quant. Biol.* **26**, 389–401.
5. S. Oehler, E. R. Eismann, H. Krämer, and B. Müller-Hill (1990) *EMBO J.* **9**, 973–979.
6. S. Oehler, M. Amouyal, P. Kolkhof, B.v. Wilcken-Bergmann, and B. Müller-Hill (1994) *EMBO J.* **13**, 3348–3355.
7. J. Müller, S. Oehler, and B. Müller-Hill (1996) *J. Mol. Biol.* **257**, 21–29.
8. W. Gilbert, J. Gralla, J. Majors, and A. Maxam (1975) In *Protein-Ligand Interactions*, (H. Sund and G. Blauer eds.) de Gruyter, Berlin, New York, pp. 193–206.
9. H. Krämer, M. Niemöller, M. Amouyal, B. Revet, B.v. Wilcken-Bergmann, and B. Müller-Hill (1987) *EMBO J.* **6**, 1481–1491.
10. A. Lwoff and A. Ullmann, eds. (1978) *Selected Papers in Molecular Biology by Jacques Monod*, Academic Press, New York.
11. J. R. Beckwith and D. Zipser, eds. (1970) *The Lactose Operon*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
12. J. H. Miller and W. S. Reznikoff, eds. (1978) *The Operon*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
13. B. Müller-Hill (1996) *The lac Operon. A Short History of a Genetic Paradigm*, de Gruyter, Berlin.

## Lac Repressor

*Escherichia coli* grown in the presence of glucose produces about four molecules of tetrameric **b-galactosidase** per cell. When grown on lactose, the number of b-galactosidase molecules increases about a thousandfold. In 1959 Pardee, Jacob, and Monod proposed that the product of the *lac* I **gene**, Lac repressor, represses the production of b-galactosidase in the absence of **inducer**. Lactose acts as an inducer. It inactivates Lac repressor transiently by binding to it. In fact b-galactosidase isomerizes lactose into allolactose (1-6-b-galactosido-glucose), which then serves as the inducer. Lac repressor may also be inactivated permanently by a mutation in the *lac* I gene (1) The repressor hypothesis was plausible, but it took seven years until Lac repressor was isolated (2). It was the first [transcription factor](#) ever isolated. Since then, Lac repressor has been analyzed extensively, so extensively that it is now one of the best-analyzed proteins.

In 1959 when the properties of Lac repressor were defined (1), it was not clear what kind of molecule it is. Is it RNA or protein? It was also unclear *how* it controls the synthesis of b-galactosidase. In 1960 Jacob and Monod proposed that Lac repressor interacts with a particular region of the DNA, which they called the **operator** (2). Mutations in this region by definition are *cis-dominant* constitutive. Mutations in just one of the genes controlled by *lac* operator (the *lac* Z gene coding for b-galactosidase or the *lac* Y gene coding for Lac **permease**) and the existence of an F' **episome** that carries the entire *lac* region made it possible to provide experimental evidence for this proposition. The authors called their model the [operon](#) model (see [Lac Operon](#)). In 1961 Jacob and Monod generalized their model with additional evidence for the *lac* system and new evidence for the *Lambda* repressor system as the general model for controlling enzyme synthesis in *E. coli* (3, 4). At that time, they were convinced that repressors are RNA, in spite of the fact that common sense indicated that they are proteins. How would inducers bind specifically to RNA? The RNA hypothesis was disproved for Lac repressor only in 1965, when two independent researchers found nonsense mutations in the *lac* I gene.

Lac repressor ceased to be a hypothetical molecule in 1966 when Gilbert and Müller-Hill announced its isolation (5). They used [equilibrium dialysis](#) with <sup>14</sup>C-labeled inducer **IPTG** as a test. They had assumed that ten monomeric repressors per cell are present and had interpreted that *in vivo* induction indicates an inhibition constant  $K_i$  of  $6 \times 10^{-6}$  M. These numbers implied that the test would never work with extracts containing wild-type amounts of Lac repressor. Thus, they set out to isolate a mutant that binds IPTG more tightly than wild-type repressor and found a mutant that binds IPTG twice as tightly. This mutation sufficed to make detecting Lac repressor possible. In retrospect is it evident that the assumptions mentioned previously were partially wrong. One year later, specific binding of Lac repressor to *lac* operator was demonstrated (6). Thus, the operon model was fully confirmed.

In 1968 a promoter mutant of the *lac* I gene was isolated and placed on potentially replicating **bacteriophage**  $\phi$ 80 DNA. This construct produced 300 times more Lac repressor than wild-type (7). This opened the way to extensive biochemical and biophysical analysis. Thus, Lac repressor was sequenced in 1973 (8). For five years its protein sequence remained the first and only one of a repressor (transcription factor). Some minor mistakes in the protein sequence were found when the DNA sequence of the *lac* I gene was determined (9). Now we know that one subunit of Lac repressor consists of 360 residues.

The modular structure of Lac repressor was recognized in 1972. All constitutive mutations that map in the region of the I gene coding for the N-terminus of Lac repressor are negative dominant, I<sup>d</sup>. This can be explained by the facts that DNA binding is located here and that all four subunits of

tetrameric Lac repressor are necessary for DNA binding (10). Otherwise loop formation with two operators is impossible (see [Lac Operon](#)). Trypsin cuts specifically at Arg 51 and Lys 59 of intact Lac repressor. The N-terminal peptides 1 to 51 and 1 to 59, called the headpieces, bind specifically to *lac* operator (11). They do so by using a **helix-turn-helix** motif. The molecular structure of the complex formed by the headpiece and *lac* operator has been worked out in detail by [NMR](#) analysis (12). Moreover, all possible variants of this protein-DNA complex have been studied by genetic analysis. Substitutions of residues 17, 18, 21, and 22 (residues 1, 2, 5, and 6 of the recognition **a-helix**) change the specific base contacts, and these residues have been replaced by all possible amino acids. Residues 1, 2, and 6 of the recognition helix essentially work additively. In contrast, residue 5 of the recognition helix cooperates with residues 1 and 2 of the recognition helix (13).

The core of Lac repressor, residues 60 to 329, forms the inducer-binding, dimeric part of Lac repressor. Here mutations are found that destroy inducer binding, that is, lead to the noninducible, dominant I<sup>s</sup> **phenotype**, or that inhibit dimer formation, that is, lead to a recessive, constitutive phenotype, I<sup>-</sup>, which still binds inducer IPTG. A most extensive genetic analysis of the *lac* I gene has been done by Jeffrey Miller. In the end he succeeded in replacing each **codon** from codon 2 to codon 329 with an **amber** codon. Each individual amber construct was tested with twelve different amber **suppressors** that insert twelve different amino acids. These mutants have been analyzed with respect to their three-dimensional structure (see later) (14).

The C-terminal tails of the four subunits of one Lac repressor molecule, residues 330 to 360, form a four-helix bundle composed of four [coiled coils](#), each of which carries two leucine [heptad repeats](#). This short tail allows Lac repressor to become a tetramer. It needs to be a tetramer to use auxiliary operators effectively by forming DNA-protein-DNA loops. If the tail is cut off, Lac repressor dimers result. Their capacity to repress is strongly reduced because they no longer can form loops with the auxiliary operators (see [Lac Operon](#)).

Lac repressor is not unique. It is a member of a family of almost a dozen similar, mostly dimeric repressors in *E. coli*. Examples with **homologous** sequences are two Gal repressors, the Cyt repressor, the Pur repressor, the Ebg repressor, the Raf repressor, the Fru repressor, and the Asc repressor. But this is not all. A half dozen monomeric **periplasmic** proteins exist that are similar in sequence and structure to the core of Lac repressor and that bind small molecules like sugars or amino acids (14).

Lac repressor was available in gram amounts in the early seventies. Yet [X-ray crystallography](#) structures of the core (15), the Lac repressor/IPTG complex, and the Lac repressor/operator complex (16) were solved only recently. These structures, the structures of related proteins, and the extensive mutant analysis, open up new ways to understand the structure and function of all possible variants of Lac repressor in depth. Finally, most recently a book appeared that describes the history and present knowledge of the paradigmatic *lac* operon (17). Lac repressor has an essential part in it.

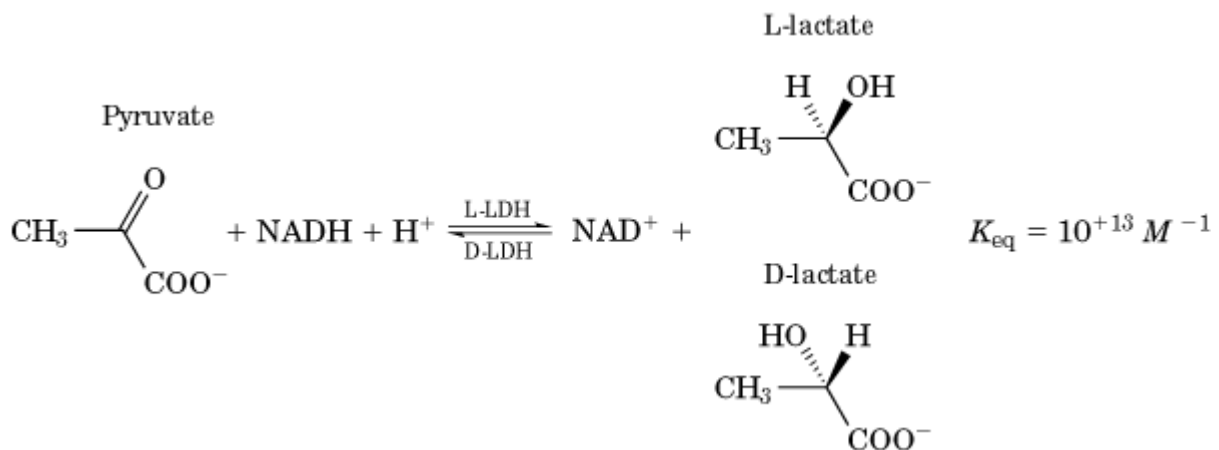
## Bibliography

1. A. B. Pardee, F. Jacob, and J. Monod (1959) *J. Mol. Biol.* **1**, 165–178.
2. F. Jacob, D. Perrin, C. Sanchez, and J. Monod (1960) *Comptes Rendues* **250**, 1727–1729.
3. F. Jacob and J. Monod (1961) *J. Mol. Biol.* **3**, 318–356.
4. F. Jacob and J. Monod (1961) *Cold Spring Harbor Symp. Quant. Biol.* **26**, 193–211.
5. W. Gilbert and B. Müller-Hill (1966) *Proc. Natl. Acad. Sci. USA* **56**, 1891–1898.
6. W. Gilbert and B. Müller-Hill (1967) *Proc. Natl. Acad. Sci. USA* **58**, 2415–2421.
7. B. Müller-Hill, L. Crapo, and W. Gilbert (1968) *Proc. Natl. Acad. Sci. USA* **59**, 1259–1264.
8. K. Beyreuther, K. Adler, N. Geisler, and A. Klemm (1973) *Proc. Natl. Acad. Sci. USA* **70**, 3576–3580.
9. P. J. Farabough (1978) *Nature* **274**, 765–769.

10. K. Adler, K. Beyreuther, E. Fanning, N. Geisler, B. Gronenborn, A. Klemm, B. Müller-Hill, M. Pfahl, and A. Schmitz (1972) *Nature* **237**, 322–327.
11. R. Ogata and W. Gilbert (1978) *Proc. Natl. Acad. Sci. USA* **75**, 5851–5854.
12. M. Slijper, A. M. J. J. Bonvin, R. Boelens, and R. Kaptein (1996), *J. Mol. Biol.* **259**, 761–773.
13. J. Lehming, J. Sartorius, B. Kisters-Woike, B. v. Wilcken-Bergmann, and B. Müller-Hill (1990) *EMBO J.* **9**, 615–621.
14. J. Suckow, P. Markiewicz, L. G. Kleina, J. Miller, B. Kisters-Woike, and B. Müller-Hill (1996) *J. Mol. Biol.* **261**, 509–523.
15. A. M. Friedman, T. O. Fischmann, and T. A. Steitz (1995) *Science* **268**, 1721–1727.
16. M. Lewis, G. Chang, N. C. Horton, M. A. Kercher, H. C. Pace, M. A. Schumacher, R. G. Brennan, and P. Lu (1996) *Science* **271**, 1247–1254.
17. B. Müller-Hill (1996) *The lac Operon. A Short History of a Genetic Paradigm*, de Gruyter, Berlin.

## Lactate Dehydrogenase (LDH)

[Enzymes](#) known as lactate dehydrogenases catalyze the interconversion of D- or L-lactate (only one of the two **isomers** by any single enzyme) and pyruvate, while donating the elements of hydrogen to an acceptor (shown as  $NAD^+$  in Eq. (1)):



Other hydrogen acceptors are:

1. Molecular oxygen (to give water) when the enzymes are known as lactate oxidases
2. Metals in the cell [membrane](#)
3. FAD linked to a [cytochrome](#)

The first two classes are least well known. The structure of the **yeast** flavocytochrome enzyme of the third class has been determined by [X-ray crystallography](#); its mechanism has been studied by [site-directed mutagenesis](#) and is well understood, but it and mandelate dehydrogenase represent an uncommon class.

The  $NAD^+$ -dependent enzymes are those usually meant by the unqualified term LDH. The biologically common L-lactate is usually assumed to be the substrate, but there are now over 22

examples of sequences of the D-LDH family from bacteria, plants, yeasts, and the arachnids, and it is increasingly necessary to specify the lactate isomer. This entry deals with NAD<sup>+</sup>-dependent L- and D-LDH.

## 1. L-LDH: Species Distribution, Functions, Genes, Proteins, and Families

L-LDH enzymes are found in the **archaea**, **prokaryia**, and **eukaryia** (Table 1). Most are relatively specific for the thermodynamically favorable reduction of pyruvate and NADH to L-lactate and NAD<sup>+</sup>. Related enzymes of the general family of (*S*)-2-hydroxyacid (*R*-CH<sub>2</sub>-CHOH-COOH) dehydrogenases are the malate dehydrogenases (where R = COOH) and the broad specificity enzymes [where R = (CH<sub>2</sub>)<sub>2-7</sub>-CH<sub>3</sub> or R = (CH<sub>2</sub>)<sub>*n*</sub>-aromatic]. Glyoxalate (H-CHOH-COOH) is a good substrate of L-LDH, but mandelate (phenyl-CHOH-COOH) is not. Some plant malate dehydrogenases will use the NADP(H) cofactor in the **glyoxysome**, but most select NAD<sup>+</sup> over NADP<sup>+</sup> by >100 : 1.

**Table 1. L-LDH Species Distribution: Genes, Enzyme Subunits, Isoenzymes, and Regulation**

| Tissue and Gene                                    | Regulation   | L-LDH-subunit  | L-LDHTetramer                                       | Isoenzyme Name and Charge       | Inhibition by NAD <sup>+</sup> + Pyruvate |
|--|--|----------------|---|---------------------------------|---|
| Anaerobic skeletal muscle: <i>lct a</i>            | Hypoxia inducible factor -1 + O <sub>2</sub> supresses gene transcription  | A or M subunit | A <sub>4</sub> or M <sub>4</sub> and hybrids with B | Muscle, just positively charged | Weakly                                    |
| Aerobic heart muscle: <i>lct b</i>                 |  | B or H subunit | B <sub>4</sub> or H <sub>4</sub> and hybrids        | Heart, high negative charge     | Strongly                                  |
| Mature testis (birds/mammals): <i>lct c</i>        |  | C or X subunit | C <sub>4</sub> or X <sub>4</sub>                    | Testis, nearly neutral          | Moderate                                  |
| Bony fishes: <i>lct e</i> and <i>lct f</i>         | Not known  | E or F subunit | E <sub>4</sub> and F <sub>4</sub>                   | Negative charge                 | Strong                                    |
| Bacterial: One gene for each of L-LDH and/or D-LDH | None or for some LDHs FBP assembles a low <i>K<sub>M</sub></i> tetramer from a high <i>K<sub>M</sub></i> (pyruvate) dimer. | L-             | Homo-tetramers or dimers                            | High negative charge            | Strongly                                  |

## 2. Metabolic Role of L-LDH, Isoenzymes, and Organ-Specific Diagnosis



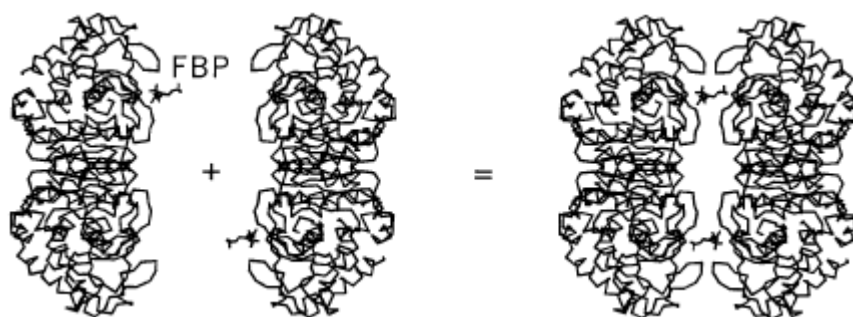
A frequent metabolic role of both L-LDH and D-LDH is to enable continuous conversion of glucose ( $C_6H_{12}O_6$ ) to two molecules of excreted lactic acid ( $2 \times C_3H_6O_3$ ) and ATP (energy) in anaerobic tissues and cells. The final oxidation of NADH to  $NAD^+$  in the pyruvate to lactate reaction regenerates the  $NAD^+$  that is consumed at the glyceraldehyde 3-phosphate dehydrogenase step earlier in glycolysis. This anaerobic role is found in higher organisms and **microorganisms**. The lactic acid is excreted. In animals, it is circulated back to the heart for further oxidation. In microorganisms, it results in acidification of the growth medium (particularly in *Lactobacilli*). In the oxygen-rich organs of higher animals, L-LDH is required to catalyze reaction (1) in the reverse direction, so that lactate from the tissues can fuel the heart.

Specialized A, B, and C-polypeptide chains (Table 1) have evolved with common catalytic mechanisms, but different susceptibilities to inhibition by  $NAD^+$ -pyruvate to suit these different metabolic roles. When newly synthesized on the ribosome, the A and B subunits can combine with each other to yield a distribution of five stable **isoenzymes**  $A_4$ ,  $A_3B$ ,  $A_2B_2$ ,  $AB_3$ , and  $B_4$ , which do not further exchange subunits without mild **denaturation**. The large difference in net charge between the A and B subunits enables the different tetramers to be easily and quickly separated by [electrophoresis](#) into five distinct bands, whose relative distribution is characteristic of each organ. When an organ is diseased, the cells die and their contents are released to the blood. A sudden increase in LDH isoenzymes in blood serum indicates the presence of disease, and the isoenzyme distribution can be used to diagnose the diseased organ.

### 3. L-LDH Regulation

In mammals, the proportions of A (also designated as M) and B (or as H) polypeptide chain made are regulated at the level of **gene expression**. An oxygen-sensitive [transcription factor](#), called hypoxia inducible factor-1, suppresses the synthesis of LDH A and the  $A_4$  isoenzyme in oxygen-rich tissues such as heart muscle (1). At low oxygen concentration, the gene for LDH A is transcribed, and  $A_4$  becomes the main isoenzyme in anaerobic tissues. In bacteria, LDH gene expression is constant, but the kinetic activity of the enzyme increases in the presence of fructose-1,6-bisphosphate (FBP) (see Fig. 1). Regulation of the D-LDH enzymes has not yet been studied extensively.

**Figure 1.** Regulation of many bacterial L-LDHs by fructose-1,6-bisphosphate (FBP). Two FBP molecules bind across the P-axis of the tetramer between the Arg-173, His-188 basic patches on each dimer. Without FBP, the mild **thermophile** enzymes dissociate to dimers, which is reversed only with a high apparent  $K_m$  for pyruvate, 5 mM; consequently, these enzymes achieve only a low kinetic flux in tissues with a low concentration of pyruvate (5). The enzymes from **mesophile** bacteria open the P-axis contact but do not dissociate. With FBP bound, the compact tetramer is assembled. The mesophile enzymes have a low  $K_m$  (0.1 mM pyruvate) and are active at cellular concentrations of pyruvate.



#### 4. D-LDH: Metabolic Role, Enzyme Family

L-LDH is the common form of the LDH enzyme, but lower organisms (arthropods, molluscs, annelids, and coelenterata) and bacteria also express D-LDH. There are now 22 gene sequences known for the D (or **R**)-2-hydroxyacid: NAD<sup>+</sup> oxidoreductase family in plants, yeasts, and bacteria. These include variants of D-LDH, where the CH<sub>3</sub>- in D-lactate is replaced, with the same chirality, by HO-CH<sub>2</sub>- (D-glycerate DH), H- (formate DH), H<sub>2</sub>O<sub>3</sub>P-O-CH<sub>2</sub>- (phosphoglycerate DH), H<sub>2</sub>O<sub>3</sub>P-O-CHOH- (erythronate-4-phosphate DH) or (CH<sub>3</sub>)<sub>2</sub>-CH<sub>2</sub>- (hydroxyisocaproate DH). With one exception, these enzymes use NAD<sup>+</sup> and not NADP<sup>+</sup> as the cofactor. Some organisms, such as *Lactobaccillus plantarum*, have both the D- and L-LDH enzymes. Enzymes with more complex substrates function in bacterial serine biosynthesis, the serine cycle, and methanol oxidation.

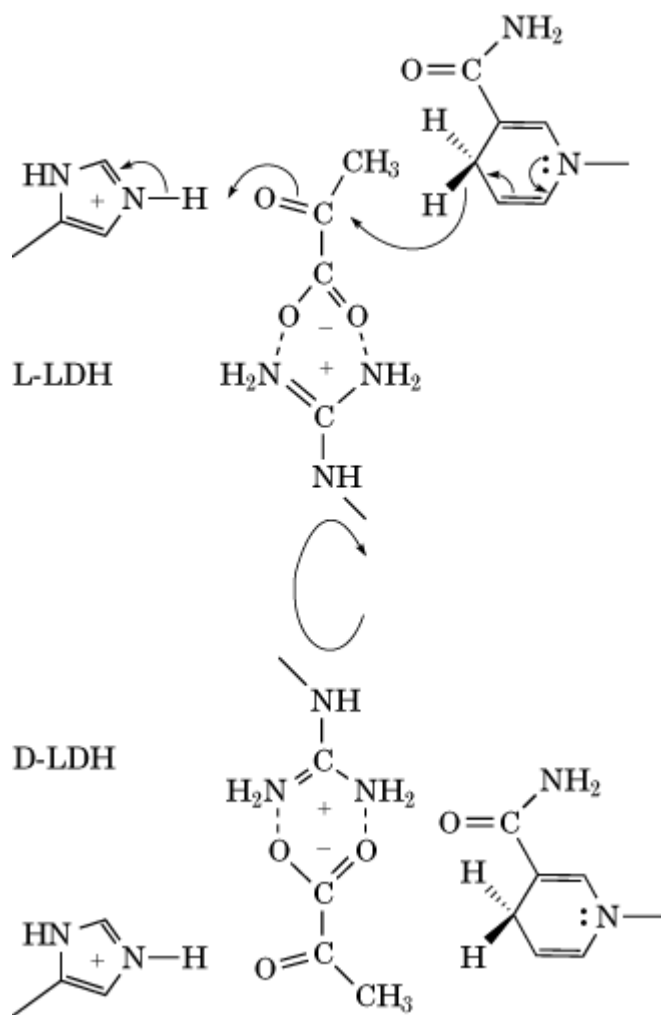
#### 5. Catalytic Mechanism

Members of both the L- and D-LDH families have analogous amino acid residues at their active sites:

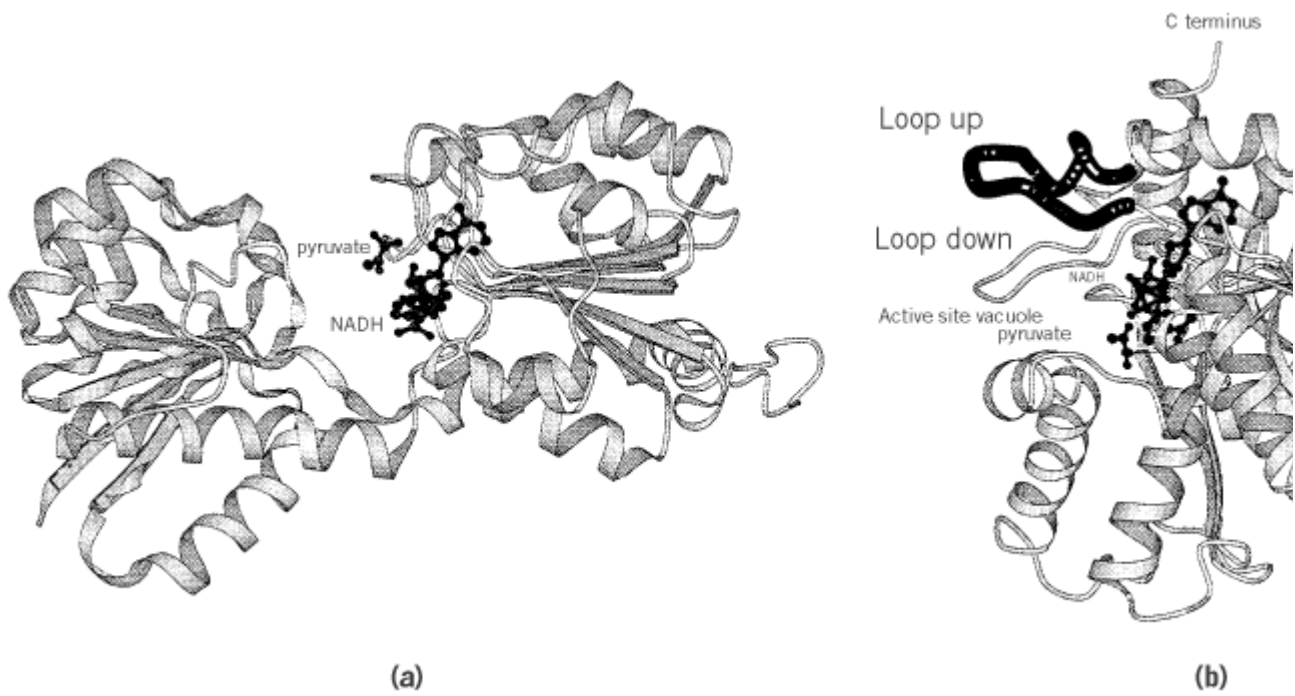
1. A catalytic histidine-aspartic (or glutamic acid) couple as an H<sup>+</sup> donor/acceptor
2. An arginine-glycine pair that binds the α-carboxylate of the pyruvate substrate
3. A glycine-rich (GxGxxG or GxAxxG) turn next to the adenine-ribose
4. An aspartic acid to select NAD<sup>+</sup> over NADP<sup>+</sup> some 25 residues further to the C-terminus

The chemical reaction involves the reversible transfer of a pair of electrons from the ring nitrogen of NADH via the >C = O of pyruvate to the protonated histidine (see Fig. 2, top). The binding of the pyruvate substrate induces a surface protein loop to fit over the protein-NADH complex (Fig. 3) to generate an internal vacuole that only accommodates small ketoacids with a single negative charge (see [Induced Fit](#)). This vacuole determines the substrate specificity; oxaloacetate becomes a good substrate if Gln102 on the underside of the mobile surface protein loop (Fig. 3) is changed by mutagenesis into positively charged arginine (2). Random variation of the bases coding for the inside of the vacuole enabled the design of almost any 2-hydroxyacid dehydrogenase (3). Kinetic **fluorescence** studies show that the slowest step in the L-LDH mechanism is the movement of this peptide loop. The two-electron chemical reduction is too fast to measure (4).

**Figure 2.** Stereochemical differences in the active sites of L- and D-LDH. The evolutionary change of an L-LDH to a D-LDH involves the change of the coenzyme-binding pocket (the Rossmann fold) and the catalytic histidine residue in constant orientation. Evolution also rotates the α-carboxyl-binding arginine residue by 180° and thus inverts the stereochemistry of the reduction.



**Figure 3.** Comparison of the structures of a single subunit of the dimeric D-LDH (a) and a single subunit of the normally D-enzyme has a bi-lobed structure in which the larger (rightmost) lobe is a cofactor-binding Rossmann fold (babab)<sub>2</sub>. The between the two lobes. The L-enzyme also has a coenzyme-binding Rossmann fold shown in a similar orientation to that substrate is bound in a small vacuole that is formed when a mobile surface loop (shown in black) moves down over a ber surface. This loop movement is the slowest event in the catalytic cycle. The movement brings a catalytic arginine residue defines the size and charge of the catalytic vacuole (6).



## 6. How the Protein Determines L- and D-Handedness of Lactate

Figure 2 shows that the same catalytic chemistry can be used to reduce a ketoacid to an L- or D-hydroxyacid if the spatial relation of the catalytic histidine and nicotinamide ring is kept constant but the ketoacid is rotated through  $180^\circ$ , so that the hydride approaches from the other face of the ketone. The L-LDH and D-LDH families have a common element of [supersecondary structure](#) (a **Rossmann fold**) that binds the cofactor along the C-terminal edge of a 6-stranded [beta-sheet](#). This fold is the core of the single domain of L-LDH and the largest domain of the bi-lobed D-LDH (Fig. 3). The reduction to either D- or L-lactate is achieved by the arginine, which binds the  $\alpha$ -carboxyl group of the substrate, being rotated in the new domain of the protein structure by  $180^\circ$ . Seven [hydrogen bonds](#) tether the pyruvate firmly and prevent no more than 1 part in  $10^5$  reduction from the incorrect carbonyl face.

### Bibliography

1. J. D. Firth, B. L. Ebert, and P. J. Ratcliffe (1995) *J. Biol. Chem.* **270**, 21021–21027.
2. H. M. Wilks, K. W. Hart, R. Feeney, C. R. Dunn, H. Muirhead, W. N. Chia, D. A. Barstow, T. Atkinson, A. R. Clarke, and J. J. Holbrook (1988) *Science* **242**, 1541–1544.
3. A. S. El Hawrani, K. M. Moreton, R. B. Sessions, and J. J. Holbrook (1996) *J. Mol. Biol.* **264**, 97–110.
4. A. R. Clarke, T. Atkinson, and J. J. Holbrook (1989) *Trends Biochem. Sci.* **14**, 101–105, 145–148.
5. A. D. Cameron, D. I. Roper, K. M. Moreton, H. Muirhead, J. J. Holbrook, and D. B. Wigley (1994) *J. Mol. Biol.* **238**, 615–625.
6. A. R. Clarke, D. B. Wigley, W. N. Chia, D. A. Barstow, T. Atkinson, and J. J. Holbrook (1986) *Nature* **324**, 699–702.

### Suggestions for Further Reading

7. J. J. Holbrook, A. Liljas, S. J. Steindel, and M. G. Rossmann (1975) *Lactate dehydrogenase, The Enzymes*, **11a**, 191–292. Covers the structure and mechanism of mammalian L-LDHs. Some amino acid sequence errors of 1975 were not corrected until the cDNAs were sequenced in the

1980s.

8. H. M. Wilks, A. Cortes, D. C. Emery, D. J. Halsall, A. R. Clarke, and J. J. Holbrook (1992) Opportunities and limits in creating new enzymes: Experiences with the NAD-dependent lactate dehydrogenase frameworks of humans and bacteria, *Annals NY Acad. Sci.* **672**, 80–93. A review of the many useful enzymes designed and synthesised on the *B. stearothermophilus* and human L-LDH framework.

## Lactoferrin

Lactoferrin is a nonheme iron-binding **glycoprotein** and a member of the [transferrin](#) gene family (1). The protein is localized in neutrophils and to exocrine secretions associated with local host defense. The biological functions of lactoferrin are associated primarily with mammalian host defense, where the protein has a multifunctional role. These functional properties include antibacterial and antiendotoxin activities, as well as direct modulation of [immune responses](#).

### 1. Structure

Lactoferrin consists of a single polypeptide chain 692 amino acid residues in length, with a molecular weight of 76 kDa. Amino acid [sequence analysis](#) and resolution of the [X-ray crystallography](#) structure of lactoferrin have demonstrated that the protein consists of two **domains** located in the N- and C-terminal halves of the protein, separated by an **a-helical** hinge region, and sharing significant amino acid [homology](#) (~40% identity). The two lobes appear to have arisen by **gene duplication** of a single primordial gene (1). Each lobe binds one atom of ferric iron, which is coordinated to six ligands; two of these are provided by a bicarbonate anion that neutralizes positively charged amino acid residues to facilitate iron coordination to four additional ligands, provided by an **aspartate**, two [tyrosine](#), and a [histidine](#) residues in each lobe (2). The protein contains two N-linked glycosylation sites provided by one [arginine](#) residue in each lobe, to which glycan residues are attached. The carbohydrate composition varies between species (3) and at different tissue sites of expression within a given species (4).

### 2. Distribution

Lactoferrin was first isolated from milk, in which it occurs as the second most abundant protein (5). However, it is expressed in all glandular epithelial cells and is a major component of exocrine secretions associated with primary host defense, including mucosal, reproductive, salivary, and lachrymal secretions (6). The protein is also produced at high levels in myeloid cells of the central immune system, specifically in the secondary granules of neutrophils (7).

### 3. Biological Functions

The primary functions of lactoferrin are associated with iron binding and mammalian host defense. The physiological role of lactoferrin with regard to iron transport and delivery to mammalian cells is still a subject of debate. Previous studies to address its role in intestinal iron absorption have provided apparently conflicting results. Although specific enterocyte receptors for lactoferrin have been identified (8), and uptake of the iron-saturated protein has been demonstrated in cultured intestinal cells, intestinal iron delivery by lactoferrin has not yet been observed *in vivo*. In contrast, recent studies indicate that depletion of lactoferrin from human milk results in an increase in iron absorption in human infants (9), suggesting that the protein may be involved in sequestration and

detoxification of free iron; this conclusion is supported by the high affinity of iron for the protein and high stability of the protein in the gastrointestinal tract (10).

Lactoferrin is an important component of both the local and central immune systems and contributes to host defense by (1) acting as a direct antibacterial agent to inhibit the growth of bacteria and (2) directly interacting with local and central immune cells to regulate immune responses. The antibacterial properties of lactoferrin are imparted by two distinct mechanisms, involving two separate domains of the protein. The first is the iron-binding domain, which causes a retardation of bacterial growth due to iron sequestration and deprivation (11). The second domain comprises a cationic region contained within the first 47 amino acids at the aminoterminal of human lactoferrin (12). When isolated as a peptide, this segment interacts with lipopolysaccharides (LPS) located in the outer cell wall of bacteria and causes a disruption of bacterial membranes, leading to a direct bactericidal effect (12). LPS acts as a potent inflammatory agent by interacting with specific immune cell **receptors** to induce expression of proinflammatory **cytokines**. The LPS-binding property of the bactericidal peptide, therefore, has been proposed to be responsible for the antiendotoxin activity of lactoferrin observed in LPS-treated mice (13) and in LPS-stimulated mononuclear cell cultures (14), resulting in inhibition of proinflammatory cytokine production.

In addition to its antibacterial activity, a significant body of evidence has accumulated to support a direct role for lactoferrin in modulation of local and systemic immune responses. Specific receptors for lactoferrin have been identified on monocytes, macrophages, platelets, and lymphocytes (15-17). These receptors are thought to mediate a direct cytokine activity of lactoferrin that results in stimulation of lymphocyte cell proliferation and regulation of myeloid cell development by altering the expression of myelopoietic factors. The protein also has been shown to have anti-tumor activity by activating natural killer cells (18). The structure and identity of the immune receptors for lactoferrin or the specific intracellular signaling mechanisms by which they mediate the cytokine functions of lactoferrin have yet to be established.

#### 4. Gene Structure and Regulation

Lactoferrin is encoded by a single-copy gene that is located on [chromosome 3q21](#) in humans and chromosome 9 in the mouse. The mouse and bovine lactoferrin genes have been characterized; both are 30- to 35-kb in size and are organized into 17 **exons** separated by 16 **introns** (19, 20). [Cis-acting](#) sequences have been identified within the **promoter** region of the lactoferrin gene that mediate its regulation by **estrogen** in the uterus (21) and may mediate its induction by **mitogens**, **cyclic AMP**-dependent signaling pathways and as an acute phase response gene (20). While a single [messenger RNA](#) species is expressed from the gene in most tissues, recent evidence suggests that a second mRNA species lacking a secretion **signal sequence**, as well as an open [reading frame](#) encoding most of the aminoterminal bactericidal segment, can be generated in a tissue-specific manner by [alternative splicing](#) and alternate promoter usage (22). However, whether this alternative RNA is translated *in vivo* has yet to be determined. **Complementary DNA** clones corresponding to both mRNAs have been isolated (22, 23). The full-length cDNA corresponding to native human lactoferrin has been expressed in mammalian cells (24), **transgenic** mice (25), and filamentous **fungi** (26). The availability of these production systems will facilitate a detailed investigation of the structural parts of lactoferrin that are required for its various functional activities.

#### Bibliography

1. M.-H. Metz-Boutigue et al. (1984) Eur. J. Biochem. **145**, 659–676.
2. B. F. Anderson et al. (1989) J. Mol. Biol. **209**, 711–734.
3. G. Spik, B. Coddenville, and J. Montreuil (1988) Biochem. **70**, 1459–1469.
4. P. Derisbourg et al. (1990) Biochem. J. **269**, 821–825.
5. P. L. Masson and J. F. Heremans (1971) Comp. Biochem. Physiol. **39**, 119–129.
6. P. L. Masson, J. F. Heremans, and C. Dive (1966) Clin. Chim. Acta. **14**, 735–739.

7. P. L. Masson, J. F. Heremans, and E. Schonke (1969) *J. Exp. Med.* **130**, 643–658.
8. S. Iyer and B. Lonnerdal (1993) *Eur. J. Clin. Nutr.* **47**, 232–241.
9. B. Lonnerdal and S. Iyer (1995) *Ann. Rev. Nutr.* **15**, 93–110.
10. G. Spik et al. (1982) *Eur. J. Biochem.* **121**, 413–419.
11. J. J. Bullen, H. J. Rogers, and E. Griffiths (1978) *Curr. Top. Microbiol. Immunol.* **80**, 1–35.
12. W. Bellamy et al. (1992) *Biochem. Biophys. Acta.* **1121**, 130–136.
13. M. Machnicki, M. Zimecki, and T. Zagulski (1993) *Int. J. Exp. Path.* **74**, 433–439.
14. S. P. M. Crouch, K. J. Slater, and J. Fletcher (1992) *Blood* **80**, 235–240.
15. H. S. Birgens et al. (1983) *Br. J. Haematol.* **54**, 383–391.
16. B. Leveugle et al. (1993) *Eur. J. Biochem.* **213**, 1205–1211.
17. J. Mazurier et al. (1989) *Eur. J. Biochem.* **179**, 481–487.
18. H. Shau, A. Kim, and S. Golub (1997) *J. Leukoc. Biol.* **51**, 343.
19. G. A. Cunningham, D. R. Headon, and O. M. Conneely (1992) *Biochem. Biophys. Res. Commun.* **189**, 1725–1731.
20. H. M. Seyfert et al. (1994) *Gene* **143**, 265–269.
21. C. T. Teng (1994) *Adv. Exp. Med. Biol.* **357**, 183–196.
22. P. D. Siebert and B. C. Huang (1997) *Proc. Natl. Acad. Sci. USA* **94**, 2198–2203.
23. M. J. Powell and J. E. Ogden (1990) *Nucleic Acids Res.* **18**, 4013.
24. K. M. Stowell et al. (1991) *Biochem. J.* **276**, 349–355.
25. G. J. Platenburg et al. (1994) *Transgenic Res.* **3**, 99–108.
26. P. P. Ward et al. (1995) *Biotechnology* **13**, 498–503.

### Suggestions for Further Reading

27. L. Ferenc and M. Viljoen (1995) *Hematologica* **80**, 252–267.
28. J. Brock (1995) *Immunol. Today* **16**, 417–419.

## Lambda CI and CRO Repressors

The [lambda phages](#) of *Escherichia coli* and *Salmonella typhimurium* encode two helix-turn-helix repressor proteins, CI and Cro, containing [helix-turn-helix motifs](#) found in other [DNA-binding proteins](#). Both of these repressors bind to a set of six similar, but nonidentical sites or **operators** in the phage [genome](#). These sites are grouped in sets of three on each side of the structural gene for the CI protein. One set, including the operators OR1, OR2, and OR3, overlaps two **promoters**; PR, which directs the synthesis of a set of genes expressed immediately after phage infection or induction, and PRM, which directs the production of the **messenger RNA** for the CI repressor itself. The second set of operators, OL1, OL2, and OL3, lies on the other side of the CI gene and overlaps the PL promoter, which directs the synthesis of another set of early mRNA species. The affinities of the CI repressor for these six sites differ, as also do those of the Cro repressor. It is this differential affinity that creates a regulatory switch determining the difference between [lysogeny](#) and lytic growth.

The structures of these two repressor proteins differ substantially. Whereas the Cro protein of phage λ contains 56 amino acid residues and exists in solution principally as a dimer, the CI repressor

contains 236 residues and can exist as a monomer, dimer, or tetramer in solution. Each CI monomer contains two domains: an N-terminal domain of some 90 residues that contains the helix-turn-helix motif and binds specifically to the operator sites, and a C-terminal domain that does not interact with DNA. Between these two regions is a flexible hinge.

When bound to DNA dimers, the CI repressor has the potential to form tetrameric complexes via their C-terminal domains. The organization of the tripartite leftward and rightward operators of phage  $\lambda$  allows such interactions between dimers bound to two adjacent sites. However, once such cooperative binding has occurred, the binding of a repressor dimer to the third element of the operator site is then independent of the other repressor dimers. For the rightward operator, the initial binding of the repressor occurs at OR1. This interaction facilitates binding to OR2, which is a slightly weaker site. Although OR2 is weaker, initial occupancy will occur at a significant frequency at this site and will facilitate binding to OR1, thereby increasing the probability of occupancy of both OR1 and OR2. By contrast, the affinity of a repressor dimer for an isolated OR3 site is 25-fold less than that for OR1, and consequently occupancy of OR3 only occurs at high repressor concentrations in the wild-type operator. If, however, the affinity of OR1 for repressor is reduced by mutation, cooperative interactions can then take place between repressor molecules bound at OR2 and OR3.

Since the occupancy of adjacent binding sites by repressor dimers is dependent on an interaction between two protein molecules, the concentration dependence of occupancy is no longer hyperbolic—as it would be for independent binding—but instead is sigmoidal. An important consequence of this dependency is that occupancy increases rapidly over a small range of repressor concentration and attains a level of >99% at concentrations significantly lower than those that would be required for the same level of independent binding. It is consequently a property of this system—and of other cooperative protein-DNA binding—that any modification of the [protein-protein interactions](#) required for cooperativity will result in a substantial change in occupancy.

In the life cycle of phage  $\lambda$ , the modification that occurs is the **proteolytic** cleavage by RecA protein of the CI repressor between the N-terminal and C-terminal domains on induction of the prophage. The result of this cleavage is that cooperative interaction between bound dimers can no longer occur and therefore the effective concentration of the repressor is reduced to a level at which the occupancy of both operator sites is very low. Repression is thus relieved, and transcription can proceed from both PR and PL, resulting in the expression of Cro from PR.

By contrast to the CI repressor, Cro binds noncooperatively to the individual operator sites and has a higher affinity for OR3 than for OR1. Cro thus preferentially represses initiation from PRM and prevents the further expression of the CI repressor after prophage induction.

In addition to its repressor function, the CI protein can activate transcription from PRM. This activation is mediated by direct protein-protein contacts between the helix-turn-helix motif and **RNA polymerase** holoenzyme, involving the exposed surfaces of helix 2 and the turn between helices 2 and 3.

#### Suggestions for Further Reading

- A. Johnson et al. (1979) Interactions between DNA-bound repressors govern regulation by the  $\lambda$  phage repressor. *Proc. Natl. Acad. Sci. USA* **76**, 5061–5065.
- A. Hochschild and M. Ptashne (1986) Cooperative binding of  $\lambda$  repressors to sites separated by integral turns in the double helix. *Cell* **44**, 681–687.

## Lambda Phage



Bacteriophages are viruses that infect and reproduce within bacterial cells. Phage  $\lambda$  has been one of the most intensively studied of all the bacteriophages. It has provided models for biological processes such as [DNA replication](#), [transcription](#), regulation of gene expression, **developmental switches**, informational suppression, homologous and site-specific [recombination](#), **restriction-modification**, and morphogenetic pathways, among others, as well as an incomparable window into the biology of its host, *Escherichia coli*. The [genome](#) of  $\lambda$  was among the earliest complete genomes to be sequenced (1). Such detailed knowledge allowed  $\lambda$  to be exploited further as a **vector** for [recombinant DNA](#) experiments, to the point that  $\lambda$  is a standard tool of the contemporary molecular genetics laboratory, as well as a subject of study in its own right in laboratories worldwide. Phage  $\lambda$  is a member of the temperate group of bacteriophages, which means that an infection may progress along either of two alternative pathways. In the lytic pathway, the host cell is destroyed (“lysed”) concomitant with the production and release to the environment of a hundred or so progeny  $\lambda$ . In the lysogenic pathway, the infected bacterium survives in the form of a carrier cell known as a lysogen. This chapter focuses on lytic development of  $\lambda$  (see [Lysogeny](#) for the alternative pathway).

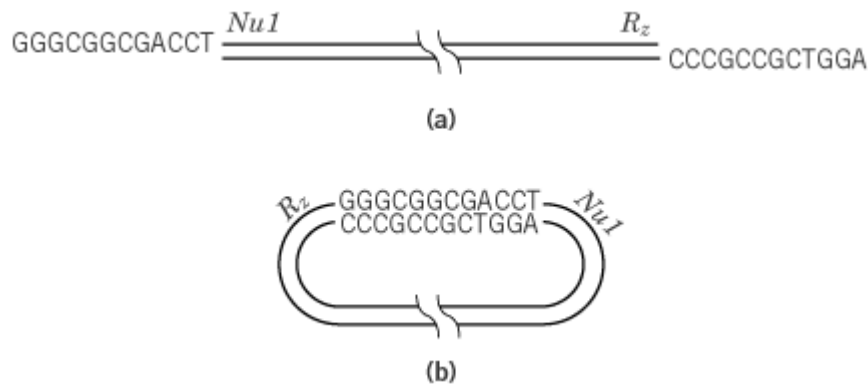
## 1. Initiation of Infection

A particle of  $\lambda$  consists of about half protein and half DNA. The protein comprises an icosahedral head, about 50 nm in diameter, attached to a flexible tubular tail about 150 nm long. The DNA, which is contained within the head, is a linear double-stranded molecule, 48,490 base-pairs in length, not counting the complementary single-stranded extensions of 12 nucleotides in length at the 5' end of each strand. These [cohesive, sticky ends](#) (cos) are generated in the course of phage morphogenesis as discussed below. Infection begins with the attachment of  $\lambda$  to the target cell, *E. coli* K-12. The cellular receptor is an outer [membrane protein](#), LamB, which is also part of the transport system for maltose.  $\lambda$  interacts with LamB through a protein located at the tip of its tail. After phage attachment, the DNA exits the phage head through the tail and enters the bacterial cytoplasm after traversing the outer membrane, periplasm, and inner membrane by an unknown mechanism. Unlike some other phages,  $\lambda$  does not use a contractile tail to assist this process.

Immediately upon DNA entry into the cytoplasm, the single-stranded extensions of the DNA pair with one another, and the DNA is converted to a covalently closed circle through the action of host-cell [DNA Ligase](#) (Fig. 1). This simple step is essential to the successful completion of lytic development for at least two reasons. First, replication of the  $\lambda$  DNA requires **supercoiling**, which is only possible in a covalently closed molecule. Second, transcription of approximately half the  $\lambda$  genome (the antiterminated  $P'_R$  transcript; see text below) traverses this site of end joining.

Incidental to the main topic here, circularization is also essential to the lysogenic pathway of development, because a circle is the substrate for integration of the phage DNA into the host [chromosome](#). From this point in the infection cycle, several biochemical activities are carried out on the DNA more or less simultaneously, including transcription, DNA replication, and site-specific recombination (lysogenic development only). For simplicity, transcription and replication are considered separately.

**Figure 1.** Circularization of the  $\lambda$  genome. (a) The linear genome as found in the virion. The sequences of the single-stranded 5' extensions are shown. Upon infection, the extensions pair and the nicks are sealed by *E. coli* DNA ligase to form the circular genome as shown in b. The positions of genes flanking the joint are shown to facilitate comparison with the transcription map of Figure 2.

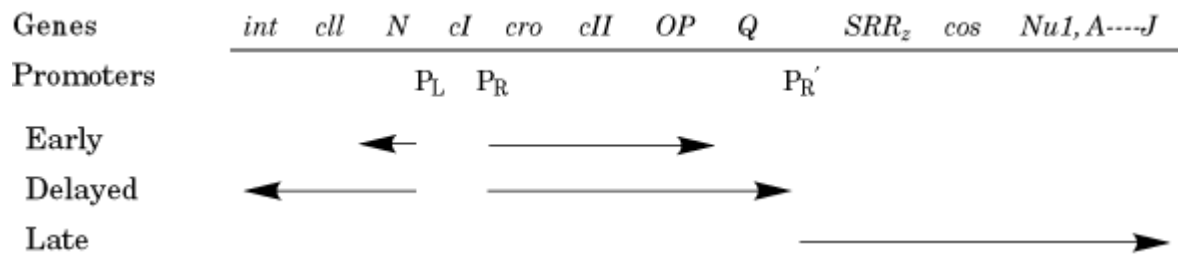


## 2. Transcription

The early transcription program is compatible with either lytic or lysogenic development. Soon, however, each  $\lambda$  infection tilts decisively one way or the other, and the programs thereafter are distinct. It therefore should come as no surprise that  $\lambda$  has deployed timing mechanisms to delay the onset of expression of some genes until the decision for lysis or lysogeny has been made. Such timing prevents the inappropriate expression of, for example, gene products that would lyse a cell that was destined to become a lysogen. The principal feature of the timing is a regulatory cascade in which products of early transcription are required for expression of late genes. These positive activators are antiterminators of transcription, agents of a mechanism of gene control that arguably has attained its most sublime refinement in  $\lambda$ .

The outlines of the lytic transcription program are depicted in Figure 2. The host **RNA polymerase**, unaided by other factors, recognizes three **promoters** within  $\lambda$  DNA, designated  $P_L$  (leftward),  $P_R$  (rightward), and  $P_R'$  (the so-called late promoter). In the absence of phage-encoded **antitermination** factors, the  $P_L$  transcript terminates after gene *N*. About half the  $P_R$  transcripts terminate after gene *cro*, the remainder after gene *P*. The  $P_R'$  transcripts terminate almost immediately to produce a short (194 nucleotide, or 6S) RNA that does not encode any protein. Thus the early proteins synthesized by  $\lambda$  are *N* and *cro* and comparatively lesser amounts of *cII*, *O*, and *P*. The *N* protein is one of  $\lambda$ 's transcription antiterminators. It is specific for transcription beginning at  $P_L$  and  $P_R$ . As the level of *N* increases, new transcription beginning at  $P_L$  ignores the terminator after gene *N* and reads genes to the left of *N*, including *cIII*, *red*, *xis*, and *int*. Similarly, new transcription beginning at  $P_R$  ignores the weak terminator after gene *cro* as well as the strong terminator after gene *P*, with the result that synthesis of *cII*, *O*, and *P* increases, and a new gene, *Q*, is transcribed for the first time. The *Q* protein is  $\lambda$ 's other transcription antiterminator. It is specific for transcription beginning at  $P_R'$ . As the level of *Q* increases, new transcription beginning at  $P_R'$  ignores its strong terminator and reads genes to the right of *Q*, including *S*, *R*, and, after traversing the *cos* site as mentioned above, some 20 additional genes encoding structural proteins of the phage particle and **enzymes** involved in morphogenesis.

**Figure 2.** Lytic transcription program. The map shows some  $\lambda$  genes (not to scale) and the locations of three promoters used in the lytic program. Note that the template is actually circular (or later, concatemeric), but for convenience it is shown as linearized between genes *J* and *int*. Early transcripts are formed by the host RNA polymerase, delayed transcripts require  $\lambda$  *N* protein, and late transcripts require the  $\lambda$  *Q* protein.



In summary then,  $\lambda$  transcripts may be classified as early, delayed, and late. Early transcription results in the synthesis of N protein, which stimulates transcription of the delayed class of genes. The latter includes *Q*, whose product is required for late gene transcription.

How is this pattern of gene expression altered in infections leading to lysogeny rather than lytic growth? Such infections are characterized by higher levels of cII protein (see [Lysogeny](#)). cII is a transcriptional activator specific for three promoters in  $\lambda$ . One of these promoters is the  $P_{RE}$  promoter, which is responsible for a burst of **lambda repressor** synthesis. The repressor turns down expression of early and delayed transcripts by preventing transcription initiation at  $P_L$  and  $P_R$  (see **Lambda repressor**). In addition, cII activates the  $P_{aQ}$  (antisense-Q) promoter. As the name implies, this promoter forms a noncoding **antisense** transcript of gene Q, which delays expression of Q protein and the late genes (2). During infections characterized by high levels of cII, these two effects conspire to delay and ultimately to prevent late-gene expression, exactly as required if the cell is to be successfully lysogenized. (The third promoter controlled by cII,  $P_{int}$ , is also essential for lysogenization, but has no effect on the lytic transcription pattern.) Just as a commitment to lysogeny entails reduction of early and delayed transcription as described above, curiously, so also does commitment to lytic development. In this case, however, the cII level is low, and the repressor is not made in significant amounts. Instead, the Cro protein predominates. This protein, like the repressor, binds near  $P_L$  and  $P_R$  and prevents initiation at those promoters. Because Cro is a product of  $P_R$ , this is an autoregulatory circuit (see **Cro protein**).

With two transcriptional antiterminators in its arsenal, how does  $\lambda$  arrange for N to be specific to  $P_L$  and  $P_R$  transcripts and Q to be specific to  $P_R'$  transcripts? This is best understood for the case of N. The  $P_L$  and  $P_R$  transcription units each contain a site called *nut* (for N-utilization; *nutL* and *nutR*, respectively). The *nut* transcript is a recognition site for N, which binds the nascent RNA at this site and transfers to the transcription complex. N remains with the transcription complex, endowing it with the ability to ignore potential termination sites encountered thereafter (3). Transcription units without *nut* are indifferent to N; conversely, hybrid transcription units in which *nut* has been inserted become competent for N-mediated antitermination. Q protein, in contrast, binds to the DNA of the  $P_R'$  promoter and interacts with the transcription complex while it is paused at nucleotide 16 or 17 of the transcript (4) (see **Antitermination**).

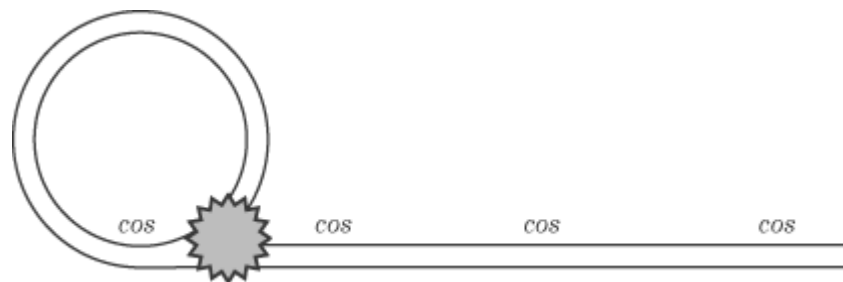
### 3. DNA Replication

For the most part, phage  $\lambda$  uses replication enzymes of the host cell (see [DNA Replication](#)). The main thing  $\lambda$  must do is divert these enzymes from their normal responsibilities and direct them to the phage DNA. For this purpose,  $\lambda$  encodes two specific replication proteins, called O and P. They are synthesized early (and also delayed), as described above; thus replication can begin soon after the initiation of infection. The O protein nucleates the assembly of a competent replication complex for  $\lambda$  by binding to the phage **origin of replication**. The P protein, a functional homologue of the host DnaC protein, binds to the host replicative [DNA helicase](#) (DnaB protein) and delivers it to the phage

origin complex by interacting with the O protein. DnaB is disassembled from P and is thus activated for its replication function with the help of three host-encoded **heat-shock** proteins, DnaJ, DnaK, and GrpE (5) (see [DnaK/DnaJ Proteins](#)).

DNA replication proceeds in two distinct phases. Initially, the substrate is a covalently closed circle, formed by *cos* joining in the early moments of infection. A pair of oppositely oriented [replication forks](#) set up at the origin form a replication bubble that grows bidirectionally until the entire circle is replicated. This is so-called  $\theta$  (theta) replication, named after the shape of the partially replicated circle. Later, replication switches to the [rolling circle DNA replication](#) mode (Fig. 3). The mechanism of this switch is not understood. In rolling circle replication, a single replication fork moves perpetually around a circle, spinning off a long linear tail, much like paper towels unwinding from a roll. Rolling circle replication is an important preparatory step to phage morphogenesis, because  $\lambda$  normally uses the concatemeric (i.e., multiple genome-length) DNA tail as the substrate for packaging.  $\lambda$  has a further adaptation to facilitate rolling circle replication. Ordinarily, double-stranded DNA ends, such as the end of the rolling circle tail, are unstable in *E. coli* because they are degraded by the RecBCD nuclease (see [Recombination](#)).  $\lambda$  encodes a RecBCD inhibitor (the *gam* gene product) to avoid this complication.

**Figure 3.** Rolling circle replication. A perpetual replication fork moving around a circular template that is monomeric, as indicated by the single *cos* site, produces a concatemeric tail. Each *cos*-to-*cos* interval of the tail carries a complete copy of the  $\lambda$  genome and reflects one passage of the replication fork around the circle. The tail is the substrate used for  $\lambda$  packaging.



#### 4. Morphogenesis and Lysis

The result of  $\lambda$  lytic infection is the production, within about 60 minutes, of approximately 100 new phage particles per infected cell. As DNA replication progresses, the replicated DNA can also be used as templates for transcription. Thus there is an accelerating pace of late protein production as infection proceeds, and this production of ingredients only stops when the cell lyses. Indeed,  $\lambda$  mutants blocked in lysis can accumulate as many as 1000 mature phage particles per cell. At some point, everything is ready for phage particle assembly to begin.

In broad outline, the construction of a  $\lambda$  particle involves the separate assembly of tails and empty, somewhat shrunken heads, called proheads. The proheads expand approximately twofold in volume as they are filled with DNA, which is excised from the concatemer by cutting at the *cos* sites by the phage-encoded terminase protein, in coordination with packaging. The first cut occurs at a randomly chosen *cos* site on the concatemer. The two ends thus produced are different; they are (i)  $\lambda$ 's "right" end (the  $R_z$  end of the genetic map) carrying the single-stranded DNA extension 5'

AGGTCGCCGCC and (ii) its "left" end (the *NuI* end of the genetic map) carrying the single-stranded DNA extension 5' GGGCGGCGACCT. After cutting, terminase remains bound to the  $\lambda$  left end, which is then brought to the portal of a prohead by interaction between terminase and the prohead. Translocation of the DNA into the prohead occurs with the consumption of an estimated 1

ATP per 2 base pairs packaged. When translocation brings the next *cos* site to the portal, it too is cleaved by the terminase, which then releases the completed head and remains bound to the  $\lambda$  left end generated at the second cleavage site. The completed head is matured by the addition of a tail, while the DNA end, bound by terminase, may begin filling another prohead. Two features of packaging are of general interest. First, *cos* is the only specific sequence required in the DNA to allow it to be packaged, a fact that has been exploited in the construction of artificial plasmid vectors that can be packaged into phage particles *in vitro* and *in vivo* (see **Cosmid**). Also, the distance along the DNA between *cos* sites determines the amount of DNA per phage, in contrast to phages that use “headful” packaging (see **P22 phage**). Luckily,  $\lambda$  is somewhat forgiving, readily giving rise to functional particles with as little as 75% or as much as 105% of the normal 48.5 kbp of DNA. This fact makes it possible to replace nonessential  $\lambda$  genes with DNA from other sources without having to make up the amount of replaced DNA precisely (see [Cloning](#)).

Phage  $\lambda$  encodes two proteins that together effect the death and lysis of the infected cell. One of these, encoded by gene *S*, is called “holin,” a protein that spans the inner membrane and is believed to form an aqueous channel to the periplasm. The holin channel permits the escape to the periplasm of endolysin, the gene *R* product, which rapidly digests the peptidoglycan layer, causing lysis. An enduring question about lysis is the timing mechanism that prevents premature lysis. In  $\lambda$ , genes *R* and *S* are transcribed from  $P_R'$ , at the same time as the genes involved in morphogenesis, and the corresponding proteins are made at the same time as other late proteins, yet lysis is delayed until about 45 minutes into infection. It has been suggested that the activity of the *S* protein may be regulated (6).

#### Bibliography

1. F. Sanger, A. R. Coulson, G. F. Hong, D. F. Hill, and G. B. Petersen (1982) Nucleotide sequence of bacteriophage  $\lambda$  DNA. *J. Mol. Biol.* **162**, 729–773.
2. B. C. Hoopes and W. R. McClure (1985) A *cII*-dependent promoter is located within the *Q* gene of bacteriophage  $\lambda$ . *Proc. Natl. Acad. Sci. USA* **82**, 3134–3138.
3. S. Barik, B. Ghosh, W. Whalen, D. Lazinski, and A. Das (1987) An antitermination protein engages the elongating transcription apparatus at a promoter-proximal recognition site. *Cell* **50**, 885–899.
4. W. S. Yarnell and J. W. Roberts (1992) The phage  $\lambda$  gene *Q* transcription antiterminator binds DNA in the late gene promoter as it modifies RNA polymerase. *Cell* **69**, 1181–1189.
5. B. A. Learn, S. J. Um, L. Huang, and R. McMacken (1997) Cryptic single-stranded-DNA binding activities of the phage  $\lambda$  *P* and *Escherichia coli* DnaC replication initiation proteins facilitate the transfer of *E. coli* DnaB helicase onto DNA. *Proc. Natl. Acad. Sci. USA* **94**, 1154–1159.
6. C. Y. Chang, K. Nam, and R. Young (1995) *S* gene expression and the timing of lysis by bacteriophage  $\lambda$ . *J. Bacteriol.* **177**, 3283–3294.

#### Suggestions for Further Reading

7. C. E. Catalano, D. Cue, and M. Feiss (1995) Virus DNA packaging: the strategy used by phage  $\lambda$ . *Mol. Microbiol.* **16**, 1075–1086.
8. R. W. Hendrix, J. W. Roberts, F. W. Stahl, and R. Weisberg (1983) “III”, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
9. A. Polissi, L. Goffin, and C. Georgopoulos (1995) The *Escherichia coli* heat shock response and bacteriophage  $\lambda$  development. *FEMS Microbiol. Rev.* **17**, 159–169.
10. R. Young (1992) Bacteriophage lysis: mechanism and regulation. *Microbiol. Rev.* **56**, 430–481.

## Laminins

Laminins are a family of large ECM glycoproteins whose name is derived from the fact that they are prominent components of the basal lamina. Laminins are heterotrimers consisting of an  $\alpha$ ,  $\beta$ , and  $\gamma$  subunit, and multiple subunit homologs can form at least 15 different laminin isoforms. The best studied member of the family, laminin-1, ( $\alpha$ -1/ $\beta$ -1/ $\gamma$ -1) is an 800,000-kDa cruciform-shaped molecule held together by disulfide bonds. Laminin-1 is widely distributed in most basal laminae, and it displays a number of important biological activities. It supports cell adhesion, migration, differentiation, polarization, and it stimulates the outgrowth of neurites by cultured neurons. In the developing nervous system, laminin-1 immunoreactivity is transiently observed in regions where axons will eventually grow out toward their synaptic partners. Laminin-1 will induce the apical-basolateral polarization of epithelial cells and cause changes in gene expression in mammary epithelial cells leading to their differentiation.

Structure-function relationships in laminin-1 have been determined using strategies similar to those applied to fibronectin. Binding sites for other matrix molecules and for cells have been identified in specific proteolytic fragments of laminin-1. The short arms of the cross, characterized by epidermal growth factor repeats, bind to collagen IV and other laminin molecules, and these interactions are thought to stabilize the basal lamina sheet. An associated 150-kDa molecule known as entactin or nidogen binds to a site near the center of the cross, and because entactin also binds collagen IV, it further solidifies the basal lamina network. Heparin sulfate proteoglycans such as perlecan bind to a globular domain at the bottom of the long arm of the cross. A neurite-promoting site at the interface of the long arm and the globular domain has been identified, which contains the amino acid sequence ile-lys-val-ala-val (IKVAV). This portion of the  $\alpha$ -1 subunit is located in a region of coiled-coil structure just above the C-terminal globular domain at the bottom of the cross. A second cell binding site near the center of the cross, containing the amino acid sequence tyr-ile-gly-ser-arg (YIGSR), appears to be recognized by melanoma cells using a non-integrin 68 kDa receptor. As many as six different integrin receptors have been reported to bind laminin-1. Integrins  $\alpha$ -1/ $\beta$ -1 and  $\alpha$ -2/ $\beta$ -1 bind to fragments from the short arms, and the  $\alpha$ -6/ $\beta$ -1 and  $\alpha$ -7/ $\beta$ -1 bind to the globular domain, near the coiled-coil region. Mouse laminin-1 contains an RGD site, but it is only exposed on proteolysis, when it can be recognized by the integrin  $\alpha$ -v/ $\beta$ -3. Whether such cryptic sites might be functional during wound repair remains to be determined. Other laminin receptors have also been identified. These include dystroglycan, which links laminin and other ECM proteins to the actin cytoskeleton via dystrophin, a large cytosolic protein associated with some forms of muscular dystrophy.

Isoforms of laminin have recently been discovered that appear to have different distributions and functions. Reduced stringency hybridization protocols and novel monoclonal antibodies have been used to identify five  $\alpha$  chains, three  $\beta$  chains, and three  $\gamma$  chains, which can assemble into at least 14 different heterotrimers. More genes may be uncovered in recently sequenced genomes. Laminin-2 ( $\alpha$ -2/ $\beta$ -1/ $\gamma$ -1; also known as merosin) is especially abundant in muscle tissue and seems to have biological activities similar to laminin-1. On the other hand, Laminin-3 ( $\alpha$ -1/ $\beta$ -2/ $\gamma$ -3; also known as S-laminin) and other  $\beta$ -2 containing laminins may affect neuronal axons in a distinct manner. Recent studies indicate that  $\beta$ -2 subunit sequences leu-arg-glu (LRE) may form attachment sites for neural cells that act as molecular stop signs and impede axon outgrowth. Laminin-3 is localized to the synaptic basal lamina, and it may act to halt motor neuron outgrowth and help initiate synapse formation. Consistent with this idea, mutant mice lacking the  $\beta$ -2 subunit fail to make normal neuromuscular synapses. Laminin-5 is localized to anchoring filaments of epithelial cells and is recognized with high affinity by integrins  $\alpha$ -3/ $\beta$ -1 and  $\alpha$ -6/ $\beta$ -4. This laminin has more of a rod shape as viewed by electron microscopy. Different laminin isoforms appear to be recognized by different subsets of integrin and nonintegrin receptors.

Although the functions of all isoforms are not yet known, genetic evidence points to crucial

biological roles. In humans, mutations in the laminin a-2 subunit gene lead to congenital muscular dystrophy, and defective a-3 has been shown to cause epidermolysis bullosa, a blistering skin disorder. In *Drosophila*, mutation of the only known a subunit gene causes many tissue defects and is embryonically lethal.

#### Additional Reading

Jones JC, Dehart GW, Gonzales M, and Goldfinger LE., *Microsc Res. Tech.* **51**(3), 211–213 (2000) and other reviews in the same volume.

Colognato H and Yurchenco PD., *Dev Dyn.* **218**(2), 213–234 (2000).

### Lampbrush Chromosomes

Lampbrush chromosomes are found in the **oocyte** and **spermatocyte** nuclei of many animals. They contain very transcriptionally active **DNA**, where loops of DNA emerging from an apparently continuous chromosomal axis are coated with **RNA polymerase**. Each RNA polymerase molecule is attached to nascent RNA and associated proteins, generating a visible “brush-like” appearance. The axes of lampbrush chromosomes, from which the loops project, consist visually of linear arrays of compacted beads, known as **chromomeres**. DNA is concentrated in the chromomeres, which represent compacted regions of **chromatin**. The axis along which the chromomeres exist consists of two distinct strands of chromatin.

Gall measured the kinetics of lampbrush chromosome breakage by **DNase I sensitivity**. These experiments led to the conclusion that there are two DNA double-helices along the chromomeric axis, but only one in each transcriptionally active loop of the chromosome that emerges from that chromomeric axis. Thus lampbrush chromosomes exist as bivalents. Lampbrush chromosomes are in the diplotene phase of the first division of **meiosis**. Each bivalent consists of two **homologous chromosomes** held together by **chiasmata**. When the chromosomes proceed towards **metaphase**, the loops contract and the metaphasic chromosomes appear.

Distinctive loops can be recognized at invariant positions of the chromosomes and depend on the DNA sequence within each loop. In *Triturus viridescens*, loops are 500  $\mu\text{m}$  wide and up to 1 mm long. Each loop contains several **transcription** units and ranges in size up to 1000 kbp. **Nucleosomes** are clearly seen within active transcription units, especially where RNA polymerases are not so densely packed. An implication of this observation is that nucleosomal structures must be able to reform very rapidly following the passage of RNA polymerase. In *Xenopus* oocytes, proteins, such as **nucleoplasmin**, associate with lampbrush loops and might facilitate nucleosomal reassembly following transit by RNA polymerase (1). Chromomeres occur in long regions of inactive chromatin that are compacted into higher order structures, resembling superbeads (see **Chromomere**). Several groups have prepared **antibodies** against amphibian oocytic nuclear proteins and have used oocytic sections or isolated lampbrush chromosomes for intranuclear localization. The large size of the chromosomes, their ease of manipulation, and the wealth of morphological detail make them ideal for such studies. **Actin**, **histone H2B**, nucleoplasmin and **RNA-binding proteins** have all been localized within loops (2). Actin filaments may be involved in extending the lampbrush chromosomal loop away from the chromomeric axis (3).

#### Bibliography

1. N. Moreau et al. (1986) *J. Cell Biol.* **103**, 683–690.

2. M. B. Roth and J. G. Gall (1987) *J. Cell Biol.* **105**, 1047–1054.
3. D. Rungger, E. Rungger-Brandle, C. Chapponier, and G. Gabbiani (1979) *Nature* **282**, 320–321.

### Suggestion for Further Reading

4. H. G. Callan (1986) *Lampbrush Chromosomes*, Springer Verlag, Berlin.

## Laue Diffraction

A Laue diffraction pattern is produced in [X-ray crystallography](#) when a stationary crystal is illuminated with a continuous spectrum of X-rays. The first X-ray diffraction pictures were taken in this way by Friedrich, Knipping, and Laue in 1912 ([1](#)), but the technique has several disadvantages and was completely superseded by data collection with monochromatic radiation. One problem with Laue diffraction is that a single diffracted beam can be composed of reflections from more than one lattice plane. This multiplet problem can easily be explained by Bragg's law (see [Bragg Angle](#)):

$$2d \sin\theta = \lambda$$

where  $d$  is the lattice plane distance,  $\theta$  the reflection angle, and  $\lambda$  the wavelength. Although  $2 \sin\theta = \frac{\lambda}{d}$ , it is also equal to  $\frac{\lambda/2}{d/2}$  and  $\frac{\lambda/3}{d/3}$ , etc. This multiplet problem is most serious for the low angle reflections.

The availability of synchrotron radiation has given new impetus to the Laue technique. Its broad, smooth spectrum, combined with high intensity, allows extremely short exposure times and the possibility of time-resolved data collection. Moreover, an enormous number of reflections is registered with one exposure. In this way, very rapid reactions can be followed in protein crystals, if a way can be found to synchronize all of the protein molecules that make up the crystal (see [Enzymes](#)).

### Bibliography

1. W. Friedrich, P. Knipping, and M. Laue (1912) *Sitzb. kais. Akad. Wiss.*, München, 303–322.

### Suggestions for Further Reading

2. D. W. J. Cruickshank, J. R. Helliwell, and L. N. Johnson (1992) *Time Resolved Macromolecular Crystallography*, Oxford Science, Oxford.
3. R. M. Sweet, P. T. Singer, and A. Smalås (1993) Considerations in the choice of a wavelength range for white-beam Laue diffraction, *Acta Crystallogr.* **D49**, 305–307.
4. J. Drenth (1995) *Principles of Protein X-ray Crystallography*, Springer, New York, Chap. "12".

## Leading and Lagging Strands

A segment of **DNA** whose replication starts from a [replication origin](#) and proceeds unidirectionally



or bidirectionally to one or two sites of [termination of DNA replication](#) is called a [replicon](#), a unit of [DNA replication](#). In each replicon, replication is continuous from the origin to the terminus and is accompanied by the movement of the replicating point, called the [replication fork](#). Both parental DNA strands are replicated concurrently at the fork. However, replication at a fork is semidiscontinuous: DNA synthesis is continuous on one strand, the **leading strand**, and discontinuous on the other, the *lagging strand* (see [Discontinuous DNA Replication](#)). This occurs because the two chains of double helical DNA are antiparallel, and **DNA polymerase** can extend a DNA chain only in the 5' → 3' direction.

The parent strand that runs 5' → 3' in the reverse direction of fork movement is termed the leading strand, and it serves as a [template](#) for the continuous DNA synthesis, in which the DNA polymerase carries out chain elongation in a highly processive manner. The other parent strand runs 5' → 3' in the direction of fork movement and is termed the lagging strand; it serves as a template for the discontinuous DNA synthesis. Short pieces of DNA, called [Okazaki Fragments](#), are repeatedly synthesized on the lagging strand; these Okazaki fragments are a few thousand nucleotides in **bacterial** cells and a few hundred in **eukaryotic** cells.

In *Escherichia coli*, Pol III holoenzyme is the major replicative DNA polymerase for both leading- and lagging-strand synthesis. The Pol III holoenzyme is a huge multiprotein complex that consists of 10 distinct polypeptide chains ([1](#), [2](#)). This enzyme extends the DNA chain with a high processivity (>500 kb of DNA can be synthesized continuously without the dissociation of polymerase from the template) and high catalytic efficiency (the velocity of chain elongation is 1000 nucleotides per second at 37°C). The catalytic core, composed of three subunits, contains the polymerase activity and a 3' → 5' exonuclease for proofreading ([3](#)). The remaining seven auxiliary subunits enhance the processivity of the core by clamping it onto the template ([4](#)). They also promote the repeated association of the polymerase necessary for discontinuous synthesis of the lagging strand. Structural analysis of the Pol III holoenzyme and studies on a reconstituted replication fork suggest that the holoenzyme is an asymmetric dimer with twin polymerase active sites: One half of the dimer has high processivity and might be the polymerase for continuous synthesis of the leading strand, whereas the other half has the recycling capacity needed for lagging-strand synthesis ([5](#)). Thus, it seems likely that a single molecule of Pol III holoenzyme acts at the replication fork catalyzing concurrently both leading- and lagging-strand synthesis ([6](#), [7](#)).

In eukaryotic DNA replication, the division of labor among the polymerases remains ambiguous. Pol a (Pol I in yeast) is apparently involved in DNA replication, since mutant cells defective in this polymerase activity are inviable. However, Pol a lacks the 3' → 5' exonuclease activity, so its DNA synthesis is inaccurate and shows a low processivity. These enzymatic characteristics make Pol a a poor candidate for the major replicative polymerase. In addition, Pol a is unique in possessing a [primase](#) activity, the only such activity thus far identified in eukaryotic cells, suggesting that Pol a may play a role in the priming of DNA synthesis ([8](#)). On the other hand, Pol d and Pol e each possess a 3' → 5' exonuclease activity and are highly processive polymerases in the presence of [proliferating cell nuclear antigen](#) (PCNA) and replication factor C ([9](#), [10](#)). Their yeast counterparts, Pol III and Pol II, respectively, are essential for cell growth and DNA replication. Therefore, these polymerases are better suited for chromosome replication. However, the specific roles of Pol d and Pol e at the replication fork remain controversial.

There is another difference in the enzymatic processes of synthesizing the leading and lagging strands. Leading-strand DNA synthesis requires RNA primer only once in the replication of each replicon, but a frequent priming process is associated with lagging-strand DNA synthesis. RNA primer must be laid down as the initiation step of each cycle of synthesis of Okazaki fragments. Therefore, the priming proteins, including the primase, are required for lagging-strand DNA synthesis, along with the replication fork.

Bibliography

1. C. S. McHenry (1991) *J. Biol. Chem.* **266**, 19127–19130.
2. C. McHenry and A. Kornberg (1977) *J. Biol. Chem.* **252**, 6478–6484.
3. H. Maki and A. Kornberg (1985) *J. Biol. Chem.* **260**, 12987–12992.
4. S. Maki and A. Kornberg (1988) *J. Biol. Chem.* **263**, 6561–6569.
5. H. Maki, S. Maki, and A. Kornberg (1988) *J. Biol. Chem.* **263**, 6570–6578.
6. C. A. Wu, E. L. Zechner, A. J. Hughes, M. A. Franden, C. S. McHenry, and K. J. Marians (1992) *J. Biol. Chem.* **267**, 4064–4073.
7. S. Kim, H. G. Dallmann, C. S. McHenry, and K. J. Marians (1996) *J. Biol. Chem.* **271**, 21406–21412.
8. R. C. Conaway and I. R. Lehman (1982) *Proc. Natl. Acad. Sci. USA* **79**, 2523–2527.
9. T. Tsurimoto, T. Melendy, and B. Stillman (1990) *Nature* **346**, 534–539.
10. T. Tsurimoto and B. Stillman (1991) *J. Biol. Chem.* **266**, 1950–1960.

### Suggestions for Further Reading

11. J. J. Blow, ed. (1996) *Eukaryotic DNA Replication*, IRL Press, Oxford.
12. M. L. DePamphilis, ed. (1996) *DNA Replication in Eukaryotic Cells*, Cold Spring Harbor Laboratory Press, New York.
13. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.
14. K. J. Marians (1992) Prokaryotic DNA replication. *Ann. Rev. Biochem.* **61**, 673–720.

## Lectins

Lectins (from Latin, *legere*, to select or choose) are proteins that bind mono- and oligosaccharides specifically and reversibly, but are devoid of catalytic activity and, in contrast to antibodies, are not products of an [immune response](#). They are typically di- or polyvalent with respect to carbohydrate binding, and therefore they agglutinate cells and precipitate polysaccharides or glycoproteins; these reactions are inhibited by the sugars for which the lectins are specific. The erythrocyte agglutinating, or *hemagglutinating*, activity of lectins is their major attribute and is used routinely for their detection and characterization; another way is by primary sequence similarities. These proteins are found in most organisms, from **viruses** and **bacteria** to **plants** and animals. Tissues or cells may contain one or more lectins of the same or different specificity. Lectins are readily obtainable in purified form, mostly by [affinity chromatography](#) on the immobilized ligands and also by [recombinant DNA](#) techniques. Hundreds of lectins have been isolated and characterized (for examples, see Table 1), some 30 to the level of their three-dimensional structure in complex with **ligands**. Generally, they are oligomeric proteins of different types with diverse combining sites. Nonetheless, many of them belong to distinct protein families that share sequence homologies and similar **tertiary** and [quaternary structures](#) (see [Protein Evolution](#)). Microbial [toxins](#) that are carbohydrate-specific are also classified as lectins.

**Table 1. Some Well-Characterized Lectins<sup>a</sup>**

---

**Specificity**

| Lectin                        | Source           | Monosaccharide | Oligosaccharide     | Structure               |
|-------------------------------|------------------|----------------|---------------------|-------------------------|
| Conanavalin A                 | Plants (legumes) | Man            | Man a 3(Mana6)      | Simple                  |
| Calnexin                      | Animals          | Man            |                     | Mosaic                  |
| <i>Dolichos biflorus</i>      | Plants (legumes) | GalNAc         |                     | Simple                  |
| Galectins                     | Animals          | Gal            | Galb4GlcNAc         | Simple                  |
| Hepatic binding protein       | Animals          | Gal/GalNAc     |                     | Mosaic                  |
| Influenza virus hemagglutinin | Influenza virus  |                | NeuAc(a2-3 / 6) Gal | Mosaic                  |
| Man-6-P receptor              | Animals          | Man-6-P        |                     | Mosaic                  |
| Peanut agglutinin             | Plants (legumes) | Gal            | GalbGalNAc          | Simple                  |
| P-type fimbriae               | <i>E. coli</i>   |                | Gala4Gal            | Macromolecular assembly |
| Selectins                     | Animals          |                | sLe <sup>x</sup>    | Mosaic                  |
| Sialoadhesins                 | Animals          |                | NeuAc(a2-3 / 6) Gal | Mosaic                  |
| Type 1 fimbriae               | <i>E. coli</i>   | Man            |                     | Macromolecular assembly |
| Wheat germ agglutinin         | Plants (cereals) | GlcNAc, NeuAc  |                     | Simple                  |

<sup>a</sup> *Abbreviations:* Man, mannose; Gal, galactose; GalNAc, *N*-acetylgalactosamine; GlcNAc, *N*-acetylglucosamine; NeuAc, *N*-acetylneuraminic acid; sLe<sup>x</sup>, sialyl-Lewis<sup>x</sup> or NeuAca(2–3)Galb(1–4)[Fuca(1–3)]GlcNAc).

Lectins are invaluable tools for the structural and functional investigation of complex carbohydrates, especially **glycoproteins**, for the study of their biosynthesis and for the examination of changes that occur on cell surfaces during physiological and pathological processes, from cell [differentiation](#) to cancer. At present, they are the focus of intense attention because of the demonstration that they act as recognition determinants in diverse biological processes.

### 1. Carbohydrate specificity

Based on their specificity, lectins are classified into five groups, according to the monosaccharide for which they have the highest affinity: mannose, galactose/*N*-acetylgalactosamine, *N*-acetylglucosamine, fucose, and *N*-acetylneuraminic acid (sugars are of the D-configuration except for fucose, which is L). Relevant for the functions of lectins is the fact that, of the numerous monosaccharides found in nature, these six alone are typical constituents of eukaryotic cell surfaces. Many lectins specific for mannose also react with glucose (the 2-epimer of mannose), but none of these reacts with galactose (the 4-epimer of glucose), nor do those specific for galactose bind

mannose. Similarly, members of the *N*-acetylglucosamine specificity group do not combine with *N*-acetylgalactosamine (or vice versa). However, most lectins that bind galactose interact also with *N*-acetylgalactosamine, in some cases preferentially. For this reason, they are classified in one specificity group, Gal/GalNAc, even though a few (eg, peanut agglutinin) do not bind *N*-acetylgalactosamine at all. Occasionally, lectins combine with apparently unrelated monosaccharides that, however, share common topological features. For instance, wheat germ agglutinin (WGA) binds both *N*-acetylglucosamine and *N*-acetylneuraminic acid, and animal mannose-specific binding proteins (MBPs) bind fucose too. Lectins of the same specificity group may combine preferentially with either the  $\alpha$ - or  $\beta$ -glycosides of the corresponding monosaccharide, whereas others lack anomeric specificity. They may also differ markedly in their affinity for other derivatives, especially oligosaccharides. Certain lectins interact with oligosaccharides only (Table 1). The **association constants** for the binding of monosaccharides and oligosaccharides to lectins are typically in the range  $10^3$  to  $5 \times 10^4 \text{M}^{-1}$  and  $10^5$  to  $10^7 \text{M}^{-1}$ , respectively; multivalent oligosaccharides bind more strongly.

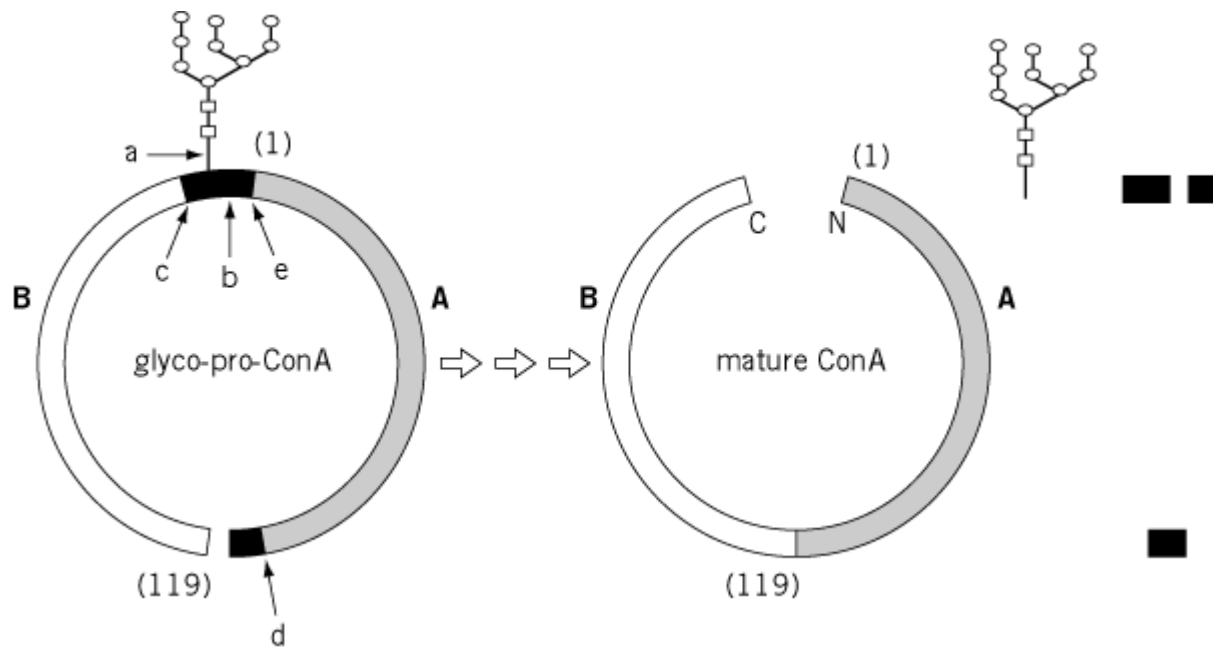
## 2. Molecular Properties

### 2.1. Plant Lectins

The largest and most thoroughly studied family of lectins is that of the legumes, with over 100 members, some 40 of which have been sequenced. These lectins, as well as others from plants, often occur as isolectins that either are coded by closely-related genes or are the product of [post-translational modification](#) such as **glycosylation** or **proteolysis**. Concanavalin A (Con A) from Jack bean, the prototype plant lectin, was first isolated in 1919 by James Sumner (of urease fame) and shown by him, in 1936, to be specific for mannose and glucose. Members of this family consist of two or four identical, or nearly identical, subunits of 25 to 30 kDa that may carry up to two *N*-linked oligosaccharides. In some cases, the subunits are fragmented. Each subunit has a single carbohydrate combining site with the same specificity. One exception is phytohemagglutinin, which occurs as a family of five tetrameric **isoforms** in all possible combinations of E and L subunits (from  $E_4$  to  $L_4$ ) that differ in their specificity and biological properties. Legume lectins also contain one tightly bound  $\text{Ca}^{2+}$  ion and one transition metal ion (usually  $\text{Mn}^{2+}$ ) per subunit, which are required for carbohydrate binding. Their sequences are about 40% **homologous**; invariant residues include several of those that participate in [hydrogen bonding](#) (an [aspartic acid](#) and an [asparagine](#) residue) and in **hydrophobic** interactions (an aromatic amino acid or [leucine](#)) with the ligand, and almost all those that coordinate the metal ions. Two animal lectins (MR60/ERGIC-53 specific for mannose, and VIP36 specific for *N*-acetylgalactosamine) are homologous to those of the legumes (1).

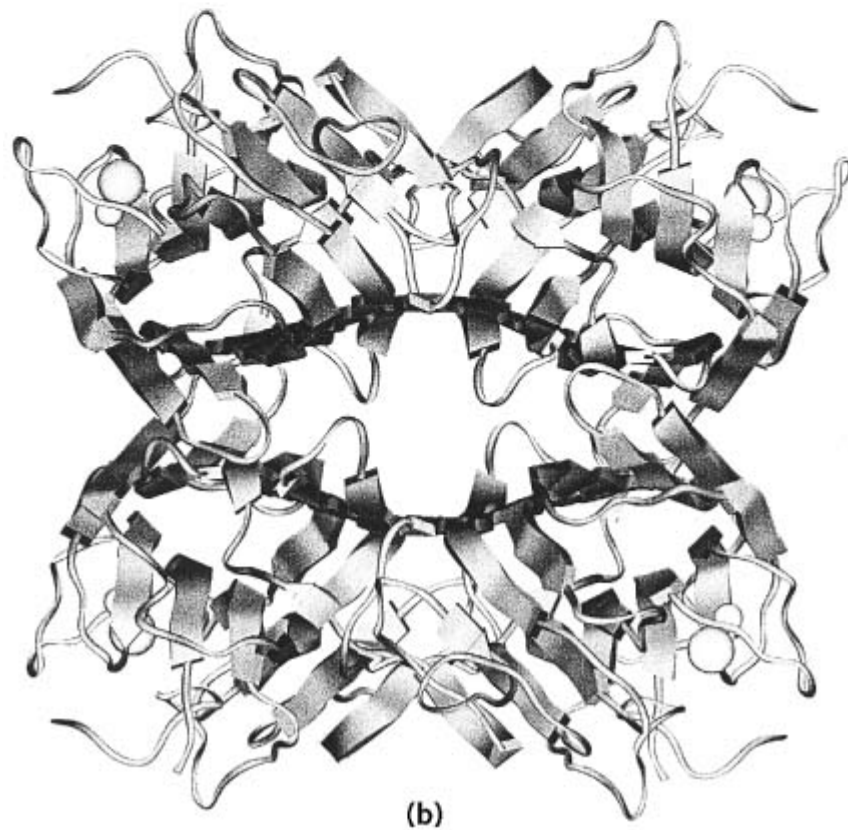
Concanavalin A occupies a special position, because it exhibits “circular homology” with the other legume lectins. This homology is obtained by aligning residue 123 of concanavalin A with the amino-terminal residue of the other lectins, proceeding to the carboxyl end of concanavalin A and continuing along its amino-terminal region. It is the result of an unusual post-translational processing of the lectin (2) (Fig. 1).

**Figure 1.** Post-translational modifications during the biosynthesis of concanavalin A; a summary of processing events converting glycosylated pro-concanavalin A to mature lectin. The amino- and carboxy-termini are indicated by N and C, respectively, and the numbers in parentheses are residue numbers in mature concanavalin A. During processing in the plant, the inactive glycosylated prolectin is deglycosylated (arrow a), resulting in appearance of lectin activity. An endopeptidase then cleaves (arrows b, c, d, and e) a carboxy-terminal nonapeptide and the glycosylated spacer (shown as solid black areas). Residues 118 (arrow d) and 119 are ligated enzymatically. Splicing thus results in a transposition of the linear arrangement of the protein sections designated B and A. (From Ref. 2, with permission.)



High-resolution [X-ray crystallography](#) analysis of the structures of legume lectins revealed that the subunits are dome-shaped, made up largely of two antiparallel [b-sheets](#) that form a **jellyroll** or *lectin* fold (3). The subunits are nearly superimposable, irrespective of the specificity of the lectins, and associate into different types of dimers or tetramers (Fig. 2). The carbohydrate-binding site is located at the top of the subunit dome, in close proximity to the metal ions, which help to position the amino acid residues that form contacts with the carbohydrate. The invariant binding site residues (aspartic acid, asparagine, and an aromatic amino acid) occupy identical spatial positions in all these lectins, and discrimination between glucose/mannose and galactose is achieved by the different orientation of the respective monosaccharides in the binding sites (4). Thus, concanavalin A binds glucose and mannose such that the Od1 and Od2 of aspartic acid are hydrogen-bonded with the 6-OH and 4-OH, respectively, and the Nd2 of asparagine is bonded with the 4-OH of the sugar. On the other hand, in lectins (such as peanut agglutinin, ECorL and SBA) that bind galactose, Od1 and Od2 of aspartic acid form hydrogen bonds with the 4-OH and 3-OH, respectively, and the Nd2 of asparagine forms a hydrogen bond with the 3-OH.

**Figure 2.** Three-dimensional model of concanavalin A. (a) Subunit dimer obtained by the antiparallel side-by-side alignment of b-sheets, leading to the formation of a contiguous 12-stranded b-sheet that extends across the dimer interface. (b) Subunit tetramer formed by the association of the central parts of both dimers; the dimer contacts are mainly through loop interactions. Prepared using program MS I/BIOSYM ING, San Diego, CA. PDB entry 5CNA. (Courtesy of Dr. Miriam Eisenstein.)



Cereal lectins, a prominent example of which is WGA, consist too of two identical subunits, but they differ markedly from those of the legumes. For instance, they are exceptionally rich in [cysteine](#) residues, which are rare in legume lectins, are devoid of metals, and possess multiple binding sites

(5). Although WGA is not glycosylated, its precursor is a glycoprotein. [X-ray crystallography](#) of the sugar complexes of WGA show that the combining site of this lectin contains several [tyrosine](#) residues that interact hydrophobically with the bound GlcNAc or NeuAc and also form hydrogen bonds with the ligand.

## 2.2. Animal Lectins

The galectins are a family of soluble animal lectins that bind exclusively b-galactosides, such as lactose and *N*-acetyllactosamine. They are found inside the **cytoplasm** and [nucleus](#) of cells, and occasionally also on the cell surface and outside the cell, and their expression is developmentally regulated. They are synthesized without a leader sequence and possess a relatively simple structure, occurring as monomers and homodimers of subunits with molecular weights of about 14 kDa, as well as larger polypeptides (30 to 35 kDa). Each galectin contains one or two copies of a homologous domain, known as the *S*-carbohydrate recognition domain (S-CRD). The tertiary structure of the galectins also exhibits the jellyroll topology found in the legume lectins, despite the absence of significant sequence homology and a different location of the combining site.

Numerous animal lectins are of the C-type ( $\text{Ca}^{2+}$ -dependent). They are **mosaic** (or multi-**domain**) molecules, characterized by an extracellular carbohydrate recognition domain (C-CRD) consisting of 115 to 130 amino acid residues, of which 14 are invariant and 18 are highly conserved (6). To the CRD is attached a variable number of domains of different kinds, which form the bulk of the molecule, and also a membrane-spanning domain. Lectins of this class are grouped into three families: *selectins*, *collectins* and *endocytic lectins*. Only three selectins are known: E-selectin, P-selectin, and L-selectin. The first two are specific for the sialylated Lewis<sup>x</sup> blood group determinant, NeuAca(2–3)Galb(1–4)[Fuca(1–3)]GlcNAc, abbreviated sLe<sup>x</sup>, and its positional isomer, sialylated Lewis<sup>a</sup> blood group determinant, NeuAca(2–3)Galb(1–3)[Fuca(1–4)]GlcNAc, or sLe<sup>a</sup>, with both L-fucose and sialic acid (or another negatively-charged group such as sulfate) required for binding. L-Selectin binds sialylated, fucosylated, and sulfated oligosaccharides on diverse mucin-like glycoproteins. Recognition of the carbohydrates is possible only when they are present on particular glycoproteins, such as cell surface mucins, pointing to the importance of the carrier molecule and carbohydrate presentation in the interaction with the lectins.

The collectins are soluble proteins; they include the mannose-binding proteins MBP-A and C, the structural unit of which is a trimer of 32-kDa subunits based on a triple helix formed by the **collagen**-like portion of the molecule. MBP-A circulates in serum of, for example, rodents and humans as a hexamer of the trimeric units. As shown by X-ray crystallography, MBP-A and MBP-C bind mannose via  $\text{Ca}^{2+}$  that serves as the nucleus of the combining site and interacts with the 3-OH and 4-OH of the ligand. Four of the five additional bonds that coordinate the metal ion are provided by the side chains of two [glutamic acid](#) and two **asparagine residues** that also are hydrogen bonded to the same (3 and 4) mannose hydroxyls. The four amino acids just mentioned are conserved in all C-type lectins specific for mannose, two of them in the sequence Glu–Pro–Asn (positions 185 to 187 in MBP-A). When Glu185 and Asn187 in MBP-A were replaced by [glutamine](#) and [aspartic acid](#), respectively, as found in galactose specific C-type lectins, galactose became the preferred ligand (7).

A prominent representative of the endocytic lectins is the rabbit hepatic asialoglycoprotein receptor [or hepatic binding protein (HBP)], the first mammalian lectin to be described. It is found on hepatocytes of different mammals and is specific for galactose and *N*-acetylgalactosamine, whereas its avian homologue is specific for *N*-acetylglucosamine. Other endocytic lectins are a fucose- and galactose-specific receptor found on Kupffer cells and the mannose-specific lectin of macrophages and hepatic endothelial cells. Except for the latter, these lectins are type II **transmembrane proteins**, consisting of a short amino-terminal cytoplasmic domain, a **hydrophobic**, membrane-spanning domain, a neck region, and a carboxy-terminal CRD. The mannose-specific macrophage surface lectin differs from the other endocytic lectins primarily in that it is a type I transmembrane protein and that its extracellular part contains a domain, closest to the membrane, with eight CRDs.

The sialoadhesins (I-type lectins) are a family of sialic acid-specific type I membrane glycoproteins with variable numbers of extracellular [immunoglobulin](#) (Ig)-like domains and are thus members of the **immunoglobulin superfamily** (8). They include the sheep erythrocyte receptor for macrophages (referred to simply as *sialoadhesin*), the lymphocyte surface antigen CD22 found only on B cells, CD33 present on early myeloid cells, and MAG, a glycoprotein associated with myelin. CD22 recognizes specifically NeuAc(a2-6)Gal(b1-4)GlcNAc. In contrast, all other known I-type lectins bind structures containing *N*-acetylneuraminic acid that is a2-3-linked. The P-type CRD has been found only in two closely related lectins, the **mannose-6-phosphate** (Man-6-P) **receptors**.

### 2.3. Microbial Lectins

Several viruses (eg, [influenza virus](#) and **polyoma virus**) contain lectins specific for *N*-acetylneuraminic acid. Many bacterial species express surface lectins, usually in the form of fimbriae (or [pili](#)). These filamentous, heteropolymeric organelles, a few nanometers in diameter and 100 to 200 nm in length, consist of helically arranged subunits (pilins) of several different types. Only one of the subunits, usually a minor component of the fimbriae, possesses a carbohydrate-binding site—for example, for mannose (in type 1 fimbriae) or galabiose, Gala(1-4)Gal (in P fimbriae).

### 3. Functions

Participation of lectins in cell recognition (Table 2) was first demonstrated in the 1940s for the influenza virus hemagglutinin and in the 1970s for bacterial surface lectins. These lectins mediate the binding of the pathogens to host cells, a step essential for the initiation of infection. Inhibitors of bacterial lectins protect animals against experimental infection by the lectin-carrying organisms, providing a basis for the development of antiadhesion therapy of microbial infections (9). Some bacterial surface lectins allow the specific binding of the bacteria to human polymorphonuclear cells and human and mouse macrophages in the absence of opsonins, which may lead to activation of the phagocytes and ingestion and killing of the bacteria (**lectinophagocytosis**).

**Table 2. Functions of Lectins**

| Lectin                      | Role in   |
|-----------------------------|---|
| <b>Microorganisms</b>       |   |
| Influenza virus             | Infection   |
| Amoeba                      | Infection   |
| Bacteria                    | Infection   |
| <b>Animals</b>              |   |
| Calnexin                    | Glycoprotein synthesis                                  |
| Galectins                   | Embryogenesis, metastasis (?)                           |
| <b>C-type lectins</b>       |   |
| Mannose-binding protein     | Host antimicrobial defense                              |
| L-selectin                  | Lymphocyte homing                                       |
| E-selectin                  | Leukocyte trafficking to sites of inflammation          |
| P-selectin                  | Leukocyte trafficking to sites of inflammation          |
| I-type lectins              | Cell-cell interactions in the immune and neural systems |
| Man-6-P receptors           | Targeting of lysosomal enzymes                          |
| Natural killer cell lectins | Cytolysis of target cells                               |



---

The galectins are postulated to function in cell adhesion. When present on the surface of metastatic murine and human cancer cells, they may be responsible for the adhesion of the cells to target organs, a step necessary for metastasis. Exposing highly metastatic cells to compounds containing lactose before injecting them into mice reduced the metastatic spread almost by half. Therefore, antiadhesive drugs may turn out to be antimetastatic.

The endocytic lectins have been assumed to facilitate clearance from the circulation of glycoproteins with complex oligosaccharide units (eg, ceruloplasmin and a  $\beta$ -acid glycoprotein) from which the terminal sialic acid has been removed, exposing the subterminal galactose. It is uncertain, however, whether this represents a physiological mechanism for regulating the turnover of serum glycoproteins (and cells), because disruption of the receptor does not result in decreased levels of desialylated forms of predominant circulating glycoproteins (10).

The **mannose-6-phosphate receptors** are responsible for targeting the appropriate enzymes to the **lysosome** subcellular compartment. The recently discovered intracellular lectins **calnexin**, MR60/ERGIC-53, and VIP-36 participate in the biosynthesis of glycoproteins, as well as in their intracellular sorting. The mannose-specific macrophage lectin has been implicated in innate antimicrobial defense (11). It binds infectious organisms that expose mannose-containing glycans on their surface, leading to their killing by lectinophagocytosis. The MBPs of mammalian serum and liver combine with oligomannosides of yeast and fungi, causing activation of **complement** in an **antibody**- and C1q-independent manner and subsequent lysis of the pathogens. They also enhance phagocytosis of the invading organisms by acting as opsinins. A mutation of a single amino acid residue in the collagen-like domain of the lectin is associated with recurrent, severe bacterial infections in infants.

The selectins provide the best paradigm for the role of sugar–lectin interactions in biological recognition (12). They mediate the initial adhesion of circulating leukocytes to endothelial cells of blood vessels, a prerequisite for the exit of the former cells from the circulation and their migration into tissues. L-Selectin is found on all leukocytes and is involved in the recirculation of lymphocytes, directing them specifically to peripheral lymph nodes. The two other selectins are expressed on endothelial cells, and only when these cells are activated by inflammatory mediators (eg, histamine, **interleukin-2**, and **tumor necrosis factor**) released from tissue cells in response to, for example, infection or ischemia. Individuals unable to synthesize the selectin ligands sLe<sup>x</sup> and sLe<sup>a</sup> suffer from recurrent bacterial infections. Prevention of adverse inflammatory reactions by inhibition of leukocyte–endothelium interactions, another application of antiadhesion therapy, has become a major aim of many pharmacological industries. As shown in animal models, oligosaccharides recognized by the selectins protect against experimentally induced lung injury and tissue damage caused by myocardial ischemia and reperfusion. In addition to their involvement in inflammation, selectins may play a role in the spread of cancer cells from the primary tumor throughout the body.

### 3.1. Plant Lectins

The role of plant lectins is still not well understood. Plant lectins may function in the establishment of **symbiosis** between nitrogen-fixing bacteria, mainly rhizobia, and leguminous plants and in defense of plants against phytopathogenic **fungi** and predatory animals.

## 4. Applications

Native lectins are used predominantly for applications that are based on precipitation and agglutination reactions or for **mitogenic** stimulation. Lectins derivatized with **fluorescent** dyes, gold particles, or **enzymes** are employed in histo- and cytochemistry. Immobilized lectins, for instance on

Sepharose, are indispensable for the isolation by affinity chromatography of glycoproteins, glycopeptides, and oligosaccharides. Mouse and human cortical (immature) thymocytes can readily be separated from the medullar (mature) ones with the aid of peanut agglutinin, making it possible to examine *in vitro* their developmental and functional relationships. SBA is used clinically for the removal, from bone marrow of haploidentical donors, of the mature T cells responsible for the lethal graft versus host reaction, affording a fraction enriched in stem cells. This fraction is employed routinely for transplantation into children born with severe combined immune deficiency (“bubble children”) with close to 70% success, and experimentally for **transplantation** into leukemic patients. Another clinical application of lectins is in blood typing—for example, to identify blood type O cells and to differentiate between M and N cells.

Lectins are potent mitogens; they are polyclonal activators stimulating lymphocytes irrespective of their antigenic specificity. Phytohemagglutinin in particular serves clinically to assess the immunocompetence of patients—for example, under chemotherapy or in those with AIDS. It is also employed for the preparation of chromosome maps for **karyotyping**, sex determination, and detection of **chromosome** defects, because the chromosomes are easily visualized in the stimulated cells. Highly toxic lectins (such as **ricin**) and moderately toxic ones (such as concanavalin A, phytohemagglutinin, and WGA) serve for the selection of lectin-resistant cell mutants that are widely employed in studies of the genetics, biosynthesis, and functions of complex carbohydrates.

### Bibliography

1. C. Itin, A.-C. Roche, M. Monsigny, and H.-P. Hauri (1996) *Mol. Biol. Cell* **7**, 483–493.
2. D. H. Jones (1995) In *Perspectives on Protein Engineering & Complementary Technologies* (M. J. Geisow and R. Epton, eds.), Mayflower Worldwide Ltd., Birmingham, U.K., pp. 70–73.
3. N. Srinivasan, S. D. Rufino, M. B. Pepys, and S. Wood (1996) *Chemtracts Biochem. Mol. Biol.* **6**, 149–164.
4. N. Sharon (1993) *Trends Biochem. Sci.* **18**, 221–226.
5. C. S. Wright and G. E. Kellogg (1996) *Protein Sci.* **5**, 1466–1476.
6. K. Drickamer and M. E. Taylor (1993) *Annu. Rev. Cell Biol.* **9**, 237–264.
7. W. I. Weis, K. Drickamer, and W. A. Hendrickson (1992) *Nature* **360**, 127–134.
8. S. Kelm, R. Schauer, and P. R. Crocker (1996) *Glycoconj. J.* **13**, 913–926.
9. D. Zopf and S. Roth (1996) *Lancet* **347**, 1017–1021.
10. S. Ishibashi, R. E. Hammer, and J. Herz J, . (1994) *J. Biol. Chem.* **269**, 27803–27806.
11. J. Epstein, Q. Eichbaum, S. Sheriff, and R. A. B. Ezekowitz, (1996) *Curr. Opin. Immunol.* **8**, 29–35.
12. L. A. Lasky (1995) Selectin-carbohydrate interactions and the initiation of the inflammatory response. *Annu. Rev. Biochem.* **64**, 113–139.

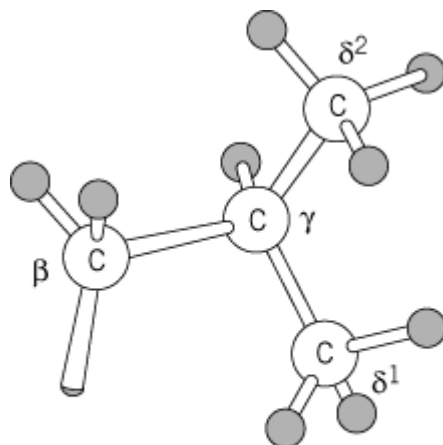
### Suggestions for Further Reading

13. K. Kasai and J. J. Hirabayashi (1996) Galectins: a family of animal lectins that decipher glycocodes. *J. Biochem.* **119**, 1–8.
14. H. Lis and N. Sharon (1997) Lectins: carbohydrate specific proteins that mediate cell recognition. *Chem. Rev.* **98**, 637–674.
15. R. M. Nelson, A. Venot, M. P. Bevilacqua, R. J. Linhardt, and I. Stamenkovic (1995) Carbohydrate–protein interactions in vascular biology. *Annu. Rev. Cell Dev. Biol.* **11**, 601–631.
16. J. M. Rini (1995) Lectin structure. *Annu. Rev. Biophys. Biomol. Struct.* **24**, 551–557.
17. N. Sharon and H. Lis (1989) *Lectins*, Chapman and Hall, London, 1989.
18. N. Sharon and H. Lis (1993) Carbohydrates in cell recognition. *Sci. Am.* **268**(1), 82–89.
19. W. Weis and K. Drickamer (1996) Structural basis of lectin–carbohydrate recognition. *Annu. Rev. Biochem.* **65**, 441–473.

## Leucine (Leu, L)

The [amino acid](#) leucine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to six **codons**—UUA, UUG, CUU, CUA, CUC, and CUG—and represents approximately 9.0% of the residues of the proteins that have been characterized. The leucyl residue incorporated has a mass of 113.16 Da, a **van der Waals volume** of  $124 \text{ \AA}^3$ , and an [accessible surface](#) area of  $180 \text{ \AA}^2$ . Leu residues are infrequently changed during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [valine](#), [isoleucine](#), [methionine](#), and [phenylalanine](#) residues.

The Leu side chain is [nonpolar](#) with no functional or reactive groups:



Leu is one of the most **hydrophobic** amino acid residues, and 45% of the residues in native [protein structures](#) are fully buried. Leu is one of the residues that favors most the **alpha-helical** conformation in model peptides. It also occurs frequently in  $\alpha$ -helices in folded proteins, somewhat more frequently than in [b-sheet](#) type of **secondary structure**.

Leucine residues are very frequently involved in interactions in **coiled coils**, and they have given their name to the so-called [leucine zipper](#).

### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Leucine Zippers

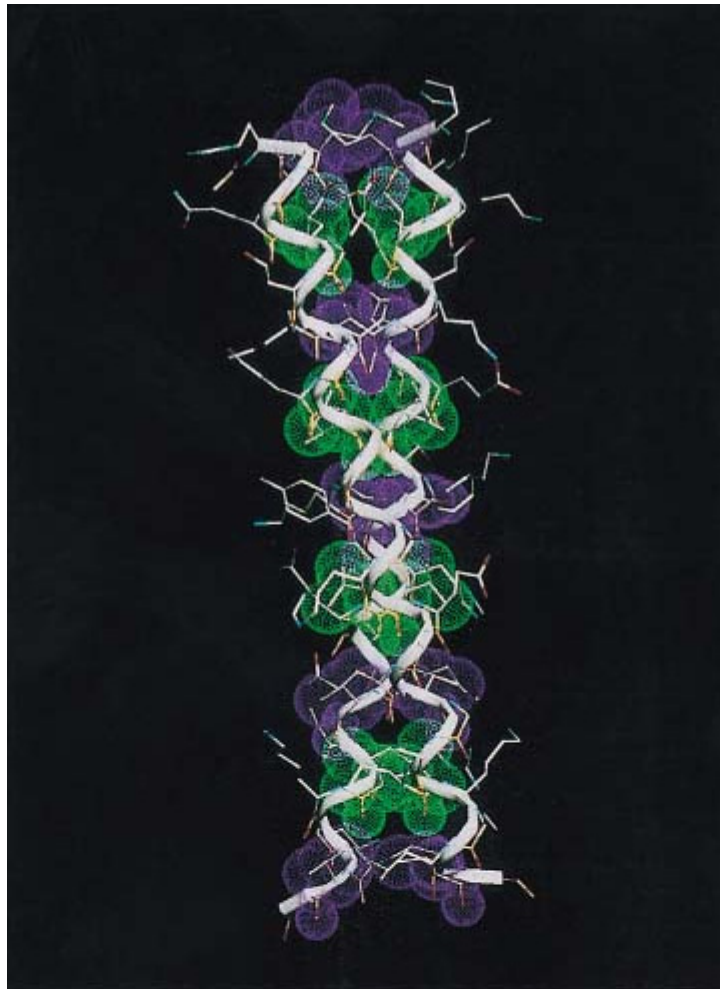
The structure of a 33-residue region of the yeast [transcription factor](#) GCN4 represents the first example of a [coiled-coil](#) conformation solved at atomic resolution ([1](#)). Like all coiled-coil structures based on the  [\$\alpha\$ -helix](#), this region has an amino acid sequence with a [heptad repeat](#) substructure of the form  $(a-b-c-d-e-f-g)_n$ , where [nonpolar](#) residues occur commonly in positions  $a$  and  $d$ . In the case of the GCN4 sequence, however, the  $d$  positions were occupied entirely by [leucine](#) residues (except for the C-terminal residue; see [Table 1](#)), hence the origin of the term *leucine zipper*. The primary function of the leucine zipper is to facilitate the dimerization of the bZIP transcription factors.

**Table 1. Heptad Distribution of Residues in the GCN4 Leucine Zipper**

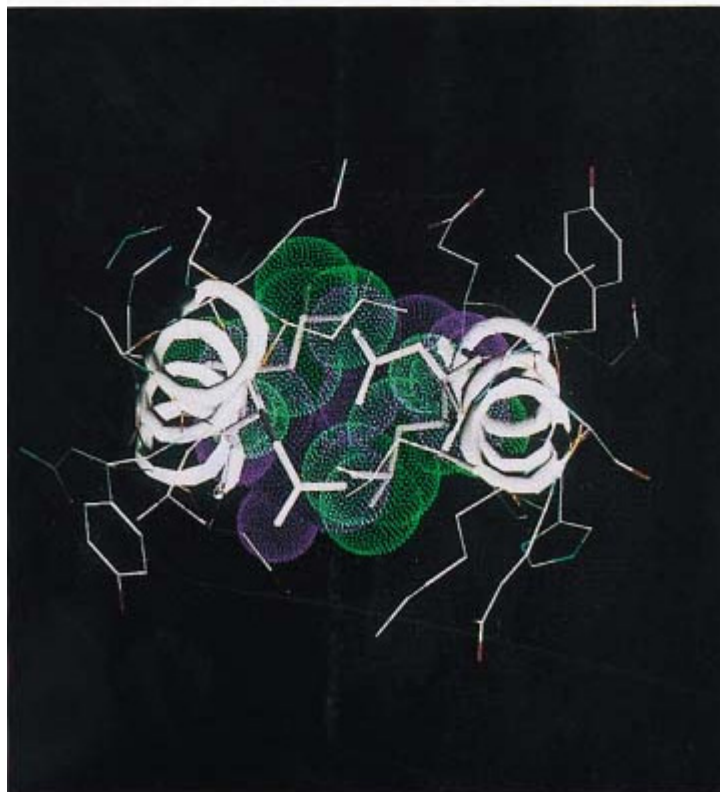
| Heptad Position | $a$ | $b$ | $c$ | $d$ | $e$ | $f$ | $g$ |
|-----------------|-----|-----|-----|-----|-----|-----|-----|
|                 |     |     |     |     |     |     | Arg |
|                 | Met | Lys | Gln | Leu | Glu | Asp | Lys |
|                 | Val | Glu | Glu | Leu | Leu | Ser | Lys |
|                 | Asn | Tyr | His | Leu | Glu | Asn | Glu |
|                 | Val | Ala | Arg | Leu | Lys | Lys | Leu |
|                 | Val | Gly | Glu | Arg |     |     |     |

The crystal structure has allowed us the opportunity of studying the detailed interactions made by various side chains in stabilizing and specifying the coiled-coil conformation ([Fig. 1](#)) and, in particular, the conformations of the side chains of the nonpolar residues that lie along the axis of the coiled-coil. The observations clearly indicated that positions  $a$  and  $d$  were structurally unique, because their apolar side chains were not equivalently directed toward the interface between the two chains. In fact, the leucine residues in position  $d$  pointed directly at the interface, whereas the valine residues that dominate at position  $a$  pointed outward from the interface. This provides a simple explanation for the observed preference at position  $a$  of b-branched apolar residues, such as valine and isoleucine, where part of their side chains can be redirected back into the interface region and away from the aqueous environment.

**Figure 1.** Structure of the GCN4 leucine zipper. **(a)** Molecular model viewed in a direction perpendicular to the long axis of the structure. Van der Waals surfaces for residues in the  $a$  and  $d$  positions are colored purple and green, respectively. **(b)** A central portion of the structure viewed in transverse section. Leucine residues packed at one level in the core of the two-stranded coiled-coil are highlighted. (Courtesy of C. L. Day and T. A. Alber, PDB coordinate reference number 2ZTA.) See color insert.



(a)



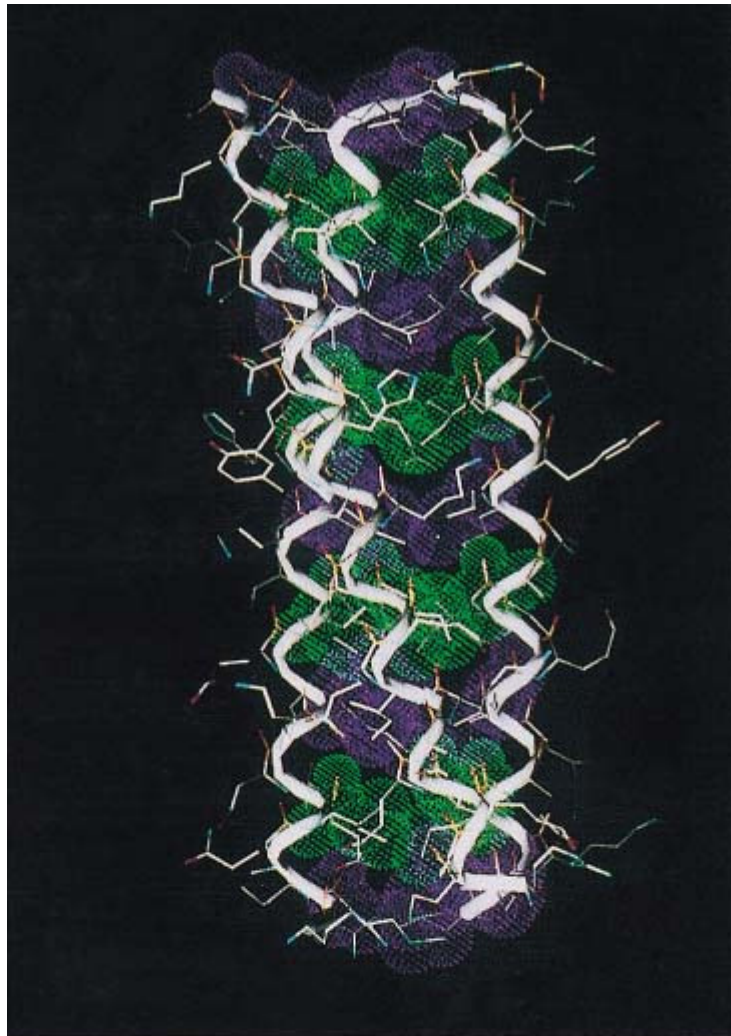
(b)

[Protein engineering](#) studies of the basic GCN4 structure have produced a number of variants:

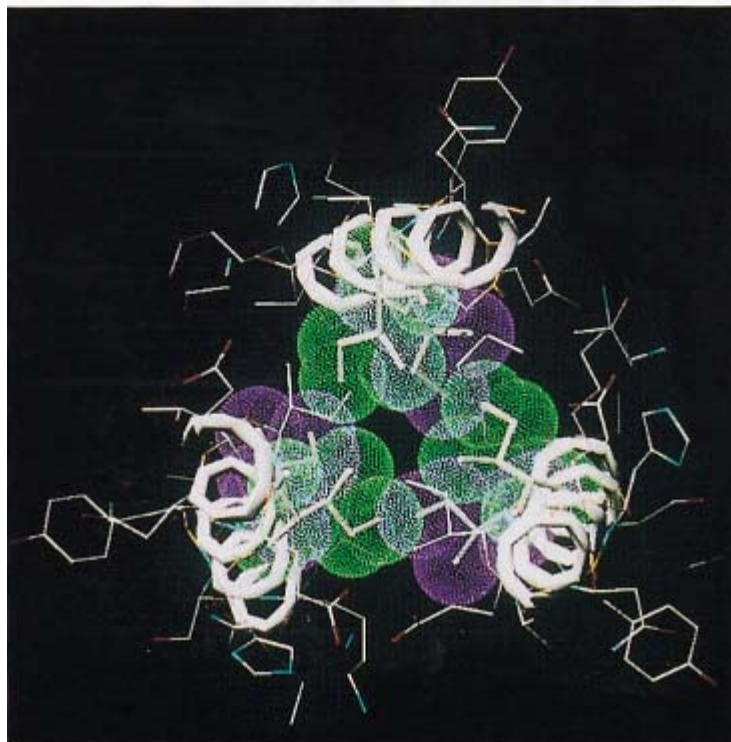
Isoleucine residues at *a* and leucines at *d* (1, 2): this gives a dimeric structure with a pitch length of 14.8 nm, a supercoil radius of 0.49 nm, and three out of a possible six interchain [electrostatic interactions](#) between oppositely charged residues in positions *e* and *g*.

Isoleucine residues at both *a* and *d* (3); this gives a trimeric structure with a pitch length of 17.5 nm, a supercoil radius of 0.67 nm, and eight out of a possible nine interchain ionic interactions between oppositely charged residues in positions *e* and *g* (Fig. 2).

**Figure 2.** Three-stranded coiled-coil conformation adopted by the GCN4-pII peptide in which residues in positions *a* and *d* of the GCN4 leucine zipper sequence have been changed to isoleucine (3). (a) The conformation seen in a direction perpendicular to the long axis of the structure. The van der Waals surfaces are colored purple for residues in position *a* and green for residues in position *d*. (b) Packing of the isoleucine residues along the axis of the coiled-coil is shown in a section of the structure. (Courtesy of C. L. Day and T. A. Alber; PDB coordinate reference number 1GCM.) See color insert.



(a)



(b)

Leucine residues at *a* and isoleucines at *d* (2); this gives a tetrameric structure with a pitch length of 20.5 nm, a supercoil radius of 0.76 nm, and five of a possible 12 interchain ionic interactions between oppositely charged residues in positions *e* and *g*. The inference from these studies is that *b*-branched residues such as valine and isoleucine tend to destabilize specific modes of assembly.

Harbury et al. (2) described “parallel” and “perpendicular” modes of packing for the apolar residues in positions *a* and *d*, where these terms refer to the relative directions of the C<sub>a</sub>–C<sub>b</sub> bonds of the residues in different chains that pack at the same axial position. In two-stranded coiled-coils, the *a* residues show parallel packing and the *d* residues display perpendicular packing. These preferences are reversed for four-stranded coiled-coils. Three-stranded coiled-coils, on the other hand, display an intermediate form of packing of the apolar residues. This is termed “acute.” Leucine zipper stability and specificity for dimerization also lies with interchain ionic interactions, particularly those occurring between residues in positions *e* and *g* (4).

Another interesting mutant of the GCN4 leucine zipper is one in which the sole asparagine residue at position *a* was replaced by a valine (5). It might have been expected that the coiled-coil would have been “improved,” because the content of apolar residues in position *a* was increased and the presence of an additional valine at *a* would have further stabilized the structure. In practice, however, a mixture of dimeric and trimeric structures was adopted. It would seem clear that the presence of a destabilizing residue such as asparagine in position *a* has the result of destabilizing the dimeric structure less than the trimeric one.

#### Bibliography

1. E. K. O'Shea, J. D. Klemm, P. S. Kim, and T. Alber (1991) X-ray structure of the GCN4 leucine zipper, a two-stranded, parallel coiled coil. *Science* **254**, 539–544.
2. P. B. Harbury, T. Zhang, P. S. Kim, and T. Alber (1993) A switch between two-, three, and four-stranded coiled coils in GCN4 leucine zipper mutants. *Science* **262**, 1401–1407.
3. P. B. Harbury, P. S. Kim, and T. Alber (1994) Crystal structure of an isoleucine-zipper trimer. *Nature* **371**, 80–83.
4. D. Krylov, I. Mikhailenko, and C. Vinson (1994) A thermodynamic scale for leucine zipper stability and dimerization specificity: *e* and *g* interhelical interactions. *EMBO J.* **13**, 2849–2861.
5. S. A. Potekhin, V. N. Medvedkin, I. A. Kashparov, and S. Y. Venyaminov (1994) Synthesis and Properties of the Peptide Corresponding to the Mutant Form of the leucine zipper of the transcriptional activator GCN4 from yeast. *Protein Eng.* **7**, 1097–1101.

#### Leucine-Rich Repeat

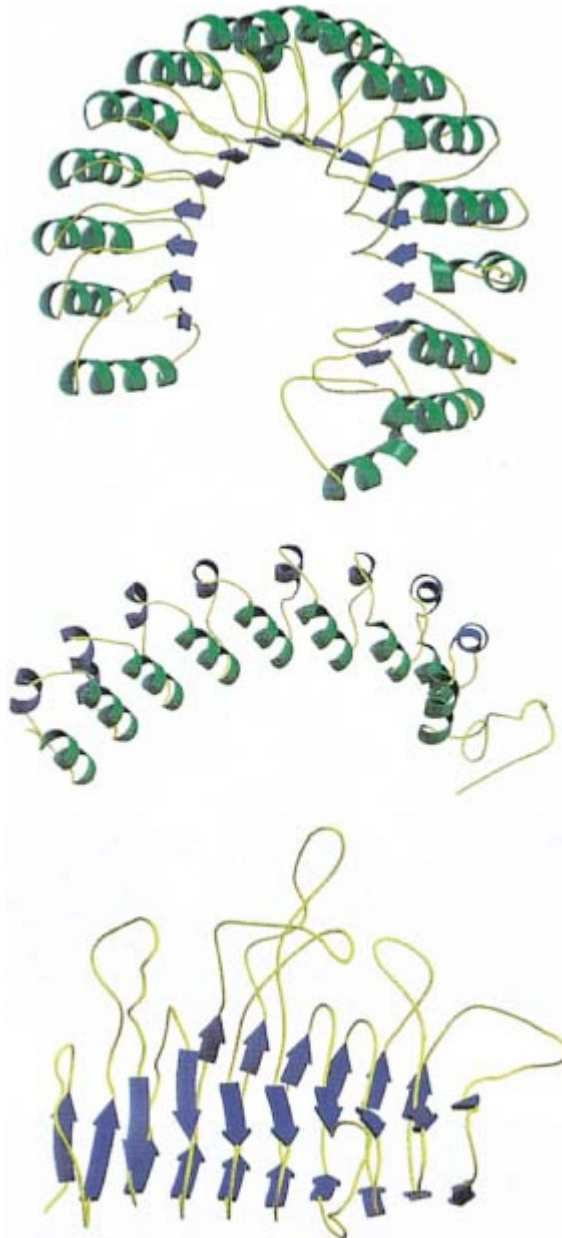
Leucine-rich repeat proteins, or LRR proteins, are a diverse group of [proteins](#) characterized by distinctive repeating sequences in their [primary structure](#) that include several [leucine](#) residues or other aliphatic residues (1, 2). For example, a typical repeating module in a LRR has the sequence, XLXXLXXNX<sub>a</sub>XX<sub>a</sub>XXXX<sub>a</sub>XXLX, where L is leucine, X is any residue, and <sub>a</sub> is an aliphatic residue. The number of residues in the repeat varies from 20 to 29, and the number of repeats in a protein varies from 1 to 41. LRR proteins are associated with functions such as [signal transduction](#),



cell adhesion, and [protein–protein interactions](#).

The structures of two LRR proteins, **ribonuclease inhibitor** (2) and leucine-rich variant (LRV) protein (3), have been determined, and they both exhibit repeating three-dimensional structural [protein motifs](#) that form a supercoil (Fig. 1). In the case of ribonuclease inhibitor, each LRR unit forms a short [b-strand](#) followed by an [a-helix](#), arranged roughly parallel to one another. Together, the repeat units generate a super helix, with the b-strands forming a parallel [b-sheet](#) that lines the interior of a horseshoe-shaped structure. In the LRV protein (3), there is no b-sheet structure; instead the repeating units form alternating a-helix and  $3_{10}$ -helix motifs arranged in a right-handed superhelix of a comma-shaped structure. In this case, it is the parallel a-helices that line the inner surface of the protein. In both ribonuclease inhibitor and LRV protein, however, the conserved leucine residues form the core of the protein.

**Figure 1.** Schematic representation of the backbone structures of two LRR proteins and of a b-helix protein. **(Top)** Ribonuclease inhibitor (2) with b-strands shown as blue arrows and a-helices shown as green coils. **(Middle)** LRV protein (3) with a-helices shown as green coils and  $3_{10}$ -helices shown as blue coils. **(Bottom)** The b-helix of the pectate lyases (4), with b-strands shown as blue arrows. In all cases, the connecting loop regions are shown in yellow. This figure was generated by Molscript (6) and Raster3D (7, 8). See color insert.



It has also been proposed that other repeating substructures could be used in LRR architecture. For example, the repeating three b-strands of pectate lyase that form a [b-helix](#) structure (4) have been suggested as a model for some LRRs (2). The molecular recognition and protein interaction roles shared by LRRs are thought to be a result of the nonglobular nature of these proteins (2) and their conformational flexibility (5).

[See also [Beta-Helix](#).]

#### Bibliography

1. N. Takahashi, Y. Takahashi, and F. Putnam (1985) *Proc. Natl. Acad. Sci. USA* **82**, 1906–1910.
2. B. Kobe and J. Deisenhofer (1994) *TIBS* **19**, 415–421.
3. J. W. Peters, M. H. B. Stowell, and D. C. Rees (1996) *Nature Struct. Biol.* **3**, 991–994.
4. M. D. Yoder and F. Journak (1995) *FASEB J.* **9**, 335–342.
5. B. Kobe and J. Deisenhofer (1995) *Nature* **374**, 183–186.
6. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.

7. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.  
 8. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

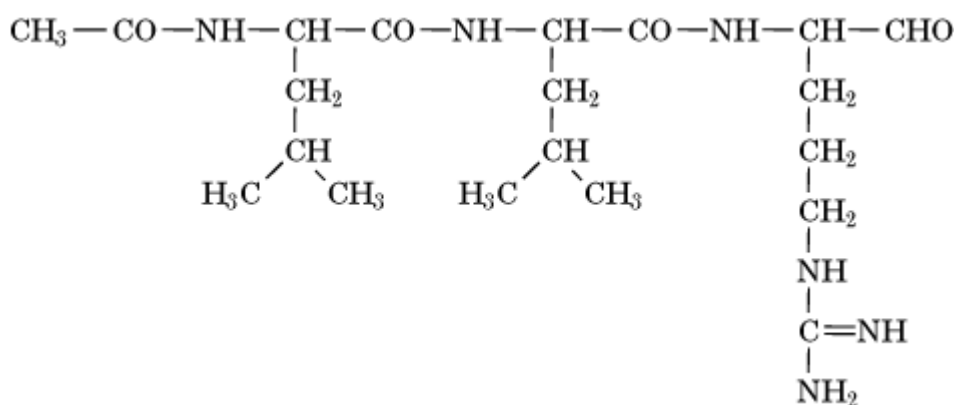
### Suggestion for Further Reading

9. S. G. St. C. Buchanan and N. J. Gay (1996) Structural and functional diversity in the leucine-rich repeat family of proteins. *Prog. Biophys. Mol. Biol.* **65**, 1–44.

## Leupeptin

Screening of actinomycetes culture filtrates for inhibition of the **proteolytic** enzyme **trypsin** led to the discovery of leupeptin (1). It is a group of water-soluble peptide analogues, and the one most commonly available is *N*-acetyl-L-leucyl-L-leucyl-L-argininal, where arginal is the α-carboxaldehyde derivative of arginine (Fig. 1). It is a potent inhibitor of not only trypsin, but also of other [proteinases](#), such as **plasmin**, **papain**, and **cathepsin B** (2). However, it is not an effective inhibitor of **chymotrypsin**, **elastase**, **thermolysin**, **pepsin**, or cathepsins A and D. Consequently, it is not useful for classifying proteinases. Specificity seems to be toward enzymes (such as trypsin) that have a strong preference for basic residues, although it also inhibits some but not all [thiol proteinases](#).

**Figure 1.** Chemical structure of leupeptin, an inhibitor of trypsin-like serine and, in some cases, thiol proteinases. It is also known as *N*-acetyl-leucyl-leucyl-argininal. The aldehyde reacts with the active-site serine hydroxyl (or cysteine thiol) group, and the guanidino side chain contributes the trypsin-like specificity.



The aldehyde group is essential for inhibition. If it is oxidized or reduced, inhibition is abolished. In solution, the aldehyde either is hydrated or forms a cyclic carbinolamine, depending on conditions. Leupeptin, like many other proteinase inhibitors, is often used to prevent unwanted degradation of proteins during the course of isolation. A 1 mM concentration of leupeptin (0.5 mg/mL) is usually recommended for this purpose. Because of its limited specificity, it should be used in conjunction with other, broad-specificity inhibitors for optimal protection. Moreover, it can be cleaved by some nonsusceptible proteinases, such as thermolysin, which further limits its effectiveness. One of the products of thermolysin cleavage can be immobilized on an insoluble polymer, which can then be employed for [affinity chromatography](#) of trypsin-family proteinases—for example, to remove trypsin from chymotrypsin samples or trypsin from trypsinogen (3).

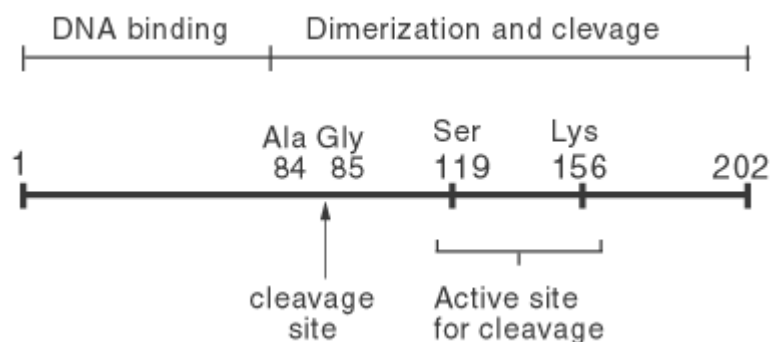
## Bibliography

1. T. Aoyagi et al. (1969) *J. Antibiot.* **22**, 558 ff.
2. H. Umezawa (1976) *Methods Enzymol.* **45**, 678 ff.
3. S. Ishii and K. Kasai (1981) *Methods Enzymol.* **80**, 842–848.

## LExA Repressor

LexA repressor controls the [SOS response](#) of *Escherichia coli*. LexA is a protein similar to  $\lambda$  repressor (see **Lambda repressor**). It has two structural and functional **domains** (Fig. 1), an N-terminal **DNA-binding** domain and a C-terminal domain that is involved in formation of dimers and in a specific **proteolytic** reaction that cuts the protein in half and inactivates its repressor function. This proteolytic event is unusual in that it is a self-cleavage reaction. In the cell, however, this reaction requires the presence of an activated form of another protein, **RecA**. Activated RecA acts indirectly to stimulate the latent self-cleavage activity of LexA. Hence this activity of RecA is termed its coprotease function. The DNA-binding domain of LexA has also been widely used in analyzing eukaryotic gene regulation by workers who fuse it to the protein of interest, thereby targeting the chimeric protein to a LexA binding site.

**Figure 1.** Organization of the LexA polypeptide chain. LexA protein is organized into two structural and functional domains separated by a hinge or connector region. The functions of the two domains are shown at the top. Specific cleavage of LexA at the site indicated by the arrow separates the two domains and inactivates the repressor function of the protein. The two residues believed to be important within the catalytic site are indicated in the second domain. Reproduced from Little (1993).



### 0.1. DNA Binding

LexA resembles  $\lambda$  repressor in its organization and function. Like  $\lambda$  repressor, LexA has a [helix-turn-helix motif](#) (1) with which it binds DNA. LexA binds to specific sites, or **operators**, near the **promoters** of about twenty SOS genes and reduces their rates of expression. The SOS gene products are involved in responding to DNA damage (see [SOS Response](#)). Unlike many other repressors, LexA does not completely repress these genes. Instead, they are expressed at low levels, typically 1 to 20% of the derepressed rate. It is thought that these low levels suffice to repair sporadic DNA damage. The binding sites for LexA (or SOS boxes) have dyad symmetry, and its **consensus sequence** is CTGTN<sub>8</sub>ACAG, where N is any base. Various SOS boxes diverge from this consensus

to one extent or another, and LexA binds with differing affinities to different sites. The relative positions of the SOS box and the promoter differ in the various regulatory regions.

The DNA-bound form of LexA is a dimer. LexA dimerizes weakly in solution, but it can dimerize on the DNA rather than in solution (2). The binding to two half-sites occurs with a high degree of **cooperativity**. The second monomer binds about  $10^6$ -fold more tightly than the first. As a result of this cooperativity, the **binding** curve is steep. Because the *in vivo* concentration of LexA is below the dimerization constant measured *in vitro*, occupancy of the operators probably changes *in vivo* over a relatively small range of LexA concentrations, giving efficient SOS induction. Cooperative DNA binding between  $\lambda$  repressor dimers has a similar role in the process of prophage induction (see **Lambda repressor**, [SOS Response](#)).

## 0.2. LexA Cleavage

This proteolytic reaction is of interest for two reasons. First, it plays a crucial role in the SOS response. Cleavage inactivates LexA by separating its DNA-binding domain from the domain required for dimerization. The isolated DNA-binding domain cannot dimerize efficiently, and hence it binds weakly.

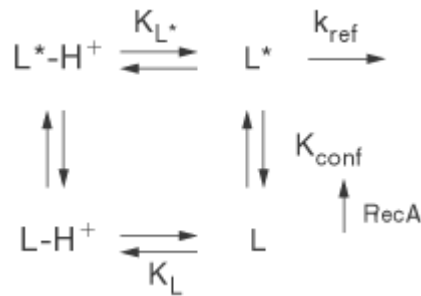
Secondly, cleavage is a self-processing reaction. Under the appropriate *in vitro* conditions, LexA cuts itself in half in an intramolecular reaction termed autodigestion. Cleavage cuts a specific bond (Ala84-Gly85, see Fig. 1) in the hinge region separating the two domains. Phage  $\lambda$  repressor undergoes this same reaction. Autodigestion is entirely parallel to an enzymatic reaction, and the C-terminal domains of LexA and  $\lambda$  repressors can act as true proteolytic **enzymes**, cleaving other substrate molecules (3). Enzymes contain **active sites**, composed of a catalytic center that carries out the chemistry of the reaction and a binding pocket that specifically binds the substrate and positions it optimally relative to the catalytic center. In LexA, the catalytic center is probably composed of a conserved serine residue, Ser119 in LexA, which acts as a nucleophile to attack the peptide bond. It is thought that this Ser is activated by a deprotonated lysine residue, Lys156 in LexA.

This Ser/Lys pair contrasts with the Ser/His/Asp **catalytic triad** in most **serine proteinases**. A similar Ser/Lys pair is present in **b-lactamase** (4). Autodigestion of LexA is stimulated *in vitro* at high pH, which is believed to titrate Lys156. The corresponding Ser and Lys residues in another cleavable protein, UmuD', lie within hydrogen-bonding distance (5) and at the end of a shallow surface groove that probably represents the binding pocket. Several mutations that alter the rate of cleavage in LexA lie in residues flanking this groove.

In the cell, LexA autodigestion is undetectable, and specific cleavage requires interaction with another regulatory protein, RecA. It is believed that RecA does not play a direct role in the chemistry of cleavage. Rather, RecA stimulates the latent self-cleavage activity of LexA. Although RecA acts as a catalyst, it acts indirectly as a coprotease. In this view, LexA cleavage is designed to be slow, but is capable of large rate increases on the order of 10,000-fold. An analogy is the **GTPase** activity of **Ras**, which is stimulated greatly by GTPase activating proteins (see [Gtpases](#)).

Analysis of mutant forms of LexA that have increased rates of self-cleavage (6) has led to a conformational model (Fig. 2) for cleavage. In this view, LexA exists in two states, a low-energy state, L in the diagram, and a higher energy state (L\*) that is competent for cleavage. Although the  $pK_a$  for Lys156 in the L form is about 10, it is reduced to  $<7$  in the L\* form. It is believed that activated RecA stabilizes the L\* form (by raising the value of  $K_{conf}$ ), leading to rapid cleavage even at physiological conditions. There is no direct physical evidence for conformational changes, however, and it is unclear whether conversion of L to L\* is binding of the cleavage site to the active site or a subtle realignment of a preformed complex.

**Figure 2.** Conformational model for LexA cleavage and the role of RecA. In this model, LexA is in two conformational states, L and L\*, coupled to the ionization of one or more groups. The ratio of these two states in the nonprotonated form of LexA is given by the equilibrium constant  $K_{\text{conf}}$ . The values of  $K_L$  and  $K_{L^*}$ ; are roughly  $10^{10}$  and  $10^5 - 10^6$ , that is, the  $\text{p}K_a$  of the titrating group is about 10 in L and 5 to 6 in L\*. Because ionization and the conformational equilibrium are **linked functions**, the conformational equilibrium in the protonated form of the protein is  $10^4$  to  $10^5$  less favorable for the L\* conformation than in the nonprotonated, high pH form. Activated RecA is believed to increase  $K_{\text{conf}}$  to a high value, perhaps by binding tightly to L\* and/or L\*-H and thereby pulling the conformation to the L\* conformation even at neutral pH. Reproduced with permission from (6).



Self-processing reactions are becoming increasingly recognized in biology. Other examples include RNA self-splicing (see **Ribozyme**), intein excision (see **Splicing**), and autophosphorylation and autophosphatase activities of sensor-receiver pairs in two-component regulatory systems (see **Phosphorylation**).

### 0.3. Use of LexA in Analysis of Eukaryotic Transcription

LexA is widely used to study the function of proteins in eukaryotic **gene regulation**. The LexA DNA-binding domain, or more commonly the entire protein, is fused to all or a portion of a eukaryotic regulatory protein (using [recombinant DNA](#) methods). Then, this chimeric protein binds to one or more LexA binding sites, located near a eukaryotic promoter. The same reporter constructs can be used with a wide variety of chimeric proteins. Although other DNA-binding proteins could probably be used for this purpose, LexA was the first to be developed in this way and continues to enjoy wide application.

Studies with such chimeras gave the first indications that eukaryotic regulatory proteins have a modular organization (7). The DNA-binding activity of GAL4 was separable from its activation function. Since then, countless studies with many other proteins have led to similar conclusions and have allowed mapping functional domains within [transcription factors](#). In some respects, however, this approach reinforces an often misleading view of protein structure and function as a linear collection of modules, and it should be applied and interpreted with caution.

### Bibliography

1. R. H. Fogh, G. Oettleben, H. Rüterjans, M. Schnarr, R. Boelens, and R. Kaptein (1994) *EMBO J.* **13**, 3936–3944.
2. B. Kim and J. W. Little (1992) *Science* **255**, 203–206.
3. B. Kim and J. W. Little (1993) *Cell* **73**, 1165–1173.
4. N. C. J. Strynadka et al. (1992) *Nature* **359**, 700–705.
5. T. S. Peat, E. G. Frank, J. P. McDonald, A. S. Levine, R. Woodgate, and W. A. Hendrickson (1996) *Nature* **380**, 727–730.
6. K. L. Roland, M. H. Smith, J. A. Rupley, and J. W. Little (1992) *J. Mol. Biol.* **228**, 395–408.
7. R. Brent and M. Ptashne (1985) *Cell* **43**, 729–736.

## Suggestions for Further Reading

8. J. W. Little (1993) LexA cleavage and other self-processing reactions, *J. Bacteriol.* **175**, 4943–4950.
9. M. Schnarr, P. Oertel-Buchheit, M. Kazmaier, and M. Granger-Schnarr (1991) DNA binding properties of the LexA repressor, *Biochimie* **73**, 423–431.

## Libraries

A molecular library is a collection of diverse molecules that can be screened for individual species exhibiting desirable properties. Molecular libraries provide a means to sample chemical space—a theoretical construct containing every possible compositional and geometric molecular variant. Molecules that are biologically relevant tend to cluster in relatively well-defined regions of chemical space, such as the space occupied by all polynucleotides or **polypeptides**. However, chemical space is by definition infinite in size; therefore, any given library can sample only a minuscule portion.

In addition to providing useful compounds, library screening provides information regarding the possible chemical solutions to a given molecular problem. For example, nucleic acid libraries may be screened for species that inhibit a **protein** target of interest, even though that protein may not normally bind nucleic acids (1). Regions of chemical space that are not representative of the chemistry of life can also be screened to derive compounds that are useful in many applications including, but not limited to, drug discovery, material sciences, and development of novel catalysts.

Libraries are usually developed for a particular purpose, such as representing all of the **genomic** DNA sequences of an organism. Thus, all members of a library are typically chemically related, with diversity operating at the level of differences in the sequences of a set of building blocks, such as nucleotides and amino acids. Polymeric libraries, assembled by stringing together monomeric building blocks, may be comprised of any type of polymer including DNA, RNA, polypeptides, carbohydrates, and lipids. In addition, nonpolymeric small molecule libraries have been developed in which the diversity operates at the level of substituents around a common scaffold, such as a heterocyclic ring.

Molecular libraries fall into two broad classes based on the source of diversity. In the first class, the library represents existing genetic diversity obtained from an organism or group of organisms. The majority of these genetic libraries sample either [messenger RNA](#) or genomic DNA sequences and store the genetic diversity as DNA copies. Depending on the context in which the DNA is presented, the stored sequences can be sampled either directly as DNA or indirectly after expression as RNA or protein species. In the second class of molecular libraries, diversity is synthesized artificially. Chemical synthesis enables the exploration of nonbiological chemical compounds, as well as nucleic acid and protein sequences not found in nature. Modern combinatorial approaches rely heavily on chemically synthesized libraries for efficient sampling of relatively short sequences, in order to derive optimal ligands for a given molecular target. In contrast, genetic libraries are heavily utilized in molecular biology and genomic projects in which the goal is to understand the biology of a given **gene**, gene product, or complex phenotype.

The number of different species within a molecular library is potentially vast. For example, [cDNA libraries](#) routinely contain multiple copies of essentially all of the  $10^4$  to  $10^5$  species of mRNA expressed within a given organism or tissue. [Phage display libraries](#) typically contain  $10^8$  to  $10^9$

different species. Oligonucleotide-based libraries can exceed  $10^{15}$  different species (2). However, libraries that contain such large numbers of different species necessarily contain small amounts of any individual species, typically far below the analytical threshold of most instruments, thereby creating enormous challenges in recovering and identifying the tiny minority of molecules that are active in a given assay. There are essentially two approaches to this detection problem. First, bead-based libraries [eg, one-bead, one-compound libraries (OBOC); (3)] can be designed to provide sufficient compound or an associated code (4) to enable detection of active molecules by highly sensitive readouts, such as [mass spectrometry](#) or [protein sequencing](#). Alternatively, the selected molecules can be amplified so that they can be detected by less sensitive assays. In practice, genetic or enzymatic amplification is currently possible only for DNA- or RNA-based libraries.

Combinatorial libraries typically contain large numbers of distinct compounds, yet with sufficient representation of each compound to allow occurrence and detection of those rare interactions that pass the selection criteria. In addition, the greater the diversity of the library (the range of chemical structures represented in the set of building blocks), the greater the probability that one or more desirable molecules will be present. For example, DNA and RNA libraries possess less chemical diversity than protein libraries because nucleic acids have fewer building blocks. However, library complexity—the number of distinct species present—is independent of the chemical composition of the library. The optimization of complexity, representation, and diversity presents numerous, often competing, challenges in the design of combinatorial libraries.

Like a library of books, a molecular library is most useful if it is addressable in some fashion, so that desirable entities within the collection can be readily identified. It is simple to address a DNA- or RNA-based library, because each species carries an inherent polynucleotide code that may be read by powerful yet routine molecular biology methods. Thus, even single copies of a DNA-based library member may be multiplied, detected, and sequenced using a combination of hybridization, primer extension, and amplification technologies. Peptide libraries based on phage display may also be addressed through the bacteriophage's DNA sequences, because the peptide moiety remains physically linked to its DNA genome within the virion. Libraries not based on nucleic acid encoding may be addressable through either spatial methods (eg, a positional array of synthetic compounds) or by encoding methods (4).

Screening is the process by which molecules with desirable attributes are isolated from a library. Nucleic acid libraries are typically screened by **hybridization** to a specific probe (5), by affinity partitioning (6-8), or by covalent modification (9). Hybridization is by far the most common approach to screening genomic DNA and cDNA libraries. For example, partial sequencing of a purified protein may yield sufficient information to design an oligonucleotide capable of annealing to the gene that encodes the target protein. The oligonucleotide can be labeled with radioactivity or a fluorophore and then used to hybridize to bacterial colonies or viral plaques formed by the library vector and host. Those rare bacteria exhibiting a signal upon hybridization to the probe are distinguished from the rest of the library and are readily isolated. In contrast, phage-display and oligonucleotide-based aptamer libraries express molecules that may adopt a complicated three-dimensional structure capable of specific binding to certain targets. The term “aptamer” refers to the property that library members are “apt” to interact with targets through their tertiary structure rather than by one-dimensional hybridization (7). Aptamer libraries are typically screened through affinity partitioning, in which binders are retained by an immobilized target, while nonbinders are removed through successive washing steps. There are also examples of oligonucleotide-based libraries that have yielded **ribozymes** capable of catalyzing the ligation of two polynucleotide chains (9). The active fraction within these ribozyme libraries is readily identified by taking advantage of the fact that ligation of the ribozyme to a unique substrate generates a unique primer binding site. PCR amplification using a 5'-ribozyme primer and a substrate-specific primer can partition any active ribozymes from the remainder of the library.

In practice, library screening encounters a fierce competition between the specific signal generated



by the few active species within the library versus the weaker, but far more prevalent, nonspecific signal contributed by the background of inactive species. In almost all types of screening, the enrichment of specific species over nonspecific background is insufficient to permit efficient isolation of the desired species in a single cycle of selection. Typically several—and often many—cycles of binding, partitioning, and reamplification are required to accomplish the desired selection. For example, if the starting library contains  $1 \times 10^9$  different species, and the per cycle enrichment is only 1000-fold, then at least three cycles of selection will be needed to force the emergence of the specific species above the background. However, pushing the selection cycle for too many rounds may be counterproductive for two reasons. First, the resultant screen may yield only a single species, thereby eliminating potentially valuable structure–activity relationship (SAR) information that would have been obtained through characterization of less optimal species from earlier rounds of selection. Second, increasing the number of selection cycles containing an amplification step often results in the emergence of clones with optimal growth and survival characteristics that may or may not interact favorably with the target. These concerns may not apply to genomic or cDNA libraries, because hybridization screening of nucleic acid libraries is usually intended to identify single clones. However, plaque- or colony-purification at each round can reduce the problem of fast growing clones.

Once active molecules are detected, they must be resynthesized in order to derive sufficient quantities for downstream studies. Amplification of DNA-based libraries is readily accomplished *in vitro* by enzymatic means or *in vivo* by introduction into an appropriate host cell. However, amplification represents a significant obstacle in the utilization of most non-DNA-based libraries. In screening small-molecule libraries, the post-screen resynthesis of the molecules found to be active against a given target can be time-consuming and expensive. Such economic liabilities restrict the use of small-molecule libraries to those with access to substantial analytical and synthetic chemistry resources. In contrast, DNA-based libraries are well within the means of most scientists.

See also [Combinatorial Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), [Peptide Libraries](#), and [Phage Display Libraries](#).

## Bibliography

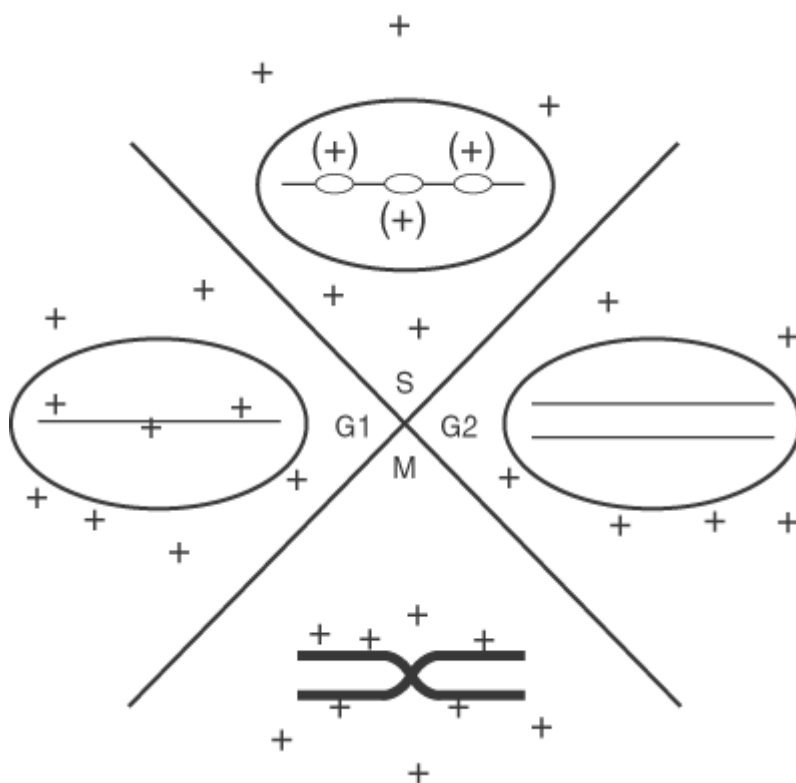
1. L. C. Bock, L. C. Griffin, J. A. Latham, E. H. Vermaas, and J. J. Toole (1992) *Nature* **355**, 564–556.
2. D. J. Kenan, D. E. Tsai, and J. D. Keene (1994) *Trends Biochem. Sci.* **19**, 57–64.
3. K. S. Lam, M. Lebl, and V. Krchnak (1997) *Chem. Rev.* **97**, 411–448.
4. S. Brenner and R. A. Lerner (1992) *Proc Natl Acad Sci USA* **89**, 5381–5383.
5. R. M. Lawn, E. F. Fritsch, R. C. Parker, G. Blake, and T. Maniatis (1978) *Cell* **15**, 1157–1174.
6. C. Tuerk and L. Gold (1990) *Science* **249**, 505–510.
7. A. D. Ellington and J. W. Szostak (1990) *Nature* **346**, 818–822.
8. D. E. Tsai, D. S. Harper, and J. D. Keene (1991) *Nucleic Acids Res.* **19**, 4931–4936.
9. D. P. Bartel and J. W. Szostak (1993) *Science* **261**, 1411–1418.

## Licensing Factor

A famous **cell fusion** experiment performed by Rao and Johnson in 1970 (1) demonstrated that

**eukaryotic** cells have a mechanism to restrict their **genome** to replicate only once per **cell cycle**. All genomic DNA in cells must replicate once in a limited period of the cell cycle, the **S phase**, but never replicate again until the cell cycle reaches the next S phase. A further study by Blow and Laskey in 1988 (2) reproduced the regulation system *in vitro* using *Xenopus* egg extracts and proposed a mechanism to license all the genomic DNA during mitosis to replicate only once in the next S phase (Fig. 1). This model predicted the presence of a hypothetical determinant, called the **licensing factor**, for the fate of replicated and unreplicated DNA. Four properties were assumed for the factor. (1) It is not able to pass through the nuclear **membrane**. (2) It associates with the **chromatin** in **M phase** when the nuclear membrane has disappeared and stays there to activate initiation of **DNA replication**. (3) The molecule binds to chromatin from the M through G1 phase, but the other population unbound to chromatin remains outside of the nucleus. (4) After DNA replication, it is quickly disrupted or inactivated. These features guarantee that the genome will replicate only once during S phase and also provide a mechanism whereby chromatin becomes competent for the next DNA replication in M phase.

**Figure 1.** Model for licensing of DNA replication.



A candidate protein that behaves like the hypothetical factor was obtained in the early 1990s (3) from the **yeast** *Saccharomyces cerevisiae* as minichromosome maintenance (MCM) gene products. The MCM genes were isolated as mutants that affect the maintenance of yeast **minichromosomes** and are essential for the initiation of chromosomal DNA replication and cell cycle progression. Studies of the localization of a subset of the MCM proteins (MCM 2, 3, and 5) revealed that they bind to the chromatin from the M through G1 phase and disappear from the nucleus according to the initiation of DNA replication, as expected for the putative licensing factor. In addition, candidates for functional licensing factor, which were composed of several polypeptide chains of around 100 kDa, were identified from the *Xenopus* egg extract as proteins that are specifically associated with the licensed chromatin (4). The genes coding for these polypeptides were highly homologous to

yeast MCM genes, and the protein complex was essential for the initiation of DNA replication in a *Xenopus* cell-free system. Another experiment, in which the licensing activity was directly purified from *Xenopus* egg extract, detected two essential components, one of which was the MCM complex (5). These results strongly suggest that the MCM complex will be the real licensing factor.

Proteins homologous to yeast MCM have also been identified in animals and **plants**, and it is thought that the proteins exist widely among all eukaryotes. In general, eukaryotes have six MCM proteins, called MCM2 to MCM7, in a hexameric complex. Here, the complex of the potential licensing factors is called MCM or the MCM complex, although there exist other types of MCM genes in yeast. Their primary structures are closely related to each other, constituting a **gene family**. The most conserved region among them exists at the middle of the proteins and has the [ATPase](#) motif, including the **DEAD box**-like sequence frequently detected in [RNA Helicases](#). Recently, it has been reported that a potential helicase activity exists in the purified MCM complex (6).

The most prominent property of MCM is its periodic association with chromatin, but no DNA-binding activities have yet been observed with isolated MCM. This probably means that MCM functions by interacting with other factors. In fact, isolation of the functional licensing factor (5) revealed that the licensing activity requires two components, one of which is MCM and the other an unidentified factor necessary for the binding of MCM to chromatin. One important property assumed for the licensing factor is its inability to pass through the nuclear membrane. However, mammalian MCM is transported into nuclei by its *nuclear localization signal* sequences (see [Nuclear Import, Export](#)) and exists in nuclei during most of the cell cycle. This means that there should be an additional factor impermeable to the nuclear membrane and involved in the chromatin binding of MCM. Thus, some of the roles of a “licensing factor” are played in collaboration with the other factor(s). Genetic studies with yeast MCM showed that they have several functional interactions with subunits of the [origin recognition complex](#) (ORC) thought to be the initiator protein of yeast chromosomal DNA replication. In this case, the role of activation to initiate DNA replication will be played by the ORC protein. Thus, MCM does not behave as the licensing factor solely but does so in association with other components. Therefore, the recent concept of the licensing factor in a precise sense represents the activities of a set of proteins. Of course, MCM has the central role in the widely defined licensing activity and can be called a licensing factor according to the narrow definition of the term.

## Bibliography

1. P. N. Rao and R. T. Johnson (1970) *Nature* **225**, 159–164.
2. J. J. Blow and R. A. Laskey (1988) *Nature* **332**, 546–548.
3. B.-K. Tye (1994) *Trends Cell Biol.* **4**, 160–166.
4. Y. Kubota et al. (1995) *Cell* **81**, 601–609.
5. M. A. Madine et al. (1995) *Nature* **375**, 421–424.
6. Y. Ishimi (1997) *J. Biol. Chem.* **272**, 24508–24513.

## Lifson–Roig Model

The Lifson–Roig (1) model for the cooperative [random coil](#) to  [\$\alpha\$ -helix](#) transition in peptides (see [Helix–Coil Theory](#)) is a statistical mechanical model similar to the [Zimm–Bragg model](#) (2), but with some important differences. Chief among these differences is the definition of a helical unit. In the Lifson–Roig treatment, the basic unit is an amino acid residue, centered on the C <sub>$\alpha$</sub>  atom and the

attached side chain, including the adjacent NH and CO moieties. A residue is defined as being either helical,  $h$ , or coil,  $c$ , based on its  $(\phi, \psi)$  angles (see [Ramachandran Plot](#)). If a residue lies in helical  $(\phi, \psi)$  space, it is defined as being helical, and all other conformations are considered nonhelical or coil. As with the Zimm–Bragg treatment, the partition function for the Lifson–Roig model is constructed with three parameters:  $n$ , the number of residues in the chain;  $w$ , the helix propagation parameter (similar to the  $s$  value in the Zimm–Bragg model); and  $v$ , a nucleation parameter for each end of the helical stretch of residues (the Zimm–Bragg  $s$  value would be akin to  $v^2$ ).

The full Lifson–Roig treatment requires a  $3 \times 3$ -correlation matrix between adjacent residues to assign the appropriate statistical weights to the partition function. Because the definition of the basic unit in this model is centered on the  $C_\alpha$  atom of an amino acid residue, it is easy to ascribe  $w$  and  $v$  values to specific residues. For this reason, most modern implementations of helix–coil theory have used the basics of the Lifson–Roig model (3). For more details, see [Helix–Coil Theory](#).

### Bibliography

1. S. Lifson and A. Roig (1961) *J. Chem. Phys.* **34**, 1963–1974.
2. B. H. Zimm and J. K. Bragg (1959) *J. Chem. Phys.* **31**, 526–535.
3. H. Qian and J. A. Schellman (1992) *J. Phys. Chem.* **96**, 3987–3994.

### Suggestions for Further Reading

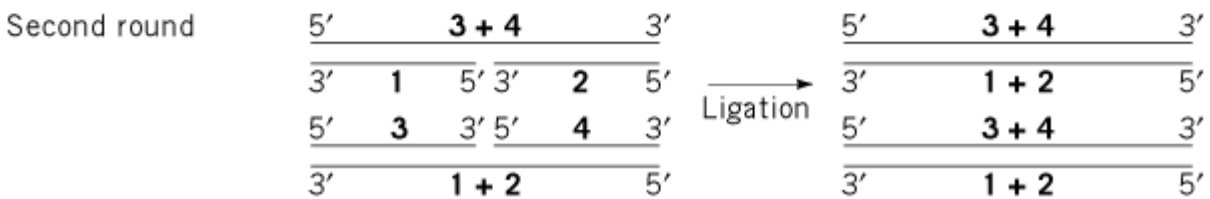
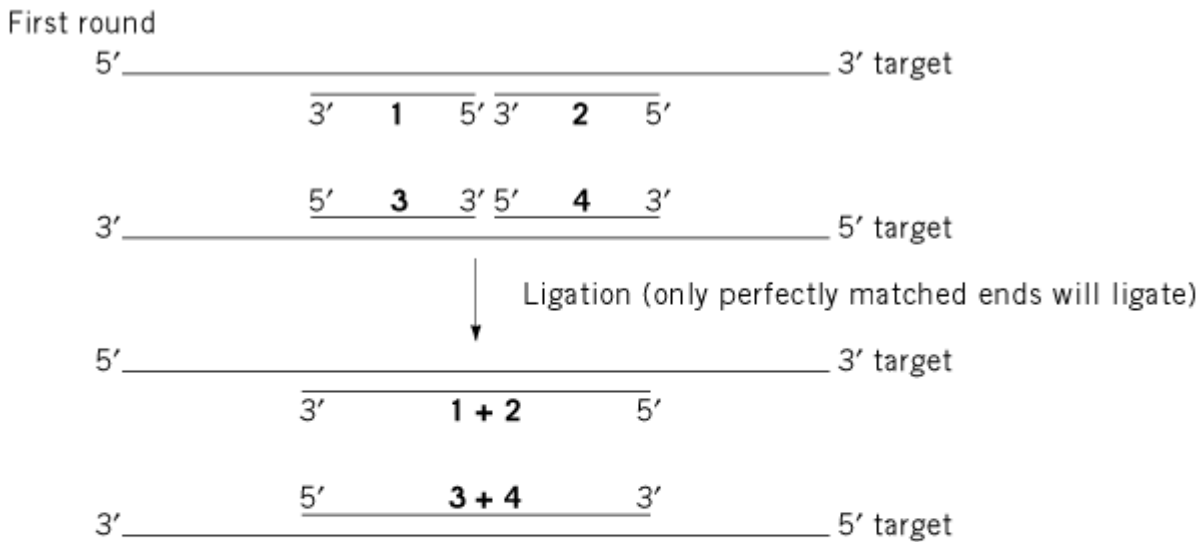
4. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, W. H. Freeman, San Francisco. Chapter 20 provides a good introduction to models for helix formation.
5. D. Poland and H. A. Scheraga (1970) *Theory of Helix-Coil Transitions in Biopolymers*, Academic Press, New York. A broad discussion of helix–coil theory. The book also reprints copies of many of the seminal articles on helix–coil theory.

## Ligase Chain Reaction

The ligase chain reaction (LCR) is another method for amplifying **DNA** and is similar to **PCR** in using a thermostable **enzyme** and specific oligonucleotide primers (1). However, LCR uses [DNA Ligase](#) to join two oligonucleotide primers that are designed to **hybridize** to the target DNA with the 3'-end of one primer adjacent to the 5'-end of the second primer. A single base-pair mismatch near the junction prevents ligation and thus prevents amplification. Therefore, LCR is an excellent method for detecting single-nucleotide changes in the target DNA.

To get amplification, four primers are used, two for each strand of the target DNA, so that the ligation product from one pair of primers is the template for ligation of the other pair of primers and vice versa (Fig. 1). LCR amplifies as few as 200 target molecules to give a detectable DNA product after 30 cycles. LCR is used to develop clinical assays (2) and is another promising tool for [recombinant DNA](#) research.

**Figure 1.** The mechanism of the ligase chain reaction. The four primers used are numbered 1 through 4.



## Bibliography

1. F. Barany (1991) Proc. Natl. Acad. Sci. USA **88**, 189–193.
2. C. S. Hill (1996) J. Clinical Ligand Assay **19**, 43–52.
3. R. D. Abramson and T. W. Myers (1993) Curr. Opinion Biotechnol. **4**, 41–47.

## Ligation

A key step in synthesizing [recombinant DNA](#) molecules and therefore in [cloning](#) is covalently joining two DNA fragments that have matching ends by the enzyme [DNA Ligase](#). This reaction, called ligation, involves forming a phosphodiester bond between the 5'-phosphate on one end and the 3'-OH on the other end of adjacent strands. Energy is required to drive the reaction, and one set of DNA ligases, represented by **bacteriophage T4** ligase, uses ATP as the source of energy, whereas the other class, represented by *Escherichia coli* DNA ligase, uses **NAD**. If the molecules to be ligated have identical **cohesive ends** that can base pair, ligation is quite efficient. However, molecules with blunt ends are more difficult to ligate (see [Blunt-End Ligation](#)). Because the cohesive ends on most DNA fragments are short, ligations are usually carried out at the relatively low temperature of 14°C. This temperature is a compromise that allows a reasonable rate of ligation, while allowing base pairing between the cohesive ends.

Cloning usually requires two ligation events. One is the joining of two linear DNA molecules, and

the other is circularization of the recombinant DNA molecule produced by the first event. The efficiency of joining increases as the DNA concentration is increased. If the concentration is too high, however, chimeric molecules are produced that contain two inserts joined to one **plasmid** molecule or two plasmid molecules joined to one insert. The second event, circularization, is a monomolecular reaction that is independent of the DNA concentration. A cloning ligation is carried out in a single step at a DNA concentration that is a compromise to allow both events, or it is carried out initially at a DNA concentration optimal for the joining step, followed by dilution and the addition of more ligase for circularization.

### Suggestion for Further Reading

A. E. T. Konson and D. S. Levin (1997) Mammalian DNA ligases, *Bioessays* **19**, 893–901; a review of DNA ligases.

## Light Scattering

Light scattering is well-suited for studying biological [macromolecules](#) in solution. The measurements are generally easy to perform, and they can be done on solutions with relatively low concentrations of the particles of interest. Information about the state of association and conformation of particles in solution can be obtained. The scattering data can be analyzed kinetically, as well as under equilibrium conditions, so information on binding mechanisms or molecular associations can be obtained. This entry focuses primarily on static light scattering. The theory and practice of [dynamic light scattering](#) is dealt with separately. Absorption and Raman light scattering techniques depend upon inelastic scattering that involves changes in energy between the incident and scattered light, and these are discussed in [Absorption spectroscopy](#) and [Vibrational Spectroscopy](#).

### 1. The Light Scattering Experiment and Underlying Theory

The static light scattering experiment is analogous to [small-angle scattering](#) experiments using neutrons or X-rays, in that a sample is irradiated with monochromatic light and the coherent, elastically scattered radiation is measured as a function of the angle between the incident and scattered beam. Static light scattering measurements monitor the total light scattering intensity averaged over time. [Dynamic light scattering](#) monitors the fluctuations in the intensity of light scattered by small volume elements of a solution, which is directly related to the Brownian motion of the solutes. Lasers are the preferred source of light, because they provide intense, coherent light of a single wavelength over a wide range of useful wavelengths, from the near ultraviolet to the infrared region of the spectrum ( $\sim 400$  to  $1000 \text{ \AA}$ ) ( $1 \text{ \AA} = 10^{-10} \text{ m}$ ). These wavelengths are  $\sim 2$  to  $3$  orders of magnitude longer than those typically provided by X-ray and neutron sources, and hence the dimensions that can be effectively probed are likewise larger. Laser light sources are extremely intense, however, and measurements can be made on solutions with much lower concentrations ( $\mu\text{M}$  as opposed to  $0.1$  to  $1 \text{ mM}$ ) of solute than is needed for the typical X-ray or neutron scattering experiment. The lower concentration limit is set by the scattering from the solvent.

Light propagates as electromagnetic radiation and is scattered predominantly through interactions between its electric field component and the electrons in a sample. The oscillating electric field causes vibrations of the electrons in the sample, turning the atoms into oscillating dipoles that in turn re-emit radiation. Light scattering depends upon refractive index fluctuations, which in polymer solutions are associated with the polarizability differences between solute and solvent. Static light

scattering is similar to X-ray scattering; however, instead of using the [scattering intensity distribution](#)  $\langle I(\mathbf{Q}) \rangle$ , the Rayleigh ratio,  $R_\theta$ , which corrects for polarization and experiment geometry, is used. For vertically polarized light and detection geometry we have

$$R_\theta = \frac{I_s r^2}{I_0} \quad (1)$$

where  $I_0$  and  $I_s$  are the intensity of the incident and scattered light, respectively,  $r$  is the sample-detector distance, and  $\theta$  is the measurement angle relative to the direction of the incident radiation. For noninteracting particles in solution (molecular weight  $M$  and solute concentration  $c$ ) with dimensions of the order of the wavelength of the light, the intraparticle structure factor,  $P_\theta = I_{s,\theta}/I_{s,\theta=0}$ , is introduced to account for the interference of light scattering from different points within an individual particle:

$$R_\theta = K c M P_\theta \quad (2)$$

$$K = \frac{4\pi^2 n_0^2}{N_A \lambda^4} \left( \frac{dn}{dc} \right)^2 \quad (3)$$

where  $n_0$  is refractive index of the solvent,  $n$  is the refractive index of the solution,  $dn/dc$  is the refractive gradient, and  $N_A$  is Avogadro's number. Equation 2 is alternatively written as (1)

$$\frac{Kc}{R_\theta} = \frac{1}{M P_\theta} \quad (4)$$

For solutions of finite concentration, interparticle interactions can give rise to an interparticle structure factor, and equation 4 is expressed in the form

$$\frac{Kc}{R_\theta} = \frac{1}{M P_\theta} + 2A_2 c \quad (5)$$

where  $A_2$  is the second virial coefficient that accounts for interparticle interactions. In cases where the interparticle interactions are concentration-dependent, scattering data can be extrapolated to  $c = 0$  (infinite dilution) to eliminate effects of interparticle interactions.

The intraparticle structure factor,  $P_\theta$ , is directly related to the [radius of gyration](#),  $R_g$ , of the scattering particle:

$$P_\theta \cong 1 - \frac{16\pi^2 n_0^2 R_g^2}{3\lambda^2} \sin^2 \frac{\theta}{2} \quad (6)$$

Equation 6 is true for  $0.05 < R_g/\lambda < 0.5$ . For very small angles, and in the limit of infinite dilution ( $c \rightarrow 0$ ), equation 6 becomes the classical light scattering equation (using the approximation,  $1/(1-x) \approx 1+x$  for small  $x$ ):

$$\frac{Kc}{R_\theta} = \frac{1}{M} \left( 1 + \frac{16\pi^2 n_0^2 R_g^2}{3\lambda^2} \sin^2 \frac{\theta}{2} \right) \quad (7)$$

and a plot of  $Kc/R_q$  versus  $\sin^2(q/2)$  yields a slope and intercept that allows calculation of  $M$  and  $R_g$  values for the particles in solution. For polydisperse polymer solutions the structural parameters evaluated in the light scattering experiment ( $M$ ,  $R_g$ , and  $A_2$ ) will be average values.

## 2. Time-Resolved Light Scattering

By combining light scattering methods with techniques such as stopped-flow, photoinduced release of **caged** compounds, photolysis of ligands and/or cofactors, or temperature-jump methods, one can monitor changes in molecular mass related to association and/or dissociation, as well as potential changes in conformation, as a function of time. Typically, experiments are set up to trigger a reaction, or change in environment, at  $t = 0$ , and to monitor the static (or dynamic) light scattering at a fixed angle. These types of measurements can be done at solute concentrations as low as 10  $\mu\text{g/mL}$  and can be used to evaluate effects such as pH, ionic strength, and so on, on the kinetic data obtained. Time scales as short as milliseconds are readily probed, and much faster time scales are accessible in principle. Using time-resolved methods, it is possible to determine whether processes are single or multiple step and, by fitting the data with single or multiple exponential terms, measure rate constants.

## 3. Examples of Light-Scattering Applications

Human [sickle cell disease](#) is characterized by the sickle shape that the patient's erythrocytes form under low oxygen tension. This effect has been traced to a point mutation in [hemoglobin](#) that results in polymerization of the hemoglobin molecules and was observed in cells using light scattering (2). Both static and dynamic light scattering have contributed to our understanding of the structural nature of the hemoglobin polymerization, its time dependence, and the nucleation processes (reviewed in Ref. 3). The hemoglobin polymerization involves the formation of 14-stranded fibers from individual hemoglobin molecules that have assumed the deoxy [quaternary structure](#). The kinetics of the aggregation continue to be studied by light scattering (4). Cells containing these aggregates clump together and have significantly reduced lifetimes as part of the disease's pathology. Another important application of light scattering to a biological system has been the study of the aggregates formed by crystallins, the major protein constituents of the eye lens responsible for the refractive properties critical to vision. These aggregates, as well as the changes they undergo as a function of lens age, cataract formation, and physicochemical changes, have been studied extensively by light scattering. A structural basis for the eye lens transparency has been proposed based on light scattering measurements (5), and kinetic studies of aggregate formation have also been studied (6). Light scattering can also be used to monitor assembly of biological structures. Lyles et al. (7) published an elegant study of the combined use of stopped-flow, classical, and dynamic light scattering applied to the analysis of matrix protein binding to nucleocapsids of the [vesicular stomatitis virus](#) that provides valuable insights into the virus assembly process. Another example of a time-resolved light scattering application is a study of the kinetics of the association of myosin subfragment-1 to [actin](#), in which a rate constant of  $9 \mu\text{M}^{-1}\text{s}^{-1}$  was determined for the association reaction, which was also determined to be anticooperative (8).

### Bibliography

1. B. H. Zimm (1948) *J. Chem. Phys.* **16**, 1093–1099 and 1099–1116.
2. M. Coletta, J. Hofrichter, F. A. Ferrone, and W. A. Eaton (1982) *Nature (London)* **300**, 194–197.
3. F. A. Ferrone (1993) *Experientia* **49**, 110–117.
4. A. Cao et al. (1997) *J. Mol. Biol.* **265**, 580–589.
5. J. Xia et al. (1996) *Biophys. J.* **71**, 2815–2822.
6. F. A. Bassi et al. (1995) *Biophys. J.* **69**, 2720–2727.
7. D. S. Lyles, M. O. McKenzie, and R. R. Hangton (1996) *Biochemistry* **35**, 6508–6518.



8. L. Blanchoin, D. Didry, M.-F. Carlier, and D. Pantaloni (1996) *J. Biol. Chem.* **271**, 12380–12386.

### Suggestions for Further Reading

9. C. R. Cantor and P. R. Schimmel (1980) In *Biophysical Chemistry, Techniques for the Study of Biological Structure and Function*, W. H. Freeman and Co., San Francisco, pp. 812–814.
10. D. L. Rousseau (1984) *Optical Techniques in Biological Research*, Academic Press, New York.
11. C. Tanford (1961) In *Physical Chemistry of Macromolecules*, Wiley New York, pp. 275–316.

## Light-Activated (Caged) Biological Ligands

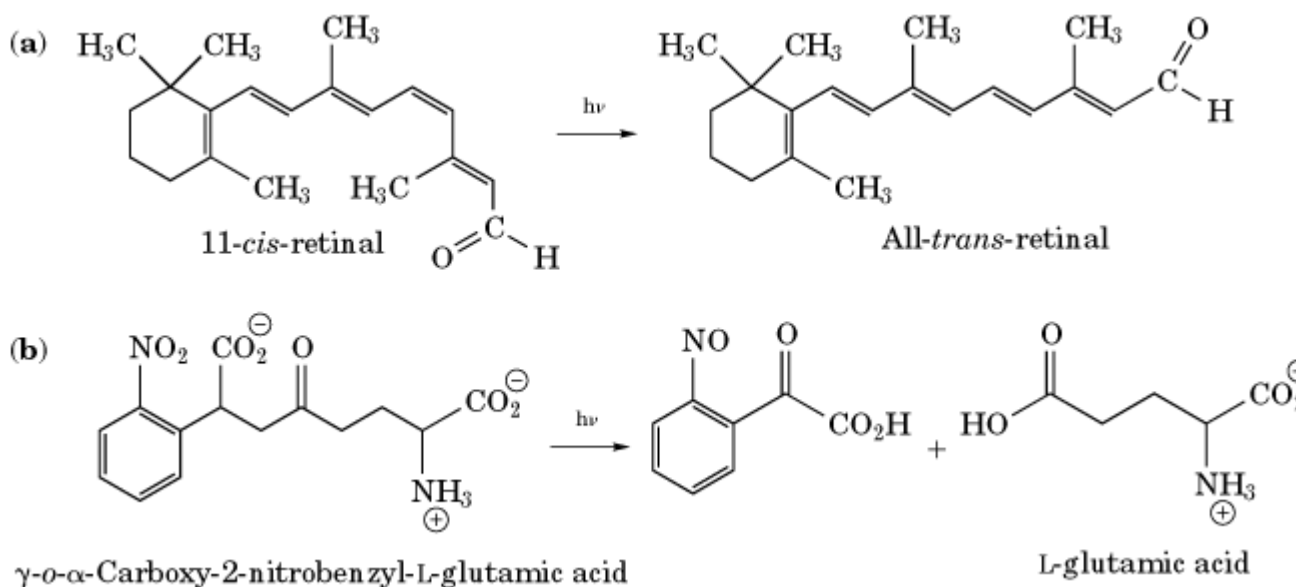
Light-activatable biological **ligands** are compounds that are biologically inert until irradiated by light. They allow one to put an inactive form of the ligand into a complex biological system, such as an intact cell or surrounding organ, and then to activate the ligand on the microsecond time scale and within areas less than  $1 \mu\text{m}^2$ . The resulting temporal and spatial resolution has opened up a vast new area of investigations of biologically interesting problems involving reactions on the surface and inside of cells. The usefulness of light-activatable ligands and new advances in their synthesis and in instrumentation, as well as their commercial availability, are expected to lead to their increasing use. The purpose of this article is to (1) mention some of the biological problems in which the use of light-activated ligands may be of special value, (2) describe the properties of light-activated ligands that make them useful in such studies, and (3) list some of the light-activated ligands already available and describe some of their photochemical characteristics and their use in biological problems.

A frequent problem with kinetic studies of the binding of ligands to their **receptors** in complex biological systems is that the diffusion of the ligand to its receptor is the rate-limiting step. Once bound to the receptor, the subsequent events usually occur much more rapidly and cannot be observed. The kinetics of such processes are observable, however, if the ligand is diffused into place in an inactive form and then activated by a rapid pulse of light.

The use of light-activatable phosphates for investigations of intracellular reactions was developed 20 years ago (1, 2). Light-activatable ligands are also referred to as *caged ligands*, a name introduced by Kaplan et al. (1). Nature uses light-activated ligands; eg, 11-*cis*-retinal (Fig. 1a) in the eyes of animals and in bacterial membranes is inactive (see [Rhodopsin](#) and [Bacteriorhodopsin](#)). On exposure to light, it isomerizes to all *trans* retinal. This initiates a chemical reaction that leads to the synthesis of chemical signals (neurotransmitters), which are released by specific neurons and then bind to neurotransmitter receptors in the membrane of the adjacent cell. Light-activatable neurotransmitters suitable for transient kinetic investigations of the chemical mechanism of neurotransmitter receptor-mediated reactions were only developed more recently (3, 4). Such light-activated neurotransmitter precursors allow transient kinetic investigations to be made of the reactions on the cell surface in the microsecond time region.

**Figure 1.** Light-activatable ligands. (a) Retinal is found in nature, both in the eye of animals and in the membrane of some bacteria. The photoisomerization of 11-*cis*-retinal to 11-*trans*-retinal initiates a cascade of biochemical reactions and occurs in the picosecond time region. (b) Glutamic acid is an important neurotransmitter in the central nervous system. A synthetically prepared and biologically inactive light-activatable glutamate derivative on exposure to near-UV

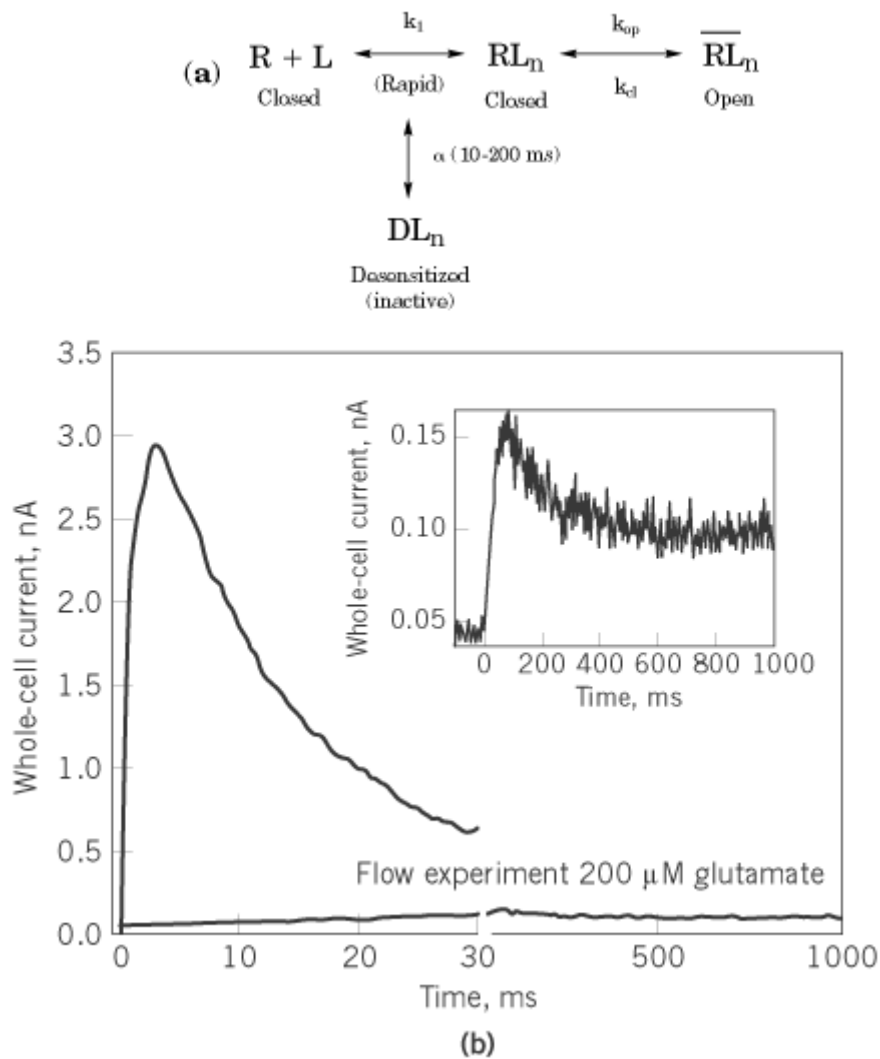
light liberates glutamic acid with a half-time of 21  $\mu$ s.



## 1. Mechanistic Investigations Using Light-Activated Ligands and Their Impact on Signal Transmission Between Cells

The experimental results of Figure 2 b demonstrate the kinetic information available with the use of photoactivated ligands. A photolabile, biologically inert precursor of the excitatory neurotransmitter glutamate (Fig. 1b) was pre-equilibrated with single rat hippocampal neurons carrying glutamate receptors. At zero time the caged glutamate was photolyzed by a laser pulse (337 nm) to liberate 100-mM free glutamate within 60  $\mu$ s (the half-time for liberation of glutamate upon irradiation is 21  $\mu$ s (5)). The free glutamate binds to glutamate receptors on the cell surface, and cation-conducting transmembrane channels consequently open. The current arising from cations passing through the channels is recorded at constant transmembrane voltage by an electrode (6) (Fig. 2b). The current first increases to a maximum value within 3 ms and then falls again, with a time constant of  $\sim$ 10 ms, due to receptor desensitization, which results in the transient closing of receptor-channels. The decrease in the current occurred in two phases: 95% of the current decreased with a time-constant of  $\sim$ 10 ms, the remainder in  $\sim$ 100 ms.

**Figure 2.** Investigations of neurotransmitter-receptor mechanisms using light-activatable receptor ligands. (a) A general mechanism for activation of neurotransmitter receptors, based on the mechanism of the acetylcholine receptor first proposed by Katz and Thesleff (46). Two (subscript  $n$  in the scheme) molecules of neurotransmitter ( $L$ ) bind to a specific receptor protein ( $R$ ) that can then form an open ion-conducting transmembrane channel ( $RL_n$ ). The neurotransmitter can also induce transiently inactive (desensitized) receptor forms ( $DL_n$ ).  $K_1$  is the receptor-neurotransmitter dissociation constant controlling channel opening,  $k_{op}$ ,  $k_{cl}$  are the rate constants for channel opening and closing, respectively, and  $a$  is the rate coefficient for receptor desensitization. (b) A laser-pulse photolysis experiment using a light-activatable biologically inert derivative of glutamate (caged glutamate) (5). A rat hippocampal neuron attached to a whole-cell current-recording electrode (6) was equilibrated with caged glutamate at pH 7.4, 22° to 23°C, and a transmembrane voltage of  $-60$  mV. The caged glutamate was irradiated by a 0.5-mJ flash of 343-nm light from a pulse dye laser (Candela SLL500, 600-ns pulse length). 100-mM glutamic acid was released, and the current due to the opening of transmembrane channels formed by the glutamate receptor recorded. The left half shows the response of the opening of transmembrane channels formed by the glutamate receptor recorded. The left half shows the response of the same neuron when it was exposed to a rapidly flowing 200-mM glutamate solution, and the inset with 100  $\mu$ M (23).



In the second experiment shown in Figure 2b, and in the inset, a 100-mM solution of free glutamate was allowed to flow rapidly over the same cell. The current thus induced reached its maximum value not in 3 ms, but in  $\sim$ 100 ms, and the maximum magnitude of this current was only about 5% of that observed in the experiment on the left. Only one falling phase of the current was observed, with a time constant of 100 ms. The receptor-mediated reaction that desensitizes with a time constant of 10 ms, and is associated with 95% of the current in the experiment shown on the left, is not seen, because the equilibration of glutamate with the receptors is rate-limiting.

Three distinct phases of the reaction could be characterized. (1) From the effect of neurotransmitter concentration on the rising phase of the current (Fig. 2b), one can determine  $K_1$ , the dissociation constant of the receptor site controlling channel opening, and the rate constants for channel opening and closing,  $k_{op}$  and  $k_{cl}$ , (Fig. 2a) (8). (2) From the effect of neurotransmitter concentration on the maximum current amplitude, one can also determine  $K_1$  and the channel-opening equilibrium constant  $F^1 = k_{op}/k_{cl}$ . (3) A falling phase of the current from which one obtains information about the rate constant for receptor desensitization  $a$  (Fig. 2a). Evaluation of all the constants controlling channel opening, which requires transient kinetic measurements over a wide range of neurotransmitter concentration, has been accomplished so far only with the nicotinic [acetylcholine receptor](#) in BC<sub>3</sub>H1 muscle cells (4, 7).

Transient kinetic techniques using light-activated ligands can have a major impact on the elucidation of the mechanism of receptor–drug interaction. Is the inhibitor **competitive** or **noncompetitive**? Do two different inhibitors bind to the same site or different sites? What is the dissociation constant of the inhibitor from the closed-channel forms and from the open-channel form ( $R$  and  $RL_n$ , and  $RL_n$ , respectively, Fig. 2a)? All these questions can now be answered (7).

The ability to perfuse cells with inactive light-activatable neurotransmitters and then to photolyze the compound in very small areas using lasers (8) or two-photon laser excitation (9) also allows one to map (1) specific receptors in specific areas of a cell (10), (2) cells containing specific receptors in brain slices (8), and to identify cells that secrete a specific neurotransmitter and cells that respond to this neurotransmitter (11). This identification can be made (11) when the cells that control an observable response are known, as in the [nematode \*Caenorhabditis elegans\*](#) (12).

## 2. General Properties of Light-Activated Ligands

Several properties of light-activated ligands are necessary to make them useful in investigations of biological reactions. (1) The major requirement for a light-activated ligand is its biological inertness in its original form. This allows one to equilibrate the compound with its reaction partners, a process that is slow on cell surfaces or when the compound must be introduced into a cell. The protecting group that is removed from a biologically active compound must also be inert in order not to interfere with the reaction to be studied. (2) A corollary of the previous requirement is that the light-activatable ligand be stable in the solutions in which the reaction is being investigated. Any thermal hydrolysis of the light-activatable ligand in these solutions to the biologically active substance violates the first requirement. (3) The light-activatable ligand must be soluble in the solutions in which it will be used and photolyzed with adequate quantum yield. This allows one to investigate the reaction over a wide concentration range of the biologically active ligand without using a light energy that will be harmful to the cell. (4) Related to this requirement, the compound must be photolyzed at a wavelength where the cell is not damaged. With muscle cells and neurons, a wavelength region above 336 nm was found to be acceptable (4,7). (5) The light-activatable ligand must be activated to the biologically active molecule in a time that is short compared to the reaction to be investigated. For transient kinetic investigations of neurotransmitter receptor-mediated reactions, caged neurotransmitters are available that are photolyzed within the 1- to 100-ms time region. Rapid photolysis and high quantum yield are also important in obtaining good spatial resolution; diffusion of the photolytically liberated biologically active compound away from the irradiated area diminishes the spatial resolution. Some light-activatable ligands that meet these requirements will now be described.

## 3. Photolabile Protecting Groups Available for the Preparation of Light Activatable Ligands

The *o*-nitrobenzyl group (Table 1A) has been used extensively in organic synthesis as a photosensitive protecting group for **carboxyl**, phosphate, hydroxyl, and **amino groups** (13-15). The use of *o*-nitrobenzyl photochemistry to achieve rapid changes in the concentration of biologically interesting molecules started with ATP (1, 2) and cyclic nucleotides (16, 17) (see [Caged ATP](#)).

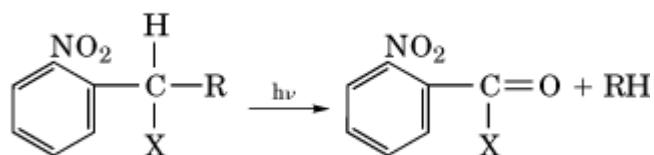
**Table 1. Some Photolabile Protecting Groups for Amino Groups, Carboxyl Groups, Thiol Groups, and Phenolic Hydroxyl Groups**

---

(Photolabile protecting groups for phosphates are discussed in [Caged ATP](#).)

---

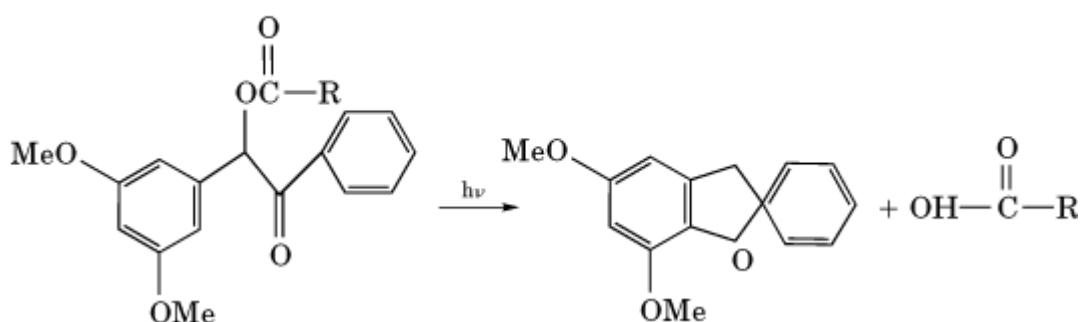
**A. *O*-Nitrobenzyl group**



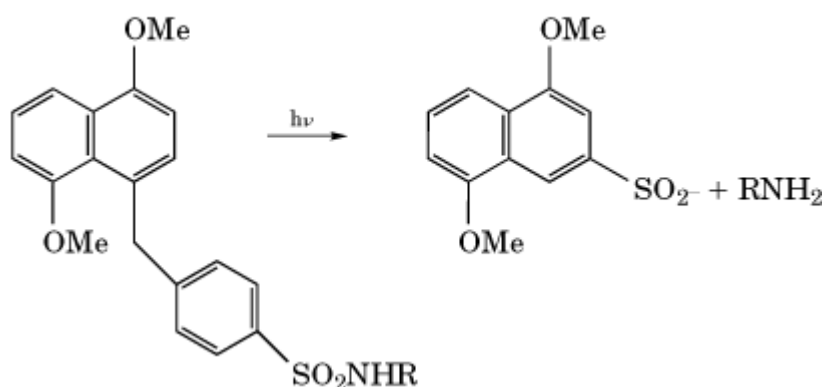
X = H, CH<sub>3</sub>, or COOH

X = H, CH<sub>3</sub>, or COOH

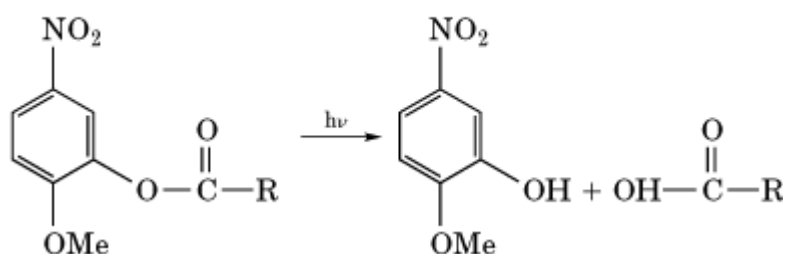
**B. 3',5'-Dimethoxy-benzoin group**



**C. Sulfonamide-protecting group**



**D. 2-Methoxy-5-nitrophenyl group**

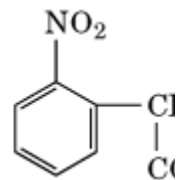
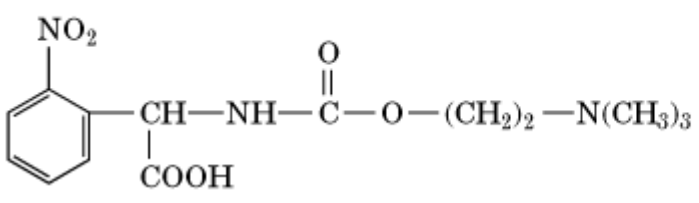
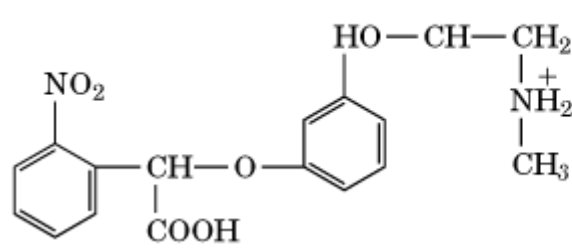
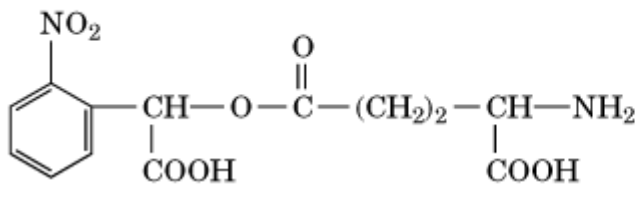
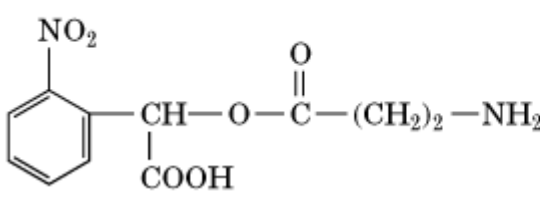


---

The  $\alpha$ -carboxy-2-nitrobenzyl group (3) Table 2A) was used to synthesize light-activatable ligands suitable for transient kinetic investigations of the excitatory acetylcholine, glutamate, and kainate receptors and the inhibitory [gamma-aminobutyric acid](#) (GABA) and glycine receptors (18). These

light-activatable ligands meet all the criteria outlined above for suitability for investigations of biological reactions.

**Table 2. Photolytic Properties of Biologically Inert, Photolabile Derivatives of Neurotransmitters**

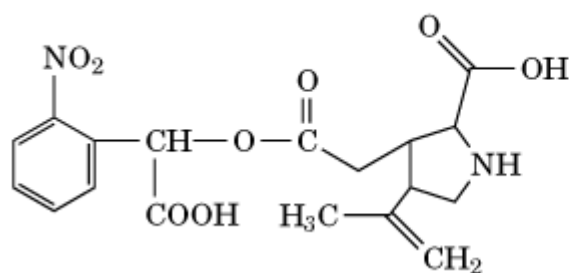
| Caged Neurotransmitter  | Group Caged        | Photolysis $t_{1/2}$ (ms) | Product Quantum Yi  |
|---|--------------------|---------------------------|---|
| <b>A. <math>\alpha</math>-Carboxy-2-nitrobenzyl caging group (neutral pH, RT)</b>   |                    |                           |   |
|   |                    |                           |  |
| Carbamoylcholine  | carbamate          | 45                        | 0.3   |
|   |                    |                           |   |
| Phenylephrine   | phenolic           | 350                       | 0.28  |
|  |                    |                           |   |
| Glutamate   | $\gamma$ -carboxyl | 21                        | 0.14  |
|  |                    |                           |   |
| $\gamma$ -Aminobutyric acid (GABA)  | $\gamma$ -carboxyl | 19                        | 0.16  |
|  |                    |                           |   |

Kainate

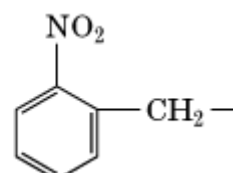
$\gamma$ -carboxyl

45

0.37



**B. 2-Nitrobenzyl caging group**

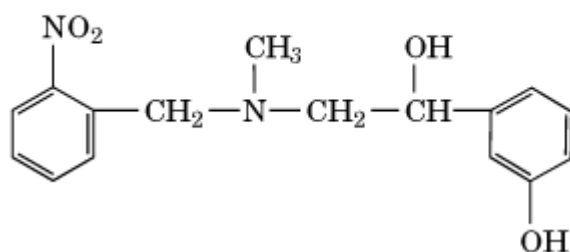


Phenylephrine

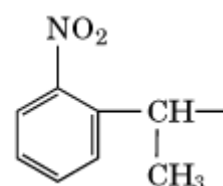
secondary  
amino group

40

0.12



**C.  $\alpha$ -Methyl-2-nitrobenzyl caging group**

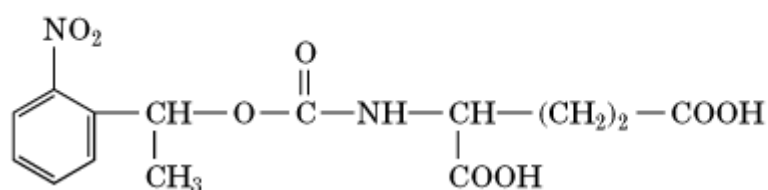


Glutamate

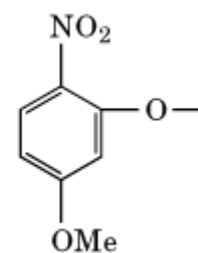
$\alpha$ -amino group

4500

0.65



**D. 2-Methoxy-5-nitrophenyl caging group**

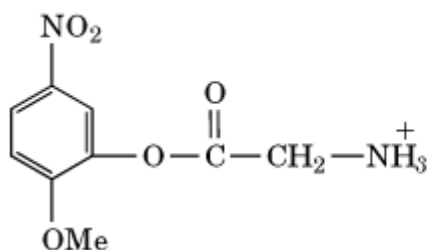


Glycine

carboxyl

&lt;1

0.2



In addition to the substitutions on the benzyl carbon of the nitrobenzyl group indicated in Table 1, nitrophenyl derivatives with different substitutions on the phenyl ring have also been synthesized [reviewed in (19), (20)]. For example, the 4,5-dimethoxy-2-nitrobenzyl group has the advantage of a much greater absorbance at 360 nm than the 2-nitrobenzyl caging group (Table 1A), but the quantum yield is much lower at this wavelength than that of the corresponding nitrobenzyl derivatives in the 300- to 337-nm wavelength region (21). A number of other light removable groups are under development for the protection of carboxylic acids and amines. Sheehan and colleagues (22) found that esters of 3',5'-dimethoxybenzoic acid (Table 1B) undergo efficient and clean photolysis. This photosensitive protecting group and its derivatives (eg, those in which the methoxy groups are lacking) have been used to protect the carboxyl group of the neurotransmitters glutamic acid and  $\gamma$ -aminobutyric acid (23) and carbamates (24).

The photodeprotection of carbamates [Table 1B (R = NH<sub>2</sub>-R'')] is of interest because it constitutes another approach for caging amino groups, which are abundant in biological effector molecules, including neurotransmitters (Table 2). The rate-limiting step in the formation of the amino-containing compound is the loss of CO<sub>2</sub> in a single exponential process with a half-time of 4.5 ms at neutral pH and room temperature (24). The use of sulfonamides for the photoprotection of amino groups (Table 1C) is also currently under investigation (25). The 2-methoxy-5-nitrophenyl group (26) has been used to synthesize compounds that on exposure to light liberate glycine (27) and  $\beta$ -alanine (28).  $\beta$ -alanine activates the inhibitory glycine receptor (29). Both compounds are photolyzed within  $\sim 1 \mu\text{s}$ , with a product quantum yield of 0.2, and are biologically inert before photolysis. The light-activatable glycine derivative (caged glycine) is considerably less stable than the  $\beta$ -alanine derivative in aqueous solutions at neutral pH and is, therefore, more difficult to use in experiments. The nitrophenyl benzyl group and its derivatives are currently the best characterized light-activatable ligands suitable for investigations of biological processes.

#### 4. The Uses of Light-Activatable Ligands

For investigations of the  $\alpha_1$ -adrenergic receptor, caged phenylephrine has been prepared by Nerbonne and her colleagues (30, 31) and by Walker and his colleagues (32). Many intracellular processes are controlled by calcium signaling. For the rapid and spatially controlled release of Ca<sup>2+</sup>, the calcium chelator ethylene diamine tetraacetic acid (EDTA) has been caged with the dimethoxy nitrophenyl group (33, 34) and the nitrophenyl group (35). Uncaging of the EDTA alters the binding affinity for Ca<sup>2+</sup> over 10,000 fold, allowing large spatially defined increases in Ca<sup>2+</sup> concentrations within 500  $\mu\text{s}$  (36). Nitric oxide plays a role in cellular function, including signaling between central nervous system neurons and controlling dilation of blood vessels (47). Several caged compounds that liberate nitric oxide upon irradiation with near-UV light are now available (37, 38). Fatty acids are essential components of cells and regulate many cellular functions (39) and a light-activatable fatty acid has recently been synthesized (40).



To produce a rapid increase in the concentration of proteins in a spatially controlled manner inside cells, functional groups of proteins and peptides have been caged. For example, the caging of the amino group of [lysine](#) residues in [G-actin](#) results in a biologically inactive molecule ([41](#)); the ability to liberate G-actin by light in a temporally and spatially controlled manner enables one to elucidate the role of the protein in muscle contraction and the formation of actin filaments. Another example is the demonstration that the activity of protein can be controlled using the  $\alpha$ -carboxy-2-nitrobenzyl group ([5](#)) to cage the [thiol groups](#) of [cysteine](#) residues ([42](#)), thus abolishing the pore-forming ability of  $\alpha$ -hemolysin, a bacterial pore-forming protein that allows relatively large molecules (up to 3000 MW) to pass through its pore ([43](#)). Irradiation with near-UV light restores the formation of pores. The caged pore-forming protein and the spatial resolution provided by laser light may be of interest for drug delivery to specific cells.

The caging of specific enzyme substrates and proteins to investigate and control intracellular reactions is still in its infancy. The quantitative use of light-activatable neurotransmitters in transient kinetic measurements has been described in greater detail in the last few years ([18](#), [44](#), [45](#)). Many receptors and their **isoforms** from various regions of the nervous system, from different species, and from animals that exhibit defects in the receptor-mediated reactions are now available. It will be possible to study the hundreds of therapeutically useful compounds and abused drugs that affect receptor function. However, only a few of the many biological problems to whose solution the use of light-activatable ligands can contribute have been mentioned.

#### Bibliography

1. J. H. Kaplan, B. Forbush, and J. F. Hoffman (1978) *Biochemistry* **17**, 1929–1935.
2. J. A. McCray and D. R. Trentham (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7237–7241.
3. T. Milburn, N. Matsubara, A. P. Billington, J. B. Udgaonkar, J. W. Walker, B. K. Carpenter, W. W. Webb, J. Marque, W. Denk, J. A. McCray, and G. P. Hess (1989) *Biochemistry* **28**, 49–55.
4. N. Matsubara, A. P. Billington, and G. P. Hess (1992) *Biochemistry* **31**, 5507–5514.
5. R. Wieboldt, K. R. Gee, L. Niu, D. Ramesh, B. K. Carpenter, and G. P. Hess (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8752–8756.
6. O. P. Hamill, A. Marty, E. Neher, B. Sakmann, and F. J. Sigworth (1981) *Pfluegers Arch.* **391**, 85–100.
7. C. Grewer and G. P. Hess (1999) *Biochemistry* (in press).
8. M. B. Dalva and L. Katz (1994) *Science* **265**, 255–256.
9. W. Denk, J. H. Stickler, and W. W. Webb (1990) *Science* **248**, 73–76.
10. W. Denk (1994) *Proc. Natl. Acad. Sci. USA* **91**, 6629–6633.
11. H. Li, L. Avery, W. Denk, and G. P. Hess (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5912–5916.
12. D. G. Albertson and J. N. Thompson (1976) *Phil. Trans. R. Soc. Lond. B* **275**, 299–325.
13. H. A. Morrison (1969) In *The Chemistry of the Nitro and Nitroso Groups* (H. Feuer, ed.) Interscience, New York, pp. 165–213.
14. V. N. R. Pillai (1980) *Synthesis* 1–26.
15. H. Schupp, W. K. Wong, and W. Schabel (1987) *J. Photochem.* **36**, 85–97.
16. J. Engels and E.-J. Schlaeger (1977) *J. Med. Chem.* **20**, 907–911.
17. J. M. Nerbonne (1996) *Curr. Opin. Neurobiol.* **6**, 379–386.
18. G. P. Hess and C. Grewer (1998) *Meth. Enzymol.* **291**, 443–473.
19. S. R. Adams and R. Y. Tsien (1993) *Ann. Rev. Physiol.* **55**, 755–783.
20. J. E. T. Corrie and D. R. Trentham (1993) *Bioorganic Photochemistry*, Vol. **2** (H. Morrison, ed.), Wiley, New York.
21. J. F. Wootton and D. R. Trentham (1989) *NATO ASI Sect. C* **272**, 277–296.
22. J. C. Sheehan, R. M. Wilson, and A. W. Oxford (1971) *J. Am. Chem. Soc.* **93**, 7222–7228.

23. K. R. Gee, L. W. Kueper, III, J. Barnes, G. Dudley, and R. S. Givens (1996) *J. Org. Chem.* **61**, 1228–1233.
24. G. Papageorgiou and J. E. T. Corrie (1997) *Tetrahedron* **53**, 3917–3932.
25. J. E. T. Corrie and G. Papageorgiou (1996) *J. Chem. Soc. Perkin Trans. I*.
26. P. Kuzmic, L. Pavlickova, and M. Soucek (1986) *Collection Czechoslovak Chem. Commun.* **51**, 1292–1300.
27. D. Ramesh, R. Wieboldt, L. Niu, B. K. Carpenter, and G. P. Hess (1993) *Proc. Natl. Acad. Sci. USA* **90**, 11074–11078.
28. L. Niu, R. Wieboldt, D. Ramesh, B. K. Carpenter, and G. P. Hess (1996) *Biochemistry* **35**, 8136–8142.
29. D. Choquet and H. Korn (1988) *Neurosci. Lett.* **84**, 329–334.
30. J. M. Nerbonne (1996) *Curr. Opin. Neurobiol.* **6**, 379–386.
31. S. M. Muralidharan, G. M. Mayer, W. B. Boyle, and J. M. Nerbonne (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5199–5203.
32. J. W. Walker, H. Martin, F. R. Schmitt, and R. J. Barsotti (1993) *Biochemistry* **32**, 1338–1345.
33. S. R. Adams, J. P. Y. Kao, G. Gryniewicz, A. Minta, and R. Y. Tsien (1988) *J. Am. Chem. Soc.* **110**, 3212–3220.
34. J. H. Kaplan and G. C. R. Ellis-Davies (1988) *Proc. Natl. Acad. Sci. USA* **85**, 6571–6575.
35. G. C. R. Ellis-Davies and J. H. Kaplan (1994) *Proc. Natl. Acad. Sci. USA* **91**, 187–191.
36. G. C. R. Ellis-Davies, J. H. Kaplan, and R. J. Barsotti (1996) *Biophys. J.* **70**, 1006–1016.
37. O. Arancio, M. Kiebler, C. J. Lee, V. Leu-Ram, R. Y. Tsien, E. R. Kandel, and R. D. Hawkins (1996) *Cell* **87**, 1025–1035.
38. N. Bettache, T. D. Carter, J. E. T. Corrie, D. Ogden, and D. R. Trentham (1996) *Methods Enzymol.* **268**, 266–281.
39. J. X. Kang, Y. F. Xiao, and A. Leaf (1995) *Proc. Natl. Acad. Sci. USA* **91**, 3997–4001.
40. Y. Xia, X. Huang, R. Sreekumar, and J. W. Walker (1997) *Bio-org. Med. Chem. Lett.* **7**, 1243–1248.
41. G. Marriott (1994) *Biochemistry* **33**, 9092–9097.
42. C. Y. Chang, B. Niblack, B. Walker, and H. Bayley (1995) *Chem. Biol.* **2**, 391–400.
43. L. Niu, K. R. Gee, K. Schaper, and G. P. Hess (1996) *Biochemistry* **35**, 2030–2036.
44. G. P. Hess (1993) *Biochemistry* **32**, 989–1000.
45. G. P. Hess, L. Niu, and R. Wieboldt (1995) *Ann. NY Acad. Sci.* **757**, 23–29.
46. B. Katz and S. Thesleff (1957) *J. Physiol. (Lond.)* **138**, 63–80.
47. E. R. Kandel, J. H. Schwartz, and T. M. Jessell (1995) *Essentials of Neural Science and Behavior*, Appelton and Lang, Norwalk, CT.

### **Suggestion for Further Reading**

48. G. Marriott, ed. (1997) *Methods in Enzymology* Vol. 291, Academic Press, San Diego, CA.

### **Light-Harvesting Complex**

Light-harvesting complexes of pigments and [proteins](#) enable **photosynthetic** organisms to use a

broad range of the solar spectrum and to grow even at low light intensities. Following absorption of a photon by an antenna pigment, the exciton moves from pigment to pigment via so-called resonance energy transfer, until becoming trapped in a [photosynthetic reaction center](#).

A variety of light-harvesting complexes exist, differing in different taxonomic groups. Plants have an inner antenna of chlorophyll-*a*-binding proteins and an outer antenna of chlorophyll-*a/b* (CAB) binding proteins. The inner antenna, typically containing 80–100 chlorophyll-*a* molecules, can be an integral part of the reaction center, as in the case of photosystem I, or be composed of two separate chlorophyll-*a* binding proteins associated with photosystem II. The CAB proteins bind half of the thylakoid chlorophyll-*a* and all of the chlorophyll-*b*, and they comprise an outer antenna of several hundred chlorophyll molecules. The various CAB proteins associated with photosystem II or photosystem I are integral membrane proteins, in the 25-kDa molecular weight range, with three transmembrane-spanning regions. Each protein has been estimated to bind up to 14 chlorophyll and one carotenoid molecules.

Cyanobacteria and red algae use water-soluble protein structures, known as *phycobilisomes*, as outer light-harvesting antennae. Purple bacteria have two types of light-harvesting proteins (LHI and LHII). They are composed of small proteins that span the membrane once and bind bacteriochlorophylls and carotenoids, and they form ring-like structures around the reaction center. Green sulfur bacteria have pigment assemblies known as *chlorosomes*. For further details, see [Photosynthesis](#).

## Lines

LINEs are *long interspersed repeated segments* (see [Interspersed DNA Elements](#)), first described in the [genomes](#) of mammals. They are over 5 kb long and are repeated approximately  $10^4$  times in the genome. A 6.4-kb member of one LINE family occurs downstream of the human **beta-globin** gene. This primate LINE family is called the *Kpn family*, because it generates [restriction fragments](#) of characteristic size when the genome is fragmented by the *Kpn* [restriction enzyme](#).

LINEs of the **mouse** genome have been prepared by digesting reassociated DNA fragments with **deoxyribonuclease** specific for single strands. Separation of the resulting fragments on the basis of size reveals many discrete bands that are assigned to 15-kbp and 6-kbp groups. The 6-kbp LINE family, called MIF-1, has been described in several species of mice and other rodents. As with the primate *Kpn* family, MIF subfamilies are distinguishable within each of several genomes by restriction-site polymorphisms (see [Restriction Fragment Length Polymorphism \(RFLP\)](#)), and related genomes contain distinct divergent subfamilies. The rodent and primate LINEs are sufficiently different not to cross-hybridize.

A highly repeated DNA element has been described in the **Xenopus laevis** genome. It is present in about 8500 copies per haploid genome and accounts for about 2.4% of the genome. The copies range in size from 6 kbp to 10 kbp because of an expandable region containing variable numbers of a tandemly repeating 183– to 204 bp unit. The element is flanked by an imperfect 18 bp inverted sequence, and inverted repeats of 180 to 185 bp are nearby. This and other indirect evidence suggest that this element may be a [transposable element](#) (1).

DNA sequence analysis of a region contained within a LINE in mice reveals a long, open reading frame (ORF) of 978 bp between two **stop codons**. This sequence is conserved among three distantly related mice species and between mouse and monkey, in a manner that is claimed to be characteristic

of regions undergoing selection for protein function (2).

The *Drosophila mobile element* jockey is similar in its structural organization and coding potential to the LINES of various organisms. A complete copy of jockey (about 5 kbp) is terminated by a segment of oligo(dA), preceded by two long ORF that overlap with a –1 frameshift. On the basis of **sequence analysis**, the first ORF codes for a nucleic acid-binding protein and the second codes for a **reverse transcriptase** that is most similar in its sequence to the putative reverse transcriptase of other LINES. The existence of a large number of *jockey* copies with a deletion in the second ORF may indicate that they can use reverse transcriptase in *trans* (3).

Sequence analysis of the rat *vasopressin* and *oxytocin gene family* reveals that the two genes are linked by a LINE containing seven ORFs encoding hypothetical proteins of 99 to 566 amino acid residues. Furthermore, both DNA strands of LINE serve as templates for **transcription**. Transcripts initiated at the 3' end are more abundant than those started from the 5' end. Transcription occurs preferentially in brain tissue, as analyzed by **Northern blot**, [in situ hybridization](#), and **nuclease** protection experiments. Most LINES are transcribed over their entire length, and a major fraction of the resulting RNA does not enter the cytoplasm, but remains in the cell nucleus (4).

#### Bibliography

1. B. K. Kay and I. B. Dawid (1983) J. Mol. Biol. **170**, 583–596.
2. S. L. Martin et al. (1984) Proc. Natl. Acad. Sci. USA **81**, 2308–2312.
3. A. F. Priimagi, L. J. Mizrokhi, and Y. V. Ilyin (1988) Gene **70**, 253–262.
4. E. Schmitz, E. Mohr, and D. Richter (1991) DNA Cell Biol. **10**, 81–91.

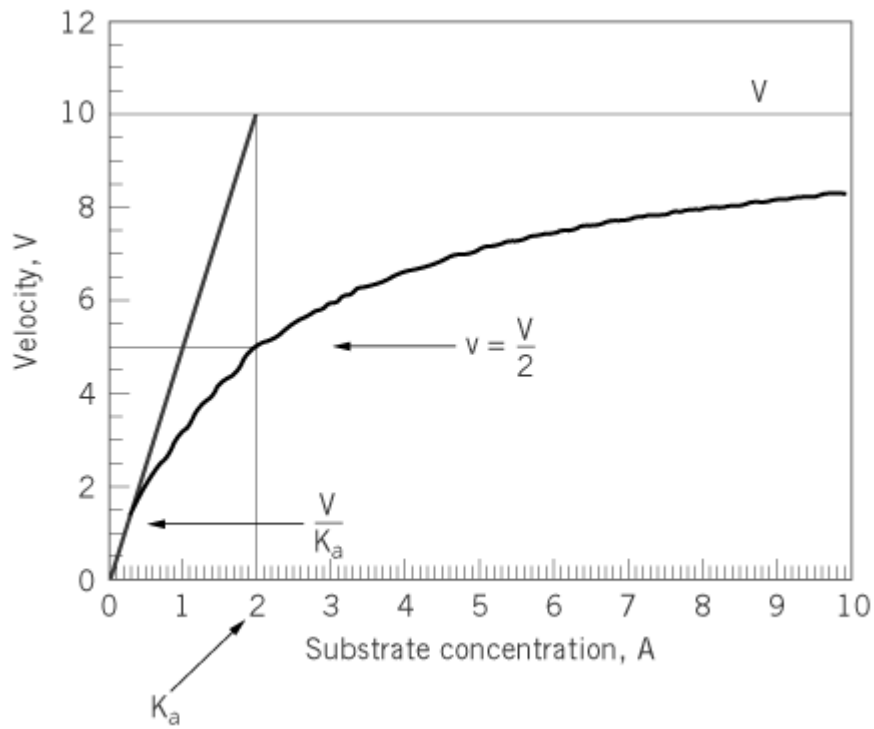
#### Lineweaver–Burk Plot

[Enzymes](#) that conform to [Michaelis–Menten kinetics](#) yield plots of initial velocity  $v$ , as a function of substrate concentration  $A$  that are rectangular hyperbolas of the form

$$v = \frac{VA}{K_a + A} \quad (1)$$

where  $K_a$  denotes the  $K_m$  (Michaelis constant) for substrate  $A$  and  $V$  is the maximum velocity of the reaction. Hyperbolas of this type pass through the origin with an initial slope of  $V/K_a$  and have an asymptote, where  $v = V$  (Fig. 1). The determination of values for  $V$  and  $K_a$  from a hyperbolic curve is difficult, especially as the asymptote cannot be well defined and error will always be associated with initial velocity data. To overcome this difficulty, the equation was rearranged to linear forms to permit the graphical determination of values for the kinetic parameters. The first of the linearizations was reported by Lineweaver and Burk (1). Subsequently, two other linearizations were suggested; one was advanced by Eadie (2) and Hofstee (3) (see [Eadie–Hofstee Plot](#)) and the other by Hanes (4).

**Figure 1.** Variation of initial velocity  $v$  as a function of the concentration of substrate  $A$  as described by Equation 1.

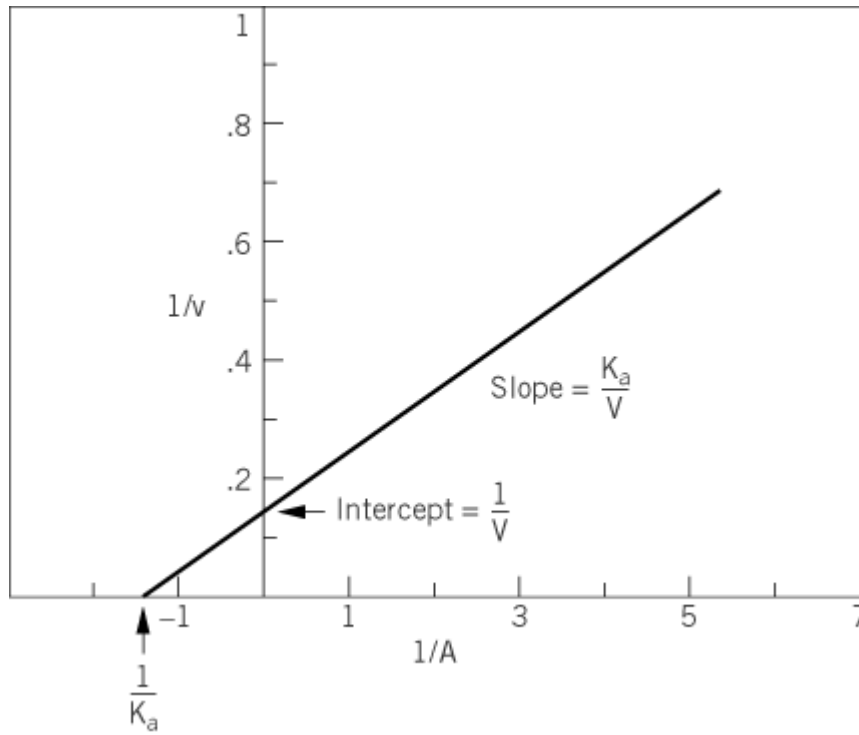


The Lineweaver–Burk rearrangement of the initial velocity equation involves taking reciprocals of each side of the equation and rearranging:

$$\frac{1}{v} = \frac{K_a}{V} \frac{1}{A} + \frac{1}{V} \quad (2)$$

A plot of  $1/v$  against  $1/A$  yields a straight line, with a slope of  $K_a/V$  and an intercept with the vertical ordinate of  $1/V$  (Fig. 2). The intersection of the line with the abscissa occurs at the point where  $1/v = 0$  and at this point  $-1/A = 1/K_a$ .

**Figure 2.** Double-reciprocal plot of the variation of the initial velocity  $v$  of an enzyme-catalyzed reaction as a function of substrate concentration  $A$ .



Determinations of values for the kinetic parameters for a two-substrate reaction are more complex. This can be illustrated by reference to the initial velocity equation for a sequential Bi–Bi reaction that conforms to Michaelis–Menten kinetics involving random binding of substrates *A* and *B* to the enzyme, to form a ternary *EAB* complex (Eq. 3):

$$v = \frac{VAB}{K_{ia}K_b + K_aB + K_bA + AB} \quad (3)$$

where  $K_a$  and  $K_b$  are the dissociation constants for the dissociation of each substrate from the ternary *EAB* complex and  $K_{ia}$  is the dissociation constant of the *EA* complex. The reciprocal form of this equation is given by Equation 4:

$$\frac{1}{v} = \frac{K_a}{V} \left[ \frac{K_{ia}K_b}{K_aB} + 1 \right] \frac{1}{A} + \frac{1}{V} \left[ \frac{K_b}{B} + 1 \right] \quad (4)$$

For a plot of  $1/v$  against  $1/A$ , Equation (4) is that of a straight line with both slope and intercept varying as a function of the concentration of substrate *B*. Data obtained for the variation of the initial velocity as a function of the concentration of *A*, at different fixed concentrations of *B*, would give a family of straight lines that intersect at a point to the left of the vertical ordinate (Fig. 3), where the  $1/v$  and  $1/A$  coordinates are  $(1/V) [1 - K_a/K_{ia}]$  and  $K_a/(K_{ia}K_b)$ , respectively. Thus, the crossover point may be above, on, or below the abscissa, depending on the relative values of  $K_a$  and  $K_{ia}$ . The intersection of each straight line with the abscissa would give only an apparent value for  $K_a$  at a particular concentration of *B*. The slope of the lines as a function of *B* is described by Equation (5):

$$\text{Slope} = \frac{K_{ia}K_b}{V} \frac{1}{B} + \frac{K_a}{V} \quad (5)$$

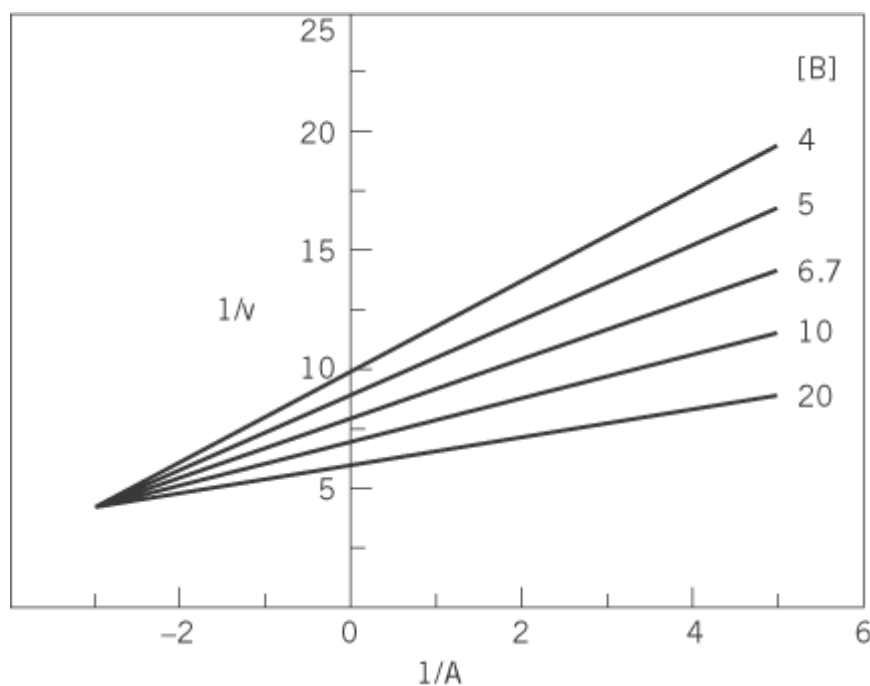
so that a replot of the slopes of the lines of the primary plot against  $1/B$  would yield a straight line

that intersects the abscissa at a point equal to  $K_a/K_{ia}K_b$ . For a *rapid-equilibrium random* mechanism, this value is equal to  $K_{ib}$ , the dissociation constant for the interaction of  $B$  with free enzyme. For an ordered mechanism, there is no  $EB$  complex. The intersection points of the lines of the primary plot with the vertical ordinate give only apparent maximum velocities at different fixed concentrations of substrate  $B$ . Variation of the intercepts with the concentration of  $B$  is described Equation (6):

$$\text{Intercept} = \frac{K_b}{V} \frac{1}{B} + \frac{1}{V} \quad (6)$$

so that a replot of the intercepts of the primary plot against  $1/B$  would be a straight line that intersects the vertical ordinate at the reciprocal of the true value for the maximum velocity and the abscissa at the reciprocal of the true value for  $K_b$ . Values for  $V$ ,  $K_{ia}$ , and  $K_a$  would be obtained in a similar manner by starting with a plot of  $1/v$  against  $1/B$  at different fixed concentrations of  $A$ .

**Figure 3.** Double-reciprocal plot of the variation of the initial velocity of an enzyme-catalyzed reaction involving the sequential addition of two substrates,  $A$  and  $B$ , as a function of the concentration of substrate  $A$  at different fixed concentrations of substrate  $B$ .



There has been considerable discussion in the past about the relative merits of the Lineweaver–Burk plot, the Eadie–Hofstee plot, and Hanes linearization procedures for obtaining the best estimates of values for kinetic parameters (3, 5). Those days have long since passed, and now computer programs are available for least-squares fitting of data, with appropriate weighting factors, to an assumed rate equation (6). Graphical methods are important for determining the form of the rate equation to which the data are to be fitted and for illustrating the results of kinetic investigations. The Lineweaver–Burk plot must be considered the most satisfactory of the three types of plot, as it shows the straightforward variation of one dependent variable as a function of the concentration of one or two independent variables.

#### Bibliography

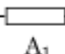
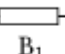
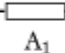
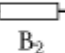
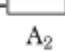
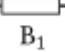
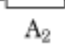
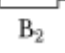
1. H. Lineweaver and D. Burk (1934) *J. Amer. Chem. Soc.* **56**, 658–666.

2. G. S. Eadie (1942) *J. Biol. Chem.* **146**, 85–93.
3. B. H. J. Hofstee (1959) *Nature* **184**, 1296–1298.
4. C. S. Hanes (1932) *Biochem. J.* **26**, 1406–1421.
5. J. E. Dowd and D. S. Riggs (1965) *J. Biol. Chem.* **240**, 863–869.
6. W. W. Cleland (1979) *Meth. Enzymol.* **63**, 103–138.

## Linkage Disequilibrium

When the frequency of **genotypes** at more than two **loci** can be expressed as the cumulative product of the respective **allele** frequencies of each locus, those genes are said to be at the state of “linkage equilibrium.” Suppose that two different loci, A and B, each have two alleles,  $A_1$  and  $A_2$ , and  $B_1$  and  $B_2$ , respectively. Let us denote frequencies of alleles  $A_1$  and  $A_2$  at locus A in a population as  $x_1$  and  $x_2$ . Similarly, let us denote the allele frequencies at locus B as  $y_1$  and  $y_2$ . The genotypic frequency,  $X_1$ , of  $A_1 B_1$  (viz., allele  $A_1$  at locus A and allele  $B_1$  at locus B) can be expressed as  $X_1 = x_1 y_1$  when those genes are at linkage equilibrium. When this relationship does not hold for some reason, those genes are said to be in a state of “linkage disequilibrium.” Therefore, the linkage disequilibrium,  $D$ , can be measured by  $D (X_1 - x_1 y_1)$ . When  $D = 0$ , linkage equilibrium exists (Fig. 1). Linkage disequilibrium usually occurs when a particular combination of alleles (a **haplotype**) is favored by or when two loci are located so closely on a chromosome that between those loci rarely takes place. It can also occur when equilibrium has not yet been reached in a population, or when there is nonrandom mating. When sampling in a finite population is nonrandom and produces an abundance of a particular haplotype, linkage disequilibrium will also be apparent. Favorable combinations of alleles will also be transmitted together at greater frequencies when they enhance reproduction and survival, leading to linkage disequilibrium.

**Figure 1.** Illustration of the calculation of linkage disequilibrium.

| Locus A   | Locus B   | Haplotype | Haplotype frequency | Haplotype frequency at equilibrium |
|---|---|-----------|---------------------|------------------------------------|
|  |  | $A_1 B_1$ | $X_1$               | $x_1 y_1$                          |
|  |  | $A_1 B_2$ | $X_2$               | $x_1 y_2$                          |
|  |  | $A_2 B_1$ | $X_3$               | $x_2 y_1$                          |
|  |  | $A_2 B_2$ | $X_4$               | $x_2 y_2$                          |

Linkage Disequilibrium:  $D = X_1 - x_1 y_1$ ,  $-D = X_2 - x_1 y_2$ ,  $-D = X_3 - x_2 y_1$ ,  $D = X_4 - x_2 y_2$



Mechanisms that restrict recombination, such as asexual reproduction, promote the continuance of linkage disequilibrium, leading to the domination of a limited number of alleles. Linkage disequilibrium is also maintained through chromosomal translocations and inversions that reduce recombination. When natural selection favors linkage disequilibrium, chromosomal rearrangements will also increase the linkage. The so-called supergenes are closely linked genes that affect one or several related traits that have arisen through linkage disequilibrium.

With the advancement of technology, it will become possible to detect single-nucleotide **polymorphisms**, (SNP; pronounced “Snip”). By detecting SNP in the noncoding region of the gene of interest, one is able to detect which gene is responsible for a particular disease by checking if the linkage disequilibrium exists between a particular disease and SNP in a target population. This is one of the latest examples of utilization of linkage disequilibrium.

#### Suggestion for Further Reading

“Linkage Disequilibrium” in , Vol. 3, pp. 1393–1394, by T. Gojobori; “Linkage Disequilibrium” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

M. Nei, (1987). *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.

## Linkage, Genetic

Two pairs of genetic markers are said to be linked when the [recombination](#) frequency between them is significantly less than 50%. The recombination frequency in a **diploid** that is heterozygous for two pairs of **alleles** is the proportion of haploid recombinants observed to the total number of haploid products examined. The genes, loci, and DNA sites affected by such alleles are also said to be linked (see [Genetic Marker](#)). Linkage is assumed to be proof of physical location on the same chromosome (see [Recombination](#); [Mitotic Recombination](#)).

Loci that exhibit recombination frequencies of 50% are said to be independent. Loci on different chromosomes give recombination frequencies of 50% because of the random assortment of chromosomes during meiosis; the same is true for random chromosome losses during the [parasexual cycle](#) and for the random distribution of genomic nucleic acid molecules in **viruses** that contain several different ones.

The recombination of alleles in two homologous DNA molecules will be observed if not just one recombination event, but any odd number of these events, occurs in the DNA segment that separates them. An even number of recombination events between them will maintain the parental configuration. The consequence is that the recombination frequency will approach 50% as the number of recombination events increases. Loci of the same chromosome will segregate independently if they are sufficiently distant and the recombination activity is sufficiently high.

Linkage disturbs Mendelian segregations (see [Mendelian Inheritance](#)) as shown in Table 1. The relationship between the recombination frequency and the observed phenotype frequencies is rather complicated in the case of an  $F_2$  generation, but relatively simple in a test cross of  $F_1$  individuals with tester individuals that are homozygous for the recessive alleles at all the loci under consideration.

**Table 1. Segregations of Linked Alleles<sup>a</sup>**

|                      |   |   |   |
|----------------------|---|---|---|
| Parents:             | $AB/AB$   | × | $ab/ab$   |
| Phenotypes:          | “AB”  |   | “ab”  |
| Their gametes:       | $AB$  |   | $ab$  |
| Generation F1:       |   |   | $AB/ab$   |
| Phenotype:           |   |   | “AB”  |
| Their gametes:       |   |   | $(\frac{1-r}{2})AB + \frac{r}{2}Ab + \frac{r}{2}aB + (\frac{1-r}{2})ab$                                       |
| Generation F2:       |   |   | $((\frac{1-r}{2})AB + \frac{r}{2}Ab + \frac{r}{2}aB + (\frac{1-r}{2})ab)^2$                                   |
| Their phenotypes:    |   |   | $(\frac{1}{2} + (\frac{1-r}{2})^2)“AB” + \frac{r(2-r)}{4}“Ab” + \frac{r(2-r)}{4}“aB” + (\frac{1-r}{2})^2“ab”$ |
| Testcross of the F1: | $AB/ab$   | × | $ab/ab$   |
| Phenotypes:          | “AB”  |   | “ab”  |
| Their gametes:       | $(\frac{1-r}{2})AB + \frac{r}{2}Ab + \frac{r}{2}aB + (\frac{1-r}{2})ab$ |   | $ab$  |
| Generation R2:       |   |   | $(\frac{1-r}{2})AB/ab + \frac{r}{2}Ab/ab + \frac{r}{2}aB/ab + (\frac{1-r}{2})ab/ab$                           |
| Their phenotypes:    |   |   | $(\frac{1-r}{2})“AB” + \frac{r}{2}“Ab” + \frac{r}{2}“aB” + (\frac{1-r}{2})“ab”$                               |

<sup>a</sup>  $A$ ,  $a$  and  $B$ ,  $b$  are two pairs of linked alleles. Capitals  $A$  and  $B$  are used to designate alleles with dominant phenotypes, “A” and “B,” respectively; and lower case  $a$  and  $b$  are used for alleles with recessive phenotypes, “a” and “b,” respectively. The coefficient  $r$  is the probability that a gamete from a doubly heterozygous diploid is a recombinant gamete. The percentage value is used in the construction of genetic maps.

Complete linkage—that is, zero recombination—is observed for loci located on the same chromosome in the absence of mechanisms to recombine chromosome pieces. This is the case in the males of *Drosophila melanogaster*, as well as in RNA viruses.

The name *linkage* was proposed by T. H. Morgan (1) for a phenomenon that had already been observed as a disturbance of Mendelian segregations.

#### Bibliography

1. T. H. Morgan (1910) Am. Nat. **44**, 449.

#### Linker DNA

Linker DNA is the DNA connecting chromatosomes in [chromatin](#). Its length can be taken to be the difference between the number of DNA base pairs in the [nucleosome](#) repeat unit and the 166 bp in the chromatosome. This varies from essentially “zero” in the yeast *Saccharomyces cerevisiae* and mammalian cortical neurons (which have a repeat length of ~166bp) to ~74bp in sea urchin sperm (repeat length ~240bp). Linker DNA is complexed with the basic C-terminal (1), and possibly N-terminal, tail of [histone](#) H1 and its variants; consequently, H1 is often designated as a “linker

histone”). The details of the interaction are unclear, but it must be largely **electrostatic** (basic protein side chains interacting with DNA phosphate groups) and is important in promoting condensation of the linker DNA in the folding of the nucleosome filament into higher-order structure (see [Chromatin](#)). The linker is relatively accessible, compared with the DNA sequestered within the body of the nucleosome, and is preferentially cleaved by endonucleases, such as micrococcal nuclease (**Staphylococcal nuclease**), that are used to release nucleosomes (and oligomers thereof) from nuclear chromatin for biochemical and biophysical study. [Transcription factor](#) binding sites in the linker region, for example between positioned nucleosomes (see [Chromatin](#)), are likewise generally more accessible than those in DNA wrapped around the octamer. The repeat length, and hence linker length, can vary throughout the [genome](#), as shown by **Southern blotting** using various probes of the “ladder” of DNA fragments extracted from a micrococcal nuclease digest of chromatin (2).

The nature of the path taken by linker DNA—whether bent or straight—is somewhat controversial, as are models for chromatin higher-order structure (see [Chromatin](#)) to which this feature is relevant. The solenoid and related models, in which the nucleosome filament is helically coiled with increasing ionic strength *in vitro*, require a bent linker to allow nucleosomes to pack together in space; other models invoke straight linkers. The models and evidence have been reviewed (3, 4). [Electron microscopy](#) suggests that the linkers may be straight in the extended nucleosome filament at low ionic strength, but linkers are not visible in the condensed structure and could well be bent. The roughly 10-nucleotide pattern of thymidine dimer formation observed for core DNA, and attributed to its curvature around the octamer, is absent for linker DNA, for which the pattern is relatively uniform (5); however, this should not be taken as definitive evidence for a straight linker rather than a differently curved linker, or a curved linker whose reactivity is altered by association with H1. Studies of the **sedimentation** properties of dinucleosomes, which become more compact with increasing ionic strength, are not readily compatible with a straight linker (6, 7). None of this evidence is absolutely conclusive, however, and definitive evidence from long polynucleosomes is badly needed. Analysis of all the available information on (average) nucleosome repeat lengths, long and short, has revealed that they are quantized, meaning that the linker DNA lengths are quantized (8). The values are related to each other by integral multiples of the helical twist of DNA (~9.5 to 10.5 bp/turn), suggesting that the basis for this is the requirement for definite protein–DNA and [protein–protein interactions](#) in a higher-order structure. Quantization would be expected for the solenoid model, but would also be compatible with the straight linker models.

Histone H1 is clearly a major factor in establishing a regular length of linker DNA in a reconstituted nucleosome array, and it presumably also plays this role *in vivo*. “Chromatin” reconstituted *in vitro* by mixing histones at high ionic strength (eg, 2 M NaCl) and then dialyzing to low ionic strength shows close packing of octamers, irrespective of the presence of H1, which has no effect because it is the last histone to bind in this procedure. Cell-free extracts from *Xenopus* eggs and oocytes and from *Drosophila* embryos will assemble plasmid DNA, irrespective of sequence, into “physiologically” spaced chromatin in an ATP-dependent fashion (see [Chromatin](#)), and the nucleosome spacing (and hence linker length) is increased by H1. Spacing may also have an electrostatic component, involving neutralization of charges on linker DNA (9), and this has been taken to suggest a connection (possibly the basis of quantization; see text above) between nucleosome spacing and the formation of higher-order structure, which is also ionic strength-dependent (10). A well-defined *in vitro* system starting with histones and DNA, and using polyglutamic acid, will also assemble “properly spaced” chromatin, in an H1-dependent, ATP-independent manner (11). In this case, formation of regular arrays is sensitive to (nucleosome positioning?) signals in the DNA.

## Bibliography

1. A. Hamiche et al. (1996) Linker histone-dependent DNA structure in linear mononucleosomes. *J. Mol. Biol.* **257**, 30–42.
2. B. Villeponteau, J. Brawley, and H. G. Martinson (1992) Nucleosome spacing is compressed in active chromatin domains of chick erythroid cells. *Biochemistry* **31**, 1554–1563.

3. V. Ramakrishnan (1997) Histone H1 and chromatin higher-order structure. *Crit. Rev. Eukaryot. Gene Expr.* **7**, 215–230.
4. J. Widom (1998) Structure, dynamics and function of chromatin *in vitro*. *Annu. Rev. Biophys. Biomol. Struct.* **27**, 285–327.
5. J. R. Pehrson (1995) Probing the conformation of nucleosome linker DNA *in situ* with pyrimidine dimer formation. *J. Biol. Chem.* **270**, 22440–22444.
6. J. Yao, P. T. Lowary, and J. Widom (1990) Direct detection of linker DNA bending in defined-length oligomers of chromatin. *Proc. Natl. Acad. Sci. USA.* **87**, 7603–7607.
7. P. J. G. Butler and J. O. Thomas (1998) Dinucleosomes show compaction by ionic strength, consistent with bending of linker DNA. *J. Mol. Biol.* **281**, 401–407.
8. J. Widom (1992) A relationship between the helical twist of DNA and the ordered positioning of nucleosomes in eukaryotic cells. *Proc. Natl. Acad. Sci. USA* **89**, 1095–1099.
9. T. A. Blank and P. B. Becker (1995) Electrostatic mechanism of nucleosome spacing. *J. Mol. Biol.* **252**, 305–313.
10. D. J. Clark and T. Kimura (1990) Electrostatic mechanism of chromatin folding. *J. Mol. Biol.* **211**, 883–896.
11. K. Liu and A. Stein (1997) DNA sequence encodes information for nucleosome array formation. *J. Mol. Biol.* **270**, 559–573.

## Linker Scanning

The study of regions involved in the regulation of gene expression typically involves analysis of nested sets of genetic deletions that enter the region of interest from both [upstream](#) and [downstream](#) directions. In many cases, however, the region to be studied is too large to be readily analyzed by substituting individual bases. In such cases, it may be preferable to use “linker scanning” mutagenesis to search within the regulatory region and locate sequences that are particularly important. This method was originally developed by McKnight and Kingsbury ([1](#)), in order to study the transcriptional control signals of the **thymidine kinase (tk)** gene of **Herpes simplex virus**. In preference to mutating separately each of the 50–100 nucleotides that composed the control sequence of interest, they developed the “linker scanning” method, to introduce clustered sets of point mutations at desired locations. In this specific example, a set of point mutations was constructed to scan systematically across a region of DNA known to harbor the control elements of the tk gene. McKnight and Kingsbury then tested the retention of transcriptional competence using a [microinjection](#) assay in individual cells and were able to identify three distinct regions upstream of the tk gene that are required for *in vivo* gene expression.

More recent refinements of this technique have been described in some detail ([2-4](#)). The most commonly used method requires a unique **restriction site** adjacent to the region being mutagenized, and it uses complementary oligonucleotides. The sequence of interest is initially **cloned** into a plasmid **vector**. The plasmid is linearized, and a nested series of 5' and 3' deletion mutations is created using [restriction enzymes](#). The fragments generated are ligated together with complementary oligonucleotides, filling in the gap between the sequence of interest and the nearby restriction site. A series of these mutants is created in order to scan the site of interest. An alternative procedure utilizes [site-directed mutagenesis](#) procedures to introduce smaller clusters of point mutations throughout the region being studied ([2](#)).

## Bibliography

1. S. L. McKnight and R. Kingsbury (1982) *Science* **217**, 316–324.
2. J. M. Green (1996) In F. M. Ausubel, R. Brent, R. E. Kingston, D. D. Moore, et al. (eds.) *Current Protocols in Molecular Biology*, Vol. **1**, Wiley, New York, Suppl. 2, pp. 8.4.1–8.4.7.
3. P. P. Harlow, G. M. Hobson, and P. A. Benfield (1996) *Methods Mol. Biol.* **57**, 287–295.
4. G. M. Hobson, P. P. Harlow, and P. A. Benfield (1996) *Methods Mol. Biol.* **57**, 279–285.

## Linking Number Of DNA

The linking number,  $Lk$ , is the defining property of DNA elementary topological domains. Generally speaking,  $Lk$  is a measure of the total number of complete revolutions that either strand makes about the other. As long as the strands remain intact,  $Lk$  is a fixed quantity. Two otherwise identical closed duplex DNAs (*cdDNA*) that differ only in  $Lk$  are termed *topoisomers*. Even though they contain exactly the same nucleotide sequence and covalent connections, the properties of the members of a family of topoisomers depend strongly on  $Lk$ . The linking number of DNA can be described formally in several different ways, all of which may be readily generalized to any two closed curves in space.

### 1. Properties of the Linking Number

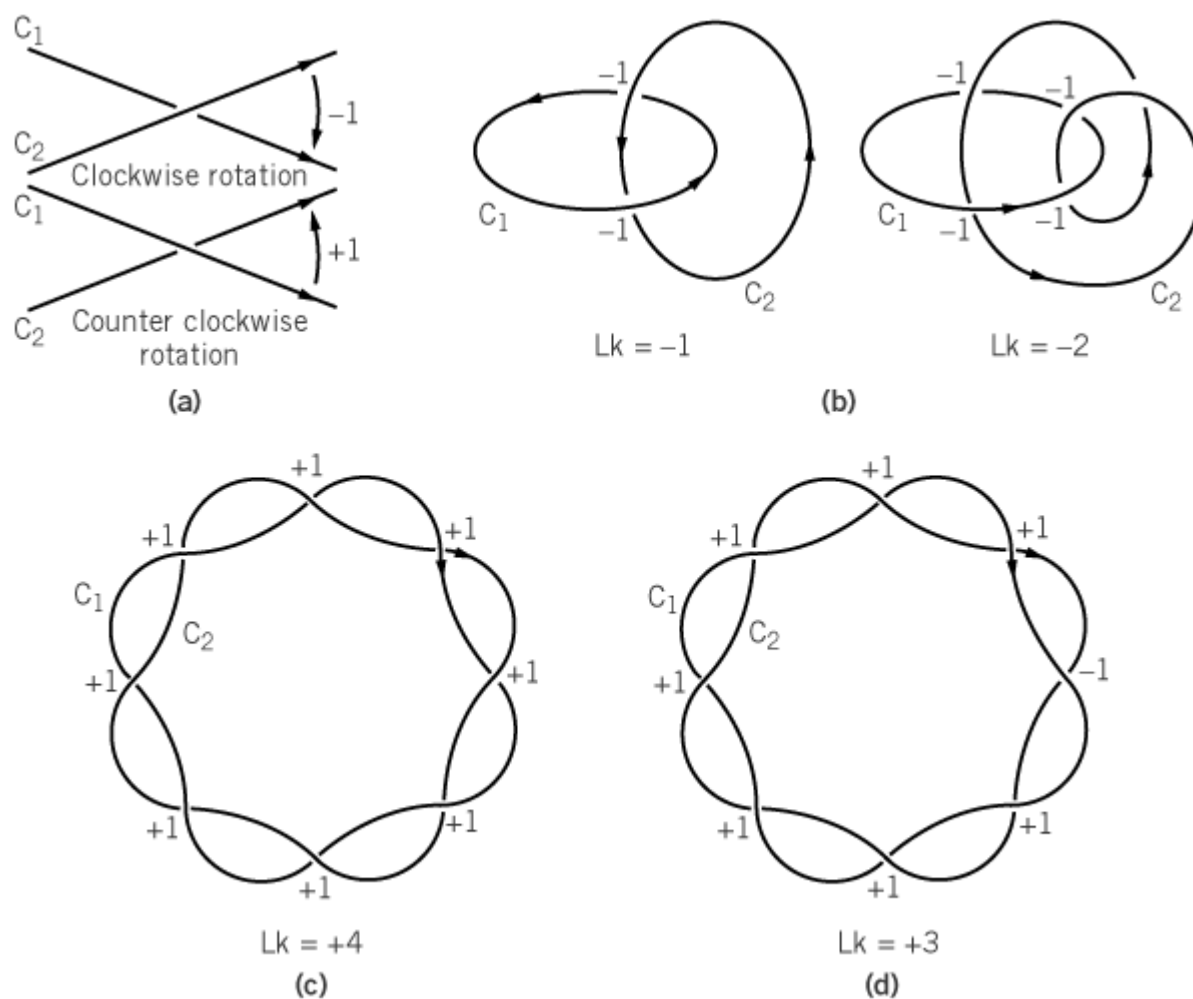
The linking number has several simple and highly useful attributes. (1)  $Lk$  is an integer for superhelical DNA (but not necessarily for protein-sealed DNA loops)—this follows from the requirement that both strands be closed curves; (2) in general,  $Lk$  is constant as long as the topological domain remains intact—a topological domain may be broken by a single- or double-stranded DNA chain scission or, if appropriate, by the disruption of links to the protein sealing the domain; (3) the linking number is a topological quantity, and its value is independent of DNA geometry—that is,  $Lk$  does not vary with deformation of the trajectory of either strand or with changes in the characteristic duplex geometric quantities (pitch, roll, twist, tilt, propeller twist, etc); (4) the linking number is independent of the ordering of the two curves; thus, for two DNA strands  $C_1$  and  $C_2$ ,  $Lk(C_1, C_2) = Lk(C_2, C_1)$ —this independence of ordering clearly distinguishes the linking number from the **twist**; and (5) finally, the double-helix structure of DNA defines a duplex axis ( $A$ ) that is considered to be a continuous line along the center of the helix. Since either of the two backbone chains can be continuously deformed into the axis, the linking number of either strand about the axis is the same as the linking number of the strands about one another. Thus, for example,  $Lk(C_1, A) = Lk(C_2, A) = Lk(A, C_1) = Lk(A, C_2)$ , and so on. The validity of these equalities requires, of course, that the duplex axis be everywhere defined. If it is not (eg, a *cdDNA* containing regions of local denaturation or of cruciform extrusion), a more extensive treatment might be required (1).

### 2. Formal Definition and Calculation of the Linking Number

#### 2.1. The Index Approach

Conceptually, the easiest method of determining  $Lk$  is the index approach. This method relies on the criteria that the two strands must wind about each other once for each contribution to  $Lk$ , and that the sign of each local contribution depends on the relative orientation of the curves. A description of how to assign index numbers is presented in Figure 1. The linking number does not depend on the perspective of the observer, so the DNA may be viewed in projection onto any convenient plane. The strands are assigned parallel orientations, in order to be consistent with the conventions used in the mathematical literature.

**Figure 1.** Calculation of the linking number using the index approach. (a) Assignment of strand orientations and index numbers at projection crossings. The curves are oriented in a parallel sense. Clockwise rotation of the tangent of the upper curve to coincide with the tangent of the lower curve at the intersection point contributes  $-1$ ; counterclockwise rotation of the tangent of the upper curve to coincide with the tangent of the lower curve at the intersection point contributes  $+1$ . (b) Two examples of the use of the index approach to calculate  $Lk$ . In each case  $Lk$  is one-half the sum of the index numbers. (c) Representation two oriented curves that intersect eight times in projection. Each index number is  $+1$ , giving  $Lk = +4$ . (d) Calculation of  $Lk$  for a case in which some of the intersections cancel. This structure, which appears similar to that of (c), has  $Lk = +3$ .



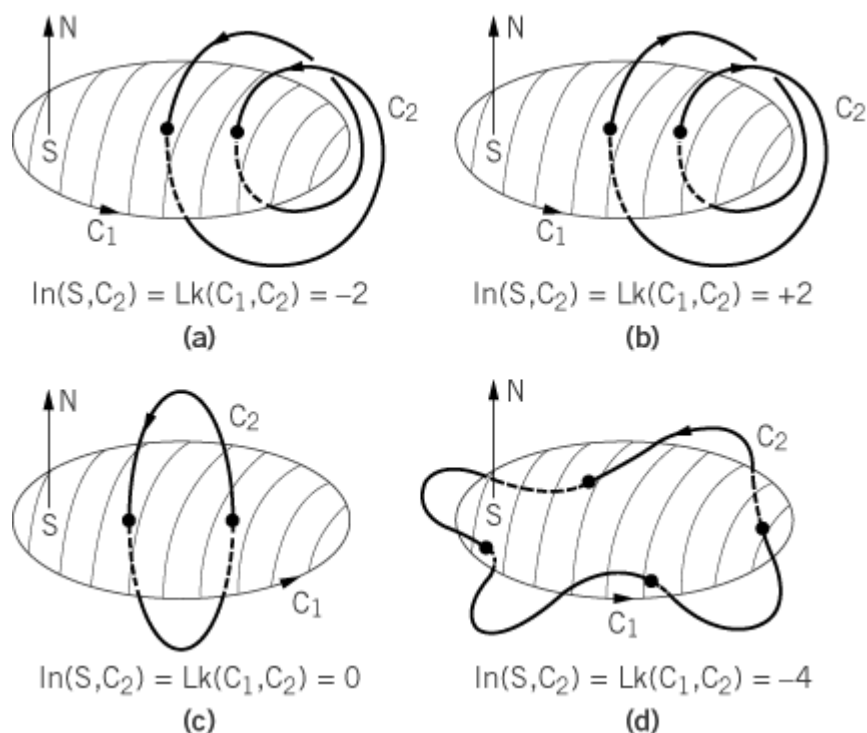
The easiest way to construct the backbone chain curves  $C_1$  and  $C_2$  is to draw a line along one strand with its arrow pointing  $5' \rightarrow 3'$ , and along the complementary strand with its arrow pointing  $3' \rightarrow 5'$ . (This parallel orientation should not be confused with the polarities of the strands, which are, of course, antiparallel.) According to this convention, each interstrand crossing in the projection is assigned an index value of  $+1$  if the locally upper strand must be rotated counterclockwise so as to make its tangent vector coincide with that of the locally lower strand. An index value of  $-1$  is assigned for the reverse case. Then  $Lk$  is one-half the sum of all these index values, since each local contribution to  $Lk$  requires two crossings to advance one turn.

## 2.2. The Surface Intersection Approach

The second method of determining  $Lk$  is the surface intersection approach (1). Here either the DNA axis or one of the two strands is considered to form a spanning surface, which the other strand then repeatedly intersects. For an illustration of how to use the spanning surface to calculate  $Lk$ , refer to

Figure 2. The surface is assigned an orientation; as shown in the figure, the surface normal is taken to be upward-pointing with the curve oriented counterclockwise, equivalent to following a right-hand rule. To each intersection is assigned an index of +1 or -1, depending on the direction of the intersection with respect to the surface orientation. By convention, the index number is positive if the tangent vector of the intersecting curve points along the direction of the surface normal and is negative for the opposite case.  $Lk$  is then the sum of these index numbers. An easily constructed spanning surface, as depicted in the figure, is that whose perimeter is formed by the axis of the closed circular DNA. For relaxed circular DNA, for example, this surface is approximately bounded by a circle. The surface intersection approach may also be employed with the other strand used to form the spanning surface. This choice of spanning surface more readily lends itself to calculation of  $Lk$  for complex cases in which the DNA axis is poorly defined locally. This occurs, for example, in the extrusion of a cruciform. As noted above, the linking number is the same for interstrand winding as for strand-axis winding.

**Figure 2.** Construction and use of the spanning surface in calculation of the linking number. The mathematical convention is followed in choosing the orientation of the curves; the arrows are pointed in the 5' → 3' direction for one strand and in the 3' → 5' direction for the other. The arrows are therefore oriented parallel to one another. One of the two strands is chosen to form the spanning surface, and direction of the normal to the surface is taken in the sense of the classic Stokes theorem; specifically, the right-hand rule is followed. Then the local contribution to the intersection number is +1 if the tangent to  $C$  at the intersection point is parallel to the surface normal. The contribution is -1 if the tangent to  $C$  at the intersection point is antiparallel to the surface normal. (a) Here the spanning surface is intersected twice in the antiparallel sense, and  $Lk = -2$ . (b) Here both intersections are in the parallel sense, and  $Lk = +2$ . (c) Here the two intersections are in the opposite sense and  $Lk = +1 - 1 = 0$ . (d) Here the surface is intersected four times in the antiparallel sense and  $Lk = 4$ .



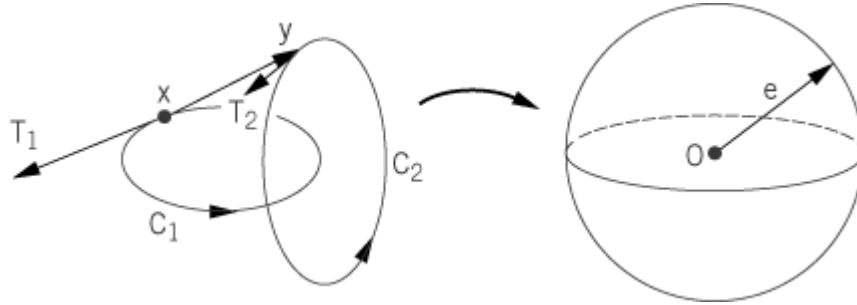
### 2.3. The Gauss Integral Approach

The third and most formal method of calculating  $Lk$  employs the Gauss integral. The integral is taken over all pairs of points on the two strands, which are again labeled  $C_1$  and  $C_2$ . A graphical description of the geometric quantities involved is given in Figure 3. The linking number is calculated from the integral (2).

$$Lk = \frac{1}{4\pi} \int_{C_1} \int_{C_2} \frac{\hat{\mathbf{e}} \cdot (\mathbf{T}_2 \times \mathbf{T}_1)}{r^2} ds_1 ds_2 \quad (1)$$

where  $x$  is an arbitrary location on strand  $C_1$ ,  $y$  is an arbitrary point on strand  $C_2$ ;  $T_1$  is the unit tangent vector to  $C_1$  at  $x$ ; and  $T_2$  is the unit tangent vector to  $C_2$  at  $y$ . The two curves are connected between any pair of these points by the vector  $y-x$ , whose scalar length is  $r = |y-x|$ . Then  $\hat{\mathbf{e}}$  is the unit vector along  $r$ , and  $\hat{\mathbf{e}} = (\mathbf{y} - \mathbf{x})/r$ . The quantities  $s_1$  and  $s_2$  are the arc lengths on  $C_1$  at  $x$  and on  $C_2$  at  $y$ . Because of its complexity, and because of the need for knowledge of the exact chain trajectory, the Gauss integral is impractical for the calculation of  $Lk$  for DNA. It is, however, useful for theoretical modeling purposes.

**Figure 3.** The Gauss map for calculation of the linking number. Here  $x$  is an arbitrary point on strand  $C_1$  and  $y$  is an arbitrary point on strand  $C_2$ .  $T_1$  is the unit tangent vector to strand  $C_1$  at  $x$ , and  $T_2$  is the unit tangent vector to strand  $C_2$  at  $y$ . The unit vector  $\hat{\mathbf{e}} = (\mathbf{y} - \mathbf{x})/r$ , where  $r = |y-x|$ . If the unit vector  $\hat{\mathbf{e}}$  is translated parallel to the origin, its terminus becomes a point on a sphere of unit radius, centered at the origin. Each pair of  $(x, y)$  points therefore maps to a unique point on this unit sphere. The Gauss integral for  $Lk$  measures how many times the  $\hat{\mathbf{e}}$  vector sweeps across the surface of this sphere in a positive or negative sense. The factor  $1/4\pi$  is the surface area of a unit sphere and normalizes the result to give the number of turns.



### 3. Linking Difference and Superhelix Density

#### 3.1. Linking Difference

The linking difference applies to an elementary topological domain (see [DNA Topology](#)) and is a measure of the deviation in the winding of a cdDNA from that of its nonclosed circular counterpart. Consequently it is  $\Delta Lk$ , rather than just the linking number,  $Lk$ , that determines the chemically and biologically interesting properties of a topological domain. The linking difference is defined by the difference between  $Lk$  and the value of the pseudo-linking number for the corresponding nicked circular or linear DNA,  $Lk_0$ . These latter species are termed, collectively, open duplex DNA:

$$\Delta Lk = Lk - Lk_0 \quad (2)$$

In contrast to  $Lk$ ,  $Lk_0$  is seldom an integer, even for completely purified cdDNA.  $Lk_0$  is related to the number of base pairs,  $N$ , and to the DNA helical repeat,  $h_0$  base pairs per turn, by

$$Lk_0 = \frac{N}{h_0} \quad (3)$$

The local value of the [B-DNA](#) helical repeat is dependent on the base composition (3) and is different for noncanonical DNA structures, such as [Z-DNA](#), H, and locally denatured DNAs. For the linear axis B form of DNA in dilute NaCl at 37°C, the average  $h_0 = 10.5$  base pairs/turn (4, 5).



### 3.2. Relaxed Duplex DNA

A nicked circular duplex DNA contains at least one backbone chain interruption. Such a DNA loses its topological domain, and free rotation about the chain scission takes place. In contrast to a nicked circular DNA, a relaxed duplex DNA (*rdDNA*) is that topoisomer whose linking number is the nearest integer to  $Lk_0$ . Since  $Lk_0$  itself is generally not an integer, the rdDNA species usually has a value of  $DLk$  that is small and fractional. This fractional displacement  $\epsilon$  is simply the difference between the exact value of  $Lk_0$  from Equation (3) and the nearest integer, and  $-0.5 \leq \epsilon \leq 0.5$ . For example, pBR322 DNA has  $N = 4363$ ,  $Lk_0 = 415.52$ , and consequently  $\epsilon = -0.48$ . The rdDNA topoisomer is therefore the one with  $Lk = 416$ . Physically,  $\epsilon$  represents the minimum fractional rotation that is required in a nicked circular DNA, whose axis initially lies in a plane, to bring the 5' and 3' ends together in order to allow covalent joining by a [DNA Ligase](#). The sign and magnitude of  $\epsilon$  can be determined experimentally from the Boltzmann distribution of topoisomers that forms on thermal equilibration in the presence of a **topoisomerase**. This distribution is well fit by a Gaussian curve in topoisomer frequency versus  $DLk$ , and  $\epsilon$  is the separation between the center of this distribution and the location of the nearest, most prominent topoisomer (6, 7) (see [Superhelical DNA Energetics](#)). A closed DNA is *underwound* if  $DLk < \epsilon$ , which is the case for all naturally occurring closed circular DNAs. The linking difference is consequently sometimes called the *linking deficiency*, although linking excess would be the appropriate term for cases in which the DNA is *overwound*, or  $DLk > \epsilon$ . The DNA is said to be *relaxed* if  $DLk = \epsilon$ , even though it might still contain a small linking difference.

Both the sign and the magnitude of  $DLk$  can be altered by changing either  $Lk$  or  $Lk_0$ . The former can be changed by treatment with DNA topoisomerase, DNA gyrase, or reverse DNA gyrase; the latter can be changed by altering the salt type or concentration (8), the temperature (6, 9), or by the addition of an intercalating drug (10). Thus, a relaxed duplex DNA can be made to **supercoil** by a change in its environment alone. For comparison purposes, it is therefore necessary to specify the conditions under which  $Lk_0$  is measured. The usual standard state is 37°C, 0.2 M NaCl, and no added reagents that change the DNA twist (11). If necessary, the value of  $Lk_0$  under standard conditions can be specified as  $Lk_0^0$ .

### 3.3. Superhelix Density

All other things being equal,  $Lk$ ,  $Lk_0$ , and  $DLk$  are proportional to the DNA length. In order to compare two DNA molecules of different lengths, it is convenient to define the associated normalized quantity, the superhelix density or specific linking difference,  $s$ . The superhelix density is defined by

$$\sigma = \frac{Lk - Lk_0}{Lk_0} \quad (4)$$

Although the term “superhelix” is traditionally used here, it should not be taken to imply any specific tertiary structure (see [DNA Topology](#)). Most naturally occurring cdDNAs have values of  $s$  between 0 and  $-0.1$  under standard conditions, but DNAs having  $s$  as great as  $-0.17$  have been prepared *in vitro* (12). cdDNA of relatively large positive  $s$  values have been prepared with the DNA reverse gyrase (13). Since both  $s$  and  $Lk_0$  are readily accessible experimental quantities, the practical way to calculate  $Lk$  is by combining  $Lk_0$  with the measured value of either  $DLk$  or  $s$ . To continue the example of pBR322 DNA, with  $N = 4363$  and  $h_0 = 10.5$  base pairs/turn, the DNA occurs naturally with an average value of  $s = -0.06$ . Taking this number to be exact for purposes of illustration, the *average* linking number for this native DNA is then 390.58. The linking number of the nearest (most prevalent) topoisomer is 391.

As with  $DLk$ ,  $s$  can be varied by changes in environmental conditions alone, even with no change in

*Lk*. For example, the temperature coefficient of *s* is  $Ds/DT = 3.1 \times 10^{-4} \text{deg}^{-1}$ . If the temperature is changed from 5 to 37°C, the change in *s* is +0.01. For pBR322 DNA, this changes *DLk* by +4.16 turns. An increase in temperature thus causes a relaxed DNA to supercoil in the positive sense (a left-handed, interwound superhelix) and reduces the supercoiling of a naturally occurring (underwound) superhelical DNA. Similar effects result from changing the cation species and concentration (8). Over the ionic strength range 0.05–0.3 M,  $Ds/DpX^+ = 4.47 \times 10^{-3}$  for the ions  $\text{Na}^+$ ,  $\text{K}^+$ ,  $\text{Li}^+$ , and  $\text{NH}_4^+$  and  $Ds/DpX^+ = 6.70 \times 10^{-3}$  for  $\text{Rb}^+$ ,  $\text{Cs}^+$ , and  $\text{Mg}^{2+}$ , where  $pX^+$  is the negative logarithm of the cation concentration. For example, if the sodium ion concentration is decreased from 1.0 to 0.01 M, the change in *s* is +0.013 and is comparable to a temperature increase of 42°C. Both these effects are clearly evident in **gel electrophoresis** experiments under the appropriate different conditions.

#### 4. Relationship of DNA Topology to DNA Geometry

In spite of its being topological (nonmetric), the linking number is equal to the sum of two very important DNA geometric properties (14): the **twist**, *Tw*, and the **writhe**, *Wr*:

$$Lk = Tw + Wr \quad (5)$$

This equation is fundamental to understanding any topological domain in DNA. An immediate consequence of this relationship is that, for any process in which *Lk* is unchanged, any changes in the twist are matched exactly by changes in the writhe (but of opposite sign):

$$\delta Lk = 0 \Rightarrow \delta Tw = -\delta Wr \quad (6)$$

The linking difference also appears in a modified version of the fundamental relationship for a topological domain, Equation (5). A nicked circular DNA has no writhe, so  $Tw_0 = Lk_0$ . Combining this condition with Equation (5) gives

$$\Delta Lk = \Delta Tw + Wr \quad (7)$$

Here  $DTw = Tw - Tw_0$  represents the deviation of the twist from the open circular (B-DNA solution) value. Equation (7) points to two differences between a relaxed duplex DNA and a nicked circular duplex DNA. The first difference is relatively minor. The nicked circular duplex DNA has  $Wr = DTw = 0$ , but for the relaxed duplex DNA  $Wr + DTw = \epsilon$ . That is, a relaxed duplex DNA may be slightly distorted by some combination of changes in *Tw* and *Wr*. Theoretical calculations indicate that all the distortion goes into twist for a perfect elastic rod (15, 16) but that it all goes to writhe for any actual case involving DNA (J. H. White, submitted). The second difference is major. In a rdDNA the twist and writhe remain coupled, such that  $dDTw = -dWr$ . For a nicked circular DNA, however, *D Lk* is not defined and Equation (7) does not apply. In this case, *DTw* and *Wr* are uncoupled and may fluctuate independently. These differences explain why nicked circular DNA often migrates more slowly than relaxed circular DNA in both gel electrophoresis and **sedimentation** experiments, even under identical solution conditions.

#### Bibliography

1. J. H. White and W. R. Bauer (1987) *J. Mol. Biol.* **195**, 205–213.
2. J. H. White (1989) in *Mathematical Methods for DNA Sequences*, M. S. Waterman, ed., CRC Press, Boca Raton, FL, pp. 225–253.
3. L. J. Peck and J. C. Wang (1981) *Nature* **292**, 375–378.

4. D. Rhodes and A. Klug (1980) *Nature* **287**, 573–578.
5. J. C. Wang (1979) *Proc. Natl. Acad. Sci. USA* **76**, 200–203.
6. R. E. Depew and J. C. Wang (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4275–4279.
7. M. Shure and J. Vinograd (1976) *Cell* **8**, 215–226.
8. P. Anderson and W. Bauer (1978) *Biochemistry* **17**, 594–601.
9. F. S. Lee and W. R. Bauer (1985) *Nucleic Acids Res.* **13**, 1665–1682.
10. W. Bauer and J. Vinograd (1971) *Prog. Mol. Subcell. Biol.* **II**, 181–215.
11. W. R. Bauer (1978) *Annu. Rev. Biophys. Bioeng.* **7**, 287–313.
12. C. K. Singleton and R. D. Wells (1982) *Anal. Biochem.* **122**, 253–257.
13. A. Kikuchi and K. Asai (1984) *Nature* **309**, 677–681.
14. J. H. White (1969) *Am. J. Math.* **91**, 693–728.
15. C. J. Benham (1989) *Phys. Rev. A* **39**, 2582–2586.
16. M. Le Bret (1979) *Biopolymers* **18**, 1709–1725.

### Suggestions for Further Reading

17. W. R. Bauer and R. Gallo (1989) *Physical and topological properties of closed circular DNA, in Chromosomes: Eukaryotic, Prokaryotic, and Viral*, K. W. Adolph, ed., CRC Press, Boca Raton, FL, Vol. I, pp. 87–126. (This review article includes descriptions of how to determine the linking number and the superhelix densities and summarizes the literature data on experimental determinations of these quantities.)
18. J. H. White (1992) Geometry and topology of DNA and DNA-protein interactions, *Proc. Symp. Appl. Math.* **45**, 17–37. (This article contains detailed descriptions of the various methods for calculation of the linking number.)

## Lipases

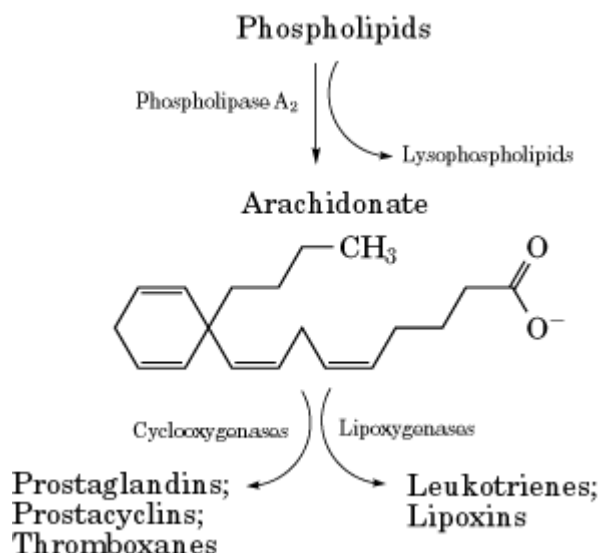
Studies on the involvement of lipids in [signal transduction](#) pathways and in cell regulatory mechanisms have demonstrated the regulation of key [enzymes](#) of [lipid metabolism](#) that produce important lipids and lipid-derived products that mediate signal transduction. This is best illustrated with the [phosphatidylinositol](#) (PI) cycle, where the activation of PI-specific [phospholipases C](#) results in the formation of [diacylglycerol](#) and [inositol trisphosphate](#), both of which function as [second messengers](#). Although this remains the best-studied example, and the anchor for the paradigm, several additional mechanisms of lipid metabolism are now appreciated to participate in significant pathways of cell regulation. Indeed, on the basis of this paradigm, all *regulated* enzymes of lipid metabolism appear to define novel and distinct pathways of lipid-mediated cell regulation. These enzymes include lipases, lipid kinases, and transacylases ([1-4](#)). Other entries deal with PI-specific phospholipases C and [sphingomyelinases](#). This entry will highlight current understanding of lipases and other regulated enzymes of lipid metabolism.

### 1. Phospholipase D

Phospholipase D catalyzes the hydrolytic cleavage of the polar head group of phospholipids (Fig. [1](#)), resulting in the formation of phosphatidic acid and, usually, a free alcohol (such as choline when the substrate is phosphatidylcholine). Mammalian phospholipase D is activated in response to a number of [growth factors](#) and [cytokines](#). Mechanistically, phospholipase D activation requires small **G**

**proteins** (such as ARF or RhoA), and its activation can also be effected through a protein kinase C-dependent mechanism (5-7).

**Figure 1.** Sites of action of phospholipases. The structure of a phospholipid is shown with the sites of action of the main phospholipases (C, D, and A<sub>2</sub>).



Phosphatidic acid, the initial product of phospholipase D action, has been shown to exert growth-promoting activities (7). Its direct targets of action are not clearly defined, but it has been shown to regulate protein [kinases](#), [Gtpases](#), and other intracellular targets. Phosphatidic acid can also be further metabolized to yield lysophosphatidic acid (through a phospholipase A<sub>2</sub> activity) or diacylglycerol (through a **phosphatase** activity). Lysophosphatidic acid has been shown to serve as an intercellular regulator, probably acting on a transmembrane **receptor**, and it has been suggested to participate in **mitogenic** activities (8). Diacylglycerol activates protein kinase C and participates in several signal transduction mechanisms.

In addition, phospholipase D has been implicated in **vesicle** formation and trafficking (6), based on the localization of its regulators (primarily the ARF family of small **G proteins**) in vesicles and their involvement in vesicle budding, transport, and fusion. In addition, yeast phospholipase D plays an important role in **sporulation** of budding yeast.

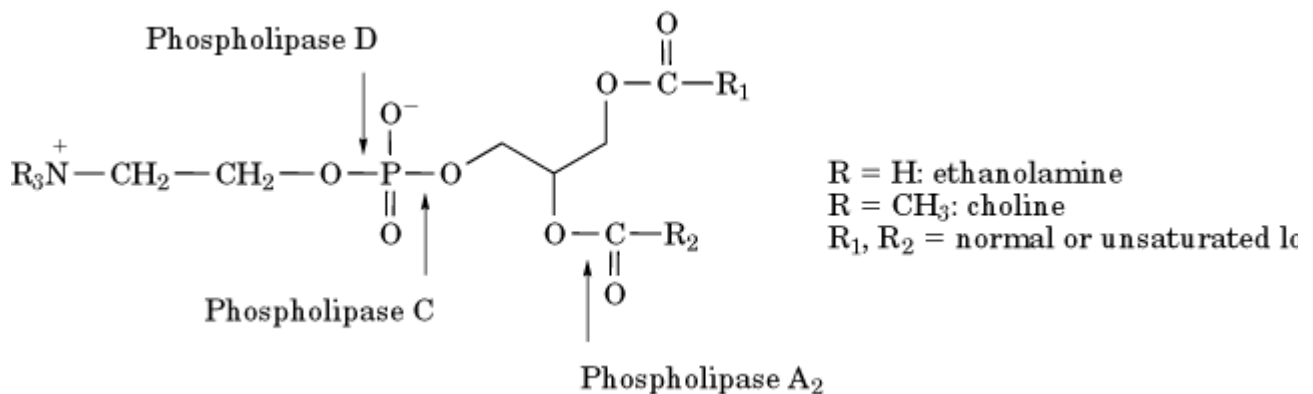
## 2. Phospholipase A<sub>2</sub> and Related Enzymes of Eiconasoid Production

Phospholipase A<sub>2</sub> attacks the acyl group at the *sn*-2 position of glycerophospholipids, which is usually occupied by arachidonate (Fig. 1). Several phospholipases A<sub>2</sub> have been purified and **cloned** (9). These include **secretory** small enzymes, as well as cytosolic regulated enzymes. The best studied of the latter is cPLA<sub>2</sub>, which is regulated by calcium and by **phosphorylation** on serine residues. This enzyme is activated in response to a number of growth factors and cytokines and displays selectivity for arachidonate in the *sn*-2 position. Therefore, this enzyme has been implicated as the major regulated PLA<sub>2</sub> involved in signal transduction. The products of this reaction are lysophospholipid and arachidonate. Lysophospholipids may have cellular functions on their own (such as the effects of lysophosphatidic acid on mitogenesis) or may serve as precursors for platelet activating factor, which has important functions in inflammation and vascular biology. The better-

studied product of the phospholipase A<sub>2</sub> reaction, however, is arachidonate. Multiple biochemical targets have been ascribed to arachidonate, including protein kinase C, **phosphodiesterases**, GTPases, and **ion channels** .

An important function for arachidonate is to serve as a precursor for the formation of eicosanoids (a collective name for products of arachidonate metabolism) (Fig. 2) (3). In particular, cyclooxygenases regulate the formation of prostaglandins, thromboxanes, and prostacyclins. These are very important intercellular mediators in platelet activation, neutrophil biology, inflammatory responses, and vascular biology. On the other hand, lipoxygenases regulate the formation of leukotrienes and lipoxins. These again serve as important intercellular messengers and regulators with important functions in inflammation, allergy, vascular biology, and blood cell activation (3).

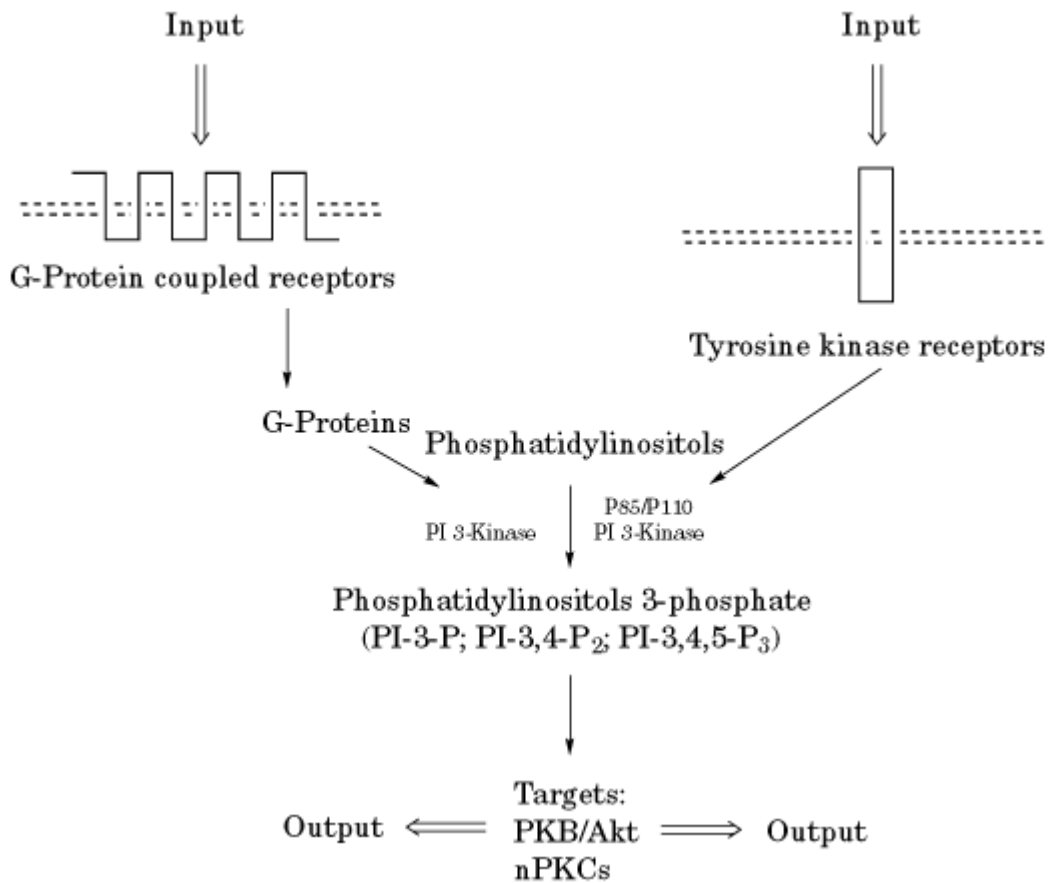
**Figure 2.** The phospholipase A<sub>2</sub> pathway and the eicosanoids. Activation of phospholipase A<sub>2</sub> results in the liberation of fatty acids from phospholipids. Arachidonate can function on its own, or it can be metabolized into two major groups of eicosanoid: either cyclooxygenases or lipoxygenases.



### 3. PI Kinases

A number of kinases that phosphorylate inositol phospholipids on unique positions of the inositol ring have been described and characterized. This results in the formation of products such as phosphatidylinositol-3-phosphate, which has important roles in prevention of cell death and regulation of growth (Fig. 3). The PI-3-kinases that catalyze and regulate this reaction are activated in response to receptor stimulation through either direct tyrosine phosphorylation of regulatory subunits or through heterotrimeric G proteins (4). The phosphorylated inositol phospholipids, in turn, regulate targets involved in signal transduction, such as certain isoforms of protein kinase C and the Akt serine threonine protein kinase. PI-4-kinases, which phosphorylate inositol phospholipids at a distinct site, may play a role in vesicle trafficking.

**Figure 3.** Scheme of activation of PI 3-kinases. G-protein-coupled receptors or tyrosine kinase receptors activate distinct forms of PI 3-kinase, which results in the formation of 3-phosphorylated derivatives of inositol phospholipids.



#### 4. Diacylglycerol Kinases

Diacylglycerol is a substrate for diacylglycerol kinases, whose action results in the formation of phosphatidic acid. These enzymes therefore serve to attenuate the levels of diacylglycerol and enhance the levels of phosphatidic acid, with consequences on pathways regulated by these lipid products.

#### 5. Transacylases

An emerging group of regulated enzymes of lipid metabolism catalyze a transacylation reaction whereby acyl groups are transferred from one lipid to another. The best studied of these transacylases are the enzymes involved in the formation of platelet-activating factor, which has important functions in vascular and endothelial biology ([10](#)).

#### Bibliography

1. E. A. Dennis, S. G. Rhee, M. M. Billah, and Y. A. Hannun (1991) *FASEB J.* **5**, 2068–2077.
2. M. Liscovitch and L. C. Cantley (1994) *Cell* **77**, 329–334.
3. C. N. Serhan (1994) *Biochim. Biophys. Acta* **1212**, 1–25.
4. C. L. Carpenter and L. C. Cantley (1996) *Curr. Opin. Cell Biol.* **8**, 153–158.
5. J. H. Exton (1994) *Biochim. Biophys. Acta* **1212**, 26–42.
6. S. Cockcroft (1996) *Chem. Phys. Lipids* **80**, 59–80.
7. S. C. Olson and J. D. Lambeth (1996) *Chem. Phys. Lipids* **80**, 3–19.
8. W. H. Moolenaar (1995) *J. Biol. Chem.* **270**, 12949–12952.
9. E. A. Dennis (1997) *Trends Biochem. Sci.* **22**, 1–2.

10. Z. Honda, M. Nakamura, I. Miki, M. Minami, T. Watanabe, Y. Seyama, H. Okado, H. Toh, K. Ito, T. Miyamoto, and T. Shimizu (1991) *Nature* **349**, 342–346.

### Suggestions for Further Reading

11. N. Divecha and R. Irvine (1995) Phospholipid signaling. *Cell* **80**, 269–278.
12. R. Graber, C. Sumida, and E. Nunez (1994) Fatty acids and cell signal transduction. *J. Lipid Med. Cell Signal.* **9**, 91–116.
13. M. Liscovitch (1992) Crosstalk among multiple signal-activated phospholipases. *Trends Biochem. Sci.* **17**, 393–399.
14. W. Moolenaar (1995) Lysophosphatidic acid, a multifunctional phospholipid. *J. Biol. Chem.* **270**, 12949–12952.
15. C. Serhan (1994) Lipoxin biosynthesis and its impact in inflammatory and vascular events. *Biochim. Biophys. Acta* **1212**, 1–25.
16. M. Venable, G. Zimmerman, T. McIntyre, and S. Prescott (1993) Platelet-activating factor: a phospholipid autacoid with diverse actions. *J. Lipid Res.* **34**, 691–702.

## Lipid Metabolism

Membrane lipids belong primarily to one of three major classes of [lipids](#): glycerolipids, sphingolipids, and sterols. These small and hydrophobic molecules (molecular weight mostly in the range of 150–2000) display a richness in structural composition that serves as the foundation for their important functions. Lipids are critical and defining components of the [membrane](#) bilayer, such that the **hydrophobic** acyl and alkyl groups constitute the hydrophobic interior of the membrane, whereas their head groups provide a [hydrophilic](#) interface with the aqueous environment on both sides of the membrane. The lipid composition of biological membranes varies between species, cell types, and subcellular compartments. Lipids also participate in the properties and functions of membranes, such as determining their fluidity or rigidity and possibly in determining microdomains in different membranes. In addition, lipids are of significance in ligand recognition, cell–cell interaction, and covalent modification of proteins.

Lipids also function as integral components in diverse [signal transduction](#) pathways, where the action of many extracellular agents is coupled to different enzymes of lipid metabolism. This results in the generation of specific lipid [second messengers](#) (such as **diacylglycerol**, platelet-activating factor, eicosanoids, and ceramide). These molecules then go on to interact with and regulate specific targets (such as protein [kinase C](#) or membrane **receptors**) and thus provide for signal transduction across biological membranes. It is this aspect of lipid function that has come under intense investigation, especially in the last two decades with the elucidation of several signal-activated [lipases](#) and other enzymes of lipid metabolism. As such, the study of regulation of lipid metabolism is at the heart of lipid-mediated signal transduction and cell regulation.

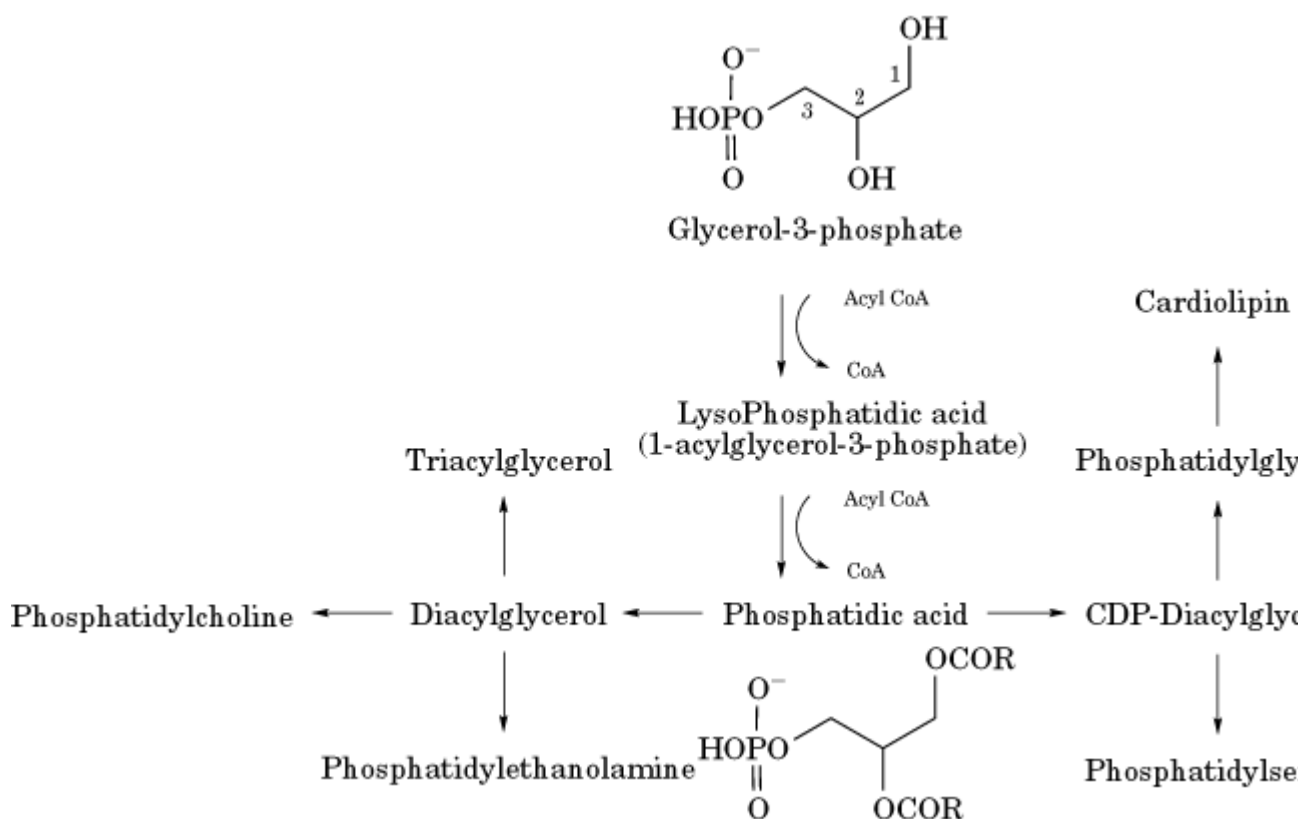
### 1. Basic Pathways of Lipid Metabolism

The metabolic pathways that regulate lipid biosynthesis and lipid composition utilize molecules that are readily available in the cell and follow simple blueprints that show significant similarity among **prokaryotes**, mammalian cells, and other **eukaryotic** systems.

#### 1.1. Glycerolipid Metabolism

In eukaryotic cells, the initial step of phospholipid biosynthesis commences with the acylation of glycerol-3-phosphate using **acetyl-CoA** as the donor and resulting in the formation of 1-acyl-glycerol-3-phosphate (lysophosphatidic acid). This is further acylated to result in the formation of phosphatidic acid, which, in turn, can serve as a precursor for either diacylglycerol or CDP–diacylglycerol. CDP–diacylglycerol serves as the precursor to phosphatidylinositol, phosphatidylserine, phosphatidylglycerol, and cardiolipin, whereas diacylglycerol can serve as the precursor for the synthesis of phosphatidylethanolamine, triacylglycerol, and phosphatidylcholine (Fig. 1). The catabolism of phosphoglycerolipids is achieved through the action of various lipases such as **phospholipase C**, which cleaves off the phospho head group and results in the formation of diacylglycerol. Other lipases also contribute to the breakdown of phospholipids to acetyl groups, glycerol, or glycerolphosphate.

**Figure 1.** The major pathways of glycerolipid biosynthesis in eukaryotes, commencing with glycerol-3-phosphate.



## 1.2. Sphingolipid Metabolism

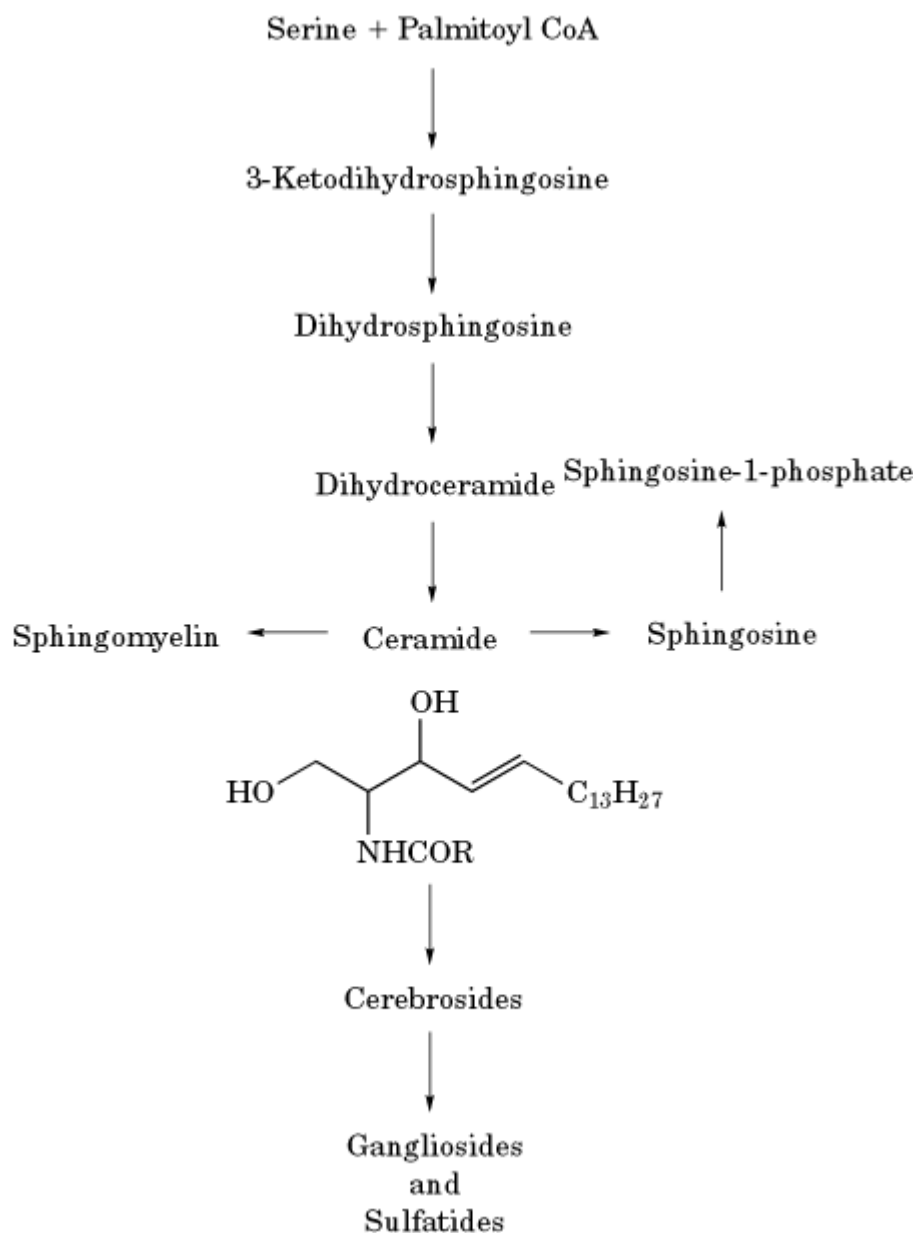
The biosynthesis of sphingolipids commences with the condensation of serine and palmitoyl CoA into 3-ketodihydrosphingosine, which, in turn, is reduced to dihydrosphingosine.

Dihydrosphingosine is acylated to form dihydroceramide, which is either oxidized to ceramide or incorporated into more complex sphingolipids (Fig. 2). The complex sphingolipids are distinguished by their specific substituents at the 1-hydroxyl position of ceramide. For example, sphingomyelin contains phosphorylcholine, whereas cerebroside contains either galactose or glucose. These glycosylated sphingolipids can serve as precursors for more complex neutral and acidic glycolipids. The catabolism of sphingolipids proceeds in a reverse and stepwise fashion through hydrolytic elimination of specific components of the head groups, eventually resulting in the formation of ceramide, which, in turn, is deacylated by the action of a ceramidase to yield sphingosine.

Phosphorylation of sphingosine results in sphingosine-1-phosphate, which is a substrate for a lyase that breaks it down into phosphoethanolamine and hexadecenal.



**Figure 2.** The major pathway of biosynthesis of sphingolipids in eukaryotes, commencing with the condensation of serine and palmitoyl CoA.



## 2. Regulation of Lipid Metabolism

The regulation of intermediary lipid metabolism follows the basic and well-established principles of intermediary metabolism, with the rate-determining steps usually associated with the initial enzymes in the biosynthetic scheme (acylation of glycerolphosphate or condensation of serine and palmitoyl CoA in the biosynthesis of glycerolipids and sphingolipids, respectively).

In addition, lipid metabolism is a subject of multiple mechanisms of regulation that are critical in the determination of the levels of individual lipid precursor substrates and lipid-derived products. This is of particular significance in the myriad pathways of signal transduction. For example, phospholipases C regulate the levels of inositol trisphosphate and diacylglycerol, whereas

phospholipases D regulate the levels of phosphatidic acid (see [Lipases](#)). Phospholipase A<sub>2</sub> regulates the levels of arachidonate and, consequently, the levels of eicosanoids that function as intracellular and intercellular messengers (see [Lipases](#)). Other regulated enzymes of lipid metabolism include lipid kinases, synthases, transacylases, and other specialized enzymes.

#### Suggestions for Further Reading

- M. Berridge (1993) Inositol trisphosphate and calcium signalling. *Nature* **361**, 315–324.
- N. Divecha and R. Irvine (1995) Phospholipid signaling. *Cell* **80**, 269–278.
- M. Liscovitch (1992) Crosstalk among multiple signal-activated phospholipases. *Trends Biochem. Sci.* **17**, 393–399.
- Y. Nishizuka (1992) Intracellular signaling by hydrolysis of phospholipids and activation of protein kinase C. *Science* **258**, 607–614.
- D. Vance (1983) in *Biochemistry*, G. Zubay, ed., Addison-Wesley, Reading, Mass., pp. 505–543.
- D. Vance and J. Vance (1991) *Biochemistry of Lipids, Lipoproteins, and Membranes*. Elsevier, New York.

## Lipids

Lipids are water-insoluble organic compounds that serve a number of essential and diverse roles including intracellular storage of metabolic fuel as triglycerides, structural elements of cell [membranes](#) as phospholipids and cholesterol, protective waxes as fatty acid esters of monohydroxylic alcohols, and substances with intense biochemical activity as the various water- and fat-soluble vitamins. They are usually classified as simple or complex lipids based on whether hydrolysis of the compound yields one or two products, in the former, or more than two, in the latter. Among the simple lipids are cholesterol and its fatty acid esters, triglycerides and [fatty acids](#). There are three major subdivisions of complex lipids: (i) **Glycerophospholipids**, which produce on hydrolysis glycerol, fatty acids, inorganic phosphate, and an organic base or polyhydroxy compound. This group includes the major phospholipids of cell membranes: phosphatidylcholine, phosphatidylethanolamine, phosphatidylinositides, phosphatidylserine, and phosphatidylglycerol. (ii) **Glycoglycerolipids**, which yield glycerol, fatty acids, and carbohydrates. These lipids are found primarily in **plants** and **bacteria**, as well as in small amounts in brain and nervous tissue of some mammals. (iii) **Sphingolipids**, which contain a long-chain base, fatty acids and inorganic phosphate, and carbohydrates. They are present in the membranes of both plants and animals, particularly in brain and nerve tissue. Sphingomyelin is the most abundant member of this group.

#### Suggestion for Further Reading

- A. L. Lehninger (1970) *Biochemistry*, Worth Publishers, New York, Chapter "10". (Present space-filling models of various lipids.)

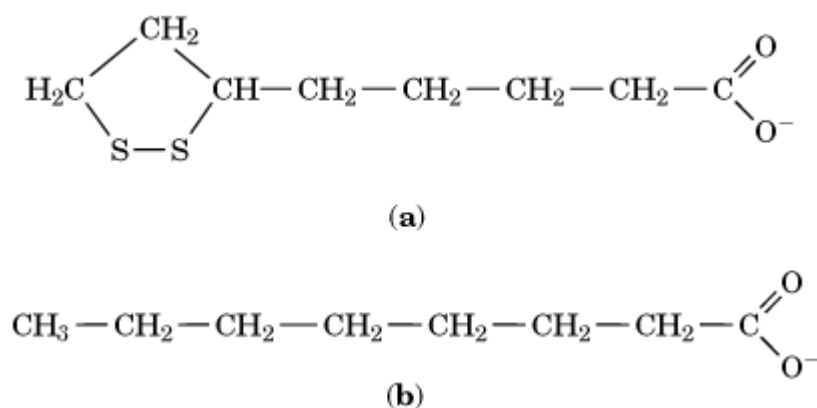
## Lipoic Acid

Lipoic acid (1,2-dithiolane-3-pentanoic acid or, trivially, thiocetic acid) is chiefly known as a protein-bound cofactor in the oxidative decarboxylation of 2-oxo acids. 2-Oxo acid dehydrogenase multienzyme complexes catalyze the oxidative decarboxylation of pyruvate, 2-oxoglutarate, and the branched chain 2-oxo acids that are derived from the transamination of leucine, isoleucine, and valine. In all these instances, the lipoyl group is found in amide linkage with the  $N^6$ -amino group of a specific [lysine](#) residue in the dihydrolipoyl acyltransferase (E2) component, where it acts as a “swinging arm” to ferry the substrate between the three active sites that successively catalyze the overall reaction (1-3). A similar lipoyl-lysine residue occurs in another widespread multienzyme system, the glycine cleavage system, which catalyzes the decarboxylation of glycine (4, 5). All lipoylated [protein structures](#) studied thus far contain an autonomously folded **domain** of about 80 amino acid residues, in which the lipoyl-lysine residue is displayed in a prominent **b-turn** (6-9). The selectivity of the lipoyl protein ligase that catalyzes the lipoylation reaction (10, 11) is unusual, depending in large part on the correct siting of the target lysine residue in the exposed b-turn of the apo-lipoyl domain (12). Free lipoic acid is inactive as a substrate in the 2-oxo acid dehydrogenase complexes and must be attached to the lipoyl domain to play its part in the systems of [active-site](#) coupling and substrate channeling that are prominent features of these complexes (13). There are strong parallels (13) between these properties of lipoic acid and those of [biotin](#) in various ATP-dependent carboxylases and 4-phosphopantothenic acid in fatty acid and polyketide synthases (see [Thiotemplate Mechanism Of Peptide Antibiotic Synthesis](#)). Lipoic acid has additionally been attributed to numerous other biological effects, among them that of a protective antioxidant against free radicals, an inhibitor of lipid peroxidation and of [HIV](#) replication, and a valuable dietary supplement (14, 15).

## 1. Structure and Biosynthesis

Lipoic acid was first discovered in the quest for the “pyruvate oxidation factor” (16). It is based on the eight-carbon [fatty acid](#), octanoic acid, modified by the insertion of sulfur atoms at C-6 and C-8 to give a dithiolane ring (Fig. 1). There is a chiral carbon atom at position 6, and generally only the R-enantiomer is active in the 2-oxo acid dehydrogenase complexes (17). There are intriguing similarities between the enzymes, which contain [2Fe-2 S] clusters (see [Iron-Sulfur Proteins](#)), responsible for the insertion of sulfur into the biotin and lipoic acid precursors (18).

**Figure 1.** Structure of lipoic acid (a) and its biosynthetic precursor, octanoic acid (b).



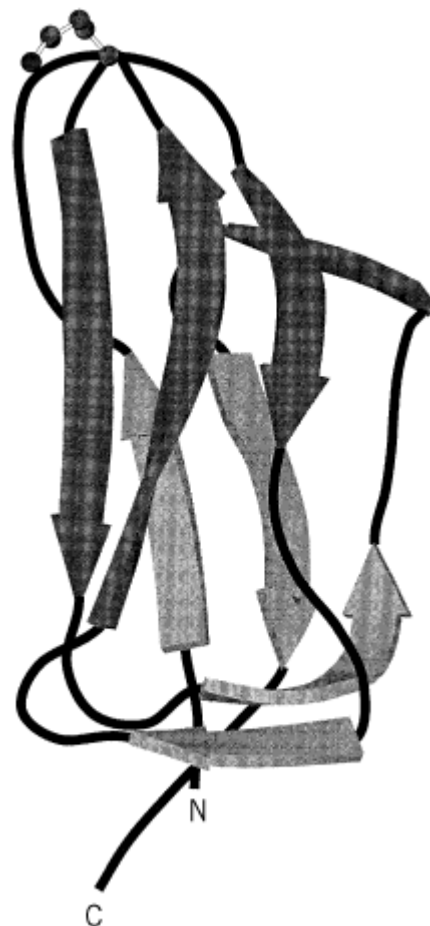
## 2. The Lipoyl Domain

The E2 polypeptide chain of the 2-oxo acid dehydrogenase complexes consists of, from the *N*-

terminus: one to three lipoyl domains (depending on the source), a peripheral subunit-binding domain, and a catalytic (acyltransferase) domain. The domains are joined together by long (20 to 30 residue) and conformationally flexible linker regions, and the catalytic domain aggregates with octahedral (24-mer) or **icosahedral** (60-mer) symmetry (again according to source), to form the structural core of the complex (1-3). Lipoic acid will not serve as a substrate for the 2-oxo acid decarboxylase (E1) component unless attached to the lipoyl domain (the value of  $k_{cat}/K_m$  is thereby raised by a factor of  $10^4$ ); moreover, the lipoyl domain restricts reductive acylation of the lipoyl group to the partner E1 of the parent complex. Thus the true substrate is the conformationally mobile lipoyl domain, an elegant molecular basis for substrate channeling and active-site coupling in these complexes (2, 13).

The structures of several lipoyl domains have been determined by means of **nuclear magnetic resonance** spectroscopy (6-8). All consist of two four-stranded **b-sheets**, with the lipoyl-lysine residue displayed in a prominent b-turn in one sheet, and the *N*- and *C*-termini close together in space in the other sheet, related by a two-fold axis of quasi-symmetry (Fig. 2). The structure of the lipoylated H-protein of the glycine cleavage system is similar (9). Interaction of the lipoyl domain with E1 is only transient, and the specificity appears to depend in large part on the nature of the amino acid residues surrounding the lipoyl-lysine residue in its b-turn and on the neighboring surface loop between strands 1 and 2 (19). Curiously, in the human autoimmune disease, primary biliary cirrhosis, the offending **antigen** is the lipoylated domain of the E2 component of the pyruvate dehydrogenase complex (20), which is normally a mitochondrial protein.

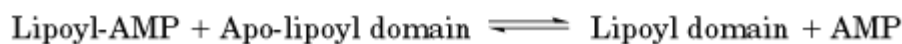
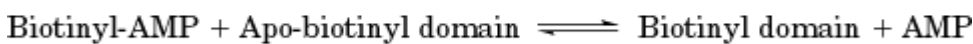
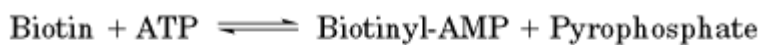
**Figure 2.** Structure of the lipoyl domain of the dihydrolipoyl succinyltransferase component of the 2-oxoglutarate dehydrogenase complex of *Escherichia coli* (based on Ref. 7). The b-sheet containing the lipoyl-lysine residue is shown in dark shading, and the b-sheet containing the *N*- and *C*-terminal residues is shown in light shading. The lipoylation site (Lys43) in the turn between strands 4 and 5 is indicated.



### 3. Post-translational Modification

The lipoyl group is attached to the target lysine residue in the lipoyl domain by an ATP-dependent enzyme, lipoyl protein ligase (10, 11). Its two-step mechanism resembles that of a fatty acyl CoA synthetase or [aminoacyl tRNA synthetase](#) (Fig. 3). The specificity of the [post-translational modification](#) depends crucially on the correct siting of the target lysine residue in the exposed b-turn of the apo-lipoyl domain, and much less on the surrounding amino acid sequence (12). In this it differs significantly from many other post-translational modifications, for which the sequence motif is dominant.

**Figure 3.** Reactions catalyzed by biotinyl and lipoyl protein ligases.



### 4. Similarities with Biotin

Biotin-lysine is the swinging arm carrying the carboxy group in multienzyme systems that catalyze carboxylation reactions. Like lipoic acid, it too is attached in amide linkage to the  $N^6$ -amino group of a specific lysine residue in the relevant enzyme. There are intriguing similarities between biotin and lipoic acid: (a) in the enzymes that catalyze their biosynthesis; (b) in the ability of both to bind tightly to avidin (although lipoic acid does so with a much higher dissociation constant,  $\sim 10^{-6}\text{M}$ ); (c) in the existence of a biotinyl domain in biotin-dependent enzymes, the structure of which closely resembles that of the lipoyl domain (13, 21); (d) in the requirement for the lipoic acid or biotin to be attached to the lipoyl or biotinyl domain before it will serve as a substrate in its parent enzyme complex; (e) in the biotinylation of the target lysine residue catalyzed by an enzyme that mechanistically resembles lipoyl protein ligase (Fig. 3); and (f) in the dependence of post-translational modification on the siting of the target lysine residue in the exposed b-turn of the apo-domain (13).

### 5. Role as a Swinging Arm

There is clear evidence for the essentially free rotation of the swinging arm on the surface of the lipoyl domain of 2-oxo acid dehydrogenase complexes (2, 3). However, the lipoyl-lysine residue in the H-protein of the glycine cleavage system is localized by interactions with the protein (9). It switches to a new position when charged with substrate, such that the aminomethylated derivative is sequestered in a surface cavity of the domain unique to the H-protein (9). In this instance, the swinging arm is fulfilling the expectation of the “hot potato hypothesis” of multienzyme complexes, protecting an unstable intermediate for presentation to the next enzyme in the sequence (13). Likewise, the biotinyl-lysine residue of the biotinyl domain of the acetyl CoA carboxylase of *E. coli* is clearly localized by interaction with the protein (21), but there is no evidence of similar biotin–

protein interactions in the 1.3 S subunit of *Propionibacterium shermanii* transcarboxylase (22). It is not known what purpose, if any, is served by the prior localization of the biotinyl-lysine residue in biotin-dependent reactions.

### Bibliography

1. L. J. Reed and M. L. Hackert (1990) *J. Biol. Chem.* **265**, 8971–8974.
2. R. N. Perham (1991) *Biochemistry* **30**, 8501–8512.
3. A. Berg, A. and A. de Kok (1997) *Biol. Chem.* **378**, 617–634.
4. K. Okamura-Ikeda, K., Y. Ohmura, K. Fujiwara, and Y. Motokawa (1993) *Eur. J. Biochem.* **216**, 539–548.
5. D. Macherel, J. Bourguignon, E. Forest, M. Faure, C. Cohen-Addad, and R. Douce (1996) *Eur. J. Biochem.* **236**, 27–33.
6. F. Dardel, A. L. Davis, E. D. Laue, and R. N. Perham (1993) *J. Mol. Biol.* **229**, 1037–1048.
7. P. M. Ricaud, M. J. Howard, E. L. Roberts, R. W. Broadhurst, and R. N. Perham (1996) *J. Mol. Biol.* **264**, 179–190.
8. A. Berg, J. Vervoort, and A. de Kok (1997) *Eur. J. Biochem.* **244**, 352–360.
9. C. Cohen-Addad, S. Pares, L. Sieker, M. Neuberger, and R. Douce (1995) *Nature Struct. Biol.* **2**, 63–68.
10. T. W. Morris, K. E. Reed, and J. E. Cronan (1995) *J. Bacteriol.* **177**, 1–10.
11. K. Fujiwara, K. Okamura-Ikeda, and Y. Motokawa (1997) *J. Biol. Chem.* **272**, 31974–31978.
12. N. G. Wallis and R. N. Perham (1994) *J. Mol. Biol.* **236**, 209–216.
13. R. N. Perham and P. Reche (1998) *Biochem. Soc. Trans.* **26**, 299–303.
14. L. Packer, E. H. Witt, and H. J. Tritschler (1995) *Free Radic. Biol. Med.* **19**, 277–250.
15. L. Packer, H. J. Tritschler, and K. Wessel (1997) *Free Radic. Biol. Med.* **22**, 359–378.
16. L. J. Reed (1957) *Adv. Enzymol.* **18**, 319–347.
17. Y.-S. Yang and P. A. Frey (1989) *Arch. Biochem. Biophys.* **268**, 465–474.
18. I. G. Serebriiskii, V. M. Vassin, and Y. D. Tsygankov (1996) *Gene* **175**, 15–22.
19. N. G. Wallis, M. D. Allen, R. W. Broadhurst, I. A. D. Lessard, and R. N. Perham (1996) *J. Mol. Biol.* **263**, 436–474.
20. M. J. Howard, C. Fuller, R. W. Broadhurst, R. N. Perham, J.-G. Tang, J. Quinn, A. G. Diamond, and S. J. Yeaman (1998) *Gastroenterology* **115**, 139–146.
21. F. K. Athappilly and W. A. Hendrickson (1995) *Structure* **3**, 1407–1419.
22. D. V. Reddy, B. C. Shenoy, P. R. Carey, and F. D. Sönnischen (1997) *Biochemistry* **36**, 14676–14682.

### Suggestions for Further Reading

23. D. B. McCormick, J. W. Suttie, and C. Wagner, eds. (1997) *Methods in Enzymology*, Vol. **279**. *Vitamins and Coenzymes*, Part I. Academic Press, San Diego.
24. M. S. Patel, T. E. Roche, and R. A. Harris, eds. (1996) *Alpha-Keto Acid Dehydrogenase Complexes*, Birkhäuser Verlag, Basel.

## Liposomes

Liposomes are synthetic vesicles comprised of bimolecular layers (bilayers) of phospholipids. They have been used as models of cell [membranes](#) (1), as drug delivery systems in which the drug is encapsulated within the liposome, and as carriers for genetic material into cells (see [Transfection](#)) (2). They are generally classified according to their size and the number of bilayers in the vesicle. Unilamellar vesicles consist of a single bilayer enclosing an aqueous core and may be small (SUVs) with a diameter of up to 300 Å, or large (LUVs) with diameters comparable to cell membranes. Multilamellar vesicles may be 1 to 50 µm in diameter; in cross section, when viewed by electron microscopy, the structure appears onion-like in which the bilayers are concentrically arranged and separated by alternating aqueous layers. SUVs generally require special procedures for their preparation, such as sonication (3); multilamellar vesicles generally form spontaneously when phospholipids are dispersed in water. LUVs have structures most closely resembling the membrane; they may be prepared, but only with special techniques such as extrusion of multilamellar vesicles under pressure through small pore membranes (4). LUVs will form spontaneously in phospholipid dispersions and in surface films, but only at a critical temperature that depends on the phospholipid composition (5). If the total lipid composition of a cell membrane is used to form the aqueous dispersion, the critical temperature for spontaneous formation of LUVs is the physiological temperature of the cell from which the lipids are removed (6) (see [Membranes](#)).

Liposomes have also been utilized to examine the nature of the fundamental forces that exist between the phospholipid bilayers of multilamellar vesicles. An approach that has been especially informative is the *osmotic stress* method, in which the aqueous spacing between the bilayers is modified by osmotic pressure. This method utilizes bilayer-impermeable, water-soluble polymers in the suspending solution to create an osmotic gradient between the internal aqueous phase of liposomes and the external solution. Equilibration results in water being driven from the liposome until the interbilayer forces balance the osmotic pressure. From measurements of the equilibrium interbilayer spacing, using X-ray diffraction methods and its dependence on osmotic pressure, Rand and Parsegian (7) have obtained the contributions of [hydration](#), [electrostatic interactions](#), and [van der Waals interactions](#) to the internal energy of the multilamellar vesicles.

#### Bibliography

1. A. D. Bangham, M. W. Hill, and N. G. A. Miller (1974) In *Methods in Membrane Biology*, Vol. 1 (E. D. Korn, ed.), Plenum Press, New York, pp. 1–68.
2. D. D. Lasic (1997) *Liposomes in Gene Delivery*, CRC Press, Boca Raton, FL.
3. C. Huang (1969) *Biochemistry* **8**, 344–352. (A systematic study of the hydrodynamic properties of sonicated vesicles; it describes the first use of gel filtration chromatography to separate SUVs from multilamellar vesicles.)
4. S. Kölschens et al. (1993) *Chem. Phys. Lipids* **65**, 1–10.
5. K. Tajima and N. L. Gershfeld (1985) *Biophys. J.* **47**, 203–209; N. L. Gershfeld (1989) *J. Phys. Chem.* **93**, 5256–5261.
6. N. L. Gershfeld (1986) *Biophys. J.* **50**, 457–461; N. L. Gershfeld L. Ginsberg (1997) *J. Membr. Biol.* **156**, 279–286.
7. R. P. Rand and V. A. Parsegian (1989) *Biochim. Biophys. Acta, Rev. Biomembr.* **988**, 315–334.

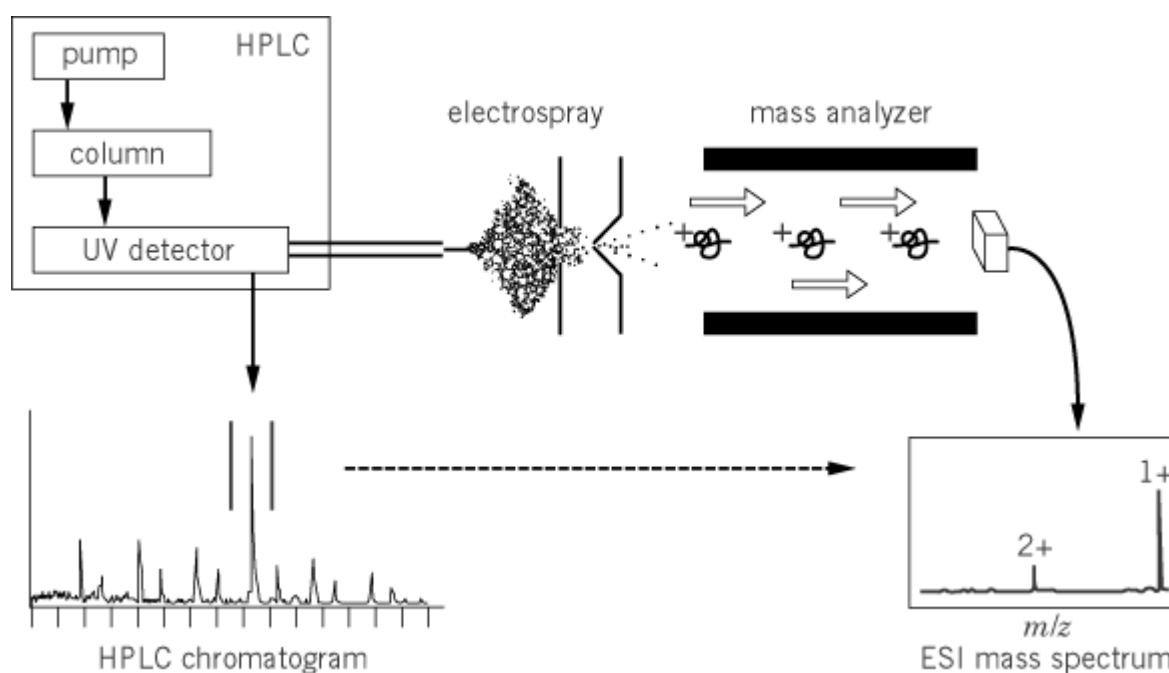
#### Suggestion for Further Reading

8. R. E. Pagano and J. N. Weinstein (1978) *Annu. Rev. Biophys. Bioeng.* **7**, 436–468. (A review of the interactions of liposomes with mammalian cells.)

## Liquid Chromatography Mass Spectrometry

Separations of macromolecules by [chromatography](#) are readily combined with [mass spectrometry](#) (MS) to monitor the molecular weights of the eluted molecules (see [Chromatofocusing](#), [Hydrophobic Chromatography](#), [Hydroxyapatite Chromatography](#), [Ion-Exchange Chromatography](#), [Isocratic Elution Chromatography](#), [Reversed-Phase Chromatography](#)). This is accomplished readily with [electrospray ionization](#) (ESI) MS, as the ions are formed directly from solution (usually an aqueous or aqueous/organic solvent system). In the past, attempts to couple liquid chromatography with mass spectrometry resulted in limited success, but ESI has made liquid chromatography mass spectrometry (LCMS) routine. Figure 1 illustrates how mass spectrometry is interfaced to liquid chromatography. Electrospray ionization also allows for MS analysis at relatively high liquid chromatography flow rates (1.0 mL/min) and high mass accuracy ( $\pm 0.01\%$ ), adding a new dimension to the capabilities of chromatographic characterization. In fact, using ESI-MS as a detector for chromatography was one of its first obvious applications.

**Figure 1.** In LCMS, the HPLC is typically interfaced with an electrospray ionization (ESI) mass spectrometer. The mass spectral data obtained on the compounds as they elute can provide valuable information in product identification.



Another use of LCMS is in the analysis of [peptide mapping](#), where the protein can be identified by the mass analysis of its constituent peptides and by comparison of this data to a protein database (Fig. 1). This approach has become increasingly useful when combined with [tandem mass spectrometry](#), where the peptides undergo fragmentation that produces more useful sequence information and thus makes it easier to identify the protein from a database (see [Proteome](#)) (1). This approach is also useful for identifying [post-translational modifications](#) of proteins (2).

#### Bibliography

1. J. R. Yates, S. Speicher, P. R. Griffin, and T. Hunkapiller (1993) *Anal. Biochem.* **214**, 397–408.
2. H. Taniguchi, M. Suzuki, S. Manenti, and K. Titani (1994) *J. Biol. Chem.* **269**, 22481–22484.



## Long Terminal Repeats

Within [retroviruses](#), the integrated, double-stranded proviral DNA genome is flanked at its 5'- and 3'-termini by identical noncoding regions designated long terminal repeats (LTR). Each consists of three “domains.” U3 and U5 are derived from unique sequences at the 3' and 5' termini of the viral RNA genome, respectively, and R denotes repeat sequences of the termini whose homology is exploited to transfer nascent DNA within or between genomes during proviral DNA synthesis. The domains are linked in the order -U3-R-U5-. Sites for initiation of **minus-** and **plus-strand** DNA (the primer binding site (PBS) and the polypurine tract (PPT)) are located immediately [downstream](#) and [upstream](#) of the 5'- and 3'-LTR, respectively. Although they provides many common control mechanisms, retroviral LTR vary considerably in size, ranging from as little as 350 bp in [Rous sarcoma virus](#) (RSV) to greater than 1700 bp in SFV-3 (Table 1). Size heterogeneity is also evident within individual LTR domains, for example, the R region of mouse mammary tumor virus (MMTV) has only 13 bp, whereas its human T-cell leukemia virus (HTLV) BLV counterpart has 235 bp. This contrasts with the U3 region of MMTV which, at almost 1200 bp, is among the largest reported. Currently, the reason for such size heterogeneity amongst retroviral LTR is unclear.

**Table 1. Size Variation among U3, R, and U5 Components of Retroviral LTR<sup>a</sup>**

| Retroviral Group | U3      | R       | U5      | Replication Primer (tRNA) |
|------------------|---------|---------|---------|---------------------------|
| B-type           | 1200    | 13–15   | 120     | Lys <sup>3</sup>          |
| Avian C-type     | 150–250 | 18–21   | 80–100  | Trp                       |
| Mammalian C-type | 450–500 | 60–70   | 75      | Pro                       |
| D-type           | 235–240 | 13–15   | 95      | Lys <sup>1,2</sup>        |
| HTLV/BLV         | 250–350 | 120–240 | 100–200 | Pro                       |
| Lentivirus       | 350–450 | 100–200 | 80–150  | Lys <sup>3</sup>          |
| Spumavirus       | 800     | 200     | 150     | Lys <sup>1,2</sup>        |

<sup>a</sup> Most notable among these are the extremely short R elements, which play a pivotal role in DNA strand transfer during replication.

Despite their identity, 5- and 3'-retroviral LTR assume different roles during the expression of proviral DNA. U3 of the 5'-LTR harbors the **promoter**, a binding site for the host-coded **RNA polymerase II** and a variety of cellular protein factors that cumulatively govern the level of [transcription](#). Viral RNA initiates at the first base of the 5'-R, and extends to the end of its 3'-counterpart, where it is **polyadenylated**. Additional levels of control have been identified within R and U5 sequences of the transcript, including (i), the [HIV](#) TAR loop, which interacts with virally coded **transactivator** proteins (eg, Tat) to augment transcription and (ii) U5-IR stem-loop structures, which interact through a variety of mechanisms with the **tRNA** replication primer to control initiation of minus-strand synthesis. Thus, the LTR can be considered a major contributor toward

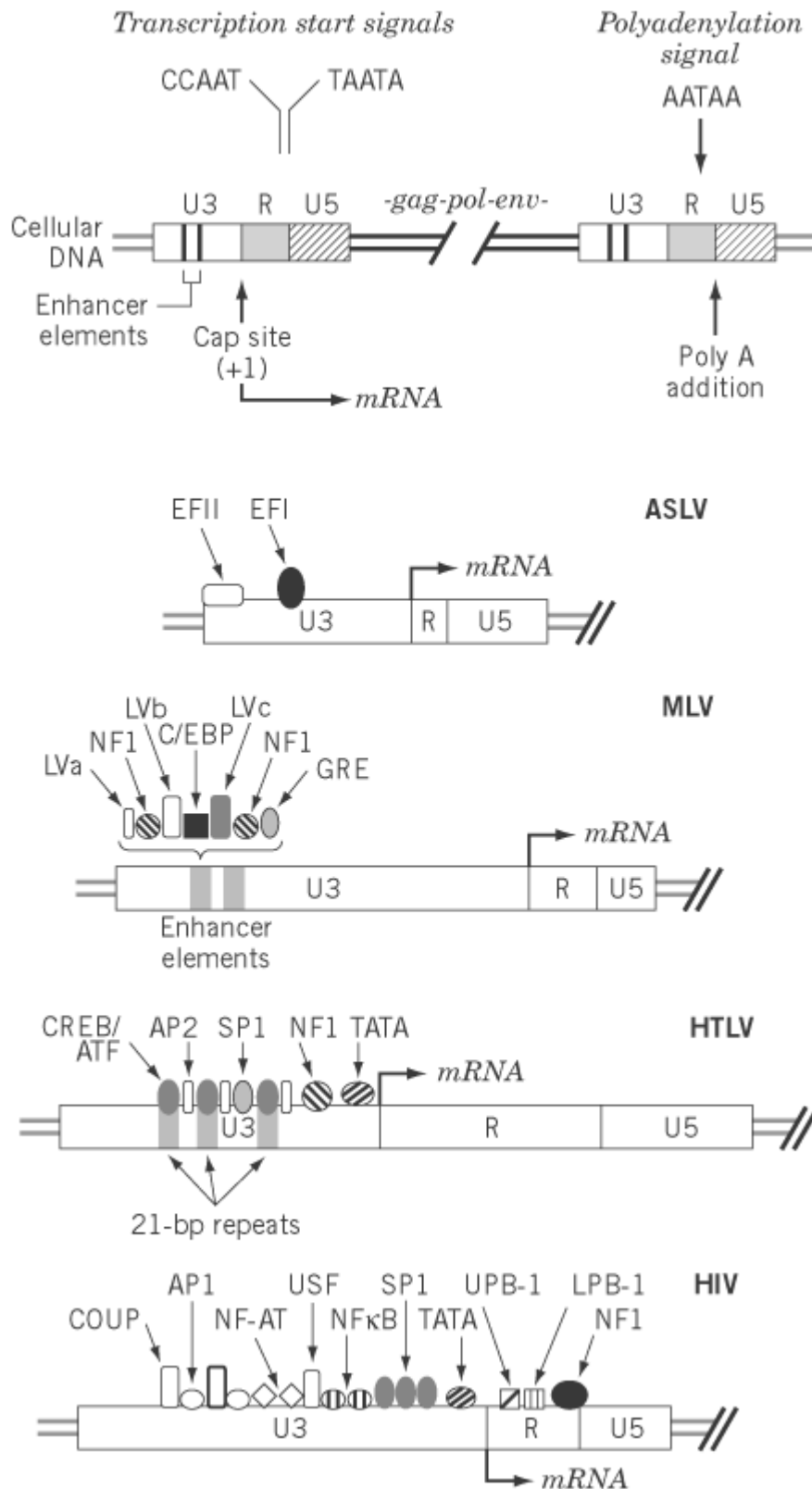
transcription and **reverse transcription**.

## 1. LTR control of viral transcription

### 1.1. The LTR Promoter

The LTR promoter provides the binding site for **RNA polymerase II** and accessory host cellular factors that interact with U3 sequences upstream from the [TATA box](#), a consensus sequence 20 to 30 bp upstream of the transcription initiation site that mediates polymerase recognition. The complexity of [cis-acting](#) sequence elements in LTR promoters, shown in Fig. 1, varies from relatively few in the case of ALSV to as many as 10 in HIV-1. Several spatially distinct binding sites for a single accessory protein may exist (eg, AP-2 sites on the HTLV promoter) in addition to multiple adjacent binding sites (eg, the NF-AT and NFkB sites on the HIV promoter). Although these recognition elements are depicted primarily upstream of the transcriptional initiation site, precedents exist for regulatory elements more distal from the promoter. An example of this is the presence of a binding site for the **C/EBP transcriptional factor** in the *gag* gene of ALSV. The TATA box and upstream elements are commonly called the basal promoter.

**Figure 1.** Transcriptional control elements of retroviral LTR. The upper portion represents the integrated form of the provirus and indicates a variety of *cis*-acting elements in the U3 and R regions of the 5'- and 3'-LTR that regulates viral transcription. Binding sites for cellular transcriptional factors in the LTR of ALSV, MLV, HTLV, and HIV are indicated.



## 1.2. Enhancers

The [enhancers](#), an important class of *cis*-acting elements are sequences that specify binding of cellular factors to enhance the basal level of transcription. Core enhancer sequences are relatively small (10 to 15 bp), but they collectively constitute a hierarchy of binding domains to generate one or more copies of an element functional independent of both position and orientation. Examples of this is are the tandem, 75-bp murine leukemia virus (MLV) enhancers depicted in Fig. 1, which

sequester seven **transcriptional factors**.

Stimulation of transcription by viral enhancer elements has a variety of biological consequences. Although avian leukosis virus (ALV)-induced B-lymphomas lack an **oncogene** they reflect enhancer activity derived from proviral integration adjacent to the promoter for *c-myc* **proto-oncogene**. Proviral integration on either side of the proto-oncogene promoter can stimulate its transcription. In contrast, Rous-associated virus type 0 (RAV-0) infections are nonpathogenic, which correlates with lack of LTR enhancer activity in this virus. Deleting an enhancer element within the MLV LTR reduces the capacity to induce thymic lymphomas in mice.

Enhancer elements also play an important role in tissue-specific expression. Wild-type Moloney murine leukemia virus (MoMLV) is not effectively expressed in the liver or brain of mice. However, when the MoMLV LTR enhancer is replaced with its counterpart from the cellular transthyretin gene and a promoter-proximal sequence that controls tissue-specific expression, the resulting recombinant virus is infectious and is expressed in previously nonpermissive tissues. These observations are clinically significant because they open the possibility of constructing tissue-specific recombinant retroviral vectors for gene therapy.

### 1.3. Regulation of LTR Transcription

MMTV, a murine retrovirus responsible primarily for mammary tumors, provides a well-studied example of LTR expression regulated by [glucocorticoid](#) hormones. The glucocorticoid **receptor** contains an N-terminal activation domain, a central DNA binding motif, and a C-terminal ligand-binding site. After hormone treatment, receptor-bound hormone is transported to the [nucleus](#), resulting in hypersensitivity of chromatin near glucocorticoid [response elements](#) (GRE). Accessibility of these exposed regions to transcription factors results in enhanced RNA polymerase II activity from promoters containing GRE motifs. The MMTV LTR contains multiple copies of the GRE. In the absence of hormone induction, a neighboring binding site in the LTR for the nuclear factor NF-1 is inaccessible because of association with [nucleosomes](#). However, after induction, transport of the hormone-receptor complex to the nucleus essentially “clears” [chromatin](#) from the retroviral LTR, making the NF-1 binding site accessible to the transcriptional activator and stimulating viral transcription.

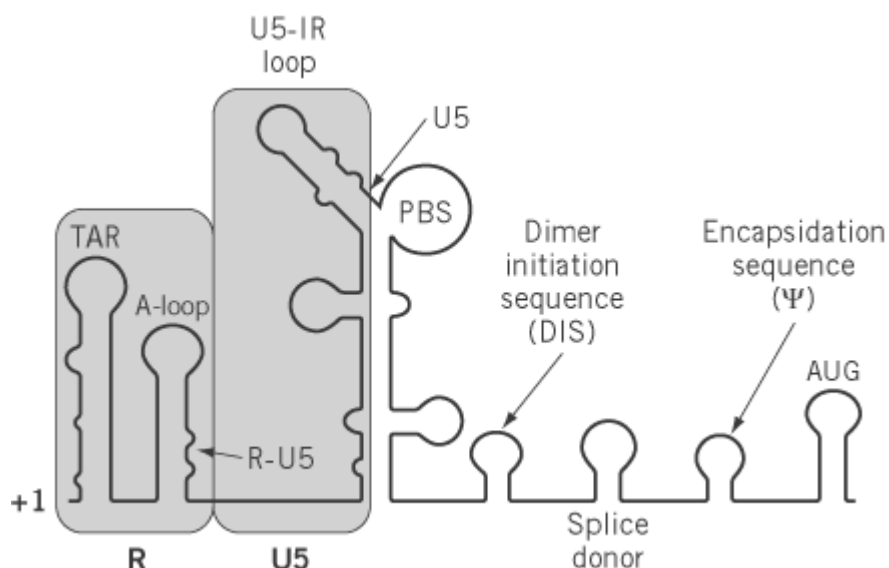
### 1.4. Transactivation

Transactivation is a feature of many complex retroviruses (spumaviruses, lentiviruses, and HTLV-related viruses) whereby, in addition to host factors, virally coded **transactivators** stimulate LTR transcription ([1](#)). Stimulation may occur through an interaction with DNA sequences in U3 (eg, the HTLV Tax protein and the *bel-1* and *taf* gene products of human and simian foamy viruses, respectively), or through binding to a particular sequence near the 5'-end of the RNA transcript (eg, HIV Tat proteins). Tax-mediated stimulation of HTLV transcription is mediated via three 21-bp repeat elements in U3 that serve as recognition sites for cyclic AMP-responsive element binding protein (CREB), a member of the cAMP response family. A direct interaction of Tax with the 21-bp repeats has not been documented and spawns the theory that Tax exerts its function through protein-protein interactions to modulate CREB. Although the transactivators of foamy viruses act through similar mechanisms, their target sequences share little in common with those of Tax.

The target of the Tat transactivator in HIV and related lentiviruses is an ordered structure near the 5' end of the viral RNA transcript, the transactivator response element or TAR (Fig. [2](#)). Through an exhaustive series of investigations, the Tat target on TAR has been located at a small U-rich bulge in the stem. Different mechanisms have been proposed for the stimulatory effect of HIV Tat. One model suggests overriding transcriptional termination, based on observations that, in the absence of Tat, transcripts that initiate in the LTR pause after synthesizing 60 to 80 nucleotides. It is envisioned that TAR-bound Tat (and possibly other cellular factors) thereby interacting with and stabilize the host transcriptional machinery. Model *in vitro* systems show that an alternative hypothesis envisages Tat acting at the stage of transcriptional initiation through an interaction with the transcription factor SP-1. This may be a fortuitous consequence of the proximity of TAR to the site of **mRNA** initiation,

hence the possibility that Tat action reflects a combination of these events also cannot be ruled out. Finally, the observation that Tat is released from infected cells and enters neighboring cells to modulate transcriptional activity suggests a further role of this protein in HIV pathogenesis.

**Figure 2.** Schematic representation of secondary structural elements in the 5'-leader RNA of retroviruses that, it is proposed, regulate transcription and reverse transcription. +1 denotes the first nucleotide of the viral transcript, that is, the 5'-end of R. RNA that corresponds to the R and U5 regions of the LTR is shaded. Adjacent elements that control dimerization (DIS) and encapsidation of the retroviral RNA genome ( $\Psi$ ) have been indicated. AUG represents the initiation codon of the *gag* open reading frame. PBS is the primer binding site.



## 2. R-U5 RNA and control of reverse transcription

The 5'-non-coding end of the viral transcript, that is, between the first nucleotide of R and the *gag* initiation codon, is defined as the leader RNA. Although the regulatory mechanisms differ among retroviruses, it is generally recognized that structural elements within the R-U5 portion of the leader play an important role in regulating the efficiency of reverse transcription. Extensive chemical and enzymatic probing has revealed a complex set of intramolecular stem-loop structures, as well as intermolecular duplexes involving different regions of the tRNA replication primer, the binding site for which lies immediately adjacent to U5 (Fig. 2). In discussing LTR control of reverse transcription, it is thus necessary to include RNA sequences that comprise the 5'-leader.

In addition to complementarity between viral PBS sequences and the 3'-terminal nucleotides of the tRNA replication primer, secondary intermolecular interactions critical to reverse transcription have been uncovered that involve U5 of RNA human and avian retroviruses. In ALSV, it has been elegantly demonstrated genetically and biochemically that bases comprising the TYC loop of tRNA<sup>Trp</sup> (the cognate replication primer of ALSV; see [Transfer RNA](#)) interact with sequences in the U5 leader stem to control initiation of minus-strand DNA synthesis (2). Although secondary structural predictions for RNA folding suggest the potential for a similar mechanism in many retroviruses, HIV-1 adopts an alternative approach. In this case (3), the U-rich tRNA<sup>Lys,3</sup> anticodon domain is implicated in controlling reverse transcription, the target of which is the A-rich U5-IR loop in the immediate vicinity of the PBS (Fig. 2). Although both models are equally plausible, they must recognize a common requirement for disruption of RNA structural elements before and accompanying initiation of reverse transcription. For example, although hybridization of tRNA<sup>Trp</sup> to

the PBS of the ALSV genome has the consequence of unwinding its TYC stem, it is necessary to disrupt the U5 leader stem to provide the type of intermolecular duplex been proposed. Furthermore, once such non-PBS intermolecular duplexes are established, they must also be disrupted by the retroviral replicating machinery. The energy to disrupt such structures may be an intrinsic feature of RT, as proposed for the avian enzyme. Alternatively, this may be supplied by way of deoxynucleoside triphosphate hydrolysis during polymerization. Finally, accessory viral proteins, such as NC, could potentially interact with RT and serve a stimulatory role during initiation of reverse transcription.

Why are such complex regulatory mechanisms of reverse transcription necessary? One clue to this may lie in observations that alternative sites on the retroviral genome with considerable sequence homology to the PBS can be identified. Although initiation of reverse transcription from such “pseudo” initiation sites might be possible, this would have severe consequences at later steps in replication when the tRNA primer is removed from nascent minus-strand DNA. Minus-strand DNA sequences immediately adjacent to the tRNA primer are destined to become terminal nucleotides of the 3'-LTR critical for recognition by the retroviral integration machinery (see Fig. 4 of [Retroviruses](#)). Thus, although the PBS provides the appropriate sequence for localizing the tRNA 3'-terminus, additional tRNA/viral base pairing between could be envisaged as “locking” or stabilizing the replication primer at the appropriate initiation site. Support for this notion comes from genetically manipulated genomes of HIV-1, where the natural PBS is replaced by a variant that specifies binding of another tRNA isoacceptor species. If the PBS alone is exchanged, the resulting virus replicates poorly. In contrast, replication kinetics are significantly improved when sequences in the U5-IR loop are simultaneously altered to complement the anticodon loop of the heterologous tRNA isoacceptor.

In addition to the PBS, two non-LTR elements of the leader RNA should be pointed out. The first of these is the dimer initiation sequence (DIS), which promotes intermolecular base pairing of the RNA genome to ensure that a dimer of identical molecules is introduced into the budding virion. In HIV-1, **palindromic** sequences in the loop of a hairpin structure mediate dimerization through a mechanism defined as the “kissing complex.” Finally, the encapsidation sequence (Y) provides a signal for packaging the retroviral genome into the budding virion that is mediated through an interaction with the NC component of the *gag* polyprotein. Location of Y between the major splice donor and the *gag* initiation codon ensures that spliced RNAs are not incorporated (see [RNA Splicing](#)). The exception to this is ALSV, whose splice donor is located in the *gag* gene, that is, downstream of the encapsidation sequence, and predicts that subgenomic RNA is transported to the virion. The observation that such species are in fact discriminated against suggests the existence of a cryptic encapsidation signal close to the 5' end of the *gag* transcript.

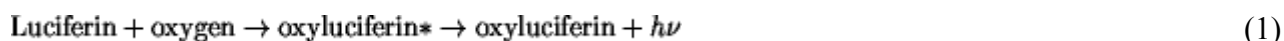
## Bibliography

1. J. Karn, M. G. Gait, M. J. Churcher, D. A. Mann, I. Mikaelin, and C. Pritchard (1994) In *RNA-Protein Interactions* (K. Nagai and I. W. Mattaj eds.), Oxford University Press, New York, pp. 192–220.
2. J. Leis, A. Aiyar, and D. Cobrinik (1993) In *Reverse Transcriptase* (A. M. Skalka and S. P. Goff, eds.) Cold Spring Harbor Monograph Series, Cold Spring Harbor, New York, pp. 33–47.
3. C. Isel, C. Ehresmann, G. Keith, B. Ehresmann, and R. Marquet (1995) *J. Mol. Biol.* **247**, 236–250.

## Luciferases And Luciferins

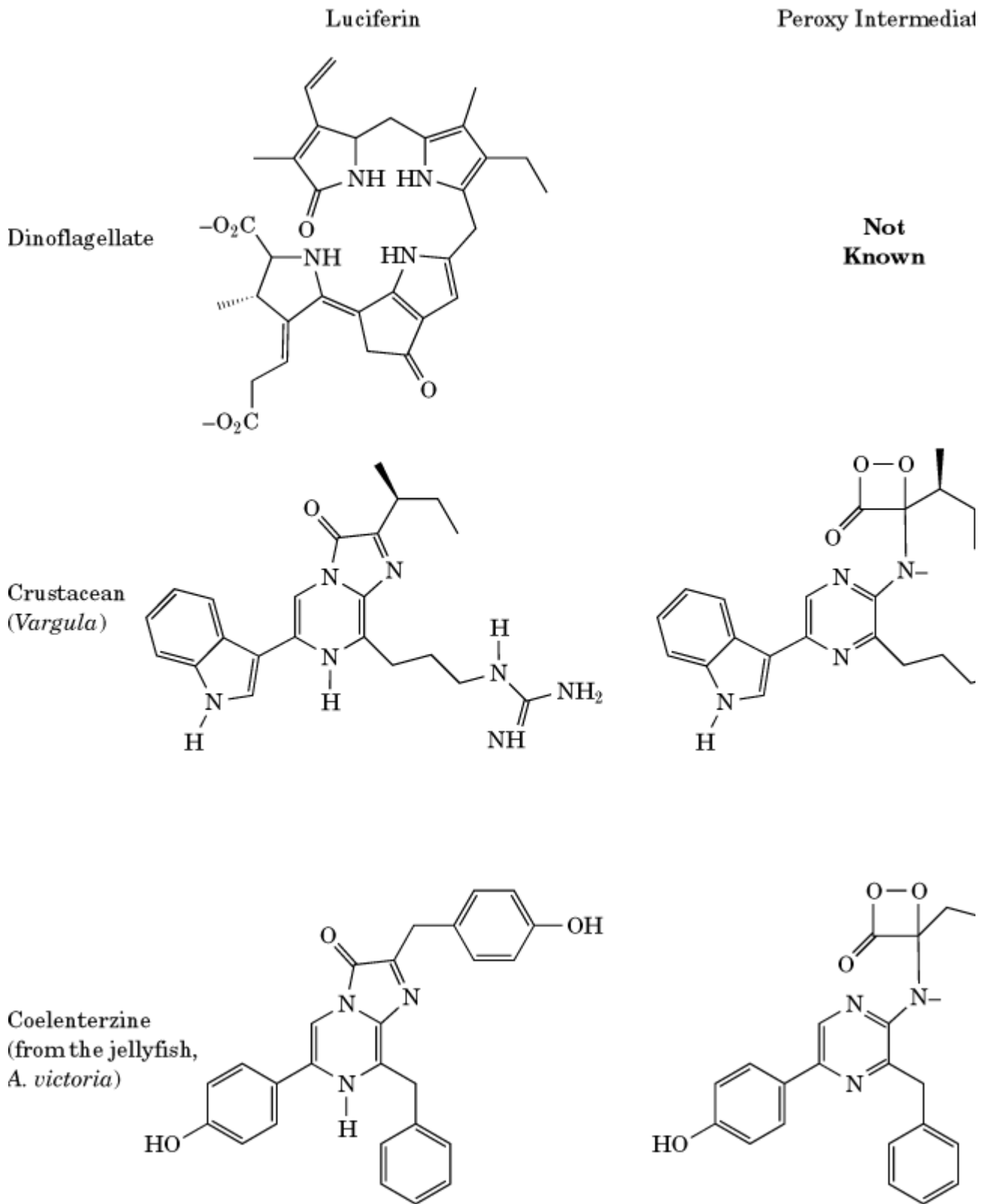
Luciferase is the generic term for an [enzyme](#) that catalyzes a reaction that produces light; luciferin is the general term for a substrate of these reactions. Luciferin and luciferase are both named from the Latin *lucifer*, referring to Venus, the morning star. (*Lucifer* literally means light-bringing.) Luciferases all catalyze oxidative reactions that give off light (and as such require oxygen), but other than this, different luciferases may have little or nothing in common. There are many different bioluminescent systems in nature, and consequently many different luciferins and luciferases.

In general, the reactions catalyzed by luciferases are special cases of **chemiluminescent** reactions. The general reaction scheme is as follows:



Thus, luciferases catalyze the oxidation of their respective luciferins, producing an excited-state molecule of oxyluciferin (the asterisk denotes an excited state). It is the decay of this molecule to the ground state that is the source of the light. Differences in luciferin molecules, or changes in the environment of the active site, may result in differences in the color of the emitted light. The variety of luciferin structures is demonstrated by the examples shown in [Figure 1](#).

**Figure 1.** The various forms of luciferin, with the corresponding peroxy intermediate and oxyluciferin.



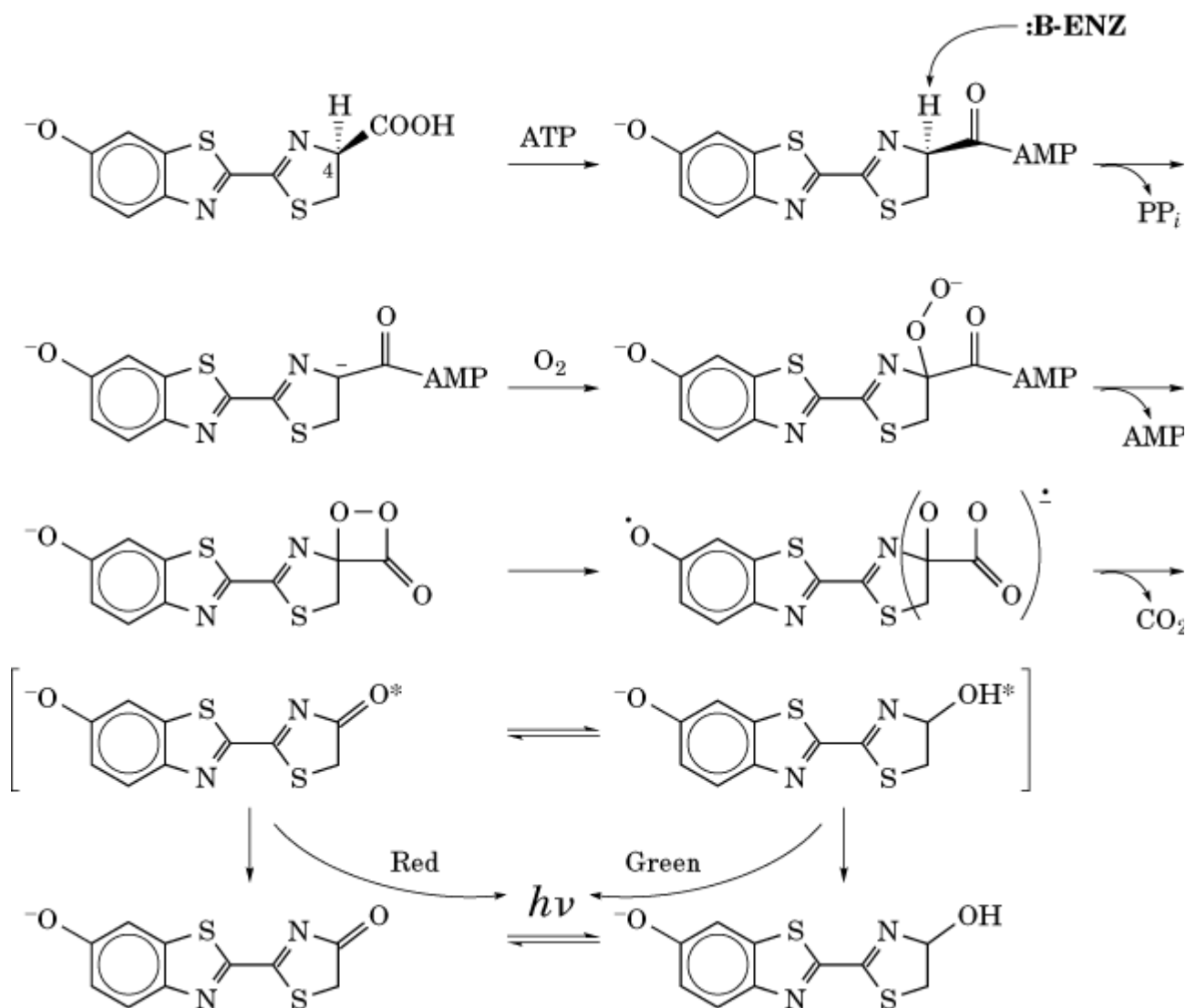
## 1. Firefly Luciferase

Perhaps the most widely known bioluminescent reaction is that of the common firefly. This reaction is exceptional in that it has a high quantum yield of about 90%, meaning that about 90% of the reacting molecules will emit a photon of light. The detailed mechanism of the firefly reaction is shown in Figure 2. In the first step of the reaction, the carboxylic acid moiety of the firefly luciferin



is activated by reaction with ATP to produce an adenylated intermediate. The C4 proton is abstracted, producing a carbanionic form of luciferin. This is the form of luciferin that reacts with molecular oxygen to produce the key intermediate, an activated dioxetanone. This dioxetanone is unstable because of the high strain energy in the four-membered ring and because of the inherently weak peroxide bond. The breakdown of this intermediate is generally believed to occur by a chemically induced electron exchange luminescence (CIEEL) mechanism (discussed below) and is accompanied by loss of carbon dioxide. The final product is the excited state of oxyluciferin, which can exist in either the keto or the enol forms of excited oxyluciferin, which are responsible for emitting red and yellow-green light, respectively.

**Figure 2.** Mechanism of bioluminescence in the firefly.



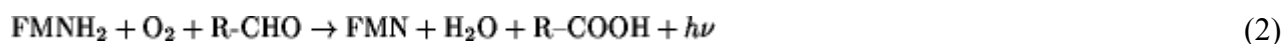
The generation of the excited state is thought to proceed by the CIEEL mechanism. This hypothesis states that the benzothiazole portion of the molecule acts as an electron-rich (easily oxidizable) donor molecule that donates an electron (going from negative to neutral in the process) to the peroxide moiety, which facilitates cleavage of the O–O peroxide bond, leaving one oxygen atom with a full negative charge, and one as a neutral radical. After loss of the carbon dioxide, the electron is returned to the donating moiety, yielding excited-state oxyluciferin.

The firefly is not the only species to make use of a dioxetanone intermediate. Other organisms such as the sea pansy *Renilla reformis*, the marine crustacean *Cypridina hilgendorffii*, and the jellyfish *Aequorea victoria* have luciferins of markedly different structures, but each is thought to proceed through a dioxetanone intermediate. The color of light emitted by these species also differs from that of the firefly reaction. It is important to note, however, that the color of light emitted is not determined solely by the structure of the luciferin. The enzyme and its microenvironment also have an important role to play. It is possible to generate mutant forms of luciferase in the laboratory that emit light of a different color than the naturally occurring reaction.

Although the reaction mechanism for firefly luciferin is known in some detail, and the structure for firefly luciferase has recently been solved (without bound luciferin), the details of the interaction between the luciferin and luciferase are not known. However, based on [sequence analysis](#) of luciferases from many firefly species, the active site of firefly luciferase has been tentatively identified. The structure is a “hammer and anvil” structure, where there is a small C-terminal “hammer” **domain** atop a large N-terminal “anvil” domain. Many of the putative active-site residues are at the interface between these domains. It is speculated that since, like all luciferase reactions, water must be excluded in order for the bioluminescence reaction to occur with a high photon yield, the hammer and anvil regions of luciferase clamp together, with the substrates in between and excluding water in the process. Of course, the enzyme must also provide general bases for abstracting protons involved in the mechanism, but no specific amino acid residues have yet been implicated.

## 2. Bacterial Luciferase

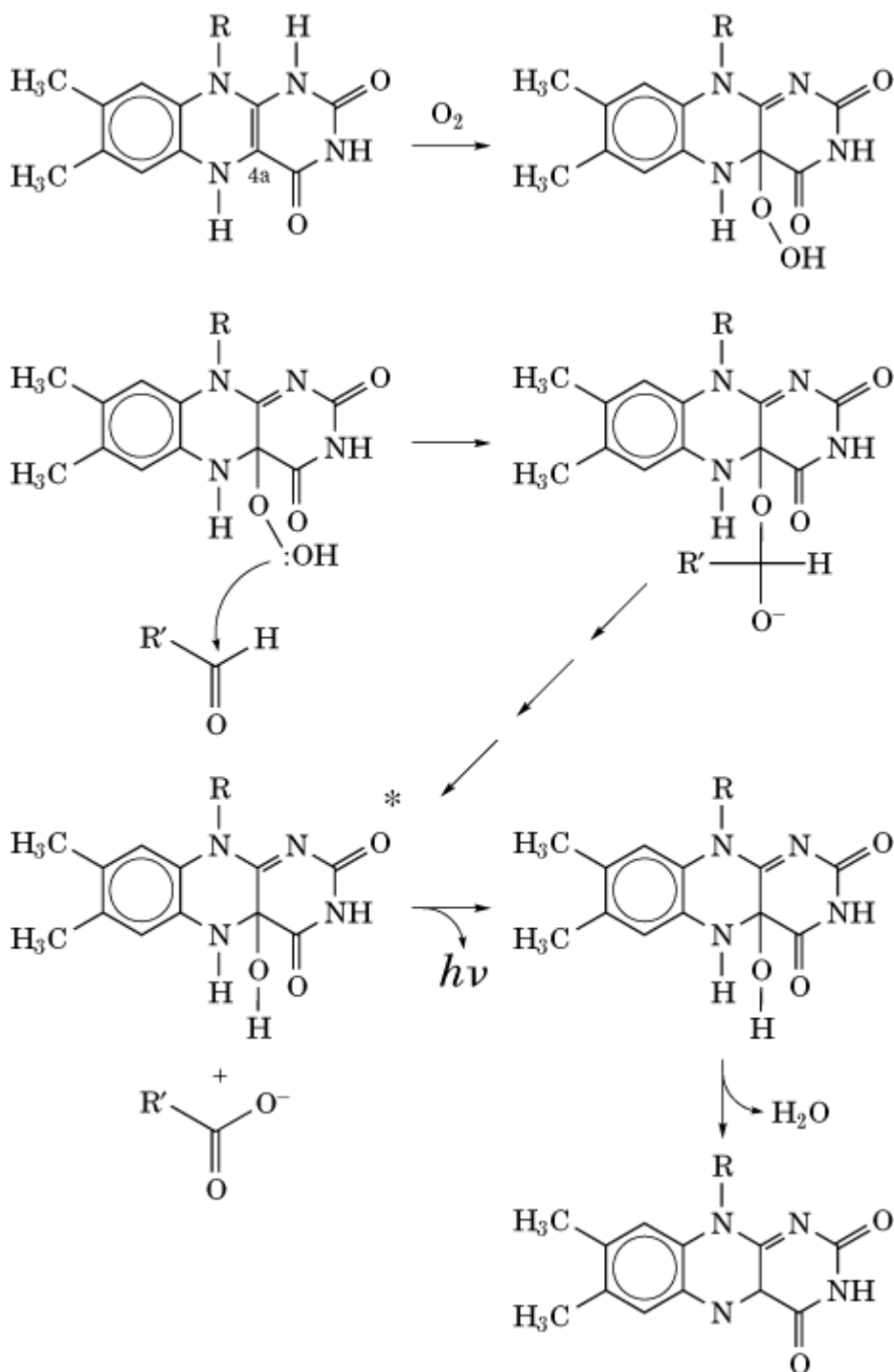
One of the best characterized bioluminescence systems is that from bacteria. Several different species of luminescent bacteria exist, but the luciferases from the various organisms are **homologous**. Unlike most bioluminescent systems, however, the bacterial luciferase system does not have a luciferin molecule that is unique to the bioluminescence reaction. Instead, these systems use reduced flavin mononucleotide (FMNH<sub>2</sub>) and a long-chain fatty aldehyde, which are converted to oxidized flavin and the corresponding carboxylic acid. Similar to other bioluminescent reactions, this oxidation reaction requires the use of molecular oxygen. Another important difference between this reaction and most bioluminescent reactions is that the peroxide species leading to the light emission is a linear rather than a cyclic peroxide. However, like firefly luciferase, the structural details of the interaction of bacterial luciferase with its substrates are not well understood. The overall chemical reaction is as follows:



where R-CHO is a long-chain aldehyde, R-COOH is the corresponding carboxylic acid, and FMN is oxidized flavin mononucleotide.

The reactants must bind to the luciferase in a specific sequence; FMNH<sub>2</sub> first, then O<sub>2</sub>, then the aldehyde. A detailed scheme of the bacterial bioluminescence reaction is shown in Figure 3. The formation of a 4a-peroxyflavin from FMNH<sub>2</sub> and O<sub>2</sub> is the first step in the reaction. This peroxyflavin may then react with the aldehyde to form a tetrahedral intermediate (the formation of an alkyl-flavin peroxide bond should be noted). Breakdown of the tetrahedral intermediate involves cleavage of the peroxide bond along with the formation of an excited-state 4a-hydroxyflavin and a carboxylic acid. The hydroxyflavin emits light as it relaxes to the ground state and then undergoes dehydration to yield FMN and H<sub>2</sub>O.

**Figure 3.** Mechanism of bioluminescence in bacteria.



Bacterial luciferase is a heterodimer, consisting of two homologous a and b subunits with a total molecular weight of about 76 kDa. The two subunits have an overall sequence identity of about 40%, and a similarity of about 80%. The crystal structure of the enzyme in the absence of substrates has recently been solved. The two subunits each form a [TIM barrel](#) structure, with a four-helix bundle as the core of the interface between the subunits. Both subunits are required for normal activity, but there appears to be only one active site per dimer. There is a narrow channel leading into a large **hydrophobic** cavity in the core of the a subunit, and this is believed to be the active site of the enzyme. Cys 106 is located at the mouth of this cavity; modifications to this residue have been shown to compromise the catalytic activity of the enzyme greatly. Further, two tryptophan residues, 194 and 250, which are postulated to be near the flavin binding site on the basis of spectroscopic

measurements, are found lining this cavity. Mutation of either of these tryptophans reduces the FMNH<sub>2</sub> binding ability of the enzyme.

In the present crystal structure, residues 262–290 are unresolved, which strongly suggests that this segment of the polypeptide chain is flexible. Given the structure of the rest of the enzyme, it is likely that the unresolved amino acids form a “lid” that covers the entrance to the active site. It is known from **proteolysis** and **chemical modification** experiments that a conformational change occurs during the catalytic cycle of the enzyme. It is speculated that this conformational change may be the flexible lid closing over the active site after substrates have been bound, excluding water from the reaction. Another interesting feature of the internal cavity centers around Asp 113. In the crystal structure, this residue forms a **hydrogen bond** network with His 44 and His 45, which mutational analysis have also shown to be important for activity and flavin binding. Although Asp 113 is not itself part of the putative active site cavity, His 44 and 45 are. Substitution at position 113 with neutral or positively-charged amino acids results in activity loss of three to six orders of magnitude. It is clear that disruption of this hydrogen bonding network interferes with the activity of the enzyme. Although the reason for this is not yet understood, it is believed to be due to distortion of the active site.

### 3. Aequorin and Calcium-Binding Photoproteins

Calcium binding proteins constitute another category of luciferases. The best known of these is aequorin from the jellyfish *Aequorea victoria*. Aequorin was originally thought to be an exception to the general rule that all bioluminescent reactions require molecular oxygen or hydrogen peroxide, because isolates of the enzyme would emit light on addition of calcium, even under anaerobic conditions. However, it was eventually determined that the cofactor for aequorin, called coelenterazine, is covalently bound to the enzyme in a reaction requiring molecular oxygen. Thus, once the coelenterazine is bound to the enzyme, it is activated and needs only calcium for the bioluminescence reaction to occur. Once the calcium is added, aequorin undergoes a conformational change and becomes capable of catalyzing the oxidation of the bound coelenterazine. The products of the reaction are light, carbon dioxide, and apoaequorin with noncovalently bound coelenteramide. Reconstitution of the activated enzyme requires molecular oxygen and fresh coelenterazine in reducing conditions. The mechanism of the bioluminescence reaction in aequorin is under some debate. It seems clear, however, that the mechanism has similarities to the firefly reaction, including a dioxetanone-like intermediate and CIEEL-type charge transfer.

### 4. Green Fluorescent Protein

*Aequorea* (and some other bioluminescent coelenterates) have a second protein in their bioluminescent systems. This noncatalytic protein is the green fluorescent protein (GFP). GFP absorbs maximally at 395 nm and has a smaller absorption peak at 470 nm, while the bioluminescence emission spectrum of aequorin has a maximum at 470 nm. Thus GFP is able to act as an antenna protein and to accept the energy from the bioluminescent reaction of aequorin. GFP then emits light with a maximum intensity at 509 nm. The presence of GFP serves to increase the quantum yield of the aequorin. The crystal structure of GFP has recently been solved; it is an 11-stranded  $\beta$ -barrel surrounding a central  $\alpha$ -helix, capped by three smaller distorted helices. The chromophore for GFP is a modified hexapeptide derived from the sequence Phe–Ser–Tyr–Gly–Val–Gln, which is found in the interior of the  $\beta$ -barrel.

### 5. Applications of Luciferase Technology

Although the study of luciferases is a fascinating academic pursuit in its own right, luciferase research has yielded many technological benefits as well. Bioluminescence is experimentally a very useful characteristic for an enzyme to have, because light production can be detected at extremely low levels and over a wide range of intensities. This, coupled with the fact that bioluminescence assays are quite rapid, has made luciferases the assay of choice in many different systems. Firefly

luciferase is sold commercially as a sensitive and accurate assay for ATP, and it is able to detect ATP down to femtomolar ( $10^{-15}$  M) concentrations. Similarly, aequorin may be used as an equally sensitive assay for calcium. Finally, bacterial and firefly luciferases are commonly used as markers for genetic engineering, in the study of gene regulation, and for many other laboratory and industrial uses. Luciferases are singularly useful in these contexts, as they allow for sensitive reporting of molecular events as they occur within living systems.

### Suggestions for Further Reading

T. O. Baldwin (1996) Firefly luciferase: the structure is known but the mystery remains, *Structure* **4**, 223–228.

T. O. Baldwin and M. M. Ziegler (1992) "The biochemistry and molecular biology of bacterial bioluminescence", in *Chemistry and Biochemistry of Flavoenzymes*, Vol. **III**, F. Müller, ed., CRC Press, Boca Raton, FL, pp. 467–530.

T. O. Baldwin, J. A. Christopher, F. M. Raushel, J. F. Sinclair, M. M. Ziegler, A. J. Fisher, and I. Rayment (1995) Structure of bacterial luciferase, *Curr. Opin. Struct. Biol.* **5**, 798–809.

E. Conti, N. P. Franks, and P. Brick (1996) Crystal structure of firefly luciferase throws light on a superfamily of adenylate-forming enzymes, *Structure* **4**, 287–298.

J. W. Hastings and D. Tu, eds. (1995) "Molecular mechanisms in bioluminescence", *Photochem. Photobiol.* **62**, 599–673.

M. Ormö, A. B. Cubitt, K. Kallio, L. A. Gross, R. Y. Tsien, and S. J. Remington (1996) Crystal structure of the *Aequorea victoria* green fluorescent protein, *Science* **273**, 1392–1395.

## Luminescence

Luminescence is the emission of light that does not result from high temperatures (cf. incandescence). In general, luminescence occurs when an atom or molecule is excited into a high energy state, and then decays to the ground state. Since electronic energy levels are quantized, the decay to the ground state is accompanied by the emission of a photon of a specific wavelength. Luminescence is categorized by the mode of excitation that produces the high energy excited state. [Chemiluminescence](#) is luminescence resulting from a chemical reaction, and bioluminescence is luminescence resulting from a biological ([enzyme](#)-catalyzed) chemical reaction, as in the firefly (see [Luciferases And Luciferins](#)). **Fluorescence** occurs when an atom or molecule is excited (either by an electric discharge or by absorption of a photon) into the singlet excited state, which then decays to the ground state. The lifetime of this excited state is very short (on the order of picoseconds), resulting in rapid emission. Some molecules may undergo intersystem crossing where the singlet excited state becomes a triplet state. Since the rules of quantum mechanics forbid a transition from a triplet excited state directly to the ground state, the triplet excited state has a long lifetime (seconds to hours), resulting in a weak but long-lived glow called phosphorescence.

## Lymphokines

The term *lymphokine* refers to soluble proteins produced by lymphoid cells that affect their

proliferation, maturation, or function. It derives from studies in the 1970s and 1980s in which lymphoid cells, upon exposure to infectious agents or mitogens, became activated and secreted proteins that acted in an autocrine and/or paracrine fashion. Since then, many of the genes for lymphokines have been **cloned**, sequenced, and characterized. Analysis of separated cell populations revealed that some of the lymphokines were produced exclusively by lymphocytes and acted primarily on lymphocytes (eg, **interleukin-2**), while other lymphokines were produced by **macrophages** and a variety of other cells and acted upon **T cells** and numerous other cells (eg, interleukin-1). Another important lymphokine studied was **interferon-g**; this protein is produced by T cells and acts on macrophages and many other cells. As the number of cloned protein factors increased to more than 40, it became apparent that many different cell types had the potential to produce a particular lymphokine; and the lymphokine, in turn, had effects on diverse cell types. This development has encouraged the use of the term “cytokine” rather than “lymphokine” for soluble protein factors, although these terms are sometimes used interchangeably. Some of these factors have been given interleukin (IL) designations and at the present time number from IL1 to IL18. In addition, some factors have retained their original designation—for example, granulocyte-macrophage [colony stimulating factor](#), **interferon-g tumor necrosis factor**, and [transforming growth factor](#)  $\beta$ . In general, lymphokines have effects on hemopoietic cells and cells of the immune system. They have short *in vivo* half-lives and act locally at very low concentrations. Some are produced constitutively, but most are up-regulated in response to infection, [antigen](#) activation, and trauma. Interleukin-2, for example, is up-regulated in T cells at least 1000-fold in response to antigenic or mitogenic stimuli. Most lymphokines are produced with a [signal peptide](#) that is cleaved at the cell membrane, thereby releasing lymphokine to the exterior. They bind to specific receptors on cell surfaces and activate (or down-regulate) clusters of genes through multiple signaling pathways.

Emphasis will be placed here on those lymphokines having effects primarily upon lymphoid cells. The lymphokines will be divided into three groups, based on their effects on T and B lymphocytes and NK (natural killer) cells. For more details about individual lymphokines refer to the following entries: [Interleukins](#), [Interferon](#), [Colony stimulating factors](#), [Transforming growth factor](#), [Tumor necrosis factor](#), and [Chemokines](#).

## 1. T-Cell Lymphokines

Many lymphokines are important in the development, differentiation, and activation of T cells. T-cell precursors arise in the bone marrow and migrate to the thymus, where they undergo maturation. The immature T cells enter the thymus as  $CD4^-CD8^-$  immature T cells. These cells develop into  $CD4^+CD8^+$  cells, and then into either  $CD4^+CD8^-$  or  $CD4^-CD8^+$  cells, which then migrate into the periphery as mature T cells. Interleukin-7 (IL-7), one of the most important cytokines in T-cell differentiation, is produced by bone marrow and thymic stromal cells and causes proliferation of pre-T cells in the bone marrow, as well as  $CD4^-CD8^-$  cells in the thymus (1). In the thymus, IL-2 is produced by  $CD4^-CD8^-$  cells and induces an autocrine cycle of stimulation and differentiation (1). IL-2 is also produced by mature T cells and stimulates NK-cell activation and B-cell growth. IL-6 is secreted by thymic stromal cells and acts as a costimulator of  $CD4^+CD8^-$  and  $CD4^-CD8^+$  cells, along with IL-2 and IL-7 (1). IL-1a and IL-1b are secreted by thymic epithelial cells and stimulate proliferation of  $CD4^+CD8^-$  and  $CD4^-CD8^+$  cells (1). IL-4 is secreted by  $CD4^-CD8^-$  cells and, much like IL-2, acts as an autocrine stimulator on these cells (1). IL-4 also induces development of Th2 cells, which are mature  $CD4^+CD8^-$  “helper” cells, so named because they play a major role in providing help to [B cells](#) (2). Th2 cells are also involved in allergic responses and regulate Th1 cell responses (see text below). The IL-4 secreted by these Th2 cells stimulates B-cell [class switching](#) and augments [IgG](#) 1 production. Th2 cells also secrete other cytokines, such as IL-5, IL-10, and IL-13 (2). IL-5 induces the production of eosinophils, primes basophils for histidine release, and provides additional help for B cells. IL-10 and IL-13 act to down-regulate the production of inflammatory cytokines by a variety of cells (3, 4). Dendritic cells produce IL-12, which drives

development of Th1 cells. Th1 cells are mature CD4<sup>+</sup>CD8<sup>-</sup> cells that are associated with inflammatory responses and regulate Th2 cell responses (2). Th1 cells also secrete interferon-gamma (IFN $\gamma$ ), which stimulates induction of **multiple histocompatibility complex** (MHC) class I and II antigen expression, enhances B and NK cell activity, and has powerful antiviral and tumoricidal effects (5). Mature CD4<sup>-</sup>CD8<sup>+</sup> cells, which are called **cytotoxic T lymphocytes** (CTLs), secrete IL-2, IFN $\gamma$ , and tumor necrosis factor b (TNFb). TNFb induces proliferation and differentiation of a wide variety of cell types, is an integral component of the inflammatory response, and also has direct cytotoxic effects on cells (6).

## 2. B-Cell Lymphokines

B cells originate from pluripotent hematopoietic **stem cells** and differentiate in the bone marrow into mature B cells (see **Hematopoiesis**). The differentiation can be divided into five steps: early pro-B, late-pro B, pre-B, immature B, and mature B cells. Upon release from bone marrow, in response to antigen activation, mature naive B cells differentiate into either plasma cells or memory B cells. Most differentiation steps appear to be regulated by lymphokines, which include interleukin 4 (IL-4), interleukin-5 (IL-5), interleukin-6 (IL-6), and interleukin-7 (IL-7).

Interleukin-7, a 25-kDa glycoprotein produced by stromal cells in the bone marrow, plays an important role in early B cell development. At late pro-B cell stage, IL-7 induces proliferation and differentiation of pro-B cells to pre-B cells by up-regulating the expression of other cytokine receptors that are necessary for the survival and growth of pre-B cells. *In vitro* studies also suggest the ability of IL-7 to support the growth of T cells. Interleukin-4, a 14-kDa glycoprotein, is a critical lymphokine for B-cell activation, proliferation, and **isotype** switching. Most IL-4 is produced by CD4<sup>+</sup> type-2 T-helper cells (TH2), and its production is used as the criterion for placement of CD4<sup>+</sup> T cells into this subset. IL-4 can exert its action on resting B cells in the absence of any costimulant. It induces an increase in class II MHC molecule expression on resting B cells, and it also causes a significant increase in the cell volume of resting B cells. Recent studies suggest that IL-4 is also an isotype “switch” factor and a regulator of allergic reactions. Interleukin-5 is a glycoprotein initially identified by its ability to support the growth and differentiation of B cells (7). IL-5 is a potent cytokine for inducing the production of natural antibody and antigen-induced **IgA** production. It can also synergize with many cytokines to enhance antibody production. In humans, the biologic effects of IL-5 are also well characterized for eosinophils. In addition to inducing terminal maturation of eosinophils, IL-5 prolongs eosinophil survival by delaying **apoptotic** death, acts as a chemotactic factor, increases eosinophil adhesion to endothelial cells, and enhances eosinophil effector function. Interleukin-6 is made by activated T cells. It is a lymphokine that induces terminal differentiation of B cells and enhances immunoglobulin production by activated B cells. The influence of IL-6 is not restricted to B cells; it induces the expression of IL-2R on IL-6R bearing T cells. Therefore IL-6 affects B cell triggering by two mechanisms: through a direct effect on B cells and through activation of T cells.

In addition to antibody production, activated B cells can produce a large number of lymphokines involved in the regulation of immune and inflammatory responses and hematopoiesis (8): (i) immunosuppressive lymphokines, such as TGF-b1 and IL-10; (ii) proinflammatory molecules: IL-1, IL-6, IL-8, and TNF; and (iii) hematopoietic growth factors: granulocyte colony-stimulating factor (G-CSF), Granulocyte-macrophage colony-stimulating factor (GM-CSF), macrophage colony-stimulating factor (M-CSF), and IL-7. IL-10, secreted at low level by activated B cells, is a key lymphokine that inhibits cell-mediated immunity and inflammation while promoting humoral responses. Overproduction of B-cell-derived IL-10 plays a role *in vivo* in B-cell lymphomagenesis.

## 3. Natural Killer Cell Lymphokines

Natural killer (NK) cells are cytotoxic lymphocytes. They differ from T and B cells in that NK cells are not triggered by specific antigens and do not require target sensitization.

IL-2, IFN $\alpha$ , IFN $\gamma$ , IL-7, IL-12, and IL-6 are able to directly activate NK cells *in vitro*. Stimulation of NK cells with IL-2 increases cytotoxic granule content, enhances surface adhesion molecule expression, induces proliferation, and can stimulate NK cells to become lymphokine-activated killer (LAK) cells. IL-2 also stimulates ADCC activity in NK cells, mRNA [transcription](#), and secretion of cytokines. IFN $\alpha$  and IFN $\gamma$  induce LAK activity in enriched/purified NK cell cultures. The use of IFN $\alpha$  and IFN $\gamma$ , in conjunction with IL-2, enhances the level of NK cytotoxicity beyond that induced by IL-2 alone. However, IFNs do not induce NK cell proliferation and have been demonstrated to inhibit IL-2-induced proliferation in peripheral blood mononuclear cells (PBMCs). IL-7 can induce LAK activity in purified NK cell populations and can increase IL-2Ra chain surface expression. IL-12 can generate high levels of LAK activity/NK cytotoxicity in purified and resting NK cells. This effect is independent of accessory cells and IL-2. IL-12 can also enhance ADCC, surface adhesion molecule expression, and granule content of NK cells. IL-6 can induce surface adhesion molecule expression, induce production of TNF $\alpha$  and TNF $\beta$ , as well as induce low levels of proliferation in purified NK cells.

NK cells can be stimulated to produce IFN $\gamma$ , TNF, IL-3, GM-CSF, M-CSF, TGF $\beta$ , and IL-8. Purified NK cells are able to produce IFN $\gamma$ , although this may be dependent on nonadherent accessory cells. In response to phorbol ester/Ca<sup>2+</sup>-ionophore, IL-2, IL-6, IL-7, IL-12, or Fc-receptor interaction, NK cells express and secrete TNF. NK cells produce IL-3, GM-CSF, and M-CSF in response to Ca<sup>2+</sup>-ionophore/phorbol ester, anti-CD16, or IL-2. IL-2 activates NK cells to produce TGF $\beta$  and IL-8. IL-7 induces NK cells to secrete low levels of TNF $\alpha$  and IL-8 and high levels of GM-CSF. IL-12 induces the expression of low levels of TNF $\alpha$  and IL-8 and high levels of IFN $\gamma$ .

#### Bibliography

1. W. He and D. Kabelitz (1994) *Int. Arch. Allergy Immunol.* **105**, 203–210.
2. A. O'Garra and K. Murphy (1996) In *Th1 and Th2 Cells in Health and Disease* (S. Romagnani, ed.), Karger, Basel, pp. 1–13.
3. N. Nicola (1994) In *Guidebook to Cytokines and Their Receptors* (N. Nicola, ed.) Oxford University Press, Oxford, U.K., pp. 83–85.
4. A. McKenzie and G. Zurawski (1994) In *Guidebook to Cytokines and Their Receptors* (N. Nicola, ed.), Oxford University Press, Oxford, U.K., pp. 92–94.
5. P. Gray (1994) In *Guidebook to Cytokines and Their Receptors* (N. Nicola, ed.), Oxford University Press, Oxford, U.K., pp. 118–119.
6. I. Millet and N. Ruddle (1994) In *Guidebook to Cytokines and Their Receptors* (N. Nicola, ed.), Oxford University Press, Oxford, U.K., pp. 105–108.
7. S. Karlen M. De Boer et al. (1998) *Int. Rev. Immunol.* **16**( 3–4), 227–247.
8. V. Pistoia (1997) *Immunol. Today* **18**(7), 343–350.

#### Suggestions for Further Reading

9. Anonymous R & D Systems Catalogue Mini-reviews, 1995 and 1997.
10. K. Takatsu (1997) *Proc. Soc. Exp. Biol. Med.* **215**(2), 121–133.
11. B. Naume and T Espevik (1994) *Scand. J. Immunol.* **40**, 128–134.
12. C. Janeway and P. Travers (1996) "Immunobiology". In *The Immune System in Health and Disease*, Vol. **8**, 2nd ed., Garland, New York, pp. 27–29.

#### Lyon Hypothesis



In 1961 Mary Lyon proposed that one of the two female [X-chromosomes](#) in eutherian (placental) animals is inactivated to equalize the expression of genes from the two X-chromosomes in female cells, relative to the single X-chromosome in male cells. It has since been shown that the assembly of [facultative heterochromatin](#) on one of the two female X-chromosomes leads to transcriptional **silencing** (2, 3), confirming the Lyon hypothesis (see [Barr Body](#), [Random X-Inactivation](#), [X-Chromosome Inactivation](#)).

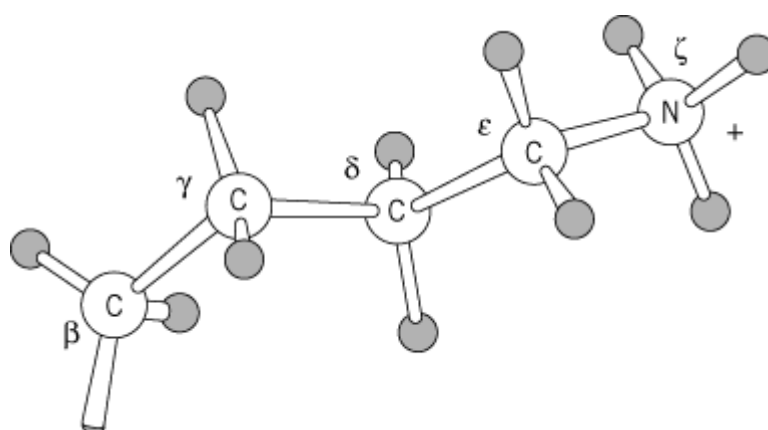
#### Bibliography

1. M. F. Lyon (1961) *Nature* **190**, 372–374.
2. M. F. Lyon (1988) *Am. J. Hum. Genet.* **42**, 8–16.
3. M. F. Lyon (1992) *Ann. Rev. Genet.* **26**, 15–27.

### Lysine (Lys, K)

The [amino acid](#) lysine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons** AAA and AAG—and represents approximately 5.7% of the residues of the proteins that have been characterized. The lysyl residue incorporated has a mass of 128.17 Da, a **van der Waals volume** of  $135 \text{ \AA}^3$ , and an [accessible surface](#) area of  $211 \text{ \AA}^2$ . Lys residues are not changed very frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [arginine](#), [asparagine](#), [threonine](#), [serine](#), and [glutamine](#) residues.

The side chain of Lys is a flexible **hydrophobic** chain of four methylene groups capped by an [amino group](#):



The amino group ionizes with an intrinsic  $pK_a$  value of about 11.1, so it is ionized under most physiological conditions. The ionized form is unreactive chemically, but there is always a finite fraction of nonionized amino groups, which are potent nucleophiles. Consequently, the amino groups of Lys residues readily undergo a typical wide variety of acylation, alkylation, arylation, and amidation reactions (see [Amino Groups](#)). These reactions can be used to measure the number of

Lys residues in a protein (see [Counting Residues](#) and [Trinitrobenzene Sulfonic Acid](#)).

The ionized amino groups of Lys residues in **protein structures** are nearly always exposed to the solvent, with the entire side chain typically exposed to the solvent and flexible, when they have **relaxation times** in the nanosecond range. Of **secondary structures**, Lys residues occur most frequently in ***α*-helices**; they also favor the helical conformation in model peptides. Virtually no Lys residues are buried within the interiors of proteins. They are occasionally used to attach prosthetic groups to proteins, such as the **Schiff Base** attachment of **pyridoxal phosphate** to some **enzymes**. In an important **post-translational modification**, Lys residues of **collagens** in the sequence –Xaa–Lys–Gly– are hydroxylated on the  $\alpha$  carbon by the enzyme lysyl 5-hydroxylase.

**Proteinases** frequently cleave polypeptide chains adjacent to Lys residues, as in the processing of pro-hormones, such as pro-**insulin**, at pairs of basic residues.

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Lysogeny

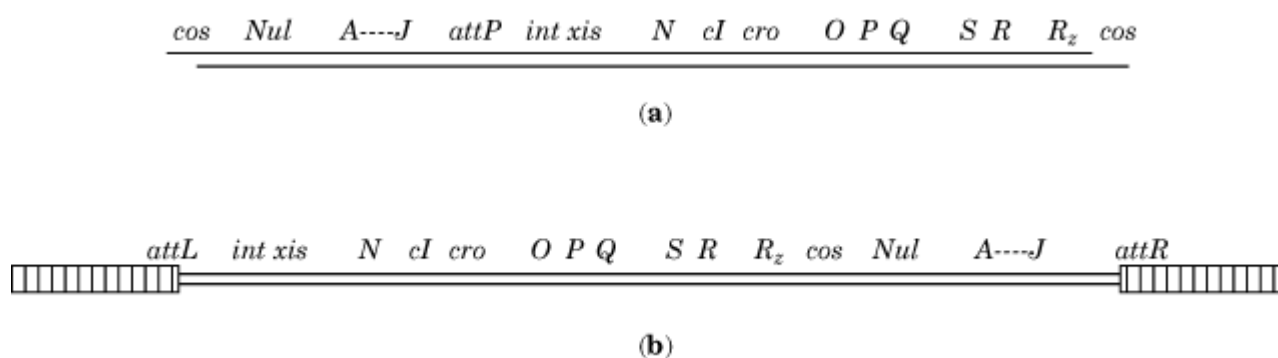
Upon infection of a bacterial cell by a bacterial virus, or bacteriophage, the **DNA** or RNA **genome** of the bacteriophage enters the cell. The encoded **RNA** and **proteins** are synthesized and replication takes place, leading to the production of about 100 progeny phages from each infected cell. Their release is usually accomplished by cell lysis, the breaking open of the cell, which is now dead. This course of events is referred to as a “lytic” infection. For the so-called “temperate” bacteriophages, whose genome invariably consists of double-helical DNA, there is another potential outcome of the infection: The expression and replication of the bacteriophage genome are inhibited, and the DNA gets integrated into the host genome, each DNA strand becoming covalently attached to, and part of, a strand of the bacterial genome. This is the “lysogenic” pathway; the phenomenon is referred to as “lysogeny,” the ensuing bacterium is a “lysogen,” and the integrated bacteriophage DNA is the “prophage.” Such a course of events provides advantages to the bacterium as well as to the bacteriophage. The bacterium becomes immune to subsequent infection by a bacteriophage like the one, whose genome it now harbors. The prophage is provided a “free ride”: Its DNA gets replicated with every round of replication of the bacterial DNA. After 8–10 generations, more copies of the bacteriophage genome will have been generated from the original one injected into the cell than if the lytic pathway had been followed. Some bacteriophages have the ability to respond to conditions that threaten the life of their host by killing the host and releasing progeny phages, in a process called “induction.” The phenomenon of lysogeny in bacteria has attracted great interest as a model system for viral latency in a eukaryotic cell, when infection results in integration of a virus genome into the cell's chromosomal DNA, as can occur with herpes or hepatitis virus.

### 1. Bacteriophage $\lambda$

**Lambda phage** ( $\lambda$ ) is the paradigm of a temperate phage. Not only is the process of establishment of lysogeny better understood for this bacteriophage than for any other, but it is also thought to be representative of the way most other bacteriophages accomplish this feat. Bacteriophage  $\lambda$  DNA relies entirely on the bacterial host's RNA and **protein biosynthesis** machinery for its expression. It contains all the appropriate signals to enable *Escherichia coli* RNA polymerase to use its DNA for

the synthesis of [messenger RNA](#), and *E. coli* [ribosomes](#) to use these mRNA to program the synthesis of bacteriophage proteins. Soon after the  $\lambda$  genomic DNA (Fig. 1A) has been injected into the bacterial cell, an irreversible decision is made, whether the infection will proceed along the lytic or the lysogenic pathways, which are mutually exclusive. Only after the lysogenic state has been well established can viable progeny be obtained by induction (see text below for additional details). In large extent due to the work carried out by Ptashne and co-workers, the establishment of lysogeny in bacteriophage  $\lambda$  is understood to a high level of molecular detail (1-4). Here, an outline of the most important steps in the process will be presented.

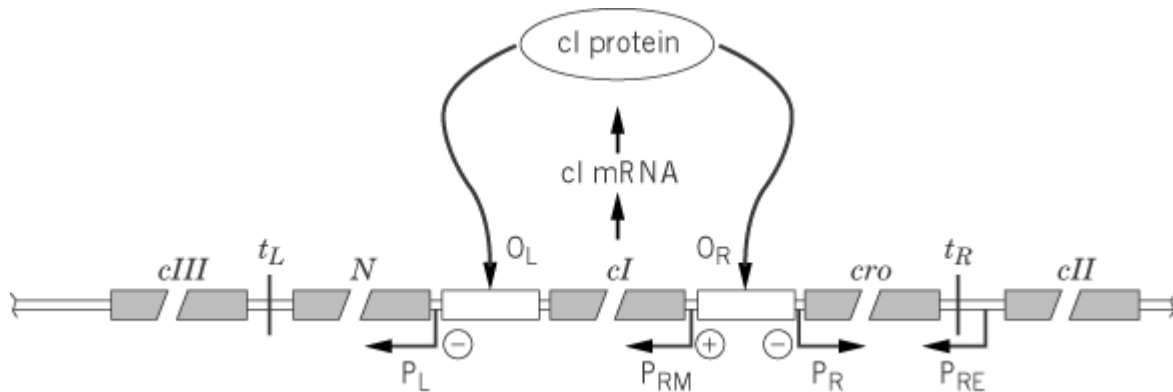
**Figure 1.** The arrangement of genes on the map of bacteriophage  $\lambda$  differs for the DNA in the phage and in the prophage: The two are circularly permuted. The map is not to scale. The double line symbolizes the two strands of the DNA double helix. (a) The linear DNA found in the phage heads (see [Lambda Phage](#)). At the ends are the complementary single-stranded regions (*cos* sites) that enable cyclization. From left to right are: *Nul* (a subunit of terminase), *A-J* (encoding tail and head proteins), *attP* (the attachment site), *int* and *xis* (involved in integration and excision), *N*, *cI*, and *cro* (regulatory genes discussed here), *O* and *P* (DNA replication), *Q* (antitermination), *S* and *R* (cell lysis) and *R<sub>z</sub>* (possibly involved in lysis). (b) Prophage DNA. At the ends are the *attL* and *attR* sites, each composed of half of the *attB* and *attP* sites originally on the bacterial and phage genomes, respectively (see text). The bacterial DNA is indicated as the striped regions on both sides of the prophage DNA; it, too, is double helical. Note that the *cos* sites are now internal. The prophage order of genes was obtained by cyclizing the bacteriophage DNA shown in (a) after injection into the cell, then reopening it at the *att* site in the process of integration into the bacterial genome.



The regulatory region of bacteriophage  $\lambda$  DNA is shown in Figure 2; this region represents approximately 10% of the entire  $\lambda$  genome (corresponding to the region between *xis* and *O* in Fig. 1A). It contains the genes encoding five regulatory proteins [*cI*, *cII*, *cIII* (pronounced “c-one,” “c-two,” and “c-three,” respectively), *N*, and *cro*], four **promoters** where RNA polymerase initiates RNA synthesis ( $P_L$ ,  $P_{RM}$ ,  $P_R$ , and  $P_{RE}$ ), two control regions or “**operators**” ( $O_L$  and  $O_R$ , each divided into three subsites), and two terminators, where RNA synthesis stops and RNA is released ( $t_L$  and  $t_R$ ). Together these affect the regulation of the life cycle of bacteriophage  $\lambda$ .

**Figure 2.** The control region of bacteriophage  $\lambda$ . This map is only intended to show the relative order of important regions and is not to scale. The DNA regions labeled *cI*, *cII*, *cIII*, *N*, and *cro* (gray “boxes”) are the genes encoding proteins *cI*, *cII*, *cIII*, *N*, and *cro*;  $O_L$  and  $O_R$  (white boxes) are binding sites for regulatory proteins *cI* and *cro*;  $P_L$ ,  $P_R$ ,  $P_{RM}$ , and  $P_{RE}$  are promoters from which transcription of RNA is initiated in the direction indicated by the horizontal arrows, and  $t_L$  and  $t_R$  are terminators where RNA synthesis stops and the transcription complex falls apart. Upward pointing arrows symbolize the synthesis of *cI* mRNA and protein. The downward arrows point to the binding sites for *cI* in a lysogen, and they also show the effect of the bound protein on transcription at the nearest promoter (shown as + for activation, – for inhibition). To establish lysogeny, *cII* binds near the  $-35$  region of  $P_{RE}$  to activate transcription from this promoter. Other important interactions and their effects: *cro* protein binds at  $O_R$  to inhibit transcription from

$P_{RM}$ ; N binds to the transcription complex that initiated at  $P_L$  to prevent termination at  $t_L$  and to the one that initiated at  $P_R$  to prevent termination at  $t_R$ . Note that names of genes are in “italics” letters and those of proteins in “roman” letters.



As discussed (see [Lambda Phage](#)), immediately upon entry of the linear I DNA (Fig. 1A) into the cell, it cyclizes by virtue of the fact that its two ends have complementary single-stranded DNA regions, and *E. coli* RNA polymerase initiates RNA synthesis at  $P_L$  and  $P_R$ , promoters for left- and rightward [transcription](#), respectively (Fig. 2). The DNA sequences of promoters  $P_{RE}$  and  $P_{RM}$  poorly resemble that of a consensus, or optimal, promoter. Therefore they are not well-utilized by RNA polymerase, unless aided by activator proteins specific to each promoter. Most of the RNAs initiated at  $P_R$  and  $P_L$  are terminated at transcription terminators  $t_L$  and  $t_R$ , but they still contain the entire coding sequences for the *cro* and *N* proteins, respectively; these two mRNA program *E. coli* ribosomes to synthesize the *cro* and *N* proteins. The **cro protein**, just like *cI* and *cII* a dimeric, sequence-specific [DNA-binding protein](#), represses initiation of RNA synthesis from  $P_{RM}$  by binding at  $O_R$ . The *N* protein is an **antiterminator** that enables RNA polymerase to continue beyond  $t_L$  and  $t_R$ , allowing additionally the mRNA for *cIII* to be synthesized from  $P_L$  and that for *cII* from  $P_R$ . These are again translated to give the *cIII* and *cII* proteins. At this point, the decision whether to proceed along the lytic or lysogenic pathways is made, in which the *cII* protein plays a major role: Accumulation of *cII* protein leads to establishment of lysogeny; otherwise the infection will be lytic. The *cII* protein activates transcription from  $P_{RE}$  (promoter for repressor establishment) (5), leading to the synthesis of *cI* mRNA, and *cI* protein. Note that even though the RNA polymerase that initiates RNA synthesis at  $P_{RE}$  traverses the *cro* gene, it does so in the wrong direction to yield functional *cro* mRNA. Other promoters activated by *cII* protein are  $P_{aQ}$  [involved in the repression of “late” genes (see [Lambda Phage](#))] and  $P_{int}$ , from which the mRNA for integrase protein (see text below) is made (6). The *cI* protein is both a repressor (to which it owes its name **lambda repressor**) and an activator of transcription.  $O_L$ -bound *cI* protein represses RNA synthesis from  $P_L$ , and  $O_R$ -bound *cI* represses utilization of the  $P_R$  promoter (and synthesis of *cro* mRNA), by denying RNA polymerase proper access to these promoters. An additional effect of *cI* bound at  $O_R$  is to facilitate formation of a functional RNA polymerase-promoter complex at  $P_{RM}$  (promoter for repressor maintenance) by favorably contacting  $P_{RM}$ -bound RNA polymerase (7). This results in greatly enhanced initiation rates for RNA synthesis from this promoter.

The mRNA synthesized from the  $P_{RE}$  promoter provides an initial burst of *cI* protein (8), needed to jump-start transcription from  $P_{RM}$ . Once  $P_{RM}$  has been activated, the amount of *cI* mRNA synthesized is sufficient to establish levels of *cI* protein in the cell sufficient for maintaining both the

activated state of  $P_{RM}$  and the repressed state of  $P_R$  and  $P_L$ . The double-helical  $\lambda$  DNA is “integrated” into the bacterial DNA with the aid of the integrase protein (encoded by the *int* gene of  $\lambda$ ) and a bacterial protein originally discovered by its required participation in this event, integration host factor (IHF). This is a very precise process, always involving the joining of the same regions of bacterial and bacteriophage  $\lambda$  DNA, known as *attB* and *attP* (attachment sites on the bacterial and phage genomes), respectively. First, *attB* and *attP* are positioned next to each other. Next, in a concerted process, both genomes are cut within the *attB* and *attP* DNA sequences, and the ends are exchanged and ligated together so that the bacteriophage becomes a prophage, embedded in the bacterial DNA. Because the *attP* site (Fig. 1A) is different from the overlapping DNA sequences found at the ends of linear  $\lambda$  DNA in bacteriophage particles (see [Lambda Phage](#)), the map of genes on the prophage is circularly permuted with respect to that of the linear bacteriophage DNA (compare Figs. 1A and 1B). The  $\lambda$ -mediated repression of  $P_L$  and  $P_R$  effectively silences the integrated bacteriophage  $\lambda$  genome, with  $\lambda$  itself the only bacteriophage protein still being synthesized (from mRNA initiated at  $P_{RM}$ ). This is also the basis for the immunity of the lysogenic cell to subsequent  $\lambda$  infections: any invading  $\lambda$  DNA will similarly be silenced by  $\lambda$  protein binding to its  $O_L$  and  $O_R$  regions, resulting in repression of transcription initiation from the promoters  $P_L$  and  $P_R$ .

The lytic/lysogenic decision is affected by the growth conditions of the infected cell; if it finds itself in a nutrient-rich environment containing the preferred carbon source for *E. coli*, glucose, conditions are favorable for the generation of progeny, and the probability is high that the infection will follow the lytic pathway. This is the case for most laboratory media in which bacteria are grown. Conversely, when in nutrient-poor medium, the bacterium will probably be able to support the production of only a limited number of progeny, and mechanisms are in place to increase greatly the probability for establishment of lysogeny. Growth conditions are sensed by a series of events, as outlined in Figure 3, ultimately affecting levels of  $\lambda$  protein in the infected cell (9). The  $\lambda$  protein is highly sensitive to **protein degradation** by an *E. coli* [proteinase](#), HflB (or FtsH) (10), which is most active in fast-growing cells. An added twist is that  $\lambda$  is so unstable that, unless it is stabilized by bound  $\lambda$  protein, it will not accumulate to any appreciable levels even when HflB activity is low. Thus the expression of  $\lambda$  protein is also required for the establishment of lysogeny. When the bacterium finds itself in nutrient-rich medium, HflB proteinase is active, and  $\lambda$  is rapidly degraded, before it can activate the  $P_{RE}$  promoter. Then the promoters  $P_L$  and  $P_R$  remain active; any residual activity of  $P_{RM}$  will be blocked by *cro* protein, so no significant  $\lambda$  levels build up, and the lytic pathway will be followed. Conversely,  $\lambda$  levels will be high in poor growth medium, setting in motion the course of events commencing with activation of the  $P_{RE}$  promoter to provide the initial burst of  $\lambda$  protein, which favors the lysogenic state. The implementation of the lytic/lysogenic decision occurs at a short region of DNA encompassing promoters  $P_R$  and  $P_{RM}$ , which has been called a “genetic switch” by Ptashne (3, 4). The molecular details of this switch ensure that it is essentially irreversible: In a lysogen,  $\lambda$  protein not only activates synthesis of its own mRNA from  $P_{RM}$ , but it also represses  $P_R$ , from which the mRNA for *cro*, a repressor of  $P_{RM}$ , would be made. Conversely, once the cell has embarked upon the lytic pathway, *cro* represses the very promoter,  $P_{RM}$ , from which mRNA for the synthesis of  $\lambda$ , the repressor of *cro* mRNA synthesis from  $P_R$ , would be initiated.

**Figure 3.** Summary of the effect of growth conditions on the lytic/lysogenic decision. The bottom two lines of the table, when read from left to right, show two possible courses of events. + indicates high cellular levels of protein; – indicates low levels, or none present.

| Medium        | Cell growth | HflB protease | cII | cI | $\lambda$ development |
|---------------|-------------|---------------|-----|----|-----------------------|
| Nutrient-rich | Fast        | Active        | -   | -  | Lytic                 |
| Nutrient-poor | Slow        | Inactive      | +   | +  | Lysogenic             |

It would not be in the best interest of bacteriophage  $\lambda$  to maintain the lysogenic relationship under conditions where the life of its host were in danger. Indeed, reversal of the lysogenic state is initiated when significant damage has been inflicted upon cellular DNA, such as may be the case upon exposure to UV or X-ray irradiation or to chemicals that covalently modify DNA (see **DNA damage**). Such conditions are much more likely to cause irreparable damage to the DNA of the bacterium than to the integrated bacteriophage DNA, because the former contains 80 times the number of base pairs and presents a correspondingly larger target. As a consequence, the prophage DNA may remain unscathed, even if the host DNA were damaged beyond repair. DNA damage results in the generation of fragments of single-stranded DNA, in the presence of which the bacterial **RecA protein** facilitates the inactivation of cI protein by a self-inflicted cleavage of its polypeptide chain (11). The resulting relief of repression at  $P_L$  and  $P_R$  sets into motion a series of events leading to the excision of the bacteriophage DNA from the bacterial genome. This is accomplished by three proteins: In addition to the  $\lambda$  integrase and bacterial IHF proteins, which were responsible for the integration reaction, the excisionase protein, encoded by the *xis* gene, is also required. Further development proceeds along a lytic-like pathway (see [Lambda phage](#)), eventually leading to cell death and the release of progeny phages.

## 2. Other Temperate Bacteriophages

Several so-called lambdoid phages, with extensive sequence [homology](#) to bacteriophage  $\lambda$ , have been characterized that behave similarly to phage  $\lambda$  in the establishment of lysogeny. The genomes of many other bacteriophages, with less homology to bacteriophage  $\lambda$  DNA, also integrate site-specifically into the bacterial genome, but a notable exception is the *E. coli* bacteriophage “mu.” Just like phage  $\lambda$ , bacteriophage mu is able to subvert the *E. coli* synthetic machinery to express and replicate its own genomic DNA. Again, there is mutually exclusive expression of the mu repressor, c (the equivalent of the phage  $\lambda$  cI protein), which silences the genome, and a protein specific for the lytic pathway, ner (which plays a role analogous to that of cro in bacteriophage  $\lambda$ ). However, the factors determining whether one or the other will predominate are not yet fully understood. Bacteriophage mu is unique among temperate bacteriophages, in that integration of its DNA into the bacterial genome is an integral part of both its lytic and its lysogenic modes of development. The initial integration is conservative, that is, the actual DNA molecule that entered the cell becomes part of the *E. coli* genome, similar to integration of bacteriophage  $\lambda$  DNA. Integration always involves the attachment of the same regions of the mu genome (near the ends of the linear bacteriophage mu DNA injected into *E. coli*) to the bacterial DNA, but in contrast to bacteriophage  $\lambda$ , the sites on the *E. coli* genome where integration takes place are almost entirely random. It is the distinct ability to integrate almost anywhere in the bacterial DNA that has given mu its name: In lysogenic strains, the integrated mu DNA can disrupt various genes of the bacterium and thus inactivate them. Bacteriophage mu has been used as a mutator agent to generate mutants useful for genetic studies.

Bacteriophage mu also stands out in that the divergence between the lytic and lysogenic pathways takes place subsequent to the integration event. If the integrated bacteriophage DNA is silenced by the c repressor, the cell becomes lysogenic. On the other hand, the ner protein may predominate, preventing accumulation of c protein; then phage expression and replication may occur to initiate lytic development. In this case, replicative integration takes place, that is, copies of the mu genome integrate at different sites on the bacterial genome, in steadily increasing numbers. Eventually they will be excised (ie, cut out of the bacterial DNA), packaged into bacteriophage particles, and released upon cell lysis and death. Mu integration has attracted considerable interest because, in

this respect, mu DNA behaves similarly to a class of mobile genetic elements called “[transposons](#).”

In addition to *E. coli*, a large number of bacterial species have been found to harbor prophages, and thus to be lysogenic for the corresponding bacteriophages. Examples include pathogenic as well as nonpathogenic bacteria, **gram positive** as well as **gram negative**, such as *Salmonella*, *Staphylococcus aureus*, and various *Streptococci* and *Mycobacteria*. With few exceptions (eg, **P22 phage**, a  $\lambda$ -like bacteriophage of *Salmonella*), however, the details of the relationship between these bacterial hosts and their bacteriophages are yet to be determined. Finally, some phages establish a relationship with their host that resembles lysogeny, in that the cell is not lysed. This is the case for the filamentous bacteriophages of *E. coli* (eg, [M13 phage](#)), whose genomes are single-stranded DNA. Subsequent to infection, progeny virus is continuously produced and released, without cell lysis, from the infected cells, which keep dividing. There is, however, no repression of genes of the bacteriophage genomic DNA, nor does this DNA integrate into the host genome, and the infected cells are not considered to be lysogens.

### Bibliography

1. A. D. Johnson, A. R. Poteete, G. Lauer, R. T. Sauer, G. K. Ackers, and M. Ptashne (1981) *Nature* **294**, 217–223.
2. M. Ptashne, A. Jeffrey, A. D. Johnson, R. Maurer, B. J. Meyer, C. O. Pabo, T. M. Roberts, and R. T. Sauer (1980) *Cell* **19**, 1–11.
3. M. Ptashne, A. D. Johnson, and C. O. Pabo (1982) *Sci. Am.* **247**, 128–140.
4. M. Ptashne (1986) *A genetic switch*, Cell Press and Blackwell Scientific Publications, Oxford, England.
5. M. C. Shih and G. N. Gussin (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6432–6436.
6. Y. S. Ho and M. Rosenberg (1985) *J. Biol. Chem.* **260**, 11838–11844.
7. M. Li, H. Moyle, and M. M. Susskind (1994) *Science* **263**, 75–77.
8. U. Schmeissner, D. Court, H. Shimatake, and M. Rosenberg (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3191.
9. M. A. Hoyt, D. M. Knight, A. Das, H. I. Miller, and H. Echols (1982) *Cell* **31**, 565–573.
10. A. Kihara, Y. Akiyama, and K. Ito (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5544–5549.
11. B. Kim and J. W. Little (1983) *Cell* **73**, 1165–1173.

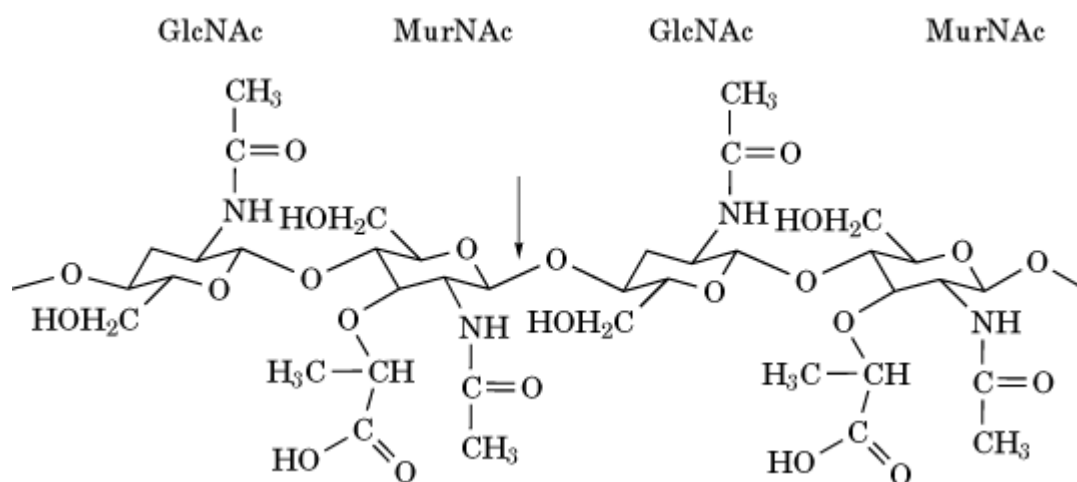
### Suggestions for Further Reading

12. L. J. Beamer and C. O. Pabo (1992) Refined 1.8 Å crystal structure of the  $\lambda$  repressor–operator complex. *J. Mol. Biol.* **227**, 177–196.
13. R. W. Hendrix, J. W. Roberts, F. W. Stahl, and R. Weisberg (1983) *Lambda IIs*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
14. S. E. Nunes-Duby, D. Yu, and A. Landy (1997) Sensing homology at the strand-swapping step in lambda excisive recombination. *J. Mol. Biol.* **272**, 493–506.
15. W. A. Rees, S. E. Weitzel, A. Das, and P. H. von Hippel (1997) Regulation of the elongation-termination decision at intrinsic terminators by antitermination protein N of phage lambda. *J. Mol. Biol.* **273**, 797–813.
16. P. A. Rice, S. Yang, K. Mizuuchi, and H. A. Nash (1996) Crystal structure of an IHF–DNA complex: a protein-induced DNA U-turn. *Cell* **87**, 1295–1306.
17. B. Stern, A. Molven, and D. Kamp (1995) Conserved sequence patterns in phages mu and  $\lambda$  DNA. *Biochim. Biophys. Acta* **1264**, 115–120.
18. N. Symonds, A. Toussaint, P. van der Putte, and M. M. Howe (1987) *Phage mu*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

## Lysozymes

The first lysozyme was discovered in 1921 by Sir Alexander Fleming, who received the Nobel Prize in 1945 for also finding the first **antibiotic**, [penicillin](#). Lysozymes are widely distributed in nature, but the term “lysozyme” most usually means the protein from hen egg white. This is because it has been studied so thoroughly, in part because it is so abundant and easy to obtain in quantity. This lysozyme (EC 3.2.1.17) is the lytic [enzyme](#) that cleaves the glycosidic bond between *N*-acetylmuramic acid and *N*-acetylglucosamine units in the polysaccharides of bacterial cell walls (Fig. 1). Consequently, lysozyme is also known as muramidase. It is a small protein ( $M_r = 14,307$ ), basic ([isoelectric point](#) = 11.5), and stable (melting temperature under physiological conditions is 78°C) (see [Protein Stability](#)). It was the first enzyme containing all the 20 usual [amino acids](#) to be sequenced (1, 2). It was the first enzyme whose precise structure was determined by [X-ray crystallography](#) and whose reaction mechanism was elucidated as a result (3, 4). Thus, hen lysozyme became a model enzyme in biochemistry and biology.

**Figure 1.** Structure of the bacterial cell wall, natural substrate for hen lysozyme, an alternating copolymer of *N*-acetylmuramic acid (MurNAc) and *N*-acetylglucosamine (GlcNAc) residues. The arrow shows the bond that is cleaved by lysozyme.



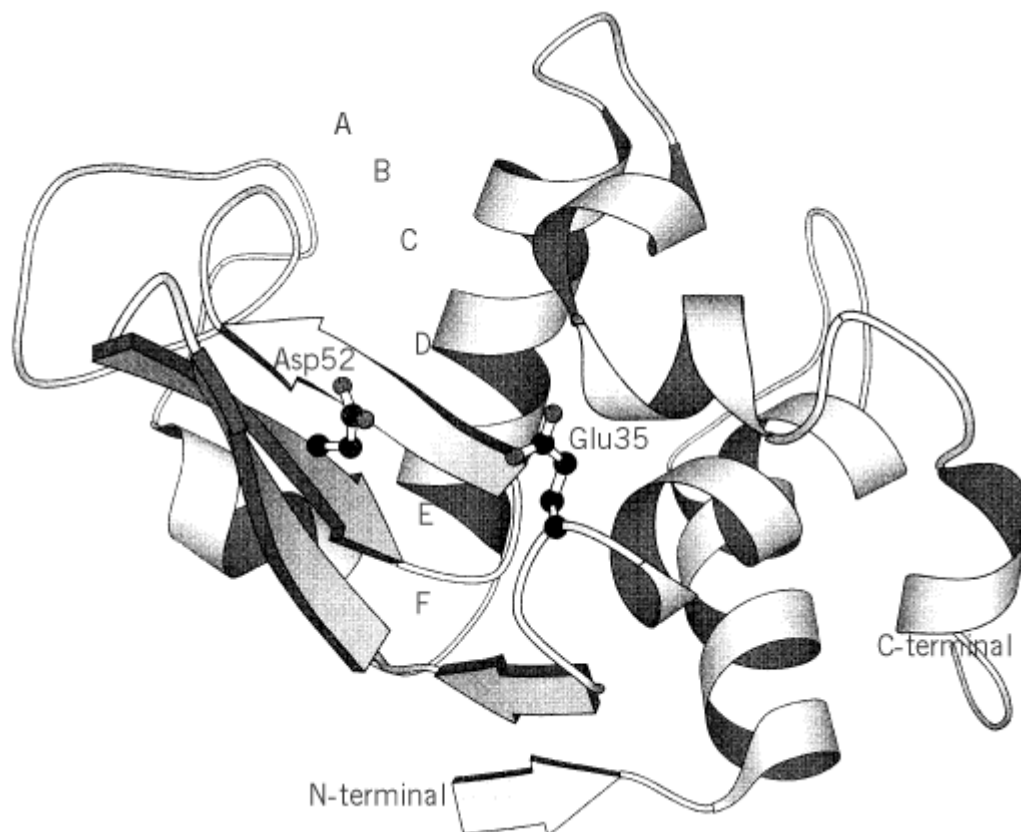
### 1. Structure

The primary structure of lysozyme was determined by Canfield (1) and Jolles et al. (2) in 1963. It consists of 129 amino acid residues, including 10 **carboxyl** and 7 **amino** groups, 11 [arginine](#) residues, 6 [tryptophan](#) residues, and 4 [disulfide bonds](#). The three-dimensional structure of lysozyme was elucidated by Sir David Phillips' group between 1961 and 1966 (3, 4) and is shown diagrammatically in Figure 2. The molecular dimensions of the protein are approximately  $3.6 \times 4.5 \times 4.2$  nm. The structure adopts a mixed  $\alpha/\beta$  fold (see [Protein Structure](#)). A deep [active site](#) cleft divides the molecule into two **domains**; one is almost entirely  $\beta$ -sheet structure (residues 40–85), the other is  $\alpha$ -helix-rich and consists of the *N*- and *C*-terminal segments (residues 1–39 and 101–129). The two domains are linked by an  $\alpha$ -helix (residues 89–99). Four disulfide bonds are formed



between [cysteine](#) residues 6–127, 30–115, 64–80 and 76–94.

**Figure 2.** The three-dimensional structure of hen lysozyme. Residues Glu35 and Asp52 are shown as the catalytic groups. The letters A–F in the active-site cleft indicate the binding subsites. The schematic figure was produced with the program MOLSCRIPT.



## 2. Enzymatic activity

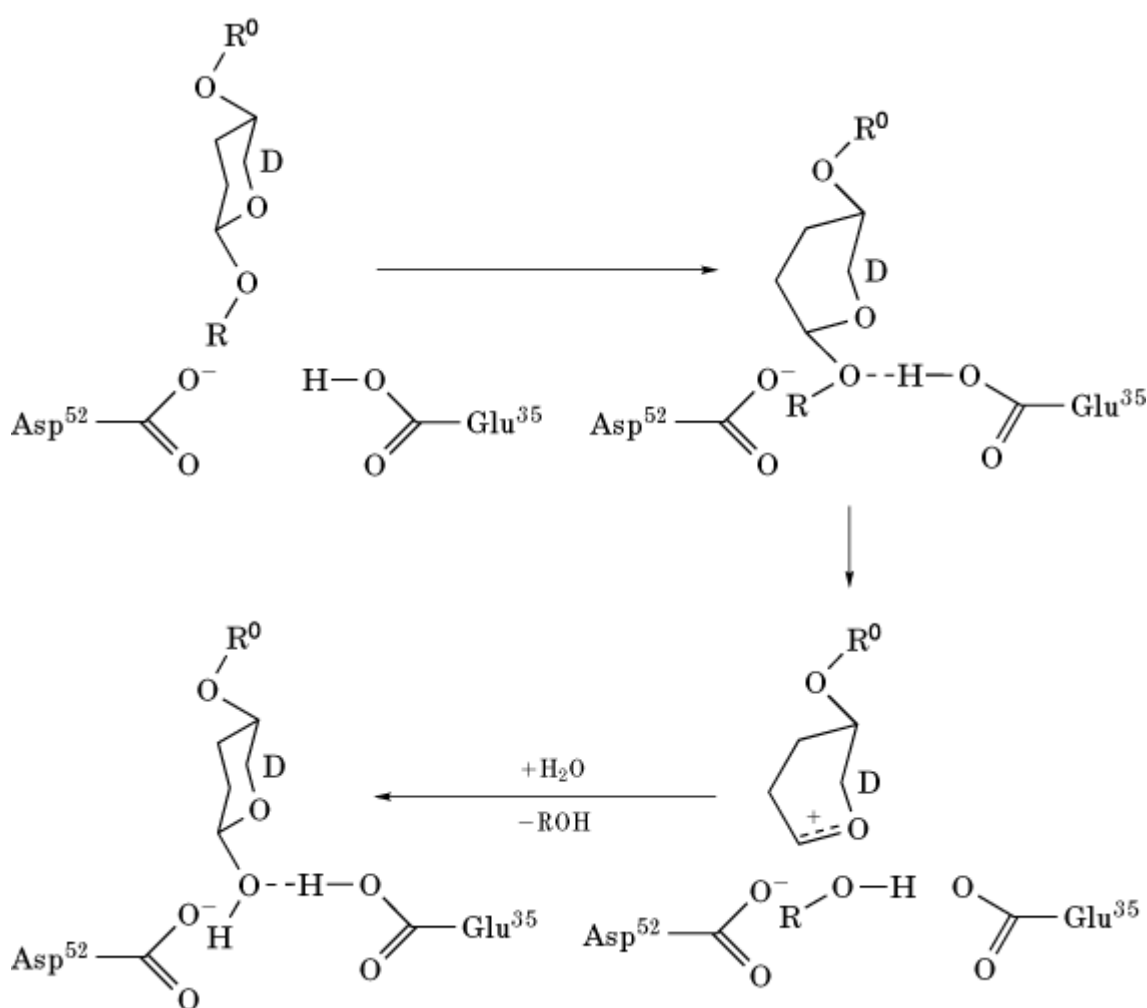
Lysozyme hydrolyses the b-1,4-glycosidic bond between *N*-acetylmuramic acid and *N*-acetylglucosamine residues in polysaccharides (Fig. 1). Thus, lysozyme activity can be determined by the decrease in the turbidity of a suspension of dried cells of UV-inactivated *Micrococcus luteus* (*lysodeikticus*) (5). A more convenient method employing cells labeled with Remazol brilliant blue has been invented (6). Lysozyme also hydrolyzes b-1,4-glycosidic bonds between *N*-acetylglucosamines, and chitin derivatives can be employed as substrate. The activity is measured colorimetrically by determining the reducing sugars produced from chitin solubilized by hydroxyethylation (glycol chitin) (7). *N*-Acetylglucosamine (GlcNAc) oligomers ((GlcNAc)<sub>*n*</sub>), or their *p*-nitrophenylacetate derivatives, are also employed as substrates.

A reaction mechanism was deduced for lysozyme from model building the enzyme–substrate complex on the basis of the structure of Figure 2 and that of its complex with (GlcNAc)<sub>3</sub>. The active-site cleft can bind six GlcNAc residues, and the binding sites for each were designated A, B, C, D, E and F, from top to bottom (from the nonreducing end of (GlcNAc)<sub>6</sub> to its reducing end).

Cleavage of the substrate occurs at the glycosidic bond between the residues bound to the D and E sites. The sugar ring at D site must be distorted from the normal chair conformation to a half-chair form to achieve good contacts with the enzyme. This distortion is thought to be used in the [transition state](#) of the reaction. Lysozyme residues Glu35 and Asp52 are in position near the cleavage site and

are considered to be the primary catalytic residues. Glu35, which lies in a hydrophobic environment, participates in catalysis in its protonated form, while Asp52, which is in a hydrophilic environment, does so in its nonprotonated form. As shown in Figure 3, Glu35 donates its proton to the oxygen linking the D and E sugars. The bond between this oxygen and C1 on the D sugar is cleaved. A carbonium ion is formed on the D sugar, and this ion is stabilized by Asp52 and by the formation of an oxocarbenium ion. Formation of the planar oxocarbenium ion is favored by the distortion of the D sugar. As a result of this stabilization of the transition state, lysozyme can rapidly hydrolyze glycosidic bonds (4). This is a good example of how enzymes generally stabilize the transition state of their enzymatic reaction (see [Catalysis](#)). The reaction is completed by one-half of the cleaved product dissociating, and water replacing it and reacting with the bound intermediate. Extraneous oligosaccharides can also bind to the vacated sites and can compete with water for the last step. As a result, lysozyme has considerable transglycosylation activity.

**Figure 3.** Reaction mechanism of lysozyme.



### 3. Physiological Function

Lysozyme is widely distributed in nature, so some biological functions (8) other than its primary anti-bacterial action are expected. Lysozyme is effective against bacterial and viral infections and also shows antiphlogistic activity in a number of pathological conditions. Immune-stimulating actions for lysozyme are widely found. Owing to these effects, lysozyme is employed as a medicine and as a food additive. Lysozymes are used as digestive enzymes in some organisms, such as

ruminants and insects.

#### 4. Comparative Biochemistry

Lysozymes are widely distributed in nature, from humans to viruses and plants. Almost 100 lysozymes have been sequenced. On the basis of their sequence homology, they are classified into four classes: chicken-type (c-type), viral-type (v-type), goose-type (g-type), and others. C-type lysozymes consist of 119–130 amino acid residues, and they include hen and human lysozymes. V-type lysozymes consist of approximately 140–350 residues and are found mainly in **bacteriophages**. G-type lysozymes are widely present in birds, and some birds have both c- and g-type lysozymes. Four g-type lysozymes have been sequenced, including those of black swan, goose, chicken, and ostrich; they have 185 amino acid residues. These three types of lysozyme are very different in size and sequence, but they are quite homologous in their three-dimensional structures.

C-type lysozymes can be subdivided into three categories: conventional, digestive, and Ca-binding lysozymes. Representative conventional lysozymes are those from hen egg white and human. Digestive lysozymes are devised to act at low pH and have low isoelectric points. Ca-binding lysozymes are interesting for their relationship to a homologous protein, **α-lactalbumin**, which also is a [calcium-binding protein](#) and has a three dimensional structure similar to that of c-type lysozymes. Ca-binding lysozymes include those from echidna, pigeon, and some milk. The Ca-binding loops are located at residues 82–91. These diversities in structure and function make lysozymes interesting proteins with which to analyze **phylogenetic** relationships.

Human lysozyme has 130 residues. Its sequence is 60% identical to that of hen lysozyme, and its properties and structure are also very similar. T4 phage lysozyme has 164 residues. It is capable of cleaving only glycosidic bonds next to *N*-acetylmuramic acid residues that are substituted with peptide side chains, and it also requires the presence of the *N*-acetamido group for activity. It is active on chloroform-treated *Escherichia coli* cells, and it shows no transglycosylation reaction. The catalytic residues corresponding structurally to Glu35 and Asp52 of hen lysozyme are Glu11 and Asp20, respectively, although the catalytic role of the latter has been questioned. Goose lysozyme has Glu73 as the catalytic residue corresponding to Glu35 in hen lysozyme, but apparently has no equivalent of Asp52. T4 lysozyme has been the subject of extensive [protein engineering](#) studies by Matthews et al. (9).

#### Bibliography

1. R. E. Canfield (1963) *J. Biol. Chem.* **238**, 2698–2707.
2. J. Jolles, J. Jauregui-Adell, I. Berner, and P. Jolles (1963) *Biochim. Biophys. Acta* **78**, 668–689.
3. C. C. F. Blake et al. (1961) *Nature* **206**, 757–761.
4. D. C. Phillips (1966) *Sci. Am.* **215**, 78–90.
5. D. Sugar (1952) *Biochim. Biophys. Acta* **8**, 302–309.
6. Y. Ito, H. Yamada, and T. Imoto (1992) *Chem. Pharm. Bull.* **4**, 1523–1526.
7. H. Yamada and T. Imoto (1981) *Carbohydr. Res.* **92**, 160–162.
8. G. Sava (1996) in *Lysozymes: Model Enzymes in Biochemistry and Biology*, P. Jolles, ed., Birkhauser, pp. 434–449.
9. B. W. Matthews (1993) *Annu. Rev. Biochem.* **62**, 139–160.

#### Suggestions for Further Reading

10. T. Imoto et al. (1972) "Vertebrate lysozymes, in" *The Enzymes*, 2nd ed., (P. D. Boyer, ed., Academic Press, New York, Vol. 7, pp. 665–868.
11. P. Jolles (ed.) (1996) *Lysozymes: model enzymes in biochemistry and biology*, Birkhauser Verlag.

## M13 Phage

From an incubation of waste water samples with several types of *Escherichia coli* bacteria, Hofschneider (1) discovered in 1963 some new types of bacteriophages, among which was one designated M(unich)13. Since then, a large number of bacteriophages that constitute the family of filamentous phages (Inoviridae) have been isolated from different [gram-negative bacteria](#). Filamentous phages belong to the smallest bacteriophages known and differ from most other bacterial viruses in that they are being reproduced continuously upon infection without lysis of their host cells; infected cells continue to grow and divide, although with a reduced rate.

Filamentous phage efficiently infect only host cells that express (sex) [pili](#), which serve as receptor sites. Based on their pilus specificity, different classes of evolutionary closely related filamentous phages can be distinguished. F-specific filamentous phages (the Ff group) only infect male *E. coli* cells carrying sex pili, encoded by an **F factor** (ie, a fertility or sex episome). Among these Ff bacteriophages, the almost identical strains M13, fd, and fl are the most extensively studied and, together with their *E. coli* host strains, the best characterized from biochemical, genetical, molecular biological, and biophysical points of view.

Because of its relative simplicity and the ease of being genetically manipulated, M13 bacteriophage has a very important role in investigating biological processes that extends far beyond its own life cycle. Apart from the still fruitful investigation of the phage structure and the studies of the processes involved in infection, replication, and assembly of the phage machinery itself, the practical use of M13 phage-derived plasmids as tools in biotechnology is especially enormous (2). The advantages of the M13 phage-derived plasmids as cloning vectors (eg, M13mp18), with their ability to generate large quantities of single-stranded DNA molecules carrying foreign DNA, should be emphasized: Single-stranded DNAs are the templates of choice for [DNA sequencing](#), labeling, or [site-directed mutagenesis](#) (3). Furthermore, the phage plays an important role in displaying foreign peptides on the phage surface, when foreign DNA is fused in-frame to a gene coding for one of the coat proteins (see [Phage Display Libraries](#)). This technique enables the construction of extensive [libraries](#) that can be easily screened to select peptides with specific affinities or activities (4).

### 1. M13 Phage Particle Architecture

As a member of the filamentous phage family, the M13 phage nucleoprotein particle is long (900 nm) and thin (6.5 nm). The filament contains the viral [genome](#): a twisted closed circular single-stranded DNA molecule of 6407 nucleotides. The DNA genome is protected by a long cylindrical protein coat, consisting predominantly of about 2700 copies of the major coat protein (gp8) and capped at the ends with about five copies each of four different minor coat proteins (5).

The strong native architecture of the phage coat seems to be maintained primarily by **hydrophobic** interactions between the individual coat proteins. The major coat proteins form a tube around the viral DNA, in an overlapping helical array, with a flexible [amphipathic](#) N-terminus located at the outside of the coat and the basic C-terminus interacting with the DNA at the inside of the coat. The hydrophobic domain of the major coat protein is located in the central section of the protein sequence, and it interlocks the coat protein with its neighboring coat proteins (6). As a result, M13 phage is very stable and resistant to most [proteinases](#), salts, extreme pH values, heating, (phospho) lipids, and [detergents](#) (but not the strong anionic detergent **SDS**). The phage particle is sensitive to mechanical stress (ultrasonication) and to chloroform, which in turn makes it susceptible to detergents (7, 8).

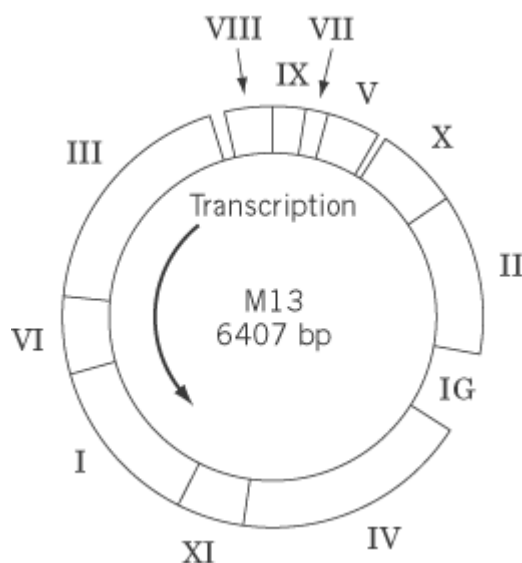
The packing of the DNA in the protein tube is determined by an [electrostatic interaction](#) between phosphate groups of the DNA and [lysine](#) residues of the major coat protein. Elongating the viral genome by adding foreign DNA to the intergenic region ( [3](#)), or reducing the number of positively charged lysine residues of the **DNA-binding** part of the major coat protein ( [9](#)), simply results in elongated phage filaments.

The end of the phage that emerges first from the host cell contains a packaging signal as an imperfect double-stranded DNA hairpin loop of the intergenic region. This end of the phage is capped by the minor coat proteins gp7 and gp9 (five copies of each). Both gp7 and gp9 are required for proper initiation of the phage assembly process. At the other end of the phage, five copies each of the other minor coat proteins gp3 and gp6 are located. These proteins are required for properly terminating the assembly process and contribute to phage particle stability ( [10](#)). The minor coat protein gp3 also contains functional domains enabling host-specific adsorption and subsequent penetration of the viral genome during a next round of infection.

## 2. Genetic Organization

The complete DNA sequence of M13 and related phages is known ( [11](#), [12](#)). The genome comprises nine open reading frames, coding for 11 identified proteins (Fig. [1](#)), and two intergenic regions. It should be noted that genes X and XI are contained completely within genes II and I, respectively. The phage-encoded proteins are organized into three functional clusters: gp2, gp10, and gp5 are involved in [DNA replication](#); the morphogenetic gp1, gp4, and gp11 are required for phage assembly and secretion; and gp3, gp6, gp7, gp8, and gp9 comprise the phage structural major and minor coat proteins. The intergenic region between genes IV and II does not code for proteins, but it contains signals for the initiation of synthesis of viral and [complementary DNA](#) strands and for phage assembly, termination of RNA synthesis, and the expression of the following genes II, X, V, VII, IX, and VIII. Furthermore, this region has been used as the site for the insertion of foreign DNA into the M13mp vectors. A smaller intergenic region between genes VIII and III also contains a [transcription](#) termination signal and a **promoter** for the expression of the following genes III, VI, I, XI, and IV. An excellent and exhaustive review on the genetic organization of M13 and related filamentous phages was written by Model and Russel in 1988 ( [13](#)) (see also Suggestions for Further Reading).

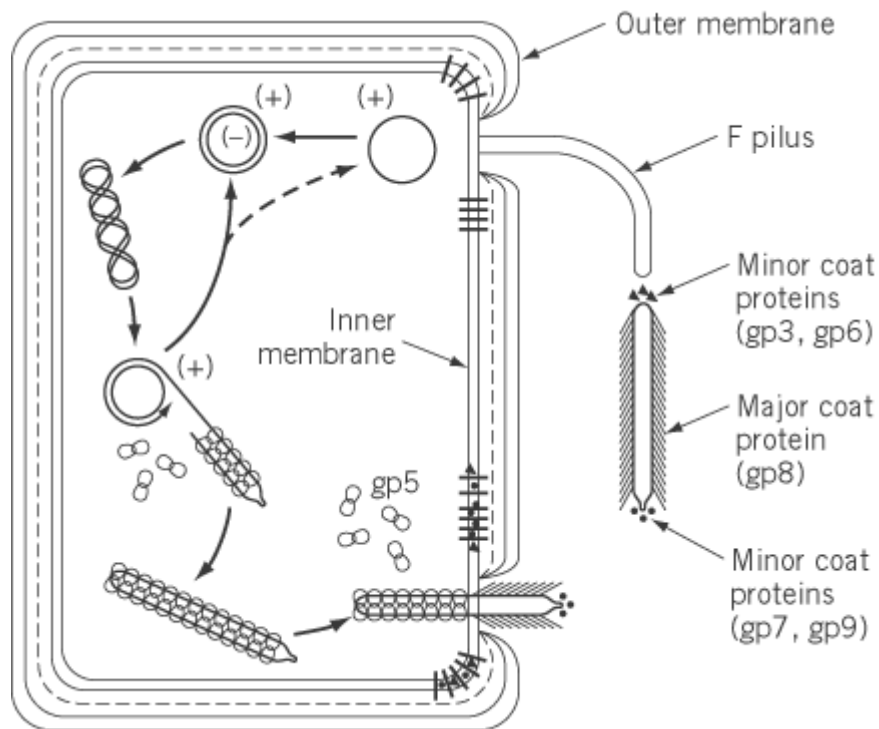
**Figure 1.** The genome of wild-type bacteriophage M13. The approximate locations of the genes are indicated by Roman numerals.



### 3. Reproductive Cycle

The reproductive life cycle of bacteriophage M13 (depicted schematically in Fig. 2) starts with the adsorption of the phage to the tip of the F pilus of the host *Escherichia coli*, specifically mediated by gp3, the viral adsorption protein (14-16). The phage is subsequently transported to the cell surface, probably as a consequence of pilus retraction. Some host-encoded outer-membrane proteins (TolQ, TolR, TolA) are reported to mediate uncoating and DNA penetration (17, 18), probably at adhesion zones between the inner and outer membranes. Also, the viral gp3 is thought to play a role in facilitating the entry of the phage genome into the cell by the formation of a pore, enabling the entire host cell envelope to pass through (19). During the infectious entry of the phage, the hydrophobic gp6, which probably functions as a sort of sealing agent to preserve phage stability (5), is lost from the phage, resulting in a destabilized nucleoprotein particle ready for further disassembly. The major coat protein, and probably also the minor coat proteins, are stripped off from the phage and subsequently deposited in the inner membrane (20). During this process, the phage genome is injected into the cell cytoplasm.

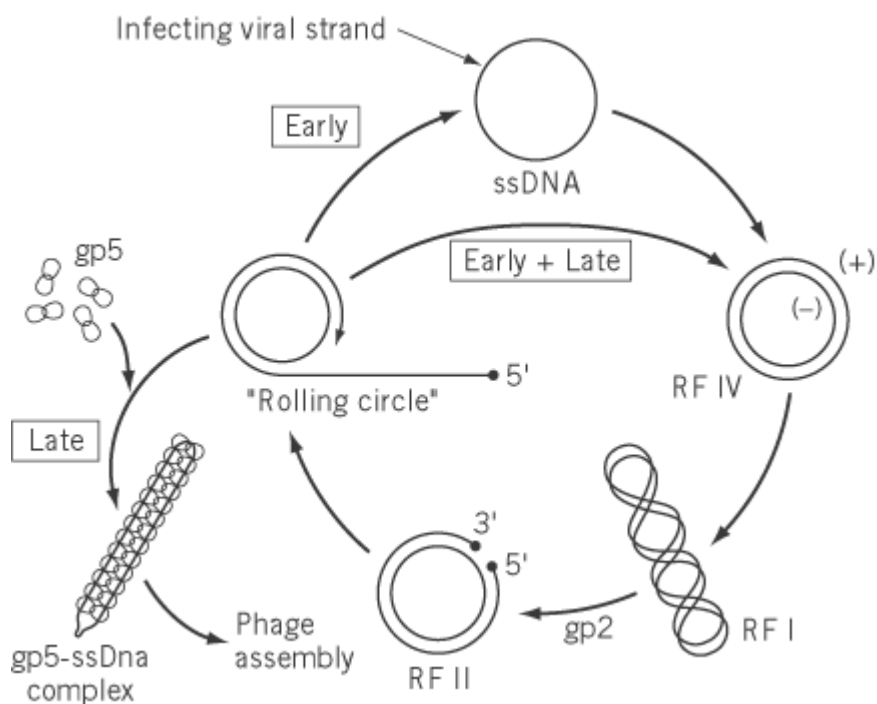
**Figure 2.** Schematic illustration of some features of the reproductive life cycle of bacteriophage M13.



Following entry into the cytoplasm, the single-stranded phage genome is converted into a double-stranded **supercoiled** covalently closed **replicative form** (RF) DNA molecule by the action of host **enzymes**. This parental RF is the template for subsequent phage DNA replication by a **rolling-circle mechanism** (21), as well as for the transcription and translation of the 11 phage-encoded proteins (Fig. 3). At an early stage of infection, a pool of progeny RF molecules is generated under control of phage-encoded gp2. At a later stage, gp2, gp5, and gp10 regulate the levels of phage RF DNA and single-stranded DNA in the cell, which determines the rate of phage **protein biosynthesis** and the rate of phage assembly and concomitant extrusion, respectively. Upon binding of gp5, the newly synthesized viral strands are separated from the DNA replication machinery. The resulting single-stranded viral genome is encapsidated by the gp5 homodimers (22), but leaving the double-stranded

packaging signal exposed, and is then ready to be assembled into progeny phage particles at the membrane-bound assembly site.

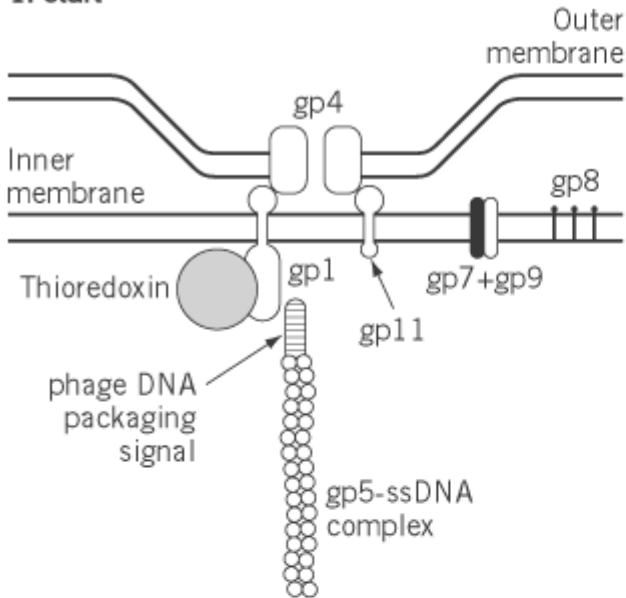
**Figure 3.** Schematic representation of the DNA replication cycle of bacteriophage M13. The early stage is characterized by the predominant synthesis of the replicative form (RF). Later in the infection cycle, the viral strand is withdrawn from the DNA replication machinery by gp5 and directed toward the phage assembly site.



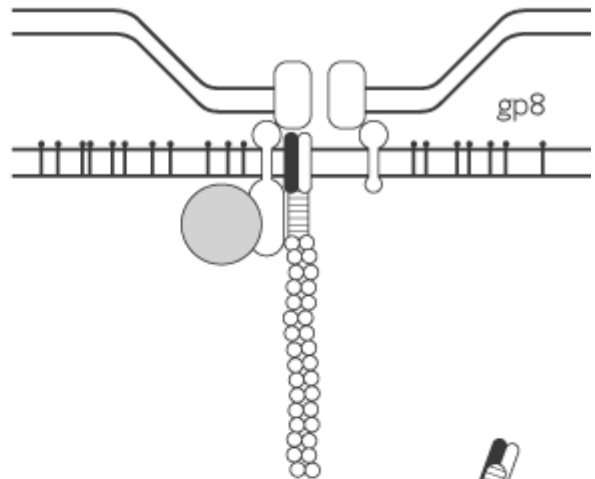
In contrast to most other bacteriophages, which are assembled intracellularly and are extruded by cell lysis, the morphogenesis of filamentous bacteriophage is concomitant with extrusion and without lysis of the host cell. Infected cells continue to grow and divide while extruding progeny phage particles into the medium. Filamentous phage emerge the infected cells at the assembly sites, which resemble adhesion zones between the inner and outer membranes (23). Three morphogenetic proteins gp1, gp4, and gp11 (the latter being an in-frame translated C-terminal part of the gp1), as well as at least one host-encoded protein ([thioredoxin](#)), are required for phage assembly. The proteins gp1 and gp11 are both associated with the inner membrane, and mixed multimers of them probably constitute the part of the extrusion channel that passes through the inner membrane (24, 25). A homomultimer of gp4 is thought to constitute the part of the extrusion channel passing through the outer membrane (26). As depicted in Fig. 4 these proteins interact and form a complex that could span the entire cell envelope.

**Figure 4.** Schematic illustration of the different stages of the membrane-associated assembly process of bacteriophage M13.

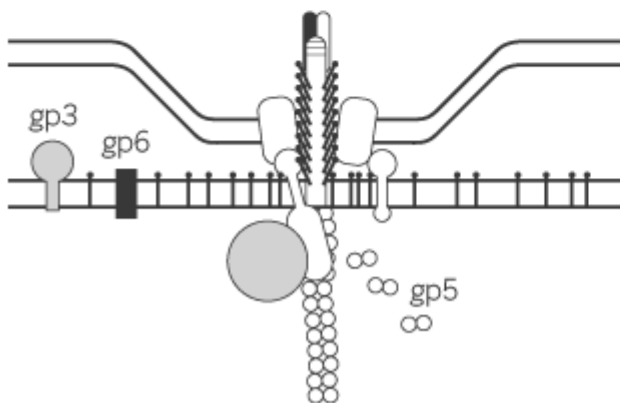
### 1: start



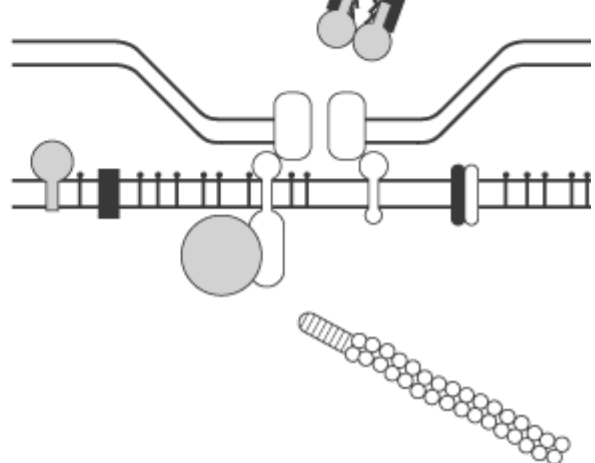
### 2: initiation



### 3: elongation



### 4: termination



On initiating phage assembly, all the constituents come together around the inner membrane-associated gp1, which seems to be the “spider in the web” (see Fig. 4). All newly synthesized coat proteins are inserted into the inner membrane, where they are stored prior to being used in the assembly process (27). The exposed double-stranded DNA packaging (or morphogenetic) signal on the gp5-ssDNA complex interacts with gp1, which in turn triggers the formation and opening of the extrusion channel. Phage assembly might also be initiated by binding of gp7 and gp9 to the exposed double-stranded DNA packaging signal, probably mediated by gp1 (28). Furthermore, the association of thioredoxin with gp1 is a prerequisite for the elongation of the nascent phage. However, the main process during assembly of the new phage particle is the replacement of each gp5 monomer in the gp5-ssDNA complex by about two gp8 major coat-protein molecules, which protects the phage genome outside the cell. For this purpose, the initially membrane-bound gp8 should specifically recognize gp1 or gp11 in the assembly site.

Upon separating from the membrane lipids, and undergoing a conformational change from an L-



shaped membrane-bound state to an aligned helical molecule, the gp8 molecules fit together with the other gp8 molecules, which all wrap around the phage DNA, thereby forming the coat of the new assembling phage. Meanwhile, the already assembled part of the new phage is extruded. The final, terminating step in assembly occurs by addition of gp3 and gp6 to close the phage filament, which is then released into the medium. Next, a voyage to another host cell can start.

#### 4. Current State of Research, Future Prospects, and Implications of Interest

The information gained about the relatively simple, and therefore extensively characterized, filamentous M13 phage structure and reproductive cycle is very important for biology, because it can be translated to other more complex biological systems; the physical and chemical principles of interactions in biological processes will be the same. On the other hand, the versatile M13 phage is the system of choice to study fundamental macromolecular interactions, because many kinds of different processes are involved in filamentous phage infection, and ample detailed background information is already available. The study of the structure–function relationships of the various phage-encoded proteins, as well as the many mutual interactions in the formation of macromolecular assemblies, is therefore an attractive and promising field of research.

Employing state-of-the-art techniques in the fields of genetics, molecular biology, biochemistry, and biophysics, M13 phage and the related filamentous phages have had an enormous impact on the current research covering molecular recognition ([4](#), [16](#)), [virus structure](#) ([6](#), [29](#), [30](#)), [protein–protein interactions](#) ([31–33](#)), protein insertion across membranes ([34](#), [35](#)), DNA-binding proteins ([36](#), [37](#)), membrane protein assembly ([38–41](#)), and [protein secretion](#) ([42](#), [43](#)).

#### Bibliography

1. P. H. Hofschneider (1963) *Z. Naturforsch. B* **18**, 203–210.
2. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning; A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
3. J. Messing (1983) *Methods Enzymol.* **101**, 20–78.
4. G. P. Smith and V. A. Petrenko (1997) *Chem. Rev.* **97**, 391–410.
5. R. E. Webster and J. Lopez (1985) In *Virus Structure and Assembly* (S. Casjens, ed.), Jones & Bartlett, Boston, pp. 235–267.
6. D. A. Marvin, R. D. Hale, C. Nave, and M. H. Citterich (1994) *J. Mol. Biol.* **235**, 260–286.
7. J. Griffith, M. Manning, and K. Dunn (1981) *Cell* **23**, 747–753.
8. D. Stopar, R. B. Spruijt, C. J. A. M. Wolfs, and M. A. Hemminga (1998) *Biochemistry* **37**, 10181–10187.
9. G. J. Hunter, D. H. Rowitch, and R. N. Perham (1987) *Nature* **327**, 252–254.
10. J. Lopez and R. E. Webster (1985) *J. Bacteriol.* **163**, 1270–1274.
11. P. M. G. F. Van Wezenbeek, T. J. M. Hulsebos, and J. G. G. Schoenmakers (1980) *Gene* **11**, 129–148.
12. D. F. Hill and G. B. Petersen (1982) *J. Virol.* **44**, 32–46.
13. P. Model and M. Russel (1988) In *The Bacteriophages*, Vol. **2** (R. Calendar, ed.), Plenum Press, New York, pp. 375–456.
14. C. W. Gray, R. S. Brown, and D. A. Marvin (1981) *J. Mol. Biol.* **146**, 621–627.
15. J. Armstrong, R. N. Perham, and J. E. Walker (1981) *FEBS Lett.* **135**, 167–172.
16. P. Holliger and L. Riechmann (1997) *Structure* **5**, 265–275.
17. R. E. Webster (1991) *Mol. Microbiol.* **5**, 1005–1011.
18. E. M. Click and R. E. Webster (1997) *J. Bacteriol.* **179**, 6464–6471.
19. G. Glaser-Wuttke, J. Keppner, and I. Rasched (1989) *Biochim. Biophys. Acta* **985**, 239–247.

20. I. Ohkawa and R. E. Webster (1981) *J. Biol. Chem.* **256**, 9951–9958.
21. D. S. Ray (1969) *J. Mol. Biol.* **43**, 631–643.
22. R. N. H. Konings, R. H. A. Folmer, P. J. M. Folkers, M. Nilges, and C. W. Hilbers (1995) *FEMS Microbiol. Rev.* **17**, 57–72.
23. J. Lopez and R. E. Webster (1982) *J. Virol.* **42**, 1099–1107.
24. J. K. Guy-Caffey, M. P. Rapoza, K. A. Jolley, and R. E. Webster (1992) *J. Bacteriol.* **174**, 2460–2465.
25. M. P. Rapoza and R. E. Webster (1995) *J. Mol. Biol.* **248**, 627–638.
26. M. Russel (1995) *Trends Microbiol.* **3**, 223–228.
27. H. Endemann and P. Model (1995) *J. Mol. Biol.* **250**, 496–506.
28. M. Russel and P. Model (1989) *J. Virol.* **63**, 3284–3295.
29. G. Kishchenko and L. Makowski (1997) *Proteins: Struct., Funct., Genet.* **27**, 405–409.
30. S. A. Overman, M. Tsuboi, and G. J. Thomas (1996) *J. Mol. Biol.* **259**, 331–336.
31. K. A. Williams, M. Glibowicka, Z. M. Li, H. Li, A. R. Khan, Y. M. Y. Chen, J. Wang, D. A. Marvin, and C. M. Deber (1995) *J. Mol. Biol.* **252**, 6–14.
32. D. Stopar, R. B. Spruijt, C. J. A. M. Wolfs, and M. A. Hemminga (1997) *Biochemistry* **36**, 12268–12275.
33. N. G. Haigh and R. E. Webster (1998) *J. Mol. Biol.* **279**, 19–29.
34. A. Kuhn (1995) *FEMS Microbiol. Rev.* **17**, 185–190.
35. M. Soekarjo, M. Eisenhawer, A. Kuhn, and H. Vogel (1996) *Biochemistry* **35**, 1232–1241.
36. A. P. M. Stassen, R. H. A. Folmer, C. W. Hilbers, and R. N. H. Konings (1995) *Mol. Biol. Rep.* **20**, 109–127.
37. R. H. A. Folmer, M. Nilges, C. H. M. Papavoine, B. J. M. Harmsen, R. N. H. Konings, and C. W. Hilbers (1997) *Biochemistry* **36**, 9120–9135.
38. F. J. M. Van de Ven, J. W. M. Van Os, J. M. A. Aelen, S. S. Wymenga, M. L. Remerowski, R. N. H. Konings, and C. W. Hilbers (1993) *Biochemistry* **32**, 8322–8328.
39. R. B. Spruijt, C. J. A. M. Wolfs, J. W. G. Verver, and M. A. Hemminga (1996) *Biochemistry* **35**, 10383–10391.
40. W. F. Wolkers, R. B. Spruijt, A. Kaan, R. N. H. Konings, and M. A. Hemminga (1997) *Biochem Biophys. Acta* **1327**, 5–16.
41. P. A. McDonnell, K. Shon, Y. Kim, and S. J. Opella (1993) *J. Mol. Biol.* **233**, 447–463.
42. N. A. Linderoth, P. Model, and M. Russel (1996) *J. Bacteriol.* **178**, 1962–1970.
43. M. Russel, N. A. Linderoth, and A. Sali (1997) *Gene* **192**, 23–32.

### Suggestions for Further Reading

44. P. Model and M. Russel (1988) In *The Bacteriophages*, Vol. **2** (R. Calendar, ed.), Plenum Press, New York, pp. 375–456. (Extensive review covering all topics mentioned here, with almost 400 references up to 1987.)
45. M. A. Hemminga, J. C. Sanders, C. J. A. M. Wolfs, and R. B. Spruijt (1993) In *New Comprehensive Biochemistry*, Vol. **25**; *Protein–Lipid Interactions* (A. Watts, ed.), Elsevier, Amsterdam, pp 191–212. (The role of the major coat protein in its *in vitro* membrane-bound state in molecular detail.)
46. M. Russel (1995) *Trends Microbiol.* **3**, 223–228. (A review about the reproductive life cycle focussing on the phage export mechanism.)
47. R. E. Webster (1996) In *Phage Display of Peptides and Proteins; A Laboratory Manual* (B. K. Kay, J. Winter, and J. McCafferty, eds.) Academic Press, New York, pp. 1–20. (A follow-up of the review of Model and Russel (1988), with more recent references. A clearly written, complete

overview describing what is presently known about the filamentous phage reproduction machinery and about what still remains to be elucidated.)

48. D. A. Marvin (1998) *Curr. Opin. Struct. Biol.* **8**, 150–158. (The most recent review about phage structure, infection, and assembly, elucidating much about the subunit organization of the phage coat.)

## Macroglobulins

Macroglobulins are large glycoproteins found in the circulation of vertebrates and invertebrates. They are also present in bird and reptile egg whites. When appropriately triggered, macroglobulins become binding proteins for proteinases belonging to all four mechanistic classes [see [Proteinase Inhibitors, Protein](#)]. In macroglobulin-proteinase complexes, the active site of the enzyme is generally left open and active. Some complexes are slightly more active, others slightly less active, toward small substrates than these enzymes. They can still be inhibited by small inhibitors. It is the interaction with large substrates that is prevented in complexes.

Complex formation starts by a specific cleavage of a peptide bond in the bait region of macroglobulin. However, in contrast to the single reactive site peptide bonds of protein serine proteinase inhibitors [see **Serine proteinase inhibitors**], one of many different peptide bonds can be broken to trigger. These accommodate the specificity of many different proteinases, so that macroglobulins are panspecific and “inhibit” most, but not all, proteinases. When elimination of action toward large substrates suffices, macroglobulins are the best choice to cut down on proteolysis by poorly understood or poorly characterized enzymes. After the bait region is triggered, the macroglobulin undergoes a conformational change that traps the enzyme molecule(s). In many macroglobulins, this noncovalent trapping is followed by formation of isopeptide bonds between the macroglobulin and the enzyme. The macroglobulins that undergo this reaction contain in them the sequence  $\frac{1}{4}\text{Cys Gly Glu Gln}\frac{1}{4}$  which leads to the formation of an internal thioester bond between the Cys and Gln residues (with the release of  $\text{NH}_3$ ). The thioester then reacts with Lys residues of the trapped enzyme to form isopeptide bonds.

The most studied of macroglobulins is human  $\alpha_2$  macroglobulin, which is a homotetramer of four chains, each of 1451 amino acid residues. It is a noncovalent homodimer of two disulfide bridged homodimers. It forms isopeptide bonds with trapped enzymes. Potentially, it has two enzyme binding sites, but whether one or two will be employed depends strongly on the kinetics of binding.

### Suggestions for Further Reading

- A. J. Barrett and P. M. Starkey (1973) The interaction of alpha 2-macroglobulin with proteinases. Characteristics and specificity of the reaction, and a hypothesis concerning its molecular mechanism. *Biochem. J.* **133**, 709–724.
- L. Sottrup-Jensen (1989) Alpha-Macroglobulins. Structure, shape and mechanism of proteinase complex formation. *J. Biol. Chem.* **264**, 11539–11542.

## Macromolecule

Macromolecules are molecules that are characterized by their very large size and high molecular weight. This diverse group of giant molecules can contain from several hundred to many thousands of atoms and have molecular weights ranging from ~1kDa to well over 1000 kDa. The term *macromolecule* is often used interchangeably with [polymer](#), though not all macromolecules are formed by polymerization of simple molecules. Macromolecules may be naturally occurring or man-made. Examples of man-made or synthetic macromolecules include plastics, fibers, and paints. Naturally occurring macromolecules include wool, cotton, wood, silk, and rubber. Of particular interest to the molecular biologist are the biological macromolecules or **biopolymers** essential for life, such as [proteins](#), **nucleic acids**, and carbohydrates.

[See also [Polymer](#) and [Biopolymer](#).]

### Suggestions for Further Reading

L. Mandelkern (1983) *An Introduction to Macromolecules*, Springer-Verlag, New York.

P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

## Macrophage

Macrophages, first known as cells very active in [phagocytosis](#), are involved in a large array of functions that confer upon them a central role in the host defense against pathogens. Observations on the importance of phagocytic cells were made as early as 1882 by Metchnikoff, who observed that macrophages could engulf and kill bacteria, thus providing a clear demonstration of the important role of cellular mechanisms in [immunity](#). He also observed that phagocytic activity was increased in **immunized** animals. This pointed to the importance of activation in macrophage physiology, a phenomenon that results from a constant interplay between macrophages and CD4<sup>+</sup> [T cells](#).

Macrophages originate from circulating monocytes, which represent a separate lineage derived from the bone-marrow [stem cells](#). They become resident in many tissues of the organism, where they function basically as scavengers for dead or **senescent** cells. In liver, where they have little contact with T cells, they are called *Kupffer cells* and are essentially devoted to these cleaning functions. When macrophages encounter microorganisms, they act quite differently. They first bind the pathogen, which implies the existence of receptors at the cell surface. Although only partially deciphered, the estimated number of different receptors expressed at the macrophage surface is probably rather large, maybe more than 100. This includes the mannose receptor, which may bind many microbial polysaccharides, but also molecules of the **immune system** such as **Fc receptors**, which bind **immunoglobulins** (Ig) (opsonization), or a variety of **cytokine** receptors, which account for possible stimulation by lymphocyte products. They express only low levels of [major histocompatibility complex](#) (MHC) class II molecules, as opposed to dendritic cells. Once bound, microorganisms are engulfed and degraded in [endosome](#) and lysosome compartments. The resulting peptides enter the cycle of [antigen processing and presentation](#) to T cells as they become expressed at the macrophage surface in association with MHC class I molecules. Interactions with T cells will involve primarily the CD4-T<sub>H</sub>1 population that will now, as an armed inflammatory T cell, produce

[interferon g](#) (IFNg) and [tumor necrosis factor a](#) (TNFa), thus ensuring amplification of the antibacterial activity of the macrophage. Fusion of phagosomes to lysosomes is increased, and substances with intense bacteriocidal activity like oxygen radicals and [nitric oxide](#) are produced. IFNg also stimulates expression of MHC class II molecules at the macrophage cell surface, resulting in a better presentation efficiency and, consequently, recruitment of more inflammatory T cells that contribute further to enhance the response. Macrophages therefore appear to be very efficient effectors to destroy intracellular pathogens, which may be bacteria, parasites, or **viruses**. If the destruction process is not complete or hindered because of a defect in the CD4 compartment (which is well known to occur in AIDS), macrophages become chronically infected and then provide a reservoir of pathogens.

Macrophages may also release soluble factors, including oxygen radicals, NO and [proteinases](#), that may also destroy extracellular pathogens. These substances are, however, also toxic for the host tissue. A down-regulation of macrophage activation must therefore be at hand so that this deleterious effect is limited. This is, in fact, provided by the CD4 T<sub>H</sub>2 T cells that will produce the cytokines **interleukin-4** (IL-4), IL-10, and IL-13. Finally, they may be targeted to specific pathogens by opsonization—that is, binding of their corresponding antibodies, which will bind their Fc receptor.

Although nonspecific as such, macrophages provide a remarkable example of efficiency of the immune system as a consequence of their balanced interaction between major helper T cell subsets.

See also entry **Antigen presentation**.

#### Suggestions for Further Reading

V. Witko-Sarsat and B. Descamps-Latscha (1994) Phagocyte-derived oxidants and proteinases as immunomodulatory mediators in inflammation. *Mediat. Inflamm.* **3**, 257–273.

E. Unanue and P. Allen (1987) The basis for immunoregulatory role of macrophages and other accessory cells. *Science* **236**, 551–557.

## Macropinosomes

Macropinosomes are large invaginations observed in cells that exhibit extensive plasma [membrane](#) ruffling, such as macrophages and dendritic cells (1). These cells engulf relatively large droplets of extracellular fluid at the site of ruffling. Macropinosome formation requires **actin**-containing **cytoskeletal** elements, and internalized macropinosomes either fuse with other organelles of [endocytosis](#) or are recycled back to the cell surface, depending on the cell type. Proteins internalized by macropinocytosis in dendritic cells can be delivered to **MHC** class II-containing vesicles, where they may be cleaved **proteolytically** to generate **peptides** that are recycled to the cell surface as MHC II-peptide complexes. **Antigens** may also be delivered to the cytosol for subsequent presentation on MHC class I antigens [see [Antigen Processing, Presentation](#) (2)]. Macropinocytosis is triggered in certain cells by treatment with phorbol esters, the ligation of [growth factor](#) receptors, or the microinjection of activated small [GTP-binding proteins](#), and it is inhibited by reagents that interfere with actin polymerization and by amiloride (1).

#### Bibliography

1. J. A. Swanson and C. Swanson (1995) *Trends Cell Biol.* **5**, 424–428.
2. C. C. Norbury, L. J. Hewlett, A. R. Prescott, N. Shastri, and C. Watts (1995) *Immunity* **3**, 783–

## MADS-Box Proteins

The [X-ray crystallography](#) structure of the complex of human **serum response factor** (SRF) with DNA revealed that the multiple functions of **DNA binding** and dimerization are integrated into a single-protein design (Fig. 1) (1). The DNA-binding region of SRF is called a MADS box, after the first four proteins identified to contain such a domain, namely MCM1 in yeast, the plant **homeotic proteins** *Agamous* and *Deficiens*, and serum response factor (SRF) in vertebrates (2). SRF binds to DNA as a homodimer, and the primary DNA-binding element is an antiparallel [coiled coil](#) of two amphipathic **a-helices** (helix-1), one from each subunit. The basic residues of these two a-helices fit into the major groove of the DNA double helix. The polypeptide chain continues from the *N*-terminal end of these a-helices to reach over the DNA backbone into the minor groove of the DNA. The rest of the dimerization interface of SRF is made up of a highly [hydrophobic](#) b-hairpin, the *C*-terminal a-helix, and the rather extended polypeptide structure connecting these two elements.

**Figure 1.** Global conformation of the SRF–DNA complex (1). The protein dimer binds from the minor groove side and bends the DNA into the minor groove. a-Helices are represented as cylinders, and b-strands are represented as arrows. Note the *N*-terminal extensions which lie along the minor groove of the DNA. The DNA sequence is depicted below.



```

5' - C C T T C C T A A T T A G G C C A T G
    G A A G G A T T A A T C C G G T A C C - 5'

```

Related SRF proteins, such as the members of the MEF-2 family of [transcription factors](#), show high sequence similarity to SRF within the regions of the *N*-terminal extension and helix-1 ([3](#), [4](#)), and the orientation of these proteins on the DNA and the details of the interactions are most probably very similar to those observed for SRF.

### Bibliography

1. L. Pellegrini, S. Tan, and T. J. Richmond (1995) *Nature* **376**, 490–498.
2. P. Shore and A. D. Sharrocks (1995) *Eur. J. Biochem.* **229**, 1–13.
3. D. Meierhans, M. Sieber, and R. K. Allemann (1997) *Nucleic Acids Res.* **25**, 4537–4544.
4. D. Meierhans and R. K. Allemann (1998) *J. Biol. Chem.* **273**, 26052–26060.

### Magnetic Force Microscopy

Magnetic force microscopy (MFM) is a [scanning probe technique](#) that is an offshoot of [atomic force microscopy](#) (AFM). Soon after the development of AFM, it was realized that AFM technology could be used to image magnetic forces. The first published MFM images were of a magnetic recording head. Since that first step, applications in the data storage industry have continued. With a resolution of better than 50 nm, AFM surpasses optical techniques and is much easier to use than electron-beam techniques with comparable resolution. Imaging is done under ambient conditions, it requires little or no sample preparation, and results are obtained in a few minutes.

As with AFM, MFM uses a flexible cantilever equipped with a sharp tip. To sensitize it to magnetic interactions, the cantilever is typically coated with a magnetic alloy. The most common imaging mode is referred to as “lift.” In this imaging mode, the sample topography is imaged along a line scan. The controlling computer “remembers” the topography and then directs the cantilever to repeat the scan along the topographical line, separated by a lift height. Because most of the short-range forces, such as [van der Waals interactions](#), diminish greatly at more than 2 to 5 nm, the remaining forces between the tip and sample are dominated by long-range electromagnetic interactions. In MFM, electrostatic forces are typically nulled, leaving only the magnetic contrast in the resulting image.

Application of MFM in the biological sciences has been limited to date, although it is possible to measure magnetic forces in a liquid environment. One application involves magnetotactic bacteria—Pyroketes that manufacture small angle crystals of ferromagnetic materials called magnetosomes and orient in the terrestrial magnetic field. The magnetosomes are small enough to support a single ferromagnetic domain. The bacteria produce chains of these particles that act as a single magnetic dipole large enough to overcome thermal randomization in the terrestrial magnetic field. The magnetosomes are produced in ambient conditions and have a narrow size distribution, so they also represent an impressive nanoengineering feat. Particles such as these may have a wide variety of applications, ranging from data storage to waste management. MFM was used to characterize the magnetic behavior of these particles in a recent study. Future applications of the MFM might include monitoring magnetic beads coated with antibodies and their interaction with various antigens.

A major application of the MFM is a new way of detecting magnetic resonance. In conventional techniques for measuring magnetic resonance, the electromagnetic signals induced in a coil or microwave cavity are detected by the collective precession of magnetic moments of the nuclei or electrons excited by an alternating magnetic field (see [NMR \(Nuclear Magnetic Resonance\)](#)). The

MFM detection technique works differently. When a paramagnetic sample is mounted on a micromechanical cantilever and placed in an inhomogeneous magnetic field, the sample is excited into magnetic resonance which produces a small oscillatory magnetic force ( $\sim 10^{14}$  to  $10^{15}$ N). The force is proportional to the volume of the sample, the magnetic field gradient, and the sample magnification. This force is detected by measuring the cantilever deflection in the optical deflection system of an AFM. The optical detection system senses the angstrom-scale vibration of a micromechanical cantilever on which the sample is mounted. Alternately, the mechanical response of the cantilever is larger if the force oscillates near the cantilever resonance frequency. Therefore, a static and sinusoidal oscillating radio-frequency polarizing magnetic field is applied to excite the magnetic spins in the sample. Then, the magnetic forces on the sample oscillate, and this oscillation of the sample on the cantilever is detected by a fiber optic interferometer. With an increase in the sensitivity of the optical detection system, it is possible to detect the magnetic resonance from single atoms. Hence it can be used for imaging microscopic samples in three dimensions. This technique of AFM-based NMR imaging (1) provides the most promising avenue of 3-D structural analysis of individual macromolecules without the need for crystallization and related complications.

#### Bibliography

1. D. Rugar, C. S. Yannoni, and J. A. Sidles (1992) Mechanical detection of magnetic resonance Nature **360**, 563–566.

#### Suggestions for Further Reading

2. P. Grutter (1994) “An Introduction to Magnetic Force Microscopy,” MSA Bulletin, **24**, 416–425.
3. R. Proksch and E. D. Dahlberg (1993) Optically stabilized, constant-height mode operation of magnetic force microscope, J Appl Phys. **73**, 5808–5810.
4. R. Proksch, T. E. Schaffer, B. M. Moskowitz, E. D. Dahlberg, and P. K. Hansma (1995) Magnetic force microscopy of the submicron magnetic assembly in magnetotactic bacteria, Appl. Phys. Lett. **66**, 2582–2584.

#### Magnetization Transfer

Magnetization is a general term used in various ways in the description of nuclear magnetic resonance (NMR) experiments. It is applied to detectable magnetic properties produced when a sample of material containing nuclei with nonzero nuclear spin is within a magnetic field.

When nuclei that have spin are placed in a magnetic field, they become distributed among their allowed energy states according to Boltzmann's law. The equilibrium distribution has different numbers of nuclei in the different levels; as a result, a net magnetization of the sample is produced. That is, all spins of the sample acting collectively produce a detectable macroscopic magnetic field that is aligned along the direction of the applied laboratory field. Magnetization aligned in this way is said to be longitudinal. Longitudinal magnetization can be converted, partly or wholly, into magnetization that is perpendicular (transverse) to the laboratory field by application of RF pulses to the sample, as is typically done in an NMR experiment. The resultant of the bulk (detectable) longitudinal and transverse magnetic field components is often referred to as the sample magnetization.

It is possible experimentally to transfer some or all of a population difference represented by the longitudinal part of a sample magnetization to energy levels other than those that were involved in



the initial creation of the population difference. A chemical reaction in which a nucleus (A) in an environment characterized by a shielding parameter  $s_A$  is transferred to another environment characterized by a different shielding parameter  $s_B$ , in the same or a different molecule, is one way this kind of magnetization transfer can take place. It is also possible to transfer transverse magnetization from one set of spins to another. Such transfers can be accomplished by a chemical reaction, or they can be the result of a process that is dependent on the spin coupling interaction between the nuclei. Magnetization transfer experiments are important ways of obtaining information about the rates of chemical reactions and have been used to measure the rates of dissociation of **enzyme-inhibitor complexes**, the rates of conformational change in proteins and nucleic acids (1), and the rates of enzymatic reactions *in vivo* (2).

Magnetization transfer by means of relaxation processes is involved in production of the [nuclear Overhauser effect](#) (NOE). (See also [NMR \(Nuclear Magnetic Resonance\)](#), [NOESY Spectrum](#), [COSY Spectrum](#).)

#### Bibliography

1. L. Y. Lian, I. L. Barsukov, M. J. Sutcliffe, K. H. Sze, and G. C. K. Roberts (1994) *Methods Enzymol.* **239**, 657–700.
2. J. R. Alger and R. G. Shulman (1984) *Quart. Rev. Biophys.* **17**, 83–124.

#### Suggestions for Further Reading

3. D. Canet (1996) *Nuclear Magnetic Resonance: Concepts and Methods*, Wiley, Chichester, Chapter "5", pp. 210–237.
4. *NMR in Drug Design*, (1996) (D. J. Craik, ed.) CRC Press, Boca Raton.

## Major Histocompatibility Complex

The major histocompatibility complex (MHC) is a **gene cluster** spanning several thousand kilobases on [chromosome 6](#) in humans and chromosome 17 in mice; similar complexes are found in most vertebrate species. The MHC was identified in the first half of the 1900s in experiments to determine the genetic region in mice responsible for the rejection or acceptance of transplanted tumors and skin. Transplant rejection subsequently was shown to be an immunological reaction (see also [Histocompatibility](#)); this and all adaptive [immune responses](#) are now known to be governed by [protein](#) molecules encoded by the highly polymorphic class I and class II loci of the MHC. While numerous genes coding for immune function-related and -unrelated proteins also reside within the MHC, the term “MHC molecules” conventionally is taken to refer to class I or class II gene products.

Vertebrate immune responses critically depend on the detection of pathogen-derived proteins. Class I and class II MHC molecules are cell-surface **glycoproteins** that bind short peptides derived from the limited **proteolytic degradation (antigen processing)** of both host and pathogen proteins. Cell-surface MHC/peptide complexes are potential ligands, or **epitopes**, for the [T-cell receptors](#) (TCRs) of circulating CD4+ and CD8+ T cells. This antigen-presenting function of MHC molecules provides the immune system with a sampling of the peptides derived from proteins residing in both the intra- and extracellular milieu of the host individual. Several mechanisms, including those responsible for T-cell maturation, typically ensure that T cells having TCR that recognize “self” MHC/peptide complexes (autoreactive T cells) are eliminated from the repertoire or otherwise

remain insensitive to self epitopes. Productive engagement of a TCR by an MHC/peptide complex results in functional activation of the responding T cell.

The **polymorphism** that exists at MHC class I and class II gene loci has no known parallel in the vertebrate [genome](#). The molecular diversity of MHC molecules is distinct from the diversity observed for [B-cell](#) receptors ( **antibodies**) and T-cell receptors, which arises through **somatic cell recombination**. MHC polymorphism resides in **germline** DNA and is observed to its fullest extent at the population level. Since the discovery of the MHC, the desire to understand the origins and functional consequences of MHC polymorphism has engaged immunologists, molecular biologists, and population geneticists alike. In recent years, extensive structural and biochemical analyses have provided molecular-level insights into MHC polymorphism and its functional concomitants. For both class I and class II molecules, structure and function are linked through the binding of peptides, and polymorphism is now understood in terms of its contribution to the peptide-binding specificities associated with various MHC allomorphs.

## 1. Class I MHC Molecules

Peptides bound by class I molecules are presented to the TCR of [cytotoxic T lymphocytes](#) (CTLs) expressing the accessory molecule CD8. Recognition of a class I/peptide complex by the TCR of a mature effector cell results in lysis (death) of the antigen-presenting cell. The pathway for peptide acquisition by class I molecules is optimized for the capture of peptides derived from proteins synthesized within the cell. The “class I-restricted” CTL response plays an important role in the detection and eradication of intracellular pathogens, such as **viruses**.

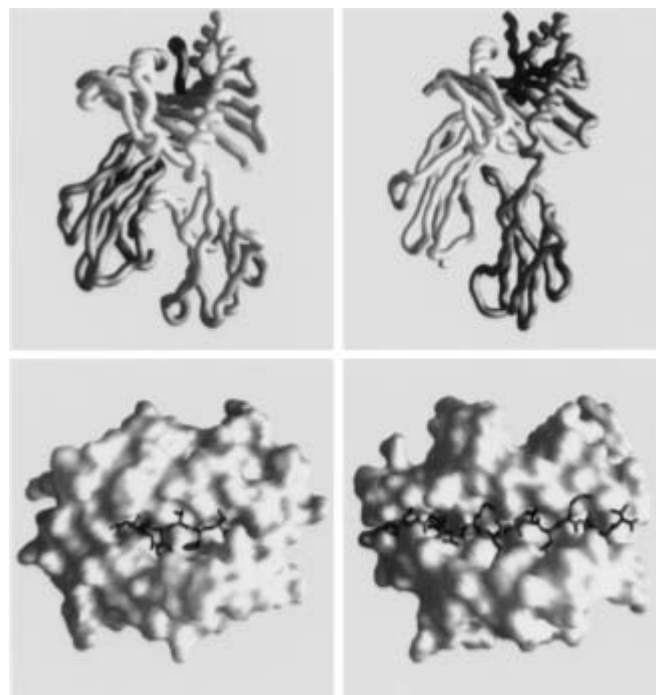
Class I molecules are heterodimers composed of a 45-kDa, glycosylated transmembrane protein that is noncovalently associated with  $\beta$ 2-microglobulin ( $\beta$ 2 m), a 12-kDa soluble protein. Both molecules are members of the **immunoglobulin superfamily**, but  $\beta$ 2 m is neither polymorphic nor encoded in the MHC; its gene resides on a different chromosome. Both  $\beta$ 2 m and the class I-bound peptide are integral components of the class I molecular structure, and the absence of either constituent results in loss of stable class I expression. Principle loci for class I molecules in the human MHC (also referred to as “HLA” for human leukocyte antigen) are designated A, B, and C. Over 50 class I **alleles** have been identified for the A and B loci combined, giving rise to cell-surface molecules that generally are serologically distinguishable. (Prior to the elucidation of class I function, these molecules were more commonly referred to as “antigens”). Nucleotide sequencing has revealed the existence of subtypes of many of these broad allelic classes. HLA alleles are formally specified, for example, as “HLA-A\*0205” or “HLA-B\*2703” to identify subtype “5” of the A2 allele or subtype “3” of the B27 allele, respectively. When subtypes are considered, the number of individual alleles existing at the A and B loci grows to at least 200, and new subtypes continue to be identified. A maternal and paternal allele of each locus (A, B, and C) is inherited, so that up to six potentially distinct class I molecules contribute to the “HLA type” of an individual.

The A, B, and C loci represent the so-called “classical” class I loci, which are characterized by a high degree of polymorphism, nearly ubiquitous tissue distribution, and documented roles in antigen presentation. Relative to A and B locus gene products, allomorphs of the C locus show more restricted polymorphism and more variable expression patterns. “Nonclassical” class I loci have also been described; in humans these include HLA-E, F, and G molecules of the MHC and the non-MHC-encoded, “class I-like” CD1 molecules. These gene products are characterized by limited polymorphism, and their expression tends to be either tissue-restricted or difficult to establish. Although functional roles have been proposed for HLA-G (1) and CD1 (2), nonclassical class I molecules remain poorly characterized relative to their classical counterparts.

The structural basis for peptide presentation by class I molecules has been made clear by [X-ray crystallography](#) (Fig. 1). Crystal structures for several human and murine class I allomorphs indicate that all class I protein molecules adopt essentially the same structure, regardless of type or origin (see [Protein Structure](#)). The membrane-distal  $\alpha$ 1 and  $\alpha$ 2 **domains** of the molecule fold to form a

peptide-binding groove; the “floor” of the groove is formed by **b- strands**, and the “walls” of the groove are composed of **a-helices**. The ends of the groove are closed, accounting for the restricted length (generally 8 to 10 residues) of peptides bound to class I molecules. Class I residues that are polymorphic are concentrated mainly in this peptide-binding region of the molecule, and they in turn dictate the antigenic peptide-binding specificity (or “motif”) of a given class I allomorph.

**Figure 1.** Structures of MHC class I and class II protein molecules with their bound peptides. Upper panels, left and right, are tracings of the main chain C<sub>α</sub> atoms of class I and class II molecules, respectively. The a chains of each molecule are shown in white; the class I light chain (b2-microglobulin) and class II b chain are shown in gray. Their bound peptides are the dark strands seen at the top of each molecule. Lower panels show views from above the peptide-binding grooves of (left) class I molecule HLA-A\*0201 complexed with a 9-residue peptide derived from matrix protein M1 of influenza A virus (3), and (right) the class II molecule HLA-DR\*0101 with a 13-residue peptide, also from influenza A virus (4). The MHC molecules are rendered as surfaces; the peptides are shown in an amino-to-carboxy, left-to-right orientation. Figures were produced with GRASP software (A. Nicholls and B. Honig, Columbia University) using Brookhaven Protein Data Bank coordinate files 1HHI (class I) and 1DLH (class II).



Biochemical analyses of peptides eluted from class I molecules have permitted the definition of the binding motifs associated with various allomorphs and have provided estimates of the heterogeneity and relative abundance of the peptides bound to cell-surface class I molecules. Consensus and individual sequences of naturally processed and presented peptides have revealed both allomorph-specific and general features of peptides associated with MHC molecules. Class I-bound peptides show pronounced biases in the amino acid types preferred at two to three key positions of the peptide (called anchor positions), while other positions exhibit various degrees of tolerance for diverse amino acids. For example, peptides associated with the allomorph HLA-A\*0201 are typically nonamers of the form “X-[Ile, **Leu, Met**, Val]-XXXXXX-[Ala, **Ile, Leu, Val]**” (where bold type indicates a strong preference and “X” denotes no clear preference); this indicates a preference for [nonpolar](#) side chains at peptide position 2 (P2) and at the carboxyl terminus (PC). In contrast, the peptide motif for HLA-B\*2705 is “[Arg, Lys]-[**Arg**]-XXXXXX-[Arg, **Leu, Lys, Tyr**].” Data also indicate that thousands of different peptides may be displayed on the cell surface; a given peptide species may be present at levels averaging from 100 to 1000 complexes per cell.

Many HLA class I molecules exhibit a binding motif characterized by anchor positions at P2 and PC, and structural data obtained for homogenous class I/peptide complexes (3) have revealed the basis for this observation. All peptides are bound to class I proteins with their amino and carboxyl termini oriented uniformly, such that their  $\alpha$ -amino group and carboxyl-terminal side chain are situated in conserved “pockets,” or subregions of the binding groove. Several class I molecules also feature a pocket that accommodates the P2 side chain of bound peptides. These conserved interactions contribute greatly to the energy of binding and account for the general P2/PC motif pattern observed for most class I molecules, while allomorph-specific features of these pockets determine the side-chain preferences at P2 and PC.

## 2. Class II Molecules

Peptides bound to class II molecules are presented to CD4<sup>+</sup> T cells, also known as “helper T cells.” Class II/peptide engagement by the TCR of these T cells triggers the release of **cytokines** that may have numerous effects on the immune response, including the promotion of **B-cell** function ([antibody](#) production) and [macrophage](#) activation. Class II molecules readily acquire peptides derived from internalized extracellular proteins, and the T-cell responses triggered by class II/peptide complexes are in turn keyed to the neutralization of extracellular pathogens, such as bacteria. In contrast to class I expression, class II molecule expression is generally limited to lymphoid cells, macrophages, dendritic cells, and endothelium, although class II expression is inducible by [g-interferon](#) in various other tissues.

Class II molecules are heterodimeric cell-surface proteins consisting of a 33-kDa  $\alpha$  chain associated noncovalently with a 28-kDa  $\beta$  chain. Prior to peptide acquisition in an **endosomal** compartment of the cell, the protein molecule is found in association with a third nonpolymorphic transmembrane protein termed the “invariant chain.” In humans, class II  $\alpha$  and  $\beta$  chains are encoded primarily by A ( $\alpha$  chain) and B ( $\beta$  chain) gene loci residing in three chromosomal regions designated DR, DQ, and DP. Gene products of all three regions function as antigen-presenting molecules. Class II genes are also highly polymorphic, especially those coding for  $\beta$  chains. Nucleotide sequencing has identified over 150 alleles of the DRB locus alone. The nomenclature for class II molecules is similar to that for class I: For example, “HLA-DRA\*0101” identifies an allele of the  $\alpha$  chain locus of the DR region.

Class II molecules and their associated peptides have also been characterized both structurally and biochemically. The  $\alpha$  and  $\beta$  chains form a peptide-binding groove highly similar to that formed by class I  $\alpha 1$  and  $\alpha 2$  domains. A conspicuous feature of the class II peptide groove is its open-ended design, which permits the accommodation of peptides longer than the nonamers typically found in the class I binding groove. Peptides bound to class II molecules may be from 12 to over 20 residues long. Class II molecules also feature well-defined pockets lined by polymorphic residues in their peptide-binding groove, and these pockets also dictate the types of peptide side chain that are permissible at a given anchor position of a class II-bound peptide. Crystallographic data (4) indicate also that several [hydrogen bonds](#) are formed between class II side chains and the peptide main-chain. While class I-peptide binding also involves conserved hydrogen bonds, these interactions appear to play a more critical role in stabilizing the class II-peptide complex.

Additional molecules encoded in the class II region include (a) the class II molecule HLA-DM and (b) subunits TAP1 and TAP2 of the TAP (transporter associated with antigen presentation) peptide transporter. HLA-DM exhibits limited polymorphism and has no known antigen-presenting function, but it plays a role in facilitating peptide acquisition by the cohort of conventional class II molecules (5). The TAP1/TAP2 heterodimer is a member of the “ATP binding cassette” family of **transporter** proteins, and is responsible for the translocation of cytosolic peptides into the [endoplasmic reticulum](#), where MHC class I-peptide assembly occurs (see [Antigen Processing, Presentation](#)).

## Bibliography

1. C. M. Schmidt and H. T. Orr (1995) *Immunol. Rev.* **147**, 53–65.
2. A. Melian, E. M. Beckman, S. A. Porcelli, and M. B. Brenner (1996) *Curr. Opin. Immunol.* **8**, 82–88.
3. D. R. Madden, D. N. Garboczi, and D. C. Wiley (1993) *Cell* **75**, 693–708.
4. L. J. Stern, J. H. Brown, T. S. Jardetzky, J. C. Gorga, R. G. Urban, J. L. Strominger, and D. C. Wiley (1994) *Nature* **368**, 215–221.
5. V. S. Sloan, P. Cameron, G. Porter, M. Gammon, M. Amaya, E. Mellins, and D. M. Zaller (1995) *Nature* **375**, 802–806.

### Suggestions for Further Reading

6. J. Klein (1986) *Natural History of the Major Histocompatibility Complex*, Wiley, New York. (Klein is a leading authority on the genetics and evolution of the MHC.)
7. V. H. Engelhard (1994) Structure of peptides associated with class I and class II MHC molecules, *Annu. Rev. Immunol.* **12**, 181–207.
8. D. N. Garboczi, P. Ghosh, U. Utz, Q. R. Fan, W. E. Biddison, and D. C. Wiley (1996) Structure of the complex between human T cell receptor, viral peptide, and HLA-A2, *Nature* **384**, 134–141. (The first high-resolution analysis of the ternary complex; a landmark.)

## Male

Male is one of the two possible [sexes](#) found in higher **eukaryotes**, the other being the [female](#). The characteristic feature of males is that, from puberty on, they produce [sperm](#), which are specialized haploid **germ cells** containing only those **organelles** that are necessary to transport genetic information and to **fertilize** a mature **oocyte**. Compared with the female, the investment to produce offspring is relatively low. After puberty, the male is able to produce sperm constantly, which contain no egg yolk, only minute amounts of cell cytoplasm, and virtually no RNA or other nutrients to promote growth of the new individual after union with the female [egg](#).

The spermatozoa are motile **gametes** specialized in the function of spreading genetic information widely throughout a population. Although there are exceptions to the rule, in higher **eukaryotes** and mammals the constant production of sperm cells and the prospect of totipotent fertilization capability cause competition for mating partners to take place primarily in the males. This is because, in many species, only some male individuals achieve reproductive success, whereas most of them suffer genetic death. This has led to genetic differentiation, and also affects the physiological behavior (1). Competition takes place not only at the level of sperm competing during the first stages of fertilization, but also during other stages of life. Evolutionary forces therefore had a great impact on sexual behavior, which can be regarded as preliminary selection to find the fittest. In order to impress the females and thus enhance the chance of copulation and prospect of reproductive success, all these factors have made male birds sing, for example, and has led to considerably different dimorphism regarding the behavior and **phenotype** in mammals.

To fulfill the task of genetic transfer, the most predominant feature of maleness is the special development of gonads and a genital tract. In chordates, male individuals share common morphological features. The organ responsible for the production of spermatozoa is the *testis* (2). Within this encapsulated organ, sertoli cells located in the ductuli support spermatogonia while they undergo several cycles of **mitosis** and then **meiosis**, until the final division makes them haploid.

During this process, germ cells move from the basal membranes to the lumen of the testis ductuli. The process of intracellular reorganization that ends with the creation of a spermatozoon is called *spermiogenesis*. Spermatozoa from the testicles are transported and then stored in the *epididymis*. In this organ, maturation of normal spermatozoa and resorption of abnormal spermatozoa occur. Maturation includes the association of binding proteins to the outer membrane of the [acrosome](#). Also, the physiological parameters of sperm change during this maturation process in the epididymis. Sperm from the epididymis head are not motile, but those that have passed the epididymis body and epididymis tail are.

Another characteristic of the male genital tract is the presence of accessory glands. The ampulla, seminal vesicles, prostate, and bulbo-urethral glands are activated on ejaculation and secrete a liquid that is, in terms of volume, the major part of the ejaculate. The secreted liquids contain buffers and sugars that serve to nourish and protect spermatozoa during their passage through the male, and subsequently the female, genital tract. The volume of the secretion of these glands can be great; for example, an ejaculate of the pig can be as much as 400 ml. The spermatid cords are united at the point where the urethra of the bladder joins in. Typically, the urethra ends in the penis.

Whether male or female genitals develop depends on the [hormone](#) status during early development of the [embryo](#). The basic plan of [development](#) in mammals is female. Without the induction of any sex hormone, the primordial cells develop into mullerian ducts and develop an ovum and uterus. The induction of sex hormones depends initially on the presence or absence of the *testis-determining factor*, which is the SRY gene found on the [Y-Chromosome](#) in mammals (3). Determination of the male sex does not, however, in all species depend on the presence of a certain [chromosome](#) or gene. In insects, for example, a number of genes must be homozygous to allow the development of male sex. In other species, sex determination depends on the relation of **autosomes** to sex chromosomes.

The interesting fact of the development of primordial germ cells is that sex determination occurs early through determinants of the female oocyte and they are not related to soma cells. In male mammals, primordial germ cells migrate during early embryonic stages of development from the vegetal pole (see [Animal Pole, Vegetal Pole](#)) to the cortex, which in the adult male is part of the testis. The primordial germ cells divide mitotically until birth. From then on, mitotic divisions of the so-called spermatogonia are stopped and they will be turned on again only during puberty. The induction of mitotic divisions and spermatogenesis is based on an increased level of testosterone. It is believed that **receptors** for testosterone lose their sensitivity to the hormone and, as a consequence, testosterone production is increased. In terms of physiology, the biggest difference with female gametogenesis is the constant production of germ cells after puberty and, in primates, the absence of a menopause.

### Bibliography

1. S. Ogawa, D. B. Lubahn, K. S. Korach, and D. W. Pfaff (1997) Proc. Natl. Acad. Sci. USA **94**, 1476–1481.
2. J. W. Foster et al. (1994) Nature **372**, 525–530.
3. C. M. Haqq, C. Y. King, E. Ukiyama, S. Falsafi, T. N. Haqq, P. K. Donahoe, and M. A. Weiss (1994) Science **266**, 1494–1500.

### **Maltose Binding Protein (MBP)**

The maltose binding protein (MBP) is a periplasmic element of the maltose transport machinery ([1-](#)

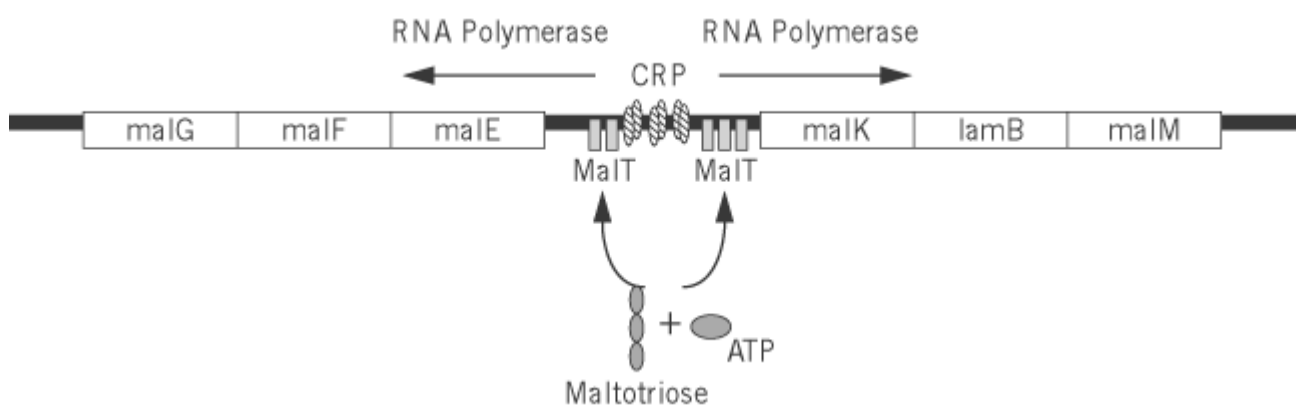
3) that belongs to the periplasmic permease family (4) and to the more ubiquitous group of ABC (ATP-binding cassette) transporters (5). MBP is also a chemoreceptor involved in sensory transduction and the **chemotactic** response toward maltose (6, 7). As a bifunctional protein, MBP is one of the most extensively studied proteins involved in transport and chemotaxis. The knowledge acquired by genetic, biochemical, and biophysical means provides important insights into the structure–function relationship of soluble globular periplasmic proteins. Its collaboration with integral membrane proteins involved in transport and chemotaxis makes MBP a choice target for studies on protein interactions, signal transmission, and elucidating the molecular mechanism of an ATP-driven transport of maltose and maltodextrins in *Escherichia coli*. It may serve as a model for a better understanding of ABC transporters, including the eukaryotic multidrug resistance (MDR) pump and the chloride channel altered in cystic fibrosis patients (CFTR) (8).

MBP is synthesized in the cytoplasm and must be exported to the periplasm (9). Its localization and folding in the periplasm are essential for its function. MBP was therefore used as a potent tool to select export mutants and to identify genes implicated in [protein secretion](#): the *sec* genes (10). MBP also serves as a favored model in studies on protein translocation and folding. Finally, the identification of “permissive” sites that could be modified or deleted without significantly perturbing MBP's functions (11) allowed the insertion in MBP of foreign epitopes for immunologic studies. MBP is also currently used, and is commercially available, as a vector protein for fusion proteins that allow expression in *E. coli* of foreign proteins and easy purification by [affinity chromatography](#) of the hybrid molecules (12).

## 1. Genetic Control, Expression, and Export

The *malE* gene encoding MBP (which is also called MalE) belongs to the *malEFG* [operon](#), which, together with the divergently transcribed *malKLM* operon, forms the *malB* region located at 91.5 min on the *E. coli* chromosomal map (13) (Fig. 1). Both operons are implicated in transport and are part of the maltose regulon, which includes at least nine genes positively controlled by the transcriptional activator MalT (14, 15). Operons of the *malB* region are also subject to dual positive control by **cyclic AMP-catabolite activator protein** (cAMP-CAP)—directly through the binding of the cAMP-CAP complex to the **promoter** region (16), and indirectly through MalT expression also being subject to CAP control (17, 18). Although MBP expression is inducible by maltose, the internal **inducer** appears to be maltotriose, which, together with ATP, binds preferentially to MalT and activates [transcription](#) (19, 20).

**Figure 1.** The *malB* region (91.5 min). Transcriptional regulation of the two divergent operons *malEFG* and *malKLM* belonging to the maltose regulon. The two operons are positively controlled by the transcriptional regulator MalT and the complex cyclic AMP-catabolite activator protein (cAMP-CAP). (See text for more detail.)



MBP, like other periplasmic proteins, is synthesized as a precursor (pre-MBP) in the cytosol before being exported to the periplasm (see [Protein Secretion](#)). [Translocation](#) is initiated when about 80% of the protein has been synthesized and requires that the pre-MBP is maintained in a “loosely” folded proteinase-sensitive state (21). This unfolded, translocation-competent conformation is maintained by the nascent pre-MBP associated with SecB, a tetrameric specialized **molecular chaperone** (see [Sec Mutants/Proteins](#)). This protein complex is then targeted to the secretion machinery in the membrane. Specific interactions between MBP and SecB seem to take place through the MBP **leader sequence** of 26 residues, which will be cleaved during translocation, and by binding motifs in mature MBP (22, 23). The initial binding of SecB to the substrate protein is probably triggered by [electrostatic interactions](#) (24, 25).

Following translocation of the polypeptide chain across the plasma membrane, mature MBP is released from the membrane and acquires its native, functional conformation. The folding process of MBP in the periplasmic space, as well as its dependence on some periplasmic chaperones, is not yet well understood (26).

In wild-type fully induced cells, the periplasmic concentration of MBP is about 1 mM, which corresponds to a 20- to 30-fold molar excess over the intrinsic membrane proteins involved in transport and chemotaxis (27). According to ultra structural observations, most of the MBP appears to converge at one pole of the cell (28).

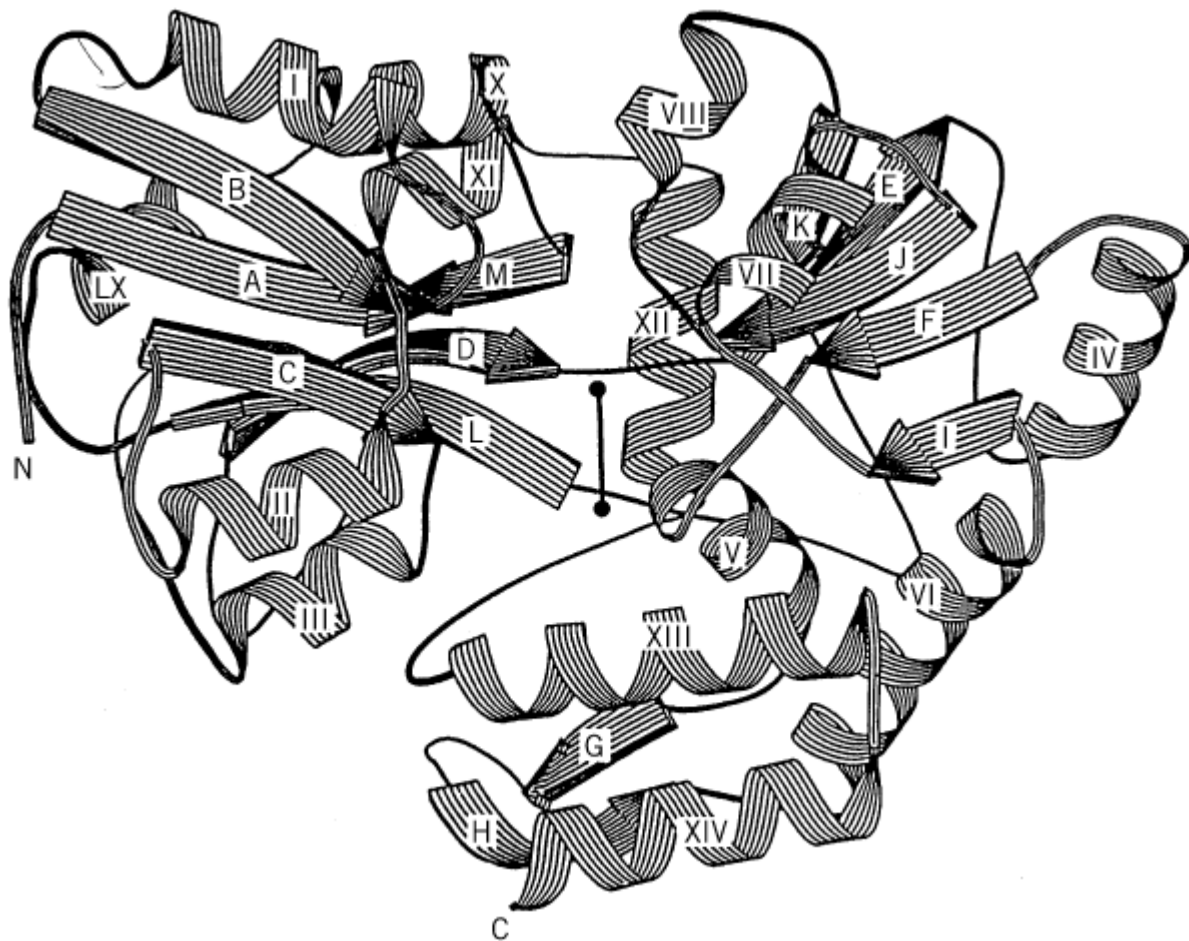
## 2. Structure and Function

Native MBP is a soluble globular molecule composed of 370 amino acid residues (29). Although it does not show significant sequence [homology](#) with other periplasmic binding proteins, it shares with them obvious structural similarities (30).

MBP is devoid of any known enzymatic activity; it shows a very stringent specificity and discriminates among several glucosides, binding selectively with high-affinity (**dissociation constant** in the micromolar range) maltose, higher  $\alpha$  (1  $\rightarrow$  4) glucose polymers, and [cyclodextrins](#) (up to maltoheptaose) (1, 31). The three-dimensional (3D) structure of MBP was determined by [X-ray crystallography](#) in the presence and absence of ligand (32, 33). In the presence of ligand, the protein appears as a bilobed, ellipsoid molecule with two distinct globular domains (designated N and C) connected by three peptide segments, forming a flexible hinge (Fig. 2). The N and C domains are each of an  $\alpha$  structural type and surround a deep cleft where maltose is sequestered and maintained by [hydrogen bonding](#) to the ligand hydroxyl groups and by **van der Waals contacts**, mainly with aromatic residues stacking their side-chains against the faces of the sugar rings. These extensive hydrogen and hydrophobic bonds provide the stereospecificity of ligand recognition and tend to stabilize the closed conformation of the protein. Dextrins longer than maltotriose have their nonreducing end protruding from the cleft.

**Figure 2.** The three-dimensional structure of maltose binding protein (MBP), which is (according to J. C. Spurlino et al. (32)) a bilobed  $\alpha$  structure. The N and C lobes surround the groove forming the maltose and maltodextrin binding site.





In the absence of ligand, MBP has an “open” conformation, where the edges of the cleft are far apart. Ligand binding induces a large conformational twisting between the N and C domains (a 35° rotation and 8° lateral twist), closing the binding site (33). This conformational modification seems essential for function and can be monitored by a characteristic change in the intrinsic **fluorescence** emission of MBP (31). One may assume that in solution and in the absence of substrate, MBP exists in a dynamic equilibrium between the open and closed form, the open form being favored energetically in the absence of ligand, the closed form when ligand is bound (33).

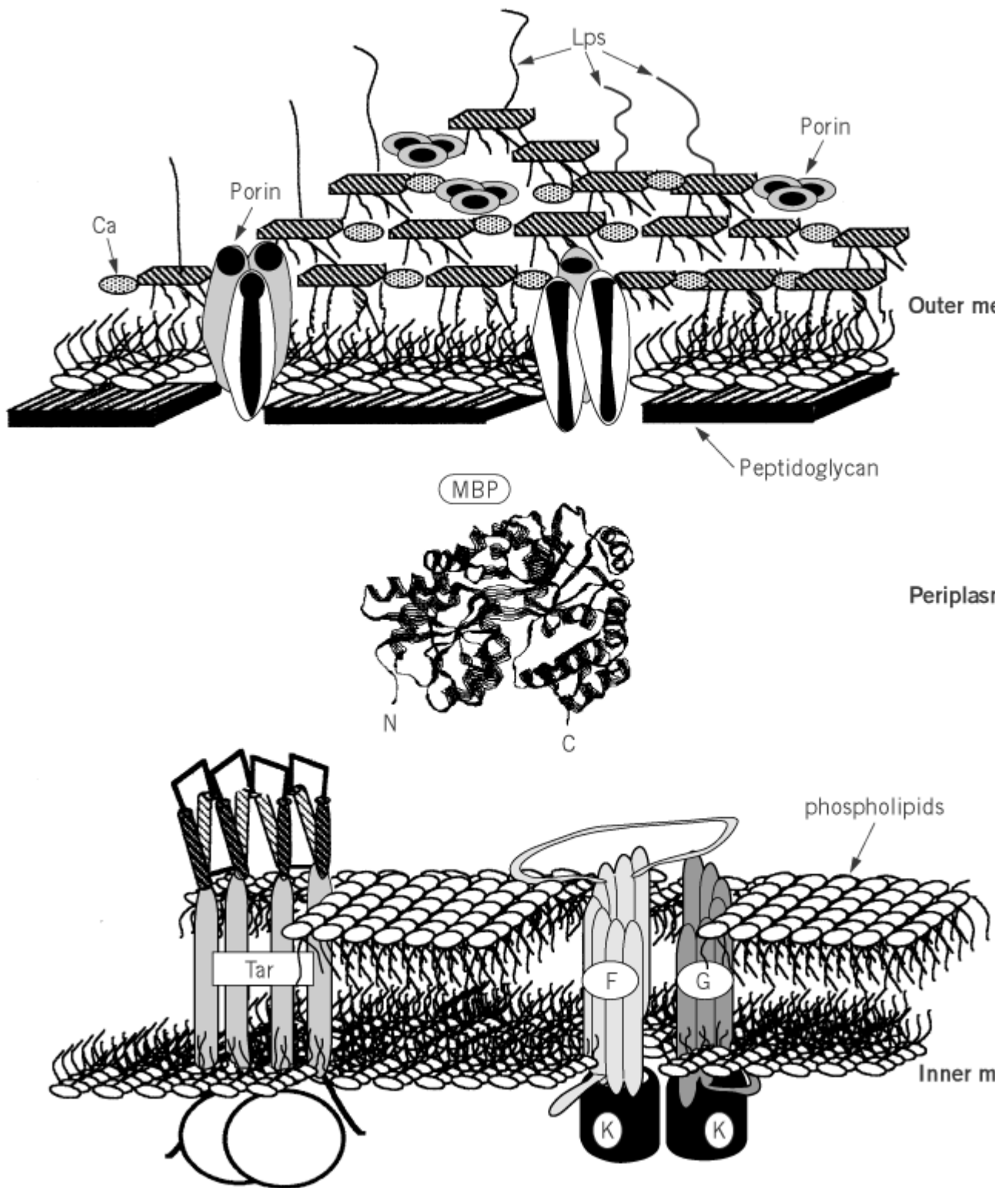
From **NMR** and UV **spectroscopic** studies of ligand binding to MBP (34, 35), and from the use of maltodextrin derivatives with modified reducing ends (36), it appears that MBP is able to bind dextrans by one of two modes: the “end on” mode, where the reducing end interacts with the binding site (as in the binding of maltose), and the “middle binding” mode, where binding to internal glucose moieties occurs along the dextrin chain. Only the conformation of MBP binding through the “end on” mode for maltose and maltodextrans seems to be productive and to trigger sensory transduction as well as the ligand translocation process (37).

### 2.1. Transport Function

Maltose and maltodextrans mobilize five proteins to cross the cell envelope of *E. coli*. The LamB porin facilitates specific ligand transfer through the outer membrane (38); the ABC transporter composed of MalF, MalG, and MalK proteins, in collaboration with MBP, ensures an energy-dependent uptake across the inner membrane (2, 39-41) (Fig. 3). Several types of deletions in the *malE* gene, as well as point mutations preventing MBP export, ligand binding, or interactions with the membrane transporter elements, result in a maltose-negative **phenotype** (1, 2, 11, 42). This strongly suggests that MBP is essential for transport and that this additional periplasmic element, acquired by the bacterial ABC transporter through evolution, must have conferred some important

advantage to the cells. When compared with transport systems that do not use periplasmic binding proteins, it appears that the presence of a periplasmic binding element increases the affinity of transport up to 100-fold. The  $K_d$  values measured *in vitro* (in the micromolar range) and the substrate concentration of half-maximal rate of transport are nearly equivalent (31, 43). This permease system, with a high rate of substrate uptake coupled to high affinity, enables the cells to grow and use maltose and maltodextrins efficiently, even when they are available at very low concentrations in the medium (the substrate can be transported against a concentration gradient of the order of  $1:10^5$ ). Interactions between MBP and the LamB porin were observed *in vitro* (44, 45), and some genetic data suggested that these associations may facilitate diffusion of dextrans through the outer membrane (46). The physiologic role of these interactions, however, remains controversial (47).

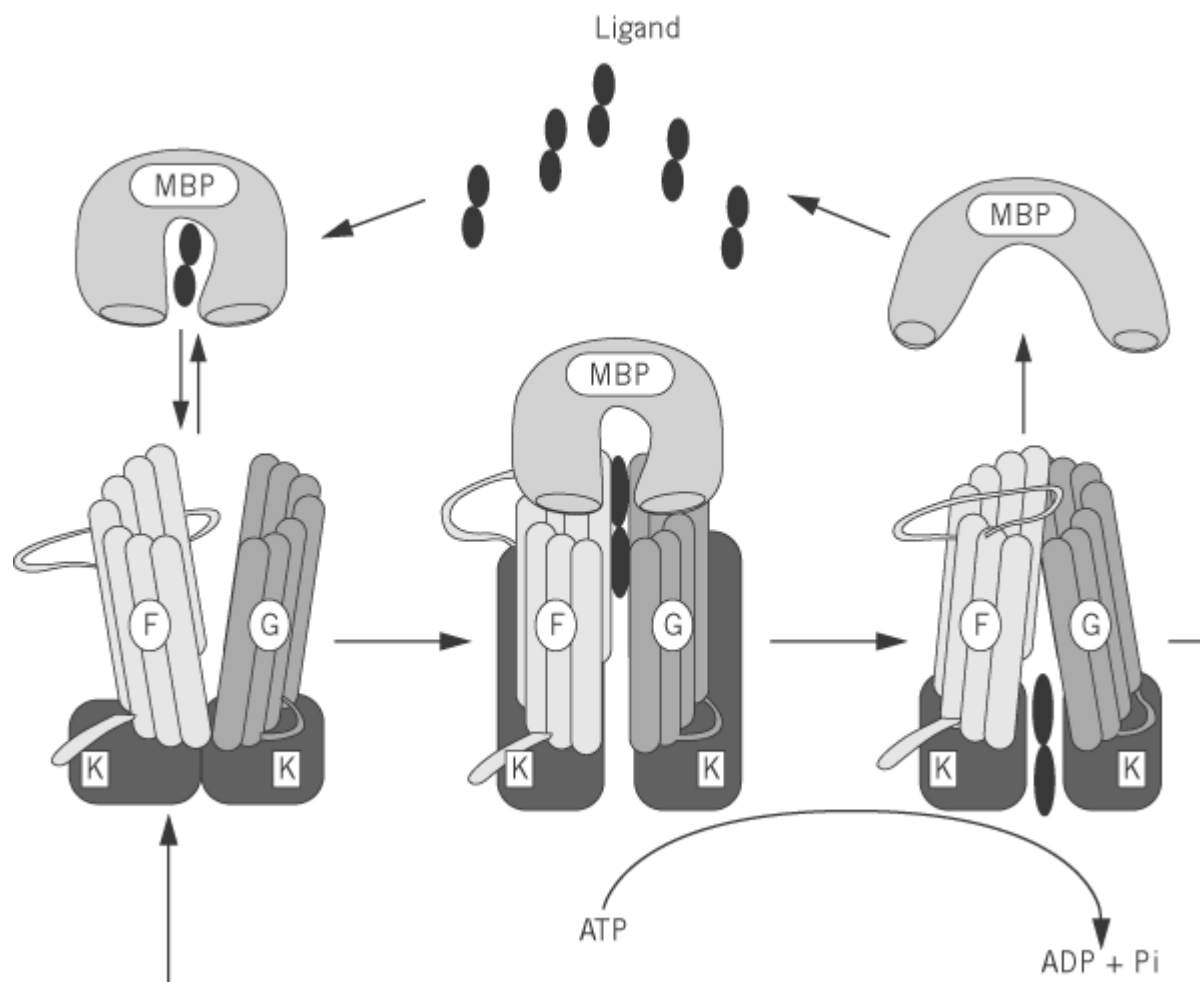
**Figure 3.** Schematic representation of the periplasmic space, the outer and inner membranes, and some constituents of the transport system (the LamB porin, MBP, and the MalFGK<sub>2</sub> membrane complex), and Tar, the chemotransducer involved in transduction. (Drawing is not to scale.) Ca, calcium.



Free substrate molecules crossing the outer membrane and reaching the periplasmic space are picked up and bound tightly by MBP. The high concentration of MBP in the periplasmic space (about 1 mM), plus a quite slow, unassisted release of the ligand ( $k_{off} = 1100s^{-1}$ ) (48), will favor retention of the ligand by the binding protein and probably prevent free ligand circulation in the periplasm (49). From this and the above genetic observations, it appears that delivery of substrate to the ABC transporter will be ensured by MBP. Translocation across the inner membrane therefore depends, on

productive interactions between the liganded MBP and the membrane complex (MalFGK<sub>2</sub>), which is composed of one copy of each of the membrane proteins MalF and MalG and of two molecules of the MalK protein, which binds and hydrolyzes ATP (40, 50). Only MBP in the presence of maltose is able to trigger ATP hydrolysis (51). It seems likely that the MBP closed conformation—essentially its liganded state—is able to activate MalK and concomitantly induce a cascade of events, leading to the release and transfer of ligands across the inner membrane (Fig. 4). It is noteworthy that a mutant MBP, with two cysteine residues introduced into each of the N and C lobes of the protein and **cross-linked** by an interdomain disulfide bond, confers a dominant-negative phenotype for growth on maltose, for maltose transport, and for chemotaxis (52). This strengthens the hypothesis that it is the closed form of MBP that interacts with the membrane transporter elements. According to a mathematical treatment of maltose transport kinetics, however, both the loaded and unloaded substrate forms (open and closed) of MBP are supposed to interact with the membrane transporter proteins (53).

**Figure 4.** The transport model for MBP. The liganded form of MBP interacts with the membrane-bound ABC transporter complex MalFGK<sub>2</sub>. This interaction triggers reversible conformational modifications of all the partners, with concomitant hydrolysis of ATP, followed by ligand translocation across the inner membrane.



Indirect but important information about the interactions between MBP and the other membrane components was provided by genetic analysis. From a transport-negative *malE* deletion mutant, revertants that did not require MBP for maltose uptake (39) were found to have double mutations in either *malF* or *malG* and to transport maltose with a  $K_m$  of 1 mM, instead of the 1  $\mu$ M for wild-type. On addition of wild-type MBP, some of the mutants became transport-negative, presumably

due to a nonproductive interaction with the MalFGK<sub>2</sub> complex. Allele-specific suppressor mutations in MBP that restored transport, as well as random mutations in MBP conferring a negative-dominant phenotype over the wild-type MBP, were isolated and mapped (54). Finally, by random insertion mutagenesis of the *malE* gene, a domain located at the surface of MBP and involved in interactions with the membrane transporter could be pinpointed (11). Deletion of this domain, which is primarily an [helix](#) (helix7), specifically inhibits transport but retains wild-type affinity and chemotactic response toward maltose. Individual residues implicated in the interaction were found to be located near the edge of the binding site in the N and C lobes (42). The present model suggests that the interactions take place between the N lobe and MalG, whereas the C lobe of MBP interacts with MalF.

## 2.2. Chemotactic Function

MBP serves as a maltose chemoreceptor for *E. coli* that enables the cell to recognize maltose and to direct its motility in a spatiotemporal gradient toward the sugar (6). This chemoreceptor function seems to have been acquired after the transport function, as the *malE* gene is coregulated with the maltose regulon and belongs to the transport specialized *malB* region. Whether MBP is free or liganded is specifically detected by the inner-membrane chemotransducer Tar, which acts as a homo-dimer (55). The Tar protein is anchored in the cytoplasmic membrane by two transmembrane helices (TM1 and TM2); TM2 connects a periplasmic and a cytoplasmic domain (Fig. 4). Tar also mediates chemoattraction to aspartate, which binds directly to the chemotransducer (56, 57). *In vitro* experiments with partially reconstituted elements have demonstrated physical contacts between MBP and Tar and shown that only the maltose-loaded form of MBP interacts with Tar (58). Mutants in *malE* that caused specific defects in the chemotactic function of MBP defined two relatively confined regions around Thr53 in the N lobe and Thr345 in the C lobe, which probably comprise a conformational interaction site with Tar (59). The distance between the two regions is supposed to change as the protein shifts from the open to the closed state on ligand binding. The regions assumed to interact with Tar appear to be distinct from those involved in productive contacts with the MalFGK<sub>2</sub> complex. Both are located on the same face of the protein. Residues in the periplasmic domain of Tar involved in liganded MBP recognition were found to be located in two antiparallel loops joining four  $\alpha$ -helices and forming the 3D structure of the periplasmic domain of Tar (60, 61). Some mutations interfere with both aspartate and maltose taxis. These genetic observations lead to the conclusion that adjacent, partially overlapping sites in the periplasmic domain of Tar interact with aspartate and with MBP loaded with maltose.

## 3. MBP as a Tool: Applications Derived from Basic Knowledge

MBP is widely used in molecular biology as a vector protein (12), because mutagenesis of the *malE* gene identified functional regions in the protein, as well as “permissive” sites where deletions or insertions do not cause deleterious effects (11).

### 3.1. Studies on Immunogenicity and Antigenicity of Peptides

Eleven permissive sites have been identified in MBP. These sites were used to test the influence of an **epitope's** surroundings on its immunologic properties (**immunogenicity** and/or **antigenicity**). For this purpose, a given peptide is inserted into different positions of MBP, and the corresponding recombinant proteins are delivered to the immune system. Insertions within the continuity of a protein limit the mobility of the inserted peptide, and the effects of such a constraint are amenable to antigenic analysis with polyclonal and [monoclonal antibodies](#). The immunogenic properties of an MBP hybrid protein may be analyzed in its purified form or within the periplasmic space of live bacteria. This type of approach, coupled to structural analysis, should help to correlate immunogenicity, antigenicity, and physical properties of the inserted peptide and to provide valuable approaches to recombinant vaccine strategy (62). Foreign epitopes, such as the C3 epitope of the [polio virus](#) and epitopes from **HIV1** coat protein or from the preS2 region of the [hepatitis B virus](#), have been subjected to this type of study (63, 64).

### 3.2. MBP as an Expression and Targeting Vector for Foreign Proteins in *E. coli*

Foreign peptides or whole proteins fused to either end of MBP (N or C fusions) can be targeted to the periplasmic space, where the disulfide bonds necessary for correct folding are formed. When necessary, hybrid proteins may be released from the cell by osmotic shock treatment and easily purified by [affinity chromatography](#) on crossed-linked amylose columns. In such fusions, both the MBP moiety and the passenger protein generally remain functional. For example, the soluble part of CD4 (V1-V2) was fused to MBP and exported to the periplasm. After affinity purification on amylose, the CD4 part bound to HIV gp120 and neutralized HIV particles *in vitro* (65).

### 4. Acknowledgments

I am greatly indebted to M. Hofnung for his constructive comments and his critical reading of the manuscript. I want to thank Ana Cova for helping to make up the final version of the manuscript.

### Bibliography

1. O. Kellerman and S. Szmelcman (1974) *Eur. J. Biochem.* **47**, 139–147.
2. H. A. Shuman (1982) *J. Biol. Chem.* **257**, 5455–5461.
3. A. L. Davidson, H. A. Shuman, and H. Nikaido (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2360–2364.
4. G. F. Ames (1986) *Annu. Rev. Biochem.* **55**, 397–425.
5. C. F. Higgins (1992) *Annu. Rev. Cell Biol.* **8**, 67–113.
6. G. L. Hazelbauer (1975) *J. Bacteriol.* **122**, 206–214.
7. M. D. Manson and M. Kossmann (1986) *J. Bacteriol.* **165**, 34–40.
8. G. F. Ames (1992) *Int. Rev. Cytol.* **137**, 1–35.
9. C. K. Murphy and J. Beckwith (1996) In *Escherichia coli and Salmonella typhimurium: cellular and molecular biology* (F. C. Neidhardt et al., eds.), ASM Press, Washington, D.C. pp. 967–975.
10. P. J. Schatz and J. Beckwith (1990) *Ann. Rev. Genet.* **24**, 215–248.
11. P. Duplay, S. Szmelcman, H. Bedouelle, and M. Hofnung (1987) *J. Mol. Biol.* **194**, 663–673.
12. M. Hofnung, A. Charbit, J. M. Clément, C. Leclerc, P. Martineau, S. Muir, D. O'Callaghan, O. Popescu, and S. Szmelcman (1992) In *Targeting of Drugs 3: the challenge of peptides and proteins* (G. Gregoriadis et al., eds.), Plenum Press, New York, pp. 109–119.
13. M. Hofnung, D. Hatfield, and M. Schwartz (1974) *J. Bacteriol.* **117**, 40–47.
14. D. Hatfield, M. Hofnung, and M. Schwartz (1969) *J. Bacteriol.* **98**, 559–567.
15. O. Raibaud, M. Débarbouillé, and M. Schwartz (1983) *J. Mol. Biol.* **163**, 395–408.
16. D. Vidal-Ingigliardi and O. Raibaud (1991) *Proc. Natl. Acad. Sci. USA* **88**, 229–233.
17. M. Débarbouillé and M. Schwartz (1979) *J. Mol. Biol.* **132**, 521–534.
18. C. Chapon and A. Kolb (1983) *J. Bacteriol.* **156**, 1135–1143.
19. O. Raibaud and E. Richet (1987) *J. Bacteriol.* **169**, 3059–3061.
20. E. Richet and O. Raibaud (1989) *EMBO J.* **8**, 981–987.
21. L. L. Randall and S. J. S. Hardy (1986) *Cell* **46**, 921–928.
22. J. R. Thom and L. L. Randall (1988) *J. Bacteriol.* **170**, 5654–5661.
23. P. M. Gannon, P. Li, and C. A. Kumamoto (1989) *J. Bacteriol.* **171**, 813–818.
24. V. A. Bankaitis, B. A. Rasmussen, and P. J. J. Bassford (1984) *Cell* **37**, 243–252.
25. L. L. Randall and S. J. S. Hardy (1995) *Trends Biochem. Sci.* **20**, 65–70.
26. J. M. Betton, D. Boscus, D. Missiakis, S. Raina, and M. Hofnung (1996) *J. Mol. Biol.* **262**, 140–150.
27. M. D. Manson, W. Boos, P. J. Bassford, and B. A. Rasmussen (1985) *J. Biol. Chem.* **260**,

9727–9733.

28. J. R. Maddock and L. Shapira (1993) *Science* **259**, 1717–1723.
29. P. Duplay, H. Bedouelle, A. Fowler, I. Zabin, W. Saurin, and M. Hofnung (1984) *J. Biol. Chem* **259**, 10606–10613.
30. F. A. Quicho and P. S. Ladvina (1996) *Mol. Microbiol.* **20**, 17–25.
31. S. Szmelcman, M. Schwartz, T. J. Silhavy, and W. Boos (1976) *Eur. J. Biochem.* **65**, 13–19.
32. J. C. Spurlino, G. Lu, and F. A. Quicho (1991) *J. Biol. Chem.* **266**, 5202–5219.
33. A. J. Sharff, L. E. Rodseth, J. C. Spurlino, and F. A. Quicho (1992) *Biochemistry* **31**, 10657–10663.
34. K. Gehring, P. G. Williams, J. G. Pelton, H. Morimoto, and D. E. Wemmer (1991) *Biochemistry* **30**, 5524–5531.
35. K. Gehring, K. Bao, and H. Nikaido (1992) *FEBS Lett.* **300**, 33–38.
36. T. Ferenci, M. Muir, K. -S. Lee, and D. Maris (1986) *Biochim. Biophys. Acta* **860**, 44–50.
37. H. Nikaido (1994) *FEBS Lett.* **346**, 55–58.
38. S. Szmelcman and M. Hofnung (1975) *J. Bacteriol.* **124**, 112–118.
39. N. A. Treptow and H. A. Shuman (1985) *J. Bacteriol.* **163**, 654–660.
40. A. L. Davidson and H. Nikaido (1991) *J. Biol. Chem.* **266**, 8946–8951.
41. E. Dassa and M. Hofnung (1985) *EMBO J.* **4**, 2287–2293.
42. S. Szmelcman, N. Sassoon, and M. Hofnung (1997) *Prot. Sci.* **6**, 1–9.
43. R. Hengge and W. Boos (1983) *Biochim. Biophys. Acta* **737**, 443–478.
44. P. Bavoil and H. Nikaido (1981) *J. Biol. Chem.* **256**, 11385–11388.
45. K. A. Stauffer, A. Hoenger, and A. Engel (1992) *J. Mol. Biol.* **223**, 1155–1165.
46. C. Wandersman, M. Schwartz, and T. Ferenci (1979) *J. Bacteriol* **140**, 1–13.
47. S. Freundlieb, U. Ehmann, and W. Boos (1988) *J. Biol. Chem.* **263**, 314–320.
48. D. M. Miller, J. S. Olson, J. W. Pflugrath, and F. L. Quicho (1983) *J. Biol. Chem.* **258**, 13665–13672.
49. T. J. Silhavy, S. Szmelcman, W. Boos, and M. Schwartz (1975) *Proc. Natl. Acad. Sci. USA* **72**, 2120–2124.
50. C. Walter, K. Höner zu Bentrup, and E. Schneider (1992) *J. Biol. Chem.* **267**, 8863–8869.
51. D. A. Dean, A. L. Davidson, and H. Nikaido (1989) *Proc. Natl. Acad. Sci. USA* **86**, 9134–9138.
52. Y. Zhang, D. E. Mannering, A. L. Davidson, N. Yao, and M. D. Manson (1996) *J. Biol. Chem.* **271**, 17881–17889.
53. E. Bohl, H. A. Shuman, and W. Boos (1995) *J. Theor. Biol.* **172**, 83–94.
54. L. -I. Hor and H. A. Shuman (1993) *J. Mol. Biol.* **233**, 659–670.
55. J. Adler, G. L. Hazelbauer, and M. M. Dahl (1973) *J. Bacteriol.* **115**, 824–847.
56. A. M. Stock and S. L. Mowbray (1995) *Curr. Opin. Struc. Biol.* **5**, 744–751.
57. D. F. Blair (1995) *Ann. Rev. Microbiol.* **49**, 489–522.
58. G. Richarme (1982) *J. Bacteriol.* **149**, 662–667.
59. Y. Zhang, C. Conway, M. Rosato, Y. Suh, and M. Manson (1992) *J. Biol. Chem.* **267**, 22813–22820.
60. P. Gardina, C. Conway, M. Kossman, and M. Manson (1992) *J. Bacteriol.* **174**, 1528–1536.
61. J. I. Yeh, H. -P. Bieman, J. Pandit, D. E. Koshland, and S. H. Kim (1993) *J. Biol. Chem.* **268**, 9787–9792.
62. A. Charbit, S. M. Newton, P. E. Klebba, J. -M. Clément, C. Fayolle, R. Loman, C. Leclerc, and M. Hofnung (1996) In *2nd International Meeting on New Approaches to Bacterial*

*Vaccine Development*. Behring Institute Mitteilungen **98**, 1–8.

63. A. Lecroisey, P. Martineau, M. Hofnung, and M. Delepierre (1987) *J. Biol. Chem.* **272**, 362–368.
64. B. Vulliez-Le Normand, P. Martineau, M. Hofnung, F. Lema, and G. A. Bentley (1997) *Prot. Engineering* **10**, 175–180.
65. S. Szmelcman, J. -M. Clément, M. Jehanno, O. Schwartz, L. Montagnier, and M. Hofnung (1990) *J. Acquir. Imm. Defic. Synd.* **3**, 859–872.

### Suggestions for Further Reading

66. H. A. Shuman and C. H. Panagiotidis (1993). Tinkering with transporters: periplasmic binding protein-dependent maltose transport in *E. coli*. *J. Bioeng. Biomemb.* **25**, 613–620.
67. E. Dassa, E. Francoz, M. Dahl, E. Schneider, C. Werts, A. Charbit, W. Saurin, and M. Hofnung (1993) "The maltose B region in *S. typhimurium*, *E. coli* and other Enterobacteriaceae". In *The Biology of Salmonella*. NATO ASI Series (F. Cabello, ed.), Plenum Press, New York, pp. 91–104.
68. W. Boos and J. M. Lucht (1996) "Periplasmic binding protein-dependent ABC transporters". In *Escherichia coli and Salmonella typhimurium: cellular and molecular biology* (F. C. Neidhart et al., eds.), ASM Press, Washington D. C., pp. 1175–1199.
69. J. J. Falke, D. F. Blair, T. J. Silhavy, and R. Schmitt (1995) BLAST 1995: international conference on Bacterial Locomotion And Signal Transduction. *Mol. Microbiol.* **16**, 1037–1050.

## Mannose 6-P Receptors

Newly synthesized soluble acid hydrolases in higher eukaryotic cells contain mannose 6-phosphate (M6P) residues on *N*-linked oligosaccharides (1). This Golgi modification allows them to be sorted from secretory proteins and delivered to the lysosomes. The M6P recognition site is formed by a two-step process involving (i) the UDP-*N*-acetylglucosamine:lysosomal enzyme *N*-acetylglucosamine-1-phosphotransferase, which forms a phosphodiester intermediate, and (ii) a-*N*-acetylglucosaminyl phosphodiester glycosidase, which generates the M6P monoester (2-4). Genetic defects in the first enzyme involved in this process cause the autosomal recessive, lysosomal storage disorders, I-cell disease (mucopolipidosis II) and pseudo-Hurler polydystrophy (ML III) (3, 5, 6). Patients with these conditions present with increased levels of lysosomal enzymes in their plasma because of a failure to generate phosphomannosyl recognition marks.

The intracellular trafficking of M6P-tagged lysosomal enzymes to the lysosomes is mediated by the mannose 6-phosphate receptors (M6PRs). Upon their delivery to the prelysosomal compartment, the lysosomal enzymes dissociate from the M6PRs, and the phosphate recognition marks are removed when the lysosomal enzymes reach the lysosome (7). The M6PRs then return to the Golgi to initiate another lysosomal enzyme delivery cycle. Approximately 10% of the M6PRs are also present on the plasma membrane where they endocytose extracellular ligands. Thus, the M6PRs function primarily as shuttle crafts for the intracellular sorting and trafficking of lysosomal enzymes. They are present in the Golgi apparatus, the *trans*-Golgi network, the prelysosomal compartment, and the plasma membrane. The highest steady-state receptor concentration, however, is in the prelysosomal compartment because the endocytic and biosynthetic routes converge at this location (8).



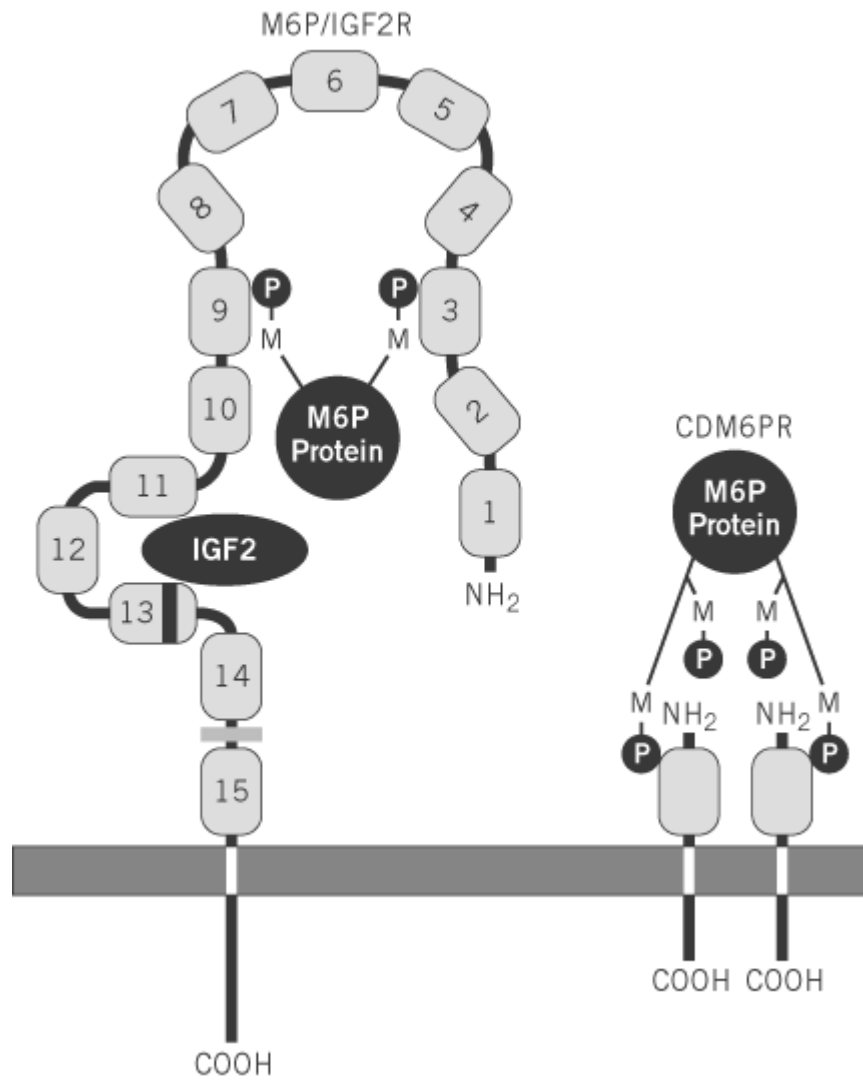
There are two M6PRs involved in the intracellular trafficking of lysosomal enzymes. The first receptor to be identified was the 275-kDa cation-independent M6PR (CIM6PR). It is also called the insulin-like growth factor 2 receptor (IGF2R) because of its additional ability to bind IGF2. In this article the CIM6PR will be referred to as M6P/IGF2R to recognize its dual function of both trafficking lysosomal enzymes and regulating extracellular IGF2 bioavailability. The 46-kDa cation-dependent M6PR (CDM6PR) was subsequently purified from bovine liver because optimal ligand binding depended upon the presence of divalent cations (9). Because these two receptors are not similar in their sequence to other lectins and appear to possess a unique class of sugar-binding domains, they have been classified as P-type lectins (10). The molecular, functional, and biological characteristics of these two unique receptors are described below.

## 1. Cation-Dependent Mannose 6-Phosphate Receptor

### 1.1. CDM6PR: Molecular Characteristics

The CDM6PR is a 46-kDa transmembrane protein that functions primarily as a homodimer (Fig. 1) (11, 12). The complete cDNA sequence for the *CDM6PR* has been determined for the bovine (13), human (14), and murine (15, 16) species (Table 1). The overall predicted structure of the CDM6PR molecule is identical in these species. It is comprised of an extracellular domain containing one M6P binding site, a single transmembrane domain, and a small cytoplasmic tail.

**Figure 1.** Structural models of the two mannose 6-phosphate receptors. The human cation-dependent mannose 6-phosphate receptor (CDM6PR) is comprised of a 159-amino-acid NH<sub>2</sub>-terminal extracellular domain, a 25-amino-acid transmembrane region (white bar), and a 67-amino-acid COOH-terminal cytoplasmic domain (14). Two different oligosaccharides of a M6P glycoprotein (M6P protein) are bound to a CDM6PR homodimer. The human mannose 6-phosphate/insulin-like growth factor 2 receptor (M6P/IGF2R) is comprised of a 2264-amino-acid NH<sub>2</sub>-terminal extracellular domain, a 23-amino-acid transmembrane region (white bar), and a 164-amino-acid COOH-terminal cytoplasmic domain (33, 42). The extracellular portion of M6P/IGF2R consists of 15 contiguous repeat regions, with a short gap (gray bar) between repeats 14 and 15. A diphosphorylated oligosaccharide of a M6P glycoprotein (M6P protein) is bound to a single M6P/IGF2R molecule at the M6P binding sites in repeats 3 and 9 (55). IGF2 is bound to its binding site in repeat 11 (36, 37, 58). The affinity of IGF2 binding is enhanced by the presence of a fibronectin, type II-like element (black bar) in repeat 13 (59). The cytoplasmic domain of both receptors contains the internalization and trafficking signals (28-30, 61, 62).



**Table 1. *CDM6PR* Species Comparisons**

| Species                   | cDNA (bp) | 5'UTR <sup>a</sup> (bp) | 3'UTR <sup>a</sup> (bp) | ORF <sup>b</sup> (bp) | Protein <sup>c</sup> (aa) | Signal Peptide (aa) | $M_r$ <sup>d</sup> |
|---------------------------|-----------|-------------------------|-------------------------|-----------------------|---------------------------|---------------------|--------------------|
| <b>Bovine</b><br>(13)     | 2243      | 37                      | 1369                    | 837                   | 279                       | 28                  | 27,963             |
| <b>Human</b><br>(14)      | 2463      | 145                     | 1487                    | 831                   | 277                       | 26                  | 27,913             |
| <b>Murine</b><br>(15, 16) | 2175      | 66                      | 1275                    | 834                   | 278                       | 27                  | 28,013             |

<sup>a</sup> Untranslated region.

<sup>b</sup> Open reading frame.

<sup>c</sup> Protein size includes amino acids in the signal peptide.

<sup>d</sup> Predicted deglycosylated molecular weight after signal peptide removal.

The complete cDNA sequence for the human *CDM6PR* is 2463 bp (14). It includes a 145-bp 5'-untranslated region, a single open reading frame of 831 bp, and a 1487-bp 3'-untranslated region. The open reading frame encodes for a protein consisting of 277 amino acids. It has a predicted deglycosylated  $M_r$  of 27,913, excluding the putative 26-amino acid signal peptide (14). The receptor consists of a 25-amino acid single membrane-spanning domain ( $M_r \approx 2588$ ), which separates the 159-amino-acid *N*-terminal region ( $M_r \approx 17,874$ ), containing five potential *N*-glycosylation sites, from the 67-amino-acid *C*-terminal region ( $M_r \approx 7486$ ). All potential *N*-glycosylation sites are utilized except Asn<sup>94</sup> (human) (17), with Asn<sup>87</sup> in the bovine receptor (human Asn<sup>113</sup>) being most important in enhancing receptor binding (18). The mature bovine, human, and murine *CDM6PR* proteins are 93% identical, with the bovine and murine sequence differing from that of human by 19 and 14 amino acids, respectively. The 67 intracellular amino acid residues are identical in these three species. Additional molecular characteristics of bovine, human, and murine *CDM6PR* genes are presented in Table 1.

The human *CDM6PR* gene has been cloned (19) and is located on chromosome 12 (14). It is distributed over 12 kb and is divided into seven exons. The mouse *CDM6pr* gene is located on chromosome 6, in a region syntenic with human chromosome 12 (20). Sequencing of the 10-kb mouse gene demonstrates that it also consists of seven exons. The 5'-untranslated region is encoded by exon 1. The extracellular, transmembrane, and cytoplasmic domains are encoded by exons 2 to 7, and exon 7 contains the 3'-untranslated region of the mRNA. The promoter region of the *CDM6PR* gene contains structural elements characteristic of promoters found in "housekeeping genes" (19).

## 1.2. *CDM6PR*: Structure and Function

The *CDM6PR* optimally binds phosphomannosyl residues in the pH range of 6.3 to 6.5, while ligand binding does not occur below pH 5 (21-23). Consequently, ligands are released from the receptor when they reach the acidic prelysosomal environment. Furthermore, the *CDM6PR* does not endocytose M6P glycoproteins because of its inability to bind ligands at neutral pH. A single *CDM6PR* molecule binds 1 mol of M6P, with a  $K_d$  of 8  $\mu$ M, and 0.5 mol of a high-mannose oligosaccharide with two phosphomonoesters, with a  $K_d$  of 0.2  $\mu$ M (23). This, coupled with the finding that the *CDM6PR* exists primarily as a homodimer (22, 24), suggests that two functional binding domains are required for optimal lysosomal enzyme recognition. The determination of the three-dimensional (3-D) structure of the bovine *CDM6PR* also revealed the presence of two molecules of *CDM6PR* per binding unit with a 40 Å linear distance between the two M6P binding sites (12).

The 3-D structure of the extracytoplasmic domain of the *CDM6PR* is similar to that of the biotin binding protein, avidin (12). The *CDM6PR* consists of two orthogonal  $\beta$ -sheets with only a single  $\alpha$ -helix at the *N*-terminal portion of the receptor. The six cysteine residues of the extracytoplasmic domain are disulfide-linked within a sheet, and they do not connect the two sheets. The disulfide bond between Cys<sup>106</sup> (human Cys<sup>132</sup>) and Cys<sup>141</sup> (human Cys<sup>167</sup>) of the bovine receptor seem to be particularly important because it brings two  $\beta$ -strand loops together to form the M6P binding pocket. Furthermore, two residues previously shown to function in M6P recognition, His<sup>105</sup> (human His<sup>131</sup>) and Arg<sup>111</sup> (human Arg<sup>137</sup>) (25), are located within this binding cavity. Asp<sup>103</sup>, Asn<sup>104</sup>, His<sup>105</sup> (human Asp<sup>129</sup>, Asn<sup>130</sup>, His<sup>131</sup>), a water molecule that is hydrogen bonded to the carboxylate of Asp<sup>103</sup> (human Asp<sup>129</sup>), and Mn<sup>+2</sup> interact with the phosphate portion of the M6P. These residues, along with those that interact with the mannose moiety, are conserved between the bovine, human, and murine *CDM6PR*s (13-16). The *CDM6PR* varies in its binding affinity for different lysosomal

enzymes (26). A model for the complex formed between the CDM6PR and b-glucuronidase suggests that these binding affinity variations result from differences in both the location and number of phosphomannosyl residues on the oligosaccharide and the spatial orientation of these sugar chains on the lysosomal enzyme molecule (12).

The signals required for appropriate intracellular trafficking of lysosomal enzymes by the CDM6PR are contained in the cytoplasmic domain. This region of the receptor is highly conserved between species, with the 67 residues in the cytoplasmic domain being identical in the bovine, human, and murine receptors (13-16). Although the CDM6PRs on the cell surface do not bind extracellular M6P glycoproteins, membrane-bound receptors are nevertheless rapidly internalized (21). One internalization signal is Phe<sup>13</sup> (human Phe<sup>223</sup>) and Phe<sup>18</sup> (human Phe<sup>228</sup>), with the latter amino acid being more important (27, 28). A second internalization signal involves Tyr<sup>45</sup> (human Tyr<sup>255</sup>). The dileucine-containing sequence, Leu<sup>64</sup>–Leu<sup>65</sup> (human Leu<sup>274</sup>–Leu<sup>275</sup>), is a third internalization signal, and it is also required for receptor sorting in the Golgi apparatus (28, 29). Another signal in the cytoplasmic tail that enables proper endosomal sorting of the CDM6PR depends on Phe<sup>18</sup>–Trp<sup>19</sup> (human Phe<sup>228</sup>–Trp<sup>229</sup>) and the modulation of this signal by Cys<sup>34</sup> (human Cys<sup>244</sup>) (30).

### 1.3. CDM6PR: Biological Functions

The CDM6PR is expressed in all cells, where it functions in intracellular lysosomal enzyme transport. Because the M6P/IGF2R can also perform a similar function, mice homozygous for a mutated *CDM6pr* gene were produced to determine the biological effects of this receptor (31, 32). The homozygous *CDM6pr* null mice are phenotypically normal and are also fertile, demonstrating that the *CDM6pr* is not required for egg fertilization or organogenesis. However, these animals exhibited defects in the targeting of multiple lysosomal enzymes, and increased levels of phosphorylated lysosomal enzymes were present in the body fluids. Thus, a physiological level of the M6p/Igf2r does not compensate for CDM6pr loss. This may result from the two M6PRs targeting distinct subsets of lysosomal enzymes because of differences in their binding characteristics (23, 26). In summary, the principal biological function of the CDM6PR is to enable efficient intracellular targeting of lysosomal enzymes.

## 2. Cation-Independent Mannose 6-Phosphate Receptor

### 2.1. M6P/IGF2R: Molecular Characteristics

The M6P/IGF2R is a 275-kDa monomeric multifunctional transmembrane-binding receptor with independent, high-affinity binding sites for phosphomannosyl glycoproteins (33-35), IGF2 (33, 34, 36, 37), and retinoic acid (38), along with a lower-affinity binding site for the urokinase-type plasminogen activator receptor (uPAR) (Fig. 1) (39). The complete cDNA sequences have been determined for the bovine (40), chicken (41), human (33, 42), murine (43, 44), and rat (34) *M6P/IGF2Rs* (Table 2). The overall structures of the M6P/IGF2R in these species is identical. It is comprised of a large extracellular domain containing the M6P, IGF2, and uPAR binding sites (33-37, 39), a single transmembrane domain, and a small cytoplasmic tail that presumably contains the retinoic acid binding site (38).

**Table 2. M6P/IGF2R Species Comparisons**

| Species               | cDNA<br>(bp)      | 5'UTR <sup>a</sup><br>(bp) | 3'UTR <sup>a</sup><br>(bp) | ORF <sup>b</sup><br>(bp) | Protein <sup>c</sup><br>(aa) | Signal<br>Peptide<br>(aa) | <i>M<sub>r</sub></i> <sup>d</sup> |
|-----------------------|-------------------|----------------------------|----------------------------|--------------------------|------------------------------|---------------------------|-----------------------------------|
| <b>Bovine</b><br>(40) | 7912 <sup>e</sup> | 152                        | 263 <sup>e</sup>           | 7497                     | 2499                         | 44                        | 269,980                           |

|                              |      |     |      |      |      |    |         |
|------------------------------|------|-----|------|------|------|----|---------|
| <b>Chicken</b><br>(41)       | 8767 | 60  | 1297 | 7410 | 2470 | 23 | 273,362 |
| <b>Human</b><br>(33, 42)     | 9090 | 147 | 1470 | 7473 | 2491 | 40 | 270,294 |
| <b>Mouse</b> (43,<br>44)     | 8894 | 132 | 1316 | 7446 | 2482 | 34 | 270,454 |
| <b>Rat</b> <sup>f</sup> (34) | 8810 | 83  | 1290 | 7437 | 2479 | 35 | 269,999 |

<sup>a</sup> Untranslated region.

<sup>b</sup> Open reading frame.

<sup>c</sup> Protein size includes amino acids in the signal peptide.

<sup>d</sup> Predicted deglycosylated molecular weight after signal peptide removal.

<sup>e</sup> 3'UTR sequence is incomplete.

<sup>f</sup> The 5' portion of the rat cDNA sequence is in GenBank (Accession No. U59809).

The complete cDNA sequence for the human *M6P/IGF2R* gene is 9090 bp (33). It is comprised of a 147-bp 5'-untranslated sequence, a large open reading frame of 7473 bp, and a 1470-bp 3'-untranslated region. The open reading frame encodes a protein of 2491 amino acid residues with a deglycosylated  $M_r$  of 270,294 after removal of the putative 40-amino-acid signal peptide. The human receptor consists of a large 2264 residue extracellular domain ( $M_r \approx 249,638$ ) comprised of 15 contiguous repeat regions averaging 150 bp, with a small gap of 27 amino acids after repeat 14. The rest of the molecule consists of a single 23-residue transmembrane region ( $M_r \approx 2346$ ) and a small 164-residue cytoplasmic domain ( $M_r \approx 18,345$ ). The receptor contains 19 potential extracellular *N*-linked glycosylation sites, the majority of which appear to be used (33, 45). The nucleotide and amino acid sequences of the *M6P/IGF2R* are highly conserved among mammalian species. The bovine, mouse, and rat amino acid sequences are approximately 80% identical with the human receptor (33, 34, 40, 43, 44); however, the chicken receptor sequence has only 60% homology with the human receptor (41). Further molecular characteristics of the bovine, chicken, human, murine, and rat *M6P/IGF2R* genes are presented in Table 2.

The mouse *M6p/Igf2r* gene has been cloned (44). It is located in a 2-megabase region on chromosome 17 (46) that shows a parental bias in replication timing, a characteristic of chromosomal regions containing imprinted genes (47). The *M6p/Igf2r* is 93 kb long and composed of 48 exons. It encodes a predicted protein of 2482 residues (Table 2). The extracellular portion of the receptor is encoded by exons 1 to 46, with each of the 15 repeat motifs being determined by 3 to 5 exons. A single fibronectin, type II-like element is found in exon 39. The transmembrane portion of the receptor is encoded by exon 46, and the cytoplasmic region is encoded by exons 46 to 48. Exon 48 also contains the 5'-untranslated region of the mRNA. A 266 bp 5'-flanking DNA region contains *M6p/Igf2r* promoter activity (48). The promoter region is highly GC-rich (70%) and is hypermethylated only on the paternal allele, because this mouse gene is paternally imprinted (49). A 54-bp segment within this flanking region (-52 to +2) contains two E-box domains and potential Sp1 and NGF-1A transcription factor binding sites (48). The genomic structure of the human *M6P/IGF2R* gene, which maps to chromosomal location 6q26 (50), has also been determined (51). It is approximately 136-kbp long, and it is similarly composed of 48 exons with all exon/intron boundaries identical to those in the mouse.

The positions of the exon/intron splice junctions are conserved between several of the repeats in the *M6P/IGF2R* and the homologous extracellular region of the *CDM6PR* gene (44, 51). These genomic

structure similarities suggest that both genes evolved from an ancestral gene that contained at least four exons. Furthermore, the sequence duplications that resulted in the formation of the *M6P/IGF2R* gene probably preceded the divergence of mammals because this receptor is also present in birds and amphibians (41, 52, 53).

## 2.2. M6P/IGF2R: Structure and Function

The M6P/IGF2R optimally binds phosphomannosyl residues in the pH range of 5.7 to 6.3, with no detectable ligand binding below pH 5 (23). This is consistent with its ability to traffic and release lysosomal enzymes into the acidic environment of the prelysosomes. In contrast to CDM6PR, M6P/IGF2R retains the capacity to bind ligands at neutral pH. This enables it to also endocytose extracellular ligands for lysosome delivery. The M6P/IGF2R only binds 2 mol of monovalent M6P ligand ( $K_d = 7 \mu\text{M}$ ) and 1 mol of a diphosphorylated oligosaccharide ( $K_d = 2\text{nM}$ ) (54) even though there is significant sequence homology between the extracellular region of CDM6PR and all 15 repeat domains in the M6P/IGF2R (13, 40). This indicates there are only two M6P binding sites per receptor, and they have been localized to repeats 3 and 9 (Fig. 1) (55).

In the bovine M6P/IGF2R, Arg<sup>435</sup> (human Arg<sup>426</sup>) in repeat 3 and Arg<sup>1334</sup> (human Arg<sup>1325</sup>) in repeat 9 are essential for high-affinity M6P binding (55). Sequence alignment of these two M6P/IGF2R repeat domains has been performed based upon the 3-D structure of the extracellular portion of the bovine CDM6PR (12). These domain alignment models demonstrate that both of these amino acids are conserved and align with Arg<sup>111</sup> (human Arg<sup>137</sup>) in the CDM6PR. They also show that the Asp<sup>103</sup> (human Asp<sup>129</sup>) that coordinates the divalent cation in the CDM6PR is absent in repeats 3 and 9 of the M6P/IGF2R, providing a possible explanation for its cation independent binding of M6P moieties. A 3-D model of the complex between the CDM6PRs and b-glucuronidase (12) further suggests that the binding affinities of these two M6PRs for lysosomal enzymes differ (26) because the location and number of M6P moieties on the oligosaccharides vary significantly between the different lysosomal enzymes.

In contrast to those of the chicken (52, 53) and frog (53), mammalian M6P/IGF2Rs bind IGF2 in addition to phosphomannosyl glycoproteins (33, 34, 36, 42, 56): 1 mol of M6P/IGF2R binds 1 mol of IGF2 with a  $K_d = 0.2\text{nM}$  at an optimum pH of 7.4 (57). IGF1 binds to the receptor with a substantially lower affinity ( $K_d = 0.4 \mu\text{M}$ ), and insulin binding is undetectable. The M6P and IGF2 binding sites of this receptor are distinct. Therefore, the M6P/IGF2R can bind both ligands simultaneously; however, binding of lysosomal enzymes impairs IGF2 binding, and vice versa (56). The minimal binding site for IGF2 is between amino acids 1508 and 1575 (human), which is located in the *N*-terminal portion of repeat 11 (36, 37, 58). Furthermore, the point mutation Ile<sup>1572</sup> → Thr<sup>1572</sup> (human) completely abolishes IGF2 binding to the M6P/IGF2R (58). This region in the chicken M6P/IGF2R has a low sequence homology (15% to 20%), with that in mammals, probably accounting for the failure of the chicken receptor to bind IGF2 (41, 58). The 43-residue fibronectin type II-like domain in repeat 13 also enhances the affinity of IGF2 binding to the M6P/IGF2R (59). Thus, high-affinity IGF2 binding requires the presence of both the IGF2 binding site in repeat 11 and the enhancer site in repeat 13 (Fig. 1). The ability of IGF2 to bind to the M6P/IGF2R is also facilitated by glypican 3 (60). This betaglycan maps to the Simpson–Golabi–Behmel syndrome locus Xq26, a disorder with a phenotype similar to that of the Beckwith–Wiedemann syndrome.

The M6P/IGF2R signals required for endocytosis and for efficient intracellular lysosomal enzyme sorting are located in the cytoplasmic tail (61, 62). The signal for rapid internalization of the M6P/IGF2R is the sequence Y<sup>24</sup>KYSKV<sup>29</sup> (human Y<sup>2351</sup>KYSKV<sup>2356</sup>), with most of the activity mediated by the four last residues (61). The single contiguous stretch of a casein kinase II site followed by two leucines in the *C*-terminal region of the cytoplasmic tail, D<sup>157</sup>DSDEDLL<sup>164</sup> (human D<sup>2482</sup>DSDEDLL<sup>2489</sup>), is the major sorting determinant for the M6P/IGF2R (62). These internalization and sorting signal sequences are identical in bovine, human, murine, and rat

## M6P/IGF2Rs.

### 2.3. M6P/IGF2R: Biological Functions

The *M6P/IGF2R* gene is expressed ubiquitously in tissue (56), and a truncated form of the receptor circulates in the blood, where it functions as an IGF2 binding protein (63-65). The *M6P/IGF2R* is also developmentally regulated (56). In mice and rats, the receptor level is highest in 16- to 20-day-old fetuses, and it declines rapidly after birth. The predominant site of *M6P/IGF2R* expression during mouse embryogenesis is in the developing myocardium, suggesting that this receptor plays an important role in heart development (66). The finding that the inheritance of homozygosity for an inactive *M6p/Igf2r* allele results in major heart abnormalities and is embryogenically lethal supports this postulate (67, 68). Fetal lethality stems from an excess of IGF2 during development because the introduction of an *Igf2* null allele rescues the *M6p/Igf2r* null mouse (68, 69). These transgenic mouse studies, and those of Kröner et al. (70), contradict the postulate that the M6P/IGF2R functions in mitogenic signaling (71). Rather, they demonstrate that the main role of the M6P/IGF2R is to control the bioavailability of IGF2 through its endocytosis and delivery to the lysosomes for degradation. To date, the majority of studies indicate that the primary mitogenic effects of IGF2 are mediated through the IGF1R (72) and the insulin receptor (73), not the M6P/IGF2R.

The *M6p/Igf2r* is also highly expressed in neurons of the forebrain, with highest expression in the pyramidal cells, the polymorphic layers of the hippocampus, and the granule cell layer of the dentate gyrus (74). These neural regions are involved in emotional behavior, information processing, and memory formation, suggesting that *M6p/Igf2r* may assist in the development of these brain functions. This postulate has been reinforced by the recent identification of the *M6P/IGF2R* as the first putative “IQ gene” (75). By comparing children with an IQ of 160 or higher to those with an average IQ, it was shown that the *M6P/IGF2R* is linked with general cognitive ability (“g”). The discovery that the *M6P/IGF2R* contains a high-affinity binding site for retinoic acid further strengthens the postulate that this receptor functions in brain and heart development (38).

Nonlysosomal enzymes containing M6P residues also bind to the M6P/IGF2R. These include epidermal growth factor receptor (76), herpes simplex glycoprotein D (77), proliferin (78), renin (79), thyroglobulin (80), and transforming growth factor b (TGFb) latent complex (81, 82). The importance of TGFb latent complex binding to the M6P/IGF2R is most clearly understood. It facilitates the proteolytic activation of TGFb by plasmin (82, 83), a reaction whose efficiency may be enhanced by the ability of uPAR to also bind directly to the M6P/IGF2R (39). Thus, in addition to the intracellular sorting of lysosomal enzymes, the M6P/IGF2R is involved in activating the growth inhibitor, TGFb, and degrading the mitogen, IGF2. The M6P/IGF2R therefore functions as a “guardian of the cell” by controlling the extracellular levels of potent growth factors and proteolytic enzymes that govern cell death, proliferation, and invasion, suggesting it also has tumor suppressor activity.

### 2.4. M6P/IGF2R and Cancer

Recent findings demonstrate that loss of heterozygosity at the *M6P/IGF2R* locus occurs in 60% of human hepatocellular carcinomas (HCCs) (84-86) and 30% of breast cancer (87). Missense mutations in the remaining allele have also been identified in repeats 9, 10, and 14 (88), and in a 6-amino-acid MARCKS (myristoylated alanine-rich protein kinase C substrates) consensus sequence present in the cytoplasmic domain (89). The *M6P/IGF2R* also contains a poly-G region in repeat 9 that is a common mutational target in human colon, gastric, and endometrial tumors with mismatch repair deficiencies and microsatellite instability (90, 91). Furthermore, *M6P/IGF2R* allelic loss is an early event in liver carcinogenesis giving rise to the clonal expansion of phenotypically normal, *M6P/IGF2R*-mutated preneoplastic cells from which HCCs often develop (86). *M6P/IGF2R* inactivation similarly occurs early in breast cancer formation (87, 92). Thus, the *M6P/IGF2R* is mutated frequently and early in a variety of human cancers.

### 2.5. M6P/IGF2R and Genomic Imprinting

A full parental complement of autosomal genes is inherited by all offspring. However, not all genes

are biallelically expressed because of genomic imprinting, an epigenetic form of gene regulation that results in the expression of only one parental allele (93-96). The *M6p/Igf2r* gene is imprinted in both mice (68, 97) and rats (98). It is expressed only from the maternal allele after embryonic implantation in all mouse tissues, except possibly the brain (99, 100). Methylation of a CpG-rich region in intron 2 of the expressed maternal *M6p/Igf2r* allele has been shown to carry the imprint signal for this gene (49, 101). Furthermore, this intron region appears normally to function as the promoter of an antisense transcript that is expressed only from the repressed paternal allele, indicating that a form of expression competition regulates imprinting of the mouse *M6p/Igf2r* gene (101).

The human *M6P/IGF2R* gene is similarly expressed monoallelically in preterm, postimplantation embryonic tissue (102). Its behavior then diverges from that in mice and rats, because monoallelic expression of the *M6P/IGF2R* becomes a polymorphic trait in humans, with most adults expressing both alleles (103-105). Thus, human cancer predisposition could result from *M6P/IGF2R* imprinting, a postulate supported by the discovery that 50% of patients with Wilms' tumor, a juvenile kidney tumor, are imprinted at this locus (106). The marked species difference in *M6P/IGF2R* imprinting between mice and humans may also have important implications for human carcinogen risk assessment based upon animal studies because only one rather than two alleles of this tumor suppressor gene would need to be inactivated in mice during carcinogenesis (88).

### 3. Conclusion

There are two receptors that bind phosphomannosyl glycoproteins with high affinity, the CDM6PR and the M6P/IGF2R. Both receptors are required for effective intracellular trafficking of lysosomal enzymes to the lysosomes, but only M6P/IGF2R can endocytose extracytoplasmic ligands. The CDM6PR contains a single M6P binding site, whereas the M6P/IGF2R possesses independent binding sites for IGF2, M6P, retinoic acid, and uPAR. Consequently, the M6P/IGF2R not only is involved in lysosomal enzyme trafficking, but also plays an important role in both embryogenesis and carcinogenesis. Species- and tissue-dependent imprinting of the *M6P/IGF2R* add further to the biological complexity of these receptors.

### 4. Acknowledgments

This work was supported by NIH grants CA25951 and ES08823 and DOD grant DAMD17-98-1-8305. I wish to thank Nancy Dahms for reading this article and making helpful suggestions. Visit the genomic imprinting website for additional genetic information on the *M6PR* genes (<http://www.geneimprint.com>).

### Bibliography

1. S. Kornfeld (1990) *Biochem. Soc. Trans.* **18**, 367–374.
2. A. Varki and S. Kornfeld (1980) *J. Biol. Chem.* **255**, 8398–8401.
3. A. Hasilik, A. Waheed, and K. von Figura (1981) *Biochem. Biophys. Res. Commun.* **98**, 761–767.
4. A. Waheed, A. Hasilik, and K. von Figura (1981) *J. Biol. Chem.* **256**, 5717–5721.
5. M. L. Reitman, A. Varki, and S. Kornfeld (1981) *J. Clin. Invest.* **67**, 1574–1579.
6. A. P. Varki, M. L. Reitman, and S. Kornfeld (1981) *Proc. Natl. Acad. Sci. USA* **78**, 7773–7777.
7. D. E. Sleat et al. (1996) *J. Biol. Chem.* **271**, 19191–19198.
8. G. Griffiths et al. (1988) *Cell* **52**, 329–341.
9. B. Hoflack and S. Kornfeld (1985) *J. Biol. Chem.* **260**, 12008–12014.
10. K. Drickamer and M. E. Taylor (1993) *Annu. Rev. Cell Biol.* **9**, 237–264.
11. M. Stein, H. E. Meyer, A. Hasilik, and K. von Figura (1987) *Biol. Chem. Hoppe Seyler* **368**,



927–936.

12. D. L. Roberts, D. J. Weix, N. M. Dahms, and J.-J. P. Kim (1998) *Cell* **93**, 639–648.
13. N. M. Dahms, P. Lobel, J. Breitmeyer, J. M. Chirgwin, and S. Kornfeld (1987) *Cell* **50**, 181–192.
14. R. Pohlmann et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 5575–5579.
15. A. Köster, G. Nagel, K. von Figura, and R. Pohlmann (1991) *Biol. Chem. Hoppe Seyler* **372**, 297–300.
16. Z. M. Ma, J. H. Grubb, and W. S. Sly (1991) *J. Biol. Chem.* **266**, 10589–10595.
17. M. Wendland et al. (1991) *J. Biol. Chem.* **266**, 4598–4604.
18. Y. Zhang and N. M. Dahms (1993) *Biochem. J.* **295**, 841–848.
19. H. J. Klier, K. von Figura, and R. Pohlmann (1991) *Eur. J. Biochem.* **197**, 23–28.
20. T. Ludwig et al. (1992) *J. Biol. Chem.* **267**, 12211–12219.
21. M. Stein et al. (1987) *EMBO J.* **6**, 2677–2681.
22. N. M. Dahms and S. Kornfeld (1989) *J. Biol. Chem.* **264**, 11458–11467.
23. P. Y. Tong and S. Kornfeld (1989) *J. Biol. Chem.* **264**, 7970–7975.
24. M. Stein et al. (1987) *Biol. Chem. Hoppe Seyler* **368**, 937–947.
25. M. Wendland, A. Waheed, K. von Figura, and R. Pohlmann (1991) *J. Biol. Chem.* **266**, 2917–2923.
26. D. E. Sleat and P. Lobel (1997) *J. Biol. Chem.* **272**, 731–738.
27. K. F. Johnson, W. Chan, and S. Kornfeld (1990) *Proc. Natl. Acad. Sci. USA* **87**, 10010–10014.
28. K. Denzer et al. (1997) *Biochem. J.* **326**, 497–505.
29. K. F. Johnson and S. Kornfeld (1992) *J. Biol. Chem.* **267**, 17110–17115.
30. A. Schweizer, S. Kornfeld, and J. Rohrer (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14471–14476.
31. A. Köster et al. (1993) *EMBO J.* **12**, 5219–5223.
32. T. Ludwig et al. (1993) *EMBO J.* **12**, 5225–5235.
33. D. O. Morgan et al. (1987) *Nature* **329**, 301–307.
34. R. G. MacDonald et al. (1988) *Science* **239**, 1134–1137.
35. B. Westlund, N. M. Dahms, and S. Kornfeld (1991) *J. Biol. Chem.* **266**, 23233–23239.
36. N. M. Dahms, D. A. Wick, and M. A. Brzycki-Wessell (1994) *J. Biol. Chem.* **269**, 3802–3809.
37. B. Schmidt et al. (1995) *J. Biol. Chem.* **270**, 14975–14982.
38. J. X. Kang, Y. Li, and A. Leaf (1997) *Proc. Natl. Acad. Sci. USA* **94**, 13671–13676.
39. A. Nykjaer et al. (1998) *J. Cell Biol.* **141**, 815–828.
40. P. Lobel, N. M. Dahms, and S. Kornfeld (1988) *J. Biol. Chem.* **263**, 2563–2570.
41. M. Zhou, Z. Ma, and W. S. Sly (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9762–9766.
42. A. Oshima et al. (1988) *J. Biol. Chem.* **263**, 2553–2562.
43. T. Ludwig et al. (1994) *Gene* **142**, 311–312.
44. G. Szebenyi and P. Rotwein (1994) *Genomics* **19**, 120–129.
45. N. M. Dahms (1996) *Biochem. Soc. Trans.* **24**, 136–141.
46. G. Laureys, D. E. Barton, A. Ullrich, and U. Francke (1988) *Genomics* **3**, 224–229.
47. D. Kitsberg et al. (1993) *Nature* **364**, 459–463.
48. Z. Liu, D. W. Mittanck, S. Kim, and P. Rotwein (1995) *Mol. Endocrinol.* **9**, 1477–1487.
49. R. Stöger et al. (1993) *Cell* **73**, 61–71.
50. P. H. Rao et al. (1994) *Cytogenet. Cell Genet.* **66**, 272–273.
51. J. K. Killian and R. L. Jirtle (1998) *Mamm. Genome* **10**, 74–77.

52. W. M. Canfield and S. Kornfeld (1989) *J. Biol. Chem* **264**, 7100–7103.
53. K. B. Clairmont and M. P. Czech (1989) *J. Biol. Chem.* **264**, 16390–16392.
54. P. Y. Tong, W. Gregory, and S. Kornfeld (1989) *J. Biol. Chem.* **264**, 7962–7969.
55. N. M. Dahms et al. (1993) *J. Biol. Chem.* **268**, 5457–5463.
56. S. Kornfeld (1992) *Annu. Rev. Biochem.* **61**, 307–330.
57. P. Y. Tong, S. E. Tollefsen, and S. Kornfeld (1988) *J. Biol. Chem.* **263**, 2585–2588.
58. F. Garmroudi et al. (1996) *Mol. Endocrinol.* **10**, 642–651.
59. G. R. Devi, J. C. Byrd, D. H. Slentz, and R. G. MacDonald (1998) *Mol. Endocrinol.* **12**, 1661–1672.
60. G. Pilia et al. (1996) *Nat. Genet.* **12**, 241–247.
61. M. Jadot, W. M. Canfield, W. Gregory, and S. Kornfeld (1992) *J. Biol. Chem.* **267**, 11069–11077.
62. H. J. Chen, J. Yuan, and P. Lobel (1997) *J. Biol. Chem.* **272**, 7003–7012.
63. W. Kiess et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7720–7724.
64. R. G. MacDonald, M. A. Tepper, K. B. Clairmont, S. B. Perregaux, and M. P. Czech (1989) *J. Biol. Chem.* **264**, 3256–3261.
65. K. J. Valenzano, J. Remmler, and P. Lobel (1995) *J. Biol. Chem.* **270**, 16441–16448.
66. K. M. McCormick, N. M. Dahms, and J. Lough (1996) *Dev. Dyn.* **207**, 195–203.
67. M. M. Lau et al. (1994) *Genes Dev.* **8**, 2953–2963.
68. Z. Q. Wang, M. R. Fung, D. P. Barlow, and E. F. Wagner (1994) *Nature* **372**, 464–467.
69. A. J. Filson, A. Louvi, A. Efstratiadis, and E. J. Robertson (1993) *Development* **118**, 731–736.
70. C. Körner, B. Nurnberg, M. Uhde, and T. Braulke (1995) *J. Biol. Chem.* **270**, 287–295.
71. I. Nishimoto et al. (1989) *J. Biol. Chem.* **264**, 14029–14038.
72. M. P. Czech, R. E. Lewis, and S. Corvera (1989) *Ciba Found. Symp.* **145**, 27–41.
73. A. Morrione et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3777–3782.
74. M. E. Couce, A. J. Weatherington, and J. F. McGinty (1992) *Endocrinology* **131**, 1636–1642.
75. M. J. Chorney et al. (1998) *Psychol. Sci.* **9**, 159–166.
76. G. Todderud and G. Carpenter (1988) *J. Biol. Chem.* **263**, 17893–17896.
77. C. R. Brunetti et al. (1994) *J. Biol. Chem.* **269**, 17067–17074.
78. S. J. Lee and D. Nathans (1988) *J. Biol. Chem.* **263**, 3521–3527.
79. P. L. Faust, J. M. Chirgwin, and S. Kornfeld (1987) *J. Cell Biol.* **105**, 1947–1955.
80. V. Herzog, W. Neumuller, and B. Holzmann (1987) *EMBO J.* **6**, 555–560.
81. A. F. Purchio et al. (1988) *J. Biol. Chem.* **263**, 14211–14215.
82. K. S. Kovacina, G. Steele-Perkins, and R. A. Rothe (1989) *Mol. Endocrinol.* **3**, 901–906.
83. P. A. Dennis and D. B. Rifkin (1991) *Proc. Natl. Acad. Sci. USA* **88**, 580–584.
84. A. T. De Souza et al. (1995) *Oncogene* **10**, 1725–1729.
85. A. T. De Souza et al. (1995) *Nat. Genet.* **11**, 447–449.
86. T. Yamada, A. T. De Souza, S. Finkelstein, and R. L. Jirtle (1997) *Proc. Natl. Acad. Sci. USA* **94**, 10351–10355.
87. G. R. Hankins et al. (1996) *Oncogene* **12**, 2003–2009.
88. A. T. De Souza, T. Yamada, J. J. Mills, and R. L. Jirtle (1997) *FASEB J.* **11**, 60–67.
89. P. J. Blackshear et al. (1992) *J. Biol. Chem.* **267**, 13540–13546.
90. R. F. Souza et al. (1996) *Nat. Genet.* **14**, 255–257.
91. H. Ouyang et al. (1997) *Cancer Res.* **57**, 1851–1854.
92. S. A. Chappell, T. Walsh, R. A. Walker, and J. A. Shaw (1997) *Br. J. Cancer* **76**, 1558–1561.

93. D. P. Barlow (1995) *Science* **270**, 1610–1613.
94. C. Sapienza (1995) *Dev. Genet.* **17**, 185–187.
95. B. Tycko (1997) *Mutat. Res.* **386**, 131–140.
96. M. A. Surani (1998) *Cell* **93**, 309–312.
97. D. P. Barlow et al. (1991) *Nature* **349**, 84–87.
98. J. J. Mills, J. G. Falls, A. T. De Souza, and R. L. Jirtle (1998) *Oncogene* **16**, 2797–2802.
99. W. Lerchner and D. P. Barlow (1997) *Mech. Dev.* **61**, 141–149.
100. J. F. Hu, H. Oruganti, T. H. Vu, and A. R. Hoffman (1998) *Mol. Endocrinol.* **12**, 220–232.
101. A. Wutz et al. (1997) *Nature* **389**, 745–749.
102. A. Wutz, O. W. Smrzka, and D. P. Barlow (1998) *Novartis Found. Symp.* **214**, 251–259.
103. V. M. Kalscheuer et al. (1993) *Nat. Genet.* **5**, 74–78.
104. O. Ogawa et al. (1993) *Hum. Mol. Genet.* **2**, 2163–2165.
105. Y. Xu, C. G. Goodyer, C. Deal, and C. Polychronakos (1993) *Biochem. Biophys. Res. Commun.* **197**, 747–754.
106. Y. Q. Xu, P. Grundy, and C. Polychronakos (1997) *Oncogene* **14**, 1041–1046.

## MAP Kinases

Mitogen-activated protein (MAP) kinases are stimulated by virtually every **mitogen** known. They are activated after the activation of [tyrosine kinase receptors](#), and they are themselves activated by **phosphorylation** on [tyrosine](#) and [threonine](#) residues (1). MAP kinases are activated by a complex series of [protein–protein interactions](#) and upstream protein kinases and phosphatases.

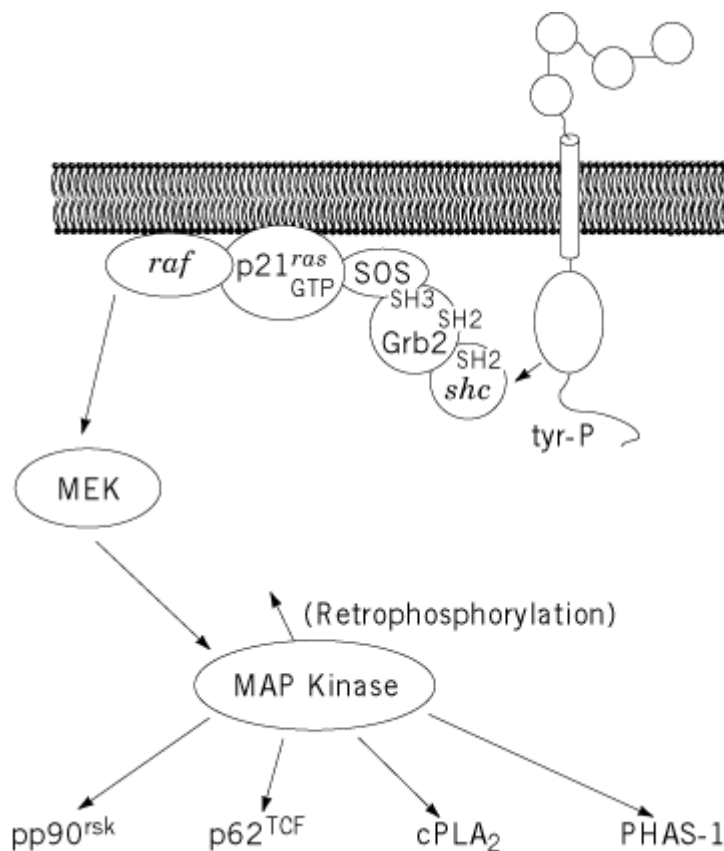
MAP kinase is directly activated by MAP kinase kinase, or (MEK) (2), a dual-specificity kinase exhibiting sequence [homology](#) to the byr1 and STE7 gene products of *Saccharomyces pombe* and *S. cerevisiae*. Although the specific pathways by which [growth factors](#) activate this enzyme remain unknown, activation of the *raf* proto-**oncogene** (3) is one of the predominant mechanisms by which MEK is activated. In some systems, however, *c-raf* is not the major MEK kinase that is stimulated by growth factors, and additional MEK kinases have been identified.

*Raf* activation is usually accompanied by enhanced phosphorylation of the enzyme on serine and threonine residues (3), although it is still unclear whether this phenomenon reflects autophosphorylation. Considerable evidence now indicates that activation of the *ras* proto- **oncogene** is a prerequisite for growth factor-dependent *raf* activation, and further that *raf* may be activated via a direct interaction with *ras* (4). Significant progress has been made in dissecting the molecular events involved in the regulation of p21<sup>ras</sup> by [tyrosine kinase receptors](#). Activation appears to involve exchange of GDP by GTP on the ras protein, catalyzed by a nucleotide exchange factor. In most mammalian cells, the growth-factor-sensitive exchange factor is presumed to be a **homologue** of the yeast protein *son of sevenless* (SOS) (5), although it is possible that other exchange factors exist. The increased binding of GTP to ras catalyzed by SOS is likely to be activated by the targeting of SOS to ras via the interaction with the adapter protein Grb2, which is mediated by the binding of one of the **SH3** domains of Grb2 with a proline-rich region in the C-terminus of SOS. This Grb2-SOS interaction is itself thought to be stimulated by the binding of Grb2 to tyrosine-phosphorylated proteins through its SH2 domain. The growth factor-induced phosphotyrosine–SH2 interaction may occur by direct binding to an autophosphorylated receptor, as is the case for the [epidermal growth](#)

[factor](#) (EGF) receptor, or it may occur through a surrogate phosphorylated protein, such as Shc or IRS-1.

Upon its activation, MAP kinase can translocate into the [nucleus](#), where it catalyzes the phosphorylation of [transcription factors](#) such as p62<sup>TCF</sup>, initiating a program of [transcription](#) that leads the cell to commit to a proliferative or differentiative cycle (Fig. 1). MAP kinase can also phosphorylate a number of other proteins involved in cellular signaling, including other kinases and [phospholipases](#). The role of the MAP kinase pathway in cell growth and differentiation has been explored with pharmacological tools, as well as by the expression of dominant interfering mutant proteins in the pathway. Thus, although numerous hormonal agents can activate MAP kinase in different cell types, the pathway is not universally used in signal transduction, even for responses produced by some tyrosine kinase receptors (6).

**Figure 1.** Regulation of the MAP kinase pathway. MAP kinases are activated by a variety of growth and differentiative signals. [Tyrosine kinase receptors](#) are thought to stimulate this pathway primarily through the activation of the *ras* proto-oncogene. Activation of growth factor receptors leads to the tyrosine phosphorylation of Shc, which then binds to Grb2 via an SH2 interaction. Shc can target the Grb2/SOS complex to the membrane, where SOS can stimulate the exchange of GDP for GTP on *ras*, thus switching this protein to the active state. Once activated, *ras* can recruit to the membrane protein kinases, such as *raf*, that can phosphorylate and activate the kinase MEK. Additional signals involving the tyrosine phosphatase Syp are also required in some cases for activation of *raf*. MEK is a dual-specificity kinase that phosphorylates MAP kinase on threonine and tyrosine residues, leading to its activation. Upon activation, MAP kinase can phosphorylate a variety of substrates, including other kinases (such as p90<sup>rsk</sup>), phospholipases, regulators of translation (such as PHAS1), and numerous transcription factors (including p62<sup>TCF</sup> and others).



Numerous studies indicate that there is a well-regulated system of feedback regulation of signals leading to MAP kinase activation (7). Additionally, the pathway may also be self-limiting regarding

the cellular consequences of activation. In the case of neuronal differentiation, the kinetics of MAP kinase activation are thought to play a major role in determining whether a cell commits to differentiate. In general, the prolonged activation and nuclear translocation of MAP kinase is associated with neurite outgrowth; other mitogenic stimuli that activate MAP kinase transiently do not produce neuronal differentiation (8). In some cells, persistent activation of MAP kinase leads to growth arrest and differentiation, which is associated with induction of the cyclin-dependent kinase inhibitor p21<sup>cip1/Waf1</sup>, and the subsequent down-regulation of cyclin-dependent kinase-4 activity. These effects can be blocked by MEK inhibition, suggesting that the sustained activation of MAP kinase leads to the phosphorylation of proteins that are perhaps the products of early immediate genes, changing dramatically the phenotype of the cell.

### Bibliography

1. A. J. Rossomondo, D. M. Rayne, M. J. Weber, and T. W. Sturgill (1991) *Proc. Natl. Acad. Sci. USA* **86**, 6940–6943.
2. C. M. Crews, A. Alessandrini, and R. Erikson (1992) *Science* **258**, 478–480.
3. U. R. Rapp (1992) *Oncogene* **6**, 495–500.
4. X. F. Zhang, J. Settleman, J. M. Kyriakis, E. Takeuchi-Suzuki, S. J. Elledge, M. S. Marshall, J. T. Bruder, U. R. Rapp, and J. Avruch (1993) *Nature* **364**, 308–313.
5. D. Bowtell, P. Fu, M. Simon, and P. Senior (1992) *Proc. Natl. Acad. Sci. USA* **89**, 6511–6515.
6. D. F. Lazar, M. J. Brady, R. J. Wiese, C. C. Mastick, S. B. Waters, K. Yamauchi, J. E. Pessin, P. Cuatrecasas, and A. R. Saltiel (1995) *J. Biol. Chem.* **270**, 20801–20807.
7. K. H. Holt, S. B. Waters, S. Okada, K. Yamauchi, S. J. Decker, A. R. Saltiel, D. G. Motto, G. A. Koretsky, and J. E. Pessin (1996) *J. Biol. Chem.* **271**, 8300–8306.
8. S. J. Decker (1996) *J. Biol. Chem.* **270**, 30841–30844.

### Marker Exchange Mutagenesis

Marker exchange describes the movement of **genes** or marker **cassettes** between different locations in [chromosomes](#). This process is often used in order to probe the effect of spatial relationships between genes on their expression. For example, bacterial genes that have been cloned on **multicopy plasmids** and subjected to mutagenesis may be returned to the chromosome by double homologous [recombination](#) events. Marker exchange occurs spontaneously in rec<sup>+</sup> strains, and the plasmid can often be readily eliminated or excluded. Marker exchange can also be used to introduce chromosomal **alleles** into a **plasmid**.

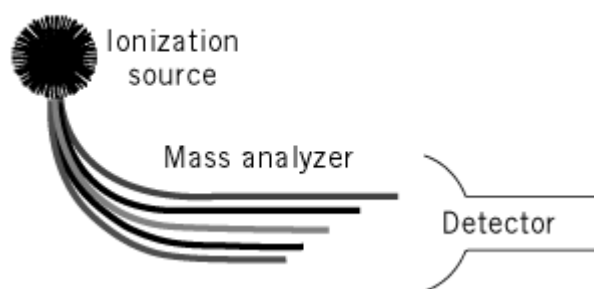
### Suggestion for Further Reading

G. B. Ruvken and F. M. Ausubel (1981) A general method for site-directed mutagenesis in prokaryotes. *Nature* **289**, 85–89.

### Mass Spectrometry

A *mass spectrometer* determines the molecular weight of chemical compounds by ionizing, separating, and measuring molecular ions according to their mass-to-charge ratio ( $m/z$ ). The ions are generated in the ionization source by inducing either the loss or the gain of a charge (eg, electron ejection, protonation, or deprotonation). Once the ions are formed in the gas phase, they can be electrostatically directed into a mass analyzer, separated according to mass, and finally detected. The result of ionization, ion separation, and detection is a mass spectrum that can provide molecular weight or even structural information (Fig. 1) (1).

**Figure 1.** The three basic components of mass spectrometers: the ionization source, the analyzer, and the detector. ESI and MALDI ionization sources have revolutionized the variety of molecules that can be analyzed and thus have broadened mass spectrometry applications.



An historically important goal in mass spectrometry has been to generate intact biomolecular ions in the gas phase to allow for their mass measurement. It has largely been the development of [electrospray ionization \(ESI\)](#) (2) and [matrix-assisted laser desorption/ionization \(MALDI\)](#) (3) sources that have efficiently accomplished this goal. As a result, mass spectrometry has become an integral part of biological research. Proteins, peptides, carbohydrates, oligonucleotides, natural products, and drug metabolites can now be analyzed routinely and examined structurally in picomole to femtomole amounts. Mass spectrometry can be used to measure the mass of biomolecules well over 200,000 Da, to provide sequence information on unknown peptides and proteins, and to detect noncovalent complexes, all with an accuracy in the measurement of molecular weight on the order of  $\pm 0.01\%$  (4).

While ESI and MALDI are fundamentally different ionization techniques, each with its respective advantages, they achieve essentially the same end result—the generation of gas phase ions via nondestructive vaporization and ionization. In both techniques, ionization typically occurs through proton addition or proton abstraction, to produce either  $[M + H]^+$  or  $[M - H]^-$  ions (where M is the molecule of interest).

## 1. Mass Analyzers

Mass analyzers separate ions according to their mass-to-charge ratio ( $m/z$ ) and are thus a crucial part of the mass spectrometer. There are numerous types and designs, yet only four analyzers have been adopted for ESI and MALDI (Table 1). ESI-MS and MALDI-MS commonly use *quadrupole mass analyzers* and *time-of-flight (TOF) mass analyzers*, respectively. The limited resolution offered by TOF mass analyzers, combined with adduct formation observed with MALDI-MS, results in accuracy on the order of 0.1% to a high of 0.01%, while ESI typically has an accuracy on the order of 0.01%. Both ESI and MALDI are now being coupled to higher-resolution mass analyzers, such as the ultrahigh-resolution ( $>10^4$ ) *Fourier-transform ion cyclotron resonance mass analyzer (FTMS)*,

with an accuracy of 0.001%. This analyzer is based on Fourier-transform mass spectrometry (FTMS). The result of the increased resolving power for ESI and MALDI mass spectrometers is an increase in accuracy for **biopolymer** analysis.

**Table 1. General Comparison of Quadrupole, Time-of-Flight, and FTMS Mass Analyzers**

| Quadrupole and Ion Trap  | Time-of-flight (TOF)                             | Fourier transform mass spectrometry(FTMS)  |
|--|--|--|
| Mass range (m/z ~4000)   | High mass range (unlimited)                      | Mass range (m/z ~10,000)   |
| Good resolution (~3000)  | Good resolution (~3000)                          | High resolution (>30,000)  |
| Relatively low cost (<\$200,000)   | Simple design, relatively low cost (<\$150,000)  | Superconducting magnet adds cost (>\$300,000)  |
| Well-suited for ESI  | Not well-suited for ESI                          | Well-suited for ESI  |
| Not well-suited for MALDI  | Well-suited for MALDI                            | Well-suited for MALDI  |
| Accuracy ~±0.01% (without internal standard)   | Accuracy ±0.2% to 0.01% (with internal standard) | Accuracy ~±0.0005% (without internal standard)   |
| Amenable to structural studies using tandem mass analysis. Ion trap especially well-suited for (MS <sup>n</sup> , typically n ≤ 4) | Not amenable to structural studies               | FTMS amenable to structural studies using multipleMS fragmentation experiments (MS <sup>n</sup> , typically n ≤ 4) |

Quadrupole and quadrupole *ion trap* mass analyzers have found utility in their capacity to interface having electrospray ionization, an interface having three primary advantages. First, quadrupoles and ion traps are tolerant of relatively high pressures (~5 × 10<sup>-5</sup> torr), which is a quality well-suited to electrospray ionization, since the ions are produced at atmospheric pressure. Secondly, quadrupoles and ion traps are now capable of routinely analyzing an *m/z* of up to 4000, which is useful because electrospray ionization of proteins and other biomolecules commonly produces a charge distribution below *m/z* = 4000. Finally, the relatively low cost of quadrupole mass spectrometers makes them attractive as electrospray analyzers. Considering these mutually beneficial features of electrospray and quadrupoles, it is not surprising that most of the successful commercial electrospray instruments thus far have been coupled with quadrupole mass analyzers.

A TOF analyzer is one of the simplest mass analyzing devices and is commonly used with MALDI ionization. Time-of-flight analysis is based on accelerating a set of ions to a detector with the same amount of energy. Because the ions have the same energy, yet different masses, the ions reach the detector at different times. The process is analogous to a pitcher throwing a golf ball and a basketball at a catcher with the same amount of energy. The golf ball will reach the catcher faster because it has a smaller mass and therefore a greater velocity. So it is with ions. The smaller ions reach the detector first because of their greater velocity, and the larger ions take longer; the analyzer is called a time-of-flight analyzer because the mass is determined by the ions' flight time through the analyzer.

FTMS is a valuable approach toward measuring ions because it offers high resolution (>30,000) and the ability to perform multiple collision events ( $MS^n$ , where  $n$  can be as high as 4). First introduced in 1974 by Comisarow and Marshall (5), the method of FTMS traps ions within a cell and then detects the image current of orbiting ions. Ions are distinguished according to their  $m/z$  by measuring their different orbiting frequencies. Coupled to ESI and to MALDI, FTMS has become an important research tool offering high accuracy with errors as low as  $\pm 0.001\%$ .

## 2. Conclusion

The historical achievement of measuring the masses of intact biomolecular ions has established MS as a valuable tool in molecular biology and biotechnology. It is the combination of attributes such as sensitivity, mass range, accuracy, and the ability to obtain structural information that has brought MS to the forefront of biotechnological research and which propels its further development as an important analytical tool for both chemists and biologists.

## Bibliography

1. G. Siuzdak (1996) *Mass Spectrometry for Biotechnology*, Academic Press, San Diego.
2. J. B. Fenn, M. Mann, C. K. Meng, S. F. Wong, and C. M. Whitehouse (1990) *Mass Spectrometry Rev.* **9**, 37–70.
3. M. Karas, D. Bachmann, U. Bahr, and F. Hillenkamp (1987) *Int. J. Mass Spectrom. Ion Proc.* **78**, 341–348.
4. G. Siuzdak (1994) *Proc. National Academy of Sciences of the U.S.A.* **91**, 11290–11297.
5. M. B. Comisarow and A. G. Marshall (1974) *Chem. Phys. Lett.* **25**, 282–283.

## Suggestions for Further Reading

6. G. Siuzdak (1996) *Mass Spectrometry for Biotechnology*, Academic Press, San Diego.
7. A. L. Burlingame and S. A. Carr, Eds. (1996) *Mass Spectrometry in the Biological Sciences*. Humana Press, Totowa, N.J.
8. R. A. W. Johnstone and M. E. Rose (1996) *Mass Spectrometry for Chemists and Biochemists*, 2nd Ed. Cambridge University Press, New York.
9. R. M. Caprioli, A. Malorni, and G. Sindona, eds. (1996) *Mass Spectrometry in Biomolecular Sciences* Kluwer Academic, Dordrecht; Boston.

## Master Chromosome

The notion that a copy of [mitochondrial DNA](#) might behave as a master [chromosome](#) to ensure transmission of identical copies of the mitochondrial [genome](#) to daughter cells was discussed early in the development of mitochondrial studies, but it has never been completely clarified. The nonrandom inheritance of multiple mitochondrial genomes to progeny cells raises problems that are still largely unsolved, although new clues will probably come from recent results concerning mitochondrial fusion, collision, motility.

In **yeast**, mitochondria can replicate the mitochondrial chromosome until 50–100 copies are present per cell. Mitochondria can then divide without DNA synthesis, eventually to contain 1–5 copies each of the mitochondrial chromosome. The population of these DNA molecules is homogeneous. A



heterogeneous population is present in **zygotes** after **conjugation**, but this is rapidly replaced by the homoplasmic situation after 10–20 generations. This nonrandom segregation had originally suggested the master copy hypothesis.

Recently very interesting fluorescence microscopic evidence has suggested the existence of a bias in the transmission of one of the parent's mitochondria to the first zygotic bud (1). The involvement of the [cytoskeleton](#) in the mitochondrial movements during the mitosis and meiosis, and hence probably in mitochondrial inheritance, has also been recently demonstrated (2). On the other hand, human cells are generally heteroplasmic and can bear varying proportions of different mitochondrial genomes. Hence the presence of inherited, or newly acquired, mutations results in mixed mitochondrial populations; this can result in partial deficiencies of mitochondrial functions and allows late-onset and characteristic tissue dependence of mitochondrial illnesses. However, as heteroplasmic cells undergo mitosis or meiosis, the proportion of mutant and **wild-type** DNA may fluctuate and eventually drift towards homoplasmy. Accumulation of mitochondrial mutations in somatic cells might play a role also in aging and in some neurodegenerative disorders (3).

In plants, the structures of mitochondrial genomes are much less clearly defined, and variable proportions of circular and linear molecules of subgenomic size are present. The latter probably arise from homologous [recombination](#) between large [inverted repeats](#) found on a large master chromosome. This master circle is very large, measuring about 570 kbp for the fertile cytoplasm of maize, for instance.

#### Bibliography

1. J. Nunnary et al. (1997) *Mol. Biol. Cell* **8**, 1233–1242.
2. K. H. Berger and M. P. Yaffe (1996) *Experientia* **52**, 1111–1116.
3. D. C. Wallace (1992) *Annu. Rev. Biochem* **61**, 1175–1212.

#### Maternal Control, Effect

[Transcription](#) from the [genome](#) of a [zygote](#) does not usually begin until several cleavage divisions have taken place. The early cleavage divisions, and often much of early embryogenesis, are accomplished by RNAs and proteins deposited in the unfertilized egg by the maternal genome during oogenesis. The early patterning of the embryo is also dependent on maternal positional information in the unfertilized egg. For example, both the anterioposterior and dorsoventral axes of *Drosophila* are determined by maternally acting genes (1-3). The extent of maternal control of early development has been most characterized in *Drosophila melanogaster* using mutations that alter development. Mutations that act in the mother to alter the development of progeny are called *maternal effect mutations*.

Maternal effect mutations that cause the death of all progeny are a special type of female sterile mutation, called *maternal effect lethal mutations*. The severity of the defects in the dead progeny can be independent of their genotype (nonrescuable by the paternal gene delivered by the sperm), or can be less severe in zygotes that receive a wild-type gene from the father (partially rescued). When the progeny are completely rescued by the wild-type paternally derived gene, the maternal-effect mutation is not classified as a female sterile. One example is the *cinnamon (cin)* mutation of *Drosophila* (4). Homozygous *cin* mutant progeny from homozygous mutant mothers die during embryogenesis. However, progeny of *cin* mutant mothers that receive a wild-type allele from the father survive. Most screens for maternal-effect mutations in *Drosophila* have not been designed to

identify mutations that are completely rescued by the paternally contributed allele, and the frequency of this type of gene is not known. Most screens would also miss maternal-effect mutations in which the progeny die later than embryogenesis. Because the zygotic genome becomes transcriptionally active relatively early in embryogenesis, few maternal-effect mutations are expected to cause death later than embryogenesis. An estimate of the proportion of the genome that can mutate to give maternal-effect lethal mutations can be made from the work of Schüpbach and Wieschaus (5). They recovered female-sterile mutations at about 8% of the frequency of lethal mutations. Of these female-sterile mutations, about 25% produced normal eggs that were fertilized but the embryos died during early development. This suggests that about 72 genes in *Drosophila* should mutate to maternal-effect lethality. As many of the female-sterile mutations are single alleles of genes that usually mutate to zygotic lethality (6), this estimate is probably somewhat low.

The study of maternal control using maternal-effect mutations is limited by the ability to recover mutations that are not lethal to the zygotes that carry them. This problem has been circumvented by the use of genetic **mosaics**. One method is to transplant germ cells from embryos homozygous for lethal mutations into wild-type embryos. These mutant germ cells are included in the developing ovary, and the resulting females have wild-type somatic cells, but mutant germ cells. This approach is limited by the difficulty and time involved in the transplantation techniques. [Mitotic recombination](#) in the germ cells of females heterozygous for recessive lethal mutations can also be used to produce homozygous mutant germ cells. The use of a dominant female-sterile mutation to prevent the formation of normal eggs from all of the germ cells except those homozygous for the mutation of interest has greatly facilitated this approach (7). Female germ cells heterozygous or homozygous for the *ovo*<sup>DI</sup> mutation are blocked early in oogenesis. Mitotic recombination in females heterozygous for this mutation produces germ cells that are homozygous for the wild-type allele (they lack the mutant allele), are no longer blocked in oogenesis, and can produce wild-type eggs. If the chromosome that carries the wild-type allele of *ovo* also carries a recessive lethal mutation, the germ cells that lose the mutant *ovo*<sup>DI</sup> allele also lose the wild-type allele for the recessive lethal mutation. Extensive experiments using this technique suggest that at least 70% of the zygotic lethal mutations in *Drosophila* have maternal effects (8).

### Bibliography

1. D. Morisato and K. V. Anderson (1995) *Annu. Rev. Genet.* **29**, 371–399.
2. C. Nüsslein-Volhard (1993) *Cancer* **71**, 3189–3193.
3. S. Roth, F. S. Neuman-Silberbert, G. Barcelo, and T. Schüpbach (1995) *Cell* **81**, 967–978.
4. B. S. Baker (1973) *Dev. Biol.* **33**, 429–440.
5. T. Schüpbach and E. Wieschaus (1989) *Genetics* **121**, 101–117.
6. N. Perrimon and A. P. Mahowald (1986) In *Gametogenesis and the Early Embryo* (J. Gall, ed.), Alan R. Liss, New York, pp. 221–235.
7. N. Perrimon, L. Engstrom, and A. P. Mahowald (1984) *Dev. Biol.* **105**, 404–414.
8. Reference 6, p. 225.

### Suggestions for Further Reading

9. N. Perrimon, L. Engstrom, and A. P. Mahowald (1989) *Genetics* **121**, 333–352.
10. T. B. Chou and N. Perrimon (1996) *Genetics* **144**, 1673–1679.
11. R. P. Ray and T. Schüpbach (1996) *Genes Dev.* **10**, 1711–1723.
12. S. Roth, F. S. Neuman-Silberberg, G. Barcelo, and T. Schüpbach (1995) *Cell* **81**, 967–978.

## Maternal Genetic Effects

The phenotype of the offspring of a genetic cross may be decided not by their own genotype, but by the genotype or the phenotype of the mother. Such *maternal effects* should be distinguished from the matrilineal inheritance of cytoplasmic organelles (see [Cytoplasmic Inheritance](#)); in this latter case the phenotype of the offspring is decided by its own genotype, which happens to coincide with that of the mother. Maternal effects are found in traits that are governed by nuclear genes, or by no genes.

### 1. Delayed Inheritance

Maternal effects governed by nuclear genes occur when the trait is decided very early in the development of an embryo by messenger RNA or proteins given by the mother to the unfertilized ovum. The phenotype is decided by the genotype of the mother. This pattern of inheritance is also called *predetermination*.

A well-known example is the determination of shell coiling in the snail *Limnaea peregra* (1). Shell coils, like DNA helices or staircases, may be twisted to the left or to the right; both phenotypes are found in *Limnaea* and other snails and are determined by two alleles of a single gene, with complete dominance of the allele for dextrality. Reciprocal crosses give different results; for example, the  $F_1$  hybrids are genetically homogeneous and will twist to the right or to the left, depending on whether the mother was the dextral or sinistral homozygote. The  $F_2$  segregants are genetically heterogeneous, but all will twist to the right. These results should be compared with those of [Mendelian inheritance](#).

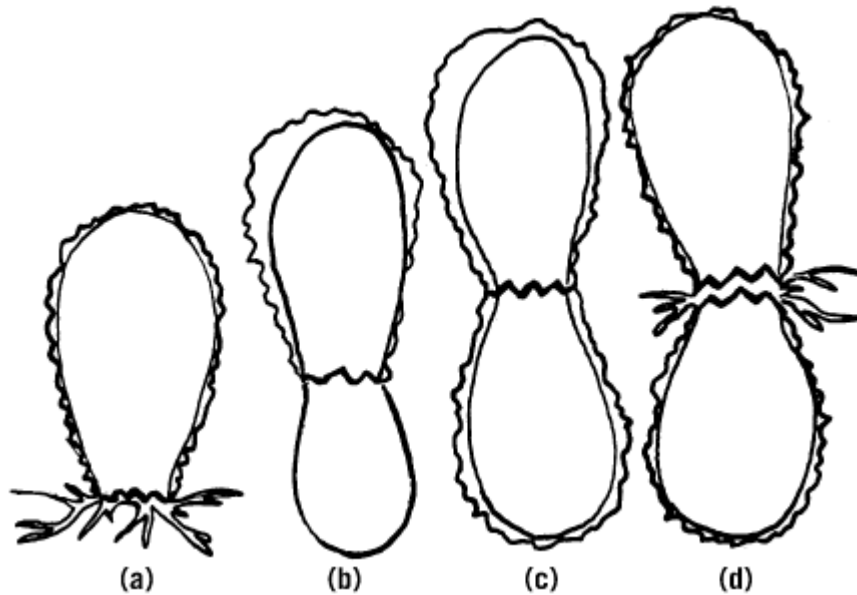
### 2. Persistent Modifications

Persistent modifications (*Dauermodifikationen*, Ref. 2) of gene expression may be induced by environmental factors, without change in the [genome](#), and may give rise to phenotypes maintained over many generations. The phenotype is usually transmitted through the mother and will fade out gradually in the absence of new stimuli, often after several generations. Persistent modifications are examples of **epigenetic effects**: A genotype may give rise to different phenotypes and express them persistently after reception of appropriate signals. Such epigenetic effects are responsible for cell differentiation in multicellular organisms.

### 3. Inherited Structural Traits Not Decided by Genes

Genomes cannot act without a preexisting cell. The structure of a zygote provides the starting conditions for the action of genes; this implies an amount of information that is difficult to estimate, but is probably vast. Some of this information may exist in nearly allelic alternatives not encoded in the genome. Thus, in *Paramecium aurelia*, the several thousand surface units that multiply by division are asymmetric and may be placed in one of two possible orientations relative to the cell axis. Both orientations exist in natural populations and breed true indefinitely. Surface units whose orientation is changed surgically maintain their new orientation and transmit it indefinitely (3). Another example of alternative phenotype not determined by the genes is represented in Figure 1. The number of teeth in the tests (a sort of shell) of some testate amoeba, like *Diffugia* and *Euglypha*, is inherited, but does not depend on the genes, as shown by changing their number mechanically (4).

**Figure 1.** Reproduction of the testate amoeba, *Euglypha*. The vegetative animal lives inside a test with a toothed rim (a). The animal grows, protrudes (b), and generates a new test (c), whose teeth are molded on those of the parental test. (From W. Schewiakoff.)



### Bibliography

1. A. E. Boycott and C. Diver (1923) Proc. R. Soc. **195**, 143.
2. V. Jollos (1921) Arch. Protistenkunde **43**, 1.
3. J. Beisson and T. M. Sonneborn (1965) Proc. Natl. Acad. Sci. USA **53**, 275–282.
4. H. S. Jennings (1937) J. Exp. Zool. **77**, 287–336.

### Maternal Inheritance

Maternal inheritance, or matrilinear inheritance, is the transmission of **genes** through the mother only. It is a common occurrence in the case of [cytoplasmic inheritance](#). The resemblance of all progeny to the mother in genetic reciprocal **crosses** is not, however, always due to maternal inheritance.

### Mating In Fungi

Mating is an essential step in the life cycle of sexually reproducing organisms. It provides a means to impose barriers on self-mating and to increase genetic variability within a population. In filamentous **fungi**, the mating type genes regulate not only sexual compatibility, but also morphogenesis and cell development and, at times, **pathogenesis**. Research into the molecular mechanisms governing mating processes of filamentous fungi impinges greatly on our understanding of many other processes central to the molecular biology of eukaryotes, such as gene regulation, [signal transduction](#), self/nonself recognition, and organelle movement.

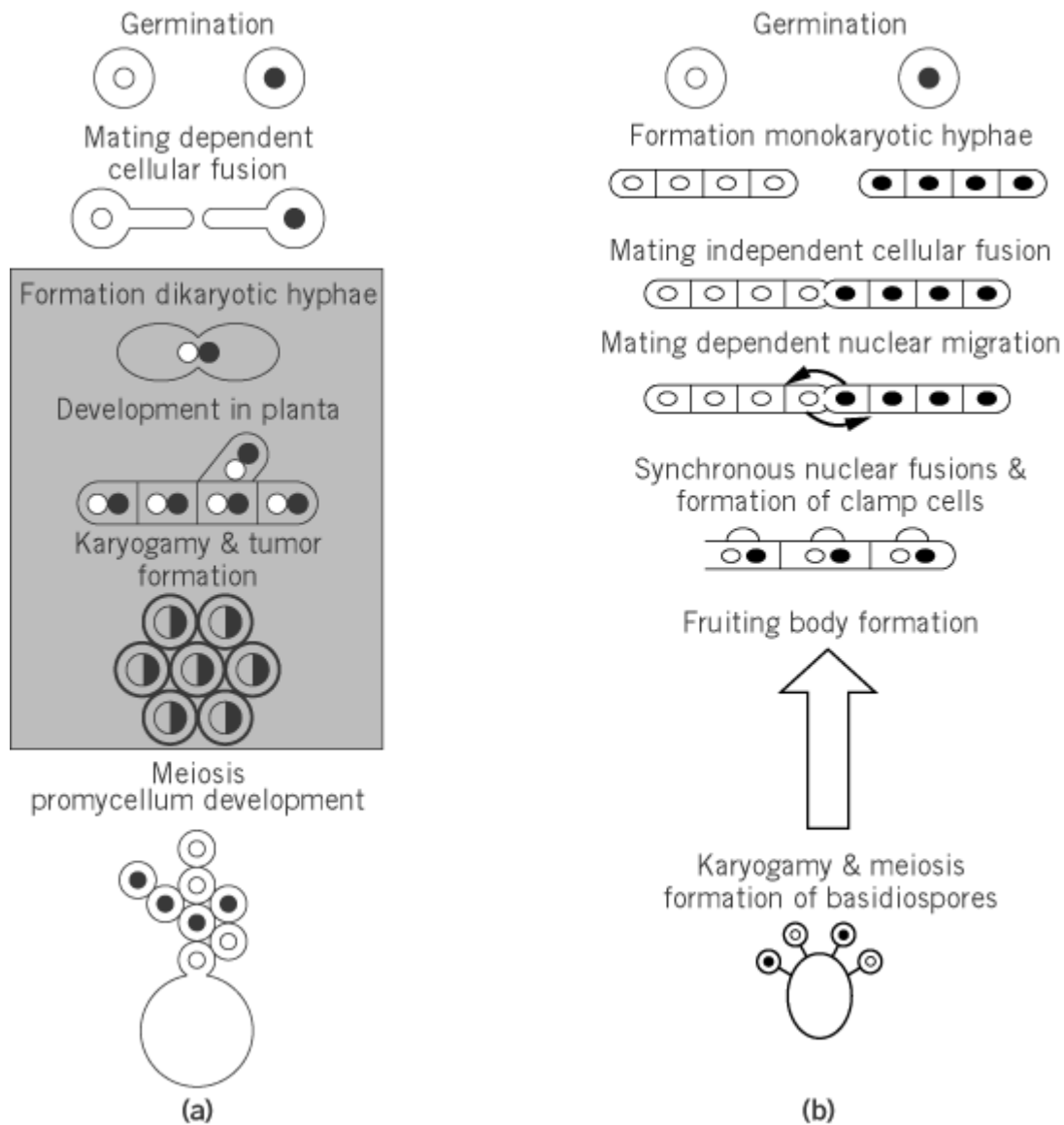
Mating type is described as either heterothallic (in which gametes come from parents of different mating type), homothallic (in which the strains are self-fertile), or, less frequently, pseudohomothallic strains that are self-fertile but in which nuclear fusion takes place between nonidentical nuclei. The switching system of *S. cerevisiae* is such an example: The mother cells switch mating type every generation so that direct descendants are now capable of partnering. No genetic variability is generated by such a system. Filamentous ascomycete species have been identified within the same species that encompass all three reproductive strategies.

The mating pathways in filamentous fungi have many features in common with those of the yeasts, *S. cerevisiae* and *S. pombe*, but they show a diversity that must reflect the unique and specialized environmental challenges faced by this widely disparate group of organisms. For example, filamentous fungi may differentiate into different vegetative forms, including asexual spores and/or infection structures. They may regulate the nuclear composition of their vegetative mycelium and produce unique sexual fruiting bodies that house the maturing meiotic spores. In order to understand this process, it is necessary to describe briefly the life cycles of several of the model organisms currently under intense study.

For the Ascomycetes, the best studied are, of course, the yeasts *S. cerevisiae* and *S. pombe*, both of which are described elsewhere in this volume. The filamentous fungi are generally represented by work on *N. crassa*, *M. grisea*, *P. anserina*, and *C. parasitica*. *N. crassa* life cycle presented in **Fungi** page is typical of this group of organisms.

The basidiomycete mating systems most often studied are *U. maydis*, the smut fungus of corn and *Coprinus cinereus* and *Schizophyllum commune*, both mushroom fungi. A brief description of their reproductive cycles are presented (see also Fig. 1): Haploid spores germinate under proper nutritional conditions to give rise to monokaryotic masses of hyphal filaments, much like that described for the ascomycetes. For *C. cinereus* and *S. commune* (both homobasidiomycete fungi), a chance encounter with another mycelium results in cellular fusion; however, nuclear fusion and subsequent development occurs only upon the determination of compatibility. For the hemibasidiomycetes (eg, *U. maydis*), conjugation tubes are extruded in response to secreted pheromones of the opposite mating type, much like the formation of shmoo in *S. cerevisiae*; cellular fusion occurs only with compatible mates.

**Figure 1.** Schematic representation of life cycles of the basidiomycetes: ( **a** ) *Ustilago maydis*, ( **b** ) *Coprinus cinereus*. In *U. maydis*, two mating type spores send out long projection tubes in response to pheromones prior to cellular fusion, thereby producing a dikaryotic infectious cell which is dependent upon plant infection in order to complete its life cycle. The dikaryotic hyphae invade and develop *in planta* whereupon plant signals result in the further development of tumors or galls on the corn. These tumors are the result of karyogamy and teliospore development. Meiosis results in the restoration of the haploid state. In *C. cinereus* two mating type spores germinate producing monokaryotic hyphae. A chance fusion of hyphae containing compatible mates allows reciprocal nuclear migrations through openings in septae called dolipores. Specialized clamp cells form at the apical cells in which a synchronous nuclear division and clamp cell septation results in reproduction of dikaryotic tissue. Appropriate stimuli induces fruiting body formation in which karyogamy and meiosis occur within the basidia located within the mushroom cap.



In the homobasidiomycetes, if a compatible mating has occurred, reciprocal nuclear migration occurs, followed by septation and the formation of a dikaryotic mycelia. The presence of two compatible nuclei within the tip cell of a growing hyphal mat, as well as additional environmental stimuli, signals further developmental processes. A unique outgrowth of the tip cell forms (called the clamp cell), followed by coordinated nuclear division of both nuclei, specific clamp cell septation and fusion with the subapical cell that results in both the tip cell and subapical cells having one nucleus from each parent. Fruiting body differentiation, followed by meiosis, restores the haploid content of the organism and completes the life cycle.

In *U. maydis*, initial cellular fusion is dependent upon having different mating specificities at the *a* locus. Following fusion at the tips of the conjugation tubes, dikaryotic hyphae are formed that are pathogenic on maize. These hyphae can grow only *in planta* and only if different alleles at the *b* mating locus are present in the paired nuclei. Thus compatibility is determined at two points in the life cycle—first, to initiate the mating cycle, and second, to complete sexual development and pathogenesis. Within the plant, tumors are formed that enclose black diploid teliospores, which ultimately burst. Meiosis occurs within the teliospore, forming a promycelial structure which then buds regenerating the haploid sporidia (1).

## 1. Mating Loci in Ascomycetes

All known heterothallic ascomycetes have a single-locus, two-allele mating system similar to that of the more familiar yeasts, *S. cerevisiae* and *S. pombe*. Mating type in these systems controls the initial fusion of compatible strains, as well as the subsequent formation of the dikaryotic hyphae that will ultimately form the zygotes. The mating type locus contains different DNA or idomorphs flanked by common sequences. In *N. crassa* there is 5.1 kbp of *A*-specific sequence and 3.2 kbp of *a*-specific sequence at the *mt* locus (2, 3); in *M. grisea* there is 2.5 kbp of *MAT1-1*-specific sequence and 3.5 kbp of *MAT1-2*-specific sequence (4); in *C. heterostrophus* there is 1.3 kbp of *MAT-1*-specific sequence and 1.2 kb of *MAT-2*-specific sequence at the mating loci (5). In contrast to the situation in the yeasts, there are no silent copies of mating-type sequences elsewhere in the genome of these fungi, thus preventing any possibility of a switching mechanism. Mating types are very stable in the heterothallic ascomycetes.

The sequences at the mating loci may encode one or more genes; the genes of alternate mating types are not related to one another by sequence. For *N. crassa*, the *mtA* locus contains three genes (*mtA-1*, *mtA-2*, *mtA-3*) while the *mta* locus contains one gene (*mta-1*) (6). In *P. anserina* the *mat*- locus contains three genes (*FMRI*, *SMR1*, *SMR2*) and the *mat+* locus only one (*FPR1*) (7). The mating loci of *C. heterostrophus* are much simpler; both contain only a single gene, one with homology to *Mata1p* of *S. cerevisiae* and the other to *Mat a-1p* of *N. crassa* (5).

The genes encoded by the mating type loci are likely to be regulatory proteins; sequence analysis suggests that some of the proteins contain potential DNA binding motifs. The *Mata-1* and *Mata-3* of *N. crassa*, *FPR1* and *SMR2* of *P. anserina*, and *Mat-2* of *C. heterostrophus* each have high mobility group (HMG) motifs and bind DNA *in vitro* (8). HMG proteins are known to bind to DNA in the minor groove and induce a bend in the DNA (9). Other known mating-related HMG box proteins include *Mat Mc* and *STE11* from *S. pombe*, and *Prf1* from *U. maydis*. The *Mata-1* (*N. crassa*) polypeptide, along with *FMRI* (*P. anserina*) and *MAT-1* (*C. heterostrophus*), are also potential transcriptional activators; they have sequence similarity to *Mata1p* of *S. cerevisiae*, which is known to bind to DNA in combination with *Mcm1p* and turn on specific genes (10). *In vitro* analysis of *Mata-1* of *N. crassa* demonstrated DNA binding while truncation of the final 188 amino acids completely abolished mating function (11). Interestingly, these three proteins have relatively acidic C-termini common to known transcriptional activation domains in other proteins. Furthermore, the C-terminal region of *FMRI* has been shown to interact with *SMR2* in a yeast two-hybrid system (12), although truncation of the carboxy-terminal 112 amino acids revealed only a minor defect in spore production and is apparently completely dispensable for mating activity (13). It is possible that this protein regulates a postfertilization step that was not assayed in these experiments. The involvement of mating type genes in postfertilization steps suggests a level of complexity in these systems that should provide new insights into interacting networks in eukaryotic organisms. While the complete picture is not yet available, it is satisfying that many of the same regulatory processes revealed in the yeasts are apparently at work in these fungi.

Only a handful of the targets that are regulated by these putative regulatory proteins have been identified. Clearly, mating specific pheromones and their corresponding receptors should be among the targets. Two mating-specific pheromones, *Vir1* and *Vir2*, have been identified in *C. parasitica* (14, 15) that have structural features in common with other known fungal mating pheromones and are down-regulated during the mycovirus infection. In addition, a number of genes that are transcriptionally up-regulated in a mating-specific manner in some of these well-studied fungi have been identified by classical genetic screens. It is expected that numerous mating-specific targets will be identified in the near future. Some will be similar between the different fungi, while others will be unique, depending upon the different developmental strategies the fungi employ following a successful mating. The methods through which this will be accomplished will likely be through the use of differential cDNA analysis, transcriptional/translational insertional libraries, and eventually the use of chip technologies which will be able to distinguish transcriptionally active genes under a wide variety of conditions.

## 2. Mating Loci in Basidiomycetes

This topic is reviewed in Refs. [16](#) and [17](#). The mating type in the basidiomycetes, while essentially built upon the same framework as in the ascomycetes, is much more complex. Instead of a one-locus, two-allele system, there is a two-locus, many (many)-allele system that provides an astounding number of possible mate choices. Furthermore, sexual development and pathogenicity in this group of fungi are intimately intertwined with one another. The best-studied model organisms in the basidiomycetes with regard to mating type are *U. maydis*, *C. cinereus*, and *S. commune*.

Mating is determined by genes at two mating type loci: *a* and *b* for *U. maydis* and *A* and *B* for the mushroom fungi. Successful mating occurs when the two partners contribute different alleles of genes at both loci (eg,  $A1B2 \times A2B3$  or  $a1b1 \times a2b2$ ). Because multiple alleles segregate at two loci, resulting in four possible combinations of mating types in the sexual progeny, these fungi are called tetrapolar whereas the ascomycetes such as *S. cerevisiae* are bipolar (only two mating types are segregated during meiosis).

In the smut fungi (*U. maydis* model) there are two different alleles of the *a* locus (*a1* and *a2*) which encode the pheromones and pheromone receptors. These peptides in turn control the initial fusion of the compatible cells and maintenance of the heterokaryon. The *b* locus of *U. maydis* encodes two divergently transcribed genes designated *bE* and *bW* that contain typical homeodomain motifs. There are an astounding 25 alleles of the *b* locus such that each of the genes is referred to as *bE1* and *bW1*, *bE2* and *bW2*, and so on. These genes regulate the expression of other genes required for filamentous growth and pathogenesis of the dikaryon. The redundancy of this mating locus was revealed in a series of deletion experiments in which it was shown that a single pair of genes, one *bE1* and one *bW2*, may be combined in a heterokaryon to produce a functionally active mating ([18](#)). With two allele specificities at *a* and 25 specificities at *b*, there are 50 possible mating types in *U. maydis*.

In *C. cinereus* the situation is a bit more complex: Each of the two mating loci have two subloci designated *Aa* and *Ab*, *Ba* and *Bb*. It is the *A* locus that encodes the homeodomain-containing regulatory proteins while the *B* locus encodes the pheromones and receptors. Although not all allele specificities are identical in structure, it is believed that the archetypical *A* locus is structured with a pair of divergently transcribed homeodomain proteins at the *a* sublocus and two pairs at the *b* sublocus. There are 160 different specificities at the *A* locus of *C. cinereus*. The *B* locus spans ~17 kb that contain three subfamilies of genes, with each subfamily containing two pheromone genes and one receptor. Each “set” of pheromone/receptor genes is functionally independent. A compatible mating occurs if only one of the subfamilies contains different allelic versions of its set of genes. Interestingly, the pheromones within a given subfamily can only stimulate receptors from the same subfamily (but clearly not its own receptor, nor other receptors from outside its subfamily). Thus, in order to generate the estimated 79 *B* mating specificities that have been described for *C. cinereus*, only four or five alleles are actually required using this modular system. Taken together, the combination of *A* and *B* loci may generate as many as 12,000 different mating types for *C. cinereus*.

For *S. commune* the nine *A* loci contain a single pair of divergently transcribed homeodomain proteins designated *Y* and *Z*, the exception being *Aa1*, which contains only a single transcript designated *Y1*. The *Ba* and *Bb* loci of *S. commune* each contain four genes: three encoding pheromone precursors and one gene encoding the receptor. Taken together, *S. commune* can accommodate 288 *A*-factor specificities (9 at *Aa*, 32 at *Ab*) and 81 *B*-factor specificities (9 at *Ba* and 9 at *Bb*) with a grand total of over 20,000 separate mating types ([19](#))!

The mating loci regulate everything from pheromone signaling, to fusion of the compatible cells or nuclei, formation and maintenance of the dikaryon, and finally morphogenesis of the fruiting body. For the mushroom fungi, the *A* locus encodes the regulatory homeodomain proteins while the *B* locus encodes the pheromones and pheromone receptors. For the smut fungi it is the *B* locus that encodes the pheromone/receptor cassette and the *A* locus that encodes the regulatory proteins.



Regardless of the apparent “mix-up” in the naming schemes, the paradigm for recognition and morphological development of the sexual fruiting body for these fungi are remarkably similar.

A pheromone response factor, *PRF1*, was identified by degenerate polymerase chain reaction (PCR) in *U. maydis* with similarity to *ROX1* of *S. cerevisiae*, a known transcription factor. *Prf1* deletion mutants are sterile and defective in their production of pheromone. Transcription of *PRF1* is itself induced 20-fold by pheromone stimulation, and deoxyribonuclease I [footprinting](#) reveals that the *Prf1* binds to the PRE elements found upstream of both *a1* and *a2* sequences ([20](#)).

The regulatory proteins of the mating loci encode two subunits of a heterodimeric protein containing a typical DNA-binding homeobox motif. Homeobox proteins are ubiquitous in eukaryotic organisms and are known to play critical roles in development. In fungi, these proteins are essential in self/nonself recognition and are required for the formation and maintenance of the dikaryon and development of the fruiting body. Importantly, **homeobox** domains have been classified into two different groups designated, HD1 and HD2. The HD1 class includes the *S. cerevisiae a2* gene, the *U. maydis bE* gene, the *S. commune Z* genes, and the *C. cinereus a-d1* genes. The HD2 class includes the *S. cerevisiae a1* gene, the *U. maydis bW* genes, the *S. commune Y* genes, and the *C. cinereus a-d2* genes. In compatible matings, each pair of HD proteins (one from the HD1 class and one from the HD2 class) forms a heterodimer that is functionally active when encoded only by genes from different specificities ([21-23](#)). Thus mating compatibility is now defined in terms of molecular recognition between two peptides.

The homeobox heterodimer that is formed is most likely a [transcription factor](#) whose function is to turn on a series of downstream targets that activates a new developmental pathway. Therefore, in order to prevent the constitutive activation of this pathway in an unmated organism, an HD1 protein must dimerize with an HD2 protein from a different specificity. In its simplest form, compatibility between mating partners occurs when an HD1 protein from one parent dimerizes with an HD2 protein from the other parent.

By analysis of chimeras between HD1 and HD2 proteins, the region responsible for discrimination between alleles of the same organism has been narrowed to the amino terminus of these proteins ([24](#)). Furthermore, in some instances, single amino acid substitutions have been shown to alter protein specificities and turn an otherwise incompatible reaction into a compatible one ([16](#), [22](#)). This model is further supported by the use of the yeast [two-hybrid system](#), where it was shown that dimerization activates transcription only when nonself combinations were assayed ([22](#)). Furthermore, mutations in one peptide that permitted dimerization between two ordinarily noncompatible proteins also permitted filamentous development in the same “illegitimate” combination. Thus self–nonself recognition and the ability to dimerize are one and the same. In *C. cinereus*, where there are multiple copies of the homeodomain proteins, specificity includes members of the same family as well as a functional HD1/HD2 dimer (*A-a1* is compatible with *A-a2* but not with *A-b2* or vice versa). Further analysis of the chimeric alleles suggested that the specificity is determined at the borders of the *N*-terminal specificity regions with the idea that these regions contain key amino acids that act as “antidimerization” points preventing dimerization of self peptides. It remains to be determined which amino acid residues are key players in this molecular recognition paradigm ([25](#)).

The targets of this transcriptional recognition system are largely unknown; yet, because many postfertilization steps are dependent upon compatible matings, an extensive alteration in the organism's gene activity is expected. *A*-regulated events include nuclear pairing and division, hook cell formation, and septation: *B* genes regulate nuclear migration and the fusion of the hook cell with the subapical cell. The pheromones and their receptors are likely candidates of this putative heterodimer transcriptional activation. They may be found at the *B* locus of the smut fungi and the *A* locus of the mushroom fungi. In *U. maydis*, the *B* locus encodes one pheromone and one receptor pair, either *mfa1* and *pra1* or *mfa2* and *pra2* whose expression is required for initial cellular fusion and proper maintenance and growth of the filamentous dikaryon ([26](#)). In *S. commune* and *C. cinereus*

the pheromones function after fusion has already taken place. The *B* locus of *S. commune* encodes three pheromone genes (*bap1*, *bap2*, and *bap3*) and a receptor (*bar1*), while the *C. cinereus* *B* locus contains six pheromone and three receptor genes. If fusion occurs between different specificities of compatible monokaryons, secreted pheromones trigger subsequent events that signal septal breakdown, thus permitting passage of the fertilizing nuclei.

### 3. Signal Transduction and the Mating Response

The mating response in fungi is controlled by a [signal transduction](#) pathway that is well-defined in *S. cerevisiae* and *S. pombe*. Upon stimulation by the appropriate mating pheromone, a [heterotrimeric G protein](#) bound to the intracellular loops of the receptor activates a [MAP kinase](#) cascade. This results in the activation of a transcription factor (Ste12p in yeast) which, in turn, binds to specific upstream sequences of a number of target genes whose transcription is altered. A Ga protein homologue from *N. crassa*, *gna-1* and a similar protein, *cpg-1* from *C. parasitica*, have been identified that lead to sterility when deleted *in vivo* ([27](#), [28](#)). Other homologues of the *S. cerevisiae* mating type signal transduction pathway have been identified in ascomycetes, either by sequence mining, by genetic complementation, or by the use of degenerate **PCR** to identify proteins of similar structure (see the supplementary materials section of the Web site at <http://www.AnnualReviews.org> from Ref. [29](#) which details signal transduction and regulatory components). Some clearly function in mating (mutations or deletions lead to sterility), while others such as *PMK1* from *M. grisea* appear to have little or no effect in the mating pathway and may be part of other parallel signal transduction pathways. It is likely, as research into mating signal transduction pathway is continued, that homologues of expected kinases, transcription factors, receptors, and **transporters** will all be identified in the filamentous fungi.

### Bibliography

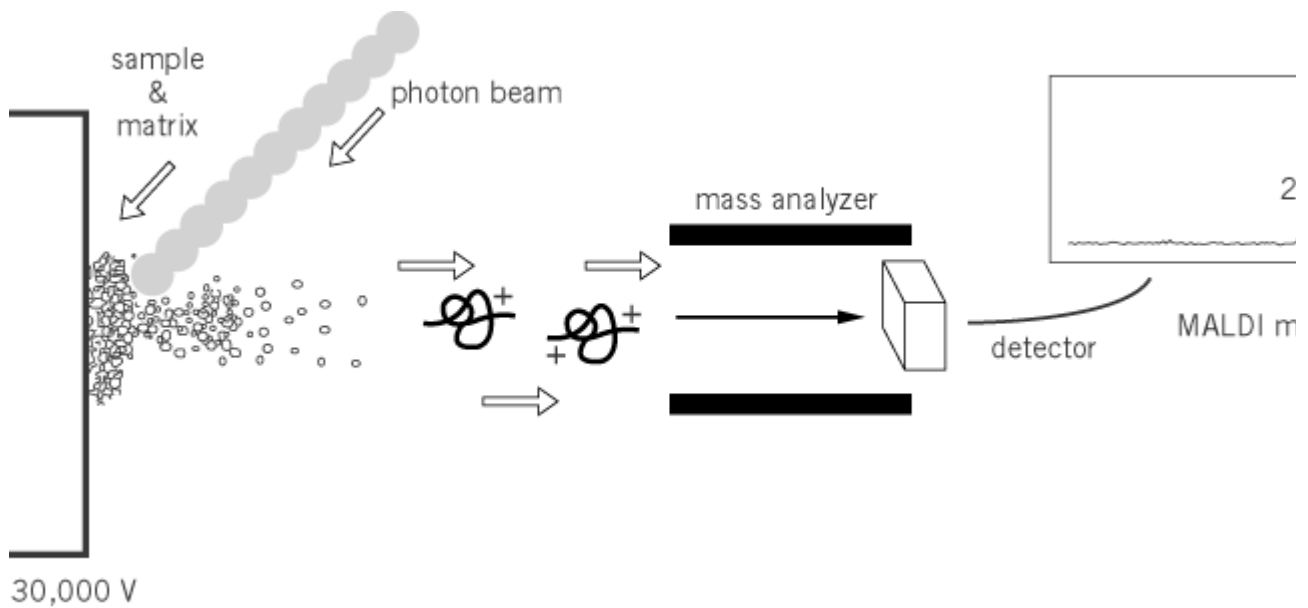
1. F. Banuett (1992) Trends Genet. **8**, 174–179.
2. N. L. Glass, J. Grotelueschen, and R. L. Metzberg (1990) Proc. Natl. Acad. Sci. USA **87**, 4912–4916.
3. C. Staben and C. Yanofsky (1990) Proc. Natl. Acad. Sci. USA **87**, 4917–4921.
4. S. C. Kang, F. G. Chumley, and B. Valent (1994) Genetics **138**, 289–296.
5. B. G. Turgeon et al. (1993) Mol. Gen. Genet. **238**, 270–284.
6. S. Chang and C. Staben (1994) Genetics **138**, 75–81.
7. R. Debuchy and E. Coppin (1992) Mol. Gen. Genet. **233**, 113–121.
8. J. W. Kronstad (1997) Trends Plant Sci. **2**, 193–199.
9. R. Grosschedl, K. Giese, and J. Pagel (1994) Trends Genet. **10**, 94–100.
10. E. J. Grayhack (1992) Mol. Cell. Biol. **12**, 3573–3592.
11. M. L. Phillely and C. Staben (1994) Genetics **137**, 715–722.
12. D. Zickler, S. Arnaise, C. E., R. Debuchy, and M. Picard (1995) Genetics **140**, 493–503.
13. E. Coppin, R. Debuchy, S. Arnaise, and M. Picard (1997) Microbiol. Mol. Biol. Rev. **61**, 411–428.
14. L. Zhang, A. C. L. Churchill, P. Kazmierczak, D. H. Kim, and N. K. Van Alfen (1993) Mol. Cell. Biol. **13**, 7782–7792.
15. P. Kazmierczak, P. Pfeiffer, L. Zhang, and N. K. Van Alfen (1996) J. Virol. **70**, 1137–1142.
16. R. Kahmann, T. Romeis, M. Bolker, and J. Kamper, (1995) Curr. Opin. Genet. Dev. **5**, 559–564.
17. L. A. Casselton and N. S. Olesnický (1998) Microbiol. Mol. Biol. Rev. **62**, 55–70.
18. B. Gillissen et al. (1992) Cell **68**, 647–657.
19. L. Casselton and U. Kues (1994) "Mating-Type Genes in Homobasidiomycetes". In *The Mycota*, (J. W. Wessels and F. Meinhardt, eds.), Vol. **1**, Springer-Verlag, Berlin.

20. H. A. Hartmann, R. Kahmann, and M. Bolker, (1996) *EMBO J.* **15**, 1632–1641.
21. A. H. Banham et al. (1995) *Plant Cell* **7**, 773–783.
22. J. Kamper, M. Reichmann, T. Romeis, M. Bolker, and R. Kahmann (1995) *Cell* **81**, 73–83.
23. Y. Magae, C. Novotny, and R. Ullrich (1995) *Biochem. Biophys. Res. Commun.* **211**, 1071–1076.
24. A. R. Yee and J. W. Kronstad (1993) *Proc. Natl. Acad. Sci. USA* **90**, 664–668.
25. R. Kahman and M. Bolker (1996) *Cell* **85**, 145–148.
26. T. Spellig, B. M., F. Lottspeich, R. W. Frank, and R. Kahmann (1994) *EMBO J.* **13**, 1620–1627.
27. F. Ivey, P. Hodge, G. E. Turner, and K. A. Borkovich (1996) *Mol. Cell Biol.* **7**, 1283–1297.
28. S. Gao and D. L. Nuss, *Proc. Natl. Acad. Sci. USA* (1996) **93**, 14122–14127.
29. J. W. Kronstad and C. W. Staben (1997) *Annu. Rev. Genet.* **31**, 245–276.

### **Matrix-Assisted Laser Desorption/Ionization**

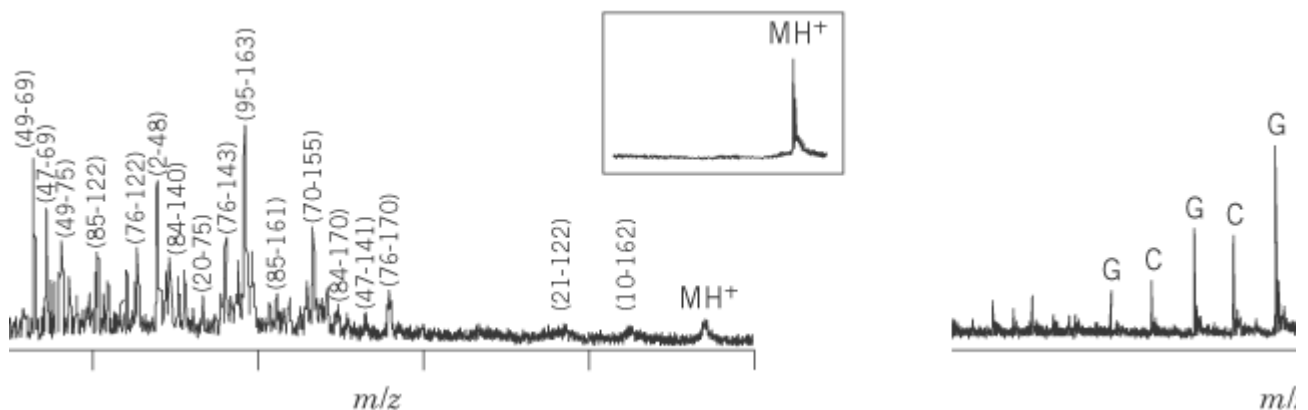
Matrix-assisted laser desorption/ionization (MALDI) is a method for producing ionic forms of molecules for analysis by [mass spectrometry](#) (MS). Gas phase ions are generated by the laser vaporization of a mixture of the molecule of interest, the analyte, in a solid *matrix* in which the matrix (usually a small crystalline organic compound) strongly absorbs the laser radiation and acts as a receptacle for energy deposition (Fig. 1) (1). This concentrated energy deposition results in the vaporization and ionization of both the matrix and the analyte ions, which contain relatively few charges. Low-molecular-weight analytes (<20 kDa) are typically ionized with only one or two charges, while larger analytes can accumulate as many as 3 to 5 charges, depending on the specific desorption conditions (ie the type of matrix material and the laser power). The relatively low number of charge states observed in MALDI makes the technique especially well suited for the analysis of multicomponent mixtures, as individual components can be easily identified by the mass spectral signal generated from their 1+ charge state.

**Figure 1.** The matrix-assisted laser desorption/ionization process.



The other commonly used ionization technique is [electrospray ionization](#) (ESI). Because ESI and MALDI are different ionization techniques, they offer different capabilities that in many cases are complementary (Table 1). For example, since MALDI is especially well suited for analyses like complex mixtures of peptides, a protein can often be identified by MS analysis of its constituent peptides produced by either chemical or enzymatic treatment of the sample (Fig. 2). This information can be especially useful when used in conjunction with protein databases or in identifying [post-translational modifications](#) (2).

**Figure 2.** Examples of data generated on a MALDI-MS. Proteins (left, inset) typically produce positive singly or doubly (inset) also produce singly and doubly charged positive ions. Proteolytic digest of the proteins (left) offers information ab identify the protein when used in conjunction with a protein database searching program. Exonuclease digestion of DNA



**Table 1. Capabilities, Limitations, and Recent Improvements of ESI-MS and MALDI-MS**

ESI Mass Spectrometry

MALDI Mass Spectrometry

|   |   |
|---|---|
| Routine femtomole to picomole sensitivity; attomolesensitivity possible                                   | Routine femtomole to picomole sensitivity; attomolesensitivity possible                               |
| Accuracy $\sim\pm 0.01\%$   | Accuracy $\pm 0.2\%$ to $0.01\%$ (with internal standard)   |
| Protein analysis to $\sim 70,000$ Da  | Protein analysis to $\sim 300,000$ Da   |
| DNA analysis to $\sim 100$ bases  | DNA analysis to $\sim 200$ bases  |
| Relatively intolerant of millimolar salt concentrations and mixture analysis                              | Tolerant of millimolar salt concentrations and mixture analysis                                       |
| HPLC-MS compatible  | Not HPLC-MS compatible  |
| Capable of observing noncovalent complexes directly from aqueous solutions                                | Not capable of routinely observing noncovalent complexes  |
| Amenable to structural studies using tandem mass analysis (see <a href="#">Tandem Mass Spectrometry</a> ) | Amenable to structural studies using enzyme digestion (see <a href="#">Tandem Mass Spectrometry</a> ) |

---

## Bibliography

1. F. Hillenkamp, M. Karas, R. C. Beavis, and B. T. Chait (1991) *Anal. Chem.* **63**, A1193–A1202.
2. E. J. Zaluzec, D. A. Gage, and J. T. Watson (1995) *Protein Expression and Purification* **6**, 109–123.

## Matrix Attachment Regions

Matrix attachment regions (MARS) are specific elements of **DNA** of the [genome](#) that associate with or attach to the [nuclear matrix](#) during the **interphase**. MARS are rich (about 70%) in adenine (A) and thymidine (T) bases, but no **consensus sequence** has yet been defined. The lengths of MARS range from a minimum of 200 to 350 base pairs (bp) to several thousand ([1](#)). Higher order [chromatin](#) structure involves discrete, and topologically independent, looped domains of 5 to 200 kbp within both mitotic [chromosomes](#) ([2](#)) and dispersed interphase chromatin ([3](#)). MARS form the bases of these looped domains and anchor the loops to the mitotic chromosome scaffolds or to the interphase nuclear matrix. In mitotic chromosomes, these elements have traditionally been referred to as SARS for [scaffold attachment regions](#) ([4](#)). In the case of the interphase matrix, these AT-rich elements are usually called MARS ([5](#)).

The chromatin loop domains are anchored to the matrix by [proteins](#) that preferentially bind AT-rich sequences. Among the proteins characterized are (1) the special [AT-rich binding protein 1](#), SATB1, that is expressed predominantly in thymocytes ([6](#), [7](#)); (2) the 95-kDa [attachment region binding protein](#), (ARBP) ([8](#)); (3) the A and B **lamins** of the nuclear envelope ([9](#), [10](#)), (4) the nuclear [intermediate filament](#)-type protein called NuMA, and (5) the **hnRNP U** protein. Most interestingly, **topoisomerase II** (ScI of Lewis and Laemmli, [11](#)) associates with both interphase MARS and mitotic SARS. Topoisomerase II may relieve torsional stress generated by [transcription](#) within

individual looped domains during the interphase (as it binds to MARS), and the enzyme may serve as a structural protein within the chromosome scaffold during mitosis (as it binds SARS).

MARS often lie close to sequences involved in **gene regulation**, such as **enhancers**. Because of this proximity, MARS have been implicated in regulating the expression of several genes including (1) the *Drosophila* [histone](#) gene cluster, (2) the human [interferon b](#) gene, (3) the [immunoglobulin m](#) heavy chain gene, (4) the human [b-globin](#) gene, and (5) the chicken [lysozyme](#) gene, to name just a few. When cells are stably transfected with genes linked to MAR elements, transcription of these genes is enhanced independently of the chromosomal integration position, but the expression depends on the number of MAR elements present. The AT-rich sequence of MARS permits base pair unwinding, which is thought to relieve negative superhelical strain introduced into the looped chromatin domain (7). AT-rich core motifs within MARS (eg, the ATATAT motif within a 300-bp MAR element downstream of the immunoglobulin heavy-chain gene enhancer) serve as the unwinding elements. The greater the number of negative **supercoils** generated, the more the motif unwinds. The SATB1 protein mentioned before preferentially binds the unwound segments of MAR elements. In fact, SATB1 does not bind the bases, but rather the particular configuration of the sugar-phosphate backbone of the DNA when the AT-rich segment is unwound (6). The current supposition is that the unwinding ability of the MARS permits them to bind tightly to the nuclear matrix.

### Bibliography

1. S. M. Gasser and U. K. Laemmli (1986) *Cell* **46**, 521–530.
2. J. R. Paulson and U. K. Laemmli (1977) *Cell* **12**, 817–828.
3. B. Vogelstein, D. M. Pardoll, and D. S. Coffey (1980) *Cell* **22**, 79–85.
4. J. Mirkovitch, M. E. Mirault, and U. K. Laemmli (1984) *Cell* **39**, 223–232.
5. P. N. Cockerill and W. T. Garrard (1986) *Cell* **44**, 273–282.
6. L. A. Dickinson, T. Joh, Y. Kowhi, and T. Kowhi-Shigematsu (1992) *Cell* **70**, 631–645.
7. T. Kohwi-Shigematsu and Y. Kohwi (1996) In *Nuclear Structure and Gene Expression* (R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 111–144.
8. J. P. von Kries, H. Buhrmester, and W. H. Strätling (1991) *Cell* **64**, 123–135.
9. M. E. E. Ludérus, A. de Graaf, E. Mattia, J. L. den Blaauwen, M. A. Grande, L. de Jong, and R. van Driel (1992) *Cell* **70**, 949–959.
10. M. E. E. Ludérus, J. L. den Blaauwen, O. J. B. de Smit, D. A. Compton, and R. van Driel (1994) *Mol. Cell Biol.* **14**, 6297–6305.
11. C. D. Lewis and U. K. Laemmli (1982) *Cell* **29**, 171–181.

### Suggestions for Further Reading

12. S. M. Gasser and U. K. Laemmli (1987) A glimpse at chromosomal order, *Trends Genet.* **3**, 16–22.
13. R. Berezney and K. W. Jeon, (eds.) (1995) "Structural and functional organization of the nuclear matrix", In *International Review of Cytology, A Survey of Cell Biology*, Vols. **162A** and **162B**, Academic Press, San Diego. Compilations of excellent reviews on matrix attachment regions and the nuclear matrix.

### Maturation Promoting Factor

In fewer than twenty years, Maturation Promoting Factor (MPF) has evolved from an obscure developmental factor, to the master regulator of the G2-to-mitosis transition in all eukaryotic cells, to a somewhat historical status remembered most fondly by an older generation of cell cycle researchers. Transfer of cytoplasm from amphibian eggs into immature oocytes causes the oocytes to proceed through meiosis and to mature into eggs. The responsible substance was termed MPF and was subsequently found to function during M-phases in all eukaryotic cells. The identification of the protein subunits of MPF in the late 1980's led to a three-way convergence of work on MPF with genetic analyses of the cell cycle in yeasts and studies of cyclins in marine invertebrates. This unification fueled the explosive growth in our understanding of the biochemical mechanisms controlling cell cycle transitions.

MPF was named and first characterized in detail during studies by Masui and Markert (1) on the role of cytoplasmic factors during cell differentiation. Their work focused on the maturation of *Rana pipiens* oocytes, although most later studies of MPF used oocytes of *Xenopus laevis*. These oocytes are naturally arrested in G2/early prophase at the start of meiosis. Secretion of progesterone by the follicle cells surrounding the oocyte initiates a signal-transduction cascade in the cytoplasm leading to progression through meiosis and arrest at the second meiotic metaphase. Fertilization of this egg triggers release from the metaphase arrest and initiates the mitotic cell cycles of the embryo. Maturation is generally scored by germinal vesicle breakdown (GVBD). The germinal vesicle is the gigantic nucleus present in oocytes whose dissolution results in displacement of pigment from the animal pole of the oocyte, resulting in the formation of a large "white spot". Injection of a number of agents can cause maturation. For example, any treatment that reduces the activity of the cAMP-dependent protein kinase (including decreases in the concentration of cAMP or injection of the regulatory (R) subunit of the kinase) leads to maturation. However, the induction of maturation by such agents as well as by progesterone requires ongoing protein synthesis. In contrast, maturation induced by injection of egg cytoplasm (containing MPF) proceeds even in the absence of protein synthesis, indicating that MPF acts late in the cascade leading from progesterone to maturation. Finally, MPF activity is "amplified" during maturation in that the injection of a small amount of MPF into an oocyte produces an egg with high levels of MPF that can be serially transferred back into oocytes. The observation that amplification occurs even in the absence of protein synthesis (2) indicates that a precursor form of MPF is stored in the oocyte.

Masui and Markert identified MPF as a developmental factor present in amphibian eggs but absent from oocytes. Subsequent more detailed studies showed that the presence of MPF activity actually follows the stages of the cell cycle, rather than the developmental stage of the egg. It is now appreciated that MPF activity—assayed as the ability of a sample to induce maturation following injection into the cytoplasm of an oocyte—is present whenever a cell is in an M phase, either meiotic or mitotic. Gerhart *et al.* showed that MPF activity cycles during *Xenopus* oocyte maturation (2). It rises as the oocyte enters the first meiotic metaphase, declines at the first meiotic division, and rises again at second meiotic metaphase. MPF was originally detected in eggs since they are stably arrested in metaphase, whereas first meiotic metaphase is only a transient state. Importantly, similar cycles of MPF activity occur during the mitotic divisions of frog embryos (2) and MPF activity is present in mitotic extracts from organisms as diverse as yeast and man (3-5), suggesting that MPF is a universal regulator of entry into M-phase.

The beginning of the molecular understanding of MPF dates to its partial purification in 1980 (6). Though the purification factor was a modest 20-30-fold, the simple demonstration of its feasibility and the rigorous quantitation of the results helped transform MPF from a "fuzzy factor" into a conventional enzyme. Purification to near homogeneity required another 8 years and indicated that MPF was composed of two proteins (7). In short order these were identified as the *Xenopus* homologs of Cdc2 and cyclin B (8-11) (see [Cell Cycle](#) and [Cyclins](#)). *cdc2*<sup>+</sup> was identified as the master regulator of the G2-to-mitosis transition during genetic studies in *Schizosaccharomyces pombe*, and its sequence indicated that it was a protein kinase. Cdc2 is the catalytic subunit of MPF,

and it is present throughout the cell cycle. Because Cdc2 is a protein kinase, contemporary assays of Cdc2 activity simply follow its ability to phosphorylate proteins such as histone H1, rather than its ability to induce oocyte maturation following microinjection. Cyclins were first identified in marine invertebrates as proteins that accumulate during interphase and that are abruptly degraded at the end of mitosis. Cyclin B is the regulatory subunit of MPF. An inactive form of the Cdc2-cyclin B complex forms the MPF precursor found in *Xenopus* oocytes.

The molecular identification of the components of MPF thus united the three major strands of cell cycle studies: MPF, yeast genetics, and cyclins. Although MPF per se is not extensively studied anymore, it has helped spawn a huge growth industry in cell cycle research. Subsequent studies have identified large families of proteins related to Cdc2 and cyclin B that regulate other cell cycle transitions (see [Cyclins](#)). These protein kinases are in turn regulated by multiple mechanisms, including phosphorylations and protein-protein interactions. So much of this work had its humble origins in simple studies of the cytoplasmic control of frog oocyte maturation by MPF.

### Bibliography

1. Y. Masui and C. L. Markert (1971) *J. Exp. Zool.* **177**, 129–146.
2. J. Gerhart, M. Wu and M. Kirschner (1984) *J. Cell. Biol.* **98**, 1247–1255.
3. H. Weintraub, M. Buscaglia, M. Ferrez, S. Weiller, A. Boulet, F. Fabre, and E.-E. Baulieu (1982) *C. R. Acad. Sci (Paris) Ser. III* **295**, 787–790.
4. P.S. Sunkara D.A. Wright, and P.N. Rao (1979) *Proc. Natl. Acad. Sci. USA* **76**, 2799–2802.
5. B. Nelkin, C. Nichols, and B. Vogelstein (1980) *FEBS Lett.* **109**, 233–238.
6. M. Wu and J.C. Gerhart (1980) *Dev. Biol.* **79**, 465–477.
7. M.J. Lohka, M.K. Hayes, and J.L. Maller (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3009–3013.
8. W.G. Dunphy, L. Brizuela, D. Beach, and J. Newport (1988) *Cell* **54**, 423–431.
9. J. Gautier, C. Norbury, M. Lohka, P. Nurse, and J. Maller (1988) *Cell* **54**, 433–439.
10. J.C. Labbé, J.P. Capony, D. Caput, J.C. Cavadore, J. Derancourt, M. Kaghad, J.M. Lelias, A. Picard, and M. Dorée (1989) *EMBO J.* **8**, 3053–3058.
11. J. Gautier, J. Minshull, M. Lohka, M. Glotzer, T. Hunt, and J.L. Maller (1990) *Cell* **60**, 487–494.

### Suggestions for Further Reading

12. J.L. Maller (1980) Regulation of oocyte maturation. *Curr. Top. Cell Regul.* **16**, 271–311.
13. J.L. Maller, (1995) Maturation-promoting factor in the early days. *Trends Biochem. Sci.* **12**, 524–528.

### Medium-Chain Dehydrogenases/Reductases

The MDRs are a large family of dehydrogenase and reductase enzymes that include the classical forms of alcohol dehydrogenase. They are usually proteins with subunits containing ~350 residue subunits and often, but not always, a zinc atom at the active site. For examples, references, and further details, see [Alcohol Dehydrogenase \(ADH\)](#).



## Melittin and Membrane-Perturbing Toxins

A very large array of protein [toxins](#) act by altering the [membrane](#) permeability barrier (1). If this appears to be a common final goal, very different strategies are adopted to reach it. Many bacteria and animal toxins are [phospholipases](#) that degrade membrane phospholipids to smaller components (fatty acids, lysolecithins, etc.) endowed with [detergent](#) properties, which increase the membrane permeability to ions and lead to osmotic lysis (2-3). Melittin (which is about 50% of the dry weight of bee venom) and *Staphylococcus aureus* d-lysin are both surface-active peptides 26 residues long. They have no sequence homology, but these two toxins show a conserved pattern of [hydrophilic/hydrophobic](#) residues that fits in an [amphipathic](#)  $\alpha$ -helix, followed by a short very hydrophilic segment (4). These peptides have a great tendency to adsorb parallel to the membrane surface, where they self-associate and perturb the bilayer structure and thickness. If there is a membrane potential, they form channels of high conductance. Melittin and d-lysin are active on most **eukaryotic** cells and, at high concentrations, cause a colloid osmotic lysis: Red blood cells first become permeable to small solutes and then leak hemoglobin (5). At low concentrations, these peptides are able to induce the activation of endogenous cell phospholipase A2, particularly in mast cells, with the release of inflammatory mediators (1).

Another family of bacterial toxins, including the  $\alpha$ -toxin of *Staphylococcus aureus* (33 kDa) and the aerolysin (48 kDa) produced by *Aeromonas sp.*, alter membrane permeability by forming pores of high conductance (6-8). These toxins concentrate on the membrane surface, by binding as monomers to various glycoproteins. They then oligomerize to form heptamers, which insert in the lipid bilayer to form  $\beta$ -barrel transmembrane pores, similar to those of mitochondrial and bacterial **porins**, with a water-filled channel with a 1- to 1.5-nm inner diameter (6, 7). These toxins can be used to permeabilize cells in a controlled way so that the cell content of small ions and metabolites, but not cytosolic proteins, can be manipulated.

Other membrane-permeabilizing toxins are specific for cholesterol-containing membranes and are activated by **thiol** compounds. These cholesterol-binding toxins (50 to 60 kDa) are produced by many different **gram-positive** bacteria; their best-known members are streptolysin O and lysteriolysin O (9). They bind very rapidly in a nonsaturable mode to membranes containing cholesterol. Two monomers then form a membrane-embedded start complex that grows by monomer recruitment to form arc- and finally ring-shaped large oligomeric complexes of 25 to 50 toxin molecules (10). The inner diameter of these cholesterol-toxin complexes is 7.5 nm, and they permit a large loss of proteins from the cell. Hence, they can be used to achieve a rather complete manipulation of the cell cytosol. All of them are active on the plasma membrane, except lysteriolysin O, whose role is that of lysing phagosomes (see [Lysozymes](#)) and hence is active only at acidic pH values.

### Bibliography

1. G. Menestrina, G. Schiavo, and C. Montecucco (1994) *Mol. Aspects Med.* **15**, 81–193.
2. H. L. Harvey (1990) In *Handbook of Toxicology* (W. T. Shier and D. Mebs, eds.), Marcel Dekker, New York, pp. 1–66.
3. J. E. Alouf and J. H. Freer, eds. (1991) *Sourcebook of Bacterial Protein Toxins*, Academic Press, San Diego, CA.
4. J. Dufourcq and S. Castano (1997) In *Guidebook to Protein Toxins and Their Use in Cell Biology* (R. Rappuoli and C. Montecucco, eds.), Sambrook and Tooze, Oxford University Press, Oxford.
5. M. T. Tosteson, S. J. Holmes, M. Razin, and D. C. Tosteson (1985) *J. Membr. Biol.* **87**, 35–44.
6. G. Menestrina (1986) *J. Membr. Biol.* **90**, 177–190.

7. L. Song et al. (1996) *Science* **274**, 1859–1866.
8. M. W. Parker et al. (1992) *Nature* **367**, 292–295.
9. J. E. Alouf and C. Geoffroy (1991) In *Sourcebook of Bacterial Protein Toxins* (J. E. Alouf and J. H. Freer, eds.), Academic Press, San Diego, CA, pp. 147–186.
10. I. Waley et al. (1995) *Infect. Immun.* **63**, 1188–1194.

## Membrane Anchors

[Proteins](#) are bound to [membranes](#) by several different mechanisms. Some proteins are bound because **hydrophobic** regions of the [polypeptide chain](#) traverse the [lipid](#) bilayer several times and, as a consequence, a substantial proportion of the protein is actually located in the hydrophobic interior of the bilayer (see [Membrane Proteins](#)). Other proteins have a much more tenuous association with the membrane and utilize a relatively small membrane anchor for attachment. Depending on the particular protein, this membrane anchor may be a single hydrophobic, membrane-spanning segment of the polypeptide or a variety of different lipids covalently attached to specific sites in the [polypeptide chain](#). Lipid anchors can confer unusual biophysical properties on membrane proteins and are believed to play an important role in protein function.

### 0.1. Membrane Proteins Can be Anchored by Polypeptide Chains or by Lipid Molecules

Biological membranes consist of a lipid bilayer to which proteins are attached by a variety of mechanisms. For most membrane proteins, attachment is due to the presence of one or more relatively hydrophobic regions of the polypeptide chain that traverse the lipid bilayer. The proteins are held in place because the hydrophobic side chains exposed on the surface of the **alpha-helical** transmembrane regions are shielded from [water](#) molecules by the lipid bilayer (see [Hydrophobic Effect](#)). Conversely, the relatively polar residues that usually border the transmembrane segments are located in the aqueous phase on either side of the bilayer. The number of transmembrane regions can vary enormously among different proteins. For example, **ion channels** or **transporters** of water-soluble molecules usually have multiple hydrophobic transmembrane regions in a single polypeptide chain, interspersed with [hydrophilic](#) regions located outside the bilayer. Because a large proportion of its structure is inside the bilayer, this type of protein is commonly referred to as an “integral” membrane protein.

Many membrane proteins, however, do not interact so extensively with the bilayer and have only one transmembrane region holding them in the membrane. Strictly speaking, these are “polypeptide-anchored” rather than “integral” membrane proteins. For some of these proteins, the transmembrane region plays a crucial role in communication between the protein domains on opposite sides of the bilayer and also in forming functional complexes with other membrane proteins. For many “polypeptide-anchored” proteins, however, the transmembrane region does not have any role other than attachment to the membrane. This is particularly common for many **cell-surface adhesion receptors** and [enzymes](#) that need to be displayed prominently on the cell surface but do not need to communicate across the bilayer. For these proteins, a polypeptide anchor with a minimal intracellular domain is sufficient.

In some membrane-bound proteins, the transmembrane region has been dispensed with altogether, and the protein is covalently linked to a lipid molecule located in the bilayer that it uses as a membrane anchor. Although the lipid-anchored proteins probably represent a minority of membrane proteins, they are widely distributed among different organisms and functional classes of protein. Several hundred have been identified, but this is undoubtedly an underestimate because it is not

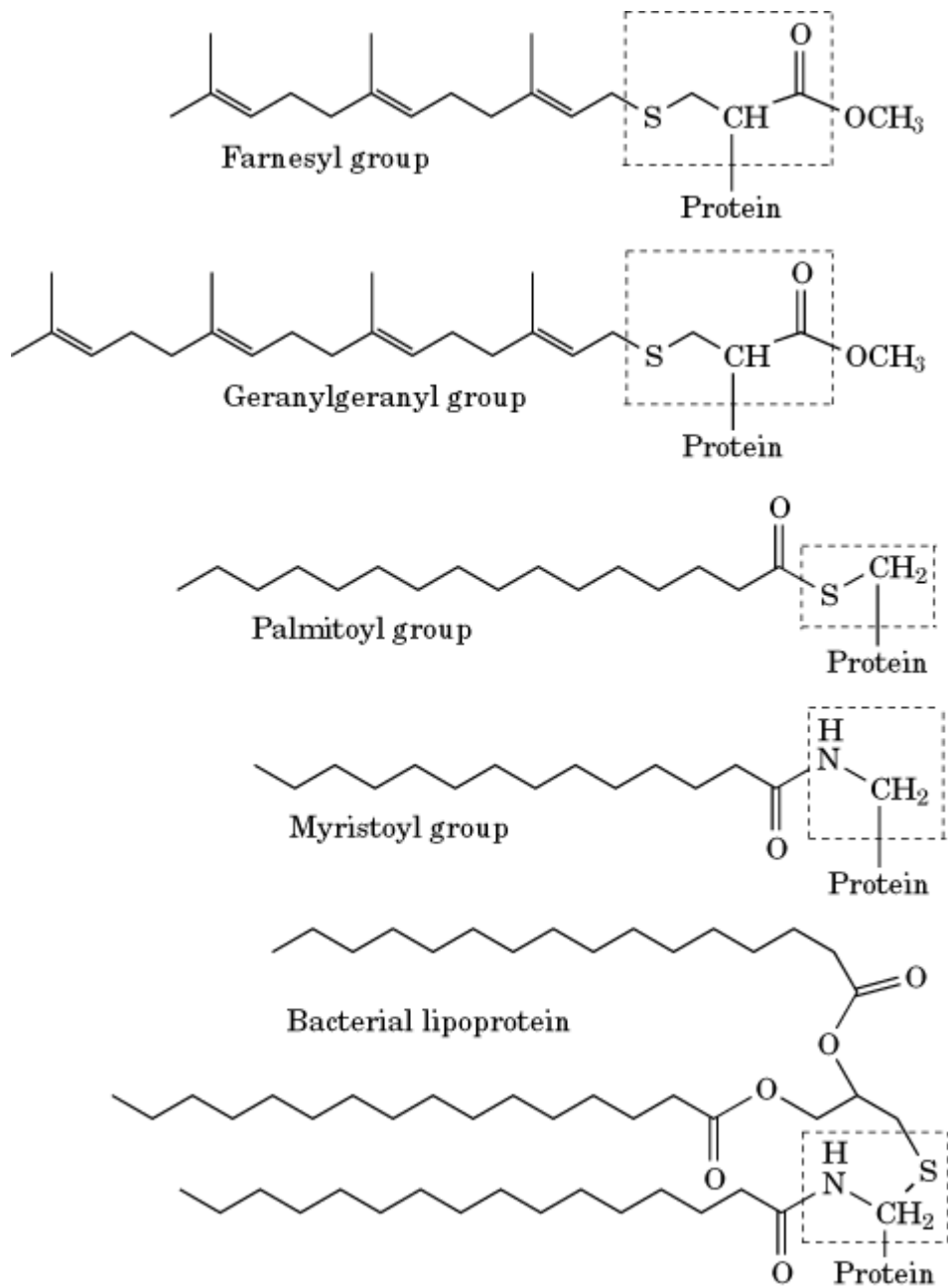
always possible to predict the occurrence of a lipid anchor from DNA sequences. The reason that lipid anchors are used for some proteins in preference to polypeptide anchors is uncertain. Lipid anchors are both structurally and metabolically varied, however, and their use provides the cell with an array of novel and versatile mechanisms for modulating the affinity of the protein for the membrane that are not possible with polypeptide anchors. The rest of this article is devoted to a general description of lipid anchors and to their potential role in regulating protein function. Additional information and examples of particular types of lipid anchor are given in individual entries.

## 0.2. Lipid Anchors Have Diverse Structures

In general, the lipid groups used as membrane anchors are attached to specific amino acid residues at or close to the N- or C-terminus of the protein. Attachment of lipid groups is also often accompanied or preceded by **proteolytic** processing at the N- or C-terminus. Five different types of lipid groups are used as membrane anchors (see also Table 1 and Figure 1):

1. An isoprenoid group that contains 15 (ie, **farnesyl**) or 20 (geranylgeranyl) carbon atoms attached to a **cysteine** residue at the C-terminus via a thioether linkage (1, 2). The carboxyl group of the cysteine may also be methylated (see **Prenylation**).
2. A myristic acid (14-carbon, saturated **fatty acid**) molecule attached to a **glycine** residue at the N-terminus via an amide linkage (see **Myristoylation**) (3, 4).
3. A palmitic acid (16-carbon, saturated fatty acid) molecule attached to a cysteine residue close to the N- or C-terminus via a thioester linkage or to internal **serine** and **threonine** residues via an ester linkage (see **Palmitoylation**) (3, 5).
4. A fatty-acid or 1,2-**diacylglycerol** molecule attached to the N-terminal cysteine via amide or thioether linkages, respectively (6, 7). This type of lipid anchoring is not found in eukaryotes but is common in eubacteria, where it anchors proteins to the inner or outer membrane (ie, facing the periplasmic space or on the cell surface) of **gram-negative** bacteria or to the cell surface of **gram-positive** bacteria, spirochaetes, **mycoplasma** and rickettsia. The proteins that use this type of anchor are functionally diverse and include structural proteins, extracellular enzymes, and proteins involved in transport and conjugation. This type of protein is commonly referred to by the generic, trivial name “*bacterial lipoprotein*” based on the first bacterial protein found covalently attached to a lipid—the murein lipoprotein of *Escherichia coli*. This is an unfortunate choice of a name because “bacterial lipoproteins” are in no sense comparable to vertebrate **lipoproteins** used for lipid transport. The latter are relatively large and heterogeneous complexes of lipids and proteins that are not covalently linked together.
5. A glycosyl phosphatidyl inositol molecule (GPI) attached to the C-terminal amino acid via an amide linkage (see **GPI Anchor**). The GPI molecule consists of a **phosphatidylinositol** molecule linked to a conserved glycan chain that contains (at minimum) a glucosamine residue, three mannose residues, and a terminal ethanolamine phosphate (8-10).

**Figure 1.** Structures of the major types of lipid groups that are used to anchor proteins to membranes. Additional details are given in the text or in the individual entries. (1) The C-terminal cysteine residue is thioether-linked to an isoprenyl group (*farnesyl* or *geranylgeranyl*), and the exposed carboxyl group is methylated. (2) The N-terminal glycine residue is amide-linked to a *myristoyl* group. (3) An internal cysteine residue (usually close to the N- or C-terminus) is thioester-linked to a *palmitoyl* group. In some proteins palmitoyl groups are ester-linked to the hydroxyl groups on serine or threonine residues (not shown). (4) In the so-called *bacteriallipoprotein* structure, the N-terminal cysteine residue is amide-linked to a fatty acid and thioether-linked to a diacylglycerol. For simplicity, all the fatty acids are shown as palmitate (the predominant fatty acid in *E. coli* murein lipoprotein), but generally the fatty acid chain length, degree of unsaturation etc. will be similar to the overall fatty acid composition of the membrane phospholipids and will vary in different bacteria. In (1)–(4) the dotted line encloses the amino acid residue modified by lipid. (5) The exposed carboxyl group on the C-terminal residue is amide-linked to the *GPI anchor*. The identity of the C-terminal residue varies for different proteins. The dotted lines enclose the major components of the GPI anchor. See **GPI Anchor**, Figure 1 for GPI anchor structure.



**Table 1. Characteristics of Lipid Anchors**

| Lipid Anchor           | Organisms           | Attachment Site                | Subcellular Location | Stability/Affinity of Membrane Binding |
|------------------------|---------------------|--------------------------------|----------------------|--|
| Isoprenyl <sup>a</sup> | Eukaryotes          | Cysteine residue at C-terminus | Intracellular        | Low–medium                             |
| Myristoyl              | Eukaryotes; Viruses | Glycine residue at             | Intracellular        | Low                                    |

|  |                     |  |                      |
|--|---------------------|--|----------------------|
| Palmitoyl                                | Eukaryotes; viruses | N-terminus<br>Cysteine residue near N- or C-terminus | Intracellular Medium |
| Fatty acyl + diacylglycerol <sup>a</sup> | Prokaryotes;        | Cysteine residue at N-terminus                       | Cell surface High    |
| GPI anchor                               | Eukaryotes          | Various residues at C-terminus                       | Cell surface High    |

<sup>a</sup> Farnesyl or geranylgeranyl.

<sup>b</sup> Also known as “bacterial lipoprotein.”

### 0.3. Lipid Anchors Provide a Mechanism for Reversible Membrane Association that Can be Regulated

The reason that lipid anchors are used by some proteins instead of polypeptide anchors is uncertain. A unique feature of lipid anchors however, is, that they do not span the membrane and consequently can be removed from and reinserted into the lipid bilayer relatively easily. This property allows the association of a lipid-anchored protein with the membrane to be reversible. In contrast, reversible membrane association of a polypeptide-anchored protein would be unlikely for two reasons: (1) During dissociation of the protein from the membrane, the energy cost of transferring even a small hydrophilic domain across the bilayer would be prohibitively high. (2) Even if it were possible to remove the protein from the membrane, it is unlikely that the transmembrane polypeptide could reinsert correctly into the bilayer without the assistance of some protein-translocating mechanism similar to that used during its biosynthesis. Correct reinsertion of an integral protein with significant protein domains on both sides of the bilayer would be essentially impossible. Although the physiological importance of reversible association has not been established directly for most lipid-anchored proteins, this property has been observed for individual proteins in all four classes of eukaryotic lipid anchors. It seems likely that reversible association is a general property of lipid anchors. An important aspect of reversible membrane association is that it also offers the possibility for rapidly changing the distribution of proteins between the membrane surface and the surrounding fluid by increasing or decreasing the membrane affinity. When a protein translocates from the three-dimensional fluid phase to the two-dimensional membrane surface, the decrease in dimensionality produces a large increase in the **effective concentration** and therefore enhances interaction with other membrane components relative to those in the fluid phase. Many lipid-anchored proteins have important roles in [signal transduction](#) and **vesicle** traffic, where the ability to bind rapidly and reversibly to membranes is essential.

The membrane affinity of lipid-anchored proteins inside intact cells is difficult to determine directly, but two experimental approaches have provided some insight into this problem.

1. The membrane affinities of model prenylated, myristoylated, or palmitoylated oligopeptides can be estimated from binding experiments with [liposomes](#) (11-13). Based on these studies, the lipid anchors can be arranged in order of increasing membrane affinity as follows:

(1)

Unfortunately, it is difficult to extrapolate from these data to proteins in living cells. The membrane affinity of the lipid anchor is offset by the loss of **entropy** resulting from restricted translational, rotational, and conformational motion. Although the entropic cost of membrane

binding undoubtedly increases with the size of the peptide, the magnitude of this increase is likely to be highly structurally dependent and therefore difficult to predict (11, 13).

2. [Site-directed mutagenesis](#) of lipid-attachment sites allows the effect of individual lipid groups on the subcellular distribution of proteins (and their response to regulatory signals) to be investigated systematically. In general, the binding affinities deduced from the model peptide/liposome studies reasonably agree with the mutagenetic studies in intact cells. Based on this information, several mechanisms for modulating the membrane affinity of lipid-anchored proteins have been proposed. Some of these mechanisms and their potential functions are described here and also in the entries on particular lipid anchors (see before).

### 0.3.1. Multiple anchors

As noted before, the lipid anchors on cytosolic proteins have low to medium affinity for membranes (compared to lipid anchors used on cell-surface proteins) as a consequence of their short length and/or branched structure (see Table 1 and Figure 1). However, the membrane affinity of a lipid-anchored protein can be increased by using multiple anchors in different combinations. Thus, the relatively weak farnesyl anchor in some **Ras proteins** can be reinforced by the adding a palmitoyl group close by. Similarly, the weak myristoyl group in some Src-related **tyrosine kinases** (p59<sup>fyn</sup> and p56<sup>lck</sup>) or some of the G<sub>α</sub> subunits of heterotrimeric **G proteins** can be reinforced by adding one or more palmitoyl residues. Small G-proteins in the Rab family use two geranylgeranyl molecules, whereas those in the Rho/Rac family only use one. When multiple anchors are used, they are generally close together. This arrangement maximizes the increase in membrane affinity due to the additional lipid group(s). Membrane binding of the second lipid group is more likely because it is restricted (or “tethered”) to a small aqueous volume close to the membrane as a result of the membrane binding of the first lipid group (12, 14). Another consequence of dual lipid anchors is that, in addition to the thermodynamic effect of increasing the membrane affinity, a second lipid anchor markedly reduces the rate of reversible association-dissociation from the membrane, causing the protein to become irreversibly trapped in the membrane (15).

### 0.3.2. Electrostatic interactions

Several phospholipids found in biological membranes have acidic head groups that confer a negative charge on the surface of the phospholipid bilayer. Consequently, regions of a polypeptide chain that are held close to the membrane (by virtue of a lipid anchor) are subject to electrostatic repulsion or attraction if that region also possesses an excess of acidic or basic residues, respectively. Electrostatic interactions might not effect proteins that have strong anchors but would have a profound influence on the membrane affinity of proteins with relatively weak lipid anchors. For example, myristoylated p60<sup>src</sup> and farnesylated K-Ras do not have any cysteine residues located close to their lipid anchors and therefore cannot increase their membrane affinity by palmitoylation (as described before). Instead, these proteins depend for their membrane association on **electrostatic** attraction resulting from polybasic regions located close to the N- or C-terminus (16, 17). The strength of the electrostatic interaction may be subject to modulation by metabolic processes, as suggested for the MARCKS protein. This myristoylated protein contains multiple sites in its polybasic region that are **phosphorylated** by [protein kinase C](#). Phosphorylation neutralizes the local positive charge and reduces the electrostatic attraction, resulting in the release of MARCKS protein from the membrane (14).

### 0.3.3. Masking or removing the lipid anchor

Masking the lipid anchor by binding it to a soluble protein or removing it with an enzyme may also provide a mechanism for regulating the membrane affinity of lipid-anchored proteins. For example, the diprenylated lipid group on Rab proteins provides a strong anchor, but it can be masked by a lipid binding site on a cytoplasmic protein (Rab GDP dissociation inhibitor, or GDI). Thus, GDI can shift the distribution of Rab proteins from the membrane to the cytoplasm (18). A similar phenomenon occurs with recoverin, a Ca<sup>2+</sup>-binding protein found in photoreceptor cells. In this case,

however, the myristate binds to a hydrophobic patch on recoverin itself, rather than to a separate protein. In the presence of  $\text{Ca}^{2+}$ , recoverin undergoes a conformational change that results in extruding the myristate from its binding site and reassociation with the membrane (19). Enzymatic cleavage of the lipid anchor may also play a role in reversing the membrane association of S-palmitoylated proteins. As noted before, palmitoylation often occurs in conjunction with other lipid anchors, and depalmitoylation may be an important mechanism for preventing these proteins becoming permanently trapped in a membrane (15). It has even been suggested that receptor activation of heterotrimeric G-proteins stimulates a cycle of palmitoylation-depalmitoylation, but the functional relevance of these claims will remain speculative until the enzymes involved in this cycle have been characterized at the molecular level (see [Palmitoylation](#)) (20, 21). There are also enzymes that can cleave the GPI anchor on cell surface proteins, and one of these occurs in mammals and other vertebrates (see [GPI Anchor](#)). In contrast to depalmitoylation, however, this process is essentially irreversible. These enzymes cleave the lipid group into two, rather than removing it altogether. Therefore, reattachment of the anchor is not a simple recapitulation of biosynthesis. Furthermore, if degradation of the GPI anchor occurs at the cell surface, the protein is lost from the cell, and it is unlikely to regain access to the site of biosynthesis. At present, the physiological purpose of GPI anchor cleavage remains obscure.

#### 0.3.4. How can lipid-anchored proteins become enriched in specific membranes?

By its very nature, the association between a lipid anchor and the lipid bilayer is highly nonspecific, whereas intracellular, lipid-anchored proteins usually have quite specific subcellular distributions. Although the mechanisms responsible for directing lipid-anchored proteins to particular membranes are not fully understood, recognition and high-affinity binding of the lipid anchor by specific membrane-bound “**receptors**” may not be involved in the majority of cases. Because membrane association is rapidly reversible, a number of less direct mechanisms might enable a lipid-anchored protein to become enriched in a particular subcellular compartment (15). First, one of the enzymes involved in attaching the lipid might be localized in a specific membrane. An example of this is the relative enrichment of S-palmitoylating enzymes in the plasma membrane of some cells.

Consequently, a myristoylated protein might bind equally well to the cytoplasmic surfaces of **Golgi**, endoplasmic reticulum, or plasma membranes, but it would become palmitoylated only in the plasma membrane. Palmitoylation would increase the membrane affinity of the protein and also because of the dual lipid anchor, the rate of dissociation would be reduced drastically, and the protein would in essence become trapped there permanently. Carboxymethylation of prenylated proteins also occurs in membranes also and may play a role in trapping proteins in particular compartments by a similar mechanism. Masking of the lipid anchor may also be regulated, so that a protein becomes localized in a particular membrane. Thus, the interaction between Rab and GDI (which masks the prenyl groups as described before) is broken by specific GDI displacement factors that are located in particular membrane compartments. Exposure of the lipid group increases the membrane affinity of Rab and, due to its proximity, results in preferential binding to the membrane that contains the displacement factor. This type of mechanism is especially useful (and possibly more energy-efficient) for Rab proteins, which play an important role in vesicular transport and must recycle repeatedly between specific membrane compartments in the endoplasmic reticulum, the Golgi, endosomes, and the plasma membrane (18).

#### Bibliography

1. S. Clarke (1992) *Annu. Rev. Biochem.* **61**, 355–386.
2. F. L. Zhang and P. J. Casey (1996) *Annu. Rev. Biochem.* **65**, 241–269.
3. D. A. Towler, J. I. Gordon, S. P. Adams, and L. Glaser (1988) *Annu. Rev. Biochem.* **57**, 69–99.
4. D. R. Johnson, R. S. Bhatnagar, L. J. Knoll, and J. I. Gordon (1994) *Annu. Rev. Biochem.* **63**, 869–914.
5. M. F. G. Schmidt (1989) *Biochim. Biophys. Acta* **988**, 411–426.
6. H. C. Wu and M. Tokunaga (1986) *Curr. Top. Microbiol. Immunol.* **125**, 127–157.
7. S. Hayashi and H. C. Wu (1990) *J. Bioenergetics Biomembranes* **22**, 451–471.

8. M. G. Low (1989) *Biochim. Biophys. Acta* **988**, 427–454.
9. M. J. McConville and M. A. J. Ferguson (1993) *Biochem. J.* **294**, 305–324.
10. M. C. Field and A. K. Menon (1993) In *Lipid Modifications of Proteins* (M. J. Schlesinger, ed.), CRC Press, Boca Raton, pp. 83–134.
11. R. M. Peitzsch and S. McLaughlin (1993) *Biochemistry* **32**, 10436–10443.
12. C. A. Buser, C. T. Sigal, M. D. Resh, and S. McLaughlin (1994) *Biochemistry* **33**, 13093–13101.
13. J. R. Silvius and F. l'Heureux (1994) *Biochemistry* **33**, 3014–3022.
14. J. T. Seykora, M. M. Myat, L.-A. Allen, J. V. Ravetch, and A. Aderem (1996) *J. Biol. Chem.* **271**, 18797–18802.
15. S. Shahinian and J. R. Silvius (1995) *Biochemistry* **34**, 3813–3822.
16. M. D. Resh (1994) *Cell* **76**, 411–413.
17. L. Alland, S. M. Peseckis, R. E. Atherton, L. Berthiaume, and M. D. Resh (1994) *J. Biol. Chem.* **269**, 16701–16705.
18. S. R. Pfeffer, A. B. Dirac-Svejstrup, and T. Soldati (1995) *J. Biol. Chem.* **270**, 17057–17059.
19. J. B. Ames, T. Tanaka, M. Ikura, and L. Stryer (1995) *J. Biol. Chem.* **270**, 30909–30913.
20. S. M. Mumby, C. Kleuss, and A. Gilman (1994) *Proc. Natl. Acad. Sci. USA* **91**, 2800–2804.
21. M. Y. Degtyarev, A. M. Spiegel, and T. L. Z. Jones (1993) *J. Biol. Chem.* **268**, 23769–23772.

### Suggestions for Further Reading

22. G. Milligan, M. Parenti, and A. I. Magee (1995) The dynamic role of palmitoylation in signal transduction, *Trends Biochem. Sci.* **20**, 181–186.
23. S. McLaughlin and A. Aderem (1995) The myristoyl-electrostatic switch: A modulator of reversible protein-membrane interactions, *Trends Biochem. Sci.* **20**, 272–276.
24. P. J. Casey and J. E. Buss, (eds.) (1995) *Lipid modifications of proteins*, In *Methods in Enzymology* **250**, Academic Press, New York.
25. R. S. Bhatnagar and J. I. Gordon (1997) Understanding covalent modification of proteins by lipid: Where cell biology and biophysics mingle, *Trends Cell Biol.* **7**, 14–20.
26. S. Ilangumaran and D. C. Hoessli (1998) *Glycosylphosphatidylinositol-Anchored Biomolecules*, R. G. Landes, Austin, TX.

## Membrane Potentials

Electrical potentials across cell [membranes](#) are a function of the electrolyte concentrations in the intracellular and extracellular solutions and of the permeabilities of the ions. For many cells, the principal diffusing ions are sodium, potassium, and chloride. The *Goldman–Hodgkin–Katz equation* (1, 2)

$$\Psi = -\frac{RT}{F} \ln \frac{P_{K^+}[K^+]_{in} + P_{Na^+}[Na^+]_{in} + P_{Cl^-}[Cl^-]_{out}}{P_{K^+}[K^+]_{out} + P_{Na^+}[Na^+]_{out} + P_{Cl^-}[Cl^-]_{in}} \quad (1)$$

provides a quantitative relation for predicting the membrane potential of resting cells, where  $R$  is the universal gas constant,  $T$  is the absolute temperature, and  $F$  is the Faraday constant; this relation requires only the permeabilities  $P$  and the concentrations of sodium, potassium, and chloride ions



inside and outside the cell. Permeabilities have been obtained by unidirectional flux measurements using radiotracers. In nerve and muscle, the intracellular concentrations of sodium and chloride are low relative to their concentration in plasma while intracellular potassium is high relative to its plasma concentration. Generally the sodium ion gradient opposed to that of potassium is maintained by the  $\text{Na}^+/\text{K}^+$  pump ATPase. In nerves, the resting potential is  $-60\text{mV}$ , with the inside of the cell at a negative potential relative to the external medium. This is close to the potassium equilibrium potential and is largely due to a much higher potassium than sodium permeability across the membrane. An action potential is generated when the resting nerve is activated by an electrical impulse; there is a rapid (0.5 ms) depolarization of the membrane potential to  $+50\text{mV}$ , which is followed by a slower (several milliseconds) repolarization to the resting potential. The rapid depolarization is due to the transient opening of sodium channels and a concomitant increase in sodium permeability, followed by an increase in potassium permeability to reestablish the resting potential. The transient character of the action potential may be accounted for by the temporal and sequential opening and closing of the sodium and potassium channels (3). Contributions to membrane potentials are also derived from electrogenic ion pumps in [mitochondria](#), **bacteria**, and **chloroplasts**. The energy released from the oxidation–reduction reactions in these systems is coupled to the generation of a proton gradient across their membranes; this gradient provides the driving force for the generation of other energy-dependent processes and to move ions against their concentration gradients (4, 5) (see [Chemiosmotic Coupling](#)).

#### Bibliography

1. D. E. Goldman (1943) *J. Gen. Physiol.* **27**, 37–60.
2. A. L. Hodgkin and B. Katz (1949) *J. Physiol.* **108**, 37–77.
3. A. L. Hodgkin and A. F. Huxley (1952) *J. Physiol.* **117**, 500–544.
4. P. D. Boyer et al. (1977) *Annu. Rev. Biochem.* **46**, 955–1026.
5. G. Hauska and A. Trebst (1977) *Current Topics in Bioenergetics* **6**, 152–211.

#### Suggestions for Further Reading

6. H. Davson (1964) *A Textbook of General Physiology*, 3rd ed., Little Brown, Boston. (A general source book for the basic physical chemistry of membrane potentials, as well as for the early experimental basis of the membrane phenomena.)
7. J. B. Finean, R. Coleman, and R. H. Michel (1978) *Membranes and their Cellular Functions*, 2nd ed., Wiley, (A readable summary of many of the properties of cell membranes.)
8. W. J. Edelman Jr.(ed.) (1971) *Biophysics and Physiology of Excitable Membranes*, Van Nostrand Reinhold, New York. (Contains a useful summary of theoretical and physical aspects of membrane potentials, along with a summary of analytical methods for action potentials.)

#### Membrane Proteins

Intrinsic or integral membrane proteins are defined as [proteins](#) that penetrate into and, most often, traverse the [lipid](#) bilayer of a biological [membrane](#). This makes them fundamentally different from proteins that are anchored to the membrane via a fatty acid or a prenyl group attached to one of their termini (see [Membrane Anchors](#)). [Protein structures](#) that partition into lipid, rather than remain in aqueous solution, have specific chemical and structural properties. They are rich in exposed **hydrophobic** amino acids and, owing to the need to satisfy the **hydrogen-bonding** potential of the peptide groups in the membrane environment, they are restricted in their patterns of **secondary**

**structure.** A consequence of the exposed hydrophobic surface is that an integral membrane protein can only be brought into aqueous solution (solubilized) in the presence of a [detergent](#). The detergent is crucial for the stability and activity of a given membrane protein in solution. In most cases, membrane proteins can be purified using the same methods as soluble proteins, but all solutions must contain a detergent at a concentration greater than the [critical micelle concentration](#).

## 1. Biological Importance

About a quarter (20% to 25%) of **genes** are predicted to encode integral membrane proteins in the **genomes** of microorganisms and multicellular organisms. This large number reflects the importance of membranes in the life of any single cell. In bacteria, the majority of membrane proteins function in **membrane transport**, secretion, and bioenergetic processes. As the importance of cell–cell communication increases in multicellular organisms, the fraction of genes linked to this membrane function also tends to increase, reflecting the presence of a large number of receptors in the cytoplasmic membrane and the sophisticated organization of intracellular traffic and cell organelles. In the light of the biological importance of membrane-bound processes, the scarcity of structural data on membrane proteins—and consequently the lack of molecular details concerning their functional mechanisms—creates frustration in present-day molecular biology.

Membrane proteins often assemble into large complexes. In particular, [photosynthetic reaction centers](#) in plants and bacteria, [ATP synthase](#), and respiratory complexes contain a large number of subunits, pigments, and bound metals. Similarly, molecular machines involved in [protein secretion](#) or ATP-driven **active transport** usually have several protein components, and **ion channels** are typically oligomeric.

## 2. Structures

Three-dimensional [protein structures](#) at atomic resolution have been determined for only about 20 membrane proteins (Table 1). This contrasts with the situation for soluble proteins, for which thousands of high-resolution structures determined by [X-ray crystallography](#) or [NMR](#) spectroscopy have already been deposited in the **structure data banks**. The small number of high-resolution structures of membrane proteins reflects the great difficulty in crystallizing membrane proteins. This appears to be caused by two factors. First, the presence of the detergent, with its own phase behavior, significantly complicates the search for the right crystallization conditions. Second, the critical protein–protein contacts in the crystal lattices of most membrane proteins occur between the polar peripheral regions. These polar parts are often small, and they may also be masked by the polar headgroups of the detergent, disfavoring the formation of a crystal lattice.

**Table 1. Structures of Membrane Proteins at Atomic Resolution (Determined by the End of 1998)**

| Protein (in Chronological Order)                           | Function          | Transmembrane Type (Helical/b-Barrel/Monotopic) | Method of Structure Determination/Current Resolution | Crystallization and a          |
|--|-------------------|---|--|--------------------------------|
| Photosynthetic reaction center from <i>Rps. viridis</i>    | Energy conversion | Helical   | X-ray/2.3 Å  | Lauryldimethoxide and he       |
| Photosynthetic reaction center from <i>Rb. sphaeroides</i> | Energy conversion | Helical   | X-ray/2.7 Å  | Lauryldimethoxide and heptanet |
| Bacteriorhodopsin from                                     | Light-driven      | Helical   | (1) Electron   | (1) and (2) N                  |

|  |  |           |   |   |
|--|--|-----------|---|---|
| <i>H. salinarium</i>   | proton translocation   |           | crystallography/3.5 Å<br>(2) 2.3 Å  | occurring 2I<br>the purple m<br>(3) Protein i<br>glucoside in<br>cubic phase                |
| Porin from <i>Rb. capsulatus</i>                                       | Facilitated diffusion across the outer membrane                  | b-Barrel  | (3) X-ray/2.5 Å<br>X-ray/1.8 Å  | Octyltetraox  |
| OmpF and PhoA porins from <i>E. coli</i>                               | Facilitated diffusion across the outer membrane                  | b-Barrel  | X-ray/2.7 Å   | (1) Octyltetra and octyl-2-hydroxyethyl<br><br>(2) Decyl dioxolane and octyl-2-hydroxyethyl |
| Porin from <i>Rb. blastica</i>   | Facilitated diffusion across the outer membrane                  | b-Barrel  | X-ray/1.96 Å  | Octyltetraox  |
| Light-harvesting complex from pea                                      | Capture of photons   | Helical   | Electron crystallography/3.4 Å  | Triton X-100<br>glucoside pl  |
| Prostaglandin H2 synthase 1 from sheep                                 | Conversion of arachidonate into a prostaglandin                  | Monotopic | X-ray/3.1 Å   | Octyl glucoside   |
| Light-harvesting complex from <i>Rps. acidophila</i>                   | Capture of photons   | Helical   | X-ray/2.5 Å   | (1) Dimethyl sulfoxide and benzamide<br><br>(2) Octyl glucoside and benzamide               |
| Maltoporin from <i>E. coli</i>   | Facilitated diffusion of maltodextrins across the outer membrane | b-Barrel  | X-ray/2.4 Å   | Decyl maltoside<br>dodecyl maltoside  |
| Cytochrome <i>c</i> oxidase from <i>Pa. denitrificans</i> <sup>a</sup> | Proton translocation coupled to the reduction of dioxygen        | Helical   | X-ray/(1) 2.8 Å (four-subunit complex)<br><br>(2) 2.7 Å (two-subunit complex) | (1) Dodecyl maltoside<br><br>(2) Undecyl maltoside<br>cyclohexyl-β-D-glucopyranoside        |
| Cytochrome <i>c</i> oxidase  | Proton   | Helical   | X-ray/2.3-2.8 Å   | Decyl maltoside   |

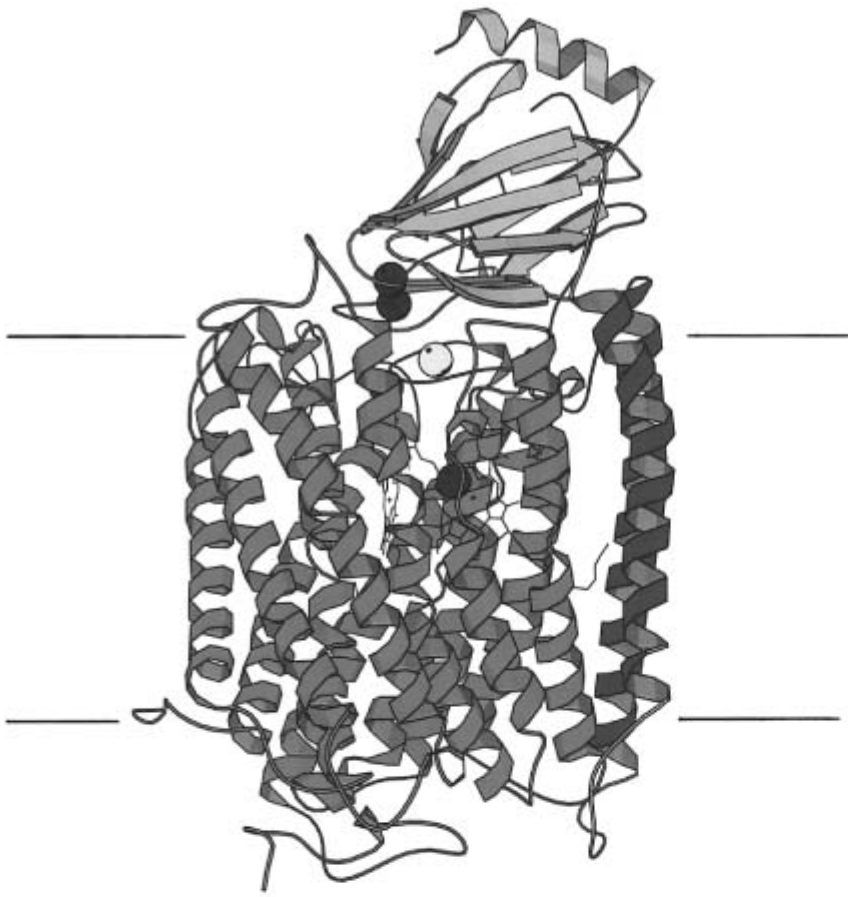
|  |   |                 |                                      |   |
|--|---|-----------------|--------------------------------------|---|
| from bovine heart  | translocation coupled to the reduction of dioxygen        |                 | (several structures)                 | some residues   |
| Prostaglandin H2 synthase 2 from man                       | Conversion of arachidonate into a prostaglandin           | Monotopic       | X-ray/2.9 Å                          | Octyl pentac  |
| Photosystem I from <i>Synechococcus elongatus</i>          | Energy conversion   | Helical         | X-ray/4.0 Å                          | Dodecyl ma  |
| $\alpha$ -Hemolysin from <i>Staphylococcus aureus</i>      | To make a pore to an eukaryotic cell membrane             | $\beta$ -Barrel | X-ray/1.9 Å                          | Octyl glucos  |
| Porin from <i>Pa. denitrificans</i>                        | Facilitated diffusion across the outer membrane           | $\beta$ -Barrel | X-ray/3.1 Å                          | Octylglucos   |
|  |   |                 | (2) 2.7 Å (two-subunit complex)      | (2) Undecyl cyclohexyl- $\beta$ -maltoside                  |
| Cytochrome <i>c</i> oxidase from bovine heart              | Proton translocation coupled to the reduction of dioxygen | Helical         | X-ray/2.3-2.8 Å (several structures) | Decyl malto<br>some residu                                  |
| Prostaglandin H2 synthase 2 from man                       | Conversion of arachidonate into a prostaglandin           | Monotopic       | X-ray/2.9 Å                          | Octyl pentac  |
| Photosystem I from <i>Synechococcus elongatus</i>          | Energy conversion   | Helical         | X-ray/4.0 Å                          | Dodecyl ma  |
| Porin from <i>Pa. denitrificans</i>                        | Facilitated diffusion across the outer membrane           | $\beta$ -Barrel | X-ray/3.1 Å                          | Octylglucos   |
| Light-harvesting complex from <i>Rs. molischianum</i>      | Capture of photons  | Helical         | X-ray/2.4 Å                          | Lauryl dimethylami  |
| Cytochrome <i>bc<sub>1</sub></i> complex from bovine heart | Respiratory electron transfer and energy transduction     | Helical         | X-ray/(1) 2.9 Å (partial structure)  | (1) Decanoy glucanamide diheptanoyl-choline (plus deoxychol |

|   |  |          |                                 |   |
|---|--|----------|---------------------------------|---|
|   |  |          | (2) 2.8–3.0 Å<br>(complete)     | (2) Dodecyl<br>dodecyl mal<br>heptyl-carba<br>methyl gluc |
| Maltoporin from <i>S. typhimurium</i>                 | Facilitated diffusion of maltodextrins across the outer membrane | b-Barrel | X-ray/2.4 Å                     | Octyltetraox<br>hexyldimeth                               |
| Cytochrome <i>bc<sub>1</sub></i> complex from chicken | Respiratory electron transfer and energy transduction            | Helical  | X-ray/3.0 Å (partial structure) | Octyl glucos  |
| Potassium channel from <i>Streptomyces lividans</i>   | Facilitated diffusion of potassium across the cell membrane      | Helical  | X-ray/3.2 Å                     | Lauryl dime   |
| Sucrose porin from <i>S. typhimurium</i>              | Facilitated diffusion of sucrose across the outer membrane       | b-Barrel | X-ray/2.9 Å                     | Mixture of c<br>hexyldimeth<br>and heptylgl               |
| OmpA protein from <i>E. coli</i>                      | Outer membrane receptor  | b-Barrel | X-ray/2.5 Å                     | Octyltetraox  |
| FhuA porin from <i>E. coli</i>                        | Ferrichrome-bound iron transport                                 | b-Barrel | X-ray/2.7 Å                     | Octyl-2-<br>hydroxyethy                                   |
| FepA porin from <i>E. coli</i>                        | Ferric enterobactin transport                                    | b-Barrel | X-ray/2.4 Å                     | Lauryl<br>dimethylami                                     |

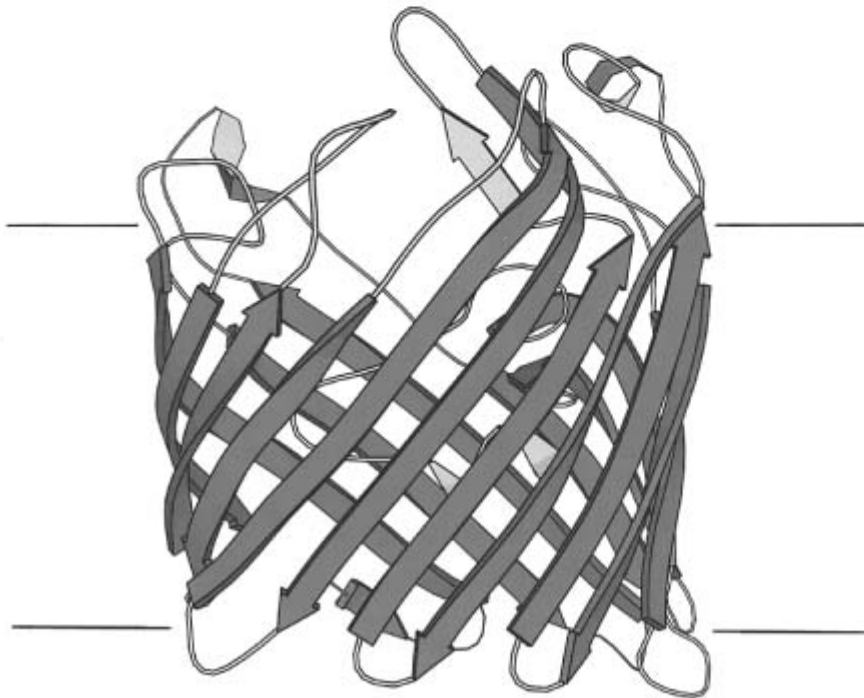
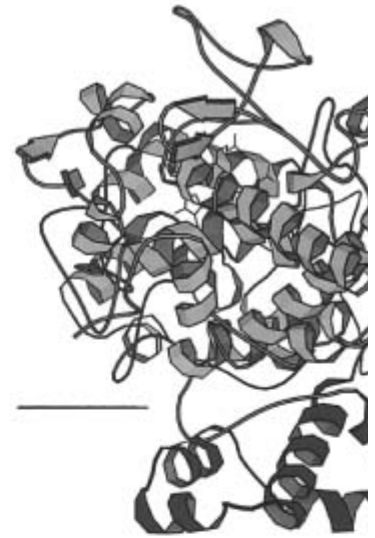
<sup>a</sup> In both crystallizations, the enzyme was cocrystallized with an antibody fragment.

The known membrane protein structures fall into three categories: a-helical, b-barrel, and monotopic membrane proteins (Table 1 and Fig. 1). The majority of the structures have been determined using three-dimensional crystals and X-ray diffraction. However, the structure of a plant [light-harvesting complex](#) has been determined using [electron crystallography](#) on two-dimensional (single-layered) crystals of the protein reconstituted with lipid. The structure of [bacteriorhodopsin](#) has been solved by both methods. [Electron microscopy](#) increasingly contributes to our knowledge of membrane protein structures. Improved image processing methods and particularly the use of [cryoelectron microscopy](#) techniques (see [Single Particle Reconstruction](#)) have already resulted in a number of low-resolution structures, either in projection through the membrane or in three dimensions after merging data from tilted two-dimensional crystals. Many of these structures are likely to evolve into a higher resolution.

**Figure 1.** Different structural types of integral membrane proteins. **(a)** A helical membrane protein complex: the two-subunit cytochrome *bc<sub>1</sub>* complex of *Paracoccus denitrificans* (PDB, file name 1AR1). **(b)** A  $\beta$ -barrel membrane protein: the OmpF porin of *E. coli* (1OPF). (c) A monomeric integral membrane protein: prostaglandin synthase 2 (5COX). Horizontal lines indicate approximately the location of the bilayer surfaces.



(a)



(b)

## 2.1. Helical Superfolds

Helical membrane proteins are characterized by having [a-helices](#) (Fig. 1a), and they appear to have superfolds characterized by the number of membrane-spanning segments and their inside/outside topology (1). In particular, there is a large family of seven-helix membrane proteins, the “7TM” superfamily, that typically act as receptors in eukaryotic organisms. Other abundant superfolds are formed by **four-helix bundles**, and some proteins have 12 predicted transmembrane helices. The latter often appear to function in **transport** of solutes or ions across membranes. A subclass of the 12-helix proteins may be formed by dimers, in which each monomer contains six a-helices.

Analysis of these superfolds may be used as a robust rule of thumb in prediction of the function of a novel integral membrane protein. It is clear, however, that no strict rules exist. For example, **membrane transport** is catalyzed by proteins that fall into different topological categories. An example is the tetrameric **potassium channels**, which can be divided into two subclasses. In one case, the monomer has two transmembrane helices, and hence the tetramer has eight, whereas the monomer of the second subclass has six membrane-spanning segments and the tetramer has 24. In both cases, only one helix of each monomer lines the ion-conducting channel (2).

The proteins in the “7TM” superfamily all have similar topologies with seven membrane-spanning a-helices and the *N*-terminus on the external side of the membrane. Many members of the family, such as [rhodopsin](#), are G-protein-coupled receptors (see [GTP-Binding Proteins](#) and [Heterotrimeric G Proteins](#)) that mediate signaling from the outside to the cellular interior and have a sensory function in processes such as hearing, seeing, smelling, and tasting. Generally speaking, however, the 7TM proteins are functionally very diverse. It seems that the seven-helix fold is an archetype that is able to adapt to different functions by accommodating diverse ligands or **prosthetic** groups.

## 2.2. Beta-Barrel Membrane Proteins

In a few membrane proteins, notably the **porins** and the membrane-bound form of some [toxins](#) such as the a-hemolysin (a heptameric toxin from *Staphylococcus aureus*), the transmembrane portion comprises an [b-barrel](#) made of 8, 14, 16, 18, or 22 antiparallel **b-strands** (Figure 1b). The b-barrel is stabilized by extensive hydrogen-bonding between the b-strands. The interior of the barrel is hydrophilic, reflecting its function of facilitating diffusion of small molecules across the membrane. Because of the extended conformation of b-strands, fewer than 10 residues are needed to traverse the bilayer. Of these residues, only every second has to be hydrophobic. This complicates the prediction of transmembrane b-strands from just the [primary structure](#). In fact, the amino acid sequences of porins and a-hemolysin are as [hydrophilic](#) as are those of soluble proteins. The monomeric a-hemolysin is a soluble protein, which rearranges to form the heptameric transmembrane pore upon interaction with lipids or detergent.

## 2.3. Monotopic Membrane Proteins

Prostaglandin H<sub>2</sub> synthases (also known as cyclo-oxygenases 1 and 2) require detergents for solubilization, but their three-dimensional structures contain no transmembrane a-helices or b-strands (Fig. 1c). These enzymes appear to be examples of monotopic membrane proteins. They associate with the membrane via four short amphipathic a-helices that run approximately parallel to the membrane plane (Fig. 1). The helices are [amphipathic](#) and contain a number of aromatic side chains, which are probably exposed to the lipid headgroup region of the membrane. Most likely, the four-helix motif is only immersed into one layer of the bilayer. This means that, although prostaglandin synthases are genuine membrane proteins, they are not transmembrane proteins. It is not known how common this type of arrangement is, because sequence-based prediction algorithms for monotopic membrane anchors have not been developed. The membrane-binding domains of cyclo-oxygenases 1 (sheep) and 2 (human) are the least conserved part of the enzymes (33% identity). However, the three-dimensional structure of the membrane anchor is well-preserved.

## 3. Structural Stability of Membrane Proteins

In contrast to soluble proteins, membrane proteins fold in an environment from which [water](#) is largely excluded. The absence of water forces the protein to form [hydrogen bonds](#) between peptide



groups within the molecule itself, driving the formation of transmembrane helices or the transmembrane  $\beta$ -barrel of **porins**. The hydrogen bonds in transmembrane helices are probably much stronger than the average bonds in soluble proteins, because there is no competition to form hydrogen bonds with water. As a result, each transmembrane helix is a highly stable and rigid structure that can be thought of as an independent folding unit. In soluble proteins, the tertiary fold at least partially results from the action of the [hydrophobic effect](#), which brings about formation of their hydrophobic core. In membrane proteins, the association of the transmembrane helices within the bilayer cannot be driven by the hydrophobic effect. Mutagenesis studies and analyses of known structures suggest that optimal [van der Waals interactions](#) are responsible for the helical packing within the bilayer. Buried residues are more conserved and more polar than those exposed to lipid, suggesting that a subtle stereochemical fit plays a key role in the formation of the [tertiary structure](#). The interiors of soluble and membrane proteins have been concluded to have similar polarities.

A consequence of this increased stability is the resistance of membrane-spanning helices toward **denaturation** with commonly used agents, such as **SDS**, [urea](#), and **guanidinium** chloride. In the routinely used [SDS-PAGE](#), integral membrane proteins with multiple transmembrane helices commonly migrate faster than expected from their true molecular weight. This may be due to (a) the maintenance of secondary and/or tertiary structure under the conditions in which soluble proteins are totally denatured and/or (b) increased binding of the anionic detergent.

#### Bibliography

1. D. T. Jones (1998) Do transmembrane protein superfolds exist? *FEBS Lett.* **423**, 281–285.
2. D. A. Doyle, J. M. Cabral, R. A. Pfuetzner, A. Kuo, J. M. Gulbis, S. L. Cohen, B. T. Chait, and R. MacKinnon (1998) The structure of the potassium channel: molecular basis of  $K^+$  conduction and selectivity. *Science* **280**, 69–77.

#### Suggestions for Further Reading

3. T. Haltia (1997) Structural features of membrane proteins. *Adv. Mol. and Cell Biol.* **22**, 229–277.
4. T. Haltia and E. Freire (1995) "Forces and factors that contribute to the structural stability of membrane proteins". *Biochim. Biophys. Acta* **1228**, 1–27.
5. D. C. Rees, L. Deantonio, and D. Eisenberg (1989) Hydrophobic organization of membrane proteins. *Science* **245**, 510–513.
6. C. Ostermeier and H. Michel (1997) Crystallization of membrane proteins. *Curr. Opin. Struct. Biol.* **7**, 697–701.

## Membranes

Membranes are the anatomical boundaries of cells and their organelles (see also [Endoplasmic Reticulum](#), [Mitochondria](#), [Nucleus](#), and [Plasmalemma](#)). The composition of membranes, which are principally lipid and protein, varies from 20% protein in the myelin sheath of nerve to 80% protein in bovine erythrocytes. The lipid is arranged as a planar bilayer (bimolecular layer) with an internal hydrocarbon interior shielded from the aqueous exterior and cytoplasmic interior of the cell by the [polar](#) moieties of the lipids. The distribution of phospholipids across the midplane of the bilayer is often asymmetric (1); in erythrocytes the external monolayer of the membrane bilayer contains mostly neutral phospholipids, such as phosphatidylcholine and sphingomyelin (see [Lipids](#)), while the internal monolayer contains a majority of the amino-phospholipids (phosphatidylserine, phosphatidylethanolamine) (2). This asymmetric distribution of lipids is maintained by metabolic

processes which, when suppressed, lead to a slow reversion toward a symmetric lipid distribution across the bilayer (3). The lipid composition of membranes varies widely within a cell, and from tissue to tissue.

The proteins that make up membranes represent virtually all the [enzymes](#) of the primary biochemical and physiological processes of living cells (eg, [photosynthesis](#), oxidative phosphorylation, nerve conduction), the centers for sensory perception (vision, taste, touch, smell), and for the selective transport of ions and molecules across membranes. Membrane-associated proteins are of two classes, those that are readily dissociated from the membrane and those that require drastic treatments for removal, a process that inevitably destroys the membrane. The former are called *peripheral* or *extrinsic* membrane proteins, while the latter are called *integral* membrane proteins. Integral proteins have been visualized in electron micrographs of **freeze-fractured** membrane preparations; they contain a stretch of 21 to 26 **hydrophobic** amino acid residues coiled into an **a-helix** that is believed to facilitate the spanning of a membrane bilayer (see [Amphipathic](#)).

The apparent morphological and compositional stability of membranes belies a system that seethes with activity. Both lipid and protein components undergo continuous synthesis and degradation, without any apparent synchronization of these processes. Thus, the various lipids and protein components in the same membrane turn over at different rates relative to each other (4). Membrane assembly is believed to occur on existing membranes, which serve as templates for insertion of newly synthesized protein and lipid. The primary centers for lipid and **protein synthesis** are also located on membranes. For **eukaryotic** cells the principal locus is the endoplasmic reticulum, and for **bacteria** the locus of synthesis is the plasma membrane. Approximately 180,000 reactions per second are performed by a bacterium to synthesize the membrane lipids for a cell population that doubles in size in 20 min. Microorganisms and many animals (eg, fish, reptiles) are *poikilothermic*, and their physiological temperatures are largely a function of the environment in which they live. Wide variations of environmental temperature are accommodated by changes in the ratio of saturated to unsaturated fatty acids in the phospholipids of their cell membranes; higher temperatures generally yield increasing amounts of saturated fatty acids.

Protein insertion into membranes requires leader sequences on the protein to convey information about the location for insertion, and to initiate entry through the lipid bilayer (5) (see [Topogenesis](#)). Cell membranes also possess the facility to undergo morphological and topological transformations that allow them to fuse with other membranes. They may also undergo [endocytosis](#), in which extracellular material is internalized by invagination of the plasma membrane, followed by formation of a pinched-off vesicle, and [exocytosis](#), where internalized material is transported out of the cell. Proteins in membranes are distributed randomly as a mosaic and diffuse laterally within the plane of the membrane with rates that are 2 to 5 orders of magnitude lower than in aqueous solution. This process has been visualized by fusing tissue culture lines of human and mouse cells with [Sendai virus](#). The fused cells have both mouse and human surface antigens, and the movement of both may be followed by attaching **fluorescent**-labeled [antibodies](#). Within 40 min after the cells are fused, the parent proteins, initially separated, are almost completely mixed and uniformly distributed over the surface of the fused cells (6). Lipid components of membranes also diffuse across the surface of the cell; [diffusion](#) coefficients of the order of  $10^{-8}\text{cm}^2/\text{s}$  have been measured (7). Exo- and endocytosis and membrane fusion may also involve processes in which specific lipids are recruited from the plasma membrane to form **vesicles** (8).

Lipid bilayer models of cell membranes (see [Monolayer](#) and [Liposomes](#)) ostensibly separate the physical and structural properties of the membranes from the biochemical processes superimposed by metabolism. The membrane bilayer is unilamellar, and within this ubiquitous biological structure integral proteins are distributed randomly, free to diffuse over the surface. These properties are embodied in the [fluid-mosaic model](#), which assumes that the lipid bilayer is unilamellar and incorporates the concept of a viscous lipid matrix for dissolved globular proteins that diffuse laterally within the membrane. A more detailed view of the fluid membrane structure is described in the

*Critical Unilamellar State Model*, which proposes that the membrane bilayer is a unique state that assembles and is stable only at a critical point, the physiological temperature,  $T_p$  (9). From this perspective, the membrane bilayer assembles spontaneously from cytoplasmic lipid metabolic pools maintained by the cell at a critical composition, and it illustrates the properties that are characteristic of critical states (10). This model is supported by findings that in large unilamellar vesicles, the membrane bilayer structure forms spontaneously only at a critical temperature that depends on lipid composition and exhibits specific heats and mechanical properties that are found only at this temperature (11). For the total lipid extracts of a wide variety of cellular systems, the critical temperature for assembly of the unilamellar structure is the physiological temperature,  $T_p$ , of the cell from which the membrane lipids are removed (9, 11). Some of the cells for which this phenomenon has been observed include bacteria ( $T_p = 20^\circ$  to  $60^\circ\text{C}$ ), human erythrocytes ( $T_p = 37^\circ\text{C}$ ), brain tissue (squid,  $T_p = 16^\circ\text{C}$ ; rat,  $T_p = 39^\circ\text{C}$ ; human,  $T_p = 39^\circ$  to  $40^\circ\text{C}$ ), and hamster synaptosomes ( $T_p = 37^\circ\text{C}$ ). The concept of a membrane bilayer that assembles and is stable only at the physiological temperature has been utilized as the basis of a theory of neurodegeneration (12).

### Bibliography

1. J. A. F. Op den Kamp. (1979) *Annu. Rev. Biochem.* **48**, 47–71.
2. A. J. Verkleij et al. (1973) *Biochem. Biophys. Acta*, **323**, 178–193.
3. A. Zachowski et al. (1986), *Biochemistry* **25**, 2585–2590.
4. P. Siekevitz (1972) *Annu. Rev. Physiol.* **34**, 117–140.
5. B. Lewin (1994) *Genes V*, Oxford University Press, Oxford, U.K., Chapter "11"–"12". (Contains an overview of the trafficking of proteins within cells, along with an overview of the mechanisms of protein insertion into membranes.)
6. L. D. Frye and M. Edidin (1970) *J. Cell Sci.* **7**, 319–335.
7. M. Edidin (1991) In *The Structure of Biological Membranes*, (P. Yeagle, ed.), CRC Press, Boca Raton, FL, pp. 539–572.
8. C. S. Chen, O. C. Martin, and R. E. Pagano (1997) *Biophys. J.* **72**, 37–50.
9. N. L. Gershfeld (1989) *Biochem. Biophys. Acta Rev. Biomembr.* **988**, 335–350; L. Ginsberg, D. L. Gilbert, and N. L. Gershfeld (1991) *J. Membr. Biol.* **119**, 65–73.
10. J. V. Sengers (1994) In *Supercritical Fluids* (E. Kiran and J. M. H. Levelt Sengers, eds.), Kluwer Academic, Netherlands, pp. 231–271.
11. K. Tajima and N. L. Gershfeld (1985) *Biophys. J.* **47**, 203–220; N. L. Gershfeld (1989) *J. Phys. Chem.* **93**, 5256–5261; N. L. Gershfeld and L. Ginsberg (1997) *J. Membr. Biol.* **156**, 279–286.
12. N. L. Gershfeld and L. Ginsberg (1995) *Rev. Neurosci.* **6**, 1–13.

### Suggestion for Further Reading

13. M. A. Yorek (1993) In *Phospholipids Handbook* (G. Ceve, ed.), Marcel Dekker, New York, Chapter "21". (Contains extensive tables of the phospholipid composition of membranes from a wide variety tissues and cells.)

### Memory Cells

Immunological memory is a basic feature of the immune system. It originates simply because when an [immunogen](#) is given for the first time, [antibody](#) production starts slowly, consisting first of

[IgM](#) that is progressively replaced by [IgG](#) after [class switching](#) of the [isotype](#). When the same [antigen](#) is given a second time, the antibody response is much faster, is maintained longer, and results exclusively in IgG production. Immunological memory is of course the basis for **vaccination**. This specific response indicates that some memory mechanism is operating. This phenomenon is characteristic of **T-cell-dependent** antigens, so the basis for immunological memory might be found in both the T- and **B-cell** populations. Furthermore, this recall effect has been reported for pure T-cell responses, such as delayed-type hypersensitivity, reinforcing the idea that the T-cell compartment also has memory. Are there T and B memory cells, or is immunological memory the result of a systemic organization of the immune system?

An obvious marker that might be expected from a memory cell is that it expresses a [repertoire](#) different from naive unstimulated cells. This is clearly the case for B cells that immediately produce IgG antibodies in secondary responses. Memory B cells arise in germinal centers, in the course of a primary immunization. Once stimulated B-cell clones have switched to IgG, they start to accumulate [somatic hypermutations](#). *In situ* selection by antigen ensures emergence of clones with the highest affinities, which may evolve either to plasma cells that will produce circulating antibodies or to memory cells, which may be endowed with a long lifespan. Memory cells are IgM<sup>-</sup>IgD<sup>-</sup> and have accumulated somatic mutations that can be identified by single-cell gene sequence determination in germinal centers. In many occasions, however, reports have claimed that memory cell survival was strictly dependent on persistence of antigen, which may be maintained as immune complexes on follicular dendritic cells in lymph nodes. Finally, it should be mentioned that the existence of separate lineages of B cells has been claimed, one of which is devoted specifically to expand as memory cells.

Identification of memory T cells is not based on a change in repertoire, because T cells do not use the class switch nor the somatic hypermutation processes. Memory T cells have nevertheless been identified, because they differ from naive T cells by an increased density of surface markers that essentially include [cell adhesion molecules](#) such as LFA-1, CD2, LFA-3, ICAM-1, and several [integrins](#) of the VLA family. Various forms of CD45, a surface phosphatase, also distinguish naive T cells (CD45RA) from memory T cells (CD45RO). Memory cells presumably result from a clonal expansion that occurs during primary immunization. Expression of adhesion molecules in increased amounts facilitates migration of memory T cells to diverse compartments of the organism, where they may encounter **antigen presenting** cells, such as Langerhans cells in the skin, thus favoring in this case the emergence of delayed-type hypersensitivity, a pure T-cell response. Memory T cells have also been reported to have a long lifespan.

Despite the arguments supporting the existence of both B and T memory cells, one cannot exclude the possibility that immunological memory might be the result of a more complex system organization that would be inspired most from the central nervous system memory. More data will be required to clarify this issue.

See also entries [Immune response](#), [Immunoglobulin](#), [Immunization](#), and [Somatic hypermutation](#).

#### Suggestions for Further Reading

- N. R. Klinman (1998) Repertoire diversification of primary vs memory B cell subsets. *Curr. Top. Microbiol. Immunol.* **229**, 133–148.
- R. W. Dutton, L. M. Bradley, and S. L. Swain (1998) T cell memory. *Annu. Rev. Immunol.* **16**, 201–223.
- T. Manser, K. M. Tumas-Brundage, L. P. Casson, A. M. Giusti, S. Hande, E. Notidis, and K. A. Vora (1998) The roles of antibody variable region hypermutation and selection in the development of the memory B-cell compartment. *Immunol. Rev.* **162**, 183–196.

## Mendelian Inheritance

Mendelian inheritance patterns occur when (i) **alleles** in heterozygous **diploids** have equal chances of being transmitted to [haploid](#) gametes, thus giving rise to a 1:1 segregation, (ii) alleles of different loci combine at random in the gametes, and (iii) gametes fuse at random to form diploids.

The name *Mendelian* honors Gregor Mendel, who observed and counted carefully the phenotypes produced in many crosses of pea plants, *Pisum sativum* (1). This plant was convenient because of its relatively short generation time, easy handling, and abundant offspring. He tried other plants as well; and he found confirmatory evidence with some of them, and results that could not be understood at the time with others—for example, *Hieracium*.

Mendelian inheritance is sometimes expressed in the form of Mendel's laws, although we owe this formulation largely to later writers, such as C. Correns and H. de Vries. The first law is the law of segregation (of alleles to gametes); the second is the law of independent assortment (of alleles of different genes). The uniformity of the hybrid generation is sometimes considered a law as well.

Two underlying assumptions of Mendelian inheritance were introduced before Mendel. One is the *individuality* and *autonomy* of inherited traits: Many traits present clear alternatives that are preserved in successive generations without blending (see [Genetic Marker](#)). The other is the *combination* of separate traits with each other. These assumptions were expressed clearly by Michel Sageret (2) following his observations with melons, *Cucumis (Cucurbita) melo*. He found “marked combinations of various characters without any mixture between them” and understood that this gives, nature “the ability to vary infinitely its products and avoid monotony.” The combinations of individual traits superseded the traditional holistic opinion for which every aspect of living beings had to be seen as a manifestation of a whole.

It is not generally known that Sageret came close to summarizing Molecular Biology, using in the figured sense terms from the arts of typography and casting, as our textbooks do today: “To what is due this ability of nature to reproduce in the descendants characters that belonged to their ascendants? We do not know, but we may well suppose that it depends on a type, on a primitive mold that contains the germ of all the organs, germ that sleeps and awakes, and develops or not according to the circumstances”.

The rules of Mendelian inheritance provide quantitative forecasts for the frequency of different traits and their combinations in the offspring of defined genetic **crosses**. Consider a pair of genetic markers, phenotypes “A” and “a,” governed by the alleles  $A$  and  $a$  of a single gene (see [Genetic Marker](#)), and two diploid parents that “breed true” [ie, are homozygous for the alleles under consideration (see [Homozygote](#), [True Breeding](#))]; the parental genotypes are written  $AA$  and  $aa$  to indicate the fact that diploids carry two copies of each allele (Table 1). Their gametes carry alleles  $A$  and  $a$ , respectively, and the  $F_1$  diploids produced by their fusion will be  $Aa$ . The phenotype of the offspring will be uniform and will depend on the dominance or recessivity of the alleles; assuming that  $A$  is dominant over  $a$ , the offspring phenotype will be “A.” The gametes of  $F_1$  individuals will have equal chances of carrying the  $A$  or the  $a$  allele; this is conveniently represented by an algebraic sum in which each allele takes its probability as a coefficient. The fusion of gametes of  $F_1$  individuals produce a heterogeneous  $F_2$  segregation: The random combination of gametes is well represented by the multiplication of their respective algebraic expressions (the squaring, in this case) and yields an algebraic expression of all genotypes each with its own probability as coefficient. The offspring of other crosses—for example, a backcross—is easily predicted by similar considerations.

**Table 1. Mendelian Segregation of a Monohybrid (a Dominant Allele  $A$  and a Recessive Allele  $a$ )**

|                      |                               |  |      |
|----------------------|-------------------------------|--|------|
| Generation P1:       | $AA$                          | $\times$   | $aa$ |
| Phenotypes:          | “A”                           |  | “a”  |
| Their gametes:       | $A$                           |  | $a$  |
| Generation F1:       |                               | $Aa$   |      |
| Phenotype:           |                               | “A”  |      |
| Their gametes:       |                               | $\frac{1}{2}A + \frac{1}{2}a$  |      |
| Generation F2:       |                               | $(\frac{1}{2}A + \frac{1}{2}a) \times (\frac{1}{2}A + \frac{1}{2}a) = \frac{1}{4}AA + \frac{1}{2}Aa + \frac{1}{4}aa$ |      |
| Their phenotypes:    |                               | $\frac{3}{4}$ “A” + $\frac{1}{4}$ “a”  |      |
| Testcross of the F1: | $Aa$                          | $\times$   | $aa$ |
| Phenotypes:          | “A”                           |  | “a”  |
| Their gametes:       | $\frac{1}{2}A + \frac{1}{2}a$ |  | $a$  |
| Generation R2:       |                               | $\frac{1}{2}Aa + \frac{1}{2}aa$  |      |
| Their phenotypes:    |                               | $\frac{1}{2}$ “A” + $\frac{1}{2}$ “a”  |      |

This algebraic treatment can be generalized to several genes, each with their respective alleles (Table 2). Again, multiplication of polynomials expresses the random combinations of the alleles of different genes.

**Table 2. Mendelian Segregation of a Multihybrid ( $n$  Pairs of Alleles,  $A^1_i$  and  $A^2_i$ )**

|                |                                       |   |
|----------------|---------------------------------------|---|
| Generation P1: | $A^1_1 A^1_1 \frac{1}{4} A^1_n A^1_n$ | $A^2_1 A^2_1 \frac{1}{4} A^2_n A^2_n$                     |
| Their gametes: | $A^1_1 \frac{1}{4} A^1_n$             | $A^2_1 \frac{1}{4} A^2_n$                                 |
| Generation F1: |                                       | $A^1_1 A^2_1 \frac{1}{4} A^1_n A^2_n$                     |
| Their gametes: |                                       | $\prod_{i=1}^n (\frac{1}{2} A^1_i + \frac{1}{2} A^2_i)$   |
| Generation F2: |                                       | $\prod_{i=1}^n (\frac{1}{2} A^1_i + \frac{1}{2} A^2_i)^2$ |

Mendelian inheritance is a consequence of the phenotypic traits being governed by individual genes located on the chromosomes. It is valid for genes located in the *autosomes* (chromosomes present in two copies in a diploid nucleus) and must be modified in the cases of imperfect diploidy, such as sex linkage and [aneuploidy](#). The random combinations of different traits apply only when the responsible genes are located on different chromosomes or at sufficient distance in the same chromosome. Genes

located sufficiently close to each other in the same chromosome tend to be transmitted to the gametes as a block.

The random segregation of alleles to the gametes does not occur in the presence of *segregation distorters*. Thus, gametes with certain changes in the structure or the number of chromosomes are formed less often than expected, particularly in the meiosis of females. Some genetic elements actively bias the segregation to their own favor. This phenomenon is called *meiotic drive*. Some meiotic-drive elements fail to impose themselves on the whole population because they are held in check by opposing selection pressures.

Independent assortment of chromosomes does not occur when there is *affinity* between some nonhomologous chromosomes, so that they tend to move together during meiosis. The result is an apparent linkage between alleles of different chromosomes.

### Bibliography

1. G. Mendel (1865, published 1866) Verh. Naturforsch. Ver. Brünn, Abhand. **4**, 3–47.
2. M. Sageret (1826) Ann. Sci. Nat. (Paris) **8**, 294–314.

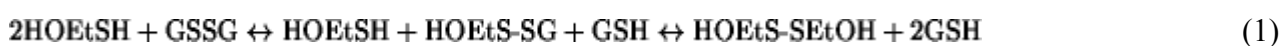
### Suggestions for Further Reading

3. V. Orel (1996) *Gregor Mendel, the First Geneticist*, Oxford University Press, New York.
4. A. H. Sturtevant (1965) *A History of Genetics*, Harper and Row, New York.
5. J. F. Crow (1991) Why is Mendelian segregation so exact? *Bioessays* **13**, 305–312.
6. T. W. Lyttle (1993) Cheaters sometimes prosper: distortion of mendelian segregation by meiotic drive. *Trends Genet.* **9**, 205–210.
7. K. G. Ardlie (1998) Putting the brake on drive: meiotic drive of *t* haplotypes in natural populations of mice. *Trends Genet.* **14**, 189–193.

## b-Mercaptoethanol

b-Mercaptoethanol has the simple structure HO—CH<sub>2</sub>—CH<sub>2</sub>—SH, abbreviated here as HOEtSH, and its [thiol group](#) is of typical reactivity and ionizes with a pK<sub>a</sub> of about 9.6. It is one of the most extensively used reagents for protecting biological materials, especially [proteins](#), from oxidation. It has the advantage of being a stable, colorless liquid that is readily miscible with [water](#). Its thiol group will react with oxidants and heavy metals, thereby protecting biological materials, and it can also be used to keep the thiol groups of proteins and other macromolecules in the reduced form. It is frequently used to reduce [disulfide bonds](#) in biological materials, especially in preparation for [SDS-PAGE](#).

Mercaptoethanol has limitations in this respect that should be recognized. Its thiol group is not a very potent reductant, very similar to that of [glutathione](#). The equilibrium constant for **thiol-disulfide exchange** between these two molecules,



is approximately unity. The equilibrium constant for the reduction of stable disulfide bonds in folded proteins,



$$K_{eq} = \frac{[P_{SH}^{SH}] [HOEtS-SEtOH]}{[P_S^S] [HOEtSH]^2} \quad (3)$$

has been measured with glutathione and will be about the same for mercaptoethanol; values in the region of less than  $5 \times 10^{-3}$  have been measured for the overall reaction of Eq. 3 (1). Therefore, a large excess of mercaptoethanol will be required to reduce stable protein disulfide bonds.

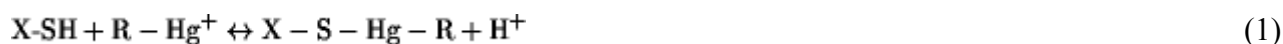
A further problem is that mercaptoethanol becomes air-oxidized readily, especially at alkaline pH values, and the disulfide form is a potent oxidant (Eq. 2 and 3). If only 1% of a 10 mM solution of mercaptoethanol becomes oxidized, it will be unable to reduce stable disulfide bonds like those previously described. In the presence of significant amounts of the disulfide form of mercaptoethanol, mixed disulfides of the mercaptoethanol with protein thiol groups will be significant. The reagents **dithiothreitol** and **dithioerythritol** are much more powerful in reducing disulfide bonds because they form a stable, intramolecular disulfide bond; the disulfide bond of HOEtSSEtOH can be considered to be intermolecular because two molecules are produced when the disulfide bond is reduced. A solution in which 99% of dithiothreitol has become oxidized will be no less potent at reducing protein disulfide bonds than the above solution of mercaptoethanol that was only 1% oxidized.

#### Bibliography

1. T. E. Creighton and D. P. Goldenberg (1984) J. Mol. Biol. **179**, 497–526.

#### Mercurials

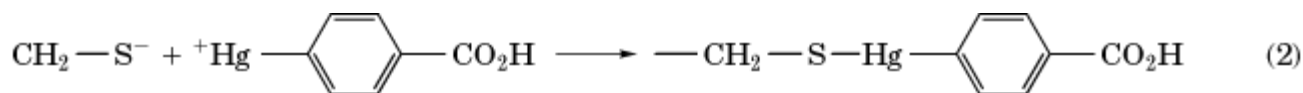
The medicinal use of mercurials can be traced back over 3000 years, although their modern therapeutic applications began with the discovery of the diuretic action of mercurous chloride in 1849. When Zeise discovered  $C_2H_5SH$  in 1834, he called it “mercaptan” (corpus mercurium captans) because the formation of mercury derivatives was a striking characteristic (1). The aryl mercurials, such as *p*-mercuribenzoate, were introduced as enzyme inhibitors in 1937, and the alkyl mercurials, such as methylmercuric chloride, in 1950. Organomercuric ions of the type  $R-Hg^+$  comprise a useful family of reagents that react specifically with thiol groups and can be used as titrants, reactivity probes, or reporter groups, depending upon the characteristics of group R. Such organomercuric ions are monovalent and react with a single thiol group:



In contrast, mercuric ion,  $Hg^{2+}$ , tends to react with two thiol groups and to cross-link them. The reaction with either is reversible, but the binding is tight, so that the equilibrium generally lies far to the right, unless an excess of another thiol reagent or of cyanide is added, when the original thiol group can be regenerated. A slight complication is that mercuric ions exist in a variety of loose complexes with components of the aqueous solution, such as hydroxide and chloride ions (2).



The classical organomercurial reagent was *p*-mercuribenzoate (3):



which is often called *p*-chloromercuribenzoate, although the chloride component is simply a counterion, and it makes little difference what salt is used. It has useful spectroscopic properties, and its absorbance increases at 250 nm upon reaction with a thiol group. Many variants have been prepared and used subsequently. Nevertheless, such reagents are not widely used now, in part because the adduct is not irreversible, and the original reagents are always present in significant quantities, due to the reverse of Eq. (1). Such mercurial reagents can always react with any other constituents subsequently added. Ions such as hydroxyl and chloride, especially if present at high concentration, can also compete with thiol groups for the mercurial. Furthermore, the C—Hg bond is relatively weak, so Hg<sup>2+</sup> can exchange rapidly with the alkyl mercurial, and inorganic mercury can be split off rather readily.

Such reactions of cysteine thiol groups with mercurials are the most obvious, and one of the more useful, ways to make heavy-atom derivatives for [X-ray crystallography](#) determination of protein structure (see [Isomorphous Replacement](#)).

#### Bibliography

1. E. M. Reid (1958) *Organic Chemistry of Bivalent Sulfur*, vol. I, Chemical Publishing Co., New York.
2. R. G. Khalifah, G. Sanyal, D. J. Strader, and W. McI. Sutherland (1979) *J. Biol. Chem.* **254**, 602–604.
3. P. D. Boyer (1954) *J. Amer. Chem. Soc.* **76**, 4331–4337.

#### Suggestion for Further Reading

4. J. L. Webb (1966) *Enzyme and Metabolic Inhibitors*, Vol. II, Chapter "7", Academic Press, New York.

## Meroblastic Cleavage

Meroblastic cleavage is the incomplete cellular division found in [embryos](#) with large amounts of yolk (1). In embryos with little yolk, the early cleavage furrows can completely separate the daughter cells, or [blastomeres](#). In many vertebrate embryos, the large amounts of yolk prevent the cleavage furrows from completely separating the daughter blastomeres. In these embryos, the cytoplasm is usually segregated to one end of the embryo and sits on top of the yolk. The cleavage furrows begin in the cytoplasm, but only proceeds part of the way through the yolk. The nuclei and associated cytoplasm are separated on one surface, but remain connected to the yolk (and to each other) at the interior. Later divisions tangential to the original cleavage furrows completely separate the cells in the dorsal and middle portions of the cytoplasmic region, while the cells more ventral and lateral continue to be connected cytoplasmically to the yolk and each other. Eventually, a disc of blastomere cells is formed that sits on top of the yolk, with the nearest cells still connected to the yolk. Meroblastic cleavage is common among birds (eg, the domestic chicken embryos that have been used extensively for embryological studies), reptiles, bony fishes, and amphibians.

## Bibliography

1. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 111–112.

## Suggestions for Further Reading

2. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 111–112.
3. J. C. Gerhart and M. W. Kirschner (1997) *Cells, Embryos, and Evolution: Toward a Cellular and Developmental Understanding of Phenotypic Variation and Evolutionary Adaptability*, Blackwell Science, Malden, MA.
4. S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
5. S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
6. J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, U.K.

## Merodiploids

A merodiploid is an essentially [haploid](#) organism that carries a second copy of a part of its [genome](#). The term is derived from the Greek, meros = part, and was originally used to describe both unstable partial diploidy, such as that which occurs briefly in recipients after mating with an Hfr strain ([1](#)), and the stable state, exemplified by F-prime strains (see [Hfr'S And F-Primes](#)). Over time, usage has tended to confine the term to descriptions of stable genetic states.

The inability of bacteria to maintain more than a single copy of their [chromosome](#) makes it easier to isolate [mutations](#). Most of the mutations useful to a geneticist are recessive, and their effects would be hidden if a second unmutated gene copy were present to provide wild-type function, as in eukaryotes. Because bacteria are haploid, mutations are expressed directly as mutant **phenotypes**, and mutants are readily isolated by straightforward selection and screening procedures, vastly expanding the range of genes attainable by classical genetic methods. Beyond this point, however, the haploid state is inconvenient, because analyzing the function of a gene identified by a mutation usually involves determining how the mutated function interacts with the wild-type one or with related mutants, for which diploidy is needed. Merodiploidy provides a solution to this problem.

The discovery of natural merodiploidy arrived happily at a moment when it would make one of its most notable contributions. Hfr strains in which the integrated F factor, or [F plasmid](#), is situated near the *lac* genes were found to give rise to F<sup>'</sup>-*lac*<sup>+</sup> **plasmids**, by excision of the *lac* region along with the F DNA ([2](#), [3](#)). When an F'*lac* transferred itself into a new cell, it rendered the *lac* gene region of the chromosome diploid. At this time, François Jacob and his colleagues had proposed that *lac* gene expression is regulated by the action of a [repressor](#) at a control locus, the **operator**, situated next to the *lac* genes. Seeking evidence for the operator, they used an F'*lac*/*lac*<sup>+</sup> strain to isolate mutants that expressed the *lac* genes even in the presence of the repressor. The mutations caused unregulated expression only of the *lac* genes to which they were linked (**cis-acting**) and not of *lac* genes in **trans**, just the property expected of the hypothetical *cis*-acting control locus ([4](#)).

The advantages offered by merodiploidy were well-illustrated by an analysis of *E. coli* mutants unable to make flagella (5). The *fla* mutants were first mapped to three clusters by complementation with F-primes carrying defined regions of the chromosome. The mutations were then crossed onto the F-primes by [recombination](#), and the resulting F-primes transferred into each mutant to determine which pairs of mutations complemented to allow synthesis of flagella. These tests provided an estimate of the minimum number of *fla* genes in each region. The literature of bacterial genetics is replete with similar examples.

The uses of merodiploidy are limited to neither *E. coli* nor plasmids. The genes for catechol metabolism in *Pseudomonas aeruginosa* were mapped by first allowing recombinational transfer of [transposon](#) insertion mutations from R-prime plasmids to the chromosome and then carrying out [complementation](#) analysis using overlapping **cosmid** clones (essentially, smaller prime plasmids) (6). Potassium transport mutants of *E. coli* were analyzed by infection with [lambda phage](#) for **transduction** and plating on potassium-deficient media: Formation of turbid plaques containing lysogens occurred only for complementing mutants (7).

The existence of natural stable merodiploid states is one reason for raising the question, Why are bacteria haploid? There are others. Partial diploidy of the F-prime type is usually well-tolerated, suggesting that doubling the gene dosage is not generally deleterious. In several bacteria, the genome is distributed over more than one DNA molecule, implying the ability to coordinate replication and gene expression of physically distinct chromosomes. Bacteria have no trouble maintaining oligo- and multicopy plasmids in a polyploid state and, when growing fast, maintain their chromosomes in a state of pseudopolyploidy owing to the formation of multiforked chromosomes. J.-P. Bouché and his colleagues (8) produced a conditionally diploid *E. coli* by reducing synthesis of the essential cell division protein FtsZ, thereby moving the replication cycle forward relative to division.

Nevertheless, haploidy is the invariable ground state. Could it be that the advantage that haploidy confers on the bacterium is the same as that offered to the geneticist? Whereas multicellular eukaryotes require the phenotypic stability of diploidy on which to build a developmental program and to prevent unregulated growth, bacteria have a different imperative: to adapt rapidly to diverse environmental and chemical challenges by generating from within their enormous populations mutations that are expressed immediately.

## Bibliography

1. E. L. Wollman, F. Jacob, and W. Hayes (1956) Cold Spring Harbor Symp. Quant. Biol. **21**, 141.
2. E. A. Adelberg and S. N. Burns (1960) J. Bacteriol. **79**, 321–330.
3. F. Jacob and E. A. Adelberg (1959) C. R. Acad. Sci. **249**, 189–191.
4. F. Jacob, D. Perrin, C. Sanchez, and J. Monod (1960) C. R. Acad. Sci. **250**, 1727–1729.
5. M. Silverman and M. Simon (1973) J. Bacteriol. **113**, 105–113.
6. C. Zhang and B. W. Holloway (1992) J. Gen. Microbiol. **138**, 1097–1107.
7. D. C. Dosch, G. L. Helmer, S. H. Sutton, F. F. Salvacion, and W. Epstein (1991) J. Bacteriol. **173**, 687–696.
8. F. Tétart, R. Albigot, A. Conter, E. Mulder, and J.-P. Bouché (1992) Mol. Microbiol. **6**, 621–627.

## Messenger RNA

Messenger RNA (mRNA) is a fully processed RNA containing a particular open reading frame (see

[Gene Structure](#)) that is **translated** by the [ribosome](#) into a specific [polypeptide chain](#). In **eukaryotes**, all mRNAs carry a [cap](#) (1), and most are **polyadenylated** (2); both of these modifications can contribute to the translational efficiency of the mRNA. Normally, mRNA have the general structure from 5' to 3': 5' cap and 5' untranslated region [5'UTR], open reading frame [ORF], 3'UTR, and [poly A](#) tail. The mRNA are synthesized from the DNA sequence of the gene as a precursor, pre-mRNA, by DNA-dependent **RNA polymerase II**, which must be processed in order to generate the mature mRNA. This processing usually involves the removal of noncoding **introns**, which are specifically removed in the [spliceosome](#) (3) (see [Gene Splicing](#)), and most mRNAs are polyadenylated cotranscriptionally (2) (see [Polyadenylation](#)). Translation of the ORF initiates normally at an AUG codon, which encodes the initiating [methionine](#) residue (see [Initiation Codon](#) and [Translation](#)). Translation of the ORF is then terminated when the ribosome encounters a [stop codon](#).

### Bibliography

1. M. Salditt-Georgieff, M. Harpold, S. Chen-Kiang, and J. E. Darnell, Jr. (1980) *Cell*, **19**, 69–78.
2. J. D. Lewis, S. I. Gunderson, and I. W. Mattaj (1995) *J Cell Sci Suppl.* **19**, 13–19.
3. M. Moore, C. Query, and P. Sharp (1993) in R. Gesteland and J. Atkins, eds., *Splicing of Precursors to Messenger RNAs by the Spliceosome*, CSH Laboratory Press, New York, pp. 303–357.

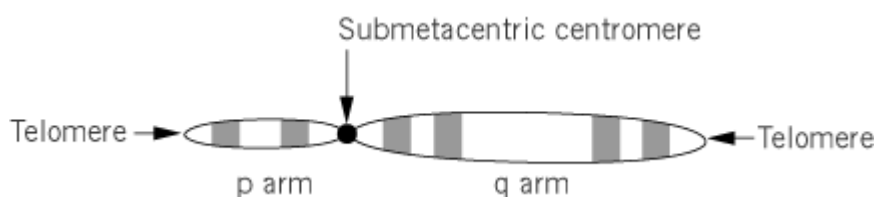
### Suggestion for Further Reading

4. J. D. Lewis and E. Izaurralde (1997) The role of the cap structure in RNA processing and nuclear export, *Eur. J. Biochem.* **247**, 461–469.

## Metacentric Chromosome

A *metacentric chromosome* has its [centromere](#) in or near its middle. Robertsonian translocations occur when two [acrocentric chromosomes](#) fuse to form a single metacentric (see Fig. 1 of [Acrocentric Chromosome](#)). Dissociation occurs when a metacentric chromosome breaks at the centromere to form two acrocentrics. This type of chromosome fusion and breakage is a major source of variation in chromosome numbers within various taxa. Species with high numbers of metacentric chromosomes tend to have fewer chromosomes than those with high numbers of acrocentrics. Metacentric and submetacentric chromosomes (where the centromere is near the middle) have two arms and appear as V- or J-shaped structures in metaphase preparations. The shorter arm of a submetacentric is called the 'p' arm, and the longer arm is the 'q' arm (Figure 1).

**Figure 1.** The basic organization of a mammalian chromosome, indicating p and q arms, telomeres, and centromeres.



### Suggestion for Further Reading

R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A synthesis*, Wiley-Liss, New York.

## Metalloproteinase Inhibitors, Protein

The classical vertebrate digestive enzymes carboxypeptidases A and B are  $\text{Zn}^{++}$ -dependent metalloproteinases. Along with many other proteinases, they are inhibited by their propeptides. In addition, carboxypeptidase A inhibitors from potatoes and from the roundworm *Ascaris lumbricoides* were intensively studied. The potato inhibitor is a single polypeptide chain of 39 residues crosslinked by three disulfide bridges. The disulfides bond pattern corresponds to that in the 30 residue squash inhibitors of serine proteinases. A small number of recombinant changes endowed the same protein molecule with trypsin inhibitory activity, located near the inhibitor's  $\text{NH}_2$  terminus and with carboxypeptidase inhibitory activity located near the COOH terminus. The three-dimensional structure of the carboxypeptidase A–potato inhibitor complex reveals that, among others, the  $P_3$ ,  $P_2$  and  $P_1$  residues make strong contacts with the enzyme. The original COOH terminal residue ( $P_1'$ ), which is a Gly, is hydrolyzed off but not fully released. It remains attached to the complex.

A 102-residue inhibitor of metalloproteinases was isolated from *Streptomyces nigrescens*. It appears to inhibit thermolysin and related endoproteinases according to the standard mechanism employed by most protein inhibitors of serine proteinases.

A large number of metalloendoproteinases, the matrix metalloproteinases, or MMPs, are known. They are involved in remodeling of tissues and invasion of cells both in normal and in pathological processes. Such enzymes are synthesized as inactive zymogens that require the release of the propeptide for activity. This release is frequently but not always autocatalytic. The activity is also controlled by a number of closely related protein inhibitors. The most studied of these are called tissue inhibitors of metalloproteinases, or TIMPs. At the time of this writing, TIMPs 1,2,3, and 4 have been described. They differ from one another by being specific for different members of the large MMP set. A striking feature of TIMPs is that they form complexes that block the actual or potential activity of both the MMP enzymes and zymogens. The inhibited zymogens are no longer capable of autoactivation. Aside from the four TIMPs, a number of other protein inhibitors of MMPs have been reported.

Recently, a three-dimensional structure of a complex of matrix metalloproteinase, MMP3, or human stromelysin and TIMP-1 was reported. The  $\text{NH}_2$  terminal Cys<sup>1</sup>, which also forms the Cys<sup>1</sup>-Cys<sup>70</sup> intramolecular disulfide serves as the  $P_1$  residue in TIMP. This residue is placed just above the catalytic  $\text{Zn}^{++}$  of MMP3. It forms two close (2.0Å) contacts with  $\text{Zn}^{++}$  involving its N and O atoms. The Thr<sup>2</sup> residue that serves as  $P_1'$  is embedded in a large hydrophobic  $S_1'$  pocket of the enzyme. There is strong reason to suppose that this mode of interaction occurs in all or most MMP - TIMP complexes.

### Suggestions for Further Reading

W. Bode and R. Huber (1992) Natural protein proteinase inhibitors and their interaction with

proteinases. *Eur. J. Biochem.* **204**, 433–451.

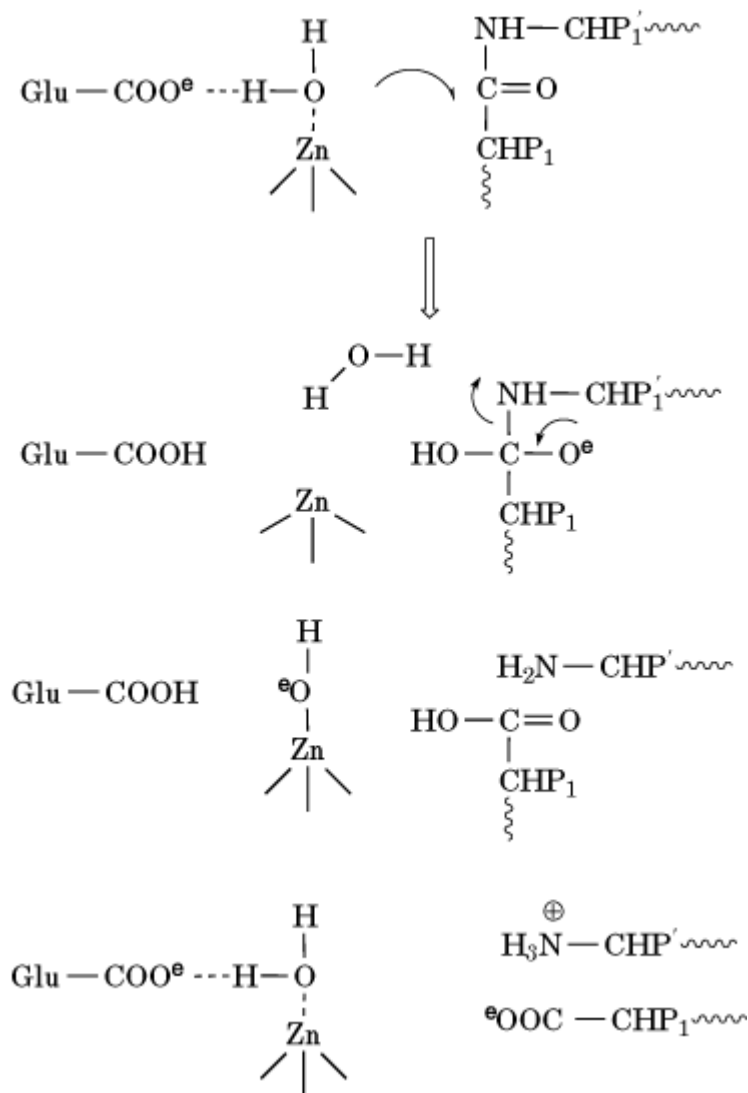
L. M. Coussens and Z. Werb (1996) Matrix metalloproteinases and the development of cancer. *Chem. & Biol.* **3**, 895–904.

F.-X. Gomis-Rüth et al. (1977) Mechanism of inhibition of the human matrix metalloproteinase stromelysin-1 by TIMP-1. *Nature* **389**, 77.

## Metalloproteinases

In addition to the **carboxyl**, **thiol**, and [serine proteinases](#), there is another large group of protein-degrading [enzymes](#) that use a zinc hydroxide ion to help catalyze the hydrolysis of [peptide bonds](#). The Enzyme Commission of the International Union of Biochemistry and Molecular Biology has assigned them the classification E.C. 3.4.24. The zinc ion is bound to the protein by interaction with three amino acid side chains, and the fourth coordination site on the zinc is occupied by a water molecule (1). As the amino acid sequences and three-dimensional structures of increasing numbers of zinc metalloproteinases became available, certain patterns of metal binding emerged. The first to be recognized was the sequence -His-Glu-X-X-His- (where X is any amino acid other than [histidine](#) or [glutamic acid](#)). The two His residues, along with a third residue—typically a glutamic acid separated from the second His by from 21 to 61 intervening residues—constitute the metal binding site (2). Later another group of metalloproteinases was found in which the third residue is another His separated from the second His by only five intervening residues (3). At least three other modes of metal binding have now been identified (4). In those enzymes where a glutamic acid residue is adjacent to the first metal-binding histidine, as in [thermolysin](#), the glutamic acid carboxyl group is thought to activate the metal-bound water molecule by removing a proton and allowing the hydroxide to attack the carbonyl group of the peptide bond that is to be hydrolyzed (Fig. 1). Similar mechanisms may pertain to other zinc metalloproteinases.

**Figure 1.** Schematic representation of the mechanism of action of zinc metalloproteinases. The carboxylate group of the active-site glutamic acid residue potentiates the attack of a zinc-coordinated water molecule on the carbonyl group of the substrate. In a substrate with more than one peptide bond, the particular peptide attacked is determined by the ability of the enzyme to recognize P<sub>1</sub> and P'<sub>1</sub>, which are the principal specificity-determining side chains. The tetrahedral intermediate rearranges to form products, which dissociate, and the active site is ready to act on another substrate.



An important subgroup of these enzymes includes the matrix metalloproteinases. Among them are the enzymes that degrade [collagen](#) (collagenases) and other constituents of the [extracellular matrix](#), hence their collective name. It is noteworthy that these enzymes bind two zinc atoms, one at the [active site](#) and a second that may have a structural role. Two zinc ions are also found in lens [aminopeptidase](#), but in this case they are close together, apparently bridged by the carboxyl side chain of a glutamic acid residue. This would seem to be an example of a co-catalytic site (5), with one metal being important for catalysis and the other for substrate binding (6).

Metalloproteinases play an important role in tissue remodeling, in degenerative diseases such as rheumatoid arthritis, and in cancer. Several snake venom [toxins](#) are metalloproteinases, and a key enzyme in blood pressure regulation, angiotensin converting enzyme, is a metalloproteinase. Classically, these enzymes are recognized by their susceptibility to inhibition by metal-binding agents (chelating agents), such as ethylenediaminetetraacetic acid (**EDTA**) or 1,10-phenanthroline. They are not inhibited by reagents that are specific for carboxyl (ie, [pepstatin](#)), thiol ([leupeptin](#)), or serine ([PMSE](#)) proteinases. A number of metal-binding substrate analogues have been found in nature (eg, phosphoramidon); and others, synthesized in the laboratory, are important pharmaceuticals.

#### Bibliography

1. B. L. Vallee and D. S. Auld (1990) *Biochemistry* **29**, 5647–5659.

2. B. L. Vallee and D. S. Auld (1993) *Acc. Chem. Res.* **26**, 543–551.
3. W. Bode, F. X. Gomis-Ruth, and W. Stocker (1993) *FEBS Lett.* **331**, 134–140.
4. N. M. Hooper (1994) *FEBS Lett.* **354**, 1–6.
5. B. L. Vallee and D. S. Auld (1993) *Biochemistry* **32**, 6493–6500.
6. H. Kim and W. N. Lipscomb (1994) *Adv. Enzymol.* **68**, 153–213.

### Suggestion for Further Reading

7. S. J. Lippard and J. M. Berg (1994) *Principles of Bioinorganic Chemistry*, University Science Books, Mill Valley, CA, pp. 257–271.

## Metalloproteins

Many [proteins](#) bind metal ions, permanently as **prosthetic groups** or more transiently as **ligands**. These metal ions play a variety of roles in metalloproteins: **electron transfer**, maintaining the [protein structure](#), **oxygen binding**, forming coordinated hydroxide radicals, substrate binding, and electrophilic catalysis (see [Metal-Requiring Enzymes](#) and [Metalloproteinases](#)).

During the past 15 years, it has become apparent that metalloproteins also play important roles in regulating the expression of genetic information. Such proteins take up stoichiometric quantities of trace metal ions and then undergo conformational changes that produce marked differences in their abilities to bind to specific sites on the **genomic** DNA or RNA of a given organism. The physiological effects of the binding of metalloproteins to, or release from, **nucleic acids** include increased resistance to heavy metals, control of iron uptake and storage pathways, recognition of packaging signals within the RNAs of [retroviruses](#), control of vertebrate [development](#) events, and recognition of [steroid hormones](#). Minerals present in high concentrations, such as  $\text{Na}^+$ ,  $\text{K}^+$ , and  $\text{Mg}^{2+}$ , can play important roles in stabilizing nucleic acid structures, but they are unlikely to be widely used in gene regulation. The remaining essential metals are normally trace elements. Hence, it is not surprising that zinc, iron, and copper are prominent metalloregulators (see also [Calcium-Binding Proteins](#), [Zinc-Binding Proteins](#), and [Iron-Binding Proteins](#)). The incorporation of a particular trace metal ion into an apoprotein is influenced by its ionic radius, thermodynamic stability, ligand-substitution kinetics, and charge (Table 1). Three to four ligands, usually the side chains of [cysteine](#) and [histidine](#) residues, typically complete the coordination sphere about a bound metal ion.

**Table 1. Coordination Environments Preferred by Common Metal Ions Found in Proteins**

| Metal Ion        | Coordination Number | Geometry    | Ligands            |
|------------------|---------------------|-------------|--------------------|
| $\text{Zn}^{2+}$ | 4                   | Tetrahedral | His, Cys           |
| $\text{Fe}^{2+}$ | 4                   | Tetrahedral | His, Cys, Glu, Asp |
| $\text{Cu}^+$    | 4                   | Tetrahedral | Cys                |
| $\text{Cd}^{2+}$ | 4                   | Tetrahedral | Cys                |



Because of their intrinsic biological functions, regulatory metalloproteins are not normally present in large quantities in cell and tissue homogenates, and most are colorless or only weakly colored. Nonetheless, the study of regulatory metalloproteins has become the fastest growth area of inorganic biochemistry. The term “metalloregulation” refers to regulating the transfer of genetic information by metal ions.

Regulatory metal-binding proteins are generally specific for a particular trace metal ion. For example, a synthetic **zinc finger** analogue (1) displays a strong preference for Zn<sup>2+</sup> over Mn<sup>2+</sup>, Fe<sup>2+</sup>, Co<sup>2+</sup>, Ni<sup>2+</sup>, and Cd<sup>2+</sup>.

## 1. Characterization

As with any metalloprotein, the metal stoichiometry is of paramount importance. It is usually determined by titration or atomic absorption spectrometry. Several thousand putative metalloregulatory proteins have been reported in the literature, usually on the basis of limited sequence patterns alone. In the absence any physical characterization, however, such reports must be viewed as speculative because they offer no real evidence regarding the metal element involved.

The nature of the metal coordination sphere can be probed by electronic absorption spectroscopy in cases involving iron or copper coordination. Zn<sup>2+</sup> presents a challenge because it is colorless, but the Zn<sup>2+</sup> can be replaced with Co<sup>2+</sup>, a chromophoric probe that possesses a comparable ionic radius, is kinetically labile, and prefers similar coordination environments.

X-ray absorption spectroscopy is particularly useful for investigating the local environment (<~5) of the metals in these proteins (2). When X rays are absorbed by metal atoms, they liberate electrons that are backscattered from neighboring (ie, ligand) atoms. This interference phenomenon provides a very precise measurement of the distance to neighboring atoms. The edge region of the spectrum provides information about the valence state and coordination geometry of the metal and the identity of neighboring atoms. The extended X-ray absorption fine structure (**EXAFS**) region of the spectrum contains information about the number and average distances of neighboring atoms and about their relative disorder (see **EXAFS**).

The structural methods do not address key issues regarding the way a regulatory protein acts as an information transducer. [Electrophoresis](#) experiments, including hydroxyl-radical [footprinting](#) and **mobility shift assays**, can indicate what region of the DNA or RNA is recognized and whether protein binding induces a structural change in the nucleic acid. More precise structural information relating to the nature of the protein–nucleic acid contacts requires their cocrystallization and **X-ray crystallographic** structural determination. Solution nuclear magnetic resonance (**NMR**) is also used in structural determination of small metalloproteins, but it has not yet been used to study complexes with DNA or RNA. Approximately 50 structures of metalloproteins have been deposited to date in the Brookhaven Protein Data Bank (see [Structure Databases](#)).

## 2. Regulatory Metalloproteins

The major regulatory metalloproteins are described in [Zinc-binding proteins](#), **Zinc fingers**, [Iron-binding proteins](#), [Iron-response elements](#), and **Metal-response elements**. In addition, organisms must frequently cope with toxic elements in the environment. Bacteria are especially rich sources of metalloregulators and encode resistance systems (3) for many toxic metal ions, including Hg<sup>2+</sup>, Tl<sup>+</sup>,

$\text{Ag}^+$ ,  $\text{AsO}_2^-$ ,  $\text{Cd}^{2+}$ , and  $\text{Cu}^{2+}$ . Heavy-metal resistance is achieved in various ways, including sequestration in proteins, such as **metallothioneins**, efflux “pumping,” and reduction to volatile metal atoms (eg,  $\text{Hg}^0$ ).

### 3. Concluding Remarks

In contrast to the surging interest in structural studies on metalloproteins, the metabolic effects of dietary metal supply on the availability of metalloregulatory proteins have not received nearly as much attention as is warranted. Sorting out these effects is complicated by the presence of trace metal ions in structural proteins and enzymes. For example, zinc is found in **RNA polymerase** and in zinc fingers.

The sensing of small molecules by metalloproteins is only beginning to be understood. How are gases of physiological relevance (eg,  $\text{O}_2$ , CO, NO,  $\text{CO}_2$ ) recognized by proteins other than the heme-containing [oxygen-binding proteins](#) and cytochrome oxidases, and what are the regulatory consequences at the genetic level? For example, the requirement for  $\text{CO}_2$  as an environmental determinant of gene expression involved in many bacterial capsule and [toxin](#) biosyntheses has long been recognized by microbiologists but is not understood. Curiously, heme proteins (4), rather than the nonheme iron-responsive proteins discussed in [Iron-response elements](#), have been implicated in most of what is presently known about the molecular mechanisms of gas sensing in biology.

The initial reports of mutant metal-binding proteins have naturally led to their consideration for use in the biotechnology industry. The development of artificial regulators of eukaryotic gene expression is a much sought after goal. Can metalloprotein-based drugs be produced that block the expression of genes whose products are important in the developing animal and plant diseases? Alternatively, could such types of regulatory proteins be used to produce large **transgenic** animals and plants? Novel zinc-finger proteins designed to function as specific transcriptional switches have been reported (5, 6). Other potential applications for engineered metal-binding proteins include highly sensitive metal-ion biosensors.

### Bibliography

1. B. A. Krizek, D. L. Merkle, and J. M. Berg (1993) *Inorg. Chem.* **32**, 937–940.
2. L. S. Powers (1982) *Bioch. Biophys. Acta* **683**, 1–38.
3. S. Silver and L. T. Phung (1996) *Annu. Rev. Microbiol.* **50**, 753–789.
4. D. Shelver, R. L. Kerby, Y. He, and G. P. Roberts (1997) *Proc. Natl. Acad. Sci. USA* **94**, 11216–11220.
5. Q. Liu, D. J. Segal, J. B. Ghiara, and C. F. Barbas III (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5525–5530.
6. J.-S. Kim and C. O. Pabo (1998) *Proc. Natl. Acad. Sci. USA* **95**, 2812–2817.

### Suggestions for Further Reading

7. I. Bertini, H. B. Gray, S. J. Lippard, and J. S. Valentine (eds.) (1994) *Bioinorganic Chemistry*, University Science Books, Mill Valley, CA.
8. L. Que Jr. (ed.) (1999) *Bioinorganic Spectroscopy and Magnetism*, University Science Books, Sausalito, CA.
9. G. L. Eichhorn and L. G. Marzilli (eds.) (1989) *Advances in Inorganic Biochemistry*, Elsevier, New York, Vol. **8**.
10. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New York, Vols. 1–35.
11. R. J. Cousins (1994) Metal elements and gene expression, *Annu. Rev. Nutr.* **14**, 449–469.

## Metallothionein

In 1957, Vallee and co-workers (1) discovered a small, **cysteine-rich protein**, which they termed metallothionein (MT). Metallothioneins have been found in most higher eukaryotes, fungi, and cyanobacteria. These proteins are produced in response to heavy-metal uptake, and they sequester these ions in protein-bound metal clusters. **Cysteine** is the sole amino acid residue that binds to metals in metallothioneins.

Regulation of metallothioneins is necessary because high metallothionein levels indiscriminately attenuate the levels of trace metals inside the cell. Metallothionein regulation is currently being studied in fungi and cyanobacteria as model metalloregulatory systems (2), where  $\text{Cu}^{1+}$  or  $\text{Zn}^{2+}$  function as the effector metal ions. Table 1 lists the key **transcription factors**, all of which dissociate from their DNA **metal-response elements** upon metal-ion binding. The **X-ray crystallographic structure** (3) of the cyanobacterial metallothionein **repressor**, SmtB, shows a classic **helix-turn-helix** motif, in common with many other **DNA-binding proteins**. The binding of DNA by SmtB, unlike other zinc metalloregulators (see **Zinc-Binding Proteins**), is inhibited by the metal. Two other transcription factors, ACE1 and AMT1 (see Table 1), are both activated by  $\text{Cu}^{1+}$ . Four copper ions congregate in an all-or-nothing manner to form a tetracopper-thiolate cluster in the N-terminus of each protein (4). The C-termini contain the **transactivation** domains. Cu- and S-EXAFS show that Cu-Cu distances of 2.7 Å and Cu-S distances of 2.26 Å characterize the clusters of both proteins and the copper-thiolate cluster in copper-loaded metallothionein. These observations point to an interesting similarity between  $\text{Cu}^{1+}$ -bound ACE1 (AMT1) and the gene product they regulate.

**Table 1. Model Transcription Factors for Metallothioneins**

| Organism                        | Protein | Metal Content                                      | Ref. |
|---------------------------------|---------|--|------|
| <i>Synechococcus</i> PCC 7942   | SmtB    | 2 $\text{Zn}^{2+}$ /monomer                        | 3    |
| <i>Saccharomyces cerevisiae</i> | ACE1    | 1 $\text{Zn}^{2+}$ ;<br>4 $\text{Cu}^{1+}$ cluster | 4    |
| <i>Candida glabrata</i>         | AMT1    | 1 $\text{Zn}^{2+}$ ;<br>4 $\text{Cu}^{1+}$ cluster | 4    |

### Bibliography

1. M. Margoshes and B. L. Vallee (1957) J. Am. Chem. Soc. **79**, 4813–4814.
2. D. J. Thiele (1992) Nucleic Acids Res. **20**, 1183–1191.
3. W. J. Cook, S. R. Kar, K. B. Taylor, and L. M. Hall (1998) J. Mol. Biol. **275**, 337–346.
4. J. A. Graden, M. C. Posewitz, J. R. Simon, G. N. George, I. J. Pickering, and D. R. Winge (1996)

Biochemistry **35**, 14583–14589.

### Suggestions for Further Reading

5. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New York, Vols. **1–35**.
6. R. J. Cousins (1994) Metal elements and gene expression, *Annu. Rev. Nutr.* **14**, 449–469.

## Metal-Requiring Enzymes

[Enzymes](#) that require metal ions for their catalytic activity fall into two classes. They are the metal-activated enzymes and the metalloenzymes. The latter contain tightly bound metals that do not dissociate during isolation or [dialysis](#) of the enzyme under conditions where activity is retained. However, such metal ions can be removed under more drastic conditions, such as low pH. Bound metal ions can be involved with the maintenance of the structural integrity of enzymes, and they can participate in electrophilic catalysis.

Metal ions that are found in metalloenzymes include those of the first transition series of elements in the periodic table:  $\text{Mn}^{2+}$ ,  $\text{Fe}^{2+}$ ,  $\text{Co}^{2+}$ ,  $\text{Ni}^{2+}$  and  $\text{Cu}^{2+}$ , as well as  $\text{Zn}^{2+}$ ,  $\text{Mo}^{2+}$  and  $\text{Ca}^{2+}$ . Examples of enzymes that contain metal ions are listed in [Table 1](#). Metal ions involved with enzymes that participate in electron transport undergo **redox reactions**. Thus, the ionic forms of iron, copper, cobalt, and molybdenum can be  $\text{Fe}^{2+}/\text{Fe}^{3+}$ ,  $\text{Cu}^+/\text{Cu}^{2+}$ ,  $\text{Co}^{2+}/\text{Co}^{3+}$ , and  $\text{Mo}^{2+}$  to  $\text{Mo}^{6+}$ , respectively. For convenience, these metal ions will be listed simply as bivalent ions.  $\text{Fe}^{2+}$  is most commonly found as a heme complex in redox enzymes such as [catalase](#) and **peroxidases** (see [Iron-Binding Proteins](#)). It also occurs as a component of **iron–sulfur** clusters in enzymes that are involved in one-electron transfer processes; NADH dehydrogenase and succinate dehydrogenase belong to this group and are flavoprotein enzymes. Like  $\text{Fe}^{2+}$ ,  $\text{Cu}^{2+}$  has multivalent oxidation states, and many  $\text{Cu}^{2+}$  enzymes are either oxidases or **hydrolases** that utilize molecular oxygen.  $\text{Co}^{2+}$  enzymes, such as methylmalonyl-CoA mutase and [ribonucleotide reductase](#), have the cobalt atom bound within a corrin ring.  $\text{Ni}^{2+}$  is rarely found as a component of metalloenzymes, but urease from jack bean is an exception. The occurrence of  $\text{Mn}^{2+}$  and  $\text{Ca}^{2+}$  in metalloenzymes is also somewhat rare (see [Calcium-Binding Proteins](#)).  $\text{Mg}^{2+}$ , the alkaline earth metal that is found so commonly in biological systems, does not play a role in the functioning of metalloenzymes, but is important in metal-activated enzymes (see below). By contrast,  $\text{Zn}^{2+}$  is an important and widely utilized metal for electrophilic catalysis (see [Zinc-Binding Proteins](#)). Not all enzymes that catalyze a particular reaction have the same requirement for a metal. Thus, fructose biphosphate aldolase from yeast and bacteria utilize Zn ions, whereas the same enzyme from muscle uses a [Schiff Base](#) intermediate to activate the substrate ([1](#)).

**Table 1. Selected Examples of Metalloenzymes**

---

| Metal Ion | Enzyme |
|-----------|--------|
|-----------|--------|

|                  |   |
|------------------|---|
| Ca <sup>2+</sup> | amylase, galactosyltransferase, <a href="#">thermolysin</a>   |
| Co <sup>2+</sup> | dioldehydrase, glycerol dehydratase, methylmalonyl-CoA mutase, <a href="#">ribonucleotide reductase</a>   |
| Cu <sup>2+</sup> | cytochrome c oxidase, dopamine-β-hydroxylase, superoxide dismutase  |
| Fe <sup>2+</sup> | <a href="#">catalase</a> , NADH dehydrogenase, nitrogenase, <a href="#">peroxidase</a> , succinate dehydrogenase, xanthine oxidase                                |
| Mn <sup>2+</sup> | arginase, histidine-ammonia lyase, pyruvate carboxylase   |
| Mo <sup>2+</sup> | nitrogenase, xanthine oxidase   |
| Ni <sup>2+</sup> | urease, Ni-Fe hydrogenase   |
| Zn <sup>2+</sup> | <a href="#">alcohol dehydrogenase</a> , <a href="#">carbonic anhydrase</a> , <a href="#">carboxypeptidase</a> , superoxide dismutase, <a href="#">thermolysin</a> |

---

The largest group of metal-activated enzymes contains the **phosphotransferases** that catalyze the transfer of the terminal phosphoryl group of ATP to an acceptor molecule that can be an alcohol, carboxylic acid, nitrogenous compound, or a phosphorylated compound (see [Kinase](#)). Their essential requirement for a bivalent metal ion is always satisfied by Mg<sup>2+</sup> or Mn<sup>2+</sup>. However, other bivalent metal ions have been shown to activate some phosphotransferases ([2](#)). The role of bivalent metal ions in the activation of phosphotransferases is to form a MgATP<sup>2-</sup> complex that then acts as the true substrate for the reaction. Thus, the binary complex formed by the interaction of the enzyme and its nucleotide substrate is an enzyme–nucleotide–metal complex. Some phosphotransferases involve a second metal ion that is liganded by the enzyme as well as the substrate. Examples are pyruvate kinase and the [biotin](#)-containing enzymes that form carboxybiotin by the initial phosphorylation of bicarbonate to carboxy-phosphate ([1](#)). Pyruvate kinase also differs from most other phosphotransferases in its requirement for K<sup>+</sup> and its inhibition, rather than activation, by Ca<sup>2+</sup>.

#### Bibliography

1. J. J. Villafranca and T. Nowak (1992) *The Enzymes* **20**, 63–94.
2. J. F. Morrison (1979) *Meth. Enzymol.* **63**, 257–294.

#### Metal Response Element

The metal response element is present in multiple copies in **genes**, such as the [metallothionein](#) genes, whose transcription is increased in response to treatment of cells with metals like zinc, copper, cadmium, and lead ([1](#), [2](#)). The metallothionein proteins are **cysteine**-rich metal-binding proteins that play a critical role in protecting cells from metal toxicity. Hence, increased levels of potentially toxic metals lead to enhanced [transcription](#) of the metallothionein genes. Inspection of the regulatory regions of the metallothionein genes reveals that they contain multiple copies of a 13- to 15-bp element with the **consensus sequence**

*CTCTGCRNCGGCCC*

where R = purine and N = any base. The core sequence of this motif (underlined above) is essential for the ability of the metallothionein gene to respond to metals. Thus mutagenesis of this sequence destroys the metal responsiveness of the metallothionein gene (3), while individual elements of this sequence can confer metal responsiveness upon a gene that is not normally induced in response to this treatment (4).

This element is thus a metal response element (MRE), which allows the transcription of specific genes to be activated in response to metals (see [Response Element](#)).

#### Bibliography

1. G. W. Stuart, P. F. Searle, H. Y. Chen, R. L. Brinster, and R. D. Palmiter (1984) Proc. Natl. Acad. Sci. USA **81**, 7318–7322.
2. H. M. Zafarulla, K. Bonham, and L. Gedamu (1988) Mol. Cell. Biol. **8**, 4469–4476.
3. A. V. Cizewski Culott and D. H. Hamer (1989) Mol. Cell. Biol. **9**, 1376–1380.
4. G. W. Stuart, P. F. Searle, and R. D. Palmiter (1985) Nature **317**, 828–831.

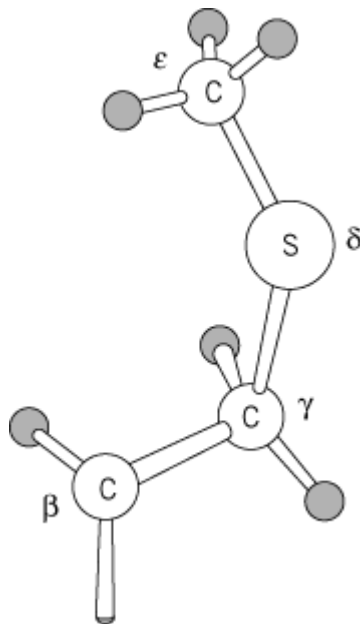
#### Suggestion for Further Reading

5. D. J. Thiele (1992) Metal-regulated transcription in eukaryotes. Nucleic Acids Res. **20**, 1183–1191.

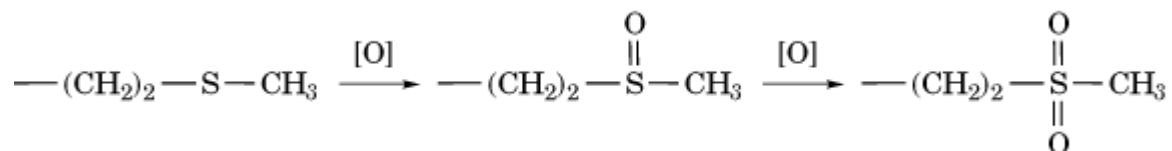
### Methionine (Met, M)

The [amino acid](#) methionine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to only one **codon**—AUG—and represents approximately 2.4% of the residues of the proteins that have been characterized. It is usually the initiating residue in [translation](#) and occurs at the *N*-terminus of the initial gene product. In **prokaryotes**, the amino group of this residue is formylated, but this group is normally removed post-translationally by a formylase [enzyme](#). Frequently, the initiating methionyl residue is removed by a **ribosome**-associated Met-aminopeptidase enzyme. An internal methionyl residue has a mass of 131.19 Da, a **van der Waals volume** of 124 Å<sup>3</sup>, and an [accessible surface](#) area of 204 Å<sup>2</sup>. Met residues are changed only moderately frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [leucine](#) residues.

The long side chain of Met residues is [nonpolar](#) and relatively unreactive chemically:



It is the only unbranched nonpolar side chain of all the natural amino acids. The sulfur atom is somewhat nucleophilic, but unlike other nucleophiles in proteins, it cannot be protonated. Consequently, it is the most potent nucleophile in proteins at acidic pH. Its reaction with **cyanogen bromide** under such conditions causes the [polypeptide chain](#) to be cleaved after Met residues, which are converted to homoserine residues; this is the most useful chemical method for cleaving polypeptide chains. The Met sulfur atom can be selectively modified under acidic conditions, readily forming sulfonium salts with alkylating agents, such as methyl iodide. This reaction can be reversed by [thiol groups](#), and the methyl group removed is equally likely to be the original one or that introduced by the methyl iodide, so this reaction offers the possibility of introducing a radioisotope label in 50% of the residues by using labeled methyl iodide. The sulfur atom of Met residues is also susceptible to oxidation by air or by more potent oxidants, such as peroxides. The sulfoxide is formed first, followed by the sulfone:



The first step, but not the second, can be reversed by sulfides and by thiol groups. Either oxidation state makes the Met residue unreactive with alkylating reagents.

In folded proteins, about 40% of the Met residues are totally buried. The Met residue favors the  $\alpha$ -helix conformation in model peptides and occurs most frequently in that type of **secondary structure** in folded [protein structures](#), although it also occurs frequently in *beta*-sheets. In some [cytochromes](#), the Met sulfur atom is one ligand to the iron atom of the heme group.

#### Suggestions for Further Reading

B. Witkop (1961) Nonenzymatic methods for the preferential and selective cleavage and modification of proteins, *Adv. Protein Chem.* **16**, 221–321.

E. Gross (1967) The cyanogen bromide reaction, *Meth. Enzymol.* **11**, 238–255.

N. Brot and H. Weissbach (1983) Biochemistry and physiological role of methionine sulfoxide residues in proteins, *Arch. Biochem. Biophys.* **223**, 271–281.

## Methionine Regulon

In addition to its role as an [amino acid](#) constituent of [proteins](#), *methionine* participates in initiating **protein biosynthesis** (as *N*-formylmethionine), as a universal methylating agent (as S-adenosylmethionine, AdoMet;), and as a precursor of **spermidine** through AdoMet. Consequently, controlling the growth rate by regulating the availability of methionine is a very effective strategy in microorganisms.

Methionine is synthesized only in **microorganisms** and **plants**, and its biosynthetic pathway is complex. The carbon skeleton derives from aspartic acid, the sulfur atom from cysteine, and the methyl group from the  $\beta$ -carbon of serine (1). In *Escherichia coli*, the genes involved in methionine biosynthesis are scattered throughout the [chromosome](#). Table 1 summarizes the names of the genes, of the protein they encode, their position on the circular chromosome, and the product of the catalyzed reaction. The conversion of aspartate to homoserine has been treated elsewhere (see [Multifunctional Proteins](#); [Enzymes](#) and [Isofunctional Proteins](#)).

**Table 1. Genes and Enzymes in *Escherichia coli* for Biosynthesizing Methionine from Homoserine**

| Genes       | Gene Product                                | Position (min) | Reaction Product             |
|-------------|---|----------------|------------------------------|
| <i>metA</i> | Homoserine transsuccinylase                 | 90.7           | <i>O</i> -Succinylhomoserine |
| <i>metB</i> | Cystathionine- $\gamma$ -synthase           | 88.9           | Cystathionine                |
| <i>metC</i> | $\beta$ -Cystathionase                      | 67.8           | Homocysteine                 |
| <i>metE</i> | Methionine synthase (cobalamin-independent) | 86.4           | Methionine                   |
| <i>metH</i> | Methionine synthase (cobalamin-dependent)   | 90.9           | Methionine                   |

The methyl group of methionine originates from the  $\beta$  carbon of serine, via methylenetetrahydrolate, in a series of enzymatic reactions listed in Table 2.

**Table 2. Biosynthesis of the Methyl Group of Methionine Used by the Cobalamin-Independent or Cobalamin-Dependent Methionine Synthases**

| Genes | Gene Product | Position (min) | Reaction Product(s) |
|-------|--------------|----------------|---------------------|
|-------|--------------|----------------|---------------------|

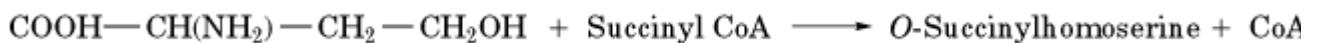


|             |  |      |                                     |
|-------------|--|------|-------------------------------------|
| <i>glyA</i> | Serine hydroxymethyl transferase                 | 57.6 | Glycine + methylenetetrahydrofolate |
| <i>metF</i> | Methylenetetrahydrofolate <sup>a</sup> reductase | 89   | Methylenetetrahydrofolate           |

---

<sup>a</sup> Methylenetetrahydrofolate is the direct methyl donor to homocysteine.

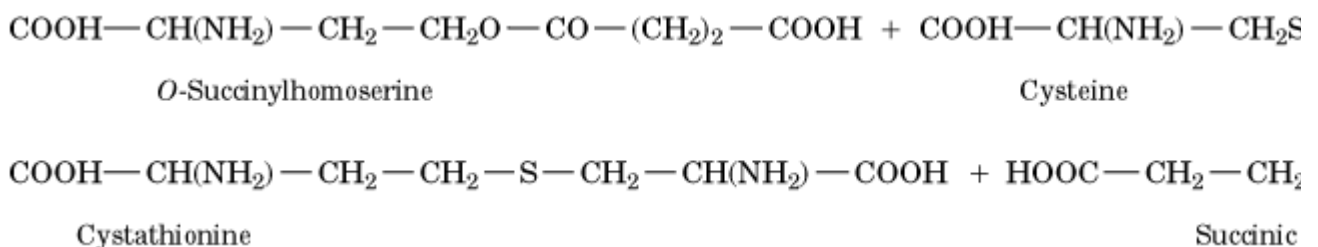
The first step of methionine biosynthesis is the succinylation of homoserine, catalyzed by *homoserine succinyltransferase*, the product of the *metA* gene (2):



In yeast, acetylCoA and acetyl- *O*-homoserine replace the succinylated products.

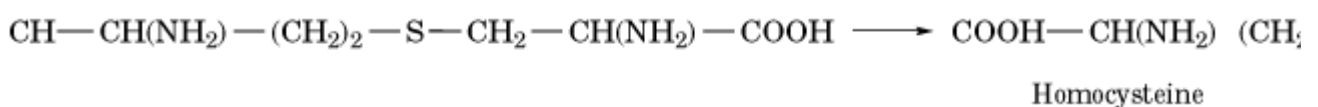
The *metA* gene codes for a polypeptide of molecular weight 35,673 (3), and the enzyme is a homodimer. It has unusual temperature sensitivity in many *Enterobacteriaceae*, which limits growth at elevated temperature (4). The enzyme appears to be a **heat shock** protein and is involved in other **stress responses**. In this respect, note that one of the genes (of unknown biochemical function) involved in methionine biosynthesis in *Saccharomyces cerevisiae* is identical to a gene responsible for high salt tolerance. Furthermore, supplementing the growth medium with methionine increases the tolerance of **yeast** to NaCl.

*O*-succinylhomoserine is transformed to cystathionine, a seven-carbon diaminoacid, in the presence of cysteine (5):



This reaction is catalyzed by cystathionine-g-synthase, encoded by the *metB* gene, and uses [pyridoxal phosphate](#) as a **cofactor**. The reaction mechanism of the enzyme is thought to proceed via a pyridoxamine derivative of vinylglyoxylate. The enzymes from *Salmonella typhimurium* and *E. coli* are homotetramers of a 40-kDa subunit, each containing one pyridoxal phosphate molecule. The *metB* gene of *E. coli* codes for 386 amino acid residues (6).

Cystathionine is cleaved to give homocysteine, ammonia, and pyruvate in a reaction catalyzed by cystathionase (cystathionine-b-lyase), the product of the *metC* gene (7).



This is also a pyridoxal phosphate-containing enzyme that has been only partially characterized. 3,3,3-Trifluoroalanine binds covalently to the enzyme and inhibits it irreversibly. The protein was first reported to be a hexamer, each subunit binding a molecule of pyridoxal phosphate, but further studies suggest that it is a homotetramer (8). Its sequence of 395 residues is strongly homologous with that of cystathionine-g-synthase, indicating a common ancestor of the two proteins (9).

A mutation in the gene *metQ* (39.5 min) allows *metC* mutants to catalyze the formation of homocysteine directly from homoserine, bypassing the normal cystathionine intermediate (10).

The last step in methionine biosynthesis involves methylation of homocysteine:



The methyl donor **R** is <sup>5</sup>N-methyl tetrahydropteroylglutamate, a member of the folate family (see [Aminopterin, Methotrexate, Trimethoprim, and Folic Acid](#)), which is derived from <sup>5</sup>N, <sup>10</sup>N-methylene tetrahydrofolate (generated by the conversion of serine to glycine catalyzed by serine hydroxymethyltransferase) by a reductase specified by the *met F* gene. This reaction produces the methyl group specific for methionine biosynthesis. Two enzymes catalyze this reaction (11):

1. a transmethylase (product of the *metH* gene) dependent on the presence of vitamin **B12** that can use either the mono- or polyglutamate forms of <sup>5</sup>N-methyl tetrahydropteroylglutamate as a methyl donor;
2. a vitamin B<sub>12</sub>-independent transmethylase (product of *metE*) that uses only the polyglutamate forms of <sup>5</sup>N-methyl tetrahydropteroylglutamate (N > 3).

The gene *metE* encodes a polypeptide chain of 84654 Da (12), and *metH* codes for an exceptionally long polypeptide of 132,628 Da showing no similarity to the *metE* product (13). *E. coli* possesses both enzymes, whereas most microorganisms possess only one. It does not synthesize vitamin B<sub>12</sub> and uses one or the other transmethylase according to the availability of vitamin B<sub>12</sub> in the medium. Strains with a mutation in *metE* require either methionine or vitamin B<sub>12</sub> for their growth. However, vitamin B<sub>12</sub> can be synthesized anaerobically. The vitamin B<sub>12</sub>-independent enzyme, coded by *metE*, is a rather inefficient enzyme.

The mechanism of the B<sub>12</sub>-dependent enzyme involves methylation of the vitamin's cobalt atom. Catalytic amounts of AdoMet are required for this priming step. After the methyl group is transferred to homocysteine, subsequent methyl groups come from <sup>5</sup>N-methylenetetrahydrofolate. In addition to AdoMet, methyl B<sub>12</sub>, and the folate derivative, the enzyme requires reduced **FAD** for activity.

The B<sub>12</sub>-independent transmethylase involves neither a methylated coenzyme nor the other requirements of the vitamin-dependent enzyme, such as AdoMet or a reducing agent. Alkylation of Cys726 of the enzyme results in complete loss of activity. This thiol might function as an intermediate methyl acceptor in catalysis, analogous to the role of cobalamin in the reaction catalyzed by the B<sub>12</sub>-dependent enzyme.

Table 3 refers to the genes responsible for synthesizing proteins involved in regulating methionine biosynthesis (see [Methionine Repressor](#)) and in synthesizing AdoMet, the corepressor of many of these proteins.

**Table 3. Genes Involved in Regulating Methionine Biosynthesis**

| Genes       | Gene product        | Position (min) | Reaction product              |
|-------------|---------------------|----------------|-------------------------------|
| <i>metK</i> | Methionine adenosyl | 66.4           | S-Adenosylmethionine (AdoMet) |

|             |                          |      |   |
|-------------|--------------------------|------|---|
|             | transferase              |      |   |
| <i>metJ</i> | Methionine repressor     | 88.9 | Methionine repressor                                      |
| <i>metR</i> | Positive regulatory gene | 86.4 | Protein regulating <i>metE</i> and <i>metH</i> expression |

---

## Bibliography

1. R. C. Greene (1996) In *Escherichia coli and Salmonella* (F. C. Neidhardt, editor-in-chief) American Society for Microbiology Press, Washington DC, vol. **1**, pp. 542–560.
2. R. J. Rowbury (1962) *J. Gen. Microbiol.* **28**, V–VI.
3. B. Duclos et al. (1989) *Nucleic Acids Res.* **17**, 2856.
4. E. Z. Ron and M. Shani (1971) *J. Bacteriol.* **107**, 397–400.
5. M. M. Kaplan and M. Flavin (1966) *J. Biol. Chem.* **241**, 4463–4471 and 5781–5789.
6. N. Duchange et al. (1983) *J. Biol. Chem.* **258**, 14868–14871.
7. M. Flavin (1975) In *Metabolic Pathways* (D. M. Greenberg, ed.) Academic Press, New York, pp. 457–503.
8. I. G. Old et al. (1991) *Prog. Biophys. Mol. Biol.* **56**, 145–185.
9. J. Belfaiza et al. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 867–871.
10. M. Simon and J. S. Hong (1983) *J. Bacteriol.* **153**, 558–561.
11. R. T. Taylor and H. Weissbach (1973) In *The Enzymes*, 3rd ed. Group Transfer (P. D. Boyer, ed.) Academic Press, New York, Vol. **IX**, part B, pp. 121–165.
12. M. E. Maxon et al. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 85–89.
13. I. G. Old et al. (1990) *Gene* **87**, 15–21.

## Methionine Repressor

### 1. Regulation of Methionine Biosynthesis by Gene Expression

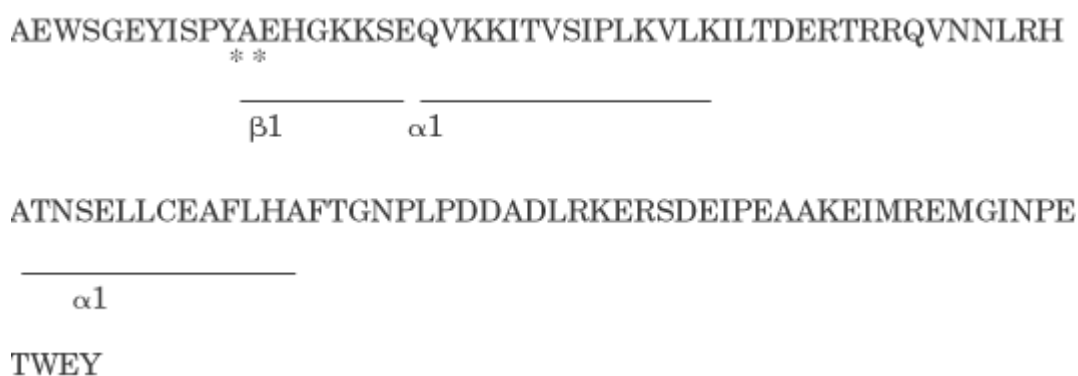
The *metA* gene codes for the first step in the methionine biosynthetic pathway, starting from homoserine (see [Methionine Regulon](#)). It has two starting sites for [transcription](#) located 74 nucleotides apart. Both **promoters** are used *in vivo*, but only one is regulated by intracellular methionine levels. The *metB*, *metC*, and *metF* genes, coding respectively for cystathionine synthase, cystathionase, and 5,10-methylene tetrahydrofolate reductase, each have a single promoter. The *metA*, *metB*, *metC*, *metE*, and *metF* genes are all **repressed** when excess methionine is added to the minimal growth medium, albeit to different degrees. The level of the enzyme varies 300-fold for homoserine transsuccinylase, 40-fold for cystathionine synthetase, 6- to 12-fold for cystathionase, 20-fold for the methylene tetrahydrofolate reductase, and 60-fold for the vitamin **B12**-independent transmethylase. Thus there is no coordinate repression, reflecting the fact that the *met* genes do not form a single [operon](#), but are scattered through the bacterial chromosome.

There is a notable exception: *metB* and *metL*, coding for cystathionine synthetase and aspartate kinase II-homoserine dehydrogenase II, constitute an operon and belong to a [gene cluster](#), **metJBLF**, located around 89 min. on the chromosome (see [Methionine Regulon](#)). Methionine also regulates the

synthesis of aspartokinase II-homoserine dehydrogenase II (*metL*) and **S-adenosylmethionine** (AdoMet) synthetase (*metK*). Although methionine affects the level of the B<sub>12</sub>-dependent transmethylase, coded by *metH*, this effect appears to be the indirect result of repression of the synthesis of an activator protein (MetR, see below). Expression of the *metE* gene is repressed by vitamin B<sub>12</sub>.

The *E. coli metJ* gene is 312 bp long and codes for the Met repressor, a molecule of 104 amino acid residues that exists as a homodimer (Fig. 1). *metJ* is transcribed divergently from *metB*. Consequently a complex 276-bp regulatory region is found between them. There is a single promoter for *metB*, whereas *metJ* is transcribed from three separate points. Only two of the three *metJ* promoters are regulated by methionine. The repressor regulates its own synthesis.

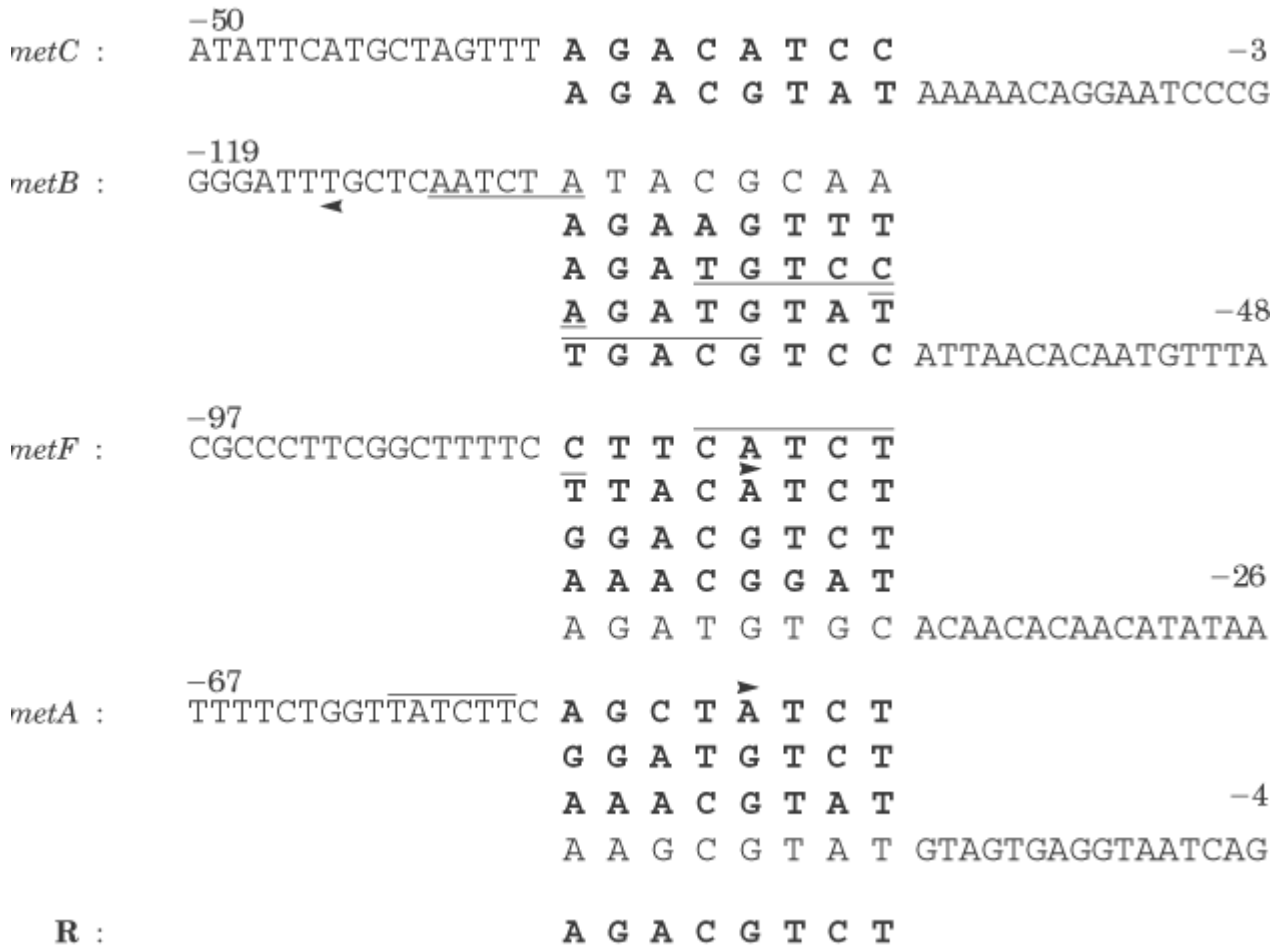
**Figure 1.** Sequence of the *E. coli* methionine repressor. The \* sign indicates the amino acid residues making contacts with DNA, inferred from the structure of the repressor/DNA complex. The underlined sections indicate the b-strand and two of the a-helices of the protein (1).



The effect of the repressor on the synthesis of the methionine biosynthetic enzymes has been assessed by measuring **Beta-galactosidase** synthesis in a cell-free system where the synthesis of this enzyme is under the control of *metF* (*metF-lacZ* fusion). Methionine had to be present in these experiments as a necessary building block for protein synthesis. Expression of b-galactosidase was progressively repressed by increasing concentrations of the aporepressor. Complete repression was attained at 600 nM. At a constant concentration of aporepressor, AdoMet enhances the repression in a concentration-dependent manner. Half-maximum inhibition is reached at 10 mM AdoMet and complete repression at 100 mM. Similar *in vitro* assay systems demonstrated that the aporepressor and AdoMet have the same effect on the expression of the *metB* and *metL* genes and on that of the *metJ* gene itself.

[Equilibrium Dialysis](#) demonstrates that the methionine repressor binds two molecules of S-AdoMet per dimer. The corresponding [Scatchard Plot](#) is linear, showing that there is no **cooperativity**. Methionine does not bind to the aporepressor. Because the repressor binding sites should be similar for all of the genes regulated by methionine, comparison of the four 5' regions of the *metC*, *metB*, *metA*, and *metF* genes reveals a repetitive unit (R) eight nucleotides long, named the “Met box.” In the alignment presented in Figure 2 of 128 positions, 89 matches, and 21 transitions are found when the repetitive units are compared to the **consensus sequence** R. This consensus sequence is a perfect [palindrome](#), AGACGTCT, which is present in an altered form two to five times in the Met boxes. The *metB* and *metJ* genes are transcribed divergently and share the same Met box.

**Figure 2.** Comparison of the upstream regions of the *metC*, *metB*, *metF*, and *metA* genes. The sequences 5' to the structural *metC*, *metB*, *metF*, and *metA* genes are presented discontinuously and have been aligned in order to focus on the presence of the underlying repetitive palindromic repetitive unit. Nucleotides matching the consensus sequence presented in line **R** are in boldface. Numbers indicate positions relative to the adenine of the respective start codon taken as +1. The -10 promoter sequences are overlined and arrowheads indicate the transcription start signals. In the case of *metB*, the overlined hexamer is the -35 box: the two underlined promoter sequences represent the -35 and -10 boxes of the first promoter of *metJ* (2).



The *metE* and *metR* genes also possess these repetitive units (3 for *metE* and 4 for *metR*). The two exceptions are *metK* and *metH* which, although regulated by methionine, do not possess the consensus sequence. The differences in the number and/or sequences of the Met boxes may be related to the different extents of repression elicited by the Met repressor. A **filter-binding assay** using [nitrocellulose](#) and **radioactive** oligodeoxynucleotides has shown that DNA fragments containing two consecutive consensus Met boxes are tightly bound by repressor in the presence of saturating AdoMet, whereas nucleotide sequences containing only one Met box are not bound. Binding is **cooperative** with respect to repressor concentration. The repressor protects the operator against nuclease digestion (see [Footprinting](#)), and the results strongly indicate the binding of an array of repressor dimers centered on the 16-bp operator site but extending into the neighboring DNA. The same was found for a fragment containing the *metF* regulatory region, where five boxes are protected. These observations are consistent with the following *in vivo* data: the level of repression of defined mutants of the *metC* operator, the smallest one known (having two Met boxes), is increased when the sequence becomes closer to the consensus sequence. Similar results were obtained with the five Met boxes of *metF*. Furthermore, the arrangement of all five boxes in tandem repeats is important for effective repression, since some point insertions within the operators lower repression 100-fold.

## 2. Three-Dimensional Structure of the Methionine Repressor of its Binary Complex with AdoMet and of its Ternary Complex with AdoMet and DNA

As in the case of the Trp repressor, (see [TRP Operon](#)) the two subunits of the Met repressor are strongly intertwined to form the dimer. AdoMet binds to the dimer at two independent but symmetrical sites, one on each monomer. The purine nucleus of AdoMet inserts into a **hydrophobic** pocket, normally occupied in the aporepressor by the side-chain of Phe65, whereas the methionine moiety lies at the protein surface. This explains why methionine itself does not bind to the aporepressor, whereas S-adenosylhomocysteine binds with about half the affinity of AdoMet. The positively charged trivalent sulfur atom lies at the C-terminus of helix B. In contrast to the Trp repressor, there is no major conformational change associated with corepressor binding. The apo and holo structures are indistinguishable, except for the absence or presence of AdoMet and small conformational changes at the N-terminus and Phe65. AdoMet does not contact DNA in the repressor-operator complex (see below), and how it regulates binding remains unclear. One possibility is a long-range electrostatic effect based on the positive charge on the sulfur atom.

Each monomer of the Met repressor is composed of three  $\alpha$ -helices and one  $\beta$ -strand, accounting for about 50% of the sequence (see Fig. 1). In the protein dimer, the  $\beta$ -strands pair to form an antiparallel [Beta-sheet](#), whereas the  $\alpha$ -helical regions pack against the sheet and against each other to stabilize the dimer. The molecule lacks a [helix-turn-helix motif](#) which makes it different from many other repressors and from the **cyclic AMP receptor protein**.

By [X-ray crystallography](#) it was found that a 19-bp oligonucleotide containing two adjacent 8-bp Met boxes binds two dimeric Met repressor molecules. One dimer binds to each half-site, each of which contains a twofold axis of symmetry that coincides with the twofold axis of the protein  $\beta$ -sheet. The  $\alpha$ -helix 1 on one dimer interacts with the same  $\alpha$ -helix 1 on the adjacent dimer to form a tetrameric protein structure. Sequence specificity is achieved by inserting the double-stranded, antiparallel, protein  $\beta$ -sheet into the major groove of B-form DNA. Direct hydrogen-bonding occurs between amino acid side-chains on the exposed face of the sheet and the base pairs. Lys23 from each  $\beta$ -strand contacts a guanine base, and the neighboring Thr-25 on each strand contacts an adenine. Residues from the N-terminus of  $\alpha$ -helix 2 make backbone contacts. The repressor also recognizes sequence-dependent distortion or flexibility of the operator phosphate backbone and confers specificity even for inaccessible base pairs.

## 3. The *metR* Gene and its Product

Surprisingly it was found that some mutants with normal *metF*, *metE* and *metH* genes express the two transmethylases from the *metE* and *metH* genes at very low levels, resulting in a growth requirement for methionine. This **auxotrophy** was overcome in strains carrying **multicopy plasmids** with either the *metE* or the *metH* gene. This implied that the methionine auxotrophy was caused by independent mutations resulting in an inability to synthesize enough homocysteine transmethylase enzyme to permit growth. It was found that these mutations are linked to *metE* but lie outside the *metE* structural gene. Their locus was called *metR*. The *metR* gene from *E. coli* encodes a 35,628-Da polypeptide 317 amino acid residues long. The native protein is a homodimer and is believed to contain a [leucine zipper](#), a motif characteristic of many eukaryotic DNA binding proteins (and of the *E. coli* [Lac repressor](#)).

Expression of *metR* is repressed by the Met repressor, as the 5' flanking region of *metR* contains four Met boxes, and by the product of *metR* itself. Conversely, the MetR protein plays no role in regulating the *metJ* gene for the Met repressor. MetR activates the expression of both *metE* and *metH*. Homocysteine, the substrate of both transmethylases, is the coactivator for *metE* but not *metH*. This positive regulatory mechanism may apply when vitamin B<sub>12</sub> is absent, ie, when there is no B<sub>12</sub>-dependent transmethylase activity (catalyzed by the *metH* product). The resulting methionine

starvation would then cause derepression of the biosynthetic enzymes and accumulation of homocysteine. The latter could then act as a signal to activate the synthesis of the vitamin B<sub>12</sub>-independent transmethylase, coded by *metE*.

In *S. typhimurium*, the MetR protein activates *metA*, the gene for homoserine succinyltransferase, but homocysteine inhibits this activation. This result explains why expression of the *metE* gene is repressed by vitamin B<sub>12</sub>. It is primarily caused by a loss of MetR-mediated activation through depletion of the coactivator homocysteine, rather than a direct repression by the methH-B<sub>12</sub> holoenzyme.

### Bibliography

1. I. Saint Girons et al. (1984) *J. Biol. Chem.* **259**, 14282–14285.
2. J. Belfaiza et al. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 867–871.

### Suggestions for Further Reading

3. I. Saint-Girons et al. (1988) “Methionine synthesis in Enterobacteriaceae: Biochemical, regulatory and evolutionary aspects,” *CRC Crit. Rev. Biochem.* **23** (suppl. 1), 1–42.
4. S. E. V. Phillips et al. (1989) “Cooperative tandem binding of methionine repressor of *Escherichia coli*,” *Nature* **341**, 711–715.
5. B. Rafferty et al. (1989) “Three-dimensional crystal structures of *Escherichia coli met* repressor with and without corepressor,” *Nature* **341**, 705–710.
6. W. S. Somers and S. E. V. Phillips (1992) “Crystal structure of the *met* repressor-operator complex at 2.8 Å resolution reveals DNA recognition by -strands,” *Nature* **359**, 387–393.
7. R. C. Greene (1996) “Biosynthesis of methionine.” In *Escherichia coli and Salmonella* (F. C. Neidhardt, editor-in-chief), American Society for Microbiology Press, Washington, DC, Vol **1**, pp. 542–560.

## Methylation, DNA

Enzymatic DNA methylation consists of the covalent attachment of a methyl group to a cytosine or adenosine residue in a defined sequence of double-strand DNA by a DNA **methyltransferases** (also called a *methylase*). DNA methyltransferases use **S-adenosyl-L-methionine** (AdoMet or SAM) as the methyl-donating cofactor and are separated into two major classes (see [Methyltransferase, DNA](#)). The first class methylates the C5 carbon of cytosine to form **5-methylcytosine** (5mC). The second class methylates exocyclic amino nitrogen atoms to form either *N*-6-methyladenine (N6mA) or *N*-4-methylcytosine (N4mC). **Eukaryotic** DNA contains 5mC and **prokaryotic** DNA contains all three methylated bases: 5mC, N6mA, and N4mC. These methyl substituents are located in the major groove of DNA and disrupt interactions with some [DNA-binding proteins](#) and [enzymes](#) by perturbing binding or [catalysis](#) (1). These methyl-induced disturbances may be due to direct steric clashes combined with longer-range conformational changes produced in the DNA–protein interface (2).

### 1. Prokaryotic Methylation

DNA methylation plays diverse roles in biology. For instance, methylation by the *N*-6-methyladenine DNA methyltransferase (Dam), found in *Escherichia coli* “marks” a DNA strand for

repair by the methyl-directed [mismatch repair](#) system. In this repair system, a newly synthesized DNA strand, which lacks methyl groups, is scanned for mispairing of its bases with the original parental DNA strand that had been methylated by Dam MTase prior to replication. The parental strand serves as a [template](#) to direct the replacement of erroneously incorporated nucleotides. Shortly after proofreading and repair, Dam methylase modifies the new daughter strand to yield fully methylated DNA. In addition, Dam methylase plays an important role in [DNA replication](#). High-efficiency *E. coli* chromosomal replication is dependent on eleven fully methylated Dam (GmATC) sites located at the *oriC* (the **origin of replication**). Although not lethal, loss of Dam methylation leads to increased spontaneous [mutation](#) rates and asynchronous DNA replication (3). Another prokaryotic N6mA methylase (CcrM) from *Caulobacter crescentus* is [cell cycle](#)-regulated and essential for viability. CcrM methylation at GATNC sites of daughter strands is temporally regulated at late stages of chromosomal replication and functions before **cell division**. Cells constitutively expressing CcrM throughout the cell cycle have abnormalities in their DNA replication and cell division. These findings indicate that the state of chromosomal DNA methylation regulates the progression of the cell cycle in this organism (4).

Prokaryotic methyltransferases, as part of [restriction–modification systems](#), methylate host DNA sequences for protection against cleavage by the partner [restriction enzyme](#). Methylation by enzymes like Dam that is involved in the normal functioning of host genomes may have been the evolutionary progenitors of restriction endonucleases and restriction–modification systems (5).

## 2. Eukaryotic Methylation

In eukaryotes, methylation functions in **gene expression** and regulation and cell division and [development](#). Eukaryotic methylation is accomplished by DNA methyltransferases (DNA MTases) to form 5-methylcytosine residues at CG sites (6). Mutations that reduce cytosine methylation in **fungi** and plants result in abnormal chromosomal segregation and stability. Methylation may also serve to compartmentalize structurally large eukaryotic genomes (such as those of plants and vertebrates) into inactive, compact, methylated regions and unmethylated regions accessible to gene [transcription](#) and regulatory factors. In vertebrate genomes, 5-methylcytosine represents 1% of bases (7). Vertebrate methylation is developmentally and tissue-specific and leads to defined, heritable patterns because of the preference of DNA MTase for hemimethylated DNA, specifically, DNA methylated on only one strand at a given methylation site (8). Li et al. (8) found that transfected mice embryos that were homozygous for a MTase mutation had 30% the normal amount of DNA methylation, major developmental abnormalities, and a recessive lethal phenotype. Methylation patterns have a role in **genomic imprinting**, where paternal and maternal **alleles** of a gene are differentially expressed. For instance, human and mouse insulin-like [growth factor](#) genes (*Igf 2*) are selectively expressed from the paternal gene, which, unlike the maternal allele, contains a methylated region upstream of the *Igf 2* promoter. This finding suggests that the methylated region may interfere with the binding of a transcriptional [repressor](#) (9). Abnormal methylation of **promoters** of [tumor suppressor genes](#) promotes tumorigenesis (10). Interestingly, the spontaneous hydrolytic deamination of 5-methylcytosine residues in DNA produces thymidine; on replication, a transition mutation occurs that is responsible for one-third of human mutations (7).

## 3. Molecular Basis of Effects of Methylation

An understanding of the effects of methylation on local DNA structure, DNA-protein interactions, and catalysis provides insight into the multiple roles of methylation. For example, thermodynamic melting experiments using GATC-containing and methylated GmATC-containing decameric oligodeoxyribonucleotide duplexes show that the  $T_m$  value (the **melting temperature** at which half of the duplex population has dissociated into single strands) was 10° lower for the methylated oligomer, indicating that methylation had a destabilizing effect on the DNA (11). Studies were also conducted with *EcoRI* endonuclease using fully methylated oligomeric duplexes methylated at the second adenine containing the canonical sequence (GAmATTC) (the biologically relevant



methylation site for *EcoRI*). Methylation destabilized the conformation of the DNA–protein interface, interrupting the precise structural geometry needed at the [transition state](#) for optimal catalysis and thereby prevented cleavage. Also, the cleavage rate constants for each strand on doubly methylated duplex were reduced ~600,000-fold compared to the unmethylated sequence. This rate reduction, combined with the lowered binding association rate caused by the methyl groups, explains why methylation effectively prevents *EcoRI* endonuclease cleavage (2).

DNA methylation is a ubiquitous, essential biological phenomenon in prokaryotes and eukaryotes. It is relevant to many areas of biology, molecular biology, biochemistry, and medicine. Examining DNA methylation at the molecular level also helps provide an understanding of the dynamic structure of DNA, protein–DNA interactions, and catalysis. Several uses for methylation of DNA in molecular cloning are discussed under [Methyltransferase, DNA](#) and [Staggered Cut](#).

### Bibliography

1. M. Nelson, E. Raschke, and M. McClelland (1993) *Nucl. Acids Res.* **21**, 3139–3154.
2. L. Jen-Jacobsen, L. E. Engler, D. R. Lesser, M. R. Kurpiewski, C. Yee, and B. McVerry (1996) *EMBO J.* **15**, 2870–2882.
3. B. R. Palmer and M. G. Marinus (1994) *Gene* **143**, 1–12.
4. C. Stephens, A. Reisenauer, R. Wright, and L. Shapiro (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1210–1214.
5. T. A. Bickle and D. H. Kruger (1993) *Microbiol. Rev.* **57**, 434–450.
6. W. C. Yen et al. (1992) *Nucl. Acids Res.* **20**, 2287–2291.
7. T. H. Bestor and G. L. Verdine (1994) *Curr. Opin. Cell Biol.* **6**, 380–389.
8. E. Li, T. H. Bestor, and R. Jaenisch (1992) *Cell* **69**, 915–926.
9. A. Razin and H. Cedar (1994) *Cell* **77**, 473–476.
10. G. L. Verdine (1994) *Cell* **76**, 197–200
11. Q. Guo, M. Lu, and N. R. Kallenbach (1995) *Biochemistry* **34**, 16359–16364.

### Suggestion for Further Reading

12. A. Razin, H. Cedar, and A. D. Riggs, eds. (1984) *DNA Methylation: Biochemistry and Biological Significance*, Springer-Verlag, New York. (Provides basic background of several aspects of DNA methylation).

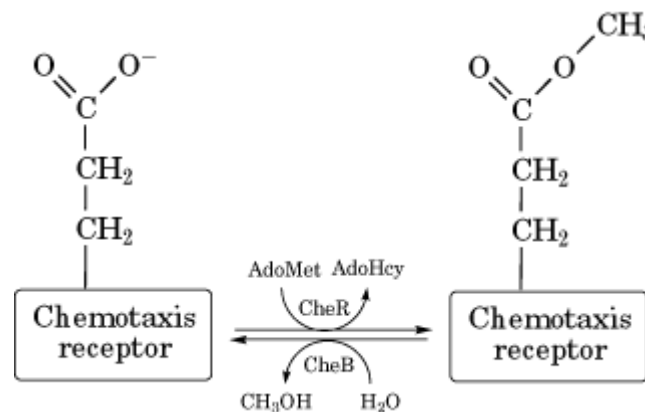
## Methylation, Protein

A number of amino acid residues are subject to methylation, catalyzed by **methyltransferases** using **S-adenosylmethionine** as the methyl donor (1). In carboxyl methylation of glutamate, aspartate, or carboxyl-terminal residues, the methylation is reversible. In methylation of sulfur or nitrogen atoms, it is probably not reversed during normal metabolism. In the former, the methylation may play a role in [signal transduction](#), or it may be connected with protein rejuvenation. In the latter, methylation is probably related to specific structures of proteins that enable them to function better (2, 3). Four distinct types of carboxyl methylation are known.

### 1. Glutamate Residues

The first is that of specific glutamate side-chains on chemoreceptors mediating [chemotaxis](#) in bacteria (Fig. 1). Increased methylation increases the activity of the CheA kinase that governs the behavior. The purpose of the receptors is to regulate the CheA kinase. For example, in *E. coli*, reduction of CheA activity caused by addition of attractant is compensated for by receptor methylation so that the stimulus is only short-lived (4) (see [Chemotaxis](#)).

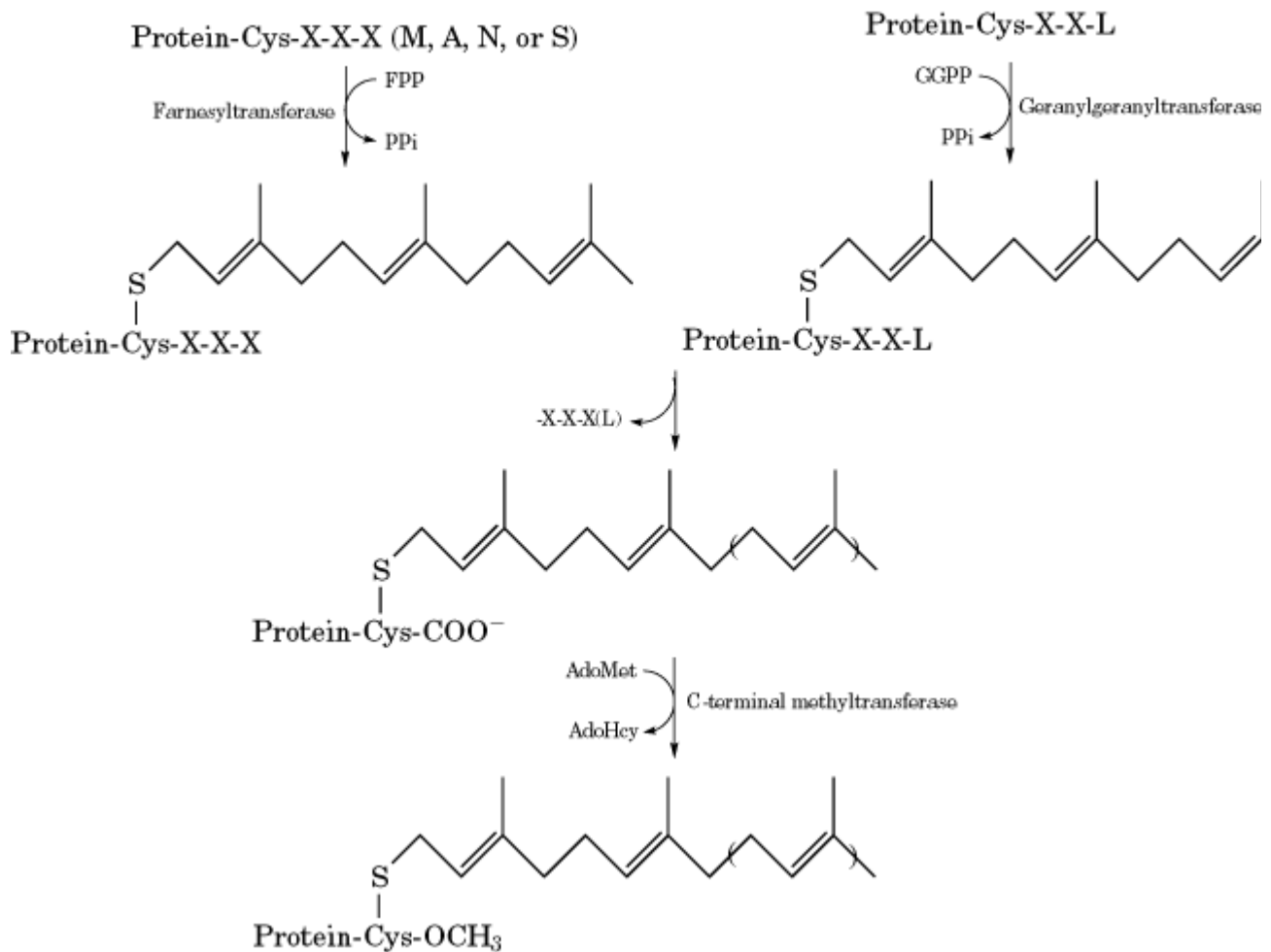
**Figure 1.** Methylation reactions on glutamate residues within the chemotaxis receptors in bacteria. The methylation reaction is catalyzed by CheR, and demethylation by CheB. Upon loss of its methyl group, the *S*-adenosyl methionine (AdoMet) is converted to *S*-adenosyl homocysteine (AdoHcy).



## 2. C-Terminal Residues

The second type involves methylesterification in eukaryotic cells of C-terminal residues, specifically [leucine](#) and isoprenylcysteine. Protein **phosphatase 2A** dephosphorylates a number of phosphorylated, regulated enzymes, and it is regulated by methylation of its C-terminal leucine residue. Following the addition of a C<sub>15</sub>-farnesyl or a C<sub>20</sub>-geranylgeranyl group to the side-chain and usually the cleavage of the three terminal amino acids (see [Prenylation](#)), certain cysteine residues now become candidates for carboxyl methylation (Fig. 2). Certain fungal **mating type** factors, **ras proteins**, analogous small **G proteins**, and the g-subunits of large G-proteins are some of the proteins regulated by methylation. The assembly and disassembly of nuclear lamins may be regulated by methylation Even eukaryotic chemotaxis and platelet aggregation may involve reversible protein methylation (3).

**Figure 2.** Reactions involved in the process of C-terminal methylation of cysteine residues in eukaryotic cells. The C-terminal target protein are shown at the top in the one-letter code. FPP is farnesyl diphosphate, and GGPP is geranylgeranyl diphosphate. The methyl donor is *S*-adenosylmethionine (AdoMet), which is released as *S*-adenosylhomocysteine (AdoHcy).



### 3. Basic Residues

Protein methylation also facilitates “permanent” structural changes in proteins to allow them to function even better. Methylation of [lysine](#) residues to mono-, di-, or trimethyl-lysine is common in bacteria and eukaryotic cells. Examples include bacterial **flagellins**, **ribosomal** proteins, [histones](#), [rhodopsin](#) and [calmodulin](#). Calmodulin plays a role in controlling many enzymes, and transgenic **plants** that have unmethylated calmodulin produce poor seed. Trimethyl-lysine has a fixed positive charge and remains charged even in a **hydrophobic** environment.

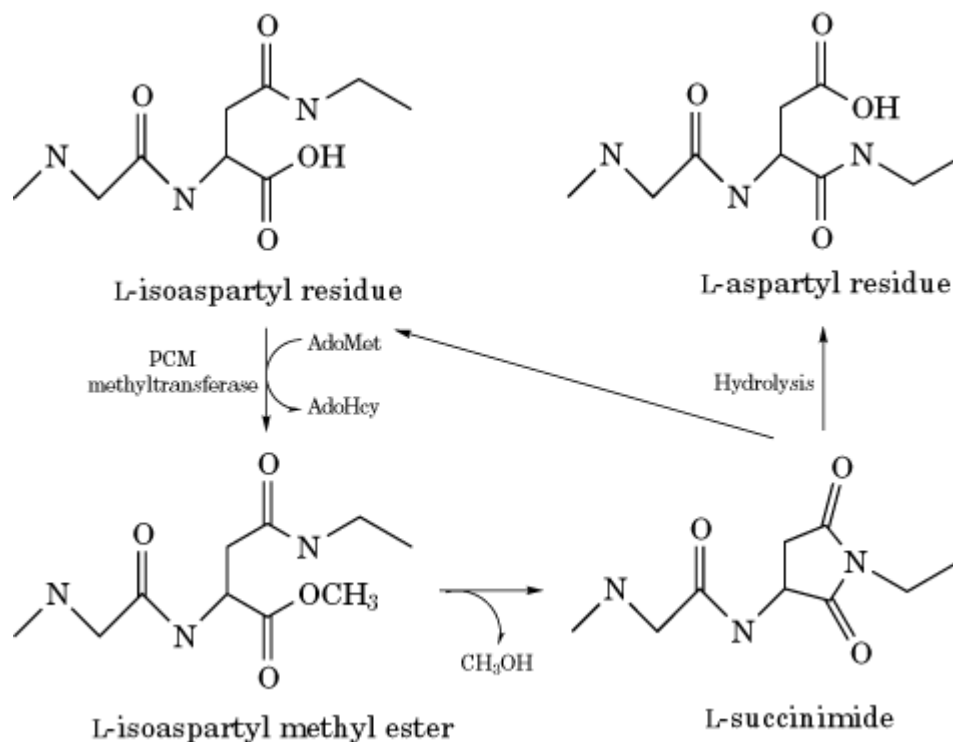
[Arginine](#) residues can be mono- or dimethylated. Examples include myelin basic proteins, **myosin**, **heat shock** proteins, and proteins of the [nucleus](#) and ribosomes. It is common to find these among [RNA-binding proteins](#), possibly because arginine methylation might disrupt certain [hydrogen bonds](#) that would occur between the protein and the nucleic acid (2).

[Histidine](#), also a basic residue, is methylated in [actin](#), myosin, [histones](#), and [rhodopsin](#), but beyond the possibility of furthering interactions with particular proteins, the purpose has not been discerned in any of the cases. One related example is the further metabolism of a particular histidine residue in ribosomal [elongation factor](#) EF-2 that is subject to [ADP-ribosylation](#) by [diphtheria toxin](#). This histidine residue is metabolized to dipthamide in a process involving trimethylation of the  $\alpha$ -amino group. Yeast unable by mutation to make this change grow more slowly. Finally, a particular amide nitrogen on an [asparagine](#) residue in certain phycocyanins and phycoerythrins, involved in [photosynthesis](#) in **cyanobacteria** and red **algae**, may be methylated, possibly to enhance efficiency of energy transfer to the [light-harvesting complex](#) (2).

#### 4. Isomerized Aspartate Residues

As proteins “age,” L-aspartyl residues become **isomerized** to L-isoaspartyl residues or **racemized** to D-aspartyl residues. Then these residues may become methylated, and their spontaneous demethylation sometimes restores the original L-isomer. Repeated cycles eventually restore the native protein (Fig. 3) (2).

**Figure 3.** Isomerization of an aspartate residue and its reversal using methylation. The succinimide intermediate can be generated by deamidation of an asparagine residue.



#### Bibliography

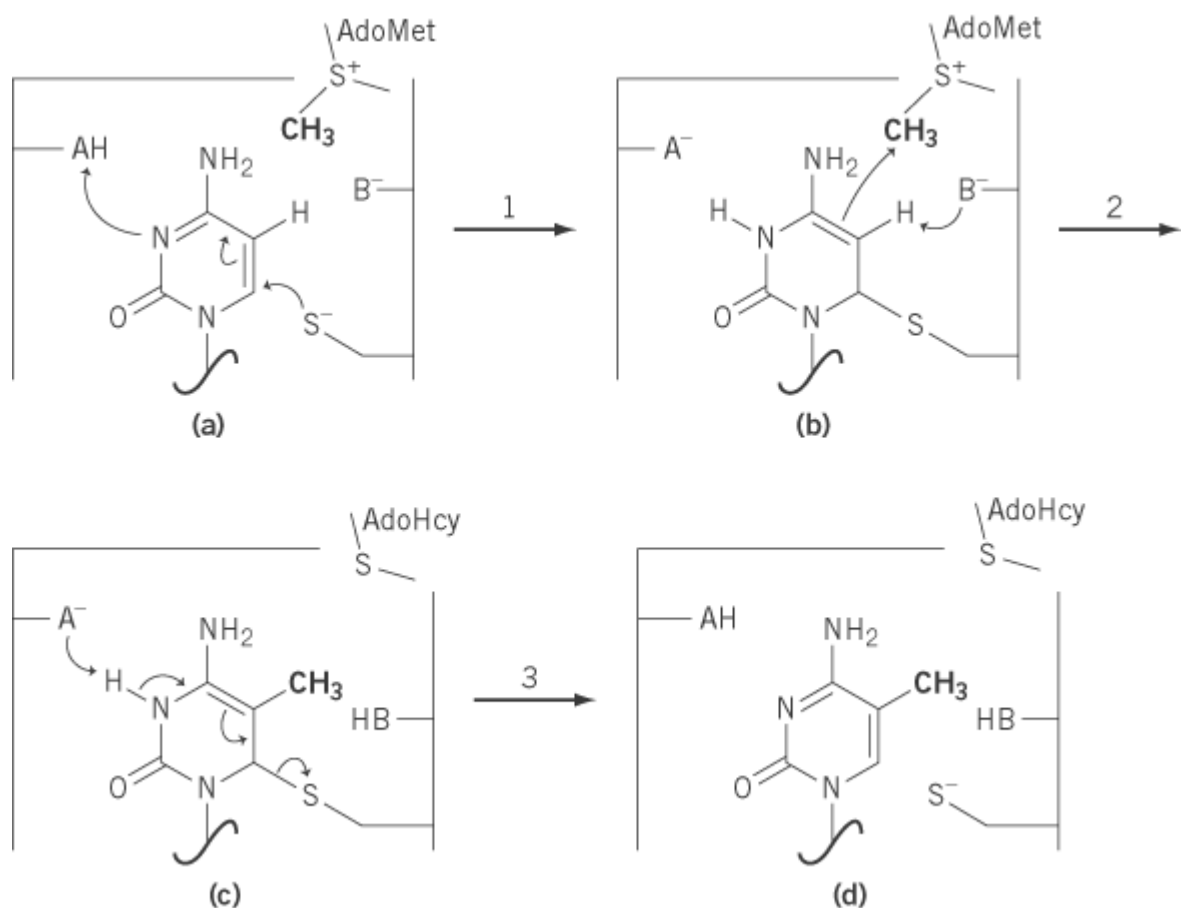
1. P. K. Chiang, R. K. Gordon, J. Tal, G. C. Zeng, B. P. Doctor, K. Pardhasaradhi, and P. P. McCann (1996) *FASEB J.* **10**, 471–480.
2. S. Clarke (1993) *Curr. Opinion Cell Biol.* **5**, 977–983.
3. C. A. Hrycyna and S. Clarke (1993) *Pharmac. Ther.* **59**, 281–300.
4. J. B. Stock and M. G. Surette (1996) In *Escherichia coli and Salmonella, Cellular and Molecular Biology*, 2nd ed. (F. C. Neidhardt, ed.), ASM Press, Washington, DC, Vol. **I**, pp. 1103–1129.

#### 5-Methylcytosine

Cytosine bases in DNA can be modified to 5-methylcytosine (5mC) by DNA **methyltransferases**

(MTases, methylases). DNA methyltransferases are monomeric enzymes that catalyze the transfer and covalent attachment of a methyl group from the cofactor **S-adenosylmethionine** (AdoMet) to DNA with the release of *S*-adenosylhomocysteine (AdoHcy). Thus, methylation occurs at cytosine or adenine bases and forms C5-methylcytosine (5mC), *N*-4-methylcytosine (N4mC) or *N*-6-methyladenine (N6mA) (see [Methyltransferase, DNA](#) for structures). DNA methylation is involved in biological phenomena ranging from prokaryotic DNA replication and host DNA protection (as part of [restriction–modification systems](#)) to eukaryotic **gene regulation** and embryonic [development](#) (see [Methylation, DNA](#)). The [X-ray crystallography](#) structures of two prokaryotic 5mC MTases, M•*Hha*I (1) and M•*Hae*III (2), have increased our understanding of the structural and mechanistic basis of 5mC methylation. The structure of M•*Hha* I (GmCGC), where mC indicates the methylated cytosine in the GCGC DNA recognition sequence, was determined in a complex with AdoMet and in a ternary complex with the cofactor and a 13-bp DNA duplex containing a **suicide substrate**, base analogue, 5-fluorocytosine (5FC). M•*Hae*III (GGmCC) was solved as a binary complex with an 18-bp DNA duplex, also containing 5FC. These complexes are consistent with a catalytic mechanism for 5mC MTases that involves a covalent DNA–protein intermediate. In a mechanism first described by Wu and Santi (3) and later modified by Erlanson, et al. (4), methylation of a cytosine base occurs in three steps (Fig. 1). The second step is prevented when 5FC is the substrate; the electronegative fluorine in place of the hydrogen at C5 cannot be eliminated, and the enzyme is locked in a dead-end, covalent complex with the DNA (5, 6). In contrast, less direct evidence for the mechanism for the amino-methylating, N4mC-forming and N6mA-forming MTases exists. Interestingly, structural comparison of the enzymes M•*Taq*I (a N6mA MTase), and M•*Pvu*II (a N4mC MTase), both crystallized without DNA, with M•*Hha*I, revealed a basic similarity in [active-site](#) architecture that allowed the identification and assignment of catalytic residues involved in a proposed catalytic mechanism for the exocyclic amino-methylating MTases (7). The proposed mechanism involves a direct transfer of the methyl group from AdoMet to the amino group, rather than through a covalent intermediate. The crystal structures of M•*Hha*I and M•*Hae*III have provided snapshots of the mechanism for both the 5mC and the exocyclic amino group methyltransferases.

**Figure 1.** Mechanism of 5mC methylation. (a) Nucleophilic attack on C6 carbon by thiol group of cysteine residue of MTase; (b) abstraction of methyl group from AdoMet and removal of hydrogen at C5; (c) regeneration of thiol group; (d) 5-methylcytosine (indicates bond attaching cytosine to the DNA, and the long straight lines indicate the active site of the protein).



Both *M•HhaI* and *M•HaeIII* are two-**domain** proteins, with the larger domain containing catalytic and cofactor binding sites and the smaller domain involved in DNA recognition. The AdoMet binding regions of *M•HhaI* and *M•HaeIII* contain [Rossman folds](#) for nucleotide binding analogous to the structures of other AdoMet-binding proteins, such as *M•TaqI* and catechol *O*-methyltransferase (2, 6), and a comparison of Ca positions of the *M•HhaI* and *M•HaeIII* DNA MTases shows structural similarity throughout the large domain and bridging region. In contrast, the smaller DNA-binding domains have little structural similarity, except for a small conserved region that serves as a scaffold for the amino acids that interact with the DNA (2). The most remarkable feature in both crystal structures containing DNA is the extrahelical cytosine base, which is flipped out of the DNA duplex and positioned in an active-site pocket in the MTases. For *M•HhaI*, the cytosine is located near AdoHcy in the AdoMet binding pocket. How does the DNA duplex compensate structurally for the cavity created by the flipped base? *M•HhaI* shows significant distortion of the phosphate backbone of the strand in which the extrahelical cytosine is located, along with interdigitation of two amino acid residues from the recognition domain into the cavity opposite the orphaned guanine base. These residues **hydrogen-bond** to the unpaired guanine residue (6). With *M•HaeIII*, on the other hand, flipping out of the base is accompanied by a reorganization of base pairs such that the lone guanine residue hydrogen bonds with a neighboring (3') cytosine base on the opposite strand. This produces another unpaired 3' guanine that hydrogen-bonds to an [arginine](#) residue, but leaves a large unfilled pocket in the DNA duplex (2). The overall structural similarity of C5-cytosine methyltransferases suggests that base flipping may be a common mechanism for this family of enzymes.

The base-flipping mechanism was first found with *M•HhaI*, but several lines of evidence suggest that the extrusion of a base from duplex DNA (or RNA) may be a prevalent occurrence among enzymes that must perform chemistry on bases. Studies using X-ray crystallography on short,

double-stranded DNA containing unpaired bases have shown that the unpaired bases can have both intra and extrahelical positions stabilized by crystal packing (2, 8). NMR studies of imino proton exchange in short, double-stranded DNA have shown the standard **free-energy** change for the base-pair opening to be 7–9 kcal/mol, comparable to protein-induced DNA distortion energies provided by the contacts made at the DNA–protein interface (8). Several DNA-binding enzymes crystallized without DNA, such as T4 b-glycosyltransferase, M• *TaqI*, and M• *PvuII* have been proposed to contain active-site pockets that would accommodate extrahelical bases (9, 10). The crystal structure of a DNA repair enzyme involved in pyrimidine dimer excision, T4 endonuclease V, complexed with duplex DNA containing a thymine dimer, has been found to have an extrahelical adenine base. The expelled adenine base, which is complementary to one of the bases in the thymine dimer, is trapped in a pocket on the surface of the enzyme (11). Although the expelled base in this structure is opposite the strand in which the thymine dimer resides, its exit from the helix renders the damaged bases accessible to the repair enzyme. The detailed mechanism of how the cytosine or adenine bases are extruded remains speculative.

Determination of the mechanism of 5mC methylation and of the crystal structures of M• *HhaI* and M• *HaeIII* provides insight into the structure–function relationships of these proteins. These findings also raise intriguing questions about the DNA dynamics resulting from DNA–enzyme interactions. Aberrant methylation patterns in mammals promote tumorigenesis and lead to developmental abnormalities (12, 13). Since eukaryotic C5-methylcytosine methyltransferases have amino acid sequence similarities to their prokaryotic relatives, the results from studies of the bacterial enzymes will likely contribute to understanding their eukaryotic counterparts.

#### Bibliography

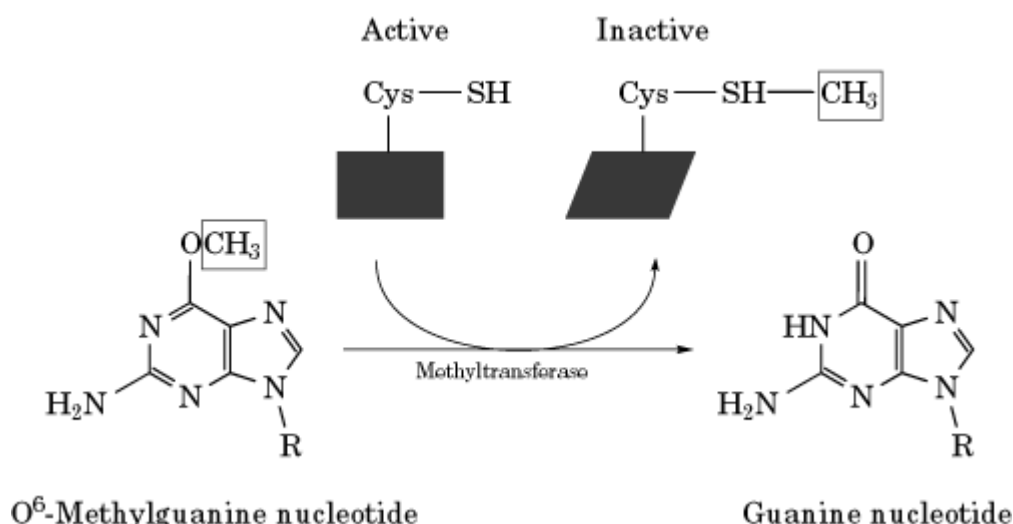
1. S. Klimasauskas, S. Kumar, R. J. Roberts, and X. Cheng (1994) *Cell* **76**, 357–369.
2. K. M. Reinisch, L. Chen., G. L. Verdine, and W. N. Lipscomb (1995) *Cell* **82**, 143–153.
3. J. C. Wu and D. V. Santi (1987) *J. Biol. Chem.* **262**, 4778–4786.
4. D. A. Erlanson, L. Chen, and G. L. Verdine (1993) *J. Am. Chem. Soc.* **115**, 12583–12584.
5. T. H. Bestor and G. L. Verdine (1994) *Curr. Opin. Cell Biol.* **6**, 380–389.
6. X. Cheng (1995) *Annu. Rev. Biophys. Biomol. Struct.* **24**, 293–318.
7. W. Gong, M. O'Gara, R. M. Blumenthal, and X. Cheng (1997) *Nucl. Acids Res.* **25**, 2702–2715.
8. D. Suck (1994) *Curr. Biol.* **4**, 252–255.
9. X. Cheng and R. M. Blumenthal (1996) *Structure* **4**, 639–645.
10. R. J. Roberts (1995) *Cell* **82**, 9–12.
11. D. G. Vassylyev, T. Kashiwagi, Y. Mikami, M. Ariyoshi, S. Iwai, E. Ohtsuka, and K. Morikawa (1995) *Cell* **83**, 773–782.
12. G. L. Verdine (1994) *Cell* **76**, 197–200.
13. P. W. Laird (1996) *Annu. Rev. Genet.* **30**, 441–64.

#### **O<sup>6</sup>-Methylguanine-DNA Methyltransferase (MGMT)**

MGMT is an important enzyme of DNA repair. It carries out direct repair of alkylated DNA by transferring the alkyl group from oxygen atoms within the DNA molecule to a cysteine residue on the enzyme. This reaction is irreversible, and thus the MGMT is said to be a “suicide” enzyme

(Figure 1). The three lesions repaired by MGMT are  $O^6$ -methylguanine,  $O^4$ -methylthymine, and methylphosphotriesters induced by methylating agents such as *N*-methyl- *N'*-nitro-*N*-nitrosoguanidine (MNNG). The enzyme also repairs other alkyl guanines, however, such as  $O^6$ -ethylguanine and  $O^6$ -butyl-guanine, albeit at a much lower efficiency.

**Figure 1.** The “suicide mechanism” reaction catalyzed by MGMT. The enzyme transfers a methyl group from the  $O^6$  position of guanine to the active-site cysteine residue of the MGMT polypeptide chain to restore the normal base and form the stable *S*-methylcysteine adduct, which inactivates the enzyme.



MGMT has been found in all free-living species tested. Enzymes from these sources exhibit considerable sequence [homology](#) and contain the [active-site](#) cysteine residue within the –Pro–Cys–His–Arg–Val–signature sequence. *Escherichia coli* contains two MGMTs: Ada and Ogt. The Ada protein was the first MGMT discovered, and it mediates the “adaptation” response in this organism (1): Exposure to subtoxic doses of an alkylating agent such as MNNG renders cells resistant to subsequent toxic doses of the same or related alkylating agents. Ada is a protein of 39 kDa, and it has two well-defined, and easily separable, NH<sub>2</sub>- and COOH-terminal domains. It performs three functions (2): (i) positive transcriptional regulator, (ii) phosphotriester methyltransferase, and (iii)  $O^6$ -methylguanine methyl transferase. The active-site residue Cys69 in the NH<sub>2</sub>-terminal domain accepts methyl residues from the phosphodiester backbone, to form *S*-methyl-cysteine. Methylation of Cys69 causes a conformational change in Ada and enables it to activate the **promoters** of *ada*, *alkA* (methyladenine DNA glycosylase), *alkB*, and *aidB* genes and thus to increase cellular resistance to alkylating agents. The COOH-terminal domain contains the active-site Cys321, which carries out direct repair by accepting methyl groups from  $O^6$ -MeGua and  $O^4$ -MeThy, to form *S*-methyl-cysteine and to restore the normal bases. In addition to Ada, *E. coli* contains another MGMT called Ogt (3). This enzyme is structurally and functionally homologous to the COOH-terminal domain of Ada; it is expressed constitutively and is not part of the adaptive response.

Like *E. coli*, most other bacteria tested exhibit the adaptive response, although the mechanistic details might vary (see [DNA Damage, Inducible Responses To](#)). In *Bacillus subtilis*, AdaA has methylphosphotriester methyltransferase and transcriptional activator activities, in a manner analogous to the NH<sub>2</sub>-terminal domain of Ada of *E. coli* (4). The AdaB of *B. subtilis*, similarly carries out the MGMT function performed by the COOH-terminal domain of the *E. coli* Ada protein. In *B. subtilis* there is also a noninducible MGMT (Ada C) comparable to the *E. coli* Ogt protein (4).



The eukaryotic MGMTs (including the human enzyme) are structural and functional homologues of Ogt-type methyltransferases ( $M_r \sim 20\text{kDa}$ ); that is, they are suicide enzymes that repair  $O^6$ -methylguanine and  $O^4$ -methylthymine (5). They have no transcriptional regulatory function, nor is there evidence that these enzymes are induced by DNA damage specifically. In humans, increased resistance to anticancer alkylating agents occurs when there is a loss of MGMT activity by gene **silencing**, accompanied by a mismatch-repair defect arising from a mutation in a mismatch-repair gene.

### Bibliography

1. L. Samson and J. Cairns (1977) *Nature* **267**, 281–283.
2. Y. Nakabeppu and M. Sekiguchi (1986) *Proc. Natl. Acad. Sci. USA* **83**, 6297–6301.
3. P. M. Potter, M. C. Wilkinson, J. Fitton, F. J. Carr, J. Brennand, and D. P. Cooper (1987) *Nucleic Acids Res.* **15**, 9177–9193.
4. F. Morohoshi, K. Hagashi, and N. Munatada (1991) *J. Bacteriol.* **173**, 7834–7840.
5. K. Tano, S. Shiota, J. Collier, R. S. Foote, and S. Mitra (1990) *Proc. Natl. Acad. Sci. USA* **87**, 686–690.

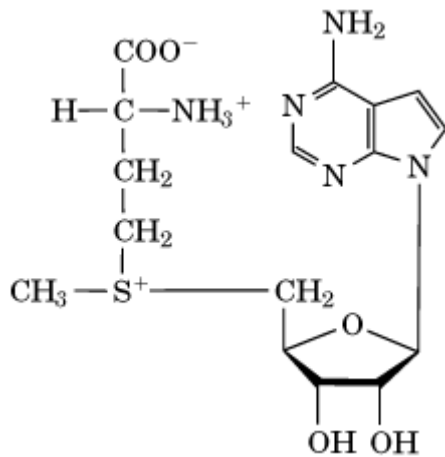
### Suggestions for Further Reading

6. T. Lindahl, B. Sedgwick, M. Sekiguchi, and Y. Nakabeppu (1988) Regulation and expression of the adaptive response to alkylating agents. *Annu. Rev. Biochem.* **57**, 133–157.
7. L. Samson (1992) The suicidal DNA repair methyltransferases of microbes. *Mol. Microbiol.* **6**, 825–831.
8. S. Mitra and B. Kaina (1993) Regulation of repair of alkylation damage in mammalian genomes. *Prog. Nucleic Acids Res. Mol. Biol.* **44**, 109–142.

## Methyltransferase

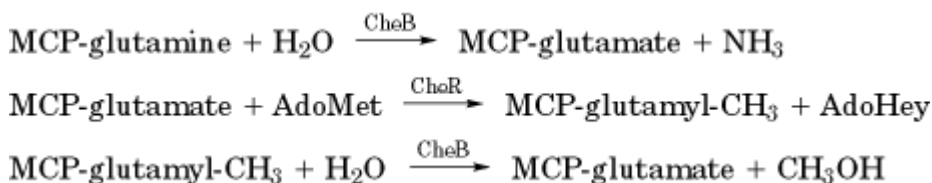
Methyltransferases use **S-adenosylmethionine** (AdoMet) as a methyl donor to catalyze methylation of functional groups on amino acids, usually on their side chains (1). Many different amino acids in proteins are methylated by a host of specific enzymes (see [Methylation, Protein](#)). Several instances of these are described in this article. The methyl group donated from AdoMet is on a sulfonium ion (trivalent sulfur atom having a “formal” positive charge) (Fig. 1) and is readily and irreversibly transferred to other atoms (1). AdoMet is an “expensive” molecule because it is synthesized from [methionine](#) and ATP with production of inorganic phosphate,  $P_i$ , and pyrophosphate,  $PP_i$ ; the latter is rapidly cleaved to  $2 P_i$  by [pyrophosphatases](#). The net cost is three “high-energy” bonds (see [Adenylate Charge](#)), which signifies that maintaining AdoMet levels is a high priority for cells generally, apparently to methylate proteins and DNA (1).

**Figure 1.** Structure of S-adenosylmethionine.



Perhaps the most well characterized methyltransferases are the CheR methyltransferases in bacteria involved in [Chemotaxis](#). These enzymes catalyze the methylesterification of a glutamate (*E*) residue within the protein **consensus sequence** (A/S)XXEEX(A/T/S)A(A/T/S) (2). These residues are located on two  $\alpha$ -helical [coiled coils](#) within the cytoplasmic region of the **receptors**, known as methyl-accepting chemotaxis proteins (2). These methylated residues are located in different places in the receptors of *Bacillus subtilis* and *Escherichia coli* (3). In *E. coli* the effect of methylation is to increase kinase activity. In *B. subtilis*, the position of the site of methylation governs whether CheA kinase activity is increased or decreased. Interestingly, in some instances the glutamate residue is produced from a genetically encoded [glutamine](#) (see [Genetic Code](#)) by the CheB methylesterase acting as a deamidase. This enzyme also hydrolyzes the methyl esters produced by the action of CheR (2) (Fig. 2).

**Figure 2.** Reactions catalyzed by CheR methyltransferase and CheB methylesterase.



In contrast with CheR, the CheB methylesterase of both *E. coli* and *B. subtilis* is activated by phosphoryl transfer from phosphorylated CheA and, hence, becomes stimulated when CheA is activated (2, 4). After autophosphorylation, the CheA kinase leads to increases in the concentrations of the phosphorylated forms of CheY and CheB. The former enhances tumbling in *E. coli* and smooth swimming in *B. subtilis*, whereas the latter causes demethylation of the receptors to reduce the activity of the CheA kinase (4). In other words, upon activation, CheA *simultaneously* causes an excitatory event (due to phosphorylated CheY) and sets in motion the adaptation process (due to the action of phosphorylated CheB). Moreover, in *B. subtilis*, both CheR and CheB are much more active immediately after the addition or removal of attractant rather than later, probably due to slow conformational changes in the receptors exposing the sites of methylation.

Other methyltransferases catalyzing methylesterifications have been identified and characterized. These include the enzyme that catalyzes the metabolically labile methylation of the carboxyl-terminal leucine residue of protein **phosphatase 2A**, an enzyme important in regulating cell

metabolism (7). The methylation and demethylation represent a novel device for regulating the phosphatase itself. The membrane-bound methyltransferase that catalyzes methylation of carboxylterminal isoprenylcysteine residues, which in yeast is the product of the *STE14* gene, has also been characterized, and strains lacking it have been studied (see [Prenylation](#)). Among other consequences, the **a-mating type** factor activity is much reduced (7).

In bacteria, the methyltransferase that catalyzes methylation of isoaspartate residues, an isomer of aspartate that occurs when proteins age, has been identified and the gene disrupted. This enzyme helps with reisolmerizing isoaspartate back to aspartate, restoring the original protein. Bacteria mutant in this gene are sensitive to **heat shock** and survive poorly in a stationary phase, results that imply the importance of this type of protein repair (7).

Arginine methylation in eukaryotes is performed by a specific subfamily of methyltransferases. This family of methyltransferases contains AdoMet binding motif similar to small molecule and nucleic acid methyltransferases (8). Many of the methyltransferases also include a C-terminal domain involved in arginine substrate recognition. The typical recognition motif for the methyltransferase is RGG, RXR, and RG, which can be methylated in three ways: N<sup>G</sup>-monomethylarginine (MMA), N<sup>G</sup>N<sup>G</sup> (asymmetric) dimethylarginine (aDMA), and N<sup>G</sup>N<sup>G</sup> (symetric) dimethylarginine (sDMA). Methylation of these conserved sites is also selective. RGG recognition sites have only been isolated as MMA and aDMA and no sDMA, while RXR and RG motifs have been isolated as aDMA-methylated but not MMA or sDMA. Smd1 and Smd3, spliceosomal small nuclear RNA binding proteins, have been isolated with sDMA methylation on the site GRG leading to the possibility that more recognition motifs exist and provide specificity as well as control of extent of methylation by the arginine methyltransferase.

Histone methylation at specific N-terminal lysine residues is believed to be associated with the epigenic inheritance of heterochromatin (transcriptionally silent DNA), and transcriptional regulation (9). This modification is catalyzed by histone methyltransferases (HMT) on K4, K9, K27, and K36 of histone H3 and K20 of histone H4 in many organisms, although the pattern may be species-specific (9). Like all methyltransferases, HMT contains an AdoMet-binding motif, and, as arginine's guanidino group is multiply methylated, a single lysine can be mono, di, and tri-substituted. Little information is available concerning HMT regulation, although, in developing rat neuron HMT, activity is enhanced, and adult rat neuronal HMT activity is reduced (9).

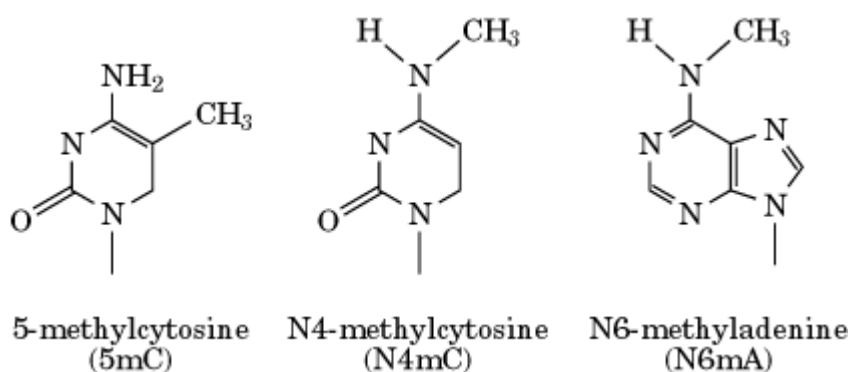
## Bibliography

1. P. K. Chiang, R. K. Gordon, J. Tal, G. C. Zeng, B. P. Doctor, K. Pardhasaradhi, and P. P. McCann (1996) *FASEB J.* **10**, 471–480.
2. J. B. Stock and M. G. Surette (1996) In *Escherichia coli and Salmonella, Cellular and Molecular Biology*, 2 ed. (F. C. Neidhardt, ed.), ASM Press, Washington, DC, Vol. **I**, pp. 1103–1129.
3. D. W. Hanlon and G. W. Ordal (1994) *J. Biol. Chem.* **269**, 14038–14046.
4. L. F. Garrity and G. W. Ordal (1995) *Pharmacol. Ther.* **68**, 87–104.
5. H. Xie and S. Clarke (1994) *J. Biol. Chem.* **269**, 1981–1984.
6. C. A. Hrycyna and S. Clarke (1993) *Pharmacol. Ther.* **59**, 281–300.
7. J. E. Visick and S. Clarke (1995) *Mol. Microbiol.* **16**, 835–845.
8. A. E. McBride and P. A. Silver (2001) *Cell* **106**, 5–8.
9. J. C. Rice and C. David Allis (2001) *Curr. Opin. Cell Biol.* **13**, 263–273.

## Methyltransferase, DNA

DNA methyltransferases [MTases, methylases] recognize specific duplex DNA sequences and catalyze the transfer of a methyl group from **S-adenosylmethionine** (AdoMet) to recipient adenine or cytosine bases within these sequences, with the release of S-adenosylhomocysteine (AdoHcy or SAH). DNA MTases are widespread in nature and serve several physiological roles (see [Methylation, DNA](#)). Prokaryotic DNA MTases function in [mismatch repair](#) and [DNA replication](#) and as part of [restriction–modification systems](#). Eukaryotic DNA MTases are involved in **gene regulation**, **development**, and **genome** compartmentalization. MTases are classified according to the base they methylate (Fig. 1). The first class methylates the C5 position of cytosine to form [5-methylcytosine](#) (5mC). The second class methylates exocyclic nitrogen atoms of adenine or cytosine to form N6-methyladenine (N6mA) or N4-methylcytosine (N4mC). The second class is further divided into different groups, a, b, or g on the basis of amino acid sequence differences and structural analysis (1, 2). Over 100 MTase genes have been **cloned** and sequenced. Analysis of several 5mC MTase amino acid sequences reveal a linear arrangement of 10 conserved motifs. One motif, a glycine-rich segment, is also found in other AdoMet-dependent MTases and is involved in **cofactor** binding (3). Another conserved segment, PC, is unique to the 5mC MTases and is a part of the [active site](#). Analysis of the primary sequences of N6mA and N4mC MTases reveals two motifs, where the arrangement and distance between the motifs defines the subgroup classifications (2). The first motif is the glycine-rich segment analogous to those in 5mC MTases and the second a tetrapeptide sequence, (Asp/Asn/Ser)–Pro–Pro–(Tyr/Phe), which has been implicated as part of the active site.

**Figure 1.** Types of methylated bases in DNA.



Discussion here will be limited to prokaryotic type II methyltransferases. For type I, IIs, III, and eukaryotic methyltransferases see [Restriction–Modification Systems](#) or **DNA methylation**.

### 1. Crystal Structures and Sequence Comparison

The crystal structures of four prokaryotic type II MTases have been determined. Two 5mC MTases, M•*HhaI* (4) with AdoMet and M•*HaeIII* (5), without AdoMet, have been solved in association with their target duplex DNAs. One N4mC MTase, M•*PvuII*, and one N6mA MTase, M•*TaqI*, have both been solved in complexes with AdoMet (6, 7). Comparison of the crystal structures of the four enzymes, as well as the results of a **proteolysis** study using M•*EcoRI* (8), reveal that all these monomeric enzymes have a two-**domain** organization, having the cofactor-binding and catalytic site and DNA recognition region associated with different domains. The structures indicate similar patterns of folding and organization at the cofactor binding and active sites (5, 9). In fact, the active site and cofactor-binding regions of M•*HhaI*, M•*PvuII*, and M•*TaqI* are found to be very similar when the central **b-sheets** of their structures are superimposed on each other (3, 6). On the other

hand, the DNA recognition domains of M• *HhaI* and M• *HaeIII* have very dissimilar structures, except for two regions of amino acids that provide a scaffold for residues contacting the DNA (5). The common catalytic domain structures led Malone, et al. (1) to do a structure-assisted amino acid sequence analysis of 42 N6mA and N4mC MTases. This analysis revealed nine variably conserved motifs analogous to those found in 5mC MTases (I–VIII and X.) (1). Motifs I–III and X were hypothesized to be responsible for forming a cofactor-binding pocket and interacting with AdoMet. Motifs IV, VI, and VIII (the conserved PC motif in 5mC MTases and (Asp/Asn/Ser)–Pro–Pro–(Tyr/Phe) in N6mA and N4mC MTases) were proposed to participate in catalysis. Motifs V and VII were postulated to provide a structural framework for the catalytic region. Motif IX may provide a structural framework for the DNA recognition region in 5mC MTases, but no analogous motif was found in the exocyclic amino-methylating MTases (1, 3). Furthermore, the N6mA MTases analyzed by Malone et al. fell into three groups, a, b, or g depending on the order of the motifs from the N to C terminus of the protein, supporting the earlier sequence analysis by Wilson (2). For Dam MTase, a representative N6mA<sub>a</sub> MTase, motifs involved in binding AdoMet are first, followed by sequences involved in DNA recognition, and finally by catalytic-site motifs. Group b, which contains most of the N4mC MTases (including M• *PvuII*), has an arrangement where catalytic-site motifs are first, followed by the DNA recognition region, and then the cofactor-binding motifs. Finally, group g, including M• *TaqI*, has an arrangement of cofactor-binding motifs, followed by the catalytic site motifs, and terminating with the DNA recognition region. In general, the solved structures support the proposals of Malone et al. (1).

## 2. DNA Recognition

MTases are monomeric enzymes that transfer one methyl group to a recipient base of one strand at the recognition sequence per binding event. Most MTases show no preference for hemimethylated or unmethylated DNA and will methylate either molecule at the same rate. MTases appear to search for their canonical recognition sites by facilitated diffusion, or sliding, along the DNA until the recognition sequence is located (10). Kinetic studies of MTases, including M• *HhaI* (a C5 MTase) (11), M• *EcoRI*, and *Escherichia coli* Dam MTase (N6mA MTases) (12, 13), show that bound cofactor is necessary for specific sequence recognition. X-ray crystallography of M• *HhaI* (14) and footprinting analysis of M• *EcoRV* (15), M• *SssI*, and M• *HhaI* (16) indicate a majority of the interactions with DNA occur with bases and phosphates in the major groove. Studies of M• *EcoRI* (17) and M• *EcoRV* (18) with base analogues indicate that some interactions may also occur in the minor groove. For some N6-methyladenine methylases, site-specific recognition may be facilitated by DNA bending. Whereas the crystal structures of M• *HhaI* and M• *HaeIII* show little bending of the DNA helix, gel-shift-mobility assays and scanning force microscopy techniques have shown that M• *EcoRI* (19) and M• *EcoRV* (20) bend DNA about 52° and 60°, respectively. A striking feature, which may be a recognition mechanism as well as a catalytic property of all MTases, was found in the X-ray crystal structures of the 5mC MTases M• *HhaI* and M• *HaeIII*. The substrate cytosine base is extrahelical and flipped from the helix into a pocket in the enzyme, where it is subsequently methylated and released. This extrahelical base flipping is considered in greater detail in 5-methylcytosine. Others have used fluorescence-based assays of M• *EcoRI* (21) or modified base interference studies with M• *EcoRV* (22) to provide evidence for a base-flipping mechanism for N6mA MTases. The structures of both M• *TaqI* and M• *PvuII*, although lacking DNA, have a “pocket” into which the target base could be flipped (6, 7). The N4mC MTases have sequence similarities to the N6mA MTases. All the DNA MTases likely require AdoMet binding and utilize phosphate and base interactions in combination with base flipping and possibly DNA bending to accomplish sequence-specific recognition.

## 3. Catalytic Properties

The steady-state kinetic analyses performed on MTases are consistent in that methylation by MTases follows an ordered bi–bi steady-state scheme, where AdoMet binding is a prerequisite to canonical DNA sequence recognition and methylation (11, 12, 23): (1) the MTase binds to nonspecific DNA

randomly, with or without AdoMet bound; and (2) in the presence of AdoMet, which confers recognition specificity, the target sequence is located and the recipient base of one strand is methylated. The product ternary complex is MTase-methylated DNA-AdoHcy. AdoHcy then dissociates, and depending on the processivity of the MTase (a measure of the ability to scan base pairs), the enzyme will dissociate from the DNA or move to the next site on the same DNA molecule (24). The overall reaction of DNA MTases is slow, with turnover rates ( $k_{\text{cat}}$ ) generally less than  $0.1\text{s}^{-1}$  (12). The average  $K_{\text{M}}$  values for AdoMet and substrate canonical DNA are in the nanomolar range. MTases are extremely efficient catalysts, with  $k_{\text{cat}}/K_{\text{M}}$  values reaching **diffusion-controlled** limits ( $10^8\text{--}10^9\text{ M}^{-1}\text{s}^{-1}$ ). The rate-limiting step for MTase catalysis occurs after the methylation step, since  $k_{\text{cat}}$  is slower than the rate of methyl group deposition,  $k_{\text{meth}}$  (11, 12, 23). The mechanism of the 5mC MTases has been elucidated and involves a covalent enzyme intermediate (see [5-Methylcytosine](#) for a detailed description of the mechanism). The mechanism of N4mC and N6mA MTases likely involves the direct methylation without a covalent enzyme intermediate (6, 25). The structure of M•*PvuII*, an N4mC-forming MTase, revealed the structural similarity of its active site to those of M•*HhaI* and M•*TaqI*, which form 5mC and N6mA, respectively (6). The similar locations of amino acids in the active sites of these three MTases that form all three kinds of methylated base allowed the assignment of specific residues to a catalytic mechanism for M•*PvuII*. This mechanism, which probably does not involve a covalent enzyme–DNA complex, can likely be extended to enzymes like M•*TaqI* that form N6mA as well as to other N4mC-forming MTases.

#### 4. Applications

Methyltransferases, or modification enzymes, can be used to modulate cleavage by restriction endonucleases (see [Restriction Enzymes](#)). Methylation of the DNA within or outside a specific restriction endonuclease recognition site prevents or inhibits cleavage by interfering with endonuclease binding or catalysis. For instance, MTases can be used decrease the number of possible restriction sites on a long DNA molecule. For example, R•*BanII* cleaves at the degenerate sequences (GAGCTC) or (GGGCCC). Premethylation of DNA with M•*HaeIII*, which is specific for (GG5mCC), reduces subsequent R•*BanII* cleavage to the single sequence, (GAGCTC), because the other sequence is rendered refractory. Also, particular restriction sites can be protected from methylation by masking the restriction site with the use of DNA-binding proteins or oligonucleotides (that form triplex DNA regions). These proteins or oligonucleotides can subsequently be removed to unmask the chosen restriction site. In another technique, methylase-limited partial digestion, larger fragments of DNA can be generated by a restriction endonuclease through prior partial methylation with the cognate MTase (26).

In summary, DNA MTases are ubiquitous and simple in composition and cofactor requirements. Thus, MTases are useful for the study of DNA–protein interactions, enzyme mechanism and kinetics, and protein–cofactor interactions. They are also useful tools for the molecular biologist.

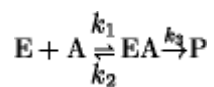
#### Bibliography

1. T. Malone, R. M. Blumenthal, and X. Cheng (1995) *J. Mol. Biol.* **253**, 618–632.
2. G. G. Wilson (1992) *Meth. Enzymol.* **216**, 262.
3. X. Cheng (1995) *Annu. Rev. Biophys. Biomol. Struct.* **24**, 293–318.
4. S. Klimasauskas, S. Kumar, R. J. Roberts, and X. Cheng (1994) *Cell* **76**, 357–369.
5. K. M. Reinisch, L. Chen., G. L. Verdine, and W. N. Lipscomb (1995) *Cell* **82**, 143–153.
6. W. Gong, M. O'Gara, R. M. Blumenthal, and X. Cheng (1997) *Nucl. Acids Res.* **25**, 2702–2715.
7. J. Labahn et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10957–10961.
8. N. O. Reich, K. A. Maegley, D. D. Shoemaker, and E. Everett (1991) *Biochemistry* **30**, 2940–2942.

9. M. O' Gara, D. McCloy, T. Malone, and X. Cheng (1995) *Gene* **157**, 135–138.
10. M. A. Surby and N. O. Reich (1996) *Biochemistry* **35**, 2201–2208.
11. J. C. Wu and D. V. Santi (1987) *J. Biol. Chem.* **262**, 4778–4786.
12. N. O. Reich and N. Mashhoon (1991) *Biochemistry* **30**, 2933–2939.
13. S. Marzabal et al. (1995) *Nucl. Acids Res.* **23**, 3648–3655.
14. S. Klimasauskas, S. Kumar, R. J. Roberts, and X. Cheng (1994) *Cell* **76**, 357–369.
15. M. D. Szczelkun, H. Jones, and B. A. Connolly (1995) *Biochemistry* **34**, 10734–10743.
16. P. Renbaum and A. Razin (1995) *J. Mol. Biol.* **248**, 19–26.
17. C. A. Brennan, M. D. Van Cleve, and R. I. Gumpport (1986) *J. Biol. Chem.* **261**, 7279–7286.
18. P. C. Newman et al. (1990) *Biochemistry* **29**, 9891–9901.
19. R. A. Garcia, C. J. Bustamante, and N. O. Reich (1996) *Proc. Natl. Acad. Sci. USA* **93**, 7618–7622.
20. S. Cal and B. A. Connolly (1996) *J. Biol. Chem.* **271**, 1008–1015.
21. B. W. Allan and N. O. Reich (1996) *Biochemistry* **35**, 14757–14672.
22. S. Cal and B. A. Connolly (1997) *Jour. Biol. Chem.* **272**, 490–496.
23. L. Szilak, A. Der, F. Deak, and P. Venetianer (1993) *Eur. J. Biochem.* **218**, 727–733.
24. M. Surby and N. O. Reich (1996) *Biochemistry* **35**, 2201–2208.
25. D. H. Ho, J. C. Wu, D. V. Santi, and H. G. Floss (1991) *Arch. Biochem. Biophys.* **284**, 264–269.
26. A. Pingould, J. Alves, and R. Geiger (1993) in *Methods in Molecular Biology*, Vol. **16**, *Enzymes of Molecular Biology*, M. M. Burrell, ed., Humana Press, Inc., Totowa, N.J., p. 167.

## Michaelis–Menten Kinetics

Reactions catalyzed by nonallosteric, classical [enzymes](#) E on substrate A to convert it to product P occur by the general mechanism



and are described by the general steady-state rate equation

$$v = \frac{VA}{K_a + A} \tag{1}$$

where  $v$  and  $V$  represent initial and maximum velocities, respectively;  $K_a$  is the [Km \(Michaelis constant\)](#) for substrate  $A$ , which is present at concentration  $A$ . This equation will apply to multisubstrate reactions, provided that only a single substrate is varied and all others are held at fixed concentrations. Reactions described by the equation are considered to obey Michaelis–Menten kinetics.

It is important to recognize that for steady-state conditions to apply, the concentration of substrate must be sufficiently greater than that of the enzyme, by at least about an order of magnitude, so that formation of the enzyme–substrate complex does not significantly reduce the concentration of free

substrate. Under these conditions, the concentration of free substrate can be equated to the total substrate concentration, and the concentration of the enzyme–substrate complex will remain essentially constant during the time over which the initial velocity of the reaction is measured.

## 1. Characteristics of the Michaelis-Menten Equation

The Michaelis–Menten equation is a saturation function identical to that which would be obtained for the interaction of an [antigen](#) with an [antibody](#) or of a **ligand** with a **receptor**. Thus, a plot of velocity as a function of substrate concentration yields a rectangular hyperbola. There are two fundamental parameters associated with the equation:

1. When  $A$  tends to infinity,  $v = V$  and the reaction is zero-order with respect to  $A$ .  $V$ , the maximum velocity of the reaction, will be dependent on the concentration of enzyme  $E_t$  and the first-order rate constant  $k_3$ . This constant is equal to  $V/E_t$  and is the [turnover number](#) for the enzyme, which gives a measure of the number of moles of product produced per mole of enzyme (or subunit) per second. The magnitude of turnover numbers varies considerably as illustrated in [Table 1](#).
2. When  $A$  is very small and tends to zero,  $v = (V/K_a)A$ . In this case, the initial velocity of the reaction is first-order with respect to  $A$ , and  $V/K_a$  represents an apparent first-order rate constant for the interaction of substrate with enzyme. This apparent rate constant has units of  $M^{-1}s^{-1}$  and can be expressed simply in terms of rate constants as

$$\frac{V}{K_a E_t} = \frac{k_1 k_3}{k_2 + k_3} = \frac{k_1}{1 + k_2/k_3} \quad (2)$$

It is only when  $k_3 \gg k_2$  that a true value for the bimolecular rate constant  $k_1$  would be observed. The rates for the formation of enzyme–substrate complexes vary considerably and can approach the diffusion rate, which is in the vicinity of  $250 \text{ mM}^{-1} \text{ s}^{-1}$  ([Table 1](#)) (see [Diffusion-Controlled Reactions](#)).

**Table 1. Values of Kinetic Parameters for Selected Enzymes**

| Enzyme                                  | Substrate                         | $K_a$<br>(mM) | $V$ ( $s^{-1}$ ) | $V/K_a$ ( $\text{mM}^{-1} \text{s}^{-1}$ ) |
|---|-----------------------------------|---------------|------------------|--|
| <b>Chymotrypsin</b>                     | acetyl-L-tryptophanamide          | 5000          | 100              | 0.02                                       |
| <a href="#">Lysozyme</a> , hen          | hexa- <i>N</i> -acetylglucosamine | 6             | 0.5              | 0.08                                       |
| Chorismate mutase                       | chorismate                        | 100           | 47               | 0.5  |
| Prephenate dehydrogenase                | prephenate                        | 65            | 95               | 1.5  |
| <a href="#">Dihydrofolate reductase</a> | dihydrofolate                     | 1             | 20               | 20   |
| <a href="#">Beta-Lactamase</a>          | benzylpenicillin                  | 50            | 2000             | 40   |
| Fumarase                                | fumarate                          | 5             | 800              | 160  |



The Michaelis constant is the ratio of the two fundamental parameters  $V$  and  $(V/K_a)$ , so it is not a fundamental parameter, but it is nevertheless an important parameter for the reasons outlined in [Km \(Michaelis constant\)](#).

## Microfilament

Microfilaments are the helical polymers found in eukaryotic cells that are formed from [actin](#). These filaments may contain only actin, or actin and associated proteins. In vertebrate striated muscle such filaments also contain tropomyosin, troponin, and nebulin (see [Thin Filament](#)). The term microfilament is largely an historical artifact, as it was used to denote the perceived morphological difference between these filaments, [intermediate filaments](#), and the thick filaments of striated muscle that are formed from myosin. While the actin-based filaments were seen by various forms of [electron microscopy](#) to be about 65 Å in diameter, in contrast to the ~100 Å diameter of the intermediate filaments, we now know that the actin filament is also approximately 100 Å in diameter ([1](#), [2](#)).

### Bibliography

1. E. H. Egelman and R. Padron (1984) *Nature* **307**, 56–58.
2. M. Lorenz, D. Popp, and K. C. Holmes (1993) *J. Mol. Biol.* **234**, 826–836.

## Microinjection

Microinjection is the technique of injecting materials directly into the interiors of individual cells. It has many applications, including the introduction of: (i) [messenger RNA](#) and [proteins](#) into frog **oocytes** and cultured mammalian cells to investigate [translation](#) and protein function; (ii) **DNA** into **zygotes** or germ cells to produce **transgenic** organisms; (iii) mouse embryonic [stem cells](#) into blastocysts to produce **chimeric** mice (see [Transgenic Technology](#)); (iv) dyes or colloids into a cell of an [embryo](#) to trace **cell division** and [differentiation](#) during [development](#); and (v) **sperm nuclei** into **oocytes** to treat human male infertility. Depending on the application, systems for microinjection vary in their degree of sophistication, but all consist of a microneedle inserted into an instrument holder, a pressure source for expelling the injectate, a microneedle micromanipulator, a microscope, and a microinjection chamber. The system may need to be placed on a vibration dampening support.

### 1. Microneedles

Microneedles are made from thin-wall borosilicate capillary glass tubing using a micropipette-puller. The dimensions of the pulled tip vary according to the application and preferences of the injector; puller settings are established according to manufacturer's guidelines and, to a significant extent, empirically. Bending of microneedles for accommodating them into microinjection chambers is achieved using the heated filament of a microforge or a small flame. The final tip diameters for injecting solutions into mammalian cells are 0.1 to 0.5 μm, and those for injecting solutions into mammalian zygote pronuclei are 0.5 to 2.0 μm. The approximate tip diameter in the lower range can

be gauged under a high-power objective (1). With this method it is difficult to resolve diameters less than  $\sim 0.5\mu\text{m}$ , but if the tip diameter can be resolved it is probably too large for microinjection of mammalian cells (2). If pullers cannot be calibrated to pull medium tip diameters, between 0.5 and  $2.0\mu\text{m}$ , tips can be chipped by bringing them at right angles against a glass bead fused to the filament of a microforge. Alternatively, they can be beveled with a micropipette-beveler. The smaller the tip diameter, the more readily it will be clogged with particulates arising from either the microneedle or the injectate. Particulates can be removed by washing and siliconizing the unpulled glass capillaries (1, 2) and by centrifuging the injectate at  $\sim 15,000\text{ g}$  for 15 min. A fiber or filament within the microneedle allows the tip to be filled from the back with the material to be injected. It is deposited inside the microneedle close to the tip, using a drawn capillary attached by tubing to a mouthpiece. It then flows along the fiber to the tip by capillary action. Alternatively, the injectate can be drawn into the tip from a small drop under the microscope if an oil-filled optical transmission system is used.

The final tip diameter for injecting mouse embryonic stem cells into blastocysts should be  $\sim 20\mu\text{m}$ . The sharpest tips are produced by breaking the pulled glass at the final diameter with a microforge, beveling the tip, then making a small spike at the tip of the bevel with a microforge (3). For simplicity, however, most injectors produce bevels by resting the pulled glass on a silicon pad and cutting at the desired diameter with a scalpel blade under a stereomicroscope (4). Similar tips may be used for the microinjection of spermatozoa under the zona pellucida or directly into the cytoplasm of mammalian oocytes, although smaller diameters are used (5-7).

## 2. Pressure Sources

### 2.1. Compressed Gas

A simple system commonly used for injecting mammalian zygote pronuclei is air-filled and consists of a 50-mL glass syringe connected by tubing to the microneedle holder. Injectate is expelled by pushing the plunger in by hand; and as long as the plunger is not pulled out, there is no risk that medium will be drawn into the microneedle and dilute the injectate. An automated system connected to a cylinder of compressed gas and providing three pressures (injection, back and purge), with a foot-operated switch between injection and back pressure, can be made in a workshop (8) or purchased commercially. Back pressure prevents entry of medium into the microneedle, and purge pressure is useful for clearing cell debris or particulates from the tip. Such a unit is desirable for microinjection of cultured mammalian cells, because the greatest control of flow of injectate can be obtained. Nevertheless, consistency of the volume in successive injections is not guaranteed, since the tip diameter alters with the collection of cell debris.

### 2.2. Hydraulic

The pressure source usually consists of a micrometer syringe connected by tubing to the microneedle holder. The whole system is filled with either light or heavy paraffin oil (the heavier oil being the most stable), silicon oil of similar viscosity, Fluorinert, or water. Injectate is deposited into microneedles as described above, and the rest of the space in the microneedle is filled with oil or Fluorinert from the back with a 30-gauge steel needle attached to a syringe. Such a system has sufficient control for injecting nuclei of mammalian cultured cells (2, 9), although control of flow is not as great as with the automated three-setting pressure system, and is also used for injecting mouse embryonic stem cells into blastocysts. For injecting solutions into cells, tubing for connecting the micrometer syringe to the microneedle holder is often hard-walled and of small diameter, although similar control of flow can be obtained with semihard wall tubing up to 6 mm in diameter connected to a micrometer syringe with a larger barrel. DNA may also be expelled from the micropipette, without the use of pressure, by iontophoresis (10).

## 3. Micromanipulators

Fully mechanical or hydraulic-mechanical combination units employ ratio reduction mechanics to reduce the movement generated by the hand. When at rest, these units are subject to gravitational

drift. Electronic units employ motors that drive in steps as small as ~150 nm and therefore have the greatest resolution of movement. Because they are driven indirectly by the hand, however, they do not respond exactly with the time and speed of hand movement. Although all three types of units are suitable for all applications, mechanical and hydraulic–mechanical combination units are most suited for use with mammalian ova, because rapid and sudden movements are often required, and electronic units are most suited for injection of mammalian cultured cells, because precise movement over small distances is important. Micromanipulators that fit onto the microscope stage are not subject to the stage being fixed in the vertical axis, but in this case the whole of the stage must be movable in the horizontal plane.

#### 4. Microscope

Inverted microscopes are usually employed for the injection of mammalian zygotes and cultured cells. The large working distance between the condenser lens and the stage accommodates a large range of designs of microinjection chambers and microneedles. Cultured cells can be injected conveniently with the microneedle entering an open dish at a steep angle from above. Upright microscopes with long working distance objectives and a stage fixed in the vertical axis are more limited for chamber design, but they can be used more conveniently for injecting nuclei of cultured mammalian cells under high-power oil-immersion objectives. Differential interference contrast (DIC) optics provide the best resolution of cellular and nuclear membranes, but are expensive and must be used with glass chambers. Hoffman Modulation Contrast optics create a similar effect to DIC providing a less expensive, but lower-resolution, alternative and can be used with plastic chambers. Phase contrast optics can also be used for injecting nuclei of mammalian ova and cultured cells but do not provide good resolution of membranes. Nevertheless, these optics are preferred by many for injecting embryonic stem cells into blastocysts (see [Microscopy](#)).

#### 5. Microinjection Chambers

This design of this component of the microinjection system varies greatly, being subject to personal preference. A versatile chamber is one in which the medium containing the target cells is sandwiched between a slide and a coverslip and is surrounded by oil ([9](#)). This chamber accepts microneedles from the side, and, if raised a few millimeters on the stage, no bending of microneedles for entry into the medium is required. Optical performance is optimal at all points, because the droplet of medium is flat on the upper and lower surfaces. Also, the medium is extremely stable with respect to movement of the microneedle. Apart from being useful for the manipulation of mammalian ova, this chamber can be used for the injection of cultured mammalian cells attached to the coverslip under high power oil immersion optics ([2](#), [9](#)). Other chambers used for the manipulation of mammalian ova include the hanging drop chamber ([11](#)) and the simple depression slide chamber ([12](#)).

#### 6. Vibration Dampening

The microneedle must be stable when viewed under the microscope, because vibration will increase the frequency of membrane lysis. If vibration occurs, microinjection systems can be placed onto heavy wood or marble tables or onto metal plates resting on top of rubber, inflated innertubes, or squash balls. Although these supports give some dampening, they do not completely eliminate vibration. Commercially available devices, such as gas cushion tables, usually do eliminate vibration and are a fraction of the cost of the microinjection system as a whole.

#### 7. Microinjection in Nonmammalian Organisms

For organisms commonly used in research, methods of microinjection have been described for the frog ([13-15](#)), [nematode](#) worm ([16](#)), **Drosophila** ([17](#)), and [zebrafish](#) ([18](#)) and are essentially the same as described above.

#### Bibliography

1. M. L. DePamphilis (1988) *BioTechniques* **6**, 662–680.
2. R. H. Lovell-Badge (1987) In *Teratocarcinomas and Embryonic Stem Cells: A Practical Approach* (E. J. Robertson, ed.), IRL Press, Oxford, pp. 153–182.
3. C. L. Stewart (1993) *Meth. Enzymol.* **225**, 823–855.
4. A. Bradley (1987) In *Teratocarcinomas and Embryonic Stem Cells: A Practical Approach* (E. J. Robertson, ed.), IRL Press, Oxford, pp. 113–151.
5. J. R. Mann (1988) *Biol. Reprod.* **38**, 1077–1083.
6. G. Palermo, H. Joris, P. Devroey, and A. C. Van Steirteghem (1992) *Lancet* **340**, 17–18.
7. Y. Kimura and R. Yanagimachi (1995) *Biol. Reprod.* **52**, 709–720.
8. W. Ansorge (1982) *Exptl. Cell Res.* **140**, 31–37.
9. E. G. Diacumakos (1973) *Methods Cell Biol.* **7**, 287–311.
10. C. W. Lo (1983) *Mol. Cell. Biol.* **3**, 1803–1814.
11. R. L. Gardner (1979) In *Methods in Mammalian Reproduction* (J. C. Daniels, ed.), Academic Press, New York, pp. 137–165.
12. B. Hogan, R. Beddington, F. Costantini, and E. Lacy (1994) *Manipulating the Mouse Embryo: A Laboratory Manual*, 2nd ed., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 189–252.
13. J. B. Gurdon and L. Wakefield (1986) In *Microinjection and Organelle Transplantation Techniques, Methods and Applications* (J. E. Celis, A. Graessmann, and A. Loyter, eds.), Academic Press, New York, pp. 269–299.
14. M. J. M. Hitchcock, E. I. Ginns, and C. J. Marcus-Sekura (1987) *Methods Enzymol.* **152**, 276–284.
15. C. J. Marcus-Sekura and M. J. M. Hitchcock (1987) *Methods Enzymol.* **152**, 284–288.
16. C. Mello and A. Fire (1995) *Methods Cell Biol.* **48**, 451–482.
17. M. Ashburner (1989) *Drosophila. A Laboratory Handbook*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 223–228.
18. M. Westerfield (1994) *The Zebrafish Book. A Guide for the Laboratory Use of Zebrafish (Brachydanio rerio)*, 2.1 ed., University of Oregon Press, Eugene, OR, pp. 5.1–5.30.

### Suggestions for Further Reading

19. N. D. Allen, S. C. Barton, M. A. H. Surani and W. Reik (1987) "Production of transgenic mice". In *Mammalian Development: A Practical Approach* (M. Monk, ed.), IRL Press, Oxford, U.K., pp. 217–233.
20. J. W. Gordon (1993) Production of transgenic mice. *Methods Enzymol.* **225**, 747–771.
21. J. R. Mann (1993) Factors influencing production frequency of transgenic mice. *Methods Enzymol.* **225**, 771–782.

### Micronucleus, Macronucleus

Ciliated protozoa (eg, *Paramecium* and *Tetrahymena*) contain a small micronucleus and a larger macronucleus. The micronucleus, although quiescent in **RNA** production, is the repository for the intact **genome**, which is usually **diploid** in ciliates. Conversely, the macronucleus is transcriptionally active, but it contains a much larger amount of DNA, all of which exists as small fragments of about

2000 base pairs (bp) that are repeated hundreds to thousands of times. All of the genes necessary for survival of the ciliate are present in and expressed from these small macronuclear DNA fragments.

The life cycle of *Tetrahymena* describes the biogenesis and fate of micro- and macronuclei. Upon conjugation between two *Tetrahymena* cells, the micronucleus in each cell undergoes **meiosis** to form four **haploid** micronuclei. Three of the haploid micronuclei degenerate, and the remaining one undergoes a single round of **mitosis**. The old macronucleus is still present at this point. The mating cells exchange a single micronucleus (pronucleus) through the cytoplasmic bridge connecting the two cells. The haploid micronucleus that enters the other cell fuses with the haploid micronucleus that remained behind there (**fertilization**), thus reestablishing the diploid state. After fusion, the diploid micronucleus of the **zygote** undergoes two rounds of mitosis, leading to four diploid micronuclei in each cell (which have now separated and are called exconjugates). Two of these diploid micronuclei become the precursors of new macronuclei, while the old macronucleus from the previous generation breaks down. One of the other two micronuclei degenerates, leaving only one of the original four. At this point in the life cycle, each exconjugate cell contains two macronuclei and the one remaining micronucleus. Then both exconjugate cells divide. In the process, the micronucleus undergoes mitosis, but the two macronuclei do not. The end result is that each cell contains one macronucleus and one micronucleus.

The conversion of a micronucleus into a macronucleus requires **chromosomal** breakage, followed by genomic DNA deletion and rearrangement (1). In *Tetrahymena*, chromosomal breakage is mediated by a specific 15-bp *cis*-acting element known as a **chromosomal breakage sequence** (Cbs). The element is found in multiple copies but only within the micronuclear genome. The Cbs elements are necessary and sufficient for breakage, and they are lost from the micronuclear genomic DNA upon cleavage. *Trans*-acting Cbs-recognition factors work in conjunction with **nucleases** to cut the DNA and then with **telomerase** to stabilize the ends. Next deletion of DNA (up to 15% in *Tetrahymena*) occurs in segments ranging from several hundred to >13kbp. The segments to be eliminated are AT-rich and are flanked on their ends by short direct repeats that do not play a role in the excisional mechanism. Instead, flanking A<sub>5</sub>G<sub>5</sub> sequences within the retained DNA (~45bp from the actual boundaries) direct where the cuts are made. Sequence elements internal to the eliminated DNA are also necessary to promote the deletion process. At least one *trans*-acting protein, Pdd1p, recognizes these internal sequences (2). Coyne *et al.* (1) present a model in which *trans*-acting factors bind the flanking A<sub>5</sub>G<sub>5</sub> sequences and then interact with Pdd1p to mediate excision.

Differences between the levels of transcription in macronuclei and micronuclei suggest differences in **chromatin** composition, particularly with regard to **histone** H<sub>1</sub>. Similar to histone H<sub>1</sub> proteins from other **eukaryotes**, that in *Tetrahymena* macronuclei associates with the linker DNA between **nucleosomes**, is soluble in 5% perchloric acid, and dissociates from chromatin at salt concentrations lower than the other histone proteins (3). It is unusual, however, in that it lacks the central **hydrophobic** globular domain characteristic of other H<sub>1</sub> histones. The micronucleus, on the other hand, has no such H<sub>1</sub> protein discernible. Instead, the micronucleus contains proteins designated a, b, and g that interact with the linker DNA, but display tryptic **peptide maps** that differ from those of the macronuclear H<sub>1</sub> species (4). Macronuclear-specific histone variants have been identified, further exemplifying fundamental biochemical differences between micro- and macronuclei, both of which exist within the same cell.

## Bibliography

1. R. S. Coyne, D. L. Chalker, and M.-C. Yao (1996) *Ann. Rev. Genet.* **30**, 557–578.
2. M. T. Madireddi, R. S. Coyne, J. F. Smothers, K. M. Mickey, M.-C. Yao, and C. D. Allis (1996) *Cell* **87**, 75–84.
3. M. A. Gorovsky (1986) In *The Molecular Biology of Ciliated Protozoa* (J. G. Gall, ed.) Academic Press, Orlando, pp. 227–261.

4. C. D. Allis, C. V. C. Glover, and M. A. Gorovsky (1979) Proc. Natl. Acad. Sci. USA **76**, 4857–4861.

### Suggestion for Further Reading

5. J. G. Gall (ed.) (1986) *The Molecular Biology of Ciliated Protozoa*, Academic Press, Orlando. Collection of comprehensive reviews of ciliate micro- and macronuclear biology.

## Microsatellite DNA

*Microsatellites* and **minisatellites** have become valuable tools in genetic mapping, forensic identity testing, and population studies as a result of their abundance, relatively uniform distribution, and high degree of **polymorphism**.

The human [genome](#) contains approximately 50,000 copies of an interspersed dinucleotide repeat with the sequence  $(T-G)_n$ , where  $n$  is between 10 and 60. Such repeats have been found in many regions that have been fully sequenced. As with minisatellite DNA regions with longer repeat elements, they often display extensive length polymorphisms (1). Other microsatellite polymorphic sequences, such as  $(C-A)_n$  and  $(G-T)_n$ , have been found in humans (2).

### 1. Use of Microsatellites for Genome Sequencing

Microsatellites are present in the genomes of all eukaryotes. They have been analyzed for size variation by the polymerase chain reaction (**PCR**) and [gel electrophoresis](#) in humans and in mice, among others. If they are sequenced and taken as the substrate for the polymerase chain reaction (see below), they become highly informative locus-specific markers. By analogy with the **sequence-tagged sites** (STS), which serve for the standardization of physical genetic maps, they are termed *sequence-tagged microsatellite sites* (STMS). Genetic maps based on STMS share the advantage with STS-based maps that mapping vocabularies are standardized to the DNA sequence base and that access to any particular locus does not require shipping or storing cloned probes (3).

Microsatellites are sequences consisting of a repeated very simple motif. The most widespread in the human genome is  $(C-A)_n$ , where  $n$  varies from 5 to about 50. They are interesting because they are very polymorphic. If such a sequence exists in a certain site of the genome, it is likely to be  $(CA)_{17}$  in one given individual,  $(CA)_{18}$  or  $(CA)_{15}$  in another one, and so on. These sequences are very widespread in the human genome because, on average, there is one every 10 kb. Therefore, they constitute perfect landmarks for a detailed genetic map. However, the polymorphisms cannot be detected by a simple [Southern blot](#), as can restriction fragment length polymorphism (**RFLP**). First, it is necessary to find a microsatellite in the region where a marker is desired, for example, by screening **cosmids** containing this region with a  $(CA)_{10}$  probe (because of the average spacing of microsatellites, one is practically certain to detect some). After identifying the fragment containing the microsatellite, it must be sequenced to define the unique DNA flanking sequences. From this information, primers must be synthesized by PCR to target the required region of the genome specifically. Analysis of the DNA from a series of individuals requires amplifying each DNA by PCR and analyzing the amplification product by [electrophoresis](#) on a highly resolving polyacrylamide gel, because one must be able to distinguish length variations of two nucleotides in a fragment of 100 to 200 nucleotides. Although more complicated than classic Southern blotting, this

method can be partially automated and is worthwhile because of the power of the method. The markers are almost always informative. Since 1992, linkage maps based on microsatellite markers are regularly put at the disposal of the scientific community (4-6). The most recent human genetic map is based on 5264 short tandem(A-C/T-G)<sub>n</sub> repeat polymorphisms, spans a genetic distance of 3,699 centiMorgan (cM), and comprises 2335 positions. The majority of the position can be ordered with a probability of at least 1000:1 against alternative orders. The average interval size is 1.6 cM; 59% of the map is covered by intervals of 2 cM at most, and only 1% remains in intervals above 10 cM (5).

A search of the DNA sequence databases revealed that the location of dinucleotide microsatellites is often conserved among mammalian species, enabling the prediction of the presence of microsatellites by comparative genetic data. In closely related species, this conservation was adequate to allow PCR primers designed for use in one species to be used to analyze microsatellite length polymorphism in the other. The ability to use heterologous PCR primers, coupled with comparative map information, facilitates the use of DNA microsatellites in gene mapping studies in closely related species, such as cattle and sheep, rat, and mouse, or primates (7).

The same type of length variation is found in nuclear genomes and also in the [chloroplast](#) genomes of many plants (8).

## 2. Possible Origin of the Microsatellites

A number of microsatellite repeats have been found associated with the **Alu** type of repeated [interspersed DNA elements](#). The association of an Alu sequence with a microsatellite repeat could result from the integration of an Alu element within a preexisting microsatellite repeat, or Alu elements could have a direct role in the origin of microsatellite repeats. A microsatellite repeat could result from errors introduced during **reverse transcription** of the primary **transcript** derived from an Alu “master gene” or from the accumulation of random mutations in the middleA- rich regions and oligo(dA)-rich tails of Alu elements after insertion and subsequent expansion and contraction of these sequences. These hypotheses have been tested by direct evolutionary comparison of the sequences of some recent Alu elements, found only in humans and absent from nonhuman primates, and older elements present at orthologous positions in a number of nonhuman primates. That the Alu elements were a source for the genesis of primate microsatellite repeats is suggested by the origin of the “young” Alu insertions, the absence of sequences that resemble microsatellite repeats at orthologous loci in chimpanzees, and the gradual expansion of microsatellite repeats in some old Alu repeats at orthologous positions within the genomes of a number of nonhuman primates (9).

## Bibliography

1. M. Litt and J. A. Luty (1989) *Am. J. Hum. Genet.* **44**, 397–401.
2. J. L. Weber et al. (1991) *Genomics* **11**, 695–700.
3. J. S. Beckmann and M. Soller (1990) *Biotechnology* **8**, 930–932.
4. J. Weissenbach et al. (1992) *Nature* **359**, 794–801.
5. G. Gyapai et al. (1994) *Nature Genetics* **7**, 246–339.
6. C. Dib et al. (1996) *Nature* **380**, 152–154.
7. S. S. Moore et al. (1991) *Genomics* **10**, 654–660.
8. W. Powell et al. (1995) *Curr. Biol.* **5**, 1023–1029.
9. S. S. Arcot et al. (1995) *Genomics* **29**, 136–144.

## Microscopy

The applications of microscopy to molecular biology are rapidly expanding as new microscopes and techniques are developed. Historically, the field of microscopy has been divided into two major realms—light microscopy and [electron microscopy](#). However, with the advent of new microscopies, such as atomic force, scanning tunneling, X-ray, NMR, and acoustic microscopies (1), the need for a new classification is evident, perhaps along the lines of the structural features observable with each kind of microscope. The first question to ask is, “How does one choose the type of microscopy to use for a particular specimen?” To answer this question, one must first know the approximate size of the specimen and the level of structural detail desired. Key aspects of important subfields of microscopy of interest to molecular biologists are treated here along with applications which illustrate the power of each technique.

### 1. Light Microscopy

#### 1.1. Transmission Light Microscopy

It is important to realize that *bright field (Kohler) microscopy* (as is true for all microscopies) allows for the observation of a magnified image of the specimen and not the specimen itself (2). The imaging process is achieved through the interference of light waves issuing ultimately from the back of the objective lens of the light microscope. Kohler illumination is produced by first positioning a lens in front of the light source so that its image is *not* in the specimen plane. A second (condenser) lens places an image of the surface of the first lens onto the specimen with as short a focal length as possible to obtain as wide a cone of light as possible to illuminate the specimen. The purpose of this illumination is first, to achieve an evenly lit field of view against which the specimen detail can be recognized and second, to attain the maximum resolution of fine detail by having a wide cone of radiation. Specimen visualization with this kind of microscopy relies on stains or contrast in the specimen itself. Fluorescent dyes offer more sensitive visualization than simple absorption contrast.

If the specimen has insufficient contrast for bright field microscopy, then *phase contrast microscopy* may be tried. Examples are unstained specimens which can be described as phase objects because they do not significantly reduce the transmitted intensity, yet differ in refractive index from the background. Hence elements in the transmitted beam can be separated according to differences in the optical path traversed (a product of the distance traveled and the index of refraction). Interference between the reunited elements will then produce detectable differences in intensity in the image. Phase contrast microscopy has remarkable sensitivity, capable of distinguishing differences in optical pathlength of  $7/360$  of a wavelength (2). However, it has a disadvantage in that the image of each feature is surrounded by a halo of bright light. This halo degrades the resolution of small objects. Phase contrast can be useful for weakly stained specimens, but is not suitable for heavily stained or thick specimens, since the interference from multiple diffractions is uninterpretable.

*Nomarski or differential interference contrast (DIC) microscopy* is superb for looking at living cells because it is noninvasive and high contrast is produced along the edges of biological structures (3). This type of microscopy detects abrupt changes of refractive index occurring over a small distance. Thus, for example, [mitochondria](#), nuclear [membranes](#), and lipid droplets are clearly visualized. The images produced with a Nomarski microscope appear three-dimensional because one side of the specimen appears lighter than the other side; this 3-D effect is actually an artifact of the polarized light used. Another useful property of Nomarski imaging is that thin optical sections can be obtained, and, unlike phase contrast imaging, thick specimens can be optically sectioned. Real-time optical sectioning allows for the observation of organelle movement. Nomarski microscopy was used to investigate spiral [chloroplast](#) movement in a filamentous alga (4).

### Bibliography



1. P. J. Duke and A. G. Michette (1990) *Modern Microscopes: Techniques and Applications*, Plenum Press, New York.
2. A. J. Lacey (1989) *Light Microscopy in Biology*, IRL Press, Oxford.
3. D. J. Rawlins (1992) *Light Microscopy*, Bios Scientific, Oxford.
4. E. G. Jordan and D. J. Rawlins (1990) *J. Cell Sci.* **95**, 343–351.

## Microtubule-Associated Proteins (MAPs)

A large number of structurally diverse microtubule-associated proteins (MAPs) are bound to the surface of [microtubules](#) in cells. It is widely accepted that many of these MAPs play critical roles in the regulation of microtubule polymerization and dynamics, in the organization of microtubule arrays, and in the functional interactions of microtubules with other cellular components. MAPs from different species, as well as from different cells and tissues of the same species, differ widely in their molecular properties and functions ([1](#)).

The best-studied MAPs to date are several neural MAPs whose functions appear to involve stimulation of microtubule polymerization during formation of neuronal cell processes and stabilization of microtubule dynamics in mature axons and dendrites. One such group of MAPs, called tau proteins, is a family of closely related 55- to 62-kDa phosphoprotein isoforms found predominantly in axons that arise by [alternative splicing](#) of the [messenger RNA](#) of a single tau gene. Tau may be involved in the etiology of Alzheimer's disease; a hyperphosphorylated form of tau is the major protein component of the paired helical filament lesions found in brains of patients with this disease ([2](#), [3](#)). A second well-studied MAP is MAP2, a 200-kDa protein that exists in two similar isoforms which, like the tau isoforms, also arise by alternate exon splicing of a single MAP2 gene. MAP2 is expressed in neurons and associates with microtubules in dendrites and cell bodies, but not in axons. Both tau and MAP2 contain remarkably similar microtubule-binding domains in their C-terminal segments and very different arm-like projection domains in their N-terminal segments. The different projection domains of MAP2 and tau presumably reflect the different functions of tau and MAP2 in axons and dendrites. MAP2 and tau are multiply **phosphorylated** on serine and threonine residues in their microtubule-binding domains by a number of **Ser/Thr protein kinases** *in vitro*, including MAP kinase, **cyclic AMP**-dependent protein kinase, and calcium-**calmodulin**-dependent protein kinase, and phosphorylation strongly influences the ability of these two MAPs to bind to microtubules and to modulate their polymerization *in vitro*. A great many other MAPs exist in different cell types and tissues, such as MAP 4, a fairly well-characterized high-molecular weight MAP prominent in mammalian cells ([4](#)). Most MAPs, however, remain to be characterized.

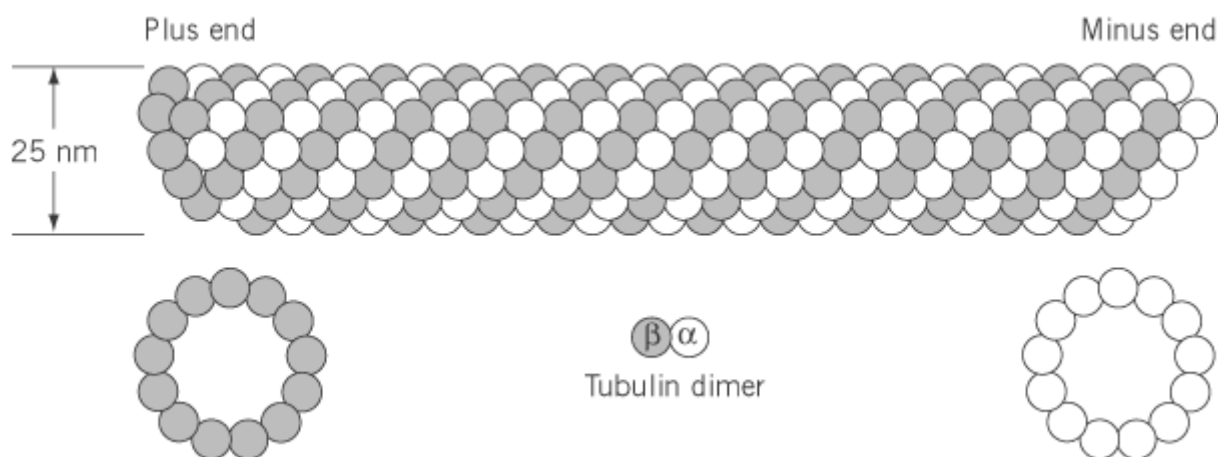
## Bibliography

1. G. Wiche, C. Oberkanins, and A. Himmler (1991) *Int. Rev. Cytol.* **124**, 217–273.
2. J. Q. Trojanowski, et al. (1993) *Clinical Neurosci.* **1**, 184–191.
3. M. Goedert, R. Jakes, M. G. Spillantini, and R. A. Crowther (1994) In *Microtubules* (J. Hyams and C. Lloyd, eds.), Wiley-Liss, New York, pp. 183–200.
4. J. C. Bulinski (1994) In *Microtubules* (J. Hyams and C. Lloyd, eds.), Wiley-Liss, New York, pp. 167–182.

## Microtubules

Microtubules are tube-shaped [protein](#) polymers, approximately 25 nm in diameter, that can vary from a few to many micrometers in length (Fig. 1). They are present in virtually all eukaryotic cells but are not found in prokaryotic cells. Microtubules are a major component of the cell's [cytoskeleton](#) and, along with [actin](#) filaments and [intermediate filaments](#), play critical roles in the determination of cell shape, in the organization of the cell's cytoplasm, and in many forms of cellular motility (1). For example, microtubules are prominent in the nervous system, where they are critical for the development, structural organization, stability, and functions of the axonal and dendritic processes of neurons. They are also critical for mitosis, the process during the [cell cycle](#) in which the previously duplicated DNA in the form of highly condensed [chromosomes](#) is distributed equally with exquisite accuracy to the daughter cells at the time of cell division. The basic building block unit of the microtubule is [tubulin](#), a 100-kDa,  $\sim 0.4\text{ nm} \times 0.5\text{ nm} \times 0.8\text{ nm}$  dimer composed of two 50-kDa subunits known as  $\alpha$ -tubulin and  $\beta$ -tubulin. A large number of additional cellular molecules, called [microtubule-associated proteins](#) (MAPs), bind to the surface of the microtubule and determine and control its many functions. Microtubules accomplish their multiple tasks in diverse kinds of cells by a combination of their polymerization dynamics, their interactions with specific MAPs, and their interactions with motor molecules, which transport cargo along the surface of the microtubule. The functions, cellular organizations, and biochemical and molecular properties that determine the ability of microtubules to carry out their many functions are summarized here.

**Figure 1.** Structure of a typical microtubule. The  $\beta$ -tubulin subunits are darkly shaded and the  $\alpha$ -tubulin subunits are lightly shaded. The  $\beta$ -tubulin subunits are exposed at the plus ends (17).



### 1. Microtubule Structures, Microtubule Organizing Centers, and the Stability and Organization of Microtubule Arrays

Microtubule polymers are composed of  $\alpha\beta$ -tubulin heterodimers arranged in a uniformly repeated head-to-tail fashion, with the long axis of the heterodimer oriented parallel to the long axis of the microtubule (Fig. 1). Thus, all of the  $\beta$ -tubulin subunits face one end of the microtubule, and the  $\alpha$ -subunits face the opposite end. The ends are designated either plus or minus, based upon their different dynamic and structural properties. Current evidence suggests that the  $\beta$ -subunits are exposed at the plus ends while the  $\alpha$ -subunits are exposed at the minus ends. Each linear array of

tubulin heterodimers is called a *protofilament*, and in cells most microtubules consist of 13 protofilaments arranged into a tube structure. However, many examples exist in which microtubules have a different number of protofilaments (from 11 to as many as 15).

Microtubules in cells are usually organized by an organelle called the *microtubule organizing center* (MTOC) or *centrosome* (2). The minus ends are tethered at the MTOC and the plus ends are away from it, often in a radial distribution. In cultured animal cells, for example, most microtubules are arranged radially in the cell, all or most emanating from the centrosome, which consists of a pair of *centrioles* surrounded by a matrix of poorly understood composition and organization. One of the known proteins of the MTOC is a form of tubulin called *g-tubulin*, which does not form microtubules itself and does not become incorporated into microtubules, but does appear to function in the nucleation of microtubule formation. Not all microtubules in different types of cells and arrays are anchored to centrosomes, and the origins of such noncentrosomal microtubules are not yet known. Interestingly, the MTOCs in cells from higher plants do not contain centrioles, but they still function as microtubule organizing centers.

Single microtubules abound in a great many cellular microtubule arrays, but many specialized microtubule arrays also exist in which microtubules are fused along their lengths into doublet or triplet microtubule structures (3). Doublet microtubules consist of 10 or 11 protofilaments of a second microtubule fused to an entire 13-protofilament first microtubule, and triplet microtubules have an additional 10 or 11 protofilaments of a third microtubule fused to the second microtubule. Doublet microtubules occur in the long shafts of motile cell appendages called *cilia* and **flagella**. Triplet microtubules occur in the basal body structures that in most animal cells are situated at the base of cilia and flagella and in centrioles, which are situated at the center of animal cell centrosomes.

One of the hallmarks of cytoplasmic microtubules is their rapid dynamics. Microtubules are assembled and disassembled from a large pool of cytoplasmic tubulin subunits, and because of their rapid dynamics, many cytoplasmic microtubules are very labile. There is a continuous exchange of tubulin subunits between the soluble tubulin pool and the microtubules. Half-times for turnover of tubulin in cytoplasmic microtubules can be remarkably short, no more than 15 to 30 s, as for example with microtubules in the mitotic spindles of cultured animal cells (4). Such microtubules polymerize and depolymerize rapidly in response to the cell's needs. Another hallmark of dynamic cytoplasmic microtubules is their sensitivity to cold temperatures; cooling cells to 0°C induces rapid disassembly of labile microtubules. However, many microtubules present in the cell's cytoplasm are relatively stable. For example, microtubules in axons of neurons serve as passive tracks for the transport of cell constituents between the nerve cell body and the nerve ending. Such axonal microtubules, stabilized by specific neuronal MAPs, exchange their tubulin with tubulin in the cytoplasmic pool very slowly. Some microtubules, such as the doublet and triplet microtubules of cilia, flagella, centrioles, and basal bodies, are so stable that, once formed, they do not detectably exchange their tubulin with tubulin in the cytoplasmic pool. It is clear, then, that cells are capable of modifying the stability of their microtubules to accomplish specific functions. Such control of stability appears to be mediated primarily by MAPs that bind to the microtubule surfaces and ends.

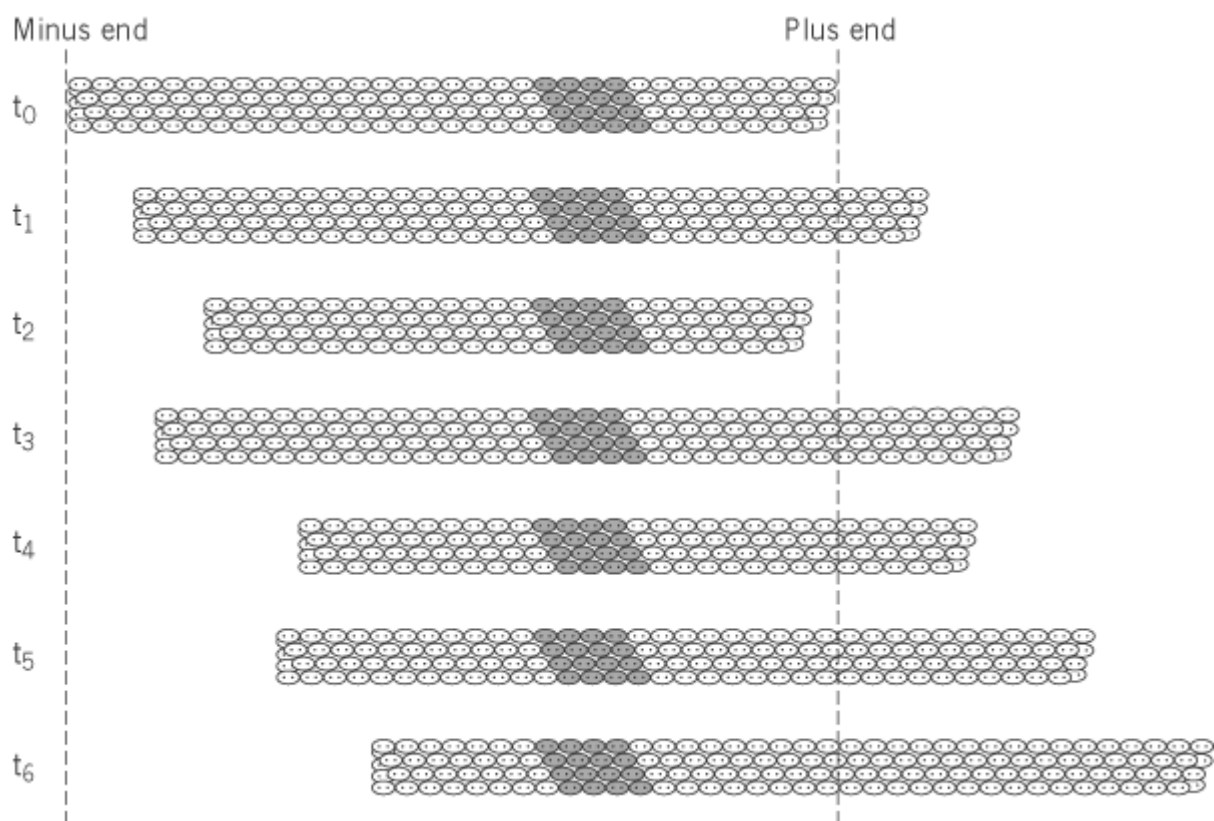
Another important characteristic of microtubules is their ability to form large organized microtubule arrays. Such arrays are especially common in protists, where bundles or sheets of parallel-aligned microtubules, held together by various MAPs, function in the creation of extremely long asymmetric cell structures, such as the dynamic microtubule array in the 400- $\mu$ m-long axopods in the unicellular organism *Echinospaerium*, or the axostyle, a remarkable motile ribbon of cross-bridged microtubules that serves as a motility organelle in the flagellate *Saccinobaculus* (5).

## 2. Microtubule Polymerization and Dynamics

Microtubule polymerization occurs by a nucleation-elongation mechanism, in which the formation of a short microtubule “nucleus” is followed by growth of the microtubule at both its ends by the

reversible noncovalent addition of tubulin subunits. It is well established that microtubules are not simple equilibrium polymers. Soluble tubulin binds GTP reversibly at a site on the b-subunit, and the GTP becomes hydrolyzed to GDP and  $P_i$  as, or shortly after, the tubulin polymerizes onto a growing microtubule end (6, 7). The irreversible hydrolysis of GTP during tubulin addition creates polymers with unique dynamics. One such behavior is called *treadmilling* (8-10), which is the net addition of tubulin at one microtubule end and net loss of tubulin at the opposite end. A second behavior is called *dynamic instability*, which is characterized by phase transitions, or switching, between phases of relatively slow growing and rapid shortening at the ends of individual microtubules (10, 11). Both kinds of dynamics appear to be highly important in specifying microtubule function in cells. Originally, dynamic instability and treadmilling were thought to be incompatible behaviors, but it is now known that at or near steady state *in vitro* and in cells, both behaviors coexist in microtubule populations. While individual microtubules may continually change length, net tubulin addition occurs at the plus end and loss at the minus within an individual microtubule, and in a balanced fashion within a microtubule population as a whole (Fig. 2).

**Figure 2.** Simultaneous treadmilling and dynamic instability of microtubules at polymer-mass steady state. Shown are consecutive “snapshots” of a microtubule exhibiting growth and shortening at its plus and minus ends, with net growth at the plus ends and net shortening at the minus ends. Individual microtubules can change length, as shown here, but net growth occurs at the plus end and net shortening at the minus end. The shaded subunits represent a marked segment.  $T_0$  = zero time;  $t_1$ ,  $t_2$ , and  $t_3$ ,  $\frac{1}{4}$  are arbitrary equal increments of time.



The mechanism responsible for the nonequilibrium dynamics of microtubules involves the gain and loss of a stabilizing “cap” at the microtubule ends, which involves in some way the irreversible hydrolysis of GTP. While the size, chemical composition, and mechanism responsible for the stochastic gain and loss of the stabilizing cap at microtubule ends are poorly understood, evidence with brain microtubules indicates that the cap may be very small, consisting of only a single layer of GTP- or GDP- $P_i$ -liganded tubulin situated at the extreme microtubule ends (see Refs. 6 and 7).

Dynamic instability is believed to be due to a stochastic loss and regain of the stabilizing cap, with growth occurring when the cap is present, and rapid shortening when the cap is lost. Treadmilling is believed to be due to the overall differences in the critical subunit concentrations at the opposite ends of the microtubule.

### 3. Mechanistic Basis of Microtubule-Mediated Cell Function

Microtubules serve two primary functions in eukaryotic cells. They function in the development and maintenance of cell shape, particularly in the creation of asymmetric cell shape, and they are used for many forms of cellular movement. Microtubule function at a mechanistic level is probably best understood in the beating of **cilia and flagella**. The microtubules of cilia and flagella are extremely stable, and their movement is not complicated by dynamics. In contrast to our understanding of microtubule-mediated movement in cilia and flagella, it has been difficult to unravel the mechanistic basis of motility mediated by labile cytoplasmic microtubules. In contrast to the stable microtubules of cilia and flagella, most microtubules in the cell cytoplasm are dynamic and it is clear that the dynamics of such microtubules, not just their presence as passive supports for motors, are integrally linked to their functions. Intriguing recent evidence indicates that microtubule motors and polymerization dynamics may work together to effect certain microtubule-dependent movements such as the complex and highly varied movements of the chromosomes during the processes of mitosis and meiosis (eg, see Ref. [12](#)).

It is well established that the dynamic instability behavior of microtubules plays an important role in many microtubule-dependent processes in cells including the attachment of mitotic spindle microtubules to the chromosomes and determination of cell polarity. In addition to dynamic instability behavior, rapid treadmilling of microtubules occurs in interphase microtubule arrays ([9](#)), and rapid flux dynamics consistent with treadmilling occurs during metaphase and anaphase of mitosis (eg, see Ref. [13](#)). During metaphase of mitosis, for example, tubulin addition occurs at the microtubule plus ends, attached at the kinetochores of the chromosomes, and balanced tubulin loss occurs at the minus ends, which are tethered at the centrosomes, while the microtubule lengths remain approximately constant. The points of microtubule attachment in the spindle are therefore not static. Instead, these regions must be mechanical protein machines capable of transiently attaching and detaching from the microtubules as the microtubules grow or shorten. The function of treadmilling in cells is not yet known. Treadmilling may be involved with maintaining tension in spindles or with translocation of molecules (signals?) from kinetochores to centrosomes. How microtubule dynamics and the dynamics-dependent functions are determined and controlled in cells at a mechanistic level, and how the dynamics and motors might work together to mediate movement, are poorly understood and are the focus of intense current research.

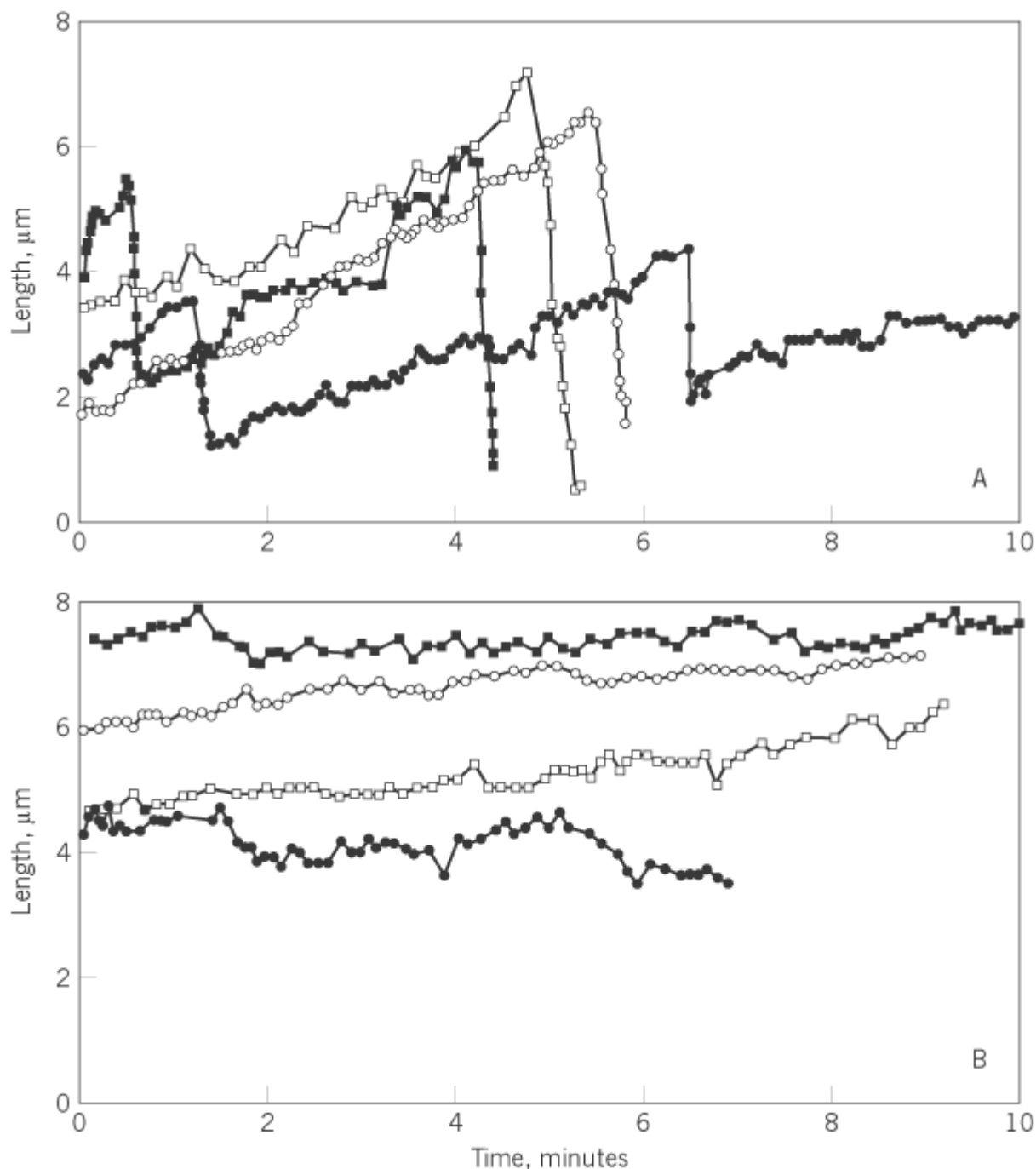
### 4. Microtubule-Targeted Drugs: The Use of Drugs for Analyzing Microtubule Function and for the Treatment of Cancer

Antimitotic drugs have become important tools in cell biology for the study of microtubule function. For example, [colchicine](#), [vinblastine](#), nocodazole, and taxol have been used extensively by investigators to determine whether specific cell processes involved the participation of microtubules. At relatively high concentrations, most of these drugs inhibit microtubule polymerization, while taxol strongly promotes microtubule polymerization. Thus, inducing microtubule depolymerization or excessive polymerization in cells has revealed whether cell processes are dependent upon microtubules.

Low concentrations of these drugs have a subtle action; they interact with microtubule ends and surfaces and powerfully suppress treadmilling and dynamic instability (eg, see Refs. [14](#) and [15](#)). Suppression of microtubule dynamics occurs at drug concentrations well below those required to depolymerize or excessively polymerize microtubules. For example, low concentrations of vinblastine, when added to bovine brain microtubules *in vitro*, powerfully suppress the rate and extent of growing and shortening at microtubule plus ends, plus the frequencies at which the ends

switch between the growing and shortening state, without appreciably affecting the polymer mass (Fig. 3). Thus, these drugs have become valuable tools for analyzing the roles of dynamics in cellular microtubule functions, such as in the complex movements of chromosomes during mitosis.

**Figure 3.** Life history plots of individual MAP-depleted bovine brain microtubules at their plus ends at steady-state *in vitro*. The traces shown are typical examples of the changes in length of the individual microtubules with time. (a) Control microtubules and (b) microtubules in the presence of 0.4  $\mu\text{M}$  vinblastine. (See Ref. 18.)



Many antimetabolic drugs that act on microtubules, including vinblastine, vincristine, vinorelbine, taxol, taxotere, and estramustine, are used in humans for the treatment of various forms of cancer. It appears that the mechanism responsible for the ability of antimetabolic compounds to inhibit cell proliferation and kill tumor cells is the kinetic stabilization of spindle microtubule dynamics. HeLa

cells, for example, are blocked at the metaphase–anaphase transition of mitosis at very low concentrations of taxol (16) or vinblastine. While such blocked cells have normal or nearly normal mitotic spindles and chromosome distributions, they cannot progress into anaphase. The powerful antiproliferative and antitumor effects of drugs that act by suppressing spindle microtubule dynamics may involve the metaphase–anaphase cell-cycle checkpoint, which prevents cells from progressing into anaphase until the spindle is fully assembled and the chromosomes are properly poised for separation. It seems that a prolonged block at the metaphase–anaphase transition in sensitive tumor cells, due to suppressed spindle microtubule dynamics, triggers [apoptosis](#), thus killing the cells.

### Bibliography

1. P. Dustin (1984) *Microtubules*, 2nd revision, Springer-Verlag, Berlin, pp. 1–482.
2. D. M. Glover, C. Gonzalez, and J. W. Raff (1993) *Sci. Am.* **June**, 62–68.
3. K. Fujiwara and L. G. Tilney (1975) *Ann. N.Y. Acad. Sci.* **253**, 27–50.
4. E. D. Salmon (1989) In *Mitosis: Molecules and Mechanisms* (J. Hyams and B. R. Brinkley, eds.), 119–181.
5. J. R. McIntosh, E. S. Ogata, and S. C. Landis (1973) *J. Cell Biol.* **56**, 304–323.
6. M.-F. Carrier (1989) *Int. Rev. Cytol.* **115**, 139–170.
7. H. P. Erickson and E. T. O'Brien (1992) *Annu. Rev. Biomol. Struct.* **21**, 145–166.
8. R. L. Margolis and L. Wilson (1978) *Cell* **13**, 1–8.
9. V. I. Rodionov and G. G. Borisy (1997) *Science* **275**, 215–217.
10. H. Hotani and T. Horio (1988) *Cell Motil. Cytosk.* **10**, 229–236.
11. T. J. Mitchison and M. Kirschner (1984) *Nature* **312**, 237–242.
12. C. E. Walczak, T. J. Mitchison, and A. Desai (1996) *Cell* **84**, 37–47.
13. T. J. Mitchison and E. T. Salmon (1992) *J. Cell Biol.* **119**, 569–582.
14. D. Panda, J. E. Daijo, M. A. Jordan, and L. Wilson (1995) *Biochemistry* **34**, 9921–9929.
15. R. Dhamodharan et al. (1995) *Mol. Biol. Cell* **6**, 1215–1229.
16. M. A. Jordan, R. J. Toso, D. Thrower, and L. Wilson (1993) *Proc. Natl. Acad. Sci USA* **90**, 9552–9556.
17. E. Nogales, S. G. Wolf, and K. H. Downing (1998) *Nature* **391**, 199–203.
18. D. Panda, M. A. Jordan, K. C. Chu, and L. Wilson (1996) *J. Biol. Chem.* **271**, 29807–29812.

### Suggestions for Further Reading

19. H. V. Goodson, C. Valetti, and T. E. Kreis (1977) Motors and Membrane Traffic. *Curr. Opin. Cell Biol.* **9**, 18–28.
20. J. M. Scholey (1996) Kinesin II, a membrane traffic motor in axons, axonemes, and spindles. *J. Cell Biol.* **133**, 1–4.
21. J. C. Waters and E. D. Salmon (1977) Pathways of spindle assembly. *Curr. Opin. Cell Biol.* **9**, 37–43.
22. L. Wilson and M. A. Jordan (1994) "Pharmacological Probes of Microtubule Function". In *Microtubules*, (J. Hyams and C. Lloyd, eds), Wiley-Liss, NY, pp. 59–84.
23. E. Nogales, M. Whittaker, R. A. Milligan, and K. Downing (1999) High-resolution model of the microtubule. *Cell* **96**, 79–88.

### Minicells

Minicells are small spherical cells that are formed in certain mutants of rod-shaped **bacteria**. In these mutant strains, the **cell division** septum is aberrantly placed adjacent to the cell pole instead of at its normal midcell site. The polar septation events give rise to spherical cells (minicells) that lack the bacterial **chromosome** but otherwise appear normal. The minicells are capable of **protein synthesis** and normal metabolic functions but are incapable of undergoing further rounds of cell division. The first description of a mutation that led to the minicell **phenotype** in *Escherichia coli* was made by Adler et al. (1). Minicell formation has since been described in a number of other **Gram-positive** and **Gram-negative** species (2-5).

The formation of polar septa, and the concomitant production of minicells, results from the use of potential division sites that are located at the cell poles. The polar division sites are probably derived from division sites that had been located at midcell during a preceding division event. In this view, components that had been located at the midcell division site are retained at the new poles of the daughter cells after septation is completed, and these sites are still capable of supporting additional septation events (6). It has also been suggested that the polar sites are induced by changes in the organization or intracellular localization of the bacterial chromosomes, rather than being residues of previous division events (7, 8).

The potential division sites at the cell poles are not used to support septum formation because they are suppressed by the gene products of the *min* genetic locus, as shown by studies of *E. coli* and *Bacillus subtilis*. The minicell **phenotype** occurs when the suppression of polar division events is lost as a result of mutations or experimental manipulations that interfere with function or expression of the *min* genes (9). The minicell phenotype can also be induced by procedures that increase the expression of the cell division gene *ftsZ*, thereby increasing the frequency of septum formation at all sites, both the aberrant sites at the poles and the normal division site at midcell (10). The minicell phenotype is also associated with certain mutations or physiological manipulations that perturb DNA structure or replication pattern (8, 11, 12).

Minicelling populations include two types of cell, spherical minicells and rod-shaped cells that are longer than the cells of wild type populations. The lengths of the rod-shaped cells range approximately from one to six times normal cell length. It has been suggested that this occurs because in each division cycle the cell has only enough “division potential” to support a single septation event (one “quantum” of division potential) (6). In this view, if the polar sites are available because of a *min* mutation, septation can occur either at a polar site or at the normal midcell site; septation at a polar site leads to one minicell and one cell that is almost two cell lengths long. Because a similarly random choice between polar and central division occurs in each generation, a portion of the cells end up being several cell lengths long before a central division occurs. The “division potential” that is postulated in this quantal model of cell division could correspond to the cellular content of FtsZ, because FtsZ appears to be the rate-limiting protein in the division process (see **Cell division**). The quantal hypothesis has recently been questioned on the basis of examination of the division pedigrees of individual cells (8) but still remains an attractive explanation of the known facts.

Minicell mutants have been used primarily for two purposes: (1) to better understand the mechanism of cell division, with emphasis on the mechanism used by the cell to identify the correct location for the division site; and (2) as factories for the synthesis and characterization of proteins.

## 1. *E. coli min* Genes

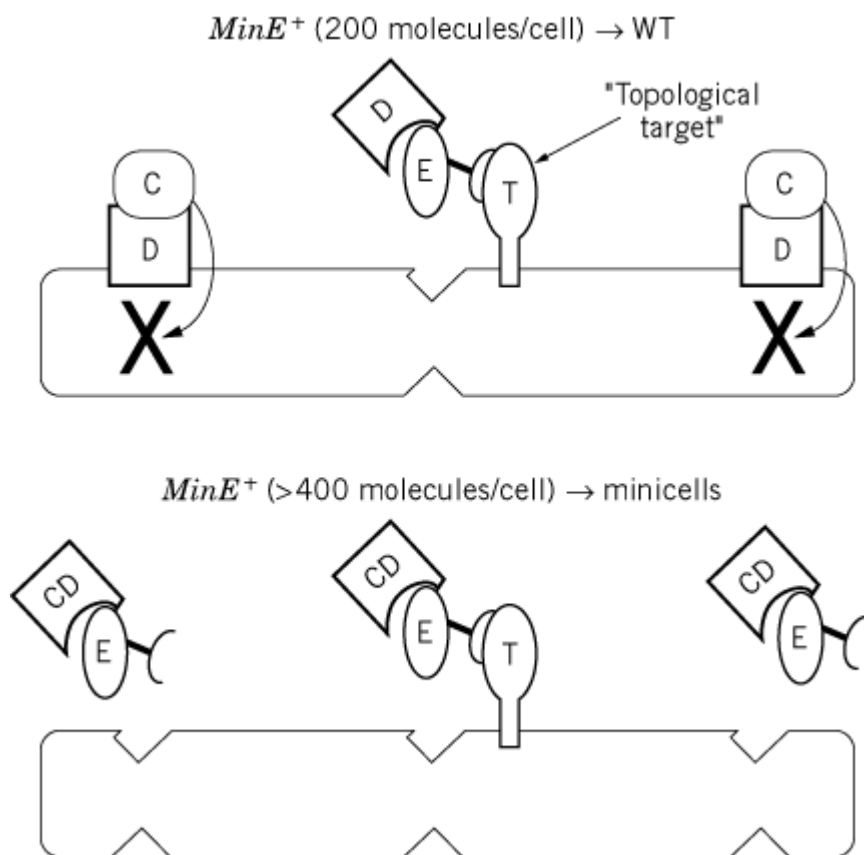
Most of what is known about the normal regulation of polar division sites has come from studies of *min* mutants of *E. coli* (13). The *min* genetic locus defines an **operon** that includes three genes, *minC*, *minD*, and *minE* (9). This gene cluster is sometimes called the *minB* locus. The term *minB* is a historical anachronism, resulting from the original misconception that mutations in two genetic loci,



called *minA* and *minB*, were required to produce the minicell phenotype. It was later discovered that the original identification of *minA* was an error and that the minicell phenotype resulted from a single mutation in *minB* (ie, in the *minCDE* gene cluster) (14, 15). This mutation later was shown to reside in *minD* (16).

The *minCDE* gene products act in concert to prevent septation at the cell poles. This is accomplished by the coordinated action of a nonspecific division inhibitor, whose activity requires the expression of *minC* and *minD*, and a topological specificity factor (MinE) that restricts the action of the division inhibitor to the polar division sites (Fig. 1). As a result, minicell formation is prevented, whereas septation at the midcell site occurs normally. Consistent with this model, loss of either MinC or MinD leads to septation at the polar sites, with production of the minicell phenotype.

**Figure 1.** Current model for function of *min* gene products. (Adapted from Ref. 13.)



MinC and MinD act in concert to produce a nonspecific inhibitor of septation, as shown by the observation that expression of *minC* and *minD* in the absence of *minE* leads to the formation of long nonseptate filaments. In the MinCD filaments, septation is prevented at the cell poles, as well as at the normal midcell site. Although at normal levels of expression MinC and MinD are both required for the division block, MinC is believed to be the actual division inhibitor protein. This conclusion is based on the observation that high levels of expression of MinC lead to the formation of nonseptate filaments even in the absence of MinD (17). MinD is believed to function as an activator of the latent or low-level division inhibition activity of MinC. Studies using the yeast two-hybrid system have indicated that MinC and MinD can interact (18).

The MinCD division inhibitor blocks septation by interfering with the earliest detectable step in cell division, the formation of the cytokinetic FtsZ ring that appears at the division site prior to septal

ingrowth (19) (see **Cell division**). Consistent with the idea that FtsZ is the MinCD target, the ability of MinCD to block division is suppressed by the overexpression of FtsZ (20, 21).

MinC also plays an essential role in a second division inhibition system, the *minC-dicB* system. In this system, DicB plays the role that MinD plays in the *minC-minD* system. The *dicB* gene is believed to be part of a cryptic **prophage** in the *E. coli* chromosome (22). Its expression is normally tightly repressed, and there are no known physiological mechanisms to relieve this repression. However, when *dicB* is induced by placing it under control of a regulatable promoter, such as  $P_{lac}$  (see **Lac Operon**), division is blocked completely, leading to formation of long nonseptate filaments (22). The DicB-mediated division block requires the presence of MinC but is independent of MinD (20). Taken together with the observation that high concentrations of MinC can block division in the absence of MinD or DicB (17), this has led to the conclusion that DicB and MinD can act as independent cofactors or activators of the MinC-mediated division inhibition reaction. It is not known how DicB or MinD function to activate the division inhibition systems.

DicB and MinD have no apparent sequence similarity, and it is likely that the two proteins use different mechanisms to activate the division inhibition reaction. The MinD protein has **ATPase** activity, consistent with the presence of a predicted **nucleotide-binding motif** in the protein sequence (23). Mutational inactivation of this site results in loss of the ability of MinD to activate the MinC-mediated division inhibition reaction. This suggests that nucleotide binding and/or hydrolysis play a role in MinD function. In contrast, the *dicB* sequence shows no similar nucleotide-binding site. A number of *minC* mutations have been identified that inactivate the division inhibition activity of the MinCD or MinC-DicB systems. Some of the mutants showed some degree of preferential resistance to either DicB or MinD (24), suggesting that the two proteins may act on different regions of MinC or on different aspects of the division inhibition system. The *minC-minD* and the *minC-dicB* division inhibition reactions also can be distinguished by differences in their response to the MinE protein (discussed below).

When *minE* is expressed coordinately with *minCD*, the MinCD-induced filamentation is prevented and the wild type division pattern is restored. Thus, MinE prevents the MinCD division inhibitor from working at midcell, without interfering with its ability to prevent minicell formation by blocking septation at the cell poles. This implies that MinE, either directly or indirectly, can distinguish between the new potential division site at midcell and the potential division sites at the cell poles.

The ability of MinE to suppress the activity of the MinC-mediated division inhibition reaction depends on the presence of MinD. Thus, although MinE can prevent MinCD-mediated filamentation, it has no effect on the division block that results from the coexpression of MinC and DicB, nor on the division block that results from high-level expression of MinC alone. Therefore, it is thought that MinD performs two roles. First, it activates the MinC division inhibitor. Second, it mediates the interaction between MinE and the division inhibition system.

To explain the site-specificity of MinE action, a model has been proposed in which MinE has a high affinity for a topological target protein that acts as a MinE docking protein or receptor (Fig. 1) (25). The target protein is presumed to be located at midcell, where it sequesters MinE. As a result, MinE only counteracts the MinCD division inhibitor at the midcell site. This model is consistent with the low abundance of MinE (approximately 200 monomers per cell) and with the observation that a twofold or greater increase in the cellular concentration of MinE leads to minicell formation. According to this model, the increased level of MinE saturates the topological target, leaving the excess MinE free to counteract the MinCD division inhibitor at the polar sites, leading to the observed minicell phenotype. As predicted from this model, the effect of the excess MinE can be counteracted by increasing the level of MinCD, which restores the wild-type division pattern, presumably due to titration of the excess MinE. Thus it is the relative concentrations of MinC, MinD, and MinE that are important, rather than the absolute concentration of the individual

components.

In a variation of this model, the topological target is located at the cell poles instead of midcell (26). This would lead to the sequestration of MinE at the polar sites, where it would play a positive role, activating the division inhibitor at the cell poles. The evidence available at this time favors the first model in which the MinE target is at midcell; but until the locations of the proteins are directly determined, both models are tenable.

Two functions are defined for MinE. First, it is an antagonist of the MinCD division inhibitor, as shown by its ability to prevent the formation of nonseptate MinCD filaments. Second, it is a topological specificity factor, limiting the action of the division inhibitor to the potential division sites at the cell poles. Genetic deletion analysis has shown that the anti-MinCD function resides entirely within the 22-residue *N*-terminal domain of MinE, MinE<sup>1-22</sup>. Thus, expression of MinE<sup>1-22</sup> prevents MinCD-induced filamentation (25). However, the MinE<sup>1-22</sup> domain lacks the topological specificity function of MinE, because expression of MinE<sup>1-22</sup> in the presence of MinCD leads to a minicell phenotype instead of a wild-type division pattern. In contrast, deletion analysis has shown that the topological specificity function of MinE, as defined by its ability to prevent minicelling, requires a domain that extends close to the carboxy-terminus of the protein (25, 26).

MinE is probably an oligomeric protein, as shown by gel filtration studies using MalE-MinE chimeric proteins. These studies, and studies using the yeast two-hybrid system, have indicated that the domain required for oligomerization lies within the carboxy-terminal half of the protein (26).

## 2. *min* Genes in *B. subtilis*

Mutations in two *B. subtilis* genetic loci, *divIVA* and *divIVB*, have been shown to lead to minicell phenotypes (2). The *divIVB* locus includes homologues of the *E. coli* *minC* and *minD* genes (27-29). The homology of the predicted amino acid sequences is strong for MinD (43.7% identity and 83.3% similarity between the *E. coli* and *B. subtilis* proteins) and less so for MinC (17.7% identity and 26.1% similarity). Similar to the results in the *E. coli* system, insertional disruption of the *B. subtilis* *minD* or *minC* gene leads to a strong minicell phenotype. However, in contrast to *E. coli*, where *minE* is located immediately downstream of *minD*, a *minE* homologue has not been identified in *B. subtilis*.

It is not known whether the *B. subtilis* MinC and MinD proteins function like their *E. coli* homologues. Similarity in function is suggested by the observation that minicelling results from mutation of *minC* or *minD* in both species, and that expression of the *B. subtilis* proteins in wild-type *E. coli* leads to minicelling (27). This suggests that *B. subtilis* MinC and/or MinD can compete with their *E. coli* counterparts and implies some similarity in function. It is possible that the Min proteins function similarly in the two species, but that the topological specificity function that is carried out by MinE in *E. coli* is incorporated into the MinC or MinD protein of *B. subtilis*, thereby eliminating the need for a separate *minE* gene. Alternatively, a *B. subtilis* *minE* homologue could be located within the adjacent gene cluster that was expressed together with *minC* and *minD* when minicelling was induced in *B. subtilis*, or a *B. subtilis* *minE* counterpart could be located elsewhere in the chromosome. This would be consistent, for example, with the observation that overexpression in *B. subtilis* of the *B. subtilis* *minC* and *minD* genes (together with adjacent genes) leads to a minicell phenotype (27). A similar result has been observed in *E. coli* when *minCDE* is overexpressed (30).

Mutations in a second *B. subtilis* locus, *divIVA*, also lead to a minicelling phenotype (2, 31). The *divIVA* gene codes for a 264-residue protein (31). The phenotype resulting from inactivation of *divIVA* differs from that of *divIVB* (*minC* or *minD*) minicelling mutants, because the *divIVA* mutants show a lower frequency of minicelling and also show extensive filamentation.

Overexpression of *divIVA* leads to a block in cell division, as shown by the formation of long

nonseptate filaments (31). The DivIVA-induced division block appears to require a functional MinC protein, since filamentation was not observed when *divIVA* was overexpressed in a MinC mutant strain. Further evidence for a functional interaction between DivIVA and MinCD comes from the observation that the phenotype of a *divIVA-minC* double mutant resembles that of a *minC* single mutant rather than that of the *divIVA* mutant alone (31). Taken together with the observation that DivIVA-mediated division inhibition requires MinC function, this implies that MinCD acts downstream of DivIVA and is consistent with the view that the *divIVA* gene product may act by modulating the activity of MinCD. This could occur, for example, if DivIVA acted as an activator of a latent division inhibitor activity of MinC, thereby functionally resembling the *E. coli* MinD and DicB proteins (discussed above).

### 3. FtsZ and Minicell Formation

Moderate overexpression of the essential cell division protein FtsZ leads to minicell formation (10). However, in contrast to the phenotype associated with mutation of the *min* genes, where the population consists of a mixture of spherical minicells and of rod-shaped cells that are longer than wild-type cell length, overexpression of FtsZ leads to a mixture of spherical minicells and of rod-shaped cells that are shorter than the cell length of wild-type cells. This occurs because the increased concentration of FtsZ leads to an increase in the frequency of septation events, with the extra septations occurring at the normally unused polar sites. As a result, minicell formation leads to the generation of one minicell and one cell that is shorter than normal. Presumably, the excess FtsZ overcomes the normal MinCD-mediated suppression of the potential division sites at the cell poles. This is consistent with the observation that overexpression of *ftsZ* suppresses the filamentation that otherwise occurs when the cellular concentration of MinC and MinD is increased above normal levels.

### 4. Effect of Minicell Mutations on Chromosomal Organization

Several observations suggest a relationship between minicell formation and chromosome location or the state of chromosomal organization. Åkerlund et al. (8) have shown that minicell formation occurs when aberrant chromosome replication takes place in cells in which *oriC* is nonfunctional and a plasmid R1 replication origin is present elsewhere in the chromosome. Minicelling also has been reported to occur when the DNA-binding Fis protein is nonfunctional, and in *gyrA* and *gyrB* mutants (11, 12), although in the *gyr* mutants there is apparently a continuum of cell lengths ranging from “classic” spherical minicells to cells of approximately normal or longer than normal cell lengths. In all of these cases, the nucleoid-free regions at the cell poles appear longer than normal. In addition, the nucleoid-free region at the poles has been reported to be longer than normal in *min* mutants (7), and it has been reported that minicell mutants produce a low but significant number of anucleate cells of normal cell length (32). An altered degree of supercoiling of some plasmids has also been reported to occur in minicell mutants (7). On the basis of these observations, it has been suggested that the *min* gene products act on chromosome organization or partition, leading in *min*<sup>-</sup> mutants to a nucleoid-free zone at the cell poles. In this view, polar septation occurs as the default mode wherever a sufficiently large nucleoid-free zone exists; and the Min proteins do not act on old division sites at the poles, but directly on the nucleoid or on the chromosome replication or partition system (7, 8).

### 5. Synthesis of Plasmid-Encoded Proteins in Minicells

Extrachromosomal genetic elements such as plasmids segregate into minicells, although the bacterial chromosome does not. Therefore, when a minicelling mutant strain is transformed with a native or recombinant plasmid, the minicells contain plasmid DNA and continue to synthesize plasmid-encoded message after the short-lived chromosomally encoded messenger RNA has decayed. Therefore, when the minicell population is labeled with a radioactive amino acid, essentially all of the labeled proteins that are synthesized are encoded by the plasmid DNA. This has been exploited as a generally applicable method to identify the gene products of desired genes by using plasmids

that include the genes of interest ([33](#)).

## Bibliography

1. H. I. Adler, W. D. Fisher, A. Cohen, and A. A. Hardigree (1967) *Proc. Natl. Acad. Sci. USA* **57**, 321–326.
2. J. N. Reeve, N. H. Mendelson, L. I. Coyne, L. L. Hallock, and R. M. Cole (1973) *J. Bacteriol.* **114**, 860–873.
3. K. Chung and L. P. Lin (1974) *Can. J. Microbiol.* **20**, 1621–1623.
4. W. G. Tankersley, J. Woodward, and A. Brown (1974) *Proc. Soc. Exp. Biol. Med.* **145**, 802–805.
5. B. Sedgwick, J. K. Setlow, M. E. Boling, and D. P. Allison (1975) *J. Bacteriol.* **123**, 1208–1217.
6. R. M. Teather, J. F. Collins, and W. D. Donachie (1974) *J. Bacteriol.* **118**, 407–413.
7. E. Mulder, M. El‘Bouhali, E. Pas, and C. L. Woldringh (1990) *Mol. Gen. Genet.* **221**, 87–93.
8. T. Åkerlund, R. Bernander, and K. Nordström (1992) *Mol. Microbiol.* **6**, 2073–2083.
9. P. A. J. de Boer, R. E. Crossley, and L. I. Rothfield (1989) *Cell* **56**, 641–649.
10. J. E. Ward and J. Lutkenhaus (1985) *Cell*, **42**, 941–949.
11. E. Orr, N. F. Fairweather, I. B. Holland, and R. H. Pritchard (1979) *Mol. Gen. Genet.* **177**, 103–112.
12. K. Hussain, E. J. Elliott, and G. P. C. Salmond (1987) *Mol. Microbiol.* **1**, 259–273.
13. L. I. Rothfield and C.-R. Zhao (1996) *Cell* **84**, p. 183–186.
14. T. H. Schaumberg and P. L. Kuempel (1983) *J. Bacteriol.* **153**, 1063–1065.
15. E. Davie, K. Sydnor, and L. I. Rothfield (1984) *J. Bacteriol.* **158**, 1202–1203.
16. C. Labie, F. Bouché, and J.-P. Bouché (1990) *J. Bacteriol.* **172**, 5852–5855.
17. P. A. J. de Boer, R. E. Crossley, and L. I. Rothfield (1992) *J. Bacteriol.* **174**, 63–70.
18. J. Huang and J. Lutkenhaus (1996) *J. Bacteriol.* **178**, 5080–5085.
19. Bi, E. and J. Lutkenhaus (1993) *J. Bacteriol.* **175**, 1118–1125.
20. P. A. J. de Boer, R. E. Crossley, and L. I. Rothfield (1990) *Proc. Natl. Acad. Sci. USA* **87**, 1129–1133.
21. E. Bi and J. Lutkenhaus (1990) *J. Bacteriol.* **172**, 5610–5616.
22. S. Béjar, F. Bouché, and J.-P. Bouché (1988) *Mol. Gen. Genet.* **212**, 11–19.
23. P. A. J. de Boer, R. E. Crossley, A. R. Hand, and L. I. Rothfield (1991) *EMBO J.* **10**, 4371–4380.
24. E. Mulder, C. L. Woldringh, F. Tétart, and J.-P. Bouché (1992) *J. Bacteriol.* **174**, 35–39.
25. C.-R. Zhao, P. de Boer, and L. Rothfield (1995) *Proc. Natl. Acad. Sci. USA* **92**, 4313–4317.
26. S. Pichoff, B. Vollrath, C. Touriol, and J.-P. Bouché (1995) *Mol. Microbiol.* **18**, 321–329.
27. A. W. Varley and G. C. Stewart (1992) *J. Bacteriol.* **174**, 6729–6742.
28. P. Levin, P. S. Margolis, P. Setlow, R. Losick, and D. Sun (1992) *J. Bacteriol.* **174**, 6717–6728.
29. S. Lee and C. W. Price (1993) *Mol. Microbiol.* **7**, 601–610.
30. P. A. J. de Boer, R. E. Crossley, and L. I. Rothfield (1988) *J. Bacteriol.* **170**, 2106–2112.
31. J.-H. Cha and G. Stewart (1997) *J. Bacteriol.* **179**, 1671–1683.
32. A. Jaffe, R. D'Ari, and S. Hiraga (1988) *J. Bacteriol.* **170**, 3094–3101.
33. M. Kuehn, F. Jacob-Dubuisson, K. Dodson Slonim, L. R. Striker, and S. Hultgren (1994) *Methods Enzymol.* **236**, 285–287.

## Minichromosome

A minichromosome represents a **DNA** molecule and associated [histones](#) that assemble as [chromatin](#) similar to that in the [chromosomes](#) of the host cell. Minichromosomes assemble chromatin most efficiently in mammalian cells if they contain **origins of replication** from viral DNA; In **yeast** cells, minichromosomes can be used that contain true chromosomal origins ([autonomously replicating sequences, \[ARS\]](#)). They can also be assembled, although less efficiently, on DNA **microinjected** or **transfected** into cell **nuclei**.

The small size of the [SV40](#) genome (5243 bp long) and the early availability of information concerning its DNA sequence and gene organization made it a convenient model for studying the structure and function of chromatin. Late in infection, the SV40 genome is organized into arrays of **nucleosomes** as a minichromosome that can be isolated from infected cells under appropriate conditions. One region of the minichromosome (ORI) contains several important recognition sites for **trans-acting** factors: the origin of replication, the binding sites for the viral regulatory protein ([T Antigen](#)), and the **promoters** driving [transcription](#) of early and late SV40 [messenger RNA](#). The chromosomal organization of this ORI region was recognized as important for [DNA replication](#) and transcription and was shown to differ from the rest of the minichromosome. The first experiments used **nucleases** to digest the minichromosome. In the nuclei of infected cells, the DNA sequence is preferentially cut in the ORI region with **DNase I** (see [Hypersensitive Site](#)). This suggested that DNA in the ORI region is more accessible to [DNA-binding proteins](#) than in the rest of the minichromosome. Chemical **carcinogens** that interact with the free DNA duplex by intercalation, such as [psoralen](#), bind preferentially to the ORI region in infected cells. Its rapid digestion with a variety of **endonucleases** in isolated minichromosomes is also consistent with preferential accessibility of the ORI region. Finally, [electron microscopy](#) reveals that 20 to 25% of the isolated minichromosomes contain a nucleosome-free region (or gap) of approximately 350 bp covering the ORI region. This gap represents the best early documentation of nucleosome positioning on regulatory DNA (1).

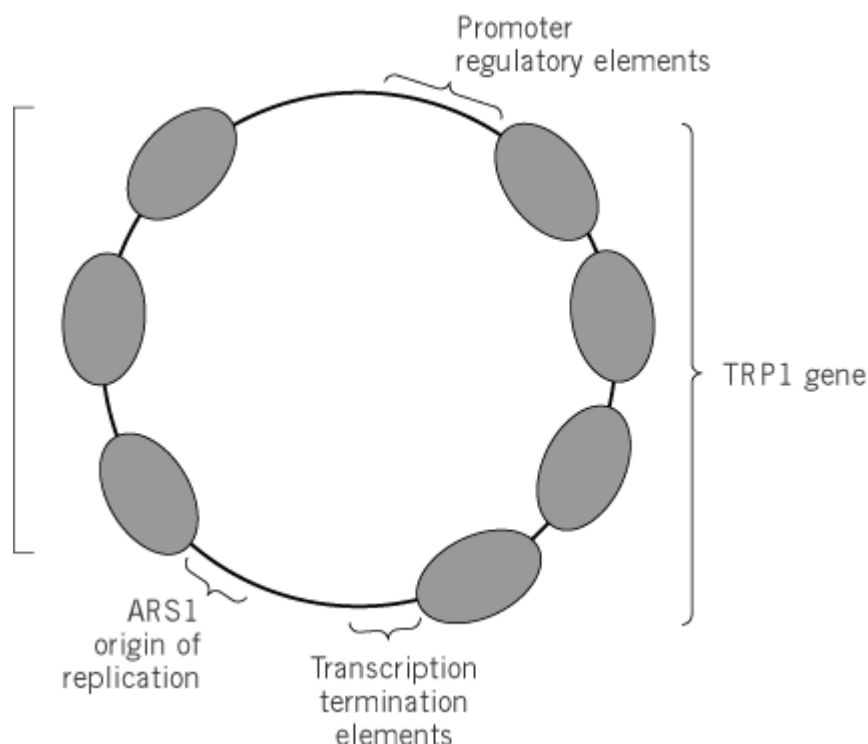
Several genes or **promoter** elements have been introduced into SV40 minichromosomes, and the influence of chromatin structure on **gene expression** was examined. Although generally repressive effects on transcription are observed, the structural basis of this repression has not yet been determined. The advantage of viral genomes for analyzing the interrelationship between chromatin structure and function is that they contain replicative origins. This means that the viral minichromosome can be studied as an **episome**, ie, without being integrated into the cellular chromosomes. Moreover, because chromatin assembly is coupled to DNA replication, the influence and access of *trans*-acting factors on nucleosomal organization are more like that within the true chromosomal context. The only reservation about studies with small viral genomes is that possible regulatory effects dependent on higher order chromatin structure are unlikely to be observed. Again, this is caused by topological constraints in folding small DNA circles in the chromatin fiber. Viral genomes other than SV40 have been very useful for the detailed analysis of chromosomal influences on transcription. These include the bovine papilloma virus (BPV)-based episomes that have been used to investigate the molecular mechanisms by which [glucocorticoid](#) receptor activates gene expression. Here, the interplay between specific chromatin structures and [transcription factors](#) has been rigorously documented (2).

Viral episomes that are DNA molecules containing a viral origin of replication have been used to establish the organization of *trans*-acting factors and nucleosomes in specific DNA sequences. However, most studies have investigated gene regulation by transiently transfecting **cloned** DNA without a viral origin of replication into eukaryotic cells. Under these conditions, such studies have rarely investigated repressive or stimulatory effects that might be attributed to chromatin. Importantly, the observed regulation due to enhancers, promoters, or other elements often does not

reflect the range of response observed in the natural chromosomal context. Extensive studies by Howard and colleagues demonstrated that the efficiency of chromatin assembly depends on the transfective conditions. The amount of DNA transfected into cells, the method of compacting DNA prior to transfection, and the efficiency with which particular DNA sequences partition to the nucleus greatly affect chromatin assembly. Over 80 to 90% of nuclear plasmid material might not be assembled into chromatin, yet some templates and protocols generate nucleosomal arrays (3). Consequently, it is not surprising that the regulation of transiently transfected DNA does not always completely reflect that found in the natural chromosomal context.

Yeast is a particularly attractive system for examining the influence of chromatin structure on the function of DNA. Experimental work with yeast has many advantages, especially for molecular biologists and geneticists. Among these are the existence of small (1500 bp) extrachromosomal minichromosomes that replicate autonomously. Particular experimental attention has been given to a plasmid present at about 100 copies/cell known as TRP1ARS1. This plasmid is 1453 bp long and consists of (1) one gene coding for N-(5'-phosphoribosyl)-anthranilate isomerase (TRP1), (2) a sequence containing a replicative origin (ARS1), and (3) a segment of unknown function (UNF) but now known to be part of the GAL3 promoter (Fig. 1). Simpson, Thoma, and colleagues have determined the chromatin structure of TRP1ARS1 in some detail (4). Nucleosomal position has been determined by the nuclease accessibility of isolated minichromosomes. Four different regions of chromatin structure have been defined. The UNF region contains three nucleosomes that have very strong DNA sequence-directed positions and are not changed by inserting DNA fragments of various lengths into the plasmid. In contrast, four loosely positioned nucleosomes are found on the TRP1 gene. These nucleosomes can easily rearrange following insertion of additional DNA. Two nucleosome-free regions also present (like the gap in the SV40 minichromosome) are hypersensitive to nuclease digestion. These sequences include the TRP1 promoter and the ARS origin of replication. Although lacking **centromeres** and **telomeres**, the TRP1ARS1 minichromosome therefore contains many elements associated with normal cellular chromosomes.

**Figure 1.** The structure of the TRP1 ARS1 minichromosome. Nucleosome (open ellipsoids) and key DNA sequences are shown. The region originally called UNF, for unknown function, is now known to be part of the GAL3 promoter.



Among the DNA sequences inserted into the TRP1ARS1 plasmid was a sea-urchin 5S RNA gene, on which base-pair resolution of the rotational positioning of DNA on the histone core was originally defined. A nucleosome formed including this sequence, with exactly the same position in the yeast minichromosome, as observed *in vitro*. This important observation shows that yeast histones *in vivo* recognize the same DNA sequence-directed nucleosome positioning elements as chicken histones *in vitro*. Similar experiments with yeast genomic sequences inserted into the TRP1ARS1 plasmid consistently reveal that the same nucleosome positioning occurs on the episome as in the chromosome. In certain instances, unstable nucleosomes, like those on the TRP1 gene, may have their positions influenced by the organization of other regions of the episome into chromatin. More recent studies have taken these observations further to dissect the contributions of specific trans-acting factors to nucleosome positioning.

### Bibliography

1. A. J. Varshavsky, O. H. Sundin, and M. J. Bohn (1978) *Nucleic Acids Res.* **5**, 3469–3477.
2. T. K. Archer, P. Lefebvre, R. G. Wolford, and G. L. Hager (1992) *Science* **255**, 1573–1576.
3. R. Reeves, C. M. Gorman, and B. Howard (1985) *Nucleic Acids Res.* **13**, 3599–3615.
4. R. T. Simpson (1991) *Prog. Nucleic Acids Res. Mol. Biol.* **40**, 143–184.

### Suggestion for Further Reading

5. A. P. Wolffe (1998) *Chromatin: Structure and Function*, 3rd ed., Academic Press, London.

## Minigene

A minigene is a shortened version of a cloned **gene** that is produced by [recombinant DNA](#) techniques. Most minigenes are made from **eukaryotic** (especially mammalian) genes because of their large size. In most cases, a sequence from the 5'-end of the gene is joined to a sequence near its 3'-end to form the minigene. The smaller size of a minigene facilitates studies of its regulatory sequences (see **Gene regulation**). It is necessary to show that the behavior of the minigene resembles that of the intact gene because internal sequences affect gene regulation ([1](#)). In some cases, only a 5'-terminal sequence is used. Such a minigene has been used to obtain high-level expression of proteins in salivary glands ([2](#)). A minigene with an intact coding region has been used to introduce [mutations](#) at specific sites in the mouse [genome](#) ([3](#)).

### Bibliography

1. L. H. Reid, R. O. Gregg, O. Smithies, and B. A. Koller (1990) *Proc. Natl. Acad. Sci. USA* **87**, 4299–4303.
2. J. Laursen and J. P. Hjorth (1997) *Gene* **198**, 367–372.
3. D. W. Melton, A. M. Ketchen, and J. Selfidge (1997) *Nucleic Acids Res.* **25**, 3937–3943.

## Minisatellite DNA



*Minisatellites* are dispersed arrays of short, tandem, direct repeats of highly variable DNA sequences found in the **genomes** of most higher **eukaryotes**. They contain core sequences resembling the [recombination](#) hotspot **Chi sequence** 5'-GCTGGTGG-3' of *Escherichia coli*, which is a binding site for proteins involved in genetic recombination. Based on this, a general function of minisatellites may be to provide binding sites for recombination proteins in eukaryotes.

Their name arose for historical reasons. Under certain **density centrifugation** conditions, their density is similar to that of the satellite peak of genomic DNA.

### 1. Use of Minisatellites for Establishing Physical Genetic Maps

The discovery of minisatellites ([1](#)) has resolved the difficulties raised by the relative lack of information provided by *restriction fragment length polymorphisms* (**RFLP**). In the vertebrate genome, the short (10- to 50-bp) tandem, direct, repeat motifs of minisatellites contain variants of a common core sequence. Their lengths are highly variable, which is expressed by their other name, VNTR, for *variable number of tandem repeats*. Interest in them comes from the fact that the number of repeats of a given minisatellite at a certain position in the genome varies from one individual to another. Thus a probe for a particular minisatellite is highly polymorphic and reveals a locus with a great number of alleles that have a high probability of being genetically informative in the great majority of cases. A minisatellite general probe of a tandem repeat of the core sequence can simultaneously detect many variable loci because of their sequence homologies and can provide an individual specific DNA fingerprint of general use in human genome analysis.

But, the real solution has been provided by the discovery of **microsatellites**, named by analogy to the minisatellites.

### Bibliography

1. A. J. Jeffreys, V. Wilson, and S. L. Thein (1985) *Nature*, **314**, 67–73.

## Mismatch Repair

Mismatches in **DNA** are noncomplementary base pairs in a double helix. There are eight such base pairs and, if the directionality of the duplex is taken into account, the number of mismatches is 12. They are corrected by the phenomenon known as mismatch repair. In addition to simple mismatches, other abnormal DNA structures that do not involve damaged bases are also substrates for mismatch repair. These include the insertion/deletion of single nucleotides, loops consisting of unpaired nucleotides in one strand, and bubbles, which consist of more than one mismatched base in a row.

Mismatches arise from three sources: [DNA replication](#) errors, [recombination](#), and base deamination.

1. *Replication errors*. **DNA polymerases** are relatively accurate enzymes that ensure replication fidelity at several stages, including base selection, insertion, and proofreading. Despite all these safeguards, DNA polymerases make replication errors at a frequency of  $10^{-3}$  to  $10^{-5}$ , depending on the polymerase.

2. *Recombination*. Mismatches also arise during recombination between **alleles** of the same gene or between closely related sequences. If these mismatches occur in meiotic chromosomes and remain uncorrected, they can give rise to anomalous segregation of the alleles, or they may lead to **gene conversion** upon correction.

3. *Base deamination*. Cytosines are often methylated at the C5 position, and these 5-methylcytosines deaminate about  $10^3$ -fold faster than does cytosine (1). The extent of methylation varies across species, ranging from a virtually undetectable level in yeast and *Drosophila*, to 1% in *Escherichia coli*. In humans, 20% of cytosines are methylated; consequently, G-T mismatches arising from deamination of cytosine are a common occurrence in humans.

Mismatches are eliminated from DNA by three general mechanisms: general mismatch repair, very short patch mismatch repair, and mismatch repair by base excision. The general mismatch repair system is the most important, and is the system most widely distributed in the biological world. It has been found in all bacteria tested, as well as in model eukaryotic organisms, such as yeast, fruit flies, and humans.

## 1. General Mismatch Repair Systems

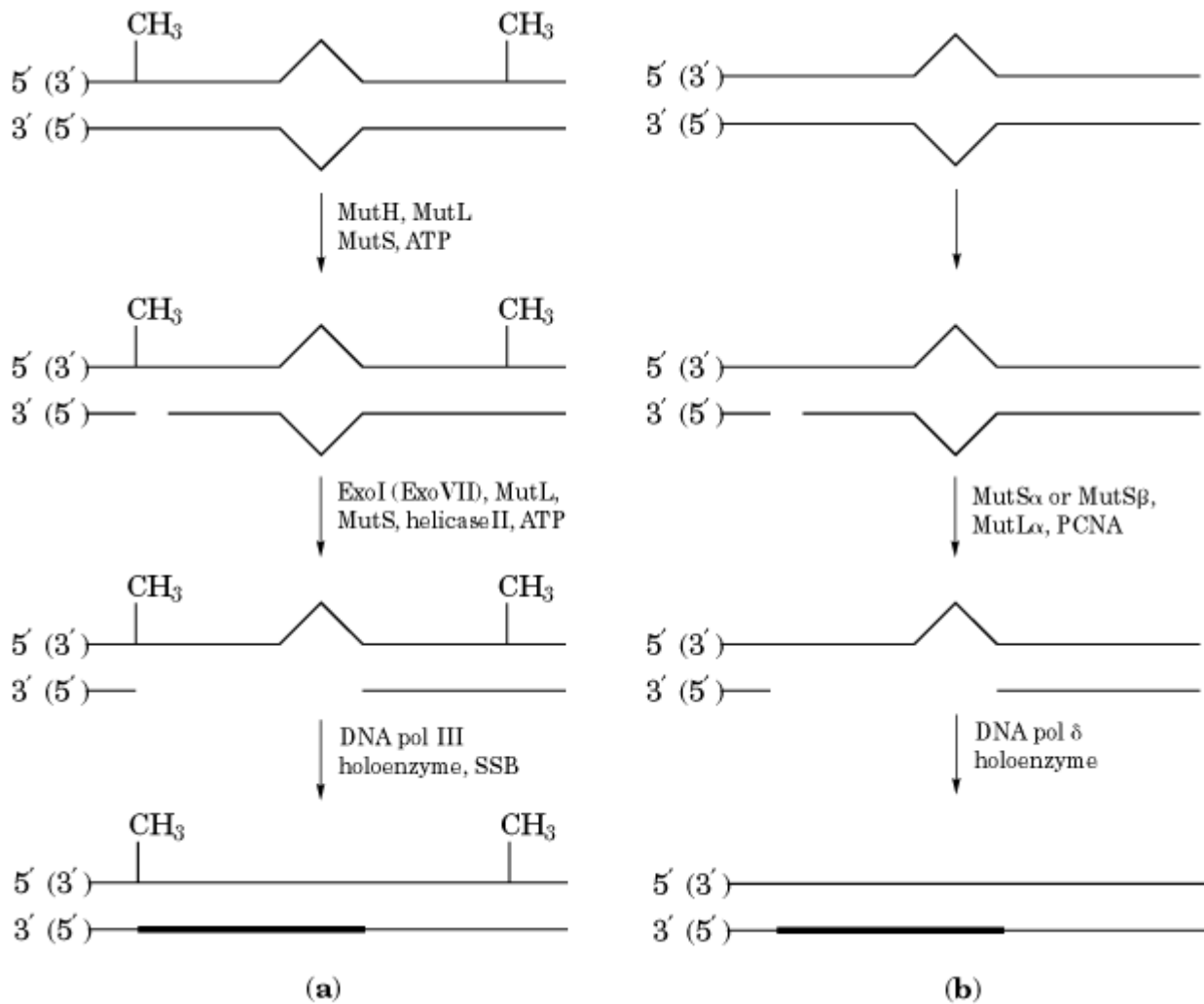
This repair system corrects all mismatches and small insertion/deletion loops in an essentially sequence-independent manner (2). The general reaction scheme is as follows: A nick is made 5' or 3' to the mismatch in the newly synthesized strand, which is by definition the “wrong” strand. A 5' to 3' or a 3' to 5' exonuclease digests the nicked strand, starting at the nick and terminating at some distance past the mismatch. Concomitant with exonucleolytic degradation, repair synthesis takes place, filling the gap generated by the exonuclease with correct nucleotides. Ligation of the repair patch completes the reaction (2). The reaction mechanism of general mismatch repair has been elucidated in considerable detail in *E. coli* and to a lesser extent in humans.

### 1.1. General Mismatch Repair in *E. coli* (Fig. 1)

In *E. coli* the three key proteins for mismatch repair are known as MutH, MutS, and MutL (2). In addition, the Dam methyltransferase plays a crucial role in strand discrimination. The Dam protein is a sequence-specific adenine DNA methyl transferase; it recognizes unmethylated or hemimethylated GATC sequences and methylates the N6 of adenine (see [Methylation, DNA](#) and [Methyltransferase, DNA](#)). Under steady-state growth conditions, both strands of the parental DNA are methylated. However, during replication a hemimethylated GATC sequence forms transiently as a result of the incorporation of a canonical (unmodified) adenine at this site by replicase. This hemimethylated GATC serves as a signal for discriminating between the parental (correct) and daughter (potentially incorrect) strands. When the newly synthesized DNA contains a mismatch, the mismatched base pair is recognized by MutS, while the hemi methylated GATC (5' or 3' to the mismatch) is recognized by MutH. Then, MutL binds to both the MutS-mismatch and MutH-GATC complexes. This binding activates the [ATPase](#) activity of MutS and the endonuclease activity of MutH, which then cleaves the phosphodiester bond 5' to the G of the GATC sequence in the unmethylated strand. The nick introduced by MutH serves as an entry site for an exonuclease (exonuclease I for 3' to 5' and exonuclease VII or RecJ for 5' to 3' digestion). One of these exonucleases digests the DNA from the site of the nick, which could be up to 1000 nucleotides away from the mismatch, to a point past the mismatch by 20 to 30 nucleotides. Digestion of DNA between the nick site and the mismatch is concurrent with resynthesis by the [DNA replication](#) machinery comprised of **DNA polymerase III**, [DNA helicase II](#), SSB ([single-stranded DNA binding protein](#)), and [DNA Ligase](#). A mutation in any of the mismatch repair-specific genes (ie, *mutS*, *mutH*, *mutL*) or in the specific methylation gene *dam* increases both the spontaneous and damage-induced mutation frequency by  $10^3$ - to  $10^4$ -fold.

**Figure 1.** Mechanisms of general mismatch repair in (a) *E. coli* and (b) humans. A nick 3' or 5' to the mismatch in the

newly synthesized strand marks that strand for correction. In *E. coli*, the newly synthesized strand is transiently undermethylated at the Dam site (GATC) and hence is nicked by MutH endonuclease. In humans, the mechanism of preferential nicking of the newly synthesized strand is not known. MutS in *E. coli* and MutSa and b in humans recognize the mismatch, and MutL or MutLa coordinate the nicking and exonucleolytic degradation of the mismatched strand to a point past the mismatch. The resulting gap is filled in by replication polymerases.



## 1.2. General Mismatch Repair in Humans (Fig. 1)

The basic steps and key enzymes of mismatch repair in humans are similar to those of *E. coli* (3-5). The mismatch is recognized by a MutS homologue, and a MutL homologue coordinates nicking with mismatch recognition. In humans, however, both MutS and MutL are encoded by several related genes (6). There are at least three *mutS* genes and proteins in humans: hMSH2, hMSH3, and hMSH6. The proteins encoded by these genes form heterodimers (hMSH2)<sub>1</sub>(hMSH6)<sub>1</sub> and (hMSH2)<sub>1</sub>(hMSH3)<sub>1</sub> which are also referred to as hMutSa and hMutSb and which are the functionally relevant forms. The MutSa is involved in recognition/repair of simple mismatches and small, one- to three-nucleotide insertion/ deletion mismatches. It appears that MutSb is involved in recognition of larger mismatch loops. There are three MutL homologues, hMLH1, hPMS1, and hPMS2, which form heterodimers (hMLH1)<sub>1</sub>(hPMS2)<sub>1</sub>, called hMutLa, and presumably (MLH1)<sub>1</sub>(PMS1)<sub>1</sub> heterodimer, which may be involved in the repair of larger loops (heterologies). The signaling mechanism that discriminates between the right and wrong strands in humans is not known. In humans, there is no enzymatic activity that methylates adenine. In contrast, up to 20% of cytosine bases are methylated. There is, however, no evidence that cytosine methylation is used as a signal in discriminating between the nascent and parental strands. It is possible that during replication the

natural nicks in the lagging strand and the [replication fork](#), or occasional nicks in the leading strand, provide both the signal for discrimination and the entry point for the mismatch repair exonuclease for mismatch removal. The helicase(s) and exonuclease(s) necessary for carrying out mismatch excision are not known, although DNA polymerases  $\alpha$ ,  $\delta$ , and  $\epsilon$  have been implicated in the repair synthesis.

A defect in mismatch repair in humans has two consequences. First, hereditary nonpolyposis colon cancer (HNPCC) is associated with a defect in mismatch repair ([6](#), [7](#)). Apparently, mutations in *MSH2* and *MLH1*, and to a lesser degree in *PMS1* and *PMS2*, that are inherited in heterozygotes are asymptomatic. However, a second mutation in the wild-type allele of an individual (by a [somatic mutation](#)) renders the cell mismatch repair defective and leads to development of colorectal cancers. HNPCC family members also develop other cancers at high frequency, including uterine, stomach, ovarian, and squamous cell skin cancers. Cancers that are caused by mismatch-repair defects are characterized by **microsatellite instability (MIN<sup>+</sup> phenotype)**, which includes expansion or contraction of simple repeat sequences, such as (A)<sub>n</sub>, (CA)<sub>n</sub> and (GGC)<sub>n</sub> ([8](#)). However, not all MIN<sup>+</sup> tumors are defective in mismatch repair. Microsatellite instability in general is a symptom, but not the cause of mutations that eventually convert a normal cell to cancerous phenotype ([8](#)). A significant number of these tumors have a mutation, presumably caused by a lack of mismatch correction, in the [transforming growth factor](#) (TGF- $\beta$ ) receptor.

The other manifestation of mismatch repair defects is resistance to anticancer drugs that exert their effect by damaging DNA ([9-11](#)). DNA damage caused by various anticancer drugs, including cisplatin and [adriamycin](#), causes misincorporation of nucleotides during replication. The resulting lesion in DNA is a mismatch superimposed upon base damage, which is referred to as a “compound lesion.” The lesions with or without a mismatch bind to mismatch recognition proteins ([12](#)). The binding activates the mismatch repair system, which removes the mismatch. Upon resynthesis, however, the mismatch is regenerated by DNA polymerase at a certain frequency, depending on the lesion. Hence, a futile cycle of excision and resynthesis ensues, which eventually leads to DNA degradation and cell death. If cells are mismatch repair-defective, this suicidal cycle is prevented ([10](#)). Although mismatch repair-defective cell lines have been isolated from tumor cell lines exposed to nitrosoguanidine or cisplatin ([11](#), [13](#), [14](#)), it is unclear at present whether or not secondary mutations in tumors contribute to development of drug resistance in clinical setting. Similarly, it is unknown at present whether or not tumors arising from HNPCC disease are resistant to cisplatin or similar drugs.

## 2. Very Short Patch (VSP) Repair

This repair system has been defined only in *E. coli*. It is a repair system specific for G-T (or G-U) mismatches. In *E. coli* the sequence CC(A/T)GG is the target for deoxycytidine methylase (Dcm), which methylates the second C in this sequence. Methylated C deaminates at a 1000-fold faster rate than does cytosine. As a consequence, G-T mismatches are produced in the *E. coli* genome at a rate that may be unacceptably high. It appears that the VSP system has evolved to deal with this specific problem. An enzyme encoded by the *vsr* gene specifically nicks 5' to the T of the mismatch in the context of a Dcm recognition site ([15](#)). However, both the Dcm methylase and the Vsr endonuclease also recognize and methylate or cleave at related sequences, albeit at a reduced efficiency. Although the Vsr protein is the only protein absolutely required for this repair pathway system, the VSR pathway is greatly stimulated by the MutL and MutS proteins by an unknown mechanism. In any event, following nicking, excision of mismatch and repair synthesis are carried out by DNA polymerase I, which generates a four- to six-nucleotide-long repair patch.

Although humans do not have a VSR system, they do have a repair system that deals specifically with G-T mismatches. In humans, the cytosines in a CG context in the so-called [CpG Islands](#) are heavily methylated and hence are subject to deamination, giving rise to G-T mismatches. In humans, this mismatch is processed by a thymine glycosylase, followed by the other steps common to all

“damaged DNA” excision repair systems. As in the case of VSR, the mammalian system for G-T mismatches also has a sequence context specificity and is more effective on G-T mismatches in CG islands.

### 3. A-G Mismatch Repair System

One of the major lesions produced in DNA by oxidative stress is 8-oxoguanine (8-oxoG). This base preferentially pairs with adenine during replication, giving rise to an 8-oxoG-adenine mismatch. A glycosylase encoded by *mutY* in *E. coli* cleaves off the mismatched adenine and the phosphodiester bond 3' to the resulting apurinic site (16). Then, the abasic sugar is eliminated, and the one- to two-nucleotide gap is filled and ligated. The mutagenic 8-oxoG base is eventually removed from the DNA by 8-oxoG glycosylase (see [Base Excision Repair](#)). Although this repair system is most active on the A-8-oxoG mismatch, it also removes the adenine residues from rare A-G mismatches that are generated by replication or recombination.

### Bibliography

1. B. K. Duncan and J. H. Miller (1980) *Nature* **287**, 560–561.
2. R. S. Lahue, K. G. Au, and P. Modrich (1989) *Science* **245**, 160–164.
3. J. J. Holmes, S. Clark, and P. Modrich (1990) *Proc. Natl. Acad. Sci. USA* **87**, 5837–5841.
4. D. C. Thomas, J. D. Roberts, and T. A. Kunkel (1991) *J. Biol. Chem.* **266**, 3744–3751.
5. W. H. Fang and P. Modrich (1993) *J. Biol. Chem.* **268**, 11838–11844.
6. R. Fishel, M. K. Lescoe, M. R. S. Rao, N. G. Copeland, N. A. Jenkins, J. Garber, M. Kane, and R. Kolodner (1993) *Cell* **75**, 1027–1038.
7. R. Parsons, G. M. Li, M. J. Longley, W. H. Fang, N. Papadopoulos, J. Jen, A. de la Chapelle, W. Kinzler, B. Vogelstein, and P. Modrich (1993) *Cell* **75**, 1227–1236.
8. Y. Ionov, M. A. Peinado, S. Malkhosyan, D. Shibata, and M. Perucho (1993) *Nature* **363**, 558–561.
9. P. Karran and M. Marinus (1982) *Nature* **296**, 868–869.
10. V. S. Goldmacher, R. A. Cuzick, and W. G. Thilly (1986) *J. Biol. Chem.* **261**, 12462–12471.
11. A. Kat, W. G. Thilly, W. H. Fang, M. J. Longley, G. M. Li, and P. Modrich (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6424–6428.
12. D. R. Duckett, J. T. Drummond, A. I. H. Murchie, J. T. Reardon, A. Sancar, D. M. J. Lilley, and P. Modrich (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6443–6447.
13. J. Drummond, A. Anthony, R. Brown, and P. Modrich (1996) *J. Biol. Chem.* **271**, 19645–19649.
14. S. Aebi, B. Kurdi-Haidar, R. Gordon, B. Cenni, H. Zheng, D. Fink, R. D. Christen, C. R. Boland, M. Koi, R. Fishel, and S. B. Howell (1996) *Cancer Res.* **56**, 3087–3090.
15. F. Hennecke, H. Kolmar, K. Bründl, and H. J. Fritz (1991) *Nature* **353**, 776–778.
16. K. G. Au, S. Clark, J. H. Miller, and P. Modrich (1989) *Proc. Natl. Acad. Sci. USA* **86**, 8877–8881.

### Suggestions for Further Reading

17. R. Kolodner (1996) Biochemistry and genetics of eukaryotic mismatch repair. *Genes Dev.* **10**, 1433–1442.
18. E. C. Friedberg, G. C. Walker, and W. Siede (1995) "DNA repair and mutagenesis". ASM Press, Washington, D.C.
19. P. Modrich (1997) Strand-specific mismatch repair in mammalian cells. *J. Biol. Chem.* **272**, 24727–24730.
20. P. Modrich and R. Lahue (1996) Mismatch repair in replication fidelity, genetic recombination, and cancer biology. *Annu. Rev. Biochem.* **65**, 101–133.

21. M. Lieb and A. S. Bhagwat (1996) Very short patch repair: reducing the cost of cytosine methylation. *Mol. Microbiol.* **20**, 467–473.

## Missense Mutation

Missense mutations comprise a class of [mutations](#) in which one **nucleotide** within the coding segment of a gene is substituted for another, changing a **codon** for one [amino acid](#) into a codon for a different amino acid. The effect of a missense mutation depends on the nature of the amino acid change (conservative versus nonconservative; see [Conservative Substitutions](#)) and on where it occurs in the **protein** encoded by the gene. Missense mutations in various locations in a protein that give a variety of phenotypes are a powerful aid in analyzing the relationship between protein structure and function. Like [nonsense mutations](#), missense mutations are **suppressed** by mutant [transfer RNAs](#), but these inevitably have deleterious effects because the amino acid sequences of most of the cell's proteins are affected.

## Mitochondria

Early **microscopic** descriptions of cells identified the existence in both animal and in plant cells of **organelles** that could be stained with vital dyes such as Janus green. Cytologists in early decades of this century were often struck by the morphological diversity of these structures, and many different names (chondriome, chondriomites) were used, in addition to the now generally used mitochondria, to indicate the same basic structure, which was only clearly identified in the 1950s by [electron microscopy](#). It then became apparent that mitochondria are surrounded by a double [membrane](#). The inner mitochondrial membrane (IM) folds inwards, thereby forming the [cristae](#).

The morphological definition of mitochondria had been preceded, in the 1940s, by the isolation of particulate mitochondrial cell fractions and by the findings that respiratory [enzymes](#) were localized in these fractions and that isolated mitochondria could respire and perform oxidative phosphorylation. The three main respiratory electron-transport complexes (NADH dehydrogenase; bc<sub>1</sub> complex; cytochrome oxidase) are now known to be localized in the IM, together with [ATP synthase](#), the complex that is responsible for the phosphorylation of ADP coupled to the electron transport. The integrity of the IM is essential for the formation of the [proton gradient](#) that drives the phosphorylation (see [Proton Motive Force](#)).

New techniques of [fluorescence microscopy](#), based on the targeting to the mitochondria of fluorescent proteins, often coupled with high-speed imaging systems, have recently revealed that mitochondria form in the cytoplasm a largely interconnected tubular network that undergoes continuous rearrangements; electron microscopy shows only sections of this network.

In 1949, Boris Ephrussi ([1](#)) in Paris observed that a spontaneous [mutation](#) occurred with high frequency (around  $10^{-2}$ ) in populations of the yeast *Saccharomyces cerevisiae*, leading to characteristic small colonies. The reduced size of the colony was related to a defect in respiration, which allowed growth only on fermentable substrates such as glucose. These mutants were called

*petite* (small), to indicate the small size of its colonies, and were later referred to as *rho*<sup>-</sup> to indicate the absence of respiration. It was soon discovered that the mutation exhibited non-Mendelian inheritance and that it was transmitted to progeny even when fusion of **nuclei** was prevented. This discovery was the origin of yeast mitochondrial genetics, which disclosed several features of the yeast mitochondrial [genome](#) [see [Mitochondrial DNA \(MTDNA\)](#)] even before sequencing techniques were available for DNA. Since most of these mutants exhibited defective **absorption spectra** of [cytochromes](#), indicating a defect in mitochondrially localized enzymes (2), it was suggested that the [cytoplasmic inheritance](#) might be due to the existence of a genetic determinant localized in mitochondria.

The presence in mitochondria of a [mitochondrial DNA](#) was demonstrated in the 1960s, and this was rapidly followed by the discovery of a mitochondrially localized machinery for **protein biosynthesis**. This is used for the *in situ* [translation](#) of [messenger RNA](#) resulting from [transcription](#) of the mitochondrial genes, which mainly code for some of the subunits of the above-mentioned respiratory enzymes. The other subunits of the respiratory enzyme complexes are encoded by the nucleus, synthesized in the cytoplasm, and imported into the mitochondrion, to be assembled with the mitochondrially synthesized subunits.

The mitochondrial protein biosynthesis apparatus has several distinctive features that differentiate it from the cytoplasmic counterpart of the same cell. The most important one is probably the sensitivity of its [ribosomes](#) to bacterial-type antibiotics, such as [chloramphenicol](#) and **erythromycin**, which are not active on **eukaryotic** ribosomes. This and other similarities with prokaryotic cells have strengthened the hypothesis of an infective origin of mitochondria in eukaryotic cells; a primitive respiring bacterium might have been trapped in a primitive cell and have become a symbiont ensuring the energy supply to the cell, while receiving from the cell metabolites and intermediates for macromolecular syntheses. The majority of its genetic information might then have passed to the nucleus. The opposite theory on the evolutionary origin of mitochondria and **chloroplasts**, called *episomic*, hypothesized a segregation to the organelles of nuclear genetic information, leading to the formation of efficient localized functions such as respiration or [photosynthesis](#). However, the recent finding of a very large protozoan mitochondrial genome (3) strongly supports the hypothesis of the infective origin of this organelle.

In the majority of mitochondrial genomes, mitochondrial ribosomal and [transfer RNA](#) are encoded by mtDNA, while other components of the mitochondrial protein biosynthesis apparatus, including enzymes, protein synthesis factors and ribosomal proteins (with few exceptions), as well as factors involved in RNA splicing, DNA transcription, and [DNA replication](#), are coded by nuclear genes. Consequently the question might be posed as to why such a complex apparatus is needed just to ensure the localized synthesis of a few mitochondrially encoded subunits of respiratory enzyme complexes.

Another very relevant aspect of mitochondrial protein biosynthesis lies in its deviations from the universal [genetic code](#). In fact the nonuniversality of the genetic code was first discovered in mitochondria. Actually, the triplet UGA, which has no sense in the universal genetic code (see **Nonsense and unassigned codons**), codes for [tryptophan](#) residues in mammalian and fungal genomes. However, deviations from universality are different in various species, or they may not exist, as is the case in the above-mentioned protozoon (*Reclinomonas americana*). This supports the notion that the translation system in the mitochondrial genomes, which code for few proteins, has more relaxed constraints and that this has resulted in the appearance of independent deviations from the universal genetic code and from the general rules that govern protein biosynthesis. Other unusual features have appeared in many mitochondrial protein synthesis apparatus, including the two-letter codon–anticodon recognition, resulting in the diminished number of [transfer RNA](#) (tRNA) necessary for protein synthesis (4), and the use of tRNA that have highly modified structures (in mammals) or that are imported from the cytoplasm (in plants and some protozoa). Another quite exceptional feature present in some mitochondria is the phenomenon of editing of RNA (5) consisting in the

addition to RNA of a number of bases not encoded by DNA.

### Bibliography

1. B. Ephrussi, H. Hottinguer, and A. M. Chimenes (1949) *Ann. Inst. Pasteur* **76**, 351.
2. P. P. Slonimski (1953) *La formation des enzymes respiratoires chez la levure*, Masson, Paris.
3. B. F. Lang et al. (1997) *Nature* **387**, 493–497.
4. S. C. Bonitz et al. (1980) *Proc. Natl. Acad. Sci.* **77**, 3167–3170.
5. J. M. Gualberto et al. (1989) *Nature* **341**, 660.

### Mitochondrial DNA (MTDNA)

The **genomes** of [mitochondria](#) have a wide range of size, structure and informational content, and the recent sequencing of several new genomes has widely contributed to understanding the extent of this diversity (Table 1).

**Table 1. Size, Structure, and Informational Content of Some Mitochondrial DNAs<sup>a</sup>**

|  | Hs     | At               | Sc            | Cr     | Pf     | Ra     |
|--|--------|------------------|---------------|--------|--------|--------|
| Size   | 16 (1) | 367              | 70–75<br>(2)  | 16-20  | 6      | 69     |
| Molecular form   | Circle | Master<br>circle | Circle<br>(3) | Linear | Linear | Circle |
| Number of genes (4)  | 37     | 52               | 35            | 12     | 5      | 92     |
| Protein-coding genes (4)   | 13     | 27               | 8             | 7      | 3      | 62     |
| NADH-dehydrogenase<br>complex  | 7      | 9                | —             | 5      | —      | 12     |
| Succinic dehydrogenase<br>bc1 complex                                | —      | —                | —             | —      | —      | 3      |
| Cytochromoxidase   | 1      | 1                | 1             | 1      | 1      | 1      |
| ATP synthase   | 3      | 3                | 3             | 1      | 2      | 3      |
| Ribosomal proteins   | 2      | 4                | 3             | —      | —      | 5      |
| RNA polymerase   | —      | 7                | 1             | —      | —      | 27     |
| Other proteins involved in<br>macromolecular synthasesor<br>assembly | —      | —                | —             | —      | —      | 4      |
| Component of RNase P   | —      | 3                | —             | —      | —      | 7      |
| RNA-coding genes   | 24     | 25               | 27            | 5      | 2      | 29     |
| Ribosomal RNAs   | 2      | 2                | 2             | 2      | 2      | 3      |
| tRNAs  | 22     | 22               | 24            | 3      | —      | 26     |
| Component of RNase P   | —      | —                | 1             | —      | —      | 1      |



---

<sup>a</sup> Key: 1—the same size and molecular form in all other studied mammals, 2—the size of *S. cerevisiae* in mtDNA is variable mainly because of strain-dependent number of introns, 3—some other yeasts have a linear molecular form, 4—not including ORFs and intronic genes; Hs—*Homo sapiens*, At—*Arabidopsis thaliana*, Sc—*Saccharomyces cerevisiae*, Cr—*Chlamidomonas reinhardtii*, Pf—*Plasmodium falciparum*, Ra—*Reclinomonas americana*.

The mtDNA of **fungi** have sizes ranging between the 18 kbp of *Schizosaccharomyces pombe* and over 100 in some other yeasts. In these genomes, the informational content usually includes three subunits of cytochrome oxidase, two or three subunits of [ATP synthase](#) and [cytochrome b](#), in addition to the RNAs of [ribosomes](#) and [transfer RNAs](#) and, in some cases, a ribosomal subunit protein. The wide differences in size are not due to substantial differences in informational content, but to the presence of large intergenic sequences, often AT-rich, and to the presence of **introns** in the genes. An important peculiarity of these introns is the presence of open reading frames (see [Gene Structure](#)), sometimes free-standing within the introns (1) or, in some other cases, also utilizing part of the sequences of the gene. These intronic proteins have various functions, including the splicing of the intron itself (2) or its [transposition](#), which is in some cases mediated by a **reverse transcriptase** activity coded by the intron (3). The above-described mtDNA are usually circular (see [Circular Chromosome](#)), but linear in some cases. In the latter case, [telomere](#) sequences ensure [DNA replication](#) of the linear molecules.

By contrast, animal mtDNA are always circular, have a very uniform size of about 16 kbp, and have a standard informational content, which includes the same proteins found in fungal genomes, plus seven subunits of the NADH dehydrogenase complex. The genes are tightly packed, and no intergenic regions are present outside the region involved in DNA replication (4).

In plants, the size and complexity of mtDNA exhibit important variation even in closely related species. The number of mitochondrial genes is greater than in animal or fungal genomes and includes genes coding for several ribosomal proteins. The DNA organization is not clear; both linear and circular molecules have been found, bearing repeated sequences that might originate multiple [recombination](#) events.

Perhaps the most interesting developments concerning the functions of mtDNA come from the two extremes of informational content found thus far. The mtDNA of the protozoan *Reclinomonas americana*, although having a size of 69 kbp, smaller than many plant and fungal genomes, contains 92 genes. These include the genes for 27 ribosomal proteins and, interestingly, for the [elongation factor](#) tufA and for four subunits of an **RNA polymerase** of bacterial type. Moreover, the genes for factors involved in [cytochrome c](#) biogenesis, in cytochrome oxidase assembly, and even the secY protein import gene, are encoded by mtDNA. At the other extreme, the mtDNA from the parasite *Plasmodium falciparum* has a genome size of 6 kbp and possesses only five genes, which code for cytochrome b, two subunits of cytochrome oxidase, and the two ribosomal RNA.

In conclusion, only four genes, namely, those coding for subunit I of cytochrome oxidase, for cytochrome b and for the two ribosomal RNA, are present in all mtDNA studied thus far and therefore might be considered to be essential for the conservation of the mitochondrial genome. It may be pertinent that subunit I of cytochrome oxidase and cytochrome b are both the heme-carrying subunits of the corresponding respiratory complexes. The small size of mammalian mtDNA results in relatively simple [restriction maps](#). Therefore, comparisons can be used to detect subtle differences among the mtDNA. Moreover, the [mutation](#) rate of mammalian mtDNA is much higher than that of the corresponding nuclear DNA (the opposite situation has been observed for fungal and plant genomes). Consequently a large number of mutations accumulate in mammalian populations, and mtDNA comparisons have been used to evaluate **phylogenetic** differences among different species

or taxa, and to study the origins and the migrations of different human or mammalian populations.

## Bibliography

1. B. Dujon (1989) *Gene* **82**, 91–114.
2. J. Lazowska, C. Jacq, and P. P. Slonimski (1980) *Cell* **22**, 333.
3. S. Zimmerly, H. Guo, P. S. Perlman, and A. M. Lambowitz (1995) *Cell* **72**, 545–554.
4. G. Attardi and J. Schatz (1988) *Annu. Rev. Cell Biol.* **4**, 289–333.

## Suggestions for Further Reading

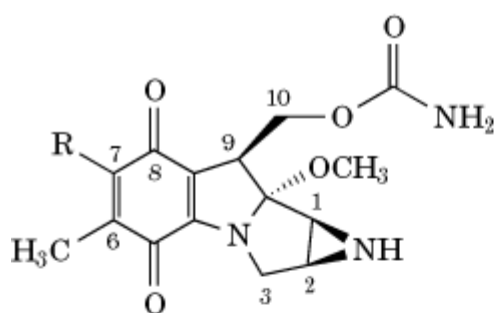
5. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, Garland Publishing Inc., Chapters "12", "14", "15".
6. N. W. Gilham (1994) *Organelle Genes and Genomes*, Oxford Univ. Press, New York.
7. J. Nunnari and D. Walter (1996) Regulation of organelle biogenesis. *Cell* **84**, 389–394 (minireview).
8. O. Poyton and E. McEwen (1996) Crosstalk between nuclear and mitochondrial genomes. *Annu. Rev. Biochem.* **65**, 563–607.
9. W. C. Wallace (1997) Mitochondrial DNA in aging and disease. *Sci. Am.* **Aug.** 22–29.
10. G. Warren and W. Wickner (1996) Organelle inheritance. *Cell* **84**, 395–400 (minireview).

## Mitomycin C

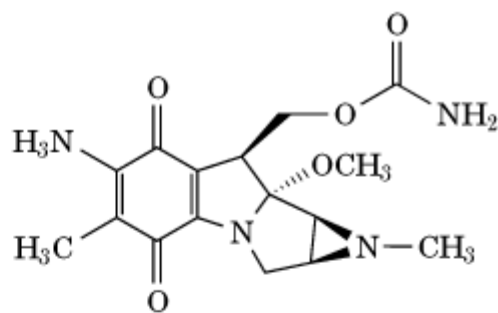
### 1. Introduction to the Anticancer Drug Mitomycin C and its Analogs

The mitomycins are a group of natural product antibiotics that were discovered in Japan in the 1950s in the fermentation cultures of *Streptomyces caespitosus* (1). Mitomycin C (MMC), mitomycin A, and porfiromycin represent three well-studied members of this family (Fig. 1). Shortly after the isolation of the mitomycins, their significant antitumor properties were discovered. MMC began use in clinical cancer chemotherapy in the 1960s and subsequently demonstrated a broad spectrum of antitumor activity against a variety of solid tumors (2). MMC continues to be an important component of combination chemotherapy for a variety of cancers, including carcinoma of the anal canal, non-smallcell lung cancer, esophageal and bladder cancer, and others (reviewed in (3)). Bone marrow suppression is the principal dose-limiting toxicity that has limited the success of MMC as a single agent. Other members of the mitomycin family, the natural product analogs, FR66979 and FR900482, were discovered in the fermentation products of *Streptomyces sandaensis no. 6897*, and the semisynthetic mitomycin derivatives, E09, KW-2149, and BMS-181174 (previously designated as BMY25067) (reviewed in (4)), have also been developed (Fig 1).

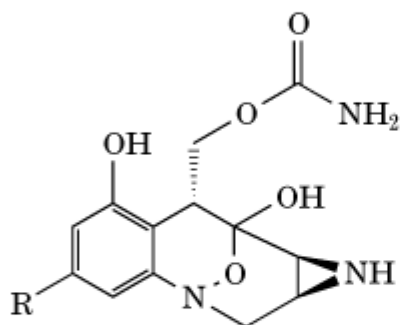
**Figure 1.** Structures of mitomycin antibiotics and common analogs. Mitomycin C, mitomycin A, and porfiromycin represent three well-studied members of the mitomycin family. More recently, two other members of the mitomycin family, the natural product analogs, FR66979 and FR900482, were discovered, and the semisynthetic mitomycin derivatives, E09, KW-2149, and BMS-181174 (previously designated as BMY25067), have also been developed.



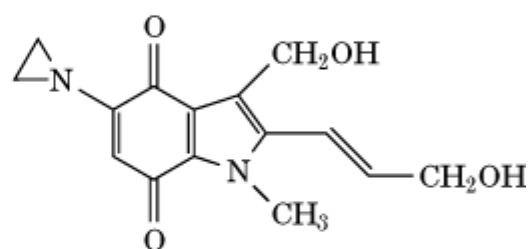
Mitomycin C R = NH<sub>2</sub>  
Mitomycin A R = OCH<sub>3</sub>



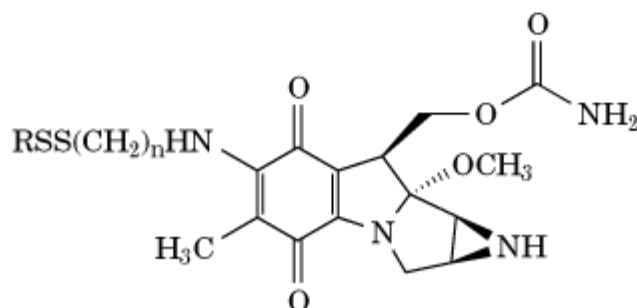
Porfiromycin



FR66979 R = CH<sub>2</sub>OH  
FR900482 R = CHO



EO9

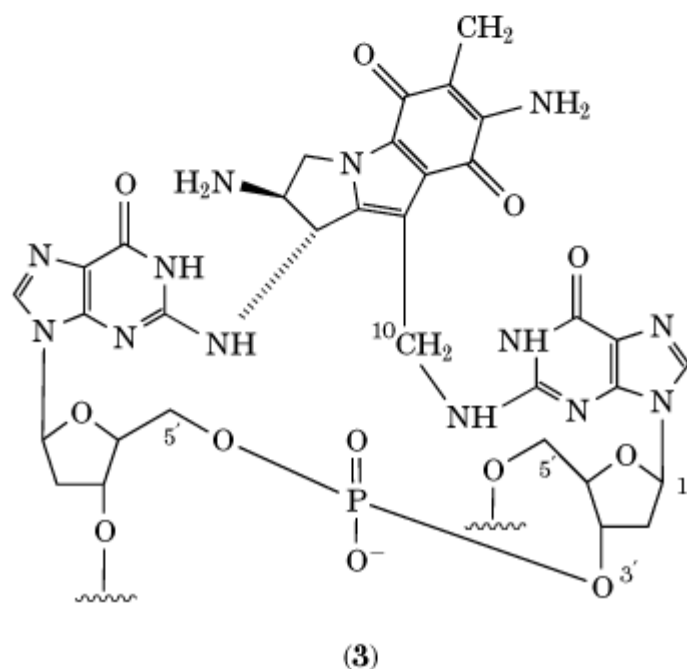
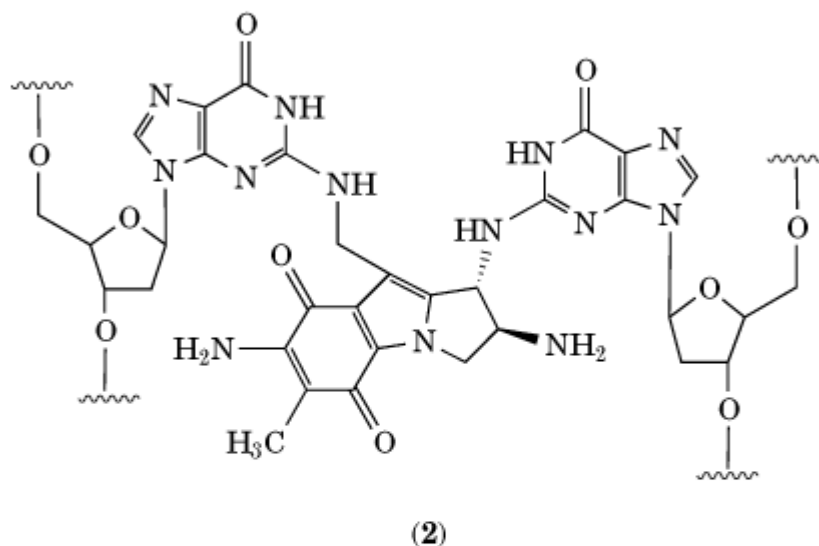
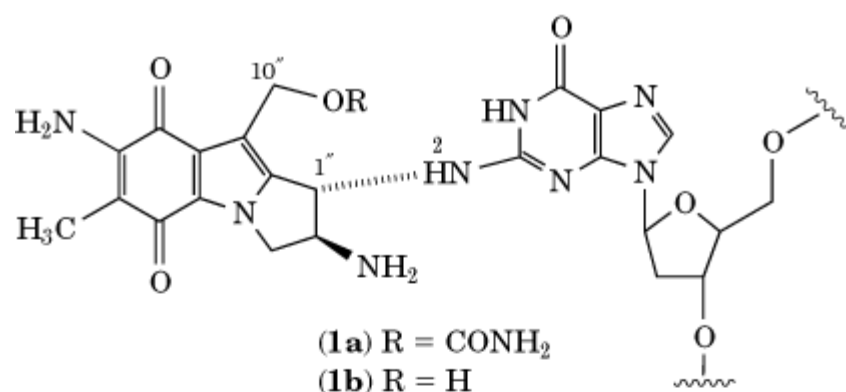


BMS-181174 R = *p*-nitrophenyl  
KW-2149 R = HOOC(NH<sub>2</sub>)CH(CH<sub>2</sub>)<sub>2</sub>CONH(CH<sub>2</sub>)<sub>2</sub>

MMC has served as a model for investigating the family of mitomycins, and the vast majority of knowledge about this group of agents is derived principally from extensive laboratory and clinical studies of MMC over the past 25 years. Like most genotoxic agents, MMC binds to DNA, RNA, and other macromolecules, and MMC–DNA adducts are believed to be principally responsible for the majority of its observed biological effects. Also like many genotoxic agents, MMC requires metabolism from the parent form to a reactive intermediate to bind to macromolecules and exert its biological effects. However, unlike most other organic genotoxins, MMC requires chemical or biochemical reduction rather than oxidation. Although the chemistry of this reduction has been well characterized *in vitro*, the specific biochemical pathways by which this occurs inside the cell remain controversial. MMC has been shown to induce many of the common biological hallmarks of other genotoxic cancer chemotherapy agents and genotoxic mutagens and carcinogens. MMC has been reported to inhibit DNA replication and RNA synthesis, induce mutations and the SOS response in bacteria, and induce mutations and sister chromatid exchanges in mammalian cells (reviewed in (5)).

Reduced MMC covalently binds preferentially to the N<sup>2</sup> position of guanine in the minor groove of duplex DNA, forming monoadducts, intrastrand crosslinks at adjacent guanines, and interstrand crosslinks at opposite guanines of adjacent basepairs, specifically in the sequence orientation of CpG, but not GpC (see Fig. 2). It is currently presumed that the CpG MMC interstrand crosslink is principally responsible for the cytotoxic and antitumor properties of MMC; however, formal evidence for this is still lacking.

**Figure 2.** Mitomycin C-DNA adducts. After reductive activation, mitomycin C (MMC) can form two monoadducts, i.e., a monofunctional monoadduct (**1a**) and a bifunctional monoadduct (**1b**), and also two types of cross-links, i.e., an interstrand cross-link (**2**) at CpG sites, and an intrastrand cross-link (**3**) at GpG sites.



## 2. MMC Structure and DNA Adduct Formation

The structures of the mitomycins are unique, containing three different functional groups thought to be important for their properties and present separately in several other known carcinogens. These

reactive groups, the aziridine (three membered ring attached to C<sup>1</sup> and C<sup>2</sup>), carbamate (–OCONH<sub>2</sub>), and quinone (five-membered ring containing ketones at C<sup>5</sup> and C<sup>8</sup>), are arranged about a pyrrolo [1,2-*a*]indole nucleus (see [1](#)). In this arrangement, these functional groups are relatively inactive; however, upon chemical or biochemical reduction of the quinone by either a one- or two-electron process, a semiquinone or hydroquinone is formed, respectively, leading to the unveiling of two electrophilic sites (C<sup>1</sup> and C<sup>10</sup>) for alkylation. According to the proposed mechanism of alkylation (reviewed in [4](#)), there are two pathways the MMC and follow: the “monofunctional” alkylation pathway, which results in a “monofunctional monoadduct,” (**1a**) and the “bifunctional” alkylation pathway, which results either in a “bifunctional monoadduct” (**1b**) or a cross linked adduct (**2** or **3**). The sequence of events leading to the activation of the C<sup>1</sup> and C<sup>10</sup> positions of MMC involves the reduction of the quinone ring, the spontaneous elimination of methanol (driving force: formation of the pyrrole), and the opening of the aziridine ring. Once the electrophilic C<sup>1</sup> and C<sup>10</sup> sites are “unmasked,” the next step is the highly specific nucleophilic attack by the N<sup>2</sup> exocyclic amino group of guanine on the C<sup>1</sup> position of MMC. If reoxidation of the quinone ring, by either O<sub>2</sub> (poor redox conditions) or by the presence of excess MMC, occurs at this stage, the compound is no longer activated and is referred to as a “monofunctional monoadduct” (**1a**), where monoadduct refers to one addition product with DNA, and monofunctional refers to a reaction at only one functional group on the MMC ring system. This reoxidation, which is autocatalytic, occurs when the initial reduction of MMC is very slow, ie, when it can be inactivated faster than it is formed.

In the “bifunctional” pathway, after the first alkylation of DNA, the C<sup>10</sup> position of MMC undergoes decarbamoylation that results in the activation of this position for further nucleophilic attack. If there is a second guanine present adjacent to or on the opposing strand of DNA (ie, at GpG or CpG sequences, respectively), then the N<sup>2</sup> position of this guanine can serve as the incoming nucleophile on the C<sup>10</sup> position of MMC, forming an N<sup>2</sup>-G-MMC-N<sup>2</sup>G intrastrand or interstrand crosslink (**2** or **3**), respectively). If there is not a second guanine in either of these two sites, H<sub>2</sub>O can serve as the incoming nucleophile and attack the C<sup>10</sup> position, forming a “bifunctional monoadduct” (**1b**). The term “bifunctional monoadduct” refers to the fact that only one addition product with DNA has formed although the mitomycin complex has undergone nucleophilic reactions at two functional groups. Other less abundant MMC–DNA adducts have also been observed and their biological activity has been investigated *in vitro* ([6](#)).

### 3. Bioreduction

Due to its requirement for reduction prior to alkylating DNA, MMC is referred to as a bioreductive alkylating agent (reviewed in [7](#)). The bioactivation of MMC can be mimicked *in vitro* by chemical reducing agents such as by H<sub>2</sub> and PtO<sub>2</sub>, which favor the formation of a monofunctional adduct with DNA; by sodium dithionite, which favors the formation of a bifunctional adduct; or by various nicotinamide adenine dinucleotide (NAD(P)H)-dependent flavoreductases. The relative proportion of adducts *vitro* or *in vivo* depends on the sequence of the DNA and on the metabolic and redox status of the cells.

Because bioreduction is key to the cytotoxic actions of MMC, a number of laboratories have dedicated a large effort to identifying specific biochemical pathways involved in the activation of MMC. A variety of bioreductive enzymes have been implicated in the reduction of MMC, such as NAD(P)H cytochrome *c* reductase, NAD(P)H quinone oxidoreductase (DT-diaphorase), xanthine oxidase, xanthine dehydrogenase, and NADH cytochrome b<sub>5</sub> reductase (reviewed in [8](#)). Despite this effort, much controversy still exists regarding the specific roles of each of these pathways in intact cells. The large number of divergent and sometimes contradictory *in vitro* and cell culture studies suggest that no single enzyme is uniquely responsible for MMC activation (reviewed in [7](#)). These studies also suggest that a number of variables can influence the activation of MMC, at least

*in vitro* and in cell culture, including the cell type and cell culture growth conditions, pH, and degree of hypoxia (7, 8). For example, *in vitro* chemical reduction studies indicate that the pH can influence both the site of the lesions on guanine (N<sup>7</sup> vs N<sup>27</sup> alkylation), as well as influence the predominating bioreduction pathway (4). However, whether such parameters also influence MMC metabolism and MMC-induced biological effects *in vivo* is not clear, because pH, oxygen concentration, and other parameters are tightly regulated over a much narrower range under physiological conditions in the intact animal. Such studies demonstrate the limitations of *in vitro* and cell culture systems in modeling the complexities of a whole organism or the microenvironment of a tumor *in vivo*.

Several studies have suggested that DT-diaphorase and cytochrome P450 reductase are often highly overexpressed in cancer cells (7, 9). The terminology “enzyme-directed drug development” - coined with mitomycin as the paradigm - refers to an approach that seeks to identify enzymes that are overexpressed in tumors when compared with normal tissues and that can be used to preferentially target drug effects to these tumor cells (reviewed in (10)). For example, the sensitivity of bladder cancer patients to MMC has been shown to be correlated with the expression of DT-diaphorase and cytochrome P450 reductase in their bladder tumors, suggesting that the lower efficacy of intravesical MMC in these patients may be due, at least in part, to their lower expression of these enzymes (9). An important area of future study will be to determine the precise biochemical pathways by which MMC is bioreduced in tumor and normal cells, and to determine the role of each pathways in the biological effects of MMC in these target and nontarget tissues. This is critical for understanding, and hopefully being able to modulate, the efficacy of MMC and its analogs as antitumor agents against various human cancers.

#### 4. Characterization and Measurement of MMC–DNA Adducts

The structures of the various MMC–DNA adducts have been extensively characterized by a variety of biophysical techniques—including high-performance liquid chromatography, ultraviolet absorbance, thermal melting curves, gel electrophoresis, circular dichroism, chemical footprinting, Fourier transform infrared spectroscopy, fast atom bombardment mass spectroscopy, and a variety of one and two-dimensional nuclear magnetic resonance techniques (4, 11-13)). There is general agreement from these studies that the predominant MMC–guanine monoadducts and the CpG MMC–DNA interstrand cross-link cause little or no perturbation of the DNA double helix, with the attached mitosene ring lying snugly within the minor groove of the DNA. In contrast, the GpG MMC–DNA intrastrand cross-link has recently been shown to bend the DNA helical axis by 14.6° (reviewed in (4)).

Covalent MMC–DNA adducts have been detected and isolated following *in vitro* reactions involving mono- or di-nucleotides, oligonucleotides, and purified DNA (eg, calf thymus DNA, *Micrococcus luteus* DNA, salmon sperm DNA, and poly(dG-dC)·poly(dG-dC) (11, 13)). MMC–DNA adducts have also been detected in various cell culture models, in intact animals, and in human cancer chemotherapy patients (14, 15). Due to the lack of readily available radiolabeled MMC, some studies have substituted the commercially available radiolabeled *N*-methyl analog of MMC, porfiromycin. Porfiromycin has been shown to have similar reactivity and biological properties as MMC, and the corresponding porfiromycin *N*-methyl analogs of the monofunctional monoadduct, bifunctional monoadduct, and cross-links of MMC have been observed (reviewed in (13)).

#### 5. Modulation of Binding Activity at Methylated CpG Sites

Sequence context plays an important role in the types and quantity of MMC adducts formed. MMC–DNA monoadducts and interstrand cross-links have been observed to show a strong preference for the 5′CpG3′ but not the 5′GpC3′ dinucleotide sequence (reviewed in (4)). There also appears to be an enhancement of MMC cross-linking at those sequences that have a purine 5′ to the CpG site and a pyrimidine (particularly thymine) 3′ to the guanine (ie, 5′PuCGPyr3′) (4). The CpG dinucleotide sequence has some unique features that are of interest with regard to its specificity for MMC cross-

linking. The CpG dinucleotide sequence is underrepresented in eukaryotic genomes compared to other G-containing dinucleotide sequences, occurring at about one fifth its expected statistical frequency (reviewed in (16)). However, clusters of CpG containing regions, referred to as “CpG islands,” are often found at much higher than predicted frequencies within many gene promoters, and 60 to 90% of these sequences are methylated at the 5' position of cytosine *in vivo*. Cytosine methylation is thought to be one mechanism for transcriptional regulation of genetic material. For example, the degree of cytosine methylation is often but not always inversely correlated with transcriptional activity. Cytosine methylation also alters the local structure of duplex DNA. Thus, an enhancement in MMC cross-linking at methylated CpG sites (reviewed in (4)) may be biologically significant because it has been proposed that DNA-damaging agents may heritably alter DNA methylation patterns. This would provide an epigenetic mechanism for heritable alterations in expressions of genes involved in carcinogenesis, and might also be one way in which MMC exerts its biological effects.

## 6. MMC–DNA Adduct Recognition and Repair

Little is known about the recognition and repair of MMC–DNA adducts *in vivo*. It is generally assumed that MMC–DNA adducts are repaired by nucleotide excision repair (NER), which is an important cellular mechanism that removes radiation-induced and chemically induced damage for DNA. Repair of the interstrand cross-link in particular is believed to pose a unique and difficult topological challenge to the repair machinery. However, the rapid kinetics of MMC cross-link removal from DNA in cell culture and *in vivo* argues for effective recognition and repair of this lesion, despite this topological problem and the relatively non distorting nature of this lesion. Several specific protein complexes from human and rodent cell nuclear extracts were found to preferentially recognize and bind to the MMC interstrand cross-link with high affinity (17). Several of these complexes appear to contain ERCC-1 (excision repair cross complementation group 1) and/or XPA (xeroderma pigmentosum complementation group A), which have both been shown to be involved in mammalian NER and to be implicated in cross-link repair (18). *In vitro* experiments have demonstrated that both purified XPA and the minimal DNA binding domain of XPA (XPA-MF122) were able to bind MMC–cross-linked DNA with a much greater specificity and higher affinity than to undamaged DNA (19). This occurred in the absence of other proteins from the NER complex (19). Such preferential binding may lead to enhanced recruitment of the NER machinery to the adduct site for subsequent efficient repair. Further studies of this type will be required to fully elucidate the biochemical mechanisms by which DNA interstrand cross-links formed by MMC are recognized and repaired in mammalian cells.

## 7. Effects of MMC on Gene Expression

In addition to its overt effects as a cytotoxic and antitumor agent, recent studies have demonstrated that lower, noncytotoxic doses of MMC can selectively modulate expression of certain target genes at doses that do not induce apoptotic or stress response, and that have no effect on expression of most other genes. For example, selective suppression of several genes associated with multidrug resistance in cancer cells was observed in several human and rodent cancer cell lines in culture (20), in a mouse tumor model *in vivo* (21), and in the tumors of human cancer patients *in vivo* (22) after a single low-dose MMC treatment. In the cell culture and rodent studies, pretreatment with MMC significantly increased the subsequent antitumor activity of a second agent such as doxorubicin or taxol (20, 21). Extensive dose–response, time course, and mechanistic studies indicated that this was principally a result of the suppression of these drug-resistance proteins by the MMC pretreatment rather than combined toxicity *per se* (20–22). Similar results were seen after pretreatment with other DNA cross-linking agents of divergent chemical structure, including cisplatin and carboplatin, suggesting that these effects were a result of formation of DNA cross-link adducts. These results suggest the basis for development of novel combination chemotherapy approaches, where an agent such as MMC is used as a modulator of the cancer cell phenotype in combination with other cytotoxic agents, rather than, or in addition to, its use as a cytotoxic agent (22).



In addition to cancer applications, MMC has been shown to selectively alter expression of other genes that may also be of benefit in noncancer clinical diseases. For example, MMC has been a useful model compound for examining effects of pharmacological agents on the gene regulation and protein biogenesis of cystic fibrosis transmembrane conductance regulator (CFTR), mutations of which are entirely responsible for the diseases, cystic fibrosis. At low doses, MMC induced a significant increase in both CFTR mRNA and protein expression. Because the problem in cystic fibrosis is essentially a lack of adequate CFTR functional expression, such a drug-induced increase in CFTR expression could potentially be of clinical benefit in treatment of cystic fibrosis patients (23). Identification of less toxic analogs of MMC that also upregulate CFTR expression is therefore of potential interest in this setting.

Low doses of MMC have also been shown to selectively modulate the expression of the xenobiotic- and drug-inducible genes, 5-aminolevulinate synthase (24, 25) and cytochrome P450 *CYP2H1* (24, 25), and the glucocorticoid hormone-inducible phosphoenolpyruvate carboxykinase gene (24, 26)). The mechanistic basis for these selective effects on gene expression is not known. Correlative evidence indicates that these effects are primarily a result of formation of MMC-DNA adducts within or near the promoter regions of these genes. However, whether this is a result of preferential adduction of these regions by MMC (ie, higher levels of adducts at these sites), or a higher sensitivity of these particular genes to nearby MMC-DNA adducts remains to be determined.

## 8. Mitosene Analogs

Preclinical and clinical trials, as well as basic mechanistic studies, with the mitomycin derivatives KW-2149, BMS-181174, EO9, FR66979, and FR900482 (Fig. 1) have suggested interesting activity profiles. Due to differences in structure, these derivatives can differ from MMC in their mechanism of bioreductive activation, cross-linking efficiency, dose-limiting toxicity, antitumor activity, and/or activity against resistant cell lines (4, 27-30). The aim of such studies with MMC analogs is to develop more effective therapies, in particular treatment regimens that are active in drug-resistant cancers, that do not induce drug resistance after prolonged treatment, and that reduce nontarget tissue toxicities. However, the clinical usefulness of these newer analogs remains to be determined. For example, BMS-181174 initially showed good promise as an antitumor agent in preclinical trials, but subsequent human phase I clinical trials revealed substantial dose-limiting toxicities that precluded its continued clinical use (28). Several other MMC analogs are currently being evaluated in preclinical and early clinical trials.

## 9. Summary

In summary, MMC and its analog represent a unique group of compounds with respect to both their interesting chemistry and biology, and their usefulness as cancer chemotherapy agents. Although much has been learned about this class of chemicals over the past 25 years, based principally on studies of MMC, there is still much to learn about these agents, including the biochemical pathways responsible for their activation and detoxification, the specific DNA lesions induced *in vivo*, their repair, and their individual roles in the biological effects observed, and the further development of these agents as gene modulators and clinical cancer chemotherapy agents.

## Bibliography

1. T. Hata et al (1956) *J. Antibiotics. Ser. A.* **9**, 141-146.
2. C. D. Doll, R. B. Weiss, and B. F. Issell (1984) *J. Clin. Oncol.* **3**, 276-286.
3. W. T. Bradner (2001) *Cancer Treat. Rev.* **27**, 35-50.
4. M. Tomasz and Y. Palom (1997) *Pharmacol. Ther.* **76**, 73-87.
5. E. W. Vogel and M. J. Nivard (1994) *Mutat. Res.* **305**, 13-32.
6. Y. Palom et al (2001) *Biochem. Pharmacol.* **61**, 1517-1529.
7. J. Cummings et al (2001) *Biochem. Pharmacol.* **56**, 405-414.

8. V. J. Spanswick, J. Cummings, and J. F. Smyth (1998) *Gen. Pharmacol.* **31**, 539–544.
9. Y. Gan et al (2001) *Clin. Cancer Res.* **7**, 1313–1319.
10. H. D. Beall and S. I. Winski (2000) *Front. Biosci.* **5**, D639–D648.
11. A. J. Warren and J. W. Hamilton (1996) *Chem. Res. Toxicol.* **9**, 1063–1071.
12. D. Norman et al (1990) *Biochemistry* **29**, 2861–2875.
13. M. Tomasz (1994) In *Molecular Aspects of Anticancer Drug-DNA Interactions*, Volume 2 (S. Neidle and M. Waring, eds.) CRC Press, Boca Raton, FL, pp. 313–349.
14. A. J. Warren, A. W. Maccubbin, and J. W. Hamilton (1998) *Cancer. Res.* **58**, 453–463.
15. A. J. Warren, D. J. Mustra, and J. W. Hamilton (2001) *Clin. Cancer Res.* **7**, 1033–1042.
16. S. B. Baylin et al (1998) *Adv. Cancer Res.* **72**, 141–196.
17. A. J. Warren et al (1998) *Environ. Mol. Mutagen.* **31**, 70–81.
18. R. D. Wood (1995) *Phil. Trans. R. Soc. Lond. B.* **347**, 69–74.
19. D. J. Mustra, A. J. Warren, and J. W. Hamilton (2001) *Biochemistry* **40**, 7158–7164.
20. M. A. Ihnat et al (1997) *Clin. Cancer Res.* **3**, 1339–1346.
21. M. A. Ihnat et al (1999) *Oncol. Res.* **11**, 303–310.
22. S. P. Anthony et al (1997) *Proc. Am. Soc. Clin. Oncol.* **16**, 230A.
23. R. Maitra et al (2001) *Cell Physiol. Biochem.* **11**, 93–98.
24. J. W. Hamilton et al (1994) *Ann. N.Y. Acad. Sci.* **726**, 343–345.
25. R. M. Caron and J. W. Hamilton (1995) *Environ. Mol. Mutagen.* **25**, 4–11.
26. R. M. Caron and J. W. Hamilton (1998) *J. Biochem. Mol. Toxicol.* **12**, 325–337.
27. L. Y. Dirix et al (1996) *Eur. J. Cancer* **32A**, 2019–2022.
28. A. S. Planting et al (1999) *Anticancer Drugs* **10**, 821–827.
29. J. H. Schellens et al (2001) *Anticancer Drugs*, **12**, 583–590.
30. D. Rohde et al (1998) *Urol. Res.* **26**, 243–247.

### Suggestions for Further Reading

31. D. Ross, H. D. Beall, D. Siegel, R. D. Traver, and D. L. Gustafson, (1996) *Enzymology of bioreductive drug activation*. *Br. J. Cancer* **74** (Suppl. 27), S1–S8
32. M. Tomasz, (1994) In *DNA Adducts: Identification and Biological Significance* (K. Hemminki, A. Dipple, D. E. G. Shuker, F. F. Kadlubar, D. Segerback, and H. Bartsch, eds.), IARC Lyon, pp. 349–357.

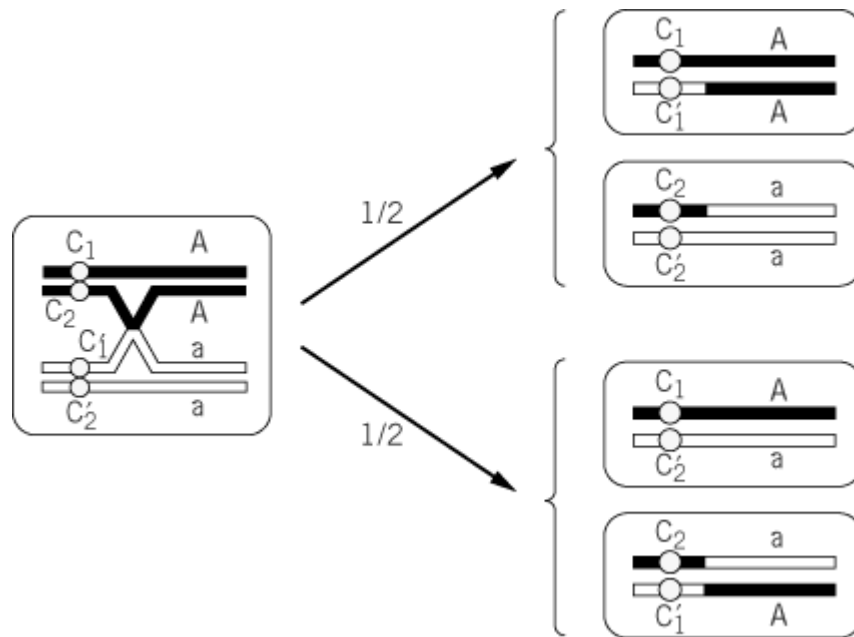
### Mitotic Recombination

Genetic [recombination](#) usually occurs in meiosis, but it can also occur in **diploid** cells that reproduce asexually by mitosis. Mitotic recombination breaks the general rule that genotypes are maintained during asexual reproduction.

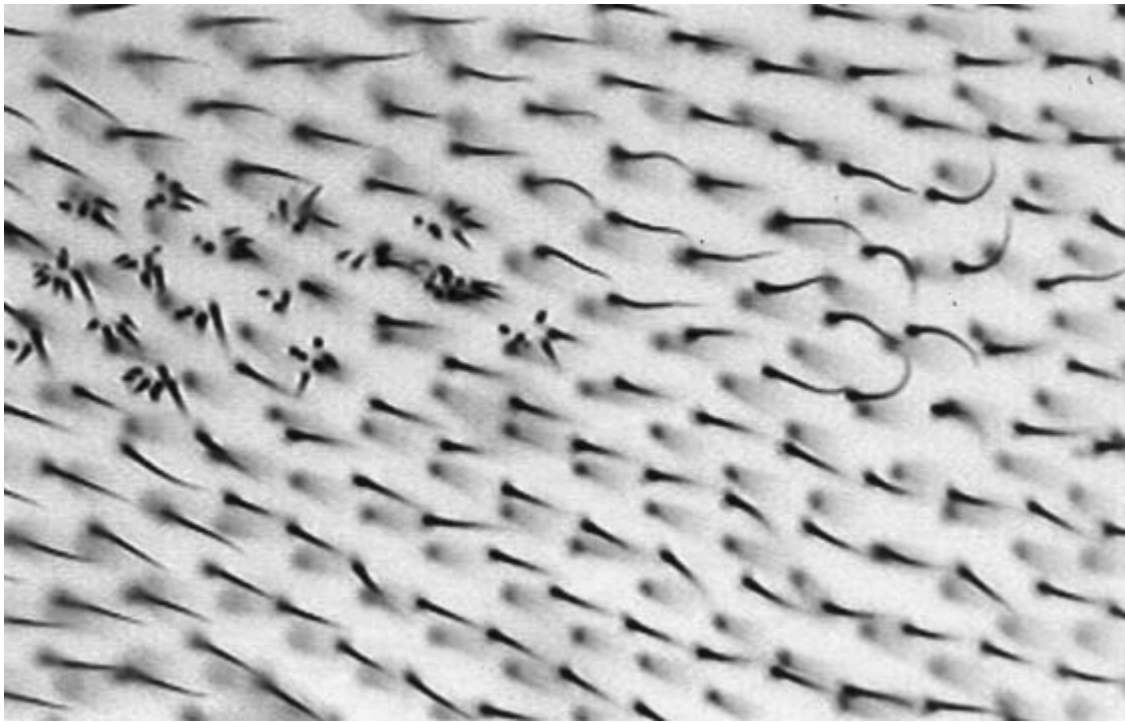
Consider a clone of heterozygous cells, genotype  $Aa$ , that carry the **alleles**  $A$  and  $a$  of a certain gene. The clone can be a colony of an unicellular organism or a tissue of a multicellular organism. A recombination event between the gene and the [centromere](#) of its [chromosome](#) may give rise to two homozygous cells,  $AA$  and  $aa$  (Fig. 1). Further mitotic growth will produce two new cell clones, genetically distinct from the original cell. In the absence of cell migration, the new clones form twin

spots in the colony or the tissue. In the case of complete dominance of a cellular genetic marker (eg, the cell color produced by the *yellow* allele of *Drosophila*), only the spot with the recessive marker will be recognized easily. The twin spots will be seen if the phenotype of the heterozygote differs from that of both homozygotes or, with complete dominance, in diploids  $Ab/aB$  when both genes are present on the same chromosome arm and recombination occurs between the centromere and the closest gene (Fig. 2). The concept of mitotic recombination was derived from the observation of twin spots in *Drosophila* (1).

**Figure 1.** Mitotic recombination between a heterozygous locus and the centromere (C) of its chromosome. Recombination occurs when each chromosome has two chromatids. Depending on the movements of the centromeres in the next anaphase, the recombination event leads to homozygosity of the alleles (above) or to an exchange between the homologous chromosomes (below).



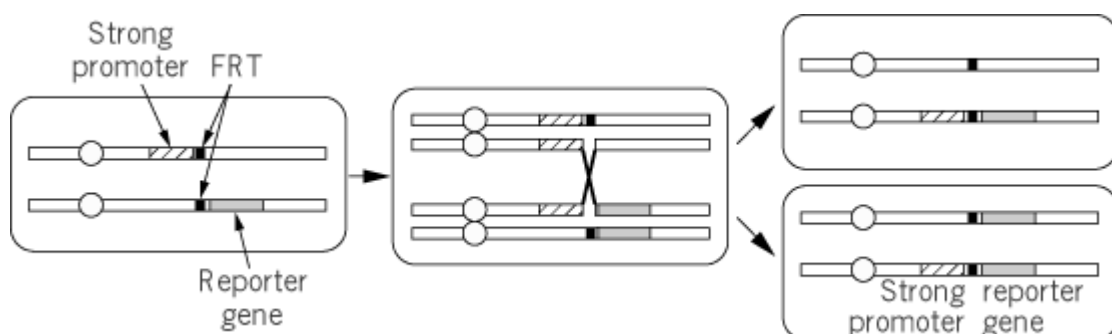
**Figure 2.** Twin spots in the surface of a *Drosophila melanogaster* adult fly heterozygous for the markers *mwh* (multiple wing hairs) and *f* (forked). Mitotic recombination in the left arm of the the third chromosome led to the formation of a spot with multiple processes per cell and another with bent cell processes, surrounded by cells of wild-type phenotype.



Mitotic recombination is rare; the frequency of alleles that become homozygous lies usually in the order of  $10^{-4}$  to  $10^{-5}$  per nuclear division, and it has been estimated that in *Aspergillus nidulans* the frequency over the whole genome is about 2% per nuclear division. The incidence can be increased by exposure to recombinogenic agents, such as X-rays, other ionizing radiations, ultraviolet radiation, and certain chemicals, including many alkylating agents. With appropriate treatments, the frequency of mitotic recombination reaches the order of magnitude of meiotic recombination.

Mitotic recombination provides a powerful tool to mark all the descendants of a single cell during development of multicellular organisms. A recombinogenic treatment of an appropriate heterozygote at a certain stage of development will produce genetically marked cell clones that can be easily recognized in later stages of development. Mitotic recombination can be made time-controlled and site-specific, as indicated in Figure 3. (see top of next page)

**Figure 3.** Time-controlled and site-directed mitotic recombination. FRT, a short DNA sequence from the 2- $\mu$ m plasmid of *Saccharomyces cerevisiae* that is the target for the action of a specific recombinase, is inserted at the site where mitotic recombination should occur. The gene for the recombinase is inserted, under an inducible promoter, elsewhere in the genome. Induction of this promoter will trigger mitotic recombination at the FRT site. Cells derived from this event may be marked by a strong promoter and a reporter gene that are brought together by recombination (see Ref. 2). See previous page.



Mitotic recombination between two genes of the same chromosome may produce new combinations of the respective alleles. The frequency of recombinants correlates with the physical distance between the genes and may be used for the construction of genetic maps. In some fungi, like *Aspergillus niger*, this is the only way to construct genetic maps by recombination. Mitotic and meiotic recombination maps conserve the order of genes, although the respective recombination frequencies are not proportional.

The [parasexual cycle](#), which includes fusion of haploid nuclei, mitotic recombination, and random chromosome loss, allows natural strains to form new combinations of genetic characters independently of the existence of a sexual cycle and meiosis.

### Bibliography

1. C. Stern (1936) *Genetics* **21**, 625–730.
2. D. Harrison and N. Perrimon (1993) *Curr. Biol.* **3**, 424–433.

### Suggestions for Further Reading

3. E. Käfer (1977) Meiotic and mitotic recombination in *Aspergillus* and its chromosomal aberrations. *Adv. Genet.* **19**, 33–131.
4. B. A. Kunz and R. H. Haynes (1981) Phenomenology and genetic control of mitotic recombination in yeast. *Annu. Rev. Genet.* **15**, 57–89.
5. G. S. Roeder and S. E. Stewart (1988) Mitotic recombination in yeast. *Trends Genet.* **4**, 263–267.
6. J. J. Panthier and H. Condamine (1991) Mitotic recombination in mammals. *Bioessays* **13**, 351–356.

### Mobile Element

**Genomes** are surprisingly plastic, containing a variety of different kinds of DNA that can move from place to place. These mobile elements generally move by site-specific [recombination](#)—that is, by the action of a specialized [recombinase](#), usually encoded by the mobile element, on special DNA sites at the ends of the element where recombination occurs; these [recombination](#) reactions can occur by a variety of mechanisms. One type of mobile DNA present in virtually all organisms examined are transposable elements. A major pathway for transposition is a DNA-based mechanism in which the transposon is excised from flanking donor site DNA by a double-strand break; alternatively, a DNA copy of the element is made by **reverse transcription** of an RNA transcribed from a chromosomal copy of the element. Both these processes result in the production of exposed 3'OHs at the ends of the mobile DNA segment. These exposed 3'OH ends are then joined to the target DNA, covalently linking the mobile element to its new site of insertion. Transposition is used by viruses to join the element to the genome of a newly infected cell. Other transposable elements move within a cell and often encode determinants, such as **antibiotic-resistance** genes that affect the cell. Nonviral transposable elements can also be dispersed between cells by translocating onto plasmids or viruses that move from cell to cell.

Mobile elements include, but are not limited to, [transposable elements](#)—that is, those elements that

move between nonhomologous sites via the recombination reaction called [transposition](#), and transposable elements may be called *mobile elements*. There are a number of mobile DNA segments that are not transposable elements and that translocate by mechanisms other than transposition.

Another reaction in which DNA segments can rearrange to make new genes occurs by a pathway distinct from transposition, but involves a recombinase that can carry out transposition-type reactions (1, 2). VDJ recombination is an ordered set of deletion reactions that assembles immunoglobulin genes from multiple gene segments, providing incredible diversity in the immune system (3) (see [Gene Rearrangement](#)). As in transposition, the breakage and joining steps also involve DNA breakage reactions to expose 3'OH ends, followed by the joining of these 3'OH ends to intramolecular targets.

Another distinct kind of element translocation that occurs between nonhomologous sites is the movement of **polyA**<sup>+</sup> elements such as the [LINES](#) and [SINES](#) of mammalian genomes that lack all [inverted terminal repeats](#) (4). Translocation involves an element-generated nick in the target DNA that serves as a primer for the *in situ* reverse transcription of an element RNA, resulting in a new DNA copy of the element at a new chromosomal location.

Another nontransposable element mobile DNA is exemplified by the [lambda phage](#) that integrates into the bacterial chromosome to become part of the chromosome during [lysogeny](#) and then excises from the chromosome to enter the lytic phase. These recombination reactions occur at specialized sites on both the viral and chromosomal genomes; these sites are related to each other by protein binding sites and a short region of [homology](#). Recombination occurs by simple breakage and joining reactions at each site (5, 6), in a reaction called *conservative site-specific recombination* (CSSR). This is a very distinct reaction from transposition, where the two transposon ends interact with a nonhomologous target site. Other CSSR-mediated reactions include the inversion of bacterial chromosomal segments to alternate between two configurations for differential gene expression and the monomerization of various dimeric DNA molecules, from plasmids to chromosomes. Yet another kind of reaction that can move different DNA sequences from place to place involves the introduction of gap at a particular site by a double-strand break that undergoes exonucleolytic trimming, followed by the filling of this gap via pairing and DNA replication, with information from a related site containing similar but distinct sequences (7).

## Bibliography

1. M. Melek, M. Gellert, and D. van Gent (1998) *Science* **280**, 301–303.
2. D. C. van Gent, K. Mizuuchi, and M. Gellert (1996) *Science* **271**, 1592–1594.
3. D. A. Ramsden (1997) *Curr. Opin. Immunol.* **9**, 114–120.
4. J. V. Moran, S. E. Holmes, T. P. Naas, R. J. DeBerardinis, J. D. Boeke, and H. H. Kazazian *Cell* in press.
5. N. D. Grindley (1997) *Curr. Biol.* **7**, R608–R612.
6. K. Mizuuchi (1997) *Genes Cells* **2**, 1–12.
7. A. Shinohara and T. Ogawa (1995) *Trends Biochem. Sci.* **20**, 387–391.

## Molecular Averaging

Molecular averaging is an important technique in [X-ray crystallography](#) for refining the structure generated (see also **R factor**). The structure determined by [crystallography](#) is the [asymmetric unit](#) of

the [unit cell](#) of the crystal lattice. If the asymmetric unit contains more than one molecule of the same kind, these molecules have the same structure, apart from small deviations at their surface due to molecular contacts. This structural similarity is extra information and imposes constraints on the phase angles of the X-ray reflections (see [Phase Problem](#)).

Bricogne proposed a very convenient method to employ this information for improving the electron density map of a protein ([1](#)). After a first electron density map has been calculated, for instance, based on the measured amplitudes of the reflections and phase angles estimated from [isomorphous replacement](#) or [molecular replacement](#), the phase angles can be improved by averaging. In this procedure, the positions of the protein molecules in the asymmetric unit are located in the electron density map as regions of relatively high density. An envelope is constructed around each of them. The region outside the envelopes is occupied by solvent.

For an ideal structure, the electron density within all envelopes should be exactly the same, and outside them it should be flat. In practice this is not true. Therefore, the electron density in the solvent region is flattened in the next step (see [Solvent Flattening](#)), and the density of the protein molecules is averaged. This new map of electron density is closer to the correct map than the first map, and structure factors are calculated from it. The phase angles of these new, calculated structure factors are combined with their measured amplitudes, and an improved map is calculated. The procedure is repeated until convergence is reached.

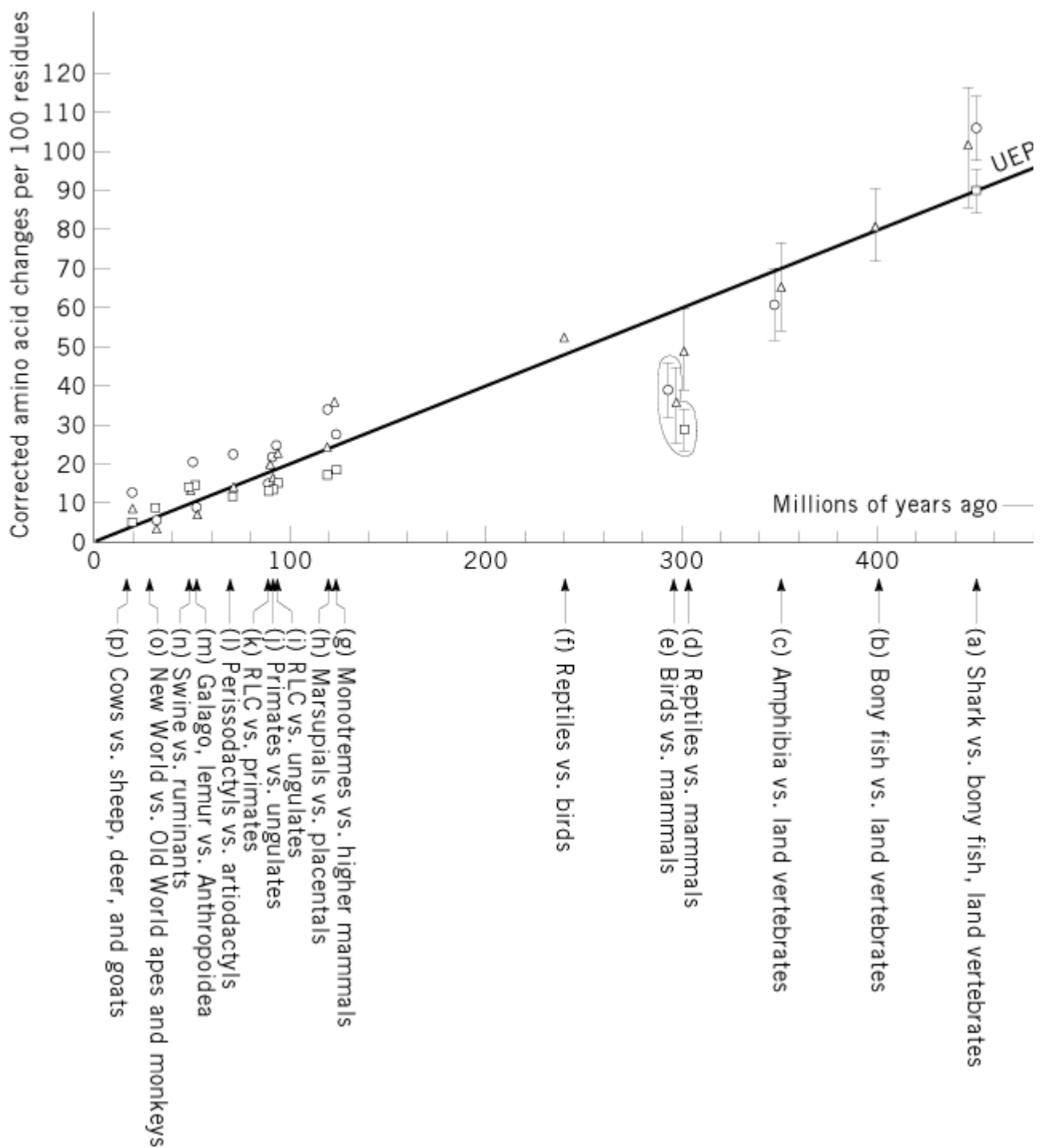
#### Bibliography

1. G. Bricogne (1974) *Acta Crystallogr.* **A30**, 395–405.

## Molecular Clock

The molecular clock is an evolutionary property of **genes** in which the number of nucleotide or amino acid substitutions between a pair of species is approximately proportional to the time since those species diverged (Fig. [1](#)), ([1](#)) which is usually given by the fossil record. Nucleotide and amino acid substitutions take place with a certain degree of regularity along evolutionary time, as if one substitution makes the needle of the evolutionary clock move ahead by a given time interval ([2](#)). Thus, this property is called “molecular clock” or “molecular evolutionary clock.”

**Figure 1.** The constant rate of divergence of the amino acid sequences of the hemoglobin and myoglobin polypeptide chains. The number of amino acid substitutions that occurred between two species is plotted against the time of divergence of the species deduced from the paleontological record. The vertical error bars for the older divergence times extend over  $\pm 2$  standard deviations, indicating the 95% confidence level. Triangles represent hemoglobin a chains, circles hemoglobin b, and squares myoglobin ([1](#)).



The molecular clock made a tremendous impact on the study of molecular evolution. In particular, the molecular clock guaranteed, to some extent, that the divergence time can be estimated only by sequence comparisons. This estimation can be made even between two species for which their morphological characters are totally different and no fossil records exist. Indeed, immediately after the discovery of the molecular clock, amino acid sequencing was used extensively for the study of long-term evolution of organisms.

The speed of the molecular clock depends on the gene of interest. For example, the rate of amino acid substitutions for [histone](#) genes is more than 1000 times slower than that for fibrinopeptides (see [Fibrinogen](#)). In principle, the stronger the functional constraints, the slower the rate of amino acid substitution. Moreover, it has become recently known that the speed of molecular clock sometimes



depends on lineages (3). Thus, a molecular clock should be used with caution, particularly when estimating the divergence time.

## Bibliography

1. M. Nei, (1987) *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.
2. E. Zuckerkandl and L. Pauling (1962) in *Horizons in Biochemistry*, M. Kasha and B Pullman, eds., Academic Press: New York.
3. S. Easteal, C. C. Collet, and D. J. Betty, (1995) *The Mammalian Molecular Clock*, Springer-Verlag and R. D. Landes, Austin, TX.

## Molecular Dynamics

[Molecular dynamics](#) (MD) simulations evaluate the motion of the atoms in a given system and provide the position or *trajectory* of these atoms as a function of time. The trajectories are calculated by solving the classical equation of motion for the molecule under consideration. This is not unlike the well-known approach by which one evaluates the speed and position of a projectile, starting from the initial velocity, the mass, and the forces by using Newton's equation of motion. However, in the case of molecules, one obtains the relevant forces on each atom from the first derivatives of the given [potential functions](#). The actual evaluation of classical trajectories is done numerically, expressing the changes in coordinates and velocities at a time increment,  $\Delta t$ , by:

$$\begin{aligned} \mathbf{r}_i(t + \Delta t) &= \mathbf{r}_i(t) + \dot{\mathbf{r}}_i \Delta t \\ \dot{\mathbf{r}}_i(t + \Delta t) &= \dot{\mathbf{r}}_i(t) + \ddot{\mathbf{r}}_i \Delta t = \dot{\mathbf{r}}_i(t) - [(\partial U / \partial \mathbf{r}_i) \mathbf{m}_i^{-1}] \Delta t \end{aligned} \quad (1)$$

where the dot designates a time derivative and we use Newton's law:

$$\mathbf{m}_i \ddot{\mathbf{r}}_i = \mathbf{F}_i = -\partial U / \partial \mathbf{r}_i \quad (2)$$

starting with a given set of initial conditions [eg, with the values of  $\mathbf{r}_i(t = 0)$  and  $\dot{\mathbf{r}}_i(t = 0)$ ], we can evaluate  $\mathbf{r}(t)$  either by numerically integrating eq. (1) or by using the somewhat more complicated but far better approximation (1).

$$\begin{aligned} \mathbf{r}_i(t + \Delta t) &= \mathbf{r}_i(t) + \dot{\mathbf{r}}_i \Delta t + [4\ddot{\mathbf{r}}(t) - \ddot{\mathbf{r}}(t - \Delta t)] \Delta t^2 / 6 \\ \dot{\mathbf{r}}_i(t + \Delta t) &= \dot{\mathbf{r}}_i + [2\ddot{\mathbf{r}}(t + \Delta t) + 5\ddot{\mathbf{r}}(t) - \ddot{\mathbf{r}}[t - \Delta t]] \Delta t / 6 \end{aligned} \quad (3)$$

This equation allows one to obtain much more accurate results than those of eq. (1) using the same  $\Delta t$ s.

The propagation of classical trajectories of the atoms of a given system corresponds to a fixed total energy (determined by the specified initial conditions). However, the evaluation of statistical mechanical averages implies that the system included in the simulation is a part of a much larger system (ensemble) whose atoms are not considered in an explicit way. Thus, in order to simulate a given macroscopic property at a specified temperature we must introduce some type of “thermostat”

in the system that will keep it at the given temperature. This can be easily accomplished by assuming equal partition of kinetic energy among all degrees of freedom. Since each atom has three degrees of freedom with kinetic energy of  $\frac{1}{2}m\dot{r}^2 = \frac{3}{2}k_B T$  (where  $k_B$  is the Boltzmann constant) we obtain:

$$T = \sum_i m_i \dot{r}_i^2 / (3nk_B) \quad (4)$$

where  $n$  is the number of atoms in our system. In general, we can adjust the temperature during the simulation by scaling the velocities. That is, when  $T$  is smaller than the target temperature we can scale  $\dot{r}$  uniformly by  $1 + \epsilon$  until the target temperature is obtained. If  $T$  is higher than the target temperature, then a scaling of  $1 - \epsilon$  is used. More sophisticated considerations for constant temperature simulations are described elsewhere (1).

MD simulations provide a powerful tool for studies of time-dependent dynamical properties and for evaluation of average properties (see below). The strength of such approaches is associated with the fact that they have the ability to simulate, at least in principle, the true microscopic behaviors of macromolecules. The weakness is associated with the fact that some properties reflect extremely long time processes that cannot be simulated by any current computer.

Early MD simulations of a condensed phase system were performed by Alder and Wainwright (2). This approach was later extended to studies of the properties of liquid water (3). The emergence of MD simulations in studies of biological systems can be traced to a simulation of the dynamics of the primary event in the visual process (4) that correctly predicted a photoisomerization process of around 100 femtoseconds. A subsequent study (5) attempted to examine the **heat capacity** of [BPTI](#) by a very short simulation of this protein in vacuum. Unfortunately, at the early stages of the development of this field, it was not possible to obtain meaningful results for average properties of macromolecules due to the need for much stronger computers to reach a reasonable convergence (the heat capacity is drastically underestimated [6], reflecting artificial diffusive motions). Nevertheless, ultrafast reactions, such as those that control the photobiological process could be simulated even at this early stage (4). Eventually, with the increase of computer power, it has become feasible to reach simulation times of nanoseconds and to start to obtain meaningful average properties of macromolecules.

It should be mentioned, however, that it is frequently possible to evaluate properties that occur in a very long time using rather short simulation times (eg, 100 ps). Examples are rate constants that can be evaluated by calculating *activation-free energies*.

The time needed to reach accurate results for average properties depends on the model used and number of local minima. For example, models that trim the protein to a sphere with the proper spherical boundary conditions (7) converge much faster than models that involve periodic boundary conditions (8) because the latter involves more molecules and more minima. As much as accuracy is concerned, the proper treatment of long-range effects is crucial, and improved convergence is usually associated with more proper treatment of long-range forces (7, 9).

MD simulation methods provide a powerful way of evaluating average properties such as **free energies**, but it must be emphasized that such properties have little to do with dynamics per se and can be evaluated by other averaging approaches such as [Monte Carlo calculations](#). Similarly, the most important factors that determine the rate constants of most biological processes do not reflect dynamical properties but rather the probability of reaching the [transition state](#) configuration (10). Nevertheless, in cases of light induced ultra-fast photobiological processes, there are probably important effects that can be considered as dynamical properties, and MD simulations provide a direct way of modeling such effects.

Despite ever-increasing computational power, which has pushed the upper limit of simulation time scales to several hundred nanoseconds for condensed-phase systems, there are many interesting

processes that are too slow to be meaningfully simulated by MD approaches. However, some of these slow processes can be handled by **Brownian Dynamics (BD)** simulations, which represent the motion of the particles as a stochastic (random) process governed by a modified Langevin equation:

$$m_i \ddot{\mathbf{r}}_i = -\gamma m_i \dot{\mathbf{r}}_i + \mathbf{f}_i(t) + \mathbf{R}_i(t) \quad (5)$$

where  $\mathbf{R}_i(t)$  is the random force on  $i^{\text{th}}$  particle,  $\gamma$  is the friction coefficient representing the dissipation of the random forces, and  $\mathbf{f}_i(t)$  is the part of the inter-particle forces exerted on particle  $i$  that are not represented by the random-force and friction terms. The average magnitude of the random forces and the corresponding friction acting on them are connected through the fluctuation-dissipation theorem, with the relationship:

$$\langle \mathbf{R}^2 \rangle = n2mkT\gamma \quad (6)$$

where  $n$  is the dimensionality of the system. The very long time behavior of the Langevin equation reduces to Einstein's diffusion law. The resulting equation of motion

$$\mathbf{r}_i(t + \delta t) = \mathbf{r}_i(t) + \frac{\mathbf{f}_i(t)}{m_i \gamma} \delta t + \boldsymbol{\rho}_i(\delta t) \quad (7)$$

no longer involves velocities and represents a coarse-grained view of the time evolution of the system.  $\boldsymbol{\rho}_i$  represents the random displacement during a time step of duration  $\delta t$ , and is related to the friction coefficient by:

$$\langle \rho^2 \rangle = \frac{n2kT}{m\gamma} \delta t \quad (8)$$

In addition to simulating slow processes (11), BD simulation strategies can also be used for simplifying the dynamics in a system in order to study general aspects of dynamics (12) or testing theories that often contain assumptions similar to those used in obtaining the Brownian equations of motion (12).

## Bibliography

1. M. P. Allen and D. J. Tildesley (1987) *Computer Simulation of Liquids*, Oxford University Press, Oxford, UK.
2. B. J. Alder and T. E. Wainwright (1957) *J. Chem. Phys.* **27**, 1208–1209.
3. F. H. Stillinger and A. Rahman (1974) *J. Chem. Phys.* **60**, 1545.
4. A. Warshel (1976) *Nature* **260**, 679–683.
5. J. A. McCammon, B. R. Gelin, and M. Karplus (1977) *Nature* **267**, 585.
6. M. Levitt (1983) *J. Mol. Biol.* **168**, 621–657.
7. J. Åqvist (1996) *J. Comp. Chem.* **17**, 1587.
8. W. F. vanGunsteren and A. E. Mark (1992) *Eur. J. Biochem.* **204**, 947.
9. F. S. Lee, Z. T. Chu, and A. Warshel (1993) *J. Comp. Chem.* **14**, 161–185.
10. A. Warshel and W. W. Parson *Quart. Rev. Biophys.* In Press.
11. T. Schlick (1995) *Curr. Opin. Struct. Biol.* **5**, 245–262.
12. A. Papazyan and M. Maroncelli (1995) *J. Chem. Phys.* **102**, 2888–2919.

## Molecular Mechanics

Examination of the properties of molecules by **computer simulation** using the corresponding [potential functions](#) is called molecular mechanics (MM). This name reflects the fact that a molecular force field considers a molecule as a collection of balls connected by springs and that examination of the mechanical properties of such a system is similar to the study of the properties of the corresponding molecule.

The term *molecular mechanics* usually includes several techniques that are aimed at determining different molecular properties. In particular, with a given set of analytical potential functions, one can evaluate the molecular equilibrium geometries and the vibrations around these configurations. This task can be accomplished in the simplest way using the Cartesian representation, with coordinates  $(x, y, z)$ . That is, the potential surface for a molecule with  $n$  atoms can be expanded formally around the equilibrium configuration  $\mathbf{r}_0$  and give

$$U(\mathbf{r}_0 + \delta\mathbf{r}) = U(\mathbf{r}_0) + \sum_{i\alpha} (\partial U / \partial r_{i\alpha}) \delta r_{i\alpha} + \sum_{i\alpha, j\beta} (\partial^2 U / \partial r_{i\alpha} \partial r_{j\beta}) \delta r_{i\alpha} \delta r_{j\beta} + \dots \quad (1)$$

where the indices  $i$  and  $j$  designate atoms while  $a$  and  $b$  run over the  $x, y,$  and  $z$  coordinates of each atom. The first term is just the energy of the molecule at the equilibrium geometry. The second and third terms can be used (see below) to evaluate the equilibrium geometry and the vibrational frequencies.

### 1. Energy Minimization

The first term in equation (1) is just the energy of the system at equilibrium. The second term represents the deviation from equilibrium and the set of  $3n$  equations (for  $i = 1, 2, \dots, n$  and  $a = x, y, z$ )

$$\partial U / \partial r_{i\alpha} = 0 \quad (2)$$

represents the condition that  $\mathbf{r}_0$  is an equilibrium configuration. This set of equations can be solved approximately by the steepest-descent method (1), where the minimum is being searched by simply going in the opposite direction to the first derivative vector. That is, this method obtains the  $n$ th step toward the minimum by

$$\mathbf{r}_{i\alpha}^n = \mathbf{r}_{i\alpha}^{n-1} - g(\partial U / \partial r_{i\alpha}) \quad (3)$$

where  $g$  is a scaling factor that is changed in an iterative way to prevent  $\mathbf{r}$  from overshooting the minimum.

The steepest-descent method converges very slowly and is not so effective in searching for minima. A much more reliable and efficient approach is the modified Newton–Raphson method (1), (2). This method is based on expanding the gradient as a Taylor series around the given  $\mathbf{r}$  and finding the  $d\mathbf{r}$  that leads to  $\mathbf{r}_0$  where the gradient is zero (ie,  $\mathbf{r}_0 = \mathbf{r} + d\mathbf{r}$ ). This gives

$$\partial U(\mathbf{r}_0) / \partial r_{i\alpha} = \partial U(\mathbf{r}_0 + \delta\mathbf{r}) / \partial r_{i\alpha} + \sum_{j\beta} F_{i\alpha, j\beta} \delta r_{j\beta} = 0 \quad (4)$$

where  $\mathbf{F}$  is the matrix of second derivatives, that is,  $F_{ia,jb} = \partial^2 U / \partial r_{ia} \partial r_{jb}$ . Then, by solving equation (4) one obtains

$$\mathbf{r}_0 = \mathbf{r} + \delta\mathbf{r} = \mathbf{r} - \mathbf{F}^+ \nabla U(\mathbf{r}) \quad (5)$$

where  $\tilde{\partial}U$  is the gradient vector [with  $(\tilde{\partial}U)_{ia} = \partial U / \partial r_{ia}$ ] and  $\mathbf{F}^+$  is the generalized inverse of  $\mathbf{F}$ , which is constructed by “filtering” the zero eigenvalues of  $\mathbf{F}$  before inverting this matrix.

It is instructive to note that both the steepest-descent and the Newton–Raphson methods lead in the direction of  $-\tilde{\partial}U$ ; however, the steepest-descent method is unable to tell us how far to go in each step, and therefore we must search for the minimum in a very ineffective way. Equation (5) requires the evaluation of the second-derivative matrix  $\mathbf{F}$ , which is quite involved. Alternatively, one can use the conjugated gradient methods, where an approximation of  $b\mathbf{F}b^+$  is being built while searching the minimum using only the first derivatives vector  $DU$  (for a description of these powerful methods and related approaches see Ref. 1).

Energy-minimization methods that exploit information about the second derivative of the potential are quite effective in the structural refinement of proteins. That is, in the process of determination of a structure by [X-ray crystallography](#), one sometimes obtains bad steric interactions that can easily be relaxed by a small number of energy minimization cycles. In fact, one can combine the potential  $U(\mathbf{r})$  with the function that is usually optimized in X-ray structure determination and minimize the sum of these functions (3) by a conjugated gradient method, thereby satisfying both the X-ray electron density constraints and steric constraints dictated by the molecular potential surface.

Although the conjugated gradient and related methods are very effective in finding local energy minima, they do not overcome the problems associated with the enormous dimensionality of macromolecules. That is, in systems with many degrees of freedom, we expect to find a very large number of local minima; and it is not clear, for example, how to find the lowest minimum in an efficient way. For this purpose, one must use much more computer time and different types of search procedures (**Monte Carlo** or [molecular dynamics](#)), generating different configurations at different regions of conformation space and locating the minimum energy structure of each region. Such methods were, however, useless until the emergence of supercomputers in the early 1980s, and even at the present time they do not allow for a complete search of a protein's conformational space.

## 2. Normal Mode Analysis of Large Molecules

Finding a local energy minimum by a convergent minimization method allows one to exploit the third term in equation (1) for evaluation of the molecular vibrations around that minimum. Although early works used internal coordinates in studies of molecular vibrations (4), it was realized that much more effectiveness is offered by using the Cartesian coordinates of the given molecule (2). All that is needed in this case is the second-derivative matrix and the atomic masses. That is, the potential for infinitesimal vibrations around the local minimum  $\mathbf{r}_0$  can be written as

$$\delta U(\mathbf{r}_0) = \frac{1}{2} \delta\mathbf{r}' \mathbf{F}(\mathbf{r}_0) \delta\mathbf{r} \quad (6)$$

where  $\mathbf{F}$  is the second-derivative matrix of equation (1). The kinetic energy for the system can be written as

$$\delta T = \frac{1}{2} \delta\dot{\mathbf{r}}' \mathbf{M} \delta\dot{\mathbf{r}} = \frac{1}{2} \sum_{i,\alpha} m_i \delta\dot{r}_{i,\alpha}^2 \quad (7)$$

where  $\mathbf{dr}$  is the vector of the atomic velocities,  $\mathbf{M}$  is the diagonal matrix of the atomic masses ( $M_{i_a,j_b} = M_i d_{ij} d_{ab}$ ) and  $\mathbf{r} = d\mathbf{r}/dt$  ( $t$  is the time). The expression for the potential energy can be greatly simplified if we transform the Cartesian coordinates to a new set of coordinates called **normal coordinates**,  $\mathbf{Q}_i$ , using

$$\delta\mathbf{r} = \mathbf{M}^{-1/2}\mathbf{L}\mathbf{Q} \quad (8)$$

where the matrix  $\mathbf{L}$  is constructed from the column vectors  $\mathbf{L}^s$  that are obtained by diagonalizing the mass-weighted  $\mathbf{F}$  matrix

$$(\mathbf{M}^{-1/2}\mathbf{F}\mathbf{M}^{-1/2})\mathbf{L}^s = (2\pi\nu_s)^2\mathbf{L}^s \quad (9)$$

Here  $\nu_s$  is the vibrational frequency of the  $s$ th molecular mode and  $\mathbf{L}^s$  is the normal mode that describes the motion of the atoms in this particular mode.

Normal mode analysis based on analytical potential functions has been found to be quite effective in studies of spectroscopic properties of various organic molecules and of molecules of biological interest (eg, Refs. [5-7](#)), as well as in refinements of force fields (eg, Ref. [2](#)). The harmonic normal mode description was also found to be quite useful for approximate evaluation of various molecular properties. For example, one can use this description in a convenient way to evaluate the average thermal atomic motions. Normal mode analysis of protein motion is also quite useful ([8](#), [9](#)).

#### Bibliography

1. R. Fletcher (1980) *Practical Methods of Optimization*, Wiley, New York.
2. S. Lifson and A. Warshel (1968) *J. Chem. Phys.* **49**, 5116.
3. A. Jack and M. Levitt (1978) *Acta Crystallogr.* **A34**, 931.
4. E. B. Wilson, J. C. Decius, and P. C. Cross (1955) *Molecular Vibrations*, McGraw-Hill, New York.
5. G. Eyring, B. Curry, R. Mathies, R. Fransen, I. Palings, and J. Lugtenburg (1980) *Biochemistry* **19**, 2410.
6. A. Warshel (1977) *Annu. Rev. Biophys. Bioeng.* **6**, 273.
7. X.-L. Li, R. S. Czernuszewicz, J. R. Kincaid, and T. G. Spiro (1989) *J. Am. Chem. Soc.* **111**, 7012.
8. M. Levitt, C. Sander, and P. S. Stern (1983) *Int. J. Quant. Chem.* **10**, 181.
9. N. Go (1990) *Biophys. Chem.* **35**, 105.

#### Molecular Replacement

If the structure of a macromolecule has been determined by [X-ray crystallography](#), it can be used to elucidate a [homologous](#) structure without the need to determine the phases of the X-ray reflections directly by methods, such as [isomorphous replacement](#) (see [Phase Problem](#)). Proteins that are homologous in their amino acid sequence have essentially the same [tertiary structures](#). Therefore, the known structure can be employed as a first model for the unknown structure.

The problem is where to put this model in the unit cell of the unknown structure. This problem is solved in two steps:

1. The orientation of the model is found by applying a rotation function.
2. It is moved translationally to the correct position using a translation function.

The method was pioneered by Rossmann and Blow (1) based on principles given by Hoppe (2). The rotational and translational processes can most easily be understood by regarding the Patterson function (see [Patterson Map](#)). With this function a vector map is calculated. Vectors between atoms in the real structure show up as vectors from the origin to maxima in the Patterson map. Most of the intramolecular vectors are short, and their end points are found in a region around the origin of the Patterson map. If intermolecular vectors were completely absent, this inner region in the Patterson map would be the same for the same protein in different crystal structures and very similar for homologous structures, apart from a rotational difference. Although the intermolecular vectors disturb the situation somewhat, in many cases the required rotation can easily be obtained by finding the maximum overlap of the observed Patterson map for the unknown structure and that calculated for the known model.

With the translation function, the positions of the model molecules in the unit cell with respect to symmetry axes are determined. The principle here is also that the calculated and observed Patterson functions should overlap to a maximum extent. An excellent computer program for molecular replacement, AMoRe, has been written by Navaza (3).

After successful application of the rotation and translation function, the preliminary structure obtained is the starting model for a refinement process.

#### Bibliography

1. M. G. Rossmann and D. M. Blow (1962) *Acta Crystallogr.* **15**, 24–31.
2. W. Z. Hoppe (1957) *Elektrochemie* **61**, 1076–1083.
3. J. Navaza (1994) *Acta Crystallogr.* **A50**, 157–163.

#### Suggestion for Further Reading

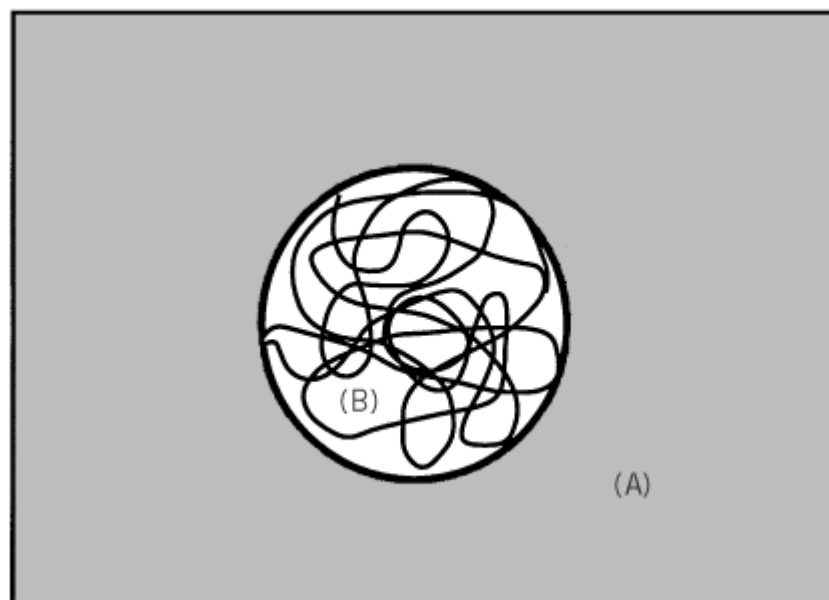
4. M. G. Rossmann (ed.) (1972) *The Molecular Replacement Method*, Gordon and Breach, New York.

## Molecular Sieve Resins

Molecular sieves are the materials used for preparing the stationary phase used in [size exclusion chromatography](#) (SEC), or gel filtration, to obtain information about the [hydrodynamic volume](#) and [Stokes Radius](#) of a [macromolecule](#). Figure 1 illustrates the essential features of molecular sieves (1). These resins are porous beads that are constructed from noninteracting or “inert” polymeric materials. When transported through a mobile phase in a column packed with this resin, a solute potentially has access to all the pores in the resin bead. The sizes of the pores in the beads are, however, not uniform but are best described as a distribution. If very small, a solute may have access to all pore sizes. A very large solute, in contrast, would not be accommodated in the pores at all and would be totally excluded. Solutes of intermediate size would partially partition into the pores of the molecular sieve. The extent of partitioning of the solute into the pores is inversely proportional to its size, with larger solutes partitioning to a lesser extent than smaller solutes. Here the properties of available molecular sieve resins are discussed, as well as phenomenological and mathematical

descriptions of the function of these resins.

**Figure 1.** Schematic representation of a molecular sieve partitioning system. The dotted area represents the mobile-phase exterior to the resin particle. The porous gel particle enclosed by the circle represents the stationary phase. The gel particle consists of two regions, the internal solvent ( **B** ) and the polymeric gel matrix or curved lines ( **A** ). Solute molecules introduced into the system partition between ( **A** ) and ( **B** ), with the extent of partitioning into ( **B** ) depending on the size of the solute.



Molecular sieve resins are prepared from a range of materials. The two major classes of stationary phase materials are derivatized glass or silica and crosslinked gels. The crosslinked gels are further subdivided into natural or synthetic; the former are based on polysaccharides and the latter on polymerized hydrophilic monomers. Advantages of the silica-based packings include their rigidity, uniform pore size, and uniform particle size (2). A summary of the types of available derivatized silica resins is provided in Table 1 of Reference 2. The major disadvantages of the derivatized silica resins are their adsorptive properties and tendency to degrade with use. The first problem results from incomplete reaction of the SiOH groups, which results in a negatively charged surface with which proteins have been found to interact (2-4). The problem with degradation of the derivatized silica based resins limits the composition of the mobile phase. The common crosslinked polymers include crosslinked dextrans or [agarose](#), the [polyacrylamide](#) resins, and composite polyacrylamide/agarose resins (2). Table 2 of Reference 2 contains a summary of the types and properties of a large number of these hydrophilic polymer gels. An inherent assumption of size exclusion chromatography is that the solute that is chromatographed does not interact with the resin. As indicated in the preceding discussion, silica-based resins are the least desirable in this respect. Polymeric resins are, however, by no means inert, and care should be exercised in choosing the composition of the mobile phase used in chromatographic experiments, in order to minimize interaction with the resin (2).

In a molecular sieving process, solute molecules partition between the solvent spaces or the pores in the resin particles (the stationary phase in the chromatographic column) and the solvent space exterior to the resin particles (the mobile phase in the chromatographic experiment). The stationary phase consists of the sieving resin equilibrated with the solvent used for the mobile phase. There are three compartments within the experimental system (1): (1) the volume exterior to the resin particles, or the *void volume*,  $V_0$ ; (2) the solvent region in the interior of the gel particles, or the *internal*



volume,  $V_i$ ; and finally (3) the volume occupied by the gel matrix,  $V_g$ . The experimental methods for measurement of these volumes are provided in [size exclusion chromatography](#). The total volume of the column is simply the sum of these three volumes:

$$V_t = V_o + V_i + V_g \quad (1)$$

The behavior of a solute, typically a protein, in gel filtration chromatography is expressed in terms of the extent of partitioning of the solute into the pores of the resin, which is measured quantitatively by the [partition coefficient](#),  $s$ . It is defined as the amount of solute distributed in the gel per unit internal volume,  $V_i$ , and solute concentration in the external space of the column,  $C$ . A partition isotherm that relates the weight of solute,  $Q_i$ , inside the gel and the solute concentration in the space exterior to the gel, or the void, can be written in terms of  $s$ .

$$Q_i = \sigma V_i C \quad (2)$$

The partition coefficient is thus a measure of the degree of partitioning of the solute into the interior space of the gel matrix. At equilibrium, the penetrable volume,  $V_p$ , of the gel that is occupied by solute molecules is given by

$$V_p = Q_i / C_p \quad (3)$$

where  $C_p$  is the solute concentration in the region of distribution. Combining Equations (2) and (3) leads to

$$\sigma = V_p / V_i \quad (4)$$

which indicates that the partition coefficient is the fraction of the internal volume of the resin that is penetrable by the solute molecule. This relationship works reasonably well at low solute concentrations. At high solute concentrations, however, corrections must be introduced to take into account the effects of nonideality. Values of  $s$  range from 0, for very large solutes that are totally excluded from the pores, to 1 for very small solutes that are completely included in the pores.

The partition coefficient can also be related to the total volume of the gel phase  $V_t - V_o$ . Using this formalism the partition isotherm is expressed as

$$Q_i = K_{av}(V_t - V_o)C \quad (5)$$

Assuming no concentration-dependence of the partition coefficient or nonideality:

$$K_{av} = \frac{V_p}{V_t - V_o} \quad (6)$$

At low solute concentrations,  $K_{av}$  is the volume fraction of the stationary phase that is occupied by the solute. This coefficient is frequently used in analysis of SEC data. The coefficient  $K_{av}$  is related to  $s$  by the following expression:

$$\sigma = K_{av}(1 + \bar{V}_g / S_r)$$

where  $ov/Vovl_g$  is the partial specific volume of the gel matrix and  $S_r$  is the solvent regain of the gel material. For a detailed discussion of these two terms, the reader is referred to Reference 1. Note that if the solvent regain and [partial specific volumes](#) for two gel materials are known, a measured value for a partition coefficient obtained in one system can be converted to the value in another. Moreover, for any sieving resin, a constant ratio exists between the two partition coefficients,  $s$  and  $K_{av}$ .

### Bibliography

1. G. K. Ackers (1975) In *The Proteins*, H. Neurath and R. L. Hill, eds., Academic Press, New York, pp. 1–94.
2. P. L. Dubin (1992) In *Advances in Chromatography*, J. C. Giddings, E. Grushka, and P. R. Brown, eds., Marcel Dekker, New York, Vol. 31, pp. 119–151.
3. K. M. Gooding and F. E. Regnier (1990) In *HPLC of Biological Macromolecules*, K. M. Gooding and F. E. Regnier, eds., Marcel Dekker, New York, pp. 47–75.
4. Y. Kato (1989) In *Size Exclusion Chromatography*, B. J. Hunt and S. R. Holding, eds., Chapman & Hall, New York, pp. 170–188.

### Suggestion for Further Reading

5. H. G. Barth, B. E. Boyes, and C. Jackson (1996) *Anal. Chem.* **68**, 445R–466R. (A comprehensive review of recent advances in SEC methodologies and of recent developments in preparation of sieving resins.)

## Molecular Surface, Volume

[Protein structure](#) can be described with the spatial coordinates of all the constitutive atoms in terms of [alpha-helices](#), [beta-sheet](#), [turns](#), irregular regions, etc. But this useful information has not yet led to a complete understanding of protein structure, particularly its dynamic nature. The efficiency of packing of atoms and residues in a protein molecule (the packing density) has important implications for its structure and its thermodynamic, mechanical, and functional properties. Related structural parameters are the surface area and the volume, which make it possible to grasp intuitively the overall structure of a protein. It is difficult, however, to define the complex surface and volume of all its atoms, including all the intramolecular cavities, so there are a number of definitions of surface and volume for protein molecules.

The geometrical surface and volume of the individual atoms of a macromolecule (the [van der Waals surface, volume](#)) can be evaluated on the basis of the structure determined by [X-ray crystallography](#). The surface and volume of this collection of atoms can be computed, but many atoms are in the interior and will normally not make contact with a molecule in the solvent. A more relevant surface is that which normally is in contact with the solvent. This is defined by rolling a spherical probe, representing a water molecule, over the van der Waals surface of the macromolecule (see Fig. 1 of [Accessible Surface](#)). Those parts of the van der Waals surface in contact with the surface of the solvent molecule are designated the contact surface. When the probe is simultaneously in contact with more than one protein atom, its interior surface defines the reentrant surface. The contact surface and the reentrant surface together make a continuous surface, which is defined as the molecular surface, and it defines a molecular volume. The surface defined by the center of the probe molecule is the [accessible surface](#).

Special thermodynamic volumes, such as the [partial specific volume](#), are introduced for understanding the physicochemical properties of proteins in solution. The pressure and temperature coefficients of the volume ([compressibility](#) and [expansibility](#)) also give useful information on the atomic packing and flexibility of protein molecule.

#### Suggestion for Further Reading

F. M. Richards (1977) Areas, volumes, packing, and protein structure. *Ann. Rev. Biophys. Bioeng.* **6**, 151–176.

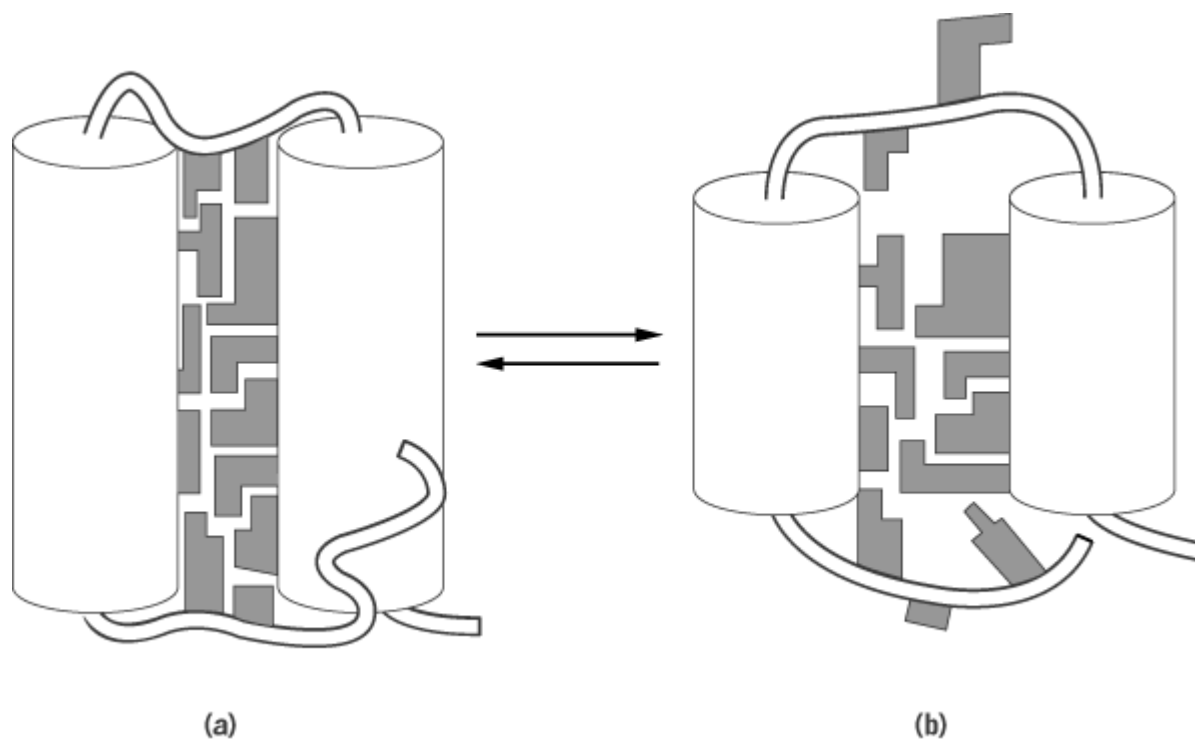
## Molten Globule

The molten globule state is an intermediate conformational state between the native and the fully unfolded states of a globular protein (see **Protein folding** and **Protein unfolding**). Many proteins can be observed in this state when partially unfolded at equilibrium, under mild **denaturation** conditions, or as a transient intermediate [kinetic](#) species, being formed rapidly from the unfolded state upon transfer to refolding conditions. The characteristics of the molten globule state are:

1. the presence of a native-like content of **secondary structure**;
2. the absence of a specific [tertiary structure](#) produced by the tight packing of amino acid side chains;
3. compactness in the overall shape of the protein molecule, with a radius 10 to 30% larger than that of the native state;
4. the presence of a loosely packed **hydrophobic** core that increases the hydrophobic surface area accessible to solvent.

Thus, in short, the molten globule is a compact globule with a “molten” side-chain structure that is primarily stabilized by nonspecific hydrophobic interactions (Fig. [1](#)).

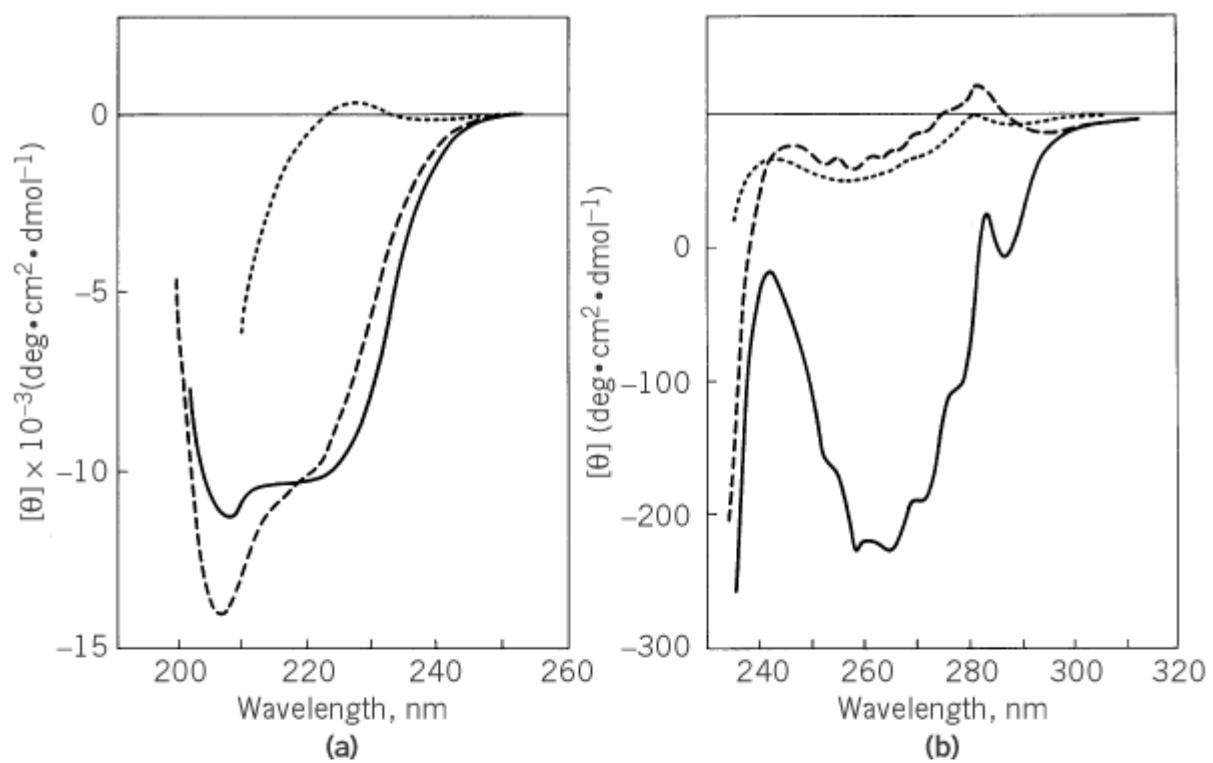
**Figure 1.** A schematic model of the native (**a**) and the molten globule (**b**) states of a protein molecule. Only two **a-helices** are presented for simplicity. According to this model, the molten globule preserves the mean overall structural features of the native protein but differs from the native state mainly by looser packing and a higher mobility for the loops and ends of the protein molecule. (From Ref. [13](#), reprinted with permission.)



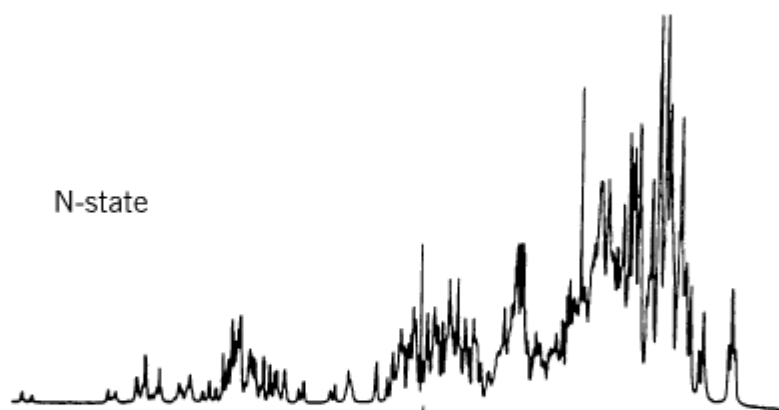
Experimentally, the molten globule state is characterized by having a native-like [circular dichroism](#) (CD) peptide spectrum below 250 nm, which arises from the secondary structure, and by an unfolded-like CD spectrum of the aromatic side chains between 250 and 320 nm, due to the absence of specific side-chain packing interactions (Fig. 2). Hydrogen atoms of the peptide backbone involved in secondary structure of the molten globule state appear to be protected from [hydrogen exchange](#) with the solvent protons, but the protection factor for the molten globule (10 to 1000) is much smaller than that for the native state, which is often greater than  $10^6$ . The nuclear magnetic resonance (NMR) spectrum of the molten globule state is closer to that of the unfolded protein, and there is little, if any, chemical shift dispersion in the spectrum, reflecting the absence of a specific tertiary structure (Fig. 2). The individual resonances in the NMR spectrum are, however, broader than those in the unfolded state, reflecting conformational fluctuations in the molten globule state. It is known that the structure of the molten globule state, as determined by hydrogen exchange and NMR spectra, is heterogeneous in proteins, including [a-lactalbumin](#), [cytochrome c](#), and apo [myoglobin](#). In these proteins, one portion of the structure is more organized and substantially protected in the molten globule state, with other portions of the structure being less organized. Solution [X-ray scattering](#) has been used to characterize the molten globule structure (1). The presence of a clear peak in the Kratky plot and the radius of gyration evaluated from the Guinier plot of the X-ray scattering curve in the molten globule state show that the protein molecule in this state is compact and globular. Limited **proteolysis** by proteolytic enzymes has also been used for probing the partly folded structures of proteins; the key result is that the molten globule can be sufficiently rigid to prevent extensive proteolysis and it appears to maintain significant native-like structure (2). A hydrophobic fluorescent dye, such as 8-anilino-naphthalene-1-sulfonate (ANS), which binds to solvent-accessible hydrophobic surfaces of a protein molecule, has also been used for characterizing the molten globule state. The molten globule binds ANS much more strongly than does either the fully folded or fully unfolded form of the protein; the latter can be generated in a concentrated solution of a strong denaturant (6 M **guanidinium** chloride or 8 M [urea](#)).

**Figure 2.** CD and NMR spectra of the native, the molten globule, and the fully unfolded states of a-lactalbumin. (a, b)

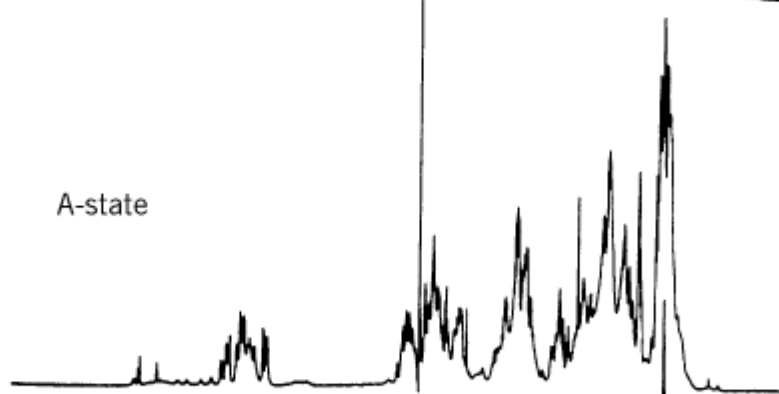
CD spectra in the peptide region (**a**) and in the aromatic region (**b**) of bovine  $\alpha$ -lactalbumin in the native state at pH 7.0 (solid line), the molten globule at pH 2.0 (dashed line), and the fully unfolded state in 6 M guanidinium chloride at pH 2.0 (dotted line), at 4.5°C. (From Ref. [14](#), reprinted with permission.) (**c,d,e**) 500-MHz  $^1\text{H}$  NMR spectra at 52°C of guinea pig  $\alpha$ -lactalbumin in the following states: (**c**) native (pH 5.4), (**d**) molten globule (pH 2.0), (**e**) fully unfolded (pH 2.0 in 9 M urea). (From Ref. [15](#), reprinted with permission.)



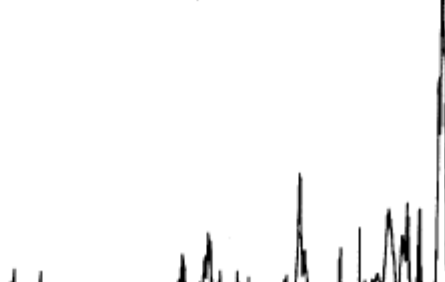
(c) N-state



(d) A-state



(e) U-state



For some proteins, the molten globule state is an equilibrium intermediate observed at an intermediate concentration of a strong denaturant (eg, 2 M guanidinium chloride) as part of a denaturant-induced unfolding transition. On the other hand, many globular proteins show a **cooperative** two-state unfolding transition without the intermediate. Whether or not the molten globule state is observed as a stable intermediate depends upon its stability relative to that of the native and the unfolded states (3). The unfolding intermediates of [carbonic anhydrase](#) and of [α-lactalbumin](#) are typical examples of a molten globule state that is stably populated at an intermediate concentration of denaturant. For these proteins, the partially unfolded states at acidic or alkaline pH are identical to the unfolding intermediate in the denaturant-induced transition, so the acidic or alkaline transitions also produce the molten globule state. For some other proteins, such as cytochrome *c*, apomyoglobin, and **b-lactamase**, the acidic or alkaline transition is known to produce a more extensively unfolded state; in these cases the addition of salt refolds the protein molecule from the unfolded to the molten globule state. The salt-induced refolding to the molten globule state is caused by counterion binding of the salt to the protein molecule, which eliminates the **electrostatic** repulsion between the charged groups. Other mildly denaturing processes that lead to the molten globule state include denaturation induced by hydrostatic pressure and by alcohols. Removal of the bound metal ion in a metal-ion binding protein sometimes results in a molten globule state, as in the case of apo-[α-lactalbumin](#) produced by removal of the bound  $\text{Ca}^{2+}$ . Covalent modification of a protein can also sometimes result in a molten globule state.

In many globular proteins, the molten globule state is observed at an early stage of the [kinetics](#) of refolding from the unfolded state. The early formation of the molten globule might be a way by which these proteins can be folded efficiently without wandering into the huge conformational space available for the proteins. Two experimental techniques, stopped-flow CD (4) and pulsed hydrogen-exchange combined with either two-dimensional NMR or electrospray ionization [mass spectrometry](#), have been used successfully to characterize the transiently formed molten globule-like states during the kinetics of refolding of many globular proteins. The stopped-flow CD studies have shown the rapid formation of the peptide secondary structure occurring within the dead-time of the stopped-flow mixing (~10ms), although how much of the secondary structure is rapidly regained depends on the protein species (3). The pulsed hydrogen-exchange technique, when combined with two-dimensional NMR, can identify the specific location of stabilized secondary structure segments in a transient intermediate of a protein. For apomyoglobin (5) and ribonuclease HI (6), comparison of the kinetic refolding intermediate and the equilibrium molten globule state has shown that the two are identical. Identification of the molten globule intermediate has also been well-established for [α-lactalbumin](#) by time-resolved CD and NMR studies (4, 7). Molten globule-like folding intermediates have also been detected and characterized in many other globular proteins. Nevertheless, this does not necessarily mean that the molten globule state must be an obligatory, universal intermediate of protein folding. Because the formation of the molten globule-like folding intermediate is usually too rapid to be coupled kinetically with the subsequent folding reactions, it is very difficult to determine whether or not the molten globule is an obligatory folding intermediate. Furthermore, several small globular proteins with approximately 60 amino acid residues are known to refold very rapidly to the native state within a few milliseconds without accumulation of the molten globule intermediate (8).

When considering the role of the molten globule state in protein folding, it is important to address the question as to whether or not the molten globule is a **thermodynamic** state. Analysis of the cooperativity parameters for denaturant-induced unfolding transitions of some proteins has suggested that the transitions from the molten globule state to the unfolded state and from the native state to the molten globule state are both all-or-none transitions, indicating that the molten globule state is a thermodynamic state. Furthermore, stability studies of mutants of apomyoglobin (9) and cytochrome *c* (10) have concluded that the molten globule states of these proteins show cooperative unfolding and are stabilized by native-like tertiary interactions, in addition to nonspecific hydrophobic

interactions, which is consistent with the proposal that the molten globule state is a distinct thermodynamic state. However, for the best-characterized molten globule of  $\alpha$ -lactalbumin, calorimetry, NMR, vibrational **Raman spectroscopy** (11), and other techniques (12) have clearly shown that the unfolding of this molten globule is not a cooperative two-state transition. Such diversity in the unfolding behavior of the molten globule state among different proteins may arise from the diversity of the molten globule structure. Because the native tertiary interactions are at least partially lost in the molten globule state, its structure must be more diverse than the native structure, and how cooperatively the molten globule unfolds may depend on how many residual native tertiary interactions are retained in this state. A good example is the molten globule state of the equine **lysozyme**, which is a **calcium-binding protein** and homologous to  $\alpha$ -lactalbumin. Although the equine lysozyme molten globule apparently resembles that of  $\alpha$ -lactalbumin, a rigorous analysis of its spectroscopic and thermodynamic properties has shown that its structure is significantly more highly organized and that its unfolding is a cooperative first-order transition accompanied by a large change in enthalpy. Because of this diversity in the intermediate conformational states of proteins, it is difficult to provide a clear structural definition of the molten globule state. Consequently, this causes some controversy, with the formation of native-like tertiary fold considered to be a characteristic of the molten globule state in certain cases, while in other cases structures with non-native tertiary folds are also molten globules. Furthermore, more than one intermediate conformational state is often observed between the native and fully unfolded states (10). Nevertheless, the four characteristics itemized at the beginning of this article are those generally accepted as the characteristics of the molten globule state.

It is now well established that not only the native state, but also non-native conformational states, play an important role in a biological cell. The protein states recognized by various **molecular chaperones** are non-native. The non-native conformation is also required for translocation of a protein across a biological membrane (see **Protein Secretion**). Various genetic diseases can be caused by the misfolding of translated polypeptides, and this misfolding results from an increased propensity of the mutant proteins to form non-native conformations. Because the molten globule state is regarded as a denatured state under physiological conditions, it definitely assumes some role in the above phenomena *in vivo*. It is also true, however, that there is a much greater diversity in the non-native conformations of proteins than in just the conformations characterized as the molten globule state.

## Bibliography

1. M. Kataoka and Y. Goto (1996) *Folding Design* **1**, R107–R114.
2. A. Fontana et al. (1997) *Folding Des.* **2**, R17–R26.
3. K. Kuwajima (1996) In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum, New York, pp. 159–182.
4. M. Arai and K. Kuwajima (1996) *Folding Des.* **1**, 275–287.
5. P. A. Jennings and P. E. Wright (1993) *Science* **262**, 892–896.
6. T. M. Raschke and S. Marqusee (1997) *Nat. Struct. Biol.* **4**, 298–304.
7. J. Balbach, V. Forge, W. S. Lau, N. A. J. Van Nuland, K. Brew, and C. M. Dobson (1996) *Science* **274**, 1161–1163.
8. T. Schindler and F. X. Schmid (1996) *Biochemistry* **35**, 16833–16842.
9. M. S. Kay and R. L. Baldwin (1996) *Nat. Struct. Biol.* **3**, 439–445.
10. W. Colón and H. Roder (1996) *Nat. Struct. Biol.* **3**, 1019–1025.
11. G. Wilson, L. Hecht, and L. D. Barron (1996) *J. Mol. Biol.* **261**, 341–347.
12. B. A. Schulman and P. S. Kim (1996) *Nat. Struct. Biol.* **3**, 682–687.
13. O. B. Ptitsyn (1992) In *Protein Folding* (T. E. Creighton, ed.), W. H. Freeman, New York, pp. 243–300.
14. K. Kuwajima, Y. Hiraoka, M. Ikeguchi, and S. Sugai (1985) *Biochemistry* **24**, 874–881.



15. J. Baum, C. M. Dobson, P. A. Evans, and C. Hanley (1989) *Biochemistry* **28**, 7–13.

### Suggestions for Further Reading

16. K. Kuwajima (1989) The molten globule state as a clue for understanding the folding and cooperativity of globular-protein structure. *Proteins* **6**, 87–103.
17. D. Barrick and R. L. Baldwin (1993) The molten globule intermediate of apomyoglobin and the process of protein folding. *Protein Sci.* **2**, 869–876.
18. H. Christensen and R. H. Pain (1994) The Contribution of the Molten Globule Model. In *Mechanisms of Protein Folding* (R. H. Pain, ed.), IRL Press, Oxford, pp. 55–79.
19. O. B. Ptitsyn (1995) Molten globule and protein folding. *Adv. Protein Chem.* **47**, 83–229. (A comprehensive review of the molten globule state with more than 450 references.)
20. O. B. Ptitsyn (1995) How the molten globule became. *Trends Biochem. Sci.* **20**, 376–379. (An interesting review of the history of the molten globule state from the perspective of Ptitsyn's group.)
21. A. L. Fink (1995) "Molten globules". In *Methods in Molecular Biology: Protein Stability and Folding* (B. A. Shirley, ed.) Humana Press, Totowa, NJ, pp. 343–360.
22. K. Kuwajima (1996) The molten globule state of  $\alpha$ -lactalbumin. *FASEB J.* **10**, 102–109.
23. A. D. Miranker and C. M. Dobson (1996) Collapse and cooperativity in protein folding. *Curr. Opin. Struct. Biol.* **6**, 31–42.
24. H. Roder and W. Colòn (1997) Kinetic role of early intermediates in protein folding. *Curr. Opin. Struct. Biol.* **7**, 15–28.
25. T. E. Creighton (1997) How important is the molten globule for correct protein folding? *Trends. Biochem. Sci.* **22**, 6–10.

## Monocentric Chromosome

[Chromosomes](#) with a single [centromere](#) are monocentric. The majority of metazoan animals and **plants** have monocentric chromosomes. Monocentric chromosomes are also called **acrocentric** if the centromere is visible at the end of the chromosome, **metacentric** if it is toward the middle, and **telocentric** if it is at the chromosome's end.

## Monoclonal Antibody

### 1. Monoclonal Antibodies

Monoclonal antibodies are of special significance in immunology because of their direct correlation with fundamental immunological principles at the cellular, genetic, and protein levels ([1](#)). Before the recognition and understanding of monoclonal antibodies, immunologists appreciated the biological benefits to the host of the extensive heterogeneity that characterized the mammalian polyclonal humoral immune response to most antigens. Heterogeneity equated with immunological diversity and represented a crucial component in a competent and protective immune system. Understanding antibody heterogeneity at the molecular level led to the hypothesis that heterogeneous immune

responses were comprised of arrays of antibodies of similar (but not identical) specificity and affinity for a given antigen. It was further surmised that each antibody was the product of a single lymphocyte, a fundamental doctrine in the Clonal Selection theory (2). Clonal selection had many biological implications that led to the following precepts: 1) fully differential lymphocytes would have a finite genetic capacity to synthesize antibody; 2) upon antigen stimulation, responding activated lymphocytes would rapidly proliferate, generating progeny-synthesizing identical antibody molecules; and 3) antibody molecules produced by a single cell and subsequent progeny would be homogeneous in specific interaction with antigen. Thus, the term “monoclonal antibody” embodies basic elements that collectively characterize the immune system. Experimental proof of these basic tenets established the foundation for modern molecular immunology.

Monoclonal antibody is by definition a homogeneous population of immunoglobulin macromolecules produced by a single lymphocytic cell and progeny. Homogeneity of antibody macromolecules can be assessed at various physiochemical levels to verify monoclonality.

First, every molecule in a monoclonal antibody population is of the same heavy (H) and light (L) chain isotype. Isotype is equated with the class of each chain (eg, IgG or IgM as H chain designations and kappa or lambda for L chains) reflecting identical primary structures in the constant domains. Thus, every lymphocyte synthesizing and secreting the Ig protein expresses identical gene products.

Second, when the variable domains associated with the H and L chains combine to form a functional active site of a given specificity in a monoclonal antibody population, all sites are idiotypically identical. Idiotypic refers to the shape, conformation, configuration, and dynamics of the antibody-active site constituted by the unique non-covalent association of the two variable domains. Idiotypic relationships are determined by eliciting anti-idiotypic antibodies specific for the unliganded active site of the monoclonal antibody. Similarly, when ligand or antigen is bound to the variable domains of a given idioform, the monoclonal antibodies assume the same unique liganded conformation as evidenced by identical metatypic properties. Epitopes associated with variable domain interactions that constitute idioform (non-liganded state) or metatype (liganded state) of a monoclonal antibody are non-linear or conformational in nature (3). Both idioform and metatype represent unique conformational properties of the active site of a monoclonal antibody and differ due to the presence or absence of antigen.

Third, homogeneous or monoclonal antibody populations display identical binding properties for the homologous antigen or ligand. Thus, in terms of affinity, the same binding constant ( $M^{-1}$ ) describes each active site within the monoclonal population. In contrast, an affinity of a polyclonal antibody population represents an average of the molecules present. In addition, identical specificity and cross-reactivities are exhibited by all molecules. Homogeneity of binding can be quantitatively measured by constructing antigen-antibody binding plots, such as Scatchard and Sips plots. In the Scatchard plot, as the binding sites become saturated, the abscissa intercept of the plot represents the valence (N) of the antibodies comprising the population (eg, N = 2.0 for IgG antibodies). In a monoclonal antibody population, the Scatchard plot will characteristically yield a straight line indicative of homogeneous reactivity with antigen. Importantly, the corresponding Sips plot yields a Sipsian heterogeneity index based on a scale ranging from 0 to 1.0. When the constant or index (slope of the line) is equivalent to 1.0, the antibody population is judged to be monoclonal or homogeneous.

Fourth, monoclonality can be ascertained by determining the amino acid sequence of the H and L chains of the antibody protein population. Since antibodies of differing specificities or affinities vary in their primary structure within the H and L chain variable domains, the ability to determine primary structures verifies population monoclonality. Heterogeneous antibody populations are not amenable to primary structure determinations reflecting significant differences.

Finally, it is instructive to mention a procedure to be used with reservation to ascertain monoclonality because the results are often misleading. In particular, isoelectric focusing yields aberrant results when analyzing monoclonal antibodies. Analyses of a monoclonal antibody population on the basis of isoelectric point determinations are often distorted by microheterogeneity that is attributable to the multiple carbohydrate moieties attached to the antibody molecules. Although the Fc-associated (ie, constant domains of H chains) carbohydrate content is of a relatively small percentage on a weight basis (eg, IgG is 2% to 3% carbohydrate, while IgM is 10% to 12% carbohydrate), the carbohydrate content varies within an antibody population that is homogeneous at the protein level. Microheterogeneity can be explained, in part, due to differential sialic acid content as the terminal sugars constituting the individual carbohydrate moieties. Microheterogeneity results in a spectrum of different isoelectric points, even in a monoclonal antibody population. Enzymatic or chemical removal of charged sialic acid residues converts heterogeneous isoelectric profiles to relatively more homogeneous patterns consistent with the uniform primary structure of the Ig protein.

When an antibody population satisfies the various criteria listed, the protein can be determined as monoclonal and the product of a single cell.

## 2. Myeloma Proteins

Before the development and universal application of hybridoma methodology to generate monoclonal antibodies of defined specificity, immunologists relied on the existence of homogeneous pathological immunoglobulins called myelomas. Myeloma proteins are produced spontaneously in humans due to transformation of a lymphocyte to the malignant state. Myeloma proteins can also be generated in specific strains of mice through administration of various reagents, such as mineral oil, administered in the peritoneal cavity. Initially, myeloma globulins produced during human and animal disease were considered devoid of functional properties and were classified along with normal immunoglobulins as non-specific and lacking antigen-binding capacity. With the advent of the previously cited Clonal Selection theory (2), homogeneous myeloma proteins become antibodies in search of antigen. Myeloma proteins in humans appear in a frequency that parallels their concentration in serum and, therefore, the respective number of lymphocytes synthesizing that class of antibody. Thus, human IgG myelomas appear most frequently (~90%) followed by IgM (10%), etc. Myelomas in the other classes, such as IgA, IgD, and IgE are quite rare in humans. In contrast, induced myelomas in mice are predominantly IgA followed in incidence by IgG and IgM. A predominance of IgA reflects the administration of mineral oil into the peritoneal cavity, which is rich in IgA-producing lymphocytes.

Myeloma proteins proved homogeneous at the protein level, and, because they were produced in high quantities, they were easily purified to homogeneity. These homogeneous proteins provided the foundation upon which the first primary structures of immunoglobulins were based. L chains of myeloma proteins in human patients exist in excess and are harvested from urine as monomers or stable dimers called Bence-Jones proteins. Since L chains were of lower molecular weight (approximately 22 kd to 23 kd) and lacked carbohydrate moieties, they were used to obtain the first primary structures. Putnan, Titani, and Whitley (4) published the first complete sequence analysis of a Bence-Jones protein called Ag, a sequence that became the standard by which all subsequent L chain primary structure analyses were compared. Roy was the second completed L chain sequence and was significant because it was the comparison of primary structures that yielded critical structure-function information. Comparison of several primary structures of H and L chains normalized for isotype revealed the existence of the variable and constant domains (5). Further comparisons within the variable domains showed that the H chain consisted of four hypervariable regions termed “complementarily determining regions” (CDRs). The L chain variable domains consisted of three CDR segments. This finding, combined with similar results from analyses of hybridoma proteins, led to delineation of the molecular basis of immunological specificity.

## 3. Hybridomas

Classical seminal studies by Köhler and Milstein in 1975 (6) led to the generation of monoclonal antibodies with defined specificity. Hybridoma methodology was based on fusing murine immune lymphocytes with a malignant mouse myeloma lymphocyte cell line that possessed certain important genetic and biochemical properties. This important methodology was necessitated by the fact that immune antibody-producing B lymphocytes were difficult to maintain as stable cell lines under tissue culture conditions. However, fusion of antibody-producing cells with a malignant myeloma cell line resulted in proliferative cells that can be maintained indefinitely.

Spontaneous fusions between lymphocytes are rare events, so, to generate large numbers of fused cells, the frequency of fusion was increased by using surface-active agents such as inactivated Sendai virus or the chemical polyethylene glycol. Upon mixing immune lymphocytes with myeloma cells and the fusing reagent, a selection system was mandatory to eliminate unfused parental cells and to assay for appropriately fused cells that synthesized antibodies with a desired specificity. The tissue culture selection process was dependent on the procurement of myeloma cell lines that no longer synthesized an immunoglobulin and lacked genes for thymidine kinase (TK) or hypoxanthine-guanine-phosphoribosyl transferase (HGPRT). These enzymes catalyzed important steps in the salvage pathway used in DNA synthesis by dividing cells. Unfused myeloma cells are unable to proliferate in the HAT selection medium (containing hypoxanthine, aminopterin, and thymidine) because aminopterin inhibited the *de novo* DNA biosynthetic pathway forcing the cells to utilize the non-functional salvage pathway. Fused myeloma cells proliferate in HAT medium because the immune lymphocytes (fusion partner) furnish the essential TK and HGPRT enzymes for the salvage pathway. Thus, cells that survived incubation with HAT medium represented a successful immune B cell-myeloma cell fusion usually on a 1:1 cellular basis. The second phase in hybridoma production was the selection of those fused cells that synthesized and secreted antibody of the desired specificity. This was generally achieved by diluting fused cell populations to a single cell and then generating clones by proliferation of the single cell. Extracellular fluid or medium from the single cell and resulting progeny was assayed for the appropriate antigen binding. Methods and technologies used to detect monoclonal antibody activity are usually based upon some variation of the solid-phase immunoassay (eg, ELISA or radioimmunoassay formats).

Hybridoma methodology and consistent generation of relatively large amounts of monoclonal antibodies with defined specificities has proven to be a source of important immunological reagents. First, the provision of adequate quantities of monoclonal antibody has facilitated the routine determination of primary structures of both the H and L chains at the protein level. Second, various monoclonal antibodies can be crystallized on a regular basis providing atomic resolutions of 2 to 3 Å by X-ray crystallographic procedures. In general, Fab fragments (four domain substructures of antibodies including the variable domains of both chains) with bound ligand crystallize most efficiently. Approximately 100 resolved structures of monoclonal antibodies (Fab fragments) have now been reported. Third, this methodology provides large homogeneous hybridoma cell populations from which mRNAs responsible for the synthesis of H and L chains can be obtained. Through development of polynucleotide primers, PCR technology can be used to generate cDNA copies for sequencing at the gene level. Fourth, gene cloning experimentation with antibody genes has been important in successful site-specific mutagenesis studies. Such studies have been critical in deciphering structure-function relationships governing antibody activity. Fifth, hybridoma methodology has led to the development of antibody derivatives such as single-chain or Fv antibodies (7). The latter represent ~25,000 structures containing only the variable domains of the H and L chains of a specific antibody molecule. Variable domains constituting single-chain antibodies are usually covalently attached through the use of flexible polylinkers (10 to 15 amino acids) encoded in the single-chain synthetic gene. Single-chain genes are subsequently incorporated into the appropriate plasmids and expressed in either procaryote or eucaryote cell lines. Sixth, antigen-antibody interactions can now be studied on a homogeneous basis yielding important information regarding the thermodynamics of antigen binding. Finally, provision of monoclonal antibodies by hybridoma technology has led to the development of standard diagnostic procedures as well as therapeutically important immunochemical reagents.

Monoclonal antibodies of defined specificity generated through hybridoma methodology have been used to solve many issues related to structure-function relationships within the antibody molecule. Based on studies with monoclonal antibodies, the property of Fab (active site) segmental flexibility within the IgG class of antibodies was solved. Homogeneous binding at both active sites within a bivalent molecule ruled out allosteric effects within antibody molecules. Thus, binding at one active site does not influence antigen binding at the adjacent site. Conformational changes transmitted throughout the molecule subsequent to binding of antigen were dismissed as an explanation for such important phenomena as complement binding and fixation. In all cases, availability of monoclonal antibodies proved important to definitively examine these significant basic questions.

Although hybridoma technology has been the method of choice to produce monoclonal antibodies of defined specificity, new techniques have now emerged. Similar to the construction of single-chain antibodies, gene segments encoding the H and L chain variable domains are genetically fused to genes encoding a bacteriophage coat protein (8). The engineered bacteriophage infect bacterial hosts, and the resulting phage particles express active antibody products on their surface. The latter can be combined with mutagenesis so that the resulting phage display library expresses many different antigen-binding domains. Phages that specifically bind antigen are selected and used to infect bacteria in a second cycle. Each selected phage produces a monoclonal antigen-reactive particle. Primary structures of the variable domains can be determined, and those genes can be fused to the antibody constant region genes to reconstruct a monoclonal antibody. Such genes transfected into myeloma cell lines are expressed and the antibody products secreted. Phage display methodology has important implications for the future.

Monoclonal antibodies proved to be the “tool of nature,” whereby immunologists first learned of the structure-function relationship in antibody molecules. However, these basic relationships extended beyond the antibody molecule when it was determined that other members of the immunoglobulin superfamily of molecules (eg, T-cell receptors, major histocompatibility molecules, T-cell accessory proteins, and adhesion molecules) utilized the same basic motifs and canonical structures first recognized in monoclonal antibodies (9).

### Bibliography

1. E. D. Day (1990) *Advanced Immunochemistry*, 2 ed., Wiley-Liss, New York, pp. 351–396.
2. F. M. Burnet (1957) *Aust. J. Sci.* **20**, 67–69.
3. J. Carrero, M. Mummert, W. D. Mallender, and E. W. Voss Jr. (1997) *Comm. Mol. Cell. Biophys.* **9**, 49–86.
4. G. Köhler and C. Milstein (1975) *Nature* **256**, 495–497.
5. R. E. Bird, K. D. Hardman, J. W. Jacobson, S. Johnson, B. M. Kaufman, S-M Lee, T. Lee, S. H. Pope, G. S. Riordan, and M. Whitlow (1988) *Science* **242**, 423–426.
6. I. Roitt, J. Brostoff, and D. Male (1996) *Immunology*, 4 ed., Mosby, London, p. 28.9.
7. F. W. Putnam, K. Titani, and E. Whitley Jr. (1966) *Proc. Roy. Soc. Lond.* **166**, 124–137.
8. E. A. Kabat and T. T. Wu (1970) *J. Exper. Med.* **132**, 211–250.
9. J. Kuby (1996) *Immunology*, 3 ed., W. H. Freeman & Co., New York, pp. 129–131.

### Suggestions for Further Reading

10. C. A. Janeway Jr. and P. Travers (1996) *Immunobiology*, Garland Publ. Co., New York.
11. T. Hunkapiller and L. Hood (1990) *Adv. Immunol.* **44**, 1.
12. R. E. Bird and B. W. Walker (1991) *Tibtech* **9**, 132–137.

## Monolayer

A *monolayer* (*monomolecular layer*) formed at aqueous interfaces by [amphipathic](#) compounds, such as [fatty acids](#) and phospholipids, is one of the early structural models for one-half of the bilayer leaflet (1) (see [Membranes](#)). Amphipathic molecules are oriented at the interface with the polar moiety embedded in the aqueous phase. Lipid monolayers form physical states and exhibit properties that are analogous to those found for these materials in bulk (2). Thus, “gaseous” and liquid-like monolayers, plus a first-order transition between these states, have been observed (3) with energies of the same magnitude as the comparable systems in bulk. Crystalline monolayers on water have not been observed, although the formation of highly viscous films has been reported.

There are two principal methods of forming monolayers on water, each of which provides a different aspect of the thermodynamic properties of films on water. For materials that are soluble in water, the surface-active solute adsorbs at the interface, and an equilibrium between the surface and the bulk solution is established. This approach provides a method for analyzing the adsorption of amphipathic materials at aqueous interfaces. The **free energy** of adsorption may be obtained from the *Gibbs adsorption isotherm*, which relates the measured surface tension  $g$  (mN/m) with the concentration  $C$  of the dissolved surface active solute as  $dg = -RTGd \ln C$ ;  $G$  represents the concentration of the solute in the surface (moles/cm<sup>2</sup>),  $R$  is the gas constant, and  $T$  is the absolute temperature. The surface tension decreases as solute adsorption increases; the maximum lowering of the surface tension is attained for saturated solutions of the solute.

The second method for forming monolayers is applied to compounds that are poorly soluble in water. The film-forming material is dissolved in a volatile organic solvent, and a dilute solution is deposited in the surface; the solvent evaporates, leaving a monolayer of the material at the interface. The films that form are considered to be insoluble, but in reality some of the material dissolves into a thin (~100  $\mu$ m) unstirred region just beneath the interface, where it slowly diffuses into the bulk of the aqueous phase. The amount of material that dissolves is governed by the Gibbs adsorption isotherm (see text above). Because the material is poorly soluble in water, the amount of lipid that dissolves is usually insignificant relative to the total amount that is deposited initially, and the film behaves as if it is insoluble (3). The surface concentration of the lipid is generally manipulated by moveable barriers to increase or decrease the area in which the film is confined. The surface tension lowering from the film-free surface tension  $g_0$  is usually represented by  $p = g_0 - g$ , as the *surface pressure*. The film balance is an instrument that combines a surface tension measuring device with a trough whose surface area, and therefore the area/mole of lipid  $A = G^{-1}$ , may be manipulated (4). The  $p$ - $A$  relationship obtained with the film balance provides a systematic experimental basis for evaluating thermodynamic properties of films, including the energetics of phase transitions, and of mixture formation with simple lipid components (3). The properties of “insoluble” films are continuous with those obtained from the soluble surface-active compounds, provided that the surface tension lowering for a saturated solution of the material is not exceeded. “Insoluble” lipid monolayers often may be supercompressed to surface pressures that exceed the values for a saturated solution of the lipid. These supercompressed films are metastable and are usually formed from lipids that contain aliphatic saturated hydrocarbon chains of 14 or more carbon atoms.

### Bibliography

1. E. Gorter and F. Grendel (1925) *J. Exp. Med.* **41**, 439–443.
2. I. Langmuir (1917) *J. Am. Chem. Soc.* **39**, 1849–1906. (Contains the first comprehensive measurements of the isotherms of various lipid compounds, and includes many of the concepts of molecular structure that still influence contemporary studies of lipid films.)
3. N. L. Gershfeld (1984) In *Cell Surface Dynamics*, (A. S. Perelson, C. DeLisi, and F. W. Wiegel,

eds.), Marcel Dekker, New York, pp. 93–131. (Includes a summary of the thermodynamic methods and results for the energetics of various surface phase transitions and for two component lipid mixtures.)

4. G. L. Gaines Jr. (1966) *Insoluble Monolayers at Liquid–Gas Interfaces*, Wiley-Interscience, New York.

### Suggestions for Further Reading

5. R. Defay, I. Prigogine, A. Bellemans, and D. H. Everett (1966) *Surface Tension and Adsorption*, Wiley, New York. (This is a rigorous treatise on the thermodynamic properties of surfaces.)
6. N. L. Gershfeld (1976) *Annu. Rev. Phys. Chem.* **27**, 349–368. (A critical survey of film balance studies.)

## Monosomy

Monosomy is a special case of [aneuploidy](#) in which one [chromosome](#) (or part of a chromosome) is represented by one copy instead of two in the normal **diploid** genome.

Partial or total monosomies of all chromosomes in humans have been described. Partial monosomies are due to chromosomal **translocations**. As an example chosen among many, a partial monosomy of chromosome 21 is caused by an unbalanced translocation between the long arms of chromosome 11 and chromosome 21. In this patient, the anomaly had been reported as a full 21 **autosomy** by using cytogenetic banding techniques. The translocation was established by a combination of a high resolution banding technique, chromosome painting, and DNA **polymorphism** analysis (1). In this respect, other descriptions of total monosomies should be reconsidered, because the occurrence of complete autosomal monosomies is now generally considered to incompatible with life. The nonlethal monosomy affecting the [X-chromosome](#) (Turner's syndrome) will be considered below.

Studies on the correlations between **karyotypes** and clinical **phenotypes** have led to the following conclusions:

1. Monosomy for an autosomal segment leads to more severe alterations to the phenotype and restricts survival more than [trisomy](#) for the same segment. When the clinical pictures of these two aneuploid states are compared, they do not go in opposite directions, as would be expected by the type-contrast approach.
2. For certain types of minor anomalies, the aneuploid segment can be narrowed down to a very short region. These regions are closer to the [telomere](#) both in trisomy and in monosomy.

### 1. Monosomy of the Sex Chromosomes

Turner's syndrome is the phenotype associated with the absence of a second sex chromosome in humans. Only about half of all patients with Turner's syndrome are monosomy 45, X on karyotyping, and there are grounds for supposing that cryptic mosaicism for at least part of the [Y-Chromosome](#) may be present in some patients. If so, this would be clinically important because of the risk to patients of gonadal neoplasm and virilization. Although cytogenetic analysis did not detect any Y-chromosomal material, specific nucleotide sequences from the **sex-determination** region of the Y-chromosome (SRY gene) were found by [Southern blot](#) analysis of **PCR** material from the patient's DNA (2).

Recent observations support the hypothesis that the Turner's syndrome **phenotype** results from a [haploid](#) dosage of genes that are common to the X- and the Y-chromosomes and escape X inactivation. A goal of current studies is identifying these putative “Turner” genes (3). Individuals with Turner's syndrome are short and present a characteristic pterygium colli. Most often, they are sterile because ovarian agenesis. Sometimes, functional ovaries allow pregnancies, but the children are often malformed.

#### Bibliography

1. B. Hertz et al. (1993) Clin. Genet. **44**, 89–94.
2. M. Kokova et al. (1993) Lancet **342**, 140–143.
3. A. R. Zinn, D. C. Page, and E. M. Fisher (1993) Trends Genet. **9**, 90–93.

#### Suggestion for Further Reading

4. A. Schinzel (1993) Karyotype-phenotype correlations in autosomal chromosomal aberrations, Prog. Clin. Biol. Res. **384**, 19–31.

### Monte Carlo Calculations

**Computer simulation** calculations of average macromolecular properties by statistical mechanical approaches require one to evaluate the Boltzmann-weighted average of the given property over the entire conformation space (see equation (1) of [Free Energy Calculations](#)). Fortunately, it is possible to estimate such an average without actually sampling the entire enormous conformational space by exploiting the fact that points with high energy do not contribute significantly to the overall average. A powerful and systematic approach of doing so is the so-called Monte Carlo method devised by Metropolis et al. (1) This approach generates random conformations by moving several atoms or torsional angles and uses the corresponding value of the potential surface to accept or reject each new conformation. If the energy of the new conformation ( $n$ ) is lower than the previous one ( $n-1$ ), the new conformation is accepted. If the energy of the new conformation is higher than that of the previous one, a random number,  $R_n$ , between 0 and 1 is generated and compared to the Boltzmann factor

$$P_n = \exp[-(V_n - V_{n-1})/k_B T] \quad (1)$$

If the random number is larger than  $P_n$ , the move is rejected. If  $R_n$  is smaller than  $P_n$ , the conformation is accepted. The desired average is then simply evaluated by

$$\langle A \rangle = \frac{1}{N} \sum_{i=1}^N A(\mathbf{r}^N) \quad (2)$$

This random procedure generates a probability that corresponds to the proper Boltzmann distribution of the entire system if enough points are evaluated. In principle, one needs an enormous number of points, but in practice one can obtain reasonable results with a reasonable amount of computer time (if the property has similar values at different conformations). In fact, Monte Carlo approaches do not involve complete random searches, because they are usually implemented with umbrella sampling or free energy perturbation approaches that bias the search toward important regions in the



conformational space (see [Free Energy Calculations](#)).

Average properties that can be evaluated by [molecular dynamics](#) (MD) techniques can also be evaluated by Monte Carlo approaches. Thus, the selection of the proper method depends on the problem at hand. The advantage of Monte Carlo approaches is the fact that they do not require the evaluation of the first derivatives of the potential functions, which is a prerequisite for MD approaches. Thus, one can use Monte Carlo approaches to nonphysical moves in studies of folding processes and related problems while using lattice models. On the other hand, MD approaches are much “smarter” in that they use derivatives that allow, in fact, much more systematic local searches. Thus, one can state that MD does a much better job in exploring local problems while Monte Carlo allows more uniform global searches. This is useful, for example, in particle insertion methods (2). Of course, one can combine the benefits of both approaches, using Monte Carlo calculations to generate starting conformations and MD for local exploration. Approaches that exploit this and related ideas have been developed (eg, Refs. 3 and 4). In fact, Brownian dynamics is a form of a smart Monte Carlo approach.

### Bibliography

1. N. Metropolis, A. W. Rosenbluth, A. H. Teller, and E. Teller (1953) *J. Chem. Phys.* **21**, 1087–1092.
2. B. Widom (1963) *J. Chem. Phys.* **39**, 2808–2812.
3. F. Guarnieri and W. C. Still (1994) *J. Comp. Chem.* **15**, 1302–1310.
4. M. Rao and B. J. Berne (1979) *J. Chem. Phys.* **71**, 129–132.

### Suggestions for Further Reading

5. M. P. Allen and D. J. Tildesley (1987) *Computer Simulation of Liquids*, Oxford University Press, Oxford, U.K.
6. M. H. Kaols and P. A. Whitlock (1986) *Monte Carlo Methods*, Vol. 1: *Basics*, Wiley, New York.

## Morphogenesis

Morphogenesis is the creation of pattern and form during [development](#). It occurs through changes in cell number, cell shape, and cell adhesion. The first morphogenetic changes occur after gastrulation has formed the three germ layers. The germ layers undergo morphogenesis to form the organs of the embryo. The ectoderm forms the epidermis and nervous system. The mesoderm forms the muscles, the vascular system, the lining of the body cavity, and the gonads. The endoderm forms the gut and digestive glands.

### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed. W. B. Saunders, Philadelphia, p. 5.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.

## Morphogens

A morphogen is a chemical which exists in a concentration or activity gradient and to whose different concentrations or activity cells respond in a qualitatively distinct manner.

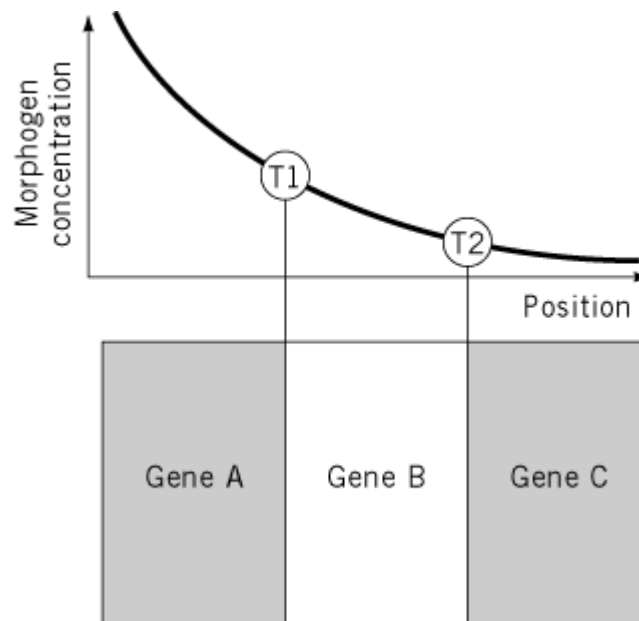
Gradient models have been proposed by developmental biologists to explain certain graded responses and regulative phenomena observed in embryos following experimental manipulations such as tissue transplantations and ablations. In the simplest gradient model, a secreted morphogen diffuses from a local source through a tissue and a concentration gradient is created as a function of the distance from the source by simple dilution. Cells would then respond according to the local morphogen concentration they experience. Gradient models are able to explain many morphogenetic phenomena but convincing molecular evidence for the existence of morphogens has only been documented in a few cases. These examples which are discussed below illustrate a variety of ways by which the concentration gradient of morphogens is created and suggest that the gradients are subject to complex cellular control.

The fundamental question of development is how a fertilized egg progresses from almost homogenous initial conditions to a fully differentiated multicellular organism. It was realized early that many cells appear to “know” their place when subjected to certain experimental conditions, as if under the influence of a morphogenetic field. The concept of a gradient field to explain such phenomena was introduced a century ago by Boveri and proved useful for many schools of embryologists working on diverse systems, most prominently on sea urchins.

However, the notion of what is in fact graded remained necessarily vague in the pre-molecular era. The term “morphogen” was introduced 1952 by Turing (1), who described the mathematical formalism according to which spatial concentration patterns can be generated if two chemicals with different diffusion rates react with each other (reaction-diffusion models). He referred to such diffusible chemicals as morphogens but as we know now such pattern-generating chemicals need not necessarily be diffusible and the term is better defined by the more general definition given on top.

The more tangible formalism of “positional information” for the gradient concept was developed by Wolpert (2). He proposed that in a gradient field cells can read their position - or distance - from a boundary. In the simplest case the gradient is generated by diffusion of a localized morphogen from a source at a cellular boundary. At distinct threshold concentrations these cells respond differently and differentiate accordingly. This is illustrated by the French flag model (Fig. 1).

**Figure 1.** The French flag model. A morphogen gradient (top) provides positional information encoded by the local concentration which can be read by cells. At distinct threshold concentrations (T1, T2) cells respond with differential activity of genes A-C. The tricolore exemplifies the cell fate pattern resulting from the morphogen read-out.



## 1. Criteria for Morphogens

Many chemicals have been proposed to act as diffusible morphogens mostly based on the observation that they can elicit qualitatively different effects in a dose-dependent manner under experimental conditions (3). However, this criterion is not sufficient evidence for a morphogen. A factor which activates a morphogen would indirectly set in motion the same range of responses via signal relay. Activin for example induces different mesodermal cell fates at distinct threshold concentrations in *Xenopus* (4) and it is instructive to study its long-range patterning effects (5, 6), but a role as a physiological mesoderm inducer is doubtful. Therefore, to provide sufficient evidence for a morphogen the following criteria should be fulfilled:

### 1. The chemical should elicit qualitatively different effects in a dose-dependent manner.

Experimentally this can be tested by applying increasing doses of the chemical or in case of gene products by overexpression. An important prerequisite is the possibility of a complex read-out of cellular responses, eg, by being able to score multiple cell fates using molecular differentiation markers.

2. Dose-dependent removal of the chemical should yield qualitatively different effects. Technically this is maybe the most difficult to achieve. Application or expression of specific inhibitors in a dose-dependent manner may be used (7). Genetically, a series of hypomorphic mutants for the [decapentaplegic](#) morphogen have been used to show a dose-dependent requirement for patterning (8, 9). However, given the redundancy of many morphogenetic processes, removal of the morphogen may be without effect unless back-up systems are also inhibited.

The chemical should naturally be present in a concentration gradient. The bicoid and wingless morphogen gradients can be detected by **antibodies** (10, 11). However, as highlighted by the [decapentaplegic](#) and BMP-4 morphogens (see below), what is graded may not be the protein itself but its activity, which may be difficult to detect by antibody.

The chemical should exert its dose-dependent effects directly and not by relay, eg, by inducing a morphogen or a cascade of other pattern-generating factors. For [growth factor](#) type morphogens that act by receptor-mediated [signal transduction](#) this question translates typically in demonstrating that intracellular activation of the signalling pathway does not lead to dose-dependent effects at long-range. This test fails for cases of nonsecreted morphogens, such as the [transcription factor](#) dorsal or

for cases of growth factors which diffuse poorly but whose activity becomes graded by a diffusible antagonist (see below).

Examples of factors in which these criteria for a morphogen have been largely met are discussed below and summarized in Table 1.

**Table 1. Identified Morphogens**

| <b>Morphogen</b>                  | <b>Protein Class</b> | <b>Process Involved</b>    | <b>Mechanism of Gradation</b>    | <b>Ref.</b>    |
|-----------------------------------|----------------------|----------------------------|----------------------------------|----------------|
| <i>Drosophila bicoid</i>          | transcription factor | a/p blastoderm             | diffusion                        | (10, 27, 12)   |
| <i>Drosophila hunchback</i>       | transcription factor | a/p blastoderm             | translational control, diffusion | (12)           |
| <i>Drosophila dorsal</i>          | transcription factor | d/v blastoderm             | nucleo-cytoplasmic distribution  | (13, 28)       |
| <i>Drosophila decapentaplegic</i> | growth factor        | d/v blastoderm, a/p wing   | activity gradation, diffusion    | (8, 9, 14, 15) |
| <i>Drosophila wingless</i>        | growth factor        | <i>Drosophila</i> d/v wing | diffusion                        | (17, 11)       |
| <i>Xenopus BMP-4</i>              | growth factor        | d/v mesoderm               | activity gradation               | (7, 22)        |

## 2. Bicoid and hunchback

The transcription factors bicoid and hunchback were the first morphogens described. They are involved in anteroposterior patterning of the *Drosophila* blastoderm, which at the time of action of these two morphogens is a **syncytium**, allowing free diffusion of cytoplasmic components. *Bicoid* mRNA encoding a *homeodomain* protein is localized and translated at the anterior pole from where the protein diffuses through the syncytium towards the posterior pole. This diffusion gradient can be detected by antibody staining. Overexpression of *bicoid* leads to dose-dependent cell fate changes. The **zinc finger** protein hunchback forms a similar anteroposterior gradient, which is created by translational control. Hunchback and bicoid protein are thought to cooperatively pattern the anterior of the embryo (12). Since this patterning process occurs in a syncytium this example may not occur widely in other multicellular organisms.

## 3. *Dorsal* in *Drosophila* blastoderm

As predicted and as was experimentally found for the transcription factor *dorsal*, a morphogen must not necessarily be diffusible. The essential feature is that the graded parameter reflects distance from a boundary and cells respond at distinct threshold concentrations. Yet, commonly, a morphogen is referred to as a diffusible molecule. *Dorsal* acts as a morphogen in dorsoventral patterning of *Drosophila* blastoderm where it is both necessary and sufficient for dose-dependent patterning (13). Antibody staining revealed a gradient of nuclear localization of dorsal protein along the dorsoventral

axis. Since only nuclear dorsal protein activates transcription this leads to a gradient of effective dorsal protein concentration. The gradient is set up by a cascade of proteases culminating in the local and probably graded activation of the toll receptor which controls the nucleo-cytoplasmic distribution of dorsal. However, it is unclear how the toll receptor signal becomes graded in the first place.

#### 4. *Decapentaplegic* and *wingless* in *Drosophila* wing development

*Decapentaplegic* functions in *Drosophila* anteroposterior wing patterning. Misexpressing *decapentaplegic* leads to dose-dependent and long-range changes of marker gene expression. The effects are direct, as they are not observed by misexpressing an activated form of the *decapentaplegic* receptor (14, 15). Partial loss of *decapentaplegic* signalling leads to changes in marker gene expression which are consistent with a requirement for dose-dependent *decapentaplegic* signalling (16). A *decapentaplegic* gradient has not been physically described in the wing but is made plausible by localized *decapentaplegic* mRNA expression.

Dorsoventral patterning in the wing is regulated by the secreted growth factor *wingless*. Similar to *decapentaplegic*, misexpressing *wingless* leads to dose-dependent and long-range changes of marker gene expression. The effects are direct, as they are not observed by misexpressing intracellular components of the *wingless* transduction pathway. Hypomorphic mutant alleles of *wingless* lead to changes in marker genes consistent with a requirement for dose-dependent *wingless* signalling (17, 11). *Wingless* mRNA expression is localized and a protein gradient has been detected by antibody staining surrounding mRNA expressing cells (11).

#### 5. *Decapentaplegic* in *Drosophila* blastoderm and *BMP-4* in *Xenopus* mesoderm

*Decapentaplegic* is a *Drosophila* homologue of vertebrate BMP proteins and both are members of the TGF- $\beta$  superfamily. In the blastoderm, microinjection of *decapentaplegic* mRNA leads to different cell fates in a dose-dependent manner. Likewise, an allelic series of *decapentaplegic* mutants with different severity leads to different fate changes. Thus, both gain- and loss-of-function of *decapentaplegic* lead to dose-dependent effects (8, 9). An evolutionarily conserved patterning process occurs during dorsoventral patterning of *Xenopus* mesoderm (18). Here, microinjection of *BMP-4* mRNA induces different cell fates in a dose-dependent manner and microinjection of a dominant-negative BMP receptor, an inhibitor of BMP signalling, leads to dose-dependent fate changes in the reverse order (7). A concentration gradient has not been detected for either *decapentaplegic* nor *BMP-4*. However, *decapentaplegic* is antagonized by a secreted inhibitor, short gastrulation, whose vertebrate homologue *chordin* has been shown to antagonize *BMP-4* and is expressed in the Spemann organizer. These antagonists bind and inactivate *decapentaplegic* and *BMP-4*. Since in *Drosophila* as well as *Xenopus* embryos the inhibitors are expressed opposite to *decapentaplegic* and *BMP-4*, an activity gradient is most likely created (19, 20-22). Thus, the antagonists create the morphogen gradient of an instructive signaling molecule. The directness of *decapentaplegic* action has not been tested in this context but *Xenopus* *BMP-4* acts in a direct and long-range fashion (7).

#### 6. Regulating gradient shape

It is evident that morphogen gradient may be of crucial importance for embryonic patterning and thus be expected to be subject to regulatory mechanisms. The *BMP-4* signal spreads about 5-10 cell diameters from the expressing cells within the mesoderm (7). However, the diffusion of *BMP-4* as well as of other TGF- $\beta$  homologues may be more limited in other cells (23, 24) and hence more tightly controlled by yet unknown mechanisms. Besides *chordin*, the *BMP-4* antagonist *noggin* is expressed in the *Xenopus* organizer which diffuses over more than 20 cell diameters (7, 22). A further complexity is introduced by the existence of *tolloid* proteases that cleave the *chordin/BMP-4* and *sog/decapentaplegic* protein complexes. Loss of *tolloid* function leads to changes in patterning both in *Drosophila* blastoderm and *Xenopus* mesoderm and hence these proteins play an important

role in regulating the shape of the decapentaplegic/BMP-4 activity gradient (25, 26).

## 7. Read-out of morphogen gradients

Little is known about how the positional information provided by morphogens is translated into distinct cellular responses. The bicoid transcription factor for example binds to and activates its target gene *hunchback* only at a certain threshold concentration. The *hunchback* promoter region contains multiple sites to which bicoid protein binds in a cooperative manner, leading to a sharp threshold concentration at which bicoid can lead to activation of *hunchback* in an on/off manner (27). Bicoid functions synergistically with *hunchback* to regulate directly further downstream target genes (12). The dorsal protein binds to regulatory sequences in the promoters of its target genes *snail*, *twist* and *decapentaplegic*, leading to activation or repression at distinct threshold concentrations. Thresholds are set by a combination of low and high affinity binding sites for dorsal protein leading to on/off responses (28). Thus, in these cases the target gene promoters serve to convert positional information into distinct cellular responses.

In the case of the diffusible morphogens *decapentaplegic*/BMP-4 and *wingless*, signal transduction pathways lead to the activation of downstream target genes encoding transcription factors. Expression of these target genes is typically not graded and occurs in an overlapping fashion indicating that the morphogen signal has been digitized. Different target genes require distinct threshold doses of morphogen signalling (7, 3). It is unclear in these cases which component of the intracellular signal transduction machinery serves to convert positional information into a more digital response.

## Bibliography

1. A. Turing (1952) *Phil. Trans. Roy. Soc. B.* **237**, 37–72.
2. L. Wolpert (1989) *Development Supplement*, 3–12.
3. C. Neumann and S. Cohen (1997) *Bioessays* **19**, 721–729.
4. J. B. A. Green, H. V. New, and J. C. Smith (1992) *Cell* **71**, 731–739.
5. J. B. Gurdon, P. Harger, A. Mitchell, and P. Lemaire (1994) *Nature* **371**, 487–492.
6. N. McDowell, A. M. Zorn, D. J. Crease, and J. B. Gurdon (1997) *Curr. Biol.* **7**, 671–681.
7. R. Dosch, V. Gawantka, H. Delius, C. Blumenstock, and C. Niehrs (1997) *Development* **124**, 2325–2334.
8. E. L. Ferguson and K. V. Anderson (1992) *Cell* **71**, 451–461.
9. K. A. Wharton, R. P. Ray, and W. M. Gelbart (1993) *Development* **117**, 807–822.
10. W. Driever and C. Nüsslein-Volhard (1988) *Cell* **54**, 95–104.
11. C. J. Neumann and S. M. Cohen (1997) *Development* **124**, 871–880.
12. R. Rivera-Pomar and H. Jäckle (1996) *Trends Genet* **12**, 478–483.
13. D. Morisato and K. V. Anderson (1995) *Annu. Rev. Genetics* **29**, 371–399.
14. T. Lecuit, W. J. Brook, M. Ng, M. Calleja, H. Sun, and S. Cohen (1996) *Nature* **381**, 387–392.
15. D. Nellen, R. Burke, G. Struhl, and K. Basler (1996) *Cell* **85**, 357–368.
16. M. A. Singer, A. Penton, V. Twombly, F. M. Hoffmann, and W. M. Gelbart (1997) *Development* **124**, 79–89.
17. M. Zecca, K. Basler, and G. Struhl (1996) *Cell* **87**, 833–844.
18. E. M. De Robertis and Y. Sasai (1996) *Nature* **380**, 37–40.
19. E. L. Ferguson (1996) *Curr. Opin. Genet. Develop.* **6**, 424–431.
20. S. Piccolo, Y. Sasai, B. Lu, and E. M. De Robertis (1996) *Cell* **86**, 589–598.
21. L. B. Zimmerman, J.-E. De Jesús-Escobar, and R. M. Harland (1996) *Cell* **86**, 599–606.
22. C. M. Jones and J. C. Smith (1998) *Develop. Biol.* **194**, 12–17.

23. C. M. Jones, N. Armes, and J. C. Smith (1996) *Curr. Biol.* **6**, 1468–1475.
24. K. M. Reilly and D. A. Melton (1996) *Cell* **86**, 743–754.
25. G. Marqués, M. Musacchio, M. J. Shimmel, K. Wünnenburg-Stapleton, K. W. Y. Cho, and M. B. O'Connor (1997) *Cell* **91**, 417–426.
26. S. Piccolo, E. Agius, B. Lu, S. Goodman, L. Dale, and E. M. De Robertis (1997) *Cell* **91**, 407–416.
27. W. Driever, V. Siegel, and C. Nüsslein-Vohlard (1990) *Development* **109**, 811–820.
28. J. Rusch and M. Levine (1996) *Curr. Opin. Genet. Dev.* **6**, 416–423.

### Suggestions for Further Reading

29. L. Wolpert (1996) *Trends. Genet.* **12**, 359–364. Discusses positional information and morphogens in general terms.
30. C. Neumann and S. Cohen (1997) *Bioessays* **19**, 721–729. Discusses up-to date morphogen candidates.
31. P. A. Lawrence and G. Struhl (1996) *Cell* **85**, 951–961. Presents a thoughtful synthesis of morphogen action.
32. H. Meinhardt, *Models of Biological Pattern Formation* (1982) Academic Press. London. Computer models simulating reaction diffusion models and their application in patterning problems.
33. L. Wolpert, *Principles of Development* (1998) Oxford University Press. London, New York. Embryology textbook discussing morphogen examples in the context of animal development.

## Morula

The morula is a stage during early embryonic development in embryos with holoblastic cleavage, where there is little yolk and the entire [zygote](#) is involved. After fertilization, the cleavage furrows in holoblastic embryos completely separate the cells. After a few divisions, the cell mass forms a small ball that resembles a mulberry (morula is Latin for mulberry). As the cleavage divisions continue, a hollow space called the *blastocoel* forms with a single layer of cells surrounding it. This stage is called the **blastula** stage.

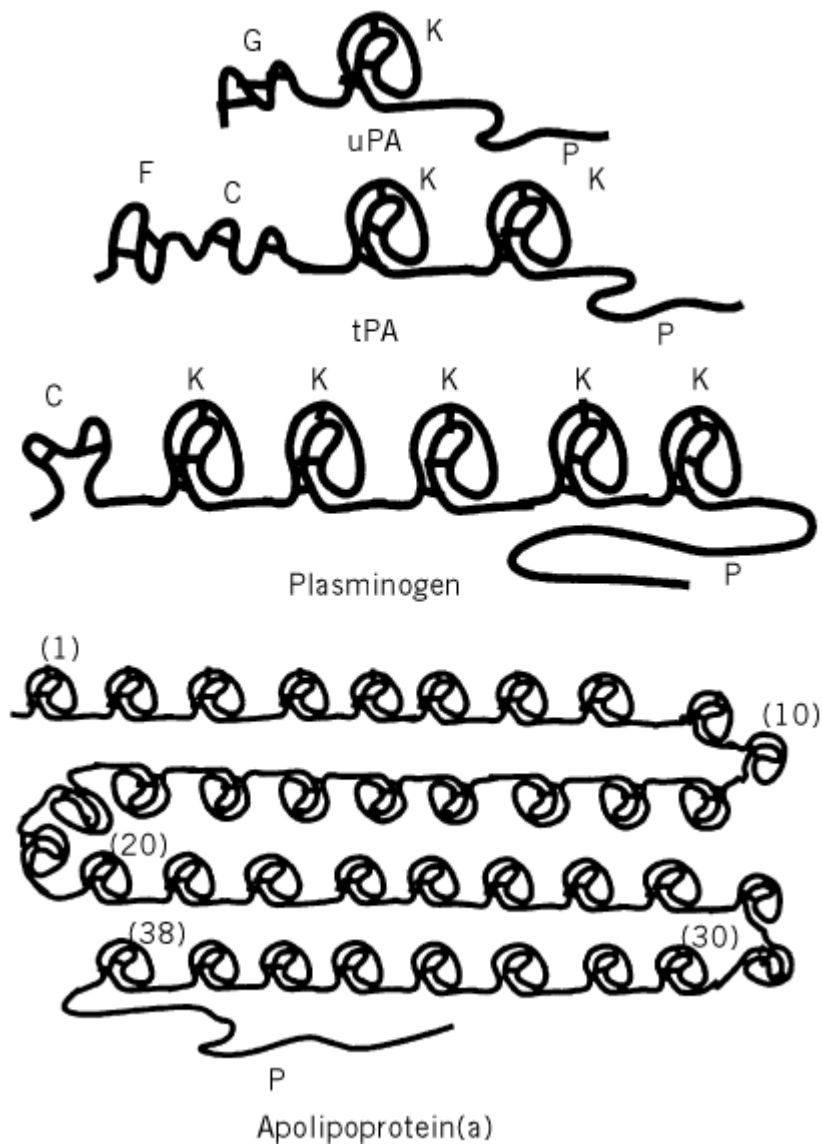
### Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed. W. B. Saunders, Philadelphia, pp. 114–115.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New York.
- S. F. Gilbert (1997) *Developmental Biology*, Sinauer Associates, Sunderland, MA.
- L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.
- J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.

## Mosaic Protein

A mosaic protein is one that is composed of more than one functional **domain**. A typical example of a mosaic protein can be seen in tissue plasminogen activator (tPA) (Fig. 1). tPA contains at least four functional domains; an [epidermal growth factor](#)-like domain (**EGF motif**), two [Kringle domains](#), and a [serine proteinase](#) domain. These DNA segments coding for these domains are believed to have been collected together into a single gene by either [domain shuffling](#) or [exon shuffling](#) during [evolution](#).

**Figure 1.** The domain structures of a number of related mosaic proteins. Kringle domains are indicated by “K” and schematically, (F—finger motif; C—Ca<sup>++</sup> binding domain; P—serine protease domain; G—EGF motif; uPA—urokinase; tPA—tissue plasminogen activator).



Most proteins possess more than one functional domain and can be considered mosaic proteins (Fig. 1). Each such domain usually carries its own function, so a mosaic protein can be considered to be a



multifunctional protein. Consequently, mosaic proteins are evolutionarily important, because they can create diverse functions of proteins. For more detail, see [Exon Shuffling](#) and [Domain Shuffling](#).

## Motor Proteins

A motor protein hydrolyzes ATP to generate movement or force along the [cytoskeleton](#) in eukaryotes. The [actin](#) filaments and **microtubular** components of the cytoskeleton have distinct plus and minus ends, which extend to different regions of the cell. Cells use motor proteins that read the filament polarity to direct specialized cargoes to specific locations. Two motor protein families move along microtubules. The **dyneins** move toward the minus ends of the microtubules, and the **kinesins** move toward the plus ends of microtubules. Flagellar dyneins generate the sliding force between axonemal microtubules that is the motive force for eukaryotic **cilia** and flagella bending. Cytoplasmic dyneins are responsible for moving membrane-bounded organelles and, **kinetochores** and for assembling the spindle. Kinesin family members are also responsible for intracellular transport of membranous organelles moving [chromosomes](#), and assembling the spindle. Myosins, the third motor protein family, interact with actin filaments and move toward the plus end of actin filaments. Myosin family members are responsible for a variety of functions in cells: contraction in muscle, cytokinesis in dividing cells, crawling locomotion, some forms of membrane-bounded organelle transport, and localization of determinants. The three motor proteins are multisubunit complexes that have specific heavy and light chains. All of the members of the three motor protein families are defined by a conserved, force-generating “head” domain that is part of the heavy chain. The heads bind the appropriate cytoskeletal polymer and hydrolyze ATP to move along the polymer. The divergent domains of the heavy chains, the “tails,” and the light chains target the motors to various cargoes or locations and in some cases specify oligomerization of the molecules.

### Suggestions for Further Reading

- J. P. Baker and M. A. Titus (1998) Myosins; matching functions with motors, *Curr. Opin Cell Biol.* **10**, 80–86.
- H. V. Goodson, C. Valetti, and T. Kreis (1997) Motors and membrane traffic, *Curr. Opin Cell Biol.* **9**, 18–28.
- N. Hirokawa (1998) Kinesin and dynein superfamily proteins and the mechanism of organelle transport, *Science* **279**, 519–552.
- T. M. Lothman, K. Thorn, and R. D. Vale (1998) Staying on track: Common features of DNA helicases and microtubule motors, *Cell* **93**, 9–12.

## Mouse

The mouse is an important tool in studying vertebrate biology, including [evolution](#), genetics, ecology, [development](#), physiology, and pathology. The mouse is increasingly important in molecular biology, and its most relevant properties are described here.

### 1. Inbred Strains

The house mouse complex consists of four commensal species, *Mus musculus domesticus*, *M. musculus musculus*, *M. musculus castaneus*, and *M. musculus bactrianus*. In addition, three aboriginal species are known, *M. spicilegus*, *M. macedonicus*, and *M. spretus*. The classic laboratory strains carry a single type of [mitochondrial DNA](#) [see the section thus titled] and primarily the nuclear genes of *M. musculus domesticus*. However, some of the nuclear genes, including the [Y-Chromosome](#), are derived from other commensal species. Although laboratory inbred strains have been established from diverse wild-caught animals by brother-sister matings and back-crossing of offspring with their parents over numerous generations, most of them are predominantly of *M. musculus domesticus* origin. Approximately 430 different mouse inbred strains are currently known. From most inbred strains, several substrains have been generated that are used for different research purposes. Table 1 lists examples of some widely used mouse inbred strains. The characteristics and molecular biology of the inbred strains and their substrains are well known. Most of the inbred strains can be traced back to the beginning of the nineteenth century. Although most differences between substrains and strains are due to mutations or residual heterozygosity, there is also evidence of genetic contamination (1).

**Table 1. Examples of Inbred Mouse Strains and Their Use in Medical Research**

| Strain | Substrain    | Features  |
|--------|--------------|---|
| 129    | 129/Re       | High incidence of spontaneous testicular teratomas  |
|        | 129/RrJ      |   |
|        | 129/Sv-ter/+ |   |
| AKR    | AKR/J        | High incidence of lymphatic leukemia, presence of the Thy 1a T-cell antigen, expression of the ecotropic retrovirus AKV |
|        | AKR/LwN      |   |
|        | AKR/FuA      |   |
|        | AKR/Cum      |   |
|        | AKR/FuRdA    |   |
|        | AKR/TIAld    |   |
| BALB/c | BALB/cHeAn   | Production of plasmacytomas on injection of mineraloil, low incidence of mammary tumors                                 |
|        | BALB/cJ      |   |
|        | BALB/cRl     |   |
|        | BALB/cWt     |   |
| C3H    | C3H/Bi       | High incidence of mammary tumors  |
|        | C3H/Fg       |   |
|        | C3H/He       |   |
|        | C3H/HeJ      |   |
|        | C3H/He-mg    |   |
|        | C3H/He-Avy   |   |
|        | C3H.RV       |   |
| C57BL  | C57BL/A      | Carries the Y-chromosome of the Asian <i>Mus musculus</i> and a LINE-1 element of <i>Mus spretus</i>                    |

|     |  |   |
|-----|--|---|
|     | C57BL/An   |   |
|     | C57BL/GrFa   |   |
|     | C57BL/KaLwN  |   |
|     | C57BL/6  |   |
|     | C57BL/10   |   |
|     | C57BL/10ScSn   |   |
|     | C57BL/10Cr   |   |
|     | C57BL/Ola  |   |
|     | C57BL/Fa   |   |
| CBA | CBA/Ca or<br>CBA/H<br>CBA/Br<br>CBA/CaN<br>CBA/J<br>CBA/St<br>CBA/H-T6 | Low incidence of mammary tumors   |
| DBA | DBA/LiA<br><br>DBA/1<br>DBA/2  | Unfostered substrains have a high incidence of mammary tumors and carry a tumor virus |

---

When some laboratory inbred strains, such as C57BL/10, are crossed with some wild mice from *M. musculus* species, the resulting F<sub>1</sub> male hybrids are sterile (hybrid sterility), whereas the females are fertile (2). Other inbred strains, such as C3H/Di, produce fertile progeny when crossed with wild-type *M. musculus*. Hybrid sterility is also observed when crossing *M. musculus domesticus* and *M. spretus*. Both species separated from a common ancestor about 3 million years ago. The male hybrids are completely sterile, whereas the females are fertile. In both cases, hybrid sterility is caused by a specific configuration of several loci on [chromosome 17](#), which have been described as t-complex or t-haplotypes. T-haplotypes are variant forms of the proximal part of mouse chromosome 17 that were introduced from an unknown species to the present-day mice about 100,000 years ago (3). During evolution approximately 2.5 million years ago, the inversion in (17)2 occurred within the ancestral chromosome. Although this inversion rose to high frequencies, it was never fixed in the wild population. Instead, the ancestral chromosome continued to evolve and accumulated three new inversions, in (17)1, in (17)3, and in (17)4. Finally, the ancestral chromosome became extinct, so that at present only the inverted forms are represented in the population. Consequently, the initial inverted chromosome in (17)2 is called the wild-type chromosome, whereas the three other inversion are called the t-haplotypes.

## 2. Genome

A haploid mouse [genome](#) consists of 20 chromosomes, 19 autosomes and either an X or a Y sex chromosome. Normal somatic cells of a mouse contain a **diploid** genome, or 40 chromosomes. As in other mammals, the physical size of a [haploid](#) genome is approximately  $3 \times 10^9$  bp. The number of genes within the mouse genome and other mammalian genomes is still a matter of discussion, but most estimates are of the order of 100,000 genes. To assess the number of murine genes, a large [expressed sequence tag](#) (EST) project is being carried out under the coordination of the Washington University School of Medicine, St. Louis, Missouri. In the course of this project, approximately

400,000 [complementary DNA](#) (cDNA) clones from different murine tissues will be analyzed. From each of these clones, the 3' end will be sequenced and the sequence put into the dbest section of genbank (site currently unavailable). More than 260,000 mouse ESTs have already been deposited into the nucleotide sequence [databases](#). This means that some sequence information is available on almost all mouse genes although most of the ESTs are not yet linked to complete genes with an assigned function. Simultaneous with the sequencing of all the genes, there are also efforts to monitor the expression of all genes in different tissues and during different stages of [development](#). With detailed data on the spatial and temporal expression patterns of specific genes, for example, [transcription factors](#), it will become feasible to understand the molecular basis for mammalian development. The expression data are collected in the gene expression database (GXD) that is maintained at the Jackson Laboratory, Bar Harbor, Maine (<http://www.informatics.jax.org/>).

As mice are genetically so well characterized through numerous studies with mouse mutants [see below], they have been chosen as one of the model organisms whose genome is to be completely sequenced within the next few years. As a prerequisite to the sequencing, an almost complete **contig** of **yeast artificial chromosome** (YAC) clones of the mouse genome has been established (<http://www-genome.wi.mit.edu/>). For many regions, there are also ordered arrays of **cosmids** available, which will probably serve as the actual sequencing startpoints. From these cloning efforts, extremely dense physical maps were generated, in addition to the longer-existing linkage maps of the mouse genome (4, 5). The combined physical and genetic map can be found in the Mouse Genome Database (MGD), which is also maintained at the Jackson Laboratory ([http://www.informatics.jax.org](http://www.informatics.jax.org/)) (6). This database features a graphic representation of each chromosome with all the mapped loci, together with information on homologous human sequences. In addition to the genetic and physical position on the chromosome, the database is linked to other databases, so that individual **DNA** or [protein](#) sequences can be quickly retrieved. Other links provide information on the **phenotypes** that are associated with mutations of certain genes.

### 3. Repetitive Sequences in the Mouse Genome

The mouse genome contains several different types of repetitive sequences. Apart from simple-sequence repetitive elements, such as **microsatellites** and **minisatellites**, the most common high-copy sequences in the murine genome belong to the [SINE](#) and [LINE](#) families of **retroposons**. The most abundant repetitive sequence is the B1 element, it occurs in about 1 million copies per mouse genome (7). This SINE sequence is derived from the 7SL RNA, like the human **Alu** sequence. B1 elements are found in all rodents, and they can be grouped into three to six subfamilies (8). Another member of murine SINEs are the B2 elements, which are present in much lower copy numbers than B1 elements (9). Intracisternal type A particle (IAP) elements, which are present also in other rodents, represent a middle repetitive class of sequences, with high similarities to [retroviruses](#). They contain *gag*, *pol*, and *env* genes, as well as [long terminal repeat](#) sequences at the 5' and 3' ends (10).

### 4. Mitochondrial DNA

The genome of murine mitochondria is a double-stranded circular DNA molecule of 16295 bp. Its complete nucleotide sequence has been determined (Genbank Accession: J01420) (11). It encodes the 12 S and 16 S **ribosomal RNAs**, 22 [transfer RNAs](#), and 13 proteins. The [genetic code](#) of murine mitochondria is different from that of other mammals in that it uses four [initiation codons](#): AUA, AUC, AUG, and AUU. The codons AGA and AGG are not used. The termination codons are UAA and UAG, but some [messenger RNA](#) transcripts end with U or UA, and a complete termination codon is formed on [polyadenylation](#). AUA codes for methionine, and UGA codes for tryptophan in the mouse mitochondria.

Each strand of the mitochondrial genome is **transcribed** from a single major **promoter**, the heavy strand and the light strand promoters (LSP). In mitochondria, [DNA replication](#) and transcription are closely related. In mouse mitochondria, all major nascent DNA strands initiate within approximately

200 nucleotides downstream of the light strand promoter. In mice, three species of LSP transcripts were observed, one that extends beyond the replication start sites, one with its 3'-terminus mapping exactly at or close to the start sites of the nascent DNA strands, and one that is covalently attached at the 5'-termini of nascent DNA from the leading-strand origin of replication ( $O_H$ ) (12). This type of transcript forms a stable RNA-DNA hybrid also known as an R-loop and serves as a mature [primer](#) for DNA synthesis after endonucleolytic processing. The stability of the RNA-DNA hybrid is probably due to the formation of intramolecular structures involving G:G base pairs stabilized by Hoogsteen base pairing (G-quartet tetraplex model).

Mitochondria, and hence mitochondrial DNA, have an almost exclusively maternal mode of inheritance (uniparental). [Sperm](#) mitochondria and their DNA are eliminated from the **oocyte** cytoplasm in intraspecies crosses at the late **pronucleus** stage (13). However, studies with hybrids between inbred strains (C57BL/6J) and *M. spretus* have shown that heteroplasmy may arise from paternal mitochondrial DNA (14). The persistence of cytoplasmic genomes is observed only in interspecies crosses that rarely occur in nature.

In general, the mutation rate of the mitochondrial DNA is higher than that of the nucleus. Compared with other mammals, however, the mouse mitochondrial DNA again seems to have an even higher rate of mutations (15); the rate per unit time is 40-fold faster in mice than in humans. Postmitotic tissues—for example, the brain, kidney, heart, or skeletal muscle—especially accumulate mutations with age. The etiology of the higher rate of mutations is not yet clear, but mitochondrial DNA is not associated with protective proteins, and it has a less efficient mechanism of DNA repair. In addition, mitochondrial DNA is exposed to reactive oxygen species generated in oxidative phosphorylation at the inner mitochondrial membrane.

## 5. Imprinting, Gametic

[Imprinting](#) is a phenomenon that has been observed in a wide range of phyla from both plant and animal kingdoms. To elucidate the mechanisms underlying imprinting, however, the mouse has been used most extensively. Currently, at least seven terms are used to describe imprinting (*chromosomal*, *parental*, *genomic*, *genetic*, and *gametic imprinting*, as well as simply *imprinting*) (16). Regarding the actual effects and processes associated with imprinting, the term *gametic imprinting* describes the phenomenon best. Gametic imprinting can be grouped into four classes, namely: (1) differential phenotypic effects of parental **alleles**, (2) monoallelic expression, (3) allele-specific **methylation**, and (4) monoallelic genome changes (17). In other words, the activity of a specific gene depends on its parental origin. But not only single genes are regulated by gametic imprinting. It has been shown using **parthenogenetic** (gynogenetic, androgenetic) mouse embryos that mammalian development depends on both parental genomes. From these findings, it is conceivable that gametic imprinting may regulate more than a single locus. Thus far, in fact, nine imprinted loci have been identified in the mouse. Six loci—Igf2, SnRPN, SP2, Ins1, Ins2, and Xist—are repressed on the maternal chromosome. Two loci, Igf2r and H19, are repressed on the paternal chromosome, and one locus, Wt1, has a yet uncharacterized monoallelic expression (16). However, the overall number of genes showing parental effects might remain small, taking into account that mice can be disomic for one parental type of six chromosomes without any obvious effect on the phenotype.

DNA methylation plays an integral role in the maintenance, and probably the establishment, of imprinting. With methylation of the cytosine in the CpG-dinucleotide, binding of transcription activators or [repressors](#) is modulated. From several analyses, it is clear that *trans*- as well as [cis-acting](#) modifiers of imprinting (imprintor loci) exist.

A well-studied imprinted gene is that of the murine insulinlike growth factor type 2 receptor (Igf2r). Igf2r has been mapped to chromosome 17 within a region that was previously shown to be responsible for early embryonal death if two overlapping deletions  $T^{hp}$  and  $t^{Lub2}$  are inherited from the mother only [see “Inbred Strains” above]. In earlier analysis, the effect of the two loci on the

development of mouse embryo was explained by a theoretical gene called Tme (T-associated maternal effect). The Tme locus is located in the same region as Igf2r on chromosome 17 at 7.35 cM. Embryos that have inherited the Tme deletion from their mother do not express the Igf2r gene and die on day 15 of gestation. Therefore, the Igf2r is a candidate gene for Tme. Later it was shown that the Igf2r gene is gametically imprinted in the oocyte primarily by methylation of a region in intron 2 (region 2), where repressor proteins can normally bind (18). The paternally derived allele of Igf2r is not methylated in this region and, therefore, the repressors can bind. At a later stage of development, the paternal allele becomes methylated at the promoter region, causing inactivation of the gene. This process is also referred to as secondary somatic imprinting. Hence, if the Igf2r gene is inherited from the father and becomes inactivated during development, and the Tme deletion is transmitted by the mother, the lack of Igf2r gene expression causes embryonal death.

## 6. Mouse Mutants

A vast number of mouse mutants have been described and characterized genetically. Mice are used as models for all types of genetic variations in mammals, as they reproduce fast and can be kept at relatively low cost per animal. Many significant findings about gene function in higher **eukaryotes** have been derived from mouse mutants. The most comprehensive collection of mutant mice is maintained and distributed by the Jackson Laboratory in Bar Harbor, Maine (<http://www.jax.org/>), which is also the administrator of the major mouse genome databases. Mutations in the mouse genome can arise spontaneously or they can be induced. To induce undirected mutations, chemical mutagens like [ethyl-nitrosourea](#) or ionizing radiation can be used. Mutants can also be created by introducing additional genetic information into the mouse genome. The resulting **transgenic** mice [see the following section] can be used to study the effects of the introduced gene construct *in vivo*. Finally, it has become possible to introduce virtually any change into the mouse genome by so-called targeted mutations. Targeted mutations rely on homologous [recombination](#) of a DNA construct with its genomic **homologue** in embryonic [stem cells](#). When the desired mutation has been achieved in a stem cell, **chimeric** mice can be produced with these cells. If the engineered stem cells become part of the **germ line**, the chimeric mice will pass the mutation on to their offspring, which are then genetically homogeneous. The availability of murine embryonic stem cells was the reason why mice were for a long period the only mammalian species in which targeted mutations could be introduced. To perform targeted mutations in other mammalian species, it is necessary to have either embryonic stem cells or the ability to **clone** the whole organism from a single differentiated cell. The technology of targeted mutations is used primarily to create **knockout** mice in which a single gene is disrupted and thereby completely inactivated. However, it is also possible with targeted mutations to modify genes in order to achieve a specific change in the amino acid sequence of the corresponding translation product or to change expression patterns by promoter mutations. Commonly used mouse mutants are listed in Table 2.

**Table 2. Examples of Widespread Mouse Mutants**

| Mutant | Affected Gene       | Phenotype   | Type of Mutation | Reference   |
|--------|---------------------|---|------------------|-------------|
| Agouti | Agouti ( <i>a</i> ) | Many alleles, for example, <i>A(y)</i> , agouti yellow; homozygous lethal; heterozygous | Spontaneous      | <i>a, b</i> |

|          |  |  |  |                      |
|----------|--|--|--|----------------------|
|          |  | yellow coat color;<br>obesity; immune<br>defects   |  |                      |
| Scid     | DNA-activated<br>protein kinase<br>( <i>Prkdc</i> )                          | Severe combined<br>immunodeficiency;<br>homozygotes lack<br>VDJ recombination<br>of immunoglobulins;<br>no functioning B and<br>T cells                                    | Spontaneous                                  | <a href="#">c</a>    |
| Obese    | Leptin ( <i>Lep</i> )  | Extreme obesity;<br>endocrine defects  | Spontaneous                                  | <a href="#">d, e</a> |
| Nude     | HNF-<br>3/forkhead<br>homologue 11<br>( <i>Hfh11</i> )                       | Homozygotes are<br>hairless because<br>keratinization is<br>defective; skin<br>abnormalities;<br>immune defects;<br>growth, viability, and<br>fertility greatly<br>reduced | Spontaneous                                  | <a href="#">f</a>    |
| Myc      | <i>c-Myc</i>   | Early development of<br>lymphomas  | Transgenic                                   | <a href="#">g</a>    |
| CF       | Cystic fibrosis<br>transmembrane<br>conductance<br>regulator ( <i>Cftr</i> ) | Similar to human<br>cystic fibrosis (CF):<br>intestinal<br>obstructions,<br>pancreatic and lung<br>symptoms less severe<br>than in human CF                                | Targeted<br>mutation<br>(knockout)           | <a href="#">h</a>    |
| Mdx      | Dystrophin<br>( <i>Dmd</i> )   | Related to human<br>Duchenne muscular<br>dystrophy; mdx mice<br>show muscle<br>weakness but normal<br>life span  | Spontaneous                                  | <a href="#">i</a>    |
| Utrn/Dmd | Utrophin ( <i>Utrn</i> )<br>and Dystrophin<br>( <i>Dmd</i> )                 | Very similar to<br>human Duchenne<br>Muscular dystrophy;<br>progressive muscle<br>weakness; mice die<br>of respiratory failure<br>before they reach 20<br>weeks of age     | Targeted<br>mutation<br>(double<br>knockout) | <a href="#">i</a>    |

<sup>a</sup> L. C. Dunn, E. C. Macdowell, and G. A. Lebedeff (1937) *Genetics* **22**, 307–318.

<sup>b</sup> S. J. Bultman, L. B. Russell, G. A. Gutierrez-Espeleta, and R. P. Woychik (1992) *Cell* **24**, 1195–1204.

<sup>c</sup> T. Blunt, D. Gell, M. Fox, G. E. Taccioli, A. R. Lehmann, S. P. Jackson, and P. A. Jeggo (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10285–10290.

<sup>d</sup> A. M. Ingalls, M. M. Dickie, and G. D. Snell (1950) *J. Hered.* **41**, 317–318.

<sup>e</sup> Y. Zhang, R. Proenca, M. Maffei, M. Barone, L. Leopold, and J. M. Friedman (1994) *Nature* **372**, 425–432, and (1995) *Nature* **374**, 479, erratum.

<sup>f</sup> M. Nehls, D. Pfeifer, M. Schorp, H. Hedrich, and T. Boehm (1994) *Nature* **372**, 103–107.

<sup>g</sup> J. M. Adams, A. W. Harris, C. A. Pinkert, L. M. Corcoran, W. S. Alexander, S. Cory, R. D. Palmiter, and R. L. Brinster (1985) *Nature* **318**, 533–538.

<sup>h</sup> J. N. Snouwaert, K. K. Brigman, A. M. Latour, N. N. Malouf, R. C. Boucher, O. Smithies, and B. H. Koller (1992) *Science* **257**, 1083–1088.

<sup>i</sup> A. S. Ryder-Cook, P. Sicinski, K. Thomas, K. E. Davies, R. G. Worton, E. A. Barnard, M. G. Darlison, and P. J. Barnard (1988) *EMBO J.* **7**, 3017–3021.

<sup>j</sup> A. E. Deconinck, J. A. Rafael, J. A. Skinner, S. C. Brown, A. C. Potter, L. Metzinger, D. J. Watt, J. G. Dickson, J. M. Tinsley, and K. E. Davies (1997) *Cell* **90**, 717–727.

## 7. Transgenic Mice

Transgenic mice carry additional genetic information that is not derived from the mouse genome. There are several ways to introduce foreign genetic material into the mouse germ line so that it is faithfully passed to succeeding generations. The most popular way is the [microinjection](#) of a DNA construct into a fertilized oocyte ([19](#)). Fertilized oocytes are isolated from hormonally superovulated females that have been mated before the procedure. The DNA construct is microinjected with a glass capillary into the male pronucleus of the fertilized oocyte. The microinjected oocyte is then implanted into a hormonally treated recipient mouse. About 50% of the microinjected oocytes will survive the microinjection and implantation process. Of the resulting offspring, up to 50% will have integrated the gene construct into their genome, so that the overall efficiency is between 20% and 25% under optimal conditions. The gene construct is usually integrated in multiple copies (up to approximately 100) in a head-to-tail [concatemer](#) arrangement. The integration site is random and, in some cases, rearrangements of the DNA, such as deletions, duplications, or translocations, have been observed in the neighborhood of the integration site. Experimental verification of an integration event is done by **PCR** or **Southern blotting**. Mouse genomic DNA for these experiments is usually extracted from a piece of tail. Whether the additional genetic material is expressed depends not only on the regulatory sequences present on the introduced DNA construct but also on the site of integration. Integration sometimes takes place at transcriptionally silent parts of the [chromatin](#), and even genes under the control of strong viral promoters will not be expressed. Whether a transgenic mouse expresses an introduced DNA construct has therefore to be determined experimentally for each individual transgenic mouse line. The integration efficiencies are higher if the DNA construct has been linearized and if all unnecessary vector sequences have been removed. A big advantage of this method is that the size of the injected DNA construct can be very large. As long as intact DNA molecules can be injected through the glass capillary, they will be integrated, and constructs of more than 100 kbp have been successfully introduced into transgenic mice by this method. For example, it is possible to introduce whole YACs into transgenic mice. This ensures that the original chromosomal structure of the transgene and all *cis*-acting regulatory elements are much better retained than in experiments in which only a cDNA construct or a smaller piece of genomic DNA is injected.

Another method for introducing genetic information into mice is to treat mouse embryos with recombinant retrovirus ([20](#)). As retroviruses integrate with high efficiencies into the host genome during their normal life cycle, they can be used as vehicles to introduce foreign DNA into the mouse genome. The constructs used for this purpose contain all necessary viral elements for their integration into the genome but are replication-deficient, so that no active retroviruses are produced within the target cells. Treatment of the mouse embryo with the recombinant retrovirus can be done either *in vitro* prior to implantation into a recipient mother or *in utero* after implantation of the



embryo, although with much lower transfection rates. The integration of retroviral DNA is very efficient, and usually one copy of the retrovirus will integrate at a random site in the genome without major rearrangements of the integration locus. Genetic transfer with retroviruses is usually achieved after the single-cell stage of mouse development, so that the resulting mice will be genetically heterogeneous. Genetically homogeneous transgenic mice can be produced in the F1 generation of these mice. With retroviruses, the size of the introduced DNA construct is limited to approximately 8 kb. Therefore, it is not possible to use this approach for very large genes or for genomic constructs with multiple introns.

The third commonly used method for producing transgenic mice is the introduction of genetic material into embryonic stem cells. With this technique, it is possible to modify specific sites in the genome (21). Additional DNA can be introduced, but endogenous DNA can also be deleted by this technique, which is mostly used to produce knockout mice.

## 8. Mutagenesis Models

Very often, the **mutagenic** potential of a chemical compound is estimated with the [Ames test](#). In this test, a strain of *Salmonella* that has a mutation in a gene of the [his operon](#) coding for an [enzyme](#) of the histidine biosynthesis pathway is treated with the putative mutagen and plated onto a medium without histidine. On this medium, only bacteria with intact histidine biosynthesis will grow, ie, bacteria that have reverted to the wild type or acquired a compensating mutation for the original defect. The number of colonies in an Ames test is therefore a direct measure of the mutagenic potential of a compound. This test is commonly used because it is very easy and cheap to perform.

However, some chemicals exhibit different mutagenic potential in bacteria and in eukaryotes. Within multicellular eukaryotes, there are also tissue-specific mutation frequencies. Two mouse models have been developed that allow a more accurate analysis of mutagenesis in mammals. The Big Blue mouse is a transgenic mouse that carries a DNA construct derived from [lambda phage](#) containing a *lacI* gene (1). The *lacI* gene codes for the [lac repressor](#), which prevents the expression of the [lac operon](#). This DNA construct can be recovered easily from high molecular weight DNA isolated from a Big Blue mouse if lambda packaging extract is added. Because of the phage lambda regulatory sequences on the transgene, the construct will be excised by the packaging extract and packaged into infectious phage particles. These phage are then plated together with a bacterial host onto X-Gal containing medium. If the phage contains an intact lac repressor, the plaques will be clear. If the *lacI* gene is mutated, however, blue plaques will be generated. It is then also very easy to investigate the exact nucleotide change of the mutation. The region containing the *lacI* gene is amplified by PCR primers, and the PCR product is sequenced and compared to the wild-type sequence.

A similar mutagenesis model is provided by the MutaMouse (2). This is also a transgenic mouse carrying a lambda-derived DNA construct. However, the additional DNA in the MutaMouse contains an intact *lacZ* gene (for b-galactosidase) instead of the *lacI* gene. The experimental procedure is the same as with the Big Blue mouse. The only difference is that mutant phage will produce clear plaques and wild-type phage will produce blue plaques in this assay. This system was further improved by incorporating a positive selection for *lacZ*<sup>-</sup> bacteria. In this system, only mutant phage can form plaques, and very large numbers of phage particles can be analyzed (3).

## Bibliography

1. E. M. Prager, H. Tichy, and R. D. Sage (1996) *Genetics* **143**, 427–446.
2. J. Forejt (1996) *TIG* **12**, 412–417.
3. L. M. Silver (1993) *TIG* **9**, 250–254.
4. W. F. Dietrich, J. Miller, R. Steen, M. A. Merchant, D. Damron-Boles, Z. Husain, R. Dredge, M. J. Daly, K. A. Ingalls, and T. J. O'Connor (1996) *Nature* **380**, 149–152.
5. W. F. Dietrich, J. Miller, R. Steen, M. A. Merchant, D. Damron-Boles, Z. Husain, R. Dredge,

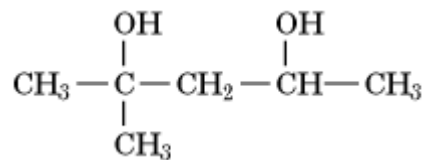
- M. J. Daly, K. A. Ingalls, and T. J. O'Connor (1996) *Nature* **381**, 172.
6. J. A. Blake, J. T. Eppig, J. E. Richardson, and M. T. Davisson (1998) *Nucl. Acids Res.* **26**, 131–138.
  7. P. L. Deininger and M. A. Batzer (1993) *Evol. Biol.* **27**, 157–196.
  8. E. Zietkiewicz and D. Labuda (1996) *J. Mol. Evol.* **42**, 66–72.
  9. W. Bains and K. Temple-Smith (1989) *J. Mol. Evol.* **28**, 191–199.
  10. S. Aota, T. Gojobori, K. Shigesada, H. Ozeki, and T. Ikemura (1987) *Gene* **56**, 1–12.
  11. M. J. Bibb, R. A. Van Etten, C. T. Wright, M. W. Walberg, and D. A. Clayton (1981) *Cell* **26**, 167–180.
  12. D. Y. Lee and D. A. Clayton (1996) *J. Biol. Chem.* **271**, 24262–24269.
  13. H. Kaneda, J.-I. Hayashi, S. Takahama, C. Taya, K. Fischer Lindhal, and H. Yonekawa (1995) *Proc. Natl. Acad. Sci. USA* **92**, 4542–4546.
  14. U. Gyllensten, D. Wharton, A. Josefsson, and A. C. Wilson (1991) *Nature* **352**, 255–257.
  15. E. Wang, A. Wong, and G. Cortopassi (1997) *Mutat. Res.* **377**, 157–166.
  16. D. P. Barlow (1994) *TIG* **10**, 194–199.
  17. S. Varmuza and M. Mann (1994) *TIG* **10**, 118–123.
  18. A. Wutz, O. W. Smrzka, N. Schweifer, K. Schellander, E. F. Wagner, and D. P. Barlow (1997) *Nature* **389**, 745–749.
  19. R. D. Palmiter and R. L. Brinster (1985) *Cell* **41**, 343–345.
  20. T. A. Stewart, P. K. Pattengale, and P. Leder (1987) *EMBO J.* **6**, 3701–3709.
  21. K. R. Thomas and M. R. Capecchi (1987) *Cell* **51**, 503–612.
  22. S. W. Kohler, G. S. Provost, A. Fieck, P. L. Kretz, W. O. Bullock, D. L. Putman, J. A. Sorge, and J. M. Short (1991) *Environ. Mol. Mutagen* **18**, 316–321.
  23. J. A. Gossen, W. J. de Leeuw, C. H. Tan, E. C. Zwarthoff, F. Berends, P. H. Lohman, D. L. Knook, and J. Vijg (1989) *Proc. Natl. Acad. Sci. USA* **86**, 7971–7975.
  24. S. W. Dean and B. Myhr (1994) *Mutagenesis* **9**, 183–185.

### **Suggestions for Further Reading**

25. J. T. Eppig and L. M. Silver (1997) Chromosome committee reports for the mouse genome. *Mamm. Genome* **7**, S1–S388.
26. P. A. Leighton, I. R. Saam, R. S. Ingram, and S. M. Tilgham (1996) Genomic imprinting in mice: its function and mechanism. *Biol. Reprod.* **54**, 273–278.
27. G. Brem (1993) "Transgenic animals". In *Biotechnology* (H.-J. Rehm and G. Reed, eds.) VCH, Weinheim, Germany, Vol. **2**, pp. 745–832.

### **MPD (2-Methyl-2,4-Pentanediol)**

MPD is commonly used to crystallize proteins for [X-ray crystallography](#) (see [Crystallization](#)). It is a relatively [nonpolar](#) molecule



although soluble in water. It decreases the solubility of proteins because it is excluded from the solvent region around their surfaces by a strong repulsion between MPD and the electric charges of the ionized groups on the protein surface. The surface of a globular protein is a mosaic of many charges, and the charge density may be very high, even if the net charge on the protein is close to zero. MPD is repulsed by the charges and concentrates in the bulk solvent. The presence of MPD lowers the solubility of the protein. When the solubility limit of the protein is reached, it precipitates or crystallizes.

The exclusion of MPD from the surface of a protein results in [preferential hydration](#) of the protein, which would be expected to stabilize the protein folded conformation (see [Stabilization And Destabilization By Co-Solvents](#)). However, unfolding reduces the charge density on the protein molecule and, consequently, the repulsion of MPD. Furthermore, MPD interacts favorably with the nonpolar surfaces of proteins that are exposed upon unfolding. Consequently, MPD interacts more favorably with the unfolded state and actually destabilizes folded proteins. So it must be used with caution.

#### Suggestion for Further Reading

S. N. Timasheff and T. Arakawa (1997) "Stabilization of protein structure by solvents", In *Protein Structure: A Practical Approach*, 2nd ed. (T.E. Creighton, ed.) IRL Press, Oxford, pp. 349–364.

## Mu Phage

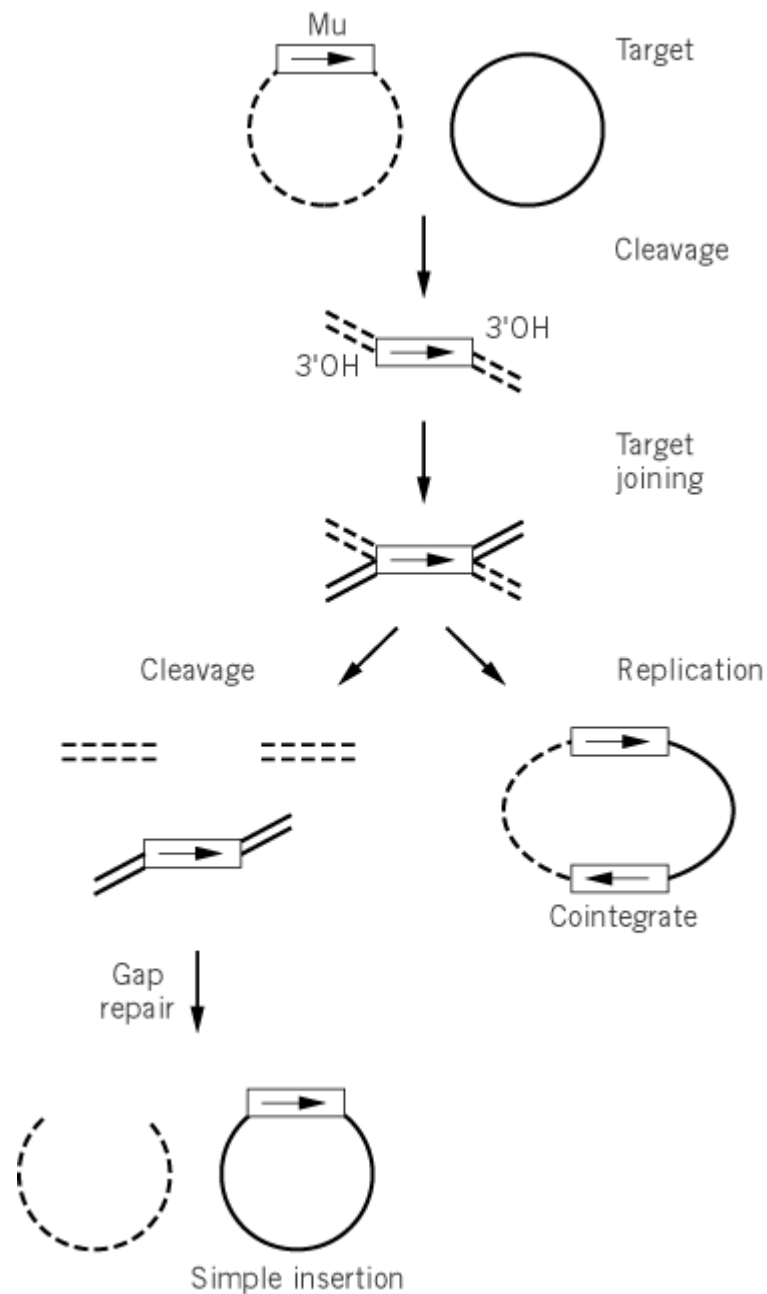
The bacteriophage Mu is a [transposable element](#) (1). It uses [transposition](#) to insert into the *E. coli* chromosome after infection and form Mu **lysogens**. The low target-site selectivity of insertion results in insertion mutations at many different locations, reflecting the disruption of many different genes, leading to a “mutator” (“Mu”) phenotype. Mu also uses transposition to replicate its DNA during lytic growth; multiple rounds of transposition result in the formation of multiple Mu-specific [replication forks](#), and [DNA replication](#) from these forks produces multiple copies of Mu. Understanding Mu has played a key role in the dissection of the control and mechanism of transposition. Mu has also been key in studying the mechanism of transposition, because Mu transposition occurs at a high frequency during lytic growth, rather than the low frequency observed for most other elements. Indeed, the first *in vitro* transposition system was established in 1983 by Mizuuchi using Mu (2).

The DNA of Mu is about 40 kbp in length and, in addition to a number of viral proteins, encodes two transposition proteins: One is MuA, the [transposase](#) that binds specifically to the ends of Mu DNA and executes the DNA breakage and joining reactions; the other is MuB, an ATP-binding protein that interacts with the target DNA and also regulates the activity of MuA. There are special transposition sequences at each end of the element, including multiple MuA binding sites near each terminus and, at about 1.2 kbp internal from the end, an [enhancer](#) site that can strongly stimulate recombination from different positions in either orientation (see Fig. 1 in [Transposable Elements](#)). The enhancer also contains several MuA sites flanking a binding site for a sequence-specific DN-bending protein called *integration host factor*. Another nonspecific DNA-bending protein, called

HU, also participates in recombination. Both of these proteins probably promote a particular architecture of the transposase-end complex. A key early step in recombination is the three-way synapsis of the MuA-bound transposon ends and enhancer. This step is critical in promoting a conformational change in MuA that converts it from a simple [DNA-binding protein](#) to an active recombinase that can execute DNA breakage and joining. Although not essential, MuB bound to target DNA can also play a role in regulating end synapsis. Thus, by influencing the assembly of the synaptic structure, MuB, via its interaction with target DNA, can play a key role in controlling transposition.

The first chemical step in Mu transposition is the introduction by transposase of a nick at the 3' ends of the transposon; no double-strand breaks are made (3) (Fig. 1). These exposed 3'OH transposon ends then attack the target DNA at staggered positions. The product of these reactions contains the transposon covalently linked to the target DNA through its 3' ends and also linked to the donor DNA by its 5' ends; it is variably called a *fusion product*, a *strand transfer product*, or a *Shapiro intermediate* [Shapiro first proposed this intermediate in a model for Mu transposition (4)]. This intermediate can undergo either of two alternative processing reactions. In one pathway, an endonuclease activity (possibly provided by the host) clips the 5' Mu ends, disconnecting the target DNA from the transposon inserted into the target DNA; thus, Mu forms a simple insertion product by a two-step mechanism. In the other processing reaction, replication across both strands of the phage initiates from the two 3'OH target ends that flank the transposon in the fusion product intermediate. Once the replication of both strands is completed, there are now two copies of the transposon linked by the donor DNA and the target DNA, in a structure called a *cointegrate*. Thus, a transposition event can provide two replication forks that mediate Mu replication; that is, two copies of transposon are derived from the transposition product formed by a single Mu substrate DNA. Formation of a cointegrate is an example of replicative transposition.

**Figure 1.** Mu transposition pathways. Mu transposition initiates with cleavages that expose the 3'OH ends of the element; these exposed ends are then joined to target DNA to make a fusion product that contains the transposon attached to the target DNA via its 3' ends and still attached to the donor DNA via its 5' ends. In replicative transposition, replication initiates from 3'OHs in the flanking target DNA at both ends of Mu. This replication results in a cointegrate, a single DNA molecule containing two copies of the transposon, the donor backbone, and the target DNA. In the nonreplicative pathway, another set of cleavages occur that disconnect the 5' ends of the transposon from the flanking donor DNA. A simple insertion results after host repair of the gaps that flank the newly inserted element (see Fig. 1 in [Transposable Elements](#)).



The key chemical steps of Mu transposition are cleavages to expose the 3' ends of the elements and the joining of the 3'OH transposon ends to the target DNA. This polarity was first described with Mu, and we now know it is true of all other transposition reactions examined that involve processing of a DNA intermediate. The [X-ray crystallography](#) structures of the MuA transposase (5) and human immunodeficiency virus (HIV) retroviral integrase (6) demonstrate that the catalytic centers of these transposases are very similar in structure, although there is little amino acid sequence [homology](#) between these proteins. A notable feature of these structures is a cluster of acidic amino acids that are not adjacent in the primary sequence. These acidic residues form a binding site for  $Mg^{2+}$ , which is an essential cofactor in all known transposition reactions and probably places a key role in the actual execution of the DNA breakage and joining steps. Many other transposable elements in bacteria and eukaryotes execute the same 3'OH steps in transposition, and it is likely that their catalytic regions will be related to the common form found in MuA transposase and HIV [integrase](#) (7, 8).

## Bibliography

1. B. D. Lavoie and G. Chaconas (1996) *Curr. Top. Microbiol. Immunol.* **204**, 83–102.
2. K. Mizuuchi (1983) *Cell* **35**, 785–794.
3. K. Mizuuchi (1992) *J. Biol. Chem.* **267**, 21273–21276.
4. J. A. Shapiro (1979) *Proc. Natl. Acad. Sci. USA* **76**, 1933–1937.
5. P. Rice and K. Mizuuchi (1995) *Cell* **82**, 209–220.
6. F. Dyda, A. B. Hickman, T. M. Jenkins, A. Engelman, R. Craigie, and D. R. Davies (1994) *Science* **266**, 1981–1986.
7. J. Kulkosky, K. S. Jones, R. S. Katz, J. P. G. Mack, and A. M. Skalka (1992) *Mol. Cell Biol.* **12**, 2331–2338.
8. P. Polard and M. Chandler (1995) *Mol. Microbiol.* **15**, 1–23.

## Multifunctional Proteins

Multifunctional proteins combine several autonomous functions on a single [polypeptide chain](#). In this sense autonomy implies that each function is assigned to a different region or **domain** of the polypeptide chain. Multifunctional proteins contrast with multienzyme complexes, in which different polypeptides are noncovalently associated. The definition excludes [enzymes](#) that catalyze different reactions using the same [active site](#), such as *asparaginase* functioning as a glutaminase, *phosphoglyceromutase*, which can catalyze three different reactions using the same reaction center, or cystathionine-g-synthase which can catalyze a whole series of analogous reactions because of its relative lack of substrate specificity. **Allosteric** enzymes are not considered multifunctional. Even though they possess binding sites for effector ligands (inhibitors or activators) situated in different regions of the polypeptide chain, they are defined only with respect to the catalytic function.

Multifunctional proteins generally have different catalytic functions residing in separate domains of the same polypeptide chain. A few examples in *Escherichia coli* are **DNA polymerase I**, phosphoribosylanthranilateisomerase-indoleglycerolphosphate synthetase (([1](#)), ([2](#))) (see [TRP Operon](#)), the aspartate kinases-homoserine dehydrogenases (([3](#)), ([4](#))), and chorismate mutase-prephenate dehydrogenase (([5](#)),([6](#))). In *Neurospora crassa*, the two functions of tryptophan synthase (see [TRP Operon](#)), namely the cleavage of indoleglycerol phosphate and the combination of indole with serine, reside on the same polypeptide chain, whereas in *E. coli* each is carried by an independent chain, noncovalently linked to the other. In this case, the *E. coli* enzyme is not a multifunctional protein but a multienzyme complex (([7](#)),([8](#))).

Other examples are the flavocytochrome  $b_2$  of *Saccharomyces cerevisiae* ([9](#)), a homotetramer that is both an l-lactate-cytochrome *c*-oxidoreductase (a flavoprotein) and a cytochrome  $b_2$ .

### 1. More than two Catalytic Functions may be Fused on the same Polypeptide Chain

DNA polymerase I from *E. coli* is a monomer endowed with three activities: DNA polymerase, 5'-3' exonuclease, and 3'-5' exonuclease ([10](#)). Carbamylphosphate synthetase from ascites hepatoma cells is a homotetramer also endowed with [aspartate transcarbamoylase](#) activity (in *E. coli*, these activities are carried by independent proteins) (([11](#)),([12](#))). In *N. crassa* a single polypeptide chain, organized as a homodimer, is responsible for five of the biochemical reactions leading to the biosynthesis of the aromatic ring ([13](#)).

There are situations where the two systems coexist. The fatty acid synthetase of *S. cerevisiae* is a multienzyme complex ( $\alpha_6\beta_6$ ) composed of two multifunctional proteins. One of the polypeptides carries the malonyl transacetylase, the  $\beta$ -hydroxyacyldehydrase, the enoyl reductase, and the palmitoyl deacylase activities, whereas the other polypeptide chain carries the  $\beta$ -ketoacylsynthase, the  $\beta$ -ketoacylreductase, and the acylcarrier protein. In combination with acetylCoA carboxylase, the complex catalyzes the condensation of acetyl subunits to the final product, palmitic acid (14). The avian fatty acid synthetase, which carries the same series of reactions, is a homodimer organized in a head-to-tail manner. Each of the identical polypeptides carry all of the above-mentioned enzymes (15).

## 2. Catalytic and Noncatalytic Functions may be Fused

In the case of cytochrome  $b_2$  from calf liver **microsomes**, one domain of the protein represents the cytochrome  $b_2$  proper, whereas another domain is responsible for an anchor function (16). The homodimeric mammalian meromyosin has a domain endowed with [ATPase](#) activity, whereas another domain has a structural function (forming **thick filaments**) (17). [Diphtheria toxin](#) is a monomer whose amino terminal part ribosylates the mammalian elongation factor, whereas the C-terminal part is responsible for transporting the catalytic part across the membrane (18). The [biotin repressor](#) is a very interesting bifunctional protein of 321 amino acid residues that acts at two different levels. In addition to its repressor function, it is endowed with acetyl CoA carboxylase holoenzyme synthetase activity (19).

## 3. Noncatalytic Functions can be Fused

Human and bovine [serum albumins](#) consist of three homologous domains. They are endowed with independent binding and transport functions for tryptophan, bilirubin, and long-chain fatty acids (20). The [Lac repressor](#), a homotetramer, has two independent domains, a small N-terminal domain that recognizes the corresponding operator DNA and a core that binds the **inducer**. The amino terminal part of the constitutive monomer mediates the assembly into tetramers (21). Other repressors also have separate domains for ligand binding and for specific DNA binding. Mammalian **immunoglobulins** have separate [antigen](#) binding and **complement** fixation sites (22).

## 4. Structural Evidence for a Multifunctional Protein

1. There should be more than one function on a single polypeptide chain.
2. Autonomy of these functions must be demonstrated by the existence of distinct domains on this polypeptide chain. Genetic analysis may produce evidence if it is possible to isolate mutants that are defective in only one function. Single-point mutations may result in the loss of more than one function. In this case, however, pleiotropic effects should be excluded. Conclusive evidence may be obtained by constructing a detailed genetic map, as done for several enzymes from *E. coli* ((23),(24) ).
3. A method providing convincing evidence is the isolation and characterization of fragments that have retained their individual function unimpaired. The N-terminal domain may be produced and isolated using [nonsense mutants](#). Limited **proteolysis** is also useful because the hinge peptide that links individual domains is often very sensitive to proteolytic attack (25).
4. Chemical modification may affect one function without affecting another. For example, *E. coli* aspartate kinase I activity is destroyed by treating the enzyme with sulfhydryl reagents, whereas the other catalytic activity of the bifunctional protein, homoserine dehydrogenase, remains intact (26).

The proper physical methods are **ultracentrifugation**, **gel filtration**, and **PAGE** in with and without

denaturing agents to determine the molecular weight of the native protein and the number of its subunits. [Isoelectric focusing](#) may provide indirect evidence that the subunits are identical. The number of autonomous functions must exceed the number of separable protein bands. The protein band must be homogeneous. Final proof is evidently provided by the total sequence of the protein or of the coding gene (**cDNA** in the case of **eukaryotes**).

### Bibliography

1. T. E. Creighton and C. Yanofsky (1966) *J. Biol. Chem.* **241**, 4625.
2. T. E. Creighton (1970) *Biochem. J.* **120**, 699.
3. F. Falcoz-Kelly et al. (1972) *Eur. J. Biochem.* **28**, 507–519.
4. J-C Patte, G. LeBras, and G. N. Cohen (1967) *Biochim. Biophys. Acta* **136**, 245–257.
5. G. L. E. Koch, D. C. Shaw, and F. Gibson (1971) *Biochim. Biophys. Acta* **229**, 795–804.
6. R. G. H. Cotton and F. Gibson (1965) *Biochim. Biophys. Acta* **100**, 76–88.
7. W. H. Matchett and J. A. DeMoss (1975) *J. Biol. Chem.* **250**, 2941–2946.
8. I. P. Crawford and C. Yanofsky (1968) *Proc. Natl. Acad. Sci. USA* **44**, 1161.
9. M. Mevel-Ninio, Y. Risler, and F. Labeyrie (1977) *Eur. J. Biochem.* **73**, 131–140.
10. A. Kornberg (1974) *DNA Synthesis*, W. H. Freeman, San Francisco.
11. M. Mori and M. Tatibana (1975) *J. Biochem.* **78**, 239–242.
12. J. C. Gerhart and A. B. Pardee (1967) *J. Biol. Chem.* **237**, 891.
13. N. W. Giles et al. (1967) *Proc. Natl. Acad. Sci. USA* **58**, 1453–1460.
14. E. Schweizer (1980) In *Multifunctional Proteins* (H. Bisswanger and E. Schminke-Ott eds.), Wiley, New York, pp. 197–215.
15. S. J. Wakil, J. K. Stoops, and V. C. Joshi (1983) *Ann. Rev. Biochem.* **52**, 537–579.
16. J. Ozols and C. Gerard (1977) *Proc. Natl. Acad. Sci. USA* **74**, 3725–3729.
17. H. E. Huxley (1976) In *The Molecular Basis of Motility* (L. M. G. Heilmeyer, J. C. Rüegg and T. Wieland, eds.), Springer, Heidelberg, pp. 9–25.
18. A. M. Pappenheimer (1977) *Ann. Rev. Biochem.* **46**, 69–94.
19. P. K. Howard, J. Shaw, and A. J. Otsuka (1985) *Gene* **35**, 321–331.
20. A. D. McLachlan and J. E. Walker (1977) *J. Mol. Biol.* **112**, 543–558.
21. T. J. Platt, J. G. Files, and K. Weber (1973) *J. Biol. Chem.* **248**, 110–121.
22. R. R. Porter (1971) *Harvey Lectures* **65**, 157–174.
23. C. Yanofsky et al. (1971) *Genetics* **69**, 409–433.
24. I. Saint Girons (1978) *Mol. Gen. Genetics* **162**, 95.
25. M. Véron, F. Falcoz-Kelly, and G. N. Cohen (1972) *Eur. J. Biochem.* **28**, 520–527.
26. P. Truffa-Bachi et al. (1968) *Eur. J. Biochem.* **5**, 73–80.

### Suggestion for Further Reading

27. H. Bisswanger and E. Schminke-Ott, eds. (1980) *Multifunctional Proteins*, Wiley, New York.

### Multigene Family

When similar **nucleotide sequences** that code for protein and have been derived from a common



ancestor (**homologous**) exist in a repeated fashion in the [genome](#), such a group is referred to as a multigene family, or gene family. For more detail, see [Gene Families](#).

## Multiple-Wavelength Anomalous Dispersion MAD

[Anomalous dispersion](#) gives information about the phases of the reflections measured in [X-ray crystallography](#) (see [Phase Problem](#)). In MAD measurements, [isomorphous replacement](#) is not necessary to solve the crystal structure of a macromolecule. Electrons in an atom are bound by the nucleus and are, in principle, not free electrons (see [Anomalous Dispersion](#)). However, they can be regarded as such if the wavelength of the incident radiation is short compared with the absorption edge wavelength. In X-ray diffraction this is normally true for the light atoms, but not for the heavy atoms. If the electrons in an atom can be regarded as free electrons, the atomic scattering factor is a real quantity  $f$  because the electron cloud is centrosymmetrical. If they are not free electrons, the atomic scattering factor becomes an imaginary quantity. This complex atomic scattering factor,  $f_{\text{anomalous}}$ , is separated into three parts:  $f_{\text{anomalous}} = f + f' + if''$  where  $f$  is the contribution to the scattering if the electrons were free electrons,  $f'$  is the real part of the correction to be applied for non-free electrons, and  $f''$  is the imaginary correction. The sum  $f + f'$  is the total real part of the atomic scattering factor. Values for  $f$ ,  $f'$ , and  $f''$  are always given in units equal to the scattering by one free electron and are listed in Ref. [1](#). Because the anomalous contribution to the atomic scattering factor is mainly due to the electrons close to the nucleus, the value of the correction factors diminishes only slowly as a function of the scattering angle, much slower than for  $f$ .

Anomalous scattering causes a violation of Friedel's law:  $I(h\ k\ \ell)$  is no longer equal to  $I(h\ k\ \bar{\ell})$ . This fact can be used profitably to elucidate the absolute configuration of the structure determined, and it can assist in the phase angle determination in the single-wavelength isomorphous replacement method. It can also be exploited for phase angle determination by combining information from different wavelengths (MAD). The principle is based on the wavelength dependence of  $f'$  and  $f''$ .

Hendrickson was the first to apply this to a protein structural determination ([2](#)). The X-ray data can be processed in two different ways:

1. the classical way, advocated by Hendrickson
  2. the simplified way in which the problem is treated as if it were isomorphous replacement.
- Hendrickson split the anomalous from the nonanomalous scattering and derived relationships between them. With a minimum of two wavelengths (but in practice at least three), he can find the protein phase angles. In the simplified way, data are collected at three different wavelengths,  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  ([3](#)):

$\lambda_1$  is regarded as the reference wavelength (remote from the absorption edge).

$\lambda_2$  has a maximum difference in  $f'$  compared with  $\lambda_1$  and this difference is used exactly as in isomorphous replacement. This results in two possible phase angles for the protein (see [Isomorphous Replacement](#)).

$\lambda_3$  has a maximum difference in  $f''$  and with this information the choice can be made between the two possible phase angles.

## Bibliography

1. International Union of Crystallography (1995) *International Tables for Crystallography*, Vol. C (A. J. C. Wilson, ed.) Kluwer, Dordrecht, Boston, London.
2. W. A. Hendrickson, J. L. Smith, R. P. Phizackerley, and E. A. Merritt (1988) *Proteins* **4**, 77–88.
3. V. Ramakrishnan et al. (1993) *Nature* **362**, 219–223.

### Suggestions for Further Reading

4. W. A. Hendrickson (1991) Determination of macromolecular structures from anomalous differences of synchrotron radiation, *Science* **254**, 51–58.
5. F. K. Athappilly and W. A. Hendrickson (1995) Structure of the biotinyl domain of acetyl-coenzyme A carboxylase determined by MAD phasing, *Structure* **3**, 1407–1419.
6. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Multivalents

A multivalent occurs following the pairing in **meiosis** of more than one [homologous chromosome](#). In meiosis, the [chromosomes](#) pair along their length during the **prophase** of the first meiotic division. This process of **synapsis** is associated with the assembly of a [synaptonemal complex](#). When the synaptonemal complex extends from end to end of each of the paired chromosomal homologues, the final structure is called a **bivalent**. Multivalents that involve the pairing of more than two homologous chromosomes occur in [polyploid](#) organisms, especially in **plants**. If all of the chromosomes do not double in number during polyploidization, then partial pairing among homologous chromosomes occurs. Some bivalents form, but multivalents, containing three or more chromosomes connected by synaptonemal complexes and univalents with a single chromosome, also form. The aberrant distribution of chromosomal material at the end of the first meiotic division generates many **aneuploid** gametophytes. Thus many polyploid plants reproduce parthenogenetically because of the large number of inviable **gametes** that would otherwise result.

## Mutagen

A mutagen is any chemical that increases the [mutation](#) rate above the spontaneous level in a cell. Few, if any, mutagens induce only one type of event, with all those mutagens that have been well-studied to date showing a spectrum of changes, sometimes considered diagnostic of that particular mutagen or mutagen class. The most extensively studied mutagens are agents that modify DNA by alkylation reactions, either directly or after metabolic conversion to reactive intermediates. This is a diverse group of chemicals for which carcinogenic potencies in rodents vary over a 10,000-fold dose range (1) (see [Dimethyl Sulfate \(DMS\)](#) and [Ethyl Methane Sulfonate \(EMS\)](#)). There is a wealth of literature on the chemistry of formation of DNA-[carcinogen](#) adducts, consideration of the potential of individual adducts for mutagenicity, and mechanisms of **DNA polymerase** fidelity and mutation in the presence of adducts. Whether or not a particular genetic lesion of DNA is mutagenic depends on the nature of the lesion, whether or not it is repaired, and, if so, whether the repair is accurate, the possibilities of replicating DNA at the lesion, and also the relative timing of [DNA replication](#). An exogenous mutagen increases the number of lesions over the background rate of endogenous lesions,

which is known to be high already. Thus, Ames et al. (2) calculated that there are  $10^6$  oxidative lesions present in a normal rodent cell.

There is increasing recognition that certain chemicals can cause mutations through other types of mechanisms. Ferguson and Baguley (3, 4) reviewed the inhibition of topoisomerase enzymes as a mutagenicity mechanism (see [DNA Topology](#)). Turner and Denny (5) reviewed the mutagenic properties of ligands that can bind noncovalently to the minor groove of DNA and span a range of DNA sequence selectivities. Kunz et al. (6) considered the role of imbalances in the pools of the DNA precursors in the genotoxicity of certain DNA-damaging agents.

### Bibliography

1. E. W. Vogel and M. J. Nivard (1994) *Mutat. Res.* **305**, 13–32.
2. B. N. Ames, M. K. Shigenaga, and L. S. Gold (1993) *Environ. Health Perspec.* **101** Suppl. **5**, 35–44.
3. L. R. Ferguson and B. C. Baguley (1996) *Mutat. Res.* **355**, 91–101.
4. L. R. Ferguson and B. C. Baguley (1994) *Environ. Mol. Mutagen.* **24**, 245–261.
5. P. M. Turner and W. A. Denny (1996) *Mutation Res.* **355**, 141–169.
6. B. A. Kunz, S. E. Kohalmi, T. A. Kunkel, C. K. Matthews, E. M. McIntosh, and J. A. Reidy (1994) *Mutation Res.* **318**, 1–64.

### Mutagenesis

Mutagenesis is the study of [mutation](#), which is classically defined as an abrupt and heritable genetic change, induced by chemical or physical agents known as **mutagens**. Practically, a mutation means any change in the **DNA** or the [chromosomes](#) that can be transmitted to descendant cells. These changes may affect a single **nucleotide**, several nucleotides in a **gene**, several genes, large chromosome segments, or whole chromosomes. Crow et al. (1) adopted the following terms:

1. *Gene mutation* affects a single gene. Heritable changes may be due to [base-pair substitution](#) mutations (*transitions* or *transversions*) or **frameshift** events in which deletions or additions of one or a few nucleotide base pairs leads to a change in the [reading frame](#) of the [genetic code](#).
2. *Chromosomal mutation* affects blocks of genes in one or more chromosomes. Heritable changes may be due to large deletions, insertions, or translocations of blocks of DNA.
3. *Genome mutation* affects the number of chromosomes without altering the chromosomal structure.

### Bibliography

1. J. F. Crow, S. Abrahamson, C. Denniston, D. Hoel, et al. (1982) *Identifying and Estimating the Genetic Impact of Chemical Mutagens*. National Academy Press, Washington, D.C.

### Suggestions for Further Reading

2. E. C. Freidberg, G. C. Walker, and W. Seide (1995) *DNA Repair and Mutagenesis*, ASM Press, Washington, D.C. An excellent and detailed modern exposition of the mutation process.
3. T. A. Lindahl (ed.) (1996) *Cancer Surveys 28: Genetic Instability in Cancer*, Cold Spring Harbor Press for the Imperial Cancer Research Fund, Cold Spring Harbor, NY. Invited reviews from

different authors on various topics, including mutation, mismatch repair, mutator genes, and endogenous DNA damage.

## Mutant

A mutant is an organism or cell with a mutation in its [genome](#). Classically, a mutant was defined as an individual with an observable **phenotype**. For example, in **microorganisms**, if the **wild-type** is **prototrophic** (that is, it grows on a defined minimal medium), a mutant is an individual that needs a growth factor (such as an [amino acid](#)) added to the medium to support growth. The observable phenotype can be morphological or biochemical changes. However, modern molecular biological techniques detect changes in the nucleotide sequence that may or may not affect the metabolism or morphology of the organism. Thus, the phenotype may be detected only as a change in the **restriction enzyme's** digestion pattern or a change in the DNA or RNA sequence. When classifying an individual as a mutant, it is necessary to have a defined wild-type for comparison. Differences in the nucleotide sequence within a population that do not result in an altered phenotype are considered **polymorphisms**. The word mutant is also commonly used as an adjective, for example, in mutant protein and mutant plasmid.

## Mutation Load

The mutation load represents the extent of genetic death in a population due to deleterious and . Genetic death can include (1) actual death of individuals before reproduction, (2) sterility, (3) lack of mating opportunities, and (4) other factors that reduce reproduction in comparison to other genotypes. Every population is exposed continuously to newly arising mutations, the average effect of which lowers the fitness of the population in the environment where it lives (1). This decrease in fitness can be considered as the price paid by a species for its capacity to evolve. Increases in mutation rates will result in an increased mutation load and higher rates of genetic death.

In general, the change of average fitness associated with maintaining the variability in a population has come to be called the “genetic load.” Muller first used the term “load” in assessing the impact of mutation on the human population (2). In population genetics and the study of molecular , it is used as a measure of the amount of natural selection associated with a certain amount of genetic variability. The mutation load is one of genetic loads. The operational definition of the genetic load ( $L$ ) is given as  $L = (W_{max} - W)/W_{max}$ , where  $W_{max}$  is the fitness of the reference population, usually taken as the one with maximum fitness, and  $W$  is the average fitness over all individuals in the population. In a random-mating population, population genetics theory has shown that the mutation load can be expressed as  $L = u$  or  $L = 2u/(1 + u)$  in the cases of complete recessive and no dominance, respectively, where  $u$  is the mutation rate.

## Bibliography

“Mutation Load” in , Vol. 3, p. 1552, by T. Gojobori; “Mutation Load” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. J. F. Crow and M. Kimura (1970) *An Introduction to Population Genetics Theory*, Harper and Row, P: New York.
2. H. J. Muller (1950) *Am. J. Hum. Gen.* **2** 111–176.

## Mutations

A mutation is a permanent change in the sequence of the heritable genetic material of a cell or organism. Any change in the sequence of the [genome](#) that is transmitted is a mutation. Transient changes to the genome that do not alter the primary sequence, such as **methylation** patterns, are transmittable, but are not mutations. To recognize a mutant, one must establish a wild-type **phenotype** or **genotype**. Forward mutations are changes from the wild-type phenotype, whereas **reversions** are changes that restore the wild-type. The organism that has a mutation is a [mutant](#), and the process that leads to the mutation is [mutagenesis](#). Mutations are rare events, but these rare events have provided scientists with the tools to dissect the inner working of cells. Mutations are usually thought of as deleterious because they cause diseases, such as birth defects and cancer, but mutations can be beneficial, and they provide the variation for [evolution](#).

Mutations are generally grouped into the categories listed in Table 1. Point mutations refer to small changes, such as the substitution of one base for another or the insertion of one or two bases. A base change that converts a **purine** to the other purine, or a **pyrimidine** for the other pyrimidine, is a [transition mutation](#). Transitions maintain the purine/pyrimidine axis of the nucleic acid **double helix**. [Transversion mutations](#) substitute a purine for a pyrimidine or a purine for a pyrimidine, thus reversing the purine/pyrimidine axis. Transitions are more common than transversions. When a point mutation occurs within the region of a gene coding for a [protein](#), the mutation has one of three effects on the gene product. It can be a [silent mutation](#), it can change the amino acid ([missense mutation](#)) or, it can cause a premature [stop codon](#) ([nonsense mutation](#)).

**Table 1. Categories of Mutations**

| Mutation Class    | Definition   |
|-------------------|--|
| Point mutations   |  |
| Base substitution | Any change in a single base of DNA (or RNA if an RNA genome)   |
| Transition        | Substitution of a purine for a purine or a pyrimidine for a pyrimidine on the same strand                        |
| Transversion      | Substitution of a purine for a pyrimidine or a pyrimidine for a purine on the same strand                        |
| Missense          | A mutation that causes a change in the codon, substituting one amino acid for another                            |
| Nonsense          | A mutation that changes a codon (sense) to a stop codon (nonsense)   |
| Frameshift        | A mutation that changes the reading frame of a transcript by insertion or deletion of $n$ bases, where $n$ is an |

|                 |  |
|-----------------|--|
|                 | integer  |
| Large mutations |  |
| Deletion        | Loss of a large number of bases; generally used for loss of more than four bases     |
| Insertion       | Gain of a large number of bases; generally used for the gain of more than four bases |
| Duplication     | Gain of an additional copy of a sequence already in the genome                       |
| Inversion       | Reversal of a segment of DNA without the gain or loss of bases                       |

---

The other classes of mutations are deletions, insertions, and inversions. Deletions are the loss of bases, as small as one or a few bases or as large as several kilobases encompassing several genes, which can have dramatic effects on the phenotype of the mutant. Small deletions and duplications can be the result of simple slippage during [DNA replication](#), most likely mediated by repeated bases. [Frameshift mutations](#) are the insertion or deletion of  $n$  bases (where  $n$  is an integer), resulting in a shift in the translational reading frame. Technically, a frameshift occurs only within the coding region of a gene, but the term is commonly used for any insertion or deletion of  $n$  bases. Large deletions result from [recombination](#) between **homologous** or nonhomologous DNA sequences. Large rearrangements are also mediated by genetic [mobile elements](#), such as **transposons**, **viruses** and **group I or II introns**. If a viral genome inserts at a specific point in the host genome (for example, at the **att site** for bacteriophage **lambda** in the *Escherichia coli* genome), the insertion is not considered a mutation. However, an illegitimate insertion of a viral genome is a mutation.

Mutations are either spontaneous or induced. One source of spontaneous mutations is replication errors. The fidelity of DNA replication is high because of the combination of base incorporation specificity, proofreading, and postreplicative repair. Most replicative **DNA polymerases** insert an incorrect base at the rate of one per  $10^4$  to  $10^5$  base pairs replicated. If a typical bacterial gene has 1000 bases, the rate of misincorporation would be 1 incorrect base per 10 to 100 genes. Most replicative DNA polymerases have associated **exonucleases** that are the “proofreading” function of the enzyme complex. If the polymerase incorporates an incorrect base, this misincorporated base is recognized by the proofreader, and the exonuclease degrades the nascent strand in the reverse order of polymerization. Then the polymerase has another chance to insert the correct base. Proofreading increases the fidelity of replication 10- to 200-fold. If misincorporated bases are not corrected by proofreading, postreplicative [mismatch repair](#) provides a third level of defense against mutation. Because DNA replication is **semiconservative**, a misincorporated base in the nascent strand creates a mismatch with the template strand. These mismatches are recognized by mismatch recognition proteins that correct the newly replicated strand in favor of the template strand. This postreplicative mismatch repair increases replication fidelity an additional 10- to 1000-fold (see [Mismatch Repair](#)). Spontaneous mutations also result from endogenous DNA damage. Cytosine spontaneously deaminates to uracil, and **5-methyl-cytosine** deaminates to thymine. If not repaired, these deaminations cause a C to T transition mutation.

It is thought that most spontaneous mutations commonly arise as unavoidable replication errors during cell growth. However, in the 1950s Ryan reported the occurrence of mutations in *E. coli* cells that were not dividing or, apparently, replicating their genomes (1, 2). More recently, Cairns et al. (3) described other instances of mutations in nondividing *E. coli* under nonlethal selective conditions and presented evidence that the only mutations recovered were those that allow the cells to grow. This phenomenon has been called “directed,” or “adaptive,” or “postselection” mutation and has

been observed now in a variety of microorganisms (4). Although the cells undergoing adaptive mutation are not dividing, the mutations appear to be caused by DNA synthesis initiated independently of **chromosomal** replication. For example, [recombination](#) functions are required in one case, suggesting that recombination intermediates initiate at least some limited DNA synthesis in nondividing cells. In other cases, the mutations are caused by replication past endogenous DNA lesions (5). The “adaptiveness” of the mutations is not so clear, however, as recent evidence indicates that with sensitive tests nonselected mutations are found. One current theory is that a small subpopulation of the cells under selection undergo a period of transient mutation and that these give rise to the apparently adaptive mutations (6).

In addition to spontaneous mutations, genetic changes are also induced by exogenous chemical and physical agents that alter or damage bases ([mutagens](#)). Some mutagens covalently modify DNA. [Alkylation](#) is a common type of covalent modification, in which an alkyl group is attached to the base. Methylation of the N7 position of guanine is generally not mutagenic, whereas methylation of its O6 position disrupts the **hydrogen bonding**, allowing T to pair with G and giving rise to a G to A transition mutation. Physical agents, such as heat, UV light and radiation are also mutagens. UV light produces cyclobutane dimers between adjacent pyrimidine bases (see **Thymine dimer**). Cells have multiple systems to deal with these dimers, including (1) white-light **photoreactivating** enzymes that split the cyclobutane ring and (2) excision-repair enzymes that recognize the distortion of the helix and remove the dimer and adjacent bases, leaving a gap that is filled in by a DNA polymerase. Mutations in excision-repair genes in humans cause xeroderma pigmentosum, a skin disorder that makes people very sensitive to sunlight.

As mentioned above, genomic changes also occur by genetic [recombination](#). Inappropriate recombination leads to deletions and duplications in the genome. Although these large alterations are devastating to the organism, [gene duplications](#) may allow evolution by providing “extra” gene copies in which a new genetic function can evolve. Homologous recombination between sister chromosomes changes the architecture of chromosomes in diploid organisms, changing the linkage of **alleles**.

Because mutations are rare, the study of mutagenic mechanisms is difficult. To increase the numbers of mutants, geneticists design selection procedures to find mutants of interests (see **Selection**). For example, if one is studying the types of mutations induced by a chemical mutagen, it would be helpful to know what change is being induced. Cupples and Miller (7, 8) developed a series of *E. coli* strains incapable of utilizing lactose as a carbon source but that revert to lactose utilization by a specific base substitution or frameshift mutation. If these strains are subjected to a mutagenic agent and then plated onto lactose minimal media, the number of mutations may be high in one strain, but low in the other strains. Because the base change needed to revert is known, the specificity of the mutagen can be determined. This is a brilliant example of how geneticists design selections to understand mutagenesis.

## Bibliography

1. F. J. Ryan (1955) *Genetics* **40**, 726–738.
2. D. Nakada and F. J. Ryan (1961) *Nature (London)* **189**, 398–399.
3. J. Cairns, J. Overbaugh, and S. Miller (1988) *Nature (London)* **335**, 142–145.
4. P. L. Foster (1993) *Annu. Rev. Microbiol.* **47**, 467–504.
5. B. A. Bridges (1997) *BioEssays* **19**, 347–352.
6. P. L. Foster (1998) *Genetics* **148**, 1453–1459.
7. C. G. Cupples, M. Cabrera, C. Cruz, and J. H. Miller (1990) *Genetics* **125**, 275–280.
8. C. G. Cupples and J. H. Miller (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5345–5349.

## Suggestions for Further Reading

9. E. C. Friedberg, G. C. Walker, and W. Siede (1995) *DNA Repair and Mutagenesis*, American Society for Microbiology, Washington, D.C.
10. A. Griffiths, J. H. Miller, D. Suzuki, R. Lewontin, and W. Gelbart (1996) *An Introduction to Genetic Analysis*, 6th ed., W. H. Freeman, New York.
11. B. Singer and D. Grundberger (1983) *Molecular Biology of Mutagens and Carcinogens*, Plenum Press, New York.
12. J. R. Beckwith and T. Silhavy (1992) *The Power of Bacterial Genetics: A Literature-Based Course*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

## Mutator Genes

Defects in mutator genes lead to organisms with spontaneous [mutation](#) rates that are significantly elevated over the normal rate. The earliest report of this type of defect was from Treffers et al. (1), who discovered the *mutT Escherichia coli* mutator strain. This mutator is quite specific in increasing the frequency of A.T → G.C transversion mutations, even though these are only one of a number of mutation pathways that occur normally. The *mutT* strain frequently produces 10<sup>3</sup>-fold increases in mutation rates. Similar genes have been described in many different organisms, each with their own spectrum of increased mutagenic activity. For example, the X-ray-sensitive yeast mutator strain described by von Borstel et al. (2) shows an increased frequency of [frameshift mutations](#). Inactivation of the yeast RAD27 gene leads to a UV-sensitive phenotype and an elevated rate of duplications of 5–108 [base pairs](#) of DNA located between repeated sequences of 3–12 base pairs (3).

It has been argued that mutator genes play an important role in adaptive [evolution](#) (4). As many as 1% of natural bacterial isolates have defective mutator genes, some of them leading to very high mutation rates. Taddei et al. (4) developed models of clonal populations adapting to a new environment. Their results suggest that strong mutator genes (such as those that increase mutation rates by 10<sup>3</sup>) will significantly accelerate adaptation to a changed environment. Although early studies on mutator genes focused on bacteria and yeast, mutator genes are increasingly being recognized as important in human cancer. Many of these mutants are defective in [DNA repair](#) systems, including DNA [mismatch repair](#), a repair pathway for oxidative lesions and mistranslation (5). For example, the human MSH2 (hMSH2) protein is involved in the recognition of mismatched nucleotides during postreplication mismatch repair. Alterations to the hMSH2 gene are found in approximately 60% of hereditary nonpolyposis colon cancer cases. Loss of the function of this gene leads to a mutator phenotype and may be involved with the multiple mutations required for multistage carcinogenesis. (see [Carcinogen](#)).

## Bibliography

1. H. P. Treffers, V. Spinelli, and N. O. Belser (1954) *Proc. Natl. Acad. Sci. USA* **40**, 1064–1071.
2. R. C. von Borstel, D. E. Graham, K. J. La Brot, and M. A. Resnick (1968) *Genetics* **60**, 233.
3. D. X. Tishkoff, N. Filosi, G. M. Gaida, and R. D. Kolodner (1997) *Cell* **88**, 253–263.
4. F. Taddei, M. Radman, J. Maynardsmith, B. Toupance, P. H. Gouyon, and B. Godelle (1997) *Nature* **387**, 700–702.
5. J. H. Miller and M. Michaels (1996) *Gene* **179**, 129–132.



## Myb Oncogene

The *myb* oncogene, the transforming gene of the avian myeloblastosis virus, was isolated in 1941 from a chicken tumor (1). The normal cellular counterpart of this **oncogene**, *c-myb*, codes for a nuclear **transcription factor** that contains an amino-terminal **DNA-binding** domain, a central transactivation domain, and a C-terminal negative-regulatory domain (2, 3). Both the v-Myb and c-Myb proteins bind to DNA in a sequence-specific manner and transactivate the **transcription** of genes that contain Myb-binding sequences in their **promoter/enhancer** elements. A comparison of the v-*myb* and *c-myb* sequences showed that the viral oncogene arose as a result of deletions in the 5' and 3' portions of the coding sequences, which results in deleting a portion of the DNA-binding domain and the entire negative-regulatory domain (2). Although the deletion in the DNA-binding domain does not affect the viral protein's ability to bind to DNA, deleting the C-terminal negative-regulatory domain results in enhanced transcriptional transactivation by the v-Myb protein, which in turn enhances the proliferative activity of the virally infected myeloid cells (4).

*c-myb* is expressed predominantly in **hematopoietic** cells and is readily induced by treating these cells with **interleukins** or **mitogens**. Studies with **antisense** oligonucleotides demonstrate that expression of the *c-myb* gene product is essential for the proliferative potential of several myeloid and **T-cell** lines. In addition, studies with hematopoietic cell lines suggest that the terminal differentiation of these cells is accompanied by down-regulation of *myb* gene expression and constitutive expression of *c-myb* blocks their terminal differentiation. Homozygous null *c-myb* mutant mice die *in utero* because of defects in fetal hepatic hematopoiesis, confirming an essential role for *c-myb* in hematopoiesis (5).

Amplification of the *c-myb* gene has been observed in several human cancers, including the acute myelogenous leukemia (AML), chronic myelogenous leukemia (CML), acute lymphocytic leukemia (ALL), T-cell leukemias, melanomas, and colon carcinomas. Recent findings suggest that virtually all **estrogen receptor**-positive breast carcinomas express high levels of *c-myb*. This gene might also play a critical role in **estrogen**-induced proliferation of mammary epithelial cells, and its overexpression may be associated with the malignant growth of these cells.

The **genomes** of higher organisms contain two *myb*-related genes, which have been named A-*myb* and B-*myb*. B-*myb* is expressed ubiquitously, but A-*myb* is expressed almost exclusively in testis, pregnant breast, and the germinal centers of the spleen. These two genes are highly **homologous** to *c-myb* and contain virtually identical **DNA-binding** domains. Like c-Myb, they contain transactivation and C-terminal regulatory domains and act as activators of **transcription**. It is possible that these two proteins function similarly to c-Myb in tissues where they are expressed. Further evidence is that A-*myb*-null mutant mice exhibit defects in spermatogenesis, breast development, and spleen germinal center development (6).

### Bibliography

1. W. J. Hall, C. W. Bean, and M. Pollard (1941) *Am. J. Vet. Res.* **2**, 272–279.
2. M. A. Baluda and E. P. Reddy (1994) *Oncogene* **9**, 2761–2774.
3. K. Weston and J. M. Bishop (1989) *Cell* **58**, 85–93.
4. G. Patel, B. Kreider, G. Rovera, and E. P. Reddy (1993) *Mol. Cell. Biol.* **13**, 2269–2276.
5. M. Introna, M. Luchetti, M. Castellano, M. Arsura, and J. Gole (1994) *Semin. Cancer Biol.* **5**, 113–124.
6. A. Toscani, R. V. Mettus, R. Coupland, H. Simpkins, J. Litvin, J. Orth, K. S. Hatton, and E. P. Reddy (1997) *Nature* **386**, 713–717.

## Myc Oncogene

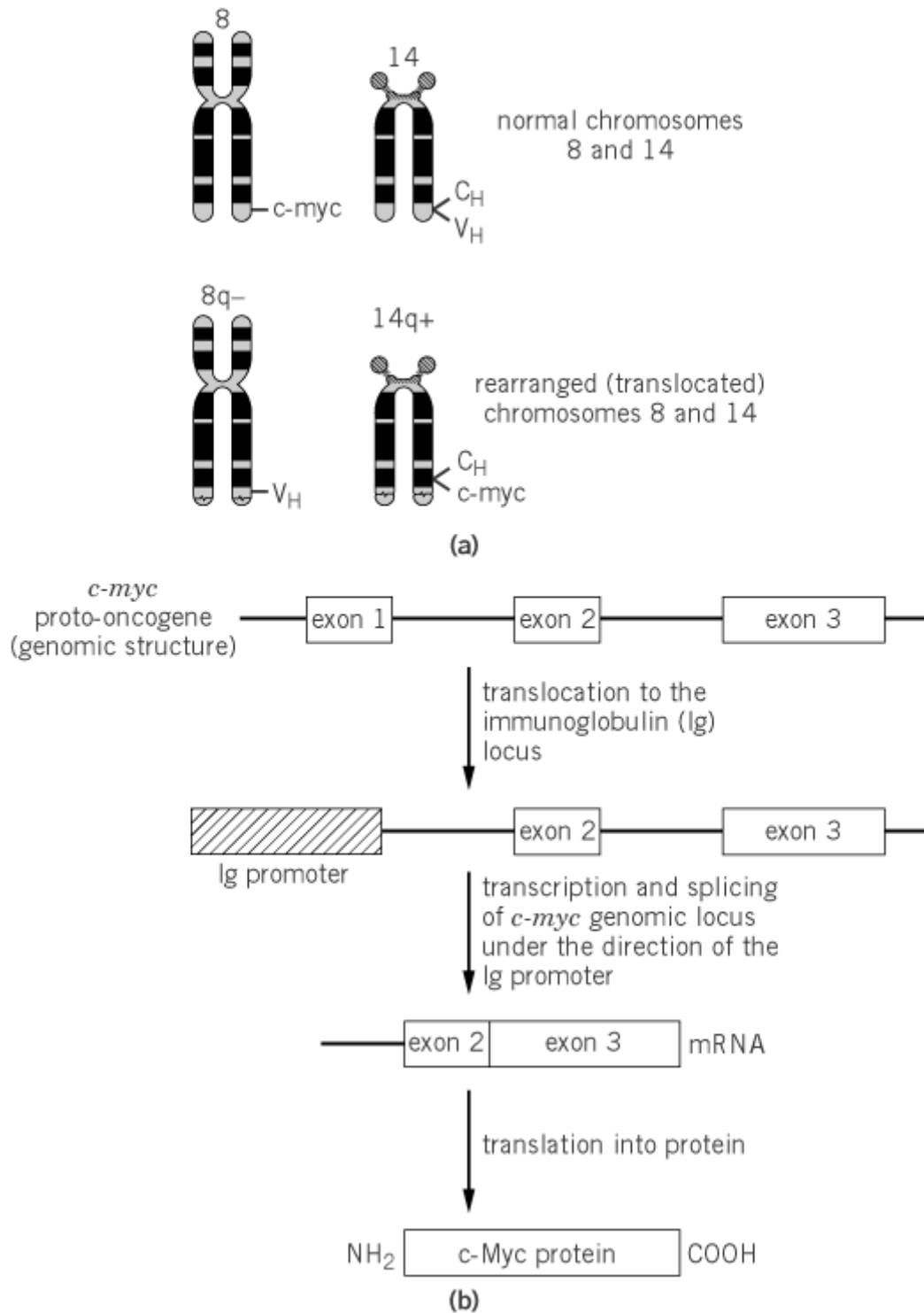
The *myc* oncogene was originally identified as the transforming gene of an acute transforming virus called MC29, which was isolated from a chicken that had spontaneous myelocytomatosis (1). The *myc*-related sequences were subsequently identified in three other independently derived chicken retroviral isolates called CMII, OK10, and MH2. Comparison of the MC29-derived v-*myc* and chicken c-*myc* sequences showed that the viral **oncogene** contains the entire coding sequence of the chicken c-*myc* gene fused to the viral gag sequence, thus producing a gag-*myc* fusion protein (2, 3). Although the fusion of the *myc* sequences to the gag gene results in elevated levels of Myc protein expression, deletion analysis of the proviral genome suggests that elevated expression of *myc* sequences alone is adequate for the transforming activity of this oncogenic virus. Thus, structural alterations play a less important role in the oncogenic activity of this gene than in other oncogenic viruses, and overexpression of the encoded protein product alone is adequate to induce transformation of appropriate target cells.

Avian and mammalian genomes code for three members of this family, which are known as c-*myc*, N-*myc*, and L-*myc* (4). All three genes code for approximately 50-kDa proteins, which belong to the **helix-loop-helix/leucine zipper** superfamily of proteins. The Myc proteins themselves do not bind to DNA with high affinity. However, they form heterodimers with a common partner termed Max, which has high affinity for the **consensus sequence** CACGTG, often referred to as the E Box. Myc-Max heterodimers bind to the E box sequence and transactivate [transcription](#) of the genes that contain this binding sequence.

c-*myc* expression in normal cells is induced by the interaction of [growth factors](#) with their cognate **receptors**, and this expression is essential for [cell-cycle](#) progression and cell proliferation. Downregulation of *myc* expression in **hematopoietic** cells, correlates with their terminal differentiation, and constitutive expression of *myc* often results in a block to the differentiation process. Studies with animal and human tumors show that the *myc* gene is frequently activated by one of three mechanisms: (1) proviral insertion; (2) gene amplification; and (3) chromosomal translocation. All three mechanisms result in elevated expression of the Myc protein. Although proviral insertion into the *myc* locus occurs in avian and murine tumors, activation of this gene in human tumors occurs predominantly by gene amplification and chromosomal translocation. Thus, several small-cell lung carcinomas, breast carcinomas, cervical carcinomas, and neuroblastomas have one of the three *myc* family members amplified. In Burkitt's lymphomas, the chromosomal [translocation](#) between [chromosomes](#) 8 and 14 leads to juxtapositioning of the c-*myc* oncogene with the regulatory sequences of the [immunoglobulin](#) (Ig) genes (Fig. 1). Immunoglobulin loci code for **antibodies**, which are produced solely by B lymphocytes. When these genes are expressed, the Ig regulatory sequences direct high levels of mRNA transcription of the Ig genes located next to them. When the c-*myc* oncogene is placed under the transcriptional control of the Ig regulatory elements, however, the B cell abnormally produces extremely high levels of c-*myc* mRNA and protein (instead of antibody), which induces a malignant phenotype.

**Figure 1.** Activation of the c-*myc* oncogene as a consequence of translocation in Burkitt's lymphoma. **(a)** During translocation, the c-*myc* gene (located on chromosome 8) translocates to chromosome 14 at a site next to the immunoglobulin heavy chain (C<sub>H</sub>) locus. **(b)** This translocation results in juxtaposing the immunoglobulin regulatory sequences upstream of the c-*myc* gene. Because the immunoglobulin promoter positively regulates high levels of

transcription in B lymphocytes, the *c-myc* gene (instead of the immunoglobulin gene) is constitutively transcribed and translated in the cell.



Recent studies have also implicated the *myc* gene in [apoptosis](#) (5). Thus, it has been found in hematopoietic cells that constitutive expression of *myc* during the withdrawal of **cytokines** leads to acceleration of apoptosis. Similarly, constitutive expression of *myc* in Rat-1 fibroblasts causes apoptosis when proliferation is inhibited by removing serum from the growth medium (see [Serum Dependence](#)) or by adding factors, such as [tumor necrosis factor](#) to the medium. Myc-mediated apoptosis is associated with elevated expression of [cyclin A](#).

## Bibliography

1. I. X. Vanov, Z. Mladenov, S. Nedyalkov, T. G. Todorov, and M. Yakimov (1964) Bull. Inst. Pathol. Comp. Anim. **10**, 5–38.
2. E. P. Reddy, R. K. Reynolds, D. K., Watson, R. A. Schultz, J. Lautenberger, and T. S. Papas (1983) Proc. Natl. Acad. Sci. USA **80**, 2500–2504.
3. D. K. Watson, E. P. Reddy, P. H. Duesberg, and T. S. Papas (1983) Proc. Natl. Acad. Sci. USA **80**, 2164–2150.
4. N. Schreiber-Agus, L. Alland, R. Muhle, J. Goltz, K. Chen, L. Stevens, D. Stein, and R. A. DePinho (1997) Curr. Top. Microbiol. & Immunol. **224**, 159–168.
5. G. T. Williams (1991) Cell **65**, 1097–1098.

## Mycoplasma

Mycoplasma is the vernacular name for the group of naturally occurring wall-less eubacteria. With one apparent exception, these are free-living microorganisms and part of the normal flora and significant pathogens of humans, animals, plants, and insects. The possible exception are the phytoplasmas (originally called “mycoplasma-like organisms”), which infect plants and insects and may be obligate intracellular parasites. Besides their clinical importance, mycoplasmas are of interest because of their small genomes: some mycoplasmas have the smallest known cellular genomes, perhaps defining the minimal number of genes a living cell can have.

The first isolation of a mycoplasma was in 1898 with the cultivation of the causative agent of bovine pleuropneumonia. This strain and subsequent similar isolates were named pleuropneumonia-like organisms (PPLO) and later renamed the genus *Mycoplasma*. Current mycoplasma taxonomy consists of a single class, *Mollicutes*, containing four orders (Table 1): *Mycoplasmatales*, sterol-requiring human and animal mycoplasmas; *Entomoplasmatales*, plant and insect mycoplasmas; *Acholeplasmatales*, non-sterol-requiring mycoplasmas that grow in a variety of habitats; and *Anaeroplasmatales*, obligate anaerobic mycoplasmas. There is also a group of plant and insect mycoplasmas, Phytoplasma, that cannot be given formal taxonomic status because it has not been possible to grow these microorganisms in axenic culture.

**Table 1. Taxonomy of the Class *Mollicutes*<sup>a</sup>**

---

|   |
|---|
| Order: <i>Mycoplasmatales</i>                     |
| Family: <i>Mycoplasmataceae</i>                   |
| Genus: <i>Mycoplasma</i> (102 recognized species) |
| Genus: <i>Ureaplasma</i> (6 recognized species)   |
| Order: <i>Entomoplasmatales</i>                   |
| Family: <i>Entomoplasmataceae</i>                 |
| Genus: <i>Entomoplasma</i> (5 recognized species) |
| Genus: <i>Mesoplasma</i> (12 recognized species)  |
| Family: <i>Spiroplasmataceae</i>                  |

Genus: *Spiroplasma* (33 recognized species)  
Order: *Acholeplasmatales*  
Family: *Acholeplasmataceae*  
Genus: *Acholeplasma* (13 recognized species)  
Order: *Anaeroplasmatales*  
Family: *Anaeroplasmataceae*  
Genus: *Anaeroplasma* (4 recognized species)  
Genus: *Asteroleplasma* (1 recognized species)  
Undefined taxonomic status: Phytoplasma<sup>b</sup>

---

<sup>a</sup> Based on data in Ref. 2.

<sup>b</sup> There is no formal taxonomic classification for Phytoplasma. They were originally called mycoplasma-like organisms (MLOs).

The defining phenotype of members of the Class *Mollicutes* is the absence of any type of cell wall structure. Cell wall genes appear to have been lost during the degenerate evolution of mycoplasmas from Gram-positive eubacteria, and mycoplasmas are bounded only by a single cell membrane. The properties of the mycoplasma genera are shown in Table 2.

**Table 2. Properties of mycoplasma genera<sup>a</sup>**

| Genus          | Effect of Oxygen     | Cholesterol Required | Genome Size            | Habitat                          |
|----------------|----------------------|----------------------|------------------------|----------------------------------|
| Acholeplasma   | Facultative anaerobe | No                   | 1,500–1,650 kb         | Animals, some plants and insects |
| Anaeroplasma   | Obligate anaerobe    | Yes                  | 1,500–1,600 kb         | Bovine and ovine rumens          |
| Asteroleplasma | Obligate anaerobe    | No                   | 1,500 kb               | Bovine and ovine rumens          |
| Entomoplasma   | Facultative anaerobe | Yes                  | 790–1,140 kb           | Plants and insects               |
| Mesoplasma     | Facultative anaerobe | No                   | 870–1,100 kb           | Plants and insects               |
| Mycoplasma     | Facultative anaerobe | Yes                  | 580–1,350 kb           | Animals and humans               |
| Phytoplasma    | — <sup>b</sup>       | — <sup>b</sup>       | 640–1,185 <sup>c</sup> | Plants and insects <sup>d</sup>  |
| Spiroplasma    | Facultative anaerobe | Yes                  | 780–2,200 kb           | Plants and insects               |
| Ureaplasma     | Facultative anaerobe | Yes                  | 760–1,170 kb           | Animals and humans               |

<sup>a</sup> Based on data in Ref. 2.

<sup>b</sup> No data available.

<sup>c</sup> These genome sizes are based on DNAs in phytoplasma-infected plant and insect tissues (21).

<sup>d</sup> Phytoplasma may be obligate intracellular parasites.

## 1. Phylogeny

Mycoplasmas arose from the *Streptococcus* branch of the Gram-positive eubacterial phylogenetic tree of species with low DNA G+C contents (1). There was an early split in the mycoplasma branch into two major branches, with the ancestral branch leading to *Acholeplasma*, *Asteroleplasma*, *Anaeroplasm*a, and Phytoplasma branches. The more recent major branch led to *Spiroplasma*, *Entomoplasma*, *Mesoplasma*, *Mycoplasma*, and *Ureaplasma* branches.

The ancestral mycoplasma branch arose from bacteria with relatively small genomes (*Streptococcus* genome sizes are 2,100 kb to 2,300 kb) and probably consisted of microorganisms with 1,500 kb to 1,700 kb genomes. It is likely that these earliest mycoplasmas were similar to *Acholeplasma*, ie, non-sterol-requiring, facultative anaerobes. Further evolution must have led to branches of obligate anaerobes (*Asteroleplasma* and *Anaeroplasm*a) and, more recently, to a Phytoplasma branch with 640 kb to 1,185 kb genomes. Phylogeny of the more recent major branch produced *Spiroplasma*, *Entomoplasma*, *Mesoplasma*, and *Mycoplasma* branches with 580 kb to 2,200 kb genomes. *Ureaplasma* is a urea-requiring branch of *Mycoplasma*. Surprisingly, the smallest genomes (about 600 kb) have evolved on several independent phylogenetic branches, meaning that mycoplasmas have been under selective pressure for genome reductions and repeatedly have explored the lower limit for the genome complexity for free-living organisms.

## 2. Viruses

Mycoplasma viruses have been found infecting *Acholeplasma*, *Spiroplasma*, and *Mycoplasma* species (2). Several viruses infecting *Acholeplasma* and *Spiroplasma* have been propagated and characterized. However, a virus infecting *My. pulmonis* is the only virus infecting a *Mycoplasma* that has been propagated and characterized.

Filamentous *Acholeplasma* and *Spiroplasma* viruses have circular, single-stranded DNA genomes. Their replication appears similar to bacterial filamentous phages, with a nonlytic infection cycle and progeny virus release by extrusion through infected cell membranes. The *Acholeplasma* viruses are short, bullet-shaped virions with genomes about 4.4 kb, and the *Spiroplasma* viruses are long, filamentous virions with genomes about 8.3 kb. The *Spiroplasma* virus SpV1 genome sequence (8,272 bases) has been determined (3).

The only known isometric mycoplasma virus is *Spiroplasma* virus SpV4, with a genome of 4,421 bases (4).

Two enveloped, quasi-spherical, *Acholeplasma* viruses have been reported: L172 with a circular single-stranded DNA genome (about 14.0 kb) and L2 with a circular double-stranded DNA genome. The L2 genome sequence (11,965 bases) has been analyzed (5). L2 is a temperate virus with an unusual infection cycle: productive infection is noncytotoxic with enveloped progeny virions released by budding through infected host cell membranes followed by establishment of lysogeny in most, if not all, infected cells.

Viruses with tailed phage morphology have been reported infecting *Acholeplasma*, *Spiroplasma*, and *Mycoplasma* species. Those that have been propagated all have isometric heads and short tails and contain linear double-stranded DNA ranging in size from 11.3 kb for *Mycoplasma* virus P1 (6) to 39.4 kb for *Acholeplasma* virus L3 (7), with circular permutation and terminal redundancy, cohesive

ends, or terminal inverted repeats and 5'-terminal proteins (2, 6, 7). The genome sequence of Mycoplasma virus P1 (11,660 bases), which infects *My. pulmonis*, has been determined (8).

In addition, the genome sequence of Mycoplasma virus MAV1 (15,644 bases), which infects *My. arthritidis*, has been analyzed (9), although the morphology of this virus is unknown.

### 3. Genome Structure and Organization

#### 3.1. Genome Structure

Mycoplasma genomes have low G + C contents and small sizes (2). The DNA G + C content ranges from 25 to 35 mol%, with the exception of the 40 mol% G + C *My. pneumoniae* genome. In general, coding regions have a higher G + C content than non-coding regions.

Mycoplasma genome sizes range from 580 kb to 2,200 kb. The larger genome sizes are found in organisms on the oldest phylogenetic branches and smaller genome sizes in organisms on more recent phylogenetic branches (Table 2). Physical and genetic genome maps have been constructed for several Mycoplasma species, *U. urealyticum*, and *S. citri* (2).

The *My. genitalium*, *My. pneumoniae*, *My. pulmonis*, and *U. urealyticum* genome sequences have been determined and are 580 (10), 816 (11), 964 (12), and 751 kb (13), respectively. Analyses of these data show that genome reductions during evolution of mycoplasmas from Gram-positive bacteria mostly involved losses of genes for cell wall synthesis and anabolic pathways.

#### 3.2. Repetitive Sequences, Transposable Elements, and Genome Rearrangements

Several types of genetic elements have been reported in *Mycoplasma* (14). The repetitive sequences that have been found are related to insertion sequences (IS) and are involved in adhesin and antigenic variation. Mycoplasma IS elements have significant sequence similarity to the IS3 class of insertion sequences, but only *My. pulmonis* IS1138 has been shown to transpose. Besides transposition, other types of genome rearrangements (ie, insertions, deletions, and DNA inversions) have been reported in *Mycoplasma* and associated with variations in antigen expression and size. A novel type of phenotypic switching in *My. pulmonis* involves regulation of restriction and modification by DNA inversion.

#### 3.3. Plasmids

Mycoplasma plasmids have been reported in *My. mycoides* (14). Two cryptic *My. mycoides* plasmids have been sequenced (1,717 and 1,875 bases) and found to have regions of similarity to each other and to a family of Gram-positive bacterial plasmids that replicate via single-stranded DNA intermediates.

Extrachromosomal DNAs have been reported in *Spiroplasma*. However, it has not been possible to determine which (if any) are plasmids because most *Spiroplasma* cultures contain viruses and many *Spiroplasma* extrachromosomal DNAs have been shown to be *Spiroplasma* virus replicative forms.

### 4. DNA Replication

#### 4.1. Biochemistry

Mycoplasma DNA replication resembles that of eubacteria: the DNA replication complex is membrane-associated (15), and replication is bidirectional from a replication origin (*oriC*) (2). Mycoplasma *oriC* sites contain *DnaA* boxes, and genes in this region have organizational and sequence similarity to Gram-positive bacterial *oriC* sites.

Three DNA polymerases have been identified in *Acholeplasma* and *Spiroplasma*, as in eubacteria, but the number of DNA polymerases in *Mycoplasma* and *Ureaplasma* is uncertain (2). One of the *Acholeplasma* and *Spiroplasma* polymerases has properties similar to eubacterial DNA polymerase I (Pol I), and another has properties similar to eubacterial DNA polymerase III (Pol III). However, it is

not known if the exonuclease activity in mycoplasma Pol III preparations is part of the DNA polymerase holoenzyme or a contaminant. In contrast, genome sequence analyses of *Mycoplasma* have identified two DNA polymerase genes, both with sequence similarity to Gram-positive bacterial Pol III. The larger encodes both DNA polymerase and 3'-5' exonuclease activities, but the smaller only encodes DNA polymerase activity. Other gene products required for DNA replication have been identified in *Mycoplasma* genome sequence analyses.

#### 4.2. Repair

*Acholeplasma*, but not *Spiroplasma* or *Mycoplasma*, has a photoreactivation (light-dependent DNA repair) system (2). In addition, *Acholeplasma* and some *Mycoplasma* species have excision (dark) DNA repair systems, but *Spiroplasma* and other *Mycoplasma* species do not. *Acholeplasma* and *Mycoplasma* are also expected to have RecA-dependent (recombinational) DNA repair because they have functional recA genes and active recombination systems (2). However, no recombinational DNA repair could be demonstrated in *Spiroplasma*, and the recA gene in *Spiroplasma* has been shown to be either truncated or to contain an internal stop codon. Some of the gene products required for DNA repair have been identified in *Mycoplasma* genome sequence analyses, but the number of DNA repair gene products appears much smaller than it does in eubacteria.

#### 4.3. Restriction and Modification

As in eubacteria, mycoplasma restriction and modification were first observed as host-controlled variation in virus plating efficiencies. Type II restriction and modification systems (with relatively simple endonuclease and modification enzymes and cleavage at or near a small palindromic site) have been reported in *Acholeplasma*, *Spiroplasma*, *Mycoplasma*, and *Ureaplasma* (14). In general, these resemble typical bacterial restriction and modification systems. A novel specificity is that of *A. laidlawii* strain JA1, in which DNA containing m<sup>5</sup>Cyt is restricted, independent of the m<sup>5</sup>Cyt-containing sequence, and host cell DNA is protected by the absence of m<sup>5</sup>Cyt nucleotides.

The only mycoplasma Type I restriction and modification system (with multimeric endonuclease and modification enzymes and cleavage at a random site far from the recognition site) has been reported in *My. pulmonis* (16). However, restriction and modification in *My. pulmonis* is regulated in a way not found in other microorganisms; regulation is via DNA inversions that change the cells' restriction and modification properties. A similar genome sequence has been found in *My. pneumoniae*, suggesting this type of system may be common in *Mycoplasma* species.

### 5. Transcription

#### 5.1. Biochemistry

Mycoplasma RNA polymerases resemble eubacterial RNA polymerases, with a core enzyme of a, b, and b' subunits and a s initiation factor (2). However, Mycoplasma genome sequence analyses show only one s gene, in contrast to the six or more s genes in eubacteria.

Mycoplasma transcription is not inhibited by rifampin, an antibiotic that binds the RNA polymerase b subunit and blocks eubacterial transcription, and purified *Acholeplasma* and *Spiroplasma* RNA polymerases are not inhibited by rifampin. Hence, the mutation to rifampin resistance must have occurred early in mycoplasma phylogeny, before the ancestral mycoplasma branch split.

#### 5.2. Promoters, Terminators, and Antiterminators

Mycoplasma promoters are similar to eubacterial promoters (2): -35 and -10 hexamers separated by a 15 to 17 base spacer, and transcription initiation 6 to 8 bases downstream from the -10 site. The mycoplasma -10 sequence is close to the eubacterial -10 consensus sequence (TATAAT), but there is significant diversity between mycoplasma -35 sequences and the eubacterial -35 consensus sequence (TTGACA). A novel type of transcriptional regulation has been found in *My. hyorhinae*, in which the spacer between the promoter -10 and -35 sites for a family of surface proteins is a poly(A) tract, and changes in the number of adenine residues in this tract correlate with changes in surface



protein gene expression and consequent antigenic variation (17).

Mycoplasma studies show transcription termination is at r-independent termination sites (2). Although several termination and antitermination factors have been identified in Mycoplasma genome sequence analyses, no r termination factor has been found, suggesting only r-independent termination is present. Antitermination sequences have been found in the spacer regions between genes in *Mycoplasma* rRNA and tRNA operons.

### 5.3. RNA processing and RNase P

*Mycoplasma* rRNA operons have putative stem-loop structures in the spacer region between 16S and 23S rRNA genes (2). In eubacteria, such structures are substrates for processing by RNase III. Mycoplasma genome sequence analyses have identified an RNase III gene, indicating processing of rRNA primary transcripts in mycoplasmas.

The RNase P gene or RNA have been identified in *Mycoplasma* (2) and found to have ribozyme activity as a tRNA-processing endonuclease. Comparative analyses of mycoplasma and bacterial RNase P RNAs are being used to investigate RNase P structure and function.

## 6. Translation

### 6.1. Biochemistry

Mycoplasma translation resembles that of Gram-positive bacteria (2). The only reported exception is the *My. genitalium tuf* gene, which uses a very different ribosome binding sequence from the Shine-Dalgarno sequences in other mycoplasma mRNAs (18). In small *Mycoplasma* genomes, about 15% of total open reading frames encode gene products needed for translation.

An interesting aspect of mycoplasma translation is the recent finding of a gene encoding 10Sa RNA in *My. pneumoniae* (11). This type of bacterial RNA appears to have both mRNA and tRNA functions and has been proposed to mediate a novel type of “trans translation” protein synthesis.

### 6.2. Codon Usage

Mycoplasma translation initiation is generally at UGA codons, although initiation at GUG and UUG codons also has been found (2). There is extensive codon bias in mycoplasmas, reflecting the low G+C contents of mycoplasma DNAs. Variations in codon usage in *My. pneumoniae* also has been found as a function of G+C content and frequency of expression (11).

There is an exception to the universal genetic code in some mycoplasmas. *Acholeplasma*, like other microorganisms, uses UGG as the tryptophan (trp) codon and UGA as a termination codon, has no tRNA with anticodon CCA, and has a trp-tRNA with anticodon UCA. However, *Spiroplasma*, *Mesoplasma*, *Mycoplasma*, and *Ureaplasma* use both UGG and UGA as trp codons and contain two trp-tRNAs, one with anticodon CCA and one with anticodon UCA. Hence, during mycoplasma phylogeny, the ancestral branch that arose from the Gram-positive bacteria retained the typical bacterial codon usage, with UGG as the only trp codon and UGA as a termination codon. The change in codon usage must have occurred on the more recent phylogenetic branch soon after it diverged from the ancestral branch and presumably involved duplication of the trp-tRNA gene with anticodon CCA and mutation of the anticodon in one of these genes to UCA.

### 6.3. Post-translational Modification

Studies of post-translational processing have involved mycoplasma membrane proteins (2). Many of these proteins have typical Gram-positive bacterial signal peptides, and cleavage processing has been confirmed in a few cases. Mycoplasma membrane proteins also may be modified by covalent addition of fatty acids, phosphorylation, and isoprenylation.

### 6.4. Heat Shock Proteins

Mycoplasmas have a heat shock response similar to that of other eubacteria: increased synthesis of a

specific subset of proteins after a shift-up in temperature (2). Two mycoplasma heat shock proteins have been identified as homologs of eubacterial DnaK and GroEL proteins, two of the most conserved heat shock proteins. The fact that a heat shock system has been conserved in mycoplasmas during their degenerate evolution from eubacteria indicates significant selective advantage for this system for adapting to environmental stress.

## 7. Genetics

### 7.1. Genetic Transfer

Mycoplasma can transfer genetic information by transformation and cell-to-cell contact (14). Many *Acholeplasma*, *Spiroplasma*, and *Mycoplasma* species have been transformed with viral and plasmid DNAs using polyethylene glycol-mediated and electroporation protocols. Both methods are successful in some mycoplasmas, but other mycoplasmas can be transformed by only one of these methods, and there are a few mycoplasmas in which transformation has not been successful.

Two types of genetic transfer in mycoplasmas requiring cell-to-cell contact have been observed: intraspecies transfer in *Spiroplasma* and *Acholeplasma*, but not *Mycoplasma*; and interspecies transfer from the Gram-positive bacterium *Enterococcus faecalis* to *Mycoplasma*. The interspecies studies involved transfer of the conjugal transposon Tn916 carrying the *tetM* selectable marker. This type of transfer probably occurs in nature, to explain the prevalence of the *tetM* gene in *Mycoplasma* and *Ureaplasma* clinical isolates. Both conjugation-like and cell fusion models have been proposed for the mechanism of cell-to-cell gene transfer in mycoplasmas.

### 7.2. Genetic Recombination

Homologous, site-specific, and illegitimate recombination have been found in mycoplasmas (14). As in bacteria, homologous recombination in *Acholeplasma* and *Mycoplasma* probably occurs via RecA-dependent pathways. These mycoplasmas have *recA* genes, and the enzymes required for homologous recombination have been identified in *Mycoplasma* genome sequence analyses. However, homologous recombination in *Spiroplasma* may involve a non-RecA-dependent pathway because *Spiroplasma* may not have complete *recA* genes. Mycoplasma site-specific recombination occurs in integration of *Acholeplasma* and *Spiroplasma* viruses during lysogenization, movement of transposable elements in *Acholeplasma* and *Mycoplasma* genomes, and DNA inversions regulating expression of restriction and modification and phase-variable surface antigens in *My. pulmonis*. Illegitimate recombination has been suggested to explain the deletions and genetic rearrangements in plasmids maintained in *Acholeplasma*.

### 7.3. Antigen and Phase Variation

*Mycoplasma* and *Ureaplasma* surface antigens can undergo high frequency variation (14). This variation may be important in pathogenesis by enabling mycoplasmas to adhere to different tissue receptors or evade a host immune response. Variation in two mycoplasmas has been studied in detail. The *vlp* gene family in *My. hyorhinae* contains three to six genes that can be expressed either individually or in combination. Regulation is via changes in the number of adenine residues in the poly(A) spacer between the promoter -35 and -10 sites for each *vlp* gene, which somehow affects *vlp* gene expression. In contrast, the *vsa* gene family in *My. pulmonis* contains at least seven genes, only one of which is expressed in any cell, and regulation among *vsa* genes is via DNA inversions. Similar DNA inversions also regulate restriction and modification in *My. pulmonis*.

### 7.4. Mycoplasmas and the Minimal Cell Genome

The small size of mycoplasma genomes led to interest in them as models for a minimal cell genome, perhaps the primordial cell genome. The increasing availability of mycoplasma and eubacterial genome sequences is enabling comparative sequence analyses to identify a putative minimal gene set (19). However, small mycoplasma genomes resulted from gene attrition during the degenerate evolution of mycoplasmas from larger eubacterial genomes in cells growing under moderate environmental conditions in niches with a variety of organic nutrient and metabolic precursors (20). In contrast, the primordial cell genome probably resulted from gene accretion in cells growing in an

environment of high temperature, reducing atmosphere, and sulfur metabolism. Therefore, a proposed minimal gene set extracted from nascent microbial genes is not expected to be a useful model for a primordial genome in early cell evolution.

### Bibliography

1. J. Maniloff (1992) In *Mycoplasmas: Molecular Biology and Pathogenesis* (J. Maniloff, R. N. McElhaney, L. R. Finch and J. B. Baseman, eds.), American Society for Microbiology, Washington, DC, pp. 549–559.
2. S. Razin, D. Yogeve, and Y. Naot (1998) *Microbiol. Mol. Biol. Rev.* **62**, 1094–1156.
3. J. Renaudin, P. Aullo, J. C. Vignault, and J.M. Bové (1990) *Nucl. Acids Res.* **18**, 1293.
4. J. Renaudin, M. C. Pascarel, and J. M. Bové (1987) *J. Bacteriol.* **169**, 4950–4961.
5. J. Maniloff, G. J. Kampo, and C. C. Dascher (1994) *Gene* **141**, 1–8.
6. N. Zou, K. Park, and K. Dybvig (1995) *Plasmid* **33**, 41–49.
7. W. Just and G. Klotz (1990) *J. Gen. Virol.* **71**, 2157–2162.
8. A.-H. T. Tu, L. L. Voelker, X. Shen, and K. Dybvig (2001) *Plasmid* **45**, 122–126.
9. L. L. Voelker and K. Dybvig (1999) *Gene* **233**, 101–107.
10. C. M. Fraser, J. D. Gocayne, O. White, et al (1995) *Science* **270**, 397–403.
11. R. Himmelreich, H. Hilbert, H. Plagens, E. Pirkl, B.-C. Li, and R. Herrmann (1996) *Nucl. Acids Res.* **24**, 4420–4449.
12. I. Chambaud, R. Heilig, and S. Ferris, et al (2001) *Nucl. Acids Res.* **29**, 2145–2153.
13. J. I. Glass, E. J. Lefkowitz, J. S. Glass, C. R. Heiner, E.Y. Chen, and G.H. Cassell (2000) *Nature* **407**, 757–762.
14. K. Dybvig and L. L. Voelker (1996) *Annu. Rev. Microbiol.* **50**, 25–57.
15. D. C. Quinlan and J. Maniloff (1972) *J. Bacteriol.* **112**, 1375–1379.
16. K. Dybvig and H. Yu (1994) *Mol. Microbiol.* **12**, 547–560.
17. D. Yogeve, R. Rosengarten, and K. S. Wise (1991) *Zbl. Bakt.* **278**, 275–286.
18. S. Loechel, J. M. Inamine, and P.-C. Hu (1991) *Nucl. Acids Res.* **19**, 6905–6911 .
19. A. R. Mushegian and E. V. Koonin (1996) *Proc. Natl. Acad. Sci. U.S.A.* **93**, 10268–10273.
20. J. Maniloff (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10004–10006.
21. H. Neimark and B. C. Kirkpatrick (1993) *Mol. Microbiol.* **7**, 21–28.

### Suggestions for Further Reading

22. K. Dybvig and L. L. Voelker (1996) Molecular biology of mycoplasmas. *Annu. Rev. Microbiol.* **50**, 25–57.
23. J. Maniloff (1996) The minimal cell genome: “On being the right size.” *Proc. Natl. Acad. Sci. U.S.A.* **93**, 10004–10006.
24. J. Maniloff, R. N. McElhaney, L. R. Finch, and J. B. Baseman (1992) *Mycoplasmas: Molecular Biology and Pathogenesis*, American Society for Microbiology, Washington, DC.
25. A. R. Mushegian and E. V. Koonin (1996) A minimal gene set for cellular life derived by comparison of complete bacterial genomes. *Proc. Natl. Acad. Sci. U.S.A.* **93**, 10268–10273.
26. S. Razin, D. Yogeve, and Y. Naot (1998) Molecular biology and pathogenicity of mycoplasmas. *Microbiol. Mol. Biol. Rev.* **62**, 1094–1156.

### Myeloma Proteins

Prior to the development and universal application of [hybridoma](#) methodology to generate [monoclonal antibodies](#) of defined specificity, immunologists relied on the existence of homogeneous pathological immunoglobulins called *myelomas*. Myeloma proteins are produced spontaneously in humans due to [transformation](#) of a lymphocyte to the malignant state. Myeloma proteins can also be generated in specific strains of mice through administration of various reagents, such as mineral oil, administered in the peritoneal cavity. However, for a long time, myeloma globulins produced during human and animal disease were considered to be devoid of functional properties and were classified along with normal immunoglobulins as nonspecific and lacking antigen binding capacity. With advent of the Clonal Selection Theory (1), homogeneous myeloma proteins become antibodies in search of antigen. Myeloma proteins in humans appear with a frequency that parallels their concentration in serum and, therefore, the respective number of lymphocytes synthesizing that class of antibody. Thus, human IgG myelomas appear most frequently (~90%) followed by IgM (10%). Myelomas in the other classes, such as IgA, IgD, and IgE, are quite rare in humans. In contrast, myelomas induced in mice are predominantly IgA, followed by IgG and IgM in incidence. The predominance of IgA reflects the administration of mineral oil into the peritoneal cavity, which is rich in IgA-producing lymphocytes.

Myeloma proteins proved homogeneous at the protein level and, because they were produced in high quantities, were easily purified to homogeneity. These homogeneous proteins provided the foundation upon which the first primary structures of immunoglobulins were based. Light chains of myeloma proteins in human patients exist in excess and are harvested from urine as monomers or stable dimers called [Bence-Jones proteins](#). Because light chains are of lower molecular weight (~22 to 23 kDa) and lack carbohydrate moieties, they were used to obtain the first primary structures. Putnam et al. (2) published the first complete sequence analysis of a Bence-Jones protein called Ag, a sequence that became the standard by which all subsequent light chains primary structure analyses would be compared. Roy was the second completed sequence and was significant because it was the comparison of primary structures that yielded critical information. The comparison of several primary structures of heavy and light chains normalized for isotype revealed the existence of the variable and constant domains (3). Further comparisons within the variable domains showed that the H chain consisted of four hypervariable regions termed **complementarity determining regions** (CDRs). The L-chain variable domains consisted of 3 CDR segments. This finding, combined with similar results from analyses of hybridoma proteins, led to delineation of the molecular basis of immunological specificity.

Monoclonal antibodies proved to be the “tool of nature” whereby immunologists first learned of the structure–function relationship in antibody molecules. However, these basic relationships extended beyond the antibody molecule when it was determined that other members of the immunoglobulin superfamily of molecules (eg, [T-cell receptors](#), **major histocompatibility** molecules, **T-cell** accessory proteins, and adhesion molecules) utilized the same basic motifs and canonical structures first determined in monoclonal antibodies (4).

#### Bibliography

1. F. M. Burnet (1957) *Aust. J. Sci* **20**, 67–69.
2. F. W. Putnam, K. Titani, and E. Whitley Jr. (1966) *Proc. R. Soc. Lond.* **166**, 124–137.
3. E. A. Kabat and T. T. Wu (1970) *J. Exp. Med.* **132**, 211–250.
4. J. Kuby (1996) *Immunology*, 3rd ed., W. H. Freeman, New York, pp. 129–131.

# Myoglobin

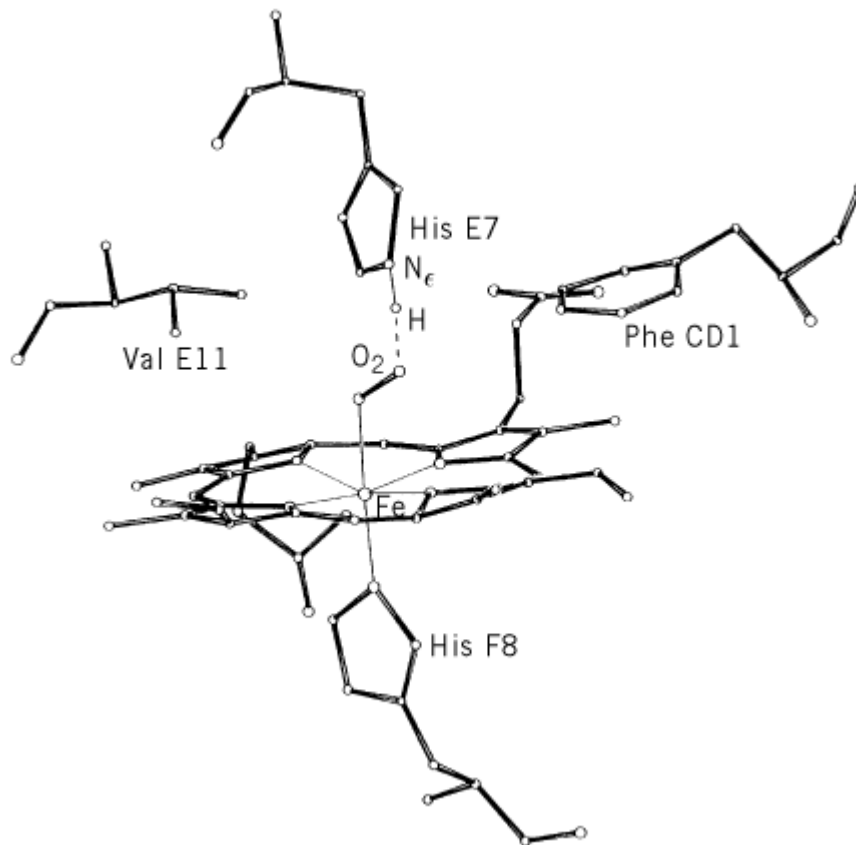
## 1. Distribution and Classification

Myoglobin (Mb) is one of hemoproteins, composed of globin and protoheme IX. It is found in cytoplasm of many vertebrate cardiac, skeletal and smooth muscles. Its fundamental function is intracellular oxygen store and transport. Vertebrate myoglobins are monomeric with a molecular mass of approximately 17 kDa. Myoglobin also occurs in invertebrates. It is found in body wall muscle of platyhelminthes, nemertea, nematoda, echiura and annelida, in radular muscle of amphineura, in buccal mass, radular muscle, pharyngeal muscle or triturative stomach of gastropods and in adductor muscle of bivalve. These invertebrate myoglobins are either monomeric or dimeric with a molecular mass of approximately 17 kDa per heme. Some species of gastropoda such as *Sulculus diversicolor* and *Nordotis madaka* (abalones) have a novel myoglobin in their radular muscle. It possesses one heme group but the polypeptide chain that has a molecular mass of approximately 40 kDa shows significant sequence homology with indoleamine dioxygenase (IDO), a tryptophan degrading enzyme carrying protoheme (1). The IDO-like myoglobin is considered to have evolved from a duplicated IDO gene through a pathway completely different from that for conventional globin (1).

## 2. Structure and Function

Skeletal muscle myoglobin of sperm whale is the first protein whose three-dimensional structure was solved by X-ray crystallography (2). Its structure is the origin for defining the “globin fold.” Several amino acid residues surrounding the heme group are important for stabilizing of the heme-globin linkage and for reversible oxygen binding (Fig. 1). His-F8 (proximal His) and Phe-CD1 are strictly conserved. In most Mbs the site E7 is occupied by His but it is occasionally replaced by other residues such as Gln, Val, Leu, or Tyr. His-E7 stabilizes the oxygen-iron bonding by forming a hydrogen bond with one of the oxygen atoms (Fig. 1). The amino acid residues at E7, E11, B10 and CD4 are known to make major contributions to control of ligand binding to Mb, as manifested in recombinant Mbs (3).

**Figure 1.** The heme group and key amino acid residues of sperm whale oxymyoglobin (14). His-F8 is the proximal His. His-E7 is the distal His, forming a hydrogen bond with the bound oxygen. Other key residues that control oxygen binding are Val-E11, Phe-CD1, Val-B10 and Leu-CD4 (the latter two are not shown).



Generally, the oxygen affinity of Mb is higher than that of hemoglobin within a single species. This enables oxygen transfer from the latter to the former. The oxygen dissociation curve of monomeric Mb is of hyperbolic shape whereas dimeric Mbs of annelida, amphineura and gastropoda often show sigmoid curves with Hill coefficient (see [Hemoglobin](#)) values greater than unity. Mb has no cofactor which modifies its oxygen binding properties. This feature is in striking contrast to that of Hb which shows the oxygen transport function allosterically regulated by several cofactors (see [Hemoglobin](#)). However, it is reported that lactate, a product of glycolysis, causes a marked lowering of oxygen affinity for sperm whale Mb and horse heart Mb (4).

The physiological functions of Mb are considered to be short-term or long-term store of oxygen (oxygen buffer) and facilitation of intracellular diffusion of oxygen from red cells to mitochondria. Mb serves as an oxygen store during temporary deficits in oxygen supply. In beating heart or exercising skeletal muscle Mb acts as a short-term oxygen store, tiding the muscle over from one contraction to the next (5). Long-term oxygen supply during diving supposed for aquatic mammals which have high concentrations of skeletal Mb is of small significance because the amount of oxygen bound to Mb is too little to ensure the long-term supply. There is a line of evidence supporting the validity of Mb-facilitated diffusion of oxygen in muscle cells. In this mechanism the amount of oxygen transported by translational diffusion of Mb makes a significant contribution to the total oxygen flux in the cell compared to spontaneous diffusion of free dissolved oxygen. However, for this mechanism to be proven to operate in vivo several discrepancies must be resolved (6).

Recently, new functional roles of Mb, most of which are not based on reversible oxygen binding but on enzymatic activities, were proposed. They include a defense mechanism against oxidative stress through a peroxidase activity of Mb (7), augmentation of mitochondrial oxidative phosphorylation (8), augmentation of signal transmittance by NO through a redox reaction with NO (9) and inhibition of the inactivation of cytochrome c oxidase by NO (i.e., protection of cellular respiration) through

NO scavenging by Mb (10). The recent observations that myoglobinless gene knockout mice are fertile and exhibit normal exercise capacity and a normal ventilatory response to low oxygen levels seemed to imply that Mb was not required to maintain the normal cardiovascular and musculoskeletal function in a terrestrial mammal (11). However, the later observations that the elimination of Mb was accompanied by the activation of multiple compensatory mechanisms in the circulatory system indicate that Mb is still important for the delivery of oxygen (12).

IDO-like Mbs have no sequence homology with conventional Mbs. Nevertheless, they show light absorption spectra similar to those of conventional Mbs and are capable of reversible oxygen binding, whereas they have completely lost IDO activity (13). Their oxygen affinity is lower than that of conventional Mbs (13). Their three-dimensional structure is not elucidated, and the amino acid residues corresponding to the proximal and distal residues of conventional Mb are not identified.

### Bibliography

1. T. Suzuki and T. Takagi, *J. Mol. Biol.* **228**, 698–700 (1992).
2. J.C. Kendrew et al., *Nature* **185**, 422–427 (1960).
3. B.A. Springer, S.G. Sliger, J.S. Olson and G.N. Phillips, Jr., *Chem. Rev.* **94**, 699–714 (1994).
4. B. Giardina et al., *J. Biol. Chem.* **271**, 16999–17001 (1996).
5. G.A. Millikan, *Physiol. Rev.* **19**, 503–523 (1939).
6. K.D. Jurgens, T. Peters, and G. Gros, *Proc. Natl. Acad. Sci. U.S.A.* **91**, 3829–3833 (1994).
7. D. Galaris, E. Cadenas, and P. Hochstein, *Arch. Biochem. Biophys.* **273**, 497–504 (1989).
8. J.E. Doeller and B.A. Wittenberg, *Am. J. Physiol.* **261**, H53–H62 (1991).
9. J.R. Lancaster, *Proc. Natl. Acad. Sci. U.S.A.* **91**, 8137–8141 (1994).
10. M. Brunori, *Trends Biochem. Sci.* **26**, 21–23 (2001).
11. D.J. Garry et al., *Nature* **395**, 905–908 (1998).
12. A. Gödecke et al., *Proc. Natl. Acad. Sci. U.S.A.* **96**, 10495–10500 (1999).
13. T. Suzuki, H. Kawamichi, and K. Imai, *Comp. Biochem. Physiol.* **B121**, 117–128 (1998).
14. S.E.V. Phillips, *Nature* **273**, 247–248 (1978).

### Additional Reading

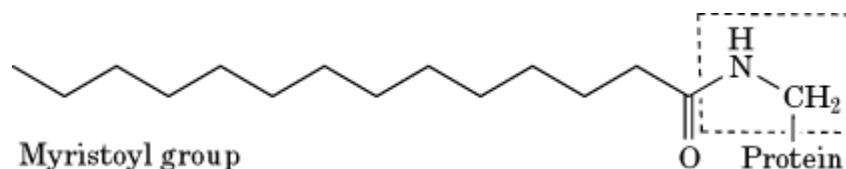
15. Suzuki T. and Imai K., *Evolution of myoglobin*, *Cell. Mol. Life Sci.* **54**, 979–1004 (1998).
16. Wittenberg J.B., *Myoglobin-facilitated oxygen diffusion: role of myoglobin in oxygen entry into muscle*, *Physiol. Rev.* **50**, 559–636 (1970).
17. Wittenberg B.A. and Wittenberg J.B., *Transport of oxygen in muscle*, *Annu. Rev. Physiol.* **51**, 857–878 (1989).

### Myristoylation

*N*-Myristoylation is a cotranslational process in which eukaryotic or viral [proteins](#) become amide-linked to the 14-carbon saturated fatty acid, myristic acid (see Fig. 1 and Table 1). Myristoylation occurs on the  $\alpha$ -amino group of the N-terminal [glycine](#) residue exposed after removal of the initiator [methionine](#) residue (see [Translation](#)). The reaction takes place in the cytosol, and is catalyzed by *N*-myristoyl transferase using myristoyl CoA as the myristic acid donor. Recent work suggests that *N*-

myristoylation may also occur post-translationally but as yet there is only one instance of this type of reaction (1). Myristoylation provides some proteins with a relatively weak [membrane anchor](#), and is important for correctly localizing them within the cell. Myristoylation also occurs on many proteins that do not appear to associate with [membranes](#). The role of the myristoyl group in these proteins is less clear, but may involve stabilizing the protein structure, facilitating [protein-protein interactions](#) by binding to a site on another protein, and viral assembly.

**Figure 1.** Modification of an N-terminal glycine residue by a myristoyl group. The N-terminal glycine residue of the protein (outlined by the dotted line) is amide-linked via its exposed  $\alpha$ -amino group to the 14-carbon saturated fatty acid, myristic acid.



**Table 1. Examples of Myristoylated Proteins**

---

|   |
|---|
| Src family tyrosine kinases (eg, p60 <sup>src</sup> , p59 <sup>fyn</sup> , p56 <sup>lck</sup> ) |
| G-protein $\alpha$ subunits (eg, $\alpha_i$ , $\alpha_o$ )                                      |
| MARCKS protein  |
| ADP ribosylation factor   |
| Endothelial nitric oxide synthase   |
| Recoverin   |
| NADH-Cytochrome $b_5$ reductase   |
| cAMP dependent protein kinase (catalytic subunit) <sup>a</sup>                                  |
| Calcineurin <sup>a</sup>  |
| Retroviral <i>gag</i> polyproteins <sup>a</sup>   |
| Viral coat proteins (eg, poliovirus VP4) <sup>a</sup>   |

---

<sup>a</sup> Not associated with membranes.

In studies with the purified N-myristoyl transferase and synthetic peptide substrates, the only amino acid that is absolutely required for myristoylation is the N-terminal glycine (2, 3). This observation agrees well with the known distribution of myristoylation among naturally occurring proteins. On the other hand, the identity of residues at positions 2–8 can also have some influence on the efficiency of myristoylation. For example, substitution of residue Ser5 (in **G-protein**  $\alpha_s$  subunits) with an Asp residue may be the reason that  $\alpha_s$  subunits are not myristoylated, in contrast to most other  $\alpha$  subunits. A search of the yeast genome database suggests that about 1% of yeast ORFs will encode myristoylated proteins (4). Palmitoyl CoA is not a substrate for N-myristoyl transferase, which accounts for the relative lack of naturally occurring N-palmitoylated proteins. The transferase is less



selective for shorter acyl chain lengths (eg, C<sub>12</sub>) or for degree of unsaturation (3). Detailed analysis of several myristoylated proteins reveals considerable acyl heterogeneity, and in some cases myristic acid is not even the most abundant molecular species.

That myristoylation occurs on proteins that do not bind to membranes suggests that the myristoyl group, unlike other lipid groups attached to proteins, may have functions other than membrane anchoring, but there is relatively little detailed information as to what these other functions might be. The amide linkage in N-myristoylated proteins is relatively resistant to hydrolysis, and is not degraded or turned over metabolically. Myristoylated and nonmyristoylated forms of the **cyclic AMP**-dependent protein kinase regulatory subunit have similar kinetic properties, and are equal in their ability to associate with the catalytic subunit. The myristate, however, does seem to increase thermal stability of the kinase, possibly by binding to an adjacent hydrophobic region of the protein (5). Myristoyl groups also seem to play a role in stabilizing interactions between protein molecules. This may be especially important for **viral** assembly, eg, the interaction between VP4 and VP3 on the inner surface of the [polio virus](#) capsid. The short length and lack of bulky side chains may make myristate particularly versatile and capable of binding to a membrane or to a protein site, a feature that is exploited for regulatory purposes as described below.

The myristoyl group has a relatively low affinity for membranes and by itself results only in weak and transient binding of proteins to membranes (6-8). It has been shown that membrane affinity of myristoylated proteins can be altered by three distinct mechanisms and may provide a way to regulate their distribution within the cell (for a general discussion of factors that affect membrane affinity of lipid-anchored proteins see [Membrane Anchors](#)):

1. Basic amino acid residues close to the N-terminus will bind to the negatively charged phospholipid bilayer and increase the membrane affinity, as proposed for the tyrosine kinase p60<sup>src</sup> and the MARCKS protein. The strength of these [electrostatic interactions](#) is reduced by **phosphorylation** of the MARCKS protein, causing it to redistribute within the cell (9).
2. [Palmitoylation](#) of a [cysteine](#) residue close to the amino terminus increases the affinity of myristoylated proteins for intracellular membranes (10, 11). This occurs in some G protein subunits (eg, a<sub>o</sub>) and Src family tyrosine kinases (eg, p59<sup>lyn</sup> and p56<sup>lck</sup>). Palmitoylation of these proteins also depends on prior [myristoylation](#), because the palmitoyltransferase activities are localized in membranes (see [Palmitoylation](#)).
3. The myristoyl group may also bind to a site on the protein itself, which prevents it from acting as a membrane anchor. The myristoyl group might be extruded by a conformational change that “closes” the binding site, causing the protein to bind to the membrane. This mechanism has been proposed for recoverin and hippocalcin (Ca<sup>2+</sup>-sensing proteins in the visual and nervous systems) and for **ADP ribosylation** factor, where the conformational changes occurs in response to Ca<sup>2+</sup> and GTP, respectively (12).

## Bibliography

1. J. Zha, S. Weiler, K. J. Oh, M. C. Wei, and S. J. Korsmeyer (2000) *J. Biol. Chem.*, in press.
2. D. A. Towler, J. I. Gordon, S. P. Adams, and L. Glaser (1988) *Annu. Rev. Biochem.* **57**, 69–99.
3. D. R. Johnson, R. S. Bhatnagar, L. J. Knoll, and J. I. Gordon (1994) *Annu. Rev. Biochem.* **63**, 869–914.
4. T. A. Farazi, G. Waksman, and J. I. Gordon (2001) *J. Biol. Chem.*, in press.
5. W. Yonemoto, M. L. McGlone, and S. S. Taylor (1993) *J. Biol. Chem.* **268**, 2348–2352.
6. S. Shahinian and J. R. Silvius (1995) *Biochemistry* **34**, 3813–3822.
7. C. A. Buser, C. T. Sigal, M. D. Resh, and S. McLaughlin (1994) *Biochemistry* **33**, 13093–13101.

8. R. M. Peitzsch and S. McLaughlin (1993) *Biochemistry* **32**, 10436–10443.
9. J. T. Seykora, M. M. Myat, L.-A. Allen, J. V. Ravetch, and A. Aderem (1996) *J. Biol. Chem.* **271**, 18797–18802.
10. M. Y. Degtyarev, A. M. Spiegel, and T. L. Z. Jones (1994) *J. Biol. Chem.* **269**, 30898–30903.
11. L. Alland, S. M. Peseckis, R. E. Atherton, L. Berthiaume, and M. D. Resh (1994) *J. Biol. Chem.* **269**, 16701–16705.
12. J. B. Ames, T. Tanaka, M. Ikura, and L. Stryer (1995) *J. Biol. Chem.* **270**, 30909–30913.

### Suggestions for Further Reading

13. D. E. Hruby and C. A. Franke (1993) Viral acylproteins: Greasing the wheels of assembly. *Trends Microbiol.* **1**, 21–24.
14. R. S. Bhatnagar and J. I. Gordon (1997) Understanding covalent modification of proteins by lipid: Where cell biology and biophysics mingle. *Trends Cell Biol.* **7**, 14–20.
15. S. McLaughlin and A. Aderem (1995) The myristoyl-electrostatic switch: A modulator of reversible protein-membrane interactions. *Trends Biochem. Sci.* **20**, 272–276.
16. M. D. Resh (1994) Myristylation and palmitylation of Src family members: The fats of the matter. *Cell* **76**, 411–413.
17. P. J. Casey and J. E. Buss (1995) Lipid modifications of proteins, *Meth. Enzymol.* **250**.

### N-End Rule

This term refers to a pathway for **protein degradation** that depends on the identity of the N-terminal amino acid residue of a protein. The N-end rule was originally elucidated in yeast by Varshavsky and coworkers, who studied the half-life of **beta-galactosidase** after genetically modifying its amino terminus. They found that if basic and large **hydrophobic** amino acids were placed, instead of the normal [methionine](#) residue at the N-terminus of b-galactosidase, the protein became very unstable. This degradative pathway requires the ubiquitin carrier protein E2<sub>14</sub> kDa (yeast UBC2/RAD6) and the Ub-protein ligase, E3a. The protein E3a (designated Ubr1 in yeast) recognizes such N-terminal protein substrates and forms a large complex with protein E2<sub>14</sub> kDa, an amidase (yeast NTA1), and an arginyl-tRNA-transferase (yeast ATE1). These additional enzymes appear to be important for modifying the N-terminus of substrates so they can be recognized by E3a. The rare cellular proteins with these abnormal amino termini can be degraded by this “N-end” pathway. The vast majority of cytosolic proteins, however, contain N-a-acetylated or methionyl amino termini, which should not be recognized by this system. Therefore, the functional importance of this pathway is uncertain. A number of recent observations suggest that this E2/E3 pair is responsible for a major portion of the protein breakdown in mammalian muscle.

### Suggestions for Further Reading

1. A. Bachmair, D. Finley, and A. Varshavsky (1986) In vivo half-life of a protein is a function of its amino-terminal residue. *Science* **234**, 179–186.
2. A. Varshavsky (1996) The N-end rule: functions, mysteries, uses. *Proc. Natl. Acad. Sci. USA* **93**, 12142–12149.

## ***n*-Octyl-β-D-Glucoside**

This [detergent](#) has a high [critical micelle concentration](#) (cmc) of 23 nM and has been used to solubilize and reconstitute many integral [membrane proteins](#) efficiently. The special advantage in using this detergent is in the efficient postsolubilization removal of the detergent by [dialysis](#). It is only mildly **denaturing**, which makes it highly efficient in solubilizing integral [membrane proteins](#), such as the nicotinic [acetylcholine receptor](#) and those for opioids, **GABA<sub>A</sub>** (1-3), and the photoreceptor [guanylate cyclase](#) (4). However, the relative inefficiency of this detergent in solubilizing plasma membrane [lipids](#) (as compared to CHAPS and dodecylmaltoside, see [Detergents](#)), makes it inappropriate for solubilizing and reconstituting the heavily lipid-dependent serotonin 1A receptor (5, 6). Nevertheless, because of its weakly denaturing property, this detergent has been efficiently used to obtain active preparations of other solubilized membrane proteins, such as the **mitochondrial** membrane protein carnitine palmitoyltransferase I (7) or the **GPI-anchored** proteins (8). It is also a popular detergent in [X-ray crystallography](#), where it is used to crystallize solubilized membrane proteins. Typical examples of such studies include the solubilization of the bifunctional enzyme fructose 6-phosphate, 2-kinase:fructose 2,6-bisphosphatase (9), and prostaglandin H synthase-1 (10).

### Bibliography

1. J. M. Gonzalez-Ros, A. Paraschos, M. C. Farach, and M. Martinez-Carrion (1981) *Biochim. Biophys. Acta.* **643**, 407–420.
2. T. Fujioka, F. Inoue, S. Sumita, and M. Kuriyama (1988) *Biochem. Biophys. Res. Commun.* **156**, 54–60.
3. S. M. J. Dunn, R. A. Shelman, and M. W. Agey (1989) *Biochemistry* **28**, 2551–2557.
4. K. W. Koch (1991) *J. Biol. Chem.* **266**, 8634–8637.
5. P. Banerjee, J. T. Buse, and G. Dawson (1990) *Biochim. Biophys. Acta.* **1044**, 305–314.
6. P. Banerjee, J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
7. M. P. Kolodziej and V. A. Zammit (1993) *FEBS Lett.* **327**, 294–296.
8. T. R. Hawn and M. Strand (1993) *Mol. Biochem. Parasitol.* **59**, 73–81.
9. E. S. Istvan, C. S. Hasemann, R. G. Kurumbail, K. Uyeda, and J. Deisenhofer (1995) *Protein Sci.* **4**, 2439–2441.
10. M. R. Garavito, D. Picot, and P. J. Loll (1996) *J. Bioenergetics Biomembranes* **28**, 13–27.

### **N-Region Diversity**

The existence of the N-region of **immunoglobulins** (Igs) was postulated when it was observed that nucleotide sequences of small, but variable, lengths identified at the VD and DJ junctions of Ig heavy chains had no counterpart in any of the **V gene**, **D gene**, or **J gene** segments in the germline. It turned out that these nucleotides were added at both the VD and the DJ junctions by the [enzyme](#) terminal deoxynucleotidyl transferase (TdT), which is specifically expressed in lymphocytes.

Direct evidence that N diversity was indeed the result of TdT activity was first obtained by

[transfection](#) of a non-Ig-producing fibroblast cell line with a DNA vector containing an insert of unrearranged D and J segments. [Gene rearrangement](#) was observed, provided that a co-transfection with the RAG1/RAG2 [recombinase](#) activating genes was made. In the absence of TdT, no additional nucleotides were inserted. When co-transfection also included the TdT gene, N diversity was observed in the rearranged substrate. The ultimate proof came from gene targeting experiments, because mice rendered TdT<sup>-/-</sup> had heavy chains that no longer contained N diversity.

The results of these experiments are also in good agreement with analysis of [B-cell](#) differentiation that takes place in fetal liver and in bone marrow. In mouse fetal liver, TdT is not expressed and, as a result, no N diversity is seen in B cells that have differentiated in this organ. By contrast, TdT expression takes place in bone marrow, but only at precise steps of B-cell differentiation—that is, at the proB stage. This is the time when DJ rearrangements first take place, which correlates nicely with the major N diversity that accumulates at the D–J joint. It is still active whenever the V to DJ recombination takes place, but is then down-regulated. As a result, no or very limited N diversity is seen in the light chains. Recent data confirmed the presence of N diversity in some light chains and suggest that this occurs later in the periphery, whenever both the TdT gene and the RAG system are reactivated transiently, especially in the course of an antigenic stimulation. This is another example of the extraordinary plasticity of the mechanisms that ensure Ig diversity.

See also entries [Gene Rearrangement](#), [Recombinase](#), and [Junctional Diversity](#).

#### Suggestions for Further Reading

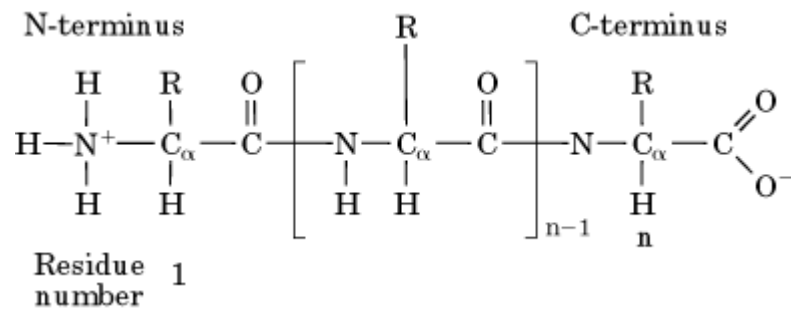
T. Komori, A. Okada, V. Stewart, and F. W. Alt (1993) Lack of N regions in antigen receptor variable region genes of TdT-deficient lymphocytes. *Science* **261**, 1171–1175.

L. A. Bentolila, G. E. Wu, F. Nourrit, M. Fanton d'Anton F. Rougeon, and N. Doyen (1997) *J. Immunol.* **158**, 715–723.

## ***N*-Terminus**

The *N*-terminus is the term for one end of the [polypeptide chain](#) of a [protein](#). Proteins are **polymers** formed by condensation of the [amino groups](#) and [carboxyl groups](#) of [amino acids](#). That end of the polypeptide chain having an uncondensed or “free” amino group is called the amino terminus or *N*-terminus (Fig. 1). By convention, the *N*-terminus is the first residue in the protein sequence. Similarly, the other end of the polypeptide chain, having an uncondensed or free carboxyl group, is called the carboxyl terminus, or ***C*-terminus**, and by convention is the last residue in the protein sequence. The *N*-terminal amino group is usually positively charged at physiological pH.

**Figure 1.** Schematic representation of the covalent structure of a protein, showing the *N*-terminus as the first residue in the polypeptide chain. The square brackets indicate the repeating part of the chain (the backbone), and the R group denotes the variable side chain of each amino acid residue. The *C*-terminus is the last residue in the chain.



The *N*-terminal protein residue can be specifically and selectively modified *in vivo* (see [Post-Translational Modifications](#)). Proteins synthesized by the cell, as opposed to chemically synthesized by [peptide synthesis](#), have an initiating **methionine residue** at the *N*-terminus. This *N*-terminal methionine residue is often removed by cellular enzymes. Also, **acetylation** of the amino group of the *N*-terminus can be catalyzed by enzymes.

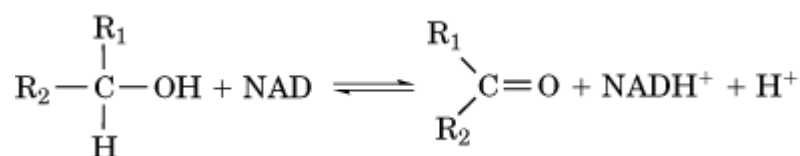
The sequence of a protein can be determined experimentally using a process called the [Edman Degradation](#). This process is based on successive removal and identification of the *N*-terminal residue of the protein.

[See also **C-Terminus** and [Polypeptide Chain](#).]

## NAD-Binding Proteins

NAD-binding proteins bind the dinucleotide NAD (nicotinamide adenine dinucleotide) and generally are [enzymes](#) that **catalyze a redox reaction** in which a proton is transferred from the carbon of a substrate alcohol group to NAD, thereby oxidizing the substrate and reducing NAD to NADH (Fig. 1). Examples of NAD-binding proteins include [lactate dehydrogenase](#), [alcohol dehydrogenase](#), and malate dehydrogenase. These different proteins catalyze the same basic reaction, but on different substrates, simply as a result of having differing substrate specificities.

**Figure 1.** The reaction catalyzed by the NAD-binding proteins.  $R_1$  and  $R_2$ , vary depending on the substrate specificity of the enzyme. For malate dehydrogenase, the substrate malate has  $R_1 = \text{COOH}$  and  $R_2 = \text{CH}_2\text{-COOH}$ . For lactate dehydrogenase, the substrate lactate has  $R_1 = \text{COOH}$  and  $R_2 = \text{CH}_3$ . For alcohol dehydrogenase, the substrate ethanol has  $R_1 = \text{CH}_3$  and  $R_2 = \text{H}$ .



The [protein structures](#) of these enzyme are composed of two **domains**; one domain binds the

substrate, and the other binds NAD. The substrate-binding domains of the different proteins are unrelated, but their NAD-binding domains have a common structural fold. The NAD-binding domain is a symmetrical 6-stranded **b-sheet** with **a-helices** on both sides, formed from two b-a-b-a-b [protein motifs](#), called **mononucleotide-binding motifs** or [Rossmann folds](#). Each of the two mononucleotide-binding motifs binds one-half of the dinucleotide; the adenine moiety of NAD binds to the first of these structural motifs, and nicotinamide binds to the second. The first three elements of secondary structure in the common adenine-binding motif, b1-aA-b2, have a similar sequence amongst the NAD-binding proteins (1). This fingerprint sequence has small **hydrophobic** residues conserved in b1, followed by a glycine-rich region (-Gly-X-Gly-X-X-Gly, where X is any residue), additional small hydrophobic residues in the aA-helix, and an Asp/Glu at the end of helix aA. The glycine residues adopt conformations forbidden to other residues and allow close packing of the strands and helix, plus a close approach between the adenine pyrophosphate and the N-terminal end of the aA-helix. The Asp/Glu residue interacts with the ribose hydroxyl. This signature sequence of ~30 residues was derived from the tertiary structures of known NAD-binding proteins and can be used to predict NAD-binding in proteins from simply their primary structure.

[See also [Proteins](#) and [Nucleotide-Binding Motif](#).]

#### Bibliography

1. R. K. Wierenga, P. Terpstra, and W. G. Hol (1986) *J. Mol. Biol.* **187**, 101–107.

#### Suggestion for Further Reading

2. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

### *Nanos* Gene

The *nanos* is a **maternal effect** gene required for the development of the fruit fly *Drosophila melanogaster*. *nanos* is required during oogenesis for the development and maintenance of the **germ line**, in the embryo for patterning of the abdomen, and in the primordial germ cells for their migration to the embryonic gonad. *nanos messenger RNA* is a model system for the study of RNA localization and [translation](#) mediated by the 3'-untranslated region (UTR), and the Nanos protein is studied for its ability to repress translation of mRNAs through sequences within their 3'-untranslated regions (see **Translational repressors**).

#### 1. Protein and RNA Structure

The *nanos* gene encodes a 3-kb primary transcript, containing two **introns**, that is processed to a mature 2.4-kb mRNA (1). The 3'-untranslated region contains sequences required for localization and translational repression of the transcript (2). *nanos* mRNA is detected at high levels in ovaries, 0- to 8-hour embryos, and adult female flies (1). Low levels of the transcript are also detected in adult male flies, but there is no genetic evidence of a role for *nanos* in males. The Nanos protein is a 400-amino-acid-residue [polypeptide chain](#) with a predicted molecular weight of 43 kDa. It contains several polyglutamine and polyasparagine repeats, which are common to *Drosophila* proteins. Analysis of *nanos* mutants reveals that a small carboxyl-terminal **domain** is required for *nanos* function *in vivo*. Within this region are two Cys-Cys-His-Cys motifs that have been shown to bind metal *in vitro* and closely resemble the retroviral nucleocapsid class of **zinc-finger** proteins (3).

## 2. Oogenesis

A *Drosophila melanogaster* ovary consists of approximately 16 ovarioles, each representing an independent egg assembly line (for review see Ref. 4). Every ovariole is a linear array of developing egg chambers beginning at the anterior with the germarium, where the progeny of the germ-line and somatic [stem cells](#) are organized into egg chambers, and progressing to more mature egg chambers as they move toward the posterior. In the germarium, a germ-line stem cell produces a cystoblast by asymmetric division. The cystoblast undergoes four rounds of division with incomplete cytokinesis, resulting in 16 cells connected by cytoplasmic bridges; one of these cells becomes the oocyte while the other 15 become nurse cells. An egg chamber that has left the germarium consists of the oocyte and nurse cells surrounded by somatically derived follicle cells. As the egg chambers mature, the nurse cells synthesize and deposit factors into the oocyte that are required for the development of the oocyte and early embryo. Late in oogenesis, the nurse cells contract and dump the contents of their cytoplasm into the oocyte. Nanos protein is expressed at high levels in the germ-line cysts of the germarium and at lower levels in the germline stem cells and dividing cystoblasts.

Females **trans**-heterozygous for strong *nanos* **alleles** lay a few eggs, but rapidly become sterile due to a loss of germ-line stem cells from the ovary. Analysis of *nanos* mutant ovaries demonstrates a number of defects (5). In most instances, ovarioles completely lack germ line, suggesting that zygotic *nanos*, like *pumilio*, is required for incorporation of germ-line cells into the developing ovarioles. Less frequently, germ-line cells appear to be incorporated into the ovariole. In some of these ovarioles, germ-line stem cells are absent and germ-line cells may develop directly into egg chambers, while in other cases stem cells seem to be established but the subsequent development of the cystoblast and maintenance of the stem cells are impaired. This suggests that *nanos* is required for the proliferation and viability of the germ-line stem cells and the developing germ-line cysts (5).

Once developing egg chambers leave the germarium, Nanos protein is no longer detectable until late in oogenesis, when Nanos protein is expressed at high levels in the nurse cells, but not the oocyte (6). The function of *nanos* at this stage is not clear, because egg chambers that lack this expression develop normally. *nanos* does not appear to play a role in spermatogenesis.

## 3. Embryonic Patterning

In *Drosophila melanogaster*, the anterior–posterior axis is defined by three independent systems: the terminal group of genes that controls the formation of the unsegmented termini at the anterior and posterior ends of the embryo; the anterior group of genes that controls development of the head and thorax; and the posterior group of genes that determines the abdomen (for review see Ref. 7). Each system requires the localized activity of a maternally supplied factor that functions as a spatial signal resulting in a **morphogenetic** gradient. The terminal and anterior systems require the activities of *torso* and [bicoid](#), respectively. *nanos* is the posterior determinant; in the absence of *nanos*, embryos fail to form abdominal segments. *nanos* mRNA is localized to the posterior pole, where its translation results in a gradient of Nanos protein from the posterior pole to approximately 50% of the egg length (8). Mislocalization of *nanos* RNA to an ectopic site leads to ectopic abdomen formation (8). The only required function for *nanos* in abdomen formation is the repression of maternal *hunchback* translation at the posterior—resulting in a gradient of Hunchback protein, with its high point at the anterior pole (9–11). Hunchback is a zinc-finger containing [transcription factor](#) that represses translation of abdominal gap genes; consequently, repressing *hunchback* in the posterior allows the ordered expression of the abdominal gap genes and abdomen formation.

### 3.1. Regulation of *nanos* mRNA Localization and Translation

Throughout oogenesis, *nanos* RNA is synthesized in the nurse cells and deposited in the developing oocyte (6). Late in oogenesis, *nanos* RNA becomes enriched at the posterior pole of the embryo (1). This process requires a localization control element located within the first 400 nucleotides of *nanos* 3'-untranslated region and the ordered activity of the posterior group genes including *oskar*, a gene necessary for organization of the pole plasm (12, 13). *oskar* mutant embryos lack pole plasm and fail

to localize *nanos* mRNA to the posterior pole (14). Ectopic expression of *oskar* at the anterior pole results in ectopic pole plasm formation and in *nanos* mRNA being mislocalized to the anterior (14). *nanos* mRNA that is localized to the posterior pole plasm is translated, whereas unlocalized *nanos* mRNA is translationally repressed, suggesting that *nanos* mRNA localization and translation are linked (2). A 90-nucleotide translational control element in the 3'UTR of *nanos* is both necessary and sufficient for translational repression of *nanos*. The translational control element can function independently of the localization element, and RNA localization can occur in the absence of the translational control element (12, 13, 15). Translational repression can be overcome by removing the 3'-UTR, localizing the transcript to the anterior pole via heterologous RNA localization sequences from the *bicoid* 3'UTR, or by specifically mutating the translational control element. A 135-kDa protein, Smaug, has been shown to bind to the translational control element *in vitro* (15). Mutational analysis of the translational-control-element indicates that failure to bind Smaug *in vitro* correlates with lack of translational repression *in vivo* (15). It is not yet understood how the translational control element and Smaug mediate repression, or how this repression is relieved at the posterior pole.

### 3.2. Function of *nanos* in Embryonic Patterning

Although the maternal *hunchback* transcript is distributed throughout the early embryo, it is only translated in the anterior half of the embryo (16). Repression of maternal *hunchback* translation in the posterior requires sequences within the *hunchback* 3'-untranslated region (termed *nanos response elements* or NREs) and the activity of *nanos* and *pumilio* (17). Mutation of the NRE, or loss of *nanos* or *pumilio* activity, results in a loss of abdominal segmentation (17). Recently it has been shown that Pumilio is an [RNA-binding protein](#) that interacts specifically with the NRE (18-20). Although Nanos can bind RNA with high affinity, it does so with low specificity (3). It has therefore been proposed that Pumilio acts as a sequence-specific adaptor for Nanos that provides the spatial aspects of *hunchback* translational regulation (19). The mechanisms by which Pumilio and Nanos regulate *hunchback* translation are unclear. Maternally supplied *hunchback* RNA is deadenylated in the presence of Nanos and Pumilio (21) which suggests that NRE-directed removal of the poly (A) tail may lead to translational repression. Ectopic expression of *nanos* in the eye [imaginal disc](#) results in disruption of the development of the adult eye (20). This activity requires endogenous *pumilio* expression. Furthermore, ectopic *nanos* can repress the expression of [reporter genes](#) bearing the NREs, suggesting that the machinery needed for *nanos* -mediated repression may be present in all cells (20). Finally, ectopic *nanos* mediates repression of a reporter gene containing [internal ribosome entry sites](#). This suggests that *nanos*-mediated translational repression may occur in a **5'-cap-independent** manner (20).

When Nanos is ectopically localized to the anterior pole of the embryo, in addition to *hunchback*, translation of the homeodomain transcription factor *bicoid* is also repressed (8). Repression of *bicoid* translation results in the loss of head structures, while repression of *hunchback* translation at the anterior causes a mirror image duplication of the abdomen, resulting in a bicaudal phenotype (8). As with *hunchback*, translational repression of *bicoid* RNA is mediated by a NRE in the *bicoid* 3'-untranslated region, and the [poly \(A\)](#) tail of *bicoid* mRNA is shortened in the presence of *nanos* and *pumilio* (21).

*Cyclin B* also contains a sequence in its 3'-untranslated region that is similar to the *nanos* response element. *cyclin B* is not translated at the posterior pole, and removal of this element is sufficient to activate translation of *cyclin B* mRNA at the posterior pole. The role of *nanos* and *pumilio* in the regulation of *cyclin B* translation has not been demonstrated (22).

## 4. Germ Cell Migration

After fertilization, the *Drosophila* embryo develops as a syncytium where nuclear rather than cellular divisions follow each other rapidly and in synchrony. The primordial germ cells are the first cells to form. They cellularize at the posterior pole of the embryo, encapsulating the pole plasm in the



process (reviewed in Ref. [23](#)). During gastrulation, the germ cells are carried along the dorsal surface of the embryo inside the invaginating posterior midgut pocket. Subsequently, the primordial germ cells migrate through the posterior midgut and move toward the somatic component of the gonad, the gonadal mesoderm. Finally, the germ cells align with the gonadal mesoderm and coalesce to form the embryonic gonad.

Nanos protein and RNA are incorporated into the germ cells as they form at the posterior pole ([6](#)). Nanos protein is present in the germ cells until gonad coalescence ([6](#)). *nanos* is not expressed in the soma, and wild-type germ cells transplanted into *nanos* mutant embryos migrate to and populate the gonad. In *nanos* mutant embryos, the germ cells form normally and are carried into the midgut pocket. However, while wild-type germ cells migrate from the midgut to the mesoderm, *nanos* mutant germ cells fail to migrate appropriately and often remain associated with the midgut ([5](#)). Analysis of *nanos* mutant germ cells demonstrates that several germ-cell markers, which are normally expressed within the gonad, are transcribed prematurely during germ-cell migration ([24](#)). These data suggest a role for *nanos* in the specification or maintenance of the germ-cell identity during migration. Although Pumilio protein is also encapsulated in the germ cells when they form, it is not known if *pumilio* and *nanos* share a common function in the germ cells as they do in abdomen formation.

## 5. Homologues

Based on similarity to the carboxyl-terminal domain, *nanos* homologues have been identified from three other Dipteran species (*Drosophila virilis*, the housefly *Musca domestica*, and the midge *Chironomus samoensis*), the leech *Helobdella robusta*, the African clawed frog *Xenopus laevis*, and the [nematode](#) *Caenorhabditis elegans* ([25-27](#)). Analysis of the distribution of the homologues that have been studied so far suggests that *nanos* is a localized factor important for establishing embryonic polarity and/or development of the germ line. The expression pattern of *nanos* RNA is virtually indistinguishable among the different Dipteran species analyzed ([25](#)). Injection of *nanos* RNA isolated from the different Dipteran species into *D. melanogaster* embryos is able to rescue the abdominal defects associated with loss of *nanos* in *D. melanogaster*, suggesting that these homologues share a common, conserved function ([25](#)).

## Bibliography

1. C. Wang and R. Lehmann (1991) *Cell* **66**, 637–647.
2. E. R. Gavis and R. Lehmann (1994) *Nature* **369**, 315–318.
3. D. Curtis, D. K. Treiber, F. Tao, P. D. Zamore, J. R. Williamson, and R. Lehmann (1997) *Embo J.* **16**, 834–843.
4. A. C. Spradling (1993) In *The Development of Drosophila melanogaster*, Vol. **1** (M. B. Arias and A. M. Bates, eds.), Cold Spring Harbor Press, Plainview, New York, pp. 1–70.
5. A. Forbes and R. Lehmann (1998) *Development* **125**, 679–690.
6. C. Wang, L. K. Dickinson, and R. Lehmann (1994) *Dev. Dyn.* **199**, 103–115.
7. D. St. Johnston and C. Nüsslein-Volhard (1992) *Cell* **68**, 201–219.
8. E. R. Gavis and R. Lehmann (1992) *Cell* **71**, 301–313.
9. M. Hulskamp, C. Schroder, C. Pfeifle, H. Jackle, and D. Tautz (1989) *Nature* **338**, 629–632.
10. V. Irish, R. Lehmann, and M. Akam (1989) *Nature* **338**, 646–648.
11. G. Struhl (1989) *Nature* **338**, 741–744.
12. E. R. Gavis, L. Lunsford, S. E. Bergsten, and R. Lehmann (1996) *Development* **122**, 2791–2800.
13. A. Dahanukar and R. P. Wharton (1996) *Genes Dev.* **10**, 2610–2620.
14. A. Ephrussi, L. K. Dickinson, and R. Lehmann (1991) *Cell* **66**, 37–50.
15. C. A. Smibert, J. E. Wilson, K. Kerr, and P. M. Macdonald (1996) *Genes Dev.* **10**, 2600–2609.

16. D. Tautz and C. Pfeifle (1989) *Chromosoma* **98**, 81–85.
17. R. P. Wharton and G. Struhl (1991) *Cell* **67**, 955–967.
18. P. D. Zamore, J. R. Williamson, and R. Lehmann (1997) *RNA* **3**, 1421–1433.
19. Y. Murata and R. P. Wharton (1995) *Cell* **80**, 747–756.
20. R. P. Wharton, J. Sonoda, T. Lee, M. Patterson, and Y. Murata (1998) *Mol. Cell* **1**, 863–872.
21. C. Wreden, A. C. Verrotti, J. A. Schisa, M. E. Lieberfarb, and S. Strickland (1997) *Dev.* **124**, 3015–3023.
22. B. Dalby and D. M. Glover (1993) *EMBO J.* **12**, 1219–1227.
23. A. Williamson and R. Lehmann (1996) *Annu. Rev. Cell Dev. Biol.* **12**, 365–391.
24. S. Kobayashi, M. Yamada, M. Asaoka, and T. Kitamura (1996) *Nature* **380**, 708–711.
25. D. Curtis, J. Apfeld, and R. Lehmann (1995) *Development* **121**, 1899–1910.
26. L. Mosquera, C. Forristall, Y. Zhou, and M. L. King (1993) *Development* **117**, 377–386.
27. M. Pilon and D. A. Weisblat (1997) *Development* **124**, 1771–1780.

### Suggestions for Further Reading

28. D. Curtis, R. Lehmann, and P. D. Zamore (1995) Translational regulation in development. *Cell* **81**, 171–178.
29. A. Forbes and R. Lehmann (1998) *nanos* and *pumilio* have critical roles in the development and function of *Drosophila* germline stem cells. *Development* **125**, 679–690.
30. D. R. Gallie (1998) A tale of two termini: a functional interaction between the termini of an mRNA is a prerequisite for efficient translation initiation. *Gene* **216**, 1–11.
31. R. Lehmann and C. Nüsslein-Volhard (1991) The maternal gene *nanos* has a central role in posterior pattern formation of the *Drosophila* embryo. *Development* **112**, 679–691.
32. C. A. Smibert, J. E. Wilson, K. Kerr, and P. M. Macdonald (1996) Smaug protein represses translation of unlocalized *nanos* mRNA in the *Drosophila* embryo. *Genes Dev.* **10**, 2600–2609.

## Natural Selection

Natural selection is defined as the differential fitness of genetically distinct individuals or genotypes within a given population. The fitness is defined as the number of offspring per parent who are capable of producing the next generation (1). Differential fitness is caused by differences among individuals in ways such as mortality, fertility, fecundity, mating success, and viability of offspring (2).

In population genetics, natural selection is classified into positive and negative selection, depending on whether the fitness is greater or less than 1. When the fitness is larger than 1, a selective advantage is conferred to the individuals or genotype. It follows that Darwinian evolution is operating; this is positive selection. On the other hand, when the fitness is less than 1, the individuals or genotype have a selective disadvantage. This means that deleterious mutations are selected out; this type of selection is negative selection, or purifying selection. When the fitness is just 1, this means selectively neutral (see [Neutral Mutation](#)).

Three major factors make important contributions to the genetic changes of a population, namely, [evolution](#): [mutation](#), natural selection, and [genetic drift](#). Selectionists believed that natural selection

is the most important because it acts like an “editor” of genetic variation; that is, it eliminates deleterious **alleles** and captures advantageous ones. On the other hand, neutralists contend that the effect of genetic drift is substantial, particularly at the molecular level (see [Genetic Drift](#)) (3). Both selectionists and neutralists agree that natural selection is the most important for morphological evolution. Both groups implicitly agree that neither natural selection nor genetic drift is effective without mutation, because mutation is the only way to produce genetic variability.

### Bibliography

1. J. F. Crow and M. Kimura (1970) *An Introduction to Population Genetics Theory*. Harper and Row Publishers, New York.
2. W. H. Li and D. Gauer (1991) *Fundamentals of Molecular Evolution*. Sinauer Associates, Sunderland, Mass.
3. M. Kimura (1983) *The Neutral Theory of Molecular Evolution*, Cambridge Univ Press, Cambridge, UK.

### Near-Field Scanning Optical Microscopy

Near-field scanning optical microscopy (NSOM) is a relatively new subfield of [microscopy](#) that allows visualization with sub-wavelength resolution, also called super-resolution (1). NSOM is a combination of optical microscopy with [scanning probe techniques](#). The principle is based on a probe consisting of an aperture smaller than the wavelength of light that is positioned in close proximity (near field; < 10nm) to the specimen. By laterally scanning the specimen with near-field optics, one can generate an image at a resolution dependent on only the probe size and the probe-to-specimen separation, each of which can be pushed into the nanometer regime. Thus, NSOM extends light microscopy to a higher spatial resolution. Since the achievable resolution degrades with increasing distance from the probe, it is possible to obtain super-resolution information in three dimensions only within a few tens of nanometers from the surface. NSOM therefore falls into the class of surface probe microscopies, eg, atomic force microscopy. However, the surface-penetrating power is sufficient to map cytoskeletal structures in addition to cell [membranes](#) (2) and [membrane proteins](#) (3). Super-resolution surface features can be obtained from thick as well as thin specimens by filtering (4), although subsurface structures are usually better studied with complementary 3-D microscopies such as confocal light or transmission electron microscopies. Resolution as fine as 30–50 nm has been obtained on biological specimens using NSOM (2). Another important aspect of NSOM is the possibility of sensitive mapping of dynamic events between two fluorophores, eg, connected by a short DNA strand, by using fluorescence resonance energy transfer (FRET; (5)).

An advantage of NSOM over microscopies that offer higher resolution, such as STEM, STM, and AFM, is the wide range of contrast mechanisms available. Contrast in NSOM can be accomplished by absorption, polarization, refractive index, fluorescence, and reflectivity (1). From the perspective of molecular biology, fluorescence may be the most powerful contrast mechanism because a wide variety of fluorescent probes have previously been developed for light [microscopy](#). A powerful application of NSOM coupled with fluorescence was demonstrated through the mapping of multiple fluorescent probes for *in situ* hybridization on a single human metaphase [chromosome](#) (6). In another study, dual-color immunofluorescence labeling coupled with NSOM was used to detect colocalized proteins and to simultaneously map membrane proteins to membrane topography (3). The power of NSOM was also shown by resolving previously unreported structural details of biological membranes including gradual domain boundaries between various [lipids](#) and a fine web structure,

which might reflect nanoscale lipid crystals (2).

## Bibliography

1. E. Betzig and J. K. Trautman (1992) *Science* **257**, 189–195.
2. J. Hwang et al. (1995) *Science* **270**, 610–614.
3. T. Enderle et al. (1997) *Proc. Natl. Acad. Sci.* **94**, 520–525.
4. R. Toledo-Crow et al. (1992) *Appl. Phys. Lett.* **60**, 2957–2962.
5. T. Ha et al. (1996) *Proc. Natl. Acad. Sci.* **93**, 6264–6268.
6. M. H. Moers et al. (1996) *J. Microscopy* **182**, 40–45.

## Necrosis

*Necrosis* was the classically accepted form of [cell death](#) and occurs in response to any *severe* physiological or environmental deviation ie, change in temperature, change in pH, or disruption of the plasma [membrane](#) (1-3). Such environmental deviations lead to a loss of homeostatic control by the cell, causing an influx of ions into the cell and loss of **mitochondrial** function, which results in a reduction of ATP levels (see [Adenylate Charge](#)) and loss of membrane integrity.

Unlike [apoptosis](#) or [programmed cell death](#), necrosis is a passive process that ends with cell rupture and induction of an [immune response](#). The latter is an undesirable outcome. The recruitment of inflammatory **macrophages** to the site of necrotic cells can also lead to damage within otherwise healthy surrounding tissues. Thus, when considering therapeutic treatments that trigger cell death, such as cancer treatment, it is desirable to trigger apoptosis and not necrosis, so as to minimize tissue damage. It is still unclear whether cells *in vivo* do necrose, or whether the necrosis observed is actually secondary necrosis in cells that have undergone apoptosis but have not been phagocytosed, due to the extent of cell death occurring. In areas where necrotic cells are prevalent—ie, centers of solid tumors, sites of infarcts, stroke damage, and rupturing of atherosclerotic plaques—apoptosis also occurs (2, 4-6). It is important to determine whether apoptotic cell death occurs initially, because this pathway, unlike necrosis, is subject to genetic regulation, making inappropriate cell death by apoptosis treatable. Moreover, if apoptosis does occur initially, it may be possible to increase the phagocytic capacity of the surrounding cells, such that all the apoptotic cells are phagocytosed prior to the onset of secondary necrosis. This would then prevent any unnecessary damage to the surrounding tissues through an inflammatory-mediated response.

The distinction between necrosis and apoptosis has been somewhat blurred by the demonstration that the anti-apoptotic gene *Bcl-2* can suppress death in response to signals that are considered to induce only necrotic death, such as cyanide (7, 8). Because cyanide damages the mitochondria, the cell is no longer able to maintain its internal environment, leading to rapid cell lysis. It has been suggested that *Bcl-2* primarily functions in regulating both ion fluxes and **cytochrome c** release in apoptotic cells (9, 10), and it may be because of its localization at the mitochondrial membrane that it has some protective effect against necrosis. Whether for the reasons mentioned above this observation is therapeutically relevant remains to be determined.

Morphologically, cells dying by necrosis show a distinct pattern of cellular breakdown [see [Apoptosis](#), Fig. 1], which consists of two stages, one reversible, the other irreversible (1). In the reversible stage, cells are able to adapt to environmental changes. Morphologically, this is characterized by the dilation of the [endoplasmic reticulum](#), along with a slight clumping of the

nuclear [chromatin](#). Other changes in the cell include cytoplasmic swelling and mitochondrial condensation due to the inner membrane shrinking from the outer. The irreversible stage is reached once the cell is unable to adapt further to the environmental change. This stage can be identified by the mitochondria undergoing “high amplitude swelling,” characterized by the appearance of dense lipid-rich aggregates within the inner membrane. The cell continues to swell, and chromatin takes on the appearance of flocculent masses, which eventually disappear, leaving a nuclear ghost (2). Finally, all membranes rupture and the cell undergoes autolysis and denaturation. Areas of tissue affected by the initial toxic insult often involve large numbers of necrotic cells, and this can result in exudative inflammation in viable tissue nearby, damaging this also. Cellular debris resulting from necrosis is phagocytosed by monocytes.

Overall, necrosis can be viewed as a relatively unhelpful process. Necrotic death results in an inflammatory response that can cause considerable damage to surrounding healthy tissues. From a clinical perspective, necrosis is difficult to prevent, unless it is the result of apoptotic cells that have not been phagocytosed.

### Bibliography

1. B. F. Trump, J. M. Valigorsky, J. H. Dees, W. J. Mergner, K. M. Kim, R. T. Jones, R. E. Pendergrass, J. Garbus, and R. A. Cowley (1973) Cellular change in human disease. A new method of pathological analysis. *Hum. Pathol.* **4**, 89–109.
2. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1980) Cell death: The significance of apoptosis. *Int. Rev. Cytol.* **68**, 251–306.
3. A. H. Wyllie, J. F. Kerr, and A. R. Currie (1972) Cellular events in the adrenal cortex following ACTH deprivation. *J. Pathol.* **106**, p. ix.
4. D. W. Choi (1996) Ischemia-induced neuronal apoptosis. *Curr. Opin. Neurobiol.* **6**, 667–672.
5. S. Bjorkerud and B. Bjorkerud (1996) Apoptosis is abundant in human atherosclerotic lesions, especially in inflammatory cells (macrophages and T cells), and may contribute to the accumulation of gruel and plaque instability. *Am. J. Pathol.* **149**, 367–380.
6. M. J. Mitchinson, S. J. Hardwick, and M. R. Bennett (1996) Cell death in atherosclerotic plaques. *Curr. Opin. Lipidol.* **7**, 324–329.
7. S. Shimizu, Y. Eguchi, W. Kamiike, H. Matsuda, and Y. Tsujimoto (1996) Bcl-2 expression prevents activation of the ICE protease cascade. *Oncogene* **12**, 2251–2257.
8. S. Shimizu, Y. Eguchi, W. Kamiike, S. Waguri, Y. Uchiyama, H. Matsuda, and Y. Tsujimoto (1996) Bcl-2 blocks loss of mitochondrial membrane potential while ICE inhibitors act at a different step during inhibition of death induced by respiratory chain inhibitors. *Oncogene* **13**, 21–29.
9. R. M. Kluck, E. Bossy Wetzl, D. R. Green, and D. D. Newmeyer (1997) The release of cytochrome c from mitochondria: a primary site for Bcl-2 regulation of apoptosis. *Science* **275**, 1132–1136 (see comments).
10. J. Yang, X. Liu, K. Bhalla, C. N. Kim, A. M. Ibrado, J. Cai, T. I. Peng, D. P. Jones, and X. Wang (1997) Prevention of apoptosis by Bcl-2: Release of cytochrome c from mitochondria blocked. *Science* **275**, 1129–1132 (see comments).

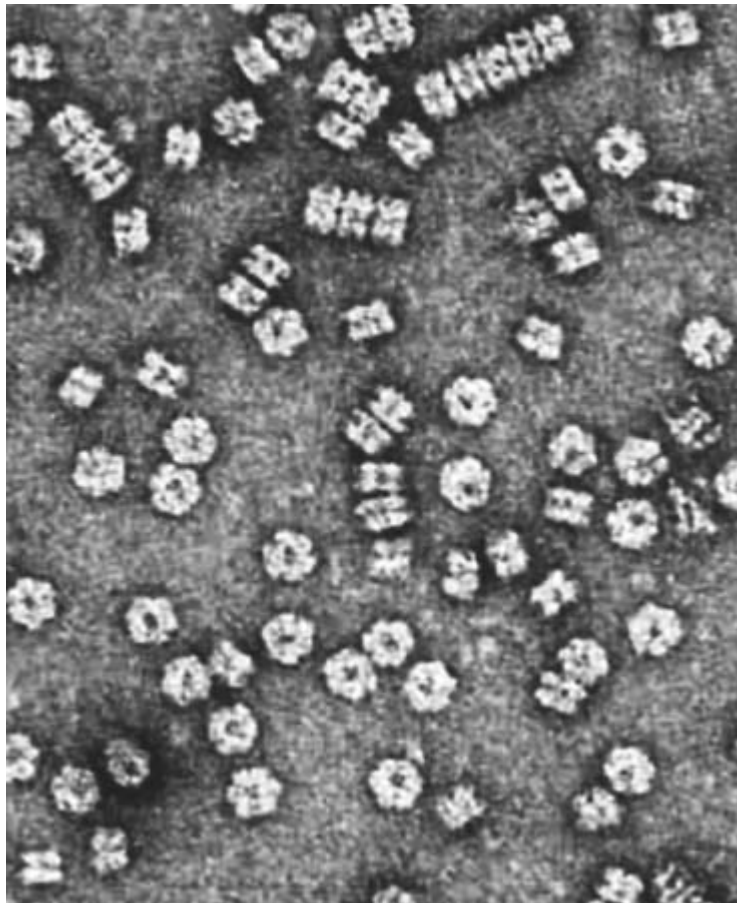
### Negative Stain

Negative staining is an excellent method for the preservation and contrast enhancement of objects

ranging in size from **macromolecules** to entire organisms, such as bacteria. This method is outstanding among [electron microscopy](#) preparation methods for its easy and rapid procedure (1). In addition to increasing the contrast, negative stain (i) also protects the specimen from distortion caused by dehydration as the stain displaces the water surrounding the specimen, and (ii) stabilizes against collapse in the vacuum of the electron microscope. However, care must be observed when choosing a stain since the preparation protocol and choice of stain can produce different results. Furthermore, the wrong stain can disorder the structure.

Staining is achieved by using heavy metal salts which have the property of drying from an aqueous solution to an amorphous “glass.” The metal salt solution is able to occupy the hydrated regions in and around the particle, and as the solution dries it forms an amorphous electron-dense replica of the particle. The image formed is light in areas occupied by biological structure (said to be stain-excluding) and dark in areas occupied by stain. Hence, the technique is called “negative” staining, since the image has the appearance of a photographic negative. The degree of permeation for a particular stain may be influenced by local charge distributions and variations in the degree of hydrophilicity of the specimen (2). A number of heavy metal stains have been employed, but uranyl acetate, phosphotungstic acid, and ammonium molybdate are the most commonly used because they exhibit excellent structural detail. The high contrast of negatively stained preparations even permits the visualization of individual protein molecules (Fig. 1) that are relatively small (>20 Kdal). Negative stain embedment also reduces the problem of radiation damage, since the negative stain replica, which is more resistant to electron radiation, supports the biological specimen. However, at higher electron doses, stain redistribution can cause a loss of resolution (3).

**Figure 1.** Molecules of *E. coli* glutamine synthetase negatively stained with uranyl acetate. This enzyme is composed of 12 identical subunits of 50,000 MW each arranged at the vertices of two eclipsed hexagons to produce a structure with 622 symmetry with dimensions of 145 Å by 95 Å. Notice the well-defined features of isolated molecules many of which show six-fold symmetry or two-fold symmetry depending on orientation. Glutamine synthetase molecules have a tendency to aggregate along the six-fold axes forming stacks of molecules which can be seen in this image.



Negative stain is most useful for revealing (i) the molecular envelope in projection and in 3-D reconstruction, (ii) the orientation of resolvable subunits in complexes, and (iii) the interactions and interfaces between molecules in assemblies and crystals. One of the downfalls of negative staining is that the resolution limit for most specimens is  $\sim 1.5$  nm, even when the structure possesses higher order. This limit is probably set by the grain size of the stain. However, there have been reports of images of negatively stained specimens with a resolution of 1.0 nm or better. This increased resolution is usually accomplished at cryotemperatures and with mixtures of negative stain and sugars or aurothioglucose (4-6), although a resolution of 0.8 nm has also been reported in negative stain without adding sugar or resorting to cryotemperature (7). In limited cases, negative stain can produce higher resolution than vitreous ice embedment, as shown for crystals of cytochrome *c* oxidase (6) in which 0.7 nm resolution was achieved in negative stain, while only 1.0 nm resolution was achieved with ice-embedded crystals. The lower resolution in ice may be ascribed to beam-induced specimen movement or specimen charging. The ease of negative stain makes it useful for surveying crystallization conditions and the degree of order of newly grown crystals, which may be used to ascertain if a higher resolution medium is warranted (8). When judging between the relative merits of negative stain and of a higher resolution medium, at times it is important to ask, "What level of structural detail will answer the research question?" For example, negative stain was chosen over ice embedment for crystals of biochemically split gap junctions as the best medium to enhance the contrast in order to visualize the extracellular domains that had been previously unseen (9).

#### Bibliography

1. M. J. Dykstra (1993) *A Manual of Applied Techniques for Biological Electron Microscopy*, Plenum Press, New York.
2. J. R. Harris (1991) *Electron Microscopy in Biology*, IRL Press, Oxford.
3. A. Bremer et al. (1992) *Ultramicroscopy* **46**, 85–111.

4. R. Rachael et al. (1986) *Ultramicroscopy* **20**, 305–316.
5. J. Lembcke and F. Zemlin (1990) *12th International Conference on Electron Microscopy*, San Francisco, Springer-Verlag, Berlin, pp. 102–103.
6. J. M. Valpuesta, R. Henderson, and T. G. Frey (1990) *J. Mol. Biol.* **214**, 237–251.
7. A. Olofsson, V. Mallouh, and A. Brisson (1994) *J. Struct. Biol.* **113**, 199–205.
8. G. A. Perkins et al. (1994) *J. Struct. Biol.* **113**, 124–134.
9. G. A. Perkins, D. A. Goodenough, and G. E. Sosinsky (1997) *Biophysical J.* **72**, 533–544.

## Nematodes

Nematodes, commonly known as roundworms, are one of the most successful groups of animals on earth. It is estimated that they represent several percent of the total biomass, and that 80% of all living metazoans are nematodes (1). Their phylum, the Nematoda, includes free-living, plant-parasitic, and animal-parasitic species (see *Ascaris*). Adult nematodes range in length from less than 1 mm for some plant and free-living species to greater than 4 m for a whale intestinal parasite.

### 1. Lifestyles

Parasitic nematodes cause extensive economic loss and human suffering. Crop damage from the plant parasites in the United States alone is estimated at more than \$5 billion annually. Livestock production is adversely affected by the animal parasites, which also are present in about 25% of the human population. Hookworm, pinworm, and filarial infections such as African river blindness and elephantiasis are common human diseases caused by parasitic nematodes, primarily in developing countries.

The majority of nematode species are free-living, found worldwide in soil as well as marine and aquatic environments, and generally harmless. One free-living soil species, *Caenorhabditis elegans*, has become one of the best understood “model organisms” for biological research, through extensive study of its genetics, anatomy, physiology, and development (see *Caenorhabditis*).

### 2. Anatomy and Physiology

In spite of their widely differing sizes and lifestyles, all nematodes are similar in their anatomy and physiology. All are simple animals with relatively few cells, and this phylum is one of only a very few that include species with fixed cell numbers and an invariant pattern of cell divisions during [development](#). The basic nematode body plan can be simply described as a tube within a tube. The outer tube is an epidermal layer of cells (hypodermis), which secretes a tough, flexible **collagenous** cuticle covering the animal's exterior. Associated with the inside of this tube are neurons and body-wall muscles. The inner tube is the digestive tract, including a muscular pharynx as well as the intestine, which runs from the posterior end of the pharynx to the anus. Between the two tubes is a space called the pseudocoelom, which is occupied by the gonad in adults. In free-living species, the pharynx is used to ingest and crush bacteria; in plant parasites it is specialized for boring into root tissue. Nematodes lack skeletal components and maintain their body shape by hydrostatic pressure.

### 3. Life Cycle and Development

Nematode embryos develop inside a chitinous egg shell, which forms shortly after fertilization. After

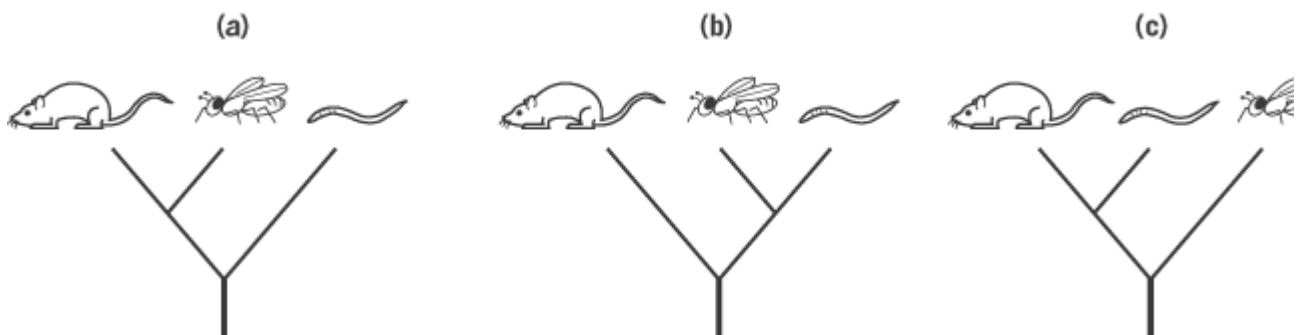


hatching as first-stage juveniles, the worms grow through three more juvenile stages, separated by molts at which a new cuticle forms and the old one is shed. The fourth-stage juveniles molt to sexually mature adults, which continue growth to varying extents in different species. The four juvenile stages are commonly referred to as larval stages (L1–L4), although this nomenclature is not strictly correct because there is no metamorphosis to adulthood. In parasitic species, particular stages are often specialized for growth in a different host or host tissue. Free-living species are capable of molting under adverse conditions to an alternative form of the L3, called a *dauer* (enduring) larva, which is resistant to environmental stress, does not feed, and can survive for much longer than the normal lifespan. When conditions improve, the dauer larva can molt to the L4 and resume the normal life cycle (see *Caenorhabditis* ).

#### 4. Evolution

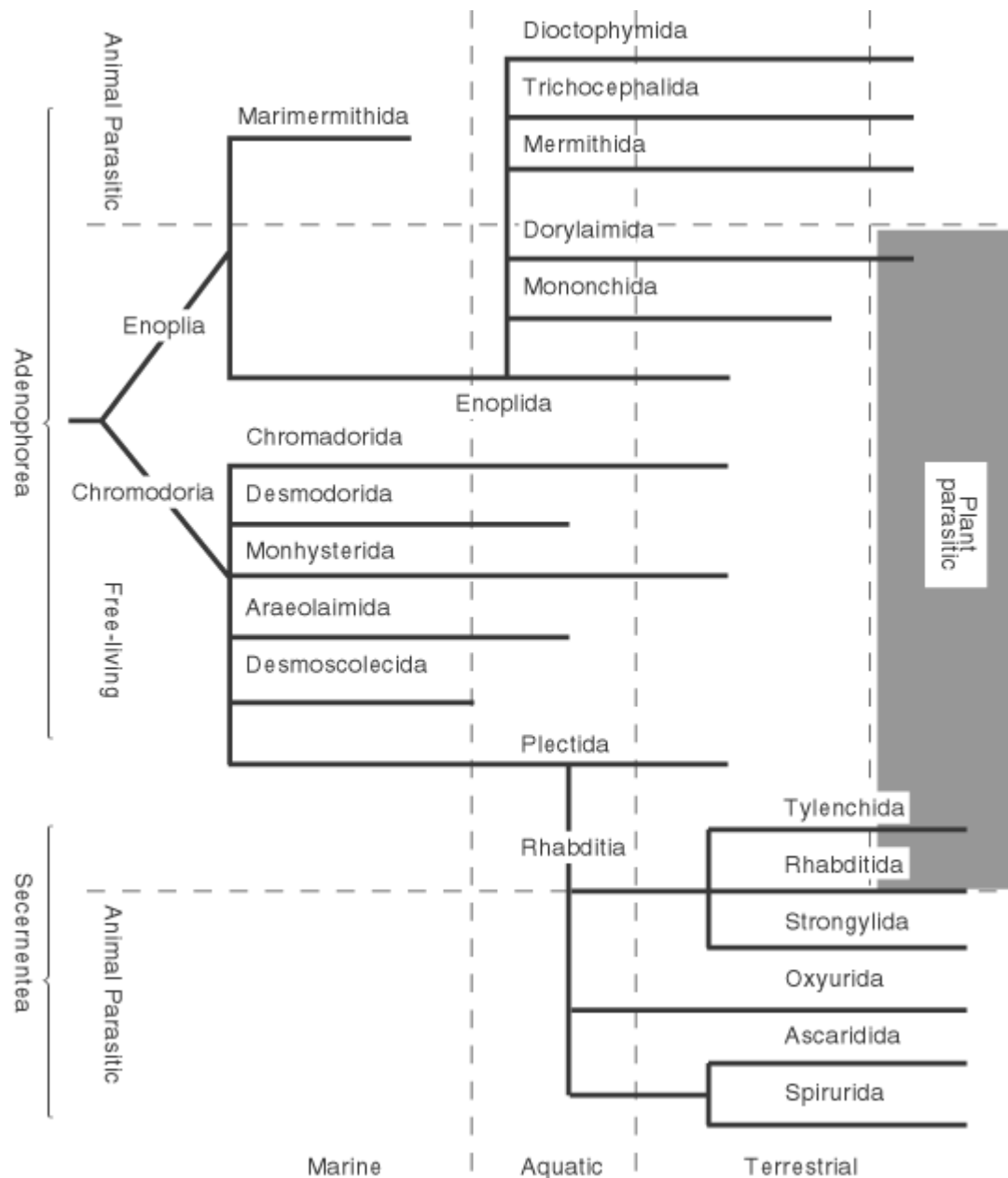
The evolutionary relationship of nematodes to other animal phyla is still a subject of debate. It has been difficult to establish, partly because there is no nematode fossil record. Figure 1 shows the three possible evolutionary relationships between nematodes, arthropods (represented by *Drosophila*), and chordates (represented by the mouse). Based on present evidence, Figure 1a seems most likely and Figure 1c least likely; however, an argument has been made for Figure 1b based on molting as a common character of nematodes and arthropods (2).

**Figure 1.** The three possible evolutionary relationships between chordates, represented by the mouse, arthropods, represented by the fruit fly *Drosophila*, and nematodes. Vertical axis represents evolutionary time. See text for further explanation. (Modified from Ref. 3.)



Within the phylum Nematoda, there are two major classes, termed Adenophorea and Secernentea. Figure 2 shows one possibility, based primarily on morphological characters, for the relationships of the various nematode orders, with their habitats and free-living or parasitic lifestyles indicated (3). A new classification based on ribosomal RNA sequences, currently in progress, suggests somewhat different relationships (M. Blaxter, personal communication).

**Figure 2.** Possible evolutionary relationships of nematode orders and their principal ecological niches. Multifurcations represent not simultaneous divergence events but rather uncertainty regarding sequences of divergence. Not indicated under “Marine” are several orders of the Secernentea that are parasites of marine vertebrates. (From Ref. 3, with permission.)



## Bibliography

1. H. M. Pratt (1994) In *The Phylogentic Systematics of Free-Living Nematodes* (S. Lorenzen, ed.), The Ray Society, London, pp. i–ii.
2. A. M. A. Aguinaldo et al. (1997) *Nature* **387**, 489–493.
3. D. H. A. Fitch and W. K. Thomas (1997) In *C. elegans II* (D. L. Riddle, T. Blumenthal, B. J. Meyer, and J. R. Priess, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 815–850.

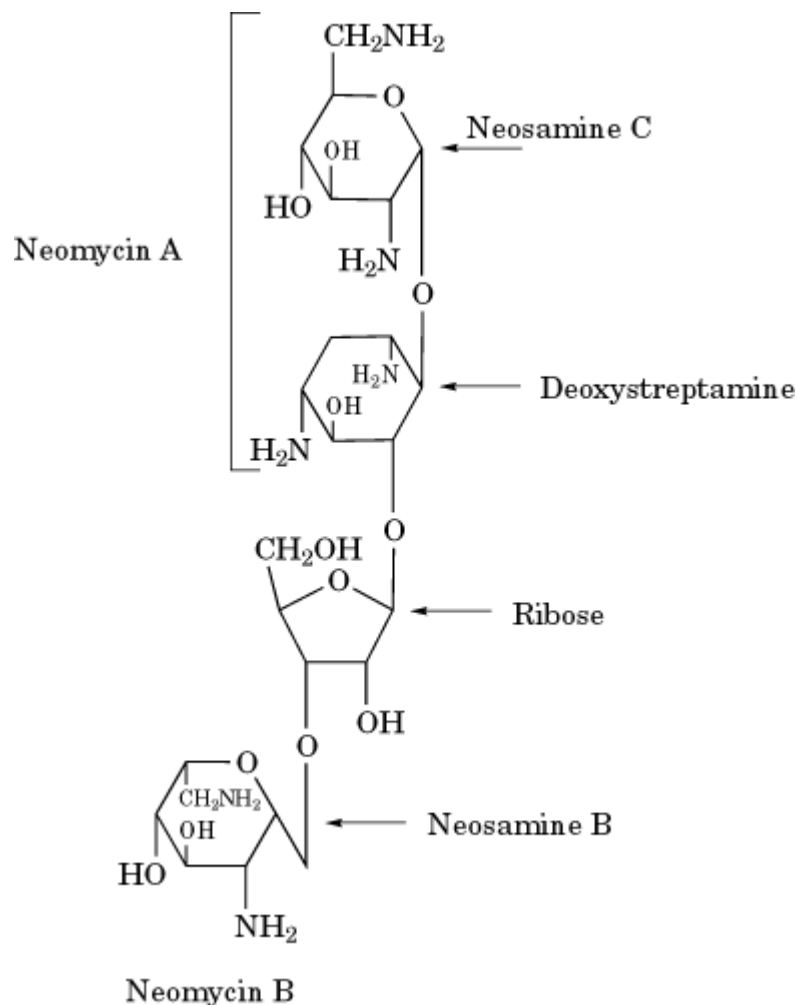
## Suggestions for Further Reading

4. A. F. Bird and J. Bird (1991) *The Structure of Nematodes*, 2nd ed., Academic Press, San Diego, CA, pp. 316.
5. W. B. Wood et al., eds. (1988) *The Nematode Caenorhabditis elegans*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 667.

## Neomycin

Neomycin is a polycationic aminoglycoside produced by actinomycetes that was first isolated in 1949 from the soil-dwelling microbe *Streptomyces fradiae* by Waksman and Lechevalier (1). As isolated, neomycin is a mixture composed primarily of two stereoisomers, neomycin B and C (Fig. 1), both of which possess antibacterial activity. In addition, there is a third component, neomycin A, that is generated by the hydrolysis of neomycins B and C. There are a number of other antibiotics which fall into the same class as neomycin, including [streptomycin](#), gentamycin, [kanamycin](#), tobramycin, and netilmicin.

**Figure 1.** Neomycins A, B, and C. The structure of neomycin B is shown in full. Neomycin C is the same, except that a second neosamine C moiety replaces the neosamine B. Neomycin A is produced by hydrolysis of neomycins B and C and consists of just the neosamine C and deoxystreptamine units.



### 1. Clinical Uses

Neomycin has been used as both an oral and a topical antibiotic and is effective against most

clinically important aerobic **gram-negative** bacteria. It has, however, only minor effects on *Streptococci* and gram-positive *Bacilli*. It was used originally to treat tuberculosis. This treatment was discontinued when it was realized that there were serious side effects (discussed below). Neomycin has also been given in combination with niacin to decrease serum cholesterol and lipoprotein A concentrations (2). Again, side effects have prevented it from being used clinically as a cholesterol-lowering agent. Neomycin sulfate is, however, widely used both as an over-the-counter topical agent for cuts and wounds and microbial skin infections and as a preparative treatment in bowel surgery.

## 2. Mechanisms of Toxicity

The side effects of neomycin administration in humans have included both nephrotoxicity and ototoxicity (involving the auditory and vestibular functions of the eighth cranial nerve). Neomycin has been found to inhibit both the renal and cochlear forms of [ornithine decarboxylase](#), enzymes important for polyamine biosynthesis (3). Use of neomycin as a cholesterol-lowering drug was discontinued due to its side effects of ototoxicity and diarrhea (2). Neomycin also induces acute neuromuscular blockade, perhaps by displacing bound calcium and inhibiting the prejunctional release of the neurotransmitter **acetylcholine** (4).

## 3. Neoresistant Gene and Use in Molecular Biology

Although limited in its clinical usefulness, neomycin has in recent years become a major tool in molecular biology. Neomycin inhibits the initiation of **protein biosynthesis** by binding to the 30S [ribosome](#), thereby blocking [translation](#). The real usefulness of neomycin as a tool in molecular biology came with the discovery of the mammalian [antibiotic resistance](#) gene (*NeoR*) (5). This gene codes for an enzyme, aminoglycoside phosphotransferase, that catalyzes the phosphorylation of neomycin, gentamycin, and kanamycin, which inactivates them as inhibitors of protein synthesis. Expression of the *NeoR* gene, therefore, blocks the toxic effects of the aminoglycosides. Growth of transfected cells in the presence of neomycin allows for selection of cells that are expressing *NeoR* coupled to other genes of interest. *NeoR* and other drug-resistant genes are being considered for use in gene therapy to protect patients from the hematological side effects of chemotherapy (6). This new approach involves the use of a chimeric **vector** that fuses the antibiotic resistance genes with the genes of therapeutic value.

## 4. Interaction with Ribosomes and Protein Synthesis

Aminoglycosides such as neomycin B are bactericidal in that they prevent translation of [messenger RNA](#) by prokaryotic ribosomes. RNA molecules that recognize neomycin share a common stem-loop structural motif (7) (see [RNA Structure](#)). Recent studies have suggested that aminoglycosides inhibit protein synthesis by direct interaction with the ribosomal RNA (8). The nuclear magnetic resonance (NMR) structure of paromycin (a member of the neomycin family) bound to the A site of *E. coli* 16S rRNA has been analyzed. Aminoglycosides such as neomycin bind to the A site of rRNA (9), which is the entry site for aminoacyl tRNAs. Binding of the aminoglycoside interferes with protein synthesis by causing **codon** misreading (10). The translational block may arise from stabilization of a unique high-affinity conformation of rRNA in the tRNA–mRNA complex, thereby perturbing the accuracy of protein synthesis (11).

## 5. Interaction with Lipids

Neomycin also appears to inhibit the [phosphatidyl inositol](#) (PI) cycle (12, 13). [Platelet-Derived Growth Factor](#) (14), [thrombin](#) (12), and a number of other molecules that bind to certain surface **receptors** mediate their cellular effects through activation of [phospholipase C](#) (PLC), leading to the cleavage of phosphatidylinositol 4,5-bis-phosphate (PIP<sub>2</sub>). Neomycin B indirectly inhibits the activity of PLC. PLC normally cleaves PIP<sub>2</sub> to generate 1,4,5-triphosphate (IP<sub>3</sub>) and 1,2-

diacylglycerol (DAG), which are involved in activation of the IP<sub>3</sub> receptor and protein kinase C (PKC), respectively. Neomycin binds to PIP<sub>2</sub> and blocks the cleavage of PIP<sub>2</sub> by PLC. Because PIP<sub>2</sub> is a cofactor in the activation of phospholipase D, neomycin also inhibits its activity and the subsequent production of phosphatidic acid (15).

## 6. Ion Channels and Calcium

Neomycin has been shown to inhibit a number of **ion channels**, including surface membrane Ca<sup>2+</sup> channels (16, 17), two intracellular Ca<sup>2+</sup> release channels, the ryanodine receptor (RyR) (18, 19), and the IP<sub>3</sub> receptor (20). K<sup>+</sup> channels, such as the ATP-sensitive K<sup>+</sup> channel (21) and the Ca<sup>2+</sup>-activated maxi-K<sup>+</sup> channel (22), also appear to be inhibited by neomycin.

The specific interaction of neomycin with an ion channel has probably been most thoroughly studied with the skeletal muscle Ca<sup>2+</sup>-release channel, also known as RYR1. Neomycin has been suggested to inhibit the Ca<sup>2+</sup>-release channel by enhancing its rate of inactivation (23), and it appears to bind to a location in the carboxy-terminal 20% of the molecule (19). There may actually be multiple binding sites for neomycin on RYR1 (18). Neomycin has been used as a tool to transfer a **fluorescent** label to its binding site on RYR1 to determine whether this region of the protein is involved in the conformational changes induced by various channel activators (24).

A specific binding motif for neomycin on the target proteins has not yet been identified. Neomycin may be binding to the binding sites for either Ca<sup>2+</sup> or Mg<sup>2+</sup>-ATP (25) on its target proteins. The identities of the actual binding sites and the specificities of these sites are unknown. Nevertheless, neomycin is thought to mediate its effect through direct interaction with ion channels and continues to be used as a tool to study ion channel properties.

## 7. Summary

Although of limited clinical value, neomycin has been an extremely valuable tool for basic science research. Its use as an inhibitor of ion channels and of the PI cycle has helped to elucidate the molecular events involved in these processes. The discovery of the antibiotic resistant genes has led to the widespread usage of neomycin and other aminoglycosides to select cells containing exogenous genes. The success of the neomycin selection technique in molecular biology and in the generation of **transgenic** mice has led to the possibility of using this approach in gene therapy (6).

## Bibliography

1. S. A. Waksman and H. A. Lechevalier (1949) *Science* **109**, 305–307.
2. S. A. Spinler and M. J. Cziraky (1994) *Ann. Pharmacother.* **28**, 343–351.
3. C. M. Henley, L. G. Mahran, and J. Schacht (1988) *Biochem. Pharmacol.* **37**, 1679–1682.
4. C. Pittinger and R. Adamson (1972) *Ann. Rev. Pharmacol.* **12**, 169–184.
5. P. J. Southern and P. Berg (1982) *J. Mol. Appl. Genet.* **1**, 327–341.
6. T. Licht et al. (1997) *Stem Cells* **15**, 104–111.
7. M. G. Wallis et al. (1995) *Chem. Biol.* **2**, 543–552.
8. D. Fourmy et al. (1996) *Science* **274**, 1367–1371.
9. J. Woodcock et al. (1991) *Eur. Mol. Biol. Organ. J.* **10**, 3099–3103.
10. P. Edelmann and J. Gallant (1977) *Cell* **10**, 131–137.
11. R. Karimi and M. Ehrenberg (1994) *Eur. J. Biochem.* **226**, 355–360.
12. D. H. Carney et al. (1985) *Cell* **42**, 479–488.
13. F. S. Vassbotn, N. Langeland, and H. Holmsen (1990) *Biochim. Biophys. Acta* **1054**, 207–212.

14. K. M. Hedberg and P. B. Bell Jr. (1995) *Exp. Cell Res.* **219**, 266–275.
15. M. Liscovitch (1996) *J. Lipid Mediators Cell Signalling* **14**, 215–221.
16. G. Suarez-Kurtz (1989) *Acta Physiol. Pharmacol. Latinoam.* **39**, 407–418.
17. K. Yamada et al. (1993) *Jpn. J. Pharmacol.* **63**, 423–432.
18. M. M. Mack, I. Zimany, and I. N. Pessah (1992) *J. Pharmacol. Exp. Ther.* **262**, 1028–1037.
19. J. P. Wang et al. (1996) *J. Biol. Chem.* **271**, 8387–8393.
20. L. G. Sayers and F. Michelangeli (1993) *Biochem. Biophys. Res. Commun.* **197**, 1203–1208.
21. X. Lin, R. I. Hume, and A. L. Nuttall (1993) *J. Neurophysiol.* **70**, 1593–1605.
22. S. Takeuchi and P. Wangemann (1993) *Hearing Res.* **67**, 13–19.
23. W. Wyskovsky et al. (1990) *Eur. J. Biochem.* **194**, 549–559.
24. J. J. Kang et al. (1992) *Biochemistry* **31**, 3288–3293.
25. S. Kocabiyik and M. H. Perlin (1992) *FEMS Microbiol. Lett.* **93**, 199–202.

### Suggestions for Further Reading

26. H. A. Lechevalier (1975) The 25 years of neomycin. *Crit. Rev. Microbiol.* **3**, 359–397.
27. J. Schacht (1993) Biochemical basis of aminoglycoside ototoxicity. *Otolaryngol. Clin. N. Amer.* **26**, 845–856.
28. H. F. Chambers and M. A. Sande (1996) "Antimicrobial agents: the aminoglycosides. In" *Goodman and Gilman's The Pharmacological Basis of Therapeutics* (J. G. Hardman, L. E. Limbird, P. B. Molinoff, R. W. Ruddon, and A. G. Gilman, eds.), McGraw-Hill Health Professions Division, San Francisco, pp. 1103–1122.

## Neoplastic Transformation

The term neoplastic transformation first came into common use in the 1950s to refer to changes in cell morphology and growth behavior in **cell culture** related to those described for tumors in animals. The changes in culture came about either spontaneously during long-term culture of mouse embryo cells or as a result of infection of chicken embryo cells with [Rous sarcoma virus \(RSV\)](#). The term was then applied retrospectively to the changes occurring in the development of cancer itself in animals and to a lesser degree in **plants**. During the 1970s, individual transforming **genes** were identified with different **RNA**-containing tumor-inducing **viruses**, including RSV, known collectively as [retroviruses](#). The RNA of these viruses was **transcribed** during infection into **DNA** and integrated into the [genome](#) of the host cell, where it brought about its transforming action. The DNA of the transforming genes, or **oncogenes**, was used to "**transfect**" normal chicken embryo cells and mouse embryo cells and to transform them into neoplastic cells in the same way as the virus itself did. The DNA of normal vertebrate cells contains nucleotide sequences similar to those of the retrovirus oncogenes, and a mutated form of one of these was isolated from the DNA of a cell culture derived from a human bladder cancer. It was inferred that oncogenes were mutated forms of cellular genes related to viral oncogenes, and they played an important part in the development of human cancer. This inference was related to much earlier results in chemical **carcinogenesis** in experimental animals, particularly mice, which suggested that the primary event in the development of a tumor was a [mutation](#) that then required a series of other events to bring the tumor to full expression. A later development was the concept of [tumor suppressor genes](#). As methods are refined for detecting genetic changes in tumors, large numbers of such changes can be found in a single tumor, particularly in solid cancers of humans, and the role of any particular genetic change in

development of the cancer is uncertain. The sum of these observations indicates that the development of cancer is an extremely complex process involving many genetic, as well as [epigenetic](#), changes in the host.

## 1. Experimental Carcinogenesis in Animals

The first strong experimental indication that mutation plays a primary role in neoplastic development came from the application of polycyclic aromatic hydrocarbons (PAHs) to the skin of mice. The earliest work in rabbits and mice indicated that repeated application of a PAH was necessary to induce a tumor. In the 1940s, it was discovered that the combined application of a PAH and a non-**carcinogenic** substance (ie, croton oil) was more effective than the PAH alone. In fact, a single application of a PAH to the skin of a mouse produced no tumors, unless it was followed by repeated applications of croton oil over many weeks. If the repeated application of the croton oil was made before the PAH, no tumor resulted. A key finding was that the interval between the single PAH application and the croton oil could be as long as 10 months and a positive result would still be obtained. It was therefore concluded that the single dose of the PAH induced a permanent genetic change, called *initiation*, that required the promoting action of croton oil or further PAH application to produce a tumor. Because the croton oil was itself noncarcinogenic, *promotion* was considered an epigenetic or physiological effect. The most active ingredient of croton oil was 12-*O*-tetradecanoyl-phorbol-13-acetate (TPA), more recently termed phorbol-13-myristate-12-acetate (PMA) (see **Phorbol esters**).

The actions of initiators and promoters are subsumed under the rubric of the two-stage theory of carcinogenesis. While the two-stage theory has gained wide acceptance, it is not without its critics who interpret results such as the seeming unidirectional action of initiators and promoters in a quite different way, that is, “carcinogenesis is most probably composed of a multifarious combination of both genetic and epigenetic phenomena” (1).

The two-stage theory of carcinogenesis is based largely on studies of carcinogenesis of mouse skin. Another carefully studied model of carcinogenesis is in the rat liver. There a carcinogenic procedure consisting of systemically inoculating a carcinogen, such as dimethylhydrazine or 2-acetylaminofluorene, that does not initiate in the same way as the PAHs do in the skin, but inhibits proliferation of the vast majority of hepatocytes and selects a few resistant hepatocytes for proliferation. When coupled with mitotic stimulation, such as removal of part of the liver (partial hepatectomy), many nodules of altered, proliferating cells are formed. If the hepatectomy is delayed for 48 h or longer, the number of nodules is greatly reduced. In any case, 90% to 98% of the nodules redifferentiate into normal-appearing liver tissue. Changes typical of cancer appear only in a small fraction of the nodules. The rapid decrease in the number of nodules with increasing delay in applying the mitotic stimulation, along with the subsequent redifferentiation of most of the nodules, has suggested that the initiating treatment does not directly result in mutational change but is a physiological adaptation resulting in hepatocytes that resist growth-inhibitory toxic material. The cells in these nodules have an increased probability of undergoing the genetic changes that eventually result in liver cancer (2). It is apparent therefore that the simple relation between initiation/mutation and promotion that arose from PAH carcinogenesis in mouse skin is inadequate as a generalization for all forms of cancer.

## 2. Mutagenic Action of Carcinogens

The two-stage model of carcinogenesis led to attempts to demonstrate the mutagenic action of known carcinogens, particularly the PAH compounds, which were the only ones that regularly caused neoplastic transformation when painted on the skin. Attempts to demonstrate their mutagenicity in the fruit fly *Drosophila melanogaster*, a favorite organism for genetic studies, were largely negative or indecisive. Beginning in the 1950s, however, genetic studies turned increasingly to the use of **microorganisms**, which would provide many millions of dividing organisms that would exhibit mutational change within a day or two. A test was developed for mutations in the

bacterium *Salmonella typhimurium* that seemed capable of detecting about 90% of the animal carcinogens (see [Ames Test](#)). The small number of presumed nongenotoxic, epigenetic carcinogens were considered a negligible problem. Broader studies later reduced the sensitivity of the assay to 53%, and they also reported that about 30% of noncarcinogens were mutagenic in the test. Further study confirmed these results and recommended dividing carcinogens into genotoxic and nongenotoxic categories (3). The picture is further complicated by the fact that a diet deficient in choline or methionine is carcinogenic. Some interpretations of these results could be used to question the significance of genetic change in the origin of cancer. However, such doubts are countered by the fact that the *Salmonella* test detects only [base pair substitutions](#) and [frameshift mutations](#) and therefore would not detect large-scale [chromosome](#) rearrangements and deletions. Such changes are, in fact, considered to be the most important ones in the origin of human cancer (4) and will be discussed further in the text below.

There is an efficient assay that detects larger-scale changes in the genome and that had identified many carcinogens that were negative in the *Salmonella* test. It uses a line of mouse lymphoma cells that are heterozygous for a gene that phosphorylates thymidine (5). When the active allele becomes inactive as a result of genetic alteration, the cell can survive under certain conditions that kill the unaltered cells containing the active allele. The inactivation can occur either by local (intragenic) changes in base sequence or by large-scale deletions or rearrangements of chromosomes. The latter types of change can be distinguished from the former because they also cause an inheritable reduction in growth rate of the altered cells. Probably because of the dual response, this assay has detected several important carcinogens (such as the [hormone](#) diethylstilbesterol) that were negative in the *Salmonella* test, and it exhibits a good quantitative correlation between mutagenicity and carcinogenicity in experimental animals (5). Many of the compounds that are positive in this assay do not interact directly with DNA, but their ultimate effect is to alter the genome regardless of their initial target. These would be termed epigenetic or *dysgenetic* agents (6).

### 3. Viral Transformation of Cells

The first clear-cut demonstration that some viruses could induce neoplastic transformation was the induction of sarcomas in chickens by inoculation of filtrates from a naturally occurring connective tissue tumor of chickens that had first been transplanted serially by intact cells in closely related chickens. The virus was later named [Rous sarcoma virus](#) (RSV) after its discoverer and was the source of most of the quantitative work done on viral carcinogenesis in animals for many years. Single infectious particles of RSV could initiate epithelial tumors on the chorioallantoic membrane of the developing chicken embryo, and progeny virus particles from the epithelial tumors initiated sarcomas or connective tissue tumors in the tissue underlying the epithelial tumors, or when inoculated into adult chickens. Every infected cell gave rise to a tumor, and unlike the cytotoxic or cell-killing viruses, all the infected cells could proliferate. It was therefore obvious that a small amount of genetic material could initiate neoplastic transformation, the first clear indication that genetic change *could* cause cancer. There was as yet, however, no evidence that some change in the host cell genome itself could induce transformation. The viral genome was made of RNA instead of DNA. However, it was discovered that the virus contained an enzyme, **reverse transcriptase**, that converted the viral RNA into DNA, which would then integrate into the genome of the cell, where the protein readout of the viral oncogene would transform the cell. This confirmed experiments on radiation of cells before infection with RSV which showed that loss of ability of cells to reproduce prevented the establishment of virus infection and suggested that the viral genome had to integrate into cellular DNA before the virus could multiply. Because there was a nucleotide sequence in the normal cellular DNA that was very similar to that of the viral RNA oncogene, the idea arose that a simple mutation in the cellular gene would lead to neoplastic transformation.

The opportunity for precise quantification of viral transformation arose when a method for infecting and transforming chicken embryo cells in culture was developed. Infection resulted in morphological transformation of spindle-shaped fibroblastic or connective tissue cells into a more rounded sarcoma cell. More significant than the morphological change, however, was the altered growth behavior of



the cell. Normal fibroblasts multiply rapidly in culture when they are sparsely distributed on the floor of the dish, but when they become crowded the cell-to-cell contact results in a marked decrease in growth rate (see [Contact Inhibition](#)). Net growth ceases when the cells multiply enough to form a confluent sheet. However, the cells transformed by RSV continue to multiply after they come in contact with other cells. If there is one transformed cell surrounded by normal cells, the latter form a single sheet of cells, but the former continues to multiply into several layers and form what is called a *transformed focus*. If the cells from the focus are dispersed from one another and transferred to a new culture dish along with a large majority of non-transformed cells, many neoplastically transformed foci result, each one like a small tumor.

Some strains of RSV are defective in the sense that they require coinfection of the same cell by a [helper virus](#) in order to produce new infectious virus. The reason for the defect is that the defective strains lack the information to produce the outer coat of the virus, which is required to adsorb to and penetrate into another cell. However, the defective virus can transform the cell to which it has gained entry even when no new infectious virus is produced. Clearly the coat protein does not participate in the transformation itself.

Just as there are a number of transforming retroviruses, there are a number of transforming viral oncogenes, and their protein translation products are associated with different functions in the cell. The transforming *src* gene of RSV codes for a protein [kinase](#) (an enzyme that **phosphorylates** protein), and considerable effort has been expended to find the cellular target of phosphorylation that causes transformation. However, the enzyme is promiscuous in the proteins it phosphorylates, and no single target protein has been found that transforms upon phosphorylation.

There is another oncogene called the *ras* gene, which stands for *rat sarcoma* virus from which it was isolated. A particular substrain, the *Harvey-* or *H-ras* gene, was the first gene from a human tumor shown to have transforming activity on [transfection](#) into a mouse cell culture line. Actually, the gene was only isolated from the tumor cells after they had been several years in culture, and the mutation could have occurred during that period. Later work showed that the *H-ras* gene, even when mutated, did not transform the target mouse cells directly but had first to recombine with a strong **promoter**, which led to a 100-fold increase in the protein product of the gene (7). Because there was no indication of a significant overproduction of the *ras* protein in the original tumor, it is unlikely the mutated gene played a major role in tumor development. It is also noteworthy that recombining the *normal* cellular *ras* gene with a strong promoter also leads to transformation, though with a lower efficiency than the recombined mutant *ras* gene does. The transfection assay in murine cells therefore leads to the questionable conclusion that the mutated *ras* gene derived from human tumor cells was the cause of the tumor, when it was in fact the artifact of [recombination](#) that produced the transformation in the test system. The results, however, do point up the potential significance of chromosome rearrangements in producing new gene combinations that are carcinogenic.

#### 4. Induction of Neoplastic Transformation by Chemical and Physical Carcinogens

Although several unsuccessful attempts were made to transform mouse embryo cells in culture with carcinogenic PAHs (see **Spontaneous transformation** below), the first success was achieved with Syrian hamster embryo cells in 1965. The untreated hamster embryo cells became large, flat, and vacuolated about 3 weeks after they were explanted to culture from the embryo, slowed down in growth, and finally stopped multiplying after 6 to 7 weeks. The main body of PAH-treated cells exhibited the same early degenerative changes as the controls, but there were also foci of criss-crossing slender, spindle-shaped proliferative cells. Upon isolation, these early-appearing proliferative foci could not be subcultured. Their number and size increased with time of incubation in the original culture. The later transformed cells had a random pattern of organization, in contrast to the flattened, degenerating nontransformed cells of the background, and they continued to multiply on subculture. After about 10 weeks in culture, the cells produced slow-growing sarcomas when inoculated into hamsters. Longer periods in culture resulted in faster growing tumors in hamsters. When single cells were treated with PAH, even if only for a few hours, a significant

fraction (2% to 26%) developed into transformed foci after they grew into full colonies. This would constitute a remarkably high frequency of transformation if it were the result of mutation in one or a few genes. There was considerable variation in cell morphology, both within and between the transformed foci. A fairly close relationship existed between the concentration of carcinogen that produced transformation and that which caused an initial inhibition of growth. Noncarcinogenic and weakly carcinogenic PAH produced neither growth inhibition nor transformation. The results suggest a relationship between a degree of cell damage and subsequent transformation, as had been proposed many years earlier. They suggest that PAH treatment introduces genetic instability into cells that results in varied forms of transformation among a fraction of the progeny. Attempts to transform primary or secondary cultures of mouse embryo cells were difficult to interpret, because even the untreated mouse cells underwent transformation. (See **Spontaneous transformation** below.)

Following the success in chemically transforming hamster embryo cells in culture, efforts were renewed with mouse embryo cells, only this time cell lines were used that had been adapted to grow in culture continuously. Single cells from a permanent line of fibroblasts derived from mouse prostate glands were exposed to a carcinogenic PAH. Under optimal conditions of treatment, 100% of the treated cells that could form a clone exhibited morphologic evidence of transformation and produced tumors upon inoculation into hamsters, while only 5% of the untreated clones did so.

Despite this remarkable result, further work with the prostate fibroblasts was discontinued, possibly because they developed higher frequencies of spontaneous transformation. Instead, a permanent line of fibroblasts was developed from embryonic body wall fibroblasts of the C3H mouse strain. These cells formed a stable, flat monolayer of cells when they reached confluence and did not appear to undergo spontaneous transformation. Brief exposure to a carcinogenic PAH led to the formation of transformed foci after a 6-week period of incubation. A peculiar feature of clonally isolated transformed cells was that they grew at a *lower* rate than their parental normal cells when seeded at low population density. This again indicates that transformation is associated with cell damage and will be considered further in the text below. The efficiency of chemically induced transformation in mass populations of the C3H cell line appeared much lower than had been obtained with clones of the prostate fibroblasts. However, later PAH treatments of low densities of the C3H cells resulted in transformation at high but variable frequencies. The same cell line was used to study the transforming effect of exposure to X-rays. These experiments gave the surprising result that the number of transformed foci produced by the X-rayed cells when they reached confluence was independent of the number of cells that was used to start the culture after the irradiation. The authors proposed that X-rays induce an ill-defined change in many or all of the surviving cells that is transmitted to their progeny. This change increases the probability of a second step, overt transformation, when the cells are maintained under the growth inhibitory conditions of confluence. The fraction of confluent cells that undergo transformation is small, but almost all the survivors of X-irradiation give rise to progeny that produce some transformed foci. As in the case of the hamster cells treated with PAH, it appears that the carcinogenic treatment induces a genetic instability in most of the cells that can result in transformation in a small fraction of their progeny.

The very high frequency of the first change, while surprising when contrasted with the rarity of specific genetic mutations, is correlated with several other findings. A similarly high proportion of X-rayed cells suffers a heterogeneous, heritable reduction in growth rate signaled by the formation of small colonies. The X-rayed cells and their progeny also display an increased frequency of mutations and an increased sensitivity to mutagenic treatment. The lesions that produce the instability appear to be related to double-strand breaks in DNA and to the elimination of chromosome fragments. Another corollary is the finding that chromosome aberrations are produced in a high proportion of cells by doses of X-rays only one-tenth the mean lethal dose. The aberrations consist mainly of chromatid and chromosome gaps and are found in several copies per cell as the X-ray dose approaches the mean lethal dose. Thus the genetic lesions at the chromosomal level are likely to be associated with the instability induced by X-rays in whole populations that increases the likelihood of neoplastic transformation.

Another example of a high-frequency transformation occurs when either thyroid or mammary cells are removed from rats, X-irradiated in culture, and reinoculated into rats made susceptible to tumor development. As many as one irradiated thyroid cell out of seven will produce a thyroid cancer, and about 1 in 100 irradiated mammary epithelia produce mammary cancer. Similar results are obtained in this procedure when the cells are treated with the chemical carcinogen and mutagen *methylnitrosourea*. It is likely that the same sequence of changes is produced in these cells as in those that produce transformed foci in culture, namely a pervasive chromosomal destabilization in a cell population that results in much rarer genetic lesions in the progeny that produce tumors.

Exposure of normal fibroblasts from baby mouse skin to white fluorescent light produced a marked increase in the number of transformed foci that appeared in the culture. A conspicuous feature of this treatment was the variation in morphology of cells between foci, contrasted with their relative similarity within a focus. Many combinations of sparse or dense growth, isometric or spindle shape of cells, and large or small, narrow or wide, flat or refractile cells were observed. Also observed were one or more characters commonly associated with spontaneous neoplastic transformation in culture, such as clumping, "cording" in linear arrays, poor spreading, criss-cross or disorderly orientation, multilayering, small cytoplasm, and variable shape. A similar effect was produced by exposing the cells to drugs, such as [colchicine](#), that combine with [tubulin](#) and prevent the formation of spindles that are necessary for chromosome segregation during cell division. In the first few divisions after removal of the antitubulin drugs, many cells became **tetraploid** (double the normal **diploid** number of chromosomes). There was also widespread asymmetric nuclear division, which creates the potential for chromosome loss. The authors felt that the resultant abnormal chromosome composition of cells ([aneuploidy](#) or heteroploidy) along with DNA damage were responsible for transformation of the cells. A large number of the transformed foci were isolated and continued for long-term culture. Although the resulting cultures grew quickly at first, this was followed by slower growth, accompanied by many dying cells. Some of the cultures died out, some grew slowly, and about one out of five became an established cell line. All of them were entirely or mostly [polyploid](#). There appeared to be a continuing but random genetic reshuffling process, with some cell death and selection of viable gene sets. The long-term survival of isolated foci of 20% was orders of magnitude greater than the bulk of the original culture, from which only about one in  $10^5$  cells survived. The surviving cells retained some of the features of neoplastic cells but did not produce tumors when inoculated into mice. The cell lines eventually became aneuploid (mainly subtetraploid), resembling in chromosome distribution the preneoplastic cells described in cultures of mouse salivary gland epithelium. Because focus formation is maximized by fluorescent light irradiation plus antitubulin treatment, the authors conclude that DNA damage, as well as chromosome reshuffling, is a significant factor in the neoplastic transformation (8).

## 5. Spontaneous Neoplastic Transformation in Newly Explanted Cells

When connective tissue is removed from the mouse embryo and the cells separated from one another in culture, they will multiply in a nutrient medium containing animal serum. They can be subcultured at 3-day intervals, but the growth rate decreases at each subculture until they reach a point of crisis. At that point, a new cell type appears that can grow indefinitely in culture. If the subcultures are kept at low population densities, the permanent cell line stops growing at confluence, behaving like normal cells. If the subcultures are made at high density, the ensuing cell line gains the capacity to multiply at high density, a characteristic of neoplastic and preneoplastic cells. Such cells produce tumors when inoculated into the same strain of mice from which the cells originated. This transforming effect of high population density also occurs in the normally behaving established lines of mouse cells. For example, if the cells are grown to confluence before subculture, transformed foci appear on a background of normal cells. These transformed cells can be maintained in culture indefinitely, and they will eventually gain the capacity to produce tumors in mice (9). This process was termed *spontaneous malignant* or *spontaneous neoplastic transformation* because it occurred in the absence of deliberate carcinogenic treatment. There was a trend toward a marked increase in

chromosome abnormalities at about the time of neoplastic transformation. The abnormalities were chromosomal breaks and minute and banded chromosomes. The frequency of transformation depended on the type of serum used in the medium, and it could be increased by exposure to fluorescent light or high concentrations of oxygen.

Chinese hamster embryo fibroblasts also undergo spontaneous neoplastic transformation in culture. A sequence of changes in growth behavior preceding transformation occurs in these cells, beginning with a crisis followed by establishment of permanent growth. The early preneoplastic changes consist of growth to higher population densities in crowded cultures and increased capacity to multiply in low concentrations of protein [growth factors](#). The capacity for isolated cells to multiply and form colonies increases, the cells gain the capacity to grow in suspension, and colonies with transformed morphology appear. After many more subcultures, the cells are able to produce tumors when inoculated into gelatin sponges that had been implanted into the subcutaneous space of the hamsters. Still later, the cells gain the capacity to produce tumors when inoculated in suspension directly into the subcutaneous space. Chromosome studies were conducted during sequential subcultures. After only six passages, about half the cells had an abnormal chromosome composition, with almost every abnormal cell having a unique karyotype (10). By the 19th passage, all the karyotypes were abnormal; about half of them were of the same type, and the others were quite variable. The presence of individual marker chromosomes increased with passage level. Despite the presence of diverse chromosome abnormalities in the early passages, the cells did not produce tumors in gelfoam sponges until about the 30th passage, nor in direct subcutaneous injection until the 54th passage. It is apparent, therefore, that only a small fraction of the total number of abnormal chromosomal combinations are associated with tumor-producing capacity, but many more combinations are associated with preneoplastic growth behavior of the cells in culture. No specific chromosome combination could be definitively linked to tumor formation, although trisomy of chromosome 5 seemed to play an important role. The sequence of changes in culture leading to tumor-forming capacity can be considered a form of progression similar to that which occurs during the development of tumors in animals.

## 6. Spontaneous Transformation in Established Cell Lines

A problem in using newly explanted or primary cultures of cells in studying spontaneous neoplastic transformation is that the growth properties of the cells are in a constant state of decline. There is an increasing rate of death with each passage, and the growth rate of the survivors steadily declines. Cells of rodent species go through a crisis period of minimal cell growth before a new cell type emerges that can grow indefinitely in culture, but normal human and chicken cells eventually expire. It is possible to study spontaneous transformation systematically in long-established cell lines of mouse embryo fibroblasts, such as the NIH 3T3 line. These can be kept in a nontransformed state, capable of forming a flat monolayer of regularly arranged cells at confluence, if they are subcultured every few days at low population density. If they are kept under the growth constraint of confluence for 2 weeks or longer, foci of transformed cells begin to appear. These foci are just like those produced in C3H10T1/2 cells after exposure to carcinogenic PAH, consisting of closely packed, multilayered criss-crossing cells. The number and size of the transformed foci depends on the length of the incubation period and on the concentration of animal serum in the medium. If a low concentration of serum is used, the confluent cells have to be subcultured and grown again to confluence for foci to appear. Each cell from a transformed focus can initiate a new focus if subcultured on a background of nontransformed cells. The continued capacity for nontransformed cells to undergo transformation depends on the type of serum and the frequency of subculture. After many low-density passages in medium containing a high concentration of calf serum, followed by growth to confluence in a low concentration of serum, the foci produced at first are small and light-staining, but they become larger and darker staining on further passages. Cells from the larger foci produce sarcomas within a few weeks when inoculated subcutaneously into immune-deficient mice. There are many types of morphologically distinct foci in a spontaneously transformed, established cell line, just as there are in primary cultures of baby mouse skin, indicating that there are many pathways to transformation in both cases. Unlike the foci from the primary cultures, however, those

from the NIH 3T3 cells are tumorigenic. There is also evidence that incipient, or unapparent, changes are occurring in the cells before they produce morphologically distinct foci (11). The progressive development of the foci has parallels to the progressively malignant stages of naturally occurring cancer in animals and humans.

The characteristic of cells in culture that correlates most closely with tumorigenesis in animals is known as *anchorage independence*, and it signifies the capacity to multiply while suspended in a soft agar or methylcellulose. Such growth in suspension demonstrates an escape from the requirement of normal cells to attach to and spread on a solid substratum in order to multiply. Oddly enough, repeated seeding of large numbers of nontransformed cells in soft agar suspension, which inhibits their growth, sometimes results in a transformation that permits a cell to multiply in suspension; this does not happen when the cells are growing at a maximal rate attached to the surface of a culture dish. Each of the transformed clones isolated from the agar is morphologically unique, just as each human cancer is unique. Inoculation of large numbers of cells from an anchorage-independent colony into immune-deficient mice results in a rapidly growing tumor within 2 to 3 weeks, but smaller numbers of inoculated cells may produce a disproportionate delay before a tumor appears. Cells from the early tumors grow as well as the original cell culture line when placed in cell culture, but very few cells from the delayed tumors grow in culture, and most of these grow at a low rate. This indicates that the cells have undergone a selection for ability to grow in the mouse during their long incubation and thereby have reduced their capacity for growth in culture. Furthermore, there is great heterogeneity in capacity to produce colonies in agar suspension among individual cells from the delayed tumors. Such variation indicates that a continuous, progressive evolution is occurring in the delayed tumor for improved capacity for continuous growth in the mouse. The great complexity of these growth dynamics has parallels in the great diversity of genetic lesions in human cancers, as will be discussed in the next section.

An unexpected feature of the work with spontaneous transformation is that the procedure of long-term confluence that evokes it also results in some cell death in the confluent culture and in a heritable reduction in growth rate at low population density of the surviving population that later produces transformed foci. It will be recalled that mouse cells transformed by treatment with carcinogenic PAHs also multiply at a lower rate at low density than the parental nontransformed cells. This appears to be true of most, but not all, clones of transformed cells, whether induced by chemicals or spontaneously. These are the same cells that continue to multiply and form transformed foci at high population density when the surrounding nontransformed cells have ceased net multiplication. The reduction in growth rate at low density of the cells is heritable, indicating that genetic damage has occurred (12). It occurs both in transformed and nontransformed cells after prolonged incubation at confluence, indicating that the genetic damage is widespread in the population, but the extent of growth reduction is highly variable, and only certain damage results in transformation. In another cell system, a wide variety of mutagenic chemicals cause a heritable reduction in growth rate of the mutated cells, and the slowdown is related to chromosomal rearrangements and deletions. Because most of these compounds are also carcinogenic, the implication is that neoplastic transformation is frequently the result of chromosomal changes rather than point mutations.

The chromosomal changes, however, are not necessarily visible at the level of the light microscope. For example, two fibroblast cell lines were isolated from a Chinese hamster embryo. Both lines have a stable diploid mode of chromosomes that are apparently normal in appearance. By the criteria of growth in suspension, colony morphology, and tumorigenicity, one of the lines is transformed and the other is not. On prolonged exposure to the cytotoxic drug methotrexate (see [Aminopterin, Methotrexate, Trimethoprim, and Folic Acid](#)), the transformed line develops resistance three times faster than the nontransformed line. A major part of the resistance to methotrexate is the result of [gene amplification](#) (increase in gene copy number) of the gene that codes for the enzyme blocked by methotrexate. In this treatment, large-scale changes in chromosome composition are seen. The implication is that chromosomal changes occur much more readily in neoplastic than in normal cells, with the further implication that those changes underlie the progressive increases in malignant

behavior (invasiveness, metastasis) that occur in tumors.

The same biochemical effect that is induced by treating cells with methotrexate can be achieved by depriving the cells of the vitamin folic acid, which is the substrate of the enzyme blocked by methotrexate. Either treatment causes single- and double-strand breaks in DNA. Folic acid deficiency also increases the metastatic capacity of cancer cells and causes a heritable decrease in the growth rate of cells. It has been suggested that dietary deficiency in folic acid may contribute to the onset of cancer in humans. It may be recalled that a deficiency of choline or methionine results in cancer in experimental animals. Therefore, a part of the damaging and transforming effect of prolonged incubation of cells at the high population density of confluence could be due to depletion of vitamins or [amino acids](#) from the medium. Heritable reduction in growth rate of cells is also a hallmark of cell aging in animals and could be the result of accumulated chromosomal damage below the resolution of the light microscope. The relation between the genetic damage, age, and neoplastic transformation will be discussed below.

## 7. Summary of Findings in Experimental Cancer and Cell Transformation

Injection of a single RNA tumor virus particle can initiate a tumor in animals or transformation in culture. Molecular analysis has shown that a single transforming gene in the viral RNA is the responsible agent. However, such transforming viruses are of little importance in nature, and the tumors they produce in adult animals usually regress. Cellular counterparts of viral transforming genes do not transform cells unless associated with strong retroviral promoters. Transformation initiated spontaneously or by physical and chemical carcinogens is a much more complex and multistaged affair than viral infection, and in many cases it involves large-scale changes at the chromosomal level. A single initiating treatment of cultured cells with a carcinogenic PAH or X-rays genetically alters most or all of the cells in a population in a manner that increases the probability of transformation in their descendants, but the transformation itself only occurs in a small fraction of descendants of each altered cell. The population-wide alteration induced by carcinogens was at first thought by many to be epigenetic or physiological in nature. However, it has been shown to be heritable, and there is evidence that treatments such as X-irradiation, chemical carcinogens, or prolonged periods of crowding and nutritional deficiency cause genetic damage that apparently differs from cell to cell, but genetically destabilizes most of them. There is a general association of transformation with chromosome changes in these treatments. Some of the carcinogens interact directly with DNA, but many do not, indicating that the genetic changes in the latter case are indirect. Such indirect effects have been called epigenetic or dysgenetic, but the ultimate effect is on DNA structure. Many of the indirectly acting agents are clastogenic or disruptive at the chromosomal level, causing recombination, deletion, or asymmetric distribution of chromosomes at **mitosis**. Such changes are not detected in the *Salmonella* test for mutagenesis, which depends on local base substitution or frameshifts in DNA. They can, however, be detected in animal cell culture. There is great diversity in the morphology of transformation, indicating a very large variety of transforming genetic changes. But there are many more chromosomal changes that do not transform cells, so the transforming changes can be considered a selection from among many possibilities. This indicates that it may not be possible to predict which chromosomal changes will cause transformation. This is all the more true because the sensitivity of cells to transformation varies with the initial state of the cells and the methods and materials used in subculturing the cells. Cells thawed from different frozen vials of the same stock of tumor cells will exhibit different growth properties, especially when those properties are repeatedly measured in serial subcultures over an extended period of time. The unpredictable details of the long-term growth behavior of cells, combined with its general dependence on the conditions of culture, are not dissimilar from the extreme heterogeneity and unpredictability of tumors in animals, and tumors are in stark contrast with the relative constancy of function and behavior of normal tissues over extended periods of time. *In vitro* transformation reveals great sensitivity to environmental conditions, such as cell density, serum type, and frequency of subculture, which therefore reflect an epigenetic or conditional aspect of transformation.

## 8. Neoplastic Transformation in Humans

The concerted application of molecular genetics to human cancer began with the transformation of NIH 3T3 cells in culture by DNA from a long-term cell population that had been cultured from a human bladder cancer and with the demonstration that the transforming gene was a mutant form of the cellular *ras* gene. (See section on **Viral transformation** .) Given the common occurrence of genetic variation in cell culture, especially that of human tumors, there is reason to question the conclusion that the mutation caused, or even existed in, the tumor. Those doubts are reinforced by the finding that the gene has to recombine with a strong mouse promoter in order to transform the NIH 3T3 cell. Nevertheless, the finding set off attempts to detect the altered gene in other human cancers. The mutation was reported in the benign tumors (adenomas) that precede colon cancers. Other genetic changes were found in early stages of colon cancer, and it was concluded that five to seven mutations were required to produce advanced stages of cancer. However, it was then found that more than 20% of the alleles of heterozygous genetic loci were lost in the average sporadic colon cancer, which casts doubt on the estimates of the number of mutations needed to produce the cancer. Indeed, the estimates of the number of altered genes in the average case of sporadic colon cancer is now estimated at greater than 25% of the entire genome (13), which implies many thousands of genetic changes during development of the cancer. It should be kept in mind, however, that *ras* mutations found in mutagen-induced rat mammary tumors actually arose from preexisting *ras* mutants. That raises the question of how many of the mutations reported in tumors arose before the tumor, and what role, if any, they had in development of the tumor. Another finding that raises specific questions about the role of *ras* mutations is that *N-ras* gene function of mice can be knocked out entirely in the earliest stages of [development](#) without affecting either the growth or development to adulthood. As methods for detecting genetic changes improve, the number of such changes in tumors increases. When the technique of comparative genomic hybridization was introduced into studies of human mammary cancer, 21 new chromosomal subregions of gene amplification were added to the five already known. The same methodology applied to prostate cancer detected about 30% of the chromosomal regions with a significant increase or decrease in gene copy number in the average case, and some of the cancers had such alterations in more than 50% of the sites (14). These changes involve large chunks of chromosomes, indicating rearrangements and large deletions or amplifications. Although some particular sites were altered in a high percentage of the prostate cancers, no two cancers had the same distribution of altered sites. This confirms what had been reported earlier, that the karyotype of glial tumors was different in every case and, furthermore, that it rapidly changed when the cells were placed in culture. The very large number of copy number changes in genes of prostate cancer seems surprising in view of evidence that nearly 50% of metastatic prostate cancers have a diploid DNA content. The diploid DNA content in so many advanced prostate cancers shows that equal proportions of the genome can be lost or gained, resulting in an overall balance of genetic material. The large number of genetic alterations in cancers indicate great instability in the tumors, but it also creates problems in determining which of the alterations bear a causal relation to development of the tumor, and how many are required to produce its progressive growth. It also raises questions about the source of such instability. One possibility is that an epigenetic change in methylation might lead to defects in chromosome disjunction and distribution at mitosis. Marked decreases are found in **methylation** among each of four genes tested in benign and malignant colonic neoplasms. Because the hypomethylation is prevalent in the benign tumors, the great majority of which have been reported to have no genetic alterations, it is possible that the altered methylation patterns drive the chromosomal changes. In support of such a possibility is the finding that experimentally induced hypomethylation induces transformation in a diploid line of Chinese hamster cells at a high frequency and that the transformation is in every case associated with DNA hypomethylation and the presence of an extra copy of part of a single chromosome. Hence, there may also be a relation between methylation and fragility of chromosome sites, which determines the likelihood of chromosome breaks in particular tissues.

A percentage of human tumors occur in dominantly inherited patterns in families. About 15% of colorectal cancers fall into this pattern (12). Hereditary nonpolyposis colorectal cancer (HNPCC) is one such familial cancer. It is related to germ-line mutations of [mismatch repair](#) genes. This results

in changes in short repeated sequences of DNA (**microsatellites**) in colorectal cancers. The question arose whether similar changes occur in normal tissues of HNPCC patients who had few tumors. Using special amplification methods, it was found that the normal tissues of these patients had a high incidence of these mutations. This shows that the existence of many mutations in cells of tissues throughout the body is not sufficient to initiate cancer, nor is the presence of the mutations in the tumor proof of their causative role. This moved a prominent geneticist to recommend a more rigorous look at the evidence that mutation induction constitutes a rate-limiting step in carcinogenesis (15). Other factors that contribute to carcinogenesis will now be considered.

## 9. The Role of the Tissue Environment in the Origin of Tumors

Teratocarcinoma is a malignant tumor derived from **germ cells** and consists of many different cell types. It can be induced experimentally in mice by transplanting 6-day-old mouse embryos under the testis capsule of adult male mice. The normal diploid tumors can be grown in culture and will produce malignant growth when a single cell is inoculated under the skin of mice. When the same cells are inoculated into the very early embryo (blastocyst stage) of mice, normal development occurs, and redifferentiated teratocarcinoma cells contribute to normal development of many tissues, including some never seen in the original tumors (16). They also contribute to the germ line, which produces normal progeny mice from **gametes** of the teratocarcinoma line. The authors of this work concluded that the conversion to malignancy in all likelihood did not involve mutational events and was completely reversible to normal function and appearance. While it has not been proven unequivocally that there were no mutations involved in causing the teratocarcinoma, it is apparent that the local environment was the determining factor in the origin and in the reversal of this malignant behavior.

A related situation exists in the newt *Triturus cristatus*. Inoculation of carcinogenic PAHs in nonregenerative regions of the newt induced metastasizing tumors that killed the host. Tumors induced in regenerating regions healed spontaneously and differentiated to normal tissues. The malignant behavior of these tumors, therefore, is dependent on the lack of regenerative potential of the surrounding tissue.

There is evidence that the local environment of tissues changes with age. For example, there is a 50-fold increase in the activation of the “inactive” **X-chromosome** in the liver of aging female mice, and there is a sixfold decrease in the transcription of a certain globulin gene in liver cells. The expression of genes is related to their state of methylation, which may serve as a surveillance mechanism for chromosome loss. The **5-methylcytosine** content of tissues decreases with age in mammals, including humans, and the rate of loss of these methylated nucleotides correlates inversely with age in two rodent species that differ in lifespan by a factor of two. This suggests a tie-in between aging, tissue organization, methylation, and chromosomal alterations underlying cancer. The incidence of cancer in humans increases sharply with age (17), and the susceptibility of tissues to carcinogenesis also increases with age. If lung tissue of mice is damaged by X-irradiation or certain cytotoxic chemicals, there is a large increase in the number of metastases that occur in the lung when metastatic cells are inoculated intravenously. Inoculation of rat liver cancer cells into the liver of old rats is much more likely to produce a progressively growing tumor there than if the cells are inoculated into the liver of young mice. In contrast, the liver cancer cells inoculated into the liver of young mice are much more likely to differentiate into normal liver cells. It is apparent that the local environment of the liver in old mice is more favorable for tumor growth than the liver of young mice.

The tissue environment in which human cancers develop appears to be different from that of normal tissues. In cancer of the bladder, there is a gradient of biochemical and cytological abnormality extending for some distance from the edge of the tumor. Loss of **alleles** of some genes occurs in morphologically normal tissue adjacent to breast cancers. Cancer of the esophagus in patients with predisposing conditions is preceded by chromosome changes in large areas of the esophagus in which cancer later arises. The transitional mucosa immediately adjacent to colorectal cancers



contains a variety of biological abnormalities but no evidence of genetic change. There is a 15-fold increase in transcripts of the [methyltransferase](#) enzyme catalyzing DNA methylation in normal colonic mucosa from patients with either polyps or cancer of the colon. It is generally believed that tumors progress by a sequence of genetic changes. As each new mutation appears that selectively favors the mutated cell, the microenvironment of altered cells in which the newly mutated cell is growing may be more favorable for further tumor development than normal tissue. These are all factors that must be taken into consideration in understanding the growth of tumors, and particularly their high degree of genetic instability. A comparable instability is seen when normal rodent cells are dissociated from one another and grown in monolayer culture. Human cells under the same conditions remain largely diploid but lose reiterated sequences of DNA, including **telomeric** DNA. Chromosomal instability persists after the cells have established themselves as a cell line, as indicated by the report that no two cells of a line of rat hepatoma cells were found to be karyotypically identical. Hence, it is possible that the high degree of genetic instability found in human cancer is due not only to the intrinsic instability of the cells, but to the loss of normal organization within and adjacent to the tumor.

## 10. Overview

There are many ways of producing neoplastic transformation: with high efficiency by RNA tumor viruses, with low efficiency by carcinogens, and by dependence on the environment, as in prolonged crowding in certain cell lines in culture. The resultant transformed cells are usually genetically altered, as indicated by the irreversibility of the altered state and the demonstration of altered gene and chromosome composition of the cells (except for exceptional cases like teratocarcinoma, which exhibits in its very origin and reversibility the importance of the surroundings in which cells grow). The idea that transformation could result from mutation in one or a few genes arose first from transformation by RNA tumor viruses and seemed to receive confirmation by transforming especially sensitive NIH 3T3 cells with an altered *ras* gene isolated from an established cell line that originated from a human bladder cancer. Attempts to generate the same kind of transformation with other genes from tumors were unsuccessful. It was then shown that the transformation by the *ras* gene depended on recombination with strong promoter elements, with no indication that it had such a relation in the original tumor. Furthermore, the NIH 3T3 target cells were shown to transform by themselves without transfection, by long-term incubation under growth constraining conditions. There is also the likelihood that the *ras* mutation in chemically induced tumors of rats occurred during normal development.

Despite these reservations, the idea grew that cellular genes with nucleotide sequences like those of RNA tumor viruses were potential “oncogenes” if activated (mutated). It was then found that some genetic loci that were heterozygous in normal tissues of an individual with cancer were homozygous in the tumor. These “loss of heterozygosity” (LOH) genes were then considered to be tumor suppressors in the heterozygous state in normal tissue, although there was no direct evidence for such a role aside from the clonal occurrence of the homozygous state in tumors. There were mathematical arguments that the development of solid cancers required five or more steps; and mutational analysis, particularly of colorectal cancers, was taken as support of the argument. As methods detecting genetic changes were refined, however, more and more of them were found in cancers. Chromosome counting evolved to banding of chromosomes, methods were developed for **allelotyping** DNA, and comparative genomic **hybridization** could detect changes in gene copy number in any region of the genome. These methods then detected changes in some common cancers at an average of 25% to 30% of the chromosomal sites per case, and in more than 50% of the sites in a few cases. This indicated a great instability in the cancers, with alterations in thousands of genes. It is well known from classic genetics that the expression of any multigenic phenomenon is very dependent on the genotypic milieu, so that a given mutation may be deleterious in one genetic milieu and advantageous in another. Thus, the combination of mutagenic changes in genotypic milieus that are different in every human, plus the sensitivity of multigenic phenotypes to the surrounding environment, account for the difficulty in predicting the likelihood of nonfamilial or sporadic cancers or their outcome once they appear. Such a high degree of complexity and the problems of

establishing causal chains in organisms were anticipated in the theoretical work of Walter Elsasser (18). Even where there is a dominant germ-line mutation that favors development of cancer with a probability approaching unity, the time of onset cannot be predicted, and only a very small fraction of the cells, all of which carry the mutation, become transformed; to do so, additional mutations are required, but they can be found in normal tissue as well. To achieve a better understanding of cancer, it will be necessary to take into account the genome of the transformed cell, the state of the surrounding tissue, the age of the organism, its diet, and the environment in which it lives.

The reader should also consult the entries on **Antioncogenes**, **Cell cycle**, **Contact inhibition**, **Oncogenes**, **Protooncogenes**, **ras genes**, **Rous sarcoma virus**, **Somatic mutation**, **Tumor necrosis factor**, **Tumor promoters**, and [Tumor suppressor genes](#) for further information on neoplastic transformation.

## 11. Acknowledgments

I wish to thank Mrs. Dorothy M. Rubin for manuscript preparation. The research was supported by the Council for Tobacco Research and the Elsasser Family Fund.

## Bibliography

1. O. H. Iverson (1995) *Crit. Rev. Oncog.* **6**, 357–405.
2. E. Farber and H. Rubin (1991) *Cancer Res.* **51**, 2751–2761.
3. J. Ashby and R. W. Tennant (1991) *Mutat. Res.* **257**, 229–306.
4. J. Cairns (1981) *Nature* **289**, 353–357.
5. C. Clive, K. O. Johnson, J. F. S. Spector, A. G. Batson, and M. M. M. Brown (1979) *Mutat. Res.* **59**, 61–108.
6. J. W. Simons (1995) *Crit. Rev. Oncog.* **6**, 261–273.
7. A. K. Chakraborty, K. Cichutek, and P. Duesberg (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2217–2221.
8. P. D. Cooper, S. A. Marshall, and G. R. Masinello (1982) *J. Cell. Phys.* **113**, 344–349.
9. K. K. Sanford and V. J. Evans (1982) *J. Natl. Cancer Inst.* **68**, 895–913.
10. L. S. Cram, M. F. Bartholdi, A. F. Ray, G. L. Travis, and P. M. Kraemer (1983) *Cancer Res.* **43**, 4828–4837.
11. H. Rubin (1994) *Proc. Natl. Acad. Sci. USA* **91**, 12076–12080.
12. M. Chow and H. Rubin (1996) *Proc. Natl. Acad. Sci. USA* **93**, 9793–9798.
13. K. W. Kinzler and B. Vogelstein (1996) *Cell* **87**, 159–170.
14. M. L. Cher, G. S. Bova, D. H. Moore, E. J. Small, P. R. Carroll, S. S. Pin, J. I. Epstein, W. B. Isaacs, and R. H. Jensen (1996) *Cancer Res.* **56**, 3091–3102.
15. B. A. Bridges (1995) *Science* **269**, 909.
16. B. Mintz and K. Illmensee (1975) *Proc. Natl. Acad. Sci. USA* **72**, 3585–3589.
17. D. Dix and P. Cohen (1980) *J. Theor. Biol.* **83**, 163–173.
18. W. M. Elsasser (1998) *Reflections on a Theory of Organisms. Holisim in Biology*, Johns Hopkins University Press, Baltimore, MD.

## Suggestions for Further Reading

19. L. Foulds (1969) *Neoplastic Development*, Vol. I, Academic Press, New York.
20. J. Cairns (1981) The origin of human cancers. *Nature* **289**, 353–357.
21. E. Farber and H. Rubin (1991) Cellular adaptation in the origin and development of cancer. *Cancer Res.* **51**, 2751–2761.
22. K. W. Kinzler and B. Vogelstein (1996) Lessons from hereditary colorectal cancer. *Cell* **87**,

159–170.

23. H. Rubin (1997) Cell aging *in vivo* and *in vitro*. *Mech. Ageing Dev.* **98**, 1–35.
24. J. B. Little, H. Nagasawa, T. Phennig, and H. Vetrus (1997) Radiation-induced genomic instability: delayed mutagenic and cytogenetic effects of x-rays and alpha particles. *Radiation Res.* **148**, 299–307.
25. G. H. Heppner and F. R. Miller (1998) The cellular basis for tumor progression. *Intl. Rev. Cytology* **177**, 1–56.

## Nerve Growth Factor and Related Neurotrophins

The neurotrophins are a family of dimeric [proteins](#) that control key neuronal behaviors, including growth, differentiation, function, and survival. These proteins are structurally similar molecules that are expressed in greatest abundance within the nervous system and various target sites of innervation. Nerve growth factor (NGF) was the first member of the neurotrophin family, having been discovered over 40 years ago by Rita Levi-Montalcini, to whom (along with Stanley Cohen) the Nobel Prize in Medicine was awarded in 1986. Studies of this protein have been at the frontier of many life sciences, particularly protein chemistry, molecular biology, structural biology, and the neurosciences.

The neurotrophins interact with membrane-spanning surface receptors of responsive cells, thereby activating [signal transduction](#) cascades that trigger distinct biological responses. The dysfunction of neurotrophic mechanisms has been implicated in the pathogenesis of a number of neurological disorders, thereby creating considerable interest in targeting them for therapeutic development. A tremendous effort is currently underway to define novel biological actions of neurotrophins and the mechanisms of their receptor interactions and signal transduction pathways, as well as to develop novel agents mimicking their biological activities.

### 1. The Discovery of NGF

The concept of a “trophic” substance guiding the growth of axons was formulated by Ramón y Cajal, at the beginning of this century. This proposal was developed to explain an exciting observation made when peripheral nerve grafts were implanted in the central nervous system (CNS). Cajal and colleagues noted that peripheral nerve grafts were able to support and direct the growth of damaged nerve fibers of the central nervous system. In the late 1940s, Rita Levi-Montalcini and collaborator Viktor Hamburger began to study the development of the nervous system, including the mechanisms whereby the correct number of peripheral neurons innervate the correct peripheral target. Using the chick embryo as a model system, these investigators demonstrated that development of the nervous system was accompanied by a steady increase in the number of neurons; as they reached their target, however, the number decreased to the final population. This discovery provided the basis for the “neurotrophic hypothesis” of neuronal development. This model states that an excess number of neurons innervate the target tissue during development. As the target produces a limited amount of neurotrophic substance, the excess neurons are deprived of trophic support and die off, presumably as a result of competitive mechanisms. It is now recognized that the excess neurons die via [apoptosis](#) (or [programmed cell death](#)), leaving the optimal number of intact connections.

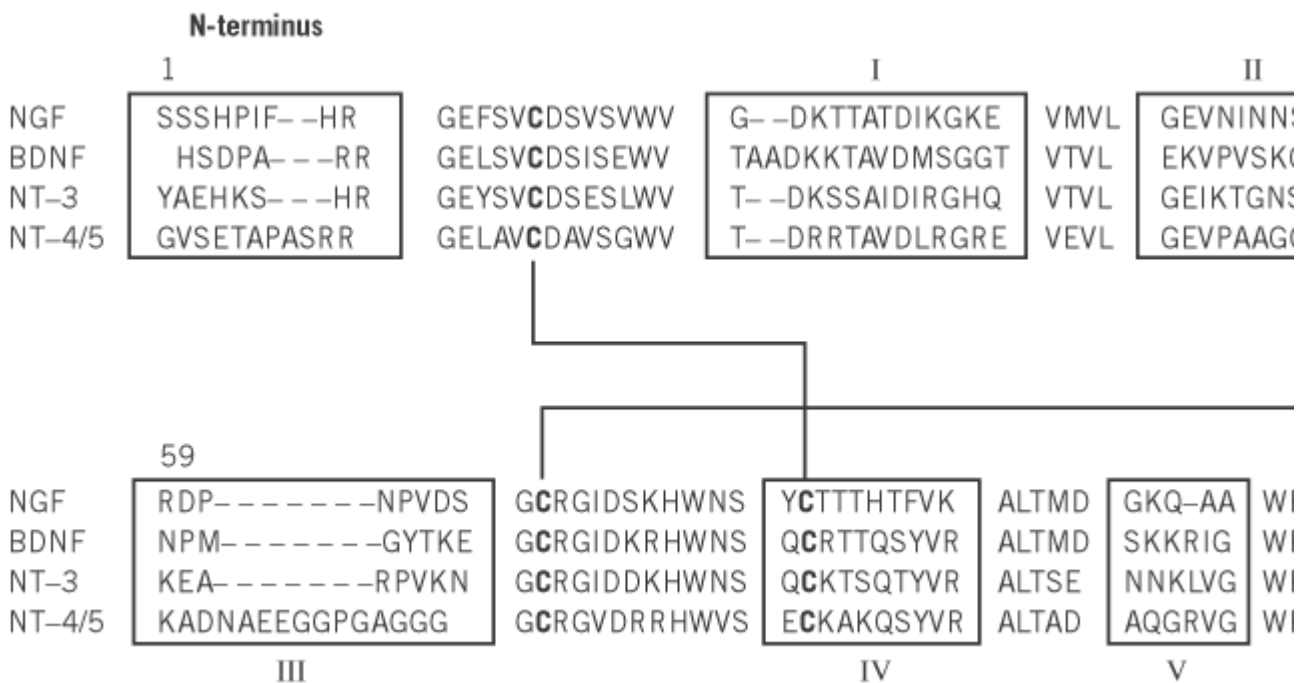
The biochemical characterization of this factor began in the early 1950s with the discovery that a sarcoma tumor [cell line](#) implanted into embryos promoted the growth of ganglia neurons, suggesting that a chemical substance was mediating the effect. Levi-Montalcini and Cohen identified a

nucleoprotein fraction obtained from the sarcoma tumor that promoted neurite outgrowth. A “nerve growth promoting factor” was purified by the late 1950s, which has subsequently become known as NGF. The biochemical characterization of this protein was greatly facilitated by the observation that the male mouse salivary gland was a very rich source of the protein. The subsequent 30 years of investigation provided the purification and amino acid sequencing of this and similar proteins, along with the [cloning](#) of the **genes** for the human neurotrophins. More recent years have witnessed the identification of additional neurotrophin family members, their receptor proteins, and a host of novel biological activities of the neurotrophins.

## 2. Structures of Neurotrophins

The mammalian family of neurotrophins is now comprised of four members. In addition to NGF, brain-derived neurotrophic factor (BDNF) was discovered using classical biochemical purification of the protein from pig brain. More recently, neurotrophin-3 (NT-3) and neurotrophin-4/5 (NT-4/5) have been identified using molecular biology techniques and [DNA sequencing](#). All four members of this family share striking sequence [homology](#) (Fig. 1). The illustrated pattern of [disulfide bonds](#) is identical in all neurotrophins. It is apparent from the alignment of the [primary structures](#) of the neurotrophins that the conserved regions of all family members are located in the same regions, leaving a series of variable regions at the amino and carboxyl termini and interspersed within the primary structure at regular intervals. It was hypothesized early in the characterization of these proteins that the conserved domains maintain the structural features of the neurotrophins, whereas the variable regions mediate the interactions with the specific receptors.

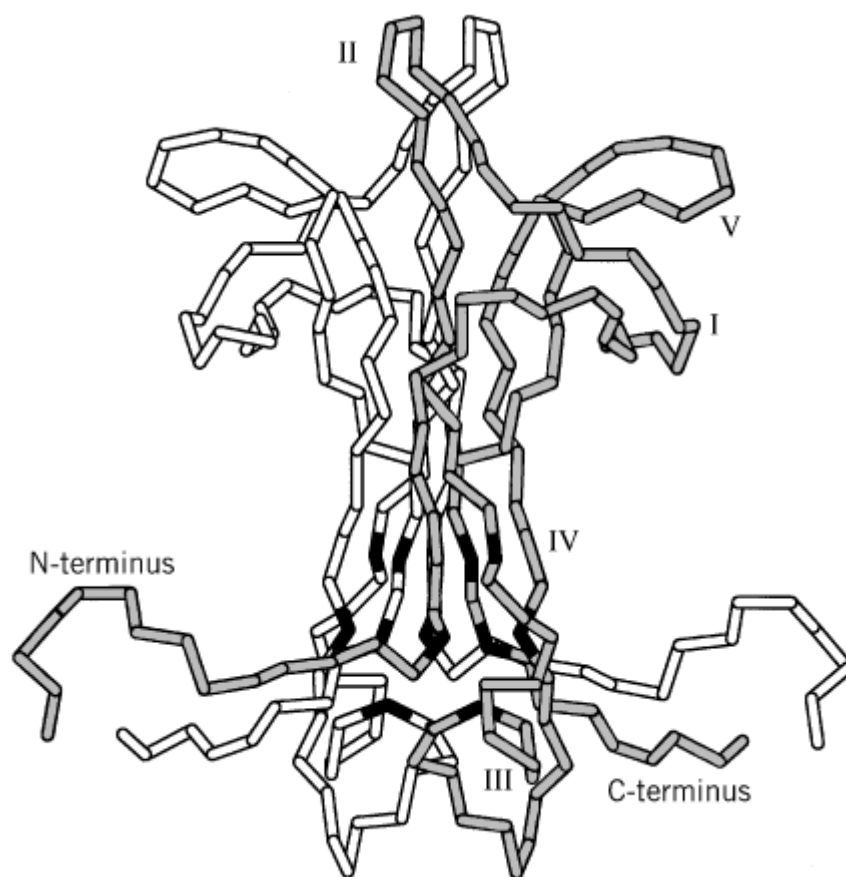
**Figure 1.** The amino acid sequences of the human neurotrophins NGF, BDNF, NT-3, and NT-4/5. The regions with the boxes; their general location within the neurotrophins is depicted in Figure 2. Cysteine residues are shown in bold, with bonds (which form a cystine knot) indicated by the single lines.



The structure of NGF was solved by [X-ray crystallography](#) in the early 1990s (1). The structure of this protein (Fig. 2) turned out to be remarkable, most notably with respect to a unique folding pattern and a novel [cystine knot](#) motif. This structure consists of three antiparallel pairs of b-strands that form a flat surface, with two polypeptide chains of NGF associating through this surface. The

three disulfide bonds that form the cystine knot link these strands and help to stabilize the unique tertiary conformation of NGF. Most of the variable regions (namely I, II, III, and V) are located within the exposed loops, which probably represent the parts that dock with receptor (Fig. 2). The crystal structure of other members of the neurotrophin family have now been determined, and the general conformation depicted for NGF appears to be common for all the neurotrophins.

**Figure 2.** The conformation of NGF as determined by X-ray crystallography (1) and molecular modeling (for the flexible amino and carboxyl termini only, ref. 11). This  $\alpha$ -carbon trace demonstrates the way two identical polypeptide chains of NGF (indicated in white and gray) associate to form a dimeric protein. Cysteine residues that define the cystine knots are illustrated in black. The variable regions of the neurotrophins are indicated by Roman numerals, as in Figure 1, for one protomer; they represent areas that are accessible, possibly for receptor binding. The amino- and carboxyl-terminal regions of NGF are involved in binding to TrkA receptor, whereas loop regions I and V are recognized by  $p75^{\text{NTR}}$ .



### 3. Neurotrophin Receptors

Each of the mammalian neurotrophins binds to two distinct receptor types. All four neurotrophins bind to the common neurotrophin receptor designated here as  $p75^{\text{NTR}}$  (which is also known in the literature as  $p75$ ,  $p75^{\text{NGFR}}$ , and  $p75^{\text{LNGFR}}$ ). The **dissociation constants** of the neurotrophins for  $p75^{\text{NTR}}$  are similar (in the nanomolar concentration range), while there are subtle differences in the rates of binding. Conversely, each of the neurotrophins interacts rather specifically with a member of the Trk **tyrosine kinase** receptor family: NGF binds TrkA, BDNF and NT-4/5 bind TrkB, and NT-3 binds preferentially to TrkC. Binding of the neurotrophins to the Trk receptors is generally of higher affinity (occurring at picomolar concentrations) than that to  $p75^{\text{NTR}}$ , so the receptors are often referred to as the high-affinity (Trk) and low-affinity ( $p75^{\text{NTR}}$ ) binding sites. Over the past decade,

significant efforts have been focused on (a) the activation of signal transduction pathways, (b) biological outcomes, and (c) the comodulation of this two-receptor system (2, 3).

#### 4. The Common Neurotrophin Receptor p75<sup>NTR</sup>

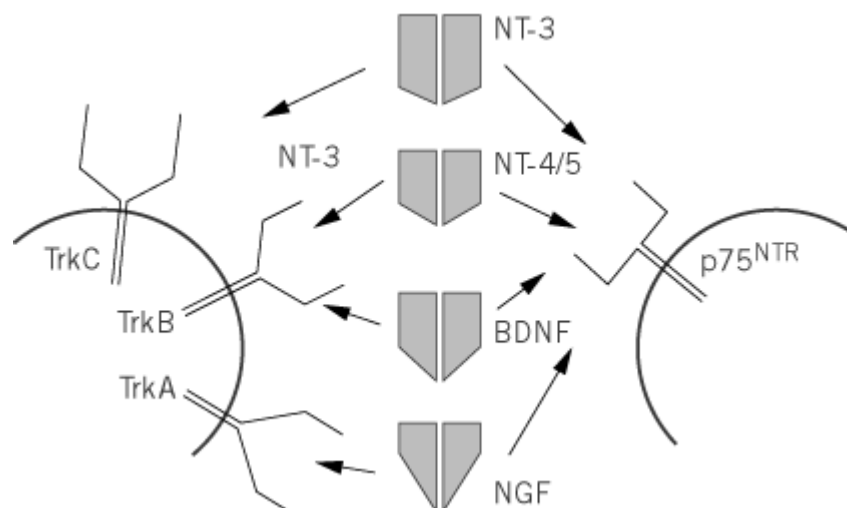
The first receptor identified for neurotrophins, p75<sup>NTR</sup>, was recognized in the early 1970s as a cell-surface receptor for NGF. The common neurotrophin receptor plays a variety of roles in neurotrophin signaling, including the enhancement of Trk activity as well as Trk-independent signaling. Signal transduction cascades induced by p75<sup>NTR</sup> include the production of ceramide as a result of activation of sphingomyelinase, activation of JNK (c-JUN amino-terminal protein kinase), and the translocation of nuclear factor kB (NFkB). Interestingly, while all of the neurotrophins bind to p75<sup>NTR</sup> with similar affinities, the resulting signaling cascades are neurotrophin-specific. Binding of NGF to the common neurotrophin receptor enhances Trk activity, while pathways that are independent of Trk activation may, under appropriate conditions, lead to an apoptotic cell death (4).

The common neurotrophin receptor is one member of what has come to be known as the “death domain” family of receptor proteins. This family (which includes Fas receptor, [tumor necrosis factor](#) receptor, and others) share specific sequence homologies, including an intracellular “death domain” sequence. In addition to the conserved intracellular domains, p75<sup>NTR</sup> and the tumor necrosis factor receptor p55<sup>TNFR</sup> also share conserved extracellular features, particularly the cysteine-rich domains (5).

#### 5. The Trk Family of Neurotrophin Receptors

A major breakthrough in the study of neurotrophins occurred in the early 1990s with the discovery that the *trk* proto-**oncogene** encoded a 140-kDa protein (p140<sup>protoTrk</sup>) that was shown to be a receptor for NGF (6). This protein was initially called Trk, but later became known as TrkA as more neurotrophin receptors of the Trk family were identified. Each of the Trk receptors is relatively specific for each neurotrophin (Fig. 3). It is now well established that the Trk receptors mediate most of the biological actions of the neurotrophins, particularly with regard to cell survival, differentiation, and growth.

**Figure 3.** Neurotrophin receptor selectivity. All neurotrophins bind to p75<sup>NTR</sup> with similar affinities. Selective receptor interactions are observed with the Trk tyrosine kinase receptors: NGF binds TrkA, BDNF and NT-4/5 bind selectively to TrkB, and NT-3 binds preferentially to TrkC. In general, the Trk receptors possess higher affinity for neurotrophins than does p75<sup>NTR</sup>.

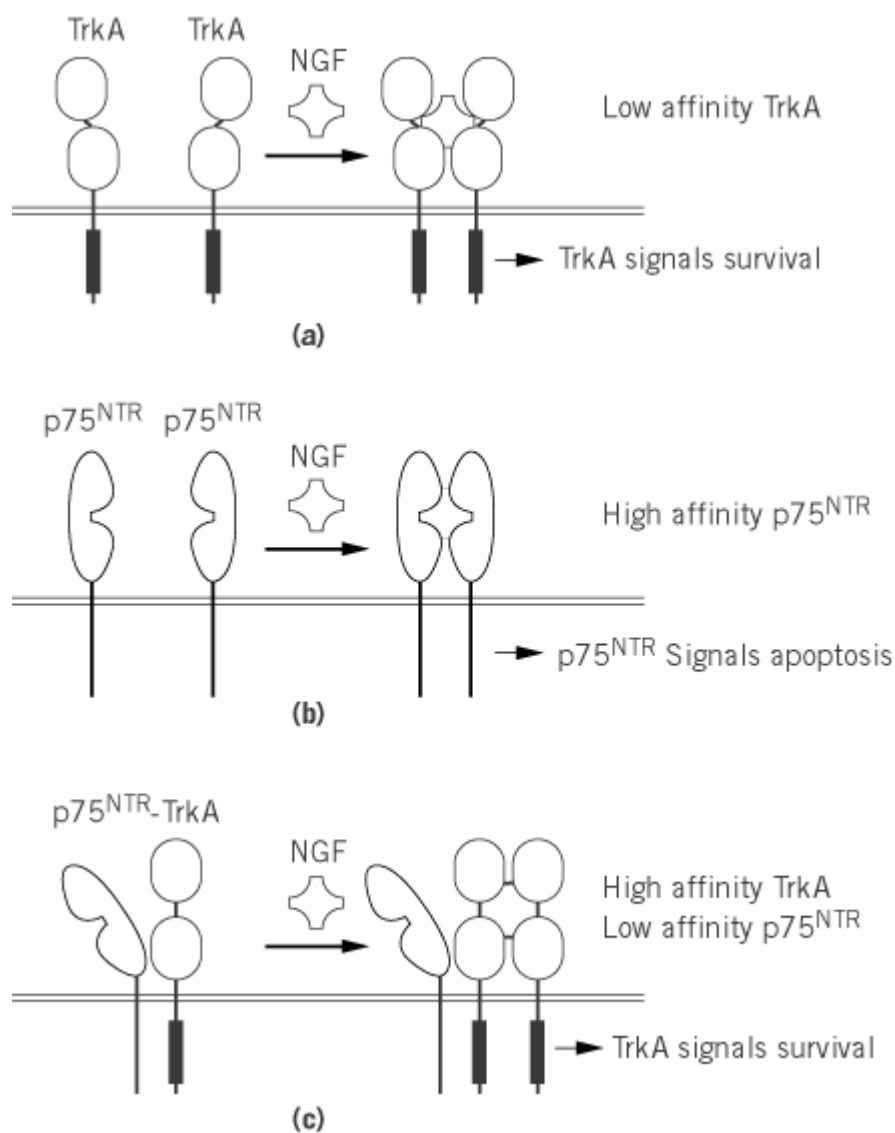


Activation of TrkA receptors is initiated by a ligand-induced dimerization that leads to activation of the receptor's tyrosine kinase activity and autophosphorylation of [tyrosine](#) residues within the intracellular domain of the receptor. Once phosphorylated, the Trk receptors interact with a number of cytoplasmic proteins. These interactions can activate multiple intracellular signaling pathways, including **phospholipase C-1**, PI3-kinase, **Ras**, MEK, and others. The signaling mechanisms of TrkB and TrkC have not been sufficiently investigated, but these receptors are also known to dimerize and to autophosphorylate upon ligand binding, and they probably share many of the transduction pathways with TrkA ([7](#)).

## 6. Interactions Between p75<sup>NTR</sup> and TrkA

Considerable experimental data accumulated over the past decade suggest that p75<sup>NTR</sup> modulates the activity of TrkA. As such, it has long been proposed that there is probably a direct interaction between these two receptors, although such an interaction has been difficult to demonstrate directly. The importance of demonstrating such an interaction has increased recently with the recognition that p75<sup>NTR</sup> can induce cell death, but only in the absence of TrkA. These observations led to the conclusion that the receptor modulation is reciprocal, and they provide further evidence that a direct interaction exists. Several recent reports provide direct evidence for an interaction between these two receptors, and it has been suggested that a transient TrkA-p75<sup>NTR</sup> complex exists in the absence of NGF. It appears that p75<sup>NTR</sup> has the ability to enhance the effects of TrkA, while the TrkA receptor can inhibit the activity of p75<sup>NTR</sup> ([8](#)). In the context of receptor expression on a neuronal cell, this situation could ensure neurotrophin-mediated survival of Trk-expressing cells, while the apoptotic death of cells expressing Trk receptors would be inhibited. The possible receptor expression patterns and examples of cellular outcomes in response to NGF are illustrated in [Fig. 4](#).

**Figure 4.** Potential cellular responses to NGF. The direct interactions of neurotrophin receptors TrkA and p75<sup>NTR</sup> is an emerging concept. Cells that express TrkA only (**a**) will undergo typical responses to NGF, including growth, survival, or differentiation. Cells expressing only p75<sup>NTR</sup> have the potential to undergo apoptosis in response to NGF (**b**). Where both receptors are present, an enhanced TrkA outcome is observed, such that cells respond to NGF more efficiently and by a TrkA-mediated mechanism (**c**). In the presence of TrkA, p75<sup>NTR</sup> does not induce an apoptotic signal. Experimental data indicate that these receptors interact allosterically, such that TrkA induces a lower affinity state in p75<sup>NTR</sup>, while p75<sup>NTR</sup> induces a higher affinity state in TrkA.



## 7. Biological Actions of Neurotrophins

Many different biological outcomes have been demonstrated to result from the treatment of a responsive cell with a neurotrophin. Most were identified originally using cultured cells obtained from the developing nervous system. With respect to nerve cells (neurons), the most studied effects of the neurotrophins are on differentiation, growth, and survival. Neuronal differentiation is a process wherein a precursor cell differentiates into a neuron (perhaps, with a specific neurotransmitter **phenotype**). The growth of neurons may include features such as the length and the number of branches of neurites. These features may be quite distinct depending on, for example, whether the neuron is developing or has been damaged (as in nerve regeneration). With regard to survival, the roles of neurotrophins are turning out to be quite complex. While the neurotrophins have classically been considered to be survival factors, it is now apparent that some cells can under certain conditions undergo NGF-mediated cell death. Apoptotic cell death appears to be mediated solely by p75<sup>NTR</sup>, because NGF-mediated cell death is generally restricted to cells that are p75<sup>NTR</sup>-positive and Trk-negative (4).

The biological actions of neurotrophins are not limited to neurons or to cells of the nervous system. Many non-neuronal nervous system cells (such as Schwann cells, oligodendrocytes, etc.) respond to neurotrophins, which may mediate such biological outcomes as cell migration and proliferation.



Outside of the nervous system, the neurotrophins influence the activity of several cellular members of the immune system, including B lymphocytes and mast cells. NGF causes degranulation and histamine release from mast cells, which can alter neuronal function by indirect mechanisms (9).

## 8. Emerging Concepts

Consistent with the recognition of expanding roles and multiple actions of the neurotrophins, new concepts of their influence on biological systems are being developed. It is now recognized that the neurotrophins can modulate specific behaviors of nerve cells. With respect to pain mechanisms, the neurotrophins can influence the phenotype of neurons, and they may play a key role in the perception of normally innocuous stimuli as painful ones. The neurotrophins have also been shown to alter the functions of neurons in acute conditions, acting to alter **ion channel** function and neurotransmitter actions within fractions of a second. These actions are quite discrete from the slower activities of their modulating functions, such as growth and survival, which generally involve the biosynthesis of new proteins and occur over a much longer time scale.

With the recognition that NGF may trigger immediate cellular responses, in contrast with active cell death versus survival, there is now considerable focus on the multiple receptors for NGF and their interactions. Specifically, the way that direct interactions of these receptor proteins may lead to such dramatically different outcomes is of tremendous interest. It is now becoming apparent that how a cell will respond to NGF is a function of the exact receptor expression pattern (including the receptor ratios) and also the signaling pathways that the cell is able to recruit. Given the complexity of a system with interacting receptors and the possibility of selective recruitment of multiple signaling pathways, it is understandable how so many different cellular outcomes could result from exposure to neurotrophins.

## 9. Therapeutic Potential of Neurotrophin Agonists and Antagonists

A number of neurological diseases are characterized by neuronal death and loss of connectivity within the nervous system. The recognition that neurotrophins represent potent survival factors and modulate neuronal plasticity has generated considerable expectations for the potential use of the agents for therapeutic development. In particular, a significant effort has been mounted to identify Trk receptor agonists that would mimic the effects of neurotrophins. The neurotrophins themselves, while the focus of much activity, will probably have limited therapeutic efficacy due to their inability to penetrate the blood–brain barrier. The identification of small organic compounds that are Trk agonists and likely to gain access to the CNS has been disappointing thus far. Clinical trials on the neurotrophins themselves have provided some encouraging results, particularly for peripheral indications, but none of them are currently approved for clinical use in humans (10). Several new approaches to modulation of neurotrophin biology are underway, including the identification of agents that alter signal transduction pathways and receptor–ligand and receptor–receptor interactions. As new roles for neurotrophin activities are being defined, most notably with regard to pain mechanisms and immune function, the clinical utility of neurotrophin antagonists may prove to be highly significant. Such compounds have been identified, and development of new therapeutic agents that block the deleterious actions of neurotrophins may emerge.

## Bibliography

1. N. McDonald, R. Lapatto, J. Murray-Rust, J. Gunning, A. Wlodawer, and T. L. Blundell (1991) *Nature* **354**, 411–414.
2. M. V. Chao and B. L. Hempstead (1995) *Trends Neurosci.* **18**, 321–326.
3. M. Bothwell (1995) *Annu. Rev. Neurosci.* **18**, 223–253.
4. G. Dechant and Y.-A. Barde (1997) *Curr. Opin. Neurobiol.* **7**, 413–418.
5. A. N. Baldwin and E. M. Shooter (1995) *J. Biol. Chem.* **270**, 4594–4602.
6. D. R. Kaplan, D. Martin-Zanca, and L. F. Parada (1991) *Nature* **350**, 158–160.

7. D. R. Kaplan and F. D. Miller (1997) *Curr. Opin. Cell. Biol.* **9**, 213–221.
8. G. M. Ross, I. L. Shamovsky, G. Lawrance, M. Solc, S. M. Dostaler, D. F. Weaver, and R. J. Riopelle (1998) *Eur. J. Neurosci.* **10**, 890–898.
9. R. Levi-Montalcini, S. D. Skaper, R. Dal Toso, L. Petrelli, and A. Leon (1996) *Trends Neurosci.* **19**, 514–520.
10. F. Hefti (1994) *J. Neurobiol.* **25**, 1418–1435.
11. I. L. Shamovsky, G. M. Ross, R. J. Riopelle, and D. F. Weaver (1996) *J. Am. Chem. Soc.* **118**, 9743–9749.

### **Suggestions for Further Reading**

12. T. Hagg (1997) "Neurotrophic factors". *Encyclopedia of Human Biology*, Vol. **6**, 2nd ed., Academic Press, New York, pp. 223–234.
13. R. M. Lindsay, S. J. Wiegand, C. A. Altar, and P. S. DiStefano (1994) Neurotrophic factors: From molecule to man. *Trends Neurosci.* **17**, 182–190.
14. A. M. Davies (1994) The role of neurotrophins in the developing nervous system. *J. Neurobiol.* **25**, 1334–1348.
15. V. Hamburger (1993) The history of the discovery of the nerve growth factor. *J. Neurobiol.* **24**, 893–897.
16. R. Levi-Montalcini (1987) The nerve growth factor 35 years later. *Science* **237**, 1154–1162.

## **Neural Crest**

The neural crest is a mass of cells that initially forms between the neural tube and the ectoderm (1). The neural crest then forms two long masses of cells along the dorsolateral side of the neural tube. Soon after forming the neural crest, the neural crest cells begin migrating laterally and ventrally. The migration spreads the neural crest cells throughout the body, coming to lie between the mesoderm and the epidermis and between the organ rudiments. During differentiation, the neural crest cells give rise to pigment cells, the visceral skeleton, subcutaneous connective tissue, and some of the peripheral nervous system.

See also [Cell Lineage](#).

### **Bibliography**

1. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 173 and 361–363.

### **Suggestions for Further Reading**

2. N. M. Shah, A. K. Groves, and D. J. Anderson (1996) Alternative neural crest cell fates are instructively promoted by TGFB superfamily members. *Cell* **85**, 331–343.
3. J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.
4. B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, pp. 173 and 361–363.

## Neural Networks and Genetic Algorithms

The information about the behavior of biological molecules can be evaluated and classified by different **computer simulation** approaches and can then be used to predict the properties of related systems. Effective ways of performing such tasks are provided by neural networks (1) and genetic algorithms. Commonly used types of neural nets contain input and output layers connected by a hidden layer. The network is first trained to perform the desired task by being presented with a set of sample input and output. During the training period, the network optimizes the strength of the connections that give the best agreement between the input and output. The trained network can now be used for predictive purposes. Neural networks have been applied to a wide range of problems including predictions of protein **secondary structure**, (2) structure–activity correlation (3), and side-chain conformations (4).

Genetic algorithms can also be used to predict molecular properties. These methods are designed to find optimal solutions to problems using strategies based on biological **evolution** (5). The first step creates a population of possible solutions whose “fitness” is calculated. Finally, a new population is generated with a bias toward a better fit. Genetic algorithms were originally developed as a technique for global optimization. However, this aim has not been accomplished, and the method is used primarily as an effective way of generating a large number of reasonable solutions (6).

### Bibliography

1. J. Zupan and J. Gasteiger (1993) *Neural Networks for Chemists*, VCH, Weinheim, Germany.
2. N. Qian and T. J. Sejnowski (1988) *J. Mol. Biol.* **202**, 865–884.
3. T. A. Andrea and H. Kalayeh (1991) *J. Med. Chem.* **34**, 2824–2836.
4. J.-K. Hwang and W.-F. Liao (1995) *Protein Eng.* **8**, 363–370.
5. D. E. Goldberg (1989) *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, Reading, MA.
6. R. S. Judson, W. P. Jaeger, A. M. Treasurywala, and M. L. Peterson (1993) *J. Comp. Chem.* **14**, 1407–1414.

## Neural Tube

The neural tube is a structure formed in vertebrate **embryos** that gives rise to the brain and spinal chord. After gastrulation, the neural folds arise as paired structures from the **blastopore** to the top of the embryo. The neural folds soon fuse to form the hollow neural tube. The anterior portion of the neural tube gives rise to the brain, the posterior portion to the spinal chord.

### Suggestions for Further Reading

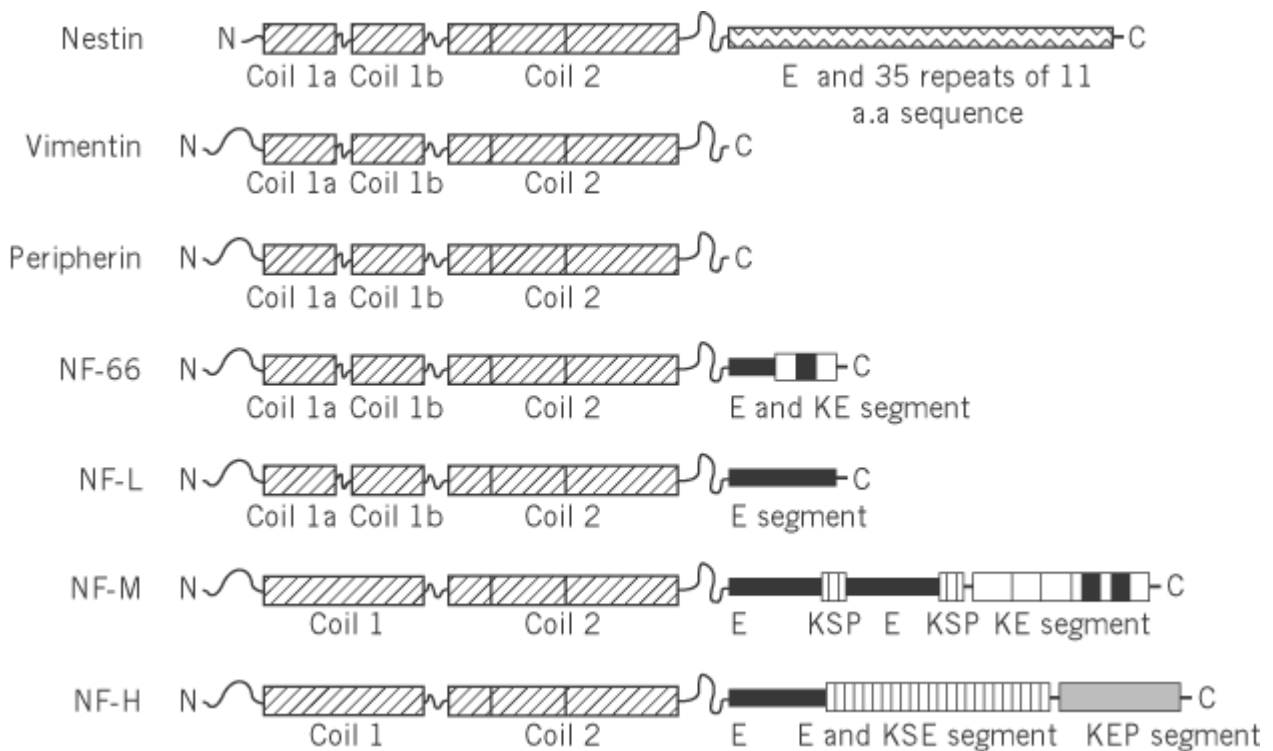
- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 172.
- J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, 2nd ed., Cambridge University Press, Cambridge, UK, 328 pages.

- A. S. Wilkins (1993) *Genetic Analysis of Animal Development*, 2nd ed., Wiley-Liss, New York, p. 546.
- A. Hemmati-Brivanlou and D. Melton (1997) Vertebrate neural induction. *Annu. Rev. Neurosci.* **20**, 43–60.
- J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.

## Neurofilaments

Neurofilaments are the intermediate filaments of mature neurons. They are obligate heteropolymers composed of three different subunits, the neurofilament “triplet proteins,” designated NF-H (“heavy”), NF-M (“medium”), and NF-L (“light”), reflecting their apparent molecular masses of 200, 145 to 160, and 70 kDa, respectively, as determined by SDS-PAGE (Fig. 1). Unique charge characteristics and extensive phosphorylation of NF-M and NF-H cause these subunits to migrate anomalously on SDS-containing gels, and their true molecular masses are smaller (1). Neurofilament triplet proteins are members of one of the five classes of intermediate filament (IF) proteins, which are tissue-specific proteins distinguished by their amino acid sequences and intron placement. Exclusively localized to neurons, the neurofilament subunit triplet proteins have distinctive amino- and carboxyl-terminal domains that characterize them as class IV IF proteins (2).

**Figure 1.** Domains of the individual polypeptide chains of neuronal intermediate filaments. Coil 1 and coil 2 regions are conserved  $\alpha$ -helical domains containing heptad repeats of hydrophobic amino acids that probably adopt coiled-coil conformations in the trimer. These regions are separated by nonhelical “links” (represented by separation of boxes) and, within coil 2, by single-residue breaks in the heptad repeats (represented by vertical lines). Note that NF-M and NF-H possess a unique, single, elongated coil 1 region, which may restrict their ability to form filaments in the absence of a “backbone” subunit, such as NF-L. The initial segment of the carboxy terminus, which constitutes the tail region of vimentin and peripherin, is followed by the following specialized domains in the other neuronal intermediate filament species. The carboxy terminus of nestin contains glutamic acid-rich segments (“E”) and about 35 repeats of an 11 amino acid (a.a.) sequence. Nestin also possesses a relatively short N-terminus. NF-66 and the neurofilament triplet proteins NF-L, NF-M, and NF-H each contain one or more glutamic acid-rich regions (“E segments”). NF-66 also contains a region rich in lysine and glutamic acid residues (“KE segment”). NF-M and NF-H contain repeating regions of the sequence Lys-Ser-Pro (“KSP2 segments”), which represent the sites of extensive phosphorylation of these two species. The extreme termini of NF-M and NF-H possess unique regions rich in lysine and glutamic acid residues (“KE segment”) and in lysine, glutamic acid, and proline (“KEP segment”), respectively. [Adapted from Ref. 7 with permission of the publisher.]



Neurofilaments are the most abundant cytoskeletal components in large myelinated axons. Like other IF proteins, individual neurofilaments are discrete, linear structures that extend for distances of at least 10  $\mu\text{m}$ . They are often curved, indicating a flexibility inherent in their construction. Arranged in parallel in axons, they maintain minimal distances of 25 to 30 nm from neighboring filaments, even in pathological states when neurofilaments accumulate. Their regular spacing and parallel alignment stem, in part, from their most unique feature, lateral sidearms that project perpendicularly from the filament backbone.

### 1. Domain Structure, Assembly, and Developmental Expression of Neurofilament Proteins

As expected for a multigene family evolving from a single ancestral gene (see [Gene Families](#)), all intermediate filaments display a similar domain organization, consisting of a central  $\alpha$ -helical rod domain of highly conserved size and a secondary structure, flanked by strongly basic domains that are highly variable in sequence. (Fig. 1) (2). The rod domain of every IF subunit is composed of four distinct tracks (IA, 1B, 2A, and 2B) that contain characteristic repeating sequences of seven hydrophobic amino acids. This “heptad repeat” strongly predicts that the rod sequence would form a parallel coiled-coil dimer—the first step in the formation of 10-nm filaments. The constant length of these four tracts across all classes of IF proteins suggests that these regions may be critical components in the assembly of the basic filament.

In contrast to the rod domain, the amino-terminal head domains of IF proteins have diverged considerably during evolution. The head domains assume predominantly  $\beta$ -sheet and  $\beta$ -turn conformations and vary in both size and charge. These regions of neurofilament subunits and other nonepithelial intermediate filaments range in molecular weight from 4.5 to 9 kDa. An excess of positively charged amino acids, particularly arginine residues, confers a highly basic charge to the head domain of neurofilament subunits and to the class III intermediate proteins. The head domain, rich in serine and threonine residues, some of which are phosphorylated or glycosylated, is believed to be important in regulating neurofilament assembly ([3-5](#), [93](#)).

The carboxy-terminal tail region of neurofilament subunits is quite distinctive ([6-8](#)) (Fig. 1).

Although the tail domains of class III IF proteins are relatively small (5 to 9 kDa) and similar in structure, neurofilament subunits have large and complex carboxyl termini. By *immunoelectron* microscopy, the NF-H and NF-M tail domains appear as long, thin, fibrous protrusions, or “sidearms,” that extend from the cores of axonal neurofilaments. They are largely random coil (6), and four subdivisions can be defined. Each subunit contains a neutral or slightly basic, proline-containing region 3.5 to 5.5 kDa in size, flanked by a 50 to 60-residue glutamate-rich region that may be the homologue of the entire tail domain of other intermediate filament proteins. In NF-H, this segment is followed by a region that consists of multiple repeats of the sequence–Lys-Ser-Pro- (“KSP repeats”). These KSP sequences, embedded within motifs of 6 or 8 residues, are arranged in nearly tandem repeats and range in number from 39 to 51, depending on the species. A common polymorphism of human NF-H generates either 43 or 44 KSP repeats (9) although the functional significance is unclear. The corresponding region of NF-M is shorter and contains 4 to 13 KSP repeats and a variable number of repeated sequences characterized by the consensus sequences KXSP or KXXSP, where X is any residue (10). Only serine residues located in KSP and KSP-like repeats are phosphorylated (11) and, at least in axonal neurofilaments, most are phosphorylated *in vivo* (8, 12, 13). NF-H terminates at the carboxy end with a variable, glutamic acid-rich region about 20 kDa in size (14, 15). Its counterpart in NF-M is a 24-kDa region rich in lysine and glutamic acid residues, but virtually no proline, which is very highly conserved in sequence among different species. This latter feature is not shared by any other known proteins, suggesting that it may serve a unique neuronal function (10).

Similar to other intermediate filaments, neurofilaments are believed to assemble from dimer, tetramer, and octamer intermediates of two or more subunits (16-21). Each NF triplet protein has an  $\alpha$ -helical rod domain and heptad repeat regions that can initiate the first stage of assembly-coiled coil dimer formation-implying that each subunit may be integral to the neurofilament. Once assembled, neurofilaments usually assume a parallel arrangement within the axonal cytoskeleton. Although single neurofilaments *in vitro* are flexible and tend to curl up (22), under appropriate conditions they can form bundles complete with cross-bridges resembling those seen *in vivo* (23). The carboxy-terminal domain of NF-M is crucial to this interaction (24).

The appearance of NF-L and NF-M in the rodent central nervous system a few days after birth is followed several days later by the appearance of NF-H, a marked up-regulation of NF-L and NF-M expression, and extensive NF-H and NF-M phosphorylation (25-28) –events associated with major axonal caliber expansion. Thus, as neurons develop, the plastic, intermediate networks, composed of the class III intermediate filaments, (e.g. vimentin and  $\alpha$ -internexin) are replaced by more stable class IV IF networks to support the highly polar morphologies of mature neurons, especially those with large caliber axons. The cDNA sequences that encode the entire protein sequences of NF-L, NF-M, and NF-H are available (15, 29), but the regulation of neurofilament gene expression is not well understood.

## 2. Posttranslational Modification

The spectrum of dynamic behaviors of neuronal cytoskeletal proteins, from local subunit assembly to the formation of stable cross-linked polymer networks, is influenced by a host of protein kinases and phosphatases (see [Phosphorylation, Protein](#)). The amino-terminal domains of NF-L and NF-M have multiple phosphorylation sites, at least 5 and 9, respectively, in the mouse (3, 4). Results from consensus sequence analysis and *in vitro* phosphorylation assays are compatible with regulation of these sites by protein kinase A, protein kinase C, and calmodulin-dependent protein kinase II (94). As neurofilaments are transported along axons, phosphate turnover continually changes the locations of phosphate groups on each subunit (4) (30). It is believed that phosphorylation/ dephosphorylation along the amino-terminal domain of NF-L, and possibly NF-M, early after subunit synthesis, modulates polymer assembly and disassembly (30, 31). The amino-terminal head domains of rat NF-L and NF-H are also posttranslationally modified by adding of O-linked *N*-acetyl glucosamine moieties (5) (see [O-Glycosylation](#)).

The NF-H carboxy-terminal domain contains more than 50 mols of phosphate. The same regions of NF-M and NF-L contain at least 10 mols and 3 mols, respectively. A number of known protein kinases, mitogen-activated kinases (ERK1) (32), (ERK2) (33), glycogen synthase kinase-3 (GSK3a, GSK3b), cyclin-dependent kinase cdk5 (21, 32, 34, 35), SAP kinase (36, 37), a neurofilament associated (115 kDa) kinase (38), a casein kinase I-like kinase (39), and several less well characterized kinases are considered candidates for regulating the function of this domain, although no kinase has yet been proven to phosphorylate neurofilaments *in vivo*. Phosphorylation begins soon after the subunits are synthesized but becomes extensive only after the neurofilaments have moved into the axon, implying that this process has a specialized role in axons (40, 41). Phosphorylation of the NF-H and NF-M tail domains induces neurofilaments to straighten, align, bundle, and develop cross-bridges *in vitro* and in cultured cells (22, 24). NF-L head phosphorylation is modulated during neuronal signal transduction (94).

### 3. Neurofilament Dynamics and Turnover

Neurons, especially those with long axons, export most neurofilament proteins into axons within 6 to 12 hours after their synthesis (40). It is believed that newly synthesized subunits are transported mainly as polymers (42, 43). It is possible, however, that some transported subunits can also incorporate into existing neurofilaments to maintain neurofilament structure or to extend filament length (44, 45, 95, 96).

Transported neurofilaments (or neurofilament proteins) advance along axons by slow axonal transport at a rate between 0.25–3 mm per day depending on subunit composition (97), the type of neuron, stage of neuronal development, and location along the nerve (46). Rates reflect the average of many intermittent rapid movements and pauses (43). Transport is more rapid in immature or regenerating axons than in mature axons and slows during aging. During transport, the addition of specific phosphate groups on the carboxy-terminal domains of NF-M and NF-H primarily at the KSP motifs (41, 47, 48) is associated with changes in neurofilament behavior and morphology. As negative charge from the introduced phosphate groups accumulates along neurofilament sidearms, the minimal spacing between neighboring neurofilaments changes from 25 to 30 nm to 50 to 55 nm (41, 48). A family of cross-linking proteins, such as BPAG and plectin (101-104), may contribute to this process. The addition of specific phosphate groups to the NF-M and NF-H tails is also associated with slowing or cessation of the movement of some filaments (47, 49, 98), facilitating their incorporation into a stationary cytoskeletal network within which individual components of the cytoskeleton associate for varying periods before advancing further along axons or being degraded locally (44, 45, 48-50). This stationary network is non-uniform along some axons, particularly at specialized sites, such as nodes of Ranvier (41, 48, 49). In addition to neurofilament protein expression levels, the equilibrium between moving and stationary neurofilaments in large part determines the size of the neurofilament network, which, in turn, influences axon caliber (51, 52). Signals from oligodendroglial cells during the earliest stages of myelination have a remarkable influence on these events (53-55). Association of oligodendroglia with regions of the axon destined to be myelinated induces phosphorylation at a specific carboxyl-terminal site on NF-H and NF-M which is closely associated with the process of neurofilament accumulation and expansion of axonal caliber (52, 55). Phosphorylation at this site may reflect important phosphorylation events at other sites on neurofilaments or other interacting proteins because mice lacking NF-H have nearly normal axon morphology (105-107).

Neurofilaments that do not incorporate along the axon reach the nerve terminals where they are presumably degraded (42, 56), although retrograde transport of neurofilaments or neurofilament fragments has not been entirely excluded (57, 58). In long peripheral axons, the radioactivity of a pulse-chase labeled pool of neurofilaments declines considerably during axonal transport (59, 60), and neurofilament breakdown products are detected along the axons (61), indicating that some neurofilaments, especially those that are stationary for long periods (49), may be turned over locally within axons. The possibility cannot be excluded, however, that neurofilaments, parts of neurofilaments, or subunits are slowly released from the cytoskeleton and then are moved into the

terminals for degradation. Phosphorylation seems to protect neurofilaments against proteolysis (62), suggesting that kinase and phosphatase activities might partly determine the turnover rate of neurofilament proteins.

#### 4. Neurofilament Function and Dysfunction

Intermediate filaments in many cell types provide a structural scaffold that supports cell shape and helps to organize cytoplasmic constituents (see [Intermediate Filaments](#)). Because of their unique sidearm extensions and ability to arrange in parallel arrays, neurofilaments are particularly well suited to fill space and support the extreme polar shape of neurons. Once axons reach their targets and form synapses, their further maturation involves up to a 10-fold increase in axonal diameter and commensurate increases in neurofilament content (63, 64). This growth is essential for achieving appropriate conduction velocity for electrical impulses. The strong correlation between neurofilament number and axon diameter during radial growth, regeneration, and neuropathological axonal atrophy or hypertrophy has suggested that neurofilaments function as a major determinant of axonal caliber (63-66). Several observations now firmly establish their importance to axonal growth. In the Japanese quail, a recessive mutation (quiverer) that interferes with NF-L synthesis (67) severely reduces radial growth, slows axonal conduction velocity, and causes the bird to quiver. Targeted disruption of mouse NF-L selectively eliminates neurofilaments and severely inhibits radial growth (68). Furthermore, expression of a mouse NF-H  $\beta$ -galactosidase fusion protein in transgenic mice selectively blocks neurofilament transport into axons and substantially reduces axonal growth (69). Studies of transgenic mice overexpressing single subunits or combinations indicate that proper subunit ratios and states of phosphorylation are important for normal radial axonal growth (66, 70, 71) but not for axonal outgrowth during development (99). NF gene deletion analysis (68, 72-75, 100) confirms the importance of NF-L and NF-M for filament assembly and caliber expansion and that the role of NF-H in this function can be partially fulfilled by compensatory phosphorylation on NF-M (52).

Given the complexity of neurofilament dynamics and regulation by multiple protein kinases, it is not surprising that neurofilaments are frequent cellular targets in neuropathological disorders, and that neurons with large caliber axons and abundant neurofilaments are particularly vulnerable. In the central nervous system, neurofilaments are a major constituent of inclusions known as Lewy bodies that develop in susceptible neuronal populations in two relatively common forms of *senile* dementia, diffuse Lewy body disease and Lewy body variant of Alzheimer's disease (76). Although  $\alpha$ -synuclein is now considered to be the principal component of Lewy bodies (76), these structures also contain all three neurofilament subunits. Neurofilaments are also major constituents of NFTs in Alzheimer brain, although the presence of neurofilaments in these tangles is overshadowed by the formation of paired helical filaments (PHF) containing tau protein (76). Neurofilament hyperphosphorylation, however, may be the earliest cytoskeletal change in neurites and neuritic plaques in preclinical stages of AD (77).

The most striking disorders involving neurofilaments are those which affect spinal motor neurons, where neurofilaments are most abundant. Abnormal accumulations of neurofilaments in the cell bodies and proximal axons of motor neurons have been established to be hallmarks of both sporadic (79) and familial (80) forms of amyotrophic lateral sclerosis (ALS) and certain less common neuroaxonal dystrophies. ALS is primarily characterized by selective degeneration of upper and lower motor neurons, beginning in mid adult life and invariably progressing to paralysis and death over a 1-5 year time period. Transgenic mouse models lend credence to the idea that these aberrant accumulations of neurofilaments contribute to the motor neuron dysfunction, rather than being simply by-products of the pathogenesis process. When wild-type mouse NF-L (81) or human NF-H are over expressed (71), mice develop selective motor neuron dysfunction associated with marked neurofilament accumulation along proximal parts of axons and atrophy more distally, as well as neurogenic atrophy of skeletal muscle. Despite these changes, large caliber neurofilament-rich axons are not lost during the disease. However, a variety of transgenic mouse constructs that alter subunit phosphorylation, polymer assembly, or neurofilament network formation can cause neurofilament



aggregations leading to neuronal dysfunction (71, 78, 81-84). Overexpression of a mouse NF-L subunit, which has a point mutation (85) within the rod domain produced pathological symptoms in mice resembling human motor neuron disease, including preferential loss of large caliber neurofilament rich axons. A possible etiological role of neurofilament protein mutations in ALS cases is suggested by analyses of the KSP repeat domains of the NF-H gene, which reveal mutations in about 1% of ALS patients. Mutations found in sporadic ALS, and very rarely in familial ALS (FALS), have included small in-frame deletions involving one or more KSP repeats and an 84-bp insertion comprising four extra KSP repeats (86-88). Recently, mutations of NF-L, which involve a single base substitution within an exon encoding a region on the rod domain, have been identified in two families afflicted with Charcot-Marie-Tooth Disease type 2, an inherited motor and sensory neuropathy (89). The relationship of these mutations to disease susceptibility is not known.

Although neurofilament gene mutations are extremely rare in FALS, about 2% of ALS patients have autosomal dominant mutations in the gene encoding the ubiquitous cytoplasmic enzyme Cu/Zn superoxide dismutase (SOD1)(90). More than 50 mutations in SOD1 have been identified in patients with FALS, the majority of which are point mutations. Neurofilament accumulations are also characteristic of the pathology in these cases, and seem to be important to the neurodegenerative mechanism. In mouse models of SOD1 – related FALS, neurologic disease develops even when neurofilaments are eliminated genetically; however, neurofilaments seem to influence the onset and pace of degeneration (91, 92, 108).

In conclusion, neurofilaments continue to figure prominently in mechanisms of human neurodegenerative disease although whether they are more important as agents of disease or as disease targets requires further investigation. Genetic evidence links mutations of neurofilament genes to human neuroaxonal disease and investigations of these mutations in mouse models should be informative. A firm understanding of neurofilament protein regulation should help to clarify the role played by neurofilaments in neurodegenerative disease.

## 5. Neurofilament Dynamics and Turnover

Neurons, especially those with long axons, export most neurofilament proteins into axons within 6 to 12 hours after their synthesis (58). It is believed that newly synthesized subunits are transported mainly as polymers or at least oligomers, based on observations that they enter axons in the same subunit stoichiometry as isolated neurofilaments (48). It is also believed that they behave identically during transport (61), even when either slow transport or the subunit itself is altered in pathological or experimental conditions (62, 63). When NFM is overexpressed by using viral **vectors** in transgenic mice that lack axonal neurofilaments, some subunits move by slow axonal transport (64). Furthermore, tagged NFL subunits **microinjected** into neurons incorporate along the length of existing neurofilaments (65, 66). These observations suggest that transported subunits can also incorporate into existing neurofilaments to maintain neurofilament structure or to extend filament length. Indeed, a large pool of NFH subunits is present during development and decreases during brain maturation (67).

Transported neurofilaments (or neurofilament proteins) advance along axons by the active process of slow axonal transport at a rate between 0.25 and 3 mm per day depending on the type of neuron, stage of development, and location along the nerve (references in 68). Transport is more rapid in immature or regenerating axons than in mature axons and slows during aging. As newly synthesized neurofilament proteins move along the axon, the carboxy-terminal domains of NFM and NFH become extensively phosphorylated, primarily at the KSP motifs (59, 60, 69). The addition of specific phosphate groups is associated with changes in neurofilament behavior and morphology. As negative charge from the introduced phosphate groups accumulates along neurofilament sidearms, they extend peripherally from the filament core and change the minimal spacing between neighboring neurofilaments from 25 to 30 nm to 50 to 55 nm (59, 69). The addition of specific phosphate groups to the NFM and NFH tails also slows or stops movement of some filaments (60, 70), facilitating their incorporation into a stationary but dynamic cytoskeletal network (64-66, 69,

[70](#)). This stationary network is nonuniform along some axons, particularly at specialized sites, such as nodes of Ranvier ([59](#), [69](#), [70](#)). In addition to neurofilament protein expression levels, the equilibrium between moving and stationary neurofilaments in large part determines the size of the neurofilament network, which in turn influences axon caliber. Signals from oligodendroglial cells during the earliest stages of myelination have a remarkable influence on these events ([71-73](#)). Association of oligodendroglia with regions of the axon destined to be myelinated induces phosphorylation at specific carboxyl-terminal sites, triggering incorporation of transported neurofilaments into a stationary network and leading to neurofilament accumulation and expansion of axonal caliber.

Neurofilaments that do not incorporate along the axon reach the nerve terminals where they are presumably degraded ([74](#), [75](#)), although retrograde transport of neurofilaments or of neurofilament fragments has not been entirely excluded ([76](#)). In long peripheral axons, the radioactivity of a **pulse-chase** labeled pool of neurofilaments declines considerably during axonal transport ([77](#), [78](#)), and neurofilament breakdown products are detected along the axons ([79](#)), indicating that some neurofilaments may be turned over locally within axons. Axons contain ample [proteinase](#) activity for this suspected purpose ([79](#), [80](#)). Neurofilaments that remain stationary along axons for long periods turn over at the same rate at proximal and distal levels of the axon ([70](#)). This pattern is consistent with a mechanism of local turnover, although the possibility cannot be excluded that neurofilaments, parts of neurofilaments, or subunits are slowly released from the cytoskeleton and then are moved into the terminals for degradation. Phosphorylation seems to protect neurofilaments against **proteolysis** by divalent cation-independent proteinases associated with cytoskeleton preparations and by calpains (see **Thiol proteinase inhibitors**) ([81](#), [82](#)). These results raise the possibility that kinase or phosphatase activity might partly determine the turnover rate of neurofilament proteins.

## 6. Neurofilament Function and Dysfunction

Intermediate filaments in many cell types provide a structural scaffold that supports cell shape and helps to organize cytoplasmic constituents (see [Intermediate Filaments](#)). Because of their unique sidearm extensions and ability to arrange in parallel arrays, neurofilaments are particularly well suited to fill space and support the extreme polar shape of neurons. Once axons reach their targets and form synapses, their further maturation involves up to a 10-fold increase in axonal diameter and commensurate increases in neurofilament content ([83](#), [84](#)). This growth is essential for achieving appropriate conduction velocity for electrical impulses. A strong correlation between neurofilament number and axon diameter during radial growth, as well as during regeneration and neuropathological axonal atrophy or hypertrophy ([83-86](#)), had suggested earlier that neurofilaments function as a major determinant of axonal caliber. Several observations now firmly establish their importance to axonal growth. In the Japanese quail, a recessive mutation (quiverer) that interferes with NFL synthesis ([87](#)) severely reduces radial growth, slows axonal conduction velocity, and causes the bird to quiver. Furthermore, expression of an NFH- **b-galactosidase** fusion protein in mice selectively blocks neurofilament transport into axons and substantially reduces axonal growth ([88](#)). Although neurofilaments are essential for achieving normal caliber, the factors that specify caliber size and neurofilament number are less clear.

Given the complexity of neurofilament dynamics and regulation by multiple protein kinases, it is not surprising that neurofilaments are commonly implicated in neuropathological states, particularly those involving neurons with large caliber axons and abundant neurofilaments. In the central nervous system, neurofilaments are a major constituent of inclusions known as Lewy bodies that develop in susceptible neuronal populations in two relatively common forms of *senile dementia*, diffuse Lewy body disease and Lewy body variant of Alzheimer's disease ([89](#)). The most striking disorders involving neurofilaments affect spinal motor neurons, where neurofilaments are most abundant. Neurofilaments accumulate markedly in the cell bodies and proximal axons of these neurons in various human motor neuron diseases, including amyotrophic lateral sclerosis (ALS) ([90](#), [91](#)). Transgenic mouse models lend credence to the idea that these aberrant accumulations of neurofilaments contribute to the motor neuron dysfunction, rather than being simply by-products of

the pathogenic process. When wild-type mouse NFL (92) or human NFH are overexpressed (93), mice develop selective motor neuron dysfunction associated with a marked neurofilament accumulation proximally and axonal atrophy distally, as well as neurogenic atrophy of skeletal muscle. Moreover, a growing number of transgenic mouse constructs that disrupt subunit phosphorylation, polymer assembly, or neurofilament network formation can cause neurofilament aggregations leading to neuronal dysfunction (92-97). Also suggesting a possible etiological role of neurofilaments in some individuals who have ALS are analyses of the KSP repeat domains of the NFH gene in DNAs from 356 sporadic ALS cases, which revealed a one-codon deletion in four patients and a 34-codon deletion in one patient, conditions that were not observed in over 300 normal individuals (98). The relationship of these mutations to disease susceptibility is unknown. Although neurofilament gene mutations have not been observed in cases of *familial* ALS (FALS), there are more than 45 mutations in the enzyme *superoxide dismutase* each of which can lead to FALS. Neurofilament accumulations are also characteristic of the pathology in these cases, and oxidative damage to the long-lived neurofilament subunits has been proposed as a mechanism that leads to aberrant filament assembly or organization, which then blocks axonal transport (99). A firm understanding of neurofilament protein dynamics should help to clarify the significance of the role played by neurofilaments in these and other neurodegenerative diseases.

## References

- “Neurofilaments” in , Vol. 3, pp. 1589–1595, by Ralph A. Nixon, Ph.D., M.D, Nathan Kline Institute, Center for Dementia Research, 140 Old Orangeburg Road, Orangeburg, NY, 10962, Tel: 845-398-5423, Fax: 845-398-5422 E-mail: Nixon@nki.rfmh.org; “Neurofilaments” in (online), posting date: January 15, 2002, by Ralph A. Nixon, Ph.D., M.D, Nathan Kline Institute, Center for Dementia Research, 140 Old Orangeburg Road, Orangeburg, NY, 10962, Tel: 845-398-5423, Fax: 845-398-5422 E-mail: Nixon@nki.rfmh.org.
1. J. P. Julien, and W. E. Mushynski, *J Neurochem*, **37**: 1579–1585 (1981).
  2. P. M. Steinert, and D. R. Roop, *Annu Rev Biochem*, **57**: 593–625 (1988).
  3. R. K. Sihag, and R. A. Nixon, *J Biol Chem*, **265**: 4166–4171 (1990).
  4. R. K. Sihag, and R. A. Nixon, *J Biol Chem*, **266**: 18861–18867 (1991).
  5. D. L. Dong, Z. S. Xu, M. R. Chevrier, R. J. Cotter, D. W. Cleveland, and G. W. Hart, *J Biol Chem*, **268**: 16679–16687 (1993).
  6. N. Geisler, E. Kaufmann, S. Fischer, U. Plessmann, and K. Weber, *Embo J*, **2**: 1295–1302 (1983).
  7. G. Shaw (1991) in *The Neuronal Cytoskeleton* (Burgoyne, R., Ed.), pp. 185–214, Wiley-Liss, New York.
  8. N. Geisler, J. Vandekerckhove, and K. Weber, *FEBS Lett*, **221**: 403–407 (1987).
  9. D. A. Figlewicz, G. A. Rouleau, A. Krizus, and J. P. Julien, *Gene*, **132**: 297–300 (1993).
  10. G. Shaw, *Biochem Biophys Res*, **162**: 294–299 (1989).
  11. E. Elhanany, H. Jaffe, W. T. Link, D. M. Sheeley, H. Gainer, and H. C. Pant, *J Neurochem*, **63**: 2324–2335 (1994).
  12. J. P. Julien, and W. E. Mushynski, *J Biol Chem*, **258**: 4019–25 (1983).
  13. S. J. Lee, U. Liyanage, P. E. Bickel, W. Xia, P. T. Lansbury, Jr., and K. S. Kosik, *Nat Med*, **4**: 730–734 (1998).
  14. J. Jordan, M. F. Galindo, R. J. Miller, C. A. Reardon, G. S. Getz, and M. J. LaDu, *J Neurosci*, **18**: 195–204 (1998).
  15. J. F. Lees, P. S. Shneidman, S. F. Skuntz, M. J. Carden, and R. A. Lazzarini, *Embo J*, **7**: 1947–1955 (1988).
  16. J. A. Cohlberg, H. Hajarian, T. Tran, P. Alipourjeddi, and A. Noveen, *J Biol Chem*, **270**: 9334–9339 (1995).
  17. C. L. Leung, and R. K. H. Liem, *J Biol Chem*, **271**: 14041–14044 (1996).

18. M. K. Lee, Z. Xu, P. C. Wong, and D. W. Cleveland, *J Cell Biol*, **122**: 1337–50 (1993).
19. G. Y. Ching, and R. K. Liem, *J Cell Biol*, **122**: 1323–1335 (1993).
20. D. A. Carpenter, and W. Ip, *J Cell Sci*, **109**: 2493–2498 (1996).
21. D. Sun, C. L. Leung, and R. K. H. Liem, *J Biol Chem*, **271**: 14245–14251 (1996).
22. J. F. Leterrier, J. Kas, J. Hartwig, R. Vegners, and P. A. Janmey, *J Biol Chem*, **271**: 15687–15694 (1996).
23. N. Hirokawa, *J Cell Biol*, **94**: 129–142 (1982).
24. T. Nakagawa, J. Chen, Z. Zhang, Y. Kanai, and N. Hirokawa, *J Cell Biol*, **129**: 411–429 (1995).
25. M. P. Kaplan, S. S. Chin, K. H. Fliegner, and R. K. Liem, *J Neurosci*, **10**: 2735–2748 (1990).
26. M. J. Carden, J. Q. Trojanowski, W. W. Schlaepfer, and V. M. Lee, *J Neurosci*, **7**: 3489–3504 (1987).
27. D. Dahl, *J Neurosci Res*, **20**: 431–441 (1988).
28. W. W. Schlaepfer, and J. Bruce, *J Neurochem*, **55**: 453–460 (1990).
29. J. P. Julien, F. Cote, L. Beaudet, M. Sidky, D. Flavell, F. Grosveld, and W. Mushynski, *Gene*, **68**: 307–314 (1988).
30. R. A. Nixon, and S. E. Lewis, *J Biol Chem*, **261**: 16298–16301 (1986).
31. S. Hisanaga, Y. Gonda, M. Inagaki, A. Ikai, and N. Hirokawa, *Cell Regul*, **1**: 237–248 (1990).
32. A. C. Pant, Veeranna, H. C. Pant, and N. Amin, *Brain Res*, **765**: 259–266 (1997).
33. Veeranna, N. D. Amin, N. G. Ahn, H. Jaffe, C. A. Winters, P. Grant, and H. C. Pant, *J Neurosci*, **18**: 4008–4021 (1998).
34. R. Starr, F. L. Hall, and M. J. Monteiro, *J Cell Sci*, **109**: 1565–1573 (1996).
35. N. P. Bajaj, and C. C. Miller, *J Neurochem*, **69**: 737–743 (1997).
36. B. I. Giasson, and W. E. Mushynski, *J Biol Chem*, **271**: 30404–30409 (1996).
37. J. Brownlees, A. Yates, N. P. Bajaj, D. Davis, B. H. Anderton, P. N. Leigh, C. E. Shaw, and C. C. Miller, *J Cell Sci*, **113**: 401–407 (2000).
38. J. Xiao, and M. J. Monteiro, *J Neurosci*, **14**: 1820–1833 (1994).
39. W. T. Link, A. Dosemeci, C. C. Floyd, and H. C. Pant, *Neurosci Lett*, **151**: 89–93 (1993).
40. R. A. Nixon, S. E. Lewis, D. Dahl, C. A. Marotta, and U. C. Drager, *Brain Res Mol Brain Res*, **5**: 93–108 (1989).
41. R. A. Nixon, P. A. Paskevich, R. K. Sihag, and C. Y. Thayer, *J Cell Biol*, **126**: 1031–1046 (1994).
42. R. Lasek, and P. Hoffman (1976) in *Cell Motility: Microtubules and Related Proteins* (Goldman, R., Pollard, T., and J. Rosenbaum, Eds.), Vol. **3**, pp. 1021–1049, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.
43. L. Wang, C. L. Ho, D. Sun, R. K. Liem, and A. Brown, *Nat Cell Biol*, **2**: 137–141 (2000).
44. S. Terada, T. Nakata, A. C. Peterson, and N. Hirokawa, *Science*, **273**: 784–788 (1996).
45. S. Takeda, S. Okabe, T. Funakoshi, and N. Hirokawa, *J Cell Biol*, **127**: 173–185 (1994).
46. R. Nixon (1991) in *The Neuronal Cytoskeleton* (Burgoyne, R., Ed.), pp. 283–307, Wiley-Liss, New York.
47. S. E. Lewis, and R. A. Nixon, *J Cell Biol*, **107**: 2689–2701 (1988).
48. S. T. Hsieh, G. J. Kidd, T. O. Crawford, Z. Xu, W. M. Lin, B. D. Trapp, D. W. Cleveland, and J. W. Griffin, *J Neurosci*, **14**: 6392–6401 (1994).
49. R. A. Nixon, and K. B. Logvinenko, *J Cell Biol*, **102**: 647–659 (1986).
50. S. Okabe, H. Miyasaka, and N. Hirokawa, *J Cell Biol*, **121**: 375–386 (1993).
51. R. A. Nixon, *Bioessays*, **20**: 798–807 (1998).
52. I. Sanchez, L. Hassinger, R. K. Sihag, D. W. Cleveland, P. Mohan, and R. A. Nixon, *Journal*

of Cell Biology, **In Press** (2000).

53. D. R. Archer, D. F. Watson, and J. W. Griffin, *J Neurochem*, **62**: 1119–1125 (1994).
54. D. F. Watson, K. P. Fittro, P. N. Hoffman, and J. W. Griffin, *Brain Res*, **539**: 103–109 (1991).
55. I. Sanchez, L. Hassinger, P. A. Paskevich, H. D. Shine, and R. A. Nixon, *J Neurosci*, **16**: 5095–5105 (1996).
56. M. Sandberg, A. Hamberger, I. Jacobson, and J. O. Karlsson, *Neurochem Res*, **5**: 1185–1198 (1980).
57. J. D. Glass, and J. W. Griffin, *J Neurosci*, **11**: 3146–3154 (1991).
58. M. Wiedau-Pazos, J. J. Goto, S. Rabizadeh, E. B. Gralla, J. A. Roe, M. K. Lee, J. S. Valentine, and D. E. Bredesen, *Science*, **271**: 515–518 (1996).
59. D. P. Stromska, and S. Ochs, *J Neurobiol*, **12**: 441–453 (1981).
60. D. F. Watson, P. N. Hoffman, K. P. Fittro, and J. W. Griffin, *Brain Res*, **477**: 225–232 (1989).
61. W. W. Schlaepfer, C. Lee, J. Q. Trojanowski, and V. M. Lee, *J Neurochem*, **43**: 857–864 (1984).
62. H. C. Pant, *Biochem J*, **256**: 665–668 (1988).
63. P. N. Hoffman, G. W. Thompson, J. W. Griffin, and D. L. Price, *J Cell Biol*, **101**: 1332–1340 (1985).
64. P. N. Hoffman, J. W. Griffin, B. G. Gold, and D. L. Price, *J Neurosci*, **5**: 2920–2929 (1985).
65. P. N. Hoffman, D. W. Cleveland, J. W. Griffin, P. W. Landes, N. J. Cowan, and D. L. Price, *Proc Natl Acad Sci U S A*, **84**: 3472–3476 (1987).
66. Z. Xu, J. R. Marszalek, M. K. Lee, P. C. Wong, J. Folmer, T. O. Crawford, S. T. Hsieh, J. W. Griffin, and D. W. Cleveland, *J Cell Biol*, **133**: 1061–1069 (1996).
67. O. Ohara, Y. Gahara, T. Miyake, H. Teraoka, and T. Kitamura, *J Cell Biol*, **121**: 387–395 (1993).
68. Q. Zhu, S. Couillard-Despres, and J. P. Julien, *Exp Neurol*, **148**: 299–316 (1997).
69. J. Eyer, and A. Peterson, *Neuron*, **12**: 389–405 (1994).
70. J. R. Marszalek, T. L. Williamson, M. K. Lee, Z. Xu, P. N. Hoffman, M. W. Becher, T. O. Crawford, and D. W. Cleveland, *J Cell Biol*, **135**: 711–724 (1996).
71. F. Cote, J. F. Collard, and J. P. Julien, *Cell*, **73**: 35–46 (1993).
72. Q. Zhu, M. Lindenbaum, F. Levavasseur, H. Jacomy, and J. P. Julien, *J Cell Biol*, **143**: 183–193 (1998).
73. G. A. Elder, V. L. Friedrich, Jr., P. Bosco, C. Kang, A. Gourov, P. H. Tu, V. M. Lee, and R. A. Lazzarini, *J Cell Biol*, **141**: 727–739 (1998).
74. G. A. Elder, V. L. Friedrich, Jr., C. Kang, P. Bosco, A. Gourov, P. H. Tu, B. Zhang, V. M. Lee, and R. A. Lazzarini, *J Cell Biol*, **143**: 195–205 (1998).
75. M. V. Rao, M. K. Houseweart, T. L. Williamson, T. O. Crawford, J. Folmer, and D. W. Cleveland, *J Cell Biol*, **143**: 171–181 (1998).
76. J. Trojanowski, M. Schmidt, R. Shin, G. Bramblett, D. Rao, and M. Lee, *Brain Path*, **3**: 45–54 (1993).
77. T. C. Dickson, C. E. King, G. H. McCormack, and J. C. Vickers, *Exp Neurol*, **156**: 100–110 (1999).
78. J. C. Vickers, J. H. Morrison, V. L. Friedrich, Jr., G. A. Elder, D. P. Perl, R. N. Katz, and R. A. Lazzarini, *J Neurosci*, **14**: 5603–5612 (1994).
79. A. Hirano, H. Donnenfeld, S. Sasaki, and I. Nakano, *J Neuropathol Exp Neurol*, **43**: 461–470 (1984).
80. A. Hirano, I. Nakano, L. T. Kurland, D. W. Mulder, P. W. Holley, and G. Saccomanno, *J Neuropathol Exp Neurol*, **43**: 471–480 (1984).
81. Z. Xu, L. C. Cork, J. W. Griffin, and D. W. Cleveland, *Cell*, **73**: 23–33 (1993).

82. P. H. Tu, M. E. Gurney, J. P. Julien, V. M. Lee, and J. Q. Trojanowski, *Lab Invest*, **76**: 441–456 (1997).
83. B. J. Gibb, J. P. Brion, J. Brownlees, B. H. Anderton, and C. C. Miller, *J Neurochem*, **70**: 492–500 (1998).
84. P. C. Wong, J. Marszalek, T. O. Crawford, Z. Xu, S. T. Hsieh, J. W. Griffin, and D. W. Cleveland, *J Cell Biol*, **130**: 1413–1422 (1995).
85. M. K. Lee, J. R. Marszalek, and D. W. Cleveland, *Neuron*, **13**: 975–988 (1994).
86. D. A. Figlewicz, A. Krizus, M. G. Martinoli, V. Meininger, M. Dib, G. A. Rouleau, and J. P. Julien, *Hum Mol Genet*, **3**: 1757–1761 (1994).
87. J. Tomkins, P. Usher, J. Y. Slade, P. G. Ince, A. Curtis, K. Bushby, and P. J. Shaw, *Neuroreport*, **9**: 3967–3970 (1998).
88. A. Al-Chalabi, P. M. Andersen, P. Nilsson, B. Chioza, J. L. Andersson, C. Russ, C. E. Shaw, J. F. Powell, and P. N. Leigh, *Hum Mol Genet*, **8**: 157–164 (1999).
89. I. V. Mersyanova, A. V. Perepelov, A. V. Polyakov, V. F. Sitnikov, E. L. Dadali, R. B. Oparin, A. N. Petrin, and O. V. Evgrafov, *Am J Hum Genet*, **67**: 37–46 (2000).
90. D. R. Rosen, T. Siddique, D. Patterson, D. A. Figlewicz, P. Sapp, A. Hentati, D. Donaldson, J. Goto, J. P. O'Regan, H. X. Deng, and et al., *Nature*, **362**: 59–62 (1993).
91. T. L. Williamson, L. I. Bruijn, Q. Zhu, K. L. Anderson, S. D. Anderson, J. P. Julien, and D. W. Cleveland, *Proc Natl Acad Sci U S A*, **95**: 9631–9636 (1998).
92. S. Couillard-Despres, Q. Zhu, P. C. Wong, D. L. Price, D. W. Cleveland, and J. P. Julien, *Proc Natl Acad Sci U S A*, **95**: 9626–9630 (1998).
93. G. Y. Ching and R. K. Liem, *J Cell Sci*, **112**: 2233–40 (1999).
94. R. Hashimoto, Y. Nakamura, S. Komai, Y. Kashiwagi, K. Tamura, T. Goto, S. Aimoto, K. Kaibuchi, S. Shiosaka, M. Takeda, *J Neurochem*, **75**: 373–82 (2000).
95. J. T. Yabe, W. K. Chan, T. M. Chylinski, S. Lee, A. F. Pimenta, T. B. Shea, *Cell Motil Cytoskeleton*, **48**: 61–83 (2001).
96. V. Prahlad, B. T. Helfand, G. M. Langford, R. D. Vale, R. D. Goldman, *J Cell Sci*, **113**: 3939–46 (2000).
97. Z. Xu and V. W. Tung, *Brain Res*, **866**: 326–32 (2000).
98. J. T. Yabe, T. Chylinski, F. S. Wang, A. Pimenta, S. D. Kattar, M. D. Linsley, W. K. Chan, T. B. Shea, *J Neurosci*, **21**: 2195–205 (2001).
99. F. Levavasseur, Q. Zhu, J. P. Julien, *Brain Res Mol Brain Res*, **69**: 104–12 (1999).
100. H. Jacomy, Q. Zhu, S. Couillard-Despres, J. M. Beaulieu, J. P. Julien, *J Neurochem*, **73**: 972–84 (1999).
101. T. M. Svitkina, A. B. Verkhovskiy, G. G. Borisy, *J. Cell Biol.*, **135**: 991–1007 (1996).
102. Y. Yang, J. Dowling, Q. C. Yu, P. Kouklis, D. W. Cleveland, E. Fuchs, *Cell*, **86**: 655–65 (1996).
103. C. L. Leung, D. Sun, M. Zheng, D. R. Knowles, R. K. Liem, *J. Cell Biol.*, **147**: 1275–1286 (1999).
104. G. Bernier, M. Pool, M. Kilcup, F. Alfoldi, Y.-E. Repentigny, R. Kothary, *Dev. Dyna.* **219**: 216–225 (2000).
105. Q. Zhu, M. Lindenbaum, F. Levavasseur, H. Jacomy, J. P. Julien, *J. Cell Biol.*, **143**: 183–93 (1998).
106. M. V. Rao, M. K. Houseweart, T. L. Williamson, T. O. Crawford, J. Folmer, D. W. Cleveland, *J. Cell Bio.*, **143**: 171–181 (1998).
107. G. A. Elder, V. L. Friedrich Jr., C. Kang, P. Bosco, A. Gourov, P. H. Tu, B. Zhang, V. M. Lee, R. A. Lazzarini, *J. Cell Bio.*, **143**: 195–205 (1998).
108. S. Couillard-Despres, J. Meier, J. P. Julien, *Neurobiol. Dis.*, **7**: 462–70 (2000).

## Suggestions for Further Reading

109. RD Burgoyne (ed.) (1991) *The Neuronal Cytoskeleton*, Wiley-Liss, New York. Contains comprehensive chapters on all aspects of neurofilaments and dynamic behavior.
110. RA Nixon. (1998) The dynamic behavior and organization of cytoskeletal proteins in neurons: Reconciling old and new findings, *Bioessays* **20**, 798–807. Recent synthetic review of historical and current data on neurofilament behavior.
111. MK Lee and DW Cleveland (1996) Neuronal intermediate filaments. *Ann. Rev. Neurosci.* **19**, 187–217. A concise, up-to-date review on all aspects of neurofilaments, including pathology.

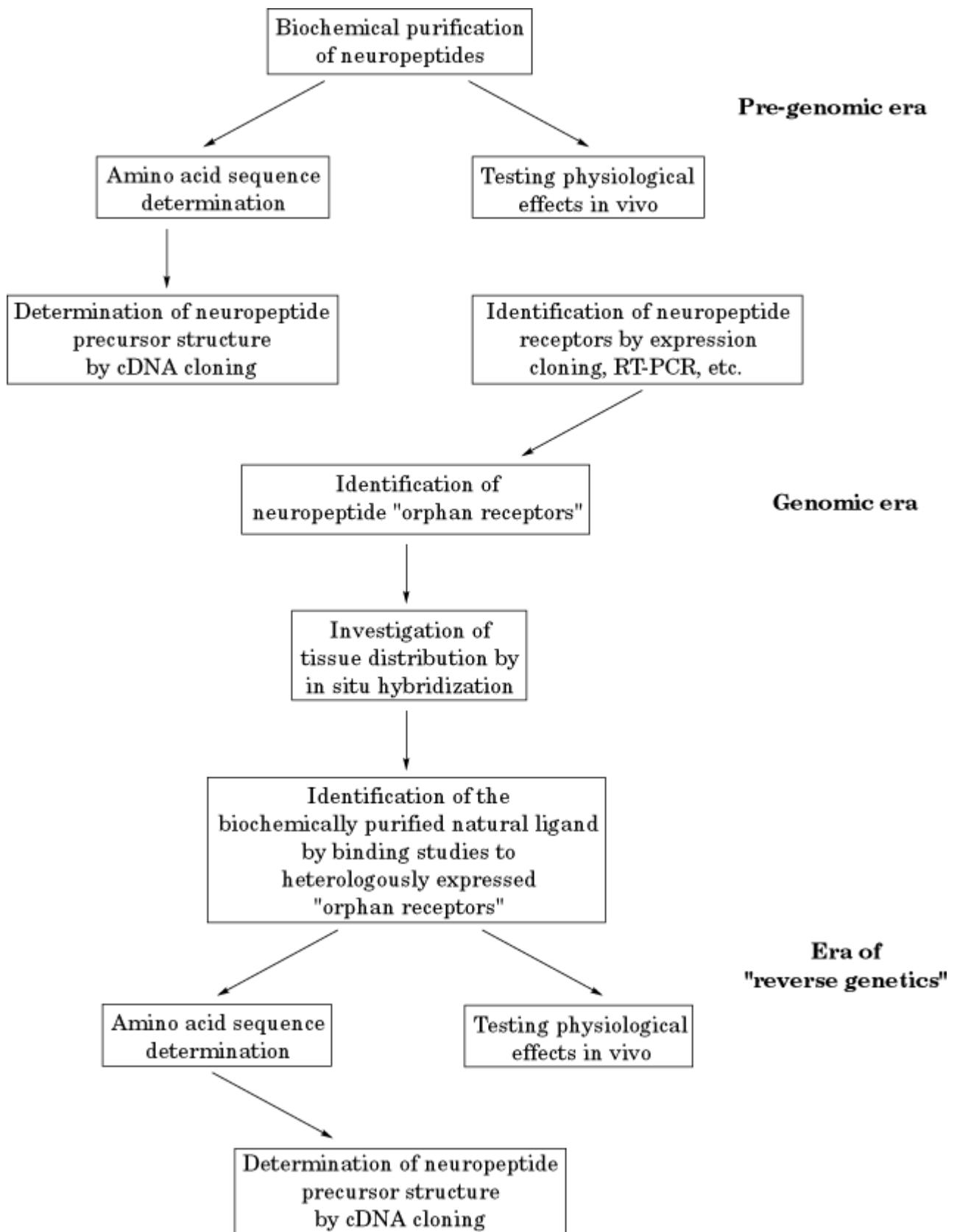
## Neuropeptides

### 1. History

By definition, neuropeptides are synthesized in nerve cells of the central and peripheral nervous systems and act as either neurohormones, neurotransmitters, or neuromodulators. They regulate and influence the behavior and homeostasis of a given organism in many and complex ways.

The concept of the peptidergic neuron in vertebrates and invertebrates was developed by Ernst and Berta Scharrer in 1928, and their work still forms the basis of contemporary neuropeptide research. Numerous neuropeptides were discovered in the late 1960s and early 1970s by work that began as a tour de force, and resulted in their biochemical isolation from huge amounts of various tissues. The advent of molecular biological techniques in the 1980s greatly facilitated the identification and characterization of novel neuropeptides in vertebrates and invertebrates. This “era of neuropeptides” was followed by extensive work aimed at identifying their cognate receptor proteins. Today, at least one type of seven-membrane-spanning-domain G-protein-coupled receptor (GPCR) has been identified for each neuropeptide. A growing number of receptors also share structural similarities with known neuropeptide receptors or neuropeptide receptor families. These putative neuropeptide receptors, referred to as “orphan receptors,” await the identification of their ligands. As a consequence, neuropeptide research has experienced a revival in the late 1990s, and now strategies called “reverse genetics” (an example is given later) are being developed to unravel the identity of the natural ligands for these orphan GPCRs and their functions within the central and peripheral nervous systems. Figure 1 schematically summarizes the history of neuropeptide and neuropeptide receptor research from the 1960s to the late 1990s.

**Figure 1.** Outlook on the history of neuropeptide and neuropeptide receptor research. During the “pregenomic” era, biochemical work prevailed. Neuropeptides were purified to homogeneity by liquid [chromatography](#) procedures. The amino acid sequence was determined by [Edman Degradation](#). Physiological effects were analyzed by *in vivo* administration of the purified neuropeptide. Identification and characterization continued during the “genomic” era. Neuropeptide precursor structures were determined by sequencing the corresponding cDNAs. A little later, the structures of the first neuropeptide receptors became known. This revealed that they belong to the large group of seven membrane-spanning-domain, G-protein-coupled receptors. Different procedures have been applied to receptor characterization, including expression in *Xenopus* oocytes (expression cloning) and polymerase chain reaction (PCR) amplification after **reverse transcription** of the mRNA (RT-PCR). In some instances, receptors (“orphan receptors”) were identified that failed to be activated by any of the known neuropeptides. In the era of “reverse genetics,” work is aimed at identifying natural ligands. This can be done by expressing “orphan receptors” in cell culture and analyzing receptor activation after applying of peptide extracts isolated from appropriate tissues. These “activities” are finally purified and analyzed, as described previously.



## 2. Neuropeptides and Neuropeptide Families

It is not feasible to list of all the neuropeptides discovered thus far in nerve cells throughout the animal kingdom. As a matter of fact, nearly every [peptide hormone](#) of a given organism is located in



nerve cells and consequently can be regarded as a neuropeptide. Table 1 shows a selection of mammalian neuropeptides, some of which can be grouped into families according to structural or functional similarities. It should be stressed, however, that this classification is somewhat arbitrary. For example, the hormones vasopressin (VP) and oxytocin (OT) (the vasopressin/oxytocin family) are synthesized in magnocellular neurons of the hypothalamo-neurohypophyseal tract and are secreted from the nerve terminals in the posterior pituitary into the systemic circulation to exert their biological functions in various peripheral organs (VP: kidney, liver, blood vessels; OT: smooth muscles of the uterus and mammary gland). They are also synthesized in parvocellular neurons of the hypothalamic paraventricular nucleus and have axonal projections to the portal blood vessel system that delivers neuropeptides to the anterior pituitary, where they act as release factors for adenohypophyseal hormones. In addition, VP and OT are produced in various extrahypothalamic neurons that have central projections and thus presumably play a role as neurotransmitters and/or neuromodulators. The same is principally true for the “hypothalamic-release and release-inhibiting hormones” thyrotropin-releasing hormone (TRH), somatostatin (SST), luteinizing-hormone releasing hormone (LHRH), etc., all of which are produced in nerve cells that have central projections. Other neuropeptides, such as cholecystokinin (CCK) and vasoactive intestinal polypeptide (VIP), were initially attributed with functions outside the brain before their synthesis and actions in the central nervous system were acknowledged.

**Table 1. Selected Mammalian Neuropeptides**

| Neuropeptide Family   | Number of Amino Acid Residues |
|---|-------------------------------|
| Oxytocin (OT)   | 9                             |
| Vasopressin (VP)  | 9                             |
| <i>Hypothalamic-releasing and release-inhibiting hormones</i> |                               |
| Corticotropin-releasing hormone (CRH)                         | 41                            |
| Growth-hormone releasing hormone (GHRH)                       | 44                            |
| Luteinizing-hormone releasing hormone (LHRH)                  | 10                            |
| Somatostatin (SST)  | 14 and 28                     |
| Thyrotropin-releasing hormone (TRH)                           | 3                             |
| <i>Opioid peptide family</i>                                  |                               |
| b-Endorphin   | 30                            |
| Dynorphin (DYN)   | 17, other forms               |
| Leu-Enkephalin (Leu-Enk)                                      | 5                             |
| Met-Enkephalin (Met-Enk)                                      | 5                             |
| Nociceptin/orphanin FQ  | 17, other forms?              |
| <i>Tachykinin family</i>                                      |                               |
| Neurokinin A, Substance K                                     | 10                            |
| Neurokinin B  | 10                            |
| Neuropeptide K  | 36                            |
| Substance P (SP)  | 11                            |
| <i>VIP/Glucagon family</i>                                    |                               |
| Glucagon-like peptide-1 (GLP-1)                               | 29                            |

|  |                        |
|--|------------------------|
| Peptide histidine-isoleucine (PHI)                     | 27                     |
| Pituitary adenylate cyclase activating peptide (PACAP) | 27 or 38               |
| Vasoactive intestinal polypeptide (VIP)                | 28                     |
| <i>NP Y family</i>                                     |                        |
| Neuropeptide tyrosine (NP Y)                           | 36                     |
| Pancreatic polypeptide (PPP)                           | 36                     |
| Peptide tyrosine-tyrosine (P YY)                       | 36                     |
| <i>Other peptides</i>                                  |                        |
| Brain natriuretic peptide (BNP)                        | 32                     |
| Calcitonin gene-related peptide (CGRP)                 | 37                     |
| Cholecystokinin (CCK)                                  | 8                      |
| Cortistatin (CST)                                      | 14 or 29, other forms? |
| Galanin (GAL)  | 29 or 30               |
| Hypocretin/orexin                                      | 14 or 29               |
| Neurotensin  | 13                     |
| Parathyroid hormone-related protein                    | 34 or 37               |

---

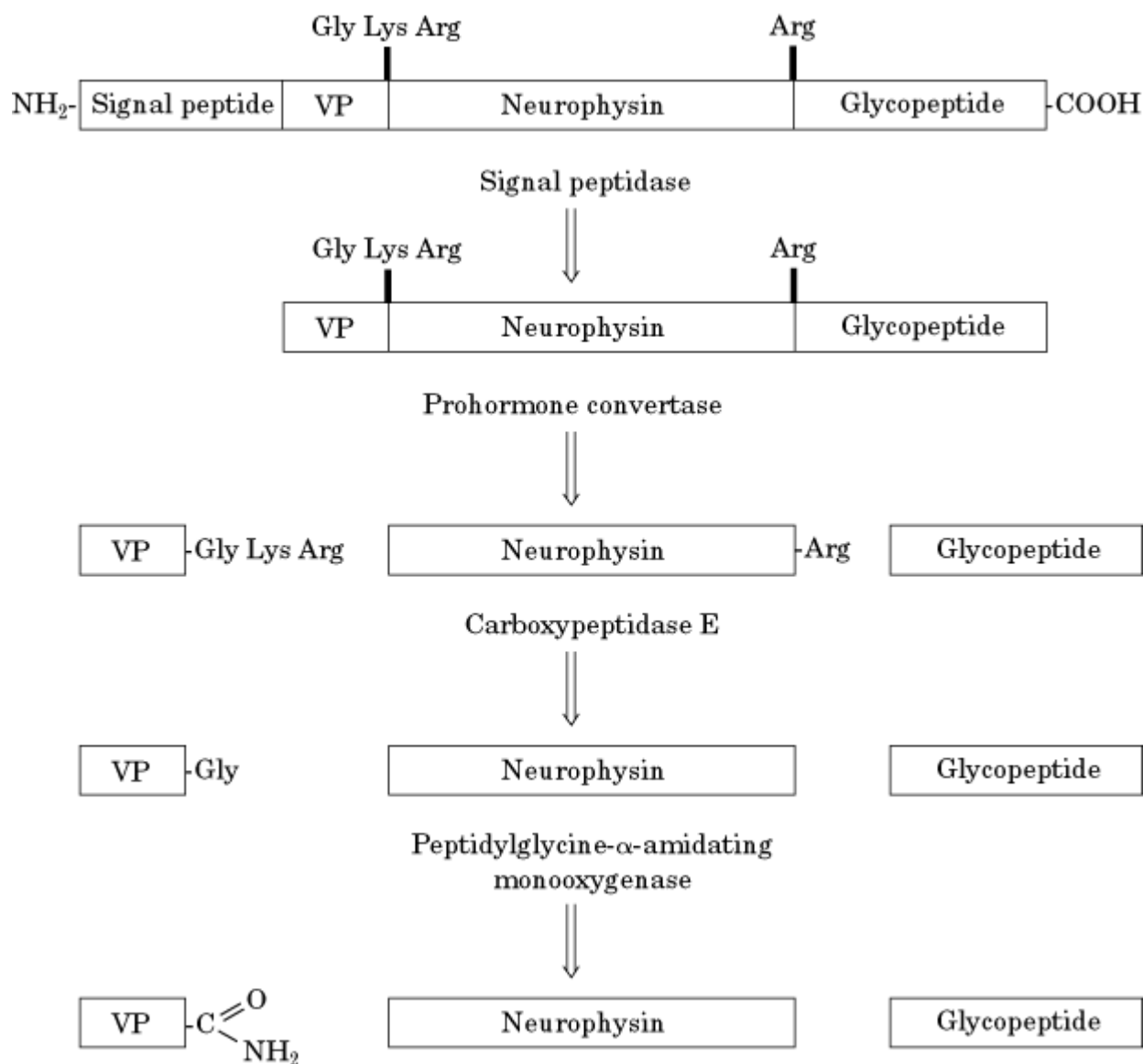
### 3. Neuropeptide Biosynthesis, Maturation, and Release

Neuropeptides are synthesized as parts of larger precursor molecules, the pre-prohormones, on the rough [endoplasmic reticulum](#) (ER). These precursors consist of a [signal peptide](#) that is cotranslationally cleaved off in the lumen of the rough ER by signal peptidases to yield the prohormones. After passage through the **Golgi complex**, where they are eventually equipped with complex carbohydrate moieties, the precursors are sorted to the regulatory secretory pathway in neurosecretory granules (also termed dense-core vesicles) by mechanisms that remain to be fully elucidated. The biologically active entities are generated in the trans-Golgi network or, more frequently, within immature **vesicles** by the enzymatic activities of endopeptidases, the prohormone convertases (PCs). PCs belong to the family of [subtilisin](#)/kexin-like [serine proteinases](#) (1). The enzymes cleave precursor molecules at the carboxyl side of either specific single or pairs of basic amino acid residues (most commonly Lys-Arg or Arg-Arg). A number of PCs have been identified. Of these, PC1 and PC2 are expressed primarily in endocrine and nerve cells. The supernumerary basic amino acids are subsequently removed by carboxypeptidase E.

The molecular composition of neuropeptide precursors is rather diverse. The VP and OT precursors, for example, harbor one copy of the respective peptide, whereas others, such as the TRH precursor, contain multiple copies of the same peptide. In a few cases, for example, in the opioid peptide precursor family, different sets of peptides arise from processing. Thus, the human [enkephalin](#) precursor contains six copies of Met-enkephalin and one copy of Leu-enkephalin (2). Before the mature neuropeptides are secreted, modification at the carboxy- and/or amino-terminal ends is frequently required for biological activity. The most common modification is the conversion of a carboxy-terminal glycine into the  $\alpha$ -amide, a two-step reaction that is catalyzed by peptidylglycine  $\alpha$ -amidating monooxygenase and requires copper, molecular oxygen and ascorbate (3, 4). Glutamyl cyclase converts the amino-terminal glutamine residue in TRH, LHRH and a few other peptides into a pyroglutamyl residue (5). Figure 2 schematically depicts maturation steps of preprohormones exemplified by the VP precursor. It should be mentioned, however, that processing of proneuropeptides in some invertebrates, particularly in cnidarians, might be even more complex and

involve enzymes that have no counterparts in mammalian species (6). Cloning of the genes for various cnidarian neuropeptide precursors has shown that sequences other than monobasic and dibasic residues must also function as processing signals. These include Thr, Ser, and Tyr residues, and X-Pro or X-Ala sequences (where X is any amino acid). These novel processing signals are located exclusively at the amino-terminal side of neuropeptides. The proteolytic enzymes remain to be fully characterized.

**Figure 2.** Structural organization, processing, and modification of pre-provasopressin. Neuropeptide precursor processing and further maturation of the neuropeptide is schematically shown for pre-provasopressin. The preprohormone consists of a signal peptide, the nonapeptide vasopressin (VP), a carrier protein called neurophysin, and a glycopeptide moiety of unknown function. The signal peptide is removed by signal peptidase, an enzyme that resides in the lumen of the rough endoplasmic reticulum. Subsequent maturation steps take place in immature dense-core vesicles, that are directed to the regulated secretory pathway. The VP precursor contains two processing signals for endoproteolytic cleavage by prohormone convertases: a dibasic sequence Lys-Arg on the amino-terminal side and a monobasic Arg residue on the carboxy-terminal side. The basic amino acids are finally removed by carboxypeptidase E. The Gly residue preceding the dibasic signal is important for subsequently modifying the neuropeptide. Peptidylglycine- $\alpha$ -amidating monooxygenase catalyzes the conversion of the carboxy-terminal Gly residue into an amide in a two-step reaction. Amidation of VP is required for its biological activity. Secretion of the mature neuropeptide takes place by exocytosis from dense-core vesicles after appropriate physiological stimulation.



Unlike classical transmitter-containing synaptic vesicles, which are organized in an active zone at the synapse, in principle, dense-core vesicles can be released from any part of the nerve cell. Initially, it was thought that release took place only from the axonal compartment. There is now compelling evidence, however, for dendritic release of at least some neuropeptides. This has been documented unequivocally for VP and OT, which can be secreted from the dendrites of hypothalamic magnocellular neurons into the central nervous system (7). In some instances, neurons synthesize both classical transmitters (for example noradrenaline, serotonin, acetylcholine, and GABA) and neuropeptides, and the parameters necessary for their release are quite distinct. Neuropeptide secretion often requires high-frequency or prolonged firing activity. Classical transmitter release, in contrast, follows the action potential in a one-to-one fashion (2, 8). At normal firing rates, neuropeptides might not be secreted and should consequently not exert any postsynaptic effects at all. High-frequency stimulation, on the other hand, allows concomitant release of classical transmitters and neuropeptides. A number of examples demonstrate either synergistic or antagonistic effects of classical transmitter/neuropeptide interactions. A synergistic interaction of acetylcholine and VIP has been demonstrated, for example, in the cat salivary gland. Acetylcholine induces saliva secretion. VIP facilitates secretion by triggering vasodilation, and it enhances acetylcholine binding to the [acetylcholine receptor](#). In some noradrenergic neurons of the mouse *vas deferens*, the coexisting neuropeptide Y (NP Y) plays a role as a functional antagonist by inhibiting noradrenaline release (8).

#### 4. Neuropeptide Inactivation in the Central Nervous System

Enzymatic degradation of neuropeptides is the major pathway for terminating the action of the signaling molecules and for adjusting neuropeptide concentrations in the extracellular fluid. To date, a small number of peptidases in the central nervous system that act as neuropeptide-degrading ectoenzymes have been characterized, and most of these also exist in the periphery, for example, in the kidney and in other tissues (9). Endopeptidase-24.11 (E-24.11) was first detected by its ability to inactivate enkephalins. The enzyme is distributed widespread within the central nervous system and is present predominantly on nerve cells at both pre- and postsynaptic sites. Besides enkephalins, other neuropeptides are effectively hydrolyzed by E-24.11 *in vitro*. Presumably, the enzyme acts at numerous peptidergic synapses. Aminopeptidase-N (AP-N) was identified as the major membrane-associated aminopeptidase engaged in the hydrolysis of enkephalins. This enzyme is also heterogeneously distributed on pre- and postsynaptic elements and, in addition, on nonneuronal cells, such as astrocytes. *In vitro*, AP-N shows broad substrate specificity and hydrolyzes the amino-terminal amino acid residues from nearly all unsubstituted peptides. Even though AP-N could play a role in inactivating all neuropeptides that have unprotected amino-termini, its action *in vivo* has only been confirmed for the enkephalins. Another aminopeptidase that has surprisingly stringent substrate specificity was identified recently. Pyroglutamyl aminopeptidase II (PGA II), a membrane-associated metallopeptidase, hydrolyzes the pyroglutamate-histidine bond of TRH. The enzyme is highly specific for TRH inactivation. It cannot degrade other neuropeptides that have pyroglutamyl-protected amino-termini, including LHRH, neurotensin, and gastrin (10). A number of neuropeptides are hydrolyzed *in vitro* by angiotensin-converting enzyme (ACE), whose main role is to convert angiotensin I to angiotensin II and inactivate bradykinin. Despite the broad distribution of ACE within the central nervous system of a variety of mammalian species, however, no clear cut evidence for its involvement in neuropeptide metabolism other than the previously mentioned peptides has been obtained.

Additional peptidases, it is claimed are involved in neuropeptide degradation/inactivation, but their precise roles have remained obscure. This is largely due to the fact that the majority of these enzymes have a predominantly cytosolic (or endosomal) localization, which contradicts a role in neuropeptide inactivation. In some instances, membrane-associated forms might exist, generated, for example, by [alternative splicing](#) of the corresponding pre-messenger RNA. More work must be done to clarify the precise roles of these enzymes. It is to be anticipated that novel ectoenzymes with neuropeptide-degrading activity will be discovered in the next few years (9).

## 5. Neuropeptide Signal Transduction via G-protein-Coupled Receptors

Defining the target molecules of neuropeptides, in particular those of opioid peptides, has been of great interest because this might pave the way for drug design (either agonists or antagonists) for treating analgesia, provided these compounds pass the blood–brain barrier. Despite extensive scientific work all over the world, it was not until the late 1980s that the structures of the first neuropeptide receptors (two tachykinin receptors and one neurotensin receptor) were elucidated by molecular biological techniques. To date, at least one receptor for each neuropeptide is known, and very often multiple forms exist, as in the case of SST receptors in rat, mouse, and humans. With rare exceptions, all neuropeptide receptors belong to the large family of GPCRs. Their subcellular pre- and/or postsynaptic localization within the target cells of neuropeptides and their regulation and desensitization are under intensive study now to clarify more precisely neuropeptide actions within the brain and in the peripheral nervous system ([11-14](#)).

## 6. Physiological Roles of Neuropeptides

The function of neuropeptides as neurohormones in the periphery is quite well defined. Their role in the central nervous system, in contrast, has largely remained elusive, despite vast literature in this field. Neuropeptide actions in the brain are extremely diverse and presumably depend on various parameters, for example, on the projection areas of peptidergic neurons and on the receptor subtype and its location at pre- and/or postsynaptic sites. It is believed that neuropeptides, either alone or in concert with classical transmitters, are engaged in aspects of learning and memory, pain perception, feeding and obesity, anxiety, aggression, reproductive behavior, sleep-wakefulness, etc. It is beyond the scope of this article to review fully the presumed central actions of all neuropeptides. Therefore, the reader is referred to the original literature in this field. Selected publications and review articles with particular reference to the roles ascribed to individual neuropeptides or neuropeptide families in the central and peripheral nervous system are listed in “Suggestions for Further Reading.”

To demonstrate the complex interplay of a variety of neuropeptides (and other signaling molecules), possible mechanisms involved in the control of feeding and energy homeostasis are briefly described. Brain lesion experiments have suggested that the hypothalamus plays an important role in energy homeostasis. Disruption of the ventromedial hypothalamus (“satiety center”) leads to obesity, whereas animals that have lesions of the lateral hypothalamus (“feeding center”) lose weight and may even starve to death in some instances. The neurochemical system in the brain that underlies the complex regulation of body weight includes a variety of classical transmitters and neuropeptides ([15](#), [16](#)). Recent genetic mouse models, molecular cloning techniques, and behavioral studies have succeeded in identifying at least some of the important positive and negative regulators of energy balance. Several neuropeptides, injected into the brain, stimulate feeding behavior. These include, for example, NP Y, galanin, melanin-concentrating hormone, and the recently identified orexins/hypocretins. A negative effect on food consumption is observed with corticotropin-releasing hormone, the closely related and even more potent urocortin, and a melanocortin-stimulating hormone (a MSH), one of the peptides derived from the complex and tissue-specific proopiomelanocortin (POMC) precursor processing. Interestingly, the action of a MSH at the melanocortin 4 receptor is antagonized by the Agouti protein, which is ectopically expressed in the central nervous system of the *lethal yellow* ( $A^Y/a$ ) mouse, an autosomal dominant obesity model, and by the Agouti-related peptide (AGRP), a newly discovered hypothalamic neuropeptide. It is believed that the resulting feeding behavior depends on the ratio of agonist to antagonist. There are almost certainly more factors that have to be identified. This is exemplified by a genetic mouse model that is deficient for an “orphan receptor,” termed bombesin receptor type-3, which also results in an obese phenotype. To date, the interplay between positive and negative regulators of food intake is far from clear. For example, targeted deletion of the NP Y precursor gene surprisingly failed to show a body weight or feeding behavior phenotype. Furthermore, the metabolic state of the periphery must be signaled to the brain by substances that cross the blood–brain barrier. Leptin, a hormone synthesized

in adipocytes, is a potentially interesting molecule that negatively and positively regulates NP Y and POMC expression, respectively. Although the physiological consequences of leptin deficiency are clear, namely, obesity, leptin levels are elevated in most human obese patients, and in rodent models, indicating that most obesity is leptin-resistant. Much work remains to be done before neuropeptide actions in the brain are precisely understood.

## 7. How Neuropeptides Are Studied

Earlier work on neuropeptide expression and distribution was based mainly on immunological procedures. **Antibodies** raised against particular peptides can be used to evaluate the tissue distribution and subcellular localization of neuropeptides. **Radioimmunoassays** continue to be used for measuring peptide secretion in different physiological (or diseased) states of a given organism. When molecular biological procedures became available, neuropeptide gene expression could be monitored qualitatively, and even quantitatively, by a variety of techniques including ***in situ* hybridization**, **Northern blot** analysis, and **solution hybridization**. In some instances, it might be interesting to measure the release of neuropeptides into the central nervous system *in vivo*, and this can be done by microdialysis and the application of push-pull systems. The biological role of neuropeptides in various aspects of behavior (learning, anxiety, feeding, reproduction, etc.) can be studied by injecting the peptide under investigation into the brain and subsequently using several behavioral paradigms. Intracerebroventricular injection of urocortin, for example, has suggested anxiogenic-like properties of this neuropeptide (17). Neuropeptide genes may also be “knocked out” by targeted deletion to investigate the physiological and/or behavioral effects of the absence of a neuropeptide. As with the NP Y “knock out” mouse mentioned before, however, the animals often do not show a phenotype deviating from that of normal individuals which carry the corresponding gene. This may be due to mechanisms operative during embryogenesis and subsequent development that ultimately compensate for the loss of function of an individual neuropeptide. To alleviate this drawback, a novel technique, called “conditional knockout,” is employed. This technique permits shutting off the expression of defined genes postnatally in a reversible manner.

Two different approaches have been chosen recently to identify novel neuropeptides, the orexins/hypocretins. In a search for the natural ligand of a cloned “orphan” receptor, the biologically active entity, termed orexin, that binds to and activates the heterologously expressed receptor was purified biochemically. After determining of the amino acid sequence of the peptide, degenerate oligonucleotide primers were designed and used to screen a [cDNA library](#). Employing this procedure, called “reverse genetics,” the structure of the orexin precursor was elucidated (18). Interestingly, the same precursor, called pre-prohypocretin, was identified by a second research group using a completely different approach. In this case, [subtractive hybridization](#) cDNA cloning techniques were employed to identify proteins specifically expressed in the rat hypothalamus (19). Other neuropeptides have been similarly identified recently. These include cortistatin, a neuropeptide structurally related to SST and possibly involved in sleep behavior (20). Another example is nociceptin/orphanin FQ, which is structurally related to the opioid peptide dynorphin A and has been characterized independently by two research groups (21, 22). It is likely that additional novel neuropeptides will be discovered as a result of the human [genome](#) project.

## Bibliography

1. N. G. Seidah and M. Chrétien (1997) *Curr. Opin. Biotechnol.* **8**, 602–607.
2. W. S. Sossin, J. M. Fisher, and R. H. Scheller (1989) *Neuron* **2**, 1407–1417.
3. A. F. Bradbury and D. G. Smyth (1987) *Biosci. Rep.* **7**, 907–916.
4. A. S. Kolhekar, R. E. Mains, and B. A. Eipper (1997) *Methods Enzymol.* **279**, 35–43.
5. W. H. Fischer and J. Spiess (1987) *Proc. Natl. Acad. Sci. USA* **84**, 3628–3632.
6. C. J. P. Grimmelikhuijzen and F. Hauser (1997) In *Neuroendocrinology. Retrospect and Perspectives* (H.-W. Korf and K.-H. Usadel, eds.), Springer-Verlag, Berlin, Germany, pp. 25–42.

7. J. F. Morris, D. V. Pow, H. W. Sokol, and A. Ward (1993) In *Vasopressin* (P. Gross, D. Richter, and G. L. Robertson, eds.), John Libbey Eurotext, Paris, France, pp. 171–182.
8. T. Hökfelt (1991) *Neuron* **7**, 867–879.
9. A. J. Turner (1997) In *Cell-Surface Peptidases in Health and Disease* (A. J. Kenny and C. M. Boustead, eds.), BIOS Scientific, Oxford, UK, pp. 275–301.
10. K. Bauer, H. Heuer, F. Iffländer, A. Peters, S. Schmitmeier, L. Schomburg, S. Turwitt, and M. Wilkens (1997) In *Cell-Surface Peptidases in Health and Disease* (A. J. Kenny and C. M. Boustead, eds.), BIOS Scientific, Oxford, UK, pp. 239–248.
11. B. Hille (1992) *Neuron* **9**, 187–195.
12. T. Gudermann, T. Schöneberg, and G. Schulz (1997) *Annu. Rev. Neurosci.* **20**, 399–427.
13. E. Grady, S. Bohm, K. McConalogue, A. Garland, J. Ansel, J. Olerud, and N. Bunnett (1997) *J. Invest. Dermatol. Symp. Proc.* **2**, 69–75.
14. J. Wess (1997) *FASEB J.* **11**, 346–354.
15. S. P. Kalra (1997) *Neuron* **19**, 227–230.
16. J. S. Flier and E. Maratos-Flier (1998) *Cell* **92**, 437–440.
17. J. L. Moreau, G. Kilpatrick, and F. Jenck (1997) *Neuroreport* **8**, 1697–1701.
18. T. Sakurai, A. Amemiya, M. Ishii, I. Matsuzaki, R. M. Chemelli, H. Tanaka, S. C. Williams, J. A. Richardson, G. P. Kozlowski, S. Wilson, J. R. S. Arch, R. E. Buckingham, A. C. Haynes, S. A. Carr, R. S. Annan, D. E. McNulty, W.-S. Liu, J. A. Terrett, N. A. Elshourbagy, D. J. Bergsma, and M. Yanagisawa (1998) *Cell* **92**, 573–585.
19. L. de Lecea, T. S. Kilduff, C. Peyron, X. B. Gao, P. E. Foye, P. E. Danielson, C. Fukuhara, E. L. F. Battenberg, V. T. Gautvik, F. S. Bartlett II., W. N. Frankel, A. N. van den Pol, E. F. Bloom, K. M. Gautvik, and J. G. Sutcliffe (1998) *Proc. Natl. Acad. Sci. USA* **95**, 322–327.
20. L. de Lecea, J. R. Criado, O. Prospero-Garcia, K. M. Gautvik, P. Schweizer, P. E. Danielson, C. L. M. Dunlop, G. R. Siggins, S. J. Henriksen, and J. G. Sutcliffe (1996) *Nature* **381**, 242–245.
21. J. C. Meunier, C. Mollereau, L. Toll, C. Suaudeau, C. Moisand, P. Alvinerie, J. L. Butour, J. C. Guillemot, P. Ferrara, B. Monsarrat, H. Mazaguil, G. Vassart, M. Parmentier, and J. Costentin (1995) *Nature* **377**, 532–535.
22. R. K. Reinscheid, H.-P. Nothacker, A. Bourson, A. Ardati, R. A. Henningsen, J. R. Bunzow, D. K. Grandy, H. Langen, F. J. Monsma Jr., and O. Civelli (1995) *Science* **270**, 792–794.

### **Suggestions for Further Reading**

23. E. R. de Kloet, V. M. Wiegant, and D. de Wied, eds. (1987) *Progress in Brain Research* **72**, "Neuropeptides and Brain Function", Elsevier, Amsterdam, Netherlands.
24. T. B. van Wiemersma Greidanus, ed. (1987) *Frontiers in Hormone Research* **15**, Karger, Basel, Switzerland.

### **Selected Publications and Review Articles**

25. A. Balasubramaniam (1997) Neuropeptide Y family of hormones: Receptor subtypes and antagonists, *Peptides* **18**, 445–457.
26. R. Barker (1996) Tachykinins, neurotrophins and neurodegenerative diseases: A critical review on the possible role of tachykinins in the aetiology of CNS diseases, *Rev. Neurosci.* **7**, 187–214.
27. K. Bedecs, M. Berthold, and T. Bartfai (1995) Galanin—10 years with a neuroendocrine peptide, *Int. J. Biochem. Cell Biol.* **27**, 337–349.
28. K. M. Braas and V. May (1996) Pituitary adenylate cyclase-activating polypeptides, PACAP-38 and PACAP-27, regulation of sympathetic neuron catecholamine, and neuropeptide Y expression through activation of type I PACAP/VIP receptor isoforms, *Ann. N. Y. Acad. Sci.* **805**, 204–216.

29. T. Darland, M. M. Heinricher, and D. K. Grandy (1998) Orphanin FQ/nociceptin: A role in pain and analgesia, but so much more. *Trends Neurosci.* **21**, 215–221.
30. R. Fischer-Colbrie, A. Laslop, and R. Kirchmair (1995) Secretogranin II: Molecular properties, regulation and processing to the neuropeptide secretoneurin, *Prog. Neurobiol.* **46**, 49–70.
31. D. R. Gehlert (1998) Multiple receptors for the pancreatic polypeptide (PP-fold) family: Physiological implications. *Proc. Soc. Exp. Biol. Med.* **218**, 7–22.
32. T. Hökfelt, and V. Mutt (1997) "Neuropeptides, In" *Encyclopedia of Neuroscience* 2nd ed., CD-ROM Version (G. Adelman and B. Smith, eds.), Elsevier, Amsterdam, Netherlands.
33. T. Hökfelt, X. Zhang, Z. Q. Xu, R. R. Ji, T. Shi, J. Corness, N. Kerekes, M. Landry, M. Rhyder, J. Kopp, K. Holmberg, and C. Broberger (1997) "The ups and downs of neuropeptides". In *Neuroendocrinology. Retrospect and Perspectives* (H.-W. Korf and K.-H. Usadel, eds.), Springer-Verlag, Berlin, Germany, pp. 5–23.
34. C. H. Hoyle (1998) Neuropeptide families: evolutionary perspectives, *Regul. Peptides* **73**, 1–33.
35. T. Jolas and G. K. Aghajanian (1997) Neurotensin and the serotonergic system, *Prog. Neurobiol.* **52**, 455–468.
36. L. Klimaschewski (1997) VIP—a “very important peptide” in the sympathetic nervous system? *Anat. Embryol.* **196**, 269–277.
37. J. C. Meunier (1997) Nociceptin/orphanin FQ and the opioid receptor-like ORL1 receptor, *Eur. J. Pharmacol.* **340**, 1–15.
38. G. A. Olson, R. D. Olson, and A. J. Kastin (1997) Endogenous opiates: 1996. *Peptides* **18**, 1651–1688.
39. W. G. Rossmanith, D. K. Clifton, and R. A. Steiner (1996) Galanin gene expression in hypothalamic GnRH-containing neurons of the rat: A model for autocrine regulation, *Horm. Metab. Res.* **28**, 257–266.
40. W. H. Rostene and M. J. Alexander (1997) Neurotensin and neuroendocrine regulation, *Frontiers Neuroendocrinol.* **18**, 115–173.
41. M. Spina, E. Merlo-Pich, R. K. W. Chan, A. M. Basso, J. Rivier, W. Vale, and G. F. Koob (1996) Appetite-suppressing effect of urocortin, a CRF-related neuropeptide, *Science* **273**, 1561–1564.
42. T. Takeda and M. Kohno (1995) Brain natriuretic peptide in hypertension, *Hypertension Res.* **18**, 259–266.
43. H. J. van Megen, H. G. Westenberg, J. A. den Boer, and R. S. Kahn (1996) Cholecystokinin in anxiety, *Eur. Neuropsychopharmacol.* **6**, 263–280.
44. J. Vaughan, C. Donaldson, J. Bittencourt, M. H. Perrin, K. Lewis, S. Sutton, R. Chan, A. V. Turnbull, D. Lovejoy, C. Rivier, J. Rivier, P. E. Sawchenko, and W. Vale (1995) Urocortin, a mammalian neuropeptide related to fish urotensin I and to corticotropin-releasing factor, *Nature* **378**, 287–292.
45. S. J. Wimalawansa (1997) Amylin, calcitonin gene-related peptide, calcitonin and adrenomedullin: A peptide superfamily. *Crit. Rev. Neurobiol.* **11**, 167–239.

## Neurotoxins

Given the fundamental importance for survival of animals having a fully functional neuromuscular apparatus, particularly in the wilderness, it is no surprise that most [toxins](#) are specific for components of nerve and muscle cells. Some of these are ion channel-specific toxins ([1-4](#)) (see



[Alpha-Bungarotoxin and Curare-Mimetic Toxins](#), [Conotoxins](#), and [Dendrotoxins](#)). Here, only toxins affecting neuroexocytosis will be considered.

Many snakes and some insects produce neurotoxins with secondary and tertiary structures closely similar to those of the pancreatic [phospholipase A2 \(PLA<sub>2</sub>\)](#). It appears that during the course of [evolution](#), the PLA<sub>2</sub> structural module has been remodeled by some animals to produce highly toxic neurotoxins (5). Some neurotoxic PLA<sub>2</sub>, such as the *Bungarus multicinctus* β-bungarotoxin and the rattlesnake crototoxins, are coupled to other proteins that are responsible for specific binding to the presynaptic membrane; in other toxins, such as notexin, neurospecificity is determined by the PLA<sub>2</sub> molecule itself. PLA<sub>2</sub> neurotoxins induce a transient release of neurotransmitter, followed by a persistent blockade of neuroexocytosis. Their molecular mechanism of action has so far escaped intense investigation, but it appears that the PLA<sub>2</sub> activity is required only in the inhibitory phase.

An opposite synaptic effect is caused by intoxication with the α-latrotoxin, or the insect-specific latroinsectotoxin, which is produced by black widow spiders of the genus *Latrodectus*. These excitatory neurotoxins are very large proteins of more than 1,400 residues that induce a massive neurotransmitter release from synaptic terminals of all types of neurons. Fusion of the small synaptic vesicles with the presynaptic membrane is not followed by [endocytosis](#), and the poisoned synaptic boutons are enlarged. Functional latrotoxin receptors have not been identified, but the very large **glycoproteins** known as neurexins appear to play a role in toxin binding (7).

Anaerobic bacteria of the genus *Clostridium* produce large neurotoxins that are the cause of the neuroparalytic syndromes of tetanus and botulism. These diseases are caused, respectively, by *tetanus* neurotoxin (TeNT) and by one of eight *botulinum* neurotoxins (BoNT, from A to G) (8). The BoNT block the release of acetylcholine at the neuromuscular junction and cause a flaccid paralysis, whereas TeNT acts on the inhibitory interneurons of the spinal cord, whose impairment results in a spastic paralysis. These toxins act in the neuronal cytosol via a zinc-[endopeptidase](#) activity that is specific for three protein components of the neuroexocytosis apparatus (8).

The clostridial neurotoxins are produced as a single inactive polypeptide chain of 150 kDa and released upon bacterial lysis. BoNT are complexed with nontoxic accessory proteins, which protect them during their transit through the stomach, and the complexes dissociate at the neutral pH of the intestine (8). Active two-chain neurotoxins, composed of disulfide-linked heavy (H, 100 kDa) and light (L, 50 kDa) chains, are generated by proteolytic cleavage at an exposed loop. TeNT and BoNT each consist of three functionally distinct domains, which play different roles in cell intoxication (8). Domain L is the catalytic part, whereas the 50-kDa amino-terminal half of the H chain governs cell penetration, and the 50-kDa carboxyl-terminal half of the H chain is responsible for the absolute neurospecificity of these toxins. Together with their intracellular enzymatic activity, this accounts for the fact that they are the most potent toxins known (mouse LD<sub>50</sub> ranging between 0.1 and 1 ng/kg).

No receptors for TeNT and BoNT are known, but there is evidence that polysialogangliosides play a role in their fixation at nerve terminals and the amino-terminus of synaptotagmin is implicated in BoNT/B binding (8).

After binding to peripheral motoneurons, TeNT and BoNT are internalized inside vesicles of unknown nature, in a temperature- and energy-dependent process (9, 10). TeNT penetrates central nervous system neurons by binding to the luminal side of synaptic vesicles during neurotransmitter release, which is followed by vesicle endocytosis (11). Refilling of the vesicle with neurotransmitter is coupled to acidification of the lumen by the vacuolar **ATPase proton pump**. Similar to [diphtheria toxin](#) (DT), this causes a conformational transition of TeNT to an acidic conformation that mediates the translocation of its L chain into the cytosol (6, 11, 12).

The L chains of TeNT and BoNT/B, /D, /F, and /G specifically recognize and cleave

VAMP/syntaxin, whereas BoNT/A and /E hydrolyze SNAP-25, and BoNT/C cleaves both syntaxin and SNAP-25 (8). This indicates that VAMP, SNAP-25, and syntaxin play key role(s) in exocytosis. With the exception of TeNT and BoNT/B, each of the different clostridial neurotoxins catalyzes the hydrolysis of different peptide bonds (8). As in the case of DT, it is likely that one molecule of toxin is capable of cleaving all target molecules present in the intoxicated cell. At variance from DT, but similar to **cholera toxin** and **pertussis toxin**, the clostridial neurotoxins do not cause neuronal death: They only paralyze the intoxicated synapse, and the poisoned animal may die because of the impairment of neuronal and muscular functions. If the animal survives, the synapse recovers its normal function within a few months. Such reversibility is at the basis of the clinical use of BoNT, which are the best available therapeutic agents for several diseases that benefit from a functional paralysis of selected neuromuscular junctions (13).

## Bibliography

1. P. N. Strong (1990) *Pharmacol. Ther.* **46**, 137–162.
2. A. L. Harvey, ed. (1991) *Snake Toxins*, Pergamon Press, New York.
3. B. M. Olivera et al. (1990) *Science* **249**, 257–263.
4. R. Rappuoli and C. Montecucco (1997) *Guidebook to Protein Toxins and Their Use in Cell Biology*, Sambrook and Tooze, Oxford University Press, Oxford.
5. R. M. Kini, ed. (1997) *Venom Phospholipase A2 Enzymes: Structure, Function and Mechanism*, Wiley, New York.
6. A. G. Petrenko (1993) *FEBS Lett.* **325**, 81–85.
7. C. Montecucco and G. Schiavo (1995) *Quart. Rev. Biophys.* **28**, 423–472.
8. M. E. Schwab, K. Suda, and H. Thoenen (1979) *J. Cell Biol.* **82**, 798–810.
9. J. O. Dolly, R. S. Black, R. S. Williams, and J. Melling (1984) *Nature* **307**, 457–460.
10. M. Matteoli et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13310–13315.
11. G. Menestrina, G. Schiavo, and C. Montecucco (1994) *Mol. Aspects Med.* **15**, 81–193.
12. P. Boquet and E. Dufloot (1982) *Proc. Natl. Acad. Sci. USA* **79**, 7614–7618.
13. C. Montecucco, G. Schiavo, V. Tugnoli, and D. De Grandis (1996) *Mol. Med. Today* **2**, 418–424.

## Neutral Mutation

can be classified into three groups, depending on their effect on the individual affected: advantageous, neutral, and deleterious. Deleterious mutations are usually eliminated quickly from the population by , namely, negative selection. It is quite possible that an advantageous mutation increases its frequency in a population by natural selection that is positive selection. In the case of a neutral mutation, its fate is determined by mere chance due to random mating (see ). Strictly speaking, both advantageous and deleterious mutations are also affected by genetic drift, unless the effective population size ( $N$ ) is infinitely large. The theory of population genetics shows that any **alleles** behave as if they are neutral when  $|Ns|$  is less than 1, where  $s$  is a selection coefficient (1). The selection coefficient,  $s$ , is defined by formulating the fitness ( $f$ ) by  $f = 1 + s$  (for an explanation of fitness, see ). Therefore, even if an advantageous mutation arises, the mutation can be treated as a neutral mutation when the effective population size is very small.

The neutral theory of molecular evolution contends that a large proportion of molecular mutations

that have become fixed in populations are neutral rather than advantageous. On the other hand, selectionists maintain that most such mutations should be advantageous, and they consider them to be the driving force of evolution. Both, however, agree that the majority of mutations are deleterious. Nevertheless, these mutations do not contribute to because they disappear quickly from the population. Therefore, the core issue of the controversy between selectionists and neutralists is what proportion of mutational changes in amino acid and nucleotide sequences is neutral when all deleterious mutations are excluded.

The controversy over the neutral theory has made great contributions to molecular evolution and population genetics in the following two aspects. Firstly, it provided general recognition that the effect of genetic drift cannot be neglected in the study of the evolutionary dynamics of molecular changes. The other was to accelerate the fusion of molecular evolution and population genetics by the recognition that molecular evolution and genetic **polymorphism** are two facets of the same phenomenon (2).

### Bibliography

“Neutral Mutation” in , Vol. 3, pp. 1602–1603, by T. Gojobori; “Neutral Mutation” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. M. Kimura (1983) *The Neutral Theory of Molecular Evolution*, Cambridge Univ. Press, Cambridge, UK.
2. W. H. Li (1997) *Molecular Evolution*, Sinauer Associates, Sunderland, MA.

## Neutron Diffraction And Scattering

Neutron diffraction and scattering yield structural information on biological molecules and their complexes at multiple levels of resolution. Neutron diffraction and scattering methods generally complement X-ray methods and other structural techniques, frequently offering unique pieces of information that enable one to put together the final jigsaw puzzle picture of how a biological molecule or complex works. Neutron diffraction analysis of crystalline samples is analogous to [X-ray crystallography](#), but at modestly high resolution ( $<3 \text{ \AA}$ ) one can locate precisely individual hydrogen atoms that can be key to enzymatic mechanisms (1) or ligand binding (2). One can also study solvent and its role in [protein stability](#) (3-5). Medium-resolution neutron diffraction (3 to 8  $\text{\AA}$ ) has been used to orient alpha-**helices** (6) and to locate specific groups, such as the retinal chromophore in two-dimensional crystals of the membrane protein [bacteriorhodopsin](#) (7), and to locate [water](#) molecules in DNA **structures** (8). Low-resolution neutron diffraction ( $>8 \text{ \AA}$ ) with [contrast variation](#) gives information on disordered regions in membrane protein crystals, or in large biomolecular complexes such as viruses, and specific examples of this applications are described in [Contrast variation. Small-angle scattering](#) of neutrons can yield information on the overall shapes of biological macromolecules in solution and, when combined with [contrast variation](#), can give information on the shapes and dispositions of individual components in biomolecular complexes.

### 1. Interactions of Neutrons with Matter and Basic Scattering Theory

Neutrons offer a number of advantages as a structural probe in biology. They are less damaging to biomolecules than X-rays. They can be produced in a range of wavelengths from tenths to tens of angstroms and hence are useful for probing many orders of magnitude of dimensions relevant in biological structures: from  $\sim 1$  to  $10^3 \text{ \AA}$ . Neutrons are neutral particles and interact principally with the atomic nuclei in a sample. Their scattering properties depend upon the complex neutron–nucleus

interaction; as a consequence, isotopes of the same element can have very different neutron scattering properties. In addition to the coherent component that can interfere and hence yield structural information the scattering of neutrons by atoms can have a significant incoherent component. The incoherent component is only significant for nonzero spin nuclei (9), and in elastic scattering experiments it gives rise to isotropic scattering that contributes to the background. The incoherent scattering is very large for the hydrogen atom, and consequently for many structural biology applications of elastic neutron scattering it is optimal to minimize the amount of hydrogen in a sample by isotopic substitution with deuterium. Alternatively, investigators have used the incoherent, inelastic scattering from neutrons to probe the dynamics of biomolecules (10). Table 1 lists the coherent scattering amplitudes and incoherent scattering cross sections for the nuclei commonly found in biomolecules.

**Table 1. Coherent Neutron Scattering Lengths,  $b_{\text{coh}}$ , and Incoherent Cross Sections,  $s_{\text{inc}}$ , for Biologically Relevant Nuclei**

| Atom       | Nucleus                | $s_{\text{inc}}$ ( $10^{-24}$ cm <sup>2</sup> ) | $b_{\text{coh}}$ ( $10^{-12}$ cm) |
|------------|------------------------|---|-----------------------------------|
| Hydrogen   | <sup>1</sup> H         | 80  | -0.3742                           |
| Deuterium  | <sup>2</sup> H         | 2   | 0.6671                            |
| Carbon     | <sup>12</sup> C        | 0   | 0.6651                            |
| Nitrogen   | <sup>14</sup> N        | ~0  | 0.940                             |
| Oxygen     | <sup>16</sup> O        | 0   | 0.5804                            |
| Phosphorus | <sup>31</sup> P        | ~0.3  | 0.517                             |
| Sulfur     | Mostly <sup>32</sup> S | ~0  | 0.2847                            |

Neutrons can be considered as plane waves with wavelengths ( $\lambda$ ); when they are scattered by atoms in a molecule, whose positions are designated by the vector  $\mathbf{r}$  from an arbitrary origin, the resultant interference of the elastic, coherent scattering is related to the spatial distribution of those atoms and can be expressed as

$$\sum_i \sum_j b_i b_j e^{-i\mathbf{Q}\cdot\mathbf{r}_{ij}} \quad (1)$$

where  $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$  and the summations are over all atoms.  $\mathbf{Q}$  is the difference between the incident and scattered wave vectors, and its amplitude is  $4\pi(\sin q)/\lambda$ , where  $2q$  is the scattering angle.  $b$  is the neutron scattering amplitude for the atom in units of length. Unlike those for X-rays, neutron scattering amplitudes show little or no dependence on scattering angle, because the dimensions of the scattering nuclei are much smaller than the wavelengths of the incident neutrons. In low-resolution solution scattering experiments, one can describe molecules as continuous density distributions, rather than sums of discrete point atom scatterers. Neutron scattering densities are calculated as the sum of the scattering lengths of atoms within a finite volume element, divided by that volume,  $Sb_i/V$ , and a molecule can be described by a scattering density distribution  $\rho(\mathbf{r})$ . The total coherent [scattering intensity distribution](#),  $I(Q)$ , for a randomly oriented molecule in a solvent with a mean scattering density  $\rho_s$  is

$$I(Q) = \left\langle \left| \int \Delta\rho(\mathbf{r})e^{-i\mathbf{Q}\cdot\mathbf{r}} d\mathbf{r} \right|^2 \right\rangle \quad (2)$$

$$= \int \int \Delta\rho(\mathbf{r}_1)\Delta\rho(\mathbf{r}_2) \frac{\sin Q|\mathbf{r}_1 - \mathbf{r}_2|}{Q|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (3)$$

The brackets  $\langle \rangle$  indicate averaging over all orientations of the particle, and the integration is over the volume of the molecule. The term  $\Delta\rho(\mathbf{r}) = \rho(\mathbf{r}) - \rho_s$  is the contrast between the scattering density of the molecule and the solvent, and [contrast variation](#) experiments involve the deliberate manipulation of this factor to alter the resultant scattering. The total [scattering intensity distribution](#) for a solution of monodisperse molecules will be directly proportional to both the number density and the square of the molecular weight of the particles. Analysis of the [small-angle scattering](#) from biological macromolecules in solution yields structural parameters, such as the [radius of gyration](#),  $R_g$ , molecular weight,  $M$ , and the vector length distribution function,  $P(r)$ , for the macromolecule.  $P(r)$  is sensitive to the symmetry of the scattering particle, and it goes to zero at its maximum linear dimension (see [Small-Angle Scattering](#)).  $P(r)$  is related to the scattering intensity distribution,  $I(Q)$ , by a Fourier transformation:

$$I(Q) = 4\pi \int P(r) \frac{\sin Qr}{Qr} dr \quad (4)$$

$$P(r) = \frac{1}{2\pi^2} \int I(Q) Q \cdot r \sin(Q \cdot r) dQ \quad (5)$$

In the case of samples that are ordered in one, two, or three dimensions, the scattered intensity is *not* spherically averaged, but rather convoluted with the repeating lattice to yield a diffraction pattern that has discrete intensity maxima that can be interpreted in terms of the repeat distances in the sample and sometimes give more detailed structural information about the ordered molecules themselves. For example, one can learn about the superhelical structures formed by protein–DNA complexes ([11-13](#)), or about the structures of [membranes](#) and the proteins they contain ([6, 7](#)). Neutron [crystallography](#) is the use of neutron diffraction data from three-dimensional crystals of macromolecules to solve the structure of the crystallized molecule at high resolution. The crystal diffraction data form a complex pattern whose maxima are indexed according to the crystal lattice indices  $h, k, l$ . The [unit cell](#) diffraction is generally described mathematically in terms of the square of the [structure factor](#)  $F_{hkl}$ , which is simply the ratio of the radiation scattering by any real sample to a point scatterer at the origin, and

$$I(Q_{hkl}) = F_{hkl}^2 = \left| \sum_j b_j e^{-2\pi i(hx_j + ky_j + lz_j)} \right|^2 \quad (6)$$

where  $x, y$ , and  $z$  are the atomic coordinates of each atom in the crystallized molecule, and the summation is over all atoms. One “solves” a crystal structure by using the measured  $I_{hkl}$  values to calculate the amplitudes of the structure factors, to which a phase must be assigned in order to extract, by Fourier transformation, the coordinates  $(x, y, z)$  for each atom or chemical group in the crystallized structure. The correct assignment of phases (which cannot be measured directly) is known as the “[phase problem](#).” Methods for solving the phase problem that are successful for X-ray crystallography, such as molecular [isomorphous replacement](#), are not useful for neutron

[crystallography](#). The reason is that they depend upon crystallizing isomorphous molecules with heavy atom labels that scatter X-rays much more strongly than other atoms in the structure and hence can be located and used to generate an initial set of model phases. Neutron crystallography does not have the equivalent of a “heavy atom,” because the neutron scattering lengths of most atoms lie in a narrow range of values. As a result, neutron crystallography generally is most useful when the X-ray crystal structure is already available and there are outstanding issues regarding for example, the locations of light atoms. Because neutron scattering amplitudes show no systematic dependence on atomic number, light atoms such as hydrogen are readily located in a neutron crystallography experiment.

One of the largest differences in coherent neutron scattering amplitudes for atoms in biological systems is between the isotopes of hydrogen ( $^1\text{H}$  and  $^2\text{H}$ ). Note from Table 1 that the scattering amplitude for  $^1\text{H}$  is negative, resulting from a  $180^\circ$  phase shift between neutrons scattered by  $^1\text{H}$  compared to the other nuclei. Thus, selective deuteration of one component of a complex in solution provides a way of altering the mean neutron scattering density of that component. Furthermore, by changing the deuterium level in the solvent, the neutron scattering contrast of each component is varied. It is these properties that facilitate [contrast variation](#) methods that, when combined with small-angle scattering, provide a means for extracting structural information on the individual components of a complex and their relative dispositions in solution. In combination with low-resolution crystallography, contrast variation can be used to locate components or features that are of too low contrast in the corresponding X-ray crystallography experiment to be distinguished.

## 2. Production of Neutrons

Neutron beams with sufficient intensity for studying the structures of biological molecules require sophisticated technology for their production. There are two classes of sources for high-intensity neutron beams for condensed matter studies operating in the world today: nuclear reactors and accelerators. Nuclear fission reactors produce energetic neutrons in quantities beyond those needed to sustain the chain reaction. The excess neutrons escape and are moderated by various substances surrounding the core of the reactor. The neutrons equilibrated with the moderator substance will have Maxwell–Boltzmann distributions of energies appropriate for the moderator temperature. High-resolution neutron crystallography instruments typically use neutrons moderated at room temperature, yielding a peak in the energy distribution around 1 Å. Small-angle neutron scattering benefits from having longer wavelengths available, and thus “cold sources” (eg, liquid hydrogen) are used to slow the neutrons further to have energy peaks in the region 4 to 10 Å. Recently, accelerator-based spallation neutron sources (14) have evolved to the point where they can be useful for biological studies. These sources use high-energy proton pulses directed onto a target of heavy nuclei. Neutrons are captured by the heavy nuclei, producing an unstable nucleus that decays, giving off pulses of neutrons in a process referred to as spallation. Pulsed neutron sources are very bright, but on average over time deliver fewer neutrons than a steady-state reactor source. Advantages in pulsed source instrumentation can be gained, however, by using time-of-flight methods to determine the wavelengths of neutrons reaching a detector (15), thus enabling the experimenter to use the entire “white” neutron beam. As spallation sources become more powerful, they are becoming more competitive with reactor sources for scattering applications. In spite of technological advances in neutron sources and instrumentation over the past few decades, however, neutron sources have intensities many orders of magnitude lower than even that of a conventional laboratory X-ray source. Some of this disadvantage can be overcome by long exposures of samples in experiments, because neutrons are less-damaging radiation. However, neutrons are generally used only when their unique properties confer an important advantage that gives specific information that cannot otherwise be obtained.

## 3. Examples of Neutron Diffraction and Scattering Applications

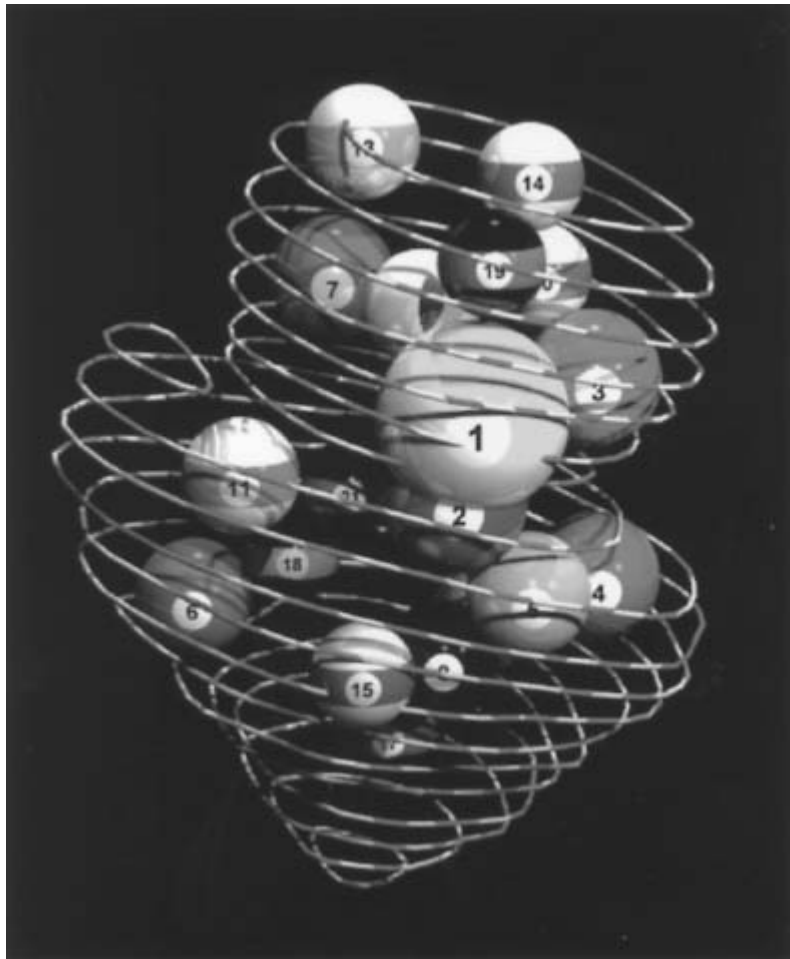
### 3.1. High-Resolution Neutron Diffraction Reveals Locations of Hydrogens in Protein Crystals

High-resolution neutron diffraction studies of three-dimensional crystals of proteins have yielded detailed information on the protonation state of individual chemical groups as well as the solvent structure that plays a role in protein folding and stability. In 1980, Kossiakoff and Spencer (1) reported the neutron crystal structure of bovine **trypsin** covalently inhibited by a **transition-state analogue**, thus revealing the protonation states of the **active-site** residues Asp102 and His57 of the **catalytic triad**. This result resolved the much debated mechanistic issue by showing conclusively that the catalytic base in the transition state of the reaction is His57 and not Asp102. High-resolution neutron diffraction data to 2.1 Å from crystals of trypsin (4) localized two-thirds of the waters of **hydration** expected from thermodynamic data. In comparison, X-ray studies at higher resolution (1.35 Å) located about half the waters of hydration, and a number of these were found in the neutron study to be spurious. In another high-resolution neutron diffraction study of carbonyl myoglobin (5), 87 water and 5 ion molecules were localized, and a number of X-ray-determined waters of hydration were identified as spurious. Neutron crystallography presents a number of difficulties as a result of the relatively low fluxes of neutron sources (compared with X-rays), and the signal-to-noise ratio is low because of the high incoherent backgrounds from hydrogenated proteins. In 1994, Kossiakoff and colleagues (16) reported obtaining 1.5-mm<sup>3</sup> crystals of perdeuterated **Staphylococcal nuclease**. They showed that there are no structural differences between the protiated and deuterated proteins at 1.9-Å resolution. Perdeuteration of proteins is thus a promising tool for gaining signal-to-noise in neutron crystallography experiments.

### 3.2. Subunit Structure of Biomolecular Assemblies: The Ribosome

**Ribosomes** are complex assemblies of proteins and RNA that act as protein factories in all living cells. They are composed of a small and large subunit. The 70S ribosome from *E. coli* is the most extensively studied ribosome, and different subunits have been the target of neutron scattering investigations. Early studies used contrast variation to locate the protein and RNA in both subunits (17-19), and more recent studies have used deuteration, triangulation (20, 22), or triple isotopic substitution (23) methods to give more structural details. In 1987, Moore and co-workers published a complete map of all 21 proteins of the small ribosomal subunit from *E. coli* using triangulation by measuring distances between pairs of proteins within the small subunit (20). In their experiment, the 30S ribosomal subunit was reconstituted from 16S ribosomal RNA and a mixture of purified 30S proteins, one or more of which was deuterated. Using the small-angle neutron scattering data from these samples, they determined the distance between the centers of mass for each pair of deuterated proteins, plus the radius of gyration for each protein. Over a 10-year period, this group obtained 105 distance data sets relating 93 different protein pairs in the 30S subunit. Using triangulation, they constructed a three-dimensional map of the 21 proteins of the 30S subunit (Fig. 1).

**Figure 1.** Representation of the neutron map of the 30S subunit of the *E. coli* ribosome (20). Each protein is represented by a sphere whose volume equals that of the protein. The numbering of the proteins adheres to the standard nomenclature for ribosomal proteins. Courtesy of Malcom Capel.



## Bibliography

1. A. A. Kossiakoff and S. A. Spencer (1980) *Nature* **288**, 414–416.
2. X. Cheng and B. P. Schoenborn (1991) *J. Mol. Biol.* **220**, 381–399.
3. J. P. Bouquiere, J. L. Finney, and H. F. J. Savage (1994) *Acta Crystallogr.* **B50**, 566–578.
4. J. S. Finer-Moore et al. (1992) *Protein Struct. Funct. Genet.* **12**, 203–222.
5. X. Cheng and B. P. Schoenborn (1990) *Acta Crystallogr.* **B46**, 195–208.
6. F. A. Samtey et al. (1994) *J. Mol. Biol.* **236**, 1093–1104.
7. J. S. Jubb, D. L. Worcester, H. L. Crespi, and G. Zaccai (1984) *EMBO J.* **3**, 1455–1461.
8. P. Langan et al. (1992) *J. Biomol. Struct. Dyn.* **10**, 489–503.
9. G. Allen and J. S. Higgins (1973) *Rep. Prog. Phys.* **36**, 1073.
10. J. C. Smith (1991) *Q. Rev. Biophys.* **24**, 227–291.
11. E. M. Bradbury and J. P. Baldwin (1986) *Cell Biophys.* **9**, 35–66.
12. V. Graziano, S. E. Gerchman, D. K. Scheider, and V. Ramakrishnan (1996) In *Basic Life Sciences*, Vol. **64: Neutrons in Biology** (E. H. Y. Chu, ed.), Plenum Press, New York, pp. 127–136.
13. G. A. Olah et al. (1995) *J. Mol. Biol.* **249**, 576–594.
14. C. G. Windsor (1981) *Pulsed Neutron Scattering*, Taylor and Francis, London, UK.
15. C. C. Wilson (1996) In *Basic Life Sciences*, Vol. **64: Neutrons in Biology** (E. H. Y. Chu, ed.), Plenum Press, New York, pp. 35–55.
16. T. R. Gamble, K. R. Clauser, and A. A. Kossiakoff (1994) *Biophys. Chem.* **53**, 15–26.



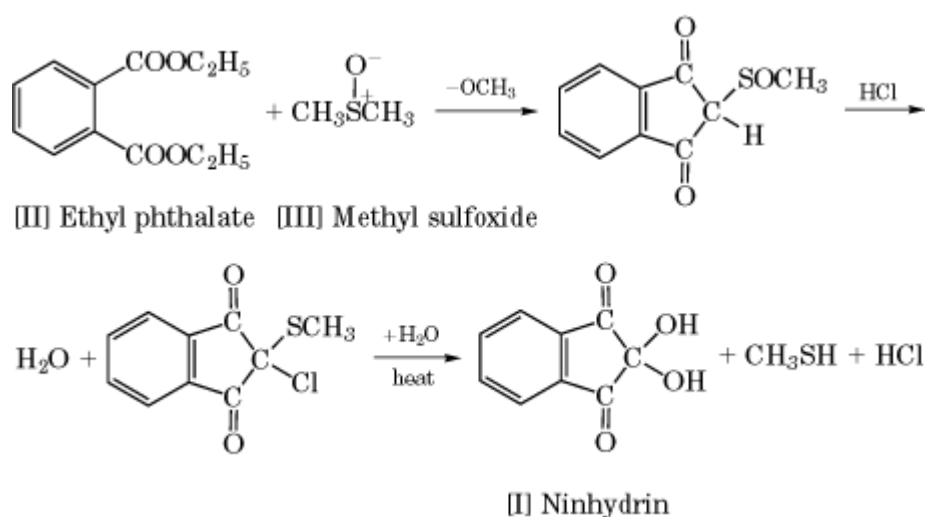
17. P. B. Moore, D. M. Engelman, and B. P. Schoenborn (1975) *J. Mol. Biol.* **91**, 101–120.
18. H. B. Stuhmann et al. (1976) *Proc. Natl. Acad. Sci. USA* **73**, 2379–2383.
19. I. N. Serdyuk and A. K. Grenader (1975) *FEBS Lett.* **59**, 133–136.
20. M. S. Capel et al. (1987) *Science* **238**, 1403–1406.
21. R. P. May, V. Nowotny, P. Nowotny, H. Voß, and K. H. Nierhaus (1992) *EMBO J.* **11**, 373–378.
22. D. H. Harrison, R. P. May, and P. B. Moore (1993) *J. Appl. Crystallogr.* **26**, 198–206.
23. I. N. Serdyuk and G. Zaccai (1996) *J. Mol. Struct.* **383**, 197–200.

### Suggestions for Further Reading

24. G. E. Bacon (1975) *Neutron Diffraction*, Clarendon Press, Oxford, UK.
25. B. Jacrot (1976) *Rep. Prog. Phys.* **39**, 911–953.
26. N. V. Raghavan and A. Wlodawer (1987) *Methods Exp. Phys.* **23**, 335–365.
27. B. P. Schoenborn and R. B. Knott, eds. (1996) *Basic Life Sciences*, Vol. **64**: *Neutrons in Biology* (E. H. Y. Chu, ed.), Plenum Press, New York.
28. H. B. Sturhmann (1987) "Molecular Biology. In" *Methods of Experimental Physics*, Vol. **23**, Part C, Academic Press, New York, pp. 367–403.

## Ninhydrin

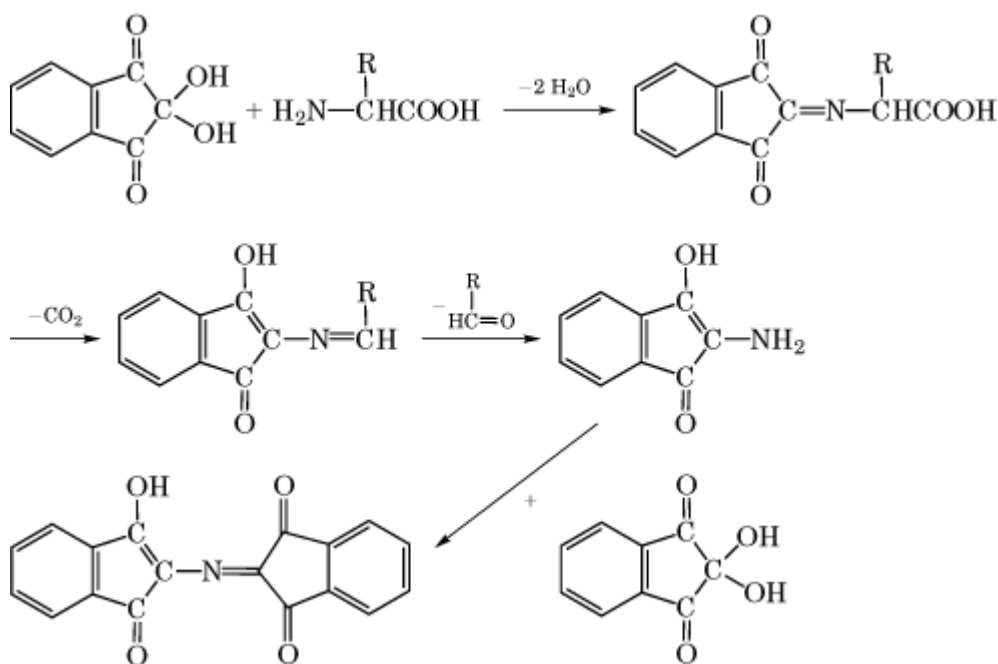
Ninhydrin [I] generates a blue-purple color on reacting with free [amino groups](#) in [amino acids](#), **peptides**, or proteins. It has a molecular weight of 178.14 and occurs as colorless or faint yellow crystals. Ninhydrin, also called 2,2-dihydroxy-1,3-indanedione, is synthesized from ethyl phthalate [II] and methyl sulfoxide [III] (Scheme 1):



### 1. Ninhydrin Reaction

When a solution containing amino acids, peptides, or proteins is boiled with ninhydrin under mildly acidic conditions, an intense blue-purple color with an absorption maximum at 570 nm develops

(Scheme 2):



[Amino acids](#), such as [proline](#) and hydroxyproline, develop color with an absorption maximum at 440 nm.

In the original instruments for automatic [amino acid analysis](#), a protein sample that has been hydrolyzed with 6 N HCl at 110°C for 24 h to amino acids is separated with an ion-exchange **chromatographic** column, and the eluate is monitored at 570 nm after being developed with the ninhydrin reaction. To detect peptides or amino acids on filter paper or other solid supports, the paper is dipped in or sprayed with ninhydrin solution in acetone.

## 2. Protein Detection by Ninhydrin Method

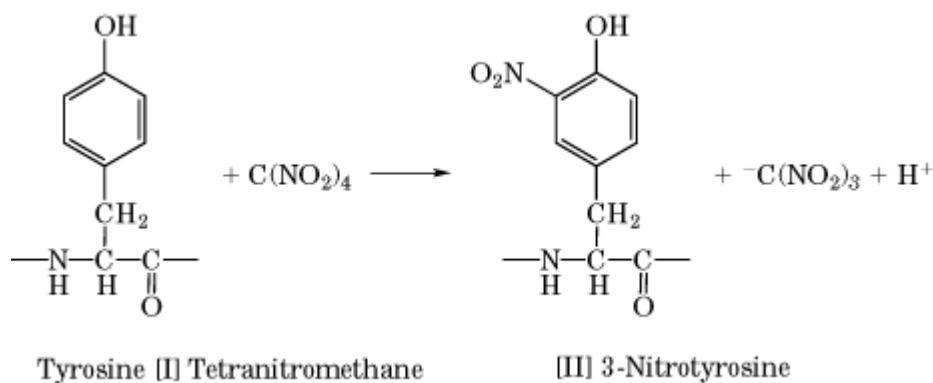
A solution containing 20 to 500 µg protein is dried on a filter paper. Nonprotein compounds that might react with ninhydrin are eluted by washing with a **trichloroacetic acid** solution. Then the proteins are hydrolyzed with Ba(OH)<sub>2</sub> solution at 120°C for 10 min, and the amino acids released are measured by the ninhydrin reaction (1). Because the amount of protein present is determined by the amino acids produced, this is an accurate method suitable for all proteins.

## Bibliography

1. R. McGrath (1972) *Anal. Biochem.* **49**, 95–102.

## Nitration

[Tyrosine](#) residues are nitrated with [tetranitromethane](#) (TNM) [I] to form 3-nitrotyrosine residues [II] under mild conditions (Scheme 1) (1):3-Nitrotyrosine is colored, with an **absorbance** maximum at 428 nm at alkaline pH, and can be utilized as a [reporter group](#). TNM also oxidizes [thiol groups](#) at low pH. The nitration reaction often causes protein polymerization, especially at high protein concentration. The modified nitrotyrosine residue is reduced with sodium hydrosulfite to an aminotyrosine residue, which then can be subjected to further modifications.



## Bibliography

1. T. Imoto and H. Yamada (1989) In *Protein Structure: A Practical Approach* (T. E. Creighton, ed.), IRL Press, Oxford, UK, pp. 247–277.

## Nitric Oxide

Nitric oxide (NO, nitrogen monoxide) is one of the smallest molecules to be synthesized in biological systems. Carrying an unpaired electron, NO is a short-lived free radical and is highly reactive in biological environments. Nitric oxide was previously well known for its roles as an air pollutant that contributes to the formation of photochemical smog, as a participant in bacterial [nitrogen fixation](#), and as a probe to study the metal-binding sites of [metalloproteins](#). Not until 1987 did scientists discover that NO is actually produced in many organs and tissues of mammals. Even more astonishing, this molecule has remarkably diverse functions, ranging from neurotransmitter and vasodilator to microbicidal mediator. It has earned a reputation as a double-edged sword, responsible for both beneficial and deleterious effects (1). Depending on the site, duration and amount of NO produced, NO can have either cytoprotective or cytotoxic effects. How can such a small and simple inorganic molecule exhibit such diverse functions? How is NO production regulated in different cell types? How can we apply our knowledge to control diseases where NO plays significant roles? Seeking answers to such questions has been the focus of intensive research for the past ten years. Our current understanding of regulation of NO biosynthesis and the roles of NO in mediating biological processes under physiological and pathophysiological conditions are the themes of this article.

## 1. Nitric oxide

The biological roles of NO were initially discovered through several lines of investigation in the fields of toxicology, immunology, cardiovascular pharmacology, and neurobiology. As early as 1916, Mitchell et al. observed that the amount of nitrate ( $\text{NO}_3^-$ ) in the urine of rats, pigs, and humans exceeded the amount present in their diets and suggested that mammals produce nitrogen oxides (2). This study was largely ignored until the late 1970s and early 1980s, when Tannenbaum and coworkers (3) demonstrated that nitrogen oxides are important products of mammalian metabolism and that infection caused a greater than tenfold increase in plasma nitrate and nitrite ( $\text{NO}_2^-$ ) levels (4, 5). By 1985, Stuehr and Marletta (6) identified **macrophages** as a source for nitrate and nitrite production. This was based on the observation that macrophages isolated from mice that were previously injected with infectious agent produced high amounts of nitrate and nitrite and that *Escherichia coli*-derived **endotoxin** induced nitrate and nitrite synthesis by cultured macrophages (6). This *in vitro* finding led to the observation that nitrogen oxides ( $\text{NO}_x$ ) are synthesized enzymatically from the amino acid L-arginine, with L-citrulline as coproduct (7). Independently, Hibbs et al. discovered that conversion of L-arginine to nitrogen oxides was required for macrophage-mediated cytotoxicity toward tumor cells (8) and that the synthesis of reactive nitrogen oxide intermediates from L-arginine, as well as the cytotoxicity of macrophages, could be blocked by a series of  $\text{N}_w$ -substituted L-arginine analogs (9). Addition of nitrite and nitrate to cultured tumor cells did not duplicate the macrophage's ability to kill, whereas NO could. These studies suggested an important role (or roles) for NO in the immune defense system (10-12). It is now known that NO is indeed a macrophage-derived cytostatic agent.

In 1980, Furchgott and Zawadzki reported that the endothelium (the innermost cell layer of blood vessel walls) releases a liable substance that diffuses into the underlying smooth muscle layer, activates soluble **guanylate cyclase**, and thus causes vascular smooth muscle relaxation (13). This liable agent was termed *endothelium-derived relaxation factor* (EDRF). In 1987, Moncada and coworkers (14) as well as Ignarro and coworkers (15), independently demonstrated that NO was produced by the endothelium and was indeed the EDRF (16, 17). These initial important discoveries in the immune and cardiovascular systems were further extended by Garthwaite and coworkers, who identified an EDRF-like activity in the brain after stimulation with N-methyl-D-aspartic acid (18). Taken together, these discoveries hinted at the wide distribution of the L-arginine:NO synthesis pathway in mammals. Studies to date have revealed that nearly every organ and tissue can produce NO.

Although NO is a simple diatomic molecule, the complexity of its **redox** chemistry in biological systems may account for its diverse functions (19, 20). Nitric oxide is a gaseous molecule with an unpaired electron ( $\cdot\text{N} = \text{O}$ ). It tends to react rapidly with other atoms or molecules that also contain unpaired electrons, for example, with oxygen ( $\text{O}_2$ ), superoxide ( $\text{O}_2^-$ ) and transition metals. Unlike most known **signal transduction** transmitters, NO is unique in being **nonpolar** and able to diffuse freely across cell **membranes**. The reaction between NO and superoxide is a **diffusion-controlled reaction**, and it generates peroxynitrite ( $\text{ONOO}^-$ ). Peroxynitrite is a powerful oxidant that can oxidize a variety of biological molecules (eg, **nitration** of **tyrosine** residues) and may therefore be responsible for certain types of NO-mediated toxicity. On the other hand, NO can effectively scavenge oxygen radicals and thus may protect against cellular damage and cytotoxicity caused by reactive oxygen radicals (21). The ultimate products of NO reacting with oxygen are nitrite and nitrate, which are largely unreactive anions and represent the biological fate of NO.

In biological environments, NO reacts primarily with transition metals, **thiol groups**, and with nitrogen and oxygen radicals generated under certain conditions. Metal- and thiol-containing proteins are well-described targets for NO. The interaction between NO and prosthetic iron groups or thiol groups leads to the formation of adducts that either activate or inactivate the target protein. These represent the major mechanisms by which NO exerts many of its biological functions. For example, binding of NO to the heme iron of soluble guanylate cyclase results in activation of the

enzyme and, thus, elevation in intracellular levels of the [second messenger](#), [cyclic GMP](#). It is through the activation of soluble guanylate cyclase that NO mediates vasorelaxation, inhibition of platelet aggregation, and neurotransmission in the central nervous system (22). Another heme-containing enzyme that is activated by NO is cyclooxygenase involved in eicosanoid synthesis (23). However, this NO-dependent activation is thought to result from S-nitrosylation of a thiol group, not from interaction with the heme moiety [reviewed by Salvemini and Masferrer (24)]. Another critical target of NO in cell signaling is p21<sup>ras</sup>, a monomeric G protein family member that plays essential roles in converting extracellular signals into intracellular biochemical events (25); see *ras* proteins). S-nitrosylation of a single [cysteine](#) residue by NO results in activation of p21<sup>ras</sup> and thus downstream signaling (26). Nitric oxide-mediated nitrosylation of non-heme iron, heme iron, or critical thiol groups of enzymes can also result in inhibition of enzyme activity. Examples of enzymes belonging to this category include aconitase, NADH:ubiquinone oxidoreductase, succinate:ubiquinone oxidoreductase, [cytochrome P450](#), glyceraldehyde-3-phosphate dehydrogenase, [ribonucleotide reductase](#), NADPH oxidase (22, 27, 28) and [caspases](#) involved in **cytokine** maturation and [apoptosis](#) (29-31); see Table 1).

**Table 1. Biological Targets of Nitric Oxide**

| <b>Target Site</b>              | <b>Examples</b>   |
|---------------------------------|---|
| Heme-containing proteins        | Soluble <b>guanylyl cyclase</b><br><br>Cyclooxygenase<br><br><b>Hemoglobin/myoglobin</b><br><br>NO synthase<br><br><a href="#">Cytochrome P450</a>  |
| Fe-S cluster containing enzymes | <i>cis</i> -Aconitase/IRE-BP<br><br>NADH: ubiquinone oxidoreductase<br><br>Succinate: ubiquinone oxidoreductase   |
| Thiol groups                    | <a href="#">Ribonucleotide reductase</a><br><a href="#">Caspases</a><br><br><a href="#">Glutathione</a><br><br>Glyceraldehyde-3-phosphate dehydrogenase<br><br><a href="#">Plasminogen</a> activating factor<br><br><a href="#">Adenylate cyclase</a> |

|                  |  |
|------------------|--|
|                  | NADPH oxidase                                    |
|                  | p21 <sup>ras</sup>                               |
|                  | Albumin  |
|                  | N-methyl-D-aspartate receptor                    |
| Free radicals    | Protein kinase C<br>Superoxide<br>Lipid radicals |
| DNA <sup>a</sup> |  |

---

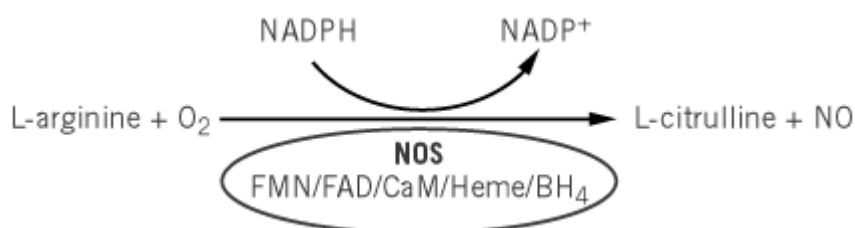
<sup>a</sup> Via deamination.

Although typically short-lived because of its reactivity with cellular targets and oxyhemoglobin of red blood cells, NO probably can also exert actions remote from the site of synthesis. Nitric Oxide has been shown to form stable adducts. Relatively long-lived nitrosothiol adducts can occur on albumin, glutathione, and hemoglobin, which all still possess the biological functions of NO, probably by its release (32-35).

## 2. Nitric oxide synthase

Our understanding of the regulation and function of NO have greatly advanced through studies on nitric oxide synthase (NOS, EC1.14.13.39), the enzyme responsible for NO synthesis. Nitric oxide is synthesized via the oxidation of one of the two chemically equivalent guanidino nitrogen atoms of L-arginine, with L-citrulline formed as the coproduct (Fig. 1). Oxidation of arginine by NOS occurs via two successive monooxygenation reactions, producing an intermediate N<sup>ω</sup>-hydroxyl-L-arginine. This is a five-electron oxidation reaction, with NADPH as the electron donor. Five bound **cofactors** and **prosthetic groups** are required for NOS activity: flavin adenine dinucleotide (FAD), flavin mononucleotide (FMN), iron-protoporphyrin IX (heme), tetrahydrobiopterin (BH<sub>4</sub>), and **calmodulin** (CaM). There are at least three NOS **isoforms**; each represents a distinct gene product and exhibits unique expression patterns and mechanisms of regulation (Table 2). The constitutive isoforms, neuronal NOS (nNOS, ncNOS, or NOS1) and endothelial NOS (eNOS, ecNOS, or NOS3), are expressed constitutively and require binding of Ca<sup>2+</sup>-CaM for catalytic function. Binding of the Ca<sup>2+</sup>-CaM complex to cNOSs is triggered by transient increases in intracellular Ca<sup>2+</sup> levels. In contrast, CaM binds tightly to the inducible NOS isoform (iNOS or NOS2) even at basal Ca<sup>2+</sup> level and thus renders iNOS activity independent of fluctuation of Ca<sup>2+</sup> levels (36). In response to Ca<sup>2+</sup> transients, cNOS intermittently produces small amounts of NO (NO “puffs”) that participate in physiological signaling processes such as vasorelaxation and neurotransmission; and on induction by endotoxin or cytokines (immunologic or inflammatory stimuli), iNOS produces large and sustained amounts of NO that exhibit not only signaling functions but also antimicrobial activity and cytotoxic effects. The roles of NO under physiological and pathophysiological conditions are typically studied via pharmacological and genetic approaches (Table 3).

**Figure 1.** The reaction catalyzed by nitric oxide synthase.



**Table 2. Characteristics of Human Nitric Oxide Synthase Isoforms**

| Isoform                                 | ncNOS   | iNOS  | ecNOS   |
|---|---|---|---|
| Predicted polypeptide size              | 61 kD/1434aa                                    | 131 kD/1153aa                                   | 133 kD/1203aa                                   |
| Cofactor/prothetic groups               | FAD, FMN, BH <sub>4</sub> , P450-like heme, CaM | FAD, FMN, BH <sub>4</sub> , P450-like heme, CaM | FAD, FMN, BH <sub>4</sub> , P450-like heme, CaM |
| Primary cellular location               | Cytosol   | Cytosol   | Golgi, caveolae                                 |
| Active enzyme                           | Dimer   | Dimer   | Dimer   |
| Ca <sup>2+</sup> dependency             | Yes   | No  | Yes   |
| Gene size                               | >200 kb   | ≈37 kb  | ≈21 kb  |
| Gene loci                               | 12q24.2   | 17q11.2–q12                                     | 7q3.5–q3.6                                      |
| Number of exon/intron                   | 29 / 28   | 26 / 25   | 26 / 25   |
| mRNA size                               | 10.0 kb   | 4.5 kb  | 4.7 kb  |
| Translation initiation/termination site | Exon 2 / 28                                     | Exon 2 / 26                                     | Exon 1 / 26                                     |
| Predominant regulation mechanism        | Ca <sup>2+</sup> transients                     | Transcription                                   | Ca <sup>2+</sup> transients <sup>a</sup>        |

<sup>a</sup> Ca<sup>2+</sup>–independent activation also occurs.

**Table 3. Tools to Study Nitric Oxide**

| Approach  | Pharmacological | Example | Genetic                               |
|-----------|-----------------|---------|---------------------------------------|
| Remove NO | Inhibitors:     |         | Knockout mice ( <a href="#">106</a> ) |

|           |              |   |                                       |
|-----------|--------------|---|---------------------------------------|
|           | nonselective | N <sup>W</sup> -substituted L-Arg analogs | ncNOS                                 |
|           | selective    | NIL <sup>a</sup>                          | iNOS<br>ecNOS                         |
|           | Scavengers   | Oxyhemoglobin                             | Antisense oligonucleotides            |
| Supply NO | NO donors    | SNAP, nitroprusside<br>nitroglycerin      | Gene transfer ( <a href="#">107</a> ) |

---

<sup>a</sup> NIL: *N*-iminoethyl lysine; SNAP: S-nitroso-*N*-acetylpenicillamine.

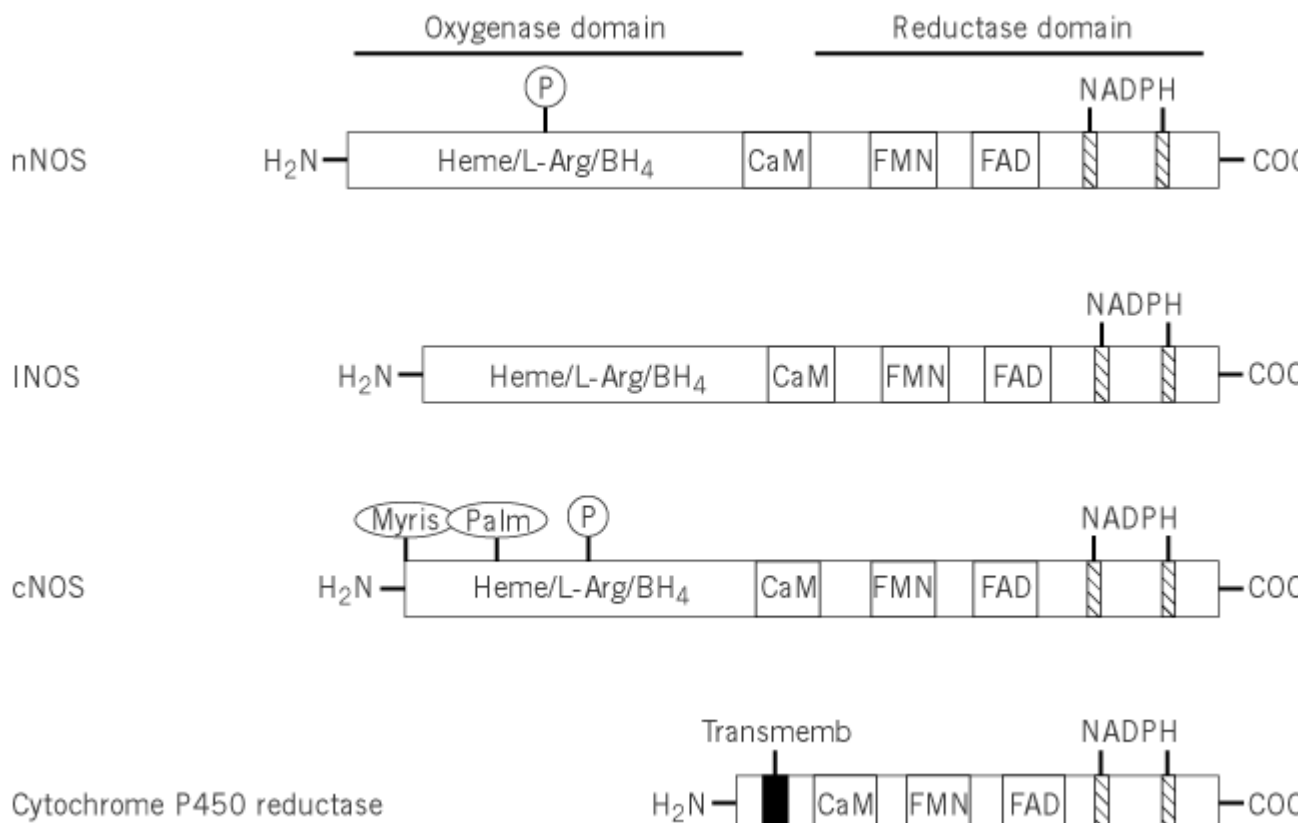
The three NOS isoforms share 50% to 60% amino acid sequence identity with each other, whereas the same NOS isoform from different species has [homology](#) ranging from 80% to 96% identity [eg, 80% sequence identity between murine macrophage iNOS and human hepatocyte iNOS at both nucleotide and amino acid levels ([37](#), [38](#)).

### 2.1. Shared Characteristics

Molecular [cloning](#) of **cDNA** of each NOS isoform has greatly advanced our understanding of the structure, functional **domains**, and potential regulatory sites of NOS. Analysis of amino acid sequences derived from NOS cDNA revealed consensus sequences for binding of NADPH, FAD, FMN, and calmodulin (Fig. [2](#)). The C-terminal half of each NOS shows significant homology to NADPH-dependent cytochrome P450 reductase ([39](#), [40](#)). The active form of all NOS isoforms appears to be a homodimer. Each subunit is composed of an N-terminal oxygenase domain and a C-terminal reductase domain, connected by a CaM-binding region (reviewed in [41](#)). Experiments utilizing limited **proteolysis** or **gene expression** of individual domains have confirmed the location of the heme, BH<sub>4</sub>, and L-arginine binding sites within the oxygenase domain and of those for FAD, FMN, and NADPH within the reductase domain. The flavins appear to transfer electrons from NADPH to the catalytic heme group, which is P450-like ([42](#)). The CaM-binding region may facilitate **electron transfer** between the reductase and oxygenase domains. Dimerization of two subunits is through the oxygenase domain and dependent on BH<sub>4</sub>, L-arginine, and heme ([43-47](#)). The three-dimensional structure of the iNOS oxygenase domain complexed with inhibitor imidazole or aminoguanidine has been most recently unraveled, allowing the first glimpse at the NOS [active site](#) ([48](#)). The structure provides a basis for designing highly selective drugs against iNOS, which may overproduce NO in a number of diseases, such as septic shock, inflammatory bowel disease, rheumatoid arthritis, multiple sclerosis, and stroke.

**Figure 2.** Schematic comparison of various nitric oxide synthase isoforms and cytochrome P450 reductase. P, phosphorylation site; Myris, myristylation site of ecNOS; Palm, palmitoylation site of ecNOS.





## 2.2. ecNOS

The distribution of endothelial cNOS is relatively limited to vascular endothelial cells, neurons and some epithelial cell types. Nitric oxide derived from vascular ecNOS exhibits many functions, including maintenance of basal vascular tone, mediation of blood vessel relaxation, regulation of blood flow and prevention of platelet adhesion and aggregation.

The ecNOS gene has been cloned from bovine (49, 50) and human (51-53) endothelial cells. The human ecNOS gene is present as a single copy in the haploid human genome, encoded by 26 exons spanning ~21 kb and located on chromosome 7. Analysis of the human ecNOS promoter reveals the presence of elements that are found in several constitutively expressed genes, namely, the SP1 and GATA motifs (51, 54). Identification of a shear stress response element in the TATA-less 5'-flanking region of the human ecNOS gene is consistent with the finding that shear stress produced by the flow of blood increases endothelial-dependent NO release and ecNOS gene expression. Hypoxia also up-regulates transcription of ecNOS (55), and chronic exercise in dogs has been shown to increase coronary artery NO production and specifically ecNOS gene expression (56). Conversely, tumor necrosis factor  $\alpha$  (TNF $\alpha$ ) down-regulates ecNOS by decreasing the stability of ecNOS mRNA (57).

ecNOS undergoes several post-translational modifications. N-myristylation (58) and reversible palmitoylation (59) are unique to ecNOS among the NOS family and account for the membrane targeting of ecNOS to the Golgi complex and plasmalemmal caveolae (small invaginations in plasma membrane that may serve as sites for signal processing) (60-63). Association of ecNOS with the Golgi region has been shown to be necessary for ecNOS to respond to intracellular signals and to produce NO in intact cells efficiently (64). Other reported post-translational modifications of ecNOS include phosphorylation (65).

ecNOS activity is regulated by the intracellular Ca<sup>2+</sup> levels (see Calcium Signaling). Agonists that

elevate the intracellular  $\text{Ca}^{2+}$  level, by either influx of extracellular  $\text{Ca}^{2+}$  following receptor activation or release of  $\text{Ca}^{2+}$  from intracellular stores, can increase NO synthesis and, hence, the NO-dependent vasculature relaxation. These agonists include acetylcholine, bradykinin, serotonin, and norepinephrine. Apart from the well-established intracellular  $\text{Ca}^{2+}$  transient dependency, accumulating evidence demonstrates that eNOS can also produce NO via a  $\text{Ca}^{2+}$ -independent pathway in response to fluid shear stress and isometric stretching (reviewed in [66](#)).

### 2.3. ncNOS

Neuronal cNOS is expressed in neurons of the central nervous system, particularly of the cerebellum and in the peripheral nervous system. In addition to the nervous system, the enzyme and/or its [Messenger RNA](#) are also detected in human skeletal muscle; spinal cord; epithelial cells of lungs, uterus, and stomach; pancreatic islet cells; male sex organ; pituitary gland; and kidney macula densa cells. In the brain, on activation of *N*-methyl- D-aspartate (NMDA) receptor by a stimulatory amino acid, such as glutamate, released from adjacent neuron (presynaptic terminus), a channel in the receptor opens, resulting in an influx of  $\text{Ca}^{2+}$  to the postsynaptic neuron, where  $\text{Ca}^{2+}$  binds to calmodulin. The  $\text{Ca}^{2+}$ -CaM complex then activates ncNOS, and NO is produced from L-arginine. Nitric oxide may then exert effects through activating soluble guanylate cyclase and, subsequently, by the elevation of the cyclic GMP content. Also, NO may diffuse back to the presynaptic terminus. A role for NO in the formation of long-term memory has been proposed. High levels of NO production following NMDA receptor stimulation may contribute to the neurotoxicity seen in strokes. In periphery nerves (referred to as nonadrenergic noncholinergic nerves), NO produced from ncNOS has been shown to be a regulator of intestinal motility and penile erection.

The ncNOS cDNA has been cloned and characterized from rat ([39](#)) mouse([67](#)), and human ([68](#)). The human ncNOS gene consists of 29 exons and 28 introns that span over 200 kb of genomic DNA and is present as a single copy in the [haploid](#) human genome ([69](#)). One characteristic of ncNOS is the remarkable diversity of the 5'-flanking regions of its mRNA isolated from different human tissues (reviewed in [70](#)). This diversity results from using different first exons, followed by alternate splicing to a common exon 2. The encoded polypeptide chain remains the same because the [translation](#) initiation site is located in exon 2. Although ncNOS is constitutively expressed, the various mRNA species are subject to different transcriptional regulation that is development- and tissue-specific, implying that other roles of NO produced by ncNOS will continue to unfold.

As with the eNOS, generation of NO by ncNOS is triggered by  $\text{Ca}^{2+}$  transients. Binding of  $\text{Ca}^{2+}$ -CaM to ncNOS acts as a switch that facilitates electron flow from the donor NADPH via flavins to the catalytic heme. The primary sequence of ncNOS also contains a consensus region (Lys-Arg-Phe-Gly-Ser) for phosphorylation by cyclic AMP-dependent kinase (protein kinase A, PKA). Although ncNOS can be phosphorylated *in vitro* by PKA, as well as by protein kinase C and CaM-dependent kinase, the lack of convincing evidence for an effect of phosphorylation on ncNOS activity raises questions about the physiological significance of this post-translational modification.

### 2.4. iNOS

Essentially every tissue and nucleated cell type has the capacity to express functional iNOS when stimulated with endotoxin lipopolysaccharide (LPS) or proinflammatory cytokines, such as TNF $\alpha$ , **interleukin-1** (IL-1), [interferon](#)  $\gamma$  (IFN- $\gamma$ ), or combinations thereof. These stimuli often act synergistically to induce iNOS expression. Hypoxia can also up-regulate iNOS expression ([71](#)). Cell types that have been shown to express iNOS on stimulation include, but are not limited to, macrophages, monocytes, hepatocytes, lymphocytes, chondrocytes, vascular smooth muscle cells, vascular endothelial cells, pancreatic islets, neutrophils, and glial cells. Under normal physiological conditions, iNOS is usually absent in cells. There is some evidence, however, that basal levels of iNOS are expressed in normal cells. For example, continuous NO is generated from iNOS in normal human lung epithelial cells *in vivo* ([72](#)), as well as from the intestinal ileal epithelium in rodents ([73](#)).

Murine macrophage was the first cell type identified to produce NO by iNOS upon microbe infection or inflammation. The large amount of NO produced by activated macrophages accounts for these cells' ability to eliminate some invading pathogens (reviewed in [74](#)). The mechanism underlying the NO-mediated antimicrobial effects is not fully understood, but it may be due to a combination of NO actions on various intracellular targets, such as inhibiting the **mitochondrial** respiratory chain (complex I and II), suppressing total **protein biosynthesis**, and inhibiting aconitase of the tricarboxylic acid cycle. As mentioned above, NO may also react with superoxide ( $O_2^-$ ) generated from activated macrophages or neutrophils to form peroxynitrite, which can generate even more deadly hydroxyl radicals. Widespread inducible NO synthesis is also seen during experimental acute sepsis and both acute and chronic inflammation. In sepsis, increased NO production helps to maximize tissue perfusion through its action on vasculature and prevent microvascular thrombosis by inhibiting platelet aggregation and neutrophil adhesion. However, overexpressed iNOS may unavoidably cause hypotension and shock or even tissue damage ([75](#)). Hence, specific iNOS inhibitors may have significant clinical implications (reviewed in [76](#)).

Unlike murine macrophages, which respond vigorously toward LPS and IFN $\gamma$ , induction of iNOS in human cells, including macrophage and monocytes, was initially difficult to demonstrate, even upon stimulation with various combinations of cytokines and LPS ([77](#)). Based on the finding that rat Kuffer cells (liver macrophages) express iNOS in response to inflammatory stimuli in much the same way as macrophages ([78](#)), subsequent coculture of rat hepatocytes with Kuffer cells or with the conditioned culture medium from LPS and IFN $\gamma$ -treated Kuffer cells led to the discovery that, besides immune systems, hepatocytes are able to produce NO in response to multiple inflammatory products ([79](#), [80](#)). More importantly, the same combination of cytokines synergistically induced high-level expression of iNOS and NO generation in human hepatocytes ([81](#)), providing conclusive evidence that put to rest any doubt about the capacity of human cells to express significant level of iNOS. The human iNOS cDNA and gene were then subsequently cloned and characterized from cytokine-stimulated human hepatocytes ([37](#), [82](#)).

To date, iNOS cDNA has been cloned from mouse ([83](#), [40](#), [84](#)), rat ([85-87](#)) and human([37](#), [88-91](#)). There is 80% sequence identity between the cDNAs of human iNOS and murine macrophage iNOS at both the amino acid and nucleotide levels. This is distinct from the two cNOS isoforms, where the homology between human and other species is consistently around 93% to 95% identity. Human iNOS cDNA has been isolated from five different cell types, namely, stimulated hepatocytes, chondrocytes, a human adenocarcinoma cell line, a human glioblastoma cell line, and cardiac monocytes. They share more than 99% identity or are nearly identical, suggesting existence of a single human iNOS gene.

The human iNOS gene is located at chromosome 17, consists of 26 exons and 25 introns spanning @ 37 kb of the genomic DNA, and encodes a polypeptide of 1153 amino acids ([82](#)). The **transcription** start site was identified 30 bp downstream of a **TATA box**, and the translational **initiation** and **termination codons** reside in exons 2 and 26, respectively. Active iNOS is a homodimer, with each subunit sharing the same **consensus sequences** for binding NADPH, flavins and calmodulin as the two cNOS isoforms (Fig. [2](#)).

Regulation of inducible NO production by iNOS occurs at multiple levels, with transcriptional regulation as the predominant mechanism. Analysis of the 5'-flanking region of human iNOS gene reveals several consensus elements for nuclear factor-kB (NF-kB), TNF **response element** (TNF-RE), and g-IRE, which are reported to be involved in LPS and IFN-g-induced gene expression ([82](#)). A @ 16-kb fragment of the 5'-flanking region of human iNOS gene was found to function as a promoter that is responsive to cytokines and capable of inducing iNOS in a human liver cell line ([92](#)). A general method used to identify the promoter region responsive to various stimuli (in this case cytokines) is to delete the complete promoter from the 5'-end systematically and to connect the shortened ones to a **reporter gene**. After transfecting target cells with various gene constructs, the cytokine inducibility of each construct can be evaluated easily by detecting the expressed reporter.

Using this technology, the human iNOS cytokine-responsive regulatory element (or elements) was mapped upstream of -3.8 kb, namely between -16 to -3.8 kb. Further studies identified three regions containing cytokine-responsive elements in the human iNOS gene: -3.8 to -5.8, -5.8 to -7.0, and -7.0 to -16 kb (92). This contrasts markedly with the murine macrophage iNOS promoter, where the LPS and cytokine-responsive regions locate within only -1 kb of the proximal 5'-flanking region (93-96), suggesting differences in transcriptional regulation between human and murine iNOS gene. Indeed, there is only about 50% sequence identity in the 5'-flanking region of human and murine iNOS genes, in contrast to the 80% identity between their cDNAs.

The cytokine-mediated up-regulation of iNOS gene transcription appears to require the [transcription factor](#) NF- $\kappa$ B in both human and rat hepatocytes (97). Compared to the murine iNOS promoter, where two functional elements reside: an NF- $\kappa$ B element at -76 to -85 bp and an IRF-E/ISRE site at -913 to -923, the human iNOS gene is regulated by a much more complex promoter-regulatory region that requires cytokine-responsive elements upstream of -3.8 kb to confer cytokine-regulated transcription. Studies of the complex mechanisms of iNOS gene regulation will undoubtedly be paramount, given the diverse roles that iNOS plays in various human disease processes.

Besides cytokines and endotoxin, a host of other agents have been shown to affect the induction of iNOS. Glucocorticoids such as dexamethasone decrease iNOS mRNA levels by down-regulating NF- $\kappa$ B activation. Steroids and induction of the **heat shock** response also inhibit iNOS gene induction (98). Hepatocellular **mitogens** and cell growth regulators, such as hepatocyte growth factor (HGF), epidermal growth factor (EGF), and transforming growth factor (TGF- $\alpha$ ), as well as p53 (99), all down-regulate iNOS gene expression. Intriguingly, NO itself suppresses its own gene expression, probably by inhibiting the binding of NF- $\kappa$ B to the iNOS promoter. This kind of negative feedback regulation may prevent cells from overgenerating NO (reviewed in 100).

Another level of regulation of inducible NO production by iNOS is the availability of substrate and cofactors. Since L-arginine is the sole nitrogen source and NADPH the only known electron donor for NO synthesis by NOS isoforms, availability of these two substrates could potentially regulate the rate of NO production. In fact, glucose-6-phosphate dehydrogenase, a regulatory enzyme acting as a valve in the pentose phosphate pathway that generates NADPH, has been shown to be coinduced with iNOS in murine bone marrow-derived macrophages (101). The rate-limiting enzyme for L-arginine biosynthesis, argininosuccinate synthetase, as well as the arginine transport system are up-regulated by the same inflammatory mediators that stimulate inducible NO production (see 102). *De novo* synthesis of BH<sub>4</sub>, another essential cofactor that is required for maximum catalytic activity and for iNOS dimerization, is also necessary for cytokine-induced NO generation in vascular smooth muscle (103). In this regard, the mRNA level of GTP cyclohydrolase-1, the rate-limiting enzyme for BH<sub>4</sub> biosynthesis, is increased along with iNOS in vascular smooth muscle and hepatocytes (104, 105).

In summary, production of NO by a variety of NOS enzymes is under stringent and complex regulation at multiple cellular levels, including cellular NOS distribution and transcriptional, post-transcriptional, translational, and post-translational regulatory mechanisms. Considering the numerous targets of NO, the role of this dual-faced molecule is likely to depend on NO chemistry in a given biological milieu. For example, NO has been recently shown to be able to activate or suppress [apoptosis](#) (programmed cell death). Detailed study into NO production and function under both normal and disease conditions will be the focus of many future investigations.

## Bibliography

1. H. H. W. Schmidt and U. Walter (1994) NO at work. *Cell* **78**, 919-925.
2. H. H. Mitchell, H. A. Schonle, and H. S. Grindly (1916) The origin of the nitrates in the urine. *J. Biol. Chem.* **24**, 461-490.
3. S. R. Tannenbaum, D. Fett, V. R. Young, P. D. Land, and W. R. Bruce (1978) Nitrite and

- nitrate are formed by endogenous synthesis in the human intestine. *Science* **200**, 1487–1489.
4. L. C. Green, S. R. Tannenbaum, and P. Goldman (1981) Nitrate synthesis and reduction in the germ-free and conventional rat. *Science* **212**, 56–58.
  5. D. A. Wagner, V. R. Young, and S. R. Tannenbaum (1983) Mammalian nitrate biosynthesis: incorporation of  $^{15}\text{NH}_3$  into nitrate is enhanced by endotoxin treatment. *Proc. Natl. Acad. Sci. USA* **80**, 4518–4521.
  6. D. J. Stuehr and M. A. Marletta (1985) Mammalian nitrate biosynthesis: mouse macrophages produce nitrite and nitrate in response to *Escherichia coli* lipopolysaccharide. *Proc. Natl. Acad. Sci. USA* **82**, 7738–7742.
  7. R. Iyengar, D. J. Stuehr, and M. A. Marletta (1987) Macrophage synthesis of nitrite, nitrate, and N-nitrosoamines: precursors and role of the respiratory burst. *Proc. Natl. Acad. Sci. USA* **84**, 6369–6373.
  8. J. B. Hibbs Jr., Z. Vavrin, and R. R. Taintor (1987) L-Arginine is required for the expression of the activated macrophage effector mechanism causing selective metabolic inhibition in target cells. *J. Immunol.* **138**, 550–565.
  9. J. B. Hibbs Jr., R. R. Taintor, and Z. Vavrin (1987) Macrophage cytotoxicity: role for L-arginine deiminase and imino nitrogen oxidation to nitrite. *Science* **235**, 473–476.
  10. J. B. Hibbs Jr., R. R. Taintor, Z. Vavrin, and E. M. Rachlin (1988) Nitric oxide: a cytotoxic activated macrophage effector molecule. *Biochem. Biophys. Res. Commun.* **157**, 87–94.
  11. M. A. Marletta, P. S. Yoon, R. Iyengar, C. D. Leaf, and J. S. Wishnok (1988) Macrophage oxidation of L-arginine to nitrite and nitrate: nitric oxide is an intermediate. *Biochemistry* **27**, 8706–8711.
  12. D. J. Stuehr and C. F. Nathan (1989) Nitric oxide: a macrophage product responsible for cytostasis and respiratory inhibition in tumor target cells. *J. Exp. Med.* **169**, 1543–1555.
  13. R. F. Furchgott and J. V. Zawadzki (1980) The obligatory role of endothelial cells in the relaxation of arterial smooth muscle by acetylcholine. *Nature* **288**, 373–376.
  14. R. M. J. Palmer, A. G. Ferrige, and S. Moncada (1987) Nitric oxide release accounts for the biological activity of endothelium-derived relaxing factor. *Nature* **327**, 524–526.
  15. L. J. Ignarro, G. M. Buga, K. S. Wood, R. E. Byrns, and G. Chaudhuri (1987) Endothelium-derived relaxing factor produced and released from artery and vein is nitric oxide. *Proc. Natl. Acad. Sci. USA* **84**, 9265–9269.
  16. M. Feelisch, M. Poel, R. Zamora, A. Deussen, and S. Moncada (1994) Understanding the controversy over the identity of EDRF. *Nature* **368**, 62–65.
  17. R. F. Furchgott (1996) The discovery of endothelium-derived relaxing factor and its importance in the identification of nitric oxide. *JAMA* **276**, 1186–1188.
  18. J. Garthwaite, S. L. Charles, and R. Chess-William (1988) Endothelium-derived relaxing factor release on activation of NMDA receptors suggests a role as intercellular messenger in the brain. *Nature* **336**, 385–388.
  19. J. S. Stamler, D. J. Singel, and J. Loscalzo (1992) Biochemistry of nitric oxide and its redox activated forms. *Science* **258**, 1898–1902.
  20. D. A. Wink, I. Hanbauer, M. B. Grisham, F. Laval, R. W. Nims, J. Laval, J. C. Cook, R. Pacelli, J. Liebmann, M. C. Krishna, M. C. Ford, and J. B. Mitchell (1996) The chemical biology of NO. Insights into regulation, protective and toxic mechanisms of nitric oxide. *Curr. Top. Cell. Regul.* **34**, 159–187.
  21. D. A. Wink, I. Hanbauer, M. C. Krishna, W. DeGraff, J. Gamson, and J. B. Mitchell (1993) Nitric oxide protects against cellular damage and cytotoxicity from reactive oxygen species. *Proc. Natl. Acad. Sci. USA* **90**, 9813–9817.
  22. F. Murad (1996) Signal transduction using nitric oxide and cyclic guanosine monophosphate. *JAMA* **276**, 1189–1192.

23. D. Salvamini, T. Miska, J. Mosferrer, K. Siebert, M. Currie, and P. Needleman (1993) Nitric oxide activates cyclooxygenase enzyme. *Proc. Natl. Acad. Sci. USA* **90**, 7240–7244.
24. D. Salvemini and J. Masferrer (1996) Interactions of nitric oxide with cyclooxygenase: in vitro, ex vivo, and in vivo studies. *Methods Enzymol.* **269**, 12–25.
25. H. M. Lander, J. S. Ogiste, S. F. A. Pearce, R. Levi, and A. Novogrodsky (1995) Nitric oxide-stimulated guanine nucleotide exchange on p21<sup>ras</sup>. *J. Biol. Chem.* **270**, 7017–7020.
26. H. M. Lander, D. P. Hajjar, B. L. Hempstead, U. A. Mirza, B. T. Chait, S. Campbell, and L. A. Quilliam (1997) A molecular redox switch on p21<sup>ras</sup>. Structural basis for the nitric oxide-p21<sup>ras</sup> interaction. *J. Biol. Chem.* **272**, 4323–4326.
27. J. S. Stamler (1994) Redox signaling: nitrosylation and related target interactions of nitric oxide. *Cell* **78**, 931–936.
28. C. Nathan (1992) Nitric oxide as a secretory product of mammalian cells. *FASEB J.* **6**, 3051–3064.
29. S. Dimmeler, J. Haendeler, M. Nehls, and A. M. Zeiher (1997) Suppression of apoptosis by nitric oxide via inhibition of interleukin-1 –converting enzyme (ICE)-like and cysteine protease protein (CPP)-32-like proteases. *J. Exp. Med.* **185**, 601–607.
30. Y. M. Kim, R. V. Talanian, and T. R. Billiar (1997) Nitric oxide inhibits apoptosis by preventing increases in caspase-3-like activity via two distinct mechanisms. *J. Biol. Chem.* **272**, 31138–31148.
31. J. Li, T. R. Billiar, R. V. Talanian, and Y. M. Kim (1997) Nitric oxide reversibly inhibits seven members of the caspase family via S-nitrosylation. *Biochem. Biophys. Res. Commun.* **240**, 419–424.
32. J. S. Stamler, O. Jaraki, J. Osborne, D. I. Simon, J. Keaney, J. Vita, D. Singel, C. R. Valeri, and J. Loscalzo (1992) Nitric oxide circulates in mammalian plasma primarily as an S-nitroso adduct of serum albumin. *Proc. Natl. Acad. Sci. USA* **89**, 7674–7677.
33. J. F. Keaney Jr., D. I. Simon, J. S. Stamler, O. Jaraki, J. Scharfstein, J. A. Vita, and J. Loscalzo (1993) NO forms an adduct with serum albumin that has endothelium-derived relaxing factor-like properties. *J. Clin. Invest.* **91**, 1582–1589.
34. B. Gaston, J. Reilly, J. M. Drazen, J. Fackler, P. Ramdev, D. Arnette, M. E. Mullins, D. J. Sugarbaker, C. Chee, D. J. Singel et al (1993) Endogenous nitrogen oxides and bronchodilator S-nitrosothiols in human airways. *Proc. Natl. Acad. Sci. USA* **90**, 10957–10961.
35. L. Jia, C. Bonaventura, J. Bonaventura, and J. S. Stamler (1996) S-nitrosohaemoglobin: a dynamic activity of blood involved in vascular control. *Nature* **380**, 221–226.
36. H. J. Cho, Q.-W. Xie, J. Calaycay, R. A. Mumford, and C. Nathan (1992) Calmodulin is a subunit of nitric oxide synthase from macrophages. *J. Exp. Med.* **176**, 599–604.
37. D. A. Geller, C. J. Lowenstein, R. A. Shapiro, A. K. Nussler, M. Di Silvio, S. C. Wang, D. K. Nakayama, R. L. Simmons, S. H. Snyder, and T. R. Billiar (1993) Molecular cloning and expression of inducible nitric oxide synthase from human hepatocytes. *Proc. Natl. Acad. Sci. USA* **90**, 3491–3495.
38. R. G. Knowles and S. Moncada (1994) Nitric oxide synthases in mammals. *Biochem. J.* **298**, 249–258.
39. D. S. Bredt, P. M. Hwang, C. E. Glatt, C. Lowenstein, R. R. Reed, and S. H. Snyder (1991) Cloned and expressed nitric oxide synthase structurally resembles cytochrome P-450 reductase. *Nature* **351**, 714–718.
40. C. R. Lyons, G. J. Orloff, and J. M. Cunningham (1992) Molecular cloning and functional expression of an inducible nitric oxide synthase from a murine macrophage cell line. *J. Biol. Chem.* **267**, 6370–6374.
41. D. J. Stuehr (1997) Structure-function aspects in the nitric oxide synthases. *Annu. Rev. Pharmacol. Toxicol.* **37**, 339–359.

42. K. A. White and M. A. Marletta (1992) Nitric oxide synthase is a cytochrome P-450 hemoprotein. *Biochemistry* **31**, 6627–6631.
43. D. K. Ghosh, H. M. Abu-Soud, and D. J. Stuehr (1996) Domains of macrophage NO synthase have divergent roles in forming and stabilizing the active dimeric enzyme. *Biochemistry* **35**, 1444–1449.
44. P. Klatt, S. Pfeiffer, B. M. List, D. Lehner, O. Glatter, H. P. Bachinger, E. R. Werner, K. Schmidt, and B. Mayer (1996) Characterization of heme-deficient neuronal nitric-oxide synthase reveals a role for heme in subunit dimerization and binding of the amino acid substrate and tetrahydrobiopterin. *J. Biol. Chem.* **271**, 7336–7342.
45. K. J. Baek, B. A. Thiel, S. Lucas, and D. J. Stuehr (1993) Macrophage nitric oxide synthase subunits. Purification, characterization, and role of prosthetic groups and substrate in regulating their association into a dimeric enzyme. *J. Biol. Chem.* **268**, 21120–21129.
46. E. Tzeng, T. R. Billiar, P. D. Robbins, M. Loftus, and D. J. Stuehr (1995) Expression of human inducible nitric oxide synthase in a tetrahydrobiopterin (H4B)-deficient cell line: H4B promotes assembly of enzyme subunits into an active dimer. *Proc. Natl. Acad. Sci. USA* **92**, 11771–11775.
47. H. J. Cho, E. Martin, Q.-W. Xie, S. Sessa, and C. Nathan (1995) Inducible nitric oxide synthase: identification of amino acid residues essential for dimerization and binding of tetrahydrobiopterin. *Proc. Natl. Acad. Sci. USA* **92**, 11514–11518.
48. B. R. Crane, A. S. Arvai, R. Gachhui, C. Wu, D. K. Ghosh, E. D. Getzoff, D. J. Stuehr, and J. A. Tainer (1997) The structure of nitric oxide synthase oxygenase domain and inhibitor complexes. *Science* **278**, 425–431.
49. K. Nishida, D. G. Harrison, J. P. Navas, A. A. Fisher, S. P. Dockery, M. Uematsu, R. M. Nerem, R. W. Alexander, and T. J. Murphy (1992) Molecular cloning and characterization of the constitutive bovine aortic endothelial cell nitric oxide synthase. *J. Clin. Invest.* **90**, 2092–2096.
50. R. C. Venema, K. Nishida, R. W. Alexander, D. G. Harrison, and T. J. Murphy (1994) Organization of the bovine gene encoding the endothelial nitric oxide synthase. *Biochim. Biophys. Acta* **1218**, 413–420.
51. P. A. Marsden, H. H. Heng, S. W. Scherer, R. J. Stewart, A. V. Hall, X. M. Shi, L. C. Tsui, and K. T. Schapper (1993) Structure and chromosomal localization of the human constitutive endothelial nitric oxide synthase gene. *J. Biol. Chem.* **268**, 17478–17488.
52. L. J. Robinson, S. Weremowicz, C. C. Morton, and T. Michel (1994) Isolation and chromosomal localization of the human endothelial nitric oxide synthase (NOS3) gene. *Genomics* **19**, 350–357.
53. W. Xu, I. G. Charles, S. Moncada, P. Gorman, D. Sheer, L. Liu, and P. Emson (1994) Mapping of the genes encoding human inducible and endothelial nitric oxide synthase (NOS2 and NOS3) to the pericentric region of chromosome 17 and to chromosome 7, respectively. *Genomics* **21**, 419–422.
54. R. Zhang, W. Min, and W. C. Sessa (1995) Functional analysis of the human endothelial nitric oxide synthase promoter. SP1 and GATA factors are necessary for basal transcription in endothelial cells. *J. Biol. Chem.* **270**, 15320–15326.
55. U. A. Arnet, A. McMillan, J. L. Dinerman, B. Ballermann, and C. J. Lowenstein (1996) Regulation of endothelial nitric oxide synthase during hypoxia. *J. Biol. Chem.* **271**, 15069–15073.
56. W. C. Sessa, K. Pritchard, N. Seyedi, J. Wang, and T. H. Hintze (1994) Chronic exercise increases coronary vascular nitric oxide production and endothelial cells nitric oxide synthase gene expression. *Circ. Res.* **74**, 349–353.
57. M. Yoshizumi, M. A. Perrella, J. C. Burnell Jr., and M. E. Lee (1993) Tumor necrosis factor downregulates an endothelial nitric oxide synthase mRNA by shortening its half life. *Circ. Res.* **73**, 205–209.

58. J. S. Pollock, V. Klinghofer, U. Forstemann, and F. Murad (1992) Endothelial nitric oxide synthase is myristoylated. *FEBS Lett.* **309**, 402–404.
59. L. J. Robinson, L. Busconi, and T. Michel (1995) Agonist-modulated palmitoylation of endothelial nitric oxide synthase. *J. Biol. Chem.* **270**, 995–998.
60. L. Busconi and T. Michel (1993) Endothelial nitric oxide synthase: N-terminal myristoylation determines subcellular localization. *J. Biol. Chem.* **268**, 8410–8413.
61. W. C. Sessa, C. M. Barber, and K. R. Lynch (1993) Mutation of N-myristoylation site converts endothelial cell nitric oxide synthase from a membrane to a cytosol protein. *Circ. Res.* **72**, 921–924.
62. G. Garcia-Cardena, P. Oh, J. Liu, J. E. Schnitzer, and W. C. Sessa (1996) Targeting of nitric oxide synthase to endothelial cell caveolae via palmitoylation: implications for nitric oxide signaling. *Proc. Natl. Acad. Sci. USA* **93**, 6448–6453.
63. O. Feron, L. Belhassen, L. Kobzik, T. W. Smith, R. A. Kelly, and T. Michel (1996) Endothelial nitric oxide synthase targeting to caveolae. *J. Biol. Chem.* **271**, 22810–22814.
64. W. C. Sessa, G. Garcia-Cardena, J. Liu, A. Keh, J. S. Pollock, J. Bradley, S. Thiru, I. M. Braverman, and K. M. Desai (1995) The Golgi association of endothelial nitric oxide synthase is necessary for the efficient synthesis of nitric oxide. *J. Biol. Chem.* **270**, 17641–17644.
65. T. Michel, G. K. Li, and L. Busconi (1993) Phosphorylation and subcellular translocation of endothelial nitric oxide synthase. *Proc. Natl. Acad. Sci. USA* **90**, 6252–6256.
66. I. Fleming, J. Bauersachs, and R. Busse (1997) Calcium-dependent and calcium-independent activation of the endothelial NO synthase. *J. Vasc. Res.* **34**, 165–174.
67. T. Ogura, T. Yokoyama, H. Fujisawa, Y. Kurashima, and H. Esumi (1993) Structural diversity of neuronal nitric oxide synthase mRNA in the nervous system. *Biochem. Biophys. Res. Commun.* **193**, 1014–1022.
68. M. Nakane, H. H. Schmidt, J. S. Pollock, U. Forstemann, and F. Murad (1993) Cloned human brain nitric oxide synthase is highly expressed in skeletal muscle. *FEBS Lett.* **316**, 175–180.
69. A. V. Hall, H. Antoniou, Y. Wang, A. H. Cheung, A. M. Arbus, S. L. Olson, W. C. Lu, C. L. Kau, and P. A. Marsden (1994) Structural organization of the human neuronal nitric oxide synthase gene (NOS1). *J. Biol. Chem.* **269**, 33082–33090.
70. Y. Wang and P. A. Marsden (1995) Nitric oxide synthase: gene structure and regulation. *Adv. Pharmacol.* **34**, 71–90.
71. G. Melillo, T. Musso, A. Sica, L. S. Taylor, G. W. Cox, and L. Varesio (1995) A hypoxia-responsive element mediates a novel pathway of activation of the inducible nitric oxide synthase promoter. *J. Exp. Med.* **182**, 1683–1693.
72. F. H. Guo, H. R. De Raeve, T. W. Rice, D. J. Stuehr, F. B. Thunnissen, and S. C. Erzurum (1995) Continuous nitric oxide synthesis by inducible nitric oxide synthase in normal human airway epithelium in vivo. *Proc. Natl. Acad. Sci. USA* **92**, 7809–7813.
73. R. A. Hoffman, G. Zhang, N. C. Nussler, S. L. Gleixner, H. R. Ford, R. L. Simmons, and S. C. Watkins (1997) Constitutive expression of inducible nitric oxide synthase in the mouse ileal mucosa. *Am. J. Physiol.* **272**, G383–392.
74. A. K. Nussler and T. R. Billiar (1993) Inflammation, immunoregulation, and inducible nitric oxide synthase. *J. Leuk. Biol.* **54**, 171–178.
75. N. Unno, H. Wang, M. J. Menconi, S. H. A. J. Tytgat, V. Larkin, M. Smith, M. J. Morin, A. Chavez, R. A. Hodin, and M. P. Fink (1997) Induction of inducible nitric oxide synthase ameliorates endotoxin-induced gut mucosal barrier dysfunction in rats. *Gastroenterology* **113**, 1246–1257.
76. T. R. Billiar (1995) Nitric oxide novel biology with clinical relevance. *Ann. Surg.* **221**, 339–349.
77. H. W. Murray and R. F. Teitelbaum (1992) L-arginine-dependent reactive nitrogen intermediates the antimicrobial effect of activated human mononuclear phagocytes. *J. Infect.*



Dis. **165**, 513–517.

78. T. R. Billiar, R. D. Curran, D. J. Steuhr, M. A. West, B. A. Bentz, and R. L. Simmons (1989) An L-arginine-dependent mechanism mediates Kupffer cell inhibition of hepatocyte protein synthesis in vitro. *J. Exp. Med.* **169**, 1467–1472.
79. R. D. Curran, T. R. Billiar, D. J. Steuhr, K. Hofmann, and R. L. Simmons (1989) Hepatocytes produce nitrogen oxides from L-arginine in response to inflammatory products of Kuffer cells. *J. Exp. Med.* **170**, 1769–1774.
80. R. D. Curran, T. R. Billiar, H. Ochoa, B. G. Harbrecht, S. G. Flint, and R. L. Simmons (1990) Multiple cytokines are required to induce hepatocyte nitric oxide production and inhibit total protein synthesis. *Ann. Surg.* **212**, 462–471.
81. A. K. Nussler, M. Di Silvio, T. R. Billiar, R. A. Hoffman, D. A. Geller, R. Selby, J. Madariaga, and R. L. Simmons (1992) Stimulation of the nitric oxide synthase pathway in human hepatocytes by cytokines and endotoxin. *J. Exp. Med.* **176**, 261–264.
82. N. A. Chartrain, D. A. Geller, P. P. Koty, N. F. Sitrin, A. K. Nussler, E. P. Hoffman, T. R. Billiar, N. I. Hutchinson, and J. S. Mudgett (1994) Molecular cloning, structure, and chromosomal localization of the human inducible nitric oxide synthase gene. *J. Biol. Chem.* **269**, 6765–6772.
83. C. J. Lowenstein, C. S. Glatt, D. S. Bredt, and S. H. Snyder (1992) Cloned and expressed macrophage nitric oxide synthase contrasts with the brain enzyme. *Proc. Natl. Acad. Sci. USA* **89**, 6711–6715.
84. Q. Xie, H. J. Cho, J. Calaycay, R. A. Mumford, K. M. Swiderek, T. D. Lee, A. Ding, T. Troso, and C. Nathan (1992) Cloning and characterization of inducible nitric oxide synthase from mouse macrophages. *Science* **256**, 225–228.
85. H. Adachi, S. Iida, S. Oguchi, H. Ohshima, H. Suzuki, K. Nagasaki, H. Kawasaki, T. Sugimura, and H. Esumi (1993) Molecular cloning of a cDNA encoding an inducible calmodulin-dependent nitric-oxide synthase from rat liver and its expression in COS 1 cells. *Eur. J. Biochem.* **217**, 37–43.
86. E. R. Wood, H. Berger, P. A. Sherman, and E. G. Lapetina (1993) Hepatocytes and macrophages express an identical cytokine inducible nitric oxide synthase gene. *Biochem. Biophys. Res. Commun.* **191**, 767–774.
87. Y. Nunokawa, N. Ishida, and S. Tanaka (1993) Cloning of inducible nitric oxide synthase in rat vascular smooth muscle cells. *Biochem. Biophys. Res. Commun.* **191**, 89–94.
88. I. G. Charles, R. M. J. Palmer, M. S. Hickery, M. T. Bayliss, A. P. Chubb, V. S. Hall, D. W. Moss, and S. Moncada (1993) Cloning, characterization, and expression of a cDNA encoding an inducible nitric oxide synthase from the human chondrocyte. *Proc. Natl. Acad. Sci. USA* **90**, 11419–11423.
89. P. A. Sherman, V. E. Laubach, B. R. Reep, and E. R. Wood (1993) Purification and cDNA sequence of an inducible nitric oxide synthase from a human tumor cell line. *Biochemistry* **32**, 11600–11605.
90. A. Hokari, M. Zeniya, and H. Esumi (1994) Cloning and functional expression of human inducible nitric oxide synthase (NOS) cDNA from a glioblastoma cell line A-172. *J. Biochem.* **166**, 575–581.
91. H. Luss, R. K. Li, R. A. Shapiro, E. Tzeng, F. X. McGowan, T. Yoneyama, K. Hatakeyama, D. A. Geller, D. A. Mickle, R. L. Simmons, and T. R. Billiar (1997) Dedifferentiated human ventricular cardiac myocytes express inducible nitric oxide synthase mRNA but not protein in response to IL-1, TNF, IFN $\gamma$ , and LPS. *J. Mol. Cell. Cardiol.* **29**, 1153–1165.
92. M. E. De Vera, R. A. Shapiro, A. K. Nussler, J. S. Mudgett, R. L. Simmons, S. M. Morris Jr., T. R. Billiar, and D. A. Geller (1996) Transcriptional regulation of human inducible nitric oxide synthase (NOS2) gene by cytokines: initial analysis of the human NOS2 promoter. *Proc. Natl. Acad. Sci. USA* **93**, 1054–1059.
93. Q.-W. Xie, R. Whisnant, and C. Nathan (1993) Promoter of the mouse gene encoding calcium-

- independent nitric oxide synthase confers inducibility by interferon gamma and bacterial lipopolysaccharide. *J. Exp. Med.* **177**, 1779–1784.
94. Q.-W. Xie, Y. Kashiwabara, and C. Nathan (1994) Role of transcription factor NF-kB/Rel in induction of nitric oxide synthase. *J. Biol. Chem.* **269**, 4705–4708.
  95. C. J. Lowenstein, E. W. Alley, P. Raval, A. M. Snowman, S. H. Snyder, S. W. Russell, and W. J. Murphy (1993) Macrophage nitric oxide synthase gene: two upstream regions mediate induction by interferon and lipopolysaccharide. *Proc. Natl. Acad. Sci. USA* **90**, 9730–9734.
  96. E. Martin, C. Nathan, and Q.-W. Xie (1994) Role of interferon regulatory factor 1 in induction of nitric oxide synthase. *J. Exp. Med.* **180**, 977–984.
  97. M. E. De Vera, B. S. Taylor, Q. Wang, R. A. Shapiro, T. R. Billiar, and D. A. Geller (1997) Dexamethasone suppresses iNOS gene expression by upregulating I-kB and inhibiting NF-kB. *Am. J. Physiol.* **273**, G1290–1296.
  98. D. A. Geller, A. K. Nussler, M. Di Silvio, C. J. Lowenstein, R. A. Shapiro, S. C. Wang, R. L. Simmons, and T. R. Billiar (1993c) Cytokines, endotoxin, and glucocorticoids regulate the expression of inducible nitric oxide synthase in hepatocytes. *Proc. Natl. Acad. Sci. USA* **90**, 522–526.
  99. K. Forrester, S. Ambs, S. E. Lupold, R. B. Kapust, E. A. Spillare, W. C. Weinberg, E. Felley-Bosco, X. W. Wang, D. A. Geller, E. Tzeng, T. R. Billiar, and C. C. Harris (1996) Nitric oxide-induced p53 accumulation and regulation of inducible nitric oxide synthase expression by wild-type p53. *Proc. Natl. Acad. Sci. USA* **93**, 2442–2447.
  100. B. S. Taylor, L. H. Alarcon, and T. R. Billiar (1997) Inducible nitric oxide synthase in the liver: regulation and function. *Biochemistry (Moscow)* **63**, 766–781.
  101. I. M. Corraliza, M. L. Campo, J. M. Fuentes, S. Campos-Portuguez, and G. Soler (1993) Parallel induction of nitric oxide and glucose-6-phosphate dehydrogenase in activated bone marrow derived macrophages. *Biochem. Biophys. Res. Commun.* **196**, 342–347.
  102. S. M. Morris Jr. and T. R. Billiar (1994) New insights into the regulation of inducible nitric oxide synthesis. *Am. J. Physiol.* **266**, E829–839.
  103. S. S. Gross and R. Levi (1992) Tetrahydrobiopterin synthesis: an absolute requirement for cytokine-induced nitric oxide generation by vascular smooth muscle. *J. Biol. Chem.* **267**, 25722–25729.
  104. D. K. Nakayama, D. A. Geller, M. Di Silvio, G. Bloomgarden, P. Davies, B. Pitt, K. Hatakeyama, H. Kagamiyama, R. L. Simmons, and T. R. Billiar (1994) Increased activity of *de novo* tetrahydrobiopterin synthesis in pulmonary artery smooth muscle cells stimulated to produce nitric oxide. *Am. J. Physiol.* **266**, L455–L460.
  105. D. A. Geller, M. Di Silvio, A. Nussler, P. D. Freeswick, D. Nguyen, N. Shah, R. L. Simmons, and T. R. Billiar (1993a) Hepatocyte nitric oxide production during endotoxemia requires expression of GTP cyclohydrolase mRNA and synthesis of tetrahydrobiopterin. *Surg. Forum* **44**, 67–69.
  106. P. L. Huang and M. C. Fishman (1996) Genetic analysis of nitric oxide synthase isoforms: targeted mutation in mice. *J. Mol. Med.* **74**, 415–421.
  107. E. Tzeng and T. R. Billiar (1996) "Genetic modulation of inducible nitric oxide synthase expression". In *Nitric Oxide Synthase Characterization and Functional Analysis* (D. Mahin and M. Maines, eds.), Academic Press, San Diego, CA, 1996, pp. 156–172.

### Suggestions for Further Readings

108. P. L. Feldman, O. W. Griffith, and D. J. Stuehr (1993) The surprising life of nitric oxide. *Chem. Eng. News* **71**, 26–38.
109. C. Nathan and Q.-W. Xie (1994) Regulation of biosynthesis of nitric oxide. *J. Biol. Chem.* **269**, 13725–13728.
110. L. Ignarro and F. Murad, eds. (1995) "Nitric oxide: Biochemistry, molecular biology, and

therapeutic implication". In *Advances in Pharmacology*, vol. **34**, Academic Press, San Diego, CA.

111. S. Moncada and E. A. Higgs, eds. (1990) *Nitric Oxide from L-arginine: A Bioregulatory System*. Elsevier, Amsterdam, the Netherlands.

## Nitrocellulose

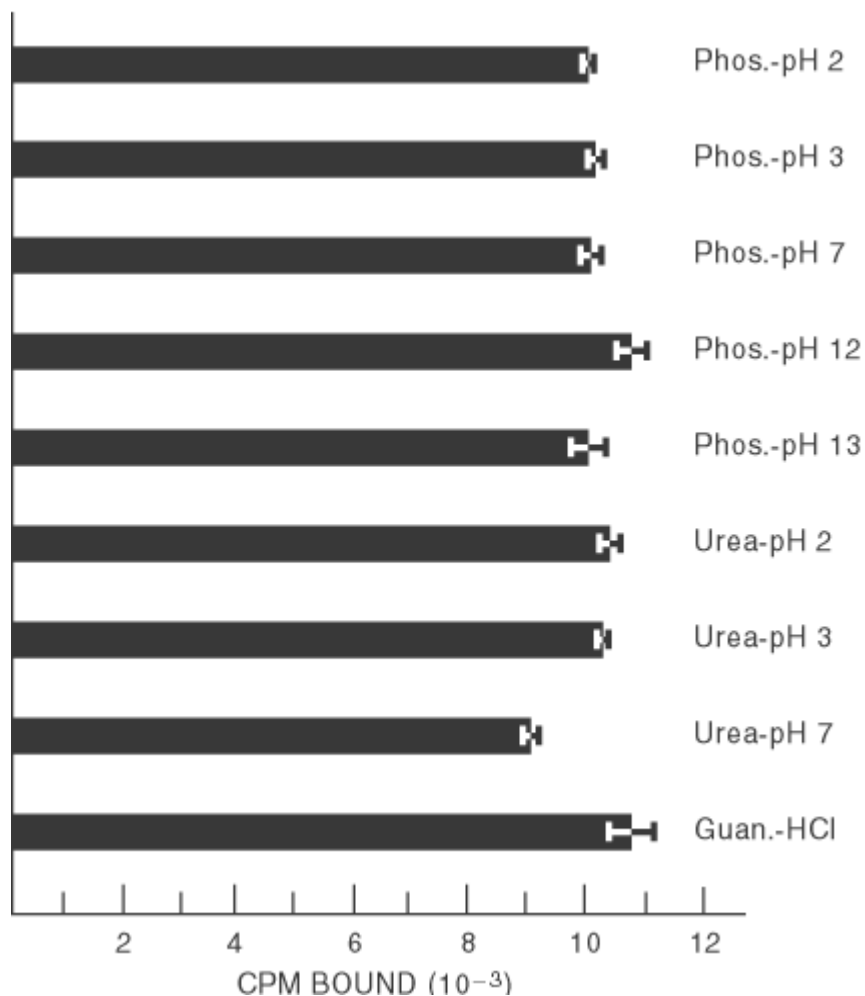
Cellulose nitrate (nitrocellulose) is produced by treating cellulose with nitric acid. Each glucose unit in the cellulose polymer is esterified with three nitrate groups, and these nitrate groups are responsible for both the negative charge of nitrocellulose at neutral pH and the unusual flammability of dry nitrocellulose. Nitrocellulose membranes are made by dissolving the nitrocellulose in organic solvents and spreading the solution as a thin film on a smooth surface. The evaporation rate of the solvent determines the porosity of the membranes, and pore sizes of 0.05 to 0.45  $\mu\text{m}$  are commercially available. Membranes with smaller pore sizes have higher binding capacities for **macromolecules** and are frequently used to bind smaller macromolecules, like **peptides**. Depending on the pore size, pure nitrocellulose membranes have a binding capacity of 80 to 150  $\mu\text{g}/\text{cm}^2$ . Because nitrocellulose membranes can become brittle with harsh treatment, they can be cast on an inert synthetic support, which are also commercially available. The tensile strength of these supported membranes is 500 times greater than ordinary nitrocellulose, but the binding capacity and background levels are the same as ordinary nitrocellulose (1). Some nitrocellulose membranes are not made from "pure" cellulose nitrate but may also contain cellulose acetate. Generally, membranes made with "pure" cellulose nitrate are preferred for most techniques in molecular biology, because "pure" cellulose nitrate has a higher binding capacity.

Nitrocellulose membranes are widely used in molecular biology because they are easy to handle, have a high binding capacity, and are compatible with a variety of assay conditions and detection systems. Although nitrocellulose membranes were originally used to filter out particles, such as **bacteria**, they are now used primarily to adsorb macromolecules throughout the filter matrix, as in [filter binding assays](#) and in [blotting](#). Procedures that frequently specify nitrocellulose membranes include: [Southern blots](#) and **Northern blots**; nucleic acid and protein dot-blot and slot-blot; immunoblots; and colony/plaque lifts. Nitrocellulose membranes are used in **Western blots** of proteins, but they cannot be used for [Edman Degradation](#) sequencing, because the membranes will dissolve in the sequencing solvents.

How macromolecules bind to nitrocellulose is not well understood, but both **electrostatic** and **hydrophobic** interactions have been suggested as possible binding mechanisms (2, 3). The following observations imply that hydrophobic interactions play a dominant role in the binding of macromolecules to nitrocellulose: (a) **SDS**-protein complexes, which are negatively charged, bind to negatively charged nitrocellulose membranes; (b) native double-stranded DNA, which is negatively charged, does not bind to nitrocellulose until the DNA is **denatured** by heat or alkali treatment; (c) nonionic [detergents](#), such as **Triton**<sup>®</sup> X-100 (4), NP-40 (5), and Tween<sup>®</sup> 20 (6), remove proteins from nitrocellulose membranes; and (d) increasing the salt concentration, which increases hydrophobic binding, strengthens the binding of proteins to nitrocellulose. Additional support for the hypothesis that hydrophobic interactions play a dominant role in the binding of proteins to nitrocellulose membranes comes from a study that compared the binding of **antibodies** and other proteins to nitrocellulose membranes in acidic, neutral, basic, and chaotropic buffers (7). In this study, **radiolabeled** proteins of various molecular weights and [isoelectric points](#) bound equally well to nitrocellulose in a variety of buffers (at pH 2, 3, 7, 12, and 13) and in the presence of 8 M [urea](#) (at

pH 2, 3, and 7) and 6 M **guanidinium chloride** (Fig. 1). This study demonstrated that the binding of most proteins to nitrocellulose was not influenced by the ionization of the proteins or by denaturing solvents.

**Figure 1.** Binding of  $^{125}\text{I}$ -labeled human anti-tetanus antibodies to a nitrocellulose membrane in different solvents. The indicated buffers and denaturing solvents were used for both membrane equilibration and protein dilution. After dilution,  $0.3 \mu\text{g}$  of  $^{125}\text{I}$ -labeled human anti-tetanus toxoid antibodies (12,000 cpm) were added to the appropriate wells of a dot-blot unit. After binding the antibodies to the membrane, each was washed four times with  $200 \mu\text{L}$  of  $0.1 \text{ M}$  phosphate buffer (pH 7). The nitrocellulose membrane was removed from the dot-blot unit, and the areas that contained the radioactive proteins were counted in a gamma counter. Each point is the mean of four different wells,  $\pm 2$  standard deviations. (Adapted from Ref. 7.)



Although nitrocellulose membranes are preferred for many nucleic acid and protein immobilization methods, there are a number of precautions that should be followed to protect the membranes, to prepare them for use and to maximize the binding of macromolecules:

1. Store dry nitrocellulose membranes away from organic solvents, oil, and dust. Vapors of some organic solvents may cause the nitrocellulose to curl, and both oil and dust may prevent binding of macromolecules.
2. Wear gloves and use forceps to handle the membranes. Fingerprints will limit binding of macromolecules to nitrocellulose.
3. Wet a nitrocellulose membrane by floating, not immersing, the membrane on the top of water to remove air that may be trapped within the matrix. Air bubbles trapped under the membrane can

be removed by using tweezers to lift a corner. After wetting, immerse the membrane in the appropriate buffer for at least five minutes. Wet nitrocellulose membrane will have a uniform darkening, and any white spots or irregularities are signs of incomplete wetting. Discard any membrane that does not wet uniformly.

4. Macromolecules bind with varying affinities to nitrocellulose. Verify that the macromolecule of interest binds to the nitrocellulose membrane and is not eluted by the detergents or other solutions used in the assay.
5. Solutions containing methanol or ethanol may cause the nitrocellulose to shrink, and acetone will dissolve it. However, xylene and toluene can be used to make the nitrocellulose transparent for scanning with a transmission densitometer (8).
6. Do not exceed 80°C in the vacuum baking oven when drying nucleic acid blots, because nitrocellulose is quite flammable.
7. Unsupported nitrocellulose membranes are brittle when dried or baked at 80°C, so handle them very carefully!

### Bibliography

1. P. Tijssen (1993) *Laboratory Techniques in Biochemistry and Molecular Biology*, Vol. 24, Elsevier, Amsterdam, p. 23.
2. S. R. Farrah, D. O. Shah, and L. O. Ingram (1981) *Proc. Natl. Acad. Sci. USA* **78**, 1229–1232.
3. J. M. Gershoni and G. E. Palade (1983) *Anal. Biochem.* **131**, 1–15.
4. B. Batteiger, W. J. Newhall, and R. B. Jones (1982) *J. Immunol. Methods* **55**, 297–307.
5. W. Lin and H. Kasamatsu (1983) *Anal. Biochem.* **128**, 302–311.
6. W. L. Hoffman and A. A. Jump (1986) *J. Immunol. Methods* **94**, 191–196.
7. W. L. Hoffman et al. (1991) *Anal. Biochem.* **198**, 112–118.
8. D. I. Stott (1989) *J. Immunol. Methods* **119**, 153–187.

### Suggestions for Further Reading

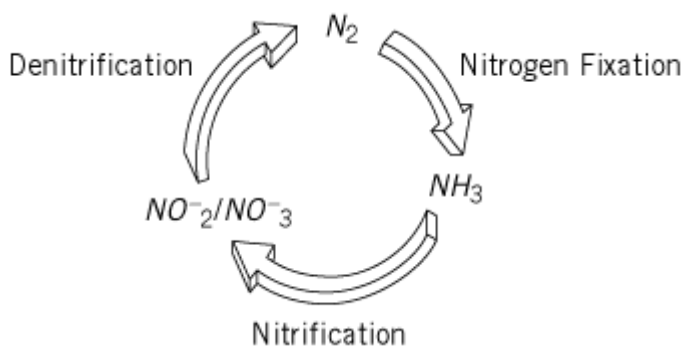
9. L. Nyholm and J. Ramlau (1988) "Nitrocellulose Membranes as Solid Phase in Immunoblotting". In *CRC Handbook of Immunoblotting of Proteins* (O. Bjerrum and N. Heegaard, eds.), CRC Press, Boca Raton, FL, pp. 101–108.
10. W. R. Brown, S. E. Dierks, J. E. Butler, and J. M. Gershoni (1991) "Immunoblotting: membrane filters as the solid phase for immunoassays". In *Immunochemistry of Solid-Phase Immunoassay* (J. E. Butler, ed.), CRC Press, Boca Raton, FL, pp. 151–172.

## Nitrogen Fixation

The conversion of dinitrogen gas ( $N_2$ ) to ammonia ( $NH_3$ ) is called nitrogen fixation. Because ammonia is necessary for the formation of biologically essential, nitrogen-containing compounds, such as **amino acids** and **nucleic acids**, a fixed nitrogen source is necessary to sustain life on earth. Furthermore, the ammonia necessary to support essential biosynthetic reactions is continually sequestered into sediments or reconverted to  $N_2$  through the combined biological processes of nitrification and denitrification (Fig. 1). So the pool of fixed nitrogen within the biosphere must constantly be replenished. Nitrogen fixation is necessary to maintain the diversity of life on earth

because most organisms cannot metabolize the abundant but relatively inert  $N_2$  molecule and must assimilate nitrogen in a “fixed” form, such as ammonia or nitrate. The three ways that nitrogen fixation occurs in the biosphere include (1) lightning and other natural combustion processes, (2) the industrial Haber–Bosch process, and (3) biological nitrogen fixation. Of these three, biological nitrogen fixation, the most significant contributor, accounts for about 65% of the total (1, 2). In addition to its pivotal role in the global nitrogen cycle, nitrogen fixation has agronomic importance because the availability of fixed nitrogen—commonly referred to as fertilizer nitrogen—usually limits to crop production. Consequently, the need to grow sufficient crops to feed the world's population requires the application of nitrogen fertilizers as a common agricultural practice. Such application of nitrogen fertilizers is expensive for several reasons: First, the Haber-Bosch process used to produce fertilizer nitrogen requires high consumption of nonrenewable resources in the form of fossil fuels. Secondly, significant transportation costs are incurred in shipping industrially produced fertilizer nitrogen to the field. Third, the application of nitrogen fertilizers often results in run-off contamination of local water systems. An alternative approach to the Haber-Bosch process for nitrogen fertilizers is to exploit the biological process, and toward this end a considerable effort has been directed at understanding the molecular mechanism of biological nitrogen fixation.

**Figure 1.** Simplified diagram of the biological nitrogen cycle.



## 1. Diazotrophs and nodulation

Although most organisms cannot fix nitrogen, a select group of **microorganisms** can. These organisms are called diazotrophs, which means “nitrogen eaters,” and they are widely distributed among and restricted to the **Archae** and **Bacterial** kingdoms. Examples of bacterial species that fix nitrogen and have been extensively studied include *Azotobacter vinelandii* (an obligate **aerobe**), *Clostridium pasteurianum* (an obligate **anaerobe**), *Klebsiella pneumoniae* (a facultative anaerobe), *Rhodospirillum rubrum* (a **photosynthetic** bacteria), *Anabaena* sp. 7120 (a heterocyst-forming **cyanobacterium**), and *Bradyrhizobium japonicum* (a symbiotic bacterium). From the agronomic perspective, symbiotic nitrogen fixers are the most important. These organisms invade leguminous plants, such as pea and alfalfa, and induce the formation of specialized structures, called nodules, on their roots. In this symbiotic association, the root nodule contains differentiated bacteria that specialize in nitrogen fixation (3). In this symbiosis, the plant provides the microbe a carbon energy source through the process of photosynthesis whereas the microbe supplies the plant host a fixed nitrogen source through the process of nitrogen fixation. Because the **enzyme** that **catalyzes** biological nitrogen fixation, nitrogenase, is extremely oxygen sensitive, another important functional feature of the root nodule is to provide an anoxic environment for nitrogen fixation. Protection of nitrogenase within the root nodule occurs by separating the oxygen-evolving process of photosynthesis and the oxygen-sensitive process of nitrogen fixation. This separation is accomplished in two primary ways: First, the plant-derived, outer cortical layer of the nodule provides a relatively strong barrier against free oxygen diffusion into the central core of the nodule,

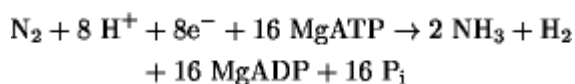
where the nitrogen-fixing bacteria are located and where photosynthesis does not occur. Secondly, the oxygen needed for bacteroidal respiration, which in turn is necessary to provide the energy and reducing equivalents to drive nitrogen fixation within the bacteroid, is delivered by an oxygen-binding plant protein called leghemoglobin.

Establishing of a symbiotic relationship between a leguminous plant host and an associated rhizobium is a complicated process that is specific between a particular host and a particular symbiont (4). During initiation of the infection process, bacteria induce the formation of curls at the tip of plant root hairs (5, 6), which then envelop the invading bacteria. Then infection threads are formed within the root hair. Infection threads are tubular structures of plant origin that penetrate the root hairs and the root cortex through which the invading bacterium traverse (7). At or about the time infection thread formation occurs, **cell division** is also induced within the root cortex, resulting in the formation of a nodule primordium (8). When the infection thread contacts a newly divided primordial cell within the root cortex, the bacteria are released from the infection thread tip. Subsequent penetration of the released bacterial cell into the cytoplasm of the primordial cell occurs through a process of [endocytosis](#) resulting in the formation of a plant cell membrane (bacteroidal membrane) that surrounds the bacterium (9). Bacterial cell division and bacteroid membrane proliferation then occur, resulting in the host plant cell cytoplasm becoming filled with bacteria. These bacteria then differentiate into cells specialized for nitrogen fixation that are called bacteroids.

To establish an effective symbiosis, developmental and metabolic cooperation between the plant and microbe are necessary. This cooperation is accomplished through the reciprocal communication and control of gene expression between the two partners. During nodule development, signaling occurs through the action of nodulation (*nod*) **gene expression** by the bacterium and expression of nodulins (ENOD sequences) by the plant. The first of these signals is provided by flavanoid molecules, three-ring aromatic compounds derived from phenylpropanoid metabolism, which are produced by the plant and found in the root exudate (10). Flavanoids act as specific **inducers** of bacterial *nodD* gene expression (11). The *nodD* gene encodes a regulatory protein that controls the expression of other *nod* genes. The concerted action of other *nod* gene products, whose synthesis is induced by *nodD* gene expression, results in forming and releasing bacterial signal molecules called Nod factors (12). Nod factors are responsible for signaling the initial root hair curling event. Their core structures are composed of b1-4 linked *N*-acetylglucosamine residues of various length (usually one to four molecules). The specificity of a particular Nod factor is determined by its level of oligomerization and by certain chemical modifications of the oligosaccharide core, such as O-acylation or N-acylation at the nonreducing end of the oligosaccharide and sulfation at the reducing end. The *nodABC* gene products are common among all the rhizobia, and they are required for catalyzing formation of the oligosaccharide core, whereas other “host-specific” *nod* gene products are responsible for oligomerizing and modifying the oligosaccharide core.

### 1.1. Nitrogenase

Nitrogenase is the enzyme that catalyzes biological nitrogen fixation. The reaction is usually depicted as follows:



All diazotrophs studied so far contain a nitrogenase that is a complex, two-component [metalloprotein](#) (13-15). The individual nitrogenase component proteins have been designated as the Fe protein and the MoFe protein, which were derived from the respective compositions of their associated metal centers. The Fe protein is a homodimer ( $M_r$  64,000) that contains two MgATP binding sites and a single [4Fe-4S] cluster (see [Iron-Sulfur Proteins](#)). The MoFe protein is an  $a_2b_2$  heterotetramer ( $M_r$  250,000) that contains two pairs of metalloclusters, called P-clusters, and iron molybdenum cofactors (FeMo cofactors). Ribbon diagrams for the three dimensional structures of the Fe protein and MoFe protein are shown in Fig. 2, and structures of the associated nitrogenase metal clusters are

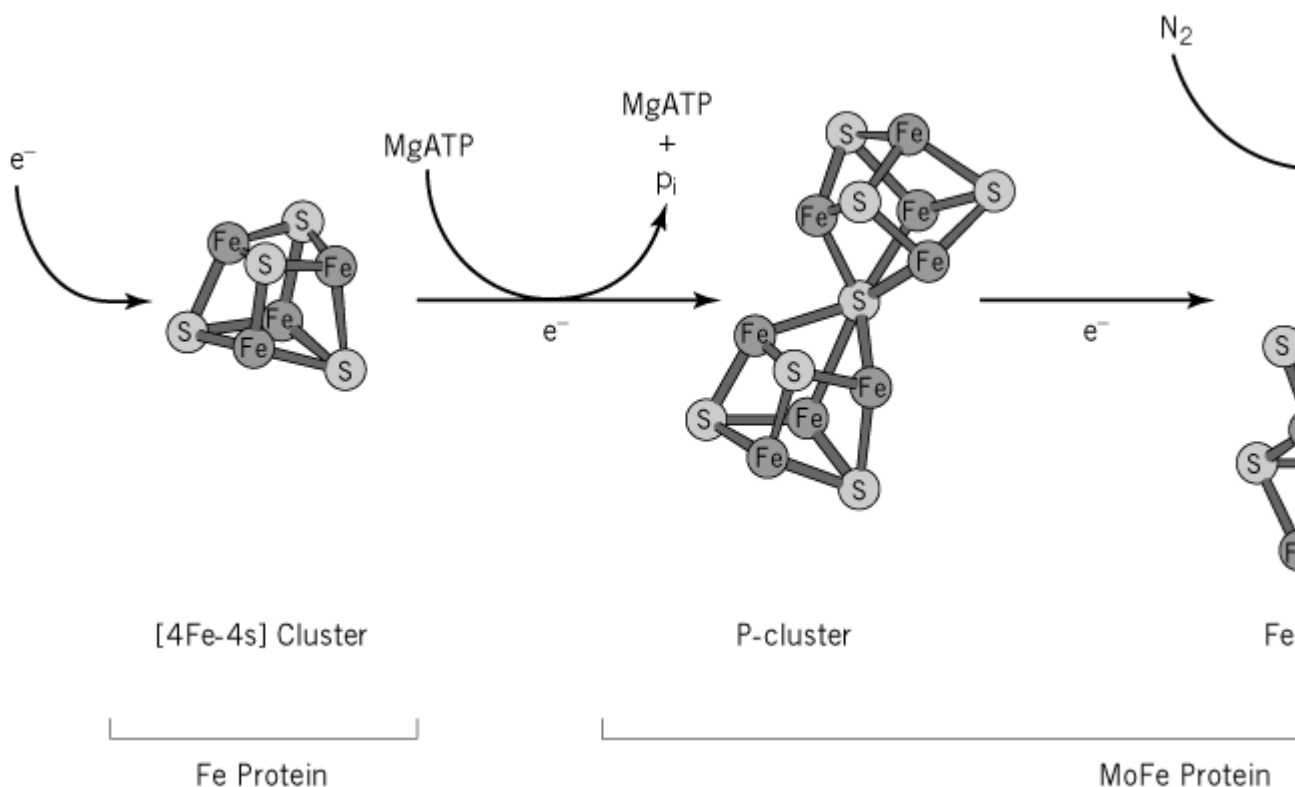
shown in Fig. 3. Because the Fe protein is a specific reductant of the MoFe protein, which in turn provides the site of substrate reduction, some investigators refer to the Fe protein as dinitrogenase reductase and the MoFe protein as dinitrogenase (16).

**Figure 2.** Ribbon diagrams of the three-dimensional structures of the nitrogenase Fe protein and MoFe protein. The view interaction with a single ab-unit of the MoFe protein (bottom).



**Figure 3.** Organization and structures of the nitrogenase metalloclusters. The path of electrons is from the [4Fe-4S] cluster and then to the FeMo cofactor. Transfer of an electron from the Fe protein to the MoFe protein is coupled to  $Mg^{2+}$  provides the  $N_2$  binding and reduction site.





### 1.1.1. Nitrogenase Mechanism and Role of the Metal Centers

During catalysis, the Fe protein delivers electrons, one at a time, to the MoFe protein in a process that couples MgATP binding and hydrolysis to the association and dissociation of the two component proteins and concomitant **electron transfer**. Both component proteins are required for MgATP hydrolysis, and neither component protein reduces any substrate in the absence of its catalytic partner. The process by which electrons are sequentially delivered to the MoFe protein and subsequently to substrate has been described by a kinetic model that involves two interlocking cycles called the Fe protein cycle and the MoFe protein cycle (17). The Fe protein cycle involves oxidizing and reducing the Fe protein's [4Fe-4S] cluster between the  $1^+$  and  $2^+$  **redox** states as it sequentially delivers electrons to the MoFe protein and is re-reduced by other electron transfer proteins (usually a [ferredoxin](#) or flavodoxin). The MoFe protein cycle involves the progressive reduction of the MoFe protein, which ultimately leads to  $N_2$  binding and reduction. Because eight electrons and eight protons are required for  $N_2$  reduction and  $H_2$  evolution, each MoFe protein cycle requires eight Fe protein cycles and storage of the electrons. Kinetic studies have also shown that  $N_2$  does not become bound to the active site until at least two, and probably more, electrons have been accumulated within the MoFe protein (18). It is not yet known where and how the electrons delivered to the MoFe protein are stored before to the binding and reduction of the substrate. However, current structural and biochemical data indicates that cluster-to-cluster electron transfer occurs as shown in Fig. 3 (13, 14, 19).

In addition to the individual structural models for the Fe protein and MoFe protein (Fig. 2), an **X-ray crystallographic** model for the docked complex has also been determined (19). In this model, docking occurs so that the two-fold symmetry axis surrounding the Fe protein's [4Fe-4S] cluster becomes paired with the surface of the MoFe protein's pseudosymmetrical ab-interface. This places the [4Fe-4S] cluster close to the P-cluster of the MoFe protein and places the P-cluster between the Fe protein [4Fe-4S] cluster and the FeMo cofactor. This arrangement is consistent with the electron transfer scheme shown in Fig. 3 and indicates that the P-cluster's role is to broker electron transfer between the Fe protein and the substrate reduction site provided by the FeMo cofactor. The P-cluster

is located at the ab interface of the MoFe protein subunits and, in its fully reduced form, is constructed from two [4Fe-4S] subclusters that share a central sulfide (Fig. 3). Upon oxidation of the P-cluster, a structural rearrangement occurs involving movement of two Fe atoms and a change in the ligand arrangement around the cluster (20). Such redox-dependent structural changes within the P-cluster might be mechanistically related to its role in accepting electrons from the Fe protein and delivering them to the FeMo cofactor.

The FeMo cofactor consists of a metal sulfur framework ( $\text{MoFe}_7\text{S}_9$ ) and one molecule of (R)-homocitrate (Fig. 3). The framework is constructed from S-bridged  $\text{MoFe}_3\text{S}_3$  and  $\text{Fe}_4\text{S}_3$  cluster subfragments. Homocitrate is coordinated to the Mo atom through its  $\beta$ -hydroxy and  $\beta$ -carboxy groups. Several lines of evidence indicate that the FeMo-cofactor provides the substrate binding and reduction site: First, mutant strains that cannot biosynthesize the FeMo cofactor also cannot catalyze nitrogen fixation (21). When the isolated FeMo cofactor is added to crude extracts prepared from such mutant strains, the *in vitro* ability to fix nitrogen is restored. Secondly, a MoFe protein that contains an altered form of the FeMo cofactor, where citrate replaces homocitrate as its organic constituent, exhibits altered catalytic properties (22, 23). Thirdly, altered MoFe proteins that have amino acid substitutions located within the FeMo cofactor's polypeptide environment also exhibit altered catalytic properties (24). Although it is not yet known how substrates interact with the FeMo cofactor during turnover, the presence of six coordinately unsaturated Fe atoms and the attachment of homocitrate to the Mo atom has invited speculation about the nature of substrate binding (25). For example, in one model it has been proposed that the carboxylate group coordinated to the Mo atom might serve as a leaving group in a mechanism that activates Mo to provide a substrate coordination site (26).

#### 1.1.2. The Role of MgATP in Nitrogenase Catalysis

The reduction of  $\text{N}_2$  to yield  $2 \text{NH}_3$  is thermodynamically favorable. Thus, the need for MgATP binding and hydrolysis during nitrogenase catalysis is kinetic so that electron transfer toward substrate reduction is favored and the flow of electrons back toward the Fe protein is prevented. One way of envisioning this process is to consider that the energy released through MgATP binding and hydrolysis could be used to open and close electron gates to ensure that multiple electrons are accumulated within the MoFe protein before they are donated to substrates. In support of this model, primary sequence and structural comparisons have revealed that the Fe protein is a member of a large class of [signal transduction](#) proteins that undergo conformational changes upon MgATP binding and hydrolysis. A consensus view (27) of the events that occur during a single turn of the Fe protein cycle and which accounts for the role of MgATP in nitrogenase catalysis is as follows. Intercomponent interaction is initiated when the reduced Fe protein binds two MgATP molecules. This elicits a conformational change in the Fe protein that makes it competent to interact with the MoFe protein. Upon complex formation, changes occur in the midpoint potentials of the respective clusters such that electron flow toward the FeMo cofactor is energetically favorable. For example, in the complexed form, the Fe protein's [4Fe-4S] cluster has a **redox potential** of about  $-620\text{mV}$ , the P-cluster a potential of about  $-390\text{mV}$ , and the FeMo cofactor a potential of about  $-40\text{mV}$  (28). In addition to eliciting redox changes that lead to electron transfer, the component docking event also triggers MgATP hydrolysis at about the same time as, or shortly after, electron transfer. Conversion of the Fe protein from the MgATP-bound state to the MgADP-bound state subsequently causes complex dissociation, which is believed to be the **rate-limiting step** in nitrogenase catalysis (29). Thus, the accumulation of electrons within the MoFe protein is a dynamic process that involves transmitting signals back and forth between the Fe protein and MoFe protein. Then the role of MgATP is to synchronize these events through sequential conformational changes induced by MgATP binding, component protein interaction, and nucleotide hydrolysis.

#### 1.1.3. Alternative Nitrogenases

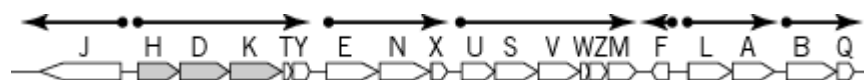
All organisms capable of nitrogen fixation that have been examined at the biochemical level have a Mo-containing nitrogenase as described previously. All such nitrogenases exhibit a high level of **primary structural** identity compared to each other. Particularly high sequential conservation is

found in the respective MgATP- and metallocluster-binding sites (30). There is, however, another class of nitrogenases that do not contain Mo, and these have been designated alternative nitrogenases (31). Two types of alternative nitrogenases have been identified so far (32, 33). One of these contains a cofactor whose Mo atom is substituted by vanadium, V (FeV cofactor), and the other contains a cofactor whose Mo atom is substituted by Fe (FeFe cofactor). The conservation in primary sequences recognized among the Mo-dependent nitrogenases also extends to the alternative nitrogenases. Thus, all nitrogenases have common structural and mechanistic features. All three nitrogenase types are found in *A. vinelandii*, and expression of each of them is under a hierarchical control that depends on the availability of Mo or V in the growth medium (34, 35). When Mo is present in the growth medium, expression of the Mo-dependent nitrogenase is stimulated, and the expression of the alternative nitrogenases is repressed. Similarly, when Mo is absent, but V is available, only expression of the V-dependent enzyme occurs. When neither Mo nor V is available, only the alternative nitrogenase that contains FeFe cofactor is expressed. Such hierarchical control by the availability of metals makes physiological sense because the Mo-dependent nitrogenase has an intrinsically higher capacity to reduce nitrogen than either of the alternative nitrogenases, and the V-dependent nitrogenase is more efficient than the Fe-only nitrogenase (36).

## 1.2. Nif-Genes and Maturation of Nitrogenase

In addition to the structural genes for nitrogenase, other genes are required for (1) coupling the reduction of the Fe protein to intermediary metabolism, (2) maturation of the nitrogenase components, and (3) regulating the expression of the nitrogenase genes. The organism that has the simplest organization of nitrogen fixation specific (*nif*) genes and is best studied at the molecular genetic level is the facultative anaerobe *Klebsiella pneumoniae*. There are 20 *nif* genes in this organism, organized into seven transcriptional units (Fig. 4). The specific designations for individual *Klebsiella pneumoniae nif* genes are also used to denote genes whose products have homologous functions in other organisms. For example, the *nif* structural genes from all diazotrophs are designated *nifH*, *nifD*, and *nifK*, and they respectively encode the Fe protein and the a- and b-subunits of the MoFe protein. A complication in the genetic nomenclature of genes, whose products are involved in nitrogen fixation, is that some organisms do not have homologues to all of the *Klebsiella pneumoniae nif*-specific genes, but have other genes related to nitrogen fixation that are not present in *K. pneumoniae*. The general convention used is that the “*nif*” designation is reserved only for those genes that have functional counterparts in *K. pneumoniae*. Genes in other organisms, whose products are involved in the process of nitrogen fixation but do not have functional counterparts in *K. pneumoniae*, have been given various other designations. One example is the designation “*fix*” which designates such genes from the rhizobia.

**Figure 4.** Organization of the nitrogen fixation specific (*nif*) genes of *Klebsiella pneumoniae*. Black arrows indicate transcription units and the direction of their transcription. Structural genes encoding the nitrogenase Fe protein (*nifH*) and MoFe protein subunits (*nifDK*) are shaded. Functions of the *nif* gene products are described in the text.



### 1.2.1. Electron Transport to Nitrogenase

A source of reducing equivalents of sufficiently low potential is required to regenerate reduced Fe protein after it has donated an electron to the MoFe protein. Both flavodoxins and ferredoxins serve this function *in vitro*. In *K. pneumoniae*, two genes, *nifF* and *nifJ*, are involved in coupling the reduction of Fe protein to intermediary metabolism (37, 38). The *nifF* gene product is a flavodoxin that, in its reduced hydroquinone form, donates an electron to the oxidized Fe protein, thereby generating the semiquinone form of the flavin moiety. The re-reduction of flavodoxin is accomplished through the catalytic activity of the *nifJ* gene product (pyruvate-flavodoxin oxidoreductase), which couples the oxidation of pyruvate, yielding acetyl-CoA and CO<sub>2</sub>, to the reduction of the semiquinone form of flavodoxin to the hydroquinone form.

### 1.2.2. Maturation of the Fe Protein and MoFe Protein

The primary translation products of the nitrogenase structural genes are not active. Instead, a consortium of other *nif*-specific genes are required to activate the immature structural components. The function of the *nif*-specific maturation gene products is to catalyze the formation and insertion of the metalloclusters into apo-forms of the Fe protein and MoFe protein. In the Fe protein, only the *nifM* gene product is specifically required for its maturation (39). The *nifM* gene product has not yet been isolated in an active form, but it is a member of a family of [peptidyl prolyl cis/trans isomerases](#) (40). Such enzymes are thought to assist in **protein folding** by catalyzing the *cis/trans isomerization* of certain prolyl peptide bonds in some proteins. The requirement for such an activity in maturation of the Fe protein is not obvious but could be related to the appropriate organization of the [cysteine](#) ligands needed for properly inserting the [4Fe-4S] cluster that is located at the Fe protein dimer interface.

Maturation of the MoFe protein, particularly the formation and insertion of the FeMo cofactor into the MoFe protein, is much more complicated. This process involves the products of the *nifH*, *nifE*, *nifN*, *nifB*, *nifV*, and *nifQ* genes (41). The nitrogenase Fe protein, a product of *nifH*, is required for both forming and inserting the FeMo cofactor (42, 43). Although the specific function of Fe protein in these processes is not known, neither its MgATP-binding or MgATP-hydrolytic properties nor its ability to transfer electrons are necessary (44, 45). Biochemical complementation experiments have shown that the FeMo-cofactor is synthesized separately and then inserted into an apo-form of the MoFe protein that contains intact P-clusters but lacks the FeMo cofactor (46). At least a portion of the biosynthetic process occurs within a complex of the *nifEN* products. The *nifEN* gene products bear primary sequential similarity to the products of the *nifDK* structural genes (47), and an  $\alpha_2\beta_2$  complex of the *nifEN* products is a molecular scaffold for FeMo cofactor assembly (48). The product of the *nifB* gene catalyzes the formation of a FeMo cofactor precursor called the B-cofactor (49). The B-cofactor probably provides the Fe-S cage necessary for FeMo cofactor construction and is inserted into the NifEN complex at an intermediate stage in FeMo cofactor formation (50). The *nifV* gene catalyzes the condensation of acetyl-CoA and  $\alpha$ -ketoglutarate to form homocitrate, the organic constituent of the FeMo cofactor (51, 52). The *nifQ* gene product has a role in inserting the Mo atom into the FeMo cofactor, but its exact role in this process is not known (53). The products of the *nifW* and *nifZ* genes might also have some role in FeMo cofactor biosynthesis because mutations in either or both of them result in lowered MoFe protein activity (54). In *K. pneumoniae*, apo-MoFe protein produced by strains lacking *nifB* or *nifEN* activity contain a low molecular weight protein encoded by *nifY*. The *nifY* gene product has a role in stabilizing a conformation of the apo-MoFe protein that is amenable to FeMo cofactor insertion (55, 56). In *A. vinelandii*, a different low molecular weight protein (called gamma), not encoded by *nifY*, serves the same function (57).

In addition to the gene products that are specifically involved in nitrogenase metallocluster biosynthesis, there are two other *nif*-genes, *nifS* and *nifU*, whose products catalyze reactions that are generally involved in mobilizing Fe and S for metallocluster assembly (58, 59). The *nifS* gene product is a cysteine desulfurase that activates S for [Fe-S] cluster formation, whereas the *nifU* gene product probably has a complementary role as an Fe carrier. Homologues to *nifU* and *nifS*, whose expression is not under *nif*-specific control, are present in many organisms. The products of these genes probably have general functions in [Fe-S] cluster formation and repair (60).

### 1.3. Regulation of Nitrogenase Expression

For most free-living, nitrogen-fixing organisms, such as *K. pneumoniae*, *nif*-gene expression should be responsive to three environmental conditions. These parameters include (1) the availability of a fixed nitrogen source, (2) whether or not oxygen is present, and (3) the energy charge status of the cell (see [Adenylate Charge](#)). The reasons that these conditions are important is that nitrogenase activity is both oxygen-sensitive and requires the intense consumption of metabolic energy. Thus, if oxygen is present, if a fixed nitrogen source is already available, and if cellular growth is limited by the availability of a carbon source rather than a nitrogen source, the

expression of *nif*-genes would represent a waste of metabolic energy. In free-living nitrogen fixers that have alternative nitrogenases, the expression of *nif*-genes is additionally controlled by the availability of Mo or V.

The organism that has been most thoroughly studied in terms of *nif*-gene expression is *K. pneumoniae*. The 20 *nif*-specific genes from *K. pneumoniae* are organized into seven transcriptional units (Fig. 4). Six of these transcription units include the *nif* structural genes, electron transport genes, and nitrogenase maturation genes, whereas the seventh transcription unit contains the *nif*-specific regulatory genes, *nifL* and *nifA* (61). Expression of the *nifLA* genes is controlled by the global regulatory elements, products of the *ntrC*, *ntrB*, and *ntrA* genes. The expression of the other *nif* transcriptional units is controlled by the products of *nifL*, *nifA*, and *ntrA*. The *ntrA* gene product (NTRA) is an alternative [sigma factor](#) that controls the expression of a wide variety of [gene families](#). The *nif* and *ntr* genes are only two examples (62). NTRA (also called  $\sigma^{54}$ ) imparts a specificity to **RNA polymerase** so that it recognizes the consensus **promoter** sequence CTGG-N<sub>8</sub>-TTGCA. This consensus sequence spans a region located 24 to 12 base pairs preceding the transcription initiation site. In contrast, normal RNA polymerase contains the abundant sigma factor, called  $\sigma^{70}$ , and this form of RNA polymerase recognizes the typical prokaryotic promoter sequence TTGACA-N<sub>17</sub>-TATACA. The so-called housekeeping RNA polymerase binding site is located in the -35, -10 regions preceding the transcription initiation site. Expression of *ntrA* is not tightly controlled. Instead, the presence of NTRA is a way for the cell to reserve a certain portion of the RNA polymerase for specialized functions, such as nitrogen fixation.

Specific regulation of the expression of the *nifLA* genes is controlled by the products of the global nitrogen regulatory genes *ntrC* (NTRC) and *ntrB* (NTRB). NTRB is a phosphatase/kinase sensor protein that controls the **phosphorylation** of NTRC in response to the ratio of  $\alpha$ -ketoglutarate to ammonia in the cell (63). When this ratio is high, NTRC is phosphorylated. When the ratio is low, NTRC is dephosphorylated. The former represents a condition of fixed nitrogen limitation and high energy charge, which signals the initiation of a two-tiered regulatory cascade that first leads to activation of the *nifLA* promoter and then results in activation of the other *nif* promoters. In its phosphorylated form, NTRC recognizes a DNA **consensus sequence** located upstream from the *nifA*  $\sigma^{54}$ -**RNA polymerase** binding site. Once bound at this upstream activator site, phosphorylated NTRC catalyzes an ATP-dependent conformational isomerization at the promoter site, which ultimately results in transcriptional initiation of the *nifLA* promoter (64). Then accumulation of the products of *nifL* (NIFL) and *nifA* (NIFA) permits specific control of the other *nif*-promoters. NIFA is structurally and functionally similar to NTRC (65, 66). Like NTRC, it binds to an upstream activator sequence, although one that differs in consensus sequence to the NTRC binding site. The NIFA binding site has a consensus sequence motif (TGT-N<sub>10</sub>-ACA) located approximately 100 bp preceding each of the *nif* promoters, except the *nifLA* promoter (67). NIFA probably activates expression of the *nif* promoters in an ATP-dependent manner analogous to the function of the phosphorylated form of NTRC.

Unlike NTRC, NIFA does not undergo phosphorylation and dephosphorylation in response to environmental signals. Instead, whether or not NIFA activates *nif* gene expression is controlled by NIFL, which acts as an antiactivator. When the oxygen level or the level of fixed nitrogen is sufficiently high, so that nitrogen fixation is either futile or not necessary, NIFL interacts with NIFA to prevent activation of the *nif* promoters. Although the details of how the complexation of NIFA and NIFL occurs are not known, it has been shown that NIFL is a flavoprotein which senses the redox status of the cell by conformational changes controlled by the oxidation state of its FAD moiety (68).

The mechanism by which *K. pneumoniae* controls *nif* gene expression is by no means universal.

For example, very little is known about how *nif* gene expression is controlled in the **clostridia** or the Archaea, although control does not appear to be through the two-tiered NTR-dependent mechanism described previously. In certain organisms another layer of nitrogenase control activity also occurs at the posttranslational level. The best described example of this type of control is the reversible **ADP-ribosylation** of the nitrogenase Fe protein that occurs in *Rhodobacter rubrum*. In this organism, nitrogenase activity is controlled at the transcriptional level and also through the antagonistic activity of ADP-ribosylation and glycohydrolase enzymes whose activities are controlled by various environmental conditions (69, 70). Another example of the control of the expression of nitrogenase activity is found in the cyanobacterium *Anabaena* 7120 (71). This organism restricts nitrogen fixation to specialized cells called heterocysts. The formation of heterocysts, and nitrogenase expression are controlled both temporally and spatially by terminal differentiation events involving rearrangements of the [genome](#).

Another specialized case of the control of nitrogenase gene expression occurs in the root nodule. It does not make sense to control nitrogenase gene expression and activity in response to the level of fixed nitrogen in the root nodule because the bacteroid functions to supply its host with fixed nitrogen. Thus, in the root nodule, *nif* gene expression is controlled by a two-tiered cascade that responds to the cellular oxygen tension. Similar to the NTRC/NTRB system, oxygen control of *nif* gene expression in rhizobia occurs through the concerted activity of an environmental sensor (FIXL) and a positive regulator of gene expression (FIXJ) (72). The sensor protein FIXL is a heme protein located within the cell membrane, and like NTRB, it has phosphatase/kinase activities. At low oxygen tension (microaerobicity), FIXL promotes the phosphorylation of FIXJ, which then activates the expression of *nifA* analogously to NTRC activation in *K. pneumoniae*. The rhizobia do not contain a protein analogous to NIFL, so control of *nif* gene expression in this case does not involve an antiactivator. Instead of the involvement of an antiactivator in the posttranslational control of certain rhizobial NIFA proteins, NIFA activity is directly sensitive to the presence of oxygen (73).

## Bibliography

1. R. C. Burns and R. W. Hardy (1975) *Nitrogen Fixation in Bacteria and Higher Plants*, Springer-Verlag, Berlin, p. 43.
2. W. E. Newton and W. H. Orme-Johnson, eds. (1980) *Nitrogen Fixation*, University Park Press, Baltimore, Maryland, 1980.
3. S. R. Long (1989) *Cell* **56**, 203–214.
4. J. Denarie, F. Debelle, and C. Rosenberg (1992) *Ann. Rev. Microbiol.* **46**, 497–531.
5. P. Y. Yao and J. M. Vincent (1969) *Aust. J. Biol. Sci.* **22**, 413–423.
6. F. B. Dazzo and A. Gardiol (1984) In *Genes Involved in Microbe Plant Interactions*, Springer-Verlag, New York, pp. 3–31.
7. D. Verma and S. Long (1983) In *Intracellular Symbiosis* (K. Jeon, ed.), Academic Press, New York, pp. 211–245.
8. K. R. Libbenga and P. A. A. Harkes (1973) *Planta* **114**, 17–28.
9. J. G. Robertson, P. Lyttleton, S. Bullivant, and G. F. Grayston (1978) *J. Cell Sci.* **30**, 129–149.
10. N. K. Peters, J. W. Frost, and S. R. Long (1986) *Science* **233**, 917–1008.
11. H. P. Spaink, C. A. Wijffelman, E. Pees, R. J. H. Okker, and B. J. J. Lugtenberg (1987) *Nature* **328**, 337–340.
12. P. Lerouge, P. Roche, C. Faucher, F. Maillet, G. Truchet, J. C. Prome, and J. Denarie (1990) *Nature* **344**, 781–784.
13. J. B. Howard and D. C. Rees (1996) *Chem. Rev.* **96**, 2965–2982.
14. B. K. Burgess and D. J. Lowe (1996) *Chem. Rev.* **96**, 2983–3011.
15. J. W. Peters, K. Fisher, and D. R. Dean (1995) *Annu. Rev. Microbiol.* **49**, 335–366.

16. R. V. Hageman and R. H. Burris (1978) *Proc. Natl. Acad. Sci. USA* **75**, 2699–2702.
17. W. H. Orme-Johnson (1985) *Ann. Rev. Biophys. Biophys. Chem.* **14**, 419–459.
18. D. J. Lowe and R. N. F. Thorneley (1984) *Biochem. J.* **224**, 903–909.
19. H. Schindelin, C. Kisker, J. L. Schlessman, J. B. Howard, and D. C. Rees (1997) *Nature* **387**, 370–376.
20. J. W. Peters, M. H. B. Stowell, S. M. Soltis, M. G. Finnegan, M. K. Johnson, and D. C. Rees (1997) *Biochemistry* **36**, 1181–1187.
21. V. K. Shah and W. J. Brill (1977) *Proc. Natl. Acad. Sci. USA* **74**, 3249–3253.
22. T. R. Hawkes, P. A. McLean, and B. E. Smith (1984) *Biochem. J.* **217**, 317–321.
23. J. Liang, M. Madden, V. K. Shah, and R. H. Burris (1990) *Biochemistry* **29**, 8577–8581.
24. D. J. Scott, H. D. May, W. E. Newton, K. E. Brigle, and D. R. Dean (1990) *Nature* **343**, 188–190.
25. I. Dance (1996) *J. Biol. Inorg. Chem.* **1**, 581–586.
26. D. L. Hughes, S. K. Ibrahim, C. J. Pickett, G. Querne, A. Laouenen, J. Talarmin, A. Queiros, and A. Fonseca (1994) *Polyhedron* **13**, 3341–3348.
27. L. C. Seefeldt and D. R. Dean (1997) *Acc. Chem. Res.* **30**, 260–266.
28. W. N. Lanzilotta and L. C. Seefeldt (1997) *Biochemistry*, in press.
29. R. N. F. Thorneley and D. J. Lowe (1984) *Biochem. J.* **224**, 887–894.
30. D. R. Dean and M. R. Jacobson (1992) In *Biological Nitrogen Fixation* (G. Stacey, R. H. Burris, and H. J. Evans, eds.), Chapman and Hall, New York, pp. 763–834.
31. P. E. Bishop, D. M. L. Jarlenski, and D. R. Hetherington (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7342–7346.
32. R. L. Robson, R. R. Eady, T. H. Richardson, R. W. Miller, M. Hawkins, and J. R. Postgate (1986) *Nature* **322**, 388–390.
33. J. R. Chisnell, R. Premakumar, and P. E. Bishop (1988) *J. Bacteriol.* **170**, 27–33.
34. M. R. Jacobson, R. Premakumar, and P. E. Bishop (1986) *J. Bacteriol.* **167**, 480–486.
35. F. Luque and R. N. Pau (1991) *Mol. Gen. Genet.* **227**, 481–487.
36. R. R. Eady (1996) *Chem. Rev.* **96**, 3013–3030.
37. D. Nieva-Gomez, G. P. Roberts, S. Klevickis, and W. J. Brill (1980) *Proc. Natl. Acad. Sci. USA*, **77**, 2555–2558.
38. S. Hill and E. P. Kavanagh (1980) *J. Bacteriol.* **141**, 470–475.
39. K. S. Howard, P. A. McLean, F. B. Hansen, P. V. Lemley, K. S. Koblan, and W. H. Orme-Johnson (1986) *J. Biol. Chem.* **261**, 772–778.
40. K. E. Rudd, H. J. Sofia, E. V. Koonin, G. Plunkett, S. Lazar, and P. E. Rouviere (1995) *Trends Biochem. Sci.* **20**, 12–14.
41. D. R. Dean, J. T. Bolin, and L. Zheng (1993) *J. Bacteriol.* **175**, 6737–6744.
42. W. A. Filler, R. M. Kemp, J. C. Ng, T. R. Hawkes, R. A. Dixon, and B. E. Smith (1986) *Eur. J. Biochem.* **160**, 371–377.
43. A. C. Robinson, D. R. Dean, and B. K. Burgess (1987) *J. Biol. Chem.* **262**, 14327–14332.
44. N. Gavini and B. K. Burgess (1992) *J. Biol. Chem.* **267**, 21179–21186.
45. D. Wolle, D. R. Dean, and J. B. Howard (1992) *Science* **258**, 992–995.
46. R. A. Ugalde, J. Imperial, V. K. Shah, and W. J. Brill (1984) *J. Bacteriol.* **159**, 888–893.
47. K. E. Brigle, M. C. Weiss, W. E. Newton, and D. R. Dean (1987) *J. Bacteriol.* **169**, 1547–1553.
48. T. D. Paustian, V. K. Shah, and G. P. Roberts (1989) *Proc. Natl. Acad. Sci. USA* **86**, 6082–6086.

49. V. K. Shah, J. R. Allen, N. J. Spangler, and P. W. Ludden (1994) *J. Biol. Chem.* **269**, 1154–1158
50. J. T. Roll, V. K. Shah, D. R. Dean, and G. P. Roberts (1995) *J. Biol. Chem.* **270**, 4432–4437.
51. T. R. Hoover, A. D. Robertson, R. L. Cerny, R. N. Hayes, J. Imperial, V. K. Shah, and P. W. Ludden (1987) *Nature* **329**, 855–857.
52. L. Zheng, R. H. White, and D. R. Dean (1997) *J. Bacteriol.* **179**, 5963–5966.
53. J. Imperial, R. A. Ugalde, V. K. Shah, and W. J. Brill (1984) *J. Bacteriol.* **158**, 187–194.
54. M. R. Jacobson, V. L. Cash, M. C. Weiss, N. F. Laird, W. E. Newton, and D. R. Dean (1989) *Mol. Gen. Genet.* **219**, 49–57.
55. T. C. White, G. S. Harris, and W. H. Orme-Johnson (1992) *J. Biol. Chem.* **267**, 24007–24016.
56. M. J. Homer, T. D. Paustian, V. K. Shah, and G. P. Roberts (1993) *J. Bacteriol.* **175**, 4907–4910.
57. M. J. Homer, D. R. Dean, and G. P. Roberts (1995) *J. Biol. Chem.* **270**, 24745–24752.
58. L. Zheng and D. R. Dean (1994) *J. Biol. Chem.* **269**, 18723–18726.
59. W. Fu, R. F. Jack, T. V. Morgan, D. R. Dean, and M. K. Johnson (1994) *Biochemistry* **33**, 13455–13463.
60. D. Flint (1996) *J. Biol. Chem.* **271**, 16068–16074.
61. W. Arnold, A. Rump, W. Klipp, U. B. Priefer, and A. Puhler (1988) *J. Mol. Biol.* **203**, 715–73862.
62. S. Kustu, E. Santero, D. Popham, and J. Keener (1989) *Microbiol. Rev.* **53**, 367–376.
63. B. T. Nixon, C. W. Ronson, and F. M. Ausubel (1986) *Proc. Natl. Acad. Sci. USA* **83**, 7850–7854.
64. D. L. Popham, D. Szeto, J. Keener, and S. Kustu (1989) *Science* **243**, 629–635.
65. M. Drummond, P. Whitty, and J. Wooton (1986) *EMBO J.* **5**, 441–447.
66. W. J. Buikema, W. W. Szeto, P. V. Lemley, W. H. Orme-Johnson, and F. M. Ausubel (1985) *Nucleic Acids Res.* **13**, 4539–4555.
67. J. Beynon, M. Cannon, V. Buchanan-Wollaston, and F. Cannon (1983) *Cell* **34**, 665–671.
68. S. Hill, S. Austin, T. Eydmann, T. Jones, and R. Dixon (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2143–2148.
69. R. G. Lowery and P. W. Ludden (1988) *J. Biol. Chem.* **263**, 16714–16719.
70. L. L. Saari, E. W. Triplett, and P. W. Ludden (1984) *J. Biol. Chem.* **259**, 15502–15508.
71. J. W. Golden, S. J. Robinson, and R. Haselkorn (1985) *Nature* **314**, 419–423.
72. J. Batut, M.-L. Daveran-Mingot, M. David, J. Jacobs, A. M. Garnerone, and D. Kahm (1989) *EMBO J.* **8**, 1279–1286.
73. H.-M. Fischer, T. Bruderer, and H. Hennecke (1988) *Nucleic Acids Res.* **16**, 2207–2224.

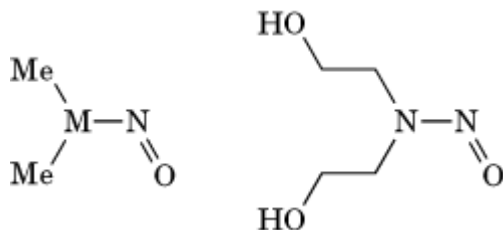
## Nitrosamines

Nitrosamines (*N*-nitrosoamines), typified by *N*-nitrosodimethylamine and *N*-nitrosodiethanolamine (Fig. 1), are **mutagens** and are reactive compounds that readily alkylate DNA via an  $S_N1$  reaction (slow formation of an alkylcarbenium ion and rapid reaction of this with nitrogen, oxygen, and



phosphorus atoms on DNA). *N*-Nitrosoamines are widespread in the environment and can be formed endogenously in humans (1-3). For example, workers in the rubber industry are exposed to high levels of *N*-nitrosodiethanolamine, which is a potent carcinogen over a considerable dose range and in more than one species. A wide range of volatile *N*-nitrosoamines are found in fried bacon and other pork products, while a different range of *N*-nitrosoamines are present in tobacco products and believed to play a significant role in human cancers associated with tobacco use. *N*-Nitrosoamines are potent carcinogens in every animal species tested; and although their role in human cancer is not definitively proved, it is considered highly likely that they are human carcinogens (2, 4).

**Figure 1.** Structures of the nitrosamines *N*-nitrosodimethylamine and *N*-nitrosodiethanolamine.



One of the most extensively studied compounds in this class is *N*-nitrosodimethylamine (NDMA). This chemical is nonreactive per se, but can be metabolized to an  $S_N1$ -type methylating agent, primarily by [cytochrome P450 IIE1](#). Although the metabolite reacts primarily to form around 60%  $N_7$ -methylguanine DNA adducts, it also generates a significant amount (approximately 6%) of  $O_6$ -methylguanine, as well as other minor *N*- and *O*-alkylated bases.  $O_6$ -Methylguanine and other *O*-alkylated adducts are directly miscoding, and they introduce base pair substitutions upon DNA replication; extensive evidence exists to link these lesions with mutagenesis and carcinogenesis by NDMA and related agents. NDMA has been described as the most powerful methylating agent known (5).

#### Bibliography

1. H. Bartsch and R. Montesano (1984) *Carcinogenesis* **5**, 1381–1393.
2. W. Lijinsky (1992) *Chemistry and Biology of *N*-Nitroso Compounds*, Cambridge University Press, Cambridge, England.
3. H. Bartsch and B. Spiegelhalder (1996) *Eur. J. Cancer Prev.* **5** (Suppl. 1), 11–18.
4. P. I. Reed (1996) *Eur. J. Cancer Prev.* **5** (Suppl. 1), 137–148.
5. V. L. Souliotis, S. K. Chabra, L. M. Anderson, and S. A. Kyrtopoulos (1995) *Carcinogenesis* **16**, 2381–2387.

#### NMR (Nuclear Magnetic Resonance)

Nuclear magnetic resonance (NMR) refers to a broad range of spectroscopic experiments intended to determine and analyze the absorption and emission of radiation by the nuclei of certain atoms. The radiation involved is in the radiofrequency range (RF) of the electromagnetic spectrum, typically

between 100 and 1000 MHz; in a strong magnetic field, it produces transitions among the allowed orientational states of nuclear magnetic moments. The exact frequencies are exquisitely sensitive to molecular structure and dynamics. Therefore, NMR experiments can provide information regarding three-dimensional structures, the rates of conformational change, and interactions between molecules. NMR is nondestructive, and samples can be recovered unchanged at the end of an experiment. NMR measurements may be made with samples in any physical state (solid, liquid, gas), although, given the dependence of the NMR phenomenon on molecular mobility, the kinds of experiments done and the information that can be obtained from them are different for different states.

## 1. Basic NMR

Atomic nuclei must possess a certain property, called angular momentum or “spin,” to be observable in an NMR experiment. At least one isotope exists of every known element that has this spin property, although this isotope may not be the most prevalent form of the element. Most experiments with biological materials use so-called spin 1/2 nuclei. Table 1 gives some of the characteristics of these nuclei. Notice the low natural abundance of some. Most carbon on the Earth's surface is the isotope  $^{12}\text{C}$ , but this isotope has zero nuclear spin and is not observable in any NMR experiment. The isotope  $^{13}\text{C}$ , however, is a spin 1/2 nucleus, and carbon NMR spectroscopy uses this isotope. There is a 1.1% chance of finding a  $^{13}\text{C}$  nucleus at any given carbon position in a molecule, although the amount of the isotope can be increased to nearly 100% by suitable synthesis using labeled precursors. Similarly, most nitrogen in biological materials is  $^{14}\text{N}$ , and the amount of  $^{15}\text{N}$  at natural abundance is low but can be increased by synthesis using isotopically enriched precursor molecules.

**Table 1. Properties of Spin 1/2 Nuclei**

| Isotope         | Natural Abundance (%) | NMR Frequency (MHz) for Magnetic Field Strength |         |         |
|-----------------|-----------------------|---|---------|---------|
|                 |                       | 11.75 T   | 17.62 T | 21.15 T |
| $^1\text{H}$    | 99.985                | 500.1   | 750.1   | 900.0   |
| $^{13}\text{C}$ | 1.108                 | 125.8   | 188.7   | 226.32  |
| $^{15}\text{N}$ | 0.37                  | 50.67   | 76.01   | 91.19   |
| $^{31}\text{P}$ | 100.0                 | 202.4   | 303.7   | 364.2   |
| $^{19}\text{F}$ | 100.0                 | 470.5   | 705.7   | 846.6   |

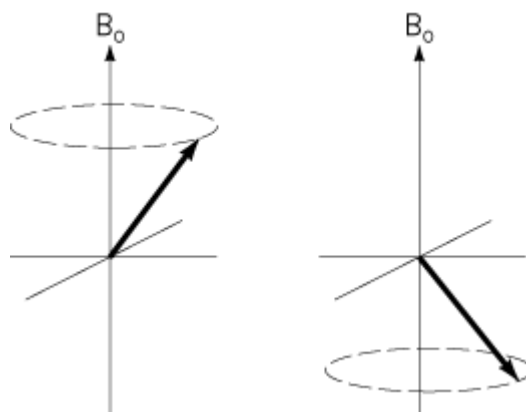
When a nucleus that has spin is placed into a magnetic field, its energy depends on its orientation in the field. A quantum mechanical description of a spin 1/2 nucleus in a magnetic field shows that the nucleus must take up one of two possible energy states, and that the difference between the energies of these states depends on the magnitude of the magnetic field. It is the transition from one of these states to the other, concomitant with the absorption or emission of a photon, that is detected in an NMR experiment.

Table 1 lists the approximate frequencies for NMR transitions of selected nuclei with various magnetic fields that are used for experiments with biological molecules. The exact frequency for any

particular nucleus in a molecule will depend on the molecular structure, as reflected by the chemical shift and any scalar (J) coupling that may be present. The basic goal of an NMR experiment is to determine the frequency or frequencies that can be associated with changes of energy state for the nucleus being studied. (See [Chemical Shift](#), [Scalar Coupling](#).)

A useful, although not completely rigorous, model of spin behavior depicts the spin (angular momentum) of a nucleus as a vector. In the case of spin 1/2 nuclei in a magnetic field, the vector has two possible orientations (“up” and “down”; Fig. 1), corresponding to the two allowed energy states. The spin vectors rotate around the direction of the laboratory magnetic field while keeping constant their angle relative to the field. As the spin vector rotates around the field in this manner, its energy stays constant. This rotation is termed Larmor precession. The frequency of precession is the same as the frequency of the radiation involved in transition of the nuclear spin from one allowed state to the other. If the laboratory field is 17.62 T, the frequency of Larmor precession for the nuclear spin vector of the hydrogen ( $^1\text{H}$ ) nucleus is  $7.5 \times 10^8$  revolutions per second.

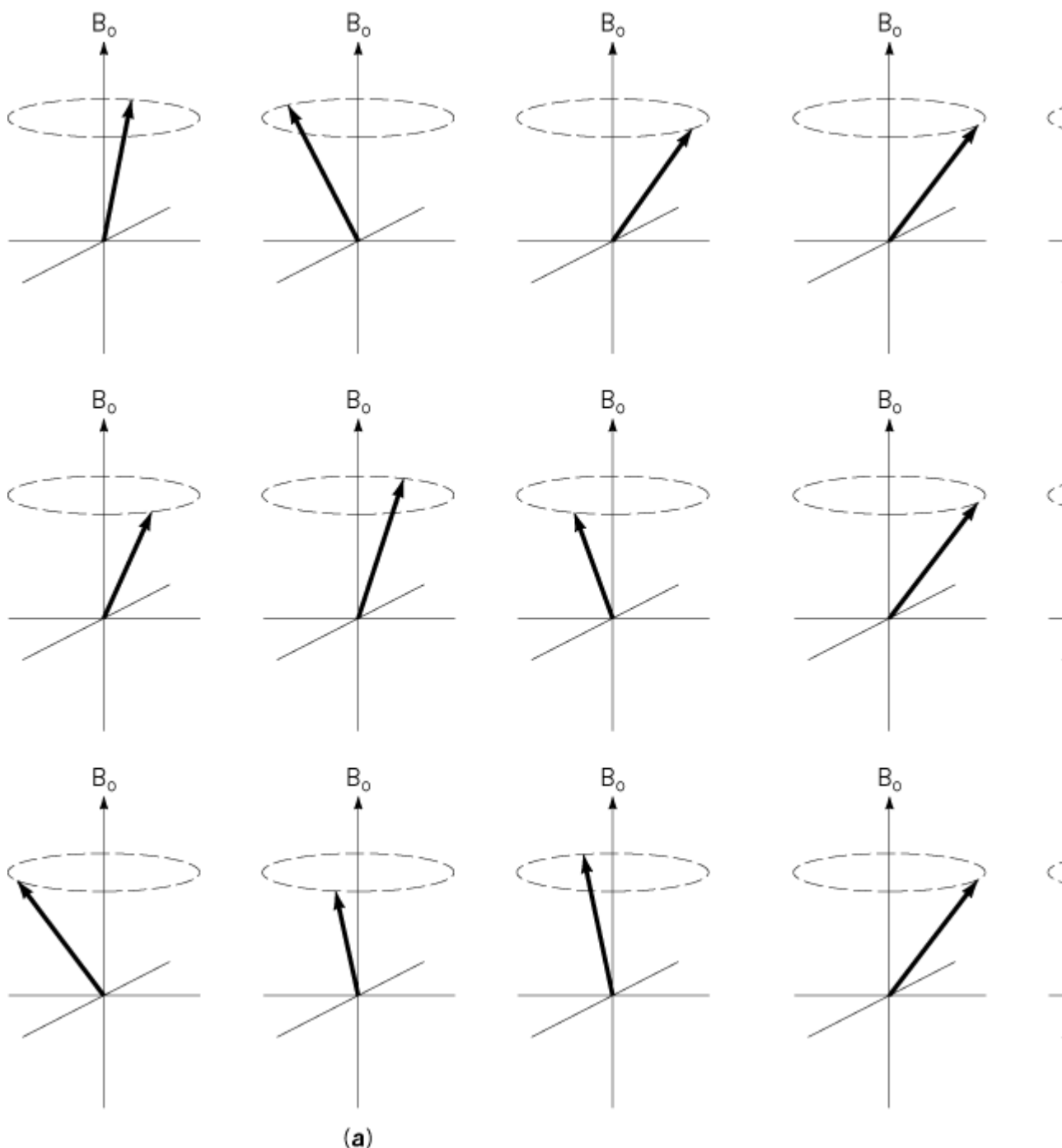
**Figure 1.** The behavior of a spin 1/2 nucleus in a magnetic field ( $B_0$ ). Such nuclei are allowed to take one of two possible orientations relative to the field. The energies of these are different. In either state, the spin undergoes Larmor precession around the direction of the field.



NMR experiments with biological molecules typically use 20–600 mL of a solution in which the material of interest is present at a concentration of 0.5 to 5.0 mM. Under these conditions, a large number ( $\sim 10^{18}$ ) of nuclei of any specific type will be present. When the sample is placed into a magnetic field, each nucleus of the sample with the spin property will be sorted into one of the allowed spin energy states and will undergo Larmor precession. The number of nuclei of a given type that is present in each particular allowed spin state is governed by the Boltzmann distribution law. Every molecule in the sample is acted on by seemingly random dynamic processes that exchange energy among the molecules of the sample, driving the system toward equilibrium. These processes are collectively known as relaxation. Any experimental operation that alters the state populations will be counteracted by these relaxation processes, so that the sample eventually returns to the Boltzmann populations that are characteristic of equilibrium.

Another feature of the equilibrium state is a random distribution of spins in their precessional motion. All of the spin-up vectors shown in Fig. 2a have the same component along the direction of the laboratory magnetic field, but their positions around the precessional path are random. (A similar situation holds for the spin-down vectors.)

**Figure 2.** Precessional behavior of spin 1/2 nuclei in a magnetic field. (a) Noncoherent precessional behavior of a collection of spins. Nuclei that have the other allowed orientation relative to  $B_0$  undergo precession at the same time that spins in the orientation illustrated are coherent.



An NMR spectrometer is only able to report the collective behavior of the nuclei of the sample, not the behavior of any single nucleus. It can produce signals only when there is a net macroscopic magnetic field from the sample that is at right angles (transverse) to the direction of the laboratory magnetic field. Although all spins of a particular type have the same Larmor precessional frequency when the sample is placed in the magnetic field, there is no defined starting point for Larmor precession (Fig. 2a). In this case, the individual transverse components of each spin vector collectively cancel, while the longitudinal components of the spins of the sample combine to create a macroscopic magnetic field component that is along the direction of the laboratory field

(longitudinal to it). To generate a signal in an NMR spectrometer clearly requires that the equilibrium state of the sample in the laboratory magnetic field be altered.

To change the sample spins away from the equilibrium state requires the absorption or emission of photons, as the nuclei of the sample change spin states. The source of the needed radiation, the transmitter of the NMR instrument, typically is switched on for only a short time (1–20  $\mu$ s). The energy of this RF pulse can alter the populations of nuclei in the allowed energy states. RF pulses of the appropriate strength and duration can also create coherence in the motions of individual spins (Fig. 2b). The resultant magnetic field produced when sample spins are precessing coherently, as suggested by Fig. 2b, has a macroscopic component in the plane that is transverse to the direction of the magnetic field—precisely what is needed to afford a detectable NMR signal. State populations altered by an RF pulse are eventually restored to their equilibrium values by relaxation processes. These processes also destroy any coherence of spin precession that has been created, and the NMR signal decays back to zero.

Except in special circumstances, RF pulses are nonselective—that is, if the radiofrequency energy source is pulsed on, all spins in the sample that have Larmor frequencies close to that frequency will respond to the pulse. If the RF pulse is at the proton Larmor frequency, all protons of the sample—those in methyl groups, those in aromatic rings, those in water molecules—will be affected by the pulse. If a pulse is designed to create coherence in proton spin motions, all spins in the sample will become coherent and, after the pulse, will contribute to the proton NMR signal detected by the spectrometer.

Spin coherence leads to NMR signals that oscillate at the Larmor frequency of the nuclei. If there are many coherences present (recognizable by their specific Larmor precessional frequencies), the detection system generates a signal that is the sum of the signals created by all coherences present. The signal is digitized and stored in the instrument computer as a series of  $10^3$ – $10^5$  numbers. Although the signal recorded is a decaying function of time, it is the frequencies present in this signal that is of interest. A Fourier transformation of the numbers representing the signal is performed to obtain these frequencies and their relative intensities.

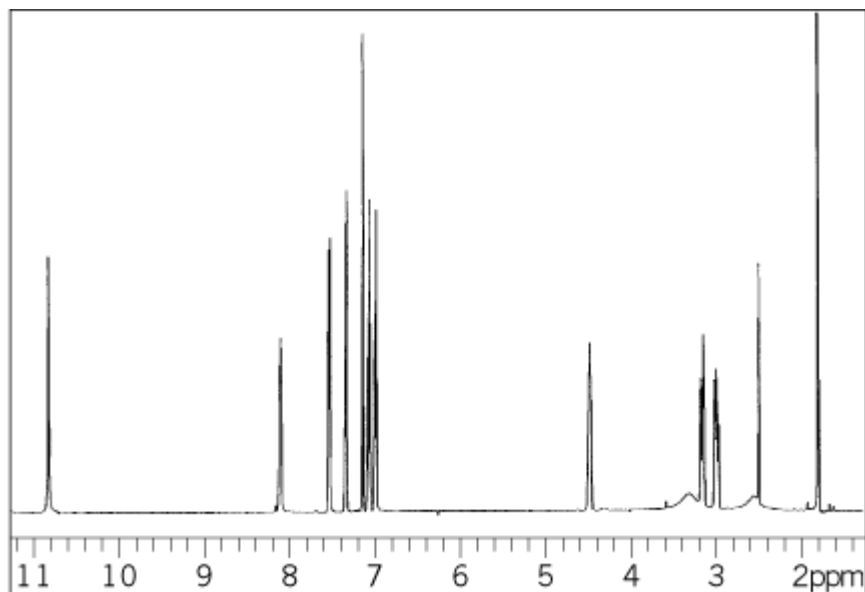
NMR signals appear at different frequencies because of [chemical shifts](#) and the effects of **scalar** (spin) **couplings**. The frequency axis of an NMR spectrum is thus usually calibrated in terms of the NMR shielding parameter, using ppm units. Because it is only practical to measure *differences* in shielding parameters, one shielding parameter (that of a signal chosen as reference) is arbitrarily set to 0. A number on a frequency axis of an NMR spectrum therefore represents the difference between the shielding parameter for the nucleus of interest and that for the reference nucleus.

## 2. Multidimensional NMR Spectra

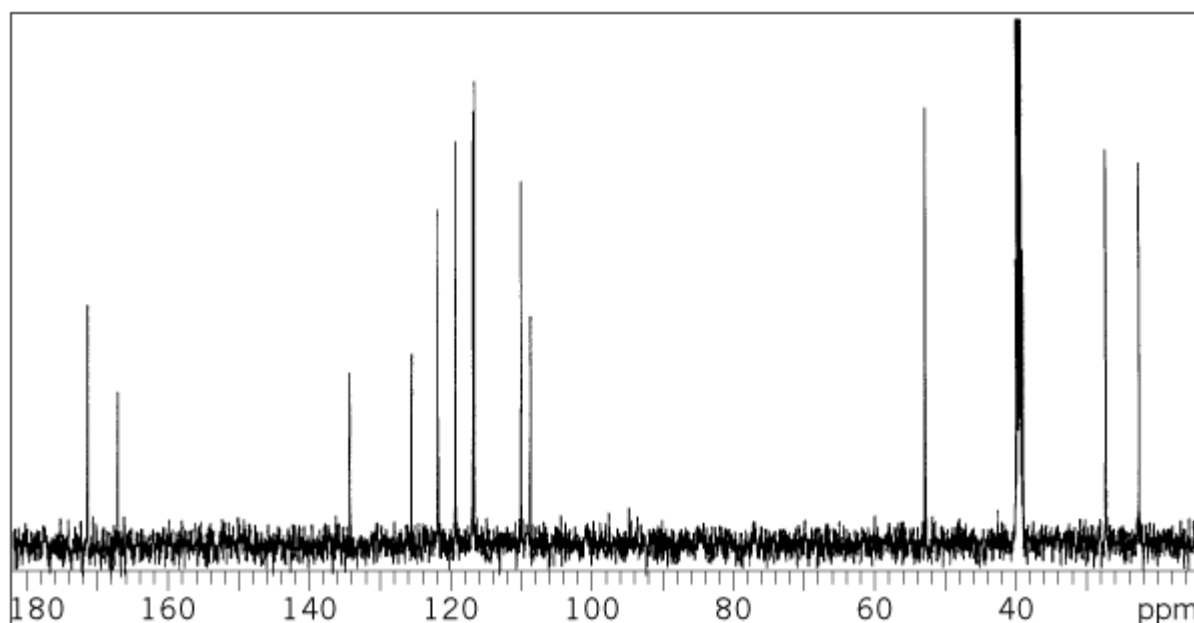
Dimension in NMR spectroscopy refers to the number of frequency or chemical shift axes that are needed to represent the results of an experiment. The most basic proton or  $^{13}\text{C}$  spectrum is a one-dimensional (1D) experiment: the intensities of spectral peaks are plotted as a function of one chemical shift axis. Figures 3 and 4 show 1D proton and  $^{13}\text{C}$  NMR spectra of an amino acid derivative. A close correspondence exists between the collections of peaks present in the spectrum and the chemical structure of the molecule examined. An experienced spectroscopist can assign these peaks to specific nuclei (hydrogens and carbons) of the molecule using any of several approaches, including analysis of the fine structure of the peaks due to scalar coupling.

**Figure 3.** The proton NMR spectrum of *N*-acetyltryptophan obtained under conditions where the proton Larmor frequency is near 500 MHz. The sample was dissolved in  $\text{d}_6$ -dimethylsulfoxide. The signal with the smallest shielding parameter (near 10.8 ppm) is readily assigned to the N-H proton of the indole ring. The signal at 2.5 ppm arises from

residual hydrogen atoms in the solvent.



**Figure 4.** A  $^{13}\text{C}$  NMR spectrum of *N*-acetyltryptophan obtained with the same sample used to acquire the proton spectrum shown in Figure 3. A signal exists for each carbon atom in the molecule. The collection of signals at 39.5 ppm is due to the  $^{13}\text{C}$  atoms in the solvent.



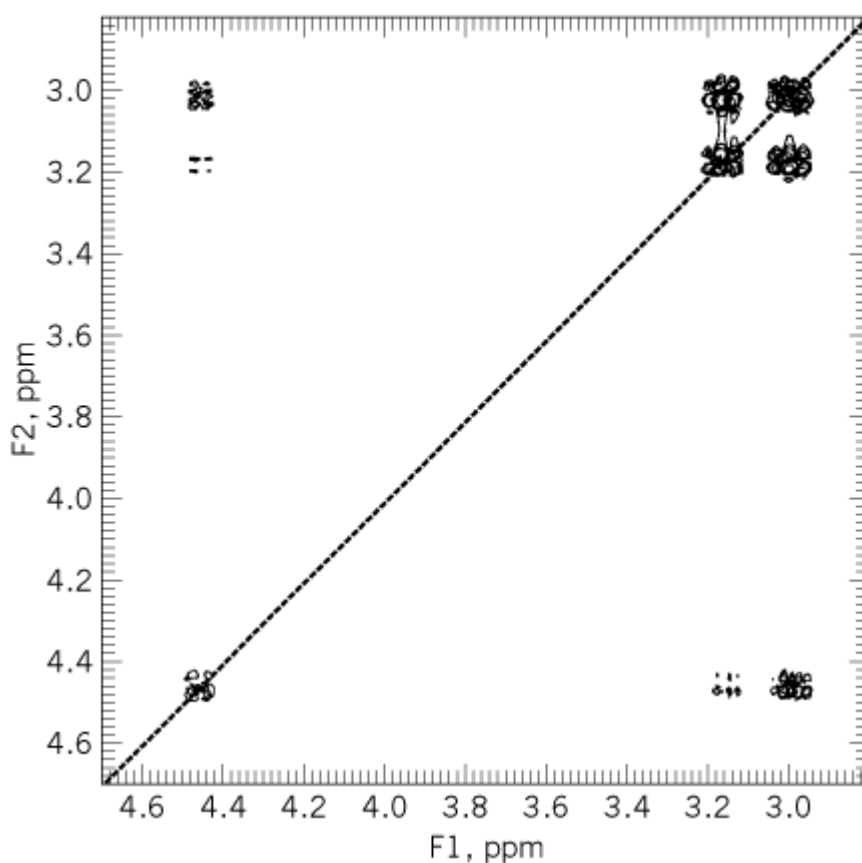
More elaborate NMR experiments can be designed in which observed intensities are dependent on two or more chemical shifts. A spectrum that involves two chemical shift axes is a two-dimensional (2D) spectrum. Reflection will show that a 2D NMR spectrum is actually a three-dimensional object with the two chemical-shift axes defining two dimensions of the object and intensity represented along the third. The most common way that 2D NMR spectra are presented is by use of a contour plot. The two chemical-shift axes appear in the plane of the page, while intensity is represented by a

series of contours that correspond to levels of intensity, in the same way that contour lines on a topographic map represent different elevations. Intensity contours in 2D experiments can be either positive or negative and therefore appear above or below the plane. One technique for indicating the sign of the intensity is to plot positive intensities with solid (continuous) lines, whereas negative intensities are plotted with dashed contours. Alternatively, contours corresponding to positive and negative intensity can be plotted in different colors.

If both chemical-shift axes of a 2D spectrum refer to the Larmor frequencies of the same nucleus, the spectrum is referred to as a homonuclear 2D experiment. Most commonly, homonuclear 2D experiments are proton-proton. In a heteronuclear 2D spectrum, the shift axes correspond to the Larmor frequencies of different types of nuclei with, perhaps,  $^{13}\text{C}$  frequencies appearing along one axis and proton frequencies along the other.

The essential features of a homonuclear 2D spectrum are abstracted in Figure 5. In this case, the range of frequencies along the axis labeled  $f_1$  is identical to the range of frequencies plotted along the other axis,  $f_2$ . There are always features along the  $45^\circ$  diagonal of a 2D map representing a homonuclear experiment. These peaks have the same coordinates (for example, S1,S1) along either axis. The coordinates of a diagonal peak correspond to the chemical shift of a peak or multiplet in the 1D spectrum of the same sample. There could be much fine structure of the diagonal peaks because of scalar coupling, but this might be obscured by the limited resolution of the manner the data are collected and displayed.

**Figure 5.** Part of a typical homonuclear proton-proton 2D NMR spectrum, using the same *N*-acetyltryptophan sample and experimental conditions that produced the 1D spectrum shown in Figure 3. This spectrum was produced by a DQF-COSY experiment. The cross peaks show that the aliphatic protons of the amino acid (with shielding parameters 2.98, 3.14, 4.45 ppm) are scalar coupled to each other.



Depending on the nature of the experiment and the structure of the molecule being studied, peaks that are off the diagonal at some coordinate (S1,S2) or (S2,S1) may appear in a 2D spectrum. These are called cross peaks, and their significance depends on the nature of the 2D experiment. Their presence, however, indicates that some relationship exists between the spins characterized by the S1 shift and spins that have the S2 shift. Figure 5 shows part of the homonuclear 2D spectrum that is obtained in a 2D NMR experiment with the same material used to obtain the 1D spectrum in Figure 3. For this particular spectrum (from a proton double quantum filtered Correlation Spectroscopy [DQFCOSY] experiment), cross peaks are present because there is scalar coupling between hydrogen atoms of the molecule. (see [COSY Spectrum](#).)

For heteronuclear experiments, only cross peaks are present in the 2D spectrum. Again, the presence of a cross peak indicates that some structural relationship exists between the spins characterized by the S1 shift along one axis and spins that have the S2 shift along the other.

NMR experiments that provide results of higher dimensionality (for example, 3D, 4D, or 5D) exist and are often applied to structural studies of biological **macromolecules**. NMR spectra of higher dimensionality than 2D experiments are impossible to represent completely on a single planar surface. These spectra are viewed and analyzed by means of computer displays, which may use color-coding to convey the intensity of an NMR absorption or emission line. Although a few 3D homonuclear experiments are done, most 3D and higher-dimension NMR experiments with biological systems are heteronuclear, usually performed with materials that have carbon or nitrogen positions that are enriched in  $^{13}\text{C}$  and  $^{15}\text{N}$ , respectively. (See Chemical shift, Scalar coupling, COSY spectrum, NOESY spectrum, ROESY spectrum, TOCSY spectrum, Magnetization transfer, Isotope editing, Triple resonance.)

#### Suggestions for Further Reading

##### Basic NMR

A. Bax (1982) *Two-Dimensional Nuclear Magnetic Resonance in Liquids*, Kluwer, Boston.

F. A. Bovey (1988) *Nuclear Magnetic Resonance Spectroscopy*, Academic, San Diego.

W. S. Brey, ed. (1988) *Pulse Methods in 1D and 2D Liquid-phase NMR* Academic, San Diego.

D. Canet (1996) *Nuclear Magnetic Resonance. Concepts and Methods*, Wiley, Chichester.

W. R. Croasmun and R. M. K. Carlson, eds. (1994) *Two-dimensional NMR Spectroscopy: applications for chemists and biochemists*, 2nd ed., V.C.H., New York.

J. D. De Certaines, W. M. M. J. Bovee, and F. Podo (1992) *Magnetic Resonance in Biology and Medicine*, Pergamon, New York.

R. K. Harris (1983) *Nuclear Magnetic Resonance Spectroscopy: a physicochemical view*, Pitman, Marshfield, Mass.

S. W. Homans (1992) *A Dictionary of Concepts in NMR*, Clarendon, Oxford.

R. Kitamaru (1990) *Nuclear Magnetic Resonance: principles and theory*, Elsevier, New York.

R. S. Macomber (1998) *A Complete Introduction to Modern NMR Spectroscopy*, Wiley, New York.

D. L. Turner (1985) *Prog. NMR Spectrosc.* **17**, 281–358.

C. H. Yoder and C. D. Schaeffer Jr. (1987) *Introduction to Multinuclear NMR*, Benjamin/Cummings, Menlo Park.

The journal *Concepts in NMR* (Wiley) provides didactic articles on many aspects of NMR spectroscopy.

##### NMR of Biological Systems



N. Beckman (1995) *Carbon-13 NMR Spectroscopy of Biological Systems*, Academic, San Diego.  
J. Cavanagh, W. J. Fairbrother, A. G. Palmer III, and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.  
G. M. Clore and A. M. Gronenborn, eds. (1993) *NMR of Proteins* CRC Press, Boca Raton.  
S. M. Cohen (1987) *N. Y. Acad. Sci.* **508** (entire volume).  
D. J. Craik, ed. (1996) *NMR in Drug Design* CRC Press, Boca Raton.  
J. N. S. Evans (1995) *Biomolecular NMR Spectroscopy*, Oxford, New York.  
R. J. Gillies (1994) *NMR in Physiology and Biomedicine*, Academic, San Diego.  
C. Ioannidou (1994) *NMR of Biological Macromolecules*, Springer-Verlag, New York.  
J. L. Markeley and S. J. Opella, eds. (1997) *Biological NMR Spectroscopy* Oxford, New York.  
P. G. Morris (1986) *Nuclear Magnetic Resonance Imaging in Medicine and Biology*, Clarendon Press, New York.  
D. G. Reid, ed. (1997) *Protein NMR Techniques Humana*, Totowa, New Jersey.  
G. K. C. Roberts, ed. (1993) *NMR of Macromolecules: a practical approach* Oxford, New York.  
K. Wuthrich (1976) *NMR in Biological Research: peptides and proteins*, American Elsevier, New York.

*Methods in Enzymology* **176**, **177** (1989) and **239** (1994) have many useful articles on deuteration and other isotope labeling methods, sample preparation, assignment methods, instrumental methods, distance geometry, molecular dynamics studies and simulations, phosphorus NMR, ligand binding, exchange effects, in vivo enzyme activity measurements, multiple quantum methods, total correlation spectroscopy (TOCSY), rotating frame Overhauser effect spectroscopy (ROESY), nuclear Overhauser effect spectroscopy (NOESY) and COSY.

## NOESY Spectrum

Proton-proton [nuclear overhauser effect](#) spectroscopy (NOESY) is a type of two-dimensional (or higher) nuclear magnetic resonance ([NMR](#)) spectrum. In it, cross peaks appear because the spins of atoms are sufficiently close to each other that their mutual interaction, through relaxation, can produce changes in the populations of the nuclear spin energy levels that are associated with those spins. The alterations of the level populations lead to changes in NMR signal intensities (NOEs). The magnitude of the effect depends on the distance between the spins; consequently, observation and calibration of NOESY cross peaks provide constraints on what can be the [tertiary structure](#) of the molecule under study.

Cross peaks in a NOESY experiment arise because magnetization that precesses at a frequency corresponding to one [chemical shift](#) during an initial phase of the experiment is transferred to a nearby spin, where it precesses at a different chemical shift during a latter phase of the experiment. The chemical-shift coordinates of a cross peak identify the shifts of the interacting sets of spins. The distance dependence of the cross-peak intensity is very great, in the simplest case varying with  $1/r^6$ , where  $r$  is the distance between the spins.

The distance between hydrogen atoms (protons) in a methylene group is about 0.179 nm, whereas the distance between adjacent protons attached to an aromatic ring is 0.248 nm. These distances are fixed by the covalent structure of these groups. Measurement of the NOEs between protons in these situations provides a possible means of calibrating NOESY cross-peak intensities and thereby

potentially obtaining estimates of other internuclear distances.

The usefulness of NOESY experiments is diminished by spin diffusion processes, in which the magnetization of a particular nucleus is transferred to nearby nuclei; this leads to NOE cross peaks between atoms that are not close in space but are linked to atoms that are. An experimental variable in NOESY-type experiments is the time during which the NOEs develop. Called the mixing time, it typically is set between 50 and 1000 ms. Spin diffusion is a time-dependent phenomenon, so the misleading cross peaks it produces can be minimized by reducing the mixing time. The heights of all cross peaks are altered to some extent by this procedure, however, and all cross-peak intensities may be perilously close to experimental noise if the mixing time is short.

The elements of NOESY experiments can be built into experiments that produce three-dimensional or higher NMR spectral results. (See also [ROESY Spectrum](#), [Distance Geometry](#), [Simulated Annealing](#).)

#### Suggestions for Further Reading

J. Cavanagh, W. J. Fairbrother, A. G. Palmer III, and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.

G. M. Clore and A. M. Gronenborn, eds. (1993) *NMR of Proteins* CRC Press, Boca Raton.

C. Ioannidou (1994) *NMR of Biological Macromolecules*, Springer-Verlag, New York.

D. Neuhaus and M. P. Williamson (1989) *The Nuclear Overhauser Effect in Structural and Conformational Analysis*, V.C.H., New York.

G. K. C. Roberts, ed. (1993) *NMR of Macromolecules: a practical approach* Oxford, New York.

K. Wuthrich (1976) *NMR in Biological Research: peptides and proteins*, American Elsevier, New York.

Methods in Enzymology 176, 177 (1989) and **239** (1994) have many useful articles on total correlation spectroscopy, (TOCSY), rotating frame Overhauser effect spectroscopy (ROESY), and NOESY.

## Noggin

*Noggin* is a **gene** expressed in Spemann's organizer of the *Xenopus* gastrula and its **homologous** structure, the node, of mammals ([1](#), [2](#)). It encodes a secreted **glycoprotein** that is a **disulfide**-linked dimer in solution ([3](#)). The protein has very high affinity (**dissociation constant** of 20 pM) for the bone morphogenetic proteins (BMPs) BMP2 and BMP4 and blocks their activity by preventing them from binding to their respective receptors, the type II BMP receptor and the type I (transducing) receptor, ALK3 ([4](#)). Noggin was the first neural-inducing molecule identified from embryos ([5](#)).

Noggin was discovered in a screen for dorsalizing activities in *Xenopus* ([1](#)). *Xenopus* embryos that are UV-irradiated shortly after fertilization develop no dorsal structures. To identify the molecules that induce axial structures, **messenger RNAs** were synthesized from pools of a **cDNA library** and injected into UV-treated embryos; those pools that induced axial structures were split, and sib-selection was used to identify the activity. Three mRNAs—*Xwnt-8*, *noggin*, and *Xnr3*—were identified in this screen ([1](#), [6](#), [7](#)). Zygotic transcripts of all three are expressed in the prospective mesoderm. Two of them, *noggin* and *Xnr3*, an unusual member of the **transforming growth factor- $\beta$**  family) are expressed in the dorsal mesoderm, Spemann's organizer. Noggin gained its name from its ability to transform the entire embryo into head structures.

The possible zygotic function of noggin was tested by applying Noggin [recombinant protein](#) to explants. When applied to prospective ventral mesoderm or epidermis, Noggin can mimic two of the known gastrula-stage activities of the organizer, namely, dorsalization of ventral mesoderm and neural induction ([3](#), [5](#)).

Neural induction by Noggin appears to be direct, in that no mesodermal intermediate is detected. Whereas mesoderm inducers can act only on blastula and early gastrula tissues, Noggin can induce neural tissue until late gastrulation. Although the treated explants are not homogeneous, only anterior types of tissue are present ([5](#)).

Expression of the mouse noggin gene is homologous to that of *Xenopus* in that noggin transcripts are found in the mouse organizer (node) and its axial mesodermal derivatives. *In situ* **hybridization** and a lacZ [reporter gene](#), integrated into and replacing the endogenous noggin-coding sequences, were used to detect later sites of noggin expression; these appear to be crucial ([2](#), [8](#)). Although the initial steps of neural induction and somite formation occur in a Noggin-deficient mouse, neural patterning and skeletal formation are abnormal. Noggin is required for ventral fates to be induced in the neural tube. During dorsal ventral patterning of the neural tube, BMPs have been shown by pharmacologic experiments to be required to induce dorsal fates ([9](#)), and sonic hedgehog has been shown from genetic experiments to be required also ([10](#)). In the absence of Noggin, excess BMP activity blocks ventral fates in the caudal neural tube. Thus, there is an increasing severity of ventralization in the neural tube, with the loss of interneurons, motor neurons, and finally the floor plate. These defects are mirrored in the mesoderm, where a rostral to caudal gradient of defects in the somites includes loss of sclerotomal and dermomyotomal derivatives. *In vitro* explant experiments show that noggin can synergize with sonic hedgehog to induce sclerotome, so endogenous BMPs are likely to be potent antagonists of somite differentiation.

Noggin has become a useful pharmacologic antagonist of BMPs. It has been used to block BMP signaling from lateral plate mesoderm in the chick as well as surface ectodermal and roofplate BMP signaling ([11-13](#)). Such experiments have established the importance of BMP signaling in not only suppressing epaxial muscle fates in the somite but also in the induction of lateral plate identities.

During skeletal development, Noggin also plays an essential role in regulating BMP activity. When BMP activity is not attenuated by Noggin, excess mesenchyme is recruited into the condensing cartilage; in addition, synovial joints do not form ([8](#)).

## Bibliography

1. W. C. Smith and R. M. Harland (1992) *Cell* **70**, 829–840.
2. J. A. McMahon et al. (1998) *Genes and Development* **15**, 1438–1452.
3. W. C. Smith, A. K. Knecht, M. Wu, and R. M. Harland (1993) *Nature* **361**, 547–549.
4. L. B. Zimmerman, J. M. De Jesus-Escobar, and R. M. Harland (1996) *Cell* **86**, 599–606.
5. T. M. Lamb et al. (1993) *Science* **262**, 713–718.
6. W. C. Smith and R. M. Harland (1991) *Cell* **67**, 753–765.
7. W. C. Smith, R. McKendry, S. Ribisi, and R. M. Harland (1995) *Cell* **82**, 37–46.
8. L. J. Brunet, J. A. McMahon, A. P. McMahon, and R. M. Harland (1998) *Science* **280**, 1455–1457.
9. K. F. Liem Jr., G. Tremml, and T. M. Jessell (1997) *Cell* **91**, 127–138.
10. C. Chiang et al. (1996) *Nature* **383**, 407–413.
11. E. Hirsinger et al. (1997) *Development* **124**, 4605–4614.
12. C. Marcelle, M. R. Stark, and M. Bronner-Fraser (1997) *Development* **124**, 3955–3963.
13. R. Reshef, M. Maroto, and A. B. Lassar (1998) *Genes Dev.* **12**, 290–303.

### Suggestion for Further Reading

14. R. M. Harland and J. C. Gerhart (1997) Formation and function of Spemann's organizer. *Ann. Rev. Cell Develop. Biol.* **13**, 611–667.

### Nonautonomous Controlling Element

“[Controlling elements](#)” was the term Barbara McClintock gave to [transposable elements](#) when she discovered them in maize in the late 1940s. This designation emphasized that these new genetic elements could control gene expression, as well as move from place to place in the [genome](#) (1). McClintock discovered two classes of controlling elements, one of which she designated *Dissociation* (Ds) and the other *activator* (Ac). She found that while Ds elements could move in the presence of Ac, they were unable to move alone, while Ac was capable of movement in the presence or absence of Ds. She hypothesized that activator made a product that could promote the movement of both Ac and Ds. We now know that this product is the Ac [transposase](#) and that Ds is a deletion derivative of Ac that contains the same special recombination sequences at its tips that the Ac transposase can act on, but is lacking a transposase. Because Ds lacks a transposase gene, it cannot promote its own movement and is dependent on this product from Ac. Ds is said to be a nonautonomous controlling element because it can't move in the absence of Ac, whereas Ac is said to be an [autonomous controlling element](#) because it encodes its own transposase that can act upon its terminal sequences to move the Ac element from place to place.

Nonautonomous elements are deleted versions of autonomous elements. Different nonautonomous versions of the same element result from different internal deletions. These deleted versions probably arise from incomplete repair events at the donor site after transposon excision (2-4); this repair involves transferring information from an intact homologue, using homologous [recombination](#) and gap repair.

In addition to Ac–Ds and several other sets of autonomous and nonautonomous elements in maize, autonomous and nonautonomous elements have been also observed in other organisms such as *Drosophila* (2), and *Caenorhabditis elegans* (3, 5). Virtually all bacterial elements that have been observed appear to be autonomous.

### Bibliography

1. B. McClintock (1956) Cold Spring Harbor Symp. Quant. Biol. **21**, 197–216.
2. W. R. Engels, D. M. Johnson-Schlitz, W. B. Eggleston, and J. Sved (1990) *Cell* **62**, 515–525.
3. E. Rubin and A. A. Levy (1997) *Mol. Cell. Biol.* **17**, 6294–6302.
4. E. S. Coen, T. P. Robbins, J. Almeida, A. Hudson, and R. Carpenter (1989) In *Mobile DNA* (D. A. Berg and M. M. Howe, eds.), American Society for Microbiology, Washington, DC, pp. 413–436.
5. R. H. Plasterk (1991) *EMBO J.* **10**, 1919–1925.

### Noncompetitive Inhibition

## 1. Linear

Linear noncompetitive inhibition of an [enzyme](#) occurs whenever an inhibitory analogue of the substrate combines with a form of enzyme other than the one with which the variable substrate combines and a reversible connection exists between the points of addition of inhibitor and the variable substrate. It should be noted that reversible connections are broken either by the release of product at zero concentration or the presence of a nonvaried substrate at an essentially infinite concentration. For a two-substrate ordered kinetic mechanism (Fig. 1), the inhibition by I would be linear competitive with respect of A. It would be linear noncompetitive with respect to B, as B combines with EA, I combines with E, and a reversible connection exists between E and EA. The general form of the equation to describe noncompetitive inhibition is shown in equation (1):

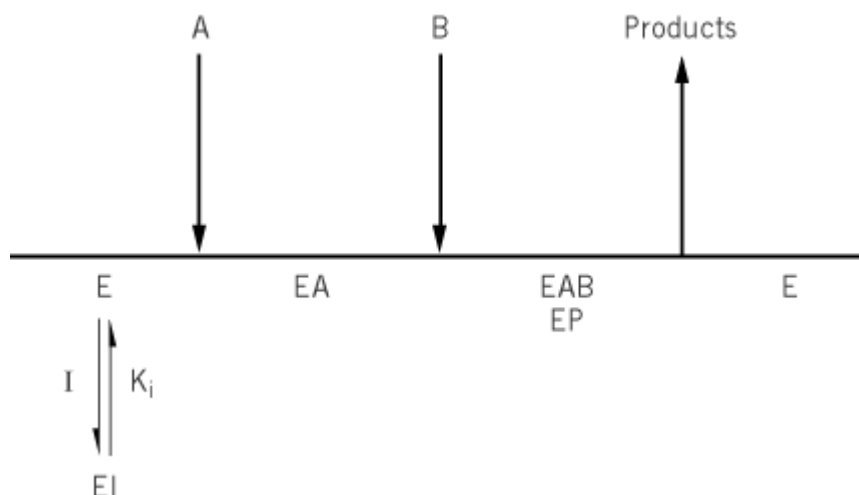
$$v = \frac{VA}{K_a(1 + I/K_{is}) + A(1 + I/K_{ii})} \quad (1)$$

where  $K_a$  denotes the  $K_m$  (Michaelis constant) for substrate A,  $V$  is the maximum velocity of the reaction, and  $K_{is}$  and  $K_{ii}$  are the inhibition constants associated with the slopes and intercepts, respectively, of the double-reciprocal plot, which has the form illustrated in equation (2):

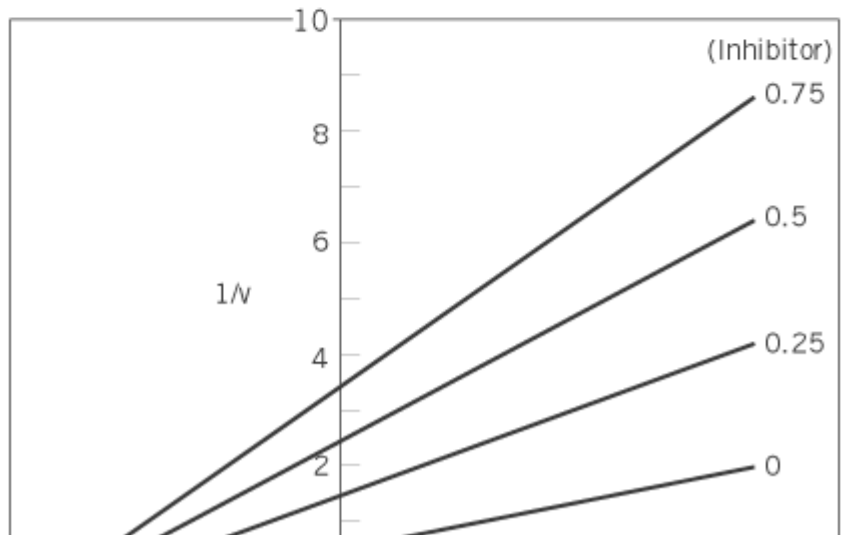
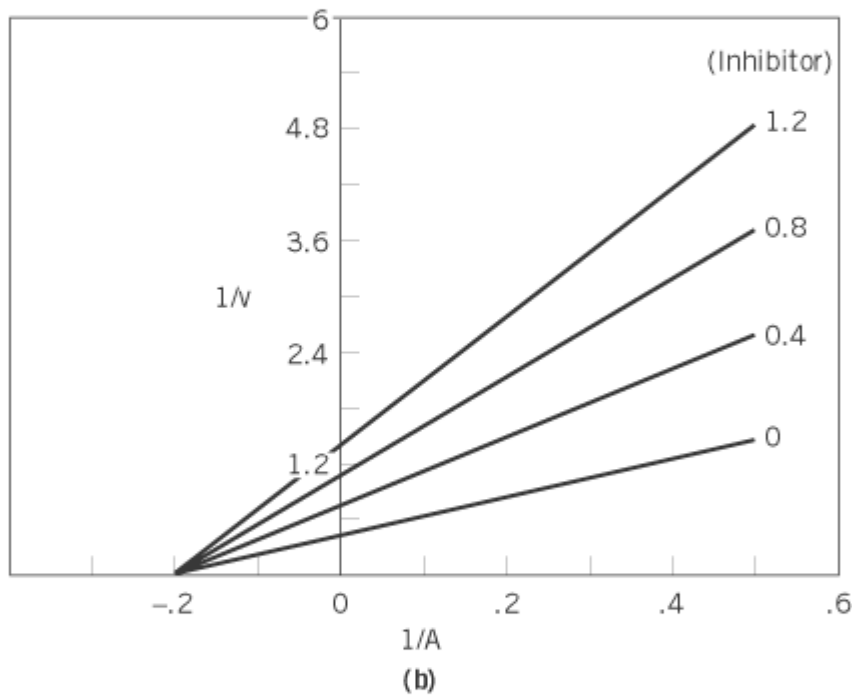
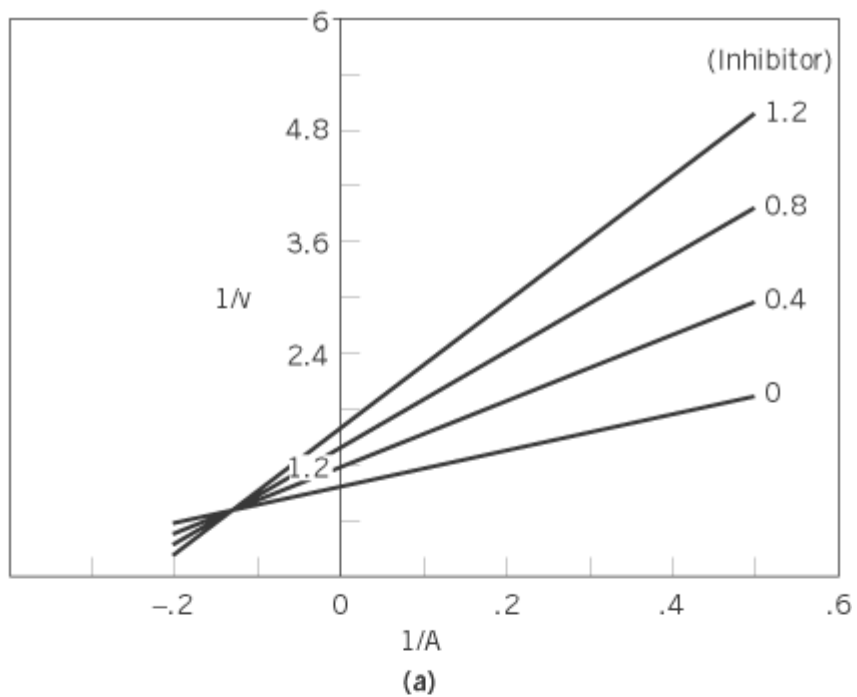
$$\frac{1}{v} = \frac{K_a}{V} \left(1 + \frac{I}{K_{is}}\right) \frac{1}{A} + \frac{1}{V} \left(1 + \frac{I}{K_{ii}}\right) \quad (2)$$

This equation shows that a plot of  $1/v$  against  $1/A$ , at different concentrations of  $I$ , would yield a family of straight lines that intersect at a common point to the left of the vertical ordinate. The  $1/v$  coordinate for the crossover point is  $(1/V)[1 - K_{is}/K_{ii}]$ , and thus, the intersection point can occur above, on, or below the abscissa, depending on whether  $K_{is}$  is less than, equal to, or greater than  $K_{ii}$ , respectively (Fig. 2). The term *noncompetitive* has been applied only for the case in which the lines of a double-reciprocal plot intersect on the abscissa; the term *mixed inhibition* is used to describe cases in which the lines intersect either above or below the abscissa. Such a distinction is unnecessary as the position of the crossover point is simply a function of the relative values of  $K_{is}$  and  $K_{ii}$ .

**Figure 1.** Kinetic mechanism for which I would give rise to linear competitive inhibition with respect to A and linear noncompetitive inhibition with respect to B.



**Figure 2.** Double-reciprocal plots for linear noncompetitive inhibition with crossover points **(a)** above, **(b)** on, and **(c)** below the abscissa. The respective values in arbitrary units for  $K_{is}$  and  $K_{ii}$  were: **(a)** 0.5 and 2.0, **(b)** 0.5 and 0.5, and **(c)** 0.3 and 0.1.



The inhibition constants associated with the slopes ( $K_{is}$ ) and intercepts ( $K_{ij}$ ) of a double-reciprocal plot are calculated and interpreted as described for [competitive inhibition](#) and [uncompetitive inhibition](#), respectively.

## 2. Hyperbolic

This type of inhibition would be observed for the scheme illustrated under hyperbolic [competitive inhibition](#), when the EA and EAI complexes give rise to product at different rates. The  $V_A$  term of the equation that describes this inhibition would be modified by the factor,  $1+I/K_{in}$ , so as to give equation (3):

$$v = \frac{VA \left(1 + \frac{I}{K_{in}}\right)}{K_a \left(1 + \frac{I}{K_{in}}\right) + A \left(1 + \frac{I}{K_{id}}\right)} \quad (3)$$

which, in reciprocal form, is described by equation (4):

$$\frac{1}{v} = \frac{K_a}{V} \left[ \frac{1 + \frac{I}{K_{in}}}{1 + \frac{I}{K_{in}}} \right] \frac{1}{A} + \frac{1}{V} \left[ \frac{1 + \frac{I}{K_{id}}}{1 + \frac{I}{K_{in}}} \right] \quad (4)$$

Values would be determined for  $K_{is}$  and  $K_{in}$  from the variation of the slope with varying concentrations of inhibitor  $I$  and for  $K_{in}$  and  $K_{id}$  from the variation of intercept with  $I$  in the same way as described for hyperbolic competitive inhibition.

## Nonpermissive Condition

The condition under which a [conditional lethal mutant](#) cannot grow is called the nonpermissive or restrictive condition. In contrast, the [permissive condition](#) allows growth. Typical nonpermissive conditions are high or low temperature, extremes of pH, osmotic pressure, oxygen concentration, and the absence of growth factors. Particular hosts can become nonpermissive or restrictive for mutant strains of **viruses**, **bacteriophages**, or other parasites (see [Conditional Lethal Mutations](#)).

## Nonpolar

A molecule or group of atoms is said to be *nonpolar* if it lacks electronegative atoms such as N, O, and S, so that it does not contain groups with substantial partial electrostatic charges. Such molecules include steroids and other [lipids](#), and the long hydrocarbon side-chains of [amino acids](#) such as **leucine**, **isoleucine**, and [valine](#). Such molecules show a tendency to move from water to oily



surroundings such as the interiors of membranes. The physical basis of this tendency, and its quantitative measurement, are discussed in [Hydrophobicity](#) and [Hydrophobic Effect](#).

## Nonrepetitive DNA

Nonrepetitive DNA is comprised of DNA sequences present in only one or in a small number of copies in a [genome](#). It contrasts with [repetitive DNA](#), especially [Alu sequences](#), [CpG Islands](#), [Sines](#), [Lines](#), [microsatellite](#) and [minisatellite](#) DNA. Nonrepetitive DNA has the reassociation kinetics expected from unique sequences and is characterized by a high  $C_0t$  value. It generally consists of protein-coding genes or a plasmid in only one copy per genome.

## Nonsense Codons and Unassigned Codons

The term *nonsense codons* is usually used for the [stop codons](#) of the [genetic code](#) that do not code for any [amino acid](#) and are used to terminate [translation](#) during **protein biosynthesis**. Unassigned codons are codons corresponding to neither amino acids nor stop codons and do not exist in the open [reading frames](#) (ORF) of genes coding for proteins. Such codons should be called the real nonsense codons because they are recognized by neither tRNA nor [release factors](#) (RF), the proteins that function to terminate translation (see [Stop Codons](#)). Such unassigned codons have been recognized only with determination of the sequences of the entire genomes of [mitochondria](#) and of some organisms with small genomes. To be an unassigned codon, the tRNA or RF corresponding to this codon probably disappeared from the translation system ([1](#)). The unassigned codons determined thus far ([2-6](#)) are shown in [Table 1](#). Recently, it has been shown that AAA is unassigned in a hemichordate ([7](#)). AAA is assigned to lysine in most metazoans and to asparagine in echinoderms. The situation in hemichordates may be an intermediate en route to the related echinoderms ([7](#)).

**Table 1. Unassigned Codons**<sup>\*a</sup>

| System   | Probable Unassigned Codon       | Cause                            |
|--|---------------------------------|----------------------------------|
| Mycoplasma capricolum                            | CGG(Arg) <sup>*b</sup>          | AT-pressure; lack of tRNAArgCCG  |
| Micrococcus luteus                               | AGA(Arg)                        | GC-pressure; lack of tRNAArg*UCU |
| AUA(Ile)   | GC-pressure; lack of tRNAIleLAU |                                  |
| <i>Torulopsis glabrata</i> mitochondriaCGN (Arg) | AT-pressure; lack of tRNAArgCCG |                                  |

|   |                               |   |
|---|-------------------------------|---|
| Prototheca<br>wickerhamii<br>(A green alga<br>mitochondria) | CGN(Arg)                      | AT-pressure; lack of<br>tRNAArgCCG                    |
| UGA(Stop)<br><br>or UAG(Stop)                               | AT-pressure;? Lack of<br>RF-2 | AT-pressure;? Lack of RF-<br>1                        |
| Drosophila* <sup>c</sup> and                                | AGG(Ser)                      | ?AT-pressure; lack of<br>tRNA <sup>Ser</sup> *GCU (2) |
| Mosquito* <sup>d</sup><br>mitochondria                      | (2)                           |   |

---

<sup>a</sup> All except for the last item from Ref. 2, p. 66.

<sup>b</sup> Parenthesis means amino acid or stop codon in the universal genetic code, or for the last item in animal mitochondrial genetic code.

<sup>c</sup> From References 3 and 4.

<sup>d</sup> From References 5 and 6.

The appearance of unassigned codons may be caused by direct mutation pressure (as a result of changes in the overall AT and GC content of the genome DNA), change of release factor, or genome economization (Table 1). It is regarded as an intermediate step in codon reassignment in producing genetic code variations (1, 2). The presence of unassigned codons is important evidence for the codon capture theory (1, 2) (see “Candida code and codon reassignment theories” section in [Genetic Code](#)).

### Bibliography

1. S. Osawa, T. H. Jukes, K. Watanabe, and A. Muto (1992) *Microbiol. Rev.* **56**, 229–264.
2. S. Osawa (1995) *Evolution of the Genetic Code*, Oxford University Press, Oxford, UK, pp. 1–205.
3. D. O. Clary and D. R. Wolstenholme (1985) *J. Mol. Evol.* **22**, 252–271.
4. D. L. Lewis, C. L. Farr, and L. S. Kaguni (1995) *Insect Mol. Biol.* **4**, 263–278.
5. C. B. Beard, D. M. Hamm, and P. H. Collins (1993) *Insect Mol. Biol.* **2**, 103–124.
6. S. E. Mitchell, A. F. Cockburn, and J. A. Seawright (1993) *Genome* **36**, 1058–1073.
7. J. Castresana, G. Feldmaier-Fuchs, and S. Pääbo (1998) *Proc. Natl. Acad. Sci. USA* **95**, 3703–3707.

### Nonsense Mutation

A nonsense mutation is a [mutation](#) that creates a translational [stop codon](#) in the coding sequence of a gene. The three chain-terminating codons, which are the normal stop signals for [translation](#), are UGA, UAA, and UGA, called **amber**, **ochre**, and **opal**, respectively (although the latter is rarely

used). Nonsense mutations are created by base substitutions at sense (coding) **codons** or by [frameshift mutations](#), which usually create a **nonsense codon** 3' to the frameshift.

Nonsense mutations are **suppressed** by mutant [transfer RNAs](#) that recognize the nonsense codon and insert their cognate amino acid into the [polypeptide chain](#). Thus, nonsense mutations are a class of [conditional lethal mutations](#). Because they occur in any gene but are suppressed by the same class of suppressors, nonsense mutations are a particularly powerful genetic tool. An elegant example was determining the life cycle of **bacteriophage T4** (1). In addition, by using mutant tRNAs with known cognate amino acids, the functional significance of amino acid replacements at various positions in a protein can be determined (2). Most suppressor tRNAs are mutant in their [anticodon](#), although mutations and modifications elsewhere in the molecule have suppressing effects. Indeed, the first UGA suppressor discovered is mutant outside of its anticodon. Mutations in other components of the translational machinery, such as the [release factors](#) and **ribosomal** subunits, also can suppress nonsense mutations (3).

Nonsense mutations can be polar, ie, can eliminate the expression of genes in the same [operon](#) encoded downstream from the mutation. The polarity of a nonsense mutation depends on its proximity to a codon that can restart translation. In *Escherichia coli* the polarity of many nonsense mutations can be suppressed by mutations in **Rho**, the transcriptional [termination factor](#). Thus in these cases polarity is probably caused by inhibition of **RNA polymerase** by Rho when it interacts with untranslated mRNA (4).

#### Bibliography

1. R. H. Epstein, A. Bolle, C. Steinberg, E. Kellenberger, E. Boy de la Tour, R. Chevalley, R. Edgar, M. Susman, C. Denhardt, and I. Lielausis (1964) Cold Spring Harbor Symp. Quant. Biol. **28**, 375–392.
2. P. Markiewicz, L. G. Kleina, C. Cruz, S. Ehret, and J. H. Miller (1994) J Mol. Biol. **240**, 421–433.
3. H. Engelberg-Kulka and R. Schoulaker-Schwarz (1996) In *Escherichia coli and Salmonella; Cellular and Molecular Biology*, 2nd ed. (F. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology Press, Washington, DC, pp. 909–921.
4. S. Adhya and M. E. Gottesman (1978) Annu. Rev. Biochem. **47**, 967–996.

#### Suggestion for Further Reading

5. T. D. Brock (1990) *The Emergence of Bacterial Genetics*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

#### Nonsense Suppression

Nonsense mutations create stop codons within the coding regions of genes. Nonsense suppression occurs when mutations in the translational apparatus allow the decoding of nonsense codons as sense codons. Most nonsense suppressor mutations alter tRNA anticodons, making them cognate the stop codon; however, a handful of suppressor mutations have also altered release factors or other translational components. Suppression requires that the amino acid inserted is compatible with protein function. Furthermore, because the nonsense mutation is decoded as “stop” at least some of the time, suppression allows only a fraction of the normal level of polypeptide synthesis. Despite

these limitations nonsense suppressors have been remarkably useful research tools. Other forms of mutational suppression are described in [Genetic Suppression](#).

There are three termination codons in the universal genetic code, UAG, UAA, and UGA. Historically, UAG nonsense mutations are called “amber,” named after the mother of a pioneering worker who first isolated such mutants in phage T4 (Harris Bernstein — Bernstein means “Amber” in German. Interestingly, the molecular nature of amber mutations was not known at the time, see Ref. [1](#)). After their discoveries, the other nonsense codons were also named according to a “color code,” with UAA as “ochre” and UGA as “opal.” tRNA suppressors have been isolated that can suppress each color, and because of U:G wobble decoding ochre suppressors will also suppress amber mutations.

Nonsense mutations are a type of conditional mutation, because the mutant phenotype is dependent on the absence of a suitable suppressor. This property has been used extensively in genetic studies of phage, bacteria, yeast, and the multicellular eukaryote *Caenorhabditis elegans* ([2](#)). For example, a large number of phage mutants have been isolated based on the ability to form plaques only on hosts that express nonsense suppressors. Such mutants are referred to as “suppressor-sensitive” and are commonly abbreviated as “sus” in the literature. *sus* mutations can serve as markers for genetic mapping and other routine genetic manipulations. In addition, studies of the molecular phenotypes during the nonproductive condition have helped illuminate the functions of the genes that carry nonsense mutations.

Nonsense suppression has also been extremely useful for studies of translational mechanisms ([3](#)). Suppressor tRNA genes are easy to manipulate, and suppression efficiencies are readily quantified; together, these properties make suppression the assay of choice for probing translational mechanisms. Suppressor tRNAs must compete with the peptide release factors that normally decode nonsense codons. Suppressors that are good competitors allow for efficient suppression of the nonsense mutation. Thus, the level of suppression can be used as a measure of the translational efficiency of tRNAs. Surprisingly, perhaps, tRNA suppressors can vary dramatically in their translational efficiency. Yarus observed that suppression efficiency was correlated with the tRNA nucleotides near the anticodon and proposed that decoding efficiency was a property of an “extended anticodon” in which the nucleotides near the anticodon contribute to its ability to translate codons ([4](#)). This hypothesis was confirmed by making measuring the effects mutations that saturate the anticodon arm of an amber suppressor. It was found that the native sequence was optimal for suppression, and that anticodon arm mutations reduce suppression efficiency by various degrees ([5](#)). Nonsense suppressors have also been used to show that nucleoside modifications near the anticodon increase translational efficiency. Mechanisms are not always clear, but a common theme is that bulky modification 3' to the anticodon increases translational efficiency, as if increased base stacking stabilizes anticodon:codon complexes ([6](#), [7](#)).

Studies of suppressors have also shown that nucleotides outside of the anticodon region can affect translational efficiency ([8-11](#)). Mutational analyses show that the central or hinge region of the tRNA is important for translation, but probably not by directly affecting anticodon:codon pairing. Instead, this region may participate in a conformational change in the tRNA that must occur during ribosomal acceptance of the aminoacyl-tRNA. Mutations that interfere with this change affect the likelihood that the tRNA will be selected, regardless of whether the anticodon is perfectly matched to the codon ([12](#)). Certain mutations in this region allow, for example, reading with a first-position U:G pair or a third-position C:A pair. It is thought the hinge mutations facilitate conformation changes leading to ribosomal acceptance despite the mismatched base pairs.

McClain and co-workers have used nonsense suppressor tRNAs to demonstrate changes in the aminoacylation specificities of several tRNAs and thus map the determinants for tRNA amino acid “identity” ([13](#)). They engineered an amber termination site near the 5' end of the gene for the easily isolable dihydrofolate reductase enzyme. They then mutagenized tRNA suppressors and determined

which amino acids the variants insert by sequencing the amino termini of isolated enzymes. They found that a small number of nucleotides are mostly responsible for defining tRNA identity. The use of suppressor tRNAs was critical for these studies because there is ambiguity about which tRNA actually decodes the amber site. Thus, a change in the amino acid inserted can be unambiguously attributed to the changes made in the amber suppressor tRNA.

Nonsense suppressors are also used to show that codon translation is affected by neighboring nucleotides (codon context). A large number of studies show that nonsense codons are more readily suppressed if the 3' neighbor nucleotide is a purine. This context effect has at least two molecular sources: the anticodon:codon complex is stabilized by base stacking with the 3' purine (14, 15), and the release factors are highly dependent on the neighbor for termination with 3' U being optimal (16, 17).

Nonsense suppressors will also be used for protein engineering. In one approach, Abelson and collaborators have constructed an extensive set of tRNA suppressors that can be used to insert a wide range of amino acids at amber mutations (18). These tRNAs may allow systematic tests of the effects of various amino acids on protein structure and function. Proteins engineered to contain amino acids with fluorescent or especially reactive side chains would be of great use in studies of protein structure and function and for biotechnological applications. Recently, it has been demonstrated that certain nonsense suppressors can be aminoacylated with novel amino acids, and that modified tRNAs can direct the incorporation into proteins that have specific nonsense mutations within their coding sequences (19). Currently, only a few novel amino acids may be used in this way. But, hopefully, a large variety of specific labels will become available.

## Bibliography

1. R. S. Edgar (1966) In *Phage and the Origins of Molecular Biology* J. Cairns, G. S. Stent, and J. D. Watson, eds., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 166–180.
2. J. Hodgkin, K. Kondo, and R. H. Waterston (1987) *Trends Genet.* **3**, 325–329.
3. E. J. Murgola (1994) In *tRNA: Structure, Biosynthesis, and Function* D. Söll and U. L. RajBhandary, eds., American Society for Microbiology Press, Washington, D. C., pp. 491–509.
4. M. Yarus (1982) *Science* **218**, 646–652.
5. M. Yarus, S. W. Cline, P. Weir, L. Breeden, and R. C. Thompson (1986) *J. Mol. Biol.* **192**, 235–255.
6. J. F. Curran (1998) In *Modification and Editing of RNA* (H. Grosjean and R. Benne, eds.), American Society for Microbiology Press, Washington, D. C., pp. 493–516.
7. G. R. Björk (1996) In *Escherichia coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt, R. Curtis III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology Press, Washington, D. C., pp. 861–886.
8. D. Smith and M. Yarus (1989) *J. Mol. Biol.* **206**, 489–501.
9. D. Smith and M. Yarus (1989), *J. Mol. Biol.* **206**, 503–511.
10. D. W. Schultz and M. Yarus (1994) *J. Mol. Biol.* **235**, 1381–1394.
11. D. W. Schultz and M. Yarus (1994) *J. Mol. Biol.* **235**, 1395–1405.
12. M. Yarus and D. Smith (1994) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, D. C., pp. 443–469.
13. W. H. McClain (1994) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, D. C., pp. 335–347.
14. M. Yarus and J. F. Curran (1992) In *Transfer RNA in Protein Synthesis* (D. L. Hatfield, B. J. Lee, and R. M. Pirtle, eds.), CRC Press, Boca Raton, FL, pp. 319–365.
15. R. H. Buckingham (1994) *Biochimie* **76**, 351–354.

16. W. T. Pedersen and J. F. Curran (1991) *J. Mol. Biol.* **219**, 231–241.
17. E. S. Poole, C. M. Brown, and W. P. Tate (1995) *EMBO J.* **14**, 151–158.
18. L. G. Kleina, J.-M. Masson, J. Normanly, J. Abelson and J. H. Miller (1990) *J. Mol. Biol.* **213**, 704–717.
19. L. Wang, A. Brock, B. Heberich, and P.G. Schultz (2001) *Science* **292**, 498–500.

## Nonsynonymous Substitution

Nonsynonymous substitution is an **evolutionary** term meaning the fixation in a population or subpopulation of a [missense mutation](#), ie, a [mutation](#) that changes a **codon** into another codon that specifies a different [amino acid](#). Nonsynonymous substitutions can be conservative, ie not greatly alter the function of the altered [protein](#) (see [Conservative Substitutions](#)). Nonsynonymous substitutions that are not conservative may indicate an adaptive change in the altered protein. This can only be confirmed, however, by demonstrating that the substitution changes the function of the protein.

## Notch Signaling

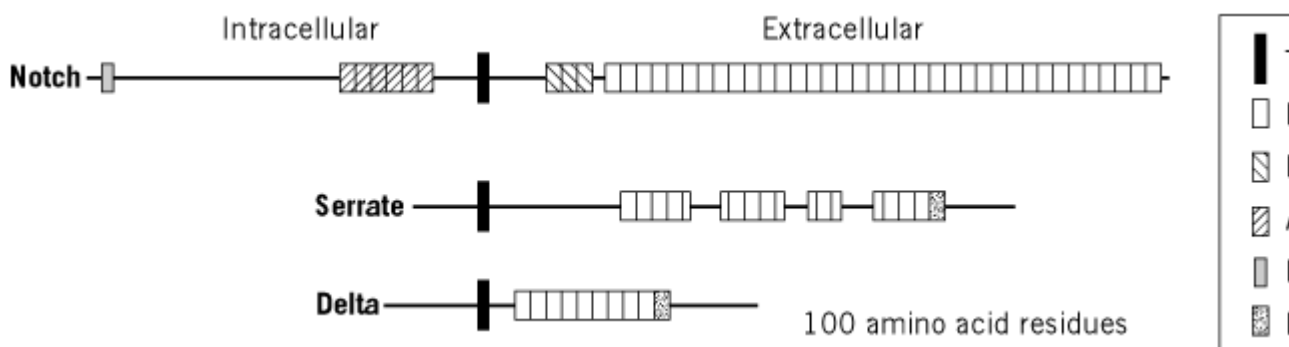
How a single cell, the fertilized egg, transforms into a complex multicellular organism is one of the most fascinating questions in biology (see [Development](#)). Two processes that are critical in this transformation are the generation of distinct cell types and the patterning of these cells into functional tissues and organs. Molecular genetic experiments in the fruit fly *Drosophila melanogaster* and the [nematode](#) worm *Caenorhabditis elegans* first revealed that similar transmembrane receptors (called the Notch receptor in *Drosophila*) play critical roles in numerous cell–cell interactions that are involved in the establishment and patterning of a wide array of cell types. Remarkably, **homologous** Notch receptors in echinoderms, amphibians, fish, birds, and mammals have subsequently been found to mediate a similar broad range of cellular interactions. This suggests that the Notch receptor is a fundamental intercellular signaling molecule used in the patterning and generation of different cell types in all metazoans. Work from a number of organisms has demonstrated that the Notch receptor is a central component of a conserved intercellular signaling pathway that has at least three additional core members. These are the transmembrane ligands Delta and Serrate and the [transcription factor](#) CSL. Consistent with the diverse interactions that the Notch pathway mediates during development, the regulation of the receptor and ligands varies a great deal in different developmental contexts. One of the most intriguing questions about the Notch signaling pathway is how this single pathway regulates the generation of such a broad range of cell types. Gaining insight into this aspect of Notch signaling has been challenging, but answers to this question promise to shed light on the fundamental question of how cells are patterned and how they assume different fates during development.

### 1. Components of the Notch Signaling Pathway

Notch homologous proteins are large [membrane proteins](#) that span the [membrane](#) once. In the

extracellular domain, all Notch proteins contain multiple, tandemly arranged [EGF motifs](#) and three Notch/Lin-12 repeats. In the intracellular domain, Notch proteins contain six [ankyrin](#) repeats, a motif involved in direct [protein–protein interactions](#), as well as a **PEST sequence** that may regulate the stability of the protein (Fig. 1). Only a single Notch protein is encoded in the *Drosophila* and sea urchin **genomes** (1). In *C. elegans*, there are two Notch-like proteins, called GLP-1 and LIN-12. In vertebrates, four Notch homologous proteins have been identified.

**Figure 1.** Schematic diagram of the [primary structures](#) of *Drosophila* Notch receptor and its ligands, Delta and Serrate.



Delta and Serrate, the two identified ligands for Notch, are also single-spanning transmembrane proteins. Much as Notch does, all ligands for Notch contain multiple EGF-like repeats in the extracellular domain, as well as a unique **cysteine-rich** motif called the DSL domain (Fig. 1). In *Drosophila*, there appears to be only one Delta and one Serrate protein although multiple homologues of each have been identified in vertebrates (eg, 2). In *C. elegans*, there are two identified ligands, LAG-2 and APX-1, which are smaller than, but similar in structure to, Delta and Serrate. Remarkably, of the 36 EGF-like repeats in *Drosophila* Notch, only two (repeats 11 and 12) are necessary for binding to either Serrate or Delta. This has raised the possibility that additional ligands may exist that bind to the other EGF-like repeats in the extracellular domain; to date, however, no other ligands have been identified.

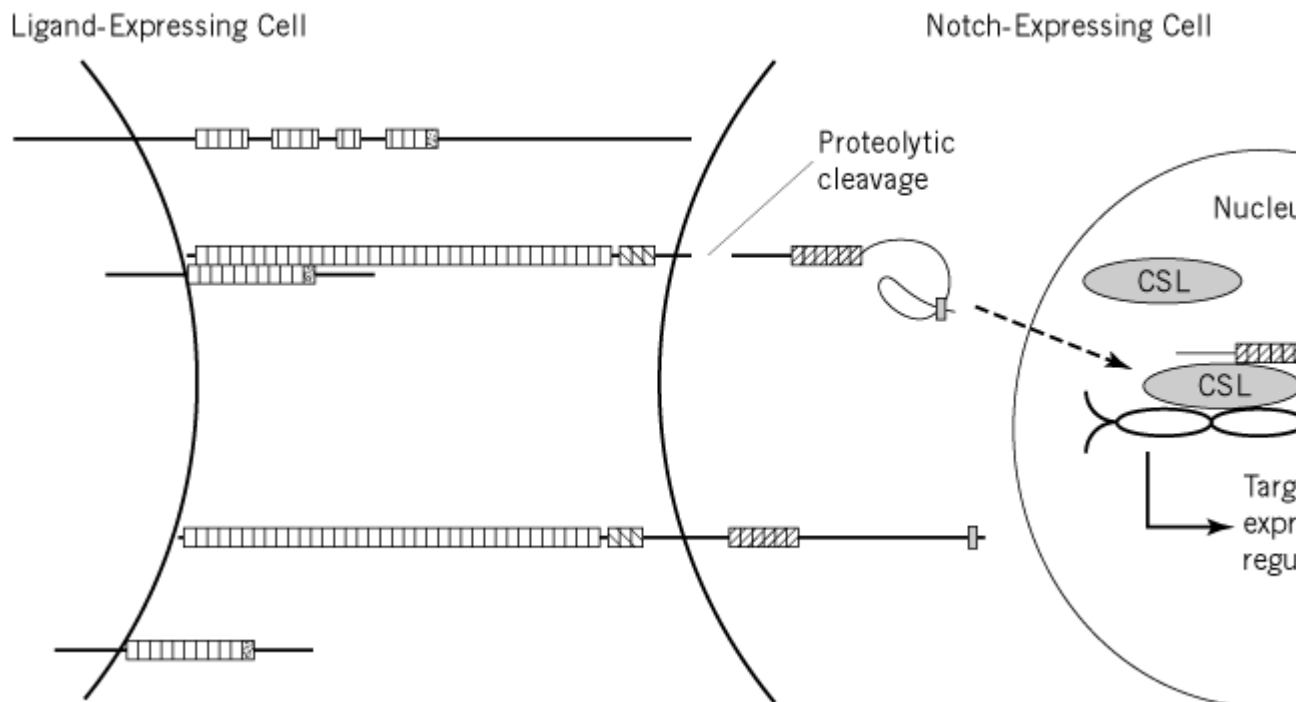
An additional shared component in the Notch pathway has been demonstrated by studies in *Drosophila*, *C. elegans*, and vertebrates. These are the homologous CSL proteins, which encode nuclear [DNA-binding proteins](#) (3, 4). CSL proteins appear to act both as [repressors](#) and activators of [transcription](#), providing a link for the Notch pathway to alterations in gene expression (5-7).

## 2. Activation of the Notch Receptor

Interestingly, almost all identified Notch proteins contain nuclear localization signals in their intracellular domain (see **Nuclear import/export**). In addition, the intracellular domains of most Notch proteins, when expressed alone in cells and untethered to the membrane, enter the nucleus and cause **phenotypes** in cells similar to the one caused by overactivation of the Notch pathway. Furthermore, the intracellular domains of some Notch proteins bind physically to CSL proteins and activate transcription of target genes. These findings have led to the proposal that the interaction of Notch with its ligands results in the proteolytic release and translocation of the intracellular domain into the nucleus (Fig. 2 and Ref. 8). The amount of intracellular domain that enters the nucleus to signal *in vivo*, however, appears to be extremely low and requires very sensitive assays to detect (9, 10). One common transcriptional target in both vertebrates and *Drosophila* for the Notch-CSL protein complex is the *Enhancer of Split* complex (7, 11). This complex encodes seven proteins related to basic helix-loop-helix proteins, as well as a protein that shares homology with **GTP-binding proteins**. Another target of Notch signaling in *Drosophila* is the *vestigial* gene, which is

activated by Notch signaling at the dorsal-ventral boundary in the wing to stimulate cell proliferation. Other targets of Notch-mediated signaling are likely to exist but have yet to be identified (eg, 12). Intriguingly, during some cell–cell interactions, Notch proteins do not require CSL proteins to signal (7, 12, 13). But how Notch proteins signal independently of CSL proteins is unknown at this time.

**Figure 2.** Hypothetical model for activation of the Notch signaling pathway. The Notch receptor can bind to either Delta neighboring cell. Ligand binding leads to proteolytic release of the intracellular domain of Notch from the membrane. A) intracellular domain is translocated into the nucleus, where it interacts with a CSL protein. This interaction leads to a cor the transcriptional activity of target genes of the Notch pathway.



### 3. Pleiotropy of Notch-mediated signaling

Notch was originally identified in *Drosophila* as a “neurogenic” gene involved in the segregation of neuroblasts and epidermal cells in the embryonic neuroectoderm. However, **temperature-sensitive mutants** in *Drosophila Notch* and *C. elegans Notch* homologues, specific mutations in *Notch* homologues, and expression and misexpression studies have uncovered an amazing diversity of cell-fate decision and patterning processes that Notch receptors mediate in *Drosophila*, *C. elegans*, and vertebrate development. For example, in *Drosophila*, Notch signaling has been implicated in the establishment of distinct cell types in the mesoderm, endoderm, Malpighian tubules, eye, and ovary. In addition, Notch signaling is crucial to the organization of the *Drosophila* dorsal-ventral wing boundary, which is a patterning process that directs wing growth and wing margin formation (eg, 14). In *C. elegans*, the two Notch-like proteins GLP-1 and LIN-12 similarly mediate numerous cell-fate decisions, and GLP-1 plays a central role in an important patterning process, the maintenance of mitotic nuclei in the distal end of the *C. elegans* ovary. In vertebrates, Notch signaling is involved in cell-fate specification during neurogenesis, [hematopoiesis](#), and muscle cell development (3, 15-17). In addition, Notch signaling is necessary for the segmentation of somites, a critical process in patterning the segmental organization of vertebrate embryos (18, 19). Furthermore, studies of human diseases have shown that altered forms of Notch signaling can lead to T-cell lymphoblastic leukemia and to a developmental disorder called Alagille syndrome, which results in developmental abnormalities in the liver, kidney, heart, eye, vertebrae, and facial structure (reviewed in 20).

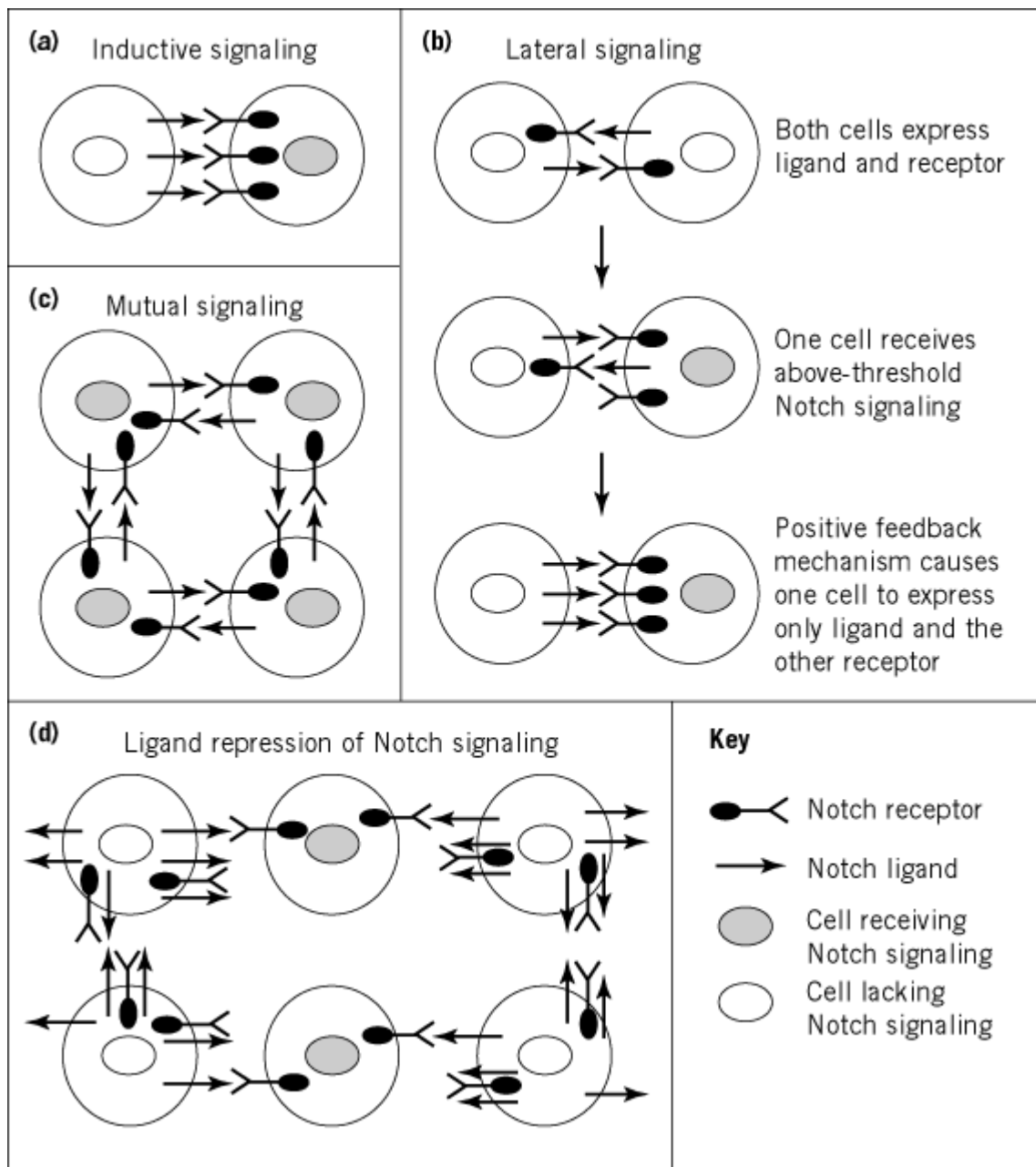


Intriguingly, particular mutations in the human Notch receptor, *Notch3*, have also been implicated in the adult onset disease CADASIL, which causes vascular alterations and leads to repeated strokes and dementia (21). The late onset of this disease suggests that Notch signaling may play an unknown function in differentiated cells of adult tissues.

#### 4. Diverse interactions between Notch and its ligands

Reflecting the numerous cell-fate and patterning processes mediated by Notch signaling, the ligands for Notch are used in diverse ways to activate and sometimes even repress the Notch receptor. In some developmental contexts, Notch and its ligands are arranged in what is known as an inductive interaction. In inductive cell–cell interactions, the signaling and receiving cells are not equivalent to one another, and one cell or tissue selectively signals to the other (Fig. 3a). An example of this is germ-line induction in *C. elegans*: the Notch-like protein GLP-1 is selectively expressed in the receiving germ-line tissue, and the ligand LAG-2 is restricted to the signaling distal-tip cell.

**Figure 3.** Diverse interactions between the Notch receptor and its ligands. (a) In inductive signaling, the signaling and receiving cells are not equivalent to one another at the onset of the interaction. The signaling cell expresses only the ligand, and the receiving cell expresses only the Notch receptor. In lateral signaling, (b) interacting cells are initially equivalent to one another, and both express ligand and receptor. A stochastic fluctuation in signaling results in one cell receiving above-threshold Notch signaling. This triggers a feedback mechanism that reinforces this signaling difference such that one cell ultimately expresses only ligand and the other only receptor. In mutual signaling, (c) a group of cells express ligand and receptor, and all signal one another. This type of signaling may be used to coordinate patterning processes or cell-fate decisions within groups of cells. High levels of Notch ligand can inhibit Notch signaling within the same cell in which they are expressed (d). The mechanism for this inhibition is not understood, but it may be mediated by direct interactions between Notch and its ligands within the same cell.



Notch and its ligands also mediate lateral signaling, where the signaling and receiving cells are initially equivalent to one another in developmental potential and signaling ability. In this type of interaction, it appears that random fluctuations in signaling activity are reinforced and strengthened through feedback mechanisms so that one cell becomes the signaler and the other cell the receiver (Fig. 3b). An example of this is the interaction between the anchor cell (AC) precursor and the ventral uterine (VU) precursor cell during *C. elegans* development. Both cells are initially equivalent in their ability to become either the AC or the VU cell, and both express the Notch-like receptor LIN-12 and the ligand LAG-2. Interactions between these cells result in the activation of a feedback mechanism, such that one cell ultimately expresses only the ligand LAG-2 (signaling cell, which becomes the AC) and the other cell expresses only the Notch-like receptor LIN-12 (receiving cell, which becomes the VU).

Notch and its ligands can also be expressed and activated in groups of cells. This type of interaction is known as mutual signaling (Fig. 3c). An example of this occurs during the early development of the *Drosophila* nervous system in the ventral neurogenic ectoderm. All cells in this region express

both Notch and the ligand Delta. Usually, only a subset of the cells delaminate from the neurogenic ectoderm and become neuroblasts; in embryos lacking either Notch or Delta, however, most cells in the neurogenic ectoderm delaminate as neuroblasts. It appears that nearly every cell in the neurogenic ectoderm has the potential to become a neuroblast but that neighboring cells interact through Notch–Delta signaling to inhibit this potential, thus restricting the number of cells that develop as neuroblasts.

In some developmental contexts, high levels of Delta or Serrate inhibit Notch signaling within the same cell in which they are expressed (Fig. 3d). An example of this type of interaction is found at the *Drosophila* dorsal-ventral wing boundary. High levels of Delta or Serrate within cells adjacent to the wing boundary inhibit these cells from receiving Notch signaling (22). This ensures that Notch signaling is activated only in cells at the boundary that do not express ligands for Notch. How ligands for Notch inhibit Notch signaling within the same cell is not fully understood, but evidence suggests that the inhibition could be mediated by direct interactions between the ligand and receptor within the same cell. Thus, the levels of ligand present on a cell may be a critical component in controlling whether neighboring cells can signal one another with the Notch pathway.

## 5. Additional regulators of the Notch pathway

Several other molecules have been identified as important regulators of the Notch pathway. A protein called Numb was first identified in *Drosophila* as playing a role in the generation of the different cell types that are required to build neuronal sensory organs located in the epidermis. Protein localization studies found that Numb is distributed asymmetrically during cell divisions in cells that construct the sensory organ. Interestingly, Numb appears to affect the generation of distinct cell types by binding to the intracellular domain of Notch, which inhibits Notch signaling. Therefore, Numb creates a signaling interaction between neighboring sister cells from a cell division very similar to inductive signaling, where one cell is the signaler and the other the receiver (see Fig. 3a).

Genetic studies have indicated that many interactions occur between the Notch pathway and another important cell–cell signaling pathway called the Wnt/**Wingless** pathway. During some interactions, these pathways appear to act synergistically; in others, however, they act antagonistically to one another. Interestingly, the *disheveled* gene product, a cytoplasmic protein and member of the Wnt/Wingless signaling pathway, may mediate direct antagonistic interactions between these two pathways. The Disheveled protein can bind to the intracellular domain of Notch and inhibit Notch signaling during the specification of *Drosophila* wing margin bristles.

Another protein that regulates the Notch pathway is Fringe, which is a secreted protein that appears to function in the extracellular space between cells. In *Drosophila*, Fringe can modulate interactions between Delta and Serrate and with Notch by blocking the ability of Serrate to activate Notch and potentiating the ability of Delta to activate the receptor (23, 24). Several vertebrate Fringe homologues have been identified as well, and they appear to mediate a similar function (25).

## 6. How does the Notch pathway control the generation of many cell types?

One of the most puzzling questions about Notch signaling is how this one pathway is involved in the formation of so many distinct cell types and patterning processes during development. Experiments in vertebrates and *Drosophila* have demonstrated that, during many Notch-mediated cell–cell interactions, Notch signaling appears to block the ability of cells to respond to differentiation signals or regulatory factors, maintaining cells in an undifferentiated state (15, 26). For example, in the *Drosophila* compound eye, inappropriate activation of Notch signaling in the presumptive R7 cell, at the time this cell is receiving the differentiation signal to become R7, inhibits this cell from differentiating as an R7 cell. Instead, after Notch signaling subsides, this cell differentiates into a cone cell, the default fate for this cell when it does not receive the R7 differentiation signal. Therefore, Notch signaling may be crucial in the establishment of many different cell types in controlling whether and when cells can respond to specific differentiation cues.

Notch signaling, however, does not always inhibit differentiation in cells. During *C. elegans* development, the Notch pathway mediates the formation of numerous cell types, yet there is currently no evidence that Notch signaling guides the generation of distinct cell types by inhibiting differentiation. Indeed, during the development of the vulva and uterus, Notch signaling appears to act as a differentiation signal in the generation of specific cells that contribute to the formation of both these tissues (27, 28). Furthermore, during *Drosophila* wing development, activation of Notch leads to the expression of growth and patterning genes (eg, 29), a finding inconsistent with the notion that Notch signaling is keeping these cells in an undifferentiated state. Therefore, although, in many cell–cell interactions, Notch signaling functions to establish distinct cell types by regulating the ability of cells to differentiate, Notch signaling does not always function in this manner. To gain a deeper understanding of how Notch signaling affects the establishment of numerous cell types and patterning process, it will be necessary to identify additional target genes regulated by Notch signaling and elucidate how these genes mediate the differentiation and patterning of distinct cell types.

### Bibliography

1. D. R. Sherwood and D. R. McClay (1997) *Development* **124**, 3363–3374.
2. C. Haddon et al. (1998) *Development* **125**.
3. J. L. de la Pompa et al. (1997) *Development* **124**, 1139–1148.
4. D. A. Wettstein, D. L. Turner, and C. Kintner (1997) *Development* **124**, 693–702.
5. A. M. Bailey and J. W. Posakony (1995) *Genes Dev.* **9**, 2609–2622.
6. J. J.-D. Hsieh and S. D. Hayward (1995) *Science* **268**, 560–563.
7. M. Lecourtois and F. Schweisguth (1995) *Genes Dev.* **9**, 2598–2608.
8. R. Kopan, E. H. Schroeter, H. Weintraub, and J. S. Nye (1996) *Proc. Natl. Acad. Sci.* **93**, 1683–1688.
9. E. H. Schroeter, J. A. Kisslinger, and R. Kopan (1998) *Nature* **393**, 382–386.
10. G. Struhl and A. Adachi (1998) *Cell* **93**, 649–660.
11. S. Jarriault et al. (1995) *Nature* **377**, 355–358.
12. C. Shawber et al. (1996) *Development* **122**, 3765–3773.
13. S. Wang, S. Younger-Shepherd, L. Y. Jan, and Y. N. Jan (1997) *Development* **124**, 4435–4436.
14. J. F. de Celis, A. Garcia-Bellido, and S. J. Bray (1996) *Development* **122**, 359–369.
15. L. A. Milner et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13014–13019.
16. R. Kopan, J. S. Nye, and H. Weintraub (1994) *Development* **120**, 2385–2396.
17. H. von Boehmer (1997) *Curr. Biol.* **7**, R308–R310.
18. R. A. Conlon, A. G. Reaume, and J. Rossant (1995) *Development* **121**, 1533–1545.
19. W.-C. Jen, D. Wettstein, D. Turner, A. Chitnis, and C. Kintner (1997) *Development* **124**, 1169–1178.
20. S. Artavanis-Tsakonas (1997) *Nature Genet.* **16**, 212–213.
21. A. Joutel et al. (1996) *Nature* **383**, 707–710.
22. C. A. Micchelli, E. J. Rulifson, and S. S. Blair (1997) *Development* **124**, 1485–1495.
23. V. M. Panin, V. Papayannopoulos, R. Wilson, and K. D. Irvine (1997) *Nature* **387**, 908–912.
24. R. J. Fleming, Y. Gu, and N. A. Hukriede (1997) *Development* **124**, 2973–2981.
25. S. H. Johnston et al. (1997) *Development* **124**, 2245–2254.
26. R. Dorsky, D. H. Rapaport, and W. A. Harris (1995) *Neuron* **14**, 487–496.
27. I. S. Greenwald, P. W. Sternberg, and H. R. Horvitz (1983) *Cell* **34**, 435–444.
28. A. P. Newman, J. G. White, and P. W. Sternberg (1995) *Development* **121**, 263–271.
29. C. J. Neumann and S. M. Cohen (1996) *Development* **122**, 3477–3485.

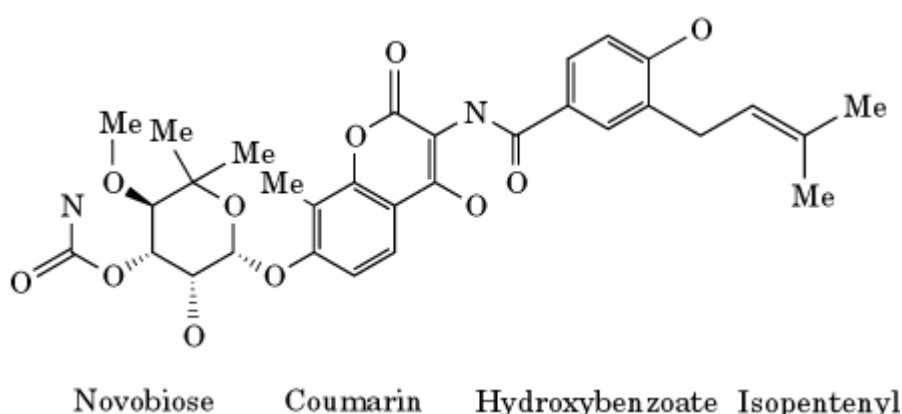
### Suggestions for Further Reading

30. S. Artavanis-Tsakonas, K. Matsuno, and M. E. Fortini (1995) Notch signaling. *Science* **268**, 225–232. A very thorough review of Notch signaling; contains an extensive bibliography.
31. C. R. Coffman, P. Skoglund, W. A. Harris, and C. R. Kintner (1993) Expression of an extracellular deletion of Notch diverts cell fate in *Xenopus* embryos. *Cell* **73**, 659–671. The first demonstration that Notch signaling maintains cells in an undifferentiated state.
32. R. G. Fehon, P. J. Kooh, I. Rebay, C. L. Regan, T. Xu, M. A. T. Muskavitch, and S. Artavanis-Tsakonas (1990) Molecular interactions between the protein products of the neurogenic loci *Notch* and *Delta*, two EGF-homologous genes in *Drosophila*. *Cell* **61**, 523–534. Elegant work demonstrating that Notch and Delta bind to each other.
33. J. Kimble and P. Simpson (1997) The LIN-12/Notch signaling pathway and its regulation. *Ann. Rev. Cell Dev. Biol.* **13**, 333–361. An outstanding recent review of Notch signaling; very good reference source.

### Novobiocin

Novobiocin is an antibiotic of the aminocoumarin class that was originally discovered in 1955 (1). It is produced by bacteria of the species *Streptomyces spheroides* and *Streptomyces niveus* (2). Novobiocin is composed of a noviose sugar linked to a substituted coumarin and a prenylated 4-hydroxybenzoic acid (Fig. 1). Because of its pronounced side effects, novobiocin is not the antibacterial treatment of choice in the clinic (3). Nevertheless, it has regained some interest recently due to its synergistic effect with some antitumor drugs (4, 5).

Figure 1. The chemical structure of novobiocin.



The antibacterial activity of novobiocin is based on its inhibition of bacterial type-II DNA topoisomerases (see [DNA Topology](#)). The type-II DNA topoisomerases are essential cellular [enzymes](#) that, together with the type-I DNA topoisomerases, determine DNA superhelicity and are involved in **chromosomal** condensation and segregation (6, 7). Using the energy derived from ATP hydrolysis, type-II topoisomerases introduce a double-stranded cut into a duplex DNA, pass a second

duplex through the opening, and then reseal the cut. In prokaryotes, the type-II topoisomerases DNA gyrase and DNA topoisomerase IV are both tetramers of two subunits (8). The A subunit of DNA gyrase (ParC in topoisomerase IV) harbors the DNA breakage-reunion activity, and the B subunit of DNA gyrase (ParE in topoisomerase IV) is responsible for ATP hydrolysis.

The bacterial-type II topoisomerases are the targets of several antibiotic compounds. Quinolones interact with the A subunit, coumarins (eg, novobiocin) and cyclothialidines (9) with the B subunit. Novobiocin (10, 11) and the cyclothialidines (12) block the [ATPase](#) activity of the B subunit of DNA gyrase by binding to the ATP-binding site. The detailed binding mode of novobiocin and of an ATP analog to a 24-kDa N-terminal fragment of the *Escherichia coli* DNA gyrase B protein has recently been determined by [X-ray crystallography](#) (13). The structure shows that ATP and novobiocin have partially overlapping binding sites. Therefore, novobiocin is a **competitive inhibitor** of ATP binding to the gyrase B subunit. An extensive **hydrogen-bonding** network is formed between novobiocin and the enzyme. Most of the bonds are contributed by the novobiose sugar group that forms hydrogen bonds with Asn46, Ala47, and Asp73. In addition there are several hydrogen bonds to ordered [water](#) molecules, which in turn hydrogen bond to Val43, Val71, Asp73, Gly77, and Thr165. The coumarin ring of novobiocin forms two hydrogen bonds to Arg136, and it is also involved in **hydrophobic** interactions with Arg76 and Pro79. The prenylated hydroxybenzoate moiety of novobiocin is not required for DNA gyrase inhibition (14) and shows few interactions with the enzyme. This group is important, however, for the antibacterial activity of novobiocin, possibly by influencing uptake of the compound by bacterial cells. The crystal structure of chlorobiocin, a close homologue of novobiocin, bound to the 24-kDa N-terminal fragment of the gyrase B subunit is also known (15). This information is very helpful in understanding differences in the affinity of various novobiocin derivatives, and it could be exploited in rational drug design.

The occurrence of resistance to novobiocin is a frequent phenomenon. [Mutation](#) of the gyrase B subunit gene that leads to an altered gyrase B subunit is the most widespread mechanism of novobiocin resistance. A summary of the known mutations leading to novobiocin resistance is given in Table 1. It is obvious that the conserved residue Arg136 (*E. coli* numbering) is most often involved in resistance development. This residue has also been mutated in two coumermycin A-resistant strains of *Borrelia burgdorferii* (16). The susceptibility of these strains to novobiocin has not been tested, but coumermycin-resistant strains are normally cross-resistant to novobiocin. Arg136 forms two hydrogen bonds with the coumarin ring in the crystal structure of the 24-kDa N-terminal fragment of *E. coli* gyrase B with novobiocin. Therefore, involvement of this residue in resistance development is not surprising (13). The effects of mutation of Arg136 on the structure of the 24-kDa N-terminal fragment and on the binding affinity and the binding mode of novobiocin have recently been described (17). Mutation of Arg136 decreases the ATPase and supercoiling activity of the DNA gyrase B subunit (18), but the residual activity of the enzyme is still sufficient to fulfill the essential task of DNA gyrase. Novobiocin-resistant mutants of gyrase B must still retain sufficient activity to carry out the essential role of DNA gyrase, so it is not surprising that the mutated residues are not directly involved in binding ATP but instead surround the ATP-binding site.

**Table 1. Mutations in the *gyrB* Gene from Novobiocin-Resistant Strains**

| Organism                | Mutation of Gyrase B <sup>a</sup> | Increase in MIC <sup>b</sup> or IC50 <sup>b</sup> | References             |
|-------------------------|-----------------------------------|---|------------------------|
| <i>Escherichia coli</i> | R136L                             | >8×   | <a href="#">27</a>     |
|                         | R136C                             | 12×   | <a href="#">18, 27</a> |
|                         | R136S                             | 13.2×   | <a href="#">18</a>     |

|                                 |   |         |                        |
|---------------------------------|---|---------|------------------------|
|                                 | R136H   | 11.6×   | <a href="#">18</a>     |
|                                 | G164V   | >3×     | <a href="#">18, 28</a> |
| <i>Staphylococcus aureus</i>    | G85S (G77)                                      | 32×     | <a href="#">29</a>     |
|                                 | I102S (I94)                                     | 32×     | <a href="#">29</a>     |
|                                 | S128L (V120)                                    | 8×      | <a href="#">29</a>     |
|                                 | R144I (R136)                                    | 64×     | <a href="#">29</a>     |
|                                 | I102V, T173N<br>(T165)                          | 32×     | <a href="#">29</a>     |
|                                 | S128L, A108S<br>(A100)                          | 64×     | <a href="#">29</a>     |
|                                 | R144I, I175T<br>(V167)                          |         |                        |
|                                 | R144I, I102S                                    |         |                        |
| <i>Streptococcus pneumoniae</i> | S127L (V120)                                    | 128×    | <a href="#">31</a>     |
| <i>Haloferax</i>                | D82G (G81),<br>S122T<br>(S121), R137H<br>(R136) | No data | <a href="#">32</a>     |

<sup>a</sup> The numbers in brackets refer to the residue in *E. coli* homologous to the mutated residue in the respective species. The amino acid residues are indicated by their one-letter abbreviations.

<sup>b</sup> MIC: Minimal inhibitory concentration; IC<sub>50</sub>: 50% inhibitory concentration.

In addition, gyrase B resistance to novobiocin can also result from an active process of efflux from the cell. Recent results have demonstrated that the baseline level of resistance of *Haemophilus influenzae* and *Pseudomonas aeruginosa* to novobiocin is caused by multidrug efflux pumps that actively extrude novobiocin ([19, 20](#)). The novobiocin-producing organism *Streptomyces sphaeroides* protects itself against the toxic effects of novobiocin by producing two gyrase B subunits, one of which is resistant to the drug ([21](#)). Two loci that confer novobiocin resistance have been identified in *Streptomyces niveus* ([22](#)). Both loci also hybridize to DNA from *S. sphaeroides* but they do not encode a gyrase because they do not hybridize to the DNA gyrase gene.

Novobiocin is widely used in basic research, to study the effects of DNA topology on **gene expression**. DNA topology is influenced by many parameters, such as osmolarity and temperature. Novobiocin is used to mimic the consequences of such parameters on DNA topology. By exposing cells to novobiocin, it was demonstrated that DNA relaxation leads to increased expression of the [sigma factor](#) that governs the **heat-shock** response and to a concomitant induction of the production of heat-shock proteins ([23, 24](#)). Similarly, it was shown that DNA **supercoiling** also regulates the expression of genes that respond to osmotic pressure ([25](#)). Furthermore, the inhibition of DNA gyrase by novobiocin leads to an increase in the expression of the genes that encode the enzyme ([26](#)).

#### Bibliography

1. H. Hoeksema, J. L. Johnson, and J. W. Hinnan (1955) *J. Am. Chem. Soc.* **78**, 6710–6711.
2. M. Steffensky, S.-M. Li, B. Vogler, and L. Heide (1998) *FEMS Microbiol. Lett.* **161**, 69–74.
3. J. C. Godfrey and K. E. Price (1972) *Adv. Appl. Microbiol.* **15**, 231–296.

4. M. J. Kennedy et al. (1995) *J. Clin. Oncol.* **13**, 1136–1143.
5. J. P. Eder, C. A. Wheeler, B. A. Teicher, and L. E. Schnipper (1991) *Cancer Res.* **51**, 510–513.
6. M. A. Gellert (1981) *Annu. Rev. Biochem.* **50**, 879–910.
7. J. C. Wang (1985) *Annu. Rev. Biochem.* **54**, 665–697.
8. R. J. Reece and A. Maxwell (1991) *Crit. Rev. Biochem. Mol. Biol.* **26**, 335–375.
9. E. Goetschi et al. (1993) *Pharmacol. Ther.* **60**, 367–380.
10. A. Ali, A. P. Jackson, A. J. Howells, and A. Maxwell (1993) *Biochemistry* **32**, 2717–2724.
11. N. A. Gormley et al. (1996) *Biochemistry* **35**, 5083–5092.
12. N. Nakada, H. Gmuender, T. Hirata, and M. Arisawa (1994) *Antimicrob. Agents Chemother.* **38**, 1966–1973.
13. R. J. Lewis et al. (1996) *EMBO J.* **15**, 1412–1420.
14. F. Reusser and L. Dolak (1986) *J. Antibiot.* **39**, 272–274.
15. F. T. Tsai et al. (1997) *Proteins* **28**, 41–52.
16. D. S. Samuels, R. T. Marconi, W. M. Huang, and C. F. Garon (1994) *J. Bacteriol.* **176**, 3072–3075.
17. G. A. Holdgate et al. (1997) *Biochemistry* **36**, 9663–9673.
18. A. Contreras and A. Maxwell (1992) *Mol. Microbiol.* **6**, 1617–1624.
19. L. Sanchez, W. Pan, M. Vinas, and H. Nikaido (1997) *J. Bacteriol.* **179**, 6855–6857.
20. R. Srikumar, T. Kon, N. Gotoh, and K. Poole (1998) *Antimicrob. Agents Chemother.* **42**, 65–71.
21. A. S. Thiara and E. Cundliffe (1993) *Mol. Microbiol.* **8**, 495–506.
22. J. I. Mitchell, P. G. Logan, K. E. Cushing, and D. A. Ritchie (1990) *Mol. Microbiol.* **4**, 845–849.
23. T. Mizushima et al. (1996) *Mol. Gen. Genet.* **253**, 297–302.
24. F. Sanchez-Lopez, J. Ramirez-Santos, and M. C. Gomez-Eichelmann (1997) *Biochim. Biophys. Acta* **1353**, 79–83.
25. A. Conter, C. Menchon, and C. Gutierrez (1997) *J. Mol. Biol.* **273**, 75–83.
26. S. Neumann and A. Quinones (1997) *J. Basic Microbiol.* **37**, 53–69.
27. I. Del Castillo, J. L. Vizán, M. C. Rodriguez-Sainz, and F. Moreno (1991) *Proc. Natl. Acad. Sci. USA* **88**, 8860–8864.
28. L. Orr, N. F. Fairweather, I. B. Holland, and R. H. Pritchard (1979) *Mol. Gen. Genet.* **177**, 103–112.
29. M. Stieger, P. Angehrn, B. Wohlgensinger, and H. Gmuender (1996) *Antimicrob. Agents Chemother.* **40**, 1060–1062.
30. B. Fournier and D. C. Hooper (1998) *Antimicrob. Agents Chemother.* **42**, 121–128.
31. R. Munoz, M. Bustamante, and A. G. de la Campa (1995) *J. Bacteriol.* **177**, 4166–4170.
32. M. L. Holmes and M. L. Dyall-Smith (1991) *J. Bacteriol.* **173**, 642–648.

### **Suggestions for Further Reading**

33. A. Maxwell (1997) DNA gyrase as a drug target, *Trends Microbiol.* **5**, 102–109.
34. R. J. Lewis (1996) Molecular mechanisms of drug inhibition of DNA gyrase, *Bioessays* **18**, 661–671.
35. D. B. Wigley (1995) Structure and mechanism of DNA topoisomerase, *Annu. Rev. Biophys. Biomol. Struct.* **24**, 185–208.
36. S. Radl (1990) Structure-activity relationships in DNA gyrase inhibitors, *Pharmacol. Ther.* **48**, 1–17.



## Nuclear Envelope

The nuclear envelope separates the [nucleoplasm](#) from the **cytoplasm**. The envelope regulates (1) the import of [proteins](#) into the [nucleus](#), (2) the export of **ribonucleoproteins** (RNP), such as those of [messenger RNA](#) and the large and small subunits of [ribosomes](#), and (3) the timing of [DNA replication](#). The nuclear envelope consists of two [lipid](#) bilayers, the inner and outer nuclear [membranes](#), which together separate the nuclear compartment from the cytoplasm. The outer membrane is continuous with the [endoplasmic reticulum](#) (ER). In fact, [ribosomes](#) are often found on the cytoplasmic surface of the nuclear envelope. The space between the outer and inner membranes is called the [perinuclear space](#). It is continuous with the ER **lumen**. In typical cells, several thousand [nuclear pore complexes](#) perforate the envelope to mediate the **nuclear import** and export of molecules that range in size from small ions to large macromolecular complexes. The inner and outer nuclear membranes fuse to form the pore membrane, a separate membrane domain of the envelope that circumvents the lateral sides of the nuclear pores complex. Just beneath the inner nuclear membrane and in contact with it is the nuclear lamina, a meshwork of **intermediate filament**-type proteins. The lamina in vertebrate somatic cells consists of **lamins** A, B, and C. The lamina provides support and shape to the envelope and is considered a component of the [nuclear matrix](#). Metabolically active cells increase their envelope surface area with folds or blebs that extend into the cytoplasm (eg, the multiple blebs within the **oocyte** nuclear envelope of [Xenopus](#)).

Nuclear pore complexes mediate bidirectional transport between the nucleus and the cytoplasm. They were first described by Callan and Tomlin in electron microscopic examination of amphibian oocyte nuclei (1). Their three-dimensional architecture and arrangement within the envelope were determined by Hinshaw et al. (2). An individual pore complex has a combined mass of approximately 124 megadaltons and consists of four ultrastructural elements: (1) two eight-membered spoke-ring assemblies on the cytoplasmic and nuclear sides; individual pore complexes display eight-fold symmetry, and they measure about 120 nm in diameter; (2) a plug that fills the center of the spoke-ring assembly; it may function as the central transporter; (3) lateral components of the nuclear pore complex that are anchored to the envelope by protein gp210 (3, 4); in addition, aqueous channels form between the pore membrane domain and the sides of the pore complexes (2); (4) filaments that extend from the pore complex into the cytoplasm, plus a basket-like structure that extends into the nucleoplasm on the other side; these peripheral structures may function in docking incoming proteins and outgoing RNP complexes (see [Nuclear Pore Complex](#)).

The continuity between the ER membrane, the outer nuclear membrane, the pore membrane, and the inner nuclear membrane suggests that proteins integrated into the ER membrane might potentially diffuse into separate membrane domains of the nuclear envelope. The outer nuclear membrane, for example, is quite similar to the ER membrane in integral membrane composition, so newly-synthesized proteins can integrate into the ER membrane and then diffuse into the outer nuclear membrane. Yet the pore membrane and the inner membrane remain unique with respect to their integral proteins (5), and how these proteins localize there remains an intriguing problem (6). The inner membrane contains a putative lamin B **receptor** (LBR) of ~58kDa that may mediate the assembly of the lamina and anchor the lamina to the underside of the inner membrane.

The LBR is a good case study of how integral proteins could make their way to the inner nuclear membrane. The LBR is synthesized in the cytoplasm without a cleavable ER **signal sequence**. It integrates into the ER membrane using eight hydrophobic transmembrane **domains** that are positioned well downstream of the amino terminus. Upon integration, its amino terminus of 204 residues (in chickens) remains in the cytoplasm. The receptor is thought to diffuse from the ER membrane to the outer nuclear membrane, past the pore membrane, and finally into the inner

membrane. The pore membrane is thought to form lateral aqueous channels with the pore complex (2). These lateral channels permit **passive transport** of water, ions, and small molecules, of course. More importantly, however, they may be large enough to permit the passage of cytoplasmic protein domains as integral membrane proteins diffuse through the pore membrane domain on their way to the inner membrane. These cytoplasmic domains may be as large as 40 to 60 kDa. Once the receptor diffuses to the inner membrane, its amino terminus faces the nucleoplasm. Then this terminus binds lamin B, DNA, and other nuclear proteins. The amino terminus of the LBR is required for proper localization because the LBR amino terminus is sufficient to localize unrelated (nonenvelope) transmembrane fusion proteins to the inner nuclear membrane (5). The first transmembrane domain of LBR can localize independently to the inner nuclear membrane (7).

Other proteins that localize to the inner nuclear envelope include the *Drosophila* otefin protein (8), the lamina-associated proteins (LAPS), which may play a role in nuclear envelope reassembly (9), and the SREBP-1 protein, a membrane-bound, cholesterol-regulated [transcription factor](#). Like LBR, SREBP-1 has an amino-terminal domain that binds DNA (10). Alternatively, some proteins may localize to the inner membrane by first associating with a cytoplasmic carrier protein via their signal sequence for **nuclear import**. Then the carrier protein delivers the protein to the nuclear pore complex. Other proteins might associate with the inner membrane by first binding the lamina or envelope-associated chromatin, but there are also other possibilities (6).

The nuclear envelope disassembles at **mitosis**, and **phosphorylation** of the lamin proteins by [maturation promoting factor](#) (MPF or [cyclin B/p34<sup>cdc2</sup> kinase](#)) may be the most significant factor. Upon disassembly, the envelope membranes form two populations of **vesicles** (see *later*), the pore complexes disassemble, and the lamina matrix depolymerizes. Lamins A and C remain associated, but soluble. On the other hand, lamin B and the LBR associate with one class of nuclear membrane vesicles, another class derives from the nuclear pore membrane domain, and it retains protein gp210 (11, 12). These disassembled materials disperse into the cytoplasm to form a pool of materials reconstructing functional envelopes in the daughter cells. Reassembly of the envelope occurs at **telophase** in stages (6, 13). First the membrane vesicles containing lamin B bind to the [chromosomes](#), most likely through ligand–receptor types of interactions. Then the vesicles fuse, which requires ATP, GTP, soluble cytoplasmic factors, and free  $Ca^{2+}$  (6). Finally, nuclear pore complexes and the lamina reassemble from pools of disassembled components. As the interphase ensues, the envelope grows by the adding of membrane from the ER and newly synthesized lamins and pore components.

The nuclear envelope may also control the timing of [DNA replication](#) (14). For example, during mitosis when the nuclear envelope normally disassembles, a cytoplasmic licensing factor is thought to gain access to the [chromatin](#). In the next interphase, DNA replication is thought to initiate at the sites where the licensing factor has bound. More licensing factor is synthesized during the interphase, but it is excluded from the nuclear compartment by the nuclear envelope. The licensing factor that originally mediated DNA replication is degraded after the **S phase**, preventing DNA replication a second time. Upon envelope breakdown, licensing factor made in the previous interphase again gains access to the chromatin to mediate the next round of DNA replication. The replication licensing factor has been fractionated into separate components. Chief among these components are the minichromosome maintenance (MCM)/P1 proteins that bind chromatin before the S-phase and are indispensable for DNA replication. As expected for a licensing factor, the MCM/P1 proteins are present at relatively few foci in G1 and the S-phase before DNA replication, and they are excluded from replicated chromatin.

Finally, the nuclear envelope has been traditionally included as a part of the interphase [nuclear matrix](#). In this capacity, the nuclear envelope organizes interphase chromatin to assist nuclear reassembly, regulation of **gene expression**, chromatin condensation, and meiotic chromosome pairing in oocytes.

## Bibliography

1. H. G. Callan and S. G. Tomlin (1950) *Proc. Roy. Soc. B* **137**, 367–378.
2. J. E. Hinshaw, B. O. Carragher, and R. A. Milligan (1992) *Cell* **69**, 1133–1141.
3. E. Hallberg, R. W. Wozniak, and G. Blobel (1993) *J. Cell Biol.* **122**, 513–521.
4. R. W. Wozniak, G. Blobel, and M. P. Rout (1994) *J. Cell Biol.* **125**, 31–42.
5. B. Soullam and H. J. Worman (1993) *J. Cell Biol.* **120**, 1093–1100.
6. C. Wiese and K. L. Wilson (1993) *Curr. Opin. Cell Biol.* **5**, 387–394.
7. B. Soullam and H. J. Worman (1995) *J. Cell Biol.* **130**, 15–27.
8. R. Padan, S. Nainudel-Epszteyn, R. Goitein, A. Fainsod, and Y. Gruenbaum (1990) *J. Biol. Chem.* **265**, 7808–7813.
9. R. Foisner and L. Gerace (1993) *Cell* **73**, 1267–1279.
10. X. Wang, R. Sato, M. S. Brown, X. Hua, and J. L. Goldstein (1994) *Cell* **77**, 53–62.
11. G. P. A. Vigers and M. J. Lohka (1991) *J. Cell Biol.* **112**, 545–556.
12. D. Lourim and G. Krohne (1993) *J. Cell Biol.* **123**, 501–512.
13. N. Chaudhary and J. -C. Courval in (1993) *J. Cell Biol.* **122**, 295–306.
14. J. J. Blow and R. A. Laskey (1988) *Nature* **332**, 546–548.

## Suggestions for Further Reading

15. S. Chevalier and J. J. Blow (1996) Cell cycle control of replication initiation in eukaryotes, *Curr. Opin. Cell Biol.* **8**, 815–821.
16. R. A. Laskey, D. Gorlich, M. A. Madine, J. P. S. Makkerh, and P. Romanowski (1996) Regulatory roles of the nuclear envelope, *Exp. Cell Res.* **229**, 204–211.
17. W. F. Marshall, A. F. Dernburg, B. Harmon, D. A. Agard, and J. W. Sedat (1996) Specific interactions of chromatin with the nuclear envelope: Positional determination within the nucleus in *Drosophila melanogaster*, *Mol. Biol. Cell* **7**, 825–842.

## Nuclear Import, Export

The [nuclear envelope](#) separates the [nucleus](#) from the **cytoplasm** and consists of two [membrane](#) bilayers that are perforated by [nuclear pore complexes](#). Two-way nucleocytoplasmic transport across this barrier is crucial in **eukaryotes** for cells to function and survive. Transport includes the export of [messenger RNA](#), the large and small subunits of [ribosomes](#), [transfer RNA](#), and small nuclear **ribonucleoprotein** particles (pre-snRNP) that continue to assemble in the cytoplasm before reentering the nucleus. Nuclear proteins synthesized in the cytoplasm must translocate into the nucleus through the pores. Passive diffusion of water, small ions, and small molecules (eg, **nucleotides**) must also occur. This article primarily addresses the import of nuclear proteins, but recent evidence suggests that nucleocytoplasmic import and export are two sides of the same coin ([1](#)). They are intimately related by the nuclear pore complex through which import and export must occur and also by the regulatory proteins involved in **energy transduction** (eg, Ran and its associated regulators) and the protein import factor known as importin a.

### 1. Mechanism

The nuclear import of small proteins (<40kDa whose diameters are no greater than 230 Å) is

generally believed to occur by diffusion through the pore complexes, and the migration rate is inversely proportional to size. The transport of [histone](#) proteins and some other small nuclear proteins may be an exception to this general rule, however, because histone H<sub>1</sub> (~21kDa) is believed to cross the nuclear envelope by a facilitated, receptor-mediated **active transport** type of mechanism (2).

Selective import of nuclear proteins larger than 40 kDa, however, requires a nuclear localization signal (NLS) and an energy-dependent mechanism of active transport at the pore complexes. Two types of NLS have been identified: (1) a simple basic motif and (2) a larger, bipartite signal. The best characterized simple NLS is the **SV40 virus** large [T Antigen](#), which has the [amino acid](#) sequence -Pro-(Lys-)<sub>3</sub>Arg-Lys-Val, -PKKKRKV- in one-letter code (3). When fused to cytoplasmic proteins, this simple NLS is sufficient to transport them into the nucleus (4). Bipartite NLS sequences include two short basic motifs that are separated by spacer segments of about 10 amino acid residues, as observed in [nucleoplasmin](#) (5). The NLS functions as a ligand. The rate of nuclear uptake of bovine [serum albumin](#) (BSA) covalently coupled to the SV40 NLS is saturable, thus indicating a ligand/receptor-mediated mechanism of nuclear protein targeting (6).

The mechanism of nuclear protein uptake depends on cytoplasmic factors, specifically, importins a and b (1, 7), also known as karyopherin a and b (8). Importin a binds the NLS-containing proteins directly through its NLS-binding **domain**. Importin a forms a heterodimer with importin b, which then mediates docking the complex to the nuclear pore or to fibers that extend from the pore into the cytoplasm. Docking is temperature-insensitive and energy-independent (9). The importins, however, are sensitive to **N-ethylmaleimide**, a thiol-alkylating compound that blocks the docking reaction. In docking, importin b binds nucleoporins, a family of proteins that reside at or within the pores. Nucleoporins p62/p58/p54 in rat cells, or Nup1 and Nsp1 in yeast, for example, may serve as docking sites for the complex of importin-NLS containing proteins (10). **Posttranslational modification** of the nucleoporins introduces O-linked *N*-acetyl-glucosamine residues (11) (see [O-Glycosylation](#)). The [lectin](#) wheat germ agglutinin (WGA) binds these sugar moieties and inhibits nuclear protein uptake. WGA does not inhibit by preventing the importins from binding nucleoporins, but by preventing the second step, namely, the actual translocation of the karyophilic protein through the pore complex.

Although understanding of the precise translocation mechanism is still in flux, translocation through the pore is known to be energy-dependent. Specifically, the GTP/GDP-binding Ran/TC4 protein is a critical player (see [GTP-Binding Proteins](#)). Interestingly, most of the Ran protein within a cell (80–90%) is contained within the nucleus as GTP-Ran. Its **guanine exchange factor** is the [chromatin](#) protein, RCC1 (regulator of chromatin condensation, Prp20 in yeast). The remaining Ran protein within the cell is cytoplasmic, existing as GDP-Ran. The GTPase activating protein, RanGAP1 (Rna1 in yeast), is also cytoplasmic. Although its precise function and associations remain unknown, this cytoplasmic GDP-Ran is necessary for karyophilic protein translocation through the pores (12). Besides Ran, the Ran-interacting factor, NTF2/p10, is another necessary factor which binds p62 of the pore complex to regulate GTP/GDP binding of Ran (13). In one model for nuclear import (14), GDP-Ran is needed to help tether the complex of importin a/importin b/NLS-bearing karyophilic protein to the pore complex by an interaction between importin b and the nucleoporins. After movement through the pore, this complex is in the nucleus and an environment rich in GTP-Ran. Now GTP-Ran competes with importin a for overlapping binding sites on importin b (8). As a consequence, importin a, with its bound NLS-bearing protein, dissociates from importin b. Importin b remains attached to the nuclear side of the pore and is thought to exit the nucleus quickly. Once within the nucleus, importin a releases the NLS-bearing protein, perhaps in exchange for RNP proteins, such as the 5'-cap binding complex that escorts capped **small nuclear RNAs** out of the nucleus (15). Therefore the export of RNAs may be directly linked to protein import via importin a.

Several nuclear proteins shuttle between the nucleus and the cytoplasm. These include the proteins Nopp140, NO38, and nucleolin of the [nucleolus](#). The significance of shuttling of nucleolar proteins

has yet to be established, but one possibility is that ribosomal subunits are escorted out of the nucleus and that the ribosomal proteins are transported into the nucleus on the return trip and eventually to the nucleolus. In addition to nucleolar proteins, some of the more abundant hnRNA-associated proteins shuttle: namely, A1, A2, D, I, and K (16) (see [Transcription](#)). HnRNP protein A1 contains a nuclear export signal (NES) sequence, a 38-residue domain near the carboxyl end of the protein sequence. It directs export of the protein (17) and also its import (18). The NES of hnRNP protein A1 and of HIV-1 Rev may function in RNA export, but those of other nuclear proteins (eg, hsc70, protein kinase inhibitor) may serve other defined nuclear shuttling activities (19).

### Bibliography

1. D. Görlich, S. Kostra, R. Kraft, C. Dingwall, R. A. Laskey, E. Hartmann, and S. Prehn (1995) *Curr. Biol.* **5**, 383–392.
2. M. Breeuwer and D. S. Goldfarb (1990) *Cell* **60**, 999–1008.
3. D. Kalderon, W. D. Richardson, A. T. Markham, and A. E. Smith (1984) *Nature* **311**, 33–38.
4. D. Kalderon, B. L. Roberts, W. D. Richardson, and A. E. Smith (1984) *Cell* **39**, 499–509.
5. J. Robbins, S. M. Dilworth, R. A. Laskey, and C. Dingwall (1991) *Cell* **64**, 615–623.
6. D. S. Goldfarb, J. Garipey, G. Schoolnik, and R. D. Kornberg (1986) *Nature* **322**, 641–644.
7. D. Görlich, S. Prehn, R. A. Laskey, and E. Hartmann (1994) *Cell* **79**, 767–778.
8. J. Moroianu, G. Blobel, and A. Radu (1996) *Proc. Natl. Acad. Sci. USA* **93**, 7059–7062.
9. M. S. Moore and G. Blobel (1992) *Cell* **69**, 939–950.
10. D. R. Finlay, E. Meier, P. Bradley, J. Horecka, and D. J. Forbes (1991) *J. Cell Biol.* **114**, 169–183.
11. G. W. Hart, G. D. Holt, and R. S. Haltiwanger (1988) *Trends Biochem. Sci.* **13**, 380–384.
12. M. S. Moore and G. Blobel (1994) *Trends Biochem. Sci.* **19**, 211–216.
13. B. M. Paschal and L. Gerace (1995) *J. Cell Biol.* **129**, 1649–1659.
14. D. M. Koepp and P. A. Silver (1996) *Cell* **87**, 1–4.
15. D. Görlich and I. Mattaj (1996) *Science* **271**, 1513–1518.
16. S. Piñol-Roma and G. Dreyfuss (1992) *Nature* **355**, 730–732.
17. W. M. Michael, M. Choi, and G. Dreyfuss (1995) *Cell* **83**, 415–422.
18. G. Dreyfuss, M. Hentze, and A. I. Lamond (1996) *Cell* **85**, 963–972.
19. L. Gerace (1995) *Cell* **82**, 341–344.

### Suggestions for Further Reading

20. F. Melchior and L. Gerace (1995) Mechanisms of nuclear protein import, *Curr. Opin. Cell Biol.* **7**, 310–318.
21. N. Pante and U. Aebi (1996) Toward the molecular dissection of protein import into nuclei, *Curr. Opin. Cell Biol.* **8**, 397–406.

### Nuclear Matrix

The nuclear matrix is a three-dimensional fibrogranular latticework of [ribonucleoprotein \(RNP\)](#) and [chromatin](#) that permeates the [nucleus](#) of **interphase** cells. Its overall function is to organize the nuclear architecture and to regulate genetic function. The precise nature of the matrix has been

elusive, perhaps because of its multiple interactions in [DNA replication](#), **gene transcription**, [RNA splicing](#), and [messenger RNA](#) (mRNA) transport. Although its precise functions remain uncertain, significant progress has been made in ultrastructural and molecular analyses to provide at least strong evidence for its existence. The first mention of the nuclear matrix was made by Fawcett (1) to describe the interchromosomal material that previously had been called simply “nuclear sap.” The RNP fibrogranular portion of the matrix was first defined by the **EDTA**-regressive staining technique developed in the late 60's and early 70's to distinguish **RNA**-containing nuclear structures from **DNA**-containing chromatin (2). Berezney and Coffey (3-5) were the first to isolate the nuclear matrix and then to characterize potential matrix proteins by **SDS-PAGE** analysis. The matrix has been defined operationally, and several modifications to the original technique, developed in the last twenty years, yield different matrix characteristics. Consequently, there are corresponding differences in descriptions of the matrix.

One particular method of matrix preparation, significant for retaining accurate nuclear morphology, was developed by Deppert (6) and by Penman (7). The nuclear matrix is prepared by treating whole cells grown on coverslips with **Triton X-100**, extracting the genomic chromatin with **DNase I** digestion, and then reextracting the cells with 0.2 to 0.6 M **ammonium sulfate** (8). Such *in situ* preparations, as they are called, retain the nuclear lamina, the residual [nucleolus](#), and the fibrogranular RNP matrix. In addition, Penman and co-workers (7) perfected whole mount and resinless thick-section techniques for [electron microscopy](#) to assess the three-dimensional matrix. Another modification included using the [detergent](#) 3,5-diiodosalicylate (LIS) to extract nuclei (9, 10). Such preparations are called nuclear scaffolds, rather than matrices, because LIS extraction replaces the high salt extraction, and minimizes protein rearrangements while efficiently removing [histone](#) proteins.

The matrix is ubiquitous in eukaryotic cells ranging from yeast to human. Ultrastructurally, the matrix consists of at least the residual [nuclear envelope](#) and its **nuclear pores** and lamina, the residual [nucleolus](#), certain regions of [euchromatin](#) that attach to the matrix, and the fibrogranular RNP matrix, which is often called interchromatin fibers and perichromatin granules. The matrix also contains 10 to 25% of the total nuclear proteins, and some tightly bound DNA and RNA remain after **nuclease** treatment and extraction (8, 11).

Using **two-dimensional gel** analyses, Fey and Penman (12) detected more than 200 potential nuclear matrix proteins. Some of these are common to all cell types. Others are specific to a particular type of cell or tissue or to particular stages of differentiation, including oncogenesis. Among the common nuclear matrix proteins are the three **intermediate filament**-like [lamin](#) proteins, A, B, and C. The nuclear mitotic apparatus protein (NuMA), of more than 200 kDa, may be an important structural component of the matrix (13-16). Like intermediate filament proteins, NuMA can form a [coiled-coil](#) structure (17, 18). Other common proteins include the **hnRNP U** protein, nuclear [actin](#), the nucleolar protein B23, and a set of related proteins called matrins (8). Protein [kinases](#) A, B, C, and CKII have been linked to the nuclear matrix, suggesting regulatory roles for the matrix in [signal transduction](#) (8). **Topoisomerase II** and an [attachment region binding](#) protein (ARBP) are matrix proteins that associate with the chromatin (see later).

Several studies have shown that different cell types share these common proteins, but that they also contain their own unique matrix proteins (12, 19-21). The protein composition of the matrix changes upon [differentiation](#) of fetal rat calvarial osteoblasts (22-25) and when they become malignant (26-28).

In their studies on mitotic chromosomal structure, Laemmli's laboratory demonstrated, after nuclease treatment and high salt extraction (29), that mitotic chromosomes consist of residual proteinacious scaffolds and that large radial chromatin loops are anchored to these scaffolds (30). As mitotic chromosomes decondense in the early interphase, these radial loops are thought to maintain their attachment to the interphase nuclear matrix. The DNA elements that attach to the mitotic scaffold and to the interphase matrix are AT-rich sequences that easily unwind with the introduction of

negative **supercoils** (see [Matrix Attachment Regions](#), or MARS). A special **AT** binding protein (SATB1) specifically binds the unwound AT-rich regions of the MARS ([31](#)), perhaps to anchor the elements to the matrix.

Several functional roles have been ascribed to the nuclear matrix. DNA replication has been linked with the matrix, because **DNA polymerase**  $\alpha$ , [primase](#), and other DNA replication factors have been found in prepared matrices ([8](#)). An intriguing hypothesis suggests that between 50,000 to 100,000 bidirectional DNA replication units ([replicons](#)) are affixed to the nuclear matrix in relatively few (150 to 300) clusters, called clustersomes ([8](#)). Bidirectional replication occurs as looped chromatin domains reel through the stationary replication machinery ([8](#), [32](#)). The clustersome model was supported by studies in which **biotin**-labeled dUTP was incorporated into newly synthesized DNA during the **S-phase** of permeabilized culture cells. Staining with [streptavidin](#) tagged with a fluorescent marker revealed distinct replicative granules whose number and size matched the values predicted for clustersomes. To demonstrate that replicons are in fact attached to the matrix, matrices were prepared *in situ* before staining with streptavidin. The staining pattern was identical to that observed in the nonextracted cells, thus indicating that the replicative factories are associated with defined matrices. Interestingly, the number and size of the stained replicative granules are again similar after DNA synthesis is carried out in matrices that retained chromatin fragments.

The nuclear matrix has also been linked to [transcription](#) ([33](#)), primarily because several [transcription factors](#) have been found associated with the nuclear matrix ([34](#)) and transcription foci have been immunolocalized to **HeLa** cell nuclear matrices ([35](#)). Analogous to the stationary DNA replicative factories described in the preceding paragraph, stationary transcription factories have also been proposed ([35](#)).

Splicing of heterogeneous nuclear RNA (hnRNA) has been intimately linked to the nuclear matrix. Interchromatin fibers and perichromatin granules constitute the bulk of the matrix of the interphase nucleus, and it was found by immunolocalization that they contain many of the splicing factors, such as **small nuclear RNAs** (snRNA) with trimethyl caps, snRNPs, and the SR family of proteins ([36](#)) (see [RNA Splicing](#)). Over 70% of the hnRNA is retained in matrices prepared in the presence of **ribonuclease inhibitors** ([11](#)). Many of the hnRNP proteins originally identified by Dreyfuss et al. ([37](#)) (eg, proteins A2, B1, C1/C2) constitute the matrix ([38](#)). The hnRNP U protein is abundant in prepared nuclear matrices, and in addition to its role in processing nascent heterogeneous nuclear RNA, it preferentially binds DNA MAR elements. For this reason, the hnRNP U protein has been called the scaffold attachment factor (SAF-A) ([38](#)).

Finally, the matrix may play a role in the actual transport of transcripts from their site of synthesis on the gene to the [nuclear pore complex](#) (another component of the matrix), where they are transported out of the nucleus and into the **cytoplasm**. **Exons** and **introns** of specific transcripts have been detected on curvilinear tracks within the nucleus ([39](#), [40](#)) and in prepared matrices ([41](#)). In all likelihood these curvilinear tracks consist of the interchromatin fibrils observed in thin-section electron micrographs.

## Bibliography

1. D. W. Fawcett (1966) *An Atlas of Fine Structure. The Cell, its Organelles and Inclusions*, W. B. Saunders, Philadelphia.
2. A. Monneron and W. Bernhard (1969) *J. Ultrastruct Res.* **27**, 266–288.
3. R. Berezney and D. S. Coffey (1974) *Biochem. Biophys. Res. Commun.* **60**, 1410–1417.
4. R. Berezney and D. S. Coffey (1975) *Science* **189**, 291–293.
5. R. Berezney and D. S. Coffey (1977) *J. Cell Biol.* **73**, 616–637.
6. M. Staufenbiel and W. Deppert (1984) *J. Cell Biol.* **98**, 1886–1894.
7. E. G. Fey, G. Krochmalnic, and S. Penman (1986) *J. Cell Biol.* **102**, 1654–1665.
8. R. Berezney (1996) In *Nuclear Structure and Gene Expression* (R. C. Bird, G. S. Stein, J. B.

- Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 35–85.
9. J. Mirkovitch, M. F. Mirault, and U. K. Laemmli (1984) *Cell* **39**, 223–232.
  10. S. M. Gasser and U. K. Laemmli (1987) *Trends Genet.* **3**, 16–22.
  11. S. Penman, B. J. Blencowe, and J. A. Nickerson (1996) In *Nuclear Structure and Gene Expression* (R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 3–36.
  12. E. G. Fey and S. Penman (1988) *Proc. Natl. Acad. Sci. USA* **85**, 121–125.
  13. C. Zeng, D. He, and B. R. Brinkley (1984) *Cell Motil. Cytoskeleton* **29**, 167–176.
  14. C. Zeng, D. He, S. M. Berget, and B. R. Brinkley (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1505–1509.
  15. D. He, J. A. Nickerson, and S. Penman (1990) *J. Cell Biol.* **110**, 569–580.
  16. D. He, C. Zheng, and B. R. Brinkley (1995) *Int. Rev. Cytol.* **162B**, 1–74.
  17. D. A. Compton, I. Szilak, and D. W. Cleveland (1992) *J. Cell Biol.* **116**, 1395–1408.
  18. C. H. Yang, E. J. Lambie, and M. Snyder (1992) *J. Cell Biol.* **116**, 1303–1317.
  19. N. Stuurman, R. van Driel, L. de Jong, A. M. Meijne, and J. van Renswoude (1989) *Exp. Cell Res.* **180**, 460–466.
  20. N. Stuurman, A. M. Meijne, A. J. van der Pol, L. De Jong, R. van Driel, and J. van Renswoude (1990) *J. Biol. Chem.* **265**, 5460–5465.
  21. R. M. Getzenberg and D. S. Coffey (1990) *Mol. Endocrinol.* **4**, 1336–1342.
  22. G. S. Stein, A. van Wijnen, J. L. Stein, J. B. Lian, J. P. Bidwell, and M. Montecino (1994) *J. Cell Biochem.* **55**, 4–15.
  23. S. I. Dworetzky, E. G. Fey, S. Penman, J. B. Lian, J. L. Stein, and G. A. Stein (1990) *Proc. Natl. Acad. Sci. USA* **87**, 4605–4609.
  24. S. Dworetzky, K. L. Wright, E. G. Fey, S. Penman, J. B. Lian, J. L. Stein, and G. S. Stein (1992) *Proc. Natl. Acad. Sci. USA* **89**, 4178–4182.
  25. J. P. Bidwell, A. J. van Wijnen, E. G. Fey, S. Dworetzky, S. Penman, J. L. Stein, J. B. Lian, and G. S. Stein (1993) *Proc. Natl. Acad. Sci. USA* **90**, 3162–3166.
  26. A. W. Partin, R. H. Getzenberg, M. J. Carmichael, D. Vindivich, J. Yoo, J. I. Epstein, and D. S. Coffey (1993) *Cancer Res.* **53**, 744–746.
  27. P. S. Khanuja, J. E. Lehr, H. D. Soule, S. K. Gehani, A. C. Noto, S. Choudhury, R. Chen, and K. J. Pienta (1993) *Cancer Res.* **53**, 3394–3398.
  28. S. K. Keesee, M. Meneghini, R. P. Szaro, and Y.-J. Wu (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1913–1916.
  29. K. W. Adolph, S. M. Cheng, J. R. Paulson, and U. K. Laemmli (1977) *Proc. Natl. Acad. Sci. USA* **74**, 4937–4941.
  30. J. R. Paulson and U. K. Laemmli (1977) *Cell* **12**, 817–828.
  31. T. Kohwi-Shigematsu and Y. Kohwi (1996) In *Nuclear Structure and Gene Expression* (R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 111–144.
  32. P. R. Cook (1991) *Cell* **66**, 627–635.
  33. G. S. Stein, A. J. van Wijnen, J. L. Stein, J. B. Lian, and M. Montecino (1996) In *Nuclear Structure and Gene Expression* (R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 177–214.
  34. A. J. van Wijnen, J. P. Bidwell, E. G. Fey, S. Penman, J. B. Lian, J. L. Stein, and G. S. Stein (1993) *Biochemistry* **32**, 8397–8402.
  35. D. A. Jackson, A. B. Hassan, R. J. Errington, and P. R. Cook (1993) *EMBO J.* **12**, 1059–1065.
  36. D. L. Spector (1993) *Ann. Rev. Cell Biol.* **9**, 265–315.
  37. G. Dreyfuss, Y. D. Choi, and S. A. Adam (1984) *Mol. Cell Biol.* **4**, 1104–1114.
  38. K. A. Mattern, R. van Driel, and L. de Jong (1996) In *Nuclear Structure and Gene Expression*



(R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds.), Academic Press, New York, pp. 87–110.

39. S. Huang and D. L. Spector (1991) *Genes Develop.* **5**, 2288–2302.

40. Y. Xing, C. V. Johnson, P. R. Dobner, and J. B. Lawrence (1993) *Science* **259**, 1326–1330.

41. Y. Xing and J. B. Lawrence (1991) *J. Cell Biol.* **112**, 1055–1063.

### Suggestions for Further Reading

42. R. Berezney and D. S. Coffey (1976) The nuclear protein matrix: Isolation, structure, and functions, *Adv. Enzyme Regul.* **14**, 63–100.

43. R. Berezney and K. W. Jeon (1995) "Structural and functional organization of the nuclear matrix", In *International Review of Cytology, A Survey of Cell Biology*, Vols. **162A** and **162B**, Academic Press, San Diego. A compilation of excellent reviews on the nuclear matrix.

44. C. R. Bird, G. S. Stein, J. B. Lian, and J. L. Stein (1996) *Nuclear Structure and Gene Expression*, Academic Press, New York.

### Nuclear Overhauser Effect (NOE)

Regardless of the complexity of the experimental design that produced them, the intensities of signals in a nuclear magnetic resonance (NMR) spectrum depend on the populations of molecules in the various allowed nuclear spin energy states immediately before the analyzing RF pulse is applied. In the simplest case, these populations are the Boltzmann populations, characteristic of the system at thermal equilibrium. Under such conditions, the relative intensities of signals from various chemically different groups of molecule spins are proportional to the relative number of spins of each type. But in a variety of circumstances, including an experimental design more complex than a single RF pulse, the number of molecules in some or all of the permitted nuclear energy states may differ from those present at thermal equilibrium. In this event, the intensity of a signal observed in the spectrum will be changed. The change in signal intensity is called a nuclear Overhauser effect, or NOE. There are several ways of expressing the NOE quantitatively. Most commonly used (1) is the equation

$$f_I\{S\} = I_p/I_0 - 1 \quad (1)$$

where  $f_I\{S\}$  indicates the NOE on the signal(s) from spin I when there is a perturbation of level populations associated with spin S.  $I_0$  is the normal intensity of the signal(s) for spin I observed from a sample at thermal equilibrium before the analyzing RF pulse, and  $I_p$  is the intensity of the same signal when there has been a perturbation of spin S before the analyzing pulse. The value of  $f_I\{S\}$  depends on the gyromagnetic ratios of spins I and S, how these spins move in the sample, the strength of the magnetic field used for the NMR experiment, and the details of how the energy level populations associated with spin S are perturbed during the course of an experiment. Importantly,  $f_I\{S\}$  depends on the distance between spins I and S. Values of  $f_I\{S\}$  ranging from 0.5 to –1.0 are possible when both I and S are protons. The first value corresponds to a 50% enhancement of the intensity associated with spin I, and the latter represents the disappearance of the signal intensity for this spin.

A change in NMR signal intensity can arise in molecular situations where two or more spins interact strongly with each other, such that perturbation of the level populations of one of the spins by RF pulses or by some other way ultimately results in the perturbation of the level populations corresponding to the other spin. It should be apparent that the NOE is strongly related to relaxation, those natural processes in any sample that tend to return the populations of the spin energy levels to their Boltzmann (thermal) equilibrium values if these level populations have been altered in some way. (See [NMR \(Nuclear Magnetic Resonance\)](#).) Consequently, the formation and disappearance of a NOE will be time-dependent.

The power of experiments that produce NOEs lies in that the bulk of the interactions that lead to relaxation are strongly dependent on the distance between interacting spins. Thus, observation of NOEs can provide internuclear distance data that ultimately can be used to deduce information about the three-dimensional structures of biopolymers. The distance dependence of the NOE is strong. In the case of  $^1\text{H}$ - $^1\text{H}$  interactions, NOEs are typically observed only when the distance between the spins is 0.22–0.55 nm, with the effect observed at 0.55 nm being less than 1% of that when the internuclear distance is 0.22 nm.

When molecules move slowly in the sample, an additional phenomenon complicates the NOE experiment. “Slow motion” conditions in this context arise when molecules have masses in excess of about 5 kDa or when solutions are highly viscous. The complication is called spin diffusion; it occurs because the magnetization transfers or level population alterations mentioned above become rapid between all protons. Under these conditions, magnetization that originates with one spin of the molecule may be transferred not only to a nearby proton but also to more distant protons. Thus, distance information that is implicit in the NOE may be obscured.

It is also possible to create NOEs because of interactions between various coherences that can be formed as a result of the application of RF pulses to the sample. These experiments, called rotating-frame NOEs, or ROEs, exhibit a different dependence on the details of molecular motions than the NOEs described above. The NOE and ROE experimental results are thus complementary. The formation of the ROE is also time-dependent and must compete with the decay of spin coherence by transverse relaxation. (See also [NOESY Spectrum](#), [ROESY Spectrum](#), [Distance Geometry](#), [Simulated Annealing](#).)

#### Bibliography

1. D. Neuhaus and M. P. Williamson (1989) *The Nuclear Overhauser Effect in Structural and Conformational Analysis*, V.C.H., New York.

#### Suggestions for Further Reading

2. J. H. Noggle and R. E. Schirmer (1971) *The Nuclear Overhauser Effect*, Academic, New York.
3. R. Freeman (1997) *Spin Choreography*, Spektrum, Oxford, Chapter "9", pp. 295–313.
4. H. Mo and T. C. Pochapsky (1997) *Prog. NMR Spectrosc.* **30**, 1–38.

## Nuclear Pore Complex

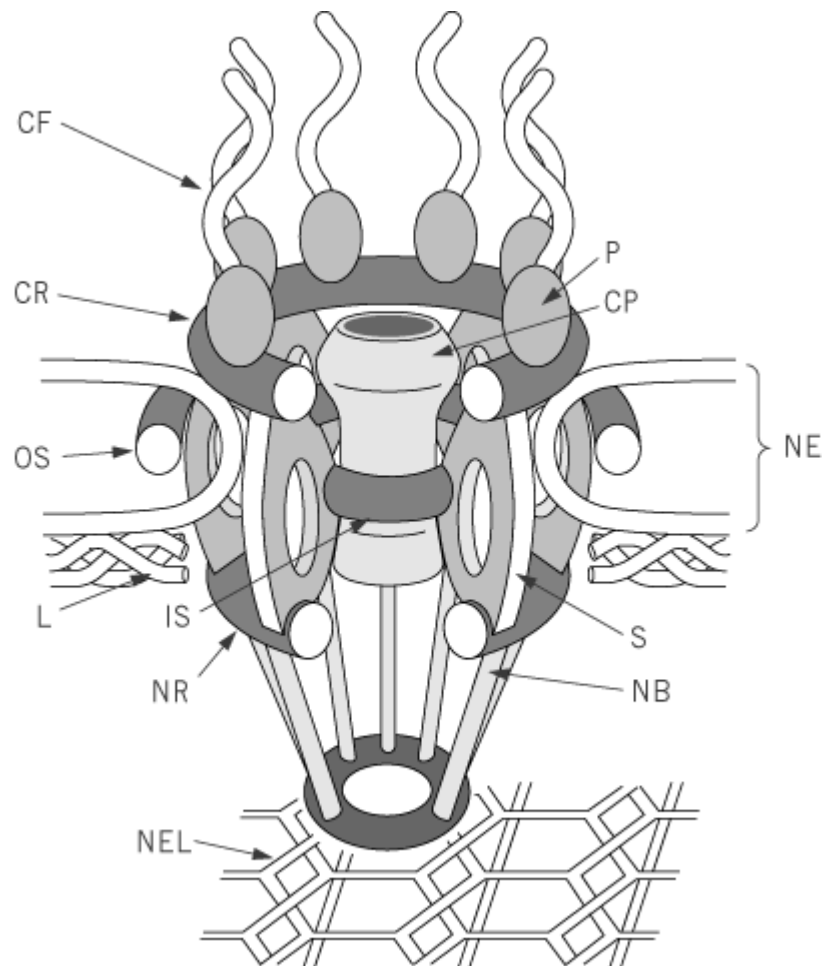
The nuclear pore complex is the **organelle** within the [nuclear envelope](#) that mediates the export of [messenger RNA](#), [transfer RNA](#), and the subunits of [ribosomes](#), which are synthesized or assembled in the [nucleus](#) but function in the **cytoplasm**. The pore complex also mediates the **nuclear import** of

[proteins](#) and protein complexes (eg, snRNP particles) that are synthesized and assembled in the **cytoplasm** but destined to function in the nucleus. The number of nuclear pore complexes varies considerably, depending on the cell type. Their number generally reflects the biosynthetic activity of the cells. For example, the nuclear envelope of the *Xenopus oocyte* has the highest known concentration of pores per unit area of nuclear envelope ( $60 \text{ pores mm}^{-2}$ ), whereas quiescent chick **erythrocytes** have two to four pores  $\text{mm}^{-2}$ .

## 1. Pore Structure

Nuclear pore complexes were first described by Callan and Tomlin (1) in their electron microscopic analysis of the nuclear envelope in amphibian oocytes. Architecturally, the nuclear pore is a huge (124 megadaltons), beautifully elaborate, symmetrical organelle within the nuclear envelope (Fig. 1). The entire complex can be graphically divided into four main component types: (1) eight-membered ring complexes, (2) the central plug, (3) cytoplasmic filaments, and (4) a filamentous nuclear basket (cage). There are two similar eight-membered ring complexes per pore, which associate through apparently **homotropic** interactions to create twofold symmetry of the complex within the plane of the envelope and are designated as the nucleoplasmic and cytoplasmic rings. The ring complexes are anchored to the nuclear envelope where the inner and outer [membranes](#) fuse together. Each ring consists of eight spoke-ring subunits of 10 to 25 nm arranged in eightfold symmetry when viewed from the surface of the envelope (3). Each spoke-ring segment consists of subcomponents that contribute to the formation of the overall complex. Some contribute to formation of the rings, and columnar and luminal portions of each spoke-ring segment help anchor the ring complexes within the nuclear envelope. A spoke protrudes from each columnar subcomponent to help form the central annulus. Spokes from the cytoplasmic ring interact with spokes from the nucleoplasmic ring. The central annulus forms an hourglass-shaped channel of 42 nm at the peripheries and between 9 and 10 nm at the narrowest point. This channel transports macromolecular complexes as large as 23 nm (see later). Eight aqueous channels form on the periphery of the complex between the columnar subcomponents and the nuclear envelope. These aqueous channels permit passive diffusion of ions and small molecules across the envelope. Extending into the cytoplasm from the cytoplasmic ring are eight long, cylindrical filaments that may serve as docking sites for proteins that must enter the nucleus. Extending into the [nucleoplasm](#) from the nucleoplasmic ring are eight filaments that join together at their ends by yet another component to form a distal annulus within the nucleoplasm. This annulus may contact the [nuclear matrix](#) to play a role in the export of nuclear products.

**Figure 1.** A schematic model of the vertebrate nuclear pore complex. The cytoplasmic face is shown uppermost, and certain facing portions have been omitted for clarity. Structures are drawn approximately to scale with the pore transverse the membrane of  $\sim 90 \text{ nm}$  thickness. The diameter of the outer spoke ring (OS) is  $\sim 120 \text{ nm}$ . Filamentous structures (CF and NB) extend from both the cytoplasmic and nuclear faces, with the former attached to particles (P) and the latter forming a basketlike structure. The overall architecture is octagonal: based on eight equally spaced spokes (S) that are sandwiched between two rings (CR and NR). The spokes are attached to an inner spoke ring (IS) that encompasses a central plug (CP). The buttresses connecting the CP with the IS are omitted. The spokes appear to transverse the pore membrane for connection to an outer ring (OS) in the nuclear envelope (NE) lumen. Well defined lamina (L) and nuclear envelope lattice (NEL) structures also seem to have connections to the nuclear pore complex. (Figure from Rout and Wentz, ref. 2, 1994. Reprinted with permission).



At least 100 different proteins (nucleoporins) are thought to constitute the nuclear pore complex, and approximately forty have been identified. They can be subdivided into pore membrane proteins (poms) and phenylalanine/glycine-rich (FG) proteins. Poms are integral nuclear membrane proteins closely associated with the pore complex, such as the 16 to 24 copies of protein gp210 that reside around the periphery of the pore complex. They anchor the pore complex within the envelope, and the bulk of the gp210 resides within the lumen of the envelope. Other poms include the yeast Pom152p and the mammalian pom121. FG proteins found in yeast, such as *NUP1*, *NUP2*, and *NSP1*, are so-called because of their repeated -X-Phe-X-Phe-Gly- (or -XFXFG-) sequences. A related -Gly-Leu-Phe-Gly- (or -GLFG-) repeated sequence is also found in yeast proteins *NUP100*, *NUP116*, *NUP49*, and *NUP145*. The functions of the repeated motifs remain unknown. Deleting them does not perturb the cell's viability nor localization of the protein to the complex. Several metazoan nucleoporins also have been identified. Some of these have been localized to the spoke-rings that face the cytoplasm and their extending filaments (nup214/CAN, p180, and p75), some to the central plug or spoke-ring structures (p62, p58, p54), and at least one to the nuclear basket (nup153). Interestingly, outside of the FG-repeat domains, no extensive homologies have been found between yeast and metazoan nucleoporins. Metazoan FG nucleoporins are modified **posttranslationally** by addition of *N*-acetylglucosamine through [O-glycosylation](#) of serine and threonine residues. The function of this glycosylation also remains unknown, but it makes it possible to probe specifically for nucleoporins and to block transport using the wheat germ agglutinin (WGA), a [lectin](#) that specifically binds O-linked *N*-acetylglucosamine.

## 2. Pore Function

The pore complex transports macromolecules into and out of the nucleus and permits diffusion of

small molecules across the nuclear envelope (see [Nuclear Import, Export](#)). Briefly, proteins that are destined for the nucleus translocate across the envelope by a **facilitated diffusion** mechanism, as in the case of small proteins such as [histones](#) (4), or an **active transport**, ATP-driven mechanism, as in the case of proteins larger than approximately 40 kDa. Large proteins usually have a short nuclear localization signal (NLS) within their linear amino acid sequence. Two types of signals are known. One is a simple basic sequence, rich in [lysine](#) residues, typified by the NLS found in the **SV40 virus** large [T Antigen](#) (namely, -Pro-Pro-Lys-Lys-Lys-Arg-Lys-Val-). Fusion proteins that consist of a normal cytoplasmic protein linked to this NLS translocate into the nucleus. Another characterized NLS is a bipartite sequence in which two short segments of basic residues are separated by approximately ten amino acid residues (5).

Transport of nuclear proteins harboring NLS sequences occurs in two steps (6). The first step is an energy-independent docking that occurs initially with the cytoplasmic filaments extending from the envelope pore complex and then with the central annulus. Cytoplasmic factors, known as importins a and b (7, 8), bind NLS sequences of proteins destined to function in the nucleus and deliver them to the filaments. Association occurs between the importins and the FG nucleoporins described previously. The docking protein may be nup214 of the filaments. Docking can be blocked by WGA, which specifically binds the O-linked *N*-acetylglucosamine residues on nucleoporins that constitute either the filaments (ie, nup214) or the central annulus of the pore complex. The second step is translocation of the docked proteins, which requires GTP hydrolysis (9-11). Translocation requires Ran/TC4, a Ras-like [GTPase](#), and its specific [guanine nucleotide exchange factor](#), called RCC1. The precise mechanism of the active transport remains uncertain, however, and other factors may be required, such as the recently described nuclear factor II (NTFII), which transiently associates with the pore complex protein p62 and Ran/TC4 (12). ATP hydrolysis is necessary for nuclear protein uptake, but again, the actual reason remains uncertain.

The nuclear pore complex is also required for nuclear export of [messenger RNA](#) (mRNA), [transfer RNA](#) (tRNA), ribosomal subunits, and prespliceosomal complexes that mature in the cytoplasm before reentering the nucleus (13, 14). The export mechanisms of these gene products are still largely unknown. Although export of mRNA must include docking, unwinding, and then translocation, it is known that the mRNA leaves the nucleus bound to protein. The nuclear cage must be intimately involved in the export process (15). Blobel's "gene gating" hypothesis (16) suggests that nuclear pores and the nuclear lamina influence (by specific binding of [chromatin](#)) the three-dimensional gene organization within the nucleus, so that transcribed genes are positioned for efficient transport of their RNA products through the pores. Tracks of **polyA**-containing mRNA have been detected in nuclei between the sites of RNA synthesis and the nuclear pores, lending support to the gating hypothesis. These tracks may be the perichromatin fibrils, which contain several **RNA splicing** components (17). Despite these uncertainties, pre-mRNA does not interact with the pore complex until all the splicing reactions are complete. HnRNP proteins (eg, A1, A2, and B1) may play a positive role in this export because a non-RNA-binding domain has been identified in their linear amino acid sequence as a nuclear export signal (NES) (14). When fused to unrelated proteins, the NES in A1 (called M9) confer bidirectional import and export (shuttling) of the fusion protein (18, 19). Other nuclear proteins contain export signals that are unrelated to the M9 sequence. For example, the activation domain of the HIV Rev protein is leucine-rich and is a strong NES that mediates the export of Rev with its bound, *unspliced* HIV RNA (14). The NES in Rev is similar to the NES in [transcription factor](#) IIIA (TFIIIA), which binds and exports 5S rRNA (20).

The nuclear pore complex disassembles into individual components during **mitosis**. Pore disassembly is concomitant with mitotic disassembly of the nuclear lamina and the vesiculation of the envelope membranes. Pore reassembly begins in **telophase** and is initiated by nucleoporins and poms aggregating within the plane of the re-forming envelope. Pore assembly itself is probably incremental, and individual components arrive at the site of reassembly from a stored mitotic pool. Fusion of the inner and outer membranes occurs at the site of pore assembly. Later in the [cell cycle](#) during the **S-phase**, pore complexes are added as the nuclear envelope grows in size. Again addition may be incremental as in telophase. Alternatively, cytoplasmic annulate lamellae, which consist of

membrane and pore-like structures, may contribute additional membrane and pores to the growing nuclear envelope.

### Bibliography

1. H. G. Callan and S. G. Tomlin (1950) *Proc. Roy. Soc.* **B 137**, 367–378.
2. M. P. Rout and S. R. Wente (1994) Pores for thought: nuclear pore complex proteins. *Trends in Cell Biol.* **4**, 357–365.
3. J. G. Gall (1967) *J. Cell Biol.* **32**, 391–399.
4. M. Breeuwer and D. S. Goldfarb (1990) *Cell* **60**, 999–1008.
5. J. Robbins, S. M. Dilworth, R. A. Laskey, and C. Dingwall (1991) *Cell* **64**, 615–623.
6. M. S. Moore and G. Blobel (1992) *Cell* **69**, 939–950.
7. D. Görlich, S. Prehn, R. A. Laskey, and E. Hartmann (1994) *Cell* **79**, 767–778.
8. M. A. Powers and D. J. Forbes (1994) *Cell* **79**, 931–934.
9. F. Melchior, B. Paschal, J. Evans, and L. Gerace (1993) *J. Cell Biol.* **123**, 1649–1659.
10. M. S. Moore and G. Blobel (1993) *Nature* **365**, 661–663.
11. D. J. Sweet and L. Gerace (1995) *Trends Cell Biol.* **5**, 444–447.
12. B. M. Paschal and L. Gerace (1995) *J. Cell Biol.* **129**, 925–937.
13. E. Izaurralde and I. W. Mattaj (1995) *Cell* **81**, 153–159.
14. G. Dreyfuss, M. Hentze, and A. I. Lamond (1996) *Cell* **85**, 963–972.
15. E. Kiseleva, M. W. Goldberg, B. Daneholt, and T. D. Allen (1996) *J. Mol. Biol.* **260**, 304–311.
16. G. Blobel (1985) *Proc. Natl. Acad. Sci. USA* **82**, 8527–8529.
17. D. L. Spector (1993) *Ann. Rev. Cell Biol.* **9**, 265–315.
18. W. M. Michael, M. Choi, and G. Dreyfuss (1995) *Cell* **83**, 415–422.
19. W. M. Michael, H. Siomi, M. Choi, S. Pinol-Roma, S. Nakielny, Q. Liu, and G. Dreyfuss (1995) *Cold Spring Harbor Symp. Quant. Biol.* **60**, 663–668.
20. U. Guddat, A. H. Bakken, and T. Pieler (1990) *Cell* **60**, 619–628.

### Suggestions for Further Reading

21. L. I. Davis (1995) The nuclear pore complex, *Ann. Rev. Biochem.* **64**, 865–896.
22. W. W. Franke, U. Scheer, G. Krohne, and E.-D. Jarasch (1981) The nuclear envelope and the architecture of the nuclear periphery, *J. Cell Biol.* **91**, 39s–50s.
23. L. Gerace and B. Burke (1988) Functional organization of the nuclear envelope, *Ann. Rev. Cell Biol.* **4**, 335–374.
24. M. W. Goldberg and T. D. Allen (1995) Structural and functional organization of the nuclear envelope, *Curr. Opin. Cell Biol.* **7**, 301–309.
25. M. Stewart and W. D. Clarkson (1996) Nuclear pores and macromolecular assemblies involved in nucleocytoplasmic transport, *Curr. Opin. Struct. Biol.* **6**, 162–165.

### Nucleoids

Nucleoids are prepared cell **nuclei** used to discern the structure and function of the [nucleus](#). They have been used to characterize nuclear [chromatin](#), nascent **RNA**, nuclear [proteins](#), and nuclear structures, such as the [nuclear matrix](#), lamina, and [nuclear pore complex](#). Originally, the term

“nucleoid” was used to describe the nuclear body within **bacteria** that contains the centrally located [chromosomes](#). These bacterial nucleoids were isolated from **spheroplasts** using solutions containing 0.1 to 2.0 *M* salt and a chelating agent (1, 2). Subsequently, the term has come to include nucleoids from **eukaryotes** prepared similarly. Eukaryotic nucleoids have been isolated by gently lysing cells with 0.4 to 2.0 *M* NaCl, the chelating agent [EDTA](#) (133 *mM*), and the nonionic [detergent Triton X-100](#) (0.67%, *v/v*). Cell lysis gently occurs over a 15-min period, and the nucleoids are enriched by **sucrose gradient centrifugation** (3). A “nucleoid cage” is formed by remnants of the [cytoskeleton](#). The plasma membrane and cytoplasmic **organelles** are absent, but within the nucleoid cage is the nuclear cage. This inner nuclear cage consists of the nuclear envelope and the entire [genome](#), although the [histone](#) proteins may be lost, depending on the salt concentrations used during lysis. Although intact, approximately 55% of the genome spills out of the nuclear material to form a shroud of intact chromatin loops that surrounds the nucleoid cage. The nuclear cage of the nucleoid contains residual **nucleoli**, nuclear particles, and fibrous aggregates that may be related to the chromosome scaffold described by Paulson and Laemmli (4).

When nucleoids are treated with low concentrations of [ethidium bromide](#) (5  $\mu\text{g/mL}$ ), negative **supercoils** within the genomic DNA relax, and the DNA diffuses from the nucleoid to form a fluorescent halo around the nucleoid periphery. When the ethidium concentration is increased to 100  $\mu\text{g/mL}$ , however, the halo of DNA is resorbed back into the nucleoid (5), because of the introduction of positive supercoils (6). If the DNA is nicked before adding excess ethidium bromide, it fails to retract back into the nucleoid, presumably because it is impossible to introduce positive supercoils. The observation that chromatin spools out and then back again as the DNA is positively supercoiled by ethidium bromide indicates that chromatin exists as closed looped domains. These early observations also suggested the presence of a [nuclear matrix](#) to which certain regions of the chromatin are anchored. Nucleoids were further used to demonstrate specific complexes within the nuclear matrix that are responsible for [DNA replication](#) and **RNA transcription**.

Experimental observations using nucleoids have received criticism in that artifactual associations between DNA or RNA with nuclear structures may arise due to the high salt conditions used in lysis. Jackson (7) provides arguments that such artifacts may be minimal: (1) the chromatin loops retain their size and the positions of genes on the loops remain constant; (2) preparation of nucleoids does not severely affect the apparent specificity of nascent RNA association with the nuclear matrix (8); (3) DNA, RNA, or [ribonucleoprotein](#) particles added before or after nucleoid purification in the presence of 2 *M* NaCl fail to associate with the nuclear cage (8); (4) no gross redistribution of DNA and RNA occurs during nucleoid preparation. Regardless of criticisms, nucleoids provided initial sources for studying nuclear matrices that many now accept as the structural support for nascent [DNA replication](#), RNA [transcription](#), pre-mRNA **splicing**, and RNA transport (9) (see [Nuclear Matrix](#)).

## Bibliography

1. O. G. Stonington and D. E. Pettijohn (1971) *Proc. Natl. Acad. Sci. USA* **68**, 6–9.
2. A. Worcel and E. Burgi (1972) *J. Mol. Biol.* **71**, 127–147.
3. P. R. Cook and I. A. Brazell (1975) *J. Cell Sci.* **19**, 261–279.
4. J. R. Paulson and U. K. Laemmli (1977) *Cell* **12**, 817–828.
5. B. Vogelstein, D. M. Pardoll, and D. S. Coffey (1980) *Cell* **22**, 79–85.
6. C. Benyajati and A. Worcel (1976) *Cell* **9**, 393–407.
7. D. A. Jackson (1986) In *Nuclear Structures Isolation and Characterization* (A. J. MacGillivray and G. D. Birnie, eds.), Butterworths, London, pp. 14–33.
8. D. A. Jackson, S. J. McCready, and P. R. Cook (1981) *Nature* **292**, 552–555.
9. C. R. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds. (1996) *Nuclear Structure and Gene Expression*, Academic Press, New York.

### Suggestion for Further Reading

10. R. Berezney and K. W. Jeon (1995) "Structural and functional organization of the nuclear matrix", In International Review of Cytology, A Survey of Cell Biology. Vols. **162A** and **162B**, Academic Press, San Diego. An excellent compilation of reviews on the nuclear matrix.

### Nucleolar Organizer

Nucleolar organizers are secondary constrictions within the mitotic [chromosomes](#) of **eukaryotes**. The primary constriction occurs at the [centromere](#). The nucleolar organizer genetic locus contains tandemly repeated **genes** that encode the 42 S to 47 S precursor of ribosomal RNA (pre-rRNA) (see [Nucleolus](#), [Ribosomes](#)). As the name implies, the nucleolar organizer gives rise to the nucleolus in the **telophase** of **mitosis**. During the late telophase, **RNA polymerase I** begins to transcribe rRNA genes once again, and the newly synthesized rRNA likely nucleates the assembly of the nucleolus by recruiting prenucleolar bodies to the [transcription](#) site.

Transcription of the 45 S rRNA genes is one of the primary events for nucleolar assembly. This was beautifully demonstrated by Karpen *et al.* (1), who inserted rRNA genes (ie, nucleolar organizers) within ectopic sites in *Drosophila* chromosomes by **P-element**-mediated integration. As a result, nucleoli were observed at ectopic sites of the polytene chromosomes. This work indicated that transcription of the ribosomal genes is necessary to initiate nucleolar assembly on the chromosome, but perhaps not sufficient. At least some of the nucleolar constituents used in the preceding interphase are retained in the cytoplasm during mitosis. In the late telophase, these constituents preassemble into extrachromosomal particles before or concomitant with the initiation of ribosomal gene transcription. These preassembled nucleolar particles, or prenucleolar bodies as they are called, contain at least the small nucleolar RNA, snoRNA U3, and the nucleolar proteins fibrillarin, nucleolin (C23), and NO38 (or B23) (2), all of which play a role in pre-rRNA processing. Although formation of the prenucleolar bodies is independent of transcription, fusion of the prenucleolar bodies to the organizer regions depends on rRNA transcription (3). Once the particles fuse with the nucleolar organizers, the resulting nucleoli grow in the early interphase because of increased rates of rRNA synthesis and the influx of other components necessary for ribosomal production (eg, ribosomal proteins). See [Nucleolus](#) for a more complete description of nucleolar structure and function.

Historically, Heitz (4) concluded that the number of secondary constrictions within certain chromosomes of *Drosophila funebris* is proportional to the number of nucleoli. He called these regions "sine acid thymonucleinico" (SAT, without DNA). It is now known that the mitotic nucleolar organizer regions retain a significant proportion of nucleolar-specific proteins relative to rDNA, as the chromosomes condense in prometaphase. These proteins include the **RNA polymerase I** complex, DNA **topoisomerase I**, and the [transcription factors](#) UBF and the SL1 complex (5). Together, these various proteins form a set that makes it possible to stain selectively for nucleolar organizers.

McClintock (6) actually coined the term "nucleolar organizer" to describe the satellite segment at the tip of **maize** chromosome 6 (the "nucleolar chromosome") that organizes the assembly of the nucleolus beginning in the telophase. If the terminal satellite element of chromosome 6 is broken into two equal portions by X-rays, her studies demonstrated that each portion gives rise to a nucleolus of normal size. McClintock's observation suggests that not all ribosomal genes within the normal nucleolus may be active. Subsequent work with *Drosophila* (7, 8) and *Xenopus* (9)



established that the tandemly repeated ribosomal genes reside at these nucleolar organizer regions. Nucleolar organizers occur on one or as many as six chromosomes, depending on the organism (10). Nucleolar gene clusters have been found on **autosomes** in some organisms and on sex chromosomes in others (eg, the X- and Y-chromosomes of *Drosophila*). Interestingly, ribosomal gene clusters are often found near **telomeres** on the shorter arm of the chromosomes that normally carry the nucleolar organizers. One possibility for this close association between the telomere and the nucleolar organizer is to reduce the adverse phenotypic affects of **unequal crossing over** between repeated rRNA genes on nonhomologous chromosomes (11, 12). Nucleolar organizers are found in humans on chromosomes 13, 14, 15, 21, and 22. All are **acrocentric** chromosomes.

### Bibliography

1. G. H. Karpen, J. E. Schaefer, and C. D. Laird (1988) *Genes Develop.* **2**, 1745–1763.
2. L. F. Jiménez-García, M. de L. Segura-Valdez, R. L. Ochs, L. I. Rothblum, R. Hannan, and D. L. Spector (1994) *Mol. Biol. Cell* **5**, 955–966.
3. P. Bell, M.-C. Dabauvalle, and U. Scheer (1992) *J. Cell Biol.* **118**, 1297–1304.
4. E. Heitz (1933) *Z. Zellforsch. Mikr. Anat.* **19**, 720–742.
5. P. Roussel, C. André, L. Comai, and D. Hernandez-Verdun (1996) *J. Cell Biol.* **133**, 235–246.
6. B. McClintock (1934) *Z. Zellforsch. Mikroskop. Anat.* **21**, 294–328.
7. F. Ritossa and S. Spiegelman (1965) *Proc. Natl. Acad. Sci. USA* **53**, 737–745.
8. F. Ritossa, K. Atwood, D. Lindsley, and S. Spiegelman (1966) *Nat. Cancer Inst. Monogr.* **23**, 449–472.
9. M. L. Birnstiel, H. Wallace, J. L. Sirlin, and M. Fischberg (1966) *Nat. Cancer. Inst. Monogr.* **23**, 431–448.
10. A. A. Hadjiolov (1985) *The Nucleolus and Ribosome Biogenesis*. Springer-Verlag. Wien.
11. A. Lima-de-Faria (1976) *Hereditas* **83**, 1–22.
12. A. Lima-de-Faria (1980) *Hereditas* **93**, 1–46.

### Suggestion for Further Reading

13. H. Busch and K. Smetana (1970) *The Nucleolus*, Academic Press, New York.

## Nucleolus

The primary function of the nucleolus is to biosynthesize **ribosomes**, both the large and small subunits. The nucleolus is the most conspicuous **organelle** within the **nucleus**. Because of its prominence, light microscopists have studied the nucleolus for more than 200 years. Fontana described the nucleolus as early as 1774, and there was a review of its literature in 1898 by Montgomery. Although much has been learned, the nucleolus is still imperfectly understood.

### 1. Overview

Assembly of large and small ribosomal subunits requires 18S, 5.8S, 28S, and 5S ribosomal **RNA** molecules (rRNA) and approximately 85 structural ribosomal **proteins**. The 18S, 5.8S, and 28S rRNA, but not 5S, are transcribed within the nucleolus by **RNA polymerase I** as parts of a large precursor transcript, pre-rRNA (47S in mammalian cells, 42S in *Xenopus*, and 35S in **yeast**) (see **Ribosomes**). RNA spacer sequences reside in front of the 18S region, between the 18S and 5.8S

regions, and between the 5.8S and 28S regions. The pre-rRNA is processed within the nucleolus by several endonuclease-catalyzed cleavages to yield mature 18S, 5.8S, and 28S rRNA. Although still imperfectly understood, these cleavage reactions require several nucleolar-specific proteins and small nucleolar RNAs (see later). Processing of the 45 S pre-rRNA occurs concomitantly with ribosomal subunit assembly. The 5 S rRNA is transcribed by RNA polymerase III from multiple copies of a **gene** that resides outside the nucleolus. After biosynthesis, the 5S RNA must migrate through the [nucleoplasm](#) to the nucleolus. Most of the ribosomal proteins translocate into the nucleus from their site of synthesis in the cytoplasm. Then they associate with the maturing rRNA in the nucleolus during assembly of the ribosomal subunits. The small ribosomal subunit consists of 18S rRNA and small subunit-specific ribosomal proteins (S-proteins), and the large ribosomal subunit consists of 5.8S and 28S rRNA, the incoming 5S rRNA, and large subunit-specific ribosomal proteins (L proteins) (1). Small and large subunits move from the nucleoli to the nuclear envelope through a [nuclear pore complex](#) and into the cytoplasm. The subunits unite to form intact ribosomes once they associate with [messenger RNA](#) to participate in synthesizing new protein (see [Translation](#)).

The number of nucleoli per nucleus varies, depending on the cell type. Yeast cells, for example, contain one nucleolus that is relatively large compared with the volume of the nucleus, and most metazoan cells contain one or a few nucleoli. On the other hand, amphibian oocytes contain hundreds of nucleoli per nucleus because of selective rDNA amplification during the early pachytene stage of meiosis I (2). The number of nucleoli also varies during the cell cycle. For example, nucleoli disassemble at the onset of **mitosis**, and they reappear beginning in the **telophase** by a process called nucleologenesis. In the late telophase/early **interphase**, a new cell displays several small nucleoli. The number reflects the number of [nucleolar organizers](#) that are characteristic of the species. As the new cells mature through the interphase, their nucleoli fuse to form larger nucleoli, and the number of nucleoli actually declines.

### 1.1. Nucleolar Ultrastructure

The interphase nucleolus consist of three ultrastructural subdivisions: the fibrillar center (FC), the surrounding dense fibrillar regions (DFR), and the peripheral granular regions (GR) (1, 3). The DFR are darker than the FC when stained with uranyl acetate and viewed by [electron microscopy](#). As the name implies, the GR consist of small granules believed to be the still immature ribosomal subunits (see later). The overall nucleolar morphology varies in different cell types because of differences in the relative sizes and distributions of these three subdivisions (3). For example, nucleoli within metabolically active cells contain DFR that are cast into nucleolonemae, which appear in the electron microscope as darkly-stained, anastomosing thread-like structures. Each thread is surrounded by lightly staining material. Conversely, nucleoli in a few cell types constitute a second morphological class and appear compact, with uniform RNA-containing structures. The third classification includes ring-shaped nucleoli in which RNA-containing structures form a peripheral ring around **chromatin-like fibrils** (3). Nucleolar morphology is influenced by the particular physiology of a given cell type. For example, the multiple nucleoli within transcriptionally active **oocytes** from *Xenopus* (stage IV) show pronounced nucleolonemae, but greatly reduced GR, perhaps because of rapid ribosomal subunit biosynthesis and export from the nucleus. Conversely, nucleoli in mature and relatively quiescent stage VI oocytes display collapsed DFR (eg, no anastomosing nucleolonemae) and very prominent GR (4).

The FCs clearly contain the DNA encoding the rRNA, and they may well be the interphase homologue of mitotic nucleolar organizers, the genetic loci that contain tandemly repeated ribosomal genes (see later). The precise function of the FC is controversial. RNA polymerase I (5, 6) and **topoisomerase I** (7) have been detected in the FC, thus suggesting a potential for transcription. Using selective *in situ* hybridization and immunolabeling techniques to detect both rDNA and rRNA within the FC, Thiry (8, 9) argued that transcription must occur within these centers. Hozák *et al.* (10), however, proposed an alternative model in which active transcription units radiate from condensed DNA within the FC. In this latter model, the transcription units occupy the borders between the FC and the DFR. Regardless of the precise site of transcription, the dense fibrillar

regions are the sites for early pre-rRNA processing and initial ribosomal subunit assembly. This assembly is a vectorial process, and maturation of the large and small ribosomal subunits continues on into the GR. Certain ribosomal proteins associate early in the DFR, and others associate in the GR. Finally, ribosomal subunits pass out of the nucleus into the cytoplasm, where final assembly with a few more ribosomal proteins takes place before they are competent to function in mRNA translation (1).

## 1.2. Ribosomal RNA Gene Organization and Expression

Genes that encode the large 45S pre-rRNA are tandemly repeated head-to-tail within genomes of most, if not all, **eukaryotic** organisms. Intergenic spacer (IGS) DNA sequences separate the individual genes. There are approximately 200 copies of ribosomal genes per haploid genome in humans, 100 in mouse, 500 to 600 in *Xenopus laevis*, 150 in *Drosophila melanogaster*, and 140 in the budding yeast *Saccharomyces cerevisiae* (1). Heitz (11, 12) and McClintock (13) demonstrated that the secondary constrictions in mitotic **chromosomes** are the genetic loci that initiate nucleolar formation as cells leave mitosis and enter the interphase. McClintock actually coined the term “nucleolar organizer” to describe these secondary constrictions. Now we know that the tandemly repeated pre-rRNA genes reside at these loci and that transcription of these genes initiates formation of nucleoli during the late telophase (see [Nucleolar Organizer](#)). Mutations that map to the nucleolar organizers reduce the number of rRNA genes or eliminate nucleoli altogether. For example, the bobbed (*bb*) mutation on the **X-chromosome** in *Drosophila* is a partial deletion of the tandem rDNA repeats. Flies homozygous or hemizygous for the mutation display pleiotropic effects of slow development, reduced fertility, and shortened (bobbed) bristles. In *Xenopus*, the anucleate mutation  $O_{nu}$  fails to synthesize rRNA because it lacks ribosomal RNA genes. Finally, McClintock (13) described a chromosomal deletion in **maize** that eliminates the nucleolar organizer and also the ability to form nucleoli.

One of the most informative techniques for elucidating nucleolar gene ultrastructure is the spreading of nucleolar chromatin in solutions with low salt and high pH for electron microscopic examination (14-16). These so-called “Miller spreads” show that nucleolar chromatin consists of nascent rRNP fibrils (rRNA associated with protein) which project from the tandemly repeated ribosomal genes. RNA polymerase I complexes pack the active ribosomal RNA genes, and the RNP fibrils remain attached to these complexes during preparation. As a result, the RNP fibrils display a gradient in length. The shortest fibrils are closest to the transcription initiation site, and the longest fibrils approach transcription termination. The gradient of RNP fibrils looks like a “Christmas tree” where the RNA polymerase I-coated gene is the tree trunk that supports the RNP branches. Interestingly, the 5'-end of each RNP fibril contains a relatively large RNP particle (see later for function). These particles have been called the “ornaments” on the Christmas trees.

The organization of a typical rRNA gene in higher eukaryotes is as follows beginning at the transcription start site: the 5'-external transcribed spacer (5'-ETS) region; the 18S rRNA region; the first internal transcribed spacer (ITS1); the 5.8S rRNA region; the second internal transcribed region (ITS2); the 28S rRNA region; and finally a short 3'-external transcribed spacer (3'-ETS). The *Xenopus* rDNA gene is described here because it is well characterized. The nascent transcript encoded by the *Xenopus* genes is 7,625 nucleotides long and has a **sedimentation coefficient** of 40S. Individual genes are separated by intergenic spacers (IGS) of variable length. For the vast majority of rRNA genes in *Xenopus*, **transcription** initiates at the **promoter** just upstream of the ETS and stops within the IGS at a site 235 bp downstream of the 28S region at the site called T2, yielding a 40S pre-rRNA transcript. The T1 site is at the very end of the 28S region, and it arises from initial posttranscriptional trimming of the 3'-ETS back to this site (17). The intergenic spacers in *Xenopus* contain several duplicated promoters and **enhancers** (18). Occasionally, transcription initiates at a distant upstream IGS promoter, reads through the other downstream promoters and enhancers (thus silencing these promoters), and then terminates at T3, a site approximately 60 bp upstream of the next rDNA proximal gene promoter (eg, 215 bp upstream of the transcription start site of the next gene). This secondary transcription unit is fully contained within the IGS. IGS transcription may

positively enhance transcription of this next gene by terminating at T3 and then handing the polymerase complex over to the adjacent downstream promoter (19, 20). Although the actual site of transcription termination has been debated, electron microscopic examination verifies that the vast majority of rDNA transcription units terminate transcription at the T2 site, at least in *Xenopus* (21).

Despite length differences in regions that encode either the mature rRNA themselves or the transcribed spacers, the overall genetic organization is similar in other eukaryotes (1). Many species of *Tetrahymena*, however, contain an **intron** within the 28S region. This is a group I intron originally characterized as self-splicing (22; see **Catalytic RNA**). Many of the *Drosophila* rRNA genes also contain one of two types of introns in the 28 S region. One type (a **transposable element**) appears only in 50% of the X-chromosome ribosomal genes. The other type appears in the ribosomal genes of both the X- and Y-chromosomes but at a frequency of 15%. Yeast has the same overall ribosomal gene structure as the higher eukaryotic cells, except that the primary transcript is shorter than in vertebrates (35S versus 40 to 45S). Yeast produces a 25S mature ribosomal RNA instead of a 28S transcript. Finally, the intergenic sequences in *S. cerevisiae* and *discoideum* contain 5S genes that are transcribed by RNA polymerase III but in the direction opposite to that of the main rDNA repeats.

Besides RNA polymerase I, several **transcription factors** (TF) are required for efficient rDNA expression (19). The **upstream binding factor** (UBF, a protein of 97 kDa in humans) contains a triple **α-helical domain** characteristic of the **HMG proteins** in **chromatin**. On its own, UBF binds promiscuously to both rDNA enhancers and promoters. UBF forms a stable complex with another factor, called SL1 in humans and rats, which then selectively binds the rDNA promoter. Transcription is initiated when an active RNA polymerase I complex recognizes the UBF-SL1 complex situated on the proximal promoter. Another transcription initiation factor, called TIF-IC in mouse cells, regulates the ability of polymerase I to recognize the UBF/SL1 complex. Topoisomerase I is relatively abundant within nucleoli. Its primary function is thought to aid RNA polymerase I transcription by relieving torsional stress generated within the DNA by the process of transcription. Low concentrations of **actinomycin D** (0.05 to 25 μg/ml) selectively inhibit RNA polymerase I transcription and strip the nascent RNP fibrils from the transcription units. This has been used to study nucleolar transcription, structure, and function (23).

Multiple copy, tandemly repeated genes encode the 5 S ribosomal RNA that assemble within the large ribosomal subunit. As mentioned above, the 5 S genes in *Saccharomyces*, *Dictyostelium*, and a few other organisms are closely linked to the rRNA genes within the same repeating units. The 5S genes in most other organisms are located on other chromosomes without a nucleolar organizer. In addition, the clusters of 5S genes are scattered between several loci in many organisms (1). The number of 5S genes within an organism usually exceeds the number of 45S rRNA genes. *Xenopus*, for example, has over 20,000 5S genes per haploid genome. The 5S genes are transcribed by RNA polymerase III and are regulated by transcription factors IIIA, IIIB, and IIIC. TFIIIA binds an intergenic control region of the 5S genes to direct the binding of TFIIIC. The 5S genes are small, 120 bp in total length. Therefore, the bound TFIIIA/C complex directs TFIIIB to bind just upstream of the transcription start site. RNA polymerase III finally associates with TFIIIB to initiate transcription. Interestingly, TFIIIB consists of three proteins, one of which is the **TATA binding protein** found in RNA polymerase II transcription initiation complexes. The 5S transcripts follow a rather circuitous route to the nucleoli. Monomeric ribosomal protein L5 binds 5S rRNA in the nucleus, and together they migrate out to the cytoplasm. They reenter the nucleus and then associate with the nucleolus, where they join in ribosomal assembly (24).

### 1.3. Nonribosomal Nucleolar Proteins

Besides RNA polymerase I, UBF, SL1, and topoisomerase I, which regulate 45 S rRNA gene transcription, many other nonribosomal nucleolar proteins are involved with pre-rRNA processing (25).

#### 1.3.1. Fibrillarin

To date vertebrate fibrillarin (Nop1 in *S. cerevisiae*) is the best characterized in terms of function. Vertebrate fibrillarin has a molecular mass of 34 kDa. It localizes mostly to the DFR, either to discrete foci within the DFR or to rings that surround the FC. This sublocalization suggests functional compartmentalization of fibrillarin within the DFR. Fibrillarin is conserved throughout eukaryotes. The amino-terminal portion contains a glycine/arginine-rich domain (called a GAR or RGG domain). Nearly half of these [arginine](#) residues are dimethylated, a **posttranslational modification** reserved for many nuclear RNA-binding proteins that contain RGG domains (26). Fibrillarin also contains a consensus RNA-binding domain (CS-RBD) in its center (27). Fibrillarin closely associates with many small nucleolar RNA (snoRNA) involved in either cleavage or base modification of the pre-rRNA. The best characterized fibrillarin association is with snoRNA U3 to form the RNP complex necessary for the first cleavage event within the ETS region of the pre-rRNA (see later).

### 1.3.2. Nucleolin

Vertebrate nucleolin is a phosphoprotein of molecular mass 90–110 kDa that is abundant in rapidly dividing somatic cells (28, 29) and in developing amphibian oocytes (30), where rates of ribosomal production are maximal. Although nucleolin shuttles between nucleoli and the cytoplasm (31), the greatest **steady state** concentration of the protein remains within the nucleolar DFR. Lesser amounts are found in the GR, and very little if any nucleolin resides in the FC. Vertebrate nucleolin consists of modular **domains** (32). The amino-terminal third contains alternating acidic domains (containing [glutamic acid](#), [aspartic acid](#), and phosphorylated [serine](#) residues exclusively) and basic domains (containing [lysine](#) residues). The serines within the acidic regions are phosphorylated by casein kinase type II enzymes during the interphase (see **Phosphorylation**). The basic domains, on the other hand, contain mitotic p34<sup>cdc2</sup>/cyclin B (MPF, [maturation promoting factor](#)) phosphorylation sites whose sequence is -Thr-Pro-(Ala or Gly)-Lys-Lys- (ie, -TP<sup>A</sup>/<sub>G</sub>KK-). Differential phosphorylation within the amino-terminal domain of nucleolin may regulate nucleolar assembly and function during interphase and also nucleolar disassembly during mitosis. The carboxyl-terminal two-thirds contains four consensus RNA-binding domains (CS-RBD), each about 80 amino acid residues long. The fourth CS-RBD is followed by a GAR (or RGG) domain of about 100 residues very near the carboxyl terminus. As in several other nuclear RNA-binding proteins (see previous discussion), the arginine residues within this GAR domain are dimethylated. Despite intense investigations of nucleolin, its precise nucleolar associations and functions remain uncertain. Nucleolin has been called a ribosome assembly factor (29). By sucrose gradient centrifugation, Herrera and Olson (33) showed that nucleolin associates with nascent rRNA, and it is believed that nucleolin binds nascent pre-rRNA *in vivo*. Recent experiments (34) showed that nucleolin decorates the RNP fibrils of Miller-spread CHO nucleolar chromatin, and it was shown by SELEX amplification that nucleolin preferentially binds a stem-loop structure consisting of a 4- to 5-bp stem and an eight-nucleotide loop containing a specific UCCCGA sequence motif. Interestingly, the six-nucleotide motif recovered by SELEX is also present in pre-rRNA, presumably at sites of nucleolin interaction (34). Binding to the motif is thought to be by one or more of nucleolin's CS-RBDs. Nucleolin's GAR domain may also function in RNA binding. It binds single-stranded DNA and RNA *in vitro*, and it displays RNA **helicase** (unwinding) activity.

Homologues of vertebrate nucleolin have been found in yeast. NSR1 in *S. cerevisiae* is an abbreviated version of the vertebrate protein. It also contains alternating acidic and basic domains at its amino terminus, but only two CS-RBD, followed by a GAR domain very near the carboxyl terminus. Gar2+ is the nucleolin homologue in *S. pombe*. Although the function of vertebrate nucleolin remains obscure, genetic **knockout** of these yeast homologues (35-37) disrupts processing of 18 S rRNA, causes a relative decrease in the abundance of 40 S small ribosomal subunits, and slows cell growth. These latter findings indicate that NSR1 and gar2+ have at least an associative role in pre-rRNA processing.

### 1.3.3. Nopp140, SRP40

Rat Nopp140 (SRP40 in yeast) is another nucleolar-specific phosphoprotein that contains alternating

acidic and basic domains (38). Its acidic domains are extremely rich in serine residues that are also phosphorylated by casein kinase II. Its basic regions are similar to those of nucleolin, and it has substantial numbers of [alanine](#), [valine](#), lysine, and [proline](#) residues. Several potential MPF sites that have the amino acid sequence -<sup>T</sup>/<sub>S</sub>PKK- are apparent in these basic domains. Similar to nucleolin, the highest steady-state concentrations of Nopp140 reside within the DFR, and to a lesser extent within the GR (39). One of Nopp140's distinguishing characteristics, however, is its ability to shuttle from nucleoli to the cytoplasm along nuclear tracks that begin within the DFR (39). Nopp140 was first identified as a nuclear localization signal (NLS)-binding protein (40). The supposition is that Nopp140 acts as a nucleolar **chaperone** by shuttling NLS-containing proteins into the nucleolus. Nopp140 specifically binds wild-type NLS sequences (not mutated versions) but only when the acidic regions are fully phosphorylated. The dephosphorylated version of Nopp140 does not bind wild-type NLS.

#### 1.3.4. NO38, B23, Nucleophosmin, and Numatrin

Similar to nucleolin, vertebrate NO38 ( $M_r$  of 38,000; also called B23, nucleophosmin, and numatrin) is a putative ribosomal assembly factor (29) and a nucleolar shuttling protein (31). Early reports showed that NO38 is enriched in the GR, but now it appears that NO38 localizes to both the DFR and GR (41). The function of NO38 also remains poorly understood. NO38 preferably binds single-stranded versus double-stranded nucleic acids. It binds RNA just as readily as it binds DNA. It binds nucleic acids cooperatively, and it destabilizes RNA helices (42). Interestingly, two **isoforms** of NO38 exist in rat (called B23.1 and B23.2), the isoforms result from [alternative splicing](#) of mRNA transcribed from one gene copy, and the B23.1 protein contains an additional 5-kDa tail at its carboxyl terminus that is not present in B23.2 (43). B23.1 is the predominant form that binds nucleic acids, and it resides within nucleoli. Conversely, B23.2 binds nucleic acids only weakly, and it resides within both the cytoplasm and the nucleoplasm but not within nucleoli. The functional significance of the two isoforms remains unknown but intriguing. Also like nucleolin, NO38 contains -Thr-Pro-Ala-Lys- (-TPAK-) sequences that are phosphorylated by MPF (44), thus suggesting that mitosis-specific phosphorylation of NO38 may help disassemble nucleoli (see later).

#### 1.3.5. Other proteins

Several other nucleolar proteins have been identified, mostly in yeast (45). The yeast Ssb1 protein contains an RNA-binding CS-RBD, and a GAR domain, followed by an acidic region. It copurifies with snR10 and snR11 small nucleolar RNA, but its role in pre-rRNA processing remains unknown. Gar1 in *S. cerevisiae* is closely associated with 18 S rRNA processing and, as its name implies, also contains a glycine/arginine-rich domain. Nop2p in *S. cerevisiae* (46) is very similar to vertebrate p120, which is closely associated with cell proliferation (47). Nop3p is another GAR-containing nucleolar protein in *S. cerevisiae*. When it is genetically depleted, yeast cells show abnormal processing of pre-rRNA, when the 27SB intermediate (precursor of the 25S mature rRNA) and the 23S intermediate (precursor of the 18S mature rRNA) both accumulate. This results in a drop in producing both small and large ribosomal subunits and in slow growth (48). Nop4p in *S. cerevisiae* is required for formation of the 60 S large ribosomal subunit (49). Nop77p is a yeast nucleolar protein that mediates interactions between Nop1 (fibrillarin) and pre-rRNA (50). Sof1p is yet another protein associated with Nop1 and U3 in 18 S rRNA processing. Several nucleolar RNA-dependent helicases in yeast (Spb4p, Drs1p, Dbp3p, and CA9) are necessary for 40S and 60S ribosomal subunit assembly (49). Pop1 and Snm1 associate with the 7-2/MRP RNA in 5.8S processing (see later). The interplay in processing between all of these proteins with the pre-rRNA and with the small nucleolar RNA (see later) remains imperfectly understood.

### 1.4. Small Nucleolar RNA and Pre-rRNA Processing

The nascent pre-rRNA transcript (47S in mammals) is processed by a series of exo- and endonucleolytic cleavage reactions to yield intermediates of discrete size before the formation of mature 18 S, 5.8 S, and 28 S ribosomal RNA (51). The transcript is modified posttranscriptionally by [methylation](#) of bases at conserved sites and by the conversion of specific **uridine** residues to pseudouridine. The assembly with in-coming ribosomal proteins to form the ribosomal subunits is

concurrent with cleavage and base modification. Several small nucleolar RNAs (snoRNA) have been identified that play either direct or associative roles in the processing reactions (45, 52). SnoRNAs in turn interact with nonribosomal nucleolar proteins (many described above) to form RNP complexes that engage in cleavage or base modification.

Initial cleavage of the mammalian 47S pre-rRNA at the 3'-T1 site yields the abundant 45S intermediate. The best characterized cleavage event, however, occurs within the 5'-ETS region of the primary transcript (53-55). U3 (56) and fibrillarin associate closely within the RNP complex known as "terminal balls" (54) that is positioned on the 5'-ETS (57). This terminal complex makes the first site-specific endonucleolytic cleavage in the processing of the 18S RNA (58). Miller spreads of nucleolar transcription units clearly show a relatively large RNP complex situated at the 5'-end of even the shortest fibrils. This RNP complex assembles on the 5'-ETS sequence specifically in the vicinity of cleavage sites at nucleotide 412 in human, 649 in mouse, or 782 in rat (57). Besides its association with the 5'-ETS cleavage, U3 is also associated with ITS1 and ITS2 cleavage reactions. Processing of 18S is inhibited by depletion of U3 (59), U14 and snR10 in yeast (60, 61), or of U14 and U22 in *Xenopus*. The 7-2/MRP snoRNA in yeast, with its associated Pop1 protein, plays a role in cleavage within ITS1 to liberate the yeast 5.8S rRNA. The U8 snoRNA in *Xenopus* mediates cleavage events within ITS2 and the 3'-ETS for proper excision of the 5.8S and 28S regions (62, 63).

SnoRNAs display conserved sequence elements that determine their conserved three-dimensional structures (45). For instance, the highly characterized U3 snoRNA contains conserved motifs designated boxes A, B, C, C', and D. Boxes A, B, and C are unique to U3. Boxes A and B interact with rRNA, and box C associates with fibrillarin. Boxes C' and D of U3 are common to many other snoRNAs that constitute a class called C/D snoRNA (box C' of U3 is simply called C in these other snoRNAs). This box C (C' in U3) is a conserved motif of six nucleotides that is six nucleotides downstream of the 5'-end of the snoRNA. Box D, on the other hand, is a conserved motif of four nucleotides that is just upstream of the 3'-end of the snoRNA. The 5'- and 3'-ends of the snoRNA form a short terminal **double helix**. Complementarity does not extend into boxes C and D, so they remain single-stranded and positioned next to the terminal helix. In the case of U14, boxes C and D are necessary for proper intron processing, stability, and accumulation (64). Interestingly, the majority of C/D snoRNAs are thought to associate with fibrillarin, again via box C.

More than fifty snoRNAs have been identified (65). Whereas snoRNA U3, U8, and U13 are encoded by their own independent genes, many of these other snoRNAs are encoded as introns of several different pre-mRNA molecules. For example, mouse U14 was the first snoRNA identified as part of an intron within the gene encoding the 70-kDa cognate **heat shock** protein hsc70 (66). Interestingly, many snoRNAs are transcribed as introns of pre-mRNA that encode proteins necessary for ribosomal assembly or factors necessary for protein translation (45). For example, snoRNA U17 in *Xenopus* is repeated in each of six introns of the pre-mRNA that encodes ribosomal protein S8 (67), and snoRNA U24 in humans and chickens is contained within an intron of the ribosomal protein L7 pre-mRNA. Both U20 and U23 in humans, mice, rats, hamsters, frogs, and fish are contained within introns of the nucleolin pre-mRNA. Expression of intronic snoRNA within these various pre-mRNA suggests a potential mechanism to coordinate synthesis of pre-rRNA processing components, ribosome structural proteins, and even translation factors. The intron processing mechanism that removes the snoRNA from the intron remains obscure. In the case of U14, the snoRNA region of the intron forms a loop whose ends are held together by a short double-helical stem of about seven base pairs. This helix actually spans what will be the final 5'- and 3'-ends of the snoRNA. An endonucleolytic cut within the stem is necessary to liberate the snoRNA from the rest of the intron. Characteristic of the intronic snoRNA, the cut is staggered, leaving a 3'-overhang of one base.

Although snoRNA U3, U8, U14, and U22 associate with fibrillarin in various pre-rRNA cleavage events, the processing roles for many of these other C/D snoRNAs have only recently been described (65, 68), and these results provide exciting discoveries about posttranscriptional methylation of the pre-rRNA. Just upstream of the D box in a particular snoRNA is a single-stranded region that is

complementary to a specific pre-rRNA segment, either within the 18S, 5.8S, or 28S regions, but not to any of the spacer regions. The D box, together with an upstream helix formed between the snoRNA and the complementary segment of the pre-rRNA, directs the methylation specifically of the fifth nucleotide upstream of the D box. Methylation occurs on the **ribose** ring as 2'-O-methylation. Approximately 105 conserved methylation sites are known within the rRNA. Therefore, the majority of the C/D snoRNAs are guide RNAs, perhaps in association with fibrillarin, for site-specific methylation of the mature ribosomal RNA. Methylation creates **hydrophobic** sites within the rRNA that may direct proper RNA folding, processing, and subsequent interaction with incoming ribosomal proteins as ribosomal subunits assemble.

### 1.5. Nucleolar Cycle

The cycle of nucleolar assembly and disassembly is intimately linked to the [cell cycle](#). The nucleolus disassembles during the prometaphase of mitosis and begins to reassemble during the telophase. The onset of mitosis is controlled by tight regulation of MPF, or cyclin B/p34<sup>cdc2</sup> kinase. Once activated, MPF phosphorylates many different cellular substrates. At least two nucleolar proteins, nucleolin and NO38, are direct substrates for MPF (44). The phosphorylation sites in hamster nucleolin are the [threonine](#) residues within nine -TP<sup>A</sup>/<sub>G</sub>KK- motifs found within the basic domains of the amino terminus. At least one other nucleolar protein of known sequence, Nopp140, displays potential MPF phosphorylation sites. The current supposition is that MPF-specific phosphorylation of key structural proteins within the nucleolus contributes to nucleolar disassembly.

Another prominent factor in nucleolar disassembly may be the shutdown of rDNA transcription. Interphase nucleoli partially disassemble when actinomycin D blocks rDNA transcription (23). Transcription of the pre-rRNA genes normally begins to decline in the **prophase**. This decline may result from tight compaction of rDNA into the nucleolar organizer regions of the mitotic chromosomes. **Immunofluorescence** studies of mitotic cells show that at least fibrillarin and nucleolin, two proteins involved with pre-rRNA processing and packaging, diffuse into the cytoplasm (the nuclear envelope has disassembled by this time) upon nucleolar disassembly. SnoRNA U8 also diffuses to the cytoplasm (69). The lack of pre-rRNA transcripts may contribute to the diffusion of RNA processing components. Interestingly, RNA polymerase I and the transcription factors UBF and SL1 all remain bound to the nucleolar organizer regions of the chromosomes (70).

Nucleolar reassembly (nucleogenesis) begins during the telophase with the appearance of prenucleolar bodies. These particles contain fibrillarin, NO38, nucleolin, Nopp140, and the snoRNA U3, but no rRNA. All are factors closely associated with the processing of pre-rRNA. As rDNA transcription by polymerase I resumes in the telophase, these prenucleolar bodies coalesce with the nucleolar organizers to reconstitute intact nucleoli by the interphase. By the early interphase, there may be several functional nucleoli, which may fuse to form fewer but larger nucleoli. Interestingly, the prenucleolar bodies share morphological and compositional similarities with coiled bodies (71), originally described in 1903 as accessory bodies of nucleoli in nerve cells (72). Coiled bodies contain fibrillarin, ribosomal protein S6, DNA topoisomerase I, and snoRNA U3, but no rRNA. In addition, coiled bodies contain the specific p80 coilin protein and snRNPs normally involved with pre-mRNA splicing. The precise relationship between prenucleolar bodies and coiled bodies remains imperfectly understood.

### Bibliography

1. A. A. Hadjiolov (1985) *The Nucleolus and Ribosome Biogenesis*, Springer-Verlag, New York.
2. H. C. MacGregor (1982) In *The Nucleolus* (E. G. Jordan and C. A. Cullins, eds.), Cambridge University Press, Cambridge, pp. 129–153.
3. H. Busch and K. Smetana (1970) *The Nucleolus*, Academic Press, New York.
4. S. B. Shah, C. D. Terry, D. A. Wells, and P. J. DiMario (1996) *Chromosoma* **105**, 111–121.
5. U. Scheer and K. M. Rose (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1431–1435.



6. G. Reimer, I. Raska, E. M. Tan, and U. Scheer (1987) *Virch. Arch. B* **54**, 131–143.
7. K. M. Rose, J. Szopa, F-S. Han, Y-C. Cheng, A. Richter, and U. Scheer (1988) *Chromosoma* **96**, 411–416.
8. M. Thiry (1992) *Nucleic Acids Res.* **20**, 6195–6200.
9. M. Thiry (1993) *J. Cell Sci.* **105**, 33–39.
10. P. Hozák, P. R. Cook, C. Schöfer, W. Mosgöller, and F. Wachtler (1994) *J. Cell Sci.* **107**, 639–648.
11. E. Heitz (1931) *Vicia. Planta* **15**, 495–505.
12. E. Heitz (1933) *Z. Zellforsch. Mikr. Anat.* **19**, 720–742.
13. B. McClintock (1934) *Z. Zellforsch. Mikr. Anat.* **21**, 294–328.
14. O. L. Miller Jr. and B. R. Beatty (1969) *Science* **164**, 955–957.
15. O. L. Miller Jr. and A. H. Bakken (1972) In *Gene Transcription in Reproductive Tissue* (E. Diczfaly, ed.), Karolinska Institute, Stockholm, pp. 155–177.
16. O. L. Miller Jr. (1981) *J. Cell Biol.* **91**, 15s–27s.
17. P. Labhart and R. H. Reeder (1986) *Cell* **45**, 431–443.
18. R. F. J. De Winter and T. Moss (1986) *Cell* **44**, 313–318.
19. R. H. Reeder (1990) *Trends Genet.* **6**, 390–394.
20. B. Sollner-Webb and E. B. Mougey (1991) *Trends Biochem. Sci.* **16**, 58–62.
21. B. Meissner, A. Hofmann, H. Steinbeiser, H. Spring, O. L. Miller Jr., and M. F. Trendelenburg (1991) *Chromosoma* **101**, 222–230.
22. T. R. Cech (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Plainview, NY, pp. 239–269.
23. U. Scheer, M. F. Trendelenburg, and W. W. Franke (1975) *J. Cell Biol.* **65**, 163–179.
24. M. Z. Barciszewska, V. A. Erdmann, and J. Barciszewski (1996) *Biol. Rev.* **71**, 1–25.
25. P. J. Shaw and E. G. Jordan (1995) *Ann. Rev. Cell Dev. Biol.* **11**, 93–121.
26. Q. Liu and G. Dreyfuss (1995) *Mol. Cell. Biol.* **15**, 2800–2808.
27. D. E. Draper (1995) *Ann. Rev. Biochem.* **64**, 593–620.
28. M. A. Lischwe, R. L. Richards, R. K. Busch, and H. Busch (1981) *Exp. Cell Res.* **136**, 101–109.
29. M. O. J. Olson (1990) In *The Eukaryotic Nucleus: Molecular Biochemistry and Macromolecular Assemblies*, Vol. **2** (P. R. Straus and S. H. Wilson, eds.), Telford Press, Caldwell, New Jersey, pp. 519–559.
30. M. Caizergues-Ferrer, P. Mariottini, C. Curie, B. Lapeyre, N. Gas, F. Amalric, and F. Amaldi (1989) *Genes Dev.* **3**, 324–333.
31. R. A. Borer, C. F. Lehner, H. M. Eppenberger, and E. A. Nigg (1989) *Cell* **56**, 379–390.
32. B. Lapeyre, H. Bourbon, and F. Amalric (1987) *Proc. Natl. Acad. Sci. USA* **84**, 1472–1476.
33. A. H. Herrera and M. O. J. Olson (1986) *Biochemistry* **25**, 6258–6264.
34. L. Ghisolfi-Nieto, G. Joseph, F. Puvio-Dutilleul, F. Amalric, and P. Bouvet (1996) *J. Mol. Biol.* **260**, 34–53.
35. K. Kondo and M. Inouye (1992) *J. Biol. Chem.* **267**, 16252–16258.
36. K. Kondo, L. R. Z. Kowalski, and M. Inouye (1992) *J. Biol. Chem.* **267**, 16259–16265.
37. M-P. Gulli, J-P. Girard, D. Zabetakis, B. Lapeyre, T. Melese, and M. Caizergues-Ferrer (1995) *Nucleic Acids Res.* **23**, 1912–1918.
38. U. T. Meier (1996) *J. Biol. Chem.* **271**, 19376–19384.
39. U. T. Meier and G. Blobel (1992) *Cell* **70**, 127–138.
40. U. T. Meier and G. Blobel (1990) *J. Cell Biol.* **111**, 2235–2245.
41. M. Biggiogera, S. H. Kaufmann, J. H. Shaper, N. Gas, F. Amalric, and S. Fakan (1991)

Chromosoma **100**, 162–172.

42. T. S. Dumbar, G. A. Gentry, and M. O. J. Olson (1989) *Biochemistry* **28**, 9495–9501.
43. J-H. Chang and M. O. J. Olson (1989) *J. Biol. Chem.* **264**, 11732–11737.
44. M. Peter, J. Nakagawa, M. Dorée, J. C. Labbé, and E. A. Nigg (1990) *Cell* **60**, 791–801.
45. E. S. Maxwell and M. J. Fournier (1995) *Ann. Rev. Biochem.* **35**, 897–934.
46. E. de Beus, J. S. Brockenbrough, B. Hong, and J. P. Aris (1994) *J. Cell Biol.* **127**, 1799–1813.
47. B. C. Valdez, L. Perlaky, Y. Saijo, D. Henning, C. Zhu, R. K. Busch, W. W. Zhang, and H. Busch (1992) *Cancer Res.* **52**, 5681–5686.
48. I. Russell and D. Tollervey (1992) *J. Cell Biol.* **119**, 737–747.
49. C. Sun and J. L. Woolford (1994) *EMBO J.* **13**, 3127–3135.
50. T. Bergès, E. Petfalski, D. Tollervey, and E. C. Hurt (1994) *EMBO J.* **13**, 3136–3148.
51. S. A. Gerbi, R. Savino, B. Stebbins-Boaz, C. Jeppesen, and R. Rivera-Leon (1990) In *The Ribosome: Structure, Function and Evolution* (W. Hill, P. Moore, D. Schlessinger, A. Dahlberg, J. Warner, and R. Garrett, eds.), American Society for Microbiology, Washington, DC, pp. 452–469.
52. B. Sollner-Webb, K. T. Tycowski, and J. A. Steitz (1996) In *Ribosomal RNA: Structure, Evolution, Processing, and Function in Protein Synthesis* (R. A. Zimmermann and A. E. Dahlberg, eds.), CRC Press, Boca Raton, pp. 469–490.
53. S. Kass, N. Craig, and B. Sollner-Webb (1987) *Mol. Cell. Biology* **7**, 2891–2898.
54. S. Kass and B. Sollner-Webb (1990) *Mol. Cell. Biol.* **10**, 4920–4931.
55. S. Kass, K. Tyc, J. A. Steitz, and B. Sollner-Webb (1990) *Cell* **60**, 897–908.
56. K. A. Parker and J. A. Steitz (1987) *Molec. Cell. Biol.* **7**, 2899–2913.
57. E. B. Mougey, L. K. Pape, and B. Sollner-Webb (1993b) *Mol. Cell. Biol.* **13**, 5990–5998.
58. E. B. Mougey, M. O'Reilly, Y. Osheim, O. L. Miller Jr., A. Beyer, and B. Sollner-Webb (1993a) *Genes Dev.* **7**, 1609–1619.
59. J. M. X. Hughes and M. Ares Jr. (1991) *EMBO J.* **10**, 4231–4239.
60. H. V. Li, J. Zagorski, and M. J. Fournier (1990) *Mol. Cell. Biol.* **10**, 1145–1152.
61. W. Q. Liang and M. J. Fournier (1995) *Genes Dev.* **9**, 2433–2443.
62. B. Peculis and J. A. Steitz (1993) *Cell* **73**, 1233–1245.
63. B. Peculis and J. A. Steitz (1994) *Genes Dev.* **8**, 2241–2255.
64. G. M. Huang, A. Jarmoloski, J. C. R. Struck, and M. J. Fournier (1992) *Mol. Cell. Biol.* **12**, 4456–4463.
65. Z. Kiss-László, Y. Henry, J-P. Bachellerie, M. Caizergues-Ferrer, and T. Kiss (1996) *Cell* **85**, 1077–1088.
66. J. Liu and E. S. Maxwell (1990) *Nucleic Acids Res.* **18**, 6565–6571.
67. F. Cecconi, P. Mariottini, F. Loreni, P. Pierandrei-Amaldi, N. Campioni, and F. Amaldi (1994) *Nucleic Acids Res.* **22**, 732–741.
68. M. Nicoloso, L-H. Qu, B. Michot, and J-P. Bachellerie (1996) *J. Mol. Biol.* **260**, 178–195.
69. A. G. Matera, K. T. Tycowski, J. A. Steitz, and D. C. Ward (1994) *Mol. Cell Biol.* **5**, 1289–1299.
70. P. Roussel, C. André, L. Comai, and D. Hernandez-Verdun (1996) *J. Cell Biol.* **133**, 235–246.
71. L. F. Jiménez-García, M. de L. Segura-Valdez, R. L. Ochs, L. I. Rothblum, R. Hannan, and D. L. Spector (1994) *Mol. Cell Biol.* **5**, 955–966.
72. S. R. Ramón y Cajal (1903) *Trab. Lab. Invest. Biol.* **2**, 129–221.

### Suggestions for Further Reading

73. E. G. Jordan and C. A. Cullis (1982) *The Nucleolus*, Cambridge University Press, Cambridge.

74. J. R. Warner (1989) The nucleolus and ribosome formation *Curr. Opin Cell Biol.* **2**, 521–527.
75. Plus references [1](#), [3](#), [16](#), [19](#), [20](#), [25](#), [29](#), [45](#), and [52](#).

## Nucleoplasm

The terms nucleoplasm and [karyoplasm](#) have been used interchangeably for more than one hundred years. According to E. B. Wilson ([1](#)), karyoplasm (nucleoplasm) is “ $\frac{1}{4}$  the nuclear substance in contradistinction to the cytoplasmic.” Today we know that the nucleoplasm includes the [nucleolus](#), [euchromatin](#) and [heterochromatin](#), a myriad of [ribonucleoprotein](#) particles such as **spliceosomes**, interchromatin granules and fibers, coiled bodies, and the [nuclear matrix](#), which has often been referred to as the “ground substance.” These nuclear contents are separated from the cytoplasm by the [nuclear envelope](#). For detailed descriptions of these components, see [Nucleus](#).

### Bibliography

1. E. B. Wilson (1925) *The Cell in Development and Heredity*, Macmillan, New York.

## Nucleoplasmin

Nucleoplasmin is an acidic, thermostable protein found in the **nuclei** of amphibian **eggs** that binds [histones](#) H2A and H2B and transfers them to DNA to form [nucleosome](#) cores *in vitro* ([1-7](#)), suggesting that nucleoplasmin is required for nucleosome assembly inside eggs. However, nucleoplasmin is not a component of the nucleosome, nor does it possess steric information required for nucleosome assembly; this information resides in the histones. These unusual properties of nucleoplasmin led to the term **molecular chaperone** ([6](#)), a usage subsequently extended to a much wider range of proteins involved transiently in assembly processes ([8](#)). Subsequent work has revealed a related function of nucleoplasmin in decondensing sperm [chromatin](#) on fertilization of the egg, resulting in the replacement of the [protamine](#) proteins of [sperm](#) chromatin by the histone proteins of the [zygote](#) ([5](#), [9-11](#)), in the reactivation of genes in somatic nuclei incubated in egg extracts ([12](#)), and in the binding of [transcription factors](#) to nucleosome cores ([13](#)).

Cell-free extracts of [Xenopus laevis](#) eggs convert added purified DNA to a regularly repeating chain of nucleosomes, using either the endogenous pool of histones or exogenous sources of purified histones. In contrast, separated DNA and histones form insoluble aggregates when incubated together directly at physiological ionic strength (0.1–0.2 M NaCl). It is the nucleoplasmin present in the cell-free extracts that prevents the formation of these aggregates, by binding transiently to histones via electrostatic interactions before the histones are transferred to DNA ([6](#)). In support of this view is the observation that complexes of histones H2A and H2B with nucleoplasmin occur in soluble extracts of [Xenopus](#) eggs and **oocyte** nuclei, and that the histones of these complexes can be transferred to DNA *in vitro* ([13](#), [14](#)). Homologues of nucleoplasmin have been detected in somatic cells ([15](#), [16](#)) but the specificity of some of the immunological methods used have been questioned

(5), and their possible role in nucleosome assembly in such cells is not established. Nucleosome assembly in adult somatic cells is linked to DNA replication, unlike the case for *Xenopus* eggs, and involves the action of other protein assembly factors (see [Nucleosome](#)). Nucleoplasmin appears to be specialized to function as a store of maternally synthesized histones for chromatin assembly during early embryogenesis.

The detailed mechanism of nucleoplasmin action is unknown, but it does not appear to require the hydrolysis of ATP, unlike the action of some other molecular chaperones such as the **chaperonins**. An important aspect is the reduction of the high positive charge density on the histone surface by the binding of the negatively charged nucleoplasmin (17). In effect, nucleoplasmin acts as an “electrostatic filter” that reduces the high positive charge density of histones and thus inhibits the tendency for unspecific aggregates to form between histones and DNA. This transient inhibitory role of nucleoplasmin allows the [self-assembly](#) properties of the histones with DNA to predominate over the incorrect interactions generated by their high densities of opposite charge. The same property is amplified to regulate the disassembly of DNA complexes with other basic proteins; massive hyperphosphorylation of nucleoplasmin during oocyte maturation allows the decondensation of the tight protamine-DNA complexes in sperm chromatin (5, 10).

### Bibliography

1. C. Dingwall and R. A. Laskey (1990) in R. J. Ellis, ed, *Molecular Chaperones*, Saunders Scientific Publications, Philadelphia, pp. 11–17.
2. R. A. Laskey and W. C. Earnshaw (1980) *Nature* **286**, 763–767.
3. S. M. Dilworth, S. J. Black, and R. A. Laskey (1987) *Cell* **51**, 1009–1018.
4. S. M. Dilworth and C. Dingwall (1988) *BioEssays* **9**, 44–49.
5. R. A. Laskey, A. D. Mills, A. Philpott, G. H. Leno, S. M. Dilworth, and C. Dingwall (1993) *Phil. Trans. Roy. Soc. Lond. B* **339**, 263–269.
6. R. A. Laskey, B. M. Honda, A. D. Mills, and J. T. Finch (1978) *Nature* **275**, 416–420.
7. J. A. Kleinschmidt, A. Seiter, and H. Zentgraf (1990) *EMBO J.* **9**, 1309–1318.
8. R. J. Ellis (1987) *Nature* **328**, 378–379.
9. A. Philpott and G. H. Leno (1992) *Cell* **69**, 759–767.
10. G. H. Leno, A. D. Mills, A. Philpott, and R. A. Laskey (1996) *J. Biol. Chem.* **271**, 7253–7256.
11. S. A. Ruiz-Lara, L. Cornudella, and A. Rodriguez-Campos (1996) *Eur. J. Biochem.* **240**, 186–194.
12. S. Dimitrov and A. P. Wolffe (1996) *EMBO J.* **15**, 5897–5906.
13. H. Chen, B. Lai, and J. L. Workman (1994) *EMBO J.* **13**, 380–390.
14. J. A. Kleinschmidt, E. Fortkamp, G. Krohne, H. Zentgraf, and W. W. Franke (1985) *J. Biol. Chem.* **260**, 1166–1176.
15. M. S. Schmidt-Zachmann, B. Hugle-Dorr, and W. Franke (1987) *EMBO J.* **6**, 1881–1890.
16. M. Cotten and R. Chalkley (1985) *EMBO J.* **6**, 3945–3954.
17. W. C. Earnshaw, B. M. Honda, and R. A. Laskey (1980) *Cell* **21**, 373–383.

### Nucleoprotein

Nucleoprotein is a complex of **DNA** and protein. In the eukaryotic [nucleus](#), DNA is organized by the [histone](#) proteins into the nucleoprotein complex known as [chromatin](#). The recognition that

nucleoprotein structures represent the functional form of genes has taken many decades to develop. Toward the end of the nineteenth century, numerous investigators formulated the theory that [chromosomes](#) determine inherited characteristics. These studies were almost entirely based on cytological observations with the light microscope. Although chromosomes are clearly only in the nucleus, the influence of components of the cytoplasm on inherited characteristics was examined by forcing embryonic nuclei into regions of the cytoplasm where they would not normally be found. These experiments and others led Morgan to propose the theory that [differentiation](#) depends on variation in the activity of genes in different cell types. The genes were clearly in the chromosomes, but their biochemical composition remained completely unknown (1).

The last quarter of the nineteenth century also saw the recognition of **RNA** (first identified as yeast **nucleic acid**) and of DNA (thymus nucleic acid) and the discovery of histones. Albrecht Kossel isolated nuclei from the erythrocytes of geese and examined the basic proteins in his preparations, which he named histones. The apparent biochemical simplicity of DNA and the obvious complexity of protein in chromosomes led investigators mistakenly to regard the latter component as the major constituent of the elusive genes. Only the gradual acceptance of experiments on the capacity of DNA alone to change the genetic characteristics of the cell led to the recognition of nucleic acid as the key structural component of a gene (2).

The elucidation of the **double-helical** structure of DNA, with its immediate implications for self-duplication, opened up the new approaches of molecular biology to clarifying the nature of genes. Although it was recognized that the double helix contains the requisite information to specify a genetic function, how this information was controlled was not understood. The apparent heterogeneity of histones due to **proteolysis** and the various **posttranslational modifications** of these proteins suggested that they might be important in regulating genes. Eventually methodological improvements for isolating and resolving the different histones demonstrated that the histones are highly conserved in eukaryotes and that only a few basic types exist. This lack of variety implied that histones themselves were unlikely to be the determinants of gene specific transcription. However, a key role for histone modification remained central to prevailing ideas of transcriptional regulation.

A major breakthrough came in the 1970s when a combination of methodologies, including **nuclease** digestion (see **DNase I sensitivity**), protein–protein **cross-linking**, [electron microscopy](#) and [sedimentation velocity centrifugation](#), determined that chromatin consists of a repetitive fundamental nucleoprotein complex, which came to be called the [nucleosome](#). The integrity of the nucleosome depends on highly specific histone–histone interactions and the recognition by the histones of DNA structural features induced as the nucleosome is assembled. The core histones are present as an octamer, consisting of two molecules of H2A, H2B, H3 and H4. Histones H3 and H4 assemble a tetramer, (H3, H4)<sub>2</sub>, that wraps DNA so that two dimers of H2A and H2B stably associate. Once two turns of DNA are wrapped around the octamer, a fifth linker histone, such as histone H1, is stably incorporated to complete the assembly process. Although all nucleosomes maintain these architectural features, there are many variations built on this common theme.

In a series of insightful experiments, Grunstein and colleagues determined that changes in nucleosomal packaging had pleiotropic effects on genetic activity. Subsequent work established that very specific modifications in histone structure could either activate or repress specific genes. This led directly to the resurgence of interest in understanding genetic activity in the natural chromosomal environment, which has characterized much of the research in eukaryotic transcriptional regulation over the past few years.

The new-found interest in the role of nucleoprotein complexes, such as chromatin, in transcriptional regulation has been fueled by progress in two specific areas. Structural studies led to the recognition that the histones are isomorphous with components of the **transcriptional** machinery. These observations provided an architectural foundation for examining the specific roles of histones and [transcription factors](#) in the assembly and function of regulatory nucleoprotein complexes. It was

shown that specific modifications to nucleosomal architecture through histone **acetylation** and removal of histones H2A/H2B or H1 alleviate the repressive effects of chromatin assembly . In certain instances, chromatin assembly also stimulates the transcription process. Thus the potential roles of nucleosomal proteins in gene control became more interesting. Biochemical purification of histone acetyltransferases and deacetylases provided an even closer link between chromatin and the transcriptional machinery. It was discovered that histone acetyltransferases are components of large macromolecular complexes, known as coactivators , which are targeted at specific promoters by transcriptional activators. Therefore a direct link was established between histone acetylation and transcriptional activation. Histone deacetylases were found within corepressor complexes that turn genes off. Once again, histone chemistry became an important variable to consider in transcriptional control.

Now it is recognized that to understand transcriptional control or any other regulated event in the nucleus, it is necessary to define the nucleoprotein structure within which the DNA is utilized. Aside from the characterization of specific architecture, we must also determine how that structure might change. Nucleoprotein complexes are not static, but dynamic. Targeted histone modifications within regulatory nucleoprotein complexes have emerged as a means of modulating the stability of repressive chromatin structures and the transcription process itself. The observations made using simple model systems are having an increasing impact on our understanding of both [development](#) and disease. It is now probable that our increasing knowledge of both nucleoprotein complex structure and function in the nucleus will provide many avenues for future advances in the biotechnological and medical fields.

#### Bibliography

1. T. H. Morgan (1934) *Embryology and Genetics*, Columbia University Press, New York.
2. O. T. Avery, C. M. MacLeod, and M. McCarty (1944) *J. Exp. Med.* **79**, 137–158.
3. M. Grunstein (1990) *Ann. Rev. Cell Biol.* **6**, 643–678.
4. K. L. Clark, E. D. Halay, E. Lai, and S. K. Burley (1993) *Nature* **364**, 412–420.
5. D. Y. Lee, J. J. Hayes, D. Pruss, and A. P. Wolffe (1993) *Cell* **72**, 73–84.
6. J. E. Brownell et al. (1996) *Cell* **84**, 843–851.

#### Suggestions for Further Reading

7. K. van Holde (1988) *Chromatin*. Springer Verlag, Berlin.
8. A. Wolffe (1998) *Chromatin: Structure and Function*. 3rd ed., Academic Press, London.

#### Nucleoside Diphosphate Kinase

For the first 35 years after its discovery in the early 1950s, nucleoside diphosphate kinase (NDP kinase) was considered the quintessential “**housekeeping** enzyme,” responsible only for the transfer of the  $\gamma$ -phosphate from a variety of nucleoside triphosphates to any common nucleoside diphosphate . Initially, NDP kinase was discovered as a partner in the succinyl-CoA synthetase reaction, using GTP synthesized by that enzyme to drive the synthesis of ATP from ADP. Later, it was recognized that the enzyme can participate in the biosynthesis of all common ribo- and deoxyribonucleoside triphosphates from the respective diphosphates and a phosphate donor (usually ATP) (see [Nucleotides, Nucleosides, And Nucleobases](#)). The past decade, however, has revealed a bewildering array of regulatory and developmental roles ascribed to NDP kinase, as well as additional catalytic capabilities.

## 1. NDP Kinase Reaction

NDP kinase functions via a phosphoenzyme intermediate. The reaction begins with transfer of the  $\gamma$ -phosphate from the nucleoside triphosphate substrate— $N_1$ TP in reaction (1) to a specific [histidine](#) residue of the enzyme. After dissociation of the resultant diphosphate ( $N_1$ DP), the second substrate ( $N_2$ DP, the phosphate acceptor) binds, and phosphate is transferred from the enzyme, yielding  $N_2$ TP, the nucleoside triphosphate product. Thus, the enzyme displays classical ping-pong kinetics (see [Kinetic Mechanisms, Enzyme](#)):



The enzyme is quite efficient, with  $k_{cat}$  values over  $1000 \text{ sec}^{-1}$ . The equilibrium constant for the reaction is close to 1. Because of this, and because ATP is the most abundant nucleoside triphosphate in most cells and organelles (see [Adenylate Charge](#)), the primary role of the enzyme has been considered to be the synthesis of the other seven triphosphates from the respective diphosphates and ATP.

Steady-state kinetic analysis revealed that the enzyme is rather nonspecific with regard to both phosphate acceptor and donor. However, recent pre-steady-state kinetic analysis of a human NDP kinase has revealed more specificity than was previously thought. This analysis depended on, first, a change in intrinsic tryptophan fluorescence accompanying phosphorylation of the enzyme and, second, the stability of the phosphoenzyme intermediate, which permitted its isolation for analysis of enzyme dephosphorylation by  $N_2$ DP. As shown in Table 1, the second-order rate constants for the formation of  $E \sim P$  varied 17-fold among the five triphosphates tested, whereas the dephosphorylation rate constants varied 145-fold among the eight phosphate acceptors tested. In general, **purines** are better substrates than **pyrimidines**, and ribonucleotides are better substrates than deoxyribonucleotides.

**Table 1. Rate Constants for Phosphorylation and Dephosphorylation of Human NDP Kinase**

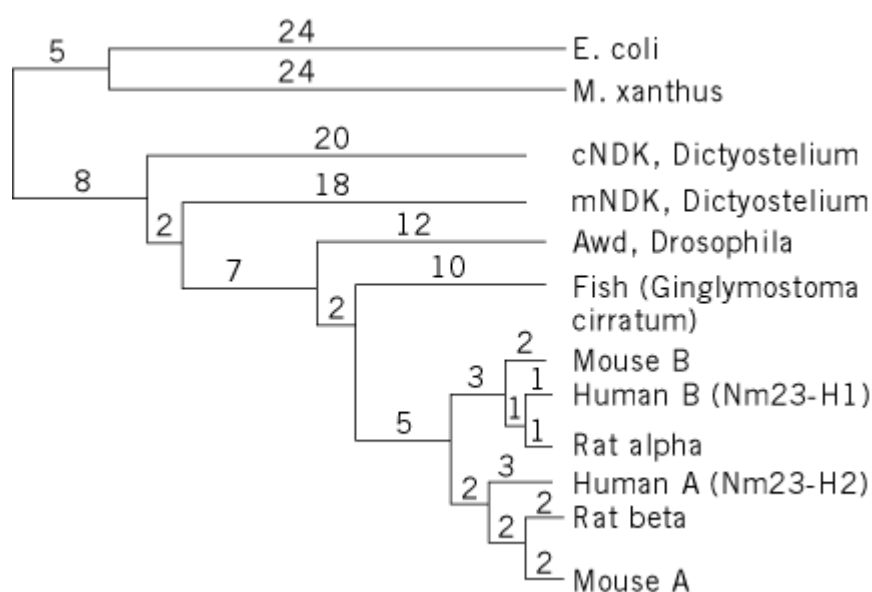
| Base     | Second-Order Rate Constant ( $\times 10^{-6} M^{-1} s^{-1}$ ) |      |      |                  |
|----------|---|------|------|------------------|
|          | rNTP  | dNTP | rNDP | dNDP             |
| Adenine  | 7.9   | —    | 17.9 | 4.5              |
| Guanine  | 12.2  | —    | 29.1 | 10.3             |
| Uracil   | 2.0   | —    | 3.8  | 1.3              |
| Cytosine | 0.7   | —    | 1.3  | 0.2              |
| Thymine  | —   | 2.3  | —    | 4.3 <sup>a</sup> |

<sup>a</sup> Data for rNTPs and dNTPs refer to enzyme phosphorylation rates; data for rNDPs and dNDPs refer to dephosphorylation of the phosphoenzyme. Source of the data: ref. 2.

## 2. Structure of NDP Kinase

The NDP kinase protein has been rather highly conserved through evolution, and there is about 43% amino acid sequence identity between the human enzyme and that of *Escherichia coli*. Figure 1 summarizes relationships among some of the known NDP kinase sequences. Although the [primary structure](#) is highly conserved, the [quaternary structure](#) is not; the enzymes of eukaryotic origin are homohexamers whereas the prokaryotic enzymes have a homotetrameric structure. Subunit molecular weights range from 15 to 18 kDa. Interestingly, preliminary data suggest that at least one mitochondrial **isoform** of a eukaryotic NDP kinase might not have a hexameric structure.

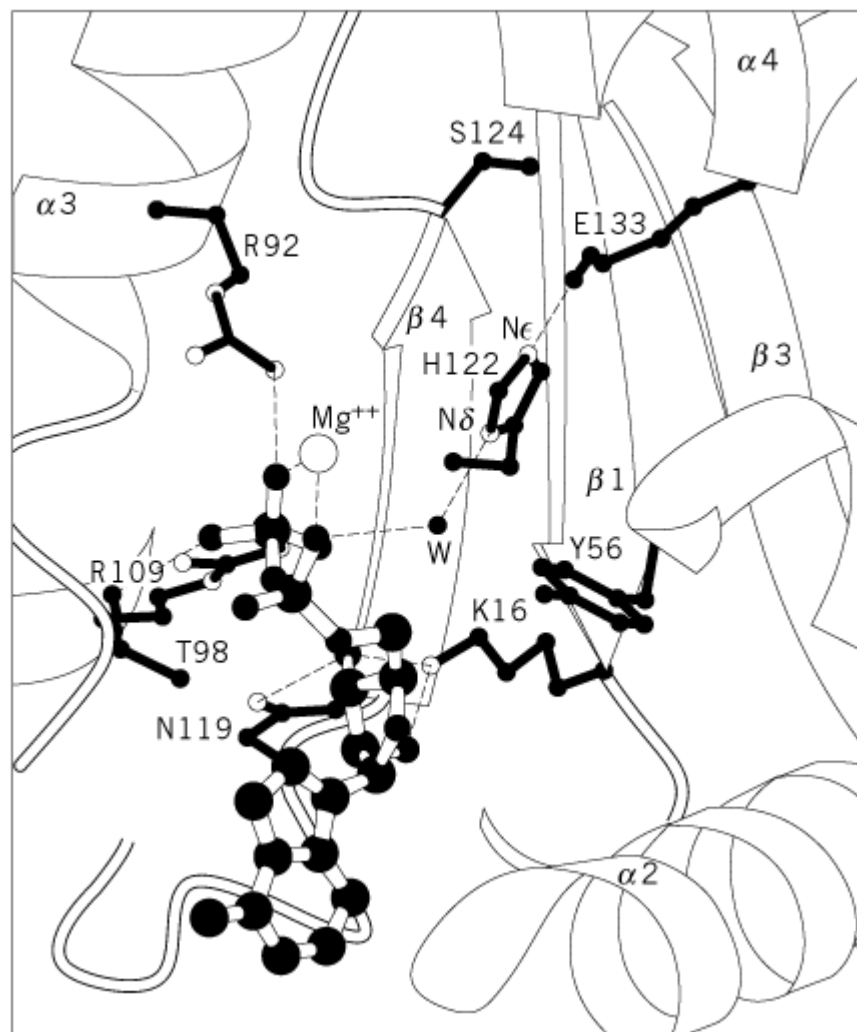
**Figure 1.** A [phylogenetic tree](#) illustrating the evolutionary relationships among NDP kinase amino acid sequences. The numbers above the horizontal lines represent evolutionary time in arbitrary units. The vertical lines indicate the positions of common ancestors. Reprinted with permission from H. Troll, T. Winckler, I. Lascu, N. Müller, W. Saurin, M. Véron, and R. Mutzel (1993) *J. Biol. Chem.* **268**, 25469–25475.



[X-ray crystallography](#) structures have been described for the NDP kinases from *Dictyostelium discoideum* (4, 5), *Drosophila melanogaster* (6), and human tissue (7). The tertiary structures are similar. The *Dictyostelium* structure has been determined complexed with ADP and with dTDP. Figure 2 shows the [active-site](#) structure of the ADP-enzyme complex. The nucleotide is bound near His122, the residue that becomes phosphorylated during the reaction. The b-phosphate interacts with two **arginine** residues and a [threonine](#) that are conserved in all known NDP kinases, and the ribose 2' and 3' hydroxyls similarly have polar interactions with Lys19 and Asn119, which are also conserved. Of considerable interest, the purine base forms no polar contacts with the enzyme, nor does thymine in the dTDP-enzyme complex. The absence of specific contacts involving the purine or pyrimidine base accounts for the relatively low nucleotide specificity of the enzyme. However, the structure explains a distinct aspect of specificity that the enzyme does display. Bourdais et al (8) reported that the diphosphates of the anti-HIV drugs azidothymidine, dideoxyadenosine, and dideoxythymidine are very poor substrates. This finding points up the importance of the interaction between the enzyme and the sugar 3' hydroxyl, which is missing in these analogues, a feature that may be significant in design of more effective antiviral nucleoside drugs.

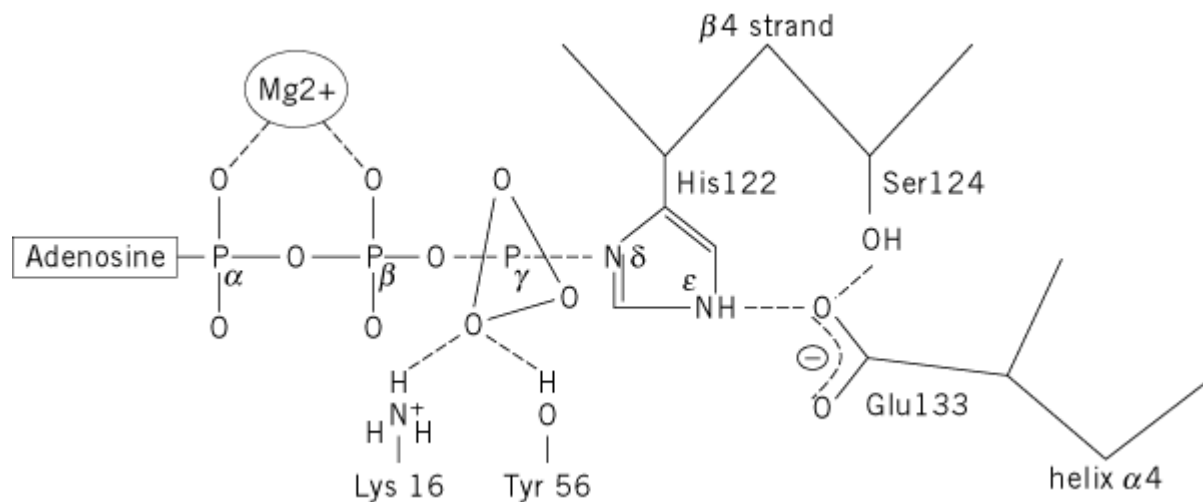
**Figure 2.** Structure of the active site of *Dictyostelium* NDP kinase with bound ADP. ADP is shown with open bonds, with the purine base at the bottom of the figure. Amino acid side chains are shown in black; W is a water molecule. Reprinted with permission from A. Tepper, H. Dammann, A. A. Bominaar, and M. Véron (1994) *J. Biol. Chem.* **269**,





The structure of the ADP-NDP kinase complex suggests a model for the transition state leading to the phosphoenzyme intermediate (4). Amino acid residues involved in both binding and catalysis, shown in Figure 3, are residues that are conserved in all NDP kinases analyzed to date.

**Figure 3.** Structure of the proposed NDP kinase [transition state](#), based on the structure of *Dictyostelium* NDP kinase complexed with ADP. Reprinted with permission from A. Tepper, H. Dammann, A. A. Bominaar, and M. Véron (1994) *J. Biol. Chem.* **269**, 32175–32180.



### 3. Additional Functions of NDP Kinase

Several recent developments have graphically demonstrated that the NDP kinase protein does more than simply synthesize nucleoside triphosphates from diphosphates. Some of these additional roles are hinted at by the existence of specific protein–protein interactions involving NDP kinase. Other functions have been traced by genetic analysis to the structural genes for NDP kinase. It is evident that NDP kinase lies at the heart of a number of developmental, genetic, and metabolic control systems.

#### 3.1. Human *nm23* Genes and Metastatic Cancer

Analysis of a closely related series of murine tumor cell lines revealed the existence of a gene, named *nm23* (nonmetastatic clone no. 23), whose action correlated with the potential of the cell lines to develop metastatic tumors after implantation in mice. Lines of low metastatic potential expressed the gene at high levels (9). Sequence analysis revealed the *nm23* gene to be the structural gene for an NDP kinase. Human cells contain two major Nm23 proteins, Nm23-H1 and Nm23-H2 (NDP kinase **isoforms A and B**, respectively). Action of the A isoform is specifically associated with metastasis. The NDP kinase enzymatic activity of the protein is not essential for this activity because some mutant forms of the protein that lack kinase activity still retain the metastatic suppression activity. Since the original cloning and analysis of the two human *nm23* genes, two additional members of the NDP kinase protein family, less closely related to each other and to Nm23-H1 and -H2, have been described in human tissues.

Important clues to the biochemical activity involved in tumor suppression are beginning to accumulate (10). NDP kinase has a protein [kinase](#) activity, being able to transfer phosphate from its catalytically important [histidine](#) residue to other proteins, including histidine residues on ATP citrate lyase and succinyl-CoA synthetase, and aspartate or glutamate residues on a number of human [membrane proteins](#). Using a cell motility assay with breast cancer cells, Wagner et al (10) explored the activities of a series of *nm23-H1* mutants. Mutations that abolished motility suppression activity also interfered with the phosphorylation of Asp and Glu residues in the membrane proteins, with minimal effects on the protein histidine kinase activity. Thus, the suppression of metastasis correlates closely with the transfer of phosphate to acidic amino acid residues. This activity is similar to that seen in two-component regulatory systems in bacteria; indeed, *E. coli* NDP kinase also can transfer phosphate from histidine to protein aspartate residues, and evidence strongly suggests regulatory roles for this activity (11).

#### 3.2. Wing Development and the *Prune* Gene in *Drosophila*

In *Drosophila*, NDP kinase is encoded by the *awd* gene (abnormal wing development). Mutations in this gene are normally not lethal, but one particular mutation, a Pro-to-Ser substitution at residue 97,

is lethal when in combination with a mutation in the *prune* gene. *Prune* controls eye color, by somehow regulating the activity of GTP cyclohydrolase, an enzyme involved in the synthesis of the pteridine eye pigments. However, the *prune* gene itself evidently encodes a small protein similar to the mammalian GTPase-activating proteins involved in **Ras**-mediated [signal transduction](#) pathways (12). *In vitro* experiments with one such protein (13) suggest that NDPK can use a protein-GDP complex as a substrate, with direct conversion of the bound GDP to GTP occurring in the absence of dissociation from the protein. Similar findings have been described for mammalian NDP kinases (14) but, to date, interactions of this type have not been unambiguously demonstrated *in vivo*.

The Pro-to-Ser mutation that makes *awd* lethal in combination with *prune* is called the *kpn* mutation (killer of *prune*). The proline residue involved lies in a loop that participates in subunit contacts. When the same mutation was engineered into the *Dictyostelium* NDP kinase, the mutant enzyme was found to undergo ready dissociation from its hexameric structure to folded monomers that retained some enzymatic activity. Lascu et al. (15) speculated that comparable dissociation occurs *in vivo* in *kpn* mutant cells and that the conditional lethality of the mutation derives from a deleterious effect of these monomers, still unknown.

Similar conclusions have been drawn from analysis of a natural mutation in human NDP kinase A (Nm23-H1), which is found in several aggressive neuroblastoma tumors (16). This mutation involves a Ser-to-Gly substitution at position 120. *In vitro*, the mutant enzyme was found to be quite unstable, and when it renatured, an intermediate accumulated that had the properties of a [molten globule](#). Whether the accumulation of such a folding intermediate at abnormal levels *in vivo* might somehow contribute toward the tumor phenotype is an exceedingly interesting, unanswered question.

### 3.3. Control of *c-myc* Gene Transcription

Postel et al (17) proposed that human NDP kinase B also serves as a [transcription factor](#) for the *c-myc* proto-oncogene. This role was identified by screening a cervical cell carcinoma [cDNA library](#) with DNA containing binding sites for PuF, the purine-binding transcription factor. Sequence analysis of a positive clone revealed the PuF gene to be identical to *nm23-H2*, the structural gene for NDP kinase B. The enzyme has been shown to activate *c-myc* transcription *in vitro* and to bind preferentially to pyrimidine-rich single-stranded DNA (18). This is a specific activity of the B isoform, whereas the metastatic tumor suppression described earlier is associated with the A isoform. Nevertheless, it is intriguing to consider whether there is a common element in the actions of these two proteins with respect to tumorigenesis and **tumor suppression**.

## 4. Protein-Protein Interactions Involving NDP Kinase

As befits a protein involved in a range of biological processes, NDP kinase has been shown to interact with numerous proteins. In *Pseudomonas aeruginosa*, NDP kinase copurifies with succinyl-CoA synthetase (19), which is consistent with the idea that one role of NDP kinase is to generate ATP from the substrate-level phosphorylation step of the citric acid cycle. The possible connection with substrate-level phosphorylation is seen also in frog heart preparations, where NDP kinase was found to copurify with five proteins (20), including glyceraldehyde 3-phosphate dehydrogenase, creatine kinase, pyruvate kinase, and vimentin.

Relatively little attention has been paid to the role of NDP kinase in DNA precursor biosynthesis. One system that has been explored is bacteriophage T4-infected *E. coli*, in which all the reactions in dNTP and DNA synthesis are catalyzed by phage-coded enzymes, except for the synthesis of dNTPs from the respective NTPs. In this system, the NDP kinase of the host bacterium, encoded by the *ndk* gene, forms part of a multienzyme complex for dNTP synthesis (21). [Affinity chromatography](#) with immobilized *E. coli* NDP kinase as the affinity ligand reveals specific interactions with several T4 phage proteins, including [dihydrofolate reductase](#), [ribonucleotide reductase](#), and [single-strand DNA-binding protein](#) (22). However, NDP kinase is dispensable, because an *E. coli* strain bearing a targeted deletion of the *ndk* gene is viable, and it supports phage T4 growth (23). Interestingly, these mutant cells have abnormal nucleoside triphosphate pools—in particular, a twentyfold elevation of

the dCTP level—and, probably as a result, an elevated spontaneous mutation rate. Lu and Inouye (24) searched in the *ndk* deletion strain for an enzyme capable of synthesizing nucleoside triphosphates from the respective diphosphates, and they made the surprising observation that adenylate kinase, the product of the *adk* gene, possesses this capacity. Because adenylate kinase was previously known only to catalyze the ATP-dependent phosphorylation of AMP or dAMP to the respective diphosphate, this is a novel activity of the *adk* protein, representing another important but incompletely understood aspect of the biochemistry of nucleoside diphosphate kinase.

## Bibliography

1. R. E. J. Parks and R. P. Agarwal (1973) *Enzymes* **8**, 307–334.
2. S. Schaertl, M. Konrad, and M. A. Geeves (1998) *J. Biol. Chem.* **273**, 5662–5669.
3. D. O. Lambeth, J. G. Mehus, M. A. Ivey, and B. I. Milavetz (1997) *J. Biol. Chem.* **272**, 24604–24611.
4. S. Moréra, I. Lascu, C. Dumas, G. LeBras, P. Briozzo, M. Véron, and J. Janin (1994) *Biochemistry* **33**, 459–467.
5. J. Cherfils, S. Moréra, I. Lascu, M. Véron, and J. Janin (1994) *Biochemistry* **33**, 9062–9069.
6. M. Chiadmi, S. Moréra, I. Lascu, C. Dumas, G. LeBras, M. Véron, and J. Janin (1993) *Structure* **1**, 283–293.
7. P. A. Webb, O. Perisic, C. E. Mendola, J. M. Backer, and R. L. Williams (1995) *J. Mol. Biol.* **251**, 574–587.
8. J. Bourdais, R. Biondi, S. Sarfati, C. Guerreiro, I. Lascu, J. Janin, and M. Véron (1996) *J. Biol. Chem.* **271**, 7887–7890.
9. A. DeLaRosa, R. L. Williams, and P. S. Steeg (1995) *BioEssays* **17**, 53–62.
10. P. D. Wagner, P. S. Steeg, and N-D. Vu (1997) *Proc. Natl. Acad. Sci. USA* **94**, 9000–9005.
11. Q. Lu, H. Park, L. A. Egger, and M. Inouye (1996) *J. Biol. Chem.* **271**, 32886–32893.
12. D. H. F. Tang, C. M. Engele, and T. R. Venkatesh (1991) *Nature* **353**, 437–440.
13. P. A. Randazzo, J. K. Northup, and R. A. Kahn (1991) *Science* **254**, 850–853.
14. S. Kikkawa, K. Takahashi, K-i. Takahashi, N. Shimada, M. Ui, N. Kimura, and T. Katada (1990) *J. Biol. Chem.* **265**, 21536–21540.
15. I. Lascu, D. Deville-Bonne, P. Glaser, and M. Véron (1993) *J. Biol. Chem.* **268**, 20268–20275.
16. I. Lascu, S. Schaertl, C. Wang, C. Sarger, A. Giartosio, G. Briand, M-L. Lacombe, and M. Konrad (1997) *J. Biol. Chem.* **272**, 15599–15602.
17. E. H. Postel, S. J. Berberich, S. J. Flint, and C. A. Ferrone (1993) *Science* **261**, 478–480.
18. M. Hildebrandt, M-L. Lacombe, S. Mesnildrey, and M. Véron (1995) *Nucl. Acids Res.* **23**, 3858–3864.
19. A. Kavanaugh-Black, D. M. Connolly, S. A. Chugani, and A. M. Chakrabary (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5883–5887.
20. A.d-S. Otero (1997) *J. Biol. Chem.* **272**, 14690–14694.
21. J. R. Allen, G. W. Lasser, D. A. Goldman, J. W. Booth, and C. K. Mathews (1983) *J. Biol. Chem.* **258**, 5746–5753.
22. L. J. Wheeler, N. B. Ray, C. Ungermann, S. P. Hendricks, M. A. Bernard, E. S. Hanson, and C. K. Mathews (1996) *J. Biol. Chem.* **271**, 11156–11162.
23. Q. Lu, X. Zhang, N. Almaula, C. K. Mathews, and M. Inouye (1995) *J. Mol. Biol.* **254**, 337–341.
24. Q. Lu and M. Inouye (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5720–5725.

## Suggestions for Further Reading

25. R. E. J. Parks and R. P. Agarwal (1973) Nucleoside diphosphate kinases. *Enzymes* **8**, 307–334. Describes the discovery of NDP kinase and all the early enzymological studies on the enzyme.

26. M. L. Lacombe and K. H. Jakobs (1992) Nucleoside diphosphate kinases as potential new targets for control of development and cancer. *Trends Pharmacol. Sci.* **13**, 46–48. An early review of the idea that NDP kinase might interact directly with protein-bound guanine nucleotides.
27. A. De LaRosa, R. L. Williams, and P. S. Steeg (1995) Nm23/nucleoside diphosphate kinase: Toward a structural and biochemical understanding of its biological functions. *BioEssays* **17**, 63–62. Reviews the relationship between NDP kinase action and tumor suppression.

## Nucleosome

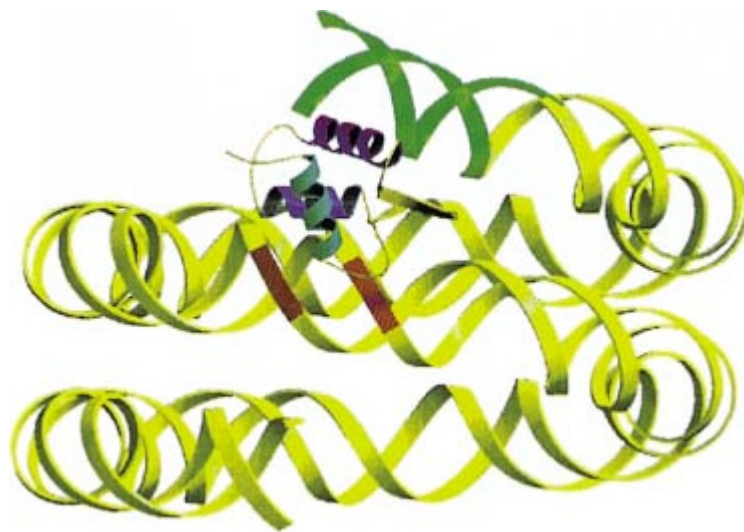
The nucleosome is the basic repeating subunit of [chromatin](#) (1), which has a regulatory role in gene [transcription](#) as well as serving to package DNA. It consists of about 166 bp of DNA wound in two left-handed superhelical turns around a wedge-shaped [histone](#) octamer (a tetramer, H<sub>3</sub><sub>2</sub>H<sub>4</sub><sub>2</sub>, of histones H3 and H4, flanked by two H2A.H2B dimers), a length of [linker DNA](#) which may vary from effectively zero (yeast, mammalian cortical neurons) to 74 bp (sea urchin sperm) and defining the repeat length of the chromatin (~170 to 40 bp), and one molecule of linker histone (H1 or one of its variants; see [Histones](#)). Nucleosomes are responsible for the familiar beaded appearance of chromatin seen in the **electron microscope**. The nucleosome is a wedge-shaped disk, about 110 Å in diameter and about 55 Å high, with a pseudo-dyad axis of symmetry close to which the entering and exiting DNA duplexes lie.

Nucleosomes (and runs of two, three, etc, nucleosomes) may be excised from chromatin by cleavage in the more accessible linker DNA using, for example, the double-stranded endonuclease micrococcal nuclease (**Staphylococcal nuclease**) and then fractionated by [sedimentation velocity centrifugation](#) through sucrose gradients. Further digestion of mononucleosomes with micrococcal nuclease, which has exonuclease as well as endonuclease activity, proceeds in two stages. First, the DNA is trimmed to ~166bp, at which point histone H1 presents a barrier to further digestion and a metastable particle termed the “chromatosome” is produced. Further digestion of the chromatosome trims the DNA further, causing loss of histone H1, presumably as its binding site close to the nucleosome core is invaded, and then trims the DNA back to 146 bp, where contacts with the octamer block further digestion. The limit product of digestion is the nucleosome core particle whose structure is now known at high resolution (see text below). Protection by histone H1 in the chromatosome was believed to be symmetric, with 10 bp protected at each end of the core particle. It now appears, however, that it may be asymmetric (2, 3), with 20 bp protected at one end of the core particle and none at the other (assuming that the addition of linker histone did not shift the dyad by one helical turn).

The exact location of H1 in the nucleosome has been controversial (summarized in Ref. 4), because there is not yet any direct structural information for a chromatosome. The general view from several lines of evidence was that the globular domain of H1 bound close to, or over, the dyad axis, on the outside of the DNA gyres, and that it contacted two duplexes through two DNA binding sites (see [Histones](#)) that were shown (5) to be required for correct binding to two duplexes and to chromatin. However, results obtained for a nucleosome reconstituted on to *Xenopus* 5S rRNA gene, which contains a strong nucleosome positioning sequence (see [Chromatin](#)), suggested a fundamentally different picture, with the globular domain bound away from the dyad axis and *inside* the DNA gyres, contacting only one duplex (6). It may be that this is a unique feature of the 5S nucleosome, although this would be a little surprising. Current evidence (summarized in Ref. 7) suggests that, for at least the majority of chromatosomes (a bulk population from chicken erythrocytes), the globular

domain of H5 makes two DNA contacts, one of which is indeed near the dyad. Helix III, the “recognition helix” in the globular domain (by analogy with the homologous structures of HNF3g and [cyclic AMP receptor protein](#); see [Histones](#)), contacts one of the exiting/entering duplexes, about one helical turn from the end of the chromatosomal DNA, and the second site makes a contact that is in the vicinity of the dyad (Fig. 1). This location is such that the basic C-terminal domain of H1 would be directed toward the linker; a linker location of this domain is suggested by electron microscopy of mononucleosomes (8). It has been argued (7) that the globular domain could well be binding to the 5S nucleosome in the same way, if assumptions about the position of the octamer in that system were to prove to be incorrect; however, the possibility remains that the 5S nucleosome is different.

**Figure 1.** The location of the globular domain of histone H5 in the 166-bp chromatosome. The recognition helix (magenta) makes a contact one double-helical turn from the end of the DNA (region shown in green) and the second binding site on the opposite face of the domain contacts the central gyre, close to the pseudodyad (3 bp around the dyad shown in red). (Based on the data in Ref. 11. Drawing courtesy of S. Muyldermans.) See color insert.



## 1. Nucleosome Core Particle

The length of DNA in nucleosome core particles produced by limit nuclease digestion as described above is relatively homogeneous, allowing particles with sequences representing essentially the entire [genome](#) to be crystallized (despite the mixed DNA sequences, histone variants, and [post-translational modifications](#)). The year 1985 saw the publication of an [X-ray crystallography](#) structure of the 206,000 Da nucleosome core particle to 7 Å resolution, which confirmed the general shape and arrangement of histones and DNA inferred from earlier studies using electron microscopy and **single-particle reconstruction** analysis. Determination of the structure of the histone octamer (free of DNA) to 3.1 Å resolution in 1991 (see [Histones](#)) allowed some of the details of protein–DNA interactions to be inferred, but the full picture has become available only in the 2.8 Å structure of the nucleosome core particle (9) published in 1997 (Fig. 2). To obtain this resolution, highly homogeneous nucleosome core particles were reconstituted (see [Chromatin](#)) from synthetic 146-bp palindromic DNA and homogeneous recombinant histones produced in *Escherichia coli*. This is also the largest piece of DNA whose structure has been determined.

**Figure 2.** Nucleosome core particle structure at 2.8 Å resolution. (a) *Left*: A view down the superhelical axis; *right*, view perpendicular to the superhelical axis. The pseudo-twofold axis is aligned vertically with the DNA center at the

top in each case. The DNA is shown in brown and turquoise; H3, H4, H2A, and H2B are shown in blue, green, yellow, and red, respectively. **(b)** The histone octamer structure in the nucleosome core, formed by removing the DNA from the nucleosome core particle image in **(a)**. **(c)** Histone tails between DNA gyres shown in a space-filling model. The colors are as in **(a)**. Note that the H2B (red) and H3 (blue) *N*-terminal tails pass through channels in the DNA superhelix (white) formed by aligned minor grooves. [**(a)** and **(c)** reproduced from ref. [9](#), with permission; **(b)** reproduced from Ref. [12](#), with permission. See color insert.





The DNA in the core particle is wound in 1.65 superhelical turns on a left-handed helical ramp on the octamer surface; the remaining 10 bp at each end are essentially straight due to contacts between neighboring molecules in the crystal, to give pseudocontinuous double helices. The core particle has a pseudodyad axis of symmetry that coincides with the dyad axis of the octamer and passes between the two H3 molecules. One base pair in the 2.1 Å structure sits on the dyad, so the DNA of the two “halves” of the particles (73 and 72 bp) are not identical, despite the initial choice of a palindromic sequence in an attempt to obtain a nucleosome core particle with exact two-fold symmetry! The DNA superhelix is not uniformly bent, as was already known, due to the nature of the local histone–DNA interactions, the sharpest distortions being at positions 15 bp and 40 to 50 bp on either side of the dyad (positions  $\pm 1.5$  and  $\pm 4$  to 5 relative to the dyad as position 0). The overall helical periodicity of the DNA in the core particle is 10.2 bp per turn, but varies in detail along the DNA.

The organization of histones is very similar to that seen in the structure of the histone octamer on its own at 3.1 Å resolution (10), with [histone folds](#) arranged in handshake motifs (Fig. 2 a). The histone folds organize the central 121 bp of DNA; the remaining 13 bp at each end are organized by *N*-terminal **a-helical** extensions to the folds in the two H3 molecules and adjacent residues from the *N*-terminal tail domain. Each histone-fold dimer in the octamer organizes 27 to 28 bp, with 4-bp stretches between contacts. Each dimer has three independent DNA-binding sites: two contributed by the loops at each end of the dimer, and one formed by the two *N*-termini of  $\alpha$ -helices of each histone fold at the dimer apex. Histone–DNA contacts include extensive [salt bridges](#) and [hydrogen bonds](#) from both main chains and side chains of the histones to the phosphate backbone, together with [nonpolar](#) contacts with DNA sugar groups, and [electrostatic interactions](#) of the positively charged *N*-termini of  $\alpha$ -helices in the histone folds with DNA phosphate groups. Such interactions might be expected for proteins that can bind to a variety of DNA sequences; the number of specific base contacts in the core particle structure is very small (one of them being a nonpolar contact of the 5-methyl group of a thymidine in the major groove). Strikingly, **arginine** side chains contact the minor groove each of the 14 times it faces inwards.

A surprising finding was that the *N*-terminal tails (see [Histones](#)), which are the most variable regions between species and which are the sites of post-translational modification, notably acetylation (see [Histone Acetylation](#)), are directed outward from the core particle, exiting between the DNA gyres; they are therefore well-placed to make contacts with other nucleosomes or other proteins. Some of the tails are visible in Figure 2a; Figure 2b shows the DNA stripped away to show the extended histone tails and the histone folds in the octamer. The four H3 and H2B tails exit through the narrow channels formed by perfectly aligned minor grooves of adjacent gyres of the DNA superhelix (Fig. 2c), the alignment being a consequence both of the change in the twist of the DNA from 10.6 bp in solution to 10.2 bp in the core particle and local adjustments to the twist. H4 and H2A tails exit in minor grooves either on the “top” or the “bottom” of the core particle. The accessible histone tails in chromatin are known to be targets for transcription-linked histone acetyltransferases (see [Histone Acetylation](#)) as well as recognition sites for proteins that lead to a more repressive chromatin structure (eg, SIR proteins, giving **silenced** regions at yeast **telomeres** and mating type loci) and for at least one chromatin remodeling complex (NURF) (see [Chromatin](#)). Histones are thus not simply packaging proteins but also have a gene regulatory function. DNA could be weakly bound to the outside of the core particle (but with low occupancy; otherwise they would have been seen in the core particle structure), and acetylation could well prevent this. Such a weak effect alone is unlikely, however, to account for the enhancing effect of acetylation on transcription *in vivo*, although it could perhaps account for results of some *in vitro* assays in which acetylation can increase the access of [transcription factors](#) to mononucleosome DNA (see [Histone Acetylation](#)). The nature of the crystal contacts hints at possible roles for the tails. The basic H4 tail, which has been implicated in a variety of biological functions, contacts an acidic patch on the surface of a neighboring octamer. Residues 16 to 25 of H4 (the most internal/proximal acetylation site is at position 16) make multiple

interactions with acidic side chains on H2A and H2B in an adjacent nucleosome in the crystal. Model building shows that a similar interaction could occur in the solenoid model for chromatin higher-order structure (9) (see [Chromatin](#)). Acetylation might then result in disruption of higher-order structure, with the relaxed structures being more accessible to transcription factors, remodeling complexes, and the transcriptional apparatus (see [Chromatin](#)). Residues 16 to 24 of H4 are also known to be necessary for interaction with Sir3 in yeast silencing and [heterochromatin](#) formation, and in this case the suggestion is that the internucleosomal contact is replaced by an interaction with Sir3 (9).

### Bibliography

1. J. O. Thomas and R. D. Kornberg (1974) Chromatin structure: oligomers of the histones. *Science* **184**, 865–868; R. D. Kornberg (1974) Chromatin structure: a repeating unit of histones and DNA. *Science* **184**, 868.
2. J. Wong et al. (1998) *EMBO J.* **17**, 520–534.
3. W. An et al. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 3396–3401.
4. C. Crane-Robinson (1997) *Trends Biochem. Sci.* **22**, 75–77.
5. F. A. Goytisolo et al. (1996) *EMBO J.* **15**, 3421–3429.
6. D. Pruss et al. (1996) *Science* **274**, 614–617.
7. A. A. Travers (1998) *Trends Biochem. Sci.* (in press).
8. A. Hamiche et al. (1996) *J. Mol. Biol.* **257**, 30–42.
9. K. Luger et al. (1997) *Nature* **389**, 251–260.
10. G. Arents et al. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10148–10152.
11. Y.-B. Zhou et al. (1998) Position and orientation of the globular domain of linker histone H5 on the nucleosome. *Nature* (in press).
12. (1997) *Nature* **389**, 232–233.

### Suggestions for Further Reading

13. K. Luger and T. Richmond (1998) DNA binding within the nucleosome core. *Curr. Opin. Struct. Biol.* **8**, 33–40.
14. K. Luger and T. Richmond (1998) The histone tails of the nucleosome. *Curr. Opin. Genet. Dev.* **8**, 140–146.

### Nucleotide-Binding Motif

The nucleotide-binding motif is a structural and functional **domain** that is frequently observed in [protein structures](#) that bind nucleotides. It has a characteristic sequence that can usually be observed in the [primary structure](#). There are two classes of the motif; the classical dinucleotide-binding fold and the classical mononucleotide-binding fold. The former binds dinucleotides such as NAD and FAD, and the latter binds mononucleotides such as ATP (and is also called the [ATP-binding motif](#)).

The dinucleotide motif was first identified by Rao and Rossmann, who identified two b–a–b–a–b units (b is a **b-strand** and a is an **a-helix**) forming a six-stranded parallel **b-sheet** surrounded by four a-helices in the structures of [lactate dehydrogenase](#) (LDH) and malate dehydrogenase (Fig. 1) (1). Both bind NAD; each half of the dinucleotide interacts with one b–a–b–a–b unit. A single b–a–b–a–b unit is called a [Rossmann fold](#). The adenine moiety of NAD (or FAD) binds to the first of the two

b–a–b–a–b– units, and the nicotinamide (or flavin) binds to the second. The first three elements of **secondary structure** in the adenine-binding motif, b1–aA–b2, are involved in the predominant interactions with adenine and have a common structure and related sequences (1). This “fingerprint” sequence consists of conserved small **hydrophobic** residues in b1, followed by a glycine-rich region (–Gly–X–Gly–X–X–Gly, where X is any residue), conserved small hydrophobic residues in helix aA, and an Asp or Glu at the end of helix aA. The first two **glycine** residues are in a loop connecting b1 and aA; the third glycine is in helix aA. The glycine residues adopt conformations forbidden to other residues and allow close packing of the b-strands and helix and a close approach between the adenine pyrophosphate and the *N*-terminal end of the aA-helix. The phosphate group is thought to interact favorably with the partial positive charge at the *N*-terminal end of the aA-helix dipole. The conserved Asp/Glu residue interacts with the ribose hydroxyl. This signature sequence of ~30 residues was derived from the tertiary structures of known NAD-binding proteins, but it can also be used to predict NAD- or FAD-binding proteins from their primary structure.

**Figure 1.** Schematic representation of the backbone structure of the nucleotide-binding domain of malate dehydrogenase (1). b-Strands are shown as arrows and a-helices are shown as coils. The bound dinucleotide (NAD) is shown as a ball-and-stick model. The two b–a–b–a–b– units that interact with the each half of the mononucleotide are shown in green and purple. The b1–aA–b2 region that interacts with the adenine portion of NAD is shown in yellow. This figure is generated using Molscrip (2) and Raster3D (3, 4). See color insert.



The fingerprint sequence is slightly different for binding the dinucleotide NADP, where the 2' ribose hydroxyl of NAD is phosphorylated. In this case, the third glycine (in the aA-helix) is replaced by Ala so that the close approach of aA and the b-strands is prevented. This provides additional space for binding the 2'-phosphate. Furthermore, the conserved Asp/Glu of the NAD/FAD sequence is replaced by Arg to interact with the 2'-phosphate.

The classical mononucleotide or **ATP-binding motif**, typified by adenylate kinase, is also an a/b

protein fold but has a five-stranded parallel  $\beta$ -sheet with a very different connectivity from the dinucleotide fold. Once again, the bound nucleotide binds to a b1–aA–b2 region that has a characteristic fingerprint sequence: Gly–X–X–Gly–X–Gly–Lys. In this case, however, the glycine-rich sequence corresponds only to the loop between b1 and aA, so that all three glycine residues form part of the loop (in contrast to the dinucleotide sequence, where the third glycine is part of helix aA). The loop forms a large pocket that accommodates the phosphate of the mononucleotide (it is also called the [P-loop](#)), and the glycine amides and the side chain of the conserved lysine residue interact with the phosphate.

[See also [Domain](#), [protein](#), [ATP-binding motif](#), and [P-loop](#).]

### Bibliography

1. C. A. Kelly et al. (1993) *Biochemistry* **32**, 3913–3922.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

5. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
6. S. T. Rao and M. G. Rossmann (1973) Comparison of supersecondary structure in proteins. *J. Mol. Biol.* **76**, 241–256 (Original description of the nucleotide-binding fold based on structural similarities between lactate dehydrogenase, malate dehydrogenase, and flavodoxin.)
7. G. E. Schulz (1992) Binding of nucleotides by proteins. *Curr. Opin. Struct. Biol.* **2**, 61–67 (Excellent review comparing and contrasting different protein folds that bind nucleotides.)

## Nucleotides, Nucleosides, And Nucleobases

Nucleotides are the monomeric constituents of **nucleic acids**, as well as structural elements of most **coenzymes** and many activated metabolic intermediates. This article describes the structure, nomenclature, and chemical and physical properties of nucleotides that serve as nucleic acid constituents.

### 1. Definitions

A nucleotide is a substance that, on hydrolysis, yields per mole at least 1 mole of a nitrogenous base, a sugar, and orthophosphate. A nucleoside yields 1 mole each of a nitrogenous base and a sugar. The term *nucleobase* refers to those nucleotides found in nucleic acids—adenine, guanine, cytosine, and uracil in RNA, and adenine, guanine, cytosine, and thymine in DNA, plus the quantitatively minor bases found in both nucleic acids. Other bases, such as the nicotinamide found in  $\text{NAD}^+$  and  $\text{NADP}^+$ , occur as parts of coenzymes or as parts of metabolically activated biosynthetic intermediates, such as uridine diphosphate glucose.

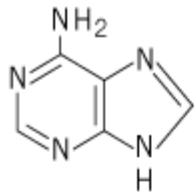
A mononucleotide is a nucleotide that, on hydrolysis, yields 1 mole each of a base and a sugar plus at least 1 mole of orthophosphate. The term nucleoside diphosphate usually refers to a mononucleotide which, on hydrolysis, yields 1 mole each of a base and a sugar plus 2 moles of orthophosphate; similarly, a nucleoside triphosphate yields 3 moles of orthophosphate per mole of

nucleotide. A dinucleotide yields, on complete hydrolysis, 2 moles each of base and sugar, plus at least 1 mole of orthophosphate; a trinucleotide yields 3 moles each of base and sugar, plus 2 or more moles of orthophosphate.

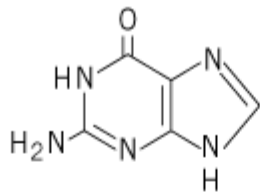
Figure 1 shows the structures of the five common nucleobases (in their most frequent **tautomeric** configuration) and several nucleosides, mononucleotides, and oligonucleotides, as well as the numbering systems used. Note that, in all nucleotides found in nucleic acids, a glycosidic bond links N-9 of a purine base or N-1 of a pyrimidine base to C-1, the carbonyl carbon, of either ribose or 2-deoxyribose. By convention, all positions in the sugar of a nucleoside or nucleotide are given primed numbers. Thus, the adenine nucleotide recovered from digestion of DNA with certain enzymes would be named 2'-deoxyadenosine 5'-monophosphate, identifying the position of the phosphate esterified to the sugar, as well as the carbon of the sugar that is in the reduced, or deoxy, configuration. Note that some modes of nucleic acid digestion will yield 2'- or 3'-nucleotides, in which the phosphate is esterified to positions 2' or 3' of the sugar, respectively. Also, in the 3',5'-cyclic nucleotides, **cyclic AMP** and **cyclic GMP**, which play multiple roles in biological [signal transduction](#), the phosphate is esterified simultaneously to carbons 3' and 5'. Because the nucleic acid biosynthetic intermediates are nucleoside 5'-phosphates, or 5'-nucleotides, these structures are shown predominantly. Table 1 gives the names and common abbreviations of the nucleotide constituents of nucleic acids and the cyclic nucleotides.

**Figure 1.** Structures of the five common nucleobases and representative nucleosides and nucleotides.

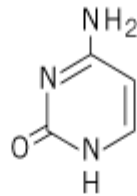
### Nucleobases



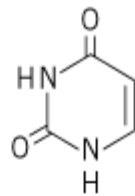
Adenine  
Ade



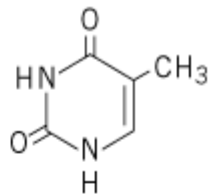
Guanine  
Gua



Cytosine  
Cyt

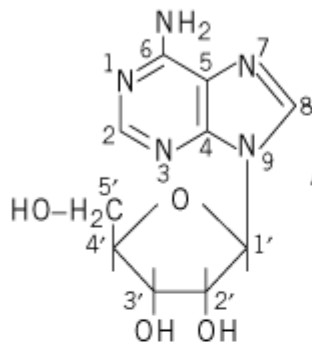


Uracil  
Ura

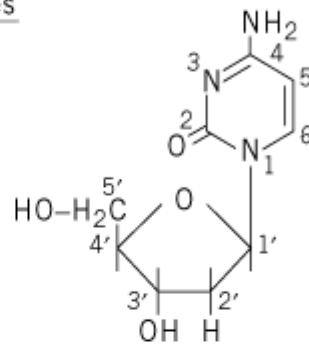


Thymine  
Thy

### Nucleosides

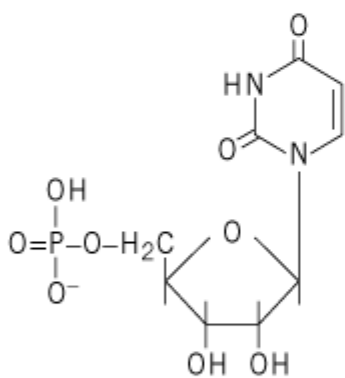


Adenosine  
Ado

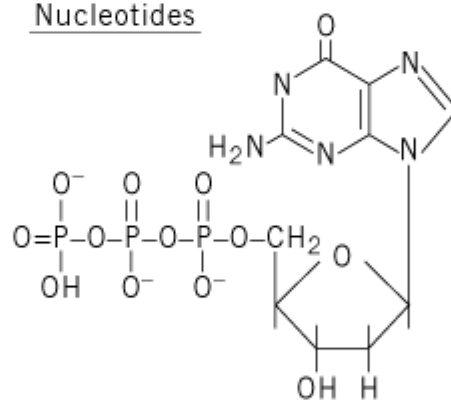


Deoxycytidine  
dCyd

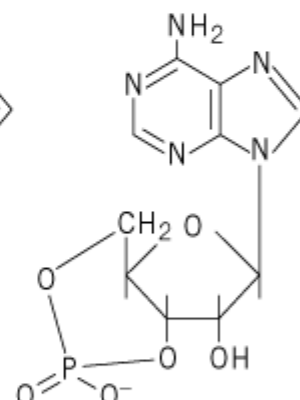
### Nucleotides



Uridine 5'-monophosphate  
UMP

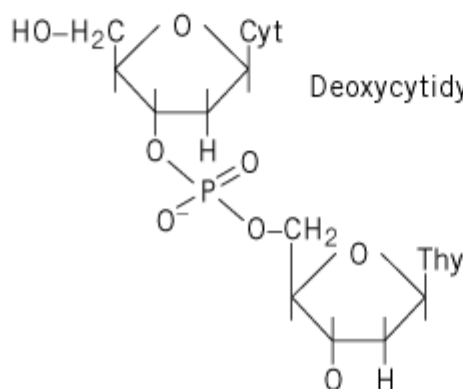


Deoxyguanosine  
5'-triphosphate  
dGTP



Adenosine  
3', 5'-monophosphate  
cyclic AMP

### A dinucleotide



Deoxycytidylylthymidine 3'-monophosphate  
dCpdTp

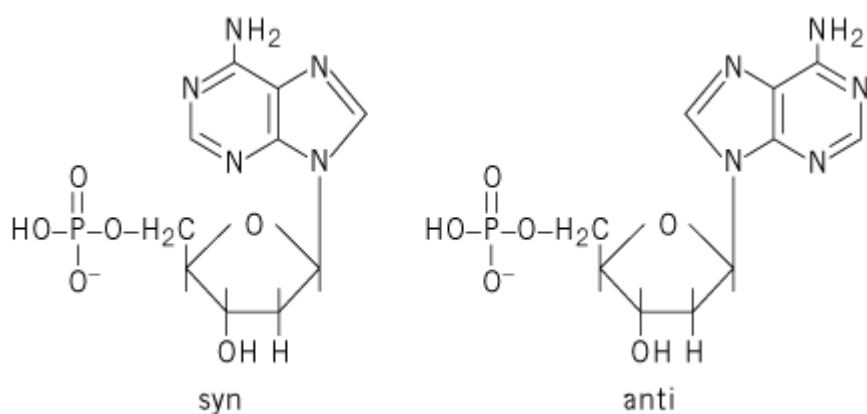
**Table 1. Names and Abbreviations of the Nucleotide Constituents of Nucleic Acids and of Cyclic Nucleotides**

| Names   | Abbreviations                |
|---|------------------------------|
| Adenosine 5'-monophosphate, 5'-adenylic acid  | AMP, Ado-5'-P, pA            |
| Cytidine 5'-monophosphate, 5'-cytidylic acid  | CMP, Cyd-5'-P, pC            |
| Guanosine 5'-monophosphate, 5'-guanylic acid  | GMP, Guo-5'-P, pG            |
| Uridine 5'-monophosphate, 5'-uridylic acid  | UMP, Urd-5'-P, pU            |
| Deoxyadenosine 5'-monophosphate, 5'-deoxyadenylic acid  | dAMP, dAdo-5'-P, pdA         |
| Deoxycytidine 5'-monophosphate, 5'-deoxycytidylic acid  | dCMP, dCyd-5'-P, pdC         |
| Deoxyguanosine 5'-monophosphate, 5'-deoxyguanylic acid  | dGMP, dGuo-5'-P, pdG         |
| Deoxythymidine 5'-monophosphate, 5'-deoxythymidylic acid, thymidine 5'-monophosphate, 5'-thymidylic acid <sup>a</sup> | dTMP, dThd-5'-P, pdT         |
| Adenosine 3',5'-monophosphate, 3',5'-cyclicadenylic acid  | cAMP, Ado-3',5'-P, 3',5'-AMP |
| Guanosine 3',5'-monophosphate, 3',5'-cyclicguanylic acid  | cGMP, Guo-3',5'-P, 3',5'-GMP |

<sup>a</sup> Thymine occurs as a minor base constituent of tRNAs. However, because the thymine ribonucleotides do not occur in nature, nucleotides containing thymine are called either *thymidine* or *deoxythymidine nucleotides*. No consistent preference for either nomenclature has emerged.

One additional structural feature of nucleotides relates to the somewhat hindered rotation about the glycosidic bond, which leads to two rather stable orientations of bases with respect to the sugar, termed *syn* and *anti* (Fig. 2). Because it is more compact to draw, most of the structures shown here are represented as *syn*, but the *anti* configuration predominates in nucleic acids (see [Z-DNA](#)).

**Figure 2.** dAMP in the *syn* and *anti* configurations.



## 2. Chemical and Physical Properties of Nucleotides

Because of the phosphate group or groups on nucleotides, these substances are highly acidic and exist as anions at physiological pH. The primary and secondary  $pK_a$  values for the phosphate groups of nucleoside 5'-monophosphates are about 1 and 6, respectively. The [amino groups](#) on adenine, guanine, and cytosine rings are extremely weak bases, undergoing protonation with  $pK_a$  values of 2 to 4.5. Because of variations in these values, and the absence of a ring amino group in uracil and thymine (see [Table 2](#) and [ref. 2](#)), one can adjust the pH to values at which different nucleotides have different net negative charges, thereby permitting their ready separation by [electrophoresis](#) or [ion-exchange chromatography](#) ([2](#)).

**Table 2. Properties of Nucleoside Monophosphates**

| Nucleotide | $pK_a$ of ring $-NH_2$ | UV absorbance maximum, nm | Measured at pH | $\epsilon_m, l^2/\text{mole-cm} \times 10^{-3}$ |
|------------|------------------------|---------------------------|----------------|---|
| AMP        | 3.8                    | 257                       | 2              | 15.0  |
| CMP        | 4.5                    | 280                       | 2              | 13.2  |
| GMP        | 2.4                    | 256                       | 1              | 12.2  |
| UMP        | —                      | 262                       | 2              | 10.0  |
| dAMP       | 4.4                    | 258                       | 2              | 14.3  |
| dCMP       | 4.6                    | 280                       | 2              | 13.5  |
| dGMP       | 2.9                    | 255                       | 1              | 11.8  |
| dTMP       | —                      | 267                       | 2              | 10.2  |
| cAMP       | N.R.                   | 256                       | 2              | 14.5  |
| cGMP       | N.R.                   | 256.5                     | 1              | 11.35 <sup>b</sup>                              |

<sup>a</sup> N.R., not reported.

<sup>b</sup> Data for cyclic nucleotides are from Ref. 10. Other data are from Ref. [1](#).

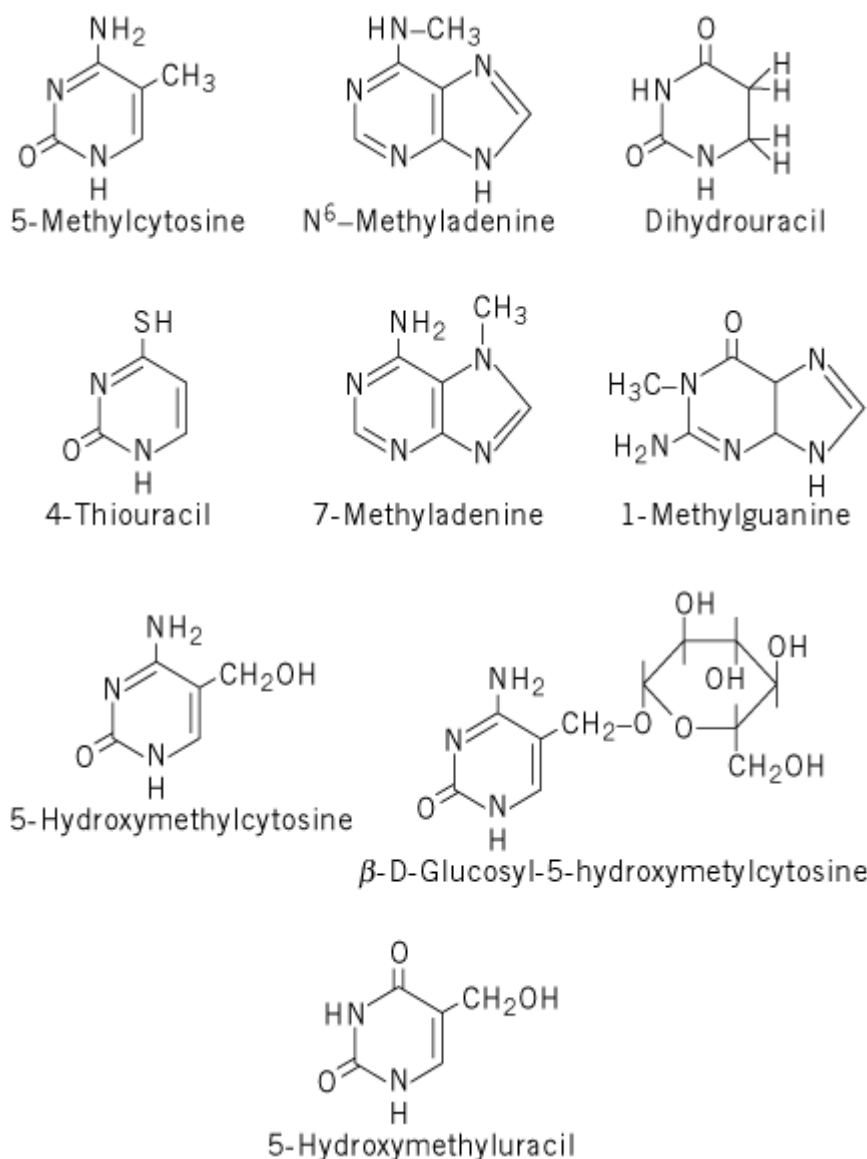


The aromatic character of the purine and pyrimidine rings leads the nucleotides to **absorb** ultraviolet light with characteristic absorption spectra, as indicated in Table 2. This allows ready detection, identification, and quantification of nucleotides and nucleotide derivatives. Because of ionization of the purine and pyrimidine rings, the ultraviolet absorption spectra are quite dependent on pH, and this has made identification and quantification methods even more specific and sensitive.

### 3. Naturally Occurring Modified Nucleotides

In addition to the five “standard” nucleobases (adenine, guanine, cytosine, thymine, and uracil), most nucleic acids contain modified bases (see [Transfer RNA](#) and [DNA Structure](#)). The most widespread is [5-methylcytosine](#), which is present at up to 5% of the cytosine level in most eukaryotic DNA (Fig. 3). The most abundant methylated base in prokaryotic DNA is N<sup>6</sup>-methyladenine (see [Methyltransferase](#), and [Mismatch Repair](#)). A variety of modified bases, some but not all of them methylated, are found in RNA, most abundantly in transfer RNA. A few of these are shown in Figure 3.

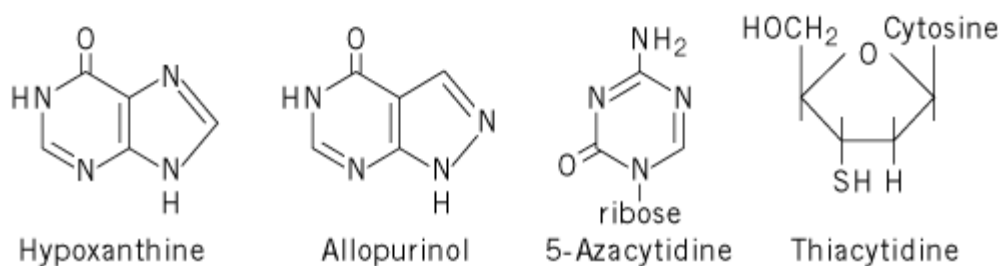
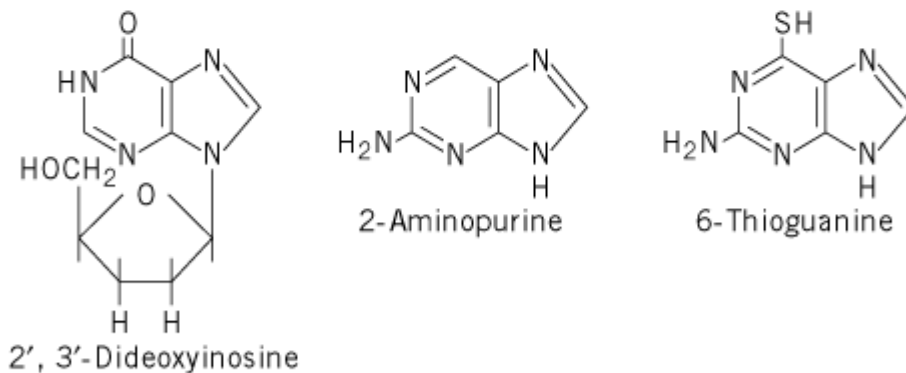
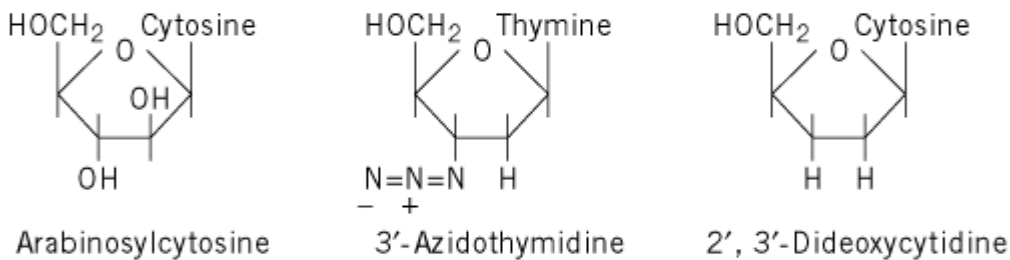
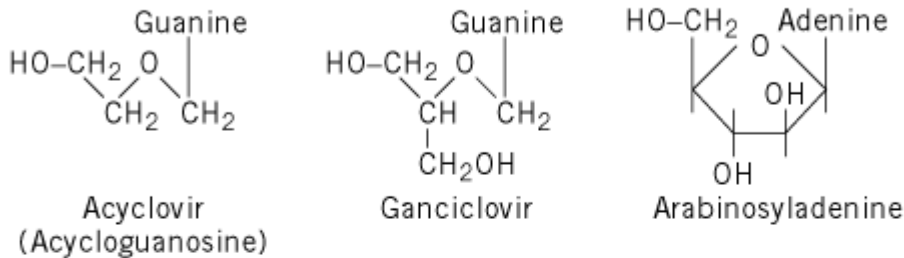
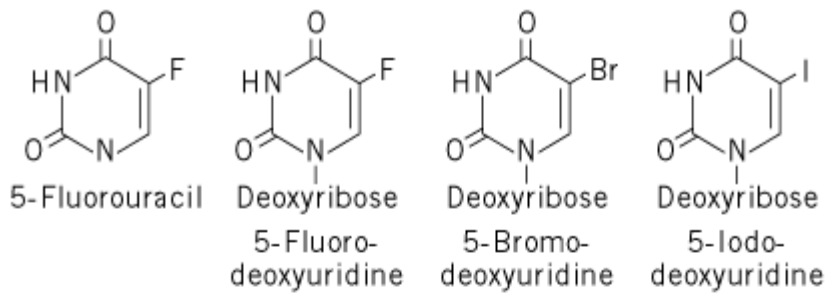
**Figure 3.** Structures of prominent modified nucleobases found in nucleic acids.



In addition to the above modified bases, which are quantitatively minor in most systems studied, there are several bacteriophages whose DNA contains modified bases as major constituents. Most prominent among these are the T-even coliphages, which contain four forms of 5-hydroxymethylcytosine, completely substituted for cytosine (Fig. 3): the base itself, with no further modifications, and the base glycosylated at the hydroxymethyl position with glucose in either the a or b configuration, or with cellobiose (glucosyl-1,4-b-D-glucose). Several phages infecting *Bacillus subtilis* contain 5-hydroxymethyluracil completely replacing thymine, and bacteriophage XP-12, which infects *Xanthomonas oryzae*, contains 5-methylcytosine completely replacing cytosine. 5-Hydroxymethyluracil has also been found as a minor constituent of DNA subjected to oxidative damage; it results from oxidation of thymine residues (3).

#### 4. Synthetic Nucleoside and Base Analogues

A large number of nucleoside and nucleobase analogues have been synthesized and developed as anticancer, antiviral, antibacterial, and antiparasitic drugs (4-6). Many of these analogues act by intracellular conversion to deoxyribonucleotide analogues, which selectively block [DNA replication](#) (because deoxyribonucleotides have no significant metabolic roles other than as DNA constituents). A few of the most prominent of these analogues are shown in Figure 4 and discussed briefly here.



#### 4.1. Halogenated Pyrimidines

The earliest-developed useful synthetic nucleobase and nucleoside analogues were 5-fluorouracil (FUra) and **5-fluorodeoxyuridine** (FdUrd), both of which were synthesized during the 1950s and shown to have potent antitumor activity; both agents are still in use today. Both the base and the nucleoside act by uptake into target cells, followed by intracellular conversion to 5-fluorodeoxyuridine 5'-monophosphate (FdUMP), an analogue of deoxyuridine monophosphate, which binds covalently to [thymidylate synthase](#), thereby blocking the synthesis of dTTP and hence

the replication of DNA.

Because of their negative charge, nucleotides are not readily transported through cell membranes. With most nucleobase or nucleoside drugs, the active form of the drug is a nucleotide analogue—FdUMP in the example under discussion. Thus, in order to use nucleic acid precursors successfully as drugs, the biochemical pharmacologist must be aware of the nucleotide biosynthetic salvage enzymes in target cells (see [Salvage Pathways To Nucleotide Biosynthesis](#)), so that a precursor is administered that can be converted to the active nucleotide. One must be concerned also with degradative enzymes. For example, nucleoside phosphorylase can cleave fluorodeoxyuridine to 5-fluorouracil plus deoxyribose-1-phosphate, a reaction that would limit the effectiveness of the drug by decreasing its availability for conversion to the active FdUMP.

Closely related to 5-fluorodeoxyuridine are the other halogenated pyrimidines, 5-bromodeoxyuridine (BrdUrd) and 5-iododeoxyuridine (IdUrd). The bromine atom of BrdUrd has a van der Waals radius similar to that of a methyl group; thus, it is readily phosphorylated by thymidine kinase and ultimately converted to the nucleoside triphosphate BrdUTP, which is a close dTTP analogue and is readily incorporated into DNA. BrdUrd has limited therapeutic utility, but it is widely used as a research reagent, for density-labeling experiments involving DNA; the heavy bromine atom imparts a significantly higher [buoyant density](#) to BrdUrd-substituted DNA than to thymidine-containing DNA. Moreover, BrdUrd is a potent chemical [mutagen](#) because bromouracil in DNA occasionally base-pairs with guanine and not with thymine in DNA (6); also, see [Mutagenesis](#)). By contrast, IdUrd is a fairly useful antiviral drug, used in treating eye infections caused by the large DNA virus, **herpes** simplex. The genome of this virus encodes a thymidine kinase of broad specificity that readily phosphorylates IdUrd to the monophosphate, and then to the diphosphate. Thus, this nucleotide is more readily incorporated into DNA of virus-infected cells, where its presence can damage the viral genome, than in uninfected cells.

#### 4.2. Arabinosyl Nucleosides

These compounds, of which the adenosine and cytidine analogues are most useful therapeutically (6), contain arabinose, which is the 2-epimer of ribose. AraA is a fairly effective antiviral agent, being used to treat viral encephalitis, among other conditions. The triphosphate of araA, araATP, is a rather specific inhibitor of DNA polymerases encoded by herpes viruses. AraC is used in chemotherapy of leukemias and some solid tumors. AraC is phosphorylated initially by deoxycytidine kinase, ultimately being converted to araCTP, a potent inhibitor of replicative DNA polymerases. A problem in the use of araC is its fairly facile deamination by cytidine deaminase, yielding the weakly effective arabinosyluracil. Therapeutic activity of araC can be potentiated by coadministration of a cytidine deaminase inhibitor. In the same sense, the efficacy of araA is often enhanced by coadministration of an adenosine deaminase inhibitor, because this enzyme converts arabinosyladenine to arabinosylhypoxanthine.

#### 4.3. 6-Thioguanine

This guanine analogue is widely used in **somatic cell** genetics. Its conversion to a toxic nucleotide is carried out by hypoxanthine guanine phosphoribosyltransferase (HGPRT). Since the structural gene for HGPRT is sex-linked, mammalian genomes contain only one copy of the gene, making this gene a convenient selectable marker for studies on mutagenesis and for generating [hybridoma](#) cell lines for [monoclonal antibody](#) production. HGPRT-negative mammalian cell mutants become resistant to the lethal effects of 6-thioguanine, because of their inability to anabolize the analogue.

#### 4.4. Allopurinol

This hypoxanthine analogue is a successful drug for treating gout (7). It is an effective inhibitor of xanthine oxidase, which converts hypoxanthine to xanthine, and xanthine to uric acid in the catabolism of purine nucleotides. Thus, allopurinol often inhibits the buildup of the insoluble uric acid, allowing excess purines to be excreted as the more soluble hypoxanthine and xanthine.

#### 4.5. 5-Azacytidine

This cytidine analogue has a nitrogen atom at position 5 of the modified pyrimidine ring. It has been useful in studies of the biological significance of **DNA methylation** because it is metabolized similarly to cytidine. When incorporated into DNA as the dCMP analogue, however, the pyrimidine ring cannot be methylated at position 5. Thus, a G-mC base pair is effectively demethylated in subsequent rounds of replication, allowing investigators to investigate the [epigenetic](#) changes that occur as a result.

### Bibliography

1. D. B. Dunn and R. H. Hall (1968) In *Handbook of Biochemistry* (H. A. Sober, ed.), Chemical Rubber Co., Cleveland, OH, pp. G-3–G-86.
2. P. R. Brown, ed. (1984) *HPLC in Nucleic Acid Research*. Marcel Dekker, New York.
3. R. A. Floyd (1990) *FASEB J* **4**, 2587–2597.
4. R. I. Christopherson and S. D. Lyons (1990) *Med. Res. Revs.* **10**, 505–549.
5. H. Mitsuya, ed. (1997) *Anti-HIV Nucleosides: Past, Present, and Future*, R. G. Landes Co., Georgetown, TX.
6. T. W. North (1991) *Encyclopedia of Human Biology* (Academic Press, San Diego) **2**, 395–402.
7. G. B. Elion (1989) *Science* **244**, 41–47.
8. S. Budavari, ed. (1989) *The Merck Index*, 11th ed., Merck & Co., Rahway, NJ, pp. 2717–2718.

### Suggestion for Further Reading

9. C. K. Mathews and K. E. van Holde (1996) *Biochemistry*, 2nd ed., Benjamin/Cummings, Menlo Park, CA, Chapters "4" and "22". Most basic information about structure and properties of nucleic acid constituents can be found in a standard biochemistry textbook.

## Nucleus

The nucleus is the center for **genetic inheritance** and **gene expression** within the cells of **eukaryotes**. It contains a myriad of reactions that occur within several macromolecular compartments and particles. This article provides an overview of nuclear structure and function. Many of the more specific topics are treated in more detail in their entries.

### 1. Genomic Inheritance

The nucleus is the most conspicuous **organelle** within eukaryotic cells. Although earlier microscopists made note of what must have been nuclei(1), historians attribute the discovery of the nucleus to Robert Brown in 1833 because of his careful descriptions of orchid cell nuclei(2). Once the nucleus was firmly established by Brown as a constant entity within cells, Hertwig and Strasburger in 1884 independently described the nucleus as the organelle of cellular inheritance. These and other late nineteenth and early twentieth century researchers theorized that inheritance is based on segregation of the nuclear [chromosomes](#) (3, 4).

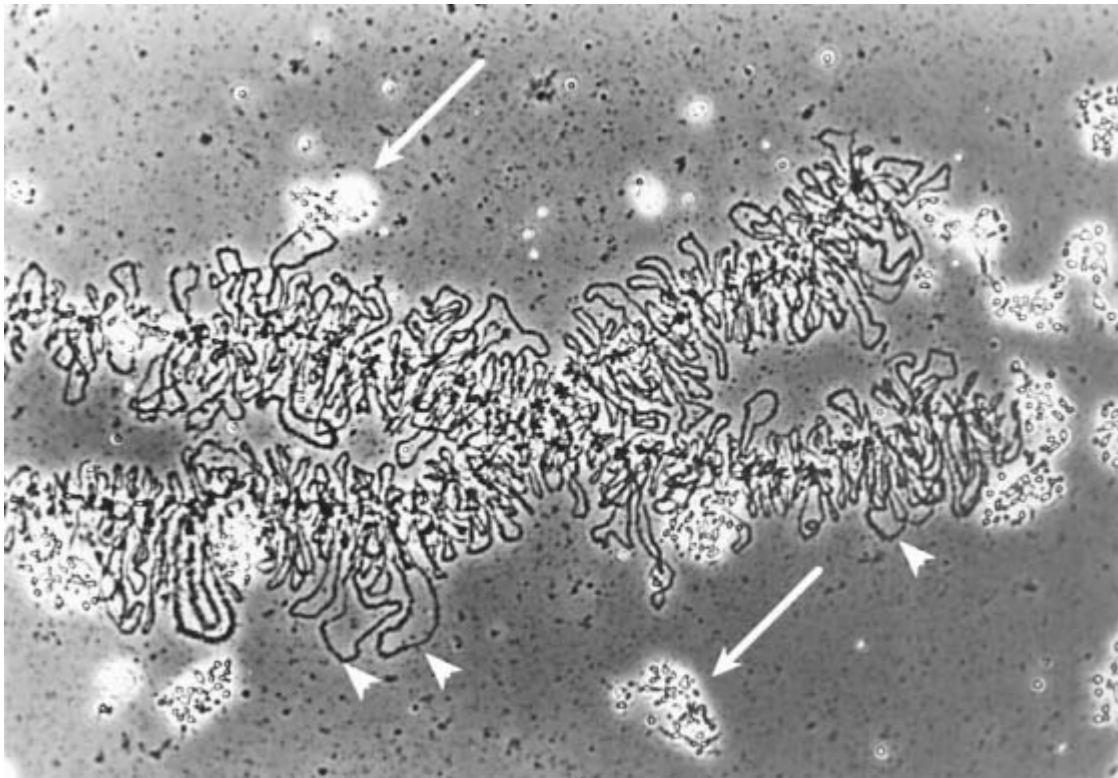
Most eukaryotic cells contain one nucleus, but there are many examples of multinucleated cells or syncytia (eg, differentiated muscle cells, the early *Drosophila* [embryo](#)). Conversely, mature mammalian red blood cells lack nuclei. Nuclei differ in size depending on the cell type: the **yeast** (*Saccharomyces cerevisiae*) nucleus is approximately 2  $\mu\text{m}$  in diameter, a typical mammalian nucleus is 4 to 6  $\mu\text{m}$ , and nuclei (germinal vesicles) in mature amphibian **oocytes** are about 600  $\mu\text{m}$

in diameter. Most nuclei are spherical, but multilobed nuclei are common, such as those found in polymorphonuclear leukocytes or mammalian epididymal cells. The nucleus is bordered by a double-membrane [nuclear envelope](#) that contains **nuclear pore complexes** for macromolecular import and export from and to the cytoplasm.

The nucleus is the repository for the complex eukaryotic [genome](#). It normally consists of unique, middle-repetitive, and highly repetitive sequences, defined by renaturation kinetics(5) (see **C<sub>0</sub>t curve**). Generally, unique DNA includes single-copy **genes** and unique but nontranscribed sequences. Middle-repetitive DNA contains genes with multiple copies, such as those found in the tandem arrays of repeated genes for [histones](#) or the 47S RNA precursor for [ribosomes](#) and the spacer sequences that separate the individual genes within these clusters. Highly repetitive sequences include relatively short DNA elements that constitute, for example, the [centromeres](#) of chromosomes. Highly repetitive, noncentromeric sequences have also been detected. Examples include the 5S ribosomal genes (over 20,000 copies per haploid *Xenopus* genome), the L1 [transposable elements](#) that constitute about 4% of primate genomes, and the *Alu* [pseudogene](#)/transposable elements that constitute about 5% of the human genome .

The eukaryotic cell must efficiently package all of its genomic DNA (which in human cells has a combined linear length of more than one meter) within a nucleus with a diameter of only 5  $\mu\text{m}$ . Yet the manner in which the DNA is packaged must also accommodate the complex processes of [DNA replication](#), [recombination](#), [DNA repair](#), and [transcription](#) . [Chromatin](#) is the term used to describe the packaged genome. It consists of genomic DNA, a nearly equal mass of [histone](#) proteins, a relatively small amount of non-histone chromosomal protein, and some nascent RNA. The most basic structural unit of chromatin is the [nucleosome](#), which consists of approximately 146 bp of genomic DNA wrapped around a complex of eight histone molecules. In the presence of histone H<sub>1</sub>, nucleosomes condense to form a solenoid with a diameter of 30 nm and six nucleosomes per turn. The solenoid is cast into looped domains that are attached to the interphase [nuclear matrix](#) by adenine/thymine-rich DNA elements called [matrix attachment regions](#) (MARS). Such compacted chromatin is then organized into one or more [chromosomes](#). The complexity of chromosomal structure within the nucleus is well illustrated by the [lampbrush chromosomes](#) of the nucleus of the amphibian **oocyte** (Fig. 1).

**Figure 1.** A phase contrast image of portions of lampbrush chromosomes isolated from an oocyte nucleus of the North American newt, *Notophthalmus viridescens*. The chromosomes were centrifuged onto a microscope slide and fixed with formaldehyde. Lateral loops of transcriptionally active DNA display nascent ribonucleoprotein matrices (arrow heads). Multiple nucleoli (arrows) within the oocyte nucleus produce enormous numbers of ribosomes. These particular nucleoli have partially fragmented during preparation. Many other small ribonucleoprotein particles are clearly evident (see text). Photomicrograph was kindly provided by Dr. Joe Gall, Carnegie Institution of Washington, Department of Embryology.



Ultrastructural examination of [glutaraldehyde](#)-fixed interphase cells reveals patches of condensed chromatin within their nuclei, usually just beneath the nuclear envelope or near the [nucleolus](#). The term '**heterochromatin**' was defined by E. Heitz in 1928 to describe chromatin that remains condensed when the chromosomes normally decondense in the late **telophase** and early **interphase** of the [cell cycle](#). Because of its condensed state during the interphase, this heterochromatin is generally believed to be transcriptionally silent. Different cell types contain different amounts of heterochromatin. Metabolically active cancer cells, for example, display very little condensed chromatin, whereas orthochromatic erythroblasts (red blood cells that will soon lose their nuclei) contain large amounts of condensed chromatin(7). There are actually two types of heterochromatin. Constitutive heterochromatin remains condensed throughout the cell cycle and [development](#). This chromatin contains highly repetitive sequences that play a structural role in chromosomal structure and mechanics (ie, centromeric DNA functioning in chromosomal movement during **mitosis**) rather than actual **gene expression**. Facultative heterochromatin has the potential for gene expression at some point in development, and it can be either condensed or decondensed, depending on the cell type. The mammalian [Barr body](#) is an excellent example of facultative heterochromatin. The Barr body is one of the two **X-chromosomes** in mammalian female cells. It remains condensed during the interphase and is almost completely silent in transcription(8), whereas the other X-chromosome in the cell is decondensed and transcriptionally very active. [Euchromatin](#) is the term used to describe decondensed interphase chromatin that is either transcriptionally active or has the potential for transcription.

Replication of eukaryotic chromosomes requires elaborate DNA polymerase complexes that replicate the entire genome within the physical confines of the nucleus (see **DNA replication**). To replicate the enormous lengths of eukaryotic chromosomes, several [autonomously replicating sequences](#) (ARS), situated along the length of each chromosome, initiate replication in both directions along the DNA, thus leading to rapid replication of the entire chromosome. Eukaryotic genomic replication occurs during the interphase of the cell cycle, specifically during the **S-phase** between growth phases 1 and 2 ( $G_1$  and  $G_2$ ). The precisely timed onset of DNA replication is controlled by elaborate [signal transduction](#) cascades that work in conjunction with [cyclin](#)-dependent

**kinases** (see [Cell Cycle](#)). Enhanced production of new histone proteins is concomitant with genomic DNA replication to ensure the timely packaging of the nascent DNA. **Telomeres** are unique structures at the very ends of chromosomes that provide for completion of their replication(9) and also interact with the nuclear envelope during **meiosis** (10, 11), perhaps to help align the chromosomes for recombination.

Fully replicated chromosomes condense during mitosis in preparation for orderly segregation into progeny cells. Each mitotic chromosome consists of two sister **chromatids** held together by their centromeres. Centromeres play a structural/mechanical function in chromatid segregation during mitosis and meiosis. They are sites for peripheral **kinetochoric** assembly during the **prophase**, one kinetochore for each chromatid. A set of microtubules from the mitotic spindle attaches to the kinetochores to align the chromosomes on the **metaphase** plate. At **anaphase**, the centromeric connection between the chromatids splits, and the kinetochoric microtubules depolymerize to “pull” the separated chromatids (now called chromosomes) toward opposite poles of the spindle.

## 2. Genomic Expression

[Transcription](#) is the initial step in **gene expression** in the nucleus. Transcription units (TU) are defined as thin-to-thick [ribonucleoprotein](#) (RNP) fibrils that remain attached to active **RNA polymerase** complexes on genes as viewed by [electron microscopy](#) of spread chromatin preparations (13) or of the transcriptionally active loops of [lampbrush chromosomes](#) (14; Fig. 1). The nascent RNP fibrils are shortest near transcription start sites, and they are longest near transcription stop sites (the so-called “Christmas tree”). The pre-[messenger RNA](#) resulting from transcription of genes that encode proteins is also known as heterogeneous nuclear RNA (**hnRNA**). Almost immediately after its synthesis by RNA polymerase II, the 5' ends of the nascent hnRNA are capped with a 7-methyl guanine nucleotide (see **5'-Cap**). The cap is important for transporting the mature mRNA from the nucleus, for protecting of the 5' end from exonucleases, and for recognizing the mRNA by translational initiation factors within the cytoplasm. [Polyadenylation](#) of pre-mRNA also occurs in the nucleus before mRNA export, and is generally believed to stabilize the transcripts in the cytoplasm.

The ribosomal RNA (rRNA) genes are transcribed within the [nucleolus](#). The precursor rRNA is processed to yield the mature 18S, 5.8S, and 28S rRNA found in cytoplasmic ribosomes.

### 2.1. Spliceosomes and Heterogeneous Nuclear RNA Processing

One of the most important aspects of eukaryotic gene expression, other than the transcription of the gene itself, is the precise splicing of premessenger RNA (pre-mRNA) transcripts, first to remove nonprotein coding **introns** and then to join (ligate) the protein-coding **exons** (see [RNA Splicing](#)). Both events must occur precisely to maintain the proper protein translational reading frame. Sequences within the nascent transcript direct the assembly of the nuclear splicing machinery. Nuclear ribonucleoprotein (RNP) particles that mediate the transesterification and ligation reactions of splicing are called **spliceosomes**. They consist of five different small nuclear RNP particles (snRNP). Each snRNP consists of a unique **small nuclear RNA** (snRNA), either U1, U2, U4, U5, or U6. The snRNP particle is referred to by its constituent snRNA (eg, the U1 snRNP). In addition to the snRNA, each snRNP consists of proteins that are either unique to the specific snRNP particle (eg, the 70-kDa U1-specific protein) or common to the different snRNP (ie, the snRNP core proteins). The common proteins initiate snRNP complex formation upon recognizing a single-stranded region of the snRNA at the Sm site.

Spliceosome assembly is a complex process involving the various snRNPs and the pre-mRNA. Auxiliary proteins are also involved in either spliceosome assembly or pre-mRNA splice site selection where alternative donor and acceptor splice site choices are possible. Upon excising the intron and ligating of the two exons, the spliceosome disassembles, and the ligated, processed mRNA is released. Individual snRNPs are conserved and reutilized for other splicing reactions.



The genes for [transfer RNA](#) in yeast, wheat germ, and vertebrates also contain small introns (6 to 80 nucleotides) just downstream of the central anticodon region. These introns are removed by endonuclease activities, and the two halves of the tRNA are ligated(15).

### 3. Subnuclear Organelles, Compartments, Particles, and their Functions

#### 3.1. The Nuclear Envelope

The nucleus is bordered by two [membrane](#) bilayers—the [nuclear envelope](#). Metabolically active cells (eg, the *Xenopus* oocyte) may display folds or blebs in the nuclear envelope as a mechanism to increase its surface area. The outer membrane is continuous with that of the [endoplasmic reticulum](#). In fact, some electron micrographs show [ribosomes](#) on the surface of the outer nuclear membrane. The space between the two nuclear membranes, the [perinuclear space](#), is continuous with the lumen of the endoplasmic reticulum. The nuclear side of the inner membrane is lined with a laminar matrix (the nuclear lamina) that consists of nuclear **lamins** A, B, and C, each an [intermediate filament](#)-type protein. The lamina lends support and shape to the nucleus and is considered part of the [nuclear matrix](#).

The inner and outer membranes fuse at the [nuclear pore complexes](#) that traverse both bilayers. Nuclear pore complexes mediate bidirectional transport of macromolecules between the nucleus and the cytoplasm (see [Nuclear Import, Export](#)). An individual pore complex displays eightfold symmetry and is a huge macromolecular assembly with a combined mass of approximately 124 megadaltons.

#### 3.2. Nuclear Matrix

The nuclear matrix is a three-dimensional fibrogranular latticework that permeates the nucleus but has been elusive to define completely. Components include the nuclear envelope with its pores, interchromatin RNP granules, perichromatin RNP fibrils, the nucleolus, tightly bound chromatin ([matrix attachment regions](#), MARS), and associated pre-mRNA. Chromatin does not simply occupy random placement within the interphase nucleus (16, 17). Instead, specific regions of euchromatin and heterochromatin associate with the nuclear envelope(18). Observations such as these suggest that specific three-dimensional positioning of genetic loci occurs within the nucleus. This is an important aspect of the gene gating hypothesis of Blobel(19), which originally proposed that nuclear pore complexes, the nuclear lamina, and components of the nuclear core (the matrix) topologically position (eg, gate) transcribable chromatin near the envelope pores for efficient export of the mRNA.

#### 3.3. The Nucleolus and Ribosomal Subunit Production

The nucleolus is the most conspicuous organelle within the interphase nucleus. Our traditional understanding of the nucleolus is that it is responsible for (1) synthesizing a large nascent precursor ribosomal RNA (pre-rRNA) that is 47S in mammals and approximately 13,000 nucleotides long in humans; (2) the processing (cleavage and base modification) of this RNA to yield mature ribosomal RNA of 18S (2,000 nucleotides), 5.8S (160 nucleotides), and 28S (5,000 nucleotides); and (3) the concomitant assembly of these RNAs with incoming ribosomal proteins to generate small and large ribosomal subunits that then pass into the cytoplasm (see [Nucleolus](#)).

Three distinct ultrastructural regions constitute the prototypical interphase nucleolus: the fibrillar center (FC), the surrounding dense fibrillar region (DFR), and the peripheral granular region (GR). Particles of preribosomal subunits are visible within the GR and are 15 to 20 nm in size. Vacuoles are sometimes found within nucleoli, but their significance is unknown. The relative sizes and distributions of the three principal nucleolar compartments vary among different cell types and among different metabolic states and cell cycle phases within the same cell type (21, 22).

The number of nucleoli per nucleus also differ. The yeast cell contains one relatively large nucleolus with respect to its nuclear volume, and most metazoan cells display one or a few nucleoli. At the other extreme, *Xenopus* oocytes contain over 1000 nucleoli per nucleus (germinal vesicle) because of selective ribosomal DNA **amplification** in the very early stages of meiosis I (pachytene). These

multiple nucleoli are independent of the lampbrush chromosomes within these oocytes (see Fig. 1), a unique feature that allows exclusive study of nucleoli without interference from nonnucleolar chromatin. These nucleoli produce an incredible 300,000 ribosomes per second during mid-oogenesis in *Xenopus*. The number of nucleoli per nucleus also varies during the cell cycle. For example, cells in the early interphase may show several small nucleoli that reflect the number of NOR unique to the species. As the interphase progresses, the small nucleoli may fuse to reduce the number of nucleoli.

Nucleoli originate from and are intimately associated with the [nucleolar organizer](#) regions (NOR) that appear as secondary constrictions within certain (“nucleolar”) mitotic chromosomes(20). The NOR are genetic loci that contain repeated 47S ribosomal RNA genes in tandem array. Each gene within the array is separated by intergenic spacer regions. The repeated ribosomal RNA genes each encode (in order) an external transcribed spacer, the 18S region, the first internal transcribed spacer, the 5.8S region, the second internal transcribed spacer, and finally the 28S region. Processing of the nascent transcript yields mature 18S, 5.8S, and 28S rRNA found in cytoplasmic ribosomes. The precise site of transcription has been controversial, although fibrillar centers (FC) of interphase nucleoli clearly contain ribosomal DNA. Good arguments place transcription within the FC, within the DFR, or on the borders between the two compartments. Nascent transcript processing (cleavage, folding, and base modification) and early ribosomal subunit assembly clearly take place within the DFR. Processing events include cleavage of the nascent transcript into smaller intermediates and then into the individual mature rRNA molecules, site-specific methylation of the pre-rRNA, and conversion of specific uridine residues to pseudouridine (see [Ribosomes](#)). Processing of the 47S transcript into the 18S rRNA (for the small ribosomal subunit) and the 5.8S and 28S rRNA (for the large subunit) occurs concomitantly with the association of ribosomal proteins that enter the nucleolus from their site of synthesis in the cytoplasm (there are approximately 85 different ribosomal proteins).

The 5S rRNA is transcribed by RNA polymerase III from multiple tandem genes that exist in genomic clusters outside the nucleolus. How the 5S transcripts make their way to the nucleolus and associate with the large ribosomal subunit remains uncertain. Assembly of large and small ribosomal subunits is a vectorial process beginning at the site of 47S RNA synthesis. By the time the particles reach the GR, the ribosomal subunits are nearly complete. Their assembly is completed once the subunits translocate to the cytoplasm, where a few more ribosomal proteins join the complexes. The large and small subunits finally associate as a complete ribosome in the recognition and translation of mRNA in the cytoplasm (see [Translation](#)).

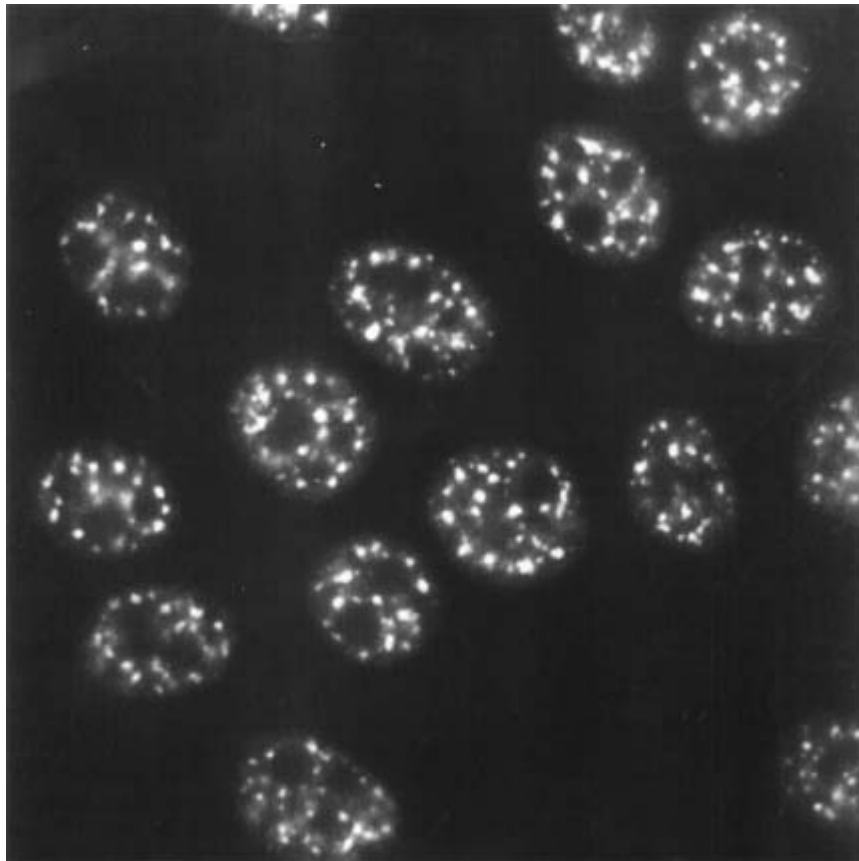
#### 3.4. Nuclear Distributions of Splicing Components

Antibodies against spliceosome components have been used in **immunofluorescence microscopy** to localize the sites of splicing and the nascent pre-mRNA within the nuclei of various cell types. The amphibian oocyte that has large lampbrush chromosomes has been particularly informative (25, 26). Antibodies directed against snRNPs or the splicing factor SC35 stain the lateral loops of giant lampbrush chromosomes within *Notophthalmus viridescens* (newt) oocytes (Fig. 1). The chromosomal loops are sites of active transcription, and their staining with antibodies directed against splicing components strongly argues that initial splicing events occur while the transcripts are still attached to RNA polymerase II. Beautiful electron micrographs of active *Drosophila* **chorion genes** clearly show the processing of nascent transcripts while still associated with RNA polymerase II(27).

Antibodies directed against splicing components brightly stain a limited number of discrete foci above a lightly and uniformly stained background in interphase nuclei. For example, anti-Sm antibodies stain 20 to 50 foci in the interphase nuclei of CHO 400 cells (Fig. 2) (24, 28). These speckles, as they are called, contain interchromatin granule clusters (IGC) interconnected by perichromatin fibrils (PF) to form a nuclear latticework. Functionally, the IGC may be centers for snRNP assembly and distribution, whereas the PF are sites of pre-mRNA synthesis and splicing as the transcripts make their way to the nuclear envelope. This latticework of clusters and

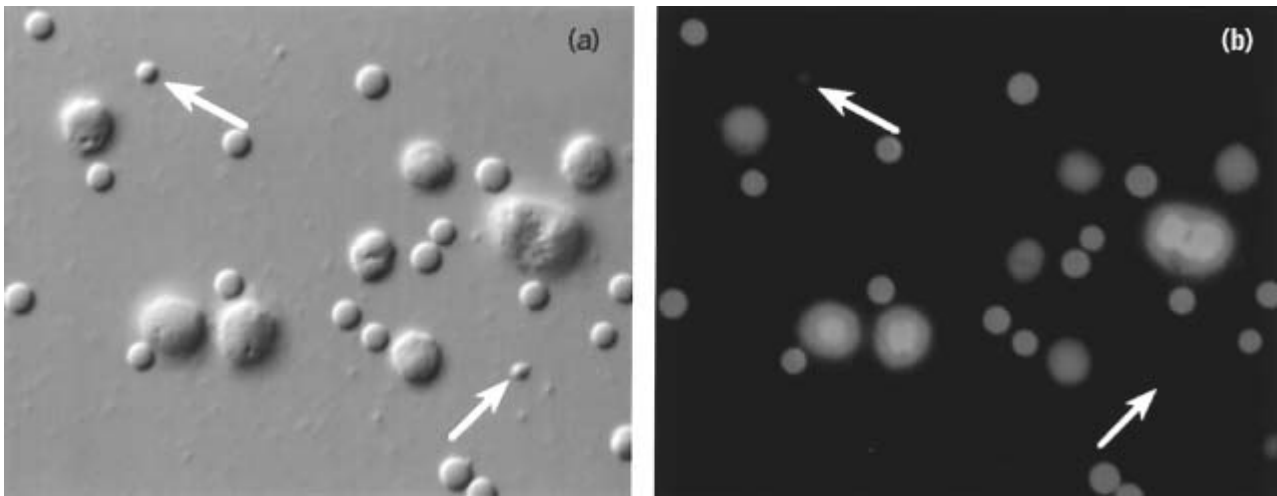
interconnecting fibrils is a dynamic structure that changes its morphology in response to changes in cell and nuclear physiology. In addition, the speckles themselves disassemble during mitosis and reassemble during the interphase. The existence and distinct distribution of the speckles within the interphase nucleoplasm clearly emphasizes the structural and functional compartmentalization of the interphase nucleus (Fig. 2).

**Figure 2.** Factors involved in pre-mRNA splicing are organized in a speckled distribution pattern in the interphase nuclei of mammalian cells. The speckles are composed of both interchromatin granule clusters (IGC) and perichromatin fibrils (PF). The PF represent sites of active transcription, and the IGC are thought to be storage and/or reassembly sites for splicing factors. Photomicrograph was kindly provided by Dr. David L. Spector, Cold Spring Harbor Laboratory.



Using antibodies directed against the Sm antigens and the unique trimethylguanosine cap of snRNA in conjunction with *in situ* hybridizations to detect snRNA, Wu et al. (26) described three types of snRNP-containing granules (A, B, and C types), which are called “snurposomes” (29). Although the antibodies also stain the lampbrush chromosomes, for the most part these granules are independent of the chromosomes. The A granules contain exclusively U1 snRNP, whereas the B granules (Fig. 3) contain all five snRNPs (26). The C granules are most interesting. They are the originally described sphere organelles (30) found near or at the histone gene clusters on the large lampbrush chromosomes of at least two amphibian species (31-33). The C granules contain the U7 snRNP (34) reserved for histone transcript processing (35). Interestingly, B particles often reside on the surface of C granules or within C granules as [inclusion bodies](#), but the functional significance of this relationship remains unknown.

**Figure 3.** Contents from a *Xenopus laevis* oocyte nucleus spread for light microscopic examination. **(a)** DIC image. The multiple nucleoli were stained with propidium iodide. Arrows show micronucleoli. **(b)** Snurposomes were stained with monoclonal antibody Y12 which specifically binds the Sm epitope of several snRNP proteins. See the text for descriptions of the nucleoli, B snurposomes, snRNP proteins, and Sm epitopes. Photomicrographs were kindly provided by Dr. Joe Gall, Carnegie Institution of Washington, Department of Embryology.



### 3.5. Coiled Bodies

Coiled bodies were first discovered by Ramón y Cajal in 1903(36) as nucleolar accessory bodies in neuronal cell nuclei, but their function remains uncertain. They were described in detail as internal coiled fibers at the ultrastructural level by Monneron and Bernhard(37), who first used the name coiled bodies. The bodies are spherical, have diameters of 0.5  $\mu\text{m}$  to 1  $\mu\text{m}$ , and are found in both plant and animal cell nuclei. One to five coiled bodies exist per nucleus. The molecular constitution of coiled bodies has been examined only recently (24, 38-40). The protein p80-coilin is highly enriched in the coiled body and is considered a marker for the coiled body, although it has also been detected throughout the nucleus in relatively low concentrations (38, 41, 42). Its function is unknown.

Evidence suggests that the coiled body is related to the C snurposome (see previous): (1) Sph-1, a resident protein of the sphere organelle, is structurally homologous to p80-coilin(43); (2) antibodies directed against human p80-coilin stain the C granules intensely; (3) when mRNA-encoding, epitope-tagged, human p80-coilin was injected into *Xenopus* oocytes, the protein quickly localized to the C granules. Therefore coiled bodies and amphibian sphere organelles are likely to be similar in function. As in the case of amphibian C granule/sphere organelles, antibodies against Sm antigens and trimethylguanosine cap structures stain coiled bodies and, like the sphere organelles, coiled bodies in human and mouse culture cells contain high concentrations of U7 snRNA(44). Significant differences exist between coil bodies and the amphibian sphere organelles, however. First, unlike the sphere organelles, where only U7 snRNP has been detected, all five mRNA-splicing snRNAs (U1, U2, U4, U5, U6) have been detected within coiled bodies by *in situ* hybridization. Parenthetically, the coiled bodies are not the speckles described before because the splicing-dependent SC35 protein has not been observed within coiled bodies. Secondly, whereas no nucleolar components have been found in amphibian sphere organelles, fibrillarin and snRNA U3 (components intimately associated with nucleoli and preribosomal RNA processing, see [Nucleolus](#)) have been detected in coiled bodies of cultured cells (41, 45). As originally reported by Ramón y Cajal(36), coiled bodies closely associate with nucleoli(46), but the physiological significance of this association remains uncertain.

Novel nuclear particles, referred to as Gemini or coiled bodies (Gems), were recently identified in HeLa cell nuclei using monoclonal antibodies against the spinal muscular atrophy (*SMN*) gene product(47). The SMN protein of 32 kDa binds several hnRNP proteins and the nucleolar protein,

fibrillarin, to form nuclear particles that closely resemble and associate with coiled bodies. Although closely related to coiled bodies in size and composition, Gems may be more closely related to the amphibian sphere organelles (C snurposomes). Their function remains unknown.

#### 4. Mitotic Regulation of Nuclear Disassembly and Reassembly

For cells to divide, the nucleus must disassemble, and the genome must be divided exactly in two. Impressive intracellular reorganizations occur with the onset of **mitosis**. Cytoplasmic [microtubules](#) disassemble into tubulin dimers that then reassemble to form the **mitotic spindle**. The [endoplasmic reticulum](#) and the [Golgi apparatus](#) vesiculate. The nuclear envelope and nucleoli disassemble, and chromatin supercondenses. The onset of mitosis is driven by the activity of [maturation promoting factor](#) (MPF), a complex between cyclin B and p34<sup>cdc2</sup> kinase(48). MPF in turn is regulated by several factors and [signal transduction](#) events, one of which is the completion of genomic DNA synthesis and any requisite [DNA repair](#). The driving mechanism behind these morphological changes is **phosphorylation** of cellular constituent proteins directly or indirectly by MPF (see **Mitosis**).

#### Bibliography

1. J. R. Baker (1949) *Quart. J. Micro. Sci.* **90**, 87–108.
2. R. Brown (1833) *Trans. Linn. Soc.* **16**, 710–711.
3. O. Hertwig (1895) *The Cell, Outlines of General Anatomy and Physiology*, Translated by M. Cambell (H. J. Stone, ed.), Swan Sonnenschein, London; MacMillan, New York.
4. E. B. Wilson (1925) *The Cell in Development and Heredity*, 3rd ed., Macmillan, New York.
5. E. H. Davidson (1986) *Gene Activity in Early Development*, 3rd ed., Academic Press, Orlando, pp. 525–540.
6. P. L. Deininger and G. R. Daniels (1988) *Trends Genet.* **2**, 76–80.
7. D. W. Fawcett (1981) *The Cell*, 2nd ed., W. B. Saunders, Philadelphia.
8. H. F. Willard (1996) *Cell* **86**, 5–7.
9. E. H. Blackburn (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 557–576.
10. J. Loidl (1990) *Genome* **33**, 759–770.
11. A. F. Dernburg, J. W. Sedat, W. Z. Cande, and H. W. Bass (1995) In *Telomeres* (E. H. Blackburn and C. W. Greider, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 295–338.
12. T. J. Yen and B. Schaar (1996) *Curr. Opin Cell Biol.* **8**, 381–388.
13. O. L. Miller Jr. and B. R. Beatty (1969) *Science* **164**, 955–957.
14. J. G. Gall, M. O. Diaz, E. C. Stephenson, and K. A. Mahon (1983) In *Gene Structure and Regulation in Development* (S. Subtelny and F. Fakatos, eds.), Alan R. Liss, New York, pp. 137–146.
15. M. Zillman, M. A. Gorovsky, and E. M. Phizicky (1991) *Mol. Cell. Biol.* **11**, 5410–5416.
16. M. Hochstrasser, and J. W. Sedat (1987) *J. Cell Biol.* **104**, 1471–1483.
17. T. Cremer, A. Kurz, R. Zirbel, S. Dietzel, B. Rinke, E. Schrock, M. R. Speicher, U. Mathieu, A. Jauch, P. Emmerich, H. Schertan, T. Ried, C. Cremer, and P. Lichter (1993) *Cold Spring Harbor Symp. Quant. Biol.* **58**, 777–792.
18. W. F. Marshall, A. F. Dernburg, B. Harmon, D. A. Agard, and J. W. Sedat (1996) *Mol. Biol. Cell* **7**, 825–842.
19. G. Blobel (1985) *Proc. Natl. Acad. Sci. USA* **82**, 8527–8529.
20. B. McClintock (1934) *Z. Zallforsch. Mikrosk. Anat.* **21**, 294–328.
21. H. Busch and K. Smetana (1970) *The Nucleolus*, Academic Press, New York.

22. A. Hadjiolov (1985) *The Nucleolus and Ribosome Biogenesis*, Springer Verlag, Wien.
23. U. Scheer and M. Dabauvalle (1985) In *Developmental Biology: A Comprehensive Synthesis*, Vol. 1 (L. W. Browder, ed.), Plenum Press, New York, pp. 385–430.
24. D. L. Spector (1993) *Ann. Rev. Cell Biol.* **9**, 265–315.
25. H. G. Callan and J. G. Gall (1991) *Chromosoma* **101**, 69–82.
26. Z. Wu, C. Murphy, H. G. Callan, and J. G. Gall (1991) *J. Cell Biol.* **113**, 465–483.
27. A. L. Beyer and Y. N. Osheim (1990) In *The Eukaryotic Nucleus, Molecular Biochemistry and Macromolecular Assemblies*, Vol. 2. (P. R. Strauss and S. H. Wilson, eds.), Telford Press, Caldwell, NJ, pp. 431–451.
28. D. L. Spector (1990) *Proc. Natl. Acad. Sci. USA* **87**, 147–151.
29. J. G. Gall (1991) *Science* **252**, 1499–1500.
30. H. G. Callan (1986) *Lampbrush Chromosomes*, Springer, Berlin.
31. J. G. Gall, E. C. Stephenson, H. P. Erba, M. O. Diaz, and G. Barsacchi-Pilone (1981) *Chromosoma* **84**, 159–171.
32. H. G. Callan, J. G. Gall, and C. Murphy (1991) *Chromosoma* **101**, 245–251.
33. J. G. Gall, A. Tsvetkov, Z. Wu, and C. Murphy (1995) *Dev. Genet.* **16**, 25–35.
34. C-H. H. Wu and J. G. Gall (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6257–6259.
35. W. F. Marzluff (1992) *Gene Expression* **2**, 93–97.
36. S. R. Ramón y Cajal (1903) *Trab. Lab. Invest. Biol.* **2**, 129–221.
37. A. Monneron and W. Bernhard (1969) *J. Ultrastruct. Res.* **27**, 266–288.
38. I. Raska, R. L. Ochs, L. E. C. Andrade, E. K. L. Chan, R. Burlingame, C. Peebles, D. Gruol, and E. M. Tan (1990) *J. Struct. Biol.* **104**, 120–127.
39. K. Brasch and R. L. Ochs (1992) *Exp. Cell Res.* **202**, 211–223.
40. A. I. Lamond and M. Carmo-Fonseca (1993) *Trends Cell Biol.* **3**, 198–204.
41. I. Raska, L. E. C. Andrade, R. L. Ochs, E. K. L. Chan, C-M. Chang, G. Roos, and E. M. Tan (1991) *Exp. Cell Res.* **195**, 27–37.
42. L. E. C. Andrade, E. K. L. Chan, I. Raska, C. L. Peebles, G. Roos, and E. M. Tan (1991) *J. Exp. Med.* **173**, 1407–1419.
43. R. Tuma, J. A. Stolk, and M. B. Roth (1993) *J. Cell Biol.* **122**, 767–773.
44. M. R. Frey and A. G. Matera (1995) *Proc. Natl. Acad. Sci. USA* **92**, 5915–5919.
45. L. F. Jiménez-García, M. de L. Segura-Valdez, R. L. Ochs, L. I. Rothblum, R. Hannan, and D. L. Spector (1994) *Mol. Biol. Cell* **5**, 955–966.
46. R. L. Ochs, T. W. Stein Jr., and E. M. Tan (1994) *J. Cell Sci.* **107**, 385–399.
47. Q. Liu and G. Dreyfuss (1996) *EMBO J.* **15**, 3555–3565.
48. S. A. MacNeill and P. A. Fantes (1995) In *Cell Cycle Control* (C. Hutchison and D. M. Glover, eds.), IRL Press, Oxford, pp. 63–105.

### **Suggestions for Further Reading**

49. G. Almouzni and A. P. Wolffe (1993) Nuclear assembly, structure, and function: The use of *Xenopus in vitro* systems, *Exp. Cell Res.* **205**, 1–15.
50. R. C. Bird, G. S. Stein, J. B. Lian, and J. L. Stein, eds. (1996) *Nuclear Structure and Gene Expression*, Academic Press, New York.
51. H. Busch (1974–1984) *The Cell Nucleus*, Vols 1–10, Academic Press, New York.
52. G. Dreyfuss, M. Hentze, and A. I. Lamond (1996) From transcript to protein. *Cell* **85**, 963–972.
53. J. G. Gall (1992) Organelle assembly and function in the amphibian germinal vesicle, *Adv. Develop. Biochem.* **1**, 1–29.
54. U. K. Laemmli and R. Tjian (1995) Nucleus and gene expression, a nuclear traffic jam:

- Unraveling multicomponent machines and compartments, *Curr. Opin Cell Biol.* **8**, 299–302.
55. J. Newport and D. Forbes (1987) The nucleus: Structure, function and dynamics, *Ann. Rev. Biochem.* **56**, 535–565.
56. P. R. Strauss and S. H. Wilson (1990) "The" Eukaryotic Nucleus, *Molecular Biochemistry and Macromolecular Assemblies*, Vols. **1** and **2**, Telford Press, Caldwell, NJ.
57. J. Strouboulis and A. P. Wolffe (1996) Functional compartmentalization of the nucleus, *J. Cell. Sci.* **109**, 1991–2000.
58. K. E. van Holde (1988) *Chromatin*, Springer-Verlag, New York.
59. A. Wolffe (1995) *Chromatin, Structure and Function*, 2nd ed., Academic Press, San Diego.

## Nude Mice

Nude mice were first described in 1963 as hairless animals. The defect is the consequence of a recessive [mutation](#) of the “nu” **gene**, located on chromosome 11. It was subsequently found that these mice were also athymic, which rendered them highly susceptible to infectious diseases and necessitated that they be maintained in a strictly controlled microbiological environment. Because of the absence of a thymus, they are unable to mount **T-cell-dependent** [immune responses](#) and are especially lacking cell-mediated immunity. It is thus possible to transplant to these mice not only **allogenic** tissues, but also tissues of **xenogenic** origin.

T cells classically develop in the thymus, because the thymus provides the necessary epitheliostromal environment. In nude mice, there is no thymus, but T cells develop very slowly and may reach more than 10% of the amount found in normal mice. These T cells may express either the TCPab or the TCRgd [T-cell receptors](#), although at a much lower density. The CD3 signaling module is also present, so the TCR is potentially functional.

[B cells](#) are present in normal amounts; but, as would be anticipated, **isotypes** that require T-cell help for [class switching](#) are absent or extremely reduced. Circulating [IgM](#) is present in normal amounts, but **IgG1**, IgG2a, IgG2b, and [IgA](#) are present only at very low levels. IgG3, which may be produced in response to T-independent **immunogens**, such as polysaccharides, is present in normal amounts. One role of the immune system that is frequently put forward is that of “immune surveillance”—that is, the elimination of potentially occurring tumors. In nude mice, however, there is apparently no impact on a possible increase in the emergence of spontaneous tumors, an observation probably explained by the presence of natural killer (NK) cells in increased numbers over normal mice. Because of the virtual absence of cell-mediated immunity, nude mice have been extensively used as recipients for tumor transplantation or for serial transfer of [hybridomas](#) as a source of [monoclonal antibodies](#).

See also entries **Antigen presentation**, [Immunogen](#), and [T cell](#).

### Suggestion for Further Reading

- T. Hunig (1983) T-cell function and specificity in athymic mice. *Immunol. Today* **4**, 84–87.

## Null Mutation

A null mutation is a [mutation](#) that causes complete loss of gene function. The mutation that causes a null mutant can be of any class, for example, point mutations, deletions, or insertions.

## O-Linked Oligosaccharides

There are three major types of [O-glycosylation](#): (1) Mucin-type O-glycosylation; (2) Glycosylation leading to the formation of proteoglycans; and (3) O-mGlcNAc glycosylation of cytoplasmic and nuclear proteins. O-GlcNAc glycosylation involves only the addition of a single hexosamine, but the other two types of glycosylation, although using only a limited number of sugars, lead to the formation of O-glycans with a wide variety of structures.

### 1. O-glycans Added to Mucin Core Proteins

The sugars commonly found in the O-glycans added to mucins in higher **eukaryotes** are the hexosamines N-acetylgalactosamine (GalNAc) and N-acetylglucosamine (GlcNAc), the monosaccharides galactose (Gal) and fucose (Fuc), and sialic acid (SA), some of these sugars, as well as tyrosine residues in the core protein, may be **sulfated**. Sugars are added individually and sequentially in the [Golgi apparatus](#), and the structure of the final O-glycan is generally thought of as being made up of three elements; (1) a core structure built on the linkage sugar, GalNAc; (2) a chain extending from the core, generally made up of repeating lactosamine units; and (3) a terminal region often containing **epitopes** recognized by naturally occurring **antibodies**:

GalNAc core—lactosamine chains—terminal epitopes (1)

O-glycans may be released from **glycoproteins** by hydrazine or by  $\beta$ -elimination in alkaline borohydride solution of the reduced form containing terminal GalNAcOH residues. After purification, analysis of the released O-glycans has been carried out using **mass spectrometric** methods, methylation chemical analysis, and nuclear magnetic resonance ([NMR](#)) techniques (1). When only a few O-glycans are added to a core protein, their structure is usually simpler than that of the O-glycans found on the mucins, which carry multiple side chains. Analysis of the O-glycans from mucins shows a high degree of structural heterogeneity, but whether the heterogeneity is between or within individual mucin molecules is not clear. These preparations have been derived from tissues that contain different cell phenotypes, which could exhibit different profiles of expression of the glycosyltransferases involved in the synthesis of the O-glycans, resulting in a population of different glycoforms where an individual molecule carries only one kind of structure. On the other hand, the addition of multiple O-glycans during transit through the Golgi apparatus may necessarily result in heterogeneity of both chain initiation and extension. The core proteins are extremely large and, after the addition of the first sugar, must be quite extended. It is easy to imagine how all the sites on the protein and on the extending chain might not be found by the relevant transferase(s).

In addition to a complete analysis of the O-glycan structure, components may be detected with specific antibodies, such as those reacting with the common blood group antigens, which recognize terminal sugars, or those reacting with the branched and unbranched poly-lactosamine chains (I and II)

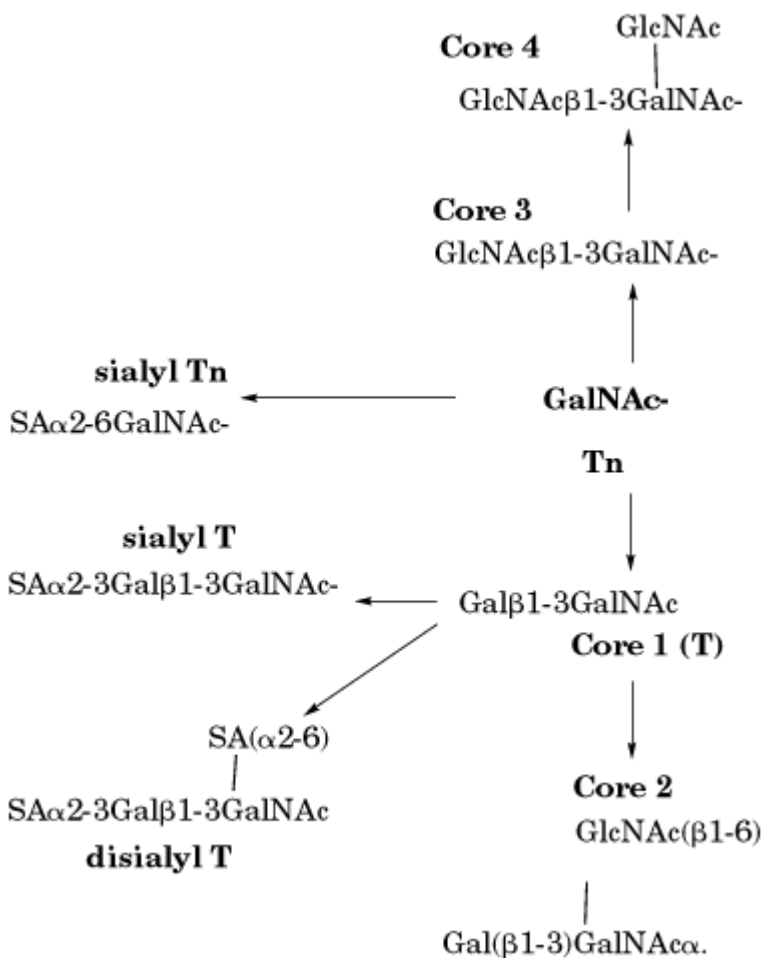


antigens, respectively).

### 1.1. Core Structures

It is possible for an O-glycan to consist only of the initial sugar GalNAc, and this has been referred to as the Tn antigen. More often, core structures are formed by the addition of GlcNAc and/or Gal units, or the chain is terminated by the addition of sialic acid, forming the sialyl Tn epitope that is specifically expressed in carcinomas. Figure 1 shows the main core structures that have been identified to date as being present on mucins. Core 1 and core 2 are by far the most common structures.

**Figure 1.** Main core structures found on mucins.



Core 1, also referred to as the T antigen, may be found as such, particularly on mucins produced by carcinomas, but it is often sialylated in normal cells to give the mono- or disialylated T antigen. The reaction catalyzed by the core1 β3-Gal-T (UDP-Gal:GalNAc-R b 1,3-Gal-transferase) is:

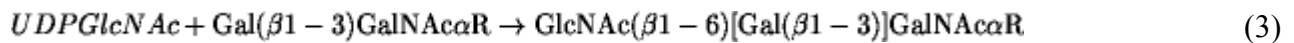


A **gene** coding for the mouse enzyme has only recently been cloned, and it remains to be seen how many enzymes catalyze this reaction, and whether the peptide sequence flanking the glycosylation site affects the specificity.

An important point is that core 1 can be a substrate not only for sialyltransferases that add sialic acid

in a-linkages and thus terminate the chain, but also for the core 2 enzyme, which is crucial for chain extension to proceed. Studies with cells transfected with these enzymes indicate that the two enzymes can compete for the core 1 substrate and therefore overlap in the Golgi apparatus. Differences in the activities of these enzymes, operating at an early stage of O-glycan synthesis, can therefore dramatically affect the final structure of the O-glycan and can explain the differences in the structure of the O-glycans added to MUC1 in breast cancer (1) or to leukosialin in T cell activation (ref 17 of [O-glycosylation.](#))

Core 2 is generated from core 1 by enzymes catalyzing the reaction:



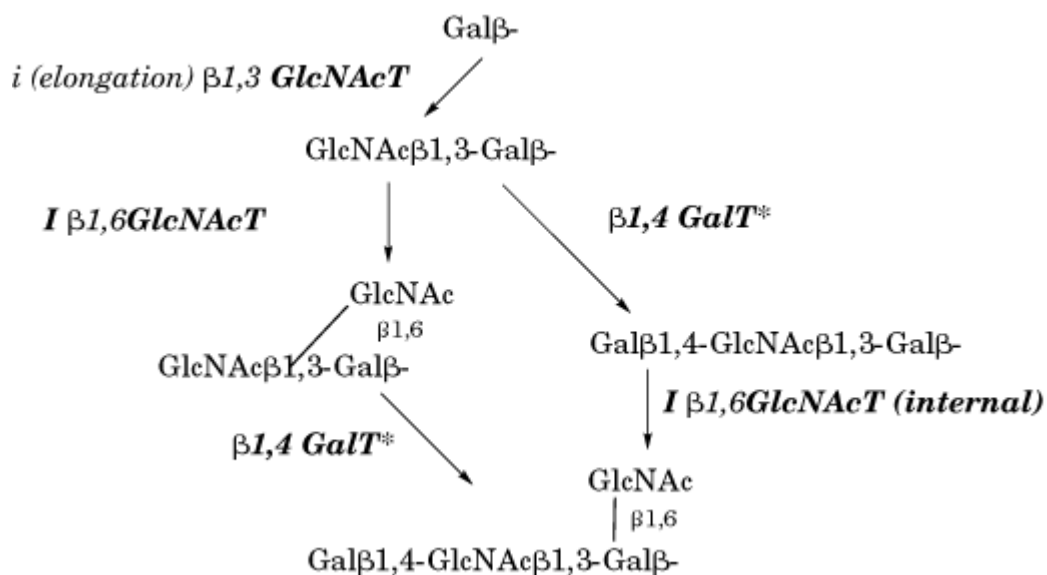
Two enzymes have been shown to catalyze the formation of core 2: the L enzyme, which catalyses only this reaction and for which the gene has been cloned (2), and the M enzyme, which has a wider substrate specificity and can also catalyze the formation of core 4 and internal chain branching. Changes in the level of activity of the enzyme(s) catalyzing the synthesis of core 2 have a profound effect on O-glycan structure in a wide variety of cell types.

Core 3 synthesis precedes Core 4 (see Fig. 1), and the enzymes catalyzing the synthesis of these structures have only been found in mucin-secreting tissues, such as those of the respiratory tract and the colon. Core 4 structures are more predominant than core 3, and because the transfer of GlcNAc to core 4 proceeds at a much faster rate than the b3-GlcNAc transferase making core 3, the activity and tissue distribution of the core 3 b3-GlcNAcT are the limiting factors in the synthesis of core 4. The enzyme catalyzing the synthesis of core 4 shows limited distribution, but it can also catalyze the synthesis of core 2 from core 1.

## 1.2. Backbone of O-Glycans

The elongation of O-glycans involves the addition of GlcNAc residues to Gal in b-1,3 or b-1,6 linkage and of Gal to GlcNAc in b-1,3 or b-1,4 linkage, to form linear and branched poly-N-lactosamine structures, as illustrated in Figure 2. Extension of the O-glycans based on core 2 can occur by the addition of sugars to either the Gal or the GlcNAc moieties. When extension is from core 4, galactose can be added to either glucosamine.

**Figure 2.** Elongation pathways forming the backbone of O-glycans.



(1) The elongation b3-GlcNAc transferase catalyses the addition of GlcNAc in b-1,3 linkage to Gal on core 1 and core 2. (See Fig. 1 for the structures of core 1 and core 2.) The enzyme shows limited distribution and may be reduced in cancers of the colon and breast. Elongation from core 1 prevents the formation of core 2, but it seems likely that core 2 is usually formed before elongation from Gal occurs.

(2) The i b3GlcNAc transferase enzyme adds GlcNAc to Gal in chains that have already been extended from cores. This enzyme is found ubiquitously and catalyzes the addition of GlcNAc alternately with Gal, producing a linear chain that is the blood group i antigen.

(3) The I b6-GlcNAc transferases are responsible for the formation of the blood group I antigen found on adult human erythrocytes, where it replaces the i antigen found in the fetal cells. The I antigen represents the branch initiated by these enzymes either immediately after the action of the b3 GlcNAc transferase or after the addition of lactosamines (see Fig. 2), and such branched O-glycans are found on many mucin-type molecules.

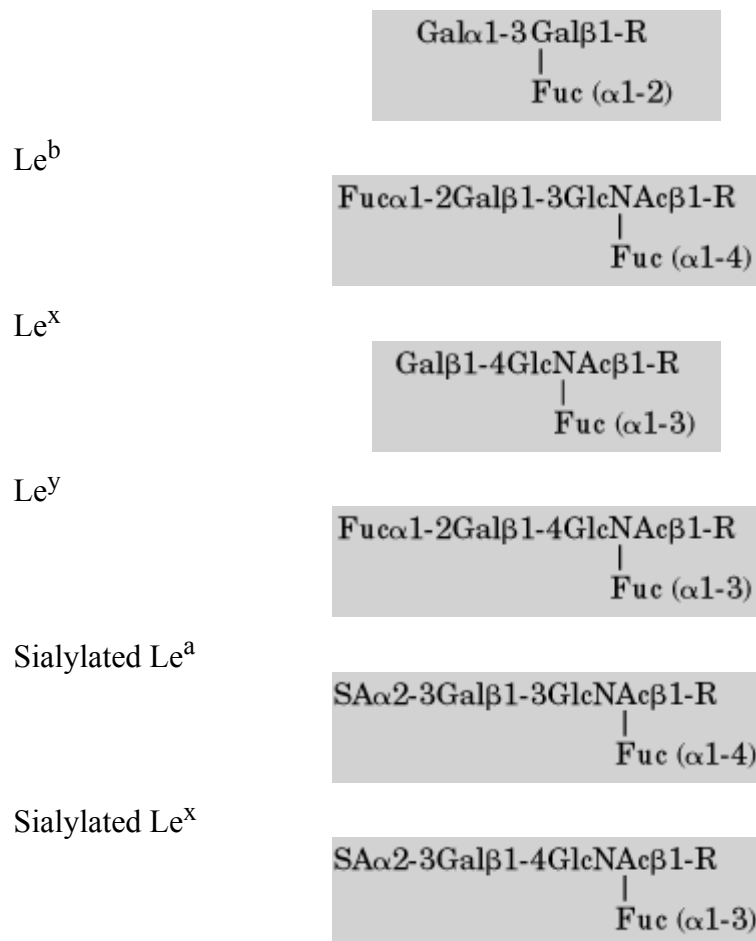
(4) The b-1,4 and b-1,3 galactosyltransferases add galactose from a UDP-Gal donor to GlcNAc in the growing polylactosamine chain in either b-1,4 or b-1,3 linkage. When the linkage is b-1,4, the chain is referred to as type 2, whereas the b-1,3 linkage is type 1. Both types of chain may be linear or branched.

### 1.3. Terminal Glycosylation

The terminal epitopes of the O-glycans on mucins are probably the most important in determining whether the molecule plays a role in cell adhesion phenomena. The epitopes recognized by antibodies related to the ABO and the Lewis blood group antigens are also found in this terminal region. Terminal sugars added in alpha linkage include sialic acid, Fuc, Gal, GalNAc and GlcNAc; Table 1 lists some of the more important structures. Some sulfation of the sugars in terminal structures may also occur.

**Table 1. Terminal Epitopes in Mucin O-glycans**

| Antigen               | Structure  |
|-----------------------|--|
| Blood group A         | $\begin{array}{c} \text{GalNAc}\alpha 1-3\text{Gal}\beta\text{-R} \\   \\ \text{Fuc} (\alpha 1-2) \end{array}$ |
| Blood group B         | $\begin{array}{c} \text{Gal}\alpha 1-3\text{Gal}\beta 1\text{-R} \\   \\ \text{Fuc} (\alpha 1-2) \end{array}$  |
| H (masked by A and B) | $\text{Fuc}\alpha 1-2\text{Gal}\beta 1\text{-R}$   |
| Le <sup>a</sup>       |  |



Sialyltransferases that are clearly specific for O-glycans are those that add sialic acid in a2-6 linkage to GalNAc(Tn) or in a2-3 linkage to core 1 (Galb1-3 GalNAc). O-glycans with an SA in a2-6 link to GalNAc(Tn) cannot be acted on by any known transferase. This link can be formed however, after the addition of SA in a2-3 linkage to Gal. The shorter O-glycans (ie, Tn, T, and their sialylated derivatives), are found on mucins expressed by some carcinomas, and the change from the normal glycosylation pathway has been analyzed in greatest detail in breast cancer (1). Several genes coding for sialyltransferases responsible for terminating the short O-glycans have been cloned (3) and show differing substrate specificities: some can synthesize the same linkage in glycolipids. Because of the multiplicity of glycosyltransferases, it is only by cloning the individual genes, thereby allowing work with the recombinant enzymes, that the biosynthetic pathways will be unambiguously clarified.

The position of glycosyltransferases in the Golgi apparatus also plays a significant role in whether the enzymes can act on a particular substrate. One of the sialyltransferases that add SA in a2-3 linkage to Gal in core 1 has been localized to the medial/trans Golgi stacks, with some found in the trans Golgi network (TGN). This is a relatively early position and allows competition with the core 2 enzyme and possibly the elongating enzymes (see ref. 12 of [O-Glycosylation](#)). It is assumed that the sialyltransferase that adds SA to Gal at the end of the extended chains may be located further down the Golgi pathway in the trans Golgi or TGN, but this has not been clarified. The location probably also relates to whether the same enzyme can add SA to lactosamine chains in both N and O-glycans. The sialylated derivatives of the Le<sup>a</sup> and Le<sup>x</sup> antigens terminating lactosamine chains are proving to be of great interest, as they appear with the change to malignancy in some tissues (4) and constitute the epitopes on the selectin ligands expressed by normal cells.

Fucosyltransferases form a large group of enzymes, and genes for several of these transferases have been isolated (5). They are involved in chain termination and are of particular interest in the synthesis of blood group antigens and in the sialylated derivatives of Le<sup>x</sup> and Le<sup>a</sup> that form the epitopes on selectin ligands. At least four α3-fucosyltransferases exist (α3 Fuc-T-III to-VI) that add fucose in α1–3 linkage to form the Le<sup>x</sup> epitope illustrated in Table 1. Of these, Fuc-T-III has the broadest substrate specificity, because it can act on type 1 or type 2 structures and also add fucose in α1–4 linkage, thus synthesizing several human blood group epitopes (Le<sup>a</sup>, Le<sup>b</sup>, Le<sup>x</sup>, and Le<sup>y</sup> as well as sialyl Le<sup>a</sup> and sialyl Le<sup>x</sup>). When added in this position, fucose terminates the chain and needs to be added after sialic acid to create the sialylated Lewis epitopes.

#### 1.4. Functions of O-Glycans

Extracellular mucins such as MUC2 form large oligomers, and although dimerization and some oligomerization occur within the cell, the interactions that continue after secretion into the mucous layer depend to a large extent on the presence of the O-glycans. Although the detailed structure of the O-glycans may not be so important for this function, it is relevant that the structures found on the mucins produced in the gastrointestinal and respiratory tracts carry large, complex O-glycans, which probably relate to the protective function that is crucial in these tissues (6). In glandular epithelia such as that in the breast, the O-glycans are shorter and simpler but still extended (7, 8). The carbohydrate side chains also serve to bind invading micro-organisms, and heterogeneity in their structure would serve to allow interactions with a variety of **receptors**. Specific interactions of defined structures present in the O-glycans involved in cell-cell adhesion, however, are now becoming clarified and are of great interest.

Selectin ligands are mucin-like glycoproteins that interact with the selectins expressed on endothelial cells, leukocytes, and platelets (ref. 15 of [O-glycosylation](#)). The interaction mediates rolling of leukocytes on blood vessels during inflammation, and, at the molecular level, the O-glycans expressed on the selectin ligand play a major role in determining the specificity of the interaction. The selectins show weak binding of sialylated fucosylated oligosaccharides, such as sialyl Le<sup>x</sup>, but bind much more strongly to glycoproteins carrying these O-glycans, which may be sulfated.

P-selectin glycoprotein ligand-1 (PSGL-1) is expressed on leukocytes and interacts with both E-selectin found on endothelial cells and P-selectin on platelets. The specificity of the ligand interaction has been studied by transfecting [complementary DNAs](#) coding for glycosyltransferases and PSGL-1 into CHO cells, which lack both the core 2 and the α1–3 fucosyltransferases required for chain extension and for the formation of the Le<sup>x</sup> epitope. Only CHO cells expressing both transferases were able to bind to P- and E-selectins (9). Because sialidase treatment also eliminates the binding of PSGL-1, sialyl Le<sup>x</sup> has been identified as the specific epitope required for interaction with the selectins. Binding to P- but not E-selectin also requires sulfation of **tyrosine residues** in the core protein. Clearly, the core protein plays a role in the presentation of the O-glycan, possibly by specifying the clustering and conformation, as well as by providing amino acid residues for sulfation.

The appearance of sialylated Le<sup>a</sup> in mucins expressed in colon cancer cells suggests that the interaction of the carbohydrate may influence the metastatic process by enhancing binding of the cancer cells to endothelial cells (10).

Sialoadhesin is a molecule expressed at high levels by macrophages (11) that interacts specifically with monosialylated core 1 (sialyl T). Its normal function is thought to relate to interactions with leukocytes. Because this structure is overexpressed on mucins expressed by cancer cells, however, the possibility exists of some interaction between the tumor cells in carcinomas and the infiltrating macrophages.

Membrane mucin MUC1 is unusual among the epithelial mucins in that it is a **transmembrane** molecule and, as such, resembles the selectin ligands. Being widely expressed on glandular epithelial

cells from which carcinomas develop, it is highly expressed by these cancers and aberrantly glycosylated, with the O-glycans being based mainly on core 1 rather than on core 2. This makes the molecule antigenically distinct, and both humoral and cellular responses to the cancer-associated mucin have been seen in breast and ovarian cancer patients. In addition to responses being generated by the whole molecule, it is becoming clear that glycopeptides can be presented by [major histocompatibility complex](#) (MHC) molecules. Moreover, peptides carrying larger core 2-based structures are not presented as well as glycopeptides carrying Tn or T ([12](#)). The role of O-glycosylation in the [immune response](#) is of interest in the wider context, as many **antigens** from infectious agents are glycoproteins carrying O-glycans, including [HIV](#), in which the *env* protein carries a sialyl Tn epitope ([13](#)).

## 2. The O-glycans of Proteoglycans

Proteoglycans are proteins that carry glycosaminoglycan side chains (GAGs), which can range from a simple linear chain of sugars to highly charged, sulfated polysaccharides. Proteoglycans can carry one to more than 100 GAGs, which consist of alternating hexosamines and hexuronic acid or galactose units and carrying sulfated substitutions at various positions.

### 2.1. Linkage Regions

The GAGs are linked to the protein core via a four-sugar bridge: glucuronic acid b1-3galactoseb1-3galactoseb1-4xylose-O-serine. The attachment of xylose to the serine residue of the protein core is catalyzed by xylosyltransferase (see [O-Glycosylation](#)). The linkage region is basically the same for most proteoglycans, and in both chondroitin sulfate and heparin sulfate (see text below), a proportion of the xylose has been shown to be phosphorylated ([14](#)). In chondroitin proteoglycans, the galactose residues have also been shown to be sulfated. Skeletal keratan sulfate is unlike the other GAG, as it is O-linked to serine or threonine residues via a GalNAc residue ([15](#)), as in mucin-type O-glycans.

### 2.2. The Glycosaminoglycan Chains

The hexosamines in the GAGs can be either D-glucosamine (GlcN) or D-galactosamine (GalN) and the hexuronic acid either D-glucuronic acid (GlcA) or L-iduronic acid (IdoA). These sugars and galactose are arranged in an alternating unbranched sequence and can carry sulfate substitutions at various positions. Table [2](#) shows the compositions of the common GAGs.

**Table 2. Composition of the Common GAG<sup>a</sup> Chains of Proteoglycans**

| Name                | Hexosamine    | Hexuronic acid                    | Galactose |
|---------------------|---------------|-----------------------------------|-----------|
| Chondroitin sulfate | Galactosamine | Glucuronic acid                   | —         |
| Dermatan sulfate    | Galactosamine | Glucuronic acid and iduronic acid | —         |
| Keratan sulfate     | Glucosamine   | —                                 | Galactose |
| Heparan sulfate     | Glucosamine   | Glucuronic acid and iduronic acid | —         |

<sup>a</sup> Key: GAG, glycosaminoglycan.

The inherent structure of the alternation of two types of monosaccharides would be expected to give simple polysaccharide units. Considerable heterogeneity exists within and among the individual

chains, however, due to modification of the repeating units that is often incomplete. This includes sulfate substitutions at differing positions and **epimerization** of carbon 5 of GlcA to form IdoA. For example the biosynthesis of heparan sulfate and heparin is initiated by the formation of [GlcAb1-4GlcNAc1-4]<sub>n</sub>, which is then N-deacetylated, N-sulfated, and undergoes C5 epimerization of GlcA to yield IdoA. The IdoA units and GlcN can then be O-sulfated. The final product is therefore very diverse, and four distinct HexA and six GlcN units have been identified, allowing 17 different HexA-GlcN and 10 different GlcN-HexA.

Which GAG is attached to which protein depends on the protein core and the cell type. For example, CD44, which can mediate cell adhesion, trafficking, and motility, can be expressed as a chondroitin sulfate or heparan sulfate proteoglycan in a cell-type specific manner (16).

### 2.3. Function of the GAGs

The biological roles of proteoglycans are many and diverse, ranging from simple mechanical support to playing an important part in cellular recognition, adhesion, motility, and proliferation. Most of these effects depend on the binding of macromolecules to the GAGs; in fact, the anticoagulant and antiproliferative activities of heparin and heparan sulfate are mediated by the free GAG. Most biological activities, however, depend on the GAG in association with the protein core of the proteoglycan. Binding of GAGs to proteins can be highly specific, as in the binding of heparin to antithrombin, or, as is often the case, the interaction can be less specific and usually electrostatic in nature. The biological activities expressed by a single GAG can often be attributed to specific carbohydrate structures. For example, the interaction between heparin and the proteinase inhibitor antithrombin is based on the occurrence in the GAG of a sequence of five specific oligosaccharides. This sequence is composed of three GlcN units, one GlcA unit, and one IdoA unit, with O-sulfate groups at various positions. The key feature of this structure is an O-sulfate group on the internal GlcN, which is essential for the high-affinity binding of the proteinase inhibitor. Another specific sequence that has been defined on GAGs is that involved in the interaction between dermatan sulfate and heparin cofactor II. As many of the genes encoding the core proteins carrying GAGs have been or are being cloned, and as the methodology improves for the sequence analysis of the GAGs, it is likely that many interactions involving highly specific recognition/binding sequences in the GAGs will be found.

### 3. Concluding Remarks

The structure of O-glycans, synthesized from relatively few building blocks, is made extremely diverse by both the order of the addition of the sugars and the large number of covalent linkages that are possible. Thus the conformation of the O-glycan differs markedly when the same sugar is added to the same substrate in a different linkage (eg, Le<sup>a</sup> vs Le<sup>x</sup>). These differences in shape can be recognized both by B cells producing specific antibodies and by cell surface receptors such as the selectins. The demonstration of such specificity has generated interest in the role of these interactions in differentiation and in cell-cell adhesion. Moreover, the very large number of genes identified that catalyze the reactions leading to the synthesis of O-glycans emphasizes the importance of their fine structure to the organism. The development of **knockout** mice defective in these genes will go some way toward identifying the crucial functions they control.

### Bibliography

1. K. O. Lloyd, J. M. Burchell, V. Kudryashov, B. W. T. Yin, and J. Taylor-Papadimitriou (1996) *J. Biol. Chem.* **271**, 33325–33334.
2. M. F. A. Bierhuizen and M. Fukuda (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9326–9330.
3. M. Chang, R. Eddy, T. B. Shows, and J. T. Lau (1995) *Glycobiology* **5**, 319–325.
4. K. Zhang, D. Baekstrom, and G. C. Hansson (1994) *Int. J. Cancer* **59**, 823–829.
5. H. Schacter (1994) In *Molecular Glycobiology* (M. Fukuda and O. Hindsgaul, eds.), IRL Press, Oxford, pp. 88–162.

6. N. G. Karlsson, H. Karlsson, and G. C. Hansson (1995) *Glycoconjugate J.* **12**, 69–76.
7. F-G. Hanisch, G. Uhlenbruck, J. Peter-Katalinic, H. Egge, J. Dabrowsli, and U. Dabrowski (1989) *J. Biol. Chem.* **264**, 872–883.
8. F -G Hanisch, J. Peterkatalinic, H. Egge, U. Dabrowski, and G. Uhlenbruck (1990) *Glycoconjugate J.* **7**, 525–543.
9. F. Li, P. P. Wilkins, S. Crawley, J. Weinstein, R. D. Cummings, and R. P. McEver (1996) *J. Biol. Chem.* **271**, 3255–3264.
10. D. Baeckstrom, G. C. Hansson, O. Nilsson, S. Johansson, S. J. Gendler, and L. Lindholm (1991) *J. Biol. Chem.* **266**, 21537–21547.
11. P. R. Crocker, S. Mucklow, V. Bouckson, A. McWilliams, A. C. Willis, S. Gordon, G. Milon, S. Kelm, and P. Bradfield (1994) *EMBO J.* **13**: 4490–4503.
12. L. Galli-Stampino, E. Meinjohanns, K. Frische, M. Meldal, T. Jensen, O. Werdelin, and S. Mouritsen (1997) *Cancer Research* **9**, 12–36.
13. J. E. S. Hansen, H. Clausen, S. L. Hu, J. O. Nielsen, and S. Olofsson (1992) *Arch Virol.* **126**, 11–20.
14. L. A. Fransson, I. Silverberg, and I. Carlstedt (1995) *J. Biol. Chem.* **259**, 1720–1726.
15. V. C. Hascall and R. J. Midura (1989) In *Keratan Sulphate* (H. Greiling and J. E. Scott, eds.), The Biochem. Soc., London.
16. S. Jalkanen and M. Jalkanen (1992) *Blood* **84**, 1802–1811.

### Suggestions for Further Reading

17. I. Brockhausen (1995) "Biosynthesis of O-glycans of the N-acetylgalactosamine- -ser/thr linkage type". In *Glycoproteins* (J. Monteuil, J. F. G. Vliegenhart, and H. Schacter, eds.), Elsevier, The Netherlands, pp. 201–259.
18. L. Kjellen and U. Lindahl (1991) Proteoglycans: structures and interactions. *Ann. Rev. Biochem.* **60**, 443–475.
19. J. Iida, A. M. L. Meijne, J. R. Knutson, L. T. Furcht, and J. B. McCarthy (1996) Cell surface chondroitin sulphate proteoglycans in tumor cell adhesion, motility and invasion. *Seminars in Cancer Biology* **7**, 155–162.

### Ochre Mutation

An ochre mutation is a [nonsense mutation](#) that changes a sense **codon** (one specifying an [amino acid](#)) into the [stop codon](#) UAA, causing premature termination of the [polypeptide chain](#) during [translation](#). The mutation, the codon, and the mutant are all called ochre. Ochre mutations arise by single base changes in the codons for seven amino acids (and in the UGA and UAG stop codons, although these would not be nonsense mutations). In principle, mutations in the anticodons of the [transfer RNAs](#) that read those seven codons could give rise to [ochre suppressors](#), but suppressors are recovered only if a second tRNA exists that reads the codon. In *Escherichia coli*, five ochre suppressors that arise by a single base change have been identified. Unlike [amber mutations](#), ochre mutations are suppressed only by ochre suppressors. Ochre suppressors have effects that are much more deleterious than amber suppressors, reflecting the fact that ochre codons are the most frequent stop codon in *E. coli* and related bacteria. Consequently, most ochre suppressors are weak (inefficient at reading ochre codons) and allow most polypeptide chains to be properly terminated. The term ochre first appears in the literature in 1965 ([1](#)), apparently coined by Sydney Brenner to



complement amber.

## Bibliography

1. S. Brenner and J. R. Beckwith (1965) *J. Molec. Biol.* **13**, 629–637.

## Ochre Suppressor

Ochre suppressors are mutant tRNAs that translate the UAA (ochre) termination codon as a sense codon. Ochre suppressors allow for protein synthesis beyond the translational block resulting in active protein. Due to wobble decoding, most ochre suppressors also decode the UAG (amber) stop codon. Ochre mutations cause protein synthesis to terminate prematurely, resulting in inactive, truncated polypeptides. Hence the term “suppressor;” these mutant tRNAs “suppress” the phenotypes of ochre mutations. These suppressors have been extensively used in prokaryotic genetic studies, and in studies of the translational apparatus and mechanisms. For complete discussions of these and other suppressors, see [Nonsense Suppression](#), [Suppressor tRNA](#), and [Genetic Suppression](#).

## Okazaki Fragments

The short stretches of DNA attached to RNA primers on the *lagging* strand during [DNA replication](#) are called *Okazaki fragments*, after their discoverer (see [Primer](#) and [Leading and Lagging Strands](#)). In other words, the Okazaki fragment is a unit of [discontinuous DNA replication](#) on the lagging strand.

Okazaki fragments were first identified during a course of study on the most recently synthesized, or nascent, DNA molecules in the replication of **bacteriophage** T4 DNA (1). Those molecules **sediment** with [sedimentation coefficients](#) of about 8 S to 10 S in alkaline sucrose gradients, corresponding to chain lengths of 1000 to 2000 residues. To make such an analysis possible, Okazaki and colleagues developed a **pulse-chase labeling** technique, by which much of the nucleotide precursor label could be captured as the nascent DNA. After the chase, with subsequent exposure to high concentrations of unlabeled precursors, the radioactive precursor was found exclusively in high-molecular-weight DNA. This means that the Okazaki fragments have a very short life-time and are connected to a long stretch of DNA in cells. RNA linked to Okazaki fragments can also be identified, although the RNA primers are removed quickly *in vivo* (2).

In a wide variety of bacteriophages and **prokaryote** cells, the size of the nascent DNA is between 1000 and 2000 nucleotides. In eukaryotic cells, however, the Okazaki fragments are between 100 and 200 nucleotides (3). This difference seems to reflect a difference between prokaryotes and eukaryotes in the rate of chain elongation on the leading-strand template. The velocity of the [replication fork](#) movement, which represents the rate of leading-strand DNA synthesis, has been estimated from the size of a replicon and the length of its replication period. In *E. coli*, the replication fork proceeds at 1000 residues per second, and one Okazaki fragment is synthesized every 1 to 2 sec. In eukaryotic cells, the rate of fork movement is slow, 10 to 100 residues per second, but the time required for completion of one Okazaki fragment is the same as in the

prokaryotic cells. The distinction between prokaryotic and eukaryotic replication may be due to a difference in the replication machinery or in the structure of the [chromosomes](#).

At least six proteins are involved in the initiation, elongation, and completion of Okazaki fragments in *E. coli*. Those are primase, [single-strand DNA binding protein](#) (SSB), DNA polymerase III, DNA polymerase I, [ribonuclease H](#) (RNaseH), and [DNA Ligase](#) (see [Discontinuous DNA Replication](#)). Primase synthesizes RNA primer in a manner dependent on DnaB protein, which acts as a major replicative [DNA helicase](#) at the replication fork (4). *In vitro* experiments suggest that the primase is recruited anew from solution for each cycle of Okazaki fragment synthesis and that association of primase with the replication fork occurs via a [protein–protein interaction](#) with the helicase, DnaB (5). Thus, the interaction between primase and DnaB at the replication fork is the primary regulator of the cycle of Okazaki fragment synthesis. SSB is required for the initiation of Okazaki fragment synthesis, as well as for chain elongation by DNA polymerase III holoenzyme. DNA polymerase I, assisted by RNaseH, removes the primers and then fills the resulting gaps, to enable the short fragments to be joined by DNA ligase (6). *E. coli* DNA ligase cannot connect RNA to DNA. This process is very important for cells to maintain a high fidelity of DNA replication, since the primase has no proofreading capacity and the [mismatch repair](#) system is not effective on errors in the primer synthesis.

*E. coli* mutant cells deficient in DNA polymerase I or DNA ligase accumulate large amounts of Okazaki fragments. If Okazaki fragments are not connected efficiently, a daughter DNA molecule synthesized on the lagging strand will suffer from double-strand DNA breaks upon the next round of DNA replication. Thus, null mutants of the genes encoding DNA polymerase I (*polA*) or DNA ligase (*lig*) cannot be isolated. Repair of the double-strand DNA breaks involves functions of **RecA** protein in *E. coli*. Temperature-sensitive *polA* or *lig* mutant cells show a lethal **phenotype** when a *recA*<sup>−</sup> mutation is introduced into the strain.

#### Bibliography

1. R. Okazaki, T. Okazaki, S. Hirose, A. Sugino, and T. Ogawa (1976) In *DNA Synthesis and Its Regulation* (M. Goulian and P. Hanawalt, eds.), Benjamin Cummings, Menlo Park, CA, pp. 832–862.
2. T. Ogawa, S. Hirose, T. Okazaki, and R. Okazaki (1977) *J. Mol. Biol.* **112**, 121–140.
3. A. B. Blumenthal and E. J. Clark (1977) *Cell* **12**, 183–189.
4. K. Arai and A. Kornberg (1981) *J. Biol. Chem.* **256**, 5267–5272.
5. K. Tougu and K. J. Marians (1996) *J. Biol. Chem.* **271**, 21398–21405.
6. B. E. Funnell, T. A. Baker, and A. Kornberg (1986) *J. Biol. Chem.* **261**, 5616–5624.

#### Suggestions for Further Reading

7. K. J. Marians (1992) Prokaryotic DNA replication. *Ann. Rev. Biochem.* **61**, 673–720.
8. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.
9. T. Ogawa and T. Okazaki (1980) Discontinuous DNA replication. *Ann. Rev. Biochem.* **49**, 421–457.

#### Oligo (dT) Cellulose

Oligo (dT) cellulose is routinely used to enrich from a total **RNA** population for RNAs that have

[poly A](#) tails, specifically, [messenger RNA](#) (1). The oligo (dT), which can be of various lengths, is attached covalently to the insoluble matrix, cellulose. This is then incubated with RNA under conditions where the poly A tail of the mRNA will hybridize to the oligo (dT). The unbound material can then be removed from the oligo (dT) cellulose by repeated rounds of washing, and the bound polyadenylated RNA eluted for further analysis, such as hybridization to a specific DNA probe in a **Northern blot**.

An alternative use of oligo (dT) is to enrich for cellular proteins that show a preference for binding to this homopolymer. A similar approach has been used with homoribopolymers in characterizing the binding characteristics of the heterogeneous **ribonucleoproteins** (2).

#### Bibliography

1. T. Maniatis, E. F. Fritsch, and J. Sambrook (1989) *Molecular Cloning: a Laboratory Manual*, Cold Spring Harbour Laboratory Press, Cold Spring Harbour, New York, pp 7.26–7.29.
2. M. Swanson and G. Dreyfuss (1988) *Mol. Cell. Biol.* **8**, 2237–2241.

#### Suggestion for Further Reading

3. T. Maniatis, E. F. Fritsch, and J. Sambrook (1989) *Molecular Cloning: a Laboratory Manual*, Cold Spring Harbour Laboratory Press, Cold Spring Harbour, New York, pp. 7.26–7.29.

## Oligomer

An oligomer is the general term for a short [polymer](#), consisting of a small number of linked, repeating monomer units. The cutoff between an oligomer and a polymer is indistinct, but an oligomer generally has between 2 and 50 repeating units. An oligomer of two repeating units may be more specifically referred to as a *dimer*, three repeating units as a *trimer*, four repeating units as a *tetramer*, and so on. Like polymers, oligomers can be synthetic or naturally occurring. For example, synthetic oligomers formed from a small number of **nucleotide** repeating units—**oligonucleotides**—are used for the polymerase chain reaction (**PCR**) in molecular biology. Many [hormones](#) and [toxins](#) are **oligopeptides**—that is, naturally occurring oligomers of [amino acids](#).

[See also [Polymer](#).]

#### Suggestion for Further Reading

P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

## Oligomeric Proteins

[Proteins](#) that contain more than one [polypeptide chain](#) are said to be *oligomeric* or *multimeric*, as opposed to *monomeric* proteins with a single chain. The Greek “oligoi” means “few” and the Latin “multi” means “many,” yet oligomeric and multimeric have become essentially equivalent. The first

term, which was introduced by J. Monod and collaborators in the context of allostery (1), will be used here for proteins with two to a few tens of polypeptide chains, as opposed to large assemblies such as muscle fibers, which have hundreds or thousands. The individual polypeptide chains in an oligomer are also called *subunits*. The terms *homo-oligomer* and *homo-multimer* apply to oligomeric proteins built from only one type of polypeptide chain, the product of the same gene. Other oligomeric proteins are *hetero-oligomers* or *hetero-multimers*.

Whereas most proteins have a [primary structure](#) (amino acid sequence), **secondary structure** ([alpha-helices](#) and [beta-sheets](#)), or a **tertiary** structure (three-dimensional), oligomers have an additional level called the *quaternary structure*. The quaternary structure of a homo-oligomer may be coded by the formula  $a_n$ , where  $n$  is the number of copies of the polypeptide chain. Many proteins are oligomers, and quaternary structures of all sorts are found in nature. Beyond formulae such as  $a_n$  that describe the number and nature of polypeptide chains, the quaternary structure is part of the three-dimensional structure and is determined along with it in [X-ray crystallography](#) or [NMR](#) studies.

## 1. Symmetry

The presence of multiple copies of the same chemical unit leads to the possibility that the three-dimensional structure has internal molecular symmetry. Molecular symmetry is easily established in crystal studies, where it can be either included in the symmetry of the crystal lattice or simply local. In [electron microscopy](#), molecular symmetry is the basis for powerful image reconstruction methods (see [Single Particle Reconstruction](#)), and it can also be used to refine X-ray crystallography data (see [Molecular Averaging](#)). X-ray data indicate that symmetry is the rule in homo-oligomeric proteins, the lack of symmetry being the exception (2). An object with internal symmetry can always be divided into smaller identical units, the smallest of which is called the [asymmetric unit](#) by crystallographers. Monod and collaborators coined the word *protomer* to designate the same entity in an oligomeric protein, and it shall be used here to distinguish the asymmetric unit of the protein from that of the crystal, which may contain more than one protomer. The protomer itself may be more than one polypeptide chain. In hetero-oligomers, it must contain each type of chains; for instance, an [immunoglobulin G](#) molecule has one heavy and one light chain in its protomer. In the discussion below, we assume for simplicity that the protein is a homo-oligomer and the protomer a single polypeptide chain, but the same conclusions are readily extended to hetero-oligomers.

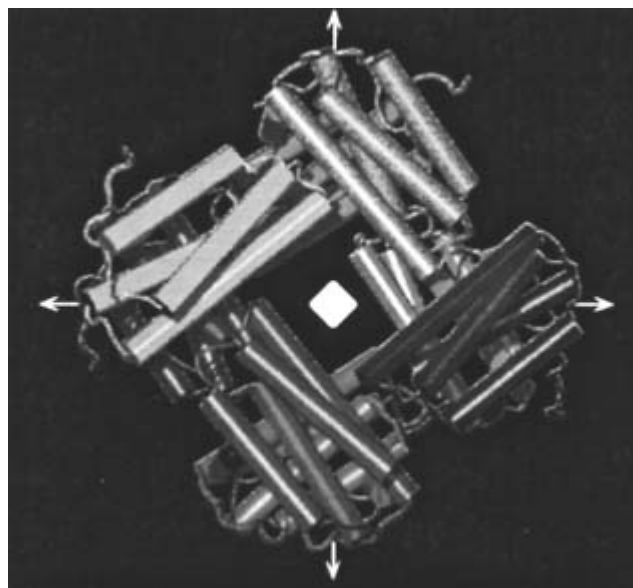
Proteins are **chiral** objects that may not have inversion centers or mirror symmetry, both of which would invert the chirality. Protomers must be related to another by a rotation, translation, or screw rotation, which is a combination of a rotation and translation. When repeatedly applied to the protomer, these operations reconstitute the whole object. Their combination constitutes a group in the mathematical sense of the word. Groups of symmetry operations that generate objects of finite size are known as *point groups*. The number of asymmetric units is a characteristic of the group called its *multiplicity*. Chiral point groups belong to one of three families:

1. *Cyclic*  $C_n$  symmetry (also noted as  $n$ ). Protomers are related by rotations of  $360/n$  degrees about a single axis  $c$ . The oligomer must therefore contain  $n$  protomers and the multiplicity is  $m = n$ .
2. *Dihedral*  $D_n$  symmetry (also noted as  $n2$ ). Protomers are related either by rotations of  $360/n$  degrees about axis  $c$ , or by a  $180^\circ$  rotation about one of  $n$  two-fold axes in the plane orthogonal to  $c$ . The oligomer then contains  $m = 2n$  protomers.
3. The *cubic symmetries* of a tetrahedron, an octahedron, or an icosahedron. All have multiplicities of 12 or a multiple of 12. Tetrahedral symmetry has nonorthogonal two-fold and three-fold axes; in addition, octahedral symmetry has four-fold axes, icosahedral symmetry five-fold axes.

The symmetry of an oligomeric protein is closely related to the number of protomers and, therefore, polypeptide chains (Table 1). The only possible point group symmetry for a homo-dimer is  $C_2$ ,

which has a single two-fold axis and  $m = 2$ . A homo-trimer must have a three-fold axis ( $120^\circ$  rotation) and cyclic  $C_3$  symmetry ( $m = 3$ ), if it is symmetric at all. On the other hand, a homo-tetramer can have two symmetries: either a four-fold axis in the cyclic point group  $C_4$ , or three orthogonal two-fold axes in the dihedral point group  $D_2$  (also noted as 222). Both point groups have  $m = 4$ , yet they yield very different quaternary structures, and  $D_2$  is much more frequently observed than  $C_4$ . Dihedral symmetry, which requires the number of subunits to be even ( $m = 2n$ ), is very common in globular soluble proteins. Homo-hexamers generally have  $D_3$  symmetry. An example is *Escherichia coli* [aspartate transcarbamoylase](#), in which each of the six protomers comprises one catalytic and one regulatory chain. Octamers have  $D_4$  symmetry, illustrated by *hemerythrin* in Figure 1. In contrast, [membrane proteins](#) often have cyclic symmetry and odd numbers of subunits. **Porins** of the bacterial outer membrane and the [bacteriorhodopsin](#) of *Halobacterium halobium* are homo-trimers with  $C_3$  symmetry, whereas the *a-hemolysin* of *Staphylococcus* is a heptamer with  $C_7$  symmetry. Cubic symmetry is less common, but it is found in [ferritin](#) and large assemblies, such as the *pyruvate decarboxylase* complex or icosahedral **viruses**. The latter have capsids made of an assembly of one or several different polypeptide chains, all present in multiples of 60, the multiplicity of the **icosahedral symmetry** group.

**Figure 1.** Octameric hemerythrin. Hemerythrin is a nonheme iron protein that transports oxygen in sipunculid worms. Sketch of the crystal structure (11) with subunits in different shades of gray and  $\alpha$ -helices as cylinders. The assembly of eight identical chains has exact dihedral  $D_4$  symmetry (also called 42 symmetry). Here it is viewed along the four-fold axis (diamond); two-fold axes (arrows) run horizontal and vertical along the diagonals of the square-shaped tetramers.



**Table 1. Symmetry of Oligomeric Proteins<sup>a</sup>**

| Protein       | Quaternary Structure Point Group |
|---------------|----------------------------------|
| <i>Cyclic</i> |                                  |

|   |                   |             |
|---|-------------------|-------------|
| HIV <a href="#">proteinase</a>              | $a_2$             | $C_2$       |
| <a href="#">Hemoglobin</a> <sup>b</sup>     | $a_2b_2$          | $C_2$       |
| <a href="#">Porin</a>                       | $a_3$             | $C_3$       |
| Neuraminidase (flu virus)                   | $a_4$             | $C_4$       |
| Pentraxins                                  | $a_5$             | $C_5$       |
| $\alpha$ -Hemolysin                         | $a_7$             | $C_7$       |
| Light harvesting complex II                 | $a_9b_9$          | $C_9$       |
| <i>Dihedral</i>                             |                   |             |
| <b>Phosphofructokinase</b>                  | $a_4$             | $D_2$       |
| <a href="#">Aspartate transcarbamoylase</a> | $a_6b_6$          | $D_3$       |
| Hemerythrin                                 | $a_8$             | $D_4$       |
| GTP cyclohydroase                           | $a_{10}$          | $D_5$       |
| GroEL <a href="#">chaperonin</a>            | $a_{14}$          | $D_7$       |
| <i>Cubic</i>                                |                   |             |
| Phaseolin                                   | $a_{12}$          | Tetrahedral |
| Apo ferritin                                | $a_{24}$          | Octahedral  |
| <b>Virus</b> coats                          | $a_{60}, a_{180}$ | Icosahedral |

<sup>a</sup> Adapted from ref. 2, where references to X-ray structures can be found.

<sup>b</sup> Hemoglobin also has approximate  $D_2$  symmetry.

## 2. Approximate Symmetry and Asymmetry

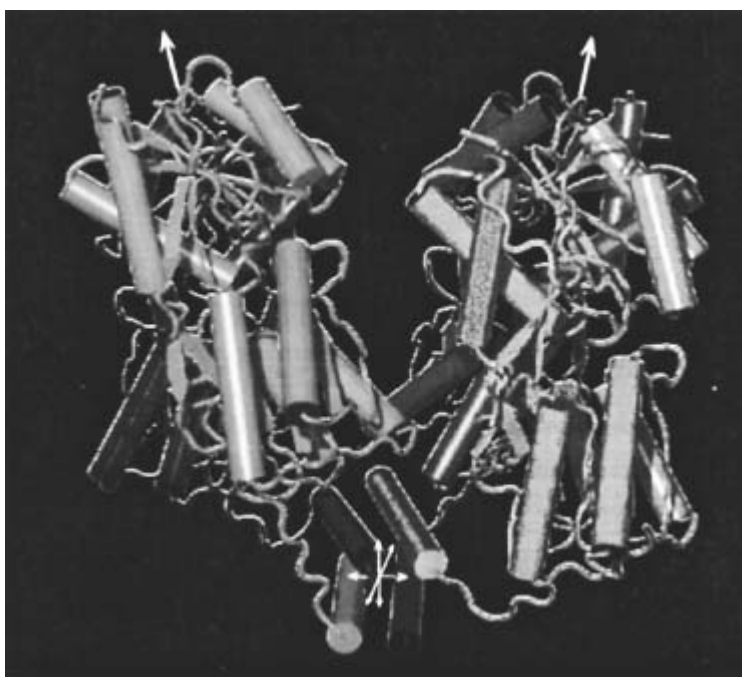
Symmetry requires the exact geometric repetition of chemically identical units, but approximate symmetry can nevertheless be observed in assemblies that do not satisfy this condition exactly. In mammalian **hemoglobins**, the  $a_2b_2$  oligomer displays the exact symmetry of point group  $C_2$  with the two-fold axis relating the two ab units; in addition, there is an approximate  $D_2$  symmetry equivalencing the very similar tertiary structures of the **homologous** a and b chains. Approximate symmetry between structural **domains** of a single-chain protein is also well documented, and it is usually interpreted as an indication that the protein derived from a symmetrical homodimer by [gene duplication](#) and [gene fusion](#).

Because symmetry is so frequent, asymmetry is remarkable when it occurs, and it is usually reflected in the function. The F1 fragment of the [ATP synthase](#) of [mitochondria](#) and [chloroplasts](#) has the formula  $a_3b_3g$ . The three-fold symmetry of the assembly of a and b chains is broken by the presence of a single g chain in the middle. Contacts with g make the three [active sites](#) carried by the three b chains nonequivalent, an essential feature of the catalytic mechanism. Simpler examples are two dimeric proteins from the human immunodeficiency virus ([HIV](#)), the HIV [proteinase](#) and HIV **reverse transcriptase**. In both, X-ray crystallography studies have demonstrated departure from two-fold symmetry. Alone, the HIV proteinase dimer does display exact  $C_2$  symmetry, but this is incompatible with the binding of a peptide substrate and with the catalytic mechanism. A peptide cannot have internal symmetry, and the mechanism of hydrolysis requires one of two intrinsically

equivalent [active site](#) aspartate residues to be protonated, the other deprotonated (see [Carboxyl Proteinase](#)). In crystalline complexes of the enzyme with substrate analogues, the symmetry is broken. Minor structural changes occur in the protein to fit the asymmetric ligand, and the complex retains approximate symmetry. In contrast, HIV reverse transcriptase displays a major asymmetry. Its two chains are the product of the same gene, but one has undergone **proteolytic** processing. They adopt grossly different folds in the dimer and play very different roles in DNA synthesis. In this case, asymmetry exists in the chemical as well as the three-dimensional structure. It is not known whether or not the asymmetry preexists the proteolytic cleavage that yields the active dimer.

Less often, oligomeric proteins display symmetries that do not belong to one of the point groups mentioned above and cannot be exact. The tetrameric *E. coli* lactose repressor is a case in point (Fig. 2). It is assembled from two dimers, each of which has  $C_2$  symmetry. The dimers are related by an approximate two-fold axis that is not orthogonal to those of the dimers as  $D_2$  symmetry would require. The result is that all four subunits have their DNA-binding headpieces on the same face of the tetramer, which they could not do in a  $D_2$  tetramer. On the opposite face, the C-terminal  $\alpha$ -helices of each subunit form most of the dimer-dimer contacts. Remarkably, the four C-terminal  $\alpha$ -helices assemble with  $D_2$  symmetry, whereas the rest of the protein does not, which a long connecting peptide makes possible by adopting different conformations in two of the subunits (3).

**Figure 2.** The lactose repressor of *E. coli*. The homo-tetramer is made of two dimers, each with two-fold symmetry. However, their two-fold axes (arrows) are not orthogonal, and only the four-helix bundle made by C-terminal  $\alpha$ -helices at the bottom of the molecule shows the usual  $D_2$  (or  $222$ ) symmetry of tetramers. Residues 1 to 61 of each 360-residue monomer form N-terminal DNA-binding domains located at the top of the molecule. In the absence of DNA, they are so mobile that they cannot be observed in the crystal structure of the protein alone (3).



### 3. Subunit Interactions

The stabilities of quaternary structures result from the contacts between subunits. Subunits form pairwise interfaces of two different characters, depending on whether or not they are related by a two-fold axis. Two-fold symmetry implies that the same surface and same amino acid residues of

both subunits are in contact. Interfaces having other cyclic symmetries, or no symmetry at all, usually involve different sets of residues on the two partners. Following the nomenclature of Monod et al. (1) once again, *isologous* interfaces refer to those with two-fold symmetry, *heterologous* the others. All interfaces are isologous in C<sub>2</sub> dimers and D<sub>2</sub> tetramers; in other systems, at least one set of interfaces must be heterologous.

When the three-dimensional structure is known from X-ray or NMR studies, the extent and nature of the interactions between polypeptide chains can be analyzed in atomic detail. A convenient estimate of the extent of the contact is the *interface area*, the area *B* of the protein surface that is buried (made inaccessible to the solvent) as a result of the interaction. In a dimer, *B* is derived from the atomic coordinates of the structure by calculating the area of the solvent-accessible surface for the dimer and for the two isolated subunits, and subtracting the first value from the sum of the other two. Table 2 quotes interface areas observed in some oligomeric proteins. *B*<sub>tot</sub> is the total area of the surface buried in the quaternary structure. In all cases, it is at least 1400 Å<sup>2</sup>, similar to the surface area buried in contacts between an antibody and a protein [antigen](#), or [proteinase](#) and a [proteinase inhibitor](#) (see [Protein-Protein Interactions](#)), and this may be the minimum required for stable association. However, most oligomeric proteins have much larger subunit interfaces than this minimum. In bovine [catalase](#), the quaternary structure buries as much as 70% of the protein accessible surface, in six large interfaces. All are isologous, and the D<sub>2</sub> symmetry of the tetramer makes them two by two equivalent. Thus, only three pairwise interface areas are listed in Table 2.

**Table 2. Area of Interfaces between Subunits of Oligomeric Proteins<sup>a</sup>**

| Protein                                     | Molecular Weight (Da) | Accessible Surface Area (Å <sup>2</sup> ) | <i>B</i> <sub>tot</sub> | Interface Area (Å <sup>4</sup> ) |      |      |
|---|-----------------------|---|-------------------------|----------------------------------|------|------|
|   |                       |   |                         | <i>B</i> <sub>pair</sub>         |      |      |
| <i>Dimers</i>                               |                       |   |                         |                                  |      |      |
| Avian pancreatic peptide                    | 8500                  | 5300                                      | 1400                    |                                  |      |      |
| Uteroglobin                                 | 15,800                | 7500                                      | 3000                    |                                  |      |      |
| Superoxide dismutase                        | 31,400                | 13,800                                    | 1350                    |                                  |      |      |
| Triosephosphate isomerase                   | 53,000                | 20,300                                    | 3180                    |                                  |      |      |
| <a href="#">Alcohol dehydrogenase</a> (ADH) | 79,800                | 29,000                                    | 3260                    |                                  |      |      |
| Citrate synthase                            | 96,000                | 28,500                                    | 9800                    |                                  |      |      |
| <i>Tetramers</i>                            |                       |   |                         |                                  |      |      |
| <b>Mellitin</b>                             | 11,400                | 6300                                      | 4160                    | 880                              | 820  | 520  |
| Glutathione peroxidase                      | 83,600                | 28,600                                    | 6280                    | 1520                             | 1520 | 180  |
| <a href="#">Phosphofructokinase</a>         | 141,200               | 40,600                                    | 14,400                  | 4500                             | 2520 | 180  |
| Catalase                                    | 231,700               | 60,900                                    | 42,300                  | 9260                             | 9140 | 4120 |
| <i>Hexamer</i>                              |                       |   |                         |                                  |      |      |
| <a href="#">Insulin</a>                     | 34,600                | 13,100                                    | 8600                    | 1280                             | 1400 | *180 |



## Octomer

|             |         |        |        |      |     |       |
|-------------|---------|--------|--------|------|-----|-------|
| Hemerythrin | 107,500 | 35,900 | 13,600 | 1780 | 260 | *1420 |
|-------------|---------|--------|--------|------|-----|-------|

---

<sup>a</sup> Adapted from ref. 4, where references to X-ray structures can be found. The molecular weight, accessible surface area, and total interface area are quoted for the whole molecule. Tetramers and larger oligomers contain several interfaces, with areas given per pair of subunits. Dimers in this table have  $C_2$  symmetry, tetramers,  $D_2$  symmetry, with isologous interfaces only. Larger oligomers also contain heterologous interfaces marked with an asterisk. ( $1 \text{ \AA} = 10^{-10} \text{ m.}$ )

In general, the buried surface area  $B_{\text{tot}}$  is distributed among several pairwise interfaces, up to  $n_{\text{pair}} = n(n-1)/2$  for an assembly of  $n$  subunits when each one is in contact with all others. In a dimer,  $n_{\text{pair}} = 1$ , in a tetramer,  $n_{\text{pair}} = 6$ . Indeed, most  $D_2$  tetramers have six pairwise interfaces like catalase. In [phosphofruktokinase](#), let us label the four subunits  $A$ ,  $B$ ,  $C$ , and  $D$ . Symmetry-equivalent  $AB$  and  $CD$  pairs make extensive contacts burying  $4500 \text{ \AA}^2$  ( $\text{\AA} = 10^{-10} \text{ m}$ ) each. As  $AC$  and  $BD$  pairs bury significantly less, and  $AD$  and  $BC$  pairs very little, phosphofruktokinase appears, based on the size of the pairwise interfaces, to be a dimer of  $AB$ -like dimers. It might be expected that  $AB$ -like dimers are actual intermediates in tetramer assembly, but this is not known. Dividing large oligomers into smaller assemblies on the basis of their three-dimensional structure is often an arbitrary decision. Hemerythrin, illustrated in Figure 1, is an octamer with  $D_4$  symmetry. Is it a dimer of cyclic tetramers, or a cyclic tetramer of dimers? Either description fits. Within a tetramer, each hemerythrin subunit makes two symmetry-equivalent heterologous interfaces with neighbors, eg, the green subunit of the top tetramer with the red and purple subunits, and two isologous interfaces with subunits of the other tetramer, eg, the green subunit with the yellow and brown bottom subunits. Both the green–yellow interface forming a dimer and the green–red or green–purple interfaces forming the tetramer are relatively large, making both types of substructures plausible.

The chemical composition of the protein surface involved in subunit contacts is generally richer in aliphatic or aromatic groups than is the surface accessible to solvent, and it therefore is less [polar](#). On average, [nonpolar](#) groups contribute 65% of the interface areas and only 57% of the accessible surface (4). From this point of view, subunit interfaces in oligomeric proteins tend to resemble protein interiors, whereas the interfaces of protein–protein complexes are more like the rest of the protein surface (see [Protein–Protein Interactions](#)). **Hydrophobic** (aliphatic or aromatic) amino acid residues are correspondingly overrepresented, and the most abundant residue at subunit interfaces is [leucine](#). The well-known [leucine zipper](#) is an example of a [coiled-coil](#) dimerization motif formed by the surfaces of two leucine-rich  $\alpha$ -helices, but there are many other situations where leucine residues contribute to interfaces. Remarkably, the second most abundant residue is [arginine](#). Its very polar guanidinium group, which is almost always a surface group in monomeric proteins, is frequently involved in the quaternary structure of oligomers. This and other polar groups buried in subunit contacts contribute to the stability of the assembly by forming intersubunit [hydrogen bonds](#). Such bonds occur at the rate of about one per  $250 \text{ \AA}^2$  of interface area (5), and the partner groups are favorite targets for [site-directed mutagenesis](#) aiming to perturb the quaternary structure.

## 4. Quaternary Structure Changes

The quaternary structure of oligomeric proteins is often as stable as the tertiary structure itself, with the subunits dissociating only under conditions that **denature** the polypeptide chains. There are nevertheless cases where a dissociation equilibrium is observed under nondenaturing conditions. In the case of human hemoglobin, the  $\alpha_2\beta_2$  oligomer is normally stable when the protein is in the deoxy-form, but when oxygen binds, it dissociates into  $\alpha\beta$  units with an equilibrium **dissociation**

**constant** of the order of 1  $\mu\text{M}$ . This is much too low for dissociation to play a physiological role in red blood cells, where the hemoglobin concentration is greater than millimolar, yet the coupling between oxygen binding and tetramer dissociation provides a powerful biochemical tool for the study of the allosteric transition (6).

While remaining tetramers, mammalian hemoglobins can still adopt more than one quaternary structure. This has a major functional significance, for the different structures have different affinities for oxygen and other ligands. The **allosteric** model of Monod et al. (1) makes quantitative predictions of how a quaternary structure change can lead to **cooperative** ligand binding in an oligomeric protein. The X-ray structures of deoxy- and oxy- or carbonmonoxy-hemoglobin give a structural basis for this model. Both forms of the protein are symmetrical  $\alpha_2\beta_2$  tetramers, yet the relative orientations of the  $\alpha\beta$  units differ, and many contacts at the interface are disrupted or rearranged. Changes within  $\alpha\beta$  units (tertiary structure changes) are of much lesser amplitude, but connect the heme sites to the interface. At least one other quaternary structure, called Y or R2, has been identified in mammalian hemoglobins, confirming the plasticity of the interface between the  $\alpha\beta$  units. In contrast, the interface between  $\alpha$  and  $\beta$  within each unit is essentially invariant.

Quaternary structure changes generally appear in X-ray structures as rotation or translation movements of the constituent subunits relative to each other, which perturb the contacts that hold together the assembly, without abolishing them. Such movements have been characterized in several allosteric proteins, for instance, in *E. coli* [aspartate transcarbamoylase](#), but they also exist in oligomeric proteins that are not usually considered allosteric and display no cooperativity in ligand binding. In immunoglobulins, the interface between the variable domains of the heavy and light chains shows significant plasticity, and this may be important for antigen recognition, since the combining site spans that interface.

## 5. Folding and Evolution

The oligomeric proteins listed in Table 2 are fairly globular, and their overall accessible surface area is correlated with their molecular weight (4). In many cases, the individual subunits are also reasonably globular and can remain folded in isolation, so these proteins assemble by the individual subunits folding, then associating. In other cases, however, the subunits themselves are often far from globular, and some parts of individual polypeptide chains may make more interactions with other subunits than with their own. In such cases, the subunit fold found in the oligomer is unlikely to preexist the assembly, and folding must therefore be tightly coupled with association. The **protein folding** of oligomeric proteins is less easily studied than that of small monomeric ones.

**Denaturation** with [urea](#) or [guanidinium salts](#) is often irreversible, as the correct reassembly of the monomers competes with their aggregation, leading to insolubility. However, it has been achieved in a number of cases, showing that complete self-assembly can be reproduced *in vitro* (7). In some cases, a folded monomer or smaller oligomer can be characterized as an intermediate. For instance, *E. coli* aspartate transcarbamoylase catalytic chains form active trimers in the absence of regulatory chains. But, intermediates are often too short-lived for detection during reassembly. This is illustrated by the possibly extreme case of the Arc repressor of bacteriophage P22, a small dimer that can refold very rapidly. No folded monomer is detected, and refolding and association appear simultaneously as a second-order reaction, even on a millisecond time scale (8).

The evolution of oligomeric proteins offers a number of similarly remarkable features. Closely related proteins usually have the same quaternary structure, and residues forming the  $\alpha\beta$  interface in mammalian hemoglobins are invariant like the internal residues. Nevertheless, the quaternary structure is much less well-conserved than the subunit fold. Globins can be monomeric ([myoglobin](#), plant and insect hemoglobins), dimeric (lamprey hemoglobin), tetrameric (mammalian hemoglobins), or larger oligomers, while preserving the same fold and the same capacity of binding heme and oxygen. Other protein families are like the globins and, in general, the quaternary structure need not be the same in homologous proteins when the amino acid sequence identity is below 50%, a

level at which tertiary structures hardly change. In globins, oligomerization is accompanied by a gain of function: Oxygen binding becomes cooperative. As mentioned above, the cooperativity results from an allosteric mechanism whereby ligation at the heme site is coupled to a change in quaternary structure. There is therefore a [natural selection](#) pressure on the quaternary structure distinct from that which maintains the tertiary structure of the subunits. It is focused on amino acid residues forming interfaces between subunits, and specific disorders in humans are associated with mutations that change these residues.

The gain of function is even more evident in many enzymes where the [active site](#) is made of residues from more than one subunit. An example is the dimeric HIV proteinase mentioned above. Its catalytic mechanism is the same as in monomeric carboxyl proteinases, such as **pepsin**, and involves two aspartate residues. In HIV proteinase, they are the same residue from each of the two subunits; in pepsin, they are different. The longer polypeptide chain of pepsin is a two-domain structure that displays approximate internal symmetry, each domain resembling the viral proteinase monomer. Thus, the single-chain enzyme reproduces the structure of the dimer, giving strong support to the hypothesis that single-chain carboxyl proteinases evolved by gene duplication and fusion from a shorter dimeric precursor. This hypothesis had been made on the basis of the sequence and structure of aspartate proteinases many years before the viral enzyme was known (9).

Dimerization is often assumed to be an acquired feature in evolution. However, in the carboxyl proteinase family, the HIV-type of dimer is the ancestral form relative to the pepsin-type of monomer and, moreover, a monomeric ancestor could not have been a proteinase, for this function requires both aspartate residues. This shows that evolution can undo quaternary structure, as well as create it. In the globin family, the more complex dimeric or tetrameric hemoglobins are considered to have evolved from single-chain globins found in lower organisms, but in many other families there is no evidence for a monomeric precursor. Nevertheless, an evolutionary scenario for dimerization can be proposed. As pointed out by Monod et al. (1), isologous association should be easier to achieve, because two-fold symmetry duplicates any favorable (or unfavorable) interaction. Thus, a small number of mutations at the protein surface might yield a  $C_2$  dimer. An alternative model for dimerization is domain swapping (10), which does not even require mutations to form novel interactions. In domain swapping, domains *A* and *B* interact either within a monomer, or with domains *A'-B'* of a second monomer, yielding a (*AB'*) (*A'B*) dimer where noncovalent *A*-to-*B'* and *A'*-to-*B* interactions are essentially identical to *A*-to-*B* interactions in the monomeric form. Domain swapping was named for the phenomenon in [diphtheria toxin](#), where it can be induced by temperature changes. There are some obvious cases in families of homologous proteins, such as the *crystallin* family of eye lens proteins in mammals (12). It includes monomeric g-crystallins and dimeric b-crystallins, all made of two domains that have very similar folds and related sequences, but swap their positions in the dimeric proteins, relative to the monomeric proteins.

Domain swapping is a plausible evolutionary source of oligomers in which large parts of the polypeptide chain wander away from their subunit to interact with others, yielding some of the largest interfaces in Table 2. It could, in principle, yield more complex assemblies with cyclic symmetry, eg, a (*AB''*) (*A''B'*) (*A'B*) trimer where each subunit donates a domain to its neighbor. However, assemblies such as hemoglobin or hemerythrin in Figure 1 show no indication of domain swapping, and many evolutionary mechanisms must coexist. Cells are very dense and complex mixtures of macromolecules. Any change on a protein surface must affect its capacity not only to form legitimate interactions with its siblings in an oligomer, but also illegitimate ones with other components of the cytoplasm, membrane, or organelle, and all these interactions compete in evolution.

## Bibliography

1. J. Monod, J. Wyman, and J. P. Changeux (1965) *J. Mol. Biol.* **12**, 88–118.
2. T. L. Blundell and N. Srinivasan (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14243–14248.

3. M. Lewis, G. Chang, N. C. Horton, M. A. Kercher, H. C. Pace, M. A. Schumacher, R. G. Brennan, and P. Lu (1996) *Science* **271**, 1247–1254.
4. J. Janin, S. Miller, and C. Chothia (1988) *J. Mol. Biol.* **204**, 155–164.
5. S. Jones and J. M. Thornton (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13–20.
6. G. K. Ackers, M. L. Doyle, D. Myers, and M. A. Daugherty (1992) *Science* **255**, 54–63.
7. R. Jaenicke (1987) *Prog. Biophys. Mol. Biol.* **49**, 117–237.
8. C. D. Waldburger, T. Jonsson, and R. T. Sauer (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2629–2634.
9. J. Tang, M. N. G. James, I. N. Hsu, J. A. Jenkins, and T. L. Blundell (1978) *Nature* **271**, 618–622.
10. M. J. Bennett, M. P. Schlunegger, and D. Eisenberg (1994) *Prot. Sci.* **4**, 2455–2468.
11. R. E. Stenkamp, L. C. Sieker, L. H. Jensen, J. D. McCallum, and J. Sanders-Loehr (1985) *Proc. Natl. Acad. Sci. USA* **82**, 713–718.
12. B. Bax, R. Lapatto, V. Nalini, H. Driessen, P. F. Lindley, D. Mahadevan, T. L. Blundell, and C. Slingsby (1990) *Nature* **347**, 776–780.

### Suggestions for Further Reading

13. C. Brändén and J. Tooze (1991) *Introduction to Protein Structure*, Garland Publishing, New York.
14. R. Dickerson and I. Geis (1983) *Hemoglobin: Structure, Function, Evolution and Pathology*, Benjamin/Cummings, Menlo Park CA.
15. R. Jaenicke and R. Rudolph (1986) Refolding and association of oligomeric proteins. *Meth. Enz.* **131**, 218–250.
16. M. Perutz (1989) Mechanisms of cooperativity and allosteric regulation in proteins, *Quar. Rev. Biophys.* **22**, 139–236.

## Oligomycin

The oligomycins are a family of macrolide antibiotics composed of three major members: oligomycins A, B, and C that are produced by *Streptomyces diastatochromogenes* (1). The family now includes rutamycin (oligomycin D), rutamycin B, and oligomycins E, F, and G (1-5). This antibiotic has antifungal activity, but it is not clinically useful because of its toxicity to mice and [HeLa Cells](#) (1, 3). However, Lardy et al (6) reported that oligomycins inhibit oxidative phosphorylation of [mitochondria](#), a finding that attracted the attention of biochemists. Since then the oligomycin family has been shown to be a specific and reversible inhibitor of two [ATPases](#): H<sup>+</sup>-ATPase that is present in mitochondrial inner membranes and the Na<sup>+</sup>, K<sup>+</sup>-ATPase that is present in animal plasma membranes. In addition, it has been reported that oligomycins inhibit some drug-translocating ATPases that confer [drug-resistance](#) to mammalian and yeast cells, although the inhibition mechanism is unknown (7, 8), and that oligomycins are suppressive agents for human B cell activation (4).

### 1. Physicochemical Properties

Oligomycins A, B, and C are commercially available, separately and in mixtures, from a few companies, including Sigma Chemical Company (USA) and Merck KGaA (Germany). Because the

three compounds are used extensively in the fields of biochemistry and molecular biology, their physicochemical properties are described.

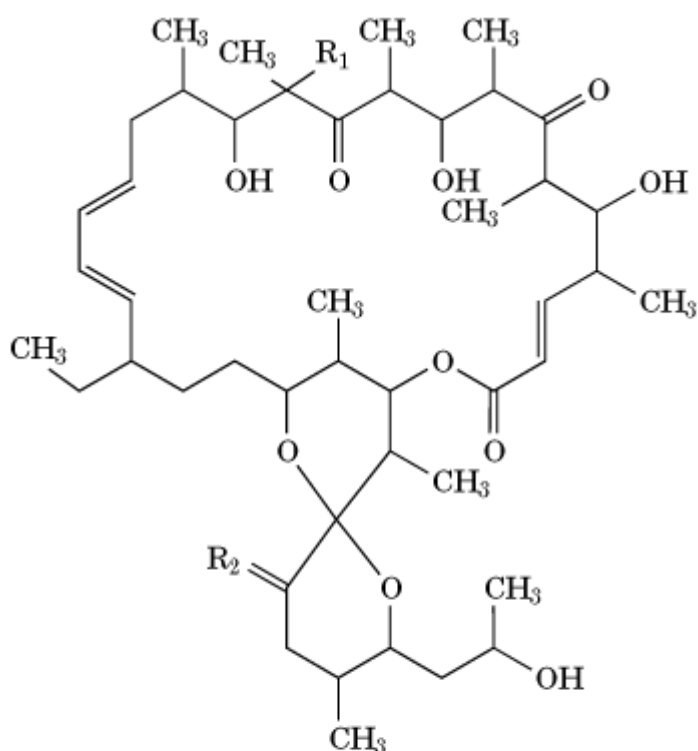
### 1.1. Molecular Formulas and Formula Weights

The molecular formulas and formula weights of oligomycins are  $C_{45}H_{74}O_{11}$  and 791.1 for oligomycin A,  $C_{45}H_{72}O_{12}$  and 805.1 for oligomycin B, and  $C_{45}H_{74}O_{10}$  and 775.1 for oligomycin C (2, 10).

### 1.2. Structures

Structures of the oligomycins are shown in Figure 1. The three-dimensional structure of oligomycin B was determined by von Glehn et al (9).

**Figure 1.** Structures of the oligomycins (2, 3, 10).



| Oligomycin | R <sub>1</sub> | R <sub>2</sub> |
|------------|----------------|----------------|
| A          | OH             | H <sub>2</sub> |
| B          | OH             | O              |
| C          | H              | H <sub>2</sub> |

### 1.3. Molecular Absorption Coefficient (e)

The maximal molecular absorbance coefficients of the oligomycins in methanol are 37,400 at 225 nm for oligomycin A, 32,200 at 224 nm for oligomycin B, and 33,500 at 224 nm for oligomycin C (2).

### 1.4. Solubilities

The oligomycins are stable compounds. Their inhibitory potency is little changed on storage at pH 3 to 10 at 37°C for 54 h (10). They are almost insoluble in water but are soluble in water-miscible solvents such as ethanol, methanol, dimethyl sulfoxide, dimethylformamide (10) acetone, glacial

acetic acid, and in ether (1). Oligomycins are usually dissolved in ethanol and stored in a refrigerator. Their maximum solubility in ethanol seems to be not more than 10 mM (about 8 mg/ml). Their solubility in water including 1% ethanol seems, from examination using [<sup>3</sup>H]-labeled oligomycin made in our laboratory and not commercially available, to be up to 10 mM (about 8 mg/ml).

### 1.5. Toxicities

The LD<sub>50s</sub> of oligomycins A, B, and C in mice injected intraperitoneally are 1.5, 2.9, and 8.3 mg/kg of body weight, respectively (1). The corresponding inhibitory concentrations (IC<sub>50s</sub>) toward HeLa cells are 0.008, 0.015, and 0.106 mg/ml, respectively (3).

### 1.6. Other Properties

The melting points, infrared absorption spectra, and optical rotations of the oligomycins, including oligomycins D, E, G, and rutamycin B are available in the literatures (1-3, 5, 11).

## 2. Biochemical Properties

Various ion (proton)-translocating ATPases are present in prokaryotic and eukaryotic organisms. These ATPase are classified into three groups: F-type, V-type, and P-type (12) (see ATPase).

Mitochondrial H<sup>+</sup>-ATPase belongs to the F-type group and is called F<sub>1</sub>F<sub>0</sub>-ATPase or [ATP synthase](#) in some cases. Other F-type H<sup>+</sup>-ATPases, present in thylakoid membranes of chloroplasts and inner membranes of bacteria, other than some photosynthetic bacteria, are oligomycin-insensitive (10).

Na<sup>+</sup>, K<sup>+</sup>-ATPase belongs to the P-type group. In some cases, this ATPase is called the **sodium pump**. Other P-type ATPases, including Ca<sup>2+</sup>-ATPase and H<sup>+</sup>, K<sup>+</sup>-ATPase, localized in the sarcoplasmic reticulum and gastric cell membranes, respectively, are oligomycin-insensitive (13, 14). V-type H<sup>+</sup>-ATPase, which is localized in membranes of subcellular particles in eukaryotic organisms, such as vacuoles and lysosomes, is also oligomycin-insensitive (15).

The mechanisms by which oligomycin inhibits mitochondria H<sup>+</sup>-ATPase and Na<sup>+</sup>, K<sup>+</sup>-ATPase are explained below.

### 2.1. Inhibition of Mitochondrial H<sup>+</sup>-ATPase by Oligomycin

Mitochondrial H<sup>+</sup>-ATPase functions as an ATP synthase coupled to the respiratory electron flow in intact mitochondria. The yeast mitochondrial ATPase is less sensitive to oligomycin than the animal mitochondrial ATPase (16). The inhibitory potency of oligomycin is often expressed as an amount of oligomycin per milligram of ATPase: 95% inhibition is achieved by 0.4-μg oligomycin/mg of bovine heart mitochondrial ATPase; 85% inhibition is achieved by 4-μg oligomycin/mg of *Neurospora crassa* ATPase; and 90% inhibition is achieved by 10-μg oligomycin/mg *Saccharomyces cerevisiae* ATPase (17). F-type H<sup>+</sup>-ATPase is composed of two components: the F<sub>1</sub> and F<sub>0</sub> portions. The F<sub>1</sub> portion is an oligomer composed of five different kinds of [hydrophilic](#) subunits named a, b, g, d, and e, and it includes the catalytic center for the synthesis or hydrolysis of ATP. The F<sub>0</sub> portion is an integral [membrane protein](#) composed of three (bacteria) or more (eukaryotic organisms) different kinds of **hydrophobic** subunits (Table 1), and it functions as a proton channel. The ATPase activity of the isolated F<sub>1</sub> region is not inhibited by oligomycin (18).

The component that controls the oligomycin sensitivity of H<sup>+</sup>-ATPase has been purified from bovine heart mitochondrial ATPase and named oligomycin sensitivity-conferring protein (OSCP) (19). However, OSCP does not contain a binding site for oligomycin; instead, it is one of the components required for correct binding between the F<sub>1</sub> and F<sub>0</sub> portions (18, 20). The amino acid sequence of OSCP has [homology](#) with that of the d subunit in the F<sub>1</sub> portion of bacterial H<sup>+</sup>-ATPase,

which is insensitive to oligomycin (21). Oligomycin-resistant mutants have been isolated from yeast and Chinese hamster ovary cell lines and analyzed genetically. The majority of the mutations were mapped at loci *oli 1* to *oli 4* (16). The *oli 2* and *oli 4* loci were found in the gene for subunit 6 of  $F_0$ , whereas the *oli 1* and *oli 3* were found in the gene for subunit 9 of  $F_0$  (Table 1) (16, 22). Some of the amino acids substituted in the mutants have been identified (23, 24). In addition, the binding of dicyclohexylcarbodiimide (DCCD), another inhibitor of  $H^+$ -ATPase, to subunit 9 is reduced by oligomycin (16). In *Escherichia coli*  $H^+$ -ATPase, which is insensitive to oligomycin, the region around *oli 4* in the gene for subunit a, which corresponds to subunit 6 (Table 1), is deleted? suggesting that the binding site for oligomycin is absent in the *E. coli* protein (22). These results suggest that subunits 6 and 9 of the  $F_0$  portion determine the oligomycin sensitivity of F-type  $H^+$ -ATPase (16, 22). OSCP apparently affects the oligomycin sensitivity indirectly.

**Table 1. Correspondence between Subunits of  $F_0$  Portions of F-Type  $H^+$ -ATPases from Various Sources (26)**

| Bacteria | Yeast  | Mitochondria | Animal Mitochondria |
|----------|--------|--------------|---------------------|
| $a^a$    | $6^b$  |              | $a^b$               |
| $b^a$    | P25    |              | PVP                 |
| —        | P18    |              | $F_6$               |
| —        | $F_1I$ |              | $F_1I$              |
| $c^a$    | $9^c$  |              | $c^c$               |
| —        | 8      |              | $A_6L$              |
| —        | —      |              | d                   |
| —        | —      |              | $F_B$               |

<sup>a</sup> Other nomenclatures of subunit a, b and c are, respectively, *uncB* (or c), *uncF* (or y), and *uncE* (or w) (25).

<sup>b</sup> Another name for subunit 6 and a is ATPase 6 (22, 24, 26).

<sup>c</sup> Another name for subunit 9 and c is ATPase 9, or Proteolipid, or DCCD-binding protein (24, 27, 28).

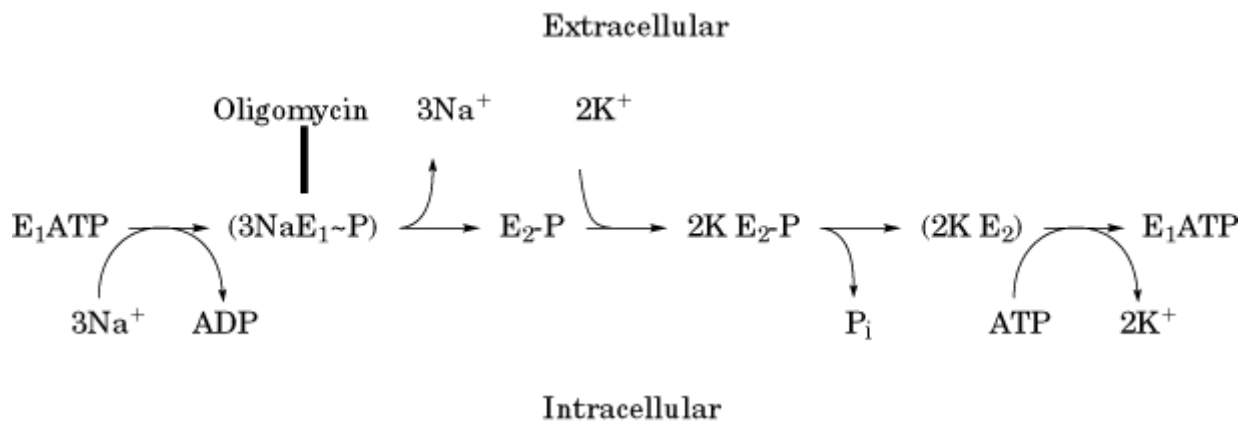
## 2.2. Inhibition of $Na^+$ , $K^+$ -ATPase by Oligomycin

$Na^+$ ,  $K^+$ -ATPase is an **antiporter**, which actively translocates  $Na^+$  from the inside to the outside of animal cells and translocates  $K^+$  in the reverse direction, using the hydrolysis of ATP as the driving force. The resultant  $Na^+$  gradient is essential for excitation of nerve tissues and for absorption in the intestine of amino acids and glucose (29). Oligomycin inhibits ATPase activity by about 80% at 10 mg/ml (30-32), and the inhibition is incomplete (32, 33).  $Na^+$ ,  $K^+$ -ATPase is composed of a and b subunits (29): the a subunit possesses the catalytic center. In the ATPase reaction sequence, known as the Post-Albers scheme (Fig. 2), oligomycin inhibits neither the formation of the  $Na^+$ -dependent phosphorylated intermediate ( $3NaE1\sim P$ ) nor the  $K^+$ -dependent  $E_22\text{-P}$  hydrolysis activity (34, 35).

The point of inhibition by oligomycin is the  $3NaE1\sim P$  to  $E_22\text{-P} + 3Na^+$  step, which is coupled to

Na<sup>+</sup> translocation (Fig. 2) (34, 35). These findings are consistent with the inhibition of only Na<sup>+</sup> translocation by oligomycin *in vivo* (32, 36). The sodium ions in the 3NaE<sub>1</sub>~P intermediate are called the occluded sodium because these sodium ions seem to be embedded transiently in the ATPase molecule (37). Homareda et al (31, 32, 38) showed that oligomycin increases the affinity of Na<sup>+</sup>, K<sup>+</sup>-ATPase for Na<sup>+</sup>, but not for K<sup>+</sup>, under nonphosphorylative conditions. This finding is in agreement with the observation that oligomycin accelerates the rate of the transition from E<sub>2</sub>K to E<sub>1</sub>Na but not in the reverse direction (39-41). These results suggest that oligomycin leads to the occlusion of sodium ions by interacting specifically with the Na-bound Na<sup>+</sup>, K<sup>+</sup>-ATPase, irrespective of whether the Na<sup>+</sup>, K<sup>+</sup>-ATPase is phosphorylated or not. Therefore, oligomycin has been used to block the reaction from 3NaE<sub>1</sub>~P to E<sub>2</sub>2-P + 3Na<sup>+</sup> and as a stabilizer of the occluded Na form.

**Figure 2.** Reaction sequence of Na<sup>+</sup>,K<sup>+</sup>-ATPase according to the Post-Albers scheme (34, 35). E<sub>1</sub> and E<sub>2</sub> denote the conformation showing high affinity for Na<sup>+</sup> and K<sup>+</sup>, respectively. E<sub>1</sub>~P and E<sub>2</sub>-P denote the phosphorylated intermediate with high and low energies, respectively. 3NaE<sub>1</sub>~P and 2K·E<sub>2</sub> denote the confirmation with occluded Na and K, respectively. P<sub>i</sub> is inorganic phosphate.



Various Na<sup>+</sup>, K<sup>+</sup>-ATPase **isoforms** have either low or high sensitivities to ouabain, a classical Na<sup>+</sup>, K<sup>+</sup>-ATPase inhibitor; both are inhibited by oligomycin (32). Oligomycin seems to act from the extracellular side of Na<sup>+</sup>, K<sup>+</sup>-ATPase molecule (42). It is suggested that the 200 N-terminal amino acids of the Na<sup>+</sup>, K<sup>+</sup>-ATPase a subunit include a binding domain of oligomycin (43).

#### Bibliography

1. S. Masamune, J. M. Sehgal, E. E. van Tamelen, F. M. Strong, and W.H. Peterson (1958) J. Am. Chem. Soc. **80**, 6092–6095.
2. G.T. Carter (1986) J. Org. Chem. **51**, 4264–4271.
3. K. Kobayashi et al (1987) J. Antibiot. **40**, 1053–1057.
4. H. Laatsch et al (1993) J. Antibiot. **46**, 1334–1341.
5. Y. Enomoto et al (2001) J. Antibiot. **54**, 308–313.
6. H. A. Lardy, D. Johnson, and W. C. McMurray (1958) Arch. Biochem. Biophys. **78**, 587–597.
7. G. D. Eytan, M. J. Borgnia, R. Regev, and Y. G. Assaraf (1994) J. Biol. Chem. **269**, 26058–26065.
8. A. Decottignies et al (1995) J. Biol. Chem. **270**, 18150–18157.



9. M. von Glehn, R. Norrestam, P. Kierkegaard, and L. Maron (1972) *FEBS Lett.* **20**, 267–269.
10. P. E. Linnett and R. B. Beechey (1979) *Methods Enzymol.* **55**, 472–518.
11. B. Arnoux, et al (1978) *J. Chem. Soc. Chem. Commun.* **1978**, 318–319.
12. P. L. Pedersen and E. Carafoli (1987) *Trends Biochem. Sci.* **12**, 146–150.
13. D. H. MacLennan (1970) *J. Biol. Chem.* **245**, 4508–4518.
14. H. E. Ives and F. C. Rector Jr. (1984) *J. Clin. Invest.* **73**, 285–290.
15. E. Uchida, Y. Ohsumu, and Y. Anraku (1985) *J. Biol. Chem.* **260**, 1090–1095.
16. W. Sebald and J. Hoppe (1981) *Curr. Top. Bioenerg.* **12**, 1–64.
17. D. S. Perlin, L. R. Latchney, and A. E. Senior (1985) *Biochim. Biophys. Acta* **807**, 238–244.
18. A. E. Senior (1979) *Methods Enzymol.* **55**, 391–472.
19. D. H. MacLennan and A. Tzagoloff (1968) *Biochemistry* **7**, 1603–1610.
20. L. Ernster, T. Hundal, and G. Sandri (1986) *Methods Enzymol.* **126**, 428–433.
21. J. E. Walker, M. J. Runswick, and M. Saraste (1982) *FEBS Lett.* **146**, 393–396.
22. A. E. Senior and J. G. Wise (1983) *J. Membr. Biol.* **73**, 105–124.
23. U. P. John and P. Nagley (1986) *FEBS Lett.* **207**, 79–83.
24. G. A. M. Breen, et al (1986) *J. Biol. Chem.* **261**, 11680–11685.
25. M. Futai and H. Kanazawa (1983) *Microbiol. Rev.* **47**, 285–312.
26. B. Hamasur and E. Glaser (1992) *Eur. J. Biochem.* **205**, 409–416.
27. Y. Hatefi (1985) *Ann. Rev. Biochem.* **54**, 1015–1069.
28. J. Hoppe, D. Gatti, H. Weber, and W. Sebald (1986) *Eur. J. Biochem.* **155**, 259–264.
29. J. B. Lingrel and T. Kuntzweiler (1994) *J. Biol. Chem.* **269**, 19659–19662.
30. A. Askari and D. Koyal (1968) *Biochem. Biophys. Res. Commun.* **32**, 227–232.
31. H. Homareda and H. Matsui (1982) *J. Biochem. (Tokyo)* **92**, 219–231.
32. T. Arato-Oshima, H. Matsui, A. Wakizaka, and H. Homareda (1996) *J. Biol. Chem.* **271**, 25604–25610.
33. J. D. Robinson (1971) *Mol. Pharmacol.* **7**, 238–246.
34. J. D. Cavieres (1977) In *Membrane Translocation in Red Cells* (J. C. Ellory and V. L. Lew, eds.) Academic Press, New York, pp. 1–37.
35. J. D. Robinson and M. S. Flashner (1979) *Biochim. Biophys. Acta* **549**, 145–176.
36. P. J. Garrahan and I. M. Glynn (1967) *J. Physiol.* **192**, 217–235.
37. I. M. Glynn and S. J. D. Karlish (1990) *Ann. Rev. Biochem.* **59**, 171–205.
38. H. Matsui and H. Homareda (1982) *J. Biochem. (Tokyo)* **92**, 193–217.
39. S. J. D. Karlish, D. W. Yates, and I. M. Glynn (1978) *Biochim. Biophys. Acta* **525**, 252–264.
40. J. C. Skou (1982) *Biochim. Biophys. Acta* **688**, 369–380.
41. M. Esmann (1991) *Biochim. Biophys. Acta* **1064**, 31–36.
42. F. Cornelius and J. C. Skou (1985) *Biochim. Biophys. Acta* **818**, 211–221.
43. H. Homareda, T. Ishii, and K. Takeyasu (2000) *Eur. J. Pharmacol.* **400**, 177–183.

## Oligopeptide

An alternative term for a [peptide](#), an oligopeptide is an [oligomer](#) or short [polymer](#) consisting of a

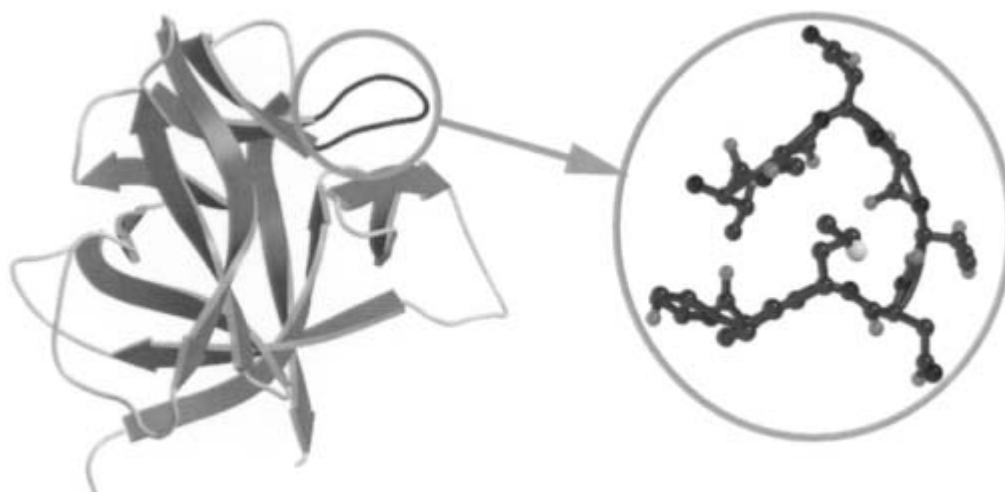
small number (usually between 2 and 50) of linked [amino acids](#). Oligopeptides have a well-defined sequence of amino acids but may or may not have a well-defined three-dimensional structure (see [Protein Structure](#)). Many oligopeptides, such as [hormones](#) and [neuropeptides](#), have important biological functions.

[See also [Peptide](#).]

## Omega Loop

The W-loop describes a common type of nonregular **secondary structure** observed in [protein structures](#) that is characterized by a loop-shaped [backbone](#) conformation resembling the Greek letter omega (W). The W-loop consists of six or more amino acid **residues** that cause a reversal in the direction of the [polypeptide chain](#) of the protein, so that the two ends of the loop are close in space (Fig. 1). Unlike the regular secondary structure types—**a-helix** and **b-strand**—the omega loop does not have repeating backbone [dihedral angles](#). Also, there are no regular repeating patterns of [hydrogen bonds](#) like those of the a-helix, although omega loops do form hydrogen bonds. Furthermore, their structures are compact and globular, with the [side chains](#) often packing into the center of the loop. In some cases, omega loops have specific functions, including substrate binding, [catalysis](#), or molecular recognition. Examples of omega loops having specific functions are (a) the phosphate-binding loop (**P-loop**) of **nucleotide-binding proteins** and (b) the **hypervariable regions** of **antibodies** that interact with **antigens**.

**Figure 1.** The structure of **interleukin-1b** (1) is depicted, with b-strands shown as arrows and connecting loops and turns as thin ribbons. An example of an W-loop, as defined by Leszczynski and Rose (2), from interleukin-1b is shown circled. **(Right)** The atomic structure of this omega loop is shown in detail. This figure was generated using Molscript (3) and Raster3D (4, 5).



## Bibliography

1. J. P. Priestle, H.-P. Schar, and M. G. Grutter (1989) Proc. Natl. Acad. Sci. USA **86**, 9667–71.
2. J. F. Leszczynski and G. D. Rose (1986) Science **234**, 849–855.
3. P. J. Kraulis (1991) J. Appl. Crystallogr. **24**, 946–950.

4. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
5. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

6. C. S. Ring, D. G. Kneller, R. Langridge, and F. S. Cohen (1992) Taxonomy and conformational analysis of loops in proteins. *J. Mol. Biol.* **224**, 685–699 (Alternative classification of loop structures in proteins.)
7. J. S. Fetrow (1995) Omega loops: nonregular secondary structures significant in protein function and stability. *FASEB J.* **9**, 708–717.

## OmpA

**Gram-negative** bacteria such as *Escherichia coli* are surrounded by two lipid bilayer [membranes](#), which are separated by the murein-containing **periplasmic** space. About 50% of the mass of the outer membrane consists of protein (1). The most abundant proteins include OmpA, Braun's lipoprotein, and the **porins**. [Transcription](#) of the *ompA* **gene** leads to a very stable [messenger RNA](#) (mRNA), whose degradation pathway is beginning to be understood. The OmpA protein is synthesized as a precursor **pre-protein** that is secreted into the periplasmic space. After proteolytic processing of the [signal peptide](#), the protein is assembled into the outer membrane. As one of the first proteins characterized with respect to secretion and outer membrane localization, OmpA has become a paradigm for the study of the underlying processes. The one-domain membrane moiety of this relatively small and monomeric protein represents a very useful model for the investigation of membrane assembly of  $\beta$ -structure membrane proteins *in vivo* as well as *in vitro*.

### 1. Discovery

The OmpA protein and *ompA* gene were discovered in the 1970s, independently by a number of different groups using different methods. Nomenclature became a problem. Although gene designations reflected different functions (see below) of the gene product [eg, *con* (2), *tolG* (3), *tut* (4)], the detection and biochemical characterization of several major proteins in outer membrane preparations led to arbitrary names in different laboratories: outer membrane protein B (5),  $\Pi^*$  (1), 7 (6), C (7), 3a (8), d (9), O-10 (10). It was finally agreed to refer to the genetic locus as the outer membrane protein A (*ompA*) since it represented the first locus of a structural gene for a major outer membrane protein to be established (11). The *ompA* gene was eventually mapped at 22 min (near *pyrD* and *fabA*) on the *E. coli* chromosome (12, 4); it was subsequently **cloned** (13, 14), and both the gene (14, 15) and the protein (16) were **sequenced**.

### 2. Expression

The *ompA* gene is constitutively expressed, leading to about  $10^5$  protein molecules per cell (1). This high expression level is reflected by the tendency to use codons that are recognized by major isoaccepting species of [transfer RNA](#) (14). Since the amount of OmpA protein in the outer membrane remains constant over a wide range of generation times, a mechanism for adjusting the rate of biosynthesis is required. In fact, there is a 2.5-fold decrease of OmpA synthesis at slow growth rates compared to the amount produced at rapid growth (17). This decrease is caused by a reduced stability of the *ompA* mRNA. In rapidly growing cells at 30°C, the *ompA* transcript has an unusually long half-life of about 15 min (18), which drops to only 4 min during slow growth (19). It

was found that a 134-nucleotide 5'-untranslated region (UTR) (20, 21) is the determinant responsible for stabilization of the mRNA (22, 23). The control exerted by this *ompA* element is under regulation by the growth rate (24, 25). The 5'-UTR serves as a target for a rate-limiting endonucleolytic cleavage event, which represents the initial step in *ompA* mRNA degradation (26, 27). Detailed studies of this element indicate that it forms a highly structured secondary structure essential for its activity as a stability determinant (28-30). Also, synonymous codon substitutions among 12 enterobacterial species do not occur at random along the *ompA* coding sequence (31), which may reflect selection for a conserved mRNA secondary structure. The *ompA* 5'-UTR can also act in *cis* to stabilize heterologous downstream RNA (23, 24, 32), suggesting that this element blocks a major mRNA degradation pathway in *E. coli*, which probably involves **ribonuclease** (RNase) E (24, 27, 33). However, the exact degradation pathway of *ompA* transcripts is still unknown. In a triple mutant that cannot synthesize **polynucleotide phosphorylase**, RNase E, and RNase II, the half-life of the *ompA* message was only barely changed (34). In contrast, a further deletion mutation in *pcnB*, encoding the poly(A) polymerase PAP I, dramatically increased the messenger half-life (2.5-fold) and influenced the decay pattern, but only in combination with the other three mutations. By itself, this *pcnB* mutation had no influence on the *ompA* transcript (34). This indicates that both termini of the transcript are important for its turnover and multiple degradation pathways exist that may or may not depend on RNase E and/or PAP I.

Cells can regulate the protein occupancy of the outer membrane. For example, **merodiploidy** for the *ompA* gene did not cause increased OmpA protein as a result of a gene dosage effect (11). Overproduction of the OmpA protein inhibited synthesis of the porins OmpC and OmpF, and its absence led to an increased concentration of these porins (and vice versa). This effect is specific for outer membrane proteins, ie, the level of periplasmic proteins remained unaffected. There are conflicting results concerning the mechanism of occupancy regulation. Initially, it was concluded that a feedback mechanism operated at the level of transcription. It was later found that overproduction of OmpC causes an almost immediate decrease (about six-fold) in the stability of *ompA* mRNA, indicating that feedback control influences mRNA decay rates. Interestingly, this cross regulation operates only for a subset of outer membrane proteins; many of those (unidentified) present at low concentration were not affected. Hence, there should be a property common to the relevant mRNA allowing them to be recognized as belonging to this subset [for more details, see ref. 35].

The *ompA* gene is adjacent to the *sulA* gene, which encodes an **SOS response**-inducible inhibitor of cell division. Both genes are arranged in tandem in the order *sulA* - *ompA*. Upon induction of the SOS response, a novel hybrid transcript encoding both the Sula and OmpA proteins is produced, resulting from transcription reading through the *sulA* terminator and the intergenic region (36). At the same time, a decrease in the abundance of the authentic *ompA* mRNA was observed as a result of **promoter** occlusion. However, the amount of the OmpA protein produced was almost the same. Finally, the turnover of the *ompA* mRNA displays a strong dependence on growth temperature; at 37°C the messenger exhibits a half-life of only 4.5 min (37).

### 3. Function

OmpA is involved in several processes, and mutations that interfere with some of them were isolated and reflected in their genetic nomenclature. It acts as a receptor for several **T-even-type bacteriophages** (*tut* mutants) (4), the best known of them being phages K3, M1, Ox2, and TuII\* (38). **Colicins** K and L need OmpA for their entry into the susceptible cell (*tolG* mutants) (3). In these cases, OmpA is parasitized by toxic agents, and thus these interactions do not, of course, represent normal physiological functions. OmpA is also essential for the formation of mating aggregates in liquid media during **F-mediated conjugation** (*con* mutants) (2). Furthermore, OmpA has been reported to form pores *in vitro*, with an estimated pore diameter of ~1nm judged by liposome swelling assays (39), or ~0.6 to 0.7 nm by black lipid bilayer measurements (40). However, using different OmpA preparations, no ion conductivity could be measured in two other laboratories

(41, 42). Only a small fraction of the molecules (2 to 3%) were observed to contain an open channel (43), and the contribution of OmpA to the overall permeability of the outer membrane of *E. coli* is expected to be insignificant, as the major porins such as OmpF and OmpC would make overwhelmingly large contributions to permeability.

The strongest evidence concerning the physiological function of OmpA came from the observation that double mutants in *ompA* and *lpp* (encoding Braun's lipoprotein, another highly abundant outer membrane protein) lacking both proteins had lost their rod shape and required increased concentrations of electrolytes for optimal growth (44). Null mutations of *ompA* or *lpp* alone do not have any phenotype under normal laboratory conditions. [Electron microscopy](#) (EM) revealed that the murein layer was no longer associated with the outer membrane in the double mutant (44), and OmpA in wild-type cells could be **crosslinked** to the murein by the **bifunctional reagent** dithio-bis-succinimidyl-propionate (45, 46). Furthermore, regions in OmpA have been identified that are **homologous** to several other proteins known or proposed to interact specifically with the peptidoglycan, eg, the peptidoglycan-associated lipoprotein or MotB protein of **flagella** motors (47, 48). Hence, OmpA is an important structural component of the bacterial cell surface, providing a physical linkage between the outer membrane and murein sacculus.

An *ompA* mutant showed attenuated virulence in two different types of *E. coli* K-1 model infection (embryonic chicken, neonatal rat), implicating OmpA in the pathogenicity of virulent *E. coli* (and possibly other enterobacteria) (49). The mutant was also more sensitive to the bactericidal effect of human **serum** by the classical pathway of **complement** activation, suggesting that OmpA contributes to *E. coli* K-1 pathogenesis by a mechanism that may involve increased serum resistance (49). The invasive capability of the **wild-type** *E. coli* K-1 toward brain microvascular endothelial cells was 25- to 50-fold greater than that of an isogenic mutant lacking OmpA, and the invasion of the wild-type *E. coli* K-1 was inhibited by purified OmpA proteins and polyclonal anti-OmpA **antibodies**. Therefore, OmpA may also be involved in causing meningitis in newborns by helping in the traversal of *E. coli* K-1 across the blood-brain barrier (50).

#### 4. Structure

The 325-residue OmpA of *E. coli* is assumed to exist as a monomer in the outer membrane; oligomers have thus far not been observed. Although microcrystalline preparations of OmpA have been available since 1984 (51), its three-dimensional structure has not yet been solved. Nevertheless, a detailed model of its membrane topology has been proposed by means of comparative [sequence analysis](#) (52), phage-mapping (53, 54), linker insertion [mutagenesis](#) in combination with [proteinase digestion](#) experiments (55, 56), spectroscopic analyses [[circular dichroism](#) (CD) and [vibrational spectroscopy](#)] (57, 58), and computer-based predictions (58). This model predicts a two-domain organization of OmpA, with the 170 N-terminal residues residing in the outer membrane, in the form of an eight-stranded antiparallel b-barrel, and the remaining 155 C-terminal residues extending into the periplasmic space. The eight antiparallel b-strands are connected by relatively large and hydrophilic surface-exposed loops (L1 to L4) and short periplasmic turns (T1 to T3). Such an up-and-down b-barrel topology appears to be a common architectural principle of integral outer membrane proteins, as observed by [X-ray crystallography](#) for the [porin](#) family (59-61). The following lines of evidence led to this model:

1. The transmembrane nature of the surface-exposed OmpA protein was first demonstrated by its chemical cross-linking to the periplasmic murein layer in intact cells (45).
2. Proteinases acting on isolated cell envelopes degraded the OmpA protein (35 kDa) to smaller proteinase-protected fragments with sizes of 19 kDa ([pronase](#)) or 24 kDa ([trypsin](#)) (1). These fragments were still active as phage receptors and in inhibiting F-mediated conjugation (which was not, however, tested for the 19-kDa fragment). They corresponded to the N-terminal region of OmpA (62), and the pronase cleavage site was mapped at amino acid residue 177 (16). Also, the *in vivo* membrane assembly of C-terminally truncated OmpA variants (up to 193 residues) was not

affected and led to functional polypeptides (receptor for phages and colicin, F-mediated conjugation). From these results, it was concluded that OmpA consists of two domains: an N-terminal transmembrane domain (residues 1 to about 180) and a C-terminal domain (residues 178–325) that is exposed in the periplasm (63).

3. Mutations that rendered cells resistant to phage mapped at four regions around residues 25, 70, 110, and 154. It was suggested that these regions are located at the cell surface where they interact with the phages (53, 54).

4. Comparative sequence analysis of five enterobacterial OmpA proteins revealed that regions with high sequence variability coincided with those regions where phage-resistance mutations had been mapped. This finding was interpreted as further evidence of their surface localization (52).

5. Spectroscopic analyses (infrared, CD, and Raman spectroscopy) indicated a high [beta-pleated sheet](#) content for the protein (55%) (57, 58). When the N-terminal domain was analyzed separately, a complete absence of [alpha-helix](#) and an even greater fraction of (mainly antiparallel)  $\beta$ -structure (68%) was observed (58). This domain is resistant to [sodium dodecyl sulfate](#) (SDS)-induced denaturation at room temperature (64), causing an aberrant migration behavior of the OmpA protein in **SDS-PAGE**; its apparent molecular mass depends on the temperature used during preparation of the sample. Prolonged heating in SDS results in a conversion from ~30 to ~35 kDa; this phenomenon has been called “heat modifiability” (65). Since OmpA adopts its heat-stable  $\beta$ -structure only during membrane assembly, demonstration of heat modifiability is sufficient to prove membrane assembly of mutant OmpA variants. However, mutants exist that assemble into the membrane but are not heat-modifiable.

6. With an antiparallel  $\beta$ -structure of the N-terminal domain, assumed, computer-based predictions suggest a detailed model of its membrane topology (58).

7. Linker insertion mutagenesis in combination with proteinase digestion experiments confirmed the location of two surface-exposed loops (55) and of all three periplasmic turns (56) of the N-terminal domain in the polypeptide chain.

8. In support of the two-domain structure of OmpA, a C-terminally truncated variant, consisting of only the N-terminal  $\beta$ -barrel domain (171 residues), was shown to assemble into the outer membrane as efficiently as the wild-type protein (56). This fragment, which was also heat-modifiable, was still active as a phage receptor (56) and in F-mediated conjugation (66).

9. Sequences homologous to the C-terminal 150 residues were found in numerous proteins, including various lipoproteins and the MotB protein of the flagella motor of both Gram-negative as well as Gram-positive bacteria, supporting its organization as a separate protein domain (47, 48).

10. Circularly permuted variants of the N-terminal domain were constructed at the DNA level and shown to assemble functionally into the outer membrane (67). This indicates the close proximity of both ends of the polypeptide chain in the folded structure and supports the  $\beta$ -barrel model.

11. Two short synthetic peptides corresponding to regions of the first two surface-exposed loops of OmpA (a hexamer, Asn27-Glu32, and a pentamer, Gly65-Asn69) significantly inhibited invasion of brain microvascular endothelial cells by *E. coli* (50), a result that could be easily explained by their surface location.

12. All four predicted surface-exposed loops (58) of the N-terminal domain could be shortened to three residues without interfering with membrane assembly (66). However, all four loops were necessary for efficient F-mediated conjugation in liquid media, and only the fourth loop was dispensable for functioning as a receptor for phage K3 (66).

Based on its transmembrane nature, OmpA was successfully used as a surface display vehicle (see [Phage Display Libraries](#)). Peptides of between 20 and 80 amino acids in length were inserted into all but the first surface-exposed loop and shown to be accessible from outside the cell ([55](#), [68](#), [69](#)). OmpA was also used for surface exposure of complete proteins, either as tripartite fusions consisting of parts of the Lpp lipoprotein, OmpA, and the target protein, eg, **b-lactamase** ([70](#), [71](#)), ScFv (single-chain antibody Fv fragment) ([72](#)), [alkaline phosphatase](#) ([73](#)), or as a sandwich fusion of influenza hemagglutinin into surface loop L3 ([74](#)). The target proteins could be detected at the cell surface. However, the exact arrangement of the polypeptide chain in the outer membrane was not analyzed, and it is possible that they adopted rather aberrant structures. This is most likely the case with the tripartite fusions, since they included only five ([70](#), [72](#), [73](#)) or fewer ([71](#)) of the proposed transmembrane strands of OmpA.

An alternative OmpA model predicts 16 transmembrane segments that are distributed throughout the polypeptide chain, forming a b-barrel-like structure ([75](#)). There are, among others, two main criticisms of this model. First, C-terminally truncated OmpA derivatives that would lack some of the predicted transmembrane segments can assemble into the membrane, which is incompatible with such an arrangement of the polypeptide chain. Second, speculations about an  $\alpha$ -helical conformation of some transmembrane segments that were put forward to match the observed secondary structure composition contradict the prediction procedure, which was based on the detection of amphipathic b-strands.

## 5. Folding

The OmpA protein serves as a well-established model for the study of membrane assembly of b-structure proteins *in vivo* as well as *in vitro*. OmpA is synthesized as a precursor protein, pre-OmpA, which is then secreted across the cytoplasmic membrane with the assistance of its 21-residue signal peptide ([76](#), [77](#)) and the **Sec** machinery. An influence of the cytosolic SecB **molecular chaperone** in this process, as detected by a decreased processing rate of the signal sequence by the Lep leader peptidase during pulse-chase experiments, was observed with a *secB* null mutant ([78](#)) and when SecB was titrated with a secretion-defective mutant of **maltose-binding protein** (MBP) ([79](#)). Furthermore, binding of pre-OmpA to SecB has been demonstrated both *in vitro* ([80](#)) and *in vivo* ([81](#), [82](#)). As one of the first secretory proteins to be identified and cloned, OmpA became a well-established model for the study of protein secretion *in vitro* ([83-88](#)) as well as *in vivo* ([89-91](#)), and its signal sequence has often been used for the export of heterologous proteins. Several secretion vectors have been developed for this purpose ([92](#)).

After being translocated across the cytoplasmic membrane, a folding intermediate, imp-OmpA (immature processed OmpA), has been identified that was not yet assembled in the outer membrane and lacked heat modifiability ([93](#)). Sucrose gradient centrifugation located this species at the cytoplasmic membrane; when overproduced, it was found to be associated with the outer membrane. For such outer membrane association, whose nature has remained unknown, a region between residues 154 and 180 was found to be responsible ([94](#), [95](#)). This region includes the last transmembrane b-strand, with a conserved aromatic residue at its C-terminal position ([96](#)). There is, however, no sequence specificity for the last b-strand to function in assembly; it only needs to be [amphipathic](#) or [hydrophobic](#), and it must not have a proline at the center ([97](#)). For incorporation into the outer membrane, OmpA does not depend on newly-synthesized lipopolysaccharides (LPS) ([98](#)). Instead, OmpA appears to associate with LPS preexisting in the outer membrane, and this association is required for phage receptor activity ([98](#)).

The b-barrel domain of OmpA is surprisingly tolerant toward changes of its amino acid sequence without preventing membrane insertion. Insertions of small peptides at all periplasmic turns ([56](#)), as well as the last three surface-exposed loops (the first loop was never tested), were well tolerated ([55](#), [68](#), [69](#)), as was drastic shortening of all surface-exposed loops to only three amino acid residues ([66](#)). Also, substitution by [leucine](#) of the turn-promoting [proline](#) residues that are present at all three periplasmic turns had no effect on membrane assembly ([99](#)). Apparently, neither turns nor loops

possess any topogenic information, which then must reside in the transmembrane b-strands. But even they exhibit a surprising tolerance toward mutagenic alterations. The last b-strand can be replaced by an arbitrary hydrophobic or amphiphilic sequence (97), and proline residues, which are generally avoided within transmembrane b-strands, can be introduced at several positions (100). Probably the information for membrane assembly is redundant and the b-barrel can accommodate some structural distortions without being affected in its general fold. The possibility of creating circularly permuted (67) and split (101) variants (by co-expressing pairs of complementary protein fragments) that assemble into the outer membrane *in vivo* underscores the robustness of the structure and of the folding process of the b-barrel domain of OmpA. In contrast, much less is known about the C-terminal periplasmic domain.

SDS-denatured OmpA could be refolded to its native conformation *in vitro* in the presence of LPS, but not synthetic dimyristoylphosphatidylcholine (DMPC) or total *E. coli* phospholipid (62). However, such denatured OmpA could readily assemble into DMPC vesicles when octylglucoside was first added to the protein. In this case, the protein was already refolded in the SDS-octylglucoside micelles and could be transferred into DMPC vesicles (102). In fact, the protein could also be assembled directly into such vesicles under conditions that tried to mimic what happens *in vivo* (103). In this case, the vesicles had to be very small, ie, highly curved, and the lipid had to be in the fluid state. When urea-denatured OmpA was added to such vesicles, it refolded and inserted in an oriented manner. This was also achieved with the tryptic fragment of OmpA that lacks most of its periplasmic domain. When the experiment was performed at a temperature below the lipid phase transition, the protein also folded and adsorbed to the vesicles but was not incorporated; raising the temperature also led in this case to incorporation. This polypeptide that is adsorbed possessed about the same content of b-structure (as judged by spectroscopic methods) as the protein incorporated into the vesicles (or present in the outer membrane) but lacked heat modifiability. Hence, this intermediate appears to correspond to the imp-OmpA observed *in vivo*. These results show that a lipid bilayer is not required for OmpA to assume its native conformation and it spontaneously refolds and inserts into a membrane, provided the latter possesses some defects (suboptimal packing of lipid in small vesicles or the presence of octylglucoside in large vesicles). From a study of the *in vitro* refolding kinetics, three transitions along OmpA's folding pathway could be distinguished (104). Their characteristic times are (1) less than a second, (2) in the range of minutes, and (3) in the range of an hour. These *in vitro* experiments appear to suggest that there are no helpers involved in the membrane assembly of OmpA. However, *in vitro* assembly is much less efficient than *in vivo* assembly; therefore, these experiments do not exclude the involvement of helper factors *in vivo*. The following proteins appear to participate in OmpA biogenesis: (1) The DsbA/DsbB system (see [Protein Disulfide Isomerase PDI](#)) is involved in the formation of the intramolecular disulfide bond that is located in the periplasmic domain (105, 106). (2) [Peptidyl-prolyl cis-trans isomerases](#) could be involved in potential [peptide bond](#) isomerizations, as OmpA contains several proline residues, although their isomeric state is unknown. A *rotA* mutant was not impaired in OmpA assembly (107), but SurA, another periplasmic isomerase, was shown to assist in the folding of OmpA (108). (3) The Skp protein, which was shown to bind selectively to OmpA and several other outer membrane proteins *in vitro*, is another candidate molecular chaperone along the route of OmpA to the outer membrane (109).

Note added in proof

The three-dimensional structure of the OmpA transmembrane domain has been solved by X-ray crystallography [A. Pautsch and G. E. Schulz (1998) *Nature Struct. Biol.* **5**, 1013–1017].

#### Bibliography

1. U. Henning, B. Höhn, and I. Sonntag (1973) *Eur. J. Biochem.* **39**, 27–36.
2. R. A. Skurray, R. E. W. Hancock, and P. Reeves (1974) *J. Bacteriol.* **119**, 726–735.
3. T. Chai and J. Foulds (1974) *J. Mol. Biol.* **85**, 465–474.
4. U. Henning, I. Hindennach, and I. Haller (1976) *FEBS Lett.* **61**, 46–48.



5. P. D. Bragg and C. Hou (1972) *Biochim. Biophys. Acta* **274**, 478–488.
6. M. Inouye and M. Yee (1973) *J. Bacteriol.* **113**, 304–312.
7. J. Koplów and H. Goldfine (1974) *J. Bacteriol.* **117**, 527–543.
8. C. A. Schnaitman (1974) *J. Bacteriol.* **118**, 442–453.
9. B. Lugtenberg, J. Meijers, R. Peters, P. van der Hoek, and L. van Alphen (1975) *FEBS Lett.* **58**, 254–258.
10. J. Uemura and S. Mizushima (1975) *Biochim. Biophys. Acta* **413**, 163–176.
11. D. B. Datta, C. Krämer, and U. Henning (1976) *J. Bacteriol.* **128**, 834–841.
12. J. Foulds (1974) *J. Bacteriol.* **117**, 1354–1355.
13. U. Henning, H.-D. Royer, R. M. Teather, I. Hindennach, and C. P. Hollenberg (1979) *Proc. Natl. Acad. Sci. USA* **76**, 4360–4364.
14. N. R. Movva, K. Nakamura, and M. Inouye (1980) *J. Mol. Biol.* **143**, 317–328.
15. E. Beck and E. Bremer (1980) *Nucl. Acid. Res.* **8**, 3011–3024.
16. R. Chen, W. Schmidmayr, C. Krämer, U. Chen-Schmeisser, and U. Henning (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4592–4596.
17. U. Lundberg, G. Nilsson, and A. von Gabain (1988) *Gene* **72**, 141–149.
18. A. von Gabain, J. G. Belasco, J. L. Schottel, A. C. Y. Chang, and S. N. Cohen (1983) *Proc. Natl. Acad. Sci. USA* **80**, 653–657.
19. G. Nilsson, J. G. Belasco, S. N. Cohen, and A. Von Gabain (1984) *Nature* **312**, 75–77.
20. N. R. Movva, K. Nakamura, and M. Inouye (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3845–3849.
21. S. T. Cole, E. Bremer, I. Hindennach, and U. Henning (1982) *Mol. Gen. Genet.* **188**, 472–479.
22. P. J. Green and M. Inouye (1984) *J. Mol. Biol.* **176**, 431–442.
23. J. G. Belasco, G. Nilsson, A. von Gabain, and S. N. Cohen (1986) *Cell* **46**, 245–251.
24. S. A. Emory and J. G. Belasco (1990) *J. Bacteriol.* **172**, 4472–4481.
25. D. Georgellis, S. Arvidson, and A. von Gabain (1992) *J. Bacteriol.* **174**, 5382–5390.
26. O. Melefors and A. von Gabain (1988) *Cell* **52**, 893–901.
27. U. Lundberg, A. von Gabain, and O. Melefors (1990) *EMBO J.* **9**, 2731–2741.
28. L. H. Chen, S. A. Emory, A. L. Bricker, P. Bouvet, and J. G. Belasco (1991) *J. Bacteriol.* **173**, 4578–4586.
29. S. A. Emory, P. Bouvet, and J. G. Belasco (1992) *Genes Devel.* **6**, 135–148.
30. V. Rosenbaum, T. Klahn, U. Lundberg, E. Holmgren, A. von Gabain, and D. Riesner (1993) *J. Mol. Biol.* **229**, 656–670.
31. J. G. Lawrence, D. L. Hartl, and H. Ochman (1991) *J. Mol. Evol.* **33**, 241–250.
32. M. J. Hansen, L. H. Chen, M. L. S. Fejzo, and J. G. Belasco (1994) *Mol. Microbiol.* **12**, 707–716.
33. E. A. Mudd and C. F. Higgins (1993) *Mol. Microbiol.* **9**, 557–568.
34. E. B. O'Hara, J. A. Chekanova, C. A. Ingle, Z. R. Kushner, E. Peters, and S. R. Kushner (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1807–1811.
35. U. Henning and R. Koebnik (1994) In *Bacterial Cell Wall* (J. M. Ghuyssen and R. Hakenbeck, eds.), Elsevier Science B.V., Amsterdam, Netherlands, pp. 381–395.
36. S. T. Cole and N. Honoré (1989) *Mol. Microbiol.* **3**, 715–722.
37. A. Deana, R. Ehrlich, and C. Reiss (1996) *J. Bacteriol.* **178**, 2718–2720.
38. H. Schwarz, I. Riede, I. Sonntag, and U. Henning (1983) *EMBO J.* **2**, 375–380.
39. E. Sugawara and H. Nikaido (1992) *J. Biol. Chem.* **267**, 2507–2511.
40. N. Saint, E. De, S. Julien, N. Orange, and G. Molle (1993) *Biochim. Biophys. Acta* **1145**, 119–123.

41. R. Benz, unpublished results.
42. P. S. Phale, unpublished results.
43. E. Sugawara and H. Nikaido (1994) *J. Biol. Chem.* **269**, 17981–17987.
44. I. Sonntag, H. Schwarz, Y. Hirota, and U. Henning (1978) *J. Bacteriol.* **136**, 280–285.
45. R. Endermann, C. Krämer, and U. Henning (1978) *FEBS Lett.* **86**, 21–24.
46. M. Leduc, K. Ishidate, N. Shakibai, and L. Rothfield (1992) *J. Bacteriol.* **174**, 7982–7988.
47. R. Demot and J. Vanderleyden (1994) *Mol. Microbiol.* **12**, 333–336.
48. R. Koebnik (1995) *Mol. Microbiol.* **16**, 1269–1270.
49. J. N. Weiser and E. C. Gotschlich (1991) *Infect. Immun.* **59**, 2252–2258.
50. N. V. Prasadarao, C. A. Wass, J. N. Weiser, M. F. Stins, S. H. Huang, and K. S. Kim (1996) *Infect. Immun.* **64**, 146–153.
51. R. M. Garavito, U. Hinz, and J.-M. Neuhaus (1984) *J. Biol. Chem.* **259**, 4254–4257.
52. G. Braun and S. T. Cole (1984) *Mol. Gen. Genet.* **195**, 321–328.
53. R. Morona, M. Klose, and U. Henning (1984) *J. Bacteriol.* **159**, 570–578.
54. R. Morona, C. Krämer, and U. Henning (1985) *J. Bacteriol.* **164**, 539–543.
55. R. Freudl (1989) *Gene* **82**, 229–236.
56. G. Ried, R. Koebnik, I. Hindennach, B. Mutschler, and U. Henning (1994) *Mol. Gen. Genet.* **243**, 127–135.
57. K. Nakamura and S. Mizushima (1976) *J. Biochem. Tokyo* **80**, 1411–1422.
58. H. Vogel and F. Jähnig (1986) *J. Mol. Biol.* **190**, 191–199.
59. M. S. Weiss, U. Abele, J. Weckesser, W. Welte, E. Schiltz, and G. E. Schulz (1991) *Science* **254**, 1627–1630.
60. S. W. Cowan, T. Schirmer, G. Rummel, M. Steiert, R. Ghosh, R. A. Pauptit, J. N. Jansonius, and J. P. Rosenbusch (1992) *Nature* **358**, 727–733.
61. T. Schirmer, T. A. Keller, Y. F. Wang, and J. P. Rosenbusch (1995) *Science* **267**, 512–514.
62. M. Schweizer, I. Hindennach, W. Garten, and U. Henning (1978) *Eur. J. Biochem.* **82**, 211–217.
63. E. Bremer, S. Cole, I. Hindennach, U. Henning, E. Beck, C. Kurz, and H. Schaller (1982) *Eur. J. Biochem.* **122**, 223–231.
64. E. Sugawara, M. Steiert, S. Rouhani, and H. Nikaido (1996) *J. Bacteriol.* **178**, 6067–6069.
65. M. G. Beher, C. A. Schnaitman, and A. P. Pugsley (1980) *J. Bacteriol.* **143**, 906–913.
66. R. Koebnik, submitted for publication.
67. R. Koebnik and L. Krämer (1995) *J. Mol. Biol.* **250**, 617–626.
68. A. Ruppert, N. Arnold, and G. Hobom (1994) *Vaccine* **12**, 492–498.
69. D. Haddad et al. (1995) *FEMS Immunol. Med. Microbiol.* **12**, 175–186.
70. J. A. Francisco, C. F. Earhart, and G. Georgiou (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2713–2717.
71. G. Georgiou, D. L. Stephens, C. Stathopoulos, H. L. Poetschke, J. Mendenhall, and C. F. Earhart (1996) *Protein Eng.* **9**, 239–247.
72. J. A. Francisco, R. Campbell, B. L. Iverson, and G. Georgiou (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10444–10448.
73. C. Stathopoulos, G. Georgiou, and C. F. Earhart (1996) *Appl. Microbiol. Biotechnol.* **45**, 112–119.
74. S. Pistor and G. Hobom (1990) *Res. Microbiol.* **141**, 879–881.
75. C. Stathopoulos (1996) *Prot. Sci.* **5**, 170–173.
76. N. R. Movva, K. Nakamura, and M. Inouye (1980) *J. Biol. Chem.* **255**, 27–29.

77. K. Gamon, W. Schmidmayr, and U. Henning (1980) *Nucl. Acid Res.* **8**, 3025–3027.
78. T. Watanabe, S. Hayashi, and H. C. Wu (1988) *J. Bacteriol.* **170**, 4001–4007.
79. D. N. Collier, V. A. Bankaitis, J. B. Weiss, and P. J. Bassford Jr. (1988) *Cell* **53**, 273–283.
80. S. Lecker et al. (1989) *EMBO J.* **8**, 2703–2709.
81. C. A. Kumamoto (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5320–5324.
82. C. A. Kumamoto and O. Francetic (1993) *J. Bacteriol.* **175**, 2184–2188.
83. M. B. Sankaram, D. Marsh, L. M. Gierasch, and T. E. Thompson (1994) *Biophys. J.* **66**, 1959–1968.
84. S. Chatterjee, D. Suci, R. E. Dalbey, P. C. Kahn, and M. Inouye (1995) *J. Mol. Biol.* **245**, 311–314.
85. K. Sato, H. Mori, M. Yoshida, M. Tagaya, and S. Mizushima (1997) *J. Biol. Chem.* **272**, 5880–5886.
86. J. Eichler and W. Wickner (1997) *Proc. Natl. Acad. Sci. USA* **94**, 5574–5581.
87. P. Fekkes, C. van der Does, and A. J. M. Driessen (1997) *EMBO J.* **16**, 6105–6113.
88. G. Matsumoto, T. Yoshihisa, and K. Ito (1997) *EMBO J.* **16**, 6384–6393.
89. R. Freudl, H. Schwarz, M. Degen, and U. Henning (1989) *J. Mol. Biol.* **205**, 771–775.
90. A. Kuhn, D. Kiefer, C. Kohne, H. Y. Zhu, W. R. Tschantz, and R. E. Dalbey (1994) *Eur. J. Biochem.* **226**, 891–897.
91. M. Pohlschröder, C. Murphy, and J. Beckwith (1996) *J. Biol. Chem.* **271**, 19908–19914.
92. J. Ghrayeb, H. Kimura, M. Takahara, H. Hsiung, Y. Masui, and M. Inouye (1984) *EMBO J.* **3**, 2437–2442.
93. R. Freudl, H. Schwarz, Y.–D. Stierhof, K. Gamon, I. Hindennach, and U. Henning (1986) *J. Biol. Chem.* **261**, 11355–11361.
94. R. Freudl, H. Schwarz, M. Klose, N. R. Movva, and U. Henning (1985) *EMBO J.* **4**, 3593–3598.
95. M. Klose, H. Schwarz, S. MacIntyre, R. Freudl, M.–L. Eschbach, and U. Henning (1988) *J. Biol. Chem.* **263**, 13291–13296.
96. M. Struyvé, M. Moons, and J. Tommassen (1991) *J. Mol. Biol.* **218**, 141–148.
97. M. Klose, F. Jähnig, I. Hindennach, and U. Henning (1989) *J. Biol. Chem.* **264**, 21842–21847.
98. G. Ried, I. Hindennach, and U. Henning (1990) *J. Bacteriol.* **172**, 6048–6053.
99. M. Klose, S. MacIntyre, H. Schwarz, and U. Henning (1988) *J. Biol. Chem.* **263**, 13297–13302.
100. M. Klose, A. Störko, Y. D. Stierhof, I. Hindennach, B. Mutschler, and U. Henning (1993) *J. Biol. Chem.* **268**, 25664–25670.
101. R. Koebnik (1996) *EMBO J.* **15**, 3529–3537.
102. K. Dornmair, H. Kiefer, and F. Jähnig (1990) *J. Biol. Chem.* **265**, 18907–18911.
103. T. Surrey and F. Jähnig (1992) *Proc. Natl. Acad. Sci. USA* **89**, 7457–7461.
104. T. Surrey and F. Jähnig (1995) *J. Biol. Chem.* **270**, 28199–28203.
105. J. C. Bardwell, K. McGovern, and J. Beckwith (1991) *Cell* **67**, 581–589.
106. D. Missiakas, C. Georgopoulos, and S. Raina (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7084–7088.
107. M. Kleerebezem, M. Heutink, and J. Tommassen (1995) *Mol. Microbiol.* **18**, 313–320.
108. S. W. Lazar and R. Kolter (1996) *J. Bacteriol.* **178**, 1770–1773.
109. R. Chen and U. Henning (1996) *Mol. Microbiol.* **19**, 1287–1294.

### Suggestions for Further Reading

110. R. Freudl, M. Klose, and U. Henning (1990). Export and sorting of the *Escherichia coli* outer

membrane protein OmpA. *J. Bioenerget. Biomembr.* **22**, 441–449.

111. U. Henning and R. Koebnik (1994). "Outer membrane proteins of *Escherichia coli*: Mechanism of sorting and regulation of synthesis". In *Bacterial Cell Wall* (J. M. Ghuyssen and R. Hakenbeck, eds.), Elsevier, Amsterdam, Netherlands, pp. 381–395.

## Oncogenes, Oncoproteins

In the past 25 years a major revolution has occurred in understanding the molecular basis of cancer. The long drive to understand how normal cells grow and contribute to the ordered development of an organism and the nature of genetic events that disrupt this ordered process and result in cancer has become a major focus of biological research. The explosion of knowledge from basic research on cancer has unified several disciplines of science and led to a clearer understanding of the mechanisms associated with normal cell growth, genetic events that produce defects in regulatory mechanisms that lead to aberrant growth and differentiation, and finally the molecular basis for disease processes, such as cancer.

The most important breakthrough that led to this scientific revolution was the discovery that a small set of **genes** (probably about 100) in mammalian and avian cells (which are estimated to contain approximately 200,000 genes) controls normal [development](#) of the organism. Aberrations introduced into these genes during the life span of an organism by **viral** infection or by exposure to chemical **carcinogens** leads to cancer. It is becoming very clear that a precise understanding of the structure of these genes and the mechanisms associated with their function will provide us with approaches to treat diseases, such as cancer, and also to correct genetically inherited forms of developmental defects. Therefore, we are at one of those exciting moments in science when many disciplines of scientific endeavor come together, and much that seemed impossibly obscure suddenly becomes clear.

It has been recognized for almost a century that cancer is a multistep process, which takes decades to develop. Because cancer cells grow at a much faster rate than normal cells and do not obey rules of ordered growth immediately suggested that cancer is an aberration of cell growth. This observation, combined with the fact that cancer is a multistep process, suggested the involvement of multiple genetic events of the growth deregulation seen in cancer cells. Research during the past decade has shown that these genetic events include activating a group of genes, termed "oncogenes," and inactivating another group of genes named "growth suppressor genes." By definition, oncogenes are genes that promote cell growth, and growth suppressor genes are those that block cell growth. A fine balance between the activities of these two groups of gene products dictates normal cell growth, and disrupting this balance provides a cell with a growth advantage that ultimately results in a neoplastic state.

### 1. Discovery of viral and cellular Oncogenes

Oncogenes were initially discovered by studying transforming [retroviruses](#) that produce tumors in animals. Several decades ago, it was observed that chickens die predominantly of cancer, suggesting that these animals have a genetic predisposition to this disease. This observation led to examining the tumors derived from these animals for a disease-transmitting agent, which led to the identification of retroviruses. In 1911 Peyton Rous at the Rockefeller Institute in New York first isolated a retrovirus from a spontaneous chicken sarcoma ([1](#)). He established the viral etiology of the tumor by demonstrating that extracts derived from the tumor could be filtered through membranes that retain bacteria and nevertheless induce cancer in animals injected with the filtered extracts. The virus

isolated by Rous has been named the [Rous sarcoma virus](#) in honor of its discoverer. Most interestingly, it was found that these viruses contain **RNA** as their genetic material and hence were called “RNA tumor viruses” or “retroviruses.” Later studies revealed that extracts from animal tumors often contain two types of RNA tumor viruses, which were termed acute transforming viruses and leukemia viruses. These viruses differ from each other in a number of properties, which are listed in Table 1. The most important biological difference between these two classes of viruses is their ability to induce tumors *in vivo* as a function of time. Although acute transforming viruses induce tumors in susceptible hosts in a very short time (1 to 2 weeks), the leukemia viruses require several months to years. A second important difference between the two classes of viruses is their ability to transform cells in tissue culture. The acute transforming viruses readily transform cells in tissue culture, whereas the leukemia viruses fail to do so. Finally, several of the acute transforming viruses (with the exception of Rous sarcoma virus) are replication-incompetent, whereas leukemia viruses replicate readily *in vitro* and *in vivo*.

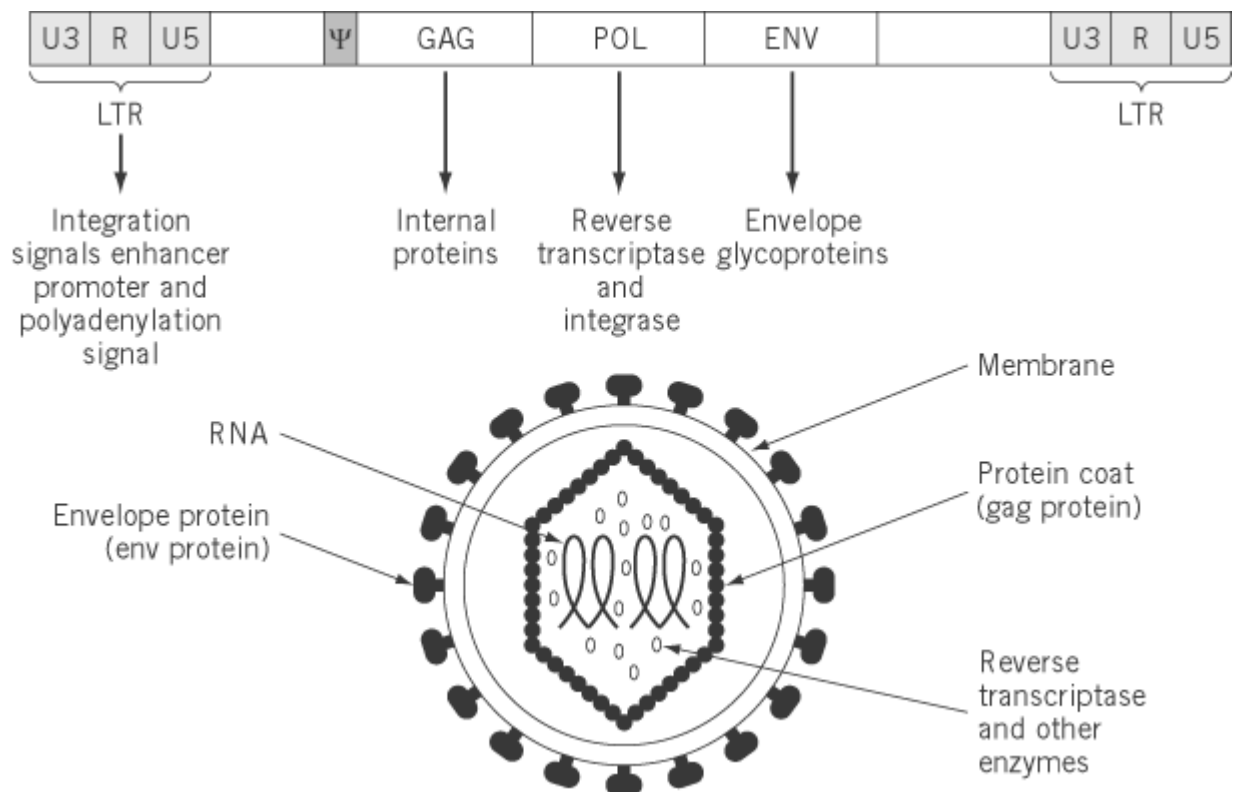
**Table 1. Retroviruses or Type C RNA Tumor Viruses**

| <b>Properties</b>                  | <b>Chronic Leukemia Viruses</b> | <b>Acute Transforming Viruses</b>                        |
|------------------------------------|---------------------------------|--|
| Induce                             | Lymphomas                       | Sarcomas, carcinomas, and hematopoietictumors            |
| Tumor latent period                | Months to years                 | Weeks  |
| Transformation of cells in culture | No                              | Fibroblasts and/orhematopoietic cells                    |
| Replication-competent              | Yes                             | Generally not  |
| Mechanism of transformation        | Integration next to oncogene    | Discrete cell-derived gene withinviral genome (oncogene) |

The advent of [recombinant DNA](#) techniques in late 1970s allowed the molecular [cloning](#) and [sequence analysis](#) of the two types of viruses, which provided a molecular explanation of their biological differences. These studies showed that the leukemia viruses encodes three polypeptides, named gag, pol and env. Of these, the *gag* gene encodes a [polypeptide chain](#) that is cleaved **proteolytically** into smaller polypeptides that form the core proteins of the virion (Fig. 1). The *pol* gene encodes a polypeptide chain that is cleaved into two polypeptides, one of which is the **reverse transcriptase** that enables the virus to convert its RNA into DNA. The second polypeptide functions as an [integrase](#) that allows the proviral DNA to integrate into the host viral genome. The *env* gene encodes the viral envelope protein that plays an important role in virus–host cell [membrane](#) interactions. In addition to these three genes, the retroviral [genome](#) contains a stretch of sequences that are duplicated at the 5' and 3' ends of the provirus, termed [Long Terminal Repeats](#) (LTRs). The LTRs contain sequences that constitute some of the most potent eukaryotic **promoter/enhancer** elements in mediating high-level [transcription](#) of viral RNA.

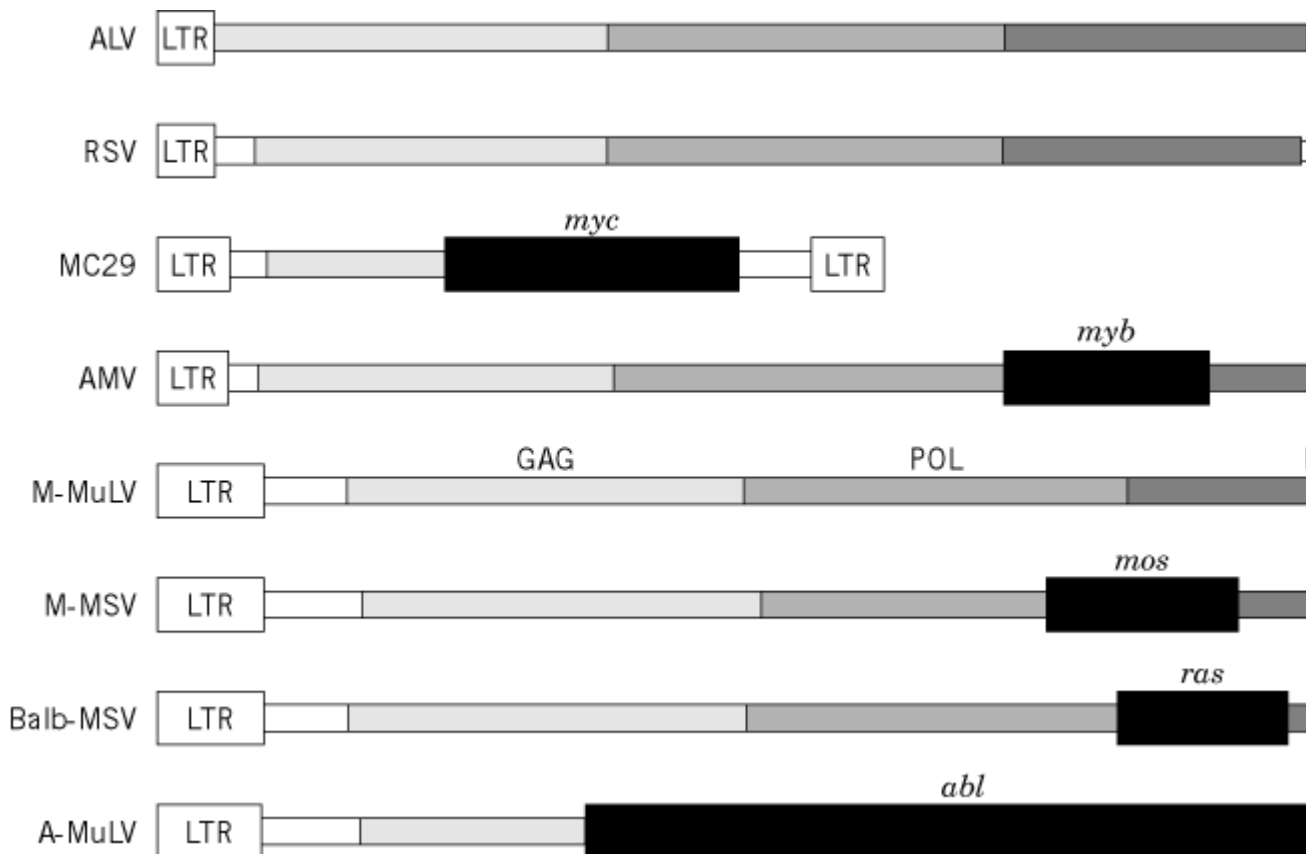
**Figure 1.** Genetic elements of retroviruses and their gene products. A typical retroviral provirus encodes three genetic elements, gag, pol, and env, that are translated into viral proteins. The gag sequences encode the viral core protein (an inner shell), the pol gene codes for reverse transcriptase (an enzyme that converts RNA into DNA) and the viral

integrase (an enzyme that facilitates the integration of viral sequences into the host genome), and the env gene product forms the viral envelope protein. The LTR sequences direct high levels of viral gene expression within the host cell.  $\Psi$  denotes the packaging signal that enables the viral proteins to recognize the viral RNA that is to be incorporated into the mature virion.



Comparison of the structures of the acute transforming and leukemia viral genomes revealed that the acute transforming viruses contain an additional segment of genomic RNA not present in leukemia viruses. Deletion of this sequence leads to loss of the transforming ability of these viruses. These observations suggested that a protein encoded by this unique piece of genetic material is responsible for inducing of tumors in animals. The first of these genes was originally found in the Rous sarcoma virus and was designated as the **src oncogene**. Structural analysis of several of the acute transforming viruses revealed that acute transforming viruses very often suffer deletions in their *env*, *pol*, and/or *gag* sequences, which explains their inability to replicate *in vitro* or *in vivo* (Fig. 2). The only exception to this rule is the Rous sarcoma virus, which contains all three structural genes of the avian leukosis virus, in addition to the *src* gene, and thus constitutes the only replication-competent acute transforming virus.

**Figure 2.** Comparison of the genomic structures of M-MuLV, ALV, and some of the acute transforming viruses. LTR: long terminal repeat.



An important breakthrough in cancer research came from the discovery that oncogenes of acute transforming viruses are in fact derived from normal cellular DNA and that this genetic information is transduced by acute transforming viruses via genetic [recombination](#). An analysis of normal cellular DNA by Stehlin and his co-workers (2), using a probe derived from the *src* oncogene (derived from Rous sarcoma virus), revealed the presence of endogenous oncogene-related sequences in normal chicken DNA. Following the lead provided by the Rous sarcoma virus, retrovirologists have isolated more than 200 different tumor-producing, acute transforming viruses from animal tumors of different species. Table 2 lists some of the acute transforming viruses and the oncogenes transduced by them. An examination of this table shows that several of the acute transforming viruses isolated from different animal species contain the same oncogene, suggesting that there are only a finite number of transforming genes in the avian and mammalian genomes and that these genes often undergo recombination with replicating retroviruses, leading to the formation of acute transforming viruses. Since the first identification of the *v-src* oncogene, a number of different approaches have been used to identify genes responsible for the altered growth properties of tumor cells. Now, the term “oncogene” is used more broadly to include any gene whose expression is associated with enhanced growth of tumor cells.

**Table 2. Retroviral Transforming Genes<sup>a</sup>**

| Gene       | Virus | Animal Origin  | Protein | Function       |
|------------|-------|----------------|---------|----------------|
| <i>src</i> | RSV   | Chicken        | p60     | Protein kinase |
|            | B77   | Chicken        | p60     | Protein kinase |
|            | rASV  | Chicken, quail | p60     | Protein kinase |
| <i>fps</i> | FSV   | Chicken        | p140    | Protein kinase |

|            |         |         |         |                      |
|------------|---------|---------|---------|----------------------|
|            | PRCII   | Chicken | p105    | Protein kinase       |
|            | UR 1    | Chicken | p150    | Protein kinase       |
|            | 16 L    | Chicken | p142    | Protein kinase       |
| <i>yes</i> | Y73     | Chicken | p90     | Protein kinase       |
|            | ESC     | Chicken | p80     | Protein kinase       |
| <i>ros</i> | UR2     | Chicken | p68     | Protein kinase       |
| <i>erb</i> | AEV     | Chicken | p75+p45 | Receptor             |
| <i>myc</i> | MC28    | Chicken | p110    | Transcription factor |
|            | CMII    | Chicken | p90     | Transcription factor |
|            | MH2     | Chicken | p100    | Transcription factor |
|            | OK10    | Chicken | p200    | Transcription factor |
| <i>jun</i> | ASV17   | Chicken | p55     | Transcription factor |
| <i>ski</i> | SK      | Chicken | p110    | Transcription factor |
| <i>rel</i> | AEV     | Turkey  | p56     | Transcription factor |
| <i>fos</i> | FBJ MSV | Mouse   | p55     | Transcription factor |
| <i>raf</i> | 3611MSV | Mouse   | p75     | Protein kinase       |
|            | MH2     | Chicken | p100    | Protein kinase       |
| <i>ras</i> | Ki-MSV  | Rat     | p21     | G-protein            |
|            | Ha-MSV  | Rat     | p21     | G-protein            |
|            | Ra-RaSV | Rat     | p29     | G-protein            |
| <i>abl</i> | Ab-MuLV | Mouse   | p20     | Protein kinase       |
| <i>fes</i> | ST-FeSV | Cat     | p85     | Protein kinase       |
|            | GA-FeSV | Cat     | p110    | Protein kinase       |
| <i>fgr</i> | GA-FeSV | Cat     | p70     | Protein kinase       |
| <i>fms</i> | MS-FeSV | Cat     | p170    | Protein kinase       |
| <i>sis</i> | SISV    | Monkey  | p28     | Growth factor        |

<sup>a</sup> Representative oncogenes, their transforming viruses, animal origin, molecular size, and function of the oncoproteins are presented here.

## 2. Mechanisms of Oncogene Activation by Retroviruses

Since this discovery of the *src* oncogene, a number of investigators have demonstrated the presence of cellular homologues to retrovirally associated oncogenes, which today number approximately about 50. These studies demonstrated clearly that all retroviral oncogenes are derived from normal cellular DNA. A detailed comparison of the structure of the viral oncogenes with their cellular counterparts revealed that several of the viral genes often represent mutant versions of their cellular homologues. To distinguish the viral and cellular versions of these genes, the oncogenes present in normal cellular DNA were termed proto-oncogenes or c-oncogenes (c = cell), and those in RNA tumor viruses were termed v-oncogenes (v = virus). Several of the proto-oncogenes have been highly conserved through evolution because some of them are readily detected even in yeast by their homology to mammalian oncogenes, suggesting an important role for these genes in cell survival.

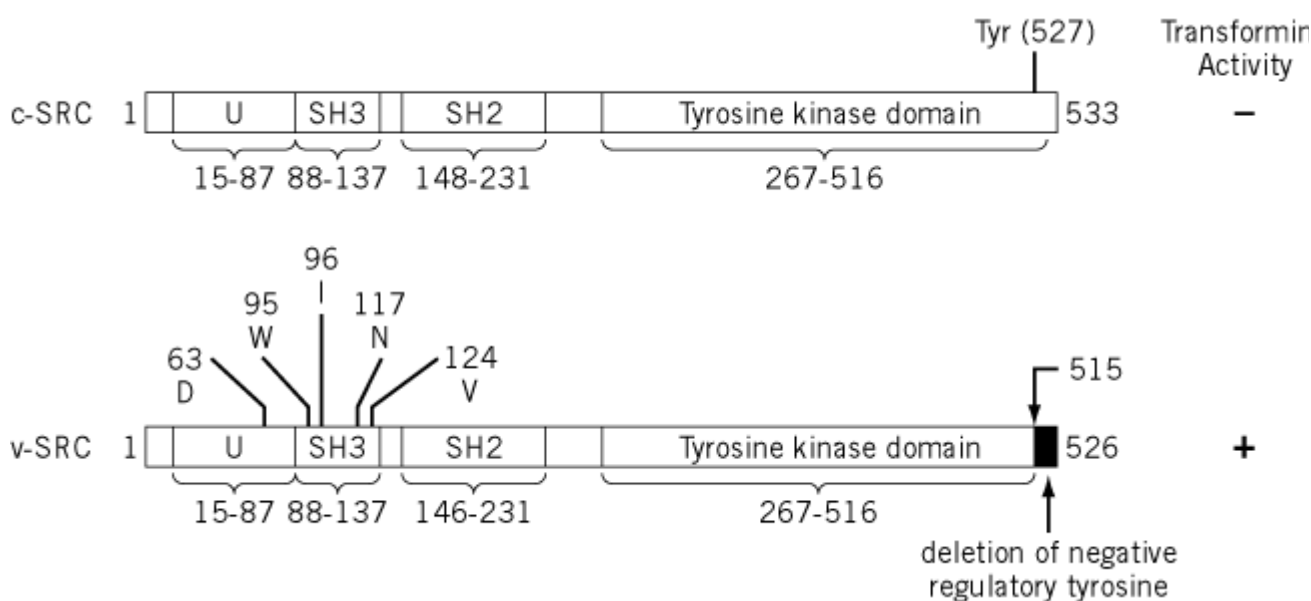
The observation that retroviral oncogenes are derived from normal cell DNA but readily transform cells in culture posed an important question how the same gene, expressed in the context of normal cell growth has no deleterious effect but readily induces transformation of cells when expressed as



an integral component of a retrovirus. A detailed structural analysis of v- and c-oncogenes revealed that retroviral oncogenes represent gain-of-function mutants of c-oncogenes. This is exemplified by the analysis of five retroviral oncogenes, *v-src*, *v-abl*, *v-ras*, *v-myb*, and *v-myc* and their normal cellular homologues, *c-src*, *c-abl*, *c-ras*, *c-myb* and *c-myc* (Fig. 2).

The *c-src* gene encodes a protein with four well-defined structural domains, termed the unique, SH3, SH2, and SH1 (or tyrosine kinase) domains (3). The kinase domain is located toward the C-terminus of the protein and contains a negative-regulatory tyrosine residue that is phosphorylated in a resting cell (3). This phosphorylated state adopts a local protein structure that renders the protein kinase inactive. Mitogenic stimulation of cells results in dephosphorylation of this tyrosine residue, which subsequently results in enhanced kinase activity of the protein. In the v-Src protein, this negative-regulatory region is deleted, resulting in constitutive activation of its kinase activity (3) (Fig. 3).

**Figure 3.** Activation of the *src* oncogene. The cellular *src* gene, designated as *c-src*, encodes a tyrosine kinase whose activity is regulated by tyrosine phosphorylation. In a normal resting cell, the Src protein is phosphorylated on a C-terminal, negative-regulatory tyrosine residue and has a low or undetectable level of enzymatic activity. Upon stimulation, the protein is dephosphorylated and active. As a consequence of transduction by RSV, the *v-src* oncogene contains a deletion at the 3' end of the gene and does not encode a negative-regulatory tyrosine residue. As a result, its kinase activity cannot be regulated. This constitutive activation of the enzyme leads to cellular transformation. *v-src* also contains a number of other mutations that might enhance its oncogenic activity, but none of the mutations alone activate *c-src*.



Abelson murine leukemia virus (A-MuLV) is a replication-defective retrovirus (4) that induces B-cell lymphomas *in vivo* and transforms both lymphoid and fibroblast cells *in vitro*. The proviral genome of A-MuLV encodes a single protein that is a fusion product of the virus-derived *gag* and cell-derived *abl* sequences (5). Like the *src* gene, the *abl* gene codes for a tyrosine kinase that also contains unique, SH3, SH2, and tyrosine kinase domains (6). But it also contains an additional C-terminal domain, whose function is not entirely clear, and the gene product of *c-abl* does not contain the negative-regulatory tyrosine residue at its C-terminus. The kinase activity of the c-Abl protein is negatively regulated by its SH3 domain, which is deleted from the v-*abl* gene product. Interestingly, the v-*abl* gene product also contains a point mutation in its C-terminal sequence, which enhances its tyrosine kinase and transforming activities (7). Thus, both the v-*src* and v-*abl* genes have alterations in their regulatory sequences that result in the constitutive activation of their tyrosine kinase activities, and this correlates with their transforming function.

The *ras* gene is an important gene that has been studied extensively because of its role in **growth-factor-associated signal-transduction** pathways (8). The viral *ras* genes code for 21-kDa proteins, designated p21, that belong to the G-protein superfamily. G-proteins bind to GTP with high affinity and hydrolyze it to GDP (see [GTP-Binding Proteins](#); [Gtpases](#); [Heterotrimeric G Proteins](#)). The GTP-bound Ras protein is biochemically active, whereas the GDP bound form is inactive. More importantly, the proportion of GTP-bound p21 Ras is tightly regulated in normal cells by feedback mechanisms, and it normally represents less than 5% of the total Ras protein. Comparison of v- and c-*ras* revealed that the viral oncogenes contain two point mutations, one in **codon 12** and a second in codon 59, both of which impair the intrinsic GTPase activity of the mutant proteins and render them resistant to negative regulation (8, 9). As a result, the v-Ras proteins remain in a constitutively activated (or GTP-bound) state, which contributes to their oncogenic activity.

The avian myeloblastosis virus, isolated in 1941 from a chicken tumor, is an acute transforming virus that causes myeloblastic leukemia in chickens and transforms myelomonocytic cells *in vitro* (10). Like most acute transforming viruses, AMV is replication-defective, having arisen by recombination between a nondefective leukemia virus and chicken cellular sequences. The normal cellular counterpart of this oncogene, c-*myb*, codes for a nuclear [transcription factor](#) that is essential for the proliferation of lymphoid, myeloid, and erythroblastoid cells. Sequence analysis of the c-*myb* gene revealed that the encoded protein contains an amino-terminal **DNA-binding** domain, a central transactivation domain, and a C-terminal negative-regulatory domain (11, 12). Comparison of the v-*myb* and c-*myb* sequences showed that the viral oncogene arose as a result of deletions in the 5' and 3' portions of the coding sequences, which results in deleting of a portion of the DNA-binding domain and the entire negative-regulatory domain. Although the deletion in the DNA-binding domain does not affect the ability of the viral protein to bind to DNA, the deletion of the C-terminal negative-regulatory domain results in enhanced transcriptional transactivation by the v-Myb protein, which in turn enhances the proliferative activity of the virally infected myeloid cells (13).

The [myc oncogene](#) was originally identified as the transforming gene of an acute transforming virus, called MC29, which was isolated from a chicken that had spontaneous myelocytomatosis (14). *myc*-related sequences were subsequently identified in three other independently derived chicken retroviral isolates. Comparison of the MC29-derived v-*myc* and chicken c-*myc* sequences showed that the viral oncogene contains the entire coding sequence of the chicken c-*myc* gene fused to the viral gag sequences, thus producing a gag-*myc* fusion protein (15, 16). Although fusion of the *myc* sequences to the gag gene results in elevated levels of Myc protein expression, deletional analysis of the proviral genome suggests that elevated expression of *myc* sequences alone is adequate for the transforming activity of this oncogenic virus. Thus, unlike with other oncogenic viruses, structural alterations play a less important role in the oncogenic activity of this gene, and overexpression of the encoded protein product alone is adequate to induce transformation of appropriate target cells.

Studies with oncogenic viruses, like those illustrated previously, provided important clues regarding the relative contributions of different mechanisms to the activation of proto-oncogenes. These comparisons between the viral and cellular oncogenes revealed a high frequency of deletions and mutations in v-oncogenes. In a vast majority of cases, these changes result in a gain of function, which contributes to the oncogenic potential of the viral oncogenes. In addition, these studies also showed that the viral LTRs play a critical role in viral oncogenesis because they direct high-level expression of v-onc gene products.

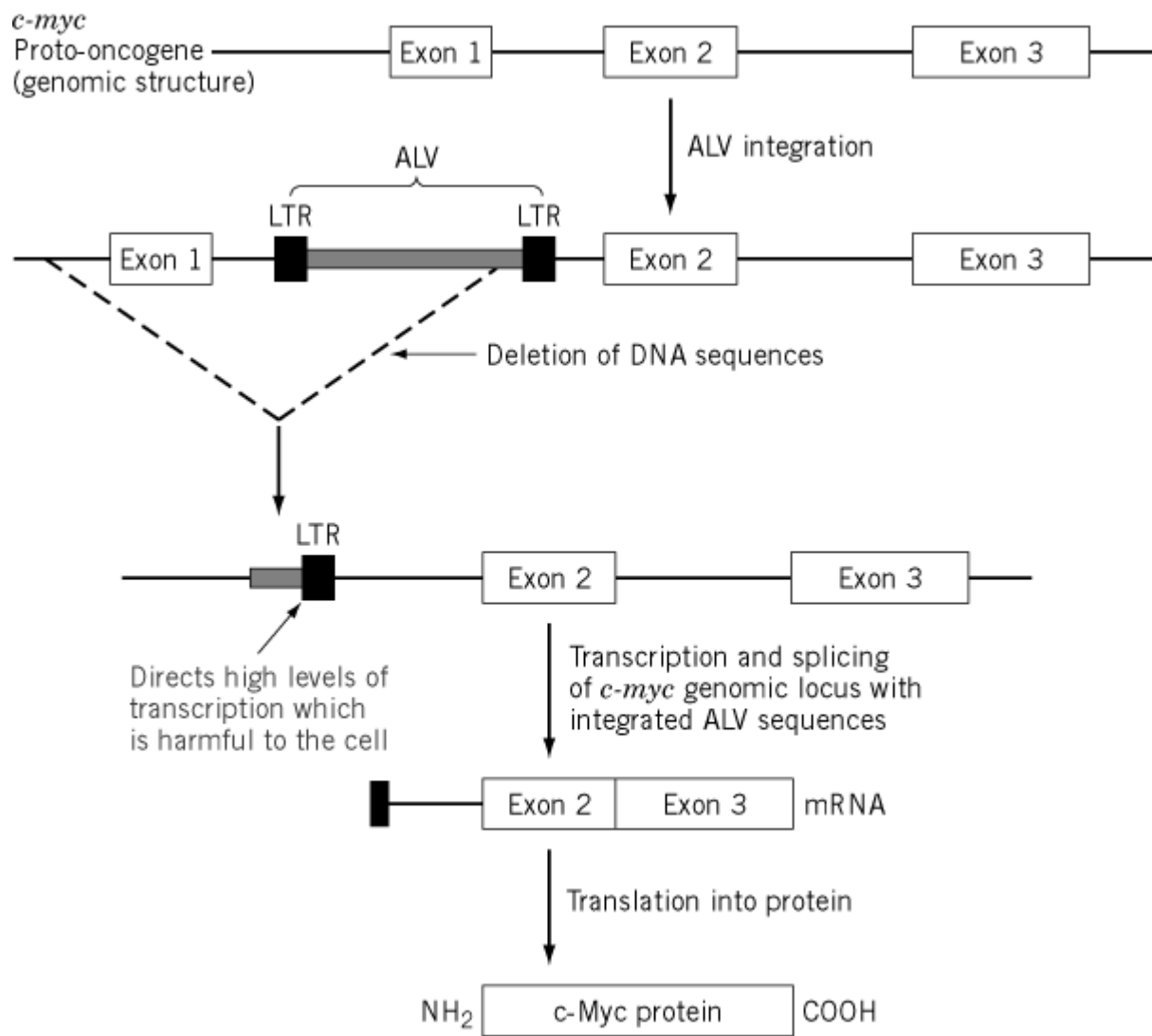
### 3. Oncogene Activation by Leukemia Viruses

The study of acute transforming viruses has provided us with a wealth of knowledge about the mechanisms of oncogene activation, but the mechanisms for the transformation by leukemia viruses, like ALV or Mo-MuLV, were not revealed by these studies. As noted in Table 1, leukemia viruses lack a cell-derived oncogene but nevertheless induce leukemias, albeit over prolonged periods of time. The long latency period for the development of leukemias suggested that genetic events in

addition to viral infection are required to induce transformation by this group of viruses. The first clues to the mechanisms for leukemia virus-induced transformation came from studies with ALV-induced bursal lymphomas in chickens (17) and Mo-MuLV-induced leukemias in mice (18).

Analysis of the chicken lymphomas induced by ALV demonstrated that these tumors contain integrated copies of the proviral genome and, most interestingly, these integration sites are often located within the *c-myc* proto-oncogene, the cellular homologue of the *v-myc* oncogene (17). Detailed analysis of the *c-myc* locus from these tumors indicated that the *c-myc* gene in normal cells contains three exons (see [Introns, Exons](#)), of which exon 1 contains the 5' noncoding sequences, whereas exons 2 and 3 contain the entire coding region of the gene (Fig. 4). In ALV-induced bursal lymphomas, this locus is disrupted by insertion of the proviral genome within intron 1 of the *c-myc* gene (Fig. 4). Several of the tumors analyzed showed that, following the integration of the provirus, the viral genome itself often undergoes rearrangements that result in the deleting most of the structural genes and one of the LTR sequences. This results in positioning a viral LTR immediately upstream of exons 2 and 3 of the *c-myc* gene, placing *c-myc* under the transcriptional control of the viral LTR. Because the viral LTR is designed to express high levels of RNA, the lymphoid cells that contain the viral LTR integration express abnormally high levels of *myc* RNA and protein, which, in combination with other genetic alterations, leads to formation of the tumor.

**Figure 4.** Activation of the *c-myc* oncogene as a consequence of ALV integration. In a normal cell, the three exons of the *c-myc* gene (designated 1, 2, and 3) are normally spliced into an mRNA that can be translated into the *c-myc* protein (only exons 2 and 3 actually code for the protein). When ALV integrates into the *c-myc* locus (usually between exons 1 and 2), the presence of the LTR sequences in the transcript results in increased expression of the *c-myc* protein, which is harmful to the cell. The integration of ALV is often followed by deletion of portions of viral sequences, leaving one of the LTRs intact.



Mo-MuLV induces thymic lymphomas, when injected into newborn mice (18), and myeloid leukemias when injected into pristane-primed BALB/C mice (19). Analysis of the thymic lymphomas induced by this virus revealed that the integration of the provirus often occurs into two common loci. As with ALV-induced lymphomas, many of the Mo-MuLV-induced lymphomas exhibit integration of the proviral genome into the *c-myc* locus and deregulated expression of this gene (18). Several of the tumors did not have the provirus integrated into the *c-myc* locus, but into a new locus, which was named *pim-1* (proviral integration MuLV) (18). *pim-1* codes for a protein kinase that acts as a transforming gene when overexpressed in lymphoid cells. Interestingly, a number of tumors had the provirus integrated in both the *c-myc* and *pim-1* loci and therefore have elevated transcriptional rates of both of these genes.

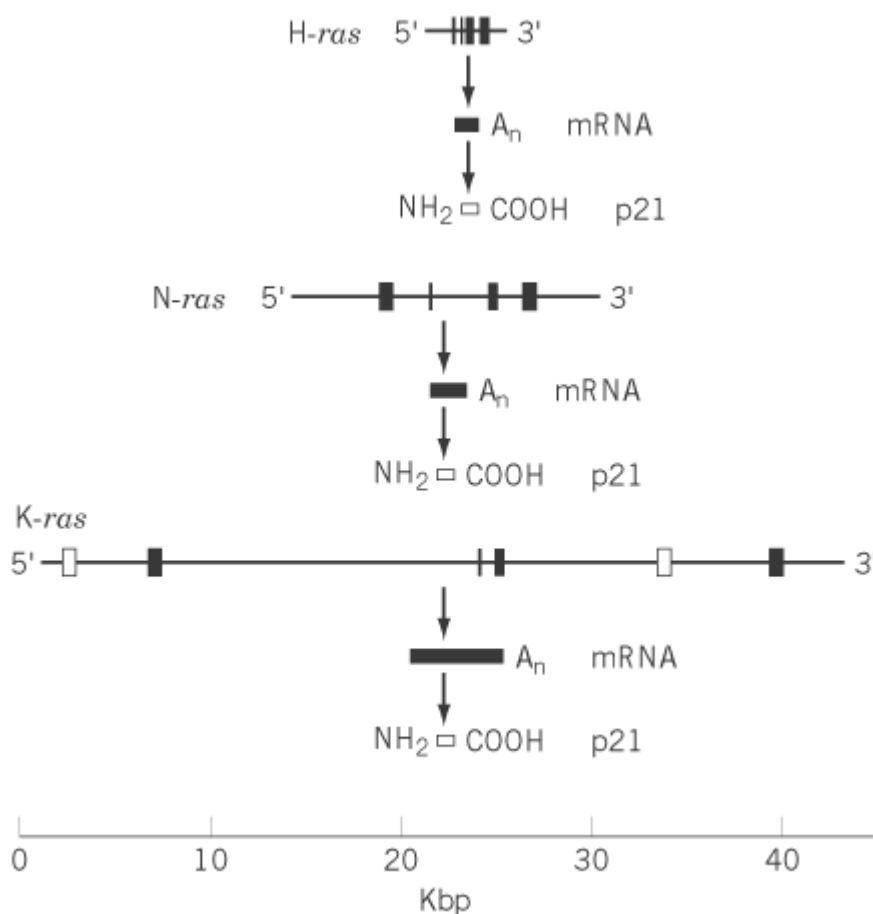
Analysis of myeloid tumors induced by Mo-MuLV revealed that the viral integrations in these tumors often occur in the *c-myb* locus, the oncogene originally identified in the avian myeloblastosis virus (AMV) (19). The mouse *c-myb* gene contains 15 exons, and viral integrations very often occur in either intron 3 or intron 9, resulting in either amino-terminal truncations or a C-terminal truncation of the protein after RNA splicing. This results in the synthesis of aberrant *myb* transcripts that lack either the 5' or 3' end. Interestingly, the aberrant messenger RNAs produced by these tumors exhibit deletions similar to those observed in the *v-myb* oncogene (20). Thus, activation of the *myb* gene is invariably accompanied by deletions in the coding regions of the gene. This finding lends support to the hypothesis that these structural alterations result in a gain of function for this oncoprotein, which in turn results in its oncogenicity.

### 3.1. Isolation of Dominant Transforming Genes from Human Tumor DNAs

Although the study of retrovirus-associated oncogenes led to the identification of a number of oncogenes, the link between these genes and human cancer was made by identifying dominant transforming genes from human tumor DNAs. The development of the NIH/3T3 cell transformation assay system, which utilizes the NIH/3T3 mouse fibroblast cell line which efficiently takes up and integrates exogenous DNA, has contributed immensely to identifying the transforming genes in human tumors. When the DNA from different human tumors is transfected and assayed for the induction of foci, about 20% of human tumors contain sequences that could transform NIH/3T3 cells (21). Molecular cloning and characterization of these sequences resulted in identifying several new oncogenes. One of the first human oncogenes identified from the T24 human bladder carcinoma cell line is related to the viral *ras* gene, previously found in the Harvey sarcoma virus that had been isolated from a rat tumor (22, 23). Detailed characterization of the genetic lesion associated with the transforming gene of the T24 bladder carcinoma cell line showed that a point mutation in codon 12 (similar to that seen in the *v-ras* oncogene), which leads to a single amino acid change in the coding sequence of this gene, results in its oncogenic activation (24, 25).

The human genome contains three closely related *ras* genes, all of which code for 21-kDa proteins and exhibit a high degree of homology (Fig. 5) (8). These three genes are called H-*ras*, K-*ras*, and N-*ras*. Of these, the H-*ras* gene is most similar to the *v-ras* genes derived from the Harvey sarcoma virus and BALB-MSV, whereas the K-*ras* gene is most similar to the oncogene encoded by the Kirsten sarcoma virus. The N-*ras* gene was originally isolated from a neuroblastoma after transfecting of the tumor DNA into NIH/3T3 cells and analyzing of the sequences responsible for inducing transformed foci. Detailed analysis of the tumor-derived oncogenic *ras* genes showed that almost all of these genes contain mutations in either the 12th or the 61st codon that result in constitutively activating of the protein products.

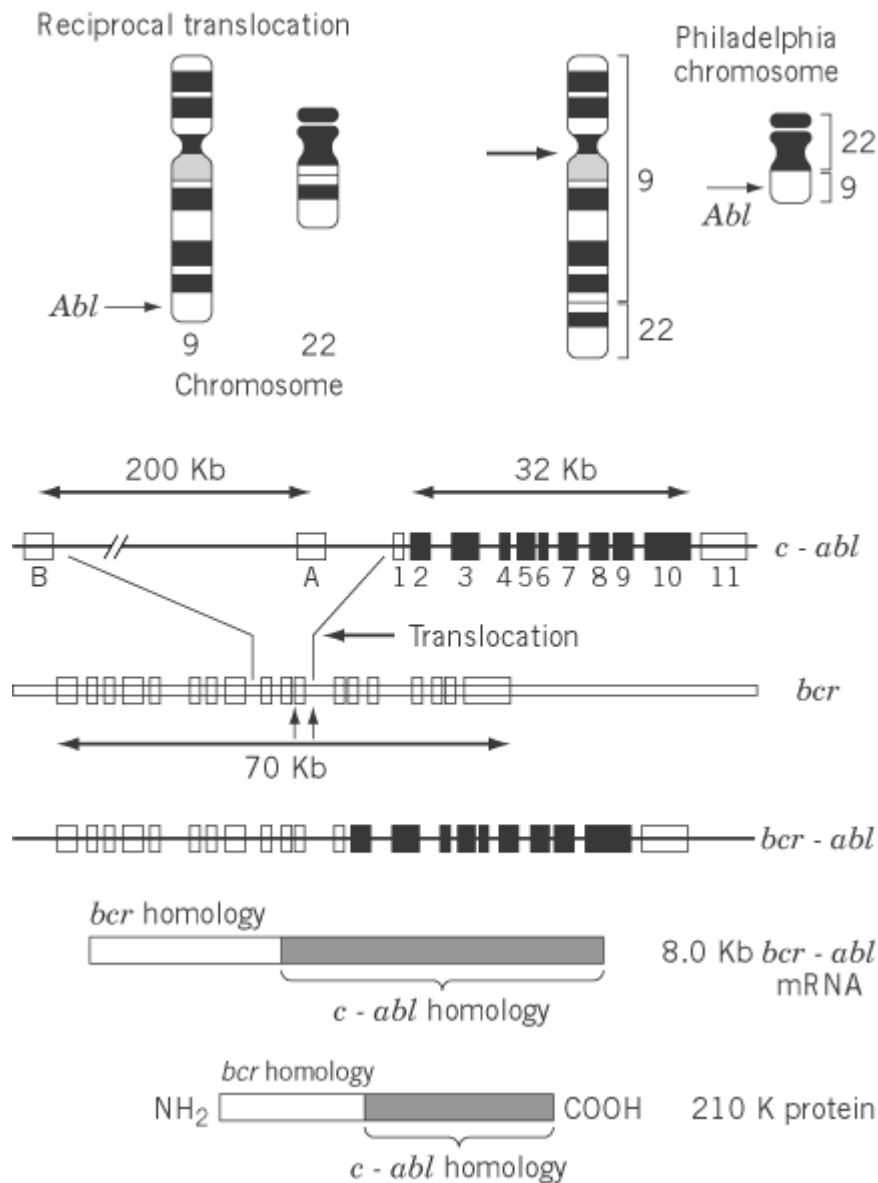
**Figure 5.** Comparison of the genomic structures of the *ras* proto-oncogenes. All of these genes encode four exons. K-*ras* encodes a fifth exon which is alternatively spliced during transcription. Although all of the encoded proteins are identical in size, the length of their transcripts and genetic loci differ considerably.



#### 4. Oncogenes and Chromosomal Abnormalities

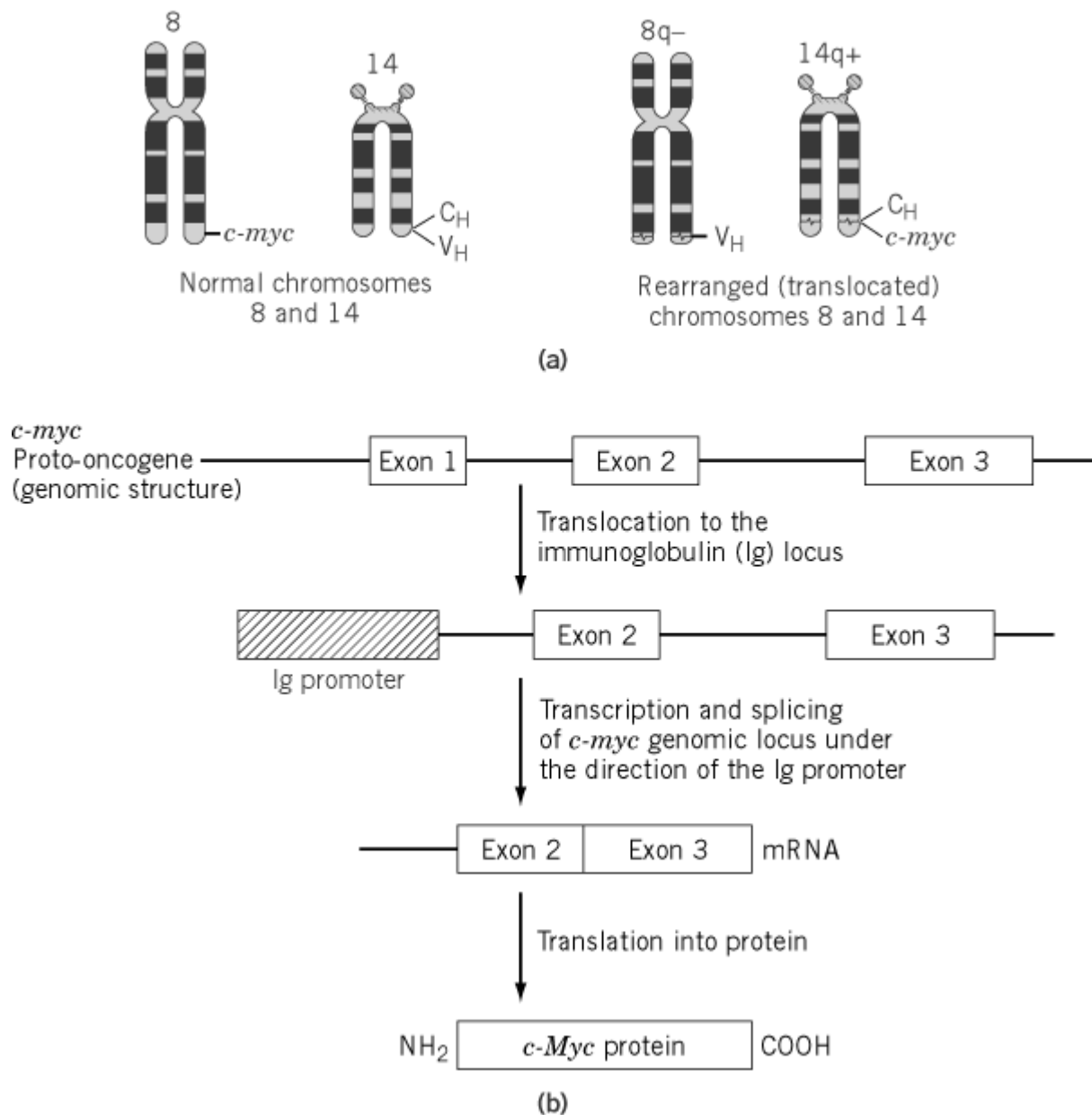
For many years, it has been known that tumor cells contain **chromosomal** abnormalities that are nonrandom. Although the hypothesis that such abnormalities might play a role in cancer was put forward in 1890 by Von Hansenmann (26), the first experimental evidence to support it was provided in 1960 by Nowell and Hungerford, who discovered the Philadelphia chromosome, a chromosomal abnormality consistently associated with human chronic myelogenous leukemia (CML) (27). This disease provides a paradigm of the genetic lesions associated with cancer because more than 90% of the cases involve a **translocation** in which a small portion of the long arm of chromosome 9 is transferred to the long arm of chromosome 22 (28) (Fig. 6). Even though later studies made it clear that a vast majority of human cancers contain one or more such chromosomal abnormalities, the actual molecular mechanisms associated with these abnormalities and their relationship to human cancer remained largely enigmatic until recently. During the past decade, a number of such chromosomal abnormalities have been studied in detail at the molecular level. All of the evidence accumulated thus far suggests that these chromosomal break points occur within or in the vicinity of an oncogene, which results in activating this gene. Thus, gene mapping studies have established that the *c-abl* oncogene is located on chromosome 9q34, the location where the break point occurs in the Philadelphia chromosome. During formation of the Philadelphia chromosome, a portion of the *c-abl* gene is translocated to chromosome 22 and is fused to a portion of the gene called *bcr*, which itself is disrupted during the translocation process (Fig. 6). This process generates a new gene, termed *bcr-abl*, that has enhanced oncogenic activity and whose expression leads the development of leukemia (29, 30). It is interesting to note that the BCR-ABL gene is structurally very similar to the GAG-ABL gene encoded by the Abelson murine leukemia virus. Both exhibit high levels of tyrosine kinase activity, which is essential for their transforming activity.

**Figure 6.** Generation of the *bcr-abl* oncoprotein by chromosomal translocation. The *abl* gene in a normal cell, is located on chromosome 9 and encodes a tyrosine kinase. During malignant transformation of myeloid cells, a portion of chromosome 9 that contains the *abl* locus translocates to chromosome 22 at the breakpoint cluster region (*bcr*) locus and generates the chimeric *bcr-abl* oncoprotein. Because the translocation results in deleting the sequences that negatively regulate *abl* tyrosine kinase activity, the fusion protein has constitutive and increased levels of enzymatic activity.



A second example of chromosomal abnormality is seen in Burkitt's lymphoma (a malignancy of the B-cell lineage), where a chromosomal exchange between chromosome 8 and 14 occurs (31, 32). This translocation results in juxtaposing the *c-myc* oncogene with the regulatory sequences of the immunoglobulin (Ig) genes (Fig. 7). Immunoglobulin loci code for **antibodies** produced solely by B lymphocytes. When these genes are expressed, the Ig regulatory sequences direct high levels of mRNA transcription of the Ig genes located next to them. When the *c-myc* oncogene is placed under the transcriptional control of the Ig regulatory elements, however, the B cell abnormally produces extremely high levels of *c-myc* mRNA and protein (instead of antibody), which induces a malignant phenotype.

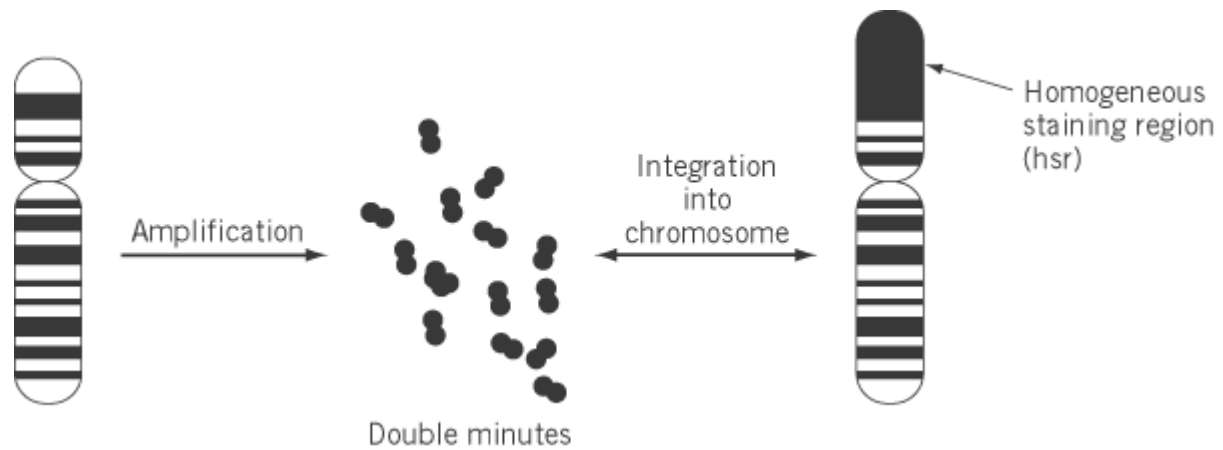
**Figure 7.** Activation of the *c-myc* oncogene as a consequence of translocation in Burkitt's lymphoma. (a) During translocation, the *c-myc* gene (located on chromosome 8) translocates to chromosome 14 at a site next to the immunoglobulin heavy chain (C<sub>H</sub>) locus. (b) This translocation results in juxtaposing the immunoglobulin regulatory sequences upstream of the *c-myc* gene. Because the immunoglobulin promoter positively regulates high levels of transcription in B lymphocytes, the *c-myc* gene (instead of the immunoglobulin gene) is constitutively transcribed and translated in the cell.



A third type of chromosomal abnormality is associated with the abnormal amplification of certain chromosomal regions that contain oncogenic sequences. The amplified DNA in these tumors can exist as a tandem copy of a set of genes on a given chromosome or as small chromosome-like structures, known as [double minute chromosomes](#) (Fig. 8). The double minute chromosomes present themselves as homogeneously staining regions (HSR) that are distributed over the entire field in a chromosome preparation. Thus, a variety of human neuroblastomas exhibit amplification of the *N-myc* gene (33), and several solid tumors are associated with the amplification of *erb-B-2* gene (34). Interestingly, the prognosis for these tumors correlates with their level of oncogenic amplification. Poor prognosis is associated with greater amplification of a given oncogenic sequence.

**Figure 8.** Formation of homogeneous staining regions as a result of DNA amplification. Once a locus is amplified by several rounds of DNA amplification, tandem repeats can be excised to form double minute fragments. These fragments can integrate into the locus of another chromosome, forming a homogeneous staining region (HSR). Because this process is reversible, double minute fragments and HSRs are amplified regions of DNA that are interchangeable.





## 5. Chemical Carcinogens often target Oncogenes

Chemical **carcinogens** have a very broad range of structures that have little or no obvious structural or biochemical similarity to each other. This diversity of their structures clouded earlier studies of possible mechanisms for their mechanism of action. A more puzzling fact about the most commonly occurring carcinogens was their relative insolubility in water. Rapid progress in this field was made when it was discovered that most chemical carcinogens do not act directly but undergo modifications inside the cell, following their uptake, that generate a new group of chemicals which act as carcinogens. The enzymes that convert procarcinogens to carcinogens are called [cytochrome P450s](#), and their major function in the cell is detoxification of foreign chemicals (35). Thus, there are two types of chemical carcinogens, one called direct-acting carcinogens that, as the name implies, act directly without any metabolic activation and the other, called indirect-acting carcinogens, which require metabolic activation for them to act as carcinogens. The terms procarcinogen and ultimate carcinogen describe the pre- and postprocessing states of indirect-acting carcinogens. Once inside the cell, both types of carcinogens interact with a wide variety of cellular components, including DNA, RNA, and proteins. It is now believed that the most important of these interactions is with DNA and results in introducing mutations. In fact, the finding that most carcinogens are mutagens allowed Bruce Ames to develop a number of assays using bacterial systems, the [Ames test](#), that permit rapid analysis of various environmental agents for their possible oncogenic activity (36).

Although it was becoming clear that chemical carcinogens may exert their influence by acting as mutagens, the identity of their targets remained unknown for quite some time. The first clues about their targets came from DNA [transfection](#) experiments using NIH/3T3 cells, where it was shown that DNA from chemically transformed cell lines contain activated *ras* genes. These studies suggested that proto-oncogenes might be the crucial targets of the chemical carcinogens (37). This concept received much needed experimental support with the demonstration that a majority of rat mammary tumors induced with a single dose of *N*-nitroso-*N*-methylurea during sexual development have point mutation in the *ras* gene that converts this proto-oncogene into a potent transforming gene (38). Following this discovery, it was shown that a number of chemical carcinogens act on proto-oncogenes and convert them into dominant transforming genes. Although the limitations of the NIH/3T3 assay system allowed detecting *ras* genes in a large number of chemically induced tumors, oncogenes other than members of the *ras* gene family have been identified in animal model systems. Thus, inducing neuroblastomas and gliomas by ethyl nitrourea or *N*-nitroso-*N*-methylurea activates the *neu* oncogene reproducibly (37). Therefore, it is increasingly evident that chemical carcinogens target cellular proto-oncogenes and bring about transformation through their action on these genes.

In summary, the evidence outlined shows that the mammalian organism contains a number of proto-oncogenes that have been highly conserved through [evolution](#). These genes can be oncogenically activated through transduction by retroviruses, mutation by chemical carcinogens, or as a result of

chromosomal abnormalities that occur during the genesis of cancer. To develop rational approaches for treating cancer, it is essential that we understand the structure and biochemical function of the proteins encoded by proto-oncogenes and the nature of biochemical changes that accompany the structural changes in naturally occurring tumors. For this reason, during the past five years, a number of investigators turned their attention to the task of delineating the function of the proteins encoded by proto-oncogenes in normal cell growth and differentiation.

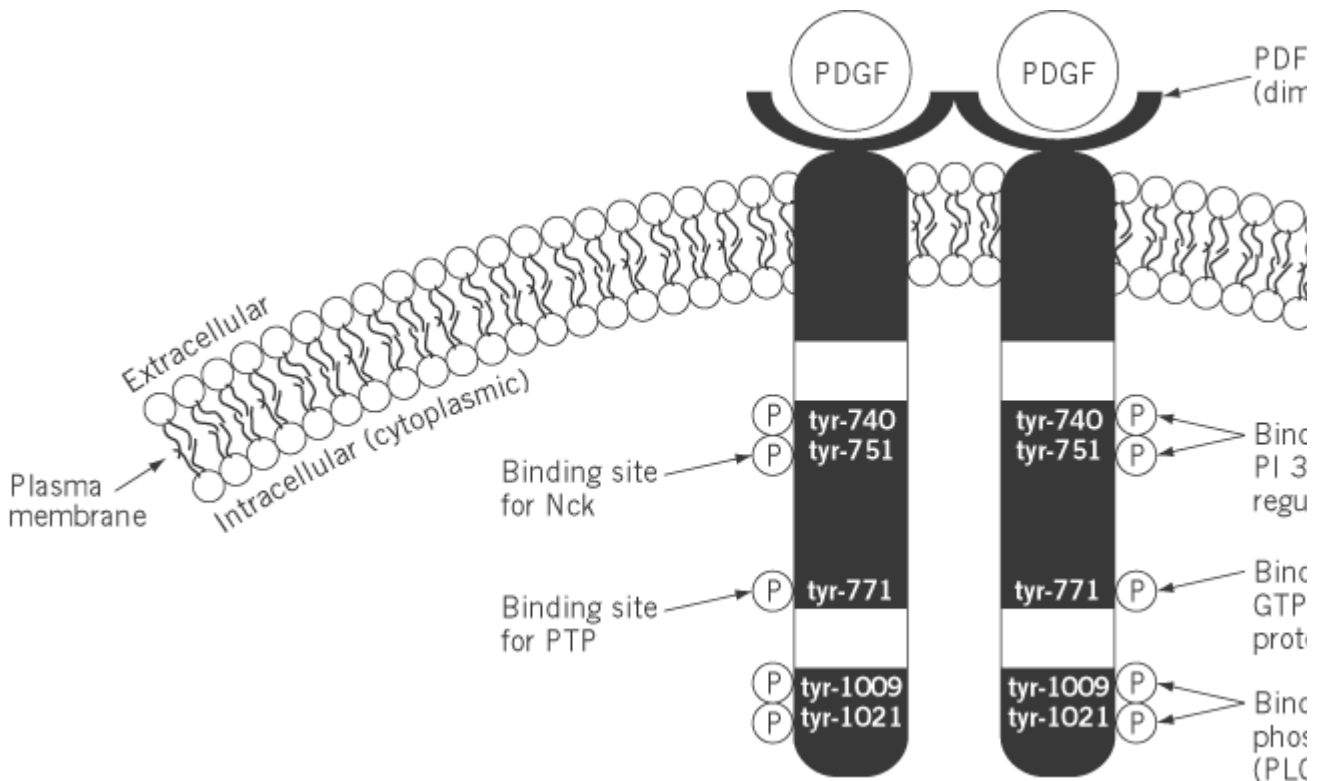
## 6. Function of Oncoproteins in the Normal Cell

The first clues about the normal function of oncogenes came from the observation that the *v-sis* oncogene product is highly related to the b-chain of the receptor for [platelet derived growth factor](#) (PDGF) (39). A second discovery that linked oncogenes with growth-factor signaling was the finding that the Erb-B2 oncogene, encoded by the avian erythroblastosis virus (Table 2), is an activated form of the [epidermal growth factor](#) (EGF) receptor (39). These two observations suggested for the first time that oncogenes might code for growth factors, their receptors, or signal-transducing proteins that relay the growth signal to the nucleus. Extensive biochemical analysis of these oncogenes has established that proto-oncogenes and their homologues are part of a stringently regulated cellular signaling network that controls cell proliferation and differentiation. In normal cells, the primary event leading to a cellular mitotic response is binding and activation of cellular receptors by growth factors (40, 41). Structural analysis of these receptors revealed that they contain a tripartite structure consisting of an external ligand-binding domain, a transmembrane domain, and an intracellular tyrosine kinase domain. Upon binding to growth factors, these receptors dimerize, which activates their intrinsic tyrosine kinase activity and leads to phosphorylation of multiple tyrosine residues in the intracellular portion of the molecule. These residues act as docking sites for a number of signal-transducing molecules that relay the signal to other proteins. An important feature of many signal-transducing molecules is the polypeptide sequences, called SH2 and SH3 domains, which were originally discovered in the *v-Src* protein (SH = Src-homology). It was subsequently discovered that peptides containing an SH2 domain recognize peptide segments containing phosphotyrosine residues in a sequence-specific manner (42). Similarly, peptides that contain an SH3 domain recognize sequence motifs that are **proline**-rich, usually containing the sequence motif Pro-X-X-Pro, where X is any amino acid residue (42). The molecules that contain these SH2 domains either have enzymatic activity or do not. The molecules that are enzymatically active include **phospholipase C $\alpha$**  (PLC), phosphatidylinositol 3-kinase (PI3K), and the GTPase-activating protein for p21 Ras, p120 GAP. The proteins that lack enzymatic activity include proteins, such as SOS and Grb-2, both of which simply function as adaptor molecules that tether proteins in a larger, more elaborate signaling complex. Collectively, these signal-transducing molecules activate multiple pathways that bring about the pleiotropic responses to growth factors. For example, PLC $\alpha$  hydrolyzes phosphoinositols to generate **diacylglycerol** and **inositol-3-phosphate**. PI3K phosphorylates phosphoinositides and generates putative [second messengers](#) that mediate rearrangements of the [cytoskeleton](#) and cellular trafficking (43). Diacylglycerol activates protein kinase C, a **serine/threonine kinase** that activates transcription factors which constitute the components of a complex called AP1.

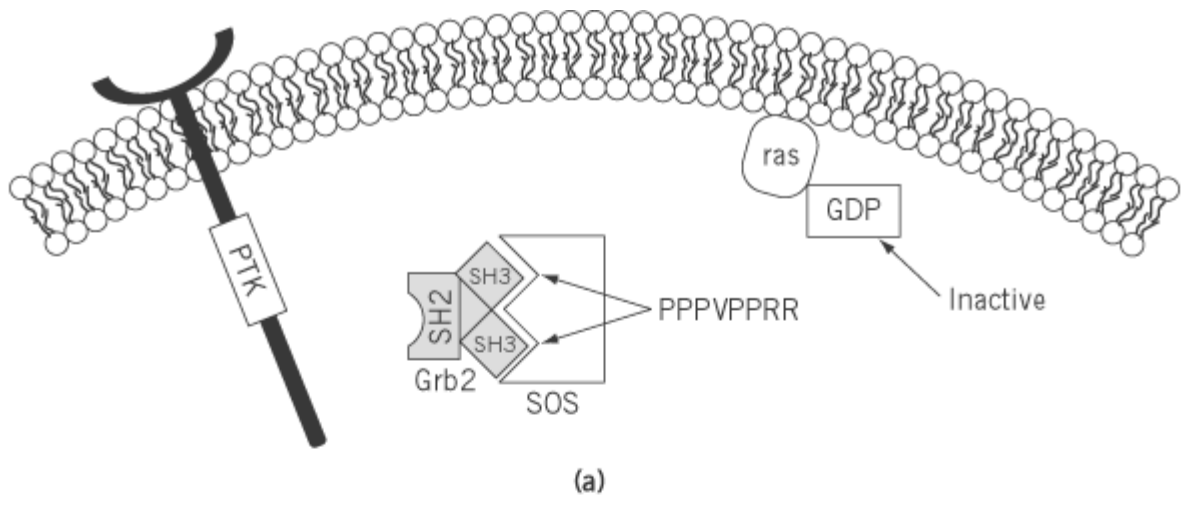
In contrast, proteins that act as adaptor molecules, such as Grb-2, simply contain SH3 and SH2 domains and bind to the receptor via their SH2 domain (Fig. 9) (39-42). Then, the SH3 domain of the protein binds to a second adaptor protein, called Sos, which contains a polyproline motif that recognizes the SH3 domain of Grb-2 and catalyzes exchange of GDP bound to Ras for GTP. Activated Ras plays a critical role in targeting Raf, another oncoprotein, to the plasma membrane. Ras contains a Cys-Ala-Ala-X box that is recognized by an enzyme called farnesyl transferase, which adds a **farnesyl** moiety to the [cysteine](#) residue, allowing Ras to bind to the plasma membrane. An adjacent polybasic domain of six [lysine](#) residues is a targeting signal for Raf. Once the Raf serine/threonine kinase has been activated by phosphorylation, it phosphorylates an enzyme called MEK. Phosphorylated MEK, in turn, initiates the activation/phosphorylation of mitogen-activated protein kinase (MAPK), which similarly activates Rsk (Fig. 10). Phosphorylation of transcription factors also results in their activation and participation in cellular proliferation pathways. The most

notable examples of the nuclear transcription factors activated by this pathway include Fos, Jun, Myc, and Myb, all of which were originally discovered as oncogenes of acute transforming viruses.

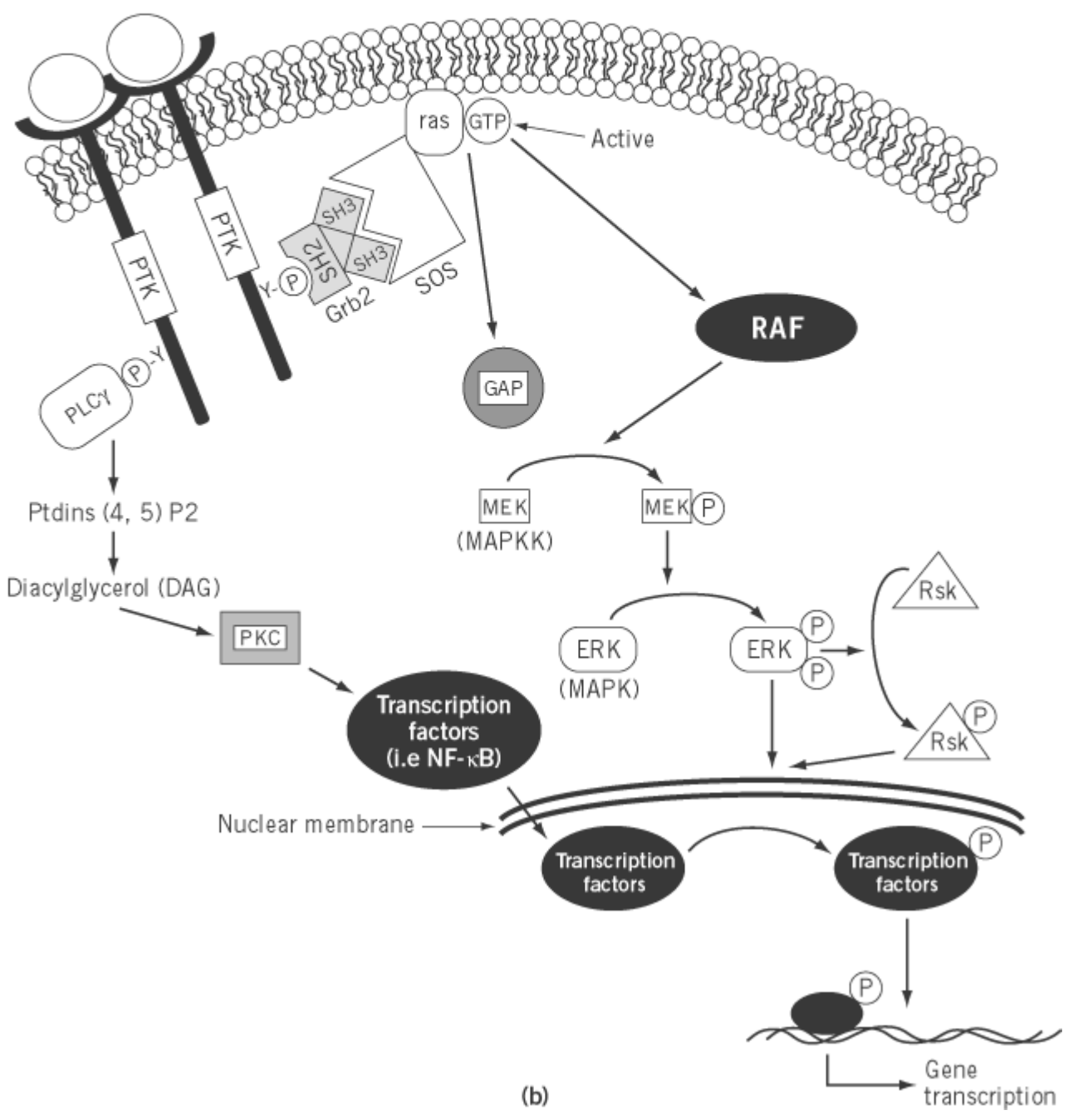
**Figure 9.** Interaction of SH2 domain signaling proteins with the PDGF receptor. The PDGF receptor contains several tyrosine residues to which SH2 domain-containing, signal-transducing proteins bind, but only after these residues have undergone autophosphorylation. The open boxes in the intracellular portion of the receptor represent a single tyrosine kinase domain.



**Figure 10.** Schematic representation of receptor tyrosine kinase (RTK) signaling. (a) In a resting cell, RTKs exist as nondimerized, nonphosphorylated transmembrane proteins. The Grb-2 and Sos adaptor proteins are bound through the two Grb2 SH3 domains and the proline-rich regions of Sos. In addition, ras remains in an inactive, GDP-bound state. (b) Once the RTK has bound to its cognate ligand, it dimerizes and undergoes autophosphorylation on multiple tyrosine residues. Now, the phosphorylated receptor can bind to the Grb2:Sos adaptor complex. Because the function of adaptor proteins is to link multiple signal-transducing proteins within a single pathway, the Grb:Sos complex triggers the formation of GTP-bound, activated ras. Subsequently, activated ras initiates the ordered, sequential phosphorylation of RAF, MEK, ERK, and Rsk which is essential for cell growth and differentiation. The phosphorylation of several classes of transcription factors enables them to bind to DNA and induce gene transcription.



(a)



(b)

The activation of multiple signaling pathways by the PDGF receptor illustrates the complexity of the signaling networks involved in regulating cell proliferation. The constitutive activation of these receptors by truncation of the amino-terminal regulatory domain, by point mutations, or by overexpression activates all of these pathways and leads to oncogenesis.

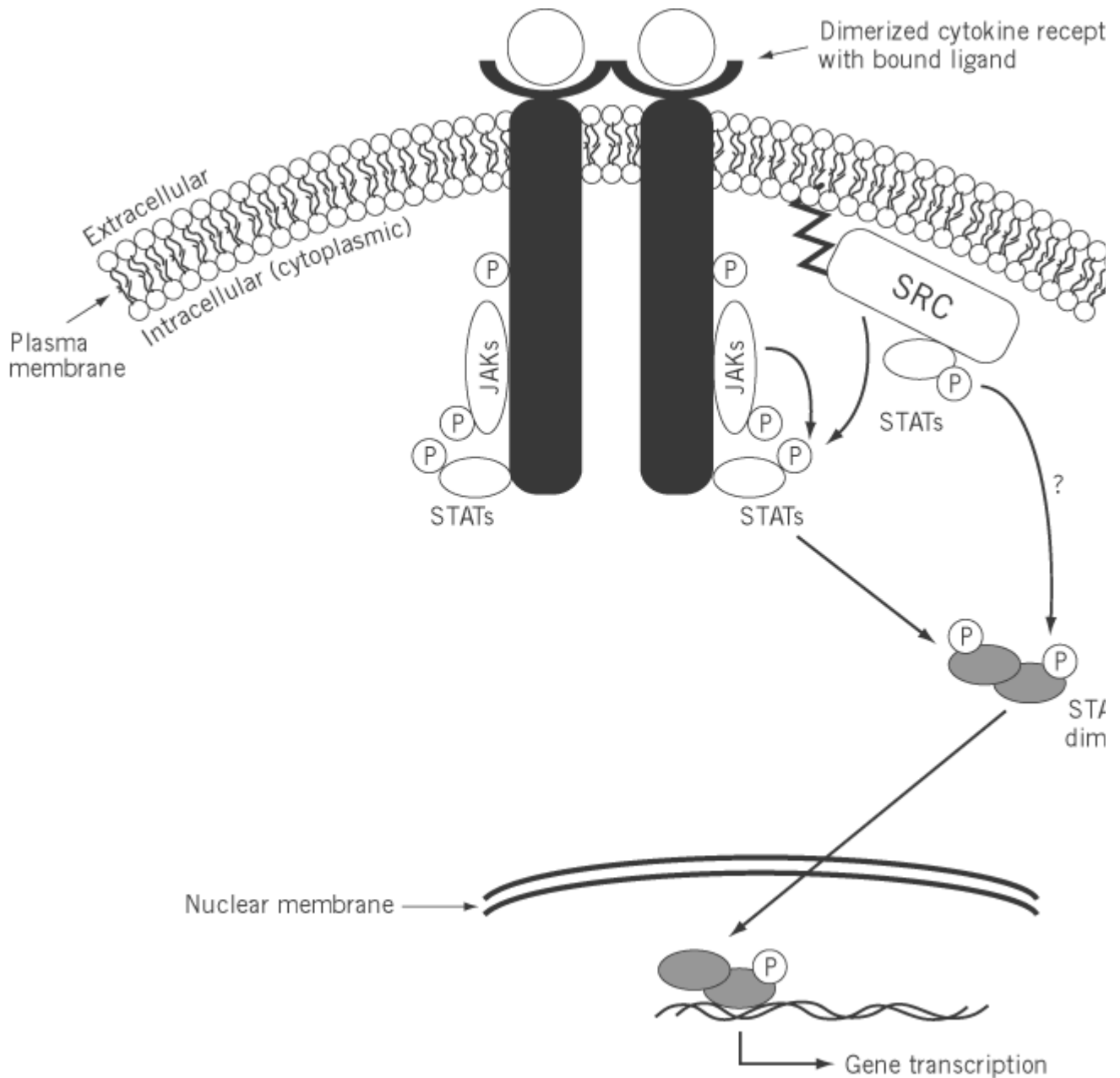
## 7. Subversion of Cytokine Signaling by Oncogenes

Cells of **hematopoietic** lineages transduce their growth signals predominantly via a group of receptors known as cytokine receptors. Molecular cloning of cytokine receptors and subsequent structure-function studies have revealed that, unlike growth-factor receptors, several of the cytokine receptors lack a cytoplasmic kinase domain. Nevertheless, interaction of a cytokine with its receptor rapidly induces tyrosine phosphorylation of the receptor and a variety of cellular proteins, suggesting that these receptors transmit their signals through cellular tyrosine kinases. Recently, new evidence has emerged to indicate that most cytokines transmit their signals via a new family of tyrosine kinases, termed JAKs, and via Src kinases ([44](#), [45](#)).

The JAK family of protein tyrosine kinases differs markedly from other classes of by the presence of an additional kinase domain. Currently, this family consists of four members, JAK1, JAK2, JAK3, and TYK2. Although JAK kinases do not contain the SH2 and SH3 domains that characterize Src family tyrosine kinases, they do possess highly conserved domains upstream of the kinase domains. These kinases, either alone or in conjunction with each other, are responsible for effects mediated by several cytokines and neurokinins, including [interleukins](#), [interferons](#), erythropoietin, prolactin, [growth hormone](#), oncostatin M, and ciliary neurotrophic factor.

Current models suggest that when cytokine receptors that contain single polypeptide chains (such as erythropoietin, growth hormone, prolactin, and GCSF) interact with cytokines, they undergo dimerization, which increases the affinity of the cytoplasmic domain of the receptor for JAK kinases and results in a ligand-dependent increase of a complex that contains the receptors and JAKs (Fig. [11](#)) ([44](#)). This association requires the membrane-proximal cytoplasmic domain of the receptor and results in activating the JAKs through an event associated with tyrosine phosphorylation. Then, the activated kinases phosphorylate the receptors and cellular substrates that regulate a wide variety of responses in cells.

**Figure 11.** Activation of STAT proteins in response to cytokine stimulation. Cytokine binding induces receptor dimerization and the activation and phosphorylation of two classes of tyrosine kinases, termed JAKs and Src kinases. Once the Src kinases and the JAKs have been activated, they phosphorylate a second class of proteins, termed STATs. Once phosphorylated, these proteins dimerize and translocate to the nucleus where they bind to DNA and stimulate transcription of genes whose promoter sequences contain STAT-responsive elements.



A surge of recent studies has demonstrated that cytokine signaling is associated with the phosphorylation of a group of transcription factors, Signal Transducers and Activators of Transcription now known as STATs, or (46). In a resting cell, these STATs remain in the cytoplasm. Stimulation of cells with cytokines results in phosphorylation of STATs mediated by Src kinases and JAKs. This, in turn, results in migration of STATs to the nucleus, where they participate in binding to DNA sequences containing the STAT [response elements](#). All of the STAT proteins are highly related and are characterized by a carboxy-terminal SH3 domain, followed by an SH2 domain. All of the proteins contain a conserved C-terminal tyrosine residue that is phosphorylated and is essential for their biochemical activity. Together, these studies suggest that, upon interaction with their receptors, cytokines activate JAKs and Src kinases, which in turn activate members of the STAT family of transcription factors.

Now, it is well established that several of the oncogenic tyrosine kinases that belong to the Src and Abl families have profound effects on the cytokine dependence of hematopoietic cell lines. Most notably, constitutive expression of the *v-src*, *v-abl*, and *bcr-abl* oncogenes in interleukin-dependent hematopoietic cell lines renders them cytokine-independent for growth. This alteration in growth-

factor dependence suggests that these oncogenes might interfere with signal-transduction pathways associated with cytokines. Recent studies have provided conclusive evidence that these oncogenes produce these transformations by constitutively activating STATs via tyrosine phosphorylation, often without activating JAK kinases (45, 47). For example, in BCR-ABL-positive CML cells, STAT 5 is phosphorylated constitutively, activating the proliferative pathways that are associated with this transcription factor in the absence of any cytokine-receptor interactions (48). These studies link chromosomal abnormalities that alter signal-transduction pathways to the onset of neoplasia.

## 8. Activation of Intracellular Phosphorylation Cascades Leads to Changes in Gene Expression

Activation of the intracellular phosphorylation cascades previously described invariably leads to two types of changes in gene expression. Early-response genes are induced with 15 minutes of growth-factor treatment, and their induction does not require **protein biosynthesis** because [cycloheximide](#), a potent inhibitor of protein synthesis, has no effect on their induction (49). By contrast, delayed-response genes are not induced until at least 1 hour after growth-factor treatment, and their induction requires protein biosynthesis. Delayed-response genes are induced by the products of early-response genes, several of which are transcription factors (50). Examples of early-response genes are *c-fos* and *c-myc*, and *c-myb* is an example of a late-response gene.

Both classes of genes are silent and not transcribed in cells in the  $G_0$  phase of the [cell cycle](#) (see later), but they are induced at high levels when growth factors are added to cells. Prolonged exposure to growth factors results in a gradual decline of gene expression, for some genes apparently to zero, and for certain others to a steady-state level above that observed during  $G_0$ .

## 9. Growth Factors and Cancer

The importance of precisely regulating receptor activation and signal transduction is further exemplified by the fact that a large variety of structural changes found in receptors lead to their constitutive activation and, consequently, to subversion of molecular control mechanisms that contributes to the onset of neoplasia. Very often, human cancers have been found when the genes for growth factor receptors are amplified. This provides a great proliferative advantage to these cells. Additional changes that result in growth advantage are point mutations in the *ras* and *raf* genes that activate this pathway constitutively and overexpression of early-response or late-response genes, such as *c-myc* or *c-myb*, which activate transition from the  $G_1$  to the S-phase of the cell cycle.

Identification of aberrations in growth-factor-induced, signal-transduction pathways in tumor cells is currently an active area of research. Detailed understanding of these mutants will provide important insights into fundamental mechanisms of normal cell growth and will also enhance our understanding of oncogenesis and open new avenues for diagnosis and therapy.

## 10. Oncogenes and the Cell Cycle

The [cell cycle](#) is a regulatory “clock” that controls the proliferative rate of a dividing cell. More importantly, many of the controls that govern the progression and length of the cell cycle and the phases that comprise it are controlled by the products of oncogenes and a second class of proteins, termed [tumor suppressor genes](#) (50). The eukaryotic cell cycle has remained highly conserved throughout evolution and therefore has been studied in organisms ranging from yeast to man. Virtually all studies conducted thus far reveal that the cell cycle is divided into four phases ( $G_1$ , S,  $G_2$ , and M; see Fig. 12) and is heavily regulated at two checkpoints: (1) a point in the  $G_1$  phase where the cell becomes committed to replicate its genetic material and; and (2) a point at the  $G_2$ /M border where the cell becomes committed to division.

Unless cells have received a stimulus to proliferate or differentiate, most remain in a resting state,

called  $G_0$ . When the organism requires additional cells, however, as is seen frequently in the hematopoietic system, extracellular stimuli induce the cells to enter the  $G_1$  phase of the cell cycle and become committed to cell division. At a late point in the  $G_1$  phase of the cell cycle a potentially dividing cell reaches the “restriction point,” a time at which the cell must determine whether the conditions are suitable for continued proliferation. Provided that conditions are conducive to proliferation, the cell proceeds past this checkpoint. An absolute prerequisite for cell growth is the duplication of its genetic material, which occurs during the S-phase. Once the DNA is replicated, the cell “ascertains” whether this process has been correctly executed at the second checkpoint during  $G_2$ . Provided that it has, the cell divides during mitosis, the M-phase. Frequently, malignant cells effectively override one or both of the intrinsic checkpoints that are normally used by their normal counterparts. Although some of the oncogenes previously discussed, such as *ras*, force progression through  $G_1$ , other genes, called [tumor suppressor genes](#), function as “gatekeepers” of these restriction points. Loss of these genes, therefore, enables a malignant cell to ignore all of the safeguards aimed at preventing unwanted cell division.

## 11. Conclusion

As summarized previously, we have come a long way during the past two decades in understanding the molecular mechanisms of cell growth, differentiation, and [neoplastic transformation](#). It is becoming increasingly clear that aberrations in oncogenes, which act as positive regulators of growth, and aberrations in tumor suppressor genes, which act as negative regulators of growth, contribute to developing the neoplastic state. Although we have accumulated a great deal of information about the role of various oncoproteins in signal-transduction pathways, we have yet to define the precise pathways that lead to cell growth versus differentiation. It is becoming apparent that cell proliferation and differentiation are interrelated and that mechanisms that lead to a block in cell differentiation result in increased proliferation of cells and ultimately into neoplasia. It is reasonable to hope that a detailed understanding of the signal-transduction pathways for cell growth and differentiation will also provide us with clues that allow us to override the block in differentiation in tumor cells. The prospect of developing such therapeutic approaches for treating cancer is very promising in the immediate future. Identifying oncogenes and delineating their function clearly has had a major impact on the development of such approaches.

## Bibliography

1. P. Rous (1911) *J. Exp. Med.* **13**, 397–411.
2. D. Stehelin, H. E. Varmus, J. M. Bishop, and P. K. Vogt (1976) *Nature*, **260**, 170–173.
3. M. T. Brown and J. A. Cooper (1996) *Biochem. Biophys. Acta* **1287**, 121–141.
4. H. T. Abelson and L. S. Rabson (1970) *Cancer Res.* **30**, 2213–2222.
5. E. P. Reddy, M. J. Smith, and A. Srinivasan (1983) *Proc. Natl. Acad. Sci. USA* **80**, 3623–3627.
6. C. Oppi, S. K. Shore, and E. P. Reddy (1987) *Proc. Natl. Acad. Sci. USA* **84**, 8200–8204.
7. S. K. Shore, S. L. Bogart, and E. P. Reddy (1990) *Proc. Natl. Acad. Sci. USA* **87**, 6502–6506.
8. M. Barbacid (1987) *Annu. Rev. Biochem.* **56**, 779–827.
9. G. Bollag and F. McCormick (1991) *Annu. Rev. Cell Biol.* **7**, 601–632.
10. W. J. Hall, C. W. Bean, and M. Pollard (1941) *Am. J. Vet. Res.* **2**, 272–279.
11. M. A. Baluda and E. P. Reddy (1994) *Oncogene*, **9**, 2761–2774.
12. K. Weston and J. M. Bishop (1989) *Cell*, **58**, 85–93.
13. G. Patel, B. Kreider, G. Rovera, and E. P. Reddy (1993) *Mol. Cell. Biol.* **13**, 2269–2276.
14. I. X. Vanov, Z. Mladenov, S. Nedyalkov, T. G. Todorov, and M. Yakimov (1964) *Bull. Inst. Pathol. Comp. Anim.* **10**, 5–38.
15. E. P. Reddy, R. K. Reynolds, D. K. Watson, R. A. Schultz, J. Lautenberger, and T. S. Papas (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2500–2504.



16. D. K. Watson, E. P. Reddy, P. H. Duesberg, and T. S. Papas (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2164–2150.
17. W. S. Hayward, B. G. Neel, and S. M. Astrin (1981) *Nature* **290**, 475–480.
18. H. T. Cuyppers, G. Selton, W. Quint, M. Zijlstra, E. Robanus-Maandag, W. Boelens, P. Van Wezenbeek, C. Meleif, and A. Berns (1984) *Cell* **37**, 141–150.
19. G. L. C. Shen-Ong, M. Potter, J. F. Mushinski, and E. P. Reddy (1984) *Science* **226**, 1077–1080.
20. R. Tantravahi, H. Dudek, G. Patel, and E. P. Reddy (1996) *Oncogene* **13**, 1187–1196.
21. G. M. Cooper (1982) *Science* **218**, 801–806.
22. E. Santos, S. R. Tronick, S. A. Aaronson, S. Pulciani, and M. Barbacid (1982) *Nature* **298**, 343–347.
23. C. Der, T. G. Krontiris, and G. M. Cooper (1982) *Proc. Natl. Acad. Sci. USA* **79**, 3637–3640.
24. C. J. Tabin, S. M. Bradley, C. I. Bargmann, R. A. Weinberg, A. G. Papageorge, E. M. Scolnick, R. Dhar, D. R. Lowy, and E. H. Chang (1982) *Nature* **300**, 143–149.
25. E. P. Reddy, R. K. Reynolds, E. Santos, and M. Barbacid (1982) *Nature* **300**, 149–152.
26. D. Von Hansemann (1890) *Virchows Arch. Pathol. Anat. Physiol.* **119**, 299–326.
27. P. C. Nowell and D. A. Hungerford (1960) *Science* **132**, 1497–1499.
28. J. D. Rowley (1973) *Nature* **243**, 290–293.
29. J. Groffen, J. R. Stephenson, N. Heisterkamp, A. De Klein, C. B. Bartam, and G. Grosveld (1984) *Cell* **36**, 93–99.
30. E. Shtivelman, R. P. Lifshitz, R. P. Gale, and E. Canaani (1985) *Nature* **315**, 550–553.
31. R. Dalla-Favera, S. Martinitti, R. C. Gallo, J. Erikson, and C. Croce (1983) *Science* **219**, 963–967.
32. R. Taub, I. Kirsh, C. Morton, G. Lenoir, D. Swan, S. Tronick, S. Aaronson, and P. Leder (1982) *Proc. Natl. Acad. Sci. USA* **79**, 7837–7841.
33. K. Alitalo, M. Schwab, C. C. Lin, H. E. Varmus, and J. M. Bishop (1983) *Proc. Natl. Acad. Sci. USA* **80**, 1707–1711.
34. W. C. Dougall, X. Qian, N. C. Peterson, M. J. Miller, A. Samanta, and M. I. Green (1994) *Oncogene* **9**, 2109–2123.
35. A. H. Conney (1982) *Cancer Res.* **42**, 4875–4882.
36. B. N. Ames (1979) *Science* **204**, 587–593.
37. M. Barbacid (1986) *Trends Genet.* **2**, 143–149.
38. S. Sukumar, V. Notario, D. Martin-Zanca, and M. Barbacid (1983) *Nature* **306**, 658–661.
39. S. A. Aaronson (1991) *Science* **254**, 1146–1153.
40. S. E. Egan, B. W. Giddings, M. W. Brooks, L. Buday, A. M. Sizeland, and R. A. Weinberg (1993) *Nature* **363**, 45–51.
41. F. McCormick (1994) *Curr. Opin. Genet. Dev.* **4**, 71–76.
42. T. Pawson and J. Schleissinger (1993) *Curr. Biol.* **3**, 434–442.
43. M. J. Berridge, J. P. Heslop, R. F. Irvine, and K. D. Brown (1984) *Biochem. J.* **222**, 195–201.
44. J. N. Ihle (1995) *Adv. Immunol.* **60**, 1–35.
45. P. Chaturvedi, M. V. Reddy, and E. P. Reddy (1998) *Oncogene* **16**, 1749–1758.
46. J. E. Darnell Jr. (1997) *Science* **277**, 1630–1635.
47. P. Chaturvedi, S. Sharma, and E. P. Reddy (1997) *Mol. Cell. Biol.* **17**, 3295–3304.
48. N. Carlesso, D. A. Fink, and J. D. Griffin (1996) *J. Exp. Med.* **183**, 811–820.
49. H. R. Herschman (1991) *Annu. Rev. Biochem.* **60**, 1347–1349.
50. X. Graña and E. P. Reddy (1996) *Oncogene* **11**, 211–220.

### Suggestions for Further Reading

51. E. P. Reddy, A. M. Skalka, and T. Curran (1988) *The Oncogene Handbook*, Elsevier, New York.
52. G. M. Cooper (1995) *Oncogenes*, Jones and Bartlett, Boston.
53. B. Vogelstein and K. W. Kinzler (1995) *The Genetic Basis of Human Cancer*, McGraw-Hill, New York.

### Opal Suppressor

Opal suppressors are mutant tRNAs that translate the UGA (opal) termination codon as a sense codon. Opal mutations cause protein synthesis to terminate prematurely, resulting in inactive, truncated polypeptides. Opal suppressors allow for protein synthesis beyond the translational block resulting in active protein. Hence the term “suppressor;” these mutant tRNAs “suppress” the phenotypes of opal mutations. These suppressors have been used in studies of the translational apparatus and mechanisms. For complete discussions of these and other suppressors, see [Nonsense Suppression](#), [Suppressor tRNA](#), and [Genetic Suppression](#).

### Operons

The term *operon* is occasionally used to mean any segment of DNA that is **transcribed** into a single RNA molecule, but in the more usual sense intended here, it denotes a set of contiguous **genes** that, though separately **translated** to yield distinct [polypeptide chains](#), are transcribed together from a common transcription start point (**promoter**). This is a common mode of gene organization in **bacteria**, but it has not been described in eukaryotic organisms, except in [nematode](#) worms and **trypanosomes**, where it takes a rather different form.

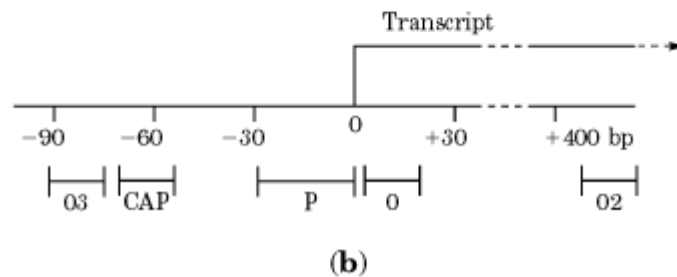
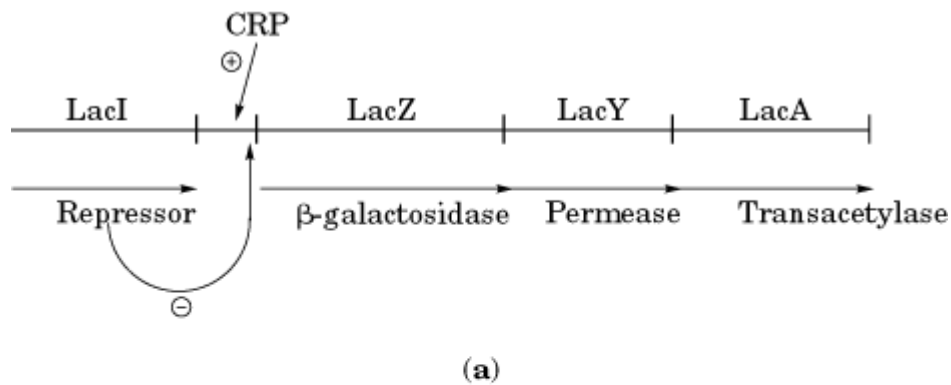
In bacteria, genes that contribute to a common function are frequently found united in an operon, an arrangement that makes it possible to regulate all their activities together by a single control of transcription. The operons of *Escherichia coli* and *Salmonella typhimurium*, by far the most thoroughly studied ([1](#)), may be considered as classified in three groups, involved, respectively, with utilizing occasionally available energy sources (eg, the sugars lactose, galactose, and arabinose; see [Lac Operon](#), [Gal Operon](#), [ara Operon](#)), with synthesis of [amino acids](#) (eg, tryptophan, histidine, phenylalanine—see leucine [TRP Operon](#)), and **protein biosynthesis** (operon organization of genes for **ribosomal** proteins). All operons are subject to some form of joint regulation of their constituent genetic activities to meet the current requirements of the cell. Thus, enzymes for utilizing sugar are produced only in the presence of the relevant sugars and the absence of glucose, the enzymes for synthesizing amino acids are not produced when the amino acids are already present in the growth medium, and ribosomal proteins are not allowed to accumulate in slowly growing cells. The joint regulation of the genes in an operon may operate through transcriptional or translational control or both.

## 1. The *lac* Operon - Transcriptional Regulation

The *lac* operon, the first to be described (2) and the one that provided the paradigm model for operon regulation, contains three genes. Two of them (*LacZ* and *LacY*) are essential for utilizing lactose as a carbon source, and code, respectively, for **beta-galactosidase**, which hydrolyses lactose to glucose plus galactose, and b-galactoside permease, a membrane protein needed to efficiently take up lactose into the cell. The third gene of the operon, *LacA*, encodes b-galactoside transacetylase, an enzyme that has defined chemistry but uncertain function and is lost by mutation without any obvious effect. The three genes are adjacent to each other in the order Z-Y-A, and are transcribed into a common [messenger RNA](#). *LacZ* is at the upstream end (“upstream” with reference to the direction of transcription and translation). All three functions of the operon are **induced** in a constant ratio (coordinately) by lactose or some other inducing b-galactoside and are repressed by glucose. The synthetic b-galactoside thiomethyl-b-D-galactopyranoside (TMG) is often used as an artificial inducer but is not itself usable as a carbon source.

The keys to understanding of induction of the *Lac* operon were mutations that resulted in the production of b-galactosidase and the other operon products constitutively, that is, without the need for induction. These mutations, which did not affect repressibility by glucose, fell into two groups. One group mapped within a closely linked gene *LacI*, outside and upstream of the operon. The other group, called operator-constitutive ( $o^c$ ), mapped within a very short segment called the **operator**, which mapped just upstream of the *LacZ* b-galactosidase-encoding sequence. A little further upstream another segment was identified through mutations that eliminate or reduce the level of expression of the operon under all conditions. This segment was called the **promoter**, the putative **RNA polymerase**-binding and transcription-initiation sequence (Fig. 1). Experiments with partial **diploids** that have duplicate *Lac* operon segments carried on [plasmids](#), showed that the operator-constitutive ( $o^c$ ) mutants are dominant over wild type, or, more precisely, [cis-dominant](#) because their mutant operators make b-galactosidase constitutive only when joined to a wild-type *LacZ*<sup>+</sup> allele, not when  $o^c$  and *Lac*<sup>+</sup> are on a separate piece of DNA. The *lacI* mutations, on the other hand, are recessive to *LacI*<sup>+</sup>, and the wild-type allele imposed the inducible state on *LacZ*<sup>+</sup> whether joined to it on the same chromosome or not, that is, it acted both in *cis* and in *trans*.

**Figure 1.** The *Lac* operon that has a single mRNA for three protein-encoding genes. Its transcription is activated by CRP protein (see text) and repressed by the repressor protein encoded by the separately transcribed *LacI*. Arrows indicate the direction of transcription/translation. The action of the repressor is cancelled by allolactose or other inducers. (b) Protein-binding sites within the *Lac* control region (dotted circle in a). The repressor binds to the three operator sites, O, O2, O3. After Ref. 3 by permission.



Based on these observations, it was postulated that *LacI* encodes a repressor protein that binds to the operator segment to block the function of the promoter and that the *LacI* protein also binds to lactose and thereby undergoes a conformational ([allosteric](#)) shift that cancels its affinity for the operator (see [Lac Repressor](#)). Later molecular analysis has confirmed this hypothesis in every particular, except that it is a lactose derivative, allolactose, not lactose itself, that binds to the repressor.

The regulation of transcription of the *Lac* operon was initially thought to be entirely through repression and release from repression, and its paradigm status was such that this was assumed to apply to gene regulation generally. That this was not the case, even for *Lac*, became clear through investigating the mechanism of repression of the operon by glucose. First, it was found that the repression is partially relieved by adding to the culture medium 2', 3' cyclic adenosine monophosphate ([cyclic AMP](#), cAMP), a metabolite that accumulates under conditions of sugar starvation. Then mutants were found that fail to respond to cAMP, and it was postulated that these fail to produce a protein with which cAMP combines to activate transcription. This protein, called CAP for catabolite activation protein (or CRP for **cAMP receptor protein**), was indeed identified. It binds to a sequence, which is a little way upstream of the promoter (Fig. 1), to activate transcription. The CRP–cAMP complex is required for transcription at all times, and when glucose is relatively abundant, cAMP levels fall to the point where activation is not effective.

Analysis of the DNA sequence of the whole *lac* operon control region combined with DNA-protein binding experiments and studies of the effects of various precisely-positioned mutations has resulted in more precise definition of the promoter, operator, and CRP-binding sequences, and also revealed further complexities (Fig. 1) (3). There are several partly overlapping sequences of 30 to 40 base pairs which conform more or less well to a consensus of *E. coli* promoters, but only one (P1) is important in the wild type. In addition to the main operator (O) there are two other short segments, quite widely spaced on either side, that are also bound by *LacI* repressor and are necessary for full repression. The highest affinity CRP-binding segment is 80 bp upstream of the transcription start point, but there is another that has lower affinity and uncertain significance, which more or less coincides with the operator.

Operon organization permits subjecting a group of functionally related genes to a common control mechanism that operates, in the case of *Lac*, at the stage of transcription. Another consequence of

common transcription was revealed by analyzing mutations that have polar effects that knock out the function of the gene in which they occur and reduce the expression of genes further downstream. Polar mutants are chain-terminating (“amber” or “ochre” [stop codons](#), or **frameshift** mutations that lead to premature stops downstream). The earlier they terminate translation, the more drastic their effect on downstream genes (4). This can be explained in the following way. Translation of mRNA in bacteria follows closely on the heels of transcription. Ribosomes attach progressively to the growing mRNA from the 5' end and begin to translate it even as it peels off the DNA template strand. There is evidence that the stability of the mRNA depends on promptly covering it with ribosomes, which become detached at termination codons. Premature polypeptide chain termination leaves mRNA unprotected by ribosomes and vulnerable to attack by **nucleases**. The severity of polar effects increases with the length of exposed mRNA, which is the distance from the premature termination point, where the ribosomes become detached, to the next initiation codon, where they reattach to the mRNA. Thus transcription and translation are interdependent. Translation obviously depends on prior transcription, and effective transcription depends on prompt translation of the transcript.

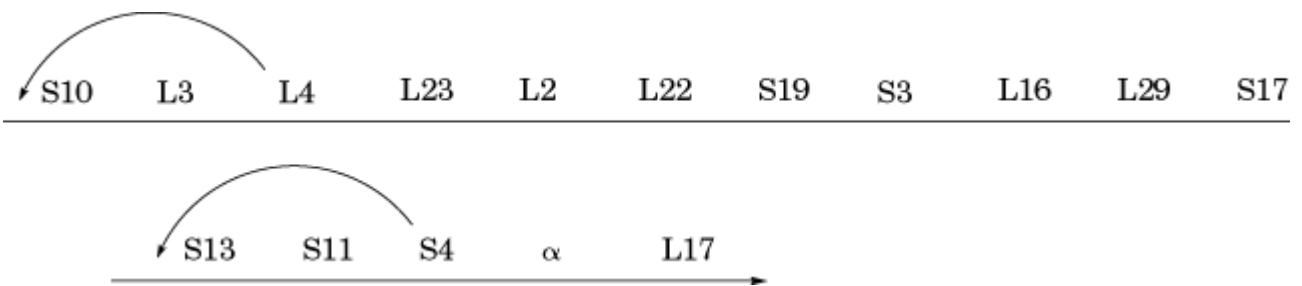
## 2. Other Bacterial Operons—Other Modes of Regulation

The combination of transcriptional activation and repression that governs the *lac* operon also applies to the [gal operon](#), which is involved in utilizing galactose. Here the repressor protein, inactivated by galactose, is encoded by a distant gene *galR*, and again glucose repression is controlled by CRP and the level of cAMP. The arabinose [ara operon](#) has its own protein (the *araR* product), which acts as an activator or a repressor depending on the presence or absence of arabinose. It was the first operon to cast doubt on the repression-only model of gene regulation.

Whereas the sugar utilization operons are regulated solely at the transcriptional level, the operons of amino acid synthesis are subject either to transcriptional or translational regulation, or both. Transcription of the tryptophan operon (see [TRP Operon](#)) is repressed by the *trpR*-encoded repressor protein in the presence of tryptophan and is derepressed when the tryptophan level drops. But there is another, more rapid way of shutting down tryptophan synthesis called **attenuation**, which operates posttranscriptionally and is mediated by the short mRNA leader sequence. Within the leader is a short **open reading frame** (ORF) that contains two successive tryptophan **codons**, a striking coincidence considering the comparative scarcity of tryptophan in proteins. Ribosomes pass through this short ORF rapidly when there is an abundance of tryptophanyl-tRNA, but they jam at the tryptophan codons when the tryptophan supply runs low. This prevents the assumption by the leader of a folding mode that would otherwise greatly reduce (**attenuate**) the translation of the rest of the operon. Similar attenuation mechanisms are found in other amino acid synthesizing operons that have runs of codons corresponding to the amino acid end product in the leader sequence in each case. Thus the **histidine operon** of *Salmonella typhimurium* (see [His Operon](#)) has eight consecutive histidine codons in its leader. Here attenuation is the sole means of regulation. There is no evidence for transcriptional repression or activation.

The six operons of *E. coli* that collectively encode all of the ribosomal proteins, together with several other components of the transcription/translation apparatus, exemplify another mode of regulation. These are regulated, at least in part, by repression of translation, and the control is exercised not by components of the growth medium (or, at least not directly) but rather by protein products of their own activity—negative feedback, the effect of which is to prevent ribosomal accumulation when growth is slowed for any reason. For each operon, one of its ribosomal protein products binds to the mRNA, generally but not always at the upstream end, to block the progress of the ribosomes (5) (Fig. 2). This stops translation of the next gene after the block and, to varying extents, that of the following genes also because translation of each gene is more or less coupled to that of the previous gene. Some of the genes, located at the downstream ends of their respective operons, escape inhibition, presumably because they recruit ribosomes independently.

**Figure 2.** The arrangements of genes within two of the *E. coli* operons that encode ribosomal proteins. L and S denote proteins of the large and small ribosomal subunit.  $\alpha$  is the  $\alpha$ -subunit of RNA polymerase. The arrowheads indicate the direction of translation of the mRNAs. The L4 and S4 proteins bind to the upstream ends of their respective mRNAs to block translation. L4 also has some effect as a transcriptional inhibitor (5).



The recently completed DNA sequence of the *Escherichia coli* genome (6), shows over 4000 genes that apparently encode polypeptide chains. About 1700 are organized into operons that have two or more genes—about 700 operons altogether. But not all sets of genes that one might think are good candidates for coregulation are linked in this way. The *E. coli* genes that encode the enzymes of arginine biosynthesis occur in scattered locations throughout the genome and, although they are all subject to a common transcriptional regulator, *argR*, they are not as tightly coupled in their activity as they would be in a common transcriptional unit. This looser type of functional grouping is sometimes called a *regulon*, and it may be appropriate in the case of the *arg* genes because the steps of arginine biosynthesis connect to a number of other metabolic pathways instead of being dedicated to a single common end product.

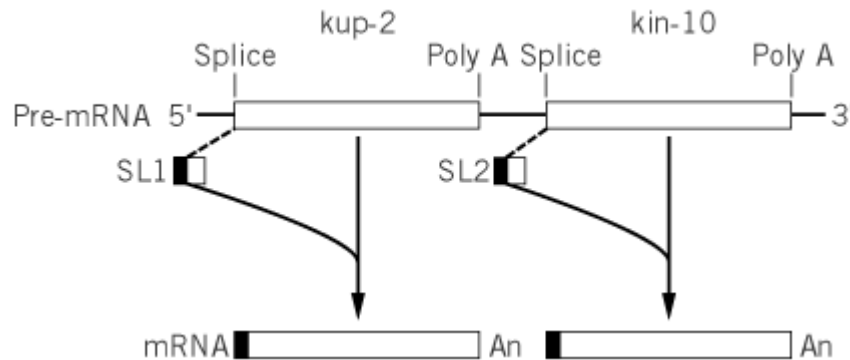
Coupling of gene activities need not be completely tight even within an operon. If the full-length multigene messenger RNA is always produced, translation is not necessarily initiated with the same efficiency at all initiation codons. The molecular yield of polypeptide chains decreases from *lacZ* to *lacY* and from *lacY* to *lacA* in the case of the *lac* operon. It seems that ribosome attachment is maximally efficient at the most upstream initiation codon and that on release from the first open reading frame the ribosomes reattach at the next initiation codon with something less than full efficiency.

### 3. Operons in *Caenorhabditis*

The generalization that operons are found only in bacteria was contradicted by the discovery that many genes in the nematode worm *Caenorhabditis elegans*, perhaps about 25% of the total, are organized in operons, in the sense of multigenic units of transcription. Another novel feature of *Caenorhabditis* gene expression, connected with the first, is that most gene messenger RNAs have leader sequences that are **spliced** after transcription (7, 8). These spliced leader sequences (SLs) are transcribed from separate genes present in multiple copies and are of several different kinds, of which SL1 is the most often used. The different gene sequences within an operon are separated by cleavage and individual splicing to yield separate mRNAs. The most upstream gene is generally spliced to SL1 but in some cases is not spliced, presumably because it has an its own adequate leader sequence. The genes further downstream are commonly spliced to SL2 or SL3, but sometimes to SL1 or one or other of at least seven less abundant SLs (8). One example is summarized in Fig. 3.

**Figure 3.** Common transcription and individual leader-sequence splicing of *Caenorhabditis elegans* genes. The individual mRNA segments are cleaved at their 5' ends and spliced to 22-nucleotide leader sequences derived from SL RNA molecules, each approximately 100 bases in length and, transcribed from multicopy genes elsewhere in the

genome, and are cleaved and polyadenylated at their 3' ends. Splicing to SL1 is characteristic of upstream gene transcripts. Those further downstream may be spliced to SL2 or some other SL. The splicing-out of any introns internal to the genes, which presumably occurs concurrently, is not depicted here.



The SL genes that occur in multicopy clusters on at least two different chromosomes have transcripts on the order of 100 nucleotides long, of which the 5' 22 nucleotides are spliced on to the mRNA precursor to provide the 5' cap site and probably essential elements of promoter sequences. The 3' ends of the cleaved operon RNA become **polyadenylated**. The splicing reaction has much in common with the mechanism used in splicing-out introns. The SL transcript contains a typical 5' splice sequence and a 3' splice sequence, typically UUUCAG, is at the cleavage sites within the operon transcript (see [RNA Splicing](#)).

The functional significance of the different SL sequences, like that of bacterial operon organization, may be to do with regulating gene expression, but here linked to development rather than to metabolism. The abundance of SL4 in mRNA, which is low in the embryo, increases sixfold during postembryonic development, mostly in hypodermal cells.

A system of very similar gene organization operates in trypanosomes (9). Many and possibly most genes in this group of parasitic protozoa are grouped in multigene transcriptional units. The transcripts are cleaved into single gene units and each of these is spliced to a common leader sequence. It is a matter of convention within the respective groups that the term operon is used for the nematodes but apparently not for the trypanosomes.

#### 4. Jointly Transcribed Nonmessenger mRNAs

Although operon-like organization of protein-encoding genes is rare in eukaryotes and is confined to nematode worms, trypanosomes, and perhaps more obscure groups yet to be investigated, it is common for several functionally distinct nonmessenger RNAs to be cut from the same primary transcript. The ribosomal RNA precursor, cleaved to yield 28S, 18S, and 5.8S molecules, is the most obvious example, but several short RNAs (snoRNAs) produced in the nucleolus and involved in processing the pre-rRNA transcripts, are themselves jointly transcribed in many eukaryotes (10). Though these multipurpose transcription units might be called operons, it is the convention to restrict the term to protein-encoding sequences.

#### Bibliography

1. E. A. Birge (1994) *Bacterial and Bacteriophage Genetics*, 3rd ed. Springer-Verlag, New York (General reference for bacterial operons).
2. F. Jacob and J. Monod (1961) *J. Mol Biol.*, **3**, 318–356.
3. W. S. Reznikoff (1992) *Mol. Microbiol.* **6**, 2419–2422.
4. W. A. Newton, J. R. Beckwith, D. Zipser, and S. Brenner (1965) *J. Mol. Biol.* **14**, 290–296.

5. J. Keener and N. Nomura (199) *Escherichia coli and Salmonella—Cellular and Molecular Biology* (F. C. Neidhardt, ed.) ASM Press, Washington D.C., pp. 1417–1431.
6. F. R. Blattner et al. (1997) *Science* **277**, 1453–1462.
7. D. A. R. Zorli, N. N. Cheng, T. Blumenthal, and J. Spieth (1994) *Nature* **372**, 270–272.
8. L. H. Ross, J. H. Freedman, and C. S. Rubin (1995) *J. Biol. Chem.* **270**, 22066–22075.
9. S. V. Graham (1995) *Parasitology Today*, **11**, 217–223.
10. D. J. Leader et al. (1997) *EMBO J.* **16**, 5742–5751.

## Optical Density

Optical density, as well as extinction, are somewhat outdated synonyms for absorbance (see [Absorption Spectroscopy](#)). The term “optical density” (OD) is still widely used as a measure of the turbidity of liquid cultures of **microorganisms**, which is a measure of their density. The increase in light scattering that accompanies growth is usually measured at 600 nm and given as OD<sub>600</sub>.

## Organizer

One of the most important ideas in developmental biology over the last 75 years was put forth by Hans Spemann based on transplantation in amphibian **embryos**. This is the idea that certain regions of the embryo, called *organizers*, are responsible for instructing other tissues to follow particular developmental pathways (1). The particular region that Spemann identified, the dorsal lip of the [blastopore](#), is usually referred to as the *Spemann organizer* in his honor.

Early in embryogenesis, when the amphibian embryo consists of a ball of cells called the **blastula**, the formation of an invagination called the *blastopore* leads to the formation of the three germ layers. These three germ layers—the endoderm, mesoderm, and ectoderm—are responsible for forming the various tissues and organs of the embryo. The ectoderm will form both the epidermis and most of the nervous system, but requires a signal from the cells of the dorsal lip of the blastopore in order to form the [neural tube](#). In a series of classic experiments, Mangold and Spemann transplanted cells from the dorsal lip of the blastopore of a donor salamander embryo to a region on the opposite side of a host embryo of a different species of salamander (the two different species were used to allow the donor and host tissues to be distinguished after the transplantation). Not only did a neural tube develop in the normal position, but a second neural tube formed near the transplanted dorsal lip cells. This second neural tube was partially formed from cells in the host embryo that had been instructed, or induced, by the transplanted cells. Similar induction events have been shown to be important in the formation of other ectodermal structures, such as the lens of the vertebrate eye (2).

Since the Spemann organizer region was first described, many experiments have been done to identify the molecules involved. A large number of molecules have been shown to induce when misexpressed, but the nature of the real signal is still in debate. Several different signaling pathways have been implicated in the Spemann organizer, including the *wingless* Wnt signaling pathway (3-5),



the [transforming growth factor](#) b pathway (6, 7), and the [fibroblast growth factor](#) pathway (8, 9). The same signaling pathways have also been shown to be involved in developmental induction in *Drosophila* (10, 11). Although the induction events in the embryo, both in amphibians and in *Drosophila*, appear to involve signaling between a number of different tissues and several interrelated cascades of gene interactions, the combination of developmental and molecular experiments has come a long way in furthering our understanding of the nature and function of organizer regions in development.

### Bibliography

1. H. Spemann (1938) *Embryonic Development and Induction*, Yale University Press, New Haven, CT.
2. W. H. Lewis (1904) *Am. J. Anat.* **3**, 505–536.
3. L. Leyns, T. Bouwmeester, S.-H. Kim, S. Picolo, and E. M. De Robertis (1997) *Cell* **88**, 747–756.
4. S. Wang, M. Krinks, K. Lin, F. P. Luyten, and M. Moos, Jr. (1997) *Cell* **88**, 757–766.
5. J.-P. Saint-Jeannet, X. He, H. E. Varmus, and I. B. Dawid (1997) *Proc. Natl. Acad. Sci. USA* **94**, 13713–13718.
6. F. Rosa, A. B. Roberts, D. Danielpour, L. L. Dart, M. B. Sporn, and I. B. Dawid (1988) *Science* **239**, 783–785.
7. A. Hemmati-Brivanlou and D. A. Melton (1992) *Nature* **359**, 609–614.
8. D. Kimelman and M. Kirschner (1987) *Cell* **51**, 869–878.
9. J. M. W. Slack, B. G. Darlington, J. K. Heath, and S. F. Godsave (1987) *Nature* **326**, 197–200.
10. P. A. Lawrence, P. Johnston, and J.-P. Vincent (1994) *Development* **120**, 3355–3359.
11. S. Morimura, L. Maves, Y. J. Chen, and F. M. Hoffmann (1996) *Dev. Biol.* **177**, 136–151.

### Suggestions for Further Reading

12. H. Spemann (1938) *Embryonic Development and Induction*, Yale University Press, New Haven, CT.
13. R. Harland and J. Gerhart (1997) Formation and function of Spemann's organizer. *Annu. Rev. Cell Dev. Biol.* **13**, 611–667.
14. R. T. Moon and D. Kimelman (1998) From cortical rotation to organizer gene expression: toward a molecular explanation of axis specification in *Xenopus*. *Bioessays* **20**, 536–545.
15. A. Hemmati-Brivanlou and D. Melton (1997) Vertebrate neural induction. *Annu. Rev. Neurosci.* **20**, 43–60.

### Origin Recognition Complex

[DNA replication](#) in bacteria or **viruses** starts at a specific region of their [chromosome](#), the [replication origin](#), and is triggered by formation of a DNA–protein complex, called the *replicator*, at that distinct site. In general, the [DNA-binding protein](#) called *initiator* recognizes the specific DNA sequence at the replication origin and nucleates formation of the DNA–protein complex in an ATP-dependent manner. Although it has been difficult to identify such a specific DNA–protein interaction for DNA replication in most **eukaryotes**, studies with a simple eukaryote, *Saccharomyces cerevisiae* (**yeast**), revealed that eukaryotes will also have the same protein–DNA interaction during initiation of DNA replication. The replication origins in yeast chromosomes have been identified as [autonomously](#)

[replicating sequences](#) (ARS), which are defined within short regions by the essential ARS consensus sequence (ACS, the A element) and its accessory elements (the B elements). Many efforts to identify the potential yeast initiator protein as the ARS binding protein have been made, and a multisubunit protein complex called the *origin recognition complex* (ORC) was finally purified from yeast nuclear extracts (1).

The purified ORC binds to the typical yeast replication origin, ARS1, in the presence of ATP, covering the 11-bp A element **consensus sequence** and neighboring B1 element. In addition to this ARS-specific binding property, ORC satisfies several other criteria as the yeast initiator protein ((2, 3): (1) ORC binds to all functional ARS. (2) The relative binding of ORC to several mutant ACS parallels their replication activities. (3) **Temperature-sensitive** mutant forms of ORC subunits result in the temporary arrest of the cell cycle in G1 phase at a **nonpermissive** temperature and make ARS **plasmids** unstable even at the permissive temperature. This instability is suppressed by increasing numbers of ARS in the plasmid, suggesting that replication initiation functions are impaired in these mutants. (4) Direct measurement of origin initiation at ARS1 by [two-dimensional gel electrophoresis](#) indicated that only a small fraction of ARS1 are active in the temperature-sensitive ORC yeast cells at the permissive temperature; on the other hand, most of the ARS1 functioned in wild-type cells. (5) The ATP-dependence of the ARS1 binding activity is consistent with the requirement of many initiator proteins for ATP in their active form at the origin sequences.

ORC is composed of six different polypeptides, called ORC1 to ORC6, with respective molecular weights of 120, 72, 62, 56, 53, and 50 kDa. All the genes encoding these subunits are essential for the viability of yeast, meaning that none of them are dispensable. In addition to DNA replication, ORC 2 and 5 were identified as factors necessary for the **silencing of mating-type** gene expression (2, 4), indicating that ORC has dual functions in DNA replication and gene silencing. The [primary structures](#) of the ORC polypeptide chains have no similarity to each other, and no typical motifs have been identified, except for strong similarity of ORC1 with one yeast replication gene product, CDC6 (Cdc18 in *Saccharomyces pombe*) (5), and purine [nucleotide-binding motifs](#) in ORC1 and ORC5. Klemm et al. (6) actually showed ATP-binding and [ATPase](#) activities in ORC regulated by ARS DNA. Binding of ORC to ARS DNA requires the presence of ATP, but not its hydrolysis, and the ARS DNA inhibits the ORC ATPase activity, indicating that ATP will stabilize ORC on the origin DNA. Since the ATP-dependent ORC binding occurs at the prereplicative stage, a mechanism for inducing the ATP hydrolysis of ARS-bound ORC is suggested to activate the replication origin.

In parallel with studies of ORC function, a dynamic protein assembly at the replication origin of yeast was studied by *in vivo* DNase I [footprinting](#) (7). The analysis demonstrated the presence of a specific DNA-protein complex at the ARS1 *in vivo*. Since the protection pattern against the **nuclease** digestion *in vivo* was nearly consistent with that obtained with purified ORC, and functional ORC is required to form the complex, it is most likely that the ORC–ACS complex observed with purified components actually exists at the replication origins in yeast cells. One important observation is a periodic change in the protection pattern in which a protection wider than the simple ORC–ACS complex appears from the end of mitosis to G1 phase and is shifted back to the narrow protection pattern at the beginning of **S phase**. Since the ORC–ACS protection pattern always exists throughout the [cell cycle](#), it was thought that ORC functions as a landing pad for additional factors, and the periodic change in the protection pattern represents the association of these factors with ORC. Since the wider protection appeared prior to initiation of replication at the replication origin and disappeared upon initiation, it represents formation of the so-called [pre-replicative complex](#), which is defined as a protein assembly to make the replication origin competent for initiation. Several replication initiation-related gene products, CDC6, CDC7, and MCM, a potential [licensing factor](#), are required to form the wider complex at the origin and are thought to be components of the pre-replication complex (8).

Accumulated knowledge about eukaryotic replication proteins inspired the idea that all eukaryotes have the same organization as replication proteins, with conserved primary structures. After the discovery of ORC in yeast, many efforts were made to identify ORC **homologues** in different

species, aiming to find their initiator proteins. Thus far, homologues have been identified in several organisms: *S. pombe*, [Arabidopsis](#), *Drosophila*, [Xenopus](#), and humans (9). The ORC2 homologue in *Xenopus* is necessary to initiate the DNA replication reaction in *Xenopus* egg extracts (10), and the homologue in *Drosophila* was genetically identified as a gene required for amplification of the **chorion** gene locus (11). In addition, ORC subunit homologues in these higher eukaryotes seem to form multisubunit complexes similar to the yeast ORC. Therefore, all eukaryotes may have a conserved mechanism to initiate DNA replication in chromosomes, in which interactions between the counterparts of ORC and ACS will be a central event.

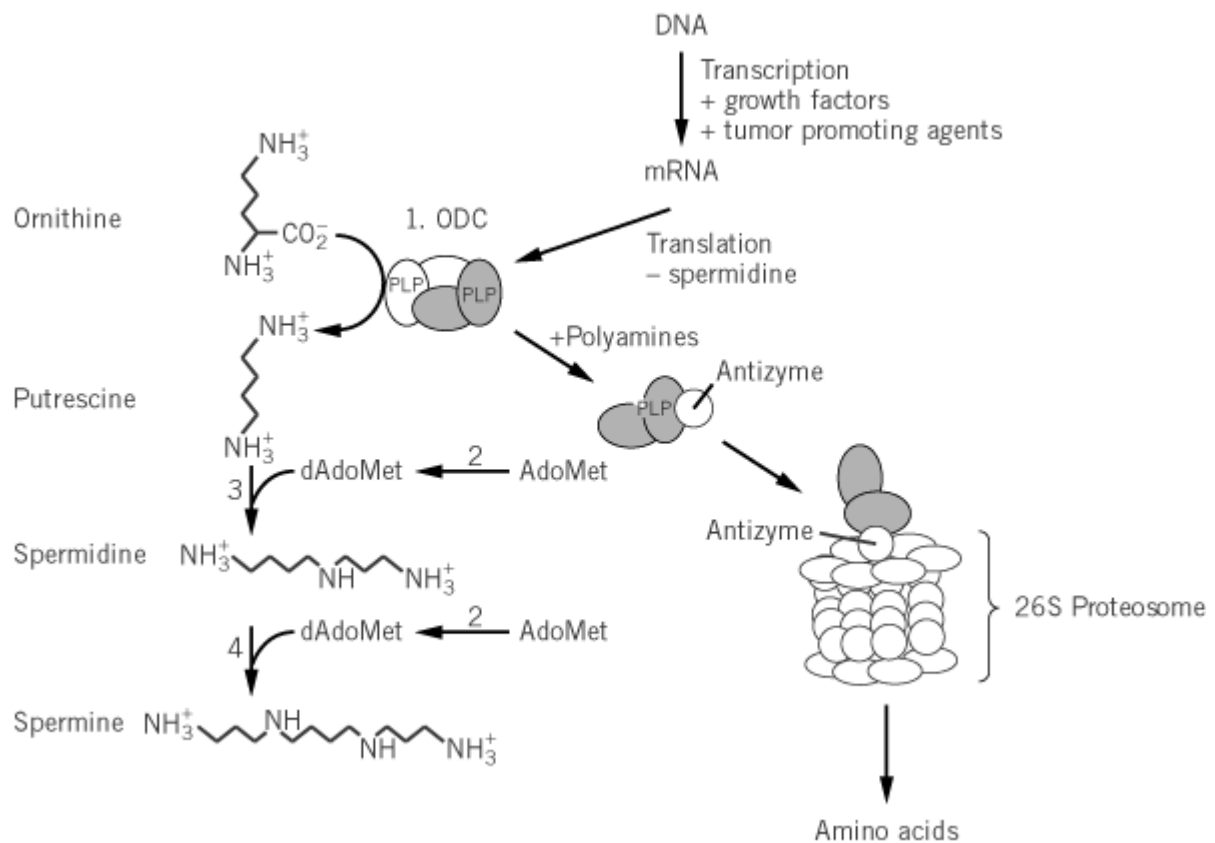
## Bibliography

1. S. P. Bell and B. Stillman (1992) *Nature* **357**, 128–134.
2. S. Loo et al. (1995) *Mol. Biol. Cell* **6**, 741–756.
3. C. A. Fox et al. (1995) *Genes Dev.* **9**, 911–924.
4. M. Foss et al. (1993) *Science* **262**, 1838–1844.
5. S. P. Bell et al. (1995) *Cell* **83**, 563–568.
6. R. D. Klemm et al. (1997) *Cell* **88**, 493–502.
7. J. F. Diffley and J. H. Cocker (1992) *Nature* **357**, 169–172.
8. J. F. Diffley et al. (1994) *Cell* **78**, 303–316.
9. K. A. Gavin et al. (1995) *Science* **270**, 1667–1671.
10. P. B. Carpenter et al. (1996) *Nature* **379**, 357–360.
11. M. F. Denissenko et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3888–3892.

## Ornithine Decarboxylase

Ornithine decarboxylase (ODC; EC 4.1.1.17), first discovered in 1968, catalyzes the **pyridoxal-5' phosphate** (PLP)-dependent decarboxylation of L-ornithine to generate putrescine (Fig. 1), the first committed step in the biosynthesis of polyamines (1-9). The polyamines putrescine, **spermidine**, and **spermine** are ubiquitous to all cells and are required for cell growth and differentiation. Consequently, ODC has been found in almost every cell type studied (4, 10-14). Overexpression of ODC or induction of its biosynthesis by growth stimuli in eukaryotic cells causes cell [transformation](#). Inhibitors of ODC arrest cell growth and are actively studied for their therapeutic potential as anticancer (1, 7, 8, 15) and antimicrobial agents (2, 5). To mediate the consequences of depletion or overexpression of ODC, higher **eukaryotes** tightly regulate levels of the enzyme through **transcriptional**, **translational**, and **posttranslational** mechanisms (reviewed in Refs. 4, 9, 16-19).

**Figure 1.** The polyamine biosynthetic pathway. Steps 1 to 4 are catalyzed by ODC, S-adenosylmethionine decarboxylase, spermidine synthase, and spermine synthase, respectively. The mechanisms of ODC regulation in mammalian cells are depicted and include transcriptional, translational, and posttranslational control. Not shown is an alternative pathway to putrescine via arginine decarboxylase present in bacteria and plants (4).



### 0.1. ODC in Tumor Promotion

Deletion of the ODC **gene** renders yeast (12), *T. brucei* (20) or mammalian (21) cells dependent on added putrescine for growth, whereas overexpression of ODC activity leads to mammalian cell transformation (22). Cell transformation caused by a number of **oncogenes**, **growth factors**, and chemical **carcinogens** is correlated with increased ODC activity. Inhibitors of ODC activity or ODC gene transcription block **oncogene**-induced transformation of rat fibroblasts (22), and 12-O-tetradecanolyphorbol-13-acetate-induced formation of tumors in epidermal cells (23, 24). Overexpression of eukaryotic **initiation factor** eIF-4E, which increases ODC translation, also causes cell transformation, and this transformation is also blocked by ODC inhibitors or by the expression of a dominant-negative mutant of ODC (25). Tissue-specific overexpression of ODC in various **transgenic** mice models causes increased tumor development in skin (26, 27), reduced fertility in testis (28), and hair loss in follicles (29), although overexpression in brain has no long-term effects on tumor incidence or neuronal degeneration (30, 31). Finally, ODC activity correlates with the acquisition of hormone-independent, poorly differentiated tumors in a rat breast cancer model (32) and of a multidrug resistance phenotype (10, 11). Thus, ODC may be useful as a clinical marker.

### 0.2. ODC as a Drug Target

A wide array of ODC inhibitors have been designed and tested, but by far the most experimentally detailed studies have been done with mechanism-based inhibitors. They and their clinical uses have been extensively reviewed (2, 4-8, 33). The best-characterized of these inhibitors, *a*-difluoromethylornithine (DFMO; eflornithine), was first described by Metcalf and co-workers. In early studies *in vitro* and in animal models, DFMO was a promising antitumor agent, but clinical studies in humans failed to demonstrate any significant impact on tumor burden. Currently, DFMO is being investigated as a chemoprevention agent for colon and cervical cancer (15, 22, 34-37).

In contrast to the poor potency of DFMO against cancer, in 1980 Bacchi and co-workers demonstrated that DFMO cures mice infected with *Trypanosoma brucei*, the causative agent of African sleeping sickness. Since this landmark discovery, DFMO has been approved for treating

Gambian trypanosomiasis in humans (reviewed in Ref. 2). The basis for the selective toxicity continues to be debated. Differences in the mammalian and *T. brucei* enzymes' ability to bind DFMO are not involved, but several metabolic differences have been found between mammalian and *T. brucei* cells that are likely to contribute to the selective toxicity:

1. Mammalian ODC is rapidly degraded intracellularly, whereas *T. brucei* ODC is stable, suggesting that differential rates of **protein degradation** are a basis for selective drug toxicity ((2), (38)).
2. Trypanosomes utilize trypanothione, a conjugate of [glutathione](#) and spermidine, to maintain cellular redox balance, and this novel role of spermidine may make the cells more sensitive to polyamine depletion (39).
3. DFMO elevates the levels of **S-adenosylmethionine** in trypanosomes, but not in mammalian cells, and may result in building up toxic S-adenosylmethionine metabolites (40).
4. Mammalian cells possess a high-throughput **transporter** for putrescine, whereas *T. brucei* cells do not (39). In contrast, the intracellular parasitic protozoa *Trypanosoma cruzi* rely on scavenged polyamines to sustain growth (39).

### 0.3. ODC Structure/Function Analysis

The genes (or **cdNA**) that encode ODC have been reported from a wide variety of organisms, including many mammalian sources, **fungi**, parasitic protozoa, and bacteria (Prosite entries PS00878 and PS00703). Interestingly, the eukaryotic and **prokaryotic** enzymes do not share significant sequence identity and differ in their three-dimensional structures. Thus, ODC activity has evolved independently at least twice (41). The eukaryotic ODC are homologous to biosynthetic arginine decarboxylase from bacteria and plants, to diaminopimelate decarboxylase from bacteria (Prosite ID PS00878) and to alanine racemase (41). In contrast, prokaryotic ODC belongs to the same family as lysine decarboxylase and biodegradative arginine decarboxylase (Prosite ID PS00703). The three-dimensional structure of bacterial ODC from *Lactobacillus*, determined by **X-ray crystallography**, has structural similarity to aspartate aminotransferase (42).

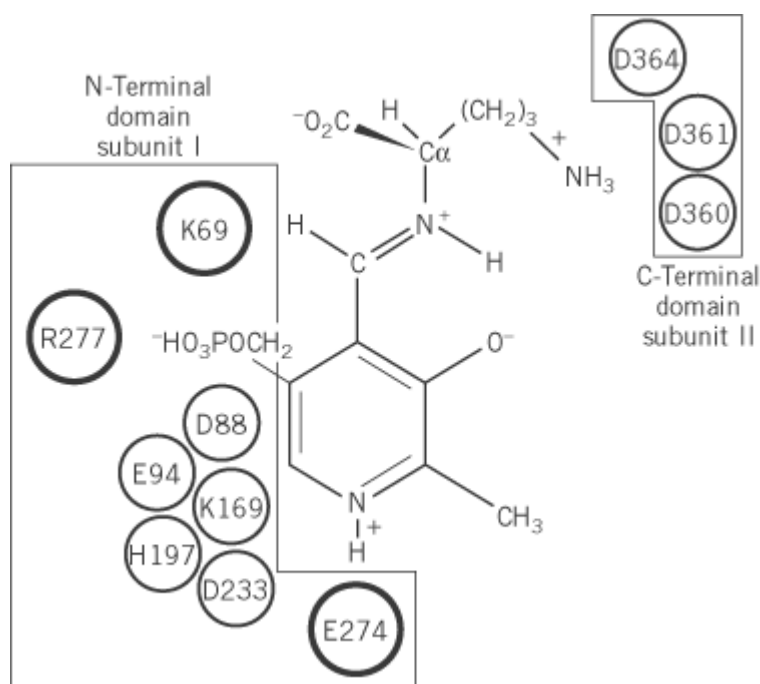
A structural model for the N-terminal domain (300 amino acids) of eukaryotic ODC predicted that it folds into a b/a barrel, and the PLP binding site formed by the C-termini of the b- **strands** is at the center of the barrel (41). This model places Lys69, His197, Asp233, Glu274, the Gly-rich loop (residues 235 to 237), and Arg277 in the PLP binding site and is strongly supported by the available biochemical data (discussed later). Now, both mouse (43) and *T. brucei* (44) ODC have been successfully crystallized, and the mouse ODC structure confirms that it folds into a b/a barrel (M. Hackert, personal communication).

ODC is an obligate homodimer, and each subunit is organized into two domains composed of a large N-terminal domain (residues 1 to 305) and a small C-terminal domain (residues 306 to 425), as demonstrated by the finding that separate polypeptides that represent these domains in *T. brucei* ODC can be coexpressed to yield an active tetramer that has kinetic properties identical to the native dimer (50). Mouse ODC may also be circularly permuted at amino acid 308 without significantly affecting activity (51). Two identical **active sites** are formed at the dimer interface and are composed of residues from the N-terminal domain of one subunit and the C-terminal domain of the other, as demonstrated by the finding that activity is restored upon mixing the inactive Lys69Ala mutant ODC with the inactive Cys360Ala ODC for the enzymes from mouse (1), *T. brucei* (45), or *L. donovani* (45). The subunits of mouse and *T. brucei* ODC are in rapid equilibrium, and heterodimers are formed immediately upon mixing two different homodimers. In contrast, the subunits from *L. donovani* ODC dissociate only after the enzyme is partially **denatured** in **urea**. Additionally, fully functional heterodimers can be formed between mouse and *T. brucei* ODC (45). This result suggests that the contacts in the dimer interface, which are essential for dimerization, are conserved between *T. brucei* and mouse ODC.

Three assays have been described to follow the decarboxylation of ornithine: a  $^{14}\text{CO}_2$  assay (1), a dye-linked **spectrophotometric** assay (45), and a **circular dichroism** assay (46). The  $k_{\text{cat}}$  for the decarboxylation of L-ornithine is typically  $10\text{s}^{-1}$ , and the  $K_m$  ranges from 90 to 400  $\mu\text{M}$  (47, 48). The enzyme is also active on lysine, but the  $K_m$  is 100-fold higher than for ornithine (48). The reaction of *T. brucei* ODC with L-ornithine proceeds through a quinoid intermediate, and product release, not decarboxylation, was determined as the rate-limiting step (49).

The roles of a number of residues in the active site of ODC have been delineated by biochemical analysis and **site-directed mutagenesis**. By reduction and peptide mapping, Lys69 was identified as the residue that forms the Schiff base with PLP, and Cys360 was identified as the site of covalent attachment of DFMO (1). Consistent with these findings, mutation of both residues significantly reduces  $k_{\text{cat}}$  (1). Mutation of Glu274 to Ala decreases  $k_{\text{cat}}$  50-fold, but wild-type activity is restored by substituting N-methyl-PLP for PLP in the reaction. This demonstrates that Glu274 interacts with the protonated pyridine nitrogen of PLP to enhance the electron-withdrawing capability of the ring (48). The Arg277 to Ala mutant has decreased PLP binding affinity and an altered  $^{31}\text{P}$ -NMR spectra of the enzyme-bound PLP (52), supporting the hypothesis that Arg277 interacts with the 5' phosphate of PLP. Mutation of Asp361 to Ala increases the  $K_m$  for ornithine 2000-fold and has little effect on  $k_{\text{cat}}$ , suggesting that Asp361 forms part of the substrate-binding pocket (48). Mutation of a number of other conserved residues also results in large decreases in enzyme activity, but the roles of these residues in catalysis remain unknown (1, 48). Taken together, these results produce a picture of the ODC active site in which PLP is bound by the large N-terminal domain and interacts with Lys69, Glu274, and Arg277, whereas the C-terminal domain of the second subunit is likely to contribute to substrate binding through interactions with residues 360 to 364. These results suggest that the substrate binds across the dimer interface (Fig. 2).

**Figure 2.** Schematic representation of the active site of eukaryotic ODC. PLP is displayed bound to ornithine via a Schiff base. Amino acids necessary for enzyme activity by site-directed mutagenesis are displayed. The three residues that have established functions are outlined in large circles.



## 0.4. ODC Regulation

### 0.4.1. Rapid Turnover of ODC and its Polyamine-Dependent Regulation

ODC from higher eukaryotes (e.g., mammals), one of the most rapidly degraded proteins known, has an intracellular half-life ranging from 10 to 60 min (the mechanism has been reviewed in Refs. [18](#) and [19](#)). Polyamines accelerate the turnover rate from the already rapid basal levels by increasing the concentration of a 26-kDa regulatory protein, termed the *antizyme*, which was first discovered for its ability to inhibit ODC activity. Polyamines regulate antizyme production by inducing ribosomal **frame-shifting** ([53](#)). The antizyme binds to the ODC monomer and targets it for ubiquitin-independent degradation by the 26 S **proteasome** ([19](#), [54](#), [55](#)). The antizyme is further regulated by the presence of antizyme inhibitor, a protein that is structurally related to ODC, but is missing many catalytic residues and cannot dimerize with native ODC ([56](#)). Antizymes that share similarity to the mammalian protein have also been found in other higher eukaryotes, for example, *Xenopus laevis* ([57](#)), where increased ODC turnover is associated with gastrulation ([58](#)). In contrast the ODC inhibitory proteins found in *Escherichia coli* are unrelated to the mammalian protein and are likely to play a very different role in ODC regulation ([59](#)).

Structural elements of ODC responsible for its rapid degradation have been delineated. Two **PEST sequences** (mouse ODC residues 298 to 333 and 423 to 461), which are associated with rapidly degraded proteins, are found in mammalian ODC ([19](#)) but not in the stable *T. brucei* enzyme ([2](#)). The C-terminal PEST sequence (423 to 416) is required for degradation and is sufficient to confer degradation on the stable *T. brucei* enzyme when expressed in mammalian cells. Mouse ODC expressed in trypanosomes is also stable, even in the presence of coexpressed rat antizyme ([60](#)), suggesting that trypanosomes lack the cellular machinery necessary to target the mouse enzyme for degradation. The requirement for the C-terminal PEST sequence for degradation in mammalian cell systems is further illustrated by the findings that mouse ODC is stabilized by truncation of as few as five residues from the C-terminus, internal deletions in the C-terminus, masking of the free C-terminus, mutation of Cys441 to Trp, or phosphorylation of ODC (reviewed in Refs. [19](#) and ([61](#))). In addition to requiring the C-terminal degradation domain, polyamine-dependent degradation requires an additional part of the sequence (residues 117 to 140 in the mouse sequence), which interacts with the C-terminal half of the antizyme (reviewed in Refs. [18](#) and [19](#)). The N-terminus of the antizyme is required to degrade mouse ODC and further confers lability when fused to the N-terminus of stable proteins in an *in vitro* assay ([62](#), [63](#)). Regulation of ODC in single-cell eukaryotes is less well studied than in mammalian cells. Yeast ODC is rapidly degraded by the proteasome ([64](#)). *Neurospora crassa* ODC is also rapidly degraded, but polyamines do not increase the degradation rate. However, in many other lower eukaryotes ODC is stable (reviewed in Ref. [19](#); see the previous discussion of *T. brucei* ODC).

### 0.4.2. Transcriptional regulation

Transcription of the mammalian ODC gene is increased in response to a number of growth stimuli. Basal levels of ODC transcription are mediated by several DNA elements in a complex cell type-specific manner, for example, the cAMP-responsive element CRE ([65](#)) and Sp1 ([66](#)). ODC transcription is stimulated above basal levels by protein kinase A through CRE in adrenal carcinoma cells ([67](#)), by c-Fos in PC12 cells ([68](#)), by c-Myc ([69](#)), and by the c-Myc.Max protein complex ([70](#), [71](#)). Overexpression of ODC in a mouse **myeloma** cell line ([72](#)) or in mouse myeloid 32D.3 cells induces **apoptosis**, and ODC is a mediator of c-Myc-induced apoptosis in the latter cells ([73](#)). The **transcription factor** AP-2 acts as a negative regulator of Myc-induced ODC gene transcription ([74](#)). The Wilms' **tumor suppressor** represses activity of the ODC promoter and suggests that tumor suppressors may have a role in maintaining normal ODC levels ([75](#)). Rotenoids ([23](#)) and tyrosine kinase inhibitors ([24](#)) inhibit 12-*O*-tetradecanolyphorbol-13-acetate-induced ODC transcription and may have utility as chemoprevention agents.

### 0.4.3. Translational regulation

Although many factors influence the levels of ODC **messenger RNA**, it is well established that the levels of mRNA do not always correspond to the ODC activity observed. For example, the growth-related changes in ODC mRNA are 10-fold less than the corresponding ODC activity upon

stimulation of Ehrlich ascites tumor cells. In quiescent cells high mRNA levels contrast with low protein levels. These observations suggest that ODC levels are controlled partially at the level of [translation](#) (76). It is predicted that the 5' untranslated region (5'UTR) of the mammalian ODC gene folds into a stable secondary structure, and deletion analysis demonstrates that the 5'UTR inhibits translation of both ODC and heterologous gene products (reviewed in Ref. 1). [Alternative splicing](#) of the 5'UTR relieves the translational suppression of ODC in tumor-derived pancreatic acinar cells and results in constitutive expression (77). The polyamines spermidine and spermine inhibit ODC translation, and this inhibition depends on the presence of the 5'UTR. Insulin-induced ODC translation also depends on the 5'UTR (1, 78). ODC activity or the activity of a reporter construct linked to the 5'UTR are greater in cells that overexpress eIF-4E, a eukaryotic transcription factor that has RNA **helicase** activity which unwinds the secondary structure of RNA. Increased ODC expression in these cells results from changes in polyamine regulation and decreased translational repression (79, 80).

ODC is an essential enzyme that promotes cell growth and controls differentiation. Consistent with its important function, the enzyme is highly regulated in mammalian cells. The continued study of both ODC structure/function and of ODC regulation should lead to a greater understanding the role of ODC in cell function and to the design of additional methods to combat cancer and infectious disease.

#### Bibliography

1. A. E. Pegg, L. M. Shantz, and C. S. Coleman (1994) *Biochem. Soc. Trans.* **22**, 846–852.
2. C. C. Wang (1995) *Annu. Rev. Pharmacol. Toxicol.* **35**, 93–127.
3. T. Oka and F. Borellini (1989) In *Ornithine Decarboxylase: biology, enzymology and molecular genetics* (S. Hayashi, ed.), Pergamon, New York, pp. 7–20.
4. P. P. McCann and A. E. Pegg (1992) *Pharmacol. Ther.* **54**, 195–215.
5. L. J. Marton and A. E. Pegg (1995) *Annu. Rev. Pharmacol. Toxicol.* **35**, 55–91.
6. C. W. Tabor and H. Tabor (1984) *Annu. Rev. Biochem.* **53**, 749–790.
7. J. Janne and L. Alhonen-Hongisto (1989) In *Ornithine Decarboxylase: biology, enzymology and molecular genetics* (S. Hayashi, ed.), Pergamon, New York, pp. 59–85.
8. J. Janne, L. Alhonen, and P. Leinonen (1991) *Ann. Med.* **23**, 241–259.
9. R. H. Davis, D. R. Morris, and P. Coffino (1992) *Microbiol. Rev.* **56**, 280–290.
10. Y. G. Assaraf, S. Drori, U. Bachrach, and V. Shaugan-Labay (1994) *Anal. Biochem.* **216**, 97–109.
11. U. Bachrach, A. Shayovitz, Y. Marom, A. Ramu, and N. Ramu (1994) *Cell. Mol. Biol.* **40**, 957–964.
12. C. W. Tabor and H. Tabor (1985) *Microbiol. Rev.* **49**, 81–98.
13. A. E. Pegg (1989) In *Ornithine Decarboxylase: biology, enzymology, and molecular genetics* (S. Hayashi, ed.), Pergamon, New York, pp. 21–34.
14. S. Merali and A. B. Clarkson (1996) *Antimicrob. Agents Chemother.* **40**, 973–978.
15. F. L. Meyskens and E. W. Gerner (1995) *J. Cell. Biochem.* **22**, 126–131.
16. L. Wen, J. K. Huang, and P. J. Blackshear (1989) *J. Biol. Chem.* **264**, 9016–9021.
17. A. E. Pegg, L. M. Shantz, and C. S. Coleman (1994) *Biochem. Soc. Trans.* **22**, 846–852.
18. S. Hayashi and Y. Murakami (1995) *Biochem. J.* **306**, 1–10.
19. S. Hayashi, Y. Murakami, and S. Matsufuji (1996) *Trends Biochem. Sci.* **21**, 27–30.
20. F. Li, S-b. Hua, C. C. Wang, and K. M. Gottesdiener (1996) *Mol. Biochem. Parasitol.* **78**, 227–236.
21. F. Svensson and L. Persson (1996) *Mol. Cell. Biochem.* **162**, 113–119.
22. A. E. Pegg, L. M. Shantz, and C. S. Coleman (1995) *J. Cell. Biol.* **22**, 132–138.



23. C. Gerhauser, M. Woongchoon, S. K. Lee, N. Suh, Y. Luo, J. Kosmeder, L. Luyeng, H. H. S. Fong, A. D. Kinghorn, R. M. Moriarty, R. G. Mehta, A. Constantinou, R. C. Moon, and J. M. Pezzuto (1995) *Nat. Med.* **1**, 260–266.
24. C. Tseng and A. Verma (1996) *Mol. Pharmacol.* **50**, 249–257.
25. L. M. Shantz and A. E. Pegg (1996) *Cancer Res.* **56**, 5136–5140.
26. T. G. O'Brien, L. C. Megosh, G. Gilliard, and A. P. Soler (1997) *Cancer Res.* **57**, 2630–7.
27. L. Megosh, S. K. Gilmour, D. Rosson, A. P. Soler, M. Blessing, J. A. Sawicki, and T. G. O'Brien (1995) *Cancer Res.* **55**, 4205–4209.
28. H. Hakovirta, A. Keiski, J. Toppari, M. Halmekyto, L. Alhonen, J. Janne, and M. Parvinen (1993) *Mol. Endocrinol.* **7**, 1430–1436.
29. A. P. Soler, G. Gilliard, L. C. Megosh, and T. G. O'Brien (1996) *J. Invest. Dermatol.* **106**, 1108–1113.
30. L. Alhonen, M. Halmekyto, V. M. Kosma, J. Wahlfors, R. Kauppinen, and J. Janne (1995) *Int. J. Cancer* **63**, 402–404.
31. R. A. Kauppinen and L. I. Alhonen (1995) *Prog. Neurobiol.* **47**, 545–63.
32. A. Manni, R. Grove, S. Kunselman, and M. Aldaz (1995) *Cancer Lett.* **92**, 49–57.
33. P. Bey, C. Danzin, and M. Jung (1987) In *Inhibition of Polyamine Metabolism. Biological Significance and Basis for New Therapies* (P. P. McCann, A. E. Pegg, and A. Sjoerdsma, eds.), Academic Press, Orlando, FL, pp. 1–32.
34. B. S. Reddy (1996) *Prev. Med.* **25**, 48–50.
35. C. L. Loprinzi, E. M. Messing, J. R. O'Fallon, M. A. Poon, R. R. Love, S. K. Quella, D. L. Trump, R. F. Morton, and P. Novotny (1996) *Cancer Epidemiol. Biomarkers Prev.* **5**, 371–374.
36. G. J. Kelloff, C. W. Boone, J. A. Crowell, V. E. Steele, R. A. Lubet, L. A. Doody, W. F. Malone, E. T. Hawk, and C. C. Sigman (1996) *J. Cell Biochem. Suppl.* **26**, 1–28.
37. M. F. Mitchell, W. K. Hittelman, R. Lotan, K. Nishioka, G. Tortolero-Luna, R. Richards-Kortum, J. T. Wharton, and W. K. Hong (1995) *Cancer* **76**, 1956–1977.
38. J. L. A. Mitchell, C. Choe, G. G. Judd, D. J. Daghfal, R. J. Kurzeja, and A. Leyser (1996) *Biochem. J.* **317**, 811–816.
39. A. H. Fairlamb and S. A. Le Quesne (1997) In *Trypanosomiasis and Leishmaniasis* (Hide, Mottram, Coombs and Holmes, eds.), CAB International, pp. 149–161.
40. C. J. Bacchi (1993) *Parasitology Today* **9**, 190–193.
41. N. V. Grishin, M. A. Phillips, and E. J. Goldsmith (1995) *Protein Sci.* **4**, 1291–1304.
42. C. Momany, S. Ernst, R. Ghosh, N. L. Chang, and M. L. Hackert (1995) *J. Mol. Biol.* **6**, 643–655.
43. A. Kern, M. A. Oliveira, N. L. Chang, S. R. Ernst, D. W. Carroll, C. Momany, K. Minard, P. Coffino, and M. L. Hackert (1996) *Proteins* **24**, 266–268.
44. N. V. Grishin, A. L. Osterman, E. J. Goldsmith, and M. A. Phillips (1996) *Proteins Struct. Function Genet.* **24**, 272–273.
45. A. Osterman, N. V. Grishin, L. N. Kinch, and M. A. Phillips (1994) *Biochemistry* **33**, 13662–13667.
46. H. B. Brooks and M. A. Phillips (1996) *Anal. Biochem.* **238**, 191–194.
47. C. S. Coleman, B. A. Stanley, and A. E. Pegg (1993) *J. Biol. Chem.* **268**, 24572–24579.
48. A. Osterman, L. N. Kinch, N. V. Grishin, and M. A. Phillips (1995) *J. Biol. Chem.* **270**, 11797–11802.
49. H. B. Brooks and M. A. Phillips (1997) *Biochemistry*, **36**, 15147–15155.
50. A. L. Osterman, D. V. Lueder, M. Quick, D. Myers, B. J. Canagarajah, and M. A. Phillips (1995) *Biochemistry* **34**, 13431–13436.
51. X. Li and P. Coffino (1993) *Mol. Cell. Biol.* **13**, 2377–2383.

52. A. L. Osterman, H. B. Brooks, J. Rizo, and M. A. Phillips (1997) *Biochemistry* **36**, 4558–4567.
53. S. Matsufuji, T. Matsufuji, N. M. Wills, R. F. Gesteland, and J. F. Atkins (1996) *EMBO J.* **15**, 1360–1370.
54. S. Elias, B. Bercovich, C. Kahana, P. Coffino, M. Fischer, W. Hilt, D. H. Wolf, and A. Ciechanover (1995) *Eur. J. Biochem.* **229**, 276–283.
55. Y. Murakami, N. Tanahashi, K. Tanaka, S. Omura, and S. Hayashi (1996) *Biochem. J.* **317**, 77–80.
56. Y. Murakami, T. Ichiba, S. Matsufuji, and S. Hayashi (1996) *J. Biol. Chem.* **271**, 3340–3342.
57. T. Ichiba, S. Matsufuji, Y. Miyazaki, and S. Hayashi (1995) *Biochim. Biophys. Acta* **1262**, 83–86.
58. U. Rosander, I. Holm, B. Grahn, H. Lovtrup-Rein, M.-O. Mattsson, and O. Heby (1995) *Biochim. Biophys. Acta* **1264**, 121–128.
59. E. S. Canellakis, A. A. Paterakis, S. C. Huang, C. A. Panagiotidis, and D. A. Kyriakidis (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7129–7133.
60. S. Hua, X. Li, P. Coffino, and C.C. Wang (1995) *J. Biol. Chem.* **270**, 10264–10271.
61. S. G. Reddy, S. M. McIlheran, B. J. Cochran, L. L. Woreth, L. A. Bishop, P. J. Brown, V. P. Knutson, and M. K. Haddox (1996) *J. Biol. Chem.* **271**, 24945–24953.
62. X. Li, B. Stebbins, L. Hoffman, G. Pratt, M. Rechsteiner, and P. Coffino (1996) *J. Biol. Chem.* **271**, 4441–4446.
63. X. Li and P. Coffino (1996) *J. Biol. Chem.* **271**, 4447–4451.
64. E. Mamroud-Kidron and C. Kahana (1994) *FEBS Lett.* **356**, 162–164.
65. J. J. Palvimo, M. Partanen, and O. A. Janne (1996) *Biochem. J.* **316**, 993–998.
66. A. P. Kumar, P. K. Mar, B. Zhao, R. L. Montgomery, D. C. Kang, and A. P. Butler (1995) *J. Biol. Chem.* **270**, 4341–4348.
67. M. S. Abrahamsen, R. S. Li, W. Dietrich-Goetz, and D. R. Morris (1992) *J. Biol. Chem.* **267**, 18866–18873.
68. C. Wrighton and M. Busslinger (1993) *Mol. Cell. Biol.* **13**, 4657–4669.
69. C. Bello-Fernandez, G. Packham, and J. L. Cleveland (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7804–7808.
70. K. E. Tobias, J. Shor, and C. Kahana (1995) *Oncogene* **11**, 1721–1727.
71. M. Selvakumaran, D. Liebermann, and B. Hoffman (1996) *Blood* **88**, 1248–1255.
72. K. E. Tobias and C. Kahana (1995) *Cell Growth and Differentiation* **6**, 1279–1285.
73. G. Packman and J. L. Cleveland (1994) *Mol. Cell. Biol.* **14**, 5741–5747.
74. S. Gaubatz, A. Imhof, R. Dosch, O. Werner, P. Mitchell, R. Buettner, and M. Eilers (1995) *EMBO J.* **14**, 1508–1519.
75. J. A. Mosheir, M. Skunca, W. Wu, S. M. Boppana, F. J. Rauscher, and J. Dosesu (1996) *Nucleic Acids Res.* **24**, 1149–1157.
76. U. M. Wallon, L. Person, and O. Heby (1995) *Mol. Cell. Biochem.* **146**, 39–44.
77. S. Pyronnet, S. Vagner M. Bouisson, A.C. Prats, N. Vaysse, and L. Pradayrol (1996) *Cancer Res.* **56**, 1742–1745.
78. M. Fogelpetrovic, S. Vujcic, J. Miller, and C. W. Porter (1996) *FEBS Lett.* **391**, 89–94.
79. L. M. Shantz, R.-H. Hu, and A. E. Pegg (1996) *Cancer Res.* **56**, 3265–3269.
80. D. Rousseau, R. Kaspar, I. Rosenwald, L. Gehrke, and N. Sonenberg (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1065–1070.

## Orphans

*Orphans* are a class of dispersed, solitary genetic elements of the [genome](#) derived from tandem multigene **families**. These displaced elements arise from genes that do or do not code for proteins, including those of [histones](#) in the sea urchin *Lytechinus pictus*; **ribosomal** genes in *Saccharomyces cerevisiae*; ribosomal genes and H3 histone genes in *Drosophila melanogaster* (1); the amylase multigene family of mice (2); the spliced-leader small RNA gene family of *Trypanosoma brucei* (3) and of the nematode *Angiostrongylus cantonensis* (4); and the human 28 S ribosomal DNA (5).

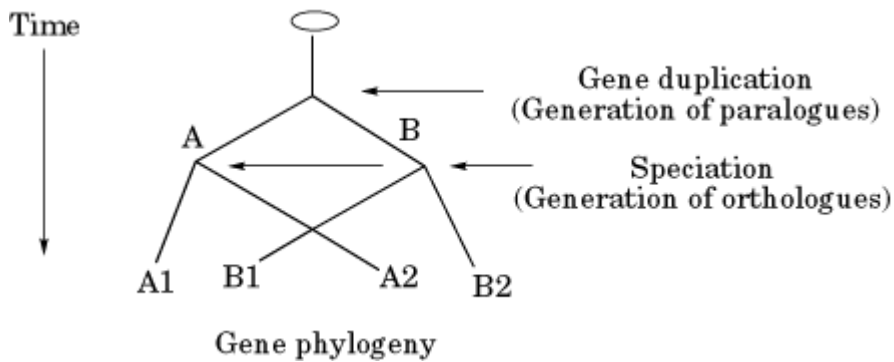
An interesting case is that of the human [immunoglobulin](#) V kappa gene regions. They have been **transposed** during [evolution](#) from the site of the kappa locus on **chromosome** 2 to chromosomes 1, 22, and other chromosomes. These orphans are very similar and may have derived from a single ancestral gene. The DNA sequences at the junction of chromosome 22 and other orphon regions are direct and inverted repeats and in one case an **Alu sequence** repeat. These unusual features may have predisposed the orphon regions to transposing by serving as target sites for enzymes involved in [recombination](#) (6).

### Bibliography

1. G. Childs et al. (1981) *Cell* **23**, 651–663.
2. S. Bodary et al. (1985) *J. Mol. Biol.* **182**, 1–10.
3. M. Parsons, R. G. Nelson, and N. Agabian (1986) *Nucleic Acids Res.* **14**, 1703–1718.
4. G. W. Joshua, F. B. Perler, and C. C. Wang (1995) *Nucleic Acids Res.* **23**, 1030–1035.
5. J. Munro, R. H. Burdon, and D. P. Leader (1986) *Gene* **48**, 65–70.
6. P. Borden, R. Jaenichen, and H. G. Zachau (1990) *Nucleic Acids Res.* **18**, 2101–2107.

## Orthologous Genes

Orthologous genes are two or more **genes** whose follows the speciation process, and they have usually the same function in each resulting species; that is, a particular gene in the ancestral will be duplicated in the two different species following speciation and will usually be maintained in each, because both species will require the function of that gene (Fig. 1) (1). These genes differ from , which diverged before or after speciation as a result of and usually have different functions. For example, the genes for the a and b polypeptide chains of hemoglobin and that for myoglobin are **homologous** but paralogous, whereas the hemoglobin a chains from different species are orthologous; similarly, the b chains and myoglobins from different species are orthologous. Only comparisons of orthologous genes can be used to reconstruct species **phylogenies**.



It is often difficult to distinguish orthologous genes from paralogous genes. However, analysis has shown that some genes have a single version in a genome. Therefore, if this is the case in two or more species, there is a strong possibility that these genes are orthologous. If the gene is a member of a known family, there is a distinct probability that it is a paralogous gene, depending on the number of genes in the family, because the gene could be a product of gene duplication before or after speciation.

Orthologous gene pairs are useful for estimating the rate of nucleotide substitution and species phylogenies. When the time since two species diverged ( $T_s$ ) is known from the fossil record, it follows that the orthologous pairs should also have diverged at the time of species divergence. Therefore, the rate of nucleotide substitution for the gene ( $\nu$ ) can be estimated by dividing the number of nucleotide substitutions ( $K_n$ ) between orthologous gene pairs by two times of the divergence time between the two species  $\nu = K_n/(2T_s)$ .

Two homologous genes sampled from the gene families between two species are often paralogous, rather than orthologous. If the two genes are paralogous, the divergence time of these genes ( $T_g$ ) must have been earlier than the speciation time ( $T_g > T_s$ ). If the divergence time of those genes is mistakenly assumed to be the divergence time between the species, the divergence time of the genes will be underestimated and will lead to overestimation of the rate of nucleotide substitution ( $K_n/(2T_g) > \nu$ ). Thus, it is important to use only orthologous genes in estimating the rate of nucleotide substitution and the times of speciation.

There have been many recent advances in genome analysis, and as a result complete genomes are known. Those of two species can be compared to determine whether the physical order of genes in the genomes is conserved between different species. To accomplish this, it is vital to distinguish orthologous genes from paralogous genes. By focusing on orthologous genes, it is possible to trace changes in gene location during evolution because of the one-to-one correspondence in gene location between genomes, unless gene duplication takes place after speciation.

#### Bibliography

“Orthologous Genes” in , Vol. 3, p. 1731, by T. Gojobori; “Orthologous Genes” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. H. Watanabe, H. Mori, T. Itoh, and T. Gojobori, (1997) *J. Mol. Evol.* **44**, 57–64.

## Oskar Gene

The *oskar* is a **maternal effect** gene in *Drosophila melanogaster* whose primary function is to organize a specialized cytoplasm, termed the germ plasm or pole plasm, at the posterior pole of the maturing egg (1). This germ plasm contains factors required for abdomen formation and germ-cell formation. The *oskar* messenger RNA is synthesized during oogenesis. It is deposited into the egg, where it becomes localized to the posterior pole through sequences in its 3'-untranslated region (UTR). Only posteriorly localized RNA is translated into Oskar protein. The site of Oskar protein synthesis determines where germ plasm is assembled, and the quantity of Oskar protein determines how many germ cells form. *oskar* is a member of the posterior group of genes (named for their effect on development of the abdomen in the posterior half of the embryo). Most of the posterior group genes facilitate localization and translation of *oskar* mRNA, while a small number of posterior genes act together with Oskar in the assembly of the germ plasm or are specifically required for either germ cell or abdomen formation.

### 1. *oskar* Protein and RNA Structure

*oskar* encodes a 2.9-kb mRNA that is expressed in the germ line of the adult female and is present at the posterior pole of early embryos. *oskar* mRNA has a 14-nucleotide 5' UTR and a 1045-nucleotide 3' UTR. Two different in-frame start codons are used to produce Oskar protein isoforms of 606 amino acids residues (69 kDa) and 468 residues (54 kDa). The short isoform has full *oskar* activity, although the long isoform is necessary to maintain localization of the *oskar* mRNA at the posterior of the oocyte (2, 3). The 3' UTR contains regulatory elements for RNA localization and translation, while the region between the two start codons contains sequences necessary for translational regulation (4). The Oskar protein is novel, in that no known sequence or structural motifs are present. The only known homologue is from *Drosophila virilis* (5).

### 2. *oskar* RNA Localization and Translation

Oskar protein synthesis is regulated at the level of mRNA localization and translation and occurs during oogenesis. Each *Drosophila* ovary consists of approximately 16 ovarioles, each of which contains a series of developing egg chambers. At the tip of each ovariole, a germ-line stem cell produces a new stem cell and a cystoblast by asymmetric cell division. Each cystoblast undergoes four rounds of division with incomplete cytokinesis, giving rise to 15 nurse cells and the oocyte. The nurse cell–oocyte cluster is interconnected by cytoplasmic bridges called ring canals and surrounded by somatic follicle cells. *oskar* RNA is produced in the nurse cells and deposited into the developing oocyte, where it accumulates from the earliest stages of oogenesis. By late stage 8 of oogenesis, the RNA is concentrated at both the anterior margin and the posterior pole of the oocyte. By stage 10, *oskar* RNA is localized exclusively to the posterior pole, where it remains until it is degraded in the early embryo after germ-cell formation [see Fig. 1 (6, 7)]. Oskar protein is translated once the mRNA becomes localized to the posterior pole, while unlocalized RNA is **translationally repressed**. This ensures that protein is made only at the posterior pole of the oocyte and embryo (2, 3).

**Figure 1.** A *Drosophila melanogaster* ovariole. *oskar* mRNA is detected by *in situ* hybridization. In the earlier stages of oogenesis (top of figure), *oskar* mRNA accumulates in the oocyte. As each egg chamber matures, the *oskar* mRNA is localized within the oocyte to the posterior pole (center and bottom of figure). Anterior is to the top, posterior is to the bottom.



The site where *oskar* RNA is localized and Oskar protein is synthesized determines where germ cells form and the posterior body region, the abdomen, develops. Therefore processes that lead to the localization of *oskar* RNA are directly linked to the establishment of anterior–posterior polarity within the oocyte. All aspects of *oskar* RNA localization are mediated by sequences within the *oskar* 3' UTR (8). Early aspects of *oskar* RNA transport into the oocyte and localization to the posterior pole of the oocyte require the correct organization of the [microtubule](#) network. Genes involved in the polarization of the early microtubule network that extends from the nurse cells to the oocyte also affect the transport of *oskar* into the oocyte. These genes include *egalitarian*, *Bicaudal-D*, ([coiled-coil](#) domain protein), *orb* (an RNA-recognition motif protein with similarity to mammalian CEBP), and cytoplasmic [dynein](#) (9). Subsequently, the establishment of anterior–posterior polarity within the oocyte is controlled by cell–cell signaling between the oocyte-bound [Gurken](#) protein, a member of the **transforming growth factor-a** (TGF-A) family, and its receptor, the *Drosophila* **EGF receptor** (*DER* or *torpedo*). Mutations in either gene disrupt localization of *oskar* RNA to the posterior pole, presumably by altering the polarity of the microtubule network within the oocyte. Indeed, inhibitors of the microtubule, but not of the [microfilament](#), cytoskeleton disrupt the localization of *oskar* RNA during oogenesis (10). Other maternal effect genes that affect *oskar* localization and affect the microtubule organization in the oocyte include *cappuccino* (for min homologue), *spire*, and *mago nashi* (a novel protein that is highly conserved from plants to humans). A role of microfilaments in *oskar* mRNA localization is suggested by the effect that mutations in the tropomyosin II gene (an **actin-binding** protein) have on *oskar* RNA localization (11). It has been proposed that the [actin](#) cytoskeleton is required for high levels of *oskar* mRNA localization or the firm anchoring of the RNA at the posterior pole (11).

Staufen, a double-stranded [RNA-binding protein](#), has been proposed as the best candidate to mediate directly the localization of *oskar* RNA. *staufen* mutants do not affect the organization of the cytoskeleton, and Staufen protein and *oskar* RNA colocalize during oogenesis (12). A direct physical interaction between Staufen and *oskar* RNA has not been demonstrated. Staufen is also required for localization of [bicoid](#) mRNA to the anterior of the oocyte and of *prospero* RNA to the apical side of dividing neuroblasts (12).

*oskar* mRNA that is not localized to the posterior pole is translationally repressed. If the RNA is translated prematurely, posterior body patterning will be specified throughout the embryo. Translational repression is mediated through specific RNA-binding sites in the 3' UTR of *oskar* mRNA. The RNA-binding protein Bruno (which corresponds to the *arrest* gene independently identified on the basis of its early arrest maternal effect phenotype) binds specifically to Bruno response elements (BREs) in the *oskar* RNA (13). The bicaudal-C protein (a KH-domain protein) has also been implicated in repressing translation of unlocalized *oskar* RNA. Only at the posterior pole, where *oskar* RNA is localized, is the RNA translationally active. Activation of translation depends on derepression mediated through sequences between the two alternative [start codons](#) of the *oskar* mRNA (4). Translation of *oskar* RNA or the stability of Oskar protein is further dependent upon the gene *aubergine* as well as upon the Vasa protein, an [RNA helicase](#) that interacts physically with Bruno protein (14, 15).

### 3. Role of *oskar* in Germ Plasm Assembly

Homozygous mutant *oskar* females are viable, but produce embryos that lack germ plasm. These embryos fail to form germ cells and do not develop an abdomen. Mislocalization of the *oskar* coding region (via the RNA localization element in the 3' UTR of [bicoid](#) mRNA) to the anterior of the oocyte results in the assembly of ectopic germ plasm, a mirror image duplication of the abdomen, and formation of ectopic germ cells (16). Thus *oskar* is both necessary and sufficient for assembly of germ plasm and for formation of both the abdomen and germ cells.

Assembly of the germ plasm is also dependent upon other factors that act downstream of *oskar*. Vasa is a **DEAD-box** helicase that physically interacts with Oskar. The products of the *tudor* and *valois*

genes are also required for germ plasm assembly. Vasa, Tudor, and Oskar are components of electron-dense particles, known as polar granules, that are found only in the germ plasm (17).

The effect of *oskar* on abdomen formation is mediated through the **nanos** gene. Oskar is responsible for localizing *nanos* mRNA to the posterior pole of the embryo, where Nanos protein is translated and acts to facilitate abdomen formation through repression of the *hunchback* mRNA. While the mechanism of *nanos* localization to the posterior pole is not clear, *oskar* is both necessary and sufficient for *nanos* localization. Overexpression of *oskar* or mislocalization of the *oskar* mRNA results in an excess of *nanos* activity at the site of *oskar* localization.

*oskar* also recruits several factors to the posterior pole that play a role in germ-cell formation. These include: *germ cell-less* mRNA, which encodes a protein that associates with the nuclear envelope; mitochondrial large ribosomal RNA (*mtlrRNA*) that is found outside of the mitochondria, where it associates with polar granules during the early stages of embryogenesis; and *polar granule component 1* (*pgc1*), which codes for an untranslated RNA. A number of different experiments suggest that *germ cell-less* and *mtlrRNA* are required for germ-cell formation (18, 19), while **antisense** experiments with *pgc1* suggest a role for this gene in germ-cell migration or development (20). Other maternal factors localized to the posterior pole in an *oskar*-dependent manner include maternal *hsp83* mRNA, whose RNA is stabilized in the posterior depending on the germ plasm; the cell cycle regulator *cyclin B*; *outspread*; and the Fat-facets protein that is a sequence-specific deubiquitinating **proteinase**. The roles of these genes in germ-cell formation, migration, or development are unclear.

## Bibliography

1. R. Lehmann and C. Nüsslein-Volhard (1986) *Cell* **47**, 141–52.
2. C. Rongo, E. R. Gavis, and R. Lehmann (1995) *Development* **121**, 2737–2746.
3. F. H. Markussen, A. M. Michon, W. Breitwieser, and A. Ephrussi (1995) *Development* **121**, 3723–3732.
4. N. Gunkel, T. Yano, F. H. Markussen, L. C. Olsen, and A. Ephrussi (1998) *Genes Dev.* **12**, 1652–1664.
5. P. J. Webster, J. Suen, and P. M. Macdonald (1994) *Development* **120**, 2027–2037.
6. J. Kim-Ha, J. L. Smith, and P. M. Macdonald (1991) *Cell* **66**, 23–35.
7. A. Ephrussi, L. K. Dickinson, and R. Lehmann (1991) *Cell* **66**, 37–50.
8. J. Kim-Ha, P. J. Webster, J. L. Smith, and P. M. Macdonald (1993) *Development* **119**, 169–178.
9. M. McGrail and T. S. Hays (1997) *Development* **124**, 2409–2419.
10. N. J. Pokrywka and E. C. Stephenson (1995) *Dev. Biol.* **167**, 363–370.
11. M. Erdelyi, A. M. Michon, A. Guichet, J. B. Glotzer, and A. Ephrussi (1995) *Nature* **377**, 524–527.
12. D. St. Johnston, D. Beuchle, and C. Nüsslein-Volhard (1991) *Cell* **66**, 51–63.
13. J. Kim-Ha, K. Kerr, and P. M. Macdonald (1995) *Cell* **81**, 403–412.
14. P. J. Webster, L. Liang, C. A. Berg, P. Lasko, and P. M. Macdonald (1997) *Genes Dev* **11**, 2510–2511.
15. J. E. Wilson, J. E. Connell, and P. M. Macdonald (1996) *Development* **122**, 1631–1639.
16. A. Ephrussi and R. Lehmann (1992) *Nature* **358**, 387–392.
17. W. Breitwieser, F. H. Markussen, H. Horstmann, and A. Ephrussi (1996) *Genes Dev.* **10**, 2179–2188.
18. T. A. Jongens, B. Hay, L. Y. Jan, and Y. N. Jan (1992) *Cell* **70**, 569–584.
19. S. Kobayashi, R. Amikura, and M. Okada (1993) *Science* **260**, 1521–1524.
20. A. Nakamura, R. Amikura, M. Mukai, S. Kobayashi, and P. F. Lasko (1996) *Science* **274**, 2075–2079.



### Suggestions for Further Reading

21. A. Ephrussi, and R. Lehmann (1992) Induction of germ cell formation by *oskar* [see comments]. *Nature* **358**, 387–392.
22. N. Gunkel, T. Yano, F. H. Markussen, L. C. Olsen, and A. Ephrussi (1998) Localization-dependent translation requires a functional interaction between the 5' and 3' ends of *oskar* mRNA. *Genes Dev.* **12**, 1652–1664.
23. J. Kim-Ha, P. J. Webster, J. L. Smith, and P. M. Macdonald (1993) Multiple RNA regulatory elements mediate distinct steps in localization of *oskar* mRNA. *Development* **119**, 169–178.
24. C. Rongo and R. Lehmann (1996) Regulated synthesis, transport and assembly of the *Drosophila* germ plasm. *Trends Genet.* **12**, 102–109.
25. P. J. Webster, L. Liang, C. A., Berg, P. Lasko, and P. M. Macdonald (1997) Translational repressor bruno plays multiple roles in development and is widely conserved. *Genes Dev.* **11**, 2510–2521.

### Ouchterlony Double Diffusion

A direct extension of the [precipitin reaction](#) to [immunoassays](#) is the Ouchterlony double diffusion technique (1). Here, two or more wells are cut in an agar slab and filled with either [antigen](#) or [antibody](#). The two materials slowly diffuse throughout the slab. As the antibody and antigen molecules encounter each other and react, a precipitate forms in the gel. Because of the highly cooperative nature of immune complex lattice formation, the precipitate forms an extremely sharp line between the two wells. The position of the precipitin line, compared to lines generated from other dilutions of antigen or antibody, gives a semiquantitative measure of the relative concentration of the reactant. However, the true value of the Ouchterlony technique is the remarkable number of details that can be inferred from the precipitin patterns formed. From the presence of multiple lines, or complex precipitin patterns formed at the intersection with multiple antigen or antibody wells in the agar slab, conclusions can be drawn about such properties as antigen heterogeneity, antigen molecular weight, shared antigenic determinants, and cross-reactivity of different antisera (2). The Ouchterlony technique is now obsolete, but it is commonly described in older literature.

#### Bibliography

1. Ö. Ouchterlony (1948) *Acta Pathol. Microbiol. Scand.* **25**, 186–191.
2. Ö. Ouchterlony and L.-Å. Nilsson (1986) "Immunodiffusion and immunoelectrophoresis", in *Handbook of Experimental Immunology*, 4th ed., Vol. **I: Immunochemistry**, D. M. Weir, ed., Blackwell, Oxford, pp. 32.1–32.50.

### Ovalbumin

The major [protein](#) in egg-white is ovalbumin. It has the typical highly ordered structure of a [serpin](#)

(see Fig. 1 in that entry), but has lost the ability to refold the loop that would correspond to its reactive center; consequently it has no [protease inhibitor](#) activity. Its smallness relative to other serpins, and its conformational stability, make it suitable as a marker for [gel electrophoresis](#). The stability is only relative, however; its reactive-center loop can be cleaved **proteolytically** to produce a form known as plakalbumin. It is likely that the induced insertion of this loop, like that that occurs in other serpins, explains a conformation, with increased stability, that forms on long-term storage: S-ovalbumin.

#### Suggestions for Further Reading

P. E. Stein, A. G. W. Leslie, J. T. Finch, and R. W. Carrell (1991) Crystal structure of uncleaved ovalbumin at 1.95 Å, *J. Mol. Biol.* **221**, 941–959.

P. G. W. Gettins, P. A. Patston, and S. T. Olson (1996) *Serpins: Structure, Function and Biology*, R. G. Landes, Georgetown, TX, pp. 135–137.

### Overlay Assay: Enzyme Zymography of Plasminogen Activators and Inhibitors

Zymography is a gel-based enzymatic assay useful for detection and semi-quantitative analysis of proteases and is performed using overlay or co-polymerization techniques (1, 2). In overlay assays, samples are fractionated by SDS-PAGE and, after removal of SDS, the polyacrylamide gel is overlaid on an agarose indicator gel containing protein substrate. Enzyme diffusion into the indicator gel occurs during incubation of the gel sandwich and results in degradation of protein substrate. The indicator gel is then stained to reveal zones of proteolysis, which appear as cleared bands on a dark background. In co-polymerization techniques, protein substrate is incorporated directly into the polyacrylamide gel matrix and, after electrophoresis and SDS removal, gels are incubated floating in buffer (2). In either method, specificity or sensitivity is tailored by modification of acrylamide gel pore size, substrate choice, incubation time, buffer conditions, and by inclusion of activators or inhibitors in gels or buffers.

Overlay assays are advantageous when substrate preference involves proteins that either bind poorly to acrylamide (albumin, casein) or are irreversibly denatured by SDS (collagen type I) and, thus, are not amenable to co-polymerization methods; gelatin, in contrast, is readily incorporated into the polyacrylamide gel matrix, and co-polymerization has been used extensively for assay of gelatinases (2). Overlay methods are also more appropriate for initial quantitative evaluations because enzyme activity is confined within the gel sandwich and, thus, not susceptible to diffusional loss to surrounding buffer, as may occur with co-polymerization assays (1). Furthermore, overlay assays are the method of choice for analysis of complex proteolytic systems in which zymogens, co-substrates, or endogenous inhibitors are relevant, as multiple components can be incorporated into indicator gels.

Representative of a complex proteolytic system is application of the overlay assay to analysis of the plasminogen activator (PA), inhibitor (PAI) system (3, 4). Briefly, samples containing PAs (urokinase [uPA], tissue plasminogen activator [tPA]) and their endogenous inhibitor (PAI) are resolved by SDS-PAGE on 10% gels. Gels are washed with Triton X-100, to remove SDS and allow protein renaturation, and then overlaid on an indicator gel containing plasminogen (zymogen) and fibrin (protein substrate). During incubation, PAs diffuse into the indicator gel, cleave plasminogen to form plasmin, a serine protease that degrades fibrin. After staining of the indicator gel, lytic zones (PA) appear as cleared bands on a dark background (zymography). In contrast, if exogenous PA is also incorporated into the indicator gel, then lysis occurs everywhere except where inhibitor is

present; thus, lytic-resistant zones (PAI) appear as dark bands on a clear background (reverse zymography).

## 1. Materials and Methods

### 1.1. Materials

Fibrinogen (plasminogen-rich) from Organon Teknika, Holland; SeaKem garose and Gelbond support backing from FMC BioProducts, Rockland, ME; human uPA and tPA from American Diagnostica, Greenwich, CT. All reagents prepared using ultra-pure water.

### 1.2. Sample Preparation and SDS-PAGE

Samples were prepared, as appropriate, in 2X or 5X Laemmli sample buffer without reducing agent or heating (4). Proteins were resolved by 10% SDS-PAGE (0.75 mm thick slab gels; 20 mA constant current) as previously described (5). After electrophoresis, gels were washed twice (2.5% Triton X-100, 200 mL/gel, 30 min each) to remove SDS and carefully overlaid on indicator gels.

### 1.3. Preparation of Fibrin Indicator Gel

Plasminogen-rich fibrinogen (5 mg/mL) was solubilized in 0.85% saline (prewarmed to 37°C) with gentle mixing (inversion several times over 1–2 h). Agarose (1%) was separately prepared in PBS (pH 7.4) by heating in a boiling water bath. Dissolved agarose was removed from heat and cooled to 60°C; 15 mL of thrombin stock (260 U/mL) and 5 mL fibrinogen solution were then added to the hot agarose (25 mL) in rapid succession and with constant swirling. The agarose-fibrin solution was quickly poured onto the hydrophilic side of a Gelbond support film (~13 × 17 cm) using a 25 mL glass pipet and a continuous, vertical motion from left to right; the Gelbond support film had been previously stabilized on a glass plate (see below). The agarose indicator gel was allowed to solidify for at least 30 min at room temperature before overlay with the acrylamide gel. Final gel composition was ~0.8% agarose, 0.08% plasminogen-rich fibrinogen, and 0.1 U/mL thrombin (4).

### 1.4. Fibrin Zymography

Washed polyacrylamide gels were carefully overlaid on fibrin indicator gels and incubated (37°C, 18–20 h) in closed containers containing water-moistened paper towels to provide a humidified atmosphere. After incubation, the acrylamide gel was carefully removed, and the indicator gel was stained (<5 min; 0.1% amido black, 70% methanol, 10% acetic acid), destained (70% methanol, 10% acetic acid), and allowed to air dry.

### 1.5. Reverse Zymography

For detection of PAIs, the indicator gel also contained uPA (0.8 U/mL).

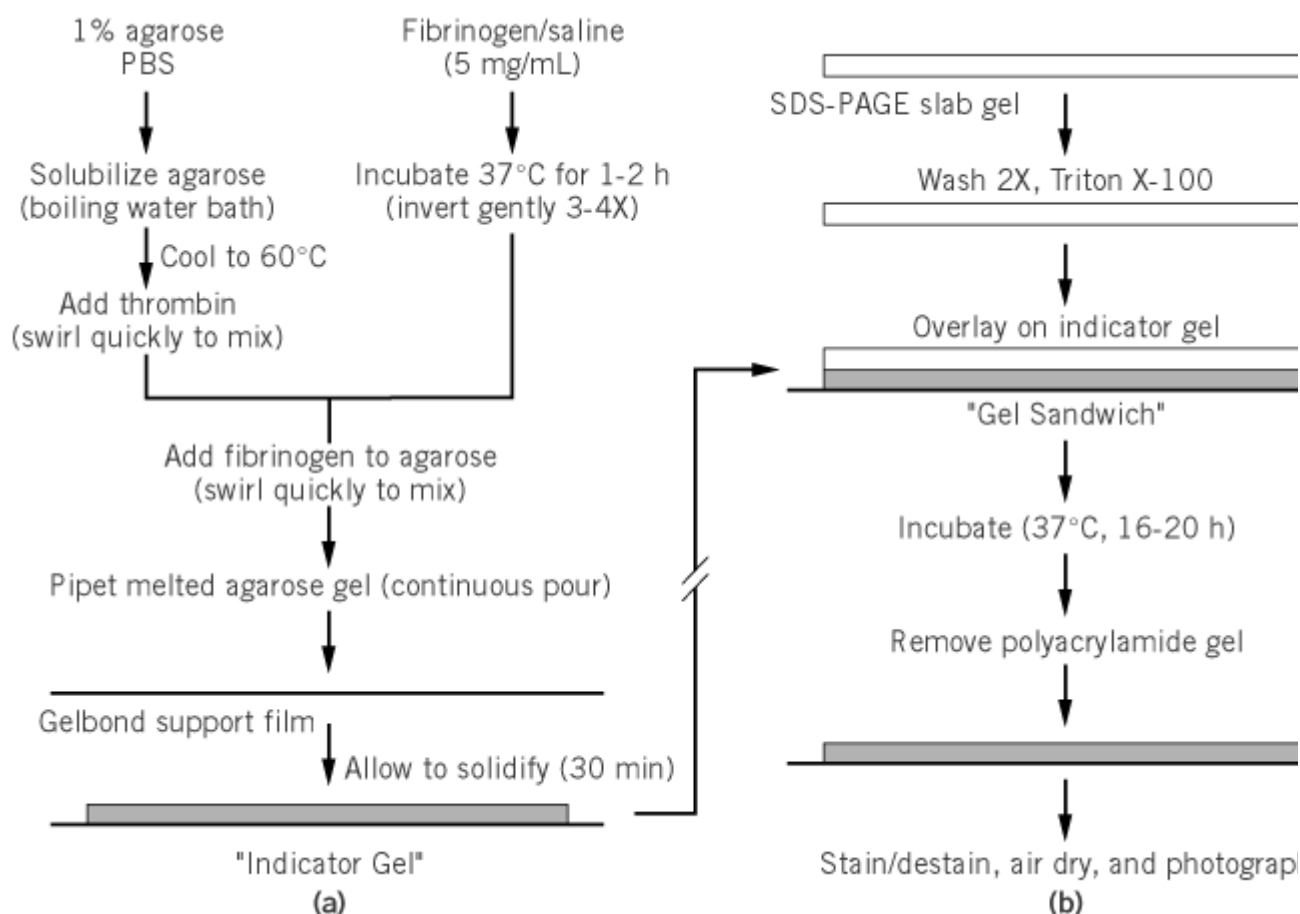
### 1.6. Notes and Tips

- 1) Excessive or too vigorous mixing during fibrinogen solubilization results in protein flocculation;
- 2) volume loss during agarose heating should be reconstituted using ultra-pure water;
- 3) Gelbond has hydrophilic (spreads water) and hydrophobic (beads water) surfaces; water placed in the center of a glass plate adheres the hydrophobic surface of the Gelbond providing a flat stable surface for gel pouring and subsequent handling of gel sandwiches;
- 4) Work quickly, fibrin polymerization and agarose solidification occur rapidly; delay in pouring produces inhomogeneities within the indicator gel;
- 5) when disassembling glass plates after electrophoresis, cut the bottom of the acrylamide gel for later orientation;
- 6) acrylamide gels enlarge and become fragile during washing; handle with care;
- 7) to position acrylamide gels carefully grasp the upper corners; overlay from bottom to top in a continuous motion to avoid air trapping;
- 8) avoid repositioning of the polyacrylamide gel when constructing the gel sandwich as artifactual lysis and streaking may result;
- 9) volumes given are for one full-sized indicator gel;
- 10) up to two indicator gels can be simultaneously prepared; extra gels can be stored at 4°C for 3–4 days in sealed plastic bags humidified with a damp paper towel.

## 2. Results and Discussion

Samples containing PAs can be readily analyzed by fibrin zymography using a gel overlay assay (Fig. 1). In such assays PAs and PAIs are resolved by electrophoresis on 10% polyacrylamide gels in the absence of reducing agent or heating. Omission of reducing agent is necessary because SDS may cause artifactual activation of thiol-dependent proteases (i.e., lysosomal thiol cathepsins) as well as inactivation of disulfide-stabilized proteases (i.e., plasminogen activators). Similarly, heating is avoided to preclude artifactual aggregation, activation, or inactivation of susceptible proteases.

**Figure 1.** Schematic of fibrin overlay assay. (a) Preparation of the fibrin indicator gel. (b) Overlay of the indicator gel (fibrin zymography). For details see MATERIALS AND METHODS section.



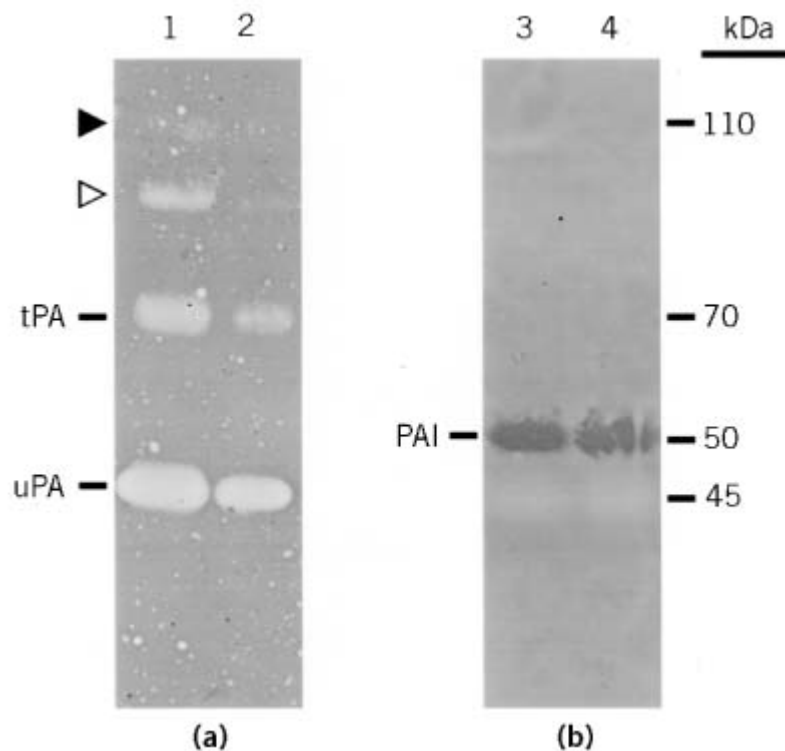
During electrophoresis the partially denatured protease migrates with respect to molecular weight. SDS removal after electrophoresis by the nonionic detergent Triton X-100 allows proteins to renature and regain activity. Wash times should be kept to a minimum and accurately monitored to avoid excessive leaching of protease and to assure reproducibility between experiments, respectively.

The gel sandwich should be constructed without repositioning of the acrylamide gel and should be incubated in a humidified atmosphere to avoid gel drying. Incubation times and conditions (temperature, pH) vary depending on the intrinsic properties of the protease being examined as well as its concentration in biological samples. For quantitative purposes, standard curves should be constructed to determine dose-time relationships with respect to size of the lytic zone. Overlay assay of PAs requires incubation for 6–18 h and is temperature-dependent (37°C).

Indicator gels require only brief staining (<5 min) and, after drying, are readily photographed by

back lighting. Overlay zymography of PAs demonstrates fibrinolytic zones at molecular weights corresponding to uPA (45 kDa) and tPA (70 kDa) as well as some high (110 kDa) molecular weight bands (Fig. 2A), presumably representing complexes of PA-PAI (3,4). In contrast, reverse zymography demonstrates fibrinolytic-resistant zones at a molecular weight (50 kDa) corresponding to PAI (Fig. 2B). Protease/inhibitor identification can be verified by parallel zymograms in which indicator gels contain PA-specific antibodies or inhibitors i.e., amiloride to inhibit uPA and erythrina to inhibit tPA, as has been described elsewhere (4).

**Figure 2.** Amido black-stained indicator gels. (a) Fibrin zymogram demonstrating uPA (45 kDa), tPA (70 kDa), and higher molecular weight PA complexes (**open and closed triangles**). (b) Reverse zymogram demonstrating PAI (50 kDa) as well as uPA shadow (45 kDa). Samples 1–4 were from conditioned media of bovine corneal endothelial cells at different growth stages (4).



In summary, overlay zymography is a versatile assay for the detection and semi-quantitative investigation of proteases, which is readily applicable to complex systems. Unique to the overlay method is the opportunity to incorporate into the indicator gel multiple components including activators, inhibitors, antibodies, zymogens, co-factors, and allosteric modifiers. This enables highly sensitive and specific assays to be developed and, thus, affords extensive biochemical investigation of diverse proteolytic systems.

#### Bibliography

1. M. S. Lantz and P. Cibrowski (1994) *Meth. Enzymol.* **235**, 563–594.
2. G. S. Makowski and M. L. Ramsby (1996) *Anal. Biochem.* **236**, 353–356.
3. D. J. Loskutoff, J. A. van Mourik, L. A. Erickson, and D. Lawrence (1983) *Proc. Natl. Acad. Sci. U.S.A.* **80**, 2956–2960.
4. M. L. Ramsby and D. L. Kreutzer (1993) *Invest. Ophthalmol. Vis. Sci.* **34**, 3207–3219.
5. G. S. Makowski and M. L. Ramsby (1997) In *Protein Structure: A Practical Approach*, (T.E. Creighton, ed.), IRL Press, Oxford, UK, pp. 1–27.

## Suggestions for Further Reading

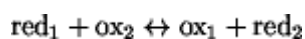
6. P. A. Andreasen, B. Georg, L. R. Lund, A. Riccio, and S. Stacey (1990) *Mol. Cell. Endocrinol.* **68**, 1–19.
7. K. Dano, P. A. Anderson, J. Grondahl-Hansen, P. Kristensen, L. S. Nielsen, and L. Skriver (1985) *Adv. Cancer Res.* **44**, 139–266.
8. L. A. Erickson, D. A. Lawrence, and D. J. Loskutoff (1984) *Anal. Biochem.* **137**, 454–463.
9. J. D. Vassalli, A. P. Sappino, and D. Belin (1991) *J. Clin. Invest.* **88**, 1067–1072.
10. M. L. Ramsby, P. C. Donshik, and G. S. Makowski (2000) *Inflammation* **24**, 45–71.

## Oxidation/Reduction Potential

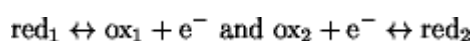
Many of the fundamental physiologic reactions in a cell, both in the metabolic pathways of degradation and in the storage of energy via production of ATP, involve the transfer of electrons from one molecule, the donor (or reducing agent), to another, the acceptor (or oxidant). They are thus called redox reactions and are catalyzed by [redox enzymes](#).

### 1. Redox Reactions and Redox Couples

The molecule donating an electron is the reductant (red) and becomes oxidized in the reaction, whereas the molecule accepting the electron is the oxidant (ox) and becomes reduced. In fact, each redox exchange:



is composed of two reactions:



The reductant and its corresponding oxidized form constitute a redox couple, in which both partners participate in the redox reaction. Each redox couple is characterized by a standard oxidation-reduction potential ( $E'_o$ ) that quantitatively defines its tendency to lose an electron. In a redox reaction, the redox couple that has the higher affinity for the electron will be the oxidant, whereas the couple having a greater tendency to donate electrons will be the reductant. The same redox couple may therefore be either a reductant or an oxidant, depending on the redox potential of the second redox couple in the reaction. Examples of such a mechanism occur in the **mitochondrial** and the anaerobic bacterial electron transfer chains and in the [photosynthesis](#) pathway in **chloroplasts**, where the reductant of the (i)th step is the oxidant of the (i-1)th step.

### 2. Electrode Potential

The redox potential can be measured electrochemically as the electromotive force generated by connecting a half-cell of the redox sample with a reference half-cell. The reference half-cell consists of a platinum electrode immersed in a 1 M H<sup>+</sup> solution and saturated with H<sub>2</sub> gas at 1 atmosphere; the sample half-cell consists of an electrode immersed in a solution of the redox couple under standard conditions (1 M of the electron donor and acceptor, 25°C, and pH 7). The electrodes are

connected to a voltmeter. A salt bridge (KCl solution) enables electrical continuity between the half-cells. The electrons flow in the direction determined by the redox potential of the sample relative to the reference. This is the voltage observed at the beginning of the experiment (standard concentration of 1 M), as the redox potential of the hydrogen electrode is arbitrarily taken as 0 V. The redox potentials ( $E_o'$ ) of redox couples are thus expressed relative to the  $H_2/2H^+$  couple.

For a redox reaction ( $ox + e^- \leftrightarrow red$ ), the Nernst equation relates the concentration of the redox species with the redox potential:

$$E = E^\circ + (RT/nF \ln[ox]/[red])$$

where  $E$  is the redox potential,  $E^\circ$  is the redox potential for components in their standard state at pH 0,  $R$  is the gas constant ( $8.3144 \text{ J K}^{-1} \text{ mol}^{-1}$ ),  $F$  is the Faraday constant ( $9.65 \times 10^4 \text{ J V}^{-1} \text{ mol}^{-1}$ ),  $n$  is the number of electrons transferred, and  $T$  is the absolute temperature. In biology, the equation is generally written as

$$E_h = E_{m,x} + (0.06/n \log[ox]/[red])$$

where  $E_h$  is the redox potential referred to the standard hydrogen electrode,  $E_{m,x}$  is the mid-point redox potential (when  $[ox] = [red] \neq \text{standard state}$ ) at a defined pH (of  $x$ ) and  $2.303 RT/F = 0.06$  for  $29.4^\circ\text{C}$ . The value of  $E_m$  for a defined redox couple depends on the relative stability of the oxidized and reduced states: the more negative the value of  $E_m$ , the more stable is the oxidized form and the stronger is the electron donor. Conversely, any factor stabilizing the reduced form makes the couple a better electron acceptor, having a more positive redox potential.

### 3. Redox Potential, Free Energy, and Equilibrium Constant

The change in **free energy** associated with a redox reaction ( $red_1 + ox_2 \leftrightarrow ox_1 + red_2$ ) is related to the difference in redox potentials of the reactants by the formula:

$$\Delta G^{o'} = -nF\Delta E^{o'}$$

in which  $n$  is the number of electrons transferred,  $F$  is the energy change as 1 M of electrons falls through a potential of 1 V ( $23.06 \text{ kcal V}^{-1} \text{ mol}^{-1}$ ), and  $\Delta E^{o'}$  is the difference in redox potentials of the reactants in volts.  $\Delta G^{o'}$  is the free energy change per mole and is expressed in kilocalories.

We know that at constant temperature and pressure a reversible reaction proceeds until an equilibrium is attained, defined by

$$[ox_1]/[red_1] \times [red_2]/[ox_2] = Keq$$

in which  $[ox_1]$ ,  $[red_1]$ ,  $[red_2]$ , and  $[ox_2]$  are the reactant concentrations, and  $Keq$  is the equilibrium constant of the reaction at a certain constant temperature.  $Keq$  is related to the  $\Delta G^{o'}$  free energy change by;

$$\Delta G^{o'} = -RT \ln Keq$$

where  $R$  is the gas constant ( $1.987 \text{ cal mol}^{-1} \text{ K}^{-1}$ ) and  $T$  is the absolute temperature ( $K$ ).

We can also say that:

$$-nF\Delta E^{of} = -RT \ln K_{eq} \text{ or } \Delta E^{of} = RT/nF \times \ln K_{eq}$$

Knowing the redox potentials of a redox couple, the last relationships make it possible to calculate the equilibrium concentrations. Conversely, the redox potential of a solution can be calculated from the relative concentrations of the redox couples at equilibrium.

#### 4. How to measure the redox potential

The oxidation of a reductant ( $red_1$ ) by an oxidant ( $ox_2$ ), at fixed temperature and pH, can be followed potentiometrically by measuring the variation of the electromotive force of a solution of the reductant as it is titrated with the oxidant. In practice, with a Pt electrode connected with a reference electrode (ie, a calomel electrode), the redox potential of the  $red_1$  solution is recorded after each addition of  $ox_2$ . The redox potential increases as the value of  $\log [ox_1]/[red_1]$  increases and is equal to  $E_m$  when the initial reductant is half oxidized,  $[ox_1] = [red_1] = 50\%$ . When  $[ox_1]$  is close to 100%, the redox potential varies rapidly, and the titration comes to an end; the point of equivalence has then been reached, i.e. an equivalent of  $ox_2$  has been added to  $red_1$ . In biological systems, the redox proteins generally carry colored groups that have substantial absorbance in the visible range of wavelengths. This allows precise determination by [absorption spectroscopy](#) of the relative concentrations of the oxidized and reduced species.

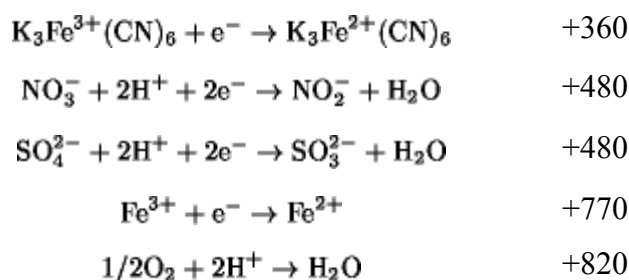
Many molecules can act as oxidizing and reducing agents in nonbiological redox reactions, whereas in the metabolic pathways of living organisms, a few molecules have been conserved during the [evolution](#) as redox agents for many different substrates.

The oxidation/reduction potentials of some biologically relevant compounds and free redox groups are presented in Table 1.

**Table 1. Oxidation/Reduction Potentials of Some Biologically Relevant Compounds and Free Redox Groups**

| Redox Couple   | $E_{m,7}(\text{mV})$ |
|--|----------------------|
| $\text{NAD(P)}^+ + \text{H}^+ + 2\text{e}^- \rightarrow \text{NAD(P)H}$            | -320                 |
| $\text{S} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{H}_2\text{S}$              | -230                 |
| $\text{Riboflavin} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Riboflavin-H}_2$  | -200                 |
| $\text{FMN} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{FMNH}_2$                 | -219                 |
| $\text{Pyruvate} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{lactate}$           | -190                 |
| $\text{Protoheme IX(Fe}^{3+}) + \text{e}^- \rightarrow (\text{Fe}^{2+})$           | -115                 |
| $2\text{H}^+ + 2\text{e}^- \rightarrow \text{H}_2$                                 | 0.0 <sup>a</sup>     |
| $\text{Dehydroascorbate} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Ascorbate}$ | +60                  |
| $\text{Ubiquinone} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Ubiquinone-H}_2$  | +100                 |
| $\text{MetBleu} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{MetBleu}$            | +110                 |
| $\text{Cu}^{2+} + 2\text{e}^- \rightarrow \text{Cu}^+$                             | +150                 |





<sup>a</sup> Taken as reference.

### Suggestions for Further Reading

J. G. Morris (1968) *A Biologist's Physical Chemistry*, Edward Arnold Ltd., London, "Chapter 12".  
 D. G. Nicholls and S. J. Ferguson (1992) *Bioenergetics 2*, Academic Press, "Chapter 3".

## Oxidoreductases

Oxidoreductases form a large class of [enzymes](#) that catalyze the oxidation of one substrate, with the concomitant reduction of another. For these donor-acceptor reactions, the electron acceptor may be a pyridine nucleotide (eg, **NAD**), [cytochrome](#), oxygen, quinone or others. The term [dehydrogenase](#) is used to describe oxidation-reduction reactions whenever possible, whereas the term *oxidase* is restricted to enzymes for which oxygen is an acceptor. Dehydrogenases act on primary and secondary alcohols, as well as hemiacetals. They also convert aldehydes to the corresponding acids, and these reactions may involve phosphorylation of the acid or acetylation of coenzyme A (CoA). Oxo groups may be oxidized by the addition of [water](#), followed by the cleavage of a **C–C**bond or dehydrogenation. Dehydrogenases can introduce double bonds by direct dehydrogenation at a single **C–C**bond. They function as oxidases by utilizing oxygen for the deamination of [amino acids](#) that can also undergo direct dehydrogenation reactions. In addition, they can dehydrogenate secondary amines to form **C=N**double bonds and oxidize NADH or NADPH as well as nitrogenous substrates and donors with sulfur or heme groups.

## Oxygen-Binding Proteins

Oxygen-binding proteins are defined as [proteins](#) that reversibly bind  $\text{O}_2$ . They are exclusively [metalloproteins](#).

### 1. Classification

The classification of oxygen-binding proteins and some of their properties and distribution are

presented in Table 1. Protoheme IX is the common oxygen-binding site for **myoglobins** (Mb), IDO-like Mb, **hemoglobins** (Hb), compact Hb, and FixL. Dioxygen is bound at the ferrous ion at the center of the heme group (see [Myoglobin](#)). The vinyl group normally at position 2 of protoheme IX is replaced by a formyl group in the chlorocruoroheme carried by chlorocruorin (Chl). Hemerythrin (Hr), myohemerythrin (myoHr), and hemocyanin (Hc) are nonheme proteins. The oxygen-binding site for Hr and myoHr is a binuclear iron center, and that for Hc is a binuclear copper center.

**Table 1. Classification of Oxygen-Binding Proteins**

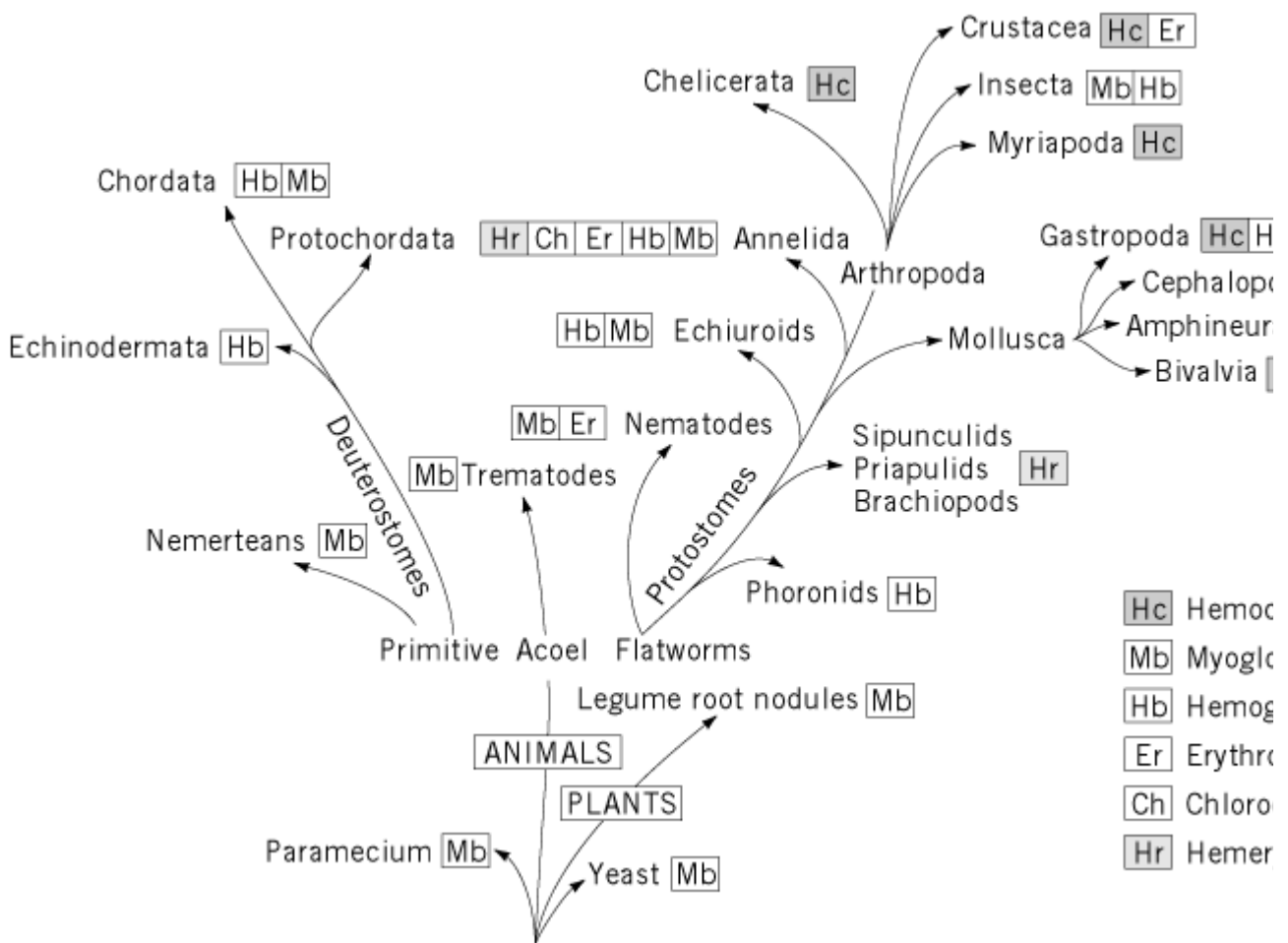
| Name               | O <sub>2</sub> Binding Site | Protein Folding                       | Molecular Mass, kDa | Aggregation State       | Color               |              | Distribution                                |
|--------------------|-----------------------------|---------------------------------------|---------------------|-------------------------|---------------------|--------------|---|
|                    |                             |                                       |                     |                         | Deoxy               | Oxy          |   |
| Myoglobin          | Protoheme IX                | Globin                                | 17 or 34            | Monomer or dimer        | Purple              | Red          | Vertebrate muscle, gastropod radularia      |
| IDO-like myoglobin | Protoheme IX                | Unknown                               | 40 or 80            | Monomer or dimer        | Purple              | Red          | Gastropod buccal radularia                  |
| Hemoglobin         | Protoheme IX                | Globin                                | 17, 34, or 44       | Monomer or dimer        | Purple              | Red          | Bacteria, unicellular eukaryotes and plants |
|                    | Protoheme IX                | Globin                                | 65                  | Tetramer                | Purple              | Red          | Vertebrate cell                             |
|                    | Protoheme IX                | Globin                                | 3500                | 144mer                  | Purple              | Red          | Annelid plasma                              |
| Compact hemoglobin | Protoheme IX                | Unknown                               | 13–15               | Monomer                 | Purple              | Red          | Cyanobacteria, protozoa, green algae        |
| FixL               | Protoheme IX                | Unknown                               | 110                 | Dimer                   | Similar to Hb or Mb |              | <i>Rhizobium</i> , <i>Bradyrhizobium</i>    |
| Chlorocruorin      | Chlorocruoroheme            | Globin                                | 3500                | 144mer                  | Greenish red        | Greenish red | Blood pigments of four polychaete families  |
| Myohemerythrin     | Binuclear Fe center         | 4- $\alpha$ -Helix bundle             | 14                  | Monomer                 | Colorless           | Violet       | Sipuncular muscle                           |
| Hemerythrin        | Binuclear Fe center         | 4- $\alpha$ -Helix bundle             | 40–110              | 3, 4 or 8mer            | Colorless           | Violet       | Red cell phyla of invertebrates             |
| Hemocyanin         | Binuclear Cu center         | $\alpha$ -Helices and $\beta$ -barrel | 450–3300            | 6, 12, 24, 36, or 48mer | Colorless           | Blue         | Arthropod plasma                            |
|                    | Binuclear Cu center         | Unknown                               | 3500–8700           | 10 or 20mer             | Colorless           | Blue         | Mollusk plasma                              |

The hemoglobins in protozoa (*Paramecium* and *Tetrahymena*), green **algae** (*Chlamydomonas*), and **cyanobacteria** are not authentic hemoglobins because their polypeptide chains are about 20% shorter and without substantial amino acid sequence similarities. Several pieces of evidence suggest that the structures of these “compact hemoglobins” differ from the “globin fold” (see [Globins](#)), and it is generally considered that they arose from a different **evolutionary** ancestor (see [Convergent Evolution](#)). Likewise, IDO-like Mb evolved from a common ancestor of indoleamine dioxygenase, a heme-containing [enzyme](#) that degrades tryptophan. Based on their gene structures and amino acid sequences, myoglobins and hemoglobins are divided into three distinct evolutionary groups: (1) “universal globins”, 17-kDa polypeptide chains, that have about 150 amino acid residues; (2) “truncated globins” or “compact globins”, 13 kDa, that have about 120 residues per chain; and (3) “IDO-like Mb” ([1](#)). The second group has been known as “truncated globin”, but “compact globin” is more appropriate because these proteins have no detectable amino acid [homology](#) with the universal globins.

## 2. Distribution, Structure, and Function

Figure [1](#) illustrates the distribution of oxygen-binding proteins in a simplified **phylogenetic** tree. Hr and Hc are distributed within certain animal phyla, such as Sipuncula, Brachiopoda, Annelida, Arthropoda, and Mollusca. Hb, the most widely distributed, occurs in **bacteria** to vertebrates, and even in **plants**, but only intermittently. The principal physiological function of the circulating proteins (vertebrate Hb, annelid Hb, Chl, Hr, and Hc) is to transport O<sub>2</sub> from the lungs or gills to peripheral tissues. The function of the noncirculating proteins (Mb and myoHr), and probably also of IDO-like Mb, is oxygen storage and intracellular transport.

**Figure 1.** Phylogenetic distribution of oxygen-binding proteins. Abbreviations undefined in the text are Er, erythrocrurin same as Chl (chlorocruorin). Erythrocrurin was the name used for what is now called annelid extracellular hemoglobin. Hemoproteins of yeast, paramecium, and legume root nodules are designated as Mb in this figure, they are sometimes designated as Hb, because they are not contained in muscle and Hb is a more comprehensive name of an oxygen carrier. From Van Ho Miller (1995) (see Reading List).



See [Myoglobin](#) for details of Mb and IDO-like Mb and [Hemoglobin](#) for details of Hb. The other oxygen-binding proteins are described here.

### 2.1. Compact Hemoglobin

Only a limited number of structure and functional studies of compact Hb have been carried out. The physiological role of these proteins is uncertain. It is proposed that Hb from the cyanobacterium *Nostoc commune* is an oxygen scavenger in [nitrogen fixation](#).

### 2.2. FixL

The FixL that occurs in *Rhizobium meliloti* and *Bradyrhizobium japonicum* is a chimeric protein consisting of a heme-containing **domain** and a protein [kinase](#) domain. Its proposed function is to sense oxygen through its heme domain and to transduce this signal by controlling the phosphorylation of transcriptional activator FixJ, which in turn induces the expression of the genes for nitrogen fixation, thereby providing oxygen-dependent control of nitrogen fixation (2). The heme domain has no significant sequence similarities to the globins, suggesting that it arose independently.

### 2.3. Chlorocruorin

Chl has only a limited phylogenetic distribution. It is found in the blood plasma of four families of polychaete annelids, the *Sabellidae*, *Serpulidae*, *Chlorhaemidae*, and *Ampharetidae*. The protein structure of Chl is essentially the same as that of extracellular annelid Hb, except for the change in the heme group (3). The vinyl to formyl replacement at position 2 of protoheme produces the characteristic greenish color and low oxygen affinity of Chl due to the enhanced electron-withdrawing properties of the formyl group. Chl shows highly cooperative oxygen binding. As with annelid Hb, its oxygen affinity is pH-dependent (see [Bohr Effect](#)) and is increased by divalent

cations like  $Mg^{2+}$  and  $Ca^{2+}$  (4).

#### 2.4. Hemerythrin (Hr) and Myohemerythrin

Hr has been found only in four phyla of marine invertebrates: Priapulida, Sipuncula, Brachiopoda, and one family of annelida, Magelonidae. It is contained in red cells called “hemerythrocytes” that float in the coelomic fluid. A noncirculating Hr is found in muscle fibers of sipunculids and is called “myoHr”, in analogy to myoglobin. Coelomic Hr is found predominantly as the octamer, but also as tetramers, trimers, and dimers. The octamer is composed of two different subunits that have different amino acid sequences (5). The octamer has a mass of 108 kDa, the shape of a square doughnut, and consists of two rings of four subunits each whose rings are associated in a face-to-face, **isologous** manner (see [Quaternary Structure](#)). Each 13.5-kDa subunit consists of four parallel **alpha-helices** twisted in a left-handed bundle (6) (see [Four-Helix Bundle Motif](#)). The [tertiary structure](#) of myoHr (13.9 kDa) is similar to that of one Hr subunit. MyoHr and each subunit of Hr contains a site that has two iron atoms, which bind one  $O_2$  molecule. The iron atoms are bound directly to seven amino-acid side chains of the protein. One of the  $O_2$  atoms binds to one of the iron atoms, and the other oxygen atom forms a [hydrogen bond](#) with the oxo bridge that connects the two iron atoms.

The oxygen affinity of myoHr is higher than that of coelomic Hr of the same species. Octameric Hr from two brachiopods (*Lingula unguis* and *Lingula reevii*) show notable cooperativity in oxygen binding and have **Hill coefficients** of 1.8 to 2.2 (see [Hemoglobin](#)), whereas other Hrs are essentially noncooperative, irrespective of their state of aggregation. It has been suggested that the allosteric unit of *L. unguis* octameric Hr is the entire molecule (7). Effects of intracellular cofactors, such as  $H^+$ ,  $CO_2$ ,  $Mg^{2+}$ ,  $Ca^{2+}$ , and  $Cl^-$ , on the oxygen affinity have been observed with some Hrs, but their physiological significance is not clear.

#### 2.5. Hemocyanin

Hc is found in the blood plasma of two phyla, arthropod and mollusk. The oxygen-binding site of Hc is composed of a pair of copper ions, each of which is liganded by three [histidine](#) residues of the protein. It was long thought that  $O_2$  binds as a peroxo ion, serially bridging the two Cu(II) ions, but now there is direct evidence that the mode of  $O_2$  binding is a  $m-h^2:h^2$  coordination, in which both oxygen atoms are bound to both of the copper ions and the O–O axis perpendicular to the Cu–Cu axis.

Arthropod and mollusk Hc exhibit similar optical and ligand-binding properties, but their molecular architectures are quite different. All arthropod Hcs are composed of hexamers. Each 75-kDa polypeptide chain carries one oxygen-binding site. These aggregate into 1, 2, 4, 6, or 8 hexamers, depending on the species, and no intermediate aggregations states are observed in each case. With this mode of assembly, one Hc molecule has from 6 to 48 oxygen-binding sites and a molecular mass ranging from 450 to 3300 kDa. Each hexamer is a “trimer of dimers”, and the shape is a trigonal antiprism. Each polypeptide chain consists of three domains. The first and second are  $\alpha$ -helical, whereas the third has an irregular structure that contains a 7-strand  $\beta$ -barrel (see [Beta-Sheet](#)). The  $O_2$ -binding site is in the second domain.

In contrast to the arthropod Hc, those from mollusks are organized entirely differently. The molecule is composed of one or more cylindrical units, each containing 10 subunits. Each subunit is a long polypeptide chain of 350 to 440 kDa composed of seven or eight globular folded regions, each of which binds one  $O_2$  molecule. It is not known whether each such region is composed of more than one domain, so it is called instead a “functional unit.” Its three-dimensional structure is not yet known.

The quaternary structures of both arthropod and mollusk Hc are stabilized by many calcium ions. The Hc oxygen affinity varies, depending on the species, especially in the Arthropoda, and reflects

adaptations to the great variety of life styles and environments of the animals. Most Hcs demonstrate a distinct [Bohr effect](#) and varying degrees of cooperativity in oxygen binding. The Hill coefficient ranges between 1 and 12. A higher degree of aggregation is not necessarily accompanied by greater cooperativity. The oxygen dissociation curves of Hcs that have large multimeric structures are not compatible with the classic allosteric **concerted model** of Monod, Wyman, and Changeux but are described best by a “nesting” model (8). This assumes that the allosteric units capable of transitions between two states,  $r$  and  $t$ , are in turn nested within larger structural units that are in equilibrium between two global states,  $R$  and  $T$ . In addition to the structural calcium ions, other calcium ions, plus magnesium ions, regulate the oxygen affinity.

There is very little sequence similarity between arthropod subunits and molluscan functional units. This and their very different molecular architectures indicate that they probably evolved independently.

### Bibliography

1. K. Fukami-Kobayashi, M. Mizutani, and M. Go (1995) In *Tracing Biological Evolution in Protein and Gene Structures* (M. Go and P. Schimmel, eds.), Elsevier, Amsterdam, pp. 271–282.
2. M. A. Gilles-Gonzalez, G. S. Ditta, and D. R. Helinski (1991) *Nature* **350**, 170–172.
3. A. N. Qabar, M. S. Stern, D. A. Walz, J.-T. Chiu, R. Timkovich, J. S. Wall, O. H. Kapp, and S. N. Vinogradov (1991) *J. Mol. Biol.* **222**, 1109–1129.
4. K. Imai and S. Yoshikawa (1985) *Eur. J. Biochem.* **147**, 453–463.
5. H. Yano, K. Satake, Y. Ueno, K. Kondo, and A. Tsugita (1991) *J. Biochem. (Tokyo)* **110**, 376–380.
6. K. B. Ward, W. A. Hendrickson, and G. L. Klippenstein (1975) *Nature* **257**, 818–821.
7. K. Imai, H. Takizawa, T. Handa, and H. Kihara (1991) In *Structure and Function of Invertebrate Oxygen Carriers* (S. N. Vinogradov and O. H. Kapp, eds.), Springer-Verlag, New York, pp. 179–189.
8. C. A. Robert, H. Decker, B. Richey, S. J. Gill, and J. Wyman (1987) *Proc. Natl. Acad. Sci. USA* **84**, 1891–1895.

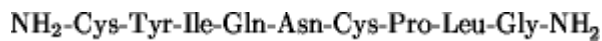
### Suggestions for Further Reading

9. M. C. M. Chung and H. D. Ellerton (1979) The physicochemical and functional properties of extracellular respiratory hemoglobins and chlorocruorins, *Prog. Biophys. Mol. Biol.* **35**, 53–102.
10. K. E. Van Holde and K. I. Miller (1995) Hemocyanins, *Adv. Protein Chem.* **47**, 1–81.
11. T. Suzuki and K. Imai (1998) Evolution of myoglobin, *Cell. Mol. Life Sci.* **54**, 979–1004.
12. D. M. Kurtz (1992) Molecular structure/function relationships of hemerythrins, *Adv. Comp. Environ. Physiol.* **13**, 151–171.
13. K. E. van Holde and K. I. Miller (1995) Hemocyanins, *Adv. Protein Chem.* **47**, 1–81.

### Oxytocin

Oxytocin (MW 1007) is a [peptide hormone](#) released by the posterior pituitary gland that displays a wide spectrum of biological activities, including a role in reproduction and social behavior (1). It stimulates uterine contraction and lactation, activates myometrial [GTPase](#) and [phospholipase C](#), and

increases sodium excretion. It consists of only nine amino acid residues, with Cys1 and Cys6 linked by a disulfide bond and with an amidated C-terminus:



The [complementary DNA](#) (cDNA) encoding the precursor for human oxytocin has been **cloned and DNA sequenced** (2). This indicated that the **genes** for oxytocin and vasopressin are linked on the same human chromosome, 20, and arose through [gene duplication](#) (3). The gene for the oxytocin precursor, pro-oxyphysin, is composed of two **introns** and three exons. The initial gene product contains oxytocin, plus a carrier protein, neurophysin.

The majority of oxytocin is produced in the neurohypophysis and magnocellular neurons of the paraventricular and supraoptic nuclei of the hypothalamus. It is processed **proteolytically** from its precursor form during transport along the axonal projections to the posterior pituitary. The enzymes and the mechanisms of the specific **posttranslational modifications** involved in the production of active peptides from the precursor protein are still not completely known. The process by which the prohormones are packaged into the membrane-bound **secretory granules** and transported down the axon supraopticohypophyseal tract is the subject of ongoing investigation. The hormone and its associated neurophysin are secreted from nerve terminals of the neurohypophysis into the portal system of the neurohypophysis by calcium-dependent [exocytosis](#), although the molecules responsible are not yet known. Oxytocin secretion is initiated by depolarizing electric impulses on the cell membranes. That is followed by calcium influx into the cell, fusion of the secretory granules with the cell membrane, and exocytotic extrusion of the hormone.

The relationship of oxytocin to its four types of receptors (4) and subsequent signaling is not completely understood (5). The human receptor has been sequenced (6), however, its structure analyzed, and its binding sites identified (7).

Oxytocin is primarily associated with lactation and parturition in females. It is used clinically to induce labor or augment labor, when appropriate, because it stimulates the contractions of uterine smooth muscle. Oxytocin release is stimulated by a baby's cry or suckling, and this subsequently releases prolactin, which causes milk to be released, by stimulating the contraction of myoepithelial cells in milk ducts of mammary glands and thereby forcing milk from alveoli in the breast. Oxytocin has been implicated in other functions. These include mating and maternal behavior, natriuresis, modulation of the cholesterol/phospholipid ratio in membranes, and even memory. However, the elimination of oxytocin in **transgenic** mice did not change their reproductive and maternal behavior (8, 9). Both male and female homozygous knockouts are fertile, and the females are capable of parturition; cross reactivity of vasopressin with the oxytocin receptors explains this apparent contradiction. Milk release was absent, but it was correctable by oxytocin injection. In males, the importance of oxytocin is even less well understood. It has been suggested that oxytocin is involved in prolactin release, penile erection, and various aggressive sexual behaviors.

## Bibliography

1. R. Ivell and J. A. Russell, editors (1995) *Oxytocin: Cellular and Molecular Approaches in Medicine and Research* (Advances in Experimental Medicine and Biology, vol. 235), Plenum, New York.
2. E. Mohr (1995) *Vitamins and Hormones* **51**, 235–266.
3. Y. Hara et al. (1990) *Mol. Brain Res.* **8**, 319–324.
4. R. Ivell et al. (1997) *Biochem. Soc. Trans.* **25**, 1058–1066.
5. C. Barberis et al. (1998) *J. Neuroendocrinology* **156**, 223–229.
6. T. Kimura et al. (1992) *Nature* **356**, 526–529.
7. B. Mouillac et al. (1995) in ref. 1, pp. 301–307.
8. K. Nishimori et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 11699–11704.

## P Element

The P element is a [transposable element](#) in *Drosophila* that can insert into many different positions in the [genome](#) (1). In addition to being an element whose [transposition](#) mechanism has been interesting to dissect, its transposition is controlled in interesting ways at several levels. Manipulation of P-element transposition has also been an immensely powerful tool for *Drosophila* genetics (2). One widespread use of P elements is to make physically tagged [mutations](#) that have allowed the ready isolation of the affected **genes** through the physical marker provided by the P element. Also, P-element transposition has provided a widely used method for introducing sequences into the *Drosophila* genome.

P elements were discovered by studying the phenomenon [hybrid dysgenesis](#) that results when males from certain strains of *Drosophila* are crossed with females of certain different strains. In the germline of the progeny of such crosses, there is a high frequency of mutations and also sterility. Characterization of mutations resulting from such crosses revealed that they resulted from the insertion of a novel 2.9-kbp piece of DNA, the P element (see Fig. 1 in [Transposable Elements](#)). Transposition of this element is usually restricted to the germline, because the transposase is assembled from four exons by [RNA splicing](#) in the P-element [messenger RNA](#), and factors in somatic cells can both inhibit [transcription](#) and block splicing. Hybrid dysgenesis results from a cross of a P-element male with a non-P-element-containing female. The use of the maternal non-P-element strain is important because, in addition to making transposase, the P element also produces a transposition inhibitor from the unspliced P-element RNA, so transposition is inhibited in eggs of P<sup>+</sup> strains.

P-element transposition is promoted by the P-element [transposase](#) (3) and a collaboration with at least one host protein that also interacts with the ends of the P element (4). Interestingly, this host protein is related to other proteins involved with [DNA repair](#). The nature of the collaboration between these proteins that promotes transposition is not yet known. P-element transposition occurs via a “cut and paste” mechanism, in which the element is excised from the donor site and then inserted into the target site. The donor gap resulting from transposon excision can be repaired by double-strand break repair using a sister homologue as a template (5).

The P element has been an immensely powerful tool in studying *Drosophila*. Not long after the discovery and molecular characterization of the P element, it was used in a landmark experiment to introduce novel DNA sequences into the *Drosophila* genome (6). A special P-element derivative was constructed that had, in addition to transposase, a further gene whose activity could be scored visually by inspection of flies. It was found that injection of this modified P element into embryos could result in its integration into the germline and expression in adult flies. Such P-element-mediated transformation is now a commonly used method in *Drosophila* analysis (2). P-element derivatives are also widely used as mutagens, so that the mutations generated are physically tagged with the element; this makes isolation of the target DNA containing the mutated gene very straightforward.

## Bibliography

1. W. R. Engels (1996) *Curr. Top. Microbiol. Immunol.* **204**, 103–124.



2. K. Kaiser, J. W. Sentry, and D. J. Finnegan (1995) In *Mobile Genetic Elements*. (D. J. Sherratt, ed.) IRL Press, Oxford, pp. 69–100.
3. E. L. Beall and D. C. Rio (1998) *EMBO J.* **17**, 2122–2136.
4. E. L. Beall and D. C. Rio (1996) *Genes Dev.* **10**, 921–933.
5. W. R. Engels, D. M. Johnson-Schlitz, W. B. Eggleston, and J. Sved (1990) *Cell* **62**, 515–525.
6. G. M. Rubin and A. C. Spradling (1982) *Science* **218**, 348–353.

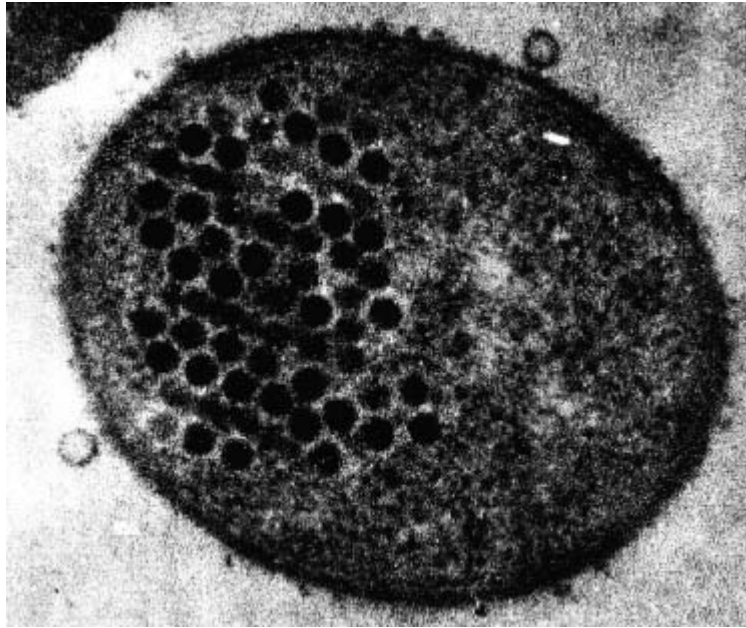
## P22 Bacteriophage

The *Salmonella typhimurium* phage P22 is a double-stranded DNA phage of the Podoviridae family that has played a key role in the development of molecular biology. Studies with P22 have contributed to our understanding of genetic exchange, control of gene expression, and biological self-assembly. It has proven itself to be a valuable laboratory tool for genetic manipulation of its host, *Salmonella*. P22 is classified as a lambdoid phage because it shares life cycle similarities with the *Escherichia coli* [lambda phage](#). Comparison of the similarities and differences between P22 and the other lambdoid phage has served to illuminate both the common themes and the range of diversity that nature employs in problem solving ([1](#)).

P22 was originally isolated as a [lysogen](#) of *Salmonella*, and it has played a key role in the development of bacterial genetics. In 1952, Zinder and Lederberg discovered the phenomenon of generalized transduction of *Salmonella* genes by P22 ([2](#)). This observation led to studies on the mechanism of the establishment of lysogeny, which contributed to our current understanding of repressors and anti-repressors. Studies of the assembly of P22 have been instrumental in framing the concepts that guide our thinking about control mechanisms in biological assembly ([3](#)). In addition, the well-defined and readily manipulable genetics of P22 have allowed it to lead the way in the development of genetic analysis of **protein folding** and assembly ([3-6](#)).

Bacteriophage P22 is a temperate virus, meaning that following infection two distinct life cycles are available to it: the lytic or the lysogenic life cycles. In the lytic life cycle, a series of viral genes are expressed that ultimately result in the production of several hundred progeny virions per cell. The progeny virions are released approximately one hour after infection, through lysis of the host cell (Fig. [1](#)). In the alternative pathway, termed the *lysogenic pathway*, the phage chromosome integrates into the host chromosome, and expression of the phage structural genes is suppressed (see [Lysogeny](#)). The integrated phage is called a *prophage*, which is replicated and passed to the daughter cells during cell division. In response to an appropriate environmental signal, such as irradiation with ultraviolet light, the prophage can excise itself, enter the lytic cycle, and generate progeny phage. The decision between whether an infecting phage enters the lytic or lysogenic pathway depends on a number of factors, the most important of which seems to be the number of infecting virions, or multiplicity of infection.

**Figure 1.** Electron micrograph of a section through a bacteriophage P22-infected *Salmonella typhimurium* cell approximately 30 min after infection. The electron-dense progeny phage within the cell contains double-stranded DNA, whereas the infecting phage that remain attached to the cell at 1 and 8 o'clock do not.



## 1. Structure

The P22 virion consists of two distinct structural elements: the capsid and the tail. The dsDNA phage genome is contained within the capsid and is thereby protected from environmental insult. The capsid of the P22 virion is a  $T = 7$  **icosahedral symmetry** lattice approximately 600 Å in diameter, composed of ~420 identical copies of the 47-kDa **gene product 5 (gp5)** coat protein (7). The capsid itself is a very stable structure, resistant to elevated temperature and drying, as well as to **nucleases** and **proteinases**. Interestingly, this stability does not arise directly from the coat protein subunits, but rather from their intersubunit contacts (8). The coat protein subunits themselves, prior to polymerization into the capsid, are only modestly stable. The tail is composed of up to six trimers of the 72-kDa gp9 tailspike protein. The tailspike protein serves as the organelle for attachment of the virus to the host cell during infection and is itself also a very stable structure (9).

The  $27 \times 10^6$ -kDa linear dsDNA **genome** is contained within the capsid of the phage in a highly condensed liquid-crystalline state. The close packing of the phosphate groups in the condensed form requires charge neutralization, which is accomplished through the binding of  $Mg^{2+}$  ions. Although the precise packing of the DNA within the capsid is not known, the B form of **DNA structure** appears to be retained (10).

The structure of the tailed vertex of the phage is surprisingly complex. The tailspike trimers are attached to the capsid through a connector region. The connector is composed of a macromolecular complex composed of 12 subunits of the gp1 portal protein, as well as the connector proteins, gp4, gp10, and gp26. The 83-kDa gp1 portal protein is arranged as a dodecameric complex, approximately 180 Å in diameter with a central channel approximately 30 Å in diameter (11). It is thought that the DNA enters and exits the head through the central channel of the portal complex. Portal complexes have been identified in nearly all the dsDNA-containing bacteriophage. Although there has been debate as to whether portal complexes are composed of 12 or 13 subunits, the existence of the portal complex at an icosahedral vertex poses an interesting structural dilemma: The portal complex exhibits 12- or 13-fold symmetry, but is located at a fivefold symmetrical vertex, resulting in a symmetry mismatch. This means that there cannot be a simple 1:1 interaction between portal and capsid subunits during assembly and function. Although the exact manner in which the portal protein interacts with the capsid is unknown, it has been suggested that such an arrangement would make it possible for the portal complex to rotate within the capsid during DNA translocation

(12).

The tailspike protein folds and assembles independently of the capsid, and then adds subsequently to the capsid to render the phage infectious. The unique biochemical properties of the tailspike protein, coupled with the well-defined genetics of the bacteriophage system, have made it an ideal subject for *in vivo* studies of protein folding. The tailspike folds in a multistep pathway, in which partially folded intermediates associate to form an SDS-sensitive protrimer (13). The protrimer then matures into a very stable SDS-insensitive trimer. The reason for this stability is that the central region of each subunit of the trimer is folded into a right-handed **b-helix** structure, and the C-terminal regions interdigitate to form **b-sheet** structures (14, 15). The ability to identify readily and quantify the ratio of partially folded and fully folded proteins by the simple technique of SDS-PAGE analysis has allowed the isolation and characterization of a series of mutations that destabilize the partially folded intermediates. Because these mutations have no effect on the stability of the final fully folded trimer, these mutants have been termed temperature sensitive for folding (*tsf*) mutants. Analysis of these mutants has lent support to the idea that protein folding and assembly takes advantage of transient interactions between amino acid residues that do not play a critical role in the stability of the final structure. At restrictive temperatures, the *tsf* mutant proteins form aggregates, indicating that aggregation occurs through the interaction of partially folded protein molecules and that alterations in the lifetime of the protein folding intermediates can result in aggregation (6, 16-18).

## 2. The Infection Process

P22 initially infects the host *Salmonella* cell through the binding of the tailspike to the O-antigenic repeating units of the bacterial lipopolysaccharide. The interaction of the tailspike protein with this receptor defines the host range. Once bound, the tailspike slowly destroys the receptor by cleaving the  $\alpha(1, 3)$ -O-glycosidic bond between rhamnose and galactose units. The virion carries up to six tailspikes, but only three are used for binding during infection (9, 19, 20). Once bound, the phage moves laterally across the cell surface through repeated cycles of binding and release, presumably until it locates a second receptor, at which point binding becomes irreversible (21).

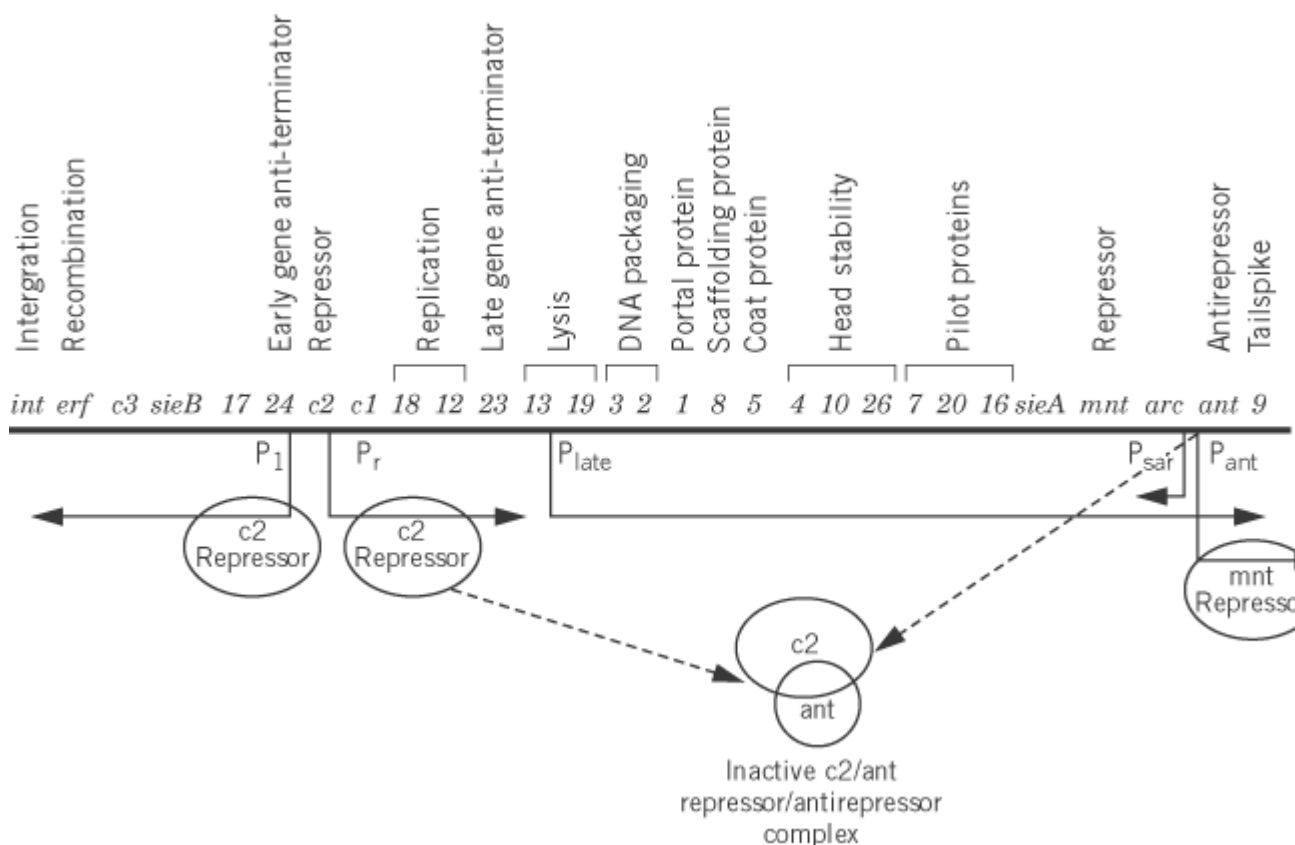
The dsDNA is then injected into the host cell in a linear fashion (22). The phage-encoded pilot proteins gp7, gp16, and gp20 are required to ensure that the DNA crosses the plasma membrane in an infectious form. Although the mechanism of their action remains unknown, they translocate from the phage into the host cell during infection (23). Following infection, the dsDNA circularizes by homologous recombination between the terminal redundancy regions, yielding a circular dsDNA that represents a single copy of the phage genome (24).

## 3. The Lysogenic Pathway

In the lysogenic state, the bacterial chromosome carries an integrated copy of the bacteriophage chromosome, termed a prophage. In the prophage, the genes responsible for growth of the phage and subsequent death of the host are repressed through the action of a specific repressor system. This repression results in a phenomenon known as *immunity*. If a host cell carrying a lysogen is subsequently superinfected by a phage carrying the same repressor system as the prophage, it too will be unable to grow. P22 maintains the prophage in a quiescent state through the action of two repressor proteins (Fig. 2). In this regard, P22 differs from  $\lambda$  phage, which has only a single repressor, *cI*. The two P22 repressor proteins are called *c2* and *mnt*, and they give rise to two immunity systems. The *immC* region is similar to the immunity region of  $\lambda$ . The P22 *c2* repressor acts at  $O_L$  and  $O_R$  to prevent transcription of the early genes. The unique *immI* region includes an antirepressor, *ant*, and its regulators *mnt* and *arc*. *Mnt* acts to repress the transcription of the *ant* antirepressor. If *mnt* is turned off, *ant* is synthesized. *Ant* binds to and inactivates the *c2* repressor protein, with the result that the prophage enters the lytic pathway. The production of *ant* itself is regulated through the action of the *arc* gene. In *arc* mutants, high levels of *ant* are produced, and these levels block phage development. The *arc* and *mnt* proteins are two members of a small family of proteins that use anti-

parallel  $\beta$ -sheet motifs rather than **helix–turn–helix** motifs to bind their operator DNA (25, 26) (see [DNA-Binding Proteins](#)).

**Figure 2.** Genetic map of bacteriophage P22. The genes are positioned above the line (not to scale), with their function indicated above. Promoter regions and the loci of repressor action are indicated below the line. The *c2* repressor protein inhibits synthesis of the early genes,  $P_1$  and  $P_r$ . The anti-terminators 24 and 23 result in expression of the early and late genes, respectively. Notice that genes with related functions are clustered.



The *ant* gene lies downstream of  $P_{late}$  and is transcribed during the lytic cycle as part of the late operon, but Ant protein is not synthesized. An additional promoter  $P_{sar}$  lies within the *ant* gene and directs the synthesis of a 69-nucleotide long antisense RNA that binds to the *ant* messenger RNA and inhibits its translation (27, 28).

In order to form a stable lysogen, the phage genome has to integrate into the host; otherwise it would be diluted out during cell division. Integration occurs preferentially at a particular site, *attB*, in the host chromosome and requires the action of the phage-encoded integrase (Int) protein and the host-encoded integration host factor (IHF). The circular phage chromosome carries a similar site, termed *attP* (29). A single crossover event catalyzed by the Int protein results in integration. For entry into the lytic cycle, excision proceeds by the reverse mechanism but, to ensure that excision does not take place prematurely, requires the additional function of the *xis* gene product, which overlaps *int* in the genetic map (30).

Entry into the lysogenic pathway is controlled by the genes *c1* and *c3*. These gene products stimulate the production of the *c2* repressor protein immediately after infection. *C1* binds to the  $P_{RE}$  (promoter repressor establishment) promoter as a tetramer and stimulates *c2* synthesis (31). High-level

expression of *c2* appears to be required for entry into the lysogenic pathway, and this is the probable reason that a high multiplicity of infection favors lysogeny.

In addition to the immunity conferred by the *immI* and *immC* repressor systems, the phage employs three other mechanisms to prevent superinfection. The first is a modification of the host O antigen, which inhibits the adsorption of P22. This function is encoded by the *a1* gene (32). The genes *sieA* and *sieB* (superinfection exclusion) block superinfection. Expression of *SieB* appears to lead to abortive infection; while the superinfecting phage enters the cell, and early functions are normally expressed, late function expression is inhibited. The *sieA* protein is an 18-kDa protein that partitions into the membrane and acts to block DNA from crossing the periplasmic space and entering the host cell (33). It therefore excludes virion DNA from entering the cells, regardless of the identity of the DNA.

#### 4. The Lytic Pathway

In the lytic pathway, P22 enters the *Salmonella* host cell and begins the process of reproduction. In P22, as in  $\lambda$ , the proteins involved in the lytic cycle are expressed sequentially. Control of the timing of the production of mRNA transcripts is achieved through the use of **anti-termination** proteins that allow the **RNA polymerase** to ignore encoded termination signals. To take advantage of this, the P22 genome is organized so that these regulatory proteins can exert their influence on a large number of genes (34). The genes themselves are clustered into functional units; for example, all the genes involved in assembly are clustered into the late operon and are under the control of a single regulatory protein. Like  $\lambda$ , the early genes of P22 are arranged into two operons,  $P_L$  and  $P_R$ , which flank the P22 *c2* repressor gene. The genes in these operons encode proteins involved in [DNA replication](#), recombination, integration, and the regulation of gene expression. Gene 24, which is transcribed off  $P_L$ , is similar in function to the  $\lambda$  N gene. It functions to prevent the termination of transcription at genetically defined sites in the chromosome and allows for complete transcription of the early genes (35, 36). Gene 23, which is transcribed off  $P_R$ , is analogous to gene Q in  $\lambda$  and allows transcription of the late genes (37). The late genes that are driven off the  $P_{late}$  promoter code for the proteins necessary for head and tail assembly.

#### 5. Generalized Transduction

Transduction is a process by which bacterial DNA is carried from one cell to another by a phage particle. Generalized transducing particles contain DNA derived entirely from the host-cell chromosome. Bacteriophage P22 transducing particles originate when the host chromosome, rather than the phage chromosome, serves as the substrate for packaging. During a lytic infection, the phage chromosome is replicated as a concatamer, and the concatamers are resolved as the double-stranded DNA is packaged by a headful mechanism. The products of phage genes 2 and 3 are required for packaging and act as a complex. Genetic evidence suggests that gp3 is responsible for recognizing a unique site or region termed a *pac* site. Packaging proceeds in an ATP-dependent reaction until the head is full, at which point the gp2/gp3 complex cuts the DNA; packaging then continues into another empty prohead. This model is the “sequential packaging model.” The P22 head can package 105% of a phage genome, accounting for the observed terminal redundancy and circular permutation of the genetic map (38-41).

Within the *Salmonella* chromosome, there are a minimum of five to six sites, called pseudo-*pac* sites, which are recognized by the gp2/gp3 complex as *pac* sites. The nonrandom distribution of pseudo-*pac* sites therefore leads to a nonrandom distribution of transduced markers (42). Mutants that map in gp3, termed HT mutants, have been isolated in which up to 50% of the phage heads carry host DNA. These mutants have proven extremely valuable for studying the genetics of *Salmonella*. The 44-kbp length of the packaged DNA corresponds to about 50 genes and has made possible linkage analysis and genetic mapping of the *Salmonella* chromosome.

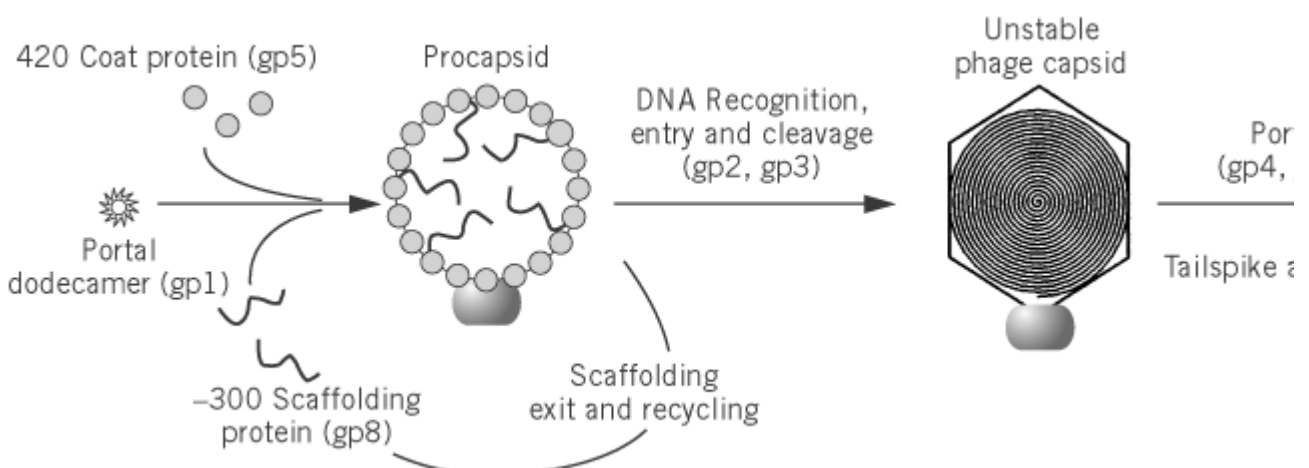
## 6. Specialized Transduction

Specialized transducing particles carry both host and phage DNA linked in a continuous stretch in a single particle. There are two mechanisms by which specialized transducing particles can arise: from aberrant prophage excision or by insertion into the phage genome of host translocatable elements. P22 has more latitude in its ability to mediate specialized transduction than does bacteriophage  $\lambda$ . Because P22 packaging does not use unique end sites, it can package extremely long stretches of DNA by segmenting them into multiple heads. If a host cell is multiply infected, an intact oversized genome can be generated by recombination.

## 7. Capsid Morphogenesis

During P22 morphogenesis, properly folded and associated tailspike trimers add to the dsDNA-containing head, to generate an infectious virion. However, the filled capsid is not the first structure that is assembled (Fig. 3). The gp5 coat protein subunits are first assembled into a structure termed a procapsid (37). Like the mature virion, the procapsid has  $T = 7$  icosahedral symmetry, and the portal protein is present at a single vertex in the form of a dodecameric complex (43, 44). However, the procapsid does not contain DNA. The interior volume is filled with approximately 300 copies of the gp8 scaffolding protein. Studies with **conditional lethal mutants** have demonstrated that the presence of the scaffolding protein facilitates coat protein polymerization and ensures a proper structure. In the absence of scaffolding protein, a variety of morphologically aberrant polymers are formed, while assembly proceeds with high fidelity (45). The procapsid is approximately 10% smaller in diameter and is less stable than the mature head. During DNA packaging, the scaffolding protein exits, and the lattice expands. Three-dimensional **single-particle reconstructions** from [cryoelectron microscopy](#) of the procapsid and mature form revealed that there are  $\sim 25$  Å diameter holes situated at the center of the hexavalent capsomeres in the procapsid and that these holes are closed during expansion. No proteins are added to the procapsid to close the holes; instead, they are closed by **domain** movements within the coat protein subunits (46). These holes are likely candidates for the exit ports for the scaffolding protein during DNA packaging (43). The freed scaffolding protein is recycled to participate in further rounds of assembly (47). Thus, the scaffolding protein functions as an “assembly **molecular chaperone**” whose presence is transiently required to provide positional information. The use of scaffolding proteins is a common theme in the assembly of the lambdoid phages and is also required for the assembly of the **herpesviruses**.

**Figure 3.** The morphogenetic pathway of the bacteriophage P22 capsid. The first structure formed is a procapsid, which coat, scaffolding, and portal proteins. The “pilot” proteins are also incorporated at this stage. The concatameric DNA is packaged with the exit of the scaffolding protein. The addition of gp4, gp10, and gp26 stabilizes the head, and tailspike binding re-



The scaffolding protein of P22, as well as all well-characterized scaffolding proteins, is a highly **α-helical** molecule that forms oligomers in solution (48). In the case of P22, dimerization of the scaffolding protein plays a key role in assembly. In both P22 and the herpes virus, the region of interaction between the coat and scaffolding protein has been mapped to the C-terminal end of the molecule. In addition to its role in assembly, it has been suggested that the scaffolding protein may function to exclude the chance incorporation of cellular proteins during head assembly, because the presence of cellular proteins within the capsid would result in a decreased internal volume and preclude the encapsidation of a complete P22 genome. P22 scaffolding protein regulates its own levels of biosynthesis, which is decreased if there is a large free pool, but up-regulated if all the scaffolding protein is incorporated into procapsid (49-51). This post-transcriptional control might be a mechanism to insure that all procapsids contain a full complement of scaffolding protein, despite the fact that a full complement is not strictly required for assembly (52).

During morphogenesis, the dsDNA is packed into the procapsid, the scaffolding protein exits, and the capsid lattice expands. The exact sequence of these events is not known, but the expansion is an exothermic process, suggesting that the capsid is “spring loaded” (53). Interestingly, the portal protein is part of the gauge that determines when the head is full; mutants in the gp1 portal protein have been identified that result in packaging a piece of P22 DNA up to 5% larger than normal (54). Following DNA packaging, the portal vertex is closed by the addition of the protein products of genes 4, 10, and 26. These proteins, which add sequentially, serve to stabilize the DNA within the capsid and also to provide the site for tail attachment (55).

There are two genes that are essential for lysis, genes 13 and 19. Gene 13 encodes an 11-kDa protein whose function is to disrupt the cell membrane by forming pores and to allow the hydrolytic enzymes access to the cell wall. Mutations in gene 13 have been isolated that delay lysis for up to several hours, despite the fact that the phage-encoded lysozyme is being produced. Gene 19 encodes the P22 lysozyme, a 16-kDa monomeric protein that attacks and degrades the peptidoglycan layer of the cell wall, leading to cell lysis. The bacteriophage P22 lysozyme is very similar to the well-studied T4 lysozyme with respect to structure and function.

Procapsid-like particles of P22 can be assembled *in vitro* from purified proteins (52). This has made it possible to perform sophisticated physical chemical studies of virus assembly. These studies have revealed that capsid assembly is a nucleation-limited reaction (56), proceeds along a well-directed pathway (56), and can be inhibited by the binding of small molecules to the coat protein subunits (57).

Over the past 30 years, the bacteriophage P22 system has contributed to our understanding of gene expression, protein folding and assembly, and basic virology. It appears that the lessons that we can learn from these systems is limited only by the creativity of the investigator.

## Bibliography

1. A. Campbell (1994) *Annu. Rev. Microbiol.* **48**, 193–222.
2. N. D. Zinder and J. Lederberg (1952) *J. Bacteriol.* **64**, 679–699.
3. J. King, R. Griffin-Shea, and M. T. Fuller (1980) *Q. Rev. Biol.* **55**, 369–393.
4. J. Jarvik and D. Botstein (1975) *Proc. Natl. Acad. Sci. USA* **72**, 2738–2742.
5. J. King and M. H. Yu (1986) *Methods Enzymol.* **131**, 250–266.
6. J. King, C. Haase-Pettingell, A. S. Robinson, M. Speed, and A. Mitraki (1996) *Faseb. J.* **10**, 57–66.
7. S. Casjens (1979) *J. Mol. Biol.* **131**, 1–14.
8. P. E. Prevelige, Jr., J. King, and J. L. Silva (1994) *Biophys. J.* **66**, 1631–1641.
9. V. Israel (1976) *J. Virol.* **18**, 361–364.

10. K. L. Aubrey, S. Casjens, and G. J. Thomas, Jr. (1992) *Biochemistry* **31**, 11835–11842.
11. C. Bazinet, J. Benbasat, J. King, J. M. Carazo, and J. L. Carrascosa (1988) *Biochemistry* **27**, 1849–1856.
12. R. W. Hendrix (1978) *Proc. Natl. Acad. Sci. USA* **75**, 4779–4783.
13. D. P. Goldenberg, D. H. Smith, and J. King (1983) *Proc. Natl. Acad. Sci. USA* **80**, 7060–7064.
14. D. P. Goldenberg and T. E. Creighton (1994) *Curr. Biol.* **4**, 1026–1029.
15. S. Steinbacher, R. Seckler, S. Miller, B. Steipe, R. Huber, and P. Reinemer (1994) *Science* **265**, 383–386.
16. C. A. Haase-Pettingell and J. King (1988) *J. Biol. Chem.* **263**, 4977–4983.
17. A. Mitraki, B. Fane, C. H. Haase-Pettingell, J. Sturtevant, and J. King (1991) *Science* **253**, 54–58.
18. M. A. Speed, T. Morshead, D. I. Wang, and J. King (1997) *Protein Sci.* **6**, 99–108.
19. V. Israel (1978) *J. Gen. Virol.* **40**, 669–673.
20. S. Iwashita and S. Kanegasaki (1976) *Eur. J. Biochem.* **65**, 87–94.
21. U. Baxa, S. Steinbacher, S. Miller, A. Weintraub, R. Huber, and R. Seckler (1996) *Biophys. J.* **71**, 2040–2048.
22. M. Rhoades, L. A. MacHattie, and C. A. Thomas, Jr. (1968) *J. Mol. Biol.* **37**, 21–40.
23. V. Israel (1977) *J. Virol.* **23**, 91–97.
24. D. Botstein and M. J. Matz (1970) *J. Mol. Biol.* **54**, 417–440.
25. M. J. Burgering, R. Boelens, D. E. Gilbert, J. N. Breg, K. L. Knight, R. T. Sauer, and R. Kaptein (1994) *Biochemistry* **33**, 15036–15045.
26. B. E. Raumann, M. A. Rould, C. O. Pabo, and R. T. Sauer (1994) *Nature* **367**, 754–757.
27. T. H. Wu, S. M. Liao, W. R. McClure, and M. M. Susskind (1987) *Genes Dev.* **1**, 204–212.
28. S. M. Liao, T. H. Wu, C. H. Chiang, M. M. Susskind, and W. R. McClure (1987) *Genes Dev.* **1**, 197–203.
29. L. Smith-Mungo, I. T. Chan, and A. Landy (1994) *J. Biol. Chem.* **269**, 20798–20805.
30. J. M. Leong, S. E. Nunes-Duby, A. B. Oser, C. F. Lesser, P. Youderian, M. M. Susskind, and A. Landy (1986) *J. Mol. Biol.* **189**, 603–616.
31. Y. S. Ho, D. Pfarr, J. Strickler, and M. Rosenberg (1992) *J. Biol. Chem.* **267**, 14388–14397.
32. M. Gough and J. V. Scott (1972) *Virology* **50**, 603–605.
33. B. Hofer, M. Ruge, and B. Dreiseikelmann (1995) *J. Bacteriol.* **177**, 3080–3086.
34. M. M. Susskind and D. Botstein (1978) *Microbiol. Rev.* **42**, 385–413.
35. S. Hilliker and D. Botstein (1975) *Virology* **68**, 510–524.
36. J. W. Roberts, C. W. Roberts, S. Hilliker, and D. Botstein (1976) in R. Losick and M. Chamberlin, eds., *RNA Polymerase*, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.
37. D. Botstein, C. H. Waddell, and J. King (1973) *J. Mol. Biol.* **80**, 669–695.
38. B. K. Tye, J. A. Huberman, and D. Botstein (1974) *J. Mol. Biol.* **85**, 501–528.
39. B. K. Tye and D. Botstein (1974) *J. Supramol. Struct.* **2**, 225–238.
40. S. Casjens and M. Hayden (1988) *J. Mol. Biol.* **199**, 467–474.
41. S. Casjens, W. M. Huang, M. Hayden, and R. Parr (1987) *J. Mol. Biol.* **194**, 411–422.
42. H. Schmieger (1982) *Mol. Gen. Genet.* **187**, 516–518.
43. B. V. Prasad, P. E. Prevelige, E. Marietta, R. O. Chen, D. Thomas, J. King, and W. Chiu (1993) *J. Mol. Biol.* **231**, 65–74.
44. P. A. Thuman-Commike, B. Greene, J. Jakana, B. V. Prasad, J. King, P. E. Prevelige, Jr., and W. Chiu (1996) *J. Mol. Biol.* **260**, 85–98.
45. W. Earnshaw and J. King (1978) *J. Mol. Biol.* **126**, 721–747.



46. P. E. Prevelige, Jr., D. Thomas, K. L. Aubrey, S. A. Towse, and G. J. Thomas, Jr. (1993) *Biochemistry* **32**, 537–543.
47. J. King and S. Casjens (1974) *Nature* **251**, 112–119.
48. M. H. Parker, W. F. Stafford III, and P. E. Prevelige, Jr. (1997) *J. Mol. Biol.* **268**, 655–665.
49. J. King, C. Hall, and S. Casjens (1978) *Cell* **15**, 551–560.
50. S. Casjens and M. B. Adams (1985) *J. Virol.* **53**, 185–191.
51. S. Casjens, M. B. Adams, C. Hall, and J. King (1985) *J. Virol.* **53**, 174–179.
52. P. E. Prevelige, Jr., D. Thomas, and J. King (1988) *J. Mol. Biol.* **202**, 743–757.
53. M. L. Galisteo and J. King (1993) *Biophys. J.* **65**, 227–235.
54. S. Casjens, E. Wyckoff, M. Hayden, L. Sampson, K. Eppler, S. Randall, E. T. Moreno, and P. Serwer (1992) *J. Mol. Biol.* **224**, 1055–1074.
55. H. Strauss and J. King (1984) *J. Mol. Biol.* **172**, 523–543.
56. P. E. Prevelige, Jr., D. Thomas, and J. King (1993) *Biophys. J.* **64**, 824–835.
57. C. M. Teschke, J. King, and P. E. Prevelige, Jr. (1993) *Biochemistry* **32**, 10658–10665.

### Suggestions for Further Reading

58. M. M. Susskind and D. Botstein (1978) Molecular genetics of bacteriophage P22. *Microbiol. Rev.* **42**, 385–413. (An outstanding review of the genetics of bacteriophage P22).
59. P. E. Prevelige Jr. and J. King (1993) Assembly of bacteriophage P22: a model for dsDNA virus assembly. *Prog. Med. Virol.* **40**, 206–221.

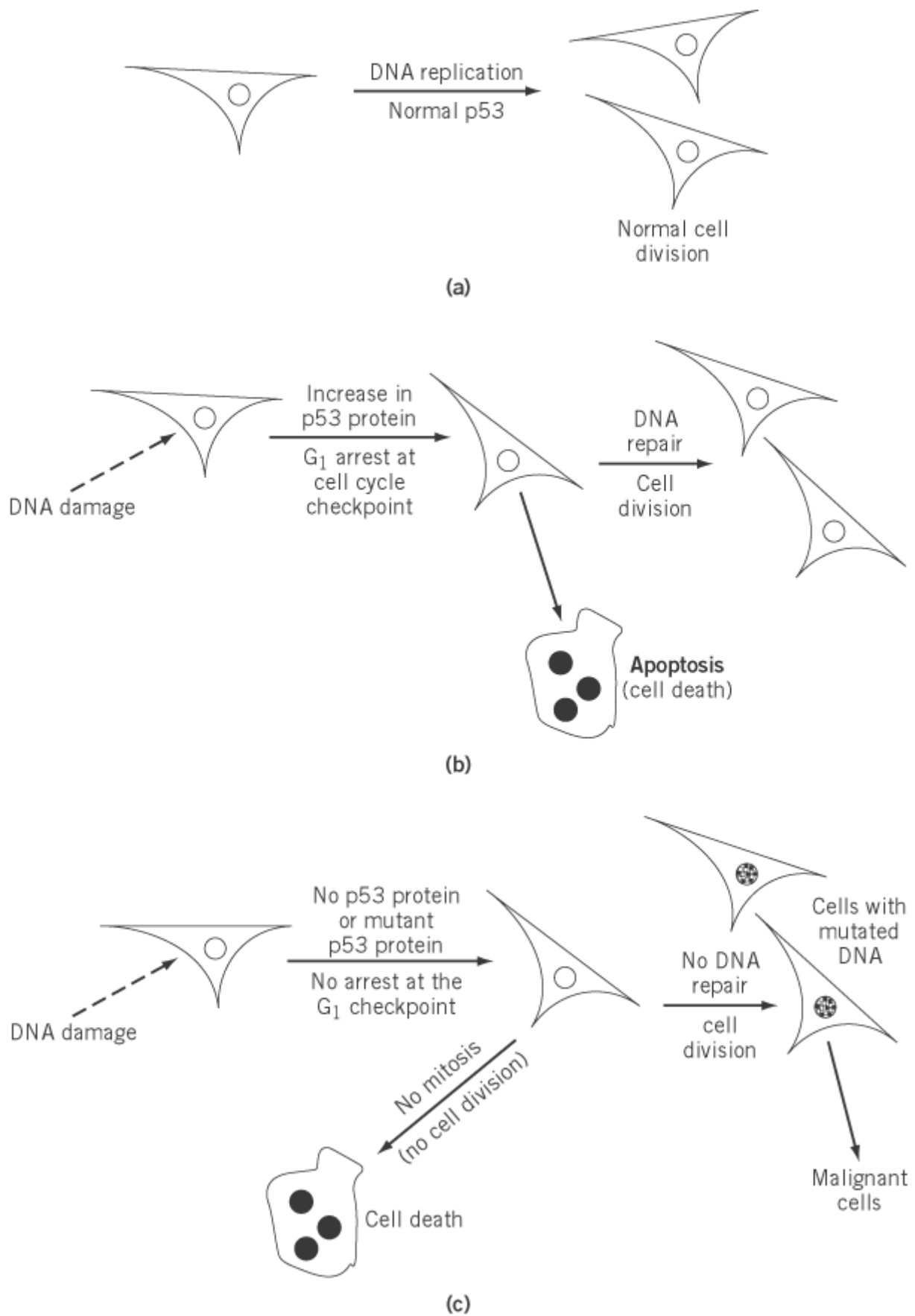
## P53

p53 is a well-characterized [tumor suppressor gene](#). It is located on human [chromosome](#) 17p13.1 and encodes a [polypeptide chain](#) of 375 amino-acid residues. Like the [retinoblastoma gene](#), loss of p53 function results in the onset of oncogenesis (see **Oncogenes** and [Neoplastic Transformation](#)). p53 is a [transcription factor](#) that regulates the expression of **cell-cycle** regulatory proteins. The p53 protein contains an N-terminal transactivation **domain**, a central **DNA-binding** domain, and a C-terminal regulatory domain. The C-terminal domain of p53 mediates oligomerization of the protein, which is essential for its DNA-binding activity. This domain also mediates its interaction with other cellular factors that enhance the DNA-binding activity of the protein.

P53 plays a pivotal role during the execution of [apoptosis \(programmed cell death\)](#) pathways induced by **DNA damage** caused by agents, such as ultraviolet or ionizing radiation (Fig. 1). Current models suggest that DNA damage results in up-regulating p53 gene [transcription](#) and accumulation of this protein in the cell (1). Accumulation of p53 appears to result in the enhanced transcription of a group of genes, such as the cell cycle-dependent kinase inhibitor p21/WAF-1/cip-1, that are responsible for stalling the cell cycle in either the G<sub>1</sub>- or G<sub>2</sub>/M-phase. This mechanism provide the cell with valuable time to make critical repairs in its genetic material. In the event that the necessary repairs cannot be made, the p53 protein initiates a “suicidal” apoptotic program, resulting in cell death, to prevent clonal expansion of a cell with a mutated genome.

**Figure 1.** Proposed function of p53. (a) P53 is not required for normal cell division. (b) In response to DNA damage,

p53 levels increase and arrest the cell at the first of two “checkpoints” in the G<sub>1</sub>-phase of the cell cycle. If all of the repairs have been made to the DNA, the cell divides normally and completes the cell cycle. However, if the cell still contains mutated or duplicated DNA sequences, it dies by a suicidal apoptotic mechanism to prevent its expansion. (c) In those cells that have mutated or lost p53 genes, the arrest at G<sub>1</sub> does not occur, and the cells that have mutated genomes proliferate and become cancerous.



Malignant cells have evolved novel strategies to destroy p53 function, and approximately 50% of human tumors contain mutations or deletions in this gene. Normal p53 protein binds to DNA and

transactivates the transcription of a distinct set of target genes. Mutated forms of p53 cannot bind DNA in this fashion, however, and therefore behave abnormally. A second mechanism by which normal p53 function is subverted in a cancerous cell is through the MDM2 protein. The MDM2 protein binds to p53, and it is amplified in many tumor cells. When these two proteins are bound, p53 cannot mediate growth arrest at the cell cycle checkpoints (as described previously), which leads to rapid clonal expansion of cells that have unstable, mutated genomes.

## Bibliography

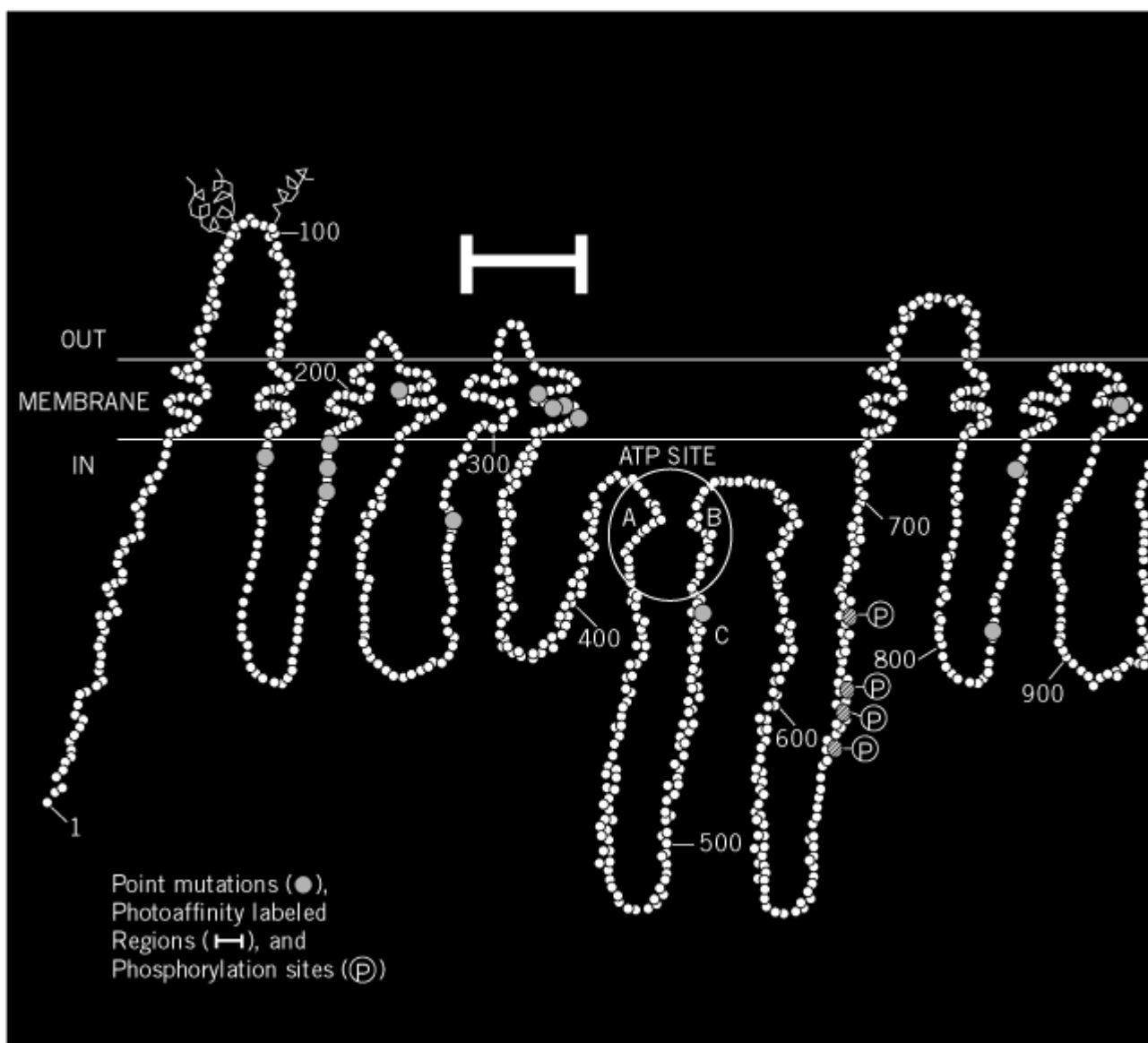
1. L. J. Ko and C. Prives (1996) *Genes Dev.* **10**, 1054–1072.

## P-Glycoproteins

The best characterized form of [drug resistance](#) in human cells is due to the overexpression of P-glycoprotein (P-gp). This 170-kDa plasma [membrane protein](#) is a member of the ATP Binding Cassette (ABC) superfamily of transporter proteins; it can extrude a range of [hydrophobic](#) anticancer drugs from the cell against a concentration gradient, and thus protect the cell. Consequently, cells that overexpress P-gp show reduced intracellular concentrations of drugs and greater resistance to them ([1](#)). Molecular probes developed to identify P-gp in cancer tissues also strongly suggest that, at least in hematological malignancies, P-gp plays a significant role in multiple drug resistance (MDR) in the clinic.

The mammalian *mdr* **gene family** consists of two members in humans and three in rodents ([2](#)). Of the two human genes, *MDR1* and *MDR2*, only *MDR1* confers drug resistance. The sequences of the mouse *mdr1* and *mdr3* are 80% identical to the *MDR1* human gene product, and they are also capable of conferring resistance to drugs, as are their hamster homologues, *pgp1* and *pgp2*. The human *MDR1* gene encodes a 1280-amino acid membrane protein, P-gp. [Hydropathy](#) plots of this protein have led to a **secondary structure** model for P-gp (see Fig. [1](#)) in which it consists of 12 putative transmembrane regions, with two **ATP-binding** domains on the cytoplasmic face of the membrane. Each half of the protein consists of six hydrophobic transmembrane putative **alpha-helices** and one ATP-binding domain; the two halves are 43% identical ([3](#)).

**Figure 1.** A hypothetical two-dimensional model of human P-glycoprotein and its functional domains. This model, with ATP sites (nucleotide binding domains), is based on hydropathy analysis of the amino acid sequence. The ATP sites are indicated. The N-linked glycosylation sites are indicated by wavy lines, and the known phosphorylation sites are shown : show the regions labeled with photoaffinity analogues. Filled circles represent amino acid residues in which mutations h transport specificity. This figure is adapted from Ref. [1](#).



Drug resistance conferred by the human and rodent P-gp's has been extensively studied. However, there continue to be controversies and alternate models vis-à-vis the mechanism by which an extremely broad range of substrates are expelled from the cells. The most widely accepted model is that the P-gp uses ATP to power a molecular pump to remove chemotherapeutic agents from the cell (1). According to this hypothesis, chemotherapeutic agents diffuse down a concentration gradient into the cell, and the pump either acts as a “flipase” to expel them as they enter the bilayer, or the drugs are first transported from the cytosol to the membrane and subsequently pumped out from the membrane. An alternative mechanism has also been suggested in which changes in the intracellular pH and membrane potential alter the transmembrane partitioning or intracellular sequestering of the drugs (4). As these are the two principal models on which most current work is based, the evidence on which each is based will be considered briefly.

The alternative partitioning model is essentially a passive diffusion model (4). The central notion is that the drugs are sufficiently hydrophobic to cross the membrane spontaneously. The subsequent asymmetric distribution of these drugs is then a consequence of the asymmetry of the chemical potential (pH and membrane potential). It is argued that because the cytosol of tumor cells is acidic, the drugs (most of which are weak bases, with  $pK_a$  values 7.4 to 8.2) get protonated and are trapped in the cytosol. Conversely, an increased cytosolic pH would decrease the intracellular accumulation

of drugs. Experimental evidence in support of this scenario comes from a number of observations. In model systems such as phospholipid vesicles, it is possible to modulate the transmembrane distribution of some of the drugs implicated in MDR by altering the pH. The cytosolic pH is higher in tumor cells exhibiting MDR and those transfected with P-gp. Verapamil, which reverses MDR, causes a partial acidic shift in the pH of cells. The model suffers, however, from a number of shortcomings. It cannot explain how drugs that are not weak bases are effluxed. A number of cells exhibiting MDR do indeed show pH shifts in agreement with the model, but a significant number of cells do not show these changes. More importantly, even where these changes in pH do occur, they are in most instances not of a magnitude sufficient to explain the several hundred-fold increase in resistance that these cells display.

The active transport model for MDR is more attractive and based on a large body of currently available evidence (5, 6). Drug binding and photoaffinity studies show a direct interaction between P-gp and many of the substrates. The drugs stimulate the [ATPase](#) activity of P-gp in proportion to its ability to transport those drugs. Specific amino acid substitutions alter the substrate specificity of P-gp. Additionally, purified P-gp reconstituted into phospholipid vesicles is capable of drug transport, even in the absence of electrochemical gradients. MDR is a complex phenomenon, and it is conceivable that the change in electrochemical gradients may be the result of epiphenomena during the prolonged selection of cells in cytotoxic drugs. Perhaps the best prospect of resolving this issue lies in the use of transient expression systems. A **vaccinia virus**-based expression system has been developed recently that allows high-level expression of P-gp in many mammalian cells (7). The use of this system clearly excludes the pleiotropic effects associated with the use of drugs in selecting cultured cells.

Studies with both cultured cells and transient expression systems provide compelling evidence that the substrates interact directly with P-gp. Table 1 lists some of the agents that interact with P-gp. This has also been demonstrated by the use of diverse substrate analogs that photoaffinity label the P-gp and yielded valuable information on the direct interaction of P-gp with its substrates: <sup>125</sup>I-labeled iodoazylazidoprazoin (IAAP), an analog of prazosin (8); <sup>125</sup>I-labeled iodomyacin, an analog of daunorubicin (9); <sup>3</sup>H-labeled azidopine (10); and 6-*O*-[2-[3-(4-azido-3-[<sup>125</sup>I]iodophenyl)propionamido]ethylcarbonyl] forskolin (AIPPF), an analog of forskolin (11). While the nature of these interactions is still unknown, photoaffinity experiments with these analogs and [site-directed mutagenesis](#) analysis indicate that the interaction is probably in the regions of transmembrane segments 5, 6, 11, and 12.

**Table 1. List of Selected Substrates and Modulators of P-Glycoprotein**

| <b>Substrates of P-glycoprotein</b> | <b>Agents that Reverse Multidrug Resistance</b> |
|-------------------------------------|---|
| <b><i>Vinca</i> alkaloids</b>       | <b>Calcium channel blockers</b>                 |
| Vinblastine                         | Verapamil                                       |
| Vincristine                         | Dihydropyridines                                |
| Anthracyclines                      | Azidopine                                       |
| Daunorubicin                        | Antiarrhythmics                                 |
| Doxorubicin                         | Quinine   |
| Epipodophyllotoxins                 | Quinidine                                       |
| Etoposide                           | Antihypertensives                               |
| Teniposide                          | Reserpine                                       |
| Antibiotics                         | Yohimbine                                       |

|                        |                            |
|------------------------|----------------------------|
| Dactinomycin           | Antibiotics                |
| Actinomycin D          | Hydrophobic cephalosporins |
| Other cytotoxic agents | Immunosuppressants         |
| Mytomycin              | Cyclosporine A             |
| Taxol                  | FK506                      |
| Topotecan              | Steroid hormones           |
| Colchicine             | Progesterone               |
| Emetine                | Megestrol acetate          |
| Gramicidin D           | HIV protease inhibitors    |
| Puromycin              | Sequinavir                 |
| Valinomycin            | Indinavir                  |
|                        | Ritonavir                  |

---

The notion that the energy of ATP hydrolysis is utilized by P-gp to pump drugs actively is central to the active pump model. Consequently, the substrate-stimulated ATPase activity of P-gp has been studied in considerable detail. P-gp consists of two transmembrane segments and two ATP-binding sites. Each ATP site is comprised of three conserved regions: Walker A and Walker B motifs and a hydrophobic dodecapeptide—a signature of ABC superfamily members, also called the “C” region. ATP hydrolysis is abolished by chemical modification with *N*-ethylmaleimide or mutations in the conserved residues of the Walker A or Walker B region. Orthovanadate, which behaves as an analogue of inorganic phosphate, inhibits the ATPase by trapping Mg-ADP at the catalytic site. Mutational analysis and chemical modification have established that both sites can bind and hydrolyze ATP. Moreover, the interaction of the two sites is essential for ATP hydrolysis and drug transport. Purified P-gp catalyzes substrate-stimulated ATP hydrolysis. The protein, however, also has a basal ATPase activity in the absence of substrate, which may be due to activation by endogenous lipids or peptides. Based on current evidence, an alternating catalytic model of ATP hydrolysis proposes that, while both sites bind ATP, only one site at any given time acts as a catalytic site. The conformation of this site is presumed to prevent the other from hydrolyzing ATP (12).

### 1. P-glycoproteins in normal cells

The wide distribution of P-gp in normal human tissues, such as the blood–brain barrier, liver, kidney, intestine, adrenal glands, and testes, has long prompted speculation on the role of P-gp in normal cells (1, 5). There is also a more urgent need to understand the normal physiological role of P-gp. Effective inhibitors of P-gp-mediated drug transport are being developed and entering clinical trials. The safe use of these inhibitors requires that the functions of P-gp in normal physiology be understood so as to anticipate, and eventually limit, the potential side effects arising from their use. In recent years, the use of knockout mice (see [Mouse](#)) with disrupted P-gp genes provided a direct method of studying these functions.

Mice with disruptions of each of the three P-gp-encoding genes (*mdr1a*, *1b*, and *mdr2*) have been obtained (13). Under laboratory conditions, all the knockout mice are healthy, fertile, normal anatomically and histologically, and have a normal life span. The most striking result obtained with the mice homozygous for a disrupted *mdr1a* gene was the role of P-gp as an active extruder of molecules that pass the blood–brain barrier. For example, ivermectin, an excellent P-gp substrate, accumulates to levels over 100-fold higher in the brains of *mdr1a* (–/–) mice than in *mdr1a* (+/+) mice. The tolerance of *mdr1a* (–/–) mice to ivermectin was also reduced 100-fold. The results from

the double knockout mice (*mdr1a + mdr1b*) are particularly important because, unlike humans, mice have two functional drug-transporting genes, and only a double knockout can tell us about the untoward effects that may be expected by the use of powerful MDR reversal agents. The double knockouts have not yet been analyzed in detail, but there are no gross physiological, anatomical, or pathological abnormalities. Particularly important is the fact that gross disturbances in corticosteroid metabolism, pregnancy, and bile formation, which would reasonably be expected to occur based on earlier speculations of the function of P-gp, are absent in these mice. These studies suggest a protective role against chemical toxicity for P-gp in humans.

### Bibliography

1. M. M. Gottesman and I. Pastan (1993) *Ann. Rev. Biochem.* **62**, 385–427.
2. W. F. Ng, F. Sarangi, R. L. Zastawy, L. Veinot-Drebot, and V. Ling (1989) *Mol. Cell. Biol.* **9**, 1224–1232.
3. C-J. Chen, J. E. Chin, K. Ueda, D. P. Clark, I. Pastan, and M. M. Gottesman (1986) *Cell* **47**, 381–389.
4. P. D. Roepe (1995) *Biochim. Biophys. Acta.* **1241**, 385–405.
5. M. M. Gottesman, C. A. Hrycyna, P. V. Schoenlein, U. A. Germann, and I. Pastan (1995) *Annu. Rev. Genet.* **29**, 607–649.
6. F. J. Sharom (1997) *J. Membrane Biol.* **160**, 161–175.
7. M. Ramachandra, S. V. Ambudkar, I. Pastan, M. M. Gottesman, and C. A. Hrycyna (1996) *Mol. Biol. Cell.* **7**, 1485–1498.
8. L. M. Greenberger (1993) *J. Biol. Chem.* **268**, 11417–11425.
9. A. Demmer, H. Thole, P. Kubesch, T. Brandt, M. Raida, R. Fislage, and B. Tummler (1997) *J. Biol. Chem.* **272**, 20913–20919.
10. E. P. Bruggemann, U. A. Germann, M. M. Gottesman, and I. Pastan (1989) *J. Biol. Chem.* **264**, 15483–15488.
11. D. I. Morris, L. M. Greenberger, E. P. Bruggeman, C. Carderelli, M. M. Gottesman, and I. Pastan (1994) *Mol. Pharmacol.* **46**, 329–337.
12. A. E. Senior, M. K. Al-Shawi, and I. L. Urbatsch (1995) *FEBS Lett.* **377**, 285–289.
13. P. Borst and A. H. Schinkle (1996) *Eur. J. Cancer* **32**, 985–990.

### Suggestions for Further Reading

14. M. M. Gottesman, I. Pastan, and S. V. Ambudkar (1996) *Curr. Opin. Genet. Dev.* **6**, 610–617. A recent review of the molecular biology, biochemistry, and energetics of P-glycoprotein, with a very helpful annotated list of references.
15. W. Stein (1997) *Physiol. Rev.* **77**, 545–590. An excellent review on the kinetics of P-glycoprotein drug transport function, which includes very detailed and often quantitative discussions.

### P-Loop

The P-loop is a fingerprint sequence that characterizes ATP- or [GTP-binding proteins](#) and is part of the [nucleotide-binding motif \(1\)](#). The characteristic P-loop sequence is Gly–X–X–Gly–X–Gly–Lys–Thr/Ser, where X is any residue. In the three-dimensional [protein structures](#) that have a P-loop sequence, such as adenylate kinase, the loop is located between a **b-strand** and an **a-helix** and forms



a binding site for the phosphate group of the nucleotide—hence the name P-loop. The loop wraps around the phosphate, with the [glycine](#) residues adopting conformations that are forbidden for other residues. The backbone amides and the conserved [lysine](#) residue interact with the phosphate of the nucleotide.

The P-loop fingerprint sequence is used to identify ATP- or GTP- binding proteins from their primary structure (1). It has also been used to classify a broad group of protein structures, the nucleoside triphosphate hydrolase fold (2), which has a P-loop structure within an a/b domain comprising a central b-sheet surrounded by a-helices. Different families within this classification have varying numbers of b-strands and different connectivities.

[See also [Nucleotide-Binding Motif](#).]

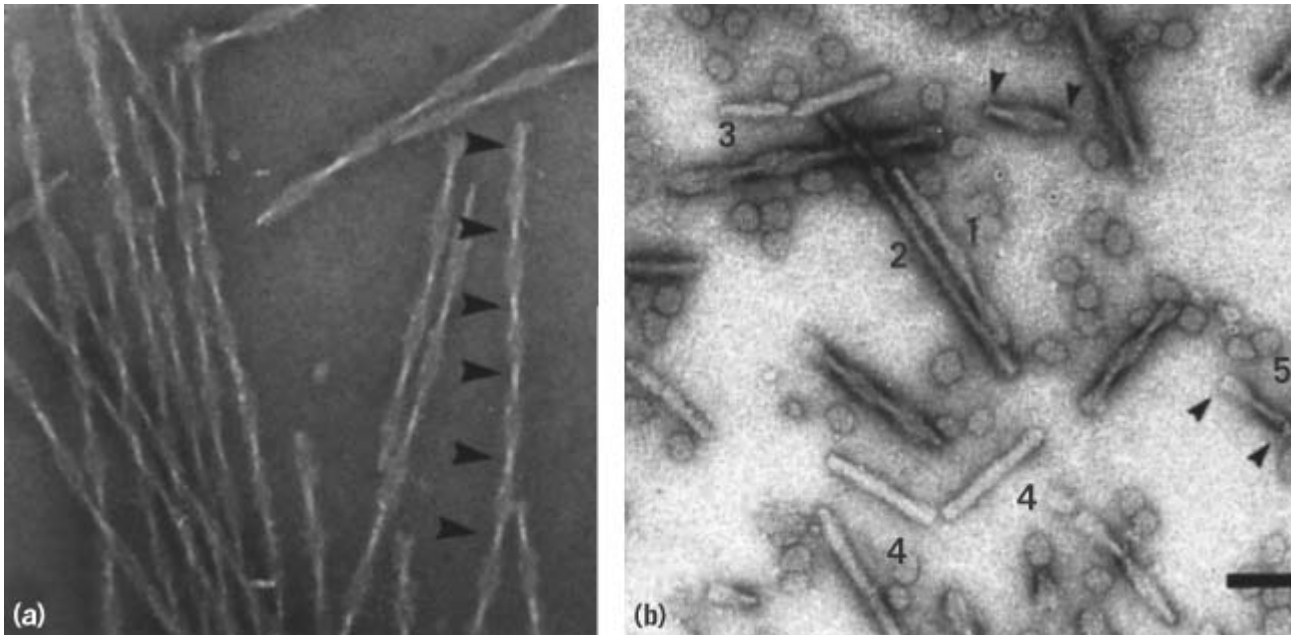
#### Bibliography

1. M. Saraste, P. R. Sibbald, and A. Wittinghofer (1990) *TIBS* **15**, 430–434.
2. A. G. Murzin, S. E. Brenner, T. Hubbard, and C. Chothia (1995) *J. Mol. Biol.* **256**, 126–143.

#### Paired Helical Filaments

Alzheimer's disease is characterized by amyloid plaques and intraneuronal inclusions or tangles in the cell body and the proximal axon and dendrites (1). The tangles are comprised of paired helical filaments (PHFs) (Fig. 1), but these bear no relationship to the [neurofilament](#) type of [intermediate filaments](#) (IFs). Indeed, PHF are composed of  $\tau$  (**tau protein**), a [microtubule-associated protein](#) that is [hydrophilic](#), rich in [proline](#), [glycine](#), [serine](#) and [threonine](#) residues, and normally highly soluble. In the PHFs, however, the  $\tau$  protein is in an abnormal state of **phosphorylation**. In adult humans,  $\tau$  is not a single protein but a family of six variants encoded by the same **gene** generated by [alternative splicing](#). All of the  $\tau$  chains have a related **domain** structure, a fundamental part of which is three or four tandem repeats, each of which is normally involved in microtubule-binding. These repeats are about 31 or 32 residues in extent and are only imperfectly preserved (2), but they have homology with other microtubule-associated proteins. The human isoforms of  $\tau$  vary also in their *N*-terminal regions, in that some contain an insert of 29 or 58 residues. Trojanowski and Lee (2) thus note that the largest  $\tau$  isoform contains 441 residues and has a 58-residue *N*-terminal insert and four tandem microtubule-binding repeats; the smallest isoform is 352 residues long and has no *N*-terminal insert and only three microtubule-binding repeats. The normal role of  $\tau$  *in vivo* is to stabilize the [microtubules](#) in the axon, thus facilitating microtubule-based axonal transport.  $\tau$  can readily be phosphorylated by the action of a variety of **kinases**; however, specific kinases acting on a serine residue in the first repeat domain cause the dissociation of  $\tau$  from microtubules. Although there is no evidence that any part of  $\tau$  adopts a regular **secondary structure**, its assembly into PHF is thought to be mediated through the repeat domains. The paired helical filaments thus formed consist of two filaments, each about 10 nm in diameter. They coil around one another with a pitch length of about 150 nm (Fig. 1). Individual  $\tau$  molecules have been observed in the electron microscope, where they appear as rodlike structures with lengths of about 35 nm. Proteolytic removal of the *N*-terminal region and part of the *C*-terminal region of  $\tau$  results in PHF with their fuzzy coats removed. Each strand in section appears to have a three-domain C-shaped structure. This also implies that the repeat domains form the core of the PHF and that the terminal ends of the  $\tau$  chains lie on the outer surface.

**Figure 1.** Electron micrographs of paired helical filaments (PHFs) from Alzheimer brains. (a) PHF from neurofibrillary et al. (5), stained with 1% phosphotungstic acid. This preparation contains homogeneous long filaments that still retain t coat.” The cross-over repeat is 75 to 80 nm, and the width varies between a minimum of about 10 nm and a maximum of Greenberg and Davies (6). This preparation results in soluble filaments of shorter length than in (a) and is more heterogeneous (2) is a straight filament with width 8 nm; (3) is a twisted filament with diameter (up to 25 nm); (4) is a straight filament with a wide diameter (18 nm); (5) is a twisted rodlike particle about 80 nm long and 18 nm wide. In many cases the particles appear to have broken apart across the filament—for example, the two rods of (3) and the short stub to the right of it, or the two straight rods above particle (3). Bar = 100 nm. (From Wille et al. (3) Mandelkow.)



A major problem in this field had been the difficulty in reassembling the entire  $\tau$  molecule *in vitro* into a structure that closely resembles native PHF, although this has now been achieved. Wille et al. (3) and Schweers et al. (4) have described the reproducible formation of PHF from recombinant  $\tau$  protein and from constructs containing three of the repeating domains. Their results showed that the state of phosphorylation of  $\tau$  was not a determinant of PHF formation, that the repeat domain in  $\tau$  can form antiparallel dimers cross-linked by [disulfide bonds](#) to their [cysteine](#) residues, and that PHF can be formed from these dimers. Inhibition of dimer formation resulted in no PHF formation, indicating that it might be the key initial step. Schweers et al. (4) have concluded that two residues of particular importance in  $\tau$  are the serine in the first repeat (residue 262 in the htau40 isoform) and the cysteine in the third repeat (residue 322 in htau40): phosphorylation of the serine residue leads to dissociation of  $\tau$  from the microtubules and allows other modes of aggregation to occur, whereas oxidation of the cysteine leads to dimer formation followed by PHF assembly.

### Bibliography

1. M. L. Shelanski and C. M. Troy (1990) "The Cytoskeleton in Neurological Disease". In *Cellular and Molecular Biology of Intermediate Filaments* (R. D. Goldman and P. M. Steinert, eds.), Plenum Press, New York, pp. 451–465.
2. J. Q. Trojanowski and V. M.-Y. Lee (1995) Phosphorylation of paired helical filament tau in Alzheimer's disease neurofibrillary lesions: focusing on phosphatases. *FASEB J.* **9**, 1570–1576.
3. H. Wille, G. Drewes, J. Biernat, E.-M. Mandelkow and E. Mandelkow (1992) Alzheimer-like paired helical filaments and antiparallel dimers formed from microtubule-associated protein tau *in vitro*. *J. Cell Biol.* **118**, 573–584.
4. O. Schweers, E.-M. Mandelkow, J. Biernat, and E. Mandelkow (1995) Oxidation of cysteine-322 in the repeat domain of microtubule-associated protein  $\tau$  controls the *in vitro* assembly of paired

helical filaments. Proc. Natl. Acad. Sci. USA **92**, 8463–8467.

5. C. Wischik, R. Crowther, M. Stewart, and M. Roth (1985) Subunit structure of paired helical filaments in Alzheimer's disease. J. Cell Biol. **100**, 1905–1912.
6. S. E. Greenberg and P. Davies (1990) A preparation of Alzheimer paired helical filaments that displays distinct tau-proteins by polyacrylamide-gel electrophoresis. Proc. Natl. Acad. Sci. USA **87**, 5827–5831.

### Suggestions for Further Reading

7. M. Goedert, R. Jakes, M. G. Spillantini, and R. A. Crowther (1994) "Tau protein and Alzheimer's disease". In *Microtubules* (J. S. Hyams and C. W. Lloyd, eds.), Wiley-Liss, New York, 183–200.

## Pair-rule Genes

Pair-rule genes play critical roles in the formation of the ***Drosophila melanogaster*** body plan, which consists of fourteen contiguous segments arranged along the anterior-posterior axis. This basic plan is established by a cascade of **segmentation gene** activities that generate spatially repeated patterns of gene expression before the completion of cellularization. Most pair-rule genes are initially expressed in patterns of seven evenly spaced stripes about four nuclei wide that represents the first evidence of a segmental body plan. The striped patterns are established by integrating crude positional cues provided by asymmetrically localized maternal factors and proteins encoded by the **gap genes**. These cues are in the form of gradients of [transcription factors](#) that interact directly with pair-rule regulatory regions, resulting in position-specific activation or repression. In several cases, pair-rule genes have been shown to contain discrete [enhancers](#) that direct the expression of single stripes of pairs of stripes. Each enhancer responds in a unique way to the asymmetric cues, and the seven stripe entire pattern represents the sum of the effects of all enhancers. Once initial stripes are established, they are refined by interactions among the pair-rule genes themselves, creating a system of overlapping patterns. This system then directs the expression of segment polarity genes in patterns of fourteen stripes about one cell wide. The fourteen stripe patterns mark cells that will form borders between segmental compartments, and, thus, establish the body plan of the mature animal.

Pair-rule genes were first identified as members of a class of zygotic recessive mutations that affect the organization of the *Drosophila* body plan ([1](#)). Loss of function mutants in pair-rule genes contain only seven segments, and close examination of specific defects in the body plan indicates that every other segment is missing. For example, embryos with reduced *even-skipped* (*eve*) gene function contain only odd-numbered segments ("even" numbered segments are "skipped"). Other pair-rule genes lack odd-numbered segments or contiguous portions of adjacent segments that overlap every other border between two segments. These mutant phenotypes suggest that the fourteen segment body plan is established by a mechanism involving genes that specify patterns of cell fate decisions with a pair-rule periodicity.

Most pair-rule genes encode proteins with well-characterized DNA-binding motifs, and, thus, function as transcription factors in *Drosophila* development (Table [1](#)). Many of these genes are expressed in patterns that appear as seven or eight transverse stripes that encircle the cellularizing blastoderm ([2-4](#)). The stripes are evenly spaced along the anterior-posterior axis in positions coincident with regions of the body plan that are disrupted in loss of function mutants. This suggests that the pair-rule genes function for the most part within their domains of expression. Interstripe

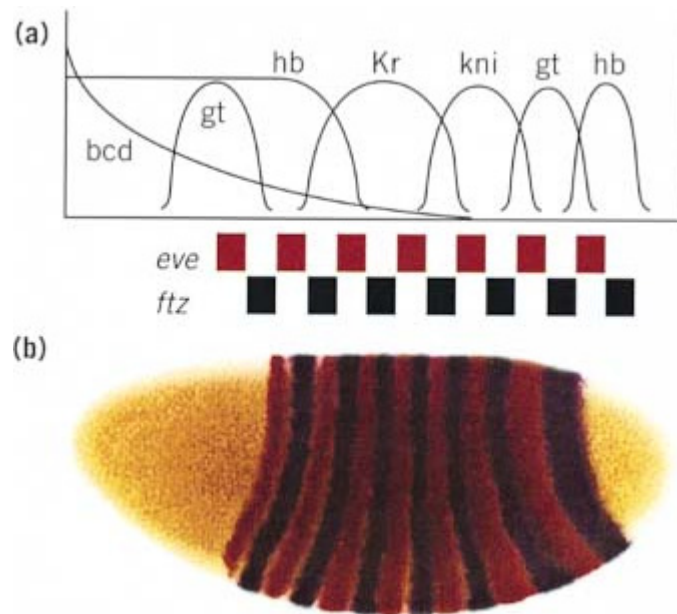
regions where the genes are not expressed are also important for normal development because near reciprocal phenotypes can be generated by ubiquitous expression of pair-rule gene products (5).

**Table 1.**

| <b>Pair-rule Gene:</b>                       | <b>Protein Motif:</b>           | <b>Reference:</b> |
|--|---------------------------------|-------------------|
| <i>hairy (h)</i>                             | helix-loop-helix                | (68)              |
| <i>even-skipped (eve)</i>                    | homeodomain                     | (4, 69)           |
| <i>fushi tarazu (ftz)</i>                    | homeodomain                     | (70)              |
| <i>runt (run)</i>                            | runt-domain                     | (71)              |
| <i>paired (prd)</i>                          | homeodomain, paired domain      | (72, 73)          |
| <i>odd-skipped (odd)</i>                     | zinc finger                     | (74)              |
| <i>odd-paired (opa)</i>                      | zinc finger                     | (75)              |
| <i>sloppy-paired (slp)</i>                   | forkhead domain                 | (76)              |
| <i>outstretched (os)</i>                     | secreted protein                | (77)              |
| <i>odd Oz (odz)</i>                          | EFG-repeats, TM tyrosine kinase | (78)              |
| <i>tenascin<sup>M</sup>(ten<sup>M</sup>)</i> | EGF-repeats, ECM protein        | (79)              |

Since the striped patterns of the pair-rule genes represent the first evidence of a metamer body plan, at least some of them must be generated *de novo*. It has been suggested that reiterated patterns such as stripes might be generated by reaction diffusion mechanisms (6-9). This model is supported by experiments in purely chemical systems that create striped patterns as the result of self-organizing processes (10). However, genetic experiments suggest that pair-rule genes are not self-organizing but are part of a hierarchy of genes that establishes polarity and segmental patterning along the anterior-posterior axis. This hierarchy is initiated by maternal effect genes that establish polarity by depositing mRNAs into the poles of the oocyte during late stages of oogenesis (11). Immediately after fertilization, the mRNAs are translated and their protein products diffuse toward the middle of the embryo creating gradients of proteins with high concentrations at the poles. These gradients direct the zygotic transcription of the gap genes, which are expressed in one or two broad domains along the anterior-posterior axis (Fig. 1). The maternal gradients and gap gene expression domains overlap, creating different combinations and concentrations of these factors according to position along the anterior-posterior axis.

**Figure 1.** (A) Spatial relationship between stripes of expression of the pair-rule genes *even-skipped (eve)* and *fushi tarazu (ftz)* and the positions of gap protein expression domains along the anterior posterior axis of *Drosophila*. Anterior is to the left. The gap genes hunchback (hb), giant (gt), Kruppel (Kr) and knirps (kni) are expressed in one or two broad domains in response to maternal gradients such as bicoid (bcd) that emanate from the poles of the embryo. (B). A double in situ hybridization experiment detects *eve* (red) and *ftz* (black) mRNAs. These genes are expressed in reciprocal patterns that establish the segmental body plan.

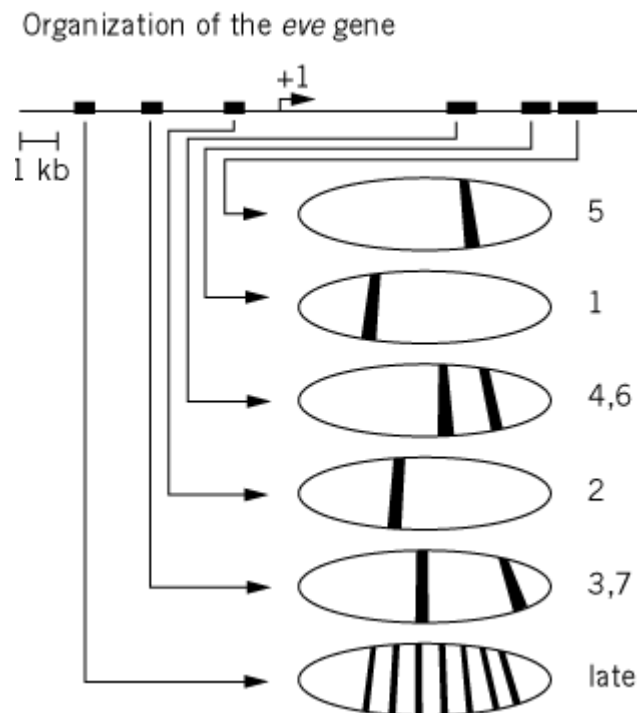


Striped patterns of the pair-rule genes are differentially disrupted in embryos lacking the functions of individual maternal effect, gap, and other pair-rule genes (12-16). These experiments are difficult to interpret, however, because removing the function of a single gene causes pattern shifts and disruptions in many other genes. Several experimental approaches have been used to unravel the mechanisms involved in pair-rule patterning. Regulatory interactions can be inferred from transgenic experiments that ubiquitously misexpress segmentation genes under the control of a heat shock inducible promoter (5, 17-20). When short heat pulses are used and the timing of pattern disruptions carefully monitored, it is possible to determine if a gene directly affects the expression of a given target or if intermediary genes must be expressed to induce the response (21, 22). Ectopic expression is also possible using enhancers that activate position-specific expression during early development (23-25). Finally, promoter truncation and P-element transformation assays (26) have identified two classes of cis-regulatory regions (enhancers) that are important for pair-rule patterning in the blastoderm. Enhancers of the first class specify the placement of individual stripes by responding primarily to gradients of maternal effect and gap proteins. Enhancers of the second class specify complete patterns of seven stripes but are active only after the first striped patterns are established. These enhancers respond to transcriptional cues provided by the pair-rule genes themselves and are thus important for mechanisms that maintain and refine initial patterns.

### 1. Regulatory Mechanisms That Control the Initiation of Individual Pair-Rule Stripes

An important breakthrough in understanding how pair-rule patterns are established arose from the analysis of *hairy* (*h*) mutations caused by translocation breakpoints in the 5' regulatory region of the gene (27). Breakpoints located near the transcription start site cause the deletion of seven segmental units (like *h* null alleles), but those further away delete only a subset of segments. Thus, the regulatory elements controlling *h* function at different positions along the anterior-posterior axis are physically separable, suggesting that individual *h* stripes are independently controlled by discrete enhancers. This hypothesis was confirmed by a series of P-element transformation experiments with *h lacZ* reporter genes, which identified discrete enhancer elements that independently control the expression of single stripes or pairs of stripes (28, 29). Similar studies led to the identification of modular enhancers in the *eve* regulatory region (30-32). Interestingly, enhancers that control stripes 2, 3, and 7 are located upstream of the *eve* coding region; those that control the other four stripes are located downstream (Fig. 2).

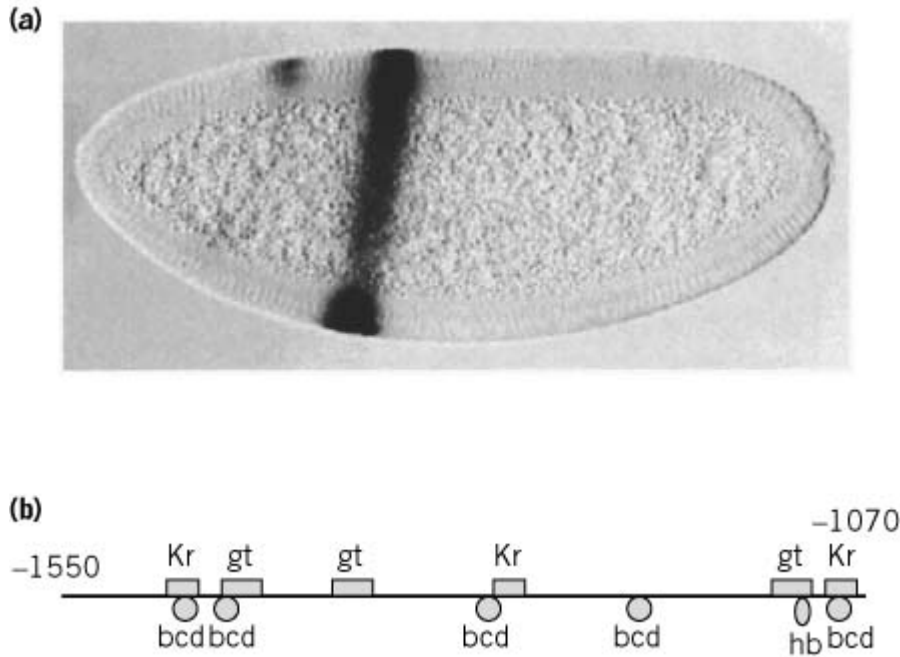
**Figure 2.** Stripe-specific enhancers in the *eve* locus. A map shows the positions of enhancers (**rectangles**) identified by promoter truncation and P-element mediated transformation experiments. Simple patterns of one or two stripes of reporter gene expression are driven five discrete enhancers. The initial patterns of seven *eve* stripes represents the sum of the activities of these enhancers. A separate enhancer directs a seven stripe pattern that appears after the initial pattern is established. This enhancer is important for maintaining and refining the pattern.



These enhancers control the positioning of individual stripes or pairs of stripes by responding to positional cues set up by the maternal effect and gap genes. The best characterized stripe enhancer controls the expression of *eve* stripe 2. Truncation analysis has narrowed this enhancer to a 480 base pair (bp) sequence that is sufficient to drive a stripe of *lacZ* reporter gene expression at the position of stripe 2(33) (Fig. 3A). Trans-acting factors involved in regulating this enhancer were identified by crossing flies containing the stripe 2-*lacZ* transgenes into embryos lacking maternal and gap-gene functions. These experiments suggest that the maternal effect gene *bicoid* (*bcd*) and the gap gene *hunchback* (*hb*) are required for activation, while the gap genes *giant* (*gt*) and *Kruppel* (*kr*) are required for setting the anterior and posterior borders of the stripe (13, 34). The expression patterns of the proteins encoded by these genes are consistent with their genetically identified functions (Fig. 1A). *bcd* and *hb* completely overlap the position of stripe 2, and *gt* and *Kr* about the anterior and posterior borders, respectively. Repressive interactions between *gt* and *Kr* are important for maintaining the spacing where stripe 2 is expressed (20, 25). All four proteins contain DNA-binding domains and have been shown to bind in vitro to multiple sites in the stripe 2 enhancer (34, 35) (Fig. 3B). Most binding sites for activator proteins are located within 50 bp of repressor sites. Mutating binding sites for *bcd* or *hb* causes a reduction in expression levels, and deleting *gt* sites causes a dramatic expansion of the stripe response into anterior regions of the embryo (33, 36, 37). These experiments suggest that the stripe 2 enhancer functions as a binary switch that controls region-specific activation by directly responding to combinations of asymmetrically localized regulatory factors.

**Figure 3. (A)** *lac Z* reporter gene expression directed by the 480 bp *eve* stripe 2 enhancer. Genetic experiments suggest that this enhancer is activated by *bicoid* (*bcd*) and *hunchback* (*hb*). The anterior and posterior borders of this stripe are

set by repression involving the gap proteins giant (gt) and Kruppel (Kr), respectively (see Fig. 1). **(B)** Positions of binding sites for proteins that regulate this enhancer. Activator sites are shown below the line, with repressor sites above. These sites enable the enhancer to make on/off decisions based on the combinations and concentrations of these proteins in a given nucleus.



Several other stripe-specific enhancers have been characterized in detail (38-40). These analyses suggest a general model for the initial establishment of individual stripes (41). In this model, factors involved in activating individual stripes are broadly distributed, and stripe borders are set by repressive interactions mediated by gradients of gap proteins. Activations and repression events involve direct binding to closely linked sites within each enhancer. The close linkage of repressor to activator sites suggests that the repressors function over short distances to interfere with the binding or activity of activator proteins. Consistent with this hypothesis, increasing the spacing between activator and repressor sites to an interval greater than ~100 bp can prevent gap protein-mediated repression from occurring (42, 43).

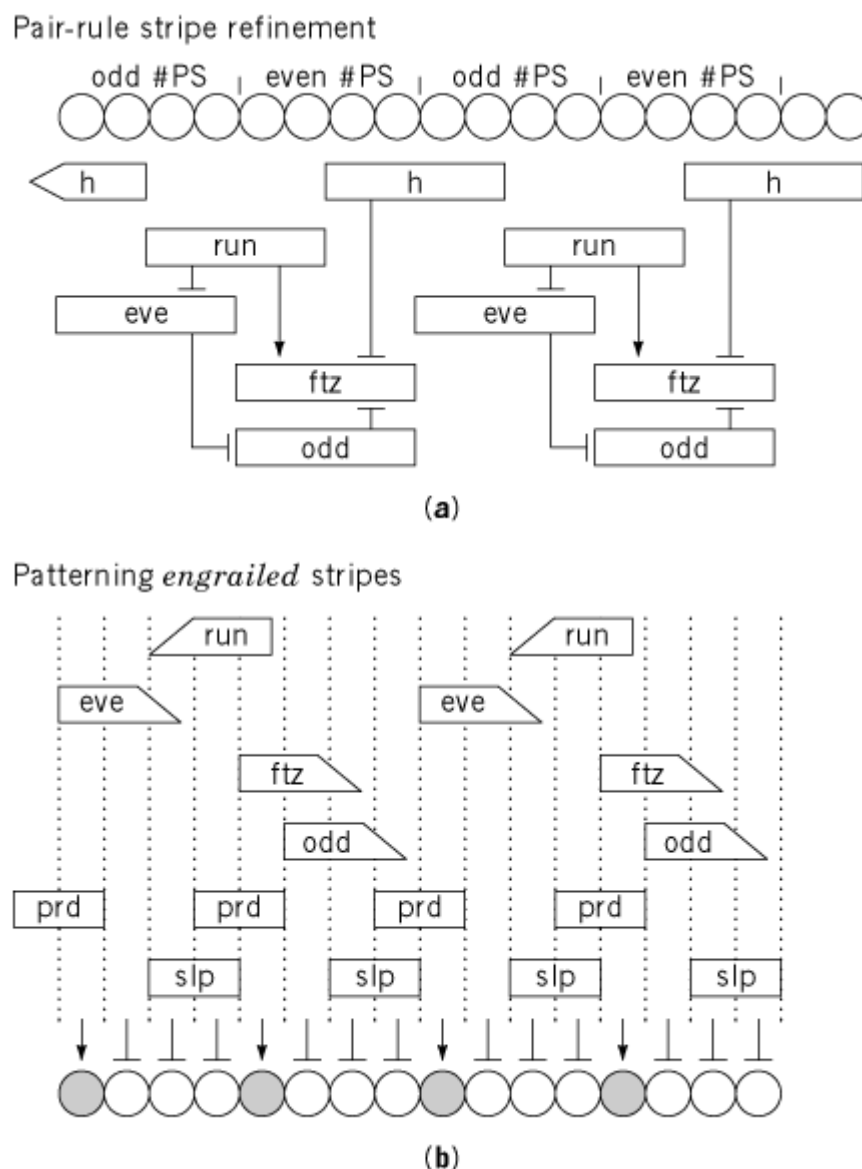
Short-range repression is also important for ensuring that individual enhancers function independently in the context of complex promoters. This has been tested by transgenic experiments using reporter genes that contain the *eve* stripe 2 enhancer and another 500 bp enhancer that regulates stripes 3 and 7 (40). In the wild-type gene, these two enhancers are separated by a ~1,700 bp sequence. When this sequence is deleted so that the enhancers are juxtaposed, there is a severe disruption of the pattern driven by the reporter gene (44). Inserting shorter spacer sequences (160 bp and 300 bp) between the enhancers restores the correct expression pattern. These results suggest that transacting factors bound to one enhancer interfere with the activity of other enhancers and that spacing between enhancers prevents such interface in the wild-type gene.

## 2. Regulatory Mechanisms That Maintain and Refine Striped Patterns

Because different mechanisms control the establishment of individual *eve* and *h* stripes, they are expressed in a specific temporal order while the embryo is going through the process of cellularization. Each enhancer contributes to a final pattern of seven stripes that are each approximately four to five nucleus diameters wide (Figure 4A). By the end of cellularization, these initial patterns are refined until individual stripes are two or three cells wide. During this process, some seven stripe patterns are replaced by fourteen stripes, either by splitting the initial stripes or by

activating new stripes in interstripe regions. Precise overlaps between specific genes are established so that each row of nuclei within a segmental unit contains a unique combination of pair-rule proteins. Genetic experiments suggest that positive interactions among the pair-rule genes maintain the stripes during this process, while negative interactions control spatial refinement. For example, there is a premature loss of *eve* stripes in *eve* and *paired* (*prd*) mutants, suggesting roles for these genes in maintaining the striped pattern (45, 46). In contrast, there is a significant posterior expansion of *eve* stripes in *run* mutants, suggesting that repression by *run* is important for refining the stripes (13, 31). Ubiquitous expression of *run* via heat shock causes a rapid repression of *eve*, consistent with this hypothesis (22). In contrast, *ftz* stripes are prematurely lost in *ftz* and *run* mutants, and expanded in *h* mutants, suggesting that these genes maintain and refine the *ftz* pattern (12, 47, 48).

**Figure 4.** (A) Interactions among pair-rule genes that maintain and refine initial striped patterns. Four parasegmental units are shown. **Arrows** represent positive interactions, **blunted lines** negative. The process of maintenance and refinement narrows stripes from four to five nuclei wide to two to three wide. (B) Spatial relationship between refined pair-rule patterns and the expression pattern of the segment polarity gene. *en*. *en* is expressed in cells that contain specific combinations of pair-rule proteins. B is colinear with A.





The molecular mechanisms involved in these interactions are only beginning to be understood. Promoter truncation and P-element transformation experiments have identified enhancers in *eve* and *ftz* that direct reporter gene expression in patterns of seven stripes (31, 32, 49). These stripes appear late in the process of cellularization, persisting through the stages of gastrulation and germ band elongation. This timing is consistent with a role for these enhancers in maintaining and refining initial striped patterns. An example is a 300 bp *eve* enhancer that directs a late pattern of seven stripes coincident with the endogenous pattern (46, 50). This enhancer is inactive in embryos that lack *eve* or *prd* function, suggesting that a combination of autoregulation and cross regulation by *prd* activates transcription of this enhancer. The enhancer contains binding sites for both proteins; mutating these sites causes a loss of expression, suggesting that direct DNA-binding is required for activation (46, 50).

Another autoregulatory element involved in pair-rule refinement is 430 bp *ftz* enhancer (AE) that contains multiple *ftz* binding sites; when these sites are mutated, the enhancer is inactivated (51). Further experiments strongly indicate that the *in vivo* activation of this enhancer involves direct binding by *ftz*. These experiments take advantage of the fact that the *ftz* protein, which contains a glutamine residue at position 50 of the homeodomain, binds *in vitro* with high affinity to a specific class of binding sites (CCATTA) (52, 53). Other homeodomain proteins, such as *bed*, contain lysine at this position and prefer a different sequence (GGATTA). Changing a single CCATTA binding site in the *ftz* AE to GGATTA causes a significant loss of AE activity (51). However, activity is restored *in trans* by a *ftz* rescue construct in which the glutamine at position 50 is replaced with a lysine residue. This experiment confirms the importance of this residue in specific DNA recognition by homeodomain proteins and demonstrates conclusively that direct DNA-binding by *ftz* is required for the activity of this enhancer.

It is worth noting, however, that the 430 bp AE represents only a small part of a complex 2,600 bp upstream element (UE) that is important for *ftz* stripe maintenance (54-57). A separate enhancer, the zebra element (ZE), also direct reporter gene expression in late cellularizing embryos (49). Together, these enhancers integrate positive inputs provided by *ftz*, the *ftz*-cofactor Ftz-F1, and *run*, as well as negative inputs mediated by *h* that refine the stripes from the posterior (48, 58, 59). Binding sites for gap proteins and other ubiquitously distributed factors such as *tramtrack* (*ttk*), have also been identified in these regions, but their exact roles in *ftz* patterning are still unclear (56, 60).

### 3. Functions of the Striped Patterns of the Pair-Rules Genes

Precise overlapping patterns of the pair-rule genes are established by the end of cellularization. These patterns comprise a coordinated system that controls the expression of the segment polarity class of segmentation genes. An example is the segment polarity gene *engrailed* (*en*), which is expressed in fourteen stripes one cell wide (61, 62). Cells that express *en* will remain segregated from neighboring cells lacking *en* throughout development. This segregation establishes a border between the anterior and posterior compartments within each parasegmental unit. Genetic and misexpression experiments suggest that odd and even-numbered stripes within the *en* pattern are regulated by different mechanisms involving pair-rule genes (21-23, 63-66) (Fig 4). Odd-numbered stripes are activated in cells that contain *eve* and *prd* protein, and the anterior and posterior borders of these stripes are set by *sloppy-paired* (*slp*) and *run* respectively. Even numbered stripes are activated in cells than contain *ftz* and *prd* protein, and the anterior and posterior borders of these stripes are probably set by *slp* and *odd-skipped* (*odd*), respectively. Very little is known, however, about the *cis* components that control the expression of either set of *en* stripes or any other segment polarity gene.

In summary, the pair-rule genes make up an integrated system that transduces asymmetrically localized positional information into fourteen stripe patterns that establish the segmental body plan of *Drosophila*. Significant progress has been made in unraveling certain aspects of the mechanisms involved, but the complexity of the system suggests that a full understanding will not possible using only a reductionist approach. Recent attempts have been made to use computer modeling to integrate

individual components of the pair-rule system into a gene circuit method (67). This method has successfully generated a complete pattern of seven stripes at the *eve* position using quantitative data representing the expression patterns of the maternal and gap proteins. In the future, this type of analysis will be important for synthesizing molecular data and predicting mechanisms that can then be tested experimentally.

## Bibliography

1. C. Nusslein-Volhard and W. Wieschaus (1980) *Nature* **287**, 795–801.
2. S. B. Carroll and M. P. Scott (1985) *Cell* **43**, 47–57.
3. E. Hafen, A. Kuroiwa, and W. Gehring (1984) *Cell* **37**, 833–841.
4. P. M. Macdonald, P. Ingham, and G. Struhl (1986) *Cell* **47**, 721–34.
5. G. Struhl (1985) *Nature* **318**, 677–80.
6. A. M. Turing, (1990) *Bull. Math. Biol.* **52**, 153–97.
7. S. A. Newman (1993) *Bioessays* **15**, 277–83.
8. T. C. Lacalli (1990) *J. Theor. Biol.* **144**, 171–94.
9. A. Hunding, S. A. Kauffman, and B. C. Goodwin (1990) *J. Theor. Biol.* **145**, 369–84.
10. Q. Ouyang and H. L. Swinney (1991) *Nature* **352**, 610–612.
11. D. St Johnston and C. Nusslein-Volhard (1992) *Cell* **68**, 201–19.
12. S. B. Carroll and M. P. Scott (1986) *Cell* **45**, 113–26.
13. M. Frasch and M. Levine (1987) *Genes. Dev.* **1**, 981–95.
14. M. Klingler and J. P. Gergen (1993) *Mech. Dev.* **43**, 3–19.
15. K. L. Hooper, S. M. Parkhurst, and D. Ish-Horowicz (1989) *Development* **107**, 489–504.
16. T. Gutjahr, E. Frei, and M. Noll (1993) *Development* **117**, 609–23.
17. G. Struhl (1989) *Ciba. Found. Symp.* **144**, 65–86.
18. E. Eldon and V. Pirrotta (1991) *Development* **111**, 367–378.
19. H. M. Krause, R. Klemenz, and W. J. Gehring (1988) *Genes. Dev.* **2**, 1021–36.
20. R. Kraut and M. Levine (1991) *Development* **111**, 611–21.
21. A. S. Manoukian and H. M. Krause (1992) *Genes. Dev.* **6**, 1740–51.
22. A. S. Manoukian and H. M. Krause (1993) *Development* **118**, 785–96.
23. M. Fujioka, J. B. Jaynes and T. Goto (1995) *Development* **121**, 4371–82.
24. D. Kosman and S. Small (1997) *Development* **124**, 1343–54.
25. X. Wu, R. Vakani and S. Small (1998) *Development* **125**, 3765–74.
26. G. M. Rubin and A. C. Spradling (1982) *Science* **218**, 348–53.
27. K. Howard, P. Ingham and C. Rushlow (1988) *Genes. Dev.* **2**, 1037–46.
28. G. Riddihough and D. Ish-Horowicz (1991) *Genes. Dev.* **5**, 840–54.
29. K. R. Howard and G. Struhl (1990) *Development* **110**, 1223–31.
30. M. Fujioka, Y. Emi-Sarker, G. Yusibova, T. Goto, and J. Jaynes (1999) *Development* **126**, 2527–38.
31. T. Goto, P. Macdonald, and T. Maniatis (1989) *Cell* **57**, 413–22.
32. K. Harding, T. Hoey, R. Warrior, and M. Levine (1989) *Embo. J.* **8**, 1205–12.
33. S. Small, A. Blair, and M. Levine (1992) *Embo. J.* **11**, 4047–57.
34. S. Small, R. Kraut, T. Hoey, R. Warrior and M. Levine (1991) *Genes. Dev.* **5**, 827–39.
35. D. Stanojevic, T. Hoey and M. Levine (1989) *Nature* **341**, 331–5.
36. D. Stanojevic, S. Small, and M. Levine (1991) *Science* **254**, 1385–7.
37. D. N. Arnosti, S. Barolo, M. Levine, and S. Small (1996) *Development* **122**, 205–14.
38. M. J. Pankratz, E. Seifert, N. Gerwin, B. Billi, U. Nauber, and H. Jackle (1990) *Cell* **61**, 309–17.

39. J. A. Langeland and S. B. Carroll (1993) *Development* **117**, 585–96.
40. S. Small, A. Blair, and M. Levine (1996) *Dev. Biol.* **175**, 314–24.
41. S. Small, and M. Levin (1991) *Curr Opin Genet Dev* **1**, 255–60.
42. A. D. Johnson (1995) *Cell* **81**, 655–8.
43. S. Gray, and M. Levine (1996) *Curr Opin Cell Biol* **8**, 358–72.
44. S. Small, D. N. Arnosti, and M. Levine (1993) *Development* **119** 762–72.
45. M. Frash, R. Warrior, J. Tugwood, and M. Levine (1988) *Genes Dev* **2**, 1824–38.
46. M. Fujioka, P. Miskiewicz, L. Raj, A. A. Gullledge, M. Weir, and T. Goto (1996) *Development* **122**, 2697–707.
47. K. Howard, and P. Ingham (1986) *Cell* **44**, 949–57.
48. C. Tsai, and P. Gergen (1995) *Development* **121**, 453–62.
49. Y. Hiromi, A. Kuroiwa, and W. J. Gehring (1985) *Cell* **43**, 603–13.
50. J. Jiang, T. Hoey, and M. Levine (1991) *Genes. Dev.* **5**, 265–77.
51. A. F. Schier, and W. J. Gehring (1992) *Nature* **356**, 804–7.
52. A. Percival-Smith, M. Muller, M. Affolter, and W. J. Gehring (1990) *Embo. J.* **9**, 3967–74.
53. J. Treisman, P. Gonczy, M. Vashishtha, E. Harris, and C. Desplan (1989) *Cell* **59**, 553–62.
54. Y. Hiromi, and W. J. Gehring (1987) *Cell* **50**, 963–74.
55. L. Pick, A. Schier, M. Affolter, T. Schmidt-Glenewinkel, and W. J. Gehring (1990) *Genes. Dev.* **4**, 1224–39.
56. W. Han, Y. Yu, N. Atlan, and L. Pick (1993) *Mol Cell Biol* **13**, 5549–59.
57. A. F. Schier, and Gehring W. J. (1993) *Embo. J.* **12**, 1111–9.
58. A. Guichet, J. W. Copeland, and M. Erdelyi et al. (1997) *Nature* **385**, 548–52.
59. Y. Yu, W. Li, and K. Su et al. (1997) *Nature* **385**, 552–5.
60. C. R. Dearolf, J. Topol, and C. S. Parker (1989) *Genes. Dev.* **3**, 384–98.
61. T. Kornberg, I. Siden, P. O’Farrell, and M. Simon (1985) *Cell* **40**, 45–53.
62. P. H. O’Farrell, C. Desplan, and S. DiNardo et al. (1985) *Cold Spring Harb. Symp. Quant. Biol.* **50**, 235–42.
63. D. E. Coulter, and E. Wieschaus (1988) *Genes. Dev.* **2**, 1812–23.
64. S. DiNardo, and P. H. O’Farrell (1987) *Genes. Dev.* **1**, 1212–25.
65. M. P. Weir, B. A. Edgar, T. Kornberg, and G. Schubiger (1988) *Genes Dev* **2**, 1194–203.
66. K. M. Cadigan, U. Grossniklaus, and W. J. Gehring (1994) *Genes. Dev.* **8**, 899–913.
67. J. Reintz, and D. H. Sharp (1995) *Mech. Dev.* **49**, 133–58.
68. C. A. Rushlow, A. Hogan, S. M. Pinchin, K. M. Howe, M. Lardelli, and D. Ish-Horowicz (1989) *Embo. J.* **8**, 3095–103.
69. M. Frash, T. Hoey, C. Rushlow, H. Doyle, and M. Levine (1987) *Embo. J.* **6**, 749–59.
70. A. Laughon, S. B. Carroll, F. A. Storfer, P. D. Riley, and M. P. Scott (1985) *Cold Spring Harb. Symp. Quant. Biol.* **50**, 253–62.
71. M. A. Kania, A. S. Bonner, J. B. Duffy, and J. P. Gergen (1990) *Genes. Dev.* **4**, 1701–13.
72. G. Frigerio, M. Burri, D. Bopp, S. Baumgartner, and M. Noll (1986) *Cell* **47**, 735–46.
73. J. Treisman, E. Harris, and C. Desplan (1991) *Genes. Dev.* **5**, 594–604.
74. D. E. Coulter, E. A. Swaykus, M. A. Beran-Koehn, D. Goldberg, E. Wieschaus, and P. Schedl (1990) *Embo. J.* **9**, 3795–804.
75. M. J. Benedyk, J. R. Mullen, and S. DiNardo (1994) *Genes. Dev.* **8**, 105–17.
76. U. Grossniklaus, R. K. Pearson, and W. J. Gehring (1992). *Genes. Dev.* **6**, 1030–51.
77. A. Levine, A. Bashan-Ahrend, O. Budai-Hadrian, D. Gartenberg, S. Menasherow, and R. Wides (1994) *Cell* **77**, 587–98.

78. S. Baumgartner, D. Martin, C. Hagios, and R. Chiquet-Ehrismann (1994) *Embo. J.* **13**, 3728–40.  
 79. D. Harrison, P. McCoon, R. Binari, M. Gilman, and N. Perrimon (1998) *Genes. Dev.* **12**, 3252–63.

### Suggestions for Further Reading

80. P. Lawrence (1992) *The making of a fly*, Blackwell Scientific Publications, Cambridge, MA.  
 81. M. Pankratz and H. Jackle (1993) "Blastoderm Segmentation. The Development of *Drosophila melanogaster*" (M. Bate and A. Martinez-Arias, eds.), Cold Spring Harbor Press.  
 82. S. Small (1997) "Mechanisms of segmental pattern formation in *Drosophila melanogaster*. *Progress in Development Biology*". (J. Collier, ed.) Oxford and IBH Publishing, pp. 137–178.  
 83. D. St Johnston and C. Nusslein-Volhard (1992). The origin of pattern and polarity in the *Drosophila* embryo. *Cell*: 201–19.

## Palindrome

A palindromic word or text reads the same in both directions. Strictly speaking, this definition applies only to the same line of the text. The term “palindrome” is traditionally applied to double-stranded **DNA** (which is “two lines”) when the **nucleotide sequence** is read forward from one strand and backward from the other, complementary strand. The reading proceeds in the same direction in terms of chemical polarity, 5' to 3', because DNA is an antiparallel duplex. A more appropriate name for such double-stranded palindromes that is also frequently used is a “sequence or element with complementary symmetry.”

The palindromes in the double-stranded sense as discussed above are not only built of two identical sequences, but they also have complementary symmetry (twofold rotational symmetry). One example, the recognition site for the Bam HI [restriction enzyme](#), is shown in Fig. 1. Such a peculiar structure is typical of restriction sites. It guarantees that both strands—not just one—will be cut by the enzyme, either simultaneously or after a second approach by the enzyme. Nucleotide sequences of various lengths with complementary symmetry are rather common. For example, they comprise the **consensus sequences** of many **repressor-binding sites**, or **operators** (1).

**Figure 1.** A nucleotide sequence with complementary symmetry (a) and a double-stranded palindrome (b)—the element with twofold rotational symmetry. Hydrogen bonds between the strands in (b) are not shown.

5' -GGATCC-3'

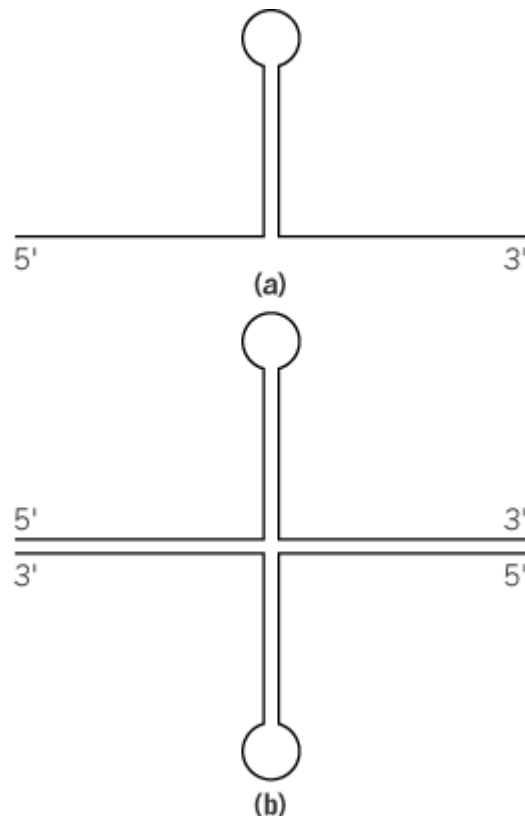
(a)

5' -GGATCC-3'  
 3'-CCATGG-5'

(b)

The condition of identity and complementarity of the two strands of the palindrome also makes the halves of each strand complementary to one another, with the potential for each strand to fold back on itself, forming a hairpin structure. In a palindromic DNA duplex, such hairpins could be made in both strands, so that a **cruciform** structure is formed (2) (see Fig. 2). The loops of the cruciforms generally could be any size and not necessarily complementary symmetric as the sequences in the stems of the cruciform. Another obvious property of palindromes is the identity of their double-stranded halves, except for their being placed in opposite polarity, like in **inverted repeats**, but without interruption between the repeats. (Various types of repeats are compared in the scheme presented in **Tandem Repeats**.) Generally speaking, any head-to-head (or tail-to-tail) fusion of identical fragments of a double-stranded DNA would generate a palindrome. For example, very long palindromes are formed in the process of **gene amplification** (3), in the size range of many kilobases. Complementary symmetries generally separated by rather long interruptions are characteristic of many **transposable elements** (4). A whole family of transposable elements also exists: the **foldback element** transposons, which are about 1 kb in size. These transposons are essentially made of inverted repeats with a very small middle loop region—almost a palindrome.

**Figure 2.** Schematic presentation of a hairpin structure in a single-stranded DNA or RNA (a) and of a cruciform structure (b). The double lines correspond to antiparallel complementary strands in the duplex parts of the structures.



Palindromes that read the same in both directions of *the same strand*, ie, proper (or true), palindromes also are known. These sequences possess a mirror symmetry, for example, ATGAGTA. Such elements are frequently encountered in intervening sequences, or **introns** (5), and may be interpreted as a device to avoid complementary folding of the sequence. Indeed, the absence of complementary sequence elements in the mirror-symmetric sequences of the true palindromes like AT $\frac{1}{2}$ TA and TG $\frac{1}{2}$ GT causes the sequences to be located in the single-stranded regions of the RNA secondary structure.

Another class of mirror-symmetry palindromes are the polypurine and polypyrimidine stretches in the H-form structure (6, 7). This partially triplex structure is formed in a normal DNA duplex at lowered pH or under torsional stress. The locally separated polypyrimidine strand folds back and binds to the duplex by Hoogsteen [hydrogen bonds](#). The sequence of this third strand can be incorporated into the structure only if it is mirror-symmetric to the sequence of the polypyrimidine strand of the Watson-Crick duplex.

### Bibliography

1. P. P. Papp, D. K. Chattoraj, and T. D. Schneider (1993) *J. Mol. Biol.* **233**, 219–230.
2. Y. Timsit and D. Moras (1996) *Q. Rev. Biophys.* **29**, 279–307.
3. J. L. Hamlin *et al.* (1991) *Prog. Nucl. Acid Res. Mol. Biol.* **41**, 203–239.
4. *Mobile DNA* (1989) (D. E. Berg and M. M. Howe, eds.) ASM, Washington, D.C.
5. J. S. Beckmann, V. Brendel, and E. N. Trifonov (1986) *J. Biomol. Struct. Dynam.* **4**, 391–400.
6. S. M. Mirkin *et al.* (1987) *Nature* **330**, 495–497.
7. M. D. Frank-Kamenetskii and S. M. Mirkin (1995) *Annu. Rev. Bioch.* **64**, 65–95.

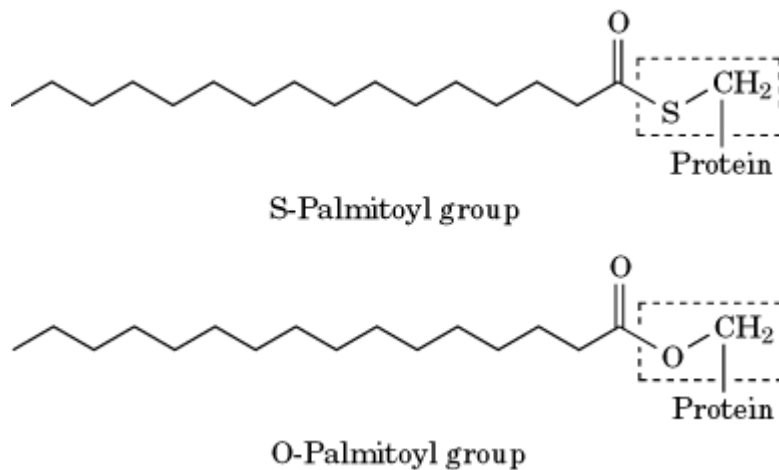
### Suggestion for Further Reading

8. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.

## Palmitoylation

Palmitoylation is a [post-translational modification](#) process in which eukaryotic and viral proteins become ester- or thioester-linked to the 16-carbon saturated [fatty acid](#), palmitic acid. Palmitoylation occurs on the [thiol group](#) of [cysteine](#) residues (S-palmitoylation) or on the hydroxyl group of [serine](#) and [threonine](#) residues (O-palmitoylation). S-palmitoylation (see Fig. 1) of cysteine residues close to the N- or C-terminus provides a [membrane anchor](#) for some cytosolic proteins, and is important for their correct localization within the cell. However, S-palmitoylation also occurs on polypeptide-anchored and integral [membrane proteins](#), and its role in these proteins is uncertain.

**Figure 1.** Modification of a cysteine or serine residue by a palmitoyl group. An internal cysteine residue in the protein is thioester-linked to the 16-carbon saturated fatty acid, palmitic acid to give an S-palmitoyl group. Palmitic acid can also be ester-linked to serine (or threonine) residues to give an O-palmitoyl group. The cysteine and serine residue in each figure is indicated by a dotted line.



In contrast to the other lipid modifications of proteins (see [Membrane Anchors](#)), palmitoylation has proven relatively difficult to study. Even though palmitoylation was discovered more than 20 years ago and occurs on at least 50 different types of proteins, this process is still poorly understood at the biochemical level [for a comprehensive description of the early work in this area, see the review articles by Schmidt (1989) and Towler et al (1988)]. Apart from detailed studies on S-palmitoylation of cysteine residues in [signal transduction](#) components (eg, Src family **tyrosine kinases**, **G-protein**  $\gamma$  subunits, Ras proteins, etc; see below), many of the initial reports of protein palmitoylation have not been substantiated in the last decade, and their functional significance is difficult to assess. There is essentially no recent information on O-palmitoylation of proteins. The lack of recent progress in these areas is due to several reasons.

1. For most of the “palmitoylated” proteins, the evidence for palmitoylation is based solely on the incorporation of  $^3\text{H}$ -labeled palmitate into cellular proteins followed by **immunoprecipitation**, **SDS-PAGE**, and [autoradiography](#). It is usually possible to distinguish ester, or O-palmitoylation of serine and threonine from thioester or S-palmitoylation of cysteine in the  $^3\text{H}$ -labeled proteins by selective chemical deacylation procedures. However, more detailed analysis of  $^3\text{H}$ -labeled palmitoylated peptides is generally not practicable, because their unpredictable chromatographic behavior makes them very difficult to purify. This technical problem has thwarted attempts to identify the precise location of the incorporated  $^3\text{H}$ -label in most of the palmitoylated proteins. Indeed, for several “palmitoylated” proteins, it was subsequently shown that the incorporated palmitate is part of a completely different structure, the [GPI anchor](#).
2. The  $^3\text{H}$ -labeled palmitate can be converted to other fatty acids before incorporation into some of the proteins. This suggests that the acyltransferase responsible for acylating the protein may be non selective and/or that different cells contain variable amounts of distinct acyltransferases with different specificities.
3. Protein S-palmitoyltransferase activities that can catalyze the transfer of a palmitoyl group from palmitoyl CoA to cysteine residues in protein or model peptide substrates have been identified in mammalian membrane fractions (1–5). Removal of the palmitoyl group by thioesterase activities has also been described. However, none of these activities has been thoroughly characterized at the molecular level, and it is difficult to tell if they are responsible for palmitoylation/depalmitoylation *in vivo*. Non enzymatic palmitoylation (autoacylation) of cysteine residues by palmitoyl CoA was also reported for several proteins, and this may account for some of the difficulties encountered in isolating and characterizing the protein S-palmitoyltransferase(s).
4. The criteria that determine which particular cysteine residue(s) in a protein will become palmitoylated are also uncertain. The presence of other lipid anchors (eg, prenyl or myristoyl groups; see below) on neighboring residues will facilitate cysteine palmitoylation presumably as

a result of transient anchoring to membranes. These observations suggest that localization of “exposed” cysteine residue close to the bilayer surface, rather than a well-defined sequence motif, is the determinant for palmitoylation (4). S-palmitoylation of cysteine residues in polypeptide-anchored and integral [membrane proteins](#) could simply be the accidental result of proximity to the bilayer surface, where palmitoyl CoA and other molecular species of acylCoA are likely to be abundant.

As a consequence of the experimental problems identified above, there is relatively little information on the functional significance of palmitoylation for most proteins. However, [site-directed mutagenesis](#) studies, on a small number of proteins that are involved in signal transduction, demonstrated that a palmitoyl group can act as a membrane anchor and may play a role in regulating distribution of proteins between membranes and the cytoplasm as well as their association with lipid rafts (6) (for a general discussion of factors that can affect membrane affinity of lipid-anchored proteins, see [Membrane Anchors](#)). In some G protein  $\gamma$  subunits and *src*-related tyrosine kinases, it has been shown that palmitoylation of cysteines near the N-terminus is dependent on prior [myristoylation](#) of the N-terminal glycine (7-9). Palmitoylation of  $\gamma$  subunits can also be increased by receptor activation (10, 11). Similarly, palmitoylation of cysteine residues near the C-terminus of some Ras proteins is dependent on prior farnesylation (see [Prenylation](#)). Because farnesyl and myristoyl groups are both relatively weak anchors, the increase in membrane affinity resulting from palmitoylation has a major effect on the amount of these proteins bound to the membrane (12, 13). Furthermore, biophysical studies indicate that addition of palmitate not only increases binding affinity, but also increases the stability of membrane anchoring by reducing the rate of dissociation from the membrane (14, 15). Palmitoylation also occurs on cysteine residues in polypeptide-anchored proteins (eg, the [transferrin](#) receptor and CD4; see Table 1) and integral proteins with multiple transmembrane regions (eg, [rhodopsin](#) and some types of **adrenergic** receptor). As expected, preventing palmitoylation of these proteins by mutagenesis of the appropriate cysteine residues has no discernible effect on membrane anchoring. However, cysteine mutagenesis often has no, or relatively subtle, effects on the functional properties of the protein either, and the biological significance of palmitoylation remains uncertain in these cases.

**Table 1. Examples of Palmitoylated Proteins**

|   |
|---|
| G protein $\gamma$ subunits   |
| Src family tyrosine kinases (e.g. p59 <sup>lyn</sup> , p56 <sup>lck</sup> etc.)                                     |
| GAP-43 (neuromodulin)   |
| H- or N-Ras proteins  |
| G protein-coupled receptors (e.g. rhodopsin, $\gamma_{2A}$ , $\gamma_{2B}$ -adrenergic receptors etc.) <sup>a</sup> |
| Transferrin receptor <sup>a</sup>   |
| CD4 <sup>a</sup>  |
| Viral envelope proteins <sup>a</sup>  |
| Myelin proteolipid protein <sup>a, b</sup>  |
| Mucus glycoproteins <sup>b</sup>  |

<sup>a</sup> Integral or polypeptide-anchored proteins.

<sup>b</sup> O-palmitoylation; all other examples are S-palmitoylation.



## Bibliography

1. L. Berthiaume and M. D. Resh (1995) *J. Biol. Chem.* **270**, 22399–22405.
2. L. Liu, T. Dudler and M. H. Gelb (1996) *J. Biol. Chem.* **271**, 23269–2376.
3. J. T. Dunphy, W. K. Greentree, C. L. Manahan, and M. E. Linder (1996) *J. Biol. Chem.* **271**, 7154–7159.
4. J. T. Dunphy, H. Schroeder, R. Leventis, W. Greentree, J. K. Knudsen, J. R. Silvius, and M. E. Linder (2000) *Biochim. Biophys. Acta* **1485**, 185–198.
5. J. A. Duncan and A. G. Gilman (1996) *J. Biol. Chem.* **271**, 23594–23600.
6. S. Moffett, D. A. Brown, and M. E. Linder (2000) *J. Biol. Chem.* **275**, 2191–2198.
7. M. Y. Degtyarev, A. M. Spiegel, and T. L. Z. Jones (1994) *J. Biol. Chem.* **269**, 30898–30903.
8. A. M. Shenoy-Scaria, L. K. Timson Gauen, J. Kwong, A. S. Shaw, and D. M. Lublin (1993) *Mol. Cell. Biol.* **13**, 6385–6392.
9. L. Alland, S. M. Peseckis, R. E. Atherton, L. Berthiaume, and M. D. Resh (1994) *J. Biol. Chem.* **269**, 16701–16705.
10. S. M. Mumby, C. Kleus, and A. G. Gilman (1994) *Proc. Natl. Acad. Sci. U.S.A.* **91**, 2800–2804.
11. M. Y. Degtyarev, A. M. Spiegel, and T. L. Z. Jones (1993) *J. Biol. Chem.* **268**, 23769–23772.
12. K. Cadwallader, H. Paterson, S. G. MacDonald, and J. F. Hancock (1994) *Mol. Cell. Biol.* **14**, 4722–4730.
13. M. H. Gelb (1997) *Science* **275**, 1750–1751.
14. S. Shahinian and J. R. Silvius (1995) *Biochemistry* **34**, 3813–3822.
15. G. Milligan, M. Parenti, and A. I. Magee (1995) *TIBS* **20**, 181–186.

## Suggestions for Further Reading

16. D. A. Towler, J. I. Gordon, S. P. Adams, and L. Glaser (1988) The biology and enzymology of eukaryotic protein acylation. *Annu. Rev. Biochem.* **57**, 69–99.
17. M. F. G. Schmidt (1989) Fatty acylation of proteins. *Biochim. Biophys. Acta* **988**, 411–426.
18. P. J. Casey and J. E. Buss (1995) Lipid Modifications of Proteins, *Meth. Enzymol.* **250**.
19. R. S. Bhatnagar and J. I. Gordon (1997) Understanding covalent modification of proteins by lipid; Where cell biology and biophysics mingle. *Trends Cell Biol.* **7**, 14–20.
20. J. T. Dunphy and M. E. Linder (1998) Signalling function of proteins palmitoylation. *Biochim. Biophys. Acta* **1436**, 245–261.

## Papain

Papain is a [thiol proteinase](#) that is obtained from the latex of the papaya tree ([1](#)) (Fig. [1](#)). It has an essential [cysteine](#) residue at its [active site](#) that must be in the reduced form in order to be active. For this reason, it is often used in the presence of a reducing agent, such as [b-mercaptoethanol](#) or **dithiothreitol**. Papain has broad specificity and will hydrolyze peptide bonds adjacent to almost any amino acid residue, although it shows a **trypsin**-like preference for [arginine](#) and [lysine](#) residues. Its activity is optimal around pH 6. It is inactivated by reagents that react with [thiol groups](#), such as **iodoacetate** or **iodoacetamide**, as well as silver and mercuric ions and *p*-chloromercuribenzoate.

Other inhibitors include [leupeptin](#), a peptide aldehyde that is a general thiol proteinase inhibitor, and E-64, an epoxide compound isolated from *Aspergillus japonicus* (Fig. 2). Leupeptin also inhibits some serine proteinases, whereas E-64 is more specific and has a strong affinity only for thiol proteinases. The epoxide group of the latter forms a thioether with the active site thiol group of the enzyme. It has little reactivity for thiol groups not in active sites.

**Figure 1.** The three-dimensional structure of the thiol proteinase papain. Only the polypeptide backbone is indicated schematically as a ribbon, with arrows for b-strands and coils for a-helices.



**Figure 2.** Chemical structure of E-64, an inhibitor of some thiol proteinases. Its chemical name is L- *trans*-epoxysuccinyl-leucylamide-(4-guanidino)-butane. The epoxy group forms a thioether linkage with the active-site thiol group of papain or cathepsin B.



interaction with cellular E6-AP protein, and they subsequently cause the **ubiquitin-mediated protein degradation** of the p53 protein. The consequent deregulation of cellular growth control by E6 and E7 genes is considered to be the molecular basis of carcinogenesis by the high-risk HPVs.

#### Suggestions for Further Reading

H. zur Hausen (1996) Papillomavirus infections-a major cause of human cancers. *Biochim. Biophys. Acta* **1288**, F55–F78.

P. M. Howley (1996) "Papillomavirinae: The Viruses and Their Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2045–2076.

### Paralogous Genes

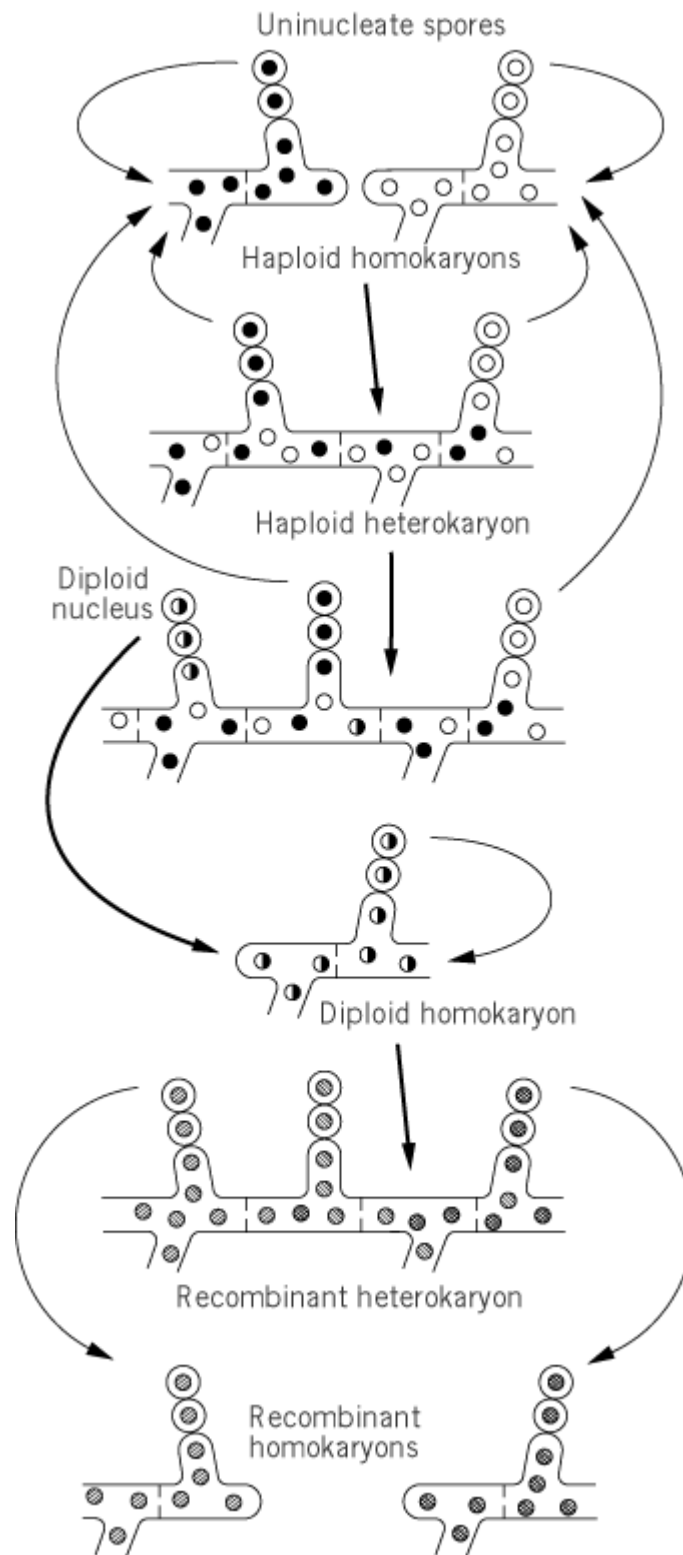
Paralogous genes are two or more **genes** that emerged during [evolution](#) from a common ancestral gene by [gene duplication](#) within a species. Usually, they acquire different functions; well-known examples are the a and b-[hemoglobin](#) polypeptide chains and that of [myoglobin](#). They contrast with **orthologous** genes, which have not been duplicated within a species, but retain the same function and have descended during evolution by the process of speciation only. Both orthologous and paralogous genes are **homologous**, but a comparison of paralogous genes from different species will not reflect the process of species divergence. Paralogous genes are usually present as members of a **gene family** within a genome. For more details, see [Orthologous Genes](#).

### Parasexual Cycle

The parasexual cycle allows new **allele** combinations to be formed independently of whether the organism can undergo a sexual cycle and meiosis. The parasexual cycle involves the fusion of [haploid](#) nuclei, [mitotic recombination](#), and random [chromosome](#) losses. Parasexuality was discovered by Giorgio Pontecorvo (1) and has been demonstrated with several **fungi** and, in a modified form, with cultured cells of multicellular animals.

The parasexual cycle is sketched in Fig. 1, taking a filamentous Ascomycete as a model. Fusion of haploid cells to form a [heterokaryon](#) occurs spontaneously in many fungi in a process called hyphal *anastomosis*. In many fungi this process occurs with certain pairs of strains said to be *vegetatively compatible* and fails with others. Compatibility requires the strains to carry certain combinations of alleles at certain genes. Heterokaryons can be prepared in the laboratory by fusion of protoplasts and other techniques in fungi without anastomosis and in many other organisms.

**Figure 1.** Parasexual cycle in filamentous *Ascomycetes*. Nuclei are followed from the haploid parents above to the haploid segregants below, through cell fusion, nuclear fusion, and haploidization. Recombinant nuclei are indicated by different shading patterns.

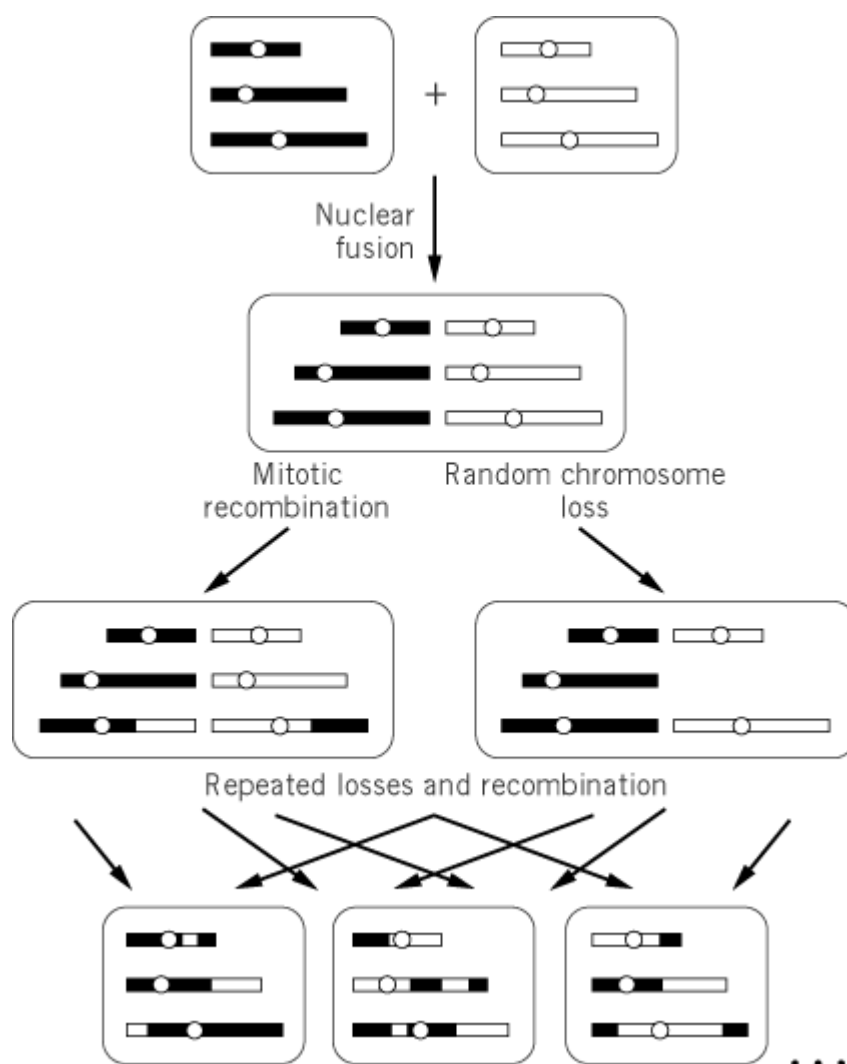


Two haploid nuclei in a heterokaryon may fuse to form a **diploid** nucleus. This process occurs seldom or not at all in vegetative heterokaryons, although it may be very quick after the fusion of gametes in the sexual cycle. Diploid cells are usually very hard to distinguish from heterokaryotic cells carrying the same genetic markers, but they may be very easily distinguished from haploid cells. In fungi in which all spores are uninucleate, colonies grown from diploid and haploid spores are distinguished easily with the help of appropriate genetic markers.

The diploid nuclei of many fungi are unstable because the loss of one chromosome of each pair is not deleterious. A slow process of haploidization takes place during mitotic growth as random chromosome losses convert the diploid into haploids through intermediates called **aneuploids**. The resulting haploid nuclei contain random combinations of the chromosomes of the two original haploid nuclei. The frequency of chromosome loss is increased by treatments with benomyl and other agents.

In addition, mitotic recombination occurs sporadically between the duplicated chromosomes of the diploid and the aneuploids and produces new combinations of their genetic markers (see [Mitotic Recombination](#)). The chromosomes in the final haploid nuclei are composed of random combinations of chromosome segments from the two original haploid nuclei (Fig. 2).

**Figure 2.** Chromosome events during the parasexual cycle. Three chromosomes are followed from the haploid parents above to the haploid segregants below, through diploidization, mitotic recombination, and chromosome loss.



As a result of the parasexual cycle, genes in different chromosomes exhibit recombination frequencies of 50 % while those that are close to each other in the same chromosome appear more or less tightly linked.

Fusion of cultured cells from different mammals gives rise first to heterokaryons and then to

*anfi*ploid nuclei, which contain all the chromosomes of the original cells. Chromosome losses reduce the number of chromosomes in these anfi

ploids, but not necessarily at random: chromosomes of *Homo sapiens* are lost more frequently than those of *Mus musculus* from their anfi

ploids. The correlation between the presence or absence of certain genetic markers and certain chromosomes in the segregants is a powerful technique for the quick assignment of genes to chromosomes (see **Cell fusion**, [Hybridomas](#)).

#### Bibliography

1. G. Pontecorvo (1956) *Annu. Rev. Microbiol.* **20**, 151–168.

#### Suggestions for Further Reading

2. G. Sermonti (1969) *Genetics of Antibiotic-Producing Microorganisms*. Wiley-Interscience, London (see Parasexuality, pp. 175–228).
3. A. J. Clutterbuck (1992) Sexual and parasexual genetics of *Aspergillus* species. *Biotechnology* **23**, 3–18.
4. C. J. Bos (1996) "Somatic recombination". In *Fungal Genetics* (C. J. Bos, ed.), Marcel Dekker, New York, pp. 73–95.

#### Parsimony

Parsimony is one of the principles of [phylogenetic tree](#) reconstruction in which the [phylogeny](#) of a group of species or genes is inferred to be the branching pattern that requires the smallest number of evolutionary changes ([1](#)). Of course, no one knows whether evolution actually occurred in this way, but the parsimony principle appears to work reasonably well, particularly for comparisons of closely related species or genes, but it tends to lose its effectiveness in comparison to more remotely related species or genes. For the latter comparisons, hidden mutational changes (due to reverse mutations) have accumulated to such an extent that the parsimony principle cannot detect all changes.

The merit of using the parsimony principle is that we can infer the amino acid or **nucleotide sequences** of all the intermediate stages in the continuous lineage between the ancestor and the present species. The intermediate stages may be trustworthy, however, only when closely-related genes are compared.

#### Bibliography

1. M. Ridley, (1997) *Evolution*, 2nd ed., Blackwell Science, USA.

#### Parthenogenesis

Parthenogenesis is the development of a [zygote](#) without [fertilization](#). Parthenogenesis can be induced experimentally by a number of procedures ([1](#)), but it occurs naturally in some species. A special case of parthenogenesis is called *arrhenotoky*. This is the mechanism of sex determination in many

insects. Fertilized eggs develop into **diploid** females, while unfertilized eggs develop parthenogenetically as **haploid** males. Arrhenotoky occurs in the four insect orders: Hymenoptera, Homoptera, Thysanoptera, and Coleoptera (2). In the Hymenoptera, arrhenotoky is the only mechanism of sex determination.

In addition to arrhenotoky, or male haploidy, is another form of parthenogenesis, called *thelytoky*, in which the unfertilized eggs develop as females. It can be either (a) cyclic, in which parthenogenetic and bisexual generations alternate, or (b) complete (3). Cyclic parthenogenesis is found among some insects, such as aphids (4), and in a few vertebrates (5). Complete parthenogenesis has appeared many times during evolution in a variety of phyla. The phylum Rotifera appears to be entirely parthenogenetic. There are even parthenogenetic species of vertebrates (6).

In arrhenotoky, the unfertilized eggs develop as haploid males, but in thelytoky the unfertilized eggs develop as diploid females (or polyploid females in some cases). In order for the unfertilized haploid eggs to develop as diploids, the genetic content of the cells must be doubled. Because parthenogenesis has arisen independently many times during evolution, this problem has been solved in a variety of ways (7). The failure of the second meiotic division during gametogenesis results in the production of diploid eggs. Meiosis can also proceed normally, with the fusion of two of the meiotic products producing a diploid egg. In other species, haploid eggs are produced, but fusion of haploid nuclei occurs during the early cleavage stages to restore diploidy.

#### Bibliography

1. Y. Fuyama (1986) *Genetics* **112**, 237–248.
2. M. J. D. White (1948) *Animal Cytology & Evolution*, Cambridge University Press, Cambridge, U.K. p. 267.
3. Reference 2, p. 280.
4. A. F. Shull (1925) *Am. Nat.* **59**, 138–154.
5. S. B. Hedges, J. P. Bogart, and L. R. Maxson (1992) *Nature* **356**, 708–710.
6. L. D. Densmore, C. C. Moritz, J. W. Wright, and W. M. Brown (1989) *Evolution* **43**, 969–983.
7. Reference 2, p. 281.

#### Suggestions for Further Reading

8. M. J. D. White (1948) *Animal Cytology & Evolution*, Cambridge University Press, Cambridge, U.K. pp. 280–302.
9. C. D. Darlington (1937) *Recent Advances in Cytology*, The Blakiston Company, Philadelphia, pp. 434–478.

#### Partial Specific (Or Molar) Volume

The volumes of macromolecules can be determined from their structures in several different ways (see [Molecular Surface, Volume](#); [Isoleucine \(Ile, I\)](#); [Accessible Surface](#)). In solution, the volumes relevant to thermodynamics are the partial specific and partial molar volumes. They involve not only the volume of the macromolecule itself, but also the contribution of solute–solvent interactions ([hydration](#)). In general, macromolecule solutions are systems with two or more components, containing the macromolecule, water, salts, etc. However, it is more convenient to reduce any such system to a binary one; one component is the solvent, ie water plus salts, etc., the other is the solute,



ie a protein. At constant temperature  $T$  and pressure  $P$ , the total volume  $V$  of such a binary solution comprising  $g_1$  grams of solvent and  $g_2$  grams of solute is expressed as follows:

$$V = g_1 \bar{v}_1 + g_2 \bar{v}_2 \quad (1)$$

where  $v_1$  and  $v_2$  are partial specific volumes of solvent and solute, respectively, as defined by

$$\bar{v}_1 = (\partial V / \partial g_1)_{T,P,g_2} \quad (2)$$

$$\bar{v}_2 = (\partial V / \partial g_2)_{T,P,g_1} \quad (3)$$

These parameters measure the change in volume upon adding further solvent or solute to the aqueous system.

The partial specific volume can be determined experimentally by measuring the densities of the solution containing the macromolecule (in units of  $\text{g cm}^{-3}$ ) as a function of its concentration. A variety of techniques have been used for density measurements: The majority of specific volume data have been obtained by pycnometry and digital densimetry (mechanical oscillator technique). Except for some rare cases, the plot of  $V$  vs  $g_2$ , which should be a straight line according to Equation (1), is not straight, due to solute–solute interactions. This means that the specific volumes of the solvent and solute are not the same as those of the pure phases of these components in most cases, and it is necessary to measure directly the values of  $v_1$  and  $v_2$ . In practice, the apparent specific volume  $f_2$  is defined instead of the partial specific volume, by assuming the specific volume of solvent in the solution to be the same as that of solvent in the pure phase  $v_1^0$ :

$$V = g_1 v_1^0 + g_2 \phi_2 \quad (4)$$

The relation between apparent and partial specific volumes is as follows:

$$\bar{v}_2 = \phi_2 + g_2 (\partial \phi_2 / \partial g_2)_{g_1} \quad (5)$$

From this equation, it is apparent that the apparent and partial specific volumes become equal at infinite dilution ( $g_2 = 0$ ). Usually, values of  $v_2$  and  $f_2$  at infinite dilution (denoted as  $v_2^0$  and  $f_2^0$ ) are used as the appropriate quantity for a protein molecule in solution not involved in intermolecular interactions with other protein molecules. When the molecular weight  $M_2$  of a solute is known, the partial and apparent molar volumes  $V_2$  and  $F_2$ , respectively, can be calculated using the equations  $V_2 = v_2 M_2$  and  $F_2 = f_2 M_2$ , respectively.

In two-component solutions, the partial specific volumes of folded proteins fall in the range of 0.69 to 0.76  $\text{cm}^3 \text{g}^{-1}$ , depending on their structure and amino acid composition (see Table 1 of [Compressibility](#)). The majority exhibit values in the range of 0.72 to 0.75  $\text{cm}^3 \text{g}^{-1}$ , which are within the range of the partial specific volumes of amino acids, which range from 0.559  $\text{cm}^3 \text{g}^{-1}$  for [aspartic acid](#) to 0.819  $\text{cm}^3 \text{g}^{-1}$  for [leucine](#) and [isoleucine](#) (Table 1). If we assume the additivity principle of volume, the apparent specific volume of a protein may be estimated from its amino acid composition by using the equation  $f_2 = \sum f_i w_i / \sum w_i$ , where  $f_i$  is the apparent specific volume of the  $i$ th amino acid residue (Table 1) and  $w_i$  its weight percent (Ref. 1). In spite of a number of simplifications and assumptions, this method yields relatively rather accurate results. It should be noted, however, that

the partial specific and apparent volumes are thermodynamic quantities, and changes in the environment (solvent composition, temperature, pressure, etc.) may cause significant changes in this volume.

**Table 1. Apparent specific volumes of amino acid residues at 25°C (Ref. 1)**

| Residue | Apparent specific volume (cm <sup>3</sup> g <sup>-1</sup> ) |
|---------|---|
| Ala     | 0.732   |
| Arg     | 0.756   |
| Asn     | 0.610   |
| Asp     | 0.573   |
| Cys     | 0.630   |
| Gln     | 0.667   |
| Glu     | 0.605   |
| Gly     | 0.610   |
| His     | 0.659   |
| Ile     | 0.876   |
| Leu     | 0.876   |
| Lys     | 0.775   |
| Met     | 0.739   |
| Phe     | 0.766   |
| Pro     | 0.748   |
| Ser     | 0.596   |
| Thr     | 0.676   |
| Trp     | 0.728   |
| Tyr     | 0.703   |
| Val     | 0.831   |

The partial specific volume of a protein in water is expressed as the sum of three contributions (Ref. 2): (1) the constitutive volume estimated as the sum of [van der Waals volumes](#) of the constitutive atoms ( $v_c$ ); (2) the volume of the cavities in the molecule due to imperfect atomic packing ( $v_{cav}$ ); (3) the volume change due to solvation or [hydration](#) ( $Dv_{sol}$ ).

$$\bar{v}_2^0 = v_c + v_{cav} + \Delta v_{sol} \quad (6)$$

The value of  $Dv_{sol}$  is usually negative, because hydration processes involve the contraction of solvent water. The volume change on transferring a protein from one solvent to another is calculated as the difference in  $v_2^0$  between the two solvents. The volume change due to any conformational change of protein is also defined as the difference in  $v_2^0$ . In such a case,  $v_c$  does not change, so the observed volume change is ascribed primarily to differences in  $v_{cav}$  and  $Dv_{sol}$ . In general, processes

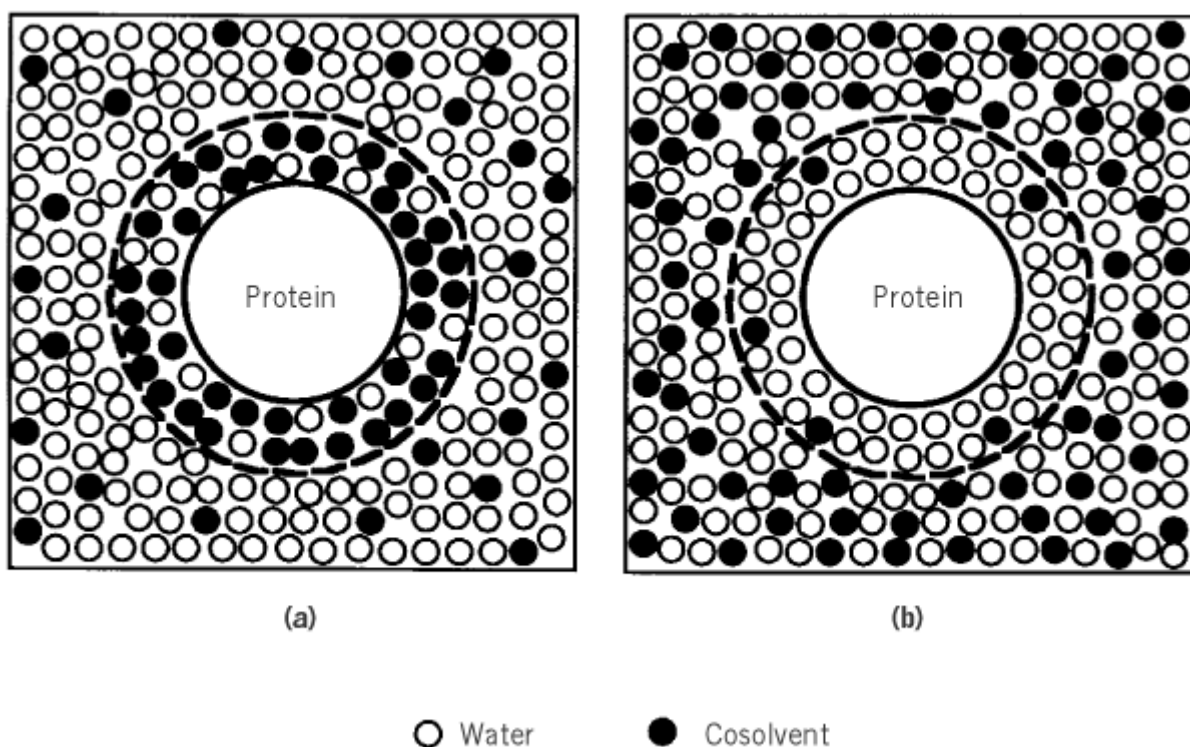
increasing the solvent accessible surface area, eg, unfolding and dissociation into subunits, are expected to produce negative volume changes due to increased hydration. Most experimental data meet this expectation, but some data show no change or even slightly positive volume changes. It is noteworthy that the volume change observed at a given protein concentration, which is based on  $v_2$ , may differ, even in its sign, from the result obtained at infinite dilution.

The partial specific volume measured at constant molality of the added cosolvent (eg, salts or organic compounds) is a useful quantity, as it is a truly thermodynamic volume. In practice, however, the measurement of this volume often accompanies difficulties in sample preparation, and most experiments (especially in **ultracentrifugation** or **small-angle X-ray scattering**) are usually done with solutions that have been subjected to [dialysis](#) to equilibrium against the solvent. At equilibrium, the chemical potentials of diffusible components are equal on the two sides of a membrane, but the molalities of the components are not necessarily identical mainly due to Donnan effects and preferential binding of solvent components to protein. The apparent isopotential specific volume  $f_2^{0'}$ , which is measured at constant chemical potential of the solvent components (equilibrium dialysis systems), can be obtained from the density differences between the protein solutions and the dialyzate by the same procedures as used for  $f_2^0$ . However,  $f_2^{0'}$  is not a partial quantity in the thermodynamic sense, since the molalities of the diffusible components are not held constant, but this specific volume is operationally very useful because it reflects preferential protein–solvent interactions. In a two-component system or a dilute buffer system where Donnan effects and/or preferential solvent interactions are absent,  $f_2^0$  and  $f_2^{0'}$  should be identical. In the presence of high concentration of a third component, however, large differences between these two specific volumes occur, as is shown in Table 2. The difference between  $f_2^0$  and  $f_2^{0'}$  was derived theoretically by Casassa and Eisenberg (3):

$$\phi_2^0 - \phi_2^{0'} = (\xi_1 1/d_0 - \bar{v}_1) = \xi_3(1/d_0 - \bar{v}_3) \quad (7)$$

where  $v_1$  and  $v_3$  are partial specific volumes of components 1 (water) and 3 (cosolvent), respectively;  $x_1$  and  $x_3$  represent preferential interaction parameters in terms of the number of grams of component 1 (or 3) that must be added per gram of component 2 to maintain components 1 and 3 at constant chemical potential. Then, the values  $x_1$  and  $x_3$  can be calculated with two specific volumes  $f_2^0$  and  $f_2^{0'}$ , which are usually measured by digital densimetry within an accuracy of  $\pm 0.002 \text{ cm}^3 \text{ g}^{-1}$  at isomolal and isopotential conditions of solvent components (Ref. 4). Positive values of  $x_3$  or negative  $x_1$  indicate the [preferential binding](#) of component 3 to protein; the opposite, negative  $x_3$  or positive  $x_1$ , indicates that water interacts more favorably with the protein, ie, there is [preferential hydration](#) (Fig. 1). Whether preferential binding or preferential hydration occurs depends on the affinity of the cosolvent for the protein. As seen in Table 2, preferential hydration is observed for cosolvents having a potential ability to stabilize and **crystallize** proteins, whereas preferential binding occurs with **denaturants**. Thermodynamic properties such as solubility (phase separation), stability, and chemical equilibrium are subject to preferential solvent interaction, because it is directly linked to the chemical potential of the protein.

**Figure 1.** Schematic illustration of solvent component distributions in mixed solvents. (a) Preferential binding: Cosolvent (additive) is present in the solvation shell of protein at a greater local concentration than in the bulk solvent. (b) Preferential hydration: Exclusion of cosolvent from the protein surface results in enrichment of water in the solvation shell. [Taken from G. C. Na and S. N. Timasheff (5).]



**Table 2. Apparent specific volumes of proteins and DNA at constant molality and constant chemical potential of various cosolvents (Ref. 6)**

| Solvent                                      | Temperature (°C) | $f_2^0$ (cm <sup>3</sup> g <sup>-1</sup> ) | $f_2^{0r}$ (cm <sup>3</sup> g <sup>-1</sup> ) | $x_3$ (g <sup>-1</sup> ) |
|--|------------------|--|---|--------------------------|
| <u>Chymotrypsinogen A</u>                    |                  |  |   |                          |
| 30% glycerol, 10 mM HCl                      | 20               | 0.727                                      | 0.745   | -0.123                   |
| 1 M sucrose, 1 mM HCl                        | 20               | 0.736                                      | 0.773   | -0.138                   |
| 10% PEG1000, pH 3                            | 20               | 0.732                                      | 0.743   | -0.07                    |
| 6 M GdmCl                                    | 20               | 0.729                                      | 0.712   | 0.15                     |
| 8 M urea                                     | 20               | 0.730                                      | 0.720   | 0.08                     |
| <br>   |                  |  |   |                          |
| <u>Bovine serum albumin</u>                  |                  |  |   |                          |
| 1 M Na <sub>2</sub> SO <sub>4</sub> , pH 5.6 | 20               | 0.735                                      | 0.788   | -0.074                   |
| 2 M glucose, pH 6                            | 20               | 0.727                                      | 0.750   |                          |
| 30% glycerol                                 | 25               | 0.729                                      | 0.744   | -0.101                   |
| 30% sorbitol                                 | 25               | 0.738                                      | 0.761   | -0.092                   |

|                               |    |       |       |       |
|-------------------------------|----|-------|-------|-------|
| 20% propylene glycol,<br>pH 2 | 25 | 0.740 | 0.726 | 0.211 |
| 1 M GdmCl                     | 20 | 0.735 | 0.729 | 0.025 |
| 6 M GdmCl                     | 20 | 0.717 | 0.724 | 0.06  |
| <u>DNA (calf thymus)</u>      |    |       |       |       |
| 0.20 M NaCl<br>(NaDNA)        | 25 | 0.503 | 0.540 |       |
| 0.20 M CsCl<br>(CsDNA)        | 25 | 0.446 | 0.503 |       |

---

The specific volumes  $f_2^0$  and  $f_2^{0r}$  of nucleic acids are very small (around  $0.5 \text{ cm}^3 \text{ g}^{-1}$ ) and remarkably influenced by solvent composition (especially salt concentration) because they are polyelectrolytes with high charge densities (Table 2).

#### Bibliography

1. A. A. Zamyatnin (1984) *Ann. Rev. Biophys. Bioeng.* **13**, 145–165.
2. W. Kauzmann (1959) *Adv. Protein Chem.* **14**, 1–63.
3. E. F. Casassa and H. Eisenberg (1964) *Adv. Protein Chem.* **19**, 287–395.
4. J. C. Lee, K. Gekko, and S. N. Timasheff (1979) *Methods Enzymol.* **61**, 26–49.
5. G. C. Na and S. N. Timasheff (1981) *J. Mol. Biol.* **151**, 165–178.
6. H. Durchschlag (1986) In *Thermodynamic Data for Biochemistry and Biotechnology* (H. J. Hinz, ed.), Springer-Verlag, Berlin, pp. 45–128. Contains details of specific volumes and their tabulation.

#### Particle Electrophoresis

As a result of the rising interest of molecular biologists in cells and **organelles**, the analytical separation and the preparative fractionation of those species by [electrophoresis](#) becomes increasingly important. Their large size makes [gel electrophoresis](#) difficult, as gels with very large pores are required. Specialized electronic devices for their separation on an analytical scale have been developed (1). Common [capillary zone electrophoresis](#) instrumentation has also been applied to particles up to  $10 \mu\text{m}$  in diameter (2). Effective resolution of particles on the basis of differences in their size and shape, using solutions of polymers of various types and molecular weights for **molecular sieving**, is possible (3), but the selection of the appropriate polymer is arbitrary at present. The resolving power is best when the particle diameter is less than the effective pore size (or “screening length”) of the polymer network in solution, and the separation of different particles is due to a differential occupation of available spaces in the network. By contrast, resolution is poor when the particle diameter exceeds that pore size, and separation is then due to a differential displacement of the network by the particle.

On preparative scales, separations in free solutions are compatible with available instrumentation, such as free-flow electrophoresis apparatus (4) or automated gel electrophoresis apparatus in its preparative mode in conjunction with polymer solutions is feasible (5). In both analytical and preparative electrophoresis, sedimentation of the particles becomes an impediment to electrophoretic resolution only when the particle size surpasses about 8  $\mu\text{m}$  in diameter (2). In those applications, electrophoretic separations in a microgravity environment are feasible (6).

### Bibliography

1. A. W. Preece and K. A. Brown (1989) *Adv. Electrophoresis* **3**, 351–404.
2. S. P. Radko, M. M. Garner, G. Caiafa and A. Chrambach (1994) *Anal. Biochem.* **223**, 82–87.
3. D. Tietz, A. Aldroubi, H. Pulyaeva, T. Guszczynski, M. M. Garner, and A. Chrambach (1992) *Electrophoresis* **13**, 614–615.
4. J. Bauer and G. Weber (1996) *Electrophoresis* **17**, 526–528.
5. E. Gombocz and E. Cortez (1995) *Appl. Theor. Electrophoresis* **4**, 197–209.
6. R. S. Snyder and P. H. Rhodes (1990) *Proc. 4th Eur. Symp. Life Sci. Res. Space*, Trieste, Italy, pp. 411–415.

### Suggestion for Further Reading

7. A. Budde, E. Knippel, G. Gruemmer, J. Treichler, H. Brockmann, E. Donath, and H. Baeumler (1996) Electrophoretic fingerprinting and multiparameter analysis of cells and particles. *Electrophoresis* **17**, 507–511.

## Partition Coefficient

A molecule's *partition coefficient* is the equilibrium constant for its transfer between two phases that are immiscible (eg, [water](#) and cyclohexane). This equilibrium constant is measured by analyzing both phases after equilibrium has been achieved by stirring or shaking, and then dividing the molar concentration in the first phase by the molar concentration in the second. To avoid (or detect) the occurrence of aggregation (eg, formation of dimers) in either phase, the experiment must be carried out at a series of concentrations and the results extrapolated to infinite dilution. Aggregation in either phase is signaled by a tendency of the apparent equilibrium constant to increase in favor of one phase, that in which aggregation is greatest, as the total concentration of solute increases. Such partition coefficients depend on the particular temperature at which they were measured and are used in determining [hydrophobicity](#). See also [Transfer Free Energy](#).

## Parvalbumin

Parvalbumin binds calcium with high affinity and appears to be involved in buffering of intracellular  $\text{Ca}^{2+}$  signals. It is a small protein (10–12 kDa) that is present in muscle, brain, kidney, testis, and adipose tissue. The various parvalbumin proteins can be classified into two evolutionarily distinct groups: the a-lineage and the b-lineage. Both types of proteins have, however, similar structures and

functions. Parvalbumins are members of the [EF-Hand Motif](#) family of [calcium-binding proteins](#).

## 1. Biological Function

Parvalbumin was first identified in skeletal muscle cells, and its function is best understood in these cells. For muscle cells to relax after a stimulus, the  $\text{Ca}^{2+}$  signal must be quenched by **active transport** of the  $\text{Ca}^{2+}$  ions into the sarcoplasmic reticulum. Parvalbumin has been implicated in skeletal muscle relaxation after a single twitch, as well as after a tetanic (continuous) contraction, and it may facilitate  $\text{Ca}^{2+}$  translocation in the sarcoplasm. There is still some controversy, however, over whether parvalbumin can exchange  $\text{Ca}^{2+}$  and  $\text{Mg}^{2+}$  ions fast enough to be involved in relaxation after a single twitch.

The function of parvalbumin in other cell types is less clear. It may act as a  $\text{Ca}^{2+}$  buffer to protect cells from damage due to calcium overload, particularly in the nervous system. It may also play an indirect role in  $\text{Ca}^{2+}$  [signal transduction](#) by affecting the nature of the incoming signal. For instance, it could modulate how quickly the signal builds, or how long the signal lasts. Parvalbumin may also be involved in cell growth and the cell cycle.

Oncomodulin is a b-lineage parvalbumin of particular interest, because it is found in abundance in certain cancer cells. It is normally found only in the preimplantation [embryo](#) and in the placenta. The function of this protein is completely unknown. However, its expression is cell cycle-dependent, and oncomodulin [antisense oligonucleotides](#) inhibit cell growth when injected into chemically transformed fibroblasts.

## 2. Ion-Binding Properties

Parvalbumin's two functional EF-hands will bind either  $\text{Ca}^{2+}$  or  $\text{Mg}^{2+}$  ions. Two main types of metal-binding loops are found in parvalbumins. The most common is the  $\text{Ca}^{2+}/\text{Mg}^{2+}$  “mixed” binding loop. These sites bind  $\text{Ca}^{2+}$  with high affinity ( $K_{\text{Ca}} = 10^7$  to  $10^9\text{M}^{-1}$ ), and  $\text{Mg}^{2+}$  with moderate affinity ( $K_{\text{Mg}} = 10^3$  to  $10^5\text{M}^{-1}$ ) (1). Binding is competitive, and at basal physiological ion concentrations ( $[\text{Ca}^{2+}] = 1 - 10 \times 10^{-8}\text{M}$ ;  $[\text{Mg}^{2+}] = 1 \times 10^{-3}\text{M}$ ), mixed sites are occupied by  $\text{Mg}^{2+}$ . Exchange between bound  $\text{Ca}^{2+}$  and  $\text{Mg}^{2+}$  is slow (2, 3). The second type of binding loop is the  $\text{Ca}^{2+}$ -specific loop. The  $\text{Ca}^{2+}$  affinity of these sites is lower than in mixed sites ( $K_{\text{Ca}} = 10^5 - 10^7\text{M}^{-1}$ ), and the  $\text{Mg}^{2+}$  affinity is significantly reduced ( $K_{\text{Mg}} = 10^1 - 10^2\text{M}^{-1}$ ) (4). These sites are metal-free under basal physiological conditions, and can therefore bind  $\text{Ca}^{2+}$  quickly. Oncomodulin is the only member of the parvalbumin subfamily that has a  $\text{Ca}^{2+}$ -specific site paired with a mixed site.

## 3. Structure

Parvalbumins are the only known members of the EF-hand family with three such motifs (a prototypical functional domain is composed of only two). However, the *N*-terminal EF-hand of parvalbumin does not bind  $\text{Ca}^{2+}$  and is therefore considered to be a nonfunctional ancestral site. The protein is largely helical, with six *alpha*-helices, labeled A through F. The functional calcium-binding loops occur between helices C and D and between helices E and F. The loop between helices A and B is shorter than the functional loops (10 residues rather than 12), and as mentioned above, does not bind metal ions.

Parvalbumin was the first EF-hand protein whose high resolution three-dimensional structure was determined (5). To date, all the known parvalbumin structures have metal bound. Structures of parvalbumin bound to both  $\text{Ca}^{2+}$  and  $\text{Mg}^{2+}$  are available (6-13). The overall structure of the protein

is very similar irrespective of which ion is bound. The two  $\alpha$ -helices flanking each functional binding loop (the C/D pair and the E/F pair) are nearly perpendicular to each other, while the A and B helices, which flank the non-functional binding site, are nearly antiparallel (Fig. 1). However, there are some differences between the binding loops with bound  $\text{Ca}^{2+}$  and those occupied by  $\text{Mg}^{2+}$ , particularly in the position of the highly conserved glutamic acid residue that is the final residue of the binding loop. In parvalbumin and all other EF-hand calcium-binding proteins studied to date, this glutamic acid is a bidentate calcium ligand. In the  $\text{Mg}^{2+}$ -bound parvalbumin structures, this glutamic acid is a monodentate ligand (14). This difference reflects the preference of the  $\text{Mg}^{2+}$  ion for six ligands, whereas the  $\text{Ca}^{2+}$  ion can accommodate seven or even eight ligands. The apo structure is of questionable physiological importance, since almost all parvalbumin binding loops are occupied by  $\text{Mg}^{2+}$  under basal physiological conditions. NMR studies indicate, however, that apo-parvalbumin has a structure very similar to that seen for the ion-bound states (15). This is consistent with its proposed role as a calcium buffer, which requires high  $\text{Ca}^{2+}$  affinity but not the conformational change associated with the transduction of a calcium signal into a mechanical or metabolic response (see [Calmodulin](#)).

**Figure 1.** Ribbon representation of  $(\text{Ca}^{2+})_2$ -parvalbumin using the coordinates of 1PVA (7).



## Bibliography

1. J. Haiech, J. Derancourt, J.-F. Pechere, and J. G. Demaille (1979) *Biochemistry* **18**, 2752–2758.
2. Y. Ogawa and M. Tanokura (1986) *J. Biochem.* **99**, 73–80.
3. Y. Ogawa and M. Tanokura (1986) *J. Biochem.* **99**, 81–89.
4. M. Milos, J.-J. Schaer, M. Comte, and J. A. Cox (1986) *Biochemistry* **25**, 6279–6287.
5. R. H. Kretsinger and C. E. Nockolds (1973) *J. Biol. Chem.* **248**, 3313–3326.
6. J.-P. Declercq, B. Tinant, J. Parello, and G. Etienne (1988) *J. Mol. Biol.* **202**, 349–353.
7. J.-P. Declercq, B. Tinant, J. Parello, and J. Rambaud (1991) *J. Mol. Biol.* **220**, 1017–1039.
8. F. R. Ahmed, D. R. Rose, S. V. Evans, M. E. Pippy, and R. To (1993) *J. Mol. Biol.* **230**, 1216–1224.



9. V. D. Kumar, L. Lee, and B. F. P. Edwards (1990) *Biochemistry* **29**, 1404–1412.
10. P. C. Moews and R. H. Kretsinger (1975) *J. Mol. Biol.* **91**, 201–228.
11. F. Roquet, B. Tinant, J.-P. Declercq, B. Rambaud, and J. Parello (1992) *J. Mol. Biol.* **223**, 705–720.
12. A. L. Swain, R. H. Kretsinger, and E. L. Amma (1989) *J. Biol. Chem.* **264**, 16620–16628.
13. C. A. McPhalen, A. R. Sielecki, and J. Santarsiero (1994) *J. Mol. Biol.* **235**, 718–732.
14. T. C. Williams, D. C. Corson, K. Oikawa, W. D. McCubbin, C. M. Kay, and B. D. Sykes (1986) *Biochem.* **25**, 1835–1846.

### Suggestion for Further Reading

15. T. L. Pauls, J. A. Cox, and M. W. Berchtold (1996). The Ca<sup>2+</sup>-binding proteins parvalbumin and oncomodulin and their genes: new structural and functional findings. *Biochim. Biophys. Acta* **1306**, 39–54.

## Pathogenesis In Fungi

### 1. Pathogenesis and Virulence

To understand how a given organism is able to infect its host, it is important to determine which aspect(s) of the biology of the fungus differs from nonpathogenic **fungi**. According to Bulmer and Fromtling, “Pathogenicity is the capacity of an organism to damage, ie, to produce disease in another animal or plant” (1). This property is the result of a direct interaction between the pathogen and the host. In many ways, pathogenicity and virulence are very broadly defined terms. If one considers all the possible gene products required for growth on or within a host, those that contribute to disease symptoms and host-range specificity, and those that govern their regulation, then we may be considering hundreds (or thousands) of potential components. Furthermore, the ability of an organism to adapt quickly to new environmental challenges is itself likely to be an effective virulence mechanism, in that it enables the pathogen to take advantage of new opportunities.

One of the most cited examples of virulence in the fungal world is the ability to switch morphological forms. While some fungi are pathogenic as yeast or conidial forms [eg, *Histoplasma capsulatum*, *Coccidioides immitis*, and *Blastomyces dermatitidis*, and *Coccidioides immitis* (2-4)], numerous correlations between hyphal growth and pathogenicity have been cited for many of the fungal pathogens. In fact, mutants incapable of switching to hyphal growth are often revealed as nonpathogenic *cpg-1* deletions in *C. parasitica* (5), deletions of *fuz7* in *U. maydis* (6), or *cph1/efg1* homozygous deletions in *C. albicans* (7). The ability to switch morphological forms is referred to as dimorphism and is governed by a number of external factors including nitrogen, pH, temperature, and intracellular biochemical pathways (eg, level of cAMP) (8). Similarly, correlations between a strain's ability to myceliate and its virulence (defined by the LD50, lethal dose to obtain 50% killing of the host) have also been cited as strong virulence factors.

Pathogen-produced products which contribute in some way to either the establishment or continuation of the disease process have also been cited as virulence factors. Historically, secretion of degradative enzymes such as [elastase](#) and certain [proteinases](#) (9, 10) have been cited as potential virulence factors using a variety of methods to correlate disease progression with levels of expression. *Paracoccidioides brasiliensis* has been shown to alter its level of virulence with the production of  $\alpha$ -1,3 glucan, suggesting that this polymer may act as a protective layer (11). Genes

that contribute to adherence of the fungus to its host have also been regarded as virulence factors (12). Numerous genetic mutations or strain isolates have been analyzed in *in vivo* experiments to identify underlying gene products causing changes in pathogenicity. In fungus–plant interactions, the pathogen may produce a substance that specifically detoxifies plant-secreted antifungal compounds. For example, *G. graminis* deletion mutants have been generated that are relatively resistant to the production of avenacin by oat roots but retain full pathogenicity to wheat, which does not synthesize saponins (13); isolates of *Nectria haematococca* are similarly able to detoxify the pea phytoalexin, pisatin, using a [cytochrome P450](#) mechanism (14). Furthermore, the biosynthetic pathway leading to melanin production in *M. grisea* correlates with altered virulence in this fungus; strains with reduced melanin content are less virulent (15). In some instances, these factors are reported to increase the pathogen's ability to penetrate host tissue (as in the case of melanin or proteinase production), while in others they may simply contribute to the overall fitness of the organism *in vivo*. For example, some auxotrophies are known to decrease virulence *in vivo* but play little or no role *in vitro*. Mutations in the *URA3* and *ADE2* genes of *C. albicans* and *C. neoformans* attenuate virulence *in vivo* (16, 17).

Alternatively, pathogen-induced factors that result in direct attacks upon the host are also known. The best-characterized system of host-selective toxin production by a plant fungus is in *Cochliobolus*. Different species of this genus are known to produce plant-specific toxins: T-toxin, HC-toxin, and victorin, of which only HC-toxin production has been deleted and shown to be nonpathogenic on maize (reviewed in Ref. 18).

One area in which there is sparse molecular information, but a great deal of inferential information, is in the interactions between the host immune system and pathogen. Research directed toward understanding the role of the unique polysaccharide capsule of *Cryptococcus neoformans* pathogenicity has uncovered a large difference in ability to phagocytose capsulized versus acapsular *C. neoformans* mutants *in vitro*, although it is not yet clear whether the ability to engulf a cell correlates with killing (4). Additionally, macrophages appear to recognize and engulf wild-type, single, and double mutants of the MAPK pathway in *C. albicans* equally well; however, only wild-type strains were able to switch morphologies, adding further support that dimorphism plays a significant role in virulence in this organism (7). It is tempting to speculate that in the *chs3* *C. albicans* cell wall mutant, prolonged survival *in vivo* (19) had more to do with a change in its interaction with the host immune system than with any other factor. If, in fact, the host defense system plays a more significant role in disease progression than has been realized, it will be interesting to determine whether other factors such as increased lipid content cited as a virulence factors in strains of *Coccidioides immitis*, *Histoplasma capsulatum*, and *Blastomyces dermatitidis* (1, 20) is a result of an alteration in the interaction rather than a simple defense on the part of the fungus.

## 2. Signal Transduction and Pathogenesis

[Signal transduction](#) pathways play an important role in cellular biology of all eukaryotic organisms, including fungi. The ability of a pathogen to receive signals from its environment and transmit those signals intracellularly to alter development or entire biochemical pathways so as best to take advantage of those changes has been recognized as a significant partner in pathogenesis of the organism. Some of the pathways have been carefully worked out in *S. cerevisiae* and *S. pombe* with regard to [cell cycle](#), filamentation, stress, and the pheromone response. It appears that the molecules required to relay a signal from the external environment to the appropriate [transcription factors](#) within the nuclei of many of the filamentous fungi are highly **homologous** to those described in other organisms. Some of the molecules have been identified and isolated on the basis of sequence homology, and some by [complementation](#) or the use of similar screens. While the precise targets of many of these molecules have yet to be discovered, it is becoming clearer that many of the pathways in yeast are similarly conserved in more distantly related organisms.

Virulence and signal transduction pathways are inescapably intertwined in the biology of fungi and

their pathogens. In the filamentous fungi, signal transduction pathways have been identified that regulate conidiation in *A. nidulans*, virulence in *C. albicans*, *C. parasitica*, *U. maydis*, and *M. grisea*, and mating. While only a few of the components have been identified in the filamentous fungi to date, there is much anticipation that most, if not all, of the players will eventually be isolated and characterized. Although their exact developmental pathways may be very different in the different fungi, it is important to keep in mind the similarities in order to appreciate the conservation of function across the biological world.

### 2.1. MAP Kinase Pathway

Signal transduction involves the activation of a receptor molecule at the surface of a cell, which transmits the signal inward toward associated proteins at the cytoplasmic interface of the plasma membrane. Changes in these proteins are relayed via a small number of known pathways (either by a cascade of kinases or by the release of a [second messenger](#)) which in turn activate a number of protein targets, some of which are themselves transcriptional activators. Subsequent alterations in transcription of whole sets of genes then reprogram the cell's "transcriptome," resulting in a particular cellular response to the initial signal. That signal may be any of a number of ligands that tells the cell something about its external environment: for example, nutritional availability, osmolarity, presence of compatible mating partners, or light. In one common pathway observed in all eukaryotes, a heterotrimeric G protein activates a mitogen-activated protein (MAP) kinase module that results in the further activation of target proteins, among which is a transcription factor. In some cells, multiple MAPK signaling pathways can be operating simultaneously, altering the organisms' transcriptome moment-by-moment. A MAP kinase pathway influence on pathogenicity has been shown for the corn smut, *U. maydis*, in which tumor production and completion of the sexual cycle is dependent upon successful compatible mating ([21](#)). Appressorium formation and pathogenic growth of *M. grisea* *in planta* also require a functional MAP kinase pathway ([22](#)). Because the components of these pathways are so highly conserved, one of the central questions that remains unanswered is how specificity is achieved and how that affects the pathogenicity of the organism.

### 2.2. Cyclic AMP Pathway

A second common pathway that impinges upon pathogenicity is the **cyclic AMP**-mediated pathway, which begins with a plasma membrane-bound sensor molecule (none of which has been isolated for the filamentous fungi to date). Upon a functional receptor–ligand interaction, a signal is transmitted to a heterotrimeric GTP-binding protein, possibly through a conformational change in the membrane protein. Second messenger cyclic AMP (cAMP) is produced via [adenylate cyclase](#), which then alters the activity of cAMP-dependent [protein kinase A](#) (PKA) by binding to the regulatory subunit of PKA. This releases the catalytic subunits, which continue to activate downstream targets (including transcriptional regulatory proteins). A few components of this pathway have been identified in filamentous fungi. Mutants in *uac1* (adenylyl cyclase of *U. maydis*) have been shown to be defective in their ability to switch from budding to filamentous growth and, in addition, were unable to cause disease on corn (reviewed in (Ref. [23](#))); disruption of *adr1* (catalytic subunit of PKA) results in nonpathogenesis, whereas inactivation of *uka1* (catalytic subunit of PKA) has apparently little or no effect on growth, morphology, or pathogenesis. Disruption of *fil1* (G a subunit of *U. hordei*) also prevented normal hyphal to bud switching. Interestingly, of the four G a proteins isolated from *U. maydis*, only one, *gpa3*, was shown to be required for pathogenic development, disruption of the other three displayed no obvious phenotypes. Furthermore, *gpa3* mutants were unable to respond to pheromone signals strongly suggesting that a single G a subunit may play a role in multiple signaling pathways as has been observed in *S. cerevisiae* and *C. albicans*. Mutants disrupted in the catalytic subunit of *M. grisea* (CPKA) were shown to be unable to form appressoria on artificial surfaces, and they were also unable to infect rice unless exogenous cAMP was administered ([23](#)).

## 3. Strategies to Identify Virulence Factors

### 3.1. Gene Knockouts

There are a number of molecular genetic strategies that can be used to identify genes responsible for virulence in fungi ([24](#)). In a standard experiment, genetic manipulations are made such that the

chromosomal copy of a gene in question is replaced by a marker gene. The growth rate of these strains is then compared to the normal wild-type strain and subsequently tested *in vivo* to determine the potential of the “knocked out” gene to cause disease. Proper controls wherein additional strains containing the chromosomal knockout, but also harboring a transformant carrying a reintroduced wild-type copy of the gene, or a heterozygote, is also compared in order to ensure that the effects observed are due to the presence or absence of a specific gene product and not some other underlying mechanism. It should be noted that due to the frequency of ectopic integration in filamentous fungi, these experiments may be difficult to interpret. They have been successfully accomplished for a wide variety of gene products in the dimorphic fungus, *C. albicans*. Using a modification of the original “ura-blaster” strategy (25), it is possible to knockout both copies of a given gene in this perpetually diploid organism. Genes in cell wall biosynthesis, signal transduction, and so on, have been identified as involved in virulence in such tests (7, 19, 26). This directed approach to virulence determination has also been successful in *A. nidulans* (27), *Gaeumannomyces graminis* (28), and *C. neoformans* (29), for example. Using this approach, there is the potential to knockout each gene in a given organism systematically and test its pathogenicity *in vivo*.

### 3.2. Gene Transfer

Virulence factors have also been identified by their ability to confer pathogenicity upon a related species or strain that was previously nonpathogenic. In this approach, a **library** is constructed and **transformed** into a nonpathogenic organism. Entire library pools are then tested for changes in related events that are suspected of having influences upon pathogenicity (adherence properties, toxin production, proteinase activity, etc.), and gene products are identified and then retested for altered virulence relative to the starting host strain. This approach has been successful in identifying a gene conferring adhesion properties from *C. albicans* to *S. cerevisiae* (30).

### 3.3. Gene Expression

This may be accomplished in a number of different ways. First is to identify all the genes expressed during an infection, by **complementary DNA** cloning, differential display analysis, or **promoter** fusions (all of which are technically feasible in the filamentous fungi). The idea is that crucial genes will be expressed sometime during or immediately prior to the establishment of an infection. These genes can then be further screened for their relative virulence, using a variety of knockout or deletion strategies. The *MPGI* protein of *M. grisea* was identified using this approach and confirmed to be important in the development of fungal infection structures (31). Using a differential display approach, two genes from *U. maydis* (32) and six genes from *C. albicans* have been identified that are differentially induced during infection (33). As additional sequence information becomes available for other fungal genomes, the use of promoter fusions as an approach will probably become somewhat more popular. Following the initial identification of clones that express a viable fusion, sequence analysis of the upstream sequences directing the fusion will reveal the target gene. Conceivably, **libraries** in which the reporter proteins are cell surface-directed will aid in the initial identification of positive clones in this type of model.

One of the most fascinating stories concerning pathogenicity in fungal biology involves the interaction between *Cryphonectria parasitica* (chestnut blight fungus) and a hypovirus. These fungi, which have wreaked havoc upon the chestnut forests of the northeastern United States, are themselves infection-susceptible to a hypovirus that causes reduced virulence of the fungus for the plant. In addition to causing reduced virulence, the virus causes altered colony morphology, reduced pigmentation, and attenuated asexual sporulation. Differential display analysis of infected and noninfected isolates identified a candidate gene *cpgI* whose reduced expression in infected isolates appears to correlate with pathogenicity (5, 34).

### 3.4. Mutagenesis

Historically, a traditional genetic approach in which mutations are isolated and further characterized has proven beneficial in the understanding of pathogenesis of fungi. One of the more recent strategies to understanding and defining virulence factors has come from alterations in **transposon**-mediated mutagenesis screens. Transposons are relatively small fragments of DNA that are capable

of recombining nearly randomly into different sites within a single [genome](#). They have proven useful in creating large number of “insertional” mutant [libraries](#) that are easily screened for loss-of-function phenotypes. Unfortunately, filamentous fungi appear to lack such elements. However, **restriction enzyme**-mediated integration (REMI) has been used to create insertional mutants in filamentous fungi. Pathogenicity genes were identified for *Cochliobolus heterostrophus* (35) and *Ustilago maydis* (36) using REMI.

A modification of this method, called signature-tagged mutagenesis, has been used successfully in bacterial pathogenesis to identify insertional mutants that fail to grow *in vivo* (37). In this method, every insertional mutant carries with it a unique oligomer tag that identifies it from all other mutants. After inoculating a mouse and recovering all survivors, DNA is prepared and hybridized to colony blots of the original input pool. Mutants not surviving are not recovered and are easily identified as not hybridizing. Identification of the gene into which the insertional tag is present is performed by sequencing. One of the only drawbacks to this approach is that only nonessential genes may be recovered, because insertions into essential genes are by definition lethal in culture.

### 3.5. Use of *In Vivo* Models

The use of *in vivo* models to confirm the role of individual genes in virulence is an important aspect of this field. However, it is important to recognize the difference between the establishment of infection and the maintenance of infection. Both are clearly critical factors in host–pathogen interactions, yet targets should be different, although somewhat overlapping. Second, the size of the inoculum is a well-established variable in disease models. It is not inconceivable that certain genes play a more significant role at low inocula in certain situations while other virulence genes are turned on under conditions of high inocula. Third, it is important to remember that virulence is multifactorial and results from a balance between the pathogen and the host immune system. Comparisons across different fungal backgrounds or different host models should proceed with caution.

## 4. Fungal Virulence Genes as Targets for Antifungal Therapy

It is not enough that we understand how an organism grows and develops or produces disease; it is our ability to use this information to circumvent those processes that will enable humanity to control and manage diseases produced by microorganisms. Compounds aimed toward inhibiting functions of known or suspected virulence factors have been suggested as likely to be effective in reducing disease even though many of the virulence genes identified to date have been shown to be nonessential for growth and reproduction of the organism *in vitro*. A compound directed toward such a nonessential target would likely prevent fungal growth until the compound were either removed or inactivated. In clinical situations where prophylaxis is appropriate, a reduction in anticipated fungal pathogen virulence could significantly influence patient outcome, because the disease will not likely gain entry initially. Unfortunately, the human patient population that often contracts fungal diseases is often immune-debilitated. From an agricultural view, inhibition of fungal growth will probably lead to increased production over a short season, making these attractive targets. Another consideration in choosing virulence factors as targets is the likelihood that parallel and/or redundant pathways may exist. With all of this, virulence genes are important to consider as potential drug targets.

Strategies for identifying fungal virulence genes as potential targets for drug development and/or vaccine development include many of the cloning strategies identified earlier. Furthermore, the use of [bioinformatics](#) and [genome](#) sequencing has made it possible to identify potential candidates based on previously known motifs, which may then be tested *in vivo* for specific disease.

## Bibliography

1. G. S. Bulmer and R. A. Fromtling (1983) "Pathogenic Mechanisms of Mycotic Agents". In *Fungi Pathogenic for Humans and Animals*, Part B, (D. H. Howard, eds.), Marcel Dekker, New

York.

2. J. E. Loyd, R. M. DesPrez, and R. A. Goodwin, Jr. (1990) *Histoplasma capsulatum*. In *Principles and Practice of Infectious Diseases* (G. L. Mandell, R. G. B. Douglas, and J. E. , eds.), Churchill Livingstone, New York, pp. 1989–1999.
3. S. W. Chapman (1990) *Blastomces dermatitidis*. In *Principles and Practice of Infectious Diseases* (G. L. Mandell, R. G. Douglas, and J. E. Bennett, eds.), Churchill Livingstone, New York.
4. D. A. Stevens (1990) *Coccidioides immitis*. In *Principles and Practice of Infectious Diseases* (G. L. Mandell, R. G. B. Douglas, and J. E. Bennett, eds.), Churchill Livingstone, New York, pp. 2008–2017.
5. S. Gao and D. L. Nuss (1996) Proc. Natl. Acad. Sci. USA **93**, 14122–14127.
6. F. Banuett and I. Herskowitz (1994) Genes Dev. **8**, 1367–1378.
7. H.-J. Lo et al. (1997) Cell **90**, 939–949.
8. S. Gold, G. Duncan, K. Barrett, and J. W. Kronstad (1994), Genes Dev. **8**, 2805–2816.
9. K. J. Kwon-Chung, D. G. Lehman, and C. P. T. Magee (1985) Infect. Immun. **49**, 571–575.
10. J. W. Rippon and G. L. Peck (1967) J. Invest. Dermatol. **50**, 54–58.
11. G. San-Blas and D. Vernet (1977) Infect. Immun. **15**, 897–902.
12. H. J. Watts, F. S. H. Cheah, V. Hube, D. Sanglard, and N. A. R. Gow (1998) FEMS Micro. Lett. **159**, 129–135.
13. P. Bowyer, B. R. Clarke, P. Lunness, M. J. Daviels, and A. E. Osbourn (1995) Science **267**, 371–374.
14. K.-M. Weltring, B. G. Turgeon, O. C. Yoder, and H. D. VanEtten (1988) Gene **68**, 335–344.
15. R. J. Howard and B. Valent (1996) Annu. Rev. Microbiol. **50**, 491–512.
16. W. A. Fonzi and M. Y. Irwin (1993) Genetics **134**, 717–728.
17. J. R. Perfect, D. L. Toffaletti, and T. H. Rude (1993) Infect. Immun. **61**, 4446–4451.
18. D. G. Panaccione (1993) Trends Microbiol. **1**, 14–20.
19. C. E. Bulawa, D. W. Miller, L. K. Henry, and J. M. Becker (1995) Proc. Natl. Acad. Sci. USA **92**, 10570–10574.
20. R. A. Cox, L. R. Mills, G. K. Best, and J. F. Denton (1972) Infect. Immun. **5**, 449–453.
21. T. Spellig, B. M. , F. Lottspeich, R. W. Frank, and R. Kahmann (1994) EMBO J. **13**, 1620–1627.
22. J.-R. Xu and J. E. Hamer (1996) Genes Dev. **10**, 2696–2706.
23. J. W. Kronstad (1997) Trends Plant Sci. **2**, 193–199.
24. L. H. Hogan, B. S. Klein, and S. M. Levitz (1996) Clin. Microbiol. Rev. **9**, 469–488.
25. E. Alani, L. Cao, and N. Kleckner (1987) Genetics **116**, 541–545.
26. L. Yaar, M. Mevarech, and Y. Koltin (1997) Microbiology **143**, 3033–3044.
27. M. A. Stringer, R. A. Dean, T. C. Sewall, and W. E. Timberlake (1991) Genes Dev. **5**, 1161–1171.
28. P. Bowyer, B. R. Clarke, P. Lunness, M. J. Daniels, and A. E. Osbourn (1995) Science **267**, 371–374.
29. Y. C. Chang and K. J. Kwon-Chumng (1994) Mol. Cell. Biol. **14**, 4912–4919.
30. M. Barki, Y. Koltin, M. Yanko, A. Tamarkin, and M. Rosenberg (1993) J. Bacteriol. **175**, 5683–5689.
31. N. J. Talbot, D. J. Ebbolle, and J. E. Hamer (1993) Plant Cell **5**, 1575–1590.
32. R. Bohlmann, F. Schauwecker, C. Basse, and R. Kahmann (1994) Plant Microbe Interact. **3**, 239–245.
33. B. Cormack and F. S. (1995) J. Cell. Biochem. (Suppl.) **19B**, abstract B4-202, 156.

34. L. Zhang, A. C. L. Churchill, P. Kazmierczak, D. H. Kim, and N. K. Van Alfen (1993) *Mol. Cell. Biol.* **13**, 7782–7792.
35. S. W. Lu et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 12649–12653.
36. M. Bolker, H. U. Bohnert, K. H. Braun, J. Gohl, and R. Kahmann (1995) *Mol. Gen. Genet.* **248**, 547–552.

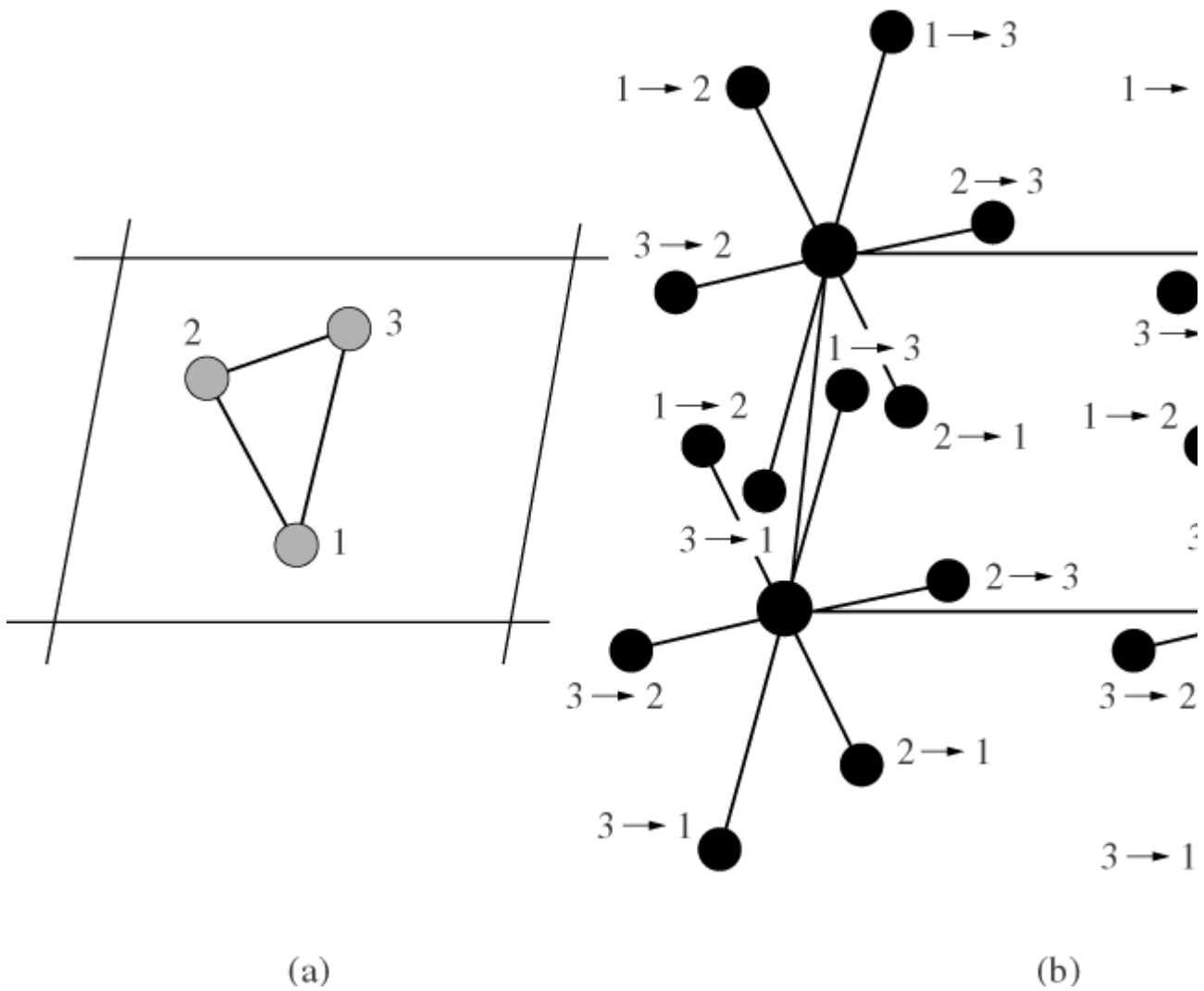
## Patterson Map

Patterson maps give structural information in [X-ray crystallography](#), even if the phases of the reflections are not known (see [Phase Problem](#)). With the Patterson function  $P(uvw)$ , a vector map is calculated. Vectors between atoms in the real structure show up as vectors from the origin to maxima in the Patterson map (Fig. 1). The size and shape of the [unit cell](#) in this map is identical to the crystal unit cell.

$$P(uvw) = \frac{1}{V} \sum_h \sum_k \sum_\ell |F(hk\ell)|^2 \cos[2\pi(hu + kv + \ell w)] \quad (1)$$

where  $u$ ,  $v$ , and  $w$  are relative coordinates in the Patterson cell;  $h$ ,  $k$ , and  $\ell$  are the indices of the reflections;  $V$  is the volume of the unit cell; and  $|F(hk\ell)|$  is the amplitude of the [structure factor](#) of reflection  $h k \ell$ . Calculation of the Patterson function from Eq. 1 does not require knowledge of the phase angles of the reflections and can always be calculated from the experimental diffraction data because, apart from correction factors,  $|F(hk\ell)| = \sqrt{I(hk\ell)}$ , where  $I(h k \ell)$  is the intensity of reflection  $h k \ell$ .

**Figure 1.** Example of a simple Patterson map. (a) A two-dimensional unit cell with three atoms. (b) The corresponding I



For very simple chemical compounds, it is possible to derive the real structure from the Patterson map. A protein Patterson map, however, has too many overlapping vectors to extract any structural information. However, protein Patterson maps are extremely useful for understanding the [molecular replacement](#) technique. They are also employed in [isomorphous replacement](#) to locate the attached heavy atoms. The difference between the amplitudes of the reflections for the native and heavy atom crystals is used in the calculation, to produce a difference Patterson map. Such a map should result only from the heavy atoms, and it should be interpretable because the number of heavy atoms attached is usually small and therefore they form a simple structure.

#### Suggestion for Further Reading

J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

#### Pax Genes

Nine *Pax* genes have been isolated from vertebrates through sequence homology to a region in the



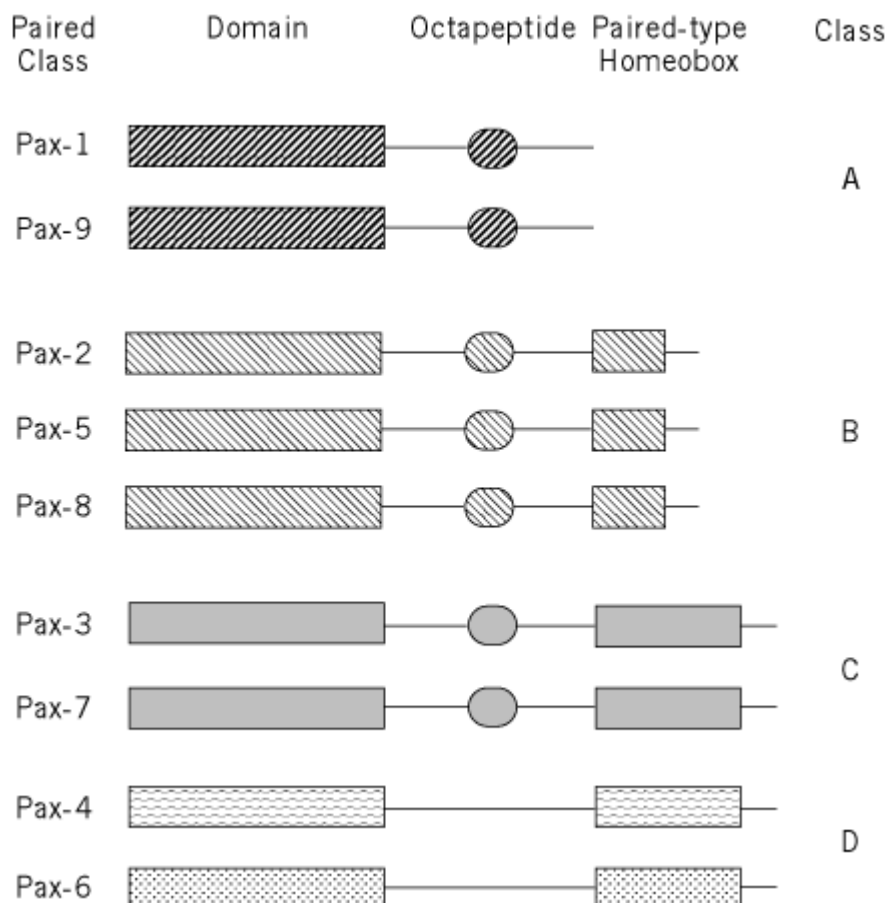
*Drosophila* gene *paired*, which codes for a DNA-binding domain. This domain of 126–amino acid residues is highly conserved and is known as the paired box. Binding sites for Pax proteins have been found in the promoter regions of genes coding for proteins as functionally varied as other transcription factors, crystallins, receptors, peptide hormones and cell adhesion proteins. Pax proteins are expressed in restricted regions during embryonic development. *Pax* genes are subdivided on the basis of the presence or absence of two further conserved sequences: a paired-type homeobox and an octapeptide coding region. These two elements are keys to the functions of Pax proteins in morphogenetic processes such as (1) establishing the competence of tissues to particular developmental fates, (2) commitment of cells to particular developmental fates, and (3) the regionalization of the central nervous system.

## 1. Structure and Biological Activity of Pax Proteins

### 1.1. Classes of Pax Proteins

Pax proteins are transcription factors that regulate gene transcription. All nine Pax proteins contain the paired box, a 136–amino acid residue domain with DNA-binding activity that has sequence similarity to a domain found in the *Drosophila* protein *paired* (1) (reviewed in (2)). As well as the paired domain, Pax proteins also carry one or both of two other conserved domains: the paired-type homeodomain and the octapeptide sequence (Fig. 1). The presence or absence of these other conserved domains is used to classify the Pax proteins into four groups (Fig. 1). The amino acid sequences of the conserved domains of a Pax protein are more similar to those within the same class than those of Pax proteins from different classes (2). Additionally, when Pax proteins have similar functions in a particular tissue, for example, Pax-2 and Pax-5 in the midbrain/hindbrain boundary, they come from the same class.

**Figure 1.** Structure of Pax proteins.



## 1.2. Pax Proteins Directly Regulate the Expression of Many Types of Genes

Pax proteins regulate gene transcription by binding to specific sequences in the enhancer and promoter regions upstream of the coding regions of the gene. This results in the initiation of transcription, which subsequently results in synthesis of the protein. The paired domain binds to DNA; thus it is crucial to the regulation of gene transcription. Changes to the amino acid sequence in this domain have been demonstrated to strongly inhibit the binding of Pax proteins to DNA (3, 4). The specific DNA sequences that Pax proteins bind to have been found in the enhancer and promoter regions of genes coding for many classes of proteins (Table 1). This is consistent with the fact that Pax genes are important regulators of organogenesis in the developing embryo.

**Table 1. Types of Genes Whose Expression is Regulated by Pax Proteins**

| <b>Gene</b>                    | <b>Protein Product Type</b>    | <b>Promoter Bound By</b> | <b>Reference</b> |
|--------------------------------|--------------------------------|--------------------------|------------------|
| <i>En1</i>                     | Transcription Factor           | Pax-2, -5, -8            | (105)            |
| <i>En2</i>                     | Transcription Factor           | Pax-2, -5, -8            | (105)            |
| <i>Crystallin</i><br>(various) | Structural Protein             | Pax-6                    | (106)            |
| <i>NCAM</i>                    | Cell Adhesion Protein          | Pax-6                    | (107)            |
| <i>Insulin</i>                 | Peptide Hormone                | Pax-6                    | (69)             |
| <i>Glucagon</i>                | Peptide Hormone                | Pax-6                    | (69)             |
| <i>Thyroxine peroxidase</i>    | Enzyme                         | Pax-8                    | (84)             |
| <i>IgH</i>                     | Immunoglobulin Subunit         | Pax-5                    | (79)             |
| <i>CD19</i>                    | Receptor                       | Pax-5                    | (76)             |
| <i>c-met</i>                   | Receptor                       | Pax-3                    | (60)             |
| <i>p53</i>                     | Cell Cycle/Apoptosis Regulator | Pax-5                    | (101)            |

## 2. Developing Processes Regulated by Pax Proteins

Pax proteins have important roles in regulating organogenesis, commitment of cells to particular developmental pathways, and the regionalization of the central nervous system (summarized in Table 2). Pax proteins help mediate these processes by (1) enabling an embryonic tissue to be competent to interact with other tissues, (2) conferring regional identity to a portion of tissue that later forms a new embryonic structure, (3) commitment of cells to a lineage with a particular differentiation pathway, (4) maintenance of boundaries between regions that have separate developmental fates, and (5) conferring multipotency to precursor cells (15, 85).

**Table 2. Summary of Pax Gene Null Mutation Phenotypes in Mouse Mutants**

| Gene               | Null Mutation  | Heterozygous Phenotype                   | Homozygous Phenotype  | Reference   |
|--------------------|--|--|---|---|
| <i>Pax-1</i>       | <i>Undulated</i> , targeted                                    | Normal                                   | Hypoplasia of vertebral bodies and discs  | ( <a href="#">52</a> , <a href="#">53</a> , <a href="#">108</a> )   |
| <i>Pax-2</i>       | <i>KrD</i> , targeted  | Small kidneys, optic coloboma            | No kidneys or genital tracts, optic oloboma   | ( <a href="#">25</a> , <a href="#">26</a> , <a href="#">28</a> , <a href="#">30</a> , <a href="#">31</a> , <a href="#">73</a> )   |
| <i>Pax-3</i>       | <i>Splotch</i>   | Lack of pigmentation                     | No limb muscle, neural crest derivatives absent or hypoplastic, persistent truncus arteriosus | ( <a href="#">45</a> , <a href="#">58-64</a> , <a href="#">22-60</a> , <a href="#">62</a> , <a href="#">109</a> , <a href="#">110</a> )                                       |
| <i>Pax-4</i>       | Targeted   | Normal                                   |   | ( <a href="#">66</a> )  |
| <i>Pax-5</i>       | Targeted   | Normal                                   | No B cells, cerebellar and inferior colliculus defects  | ( <a href="#">26-28</a> , <a href="#">77</a> , <a href="#">78</a> )   |
| <i>Pax-6</i>       | <i>Sey</i> , <i>Sey<sup>Neu</sup></i> , <i>rSey</i> , targeted | Microphthalmia, corneal and iris defects | No eyes or olfactory bulbs, forebrain and motoneuron/interneuron specification defects        | ( <a href="#">7-9</a> , <a href="#">14</a> , <a href="#">15</a> , <a href="#">30</a> , <a href="#">31</a> , <a href="#">35-38</a> , <a href="#">47</a> , <a href="#">67</a> ) |
| <i>Pax-7</i>       | Targeted   | Normal                                   | No muscle satellite cells, maxillary process, nasal capsule                                   | ( <a href="#">22</a> , <a href="#">46</a> , <a href="#">85</a> )  |
| <i>Pax-8</i>       | Targeted   | Normal                                   | Absence of thyroid  | ( <a href="#">83</a> )  |
| <i>Pax-9</i>       | Targeted   | Normal                                   | Absence of teeth  | ( <a href="#">55</a> , <a href="#">81</a> )   |
| <i>Pax-2/Pax-5</i> | Double homozygous  | —  | Absence of isthmic organizer, cerebellum and inferior colliculus                              | ( <a href="#">26</a> , <a href="#">28</a> )   |
| <i>Pax-3/Pax-7</i> | Double homozygous  | —  | Neuronal specification defects (somitic derivatives not reported)                             | ( <a href="#">22</a> )  |
| <i>Pax-1/Pax-9</i> | Compound heterozygous, double homozygous                       | —  | Absence of vertebral bodies and discs, teeth  | ( <a href="#">51</a> )  |

The roles of Pax proteins in these processes have been studied in mice having natural or targeted mutation of the *Pax* genes (Table 2). One property of *Pax* genes is that, for some developmental processes, two fully functional copies of the genes are required for normal differentiation of the tissue. Thus, developmental abnormalities in the kidney, optic nerve, eye and pigmentation are found in mice heterozygous for mutations in particular Pax genes. These defects are broadly consistent with human syndromes in which one copy of the corresponding *PAX* genes are mutated. This characteristic, known as haploinsufficiency, results in autosomally-dominant phenotypes that are not the result of novel gene function but, rather, the consequence of depleted levels of functional protein. Increased levels of Pax protein can also alter developmental fate; extra *Pax-6* genes results in developmental defects of the eye (5).

## 2.1. Pax-6 is Necessary for Early Eye Development

The embryonic eye is formed from an extension of the developing forebrain known as the optic vesicle, which is the precursor tissue for the retina, and a portion of the ectoderm, the lens ectoderm, that overlies this optic vesicle (reviewed (6)). A crucial early process of eye development is the induction of the lens ectoderm by the optic vesicle to form firstly a lens placode, then a vesicle and finally the mature lens. In animals homozygous for the natural Pax-6 null mutations *Sey* (7-9), *Sey<sup>Neu</sup>* (9) and rat r*Sey* (10), lens development fails to occur. Pax-6 is expressed in the presumptive lens ectoderm and in the induced lens placode (11, 12). This pattern of expression strongly suggested that Pax-6 proteins is part of the mechanism that allows the ectoderm to respond to a lens induction signal from the optic vesicle. This notion was supported by explant recombination experiments between homozygous r*Sey* and wild-type tissue which showed that the homozygous ectoderm failed to respond to optic vesicle interactions but that homozygous mutant optic vesicle could still induce lens ectoderm (13). A conclusive study of this interaction employed conditional knockout in mice to specifically ablate functional *Pax-6* genes in the lens while leaving fully functional *Pax-6* in the optic vesicle and other tissues (14). Lens ectoderm with no functional Pax-6 underwent an initial response resulting in increased expression of the transcription factor Sox-2, then its development arrested. Although optic vesicle development is arrested in the natural *Pax-6* mutants, optic vesicles in the conditional knockout mice underwent development to form not one but several retinas. Thus, the optic vesicle is also required for correct placement of the retina, but not further development of the retina.

Development of the optic vesicle into the retina is characterized by the sequential generation of mature retinal cell types from stem cells in the original neuroectodermal layer (reviewed in (6)). An important feature of this process is the prevention of differentiation of all the stem cells into differentiated cells, so that a stem cell population is retained for the generation of later appearing cell types. However, specific ablation of *Pax-6* in the retina by conditional knockout, results in all the stem cells of the retina differentiating into the first cell type to appear, the amacrine cells, and the consequent arrest of further retinal development (15). Thus, *Pax-6* is necessary for maintenance of this stem cell population.

## 2.2. Pax Proteins and Patterning in the Central Nervous System

Pax proteins are one of many classes of regulator proteins that control the generation of neural tissue. They are part of the mechanism that specify certain neuronal populations. Moreover, some structures in the CNS require Pax proteins for their survival. Pax proteins also modulate the extracellular environment of certain structures resulting in well-demarcated boundaries with neighboring tissues.

The neural tube and hindbrain exhibit particularly striking patterns of Pax gene expression. *Pax-6* is expressed in most of the neural tube but is missing in the most dorsal and ventral margins (11, 16). *Pax-3* and *Pax-7* are expressed in the dorsal neural tube (17, 18). Combinations of transcriptional regulators including the Pax proteins are specified in individual cells on the basis of positional information provided by opposing concentration gradients of secreted peptide factors (reviewed in (19)). Sonic hedgehog, is secreted from the ventral margin of the neural tube, known as the floorplate, and from the notochord while BMP proteins diffuse from the dorsal margin of the neural tube (the roofplate) and the overlying ectoderm is the second group of signals. Additionally, transcriptional regulator expression can be modulated by other transcriptional regulators. Gain or loss of function for individual transcriptional regulators can change the specification of restricted neuronal types. Deletion of functional *Pax-6*, which is normally expressed in the middle of the neural tube, results in the respecification of the motor neurons supplying the hypoglossal nerve and some ventral interneurons (16, 20). Inappropriate expression of Pax-6 in the ventral neural tube respecifies a population of Nkx-2.2-expressing neurons by repressing expression of the transcriptional regulator (21). Loss of functional Pax-3 and Pax-7 genes results in one population of neurons inappropriately expressing the transcriptional regulator En-1 and the respecification of some ventral commissure neurons to other fates (22). Ectopic expression of Pax-3 in the ventral neural tube prevents the formation of the floorplate (23).

Pax-2 and Pax-5 are required for the maintenance of the isthmus organizer, which forms the

boundary between the midbrain and hindbrain. This structure is required for the development of the inferior colliculus and the cerebellum (reviewed in (24)). Development of regions of the optic tectum and cerebellum were hypoplastic in both Pax-2 (25, 26) and Pax-5 null (27) mutant homozygous mouse embryos. These results suggested that Pax-2 and Pax-5 had similar, but not redundant, functions in the maintenance of the isthmus organizer. Generation of Pax-2/Pax-5 double null mutant homozygotes led to the complete absence of the isthmus organizer, colliculus and cerebellum (26, 28, 29) and the extension of tissues adjoining the missing structures into the consequent space.

Pax proteins have an instructive role in the establishment of the boundary between the developing optic nerve tract and the retina. Both these structures form from the optic vesicle. The more medial portion, the optic stalk, is the optic nerve tract anlage and the more lateral, the optic cup, forms the retina. The developing optic nerve tract is characterized by expression of Pax-2 (25, 30). Absence of functional *Pax-2*, results in the extension of *Pax-6* expression into the more medial portion of the optic vesicle and ultimately the generation of retinal epithelial cells (25, 30, 31). When functional *Pax-6* is missing, the *Pax-2* expression domain is shifted into the lateral optic vesicle (30). Furthermore, binding of Pax-6 protein to the *Pax-2* promoter represses transactivation and binding of Pax-2 protein to *Pax-6* promoter also represses transactivation (30).

In the forebrain, Pax-6 has been shown to establish the correct boundaries between the structures in the forebrain and between the forebrain and the midbrain. In Pax-6 null mutant homozygotes, the caudal forebrain develops markers characteristic of the midbrain (32-34). The dorso-anterior aspect of the forebrain, the telencephalon, develops into the cerebral cortex and the striatum. Lack of functional Pax-6 leads to alterations in its boundary with the thalamus (35), defects in the developing cortical-striatal boundary (36, 37) and altered radial glia differentiation (38). The telencephalon phenotype is, in part, a function of modulation of the levels of the extracellular matrix proteins R-cadherin and tenascin (36) resulting in the loss of a well demarcated boundary with neighboring tissues.

### 2.3. Pax Genes Regulate Migration of Neural Crest Cells

The neural crest lays a crucial role in vertebrate development. Structures derived from neural crest cells include many peripheral neurons, including those in the spinal ganglia (dorsal root, sensory and autonomic), some of the bones of the skull, elements of the inner ear, cardiovascular structures and melanocytes (the pigmented cells) (reviewed in (39)). Neural crest cells migrate from the dorsal region of the neural tube and brain vesicles just after the fusion of their neural folds. This migratory process is influenced by three Pax proteins: Pax-3, Pax-6 and Pax-7. The migration of neural crest cells is blocked at thoracic and lumbar levels in homozygous splotch Pax-3 null mutant embryos and is severely impaired in cervical regions (40-42). This results in small or absent spinal ganglia, deficient ensheathing of peripheral nerves and the absence of melanocytes (43). Lack of neural crest migration into cardiac tissue results in the failure of the separation of the aorta and the pulmonary artery into two separate blood vessels (44, 45). Pax-7 and Pax-6 play similar roles to Pax-3 in the generation of more anterior neural crest-derived craniofacial structures. The maxillary process and nasal capsule are missing in Pax-7 null homozygotes (46) and Pax-6 null allele homozygotes lack lateral nasal processes (47).

### 2.4. Pax-1 and Pax-9 and Vertebral Body Formation

The vertebral column is derived from the sclerotome compartment of somites (48). *Pax-1* and *Pax-9* are expressed in the sclerotomes of the differentiated somite (1, 49), as a result of sonic hedgehog protein signaling from the notochord and the neural tube (50). Some *Pax-1* expressing cells then migrate ventromedially to surround the notochord and differentiate to form the vertebral bodies and intervertebral discs (1, 49). The cells fated to form the vertebral bodies condense to form cartilage in a process known as chondrogenesis as evidenced by *Sox-9* and *collagen II* expression (51). It is the fate of this cell population that Pax-1 and Pax-9 regulate. Mouse embryos homozygous for the natural null *Pax-1* mutation *undulated* (52, 53) and the targeted *Pax-1* null allele (54) exhibit vertebral body hypoplasia and fusions due to inhibited chondrogenesis and agenesis of the intervertebral discs. Hypoplasia is particularly pronounced in the posterior portion of each vertebra,

and heterozygotes for *Pax-1* null alleles have hypoplastic cervical and lumbar vertebrae (54). When embryos are homozygous for targeted null alleles of both *Pax-1* and *Pax-9*, chondrogenesis does not occur because decreased proliferation and increased apoptosis leads to the disappearance of the vertebral body fated cell population (51). However, homozygotes for the *Pax-9* targeted null allele have a normal vertebral column (55). These studies show that *Pax-1* and *Pax-9* have partially redundant roles in the vertebral bodies, but that *Pax-1* acts uniquely in intervertebral disc formation and in normal chondrogenesis of the posterior of each vertebra.

#### 2.5. Pax-3 Requirement for the Migration of Limb Muscle Precursor Cells

The skeletal muscles of the trunk and limbs are derived from the dermomyotome of the somatic mesoderm (reviewed in (48)). The myotome, which develops from the dermomyotome, forms the muscles of the back and ribcage. Precursors of the limb muscle migrate from the leading lateral edge of the dermomyotome into the developing limb bud and coalesce into dorsal and ventral pre-muscle masses. These pre-muscle masses then form the mature musculature. Both *Pax-3* (56, 57) and *Pax-7* (18, 46) are expressed in the dermomyotome with *Pax-3* expressed more strongly in the lateral portion and *Pax-7* in the medial portion. Loss of functional *Pax-3* results in the failure of limb muscle precursors to make the transition from epithelium to mesenchyme and migrate into the limb (58, 59). This failure in migration is due in part because *Pax-3* transactivation is required for the expression of the SF/HGF receptor *c-met* (60-62). SF/HGF protein is expressed in the limb buds and limb muscle precursor cells do not migrate in embryos homozygous for the null allele for *c-met* (63). Additionally, apoptotic cells are seen in the somites of mice homozygous for the *Pax-3* null allele *spotch* at the developmental stages that limb muscle precursors emigrate from the somite (64).

#### 2.6. Pax-4 and Pax-6 Mediate Endocrine Cell Specification in the Pancreas

The endocrine cell populations of the pancreas are found in the Islets of Langerhans (reviewed in (65)). The a, b, d and PP cell lineages synthesize glucagon, insulin, somatostatin and pancreatic polypeptide, respectively. *Pax-4* is expressed in the b (insulin-producing) and d (somatostatin-producing) cells (66), and *Pax-6* is expressed in all four lineages (67). Inactivation of both copies of the *Pax-4* gene resulted in the failure of insulin- and somatostatin-producing cell populations to differentiate and resulted in increased numbers of glucagons-producing cells (66). The latter finding is consistent with the discovery that *Pax-4* protein binding to the enhancer of the glucagon gene represses transcription (68). *Pax-6* homozygous null mutants lack glucagon-expressing cells and fail to form islets of Langerhans (67, 69). Moreover, *Pax-6* protein transactivates the glucagon and insulin promoters and also binds to the somatostatin promoter (69). In embryos lacking both functional *Pax-4* and *Pax-6* genes, endocrine pancreatic cells fail to develop altogether (67). Hence the actions of *Pax-4* and *Pax-6* are required to generate the full complement of endocrine pancreas cell types from a single endocrine pancreatic progenitor population.

#### 2.7. Pax-2 is Required for Mesenchymal-to-Epithelial Transitions in the Urogenital System

The development of the definitive kidney commences when a mesenchymal structure, known as the metanephric mesenchyme, induces the ureter to bud out from the caudal Wolffian duct (reviewed in (70)). Contact by the ureter induces the metanephric mesenchyme to differentiate into the glomerulae which are the proximal and distal tubular epithelial structures seen in the nephrons of the kidney. The metanephric mesenchyme expresses *Pax-2* (71, 72). *Pax-2* is necessary to allow the metanephric mesenchyme to undergo mesenchymal-to-epithelial transitions since mice that are homozygous for either the targeted *Pax-2* null mutation (73) or the natural *KrD* mutation (31). Moreover, genital structures of epithelial morphology such as the vas deferens and seminal vesicles in the male and the uterus and vagina in the female are also missing in mice homozygous for these mutations. *Pax-2's* role in kidney formation is dosage dependent, kidneys in heterozygous mice are hypoplastic.

#### 2.8. Pax-5 and the B-Cell Lineage

The B-cell lineage is one of the two cell types that mediate the immune response. *Pax-5* is expressed in B cells until terminal differentiation (74) (reviewed in (75)). In *Pax-5* null allele homozygotes, differentiation is arrested at a very early stage prior to the initiation of expression of the early marker, CD19 (27). In fact, the *CD19* gene's promoter contains *Pax-5* binding sites (76). However,

the *Pax-5* (-/-) pre B cells can populate the thymus when injected into *Rag-2* (-/-) mutants (77), and can be diverted to other fates in vitro (78). *Pax-5* is thus required for the generation and maintenance of an immature B-cell population. The finding that *Pax-5* represses expression from the *Immunoglobulin Heavy Chain* gene (*IgH*) enhancer (79) suggests that downregulation of *Pax-5* expression may be an important part of terminal differentiation.

## 2.9. Roles of Pax Proteins in Other Tissues

*Pax* genes have been found, from developmental defects observed in null-mutant homozygote embryos, to participate in the regulation of other developmental processes that are less studied. Nevertheless, the roles *Pax* genes play here are similar to the others described above.

*Pax-1* and *Pax-9* are required for full development of tissues derived from the third and fourth pair of pharyngeal pouches such as the thymus, parathyroid gland and the ultimobranchial bodies. These structures are completely missing from *Pax-9* null allele homozygotes (55) and mice with no functional *Pax-1* have a reduced thymus with 50 to 80% fewer T-cells (80).

*Pax-9* plays a crucial role in the generation of teeth at the tooth germ stage. It is expressed at discrete sites in the mesenchyme along the developing mandible due to inductive signaling by FGF-8 and antagonistic signaling by BMP-2 and -4 (81). Only *Pax-9* expressing portions of the mesenchyme can initiate tooth formation (81) in an in vitro assay; its absence in *Pax-9* null mutant homozygotes prevents further development of tooth primordia (55).

The thyroid gland is formed from two cell populations: endodermal cells that express *Pax-8* and pharyngeal-derived mesenchymal cells that form the parathyroid (82). Ablation of functional *Pax-8* prevents the endodermal cells participating in thyroid development resulting in thyroid agenesis (83). Moreover, *Pax-8* protein can transactivate the promoters of thyroglobulin and thyroperoxidase (84).

The olfactory bulbs differentiate from the nasal placode, which is induced in a *Pax-6* dependent process. These structures are missing in *Pax-6* null allele homozygotes (7-9, 47).

Finally, satellite cells are missing in the muscle tissue of *Pax-7* null mutant homozygotes (85). However, the muscle stem cell population is unaffected.

## 3. Pax Genes and Human Disease

### 3.1. Human Syndromes Resulting from Mutation of Pax Genes

*PAX* gene mutations have been implicated in a number of syndromes over the last decade. The symptoms of these conditions are broadly consistent with the observations in heterozygous null mutant mice.

Two syndromes are known to result from mutations at the *Pax-6* locus: Aniridia and Peter's anomaly. Aniridia is a syndrome in which the iris fails to develop fully (86-88). This is consistent with the iris malformations in heterozygous Sey mice (7, 8). Significantly, the posterior layer of the iris is derived from the optic cup, which expresses *Pax-6*. Peter's anomaly is a defect in the formation of the cornea that renders it opaque (89). The cornea develops from ectodermal tissue that expresses *Pax-6* and forms over the lens after its formation.

The symptoms of Waardenburg's Syndrome 1 are pigmentation defects and hearing loss, which occurs as a result of *PAX-3* mutation (90, 91). These defects occur as a result of the impaired migration or differentiation of neural crest cells, which give rise to structures of the inner ear and to the melanocytes. Sufferers of this syndrome are deaf, however mice heterozygous for the *Spotch* null mutation of *Pax-3* are not (91).

In humans, disruption of *PAX-2* function leads to optic nerve colobomass, renal hypoplasia and vesicoureteral reflux (92, 93). This syndrome is consistent with the optic nerve and kidney phenotypes

in mice heterozygous for the targeted Pax-2 null mutation (73).

### 3.2. Pax Genes and Cancer

Transcription factors that are inappropriately expressed or mutated play a significant role in tumorigenesis. Some studies have identified *PAX* genes as tumor promoters in certain cancers. The potential for *PAX* genes to promote neoplasia was shown in a study where *Pax-1,-2,-3* and *-6* transformed 3T3 cells could form tumors when injected subcutaneously into nude mice (94). Moreover, the cells transformed in focus assays.

Chromosome translocations leading to fused genes that contain some PAX sequences is one way by which *PAX* genes may promote tumorigenesis. Rhabdomyosarcomas are soft-tissue tumors with a muscle morphology found in the young. Characteristic of this sarcoma is a translocation from chromosome 13 to the *PAX-3* locus in chromosome 2 GALILI, (95). This results in a fusion of sequences of a forkhead-class gene to the homeobox of *PAX-3*, which expresses a chimeric protein that can transactivate gene expression via PAX-3 binding sites in their enhancers (96), (97)s. More recently a *PAX-7* to *forkhead* translocation was isolated from a rhabdomyosarcoma (98). The effect of a translocation like this may put expression of the chimeric gene under a different, looser regulation. Alternatively, the resultant fusion proteins may simply have better transactivator potential than their parent proteins. A recent study found that the PAX-3/forkhead fusion protein was unresponsive to the PAX-3 repressor protein, hDaxx (99).

In other cases, the genetic lesion does not change the coding sequences of the PAX gene but only the promoter-enhancer region. An example of this is the translocation of the Em enhancer of the *IgH* gene to the *PAX-5* locus in diffuse large-cell lymphomas (100). PAX-5 has been shown to directly repress the transcription of p53 (101) whose protein is crucial for initiation of apoptosis and cell cycle control. Another is the elevated *PAX-2* and *PAX-8* expression in Wilm's tumor (102-104).

## 4. Conclusion

Over the last twelve years, the *Pax* genes have been identified, their patterns of expression characterized, and many of their functions analyzed. Functional analyses have been helped by the fortuitous discovery of natural *Pax* rodent and zebrafish mutants, which complemented the targeted mutational studies. All *Pax* genes display restricted patterns of expression in both neural and nonneural structures. Some *Pax* genes have very similar functional roles in the development of certain tissues. The best example of this is *Pax-2* and *Pax-5* in the morphogenesis of the midbrain-hindbrain boundary. *Pax* protein functions in some morphogenetic processes are sensitive to gene dosage, and the malformations seen in heterozygous mouse mutants and in the corresponding human syndromes resulting from mutation of the correlating PAX gene are strikingly similar. Thus haploinsufficiency is a conserved feature of *Pax* function in different species.

*Pax* proteins play important roles in certain morphogenetic processes. One of these is in conferring competence to respond to inductive influences. Examples include the competence of ectodermal tissue to be induced to form lens or nasal placode tissue, metanephric mesenchyme to form kidney epithelial structures, and the competence of dermomyotome cells to form a migratory myogenic cell type. *Pax* proteins also function in cell specification, examples being the B cells, pancreatic endocrine cells, the neural e developing crest, and certain neuronal cell types. In these cases, *Pax* proteins also help to mediate the survival of these cell types, since null mutations of *Pax* genes also lead to the disappearance of the cell types.

Studies in coming years will identify more structures whose genesis requires the actions of *Pax* genes particularly in the developing brain. Studies will also focus on *Pax* protein interactions with other nuclear proteins in forming transcriptional complexes, which is poorly characterized at present. Yet other studies will investigate the control of the expression of *Pax* genes.

## Bibliography



1. U. Deutsch, G.R Dressler, and P. Gruss, *Cell* **53**, 617–625 (1988).
2. C. Walther et al., *Genomics* **11**, 424–434 (1991).
3. G. Chalepakis et al., *Cell* **66**, 873–884 (1991).
4. K.J Vogan, D.J Epstein, D.G Trasler, and P. Gros, *Genomics* **17**, 364–369 (1993).
5. A. Schedl et al., *Cell* **86**, 71–82 (1996).
6. D. Jean, K. Ewan, and P. Gruss, *Mech. Dev.* **76**, 3–18 (1998).
7. B.L Hogan et al., *J. Embryol Exp. Morphol.* **97**, 95–110 (1986).
8. B.L Hogan, E.M Hirst, G. Horsburgh, and C.M. Hetherington, *Development* **103**, 115–119 (1988).
9. R.E Hill et al., *Nature* **354**, 522–525 (1991).
10. T. Matsuo et al., *Nat. Genet.* **3**, 299–304 (1993).
11. C. Walther and P Gruss, *Development* **113**, 1435–1449 (1991).
12. H.S Li et al., *Dev. Biol.* **162**, 181–194 (1994).
13. M. Fujiwara, T. Uchida, N. Osumi-Yamashita, and K Eto, *Differentiation* **57**, 31–38 (1994).
14. R. Ashery-Padan, T. Marquardt, X. Zhou, and P. Gruss, *Genes Dev.* **14**, 2701–2711 (2000).
15. T. Marquardt et al., *Cell* **105**, 43–55 (2001).
16. J. Ericson et al., *Cell* **90**, 169–180 (1997).
17. M.D Goulding et al., *Embo J.* **10**, 1135–1147 (1991).
18. B. Jostes, C. Walther, and P. Gruss, *Mech. Dev.* **33**, 27–37 (1990).
19. T.M. Jessell, *Nat. Rev. Genet.* **1**, 20–29 (2000).
20. N. Osumi et al., *Development* **124**, 2961–2972 (1997).
21. J. Briscoe, A. Pierani, T.M Jessell, and J. Ericson, *Cell* **101**, 435–445 (2000).
22. A. Mansouri and P. Gruss, *Mech. Dev.* **78**, 171–178 (1998).
23. P. Tremblay, F. Pituello, and P. Gruss, *Development* **122**, 2555–2567 (1996).
24. H Nakamura, *Trends Neurosci.* **24**, 32–39 (2001).
25. M. Torres, E. Gomez-Pardo, and P. Gruss, *Development* **122**, 3381–3391 (1996).
26. M. Schwarz et al., *Proc. Natl. Acad. Sci. U.S.A.* **94**, 14518–14523 (1997).
27. P. Urbanek et al., *Cell* **79**, 901–912 (1994).
28. P. Urbanek, I. Fetka, M.H Meisler, and M. Busslinger, *Proc. Natl. Acad. Sci. U.S.A.* **94**, 5703–5708 (1997).
29. M. Schwarz et al., *Mech. Dev.* **82**, 29–39 (1999).
30. M. Schwarz et al., *Development* **127**, 4325–4334 (2000).
31. S.A. Keller et al., *Genomics* **23**, 309–320 (1994).
32. J.C. Grindley et al., *Mech. Dev.* **64**, 111–126 (1997).
33. G.S. Mastick, N.M. Davis, G.L. Andrew, and S.S. Easter, *Development* **124**, 1985–1997 (1997).
34. N. Warren and D.J. Price, *Development* **124**, 1573–1582 (1997).
35. A. Stoykova, R. Fritsch, C. Walther, and P. Gruss, *Development* **122**, 3453–3465 (1996).
36. A. Stoykova, M. Gotz, P. Gruss, and J. Price, *Development* **124**, 3765–3777 (1997).
37. A. Stoykova, D. Treichel, M. Hallonet, and P Gruss, *J. Neurosci.* **20**, 8042–8050 (2000).
38. M. Gotz, A. Stoykova, and P Gruss, *Neuron* **21**, 1031–1044 (1998).
39. C. LaBonne and M. Bronner-Fraser, *Annu. Rev. Cell Dev. Biol.* **15**, 81–112 (1999).
40. T. Franz and R. Kothary, *Brain Res. Dev. Brain Res.* **72**, 99–105 (1993).
41. T. Franz, *Anat. Embryol. (Berlin)* **187**, 371–377 (1993).
42. G.N. Serbedzija and A.P. McMahon, *Dev. Biol.* **185**, 139–147 (1997).

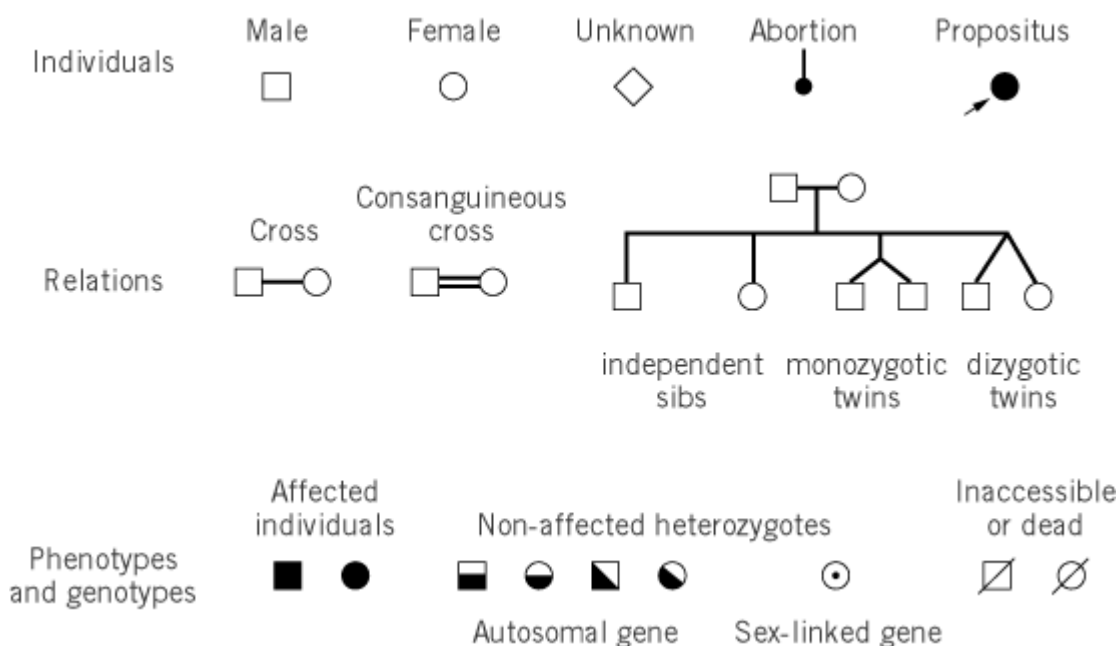
43. T. Franz, *Acta Anat. (Basel)* **138**, 246–253 (1990).
44. T. Franz, *Anat. Embryol. (Berlin)* **180**, 457–464 (1989).
45. S.J. Conway, D.J. Henderson, and A.J. Copp, *Development* **124**, 505–514 (1997).
46. A. Mansouri, A. Stoykova, M. Torres, and P. Gruss, *Development* **122**, 831–838 (1996).
47. J.C. Grindley, D.R Davidson, and R.E. Hill, *Development* **121**, 1433–1442 (1995).
48. B. Brand-Saberi and B. Christ, *Curr. Top. Dev. Biol.* **48**, 1–42 (2000).
49. C. Ebensperger et al., *Anat. Embryol. (Berl)* **191**, 297–310 (1995).
50. A.G Borycki, L. Mendham, and C.P. Emerson, *Development* **125**, 777–790 (1998).
51. H. Peters et al., *Development* **126**, 5399–5408 (1999).
52. R. Balling, U. Deutsch, and P. Gruss, *Cell* **55**, 531–535 (1988).
53. J. Wallin et al., *Development* **120**, 1109–1121 (1994).
54. B. Wilm et al., *Proc. Natl. Acad. Sci. U.S.A.* **95**, 8692–8697 (1998).
55. H. Peters, A. Neubuser, K. Kratochwil, and R. Balling, *Genes Dev* **12**, 2735–2747 (1998).
56. B.A. Williams and C.P. Ordahl, *Development* **120**, 785–796 (1994).
57. M.D. Goulding, A. Lumsden, and P. Gruss, *Development* **117**, 1001–1016 (1993).
58. E. Bober et al., *Development* **120**, 603–612 (1994).
59. M. Goulding, A. Lumsden, and A.J. Paquette, *Development* **120**, 957–971 (1994).
60. J.A Epstein et al., *Proc. Natl. Acad. Sci. U.S.A.* **93**, 4213–4218 (1996).
61. G. Daston, E. Lamar, M. Olivier, and M. Goulding, *Development* **122**, 1017–1027 (1996).
62. X.M Yang, K. Vogan, P. Gros, and M. Park, *Development* **122**, 2163–2171 (1996).
63. F. Bladt et al., *Nature* **376**, 768–771 (1995).
64. A.G. Borycki et al., *Development* **126**, 1665–1674 (1999).
65. H. Edlund, *Diabetes* **50**, S5–9 (2001).
66. B. Sosa-Pineda et al., *Nature* **386**, 399–402 (1997).
67. L. St-Onge et al., *Nature* **387**, 406–409 (1997).
68. S.B. Smith, H.C. Ee, J.R. Connors, and M.S. German, *Mol. Cell Biol.* **19**, 8272–8280 (1999).
69. M. Sander et al., *Genes Dev.* **11**, 1662–1673 (1997).
70. M. Horster et al., *Pflugers Arch.* **434**, 647–660 (1997).
71. G.R Dressler et al., *Development* **109**, 787–795 (1990).
72. H.O. Nornes et al., *Development* **109**, 797–809 (1990).
73. M. Torres, E. Gomez-Pardo, G.R. Dressler, and P. Gruss, *Development* **121**, 4057–4065 (1995).
74. B. Adams et al., *Genes Dev.* **6**, 1589–1607 (1992).
75. S.L. Nutt, A.G. Rolink, and M. Busslinger, *Cold Spring Harb. Symp. Quant. Biol.* **64**, 51–59 (1999).
76. Z. Kozmik et al., *Mol. Cell Biol.* **12**, 2662–2672 (1992).
77. A.G. Rolink, S.L. Nutt, F. Melchers, and M. Busslinger, *Nature* **401**, 603–606 (1999).
78. S.L Nutt, B. Heavey, A.G. Rolink, and M. Busslinger, *Nature* **401**, 556–562 (1999).
79. M. Singh and B.K Birshstein, *Mol. Cell Biol.* **13**, 3611–3622 (1993).
80. J. Wallin et al., *Development* **122**, 23–30 (1996).
81. A. Neubuser, H. Peters, R. Balling, and G.R Martin, *Cell* **90**, 247–255 (1997).
82. D. Plachov et al., *Development* **110**, 643–651 (1990).
83. A. Mansouri, K. Chowdhury, and P. Gruss, *Nat. Genet.* **19**, 87–90 (1998).
84. M. Zannini, H. Francis-Lang, D. Plachov, and R. Di Lauro, *Mol. Cell Biol.* **12**, 4230–4241 (1992).

85. P. Seale et al., *Cell* **102**, 777–786 (2000).
86. C.C Ton et al., *Cell* **67**, 1059–1074 (1991).
87. T. Glaser, D.S Walton, and R.L Maas, *Nat. Genet.* **2**, 232–239 (1992).
88. T. Jordan et al., *Nat. Genet.* **1**, 328–332 (1992).
89. I.M Hanson et al., *Nat. Genet.* **6**, 168–173 (1994).
90. M. Tassabehji et al., *Nature* **355**, 635–636 (1992).
91. K.P. Steel and R.J. Smith, *Nat. Genet.* **2**, 75–79 (1992).
92. P. Sanyanusin et al., *Nat. Genet.* **9**, 358–364 (1995).
93. L.A. Schimmenti et al., *Am. J. Med. Genet.* **59**, 204–208 (1995).
94. C.C. Maulbecker and P Gruss, *EMBO J.* **12**, 2361–2367 (1993).
95. D.N. Shapiro et al., *Cancer Res.* **53**, 5108–5112 (1993).
96. W.J. Fredericks et al., *Mol. Cell Biol.* **15**, 1522–1535 (1995).
97. J.E. Sublett, I.S. Jeon, and D.N. Shapiro, *Oncogene* **11**, 545–552 (1995).
98. R.J. Davis and F.G. Barr, *Proc. Natl. Acad. Sci. U.S.A.* **94**, 8047–8051 (1997).
99. A.D. Hollenbach, J.E. Sublett, C.J. McPherson, and G. Grosveld, *EMBO J.* **18**, 3702–3711 (1999).
100. M. Busslinger et al., *Proc. Natl. Acad. Sci. U.S.A.* **93**, 6129–6134 (1996).
101. E.T. Stuart, R. Haffner, M. Oren, and P. Gruss, *EMBO J.* **14**, 5638–5645 (1995).
102. G.R. Dressler and E.C. Douglass, *Proc. Natl. Acad. Sci. U.S.A.* **89**, 1179–1183 (1992).
103. M.R. Eccles et al., *Cell Growth Differ.* **3**, 279–289 (1992).
104. A. Poleev et al., *Development* **116**, 611–623 (1992).
105. D.L. Song, G. Chalepakis, P. Gruss, and A.L Joyner, *Development* **122**, 627–635 (1996).
106. A. Cvekl and J. Piatigorsky, *Bioessays* **18**, 621–630 (1996).
107. B.D. Holst, Y. Wang, F.S. Jones, and G.M. Edelman, *Proc. Natl. Acad. Sci. U.S.A.* **94**, 1465–1470 (1997).
108. H. Koseki et al., *Development* **119**, 649–660 (1993).
109. T. Franz et al., *Anat. Embryol. (Berl)* **187**, 153–160 (1993).
110. M. Goulding et al., *Genomics* **17**, 355–363 (1993).

## Pedigree

A pedigree is a family tree, a graphical representation of biological descent (Fig. 1). Each generation is represented in a separate row. Males are represented by squares and females are represented by circles; alternative respective symbols are ♂, representing the shield and spear of Mars, and ♀, the mirror of Venus. Horizontal lines join parents and sets of sibs; vertical lines indicate descent. Sibs are ordered according to birth from left to right. Double horizontal lines call attention to parents that are closely related (inbreeding); sometimes they are used to join monozygotic twins as well. Filled, colored, or shaded symbols indicate affected individuals in pedigrees that represent the inheritance of certain traits; partly modified symbols indicate unaffected **heterozygous** carriers. Several sibs can be pooled together by indicating their number within a single symbol. A pedigree may be described in words, but this is usually impractical.

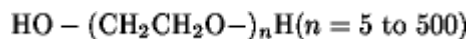
**Figure 1.** Conventional symbols used in the representation of pedigrees.



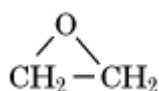
In the analysis of pedigrees, attention must be given to inconstant phenotypes due to effects of the environment or the rest of the genotype of the affected individuals on the expression of the **allele**. Some alleles are unstable and change from one generation to the next (see [Trinucleotide Repeats](#)). Because of the small size of human families, linkage has to be deduced from the pooled results of many groups of sibs (see **Lod scores**).

## PEG (Polyethylene Glycol)

PEG is a polyether polymer commonly used to induce the [precipitation](#) or [crystallization](#) of [macromolecules](#) and [viruses](#). In spite of its name, polyethylene glycol, it has hydroxyl groups only at each end of the polymeric chain:



Except for these ends, the polymer is essentially [nonpolar](#). The ether oxygen atom is only slightly polar and [hydrogen bonds](#) with [water](#) just sufficiently to make PEG water-soluble. PEG is produced by the simple polymerization of ethylene oxide,



and is available in different lengths, whose molecular weights vary from 200 to 20,000, depending on the polymerization process. The average number of monomers in each molecule may be estimated by dividing the molecular weight of the polymer by 44, the molecular weight of the monomer.

PEG precipitates proteins because it is a large, water-soluble polymer that is excluded from the surfaces of macromolecules simply because of its large physical size, a phenomenon known as steric exclusion (see [Excluded Volume](#)). A large polymer cannot penetrate the interior of a folded protein or other macromolecule, so it is effectively excluded from a shell around the macromolecule. The thickness of the shell is defined by the center of the polymer at its point of closest approach to the macromolecule. The presence of the polymer in effect increases the concentration of the macromolecule because the macromolecule is excluded from the volume of the solution occupied by the polymer. When the solubility limit of the macromolecule is reached, it precipitates or crystallizes.

Water molecules are smaller and penetrate within the shell of solvent around the macromolecule from which PEG is excluded. As a result, a zone enriched in water is formed around the protein molecule. The end result is [preferential hydration](#) of the macromolecule and preferential exclusion of the polymer (1). Solely on this basis, PEG would be expected to stabilize the folded conformation of the macromolecule (see [Stabilization And Destabilization By Co-Solvents](#)). However, PEG is a relatively nonpolar molecule that interacts favorably with the nonpolar surfaces exposed when a macromolecule unfolds, and the net effect is that PEG destabilizes folded proteins (1, 2).

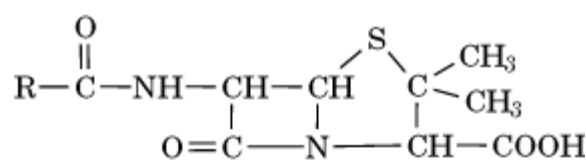
#### Bibliography

1. T. Arakawa and S. N. Timasheff (1985) *Biochemistry* **24**, 6756–6762.
2. G. G. Hammes and P. R. Schimmel (1967) *J. Am. Chem. Soc.* **89**, 442–446.

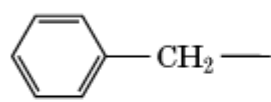
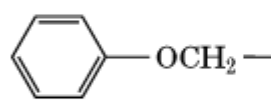
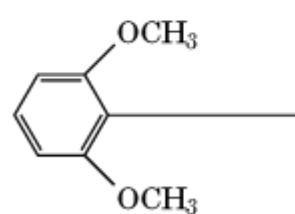
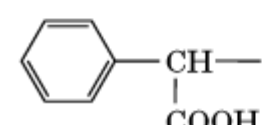
## Penicillin

The penicillins are antibacterial agents belonging to the group known as the  $\beta$ -lactam antibiotics. Penicillin was first discovered about 70 years ago by Alexander Fleming as a metabolic by-product of the fungus *Penicillium notatum* (1). Since then, other fungi have been shown to produce several different structural types of penicillins. All penicillins share a common bicyclic nucleus that is composed of fused  $\beta$ -lactam and thiazolidine rings (Fig. 1). This compound, known as 6-aminopenicillanic acid, is a cyclic [peptide](#) formed by the condensation of L-cysteine and D-valine. An important structural feature of the penicillins is the acyl side-chain attached to 6-aminopenicillanic acid (represented by R in Fig. 1). The different members of the penicillin family are distinguished by their acyl side chains, some examples of which are shown in Figure 1b. The derivatives known as penicillin G and penicillin V are natural penicillins, namely, fungal metabolites. The medical application of penicillin in the 1940s highlighted the birth of the antibiotic era, and subsequent misuse of this and other antibiotics has led to the current [drug resistance](#) crisis (2) (see [Antibiotic Resistance](#)).

**Figure 1.** (a) The structure of the penicillin nucleus; R represents a hydrogen atom in the case of 6-aminopenicillanic acid. (b) Examples of acyl side chains of natural and semisynthetic penicillins.



(a)

| <u>R</u>  | <u>Derivative</u> |
|---|-------------------|
|    | Penicillin G      |
|    | Penicillin V      |
|    | Methicillin       |
|  | Carbenicillin     |

(b)

Comparative studies indicated that the structure of the acyl side chain significantly influenced the antibacterial spectra and physical properties of penicillin. The development of semisynthetic penicillins was an important event in penicillin development. Nayler (3) has presented a historical account of the discovery of semisynthetic penicillins. Semisynthetic antimicrobial agents are derivatives of a natural product that have been deliberately modified chemically. In the case of the semisynthetic penicillins, the natural product is the penicillin nucleus, 6-aminopenicillanic acid. The discovery of methods to produce large amounts of 6-aminopenicillanic acid in the late 1950s was the pivotal event that led to the development of semisynthetic penicillins. The subsequent synthetic chemical acylation of 6-aminopenicillanic acid gave rise to numerous new penicillins with novel properties, and this approach was essential to keep up with the growing demands to solve problems faced by medical applications. For example (see Fig. 1), this technology gave rise to agents that exhibited broad-spectrum antibacterial activity (eg, ampicillin and carbenicillin), agents that could be administered orally (eg, ampicillin), and agents that were resistant to the action of *b-lactamase* (eg, methicillin).

Penicillin inhibits the synthesis of the bacterial cell wall polymer known as peptidoglycan or murein. This topic has been reviewed recently by Holtje (4, 5). Peptidoglycan is a single large macromolecule that completely envelopes bacterial cells. It represents an exoskeleton that supports the underlying cytoplasmic membrane. In this regard, it serves an absolutely essential function in preventing osmotic lysis of the bacterial cells. The actual targets of penicillin are membrane-bound enzymes, known as *penicillin-binding proteins* (PBPs), that catalyze the transpeptidation reaction in

the final stage of peptidoglycan polymerization. Transpeptidation creates crosslinkages between adjacent peptide chains in the peptidoglycan polymer. The amide group in the  $\beta$ -lactam ring of penicillin represents a structural analogue of the D-alanyl-D-alanine peptide bond involved in the transpeptidation reaction. Penicillin is a **suicide substrate** for PBPs, and the interaction of penicillin with its PBP target results in the formation of a stable penicilloylated enzyme derivative. For details of the mechanism of action of penicillin, see [Penicillin-Binding Proteins](#) (see Fig. 1 of that article).

Penicillin is bactericidal. The lethal effect of penicillin treatment is normally associated with bacterial cell lysis (4, 5). This lysis is catalyzed by one or more bacterial enzymes known as peptidoglycan hydrolases. These enzymes cleave chemical bonds in peptidoglycan. Their normal physiological functions are poorly understood. It is believed that the activities of at least some of these hydrolases create openings in the completely closed peptidoglycan molecule, to permit the insertion of new subunits during the growth of bacterial cells. The terminal stage of peptidoglycan synthesis is consequently proposed to be a coordinated process catalyzed by two classes of enzymes: the PBPs and the peptidoglycan hydrolases. These enzymes are further proposed to form a multienzyme complex. Therefore, bacteria are killed by penicillin in a two-step process. In the first step, penicillin inhibits activities of the PBPs involved in peptidoglycan polymerization. In the second step, the peptidoglycan hydrolases become uncoupled from the peptidoglycan biosynthetic process and cause lysis of the bacteria.

It has long been recognized that penicillin kills only actively growing bacteria (6). More recent studies have confirmed that the bactericidal effect of penicillin is directly dependent on the bacterial growth rate (7). The first application of penicillin in genetics was based on this concept. Lederberg and Zinder (8) and Davis (9) independently developed a technique, known as penicillin selection, for the isolation of auxotrophic mutants of *Escherichia coli*. The procedure involves inoculating a sample of *E. coli* containing rare auxotrophic mutants into a minimal medium lacking the nutritional factor(s) required by the mutants. Only the wild-type cells in the population are capable of growing in this medium, and the majority of these are subsequently killed by the addition of a lethal dose of penicillin. The auxotrophic mutants are tolerant to penicillin and are readily recovered by screening the survivors of the penicillin treatment.

The basis for the penicillin tolerance of amino acid-deprived cells has been investigated. Amino acid deprivation of *E. coli* results in the coordinate inhibition of a variety of metabolic activities in a phenomenon that has been called the “stringent response” (10) (see [Stringent Control](#)). The stringent response is probably a strategy to restrict energy consumption in an effort to promote the survival of bacteria during periods of starvation. Amino acid-deprived *E. coli* accumulate a novel nucleotide, guanosine 3',5'-bispyrophosphate (ppGpp), which is believed to mediate the stringent response. Membrane phospholipid synthesis, cell wall peptidoglycan synthesis, and penicillin-induced lysis are among the many metabolic activities that are inhibited during the stringent response (11, 12). Moreover, the activities of the membrane peptidoglycan polymerases and peptidoglycan hydrolases have been shown to be dependent on ongoing phospholipid synthesis. Therefore, the inhibition of phospholipid synthesis by ppGpp results in the inhibition of the membrane-associated reactions in peptidoglycan metabolism and is the direct cause of the penicillin tolerance of starved *E. coli*.

The concept of penicillin selection can be applied more broadly. It can be used in any situation where it is possible to design culture conditions that do not permit the growth of the cells to be selected. For example, penicillin selection can be used to select cells that have been spontaneously cured of [transposon](#) insertions. Therefore, cured cells can be selected from a population of cells that carry the transposon Tn 5 (which encodes [kanamycin](#) resistance) by growing the culture in a medium containing both kanamycin and ampicillin. The ampicillin will kill only those cells that carry Tn 5, and cured cells can be isolated by screening the surviving population.

Genetic elements encoding  $\beta$ -lactamases (see [Penicillin-Binding Proteins](#)) are currently widely employed as selective markers in molecular biological applications. Consequently, penicillin derivatives, especially [ampicillin](#), are among the most common chemicals encountered in a

molecular biology laboratory. Penicillin selection is attractive because the antibiotics are relatively inexpensive and are highly effective because of their bactericidal action. In medical applications, penicillin derivatives are considered to be of low toxicity. They are allergenic, however, and the development of either immediate or delayed hypersensitivities are the most frequent and serious problems associated with penicillin therapy. It is notable that anaphylaxis to penicillin has been reported after nontherapeutic exposure (13). It is therefore conceivable that exposure to penicillin in the laboratory could result in hypersensitive reactions in sensitized individuals, and those who are known to be allergic to penicillin should be alerted to this possibility.

### Bibliography

1. A. Fleming (1929) *Br. J. Exp. Pathol.* **10**, 226–236.
2. S. B. Levy (1992) *The Antibiotic Paradox. How Miracle drugs are Destroying the Miracle*, Plenum Press, New York.
3. J. H. C. Nayler (1991) in *50 Years of Penicillin Application. History and Trends*, H. Kleinkauf and H. von Döhren, eds., Public Ltd., Czech Republic, pp. 64–74.
4. J.-V. Höltje (1995) *Arch. Microbiol.* **164**, 243–254.
5. J.-V. Höltje (1996) *Microbiology* **142**, 1911–1918.
6. G. L. Hobby, K. Meyer, and E. Chaffee (1942) *Proc. Soc. Exp. Biol. Med.* **50**, 281.
7. E. Tuomanen, R. Cozens, W. Tosch, O. Zak, and A. Tomasz (1986) *J. Gen. Microbiol.* **132**, 1297–1304.
8. J. Lederberg and N. Zinder (1948) *J. Am. Chem Soc.* **70**, 467–468.
9. B. D. Davis (1949) *Proc. Natl. Acad. Sci. USA* **35**, 1–10.
10. M. Cashel, D. R. Gentry, V. J. Hernandez, and K. E. Rudd (1996) In *Escherichia coli and Salmonella typhimurium. Cellular and Molecular Biology*, F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds., ASM Press, Washington, DC, pp. 1458–1496.
11. D. G. Rodionov, A. G. Pisabarro, M. A. De Pedro, W. Kusser, and E. E. Ishiguro (1995) *J. Bacteriol.* **177**, 992–997.
12. D. G. Rodionov and E. E. Ishiguro (1995) *J. Bacteriol.* **177**, 4224–4229.
13. M. Blanca, J. Garcia, J. M. Vega, A. Miranda, M. J. Carmona, C. Mayorga, F. Moreno, and C. Juarez (1996) *Clin. Exp. Allergy* **26**, 335–340.

### Suggestions for Further Reading

14. J.-V. Höltje. (1998) Growth of the stress-bearing and shape-maintaining murein sacculus of *Escherichia coli*, *Microbiol. Mol. Biol. Rev.* **62**, 181–203.
15. J. M. Ghuyssen (1991) Serine-lactamases and penicillin-binding proteins, *Annu. Rev. Microbiol.* **45**, 37–67.

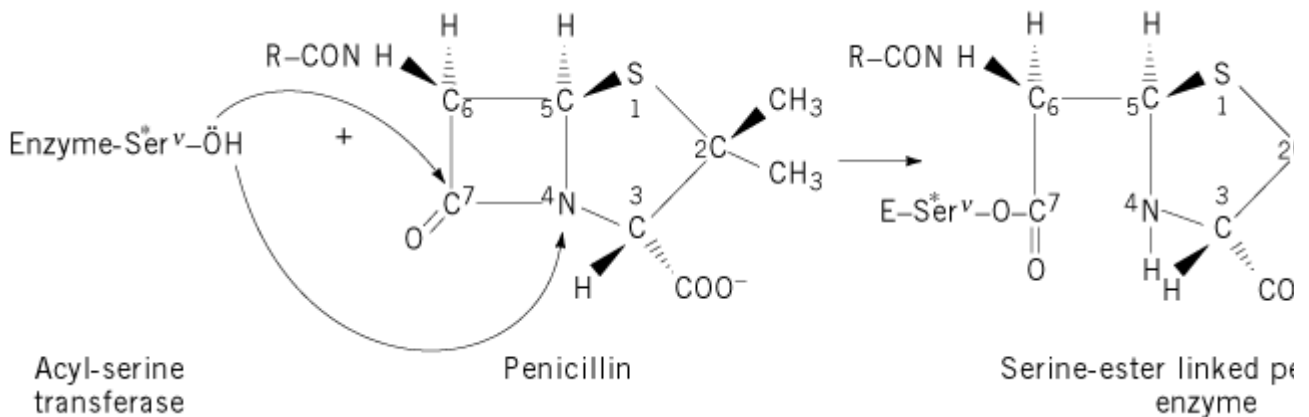
### Penicillin-Binding Proteins

The concept of antibiosis originated with Ehrlich and the sulfa drugs. It was firmly established by Fleming, Florey, and Chain, who discovered the uses of penicillin. The chemical structure of penicillin has some remarkable features: its  $\beta$ -lactam amide bond is equivalent to a [peptide bond](#) extending between two carbon atoms that have the D configuration (see **Stereochemistry**), in the  $\alpha$  position to a free [carboxyl group](#). Penicillin kills bacteria because it inactivates specialized acyl



serine transferases that are involved in the assembly and metabolism of the cell-wall peptidoglycan. It immobilizes the catalytic center in the form of a stable serine ester-linked penicilloyl enzyme derivative (Fig. 1). Variants of penicillin, such as [ampicillin](#), also exist.

**Figure 1.** The structure of penicillin and its reaction with the active-site serine residue of acyl-serine transferases. The mechanism of the reaction is shown at left, and the resulting penicilloyl enzyme is shown at right. If the penicilloyl enzyme is stable, the enzyme is a penicillin-binding protein. Alternatively, if the penicilloyl enzyme is hydrolyzed rapidly, the enzyme is a  $\beta$ -lactamase.



The bacterial cell-wall peptidoglycan is a covalently closed, net-like **polymer** (1). Linear glycan strands are made up of alternating  $\beta$ -1,4-linked units of *N*-acetylglucosamine and *N*-acetylmuramic acid. The carboxyl groups of the *N*-acetylmuramic acid residues are amide linked to the amino terminal L-alanine residues of peptide units L-alanyl-g-D-glutamyl-L-diaminoacyl-D-alanine. Neighboring peptide-substituted glycan strands are crosslinked by peptide bridges that extend from the D-alanine at the carboxyl end of one peptide to the side-chain [amino group](#) of the diamino acid residue of another peptide. The nature of the diamino acid residue and the composition of the peptide bridge vary according to the bacterial species.

In *Escherichia coli*, the diamino acid residue is *meso*-diaminopimelic acid, and the peptide bridges are direct D-alanyl-(D)-*meso*-diaminopimelic acid bonds. The immediate biosynthetic precursor of the cell-wall peptidoglycan is lipid II (2), in which a disaccharide peptide unit is linked to a C<sub>55</sub>H<sub>89</sub> undecaprenyl lipid carrier via a pyrophosphate bridge involving C<sub>1</sub> of *N*-acetylmuramic acid, and the peptide borne by *N*-acetylmuramic acid is a pentapeptide that terminates with a D-alanyl-D-alanine sequence. From this precursor, glycosyl transferases catalyze glycan chain elongation by displacement of the pyrophosphate linked to C1 of *N*-acetylmuramic acid by the 4-hydroxyl group of *N*-acetylglucosamine. Then acyl serine transferases catalyze peptidoglycan crosslinking at the expense of the D-alanyl-D-alanine bond of the pentapeptide units. This reaction proceeds via formation of a peptidyl ( $\sim$ L-alanyl-g- D-glutamyl-L-diaminoacyl- D-alanyl) enzyme intermediate linked as an ester to a serine residue at the [active site](#), with the concomitant release of the carboxy terminal D-alanine residue of the pentapeptide. It is achieved by the transfer of the peptidyl moiety to the side-chain amino group of the diamino acid residue of another peptide. Because the reaction involves breaking a D-alanyl-D-alanine bond, the acyl serine transferases are classified as DD-transpeptidases.

Acyl serine transferases are also involved in the control of the extent of peptidoglycan crosslinking. DD-carboxypeptidases hydrolyze the D-alanyl-D-alanine peptide bonds of pentapeptide units by transferring the peptidyl moiety of the acyl-enzyme intermediate to water. DD-endopeptidases hydrolyze the interpeptide bonds, which in some bacteria extend between two D centers in a position to a free carboxylate, as, for example, the D-alanyl-(D)*meso*-diaminopimelic acid bonds in the *E. coli*

peptidoglycan.

Penicillin has similarities to the above substrates and is a **suicide substrate** of the DD-(trans, carboxy, endo)peptidases. These DD-peptidases react with and rupture the b-lactam amide bond of penicillin but produce a serine ester-linked penicilloyl enzyme that is almost completely inert and inactive. The catalytic center turns over very slowly—one time or less per hour—and the inactivated enzymes are detectable as penicillin-binding proteins (PBPs).

In some cases, proteins acquired the ability to catalyze the hydrolysis of the serine ester-linked penicilloyl enzyme, which resulted in a remarkable defensive mechanism by producing enzymes, (the [b-lactamases](#)), with the ability to hydrolyze penicillin and other b-lactam antibiotics. On good b-lactam substrates, the catalytic center of the serine b-lactamases can turn over 1000 times or more per second.

The DD-peptidases and serine b-lactamases form a superfamily of penicilloyl serine transferases (3). They are an exemplary model of molecular [evolution](#). Numerous changes in their amino acid sequences, and gene fusion with the genes for several polypeptides of different origins, resulted in a prolific expansion of groups, classes, and subclasses of these enzymes, with multiple personalities. They share structural motifs. Some have similar enzymatic activity but varying substrate specificities. Others catalyze distinct reactions and fulfill different functions.

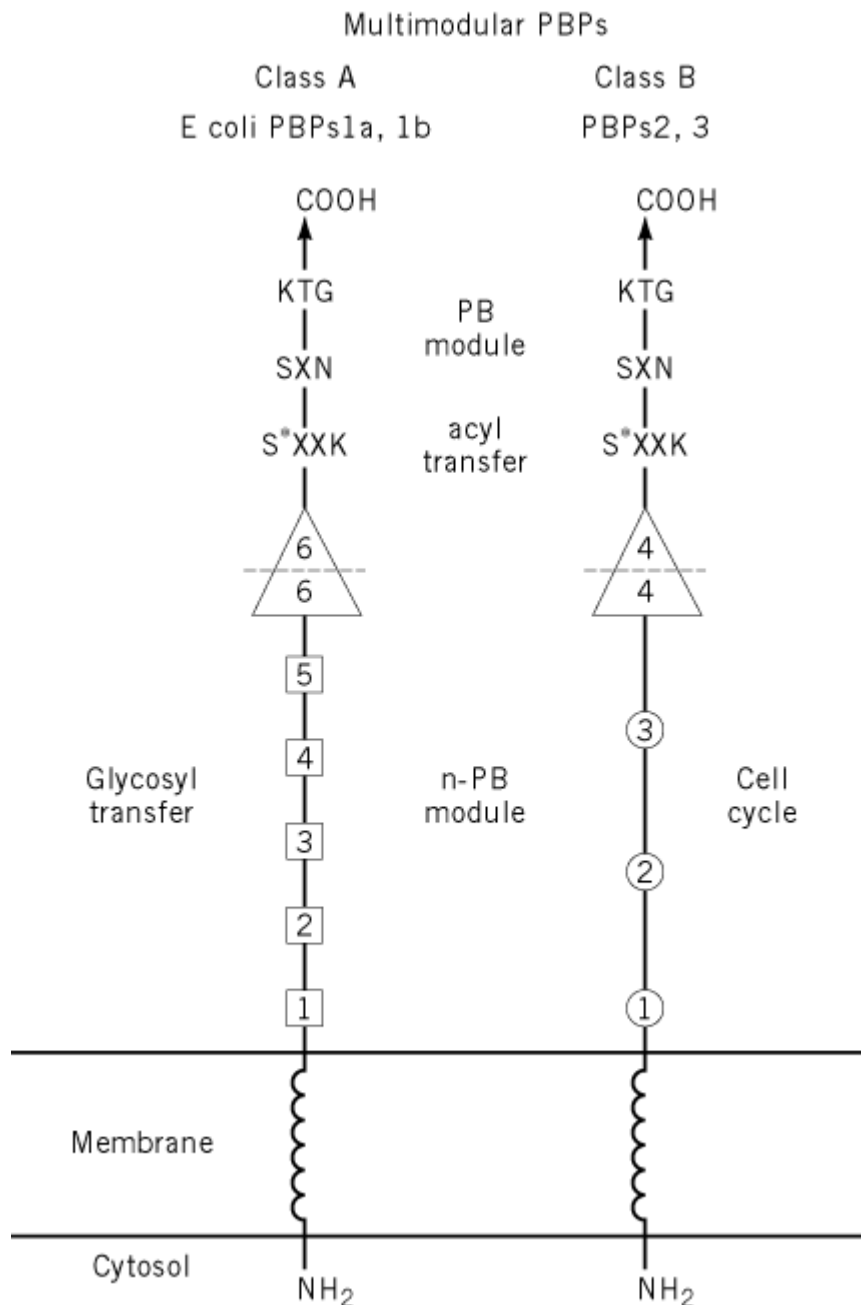
## 1. Protein Structure

The DD-peptidases and b-lactamases presumably diverged from a common ancestor while retaining the same basic [protein structure](#) (4). Similarities in their overall amino acid sequences are almost nonexistent, but folding invariably brings three common motifs of the polypeptide chains close to each other in the [tertiary structure](#), forming the catalytic center at the boundary between an all- $\alpha$  **domain** and an  $\alpha/\beta$  domain, which consists of a five-stranded [beta-sheet](#) covered by **alpha-helices** on both faces. With “x” denoting a variable amino acid residue, motif 1 has the sequence -Ser-x-x-Lys- (where Ser is the active-site serine), is at the amino end of an  $\alpha$ -helix, and occupies a central position in the cavity. Motif 2 has the **consensus sequence** [Ser/Tyr]-x-[Asn/Asp/Cys] and is on a loop between two  $\alpha$ -helices on one side of the cavity. Motif 3 has the sequence [Lys/His]-[Thr/Ser]-Gly and is on the innermost strand of the  $\beta$ -sheet on the other side of the cavity. The relative positions of these three motifs are well conserved in the various proteins.

## 2. Penicillin targets in *Escherichia coli*

*E. coli* possesses nine PBPs. PBPs 4, 5, 6, and 7 are monofunctional DD-carboxypeptidases that are not essential for viability. PBPs 1a, 1b, 1c, 2, and 3 are multidomain proteins (3) in which a transmembrane spanner is linked to the amino end of a non-penicillin-binding (n-PB) module that is linked to the amino end of an acyl serine transferase, penicillin-binding (PB) module. The two modules comprise a single polypeptide chain that folds on the exterior of the plasma membrane (Fig. 2). There are two classes: A and B (see text below).

**Figure 2.** Design and membrane topology of multimodular penicillin-binding proteins (PBPs) of class A and class B. Indicated are the five motifs of the n-PB modules of class A PBPs *f*, the three motifs of the n-PB modules of class B PBPs  $\circ$ , and the intermodule junction sites D. For more details on the class-specific motifs, see text. The PB modules bear the characteristic motifs of the penicilloyl serine transferase superfamily. Some PBPs have inserts and N- and C-terminal extensions, which are not shown.



PBP1a and PBP1b of class A (850 and 844 amino acid residues, respectively) combine the transglycosylase (n-PB module) and the DD-transpeptidase (PB module) activities that are required for peptidoglycan assembly. Each catalyzes the conversion of lipid II into polymeric peptidoglycan in *in vitro* assays (5). Loss of either PBP1a or PBP1b is tolerated, but loss of both is fatal (6). Under these extreme conditions, PBP1c (770 amino acid residues) does not rescue *E. coli* from death.

PBP2 and PBP3 of class B (633 and 588 amino acid residues, respectively) are inert on lipid II (7). They combine in a single polypeptide chain a [cell cycle](#) (n-PB) module and an acyl serine transferase (PB) module. Both modules are essential, and loss of neither PBP2 nor PBP3 is tolerated. The genes for PBP2, the DD-carboxypeptidase PBP5, and the integral membrane protein RodA cluster at the 14-min region of the chromosome; these proteins prescribe traits as complex as cell-wall expansion and cell-shape maintenance (8). PBP3, MraW, FtsL, FtsQ, FtsA, the cytoskeletal FtsZ, and FtsW (an analog of RodA)—the encoding genes of which cluster at the 2-min region of the chromosome—prescribe traits as complex as septum formation (9).

### 3. The Multimodular PBP Family

Sixty-three multimodular PBPs of various bacterial species (10) comprise two classes, A and B, with distinct amino acid sequence signatures in their n-PB modules but conserved motifs in their acyl serine transferase-PB modules. Their junction sites between the n-PB and PB modules usually have the carboxy terminal sequence Asp-x-x-x-Gln (Fig. 2). There is a clustering into several subclasses of A and B, corresponding to **Gram-negative** and **Gram-positive** bacteria.

Based on these observations, the cores of the n-PB modules (n-PB cores) of the multimodular PBPs start at the amino end of the first motif and terminate at the intermodule junction site. The cores of the associated PB modules (PB cores) start about 70 residues N-terminal of the Ser-x-x-Lys motif— at the carboxy end of the intermodule junction site (where insertions occasionally occur), and they terminate about 70 residues C-terminal of the Lys-[Thr/Ser]-Gly motif; PBPs also have C-terminal extensions. In *E. coli* PBP1b, motif 3 of the PB core has the sequence Lys-Thr-Gly at residues 698–700, residues 762–780 are essential, but residues 781–844 are dispensable (11). In *Streptococcus pneumoniae* PBP2x (the structure of which is known at low resolution [(12)]), motif 3 of the PB core has the sequence Lys-Ser-Gly at residues 547–549, and residues Ser617–Asp750 form a carboxy terminal module with its own fold.

The multimodular PBPs are fully integrated hybrids. In contrast to the monofunctional PBPs, which are autonomous folding entities, the PB modules of multimodular PBPs have lost the ability to fold by themselves (13-15). They require the assistance of the associated n-PB modules, which implies precise and specific module-module interactions. Correct folding (in terms of penicillin binding) of the associated n-PB and PB modules of *E. coli* PBP3 is independent of the transmembrane spanner (13). In contrast, the cell septation activity of *E. coli* PBP3 is dependent on the membrane-anchor module, indicating that the membrane anchor and its cytosolic tail have more sophisticated functions than that of a simple anchoring device (16).

### 4. Evolution of Multimodular Class A PBPs

The n-PB cores of the class A PBPs form a continuum of homologous sequences, and they have an extended amino acid sequence signature in the form of five motifs, indicative of conserved transglycosylase activity. In *E. coli* PBP1a, motif 1 starts at residue 86 with the sequence Glu-Asp-Ser-Arg-Phe-x-Glu-His-x-Gly, motif 2 starts at residue 117 with the sequence Gly-Ala-Ser-Thr-Ile-Thr-Gln-Gln, motif 3 starts at residue 139 with the sequence Arg-Lys-x-x-Glu, motif 4 starts at residue 156 with the sequence Lys-x-Glu-Ile-Leu-Glu-x-Tyr-x-Asn, motif 5 starts at residue 221 with the sequence Arg-Arg-x-x-x-Val-Leu-, and Gly284 is at the amino end of the intermodule junction site. Residues Glu86 of motif 1 and Glu143 of motif 3 are almost certainly important components of the transglycosylase catalytic center.

The associated PB cores of the class A PBPs of Gram-negative bacteria cluster into two subclasses: A1 (prototype: *E. coli* PBP1a) and A2 (prototype: *E. coli* PBP1b), which form a supercluster. Those of Gram-positive bacteria cluster into three subclasses: A3 (prototype: *S. pneumoniae* PBP1a), A4 (prototype: *S. pneumoniae* PBP2a), and A5 (prototype: *S. pneumoniae* PBP1b).

Although performing the same transglycosylase-transpeptidase activities, PBPs of subclasses A1 and A2 in Gram-negative bacteria and PBPs of subclasses A3, A4, and A5 in Gram-positive bacteria are almost certainly **paralogous**. They prescribe distinct, subtle traits in peptidoglycan assembly in these two groups of bacteria that sometimes are difficult to detect.

### 5. Evolution of Multimodular Class B PBP

The class B PBPs evolved differently than did class A PBPs (10). The n-PB and PB cores underwent adaptive evolution in such a concerted manner that, in Gram-negative bacteria, an n-PB core of

subclass B2 (prototype: *E. coli* PBP2) or B3 (prototype: *E. coli* PBP3) is linked to a PB core of subclass B2 or B3, respectively. In Gram-positive bacteria, an n-PB core of subclass B1 (prototype: *Enterococcus faecium* PBP3), B4 (prototype: *S. pneumoniae* PBP2x), or B5 (prototype: *S. pneumoniae* PBP2b) is linked to a PB core of subclass B1, B4, and B5, respectively.

This unique combination of the n-PB and PB cores indicates that they have different functions in cell morphogenesis (n-PB module) and peptidoglycan crosslinking (PB-module). The **Gram-negative** PBPs of subclasses B2 and B3 are paralogous, and the **Gram-positive** PBPs of subclasses B1, B4, and B5 are also paralogous. In contrast, the **Gram-negative** PBPs of subclasses B2 or B3 may be orthologous to the **Gram-positive** PBPs of subclass B5 or B4, respectively.

The **Gram-positive** PBPs of subclass B1 have no equivalent in the **Gram-negative** bacteria. They are endowed with unique properties, and they represent an important mechanism of resistance to penicillin. These PBPs have a very low affinity for the drug. They allow the strains that (over) produce them to grow in the presence of penicillin at concentrations sufficient to inactivate the PB modules of all the other PBPs of class A and class B. To all appearances, they are able to perform the basic functions required for cell-wall peptidoglycan assembly in a cell cycle-dependent fashion in conjunction (presumably) with the transglycosylase-n-PB module of class A PBPs (or with monofunctional transglycosylases).

The question of how the class B PBPs function in cell morphogenesis is still open. Yet the unique amino acid sequences of their n-PB modules (10) and the three-dimensional structure of *S. pneumoniae* PBP2x (12) are worthy of reflection. The n-PB cores of all class B PBPs have three conserved motifs. In *S. pneumoniae* PBP2x, motif 1 starts at residue 76 and has the sequence Arg-Gly-x-x-x-Asp-Arg-Asn-Gly, motif 2 starts at residue 186 and has the sequence Arg-x-Tyr-Pro-x-Gly, motif 3 starts at residue 215 and has the sequence Gly-x-x-Gly-x-Glu-x-x-x-Asn, and Gly258 is the amino end of the intermodule junction site. The n-PB module of PBP2x is shaped like a pair of sugar tongs; motifs 1, 2, 3 and the amino half of the intermodule junction site are located in the head of the sugar tongs (which fits into a noncatalytic groove of the PB module), and the long polypeptide stretch that extends between motifs 1 and 2 is well exposed at the surface of the protein.

Other class B PBPs probably adopt the same basic folded structure but with subclass- and species-specific variations. The acyl serine transferase-PB module of class B PBPs probably prescribes traits related to peptidoglycan crosslinking, and this activity might be regulated by the associated n-PB module itself, in conjunction with other components of the morphogenetic networks.

## 6. Penicillin-Oriented Evolution

Evolution is occurring before our eyes. The massive use of b-lactam antibiotics functions to fuel the emergence of b-lactamases of widely varying specificities by the alteration of a limited number of amino acid residues in the wild-type enzymes (17). Oriented evolution also results in the spread of low-affinity PBPs of subclass B1 among enterococcal and staphylococcal pathogens (18, 19). Strains close to the squirrel *Staphylococcus sciuri* might be the source of the PBPqa-encoding gene in human MRSA (methicillin-resistant *S. aureus*) strains (20). Multimodular PBPs also evolve into variants having a decreased affinity for the drug. Penicillin-resistant *S. pneumoniae* strains possess PBPs that are reduced-affinity variants of PBP1a of subclass A3, PBP2x of subclass B4, and PBP2b of subclass B5 (21). These altered PBPs are the products of **mosaic genes** in which **sensitive** sequence blocks are replaced by homologous, **resistant** blocks from related, naturally penicillin-resistant species by [transformation](#) and [recombination](#) events (22). The origin of the **foreign** DNA sequences is unknown, but several sources are probably involved.

## 7. Penicillin Sensory Transducers

BlaR is a bipartite protein involved in the inducibility of b-lactamase synthesis in *Bacillus licheniformis* (23). An amino terminal module is embedded in the plasma membrane via a **four-helix**

**bundle.** Loops connecting the transmembrane segments are exposed on the outer and inner faces of the membrane, respectively; the carboxy end of the fourth transmembrane segment is fused to an extracellular penicilloyl serine transferase module. This module possesses the three typical motifs described earlier (Ser402-x-x-Lys, Tyr476-x-Asn, and Lys539-Thr-Gly), and it shares 32% identity with the Oxa-2 class D b-lactamase.

As an independent entity, the penicilloyl serine transferase module of BlaR is a high-affinity PBP. As a component of the full-size BlaR, it adopts a different conformation, and reception of the penicillin-induced signal does not involve penicilloylation of the active-site Ser402. These properties suggest that a class D b-lactamase acquired a new property—penicillin binding—via local structural changes. Fusion of this PBP to another polypeptide then resulted in a hybrid that performs another function: gene regulation.

Signal emission in the cytosol is probably mediated by the intracellular loop that connects transmembrane segments 3 and 4. This loop has no recognizable site of **methylation**/demethylation and no recognizable site that could be associated with a histidine kinase, but it possesses the His-Glu-Leu-Tyr-His **consensus sequence** of a neutral zinc [peptidase](#) (see [Metalloproteinases](#)). This putative peptidase might adopt an active conformation in response to the signal transmitted via the four  $\alpha$ -helix bundle. Following this hypothesis, BlaR would be another example of control of gene activity via proteolytic degradation of key regulatory proteins.

Penicillin receptors similar to the *B. licheniformis* BlaR are involved in the inducibility of the synthesis of b-lactamase and low-affinity PBP2a of subclass B1 in *S. aureus* (24). Gram-negative bacteria have developed a different mechanism. b-Lactamase synthesis is the result of the b-lactam-induced, peptidoglycan hydrolase-mediated deregulation of the peptidoglycan recycling process (25).

## Bibliography

1. J. M. Ghuysen (1968) *Bact. Rev.* **32**, 425–464.
2. J. van Heijenoort (1996) In: *Escherichia coli and Salmonella*, 2nd ed., ASM Press, Washington, D C, pp. 1025–1034.
3. J. M. Ghuysen (1991) *Annu. Rev. Microbiol.* **45**, 37–67.
4. J. A. Kelly, A. P. Kuzin, P. Charlier, and E. Fonzé (1998) *Cell. Mol. Life Sci.* **54**, 353–358.
5. M. Matsushashi (1994) *New Comprehensive Biochem.* **27**, 55–72.
6. J. I. Kato, S. Hideho, and Y. Hirota (1985) *Mol. Gen. Genet.* **200**, 272–277.
7. M. Adam, C. Fraipont, N. Rhazi, M. Nguyen-Distèche, B. Lakaye, J. M. Frère, B. Devreese, J. Van Beeumen, Y. van Heijenoort, J. van Heijenoort, and J. M. Ghuysen (1997) *J. Bacteriol.* **179**, 6005–6009.
8. D. Joseleau-Petit, D. Thévenet, and R. D'Ari (1994) *Mol. Microbiol.* **13**, 911–917.
9. M. Vicente, M. J. Gomez, and J. A. Ayala (1998) *Cell. Mol. Life Sci.* **34**, 317–324.
10. C. Goffin and J. M. Ghuysen (1998) *Microbiol. Mol. Biol. Rev.*, **62**, 1079–1093.
11. F. Lefèvre, M. H. Rémy, and J. M. Masson (1997) *J. Bacteriol.* **179**, 4761–4767.
12. S. Pares, N. Mouz, Y. Pétilot, R. Hakenbeck, and O. Dideberg (1996) *Nature Struct. Biol.* **3**, 284–289.
13. C. Goffin, C. Fraipont, J. Ayala, M. Terrak, M. Nguyen-Distèche, and J. M. Ghuysen (1996) *J. Bacteriol.* **178**, 5402–5409.
14. E. C. Y. Wu, W. E. Alborn Jr., J. E. Flokowitsch, J. Hoskins, S. Unal, L. C. Blaszcak, D. A. Preston, and P. L. Skatrud (1994) *J. Bacteriol.* **176**, 443–449.
15. M. Mollerach, P. Partoune, J. Coyette, and J. M. Ghuysen (1996) *J. Bacteriol.* **178**, 1774–1775.
16. L. M. Guzman, D. S. Weiss, and J. Beckwith (1997) *J. Bacteriol.* **179**, 5094–5103.
17. K. Bush and G. Jacoby (1997) *J. Antimicrob. Chemother.* **39**, 1–3.

18. B. Bergi-Bächli (1994) *Trends Microbiol.* **2**, 389–393.
19. D. Raze, O. Dardenne, S. Hallut, M. Martinez-Bueno, J. Coyette, and J. M. Ghuysen (1998) *Antimicrob. Ag. Chemother.* **42**, 534–539.
20. S. Wu, C. Piscitelli, H. de Lencastre, and A. Tomasz (1996) *Microb. Drug Resist.* **2**, 435–441.
21. R. Hakenbeck and J. Coyette (1998) *Cell. Mol. Life Sci.* **54**, 332–340.
22. C. G. Dowson, J. T. Coffey, and B. C. Spratt (1994) *Trends Microbiol.* **2**, 361–366.
23. K. Hardt, B. Joris, S. Lepage, R. Brasseur, J. O. Lampen, J. M. Frère, A. L. Fink, and J. M. Ghuysen (1997) *Mol. Microbiol.* **23**, 935–944.
24. B. Joris, K. Hardt, and J. M. Ghuysen (1994) *New Comprehensive Biochem.* **27**, 505–516.
25. C. Jacobs, J. M. Frère, and S. Normark (1997) *Cell* **88**, 823–832.

### Suggestions for Further Reading

26. J. V. Höltje (1998) Growth of the stress-bearing and shape-maintaining murein sacculus of *Escherichia coli*. *Microbiol. Mol. Biol. Rev.* **62**, 181–203.
27. J. Lutkenhaus and S. G. Addinall (1997) Bacterial cell division and the Z ring. *Annu. Rev. Biochem.* **66**, 93–116.
28. N. Nanninga (1998) Morphogenesis of *Escherichia coli*. *Microbiol. Mol. Biol. Rev.* **62**, 110–129.

## Pepsin, Pepsinogen

*Pepsin* is a [carboxyl proteinase](#) that is secreted from the chief (or peptic) cells that are part of the oxyntic glands, located on the lining of the stomach (1). It is synthesized as a 42.5-kDa [zymogen](#) precursor, *pepsinogen*, which is stored in the chief cells and released in response to neural and hormonal signals activated by eating. Once pepsinogen enters the acidic environment of the stomach, it undergoes a change in conformation that exposes its [active site](#). This allows pepsinogen to catalyze the hydrolysis of the [peptide bond](#) between residues 16 and 17 of its own polypeptide chain. Subsequently, residues 17 to 44 are also removed proteolytically, and the resultant 35-kDa protein is pepsin. It is active at acidic pH and has broad specificity, with preference for hydrolysis of peptide bonds involving [tryptophan](#), [phenylalanine](#), [leucine](#), and [methionine](#) residues, especially when they are adjacent to other **hydrophobic** amino acids. It is often used to digest proteins under acidic conditions in order to obtain peptides linked by [disulfide bonds](#). Disulfide rearrangement, which can sometimes occur at neutral or alkaline pH, is minimized at acidic pH. Pepsinogen is stable at neutral pH and becomes activated at acidic pH, whereas pepsin is stable at acid pH but becomes inactivated at neutral pH.

Pepsin is inhibited by [pepstatin](#), *p*-bromophenacyl bromide, diazoketones, and various epoxides. It has been reported that pepsin can catalyze transpeptidation reactions, which could, in principle, cause confusion when trying to sequence proteins. This has not been a major problem in practice, however.

### Bibliography

1. M. Plebani (1993) *Crit. Rev. Clin. Lab. Sci.* **30**, 273–328.

## Pepstatin

This naturally occurring peptide analogue was discovered by screening culture filtrates of the mold actinomycetes for inhibitors of the **proteolytic** enzyme **pepsin** (1). It also inhibits other [carboxyl proteinases](#), but not [serine](#), [thiol](#), or [metalloproteinases](#). Its chemical structure is: isovaleryl-valyl-valyl-statine-alanyl-statine [where statine is the unusual residue (3*S*, 4*S*)-4-amino-3-hydroxy-6-methylheptanoic acid] (Fig. 1, see top of next page). Analogues of pepstatin (also called pepstatin A) have been obtained from *Streptomyces* strains with variable inhibitory properties toward different carboxyl proteinases, and a range of synthetic pepstatins have been prepared that have quite remarkable affinities for specific carboxyl proteinases (2). Pepstatin is often used to prevent unwanted proteolysis when isolating proteins under acidic conditions. It is effective at a concentration of about 1 mg/mL.

### Bibliography

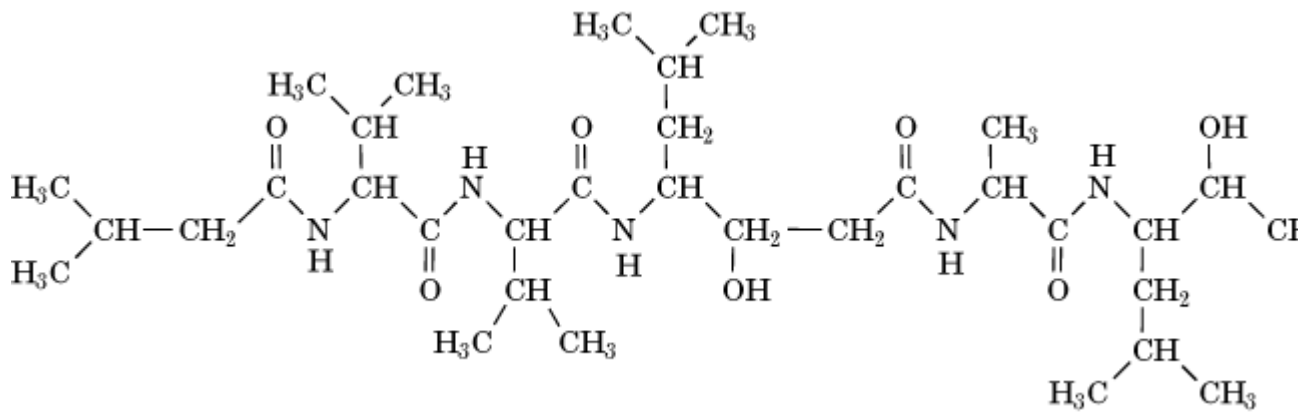
1. H. Umezawa et al. (1970) *J. Antibiot.* **23**, 259–262.
2. T. Aoyagi and H. Umezawa (1975) In *Proteases and Biological Control* (E. Reich, D. B. Rifkin, and E. Shaw, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 429–454.

## Peptidases

Peptidases are [enzymes](#) that catalyze the hydrolysis of [peptide bonds](#)—that is, the -CONH- bonds that link [amino acids](#) in [peptides](#), [polypeptides](#), and [proteins](#). According to the Enzyme Commission of the International Union of Biochemistry and Molecular Biology, they are Class 3, subgroup 4 (peptide hydrolase) enzymes (E.C. 3.4). Some peptidases only act on small substrates—di-, tri-, or oligopeptides—while others act preferentially on protein substrates. The latter are more properly referred to as [proteinases](#) (also known as proteases, proteolytic enzymes). Some peptidases cleave peptide bonds involving the amino acids at either the beginning or end, i.e., the amino terminal (N-terminal) or carboxy terminal (C-terminal)—of a peptide chain. They are therefore called *exopeptidases* and, depending on their particular specificity, are either [aminopeptidases](#) or [carboxypeptidases](#). Other peptidases cleave internal bonds in a polypeptide chain and are called [endopeptidases](#). Most peptidases are relatively specific and will only act at peptide bonds in which a particular amino acid or type of amino acid (acidic, basic, neutral, **hydrophobic**, etc.) contributes either the -CO or -NH group. Some have broad specificity, especially when acting on unfolded, **denatured** proteins, but they can be highly specific or not active at all with folded, native proteins.

**Figure 1.** Chemical structure of pepstatin [isovaleryl-valyl-valyl-(3*S*, 4*S*)-4-amino-3-hydroxy-6-methylheptanoyl-alanyl-(3-amino-3-hydroxy-6-methylheptanoic acid)], a transition-state analogue inhibitor of certain carboxyl proteinases, such as pepsin, cathepsin D, and many microbial carboxyl proteinases. Because its structure resembles the tetrahedral intermediate that forms during substrate hydrolysis, it forms a tightly bound complex with the active site of each of these enzymes. See [Pepstatin](#), previous page.





Peptidases have myriad functions. They activate or inactivate peptide [hormones](#); they degrade dietary proteins and peptides; they participate in structural growth and remodeling by degrading [collagen](#) and other structural proteins; they protect against infectious agents; and they induce [blood clotting](#). They regulate blood pressure, play an important role in fertilization, and control the [cell cycle](#). They are ubiquitous and enormously diverse. As might be expected, inappropriate peptidase activity could have devastating consequences, so there are numerous peptidase inhibitors and other biological means of curtailing peptidase activity (see [Proteinase Inhibitors](#)).

While all peptidases catalyze peptide bond hydrolysis, they do not all do this the same way. In order to act under various conditions, four different mechanisms have evolved, and peptidases are classified as [carboxyl proteinase](#), [thiol proteinase](#), [serine proteinase](#), and metalloproteinase, depending on the particular mechanism they employ. Even within these classes, there are many thematic variations such that it seems that there are almost as many peptidases as there are peptides.

Peptidases have a host of commercially significant uses ranging from “clot busters” (both in blood vessels and in drain pipes) to meat tenderizers. Peptidase inhibitors represent a multibillion dollar pharmaceutical market where they are important for controlling high blood pressure, as well as preventing the progress of acquired immune deficiency syndrome (AIDS).

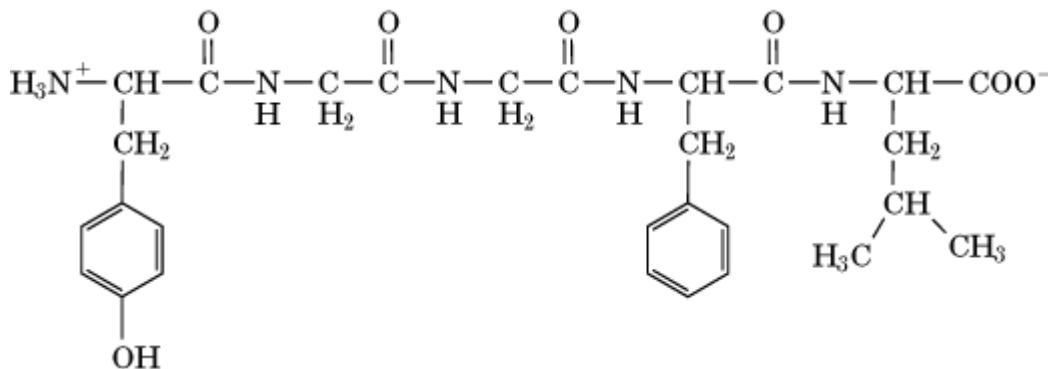
## Peptide

Peptides, sometimes referred to as [oligopeptides](#), are [oligomers](#) in which the repeating units are [amino acids](#). Peptides have a defined sequence of amino acids that are linked together by formation of [peptide bonds](#). In contrast to **polypeptides** and [proteins](#), peptides consist of a small number of amino acids. The distinction between a peptide and a polypeptide is somewhat arbitrary, but generally a peptide has between 2 and 50 amino acid residues. A peptide of two amino acids may be more specifically referred to as a *dipeptide*, three repeating units as a *tripeptide*, four repeating units as a *tetrapeptide*, and so on.

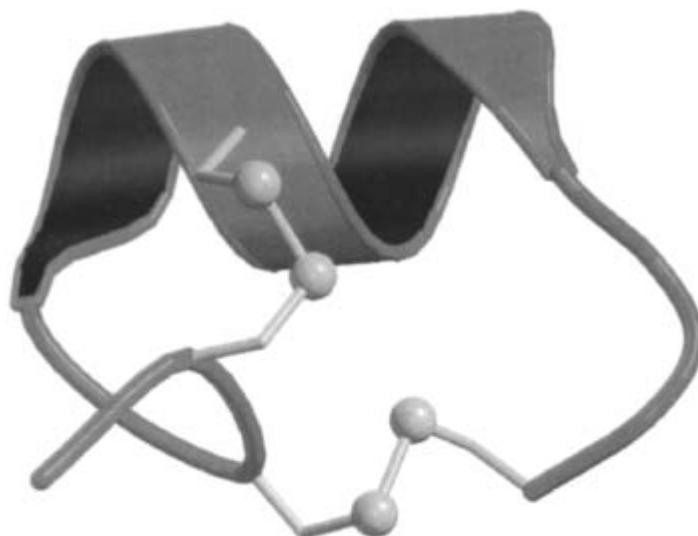
Most peptides are unstructured, described as having a **random coil conformation**, but others have highly ordered **secondary** and [tertiary structure](#) similar to that observed in larger proteins. The more structured peptides often incorporate [disulfide bonds](#). Many peptides have important biological functions. For example, the pentapeptide leucine-enkephalin (Fig. 1) is a neurotransmitter with opiate-like activity; the hormone [insulin](#) is formed from two disulfide cross-linked peptide chains;

and the major [toxins](#) from venomous spiders and cone shells are disulfide-rich peptides (Fig. 2). Naturally occurring peptides are usually produced by hydrolysis or cleavage of precursor proteins, through the action of [peptidases](#) or [proteinases](#).

**Figure 1.** Chemical structure of the pentapeptide Leu-enkephalin, which has the sequence Tyr–Gly–Gly–Phe–Leu.



**Figure 2.** Schematic representation of the structure of the 16-residue  $\alpha$ -conotoxin PnIA (1), showing the backbone and disulfides. The two turns of  $\alpha$ -helix are depicted as a coil. The two disulfide bonds are shown in ball-and-stick representation. This figure was generated using Molscript (2) and Raster3d (3, 4).



[See also **Polypeptide chain.**]

#### Bibliography

1. Hu et al. (1996) *Structure* **4**, 417–423.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

5. M. Bodanszky (1988) *Peptide Chemistry*, Springer-Verlag, New York.  
 6. P. D. Bailey (1990) *An Introduction to Peptide Chemistry*, Wiley, Chichester.

## Peptide Antibiotic Biosynthesis

Peptide antibiotics include a variety of [peptide](#) structures with diverse biological properties (1). Historically, the first antibacterial compounds identified were peptides ([penicillin](#), [gramicidin](#), [tyrocidine](#)). New such peptides have been characterized steadily, including [actinomycin](#) (antitumor); [cyclosporin](#) (immunosuppressant); destruxin (insecticide); ferrichrome (siderophore); [pepstatin](#), antipain, and bestatin ([proteinase inhibitors](#)); microcystin ([phosphatase inhibitor](#)); and defensin and magainin (anti-infectives).

Their biosynthesis may occur either (i) on [ribosomes](#) by direct [transcription](#), [translation](#), and **post-translational processing** of peptide antibiotic **genes** (see **Protein biosynthesis**)—for example, nisin—or (ii) without ribosomes but on a multienzyme system, using the **thiotemplate mechanism**. The structural complexity of peptide antibiotics ranges from modified [amino acids](#) to **polypeptide** chains of up to 50 amino acid residues (2). Ribosomal systems can incorporate only the L isomer of the 20 normal [amino acids](#), plus [selenocysteine](#), while there are no such limitations to the nonribosomal systems, which can use hydroxy acids, fatty acids, and amines. Nevertheless, peptides of ribosomal origin are known to include D amino acid residues, cyclic side chains, and [disulfide bonds](#). All nonribosomal peptides also originate from linear precursors, but a variety of types of cyclization are known. Identification of the biosynthetic origin is crucial for subsequent genetic analysis, and the structural properties of compounds arising from the two types of pathway are compared in Table 1.

**Table 1. Structural Features of Peptide Antibiotics of Ribosomal and Nonribosomal Origin**

| Feature                     | Ribosomal Path   | Nonribosomal Path   |
|-----------------------------|--|---|
| Size (amino acids)          | No size limitation                                     | 2 to about 50, 4 to 10 dominating   |
| amino acid constituents     | 21 protein amino acids and modified amino acids        | Various types of amino acids modified amino acids including 2-, 3-, and 4-amino compounds (more than 300 known) |
| D-Amino acids               | Not more than one, epimerized post-translationally     | Often several, either incorporated directly, or epimerized during synthesis                                     |
| Non-amino acid constituents | Acyl residues, amines originating from decarboxylation | Various acyl residues, including aromatic acids, hydroxy acids  |
| Cyclic structures           | Rare; frequent disulfide cycles, thioether cycles      | More frequent than linear structures, various peptide bond cyclizations, but also lactones                      |
| Modifications               | Hydroxylation,   | N-Methylation,  |

|                      |   |   |
|----------------------|---|---|
|                      | dehydration (Ser, Thr), side-chain cyclization (Cys to thiazoles, Thr to oxazoles, Glu to pyroGlu)                    | hydroxylation, side-chain cyclization (Cys to thiazole), glycosylation, side-chain crosslinking (aromatic rings)                            |
| Unusual constituents | Not known   | Urea type of peptide bond, phospho-amino acids, amino-modified fatty-acid-derived components (lipopeptides), mixed polyketide structures    |
| Biosynthesis         | Gene can be identified; consider splicing, processing, and posttranslational modification; often prepeptides detected | Nonribosomal enzyme systems present; peptide families are frequent in the same or related organisms, biosynthesis of rare precursors needed |
| Sources              | Various animals and plants, sometimes bacteria  | Mainly bacteria and lower fungi, occasionally plants and insects  |

---

## 1. Ribosomal Biosynthesis

Peptide antibiotics that are synthesized initially on ribosomes comprise a rapidly expanding field of research (3, 4). The compounds synthesized in this way range from linear polypeptide chains to compact multicyclic peptides with several disulfide bonds or thioether linkages (lantibiotics). Their genes are usually isolated by **reverse genetics**, and they usually reside within gene clusters, together with genetic information for their post-translational modification and export, and for resistance of the host to their actions. Genes for typical animal peptides, like defensins, are present in multiple copies with minor differences in sequence (5). The signals for their targeting and processing are critical for their functions in innate immunity. A peptide from the venom from the spider *Angelopsis aperta* is activated by **epimerization** of a specific residue (6).

Internal cyclizations other than disulfide bond formation have been restricted to bacterial sources, and the enzymes involved in formation of thiazolidines (microcin B) and lantibiotic ethers (dehydration of [serine](#) and [threonine](#) residues to dehydroalanine and dihydrobutryne, respectively, followed by addition of the [thiol group](#) of cysteine) are being characterized (7). The post-translational modification of these antibiotics takes place on a membrane-attached multienzyme complex, which also facilitates their export and their role in intercellular communication.

## 2. Synthesis by Peptide Synthetases

The biosynthesis of peptide antibiotics by nonribosomal systems, on peptide synthetases, requires a considerable amount of information, and these multienzyme complexes are among the largest protein structures known, with masses up to 1700 kDa, corresponding to a 45.5 kbp open reading frame, to synthesize cyclosporin (see [Gene Structure](#)). Synthesis of each amino acid residue of such a peptide antibiotic requires a minimum of three catalytic **domains**: (i) the activating adenylate domain, (ii) the carrier protein, and (iii) the condensing domain, which collectively are referred to as a module (see [Thio-template Mechanism Of Peptide Antibiotic Synthesis](#)). The nonribosomal code for the amino acid sequence of the peptide is determined by the substrate specificity of the activating domain. Selection and activation of the appropriate amino acid is followed by its attachment as a thioester to the 4'-phosphopantetheine prosthetic group attached to the adjacent carrier domain (8); it is then linked covalently to the next amino acid, on the adjacent module. The thioester intermediates

may be subjected to various modifications, such as epimerization, *N*-methylation, and hydroxylation, catalyzed by additional domains introduced into the module structure. The intermediates in the reaction are covalently attached to the peptide synthetase, and only the final product is released. Examples of known peptide synthetases are compiled in Table 2.

**Table 2. Nonribosomal Peptide Synthetase Systems**

| Peptide              | Organism   | Structural Type <sup>a</sup> | Gene Cloned | Enzymology |
|----------------------|--|------------------------------|-------------|------------|
| <b>Linear</b>        |  |                              |             |            |
| Bacilysin            | <i>Bacillus subtilis</i>   | P-2-M                        | (+)         | (+)        |
| ACV                  | <i>Streptomyces clavuligerus</i><br><i>Aspergillus nidulans</i><br><i>Penicillium chrysogenum</i><br><i>Acremonium chrysogenum</i> | P-3                          | +           | +          |
| Bialaphos            | <i>Streptomyces hygroscopicus</i>  | P-3                          | +           | +          |
| Anguibactin          | <i>Vibrio anguillarum</i>  | R-P-2-M                      | +           | –          |
| Phaseolotoxin        | <i>Pseudomonas syringae</i> pv. <i>ph.</i>   | P-4-M                        | (+)         | –          |
| Ardacin              | <i>Kibdelosporangium aridum</i>  | P-7-M                        | (+)         | –          |
| Pyoverdin            | <i>Pseudomonas aeruginosa</i>  | R-P-8-M                      | +           | –          |
| <b>Cyclopeptides</b> |  |                              |             |            |
| Enterobactin         | <i>Escherichia coli</i>  | P-C-E-3                      | +           | +          |
| HC-toxin             | <i>Cochliobolus carbonum</i>   | C-4                          | +           | +          |
| Tentoxin             | <i>Alternaria alternata</i>  | C-4                          | –           | +          |
| Echinocandin         | <i>Aspergillus nidulans</i>  | R-C-6                        | –           | (+)        |
| Microcystin          | <i>Microcystis aeruginosa</i>  | C-7                          | (+)         | –          |
| Iturin               | <i>Bacillus subtilis</i>   | C-8                          | (+)         | –          |
| Gramicidin S         | <i>Bacillus brevis</i>   | C-(P-5) <sub>2</sub>         | +           | +          |
| Tyrocidin            | <i>Bacillus brevis</i>   | C-10                         | +           | +          |
| Cyclosporin          | <i>Tolypocladium niveum</i>  | C-11-M                       | +           | +          |
| Mycobacillin         | <i>Bacillus subtilis</i>   | C-13                         | –           | +          |
| <b>Lactones</b>      |  |                              |             |            |
| Actinomycin          | <i>Streptomyces chrysomallus</i>   | R-(L-5) <sub>2</sub> -M      | +           | +          |
| Destruxin            | <i>Metarhizium anisopliae</i>  | L-6                          | +           | (+)        |
| Etamycin             | <i>Streptomyces griseus</i>  | R-L-7                        | –           | (+)        |
| Surfactin            | <i>Bacillus subtilis</i>   | L-8                          | +           | +          |
| Quinomycin           | <i>Streptomyces echinatus</i>  | (R-P-4) <sub>2</sub>         | –           | +          |

|                                 |                                 |                                     |     |     |
|---------------------------------|---------------------------------|-------------------------------------|-----|-----|
| R106                            | <i>Aureobasidium pullulans</i>  | L-9                                 | –   | (+) |
| Syringomycin                    | <i>Pseudomonas syringae</i>     | R-L-9                               | (+) | –   |
| Syringostatin                   |                                 |                                     |     |     |
| SDZ90-215                       | <i>Septoria sp.</i>             | L-10                                | –   | +   |
| <b>Depsipeptides</b>            |                                 |                                     |     |     |
| Enniatin                        | <i>Fusarium sp.</i>             | C-(P <sub>2</sub> ) <sub>3</sub> -M | +   | +   |
| Beauvericin                     | <i>Beauveria bassiana</i>       | C-(P <sub>2</sub> ) <sub>3</sub> -M | –   | +   |
| <b>Branched polypeptides</b>    |                                 |                                     |     |     |
| Bacitracin                      | <i>Bacillus licheniformis</i>   | P-12-C-7                            | +   | +   |
| Nosiheptide                     | <i>Streptomyces actuosus</i>    | R-P-13-C-10-M                       | –   | +   |
| Thiostrepton                    | <i>Streptomyces laurentii</i>   | R-P-17-C-10-M                       | –   | +   |
| <b>Branched peptidolactones</b> |                                 |                                     |     |     |
| Lysobactin                      | <i>Lysobacter sp.</i>           | P-11-L-9                            | (+) | +   |
| A21798A                         | <i>Streptomyces roseosporus</i> | R-P-13-L-10                         | (+) | +   |
| A54145                          | <i>Streptomyces fradiae</i>     | R-P-13-L-10                         | –   | +   |
| Tolaasin                        | <i>Pseudomonas tolaasii</i>     | R-P-18-L-5                          | (+) | –   |

<sup>a</sup> The abbreviations used are: P, peptide; C, cyclopeptide; L, lactone; E, ester; R, acyl; M, modified. The structural types are defined by the number of amino-, imino-, or hydroxy acids in the precursor chain. The ring sizes of cyclic structures are indicated in the number following C, L, or E, defining the type of ring closure. The abbreviations used for unusual amino acids and other compounds are listed in the abbreviations footnote on the first page.

Current applications of such enzyme systems are the enzymatic synthesis of peptide analogues, taking advantage of the relatively low stringency of the nonribosomal code—that is, the specificities of the peptide synthetases (9). Efforts are being made to alter the modular construction of peptide synthetases so as to generate new enzyme systems capable of synthesizing novel peptides.

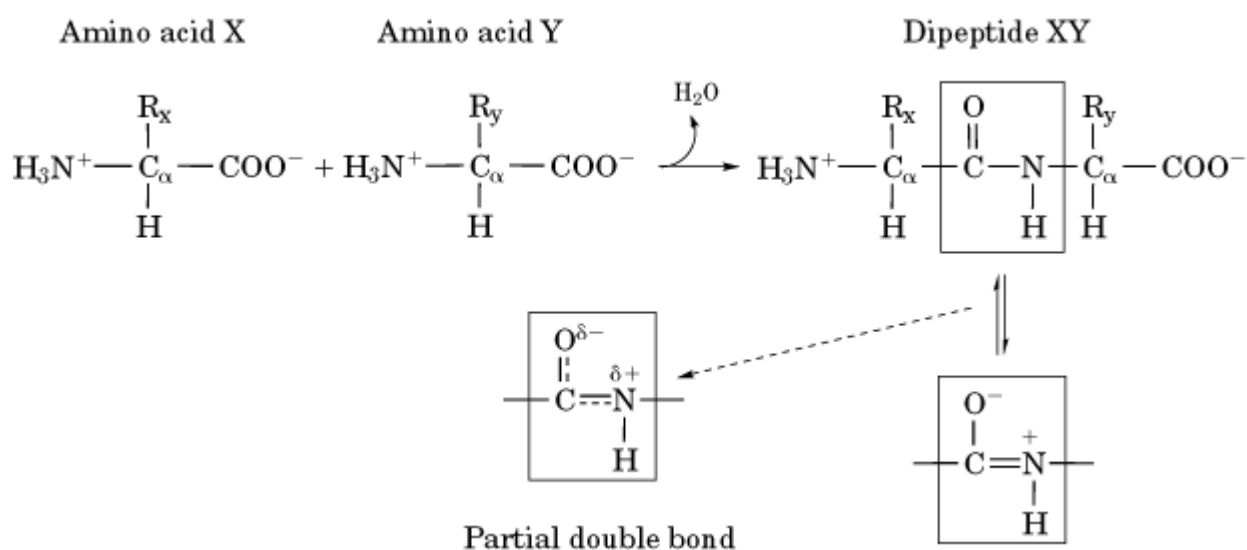
## Bibliography

1. H. Kleinkauf and H. von Döhren (eds.) (1990) *Peptide Antibiotics*, de Gruyter, Berlin.
2. H. Kleinkauf and H. von Döhren (1997) "Peptide Antibiotics". In H. Kleinkauf and H. von Döhren, (eds.), *Products of Secondary Metabolism*, Vol. 7 of H. Rehm, L. G. Reed, A. Pühler, P. Stadler (eds.), *Biotechnology*, 2nd ed., Verlag Chemie, Weinheim, pp. 277–322.
3. R. E. W. Hancock and R. Lehrer (1998) *Trends Biotechnol.* **16**, 82–88.
4. W. F. Broekaert et al. (1997) *Crit. Rev. Plant Sci.* **16**, 297–323.
5. C. L. Bevins et al. (1996) *Genomics* **31**, 95–106.
6. S. D. Heck et al. (1995) *Proc. Natl. Acad. Sci. USA* **93**, 4036–4039.
7. R. W. Jack, F. Götz, and G. Jung (1997) "Lantibiotics". In H. Kleinkauf and H. von Döhren (eds.), *Products of Secondary Metabolism*, Vol. 7 of H. Rehm, L. G. Reed, A. Pühler, P. Stadler (eds.), *Biotechnology*, 2nd ed., Verlag Chemie, Weinheim, pp. 323–368.
8. M. A. Marahiel, T. Stachelhaus, and H. D. Mooz (1997) *Chem. Rev.* **97**, 2651–2674.

## Peptide Bond

The peptide bond is the chemical link that connects [amino acids](#) to form the [polypeptide chains](#) of [peptides](#) and [proteins](#). The repeating units of these linear [polymers](#), the 20 amino acids, are naturally occurring chemical entities comprising an [amino group](#), a [carboxyl group](#), a hydrogen, and a [side chain](#) of variable chemistry. These four groups are connected to a central carbon atom, denoted  $C_a$ ; (or  $C_{\alpha}$ ) (Fig. 1). Peptide bonds are formed by condensation of the carboxyl group of one amino acid with the amino group of another amino acid, giving rise to a dipeptide. Repeated condensation of amino acids produces tripeptides, tetrapeptides, and so on. Proteins can incorporate many thousands of peptide-linked amino acids.

**Figure 1.** General scheme for formation of a dipeptide XY from two amino acids, X and Y, by reaction of the carboxyl group of amino acid X with the amino group of amino acid Y. The R group attached to  $C_a$ ; is variable, depending on the type of amino acid ( $R_x$  and  $R_y$ ). The peptide bond of the dipeptide product is boxed, and its partial double bond character is depicted.

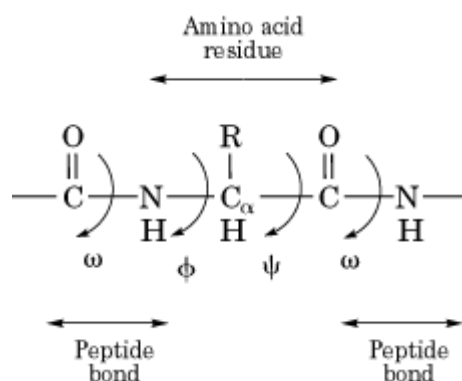


*In vivo*, peptide bond condensation occurs at the [C-terminus](#) of the growing polypeptide chain to produce proteins having specific sequences defined in the [genetic code](#). This reaction is catalyzed in a process called [translation](#), by [enzymes](#) of the cell's **protein biosynthesis** machinery. Synthetic peptides and proteins can also be produced by chemical [peptide synthesis](#). In contrast to *in vivo* protein biosynthesis, **solid-phase synthesis** usually occurs by condensation of amino acids to the [N-terminus](#) of the growing polypeptide chain. Peptide bonds can be broken down, or hydrolyzed, by enzymes such as [peptidases](#) and [proteinases](#).

Each peptide unit in a polypeptide chain has three chemical bonds. The central peptide bond is not

freely rotatable, because it has partial double bond character (Fig. 1). This imposes planarity on the peptide bond, so that only two discrete conformations are available. These two conformations, defined by the angle  $\omega$  (Fig. 2), are called the [cis](#) and [trans configurations](#). Of these two, the *trans* configuration is energetically preferred, and is much more commonly observed, unless the following residue is [proline](#) (see [Cis/Trans Isomerization](#)).

**Figure 2.** Section of a polypeptide structure showing two peptide bonds ( $-\text{CONH}-$ ) and one amino acid residue ( $-\text{NH}-\text{C}_\alpha;(\text{R})(\text{H})-\text{CO}-$ ). The dihedral angles  $\omega$ ,  $\phi$ , and  $\psi$  of the peptide backbone are labeled.



The single bonds on either side of the peptide bond are rotatable. The angle around the  $\text{N}-\text{C}_\alpha$  bond of a peptide is called phi ( $\phi$ ), and the angle of the  $\text{C}_\alpha-\text{C}$  bond is called psi ( $\psi$ ). Certain combinations of  $\phi$  and  $\psi$ , defined as the “allowed” regions of a  $\phi/\psi$  or [Ramachandran Plot](#), are preferred and energetically favored over others. The allowed  $\phi/\psi$  angles define the [conformations](#) of regular **secondary structure** such as [a-helix](#) or [b-strand](#) that are stabilized by backbone [hydrogen bonds](#) between peptide bond units.

[See also [Polypeptide Chain](#) and [Protein Structure](#).]

#### Suggestions for Further Reading

A. G. Walton and J. Blackwell (1973) *Biopolymers*, Academic Press, New York.

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

## Peptide Hormones

The first identified peptide hormone, secretin, was found in 1902 by Baylis and Starling. They discovered that the bioactive compound stimulated the secretion of saliva, and thus they named the factor secretin. Because the compound was produced by one tissue and transported by the blood to a distant organ, where it stimulated the saliva secretion, it was called a hormone based on the Greek word *hormao*, which means “to stimulate”. Since then it has been found that many physiological functions in the body are controlled by such communication with the endocrine and neural systems (1). That is, the regulation, integration, or coordination of various metabolic activities and bodily functions of tissues and organs involve the synthesis of low concentrations (micromolar to



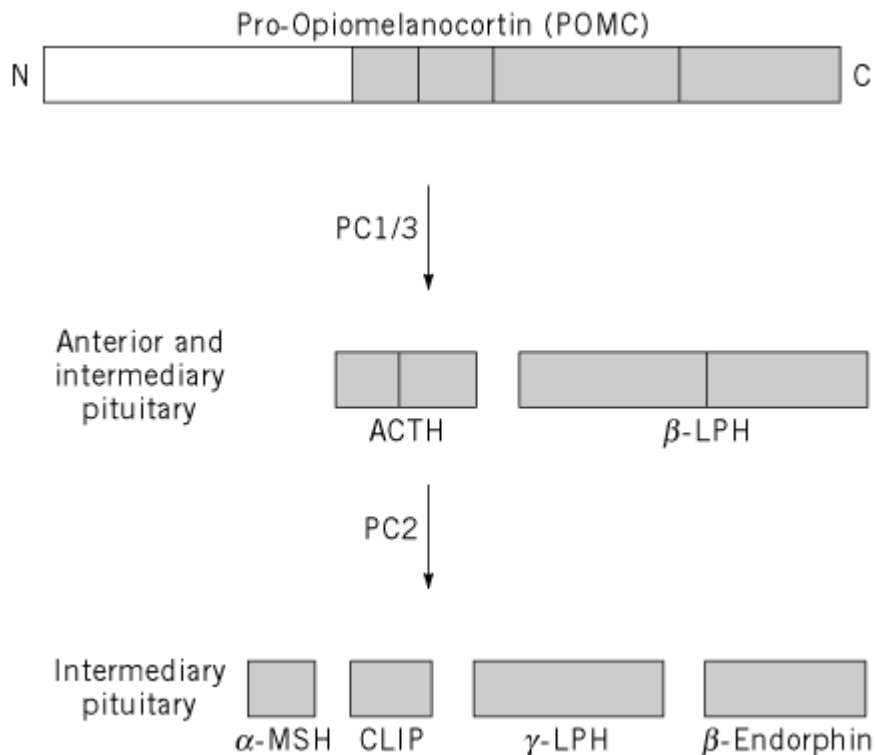
picomolar) of biochemical signals or messengers that can be disseminated and transported throughout the organism, to effect the required changes through **receptors** at distant target cells. Such regulation of function also requires that the production of the hormone be self-limiting; consequently, it requires both negative and positive feedback loops and interactive modulation of G-coupled systems (see [GTP-Binding Proteins](#)) of neurotransmitter (see [Neuropeptides](#)) systems to affect functional balance, or [homeostasis](#). Through modulation, downregulation, desensitization, upregulation, and receptor cross talk, hormones in some cases interact with the central nervous system (CNS) to influence neuronal function, to alter behavior, and to control other hormones. In other cases, one function can be influenced by several different hormones. All operate primarily under the hierarchical control of the central nervous, endocrine, and **immune systems**. Although hormones were originally defined as material that is secreted by endocrine glands, the term is now often applied to hormones that are expressed in other tissues or are formed from precursors located in other tissues. For example, angiotensin, which maintains the water and electrolyte balance and controls blood pressure, is produced by converting angiotensinogen in the blood with renin secreted by the juxtaglomerular cells in the kidney.

## 1. Biosynthesis, Secretion, and Interaction with Receptors

A peptide hormone is first synthesized on the [ribosome](#) as a pre-pro-hormone following [transcription](#) and [translation](#) of the genetic information (2) of the [messenger RNA](#). This precursor of high molecular weight contains an *N*-terminal [signal peptide](#), which guides the protein to the rough [endoplasmic reticulum](#) (ER), where the signal peptide is removed proteolytically by a signal peptidase, in a form of co-translational processing. The pro-hormone, or hormone precursor (see [Polyproteins](#)), that is produced usually does not possess much biological activity. Most pro-hormones and pro-neuropeptides contain within their polypeptide chains several smaller peptides with various biological activities. These pro-hormones normally undergo **posttranslational processing**, **protein folding**, [disulfide bond](#) formation, or chemical modification, such as [N-glycosylation](#), [O-glycosylation](#), **acetylation**, or [amidination](#), when they are processed through the [Golgi apparatus](#). The pro-hormones undergo specific proteolytic cleavage by endoproteinases or processing enzymes and then are sorted (3) into **secretory granules**. The polypeptide hormone is stored in the secretory vesicles and, in some cases, may undergo further processing by enzymes.

Almost all of the cleavage sites in the precursor protein contain a dibasic sequence, Lys–Arg or Arg–Arg, that is recognized by the processing enzymes. The **kex-2** proteinase of yeast, kexin, was the first enzyme found to cleave at Lys–Arg and Arg–Arg sequences. Subsequent use of the *kex-2* gene as a probe (see [Probe Hybridization](#)) led to the [cloning](#) of similar gene products from mammalian cells (4). Investigations demonstrated that these convertases, which cleave at the *C*-terminal side of Lys–Arg and Arg–Arg sequences, play an important physiological role. Differential processing of a single precursor protein in two different tissues can now be explained by the presence or relative concentrations of different posttranslational enzymes in tissues. For example, pro-opiomelanocortin (POMC) is differentially processed in the anterior and intermediary pituitaries (5). In the anterior and intermediary pituitary, ACTH and b-LPH are the predominant products because PC-1/3 (pro-protein convertase) is most prevalent. In the intermediary pituitary, however, alpha-melanocyte-stimulating hormone (a-MSH), corticotropin-like intermediary peptide (CLIP), and b-endorphin are secreted because PC-2 is more prevalent there (see Fig. 1).

**Figure 1.** Processing of the pro-opiomelanocortin (POMC), to produce ACTH (adrenocorticotropic hormone), B-LPH (b-lipotropin), A-MSH (a-melanocyte stimulating hormone), CLIP (corticotropin-like intermediary peptide), g-LPH, and b-endorphin. The pro-protein convertase PC-1 / 3 carries out the initial cleavages in the anterior and intermediary pituitary, whereas convertase PC-2 is present primarily in the intermediary pituitary.



The hormone within the secretory vesicles is released from the vesicles by [exocytosis](#) upon stimulation by extracellular signals provided by neurotransmitters or other hormones, whereupon the hormone is transported by the blood or extracellular fluid to the target cells. Once a hormone is secreted, it usually binds to a specific plasma protein carrier, which **prolongs** its half-life because the complex reduces the rate of hormone destruction by plasma proteinases.

The release of the peptide is normally controlled by a feedback mechanism. In the case of ACTH, for example, the secretion of corticoid is stimulated by ACTH, but on the other hand, corticoid suppresses CRF and ACTH secretion from the hypothalamus. In turn, ACTH suppresses CRF secretion.

The type of endocrine hormone secretion is classified by the distance between the cells that secrete the peptide and its target cells. Classically, endocrine hormones act at a distant target, while paracrines act upon nearby cells, and an autocrine acts upon the cell that secreted it. Many hormones often function in all three ways, however, making such distinctions difficult. Hormones secreted by neural cells and neurons are called neuroendocrines.

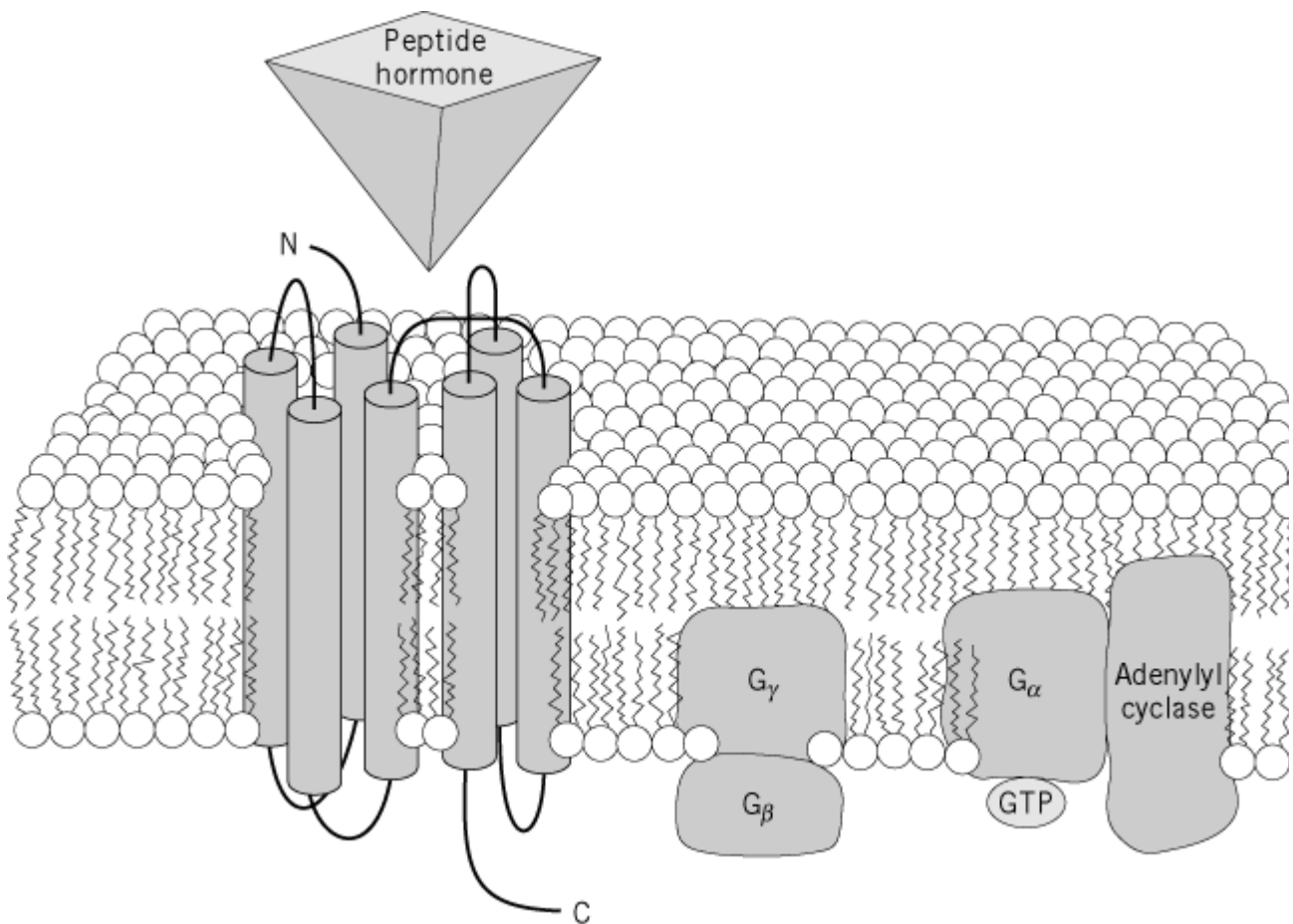
When the hormone reaches its target cell, it binds to its specific receptor in the cell [membrane](#). Each receptor is a membrane protein or complex and consists of three domains, (1) extracellular, (2) transmembrane, and (3) intracellular domains. The [tertiary structure](#) of the receptor is essential for specific binding of the ligand. This interaction induces structural changes in the receptor that result in activation of an intracellular enzyme that generates a transient increase in concentration of a signaling compound (see [Signal Transduction](#)), such as [cyclic AMP](#) (cAMP), [cyclic GMP](#) (cGMP), **diacylglycerol** (DAG), or inositol triphosphate, IP<sub>3</sub>. The elevated concentration of these [second messengers](#) and of the **phosphorylated** forms of proteins triggers the appropriate metabolic response. Proteins are phosphorylated on [serine](#), [threonine](#), or [tyrosine](#) residues by cAMP-dependent protein kinase (PKA), DAG-activated protein kinase C (PKC), or **tyrosine kinases**.

Many hormone receptors have a common structure in which the single polypeptide chain passes back and forth through the membrane seven times, forming seven transmembrane [alpha-helices](#). For these

receptors, the hormone-binding signal is transmitted or transduced to the interior of the cell by binding the receptor-ligand complex to a series of G-proteins (see [Heterotrimeric G Proteins](#)). G-proteins bind GDP or GTP. When G-proteins bind to their receptors, GTP exchanges with GDP bound to the  $\alpha$ -subunit of the G-protein. The  $G\alpha$ -GTP complex binds [adenylate cyclase](#), activating the enzyme. The activation of adenylate cyclase produces cAMP in the cytosol and leads to the activation of PKA and the phosphorylation of enzymes. Many classes of G-proteins exist, such as  $G_s$  and  $G_i$ , stimulatory and inhibitory G-proteins, respectively. A second class of G-protein coupled-receptor synthesizes DAG and IP3 as second messengers. The hormone-receptor complex interacts with G-proteins, and this is followed by the G-protein activation of phospholipase C-g (PLCg). PLCg hydrolyzes phosphatidylinositol bisphosphate to IP3 and DAG. IP3 binds to the endoplasmic reticulum and opens the calcium channels so that stored calcium flows into the cytosol and activates enzymes. On the other hand, the DAG binds and activates PKC and also opens **calcium channels**. PKC phosphorylates serine and threonine residues, thereby modulating the catalytic effect of the corresponding enzymes.

One receptor, that of ANF (atrial natriuretic factor), is coupled to the production of cGMP by guanylate cyclase. Another ANF receptor lacks guanylate cyclase activity and is thought to play a role in the clearance of ANF. Yet another type of receptor is the heterotetrameric [insulin](#) receptor. It is composed of two  $\alpha$ - and two  $\beta$ -subunits. The  $\alpha$ -subunits protrude outside the membrane and contain the binding sites, while the  $\beta$ -subunits penetrate through the membrane in a single transmembrane  $\alpha$ -helix and contain the tyrosine-kinase activity. When insulin binds to the receptor, intramolecular phosphorylations are triggered, and this autophosphorylation enhances its ability to phosphorylate enzymes that induce the physiological response (Fig. 2) (see also [Heterotrimeric G Proteins](#)).

**Figure 2.** General mechanism of action of peptide hormones that bind to receptors and thereby alter the activity of adeny cyclase through heterotrimeric G-proteins. The peptide hormone binds to the extracellular part of the receptor embedded plasma membrane. This is transmitted to the  $\alpha$  subunit of the heterotrimeric G-protein, which dissociates in the  $G\beta\gamma$ -dimer and the  $G\alpha$  subunit, which binds to adenylate cyclase and alters its activity accordingly.



## 2. Methods of Isolation and Structure Determination

The structure of a peptide hormone was first demonstrated by de Vigneaud in 1953, when he determined the structure of oxytocin and confirmed it by total chemical synthesis. In 1954, the structure of insulin was determined by Sanger. The chemical structure of secretin was not determined until 1966, 60 years after its initial discovery. Because peptide hormones exist in such small concentrations, it was difficult then to isolate, purify, and analyze them. Subsequent development of better detection methods with higher sensitivities, such as [radioimmunoassay](#) (RIA), made possible the rapid isolation and structure determination of a large number of peptide hormones and neuropeptides. Later gene [cloning](#) provided the sequence of the [complementary DNA](#) of hormone precursors by using a DNA probe (see [Probe Hybridization](#)) based on the known peptide sequence. This also provides the structure of previously undetected protein **isoforms**. For example, in the case of endothelin-1 (ET-1), which had been isolated from cultured endothelium cells of blood vessels, screening of the human chromosomal DNA library revealed the existence of genes for three isoforms: ET-1, ET-2, and ET-3. Each isoform consists of 21 amino acid residues, and they have high sequence [homology](#). Another method of detecting new peptide hormones is by screening the DNA sequence of the precursor protein. In this case, another new peptide hormone that is linked in the precursor molecule can be detected. In the precursor of [glucagon](#), two amino acid sequences that resemble glucagon were found, GLP-1 and GLP-2. Another example is the case of calcitonin gene-related peptide (CGRP), which stimulates the formation of cAMP. Analysis of the gene of the precursor protein of calcitonin (CT), which modulates serum calcium concentrations, revealed that the CT gene consists of five **introns** and six exons. As a consequence of [alternative splicing](#) of the RNA transcripts (see [Transcription](#)), the mRNA for the CT precursor is found in the thyroid but the mRNA for CGRP is formed in the nervous system.

As a result of the progress in gene technology, the chemical structures of a large number of hormone receptors have now been determined. Most of their intrinsic ligands also have been identified. However, there are still receptors for which no ligands are known. New techniques are being developed to search for new bioactive peptides that serve as ligands for such orphan receptors. In 1998 (6), a prolactin-releasing factor or peptide was found in the search for the ligand for one such orphan receptor that is expressed in the human pituitary. The 31-residue peptide stimulates the pituitary to release prolactin.

### 3. Synthesis and Conformational Analysis

After the isolation and structure determination of a peptide hormone has been achieved, the peptide is often chemically synthesized to confirm the structure. With the synthesized material, structure-activity relationships can be explored to design drugs for treating individuals with peptide-hormone deficiencies or with malfunctioning peptides. The technique of **combinatorial chemistry** has recently made it possible to synthesize millions of compounds and to screen them for activity rapidly and efficiently (7). In addition, gene technology has provided a means of producing human peptide sequences, such as insulin, in *Escherichia coli*, yeast, or other [expression systems](#). This has afforded the opportunity to produce peptide hormones and receptors that are uniformly labeled with  $^{13}\text{C}$  and  $^{15}\text{N}$  isotopes for nuclear magnetic resonance ([NMR](#)) analysis. These labeled compounds can then be used to investigate the interaction between ligand and receptor, as well as to study the conformational changes that occur upon binding. Consequently, the ongoing study of the tertiary structures by NMR and [X-ray crystallography](#) is now providing more details regarding the mechanisms of hormonal action.

#### Bibliography

1. A. W. Norman and G. Litwack, editors (1987, 1997) *Hormones*, 1st and 2nd eds., Academic Press, New York.
2. W. W. Chin (1995) In *Principles and Practice of Endocrinology and Metabolism*, 2nd ed., J. B. Lippincott, Philadelphia, pp. 8–41.
3. M. J. Perone, S. Windeatt, and M. G. Castro (1997) *Histology and Histopathology* **12**, 1179–1188.
4. S. P. Smeeckens (1993) *Bio/Technology* **11**, 182–186.
5. M. G. Castro and E. Morrison (1997) *Critical Rev. Neurobiol.* **11**, 35–57.
6. S. Hinuma (1998) *Nature* **393**, 272–276.
7. C. Pinilla et al. (1995) *Biopolymers* **37**, 221–240.

#### Peptide Libraries

Peptide libraries were first introduced during the mid-1980s as a means to elucidate the details of **antibody–antigen** recognition (1, 2). Prior to that time, [epitope](#) mapping required fragmentation analysis, either at the protein or DNA level, which was a laborious and time-consuming process. Geysen et al. (1) demonstrated that synthetic peptide libraries could be prepared and screened by [enzyme-linked immunosorbent assay](#) (ELISA) to elucidate antibody–epitope recognition to the resolution of a single amino acid residue. This was the first case in which a problem in molecular recognition was elucidated using a synthetic library and is often cited as a critical development in the emergence of [combinatorial libraries](#).

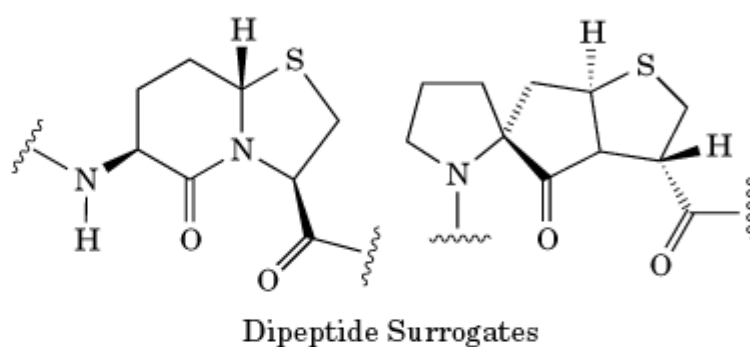
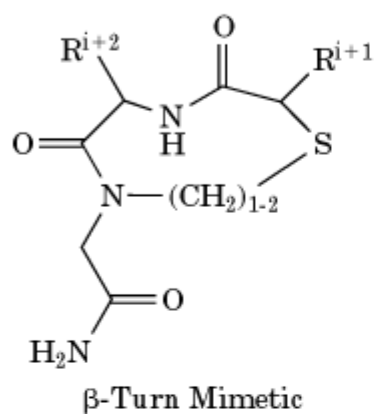
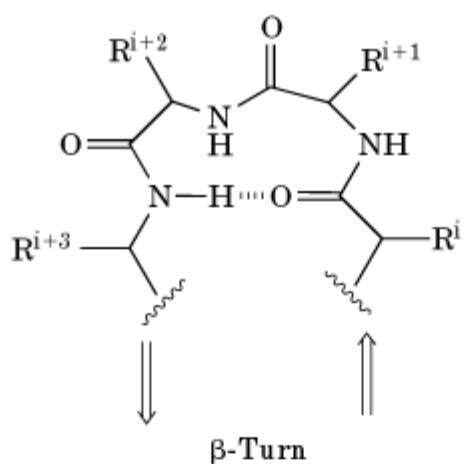
The first generation of synthetic peptide libraries was limited in size to a few hundred peptides, due to the need for each library member to be prepared as a discrete compound. However, the combination of mixture-synthesis techniques (14), reliable strategies for identifying active species from complex mixtures (3, 4), and advances in peptide sequencing methods quickly led to the routine preparation and screening of libraries containing millions of distinct peptide sequences. Several reports of ten-residue peptide libraries containing more than  $10^{12}$  peptides have even appeared (5, 6). When coupled with an increasing availability of unnatural amino acid building blocks, peptide libraries of this scale and complexity became valuable tools for identifying novel peptide agents for a wide range of applications (eg, see Refs. 7 and 8). Soluble and solid-phase peptide libraries continue to find applications in the identification of **antigenic** determinants. In addition, they have been used to identify receptor ligands (5, 9-14), inhibitors of enzymatic activity (15), and novel antibacterial agents (3, 7, 14). Many of these applications have employed standard assay formats (eg, soluble activity or binding assays), requiring complex deconvolution or iterative synthesis and screening cycles to identify the active species from complex synthetic pools (see [Combinatorial Synthesis](#)).

Split synthesis methods for preparing one-bead, one-compound (OBOC) peptide libraries have overcome many of the limitations of standard assay formats and pooling strategies. Because each bead carries a single compound in amounts approaching several nanomoles, assays can be performed directly on the resin beads that detect interactions of bead-linked compound with a soluble target. A typical assay utilizes ELISA methodologies to indicate beads bearing compounds that interact with the target (16). Positive beads may be selected manually, or more sophisticated instrumentation may be employed to detect and isolate active beads (17). Alternatively, the OBOC format is amenable to off-bead assays through the use of cleavable linkers inserted between the peptide and the solid-phase support (18, 19). In an off-bead assay, the peptide is released from the resin bead, but proximity to the bead is maintained through use of a diffusion-limiting matrix or a positional array. For example, distribution of the library and a test cell line within a soft agar matrix, followed by partial peptide release, leads to the formation of locally high concentrations of peptide in the agar surrounding each bead. Zones of reduced cell growth (or any other reporting mechanism such as color change) surrounding single beads are used to identify peptides having the desired activity. Beads producing a zone of activity are removed and analyzed by standard [Edman Degradation](#) or **mass spectrometric** analysis to determine the sequences of the active peptides. This approach has been used for identifying peptide ligands for [G-protein-coupled receptors](#) (20), as well as cytotoxic peptides that may have application as anticancer agents (19). This represents but a few examples of the unique screening methods that have been developed using peptide chemistry coupled with the OBOC library format. Other notable examples include the use of a fluorescence-activated cell sorter (FACS) (17) (see [Flow Cytometry](#)) to identify active beads, as well as direct visualization of cancer cells binding to peptides on beads (21).

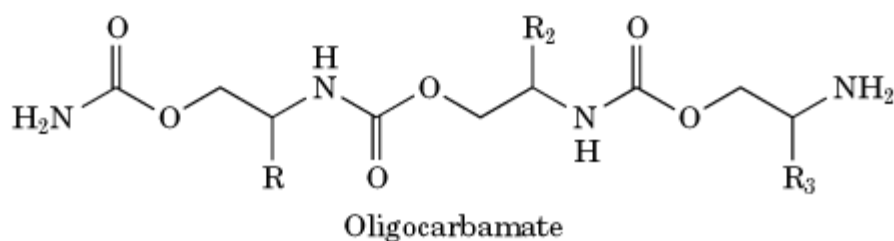
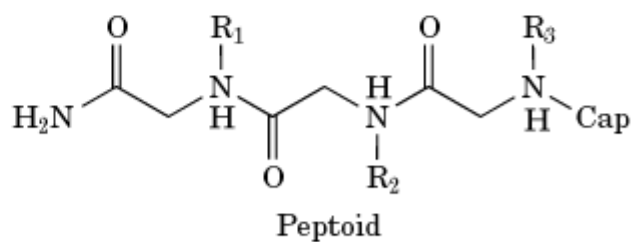
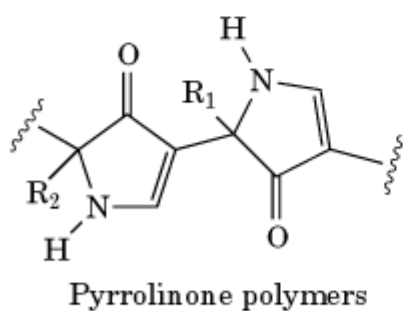
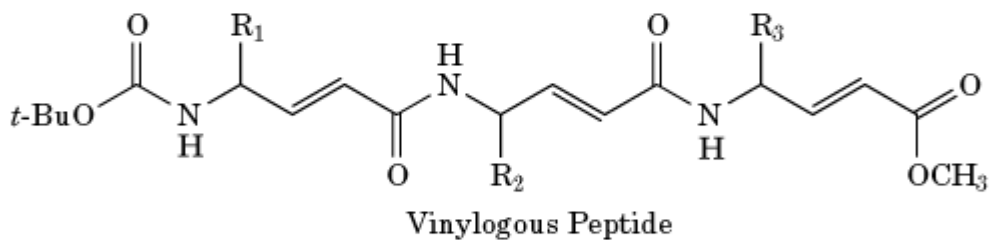
The initial excitement surrounding synthetic peptide libraries has been tempered in recent years. Peptides are linear, highly flexible polymers with fixed side-chain spacing. In contrast, most proteins adopt complex three-dimensional **structure** that place active chemical groups in positions that cannot be approximated by linear peptides. As a result, many targets fail to yield hits from short peptide libraries. The introduction of constraints into the peptide chain, including cyclization or use of domain scaffolds with defined structural propensities (Fig. 1a), may increase the success rate for some targets. In addition, altering the polymeric backbone can change the spacing of the side chains (22-25). Although unnatural amino acids can be used to shift the side-chain spacing (eg,  $\beta$ -amino acids), the problems enumerated above associated with a linear and flexible polymer remain. A wide range of peptidomimetics that employ novel backbones have emerged to address this issue (Fig. 1b) (26-28). Such compounds also address a second major criticism of peptides: They are ill-suited for *in vivo* applications due to poor pharmacokinetic properties. By utilizing building blocks with physicochemical properties compatible with uptake and adsorption, peptidomimetics have moved a step closer to lead compounds for drug discovery. The true utility of peptidomimetics for drug

discovery remains unknown, but their potential for use as research tools is high.

**Figure 1.** Structural aspects of peptide libraries. (a) Illustration of one structural bias, a **b-turn** (top), that can be introduced in a peptide library, and two approaches for creating libraries based on this structural scaffold. They are cyclization (bottom left) and the use of constrained dipeptide surrogates (bottom right). (b) Examples of various peptidomimetics.



(a)





See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#); and [Phage Display Libraries](#).

## Bibliography

1. H. M. Geysen, R. H. Meleon, and S. J. Barteling (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3998–4002.
2. R. A. Houghten (1985) *Proc. Natl. Acad. Sci. USA* **82**, 5131–5135.
3. R. A. Houghten, C. Pinilla, S. E. Blondelle, J. R. Appel, C. T. Dooley, and J. H. Cuervo (1991) *Nature* **354**, 84–86.
4. C. Pinilla, J. R. Appel, P. Blanc, and R. A. Houghten (1992) *Biotechniques* **13**, 901–905.
5. R. A. Houghten and C. T. Dooley (1993) *NIDA Res. Monogr.* **134**, 66–74.
6. C. Pinilla, J. R. Appel, and R. A. Houghten (1994) *Biochem. J.* **301**, 847–853.
7. S. E. Blondelle, E. Takahashi, P. A. Weber, and R. A. Houghten (1994) *Antimicrob. Agents Chemother.* **38**, 2280–2286.
8. K. S. Lam, M. Lebl, V. Krchnak, S. Wade, F. Abdul-Latif, R. Ferguson, C. Cuzzocrea, and K. Wertman (1993) *Gene* **137**, 13–16.
9. C. T. Dooley, N. N. Chung, B. C. Wilkes, P. W. Schiller, J. M. Bidlack, G. W. Pasternak, and R. A. Houghten (1994) *Science* **266**, 2019–2022.
10. R. A. Houghten (1993) *Gene* **137**, 7–11.
11. C. T. Dooley, R. A. Kaplan, N. N. Chung, P. W. Schiller, J. M. Bidlack, and R. A. Houghten (1995) *Pept. Res.* **8**, 124–137.
12. C. T. Dooley, N. N. Chung, P. W. Schiller, and R. A. Houghten (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10811–10815.
13. C. T. Dooley and R. A. Houghten (1993) *Life Sci.* **52**, 1509–1517.
14. R. A. Houghten, J. R. Appel, S. E. Blondelle, J. H. Cuervo, C. T. Dooley, and C. Pinilla (1992) *Biotechniques* **13**, 412–421.
15. J. Vagner, G. Barany, K. S. Lam, V. Krchnak, N. F. Sepetov, J. A. Ostrem, P. Strop, and M. Lebl (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8194–8199.
16. K. S. Lam, M. Lebl, and V. Krchnak (1997) *Chem. Rev.* **97**, 411–448.
17. M. C. Needels, D. G. Jones, E. H. Tate, G. L. Heinkel, L. M. Kochersperger, W. J. Dower, R. W. Barrett, and M. A. Gallop (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10700–10704.
18. M. Lebl, M. Patek, P. Kocis, V. Krchnak, V. J. Hruby, S. E. Salmon, and K. S. Lam (1993) *Int. J. Pept. Protein Res.* **41**, 201–203.
19. S. E. Salmon, R. H. Liu-Stevens, Y. Zhao, M. Lebl, V. Krchnak, K. Wertman, N. Sepetov, and K. S. Lam (1996) *Mol. Diversity* **2**, 57–63.
20. C. K. Jayawickreme, G. F. Graminski, J. M. Quillan, and M. R. Lerner (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1614–1618.
21. M. E. Pennington, K. S. Lam, and A. E. Cress (1996) *Mol. Diversity* **2**, 19–28.
22. R. J. Simon, R. S. Kania, R. N. Zuckermann, V. D. Huebner, D. A. Jewell, S. Banville, S. Ng, L. Wang, S. Rosenberg, C. K. Marlowe, et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9367–9371.
23. A. B. Smith III, T. P. Keenan, R. C. Holcomb, P. A. Sprengeler, M. C. Guzman, J. L. Wood, P. J. Carroll, and R. Hirschmann (1992) *J. Am. Chem. Soc.* **114**, 10672–10674.
24. C. Y. Cho, E. J. Moran, S. R. Cherry, J. C. Stephans, S. P. A. Fodor, C. L. Adams, A. Sundaram, J. W. Jacobs, and P. G. Schultz (1993) *Science* **261**, 1303–1305.

25. M. Hagihara, N. J. Anthony, T. J. Stout, J. Clardy, and S. L. Schreiber (1992) *J. Am. Chem. Soc.* **114**, 6568–6570.
26. M. A. Gallop, R. W. Barrett, W. J. Dower, S. P. A. Fodor, and E. M. Gordon (1994) *J. Med. Chem.* **37**, 1233–1251.
27. E. M. Gordon, R. W. Barrett, W. J. Dower, S. P. A. Fodor, and M. A. Gallop (1994) *J. Med. Chem.* **37**, 1385–1404.
28. L. A. Thompson and J. A. Ellman (1996) *Chem. Rev.* **96**, 555–600.

## Peptide Mapping

Peptide mapping allows rapid overall evaluation of a **protein**'s [primary structure](#) from limited analysis of the **peptides** produced by **proteolysis**. Only a single [amino acid](#) difference between two peptides caused by [mutation](#) or **posttranslational modification** alters the physical properties of the peptide in which that difference occurs. So peptide mapping is a very sensitive way of detecting differences between closely related proteins. The peptides produced by a specific [proteinase](#) are separated into a pattern ("peptide map"), and the maps of related proteins are compared. Single differences, such as the Glu/Val difference between normal/**sickle cell hemoglobins** that was early recognized by peptide mapping (1) or the relationships between **isozymes** are found rapidly (2). Subsequent refinements in the methods and resolution of mapping produced excellent one- and two-dimensional (2-D) techniques that have been reviewed (3) and much used. The classical peptide mapping methods of two-dimensional **fingerprinting** and [Cleveland maps](#) are described in those articles. These classical approaches are still practical when the amounts of protein are sufficient, but molecular biology now centers on different approaches.

One extension increases the sensitivity, as required by present-day and future needs, to evaluate proteins available only in very small amounts. This has become possible by modern separation techniques, such as capillary **HPLC** (4), **capillary electrophoresis** (5), and high-resolution [two-dimensional gel electrophoresis](#) (6, 7). A further extension is to rely on [mass spectrometry](#), using laser desorption instruments, triple quadrupoles, ion-trap devices, or other techniques, and perform the mapping, detection, and analysis in one step. In this manner, using entire peptide mixtures or preselected fragments, direct determination of the total masses of peptide fragments from a protein gives a unique "fingerprint" and is expected soon to be the mapping method of general choice. Combined with [databases](#) of peptides from all known proteins, it allows direct identification of any protein that has been characterized from the same or a closely related source (8-10). Then, mass spectrometric peptide mapping and simultaneous data bank screening will constitute a powerful technique as a descendant of the classical peptide mapping idea. In the near future, when crucial **genomes** are completely sequenced, the work remaining will center on identifying of each gene product and on [proteome](#) analysis to correlate **gene expression** with predicted protein structures stored in data banks.

### 1. Liquid chromatography

When reverse-phase HPLC was introduced in the mid-1970s, column **chromatographic** methods largely replaced paper and **thin-layer** plate separations. Further miniaturization and additional separation techniques have since continuously altered the separation methods (4), but the principle of peptide mapping by comparing efficient separations of peptides remains. In many cases, the resolution now is so high that complex peptide mixtures are separated completely by single-column chromatographic steps, thereby reintroducing one-step mapping methods. Instruments are also so

reproducible that repeated runs, even at different times, are virtually identical, thereby introducing mapping comparisons on the same instrument on different occasions. The principle of peptide mapping has gained further momentum which is expected to continue in present and future stages of proteome research with all of its gene product characterizations.

The principle is still to compare two or more expressed proteins by suitable proteolytic digests, followed by identical instrument runs, and subsequently evaluating all peak positions and heights of the resulting patterns. The era of peptide mapping has moved from spots or bands on solid supports and multiple but simultaneous separations to peaks and positions on microcolumns after successive separations. In addition to detecting structural differences, many other applications of this technique are now in use, for example, [epitope](#) mapping, determination of **glycosylation**, **phosphorylation**, or oxidation sites, and assignment of [disulfide bonds](#). Peptide mapping is used extensively for quality control and characterization of **recombinant DNA** -derived products.

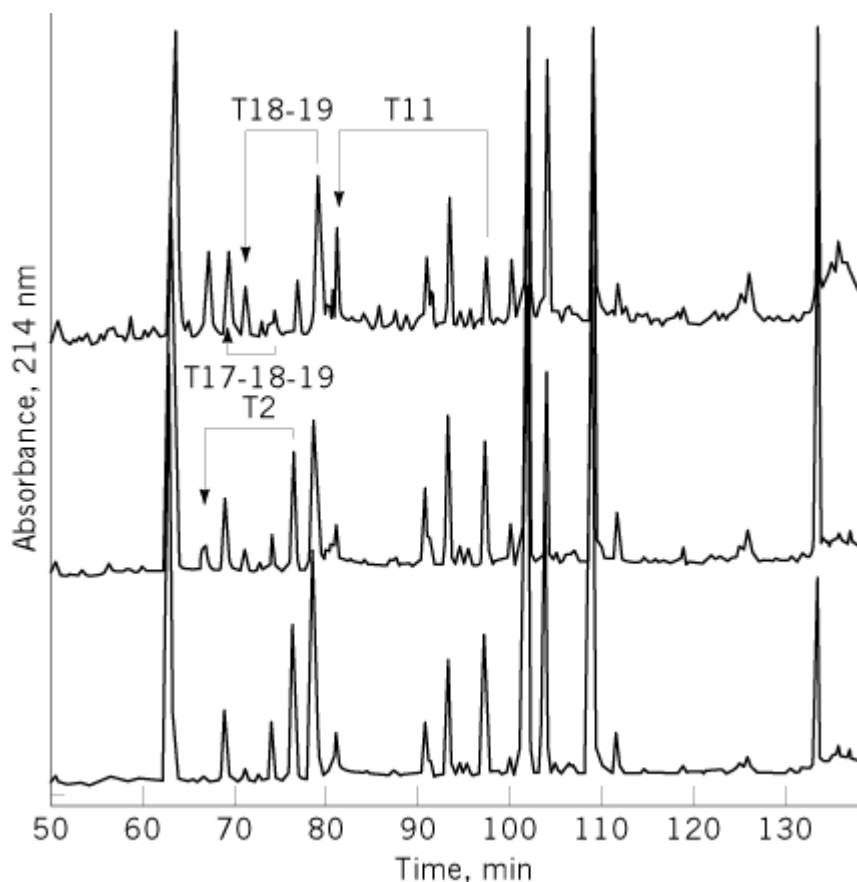
The major advantages in using microcolumn methods for peptide mapping are speed (often minutes), sensitivity (femtomole amounts, even with detection by **UV absorbance** ), and ease of recovery for further characterization by [amino acid analysis](#), [Edman Degradation](#), and mass spectrometry. Several high-performance separation media have been developed that allow rapid HPLC resolution of peptides and proteins. The classical silica-based stationary phase HPLC column has been used extensively for nearly twenty years in numerous applications, and alternative materials, such as divinylbenzene cross-linked polystyrene ([11](#)), have since emerged. Frequently, the stationary phase is packed in microbore columns with a narrow inner diameter (on the order of 1 to 2 mm), with low flow rates (down to 40  $\mu\text{L}/\text{min}$ ), small elution volumes (40 to 60  $\mu\text{L}$ ), better detector response (5- to 20-fold higher) than conventional HPLC columns, and requiring only picomole quantities of protein. However, the current trend is toward further miniaturization and smaller inner diameters, leading to still lower detection limits, but also toward implementation of powerful detection techniques, such as Z-shaped flow cells for efficient UV absorbance, laser-induced **fluorescence**, electrochemical detection, and mass spectrometry.

Capillary reverse-phase HPLC is conducted with column inner diameters of between 350 and 75  $\mu\text{m}$  and at flow rates of 4 to 0.2  $\mu\text{L}/\text{min}$ . The low flow rates necessitate special technical solutions for recovering separated peptides. Such a system has recently been developed [[MicroBlotter](#) ([12](#), [13](#))] for subpicomole peptide mapping after digestion of proteins in gels ([14](#), [15](#)) or on membranes after [blotting](#) ([16](#)) and for isolating peptides for microsequence analysis. The capillary outlet intermittently touches a strip of poly(vinylidene difluoride) (PVDF) membrane while moving along the strip. After sample collection, the PVDF strip is aligned with the chromatogram to excise spots that correspond to peptide candidates for structural analysis.

The protein to be analyzed is still digested with **trypsin**, Lys-C proteinase, or another suitable proteolytic enzyme. Protocols for microdigestion in small volumes (1 to 10  $\mu\text{L}$ ) and for recovering peptides from small samples (a few picomoles or less) have been developed ([17](#)). After digestion, the proteolytic peptides are derivatized with fluorescent tags ([18](#)) for detection. Generally, however, separation problems may arise from precolumn derivatization, and the digest is often injected directly onto the column after adjusting the pH to that of the mobile phase at the start of chromatographing. The most widely used solvent systems still consist of gradients of acetonitrile (often 0 to 60%, v/v) in aqueous trifluoroacetic acid (TFA) (often 0.1% v/v), but several other chromatographic solvents have also been suggested, for example, heptafluorobutyric acid instead of TFA, hydrochloric acid, sodium [phosphate buffer](#) at acidic or neutral pH, and ammonium bicarbonate, ammonium acetate, triethylammonium phosphate, or 1.0% sodium chloride at neutral pH. The most common detection technique is still UV absorbance, normally at both 214 and 280 nm, the latter for identification of peptides that contain [tryptophan](#) and [tyrosine](#) residues. Alkylation of [cysteine](#) with [ $^{14}\text{C}$ ]- **iodoacetic acid** or 4-vinylpyridine is common to prevent disulfide formation and peptide aggregation and to promote convenient detection of cysteine-containing peptides in the resulting map. A practical example of capillary HPLC peptide mapping (Fig. [1](#)) illustrates how *in*

*in vivo* oxidation of a recombinant protein hormone after intravenous administration is detected (17).

**Figure 1.** Capillary HPLC tryptic peptide maps of recombinant human [growth hormone](#) (lower profile) and of the same protein recovered from rat serum taken 15 (middle profile) and 45 (upper profile) min after an intravenous dose of the preparation. Arrows link methionine-containing peptides (right in each pair) with their corresponding oxidized forms (left in each pair), showing the increase in oxidized form with time *in vivo*. The peptides were separated with a 320  $\mu\text{m}\times 15\text{cm}$  reverse-phase ( $\text{C}_{18}$ ) column. From Ref. 17 with permission.

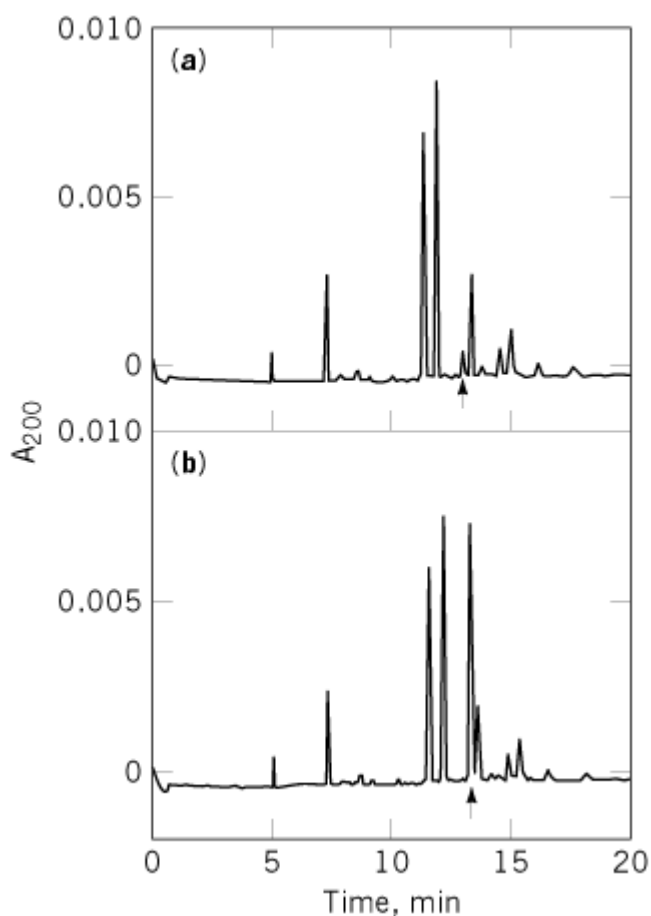


## 2. Capillary electrophoresis

Capillary [electrophoresis](#) in different separation modes, often [capillary zone electrophoresis](#) (CZE), provides a further powerful approach to peptide mapping (19), that has high sensitivity (low femtomole level) and short separation time (minutes). More than  $10^6$  theoretical plates are achieved and used to generate reproducible peptide patterns suitable for mapping experiments (Fig. 2). Capillary electrophoresis is also a suitable complement to reverse-phase HPLC for identifying peptides because the separation mechanisms are different, based on mass-to-charge ratio and **hydrophobic** interaction, respectively. Consequently, peptide fragments not resolved by reverse-phase HPLC are often completely resolved by CZE. Similarly, reverse-phase HPLC may resolve peptides not separable by CZE (21).

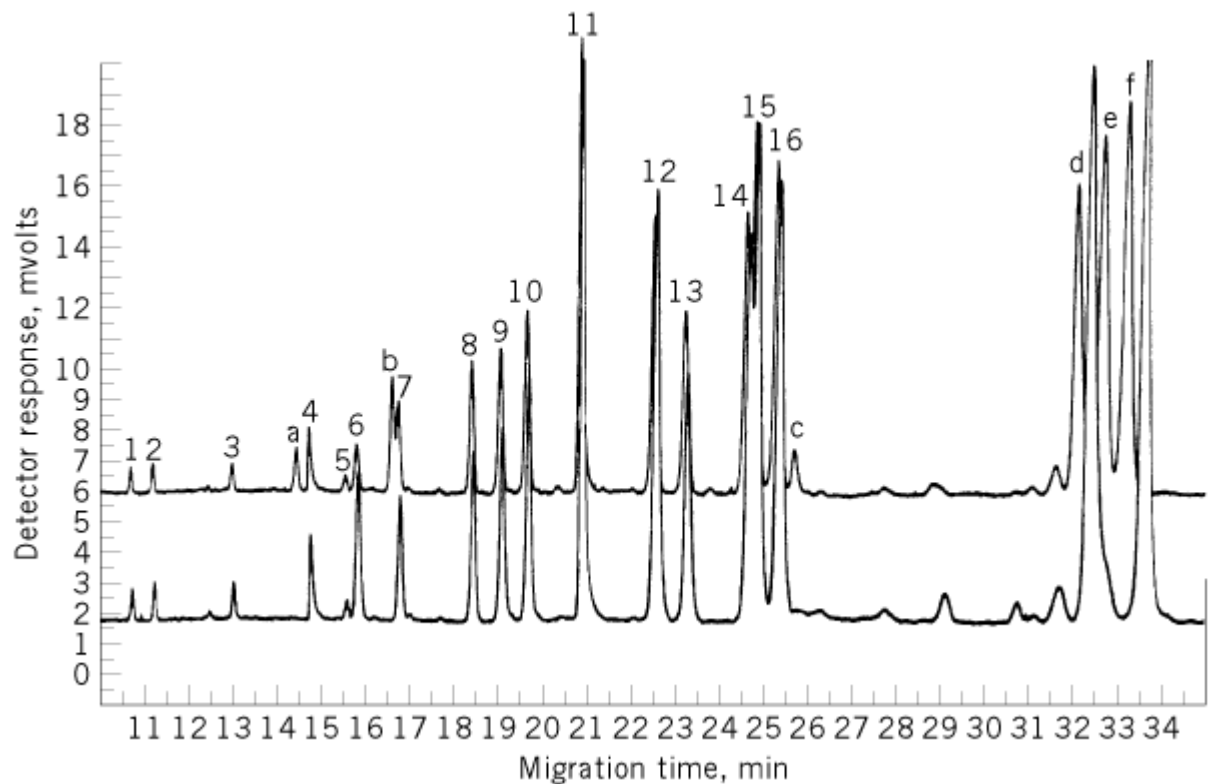
**Figure 2.** Peptide mapping of a fraction from a brain extract by capillary electrophoresis, to illustrate reproducibility with time and use of mapping for peptide identification. The separation in (a) reveals that the fraction contains a mixture of short peptides, the one at the arrow is suspected to be a neuroactive tripeptide, Gly-Pro-Glu. After adding of the corresponding synthetic peptide, reanalysis in (b) reveals that the synthetic peptide (arrow) comigrates with the

biological peptide component. Considering that the time between separations ( **a** and **b** ) is >2 weeks, the repetitive pattern confirms the reproducibility of the system on a sensitive scale (0.1 to 2 pmol analyzed). From Ref. [20](#) with permission.



Fused silica capillaries that have charged surfaces in which [electroendosmosis](#) takes place are commonly employed, typically with inner diameters of 20 to 200  $\mu\text{m}$  and lengths of 20 to 100 cm. Small inner diameters are best suited for peptide mapping because the more narrow the capillary, the better the resolution. For peptide separations, sodium phosphate buffer, pH 2.5, is a good electrolyte that allows sensitive detection at short wavelengths (190 to 200 nm) because of its high UV transparency. Including short-chain, ion-pairing reagents (hexanesulfonic or heptanesulfonic acid) in the [buffer](#) improves resolution, in particular of hydrophobic peptides, and has been tested in peptide mapping (Fig. [3](#)). The voltages applied range from 5 to 30 kV, and electric field strengths range from 200 to 500 V/cm. Separations are influenced by a combination of electrophoretic mobility and electroendosmosis. CZE can be used preparatively, and the components separated in the peptide map can be further characterized by [amino acid analysis](#), Edman degradation, or mass spectrometry ([5](#), [23](#)). In particular, electrospray mass spectrometry has gained recent popularity as an on-line detection technique with simultaneous tandem mass spectrometric capability for peptide structure information at the low femtomole level ([24-26](#)).

**Figure 3.** Tryptic peptide mapping of recombinant human erythropoietin in phosphate buffer, pH 2.5, with 100 mM heptanesulfonic acid by capillary electrophoresis. The maps from two different expression systems are shown, the upper without glycosylation (*E. coli*), the lower with (Chinese Hamster Ovary cells) but subsequently treated with N-glycanase. Although 16 peaks have identical migrations, indicating identical primary gene products, peaks a–f differ and can be attributed to differences at specific sites in the eukaryotic form. From Ref. [22](#) with permission.



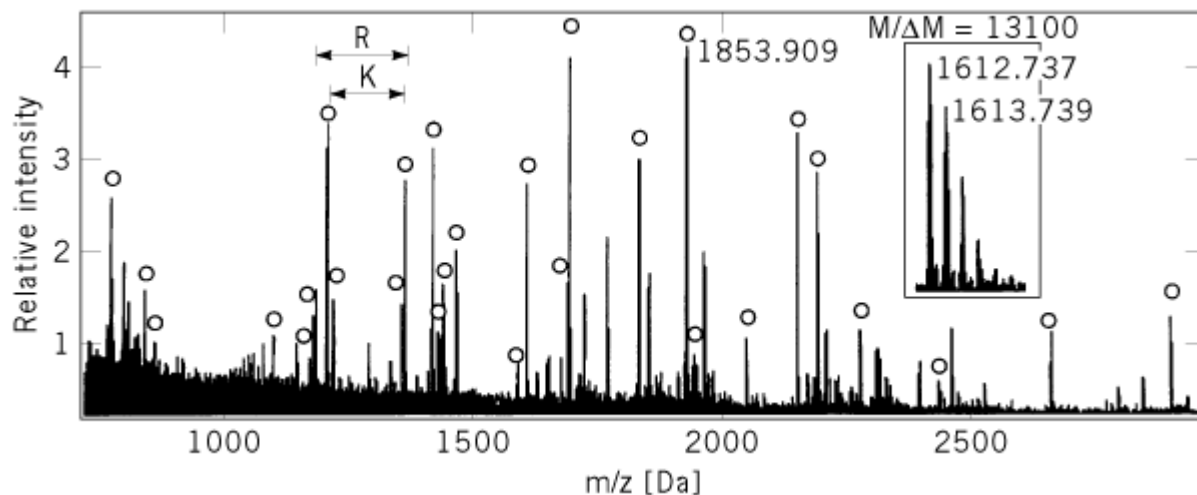
### 3. Mass spectrometry

Protein fingerprinting by determining the masses of all the peptides generated by proteolysis of polypeptide chains separated by [gel electrophoresis](#) has recently emerged as an extremely powerful tool for protein identification. It combines sensitivity and high throughput (8). Complex protein mixtures of crude cell lysates can be separated with 2-D gel electrophoresis systems (27). For identification, the amount of protein in many 2-D gel spots is well below 1 picomole and thus too small for the commonly used sequencer technology, but it is accessible to recently developed sequencer instruments (down to 100 femtomole by cartridge miniaturization and capillary HPLC for phenylthiohydantoin detection; see [Edman Degradation](#)) using matrix-assisted laser desorption/ionization (MALDI) (28) or electrospray (29) mass spectrometry, recently extended by nanoelectrospray ionization (30, 31). As a consequence, a novel type of “peptide mapping” of proteins has emerged, using their fragment masses and efficient algorithms to search peptide-mass databases, which permits structural assignment, protein identification, and determination of minor differences (32-34).

As in traditional peptide mapping, the first step in this procedure is proteolytic cleavage of the protein (commonly, still with trypsin), either in solution, in a gel (14, 15), or on a membrane after blotting on nitrocellulose or PVDF membranes (16). The MALDI mass analysis and the database search are highly automated (35), which increases the throughput and makes multiple protein identifications efficient (Fig. 4). If the protein is still unknown after the database search, the peptides are further analyzed by tandem mass spectrometry (36, 37). Using the nanoelectrospray ion source, sequence stretches up to 16 residues obtained from as little as 5 ng of protein digested with trypsin in a gel have been reported (10).

**Figure 4.** Peptide mapping by automated MALDI mass spectrometry and database searching after in-gel digestion with trypsin of a yeast protein recovered from a 2-D gel system. Ion signals whose measured masses match the

calculated masses of protonated tryptic peptides are indicated with circles. Sequence coverage is greater than 70%. Terminal Arg (R) and Lys (K) sequence tags are marked by arrows. Also shown is a magnification of one of the tryptic peptide peaks that reveals the resolution of molecules containing different natural isotopes. From Ref. [35](#) with permission.)



#### 4. Conclusion

Peptide mapping, complemented with Edman degradation and mass spectrometry, involves principles from the 1950s and earlier. But they have survived in the newest phases of modern molecular biology and now position themselves at the forefront of developments by providing rapid and sensitive identification of minor differences among related, alternatively synthesized, and posttranslationally modified gene products at femtomole and soon attomole levels.

#### Bibliography

1. V. M. Ingram (1956) *Nature* **178**, 792–794.
2. H. Jörnvall (1970) *Eur. J. Biochem.* **16**, 41–49.
3. W. J. Gullick (1986) In *Practical Protein Chemistry: A Handbook* (A. Darbre, ed.), Wiley, Chichester, U.K., pp. 207–225.
4. R. L. Moritz et al. (1994) *Methods* **6**, 213–226.
5. T. Bergman, B. Agerberth, and H. Jörnvall (1991) *FEBS Lett.* **283**, 100–103.
6. A. Görg, W. Postel, and S. Günther (1988) *Electrophoresis* **9**, 531–546.
7. J. E. Celis et al. (1996) *FEBS Lett.* **398**, 129–134.
8. S. D. Patterson and R. Aebersold (1995) *Electrophoresis* **16**, 1791–1814.
9. D. F. Hunt et al. (1992) *Science* **255**, 1261–1263.
10. M. Wilm et al. (1996) *Nature* **379**, 466–469.
11. N. B. Afeyan et al. (1990) *J. Chromatogr.* **519**, 1–29.
12. M. L. Kochersperger, K.-L. Hsi, and P.-M. Yuan (1994) *Protein Sci.* **3**, Suppl. 1, 98, 265-M.
13. K.-L. Hsi et al. (1995) *Protein Sci.* **4**, Suppl. 2, 150, 540-M.
14. J. Rosenfeld et al. (1992) *Anal. Biochem.* **203**, 173–179.
15. U. Hellman et al. (1995) *Anal. Biochem.* **224**, 451–455.
16. J. Fernandez et al. (1992) *Anal. Biochem.* **201**, 255–264.
17. J. E. Battersby et al. (1995) *Anal. Chem.* **67**, 447–455.
18. J.-Y. Chang et al. (1982) *Eur. J. Biochem.* **127**, 625–629-M.

19. J. P. Landers (1993) *Trends Biochem. Sci.* **18**, 409–414.
20. T. Bergman (1993) In *Methods in Protein Sequence Analysis* (K. Imahori and F. Sakiyama, eds.), Plenum, New York, pp. 21–28.
21. J. Bullock (1993) *J. Chromatogr.* **633**, 235–244.
22. R. S. Rush et al. (1993) *Anal. Chem.* **65**, 1834–1842.
23. A.-C. Bergman and T. Bergman (1996) *FEBS Lett.* **397**, 45–49.
24. J. F. Banks, Jr. and T. Dresch (1996) *Anal. Chem.* **68**, 1480–1485.
25. J. F. Kelly, L. Ramaley, and P. Thibault (1997) *Anal. Chem.* **69**, 51–60.
26. H. R. Morris et al. (1996) *Rapid Commun. Mass Spectrom.* **10**, 889–896.
27. W. F. Patton et al. (1990) *Biotechniques* **8**, 518–527.
28. M. Karas and F. Hillenkamp (1988) *Anal. Chem.* **60**, 2299–2301.
29. J. B. Fenn et al. (1989) *Science* **246**, 64–71.
30. M. S. Wilm and M. Mann (1994) *Int. J. Mass Spectrom. Ion Process.* **136**, 167–180.
31. M. Wilm and M. Mann (1996) *Anal. Chem.* **68**, 1–8.
32. W. J. Henzel et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5011–5015.
33. D. J. C. Pappin, P. Hojrup, and A. J. Bleasby (1993) *Curr. Biol.* **3**, 327–332.
34. M. Mann, P. Hojrup, and P. Roepstorff (1993) *Biological Mass Spectrom.* **22**, 338–345.
35. A. Shevchenko et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14440–14445.
36. D. F. Hunt et al. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 6233–6237.
37. K. Biemann and H. A. Scoble (1987) *Science* **237**, 992–998.

### Suggestions for Further Reading

38. *Methods in Protein Structure Analysis*. Proceedings from the conference series with the same name, usually published biannually. Latest issue now in press (*J. Prot. Chem.*), from the XIth conference at Annecy, 1996.
39. *Techniques in Protein Chemistry. Proceedings from the annual meetings of the Protein Society, usually published yearly by Academic Press.*
40. *ABRF News. Proceedings from the workshops of the Association of Biomolecular Resource Facilities.*

### Peptide Nucleic Acids

Peptide nucleic acids (PNA) are DNA mimics in which the phosphoribose backbone of DNA has been substituted by a pseudopeptide backbone to which the normal nucleobases are tethered. An 8-mer thymine oligoamide was first described in 1991. The name PNA was coined when it was demonstrated that this molecule forms a rigid **triple helix** with a polyadenine DNA strand and that a PNA of mixed sequence that contains all four DNA bases forms duplexes with DNA and RNA that obey the Watson–Crick [base-pair](#) rules. In 1994, it was reported that two PNAs that have complementary base sequences hybridize into a double-helical structure. PNA has attracted great attention within medicinal chemistry and molecular biology, and also in fields, such as organic and physical chemistry because of its interesting chemical and physical properties. Because it binds efficiently and sequence-specifically to single-stranded RNA and DNA and, under certain conditions, also to double-stranded DNA, PNA has great potential for diagnostic and pharmaceutical

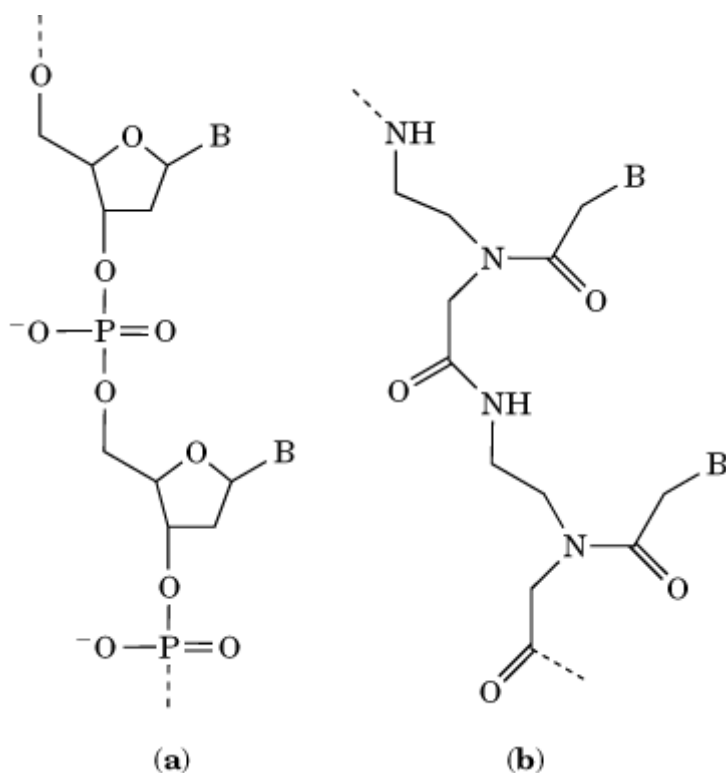


applications.

## 1. Chemistry and Synthesis

The backbone of PNA is based on 2-aminoethyl glycine linkages that mimic the normal phosphodiester backbone of DNA (1). A methylene carbonyl linker connects the standard nucleotide bases to this backbone at the amino nitrogens (Fig. 1). This chemistry has some interesting consequences. In contrast to DNA or phosphorothioate oligonucleotides, PNA is a nonionic molecule. In addition, it is achiral, which avoids diastereomeric structures and the need to develop stereoselective synthesis reactions. Schemes for protecting reactive groups on the PNA monomers are based on either Boc or Fmoc chemistry (see [Peptide Synthesis](#)). In both cases, oligomers can be assembled using established, **solid-phase** peptide synthesis procedures and a (methylbenzhydryl) amine polystyrene resin as solid support (2). PNA oligomers are cleaved from the support by treatment with either anhydrous hydrogen fluoride or trifluoromethanesulfonic acid, as in traditional peptide chemistry, followed by purification using **HPLC**. Modifications of PNA are obtained by coupling amino acids during solid-phase synthesis or by attaching other compounds that contain a carboxylic acid group to the exposed N-terminal amino group. The ease with which PNAs are modified should facilitate optimization of specific functional requirements.

**Figure 1.** Comparison of the chemical structures of (a) DNA and (b) PNA. B = nucleobase.



## 2. Physical-Chemical Properties

PNA has proven very potent as a DNA mimic capable of participating in both duplex and triplex formation. Moreover, it displays various interesting properties beyond those of natural DNA.

### 2.1. PNA-DNA/RNA Duplexes

The properties of PNA-DNA duplexes have been extensively investigated using the mixed-sequence

pentadecamer H-TGTACGTCACA ACTA-NH<sub>2</sub> (3). This PNA oligomer forms a duplex with complementary DNA that is antiparallel, has the amino-terminus of PNA facing the 3' end of DNA, and has a  $T_m$  (melting temperature) of 70°C. The corresponding DNA-DNA complex has a  $T_m$  of only 53°C. PNA also binds to the parallel, complementary, DNA target, although it has lower affinity ( $T_m = 56^\circ\text{C}$ ). The preference for antiparallel orientation is a general feature of mixed-sequence PNA-DNA duplexes. Kinetic binding studies have shown that antiparallel PNA-DNA duplex formation is very fast, faster than formation of a corresponding DNA-DNA duplex, whereas considerably slower kinetics were observed for formation of parallel PNA-DNA complexes (4). The thermal stability of PNA-RNA duplexes is generally higher than that of PNA-DNA duplexes, and the rate of association between PNA and RNA is somewhat faster than that between PNA and DNA.

The greater thermal stability of PNA-DNA/RNA duplexes compared to the corresponding DNA-DNA duplexes is predominantly ascribed to the lack of electrostatic repulsions between the two strands. This hypothesis is based on the observation that PNA-DNA and DNA-DNA duplexes have equal thermal stability (melting temperatures) at ionic strengths above 1 M Na<sup>+</sup>. As expected, the stability of PNA-DNA hybrids is little affected by changes in ionic strength, except in the limit of low ionic strength, where the stability increases. The higher stability of PNA-DNA compared with DNA-DNA duplexes at low ionic strength is due to a more favorable **entropy** change linked to duplex formation, which is attributed to a release of counterions from the DNA polyanion as PNA binds (5).

The sequence-specificity of PNA binding to nucleic acids is manifested in the considerably lower stabilities of heteroduplexes that contain mismatches. The effect of base-pair mismatches in the middle of a pentadecamer PNA-DNA duplex revealed decreases in  $T_m$  of 8 to 20°C for single mismatches (3). Similar decreases in  $T_m$  were observed in PNA-RNA for the corresponding mismatches. Thus, the sequence discrimination of a PNA oligomer is equal to, or even better than, that of a DNA oligomer.

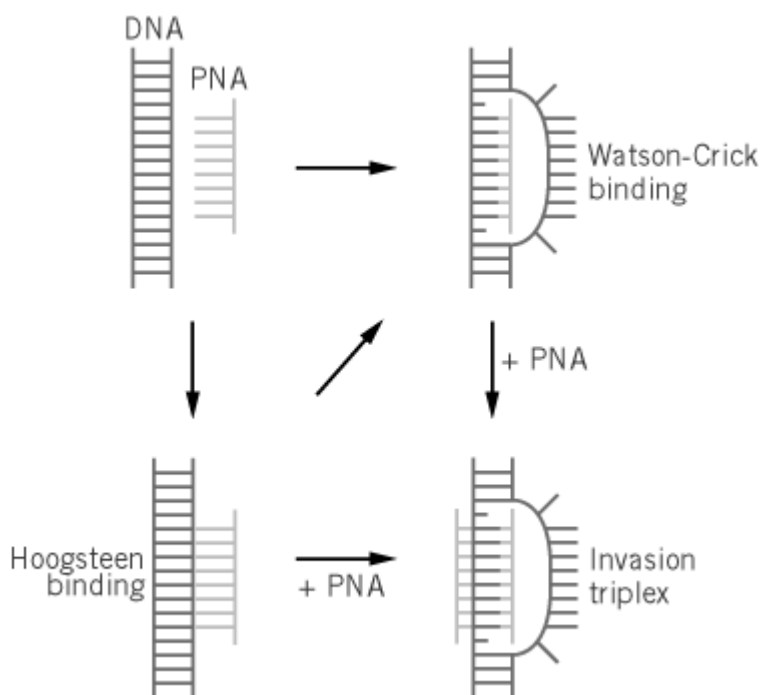
## 2.2. (PNA)<sub>2</sub>-DNA Triplexes

Homopyrimidine PNA or PNA oligomers that have a high pyrimidine content bind complementary purine DNA by forming (PNA)<sub>2</sub>-DNA triplexes (6), virtually without forming of any duplex. Monophasic, well-defined melting curves are obtained for these very stable complexes ( $T_m > 70^\circ\text{C}$  for decamer PNAs), where the binding is most probably governed by standard Watson-Crick and Hoogsteen base pairing. Triplex formation involving cytosine in the homopyrimidine PNA is pH-dependent (lower pH increases hybrid stability), as expected if C<sup>+</sup>GC triplets with one cytosine protonated are being formed. It is interesting to note that experiments still indicate a (PNA)<sub>2</sub>-DNA stoichiometry even at pH 9, where cytosine is not expected to be protonated. This suggests that only a single hydrogen bond of cytosine is sufficient or that its  $pK_a$  is drastically changed.

## 2.3. Binding to Double-Stranded DNA: Strand Invasion

The high stability of (PNA)<sub>2</sub>-DNA triplexes enables the PNA to dissociate double-stranded DNA (dsDNA). For example, homopyrimidine PNA oligomers displace the pyrimidine strand of dsDNA at the target site and form a (PNA)<sub>2</sub>-DNA triplex with the DNA homopurine target. The expelled DNA pyrimidine strand forms a single-stranded loop (Fig. 2). PNA-DNA complexes formed by strand invasion are characterized by high thermodynamic stability and high specificity, which seems contradictory. The strand invasion reaction, however, is kinetically controlled, which allows screening by the PNA for the target sequence, without becoming trapped in thermodynamically stable but incorrectly matched complexes. Although restricted to low ionic strength, this is an interesting binding property of PNA that is being pursued extensively, and several applications have been demonstrated.

**Figure 2.** Schematic general model of the possible mechanisms of PNA<sub>2</sub>-DNA triplex formation by strand invasion of a DNA duplex.



The mechanism of formation of a strand-displacement complex involving a PNA(pyrimidine)<sub>2</sub>-DNA (purine) triplex has been investigated by [gel retardation assay](#) using DNA [restriction fragments](#) that contain single PNA targets (7). The reaction obeys pseudo-first-order kinetics, and a single mismatch in a 10mer PNA reduces the binding rate constant by more than a factor of 100. For the polynucleotide system poly(dA):poly(dT) and PNA-thymine oligomers, the PNA binding reaction was studied in real time by **circular dichroic** (CD) spectroscopy (8). This revealed somewhat complex reaction kinetics and also demonstrated the presence of at least one intermediate species. The rate-limiting step is a transient opening of a few base pairs in the DNA duplex (as reflected by the size of the activation barrier). Quite remarkably, the reaction is second-order in PNA concentration, indicating that both PNA strands are involved in the precursor complex. Possibly, the two PNAs are associated loosely with the DNA duplex before the base-pair opening step. Interestingly, intercalators bound to the starting duplex DNA increase the PNA binding rate, even at higher salt concentration, presumably by facilitating initiation of PNA invasion at the intercalation site. Even if the formation of strand-displacement complexes is highly dependent on the ionic strength, complexes formed at low ionic strength are generally stable kinetically at salt concentrations up to 0.5 M NaCl.

PNA targeted to dsDNA sequences other than purine-rich stretches was recently investigated, but such complexes have reduced stability (9). With purine-rich PNA sequences, duplexes formed by the invasion into dsDNA at a complementary pyrimidine target (making a PNA-DNA duplex) have been detected by CD spectroscopy, and by a gel electrophoresis **mobility-shift assay**. Somewhat surprisingly, cytosine-rich and alternating thymine-cytosine PNA oligomers bind to their corresponding duplex DNA sequence of intact dsDNA by forming PNA-(DNA)<sub>2</sub> triplexes. The pH-dependence indicates that this probably involves Hoogsteen base pairing. These recent findings suggest that there are vast possibilities for using PNA to target a variety of base sequences of dsDNA, although more studies are necessary to assess the details of the various binding mechanisms.

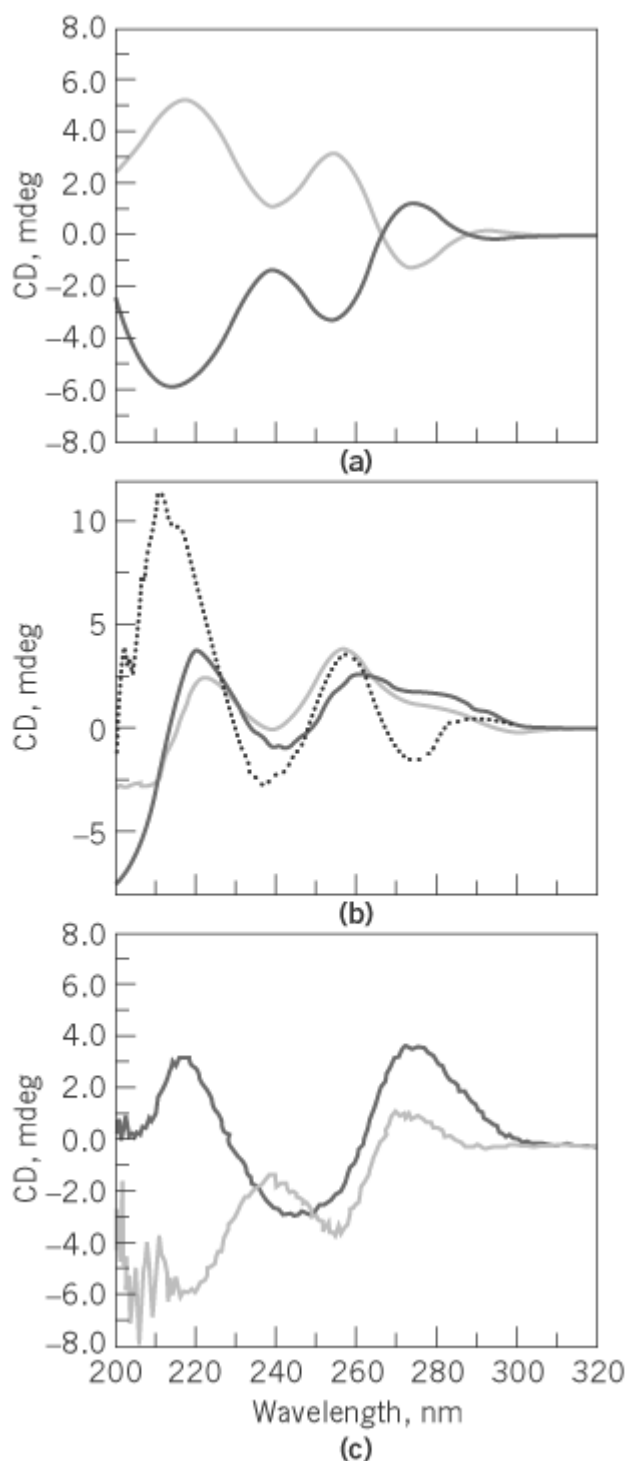
#### 2.4. PNA-PNA Duplexes and PNA-PNA-PNA Triplexes

Complementary PNA decamers that contain all four nucleobases form stable antiparallel duplexes very high in thermal stability (PNA–PNA  $T_m = 67^\circ\text{C}$ ; PNA–DNA  $T_m = 51^\circ\text{C}$ , DNA–DNA  $T_m = 33.5^\circ\text{C}$  for identical decamer sequences) that have a high sequence specificity that demonstrates Watson–Crick base pairing (10). As observed for PNA–DNA duplexes, the antiparallel PNA–PNA orientation is preferred, although the parallel PNA–PNA complex is still considerably stable ( $T_m = 47^\circ\text{C}$ ).

Moreover, a PNA<sub>3</sub> triplex that has two pyrimidine decamers and one purine decamer was detected by optical [spectroscopy](#) and [mass spectrometry](#), and the CD spectrum indicates a helical structure similar to (DNA)<sub>3</sub> and (PNA)<sub>2</sub>–DNA triplexes (11).

The base pairing in PNA–PNA duplex formation occurs rapidly on a subsecond timescale. Once formed, some PNA–PNA duplexes undergo a conformational transition (10). Although PNA in complexes with DNA or RNA adopts right-handed helical structures as a result of the hybridization, PNA–PNA duplexes, which are devoid of intrinsic chiral constituents, are believed to exist as a racemic mixture of right- and left-handed helices. However, the attachment of a chiral amino acid to the C-terminus of one strand can induce inversion into one predominating helical sense among the PNA–PNA duplexes, which was detectable as a slowly growing CD signal upon mixing. Such a chiral bias depends on the PNA sequence at the C-terminus (a G–C base pair is preferred) and on the nature of the terminal chiral amino acid. The D and L forms of amino acids give rise to mirror image CD spectra of PNA (Fig. 3a), indicating that helices of opposite handedness are formed. The nature of the side chain (**hydrophobic** or **hydrophilic**) also decides the preferred helical sense. Helix reorganization occurs on a timescale of minutes at room temperature and exhibits first-order kinetics. The reaction shows a high entropic activation barrier, as would be expected for a transition from one helical sense to another. PNA–PNA duplex helicity that originates from a chiral amino acid at the PNA terminus can propagate about 10 base pairs (chiral persistence length), or approximately one helical turn.

**Figure 3.** CD-spectra of GTAGATCACT·AGTGATCTAC duplexes. (a) PNA–PNA duplex with terminal D-lysine (light gray) and L-lysine (dark gray). (b) PNA–DNA antiparallel (light gray), PNA–RNA antiparallel (dark gray), and PNA–DNA parallel ( ). (c) PNA–PNA (light gray) and DNA–DNA (dark gray).



### 3. Structure of PNA Complexes

The unique chemical and physical properties of PNA have inspired investigations of structural details of PNA-nucleic acid hybrids by various techniques, including [NMR](#), [X-ray crystallography](#), CD, and linear dichroic (LD) spectroscopy.

#### 3.1. Dichroic Spectroscopy of PNA-DNA Hybrids

The CD of PNA-DNA and PNA-RNA duplexes provided the first indications of the overall base-pair geometry in these complexes (3). The CD spectra of DNA-DNA, antiparallel PNA-DNA, and antiparallel PNA-RNA duplexes are all similar, suggesting that PNA participates in forming right-

handed helices that have base-pair geometry not much different from that found in [B-DNA](#) or [A-DNA](#) helices (Fig. [3b](#)). In contrast, the spectra of parallel PNA-DNA and PNA-RNA complexes deviate more from the DNA-DNA spectrum, suggesting different base stacking (Fig. [3b](#)).

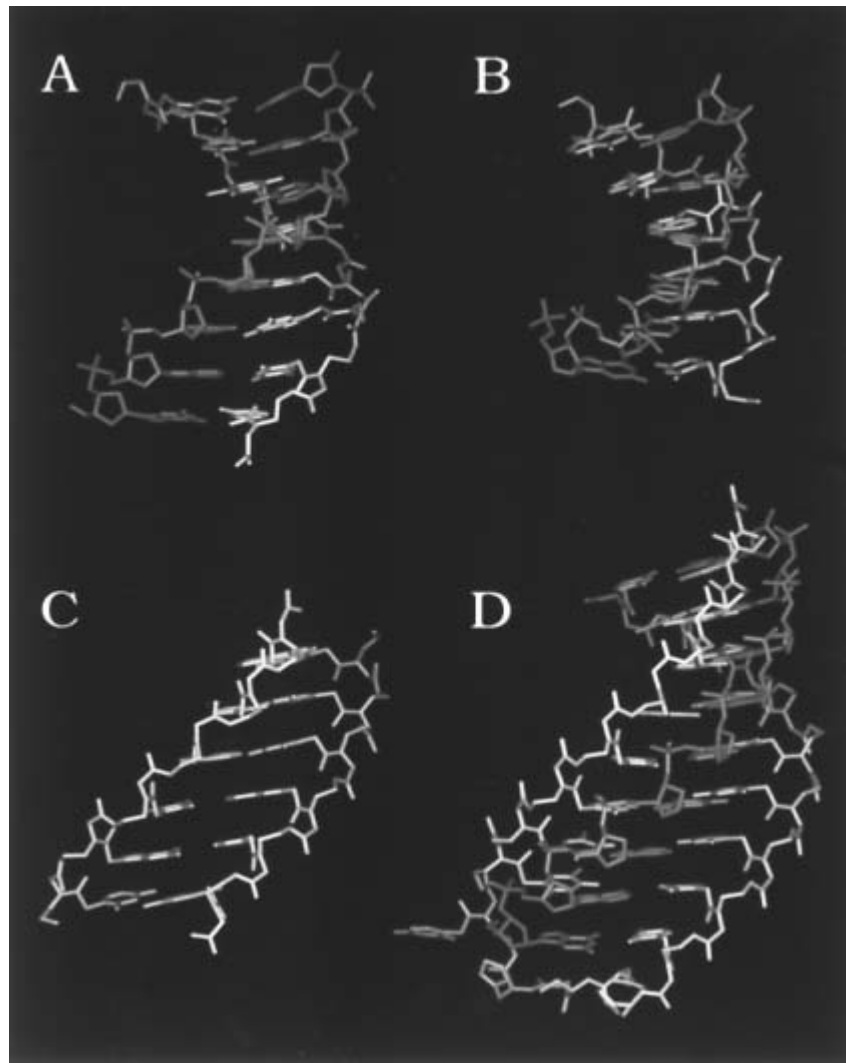
Earlier, LD and CD showed that PNA- $T_8$  with poly(dA) forms a right-handed PNA<sub>2</sub>-DNA triplex, in which the base triplets are arranged with their planes approximately perpendicular to the helix axis, similar to the standard poly(dT)<sub>2</sub>-poly(dA) triplex ([6](#)). Further, the PNA<sub>2</sub>-DNA triplex has a remarkably rigid structure, indicating efficient stacking interactions. The high thermal stability of (PNA)<sub>2</sub>-DNA triplexes is expected because of the absence of interstrand repulsions, and it is assisted by hydrophobic and dipolar interactions, and possibly also by intrabackbone [hydrogen bonds](#) (see later). Titration experiments with PNA- $T_8$  added to poly(dA) show by LD and CD that no duplex stretches are present at subsaturation ratios but that segments of (PNA)<sub>2</sub>-DNA triplexes are distributed evenly on the poly(dA) strand separated by flexible single-stranded regions. Only close to the saturation point, 2:1 PNA to DNA bases, does the complex become rigid due to complete formation of triplex and stacking between each PNA triplex segment.

A helical structure of PNA-PNA complexes has also been suggested. The CD spectrum of PNA-PNA resembles that of the corresponding DNA-DNA duplex (Fig. [3c](#)) in shape and amplitude, so it is inferred that the structures similar in their base-stacking interactions and degree of organization ([10](#)).

### 3.2. NMR and X-Ray Crystal Structures of PNA-Containing Complexes

Three-dimensional molecular structures have recently been reported for several PNA complexes (Fig. [4](#)). The structures of PNA-RNA hexamer ([12](#)) and PNA-DNA octamer ([13](#)) duplexes have been determined by NMR. The structures of a nonamer (PNA)<sub>2</sub>-DNA ([14](#)) triplex and a hexamer PNA-PNA ([15](#)) duplex were determined by X-ray crystallography. Some general conclusions can be drawn from these structures. It is clear from the heteroduplexes, that binding of PNA still allows the backbone of the nucleic acid partner to retain a nearly “normal” conformation. The RNA strand in the RNA-PNA duplex, has approximately an A-conformation with C3-*endo* sugar pucker, whereas the DNA in the PNA-DNA duplex is closer to B-form with C2-*endo* sugar pucker (see [DNA Structure](#)). It is striking, however, that the base-pair positions in both duplexes are more A-like in their helical displacement, but, in contrast to a canonical A-form helix, are more perpendicular to the helical axis. The two duplexes also show a larger pitch—about 13 base-pairs for the PNA-DNA duplex—than the usual 10 to 11 base pairs observed in nucleic acid duplexes.

**Figure 4.** Structures of PNA duplexes and triplexes. (A) and (B). Solution structures of PNA-DNA octamer duplex (A) and PNA-RNA hexamer duplex (B). (C) and (D). X-ray crystallographic structures of a PNA-PNA duplex (C) and of PNA<sub>2</sub>-DNA nonamer triplex (D). Figure kindly provided by Dr. Scott Carter.



The crystal structure of the  $(\text{PNA})_2\text{-DNA}$  triplex reveals hydrogen bonds between phosphate groups of the DNA backbone and amide protons in the PNA backbone of the Hoogsteen strand. These contacts could be important for the ability of homopyrimidine PNA oligomers to form very stable triplexes with complementary oligonucleotides. Moreover, the two PNA strands in the triplex structure are nearly identical and resemble the PNA strand in the PNA-DNA duplex. Both the  $(\text{PNA})_2\text{-DNA}$  triplex and especially the PNA-PNA duplex show very wide (26 Å and 28 Å, respectively) helices with a large pitch (16 and 18 bases per turn for the triplex and duplex, respectively), strongly suggesting that this latter type of helix is the natural conformation of helices that have PNA backbones. It is further notable that the base pairs in the PNA-PNA duplex exhibits interbase-pair stacking overlaps close to A-form. Thus, it is clear that PNA helices have a structure of their own that resembles, but is also distinctly different from, known nucleic acid helices.

#### 4. Structure–Activity Studies

Because PNA has demonstrated unique hybridization properties toward complementary oligonucleotides, it has been of the utmost interest to understand what structural elements are responsible for its behavior. Furthermore, it is important to know what kind of chemical modifications can be accomplished and accommodated while retaining binding to DNA and RNA.

A large number of studies on backbone-modified PNAs show that only minor deviations from the

original aminoethylglycine backbone can be tolerated. Extensions in any of the three possible linking directions impair stability of its complexes, and the restricted flexibility imposed by the secondary amido group in the nucleobase linker is also necessary because reducing this to a tertiary amine is deleterious to the hybridization potency (16). Isomerization to the retroinverse structure, which essentially moves a methylene group from the ethyl to the glycine moiety, also results in PNAs of low hybridization efficiency. However, there is much freedom in the way substituents may be placed in the  $\alpha$ -position of the glycine moiety of the PNA backbone. Backbones based on natural amino acids other than glycine result in fair to good DNA mimicking properties. Even cyclic substituents, such as cyclohexyl substitution at the aminoethyl linker, are acceptable, provided the right stereoisomer is chosen. By using PNA monomers that have incorporated natural amino acids, chemical functionality can be introduced into the backbone and the physical properties of the PNAs, such as hydrophilicity, hydrophobicity, ionic character, etc. might be finely tuned and controlled. It is believed that this greatly facilitates the optimization of pharmacokinetic behavior, for example (see later).

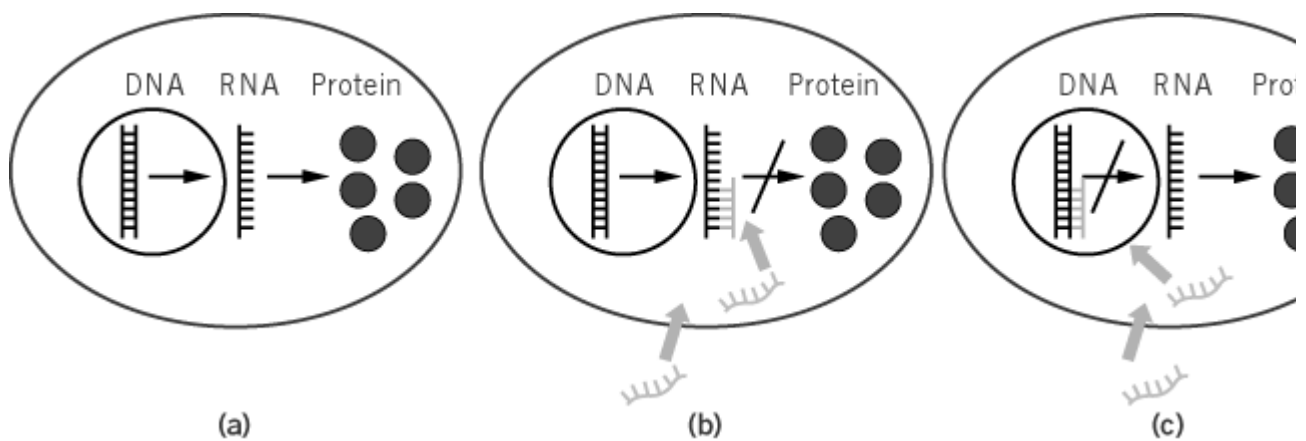
PNA-DNA chimeras are potentially interesting because they may allow DNA-recognizing enzymes (such as [ribonuclease H](#) and **DNA polymerase**) to interact with the DNA part, while the PNA part would provide high binding affinity and specificity. For this purpose, a linker is necessary that connects the two oligomers to allow stacking between the neighboring bases in the PNA and DNA parts of the chimera. Successful reports describe the synthesis of PNA-DNA chimeras of both DNA-3'-N-PNA and DNA-5'-C-PNA type junctions (17). Such chimeras hybridize to complementary oligonucleotides less efficiently than pure PNAs, but they are able to activate ribonuclease H.

## 5. A Possible Genetic Drug?

The potential use of PNA as an **antisense** or antigene drug for sequence-specific modulation of gene expression has bright prospects. Oligonucleotide analogs can be designed to recognize and hybridize to complementary sequences in a particular [messenger RNA](#) and thereby inhibit its translation (see also [Antisense Oligonucleotides](#)). Alternatively, oligonucleotides can be targeted to dsDNA and interfere there with [transcription](#) of a particular gene (the antigene strategy). These two approaches are depicted in Figure 5. *In vitro* assays examining the effect of PNA on [DNA replication](#), transcription, and translation all look quite promising. The sequence specificity is high, and the biological stability of PNA is sufficient (18). Drawbacks include the poor cellular uptake of PNA and the sensitivity of strand invasion complexes to salt concentrations, which must be addressed before general applications can be developed. Moreover, the pharmacological properties of PNAs have not yet been thoroughly investigated.

**Figure 5.** Principles of the antisense and antigene gene-targeting strategies. In normal cells (a), DNA is transcribed into mRNA which then is translated into protein. When the cells are treated with an antisense oligomer (b), hybridization to a sequence specific mRNA inhibits expression of the protein at the level of translation. Treatment of the cells with an antigene oligo complementary to a sequence in the DNA leads to inhibition of transcription of the gene into mRNA.

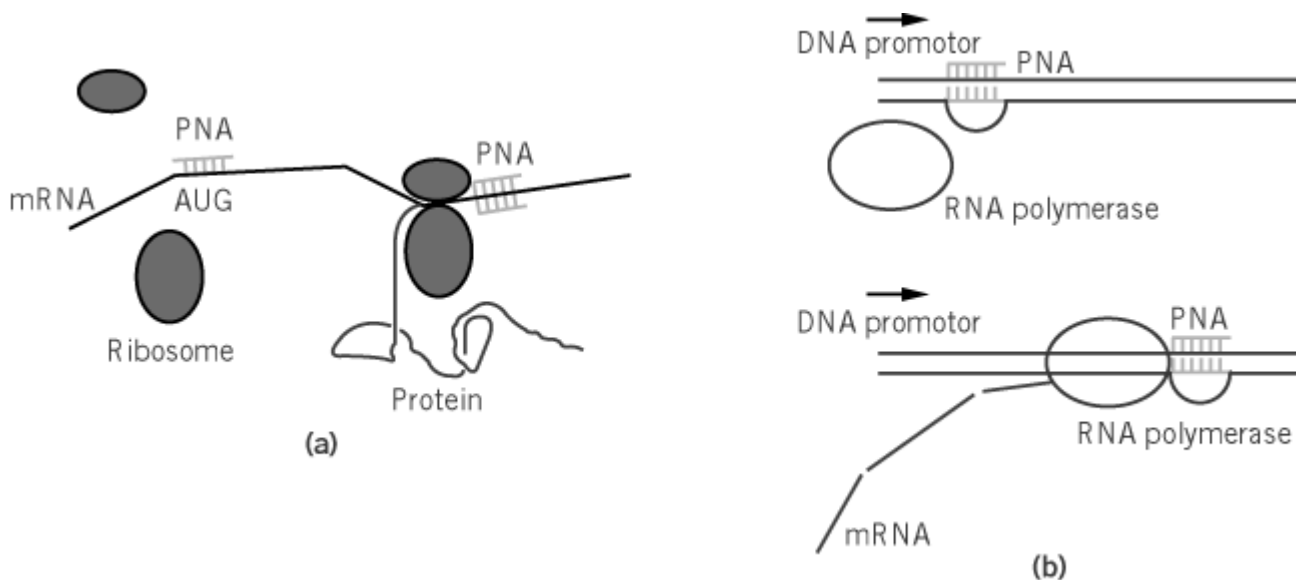




### 5.1. Transcription Arrest by PNA

When targeting genomic DNA by the antigene strategy, inhibition should in theory be obtained by the formation of just one PNA-DNA complex in each haploid [genome](#). Because pyrimidine PNAs form extremely stable, strand-invasion complexes upon binding to dsDNA, they should be excellent antigene agents (18). It has been shown that such complexes indeed arrest transcription by both prokaryotic and eukaryotic **RNA polymerases** *in vitro*. PNA bound to the **promoter** clearly blocks the access of the polymerase protein to the DNA, whereas PNA bound within the coding region of a gene sterically blocks polymerase progression (Fig. 6a). In therapeutic applications, several issues are, however, still unclear. First, to bind dsDNA, a poly-purine target must be present in the gene of interest. This, of course, limits the number of possible targets. Second, and this could be more serious, the formation of strand-displacement complexes occurs very slowly at physiological salt concentration. Modifications of PNA by intercalators or positively charged lysine residues may, however, improve the binding rates. Furthermore, chemically linking the end of the Watson-Crick and Hoogsteen PNA strands to each other or incorporating pH-independent pseudoisocytosines into the Hoogsteen strand can increase the association rates with dsDNA. It was found that binding of PNA to **supercoiled** plasmid DNA is much faster than to linear DNA, indicating that the [DNA topology](#) is an important parameter for the rate of binding. Because transcriptionally active chromosomal DNA is usually negatively supercoiled, this could be of relevance for PNA binding *in vivo*. It has also been shown that binding of PNA to dsDNA is enhanced when the DNA is being transcribed, probably because of partly ssDNA in the transcription bubble.

**Figure 6.** Translational arrest and transcriptional inhibition. (a) Translation initiation is blocked by either a duplex- or a triplex-forming PNA, whereas a triplex-forming PNA is required for translation elongation arrest of the ribosome. (b) Transcriptional inhibition is accomplished by homopyrimidine PNAs that form invasion triplexes either at the promoter or the coding region of the gene.



## 5.2. PNA as an Artificial Promoter

Even though not considered an antigene effect, the ability of PNA to activate transcription is worth mentioning in this context. The resemblance between PNA strand-invasion hybrids and transcriptional initiation complexes allows transcription to start on the looped-out strand of the PNA-DNA complex. Experiments with *Escherichia coli* RNA polymerase show that the looped-out DNA strand acts as a template, and the length of the resulting transcripts corresponds exactly to run-off transcripts initiated at the bound PNA. Such a “PNA promoter” is very strong and might have the potential to activate genes *in vivo*, provided that strand invasion is accessible inside cells. As mentioned, PNA also arrests transcription, but the two processes differ as to the size of the displacement loop. Although arrest works for decamer targets, transcriptional initiation requires a larger loop.

## 5.3. Translational Arrest by PNA

Because PNA-RNA complexes are not recognized by ribonuclease H, PNA-mediated antisense effects depend on mechanisms involving steric blocking of RNA processing, transport into the cytoplasm, or translation. Studies performed *in vitro* show that translational initiation is blocked by a PNA that forms either a duplex (mixed sequence) or a triplex (pyrimidine-rich), whereas arrest of translational elongation by the ribosome requires a triplex-forming homopyrimidine PNA (Fig. 6b). Antisense properties in a cellular context have been possible to assess by [microinjection](#) of PNA directly into cultured cells (circumventing the cellular uptake problem), and sequence-specific effects have been confirmed (18).

Bacteria have attracted relatively little attention as targets for antisense agents. However, the demonstration of gene- and sequence-specific antisense inhibition by PNA in *E. coli* opens possibilities for antisense antibacterial drugs and gene-function analysis in bacteria (19). In view of the resistance to current antibiotics displayed by an increasing number of strains of pathogenic bacteria, the possibility of novel antibiotics seems especially important.

## 5.4. Delivery and Biological Stability of PNA

The prospect of using PNA as a genetic drug makes sense only if PNA is delivered efficiently to cells and distributed so that it interacts with the target mRNA or DNA. Studies conducted with cells in culture indicate that PNA uptake is slow and that PNA ends up in **endosomal** compartments (18). Using a liposome system as an *in vitro* cell model, the PNA penetration properties based on passive [diffusion](#) have been investigated in detail and are slow (requiring days) (20). Therefore, PNA molecules need to be modified somehow, either by conjugation to ligands so that active receptor-

mediated cellular uptake **channels** might be exploited or by attachment to moieties that make them attracted to the membrane, so that the local concentration of PNA near or within the cell membranes increases. It has been found that a PNA conjugated to an [antibody](#) for [transferrin](#) transports across the blood-brain barrier, possibly via the transferrin receptor (21). Most other DNA analogs also show poor uptake properties, similar to those observed with PNA. It is always important, however, to extend uptake studies to *in vivo* situations. Surprisingly, phosphorothioate oligonucleotides, which are taken up only inefficiently in cell culture, have favorable pharmacokinetic behavior in animal models (18).

Agreeing with its chemical remoteness to natural nucleic acids or peptides, PNA exhibits resistance to degradation by cellular **nucleases** and [proteinases](#) and remains intact when incubated in serum or cellular extracts. Thus, PNA has ample biological stability for drug development. Furthermore, no signs of general toxicity by PNA have been observed thus far.

## 6. PNA as a Diagnostic Tool

An increasing number of applications have been described that confirm the potential for PNAs to facilitate the development of novel methodologies in molecular biology, diagnostics, and related fields. A few of these applications are described here.

### 6.1. PNA-Directed PCR Clamping

The high stability and sequence selectivity of PNA-DNA complexes have been exploited in analyzing single base-pair mutations in DNA by employing **PCR** amplification (22). PNA directed toward the primer binding site, the region adjacent to the primer site, or toward the middle of the PCR region, decreases amplification because PNA cannot function as a primer for **DNA polymerase**. The most efficient and discriminative clamping was observed, as might have been expected, with homopyrimidine PNA decamers. Formation of such PNA<sub>2</sub>-DNA triplexes inhibit the polymerase at any PNA binding site. With this technique, a single mismatch in the primer site is detectable.

### 6.2. PNA and S1 Nuclease as Artificial Restriction Enzymes

PNAs in combination with [S1 nuclease](#) have been used as “artificial restriction enzymes” (23). Hybridization of a homopyrimidine PNA oligomer to the complementary target on dsDNA leads to strand displacement of the noncomplementary DNA strand. Upon treatment with nuclease S1, the single-stranded DNA is cut into well-defined fragments. The binding of two PNA molecules to two adjacent targets on the same or opposite DNA strands led to an opening of the entire region, including the intervening base pairs, and thereby provided an accessible substrate for the nuclease and considerably increased the sensitivity to cleavage.

### 6.3. PNA-Mediated Nucleic Acid Purification

PNA oligomers that carry six [histidine](#) residues can be used to purify target nucleic acids by nickel [affinity chromatography](#) (24). Moreover, PNAs linked to [biotin](#) in combination and **streptavidin**-coated magnetic beads can be used to purify genomic DNA of *Chlamydia trachomatis* directly from urine samples. This “purification by hybridization” approach offers a number of advantages over traditional purification methods. One drawback, however, is that it requires knowledge of a target sequence and a dedicated capture oligomer must be synthesized for each different target nucleic acid. Short pyrimidine PNAs form unprecedentedly stable triplexes with complementary nucleic acids. Because being short means that their target sequence is prevalent in large nucleic acids, such short PNAs can be used as generic capture probes to purify large nucleic acids. It has been shown that a biotin-tagged PNA-thymine heptamer purifies human genomic DNA efficiently from whole blood by a simple and rapid procedure.

## 7. Mimic of DNA

Studies of PNA and its complexes with nucleic acids add new perspectives to our understanding of

DNA and also RNA. The base-pair stacking in the PNA-DNA and PNA-RNA duplexes is similar to that of the pure nucleic acid duplexes, showing that neither the absence of backbone charge nor absence of chirality in one of the strands markedly affects the helical stacking. Thus, the helicity can be considered a result of the stacking interactions (mainly dipolar interactions between the nucleobases) in combination with strains from the backbone, rather than arising from phosphate–phosphate repulsion that strives to increase the separation of negative charges along the backbone. An important manifestation of the preferred helical structure is also the helical stacking in PNA-PNA duplexes (10, 11). It is speculated that the helical shape of nucleic acids was selected in evolution because it allows better exposure of signature features of the bases on the helical surface without loss in secondary structure stability and thereby facilitates protein recognition.

The structural similarity between PNA-DNA and DNA-DNA duplexes offers the opportunity to investigate the forces responsible for binding different ligands to dsDNA (25). Assuming that the main binding forces can be divided into [electrostatic interactions](#) and hydrophobic contributions, the effect of replacing the negatively charged DNA backbone with a neutral peptide backbone could help assess the role of the phosphate groups in binding a given ligand. Typical molecules that intercalate or bind to grooves exhibit strongly reduced affinities for PNA-DNA complexes. None of the intercalators studied exhibited any significant binding, whereas classical minor groove binders bind weakly to PNA-DNA complexes. No simple ligand has been found that binds to PNA-PNA duplexes. These findings illustrate the importance of electrostatic attractions for ligand binding to DNA.

## 8. PNA and Prebiotic Chemistry

Identifying the first genetic molecule could provide the key to understanding the origin of life. The difficulties in spontaneously forming nucleotides under prebiotic conditions might be used to argue that nucleic acids (DNA or RNA) were not the first genetic materials on the primitive earth. Therefore, it has been proposed that one or more genetic systems composed of simpler monomers may have preceded the [RNA world](#), an era in which RNA molecules would have acted as both the genetic material and the catalysts. Then, genetic takeover would have allowed the development of RNA-controlled life. PNA is an example of an oligomer that could have preceded RNA because it is simpler than RNA but exhibits all of the important properties of nucleic acids with respect to information storage. PNA oligomers can act as templates for polymerizing of the complementary activated RNA mononucleotides, and an RNA oligonucleotide template facilitates the synthesis of complementary PNA strands. These experiments suggest that transitions between different genetic systems may occur without loss of information. A continuous transition from one genetic system to another is possible only if it is possible to form chimeras, that is, mixed molecules that contain building blocks of both systems. In fact, template-directed formation of PNA-DNA chimeras has been demonstrated (26).

PNA is attractive as a primordial [macromolecule](#) because it is simple and stable, but, not least, also because PNA is achiral. Introducing and maintaining homochirality on the primitive earth is often considered a crucial issue. How would the first RNA and DNA molecules (which are chiral) have formed from a random racemic mixture of chemically related subunits, including all the possible optical isomers? One possibility is that some physical selection mechanism resulted in forming a chiral environment before the origin of life, or, which may be more likely, that a general chiral bias was active during evolution and led to a final complete selection of handedness (27). Because homochirality is required for the function of enzymes and nucleic acids, an adaptive process might have selected for homochirality at some stage in prebiotic chemical evolution. The discovery of PNA-PNA helices with the handedness determined by a terminal chiral seed suggests that a molecular precursor of RNA could have provided a mechanism for chiral amplification (28).

## 9. Concluding Remarks

The development of PNA has had implications for many areas of DNA chemistry, biology, and

technology, ranging from basic molecular recognition, self-assembly, chiral induction, and genetic therapeutic drugs, to our understanding of the structure and function of Nature's genetic material, DNA, and speculation on its prebiotic predecessors and origin. Even some novel materials have their origin in PNA (29). Therefore, PNA should be viewed as a DNA mimic and also as a structural and self-recognizing system in its own right.

### Bibliography

1. P. E. Nielsen (1991) *Science* **254**, 1497–1500.
2. K. L. Christensen et al. (1995) *J. Peptide Science* **3**, 175–183.
3. M. Egholm et al. (1993) *Nature* **365**, 566–568.
4. K. Jensen, H. Orum, P. E. Nielsen, and B. Norden (1996) *Biochemistry* **36**, 5072–5077.
5. S. Tomac et al. (1996) *J. Am. Chem. Soc.* **118**, 5544–5552.
6. S. K. Kim et al. (1993) *J. Am. Chem. Soc.* **115**, 6477–6481.
7. V. V. Demidov et al. (1995) *Proc. Natl. Acad. Sci.*, **92**, 2637–2641.
8. P. Wittung, P. E. Nielsen, and B. Norden (1996) *J. Am. Chem. Soc.* **118**, 7049–7054.
9. P. Wittung, P. E. Nielsen, and B. Norden (1997) *Biochemistry* **36**, 7973–7979.
10. P. Wittung et al. (1994) *Nature* **368**, 561–563.
11. P. Wittung, P. E. Nielsen, and B. Norden (1997) *J. Am. Chem. Soc.* **119**, 3189–3190.
12. S. C. Brown, S. A. Thomson, J. M. Veal, and D. G. Davis (1994) *Science* **265**, 777–780.
13. M. Leijon et al. (1994) *Biochemistry* **33**, 9820–9825.
14. L. Betts, J. A. Josey, J. M. Veal, and S. R. Jordan (1995) *Science* **270**, 1838–1841.
15. H. Rasmussen et al. (1997) *Nat. Struct. Biol.* **4**, 98–101.
16. B. Hyrup et al. (1994) *J. Am. Chem. Soc.* **116**, 7964–7970.
17. K. H. Petersen et al. (1995) *Bioorg. Med. Chem. Lett.* **5**, 1119–1124.
18. H. Knudsen and P. E. Nielsen (1997) *Anticancer Drugs* **8**, 113–118.
19. L. Good and P. E. Nielsen (1998) *Nat. Biotechnol.* **16**, 355–358.
20. P. Wittung et al. (1995) *FEBS Lett.* **365**, 27–29.
21. W. M. Partridge, R. J. Boado, and Y.-S. Kang (1995) *PNAS*, **92**, 5592–5596.
22. H. Orum et al. (1993) *Nucleic Acid Res.* **21**, 5532–5336.
23. V. Demidov et al. (1993) *Nucleic Acid Res.* **21**, 2103–2107.
24. H. Orum et al. (1995) *BioTechniques* **19**, 472–480.
25. P. Wittung et al. (1994) *Nucleic Acid Res.* **22**, 5371–5377.
26. S. L. Miller (1997) *Nat. Struct. Biol.* **4**, 167–169.
27. B. Norden (1978) *J. Mol. Evol.* **11**, 313–332.
28. A. W. Schwartz (1994) *Curr. Biol.* **4**, 758–760.
29. R. H. Berg, S. Hvilsted, and P. S. Ramanujam (1996) *Nature* **383**, 505–507.

### Suggestion for Further Reading

30. M. Eriksson and P. E. Nielsen (1996) PNA-nucleic acid complexes. Structure, stability and dynamics, *Q. Rev. Biophys.* **29**, 369–394.

## Peptide Synthesis

The chemical synthesis of small [proteins](#) has only recently become feasible with reliability. Progress in this field has depended on the development of the solid phase procedure and the automated synthesizer by Merrifield to reduce the labor and errors associated with repetitive manual procedures. Difficulties in generating pure peptides arose from accumulated byproducts resulting from synthetic difficulties in chain assembly and deprotection of side-chain functionality. Advances in analytical and purification technologies have provided tools for addressing these problems. A variety of strategies have been developed to overcome the synthetic difficulties, and several laboratories have invested sufficiently in the optimization of the technology to produce synthetic small proteins of sufficient purity as to remove any doubts regarding their biological effects (1).

One of the first examples where a small synthetic protein had a major impact on a biological problem of clinical significance was the development of human immunodeficiency virus (HIV) [proteinase inhibitors](#). Although the gene sequence for the 99-residue protein which dimerized to give the active proteinase, was identified when the DNA sequence of the gene for HIV was determined, [cloning](#) and expressing of HIV proteinase proved to be exceptionally difficult for pharmaceutical groups, even for those with records of considerable previous success in biotechnology. Synthetic HIV proteinase, a 99-residue homodimer prepared in the laboratory of Stephen Kent (2), was used in the initial determination of HIV proteinase specificity and inhibitor screening. The availability of this synthetic protein certainly expedited the development of HIV proteinase inhibitors used for the treatment of acquired immune deficiency syndrome (AIDS). This example and others would argue that chemical synthesis of small proteins may be more rapid than [cloning](#) and expression for some proteins, and it can provide ready access to sufficient quantities (hundreds of milligrams) of protein for screening or biophysical studies. In fact, the first X-ray [crystallography](#) structure (3) of HIV proteinase complexed with an inhibitor, MVT-101, used synthetic HIV proteinase and set the record for the largest synthetic molecule whose structure was confirmed by crystallography.

## 1. Solution Versus Solid Phase

Traditional approaches to peptide synthesis were based on solution methods using convergent synthetic schemes, with isolation and characterization of each intermediate. This provided well-defined products with confidence in their final structures, but at the price of increased labor and losses at each synthetic step in separating product from reagents and byproducts. Until recently, the synthesis of small proteins (60 to 100 residues) has been dominated by fragment condensation in solution with maximal protection of side chains. The desired protein sequence is divided into peptide segments of about 10 residues, which is about the maximum length of peptide readily prepared by stepwise addition in solution. After each reaction, the product is isolated and fully characterized before proceeding with the next reaction. Once the set of fragments has been prepared, they are combined pairwise to generate segments of approximately 20 residues. These are then combined pairwise to generate fragments of approximately 40 residues, and this fragment condensation continues until the desired sequence is obtained. The protecting groups on side chains are then removed to give the fully unprotected peptide chain, which is allowed to fold, resulting in the desired protein. First, the 80-residue protein is divided into eight 10-residue segments that are prepared by stepwise elongation requiring nine coupling and nine deprotection steps, with isolation of each intermediate to give a completely protected fragment with both amino and carboxyl termini protected. Depending on the role of the segment, either its *N*- or *C*-terminus is deprotected and reacted with its adjacent fragment. This process is continued until the entire protein is assembled, deprotected, purified, and allowed to fold.

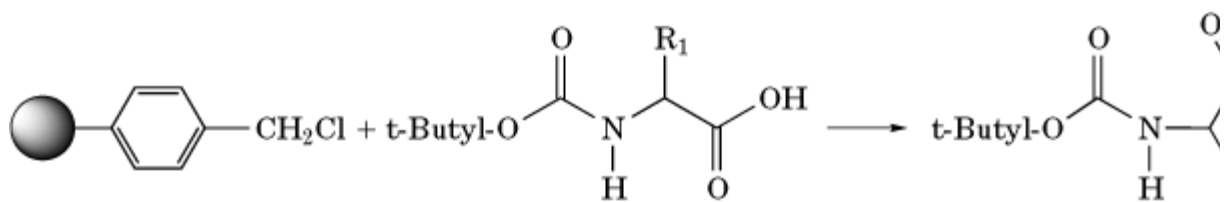
The solid-phase approach utilizes a polymeric protecting group that allows the use of excess reagents to force reactions to near completion by the mass action law and trivial isolation of polymeric product by filtration and washing. Intermediates are not isolated, and purity of the final product depends on complete reaction at each synthetic step and minimization of side reactions during the buildup of the oligomeric peptide and subsequent removal from the polymer with deprotection, to

give the desired product. The advent of high-performance liquid chromatography (HPLC) with more sophisticated techniques such as nuclear magnetic resonance ([NMR](#)), **capillary electrophoresis**, and [mass spectrometry](#) for purification and characterization of the intermediates and final products allows routine synthesis of peptides in the 50- to 100-residue range. Because of the difficulties in purification of larger peptides and small proteins with only minor differences in structure from byproducts, the unambiguous synthesis of larger peptides and small proteins is best accomplished by assembly of fragments that have been purified and fully characterized. This prevents accumulation of side products with only minor structural differences that can be difficult to remove in the final mixture. Initially, chain assembly was relatively easy to optimize, and the majority of undesirable side products in the final cleavage were due to incomplete deprotection of side chains. Considerable effort over two decades was devoted to understanding the sequence-dependent problems leading to truncated sequences or those missing a residue, but these efforts were hampered by the polymeric support itself, which limited application of the normal methods for characterizing intermediates. This effort has led to the current state of the technology, in which average reaction yields are estimated to be greater than 99.5%. Such yields are essential if multiple sequential chemical reactions are performed without isolation and purification of intermediate products.

Automation of the solid-phase reaction was initiated almost immediately once a viable synthetic scheme for peptides was evolved, and the first automated synthesizer was announced by Merrifield and Stewart in 1965 ([4](#)). Continuous development of synthesizers and the associated chemistry allows the automated addition of 75 residues per day to a growing peptide chain ([5](#)). In many cases, the repetitive yields are sufficiently high that useful products can be isolated from the synthetic mixture by HPLC when small proteins are prepared. If one desires to be more confident that the observed properties are uniquely determined by the targeted sequence, then a more conservative approach utilizing fragment condensation is still required.

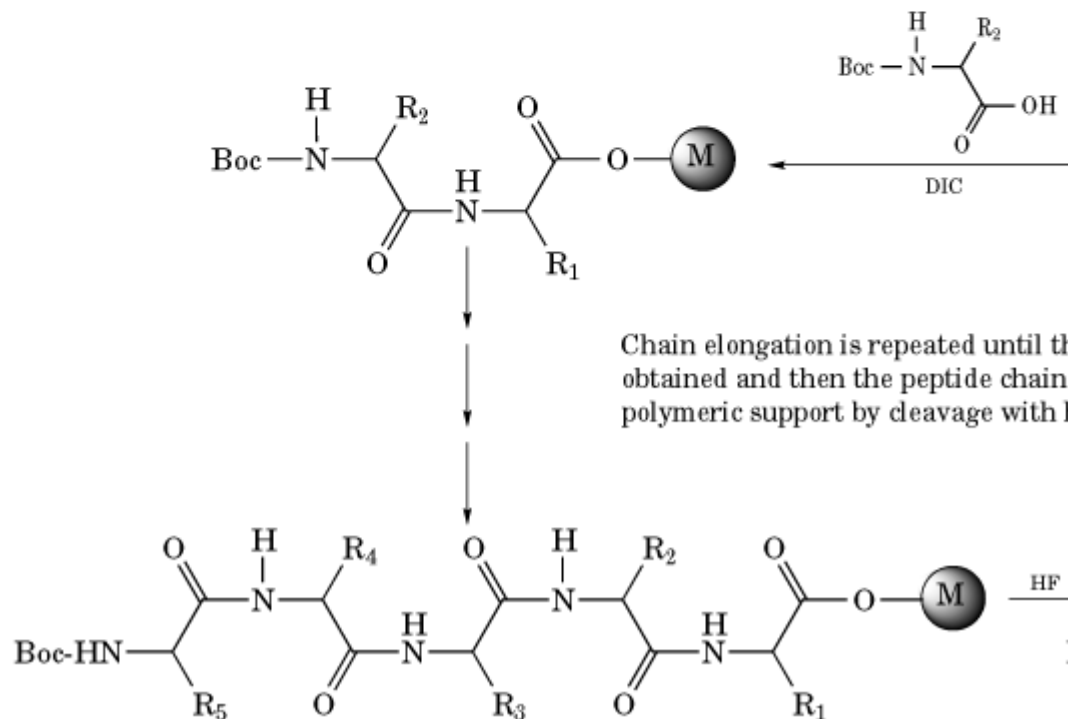
Illustrative are the solid-phase protocols for the two strategies (Boc and Fmoc) commonly used for the synthesis of peptides. The Boc strategy ([Fig. 1](#)) is often combined with a 1% to 2% cross-linked polystyrene support and a benzyl ester linkage to the polymer, requiring strong acid such as hydrogen fluoride for deprotection. The procedure favored by most synthetic laboratories uses an acid-labile linkage similar to the *p*-methoxybenzyl ester linkage of the Wang resin and a base-labile amino protecting group, the fluorenylmethyloxycarbonyl (Fmoc), on the added amino acids ([Fig. 2](#)). One can use side-chain protection with similar acid lability to the Wang linkage to give free peptide upon cleavage, or use more stable side-chain protection to give the protected peptide for fragment condensation after purification and characterization. In the latter case, a final deprotection with strong acid such as HF is required.

**Figure 1.** A common implementation of the BOC strategy for peptide synthesis.



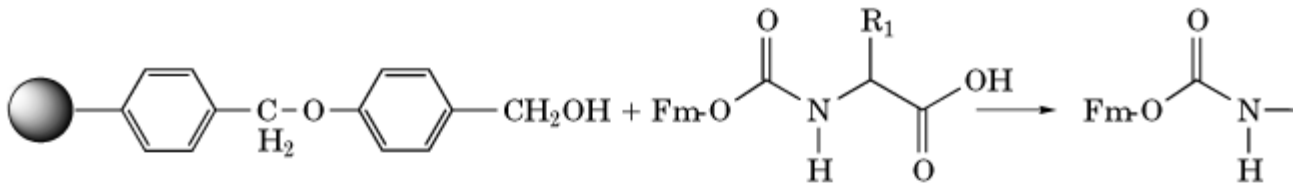
Loading of Merrifield polymer with Boc-protected C-terminal residue

The cycle of chain elongation consists of Boc removal with trifluoroacetic acid (TFA), neutralization with a base such as triethylamine (TEA), and coupling of the next Boc-amino acid (R<sub>2</sub>) in the presence of a suitable coupling reagent such as diisopropylcarbodiimide (DIC).



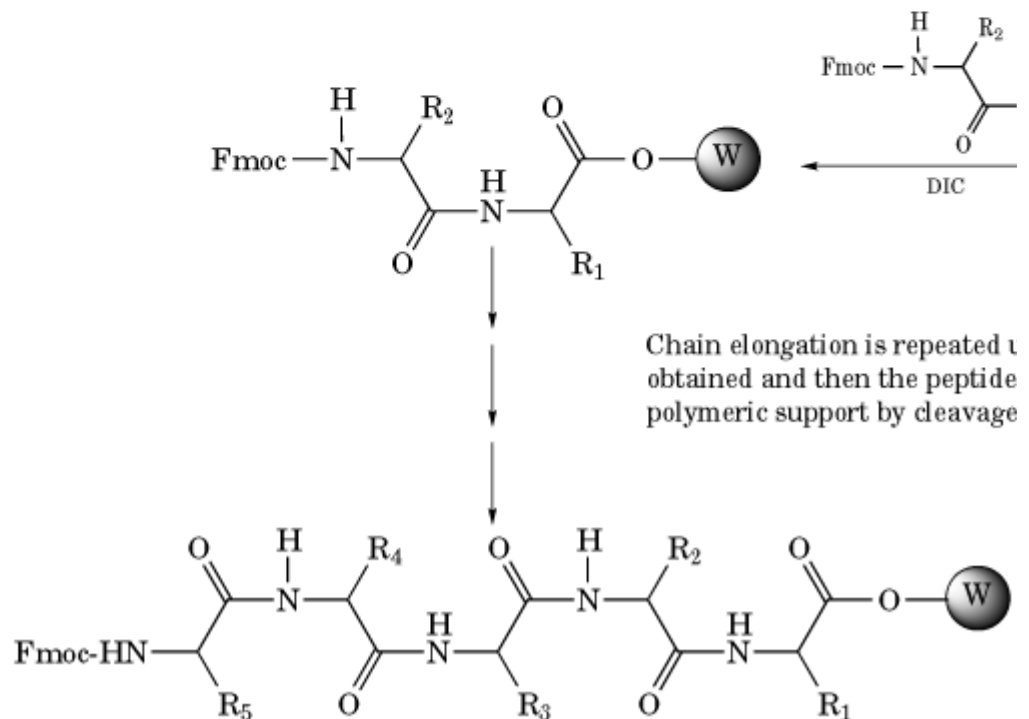
**Figure 2.** The alternative Fmoc/Wang support strategy which allows ready synthesis of protected peptide fragments as





Loading of Wang polymer with Fmoc-protected C-terminal residue

The cycle of chain elongation consists of Fmoc removal with piperidine (Pip), neutralization with a base such as triethylamine (TEA), and coupling of the next Fmoc-amino acid ( $R_2$ ) in the presence of a suitable coupling reagent such as diisopropylcarbodiimide (DIC).

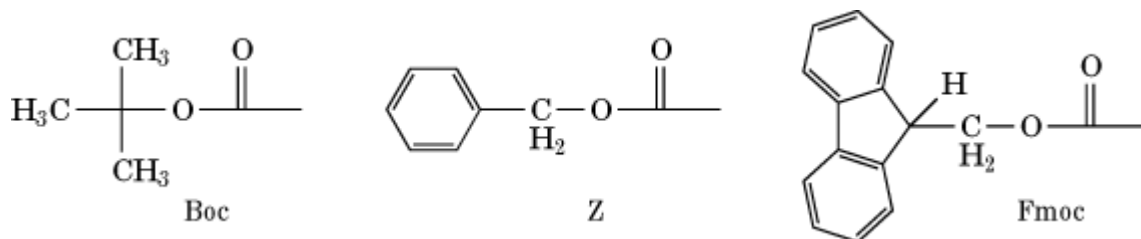


## 2. Orthogonal and Selective Protection

Because many of the 20 common amino acids have side chains containing chemically reactive groups (**carboxyl**, hydroxyl, imidazole, indole, phenol, **amine**, guanidinium) which could compete for selective **peptide bond** formation, they must be protected temporarily during assembly to make sure that only the desired  $\alpha$ -amino and  $\alpha$ -carboxyl groups react to join the protected amino acid to the growing peptide chain with the desired linkage. It is essential that these side-chain protecting groups are stable (orthogonal) to the conditions used for removal of the amino protecting group (Fig. 3) in chain elongation and yet be readily removable at the end of the synthesis. Much effort has been expended developing compatible sets of protecting groups for the 20 naturally occurring amino acids. The peptide chain is routinely assembled stepwise from the *C*- to *N*-terminus because one can

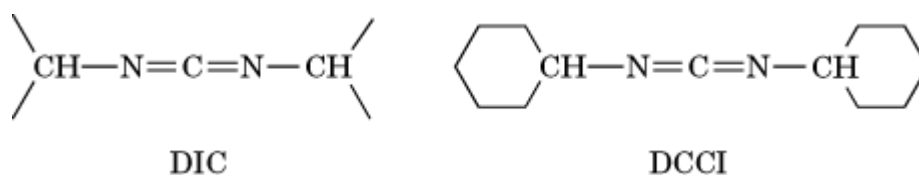
best preserve the chiral integrity of the component amino acids by activating the carboxyl group of a urethane-protected amino acid [ROCO-, where R = *t*-butyl (Boc), R = benzyl (Z), and R = fluorenylmethyl (Fmoc)] for coupling with the *N*-terminal amine of the growing peptide chain. After the addition of a protected amino acid, the amino-protecting group must be removed in order to provide the amine nucleophile for the next reaction. This must be accomplished while retaining the side-chain protecting groups of the growing peptide chain. This leads to the concept of orthogonal protection, in which the  $\alpha$ -amino group is protected with a group that can be removed selectively by a mechanism that does not affect the side-chain protecting groups. A common example in solid-phase synthesis would be the use of the Fmoc group, which can be removed repetitively from the  $\alpha$ -amino groups under basic conditions, and the use of acid-labile protecting groups, such as Boc or OtBu, on the side chains that must remain throughout chain elongation. If one attempts to couple a peptide to another peptide by activation of the carboxyl group of one of the two, the chiral integrity of the amino acid whose carboxyl group is activated (which is analogous to coupling an acetyl amino acid, which can form a cyclic intermediate, leading to rapid exchange and inversion of the alpha carbon) is often compromised by **racemization**, consequently, a mixture of peptides with *R* and *S* configurations at the activated amino acid results. This leads to the use of special procedures for coupling peptides, as well as a strong proclivity to favor glycine or proline as the *C*-terminal residues in fragment coupling, due to their resistance to racemization. Chemical or enzymatic ligation of peptide fragments offers an alternative approach to simple chemical activation of the *C*-terminal carboxyl group of one fragment and nucleophilic attack by the *N*-terminal amino group of the other fragment.

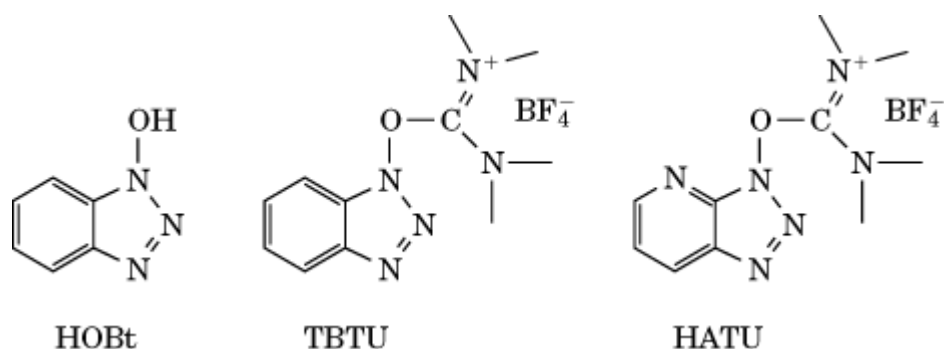
**Figure 3.** Urethane protecting groups for amino groups in common usage; Boc, *t*-butyloxycarbonyl; Z, carbobenzoxy (or benzyloxycarbonyl); Fmoc, fluorenylmethyloxycarbonyl.



### 3. Coupling Reagents

The choice of one of the numerous reagents to activate the carboxyl group for nucleophilic attack by the amine function, to produce the peptide bond, depends on numerous factors. When coupling peptide fragments, reagents that minimize racemization are preferred; a traditional favorite in solution chemistry for coupling fragments has been the azide group. A common procedure is activation by a combination of a diimide, such as 1,3-diisopropylcarbodiimide (DIC) or 1,3-dicyclohexylcarbodiimide (DCCI), with hydroxybenzotriazole (HOBt), or through the use of an activated derivative of HOBt, such as TBTU, to generate the activated HOBt ester *in situ*. Incorporation of multiple sequential sterically hindered amino acids (**valine, isoleucine, aminoisobutyric acid, etc.**) often require the use of special coupling reagents, such as acid fluorides, HATU, etc.) to provide reasonable reaction rates and efficient coupling yields.





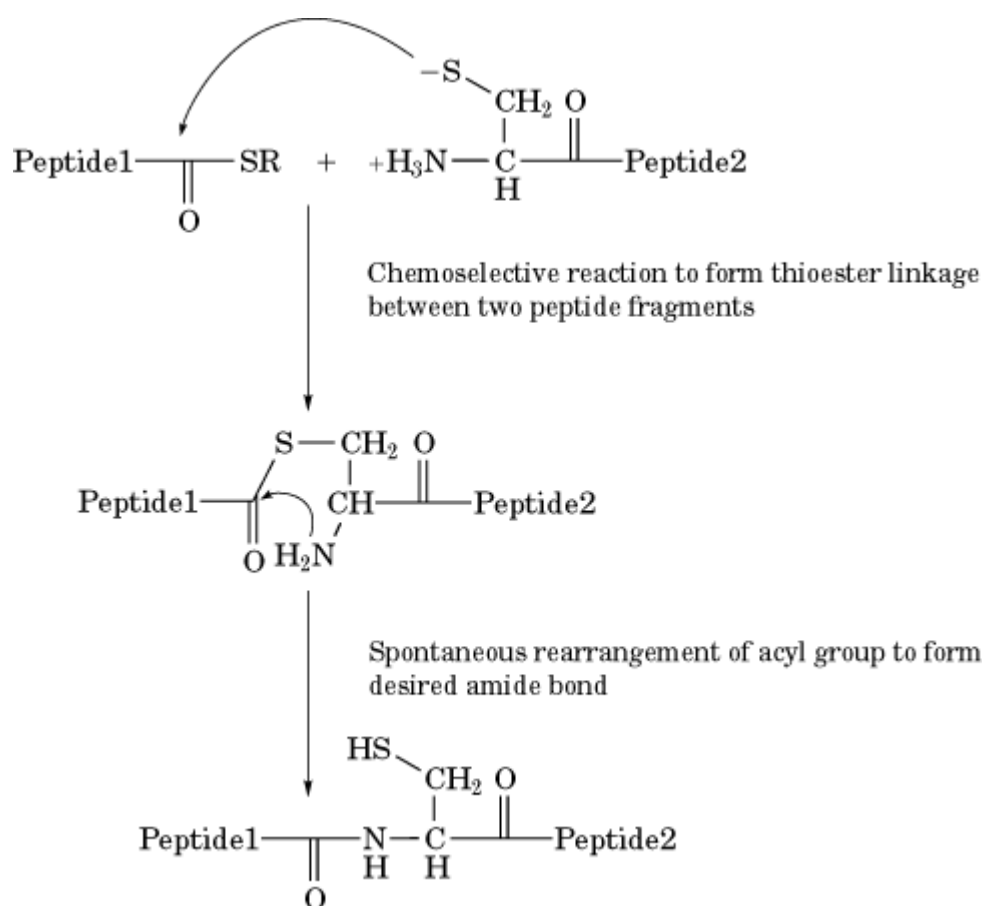
#### 4. Minimal Versus Maximal Protection

In order to get selectivity in the reaction, all reactive groups except the two that are desired to react can be protected transiently with chemically stable protecting groups. This is referred to as *maximal protection* and would be preferred, except for incomplete removal and side reactions associated with deprotection of the fully protected peptide product containing a heterogeneous set of protecting groups. In solid-phase synthesis, assembly of the correct sequence is often quite facile, but incomplete removal of the protecting groups can lead to an intractable mixture. By careful control of reaction conditions (aqueous, pH, etc.), one can eliminate some of the side-chain protection to minimize problems associated with protecting group removal. This is referred to as *minimal protection* and is preferred if selective acylation of the amino terminus by the carboxyl of the *N*-terminal peptide can be achieved.

#### 5. Chemical Ligation

Certain groups of chemical functionality have exceptional affinities for selective reaction, particularly in the formation of small rings. Recent developments have led to the synthesis of large peptide fragments by solid-phase synthesis, removal from the polymeric support and of most, if not all, side-chain protecting groups, and then selective coupling of the unprotected fragments based on chemical ligation. The most developed strategy for this approach is the use of an *N*-terminal [cysteine](#) residue with a special reactive *C*-terminal group on the other peptide. This approach was initially conceived by Kemp in his thiol capture strategy (7) and was reduced to a practically successful general approach by the groups of Kent (8-10) and Tam (11-13). For illustrative purposes, the native chemical ligation procedure of Dawson et al. (9) is described (Fig. 4), because it seems to have had the most practical impact. In this case, solid-phase synthesis is used to prepare two unprotected peptide segments that are combined in aqueous solution. The *C*-terminal fragment contains an *N*-terminal cysteine residue, and the *C*-terminal peptide fragment is prepared as the thioester. The thioester is displaced by the thiolate anion of a Cys residue. If the Cys is *N*-terminal, an acyl migration through formation of a five-membered ring occurs to generate the desired stable amide bond. If the sulfur atom of a Cys residue within the peptide chain is involved, then the thioether formed is capable of being displaced by other thiols, until the fragment migrates to the *N*-terminal Cys, when the stable rearrangement can occur. An alternative strategy, using an *N*-terminal *b*-bromoalanine of fragment two and the *C*-terminal thioester of fragment one to give the same covalent thioester intermediate by thioesterification, has been explored by Tam et al. (12)

**Figure 4.** Native chemical ligation of two unprotected peptide fragment using thiol–thioester exchange with ligation at the *N*-terminal cysteine residue (9).



While many proteins contain cysteine residues, and some have them optimally spaced for fragment assembly, many do not, and one of the research goals in this area is the extension of these concepts to allow coupling at residues other than the *N*-terminus of cysteine residues. Some success has been achieved with ligation at X-Gly and Gly-X sites (14), as well as at X-His sites (15). Another development that bodes well for the chemical synthesis of larger proteins is the adaptation of chemical ligation to solid-phase using [agarose](#) as the polymeric support. This should allow repetitive ligation of smaller fragments to generate the larger protein stepwise on the polymer.

## 6. Applications

One might question why one would want to prepare small proteins when one can easily express the genetic message in a variety of biological systems, including eukaryotic systems that include [post-translational modification](#). Limitations of normal expression systems to the 20 common amino acids restricts chemical modification of the parent protein for a variety of biophysical and pharmaceutical applications. Incorporation via molecular biology of novel amino acids as spectroscopic probes, or to introduce conformational constraints, is feasible using the [suppressor mutation](#) approach pioneered by the Schultz group (16-18). In this case, however, steric limitations of the **protein biosynthetic** machinery preclude the use of certain unusual amino acids or constrained dipeptides (19), which are readily accessible via organic synthesis. The quantities of protein obtained by this approach are generally small, requiring assay methods (enzymatic or spectroscopic) that are sensitive as well as specific.

### 6.1. Peptide Libraries

The explosion of combinatorial chemistry in the pharmaceutical industry as a paradigm for drug discovery was foreshadowed in the use of [peptide libraries](#) by Geysen (20) to map peptide **epitopes** or **antigenic** sites on proteins. Numerous strategies (21-23) to synthesize mixtures of thousands to

millions of peptides and allow selection of those with the desired activities (24) have developed over the past 20 years. [Combinatorial synthesis](#) of individual compounds for assay has developed rapidly, because it eliminates the problem of deconvolution of the mixture to identify the individual compound responsible for the observed activity. It also simplifies the pharmacological evaluation of compounds, due to concerns over synergy or the inhibition of activity that is potentially possible in the bioassay of mixtures. Nevertheless, the most effective strategy for lead generation and optimization will be determined by the cost and efficiency of the biological screen versus that of synthesis.

## Bibliography

1. T. W. Muir and S. B. Kent (1993) *Curr. Opin. Biotech.* **4**, 420–427.
2. J. Schneider and S. B. Kent (1988) *Cell* **54**, 363–368.
3. M. Miller, B. K. Sathyanarayana, A. Wlodawer, M. V. Toth, G. R. Marshall, L. Clawson, L. Selk, J. Schneider, and S. B. H. Kent (1989) *Science* **246**, 1149–1152.
4. R. B. Merrifield and J. M. Stewart (1965) *Nature* **207**, 522–523.
5. M. Schnolzer, P. Alewood, A. Jones, D. Alewood, and S. B. Kent (1992) *Int. J. Peptide Protein Res.* **40**, 180–193.
6. L. A. Carpino, D. Sadat-Aalae, H. G. Chao, and R. H. DeSelms (1990) *J. Am. Chem. Soc.* **112**, 9651–9652.
7. D. S. Kemp and R. I. Carey (1993) *J. Org. Chem.* **58**, 2216–2222.
8. M. Schnolzer and S. B. Kent (1992) *Science* **256**, 221–225.
9. P. E. Dawson, T. W. Muir, I. Clark-Lewis, and S. B. Kent (1994) *Science* **266**, 776–779.
10. T. W. Muir, P. E. Dawson, and S. B. Kent (1997) *Methods Enzymol.* **289**, 266–298.
11. C.-F. Liu and J. P. Tam (1994) *J. Am. Chem. Soc.* **116**, 4149–4153.
12. J. P. Tam, Y. A. Lu, C. F. Liu, and J. Shao (1995) *Proc. Natl. Acad. Sci. USA* **92**, 12485–12489.
13. C. F. Liu, C. Rao, and J. P. Tam (1996) *J. Am. Chem. Soc.* **118**, 307–312.
14. L. E. Canne, S. J. Bark, and S. B. H. Kent (1996) *J. Am. Chem. Soc.* **118**, 5891–5896.
15. L. S. Zhang and J. P. Tam (1997) *Tetrahedron Lett.* **38**, 3–6.
16. C. J. Noren, S. J. Anthony-Cahill, M. C. Griffith, and P. G. Schultz (1989) *Science* **244**, 182–188.
17. D. Mendel, J. Ellman, and P. G. Schultz (1993) *J. Am. Chem. Soc.* **115**, 4359–4360.
18. J. A. Ellman, D. Mendel, and P. G. Schultz (1992) *Science* **255**, 197–200.
19. S. Hanessian, G. McNaughton-Smith, H.-G. Lombart, and W. D. Lubell (1997) *Tetrahedron* **53**, 12789–12854.
20. H. M. Geysen (1985) *Immunol. Today* **6**, 364–369.
21. R. A. Houghten, C. Pinilla, S. E. Blondelle, J. R. Appel, C. T. Dooley, and J. H. Cuervo (1991) *Nature* **354**, 84–86.
22. K. S. Lam, S. E. Salmon, E. M. Hersh, V. J. Hruby, W. M. Kazmierski, and R. J. Knapp (1991) *Nature* **354**, 82–84.
23. A. Furka (1995) *Drug Dev. Res.* **36**, 1–12.
24. G. Koppel, C. Dodds, B. Houchins, D. Hunden, D. Johnson, R. Owens, M. Chaney, T. Usdin, B. Hoffman, and M. Brownstein (1995) *Chem. Biol.* **2**, 483–487.

## Suggestions for Further Reading

25. B. Merrifield (1993) *Life During a Golden Age of Peptide Chemistry*; American Chemical Society, Washington, D.C.
26. B. Gutte, ed., (1995) *Peptides: Synthesis, Structures and Applications*, Academic Press: New York.

27. M. Bodansky (1984) *Principle of Peptide Synthesis*, Springer-Verlag, New York.
28. M. Bodansky and A. Bodansky (1984) *The Practice of Peptide Synthesis*, Springer-Verlag, New York.

## Peptidyl Prolyl *Cis/Trans* Isomerases

Only a few [cis/trans isomerizations](#) are catalyzed by enzymes. Most of these enzymes catalyze isomerization of C—C double bonds. However, the highly abundant and ubiquitously distributed peptidyl prolyl *cis/trans* isomerases (PPIases) have evolved to accelerate the rotation about a formal single bond. They catalyze the *cis/trans* isomerization of **imide** peptidyl-proline bonds and are inactive toward the -CONH- peptide moiety. In enzyme nomenclature, PPIases are classified under EC number 5.2.1.8.

Initially, PPIases were identified in the cytosol of kidney cells from their ability to accelerate the *cis* to *trans* isomerization of chromogenic tetrapeptides of the type succinyl-Ala-Xaa-Pro-Phe-4-nitroanilides (where Xaa is any natural amino acid) by an assay that utilizes the conformational specificity of **chymotrypsin** (1). Briefly, in this *cis/trans* assay, the *trans* prolyl bond isomer of the substrate succinyl-Ala-Ala-Pro-Phe-4-nitroanilide (about 92% of the total) is cleaved completely at the anilide bond in a few seconds when high concentrations ( $\gg$ mM) of the proteinase are used. The remaining *cis* isomer of the peptide (8%) is resistant to chymotrypsin until it undergoes *cis* to *trans* isomerization. When monitored by the absorbance of the released product (4-nitroaniline), the isomerization is the **rate-limiting** step for proteolyzing the *cis* isomers indirectly. It can be quantified as a first-order reaction with a half-time of about 90 s at 10°C. This kinetic phase, which is equivalent to a quasi-irreversible *cis* to *trans* isomerization, has an accelerated rate in the presence of catalytic amounts of a PPIase. The calculated rate constants depend linearly on the PPIase concentration. This assay could not demonstrate catalysis in both directions, but this was possible with NMR (2, 3). The standard assay was greatly improved by dissolving the assay peptides in LiCl/trifluoroethanol to increase the fraction of *cis* isomer (4).

PPIases were initially anticipated to be folding/rearrangement catalysts for proteins bearing *cis* prolyl bonds in the native state because the slow *trans* to *cis* interconversion following **protein synthesis** on the [ribosome](#) will seriously limit the rate of *in vivo* **protein folding**. Indeed, PPIases increase the rate of the slow kinetic phases in *in vitro* refolding of many denatured proteins known to be limited in rate by prolyl isomerization. In support of a physiological role in protein folding, the involvement of a PPIase in triple-helix formation in type I pro-[collagen](#) was shown in chick embryo fibroblasts *in cellulo* (5). However, the biological functions of PPIases are not yet well understood, and they may have more general roles than affecting only *de novo* protein folding.

According to our current knowledge, the enzyme class of PPIases comprises three families with unrelated amino acid sequences: **cyclophilins** (abbreviated Cyp), **FK506-binding proteins** (abbreviated FKBP), and parvulins. An additional subfamily, the prokaryotic [trigger factors](#), exhibits weak sequence similarity to FKBP proteins but lacks FK506-binding ability (6). The amino acid sequences of several dozen different PPIases are known, encompassing molecular masses from 10.1 kDa (92 amino acids in length) for *E. coli* parvulin (7) to 158 kDa for the NK-TR1 cyclophilin of human, large granular lymphocytes (8). Typically, *Escherichia coli* cells harbor two cyclophilins, four FKBP-like PPIases, the archetype of parvulins, along with the SurA protein containing two parvulin domains, and the ribosome-bound trigger factor. Among the 470 predicted coding regions in the free-living bacterium [Mycoplasma genitalium](#), which has the smallest known genome of any

free living organism, a trigger factor homologue is the sole PPIase present (9).

Members of the PPIase families resemble perfectly evolved enzymes for catalytic efficiency (10). For oligopeptide substrates, both the Michaelis constant,  $K_m$ , and the turnover number,  $k_{cat}$ , are large and yield values of  $k_{cat}/K_m$  of up to  $>10^7 M^{-1} s^{-1}$ . These bimolecular rate constants approach the **diffusion-controlled** limits for enzyme reactions (see [Enzymes](#)). With  $k_{cat}/K_m = 1.1 \times 10^6 M^{-1} s^{-1}$ , the *E. coli* trigger factor catalyzing refolding of a denatured **ribonuclease T1** variant exemplifies a particularly efficient catalysis in protein folding (11). This efficacy results from having extended secondary binding sites to bind the unfolded protein substrate favorably. These subsites encompass the two protein domains flanking the central PPIase domain of the enzyme. Many of the PPIase enzymes are specific only for certain amino acid residues flanking the proline in the primary structure, but this is less important with the cyclophilins. Conversely, a **hydrophobic** side chain of the amino acid preceding the proline enhances catalysis by FKBP.

PPIases applied in catalytic amounts cannot alter the *cis/trans* equilibrium of the substrates, as expected. When protein refolding is hindered by misfolding and protein aggregation, an increase of the yield of native protein is not obtained by supplementing the refolding buffer with a PPIase. Thus, an important aspect of **chaperones** is missing in PPIases (12).

Immunosuppressive compounds, like cyclosporin A, FK506 and [rapamycin](#), bind tightly to the respective active site of cyclophilins and FKBP and inhibit their enzymatic activities reversibly and competitively and have inhibition constants in the picomolar to micromolar range. Cross-inhibition of the different classes does not occur. At the present time, no specific reversible inhibitors for parvulins and trigger factors have been reported.

Definitive identification of the natural substrates of PPI-ases has not been reported yet. Nevertheless, a number of binding proteins, which may include the putative cellular substrates, were identified by chemical **cross-linking**, affinity chromatography, the yeast two-hybrid screen, and by co-purification. For example, the nuclear parvulin-like PPIase Pin1 has affinity for NIMA protein [kinase](#) in the two-hybrid system. The catalytic activity of this isomerase is involved in **cell-cycle** control in **eukaryotes** (13). For the trigger factor family, association with nascent chains derived from secreted and cytosolic proteins and specific affinity for the 50S subunit of the *E. coli* ribosome are indicated by cross-linking and fractionation of enzyme activity (6, 14, 15). On the other hand, the same PPIase has been detected as an indispensable component of the heterooligomeric complex of the chaperone GroEL with the degradable fusion protein CRAAG. This system has been used to determine the requirements for rapid degradation of abnormal, misfolded proteins by the proteinase ClpP in *E. coli* cells (16). When a trace of copurifying protein of the ryanodine receptor was analyzed, it became obvious that intracellular calcium release channels on the **endoplasmic** or sarcoplasmic reticulum (ryanodine receptor, RyR) and the **inositol** 1,4,5-triphosphate receptor constitute heterooligomeric complexes of the type  $((RyR)_4/(FKBP12)_4)$  with either cytosolic FKBP12 or FKBP12.6 (17, 18). A dysfunctional state, associated with the impairment of the receptor-mediated  $Ca^{2+}$  retardation of the channel, arises from inhibiting PPIase activity with both FK506 and rapamycin (19). Furthermore, FKBP and cyclophilins in the range of 40 to 60 kDa were found as components of the unactivated **steroid receptor** complex, associated with the **hsp90** component of the receptor (20). Interestingly, two members of different families of proteins that assist protein folding, an enzyme and a chaperone, coexist as binding partners in this receptor complex.

## Bibliography

1. G. Fischer, H. Bang, and C. Mech C. (1984) Biomed. Biochim. Acta **43**, 1101–1111.
2. D. Kern et al. (1995) Biochemistry **34**, 13594–13602.
3. R.M. Justice, Jr. et al. (1990) Biochem. Biophys. Res. Commun. **171**, 445–450.

4. J.L. Kofron et al. (1991) *Biochemistry* **30**, 6127–6134.
5. B. Steinmann, P. Bruckner, and A. Superti-Furga (1991) *J. Biol. Chem.* **266**, 1299–1303.
6. G. Stoller et al. (1995) *EMBO J.* **14**, 4939–4948.
7. J.U. Rahfeld et al. (1994) *FEBS Lett.* **352**, 180–184.
8. S.K. Anderson et al. (1993) *Proc. Natl Acad. Sci. USA* **90**, 542–546.
9. C.M. Fraser et al. (1995) *Science* **270**, 397–403.
10. J.J. Burbaum, R.T. Raines, W.J. Albery, and J.R. Knowles (1989) *Biochemistry* **28**, 9293–9305.
11. C. Scholz et al. (1997) *EMBO J.* **16**.
12. G. Kern, D. Kern, F.X. Schmid, and G. Fischer (1994) *FEBS Lett.* **348**, 145–148.
13. K.P. Lu, S.D. Hanes, and T. Hunter (1996) *Nature* **380**, 544–547.
14. Q.A. Valent et al. (1995) *EMBO J.* **14**, 5494–5505.
15. T. Hesterkamp, S. Hauser, H. Lutcke, and B. Bukau (1996) *Proc. Natl Acad. Sci. USA* **93**, 4437–4441.
16. O. Kandror, O.M. Sherman, M. Rhode, and A.L. Goldberg (1995) *EMBO J.* **14**, 6021–6027.
17. H.B. Xin et al. (1995) *Biochem. Biophys. Res. Commun.* **214**, 263–270.
18. A.P. Timerman et al. (1996) *J. Biol. Chem.* **271**, 20385–20391.
19. A.R. Marks (1996) *Trends Cardiovasc. Med.* **6**, 130–135.
20. T. Ratajczak and A. Carrello (1996) *J. Biol. Chem.* **271**, 2961–2965.

### Suggestions for Further Reading

21. J.E. Kay (1996) Structure-function relationships in the FK506-binding protein (FKBP) family of peptidyl-prolyl *cis/trans* isomerases, *Biochem. J.* **314**, 361–385; an important reference source about FKBP.
22. G. Fischer (1994) Peptidyl-prolyl *cis/trans* isomerases and their effectors, *Angew. Chem., Int. Ed. Engl.* **33**, 1415–1436.
23. F.X. Schmid (1993) Prolyl isomerase: Enzymatic catalysis of slow protein-folding reactions, *Ann. Rev. Biophys. Biomol. Struct.* **22**, 123–143; summarizes PPIase catalysis in protein folding.
24. A. Galat and S. M. Metcalfe (1995) Peptidyl proline *cis/trans* isomerases, *Prog. Biophys. Mol. Biol.* **63**, 67–118; contains an exhaustive bibliography of many aspects of PPIases.

## Peptidyl Transferase

Peptidyl transferase is an integral part of the large subunit of [ribosomes](#) and catalyzes **peptide-bond** formation in the elongation step of [translation](#) during protein biosynthesis. Peptide-bond formation involves an acyl group *O*-to-*N* migration and converts an ester to an amide. Peptidyl-[transfer RNA](#) in the P site is the ester bond of high energy content that donates its peptidyl group to the amino group of aminoacyl-tRNA bound in the A site of the ribosome. The peptidyl transferase reaction is often measured by the [puromycin](#) reaction, in which puromycin acts as the acceptor substrate to form peptidyl-puromycin. In bacteria, peptidyl transferase activity can be detected in the 50 S ribosomal subunit, and several ribosomal proteins are thought to play a role in the activity. The primary activity of peptidyl transfer, however, is a **ribozyme** activity encoded in the 23 S rRNA. The peptidyl transferase activity of *Escherichia coli* and *Thermus aquaticus* ribosomes are resistant to



conventional protein extraction procedures (1), and the formation of a G-C pair between G2252 in the conserved hairpin loop of *E. coli* 23 S rRNA and C74 in the 3'-terminal region of tRNA is a prerequisite for the peptidyl transfer reaction. These findings suggest that not all the large subunit proteins are necessarily required for peptidyl transferase activity. It has been shown recently that *E. coli* 23 S rRNA synthesized by *in vitro* transcription with T7 **RNA polymerase** exhibits the primary peptidyl transfer activity in the absence of any ribosomal proteins *in vitro*, and that the truncated domain V transcript alone has this activity (2).

Peptidyl transferase is also involved in the termination reaction of protein synthesis and hydrolyzes the ester bond between the peptide chain and tRNA at the P site of the ribosome (see [Release Factor](#)). Peptidyl transferase is inhibited by various antibiotics, such as [chloramphenicol](#), lincomycin, carbomycin (prokaryotes), and [cycloheximide](#) (eukaryotes).

### Bibliography

1. H. F. Noller, V. Hoffarth, and L. Zemniak (1992) *Science* **256**, 1416–1419.
2. I. Nitta, Y. Kamada, H. Noda, T. Ueda, and K. Watanabe (1998) *Science* **281**, 666–669.

### Suggestion for Further Reading

3. K. S. Wilson and H. F. Noller (1998) *Cell* **92**, 337–349.

## Perinuclear Space

The perinuclear space is the luminal space between the inner and outer [nuclear envelope](#) bilayers that separate the [nucleus](#) from the **cytoplasm**. The width of the space is between 100 Å and 500 Å. Because the outer nuclear membrane is continuous with the [endoplasmic reticulum](#), the perinuclear space is continuous with the **lumen** of the endoplasmic reticulum.

## Permissive Condition

The condition under which a **conditional lethal mutant** grows is called the permissive condition, in contrast to the [nonpermissive condition](#) or restrictive condition, which restricts growth (see [Conditional Lethal Mutations](#)).

## Peroxidase

Peroxidases occur throughout the biosphere and **catalyze** the oxidation of various substrates at the expense of peroxide. Peroxidases are involved in **plant** growth hormone metabolism, the production

of thyroxine, neutrophil-mediated detoxification reactions, prostaglandin biosynthesis, and lignin degradation, to name a few. Although peroxidases are found in plants, animals, and **microorganisms**, the most thoroughly understood are the nonmammalian heme-containing peroxidases. They all are single [polypeptide chains](#) in the range of 30 to 40 kDa and contain a single ferric protoporphyrin IX as the **prosthetic group**. Based on [sequence analysis](#) and now **X-ray crystallographic** structures, the nonmammalian peroxidases have been divided into three classes (1): intercellular (Class I), extracellular fungal (Class II), and extracellular plant peroxidases (Class III). The extracellular enzymes are glycosylated, and in the case of Class III enzymes there often are many **isozymes** with varying degrees of glycosylation. The extracellular peroxidases contain two **calcium-binding** sites and at least four [disulfide bonds](#) presumably for the additional stability required of secreted enzymes.

Now a wealth of structural information is available about peroxidases because of recent crystal structure determinations. Coupled with the cloning and expression of several peroxidases in recombinant systems, this has opened the way for using protein engineering methods to study structure–function relationships.

## 1. Catalytic Cycle

The overall peroxidase cycle occurs in three distinct steps:



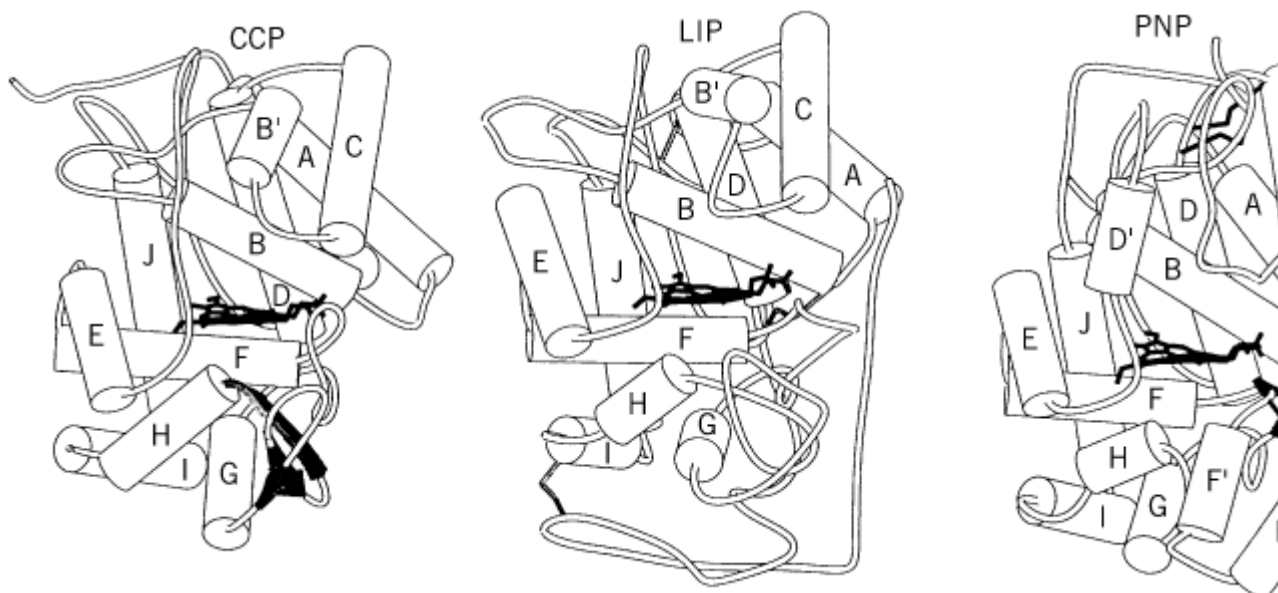
where P is the enzyme porphyrin whose heme iron is indicated, and S is the substrate. In Step 1, the peroxide removes two electrons from the enzyme to give what is called Compound I. In doing so, the peroxide O–O bond is broken, produces water, and the second peroxide-derived oxygen atom remains coordinated to the iron. One electron is removed from the iron to give the oxy-ferryl intermediate ( $\text{Fe}^{4+} = \text{O}$ ), and a second electron is removed from the porphyrin to give a porphyrin p cation radical. In some peroxidases, an amino acid side chain is oxidized to a radical rather than the porphyrin. Normally, the formation of Compound I is quite easy to follow because the resting enzyme is red/brown in color, whereas Compound I is green. In Step 2, the porphyrin radical is reduced, producing a substrate radical,  $\text{S} \cdot$ . Because of porphyrin reduction, Compound II is no longer green, but red. Finally in Step 3, Compound II is reduced by a second substrate molecule. In plant peroxidases, the substrates are normally small aromatic molecules. Once substrate radicals form, they dimerize or disproportionate in nonenzymatic reactions.

## 2. Three-Dimensional Structures

From 1980 to 1993, only one nonmammalian heme peroxidase crystal structure was known, [cytochrome c](#) peroxidase. Since then an additional six structures have been determined (2), and there may well be more on the way. Representative examples of Class I, II, and III peroxidases are shown in Fig. 1. The sequence identity can be as low as 15% between peroxidases, yet the crystal structures have a similar fold independent of sequence [homology](#) or **phylogenetic** origin. This underscores the important “rule” in structural biology that [tertiary structure](#) is far more conserved than [primary structure](#). All of the various peroxidases consist of a core of 10 alpha helices, two of which, the B

and F helices, sandwich the heme in place. The various peroxidases differ in structure primarily on the surface. For example, Class III peroxidases have two additional  $\alpha$ -helices on the surface (F' and F''), where cytochrome *c* peroxidase has a **beta**-sheet structure. It is thought that these additional helices in the Class III enzymes may play a role in substrate binding.

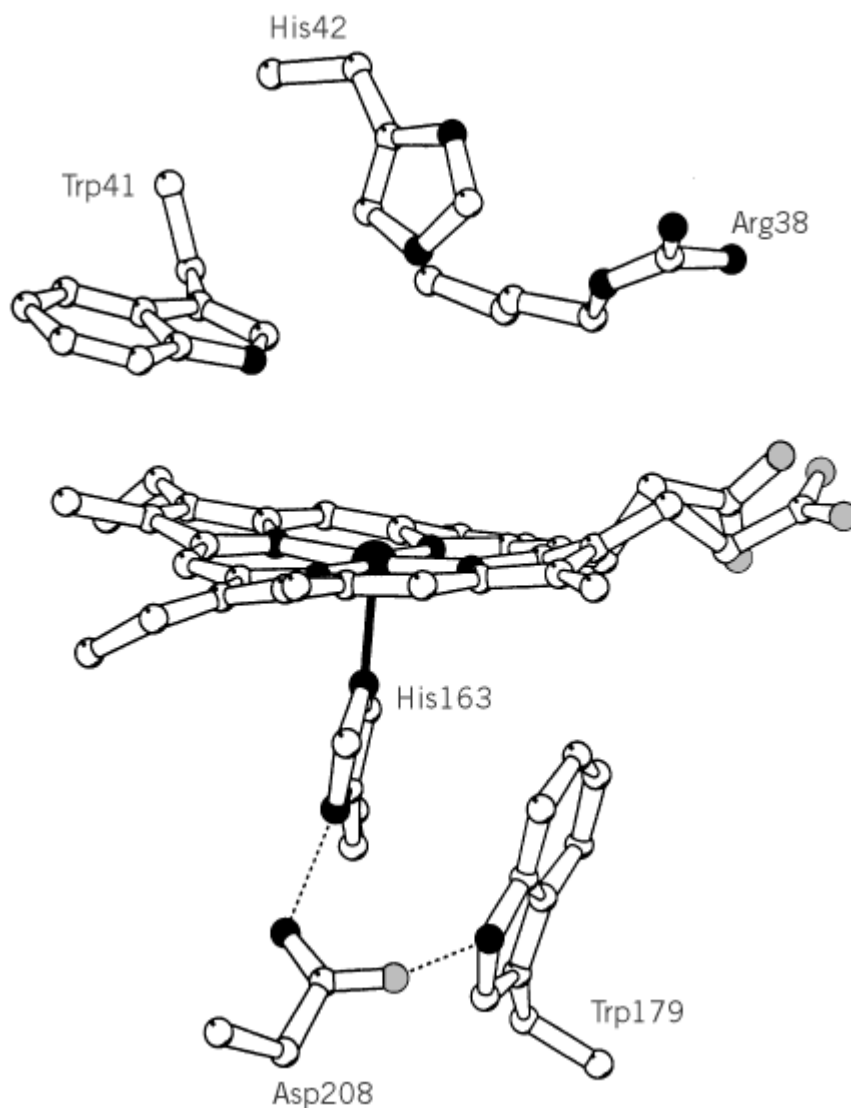
**Figure 1.** Representative crystal structures from each of the three primary classes of nonmammalian heme peroxidases: I (CCP, cytochrome *c* peroxidase); Class II (LIP, lignin peroxidase); Class III (PNP, peanut peroxidase). The  $\alpha$ -helices are as labeled cylinders, and the  $\beta$ -sheet structure as shaded ribbons.



### 3. Active-Site Structures

Not surprisingly, the [active sites](#) of heme peroxidases are highly conserved (Fig. 2). The proximal pocket contains the His residue coordinated to the heme iron. This His ligand donates a [hydrogen bond](#) to the conserved Asp residue. It is thought that this interaction imparts a greater imidazolate character to the His ligand, resulting in extra negative charge on the His residue that stabilizes the  $\text{Fe}^{3+}$  state of the heme. This is one reason that peroxidases exhibit a lower redox potential in the  $\text{Fe}^{3+}/\text{Fe}^{2+}$  couple than the [globins](#), even though the globins and peroxidases both use a His ligand. Directly adjacent to the His ligand, Most peroxidases have an aromatic residue stacked parallel to and in contact with the His ligand. This residue is Phe in most peroxidases, but Trp in the two known Class I peroxidase structures. For the class I peroxidases, the Trp indole ring N atom donates a hydrogen bond to the conserved proximal Asp residue (Fig. 2).

**Figure 2.** The active site structure of ascorbate peroxidase, a representative example of a nonmammalian heme peroxidase.



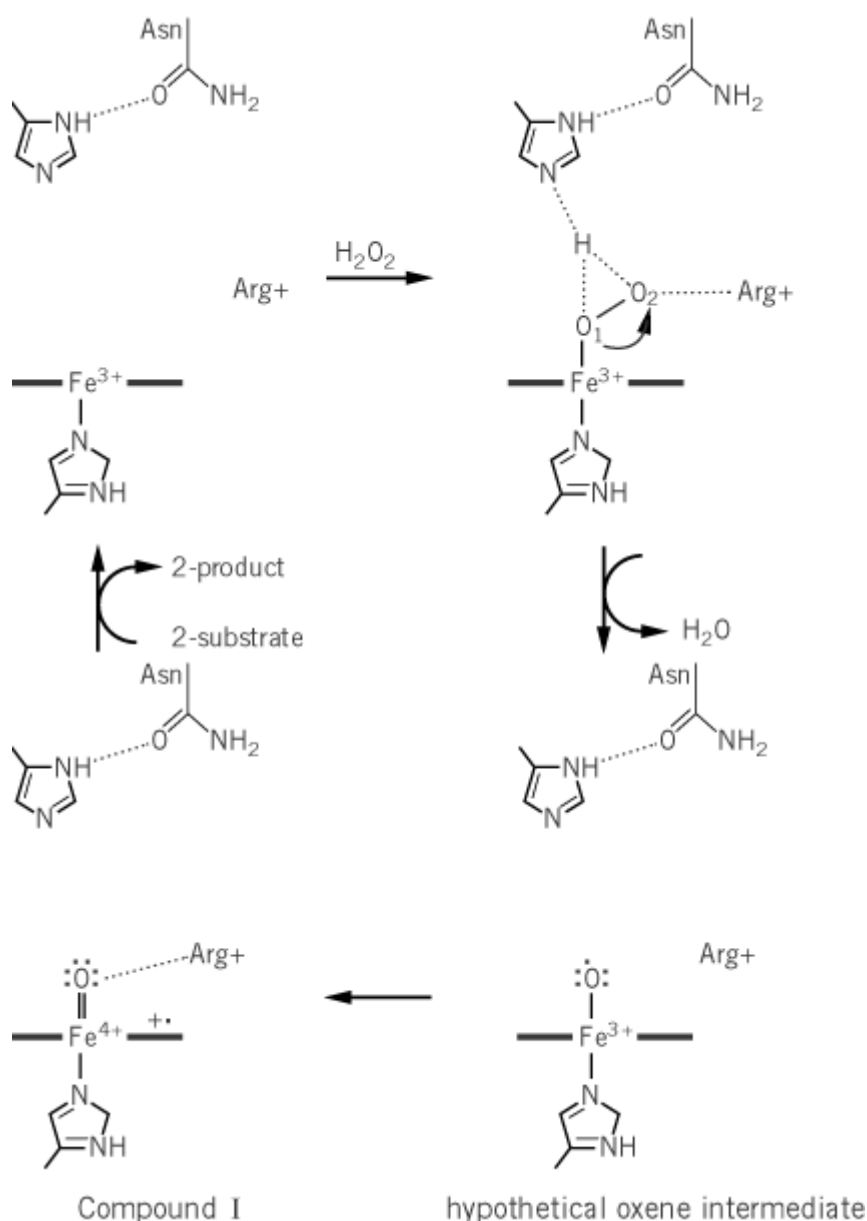
The B helix provides key catalytic residues on the opposite, distal side of the heme. It is thought that the distal His and Arg residue operate in concert to help cleave the peroxide O–O bond. Stacked parallel to the distal heme surface is a conserved aromatic residue, which is Phe in the known crystal structures, except for the Class I peroxidases, which use Trp.

#### 4. Catalytic Mechanism

Most of the effort in understanding the peroxidase mechanism has centered on the formation of Compound I. It is thought that the distal His and Arg residues operate in concert to ensure that the peroxide O–O bond is cleaved heterolytically (3; Fig. 3). The distal His operates as an acid/base catalyst by removing a proton from the incoming peroxide and delivering it to the peroxide O<sub>2</sub> atom. As shown in Fig. 3, the distal His donates a hydrogen bond to an invariant Asn residue, which ensures that the N<sub>δ2</sub> of the distal His is free to accept a proton from the incoming peroxide. The Arg residue may help to stabilize a developing negative charge on the leaving O<sub>2</sub> atom. The net result is heterolytic fission of the peroxide O–O bond. This leaves behind an “oxene” O atom that contains only six valent electrons, a potent oxidant. The oxene center removes one electron from the iron atom, to give Fe<sup>4+</sup> and one from a neighboring organic group, normally the porphyrin. The crystal structure of Compound I (4) shows that the distal Arg residue swings in to hydrogen-bond with the

iron-linked oxygen atom. This may account in part for the stability of the  $\text{Fe}^{4+} = \text{O}$  center in Compound I.

**Figure 3.** The catalytic mechanism of Compound I formation.



The formation of Compound I is stoichiometric and has a second-order rate constant of between  $10^7$  and  $10^8 \text{M}^{-1}\text{s}^{-1}$  for most heme peroxidases. [Site-directed mutagenesis](#) has been used to probe this mechanism by amino acid substitution of key catalytic groups in the active site. Changing the distal His to Leu (5) lowers the rate of Compound I formation by  $10^5$ , demonstrating that the distal His is indeed a critical residue. Somewhat surprisingly, the nature of the proximal ligand is not important for Compound I formation because a Gln that replaces the proximal His residue affects only the stability of Compound I and not the rate of formation (6).

## 5. Enzyme Substrate Complexes

All of the known heme peroxidase structures show that one heme edge is accessible to small aromatic molecules that are typical peroxidase substrates. It is generally thought that small substrates interact at this site and deliver electrons directly to the heme, resulting in reduction of the Compound I porphyrin radical or Compound II  $\text{Fe}^{4+} = \text{O}$  center. Nevertheless, there is no direct structural evidence for this view. To date, only the structures of two heme peroxidase-substrate complexes are known. Manganese peroxidase oxidizes  $\text{Mn}^{2+}$  to  $\text{Mn}^{3+}$ , and the  $\text{Mn}^{3+}$  operates in a complex with dicarboxylic acids as a diffusible oxidant of lignin. The crystal structure of manganese peroxidase (7) shows that  $\text{Mn}^{2+}$  coordinates with protein carboxylates and one heme propionate. Cytochrome *c* peroxidase is unusual among the known heme peroxidase structures because its substrate is another protein, cytochrome *c*. The crystal structure of the cytochrome *c* peroxidase–cytochrome *c* complex (8) has provided important insights into the nature of interprotein **electron transfer** reactions.

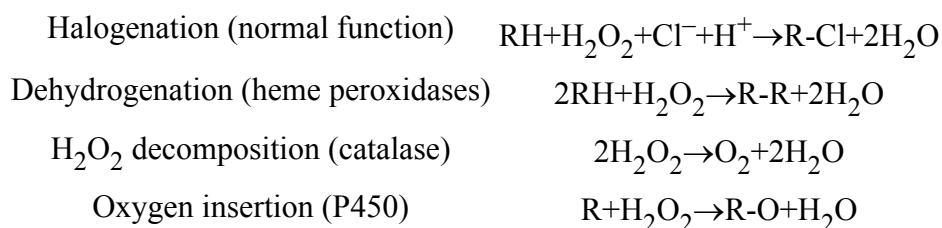
## 6. Heme Halo Peroxidases

The crystal structure of one heme halide peroxidase is known (9). Chloroperoxidase (CPO) is secreted by the mold *Caldariomyces fumago*, where it functions to chlorinate 1,3-cyclopentanedione to the natural product 2,2-dichloro-1,3-cyclopentanediol (caldariomycin). The catalytic cycle of CPO is similar to other heme peroxidases. However, in CPO the oxidizing power of Compound I is used to oxidize  $\text{Cl}^-$  ions to HOCl, and then the HOCl operates as an organic chlorinating agent. Traditional heme peroxidases cannot catalyze this reaction probably because the redox potential of  $\text{Cl}^-$  is too high. However, CPO operates at very low pH ( $\approx 3$ ) where the redox potential and hence the reactivity of Compound I may be significantly increased over other peroxidases that operate closer to physiological pH.

CPO consists of a single 30-kDa polypeptide chain, one heme group, and one disulfide bond. The crystal structure shows at least 14 different sites of **glycosylation**. The fold is unlike that of other heme peroxidases, and the proximal ligand is a Cys residue not His. This is very similar to [cytochromes P450](#), which also use a Cys ligand. As a result, CPO shares spectral features characteristic of P450s. The most significant is the 450-nm band in the visible absorption spectrum when the enzyme is reduced in the presence of carbon monoxide to give the  $\text{Fe}^{2+}$ –CO complex.

### 6.1. Reactions Catalyzed by Chloroperoxidase

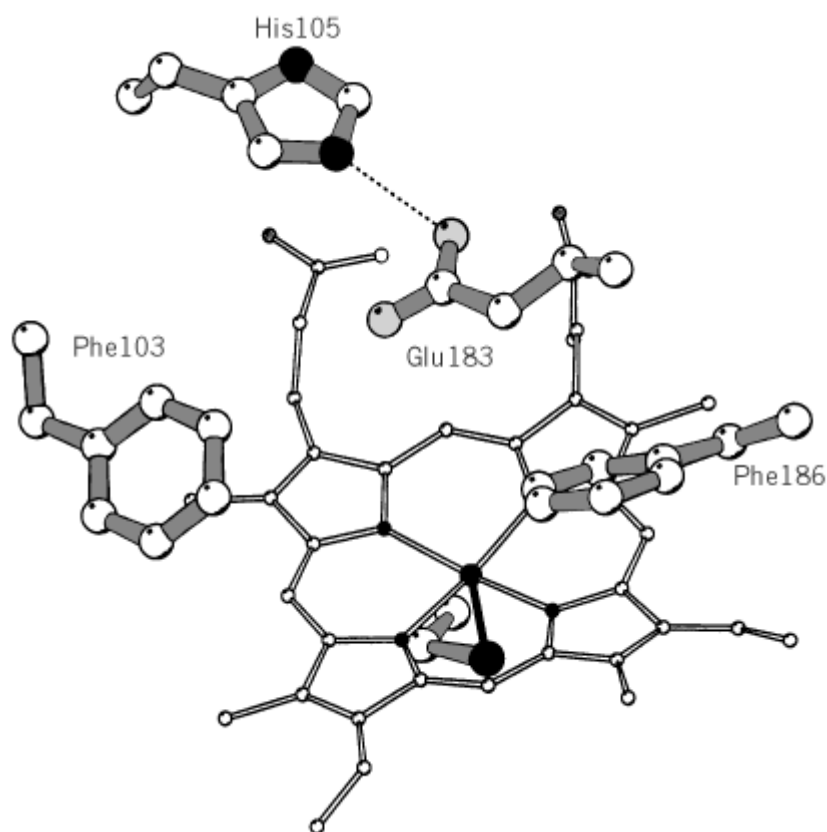
CPO is one of the most diverse of the known heme enzymes. In addition to operating as a chloroperoxidase, CPO exhibits activities characteristic of traditional heme peroxidase, [catalase](#), and cytochrome P450s:



Sulfoxidation and epoxidation reactions catalyzed by CPO proceed with a high degree of enantioselectivity (10, 11). Even the typical peroxidase dehydrogenation reactions proceed with enantioselectivity, something other peroxidases cannot do. The CPO crystal structure provides some insights into these unusual functional properties. Although CPO is a helical protein, the overall fold is unique not the typical peroxidase fold. In addition, the heme is accessible from the top distal surface, not the heme edge. Bracketing the distal heme pocket near the  $\text{Fe}^{4+} = \text{O}$  center are two key

Phe residues and other **hydrophobic** groups that form a shallow pocket that is the most probable site for substrate binding (Fig. 4). The asymmetrical nature of this pocket can account for the observed enantioselectivity.

**Figure 4.** The active site structure of chloroperoxidase. Glu183 most likely serves the same acid/base catalytic role as the distal His in traditional heme peroxidases. The two Phe residues shown help form the substrate binding site that enables this enzyme to catalyze enantioselective sulfoxidation, epoxidation, and electron transfer reactions.



Like other peroxidases, CPO forms Compound I with an  $\text{Fe}^{4+} = \text{O}$  center and porphyrin radical. The catalytic groups, however, are quite different in CPO. Directly adjacent to the peroxide binding site, CPO has a Glu residue rather than His, and there is no active-site Arg in CPO.

## 7. Other Heme Peroxidases

Although the nonmammalian heme peroxidases have received most of the attention, the heme-containing peroxidases are very diverse. Mammals produce well-known peroxidases like myeloperoxidase, **lactoperoxidase**, and thyroid peroxidase. Although these peroxidases contain heme and have active sites similar to the nonmammalian heme peroxidases, these enzymes are larger and have different three-dimensional structures (12). The best understood is myeloperoxidase. The function of myeloperoxidase in white blood cells and the related eosinophil peroxidase is to oxidize halides and thiocyanate ions to produce HOCl, HOBr, and HOSCN, which are potent antipathogen oxidizing agents. The marine worm *Amphitrite ornata* produces a dehalogenating heme-containing peroxidase (13). *Mycobacterium tuberculosis*, the causative agent in tuberculosis, produces a heme catalase/peroxidase that metabolizes the most widely used antituberculosis drug, isoniazid. Strains resistant to isoniazid have a point mutation in the catalase/peroxidase, which implicates the peroxidase as a drug activator (14). Other bacteria, such as *Pseudomonas aeruginosa*, produce a

cytochrome *c* peroxidase that contains two heme groups (15).

## 8. NonHeme Peroxidases

A number of peroxidases do not have heme as the prosthetic group. The fungus *Curvularia inaequalis* produces a chloroperoxidase that utilizes vanadium (VO<sub>3</sub>) as a cofactor (16). Vanadium bromoperoxidases are also well known in **algae** (17). NADH peroxidase uses a flavin (18), and glutathione peroxidase is unique in utilizing an essential **selenocysteine** residue at the active site to oxidize **glutathione** and reduce hydroperoxides (19). There also are peroxidases that contain no metal center or other cofactors, such as the haloperoxidase of *Pseudomonas pyrrocinia* (20). In many of these, a detailed understanding is just developing based on recent advances in **cloning** and expression of these unusual and interesting peroxidases.

## Bibliography

1. K. G. Welinder (1992) *Curr. Biol.* **2**, 388–393.
2. M. Gajhede, D. J. Schuller, A. Henriksen, A. T. Smith, and T. L. Poulos (1997) *Nat. Struct. Biol.* **4**, 1032–1038.
3. T. L. Poulos (1987) *Adv. Inorg. Biochem.* **7**, 1–36.
4. S. L. Edwards, H. X. Nguyen, R. C. Hamlin, and J. Kraut (1987) *Biochemistry* **26**, 1503–1511.
5. J. E. Erman, L. B. Vitello, M. A. Miller, A. Shaw, K. A. Brown, and J. Kraut (1993) *Biochemistry* **32**, 9798–9806.
6. K. Choudhury, M. Sundaramoorthy, A. Hickman, T. Yonetani, E. Woehl, M. F. Dunn, and T. L. Poulos (1994) *J. Biol. Chem.* **269**, 20239–20249.
7. M. Sundaramoorthy, K. Kishi, M. H. Gold, and T. L. Poulos (1994) *J. Biol. Chem.* **269**, 32759–32767.
8. H. Pelletier and J. Kraut (1992) *Science* **258**, 1748–1755.
9. M. Sundaramoorthy, J. Ternier, and T. L. Poulos (1995) *Structure* **3**, 1367–1377.
10. S. Colonna, N. Gaggero, L. Casella, G. Carrea, and P. Pasta (1992). *Tet. Asymm.* **3**, 95–106.
11. E. J. Allain, L. P. Hager, L. Deng, and E. Jacobsen (1993). *J. Am. Chem. Soc.* **115**, 4415–4416.
12. R. Fenna, J. Zeng, and C. Davey (1995) *Arch. Biochem. Biophys.* **316**, 653–656.
13. Y. P. Chen, S. A. Woodin, D. E. Lincoln, and C. R. Lovell (1996) *J. Biol. Chem.* **271**, 4609–4612.
14. Y. Zhang, B. Heym, B. Allen, D. Young, and S. Cole (1992) *Nature* **358**, 591–593.
15. V. Fulop, C. J. Ridout, C. Greenwood, and J. Hajdu (1995) *Structure* **3**, 1225–1233.
16. A. Messerschmidt and R. Wever (1996) *Proc. Natl. Acad. Sci. USA* **93**, 392–396.
17. B. E. Krenn, M. G. Tromp, and R. Wever (1989) *J. Biol. Chem.* **264**, 19287–19292.
18. T. Stehle and G. E. Schulz (1992) *J. Mol. Biol.* **224**, 1127–1141.
19. O. Epp, R. Ladenstein, and A. Wendel (1983) *Eur. J. Biochem.* **133**, 51–69.
20. C. Wolfframm, F. Lingens, R. Mutzel, and K. H. van Pee (1993) *Gene* **130**, 131–135.

## Suggestion for Further Reading

21. T. L. Poulos and R. Fenna (1994) *Met. Biol.* **30**, 25–75. Review that leads back to primary literature.

## Pertussis Toxin



Among the several virulence factors implicated in the pathogenesis of whooping cough, a major role is played by pertussis toxin (PT), which is also a major determinant in the development of antipertussis protective [antibodies](#) (1). PT is active on a variety of cells and causes a wide range of effects, due to the **ADP-ribosylation** of different trimeric cytosolic **G proteins** involved in cell [signal transduction](#) (2).

PT is composed of five different subunits: S1 (21 kDa) is the catalytic part, and S2, S3, two copies of S4, and S5 form the cell-binding pentamer B. Despite the heterogeneity of the components of B and lack of an appreciable primary structure similarity, the overall structural organization of PT is closely similar to those of other oligomeric toxins having B units (3-7). The B pentamer is disk-shaped, and its subunits fold in six antiparallel  $\beta$ -strands with  $\alpha$ -helices lining the central pore (6). S2 and S3 possess additional domains with a [lectin-like](#) fold, which were suggested to determine the different sugar-binding properties of PT being different from those of [cholera toxin and enterotoxins](#) (6). In fact, PT binds to **glycoproteins**, such as fetuin, having an Asn-linked branched mannose core with terminal N-acetylglucosamine (8).

At variance from cholera toxin, the carboxyl terminus of S1 enters the B oligomer central pore only halfway, and there are fewer protein-protein contacts (4, 6). The similarity with cholera toxin is very high in the [active site](#), with overlapping glutamate and [tyrosine](#) residues essential for activity. His35 occupies a central position in the **NAD-binding** cavity of PT, as is also the case with the active sites of [diphtheria toxin](#) and exotoxin A from *Pseudomonas*. Little information is available on the internalization of PT or the mechanism of membrane translocation of S1 (8). That the structural and functional organization of PT is similar to that of cholera toxin, enterotoxins, and shiga toxins suggests there are further similarities in their mode of cell penetration.

Reduction of the S1 subunit in the presence of [detergents](#) frees an ADP-ribosyltransferase activity specifically directed toward the  $\alpha$  subunit of several trimeric G proteins. At variance from cholera toxin and enterotoxin, PT catalyzes the ADP-ribosylation of a *cysteine* residue present in a Xaa-Cys-Gly-Leu-Xaa motif contained in the carboxyl-terminal part of the  $\alpha$  subunit of several trimeric G proteins, including  $G_i$ ,  $G_o$ , and  $G_{\text{gust}}$ .  $G_s$  has a [tyrosine](#) residue in the position of the cysteine and is not modified by PT. Another difference from cholera toxin is that the S1 enzymic activity is greatly enhanced by ATP and other adenine nucleotides (9). The result of the S1-catalyzed ADP-ribosylation of  $G_i$  and  $G_t$  is different from the cholera toxin or enterotoxin modification of  $G_s$  because the carboxyl-terminal cysteine residue of  $G_i$  and  $G_t$  is centered in the area of interaction between the G protein and plasma membrane receptor. As a consequence, the inhibitory physiological effects normally mediated by these G proteins are not transmitted. The PT oligomer B has biological functions independent of the ADP-ribosyltransferase activity, such as that of inducing the proliferation of T lymphocytes (10).

## Bibliography

1. R. Sekura, J. Moss, and M. Vaughan (1985) *Pertussis Toxin*, Academic Press, New York.
2. M. Domenighini, M. Pizza, and R. Rappuoli (1995) In *Bacterial Toxins and Virulence Factors in Disease* (J. Moss et al., eds.), Marcel Dekker, New York, pp. 59-80.
3. T. Sixma et al. (1993) *J. Mol. Biol.* **230**, 890-918.
4. R. G. Zhang et al. (1995) *J. Mol. Biol.* **251**, 563-573.
5. M. E. Fraser, M. M. Chernaiia, Y. V. Kozlov, and M. N. G. James (1994) *Nature Struct. Biol.* **1**, 59-64.
6. P. E. Stein et al. (1994) *Structure* **2**, 45-57.
7. A. G. Murzin (1993) *EMBO J.* **12**, 861-867.

8. H. R. Kaslow and D. Burns (1992) *FASEB J.* **6**, 2684–2690.
9. J. Moss et al. (1986) *Biochemistry* **25**, 2720–2725.
10. M. Tamura et al. (1983) *J. Biol. Chem.* **258**, 6756–6761.

## Pest Regions

One of several features of [proteins](#) that have been proposed lead to their rapid **protein degradation** is the presence of so-called PEST regions: protein domains that are enriched in **proline, glutamic acid, serine,** and [threonine](#) residues. These regions were first described by Rechsteiner and co-workers in the mid-1980s after searching for structural features of proteins that correlate with short half-lives. They developed a computer-based algorithm locating the presence of at least one P, E, S, or T residue in a [hydrophilic](#) stretch of at least 12 amino-acid residues and flanked by a positively charged residue. Consequently, PEST regions do not represent a specific sequence per se, but they are instead a general feature of the relevant protein domains. Less than 10% of mammalian proteins of known amino acid sequence have such PEST regions, but a high percentage of short-lived proteins do. In some cases, proteins become more stable after removal of their PEST regions, and a few examples have been shown where transplantation of PEST regions into long-lived proteins causes their rate of degradation to be increased.

The mechanism by which PEST sequences may influence degradative rates has remained elusive. Some PEST-containing proteins are degraded by **proteasomes** in a **ubiquitin**-dependent manner, and they require different E2 and E3 enzymes (see [Ubiquitin](#)). On the other hand, PEST regions may simply correlate with rapid hydrolysis because they often contain **consensus sequences** for known protein **kinases**, and **phosphorylation** often signals ubiquitin conjugation and degradation. Also, proline-rich hydrophilic domains are likely to exist as unfolded loops on the surface of [protein structures](#) and to be more susceptible to interactions with **molecular chaperones**, kinases, and other [proteinases](#). Theories involving degradation by calpains and proline endopeptidase have, however, proved false.

### Suggestions for Further Reading

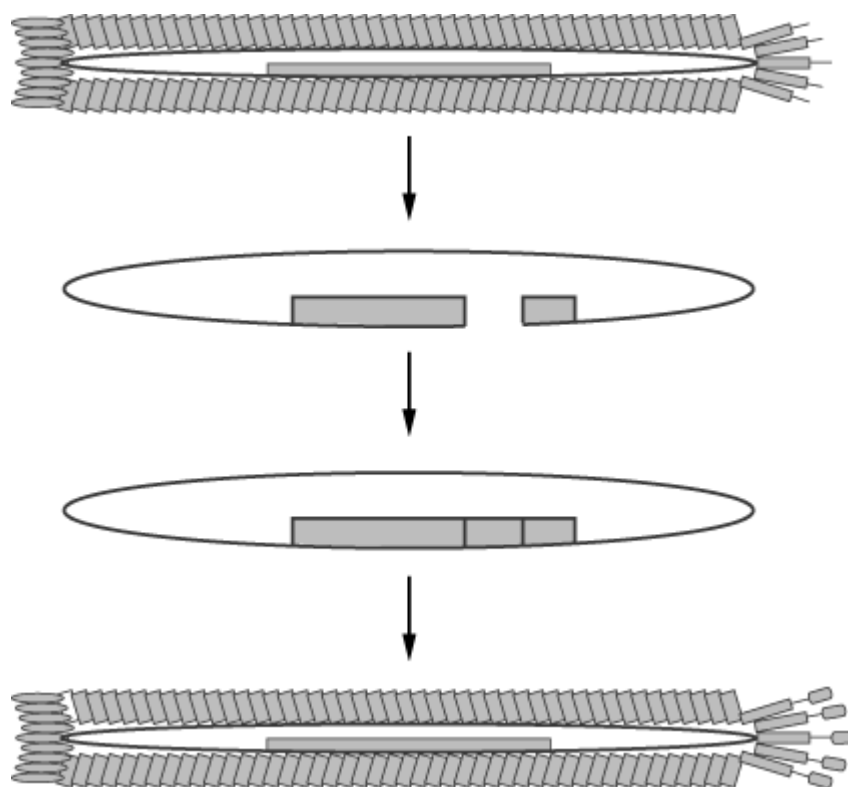
- S. Rogers, R. Wells, and M. Rechsteiner (1986) Amino acid sequences common to rapidly degraded proteins: the PEST hypothesis, *Science* **234**, 364–368.
- M. Rechsteiner and S. W. Rogers (1996) PEST sequences and regulation by proteolysis, *Trends Biochem. Sci.* **21**, 267–271.

## Phage Display Libraries

Phage display libraries are a variant of [expression libraries](#) in which **cloned** DNA inserts are expressed as polypeptide **fusion proteins** with bacteriophage capsid proteins, so that the foreign proteins are “displayed” on the exterior of the modified **virus** particle (Fig. 1) (1-3). Because the virus also contains a DNA [genome](#), a physical linkage is maintained between the displayed peptide

and its encoding DNA sequence. This linkage is the core of the power of phage display. Although the displayed proteins behave in many respects as soluble factors, they remain associated with their DNA code as part of an infectious particle. Thus, even trace amounts of display phages with desired activity can be amplified by infection, and the active peptide sequences can be decoded by DNA sequencing.

**Figure 1.** Schematic illustration of phage display. A double-stranded copy of the phage genome is cleaved enzymatically, and foreign DNA is inserted in-frame into a capsid gene. After introduction of recombinant phage DNA into bacterial cells, novel fusion phage are generated with foreign peptides displayed as capsid fusion proteins on the phage surface.



Several different bacterial viruses have been modified to serve as phage display vectors. By far the most common type of vector is derived from the filamentous Ff bacteriophage (4). The popularity of this virus is due first to the fact that it packages a relatively small single-stranded DNA genome, providing some advantages in [mutagenesis](#) and [DNA sequencing](#), and second to the robust technologies that have been developed for displaying both short peptides and larger proteins. Other viruses that have been adapted for phage display include [lambda phage](#) and **T4 phage**. Future directions in this field will likely include vector construction using human enteric bacteriophages for the purpose of mucosal vaccine development. In addition, certain human viruses could conceivably be modified as viral display vectors for human gene therapy applications.

Filamentous phage display vectors fall into two categories. First are the phage vectors, which accommodate foreign DNA fused to one of several capsid genes within the phage genome. Phage vectors are adequate for peptide display, but larger polypeptide–capsid fusion proteins tend to interfere with phage assembly or infectivity. This interference can be tolerated if the phage can incorporate a preponderance of normal capsid proteins, with only a small proportion of fusion capsid proteins present. Such a situation occurs in the second category of vectors, called “phagemids,” in which the foreign DNA and the capsid gene are both cloned into a special type of plasmid containing

the phage DNA-packaging signal. Foreign DNA is thus cloned as a gene fusion with the phagemid capsid protein. A phagemid is not infectious by itself, but in the presence of a “helper” phage the phagemid DNA can be packaged into an infectious virus particle that displays a mixture of normal and fusion capsid proteins. Using phage and phagemid vectors, a great number of foreign sequences have been expressed as viral capsid fusion proteins, including peptides, **antibodies**, **enzymes**, **DNA-binding proteins**, [RNA-binding proteins](#), and regulatory proteins (reviewed in Ref. 4). Furthermore, any of these proteins can be diversified to form combinatorial phage display libraries, which can then be screened to identify novel proteins with desirable properties.

Phage display screening is a robust technology that typically operates through [affinity selection](#) of phage-encoded molecules that interact with target proteins (1). A phage display library typically expresses billions of different peptides or polypeptides. The active minority of the library can be fractionated from the noninteracting majority by affinity partitioning over a solid-phase target. Phages that bind are then readily eluted and amplified by infection of a suitable bacterial host, and further rounds of selection can be performed. This iterative selection process provides a powerful means to obtain specific peptide aptamers, as well as sequence–activity relationship (SAR) data for any target of interest.

Phage display has proven to be a powerful means of obtaining functional probes to study any target of interest. Thus, peptide aptamers can be used to characterize target function in much the same way that antibodies have been used historically to detect, trap, and inhibit their targets. However, unlike animal-based antibody technologies, phage display can yield peptide or antibody functional probes without requiring immunogenicity or [immunization](#) by the target. Furthermore, the complexity of modern phage display libraries exceeds the genetically encoded complexity of the immune system. Thus, phage display technologies are likely to continue to find increasing applications throughout the biological and medical sciences.

See also [Combinatorial Libraries](#), [Libraries](#), [Combinatorial Synthesis](#), [Affinity Selection](#), [DNA Libraries](#), [Genomic Libraries](#), [cDNA Libraries](#), [Expression Libraries](#), and [Peptide Libraries](#).

#### Bibliography

1. G. P. Smith (1985) *Science* **228**, 1315–1317.
2. S. F. Parmley and G. P. Smith (1988) *Gene* **73**, 305–318.
3. W. Markland, B. L. Roberts, M. J. Saxena, S. K. Guterman, and R. C. Ladner (1991) *Gene* **109**, 13–19.
4. B. K. Kay, J. Winter, and J. McCafferty (1996) *Phage Display of Peptides and Proteins*, Academic Press, San Diego.

#### Suggestions for Further Reading

5. J. D. Watson (1987) *Molecular Biology of the Gene*, 4th ed., Benjamin-Cummings, Menlo Park, CA.
6. J. Sambrook, E. F. Fritsch, and T. Maniatis (1989) *Molecular Cloning: A Laboratory Manual*, 2nd ed., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
7. D. J. Kenan, D. E. Tsai, and J. D. Keene (1994) *Trends Biochem. Sci.* **19**, 57–64.
8. (1996) *Combinatorial Chemistry, Methods in Enzymology*, vol. **267**, Academic Press, San Diego.
9. L. A. Thompson and J. A. Ellman (1996) *Chem. Rev.* **96**, 555–600.
10. B. K. Kay, J. Winter, and J. McCafferty (1996) *Phage Display of Peptides and Proteins*, Academic Press, San Diego.

## Phagocytosis

Phagocytosis is the ability of a cell to internalize large particles, as exemplified by the feeding of phagocytic amoebae, which is one type of [endocytosis](#). Phagocytosis is initiated by the binding of large extracellular particles (eg, bacteria) to specific cell-surface receptors, either directly or indirectly. Indirect recognition involves host-derived opsins, such as [antibody](#) or **complement** fragments, that coat the particles and serve as ligands for some phagocytic receptors. The best-characterized opsins are antibodies that bind to surface **antigens** on infectious microorganisms and form an antibody coat in which the tail region of each antibody molecule (the **Fc region**) is exposed. Binding of the antibody-coated particles to specific Fc receptors on phagocytic cells initiates phagocytosis. Ligands for other receptors, including complement receptors (see [Complement System](#)), [integrins](#), cell-surface [lectins](#), and certain glycosylphosphatidylinositol (**GPI-anchored** proteins [eg, CD14 (1)], also activate phagocytosis.

A common feature of the phagocytic events mediated by these various receptors is that particle internalization depends on the polymerization of [actin](#). Activated cell-surface receptors transmit signals to the cytoplasm that induce actin polymerization (2, 3). During Fc receptor-mediated phagocytosis at least, these activation signals involve the small [GTP-binding proteins](#) Rac/Rho/CDC42 and **tyrosine kinases**. The localized polymerization of actin at the site of particle attachment to the cell pushes plasma membrane around the particle, allowing further opsin-receptor interactions (zippering) and the progressive envelopment of the particle within a membrane invagination tightly apposed to the particle surface. When engulfment is completed, the forming phagocytic vesicle, or phagosome, is sealed and detaches from the cell surface.

In most cases the phagosome fuses with resident compartments of the endocytic pathway ([endosomes](#) and **lysosomes**) to acquire the spectrum of hydrolases that digest the phagosome content. The resulting digestion products are transferred to the cytosol and used to synthesize cellular constituents or are secreted from the cell. In multicellular organisms, phagocytosis is primarily associated with defense against infection by ingesting invading microorganisms and the scavenging of **senescent**, damaged, or **apoptotic** cells (1) and cellular debris. These functions are generally performed by specialized white blood cells, such as macrophages and neutrophils, that act as “professional” phagocytes. These cells are extremely active. In humans, for example, macrophages remove around  $10^{11}$  senescent red blood cells every day. However, most cells are capable of phagocytosis when challenged with appropriate ligands.

Many invasive bacterial, fungal and protozoan pathogens use phagocytosis to enter cells. Some, such as *Listeria* and *Yersinia*, use [cadherins](#) or [integrins](#), respectively, as receptors to initiate endocytic events similar to those described previously. Others, such as *Shigella*, secrete proteins that activate phagocytosis directly. Then the fate of the phagocytic vesicle varies according to the nature of the pathogen internalized. *Listeria*, *Shigella*, and *Trypanosoma cruzi* induce lysis of the phagosome membrane to gain access to the cytoplasm. Other agents (*Legionella*, *Salmonella*, *Mycobacterium avium*) reside and replicate within phagocytic vesicles and produce components that prevent fusion with lysosomes, thus avoiding the hydrolytic consequences of endocytosis, or modify the environment within the phagosome to minimize the effects of the hydrolytic environment (4, 5).

(See also [Endocytosis](#); [Pinocytosis](#); [Macropinosomes](#).)

### Bibliography

1. A. Devitt, O. D. Moffatt, C. Raykundalia, J. D. Capra, D. L. Simmons, and C. D. Gregory (1998) *Nature* **392**, 505–509.

2. S. Greenberg (1995) Trends Cell Biol. **5**, 93–99.
3. K. Ireton and P. Cossart (1998) Curr. Opin Cell Biol. **10**, 276–283.
4. Y. K. Oh, C. Alpuche Aranda, E. Berthiaume, T. Jinks, S. I. Miller, and J. A. Swanson (1996) Infect. Immun. **64**, 3877–3883.
5. S. Sturgill-Koszycki et al. (1994) Science **263**, 678–681.

## Phalloidin

Phalloidin is a cyclic heptapeptide found in the highly poisonous mushroom *Amanita phalloides*. Phalloidin binds with very high affinity to F-[actin](#) and stabilizes these filaments by preventing depolymerization. It contains several uncommon [amino acids](#) and is resistant to all proteolytic enzymes tried thus far. While phalloidin has the ability to be highly toxic to a cell due to the changes that it will induce in actin polymerization, it is now known that the toxicity of *Amanita phalloides* is due mainly to the a- and b-amanitins (potent inhibitors of **RNA polymerase II**), since phalloidin is not absorbed from the gastrointestinal tract ([1](#)). Nevertheless, phalloidin has been an extremely important drug in answering structural, biochemical, and cell biological questions about actin.

Phalloidin does not bind to monomeric actin (G-actin), but binds with a 1:1 stoichiometry to protomers within the F-actin filament. This is the opposite mode of binding of DNase I (see [DNase 1 Sensitivity](#)), which binds with high affinity to G-actin but not to F-actin. As a result of this selective binding, both may be used to assay the concentrations of G- vs F-actin within a fixed cell. The phalloidin-stabilized filament resists depolymerization that would otherwise be induced by a reduction of the actin concentration, by chaotropic agents, **cytochalasins**, heat denaturation, and proteolysis ([2](#)). The critical concentration is the concentration of the monomeric species of actin at which it is in equilibrium with **polymers**, and this is about  $10^{-6}$  M for muscle actin. Phalloidin reduces this critical concentration by a factor of about 25. These properties of phalloidin have made it a very useful tool for labeling actin filaments within a cell (using fluorescent derivatives of phalloidin) and for stabilizing actin filaments used in structural and biophysical studies.

## Bibliography

1. A. Jaeger, F. Jehl, F. Flesch, P. Sauder, and J. Kopferschmitt (1993) J. Toxicol. Clin. Toxicol. **31**, 63–80.
2. T. Wieland and H. Faulstich (1978) CRC Crit. Rev. Biochem. **5**, 185–260.

## Phase Problem

To determine the structure of a molecule by [X-ray crystallography](#), after collecting and processing the data, the electron density  $r(x, y, z)$  in the crystal [unit cell](#) can be calculated at every position  $x, y, z$ :

$$\rho(xyz) = \frac{1}{V} \sum_h \sum_k \sum_\ell |F(hk\ell)| \times \exp[-2\pi i(hx + ky + \ell z) + i\alpha(hk\ell)] \quad (1)$$

To use this equation requires that the values of both  $F(hk\ell)$  and  $\alpha(hk\ell)$  for each reflection  $h k \ell$  are known. The  $F(hk\ell)$ 's are easily found because (apart from corrections) they are proportional to the square root of the measured scattered intensities  $I(hk\ell)$ . The phase problem is that the phase angles  $\alpha(hk\ell)$  are not provided directly by the experimental data. This fundamental problem greatly delayed the use of X-ray crystallography initially, but it is hardly a problem anymore. For small molecules, it has long been solved by [direct methods](#). In the crystallography of macromolecules, several methods exist to determine the phase angles of the reflections:

1. [isomorphous replacement](#)
2. [molecular replacement](#)
3. **MAD** (multiple-wavelength anomalous dispersion)

Classical direct methods are applicable to macromolecular crystals, only if the protein is small and the X-ray data can be collected to high resolution (1).

#### Bibliography

1. G. M. Sheldrick et al. (1993) *Acta Crystallogr.* **D49**, 18–23.

#### Suggestion for Further Reading

2. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

### Phase Variation, Pili

Certain enterotoxigenic strains of *Escherichia coli* produce pilus-like appendages called *colonization factor antigens* (CFAs) that are longer than common type I [pili](#). Production of type I pili or CFA is mutually exclusive in these bacteria and thus represents a specialized example of bacterial phase variation. That displayed by type IV pili of *Neisseria gonorrhoeae*, is the best-characterized system. The N-terminal domain of *N. gonorrhoeae* pilin is **hydrophobic** and highly conserved among different variants. Other regions of the pilin molecule, however, known as *minicassettes*, vary widely in amino acid sequence. There is also wide variation at the DNA level between the various copies of the **genes** that encode a particular pilin protein. In *N. gonorrhoeae* the expression of the *pilE* gene is controlled at [transcription](#) by a two-component regulatory system; this comprises the transcriptional activator PilA and the sensor component PilB, which **phosphorylates** PilA. If  $\text{Pil}^+$  and  $\text{Pil}^-$  respectively represent the phenotypes when functional pili are assembled or not on the bacterial cell surface, then phase variation represents the fluctuation that occurs between  $\text{Pil}^+$  and  $\text{Pil}^-$ . There are several types of phase variation (1).

1. The first arises because the expressed protein has two possible sites ( $\text{P}^+$  and  $\text{P}^s$ ) at which subsequent processing can take place. When cleaved **proteolytically** at  $\text{P}^+$ , the pilin protein assembles into normal pili ( $\text{Pil}^+$ ). In contrast, when it is cut at  $\text{P}^s$  it no longer contains the hydrophobic region required for assembly. The protein is secreted and lost from the cell, so there is no pili formation ( $\text{Pil}^-$ ).

2. The second type of phase variation is similar to the first in that it again gives rise to a pilin protein that is aberrant and not capable of forming pili *in vivo*. Although the pilin gene is present in multiple copies in the [chromosome](#), generally only one of these has a **promoter** and is expressed. This gene is designated *pilE*, where *E* represents expression. Those gene copies that are not expressed are termed *pilS*, where *S* represents silent. Both the *pilE* and *pilS* genes contain a number of internal repeated DNA sequences. Homologous [recombination](#), which involves exchanges of the interior regions of the genes, can lead to unequal exchange of DNA because of the repetitive nature of the sequences. As a result, the PilE protein produced may be larger in size than normal. As before, this cannot be processed in such a way that pili formation is possible, hence the **phenotype** is Pil<sup>-</sup>.

3. The third mechanism of phase variation concerns the *pilC* gene. The expressed protein PilC, which has been detected by [immunoelectron microscopy](#) at the tip of the pili, is an important protein in the assembly and functioning of intact mature pili. Towards the 5' end of the *pilC* coding sequence there is a long run of G bases. Slipped-strand mispairing during synthesis alters the number of G bases and generates a [frameshift mutation](#). This again gives rise to Pil<sup>-</sup> phenotype. Note, however, that pili of normal appearance are produced in this case and that another phase-variable pilus-associated protein replaces PilC in its assembly function, but it lacks the ability to bind to epithelial cells (2). Phase variation probably evolved to allow the bacterium to escape the [immune response](#) of its host and to allow it to adapt to new cell types (1).

#### Bibliography

1. A. A. Salyers and D. D. Whitt (1994) *Bacterial Pathogenesis: A Molecular Approach*, ASM Press, Washington, D.C.
2. T. Rudel, D. Facius, R. Barten, I. Scheuerpflug, E. Nonnenmacher, and T. F. Meyer (1995) Role of pili and the phase-variable PilC protein in natural competence for transformation of *Neisseria gonorrhoeae*. *Proc. Natl. Acad. Sci. USA* **92**, 7986–7990.

#### Phenotypic Lag

After a cell sustains a [mutation](#), there may be a delay of one or more generations before the mutant **phenotype** is expressed. This delay has two causes. First, a recessive mutation may be **complemented** by the wild-type **allele** on another [chromosome](#). In mitotically dividing **diploid** cells, at least a one-generation delay occurs before the phenotype is expressed. Under most laboratory conditions, [haploid](#) organisms also have more than one chromosome per cell and show a lag of one or two generations. Secondly, when a mutation that gives rise to a mutant protein first appears, the cell may have many copies of the **wild-type** protein. Expression of the mutant phenotype may be delayed until these preexisting proteins are diluted or degraded. A classical example from *Escherichia coli* genetics is resistance to **bacteriophage** T1. Resistance is conferred by mutations that alter the phage **receptor** on the cell surface, but several generations are needed before the preexisting sensitive receptors are diluted. Even dominant phenotypes may fail to be expressed until the cellular concentration of the wild-type protein falls. Phenotypic lag complicates the determination of mutation rates from [fluctuation tests](#) (1-4).

#### Bibliography

1. S. E. Luria and M. Delbrück (1943) *Genetics* **28**, 491–511.
2. P. Armitage (1952) *J. R. Stat. Soc.* **B 14**, 1–40.
3. A. L. Koch (1982) *Mutat. Res.* **95**, 129–143.

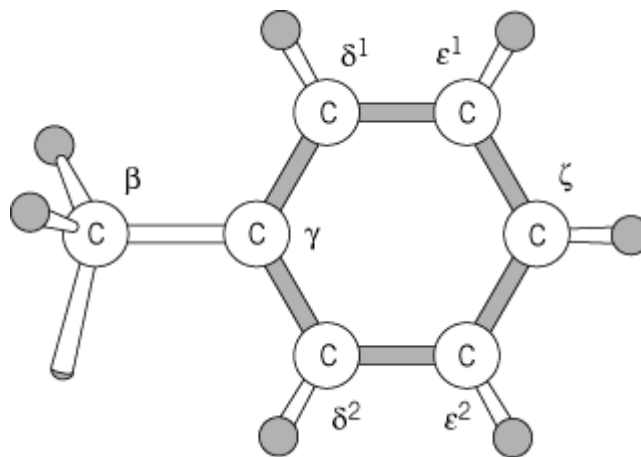


4. F. M. Stewart, D. M. Gordon, and B. R. Levin (1990) *Genetics* **124**, 175–185.

## Phenylalanine (Phe, F)

The [amino acid](#) phenylalanine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—UUU and UUC—and represents approximately 3.9% of the residues of the proteins that have been characterized. The phenylalanyl residue incorporated has a mass of 147.18 Da, a **van der Waals volume** of  $135 \text{ \AA}^3$ , and an [accessible surface](#) area of  $218 \text{ \AA}^2$ . Phe residues are changed infrequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [tyrosine](#) and [leucine](#) residues.

Phe is one of the three amino acids with an aromatic ring in its side chain:



which has **absorbance** and **fluorescence** spectral properties that are very useful in quantifying proteins and analyzing their structural properties. The other two—[tyrosine](#) and [tryptophan](#)—predominate, however, so Phe residues are useful only in their absence. The aromatic ring of Phe residues is comparable chemically to that of benzene or toluene; consequently, it is largely [nonpolar](#) (except for the tendency of the electrons of the peripheral H atoms to be drawn into the aromatic ring) and is chemically reactive only under extreme conditions that are not applicable to proteins.

Tyr residues occur primarily within the interiors of folded [protein structures](#), and approximately 50% are totally buried; normally, they are still rotating rapidly by  $180^\circ$  rotations about the  $C_b-C_g$  single bond. They tend to occur primarily in **alpha-helices**, although also in **beta-sheets**, and they favor the  $\alpha$ -helical conformation in model peptides.

### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties* 2nd ed., W. H. Freeman, New York.

## Phenylketonuria

The metabolic disorder phenylketonuria (PKU) was first described in 1934 by A. Følling, who was investigating the excretion of ketone bodies in the urine of two mentally retarded children. The reagent in use, acidified  $\text{FeCl}_3$ , typically turns a purple-red color with ketones. The urine of these children turned bright blue-green instead. Examination of additional mentally retarded children uncovered eight more who were similar, both clinically and in their urine reactions. Følling subsequently demonstrated that the responsible urinary constituent was phenylpyruvic acid, thus the name phenylketonuria. In addition to mental retardation, affected children share a number of clinical features, including irritability, seizures, eczema, and a “mousy odor.” They tend to be more lightly pigmented than others in the same family.

The rarity of this disorder in the general population and the occurrence of more than one affected person in a sibship suggested recessive inheritance. The equal occurrence of males and females and the fact that the parents and other relatives were not affected was consistent with autosomal recessive inheritance.

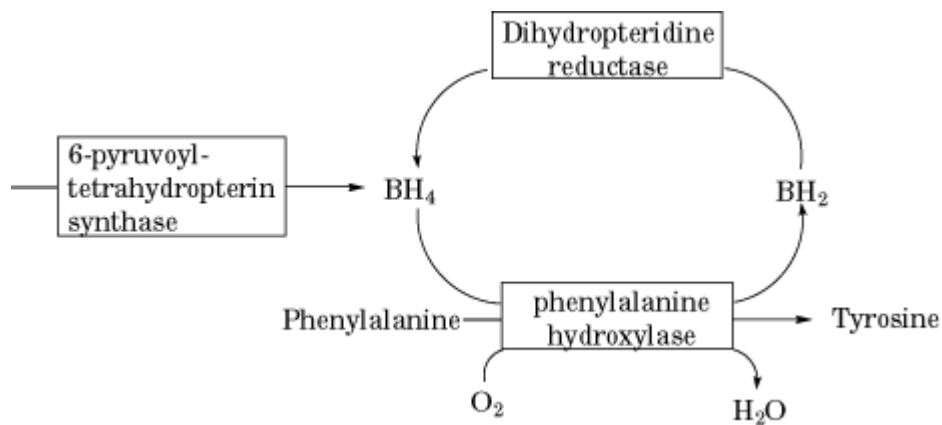
## 1. The Biochemistry of PKU

Følling suggested that patients who have PKU have a block in metabolism of phenylalanine, and in 1947 G.A. Jervis proposed that the metabolic error is in the conversion of phenylalanine to tyrosine. He subsequently proved this hypothesis by showing that liver from a PKU patient cannot make the conversion. Thus, PKU was among the earliest examples of a *metabolic error*, a term proposed by A.E. Garrod in 1909, and it has served as a model for other metabolic defects.

Phenylalanine is efficiently reabsorbed in the kidneys, leading to maintenance of very high blood levels. In normal persons, the blood phenylalanine is ca.  $60 \mu\text{M}$ . The level in PKU patients can be 40 times higher. This leads to formation and excretion of otherwise rare metabolites of phenylalanine, such as phenylpyruvic acid, phenylacetic acid, and *o*-hydroxyphenylacetic acid. These metabolites very likely contribute to the clinical features of PKU by interfering with various metabolic and developmental processes.

Phenylalanine hydroxylase (PAH; L-phenylalanine-4-monooxygenase; EC 1.14.16.1) is the enzyme that is ordinarily nonfunctional in PKU. The predicted polypeptide molecular weight is 51,862 daltons, but the enzyme is usually a homotetramer. It catalyzes the oxidation of L-phenylalanine, using tetrahydrobiopterin as a cofactor (Fig. 1). Production of PAH occurs almost entirely in the liver. A large number of PAH alleles are known. Some produce completely inactive PAH; others produce PAH that has reduced or normal activity. Persons whose allelic combinations produce PAH activity  $<1\%$  of normal have classical PKU. If there is some activity, a mild form of the disorder may result. Heterozygotes, with one normal and one inactive allele, are phenotypically normal even though the total enzymic activity is reduced.

**Figure 1.** Three enzymes that are important in the conversion of phenylalanine to tyrosine. Additional enzymes and cofactors are required to form  $\text{BH}_4$ . Reduced activity of phenylalanine hydroxylase (PAH) can result from structural modifications of PAH or from deficiency of  $\text{BH}_4$ .



## 2. Treatment of PKU

Humans cannot make phenylalanine or tyrosine and must depend on diet for these two essential amino acids. The phenylalanine content of foods is in excess of body requirements for protein synthesis, and the surplus is metabolized via tyrosine. This suggested the possibility of reducing the accumulation of phenylalanine by reducing the dietary intake. Such efforts, initiated in 1953 by H. Bickel, ultimately proved successful. The low-phenylalanine diet must be introduced soon after birth of an affected infant in order to avoid permanent damage to the developing central nervous system.

With the development of an effective method of treatment, the urgency of neonatal diagnosis became apparent. In the early 1960s, R. Guthrie developed a microbial test using whole blood from a heel stick. This rapid and inexpensive test made large scale screening possible, and required screening of newborns for elevated blood phenylalanine has become widespread in western countries. Once the central nervous system has matured fully, persons with PKU are somewhat less sensitive to increased blood levels of phenylalanine. This includes not only regulation of protein intake but also avoidance of the artificial sweetener aspartame, which is a dipeptide derivative of aspartic acid and phenylalanine. However, long-term studies have shown benefits of continued treatment of adults.

Although dietary regulation has been highly successful, it is not a perfect treatment. Treated patients show mildly depressed IQ compared to non-PKU controls. One suggested cause is increased dysmyelination with age in the presence of elevated blood phenylalanine. Also, tyrosine is lower in PKU patients, treated or untreated, and this may lead to deficiency of neurotransmitters.

In view of the single gene and its single protein product that are defective in PKU, this disorder would seem a good candidate for gene therapy. It has not been considered a high priority, however, because of the effectiveness of dietary restriction of phenylalanine. When procedures for gene therapy have been sufficiently developed and tested, it is likely that they will be applied to PKU. Introduction of functional PAH genes during the neonatal period should cure the metabolic defect. This, in effect, was the case with one young boy who received a liver transplant for reasons unrelated to his PKU status.

The successful treatment of PKU by diet has led to the problem of maternal PKU, which occurs in successfully treated women who are homozygous for PKU and whose offspring therefore are heterozygous. The developing fetus, although genotypically normal (but heterozygous), depends on the mother to regulate blood phenylalanine. The fetus is also highly susceptible to moderate increases in blood phenylalanine. Prior to conception and during pregnancy, a PKU mother must regulate her blood phenylalanine levels rigorously to avoid damage to the fetus. The recommended maximal level is 360  $\mu$ M.

One consequence of introducing screening was the recognition that some children have levels of phenylalanine that are above normal but that do not reach the levels that lead to pathology and that do not require dietary management. These cases of “hyperphenylalaninemia,” defined as plasma levels  $>120 \mu\text{M}$ , are often due to homozygosity for PAH alleles that have low but nonzero enzyme activity. In other instances, the hyperphenylalaninemia is transient in neonates and disappears with time.

### 3. Other Defects That Interfere with Metabolism of Phenylalanine

In rare instances, high levels of blood phenylalanine are due to deficiencies of PAH activity caused by other loci. The cofactor 5,6,7,8-tetrahydrobiopterin ( $\text{BH}_4$ ) is produced from 7,8-dihydrobiopterin ( $\text{BH}_2$ ) by the enzyme dihydropteridine reductase (DHPR). The gene that codes for DHPR is on chromosome 4 at 4p15.31. In one form of PKU, loss of activity of this gene leads to absence of  $\text{BH}_4$  and therefore to inactivity of PAH (Fig. 1). Because  $\text{BH}_4$  is a cofactor in several other reactions that are important to brain function, dietary management of blood phenylalanine levels does not restore a normal phenotype. Such cases occur in 1 to 2 births per million.

A second rare variant of PKU is due to deficiency of 6-pyruvoyl tetrahydropterin synthase (Fig. 1), one of several enzymes necessary for  $\text{BH}_2$  synthesis. This enzyme is coded by a gene on chromosome 11 at 11q22.3-q23.3. As in the case of DHPR deficiency, the neurological effects are not prevented by management of diet.

### 4. Molecular Genetics of PAH

The gene that codes for PAH is on chromosome 12 at 12q24.1. It is some 90 kb long with 13 exons. The messenger RNA is 2,677 nt long, consisting of 472 nt in the 5' untranslated region, 1,356 nt (452 codons) in the coding sequence, and 849 nt in the 3' untranslated region. Exons 1 to 12 are relatively small, varying in size from 57 nt to 197 nt. Exon 13 has 890 nt, of which 41 are coding. Introns vary in size from 1.0 kb to 23.5 kb.

Some 200 distinct mutations associated with PKU have been identified in the PAH gene. There are presumed to be many more that are associated with unrecognized hyperphenylalaninemia or that are normal variants. Some normal variants are known, and in a few instances they are polymorphic. Of the mutations that cause PKU, approximately three-fourths are nucleotide substitutions that lead to amino acid substitutions or premature chain termination. Another 25 or more are nucleotide substitutions that interfere with normal splicing of the RNA transcript. A number of small gene deletions have also been identified. A high frequency of mutations is found in exon 7, thought to be due to ascertainment bias rather than mutational hotspots. Almost any change in amino sequence in this region of the PAH molecule is likely to interfere with function.

A relatively small number of these mutations account for the majority of the cases of PKU, and particular mutations may be prevalent only in a limited population. Because several different mutations may occur with relatively high frequencies in a population, most so-called homozygotes are in fact compound heterozygotes, that is, the two alleles involve different mutations. Both, however, are nonfunctional. Comparison of mutations in European and Asian populations shows quite different spectra, indicating that most of the extant mutations arose after separation of the major races.

### 5. Population Genetics of PKU

The frequency of patients with PKU varies considerably among populations. In Norway, where PKU was first discovered, the frequency is ca. 1 in 14,000 births. Among Northern Europeans in general, the frequency is ca. 1 in 10,000 births, with a higher incidence of 1/4,500 in Ireland. Among U.S.

whites, the frequency is 1/8,000. The frequency among African-Americans is ca. 1/50,000. The frequency is very low in Ashkenazi Jews and Asians. A frequency of homozygotes of 1/10,000 corresponds to a gene frequency of 0.01 and a frequency of heterozygotes of 2%.

Investigation of PKU has been aided by the discovery of eight polymorphic DNA variants within the PKU gene. These were originally recognized as restriction fragment length polymorphisms, although now PCR is used for some sites. One of the sites involves different numbers of short tandem repeats, providing several alleles. Altogether, there are about 800 possible combinations (*haplotypes*) of “alleles” at these sites, of which only about 75 have been observed. Within a population, many fewer are observed, and usually only a few reach substantial frequencies.

The study of mutant-allele—haplotype associations in a variety of populations has provided insight into the origin and spread of the genes in historic and prehistoric periods. For example, the mutation R261Q is widespread in Europe but is always associated with haplotype 1. In general, haplotype 1 has a normal PAH allele. The interpretation is that this particular mutation arose on a haplotype 1 chromosome, and all R261Q mutations are descended from that single mutational event. In contrast, the R408W mutation occurs in northwest Europe on a haplotype 1 background and in eastern Europe on a haplotype 2 background. This suggests that two independent mutations occurred, one in a Celtic population and the other in a Slavic population.

#### Additional Reading

Scriver C.R., Kaufman S., Eisensmith R.C., and Woo S.L.C., The hyperphenylalaninemias, in C.R. Scriver, A.L. Beaudet, W.S. Sly, and D. Valle, eds., *The Metabolic and Molecular Bases of Inherited Disease*, 7th ed., vol I McGraw-Hill, New York, 1995, pp. 1015–1075.

*On-Line Mendelian Inheritance in Man*: <http://www3.ncbi.nlm.nih.gov/htbin-post/Omim/dispmim?261600>

*PAH locus database* : <http://www.mcgill.ca/pahdb/>

## Phosphate Buffers

Phosphate [buffers](#) are commonly used in biological work. Phosphoric acid,  $\text{O}=\text{P}(\text{-OH})_3$ , has  $\text{p}K_a$  values of about 2, 7, and 12. All three hydroxy groups are identical and, indeed, all four oxygen atoms are identical on the NMR time scale as a result of rapid hydrations and dehydrations. Thus, 2 is the  $\text{p}K_a$  for the first  $\text{H}^+$  to dissociate, even though it is coming equally from each of three groups. The next  $\text{p}K_a$  is raised to 7 because it is much harder for  $\text{H}^+$  to dissociate from the charged  $\text{H}_2\text{PO}_4^-$  ion than from the neutral  $\text{H}_3\text{PO}_4$  molecule and, similarly, for the third  $\text{p}K_a$ . Hence, the Henderson–Hasselbalch equation for orthophosphate is

$$\begin{aligned}\text{pH} &= 2 + \log \frac{[\text{H}_2\text{PO}_4^-]}{[\text{H}_3\text{PO}_4]} = 7 + \log \frac{[\text{HPO}_4^{2-}]}{[\text{H}_2\text{PO}_4^-]} \\ &= 12 + \log \frac{[\text{PO}_4^{3-}]}{[\text{HPO}_4^{2-}]}\end{aligned}\tag{1}$$

The large separation of the  $\text{p}K_a$  values means that, for a buffer of  $\text{H}_2\text{PO}_4^-$  and  $\text{HPO}_4^{2-}$ , the

concentrations of  $\text{H}_3\text{PO}_4$  and  $\text{PO}_4^{3-}$  are negligible, and the  $p_a$  values of 2 and 12 become irrelevant.

An advantage of phosphate as a buffer is that salts such as  $\text{NaH}_2\text{PO}_4$  and  $\text{Na}_2\text{HPO}_4$  are solid and both can be weighed out in the appropriate quantities; hence, the reproducible preparation of phosphate buffers need not involve adjustments of pH. Care must be taken, however, that the salt used, hydrate or anhydrous, does not lose or gain water on storage. Another advantage is that phosphate is transparent in the UV, down to wavelengths at which proteins may be estimated by their [peptide bond](#) absorbance, where carboxylate ions also absorb.

The main disadvantage is that with a doubly charged buffering species, the  $pK_a$  of 7.2 is three times more sensitive to changes in ionic strength than a buffer with only singly and zero charged components. Another disadvantage is that phosphate can support algal or fungal growth as a nutrient. Phosphate will also sequester and precipitate many cations, especially  $\text{Ca}^{2+}$  (see [Calcium Signaling](#)).

### 1. Diphosphate (pyrophosphate)

A variant of phosphate buffer is diphosphate (pyrophosphate). The acid,  $\text{HO}-\text{P}(=\text{O})(\text{OH})-\text{O}-\text{P}(=\text{O})(-\text{OH})-\text{OH}$ , has two strongly acidic groups, with  $pK_a$  values near 2, as with orthophosphate. The  $-\text{P}(=\text{O})(-\text{OH})-\text{O}^-$  groups generated are electron-withdrawing, and so the next  $-\text{OH}$  group has its  $pK_a$  lowered from the 7.2 of orthophosphate to 6.2 (at an ionic strength of 0.1 M; with its more strongly charged ions, the dissociation constants of diphosphate are even more sensitive than orthophosphate to changes in ionic strength). The  $-\text{P}(=\text{O})(-\text{O}^-)_2$  group, however, is slightly electron-donating, so that the final  $pK_a$  is raised to 8.4. Hence, the Henderson–Hasselbalch equation for pyrophosphate in the neutral pH region is

$$\text{pH} = 6.2 + \log \frac{[\text{HP}_2\text{O}_7^{3-}]}{[\text{H}_2\text{P}_2\text{O}_7^{2-}]} = 8.4 + \log \frac{[\text{P}_2\text{O}_7^{4-}]}{[\text{HP}_2\text{O}_7^{3-}]} \quad (2)$$

Diphosphate can be hydrolyzed to orthophosphate, and the equilibrium strongly favors orthophosphate, but diphosphate is stable for long periods of time in neutral and alkaline conditions, even on heating.

## Phosphatidylinositol

Phosphatidylinositol, generally abbreviated as PI or PtdIns, is a glycerophospholipid that has *myo*-inositol (one of nine stereoisomers of hexahydroxycyclohexane) as the polar head group. It is found in the [membranes](#) of eukaryotes and some eubacteria. Because of the anionic phosphodiester group, it contributes to the negative charge of the phospholipid bilayer. The free hydroxyl groups of PI can be phosphorylated or glycosylated to produce a wide variety of inositol phospholipids with specific functions. Two distinct types of inositol phospholipid have particularly important functions in eukaryotes that are described in separate entries:

1. Phosphatidylinositol phosphates are key components of several intracellular signaling pathways, either as the intact molecules or after they are cleaved by [phospholipase C](#) to **1,2-diacylglycerol** and an **inositol phosphate**.

2. Glycosylphosphatidylinositols (or GPIs) can be covalently attached to the C-terminus of many proteins and are responsible for anchoring them to the membrane (see [Membrane Anchors](#) and [GPI Anchor](#)). GPIs may also occur in a “free form” without attached protein. Although the role of these free GPIs in mammalian cells is uncertain, they are an essential part of the protective coat used by many protozoal parasites.

## Phosphatidylinositol Kinases

Phosphoinositides serve important functions in controlling the regulation of growth and differentiation in cells, as well as serving as substrates for the generation of inositol trisphosphate,  $IP_3$  (see [Inositol Lipids and Phosphates](#)). In response to a variety of [growth factors](#) and **cytokines**, phosphatidylinositol (PI) is phosphorylated sequentially, to form the polyphosphoinositide PI-3',4',5'- $P_3$  ( $PIP_3$ ). These phosphorylations are catalyzed by the enzyme PI 3'-kinase. The lipid products of this enzyme are not substrates for [phospholipases](#), but instead have direct regulatory properties themselves.

PI 3'-kinase has a number of **domains** for [protein-protein interactions](#) that contribute to its regulation and localization. The enzyme consists of a catalytic (p110) and regulatory (p85) subunit. p85 has two **SH2 domains** and one SH3 domain, which play a crucial role in regulation of the activity and targeting of the enzyme. Upon growth-factor-receptor activation, p85 is recruited to receptors phosphorylated on tyrosine residues, or their substrates. The occupancy of these SH2 domains in p85 causes activation of the enzyme, while also serving to target it to the plasma membrane, or in some cases to intracellular membranes, where its lipid products can bind to proteins (1).

The generation of 3' phosphoinositides appears to play an important role in [signal transduction](#). Activation of PI 3'-kinase is necessary for the full expression of a number of [tyrosine kinase receptors](#), including those for [insulin](#), **nerve growth factor**, and others (1). Activation of the enzyme has been shown to prevent [apoptosis](#) in a number of cell types. Moreover, a [retrovirus](#) encoding PI 3'-kinase was found to induce hemangiosarcomas in chickens, and mutations in the enzyme were associated with increased life span in [Caenorhabditis elegans](#). Many advances in identifying the role of PI 3'-kinase have come from the use of pharmacological inhibitors, such as wortmannin (1).

The targets of the lipid products of PI 3'-kinase have received a great deal of attention. PI 3'-kinase appears to play a universal role in membrane trafficking, perhaps explaining its wide importance in signal transduction. [Cell adhesion molecules](#) can also cause its activation, linking these events to changes in the actin [cytoskeleton](#). Although there have been numerous reports on the essential role that PI 3'-kinase plays in most aspects of signal transduction, most attention has focused on downstream changes in **phosphorylation**. There have now been a number of protein kinases that are activated after PI 3'-kinase activation. One group belongs to the calcium-independent [protein kinase C](#) family members, especially  $\alpha$  and  $\eta$ . Additionally,  $PIP_3$  can induce a phosphorylation cascade.

This involves the activation of the **serine/threonine kinases** PDK and Akt, which is also known as protein kinase B (PKB). It is now thought that  $PIP_3$  can bind to the **pleckstrin homology** (PH) domain of these kinases, effectively recruiting them to the plasma membrane or intracellular membranes for activation. While PDK appears to be controlled mainly by targeting, it can phosphorylate and activate Akt, which in turn can phosphorylate a variety of intracellular substrates,

including [transcription factors](#), other kinases, proteins controlling apoptosis, and others (1).

## Bibliography

1. A. Toker and L. C. Cantley (1997) *Nature* **387**, 673–676.

## Phosphofructokinase

Phosphofructokinase, or PFK (EC 2.7.1.11), is one of the best characterized allosteric enzymes. It catalyzes the step of glycolysis in which fructose-6-phosphate (FRU-6P) is phosphorylated to fructose-1,6-bisphosphate (FRU-1,6P<sub>2</sub>) using either ATP or pyrophosphate (PP<sub>i</sub>) as a phosphate donor. This enzyme is present in most types of eukaryotes and eubacteria but occurs rarely in archaeobacteria. PFKs from animals, yeast, and most bacteria use ATP as a phosphate donor. Those that use PP<sub>i</sub> are found in plants, protozoa, and some bacteria, often in addition to ATP-dependent enzymes. The amino acid sequences of most of the ATP-dependent PFKs are homologous, indicating that they form an evolutionary gene family. This family includes both “small” bacterial enzymes and “large” enzymes from mammals and yeast. It is a good model of enzyme evolution and of the appearance of “new” regulatory functions. PFK in most organisms is a highly regulated enzyme that has several allosteric effectors, and it controls (at least in part) the rate of glycolysis. The PFK from *Escherichia coli* has become a paradigm for the concerted allosteric model (1). X-ray crystallography, site-directed mutagenesis, and enzyme kinetic studies on bacterial ATP-dependent PFKs have also explained their catalytic and regulatory properties in some detail.

### 1. Isoenzymes

In humans (and probably in other mammals), there are three homologous types of PFK polypeptide chains, the muscle M, liver L, and platelet P isoenzymes. They are encoded by different genes, located on chromosomes 21, 1, and 10, respectively. All have several introns. The few cases of alternative splicing reported always lead to deletion of one exon and formation of a nonfunctional enzyme. Even though the M, L, and P type proteins are different, they share enough similarities to assemble into active, multiple, heterologous oligomers, usually tetramers. The many different PFK isoforms generated in this way have somewhat different catalytic and regulatory properties, such as saturation by ATP and Fru-6P or sensitivity to allosteric effectors (2). Some species, such as *E. coli*, also have other minor proteins with PFK activity, but have unrelated amino acid sequences indicating that they do not belong to the main family of PFKs.

### 2. Metabolic Role

The reaction catalyzed by PFK, using either ATP or PP<sub>i</sub>, is irreversible under physiological conditions. Its complex regulation by adenine nucleotides and by various metabolites indicates that PFK controls the production of cellular energy and the carbon flux between glucose and pyruvate. In several instances, the activity of PFK is coupled to that of pyruvate kinase because Fru-6P, the substrate of PFK, is an allosteric activator of one form of pyruvate kinase. This coupling is stronger in bacteria, where phospho *enol*pyruvate (PEP), the substrate of pyruvate kinase, is the allosteric inhibitor of PFK. In some microorganisms (Gram-positive bacteria and mycoplasmas, but not *E. coli*), the genetic expression of these two enzymes is coordinated. They catalyze reactions seven steps apart in glycolysis. In these cases, the genes that code for PFK and pyruvate kinase are adjacent on the chromosome and constitute an operon under transcriptional control by the same promoter (3).

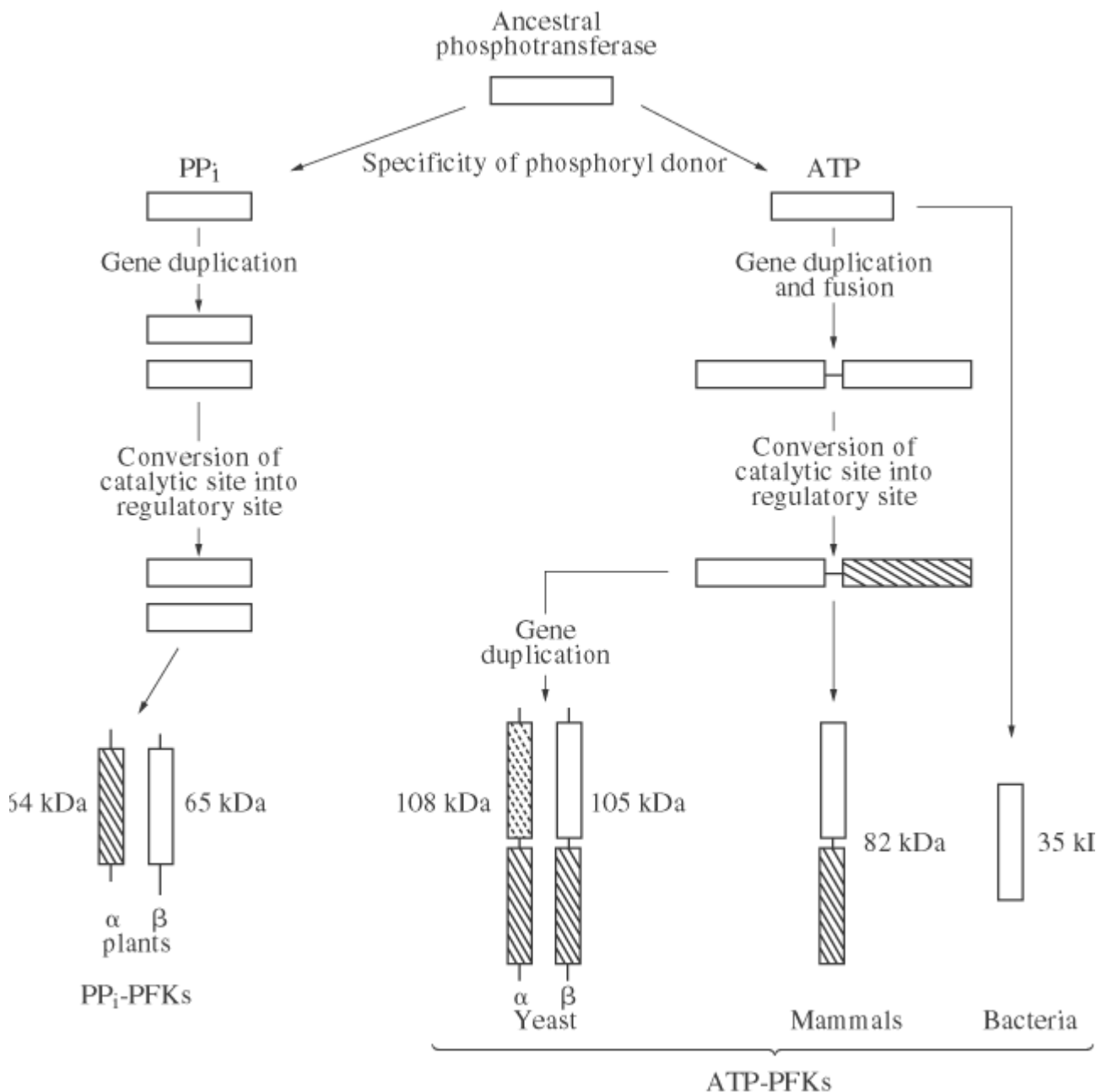


The enzymatic (and sometimes genetic) coupling between two enzymes that each catalyze irreversible phosphoryl transfer steps is probably crucial for controlling the glycolytic flux, especially in anaerobes that use glycolytic fermentation as their main energy-producing pathway.

### 3. Primary Structures and Evolution

Most of the ATP-dependent PFKs, whether from animals, yeasts, or bacteria, have homologous sequences and belong to an evolutionary gene family. The simplest members of this family are bacterial PFKs whose polypeptide chains are about 35 kDa. Mammalian PFKs are approximately twice as large, about 80 kDa, and have an internal sequence repeat. Their N- and C-terminal halves are homologous and both resemble bacterial PFKs. This suggests that they have been produced by tandem gene duplication and fusion (4). All of the important active site residues are conserved in the N-terminal half, indicating a catalytic function, whereas replacement of several crucial residues, notably that of Asp127 by Ser, suggests that the C-terminal half is inactive. Yeast PFK is composed of two types of related chains, a of 108 kDa and b of 105 kDa. Each has a segment homologous to mammalian PFKs and unrelated N- and C-terminal extensions. These chains have formed by duplication of a mammalian-like gene followed by independent acquisition of terminal extensions (Fig. 1). The distribution of residues crucial for catalysis, plus various kinetic studies, suggest that the active site is carried by the N-terminal half of the b-chain and that the a-chain has only a regulatory role. Homologous residues are referred to here by their numbers in the *E. coli* enzyme.

**Figure 1.** Evolution of PFKs suggested by the comparison of their primary sequences. The unrelated segments are represented by thin lines. Among the homologous segments, those that have retained a catalytic function are open, and those that have lost their catalytic abilities and evolved to acquire a regulatory function are crossed. The special weak crossing on the N-terminal segment of the yeast a chain reflects that the evidence in favor of a pure regulatory role is weaker than for the other eukaryotic PFKs.



PP<sub>i</sub>-dependent PFKs are made of two homologous chains, each of which has a segment distantly related to bacterial PFKs. This suggests that all PFKs might derive from the same ancestral protein (Fig. 1). The conservation of critical residues indicates that the catalytic site is on the b-chain and that the a-chain is regulatory.

#### 4. Quaternary Structure

Bacterial and mammalian PFKs are active as tetramers and are inactivated upon dissociation into dimers (see [Quaternary Structure](#)). The allosteric inhibition of the PFK from *Thermus thermoaquaticus* by PEP involves such a dissociation. The crystal structures of bacterial PFKs explain why an oligomeric state is required for activity. The Fru-6P binding site is at the interface between two subunits. The fructose moiety interacts with one subunit and the 6-phosphate group with the other (5). Under some conditions, tetrameric mammalian PFKs associate further into octamers and even more complex structures. They are active, but the role of this aggregation is not

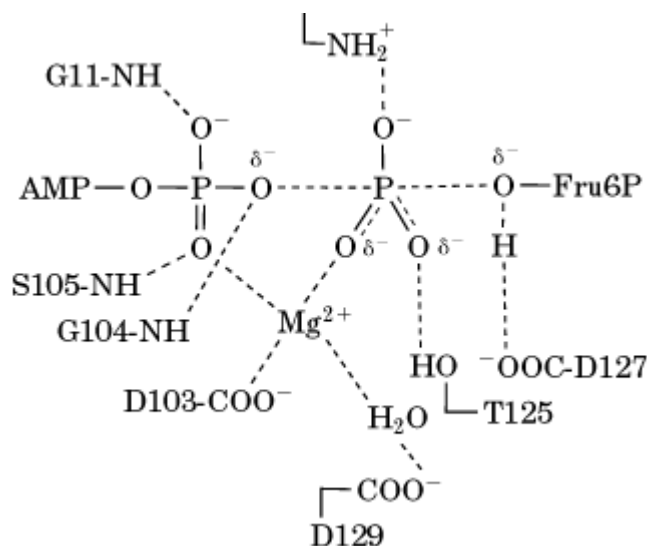
clear.  $PP_i$ -dependent PFKs also have a tetrameric structure that contains two a- and two b-chains. Yeast PFK is active only as an octamer that has four a- and four b-chains.

## 5. Catalytic Mechanism of Phosphoryl Transfer

Experiments with chiral ( $^{16}O$ ,  $^{17}O$ , and  $^{18}O$ )-ATP have shown that transfer of the  $\gamma$ -phosphate from ATP to Fru-6P occurs directly within a ternary complex (6) without a phosphoryl-enzyme intermediate. Steady-state kinetics are compatible with a random order of substrate binding, and the pH-dependence of the maximum velocity suggests that  $Mg^{2+} \cdot ATP^{4-}$  is the most active ionic form (7).

X-ray crystallography of bacterial PFKs and site-directed mutagenesis suggest that several PFK active site groups are involved in catalyzing phosphoryl transfer (Fig. 2). The deprotonated side chain of Asp127 acts as a base and abstracts a proton from the 1-OH group of Fru-6P so as to increase its nucleophilicity. Mutation of Asp127 to Ser decreases the catalytic rate constant  $10^4$ -fold (8). The transferred phosphoryl group is stabilized by interactions with the magnesium ion, by the positive charge of Arg72, and by a hydrogen bond with Thr125. The deprotonated side chains of Asp103 and Asp129 hold the magnesium ion that contributes electrostatic catalysis. The amide NH of Gly104 makes a hydrogen bond with the  $\beta\gamma$ -oxygen bridge to facilitate cleavage of the O-P bond (5). X-ray crystallography also shows that the active site of *E. coli* PFK takes two conformations, one “open” that binds substrates and releases products and one “closed” in which catalysis occurs (5). Rapid kinetics have shown that conformational changes indeed occur in *E. coli* PFK upon substrate binding (and/or product release) and that these changes are slow enough to be rate-limiting in some conditions for the catalytic cycle (9).

**Figure 2.** The transition state of the reaction catalyzed by PFK with the main interactions within the active site as suggested by X-ray crystallography (5) and site-directed mutagenesis. The residues shown are conserved in all of the active PFKs except for Thr125, replaced by a serine in eukaryotic PFKs, which conserves an OH group. The partial covalent bonds are indicated by broken lines, and the hydrogen and electrostatic bonds by dotted lines.



## 6. Purification and assay of activity

Most PFKs are purified by affinity chromatography using a nucleotide analog covalently coupled to a resin. The usual assay measures the formation of Fru-1,6P<sub>2</sub> by coupling it to the disappearance of

NADH using aldolase, triosephosphate isomerase, and glycerol-phosphate dehydrogenase (10, 11).

## 7. Regulation of activity

The catalytic activities of all PFKs are sensitive to pH, ionic strength, divalent metal ions and/or their chelators (the real substrate is the ATP – Mg<sup>2+</sup> complex), and phosphate, NH<sub>4</sub><sup>+</sup>, and K<sup>+</sup> ions (10, 11). In addition, PFKs often have complex regulatory properties, such as sensitivity to allosteric effectors and/or a cooperative saturation by the substrate Fru-6P. Only a few (eg, that from *Dictyostelium discoideum*) follow Michaelis–Menten kinetics under all conditions studied.

### 7.1. Allosteric Effectors

The activities of PFKs are modulated by several allosteric effectors that differ among various organisms (Table 1). Most of these effectors do not modify the maximum velocity, but affect only the saturation by Fru-6P, influencing the cooperativity and/or half-saturating concentration. All PFKs are activated by an increase in levels of ADP and/or AMP, which would express a shortage of cell energy (see [Adenylate Charge](#)). In addition, most eukaryotic PFKs are strongly inhibited allosterically by ATP, and show nonhyperbolic ATP saturation. Because of this inhibition by ATP, PFK in the liver is almost inactive with the *in vivo* concentrations of Fru-6P and ATP, and the total PFK activity is much too low to account for the glycolytic flux measured *in vivo*. This discrepancy between the expected and measured glycolytic flux led to the discovery of the most potent allosteric activator of eukaryotic PFKs, Fru-2,6P<sub>2</sub>, which relieves allosteric inhibition by ATP (12). The concentration of Fru-2,6P<sub>2</sub> in mammalian cells is controlled by hormones that act on cell metabolism through cAMP and protein kinase A. Fru-2,6P<sub>2</sub> decreases when cAMP increases. This hormonal control of PFK activity via Fru-2,6P<sub>2</sub> regulates the balance between the degradation of glycogen and the synthesis of glucose and/or lipids. Fru-1,6P<sub>2</sub> also activates PFKs, but with a much lower affinity than Fru-2,6P<sub>2</sub> (12).

**Table 1. The Main Effects of PFK in Different Organisms**

| PFKs from            | Activation                       | Inhibition     |
|----------------------|----------------------------------|----------------|
| Bacteria             | ADP or GDP <sup>a</sup>          | PEP            |
| Mammals <sup>b</sup> | AMP.ADP<br>Fru-2.6P <sub>2</sub> | Citrate<br>ATP |
| Yeast                | AMP.ADP<br>Fru-2.6P <sub>2</sub> | Citrate<br>ATP |

<sup>a</sup> ADP is probably the physiological effector, but GDP has been used to activate bacterial PFKs in many kinetic studies because ADP is also a competitive inhibitor of the substrate ATP.

<sup>b</sup> The three different isotypes of PFK in mammalian cells are sensitive to the same effectors but to varying extents.

Most eukaryotic PFKs are also allosterically inhibited by citrate, the first metabolite in the Krebs cycle. Citrate inhibition is responsible for the “Pasteur effect” in yeast, a decrease of glucose consumption upon shift from anaerobic to aerobic growth. Indeed, the reoxidation of NADH by

oxygen via the respiratory chain raises the level of the Krebs cycle intermediates, and the citrate inhibition of PFK decreases the rate of glycolysis.

Bacterial PFKs are insensitive to citrate, but are sensitive to feedback inhibition by PEP, the penultimate end product of glycolysis.

## 7.2. Appearance of “New” Regulatory Functions upon Evolution of Sites

The variety of allosteric effectors of PFKs can be explained by their evolution. Bacterial PFKs originally possessed an active site that binds both ATP and Fru-6P, plus an effector site that binds ADP or PEP. After the gene duplication and fusion that doubled the polypeptide chain, the active site on the N-terminal half retained its function, whereas that on the C-terminal active site evolved into a regulatory site (Fig. 1). The Fru-6P subsite became the activation site for Fru-2,6P<sub>2</sub> and the ATP subsite became an inhibitory site. Similarly, evolution could have modified an inhibitory site for PEP in bacterial PFKs into one for citrate in eukaryotic PFKs. The evolution from bacterial to eukaryotic PFKs has also created the structural elements needed for communication between the active site and the “new” regulatory sites.

## 7.3. The Cooperativity of *E. coli* PFK toward Fru-6P and the Concerted Model of Monod, Wyman, and Changeux

The steady-state kinetics of *E. coli* PFK were the first to be analyzed according to the concerted allosteric model. Quantitative agreement between experimental data and predictions from the model was remarkable (1). According to the model, the protein is in equilibrium between two conformational states, R and T. The active R state has a much higher affinity for Fru-6P and for the activator ADP (or GDP), whereas the inactive T state has a much higher affinity for the inhibitor PEP. In the absence of ligand, the equilibrium between the R and T states of free PFK largely favors T, and the a ratio  $T_0/R_0 = 4 \times 10^6$  (1). The simplicity of the concerted model is that a unique transition between the two states R and T is involved in both the cooperativity toward Fru-6P and the influence of allosteric effectors. Some bacterial PFKs show sigmoidal saturations by Fru-6P only in the presence of the inhibitor PEP, confirming that the transition into an inactive state induced by PEP and the cooperativity toward Fru-6P are related.

The crystal structure of the complex between a bacterial PFK and an analog of the inhibitor PEP shows that the latter is bound to a regulatory site remote from the active site and that the protein has changed its quaternary conformation. The change induced by PEP closes the Fru-6P binding site (13). This is consistent with the observation that PEP inhibits by decreasing PFK's apparent affinity for Fru-6P. Some residues that bind PEP are also involved in binding the allosteric activator ADP when PFK is in its active conformation. The opposite effects of binding activators and inhibitors to a single effector site result because it has two mutually exclusive conformations, that agree with the concerted allosteric model. The quaternary structure of the inactive conformation is less stable, which explains the PEP-induced dissociation of the PFK from *T. thermoaquaticus*.

However, the transition between the two states R and T seen by X-ray crystallography cannot explain entirely the cooperativity of *E. coli* PFK toward Fru-6P. In the absence of ligand, free PFK is not in the inactive T<sub>0</sub> state because (1) its crystal structure is that of the R state (14), and (2) it binds Fru-6P with high affinity and without any cooperativity (15). Presteady-state kinetics show that binding substrates and/or effectors to *E. coli* PFK is accompanied by conformational changes within the R state, which are slow enough to be rate-limiting for the catalytic cycle. Thus, the cooperativity of *E. coli* PFK is in part independent of the transition between the crystallographic R and T states and is due to a change in the rate-limiting step from a conformational change at low Fru-6P concentrations to phosphoryl transfer at high Fru-6P (9). Such a kinetic origin of cooperativity explains observations with some mutants that are difficult to fit with the concerted model, such as values of the Hill coefficient larger than the number of Fru-6P binding sites (16) and reversal of the influence of effectors, where PEP becomes an activator (17) and GDP an inhibitor.

#### 7.4. Phosphorylation of Mammalian PFKs by Protein Kinase A

Mammalian PFKs are phosphorylated on specific sites by protein kinase A. This produces changes in kinetic properties and in the tendency to aggregate into octamers, but the physiological role of this phosphorylation is unknown.

#### Bibliography

1. D. Blangy, H. Buc, and J. Monod (1968) *J. Mol. Biol.* **31**, 13–35.
2. G. A. Dunaway, T. P. Kasten, T. Sebo, and R. Trapp (1988) *Biochem. J.* **251**, 677–683.
3. P. Branny, F. De la Torre, and J.-R. Garel (1993) *J. Bacteriol.* **175**, 5344–5349.
4. R. A. Poorman, A. Randolph, R. G. Kemp, and R. L. Heinrikson (1984) *Nature* **309**, 467–469.
5. Y. Shirakihara and P. R. Evans (1988) *J. Mol. Biol.* **204**, 973–994.
6. R. L. Jarvest, G. Lowe, and B. V. L. Potter (1981) *Biochem. J.* **199**, 427–432.
7. I. Auzat, G. Le Bras, and J.-R. Garel (1994) *Protein Peptide Lett.* **1**, 179–182.
8. H. Hellinga and P. R. Evans (1987) *Nature* **327**, 437–439.
9. I. Auzat, E. Gawlita, and J.-R. Garel (1995) *J. Mol. Biol.* **249**, 478–492.
10. D. Kotlarz and H. Buc (1982) *Methods Enzymol.* **90**, 60–70.
11. K. Uyeda (1979) *Adv. Enzymol.* **48**, 193–244.
12. E. Van Schaftigen (1987) *Adv. Enzymol.* **59**, 315–395.
13. T. Schirmer and P. R. Evans (1990) *Nature* **343**, 140–145.
14. W. R. Rypniewski and P. R. Evans (1989) *J. Mol. Biol.* **207**, 805–821.
15. D. Deville-Bonne and J.-R. Garel (1992) *Biochemistry* **31**, 1695–1700.
16. I. Auzat, G. Le Bras, and J.-R. Garel (1995) *J. Mol. Biol.* **246**, 248–253.
17. F. T. K. Lau and A. R. Fersht (1987) *Nature* **326**, 811–812.

#### Suggestions for Further Reading

18. P. R. Evans (1992) "Activity and allosteric regulation in bacterial phosphofructokinase", *Proceedings from The Robert A. Welch Foundation Conference on Chemical Research XXXVI*, pp 39–54. A reference that may be difficult to find but has an outstanding short and clear synthesis of the structure, catalytic mechanism, regulation, and evolution of phosphofructokinases.
19. L. A. Fothergill-Gilmore and P. A. M. Michels (1993) Evolution of glycolysis, *Prog. Biophys. Mol. Biol.* **59**, 105–235. An exhaustive comparison of glycolytic enzymes from various species with fine sections on the structure, evolution, and expression of the genes that code for phosphofructokinases and a very complete list of references.

#### Phospholipases

Phospholipases hydrolyze membrane phospholipids (see [Lipases](#)). The lipid substrates for these enzymes are usually quite specific. Hormone-regulated [phospholipase C](#) (PLC) can hydrolyze a number of membrane phospholipids, including phosphatidylcholine and phosphatidylethanolamine. Phosphatidylinositol (PI) and its phosphorylated derivatives (polyphosphoinositides) are the most common substrates involved in [second messenger](#) generation, producing soluble **inositol phosphates**, such as inositol trisphosphate (IP<sub>3</sub>), that are involved in [calcium signaling](#).

Phospholipase D enzymes usually prefer phosphatidylcholine as substrate, generating phosphatidic

acid and choline. Phosphatidic acid has been proposed to mediate a variety of biological responses involved in mitogenesis and protein trafficking (1).

The turnover of polyphosphoinositides occurs in a cycle of hydrolysis and resynthesis, called the “PI cycle.” PI is phosphorylated sequentially on the 4' and 5' hydroxyl groups of the inositol ring, to produce PI-4',5'-P<sub>2</sub>. Upon hormonal activation, PLC subsequently catalyzes the phosphodiesteratic cleavage of this lipid, resulting in the generation of diacylglycerol and IP<sub>3</sub> in the cell. After its degradation, PI can be resynthesized by the phosphorylation of diacylglycerol to PA, via the activity of a diacylglycerol kinase and the subsequent resynthesis of PI via PI synthase. Thus, this cycle results in the generation of two potential second messengers, the membrane-associated diacylglycerol, which can activate [protein kinase C](#), as well as IP<sub>3</sub>, which can lead to the mobilization of intracellular calcium levels.

Several forms of PI-PLC have been identified (1). These proteins all exhibit an absolute requirement for calcium and are equally capable of hydrolyzing PI and the polyphosphoinositides PIP and PI-4',5'-P<sub>2</sub>. These proteins exist in three subfamilies, b, g, and d, each with multiple members. They have a number of conserved and divergent domains, reflecting their differential regulation. All three subclasses have a pleckstrin homology (PH) domain, which probably mediates some type of interaction with lipids in the membrane, and an [EF-hand motif](#), presumably involved in **calcium binding**. Moreover, there are two additional conserved domains, designated X and Y, that contain the catalytic core of the enzyme; there is also a conserved C2 domain, also likely to be required for calcium interaction. PLCg, which is regulated by [tyrosine kinase receptors](#), contains two **SH2 domains** and one SH3 domain and has been shown to undergo **phosphorylation** of [tyrosine](#) residues. PLCb is responsible for regulation through [G-protein-coupled receptors](#) and is generally regulated by a subunits of the Gq family of [heterotrimeric G proteins](#), linked to a subset of such receptors. Additionally, their bg subunits have also been shown to regulate PLCb, although the efficacy of this activation is not comparable to that seen with Gqa. PLCb can also act as an activator of the [GTPase](#) activity of Gq, effectively behaving as a GTPase activating protein (GAP).

## Bibliography

1. W. D. Singer, H. A. Brown, and P. C. Sternweis (1997) *Annu. Rev. Biochem.* **66**, 475–509.

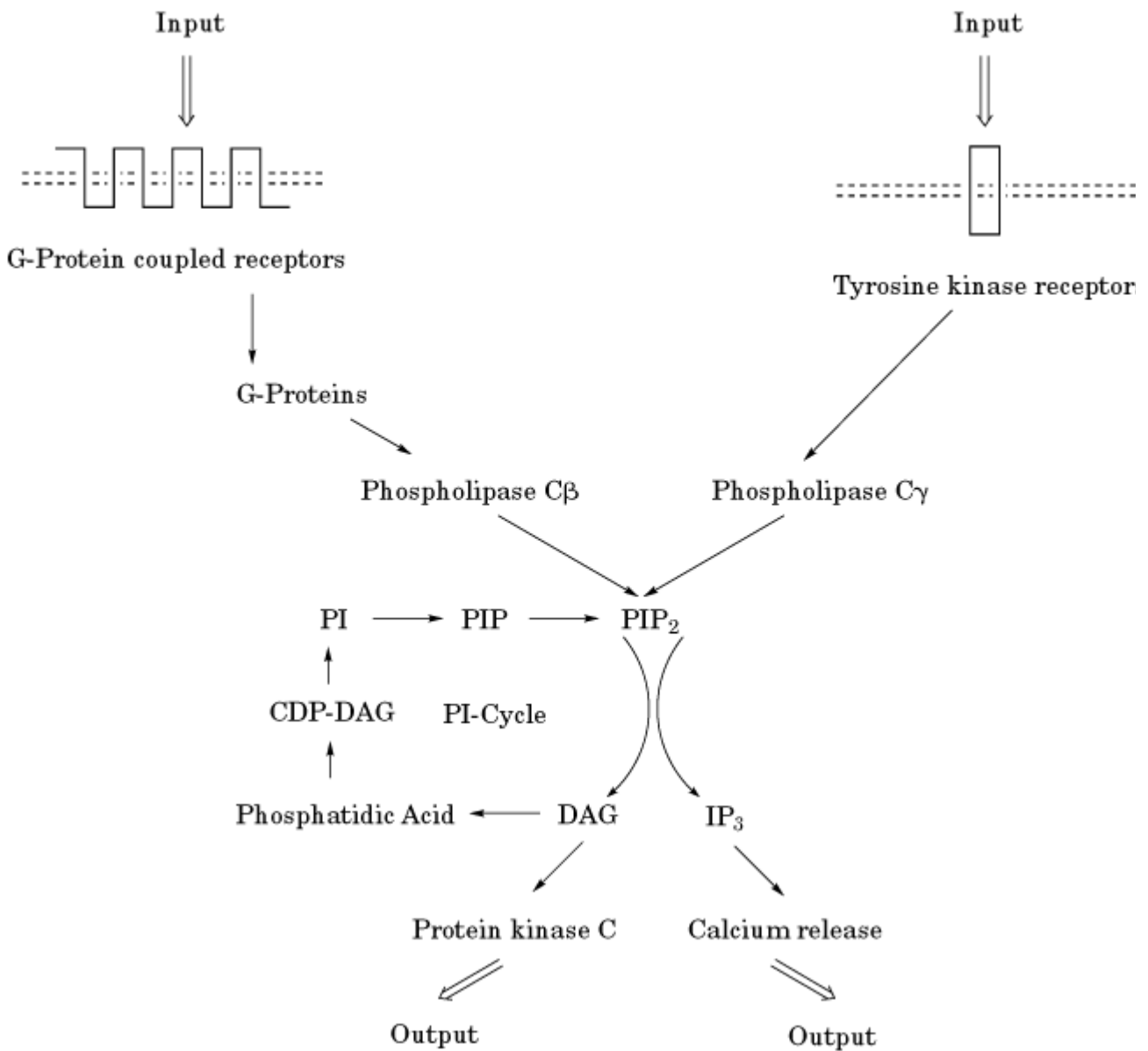
## Phospholipases C

### 1. Historical Perspective

The phosphatidylinositol (PI)-specific phospholipases C have emerged as key and determining components of the PI cycle of [signal transduction](#). In this pathway, the action of any of several extracellular agents (including [hormones](#), [cytokines](#), neurotransmitters, and xenobiotics) results in activation of PI-specific phospholipases C that cleave phosphatidylinositol bisphosphate (PIP<sub>2</sub>) (1-3). This results in the formation of inositol trisphosphate (IP<sub>3</sub>) and **diacylglycerol** (DAG), which act as important [second messengers](#) involved in the regulation of [calcium signaling](#) and protein kinase C, respectively (Fig. 1).

**Figure 1.** Regulation and function of the phosphatidylinositol (PI) cycle. The scheme illustrates two mechanisms for

activation of the PI cycle through (1) direct phosphorylation of phospholipase C  $\gamma$  by tyrosine kinase receptors, and (2) C proteins in the case of G-protein-coupled receptors. The two products of  $\text{PIP}_2$  hydrolysis,  $\text{IP}_3$  and diacylglycerol (DAG), function as second messengers in activating calcium release and protein kinase C, respectively. Resynthesis of  $\text{PIP}_2$  from DAG completes the PI cycle.



Early insight into the PI cycle came from studies by Hokin and Hokin, who identified turnover of inositol phospholipids in cells stimulated with acetylcholine. Work over the ensuing decades resulted in the delineation of the biochemical components of the PI cycle and the realization of its intimate connection to changes in calcium levels in the cells. Thus, it became appreciated that specific phospholipases of the C type were activated in response to acetylcholine and other hormones. It was also appreciated that  $\text{PIP}_2$  was the preferred substrate in this pathway and that this was coupled to changes in calcium homeostasis (3). It was then determined that the  $\text{IP}_3$  product of this reaction acted directly on specific **receptors** internal to the cell to cause changes in the levels of free calcium. Contemporaneously, a direct connection was made between DAG and protein kinase C, demonstrating a second messenger function for DAG (4, 5).

## 2. Regulation of Phospholipase C



PI-specific phospholipases C have been purified and cloned. They exist as a family of related phospholipases with close similarity in catalytic activity and preference for  $\text{PIP}_2$  as a substrate, but with distinct mechanisms of regulation (1). Notably, the two best-studied members of the family, phospholipase C $\beta$  and phospholipase C $\gamma$ , are coupled through distinct mechanisms to two different classes of membrane receptors. **G protein**-coupled receptors regulate phospholipase C $\beta$  through direct activation by hetero-trimeric G proteins. On the other hand, phospholipase C $\gamma$  is tyrosine phosphorylated (see [Phosphorylation, Protein](#)) by transmembrane tyrosine kinase receptors (and also by nonreceptor tyrosine kinases) through interactions with specific **SH2** modules in phospholipase C $\gamma$  that are recognized by specific SH2-binding **domains** in the tyrosine kinases (Fig. 1).

### 2.1. Role of Phospholipases C in Signal Transduction and Cell Biology

Given their intimate coupling to two major classes of transmembrane receptors, phospholipases C serve a critical role in allowing the ligands of these receptors to regulate the intracellular levels of calcium and the activity of protein kinase C. This is accomplished by the activity of the two products of the phospholipase C reaction.  $\text{IP}_3$  is a hydrophilic and diffusible small molecule that interacts specifically and with high affinity with the  $\text{IP}_3$  receptor in the [endoplasmic reticulum](#). This results in release of calcium from stores in the endoplasmic reticulum into the cytosol. In turn, free calcium appears to coordinate a variety of cellular responses through the activation of various enzymes, including protein kinases and **phosphodiesterases** (see [Calcium Signaling](#)). Simultaneously, DAG, which is formed in the membrane, activates several members of the protein kinase C family. Protein kinase C is capable of phosphorylating a variety of substrates that result in the regulation of diverse biological responses. Protein kinase C is also activated potently by **phorbol esters**, which are plant toxins that cause a myriad of cellular responses, including effects on blood cell activation, cell growth regulation, tumor promotion, memory formation, immune responses, inflammatory reactions, hormone release, and neurotransmission, most likely through activation of protein kinase C (4). Therefore, the two products of the phospholipase C reaction serve as critical second messengers, often acting in synergy to regulate a variety of cellular responses.

### Bibliography

1. S. G. Rhee, P. G. Suh, S. H. Ryu, and S. Y. Lee, (1989) *Science* **244**, 546–550.
2. P. W. Majerus, T. M. Connolly, H. Deckmyn, T. S. Ross, T. E. Bross, H. Ishii, V. Bansal, and D. Wilson (1986) *Science* **234**, 1519–1526.
3. M. J. Berridge, P. H. Cobbold, and K. S. R. Cuthbertson, (1988) *Phil. Trans. Roy. Soc. Lond.* **320**, 325–343.
4. Y. Nishizuka (1992) *Science* **258**, 607–614.
5. R. M. Bell (1986) *Cell* **45**, 631–632.

### Suggestions for Further Reading

6. M. Berridge (1993) Inositol trisphosphate and calcium signalling. *Nature* **361**, 315–324.
7. N. Divecha and R. Irvine (1995) Phospholipid signaling. *Cell* **80**, 269–278.

## Phosphorous Isotopes

Phosphorus is element number 15 in the periodic table and has valence states of 3 or 5 (1). Seventeen isotopes have been identified (2), ranging in atomic mass number from  $^{26}\text{P}$  to  $^{42}\text{P}$ , most of which

have physical half-lives of seconds or milliseconds and do not have important applications. Only one stable isotope of phosphorus is found in nature ( $^{31}\text{P}$ ).

The most important **radioactive** isotope of phosphorus is  $^{32}\text{P}$  (half-life = 14.29 days) (see [Radioisotopes](#)). It undergoes beta-minus decay to  $^{32}\text{S}$ , which is stable. Phosphorus-32 yields one beta particle per decay, with an energy of 1.71 MeV maximum, 0.695 MeV on average. Another useful radioisotope of phosphorus is  $^{33}\text{P}$  (half-life = 25.4 days). It yields one beta particle per decay, with an energy of 0.1667 MeV maximum, 0.0486 MeV on average.

Phosphorus-32 is usually prepared in an accelerator by bombarding a stable sulfur target with deuterons according to the reaction  $^{34}\text{S}(d,a)^{32}\text{P}$ . It may also be prepared by accelerator deuteron bombardment of  $^{31}\text{P}$  according to the reaction  $^{31}\text{P}(d,p)^{32}\text{P}$  or by reactor neutron-activation according to  $^{31}\text{P}(n,g)^{32}\text{P}$ . Carrier-free radioisotope is prepared by reactor production according to the reaction  $^{32}\text{S}(n,p)^{32}\text{P}$ .

Phosphorus-32 was first made by Fermi in 1934 using neutrons from a radium–beryllium source (3). In 1935, Hevesey made  $^{32}\text{P}$  by the same method and studied the metabolism of phosphorus compounds in rats (3). Hevesey also studied blood volume using  $^{32}\text{P}$ -labeled red blood cells. When it was found that  $^{32}\text{P}$  killed blood cells, Lawrence produced it in a cyclotron and in 1938 was the first to use it as a therapeutic radioisotope for treatment of leukemia in humans (3). Phosphorus-32 is presently used as an effective treatment for polycythemia vera (4, 5), for treatment of painful bone metastases from prostate and other cancers (6), for restenosis stents in treating heart disease, and in colloidal form for direct intratumoral injection, or *infusional brachytherapy* (7). The systemic use of  $^{32}\text{P}$  has been limited by concerns over potential leukemogenic effects.

As a diagnostic agent,  $^{32}\text{P}$  has been used to study skeletal growth and bone tumor kinetics. In molecular biology,  $^{32}\text{P}$  is used primarily to label **nucleic acids**, which have one phosphorous atom per nucleotide, and for **phosphorylation** of proteins. Phosphorus-32 is readily detected by beta-particle liquid scintillation counting, **Cerenkov radiation detectors**, [autoradiography](#), and [fluorography](#).

#### Bibliography

1. D. R. Lide and H. Pr. Frederikse, eds. (1995) *CRC Handbook of Chemistry and Physics*, CRC Press, Boca Raton, Fla.
2. Knolls Atomic Power Laboratory (1966) *Chart of the Nuclides*, 15th ed., available from General Electric Company San Jose, Calif.
3. M. Brucer (1990) *A Chronology of Nuclear Medicine*, Heritage Publications, Inc., St. Louis, Mo.
4. J. Lazlo (1974) *Postgrad. Med.* **55**, 168–173.
5. C. Parmentier and P. Gardet (1994) *Nouv. Rev. Fr. Hematol.* **36**, 189–192.
6. E. B. Silberstein (1993) *Semin. Oncol.* **20** (suppl. 2), 10–21.
7. S. E. Order et al. (1996) *Ann. Acad. Med. Singapore* **25**, 347–351.

#### Suggestion for Further Reading

8. H. N. Wagner, Jr. (1968) *Principles of Nuclear Medicine*, W. B. Saunders Co., Philadelphia.

## Phosphorylase Kinase

Phosphorylase kinase (E.C. 2.7.1.38) catalyzes the **phosphorylation** of residue Ser14 in relatively inactive [glycogen phosphorylase b](#), converting it to the much more active phosphorylase *a*. Its discovery in 1955 by Fischer and Krebs, and their later demonstration that the [enzyme](#) converting phosphorylase *a* to *b* is a protein **phosphatase** that removed the phosphate, provided a firm basis for further investigation of the first metabolically interconvertible enzyme system to be discovered. These investigators and their colleagues characterized phosphorylase kinase and demonstrated that it was itself activated by phosphorylation by [cyclic AMP](#) (cAMP)-dependent **protein kinase** (cAPK), thus suggesting the famous cascade system for the regulation of glycogen metabolism by [hormones](#). In addition, it was discovered in this same laboratory that  $\text{Ca}^{2+}$  is required for activity, and later by Drummond et al. that the activation of glycogenolysis is coordinated to muscular contraction by the release of calcium from the [endoplasmic reticulum](#) in response to nervous stimulation (see [Calcium Signaling](#)). It was not until 1978 that the small calcium-binding protein [calmodulin](#) was found to be an integral part of phosphorylase kinase and responsible for the allosteric activation by  $\text{Ca}^{2+}$ . Calmodulin was designated subunit d because three other subunits, the regulatory a and b and the catalytic g, were already known. Each is present in four copies, so the **holoenzyme** may be represented as  $a_4b_4g_4d_4$  with an overall molecular weight of 1.3 million. This enormous protein presents a formidable challenge to investigators, but the subunits have been isolated and their primary sequences determined, while progress is being made on their functions and on the complex regulation of this enzyme. The gene for the catalytic domain of subunit g has been **cloned**, and the structure of the subunit determined by [X-ray crystallography](#), so the molecular basis of the catalytic activity has been elucidated. The comprehensive review by Picket-Gies and Walsh (1) covers the literature until 1986 and can be used as the essential reference for this article. It is suggested for further reading.

### 1. Structure and Function

The molecular weights of the subunits have been determined by sequence analyses (2) and are  $a = 138,422$ ,  $b = 125,205$ ,  $g = 44,673$ , and  $d = 16,680$ . Various analyses by [electron microscopy](#) of  $a_4b_4g_4d_4$  show a bilobed butterfly-like structure that may be cleaved by **trypsin** into two halves, suggesting that the subunits are arranged as a bridged dimer of octamers ( $a_2b_2g_2d_2$ ), while another study found a “chalice” form. These observations have been reconciled by a three-dimensional model that can account for all views (3). Essentially there are three twofold axes of symmetry relating the four abgd protomers, resulting in a tetrahedral arrangement. Proteolytic digestion suggests that the b subunits play an essential role in holding the holoenzyme together, but there are additional interactions of other subunits. Lithium bromide causes dissociation into two complexes, agd\_ and gd\_ both of these have activity toward phosphorylase *b* equal to that of the holoenzyme.

The determination of the amino acid sequences of all four subunits provided valuable insights into their various functions. The g subunit proved to have substantial homology with the catalytic subunit of cAPK. There is an additional C-terminal region (approximately residues 286 to 386) that can be removed by **chymotrypsin**, to yield a catalytic domain (residues 1 to 285) with full activity and no calmodulin/ $\text{Ca}^{2+}$  sensitivity. The C-terminal or regulatory **domain** contains two calmodulin-binding domains, suggesting a double binding mode that accounts for the unusual phenomenon that calmodulin binds to the isolated g subunit in the absence of calcium, becoming subunit d, and suppresses its activity. **Troponin C**, which is **homologous** to calmodulin, also binds to the g subunit and may provide a more direct link between the calcium activation of skeletal muscle contraction and glycogenolysis. The a and b subunits are highly homologous, suggesting a common ancestral protein. Both have calmodulin-binding domains, but apparently only one extrinsic calmodulin binds per pair of a/b units, known as d' subunits, in the intact holoenzyme in the presence of excess

calcium, resulting in four d' subunits in all. The b subunit contains two sites for phosphorylation by cAPK, on Ser26 and Ser700. The two sites appear to be phosphorylated equally rapidly, and either one results in an increase in activity of the holoenzyme, but a total of only one phosphate per subunit is found. The rate of phosphate incorporation into the b subunit is 5- to 10-fold faster than that into the a subunit and is a prerequisite for the latter, suggesting a conformational change. The a subunit is phosphorylated at Ser1018 by cAPK, with resultant increases in activity. Both the a and b subunits have several sites for autophosphorylation, but the evidence that this phenomenon has significance for *in vivo* regulation is ambiguous.

Kinetics of the phosphorylase kinase reaction were studied using a tetradecapeptide from phosphorylase *b* (residues 5 to 18) that contained the phosphorylatable Ser14 in order to avoid substrate-directed effects. The **kinetic mechanism** was found to be rapid equilibrium random Bi-Bi, indicating that either ATP or the protein can bind first and that the phosphoryl transfer step in the ternary complex is rate-limiting. Several studies on the effects of deleting or substituting various amino acid residues of the peptide permit one to conclude that the kinase shows a marked preference for peptides having the natural sequence around Ser14, while the latter is also important because a Thr14 peptide has a very low rate of phosphorylation. Because the first 18 residues of phosphorylase *b* are disordered, it is not surprising that the corresponding peptide, presumably in a random conformation, should be a good substrate for the kinase. The kinase must, however, recognize more than just the *N*-terminal peptide, because the intact protein substrate has a  $K_m$  at least fivefold lower and a maximum velocity ( $V_{max}$ ) fivefold greater. Phosphorylase kinase will also phosphorylate Ser7 of glycogen synthase, but this is also the target for other kinases, including cAPK (which does not act on phosphorylase), so the physiological significance of its action on this and several other proteins is not known.

The high-resolution structure of the crystallized catalytic domain of the g subunit has been determined by Johnson et al. (4) and compared to other protein kinases (5). The structure is very similar to that of cAPK and contains the same catalytic residues in equivalent positions. These include Lys57 to interact with the phosphates of ATP and Asp167 to interact with the  $Mg^{2+}$  of ATP/ $Mg^{2+}$ . An *N*-terminal lobe, mostly **b-sheet**, is connected to the  $\alpha$ -helical *C*-terminal lobe by polypeptide chain, allowing opening and closing of the **active site** cleft between the two lobes. The suggested catalytic mechanism is that Asp149 abstracts a proton from the substrate serine residue and the resultant alcoholate ion attacks the terminal phosphate of ATP. In cAPK, an activation segment has Thr197 phosphate, which interacts with Arg165, next to the catalytic Asp166, and other basic residues, but in phosphorylase kinase this is replaced by Glu182, which forms an ionic link with Arg148 of the catalytic loop and serves the same function of maintaining Asp149 and other catalytic residues in the correct orientation.

## 2. Regulation

Because so many results from *in vitro* studies may have no physiological significance, only those most likely to apply *in vivo* will be considered here. Because nearly all the phosphorylase and 20% to 40% of the kinase are associated with the glycogen particle in muscle, the kinase action on phosphorylase must occur with the latter bound to glycogen by its catalytic face, leaving the control face with Ser14 exposed to the cytosol (see [Glycogen Phosphorylase](#)). It is relevant, then, that glycogen decreases the  $K_m$  for phosphorylase *b* 10-fold; and this is likely to be via conformational changes in the substrate, because glycogen has no effect on phosphorylation of the tetradecapeptide containing residues 5 to 18. The **allosteric** inhibitor of phosphorylase *b*, glucose-6-phosphate, also inhibits kinase activity on the protein, but not on the peptide, suggesting augmentation of the metabolite feedback. There is a good correlation between the increase of  $Ca^{2+}$  levels due to nervous or electrical stimulation, muscle contraction, and the activation of both phosphorylase kinase and its substrate, and this occurs without changes in either the level of cAMP or the phosphorylation state of the kinase. The molecular basis is the allosteric activation due to binding of  $Ca^{2+}$  to the d subunit

and, at higher  $\text{Ca}^{2+}$  levels, of the binding of  $\text{Ca}^{2+}$ -calmodulin to the a/b\_subunits (the d' subunits). In contrast, the administration of adrenaline can cause coordinated changes in cAMP levels, activation of cADK and phosphorylase kinase, and formation of phosphorylase *a*. A point still in contention is whether or not  $\text{Ca}^{2+}$  is still an obligate requirement for kinase catalysis, as *in vitro* experiments suggest, but the increased affinity for the metal ion after phosphate incorporation into the b and\_a\_subunits may allow low endogenous levels to suffice. Heilmeyer (2) summarizes current knowledge of regulation, including a possible role for an allosteric ATP/Mg<sup>2+</sup> site on the a subunit that may stimulate binding of the same substrate to the g subunit by displacing inhibitory ADP. He suggests that an ancestral gd phosphorylase kinase, responsive only to  $\text{Ca}^{2+}$ , could not handle all the signals needed to regulate the conversion of phosphorylase *b* to *a*. The addition of the two regulatory subunits b and\_a permitted regulation by cAMP-dependent phosphorylation, extrinsic calmodulin, and substrate binding (ATP/Mg<sup>2+</sup>).

## Bibliography

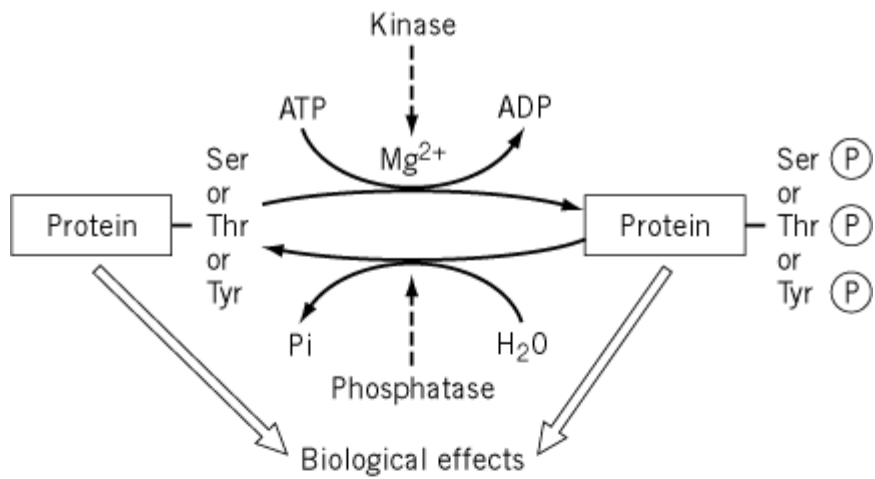
1. C. A. Pickett-Gies and D. Walsh (1986) *The Enzymes*, Vol. 17, 3rd. ed. pp. 395–459.
2. L. M. G. Heilmeyer, (1991) *Biochim. Biophys. Acta* **1094**, 168–174.
3. M. T. Norcum, et al. (1994) *J. Mol. Biol.* **241**, 94–102.
4. D. J. Owen, et al. (1995) *Structure* **3**, 467–482.
5. L. N. Johnson, M. E. M. Noble, and D. J. Owen (1996) *Cell* **85**, 149–158.

## Phosphorylation, Protein

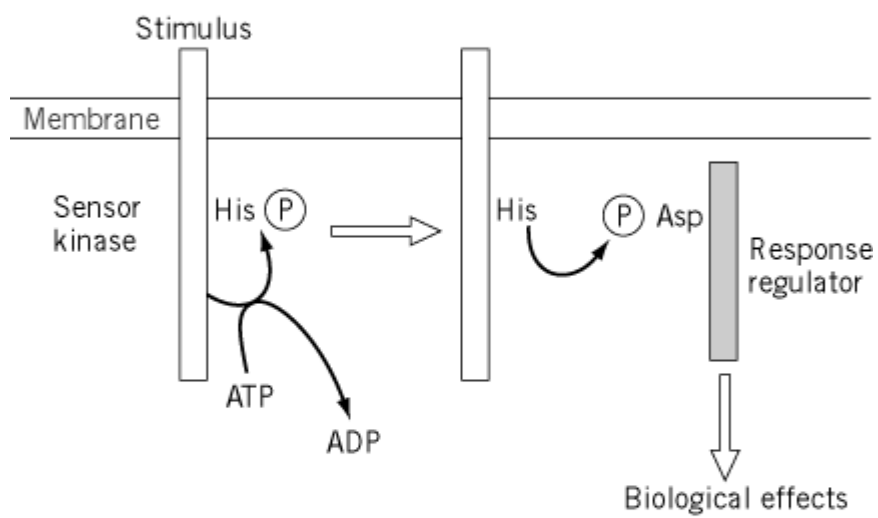
Phosphorylation is a common, and usually reversible, [post-translational modification](#) of [proteins](#). The amino acid residues that can be phosphorylated are listed in Table 1. Protein phosphorylation is either an intermediate step in the course of an **enzymatic** reaction or a more durable modification that regulates protein properties. This entry focuses on regulatory protein phosphorylation. In this context, the phosphorylation reaction is catalyzed by protein **kinases**, which form a specific class of phosphoryl-transferases, generally transferring the phosphoryl group from the g position of ATP (Fig. 1a, see exceptions below) to the side chain of an amino acid residue. The dephosphorylation reaction corresponds to the transfer of the phosphoryl group to a [water](#) molecule (ie, the hydrolysis of the phosphoester bond). Both reactions are energetically favorable and occur readily in the presence of the appropriate catalyst, a protein kinase or **phosphatase**. In a living cell, the level of phosphorylation of a protein is the result of a dynamic equilibrium between its rates of phosphorylation and dephosphorylation. As mentioned above, the phosphoryl donor is usually ATP, in complex with Mg<sup>2+</sup>. *In vitro*, for some protein kinases, ATP or Mg<sup>2+</sup> can be replaced by GTP or Mn<sup>2+</sup>, respectively. In somewhat artificial conditions (ie, very low levels of ATP and high levels of ADP), protein kinases can catalyze the transfer of phosphoryl group from the phospho-amino acid of a protein to ADP, generating ATP. Although secreted proteins or peptides are known to be phosphorylated, the regulatory phosphorylation reactions that have been well-characterized to date take place within cells.

**Figure 1.** Major pathways of regulatory protein phosphorylation/dephosphorylation. (a) The most general case, observed in all types of living organisms, corresponds to phosphorylation of serine, threonine, or tyrosine residues by protein kinases using ATP/Mg<sup>2+</sup> as phosphoryl donor. The dephosphorylation is catalyzed by a protein phosphatase,

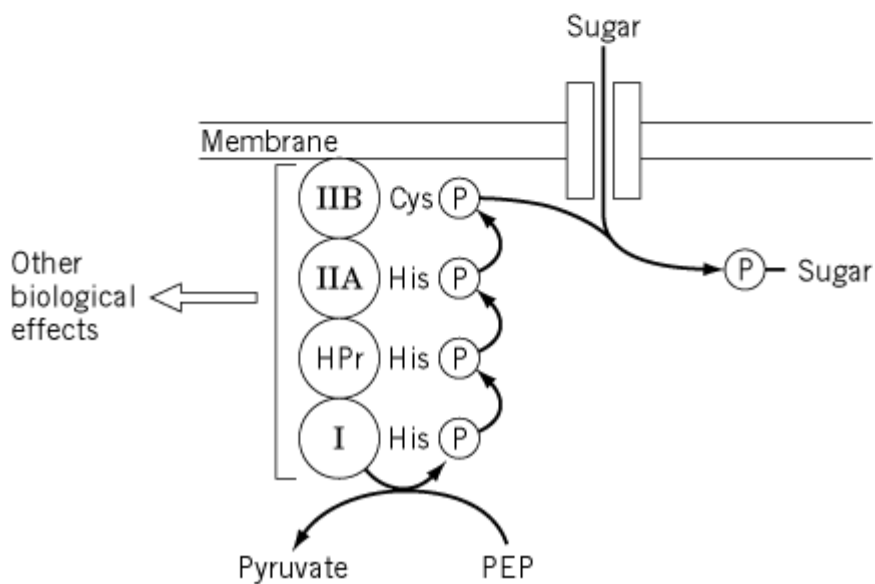
usually distinct from the protein kinase. The activity of both the protein kinase and phosphatase can be regulated by intracellular or extracellular stimuli. The phosphorylated or the dephosphorylated forms of the protein substrate, or both, may be responsible for the biological response observed. **(b)** Although not as common as the previous phosphorylation system, the two-component system is widespread in eubacteria and has also been found in archaeobacteria and eukaryotes. A “sensor protein kinase” autophosphorylates histidine residues in response to an environmental stimulus, using ATP/Mg<sup>2+</sup> as phosphoryl donor. The phosphate is then transferred to an aspartate residue of a “response regulator,” which may be located on the same polypeptide chain as the sensor kinase, or on a distinct one. The phosphorylated response regulator triggers the biological response. **(c)** The third phosphotransferase system has been described in bacteria only. In this system coupled to sugar transport, the phosphoryl group is transferred successively from phosphoenolpyruvate (PEP) to histidine and cysteine residues of several proteins (indicated with roman numerals) and to an incoming sugar. In addition to its function in sugar metabolism, this system has been shown to be coupled to several biological responses. HPr, histidine-containing phosphocarrier protein (see the text for references); Pi, inorganic phosphate;  $\text{P}$ , phosphate.



(a)



(b)



(c)

**Table 1. Amino Acid Residues That Can Be Phosphorylated in Proteins**


| Acceptor group   | Amino acid                  | Product                                    |
|------------------|-----------------------------|--|
| Alcohol          | Serine, threonine           | Phosphoester                               |
| Phenol           | Tyrosine                    | Phosphoester                               |
| Basic amino acid | Histidine, arginine, lysine | Phosphoramidate                            |
| Thiol            | Cysteine                    | Phosphate thioester                        |
| Acyl             | Aspartate, glutamate        | Mixed phosphate-carboxylate acid anhydride |

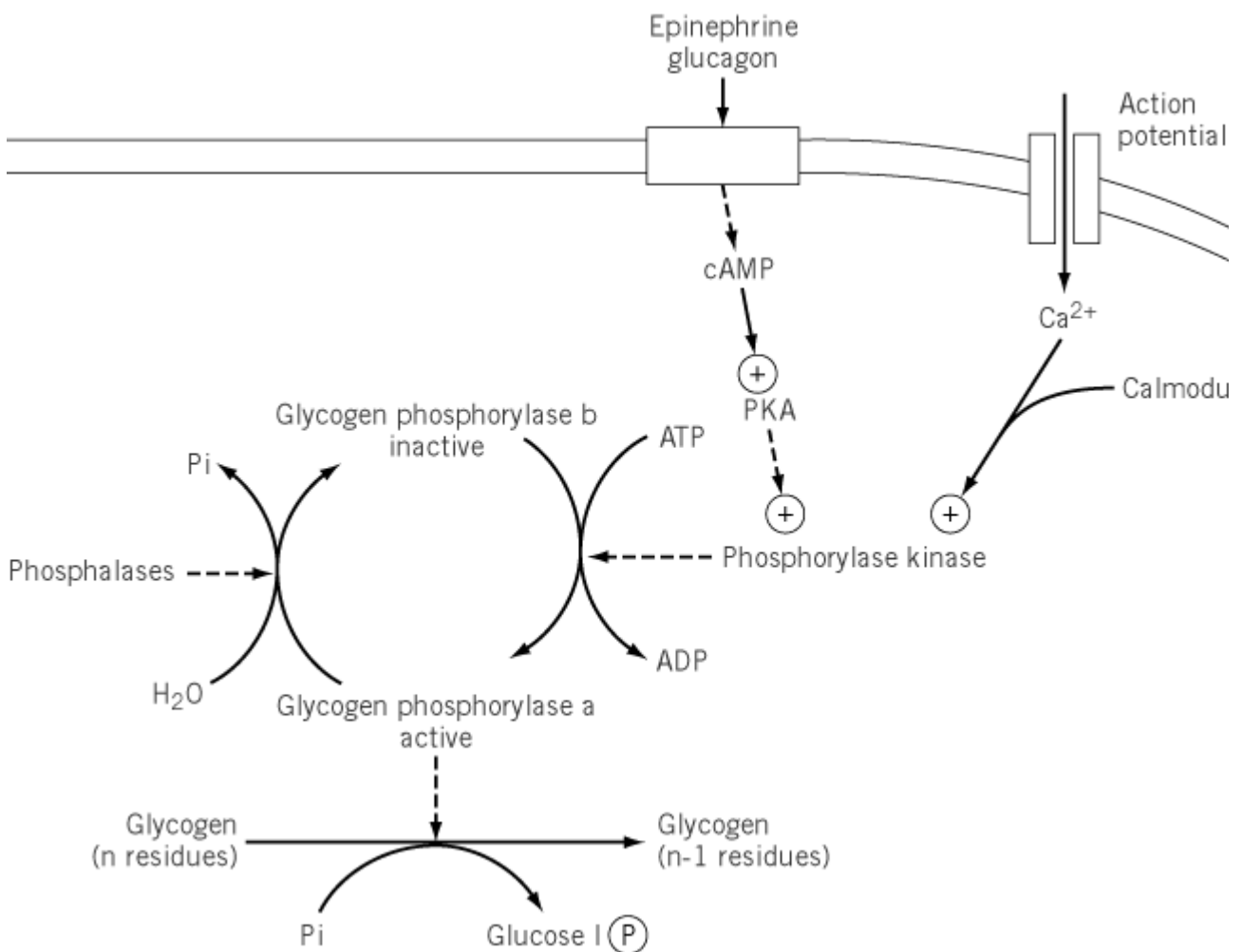
## 1. Historical Perspective

Regulation by protein phosphorylation was discovered by Edmund Fischer and Edwin Krebs in the course of their studies of glycogen metabolism between 1955 and 1958 (1, 2). It was known from the work of Carl and Gerty Cori that [glycogen phosphorylase](#), the enzyme that degrades glycogen, exists in two forms, an active form (phosphorylase a) and an inactive form (phosphorylase b). Fischer and Krebs showed that the only difference is that phosphorylase a is phosphorylated and they identified [phosphorylase kinase](#), the protein kinase that phosphorylates phosphorylase b on Ser14, transforming it into the active phosphorylase a (Fig. 2). The same authors later showed that [cyclic AMP](#), the intracellular [second messenger](#) formed under the action of epinephrine or [glucagon](#), exerts its effects on glycogen catabolism by activating a protein kinase. This protein kinase, cyclic AMP-dependent protein kinase (or [protein kinase A](#), PKA), phosphorylates and activates phosphorylase kinase. In skeletal muscle, they demonstrated that phosphorylase kinase is also activated by  $\text{Ca}^{2+}$ -[calmodulin](#), allowing the coupling of metabolic events with contraction. Although it was originally thought that these reactions might be specific for the regulation of glycogen, the subsequent work of Paul Greengard and others showed that phosphorylation of proteins accounts for the action of many extracellular signals, and several protein kinases activated by [cyclic GMP](#),  $\text{Ca}^{2+}$ -calmodulin and  $\text{Ca}^{2+}$ -[lipids](#) were discovered (3-5). During the same period, the protein phosphatases involved in the regulation of glycogen metabolism started to be studied methodically (6). This opened a new field, because it is now apparent that regulation of protein phosphatases is as important as regulation of protein kinases in the control of cell function. Another major step in the study of protein phosphorylation was the discovery of phosphorylated [tyrosine](#) residues by Tony Hunter in 1979 (7). Phosphotyrosine raised considerable interest when it was discovered that it is increased in cells transformed by [viruses](#), and that v-*src*, and v-*abl*, two major viral [oncogenes](#), are tyrosine kinases. Shortly thereafter, it was shown that [growth factor](#) receptors, which are [transmembrane](#) proteins, have an intrinsic tyrosine kinase activity associated with their intracellular domain (8) (see [Tyrosine Kinase Receptors](#)). This tyrosine kinase activity is stimulated upon binding of the growth factor to the extracellular domain of the receptor. Thus, it became evident that, in spite of its low abundance in normal cells, phosphotyrosine residues are a key step in the control of cell growth and [differentiation](#). More recently, molecular [cloning](#) of [genes](#) has identified a very large number of protein kinases and phosphatases in all living organisms that have been studied. It is now apparent that protein phosphorylation is a fundamental mechanism of [signal transduction](#) used by all cells to regulate the properties of their proteins in response to external or internal clues. The list of fundamental biological processes regulated by protein phosphorylation would be nearly endless, but would include gene expression, energy production and storage, cell division and differentiation, as well as synaptic plasticity in the nervous system, which is thought to underlie learning and memory.



**Figure 2.** Regulation of glycogen degradation. The main pathways regulating glycogen degradation, the study of which led to the discovery of protein phosphorylation, are illustrated. Glycogen is degraded stepwise by glycogen phosphorylase. Schematically, glycogen phosphorylase exists in either a dephosphorylated inactive form (b) or a phosphorylated active form (a). The transformation of b-form to a-form is catalyzed by phosphorylase kinase. Phosphorylase kinase itself is activated by two independent pathways. In skeletal muscle,  $\text{Ca}^{2+}$  entering the cell during the action potential binds to calmodulin and activates phosphorylase kinase directly. In addition, phosphorylase kinase is activated by phosphorylation by cAMP-dependent protein kinase (PKA). PKA is directly stimulated by cyclic AMP (cAMP), which is generated in response to the stimulation of adenylyl cyclase by various seven-transmembrane-segment [G-protein-coupled receptors](#) (e.g. receptors for epinephrine or glucagon). Dephosphorylation and inactivation of glycogen phosphorylase is catalyzed by several enzymes, the study of which led to the description and classification of protein phosphatases. See the text for

references. Pi, inorganic phosphate; , phosphate;  $\oplus$ , activation.



## 2. Phosphorylation in Prokaryotes

Studies of protein phosphorylation in prokaryotes initially lagged behind those in eukaryotic cells. However, the paramount role of regulation by protein phosphorylation is now recognized in both Archaea and Bacteria (9, 10). Bacteria possess ATP-dependent phosphorylation systems, using protein kinases that have sequence **homologies** with those of eukaryotes. These protein kinases phosphorylate [serine](#) or [threonine](#) residues and are, in many cases, regulated by cellular metabolites, rather than by second messengers. An original enzyme of bacteria is isocitrate dehydrogenase kinase/phosphatase of enteric bacteria, which has combined phosphatase and kinase activities located

on the same polypeptide chain (11). A second system of protein phosphorylation in eubacteria corresponds to the sensor kinase/response regulator, often termed the two-component system (12). In such systems, a sensor kinase, generally membrane-associated, phosphorylates itself on a [histidine](#) residue in response to external stimuli [eg, chemoattractant, nutrient, osmotic pressure, etc. (Fig. 1b)]. The phosphoryl group is then transferred to an [aspartic acid](#) residue located on a distinct protein or on the same polypeptide chain, the response regulator, triggering the appropriate biochemical response. This type of signal transduction seems to be evolutionarily very ancient, because it has now been also identified in Eucarya and Archea (10). In the third phosphorylation system, described thus far only in bacteria, the phosphoryl group donated by phosphoenolpyruvate (PEP) is passed down a chain of phosphorylatable proteins (9). Phosphorylation of these proteins occur on histidine and [cysteine](#) residues (Fig. 1c). Such PEP-dependent systems are implicated in a wide variety of regulations, including carbohydrate metabolism, [transcription](#), and [adenylate cyclase](#). Protein phosphatases are less well-characterized in prokaryotes. Relatives of the PPP phosphoserine/threonine protein phosphatases gene family have been identified in Archea and in eubacterial bacteriophages (10). Finally, it should be noted that bona fide tyrosine phosphorylation has been demonstrated in bacteria, as well as the presence of potential tyrosine phosphatases that bear a similar sequence signature to that of eukaryotic tyrosine phosphatases in their [active site](#) (10).

### 3. Phosphorylation in Eukaryotes

Regulatory protein phosphorylation in eukaryotes occurs primarily on serine, threonine, and tyrosine residues (Fig. 1a). The great majority of the protein kinases responsible for the phosphorylation of these residues have a conserved catalytic **domain** that belongs to a single gene family. It has been estimated that more than a thousand different protein kinases can be coded in the [genome](#) of higher eukaryotes (13, 14). The genome of *Saccharomyces cerevisiae* contains genes for 113 such conventional protein kinases (15). Interestingly, a number of lipid kinases have a catalytic domain related to that of conventional protein kinases. Some of these enzymes are capable of autophosphorylation; and at least one of them, DNA PK, is a genuine protein kinase, stimulated by the ends of double-stranded DNA (16). In addition, histidine kinases related to the bacterial two-component system have been identified in several eukaryotes (Fig. 1b). They are, for example, involved in the responses of yeast to osmotic shock (17) and of plants to [ethylene](#) (18). An original group of serine kinases related to histidine kinases corresponds to the mitochondrial protein kinase that phosphorylates  $\alpha$ -ketoacid dehydrogenases (11).

The catalytic domains of regulatory protein phosphatases in eukaryotes belong to several gene families: Serine/threonine phosphatases belong to two gene families of metal-requiring enzymes, PPP and PPM (19), whereas tyrosine and dual-specificity phosphatases form a superfamily of genes that have a common signature sequence in their active site (Cys-X<sub>5</sub>-Arg) and a common catalytic mechanism. These enzymes are grouped in four families: classical phosphotyrosine phosphatases (PTPs), low-molecular-weight PTPs, VH1-like dual-specificity phosphatases, and cdc25-like dual-specificity phosphatases (20).

### 4. Specificity of Protein Phosphorylation

Protein kinases display a high degree of substrate specificity. First, the vast majority phosphorylate either serine/threonine or tyrosine side chains, whereas only a few kinases have a dual specificity, for both threonine and tyrosine. In addition, protein kinases are highly sensitive to the immediate environment of the phosphorylated residue (21). Thus, specific **consensus sequences** have been identified for many of them (Table 2). In some cases, the specificity also depends on higher-order determinants, such as interactions between the protein substrate and the protein kinase at sites located at distance from the phosphorylated residue. In living cells, the proximity between the protein kinase and its substrate plays an important role. This is achieved by specific targeting mechanisms, which enrich the kinase concentration at particular sites of the cell (22, 23). There are also a number of scaffolding proteins that associate several enzymes, including kinases and/or

phosphatases and their substrates (24). In the case of protein phosphatases, the nature of the sequence surrounding the phospho-amino acid to be dephosphorylated seems generally less critical than in the case of protein kinases; and the role of higher-order determinants, as well as that of targeting processes, appears fundamental (25, 26). It should be pointed out that, with the exception of some monosubstrate systems, there is no perfect match between the substrate specificity of most protein kinases and phosphatases, implying a complex intertwining of the regulation by these enzymes.

**Table 2. Examples of Consensus Phosphorylation Sites for Protein Kinase**

|  |  |
|--|--|
| cAMP-dependent protein kinase                    | Arg/Lys-Arg-Xxx-Ser*                   |
| Ca <sup>2+</sup> -calmodulin-dependent kinase II | ArgXxxXxxSer*/Thr*                     |
| Protein kinase C                                 | Arg/Lys (Xxx)1-2 Ser*/Thr* Xxx-Arg/Lys |
| Casein kinase 1                                  | Acid-Xxx-Xxx-Ser*                      |
| Casein kinase 2                                  | Ser*/Thr* Xxx-Xxx Asp/Glu              |
| MAP-kinases, cyclin-dependent kinases            | Ser*/Thr* Pro                          |
| EGF-receptor tyrosine kinase                     | Glu/Asp-Tyr*-Hyd                       |

## 5. Methods to Study Protein Phosphorylation

### 5.1. General Methods to Study the State of Phosphorylation of Proteins

(detailed descriptions of methods can be found in Ref. 27-30). The most commonly used approach is metabolic **radiolabeling** of cells of interest with <sup>32</sup>P-orthophosphoric acid (sometimes called “front-phosphorylation”). <sup>32</sup>Pi is taken up by cells, incorporated into endogenous ATP (as well as in many other phosphorylated molecules), and used to phosphorylate proteins, lipids, and nucleic acids. Labeled phosphoproteins can be studied by [electrophoresis](#) (usually [two-dimensional gel electrophoresis](#)) or following **immunoprecipitation** with specific **antibodies**. However, proteins are often phosphorylated on several sites, which can have different physiological meanings, and other approaches are required to analyze independently the phosphorylation of each of them (see below). Moreover, the amount of <sup>32</sup>P incorporated in a given site does not reflect directly the stoichiometry of phosphorylation of this site, but rather the turnover rate of the phosphoryl group. That is, a site phosphorylated at a high stoichiometry but with a low turnover rate (eg, “constitutive phosphorylation”) will appear poorly labeled, whereas a residue phosphorylated at a low stoichiometry, but with a high turnover rate, will appear intensely labeled. The levels of labeling are also dependent on the specific activity of the intracellular ATP pools, which may vary, depending on the cell population or physiological status. Therefore, other methods have been developed to avoid these pitfalls. In “back-phosphorylation” assays, proteins of interest are isolated in conditions in which their phosphorylation state is preserved. They are then phosphorylated stoichiometrically *in vitro* with g<sup>32</sup>P-ATP and a purified kinase, active on the phosphorylation site of interest. Thus, the higher the level of phosphorylation of the site *in vivo*, the lesser radioactive phosphate is incorporated *in vitro*. The use of this method is limited, however, by the requirement for a highly active purified protein kinase and by its low sensitivity. An alternative approach that is enjoying considerable success is the use of phosphorylation-state-specific antibodies. Some of these

antibodies react with a phosphorylated amino acid, almost independently of its sequence environment. For example, antibodies specific for phospho-tyrosine are extremely useful tools in the study of protein tyrosine phosphorylation. Moreover, specific antibodies can be raised against synthetic phosphopeptides encompassing the phosphorylated site of a protein of interest. Such antibodies are specific for a given protein phosphorylated at a particular site, and they provide powerful tools that can be used for immunoblotting, immunocytochemistry, and other approaches.

Phosphoproteins are often studied by one- or two-dimensional [electrophoresis](#). In nondenaturing conditions (eg, **isoelectrofocusing**), the presence of a phosphorylated residue shifts the electrophoretic mobility of the protein toward a more acidic form. In the presence of **SDS**, phosphorylation of a protein may have no effect on its electrophoretic mobility, or it may increase its mobility and its apparent size. The slowing of the protein mobility by phosphorylation during [SDS-PAGE](#) is often attributed to a decrease in the charge density due to decreased binding of SDS to the phosphorylated protein. In contrast, it is noteworthy that phosphorylation at specific sites may increase the mobility of some rare proteins in SDS-PAGE.

## 5.2. Methods for Identifying the Phosphorylated Amino Acid(s)

Following prelabeling with  $^{32}\text{P}$  or *in vitro* phosphorylation with  $^{32}\text{P}$ -ATP, the nature of the phosphorylated residues can be determined by acid hydrolysis of the [peptide bonds](#) and separation of phospho-serine, phospho-threonine, and phospho-tyrosine by [thin-layer electrophoresis](#). Phospho-histidine and phospho-aspartate are extremely acid-labile and cannot be identified by this approach. On the other hand, the resistance of phospho-tyrosine to alkaline pH is used to study this phospho-amino acid preferentially. The phosphopeptides generated by partial **proteolysis** of a  $^{32}\text{P}$ -labeled phosphoprotein with specific endoproteinases can be studied by two-dimensional phosphopeptide maps (usually a combination of thin-layer electrophoresis and chromatography). Such two-dimensional phospho-**peptide maps** are very useful for comparing the sites phosphorylated on the same protein under various circumstances, *in vivo* or *in vitro*. Phosphopeptides can also be separated by high-performance liquid chromatography (HPLC) and, if sufficient amounts of material are recovered, submitted to automatic [protein sequencing](#), allowing the identification of the phosphorylated residue and its surrounding protein sequence. The use of powerful technologies of [matrix-assisted laser desorption/ionization](#) (MALDI) [mass spectrometry](#) now provides an interesting alternative approach to identify phosphorylation and other post-translational modifications. [Site-directed mutagenesis](#) of individual residues (replacement of a “phosphorylatable” residue by a “nonphosphorylatable” one—for example, replacement of a serine or a threonine by an [alanine](#), or of a tyrosine by a [phenylalanine](#)) is often used to test the phosphorylation of a precise residue. Although this approach is extremely useful, it should be kept in mind that the information obtained is indirect and that such mutations may have consequences on other properties of the protein other than preventing phosphorylation of the mutated residue.

## 6. Effects of Phosphorylation on the Properties of Proteins

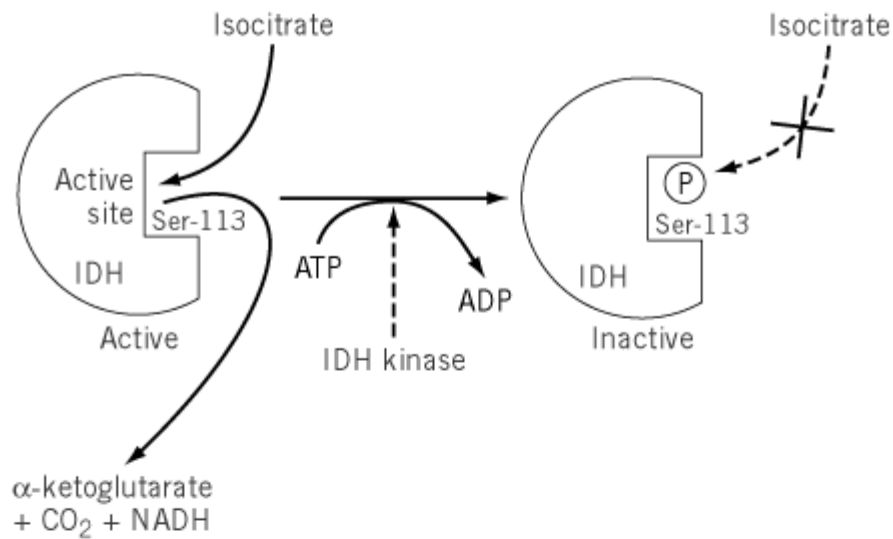
The effects of phosphorylation on the properties of many proteins are known to be functionally important. They include the activation or inhibition of enzymes, the opening or closing of **ion channels**, the increase or decrease in the activity of [transcription factors](#), the aggregation or disassembly of **cytoskeletal** components, and many other effects. There are only a few cases, however, in which the precise molecular mechanism by which phosphorylation of a specific residue brings about the functional changes observed is known in detail. The consequences of the phosphorylation of an amino acid residue can be only local, or they can be amplified in the three-dimensional [protein structure](#). Phosphorylation can also alter the interactions of a protein with others. A phosphate group being bulky and highly charged (two negative charges at physiological pH), its presence alters dramatically the properties of the amino acid side chain to which it is bound. It is often attempted to mimic the consequence of phosphorylation of a residue by its mutation to a negatively charged amino acid (**aspartate** or **glutamate**). In some cases, such mutations reproduce some of the effects of phosphorylation, while they prevent the normal phosphorylation of the

corresponding residues. However, the size and the charge of the carboxyl group are less than those of a phosphoryl group, accounting for the many failures of this approach.

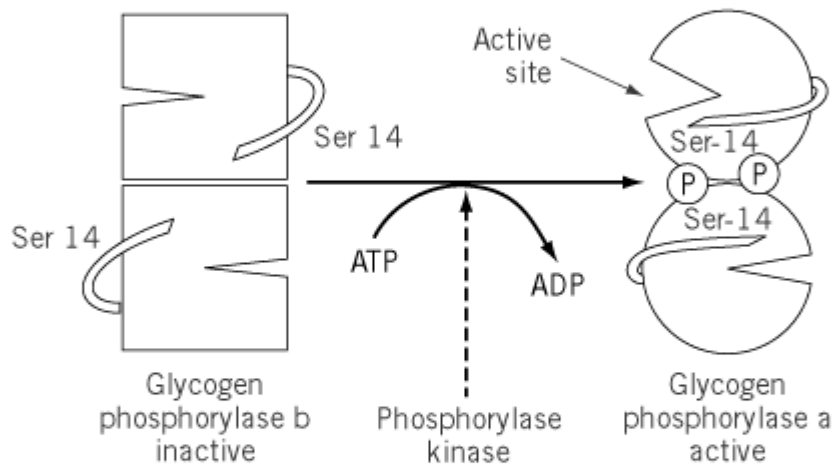
### 6.1. Local Effects of Phosphorylation on the Properties of Proteins

One example of a protein in which phosphorylation modifies dramatically the activity of an enzyme by a local action on the active site is isocitrate dehydrogenase (IDH), an enzyme regulated by phosphorylation in plants and bacteria. In the unphosphorylated state, IDH is active and catalyses the oxidative decarboxylation of isocitrate to  $\alpha$ -ketoglutarate, a critical step in the Krebs citric acid cycle. In the presence of high levels of ATP, IDH is switched off by phosphorylation, and citrate is funneled to the glyoxylate cycle, a biosynthetic pathway that allows plants and bacteria to grow on acetate. *Escherichia coli* IDH is phosphorylated on Ser113 by an enzyme that is also able to dephosphorylate the same residue (a tandem kinase/phosphatase; see above). Ser113 is located in the substrate binding site, near the active site of the enzyme (31). Phosphorylation of this serine prevents isocitrate binding by steric hindrance and electrostatic repulsion (Fig. 3a).

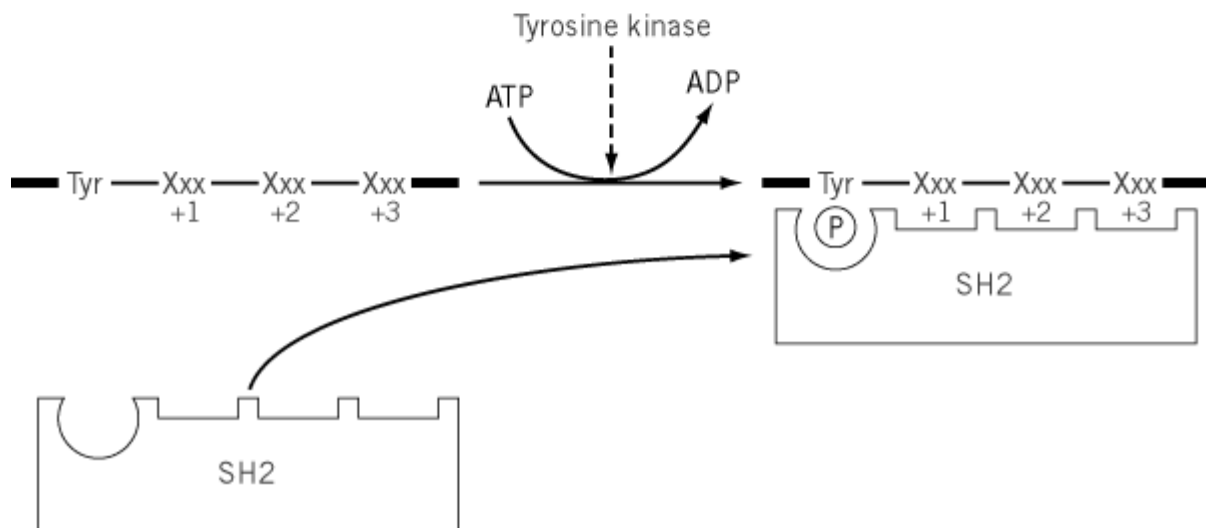
**Figure 3.** Effects of phosphorylation on proteins. (a) Local effects of phosphorylation on isocitrate dehydrogenase (IDH). IDH transforms isocitrate into  $\alpha$ -ketoglutarate, with concomitant production of  $\text{CO}_2$  and reduction of  $\text{NAD}^+$  in NADH. In *E. coli* IDH, phosphorylation of serine-113 in the active site prevents the binding of the substrate and inactivates the enzyme. (b) Allosteric activation of glycogen phosphorylase by phosphorylation. Phosphorylase is a homodimer; phosphorylation of Ser14 in both subunits leads to rearrangement of the amino-terminal region, modification of the subunit interaction, and changes in the tertiary and quaternary structure of the enzyme, making the active site more accessible to the substrate. (c) Regulation of SH2-domain binding by tyrosine phosphorylation. When particular protein sequences are phosphorylated on tyrosine residues, they bind with high affinity to specific SH2 domains. The presence of phosphotyrosine is necessary for the binding, but the specificity of the interaction is achieved by the recognition by the SH2 domain of the residues on the carboxy-terminal side of the phosphorylated tyrosine (Xxx +1, +2, +3). See the text for references.  $\text{P}$ , phosphate.



(a)



(b)



(c)

Another example of a well-studied local effect of phosphorylation is provided by the structural studies of protein kinases themselves (32). For the protein kinases to be active, a peptide loop in the

vicinity of the active site, termed the *activation loop* or *T loop*, must be correctly positioned. In many protein kinases, the correct positioning of the activation loop requires its phosphorylation on one or two residues. Interestingly, in some cases the phosphorylation appears constitutive (eg, protein kinase A), whereas in others (eg, mitogen-activated protein kinases, [MAP kinases](#)) the phosphorylation of the activation loop is a fundamental regulatory step, catalyzed by a specific activating protein kinase.

## 6.2. Effects of Phosphorylation on the Overall Structure of Proteins

Phosphorylation of a single residue can have dramatic consequences on the [tertiary structure](#) and [quaternary structure](#) of a protein, as demonstrated by the study of the [X-ray crystallography](#) structure of glycogen phosphorylase in the phosphorylated and nonphosphorylated states ([33](#)). Phosphorylase is a homodimer, and phosphorylation occurs on a serine in the amino-terminal region of each peptide chain (Ser14), located strategically at the subunit interface. Phosphorylation of this serine leads to the reorganization of the amino terminus of each subunit and the modification of its interaction with the other subunit. As a consequence, the enzyme is stabilized in a form in which the active site, located at a distance from the phosphorylated serine, becomes more accessible to the substrate (Fig. [3b](#)). In this case, the effects of phosphorylation on the structure of the enzyme are comparable to those occurring during classical **allosteric** activation, except that the effects of phosphorylation are due to a covalent modification of the enzyme, and not to ligand binding.

## 6.3. Effects of Phosphorylation on Protein–Protein Interactions

Protein phosphorylation can have dramatic regulatory effects on [protein–protein interactions](#), as disclosed by the study of protein tyrosine phosphorylation. Indeed, a 100-residue domain first identified in Src, after which it was termed Src-homology 2 (**SH2**), and then found in many other proteins, was shown to bind specific peptides phosphorylated on tyrosine ([34](#), [35](#)). Phosphorylation of the tyrosine residue is necessary for binding, but the specificity of the interaction is achieved by the recognition of additional residues, located on the carboxy-terminal side of the phosphorylated tyrosine, by the SH2 domain (Fig. [3c](#)). Binding of SH2 domains to specific peptide sequences phosphorylated on tyrosine residues allows the clustering of enzymes directly or via adapter proteins, around activated growth-factors receptors. Such enzymes, including **phosphatidylinositol-3-kinase**, **phospholipase C $\alpha$** , or **guanine-nucleotide exchange factors**, are thus brought into the vicinity of their membrane-associated substrates. SH<sub>2</sub>-mediated clustering triggers the cascades of reactions responsible for the effects of growth factors. SH2 domains can also have very different functions, such as maintaining a protein in an inactive state, as in the case of Src-family tyrosine kinases ([36](#), [37](#)). A carboxy-terminal phosphorylated tyrosine is involved in an intramolecular interaction with the SH2 domain, placing the enzyme in a closed conformation unable to reach its substrates. Activation results from dephosphorylation of the carboxy-terminal tyrosine, or by displacement of the SH2 domain by a competing phosphopeptide. Additional domains that interact specifically with phosphorylated proteins have been identified, including PTB domains (phosphotyrosine-binding) ([38](#)) and 14–3–3 proteins, whose interaction with their partner proteins appears to require phosphorylation of the latter on a serine residue ([39](#)).

## 7. The Place of Protein Phosphorylation in the Signal Transduction Networks

Many signal transduction systems in cells use several modules whose hierarchical organization allows powerful amplification and precise modulation. Protein phosphorylation is involved in most of these transduction systems, either as an output mechanism, by altering the properties of target proteins, or as a module of information processing. Indeed, phosphorylation reactions in cells often occur in cascades in which activation of a first protein kinase phosphorylates and activates a second protein kinase, which, in turn, phosphorylates and activates a third kinase, and so on. This situation may account in part for the very large number of protein kinases, and it apparently has a number of interesting properties ([40](#), [41](#)). First, phosphorylation cascades allow amplification of the signal: One molecule of the first activated kinase can phosphorylate several secondary kinases, which can each phosphorylate several tertiary kinases, and so on. Another property of these cascades is to acquire kinetic properties, such as positive or negative **cooperativity**, similar to those of allosteric enzymes,

even if the kinases composing the cascade themselves follow classical [Michaelis–Menten kinetics](#) . Thus, protein kinase cascades can convert graded inputs into switch-like outputs (42). Finally, the kinase cascades provide many levels of regulation and possible cross-talks between different regulatory pathways. Thus, it is likely that the organization of protein kinases in cascades confers strong evolutionary advantages, as attested by the very high degree of conservation of many of these cascades.

## 8. Examples of Phosphorylation Dysfunction in Human Diseases

As mentioned above, virtually all intracellular biological processes are regulated to some extent by protein phosphorylation. Thus, it is not surprising that dysfunction of these reactions may be responsible for some pathological states, including in humans. Protein kinases and phosphatases regulate cell growth, and a number of mutations leading to the uncontrolled activation of phosphorylation pathways are oncogenic (43, 44). Several oncogenes are tyrosine kinases (eg, v-Src, v-Fyn, v-ErbB, Met) or serine/threonine kinases (eg, Raf1, Mos) that are activated by mutation, whereas it is the loss of function of protein tyrosine and phospholipid phosphatase PTEN that is oncogenic (45). Loss of function of protein kinases is also responsible for various genetic diseases, including immune deficiency by the absence of  $\gamma$ -globulins [mutation of a nonreceptor tyrosine kinase in Bruton's disease (46)], intestinal palsy [mutation of a receptor tyrosine kinase Ret in Hirschprung's disease (47)], or myotonic dystrophy [mutation of a serine/threonine kinase (48)]. Conversely, protein phosphatases are used as weapons by some pathogenic organisms. For example, one of the virulence genes of *Yersinia* bacteria, including the agent of bubonic plague, is a phosphotyrosine phosphatase (49). Some environmental [toxins](#) exert their effects by stimulating protein kinases irreversibly (eg, stimulation of [protein kinase C](#) by oncogenic **phorbol esters**) or inhibiting protein phosphatases (inhibition of serine/threonine phosphatases by okadaic acid or microcystins, which are responsible for food poisoning) (50). On the other hand, the immunosuppressant drugs that have allowed the recent progress in organ grafting, [cyclosporin A](#) and [FK506](#), act by inhibiting calcineurin, a  $\text{Ca}^{2+}$ /calmodulin-dependent protein phosphatase (51). Thus, the large number of protein kinases and phosphatases, and their involvement in numerous biological functions, has prompted many investigators to try to develop specific inhibitors or activators of these enzymes as potential therapeutic tools.

## Bibliography

1. E. G. Krebs (1993) *Angew. Chem. (Engl.)* **32**, 1122–1129.
2. E. H. Fischer (1993) *Angew. Chem. (Engl.)* **32**, 1130–1137.
3. J. F. Kuo and P. Greengard (1970) *J. Biol. Chem.* **245**, 2493–2498.
4. M. B. Kennedy and P. Greengard (1981) *Proc. Natl. Acad. Sci. USA* **78**, 1293–1297.
5. Y. Takai, A. Kishimoto, Y. Iwasa, Y. Kawahara, T. Mori, and Y. Nishizuka (1979) *J. Biol. Chem.* **254**, 3692–3695.
6. P. Cohen (1994) *Bioessays* **16**, 583–588.
7. T. Hunter (1996) *Biochem. Soc. Trans.* **24**, 307–327.
8. H. Ushiro and S. Cohen (1980) *J. Biol. Chem.* **255**, 8363–8365.
9. M. H. Saier, Jr., L. F. Wu, and J. Reizer (1990) *Trends Biochem. Sci.* **15**, 391–396.
10. P. J. Kennelly and M. Potts (1996) *J. Bacteriol.* **178**, 4759–4764.
11. K. M. Popov, J. W. Hawes, and R. A. Harris (1997) *Adv. Second Messenger Phosphoprotein Res.* **31**, 105–111.
12. L. A. Alex and M. I. Simon (1994) *Trends Genet.* **10**, 133–138.
13. T. Hunter (1987) *Cell* **50**, 823–829.
14. S. K. Hanks and T. Hunter (1995) *FASEB J.* **9**, 576–596.
15. T. Hunter and G. D. Plowman (1997) *Trends Biochem. Sci.* **22**, 18–22.
16. K. O. Hartley, D. Gell, G. C. M. Smith, et al. (1995) *Cell* **82**, 849–856.



17. T. Maeda, S. M. Wurgler-Murphy, and H. Saito (1994) *Nature* **369**, 242–245.
18. C. Chang, S. F. Kwok, A. B. Bleecker, and E. M. Meyerowitz (1993) *Science* **262**, 539–544.
19. D. Barford (1996) *Trends Biochem. Sci.* **21**, 407–412.
20. E. B. Fauman and M. A. Saper (1996) *Trends Biochem. Sci.* **21**, 413–417.
21. B. E. Kemp and R. B. Pearson. (1990) *Trends Biochem. Sci.* **15**, 342–346.
22. D. Mochly-Rosen (1995) *Science* **268**, 247–251.
23. M. C. Faux and J. D. Scott (1996) *Trends Biochem. Sci.* **21**, 312–315.
24. T. Pawson and J. D. Scott (1997) *Science* **278**, 2075–2080.
25. P. J. Kennelly and E. G. Krebs (1991) *J. Biol. Chem.* **266**, 15555–15558.
26. M. J. Hubbard and P. Cohen (1993) *Trends Biochem. Sci.* **18**, 172–177.
27. T. Hunter, ed. (1991) *Methods Enzymol.* **200**.
28. T. Hunter, ed. (1991) *Methods Enzymol.* **201**.
29. A. C. Nairn and S. Shenolikar (1995) *Neuroprotocol*, **6**, 2–111.
30. A. A. Boulton, G. B. Baker, and H. C. Hemmings, Jr., eds, (1996) *Neuromethods: Post-translational Modifications* Humana Press, Totowa, NJ.
31. J. H. Hurley, A. M. Dean, J. L. Sohl, D. E. Koshland, Jr., and R. M. Stroud (1990) *Science* **249**, 1012–1016.
32. L. N. Johnson, M. E. M. Noble, and D. J. Owen (1996) *Cell* **85**, 149–158.
33. L. N. Johnson (1992) *FASEB J.* **6**, 2274–2282.
34. L. E. M. Marengere and T. Pawson (1994) *J. Cell Sci.* **107**, 97–104.
35. T. Pawson (1995) *Nature* **373**, 573–580.
36. W. Q. Xu, S. C. Harrison, and M. J. Eck (1997) *Nature* **385**, 595–602.
37. F. Sicheri, I. Moarefi, and J. Kuriyan (1997) *Nature* **385**, 602–609.
38. M. M. Zhou, K. S. Ravichandran, E. T. Olejniczak, et al. (1995) *Nature* **378**, 584–592.
39. A. J. Muslin, J. W. Tanner, P. M. Allen, and A. S. Shaw (1996) *Cell* **84**, 889–897.
40. E. Shacter, E. R. Stadtman, S. R. Jurgensen, and P. B. Chock (1988) *Methods Enzymol.* **159**, 3–19.
41. S. Swillens, J. M. Boeynaems, and J. E. Dumont (1988) *Methods Enzymol.* **159**, 19–27.
42. J. E. Ferrell, Jr. (1996) *Trends Biochem. Sci.* **21**, 460–466.
43. G. A. Rodrigues and M. Park (1994) *Curr. Opin. Genet. Dev.* **4**, 15–24.
44. T. Hunter (1997) *Cell* **88**, 333–346.
45. J. Li, C. Yen, D. Liaw, et al. (1997) *Science* **275**, 1943–1947.
46. S. Tsukada, D. J. Rawlings, and O. N. Witte (1994) *Curr. Opin. Immunol.* **6**, 623–630.
47. G. Romeo, P. Ronchetto, Y. Luo, et al. (1994) *Nature* **367**, 377–378.
48. Y.-H. Fu, D. L. Friedman, S. Richards, et al. (1993) *Science* **260**, 235–238.
49. K. Guan and J. E. Dixon (1990) *Science* **249**, 553–556.
50. C. MacKintosh and R. W. MacKintosh (1994) *Trends Biochem. Sci.* **19**, 444–448.
51. J. Liu, J. D. Farmer, W. S. Lane, J. Friedman, I. Weissman, and S. L. Schreiber (1991) *Cell* **66**, 807–815.
52. R. B. Pearson and B. E. Kemp (1991) *Methods Enzymol.* **200**, 62–81.
53. Z. Songyang, K. P. Lu, Y. T. Kwon, et al. (1996) *Mol. Cell. Biol.* **16**, 6486–6493.

## Phosphotransferase

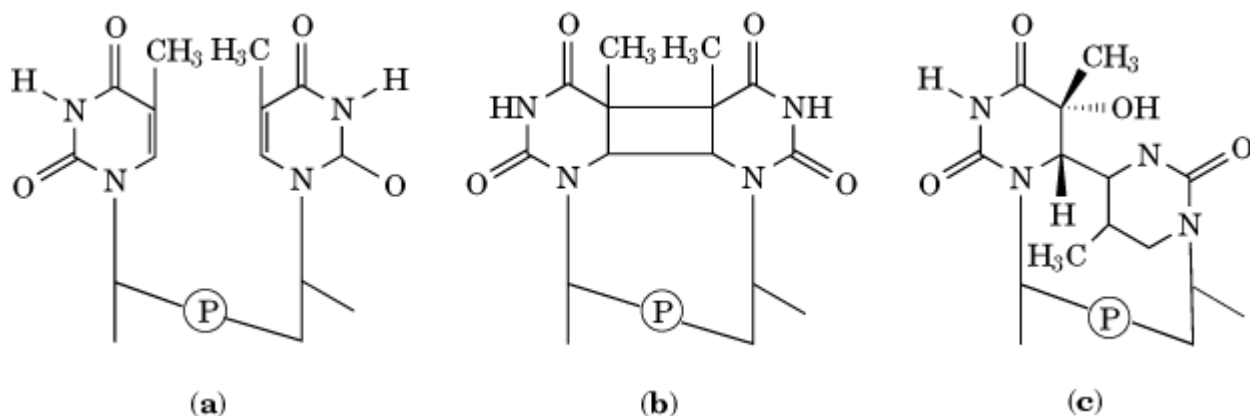
Phosphotransferases form a large group of [enzymes](#) that are involved with the transfer of a  $\gamma$ -phosphoryl group from ATP, or another nucleoside triphosphate, to an acceptor molecule (see [Kinase](#)). They are divided into groups according to the identity of the phosphoryl group acceptor. These acceptors can be alcohols and carboxylic acids, as well as nitrogenous and phosphorylated compounds. All phosphotransferases have an essential requirement for a bivalent metal ion (see [Metal-Requiring Enzymes](#)).

## Photolyase/Photoreactivation

Photoreactivation is the reversal of the effects of ultraviolet (UV) light (200 to 300 nm) by simultaneous or subsequent exposure of the organism to near UV-visible (300 to 500 nm) light ([1](#), [2](#)). Photoreactivation prevents the induction by UV light of growth delay, [mutation](#), and [cell death](#). Although several factors contribute to the phenomenon of photoreactivation, the main mechanism involved is the reversal of UV-induced pyrimidine dimers in DNA by the enzyme called photoreactivating enzyme, or DNA photolyase ([3](#)). DNA photolyase is the only known enzyme to employ light energy to catalyze a chemical reaction ([3](#)).

The two major lesions induced in DNA by UV light are cyclobutane pyrimidine dimers (Pyr $\leftrightarrow$ Pyr) and [6–4] photoproducts (Pyr[6–4]Pyr) (Fig. [1](#)). There are separate photolyases that are specific for Pyr $\leftrightarrow$ Pyr and for [6–4] photoproducts ([4](#), [5](#)), but no known enzyme can act on both lesions. The structures and reaction mechanisms of the two types of photolyases are similar and may be summarized as follows. Each consists of a single polypeptide chain of 55 to 65 kDa containing a chromophore that acts as a photoantenna (folate or deazaflavin) and a catalytic cofactor (FADH<sub>2</sub>) ([6-9](#)). It binds to the photodamage in DNA in a light-independent manner, producing a relatively stable enzyme-substrate complex. Upon exposure to light, the enzyme breaks the bonds that link the two pyrimidines to one another, generates standard canonical bases as products, and then dissociates from the repaired DNA. In contrast to many other enzymatic reactions that involve **electron transfer**, catalysis by photolyase involves a cyclic electron transfer and is not a **redox** reaction, in that upon completion of the enzymatic reaction the redox states of the enzyme and the substrate remain the same ([6](#)).

**Figure 1.** The two major photoproducts produced in DNA by ultraviolet light from (a) thymine dinucleotide: (b) the cyclobutane thymine dimer, and (c) the [6–4] photoproduct. Similar photoproducts arise from other adjacent dipyrimidines in DNA.



## 1. Evolution

Photolyases are flavoproteins, yet the primary sequences of the apoenzymes do not reveal any obvious [homology](#) to the signature sequences found in other flavoproteins (4). Instead, photolyases have remarkable sequence homology to blue-light photoreceptors (10), which regulate circadian rhythms in organisms ranging from mustard to humans (11-13). It appears that the necessity for the flavin cofactor to function from the light-excited state has imposed unique constraints on the flavin binding site of photolyase/blue-light photoreceptor family of proteins, such that the active-site geometry of these proteins is different from all other flavoproteins that carry out catalysis with ground-state flavin.

The distribution of the members of the photolyase/blue-light photoreceptor family in the biological world does not fit a readily defined pattern of evolutionary selection. The coliform bacterium *Escherichia coli* contains photolyase, but the soil bacterium *Bacillus subtilis* does not. The budding yeast *Saccharomyces cerevisiae* has the enzyme, but the fission yeast *S. pombe* does not. *Drosophila melanogaster* possesses cyclobutane pyrimidine dimer photolyase, [6-4] photolyase, and blue-light photoreceptor (14), while humans lack photolyase but express two blue-light photoreceptor proteins. Table 1 shows the expression pattern of this class of proteins in select members of the “three kingdoms of life”: bacteria, archaea, and eukarya. No general rules emerge from this table. The only rules that appear to have arisen from investigations of these groups of proteins are the following. (i) Cyclobutane photolyase has been detected in bacteria, plants, and animals through marsupial animals (15), but it has not been found in placental mammals such as human (16) and mouse. (ii) The [6-4] photolyase has been found in plants, insects (*Drosophila*), and some vertebrates (*Xenopus*, rattlesnake), but not in mammals such as human or mouse, nor in bacteria. (iii) Blue-light photoreceptors have been found in plants (*Arabidopsis*, mustard) and animals (*Drosophila*, human, and mouse), but not in bacteria. (iv) Of all bacterial and eukaryotic **viruses** investigated, photolyase has, so far, been found only in a grasshopper virus.

**Table 1. Distribution of the Photolyase/Blue-Light Photoreceptor Family of Proteins in the Biological World**

| Bacteria       |                    | Archaea              |                     |                      | Eucarya                |                   |
|----------------|--------------------|----------------------|---------------------|----------------------|------------------------|-------------------|
| <i>E. coli</i> | <i>B. subtilis</i> | <i>M. thermoauto</i> | <i>M. janaschii</i> | <i>S. cerevisiae</i> | <i>D. melanogaster</i> | <i>H. sapiens</i> |

|              |     |    |     |    |     |     |     |
|--------------|-----|----|-----|----|-----|-----|-----|
| Photolyase   | Yes | No | Yes | No | Yes | Yes | No  |
| [6–4]        | No  | No | No  | No | No  | Yes | No  |
| Photolyase   |     |    |     |    |     |     |     |
| Cryptochrome | No  | No | No  | No | No  | Yes | Yes |

---

The tissue distribution of photolyase in animals that possess the gene is rather surprising. It is highly expressed in *Drosophila* ovaries, where it is expected to protect germ cells from UV damage. However, it is also expressed in all internal organs of opossum, including high levels of expression in the brain. Considering the nocturnal nature of the opossum, and the low probability even in bright daylight of UV and photoreactivating light penetrating the skull of the opossum, this expression pattern remains enigmatic.

## 2. Structure

Photolyases consist of a single polypeptide of 500 to 700 amino acid residues and two noncovalently bound cofactors/chromophores. One of the chromophores is always flavin adenine dinucleotide (FAD), and the other is either a pterin (methenyltetrahydrofolate in the case of *E. coli* and *S. cerevisiae* photolyases) or deazariboflavin. Deazariboflavin is considered an ancient molecule, and photolyases with this cofactor have been found only in a few organisms that are capable of synthesizing this uncommon cofactor: *Aspergillus nidulans*, *Streptomyces griseus*, and *M. thermoautotrophicum*. Most other photolyases characterized to date, including those from *E. coli*, *S. cerevisiae*, *D. melanogaster*, and marsupials, contain a pterin as the second chromophore. Only two [6–4] photolyases have been characterized in detail, those of *D. melanogaster* and *X. laevis*, and both contain a pterin as the second chromophore (17). Blue-light photoreceptors from mustard, *A. thaliana*, and humans contain FAD and a pterin (11, 12).

The structures of *E. coli* and *A. nidulans* photolyases have been determined (18, 19). Even though the two enzymes are only 39% identical in sequence, they have remarkably similar structures, such that the traces of C<sup>α</sup> atoms are superimposable. The enzymes consist of two well-defined domains: an N-terminal a/b domain, which is connected to the C-terminal **alpha-helical** domain through a long inter-domain loop. The folate or the deazariboflavin chromophores are located in a crevice separating the two domains. In the *E. coli* enzyme, the folate is close to the surface of the protein and hence is not tightly bound to the enzyme; furthermore, its orientation is not optimal for energy transfer to the FAD cofactor. In contrast, the deazariboflavin chromophore in *A. nidulans* photolyase is deeply buried between the two domains, and it has a more favorable orientation relative to FAD for efficient energy transfer between the chromophores. Thus even though center-to-center distances between the two chromophores in the two photolyases are about 17 Å, the energy transfer from the second chromophore to FAD is more efficient (97% to 70% ) in *A. nidulans* photolyase than in that of *E. coli* (20). The helical domain is made up exclusively of  $\alpha$ -helices that are packed in two rather compact clusters. The FAD is buried between these two clusters, and the flavin and adenine rings of the cofactor are in the unusual *cis* conformation in both photolyases.

The crystal structures of both enzymes also reveal how they bind to DNA. A plot of the surface potential of the enzyme reveals a longitudinal groove traversing both the a/b and the  $\alpha$ -helical domains that is paved with positive charges. At its center is a hole that has the shape and size appropriate for cyclobutane pyrimidine dimer. The bottom of the hole is made up by the flavin ring. These structural features suggest that the DNA lies in the positively charged groove through ionic interactions between the DNA backbone and the positive charge of the side chains and that the dimer “flips out” of the duplex into the hole within the enzyme. Following repair, the flipped-out dinucleotide would undergo a conformational change upon transformation of the dimer into two

monomers, which would force the bases out of the hole and assist the repaired DNA to dissociate from the enzyme. Biochemical data indicate that [6–4] photolyase also employs a “flip-out” mode of complex formation during catalysis.

### 3. Reaction Mechanism

Photolyase works by the classical [Michaelis–Menten kinetic](#) mechanism; that is, the enzyme and substrate form a noncovalent complex that is converted to an enzyme–product complex after catalysis and eventually dissociates to yield product and free enzyme. The light-dependence of this catalysis step is unique to photolyase.

#### 3.1. Binding

Photolyase, unlike some other high-specificity [DNA-binding proteins](#), does not rely on diffusion along the linear DNA molecule to recognize its substrate. Similarly, the enzyme, in contrast to most proteins that bind to DNA in a sequence-specific manner and hence require double-stranded DNA, binds to its cognate lesions in single- and double-stranded DNA with almost equal affinities (21). This property of the enzyme, combined with the fact that the enzyme can bind the photoproduct even in a dinucleotide form and repair it with high affinity, are consistent with a model in which the enzyme binds mostly to the backbone of the damaged strand and flips-out the photolesion into the catalytic hole for repair.

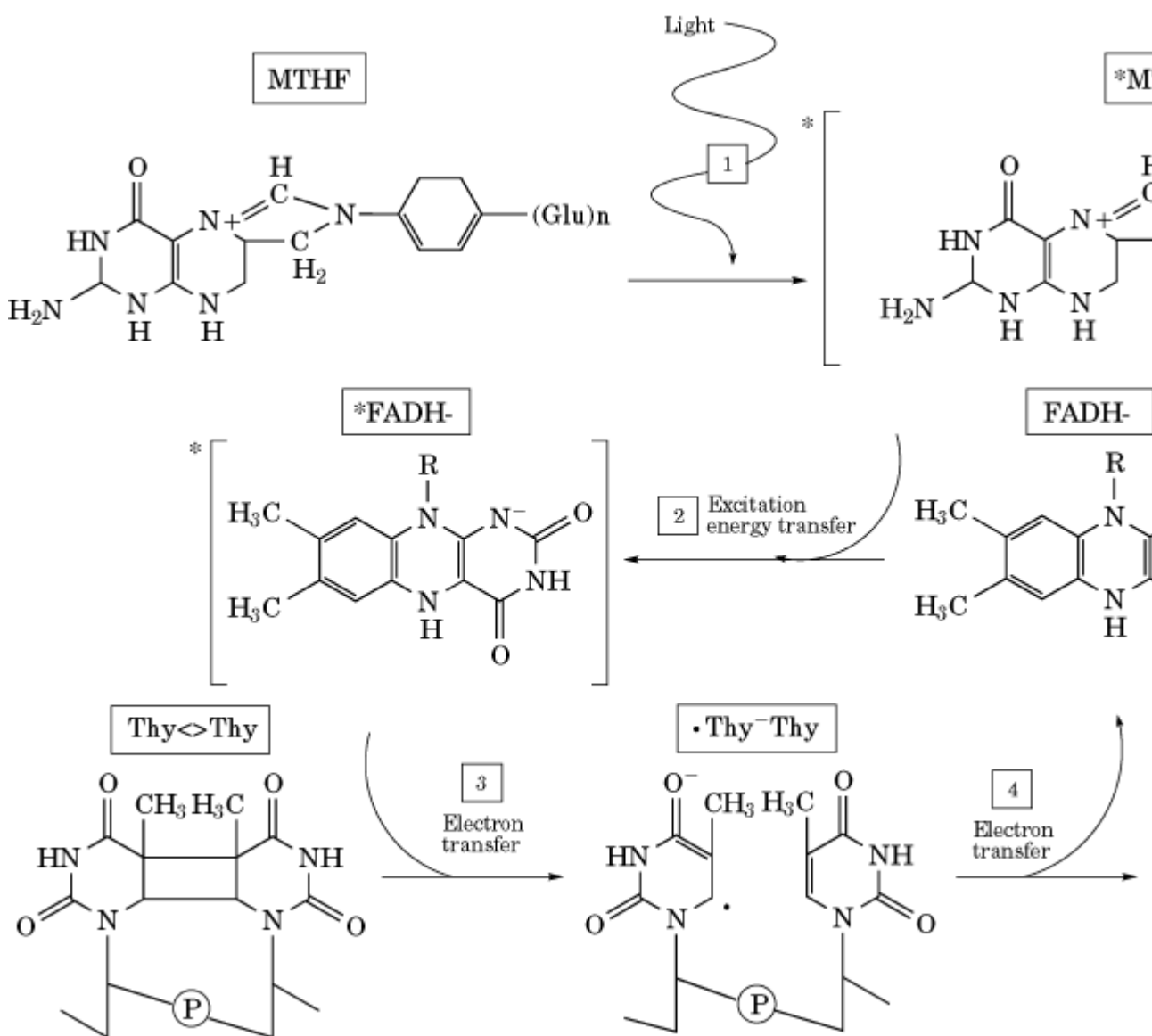
#### 3.2. Catalysis

A most extraordinary aspect of catalysis by photolyase is its absolute dependence on light. According to the Woodward–Hoffman rule on symmetry in molecular orbitals, the symmetry of molecular orbitals in a chemical reaction is conserved between reactants and products. A corollary of this rule is that photochemically allowed reactions that proceed from the lowest unoccupied molecular orbitals (LUMO) are thermally forbidden; and conversely, thermally allowed reactions that proceed from the highest-occupied molecular orbitals (HOMO) are photochemically forbidden. In line with these rules, the cyclobutane ring of Pyr◊Pyr is stable at pH 1 and 160 °C for 6 h, but it is reversed to pyrimidine monomers with 240-nm light with a quantum yield of about 1.0. Surprisingly, however, even though photolyase utilizes light energy to catalyze the repair, the reaction is not a photochemical reaction in the strict sense, because it does not involve photoexcitation of the pyrimidine dimer, either directly or indirectly. Consequently, splitting of Pyr◊Pyr by photolyase does not break the rule of conservation of orbital symmetry. The enzyme accomplishes this feat by converting the Pyr◊Pyr into a radical, which then readily proceeds, thermally, to monomers. Light (300 to 500 nm) is used simply to excite the cofactor of the enzyme into a strong reductant capable of donating an electron to the relatively inert Pyr◊Pyr.

Having thus stated the general features of catalysis, the specifics for Pyr◊Pyr photolyase can be summarized as shown in Figure 2. The enzyme binds to its substrate independent of light, and the reaction proceeds upon exposure to light as follows (6). First, a photon is absorbed by the photoantenna (pterin or deazaflavin) to excite the chromophore to the excited singlet state. Second, the photoantenna transfers the energy by dipole–dipole interactions to FADH<sub>2</sub> with either 70% (pterin) or 97% (deazaflavin) efficiency. Third, the excited singlet state of FADH<sub>2</sub> transfers an electron to Pyr◊Pyr with about 100% efficiency, generating flavin neutral radical and Pyr◊Pyr radical anion. Fourth, the Pyr◊Pyr radical anion splits to pyrimidine monomers, concomitant with electron transfer back to restore the flavin to its catalytically active state. The reaction is a photoinitiated cyclic electron transfer, which results in bond rearrangement without changing the redox status of the reactants.

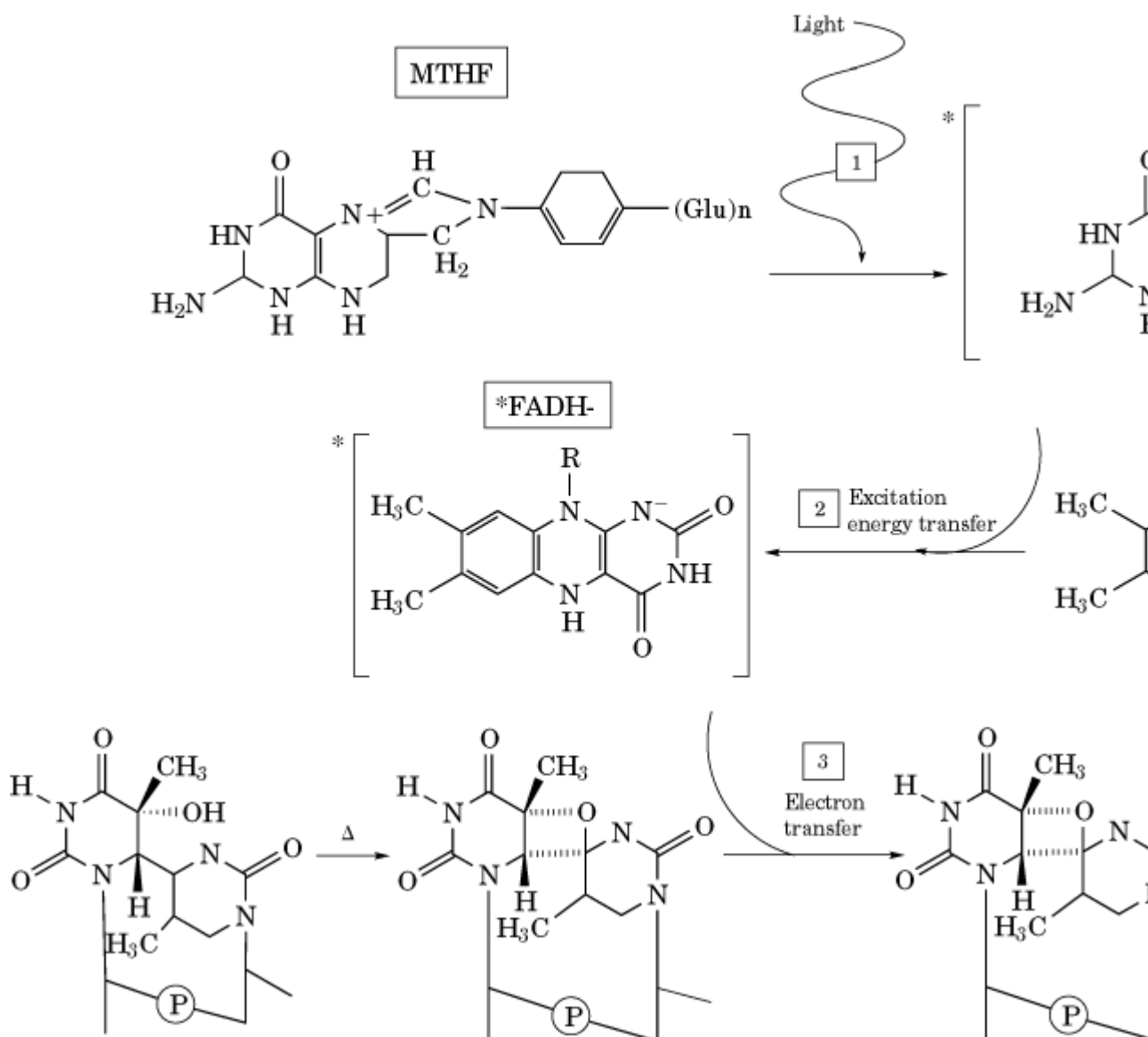
**Figure 2.** Photoreactivation photocycle for photolyase. The active site of the enzyme complexed to substrate is shown. (1) Pterin absorbs a near-UV photon. (2) The excited singlet-state pterin transfers energy to FADH<sup>•</sup> by dipole–dipole intera

FADH<sup>-</sup> transfers an electron to the thymine dimer (4). The thymine dimer anion radical collapses to two thymines, concomitant with the flavin neutral radical.



The [6-4] photolyase functions in a similar manner (17, 22), with one important exception. Formation of [6-4] photoproduct encompasses joining of C6 of the 5' base to the C4 of the 3' base via a sigma bond, accompanied by transfer of the group at the C4 position (-NH or -OH) of the 3' base to the C5 position of the 5' base. Thus, simple cleavage of these new bonds would not restore the bases to normal. Hence, before the photoinduced cleavage of these bonds, the enzyme converts the "open form" of [6-4] photoproduct to an oxetane intermediate in which the NH or OH that was transferred to the C5 position of the 5' base is shared by the C5 of the 5' base and the C4 of the 3' base (22). When the oxetane intermediate receives an electron from the active-site flavin, the resulting radical intermediate collapses to the two original bases, with back electron transfer to restore the flavin radical to its catalytically active form (Figure 3).

**Figure 3.** [6-4] Photolyase photocycle. Upon binding to the enzyme, the open form of the [6-4] photoproduct is converted to an oxetane intermediate. The photocycle proceeds as follows: (1) The photoantenna absorbs a photon and (2) transfers energy to the catalytic cofactor; (3) the cofactor transfers an electron to the oxetane intermediate; (4) the photoproduct radical splits into two normal bases, concomitant with electron transfer back to the



The blue-light photoreceptors that exhibit 40% to 60% sequence identity to photolyases (10) have absolutely no repair activity (11). They are photoreceptors for regulating plant growth or for photoentrainment of the circadian clock in animals (13). However, the reaction mechanism of these molecules is not known at present.

#### Bibliography

1. A. Kelner (1949) *Proc. Natl. Acad. Sci. USA* **35**, 73–79.
2. R. Dulbecco (1949) *Nature* **163**, 949–950.
3. C. S. Rupert, S. Goodgal, and R. M. Herriott (1958) *J. Gen. Physiol.* **41**, 457–471.
4. G. B. Sancar, F. W. Smith, M. C. Lorence, C. S. Rupert, and A. Sancar (1984) *J. Biol. Chem.* **259**, 6033–6038.
5. T. Todo, H. Takemori, H. Ryo, M. Ihara, T. Matsunaga, O. Nikaido, K. Sato, and T. Nomura (1993) *Nature* **361**, 371–374.
6. G. Payne and A. Sancar (1990) *Biochemistry* **29**, 7715–7727.
7. M. S. Jorns, G. B. Sancar, and A. Sancar (1984) *Biochemistry* **23**, 2673–2679.
8. J. L. Johnson, S. Hamm-Alvarez, G. Payne, G. B. Sancar, K. V. Rajagopalan, and A. Sancar

- (1988) Proc. Natl. Acad. Sci. USA **85**, 2046–2050.
9. A. P. M. Eker, J. K. C. Hessels, and J. O. van de Velde (1988) Biochemistry **27**, 1758–1765.
10. M. Ahmad and A. R. Cashmore (1993) Nature **366**, 162–166.
11. K. Malhotra, S. T. Kim, A. Batschauer, L. Dawut, and A. Sancar (1995) Biochemistry **34**, 6892–6899.
12. D. S. Hsu, X. Zhao, S. Zhao, A. Kazantsev, R. P. Wang, T. Todo, Y. F. Wei, and A. Sancar (1996) Biochemistry **35**, 13871–13877.
13. Y. Miyamoto and A. Sancar (1998) Proc. Natl. Acad. Sci. USA **95**, 6097–6102.
14. T. Todo, H. Ryo, K. Yamamoto, H. Toh, T. Inui, H. Ayaki, T. Nomura, and M. Ikenaga (1996) Science **272**, 109–112.
15. S. Yasuhira and A. Yasui (1992) J. Biol. Chem. **267**, 25644–25647.
16. Y. F. Li, S. T. Kim, and A. Sancar (1993) Proc. Natl. Acad. Sci. USA **90**, 4389–4393.
17. X. Zhao, J. Liu, D. S. Hsu, S. Zhao, J. S. Taylor, and A. Sancar (1997) J. Biol. Chem. **272**, 32580–32590.
18. H. W. Park, S. T. Kim, A. Sancar, and J. Deisenhofer (1995) Science **268**, 1866–1872.
19. T. Tamada, K. Kitadokoro, Y. Higuchi, K. Inaka, A. Yasui, P. E. de Ruyter, A. P. M. Eker, and K. Miki (1997) Nature Struct. Biol. **4**, 887–891.
20. S. T. Kim, P. F. Heelis, T. Okamura, Y. Hirata, N. Mataga, and A. Sancar (1991) Biochemistry **30**, 11262–11270.
21. I. Husain, G. B. Sancar, S. R. Holbrook, and A. Sancar (1987) J. Biol. Chem. **262**, 13188–13197.
22. S. T. Kim, K. Malhotra, C. A. Smith, J. S. Taylor, and A. Sancar (1994) J. Biol. Chem. **269**, 8535–8540.

### Suggestions for Further Reading

23. G. B. Sancar (1990) DNA photolyases: physical properties, action mechanism, and roles in dark repair. Mutation Res. **236**, 147–160.
24. A. Sancar (1994) Structure and function of DNA photolyase. Biochemistry **33**, 2–9.
25. P. F. Heelis, R. F. Hartman, and S. D. Rose (1995) Photoenzymatic repair of UV- damaged DNA: a chemist's perspective. Chem. Soc. Rev. **24**, 289–297.
26. M. Ahmad and A. R. Cashmore (1996) Seeing blue: the discovery of cryptochrome. Plant Mol. Biol. **30**, 851–861.
27. A. Sancar (1996) "No "End of History" for photolyases". Science **B272**, 48–49.

### Photon Correlation Spectroscopy

Dynamic light scattering is a method for measuring the [diffusion](#) of macromolecules within a solution due to Brownian motion. It differs from "classical" or "static" [light scattering](#), which permits the measurement of the molecular weight and radius of gyration of macromolecules from the light they scatter at varying scattering angles. Dynamic light scattering measures instead the intensity fluctuations of the scattered light on the timescale of  $10^{-3}$  to  $10^{-9}$ s. The frequencies of these fluctuations depend on how rapidly the molecules are moving by diffusion. These measurements require lasers, which provide light of high intensity, collimation, monochromaticity, and, most



importantly, spatial and time coherence; the last attributes mean that the light is emitted from the laser as a continuous wave rather than as short bursts.

The physics underlying dynamic light scattering is complicated (1). In simple terms, the moving macromolecules will *Doppler-broaden* the otherwise monochromatic incident radiation. The scattered intensity will fluctuate because of *beating interference* of scattered waves of different, but similar, wavelengths; the situation is analogous to the fluctuations in intensity of a radio channel caused by interference from another radio channel of a very close wavelength. This broadening of the originally monochromatic light is why the technique is often referred to as *quasi-elastic light scattering* (QLS). The detector sends the intensity signal to a special computer, called an *autocorrelator*, which compares or correlates the intensity at different times; for this reason, the technique has another term, *photon correlation spectroscopy* (PCS). How rapidly the intensity fluctuates over short periods of time, or delay times,  $t$ , is represented by how a parameter known as the *normalized intensity autocorrelation function*,  $g^{(2)}(t)$ , decays as a function of  $t$ . The superscript “(2)” is used to indicate that it is an intensity, as opposed to an electric field “(1),” autocorrelation function. Many data sets of  $g^{(2)}(t)$  as a function of  $t$  are accumulated and averaged. The degree of averaging necessary depends on the incident laser intensity and the size and concentration of the scattering macromolecule (at a given concentration; larger molecules scatter more). For globular proteins, sufficient data can usually be acquired within one to several minutes.

Analysis of how the normalized intensity autocorrelation function  $g^{(2)}(t)$  decays as a function of  $t$  can be used to evaluate the translational diffusion coefficient,  $D_t$ . For dilute solutions of spherical or near-spherical (ie, globular) macromolecules, the variation of  $g^{(2)}(t)$  with  $t$  can be represented by the simple logarithmic equation

$$\ln[g^{(2)}(\tau) - 1] = -2D_t k^2 \tau \quad (1)$$

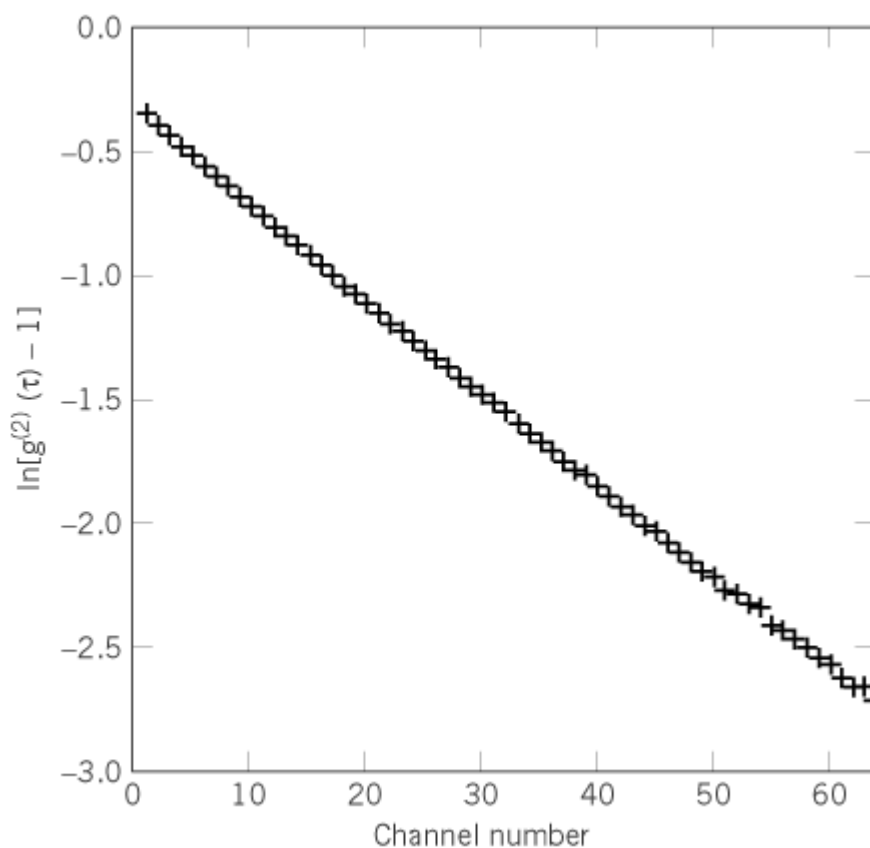
where the parameter  $k$  is known as the *Bragg wave vector*, whose magnitude is defined by

$$k = (4\pi n/\lambda) \sin(\theta/2) \quad (2)$$

where  $n$  is the refractive index of the medium,  $q$  is the scattering angle, and  $l$  is the wavelength of the incident light. Consequently, the value of  $D_t$  can be found from a plot of  $\ln[g^{(2)}(t) - 1]$  versus  $t$ .

Figure 1 illustrates an example for the motility protein [dynein](#).

**Figure 1.** Dynamic light scattering of a solution of dynein. The channel number is a measure of the delay time,  $q$  (see Eq. (1) of the text).



Two important practical conditions need to be noted: ( 1) the value of  $D_t$  is sensitive to the temperature, which needs to be accurately controlled or, at the very least, monitored during the measurement; and ( 2) the scattering signal is very sensitive to the presence of trace amounts of dust or aggregation products; solutions and the scattering cell or cuvette need to be scrupulously clean and free of particulates. The samples must be filtered and centrifuged and the vessels washed. Special filling devices have been constructed to minimize this dust problem (2).

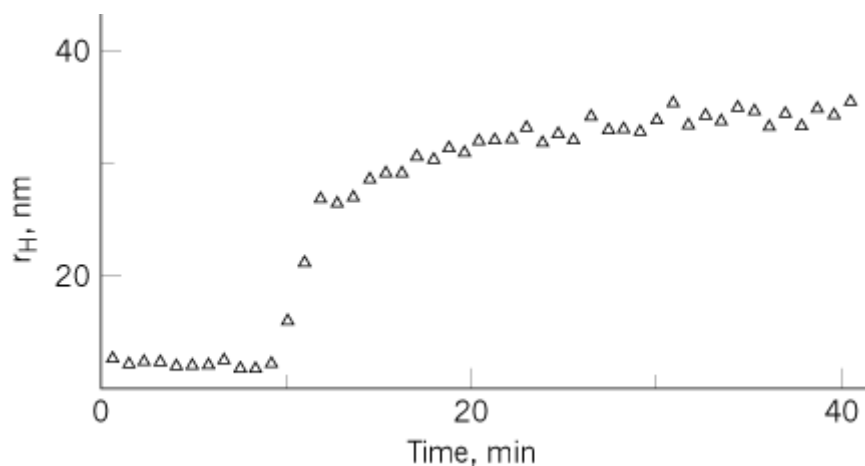
The value of  $D_t$  can be converted to the standard conditions of water at 20.0°C and extrapolated to zero sample concentration (see [Diffusion](#)). An additional extrapolation is necessary if the macromolecule is not globular. There is an extra term on the right-hand side of equation 1 that arises from rotational diffusion at finite angles  $q$ . This term approaches zero as  $q$  does. Therefore, it is necessary to make the measurements of  $D_t$  at a number of angles  $q$  and to extrapolate to zero angle (or Bragg vector  $k$ ). The extrapolations to zero sample concentration and to zero angle  $q$  (or  $k$ ) can be performed simultaneously on a biaxial extrapolation plot known as a *dynamic Zimm plot* (3). If an angular extrapolation is not necessary, a scattering angle of 90° is usually chosen, and fixed-angle instruments are usually set at this angle. At lower angles, the problems due to contamination with large particles (see (2) above) are accentuated.

If a sample is heterogeneous, it is possible, at least in principle, to obtain a distribution of diffusion coefficients after various assumptions and mathematical manipulation of the autocorrelation data; these methods have been reviewed (4), and several commercially available computer programs are available. A more simple way of representing heterogeneity is with the *polydispersity factor* (PF), which is obtained by comparing linear with quadratic or more complex fits of the normalized autocorrelation function decay data (5).

Measurements can be made very rapidly with dynamic light scattering, and it is possible to follow

biomolecular assembly–disassembly processes that occur on a timescale of minutes. An example of the swelling of southern bean mosaic virus in response to the removal of bound  $\text{Ca}^{2+}$  ions is given in Fig. 2 (6).

**Figure 2.** Following the dynamics of a process using dynamic light scattering: effect of removal of calcium ions (by adding EGTA at 10 min) on the hydrodynamic radius ( $r_H$ ) of southern bean mosaic virus (adapted from Ref. 6).



#### Bibliography

1. W. Brown, ed. (1993) *Dynamic Light Scattering. Theory and Applications*, Oxford Univ. Press, U.K.
2. A. H. Sanders and D. S. Cannell (1980) in *Light Scattering in Liquids and Macromolecular Solutions*, V. Degiorgio, M. Corti, and M. Giglio, eds., Plenum Press, New York, pp. 173–182.
3. W. Burchard (1992) in *Laser Light Scattering in Biochemistry*, S. E. Harding, D. B. Sattelle, and V. A. Bloomfield, eds., Royal Society of Chemistry, Cambridge, U.K., pp. 3–22.
4. R. M. Johnsen and W. Brown (1992) in *Laser Light Scattering in Biochemistry*, S. E. Harding, D. B. Sattelle, and V. A. Bloomfield, eds., Royal Society of Chemistry, Cambridge, U.K., pp. 77–91.
5. P. N. Pusey (1974) in *Photon Correlation and Light Beating Spectroscopy*, (H. Z. Cummings and E. R. Pike, eds.) Plenum Press, New York, pp. 387–428.
6. M. Brisco, C. Haniff, R. Hull, T. M. A. Wilson, and D. B. Sattelle (1986) *Virology* **148**, 218–220.

#### Suggestions for Further Reading

7. S. E. Harding, D. B. Sattelle, and V. A. Bloomfield, eds. (1992) *Laser Light Scattering in Biochemistry*, Royal Society of Chemistry, Cambridge, U.K.
8. K. S. Schmitz (1990) *An Introduction to Dynamic Light Scattering by Macromolecules*, Academic Press, New York (a comprehensive introductory treatise, although demanding some familiarity with mathematics).
9. K. E. Van Holde (1985) *Physical Biochemistry*, 2nd ed., Prentice-Hall Englewood Cliffs, N.J.

## Photosynthesis

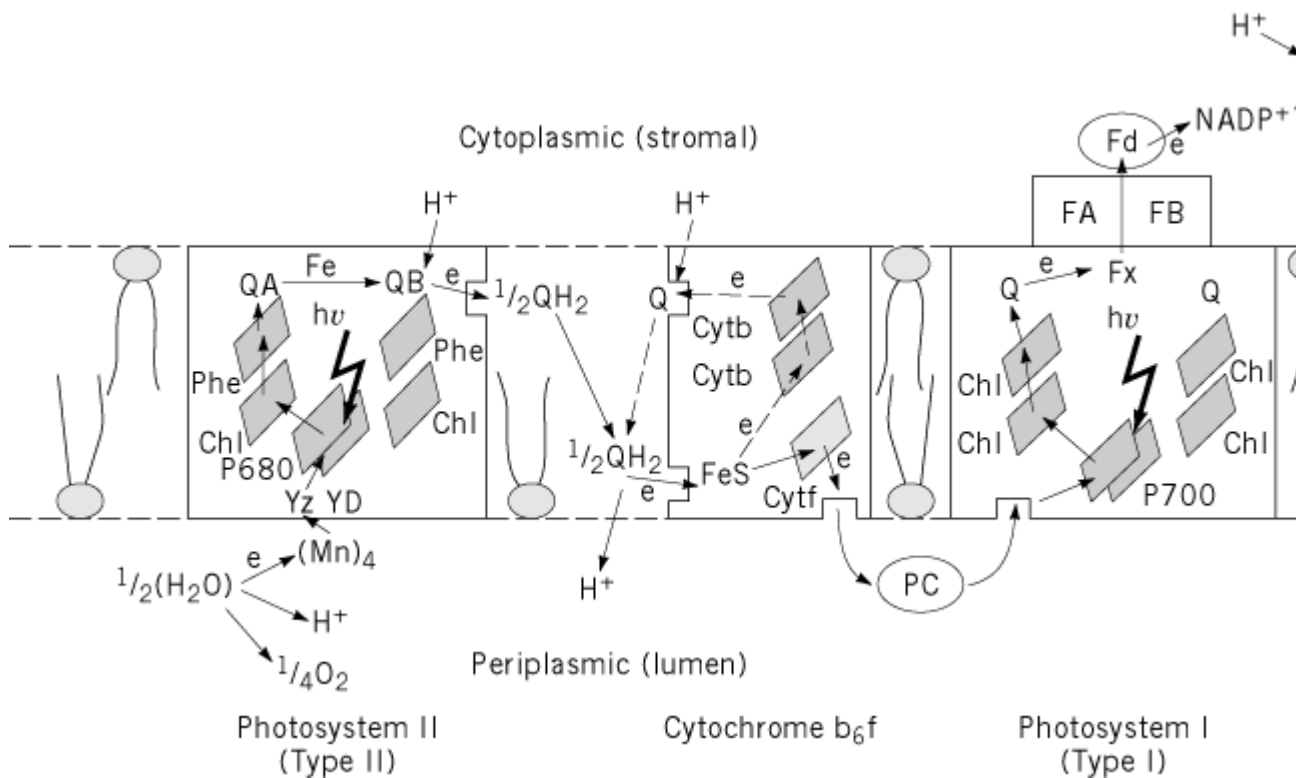
Photosynthesis is the **energy transduction** process whereby solar light energy is converted into a chemical form and used for the synthesis of organic molecules. Photosynthesis takes place in higher plants, algae, and certain bacteria. It involves light-induced **redox** processes in which carbon dioxide becomes reduced. The ability to use [water](#) as an electron source distinguishes the oxygenic photosynthetic organisms—eukaryotic plants and algae, plus prokaryotic cyanobacteria and prochlorophytes, from the nonoxygenic photosynthetic bacteria. Photosynthetic organisms evolved more than 3.5 billion years ago, and they are the origin of current biomass, as well as of fossil fuels.

In all types of photosynthetic organisms, the photosynthesis reaction can be divided into two major phases: (1) the light reactions and (2) the carbon fixation reactions (the Calvin cycle). The light reactions lead to the production of ATP and NADH or NADPH. These molecules are required for the subsequent CO<sub>2</sub> fixation and the formation of trioses, glucose, and more complex carbohydrates.

The primary energy conversion is driven by redox-active pigment-protein complexes known as [photosynthetic reaction centers](#). Together with [light-harvesting complexes](#) and electron transport components, such as [cytochromes](#), [iron-sulfur proteins](#), and quinones, the reaction centers catalyze electron and proton transport across a biological [membrane](#) to create a [proton motive force](#). In analogy with the respiratory process in [mitochondria](#) and bacteria, this proton motive force drives the synthesis of ATP (adenosine triphosphate) by the enzyme [ATP synthase](#).

In oxygenic photosynthesis, NADP<sup>+</sup> is reduced by **electron transfer** from two reaction centers, photosystem II and photosystem I, which are coupled in series (Fig. 1). The electrons are supplied by photosystem II through its unique ability to oxidize water, leading to release of molecular oxygen. The nonoxygenic photosynthetic bacteria contain only one reaction center and typically use organic or inorganic compounds as reductant to reduce NAD(P), or they perform cyclic electron transport to produce ATP. Two groups of reaction centers (type I and type II) are present among the various photosynthetic bacteria and have been shown to be highly analogous to photosystem I and photosystem II of oxygenic photosynthesis, respectively (1).

**Figure 1.** Light-driven, noncyclic electron transport in oxygenic photosynthesis. Two types of photosynthetic reaction centers are linked in series, connected by a cytochrome *b<sub>6</sub>f* complex, to extract electrons and protons from water and to make them available for the reduction of NADP<sup>+</sup> and ultimately CO<sub>2</sub>. The quinone reduction is driven by a type II reaction center called *photosystem II*, the reaction center responsible for creating sufficient reduction potential to reduce NADP<sup>+</sup> via FeS centers is a type I reaction center called *photosystem I*. Photosystem II and photosystem I are evolutionarily related to the type II reaction center of purple bacteria and the type I reaction center of green sulfur bacteria, respectively (see Fig. 2). Both types of reaction center are depicted as having cofactors arranged around a two-fold symmetry axis as revealed by X-ray crystallography (3, 6). *Chl* is chlorophyll-*a*; *P680* (in the case of photosynthetic bacteria these cofactors are either bacteriochlorophyll or bacteriopheophytin). A “special” pair of interacting chlorophylls, which act as the primary electron donor (P<sub>D</sub>), is called *P680* for photosystem II and *P700* for photosystem I. *Fe* is a nonheme iron; *F<sub>x</sub>*, *F<sub>A</sub>*, *F<sub>B</sub>* are low potential electron acceptors, and *F<sub>D</sub>* is ferredoxin. (Mn)<sub>4</sub> is the cluster of manganese atoms involved in the oxidation of water; *Y<sub>Z</sub>* and *Y<sub>D</sub>* are tyrosine residues—the former are on the main route of electron transfer. Arrows show electron transfer paths. The two photosystems are functionally connected by the cytochrome *b<sub>6</sub>f* complex, which oxidizes the reduced quinone (QH<sub>2</sub>, plastoquinol) and copper-containing extrinsic protein, plastocyanin (*PC*). Because quinones are normally 2 electron/2 proton acceptors, the cytochrome complex accepts two electrons, one reducing an iron-sulfur center (*FeS*) and the other reducing a low potential electron acceptor (cytochrome *cyt b*). The reduced FeS center donates an electron to cytochrome (*cyt f*). The other electron participates in independent cyclic electron transport, involving high and low potential b-type cytochromes, which functions as a proton pump. Developed from reference 1.



In eukaryotic organisms, photosynthesis occurs in specialized cell organelles (**chloroplasts**) enclosed by a double outer membrane (envelope) and an internal membrane system (thylakoid membrane). This organelle contains, like the mitochondria, its own DNA. The thylakoid membrane harbors the two photosystems, the light-harvesting antenna, electron transport components, and the ATP synthase (CF<sub>0</sub>-CF<sub>1</sub>). Despite their prokaryotic nature, cyanobacteria have a photosynthetic thylakoid membrane that is differentiated from the plasma membrane. In nonoxygenic photosynthetic bacteria, various forms of invaginations of the cell membrane give rise to so-called *chromatophores*, to which the photosynthetic light reactions are confined. The reduction of CO<sub>2</sub> and the subsequent enzymatic reactions of the Calvin cycle take place in the cytoplasm of prokaryotic cells or in the soluble stromal compartment of chloroplasts. The photosynthetic apparatus is remarkably adaptable to various stress conditions, in particular to short- and long-term changes in light intensity. Under severe light stress, however, photosystem II can suffer from oxidative damage leading to (photo)inhibition of photosynthesis.

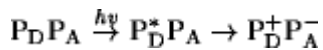
### 1. Photosynthetic Reaction Centers

Photosynthetic bacteria are rudimentary in a general biological perspective, but they have been extremely useful as simple experimental model organisms to elucidate the structure and function of the photosynthetic light reactions and to understand the evolution of photosynthesis (2). These bacteria can be classified into two major groups based on their types of reaction centers: (1) green-sulfur bacteria (eg, *Chlorobium limicola*) and heliobacteria (eg, *Heliobacterium chlorum*, *Heliobaccillus mobilis*) use low potential iron-sulfur clusters as electron acceptors, and their reaction centers are designated FeS type or type I; and (2) purple bacteria (eg, *Rhodobacter sphaeroides*, *Rhodospseudomonas viridis*, *Rhodospirillum rubrum*, *Chromatium vinosum*) and green gliding bacteria (*Chloroflexus aurantiacus*) have type II reaction centers with pheophytin and two quinone molecules as electron acceptors.

**X-ray crystallographic**, functional, and genetic analyses have revealed that photosystem I and

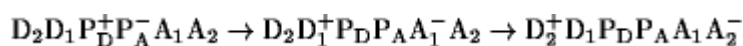
photosystem II of oxygenic photosynthesis are closely related to the bacterial type I and type II reaction centers, respectively (Fig. 1).

The reaction center is the most vital part of the photosynthetic apparatus in which the primary light-induced redox processes take place (1). The type I and II reaction centers have a similar basic architecture, and they function using essentially the same mechanistic principle. Both use a “special pair” of chlorophyll or bacteriochlorophyll molecules as the primary electron donor ( $P_D$ ) and a (bacterio-) chlorophyll (type I) or a (bacterio-)pheophytin (type II) as primary electron acceptor ( $P_A$ ). When excited by light,  $P_D$  donates an electron to  $P_A$  within a few picoseconds, and the radical pair  $P_D^+P_A^-$  is formed:



The primary donor and acceptor are arranged in the reaction center in such a way that the primary charge separation occurs across the lipid bilayer of the photosynthetic membrane.

The stability of this primary charge separation is ensured by rapid transfer of the electron from  $P_A^-$  to the secondary electron acceptors  $A_1$ ,  $A_2$ , and so on and by passing an electron to  $P_D^+$  from secondary electron donors  $D_1$ ,  $D_2$ , and so forth:



In type I reaction centers,  $A_1$  is a quinone, while  $A_2$  is a low potential 4Fe–4 S cluster ( $F_x$ ). The donors are typically *c*-type cytochromes or the copper protein plastocyanin (in oxygenic electron transport only) (see [Electron Transfer Proteins](#)). In the type II reaction center, quinones constitute both the  $A_1$  ( $Q_A$ ) and  $A_2$  ( $Q_B$ ) acceptors. These are normally plastoquinone or, in purple bacteria, ubiquinone. In the bacterial type II centers, cytochromes are often the electron donors, while in oxygenic photosynthesis the electrons are provided by the unique system of water oxidation. The primary donor of photosystem II (P680) is a powerful oxidant in its  $P680^+$  form. In fact,  $P680^+$  is the most oxidizing species found in living organisms ( $E_m = +1.1V$ ), enabling extraction of electrons from water (see [Electron Transfer Proteins](#) and [Redox Enzymes](#)). Oxidation of water requires the accumulation of four oxidizing equivalents, which are necessary to abstract four reducing equivalents from two water molecules to produce one molecule of dioxygen ( $O_2$ ). The water-oxidation system involves a cluster of four manganese atoms that are converted to higher valence in a stepwise manner by four separate photochemical events (the *S*-state cycle) involving P680 and a redox-active [tyrosine](#) residue in the reaction center.

### 1.1. Type II Centers

Determination of the high-resolution three-dimensional structure of the purple bacterial reaction center was a landmark in photosynthesis research (3). The structure revealed an approximately twofold symmetrical relationship of the prosthetic groups (Fig. 1) and the two protein subunits (L and M) that bind them. These two subunits each possess five transmembrane helices and have very similar structures, forming a heterodimeric structure. The two bacteriochlorophylls of the primary electron donor,  $P_D$ , form a special pair at the periplasmic side of the chromatophore membrane.

Following the structure across the membrane, the next cofactors are two additional bacteriochlorophylls, then two bacteriopheophytins, followed by two quinones in the  $Q_A$  and  $Q_B$  sites near the cytoplasmic surface. All the duplicated cofactors (Fig. 1) are related by a twofold axis passing through the special  $P_D$  pair and a nonheme iron located midway between the two quinones. Notably, despite this approximate twofold symmetry, the electron transfer to  $Q_B$  via  $Q_A$  occurs

exclusively along one branch, while the other branch is inoperational (Fig. 1). It is still not clear whether the bacterial chlorophyll located between the primary donor and the bacteriopheophytin is directly involved in the charge separation process as a redox-active species. The symmetry of the reaction center is broken by the H subunit, which has a single membrane-spanning [alpha-helix](#) and a globular **domain** at the cytoplasmic side of the membrane. Furthermore, there is a four-heme cytochrome present on the periplasmic side in most purple bacteria.

Functional and compositional analyses of oxygenic photosynthesis have revealed that photosystem II carries out electron transfer reactions similar to those occurring in the reaction centers of purple bacteria (3). On the acceptor side, photosystem II has pheophytin,  $Q_A$  and  $Q_B$  plastoquinone molecules, as well as the nonheme iron (Fig. 1), and electron transfer occurs with similar rate constants. Furthermore, the two protein subunits of the photosystem II reaction center ( $D_1$  and  $D_2$ ) show homology with the L and M proteins and form a heterodimeric structure similar to the bacterial reaction center. This has been confirmed by the structure of photosystem II obtained at 0.8 nm resolution by [electron microscopy](#). In contrast to the bacterial reaction center, two photosystem II complexes are tightly associated to form a dimeric structure (4).

P680 is probably a dimer of chlorophyll molecules, but the relative orientation of the chromophores and their degree of exciton interaction seem to differ somewhat from the bacterial primary donor. Most significantly, photosystem II differs in its donor properties, allowing the extraction of electrons from water (Fig. 1). Most experimental observations suggest a close association between the active manganese cluster and the reaction center heterodimer.

Despite all the redox components required for the light-induced electron transfer from water to plastoquinone being associated with the  $D_1/D_2$  reaction center heterodimer, photosystem II contains as many as 25 different subunits (1). These are encoded either by chloroplast or nuclear genes designated *psbA-X*. Closely associated with the reaction center is cytochrome b559, which is not a part of the main electron transfer chain but instead protects against photodamage. Most subunits are integral [membrane proteins](#) and span the membrane, but three extrinsic proteins are located at the inner thylakoid surface to stabilize the manganese cluster and to sequester the  $Ca^{2+}$  and  $Cl^-$  ions that are also required for water oxidation. A peculiar feature of the complex is the relatively large number (up to 10) of small proteins with only one membrane-spanning segment. Two chlorophyll-*a*-binding proteins, CP47 and CP43, serve as inner light-harvesting antenna closely associated with the reaction center. Each contains 40 chlorophyll-*a* molecules and possesses 6 transmembrane domains; they are complemented by an outer antenna constructed of various pigment-protein complexes, depending on the taxonomic group (see text below).

## 1.2. Type I Centers

Of the type I reaction centers, photosystem I is the best characterized, partly due to the 0.4 nm-resolution structure of the complex obtained from a thermophilic cyanobacterium (6). Notably, the basic three-dimensional structural arrangement of photosystem I is the same as for the type II center, despite its different content of redox-active groups. Photosystem I has the twofold symmetry with the axis passing through the chlorophyll-*a* dimer of the primary donor (P700) and a 4Fe-4 S center (Fig. 1). Two additional pairs of acceptor chlorophyll molecules occur along this axis. A fourth pair of chlorophyll molecules appears to connect those of the reaction center and of the light-harvesting antenna. The photosystem I reaction center is a heterodimer composed of two large reaction center subunits (*PsaA* and *PsaB*), each possessing 11 membrane-spanning segments. In contrast to the type II center, the two reaction center subunits also harbor the inner light-harvesting antenna of about 100 chlorophyll-*a* molecules. The 22  $\alpha$ -helices of the *PsaA* and *PsaB* proteins are arranged in such a way that 5 helices from each subunit form a central reaction center portion, in analogy with the  $D_1/D_2$  and L/M heterodimers. The additional 12  $\alpha$ -helices are located on either side of this central portion and harbor the antenna chlorophylls. This arrangement should be compared to the inner chlorophyll-*a* antenna of photosystem II composed of two distinct proteins, CP47 and CP43, together comprising

12 membrane spanning helices (4).

Photosystem I is also a multisubunit complex of both nuclear and chloroplast genetic origin, with up to 15 different subunits (*PsaA - PsaM*) (5). *PsaC*, an extrinsic subunit at the outer stromal surface (Fig. 1), provides ligation for two acceptor FeS clusters ( $F_A$  and  $F_B$ ). There are also subunits providing docking sites for [ferredoxin](#) and plastocyanin (see [Azurin](#)) at the acceptor and donor sides of the reaction center, respectively. The eukaryotic photosystem is complemented by an outer antenna of chlorophyll-*a/b* -binding proteins (LHCI). In cyanobacteria, but probably not in plants, the photosystem I complex is trimeric.

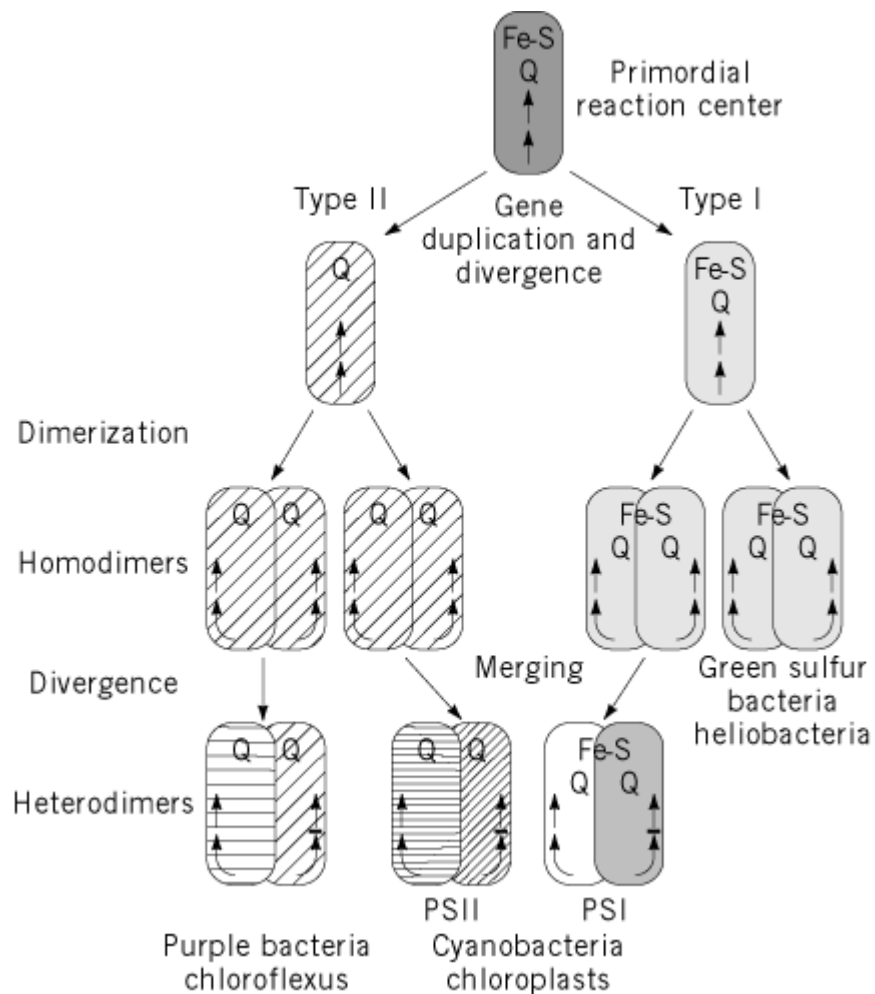
No structural information is yet available for the nonoxygenic type I reaction centers of green sulfur bacteria or heliobacteria. There is compelling evidence that these photosynthetic bacteria have homodimeric reaction centers despite the same overall arrangement of the redox active groups as in photosystem I. It is therefore assumed that these bacterial reaction centers are evolutionarily more primitive and ancestral to the heterodimeric reaction centers (7, 8).

### 1.3. Evolutionary Aspects of Photosynthesis

The occurrence of two different but analogous reaction centers, composed of two related, yet distinct protein subunits, provides insight into the evolution of photosynthesis (1, 7, 8). In the beginning, there was probably a primordial reaction center composed of one single subunit (Fig. 2), which may have evolved to a homodimeric structure with both electron transfer branches operational. Then a [gene duplication](#) took place in early evolution, the heterodimeric reaction centers were generated by [divergent evolution](#), and one of the electron transfer branches was closed. It is assumed that the ancestral reaction center contained both quinone and iron–sulfur acceptors, but that the latter were lost in the type II centers. Later the two types of reaction centers were linked to pave the way for current oxygenic photosynthesis (Fig. 2). It is not clear when the ability to oxidize water was acquired, but no single type II reaction center able to split water has been identified. The chloroplast of eukaryotic photosynthetic organisms is generally assumed to have its evolutionary origin in symbiotic ancestors of cyanobacteria.

**Figure 2.** Evolution of photosynthetic reaction centers. A primordial reaction center containing both quinones and Fe-S centers as electron acceptors, evolved into type I and type II centers, where the latter lost their FeS component. Gene duplications and divergence gave rise to homodimeric, and finally heterodimeric, reaction centers where only one electron transfer route is operational. In cyanobacteria, the two types of reaction centers merged to pave the way for oxygenic photosynthesis (see Fig. 1). The scheme is developed from References 6 and 7.





## 2. Light-Harvesting Protein Complexes

The evolution of type I and type II reaction centers has been matched by the development of a multitude of different light-harvesting pigment–protein systems. This enables photosynthetic organisms to use the entire solar spectrum and provides versatility with respect to the surrounding light environment. In bacteria there are typically 40 chlorophyll-containing molecules per reaction center, while in plants this number could be several hundred (2). The correct organization of the various pigments within the antenna systems is crucial in order to ensure efficient energy transfer. Following absorption of a photon by one of the antenna pigments, the excitation moves from pigment to pigment via so-called resonance energy transfer, until being trapped in the reaction center. These ultrafast processes can now be followed by femtosecond ( $10^{-15}$ s) flash spectroscopy.

Purple bacteria have two types of light-harvesting protein complexes, designated LH1 and LH2 (9). They are composed of two small proteins (a and b) that each span the membrane once and bind bacteriochlorophyll-*a* and carotenoids. LH1 is closely associated with the reaction center (L, M, and H subunits). LH2 is located more peripherally and functions as an outer and variable antenna. The structure of LH2 has been determined to high resolution by X-ray crystallography (9). The complex has a donut shape with 9 a-subunits forming a hollow cylinder with 1.8 nm inner and 3.4 nm outer radii. All the pigments are located between these cylinder walls. The bacteriochlorophylls are ligated to [histidine](#) residues. Eighteen molecules form an overlapping ring located toward the outer membrane surface, while nine are located more toward the inner portion of the membrane. The latter are perpendicular to the other bacteriochlorophylls and with the bacteriochlorophyll ring parallel to the membrane surface. The carotenoids are arranged in close association with the bacterial

chlorophylls. This ring of chlorophylls allows the excitation to become delocalized over a large area, thereby facilitating the energy transfer to LHI, which is likely to have a structure similar to LHII. Purple bacteria can also contain bacteriochlorophyll-*b*, in addition to bacteriochlorophyll-*a*. Typically, purple bacteria can harvest light in a very broad interval of the spectrum, even above 1000 nm.

Green sulfur bacteria have a peculiar antenna system in *chlorosomes*, which appear as bag-like structures on the inside of the cytoplasmic membrane (2). They contain a complex mixture of bacteriochlorophylls *c*, *d*, and *e* assembled into rod elements approximately 5–10 nm in diameter. Notably, chlorosomes contain very few proteins, and the polymeric structures are maintained primarily by pigment–pigment interactions. A water-soluble baseplate bacteriochlorophyll-*a* protein complex sandwiched between the chlorosomes and the cytoplasmic membrane provides the link between the chlorosome antenna and the reaction center chlorophylls. Green sulfur bacteria can absorb light up to 850 nm. Heliobacteria do not contain an outer antenna but manage with the inner bacteriochlorophyll-*a* antenna of its type I reaction center.

Cyanobacteria use water-soluble protein structures known as *phycobilisomes* as outer light-harvesting antenna (10). A phycobilisome contains several hundreds of phycobilin pigments (phycocyanin, phycoerythrin, or allophycocyanin) covalently bound to the phycobiliproteins. The organization of a phycobilisome involves rod-like structures composed of 5 to 6 elements that are connected to a tricylindrical core of phycobiliproteins carrying allophycocyanin. The excitation energy migrates, via the core, from the peripheral rod sections to the CP47/CP43 inner chlorophyll-*a* antenna, and finally to the D<sub>1</sub>/D<sub>2</sub> reaction center. The phycobilisomes are located at the outer thylakoid membrane surface facing the cytoplasm and connected to photosystem II via nonpigmented linker proteins.

The plant photosystems have chlorophyll-*a/b*-binding (CAB) proteins as their outer light-harvesting antenna (10). Altogether there are ten different but related nucleus-encoded CAB proteins: 6 associated with photosystem II and 4 with photosystem I. These are all integral membrane proteins, each with a molecular weight of about 25 kDa and with three transmembrane-spanning regions. The CAB proteins noncovalently bind approximately half of the chlorophyll-*a* and all of the chlorophyll-*b*, as well as a majority of the carotenoids. Most attention has been paid to the major CAB protein complex of photosystem II, designated LHCII. It is the most abundant thylakoid protein complex and, apart from a crucial function in light harvesting, is essential for maintaining the thylakoid membrane organization and to adjust the antenna of photosystem II in response to short- and long-term fluctuations in the light environment. LHCII normally exists as a trimer. Each monomer has been estimated to bind up to 14 chlorophyll molecules, one chlorophyll molecule for every 15 amino acids, indicating a very tight packing. The structure of LHCII has been determined to 0.34 nm resolution by [electron crystallography](#) (11).

Green algae have basically the same outer light-harvesting antenna as plants. In brown algae, the CAB proteins are homologous to the plant CAB proteins. Interestingly, the CAB proteins are not unique to eukaryotic organisms, since the prokaryotic *prochlorophytes* have a similar kind of light-harvesting proteins. On the other hand, eukaryotic red algae have no CAB proteins, but phycobilisomes similar to those of cyanobacteria.

The carotenoids of the various light-harvesting complexes and reaction centers not only play a role in light-harvesting but are also essential as protectants against oxidative damages.

### 3. Electron Transport

#### 3.1. Oxygenic Photosynthesis

Photosystem I and photosystem II are linked in series via plastoquinone, the cytochrome *b<sub>6</sub>f* complex, and plastocyanin (Fig. 1). When the secondary plastoquinone acceptor of photosystem II

( $Q_B$ ) has received two electrons, it is protonated and leaves its binding site on the  $D_1$  protein. This  $Q_B$  site is the major target for several commercial herbicides, such as atrazine and diuron. The reduced plastoquinone interacts with the  $Q_0$  site of cytochrome  $b_6f$ . This cytochrome complex is analogous to the mitochondrial cytochrome  $bc_1$ , whose structure has been solved to high resolution. The photosynthetic cytochrome complex contains four redox centers: one heme  $c$  ( $f$ ), two hemes  $b$ , and a high potential Rieske  $2Fe-2S$  cluster. The redox components are ligated to three protein subunits: cytochrome  $f$  and cytochrome  $b_6$ , which are both plastid-encoded, and the nucleus-encoded Rieske FeS protein. Altogether, the complex contains at least five subunits. It performs quinol oxidation via cytochrome  $b_6f$ , involving a so-called Q cycle, which translocates twice the number of protons that could be translocated by the plastoquinone/plastoquinol oxidoreduction alone. The electrons are eventually linked via cytochrome  $f$ , the FeS center and plastocyanin to photosystem I, reducing  $P700^+$ . The final step in noncyclic electron transfer is the reduction of  $NADP^+$ , which is mediated by ferredoxin and the ferredoxin-  $NADP^+$  reductase, a flavoprotein.

Noncyclic electron transport gives rise to proton transport into the inner (lumen) thylakoid compartment, via cytochrome  $b_6f$  and the plastoquinone/plastoquinol oxidoreduction (Fig. 1). In addition, protons are translocated into the lumen as a consequence of the oxidation of water. Cyclic electron transfer leading to transmembrane proton transport can occur around photosystem I and cytochrome  $b_6f$ , although the precise mechanism is still under debate. The proton motive force created by non-cyclic or cyclic electron transport substantiate the [chemiosmotic coupling](#) theory of ATP formation.

### 3.2. Nonoxygenic Photosynthesis

The simplest form of electron transport occurs in purple bacteria (2). The electrons from the quinol formed at the reaction center are transported back to the primary donor via a cytochrome  $bc_1$  complex and at least one other cytochrome. This light-driven cyclic electron transfer is coupled to ATP production via a proton motive force. Purple bacteria do not possess noncyclic electron transport in the classical sense. A reductant can donate electrons into the quinone pool or to a cytochrome, but these reductants are not sufficiently strong to reduce  $NAD^+$  directly. Instead, this is performed via “reversed” electron transport involving a NADH dehydrogenase enzyme. Some of the energy from the light-mediated cyclic electron transfer is used to drive these reactions.

Photosynthetic bacteria containing type I reaction centers are able to perform non-cyclic electron transport using other compounds as reductants, for example, simple sulfur components (2). The iron-sulfur centers on the acceptor side have redox potentials that are sufficiently low to allow reduction of  $NAD^+$ . In addition, sulfur bacteria and heliobacteria can perform cyclic electron transport coupled to ATP formation.

## 4. Carbon Dioxide Fixation and Photorespiration

The conversion of carbon dioxide to carbohydrates occurs in a cyclic process, called the *Calvin cycle* or the *photosynthetic carbon reduction cycle*. These reactions take place in the cytoplasm of prokaryotic cells or in the stroma of chloroplasts, mediated by at least a dozen enzymes. In the first step, ribulose-1,5-bisphosphate (RuBP) reacts with a molecule of carbon dioxide, thereby forming two identical C3 molecules of 3-phosphoglycerate, 3-PGA. This reaction is catalyzed by **ribulosebisphosphate carboxylase** (Rubisco), which is the most abundant protein on earth, and in eukaryotes consists of eight copies each of a small and a large subunit. In subsequent reactions, 3-PGA is converted into hexose and pentose phosphates, which requires 3 mol of ATP and 2 mol of NADPH per mole of  $CO_2$  fixed. Hexose phosphates are converted to starch or are exported from the chloroplast via the phosphate translocator as triose phosphates. RuBP is regenerated from 3-PGA in reactions requiring further input of NADPH and ATP. The activities of some of the Calvin cycle

enzymes are regulated and coupled to the redox state of the electron transport chain via [thioredoxin \(12\)](#).

A proportion of RuBP is oxygenated by Rubisco, which also can function as an oxygenase, in a potentially wasteful reaction (in terms of quantum efficiency of CO<sub>2</sub> fixation), to form phosphoglycolate and 3-phosphoglycerate. This reaction is called *photorespiration*. The major part of the carbon flow to glycolate is returned, via the photosynthetic carbon oxidation cycle, back into the Calvin cycle. Production of 1 mol of CO<sub>2</sub> results in the net consumption of 7 mol of ATP and 4 mol of NADPH. Although photorespiration lowers the quantum efficiency of CO<sub>2</sub> fixation, it is an important side reaction because it can dissipate excess energy under unfavorable conditions, thereby preventing oxidative damage to the photosynthetic apparatus. In addition, photorespiration provides a link between the carbon and nitrogen cycles and can provide additional ATP to the mitochondria and the cytosol. Under conditions of a short supply of carbon dioxide, such as under conditions of water stress when the leaves close their stomatas to reduce water loss or at elevated temperatures, photorespiration is enhanced.

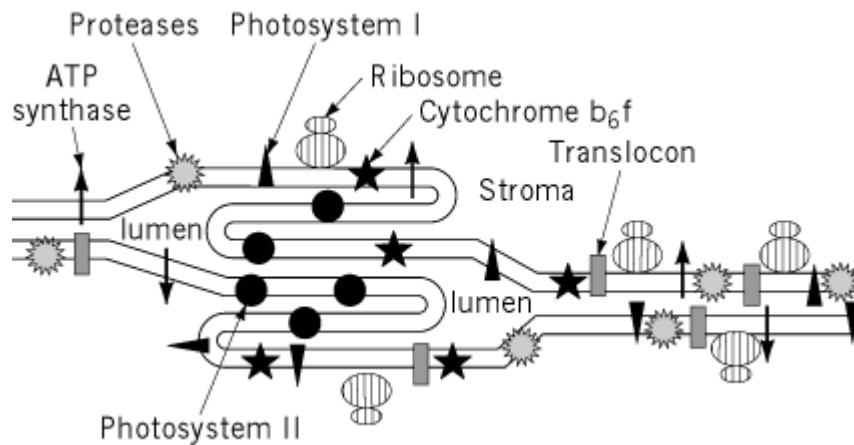
In the so-called C<sub>4</sub> plants, photorespiration is strongly reduced ([12](#)), and the primary carboxylation reaction is of phosphoenolpyruvate (PEP), involving an enzyme that has a much higher affinity for CO<sub>2</sub> than Rubisco. In this carboxylation reaction, 4-carbon carboxylic acids are formed (hence the name C<sub>4</sub> plants) in specialized mesophyll cells and are then transported to the bundle sheath cells. There they are decarboxylated, and the concentrated CO<sub>2</sub> is assimilated into the Calvin cycle via Rubisco. The photorespiration is kept low because of both the elevated CO<sub>2</sub> levels and the very low proportion of oxygen-evolving photosystem II centers in the bundle sheath cells.

## 5. Thylakoid Membrane Organization

### 5.1. Thylakoid Membrane Heterogeneity

The thylakoid membrane is a lipid bilayer that is composed mainly of galactolipids with an unusually high level of unsaturation and encloses an aqueous compartment, the lumen. The thylakoid membrane shows a strong lateral heterogeneity ([5](#)), in terms of both protein and lipid compositions (Fig. [3](#)). A substantial amount of the membrane is folded into stacks (“grana stacks”), where the outer membrane surfaces are in close contact. The extent of stacking varies greatly with plant species and light conditions. The stacked regions contain mainly photosystem II, whereas the unstacked, stroma-exposed membranes contain photosystem I and the ATP synthase (Fig. [3](#)). The intermediate electron carrier cytochrome *b<sub>6</sub>f* is distributed over both membrane regions. Destacking of the membrane and randomization of the protein complexes occurs when divalent cations (mainly Mg<sup>2+</sup>) are depleted. In cyanobacteria, there is no thylakoid stacking or lateral segregation between the photosystems.

**Figure 3.** Schematic model of thylakoid membrane organization and its lateral heterogeneity. The thylakoid membrane is a lipid bilayer that encloses an aqueous compartment, the lumen. The proteins and protein complexes within the thylakoid membrane are unequally distributed over the stacked and stroma-exposed membrane regions.



## 5.2. Biogenesis and Maintenance of the Thylakoid Membrane

Chloroplasts develop from structurally very simple proplastids, which consist of an envelope, a nucleoid with DNA fibrils, a small number of ribosomes, and some flattened thylakoid membrane sacs that appear to be continuous with the inner membrane of the envelope. Development of the proplastid requires light, several chromophores, and multiple nuclear genes, all acting together through several [signal transduction](#) chains in a tightly regulated process. When leaves develop in total darkness, so called etioplasts are formed, which contain no chlorophyll and no thylakoid membrane system. These etioplasts consist of lipid-rich cube-like structures called prolamellar bodies, containing mainly the light-regulated enzyme protochlorophyllide reductase and its pigment substrate protochlorophyllide. Within a day of illumination of the etioplast, chlorophyll is synthesized and an active thylakoid membrane system is formed.

Biogenesis of the thylakoid membrane system involves the synthesis of at least 70 to 85 proteins of the four main photosynthetic protein complexes (5). Many of these proteins are encoded in the nucleus; after synthesis in the cytosol, they are imported into the chloroplast and targeted and inserted into the thylakoid membrane. To increase targeting efficiency to the chloroplast and to prevent missorting to other organelles, these nucleus-encoded proteins are synthesized with a **signal sequence** that is recognized by a receptor on the outer envelope surface (13). After nucleotide triphosphate-driven import into the chloroplast, this signal sequence is cleaved off by a stromal peptidase. Targeting to the thylakoid is facilitated by a second signal sequence, or via information within the mature protein. Different pathways for targeting and insertion into the membrane have been found, and a number of structurally important components, such as the **Sec** machinery and a [signal recognition particle](#) (SRP) have been identified (14). It should be noted that most chloroplast-encoded proteins do not contain a signal sequence. The chloroplast-encoded thylakoid proteins are inserted into the membrane either co-translationally or post-translationally, possibly by mechanisms similar to those in *Escherichia coli*. Most proteins that form the central functional core of each photosynthetic complex are chloroplast-encoded, whereas the more peripheral proteins are nuclear-encoded. To achieve the proper stoichiometry of the different proteins, regulation of synthesis takes place at transcriptional, translational, and post-translational levels, via various signal transduction chains and through feedback mechanisms.

The two photosystems and the cytochrome *b<sub>6</sub>f* complex contain in total at least 15 different cofactors and many different pigments. The biosynthesis of these cofactors and pigments is tightly linked to the synthesis and accumulation of the proteins. Very little is understood, however, about the mechanisms of this coordination.

It has become clear that, in addition to the proteins of the photosynthetic complexes, there must be a fairly large number of proteins involved in the biogenesis and maintenance of the complexes.

Several functions must be performed during these processes: biosynthesis and ligation of cofactors, membrane insertion, folding, and degradation of the proteins. Notably, these auxiliary proteins are mostly located in the stroma-exposed thylakoid regions (Fig. 3).

Recently, members of several classes of these enzymes have been discovered, and most of these auxiliary proteins are encoded in the nuclear genome. Some of these genes have an eukaryotic origin, whereas many originate from prokaryotic chloroplast predecessors. Sequencing of various bacterial genomes, such as those of *Synechocystis* 6803 and *E. coli* and the ongoing sequencing of the higher plant genome of *Arabidopsis thaliana* and other plant species, will prove helpful in identification of the components involved in biogenesis of the photosynthetic apparatus.

### 5.3. Short- and Long-Term Adaptation to the Environment

The photosynthetic apparatus is constantly regulating its activity, in response to both the quantity and the quality of light, as well as the varying demands for metabolic end products of the light and dark reactions (NADPH, ATP, glucose, etc). The two photosystems operate in series but do not have the same action spectrum. Photosystem II can use light at wavelengths below 680 nm, whereas photosystem I can use wavelengths up to 700 nm. Maximal absorption of the outer antenna of photosystem I (LHCI) is 5 nm red-shifted (thus to lower energy level) compared to that of photosystem II (LHCII). Thus, changes in color of the light will affect the relative activities of the two photosystems. For optimal quantum use, it is important to balance the activities of the two photosystems, and several short- and long-term adaptive and regulatory processes must occur. The demand for ATP and reducing equivalents by the dark reaction shows several-fold variations on changing environmental conditions, such as cold, heat, or water stress, and also during different stages of plant development (eg, vegetative growth, senescence).

On a short-term timescale (milliseconds to minutes), light energy can be redistributed via disconnection of the major CAB protein (LHCII) from photosystem II (5). This process is driven by redox-controlled phosphorylation of LHCII. The disconnected LHCII can migrate along the thylakoid membrane to photosystem I and thereby redistribute absorbed light energy. This reversible process is called the *state transition*. Some of the light absorbed by CAB proteins can also be dissipated in the form of heat. This heat dissipation is probably occurring through formation of exciton traps through the reversible conversion of xanthophyll pigments. Furthermore, cyclic electron flow from photosystem I back into the plastoquinone pool could help to balance the activities of the two photosystems and to alter the ATP/NAPDH ratio.

If the short-term adaptations are not sufficient, the quantities of the structural components of the photosynthetic apparatus are changed through their regulated biosynthesis and proteolysis. Various signal transduction chains that respond to different environmental conditions exist to make these dynamic changes possible.

### 5.4. Light-Induced Stress, Photoinhibition and Repair

If the adaptation and regulatory processes described above are not sufficient, damaging oxidative reactions can occur. Several scavenging mechanisms, such as superoxide dismutases and exciton transfer to carotenoids, are present in the thylakoid to quench toxic oxygen species. The main target for photoinhibitory damage of the photosynthetic apparatus is photosystem II, especially its D<sub>1</sub> reaction center protein. Consequently, the D<sub>1</sub> protein has a much shorter lifetime than any other photosynthetic protein. Photoinactivation of photosystem II function can occur either on its donor side, when electron donation to P680 is limiting, or on its acceptor side, when the quinones become overreduced. Since photosystem II is not functional without the D<sub>1</sub> protein, plants have evolved a very efficient repair mechanism. Two important features of this repair mechanisms are (1) a controlled disassembly of photosystem II, to ensure that the other, undamaged, protein subunits of photosystem II are not proteolyzed and can be recycled; and (2) tightly coupled proteolysis and synthesis of the D<sub>1</sub> protein through an as-yet unidentified mechanism.

## Bibliography

1. J. Barber and B. Andersson (1994) *Nature* **370**, 31–34.
2. R. E. Blankenship, M. T. Madigan, and C. E. Bauer (1995) *Anoxygenic Photosynthetic Bacteria*, Kluwer Academic Publishers, Dordrecht.
3. J. Deisenhofer, O. Epp, K. Miki, R. Huber, and H. Michel (1985) *J. Mol. Biol.* **180**, 385–398.
4. K.-H. Rhee, E. P. Morris, J. Barber, and W. Kühlbrandt (1998) *Nature*, **396**, 283–286.
5. B. Andersson and J. Barber (1994) in *Advances in Molecular and Cell Biology*, E. E. Bittar, ed., JAI Press, Greenwich, CT, Vol. **10**, pp. 1–53.
6. W.-D. Schubert, O. Klukas, N. Krauss, W. Saenger, P. Fromme, and H. T. Witt (1997) *J. Mol. Biol.* **272**, 741–769.
7. R. E. Blankenship (1992) *Photosynth. Res.* **33**, 91–111.
8. W. F. J. Vermaas (1994) *Photosynth. Res.* **41**, 285–294.
9. G. McDermott, S. M. Prince, A. A. Freer, A. M. Hawthornthwaite-Lawless, M. Z. Papiz, R. J. Cogdell, and N. W. Isaacs (1995) *Nature* **374**, 517–521.
10. D. J. Simpson and J. Knoetzel (1996) in *Advances in Photosynthesis*, D. R. Ort and C. F. Yocum, eds., Kluwer Academic Publishers, Dordrecht, Vol. **4**, pp. 493–506.
11. W. Kühlbrandt, D. N. Wang, and Y. Fujiyoshi (1994) *Nature* **367**, 614–621.
12. G. Edwards and D. Walker (1983) *C<sub>3</sub>, C<sub>4</sub>: Mechanisms, and Cellular and Environmental Regulation of Photosynthesis*, Blackwell Scientific Publications, Oxford.
13. B. D. Bruce and K. Keegstra (1994) in *Advances in Molecular and Cell Biology*, E. E. Bittar, ed., JAI Press, Greenwich, CT, Vol. **10**, pp. 389–430.
14. D. J. Schell (1998) *Annu. Rev. Plant. Physiol. Plant Mol. Biol.* **49**, 97–126.

## Suggestions for Further Reading

15. B. Andersson and J. Barber (1992) Too much of a good thing; light can be bad for photosynthesis, *Trends Biochem. Sci.* **17**, 61–66.
16. W. Nitschke and A. W. Rutherford (1991) Photosynthetic reaction centres: variations on a common structural theme? *Trends Biochem. Sci.* **16**, 241–245.
17. R. Bassi, D. Sandonà, and R. Croce (1997) Novel aspects of chlorophyll a/b-binding proteins, *Physiol. Plant.* **100**, 769–779.
18. Y. Cohen, S. Yalovsky, and R. Nechushtai (1995) Integration and assembly of photosynthetic protein complexes in chloroplast thylakoid membranes, *Biochim. Biophys. Acta* **1241**, 1–30.
19. K. J. van Wijk, B. Andersson, and E.-M. Aro (1996) Kinetic resolution of the incorporation of the D<sub>1</sub> protein into PSII and localisation of assembly intermediates in the thylakoid membranes of spinach chloroplasts, *J. Biol. Chem.* **271**, 9627–9636.

## Photosynthetic Reaction Center

The photosynthetic reaction center is a **redox**-active complex of pigments and proteins that is the most vital part of the **photosynthetic** apparatus. All reaction centers have similar basic architectures, and their functions follow essentially the same mechanistic principles. They use a “special pair” of chlorophyll molecules as a primary electron donor (P<sub>D</sub>) and a chlorophyll or pheophytin as the primary electron acceptor (P<sub>A</sub>). On excitation by light, P<sub>D</sub> donates an electron to P<sub>A</sub>, and a radical

pair  $P_D^+P_A^-$  is formed. This primary charge separation is stabilized by transfer of the electron to secondary acceptors and by passing an electron to  $P_D^+$  from secondary donors. The donors and acceptors are arranged in the reaction center in such a way that the primary charge separation occurs across the lipid bilayer of a photosynthetic [membrane](#). The redox components are ligated to two distinct, but related, membrane-spanning proteins. In oxygenic photosynthesis in plants, algae and cyanobacteria, two reaction centers, photosystem I and photosystem II, are coupled in series. Photosystem I belongs to the type I reaction centers, which contain FeS clusters as electron acceptors and are also found in anoxygenic green sulfur photosynthetic bacteria. Photosystem II belongs to the type II reaction centers, which contain quinone acceptors and are also present in purple bacteria. The two types of reaction centers are thought to be related evolutionarily. The structures of photosynthetic reaction centers are known primarily through X-ray crystallography, in particular from the elucidation of the reaction center of purple bacteria. See [Photosynthesis](#) for further details.

## Phylogenetic Tree

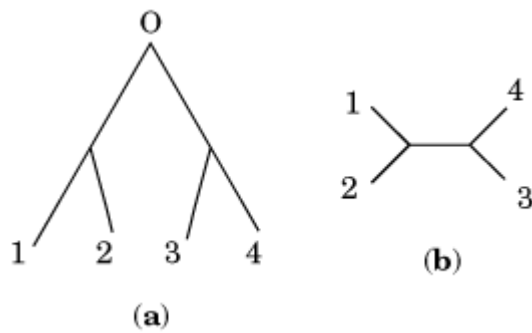
A phylogenetic tree is a tree-like presentation of the **phylogenetic** relationships among species or **genes**. In particular, a phylogenetic tree can be constructed by using evolutionary distances between amino acid or **nucleotide sequences** (see [Sequence Analysis](#)). The tree obtained by this construction is called a “molecular phylogenetic tree” or simply a “phylogenetic tree.”

The methods for constructing a phylogenetic tree can be separated into two major categories, depending on the target traits. These are character-based methods and distance-based methods ([1](#), [2](#)). For example, the maximum-likelihood (ML) method ([3](#)) for constructing a phylogenetic tree is a character-based method, whereas the neighbor-joining (NJ) method ([4](#)) is a distance-based method. In most instances, the following procedures are used to construct a phylogenetic tree: First, the nucleotide (or amino acid) sequences to be compared are aligned simultaneously with each other. When the number of sequences to be compared is more than two, the alignment is called a “multiple alignment.” A character-based method, such as the ML method, can then be used immediately. For the distance-based methods, on the other hand, the number of nucleotide or amino acid substitutions is estimated for all possible combinations of sequence pairs. Using the substitution numbers as [evolutionary distances](#), one is able to use distance-based methods, such as the NJ method.

Phylogenetic trees are either rooted and unrooted trees (see Fig. [1](#)). The rooted tree has an ancestral node, but the unrooted tree is just a network without any ancestral node. The ML methods give a rooted tree, whereas the NJ methods produce a network or unrooted tree. It is, however, possible to identify the ancestral node in a phylogenetic tree that has been constructed by the NJ method. If a particular gene sequence is known to be apparently the most remotely related one to other sequences, it is included as an outgroup into a set of sequences to be compared when a phylogenetic tree is constructed. Then, the branching point between the outgroup and all other sequences can be identified as an ancestral point. Note that if the outgroup is related very remotely to other genes, the substitution numbers may not be estimated correctly because of the so-called saturation effect.

**Figure 1.** Schematic representation of (a) rooted and (b) unrooted phylogenetic trees.





When a phylogenetic tree is constructed by using gene sequences, the tree obtained (the “gene tree”) may be different from the tree of species (the “species tree”) that has been constructed utilizing information other than sequences, such as the fossil record. In particular, if [paralogous genes](#) are included in the analysis, the phylogenetic tree will not be the same as the species tree. This is because paralogous genes, by definition, do not follow the speciation process. Thus, special attention to paralogous genes is required when the species tree is inferred from the gene tree. A species tree can be constructed only by comparing [orthologous genes](#).

#### Bibliography

1. M. Nei, (1987) *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.
2. R. F. Doolittle, (1996) *Methods Enzymology*. **183**, 659–669.
3. J. Felsenstein, (1981) *J. Mol. Evol.* **17**, 368–376.
4. N. Saitou, and M. Nei, (1987) *Mol. Biol. Evol.* **4**, 406–425.

#### Phylogeny

Phylogeny is the evolutionary relationship among organisms or **genes** that takes into account the time and the divergence process from a common ancestor. When evolutionary relationships among organisms are studied by the techniques of molecular biology, the subject and the phylogeny obtained are called “molecular phylogeny” (1). Molecular phylogeny can be studied by immunological methods, protein [electrophoresis](#), DNA-DNA **hybridization**, and amino acid and **nucleotide sequences**. For example, the immunological method utilizes the fact that cross-reactions in [antibody–antigen interactions](#) are stronger for closely related organisms than for distantly related ones. With this method, however, quantitative analysis has been difficult. Now the numbers of nucleotide or amino acid substitutions are used, as evolutionary distances, primarily to reconstruct the phylogeny of genes or proteins (see [Phylogenetic Tree](#)).

Previously in the study of phylogeny, controversy existed between cladistics and phenetics. Cladistics is the study of the pathway of evolution, in which the ancestor–descendant relationship is always discussed with the topology of a rooted phylogenetic tree, called a “cladogram” (see [Phylogenetic Tree](#)). On the other hand, phenetics is the study of the relationship among a group of organisms on the basis of the degree of similarity between their molecular, phenotypic, and anatomic traits. A tree-like network presenting phenotypic relationships is called “phenogram.” When a [molecular clock](#) exists, it is generally considered that a phenogram and a cladogram at the molecular level become identical to each other.

When the phylogeny of closely related organisms or genes is studied, for example, for individuals or **alleles** within a given population, it may be better to call it genealogy rather than phylogeny. Molecular genealogy is important for the study of genetic relationships among individuals or their alleles. The coalescence time is the shortest time for individuals or alleles to share the common ancestor in a given genealogy, and this can give useful information about the time and degree of genetic diversity in the ancestral population (2).

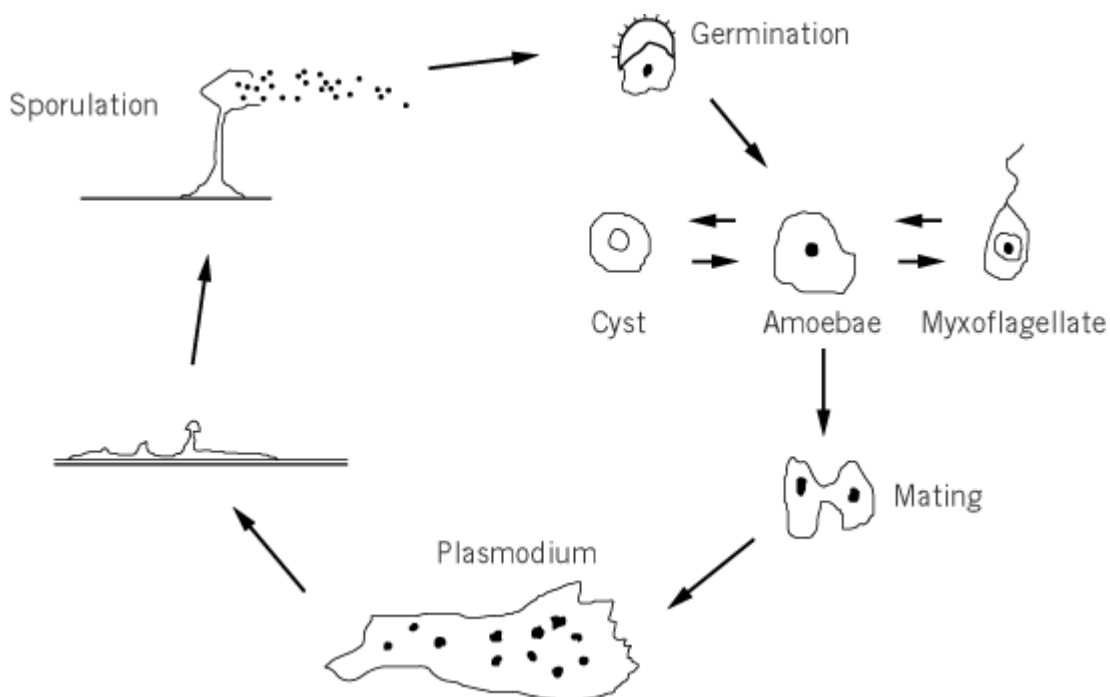
### Bibliography

1. W. H. Li and D. Graur, (1991). *Fundamentals of Molecular Evolution*, Sinauer Associates, Sunderland, MA.
2. M. Nei, (1987) *Molecular Evolutionary Genetics*, Columbia Univ. Press, New York.

### Physarum

The acellular slime mold *Physarum polycephalum* has been studied in the laboratory extensively as a model eukaryotic organism. This soil amoeba has provided insights into many basic cellular functions. It is most famous for its complex developmental and sexual cycle. In the laboratory, it can be maintained in several forms, including a plasmodium that is basically an enormous single cell containing several thousands of synchronously dividing **nuclei**. A simplified schematic of its developmental cycle is shown in Figure 1.

**Figure 1.** Schematic diagram of the *Physarum* developmental cycle.



In conditions that are not limiting in terms of food or moisture, *Physarum* is usually found in the plasmodium form. When the food supply begins to dwindle, the cell can either migrate to find more or begin the process of sporulation. In early stages of sporulation, cytoplasmic projections are formed that are filled with haploid nuclei. As the process continues, structures called *fruiting bodies* are constructed, and the nuclei packed in a protective shell and become spores. A mature fruiting body can rupture and spread the spores over the immediate area, where they remain dormant until conditions become right for germination or they can get carried to a more favorable place by wind or insects.

When conditions are optimal, the spore germinates, releasing a haploid cell. This amoeba can search for food and can transform into a flagellate swimming form if it encounters water. Interestingly, this transition between the amoeba and myxoflagellate forms occurs quickly and does not require additional protein biosynthesis, only rearrangement of existing components. If the amoeba encounters unfavorable conditions, it can change into a dormant and resistant form called a *cyst*. Amoebae can also mate or fuse with others, and the resulting cell grows into the plasmodium.

The *Physarum* [genome](#) is complex. The nuclei in the plasmodium form are a mixture of haploid and **diploid**. There has also been reported extrachromosomal linear DNA that can be maintained even upon [transformation](#) into yeast cells (1). There appear to be substantial differences between strains and under different growth conditions. One strain was reported to have about 70 [chromosomes](#), containing a total of  $1$  to  $5 \times 10^9$  base pairs of DNA (2).

A wide variety of genetic tools have been developed for manipulation of *Physarum*. Mutations can be created, and the resulting phenotypes analyzed in many ways. Extrachromosomal plasmids have been constructed into a family of vectors useful in stable or transient expression of heterologous genes. Multiple selectable markers systems are available, allowing the construction of strains in which complementation can be used to analyze the functions and interactions of proteins. Homologous recombination has been shown to occur, allowing gene knock outs and gene replacements to be made.

Of particular interest in *Physarum* is a complex system of photoreceptors and signaling that controls the onset of sporulation (3). The naturally synchronous cell cycle of *Physarum* plasmodium has been especially interesting, because fluctuations of inositol phosphate levels appear to be involved in its control (4).

Studies using the acellular slime mold *Physarum polycephalum* as a model eukaryotic organism have made contributions to our understanding of basic cell functions. Interesting insights will continue to be made using this soil amoeba in the laboratory.

## Bibliography

1. P. Kunzler (1985) *Nucleic Acids Res.* **13**, 1855–1869.
2. J. Mohberg (1977) *J. Cell Sci.* **24**, 95–108.
3. C. Starostzik and W. Marwan (1994) *J. Bacteriol.* **176**, 5541–5543.
4. M. Belyavskiy and H. W. Sauer (1992) *Eur. J. Cell Biol.* **58**, 371–376.

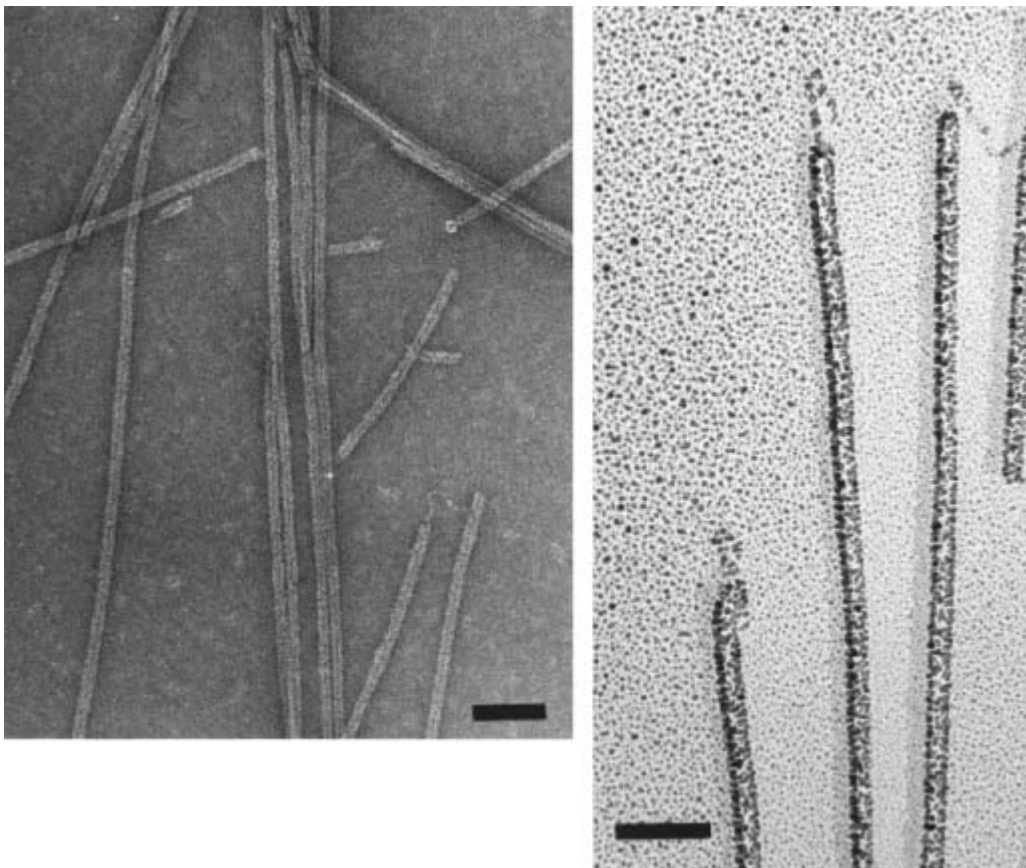
## Suggestions for Further Reading

5. H. C. Aldrich and T. W. Daniel, eds. (1982) *Cell Biology of Physarum and Didymium*, two volumes, Academic Press, New York.
6. W. F. Dove, J. Dee, S. Hatano, F. B. Hugi, and K. E. Wohlfarth-Bottermann, eds. (1986) *The Molecular Biology of Physarum polycephalum*, Plenum Press, New York.
7. W. Sauer (1982) *Developmental Biology of Physarum*, Cambridge University Press, Cambridge, U.K.

## Pili

Adhesion pili are fairly stiff, filamentous structures, typically about 6.8 nm in diameter and about 1  $\mu\text{m}$  in length, that protrude from a **bacterium** and provide specific attachment sites between the bacterium and some target cell (Fig. 1). The pili (or *fimbriae* as they are also known) terminate in a flexible fibrillum about 2.5 nm in diameter. The link between the bacterium and the cell is maintained during mechanically challenging situations *in vivo*, but the pili undergo significant structural rearrangement as a consequence. **Single-particle reconstruction** work has shown that the pilus may unwind without depolymerizing to produce a fibrillum-like structure, thus generating a filament as much as five times its original length whilst maintaining a high degree of structural coherence (1). A wide variety of pili exist *in vivo*, and these include the mannose-binding type I pili of the *Enterobacteriaceae* family and the gonococcal type IV pili. The data suggest that pili from all sources have a conserved structure and a similar morphology.

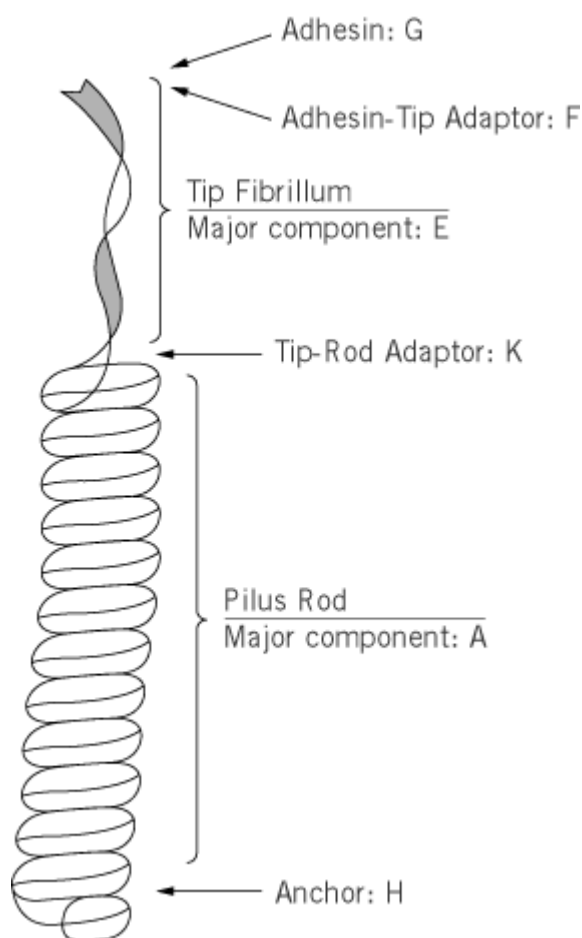
**Figure 1.** Electron micrographs of (left) negatively stained and (right) metal-shadowed P pili from *E. coli* that cause pyelonephritis. The 2.5-nm-diameter fibrillum extending from the end of the pilus is clearly shown at right. Bar = 50 nm. (Courtesy of E. Bullitt and L. Makowski.)



P pili, the surface appendages of *Escherichia coli*, have been particularly well characterized by X-ray [fiber diffraction](#), [scanning transmission electron microscopy](#) (STEM), and three-dimensional reconstruction methods (1, 2), and for that reason they are described here in more detail. The pilus of

diameter 6.8 nm contains approximately 1000 copies of a single protein (PapA) of molecular weight 16.5 kDa (Fig. 2). Pili-terminating fibrillae are composed of proteins known as PapE, PapF, and PapK. The former two proteins position PapG at the tip of the fibrillus: PapG is a digalactose-binding molecule and is responsible for binding to the host (1). PapA, PapE, and PapF are similar in sequence, especially at their C-termini where they bind to PapD, the periplasmic **molecular chaperone**. The crystal structure of PapD (3) reveals a pair of **domains**, each comprising an antiparallel  $\beta$ -barrel, linked by a hinge region. PapH appears to stabilize the attachment of the pilus to the outer membrane of the bacterium. The assembly of P pili also requires PapC, an outer membrane protein.

**Figure 2.** Schematic of the position of proteins in P pili. (From Ref. 4, with permission.)



Fiber X-ray diffraction patterns from orientated pili are consistent with a helical structure having about 33 subunits in 10 turns (2). This helix is characterized by an axial repeat of 24.45 nm, a pitch length of 2.445 nm, a rise per subunit of 0.742 nm, and 3.30 subunits per turn. Measurements of the mass-per-unit length determined using STEM were consistent with each subunit comprising a single PapA molecule. Three-dimensional reconstruction methods yielded similar helical parameters based on a structure with 23 subunits in seven turns (an axial repeat of 17.43 nm, a pitch length of 2.49 nm, a rise per subunit of 0.758 nm, and 3.28 subunits per turn: Fig. 3). These values are also very similar to those reported for type I pili from *E. coli* (pitch length 2.32 nm and 3.125 subunits per turn), in which a single pilin (FimA) predominates. This protein has about 25% sequence identity with PapA and would be expected to have a similar three-dimensional conformation. Using rotary shadowing, the handedness of the P pili helix has been determined. The three-start helix with a pitch length of about 9 nm has been visualized directly and is left-handed. Hence, the fundamental one-start helix is

right-handed. The pili have an off-centered helical cavity (the radial coordinate of the center of the cavity is about 0.5 nm) with approximate dimensions 2.5 by 1.5 nm (1). Channels emanating out from the cavity link it to the outside of the pilus. The shape of the PapA molecule is not yet known, but it has been suggested that its long axis may lie along the same direction as the one-start helix, thus facilitating a possible structural transition to the more extended helical form arising from the response of the pili to mechanical stress. Bullitt and Makowski (1) have proposed that the relatively stiff pilus can be transformed to a flexible fibrillum if the interactions between the PapA molecules on adjacent turns of the one-start helix are broken and if there is a rotation of the domains along the direction of the one-start helix. PapA, as well as some of the other Pap proteins, would thus be expected to be elongated under these conditions. Because the mass-per-unit-length of a fibrillum is about one-fifth of that of the pilus, it is possible that the pilus could be extended fivefold. As a result, many pili with different lengths would be able to bind the cell simultaneously, thus providing stronger adhesion than if the pili were of fixed length. In that case they would have a greater chance of being broken individually when under shear.

**Figure 3.** Three-dimensional reconstruction of P pili from electron microscopy and image analysis. P pili are 6.8 nm in diameter and have a helical symmetry with 23 subunits in seven turns. The pitch length of the helix is 2.49 nm, and the axial rise per residue is 0.758 nm. *E. coli* may utilize plastic deformation of P pili to withstand shear forces in the host cell environment, extending the pilus to a thin fibrillum-like structure to maintain attachment to the target cell. Bar = 2.5 nm. (Courtesy of E. Bullitt and L. Makowski.)



## Bibliography

1. E. Bullitt and L. Makowski (1995) Structural polymorphism of bacterial adhesion pili. *Nature (London)* **373**, 164–167.
2. M. Gong and L. Makowski (1992) Helical structure of P pili from *Escherichia coli*: evidence from X-ray fiber diffraction and scanning transmission electron microscopy. *J. Mol. Biol.* **228**, 735–742.
3. A. Holmgren and C. I. Branden, (1989) Crystal structure of chaperone protein PapD reveals an immunoglobulin fold. *Nature (London)* **342**, 248–251.
4. M. J. Kuehn, F. Jacob-Dubuisson, K. Dodson, L. Slonim, R. Striker, and S. J. Hultgren (1994) Genetic, biochemical, and structural studies of biogenesis of adhesive pili in bacteria. *Methods Enzymol.* **236**, 282–306.

## Suggestion for Further Reading

5. M. J. Kuehn, F. Jacob-Dubuisson, K. Dodson, L. Slonim, R. Striker, and S. J. Hultgren (1994) Genetic, biochemical, and structural studies of biogenesis of adhesive pili in bacteria. *Methods Enzymol.* **236**, 282–306.

## Pinocytosis

Pinocytosis is a form of [endocytosis](#) that is mediated by several different vesicular mechanisms and is responsible for both receptor-mediated and fluid-phase endocytosis. **Clathrin**-mediated endocytosis is the best-characterized pinocytic pathway and, in some cell types at least, is the major pathway for pinocytosis.

### 0.1. Clathrin-Dependent Endocytosis

The uptake of most receptor-bound ligands and extracellular fluid is mediated by endocytic clathrin-coated vesicles, typically 100 to 150 nm in diameter ([1](#)). These form constitutively from plasma membrane clathrin-coated pits during most of the [cell cycle](#), but their formation, in common with other membrane trafficking events, does not occur during mitosis. Clathrin-coated vesicle assembly is initiated by the multifunctional AP-2 adaptor protein complex. The recruitment of AP-2 from the cytosol is poorly understood, but it may involve binding to inositol phospholipids and possibly to a polypeptide receptor [two related adaptor complexes, AP-1 and AP-3, are involved in vesicular transport from intracellular compartments ([2](#))]. Membrane-bound AP-2 provides sites for attachment of cytosolic clathrin triskelions, which assemble either as flat lattices or invaginating lattice structures (see [Clathrin](#)). Clathrin assembly is required to drive membrane invagination at the cell surface, although local concentrations of specific lipids also contribute to membrane deformation.

AP-2 also provides binding sites for a spectrum of other proteins that are required structurally or catalytically for the formation of coated vesicles. These include amphiphysin, **dynamin**, EPS-15, synaptojanin, AP180/CALM, and the endophilins (SH3p4, SH3p8, SH3p13) ([3](#)). The precise roles of these proteins are still being established. Dynamin is currently the best characterized. This 100-kDa [GTPase](#) is recruited to coated pits *via* the **SH3** domain of amphiphysin and is required for pinching off (scission of) invaginated coated vesicles ([1](#)). GTP hydrolysis is required for a conformational change in dynamin that leads to fusion of the membrane at the neck of an invaginated vesicle. Cells that express dominant-negative mutant forms of dynamin do not internalize [transferrin](#) and other ligands taken up by receptor-mediated endocytosis. Moreover, the *Shibire* mutant of *Drosophila*

*melanogaster* has a temperature-sensitive defect in dynamin. At the nonpermissive temperature, *Shibire* mutants are paralyzed because they cannot release clathrin-coated vesicles, which is required to regenerate synaptic vesicles. Dynamin has also been implicated in forming of caveolae and non-clathrin vesicles (4, 5), but these processes are not well understood, and the full extent of dynamin function is still to be clearly established. Rab5 and RabGDI are also required for coated vesicle budding (6), and this may reflect some quality-control mechanism to ensure that endocytic-coated vesicles also include the correct SNARE proteins and allow the subsequent targeting and fusion of the vesicles with early endosomes.

Clathrin-dependent endocytosis is inhibited by treatments that interfere with the functions of the regulatory proteins, such as cytosolic acidification, hypertonic medium, or  $K^+$ -depletion, or by overexpressing modified forms of key protein components (clathrin hub domains, GTP-bound dynamin, amphiphysin-SH3 domains, Eps15 AP-2-binding domains), all of which act as dominant negatives (7). Under these conditions clathrin-dependent, receptor-mediated endocytosis (e.g., transferrin uptake) is blocked. Although the overexpression, for example, of dominant-negative dynamin blocks clathrin-dependent endocytosis, fluid-phase endocytosis is unaffected or only partially affected, suggesting that there is also a clathrin-independent pathway for constitutive endocytosis. Although this clathrin-independent pathway may be active to some extent under normal conditions (its contribution to total endocytosis may vary in different cell types ranging from <10% in baby hamster kidney cells to  $\approx 50\%$  in hepatocytes), its activity may be up-regulated when clathrin-dependent endocytosis is effectively blocked (8). The molecular basis of this clathrin-independent pathway is not known, but there are several other pathways that could mediate fluid-phase uptake (see later).

Following endocytosis, clathrin-coated vesicles must be uncoated to allow fusion with early endosomes. Uncoating, which is ATP-dependent and requires the proteins auxilin and hsc70, releases clathrin, AP-2, and the other vesicle components, allowing them to be reused (9).

## 0.2. Clathrin-Independent Endocytosis

Caveolae, or plasmalemmal vesicles, are flask-shaped (60 to 80 nm in diameter) invaginations seen on the cytoplasmic aspect of the plasma membrane in many cell types. It is believed that they form by the oligomerization of a coat protein, caveolin (VIP-21). The function of caveolae is unclear. In endothelial cells, they mediate transcytosis, but in other cells they may provide an alternative endocytic route, although their net contribution to total endocytosis is likely to be small. Caveolae function in the concentration of a subset of cell-surface proteins, including **GPI-anchored** molecules, such as the folate receptor, and may internalize these components. Caveolae have also been implicated in the entry of **SV40 virus**. The caveolar membrane has a characteristic lipid composition distinct from the bulk of the plasma membrane and includes a relatively high cholesterol content. Thus drugs that preferentially bind cholesterol, such as filipin or nystatin, disrupt caveolae structurally and functionally. This distinctive membrane lipid composition also includes a concentration of the ganglioside GM1. GM1 functions as the receptor for **cholera toxin** and, in keeping with this, the endocytic uptake of cholera toxin is caveolae-mediated. Dynamin (see previous) has also been implicated in the pinching off required to release caveolae into the cytosol, presumably using the same mechanism as in clathrin-mediated endocytosis (4, 5). Caveolae also function in **calcium signalling** and **signal transduction**, and purified caveolae contain **heterotrimeric G proteins**, **Src** family protein kinases, and protein kinase C. Certain cell types, for example lymphocytes, lack VIP-21/caveolin and identifiable caveolae, suggesting that caveolar function is not essential for cell viability. The fate of internalized caveolae is unclear, and it has been suggested that they may fuse with early endosomes.

## Bibliography

1. S. L. Schmid (1997) *Annu. Rev. Biochem.* **66**, 511–548.
2. G. Odorizzi, C. R. Cowles, and S. D. Emr (1998) *Trends Cell Biol.* **8**, 282–288.
3. O. Cremona and P. De Camilli (1997) *Curr. Opin Neurobiol.* **7**, 323–330.



4. J. R. Henley, E. W. Krueger, B. J. Oswald, and M. A. McNiven (1998) *J. Cell Biol.* **141**, 85–99.
5. P. Oh, D. P. McIntosh, and J. E. Schnitzer (1998) *J. Cell Biol.* **141**, 101–114.
6. H. McLauchlan, J. Newell, N. Morrice, A. Osborne, M. West, and E. Smythe (1998) *Curr. Biol.* **8**, 34–45.
7. H. Riezman, P. G. Woodman, G. van Meer, and M. Marsh (1997) *Cell* **91**, 731–738.
8. H. Damke, T. Baba, D. E. Warnock, and S. L. Schmid (1994) *J. Cell Biol.* **127**, 915–934.
9. S. E. Holstein, H. Ungewickell, and E. Ungewickell (1996) *J. Cell Biol.* **135**, 925–937.
10. R. G. Parton (1996) *Curr. Opin. Cell Biol.* **8**, 542–548.
11. E. Stang, J. Kartenbeck, R. G. Parton (1997) *Mol. Cell. Biol.* **8**, 47–57.

## Plant Genetic Engineering

Plant genetic engineering is defined as the creation of a complete **plant** containing a fragment of foreign **DNA** stably inserted into its nuclear **genome**. During the last 15 years, the majority of the world's most significant crop plants have been transformed. Apart from a significant impact on plant research in general, the products of **transgenic** plants are currently entering the marketplace and, indeed in some cases, are achieving market dominance.

Plant **transformation** relies on inserting foreign DNA into a single plant cell, followed by regenerating that cell into a complete plant. Often this process relies on **totipotency**, the unique ability of isolated plant cells, when cultured in the presence of plant **hormones**, to divide to form an undifferentiated callus that, in turn, is induced to form shoots and roots and thus a complete plant. Two general means are available for transferring foreign DNA to plant cells, *Agrobacterium*-mediated transformation and naked DNA uptake.

*Agrobacterium* is a soil **bacterium** that contains a tumor-inducing **Ti plasmid**. During *Agrobacterium*-mediated transformation, a defined fragment of the Ti plasmid, the so-called transferred or T-DNA (see **T-Complex, -DNA, -Region, -Strand**), is transferred from the bacterium to the infected plant cell and integrated stably into the plant nuclear genome. The T-DNA is delimited by 25-bp imperfect border repeats. Transformation **vectors** based on the T-DNA have been developed and can be grouped into two types, cointegrative and binary. **Cointegrative vectors** are based on a Ti plasmid in which the body of the T-DNA within the border sequences has been replaced by a defined sequence of DNA. Foreign DNA contained on an appropriate **plasmid** is introduced into *Agrobacterium* and integrated into the cointegrative vector by **recombination**. Binary vectors are based on plasmids stable in both *Escherichia coli* and *Agrobacterium*. They contain T-DNA border sequences between which foreign DNA can be **cloned** to be transferred to the plant genome. Then, the binary vector is transferred to a host *Agrobacterium* that contains a Ti **helper plasmid** which contains functions necessary for gene transfer to plant cells, but lacks a T-DNA. Because cointegrative vectors are based on the Ti plasmid, they are very stable in *Agrobacterium*, but they can be difficult to work with in practice. On the other hand, binary vectors can be handled with relative ease, although they are unstable in *Agrobacterium* and continued **antibiotic resistance** selection is required to maintain them (1).

Transformation using *Agrobacterium* requires direct contact between the plant cell and the infecting bacteria. In practice, this is achieved by cocultivation of *Agrobacterium* with isolated plant cells or explants, such as leaves, roots, or tubers. In this case, bacteria are incubated with plant material, and transformation occurs over a period of about two days. Then, the plant material is washed to remove

excess bacteria and cultured further to allow plant regeneration (1). Recently, a novel nontissue culture approach of transformation has been described using the model plant [Arabidopsis](#). Here the developing floral buds of intact plants are vacuum infiltrated with the *Agrobacterium*. The plants are washed, allowed to recover, and selfed. From progeny seeds, transgenic plants can be selected by growth on antibiotics to which resistance is encoded by the foreign DNA contained on the transformation vector (2). It has been generally assumed that *Agrobacterium*-mediated transformation is limited to the obvious natural hosts of *Agrobacterium*, the dicotyledonous plants. Recently, however, it has been established that *Agrobacterium* can be used effectively to create transgenic monocotyledonous plants, such as rice (3).

Naked DNA uptake was developed for two reasons: (1) as a means of transforming plants whose cells, or tissues, are not amenable to culture and to high-frequency *Agrobacterium* transformation; (2) because of the possibility of using nonspecialized plant transformation vectors in transformation. Naked DNA transformation uses standard bacterial plasmids as vectors. The only requirement is that the plasmid contains a gene linked to a plant-specific promoter that is functional in transgenic plants. Several methods are available to transfer foreign DNA into plants. The two techniques used most frequently with isolated cells, treatment with **polyethylene glycol** and electric discharge or electroporation (see [Transfection](#)), achieve relatively high frequencies of transformation, which generally result in measurable activity of the transferred gene constructs after approximately one day. This allows transient expression assays, for example, to test the expression of specific promoter/gene constructs. Cells transformed in this manner can also be regenerated to intact plants as described previously. Naked DNA can also be introduced into plant cells by particle bombardment. This involves coating gold or tungsten particles with DNA and using a gas or electric discharge to project the particles into target cells. This can be used to study transient expression in intact tissue, but its greatest practical application has been to insert foreign DNA into explants from which shoot or root organogenesis can occur (4).

In both *Agrobacterium*-mediated transformation and naked DNA uptake, the foreign DNA stably integrates into the nuclear genome of the plant. *Agrobacterium*-mediated transformation can result in inserting a single unrearranged fragment of DNA that represents the T-DNA. In transgenic cells resulting from naked DNA uptake, however, the DNA might be rearranged and can occur as multimers.

Upon insertion into the genome, a gene linked to a plant-specific promoter is expressed in a pattern determined by that promoter. A variety of [reporter genes](#) have been developed to allow direct selection for the growth of transgenic material (for example, resistance to [kanamycin](#), hygromycin, or herbicides) or for screening for expression (for example, histochemical staining with **b-glucuronidase** or **fluorescence** resulting from expression of the jellyfish [green fluorescent protein](#)). An increasing variety of promoters are available to produce expression in most plant tissues (for example, the 35 S-RNA promoter derived from [cauliflower mosaic virus](#)) or expression in a tissue-specific manner (for example, tapetum-specific expression). Recently promoter constructs have been developed so that the expression of specific genes can be induced in defined tissues following the external application of an inducing substance, such as [glucocorticoids](#) (5).

Gene expression can also be inactivated in transgenic plants, by the expression of either [antisense oligonucleotides](#) (6) or sense constructs of the target gene, a phenomenon known as cosuppression. The exact mechanism by which cosuppression occurs is not known nor is the reason for the instability of transgenic expression that often occurs in some plant lines (7).

In addition to transferring foreign genes to the genomes of plants, transformation in plants can be considered a form of insertional [mutagenesis](#). This is the basis of *gene tagging*, where a defined sequence of DNA inserts into the genome and knowledge of this sequence allows recovering the flanking plant sequence from mutant genomes by **hybridization**, **PCR**, or plasmid rescue. Either [transposable elements](#) transferred to a novel plant host or T-DNA can be used in gene tagging. The latter has the advantage of producing single insertions, simplifying genetic analysis, but requires

screening large numbers of mutant individuals following a transformation experiment. The transposable element approach has the advantage of using populations of mutant individuals produced by genetic crossing. Each individual contains multiple transposon insertions, thus reducing the number of mutants that need to be screened for a mutant phenotype. The construction of gene tags can vary. For example, at their borders they can contain a reporter gene lacking a promoter, so that expression of the reporter gene occurs only if it is inserted into the genome near a promoter (8, 9). Another example is the use of a tag that contains transcriptional **enhancers**, so that expression of flanking genes becomes activated following insertion of the tag into the plant genome (10).

Plant genetic engineering has had a profound effect on basic plant biology. Transgenic plants have been used to study the mechanism of gene expression and the developmental and biochemical consequences of the expression or inhibition of the expression of various genes (11).

The creation of transgenic plants has begun to have a profound effect on agriculture because of the steady appearance of genetically engineered crop plants. These include plants tolerant of herbicides or resistant to insect predation and viral infection. Interestingly, plants altered in developmental events or biochemical processes have also gained a place in the agronomic marketplace. Possibly the most widely known example of this is the tomatoes engineered, using antisense technology, so that the polygalacturonidase gene, which is normally involved in the softening of fruit, is not expressed during fruit development. These tomatoes have delayed ripening (12). Another interesting example is the potatoes engineered so that they contain a changed starch content and quality for industrial processes (13). The acreage of transgenic crops planted is steadily increasing worldwide (14). This technology is, however, not without its opponents. Concerns have been raised as to transfer of the transgenes from transgenic plants to local plant relatives, increasing the weediness of certain plant species, and the creation of products that may induce unexpected allergies in consumers.

### Bibliography

1. J. Draper, R. Scott, P. Armitage, and R. Walden (1989) *Plant Genetic Transformation and Gene Expression. A Laboratory Manual*, Blackwell Scientific, Oxford.
2. F. F. White (1993) *Transgenic Plants*, Vol. 1 (D. Kung and R. Wu, eds.), Academic Press, San Diego, pp. 15–48.
3. Y. Hiei, S. Ohte, T. Komari and T. Kumashiro (1994) *Plant J.* **6**, 271–282.
4. N. Bechtold, J. Ellis, and G. Pelletier (1993) *C. R. Acad. Sci. Paris, Life Sciences* **316**, 1194–1199.
5. T. Aoyama, C. H. Dong, Y. Wu, M. Carabelli, G. Sessa, I. Ruberti, G. Morelli, and N-H Chua (1995) *Plant Cell* **7**, 1775–1785.
6. A. R. van der Krol, J. N. Mol, and A. R. Stuitje (1988) *Gene* **72**, 45–50.
7. S. P. Kumpatla, M. B. Chandrasekharan, L. M. Iyer, G. Li, and T. C. Hall (1998) *Trends Plant Sci.* **3**, 97–104.
8. M. A. Haring, C. M. Rommers, H. J. Nijkamp, and J. Hille (1991) *Plant Mol. Biol.* **16**, 449–461.
9. J. F. Topping, W. Wei, M. C. Clark, P. Muskett, and K. Lindsey (1995) *Methods Mol. Biol.* **49**, 63–79.
10. H. Hayashi, I. Czaja, H. Lubepow, J. Schell, and R. Walden (1992) *Science* **258**, 1350–1353.
11. T. ap Rees (1995) *Trends Biotech.* **13**, 375–378.
12. M. G. Kramer and K. Redenbaugh (1994) *Euphytica* **79**, 293–297.
13. R. G. F. Visser et al. (1991) *Mol. Gen. Genet.* **225**, 289–296.
14. USDA Animal and Plant Health Inspection Service. Internet <http://www.aphis.usda.gov/bbep/bp>.

### Suggestions for Further Reading

15. R. Walden (1989) *Genetic Transformation in Plants*, Prentice-Hall, Englewood Cliffs, NJ.

16. M. J. Chrispeels and D. E. Sadava (1994) *Plants, Genes and Agriculture*, Johns and Barlett.
17. S. D. Kung and R. Wu (1992) *Transgenic Plants*, Vols. **1** and **2**, Academic Press, San Diego.
18. C. Tudge (1988) *Food Crops of the Future*, Blackwell Scientific Press.

## Plant Hormones

Development of higher organisms—plants as well as animals—is largely based on a coordinate division, expansion, and [differentiation](#) of cells. This coordination is primarily based on cell-to-cell communication over short and long distances and is orchestrated by messengers of various chemical nature. One class of signaling molecules contains the [hormones](#), characterized by their ability to exert, most often, specific effects over a distance and by their capacity to act at very low concentrations. This entry presents a brief overview of the current state of the art in the field of plant hormones, or phytohormones.

Since the first discovery of phytohormones, approximately 70 years ago, dedicated efforts of plant physiologists and biochemists have produced a wealth of data on their chemical nature, pathways of biosynthesis and metabolism, possible conjugates, transport mechanisms throughout the plant, and physiological functions ([1](#)). Studies have been based traditionally on the correlation of the effects observed upon exogenous application of hormones with changes in their endogenous concentration. A major disadvantage of this approach is that it remains uncertain whether the effect reflects the normal physiological role of the hormone. Another problem arises from the fact that the distribution of the growth regulator over the entire plant body is largely unpredictable. Moreover, problems of uptake and rapid metabolism were not uncommon. The use of inhibitors had the additional drawback of potential lack of specificity. In the past decade, our knowledge of the molecular basis of plant hormone biology has made a giant leap forward, mainly owing to molecular genetic studies of mutants in thale cress ([Arabidopsis thaliana](#) (L.) Heynh.). Molecular techniques have allowed the unraveling of mechanisms by which phytohormones control tissue-specific and developmentally regulated responses. Genetic engineering enables changes in internal concentration or perception of a hormone in a temporally and spatially controlled manner. Moreover, several applications in agriculture and horticulture are being explored ([2-4](#)). Introduction of bacterial hormone biosynthesis genes (mainly from [Agrobacterium tumefaciens](#)) into plants has confirmed much of the classical view of hormone action ([5](#)).

Our current understanding of the biosynthesis, perception, [signal transduction](#), control of **gene expression**, and effects of each class of phytohormones are presented in the separate entries. Six classes of compounds are recognized as true phytohormones: **auxins**, **cytokinins**, **gibberellins**, **abscisic acid**, **ethylene**, and [brassinosteroids](#). Other growth regulators, such as polyamines, oligosaccharides, and salicylic and jasmonic acids, also influence plant growth and development. For a description of their effects, see recent reviews ([6-8](#)).

It is important to stress that plant hormones do not play exclusive roles in development, but that they instead act coordinately at different stages. The most direct indication for this characteristic is their overlapping physiological roles ([1](#)). This cross-talk is most certainly based on molecular interactions between factors regulating signal transduction and feedforward or feedback loops of different hormone pathways. Moreover, environmental cues can exert an additional level of control on plant hormone abundance and activity. Thus, hormone signaling pathways are probably integrated with those participating in phototransduction, control of the biological clock, nutrient sensing, and stress ([9](#)). Given the redundancy of involvement of signaling intermediates, such as **G proteins**,

**phosphorylation** cascades, and [calcium signaling](#), mechanisms should exist to maintain accuracy of individual signaling chains (10). Much of the specificity might derive from a restricted spatial distribution of signaling components together with their assembly in complexes, called *transducons*. The latter would consist of a module of receptors, receptor targets, and downstream signaling factors, presumably coupled by 14–3–3 proteins (11). Such scaffolding proteins have been demonstrated in signaling pathways in bacteria and in yeast (12, 13). Interactions between different transduction cascades could be exerted by components that navigate between them—for example, [second messengers](#) or mobile protein **kinases** (14).

With the current progress in sequencing of the [genome](#) of *Arabidopsis*, plus the technical advances that allow rapid map-based cloning, creation of large tagged-mutant collections (15), and production of **knockout** mutants by homologous [recombination](#) (16), it can be predicted that the molecular controls of hormonal pathways and their target genes will be unraveled at a dazzling pace in the next few years. The puzzle of inter- and intracellular hormonal networking will soon come of age.

### Bibliography

1. P. J. Davies (1995) *Plant Hormones: Physiology, Biochemistry and Molecular Biology*, Kluwer, Dordrecht, The Netherlands.
2. C. Mariani, V. Gossele, M. De Beuckeleer, M. De Block, R. B. Goldberg, W. De Greef, and J. Leemans (1992) *Nature* **357**, 384–387.
3. R. Ayub, M. Guis, M. Ben Amor, L. Gillot, J.-P. Roustan, A. Latché, M. Bouzayen, and J.-C. Pech (1996) *Nature Biotechnol.* **14**, 862–866.
4. G. L. Rotino, E. Perri, M. Zottini, H. Sommer, and A. Spena (1997) *Nature Biotechnol.* **15**, 1398–1401.
5. H. J. Klee and C. P. Romano (1994) *Crit. Rev. Plant Sci.* **13**, 311–324.
6. J. A. Ryals, U. H. Neuenschwander, M. G. Willits, A. Molina, H.-Y. Steiner, and M. D. Hunt (1996) *Plant Cell* **8**, 1809–1819.
7. R. A. Creelman and J. E. Mullet (1997) *Plant Cell* **9**, 1211–1223.
8. R. Walden, A. Cordeiro, and A. F. Tiburcio (1997) *Plant Physiol.* **113**, 1009–1013.
9. T. H. Thomas, P. D. Hare, and J. van Staden (1997) *Plant Growth Regul.* **23**, 105–122.
10. A. J. Trewavas and R. Malhó (1997) *Plant Cell* **9**, 1181–1195.
11. P. C. Sehnke and R. J. Ferl (1996) *Curr. Biol.* **6**, 1403–1405.
12. T. W. Grebe and J. Stock (1998) *Curr. Biol.* **8**, R154–R157.
13. I. Herskowitz (1995) *Cell* **80**, 187–197.
14. D. M. Braun and J. C. Walker (1996) *Trends Biochem. Sci.* **21**, 70–73.
15. N. Bechtold, J. Ellis and G. Pelletier (1993) *C. R. Acad. Sci. Paris (Life Sci.)* **316**, 1194–1199.
16. K. Yoon, A. Cole-Strauss, and E. B. Kmiec (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2071–2076.

### Suggestions for Further Reading

17. C. H. Bornman, ed. (1997) Special volume: Hormones, regulating and signalling substances in plant growth and development. *Physiol. Plant.* **100**, 407–738.
18. J.-P. Kader, ed. (1998) Special issue: *Arabidopsis thaliana*. *Plant Physiol. Biochem.* **36**, 1–197.
19. K. Palme, ed. (1994) *Signals and Signal Transduction Pathways in Plants*. Kluwer, Dordrecht, The Netherlands.
20. F. B. Salisbury and C. W. Ross (1992) *Plant Physiology*, 4th ed. Wadsworth Publishing Company, Belmont (CA).
21. J. van Staden, ed. (1997). Special issue: Opinions on current research in cytokinin physiology. *Plant Growth Regul.* **23**, 1–134.

## Plasmalemma

*Plasmalemma* (plasma+Gr. lemma = pellicle) is a membrane structure that is synonymous with the plasma [membrane](#) of cells. There are other usages of the term. In sea urchin eggs, it is the vitelline membrane covering the hyaline plasma layer; upon fertilization the vitelline membrane thickens and becomes the fertilization membrane, which separates from the plasma membrane lying just beneath. In **amoeba**, it is the outermost membrane, which also contains many filamentous structures. In **plants**, it is the external membrane of the cytoplasm; under osmotic stress in hypertonic salt solutions, plant cells undergo *plasmolysis*, where the plasmalemma separates from the cellulose wall of the cell. (See also [Membranes](#).)

### Suggestion for Further Reading

H. Davson (1964) *A Textbook of General Physiology*, 3rd ed. Little Brown, Boston. (Contains a concise description of the various cell membranes and their physiological characteristics.)

## Plasminogen

Plasminogen is one of the most abundant **zymogens** of **serine proteinases** in mammalian species; it is expressed constitutively in the liver and is present in blood at 200 µg/mL. The activated [proteinase](#), plasmin, is one of the most active, because of its broad substrate specificity. The activation system for plasminogen is therefore a powerful and potentially promiscuous **proteolytic** system, which consequently has mechanisms for its strict regulation. The most fundamental of these is the ability of plasminogen to bind to various insoluble matrices (ie, fibrin and certain cell-surface components), a process regulating both the activation of plasminogen (see [Plasminogen Activators](#)) and the inhibition of plasmin.

The binding properties of plasminogen are mediated by its various [kringle domains](#). Plasminogen is a [mosaic protein](#) composed of a 77-residue *N*-terminal **domain**, five kringle modules, and the serine proteinase domain. The kringle modules each contain a binding site for [lysine](#) or other  $\omega$ -amino carboxylic acids. Kringle 1 binds these ligands with the highest affinity (**dissociation constant**  $K_d = 10 \mu M$  for 6-aminohexanoic acid) and represents the archetype for this binding site. It contains both anionic and cationic centers, with an intervening **hydrophobic** depression ([1](#), [2](#)). The cationic center is not present in all kringles, such as kringle 5, and these have a preference for ligands lacking a carboxylate function. Therefore, plasminogen can interact with the side chains of internal lysine residues, as well as those at the *C*-termini of proteins degraded by proteinases with a specificity for basic residues, including **trypsin**, or plasmin itself. The latter mechanism may act to amplify the plasminogen activation system.

The gross conformation of plasminogen is dramatically affected by kringle-binding ligands ([3](#)). In the absence of such ligands, it adopts a compact conformation, due to poorly understood interactions between the *N*-terminal domain and the kringle modules. Exogenous kringle-binding ligands disrupt these interactions and cause an extension of the molecule to approximately double its length. This conformational change can also be accomplished by proteolytic removal of the *N*-terminal peptide,

generating Lys78-plasminogen (to discriminate between them, the native molecule is often referred to as *Glul-plasminogen*). The conformation of plasminogen has a very marked effect on its activation, with the extended or “open” form being activated at least 10-fold more efficiently. The absolute difference in activation rates may be much larger, as the “closed” form is in equilibrium with the “open” form in the absence of exogenous ligands. Thus the physiologically relevant substrate form of plasminogen is thought to be that bound to various surfaces in the “open” conformation.

Plasmin free in solution is very rapidly inhibited by the [serpin](#)  $\alpha_2$ -antiplasmin, which is dependent largely on an initial interaction between the C-terminal region of the inhibitor with the lysine-binding site of plasmin kringle 1 (4). Hence, plasmin inhibition is decreased when these sites are occupied by binding to fibrin or cell surfaces (5). Plasmin and plasminogen have moderate affinities for these surfaces ( $K_d \sim 1 \mu M$ ), with dissociation half-lives of the order of minutes. Thus binding of plasminogen and subsequent dissociation of plasmin, respectively, regulate the activation and inhibition of this proteolytic system.

One final notable feature of plasminogen is its potential to attain proteolytic activity in the absence of the activation cleavage at Arg561–Val562. The non-enzymatic bacterial protein streptokinase binds to the serine proteinase domain of plasminogen and induces a conformational change, with the formation of a functional [active site](#). This has full activity against peptide substrates, but the specificity of the complex against macromolecular substrates is changed from that of plasmin to that of a plasminogen activator. Therefore, “catalytic” amounts of streptokinase can rapidly lead to full activation of plasminogen. Although this property is used clinically to generate plasmin activity *in vivo*, an equivalent endogenous mechanism has not been detected.

#### Bibliography

1. I. I. Mathews, P. Vanderhoff-Hanaver, F. J. Castellino, and A. Tulinsky (1996) *Biochemistry* **35**, 2567–2576.
2. M. R. Rejante and M. Llinás, (1994) *Eur. J. Biochem.* **221**, 939–949.
3. W. F. Mangel, B. H. Lin, and V. Ramakrishnan (1990) *Science* **248**, 69–73.
4. W. E. Holmes, L. Nelles, H. R. Lijnen, and D. Collen (1987) *J. Biol. Chem.* **262**, 1659–1664.
5. V. Ellis, N. Behrendt, and K. Danø, (1991) *J. Biol. Chem.* **266**, 12752–12758.

#### Plasminogen Activators

Mammalian species have two plasminogen activators, urokinase- or urinary-type plasminogen activator (uPA) and tissue-type plasminogen activator (tPA), that are the products of separate genes and members of the chymotrypsin family of serine proteinases. Both are mosaic proteins with a modular structure similar to that of the blood coagulation proteinases to which they are closely related [[Blood Clotting](#)]. The two plasminogen activators provide an excellent example of the regulatory potential made available by the modular construction of these proteinases. Although catalyzing the same reaction, *i.e.* specific hydrolysis of Arg560-Val561 of plasminogen, gross differences in the organization of the N-terminal modules of the plasminogen activators, together with subtle differences in the serine proteinase domain, lead to remarkably different functional properties. These are reflected in their biological roles with tPA being predominantly responsible for plasmin-catalyzed fibrin dissolution, whereas uPA is thought to be responsible for generating plasmin activity in the context of extracellular matrix degradation and thus to be involved in tissue

remodeling and invasive cell migration (e.g. tumor invasion and metastasis; vascular remodeling) (see [Urokinase](#)).

### 1. *tPA*

tPA is composed of five independent domains; an N-terminal fibronectin type-I (finger) module, an EGF-like module, two kringle modules and the serine proteinase domain. It is secreted primarily by vascular endothelial cells as a single-chain glycoprotein of 527 residues ( $M_r$ , 66,000). This single-chain form also undergoes a plasmin-catalyzed “activation” cleavage (Arg275-Ile276) to give a disulfide-bridged two-chain molecule. However, uniquely among the serine proteinases, the catalytic activities of these two forms are very similar. They vary by approximately 10–20-fold in the absence of cofactors, but in the presence of fibrin [see “[Fibrinogen](#)”] as a cofactor their activities are essentially indistinguishable. Molecular modeling and mutagenesis studies suggest that the “active zymogen” nature of tPA is largely due to Lys429 (Lys156 in chymotrypsin numbering); in the absence of a proteolytically exposed N-terminal  $\alpha$ -amino group at Ile276, it can form a salt-bridge with Asp477, substituting for the interaction between the  $\alpha$ -amino group of Ile16 with Asp194 in chymotrypsin (1). The crystal structure of two-chain tPA is compatible with this hypothesis (2).

Fibrin has a dramatic effect on the activities of both forms of tPA. The catalytic efficiency increases up to 1000-fold, as a result of a large reduction in the  $K_m$  (Michaelis constant) for plasminogen and a more modest increase in the value of the catalytic rate constant  $k_{cat}$  (3). More than one mechanism appears to be responsible for these kinetic effects. Fibrin principally enhances plasminogen activation by provides a template for the coincident binding of enzyme and substrate in a ternary complex. In contrast to the binding of uPA to uPAR, the interactions between tPA and fibrin responsible for this binding are complex and not fully elucidated. The fibronectin type-I and the lysine-binding properties of the second kringle module are involved, with sites in the serine proteinase domain also implicated. Secondary to this template effect, fibrin also appears to induce conformational changes in the serine proteinase domain of single-chain tPA, bringing it to catalytic equivalence with two-chain tPA, and increasing its activity against small peptide substrates. Whether these are direct effects or result from the significant interdomain interactions that exist in tPA is presently undetermined.

In addition to its fibrinolytic role, analysis of tPA<sup>-/-</sup> knock-out mice has revealed a role for tPA in the brain, in both normal and neurodegenerative situations. The proteolytic activity of tPA is implicated in the plasticity and reorganization of synapses involved in various learning processes, although it is not known whether plasminogen activation is required in these processes or whether tPA acts on other substrates. tPA is also involved in neuronal cell death following excitotoxin stimulation of the hippocampus and in ischemic stroke. The former process is known to require plasminogen activation, with plasmin acting to degrade laminin, an important extracellular matrix component in the central nervous system (4).

### Bibliography

“Plasminogen Activators” in , Vol. 3, pp. 1865–1866, by Vincent Ellis, Ph.D., University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ; “Plasminogen Activators” in (online), posting date: January 15, 2002, by Vincent Ellis, Ph.D., University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ.

1. K. Tachias and E.L. Madison (1997) *J. Biol. Chem.* **272**, 28–31.
2. D. Lamba, M. Bauer, R. Huber, S. Fischer, R. Rudolph, U. Kohnert and W. Bode (1996) *J. Mol. Biol.* **258**, 117–135.
3. M. Hoylaerts, D.C. Rijken, H.R. Lijnen and D. Collen (1982) *J. Biol. Chem.* **257**, 2912–2919.
4. Z.L. Chen and S. Strickland (1997) *Cell* **91**, 917–925.



### Suggestions for Further Reading

5. K. Danø, P.A. Andreasen, J. Grøndahl-Hansen, P. Kristensen, L.S. Nielsen and L. Skriver. (1985) Plasminogen activators, tissue degradation and cancer. *Adv. Cancer Res.* **44**, 139–266.
6. H.R. Lijnen and D. Collen (1998) "t-Plasminogen activator", in *Handbook of Proteolytic Enzymes*, A.J. Barrett, F. Woessner and N. Rawlings, eds., Academic Press, London, pp 184–190.

### Plasmogamy

Plasmogamy occurs when two different cells fuse with each other and the **cytoplasm** of both is joined. During **fertilization**, plasmogamy occurs when the **sperm** has entered the cell [membrane](#) of the **egg** and a male and a female **pronucleus** develop. However, plasmogamy in animal **reproduction** is unusual in that the amount of cytoplasm of the two **gametes** is disparate. The **oocyte** is one of the biggest cells in the body and contains most of the cytoplasm, whereas the spermatozoon is a very specialized cell that contains only very small amounts.

The length of time it takes the sperm head to fuse with the oocyte cytoplasm depends on the species. In general, fusion takes place during **meiosis** of the oocyte. These divisions are necessary, as the primordial germ cells remain **diploid** until ovulation commences at the beginning of puberty under the influence of follicle-stimulating hormone (FSH) and luteinizing hormone (LH). During growth and oogenesis, the oocyte becomes temporarily **tetraploid** and undergoes two meiotic divisions, resulting in the extrusion of polar bodies. Formation of the first polar body involves a cytoplasmic protrusion at the **animal pole**. The cytoplasm is not divided equally among the two cells, however, and the polar bodies contain less material than the remaining oocyte. The first polar body in a number of different species appears to be a miniaturized cell. The second polar body appears after the second meiotic division and contains a variable amount of cytoplasmic contents. In most cases, however, the fate of the polar bodies is a degenerative process. During this stage of fertilization, the formation of crenations along the vegetal pole occurs, along with redistribution of cytoplasmic compounds.

Upon sperm entry, which is the first step that can be considered a unification in which the cytoplasm are joined, a number of events are triggered. Sperm entry itself is mediated by an adhesion protein named *fertilin* (previously known as PH-30). The  $\alpha$ -subunit of fertilin is involved in mediating the binding of sperm and egg plasma membranes (1); there is also evidence for adhesive activity of an egg b1 [integrin](#) -mediated interaction (2). These gene products are expressed by the ADAM **gene family** and are processed during sperm maturation (3, 4).

The first measurable event during plasmogamy is a rapid depolarization of the plasma membrane on sperm entry. This process is used in some species to prevent more than one spermatozoon from entering the egg. In sea urchins, eggs with a fertilization potential greater than 0 mV are not polyspermic. Maintaining a potential of +5 mV causes a block of fertilization. If the applied voltage is lowered below 0 mV, fertilization will occur and, if the potential is maintained below –30 mV, the eggs will become polyspermic. These observations indicate that a rapid change in the potential of the plasma membrane is a first reaction of the egg to sperm entry. The change in polarity occurs concomitantly with the influx of calcium (see [Calcium Signaling](#)). The release of calcium and of the cortical granules occurs in a sort of wave, starting from the point of sperm entry. Within the cell, the pH is raised, which is assumed to be another factor of egg activation, along with enhancing or

initiating **protein biosynthesis**, **DNA replication**, **messenger RNA** polyadenylation, glucose-6-phosphate dehydrogenase activity, potassium conductance, glycogenolysis, and thymidine uptake. These events go together with the **pronucleus** development and **chromatin** condensation. The cortical granule reaction does not seem to be affected by the raised pH.

**Exocytosis** of the cortical granules is another event that occurs during plasmogamy. The contents of the granules trigger the conversion of the ZP2 and ZP3 **glycoprotein** receptors for sperm to ZP2f and ZP3f, respectively, which are the inactivated forms. In the mouse, rabhlin-3A is involved in the process of calcium-dependent exocytosis on fertilization (5).

There are several explanations for the calcium release mechanism: It can be mediated by **inositol triphosphate** ( $IP_3$ ) or be  $IP_3$ -independent, or it can be mediated by **G-protein** or **tyrosine kinase**-linked enzymes. It has been shown that  $IP_3$  and **diacylglycerol** are involved in a number of **signal transduction** processes and that  $IP_3$  is involved in the autocatalytic cycle of calcium release at fertilization (6). These observations have led to the assumption that the sperm induces an initial production of  $IP_3$ , but the mechanism of  $IP_3$  generation has not been elucidated. In addition to this theory, the production of  $IP_3$  and diacylglycerol is assumed to be linked to proteins, as their involvement was shown in several species in which G-proteins triggered the formation of  $IP_3$ . Up until now, however, no corresponding molecules were found that could support this kind of sperm-induced egg activation. It was also shown, in sea urchins, that ryanodine is able to induce calcium release through ryanodine receptors, which are distantly related to the  $IP_3$  receptors. In experiments with agonists and antagonists of the  $IP_3$  receptor, the involvement of both systems in calcium release was observed. When the  $IP_3$ -induced calcium release was blocked, however, other systems were not able to compensate. At this time, it must be assumed that a number of mechanisms act together to induce calcium release, but the initial triggering mechanism is still not clear.

#### Bibliography

1. J. P. Evans, R. M. Schultz, and G. S. Kopf (1997) *Dev. Biol.* **187**, 94–106.
2. J. P. Evans, G. S. Kopf, and R. M. Schultz (1997) *Dev. Biol.* **187**, 79–93.
3. R. Yuan, P. Primakoff, and D. G. Myles (1997) *J. Cell. Biol.* **137**, 105–112.
4. D. G. Myles, and P. Primakoff (1997) *Biol. Reprod.* **56**, 320–327.
5. N. Masumoto, T. Sasaki, M. Tahara, A. Mammoto, Y. Ikebuchi, K. Tasaka, M. Tokunaga, Y. Takai, and A. Miyake (1996) *J. Cell. Biol.* **135**, 1741–1747.
6. M. J. Whittaker, and R. F. Irvine (1984) *Nature* **312**, 636–639.

#### Suggestion for Further Reading

7. F. J. Longo, ed. (1997), *Fertilization*, 2nd ed., Chapman & Hall, London, Weinheim, New York, Melbourne, Madras, pp. 49–65; provides further information concerning early events during plasmogamy.

#### Plastocyanin

Plastocyanins belong to the family of cupredoxins [ie, relatively small (~10 kDa, 97–104 amino acid) blue single copper proteins, including **azurins**] involved in **electron transfer** (1). They are

found in the thylakoid lumen of **chloroplasts**, where they function as electron carriers in oxygenic [photosynthesis](#) from cytochrome *f* in the  $b_6f$  complex to P700<sup>+</sup> of photosystem I (2).

The three-dimensional (3-D) [protein structure](#) of plastocyanin isolated from the poplar tree was the first one of the family to be determined, in 1978 (3). Since then, structures of other (algae and plant) plastocyanins have been determined, either by [X-ray crystallography](#) or by nuclear magnetic resonance (NMR) methods. In parallel, the oxidized poplar plastocyanin structure has been refined to 1.22 Å resolution (4), and its reduced form structure has also been determined (5). Despite sequence **divergence** among plastocyanins of algae and vascular plants, their 3-D structures are remarkably conserved. It is a barrel-shaped molecule, constructed of eight antiparallel **b-strands** (Fig. 1). The metal-binding site is at one end of the b-barrel and consists of (a) a copper ion coordinated in a trigonal bipyramidal geometry with the thiolate of a [cysteine](#) residue and two imidazole nitrogens of [histidine](#) residues in the trigonal plane and (b) a methionine sulfur and a remote carbonyl oxygen atom occupying axial positions. This coordination sphere is typical of the blue or Type 1 (T1) copper site in all cupredoxins (1). The protein's intense blue color is due to a  $\pi S \rightarrow Cu(d_{x^2-y^2})$  ligand-to-metal charge transfer involving the thiolate ligand as the electron donor, and it has a molar extinction coefficient of  $\sim 3000 \text{ M}^{-1}\text{cm}^{-1}$ . This is more than 100 times larger than that found for simple Cu(II) complexes (6, 7). A second characteristic property associated with the blue T1 Cu(II) site is an EPR spectrum (see **Electron paramagnetic resonance**) displaying a hyperfine splitting in the  $g_{\parallel}$  region, due to interaction of the copper nuclear and electron spins, which is exceptionally narrow ( $\sim 0.008 \text{ cm}^{-1}$  and approximately 50% smaller than those of ordinary copper complexes. This is attributed to delocalization of the unpaired  $Cu(d_{x^2-y^2})$  electron onto the Cys(S) *pp* orbital, thus reducing the nuclear–electron interaction. The unique spectroscopic properties are due to the T1 Cu (II) site structure and are assumed to be related to the electron transfer function of cupredoxins in general (6-9). The relationship between these unusual properties and the electron transfer reactivity is a central issue in studies of biological [electron transfer proteins](#).

**Figure 1.** The three-dimensional structure of plastocyanin. The copper center with its ligands is seen near the top of the molecule, just below the “Northern” hydrophobic patch. The negatively charged residues are centered at the right-hand side of the figure.



The **oxidation–reduction potentials** of Cu(II)/Cu(I) in plastocyanins (~350 mV) are higher than generally observed for copper complexes. The metal–protein interactions within this site maintain the Cu(II) away from a square-planar geometry, toward a distorted tetrahedral geometry. This explains the relative stabilization of the Cu(I) state relative to Cu(II). Structure-imposed Cu → L p-backbonding has been proposed to account further for stabilization of the cuprous state, because strong p-interaction with the *dp* orbitals results in an increase of the ligand field strength (9-11).

A flat **hydrophobic** surface present at one end (“north”) and a patch with pronounced negatively charged residues on one side of the protein barrel (east) characterize the plastocyanin surface. Both the negative and the hydrophobic patches were proposed to be involved in the interactions of plastocyanin with its physiological reaction partners (12, 13). In eukaryotic plastocyanins, this negatively charged patch surrounds a nearly conserved residue, Tyr83, that has been implicated to act in one proposed electron transfer pathway. The shape and charge distribution around the negative patch make it a unique region on the surface of the molecule in both plant and algae plastocyanins. The hydrophobic surface located at the “north” end of the molecule surrounds His87, the only solvent-exposed copper ligand and hence an obvious site for electron transfer (14). This patch consists of at least eight amino acid residues that are conserved in eukaryotic plastocyanins, with no charged or polar residues in this region. The solvent-accessible surface of this hydrophobic patch is approximately 550 Å<sup>2</sup>. The side chain of Phe35, the only large residue in this area, lies along the surface of the molecule, without extending significantly into the solvent. The flatness of this patch suggests that it is important for the interaction of plastocyanin with its reaction partners. Significantly, all the cupredoxins of known structure contain a hydrophobic area surrounding their solvent-exposed copper-liganding histidine residue.

Considerable details of the structure–function relationship of plastocyanin were attained in recent years. These have established that the electron transfer activity is indeed confined to the exposed

histidine residue in the flat patch, while the tyrosine residue surrounded by the negative patch has apparently a role in formation and proper alignment of plastocyanin in its electron transfer complex formed with P700 in the [photosynthetic reaction center](#) of photosystem I (15, 16).

## Bibliography

1. E. T. Adman (1991) *Adv. Protein Chem.* **42**, 145–198.
2. M. Haehnel (1984) *Annu. Rev. Plant Physiol.* **35**, 659–693.
3. P. M. Colman, H. C. Freeman, J. M. Guss, et al. (1978) *Nature* **272**, 319–324.
4. J. M. Guss, H. D. Martunic and H. C. Freeman (1992) *Acta Crystallogr.* **1348**, 790–811.
5. J. M. Guss, P. R. Harrowell, M. Murata, V. A. Norris, and H. C. Freeman (1986) *J. Mol. Biol.* **192**, 361–387.
6. B. G. Malmström, B. Reinhammar, and T. Vänngård (1968) *Biochim. Biophys. Acta* **156**, 67–76.
7. E. I. Solomon, J. W. Hare, and H. B. Gray (1976) *Proc. Natl. Acad. Sci. USA* **73**, 1389–1393.
8. E. I. Solomon, M. J. Baldwin, and M. D. Lowery (1992) *Chem. Rev.* **92**, 521–542.
9. H. B. Gray and E. I. Solomon (1981) In *Copper Proteins*, Vol. **3** (T. G. Spiro, ed.), Wiley, New York, pp. 1–39.
10. T. Pascher, B. G. Karlsson, M. Nordling, B. G. Malmström, and T. Vänngård (1993) *Eur. J. Biochem.* **212**, 289–296.
11. C. S. St. Clair, W. R. Ellis, and H. B. Gray (1992) *Inorg. Chim. Acta* **191**, 149–155.
12. O. Farver, Y. Shahak, and I. Pecht (1982) *Biochemistry* **21**, 1885–1890.
13. D. J. Cookson, M. T. Hayes, and P. E. Wright (1980) *Biochim. Biophys. Acta* **591**, 162–176.
14. K. Sigfridsson, M. Sundahl, M. J. Bjerrum, and O. Hansson (1996) *J. Bioinorg. Chem.* **1**, 405–414.
15. W. Haehnel, T. Jansen, K. Gause, R. B. Klösgen, B. Stahl, D. Michl, B. Huvermann, M. Karas, and R. G. Herrmann (1994) *EMBO J.* **13**, 1028–1038.
16. M. Hippler, J. Reichert, M. Sutter, E. Zak, L. Altschmied, U. Schröer, R. G. Herrmann, and W. Haehnel (1996) *EMBO J.* **15**, 6374–6384.

## Platelet-Derived Growth Factor

Platelet-derived growth factor (PDGF) is a [growth factor](#) that exerts multiple biological activities in many cell types and tissues. PDGF exists in three forms. Each form consists of a homo- or heterodimeric combination of two genetically distinct but structurally related polypeptide chains, A and B. PDGF chains play important roles in connective tissue development. Mice deficient in B chain lack mesangial cells in their kidneys and develop cardiac and other microvascular abnormalities. PDGF A chain-deficient mice manifest generalized lung emphysema due to the absence of alveolar myofibroblasts. PDGF isoforms exert their effects on target cells by binding to two structurally related protein [tyrosine kinase receptors](#), alpha (α) and beta (β). The alpha receptor binds to the A or B chain, whereas the beta receptor binds only to the B chain. Differences in their biochemical properties and the relative abundance of receptor subunits and or PDGF chains determine the biologic responses to PDGF.

Upregulation of PDGF chains and receptors has been implicated in the pathogenesis of diverse conditions, including carcinogenesis, atherosclerosis, glomerulonephritis, and pulmonary and hepatic

fibrosis. Molecular and genetic approaches have demonstrated an essential role for PDGF and its receptors in the development of connective tissue in the kidney, lung, cardiovascular, skeletal, and nervous systems.

## 1. Historical background

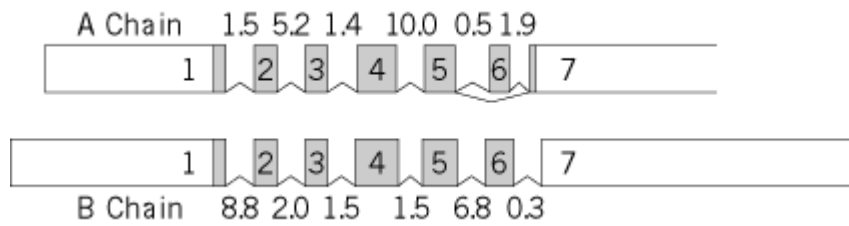
PDGF was first recognized as serum-associated growth-promoting activity for smooth muscle cells (1), fibroblasts (2), and glial cells (3). It was subsequently purified by different groups from human platelets (reviewed in Ref. 4). The 30-kDa protein was found to consist of two separate polypeptide chains linked by [disulfide bonds](#). The two chains were eventually separated by reverse-phase HPLC and referred to as A and B chains, based on their order of elution. On sequencing, the B chain turned out to be identical to p28<sup>SIS</sup>, the transforming protein of simian sarcoma virus (SSV). This finding provided the first link between growth factors and **oncogene** products and suggested a mechanism by which oncogene products subvert the mitogenic pathway and stimulate cell growth and transformation. Platelet PDGF now appears to be a mixture of PDGF-AB, -BB, and -AA (5-7). The [cloning](#) and availability of [complementary DNAs](#) encoding for both the A and B chains have allowed the production and isolation of the [recombinant proteins](#) using different expression systems. The subtle differences in the biological activities of the various purified PDGF isoforms could later be explained by their receptor-binding specificity and subtle differences in their [signal transduction](#) pathways.

Ligand-binding and **affinity-labeling** techniques demonstrated the presence of a specific 185-kDa PDGF receptor (4). Furthermore, PDGF stimulated an intrinsic protein **tyrosine kinase** activity of the receptor and resulted in the **phosphorylation** of specific cytoplasmic substrates. The murine beta receptor was cloned by Yarden et al. (8) and by Rönstrand et al. (9), followed by the alpha receptor (10, 11).

## 2. Structure and Function of PDGF Isoforms

PDGF A and B are hydrophilic soluble proteins that share approximately 51% amino acid identity in the mature human proteins, and there is 35% similarity between *Xenopus* PDGF A and mouse PDGF B. Each chain is, however, more highly conserved between species, with *Xenopus* and mouse PDGF A sharing 71% similarity and mouse and human PDGF B, 89%. Their gene structure and organization are also remarkably similar (Fig. 1) (12). Both chains are synthesized as larger precursor molecules that are extensively processed. The A chain is the product of a single gene located on chromosome 7 that contains seven exons and gives rise to one of two differentially spliced variants (see [Alternative Splicing](#)). The “long” variant, a **prepropeptide** of 211 residues, is synthesized with a [signal peptide](#) of 20 residues, a propeptide sequence of 66 residues, and a mature chain of 125 residues. The “short” 196-residue variant shows a 20-residue signal sequence, a 66-residue propeptide, and a 110-residue mature form (13-15). The difference between the long and short variants is the result of alternative exon usage, with the long form using exon 6 (18 residues), but not exon 7, and the short form using exon 7 (3 residues), but not exon 6 (14). Within exon 6 lies an approximately 10-residue sequence that signals cell retention (14, 16, 17). The short variant contains no retention sequence and is secreted (18, 19). No mechanism for C-terminal processing of the long form of the A chain has been reported, and it is not clear whether this variant is secreted (20, 21). This sequence mediates interaction with heparan sulfate proteoglycans. Therefore differential compartmentalization of the chains may have important biologic implications. PDGF A is probably subject to both O-glycosylation and N-glycosylation, whereas PDGF B has only O-glycosylation. Recombinant PDGF binds and is biologically active, however, suggesting that glycosylation is not required for its activity (22).

**Figure 1.** Schematic illustration of the organization of the PDGF A- and B-chain genes.



The B chain is also the product of a seven-exon gene located on chromosome 22. The B-chain gene is now known to be identical to the human *c-sis* gene, the normal human cell counterpart to the monkey *v-sis* (simian sarcoma) virus gene (23). The protein coded for by *c-sis* is a 27-kDa, 241-residue prepropeptide with a 20-residue signal sequence, 61-residue propeptide, and 160-residue “mature” polypeptide (13, 24). C-terminal cleavage of the mature B chain occurs in two stages, with a **trypsin-like** cleavage after residues 111–160, followed by a [carboxypeptidase](#) cleavage of the remaining arginine at residue 110 (20, 23, 25), resulting in a final mature product of 109 residues and 12 kDa. A retention sequence (residues 212–216) has also been identified in the B-chain C-terminus. As with the long variant of chain A, failure to remove this peptide also results in B-chain glycosaminoglycan retention (16, 17, 26). Dimerization of the A and B chains involves two interchain disulfide bonds mediated by eight conserved cysteine residues; four of these residues are essential for transforming activity (15, 27). The [X-ray crystallography](#) structure of PDGF-BB dimer (28) consists of four irregular antiparallel **beta-strands** and a 17-residue amino-terminal tail. Of the eight disulfide-bonded cysteine residues, six form a knotted arrangement and two form two interchain disulfide bonds. The edges of the four-stranded **beta-sheets** form the dimer, which results in the majority of intersubunit contacts being between the first two strands of the b-sheet and with the N-terminal tail. The total [accessible surface](#) buried is estimated to be 2200 Å<sup>2</sup>, and most of the buried residues are **hydrophobic** in nature.

There is little to modest expression of PDGF or PDGF receptors in normal adult tissue, but expression is markedly upregulated during disease processes and development. PDGF expression in cultured cells is regulated by the culture conditions, as well as soluble factors and cell matrix organization. The expression of PDGF chains can be regulated at the levels of [transcription](#), [translation](#), and [post-translational modification](#), as well as [protein secretion](#). Released PDGF is rapidly sequestered by binding to  $\alpha$ -2 [macroglobulin](#), a circulating 720-kDa glycoprotein that binds B-chain-containing PDGF isoforms, BB and AB (29). PDGF binding results in its rapid clearance via  $\alpha$ -2 macroglobulin receptors (30).

Transcription of the PDGF genes can be modulated by diverse extracellular stimuli, such as mitogens, cytokines, tumor promoters, shear stress, and oxygen concentration. The response of the PDGF genes to these stimuli is complex, as is the cell type- and developmental stage-dependent transcription, depending on the interaction between specific DNA elements and regulatory [transcription factors](#). PDGF A and B promoters contain a [TATA box](#) located at approximately 30 bp upstream of the transcription initiation site. The PDGF B promoter contains a phorbol ester [response element](#) and sequences that can bind Ets-, AP1-, and Sp1-like transcription factors. It also contains a shear stress response element that interacts functionally with the transcription factor NF- $\kappa$ B. It is noteworthy that these sites often overlap, making regulation of transcription of genes even more complex. Transcription of PDGF B promoter also requires the presence of the gene's entire first intron. The PDGF B-chain gene is subject to negative transcriptional regulation by a mechanism that has not yet been elucidated.

Contrary to the PDGF B promoter, the PDGF A promoter is exceptionally G/C-rich. Eight consensus binding sites for the transcription factor Sp1 are located within the first 1000 bp upstream of the transcription initiation site. Moreover, the first 600 bp of the promoter contain at least six binding sites for transcription factors Egr-1 and Wilm's tumor-associated WT-1. Depending on the cell type,

Egr-1 can either activate or repress the activity of the PDGF A promoter. A PMA-responsive element and [serum response element](#) are also present in this promoter. A negative transcription regulatory element is located between  $-1.9$  and  $-0.9$  kb (31). **Zinc-finger** transcription factors, including Sp1, Sp3, and Egr-1, have recently been shown to mediate high B-chain expression in aortae of newborn rats (32).

Two transcripts of the PDGF B gene are induced in endothelial cells in response to [cycloheximide](#) and [transforming growth factor](#) b1 (TGF-b1). Transcription initiation of the 2.8-kb messenger RNA starts 15 nucleotides upstream of the translation start site, while a 3.0-kb transcript starts at 200 nucleotides further upstream. Similar truncation of the 5' untranslated region contributes to the marked increase in the 2.6-kb transcript in the developing rat brain (33). The expression of a truncated PDGF B mRNA species provides a mechanism of escape from the potent translation inhibitory action of the exon 1-derived 5' untranslated sequence. Choriocarcinoma cell lines express a 2.6 kb PDGF mRNA species; transcription is initiated at an alternative promoter located within intron 1 and lacks the coding sequence for the signal peptide of the PDGF B-chain precursor. The resulting protein product may be targeted directly to the nucleus in the presence of a nuclear localization signal encoded by exon 6 of the PDGF B gene.

The PDGF A gene uses two **promoters**. The regular transcription start site is located approximately 35 bp upstream of the TATA box, whereas an alternative initiation site is located 470 bp downstream of the first site, but still 380 bp upstream of the translation initiation site.

PDGF chain expression may be also regulated at the level of mRNA stability, as well as at the level of translation. In general, the varying levels of PDGF B mRNA result from differences in transcription rate of the gene, rather than from differences in mRNA stability. The PDGF B mRNA half-life varies between 1 and 3.5 h. The 1-kb-long leader of the 3.5-kbp PDGF B transcript can potentially inhibit translation *in vitro*. The inhibitory activity is partly relieved during megakaryocytic differentiation of K562 cells, when transcription of the PDGF B gene is strongly induced. The 179-nucleotide-long sequence immediately upstream of the PDGF B-gene open reading frame is important for the relief of translation repression. The PDGF A 5'-untranslated sequence contains three ATG codons, potential [start codons](#), the function of which are unknown. The region between +99 and 184 bp relative to the PDGF A transcription start site, which lacks an ATG codon, inhibits translation in rhabdomyosarcoma cells, an effect overcome by sequences located between +184 and +338 bp. In summary, PDGF chains are regulated in a highly complex fashion at multiple levels in a cell- and tissue-specific manner.

Alternative splicing of PDGF mRNA was first recognized as a regulatory mechanism of PDGF A-chain transcripts. The resulting PDGF A mRNA species are expressed in a variety of normal tissues and transformed cells and in diverse mammalian species (34). Alternative splicing of PDGF B mRNA occurs in a recombinant retrovirus that expresses the wild-type PDGF B chain and induces fibrosarcoma in mice. Provirus derived from the tumors lack a 149-base sequence derived from exon 7 as a result of alternative splicing (35).

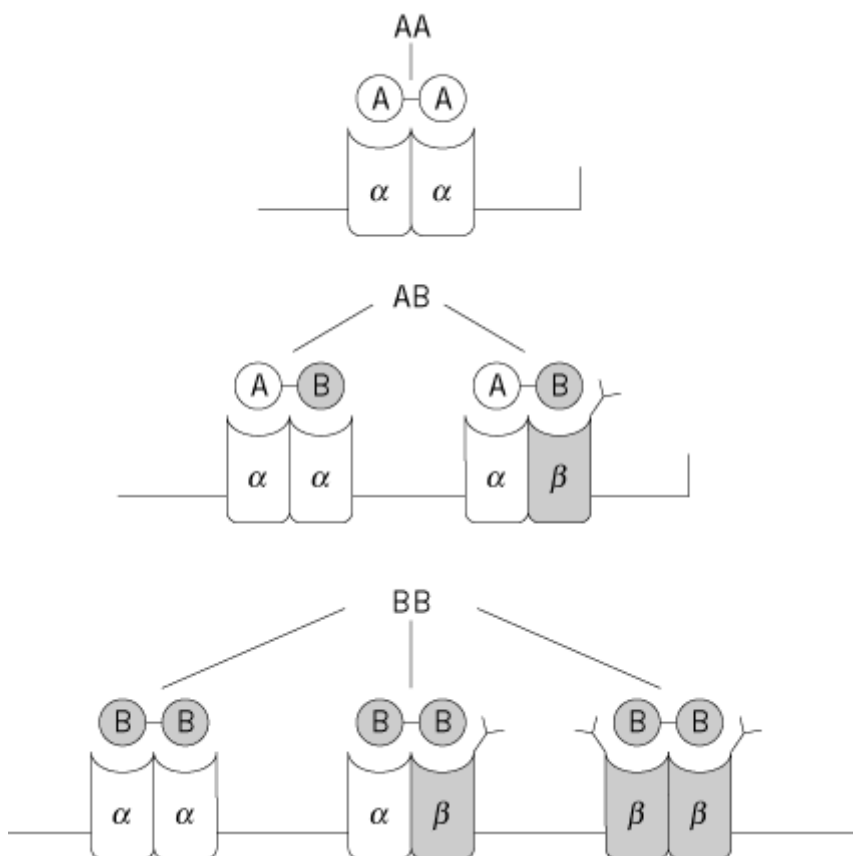
There is evidence for posttranslational modulation of PDGF. PDGF protein activity may be regulated by an extracellular glycoprotein, a secreted protein, acidic and rich in cysteine (SPARC), which binds to the B chain-containing isoforms and inhibits binding to the receptors (12). [Thrombin](#) stimulates release of PDGF B by selective **proteolytic** cleavage of the membrane-associated PDGF precursor protein (36). It is evident that regulation of PDGF chain expression is complex and that the regulation of the two chains is discordant and cell- and tissue-specific. There is very little information available as to the transcriptional machinery that regulates PDGF A or b chain *in vivo* in disease or development.

### 3. PDGF Receptors



Two distinct human transmembrane binding proteins for PDGF have been identified, a 170-kDa a-receptor (PDGF-Ra) (37) and a 190-kDa b-receptor (PDGF-Rb) (38). The two receptor proteins are structurally related and consist of an extracellular portion containing five **immunoglobulin**-like domains, a single transmembrane region, and an intracellular portion with a protein tyrosine kinase domain. The two chains have 44% overall sequence identity (39, 40). A functional PDGF receptor is formed when the two chains of a dimeric PDGF molecule each bind one of the above receptor molecules, resulting in bring them into proximity and their dimerization, and activation (Fig. 2). In addition, a 90-kDa soluble form of PDGF-Ra, consisting of the extracellular domain, has been found in cell culture medium and in human plasma (41). Although each subunit of dimeric PDGF binds to one receptor monomer, it is unclear whether these PDGF subunits need to be covalently linked. Recent evidence suggests that noncovalently linked B chains are able to activate the PDGF receptor (42). PDGF-Ra binds each of the three forms of PDGF dimers with high affinity. Although PDGF-Rb binds both PDGF-BB and PDGF-AB with high affinity ( $K_d = 0.5\text{pM}$  and  $1\text{--}2.5\text{ nM}$ , respectively), it has not been reported to bind to PDGF-AA (43, 44). The apparent high-affinity binding of the AB dimer to the b receptor must be interpreted with caution, however. Although PDGF-AB can bind to mutant 3T3 cells expressing only b receptors, it requires hundred-fold more PDGF-AB to dimerize the b receptors and activate the cells than is required for cells also expressing a receptors. This required concentration is probably not physiologically relevant (44).

**Figure 2.** Dimerization of PDGF receptor subtypes (a and b) by PDGF isoforms.



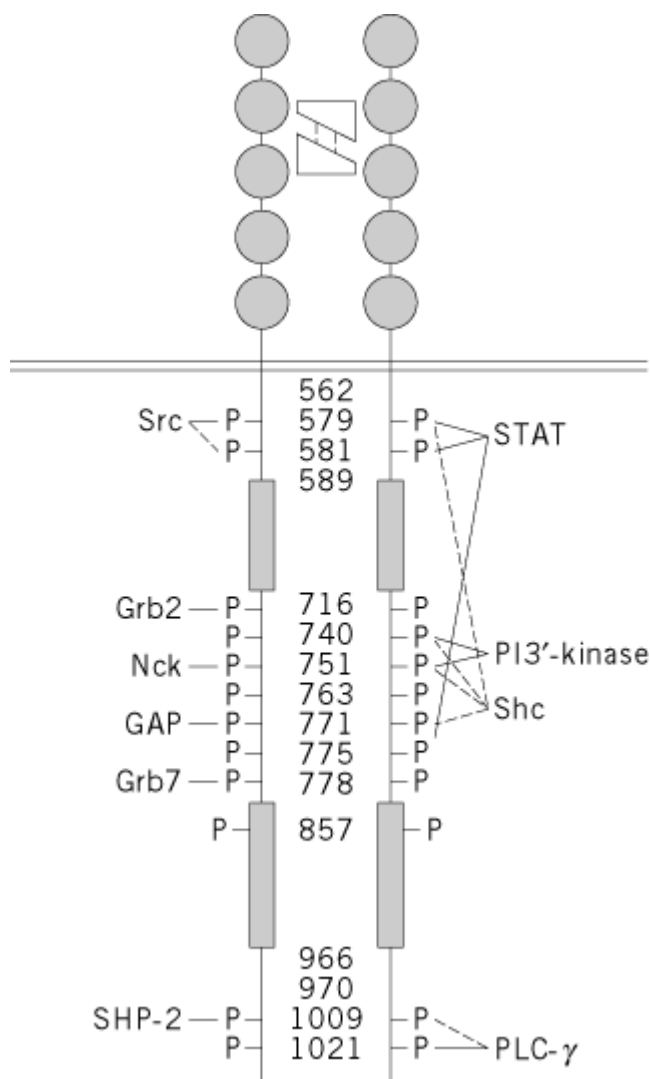
PDGF receptors are differentially expressed and regulated in a cell- and tissue-specific manner. Cells known to express only a receptors include oligodendroglial progenitors, liver endothelial cells and mesothelium (43), and platelets (45). Cells expressing only b receptors include central nervous system capillary endothelium, neurons, and monocytes/macrophages (46). Cells showing coincident expression of a and b receptors include smooth muscle cells, fibroblasts (43), Schwann cells (47),

and mesangial cells (48). The relative abundance of the a and b receptors dictate the biologic activities of the PDGF isoforms. PDGF receptors, particularly the b receptor, are upregulated in inflammatory, proliferative, and fibrotic diseases. In [tissue culture](#) studies, very few factors that regulate the a or b receptor have been identified (49). The molecular mechanisms that regulate receptor gene expression are not known. Regulatory sequences in the 5'-flanking region of the murine b receptor has been cloned. A 2.0-kb sequence appears to confer constitutive b-receptor expression in transgenic animals using reporter genes (50-52).

The signals transduced by the two receptors and their biological effects share similarities, but are not identical (53). Some effects are specific, however; for example, cell-surface [actin](#) reorganization or ruffling is mediated solely by the b receptor. PDGF binding to cells induces receptor internalization and downregulation (54). PDGF stimulation follows activation of its intrinsic tyrosine kinase activity, with phosphorylation of several substrate proteins, including the receptor itself (55). This ligand-induced autophosphorylation is an intermolecular mechanism (56). After ligand stimulation, human PDGF-Rb has at least eleven phosphorylated tyrosine residues (Fig. 3). Two of these are in the juxtamembrane region of the receptor (Tyr579 and Tyr581), seven are in the kinase insert sequence of the receptor (Tyr740, Tyr751, Tyr762, Tyr763, Tyr771, Tyr775, and Tyr778), another is in the second tyrosine kinase domain (Tyr857), and the last two are in the cytoplasmic tail downstream of the tyrosine kinase domain (Tyr1009 and Tyr1021) (34, 57-60). The autophosphorylated tyrosine residues serve as binding sites for several signaling proteins, which contain an **Src** homology-2 (**SH2**) domain. One of these proteins is **phospholipase C-g1** (PLC-g1), which binds directly to phosphorylated directly to Tyr1009 and Tyr1021. After binding, PLC-g1 is tyrosine-phosphorylated by the receptor, which is required for enzymatic activity. Activation of PLC-g1 by PDGF-R stimulates hydrolysis of PtdInsP<sub>2</sub> (phosphatidylinositol diphosphate) to generate InsP<sub>3</sub> (inositol trisphosphate) and DAG (diacylglycerol), which mobilize intracellular calcium and activate protein kinase C, respectively (61) (see [Inositol Lipids and Phosphates](#)). Mutation of tyrosine residues that specifically bind PLC-g1 in the PDGF receptor results in loss in phosphoinositide turnover, as well as loss of mitogenesis (62, 63). Phosphatidylinositol-3 kinase (PI-3K) is a heterodimer of a 110-kDa catalytic and a 85-kDa regulatory subunit that consists of only SH2 and SH3 domains (64). The SH2 domain of PI-3K binds to phosphorylated Tyr740 and Tyr751 in the PDGF receptor-kinase insert (57, 63). After binding, PI-3K becomes activated and produces D3-phosphorylated inositides. Initial reports suggested that protein kinase C-z could be the downstream target of PI-3K, since D3 inositides can activate this enzyme *in vitro* (65). It is now established that these lipids serve to recruit another **serine/threonine protein kinase**, PKB/Akt, to the membranes, where it is phosphorylated and activated by PDK1 and PDK2. Akt has been shown to mediate the effects of PI-3K on cell survival and differentiation, glycogen metabolism, and protein synthesis (66). Mutation of the two tyrosine residues that bind PI-3K results in loss of PDGF-induced cell proliferation (67). PI-3K has also been implicated in receptor internalization (68) and in PDGF-induced chemotaxis and actin reorganization (69), but some responses are cell-specific. PDGF also stimulates the **Ras** pathway. A protein tyrosine phosphatase, SHP-2 (PTP-1D), can bind to the phosphorylated Tyr1009 at the carboxy-terminal end of the receptor, and serve as the docking site for a Grb2/SOS complex. Grb2 is an adaptor molecule, composed of only SH2 and SH3 domains and devoid of catalytic activity, and SOS is a [guanine nucleotide exchange factor](#) for Ras, promoting the exchange of GDP to GTP on Ras, thereby activating it (70). Ras, in turn, recruits and activates another serine/threonine protein kinase, Raf-1 (71), whose activity is increased in PDGF-stimulated cells (72, 73). Activated Raf-1 activates the [MAP kinase](#) cascade. Ras is important in the mitogenic signal of PDGF (74). It is noteworthy that Ras-GAP (for GTPase-activating protein), the enzyme that terminates the Ras signal, can also be recruited on the activated PDGF receptor (on Tyr771), probably allowing a fine tuning of the Ras signal. The two tyrosine residues in the juxtamembrane region of the PDGF receptor are a docking site for some members of the Src family of protein tyrosine kinases, namely, *c-src*, *c-fyn*, and *c-yes* (75). The role of Src family of tyrosine kinases in regulating phases of the [cell cycle](#) is the subject of intense investigation. For instance, Src is not required for PDGF-Rb to induce mitogenesis (76), but it is recruited and activated by PDGF-Ra, and apparently phosphorylates Shc (77). Nck is another adaptor molecule, like Grb2, made of SH2 and

SH3 domains, that binds to the activated PDGF receptor (78). Again, its precise function in PDGF signaling is unknown. Recently, a novel, Ras-independent pathway has been described. Janus kinase-1 (JAK-1) as well as signal transducer and activator of transcription (STAT)-1 and -3, are also activated by PDGF (see [JAK/STAT Signaling](#)). STAT proteins contain SH2 and SH3 domains and are tyrosine-phosphorylated; they dimerize and translocate to the nucleus, where they may induce specific subset of genes. The JAK family of tyrosine kinases are involved in the activation of STAT proteins. However, STAT1a is recruited and can be directly phosphorylated by PDGF-R (79).

**Figure 3.** Schematic illustration of the interaction between the autophosphorylated PDGF-Rb and downstream signal transduction molecules. A ligand-induced dimeric receptor complex is shown. Ig-like domains in the extracellular region are represented by filled circles. The numbers of all intracellular tyrosine residues are indicated. The kinase domain is represented by filled rectangles. The specificity in the interaction between SH2 domain molecules and autophosphorylated tyrosine residues is indicated; solid lines indicate high affinity interactions, and broken lines interactions of lower affinity.



PDGF signaling is negatively regulated by a class of protein-tyrosine phosphatases (PTP) called low molecular weight (LMW)-PTP. These LMW-PTP are cytosolic enzymes that translocate to the membrane fraction on stimulation by PDGF, where they dephosphorylate PDGF-R (80). LMW-PTP interfere mainly with PDGF activation of the Src and JAK-STAT pathways (81).

The role of PDGF receptor signaling in mediating specific biologic responses in different cells or different environments remains an enigma. Much work is still needed to understand the [protein-protein interactions](#) that may be involved, as well as identification of additional molecules.

#### 4. Biology of PDGF

Aberrant expression of PDGF or its receptors is likely to be involved in the stimulation of the growth of certain tumors. However, it also plays a part in the development of certain nonmalignant disorders involving excess of cell migration and proliferation. In considering the role of PDGF in different disease processes, it is important to take into account the fact that the level of the ligands and receptors can be modulated by a network of other cytokines or soluble factors.

##### 4.1. Strategies for Genetic Analysis of PDGF Function

Targeted disruption of the genes coding for PDGF A and B chains, as well as the a and b receptors, has been accomplished. The PDGF-B and -Rb null phenotypes are embryonically lethal and strikingly similar ([82](#), [83](#)). At late gestation (E17-19), homozygous mutants develop general hemorrhaging and edema. At the time of birth, many homozygotes still display life signs. The data thus far indicate the existence of a developmental defect of the blood vessel wall that appears to be the major cause of bleeding. Both PDGF-B and PDGF-Rb null mutant embryos display a developmental defect of the kidney glomerulus that consists of a complete lack of mesangial cells, smooth-muscle-like vascular pericytes that regulate blood flow in glomeruli, filtering units of the kidney. A capillary tuft is not formed. Rather, a few, or a single, microaneurysm-like capillary loop fill out the Bowman's space. Basement membrane, endothelial and epithelial cells are present. Heart defects are also seen in the PDGF B null mutant embryos but not in the PDGF-Rb null mutants, suggesting that PDGF-Ra, which is widely expressed in the heart, may compensate for the absence of PDGF-Rb. Cerebral vessel microaneurysms are also present in the brain of the PDGF B chain null mutants ([84](#)). Defect or absence of vascular pericytes derived from connective tissue, which are essential for the normal development and integrity of the capillary wall, seems to be the mechanism underlying cutaneous bleeding (purpura) and microaneurysm formation in the mutant embryos. Recent studies utilizing chimeric mice from b-receptor-null and wild-type cells indicate that their receptor plays a role in recruitment of smooth muscle cell progenitors to sites of epithelium and endothelium ([85](#)).

PDGF-A null mice were found to develop until birth and live for a few weeks postnatally, at which time they appeared severely growth retarded. These mice suffer from generalized lung emphysema, with resulting atelectasis mixed with hyperinflated regions ([86](#)). Genotyping of large numbers of offspring at various embryonic ages reveals that there is also a loss of PDGF-A null embryos during early gestation ([86](#)). There is a loss of approximately 50% of the PDGF-A null homozygotes before E10, but the causes of this loss and the mutant phenotype at this early stage are currently unknown. Following implantation, PDGF-A and -Ra soon adopt an appositional expression pattern; later in development, PDGF-A appears generally expressed in epithelial tissues and PDGF-Ra in mesenchymal tissues ([87](#), [88](#), [89](#)). It is likely that the early embryonic loss of PDGF-A null homozygotes interferes with epithelium-to-mesenchyme signaling. The PDGF-A null embryos surviving beyond E10 appear slightly growth-retarded at E16-P2, after which point the growth retardation becomes increasingly severe. From approximately 2 weeks of age, the PDGF-A null homozygotes are invariably less than 50% the size of heterozygotes or wild-type litter mates. At 3 weeks of postnatal age, the PDGF-A null mice have invariably developed a severe generalized lung emphysema. Histologically, grossly enlarged distal air spaces are seen without signs of subcompartmentalization by alveolar walls. The mature alveolar septa contain alveolar smooth muscle cells (also called contractile interstitial cells) embedded in deposits of elastin. In the PDGF-A null lungs, several defects associated with the alveolar smooth muscle cells were found. A population of PDGF-Ra-positive cells was identified in the normal E18.5 lung, which had the expected distribution of a putative alveolar smooth muscle cell progenitor. Such cells were absent from the mutant lungs. However, many other sites of PDGF-Ra expression, in both lung and other tissues, remained unaffected in the mutant mice. It was therefore proposed that alveolar deficiency in

PDGF-A null lungs reflects the failure of a population of PDGF-Ra-positive smooth muscle cell progenitors to develop during late stage gestation of the mouse (86).

The phenotype of the homozygous PDGF-Ra mutants differs somewhat from the patch mutant mice, which harbor several deletions, including the PDGF-Ra gene. Homozygotes of the PDGF-Ra mutants die during embryonic development and exhibit incomplete cephalic closure similar to that observed in a subset of the patch mutants. There is also increased [apoptosis](#) of migrating neural crest cells and deficiency of myotome formation, which leads to alterations in mutant vertebrae, ribs, and sternum. This phenotype also differs from that of PDGF-A null mutants, suggesting the existence of compensatory mechanisms probably mediated by the PDGF B chain and b receptor.

## 5. Concluding Remarks

PDGF and its receptors play an integral role in inflammatory, proliferative, and fibrotic disorders. Its role in carcinogenesis is based on circumstantial evidence and has received a new impetus with the recent discovery of the beta receptor translocation in leukemia.

The role of signal transduction pathways in mediating specific biologic activities is a subject of intense investigation, and the development of specific inhibitors of the receptor and of the signaling molecules may open new avenues of therapeutic potential. Therapeutic trials utilizing selective inhibitors of PDGF-R signaling are currently underway in humans with soft tissue tumors and should provide valuable information as to the contribution of PDGF to tumor growth.

Analysis of PDGF ligands and receptor-deficient mice has helped to clarify their role in development. One major interesting finding in the single gene targeting approach, ie. ligand or receptor, is the length of survival of embryos lacking one such gene, considering the widespread distribution of PDGF chains and their receptors at early stages of development. Moreover, the phenotype of PDGF A mutant mice differs from that of the a-receptor mutants. The phenotype of the B chain mutant mice is also somewhat different from that of the b-receptor mutants. This is most likely due to the fact that signaling by the remaining ligand or receptor can still occur in any null embryo. In the PDGF B null mutant embryo, PDGF A can signal through the a and b receptors. Precise localization of the coproduction patterns of ligands and receptors in the normal and abnormal tissue should help interpret the data obtained in the single-gene targeting approach. Further analysis of double-receptor and double-ligand mutants is underway and should help further understanding the requirement and precise role of PDGF during embryogenesis. Quantitative chimeric mice studies should further help determine its role in development and disease.

## Bibliography

1. R. Ross, J. Glomset, B. Kariya, and L. Harker (1974) *Proc. Natl. Acad. Sci. USA* **71**, 1207–1210.
2. N. Kohler and A. Lipton (1974) *Exp. Cell Res.* **87**, 297–301.
3. B. Westermark and A. Wasteson (1976) *Exp. Cell Res.* **98**, 170–174.
4. C. H. Heldin (1992) *EMBO J.* **11**, 4251–4259.
5. C. H. Heldin, A. Johnsson, S. Wennergren, C. Wernstedt, C. Betsholtz, and B. Westermark (1986) *Nature* **319**, 511–514.
6. Y. Yarden, J. A. Escobedo, W. J. Kuang, T. L. Yang-Feng, T. O. Daniel, P. M. Tremble, E. Y. Chen, M. E. Ando, R. N. Harkins, U. Francke et al. (1986) *Nature* **323**, 226–232.
7. C. E. Hart, M. Bailey, D. A. Curtis, S. Osborn, E. Raines, R. Ross, and J. W. Forstrom (1990) *Biochemistry* **29**, 166–172.
8. A. Hammacher, U. Hellman, A. Johnsson, A. Ostman, K. Gunnarsson, B. Westermark, A. Wasteson, and C. H. Heldin (1988) *J. Biol. Chem.* **263**, 16493–16498.
9. L. Rönstrand, M. P. Beckmann, B. Faulders, A. Ostman, B. Ek, and C. H. Heldin (1987) *J. Biol. Chem.* **262**, 2929–2932.

10. C. E. Hart, J. W. Forstrom, J. D. Kelly, R. A. Seifert, R. A. Smith, R. Ross, M. J. Murray, and D. F. Bowen-Pope (1988) *Science* **240**, 1529–1531.
11. C. H. Heldin, G. Backstrom, A. Ostman, A. Hammacher, L. Ronnstrand, K. Rubin, M. Nister, and B. Westermark (1988) *EMBO J.* **7**, 1387–1393.
12. E. W. Raines, D. F. Bowen-Pope, and R. Ross (1991) "Platelet-Derived Growth Factor", in *Peptide Growth Factors and Their Receptors*, M. B. Sporn and A. B. Roberts, Springer-Verlag, New York, pp. 173–262.
13. C. Betsholtz, A. Johnsson, C. H. Heldin, B. Westermark, P. Lind, M. S. Urdea, R. Eddy, T. B. Shows, K. Philpott, A. L. Mellor et al. (1986) *Nature* **320**, 695–699.
14. D. T. Bonthron, C. C. Morton, S. H. Orkin, and T. Collins (1988) *Proc. Natl. Acad. Sci. USA* **85**, 1492–1496.
15. F. Rorsman, M. Bywater, T. J. Knott, J. Scott, and C. Betsholtz (1988) *Mol. Cell Biol.* **8**, 571–577.
16. W. J. LaRochelle, M. May-Siroff, K. C. Robbins, and S. A. Aaronson (1991) *Genes Dev.* **5**, 1191–1199.
17. A. Ostman, M. Andersson, C. Betsholtz, B. Westermark, and C. H. Heldin (1991) *Cell Regul.* **2**, 503–512.
18. M. P. Beckmann, C. Betsholtz, C. H. Heldin, B. Westermark, E. Di Marco, P. P. Di Fiore, K. C. Robbins, and S. A. Aaronson (1988) *Science* **241**, 1346–1349.
19. R. A. Pollock, and W. D. Richardson (1992) *Growth Factors* **7**, 267–277.
20. A. Ostman, J. Thyberg, B. Westermark, and C. H. Heldin (1992) *J. Cell Biol.* **118**, 509–519.
21. M. Andersson, A. Ostman, B. Westermark, and C. H. Heldin (1994) *J. Biol. Chem.* **269**, 926–930.
22. A. Johnsson, C. H. Heldin, A. Wasteson, B. Westermark, T. F. Deuel, J. S. Huang, P. H. Seeburg, A. Gray, A. Ullrich, and G. Scrace (1984) *EMBO J.* **3**, 921–928.
23. C. E. Heldin (1993) in *Cytokines, Vol. 5: Biology of Platelet-Derived Growth Factor*, B. Westermark et al., eds., Basel and Karger, pp. 1–30.
24. T. Collins, D. Ginsburg, J. M. Boss, S. H. Orkin, and J. S. Pober (1985) *Nature* **316**: p. 748–750.
25. J. L. Kelly, A. Sanchez, G. S. Brown, C. N. Chesterman, and M. J. Sleight (1993) *J. Cell Biol.* **12**, 1153–1163.
26. P. Ataliotis and M. Mercola (1997) *Int. Rev. Cytol.* **172**, 95–127.
27. C. Oefner, A. D'Arcy, F. K. Winkler, B. Eggiman, and M. Hosang (1992) *EMBO J.* **11**, 3921–3926.
28. J. C. Bonner and A. R. Osornio-Vargas (1995) *J. Biol. Chem.* **270**, 16236–16242.
29. K. P. Crookston, D. J. Webb, B. B. Wolf, and S. L. Gonias (1994) *J. Biol. Chem.* **269**, 1533–1540.
30. X. Lin, Z. Wang, L. Gu, and T. F. Deuel (1992) *J. Biol. Chem.* **267**, 25614–25619.
31. L. A. Rafty and L. M. Khachigian (1998) *J. Biol. Chem.* **273**, 5758–5764.
32. M. Sasahara, S. Amano, H. Sato, J. G. Yang, Y. Hayase, M. Kaneko, M. Suzaki, and F. Hazama (1998) *Oncogene* **16**, 1571–1578.
33. A. Kashishian, A. Kazlauskas, and J. A. Cooper (1992) *EMBO J.* **11**, 1373–1382.
34. M. Pech, A. Gazit, P. Arnstein, and S. A. Aaronson, (1989) *Proc. Natl. Acad. Sci. USA* **86**, 2693–2697.
35. A. A. Soyombo and P. E. DiCorleto (1994) *J. Biol. Chem.* **269**, 17734–17740.
36. J. C. Bonner (1994) *Ann. NY Acad. Sci.* **737**, 324–338.
37. R. G. Gronwald, F. J. Grant, B. A. Haldeman, C. E. Hart, P. J. O'Hara, F. S. Hagen, R. Ross, D. F. Bowen-Pope, and M. J. Murray (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3435–3439.
38. L. Claesson-Welsh, A. Eriksson, B. Westermark, and C. H. Heldin (1989) *Proc. Natl. Acad. Sci.*

USA **86**, 4917–4921.

39. L. Claesson-Welsh, A. Hammacher, B. Westermark, C. H. Heldin, and M. Nister (1989) *J. Biol. Chem.* **264**, 1742–1747.
40. J. Tiesman and C. E. Hart (1993) *J. Biol. Chem.* **268**, 9621–9628.
41. W. C. Kenney, M. Haniu, A. C. Herman, T. Arakawa, V. J. Costigan, J. Lary, D. A. Yphantis, and A. R. Thomason (1994) *J. Biol. Chem.* **269**, 12351–12359.
42. L. Claesson-Welsh (1993) in *Cytokines, Vol 5: Biology of Platelet-Derived Growth Factor*, B. Westermark et al., eds., Basel and Karger, pp. 31–43.
43. R. A. Seifert, A. van Koppen, and D. F. Bowen-Pope (1993) *J. Biol. Chem.* **268**, 4473–4480.
44. F. S. Vassbotn, O. K. Havnen, C. H. Heldin, and H. Holmsen (1994) *J. Biol. Chem.* **269**, 13874–13879.
45. T. Inaba, H. Shimano, T. Gotoda, K. Harada, M. Shimada, J. Ohsuga, Y. Watanabe, M. Kawamura, Y. Yazaki, N. Yamada et al. (1993) *J. Biol. Chem.* **268**, 24353–24360.
46. P. A. Eccleston, K. Funa, and C. H. Heldin (1993) *Dev. Biol.* **155**, 459–470.
47. H. E. Abboud, G. Grandaliano, M. Pinzani, T. Knauss, G. F. Pierce, and F. Jaffer, (1994) *J. Cell Physiol.* **158**, 140–150.
48. T. B. Barrett, C. M. Gajdusek, S. M. Schwartz, J. K. McDougall, and E. P. Benditt (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6772–6774.
49. A. E. Ballagi, A. Ishizaki, J. O. Nehlin, and K. Funa (1995) *Biochem. Biophys. Res. Commun.* **210**, 165–173.
50. A. Ishisaki, T. Murayama, A. E. Ballagi, and K. Funa (1997) *Eur. J. Biochem.* **246**, 142–146.
51. E. Shinbrot, X. Liao, and L. T. Williams (1997) *Dev. Dyn.* **208**, 211–219.
52. C. H. Heldin (1997) *FEBS Lett.* **410**, 17–21.
53. C. H. Heldin, A. Wasteson, and B. Westermark (1982) *J. Biol. Chem.* **257**, 4216–4221.
54. C. H. Heldin, A. Ernlund, C. Rorsman, and L. Ronnstrand (1989) *J. Biol. Chem.* **264**, 8905–8912.
55. J. D. Kelly, B. A. Haldeman, F. J. Grant, M. J. Murray, R. A. Seifert, D. F. Bowen-Pope, J. A. Cooper, and A. Kazlauskas (1991) *J. Biol. Chem.* **266**, 8987–8992.
56. A. Kazlauskas and J. A. Cooper (1989) *Cell* **58**, 1121–1133.
57. W. J. Fantl, J. A. Escobedo, G. A. Martin, C. W. Turck, M. del Rosario, F. McCormick, and L. T. Williams (1992) *Cell* **69**, 413–423.
58. L. Ronnstrand, S. Mori, A. K. Arridsson, A. Eriksson, C. Wernstedt, U. Hellman, L. Claesson-Welsh, and C. H. Heldin (1992) *EMBO J.* **11**, 3911–3919.
59. S. Mori, L. Ronnstrand, K. Yokote, A. Engstrom, S. A. Courtneidge, L. Claesson-Welsh, and C. H. Heldin (1993) *EMBO J.* **12**, 2257–2264.
60. Y. Nishizuka (1992) *Science* **258**, 607–614.
61. M. Valius, C. Bazenet, and A. Kazlauskas (1993) *Mol. Cell Biol.* **13**, 133–143.
62. M. Valius and A. Kazlauskas (1993) *Cell* **73**, 321–334.
63. E. Y. Skolnik, B. Margolis, M. Mohammadi, E. Lowenstein, R. Fischer, A. Drepps, A. Ullrich, and J. Schlessinger (1991) *Cell* **65**, 83–90.
64. H. Nakanishi, K. A. Brewer, and J. H. Exton (1993) *J. Biol. Chem.* **268**, 13–16.
65. J. Downward (1998) *Curr. Opin. Cell Biol.* **10**, 262–267.
66. J. A. Escobedo and L. T. Williams (1988) *Nature* **335**, 85–87.
67. M. Joly, A. Kazlauskas, F. S. Fay, and S. Corvera (1994) *Science* **263**, 684–687.
68. V. Kundra, J. A. Escobedo, A. Kazlauskas, H. K. Kim, S. G. Rhee, L. T. Williams, and B. R. Zetter (1994) *Nature* **367**, 474–476.
69. P. Chardin, J. H. Camonis, N. W. Gale, L. van Aelst, J. Schlessinger, M. H. Wigler, and D. Bar-

- Sagi (1993) *Science* **260**, 1338–1343.
70. S. A. Moodie, B. M. Willumsen, M. J. Weber, and A. Wolfman (1993) *Science* **260**, 1658–1661.
71. D. K. Morrison, D. R. Kaplan, U. Rapp, and T. M. Roberts (1988) *Proc. Natl. Acad. Sci. USA* **85**, 8855–8859.
72. D. K. Morrison, D. R. Kaplan, J. A. Escobedo, U. R. Rapp, T. M. Roberts, and L. T. Williams (1989) *Cell* **58**, 649–657.
73. L. S. Mulcahy, M. R. Smith and D. W. Stacey (1985) *Nature* **313**, 241–243.
74. R. M. Kypta, Y. Goldberg, E. T. Ulug, and S. A. Courtneidge (1990) *Cell* **62**, 481–492.
75. K. A. DeMali and A. Kazlauskas (1998) *Mol Cell Biol.* **18**, 2014–2022.
76. J. A. Gelderloos, S. Rosenkranz, C. Bazenet, and A. Kazlauskas (1998) *J. Biol. Chem.* **273**, 5908–5915.
77. R. Nishimura, W. Li, A. Kashishian, A. Mondino, M. Zhou, J. Cooper, and J. Schlessinger (1993) *Mol. Cell Biol.* **13**, 6889–6896.
78. G. G. Choudhury, N. Ghosh-Choudhury, and H. E. Abboud (1998) *J. Clin. Invest.* **101**, 2751–2760.
79. S. Rigacci, M. Bucciantini, R. Marzocchini, and A. Berti (1998) *FEBS Lett.* **432**, 145–149.
80. P. Chiarugi, P. Cirri, F. Marra, G. Raugei, T. Fiaschi, G. Camici, G. Manao, R. G. Romanelli, and G. Ramponi (1998) *J. Biol. Chem.* **273**, 6776–6785.
81. P. Leveen, M. Pekny, S. Gebre-Medhin, B. Swolin, E. Larsson, and C. Betsholtz (1994) *Genes Dev.* **8**, 1875–1887.
82. P. Soriano (1994) *Genes Dev.* **8**, 1888–1896.
83. P. Lindahl, B. R. Johansson, P. Leveen, and C. Betsholtz (1997) *Science* **277**, 242–245.
84. J. R. Crosby, R. A. Seifert, P. Soriano, and D. F. Bowen-Pope (1998) *Nat. Gen.* **18**, 385–388.
85. H. Bostrom, K. Willetts, M. Pekny, P. Leveen, P. Lindahl, H. Hedstrand, M. Pekna, M. Hellstrom, S. Gebre-Medhin, M. Schalling, M. Nilsson, S. Kurland, J. Tornell, J. K. Heath, and C. Betsholtz (1996) *Cell* **85**, 863–873.
86. A. Orr-Urtreger, and P. Lonai (1992) *Development* **115**, 1045–1058.
87. S. L. Palmieri, J. Payne, C. D. Stiles, J. D. Biggers, and M. Mercola (1992) *Mech. Dev.* **39**, 181–191.
88. G. C. Schatteman, K. Morrison-Graham, A. van Koppen, J. A. Weston, and D. F. Bowen-Pope (1992) *Development* **115**, 123–131.

### Suggestions for Further Reading

89. R. P. H. Dirks and H. P. J. Bloemers (1996) Signals controlling the expression of PDGF. *Mol. Biol. Rep.* **22**, 1–24. (Comprehensive review of regulation of PDGF expression; contains exhaustive bibliography.)
90. C. H. Heldin (1992) Structural and functional studies on platelet-derived growth factor. *EMBO J.* **11**, 4251–4259. (Lucid chronologic review of the topic.)
91. C. H. Heldin (1997) Simultaneous induction of stimulatory and inhibitory signals by PDGF, *FEBS Lett.* **410**, 17–21. (Updated review on PDGF receptor signaling.)
92. P. Ataliotis and M. Mercola (1997) Distribution and functions of platelet-derived growth factors and their receptors during embryogenesis, *Int. Rev. Cytol.* **172**, 95–127. (Up-to-date review of the role of PDGF in embryogenesis.)
93. B. Westermark and C. Sorg (1993) in *Cytokines Vol. 5: Biology of Platelet-Derived Growth Factor*, B. Westermark et al., eds., Basel and Karger, 1993, pp. 1–167. (Comprehensive review of the cellular and molecular biology of PDGF.)



## Pleckstrin Homology (PH) Domains

The pleckstrin homology or PH domain is a ~100 amino acid protein module found in numerous different types of human protein and is one of the most abundant non-enzymatic protein modules found in the human genome. The draft human genomes available in early 2001 claim in the one case 252 and the other 193 genes encode proteins encoding PH domains (1, 2). Both draft genomes agree that genes encoding proteins containing SH2 and SH3 domains are significantly less numerous than this. PH domains are also very abundant in the proteins of *Drosophila*, *C. elegans*, yeast and protists, suggesting a very early evolutionary origin and therefore presumably fundamental importance in eukaryotic biology. Proteins containing PH domains are associated with cellular signaling, cellular membrane dynamics and the cytoskeleton, and defects in some of these proteins are associated with a variety of cancers and other serious human diseases. These domains generally mediate the membrane localization of the proteins in which they are found, although individual PH domains appear to do this in quite different ways as a result of their surprisingly variable binding and functional properties. Some PH domains bind primarily to membrane phospholipids, others bind both lipids and membrane associated proteins while still others bind primarily to membrane associated proteins. Recent evidence suggests that the PH domain superfamily is even larger than originally thought, since protein tyrosine binding (PTB) and Enabled/Vasp homology 1 (EVH1) domains, which have virtually no sequence identity with the “classical” PH domain, nevertheless have the same 3D structure. Apparently the PH domain structural prototype has been widely used and heavily modified during evolution and other families of the PH domain superfamily may be discovered as more protein 3D structures are determined.

**Table 1. Brief Overview of Proteins with PH Domains**

---

|   |  |
|---|--|
| “ Classical ” PH domains  | <b>Kinases.</b> Tec/Btk family, bARK family, Rac/Akt/Protein kinase B family.  |
|   | <b>G-protein related.</b> Ras-GAP, Ras-GRF, Trio, ARNO.  |
|   | <b>Cytoskeletal.</b> Syntrophin, AFAP110, Kinesin Kif-1B, b-spectrin isotopes.   |
|   | <b>Signaling adapters.</b> Grb-7, Grb-10, IRS-1, Gab-1, DOS.   |
|   | <b>Membrane associated.</b> Dynamin family, Oxysterol binding protein.   |
| Dbl subfamily of classical PH domains; PH domains adjacent to Dbl homology domains. | <b>Dbl Family.</b> Dbl, Bcr, Vav, SOS, FGD1, Bcr and several others contain PH domain next to a Dbl homology (DH) domain, which activates Rac/Rho family small G proteins by stimulating GTP binding. Several (e.g. Tiam-1, FGD1) also contain classical PH domains. |

|   |   |
|---|---|
| Phospholipase C sub family of classical PH domains. | <b>Phospholipase C family.</b> All have a N-terminal PH domain which has been shown to bind strongly to PIP2 and IP3 in the case of PLC d1. PLC d and g isoforms also contain a second classical PH domain.   |
| PTB domains   | <b>Shc family.</b> Shc isoforms and several more distantly related proteins.<br><br><b>IRS-1 family.</b> IRS-1, IRS-2, 3 and 4, Dok, FRS-2 and others.<br><br><b>Others.</b> Numb, X11 and relatives. Many PTB domains bind to tyrosine phosphorylated peptides in activated receptor tyrosine kinases (Shc, IRS-1 families), others bind non-phosphorylated peptides in other membrane proteins (X11, Numb). |
| Enabled/Vasp Homology 1 (EVH1) domains              | <b>EVH1 family.</b> Enabled ( <i>Drosophila</i> ena and mammalian Mena), EVL, VASP, Wiskott-Aldrich syndrome protein (WASP) family. All bind to Pro rich peptides in various actin associated proteins. Invariably at N-terminus of molecules in which they are found.  |

---

PH domains were first described in a total of less than 20 molecules in 1993 (3, 4). These reports generated considerable interest partly because many of these molecules, such as for example mammalian son of sevenless (mSOS), p120 ras-GAP, ras-GRF and dynamin, had already been widely studied. The reason the PH domain had been overlooked in these and other molecules is that the defining amino acid sequences are a series of poorly conserved peptides interspersed with linkers which are even less well conserved in sequence and length (figure 1). Sequences of this type were difficult to identify with BLAST, FASTA and most other available sequence analysis programs, which worked best with relatively long regions of alignment. However several PH domains somewhat more divergent in primary sequence were quickly described using various kinds of sequence profiling and more specialized search methods. These included PH domains found in the Bruton's tyrosine kinase (Btk) family of src related non receptor tyrosine kinases, the signaling adapter IRS-1 and in the b-adrenergic receptor kinases 1 and 2 (bARK) (5, 6). More divergent but still convincing PH domains were described in the Dbl family of proto-oncogenes (5) and at the N-termini of all phospholipase C (PLC) isoforms (7). The Dbl and PLC families of sequences contain a C-terminal region with key conserved amino acids convincingly similar to that of the other PH domains, but are harder to align in the region corresponding to the N-terminal ~85 amino acid (figure 1). However they do have a similar content of charged, hydrophobic and turn promoting amino acids throughout the sequence alignment, and, for the purposes of this article, the original PH domains and those found in the Dbl and PLC molecules will be referred to as "classical" PH domains.

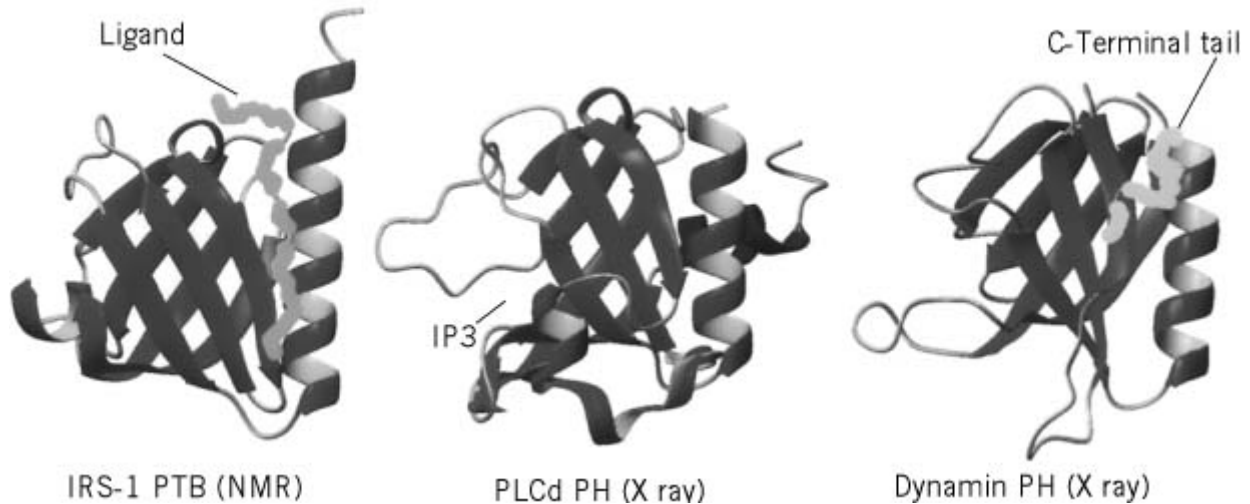
**Figure 1.** Alignment of PH domain structures current determined along with consensus structure. Large hydrophobic and charged amino acids (D, E, K, R, H) are all bold and shaded.

|                   |  |                                   |                            |                 |
|-------------------|--|-----------------------------------|----------------------------|-----------------|
| Pleckstrin        | MEPKRIREGYLVKK   | GSVFNTWKPMWVVLLED                 | GIEFYKKKSDN                |                 |
| $\beta$ -spectrin | MEGFLNRKHEWEAHNKKASSRSWHNVYCVINNQ  |                                   | EMGFYKDAKSA                |                 |
| Dynam in          | MKTSGNDEILVIRKGWLTINNIGIMKG  | GSKEYWVFLTA                       | ENLSWYKDDKEI               |                 |
| $\beta$ ARK       | MHGYSKMGNPFL   | TQWQRRYFYLP                       | NRLEWRGEGEAPI              |                 |
| hSOS              | CNEFIMEGTLTRVGAK   | HERHIFLFD                         | GLMICCKSNHGQPI             |                 |
| PLC- $\delta$ 1   | DPDLQALLKGSQLLKV   | KSSSWRRERFYKLMRSPESQEDCKTIWQESRKV |                            |                 |
| Shc               | GPGVSYLVRYMGCV (51)  | NLKFAGMPITLTVSTS                  | SLNLMAADCK                 |                 |
| IRS-1             | KEVWQVILKPKGLG   | QTKNLIGIYRLCLTSK                  | TISFVKLNSEA                |                 |
| Mena              | MSEQSICQARAAVMVYDDANKKWVPAGGSTGFSRVHIYHHT  |                                   | GNNTFRVVGRKIQD             |                 |
|                   | ----- $\beta\beta\beta\beta$ ----- $\beta\beta\beta\beta$ -----  |                                   |                            |                 |
|                   | $\beta$ 1  | L1                                | $\beta$ 2                  |                 |
|                   |  |                                   | L2                         |                 |
| Pleckstrin        | GMIPLK   | GSTLTSPCQDFGKR                    | MFVFKITTTKQ                | QDFFFQAAFLI     |
| $\beta$ -spectrin | VPVSLK   | EAICEVALDYKKKKH                   | VFKLR LSD                  | GNEYLFQAKDD     |
| Dynam in          | KYMLSVVDN  | LKLRDVEGFMSSKH                    | IFALFNTEQRNVYKDYRQLELACETQ |                 |
| $\beta$ ARK       | SLLTME   | EIQSVEETQIKERK                    | CLLLKIRGGK                 | QFVLQCDSI       |
| hSOS              | SNYRLKEKFFMRVQINDKDDTNEYKH   |                                   | AFEIILK DEN                | SVIFSAKSA       |
| PLC- $\delta$ 1   | QLFSIE   | DIQEVRMGHRTEGLEKFARDIPEDRCFSIVFK  |                            | DQRNTLDLIAPSP   |
| Shc               | QIIANHHMQS   | ISFASGGDPDT                       | AEYVAYVAKDP                | VNQRACHILECPEI  |
| IRS-1             | AVVLQLMN   | IRRCGHCEN                         | FFFIEVGRSA                 | VTGPGEFWMQVDDSI |
| Mena              | HQVVINCAIKGLKYNQATQ  |                                   | TFHQWRDAR                  | QVGLNFGSKI      |
|                   | - $\beta\beta\beta$ ----- $\beta\beta\beta\beta$ ----- $\beta\beta\beta\beta$ ----- $\beta\beta\beta\beta$ ----- |                                   |                            |                 |
|                   | $\beta$ 4  | L4                                | $\beta$ 5                  | L5              |
|                   |  |                                   | $\beta$ 6                  | L6              |
|                   |  |                                   |                            | $\beta$ 7       |

At the time of writing the 3D structures of numerous classical PH domains have been determined, including the N-terminal PH domain of pleckstrin, PH domains of  $\beta$ -spectrin from two species, dynamin, bARK1, mSOS, phospholipase C-d1, Btk and Grp1. Despite the lack of sequence conservation, all of these structures are very similar, composed of an N-terminal sandwich containing of 4 antiparallel  $\beta$ -strands overlaying 3 antiparallel strands which together form a compact wedge shape, the wide end of which is capped by a single C-terminal  $\alpha$ -helix (see figure 2). The regions of  $\beta$ -strand and  $\alpha$ -helix correspond to the better conserved regions of the sequence alignment shown in figure 1, while the loops in the structure correspond to the less conserved regions. The compact PH domain structure has the N and C termini very close together, a feature also seen in SH2, SH3 and many other modular domains.

**Figure 2.** Structure of IRS-1 PTB domain, PLC-d1 PH domain and dynamin PH domains, all in complex with ligands. The basic structure of 7 antiparallel  $\beta$ -strands capped by a C-terminal  $\alpha$ -helix is preserved in each case, although the length, amino acid sequence and conformation of the loops between these segments are not conserved. The peptide ligand for the IRS-1 PH/PTB domain and the binding site for IP3 on PLC-d1 are also indicated, showing the diversity

of type and location of ligand binding sites on these domains. Figure taken from reference [57](#), with permission.



This basic PH domain structure is shared by at least two other classes of modular binding domain. The first class is typified by a region of the adapter molecule Shc and relatives, which was shown to include a novel type of phosphotyrosine binding domain (8). Independently a region of the adapter molecule IRS-1 C-terminal to the classical PH domain was shown to contain a phosphotyrosine binding site, and represents the prototype for a further family of domains found in IRS-2, IRS-3, IRS-4, Dok and FRS-2 (9, 10). Shc and IRS-1 bind phosphotyrosine containing peptides of similar consensus sequence, though binding specificity resides in the amino acids N-terminal to the Tyr residue and so is functionally quite different from the mechanisms by which SH2 domains bind phosphotyrosine containing peptides. These regions became known as PTB (i.e. protein tyrosine binding) domains. Since the amino acid sequence of the Shc and IRS-1 PTB domains show no superficial similarities with the classical PH domains it was surprising that both domains were found to be structurally similar to classical PH domains (11, 12), and both can be considered part of the PH domain superfamily (Figure 2). Similarly sequence similarities defined the Enabled/Vasp Homology 1 (EVH1) domain found in *Drosophila* Enabled, mammalian Enabled (Mena), VASP, EVL, Wiskott-Aldrich syndrome protein (WASP), Homer and others (13). EVH1 domains contain a specific binding site for specific Pro rich sequences found in a variety of actin associated proteins (13). Surprisingly the structure of two different EVH1 domains were recently shown to be similar to that of classical PH and PTB domains (14, 15), so that EVH1 domains also belong to the PH domain superfamily. When the amino acid sequences are aligned structurally, PTB and EVH1 domains show few identities and no absolutely conserved residue compared to PH domains. However all three types of domain are rich in large hydrophobic amino acids such as Trp, Phe and Tyr and also rich in small turn promoting amino acids Gly, Ser, Ala and Pro. The short C-terminal region of all three corresponds to an amphipathic  $\alpha$ -helix which never contains any Pro residues. This C-terminal region contains a large hydrophobic amino acid, almost always a Trp residue in classical PH domains, followed four amino acids C terminal by another large hydrophobic amino acid, frequently a Leu or Ile residue (these two positions are indicated by asterisks in figure 1). This region is also always rich in acidic amino acids, particularly in the linker between the seventh b strand and the C-terminal  $\alpha$ -helix. While it is possible that the structural and sequence similarities between classical PH domains, PTB domains and EVH1 domains are fortuitous, it seems much more likely that the basic PH domain fold has been elaborated greatly during evolution, producing a peculiarly versatile family of binding domains. This hypervariability would also account for the difficulty of identifying member of the PH domain superfamily by sequence relationships alone.

Proteins with classical PH domains frequently contain enzymatic modules such as protein kinase domains (e.g. Btk family Tyr kinases, Protein kinase B (PKB, also known as Akt or Rac kinase)

Ser/Thr kinases, the  $\beta$ -adrenergic receptor Ser/Thr kinases), various phospholipid modifying domains (Phospholipase C $\beta$ ,  $\gamma$  and  $\delta$ , PI-4 kinase) or domains which regulate small G proteins (e.g. Dbl family proto-oncogenes, p120 Ras-GAP, Ras-GRF, mSOS, Sec7 family proteins). Many other PH domain containing proteins are components of the cytoskeleton (e.g. certain  $\beta$ -spectrin isotopes, the syntrophins, actin filament associated protein AFAP-110, the kinesin family protein Kif1b) or belong to the family of so-called “adapter” molecules (e.g. IRS-1, Shc, Grb7, Gab1), which contain binding sites of various types but have no intrinsic enzymatic activity. The types of molecules which contain classical PH domains are generally similar to those which contain SH2 and SH3 domains and in several cases (e.g. Phospholipase C $\gamma$ , Btk family tyrosine kinases, the Vav Dbl family proto-oncogene) all three domain types are found in the same molecule. Proteins with classical PH domains frequently also contain Pro rich peptides which can potentially bind specific SH3, WW or other Pro rich peptide binding domains (e.g. mSOS, dynamin, Grb-7 and Grb-10) or Tyr phosphorylation sites which can bind to specific SH2 domains (e.g. AFAP-110, IRS-1, Shc). A few molecules contain multiple classical PH domains (e.g. pleckstrin, PLC  $\gamma$ , Tiam-1), and one family of proteins, typified by ARAP3, has no less than 5 of these domains (16). PH domains often neighbor enzymatic domains, and in some cases this relationship appears to be obligatory. For example every known Dbl homology (DH) domain is flanked at the C-terminus by a PH domain (17), and this PH domain is vital for the normal cellular function of the DH domain (18, 19). Finally PH domains are found somewhat more frequently than expected by chance at the extreme N or C-termini of proteins. PTB domains are found in molecules such as IRS-1, Shc, Numb and X11 which generally seem to function as membrane associated signaling adapters (10). In the IRS-1 and Dok families of molecules PTB domains are found immediately C-terminal to a classical PH domain, while in most other molecules the PTB domain is at the extreme N-terminus (10). EVH1 domains are found in non enzymatic molecules which associate with membranes and the actin cytoskeleton and which might be referred to as cytoskeletal adapters (13). EVH1 domains are invariably at the extreme N-terminus of molecules in which they are found, as are many classical PH domains. PTB and EVH1 domains, like some PH domains, are therefore found in adapter molecules, in support of the view that all three domain types are evolutionarily related.

The exact function of PH domains has been controversial, and much of the controversy no doubt occurred because it was expected that these domains, like SH2 and SH3 domains, would have one predominant class of ligand. However current evidence suggests that PH domains have a surprising variety of binding partners. Some bind directly to membrane phospholipids, some to other generally membrane associated proteins and some to both membranes and proteins. The affinity and specificity of these interactions are quite variable. Many PH domains bind phosphatidylinositol 4,5 bisphosphate (PIP<sub>2</sub>) and in some cases inositol 1,4,5, trisphosphate (IP<sub>3</sub>) (20-24). While some of these affinities are high, in the nM range, most are of relatively low affinity and specificity but may nevertheless be biological significant. Other PH domains such as those in PKB and Btk family kinases show preferential and high affinity and specificity binding for phosphatidylinositol 3,4,5 trisphosphate (PIP<sub>3</sub>), phosphatidylinositol-3,4-bisphosphate (PI-3,4-P<sub>2</sub>) or the PIP<sub>3</sub> head group inositol-1,3,4,5-tetrakisphosphate (IP<sub>4</sub>) (24-26). PKB and Btk kinases become membrane localized by their PH domains following activation of PI-3 kinases which generate PIP<sub>3</sub> and PI-3,4-P<sub>2</sub>. The binding sites for PIP<sub>2</sub> on pleckstrin,  $\beta$ -spectrin and PLC-d1 PH domains involve primarily the L1 and L5 linkers (23) (27, 28). However the phosphoinositide binding sites the dynamin and Btk PH domains appear to involve primarily the L1 and L3 linkers and would therefore be on different faces of the PH domain (24, 29). A consensus from these studies is that PH domains which are rich in basic residues in the N-terminal  $\beta$ -strand region generally bind lipids with higher affinity and specificity than those with fewer basic amino acids. The recent inclusion of PTB domains into the PH domain superfamily is supported by the finding that certain of these, namely those of Shc and mDab1, bind phospholipids through interactions mediated by the basic N terminal  $\beta$ -strand regions (30-32). EVH1 domains however have not so far been shown to bind directly to lipids and generally lack a concentration of basic residues in this N-terminal segment.

Protein binding by PH domains is if anything even more variable. The well studied bARK1 PH domain binds to the  $\beta$   $\gamma$  subunits of trimeric G-proteins (G $\beta$   $\gamma$ ), albeit with higher affinity if the

peptide sequence immediately following the PH domain is included, and PH domains from several other proteins have been reported to bind Gb g *in vitro*, though the binding was rather variable and in most cases not as robust as the Gb g-bARK interaction (33). Much evidence suggests that some of these may also bind *in vivo* (34-38). PH domains of Rac/Akt/PKB and Btk protein kinases have also been reported to bind to protein kinase C isotypes (39). More recently several PH domains have been reported to bind actin directly, though strong actin binding, like Gb g binding, was not seen with every PH domain tested (40). In addition the PH domain from a member of the Btk family was shown to bind the FERM domain of FAK, the focal adhesion kinase (41), and the PH domain of IRS-1 was shown to interact with PHIP, a protein required *in vivo* for efficient phosphorylation of IRS-1 by the activated insulin receptor complex (42). The theme that emerges is that PH domain binding is generally to a protein which is associated with the membrane, so that these protein-protein interactions may be functionally similar to direct lipid binding as described above. PTB domains certainly fit into this pattern since many of them bind to tyrosine phosphorylated peptide sequences found in activated receptor tyrosine kinases. However, in a similar fashion to classical PH domains, PTB domains show distinctly and somewhat surprisingly different types of binding. For example the Dab PTB domain, paradoxically, binds to peptides similar in sequence to those of other PTB domains but only if they are not tyrosine phosphorylated (31). The PTB domains of Numb and X11 also recognize uniquely different peptide consensus sequences but only if they are not tyrosine phosphorylated (10). Recent studies of Numb suggest that this PTB domain has at least three not obviously related *in vivo* binding sites for other proteins (43). Finally the various EVH1 domains bind to Pro rich sequences in their target molecules, although the exact peptide binding site on the domain is slightly different in the two currently known structures, the EVH1 domains of Mena and Homer (14, 15). We can only conclude that the various different types of classical PH domain, PTB and EVH1 domains bind a bewildering variety of unrelated protein ligands, and may bind them using different parts of the 3D structure. Nevertheless it is noteworthy that, as noted for classical PH domains, almost all of the described PTB and EVH1 domain ligands are associated with cellular membranes either as integral lipid components or as membrane associated proteins.

Probably the best studied classical PH domain is that of bARK1, which binds both PIP2 and Gb g. PIP2 binding is of low affinity and is mediated by the N-terminal b-strand region of the PH domain, while Gb g binding is mediated by an interaction mapped to the C-terminal a-helix and some following sequence (44). This region is now known to correspond to an a-helical extension from the PH domain (45). The peptide following the PH domain contains two serine residues, Ser 670 and Ser 685, which may become phosphorylated *in vivo* and which regulate Gb g binding. Ser 670 can be phosphorylated by MAP kinase and inhibits Gb g binding (46), while Ser 685 can be phosphorylated by Protein Kinase A and enhances Gb g binding (47). Clearly there is considerable potential for regulatory complexity in the regulation of PH domain binding interactions if bARK1 is typical, and similar regulation of membrane binding may also occur. The presence of both a phospholipid and a protein binding site on the same small protein module naturally suggests allosteric regulation of binding. The bARK1 PH domain responds to both PIP2 and Gb g, and both ligands function in a coordinated manner apparently to regulate the membrane binding of the bARK PH domain (48, 49). To date the significance of phosphorylation on this cooperativity has not been examined. Studies of dynamin, an abundant protein important in endocytosis which contains a PH domain adjacent to a GTPase domain shows that PIP2 activates the GTPase while Gb g inhibits it and that the affinity of dynamin for Gb g is increased in the presence of PIP2 (50). These effects are presumably mediated by the dynamin PH domain and, if so, provide another glimpse of the potentially complex regulation of binding interactions by this domain. More recently Gb g was shown to act through the PH domain of PLC-b2 and PLC-b1 to regulate these enzymes *in vivo* (37, 38), and Gb g was shown to bind to the PH domain of Btk and also to regulate Btk enzymatic activity but by binding to a site independent of the PH domain (51).

Making functional sense of this variety of binding sites has been controversial. However in almost every case the domain binding involves direct or indirect membrane association, and it seems reasonable to conclude that PH, PTB and EVH1 domains all function in this way, to target molecules, by a variety of different mechanisms, to the membrane. It may be that the multiple

binding properties attributed to several PH domains are typical, and similar multiple binding partners remain to be discovered for other PH, PTB and EVH1 domains. It has been pointed out that the compact 7 b-strand, 1 a-helix structural template presenting a set of hypervariable peptide loops is an arrangement somewhat similar to that seen in the antibody variable region (52). These authors proposed that the PH domain may actually function like an antibody variable region, defining a binding activity which can be quite different in different molecules.

Proteins with PH domains are involved in many serious human diseases. One striking example is the Dbl family of proteins, many of which were discovered in gene transfer experiments using tumor DNA (17). Oncogenic forms of these proteins frequently have N-terminal deletions but always retain the Dbl homology (DH) domain and the adjacent PH domain. Since the DH domain is a guanine nucleotide exchange factor for Rac/Rho/Cdc42 family small G proteins, defects in these molecules presumably result in aberrant regulation of these important small G proteins. Point mutations in the non-receptor tyrosine kinase Btk are causative of Bruton's X-linked agammaglobulinaemia, a serious human immunodeficiency in which B cells fail to develop, resulting in a severely compromised immune system (53). Some of these mutations directly affect the binding of the Btk PH domain to PIP3 and IP4 (24, 54). The ubiquitous PKB Ser/Thr kinases may be virally transduced and rendered oncogenic by a combination of overexpression and fusion to viral GAG protein, which renders the PKB component constitutively membrane localized (55, 56). The EVH1 domain containing protein WASP is mutated in patients with Wiscott-Aldrich syndrome (WAS), and the majority of these mutations are found in the EVH1 domain. WAS patients have severe defects in platelets and T cells due to aberrant organization of the actin cytoskeleton. There are many other examples of human diseases associated with particular defects in PH domain superfamily molecules, underlining the importance of these domains in the normal regulation of cellular function.

In summary, PH domains are a very large family of protein modules based on a conserved structural module of 7 antiparallel b strands followed by a single a-helix. This family is so variable as to be virtually unrecognizable by primary amino acid sequence alone, so that definitive recognition of new classes of PH domain will only come from further structural work. This family of modules is also surprisingly variable in terms of binding properties, and the best current functional explanation for this is that PH domains are the cellular equivalent of the antibody molecules, a structurally conserved unit with hypervariable loops which can be evolutionarily modified to produce binding sites for virtually any ligand. The central theme seems to be the regulation of membrane binding, though this appears to be achieved in a bewildering variety of different ways. In conclusion, the fundamental role of PH domains is probably to direct and regulate the formation of protein-membrane and specific protein-protein interactions, allowing the formation of a wide variety of specific and membrane associated signaling complexes.

## References

- “Pleckstrin Homology Domains” in , Vol. 3, pp. 1875–1879, by Dr. Gerry Shaw, McKnight Brain Institute, Professor of Neuroscience, Anatomy and Cell Biology, Gainesville, Florida, 32610, Tel (352) 294 0037 (lab), (352) 294 0038 (office), Fax (352) 392 8347, Email shaw@mbi.ufl.edu;
- “Pleckstrin Homology (PH) Domains” in (online), posting date: January 15, 2002, by Dr. Gerry Shaw, McKnight Brain Institute, Professor of Neuroscience, Anatomy and Cell Biology, Gainesville, Florida, 32610, Tel (352) 294 0037 (lab), (352) 294 0038 (office), Fax (352) 392 8347, Email shaw@mbi.ufl.edu.
1. Lander, E. S. et al. (2001) *Nature* **409**, 860–921.
  2. Venter, J. C. et al. (2001) *Science* **291**, 1304–1351.
  3. Mayer, B. J., Ren, R., Clark, K. L. & Baltimore, D. (1993) *Cell* **73**, 629–30.
  4. Haslam, R. J., Koide, H. B. & Hemmings, B. A. (1993) *Nature* **363**, 309–10.
  5. Musacchio, A., Gibson, T., Rice, P., Thompson, J. & Saraste, M. (1993) *Trends Biochem Sci* **18**, 343–8.

6. Shaw, G. (1993) *Biochem Biophys Res Commun* **195**, 1145–51.
7. Parker, P. J., Hemmings, B. A. & Gierschik, P. (1994) *Trends in Biochemical Sciences* **19**, 54–55.
8. Bork, P. & Margolis, B. (1995) *Cell* **80**, 693–4.
9. O'Neill, T. J., Craparo, A. & Gustafson, T. A. (1994) *Mol Cell Biol* **14**, 6433–42.
10. Margolis, B., Borg, J. P., Straight, S. & Meyer, D. (1999) *Kidney Int* **56**, 1230–7.
11. Zhou, M. M., Ravichandran, K. S., Olejniczak, E. F., Petros, A. M., Meadows, R. P., Sattler, M., Harlan, J. E., Wade, W. S., Burakoff, S. J. & Fesik, S. W. (1995) *Nature* **378**, 584–92.
12. Eck, M. J., Dhe-Paganon, S., Trub, T., Nolte, R. T. & Shoelson, S. E. (1996) *Cell* **85**, 695–705.
13. Reinhard, M., Jarchau, T. & Walter, U. (2001) *Trends Biochem. Sci.* **26**, 243–9.
14. Prehoda, K. E., Lee, D. J. & Lim, W. A. (1999) *Cell* **97**, 471–80.
15. Beneken, J., Tu, J. C., Xiao, B., Nuriya, M., Yuan, J. P., Worley, P. F. & Leahy, D. J. (2000) *Neuron* **26**, 143–54.
16. Krugmann, S., Anderson, K. E., Ridley, S. H., Risso, N., McGregor, A., Coadwell, J., Davidson, K., Eguinoa, A., Ellson, C. D., Lipp, P., Manifava, M., Ktistakis, N., Painter, G., Thuring, J. W., Cooper, M. A., Lim, Z. Y., Holmes, A. B., Dove, S. K., Michell, R. H., Grewal, A., Nazarian, A., Erdjument-Bromage, H., Tempst, P., Stephens, L. R. & Hawkins, P. T. (2002) *Mol Cell* **9**, 95–108.
17. Cerione, R. A. & Zheng, Y. (1996) *Curr Opin Cell Biol* **8**, 216–22.
18. Whitehead, I., Kirk, H. & Kay, R. (1995) *Oncogene* **10**, 713–21.
19. Chardin, P., Paris, S., Antonny, B., Robineau, S., Beraud-Dufour, S., Jackson, C. L. & Chabre, M. (1996) *Nature* **384**, 481–4.
20. Harlan, J. E., Hajduk, P. J., Yoon, H. S. & Fesik, S. W. (1994) *Nature* **371**, 168–70.
21. Cifuentes, M. E., Delaney, T. & Rebecchi, M. J. (1994) *J. Biol. Chem.* **269**, 1945–8.
22. Yagisawa, H., Hirata, M., Kanematsu, T., Watanabe, Y., Ozaki, S., Sakuma, K., Tanaka, H., Yabuta, N., Kamata, H., Hirata, H. & Nojima, H. (1994) *J. Biol. Chem.* **269**, 20179–88.
23. Harlan, J. E., Yoon, H. S., Hajduk, P. J. & Fesik, S. W. (1995) *Biochemistry* **34**, 9859–64.
24. Salim, K., Bottomley, M. J., Querfurth, E., Zvelebil, M. J., Gout, I., Scaife, R., Margolis, R. L., Gigg, R., Smith, C. I., Driscoll, P. C., Waterfield, M. D. & Panayotou, G. (1996) *Embo J* **15**, 6241–50.
25. James, S. R., Downes, C. P., Gigg, R., Grove, S. J., Holmes, A. B. & Alessi, D. R. (1996) *Biochem J* **315**, 709–13.
26. Fukuda, M. & Mikoshiba, K. (1996) *J Biol Chem* **271**, 18838–42.
27. Hyvönen, M., Macias, M. J., Nilges, M., Oschkinat, H., Saraste, M. & Wilnanns, M. (1995) *EMBO J.* **14**, 4676–85.
28. Ferguson, K. M., Lemmon, M. A., Schlessinger, J. & Sigler, P. B. (1995) *Cell* **83**, 1037–46.
29. Fushman, D., Cahill, S. & Cowburn, D. (1997) *J Mol Biol* **266**, 173–94.
30. Rameh, L. E., Arvidsson, A., Carraway, K. L., 3rd, Couvillon, A. D., Rathbun, G., Crompton, A., VanRenterghem, B., Czech, M. P., Ravichandran, K. S., Burakoff, S. J., Wang, D. S., Chen, C. S. & Cantley, L. C. (1997) *J Biol Chem* **272**, 22059–66.
31. Howell, B. W., Lanier, L. M., Frank, R., Gertler, F. B. & Cooper, J. A. (1999) *Mol Cell Biol* **19**, 5179–88.
32. Ravichandran, K. S., Zhou, M. M., Pratt, J. C., Harlan, J. E., Walk, S. F., Fesik, S. W. & Burakoff, S. J. (1997) *Mol Cell Biol* **17**, 5540–9.
33. Touhara, K., Inglese, J., Pitcher, J. A., Shaw, G. & Lefkowitz, R. J. (1994) *J Biol Chem* **269**, 10217–20.
34. Tsukada, S., Simon, M. I., Witte, O. N. & Katz, A. (1994) *Proc Natl Acad Sci U S A* **91**, 11256–60.



35. Luttrell, L. M., Hawes, B. E., Touhara, K., van Biesen, T., Koch, W. J. & Lefkowitz, R. J. (1995) *J Biol Chem* **270**, 12984–9.
36. Langhans-Rajasekaran, S. A., Wan, Y. & Huang, X. Y. (1995) *Proc Natl Acad Sci U S A* **92**, 8601–5.
37. Wang, T., Dowal, L., El-Maghrabi, M. R., Rebecchi, M. & Scarlata, S. (2000) *J Biol Chem* **275**, 7466–9.
38. Razzini, G., Brancaccio, A., Lemmon, M. A., Guarnieri, S. & Falasca, M. (2000) *J Biol Chem* **275**, 14873–81.
39. Konishi, H., Kuroda, S., Tanaka, M., Matsuzaki, H., Ono, Y., Kameyama, K., Haga, T. & Kikkawa, U. (1995) *Biochem Biophys Res Commun* **216**, 526–34.
40. Yao, L., Janmey, P., Frigeri, L. G., Han, W., Fujita, J., Kawakami, Y., Apgar, J. R. & Kawakami, T. (1999) *J Biol Chem* **274**, 19752–61.
41. Chen, R., Kim, O., Li, M., Xiong, X., Guan, J. L., Kung, H. J., Chen, H., Shimizu, Y. & Qiu, Y. (2001) *Nat Cell Biol* **3**, 439–44.
42. Farhang-Fallah, J., Yin, X., Trentin, G., Cheng, A. M. & Rozakis-Adcock, M. (2000) *J Biol Chem* **275**, 40492–7.
43. Zwahlen, C., Li, S. C., Kay, L. E., Pawson, T. & Forman-Kay, J. D. (2000) *Embo J* **19**, 1505–15.
44. Koch, W. J., Inglese, J., Stone, W. C. & Lefkowitz, R. J. (1993) *J. Biol. Chem.* **11**, 8256–8260.
45. Fushman, D., Najmabadi-Haske, T., Cahill, S., Zheng, J., LeVine, H., 3rd & Cowburn, D. (1998) *J Biol Chem* **273**, 2835–43.
46. Pitcher, J. A., Tesmer, J. J., Freeman, J. L., Capel, W. D., Stone, W. C. & Lefkowitz, R. J. (1999) *J Biol Chem* **274**, 34531–4.
47. Cong, M., Perry, S. J., Lin, F. T., Fraser, I. D., Hu, L. A., Chen, W., Pitcher, J. A., Scott, J. D. & Lefkowitz, R. J. (2001) *J Biol Chem* **276**, 15192–9.
48. Pitcher, J. A., Touhara, K., Payne, E. S. & Lefkowitz, R. J. (1995) *J Biol Chem* **270**, 11707–10.
49. DebBurman, S. K., Ptasienski, J., Boetticher, E., Lomasney, J. W., Benovic, J. L. & Hosey, M. M. (1995) *J Biol Chem* **270**, 5742–7.
50. Lin, H. C. & Gilman, A. G. (1996) *J Biol Chem* **271**, 27979–82.
51. Lowry, W. E. & Huang, X. Y. (2002) *J Biol Chem* **277**, 1488–92.
52. Cohen, G. B., Ren, R. & Baltimore, D. (1995) *Cell* **80**, 237–48.
53. Mattsson, P. T., Vihinen, M. & Edvard Smith, C. I. (1996) *BioEssays* **18**, 825–34.
54. Fukuda, M., Kojima, T., Kabayama, H. & Mikoshiba, K. (1996) *J Biol Chem* **271**, 30303–6.
55. Cheng, J. Q., Godwin, A. K., Bellacosa, A., Taguchi, T., Franke, T. F., Hamilton, T. C., Tsichlis, P. N. & Testa, J. R. (1992) *Proc. Natl. Acad. Sci. USA.* **89**, 9267–9271.
56. Bellacosa, A., Franke, T. F., Gonzalez-Portal, E., Datta, K., Taguchi, T., Gardner, J., Cheng, J. Q., Testa, J. R. & Tsichlis, P. N. (1993) *Oncogene* **8**, 745–754.
57. Cowburn, D. (1996) *Structure* **4**, 1005–1008.

### **Suggestions for Further Reading**

58. Forman-Kay, J.D., and T. Pawson. 1999. Diversity in protein recognition by PTB domains. *Curr Opin Struct Biol.* **9**: 690–5.
59. Hurley, J.H., and S. Misra. 2000. Signaling and subcellular targeting by membrane-binding domains. *Annu Rev Biophys Biomol Struct.* **29**: 49–79.
60. Maffucci, T., and M. Falasca. 2001. Specificity in pleckstrin homology (PH) domain membrane targeting: a role for a phosphoinositide-protein co-operative mechanism. *FEBS Lett.* **506**: 173–9.
61. Reinhard, M., T. Jarchau, and U. Walter. 2001. Actin-based motility: stop and go with

## Pleiotropy

A single [mutation](#) can cause more than one change in phenotype, in which case it is called pleiotropic. The mutation is also be said to have pleiotropic effects. An example is a mutation in a gene encoding a global regulatory protein.

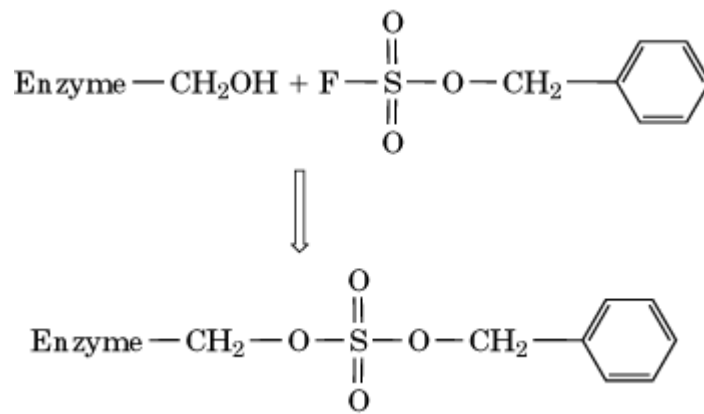
## Ploidy

Ploidy is the number of complete sets of [chromosomes](#) (or **autosomes**) in a cell. [Haploid](#) (or [euploid](#)) cells have a ploidy of one, **diploid** cells a ploidy of two, **triploid** cells a ploidy of three, etc. [Polyploid](#) cells have more than the normal chromosome content. Deviations from complete sets of chromosomes are described by [aneuploidy](#).

## PMSF (Phenylmethylsulfonyl Fluoride)

A number of *sulfonyl halides* have been found to be helpful in the study of [proteins](#) and other biological molecules. One of the most useful is phenylmethylsulfonyl fluoride (PMSF). It is a crystalline solid of low toxicity that is a potent, irreversible inhibitor of **chymotrypsin** and other [serine proteinases](#) (1). Aqueous solutions of PMSF undergo slow hydrolysis, but it is reasonably stable in anhydrous alcohol (2). When used to minimize protein degradation during isolation procedures, it is added to buffer solutions prior to use by 100-fold dilution from a stock 50 mM solution in isopropanol (3). The product of reaction with chymotrypsin is the phenylmethylsulfonyl ester of the active site serine (Figure 1). Although this ester is quite stable, it is possible to regenerate the active enzyme by incubation under acidic conditions (pH 2) for several hours (4). It is also possible, under alkaline conditions, to convert the product to *anhydrochymotrypsin* in which the [active site](#) serine residue becomes a dehydroalanine residue (4).

**Figure 1.** Interaction of PMSF (phenylmethylsulfonyl fluoride) with the active-site serine hydroxyl group of a serine proteinase.



## Bibliography

1. D. E. Fahrney and A. M. Gold (1963) *J. Am. Chem. Soc.* **85**, 997 ff.
2. G. T. James (1978) *Anal. Biochem.* (1978) **86**, 574–579.
3. G. H. Goodwin (1989) In *Protein Purification Methods: A Practical Approach* (E. L. V. Harris and S. Angel, eds.), IRL Press, Oxford, U.K., pp. 97–101.
4. A. M. Gold (1967) *Methods Enzymol.* **11**, 706–711.

## Point Accepted Mutation

The point accepted mutation (PAM) is a measure of the rate at which point [mutations](#) that substitute one [amino acid](#) residue for another have been incorporated in a **gene** lineage during [evolution](#). A PAM is sometimes used as a unit of **evolutionary divergence**. In practice, one PAM is equivalent to a unit of evolutionary divergence in which 1% of the amino acid residues of the protein product of a gene have been changed. On the basis of this PAM model of evolution, Dayhoff et al. ([1](#)) invented the amino acid substitution matrices in which the substitution scores between a pair of different amino acids are proportional to the natural log of the ratio of the target frequencies to the background frequencies ([1](#)). The background frequencies are the frequencies of each possible substitution, which would be determined simply by the overall frequencies of the different amino acids if the changes were purely random. The target frequencies are the observed substitution frequencies, which are biased toward substitutions of similar amino acids that do not seriously undermine the protein function. To estimate the target frequencies, pairs of very closely related sequences were used to collect mutation frequencies corresponding to 1 PAM, followed by extrapolation of the data to a distance of 250 PAMs. The PAM250 matrix is given in Figure [1](#) (see top of next page). This matrix is often used for both [aligning sequences](#) and conducting [homology](#) searches.

**Figure 1.** The PAM250 score matrix. The amino acid residues are indicated by their one-letter abbreviations. The values given are the natural logarithms of the frequencies with which the pairs of amino acid residues have been changed during evolution in proteins of known sequence. This value reflects both the probability of the codons for the two residues being interconverted by mutation and the effect that such a mutation will have on the function of the protein. Negative numbers indicate changes that are disfavored; positive numbers correspond to changes that occur

more frequently, and they are shaded.

|   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|
| A | 2  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| R | -2 | 6  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| N | 0  | 0  | 2  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| D | 0  | -1 | 2  | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| C | -2 | -4 | -4 | -5 | 12 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| Q | 0  | 1  | 1  | 2  | -5 | 4  |    |    |    |    |    |    |    |    |    |    |    |    |    |   |
| E | 0  | -1 | 1  | 3  | -5 | 2  | 4  |    |    |    |    |    |    |    |    |    |    |    |    |   |
| G | 1  | -3 | 0  | 1  | -3 | -1 | 0  | 5  |    |    |    |    |    |    |    |    |    |    |    |   |
| H | -1 | 2  | 2  | 1  | -3 | 3  | 1  | -2 | 6  |    |    |    |    |    |    |    |    |    |    |   |
| I | -1 | -2 | -2 | -2 | -2 | -2 | -2 | -3 | -2 | 5  |    |    |    |    |    |    |    |    |    |   |
| L | -2 | -3 | -3 | -4 | -6 | -2 | -3 | -4 | -2 | 2  | 6  |    |    |    |    |    |    |    |    |   |
| K | -1 | 3  | 1  | 0  | -5 | 1  | 0  | -2 | 0  | -2 | -3 | 5  |    |    |    |    |    |    |    |   |
| M | -1 | 0  | -2 | -3 | -5 | -1 | -2 | -3 | -2 | 2  | 4  | 0  | 6  |    |    |    |    |    |    |   |
| F | -3 | -4 | -3 | -6 | -4 | -5 | -5 | -5 | -2 | 1  | 2  | -5 | 0  | 9  |    |    |    |    |    |   |
| P | 1  | 0  | 0  | -1 | -3 | 0  | -1 | 0  | 0  | -2 | -3 | -1 | -2 | -5 | 6  |    |    |    |    |   |
| O | 1  | 0  | 1  | 0  | 0  | -1 | 0  | 1  | -1 | -1 | -3 | 0  | -2 | -3 | 1  | 2  |    |    |    |   |
| T | 1  | -1 | 0  | 0  | -2 | -1 | 0  | 0  | -1 | 0  | -2 | 0  | -1 | -3 | 0  | 1  | 3  |    |    |   |
| W | -6 | 2  | -4 | -7 | -8 | -5 | -7 | -7 | -3 | -5 | -2 | -3 | -4 | 0  | -6 | -2 | -5 | 17 |    |   |
| Y | -3 | -4 | -2 | -4 | 0  | -4 | -4 | -5 | 0  | -1 | -1 | -4 | -2 | 7  | -5 | -3 | -3 | 0  | 10 |   |
| V | 0  | -2 | -2 | -2 | -2 | -2 | -2 | -1 | -2 | 4  | 2  | -2 | 2  | -1 | -1 | -1 | 0  | -6 | -2 | 4 |
|   | A  | R  | N  | D  | C  | Q  | E  | G  | H  | I  | L  | K  | M  | F  | P  | S  | O  | T  | Y  | V |

## Bibliography

1. M. O. Dayhoff, R. M. Schwartz, and B. C. Orcutt (1978) in *Atlas of Protein Sequence, and Structure*, M. O. Dayhoff, ed., National Biomedical Research Foundation, Washington, DC.

## Poison Sequence

Poison sequence is a general term used to describe **nucleic acid** or [amino acid](#) sequences whose presence in a molecule prevents that molecule from carrying out a normal function, such as replication of **DNA**, **intron** splicing, or translocation of a [membrane protein](#). Most often, however,

the term “poison sequence” describes a sequence of DNA that cannot be **cloned**. In certain cases, poison sequences are used as experimental tools.

**Plasmid** pBR322 poison sequences have been identified that prevent efficient replication of recombinant SV40 virus-pBR322 **shuttle vectors** constructed to move DNA between bacterial and mammalian cells for **mutagenic** studies. Deletions near the pBR322 **origin of replication** increase vector replication efficiency, indicating that sequences near the pBR322 origin inhibit (poison) SV40 replication in mammalian cells (1).

Some poison sequences make it extremely difficult to clone complete sequences of DNA, cause the deletion or rearrangement of inserts, and frequently result in the loss of the poison sequence (2, 3). The presence of a poison sequence in the *gag* gene of the mouse mammary tumor **virus** (MMTV), for example, makes it difficult to study the pathogenicity of the virus, because for many years the full sequence could not be cloned in an intact state for further studies (3). Other examples of poison sequences are those that prevent **splicing** of premessenger RNA *in vitro* (4) and the highly polar cytoplasmic **domain** of *Escherichia coli* **leader peptidase**, which acts as a “translocation poison” sequence that prevents membrane insertion of hydrophobic domains of the protein (5).

Various uses are made of poison sequences. For example, poison sequences are used to create positive-selection cloning **vectors**. A toxic sequence or gene is cloned downstream from multiple cloning sites, and the poison sequence prevents bacterial growth unless it is inactivated by inserting of a DNA fragment into one of the cloning sites (6). In a different type of experiment, a poison sequence that interfered with gene expression when placed into an intron was used to study illegitimate **recombination** of mammalian cell **chromosomes** by selecting for **mutations** that eliminate or alter the poison sequence (7).

#### Bibliography

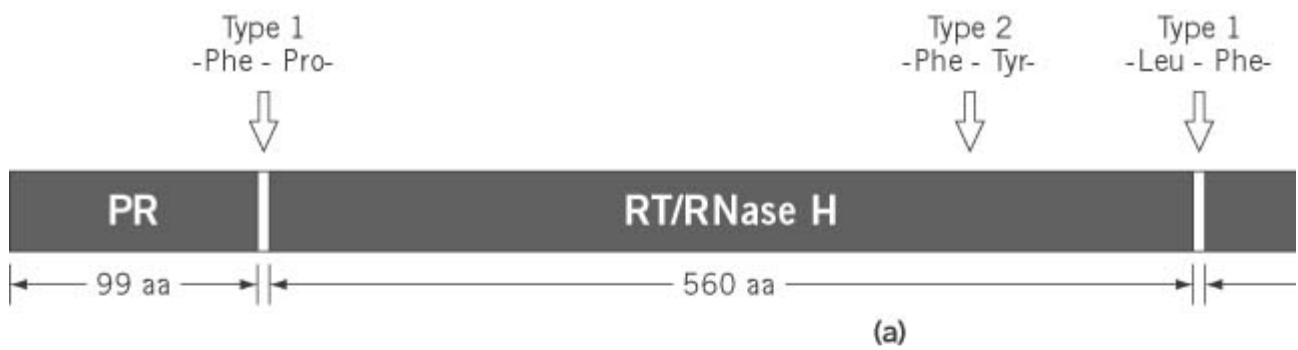
1. M. Seidman (1989) *Mutation Res.* **220**, 55–60.
2. E. M. Karawya (1988) *Ital. J. Biochem.* **37**, 365–379.
3. D. W. Morris, H. D. Bradshaw, Jr., H. T. Billy, R. J. Munn, and R. D. Cardiff (1989) *J. Virol.* **63**, 148–158.
4. P. J. Furdon and R. Kole (1988) *Mol. Cell. Biol.* **8**, 860–866.
5. G. von Heijne, W. Wickner, and R. E. Dalbey (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3363–3366.
6. P. Bernard, P. Gabant, E.M. Bahassi, and M. Couturier (1994) *Gene* **148**, 71–74.
7. T. Porter, S. L. Pennington, G. M. Adair, R. S. Nairn, and J. H. Wilson (1990) *Nucleic Acids Res.* **18**, 5173–5180.

#### **Pol Gene**

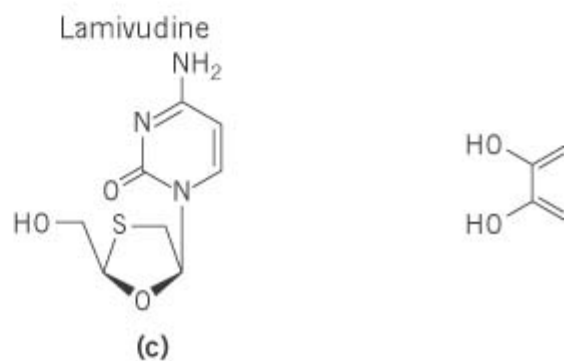
The *pol* gene of **retroviruses** contains coding sequences that specify retroviral **enzymes** involved in (i) **replication** of the RNA **genome**, **reverse transcriptase** (RT) (ii) integration of double-stranded proviral DNA, **integrase** (IN), and (iii) **proteolytic** maturation of precursor **polyproteins** that accompany virion budding, **protease** (PR) (Fig. 1). In many lentiviruses, a deoxyuridine triphosphatase (DU or dUTPase) **domain** has also been identified between the RT and IN **open reading frames**, and it may play a role in viral expression in different cell types. Over the last decade, these enzymes have been the subject of intense scrutiny due to the emergence of **HIV** and AIDS and the continuing search for therapeutic agents. Except for human foamy virus (and perhaps

other spumaviruses), where an independent **mRNA** for *pol* enzymes has been established, the *pol* gene is translated as a larger *gag-pol* [polyprotein](#) that constitutes ~5% of the total virion *gag* content. Mechanisms for differential synthesis of virion structural and enzymatic proteins (translational **suppression** and ribosomal [frameshifting](#)) are outlined in [Retroviruses](#). For unknown reasons, ribosomal frameshifting occurs after the PR coding region of avian viruses, with the consequence that this enzyme is synthesized at about a 20-fold molar excess over RT and IN. The possibility that PR plays an accessory structural role in avian retroviruses has been considered but is as yet unproven. Because the avian PR is considerably less active than its counterpart in other retroviruses, overexpression may be a means of correcting this defect. General features of each of these enzymes are presented in the following sections.

**Figure 1.** The *pol* open reading frame of HIV-1. The sizes of the three *pol*-derived enzymes, PR, RT/RNase H, and IN, & recognition sites for PR-mediated maturation of the *pol* polyprotein (**a**). Ribbon diagrams represent the structure of each fragment containing the catalytic core of IN has been crystallized. In the absence of supportive evidence, this is represented as a model. Antivirally active that are in clinical use (Sequinavir, Lamivudine) or under investigation (chicoric acid).



(b)



### 0.1. Protease

Viral PR controls events whereby *gag* and *pol* proteins are released from each other during budding. Loss of PR activity through [site-directed mutagenesis](#) does not impair virion assembly, but the resulting particles are noninfectious. In the strictest sense, PR might be considered a [zymogen](#) because it must be released autocatalytically from the precursor polyprotein to fulfill these events. Because the active form of all retroviral proteases is a dimer, therefore it follows that maturation events must be preceded by association of larger *gag-pol* precursors into dimers. The notion that precursors can dimerize is not unreasonable because most, if not all, of the components they encode exist in either dimeric or multimeric states. However, within this dimer, the cleavage sites that dictate PR excision lie on either side of the catalytic site, suggesting that its removal from the precursor occurs via *trans* cleavage (2).

Before their three-dimensional structures, were elucidated amino acid [sequence analysis](#) indicated

that the retroviral enzyme is distantly related to a class of cellular [carboxyl proteinases](#) (renin, [pepsin](#), [cathepsin D](#), and the erythrocyte membrane aspartyl proteinase), a common feature of which is the sequence -Asp-Thr/Ser-Gly- at the [active site](#). However, the two types of [enzymes](#) have major structural differences. Although the bilobed nature of the cellular enzyme provides a mechanism for duplicating the active site -Asp-Thr/Ser-Gly- motif, retroviral PR is a dimer of identical subunits, each of which contributes one copy of the active site. Subunits of the PR monomer (which vary in size from 99 to 126 residues), intertwine at their N- and C-termini to form a four-stranded, antiparallel  **$\beta$ -sheet**. A binding cleft above the active site accommodates the peptide substrates of precursor polyproteins. Detailed crystallographic analysis of the PR from HIV and ALSV has indicated that four residues (denoted P4 to P1 and P1' to P4' (see [Proteinase Inhibitors](#), [Proteins](#)) on either side of the scissile bond (P1-P1') are accommodated in the binding cleft of the homodimer at their respective enzyme subsites (S4-S1 and S1'-S4'). The entrance to the binding cleft is covered by "flaps" that undergo considerable reorganization after substrate binding.

N-terminal amino acid **sequencing** of mature *gag* and *pol* proteins, recently superseded by highly accurate determination of their molecular weights via **mass spectroscopy**, has been used to determine the PR cleavage sites in precursor polyproteins of many retroviruses. Most cleavage sites are **hydrophobic** and can be grouped into two classes based on the nature of the residue downstream of the scissile bond (P1'). Type I sites prefer a [proline](#) residue at P1', whereas Type II sites prefer [alanine](#), [leucine](#), or [valine](#) at this position. Although there is significant sequence heterogeneity when cleavage sites are compared, the various retroviral PR are not readily interchangeable with respect to heterologous peptide substrates, indicating that sequence and structure combine to provide specificity. For HIV, an exhaustive substrate analysis has provided an invaluable framework for developing potent and highly selective peptide-based inhibitors containing a nonhydrolyzable P1-P1' bond. The structure of Sequinavir, developed by Hoffmann-La Roche, is illustrated in Fig. [1](#). However, the success of this strategy is dampened by the rapid rate at which resistance is achieved by the virus, caused in part by the promiscuity of the retroviral [polymerase](#).

## 0.2. Reverse Transcriptase

Converting the single-stranded viral RNA genome into double-stranded proviral DNA, uninterrupted by ribonucleotides, is the responsibility of a single viral enzyme, reverse transcriptase (RT). These events require a combination of RNA- and DNA-dependent DNA synthesis, together with a degradative activity to remove RNA from the RNA-DNA replicative intermediate ([ribonuclease H](#) or RNase H). The two enzymatic activities reside in different domains of the enzyme, and in certain circumstances they can be physically separated into enzymatically active polypeptides. Structural similarity between the retroviral RNase H domain and its bacterial counterpart, coupled with the observation that RT encoded by the mitochondrial plasmid of **Neurospora** lacks this domain (which is supplied in the form of a mitochondrial RNase H), suggests that the retroviral polymerase evolved from an ancestral **DNA polymerase** by sequestering an RNase H domain from its host. In doing so, this provides a convenient means by which the polymerizing and degradative activities are coordinated during replication, avoiding the dependence on a host-coded function.

Despite common catalytic mechanisms, retroviral RTs differ significantly in [quaternary structure](#). The isolated form of MLV RT is an ~80-kDa monomer, although recent evidence suggests that DNA synthesis is catalyzed by a homodimer. In contrast, HIV and related lentiviral RT are heterodimers of 66 and 51 kDa subunits (p66 and p51, respectively). These enzymes are derived from the same gene but differ inasmuch as partial PR-mediated maturation of the p66/p66 homodimer removes a copy of the C-terminal RNase H domain from one subunit. The heterodimeric avian enzyme represents an unusual scenario whereby the smaller  $\alpha$  subunit contains the DNA polymerase and RNase H domains and the larger  $\beta$  subunit contains these and the entire integrase (IN) polypeptide. Whether it is advantageous to retain a second copy of the RNase H domain (e.g., MLV and ALSV) or a copy of IN (ALSV) remains to be established.

A major advance in understanding the retroviral polymerase has been elucidation of the three-



dimensional structure of HIV-1 RT cocrystals containing either the nonnucleoside inhibitor Nevirapine (3) or double-stranded DNA (4). Based on the anatomical resemblance to a right hand, subdomains of the larger or p66 subunit have been designated “fingers,” “palm,” “thumb,” and “connection,” the last of which connects the DNA polymerase and RNase H catalytic centers (Fig. 1). These four subdomains are duplicated in the smaller p51 subunit. However, alternative folding has the consequence that the p51 connective subdomain masks its DNA polymerase catalytic center, that is, catalytic activities in the parental heterodimer, are mediated by p66. The necessity for asymmetrical organization of identical subunits is not immediately clear, although the smaller HIV polypeptide has been implicated in both accommodating the tRNA replication primer and resistance to nonnucleoside drugs. Subdomains of both subunits contribute to the nucleic acid binding cleft, which is capable of accommodating ~18 bp of duplex. The p66 thumb participates in binding and translocating the template–primer duplex, whereas its N-terminal finger subdomain interacts with the single-stranded template. Catalytic events occur within the palm subdomain, mediated in part by an aspartic acid triad (Asp110, Asp185, and Asp186).

During replication, RT must accommodate three structurally diverse nucleic acid substrates. **Minus-strand** synthesis initiates from an A-form RNA duplex, thereafter continuing along a non-A, non-B RNA-DNA hybrid. In contrast, B-form duplex DNA occupies the nucleic acid binding cleft during **plus-strand** synthesis. Surprisingly, the cocrystal of RT and duplex DNA indicates that the nucleic acid is more A-form at the polymerase catalytic center, B-form at the RNase H catalytic center, and bent over an angle of 45° between the two (see [DNA Structure](#)). Whether all nucleic acid duplexes adopt a similar geometry within RT and other related nucleic acid-polymerizing enzymes remains to be established. In addition to the HIV-1 enzyme, a fragment of MLV RT that contains the finger and palm subdomains has been crystallized and shows the same general features.

A second and equally important function of the retroviral polymerase is mediated by its C-terminal RNase H domain, which degrades the RNA moiety of the RNA–DNA replication intermediate relatively nonspecifically (under appropriate conditions, this domain hydrolyzes duplex RNA, although the biological significance of this has not been established). At later stages in replication, a more precise role for RNase H involves generating the appropriate polypurine tract (PPT) 3'-terminus for initiating of plus-strand synthesis and removing minus- and plus-strand primers from nascent DNA. The use of highly purified recombinant enzymes and model heteropolymeric template–primer combinations has allowed a detailed analysis of RNase H function. Crystallographic data indicates that accommodating an RNA–DNA hybrid in the DNA polymerase catalytic center positions the RNA strand ~18 nucleotides upstream for endonucleolytic cleavage, a notion that has been validated experimentally. However, RNase H-mediated events do not cease here but continue along the template in a 3'- to 5'- direction to within eight nucleotides of the primer terminus. The two hydrolytic functions, designated “polymerization-dependent” and “polymerization-independent,” are clearly necessary for replication steps where nascent DNA is transferred within or between strands of the RNA genome. In such model systems, the accepted notion is that RNase H activity is directed by RT, whose polymerase catalytic center is appropriately positioned over the 3'-terminus of a recessed DNA primer. During minus-strand synthesis, however, short RNA fragments remain hybridized to a considerably longer DNA, making a DNA 3'-terminus unavailable to “direct” RNase H activity. Under these conditions, it is now clear that the replicating enzyme positions itself over an RNA 5'-terminus to catalyze the same sequence of hydrolytic events.

Structural similarity between the RNase H of bacterial and retroviral origin has prompted several biochemical and biophysical studies of the isolated retroviral domain. Although it folds correctly, this HIV-1 domain lacks activity, attributable in part to the observation that the RNase domain is not released intact during maturation of the p66/p51 heterodimer, but lacks important N-terminal structural elements. Secondly, an  $\alpha$ -helix of the bacterial enzyme implicated in nucleic acid binding ( $\alpha$ -helix C) has been eliminated from the HIV-1 domain. This most likely reflects structural modifications that accompanied acquisition of RNase H activity by an ancestral RT because the majority of nucleic acid binding would be assumed by the polymerase domain. Unlike HIV RT, the

RNase H domain of the MLV enzyme retains  $\alpha$ -helix C and can be removed from the polymerase domain in an active form. When administered in combination with anti-PR drugs, nucleoside-based RT inhibitors such as Lamivudine (Fig. 1) are currently proving highly effective as antiviral agents.

### 0.3. Integrase

Insertion of the provirus into the host genome is accomplished by a two-step mechanism, for which a dimeric or higher order form of the C-terminal *pol* component, IN, suffices (1). The first event (in the cytoplasm) involves site-specific cleavage at the 3'-end of each viral LTR 5' to the conserved -C-A- dinucleotide to expose new 3'-OH ends. Subsequently, a concerted cleavage-ligation reaction, involving these OH ends as attacking **nucleophiles** in a transesterification reaction, generates staggered cuts in the host DNA, allowing covalent linkage of the viral 3'-ends. The 5'-ends of viral RNA are flanked by short gaps, which are repaired by cellular enzymes. *In vitro* studies with model duplex DNA substrates and recombinant proteins from human and avian retroviruses have been instrumental in uncovering mechanisms of integration. Structurally, IN can be divided into three distinct regions: (i) an N-terminal domain of ~50 amino acids; (ii) an ~150-residue catalytic core; and (iii) a C-terminal domain of some 100 residues.

The N-terminal domain is characterized by a constellation of [histidine](#) and [cysteine](#) residues commonly associated with the DNA binding “**zinc finger**” motifs of transcriptional regulators, although the spacing between residues of the IN motif is slightly larger. Despite a high degree of conservation between this domain of IN proteins, its removal does not have a deleterious effect on its processing and joining activities. The isolated N-terminal domain also has no affinity for DNA, although this may not reflect its function in the parental enzyme. However, the N- and C-terminal domains may associate to establish higher order forms of IN necessary for activity. The central portion, or catalytic core, of IN is characterized by a constellation of invariant amino acids -Asp-Asp-(35 residues)-Glu- which is a hallmark of the [superfamily](#) of **nucleases** and **phosphotransferases** (one member of which is the RNase H domain of retroviral RT). Simultaneous joining of both LTR ends 4- to 6-bp apart to target DNA implies that IN should function minimally as a dimer. This has in fact been observed in crystals of the catalytic core. The least conserved portion of retroviral IN is the C-terminal domain, with which a relatively nonspecific DNA binding function is associated. As determined for the catalytic core, the isolated C-terminal polypeptide exists as a dimer, which has spawned the theory that duplex DNA may be accommodated within a dimer of this domain. Supporting crystallographic data of intact IN that contains its DNA substrate(s) are currently unavailable.

The absence of a cellular counterpart makes HIV IN an especially attractive target for therapeutic intervention. Unfortunately, although a substantial number of inhibitors have been reported by using *in vitro* biochemical assays, these have little effect *in vivo*. A major barrier is their inability to penetrate target cells. However, chicoric acid (Fig. 1) and derivatives do not suffer from these problems, and they may herald a first generation of potent antiviral agents.

### 0.4. Deoxyuridine Triphosphatase

In lentiviruses, such as FIV, EIAV, CAEV and visna maedi, and the D-type retroviruses MPMV and SRV a fourth open reading frame that encodes a protein of ~20 kDa has been identified between the RT and IN genes. Originally designated a pseudoprotease from amino acid sequence analysis, more recent analysis of virion-associated and recombinant proteins has unambiguously demonstrated that this protein contains deoxyuridine triphosphatase (dUTPase) activity. Cellular dUTPase plays an important role in **nucleic acid** biosynthesis, where it hydrolyzes dUTP to dUMP and pyrophosphate and thereby provides a substrate for [thymidylate synthase](#). In doing so, this provides a mechanism for minimizing **uracil** incorporation into DNA, which would be potentially **mutagenic**. Biochemical characterization indicates that the active species of the cellular and retroviral dUTPase is a trimer.

Why this activity is acquired by a limited number of retroviruses is not immediately clear. One clue may come from studies with **herpes simplex virus** (HSV), a DNA virus whose genome also encodes a dUTPase. Although eliminating dUTPase activity does not result in loss of viral infectivity, such

mutants have reduced neurovirulence and neuroinvasiveness *in vivo*. Because dUTPase levels are low in terminally differentiated or nondividing cells, it has been hypothesized that certain retroviruses retain this activity to permit high-level replication in macrophages. In support of this, equine macrophages support growth of wild-type EIAV but not a variant genetically manipulated to remove dUTPase activity, although this mutant virus replicates efficiently in nonmacrophage cell lines.

## Bibliography

1. A. Wlodawer and J. W. Erickson (1993) *Ann. Rev. Biochem.* **62**, 543–585.
2. L. A. Kohlstaedt, J. Wang, J. M. Friedman, P. A. Rice, and T. A. Steitz (1992) *Science* **256**, 1783–1790.
3. A. Jacobo-Molina, J. Ding, R. G. Nanni, A. D. Clark, X. Lu, C. Tantillo, R. L. Williams, G. Kamer, A. L. Ferris, P. Clark, A. Hizi, S. H. Hughes, and E. Arnold (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6320–6324.
4. M. D. Andrade and A. M. Skalka (1996) *J. Biol. Chem.* **271**, 19633–19636.

## Polar

Molecules or groups of atoms bearing significant positive or negative charges, even if only partial, are said to be *polar*. These molecules or groups often contain O, S, or N atoms that can act as [hydrogen bond](#) acceptors or donors, allowing them to interact with other groups possessing charges of opposite sign or with solvent [water](#), which is itself highly polar because it bears partial positive and negative charges. Molecules or groups that bear full electrostatic charges, such as the  $\text{-NH}_3^+$  and  $\text{-COO}^-$  groups of amino acids in neutral solution, are especially polar and difficult to remove from water. In a strict sense, polarity is best measured by determining the *dipole moment*. In a more practical sense, [hydrophilicity](#) or [hydrophilic](#) character provides a useful index of a molecule's polarity, by measuring experimentally the difficulty of removing it from solvent water.

## Polar Plasma

Although, formerly, it was hypothesized that **germ cells** share the same origin as **somatic cells** and that these develop to primordial germ cells during embryogenesis, it is now known that this is not the case. During the first stages of embryonic development, an event takes place that determines which of the cells will become somatic and which will become germ cells. It is amazing that the decision as to which of the cells are to transfer the genetic information to the next generation in later life is made so shortly after [fertilization](#). Some somatic cells are certainly involved in forming the genitalia, the urogenital tract, and the Wolfian and the Mullerian ducts, and also in developing the embryonic gonads. However, these allow a later invasion by the primordial germ cells that have a different origin and are located at the posterior end of the fertilized [egg](#), as has been demonstrated in insect cells. Thus, it is of interest to determine what factors make cells become germ cells.

In a number of species, there are factors that are necessary for germ line formation, and these are

accumulated in a certain region of the egg, the so-called [polar plasma](#) (or pole plasm). In *Drosophila* embryos, germ plasma is located in the posterior pole region. In transplantation experiments, it was shown that the pole plasma contains the determinants to form a germ line. This finding was supported by experiments with eggs in which the polar plasma was inactivated by ultraviolet (UV) irradiation: Development of an abdomen and of germ cells failed (1). Thus, it was assumed that cytoplasmic compounds determine the fate of the future germ line cells. Given that there is virtually no [gene transcription](#) at the very early stages of cell division from the genome of the embryo, it is likely that maternal **genes** direct early cell [differentiation](#) and the development of germ cells. These gene products were identified by selecting females with homozygous maternal-effect mutants. The genes and gene products found thus far include: tudor, *oskar*, *staufen*, *vasa*, *nanos*, *valois*, *cappuccino*, and *spire*. Deletion of these genes causes failure to produce germ cells and/or abdomen because either a direct effect on a subsequent **gene expression** is missing or the localization of gene products to the posterior end is lost on egg activation.

The protein *nanos* is essential for germ line formation in *Drosophila*. Pole cells lacking activity of this protein do not become functional germ cells (2), as they fail to migrate into the gonad. *Nanos* is therefore a protein necessary for transport of germ cells.

The proteins *oskar*, *staufen*, *vasa*, and *tudor* are necessary for the assembly of the pole plasm (3, 4). *Staufen* is required for localization of *oskar* RNA and translation. The localizations of *vasa* and *tudor* are dependent on *oskar* protein, a [transcription factor](#), and are required for accumulation of *oskar* protein at the posterior pole. Mislocation of *oskar* protein to the anterior pole results in the formation of germ cells at the anterior pole, thereby indicating the importance of correct transport and localization of transcription factors. *Tudor* is concentrated in the polar granules and is involved in two distinct determinative processes in embryogenesis. Mutations of the *tudor* gene affect the later segmentation of the abdomen and the determination of the primordial germ cells (5). The *mago nashi* locus encodes an essential product for germ plasma assembly. The gene product is necessary for the specification of the anterior-posterior axis and the dorsoventral coordinates (6). *Vasa* is also necessary for the formation of polar granules and germ cells (7). Further, mitochondrial large ribosomal RNA was identified in the germ plasma to act as a cytoplasmic factor that induces pole cell formation; it is tightly associated with polar granules, the distinctive organelles of the germ plasma. Also, Hsp83, a member of the Hsp90 **molecular chaperone** family, seems to be a component of the posterior polar plasma.

The polar plasma and the germ line cells may be of interest for attempts to clone individuals. Although fully differentiated cells are probably not able to return to their totipotent stage, primordial cells appear to retain omnipotence for at least 30 days, as has been shown in cattle.

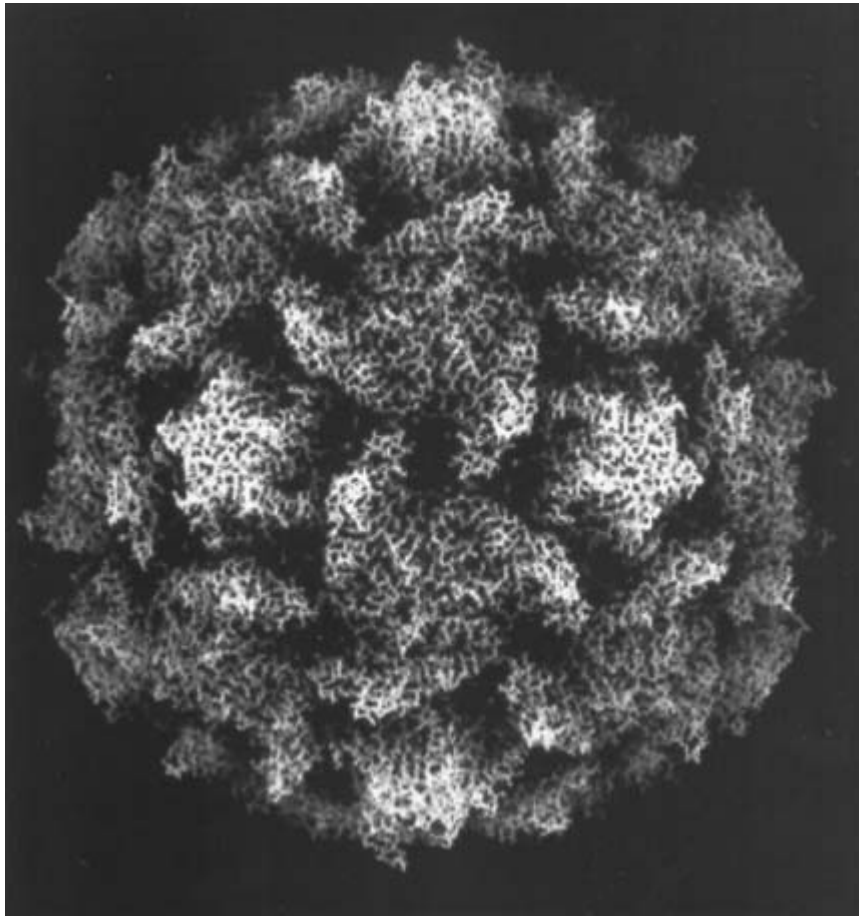
## Bibliography

1. R. M. Savage, and M. V. Danilchik (1993) Dev. Biol. **157**, 371–382.
2. S. Kobayashi, M. Yamada, M. Asaoka, and T. Kitamura (1996) Nature **380**, 708–711.
3. A. Ephrussi, and R. Lehmann (1992) Nature **358**, 387–392.
4. A. Ephrussi, L. K. Dickinson, and R. Lehmann (1991) Cell **66**, 37–50.
5. A. Bardsley, K. McDonald, and R. E. Boswell (1993) Development **119**, 207–219.
6. P. A. Newmark, and R. E. Boswell (1994) Development **120**, 1303–1313.
7. B. Hay, L. Y. Jan, and Y. N. Jan (1988) Cell **55**, 577–587.

## Poliovirus

Poliovirus, the causative agent of poliomyelitis, is a human enterovirus that belongs to the family *Picornaviridae* (see [Rhinovirus](#)) and contains a single-stranded RNA [genome](#). The precise three-dimensional structure of the icosahedral virion particle was elucidated by [X-ray crystallography](#) studies (Fig. 1). These structural investigations indicated that the virion particle is composed of 60 copies each of the viral capsid proteins VP1, VP2, VP3, and VP4 and indicated the structures of the neutralizing **epitopes** and viral attachment sites, called the “canyon” (see [Rhinovirus](#)), to the poliovirus receptor (PVR) on the surface of permissive cells.

**Figure 1.** Three-dimensional structure of poliovirus. Courtesy of Drs. James M. Hogle and Robert Grant, Harvard Medical School.



The genome of poliovirus has [messenger RNA](#) polarity, is approximately 7500 nucleotides in length, is **polyadenylated**, and is linked covalently at its 5' end to a small protein called VPg. The RNA itself is infectious; cells transfected with the RNA produce progeny virions that are infectious. The [complementary DNA](#) corresponding to the entire genome is also infectious. This infectious cDNA clone blazed a trail in the field of the molecular genetics of poliovirus.

The nucleotide sequence of viral mRNA in infected cells is essentially identical to that of the virion RNA, but the VPg is missing in the mRNA. Poliovirus mRNA harbors a long 5' untranslated region (UTR) of approximately 740 nucleotides in length that is important for synthesis of viral RNA and proteins. A possible cloverleaf-like structure formed by the 5'-proximal end of the RNA (approximately 90 nucleotides) is a probable *cis*-element that regulates the synthesis of the plus-

strand RNA. The following sequence of 400 to 500 bases makes up the internal ribosomal entry site (IRES), which directs the viral [translation](#) initiation step in a manner independent of whether or not there is a **5'-cap**. The IRES is assumed to carry a number of secondary structures, and multiple host cellular factors are required for its functions. Polypyrimidine-tract binding protein (PTB), La protein, and poly-riboC binding protein (PCBP)-2 are thought to be involved in the translation initiation process, although these cellular proteins are usually distributed mainly in the [nucleus](#).

Attenuated strains were developed in 1950s and have been used efficiently as oral live vaccines. The vaccine strains can grow well in the human alimentary tract and elicit neutralizing antibodies against the virus, but they have a very poor replicating capacity in the central nervous system (CNS). Thus, the virulent and attenuated strains differ in their CNS tissue responses. Molecular genetic analysis of the CNS response to poliovirus revealed that the phenotype was determined primarily by the IRES region. Thus, it is possible that the CNS response of the virus is the reflection of CNS-specific expression of the IRES function.

The mRNA has a unique long open reading frame that encodes the large viral precursor [polyprotein](#). The large polyprotein is cotranslationally cleaved by viral [proteinases](#) 2A and 3C to form the viral capsid proteins VP1-4 and the viral replication proteins 2A, 2B, 2C, 3A, 3B, 3C, and 3D. Among the latter, 2A and 3C are proteinases, 3B is VPg, and 3D is RNA-dependent **RNA polymerase**. VPg is importantly involved in the initiation of the viral RNA synthesis. In addition to the final processed products, intermediate proteolytic fragments such as 3AB (3A + 3B) and 3CD (3C + 3D) are also considered to play important roles in the initiation step of viral RNA synthesis.

Poliovirus infects only primates, and other animal species are generally not susceptible to the virus. This species-specificity of poliovirus has long been considered to be determined by PVR protein. PVR is known to be a member of [immunoglobulin](#) superfamily. Transgenic mice carrying the human *PVR* gene were actually susceptible to poliovirus. A line of such transgenic mice is now approved to be an animal model for poliomyelitis, in addition to a monkey model.

#### Suggestions for Further Reading

C. Mirzayan and E. Wimmer (1994) "Polioviruses: Molecular Biology". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1119–1132.

R. R. Rueckert (1996) "Picornaviridae: The Viruses and Their Replication". In *Fields Virology*, 3rd ed., (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 609–654.

E. Wimmer, C. U. T. Hellen, and X. Cao (1993) Genetics of poliovirus. *Annu. Rev. Genet.* **27**, 353–436.

## Poly A

The 3' ends of almost all [messenger RNA](#) from **eukaryotes** consist of a homopolymer of adenosine residues that is added post-transcriptionally by [polyadenylate polymerase](#) in the [polyadenylation](#) reaction (1, 2). The length of this [poly A](#) tail is variable and ranges between approximately 20 and 250 residues, depending on the mRNA. The poly A tail has also been implicated in RNA stability and in the regulation of [translation](#) (3, 4). It is important to note that the poly A tail is not a static modification, but a dynamic one, which can be shortened or lengthened. This can be exemplified by analyzing the length of the poly A tail of a specific mRNA during different times of [development](#). This may be important for developmental regulation of **gene expression** (2, 3).

## Bibliography

1. E. Wahle (1995) *Biochim. Biophys. Acta*, **1261**, 183–194.
2. D. F. Colgan and J. L. Manley (1997) *Genes Devel.* **11**, 2755–2766.
3. L. E. Hake and J. D. Richter (1997) *Biochim. Biophys. Acta-Rev. Cancer* **1332**, M31–M38.
4. M. Muckenthaler, N. Gunkel, R. Striepecke, and M. W. Hentze (1997) *RNA* **3**, 983–995.

## Suggestion for Further Reading

5. W. Keller and L. Minvielle-Sebastia (1997) A comparison of mammalian and yeast pre-mRNA 3'-end processing, *Curr. Opin. Cell Biol.* **9**, 329–336.

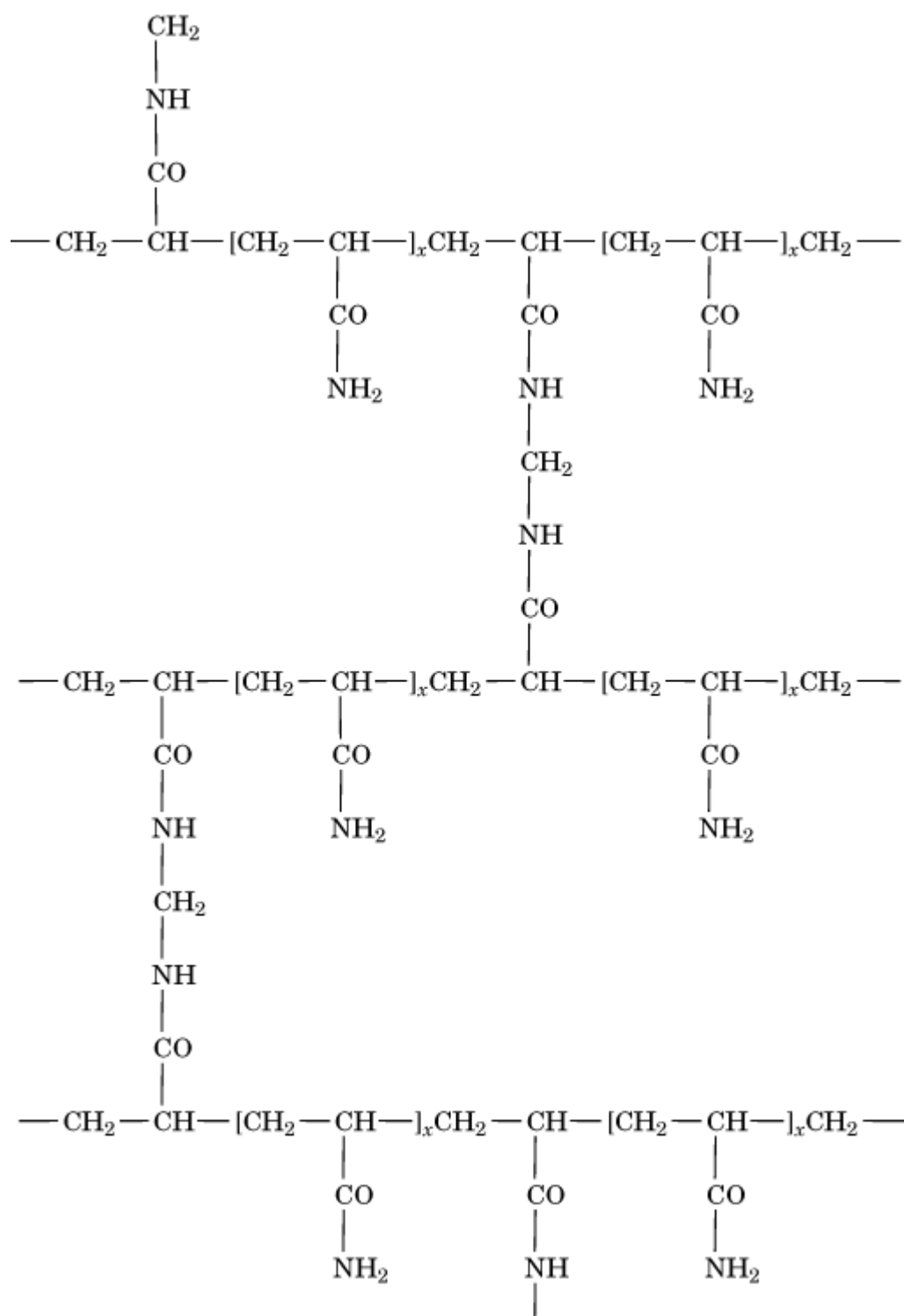
## Polyacrylamide

For over 30 years, polyacrylamide has been, and still is, the building material for [gel electrophoresis](#) with the widest range of pore sizes. It is suitable for resolving by [electrophoresis](#), on the basis of size, shape, and conformation, species with molecular weights ranging from a few hundred daltons to those of large particles the size of a **microsome**. Gel concentrations in the range of 3% to 50% (w/v) acrylamide are practical. Polyacrylamide gel electrophoresis is commonly abbreviated as PAGE .

Polyacrylamide (Fig. [1](#)) is formed by polymerization of the monomer acrylamide by a complex polymerization reaction that is catalyzed by free radicals. Thus, to obtain a reproducible polymer, along with a high degree of conversion of monomers to polymer, a considerable number of parameters need to be controlled, especially the following:

1. the purity of the monomers, especially the absence of acrylamide polymers and of acrylic acid;
2. purity and concentrations of the free-radical donors that initiate the polymerization;
3. levels of inhibitors (eg, free-radical captors, especially oxygen);
4. temperature;
5. pH;
6. divalent metal concentrations;
7. the reaction time, which roughly determines the degree of polymerization.

**Figure 1.** Elements of a Bis-crosslinked polyacrylamide gel network. The linear polyacrylamide chains are depicted horizontally, the Bis crosslinks vertically.



A large number of polymerization initiator species are available, but only three are commonly used: persulfate, tetramethylethylenediamine (TEMED), and riboflavin activated by light. Persulfate and TEMED together are effective polymerization catalysts in the alkaline pH range. Acidic gels are more effectively polymerized in the presence of all three catalysts combined in optimal ratios. Even in the alkaline range, photopolymerization of riboflavin in the presence of persulfate will produce gelation within the desired time, about 10 min for analytical-scale gels, coupled with precise control of the time at which polymerization starts by the start of illumination.

Polyacrylamide is normally crosslinked with N,N'-bisacrylamide (Bis) (Fig. 1). The crosslinking



retards progressively the rate of migration of macromolecules as the fraction of Bis is increased from 0% to 5% of the total acrylamide monomers; beyond that, the retardation diminishes with further increased crosslinking, and it becomes negligible when acrylamide and Bis are copolymerized in equal amounts (1). Polyacrylamide that has not been crosslinked is able, however, to separate some DNA conformations that crosslinked polyacrylamide cannot. More generally, the pore architecture of crosslinked vinyl polymers like polyacrylamide can be varied greatly through choice of a large number of crosslinking agents (2) or initiator conditions (3, 4). Lateral aggregation of gel fibers brought about by inclusion in the polymerization mixture of hydrophilic polymers, such as polyethylene glycol, produces large pores (5). Furthermore, the retardation of molecules in a polyacrylamide gel can be modulated by the average length of the polyacrylamide chains. The chain length and the effective pore size are inversely related, and the average chain length can be controlled through the polymerization conditions. The effective “pore size” of the network is defined by the ratio of the diameter of one covalent polyacrylamide particle to the average distance in three dimensions between polymer chains. If this ratio is too small, or if it exceeds unity, molecular sieving ceases or, at least, diminishes. The average chain length is directly related to the molecular sieving effect of the network. When the average chain length is decreased through excessive initiator concentrations or by the presence of inhibitors during polymerization, the effectiveness of sieving decreases. Also important for molecular sieving is the distribution of crosslinks. Crosslinking agents that enter into the growing polymer chain more slowly than the monomer tend to give rise to an uneven distribution of crosslinks (“nodules”), which increases the effective pore size. Polyacrylamide chains lacking crosslinks form gels above 10% monomer concentration, and those with crosslinks form gels above 3% monomer concentration. Uncrosslinked polyacrylamide excels in network homogeneity but is less effective than crosslinked polyacrylamide in the molecular sieving of proteins.

Freshly prepared polyacrylamide is an oxidative medium containing, in part, peroxide-type uncharged oxidizing agents that cannot be removed by pre-electrophoresis alone. Electrophoresis of charged reducing agents, such as thioglycolate, loaded prior to, or with, the protein sample, is often needed to overcome the oxidative effect of polyacrylamide on proteins.

Polyacrylamide has good adherence to glass walls, and it can even be prepared covalently bound to them, rendering it compatible with vertically oriented apparatus. In horizontally oriented apparatus, polyacrylamide can be firmly bonded to thin plastic sheets (GelBond-PAG); this allows for the preparation of thin gels, compatible with the high field strengths that enhance resolution.

#### Bibliography

1. A. Chrambach, T. M. Jovin, P. J. Svendsen, and D. Rodbard (1976) In *Methods of Protein Separation*, Vol. 2 (N. Catsimpoalas, ed.), Plenum Press, New York, pp. 49–67.
2. C. Gelfi and P. G. Righetti (1981) *Electrophoresis* 2, 213–227.
3. T. Lyubimova and P. G. Righetti (1993) *Electrophoresis* 14, 191–201.
4. M. Chiari, C. Micheletti, P. G. Righetti, and G. Poli (1992) *J. Chromatogr.* 598, 287–297.
5. P. G. Righetti, S. Caglio, M. Saracchi, and S. Quaroni (1992) *Electrophoresis* 13, 587–594.

#### Suggestion for Further Reading

6. A. Chrambach and D. Rodbard (1972) Polymerization of polyacrylamide gels: Efficiency and reproducibility as a function of catalyst concentration. *Separation Sci.* 7, 663–723.

#### Polyadenylate Polymerase

Addition of the [poly A](#) tail to a [messenger RNA](#) precursor requires the action of polyadenylate polymerase, a single-subunit [enzyme](#) that catalyzes addition of the adenosine homopolymer to the 3' end of the pre-mRNA (1). The sequence of the catalytic **domain** of the enzyme is conserved within **eukaryotes**. Under normal conditions, the enzyme is highly specific, polyadenylating only those pre-mRNA that have gone through the cleavage step of cleaving the precursor near its 3' end, to generate the site for poly A addition, and still have the cleavage and polyadenylation stimulatory factors bound to the sequence signaling polyadenylation, AAUAAA (2, 3) (see [Polyadenylation](#)). Under appropriate conditions *in vitro*, however, it can polyadenylate any RNA that has a free 3' hydroxyl group; this is commonly referred to as the nonspecific reaction.

#### Bibliography

1. E. Wahle (1995) *Biochim. Biophys. Acta* **1261**, 183–194.
2. S. Bienroth, W. Keller, and E. Wahle (1993) *EMBO J.* **12**, 585–594.
3. K. Murthy and J. Manley (1992) *J. Biol. Chem.* **267**, 14804–14811.

#### Suggestion for Further Reading

4. W. Keller and L. Minvielle-Sebastia (1997) A comparison of mammalian and yeast pre-mRNA 3'-end processing, *Curr. Opin. Cell Biol.* **9**, 329–336.

## Polyadenylation

Many [messenger RNAs](#) (mRNAs) are modified after [transcription](#) by the addition of multiple AMP moieties at the 3' terminus in both prokaryotic and eukaryotic organisms, but the biochemical mechanism of poly(A) addition and the function of the tails are significantly different in the two types of organisms. In eukaryotes, a large multiprotein complex recognizes specific sequences in the 3' untranslated region (3' UTR) of the transcript leading to the cleavage of the mRNA. The combined action of poly(A) polymerase, a poly(A)-binding protein, and a cleavage/polyadenylation specificity factor (CPSF) subsequently leads to the addition of poly(A) tracts that vary in length from organism to organism (1). Poly(A) tails have been implicated in conferring stability to mRNA transcripts, promoting their efficiency of [translation](#) and having a role in the transport of processed mRNA from the [nucleus](#) to the cytoplasm (2, 3). In addition, recent studies suggest a relationship between mRNA splicing and polyadenylation (1). Taken together, polyadenylation in eukaryotic organisms apparently plays an important role in the regulation of gene expression and possibly [development](#).

In contrast, poly(A) tail addition in prokaryotic organisms only seems to require an available 3' terminus and poly(A) polymerase (4). Furthermore, unlike in eukaryotes, only a small proportion of each prokaryotic mRNA is polyadenylated at any given time (5). Because the presence of a poly(A) tail on a prokaryotic transcript appears to target the molecule for rapid decay (4, 6), poly(A) tails are also most likely involved in gene regulation, but primarily from the perspective of controlling mRNA decay rates (4, 7). Another distinction between prokaryotes and eukaryotes is that 23S rRNA is the most prevalent polyadenylated species in *Escherichia coli* (4).

The discovery of enzymes in both prokaryotes and eukaryotes that could add poly(A) tails to RNA molecules occurred in the early 1960s (8, 9). However, it was another 10 years before poly(A) tails

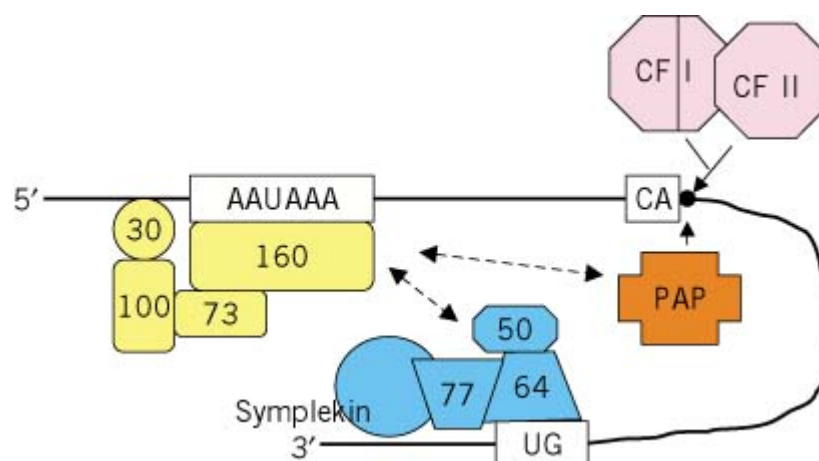
were actually identified as a post-transcriptional modification of eukaryotic mRNAs (10-12). Although the possibility of poly(A) tails in the bacteria *E. coli* was raised in 1975 (13), it was not until 1978 that solid evidence was presented for their existence (14). Subsequently, it has been shown that the length of the actual poly(A) tails varies from organism to organism [*E. coli*, 10–40 residues (6); yeast, ≈90 residues (15); vertebrates, ≈300 residues (16)].

## 1. Eukaryotes

The enzyme responsible for the actual addition of the poly(A) tail, poly(A) polymerase, appears to be ubiquitous in eukaryotic organisms, **archaea**, and both *Gram-negative* and **Gram-positive** bacteria. The prokaryotic and eukaryotic enzymes form part of a large superfamily of nucleotidyltransferase proteins that share significant amino acid **homology** in their amino termini (17, 18). The enzymes employ ATP to add AMP moieties to the 3' termini of RNA substrates.

In the case of the mammalian poly(A) polymerase, the length of the tail is directly related to the processivity of the addition reaction (19). Although poly(A) polymerases are single-subunit enzymes, the process for the selection of polyadenylation sites in eukaryotic transcripts involves both a large number of proteins and in most cases specific polyadenylation signals. The site recognition systems have been studied extensively in both mammals and yeast (20). In the mammalian system, there are three essential features of the polyadenylation element (Fig. 1). The polyadenylation signal sequence (AAUAAA) is normally located 10 to 30 bases upstream of the cleavage/polyadenylation site. A second motif that is GU rich is located 20 to 40 nucleotides downstream of the cleavage site (Fig. 1). Finally, a CA dinucleotide immediately 5' to the cleavage site is preferred but not absolutely required (21).

**Figure 1.** Outline of the mammalian polyadenylation machinery. The three *cis*-acting RNA elements associated with the selection of a polyadenylation site are shown in rectangular boxes. The solid black circle denotes where the mRNA is cleaved to generate the polyadenylation site. The numbers inside the yellow and blue boxes indicate the subunit sizes in kDa. Yellow boxes are the components of the cleavage/polyadenylation specificity factor (CPSF). Blue boxes indicate the subunits of the cleavage stimulation factor (CstF). The pink octagons are the two cleavage factors (CF I and CF II). The vertical line in the CR I octagon indicates a dimeric protein. PAP is poly(A) polymerase. The double-headed dashed arrows indicate protein–protein interactions. This figure is loosely adapted from Refs. 20 and 25, and is not drawn to scale.



In addition to the **cis-acting** sequences, at least four protein factors are required to generate a 3' terminus for polyadenylation in eukaryotes. These include the cleavage/polyadenylation specificity factor (CPSF), which consists of four distinct polypeptides (160, 100, 73, and 30 kDa) (1). The 160-kDa subunit appears to bind directly to the AAUAAA recognition signal, but its specificity appears

to be increased in the presence of the other CPSF subunits (1). Various data also suggest that the 160-kDa CPSF subunit interacts with poly(A) polymerase.

A second recognition protein is the so-called cleavage stimulation factor (CstF) that consists of three subunits (77, 64, and 50 kDa) and interacts at the downstream GU-rich element (Fig. 1) (1). The 64-kDa protein contains an **RNA-binding** domain (22); the 77-kDa subunit has been shown to associate with the 160-kDa CPSF subunit, indicating physical contact between CPSF and CstF. The nuclear protein **symplekin** has recently been shown to assist in the assembly of the CstF complex (23).

Although the CPSF and CstF protein complexes have been extensively studied, the multisubunit cleavage factors, CF I and CF II (24), associated with the actual mRNA cleavage reaction, remain poorly characterized. The CF I complex consists of two subunits (68 and 25 kDa). The exact composition of CF II remains to be determined. At the present time it has not been determined which protein factor(s) contain the actual endonuclease activity (25).

Once the polyadenylation site has been generated through the interaction of CPSF, CstF, CF I, and CF II, polyadenylation is initiated by poly(A) polymerase. However, another protein, called poly(A)-binding protein II (PAB II) plays an integral role in mediating the length of the tail (26). PAB II binds to the poly(A) tail early in its synthesis, and helps form a quaternary complex between the substrate RNA, CPSF, and poly(A) polymerase. This complex helps to provide the processivity of the polyadenylation reaction, because dissociation of the complex leads to a cessation in poly(A) tail elongation.

All eukaryotic polyadenylation systems are not identical. There are significant differences between the poly(A) addition process in yeast and in mammalian cells, but many of the proteins involved do share considerable sequence homology (20, 25). Yeast homologues of all four CPSF subunits have been identified along with several additional factors (1). Furthermore, it appears the mammalian AAUAAA and GU-rich elements are replaced by an A-rich positioning element along with an AU efficiency element (27). This would suggest that although there may be differences in the actual factors involved in 3' end formation, the basic process of polyadenylation in eukaryotic organisms has been functionally conserved (25).

## 2. Prokaryotes

In contrast to the well-characterized mammalian and yeast polyadenylation systems, polyadenylation in bacteria is still only partially understood. Recent experiments have demonstrated that the poly(A) levels in *E. coli* are regulated by controlling the amount of **poly(A) polymerase I** (4). In addition, there is a competition between the synthesis of poly(A) tails by poly(A) polymerase I and their degradation by 3' → 5' exonucleases such as **RNase II** and **polynucleotide phosphorylase** (28). Interestingly, poly(A) tails in *E. coli* can contain non A residues (4), which are added by polynucleotide phosphorylase working biosynthetically (29). In this reaction, nucleoside diphosphates are employed instead of nucleoside triphosphates. In the absence of poly(A) polymerase I, all of the residual poly(A) tails are synthesized by polynucleotide phosphorylase (29). Thus, unlike in eukaryotes, there are at least two distinct enzymes that can synthesize poly(A) tails in *E. coli*. Interestingly, it has just been reported that in spinach chloroplasts, all the poly(A) tails are synthesized by polynucleotide phosphorylase (30).

Although the *E. coli* poly(A) polymerase I is both functionally and structurally related to eukaryotic poly(A) polymerases (17, 18), at the current time there is no evidence for a specific polyadenylation signal, because poly(A) tails can be located at many different sites within a particular transcript (4, 31). The poly(A) tails synthesized in the bacterium appear to target specific mRNAs for more rapid degradation (4, 6, 32), but it is not clear whether polyadenylation in *E. coli* has any other functional role beyond its involvement in mRNA decay. Amino acid sequence comparison among the many newly sequenced bacterial genomes has demonstrated that poly(A) polymerases are found in most

prokaryotes.

### 3. Use in Purification of mRNA

From a practical perspective, polyadenylation has provided scientists with a powerful method for selecting specific mRNAs that can be converted into DNA and subsequently cloned (see cDNA Libraries). In this fashion, it has been possible to isolate the coding sequences for a large number of eukaryotic genes. The method involves initial enrichment of polyadenylated RNA by passage of total cellular RNA over a solid support matrix containing oligo(dT) molecules. In eukaryotes, the polyadenylated RNA will bind to the column, whereas the bulk of the RNA (ribosomal and transfer RNAs) will pass through the column. With bacterial RNAs, some rRNA will be associated with the column because they contain poly(A) tails. The polyadenylated RNA can then be eluted and used as a template for the enzyme **reverse transcriptase** using an oligo(dT) primer. The RNA/DNA hybrid can then be converted to fully double-stranded DNA in a second reaction involving **RNase H** and reverse transcriptase or a **DNA polymerase**. These fully double-stranded DNAs can be cloned into an appropriate vector for further analysis. Since poly(A) tails exist in both eukaryotes and prokaryotes, this method can be used with almost any organism.

### Bibliography

1. D. F. Colgan and J. L. Manley (1997) *Genes Dev.* **11**, 2755–2466.
2. A. B. Sachs, P. Sarnow, and M. W. Hentze (1997) *Cell* **89**, 831–838.
3. M. Wickens, P. Anderson, and R. J. Jackson (1997) *Curr. Opin. Genet. Dev.* **7**, 220–232.
4. B. K. Mohanty and S. R. Kushner (1999) *Mol. Microbiol.* **34**, 1094–1108.
5. G. -J. Cao and N. Sarkar (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 10380–10384.
6. E. B. O'Hara, J. A. Chekanova, C. A. Ingle, Z. R. Kushner, E. Peters, and S. R. Kushner (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 1807–1811.
7. S. R. Kushner (1996) In *Escherichia coli and Salmonella: Cellular and Molecular Biology*, Second Edition (F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, J. Low, K.B., B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), ASM Press, Washington, DC, 1996.
8. J. August, P. J. Ortiz, and J. Hurwitz (1962) *J. Biol. Chem.* **237**, 3786–3793.
9. M. Edmonds and R. Abrams, (1960) *J. Biol. Chem.* **235**, 1142–1148.
10. J. E. Darnell, R. Wall, and R. J. Tushinski (1971) *Proc. Natl. Acad. Sci. U.S.A.* **68**, 1321–1325.
11. S. Lee, J. Mendecki, and G. Brawerman (1971) *Proc. Natl. Acad. Sci. U.S.A.* **68**, 1331–1335.
12. M. Edmonds, M. H. Vaughan, and H. Nakazato (1971) *Proc. Natl. Acad. Sci. U.S.A.* **68**, 1336–1340.
13. P. R. Srinivasan, M. Ramanarayanan, and E. Rabbani (1975) *Proc. Natl. Acad. Sci. U.S.A.* **72**, 2910–2914.
14. N. Sarkar, D. Langley, and H. Paulus (1978) *Biochemistry* **17**, 3468–3474.
15. B. Groner and S. L. Phillips, (1975) *J. Biol. Chem.* **250**, 5640–5646.
16. G. Brawerman (1981) *CRC Crit. Rev. Biochem.* **10**, 1–38.
17. G. Martin and K. W. (1996) *EMBO J.* **15**, 2593–2603.
18. D. Yue, N. Maizels, and A. M. Weiner (1996) *RNA* **2**, 895–908.
19. E. Wahle (1995) *J. Biol. Chem.* **270**, 2800–2808.
20. W. Keller and L. Minvielle-Sebastia (1997) *Curr. Opin. Cell Biol.* **9**, 329–336.
21. F. Chen, C. Macdonald, and J. Wilusz (1995) *Nucl. Acids Res.* **23**, 2614–2620.
22. Y. Takagaki, C. C. Macdonald, T. Shenk, and J. L. Manley (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 1403–1407.
23. Y. Takagaki and J. L. Manley (2000) *Mol. Cell Biol.* **20**, 1515–1525.

24. Y. Takagaki, L. C. Ryner, and J. L. Manley (1989) *Genes Dev.* **3**, 1711–1724.
25. A. J. Shatkin and J. L. Manley (2000) *Nat. Struct. Biol.* **7**, 838–842.
26. E. Wahle, A. Lustig, P. Jenö, and P. Maurer (1993) *J. Biol. Chem.* **268**, 2937–2945.
27. Z. Guo and F. Sherman (1996) *Trends Biochem.* **21**, 477–481.
28. B. K. Mohanty and S. R. Kushner (2000) *Mol. Microbiol.* **36**, 982–994.
29. B. K. Mohanty and S. R. Kushner (2000) *Proc. Natl. Acad. Sci. U.S.A.* **97**, 11966–11971.
30. S. Yehudai-Resheff, M. Hirsh, and G. Schuster (2001) *Mol. Cell. Biol.* **21**, 5408–5416.
31. J. Haugel-Nielsen, E. Hajnsdorf, and P. Regnier (1996) *EMBO J.*, **15**, 3144–3152.
32. E. Hajnsdorf, F. Braun, J. Haugel-Nielsen, and P. Regnier (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 3973–3977.

### **Suggestions for Further Reading**

33. G. A. Coburn and G. A. Mackie (1999) *Prog. Nucl. Acid Res.* **62**, 55–108.
34. D. F. Colgan and J. L. Manley (1997) *Genes Dev.* **11**, 2755–2766.
35. N. Sarkar (1997) *Annu. Rev. Biochem.* **66**, 173–197.
36. W. Keller and L. Minvielle-Sebastia (1997) *Curr. Opin. Cell Biol.* **9**, 329–336.
37. A. J. Shatkin and J. L. Manley (2000) *Nat. Struct. Biol.* **7**, 838–842.

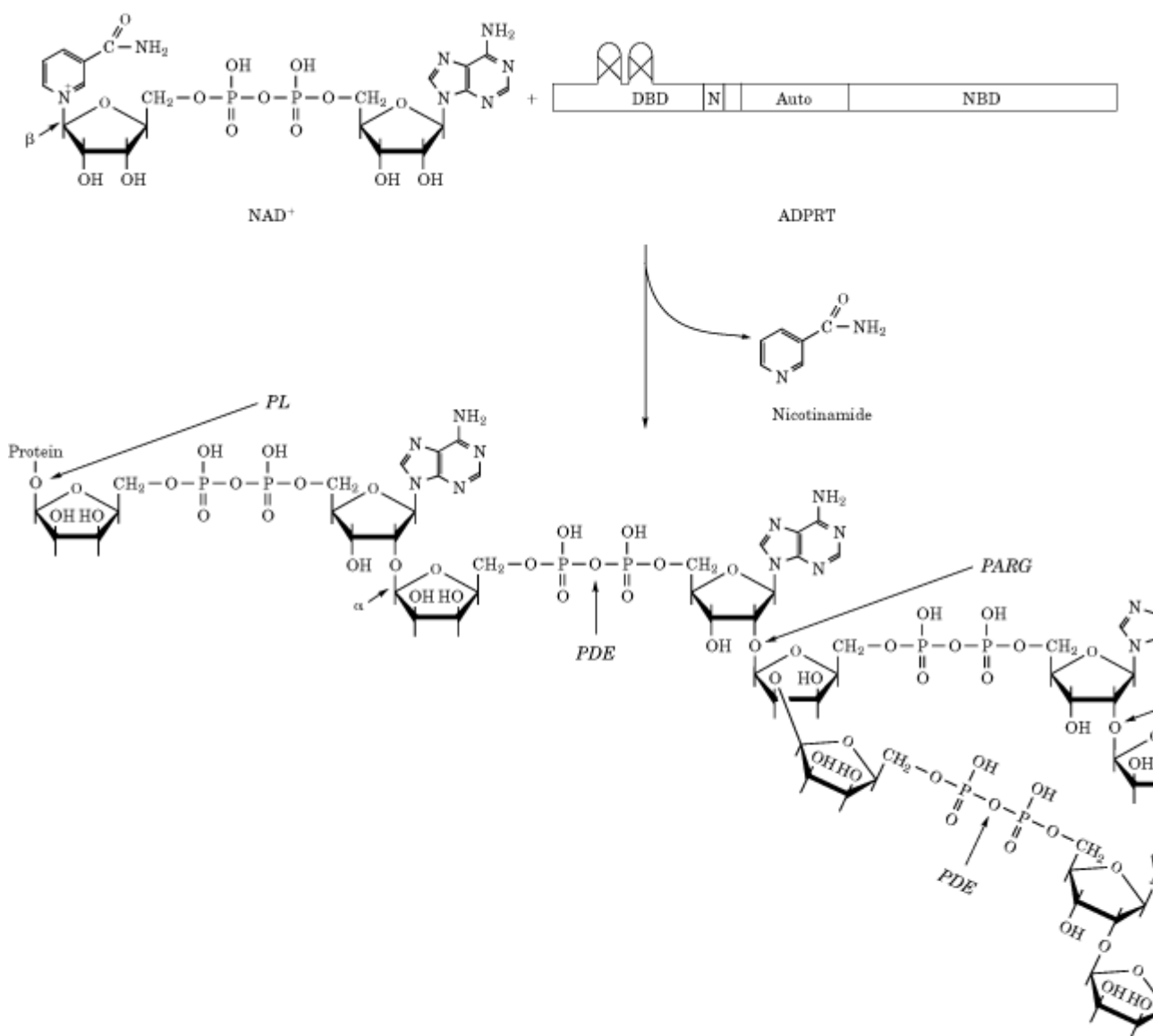
### **Poly(ADP-ribosyl)ation**

Poly(ADP-ribosyl)ation is a post-translational modification, and plays a role in regulating central cellular processes in eukaryotes. Poly(ADP-ribosyl)ation is catalyzed by the enzyme AND<sup>+</sup> protein (ADP-ribosyl) transferase polymerizing (EC.2.4.2.30). The specific inhibition of poly(ADP-ribosyl)ation in cells enhances the cytotoxic effects of carcinogens consistently. Therefore, poly(ADP-ribosyl)ation has long been considered to function in cellular surveillance of genotoxic stress. The most abundant enzyme responsible for this activity is poly(ADP-ribose) polymerase-1 (PARP-1). PARP-1 has been isolated and extensively characterized. PARP-1 has been demonstrated to specifically recognize and bind to DNA single-strand breaks and is, therefore, regarded as a “molecular nick sensor.” In an immediate response to the occurrence of DNA lesions, PARP-1 synthesizes polymers of ADP-ribose. Although not essential, PARP-1's activity has been shown to greatly facilitate DNA base excision repair. Recently, several other enzymes catalyzing this kind of modification have been identified. The discovery and characterization of further PARP enzymes suggest new functions for poly(ADP-ribosyl)ation, such as in regulation of telomere elongation and cellular transport.

The first report that led to the discovery of (ADP-ribosyl)ation was published in 1963 (1). It was observed that incorporating radioactively labeled ATP into the acid-insoluble fraction of a nuclear preparation from chicken liver was stimulated 1000 times by nicotinamide mononucleotide (NMN). The polymer product was later identified as poly(ADP-ribose). The synthesis of poly(ADP-ribose) depends on the availability of AND<sup>+</sup>, which is synthesized from NMN and ATP. PARP enzymes use NAD<sup>+</sup> as substrate and form poly(ADP-ribose), with a concomitant release of nicotinamide (Fig. 1). Multiple units of ADP-ribose are transferred onto glutamic acid. Therefore, it was suggested that PARP-1 could be an element of the G<sub>2</sub> checkpoint. Under conditions of severe damage, PARP-1 activation may cause depletion of the cellular NAD<sup>+</sup> pool and thereby can trigger cell death. During

the onset of apoptosis, specific interleukin-converting enzyme proteinases are activated, which results in the appearance of specific proteolytic PARP-1 fragments (14). These proteolytic products are therefore considered as markers of apoptosis.

**Figure 1.** Poly(ADP-ribosylation) reaction: The functional domains of mammalian PARP-1 are schematically represented including two Zn-finger motifs (indicated by the x loops) and a nuclear localization signal (N), the automodification (Au) domains are indicated. A PARP-enzyme catalyzes the polymerization of (ADP-ribose) moieties and the transfer to protein with concomitant release of nicotinamide. The conformation (a or b) of glycosidic bonds in NAD<sup>+</sup> and poly(ADP-ribose) involved in polymer catabolism are indicated. PL, protein lyase; PDE, phosphodiesterase; PARG, poly(ADP-ribose) gly



In living cells poly(ADP-ribose) polymers possess a half-life of only a few minutes, due to the activity of specific catabolizing enzymes. The cleavage sites are indicated in Figure 1. Poly(ADP-ribose) glycohydrolase is the physiological counterpart of PARP enzymes, and cleaves the polymers

with high specificity for the glycosidic bounds, generating free ADP-ribose (15). The transient nature of poly(ADP-ribose) synthesis makes this reaction also well-suited for the putative regulation of other cellular processes. At about the same time, resulting from independent research approaches, several new proteins with partial homology to the catalytic domain of PARP-1 were identified: PARP-2, PARP-3, vault-PARP (VPAAP), and tankyrase (see Table 1). Interestingly, also an alternative catalytic active form of PARP-1 was isolated (sPARP-1), which appears to be a short product resulting from an alternative initiation event in the PARP-1 gene (16). Of the new PARP members PARP-2 bears the strongest resemblance to PARP-1, and might also be functionally similar (17). VPAAP was isolated from vault complexes and is a 193-kDa protein containing a catalytic-active PARP domain (18). Only little is known about vault particles; however, it is speculated that these largest known ribonucleoprotein complexes serve a carrier function. Therefore, poly(ADP-ribosylation) by VPAAP might play a role in regulation of cellular transport. Another nuclear PARP enzyme, tankyrase was identified through its interaction with the telomere-specific DNA-binding protein TRF1 (7). TRF1 is a negative regulator of telomere elongation, and if poly(ADP-ribosylated), is released from binding to telomeres (7). Thus, tankyrase may function as a telomere-length regulator by promoting telomere elongation via poly(ADP-ribosylation) of TRF1 (19). Database search analyses indicate that the number of PARP-related proteins still is growing (Table 1). For example, recent studies have revealed a tankyrase homolog, termed tankyrase 2 (20).

**Table 1. Enzymes Involved in Poly(ADP-ribose) Metabolism**

| <b>Name of Enzyme</b> | <b>Alternative Designations</b> | <b>Chromosomal Location</b> | <b>Cellular Distribution</b> | <b>Putative Functions in Regulation of</b> |
|-----------------------|---------------------------------|-----------------------------|------------------------------|--|
| <b>PARP-1</b>         | ADPRT                           | 1q41-q42                    | Nucleoplasm                  | DNA base excision repair                   |
|                       | PARS                            |                             | Damaged DNA                  | Recombination                              |
|                       | ADPRT1                          |                             | Centromeres                  | Transcription                              |
|                       | pADPRT-1                        |                             |                              | Cell cycle                                 |
| <b>PARP-2</b>         | ADPRT2                          | 14q11.2-q12                 | Nucleoplasm                  | Cell death<br>DNA base excision repair     |
|                       | pADPRT-2                        |                             | Damaged DNA                  |  |
|                       | ADPRTL2                         |                             |                              |  |
| <b>PARP-3</b>         | ADPRTL3                         | 3q22.2-q21.1                |                              |  |
| <b>PARP-4</b>         | ADPRTL4                         | 22q11.1                     |                              |  |
| <b>VPAAP</b>          | VAULT3                          | 13q11                       | Cytoplasm                    | Cellular transport                         |



|                    |             |          |                     |                                 |
|--------------------|-------------|----------|---------------------|---------------------------------|
|                    | ADPRTL1     |          | Vault particles     |                                 |
|                    |             |          | Nucleoplasm         |                                 |
|                    |             |          | Mitotic spindle     |                                 |
| <b>Tankyrase</b>   | PARP-5      | 8q       | Telomeres           | Telomere length                 |
|                    | Tankyrase-1 |          | Mitotic centrosomes |                                 |
|                    | TNKS        |          | Nuclear pore        |                                 |
|                    |             |          | Golgi complex       |                                 |
| <b>Tankyrase-2</b> | TNKS-2      | 10q23.1  | Microsome fraction  | Signal transduction pathways    |
| <b>PARG</b>        |             | 10q11.23 | Nucleoplasm         | Catabolysis of poly(ADP-ribose) |
|                    |             |          | Cytoplasm           |                                 |

---

In resting cells, the natural content of detectable poly(ADP-ribose) is rather low, therefore the most important questions continue to be how and when PARP enzymes are activated and what are their relevant substrates of modification. The activation of PARP-1 in response to DNA damage is well understood. The characterization of the mechanisms for activation of the other PARPs will gain further insights into the biological roles of poly(ADP-ribosyl)ation.

#### Bibliography

1. P. Chambon, J. D. Weil, and P. Mandel (1963) *Biochem. Biophys. Res. Commun.* **11**, 39–43.
2. I. Kameshita, Z. Matsuda, T. Tanigushi, and Y. Shizuta (1984) *J. Biol. Chem.* **259**, 4770–4476.
3. F. Simonin, L. Höfferer, P. L. Panzeter, S. Müller, G. De Murcia, and F. R. Althaus (1993) *J. Biol. Chem.* **268**, 13454–13461.
4. A. Ruf, J. Menissier-de Murcia, G. De Murcia, and G. E. Schulz (1997) *Proc. Natl. Acad. Sci. U.S.A.* **93**, 7481–7485.
5. R. C. Benjamin and D. M. Gill (1980) *J. Biol. Chem.* **255**, 10502–10508.
6. S. L. Oei, J. Griesenbeck, M. Schweiger, and M. Ziegler (1998) *J. Biol. Chem.* **273**, 31644–31647.
7. S. Smith, I. Gariat, A. Schmitt, T. de Lange (1998) *Science* **282**, 1484–1487.
8. Z. Q. Wang et al (1995) *Genes Dev.* **9**, 509–520.
9. B. W. Durkacz, O. Omidiji, D. A. Gray, and S. Shall (1980) *Nature*, **283**, 593–596.
10. M. Molinete et al (1993) *EMBO J.*, **12**, 2109–2117.
11. M. Masson et al (1998) *Mol. Cell. Biol.* **18**, 3563–3571.

12. M. S. Satoh, G. G. Poirier, and T. Lindahl (1993) *J. Biol. Chem.* **268**, 5480–5487.
13. S. L. Oei and M. Ziegler (2000) *J. Biol. Chem.* **275**, 23234–23239.
14. S. H. Kaufmann, S. Desnoyer, Y. Ottaviano, N. E. Davidson, and G. G. Poirier (1993) *Cancer Res.* **53**, 3976–3985.
15. W. Lin, J. C. Ame, N. Aboul-Ela, E. L. Jacobson, and M. K. Jacobson (1997) *J. Biol. Chem.* **272**, 11895–11901.
16. F. R. Sallmann, M. D. Vodenicharov, Z. Q. Wang, and G. G. Poirier (2000) *J. Biol. Chem.* **275**, 15504–15511.
17. J. C. Ame et al (1999) *J. Biol. Chem.* **274**, 17860–17868.
18. V. A. Kickhoefer et al (1999) *J. Cell. Biol.* **146**, 917–928.
19. S. Smith and T. de Lange (2000) *Curr. Biol.* **10**, 1299–1302.
20. A. N. Kuimov et al (2001) *Genes Immun.* **2**, 52–55.

### Suggestions for Further Reading

21. D. D'Amours, S. Desnoyers, I. D'Silva, and G. G. Poirier (1999) Poly(ADP-ribosylation) reactions in the regulation of nuclear functions. *Biochem. J.* **342**, 249–268.
22. F. Dantzer et al (1999) Involvement of poly(ADP-ribose) polymerase in base excision repair. *Biochimie* **81**, 69–75.
23. M. K. Jacobson and E. L. Jacobson (1999) Discovering new ADP-ribose polymer cycles: Protecting the genome and more. *Trends Biochem. Sci.* **24**, 415–417.
24. A. A. Pieper, A. Verma, J. Zhang, and S. H. Snyder (1999) *Trends Pharmacol. Sci.* **20**, 171–181.
25. S. Shall and G. de Murcia (2000) Poly(ADP-ribose) polymerase-1: What have we learned from the deficient mouse model? *Mutat. Res.* **460**, 1–15.
26. S. Smith (2001) The world according to PARP. *Trends Biochem. Sci.* **26**, 174–179.
27. M. Ziegler and S. L. Oei (2001) A cellular survival switch: Poly(ADP-ribosylation) stimulates DNA repair and silences transcription. *Bioessays* **23**, 543–548.

## Polyamino Acids

Polyamino acids are **polymers** of varying lengths that are formed from only one (or sometimes a few) [amino acid\(s\)](#). Both polyamino acids and [proteins](#) are polymers formed from amino acids, but the difference is that proteins have a defined and complex sequence usually including all 20 different natural amino acids. The first peptides to be formed by chemical [peptide synthesis](#) were polyamino acids such as polyglycine, polyserine, and polylysine. These served as important models for structural and conformational studies. The two polyamino acids polyglycine and [polyproline](#) form an unusual helical backbone **secondary structure**.

[See also [Polymer](#) and [Proteins](#).]

## Polyclones

The term *polyclone* was introduced by Crick and Lawrence (1) and is intimately linked to the concept of compartments in [development](#) of the fruit fly *Drosophila* (2, 3). Compartments are well-defined parts of the body characterized by the lineage of their constituent cells; they are formed exclusively by the descendants of a small group of neighbor cells. The term *polyclone* is used to emphasize that each compartment always arises from more than one cell; the initial number varies from 5 to 50 cells. The descendants of each cell of the original group (the *founder polyclone*) will populate different regions of the adult compartment in different individuals but, together, the polyclone will always, in all individuals, construct the same region of the fly.

The polyclone is a unit of cell lineage in development. Compartmentalization is a reiterative process; during development, a polyclone can be subdivided into daughter polyclones, which will form specific subregions of the original compartment. The compartment/polyclone work provided, for the first time, a description of developmental segregations in development, that is, how to subdivide an organism into its constituent parts.

Although there were previous hints from lineage analyses in *Drosophila* and in the milkweed *Oncopeltus* (4), compartments were demonstrated as a result of the development of the *Minute technique* (5). This method allowed the production of marked clones that were able to proliferate more rapidly than surrounding cells and to reach very large size. These clones could nearly fill entire adult regions but would never transgress the boundaries of the compartment; thus, these fast-growing clones delineated compartments. These were originally described for the thoracic and cephalic [imaginal disks](#), but they were later found for the rest of the body also. The first compartmentalization events take place in early embryogenesis and affect all the germ layers, indicating that the whole body is compartmentalized.

Compartmentalization is an [epigenetic](#) subdivision of the body into parts—polyclones. These are not only fundamental units of cell lineage but also units of genetic control of development (6) and of growth and proliferation (7). Moreover, recent results (8) have demonstrated that compartment borders play a critical role in the signaling mechanisms involved in pattern formation.

## 1. Units of genetic control of development

A principal property of polyclones, and a basic tenet of the compartment hypothesis, is that polyclones are the realm of action of some key regulatory **genes** that establish developmental programs in groups of cells. For example, the identity of the body segments along the anteroposterior axis is specified by the *Hox* genes, whose domains of function and expression are delimited by compartment boundaries, indicating that *Hox* genes recognize polyclones as units of their expression.

Similarly, other genes involved with the specification of more discrete body regions become activated in specific polyclones. For example, the subdivision of embryonic metameres into anterior (A) and posterior (P) polyclones is followed by the activation of the **homeobox** gene *engrailed* in each P polyclone, whereas it is permanently turned off in the A polyclones. Thus, the P polyclone is the developmental unit of *engrailed* function, which gives P cells their specific identity.

A similar phenomenon occurs later, during the development of the wing (and presumably the haltere) disk, when a compartment boundary appears separating dorsal versus ventral polyclones (2, 3). All the cells of the dorsal polyclone, and none of the ventral, acquire activity of the homeobox gene *apterous*, which gives them specific dorsal identity (9).

## 2. Units of growth

Polyclones appear to be units of size control in development. This is indicated by the Minute

experiments (5), in which a fast-proliferating clone can fill as much as 80% to 90% of the compartment, yet it is of normal size. There must be a mechanism restricting the proliferation of the other cells of the polyclone in order to build a compartment of normal size. This mechanism has been called “cell competition” and operates within compartments (7).

### 3. Compartment borders as sources of morphogens

Recent work (8, 10) has shown that compartment borders play a key role in patterning processes. In the development of the wings and legs, signaling (through the Hedgehog product) from P polyclone cells across the anteroposterior compartment border to the adjacent A polyclone cells triggers the production of the [morphogen](#) *decapentaplegic*, which then diffuses to both polyclones. Similarly, the dorsoventral compartment border in the wing is the source of the signaling molecule *Wingless*.

### Bibliography

1. F. H. C. Crick and P. A. Lawrence (1975) *Science*, **189**, 340–347.
2. A. García-Bellido, P. Ripoll, and G. Morata (1973) *Nature New Biol.* **245**, 251–253.
3. A. García-Bellido, P. Ripoll, and G. Morata (1976) *Devel. Biol.* **48**, 132–147.
4. P. A. Lawrence (1973) *J. Embryol. Exp. Morph.* **30**, 681–699.
5. G. Morata and P. Ripoll (1975) *Devel. Biol.* **42**, 211–221.
6. G. Morata and P. A. Lawrence (1975) *Nature*, **255**, 614–617.
7. P. Simpson and G. Morata (1981) *Develop. Biol.* **85**, 299–308.
8. K. Basler and G. Struhl (1994) *Nature* **368**, 208–214.
9. F. Diaz-Benjumea and S. Cohen (1993) *Cell* **75**, 741–752.
10. T. Tabata and T. B. Kornberg (1994) *Cell* **76**, 89–102.

### Suggestions for Further Reading

11. G. Morata and P. A. Lawrence (1977) Homeotic genes, compartments and cell determination in *Drosophila*. *Nature* **265**, 211–216.
12. A. Garcia-Bellido, P. A. Lawrence, and G. Morata (1979) Compartments in animal development. *Sci. Amer.* **241**, 102–110.
13. P. A. Lawrence (1992) *The Making of a Fly*. Blackwell Scientific Publications, Oxford, U.K.
14. P. A. Lawrence and G. Struhl (1996) Morphogens, compartments and pattern: lessons from *Drosophila*? *Cell* **85**, 951–961.

### Polycomb Group

The *Drosophila* Polycomb group (PcG) comprises approximately 13 identified **genes** whose products are involved in **silencing** the [transcription](#) of target genes. Additional genes that contribute to PcG-dependent silencing are postulated to exist (1, 2). Although the list of genes that are regulated by the PcG continues to grow, PcG genes are defined primarily as negative regulators of the [homeotic genes](#) of the **Antennapedia** (ANT-C) and **bithorax** (BX-C) gene complexes. The homeotic genes must be correctly expressed throughout embryonic, larval, and pupal [development](#) in order to instruct the cells within each body segment correctly as to their respective segmental identities. Misexpression of homeotic genes results in the inappropriate development of specific body parts in place of normal structures. For example, ectopic expression of the *Antennapedia* gene

in head cells causes legs to develop in place of antennae. Expression of homeotic genes during early embryogenesis is controlled initially by [transcription factors](#) that are encoded by the segmentation genes. However, the segmentation proteins are degraded shortly after the expression patterns of homeotic genes are established. From that point onward, continuing through larval and pupal development, PcG proteins are responsible for maintaining the silenced states of homeotic genes in cell lineages in which they are initially repressed by segmentation proteins. A second group of proteins, encoded by the **trithorax** group (*trxG*) genes, are required to maintain the fully active states of homeotic genes in those cell lineages in which they are meant to be expressed. It is generally believed that control of target gene expression by these opposing groups of proteins involves modification of [chromatin](#) structure.

All PcG proteins that have been studied sufficiently are nuclear proteins that bind to [chromosomes](#) in a DNA-sequence-dependent manner and are either largely or completely co-localized on [polytene chromosomes](#). Multiple PcG proteins have been shown to **immunoprecipitate** together from fly embryo extracts, further demonstrating their *in vivo* association. Several have been shown to interact directly *in vitro*. Only one identified PcG protein possesses **DNA-binding** activity. Thus, the PcG proteins appear to function as components of multimeric protein complexes, and the association of most PcG proteins with target loci must be mediated by [protein-protein interactions](#). The varied and pleiotropic **phenotypes** produced by most PcG mutations, however, and recent biochemical and immunoprecipitation studies collectively suggest that multiple PcG complexes composed of different combinations of proteins may exist at different loci *in vivo* (3-5).

PcG homologues have been identified in [nematodes](#), vertebrates, and even plants. Molecular, biochemical, and genetic analyses of these **homologues** suggest that PcG-dependent silencing has been conserved throughout [evolution](#) across a wide **phylogenetic** spectrum.

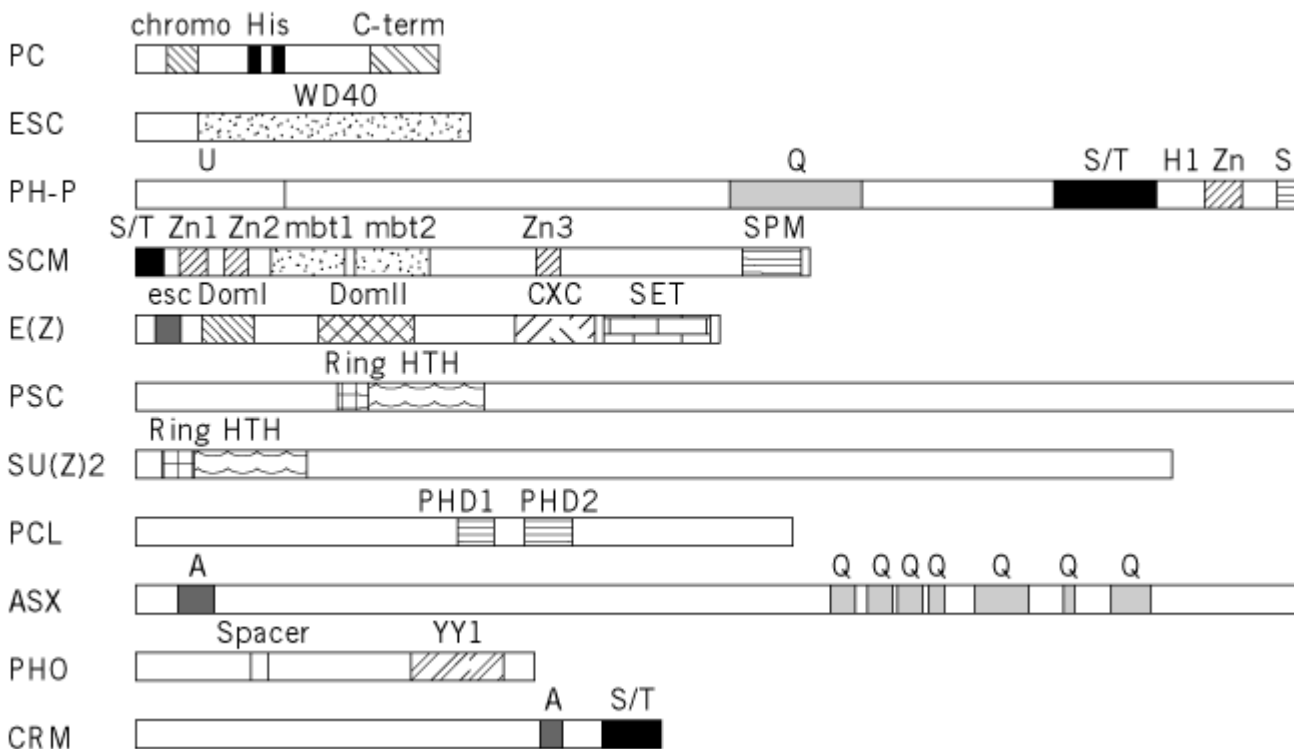
## 1. Members of the Polycomb group

This section includes a description of each individual *Drosophila* PcG gene, its protein product, and nonfly homologues (if identified).

### 1.1. Polycomb (Pc)

Mutant alleles of *Pc* were first identified on the basis of their dominant adult homeotic phenotypes, which were later found to be due to derepression of homeotic genes. For example, many *Pc* heterozygotes display partial transformation of second legs into first legs—that is, mesothoracic legs into prothoracic (6); this is caused by ectopic expression of the *Sex combs reduced* gene (of the ANT-C) during larval development (7). This is most easily apparent by the ectopic development of specialized bristles called sex combs, which normally are only found on the first legs of males, on the second legs. It is on the basis of this phenotype that *Pc*, and subsequently some of the other PcG genes, were named. Homozygous *Pc* mutants die as embryos and display posteriorly directed homeotic transformations of most body segments (8). The function of *Pc* in regulating homeotic genes is conserved in its mammalian homologue *M33*. Knockout mutations of *M33* causes posterior transformations of the axial skeleton of mouse embryos (9). Furthermore, expression of *M33* in *Drosophila* partially rescues *Pc* mutant phenotypes, demonstrating conservation of the biochemical activity of PC, and probably that of other PcG proteins with which it interacts (10).

**Figure 1.** Diagrams of the primary structures of PcG proteins and likely domains identified by sequence analysis. PC: ch domain (chromo), His-rich region (His), conserved C-terminal region (C-term); ESC: WD40 repeats (WD40); PH-P: Un region not found in PH-D (U), Gln-rich region (Q), Ser/Thr-rich region (S/T), H1 domain, Zn finger (Zn), SPM domain (SCM: Ser/Thr-rich region (S/T), three Zn fingers (Zn1, Zn2, Zn3), mbt repeats (mbt1 and mbt2), SPM domain (SPM); E binding domain (esc), Domain I (DomI), Domain II (DomII), CXC domain (CXC), SET domain (SET); PSC and SU(Z): fingers (Ring), helix-turn-helix regions (HTH); PCL: two PHD fingers (PHD1 and PHD2); ASX: Ala-rich region (A), G regions (Q), Cys cluster (Cys); PHO: Spacer region of homology with YY1 (Spacer), Zn-finger-containing region with h to YY1 (YY1); CRM: Ala-rich region (A), Ser/Thr-rich region (S/T).



PC protein binds to about 100 [euchromatin](#) sites on polytene chromosomes, including the locations of known target genes, such as the ANT-C and BX-C ([11](#), [12](#)). This strongly suggests that the homeotic genes are only a few among the many genes that are regulated by PcG proteins. PC is 390 amino acid residues in length, with a predicted relative mass of 44 kDa, although it migrates at a higher molecular weight on [SDS-PAGE](#) gels, and it is composed of at least two functional **domains**. The N-terminal chromo domain is 58% identical to a 48-residue region of the heterochromatic *Drosophila* HP1 protein ([13](#)), which is involved in position-effect variegation (PEV) ([14](#)). The chromo domain is required for polytene chromosome binding by PC, whereas the evolutionarily conserved C-terminal domain is not ([15](#)). However, deletions of this region do alter the subnuclear distribution of PC in diploid embryonic nuclei, but not as dramatically as do mutations in the chromo domain. ([16](#)). A chimeric HP1-PC protein, in which the HP1 chromo domain is replaced by that of PC, binds to both [heterochromatin](#) and to the normal euchromatic PC sites, and it recruits other PcG proteins, including endogenous PC, polyhomeotic (PH), and Posterior sex combs (PSC), to heterochromatin ([17](#)). Consistent with these observations, biochemical studies suggest that the PC chromo domain is capable of forming multimers *in vitro* ([18](#)). Conflicting results have been reported, however, regarding the ability (or inability) of PC to self-associate ([19](#), [20](#)). The function of a histidine-rich region in the middle portion of PC has yet to be defined.

PC is a component of large (~2MDa) multiprotein complexes and co-fractionates with PH in *Drosophila* embryo extracts ([21](#)). *In vivo* association of PC and PH is confirmed by their co-immunoprecipitation and co-localization on polytene chromosomes ([21](#)). PC and PH both also co-immunoprecipitate with PSC ([19](#)). PC does not directly interact with PH, but it does bind PSC *in vitro* and in yeast **two-hybrid assays**. PSC binding does not require the PC chromo domain ([19](#)), but *in vivo* PC-PH complexes are disrupted by mutations in the chromo domain ([4](#)). *Xenopus* homologues of PC and PSC, XPC, and XBMI-1, respectively, interact directly, demonstrating evolutionary conservation of this relationship ([20](#)). Conservation of this multiprotein complex is further demonstrated by *in vivo* association of M33, BMI-1, MEL-18 [BMI-1 and MEL-18 are mouse homologues of PSC ([22](#), [23](#))], and Rae28/MPH-1 [a mouse homologue of PH ([24](#), [25](#))] ([24](#)).

## 1.2. extra sex combs (esc)

Unlike other members of the PcG, *esc* appears to be involved only in regulation of the homeotic genes and is primarily only required during embryogenesis, during the window of time when PcG-dependent silencing is established (26, 27). Therefore, it has been proposed that ESC protein may somehow recognize the repressed state of these target genes and, through interactions with other PcG proteins, contribute to the assembly of silencing complexes (26, 28, 29). The 425-residue ESC protein has a relative molecular mass of 48 kDa. Apart from a short N-terminal region that includes a **PEST sequence**, ESC is composed primarily of seven WD40 repeats (26, 28, 29), which are likely to fold into a b-propeller toroid [protein structure](#) (see [Beta-Sheet](#)) (30). The primary activity of WD40 repeat domains, which are found in a wide variety of proteins with unrelated biological functions, is to interact with other proteins (31). One ESC-binding partner is another PcG protein, Enhancer of zeste [E(Z)] (32). Evolutionarily conserved residues within two adjacent ESC loops that project above the toroidal plane are required for E(Z) binding *in vitro* and are essential for ESC activity *in vivo* (30, 32). The physical relationship between ESC and E(Z) is conserved in their mouse and human homologues (3, 5, 32). The mouse embryonic ectoderm development (EED) protein shares 55% amino acid identity with ESC, although their N-termini are not well conserved (33). Human and mouse EED are identical (32). Homozygous *eed* mutant mouse embryos display posteriorly directed homeotic transformations, which is consistent with conservation of *esc* function in mammals (33). The [Caenorhabditis elegans](#) maternal-effect sterile genes, *mes-2* and *mes-6*, are required to maintain germline-specific transcriptional silence (34) and are homologues of *E(z)* and *esc*, respectively (35, 36). No other PcG proteins have been identified in [sequence database](#) searches of the *C. elegans* genome (36). This suggests that the function of ESC–E(Z) complexes may be at least partially independent of other PcG proteins.

### 1.3. polyhomeotic (ph)

The *ph* locus is tandemly duplicated and composed of two transcription units with largely redundant functions (*ph<sup>P</sup>* and *ph<sup>D</sup>*). Homozygotes with mutations in either gene alone are viable. However, mutations that disrupt both *ph<sup>P</sup>* and *ph<sup>D</sup>* are embryonic lethal when homozygous and produce posteriorly directed homeotic transformations. Embryos that also lack maternally produced *ph* fail to develop beyond the blastoderm stage (37). *ph<sup>P</sup>* and *ph<sup>D</sup>* encode 167 kDa and 149 kDa proteins, respectively (38, 39). The PH<sup>P</sup> protein contains 193 N-terminal residues that are not found in PH<sup>D</sup>. Otherwise, the two proteins are extremely similar. Both contain glutamine repeats, a Ser/Thr-rich region, an unusual Cys<sub>4</sub> **zinc finger**like repeat, and two regions, H1 and SPM, that are highly conserved in Rae28/MPH1 protein, a mouse homologue (24, 25). The H1 region binds to PSC protein, which may account for the indirect association of PH and PC *in vivo* (19). Interaction of BMI-1 with the human homologues of PH, HPH1, and HPH2, demonstrates conservation of this relationship (40). The largely **alpha-helical** SPM domain, which is related to a more divergent motif called the SAM domain that is found in proteins with a wide variety of functions (41), is most similar to 60-residue regions of the fly Sex combs on midleg (SCM) PcG and the **tumor suppressor** l(3)mbt proteins (42). Both the PH and SCM SPM domains are capable of forming homodimers, as well as heterodimers, with each other (43). *In vivo* association of PH and SCM is demonstrated by their complete co-localization on polytene chromosomes.

### 1.4. Enhancer of zeste [E(z)]

In addition to its role in maintaining the silence of homeotic genes, *E(z)* is also involved in the repression of the *white* gene by *zeste*<sup>1</sup> (44-46). This latter function is shared by *Scm*, *Psc*, and *Suppressor 2 of zeste* [*Su(z)2*] (46). Complete lack of *E(z)* activity blocks cell proliferation (47, 48). Consistent with a role in general chromatin organization, 45 identified binding sites almost completely co-localize with those of other PcG proteins, plus low levels of E(Z) protein appear to be ubiquitously distributed along the euchromatic arms of polytene chromosomes (49). The 760-residue (87 kDa) E(Z) protein is composed of multiple functional domains. The 116-residue C-terminal SET domain is a region of homology shared by the fly trithorax (TRX) (50), SU(VAR)3-9 (a suppressor of PEV) (51), and *trxG* ASH1 proteins (52). Numerous other SET domain proteins have been identified in organisms ranging from yeast to humans (53, 54). The E(Z) SET domain has recently

been shown to bind mammalian myotubularin, a dual-specificity **phosphatase**, and the myotubularin-related protein Sbf, which functions as an antiphosphatase due to lack of phosphatase activity (55). Presumably, by tethering either phosphatase or antiphosphatase to target loci, one function of SET domain proteins, including E(Z), may be to modulate the state of phosphorylation of other chromatin proteins, and thus their activities. The E(Z) CXC domain is a **cysteine**-rich region that flanks the SET domain on the N-terminal side and is required for chromosome binding (49). Domains I and II are two evolutionarily conserved regions in the N-terminal half of E(Z), whose functions are yet to be defined (53). A 33-residue domain near the N-terminus binds ESC protein (32). E(Z) is required throughout development in order to maintain the silence of homeotic genes (45, 48, 56), whereas ESC is primarily only required during early embryogenesis (26, 27). Therefore, ESC-E(Z) complexes may contribute to the establishment of homeotic gene silencing, and E(Z) may play a transitional role in maintaining their permanent silenced states. *E(z)* is also involved in the regulation of a wide variety of genes independently of *esc* (46, 57, 58) and, indeed, also appears to contribute to *trxG*-dependent transcriptional activation during larval development (59). Since *esc* is only involved in the early silencing of homeotic genes, E(Z) must interact with proteins other than ESC at other target loci and at other developmental stages.

*E(z)* homologues have been identified in mice [*Ezh1/Enx2* (53); *Ezh2/Enx1* (60)], humans [*EZH1/ENX2* (61); *EZH2/ENX1* (53, 62, 63)], *C. elegans* [*mes-2* (35)], and *Arabidopsis thaliana* [*curly leaf* (*clf*) (64) and *medea* (65)]. Most strikingly, like *E(z)*, *Arabidopsis clf* is involved in maintaining, not initiating, the silenced state of a homeotic target gene. Thus, the function of E(Z) has been conserved across kingdoms.

### 1.5. Posterior sex combs (PSC)

*Psc* mutations produce dominant anterior transformations of adult legs and suppression of the *zeste*<sup>1</sup> eye color (46, 66). Homozygous *Psc* mutant embryos exhibit posteriorly directed homeotic transformations and ectopic expression of homeotic genes (1, 67, 68). Mouse homologues of *Psc* [and *Su(z)2*], *bmi-1*, and *mel-18* are a proto-**oncogene** and tumor suppressor gene, respectively (69-71). Knockout mutations of the mouse *bmi-1* or *mel-18* genes produce posteriorly directed homeotic transformations of the axial skeleton and ectopic expression of homeotic genes (72, 73). These mutants also exhibit defects in T- and B-lymphocyte development, in addition to a variety of other maladies. Overexpression of *bmi-1* in transgenic mice produces opposite, anteriorly directed transformations (74).

PSC protein is 1603 residues in length, with a relative mass of 170 kDa (22). It contains an ~200-residue region (the HR domain) with 37.4% identity to a similar region in the SU(Z)2 protein. The HR domain includes a Cys<sub>3</sub>-His-Cys<sub>4</sub> RING finger, followed by a **helix-turn-helix** (HTH) motif, which binds to the H1 domain of PH protein (19). The HTH and RING finger regions bind to PC protein (19). PSC is largely, but not completely, co-localized with PC and PH on polytene chromosomes (75, 76), which is consistent with both their *in vivo* association and the involvement of PSC in regulation of genes independently of PC and PH. The BMI-1 RING finger is required for subnuclear localization (24). Although the MEL-18 RING finger is able to bind DNA (77), this activity has not been demonstrated for PSC.

### 1.6. Polycomblike (Pcl)

Originally identified as a dominant **enhancer** of *Pc* (78), homozygosity for *Pcl* alleles results in embryonic lethality and moderate posteriorly directed homeotic transformations (1). *Pcl* encodes an 857-residue protein that is ubiquitously expressed in embryonic nuclei and completely co-localized with PC on polytene chromosomes (12). The sequence of PCL is unique, with the exception of two adjacent PHD fingers near the middle of the protein. The PHD finger motif has a Cys<sub>4</sub>-His-Cys<sub>3</sub> pattern and spans approximately 50 to 80 amino acid residues (79). Interestingly, three other *Drosophila* PHD finger-containing proteins are members of *trxG*: TRX, ASH1, and ASH2. Thus, although the function of PHD fingers has yet to be elucidated, this suggests that these *trxG* proteins and PCL may possess common biochemical activities. For example, if the PHD finger is a protein-



protein interaction motif, this shared domain may permit interaction with a common binding partner.

### 1.7. Sex combs on midleg (Scm)

Individuals that are homozygous for **null *Scm* alleles** die as embryos and display mild posteriorly directed homeotic transformations (1) and ectopic expression of *abd-A* (42). The 877-residue SCM protein, which completely co-localizes with PH on polytene chromosomes, contains several motifs shared by PH, *Rae28/MPH1*, and *l(3)mbt* (42). In addition to the SPM dimerization domain, SCM contains two Cys<sub>2</sub>-Cys<sub>2</sub> zinc fingers that have spacing unlike other Zn fingers [*PH*, *Rae28/MPH1*, and *l(3)mbt* each contain a single copy of similar zinc fingers], and two copies of ~100-residue repeat, three of which are present in *l(3)mbt* but absent from *PH* and *Rae28/MPH1*. The null or strongly hypomorphic *Scm*<sup>XF24</sup> allele contains a deletion of most of the C-terminal SPM domain, demonstrating the importance of this domain for SCM function (42).

### 1.8. Additional sex combs (Asx)

Mutant *Asx* alleles were initially identified on the basis of their PcG-like embryonic homeotic phenotypes (1). *Asx* alleles also increase the severity of homeotic transformations produced by *trx* mutant alleles (80). In addition, *Asx* is a dominant enhancer of PEV, which is more typical of *trxG* genes (81). Therefore, like *E(Z)*, *ASX* protein may be involved in both PcG-dependent silencing and *trxG*-dependent activation. *Asx* encodes a novel 1668-residue protein with a relative mass of 182 kDa (80). Noted features include two alanine-rich regions near the N-terminus, multiple glutamine repeats, and a cysteine cluster near the C-terminus with an unusual spacing of Cys residues. *ASX* binds to about 90 sites on polytene chromosomes; 70% of these overlap with *PC*, *PH*, and *PSC* binding sites, while 30% are unique. This reinforces the idea that, although *ASX* is required for PcG-dependent silencing of many target genes, including the homeotic genes, it may function independently of PcG proteins at many others.

### 1.9. pleiohomeotic (pho)

Individuals that are homozygous for *pho* null alleles die during pupal development. Those reaching pharate adult stage display posteriorly directed transformations of legs and partial antenna to leg transformations, characteristic of PcG mutations (82, 83). Lack of maternally produced *pho* results in embryonic lethality and severe defects (84). *pho* encodes the *Drosophila* homologue of the mammalian transcription factor *YY1*, which is a DNA-binding zinc-finger transcription factor (85). Binding sites for the *PHO* protein (520 residues in length) have been identified in multiple Polycomb-response elements (PREs), including each of the *BX-C* PREs, an *Scr* PRE, and the *engrailed* *cis*-regulatory region (85, 86), which is also regulated by the PcG (57). Thus far, *PHO* is the only PcG protein shown to possess sequence-specific DNA-binding activity. Therefore, *PHO* may play a central role in assembling PcG protein complexes at PREs.

### 1.10. cramped (crm)

Males that are **hemizygous** for mutant alleles of the X-linked *crm* locus exhibit proximodistal transformations of leg segments, in addition to anterior transformations of legs that are typical of PcG mutations (87). These mutant phenotypes are enhanced by mutations in *Pc* and *mus209*, which encodes [proliferating cell nuclear antigen](#) (PCNA) (88). PCNA is essential for [DNA replication](#), suggesting a possible connection between PcG-dependent silencing and DNA replication. Both *crm*<sup>sa</sup> and *mus209*<sup>B1</sup> also suppress PEV (88, 89). *crm* encodes a novel 693-residue protein with a relative mass of 75 kDa (88). It contains one alanine-rich region and a serine/threonine-rich region near the C-terminus. *CRM* is a nuclear protein that is ubiquitously expressed and associated with embryonic diploid chromosomes during interphase, but not during mitosis. *CRM* and *PCNA* appear to be co-localized in both embryonic and polytene nuclei, and both display subnuclear patterns different from those of *PC*.

### 1.11. multi sex combs (mxc)

Adult viable and pupal lethal alleles produce homeotic phenotypes typical of other PcG mutations. Careful analysis of recessive lethal alleles has also revealed a tumor suppressor function (90). In

addition, lack of *mxo* activity specifically blocks cell proliferation in the germline (91). Cloning of the *mxo* gene has not been reported.

#### 1.12. Sex combs extra (*Sce*)

*Sce* homozygotes derived from heterozygous mothers die as first instar larvae, with weak posteriorly directed homeotic transformations. *Sce* homozygotes derived from *Sce* mutant female germlines display extreme homeotic transformations (84). *Sce*'s identity as a PcG gene is further supported by genetic interactions with other PcG mutations (67, 92). Cloning of the *Sce* gene has not been reported.

#### 1.13. super sex combs (*sxc*)

Putative null *sxc* homozygotes die as pharate adults and display homeotic transformations. Absence of maternal and zygotic expression results in embryonic lethality and mild posteriorly directed homeotic transformations (93). Cloning of the *sxc* gene has not yet been reported.

#### 1.14. Enhancer of Polycomb [E(Pc)]

Mutant *E(Pc)* alleles do not by themselves produce homeotic phenotypes, making it questionable whether to categorize *E(Pc)* as a PcG gene. They do, however, enhance the severity of phenotypes produced by other PcG mutations (94, 95). Therefore, it is likely that E(PC) protein contributes some biochemical activity that is associated with PcG-dependent silencing. Unlike most PcG genes, *E(Pc)* is also a suppressor of PEV, suggesting that it may provide a connection between these two groups of proteins that are involved in chromatin-mediated transcriptional regulation (81). Cloning of the *E(Pc)* gene has not yet been reported.

#### 1.15. Suppressor 2 of zeste [*Su(z)2*]

Like *E(Pc)*, *Su(z)2* alleles do not by themselves produce homeotic phenotypes. Mild interactions have been observed between *Su(z)2* alleles and a few PcG loci. For example, the recessive embryonic lethality of *Su(z)2* null alleles is partially suppressed by *pho* alleles, and dominant gain-of-function *Su(z)2* phenotypes are partially suppressed by *Sce* (96). SU(Z)2 protein binds to many of the same polytene chromosome sites as other PcG proteins (76). This observation, in addition to the presence of the HR domain that it shares with PSC, suggests that SU(Z)2 probably contributes to a similar aspect of chromatin-mediated transcriptional regulation but may not be involved in regulation of some PcG target genes, such as the homeotic genes. The *Su(z)2* gene is located adjacent to the *Psc* gene. The RING finger-containing HR domain of SU(Z)2, which is conserved in PSC and homologous proteins from other organisms, is responsible for locus-specific polytene chromosome binding (97).

## 2. Mechanisms of PcG-dependent silencing

PcG-dependent silencing requires the presence of DNA sequences that are referred to as Polycomb-response elements (PREs) (98), which are distinct from both promoters and enhancers of target genes and can control their activities over distances of 20 to 30 kbp. When placed within P element transgenes, PREs convey PcG-dependent silencing upon reporter genes (98) and create new PcG protein binding sites at their sites of insertion (99). Such studies have narrowed the sizes of PREs to a few hundred bp; however, it is likely that additional weak PREs are located close to the primary PREs, since larger PRE-containing fragments exhibit more consistent silencing and PcG protein binding (100). Comparison of the sequences of different PREs has revealed little concerning their mechanism of action. PREs contain binding sites for GAGA factor (101), the product of the *trxG* *Trithorax-like* gene (102), suggesting that GAGA factor chromatin remodeling activity may contribute to both *trxG*-dependent activation and PcG-dependent silencing. However, GAGA alone cannot be responsible for PRE activity, since GAGA binding sites are present at many genomic locations that do not function as PREs. Recently, binding sites for the PcG PHO protein have been identified in multiple PREs (86), suggesting a central role for PHO in PcG complex assembly.

Although we are beginning to learn more about the activities of individual PcG proteins, and it

appears that they function as components of multimeric protein complexes and require the presence of a PRE in order to silence a target gene, a great deal about PcG-dependent silencing remains mysterious. What is unknown or unclear may be broken down into three major questions.

### 2.1. 1. How do PcG proteins recognize the repressed state of target genes?

Most PcG proteins are ubiquitously distributed, at least in embryos, yet they only silence target genes in those cell lineages in which the genes are initially repressed by gene specific transcription factors, such as the products of segmentation genes. So far, there is no solid evidence indicating how this may happen. It is possible that one or more PcG proteins may interact with a repressor, such as HB or Kruppel (KR), leading to the recruitment of additional PcG proteins, but no such interaction has yet been detected. Alternatively, one or more PcG proteins may somehow recognize the repressed conformation of target genes (103). Such distinguishing features could include (1) lack of enhancer-promoter loop formation; (2) repressor-dependent availability of PcG binding site(s), or that of a protein that assists in the assembly of PcG silencing complexes; or (3) some other characteristic of repressed chromatin that distinguishes inactive from active genes. Whatever the mechanism, it may only be possible to establish PcG silencing *de novo* during embryogenesis. Reporter genes within a P element that also contains a GAL4 upstream activating sequence and a PRE from the BX-C exhibit PcG-dependent silencing. If a pulse of GAL4 is expressed from a second transgene during larval development, this silencing can be overridden, but only transiently. If, however, the pulse of GAL4 expression occurs in embryos, activation is maintained throughout development (104). Therefore, some protein or activity that is necessary for establishment of PcG silencing must exist in embryos and not in larvae. One PcG protein that fits this description is ESC, which is only required during embryogenesis. Given the specificity of *esc* for homeotic gene repression, however, other protein(s) would be needed to provide that function at other PcG target genes.

### 2.2. 2. What is the mechanism of PcG-dependent silencing?

Several alternative models for PcG-dependent silencing have been proposed. They fall into three general categories.

#### 2.2.1. Chromatin accessibility models

The basic idea of chromatin accessibility models is that PcG proteins modify the chromatin structure of target genes in some way that sterically blocks activators from binding to their DNA targets in enhancer and/or promoter regions. On the basis of the shared presence of the chromo domain in PC and in the heterochromatic HP1, which is involved in PEV, it was first proposed that PcG complexes fold target genes into compacted heterochromatin-like states (13). PEV is another example of transcriptional silencing, which results when a normally euchromatic gene is juxtaposed with heterochromatin (105). Expression of that gene in some cell lineages, and its repression in others, may depend on whether it is packaged into a condensed heterochromatic state. Additional sequence motifs, such as the SET domain, subsequently have been identified that are shared by PcG/trxG and heterochromatic proteins that are involved in PEV (50-52). Therefore, it is likely that there exist some mechanistic similarities between heterochromatin and PcG complexes, at least at the level of the biochemical activities of individual proteins. PcG-dependent silencing shares some phenomenological similarities with PEV. For example, within P element constructs, PREs convey variegated expression patterns of the *white* reporter gene that are reminiscent of PEV repression of the *white<sup>mottled4</sup>* allele (106, 107). Zink and Paro (108) have provided evidence suggesting that such variegated expression results from competition for target-site binding by PcG proteins and activators, and that binding of PC or an activator (in this specific study, GAL4) to a target gene is mutually exclusive. *In vivo* **cross-linking** studies originally suggested that PcG proteins coat target genes, such as the BX-C, supporting the heterochromatin-like model (109). However, refined experimental techniques now suggest that they are limited to PRE-containing regions and that different combinations of proteins may associate with different sites (4).

Selective exclusion is a variation of the chromatin accessibility model, in which PcG complexes may screen proteins on the basis of size, shape, or some other quality. For example, the PcG is able to

block GAL4-activated transcription of a reporter gene by **RNA polymerase II** in *Drosophila* embryos, but cannot block transcription by T7 RNA polymerase (110). McCall and Bender (110) have proposed several alternative mechanisms of selective exclusion that might explain these observations. The PcG may modify the distribution or structure of **nucleosomes** in some way [e.g., by causing the deacetylation of core **histones** (111)] that inhibits DNA binding by some proteins, but not others. Alternatively, PcG complexes may themselves directly exclude binding by some other proteins on the basis of size or shape. Or PcG proteins could interact directly with and inhibit binding by specific proteins that are required for transcription by RNA polymerase II (e.g., activators, basal transcription factors, or accessory proteins such as the SWI/SNF complex). The distinguishing feature of all these variations of the selective exclusion model is that they propose relatively subtle modifications of chromatin structure compared to the heterochromatin model.

### 2.2.2. Enhancer interference model

It has been proposed that the PcG does not necessarily prevent transcription factors from binding to target genes but prevents them from activating transcription (112, 113). This could involve blocking enhancer-promoter loop formation or in some other way preventing activators from interacting with the basal transcription machinery. In support of this model, PcG proteins that are tethered to reporter plasmids in mammalian tissue culture cells have been shown to repress transcriptional activation in transient expression assays (112). This ability to repress varies with transcription factors that have different activation domains, suggesting an effect on activation but not access. These observations are also consistent with some versions of the selective exclusion model, in that the effects on different activators may be due to the exclusion of different accessory proteins that may be required for activation by some, but not other, transcription factors. The ability of the PcG to block GAL4-activated transcription by RNA polymerase II, but not transcription by T7 RNA polymerase (110), may also be explained by enhancer interference. In this scenario, the PcG might not prevent GAL4 from binding to DNA but might block its interaction with basal transcription factors or accessory proteins that are needed for RNA polymerase II transcription. This interpretation, however, would be in conflict with the mutually exclusive binding of PC and GAL4 to polytene chromosomes (108).

### 2.2.3. Subnuclear compartmentalization

It has been suggested that PcG-dependent silencing may sequester target genes in transcriptionally inert nuclear compartments (114). This model is supported by two lines of observations. First, in **restriction enzyme** digestions of intact chromatin, PcG-dependent silencing does not have a detectable effect on access of the enzymes to recognition sites located within target genes (114). The second comes from analogy to PEV, in which silenced genes are associated with centromeric heterochromatin (115, 116). However, intranuclear distribution of PcG proteins does not correlate with high concentrations of DNA, demonstrating that silenced genes are not sequestered to heterochromatic regions (117). In addition, in some cells, one adjacent gene may be active, the other silent. Therefore, sequestering, if it occurs, must be on a microcompartmental level.

### 2.3. 3. How is the silenced state propagated through the cell cycle?

Once the silenced state of target genes is established during embryogenesis, it must be maintained through many cycles of cell division. This requires the continuous activity of all PcG proteins, with the exception of ESC, and continued presence of a PRE. Excision of a PRE results in loss of silencing, demonstrating that protein–protein interactions and/or chromatin modification alone are not sufficient to propagate the silenced state (118). Maintenance is particularly perplexing in light of recent observations that PC, PH, and PSC dissociate from chromatin during mitosis (117). It is not known if all PcG proteins behave in a similar manner, or if residual levels of one or more PcG protein remain associated with silenced loci. Should any PcG protein(s) not dissociate from chromatin during mitosis, it might provide nucleation sites for the reformation of silencing complexes. Given the DNA-binding activity of PHO and the presence of PHO binding sites in all of the major BX-C PREs, PHO may bind PREs (either remain bound throughout mitosis or rebind following telophase) and, through protein–protein interactions, direct the binding of other PcG proteins. This would be consistent with the continued requirement of PREs to maintain silencing. If PHO, or some other PcG protein, remains bound to PREs through mitosis, this may be sufficient to

reassemble silencing complexes at the correct loci. If, however, a protein such as PHO dissociates from DNA during mitosis, some other mechanism must distinguish silenced from active PcG target genes. This could involve chromatin modification. For example, histone deacetylation could mark the chromatin of silenced genes. In this model, a combination of PHO binding sites and local chromatin modification might be sufficient for PcG complex reassembly.

A remarkable feature of the aforementioned GAL4 activation of a PRE-containing transgene is the transmission of the active state to a fraction of the progeny of the treated individual (104). Thus, the GAL4-induced derepressed state may be propagated through meiosis, fertilization, and many additional mitotic divisions during development. Two alternative explanations have been proposed. The active state may be maintained by trxG proteins (104). Alternatively, this may reflect disruption of the normally silent state of germline chromatin by the pulse of GAL4 expression in germline as well as somatic cells, which would then permit the zygote to receive the transgene in the exceptional depressed state (104, 111).

### 3. Nonhomeotic functions of PcG proteins

*Drosophila* PcG proteins are involved in the regulation of many genes other than the homeotic genes, although these interactions are not as well characterized. Mammalian PcG proteins are required for proliferation and activation of **hematopoietic** cells (9, 73, 119, 120). A plant homologue of E(Z) (*medea*) restricts cell proliferation during embryogenesis (65), and worm homologues of E(Z) and ESC are involved in germline-specific silencing (35, 36). Thus, PcG proteins appear to perform an evolutionarily ancient transcriptional silencing function that has been adapted for the regulation of different genes and/or other developmental purposes in phylogenetically dispersed organisms.

### Bibliography

1. G. Jurgens (1985) *Nature* **316**, 153–155.
2. H. L. Landecker, D. A. R. Sinclair, and H. W. Brock (1994) *Dev. Genet.* **15**, 425–434.
3. R. G. A. B. Sewalt, J. van der Vlag, M. J. Gunster, K. M. Hamer, J. L. den Blaauwen, D. P. E. Satijn, T. Hendrix, R. van Driel, and A. P. Otte (1998) *Mol. Cell. Biol.* **18**, 3586–3595.
4. H. Strutt and R. Paro (1997) *Molec. Cell. Biol.* **17**, 6773–6783.
5. M. van Lohuizen, M. Tijms, J. W. Voncken, A. Schumacher, T. Magnuson, and E. Wientjens (1998) *Mol. Cell. Biol.* **18**, 3572–3579.
6. I. Duncan and E. B. Lewis (1982), *Developmental Order: Its Origin and Regulation*, Alan R. Liss, New York, pp. 533–554.
7. A. Busturia and G. Morata (1988) *Development* **104**, 713–720.
8. E. B. Lewis (1978) *Nature* **276**, 565–570.
9. N. Core, S. Bel, S. J. Gaunt, M. Aurrand-Lions, J. Pearce, A. Fisher, and M. Djabali (1997) *Development* **124**, 721–729.
10. J. Muller, S. Gaunt, and P. A. Lawrence (1995) *Development* **121**, 2847–2852.
11. B. Zink and R. Paro (1989) *Nature* **337**, 468–471.
12. A. Lonie, R. D'Andrea, R. Paro, and R. Saint (1994) *Development* **120**, 2629–2636.
13. R. Paro and D.S. Hogness (1991) *Proc. Natl. Acad. Sci. USA* **88**, 263–267.
14. J. C. Eissenberg, T. C. James, D. M. Foster-Hartnett, T. Hartnett, V. Ngan, and S. C. R. Elgin (1990) *Proc. Natl. Acad. Sci. USA* **87**, 9923–9927.
15. S. Messmer, A. Franke, and R. Paro (1992) *Genes & Devel.* **6**, 1241–1254.
16. A. Franke, S. Messmer, and R. Paro (1995) *Chromosome Res.* **3**, 351–360.
17. J. S. Platero, E. J. Sharp, P. N. Adler, and J. C. Eissenberg (1996) *Chromosoma* **104**, 393–404.
18. I. G. Cowell and C. A. Austin (1997) *Biochim. Biophys. Acta* **1337**, 198–206.

19. M. Kyba and H. W. Brock, H. W. (1998) *Mol. Cell. Biol.* **18**, 2712–2720.
20. M. J. Reijnen, K. H. Hamer, J. L. den Blaauwen, C. Lambrechts, I. Schoneveld, R. van Driel, and A. P. Otte (1995) *Mech. Dev.* **53**, 35–46.
21. A. Franke, M. DeCamillis, D. Zink, N. Cheng, H. W. Brock, and R. Paro (1992) *EMBO J.* **11**, 2941–2950.
22. B. P. Brunk, E. C. Martin, and P. N. Adler (1991) *Nature* **353**, 351–353.
23. M. van Lohuizen, M. Frasch, E. Wientjens, and A. Berns (1991) *Nature* **353**, 353–355.
24. M. J. Alkema, M. Bronk, E. Verhoeven, A. Otte, L. J. van't Veer, A. Berns, and M. van Lohuizen (1997) *Genes & Dev.*, **11**, 226–240.
25. M. Nomura, Y. Takihara, and K. Shimada (1994) *Differentiation* **57**, 39–50.
26. J. Simon, D. Bornemann, K. Lunde, and C. Schwartz (1995) *Mech. Dev.* **53**, 197–208.
27. G. Struhl and D. Brower (1982) *Cell* **31**, 285–292.
28. T. Gutjahr, E. Frei, C. Spicer, S. Baumgartner, R. A. H. White, and M. Noll (1995) *EMBO J.* **14**, 4296–4306.
29. S. S. Sathe and P. J. Harte (1995) *Mech. Dev.* **52**, 77–87.
30. J. Ng, R. Li, K. Morgan, and J. Simon (1997) *Mol. Cell. Biol.* **17**, 6663–6672.
31. E. J. Neer, C. J. Schmidt, R. Nambudripad, and T. F. Smith (1994) *Nature*, **371**, 297–300.
32. C. A. Jones, J. Ng, A. J. Peterson, K. Morgan, J. Simon, and R. S. Jones (1998) *Mol. Cell. Biol.* **18**, 2825–2834.
33. A. Schumacher, C. Faust, and T. Magnuson (1996) *Nature* **383**, 250–253.
34. W. G. Kelly and A. Fire (1998) *Development* **125**, 2451–2456.
35. R. Holdeman, S. Nehrt, and S. Strome (1998) *Development* **125**, 2457–2467.
36. I. Korf, Y. Fan, and S. Strome (1998) *Development* **125**, 2469–2478.
37. J.-M. Dura, J. Deatrick, N. B. Randsholt, H. W. Brock, and P. Santamaria (1988) *Roux's Arch. Devl. Biol.* **197**, 239–246.
38. M. DeCamillis, N. Cheng, D. Pierre, and H. W. Brock (1992) *Genes and Dev.* **6**, 223–232.
39. J. W. Hodgson, N. N. Cheng, D. A. Sinclair, M. Kyba, N. B. Randsholt, and H. W. Brock (1997) *Mech. Dev.* **66**, 69–81.
40. M. J. Gunster, D. P. E. Satijin, K. M. Hamer, J. L. den Blaauwen, D. de Bruijn, M. J. Alkema, M. van Lohuizen, R. van Driel, and A. P. Otte (1997) *Mol. Cell. Biol.* **17**, 2326–2335.
41. C. Ponting (1995) *Protein Sci.* **4**, 1928–1930.
42. D. Bornemann, E. Miller, and J. Simon (1996) *Development* **122**, 16212–1630.
43. A. J. Peterson, M. Kyba, D. Bornemann, K. Morgan, H. W. Brock, and J. Simon (1997) *Mol. Cell. Biol.*, **17**, 6683–6692.
44. W.-E. Kalisch and B. Rasmuson (1974) *Hereditas* **78**, 97–104.
45. R. S. Jones and W. M. Gelbart (1990) *Genetics* **126**, 185–199.
46. C.-t. Wu, R. S. Jones, P. F. Lasko, and W. M. Gelbart (1989) *Mol. Gen. Genet.* **218**, 559–564.
47. M. Gatti and B. S. Baker (1989) *Genes & Devel.* **3**, 438–453.
48. M. D. Phillips and A. Shearn (1990) *Genetics* **125**, 91–101.
49. E. C. Carrington and R. S. Jones (1996) *Development* **122**, 4073–4083.
50. R. S. Jones and W. M. Gelbart (1993) *Mol. Cell. Biol.* **13**, 6357–6366.
51. B. Tschiersch, A. Hofmann, V. Krauss, R. Dorn, G. Korge, and G. Reuter (1994) *EMBO J.* **13**, 3822–3831.
52. N. Tripoulas, D. Lajeunesse, J. Gildea, and A. Shearn (1996) *Genetics* **143**, 913–928.
53. G. Laible, A. Wolf, R. Dorn, G. Reuter, C. Nislow, A. Lebersorger, D. Popkin, L. Pillus, and T. Jenuwe in (1997) *EMBO J.* **16**, 3219–3232.

54. M. J. Stassen, D. Bailey, S. Nelson, V. Chinwalla, and P. J. Harte (1995) *Mech. Dev.* **52**, 209–223.
55. X. Cui, I. De Vivo, R. Slany, A. Miyamoto, R. Firestein, and M. L. Cleary (1998) *Nature Genet.* **18**, 331–337.
56. A. Shearn, G. Hersperger, and E. Hersperger (1978a) *Genetics* **89**, 341–353.
57. D. Moazed and P. H. O'Farrell (1992) *Development* **116**, 805–810.
58. F. Pelegri and R. Lehmann (1994) *Genetics* **136**, 1341–1353.
59. D. LaJeunesse and A. Shearn (1996) *Development* **122**, 2189–2197.
60. O. Hobert, B. Jallal, and A. Ullrich (1996) *Mol. Cell. Biol.* **16**, in press.
61. K. J. Abel, L. C. Brody, J. M. Valdes, M. R. Erdos, D. R. McKinley, L. H. Castilla, S. D. Merajver, F. J. Couch, L. S. Friedman, E. A. Ostermeyer, E. D. Lynch, M.-C. King, P. L. Welsh, S. Osborne-Lawrence, M. Spillman, A. M. Bowcock, F. S. Collins, and B. L. Weber (1996) *Genomics* **37**, 161–171.
62. H. Chen, C. Rossier, and S. E. Antonarakis (1996) *Genomics* **38**, 30–37.
63. O. Hobert, I. Sures, T. Ciossek, M. Fuchs, and A. Ullrich (1996) *Mech. Dev.* **55**, 171–184.
64. J. Goodrich, P. Puangsomlee, M. Martin, D. Long, E. M. Meyerowitz, and G. Coupland (1997) *Nature* **386**, 44–51.
65. U. Grossniklaus, J.-P. Vielle-Calzada, M. A. Hoepfner, and W. B. Gagliano (1998) *Science* **280**, 446–450.
66. C.-t. Wu and M. Howe (1995) *Genetics* **140**, 139–181.
67. P. N. Adler, E. C. Martin, J. Charlton, and K. Jones (1991) *Dev. Genet.* **12**, 349–361.
68. J. Simon, A. Chiang, and W. Bender (1992) *Development* **114**, 493–505.
69. Y. Haupt, W. S. Alexander, G. Barri, S. P. Klinken, and J. M. Adams (1991) *Cell* **65**, 753–763.
70. M. van Lohuizen, S. Verbeek, and B. Scheijen (1991) *Cell* **65**, 737–752.
71. M. Kanno, M. Hasegawa, A. Ishida, K. Isono, and M. Taniguchi, M. (1995) *EMBO J.* **14**, 5672–5678.
72. T. Akasaka, M. Kanno, R. Balling, M. Antonio Mieza, M. Taniguchi, and H. Koseki (1996) *Development* **122**, 1513–1522.
73. N. M. T. van der Lugt, J. Domen, K. Linders, M. van Roon, E. Robanus-Maandag, H. T. Riele, M. van der Valk, J. Deschamps, M. Sofroniew, M. van Lohuizen, and A. Berns (1994) *Genes Dev.* **8**, 757–769.
74. M. J. Alkema, N. M. T. van der Lugt, R. C. Bobeldijk, A. Berns, and M. van Lohuizen (1995) *Nature* **374**, 724–727.
75. E. C. Martin and P. N. Adler (1993) *Development* **117**, 641–655.
76. L. Rastelli, C. S. Chan, and V. Pirrotta (1993) *EMBO J.* **12**, 1513–1522.
77. M. Tagawa, T. Sakamoto, K. Shigemoto, H. Matsubara, Y. Tamura, T. Ito, I. Nakamura, A. Okitsu, K. Imai, and M. Taniguchi (1990) *J. Biol. Chem.* **265**, 20021–20026.
78. I. Duncan (1982) *Genetics* **102**, 49–70.
79. R. Aasland, T. J. Gibson, and A. F. Stewart (1995) *Trends Biochem. Sci.* **20**, 56–59.
80. D. A. R. Sinclair, T. A. Milne, J. W. Hodgson, J. Shellard, C. A. Salinas, M. Kyba, F. Randazzo, and H. W. Brock (1998) *Development* **125**, 1207–1216.
81. D. A. R. Sinclair, N. J. Clegg, J. Antnchuk, T. A. Milne, K. Stankunas, C. Ruse, T. A. Grigliatti, J. A. Kassis, and H. W. Brock (1998) *Genetics* **148**, 211–220.
82. W. J. Gehring (1970) *Dros. Inform. Serv.* **45**, 103.
83. J. R. Girton and S. H. Jeon (1994) *Genetics* **161**, 393–407.
84. T. R. Breen and I. M. Duncan (1986) *Dev. Biol.* **118**, 442–456.
85. J. L. Brown, D. Mucci, M. Whiteley, M.-L. Dirksen, and J. A. Kassis (1998) *Molec. Cell* **1**,

1057–1064.

86. J. Mihaly, R. K. Mishra, and F. Karch (1998) *Molec. Cell* **1**, 1065–1066.
87. R. E. Rayle and M. M. Green (1968) *Genetica* **39**, 497–507.
88. Y. Yamamoto, F. Girard, B. Bello, M. Affolter, and W. J. Gehring (1997) *Development* **124**, 3385–3394.
89. D. S. Henderson, S. S. Banga, T. A. Grigliatti, and J. B. Boyd (1994) *EMBO J.* **13**, 1450–1459.
90. P. Santamaria and N. B. Randsholt (1995) *Mol. Gen. Genet.* **246**, 282–290.
91. F. Docquier, O. Saget, F. Forquignon, N. B. Randsholt, and P. Santamaria (1996) *Roux's Arch. Dev. Biol.* **205**, 203–214.
92. R. B. Campbell, D. A. R. Sinclair, M. Couling, and H. W. Brock (1995) *Mol. Gen. Genet.*, **246**, 291–300.
93. P. Ingham (1984) *Cell* **37**, 815–823.
94. T. Sato, M. A. Russell, and R. E. Denell (1983) *Genetics* **105**, 357–370.
95. T. Sato, P. H. Hayes, and R. E. Denell (1984) *Dev. Genet.* **4**, 185–198.
96. P. N. Adler, J. Charlton, and B. Brunk (1989) *Dev. Genet.* **10**, 249–260.
97. E. J. Sharp, N. S. Abramova, W. J. Park, and P. N. Adler (1997) *Chromosoma*, **106**, 70–80.
98. J. Simon, A. Chiang, W. Bender, M. J. Shimell, and M. O'Connor (1993) *Dev. Biol.*, **158**, 131–144.
99. B. Zink, Y. Engstrom, W. J. Gehring, and R. Paro (1991) *EMBO J.* **10**, 153–162.
100. A. Chiang, M. B. O'Connor, R. Paro, J. Simon, and W. Bender (1995) *Development* **121**, 1681–1689.
101. K. Hagstrom, M. Muller, and P. Schedl (1997) *Genetics* **146**, 1365–1380.
102. G. Farkas, J. Gausz, M. Galloni, G. Reuter, H. Gyurkovics, and F. Karch (1994) *Nature* **371**, 806–808.
103. S. Poux, C. Kostic, and V. Pirrotta (1996) *EMBO J.* **15**, 4713–4722.
104. G. Cavalli and R. Paro (1998) *Cell* **93**, 505–518.
105. G. Reuter and P. Spierer (1992) *Bioessays* **14**, 605–612.
106. C. S. Chan, L. Rastelli, and V. Pirrotta (1994) *EMBO J.* **13**, 2553–2564.
107. J. G. Grindhart and T. C. Kaufman (1995) *Genetics* **139**, 797–814.
108. D. Zink and R. Paro (1995) *EMBO J.* **14**, 5660–5671.
109. V. Orlando and R. Paro (1993) *Cell* **75**, 1187–1198.
110. K. McCall and W. Bender (1996) *EMBO J.* **15**, 569–580.
111. V. Pirrotta (1998) *Cell* **93**, 333–336.
112. C. A. Bunker and R. E. Kingston (1994) *Mol. Cell. Biol.* **14**, 1721–1732.
113. V. Pirrotta and L. Rastelli (1994) *BioEssays* **16**, 549–556.
114. J. Schlossherr, H. Eggert, R. Paro, S. Cremer, and R. S. Jack (1994) *Mol. Gen. Genet.* **243**, 453–462.
115. A. K. Csink and S. Henikoff (1996) *Nature* **381**, 529–531.
116. A. F. Dernburg, K. W. Broman, J. C. Fung, W. F. Marshall, J. Philips, D. A. Agard, and J. W. Sedat (1996) *Cell* **85**, 745–759.
117. P. Buchenau, J. Hodgson, H. Strutt, and D. J. Arndt-Jovin (1998) *J. Cell Biol.* **141**, 469–481.
118. A. Busturia, C. D. Wightman, and S. Sakonju (1997) *Development* **124**, 4343–4350.
119. T. Akasaka, K.-I. Tsuji, H. Kawahira, M. Kanno, K.-I. Harigaya, Y. Ebihara et al. (1997) *Immunity* **7**, 135–146.
120. Y. Takihara, D. Tomotsune, M. Shirai, Y. Katoh-Fukui, K. Nishii, A. Moraleb et al. (1997) *Development* **124**, 3673–3682.



### Suggestions for Further Reading

121. V. Pirrotta (1998) Polycomb the genome: PcG, trxG, and chromatin silencing, *Cell* **93**, 333–336. (A mini-review that discusses the implications of recent reports for the possible mechanisms of PcG-dependent silencing.)
122. M. van Lohuizen (1998) Functional analysis of mouse Polycomb group genes, *Cell Mol. Life Sci.* **54**, 71–79. (A discussion of the pleiotropic functions of mammalian PcG genes with an extensive list of primary references.)

### Polyglycine

Polyglycine is the name given to the [polyamino acid](#) formed from [glycine](#). Like [proteins](#), the polyglycine [polymer](#) is formed by condensation of the [amino group](#) of one [amino acid](#) and the [carboxyl group](#) of another amino acid. On the other hand, proteins have a defined and complex sequence that includes all 20 of the amino acids, but polyglycine is a homopolymer formed only from glycine residues.

Glycine is the simplest amino acid residue, having a hydrogen atom as its [side chain](#). Furthermore, unlike the other amino acids, glycine is not **enantiomeric**. The relatively small size of the glycine side chain confers greater conformational flexibility to the polyglycine backbone compared with the backbone of other amino acid polymers. The polyglycine backbone adopts primarily two **conformations**, called poly(Gly) I and poly(Gly) II. Poly(Gly) I is an antiparallel rippled (rather than pleated) **b-sheet** (1). Poly(Gly) II is a helical conformation with three residues per turn and no intrachain [hydrogen bonds](#) (2) and is similar to the poly(Pro) II conformation of [polyproline](#). The backbone of [collagen](#), a structural protein that has a distinctive repeating unit in its **polypeptide** sequence corresponding to Gly–X–Y, where X is often proline, adopts a conformation like that of poly(Gly) II or poly(Pro) II.

Polyglycine may also represent an unusual [post-translational modification](#); up to 34 glycyI units have been observed covalently bound to the g-carboxyl group of C-terminal **glutamic acid residues** in [tubulin](#) (3).

### Bibliography

1. B. Lotz (1974) *J. Mol. Biol.* **87**, 169–180.
2. F. H. C. Crick and A. Rich (1955) *Nature* **176**, 780–781.
3. V. Redeker et al. (1994) *Science* **266**, 1688–1691.

### Polylinker

A polylinker is a short **DNA** sequence containing two or more different sites for cleavage by [restriction enzymes](#). Polylinkers are introduced into **vectors** to make [cloning](#) easier by providing

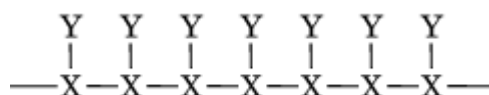
sites that allow cloning DNA, cut with any of a number of different restriction enzymes, into a single **plasmid**. They are useful for cloning a DNA fragment that has a different **cohesive end** at each end (ie, was produced by cutting out a DNA fragment from a molecule with two different restriction enzymes). When a plasmid cut with two different enzymes is **ligated** to a DNA fragment containing the two cohesive ends present in the plasmid and then the ligation mixture is **transformed** into **bacteria**, most of the transformants contain the desired recombinant plasmid in which one molecule of the fragment is precisely oriented in the plasmid. In contrast, when a DNA fragment is cloned into a plasmid cut with a single enzyme, so that both cohesive ends are the same, it recircularizes, and this produces a high background of transformants lacking an insert. This background is not present when the two ends of the fragment are different. Furthermore, the insert has a single orientation, whereas a fragment with identical cohesive ends inserts in either orientation with equal probability. Pairs of polylinkers are available that have the same set of sites, but in opposite orientation, to allow cloning the fragments in either orientation (i.e., pUC11 and pUC12). Polylinkers became somewhat less important once **PCR** became available because it is possible to introduce any desired restriction site into the PCR product by designing primers that incorporate the desired restriction sites.

## Polymer

A polymer is a large molecule, or [macromolecule](#), composed of many copies of repeating units joined together to form a long chain. Polymers can incorporate many thousands of atoms and have extremely high molecular weights. Such large molecules are ubiquitous in nature, and many are essential for life. Naturally occurring polymers include wool, cotton, wood, silk, and rubber. Of particular interest to the molecular biologist are **biopolymers** such as [proteins](#), **nucleic acids**, and carbohydrates. A diverse range of man-made polymers has been developed, including plastics, synthetic fibers, and paints.

The long chain, or [backbone](#), is the constant or repeating part of the polymer (Fig. 1). The [side chains](#), or functional groups, attached to the backbone can be variable. Homopolymers are polymers in which the repeating units are chemically and stereochemically identical. Both the backbone and side chains of homopolymers are constant throughout the length of the polymer. **Polyethylene glycol** and [polyacrylamide](#) are homopolymers commonly used in molecular biology. The repeating monomer units of these two reagents are ethylene glycol and acrylamide, respectively.

**Figure 1.** Schematic representation of a polymer, showing the repeating nature of the backbone (X). The side chains (Y) connected to the backbone are constant in homopolymers but variable in copolymers.



Heteropolymers, or copolymers, are formed from two or more kinds of monomer unit. The arrangement or order of the repeating units in copolymers can be either sequential or random. Proteins and nucleic acids are examples of biological copolymers (**biopolymers**) that have a specific sequence. The repeating units of proteins are the 20 different naturally occurring [amino acids](#), and the repeating units of nucleic acids are the [nucleotides](#).

[See also [Biopolymer](#), [Oligomer](#), and [Macromolecule](#).]

### Suggestions for Further Reading

- A. G. Walton and J. Blackwell (1973) *Biopolymers*, Academic Press, New York.  
L. Mandelkern (1983) *An Introduction to Macromolecules*, Springer-Verlag, New York.  
P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

## Polymerase Chain Reaction

The polymerase chain reaction (PCR) was originally described in 1986 (1) and since then has become one of the most frequently used and widely applied techniques in molecular biology. In its simplest form, two oligodeoxyribonucleotides (oligos) are used to select a specific target sequence from a mixture of DNA molecules and to act as primers for DNA extension using the target DNA as [template](#). Multiple rounds of template **denaturation**, primer **annealing**, and extension by **DNA polymerase** result in specific amplification of the target sequence. Because PCR is sufficiently sensitive to detect and amplify even one molecule of target DNA from a complex mixture of DNA molecules, it has proved a useful technique in many diverse fields. Various improvements and numerous variations on the original technique have enabled PCR to be used for a wide variety of applications, including [cloning](#) known and unknown sequences of DNA, [genome](#) and **genetic mapping**, **gene expression** studies, diagnostic purposes, analysis of ancient DNA, **sequencing**, and **mutational** analysis. The following is a survey of many of the common PCR techniques and applications.

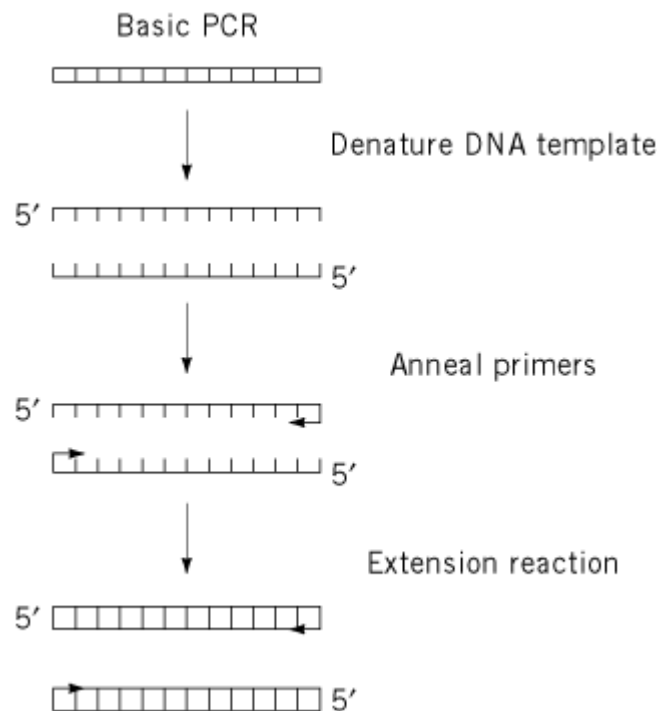
### 1. The Standard PCR Technique

In the standard PCR technique (Fig. 1), two oligos are designed to be complementary to the end sequences of the DNA to be amplified (the target sequence, or amplicon). These oligos, or PCR “primers,” are complementary to opposite strands of the double-stranded DNA and are generally 15 to 25 nucleotides in length. DNA synthesis primed from the two oligos results in selective amplification of the complementary strands of DNA contained between the two primers. Amplification occurs during repetitive cycles of synthesis that repeat the following three steps: (1) denaturation of the double-stranded DNA template; (2) annealing of the two primers to the single-stranded template; and (3) extension of the primers by a thermostable DNA polymerase. In theory, every cycle doubles the number of target sequences, although in practice PCR is not quite that efficient. The reaction mix includes the following: the DNA sample, two primers that “select” DNA to be amplified; the thermostable DNA polymerase to extend the primers; an excess of the four deoxynucleotides A, T, G, and C needed for DNA synthesis; and generally 10 mM Tris-HCl buffer, 50 mM KCl to provide the proper ionic strength, and  $Mg^{2+}$  as a cofactor (optimal  $Mg^{2+}$  concentrations may need to be determined for each reaction). Each step in the cycle is performed at a different temperature, and PCR machines are available to perform the temperature cycling automatically. Denaturation occurs at 94 to 96°C, and extension by polymerase is generally at 72°C. The temperature of the second step, primer annealing, is crucial and depends on the sequence of the primers, the concentration of the components of the PCR mix and the salt concentration. If the temperature is too high, not enough of the primers anneal to template DNA to have an efficient PCR, but a low annealing temperature allows nonspecific priming and results in spurious amplification products. The  $T_m$  (temperature when half the molecules are double-stranded and half are single-stranded) is often used for the annealing temperature and then adjusted up or down if necessary. If

the primer is 18 bases long or less, the  $T_m$  can be estimated by the following formula:

$$T_m = 4^\circ\text{C}(\text{number of GC base pairs}) + 2^\circ\text{C}(\text{number of AT pairs})$$

**Figure 1.** Basic PCR. The polymerase chain reaction uses multiple rounds of DNA denaturation, primer annealing, and primer extension to amplify a specific target sequence. First, double-stranded template DNA is denatured by brief incubation at 95°C. Secondly, primers that select the target sequence are annealed to single-stranded template DNA at a temperature determined by the length and sequence of the primers. Finally, the temperature is raised to 72°C to permit primer extension by a thermostable DNA polymerase, thus producing double-stranded target sequence. Each cycle of PCR repeats the same three steps, using the newly synthesized target DNA as additional PCR template. Thus, each round of PCR theoretically doubles the previous number of target sequences. Millions of copies of the target sequence are synthesized after multiple rounds of PCR. Horizontal arrows represent PCR primers. Double-stranded DNA is represented by horizontal lines joined by short vertical lines. The vertical lines represent the Watson–Crick base pairs.



### 1.1. Template

Standard PCR can be performed on any single-stranded or double-stranded DNA starting material: **genomic** DNA, **viral** DNA, cloned sequences, crude cell lysates, or ancient DNA, for example. Because the target sequence is amplified, only small amounts of starting material are needed (the amounts can vary from as little as a single molecule up to nanogram amounts of cloned DNA or microgram amounts of genomic DNA). Because of the sensitivity of PCR, even a small amount of contaminating DNA is amplified if it is complementary to the PCR primers. Therefore, it is important to perform proper controls and to use special precautions when contamination can interfere with interpretation of the data (see **Analysis of ancient DNA** below).

The DNA polymerase used in PCR is specific for DNA, but it is possible to analyze [messenger RNA](#) using PCR methods by incorporating a **reverse transcriptase** (RT) step to convert mRNA into double-stranded [complementary DNA](#), known as RT-PCR. The reverse transcription and PCR can be performed by a single enzyme, *rTth*, because it has both RT and polymerase activities.

### 1.2. Primers

There are various considerations when designing PCR primers. Although the 3' end of the primer, especially the final two bases, is important for specificity and should be complementary to the target sequence, it is possible to attach noncomplementary sequences to the 5' end, for example, [restriction enzyme](#) sites (to facilitate cloning of PCR amplicons) or **promoters**. GC pairs are more stable than AT pairs (GC pairs have three [hydrogen bonds](#) and AT pairs have two), and it is preferable to have either G's or C's at the 3' end of the primer and an overall GC content of 45 to 55%. Computer programs are available to help design the most efficient PCR primers for known sequences. Ideally the  $T_m$  of both primers should be the same within 1 to 2°C so that they both anneal specifically at the same annealing temperature. To prevent primer–primer binding and a decrease in efficiency of the amplification, primers must not be complementary to themselves or to each other.

Nested primers are used to increase the specificity of the reaction and to obtain larger amounts of amplified DNA. A second PCR is carried out using two “nested” primers that are complementary to sequences located 3' to the original primers, within the amplified DNA. The nested primers are used to reamplify a product, using a small amount of DNA from the original PCR product as a template.

For some experiments, it is necessary or more efficient to amplify a number of different products simultaneously in a single reaction, using different primers. This is termed *multiplex* PCR. For multiplex PCR, it is important to optimize primer choice and reaction conditions to eliminate spurious amplification products.

### 1.3. Polymerase

Various thermostable DNA polymerases are available that withstand the repeated high temperatures used for the denaturation step. *Thermus aquaticus* DNA polymerase (Taq) is commonly used and, although it has a low error rate, it is sufficiently accurate for most applications. If accuracy is crucial to the experiment, it may be necessary to use a polymerase having an inherent 3'–5' exonuclease (“proofreading”) activity, such as *vent* polymerase (2). For amplifying targets longer than 10 kb, extralong PCR (XL PCR) conditions must be used, and it is important to use a polymerase that has proofreading activity, because it is thought that errors result in termination of synthesis (3).

## 2. Cloning PCR Products

*Taq* polymerase frequently adds a single adenosine nucleotide to the 3' end of double-stranded DNA, thus making the products of a PCR reaction unclonable because of their “ragged” ends. These ends can be **filled-in** with T4 DNA polymerase and then **blunt-end ligated** into a vector. Alternatively, it is possible to clone the PCR products directly into special vectors having 3' T overhangs flanking the insertion site. A third possibility is to incorporate a restriction site into the primers that then is cut to yield **cohesive ends** for cloning.

## 3. Applications of the Polymerase Chain Reaction

### 3.1. Amplification of DNA: Cloning Substrate or Diagnostic Tool

PCR is useful because it allows identifying and isolating one or a few copies of a specific target sequence from a complex background. This permits rapid cloning of known sequences without the preparation of a **library** for gene isolation. PCR is also useful for screening for the presence of pathogens, even when only a small amount of sequence information is available for the pathogen (4). Production of amplified product using pathogen-specific primers signals the presence of the pathogen in the sample.

### 3.2. Manipulation of DNA

Because PCR primers can be designed to have a restriction site incorporated at the 5' end, it is possible to use PCR to amplify entire genes or portions of genes for **genetic engineering** purposes. This is helpful when restriction sites in the sequence of interest are too frequent or limiting and DNA must be moved from one vector to another or DNA sequences must be deleted or juxtaposed.

### 3.3. Gene Cloning with Degenerate Primers

PCR is used to clone genes even when the exact sequence of the gene of interest is not known, for example, when a gene sequence has been determined by [reverse translation](#) of a protein sequence. In this case, a mixture of oligos (a degenerate primer) is used that represents all of the possible coding sequences. Then the resulting amplified DNA is labeled and used as a **hybridization** probe for screening a [cDNA library](#) to isolate the actual coding sequence.

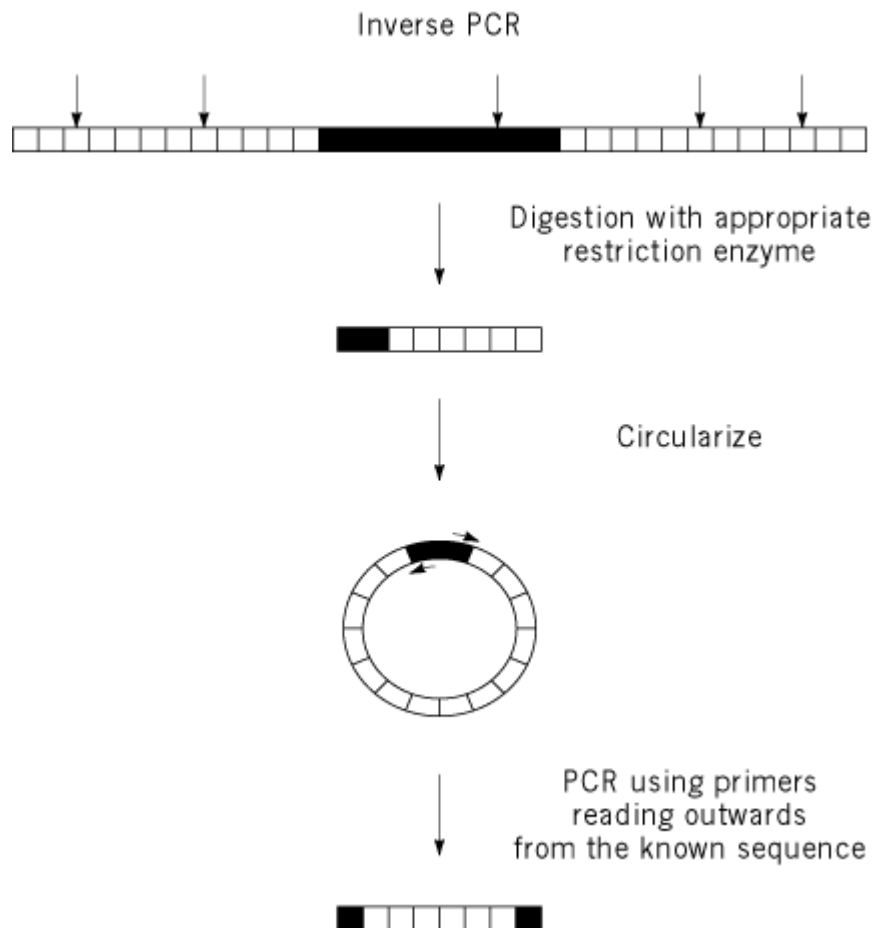
### 3.4. Recombination and Site-Directed Mutagenesis

Linear DNA molecules having homologous ends recombine *in vivo* in *Escherichia coli*, and two linear products that recombine to form a circle containing an **origin of replication** and a selectable marker are used to **transform** *E. coli*. Primers with appropriately homologous 5' ends are used to amplify the DNA sequence of interest from a plasmid or a genome and to amplify a cloning vector. If the linear molecules are used to transform *E. coli*, [recombination](#) between the two molecules results in the DNA being inserted into the vector ([5](#)). Specific [mutations](#) are introduced into insert DNA sequences by using mismatched primers to amplify the sequence of interest, but this method is useful only for creating mutations at the ends of sequences. Other methods have been developed to create mutations elsewhere in the amplified DNA, for example, mutagenesis by overlap extension ([6](#)). First, the DNA sequence of interest is amplified as two separate but overlapping fragments by using two different sets of primers. The region of overlap includes a mismatch contained in both of the two overlapping (but oppositely oriented) central primers. Then the two PCR products are mixed, denatured, and reannealed. The annealed products that have their overlap at the 3' end of each molecule act as primers extending each other to produce the full-length target DNA containing the new mutation. Using this method, a mutation is engineered anywhere in the target DNA.

### 3.5. Cloning Unknown Flanking Sequences: I-PCR and RACE

Inverse-PCR (IPCR) is often used to clone unknown DNA that flanks known sequences, such as DNA adjacent to **yeast artificial chromosome (YAC)** ends or the DNA flanking a [transposon](#) insertion (see [Transposon Tagging](#)). The starting DNA is digested with a restriction enzyme that cuts near the end of the known sequences and also at an undefined site in the flanking DNA (Fig. [2](#)). [Ligation](#) is performed under conditions that encourage intramolecular events, so that a circular DNA molecule is formed containing a small amount of known DNA adjacent to a small amount of unknown DNA. Two PCR primers are designed to prime outward from the ends of the known DNA, so that the unknown DNA between them is amplified and then can be cloned and used in further studies.

**Figure 2.** Inverse PCR. Inverse PCR is used to clone unknown DNA flanking known DNA, for example, genomic DNA flanking a transposon insertion. Total genomic DNA is isolated and digested with a restriction enzyme that releases a small piece of one transposon end and which cuts fairly frequently in the genomic DNA. Restriction enzyme digestion results in a DNA fragment containing the transposon end joined to a small piece of flanking DNA. The digested DNA is ligated under conditions that favor self-ligation, resulting in the formation of circular molecules. PCR is carried out using primers that read outward from the known transposon sequence into the unknown flanking DNA. Multiple rounds of PCR result in specific amplification of the unknown genomic DNA flanking the transposon. Vertical arrows represent restriction enzyme cleavage sites in the genomic DNA and the inserted transposon. The solid black rectangle represents the transposon insertion sequence. Horizontal lines joined by short vertical lines represent double-stranded genomic DNA. Horizontal arrows represent transposon-specific PCR primers.



The *rapid amplification of cDNA ends* (RACE) procedure is a method for obtaining full-length cDNAs from a known partial cDNA sequence, and it is used to obtain either the 3'- or 5'-ends of cDNAs. For isolation of 3'-ends, mRNA is reverse translated into cDNA using a primer (the RT primer) consisting of oligo(dT) followed by a unique sequence. Then PCR amplification is carried out using a primer specific to the known sequence and a primer (Primer U) specific to the unique sequence of the RT primer. Then nested PCR is performed to improve specificity. To clone 5' ends, reverse translation is carried out using a primer (Primer K) specific to the known cDNA sequence, reading toward the 5' end. A [poly A](#) tail is appended to the first strand cDNA, and the RT primer utilized for 3' RACE is used to synthesize the second strand of cDNA. Then PCR amplification is carried out using a gene-specific primer upstream of Primer K in conjunction with Primer U. Again this amplification is followed by nested PCR. 5'- and 3'-RACE are used to isolate full-length cDNAs in a few days, compared to the weeks needed to screen cDNA libraries and analyze library clones. This is done by coamplifying the target with an internal control and then comparing the ratio of target to control. The control is either an endogenous or exogenous gene or mRNA.

### 3.6. Genomic and Genetic Mapping Using PCR

Many different PCR-based mapping techniques have been described. PCR-based techniques are used to map sequences to specific segments of cloned DNA, such as yeast or bacterial artificial chromosomes (YACs or BACs), **cosmid** contigs, or segments of chromosomes (using radiation hybrid mapping techniques). All that is needed is a set of sequence-specific PCR primers. For example, a sequence-tagged-site (STS) consists of two PCR primers that specifically amplify the 3' untranslated region of a gene. In humans, YAC pools and **somatic cell** human-rodent hybrids are screened for the presence of the amplified STS PCR product to determine the STS map location (8). The presence of a PCR product indicates that the YAC or somatic cell hybrid contains the STS

sequence.

PCR-generated markers are also useful for creating genetic maps, rapidly producing markers for fine-mapping, and integrating genomic and genetic maps. PCR markers are used like any other **phenotypic** or DNA marker [eg, **restriction fragment length polymorphism** (RFLP) marker] to analyze differences among individuals in a population. DNA polymorphisms amplified by using repeat-specific primers [eg, [Alu sequences](#) in humans (9), **microsatellite** markers in plants (10)] have proved especially useful markers for mapping or for fingerprinting individuals. Arbitrary primers used to detect a large number of DNA differences between parental strains rapidly, and these random amplified polymorphic DNAs (RAPDs) are also used for mapping (11). RAPDs are detected between a pool of DNA from mutant progeny and a pool of DNA from wild-type progeny, where the progeny result from a genetic cross between two polymorphic strains, one of which carries the mutation. This allows rapid identification of a large number of RAPD markers in the region of the mutant gene and greatly facilitates cloning of genes located in regions of the genome that lack physical markers.

Because PCR markers are amplified by using sequence-specific primers, they can be used to position genes on both genetic and genome maps. For example, PCR methods have been used to map 47 [expressed sequence tag](#) (EST) sequences that are similar to plant disease-resistance genes (12). YACs containing the sequences were identified, and genetic map positions were obtained for those ESTs located on previously unmapped YACs. Several ESTs mapped to regions harboring genetically mapped disease-resistance loci, thus providing candidate disease-resistance genes.

### 3.7. Analysis of Gene Expression

A large number of PCR-based methods have been developed for analyzing gene expression. The sensitivity of PCR makes it especially useful in analyzing rare transcripts that cannot be analyzed by **Northern blotting** techniques. For known sequences, quantitative PCR is used to analyze relative levels of gene expression in different tissues or after different treatments.

Various PCR-based methods have been developed to identify and isolate differentially expressed genes (see [Subtractive Hybridization](#)). Two of the most commonly used procedures are representational difference analysis (RDA) (13) and differential display (14). RDA is used to select for genes expressed in only one mRNA population (the tester mRNA) compared to a second mRNA population (the driver). After cDNA synthesis and amplification of both populations, adapters are ligated only to the tester cDNA population (T-adapters). The tester and driver are mixed, denatured, and hybridized so that common sequences between the populations form tester–driver hybrids. Because of the excess of driver in the hybridization mix, only tester-specific sequences form tester–tester molecules. These are amplified using T-adapter-specific primers and used for further studies. RDA results in identifying a set of tissue- or treatment-specific cDNAs.

Differential display uses an arbitrary primer to amplify cDNAs obtained from different mRNA samples randomly. One primer (5'-T<sub>11</sub>NN, where NN are any two specific nucleotides) selects only cDNAs that have the nucleotides NN immediately adjacent to the poly A tail. When PCR is carried out using this primer in conjunction with a random 10-mer primer, the same subset of cDNAs is selectively amplified in each sample analyzed. PCR reactions from the different samples are run side by side on sequencing gels, so that gene expression differences can be visualized as bands present in one lane and absent in another. The bands of interest are cut out of the gel and the DNA eluted, cloned, and used for further analysis. This method is useful for analyzing many different tissues or treatments at once, but a large number of different primers is needed to survey for differences in all of the cDNAs in a sample.

PCR *in situ* methods have also been developed (see [In situ hybridization](#)). This method uses labeled PCR probes to amplify sequences *in situ* in cells or tissue samples. *In situ* PCR is useful for visualizing the location of rare transcripts that are difficult to visualize by normal *in situ*



hybridization methods.

### 3.8. Sequencing (see [DNA Sequencing](#))

PCR is used to amplify DNA before cloning and sequencing, or the amplified DNA is sequenced directly without prior cloning. The sequence results differ between the two strategies. If a cloned piece of DNA is sequenced, only one sequence is determined, even if there are several different sequences in the PCR product. Direct PCR sequencing, however, surveys the entire mixture of amplified sequences in the sample. Thus, for example, if the sequence being analyzed has two allelic forms, the sequence will be unambiguous at the positions where both alleles are identical but will represent both bases equally at the sites of difference. Double-stranded PCR products are sequenced directly, or single-stranded products are produced by asymmetric PCR. Several asymmetric PCR methods have been described. For example, PCR amplification is carried out with an excess of one primer over the other. Alternatively, the template is amplified first by using equal amounts of both primers, and then a second reaction is carried out using one primer in excess over the other. Thermal asymmetric PCR utilizes primers that have a  $T_m$  difference of 10°C to create an asymmetric PCR reaction. Double-stranded product is produced during 20 to 25 rounds of PCR using the lower  $T_m$  as the annealing temperature. Then subsequent rounds of PCR are carried out at the higher  $T_m$  to produce a single-stranded sequencing product.

### 3.9. Special Uses: Analysis of Ancient DNA and Forensics

The sensitivity of PCR (its ability to detect even one target molecule in a complex mixture) and the fact that PCR is used to amplify products from unclonable DNA samples (when the DNA has not been highly purified or when it is partially degraded) have made PCR a valuable tool when the starting material is limited or of poor quality. For these reasons, PCR has proven especially useful in analyzing ancient DNA isolated from, for example, museum specimens, slides, frozen mammoths, or insects in amber. In dealing with such variable and limited starting material, it is essential that special controls (eg, no template and control extract amplifications) and precautions (eg, multiple extracts from independent tissue samples, phylogenetic inferences) be taken to ensure the authenticity of the ancient DNA sequences analyzed.

PCR is used in forensic science for much the same reasons that it is applied to analyzing of ancient DNA. It enables the production of sufficient high-quality sample material for further analysis. Additionally, analysis using PCR-based markers is much more rapid than earlier methods, such as RFLP analysis, for matching samples (such as blood samples or sperm) to individuals.

## 4. Conclusion

This brief survey of various PCR methods and applications indicates the wide variety of uses for this simple technique and the large number of variations derived from the initial PCR procedure. Because of the sensitivity, flexibility, and simplicity of PCR, the uses and variations of the polymerase chain reaction will undoubtedly continue to multiply in the future.

## Bibliography

1. K. Mullis, F. Falcoma, S. Scharf, R. Snikl, G. Horn, and H. Erlich (1986) Cold Spring Harbor Symp. Quant. Biol. **51**, 260.
2. K. A. Eckert and T. A. Kunkel (1991) In *The Polymerase Chain Reaction I: A Practical Approach*. (M. J. McPherson, P. Quirke, and G. R. Tayler, eds.), IRL Press, Oxford, pp. 17–24.
3. S. Chang, C. Fockler, W. M. Barnes, R. Higuchi (1994) Proc. Natl. Acad. Sci. USA **91**, 5695–5699.
4. W. J. Martin (1994) In *The Polymerase Chain Reaction* (K. Mullis, F. Ferrè, and R. A. Gibbs, eds.), Birkhäuser, Boston, pp. 406–417.
5. D. H. Jones and B. H. Howard (1991) BioTechniques **10**, 62–66.

6. S. N. Ho, H. D. Hunt, R. M. Horton, J. K. Pullen, and L. R. Pease (1989) *Gene* **77**, 51–59.
7. F. Ferrè, A. Marchese, P. Pezzoli, S. Griffin, E. Buxton, and V. Boyer (1994) In *The Polymerase Chain Reaction* (K. Mullis, F. Ferrè, and R. A. Gibbs, eds.), Birkhäuser, Boston, pp. 67–88.
8. R. Berry et al. (1995) *Nature Genet.* **10**, 415–423.
9. M. Orita, T. Sekiya, and K. Hayashi (1990) *Genomics* **8**, 271–278.
10. S. R. McCouch, X. Chen, O. Panaud, S. Temnykh, Y. Xu, Y. G. Cho, N. Huang, T. Ishii, and M. Blair (1997) *Plant Mol. Biol.* **35**, 89–99.
11. M. McClelland, F. Mathieu-Daude, and J. Welsh (1995) *Trends Genet.* **11**, 242–246.
12. M. A. Botella, M. J. Coleman, D. E. Hughes, M. T. Nishimura, J. D. Jones, and S. C. Somerville (1997) *Plant J.* **12**, 1197–1211.
13. H. Hubank and D. G. Schatz (1994) *Nucleic Acids Res.* **22**, 5640–5648.
14. P. Liang and A. B. Pardee (1992) *Science* **257**, 967–971.
15. M. Höss, O. Handt, and S. Pääbo (1994) In *The Polymerase Chain Reaction* (K. Mullis, F. Ferrè, and R. A. Gibbs, eds.), Birkhäuser, Boston, pp. 257–264.
16. B. Budowle, A. Sajantila, M. N. Hochmeister, and C. T. Comey (1994) In *The Polymerase Chain Reaction* (K. Mullis, F. Ferrè, and R. A. Gibbs, eds.), Birkhäuser, Boston, pp. 244–256.

### Suggestions for Further Reading

17. K. Mullis, F. Ferrè, and R. A. Gibbs, eds. (1994) *The Polymerase Chain Reaction*, Birkhäuser, Boston. Various interesting articles on PCR use and methodology.
18. B. A. White, ed. (1997) *PCR Cloning Protocols, from Molecular Cloning to Genetic Engineering. Methods in Molecular Biology*, Vol. **67**, Humana Press, Clifton, NJ. Excellent methods manual with short descriptions of the techniques.
19. O. Bagasra and J. Hansen (1997) "In Situ" *PCR Techniques*, Wiley, New York.
20. R. Rapley, ed. (1996) *PCR Sequencing Protocols: Methods in Molecular Biology*, Vol. **65**, Humana Press, Clifton, NJ

## Polymerases

Polymerases represent a broad class of [enzymes](#) whose roles are central to **replicating** and expressing **genes**. This family of enzymes functions during [DNA replication](#) and RNA [transcription](#) to synthesize a **nucleic acid** from the genetic information encoded by the [template](#) strand. Polymerases are unique because these enzymes take direction from another molecule, a polynucleotide substrate. The nucleic acid template is either DNA or RNA, as is the strand synthesized, so polymerases can be divided into (1) DNA-dependent DNA polymerases; (2) RNA-dependent DNA polymerases; (3) DNA-dependent RNA polymerases; and (4) RNA-dependent RNA polymerases. Most polymerases involved in DNA replication and RNA transcription function in a **multienzyme complex**. During replication this complex (known as the replisome ) must be proficient in unwinding double-stranded DNA, initiating synthesis encoded by a single-stranded template (priming), incorporating deoxynucleoside triphosphates (synthesis), excising improperly templated nucleosides (**proofreading**), reducing DNA secondary structures on the template strand, unwinding the [DNA topology](#), and joining up the synthesized strand. Proteins involved at the **replication fork** include **helicase**, [primase](#), polymerase, 3'-5' exonuclease, [single-stranded DNA binding proteins](#), **topoisomerase**, and **ligase**. During gene transcription, the multiprotein complex

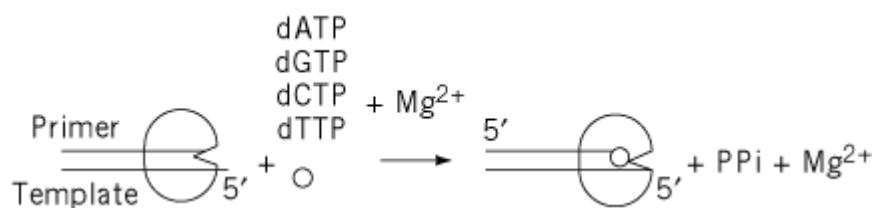
also contains components that have diverse function necessary for efficient initiation, elongation, and termination of RNA transcripts.

The structure of DNA discovered by Watson and Crick in 1953 contains complementary bases that interact via [hydrogen bonds](#). It suggested a mechanism of DNA replication in which each of the two strands of DNA is a template for the synthesizing a new DNA strand. This **semiconservative replication** mechanism involves generating two daughter duplexes, each containing one parent strand and one newly synthesized strand. DNA replication is an enzymatic process as discovered three years later by Arthur Kornberg and colleagues by experiments in which extracts of *Escherichia coli* were incubated with nucleic acid in the presence of **radiolabeled** deoxynucleoside triphosphates (dNTP). Almost fifty years after this discovery, we continue to search for the precise mechanism of nucleic acid polymerization and the enzymes involved during transcription and replication in various organisms.

## 1. Basic Enzymology

A general function of polymerases is to incorporate a nucleotide directed by a template onto a growing primer terminus. The generalized mechanism of DNA polymerization involves sequentially (1) binding the polymerase to its template primer; (2) binding the appropriate dNTP to the pol-DNA complex; (3) nucleophilically attacking the 3'-OH of the primer on the alpha phosphate of the dNTP, resulting in phosphodiester bond formation (Fig. 1); (4) releasing the pyrophosphate generated; and (5) translocating the polymerase toward the new 3'-OH of the primer. Polymerases catalyze these steps, and the products generated after each step are primers extended by one nucleoside in a 5'-3' direction and a pyrophosphate moiety. This phosphoryl transfer reaction is typically template-directed, such that adenine is incorporated opposite a template cytosine, etc. This template-directed reaction is highly accurate, and mistakes (such as misincorporation of thymine across from guanine) occur only once per  $5 \times 10^3$  to  $1 \times 10^6$  nucleotides polymerized, depending on the polymerase (1).

**Figure 1.** General mechanism for nucleotide incorporation by polymerases. Upon binding to a DNA or RNA template, polymerases incorporate complementary nucleotides. Selection of the proper nucleotide is governed by stable Watson-Crick hydrogen bond interactions, proper base stacking interactions, and the geometry of the polymerase active site, which can exclude improperly paired bases. DNA polymerases incorporate 2'-deoxynucleotides, whereas RNA polymerases incorporate nucleotides containing a 2'-OH group within the ribose. The product of one such phosphoryl transfer reaction is a DNA or RNA primer extended by one nucleotide and a pyrophosphate group, which is displaced from the polymerase active site.



## 2. Evolution

Because polymerases play a pivotal role in replicating and expressing genetic material, the amino acid sequences of the same polymerase from related organisms are conserved, and their sequence alignments are often used for **phylogenetic** purposes. Protein [sequence analysis](#) of diverse polymerases has identified three regions of homology, motifs A, B, and C (2). Structural and mutagenetic data have been used to attribute functional roles to these regions. It is believed that both motifs A and B bind to the base and sugar rings of the incoming dNTP, motif B also binds portions

of the template, and motif C binds to the divalent cation **cofactors** (Fig. 2). Other generalizations from sequence alignment studies include the following: (1) All classes of polymerases (excluding complex **RNA polymerases**) have motifs A and motif C in common, and both of these motifs contain an invariable **Asp** residue that binds a divalent metal cofactor (3). (2) Mammalian DNA polymerase  $\beta$  and **terminal deoxynucleotidyl transferase** represent unique polymerases, evolutionarily distinct from all other known polymerases (2, 4).

**Figure 2.** Amino acid sequence alignment of conserved motifs in polymerases. The amino acid residues are given with one-letter abbreviations. Bold letters represent amino acids that are identical across all families of polymerases. Positions generally occupied by hydrophobic residues are represented by “h,” and those typically containing hydrophilic residues are represented by “p.” Hyphens indicate nonconserved loci. The listed motifs are from Ref. 2.

|                | A                     | B                          | C                  |
|----------------|-----------------------|----------------------------|--------------------|
| Pol I          | <b>DYSQIELRVL</b> AHL | Rp- <b>AKh-FGhhYG</b>      | VHDEhV             |
| Pol $\alpha$   | h <b>DF</b> -SLYPS    | <b>Kh</b> - -N <b>ShYG</b> | YGD <b>TDS</b> hFh |
| Pol $\beta$    | T <b>Dh</b> Lhp       |                            | <b>DhDh</b> LIT    |
| DNA-dep RNAPol | <b>DG-C-Gh</b> QHh    | <b>Kp-VMT-hYG</b>          | h <b>HDSh</b> -T   |
| RNA-dep DNAPol | <b>DSp</b> -AY        |                            | Yh <b>DDhhh</b>    |
| RNA-dep RNAPol | h <b>Dp</b>           |                            | <b>GDD</b> -hh     |

### 3. X-Ray Crystallographic Structures

The three-dimensional structures of all polymerases solved to date (ie, Klenow fragment, **HIV-1 reverse transcriptase**, **bacteriophage T7 RNA polymerase**, pol  $\beta$ , **Taq DNA polymerase**, *Bacillus* DNA polymerase and bacteriophage RB69 DNA polymerase) morphologically resemble a “cupped human right hand.” In the HIV-1 reverse transcriptase, DNA binds across the palm subdomain, the finger subdomain interacts with single-strand portions of the template, and the thumb subdomain “grips” the duplex portion of the DNA (5). In each of these structures, the polymerase **active site**, where the incoming dNTP binds, is located in the palm subdomain (6). The active sites of all of these polymerases, including a  $\beta$  hairpin loop region where a conserved Asp residue is located, are superimposable. These polymerases have a **catalytic triad** composed of either three Asp residues or two Asp and one **Glu**, which function initially to bind the phosphate of the incoming dNTP/divalent cation complex and subsequently to assist the phosphoryl transfer catalysis (7). These diverse polymerases preserve similar architectures suggesting that they probably also exhibit similar mechanisms.

### 4. Properties

Among the most important properties of polymerases that guarantee the accurate replication of DNA or RNA are (1) fidelity (a measure of the accuracy of the polymerase) and (2) processivity (a measure of how many nucleotides are incorporated before the substrate dissociates). Polymerases exhibit varying degrees of fidelity, ranging from one misincorporation per 5000 to one per  $10^7$  nucleotides polymerized (8, 9). Those that incorporate the proper templated nucleotide at high efficiency are termed high-fidelity enzymes, and those that frequently misinsert a nucleotide are termed low-fidelity. Several polymerases contain a 3'-5' exonuclease subdomain (ie, a **proofreading** subunit) which increases the fidelity of the enzyme by approximately 10- to 100-fold.

The fidelity of polymerases is determined by one of several procedures. Fidelity of DNA synthesis was initially measured by utilizing polynucleotide templates consisting of only one or two types of nucleotides, such as an alternating poly d(A-T) template, and measuring the extent of misincorporation of radioactive cytosine or guanine nucleotides (10, 11). Greater sensitivity has been obtained with biological reversion assays, in which misincorporation by DNA polymerase results in the converting an **amber** mutation (ie, [stop codon](#)) in a **plasmid** into one that encodes an active, full-length protein (12). The forward mutational assays developed more recently offer the additional advantage of determining the mutational spectrum, that is, the types of misincorporated nucleotides catalyzed by the polymerase (13). *LacZ* has been most extensively utilized in these forward mutational assays as a [reporter gene](#) for studies on the mutational spectrum of DNA polymerases. Upon transformation of the copied plasmid (which encodes the *LacZ* gene) into *E. coli* and plating the transfected bacteria in the presence of **X-gal** (which is converted to a blue staining metabolite by the protein encoded by the *LacZ* gene, ***β*-galactosidase**), the fidelity is determined simply by counting the number of blue and white colonies resulting from functional (or nonmutated) or nonfunctional (or mutated) *LacZ* gene, respectively (see **Blue-white selection**). Sequencing the *LacZ* gene mutants determines the mutational spectrum. The fidelity of incorporation is also determined kinetically by comparing the ratio  $k_{cat}/K_m$  of the incorrect nucleotide to that of the correct nucleotide, this ratio directly reflects the efficiency of nucleotide incorporation (14). As a second step, the same assay measures the fidelity of extension by using primers that terminate in a noncomplementary nucleotide and measuring the incorporation of complementary nucleotides onto the end of this primer (15).

Processivity refers to the number of nucleotides incorporated per binding event of the polymerase with the template-primer complex. The processivity values of different polymerases range from one nucleotide to about ten thousand. The processivities of several polymerases involved in genomic replication are enhanced upon binding to a second protein, termed the processivity factor. For example, to fulfill their roles efficiently during DNA replication in **eukaryotes**, DNA polymerases  $\delta$  and  $\epsilon$  associate with a homotrimer that has 36-kDa subunits of proliferating cellular nuclear antigen (PCNA) which form a “**sliding clamp**” (16). Phage T4 gene 45 protein (17) and *E. coli* beta (18) similarly augment the processivity of T4 DNA pol and pol III, respectively, by acting as “sliding clamps” bound to the polymerase, thus preventing its dissociation from DNA.

## 5. General mechanism

The kinetic pathway, including the rates of association and dissociation of DNA polymerase with the template-primer and with dNTP, has been established for T7 DNA polymerase, the Klenow fragment of *E. coli* polymerase I (19), pol  $\beta$ , and HIV reverse transcriptase (20). DNA polymerases bind DNA tightly and have **dissociation constants** ( $K_d$ ) ranging from 1 to 50 nM. The apparent  $K_m$  (Michaelis constant) of the incoming dNTP is approximately 5 to 50  $\mu$ M. Subsequent to dNTP binding, polymerases undergo a conformational change that brings the alpha phosphate of the incoming nucleotide into proximity of the 3'-OH of the primer. It is thought that this conformational change enhances the discrimination between properly paired and mispaired nucleotides, thus ensuring that even relatively error-prone polymerases, such as HIV reverse transcriptase, incorporate the correct nucleotide 99.98% of the time. Polymerases that have higher fidelity presumably contain a more constrained active site during the conformational change step. The **rate-limiting step** during the incorporation of a single nucleotide in **steady-state** reactions is polymerase dissociation from the DNA, that is, the rate of dissociation is slow (20).

## 6. *In Vivo* Roles of Various Polymerases

DNA-dependent DNA polymerases are primarily involved in DNA replication and [DNA repair](#) in organisms that contain a DNA genome (e.g., **eukaryotes**, **prokaryotes**, bacteriophages including T3, T4, T7, FX 174, and DNA **viruses**, including **Epstein–Barr**, [cytomegalovirus](#), and other **Herpes**-related viruses). In eukaryotes, it is believed that polymerases  $\delta$  and  $\epsilon$  conduct much of the

synthesis of the leading and lagging strands, and polymerase  $\alpha$  synthesizes the RNA primers and extends these primers by adding a limited number of nucleotide residues. Mutational analysis in **yeast** has shown that each of these three proteins is essential for growth and DNA replication.

In addition to DNA replication, DNA-dependent DNA polymerases are also involved in DNA repair. Specifically, these polymerases bridge gaps in DNA resulting from excising damaged or mispaired bases. It is thought that polymerases  $\delta$  and  $\epsilon$  are essential for filling large DNA gaps during [mismatch repair](#) (which leaves a gap of thousands of nucleotides), nucleotide [excision repair](#) (which leaves about a 30 nucleotide gap), and a subset of [base excision repair](#) pathways that leave large patches of 5 to 10 nucleotides (21). Polymerase  $\beta$  is highly efficient at filling gaps of six nucleotides or less, making it the ideal candidate for a subset of base excision repair that leaves short patches of less than 5 nucleotides (22). The newly discovered pol  $\zeta$  in yeast (the pol domain is the product of the Rev3 gene that associates tightly with the product of the Rev7 gene, forming the functioning Rev3-Rev7 complex) is efficient at bypassing bulky lesions (23). This polymerase might be involved in synthesizing DNA across damaged regions that have yet to be repaired.

In addition, eukaryotes contain a unique polymerase, **terminal deoxynucleotidyl transferase**, which extends a DNA primer in the absence of a template (24). Considering that terminal deoxynucleotidyl transferase is found only in **immune**-competent cells, it is a reasonable conclusion that this enzyme is involved in generating immunological diversity (25). *E. coli*, the most studied prokaryote, contains three polymerases: Pol I, II, and III (26). Pol III is the primary replicative polymerase for both leading and lagging strands, whereas Pol I is involved in filling gaps during base excision and nucleotide excision repair. The function of Pol II is unknown.

DNA-dependent RNA polymerases are generally involved in expressing genetic material (see [Transcription](#)). Eukaryotic cells contain four types of such RNA polymerases, which interact with numerous [transcription factors](#) during initiation, elongation, and termination of transcripts. Eukaryotic RNA pol I is located in the [nucleolus](#) at 40,000 copies per cell. It synthesizes the large (35 to 47 S) RNA of [ribosomes](#). RNA polymerase II is located in the nucleoplasm at 40,000 copies per cell and functions primarily to synthesize [messenger RNA](#) precursors and several types of **small nuclear RNA** (snRNA). RNA pol III is located in the nucleoplasm in 20,000 copies per cell and functions in the synthesis of [transfer RNA](#) and small (5 S and 7 S) ribosomal RNA. [Mitochondria](#) of eukaryotes also contain a mitochondrial RNA polymerase, a monomeric enzyme encoded by the nuclear genome that is responsible for transcribing mitochondrial genes.

Prokaryotic DNA-dependent **RNA polymerase** is typified by *E. coli* RNA polymerase, which is composed of six subunits of five unique proteins:  $\alpha$  (which functions as a backbone onto which the other proteins bind);  $\beta$  (assists in DNA binding);  $\beta'$  (contains part of the active site); and  $\sigma$  (contributes to initiation). The  $\sigma$  subunit dissociates after initiation of the transcript, leaving behind the core polymerase. Bacteriophage DNA-dependent RNA polymerases are unique because they are single-subunit polymerases that do not require transcription factors. Some examples of these polymerases include T7, T3, and SP6 RNA polymerases. These polymerases are highly processive (>1000 nucleotides synthesized per binding event) and are used widely in laboratories to generate RNA transcripts *in vitro*.

RNA-dependent DNA polymerases and RNA-dependent RNA polymerases are generally encoded by **viruses** and **transposons**. RNA-dependent DNA polymerase (which also includes **reverse transcriptase**) is encoded by **retroviruses**, **LINEs**, **copia elements**, retroelements, and **hepatitis B virus**. It functions in the complicated life cycles of these organisms to convert their RNA genomes to DNA.

RNA-dependent RNA polymerases are encoded primarily by RNA viruses, such as **polio**, **hepatitis virus type C**, [rhinovirus](#), and Colxsackie. They function to replicate the viral genome and potentially to transcribe viral genes.

## Bibliography

1. T. A. Kunkel and L. A. Loeb (1981) *Science* **213**, 765–767.
2. M. Delarue, O. Poch, N. Tordo, D. Moras, and P. Argos (1990) *Protein Eng.* **3**, 461–467.
3. L. Blanco, A. Bernad, M. A. Blasco, and M. Salas (1991) *Gene* **100**, 27–38.
4. P. H. Patel, A. Jacobo-Molina, J. Ding, C. Tantillo, A. D. Clark, Jr., R. Raag, R. G. Nanni, S. H. Hughes, and E. Arnold (1995) *Biochemistry* **34**, 5351–5363.
5. A. Jacobo-Molina, J. Ding, R. G. Nanni, A. D. Clark, Jr., X. Lu, C. Tantillo, R. L. Williams, G. Kamer, A. L. Ferris, P. Clark, A. Hizi, S. H. Hughes, and E. Arnold (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6320–6324.
6. H. Pelletier, M. R. Sawaya, A. Kumar, S. H. Wilson, and J. Kraut (1994) *Science* **264**, 1891–1903.
7. G. P. Mullen, E. H. Serpersu, L. J. Ferrin, L. A. Loeb, and A. S. Mildvan (1990) *J. Biol. Chem.* **265**, 14327–14334.
8. T. A. Kunkel and L. A. Loeb (1979) *J. Biol. Chem.* **254**, 5718–5725.
9. M. Fry and L. A. Loeb (1986) *Animal Cell DNA Polymerases*, CRC Press, Boca Raton, FL.
10. J. F. Speyer (1965) *Biochem. Biophys. Res. Commun.* **21**, 6–8.
11. N. Battula and L. A. Loeb (1974) *J. Biol. Chem.* **249**, 4086–4093.
12. L. A. Weymouth and L. A. Loeb (1978) *Proc. Natl. Acad. Sci. USA* **75**, 1924–1928.
13. T. A. Kunkel (1985) *J. Biol. Chem.* **280**, 5787–5796.
14. M. S. Boosalis, J. Petruska, and M. F. Goodman (1987) *J. Biol. Chem.* **262**, 14689–14699.
15. F. W. Perrino and L. A. Loeb (1989) *J. Biol. Chem.* **264**, 2898–2905.
16. T. S. Krishna, X. P. Kong, S. Gary, P. M. Burgers, and J. Kuriyan (1994) *Cell* **79**, 1233–1243.
17. T. C. Jarvis, L. S. Paul, J. W. Hockensmith, and P. H. von-Hippel (1989) *J. Biol. Chem.* **1989**, 12717–12719.
18. X. P. Kong, R. Onrust, M. O'Donnell, and J. Kuriyan (1992) *Cell* **69**, 425–437.
19. R. D. Kuchta, V. Mizrahi, P. A. Benkovic, K. A. Johnson, and S. J. Benkovic (1987) *Biochemistry* **26**, 8410–8417.
20. K. A. Johnson (1993) *Annu. Rev. Biochem.* **62**, 685–713.
21. A. Blank, B. Kim, and L. A. Loeb (1994) *Proc. Natl. Acad. Sci. USA* **91**, 9067–9051.
22. R. W. Sobol, J. K. Horton, R. Kuhn, H. Gu, R. K. Singhal, P. K. Rajewsky, and S. H. Wilson (1996) *Nature* **379**, 183–186.
23. J. R. Nelson, C. W. Lawrence, and D. C. Hinkle (1996) *Science* **272**, 1646–1649.
24. C. Penit, M. J. Gelibert, C. Transy, and P. Rouget (1982) *Adv. Exp. Med. Biol.* **145**, 61–73.
25. T. A. Kunkel, K. P. Gopinathan, D. K. Dube, E. T. Snow, and L. A. Loeb (1986) *Proc. Natl. Acad. Sci. USA* **83**, 1867–1871.
26. A. Kornberg and T. Baker (1992) *DNA Replication*, Freeman, New York.

### Suggestion for Further Reading

27. A. Kornberg and T. A. Baker (1992) *DNA Replication*, Freeman, New York.

## Polynucleotide Phosphorylase

Polynucleotide phosphorylase (PNPase, polyribonucleotide: orthophosphate nucleotidyltransferase, EC 2.7.7.8) was the first [enzyme](#) to be discovered that catalyzes the synthesis of polyribonucleotides with a 3',5'-phosphodiester bond [reviewed in (1, 2)]. In the forward reaction, long polyribonucleotides,  $[{}_pN]_n$ , are synthesized from various nucleoside diphosphates ( $ppN$  or NDP), and inorganic phosphate,  $P_i$  is eliminated. In the reverse reaction, PNPase is a processive 3' → 5' exoribonuclease that catalyzes the phosphorolysis of polyribonucleotides, liberating nucleoside diphosphates:



In a processive reaction, the enzyme catalyzes the reaction completely on one RNA molecule before commencing on another, so that the substrate does not leave the enzyme during the reactions.

PNPase was discovered by Grunberg-Manago and Ochoa during a study of the mechanism of biological phosphorylation in *Azotobacter vinelandii* (*A. agilis*) (3-5). Studies of the nature of ribonucleotide incorporation into nucleic acids led to the recognition of the same enzyme in *Escherichia coli* cell-free extracts (6, 7). The enzyme was also isolated from *Micrococcus luteus* (formerly *M. lysodeikticus*) (8) and subsequent work has shown that it is widely distributed among **bacteria**.

PNPase has been the subject of many studies. It was employed as a tool for synthesizing model nucleic acids and for solving many important biological problems. Thus, establishing the [genetic code](#) was facilitated by the ability of PNPase to catalyze the synthesis of heteropolymers and triplet nucleotides. The advances made in understanding the physicochemical properties of **polynucleotide chains** and their **hybridization** reactions and the synthesis of polynucleotide **inducers** of [interferon](#) are additional examples of the role played by the enzyme. Renewed interest in PNPase arose with the finding that the enzyme is associated in a multiprotein complex that is involved in [messenger RNA](#) (mRNA) degradation, and it may serve other important physiological functions.

## 0.1. Properties

### 0.1.1. Distribution and Purification

PNPase is widely distributed among a variety of **aerobic**, **anaerobic**, **halophilic**, and **thermophilic** bacteria (2, 9). The properties of the enzyme differ somewhat in various bacterial species. About 10% of the total activity is found attached to washed *E. coli* [ribosomes](#), probably bound to mRNA (10), but [immunoelectron microscopy](#) indicates that most is located in the **cytoplasm** (11). Some enzyme activity is also found in membrane **vesicles** isolated from *E. coli* cells (12). In several other bacterial species, however, the enzyme is primarily found associated with the cell [membrane](#) (2, 13). PNPase was also isolated from **plants** (14, 15), where it is localized in the **chloroplasts** (16, 17). The enzyme is scanty in animal cells. Some activity was reported to be associated with the [endoplasmic reticulum](#) ribosomes (18, 19) although, the results could result from a combination of other enzymes (20). PNPase was purified from a wide variety of microorganisms (1, 2). Essentially homogenous preparations have been obtained from *E. coli* (21-25), *A. vinelandii*. (5, 26), *M. luteus* (27, 28), *V. costicola* (29), *S. antibioticus* (30) and many other bacterial species (2).

Most of the bacterial PNPase preparations are primer-independent forms, which catalyze *de novo* polymerization processively. Primer-dependent forms are derived from limited digestion by [trypsin](#) of PNPase from *M. luteus* (31). PNPase preparations from *Bacillus stearothermophilus* (21), *Thermus thermophilus* (33), and cucumber leaves (15) show primer dependency. More recently, amplification of the enzyme with [recombinant DNA](#) techniques in *E. coli* cells allowed its efficient purification and was used to raise polyclonal antibodies (11, 34).

### 0.1.2. Metal Ion Requirement and Properties of PNPase

Magnesium ions are required for the reactions catalyzed by PNPase but can be partially replaced by



Mn<sup>2+</sup> (1, 7, 22, 35) or other divalent cations, such as Co<sup>2+</sup>, Ni<sup>2+</sup>, Cd<sup>2+</sup>, Cu<sup>2+</sup>, and Zn<sup>2+</sup>, although with much lower efficiencies (2, 35, 36). Polymerization of GDP with *E. coli* PNPase, however, proceeds efficiently only in the presence of Mn<sup>2+</sup> at 60°C (37). A mutant PNPase from *E. coli* Q13 requires Mn<sup>2+</sup> rather than Mg<sup>2+</sup> (38). Purified PNPase is unstable above 55°C and is rapidly and irreversibly inactivated at 65°C (1, 22, 25, 39). The *M. luteus* (40) and the *C. perfringens* (41) enzymes are less stable. The *E. coli* PNPase is stabilized against heat inactivation by the presence of NDPs or oligoribonucleotides with a free 3'-OH terminus (39).

PNPase is sensitive to **proteolytic** digestion, which may account for differences in subunit structure and catalytic properties of PNPase from various sources [see (2)]. Thus, native PNPase from *M. luteus* is primer-independent, catalyzes *de novo* polymerization processively, and is only slightly stimulated by oligoribonucleotides. Limited trypsin digestion of the native enzyme alters its polymerization activity without affecting its ability to phosphorylate polynucleotides. The trypsinized enzyme (form T) catalyzes elongation of primers by a random mechanism and is stimulated by oligonucleotides (32, 42-44).

The *de novo* polymerization of NDPs, particularly at low Mg<sup>2+</sup> concentrations, is preceded by a lag period, which may be overcome by adding polynucleotides or short oligonucleotide primer molecules with a free 3'-hydroxyl group (21, 22, 45). The oligonucleotide primers are incorporated into the polymer synthesized by PNPase from *M. luteus* (46), and their effect is maximal in the polymerization of GDP, which proceeds with a slow rate in the absence of such primers (47). These primers also accelerate the NDP↔P<sub>i</sub> exchange reaction (22, 45, 48). The ADP↔P<sub>i</sub> exchange reaction is inhibited by poly(U), but not by poly(A), poly(G) or poly(C) (49). When blocked with a 3'-terminal phosphate moiety, oligonucleotides inhibit the polymerization reaction and exchange reactions (48, 50, 51). Several modified NDPs, such as 6-azauridine diphosphate, block the catalytic activity of PNPase (52). A heat-stable polypeptide activates the ADP↔P<sub>i</sub> exchange reaction (7, 22). Depending on the source of the enzyme, poly- L-lysine and other polyamines stimulate or inhibit enzyme activity (2). The actions of several other chemical agents that inhibit the catalytic activity of PNPase have been extensively summarized (1).

## 0.2. Structure

### 0.2.1. Molecular Weight

PNPase from *E. coli* is composed of three identical subunits (25, 53). Each contains 711 amino acids and has a molecular mass of 77,122 Da (54). A higher molecular weight of 84 to 88(±5)kDa is observed by **SDS-PAGE**, which is attributed to the acidic character of PNPase (25, 53). PNPase from *M. luteus* is also active as a homotrimer (25, 55). The molecular weight of the entire enzyme from *E. coli* has been determined by **sedimentation equilibrium**, **gel filtration**, **sucrose gradient centrifugation**, and **gel electrophoresis** under nondenaturing conditions. In the latter method, enzyme activity can be visualized after electrophoresis by incubating the gels in the presence of ADP and Mg<sup>2+</sup>, followed by staining the poly(A) formed *in situ* with **acridine** orange in the presence of lanthanum chloride (56). Other *in situ* methods for visualizing active enzyme molecules have also been published (2). The molecular weight of purified *E. coli* PNPase determined by sedimentation equilibrium ranges from 214 to 252(±20)kDa [see (2)]. The molecular weight and the physical and chemical properties of *M. luteus* enzyme are very similar (28). In crude extracts *E. coli* PNPase displays microheterogeneity, and higher level forms are detected, which may arise from the association of an additional polypeptide subunit (53). In some other bacterial species, PNPase shows considerable heterogeneity in the composition of its subunits, which may result from limited proteolysis of the enzyme. PNPase from cucumber leaves is a complex of three subunits, possibly not identical, of about 50 kDa each (15).

### 0.2.2. Complexes with Other Proteins

Sucrose gradient centrifugation and gel filtration detect the presence of PNPase in higher molecular weight complexes. These PNPase forms arise from association with enolase (53), an essential

glycolytic enzyme. Furthermore, it was found that PNPase is associated in a multicomponent complex called the degradosome, together with ribonuclease E (RNase E, an endoribonuclease). Association between the two enzymes requires the presence of an intact C-terminal region in RNase E, whereas this region is not required for the catalytic activity of RNase E (57). Several additional proteins are found in this complex, including a “**DEAD box**” RNA helicase and enolase (58-60), a **heat shock** protein, DnaK (61), and sometimes GroEL, a **chaperonin** that mediates protein folding and assembly (62). The average molar ratios of the proteins in the complex were estimated from silver-stained gels at 1.0: 0.9: 0.93: 1.8 for RNaseE: PNPase: RNA helicase: enolase (61). Only 5 to 10% of the cellular enolase is thought to be in the complex. Most of the PNPase and enolase in a cell are found in a binary complex, not in the degradosome (59). Chloroplast PNPase is also associated with an endonuclease that is homologous to RNase E (17). The physiological roles of this multiprotein complex is discussed in the section dealing with the function of PNPase.

### 0.2.3. The PNPase gene

The *pnp* gene for *E. coli* PNPase (54) is part of an operon together with the ribosomal protein S15 gene (*rpsO*). These two proteins are expressed differentially following processing of the *rpsO-pnp* transcript by RNase E and RNase III (63, 64). Two RNA-binding **domains** have been characterized in the C-terminal amino acid sequence of PNPase. The first domain, the KH module, is also found in many hnRNP proteins, in a number of other RNA-associated proteins from human and yeast cells, and in *E. coli* ribosomal protein S3 and the **transcription factor** NusA. The structure of the KH module determined by NMR consists of an antiparallel, three-stranded **b-sheet** connected by two **a-helical** regions (65). The second RNA-binding domain is strongly homologous to each of the four homologous stretches in the middle and C-terminal parts of ribosomal protein S1 (54). The structure of this S1 motif of PNPase consists of a five-stranded, antiparallel b-barrel (66). Conserved residues on one face of the barrel and adjacent loops form the putative RNA-binding site. The structure is very similar to that of **cold shock** proteins, suggesting that all these proteins are derived from an ancient nucleic acid-binding protein. The RNA-binding motifs of PNPase are believed to be involved in the processive nature of the phosphorolysis reaction. The three RNA-binding domains in trimeric PNPase might allow the enzyme to associate with adjoining regions of the same polynucleotide (54). This suggestion is supported by the observation that proteolytic cleavage of PNPase leads to 62-kDa subunits that still associate to form an active enzyme but have lost the ability to bind oligonucleotides more than 11 nucleotides long. Moreover, the processive nature of the reaction of the modified enzyme was also reduced (67).

The sequences of the *pnpA* genes from *E. coli* (54), *B. subtilis* (68), *H. influenzae* (69), and *Photobacterium luminescens* (70) indicate polypeptide chains of 705 to 711 residues of similar sequence. For example, those from *B. subtilis* and *E. coli* are 50% identical. That from *S. antibioticus* codes for an **open reading frame**, corresponding to 740 amino acids, and is 44% identical to the *E. coli* enzyme (30). There are two highly conserved domains near the C- and N-termini in these sequences, but overall homology is apparent along the entire molecule. Spinach chloroplast PNPase (17) is strikingly homologous (43% identity, 63% similarity) to *E. coli* PNPase (54). The chloroplast PNPase is longer than *E. coli* PNPase by at least 75 amino acids at the N-terminus, which probably represents a unique domain of the chloroplast protein. The content of acidic residues (14.2%), together with its low content of basic residues (12.65%), predicts an acidic protein with an isoelectric point of 5.07 for the *E. coli* enzyme, as compared to a value of 6.1 estimated by **isoelectric focusing** (25).

### 0.2.4. Reactions Catalyzed

PNPase catalyzes the reversible polymerization of NDPs to polyribonucleotides and inorganic phosphate (Eq. 1). Each of the four common NDPs can serve separately as a substrate for the polymerization reaction, leading to formation of homopolymers. Polymerization of a mixture of NDPs that contains different bases results in producing a random copolymer. The enzyme does not require a template and cannot copy one. The polymerization of GDP by *E. coli* PNPase is very slow (37), but copolymerization of GDP with the other NDPs is readily achieved (1). Polynucleotides are synthesized *de novo* in the 5' → 3' direction *via* a processive mechanism.

The mechanism of chain initiation and formation of the first internucleotide bond involves the reaction between two NDP molecules, one of which serves as the accepting 5'-terminus. Analysis of the newly synthesized polynucleotides reveals that they contain a monophosphate group at the 5'-terminus, rather than the expected 5'-pyrophosphate group (2, 55, 71, 72). The pyrophosphate linkage may be broken before or concomitant with the formation of the first internucleotide linkage, although other mechanisms cannot be ruled out (55). In addition, ApA and pApA do not undergo phosphorolysis but accumulate as resistant end-products of poly(A) degradation (73, 74). This implies that the initiation of *de novo* polymerization involves an initial irreversible step.

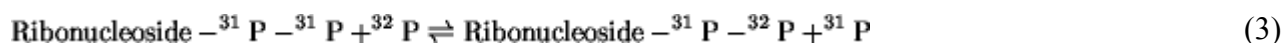
Under suitable conditions, the enzyme also catalyzes the elongation of a primer oligonucleotide with a free 3'-terminal hydroxyl group:



where R represents the oligonucleotide primer with at least two nucleoside residues and a free 3'-terminal hydroxyl group. Elongation of oligonucleotide primers may occur by a nonprocessive mechanism (1, 43). Thus, the trypsinized T form of PNPase from *M. luteus* has a reduced *de novo* polymerization activity but readily elongates the 3'-end of an oligonucleotide primer. The elongation is nonprocessive, and only small oligonucleotides are formed (43).

In the reverse reaction, PNPase serves as a 3' → 5' exoribonuclease and releases NDPs sequentially from the 3'-OH terminus of polynucleotide substrates. The phosphorolysis reaction proceeds stepwise, starting from the 3'-OH terminus of polyribonucleotides, to liberate nucleoside diphosphates. PNPase readily degrades single-stranded polyribonucleotides (75). The progress of the phosphorolytic reaction, however, depends on the secondary structure of the RNA and is impeded or stopped when the enzyme encounters stable stem-loop structures (11, 76). Thus, the enzyme acts more slowly on **transfer RNA**, ribosomal RNA (67, 75-79), or messenger RNA [except for the **poly (A)** tail, which is degraded rapidly (80)]. The rate and extent of phosphorolysis is increased by raising the temperature (75, 80) or by adding 2 M **urea** (78). PNPase phosphorolyzes short oligonucleotides by a nonprocessive mechanism (1, 81-84) but long polynucleotides by a processive mechanism (44, 67, 81). Arsenate can replace phosphate ions in degrading polynucleotides by PNPase, which produces labile 5'-phosphorylarsenate nucleotides (85, 86). Aminoacylated tRNA chains are substrates of PNPase (87), but polyribonucleotides with 3'-terminal phosphate groups are not (9). Kinetic analysis indicates that multiple subsites exist for the interaction of the enzyme with polynucleotides and oligonucleotides (73).

PNPase also catalyzes an exchange reaction between <sup>32</sup>P-labeled inorganic phosphate and the β-phosphate of nucleoside diphosphates:



This reaction is apparently a result of combined polymerization and phosphorolytic reactions that occur under approximate equilibrium conditions. The kinetic parameters of the exchange reaction are similar to the polymerization reaction. No evidence was found for reversible formation of a covalent, nucleoside monophosphate–enzyme complex (2, 9, 22).

### 0.3. Reaction Mechanisms

The processive phosphorolytic degradation of long polyribonucleotides and of tRNA by PNPase (67, 88) implies that the affinity of the enzyme for its polymeric substrate must be very high, presumably resulting from the presence of multiple binding subsites on the enzyme. On the other hand, the degradation of oligonucleotides, at least up to a chain length of six nucleotide residues, proceeds by a nonprocessive, random mechanism. The transition between the two mechanisms is still not clear.

One would expect that, once the enzyme has processively degraded a polymer to its minimum chain length of six nucleotide residues, the oligonucleotide should be released and accumulate in the reaction mixture until all of the polymer chains are degraded. Such an accumulation of oligonucleotides, however, has never been observed. The only oligonucleotide found in sizable amounts during phosphorolysis is the dimer, which cannot be phosphorolyzed. The  $K_m$  (Michaelis constant) of the *E. coli* PNPase for oligo(A)s of chain length three to six decreases from 210 mM to 50 mM as chain length increases. The presence of a 5'-phosphate residue only affects the  $K_m$  values of the tri- and tetranucleotides. For long polymers, the  $K_m$  is much smaller (of the order of 10 nM polymer). The transitional length to this very low  $K_m$  is greater than 10 and fewer than 40 nucleotide residues (89). The diversity of mechanisms probably results from the existence of different binding sites on the PNPase molecule, but the details need to be elucidated.

#### 0.4. Physiological Function

Upon its discovery, the enzyme was named polynucleotide phosphorylase, suggesting that it is involved in RNA degradation rather than synthesis (4). However, structured RNAs, such as rRNA and tRNA, are attacked slowly, whereas nonstructured homopolymers are rapidly phosphorolyzed (75, 77). Because mRNA was believed to be nonstructured then, it was suggested that the role of the enzyme is to degrade mRNA. An unequivocal demonstration of a biological function was lacking, however, for a number of years. *E. coli* contains another 3' → 5' exonuclease, RNase II. This enzyme hydrolyzes single-stranded RNA processively by releasing nucleoside 5'-monophosphates from the 3'-end. (90-93). The activities of these two exonucleases vary in different bacterial species. In *E. coli* extracts, about 90% of the degradative activity is hydrolytic and due to RNase II, 10% due to PNPase, whereas in *B. subtilis* most of the degradative activity is phosphorolytic (94).

Studies of mutants defective in these enzymes indicate that the final breakdown of mRNA *in vivo* depends on both PNPase and RNase II. Single mutant strains that lack only one of these enzymes are viable, whereas the absence of both enzymes leads to inviability in *E. coli* (95, 96) and in *B. subtilis* (97). The physiological effects of these exonucleases are influenced by the stability of the secondary structure at the 3'-end of mRNA. PNPase is less sensitive to such an obstacle than RNase II. Thus, PNPase completely degrades strongly structured tRNA molecules (67, 88), whereas RNase II cannot (98). RNase II also stalls at 3'-terminal hairpin structures of mRNA (99-101). It is believed that the decay of mRNA in *E. coli* is initiated by endonucleolytic cleavage that creates a so-called entry site for the two exonucleases (101).

Endonucleolytic cleavage is the rate-limiting step for mRNA degradation and is catalyzed by RNase E (102-105) or occasionally by RNase III (63, 64, 106-109). The observation that RNase E and PNPase form a large complex *in vivo* with at least two and perhaps more proteins, the RNA degradosome (see previous), raises the intriguing possibility that these enzymes cooperate in the coordinated regulation of RNA processing and degradation. Following an endonucleolytic cleavage by RNase E, the newly formed 3'-end would be directly attacked by PNPase. In this way, the PNPase-RNase E association in the degradosome would eliminate a slow step in which a free PNPase molecule must find the 3'-end of the mRNA. The degradosome also contains an ATP-dependent RNA helicase that aids in degrading structured RNA molecules by PNPase by unwinding RNA structures that impede its processive activity.

Degradation of RNA stem-loop structures by the degradosome has an absolute requirement for ATP and is inhibited by **antibodies** to the helicase. Interestingly, inactivation of RNase E does not inhibit the ATP-promoted degradation of RNA by the degradosome. Thus, PNPase and helicase can degrade structured RNA without endonuclease cuts (59). At least two more proteins, DnaK and enolase, are present in this complex. DnaK is also part of the degradosome and is involved in many cellular processes including the **heat shock** response. The association of the heat shock protein, DnaK, and the cold shock protein, PNPase, with RNase E raises the possibility that these proteins are involved in folding this very large enzyme complex in response to changes in the environment (61).

Enolase is also part of the degradosome but is an essential glycolytic enzyme, so that a role in RNA metabolism would not be obvious. It is possible that phosphoenolpyruvate is required to generate deoxynucleoside triphosphates (see following).

The secondary structure of RNA-OUT, the **antisense RNA** that regulates Tn10/IS10 transposition, provides greater resistance to degradation by RNase II than to PNPase (110). Moreover, the presence of RNase II protects RNA-OUT against PNPase degradation, and *rpsO* mRNA, coding for ribosomal protein S15, is also stabilized significantly by the presence of RNase II (111). The protection observed was attributed to the formation of a stable RNA-RNase II complex that sequesters the 3'-end of the transcript, making it no longer accessible to PNPase. Further studies, however, have provided a somewhat different explanation (112). Using a synthetic partially duplexed RNA as substrate, RNase II disengages from its substrate upon stalling at a region of secondary structure and reassociates with a new 3'-end. The stabilization of antisense RNA-OUT and of *rpsO* mRNA by RNase II are thus thought to be due to the removal of the 3'-overhang, rather than to the formation of a stable RNase II complex. Removal of the single-stranded 3'-overhang reduces the affinity of RNase II for the transcript and the ability of PNPase to bind and initiate phosphorolysis of RNA. In support of this theory, *in vitro* synthesized RNA that is trimmed back to a region of secondary structure with purified RNase II is not a substrate for PNPase (112). This also explains the role of poly(A) polymerase in mRNA degradation. Although polyadenylation has commonly been regarded as a special feature of eukaryotic mRNA (113-115), there are many reports of poly(A) tails on bacterial RNA [see (116-118)]. Recent work has implicated PNPase and the product of the *pcnB* gene, coding for poly(A) polymerase, in the degradation of RNA I, a highly structured, small, untranslated, antisense inhibitor of the replication of ColE1 replicons (119). Polyadenylation also destabilizes the mRNA encoding ribosomal protein S15 (118, 120) and regulates the decay of mRNA in general (121). These results support the idea that polyadenylation facilitates the exonucleolytic degradation of structured RNA, particularly by PNPase. Degradation of the mRNA encoding ribosomal protein S20 also depends on both oligoadenylation and PNPase activity. It is not sensitive to attack by RNase II and, in contrast to its stimulation of PNPase activity, polyadenylation of the 3'-end of S20 mRNA cannot overcome its natural resistance to RNase II action (99). It also appears that the degradation of polyadenylated RNA by PNPase takes place without endonucleolytic cleavage by RNase E. It seems probable that the poly(A)-dependent degradation of RNA in *E. coli* exists largely to facilitate exonucleolytic degradation of highly folded RNA fragments.

It was suggested that PNPase plays a role in supplying precursors for DNA synthesis. The NDPs released by RNA phosphorolysis can be later reduced to the deoxy versions (see [Ribonucleotide Reductases](#)) and incorporated into DNA after phosphorylation. Thus, the proposed key role for PNPase is to couple mRNA turnover to DNA synthesis (92, 122, 123).

An intriguing function of PNPase is its requirement for the development of competence in *B. subtilis* (68), which is the ability to take up exogenous DNA. It is a highly regulated phenomenon in *B. subtilis* that involves many genes. PNPase is necessary for the expression of the late competence genes, possibly functioning at a posttranscriptional step. A *B. subtilis* *pnpA* deletion strain also exhibits unusual morphological characteristics (97). PNPase is required to establish **bacteriophage P4** immunity (124). PNPase is also involved in the maturation of chloroplast mRNA, which is an essential step in synthesizing stable chloroplast RNA in higher plants (17). Finally, it was shown that guanosine pentaphosphate synthetase I from *S. antibioticus* is a bifunctional enzyme capable of both pppGpp synthesis and PNPase activities (30). Many additional roles for PNPase have recently been observed and because it has such a complicated structure, it is quite possible that other important functions are yet to be discovered.

#### 0.5. Regulation of Gene Expression

The gene encoding PNPase (*pnp*) starts with an unusual UUG triplet. It is situated at 69 minutes on the *E. coli* [chromosome](#) (125), downstream of the *rpsO* gene that codes for ribosomal protein S15. The *pnp* operons of *B. subtilis* and *Photobacterium* (70) are organized similarly (68). In *E. coli*, each

gene has its own **promoter**, but there is also co-expression of *rpsO* and *pnp* from the P1 promoter, located 108 nucleotides upstream from *rpsO* (126). The P1 transcripts are attenuated by a **rho**-independent terminator located 38 nucleotides downstream from the termination codon of *rpsO*. A second promoter, P2, is located between the *rpsO* and the *pnp* genes, 158 nucleotides upstream of *pnp*. The P1 polycistronic transcript and the P2 transcript are both processed by RNase III at two sites located 81 and 119 nucleotides upstream of the *pnp* initiation codon (63). The expression of *pnp* is under translational control, and the first step of this control is cleavage of the *pnp* message by RNase III (63, 127). Several partially deregulated mutants map in a region that overlaps the ribosome loading site, located between nucleotides -40 and +87. This indicates that the translational initiation region and the beginning of the *pnp* coding region are important for regulation. Accordingly, the stability of the message strongly depends on translation (128). One explanation for these observations is that PNPase also binds specifically to the newly formed 5'-region of the message generated by RNase III cleavage. This binding interferes with the translation process and thereby allows the enhanced degradation of its mRNA. Indeed, scanning transmission [electron microscopy](#) supports the idea that both the 5'- and 3'-ends of RNA are bound by PNPase (55). What is the role of RNase III in generating 5'-ends that are specifically recognized by PNPase? Because PNPase does not recognize specific sequences, one may suggest that the distance between the 5'-end of the mRNA and the ribosome loading site is critical in determining the capacity of PNPase to inhibit the translation of its message. In addition, expression of RNase II and PNPase might be interregulated (129). *E. coli*, *B. subtilis* and *Photothabdus* PNPases are cold shock proteins (70, 97, 130). In *Photothabdus* the increase in the level of *pnp* mRNA may be a result of a cold-inducible promoter, P2, in the intergenic region between *rpsO* and *pnp* (70). It was found that *pnp* mutants of *B. subtilis* and *E. coli* do not grow at low temperatures of 15°C and 23°C (97); C. Portier, personal communication) indicating that, under these conditions, the 3' → 5' exonuclease activity of PNPase cannot be compensated for by other nucleases.

## 0.6. Research Applications

PNPase is a useful tool for the synthesizing of homopolymers, heteropolymers, and oligonucleotides that have defined sequences. The enzyme is also employed to probe the structural conformational aspects of duplex polynucleotides involved in interferon induction (131, 132). Additional applications of the enzyme involve the synthesizing **radiolabeled** polynucleotides, the **radiolabeling** of NDPs and NTPs at their 3'-position (133-135), and sequence analysis of short oligonucleotides (84). The 3'-exonucleolytic activity of PNPase was used to analyze the size and composition of the 3'-terminal sequence of RNA molecules and their functions. With a molar excess of PNPase over the substrate, a synchronous mode of phosphorolysis is established, in which NDP molecules are sequentially released from the 3'-terminus of the RNA chains (80, 136-140). At 0°C the poly(A) tails of mRNA molecules are readily phosphorolyzed, whereas the rest of the RNA chains remain intact. Under these conditions, it was possible to determine the length of the poly(A) tails from various mRNAs and to establish their function (80, 114).

The enzyme displays low specificity for the side chains on the purine or pyrimidine moieties, whereas it shows high specificity for the number of phosphate groups on the nucleoside and the nature of the sugar moiety on the NDP substrate. The slow polymerization rate observed with some modified NDPs is enhanced by replacing  $Mg^{2+}$  with  $Mn^{2+}$ , raising the temperature, and by the adding organic solvents, such as methanol or dimethylsulfoxide [see (2)]. Blocking of NDPs at the 3'-position yields "monofunctional" substrates, of which only one residue may be added to an oligonucleotide primer, thus serving as chain terminators. The blocking group can be subsequently removed chemically from the oligonucleotide product, permitting a succession of single addition reactions. This procedure has been utilized for the stepwise synthesis of polynucleotides of defined sequence (87, 141-144).

Combining these reactions and using T4 RNA ligase to ligate the synthesized oligonucleotides permits synthesizing appreciably long oligonucleotides (145, 146). PNPase cannot phosphorolyze DNA (9). However, the enzyme can direct the reversible addition of a single deoxynucleotidyl

residue to ribooligonucleotide primers. Further addition of oligodeoxynucleotide residues to the resulting product is very sluggish (147-152). In the presence of  $Mn^{2+}$ , *E. coli* PNPase catalyzes the transfer of deoxyribonucleotide residues from the deoxyribonucleoside diphosphates to the 3'-OH end of an oligodeoxyribonucleotide primer that has a minimal length of three nucleoside residues. The repeated addition of single residues allows synthesizing oligodeoxyribonucleotides of defined sequence (153, 154).

We thank Drs. Claude Portier, Herald Putzer, and Philippe Regnier for their helpful suggestions.

### Bibliography

1. T. Godefroy-Colburn and M. Grunberg-Manago (1972) In *The Enzymes* (P. D. Boyer, ed.), Academic, New York, Vol. 7, pp. 533–574.
2. U. Z. Littauer and H. Soreq (1982) In *The Enzymes* (P. D. Boyer, ed.), Academic, New York, Vol. 17, pp. 517–553.
3. M. Grunberg-Manago and S. Ochoa (1955) Fed. Proc. **14**, 221.
4. M. Grunberg-Manago and S. Ochoa (1955) J. Am. Chem. Soc. **77**, 3165–3166.
5. M. Grunberg-Manago, P. J. Ortiz, and S. Ochoa (1956) Biochim. Biophys. Acta **20**, 269–285.
6. U. Z. Littauer (1956) Fed. Proc. **15**, 302.
7. U. Z. Littauer and A. Kornberg (1957) J. Biol. Chem. **226**, 1077–1092.
8. R. F. Beers, Jr (1956) Fed. Proc. **15**, 13.
9. M. Grunberg-Manago (1963) In *Progress in Nucleic Acid Research* (J. N. Davidson and W. E. Cohn, eds.), Academic, New York, Vol. 1, pp. 93–133.
10. Y. Kimhi and U. Z. Littauer (1967) Biochemistry **6**, 2066–2073.
11. B. Py, H. Causton, E. A. Mudd, and C. F. Higgins (1994) Mol. Microbiol. **14**, 717–729.
12. P. Owen and R. Kaback (1979) Biochemistry **18**, 1413–1422.
13. A. Abrams and P. McNamara (1962) J. Biol. Chem. **237**, 170–175.
14. S. Brishammar and N. Juntti (1974) Arch. Biochem. Biophys. **164**, 224–232.
15. Z. Kahn and H. Fraenkel-Conrat (1985) Proc. Natl. Acad. Sci. USA **82**, 1311–1315.
16. R. I. Feldman and D. S. Sigman (1984) Eur. J. Biochem. **143**, 583–588.
17. R. Hayes *et al.* (1996) EMBO J. **15**, 1132–1141.
18. A. A. Del'vig (1975) Biokhimiya **40**, 824–832.
19. A. A. Del'vig (1978) Biokhimiya **43**, 579–591.
20. R. M. S. Smillie (1963) In *Progress in Nucleic Acid Research*, (J. N. Davidson and W. E. Cohn, eds.), Academic, New York, Vol. 1, pp. 27–58.
21. F. R. Williams and M. Grunberg-Manago (1964) Biochim. Biophys. Acta **89**, 66–89.
22. Y. Kimhi and U. Z. Littauer (1968) J. Biol. Chem. **243**, 231–240.
23. Y. Kimhi and U. Z. Littauer (1968) In *Methods in Enzymology* (L. Grossman and K. Moldave, eds.), Academic, New York, Vol. **XIIB**, pp. 513–519.
24. C. Portier, R. v. Rapenbusch, M. N. Thang, and M. Grunberg-Manago (1973) Eur. J. Biochem. **40**, 77–87.
25. H. Soreq and U. Z. Littauer (1977) J. Biol. Chem. **252**, 6885–6888.
26. A. T. Gajda, G. Z. d. Behrens, and P. S. Fitt (1970) Biochem. J. **120**, 753–761.
27. C. H. Letendre and M. F. Singer (1975) Nucleic Acids Res. **2**, 149–164.
28. E. K. Barbehenn, J. E. Craine, A. Chrambach, and C. B. Klee (1982) J. Biol. Chem. **257**, 1007–1016.
29. K. Harry, N. Sharma, and P. S. Fitt (1985) Biochim. Biophys. Acta **828**, 29–38.

30. G. H. Jones and M. J. Bibb (1996) *J. Bacteriol.* **178**, 4281–4288.
31. C. B. Klee and M. F. Singer (1967) *Biochem. Biophys. Res. Commun.* **29**, 356–361.
32. C. B. Klee (1969) *J. Biol. Chem.* **244**, 2558–2566.
33. F. Hishinuma, K. Hirai, and K. Sakaguchi (1977) *Eur. J. Biochem.* **77**, 575–583.
34. G. Marumo, T. Noguchi, and Y. Midorikawa (1993) *Biosci. Biotechnol. Biochem.* **57**, 513–514.
35. G. Soe and J. Yamashita (1980) *J. Biochem. (Tokyo)* **87**, 101–110.
36. C. Babinet, A. R. J. M. Dubert, M. N. Thang, and M. Grunberg-Manago (1965) *Biochem. Biophys. Res. Commun.* **19**, 95–101.
37. M. N. Thang, M. Graffe, and M. Grunberg-Manago (1965) *Biochim. Biophys. Acta* **108**, 125.
38. W. T. Hasie and J. M. Buchanan (1967) *Proc. Natl. Acad. Sci. USA* **58**, 2468–2475.
39. J. M. Lucas and M. Grunberg-Manago (1964) *Biochem. Biophys. Res. Commun.* **17**, 395–400.
40. F. N. Brenneman and M. F. Singer (1964) *Biochem. Biophys. Res. Commun.* **17**, 401–411.
41. P. S. Fitt, F. W. Dietz, Jr., and M. Grunberg-Manago (1968) *Biochim. Biophys. Acta* **151**, 99–113.
42. P. S. Fitt and E. A. Fitt (1967) *Biochem. J.* **105**, 25–33.
43. R. E. Moses and M. F. Singer (1970) *J. Biol. Chem.* **245**, 2414–2422.
44. C. B. Klee and M. F. Singer (1968) *J. Biol. Chem.* **243**, 923–927.
45. M. F. Singer, L. Heppel, and R. J. Hilmoe (1957) *Biochim. Biophys. Acta* **26**, 447–448.
46. M. F. Singer, L. A. Heppel, and R. J. Hilmoe (1960) *J. Biol. Chem.* **235**, 738–750.
47. F. N. Brenneman and M. F. Singer (1964) *J. Biol. Chem.* **239**, 893–901.
48. R. F. Beers, Jr (1961) *J. Biol. Chem.* **236**, 2703–2709.
49. L. A. Heppel (1963) *J. Biol. Chem.* **238**, 357–366.
50. S. Mii and S. Ochoa (1957) *Biochim. Biophys. Acta* **26**, 445–446.
51. S. Ochoa and S. Mii (1961) *J. Biol. Chem.* **236**, 3303–3311.
52. J. Skoda, J. Kara, Z. Sormova, and F. Sorm (1959) *Biochim. Biophys. Acta* **33**, 579–580.
53. C. Portier (1975) *Eur. J. Biochem.* **55**, 573–582.
54. P. Regnier, M. Grunberg-Manago, and C. Portier (1987) *J. Biol. Chem.* **262**, 63–68.
55. M. Sulewski, S. P. Marchese-Ragona, K. A. Johnson and S. J. Benkovic (1989) *Biochem.* **28**, 5855–5864.
56. M. N. Thang, D. C. Thang and J. Leautey (1967) *C. R. Acad. Sci. (Paris)* **265**, 1823–1826.
57. M. Kido *et al.* (1996) *J. Bacteriol.* **178**, 3917–3925.
58. R. S. McLaren, S. F. Wewbury, G. S. C. Dance, H. C. Causton and C. F. Higgins (1991) *J. Mol. Biol.* **221**, 81–95.
59. B. Py, C. F. Higgins., H. M. Krisch, and A. J. Carpousis (1996) *Nature* **381**, 169–172.
60. A. J. Carpousis, G. Van-Houwe, C. Ehretsmann, and H. M. Kirsch (1994) *Cell.* **76**, 889–900.
61. A. Miczak, V. R. Kaberdin, C. L. Wei, and S. L. Chao (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3865–3869.
62. B. Sohlberg, U. Lundberg, F. U. Hartl, and A. V. Gabain (1993) *Proc. Natl. Acad. Sci. USA* **90**, 277–281.
63. C. Portier, L. Dondon, M. Grunberg-Manago, and P. Regnier (1987) *EMBO J.* **6**, 2165–2170.
64. P. Regnier and E. Hajnsdorf (1991) *J. Mol. Biol.* **217**, 283–292.
65. M. A. Castiglione-Morelli *et al.* (1995) *FEBS Lett.* **358**, 193–198.
66. M. Bycroft, T. J. Hubbard, M. Proctor, S. M. Freund, and A. G. Murzin (1997) *Cell* **88**, 235–242.
67. M. N. Thang, W. Guchbauer, H. G. Zachau, and M. Grunberg-Manago (1967) *J. Mol. Biol.* **26**,



403–421.

68. A. Luttinger, J. Hahn, and D. Dubnau (1996) *Mol. Microbiol.* **19**, 343–356.
69. R. D. Fleischmann *et al.* (1995) *Science* **269**, 496–512.
70. D. J. Clarke and B. C. Dowds (1994) *J. Bacteriol.* **176**, 3775–3784.
71. R. A. Harvey and M. Grunberg-Manago (1966) *Biochem. Biophys. Res. Commun.* **23**, 448–452.
72. J. F. Marlier and S. J. Benkovic (1982) *Biochem.* **21**, 2349–2356.
73. J. Y. Chou and M. F. Singer (1970) *J. Biol. Chem.* **245**, 1005–1011.
74. M. Singer (1958) *J. Biol. Chem.* **232**, 211–228.
75. M. Grunberg-Manago (1959) *J. Mol. Biol.* **1**, 240–259.
76. H. Causton, B. Py, R. S. McLaren, and C. F. Higgins (1994) *Mol. Microbiol.* **14**, 731–741.
77. U. Z. Littauer (1961) In *Protein Biosynthesis*, (R. J. C. Harris, ed.), (Academic, New York), pp. 143–162.
78. U. Z. Littauer and V. Daniel (1962) In *Acides Ribonucléiques et Polyphosphates Structure, Synthèse Fonctions* M. Grunberg-Manago and M. J. P. Ebel, eds., C. N. R. S., Strasbourg, pp. 277–294.
79. U. Z. Littauer and H. Eisenberg (1959) *Biochim. Biophys. Acta* **32**, 320–327.
80. H. Soreq, U. Nudel, R. Salomon, M. Revel, and U. Z. Littauer (1974) *J. Mol. Biol.* **88**, 233–245.
81. J. Y. Chou and M. F. Singer (1970) *J. Biol. Chem.* **245**, 995–1004.
82. M. F. Singer, R. J. Hilme, and M. Grunberg-Manago (1960) *J. Biol. Chem.* **235**, 2705–2712.
83. A. Guissani (1977) *Eur. J. Biochem.* **79**, 233–243.
84. G. Kaufmann, H. Grosfeld, and U. Z. Littauer (1973) *FEBS Lett.* **31**, 47–52.
85. M. F. Singer and B. M. O'Brien (1963) *J. Biol. Chem.* **238**, 328.
86. M. F. Singer (1963) *J. Biol. Chem.* **238**, 336–343.
87. G. Kaufmann and U. Z. Littauer (1970) *Eur. J. Biochem.* **12**, 85–92.
88. M. N. Thang, B. Beltchev and M. G. Manago (1971) *Eur. J. Biochem.* **19**, 184–193.
89. T. Godefroy (1970) *Eur. J. Biochem.* **14**, 222–231.
90. P. F. Spahr and D. Schlessinger (1963) *J. Biol. Chem.* **238**, 2251–2253.
91. P. F. Spahr (1964) *J. Biol. Chem.* **239**, 3716–3726.
92. M. Sekiguchi and S. S. Cohen (1963) *J. Biol. Chem.* **238**, 349–356.
93. M. F. Singer and G. Tolbert (1965) *Biochem.* **4**, 1319–1330.
94. M. P. Deutscher and N. B. Reuven (1991) *Proc. Natl. Acad. Sci. USA* **88**, 3277–3280.
95. T. G. Kinscherf and D. Apirion (1975) *Mol. Gen. Genet.* **139**, 357–362.
96. W. P. Donovan and S. R. Kushner (1986) *Proc. Natl. Acad. Sci. USA* **83**, 120–124.
97. W. Wang and D. H. Bechhofer (1996) *J. Bacteriol.* **178**, 2375–2382.
98. R. S. Gupta, T. Kasai, and D. Schlessinger (1977) *J. Biol. Chem.* **252**, 8945–8949.
99. G. A. Coburn and G. A. Mackie (1996) *J. Biol. Chem.* **271**, 15776–15781.
100. G. Guarneros and C. Portier (1990) *Biochimie* **72**, 771–777.
101. F. Braun, E. Hajnsdorf, and P. Regnier (1996) *Mol. Microbiol.* **19**, 997–1005.
102. B. K. Ghora and D. Apirion (1978) *Cell* **15**, 1055–1066.
103. C. P. Ehretsmann, A. J. Carpousis, and H. M. Krisch (1992) *Genes Dev.* **6**, 149–159.
104. C. F. Higgins, S. W. Peltz, and A. Jacobson (1992) *Curr. Opin. Genet. Dev.* **2**, 739–747.
105. O. Melefors, U. Lundberg, and A. von-Gaba in (1993) In *Control of mRNA Stability* (J. G. Belasco and G. Brawerman, eds.), Academic, New York, pp. 53–70.
106. J. C. Bardwell *et al.* (1989) *EMBO J.* **8**, 3401–3407.

107. D. Court (1993) In *Control of mRNA Stability* (J. G. Belasco and G. Brawerman, eds.), Academic, New York, pp. 71–116.
108. E. Hajnsdorf, A. J. Carpousis, and P. Regnier (1994) *J. Mol. Biol.* **239**, 439–454.
109. P. Regnier and M. Grunberg-Manago (1990) *Biochimie* **72**, 825–834.
110. C. N. Pepe, S. Maslesa-Galic, and R. W. Simons (1994) *Mol. Microbiol.* **13**, 1133–1142.
111. E. Hajnsdorf, O. Steier, L. Coscoy, L. Teyssset, and P. Regnier (1994) *EMBO J.* **13**, 3368–3377.
112. G. A. Coburn and G. A. Mackie (1996) *J. Biol. Chem.* **271**, 1048–1053.
113. G. Brawerman (1981) *CRC Crit. Rev. Biochem.* **10**, 1–38.
114. U. Z. Littauer and H. Soreq (1982) In *Progress in Nucleic Acid Research* (W. E. Cohn, ed.), Academic, New York, Vol. **27**, pp. 53–83.
115. E. J. Baker (1993) In *Control of Messenger RNA Stability* (J. G. Belasco and G. Brawerman, eds.), Academic, New York, pp. 367–415.
116. G.-J. Cao and N. Sarkar (1992) *Proc. Natl. Acad. Sci. USA* **89**, 10380–10384.
117. L. He *et al.* (1993) *Mol. Microbiol.* **9**, 1131–1142.
118. E. Hajnsdorf, F. Braun, J. Haugel-Nielsen, and P. Regnier (1995) *Proc. Natl. Acad. Sci. USA* **92**, 3973–3977.
119. F. Xu and S. N. Cohen (1995) *Nature* **374**, 180–183.
120. J. Haugel-Nielsen, E. Hajnsdorf, and P. Regnier (1996) *EMBO J.* **15**, 3144–3152.
121. O'Hara *et al.* (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1807–1811.
122. J. Fricke, J. Neuhard, R. A. Kelln, and S. Pedersen (1995) *J. Bacteriol.* **177**, 517–523.
123. A. Danchin (1996) *Mol. Microbiol.* **20**, 895–897.
124. F. Piazza, M. Zappone, M. Sana, F. Briani, and G. Deho (1996) *J. Bacteriol.* **178**, 5513–5521.
125. C. Portier, C. Migot, and M. Grunberg-Manago (1981) *Mol. Gen. Genet.* **183**, 298–305.
126. C. Portier and P. Regnier (1984) *Nucleic Acids Res.* **12**, 6091–6102.
127. M. R.-L. Meur and C. Portier (1992) *EMBO J.* **11**, 2633–2641.
128. M. R.-L. Meur and C. Portier (1994) *Nucleic Acids Res.* **22**, 397–403.
129. R. Zilhao, J. Plumbridge, E. Hajnsdorf, P. Regnier, and M. Arraiano (1996) *Microbiology* **142**, 367–375.
130. P. G. Jones, R. A. VanBogelen, and F. C. Neidhart (1987) *J. Bacteriol.* **169**, 2092–2095.
131. P. F. Torrence (1981) In *Methods in Enzymology* (S. Pestka, ed.), (Academic, New York), Vol. **78**, pp. 326–331.
132. M. N. Thang and W. Guschlbauer (1992) *Pathol. Biol. Paris* **40**, 1006–1014.
133. U. Z. Littauer, Y. Kimhi, and M. Avron (1964) *Anal. Biochem.* **9**, 85–93.
134. E. Gilboa, H. Soreq, and H. Aviv (1977) *Eur. J. Biochem.* **77**, 393–400.
135. B. Vennstrom, U. Pettersson, and L. Philipson (1978) *Nucleic Acids Res.* **5**, 205–219.
136. U. Nudel, H. Soreq, and U. Z. Littauer (1976) *Eur. J. Biochem.* **64**, 115–121.
137. R. Salomon, I. Sela, H. Soreq, D. Givon, and U. Z. Littauer (1976) *Virology* **71**, 74–84.
138. H. Grosfeld, H. Soreq, and U. Z. Littauer (1977) *Nucleic Acids Res.* **4**, 2109–2122.
139. H. Soreq, A. D. Sagar, and P. B. Sehgal (1981) *Proc. Natl. Acad. Sci. USA* **78**, 1741–1745.
140. P. B. Sehgal, H. Soreq, and I. Tamm (1978) *Proc. Natl. Acad. Sci. USA* **75**, 5030–5033.
141. G. Kaufmann, G. Fridkin, M. Zutra, and U. Z. Littauer (1971) *Eur. J. Biochem.* **24**, 4–11.
142. G. N. Bennett, J. K. Mackey, J. L. Wiebers, and P. T. Gilham (1973) *Biochem.* **12**, 3956–3962.
143. G. C. Walker and O. C. Uhlenbeck (1975) *Biochem.* **14**, 817–824.
144. Y. Kikuchi, K. Hirai, and K. Sakaguchi (1975) *J. Biochem. (Tokyo)* **77**, 469–472.
145. G. Kaufmann and U. Z. Littauer (1974) *Proc. Natl. Acad. Sci. USA* **71**, 3741–3745.

146. P. J. Romaniuk and O. C. Uhlenbeck (1983) In *Methods in Enzymology* (L. Grossman and K. Moldave, eds.), (Academic, New York), Vol. **100**, pp. 52–59.
147. G. Kaufmann and U. Z. Littauer (1969) *FEBS Lett.* **4**, 79–83.
148. S. Bon, T. Godegroy, and M. Grunberg-Manago (1970) *Eur. J. Biochem.* **16**, 363–372.
149. J. Y. Chou and M. F. Singer (1971) *J. Biol. Chem.* **246**, 7486–7496.
150. J. Y. Chou and M. F. Singer (1971) *J. Biol. Chem.* **246**, 7497–7504.
151. J. Y. Chou and M. F. Singer (1971) *J. Biol. Chem.* **246**, 7505–7513.
152. J. Y. Chou and M. F. Singer (1971) *Biochem. Biophys. Res. Commun.* **42**, 306–311.
153. S. Gillam and M. Smith (1980) In *Methods in Enzymology* (L. Grossman and K. Moldave, eds.), (Academic, New York), Vol. **65**, pp. 687–701.
154. M. Singer and P. Berg (1991) *Genes and Genomes: A changing perspective* University Science Mill Valley, CA.

## Polyomavirus

Polyomaviruses include BKV, JCV, [SV40](#), LPV, murine polyomaviruses, hamster polyomavirus, and others. BKV and JCV grow permissively in human cells, while SV40 and LPV grow permissively in monkey cells. The murine polyomavirus will here be simply referred to as polyomavirus. Because polyomavirus produces no detectable diseases when it multiplies in adult mice, it is regarded as a harmless passenger. However, polyomavirus inoculated into newborn mice induces parotid gland carcinoma, epidermal carcinoma, hemangioma, mammary carcinoma, thymic epithelioma, and many others, although productive virus infection does not establish in the mice. Polyomavirus replicates well in both primary cultures and established murine cell lines and produces progeny viruses (permissive infection). Primary cultured cells from other rodent origins are abortively infected with polyomavirus (nonpermissive infection). Some, but not all, of the infected cells are immortalized. Some immortalized cells are released from the tightly controlled cellular growth-regulatory system and become transformed. In contrast to SV40-transformed cells, polyomavirus-transformed cells are highly **oncogenic** in animals because of the function of the polyomavirus-encoded middle [T Antigen](#).

The virion of polyomavirus has a diameter of about 45 nm and a [sedimentation coefficient](#) of 240S (see [Papovavirus](#)). The [genome](#) of polyomavirus is double-stranded circular DNA, 5292 base pairs in length, and encodes early and late genes. Three proteins with molecular weights of 100 kDa, 56 kDa, and 22 kDa, which are named large, middle, and small T antigens, respectively, are translated from the early [messenger RNA](#) that are transcribed before the onset of viral [DNA replication](#). The large and middle T antigens are localized exclusively in nuclei and on the inner and plasma membranes, respectively. The small T antigen exists in all fractions as a soluble form. Viral capsid proteins VP1, VP2, and VP3 are encoded by the late gene.

Considerably less is known about the molecular basis of polyomavirus DNA replication than SV40. The large T antigen can bind to **retinoblastoma** family proteins but, unlike SV40 large T antigen, fails to associate with [p53](#). Inactivation of p53 might not be necessary for making murine cells enter S phase. The large T antigen binds to the origin of [DNA replication](#) and unwinds the adjacent region by its ATPase and DNA helicase activities (preinitiation complex). Formation of the initiation complex and the chain elongation reaction are considered to occur in a similar mode as described in the SV40 system (see [SV40 \(Simian Virus 40\)](#)). The phosphorylation at Thr278 of the large T

antigen by **cyclin**-dependent kinase is required for viral DNA replication.

The small T antigen forms a complex with protein **phosphatase** 2A and inhibits its enzymatic activity. Thereby, the enhanced mitogen-activated protein (MAP) kinase pathways stimulate cell proliferation. This function is dispensable for viral DNA replication, however, because a mutant lacking the small T antigen still produces progeny viruses at the same efficiency as wild type.

Transfection experiments using [complementary DNAs](#) encoding the large, middle, and small T antigens separately have produced the following results. First, the large T antigen alone can immortalize primary-cultured cells from both permissive and nonpermissive animals. Second, the middle T antigen alone can malignantly transform established cells, but not primary-cultured cells. Third, in addition to the middle T antigen, coexpression of the large or small T antigen is necessary for both transformation of primary-cultured cells and induction of tumors in newborn animals.

The oncoproteins c-src, c-yes, and c-fyn are associated with the region from residues 185 to 210 of the middle T antigen. The middle T antigen activates c-src by repressing **phosphorylation** at Tyr527 of c-src or by making the phosphorylated Tyr527 accessible to protein **phosphatase**. The middle T antigen can similarly activate c-yes, but not c-fyn. Once the activated c-src kinase phosphorylates Tyr250, 315, and 322 of the middle T antigen, the [SH2 domains](#) of SHC, the regulatory subunit of PI3 kinase, and **phospholipase** C-g1 recognize the amino acid sequence containing the respective phosphotyrosine residues. Then, the tyrosine residues of the three proteins are phosphorylated by the activated c-src associated with the middle T antigen. The other SH2 domain-containing proteins, such as Grb2 and IRS-1, form complexes with the SHC, phosphoinositide-3-kinase subunit, and phospholipase C-g1 through the phosphotyrosine residues. The subsequent unusual activation of the [signal transduction](#) pathways mediated by the Raf and/or S6 kinase cascade is considered to be the molecular basis of the transformation by the middle T antigen.

#### Suggestions for Further Reading

S. M. Dilworth (1995) Polyoma virus middle T antigen: meddler or mimic? *Trends Microbiol.* **3**, 31–35.

F. Kiefer, S. A. Courtneidge, and E. F. Wagner (1994) Oncogenic properties of the middle T antigen of polyomaviruses. *Adv. Cancer Res.* **64**, 125–157.

K. V. Shah (1996) "Polyomaviruses". In *Fields Virology*, 3rd edition (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2027–2043.

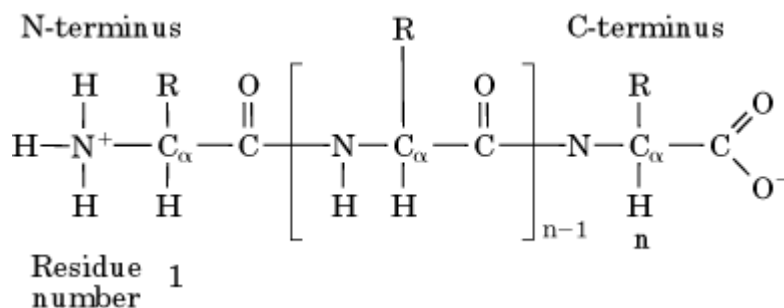
## Polypeptide Chain

Polypeptide chain is a term that describes the basic covalent structure of [proteins](#). That is, proteins are **polymers** comprising [peptide bond](#) units (hence, *polypeptide*) that connect together a large number of [amino acids](#) into a linear chain. The polypeptide chain thus consists of a regularly repeating part, the [backbone](#) or main chain, and a variable part, the [side chains](#), corresponding to the functional groups of the amino acids. The term is most often used for proteins, which can consist of one or more polypeptide chains, but can also be used more generally for all amino acid polymers including **peptides**, [polyamino acids](#), and chemically synthesized polymers of amino acids.

The linear polypeptide chain is formed from amino acids by a condensation reaction between the acidic or [carboxyl group](#) of one amino acid and the [amino group](#) of another. Amino acid units within the polypeptide chain are referred to as *residues*. The backbone of a polypeptide chain is a repeating

sequence of three atoms corresponding to the amide nitrogen (N) of the peptide bond, the central tetrahedral carbon ( $C_{\alpha}$ ), and the peptide carbonyl carbon (C) (Fig. 1). The variable side chains or functional groups (R) are covalently linked to  $C_{\alpha}$ . The amino acid sequence defines the order of the side chains throughout the polypeptide chain.

**Figure 1.** Schematic representation of the polypeptide chain of a protein, showing in square brackets the atoms of the peptide repeating unit (N,  $C_{\alpha}$ ; and C) that form the backbone. The functional group or side chain, R, of each residue varies, depending on the type of the amino acid residue. By convention, numbering of amino acid residues in the polypeptide chain is from the *N*-terminus (residue 1) to the *C*-terminus (residue *n*).



The two ends of the polypeptide chain differ from each other. At one end there is an uncondensed or free amino group, and at the other end there is a free carboxyl group. The ends of the chain are therefore referred to as the amino terminus or ***N*-terminus** and the carboxyl terminus or ***C*-terminus**. By convention, numbering of the amino acid residues in the polypeptide chain is from the *N*-terminus to the *C*-terminus (Fig. 1).

The chemical nature of the peptide bond in the polypeptide chain allows [hydrogen bond](#) formation between peptide bond units. The propensity of the polypeptide chain to form these intrachain interactions gives rise to the several types of regular backbone structure, or **secondary structure**, observed in [protein structures](#) including **a-helices**, **b-strands**, and [turns](#).

[See also [Peptide Bond](#) and [Protein Structure](#).]

Suggestion for Further Reading

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

## Polyploid

Cells with more than the normal complement of DNA in their **genomes**, usually with increased numbers of the standard [chromosomes](#), are called polyploid. In many plants, but more rarely in animals, mutations occur that increase the normal number of chromosomes. For example, if the [haploid](#) set of chromosomes of a given plant is 7, the diploid set is 14. Mutants appear that contain 3, 4, 5, 6, 7, 8, etc. genomes. There are polyploid forms called **triploid** (three genomes), tetraploid (four genomes), pentaploid (five genomes), etc. Each of the chromosomes is repeated by the same

integer. During meiosis, polyploidy is maintained in subsequent generations if there is an even number of genomes.

## Polyproline

Polyproline is the **polyamino acid** formed from [proline](#) residues. Like [proteins](#) and other **polypeptide chains**, the polyproline [polymer](#) is formed by condensation of the [amino group](#) of one [amino acid](#) and the [carboxyl group](#) of another amino acid. Proteins, however, have a defined and complex sequence that usually includes all 20 of the amino acids, but polyproline is formed from proline residues only. Polyproline adopts an unusual [protein structure](#) that is often used in **molecular recognition**.

Unlike the polypeptide chain formed from other amino acid residues, the [peptide bond](#) of the proline [backbone](#) cannot take part in [hydrogen bond](#) interactions because there is no hydrogen atom on the backbone nitrogen. Furthermore, the cyclic [side chain](#) of proline restricts the backbone [conformation](#) of polyproline. Therefore, the **secondary structure** elements observed in most protein structures, such as **a-helices** or **b-strands**, are not observed in the structure of polyproline. Instead, the polyproline backbone adopts one of two possible conformations, called poly(Pro) I (backbone  $\phi$  and  $\psi$  angles of  $-83^\circ$  and  $158^\circ$ , respectively) and polyProII (backbone  $\phi$  and  $\psi$  angles of  $-78^\circ$  and  $149^\circ$ , respectively) (see [Ramachandran Plot](#)). The difference between the two is that poly(Pro) I has all its peptide bonds in the *cis* [configuration](#) (backbone  $\omega$  angle  $0^\circ$ ) and poly(Pro) II (Fig. 1) has the all *trans* [configuration](#) (backbone  $\omega$  angle  $180^\circ$ ).

**Figure 1.** Schematic representation of the polyPro(II) helix for a three-residue sequence of polyproline. The first residue is at the top of the figure, and the third residue is at the bottom. The cyclic side chains of the proline residues are shown in light gray. The backbone nitrogen atoms, and backbone oxygen atoms are shown as dark spheres. The curve of the helix is shown by the dark line that traces the backbone of the polyproline sequence. This figure is generated using Molscript and Raster3D .



The polyamino acid formed from [glycine](#) (polyglycine) adopts a conformation called poly(Gly) II, which is similar to the poly(Pro) II conformation. Similarly, the backbone of each polypeptide chain of [collagen](#), a structural protein that has a distinctive repeating unit in its polypeptide sequence corresponding to Gly–X–Y, where X is often proline, also adopts a conformation like that of poly(Pro) II.

Regions of protein sequence that are rich in proline residues can also adopt the poly(Pro) II conformation. These proline-rich regions are important recognition elements for protein–protein interactions. For example, **SH3 domains** of proteins recognize proline-rich regions in other proteins.

### Bibliography

1. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
2. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
3. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

4. V. Sasisekharan (1959) Structure of poly-L-proline II. *Acta Crystallogr.* **12**, 897–903. (Early description of the polyPro(II) conformation.)
5. A. A. Adzhubei and M. J. E. Sternberg (1993) Left-handed polyproline II helices commonly occur in globular proteins. *J. Mol. Biol.* **229**, 472–493.
6. H. Yu, J. K. Chen, S. Feng, D. C. Dalgarno, A. W. Brauer, and S. L. Schreiber (1994) Structural basis for the binding of proline-rich peptides to SH3 domains. *Cell* **76**, 933–945.

## Polyproteins

Large precursor [proteins](#) encoding one or more active **peptides** or **polypeptide chains** are referred to as polyproteins. Polyproteins are synthesized on [ribosomes](#) in the normal method of **protein biosynthesis**, transported into the lumen of the rough [endoplasmic reticulum](#) (ER) and then transported to the [Golgi apparatus](#), where they are sorted into [secretory vesicles](#). In this process, the proteins undergo a series of **posttranslational modifications** and are processed into multiple copies of the same peptide or into different polypeptides or proteins. Thus many peptide [growth hormones](#) (see [Peptide Hormones](#)) and signaling peptides are cleaved proteolytically from their precursor pro-hormones by endoproteinases at basic residues or pairs of basic residues at the recognition sequence. This often occurs in the secretory granules or vesicles, prior to [exocytosis](#). Furthermore, the differential processing of the polyprotein in various tissues often leads to different sets of products (1).

For example, pro-opiomelanocortin (POMC) is a 32-kDa polyprotein that is processed to adrenocorticotrophic hormone (ACTH) and b-lipotropin (b-LPH) in the corticotrophic cells of the anterior pituitary under the control of corticotropin-releasing hormone (CRH). However, alpha-melanocyte-stimulating hormone (α-MSH), corticotropin-like intermediary peptide (CLIP), g-lipotropin, and b-endorphin are the products generated in the intermediary pituitary under the control of dopamine (2) (see Fig. 1 of [Peptide Hormones](#)). Like many [translation](#) products, polyproteins also contain a [signal peptide](#) for transport into the ER that is cleaved before processing. In another type of polypeptide processing, multiple copies of [enkephalin](#) are contained within the enkephalin precursor

(3).

Other examples of polyproteins are found in the translation products of **viruses**. Production of the infectious human immunodeficiency virus ([HIV](#)), for example, requires proper polyprotein processing by the dimeric HIV viral proteinase to activate it. Other viral polyproteins include those of the hepatitis A, C, and G, rubella, dengue, poliovirus, foot and mouth disease, and yellow fever viruses. These polyproteins not only play a role in infection and disease, but they may be of use in gene therapy. Viral **vectors** are now being investigated as a means of expressing and delivering peptide drugs to patients with deficiencies or genetic mutations (4).

Likewise, hybrid genes in **transgenic** animals can control synthesis of human polyproteins that can be isolated from milk, subsequently processed into human proteins, and used as therapeutic drugs (5).

### Bibliography

1. W. W. Chin (1995) *Principles and Practices of Endocrinology and Metabolism*, 2nd ed. (K. L. Becker, ed.), JB Lippincott, Philadelphia, pp. 8–20.
2. M. G. Castro and E. Morrison (1997) *Crit. Rev. Neurobiology*, **11**, 35–57.
3. R. B. Parekh and C. Rohlff (1997) *Curr. Opin. Biotechnol.* **8**, 718–723.
4. D. A. Brooks (1997) *FEBS Lett.* **409**, 115–120.
5. M. Nakanishi (1997) *Crit. Rev. Ther. Drug Carrier Syst.* **12**, 263.

### Polyribosome, Polysome

The [ribosome](#) binds to one end of a [messenger RNA](#) and initiates protein biosynthesis, moving along the mRNA. mRNAs that are being translated actively are associated with multiple ribosomes, forming polysomes. Because the length of mRNA covered by one ribosome is 30 to 35 nucleotides, the size of the polyribosome varies, depending on the length of the mRNA as well as on its translatability.

In prokaryotic cells, where translation and [transcription](#) are coupled in the same cellular compartment, mRNA is bound with ribosomes to start translation even before the transcript is complete, forming polysomes. The eukaryotic translation machinery has been found to be localized near the [nuclear pore complex](#) and to bind to the [Cap](#) signal of mRNA when translocated from the nucleus to the cytoplasm. During translation of eukaryotic secretory or membrane-associated proteins, ribosomes are bound to the [endoplasmic reticulum](#) (ER) through the [signal peptide](#) of the nascent polypeptide. The polysomes are visible there, giving this part of the endoplasmic reticulum the name *rough endoplasmic reticulum*.

### Polysomy

In cases of *polysomy*, in contrast to **polyploidy**, there is only one kind of [chromosome](#) in the triple



or quadruple state in the [genome](#), whereas all others have the normal diploid number. Many cases of polysomy have been described in the **plants** *Datura stramonium*, *Oenothera lamarckiana*, *Matthiola incana*, *Nicotiana tabacum*, *Lycopersicon esculentum*, and in the fly ***Drosophila melanogaster***. Polysomy is generally constant over generations. A case worth mentioning is that of **trisomies** in man.

### 1. Klinefelter's Syndrome in Humans

*Klinefelter's syndrome* is characterized by the [karyotype](#) 47,XXY, in other words by the presence of an additional [X-chromosome](#) in the male genome. It is associated with high stature, hypogonadism (small firm testes), azoospermia and infertility, gynecomastia, bilateral cryptorchidism, obesity and in some cases, by minor mental retardation. Testicular biopsy has revealed hyperplasia of the Leydig cells, which contain microcrystalline formations, absence of annular lamellae in the Sertoli cells, and hyalinization of the seminiferous tubular membrane. A certain number of patients are diagnosed only by their sterility. Prolonged testosterone therapy is beneficial. Secondary sex characteristics increase and gynecomastia generally disappears.

In addition to the pure 47,XXY karyotype, patients having 48,XXXY and 49,XXXXY have been described. Hormonal examination of the patients of these two groups reveals hypoandrogenia and a relative hyperestrogenuria. The occurrence of Klinefelter's syndrome occurs in one of every 600 live births.

XX-males have also been encountered, without a Y-chromosome in the karyotype. From a study using restriction fragment length polymorphisms (**RFLP**), these XX-males seem to have inherited an X-chromosome from each parent. The results argue that this is not a genetic **mosaic** involving a Y-containing cell line in males, but they do not exclude an X-Y (or Y-autosome) chromosomal [translocation](#) during paternal **meiosis** (1).

### Bibliography

1. D. C. Page and A. de la Chapelle (1984) Am. J. Hum. Genet. **36**, 565–575.

## Polytene

In certain types of secretory cells of fly larvae, all copies of [homologous chromosomes](#) remain side by side, creating a single, giant *polytene* chromosome. In some cases, large cells undergo a polytene to [polyploid](#) conversion, in which the chromosomes separate. This demonstrates that the basic structure of a polytene chromosome is similar to that of a normal chromosome.

## Polytene Chromosome

Polytene chromosomes result from the duplication of **DNA** during the [cell cycle](#) without a following **mitosis**, so that the **chromatids** do not separate. Each **chromosomal** homologue (see [Homologous Chromosomes](#)) contains as many as 1,024 chromatids in *Drosophila melanogaster*. The process of

repeated rounds of [DNA replication](#) in a single **S-phase** is called Endoreduplication. If the chromatids separate and remain in the same nucleus, the cells are called Endopolyploid. Giant polytene chromosomes result when the chromatids do not separate. Polytene chromosomes are immense compared to the normal chromosomes found in a **diploid** somatic [nucleus](#). All of the homologous pairs of chromosomes remain side by side, forming a single giant chromosome. Like [lampbrush chromosomes](#), it is possible to isolate polytene chromosomes under physiological ionic conditions with their higher order [chromatin](#) structure preserved.

This integrity of structure has facilitated a great deal of informative cytological and **immunofluorescent** analysis of chromosomal structure. Immunologic staining has revealed the specific distribution of several non-[histone](#) proteins, including **RNA polymerase II** and proteins apparently responsible for the generation of inactive [heterochromatin](#). Grossbach and colleagues used **antibodies** against variants of histone H1 to demonstrate that the localization of particular linker histones can be highly specific for individual chromosomal **domains** within polytene chromosomes (1). In contrast, a similar approach with antibodies raised against **HMG14**-related proteins leads to the general immunofluorescent staining of transcriptionally active domains. Turner found that individual chromosomal domains within polytene chromosomes are significantly enriched with forms of histone H4 that have particular states of **posttranslational modification** (2). In related experiments, it has been discovered that hyperacetylation of histone H4 on the male [X-chromosome](#) of *Drosophila* correlates with the increased transcriptional activity necessary for the phenomenon of **dosage compensation**. In this phenomenon, the transcriptional activity of all genes on the single male X-chromosome increases relative to genes on the other chromosomes, which are present in two copies per cell. These studies strongly suggest that chromosomes possess highly selective microheterogeneity in protein composition in which individual chromosomal domains contain particular histone variants or posttranslational modification states. Now important questions exist about how and why a particular protein or enzymatic activity leading to histone modification is targeted at an individual chromosomal domain. More detailed ultrastructural analysis, especially of the [Balbiani ring](#) genes in the polytene larval salivary glands of *Chironomus tentans*, has yielded much information concerning transcription, chromatin structure, and RNA processing.

#### Bibliography

1. E. Mohr, L. Trieschmann, and U. Grossbach (1989) Proc. Natl. Acad. Sci. USA **86**, 9308–9312.
2. B. M. Turner, A. J. Birley, and J. Lavender (1992) Cell **69**, 375–384.

#### Suggestion for Further Reading

3. D. DePomerai (1990) *From Gene to Animal*, 2nd ed., Cambridge University Press, Cambridge, UK.

## Ponceau S

Ponceau S is a dye commonly used to stain [proteins](#) on [blotting](#) membranes (1). The formal nomenclature of Ponceau S [I] is 3-hydroxy-4-[2-sulfo-4-(4-sulfophenylazo)phenylazo]-2,7-naphthalenedisulfonic acid, tetrasodium salt. Its molecular weight is 760.58. Ponceau S is used primarily because it is simple and rapid. It is not very sensitive, however, requiring 200 ng of protein in a normal spot to be visible. Proteins appear red on a pink background. The staining is reversible, so the blot can be used after initial protein staining for a second detection method.

#### Bibliography

1. S. Best and D. W. Speicher (1995) In *Current Protocols in Protein Science* (J. E. Coligan et al. eds.), Wiley, pp. 10.8.3.

## Pore Gradient Electrophoresis

It is often desirable to have a gradient of pore sizes for [gel electrophoresis](#), to provide varying degrees of **molecular sieving** to mixtures of [macromolecule](#) differing in size. This can be accomplished by generating a gradient of gel concentrations while preparing the gel. During electrophoresis, the sample molecules will encounter varying pore sizes as they migrate in the gel, which will affect their electrophoretic mobilities to an extent depending on their sizes and shapes.

Gel concentration gradients are usually formed using [polyacrylamide](#), but they can also be formed with [agarose](#). Three kinds of gradient gels need to be distinguished, each serving a different purpose:

1. Pore gradients formed parallel to the direction of electrophoresis, using a gel concentration range that allows for the migration of the sample at all gel concentrations. The decreasing mobility of the faster-migrating species results in a compaction of the gel pattern of the sample, so that molecules differing drastically in their mobilities can be resolved on the same gel. This is especially important for **SDS-PAGE**, where a greater range of molecular weights can be resolved on a single gradient gel than on a normal gel.
2. Pore limit gels are similar to those in gradient gel 1, but the samples encounter gel concentrations high enough to virtually halt their migration (2). Such pore limit gels allow the particle size of the arrested species to be estimated by comparison of its “arrested” position with that of molecular size standards of the same shape.
3. The third variety has the gel gradient perpendicular to the direction of migration, for [transverse gradient gel electrophoresis](#). A uniform sample applied across the top of the gel migrates at continually varying gel concentrations and pore sizes; a single gel can give an entire [Ferguson plot](#) of the electrophoretic mobility as a function of gel concentration.

The reproducibility of polyacrylamide gel gradients depends critically on the control of the polymerization reaction over the entire range of gel concentrations; such control may require more than the usual type of gradient maker. Suitable gradient makers and computer programs predicting pore gradient profiles of any shape and any slope are available (3).

### Bibliography

1. D. Rodbard, G. Kapadia, and A. Chrambach (1971) *Anal. Biochem.* **40**, 135–157.
2. K. Felgenhauer (1979) *J. Chromatogr.* **173**, 299–311.
3. K. Altland and A. Altland (1984) *Electrophoresis* **5**, 143–147.

### Suggestion for Further Reading

4. G. H. Weiss and D. Rodbard (1972) Diffusion dependent peak broadening in pore gradient electrophoresis. *Separation Sci.* **7**, 217–232.

## Porin

Porins are [membrane proteins](#) of the outer membrane of [Gram-negative](#) bacteria, with an [b-barrel](#) architecture (1, 2). Although much more is known about bacterial porins, there are functionally and structurally similar proteins, known as voltage-dependent anion channels (VDACs), that are believed to be present in the outer membrane of [mitochondria](#) (3). Known porins are composed of 250 to 450 amino acid residues. Their function is to facilitate [diffusion](#) of small molecules across the membrane by allowing solutes to pass through an aqueous channel in the middle of the transmembrane b-barrel. Some porins are nonspecific and permeable to any solutes smaller than 600 Da, but they can be cation- or anion-selective, favor polar solutes over [nonpolar](#) ones, or have specific substrates, such as maltodextrins (α1–4-polyglucose) or sucrose. The rate of transport in the nonspecific porins is a linear function of the concentration gradient of the solute. In contrast, the specific porins follow [Michaelis–Menten kinetics](#), indicating initial binding of the solute.

The porins for which three-dimensional [protein structures](#) have been determined are homotrimers containing three identical transmembrane channels. The OmpA protein of *Escherichia coli* is thought to be a monomeric porin with a barrel comprised of eight **b-strands**. The membrane-bound form of α-hemolysin, a bacterial [toxin](#), belongs to the same structural class, because its transmembrane region comprises a 14-stranded antiparallel b-barrel (4). A similar structure has been predicted for the channel assembled by the bacterial toxin aerolysin in its heptameric membrane-bound form.

The size of the pore is determined by one or more extracellular surface loops of polypeptide chain that fold back into the channel. Although the 18-stranded b-barrel of glycoporin is clearly wider than the 16-stranded barrel of general porins, the latter have a wider channel because their intrachannel loop structures are less extensive. The width of the channel mouth (the eyelet) of the general porins is determined by one long and structured loop between adjacent b-strands, which controls access to the channel. In some porins, additional control is provided by a transverse electric field, which is generated by an uneven distribution of positively and negatively charged residues inside the channel mouth. In maltoporins and sucrose porin, the channel mouth is constricted by three surface loops. For this reason, it is much narrower than the eyelet in the nonspecific porins. There is an extensive binding site or an aromatic path for the sugar rings to be translocated (5-8). Suitably positioned partners for [hydrogen bonds](#) to the sugar hydroxyl groups assist the translocation process.

Porins have a number of peculiar properties. First, their sequences are at least as [hydrophilic](#) as are those of soluble proteins. Yet, the lipid-exposed surface of the b-barrel is highly **hydrophobic**, and the proteins are only soluble in the presence of a [detergent](#). Second, the primary structures of porins form at least 10 families that show no clear sequence [homology](#) with each other. It appears that highly diverse [primary structures](#) can fold into very similar three-dimensional structures, as shown by the eleven porin structures presently known at atomic resolution (see Table 1 of [Membrane Proteins](#)). The sequences of the mutually homologous nonspecific porins of *E. coli* (OmpF and PhoE) are not related to those of the general diffusion porins from the phototrophic bacteria *Rhodobacter capsulatus* and *Rhodospseudomonas blastica*, or from *Paracoccus denitrificans*, despite the fact that all five fold into 16-stranded b-barrels. The maltoporin family has an 18-stranded topology and is not related to the nonspecific porins. Third, porins are unusually stable, forming trimers that are resistant toward **denaturation** by [SDS](#), even at elevated temperatures. Perhaps owing to this high stability, porins appear to be relatively easy to crystallize, as shown by the existence of a number of high-resolution structures.

In fact, the first well-ordered crystals of a membrane protein were those of the *E. coli* OmpF porin grown in the late 1970s in Jürg Rosenbusch's laboratory in Basel. Owing to an unfavorable crystal symmetry, however, this crystal structure was not determined until 1995.

The diverse nature of porin primary structures, along with their hydrophilic nature, makes the assignment of a novel sequence as a porin problematic (see **Homology modeling**). The hydrophobic transmembrane b-barrel is composed of b-strands whose lengths range from 6 to 17 residues. Moreover, only every second residue of each strand faces the lipid core of the membrane and is consistently hydrophobic. Owing to the fluid nature of the bilayer core, only the hydrophobic nature of the lipid-facing residues needs to be conserved. The rest of the residues in the barrel can be either hydrophilic or hydrophobic, depending on whether they are exposed to the aqueous pore or buried in the protein structure. As with helical and monotopic membrane proteins, the strands are rich in aromatic residues in the region that interacts with the lipid headgroups (see [a-Helix](#)). A [protein structure prediction](#) method based on the detection of these interfacial aromatic residues and on the “every second residue hydrophobic” pattern correctly assigned 16 out of the 18 strands in maltoporin (9).

Because porins reside in the outer bacterial membrane, they must first be translocated through the bacterial inner membrane to the **periplasm**. To this end, the porin sequences carry an *N*-terminal **signal sequence**. The subsequent folding pathway is not known with certainty, but it has been proposed that the folding of porins begins in the aqueous periplasm by the formation of the trimer interface (2). The interface resembles the hydrophobic core of soluble proteins and could, therefore, form spontaneously before insertion into the membrane. The lipid phase would then induce the formation of the three b-barrels, because the barrel both satisfies the hydrogen-bonding requirement of the peptide groups and uses the [hydrophobic effect](#) in the interaction with the oily core of the bilayer. A trimeric porin that is indistinguishable from the natural molecules has been produced by heterologous expression by [protein engineering](#), followed by [protein folding in vitro](#) from [inclusion bodies](#) in the presence of detergents.

Much of the biological interest in porins stems from the observation that in pathogens porins often are the **antigens** against which the host's **antibodies** are directed. Thus, porins from pathogenic bacteria could be used to raise antibodies and make vaccines. Maltoporin of *E. coli* is also called LamB, because it is the receptor for [lambda phage](#).

## Bibliography

1. T. Schirmer (1998) General and specific porins from bacterial outer membranes. *J. Struct. Biol.* **121**, 101–109.
2. G. E. Schulz (1996) Porins: general to specific, native to engineered passive pores. *Curr. Op. Struct. Biol.* **6**, 485–490.
3. C. A. Mannella (1997) On the structure and gating mechanism of the mitochondrial channel, VDAC. *J. Bioenerg. Biomembr.* **29**, 525–531.
4. L. Song, M. R. Hobaugh, C. Shustak, S. Cheley, H. Bayley, and E. J. Gouaux (1996) Structure of staphylococcal  $\alpha$ -hemolysin, a heptameric transmembrane pore. *Science* **274**, 1859–1866.
5. R. Dutzler, Y-F. Wang, P. J. Rizkallah, J. P. Rosenbusch, and T. Schirmer (1996) Crystal structures of various maltooligosaccharides bound to maltoporin reveal a specific sugar translocation pathway. *Structure* **4**, 127–134.
6. J. E. W. Meyer, M. Hofnung, and G. E. Schulz (1997) Structure of maltoporin from *Salmonella typhimurium* ligated with a nitrophenyl-maltotrioxide. *J. Mol. Biol.* **266**, 761–775.
7. Y.-F. Wang, R. Dutzler, P. J. Rizkallah, J. P. Rosenbusch, and T. Schirmer (1997) Channel specificity: structural basis for sugar discrimination and differential flux rates in maltoporin. *J. Mol. Biol.* **272**, 56–63.
8. D. Forst, W. Welte, T. Wacker, and K. Diederichs (1998) Structure of the sucrose-specific porin ScrY from *Salmonella typhimurium* and its complex with sucrose. *Nature Struct. Biol.* **5**, 37–46.
9. T. Schirmer and S. W. Cowan (1993) Prediction of membrane-spanning  $\beta$ -strands and its application to maltoporin. *Protein. Sci.* **2**, 1361–1363.

## Position Effect

Early studies by cytologists led to the realization that some **chromosomal** regions have properties distinct from the rest of the chromosome. Large segments of **chromatin** are highly condensed and replicate late in **S-phase**. Geneticists determined that these chromosomal regions, which they called **heterochromatin**, did not participate in meiotic **recombination**. However, heterochromatin does have significant genetic effects. The most common observed consequence of heterochromatin formation is repression of **transcription** in the heterochromatin itself and in regions of chromatin adjacent to the heterochromatin domain. The variability in **gene expression** at the border of the heterochromatin is described as “position effect variegation”. The functional organization of the chromosome into discrete **domains** has been increasingly recognized through experiments in *Drosophila* that use the phenomenon of position effect. These experiments employ a powerful combination of techniques including genetic analysis and cytological observation of the large **polytene chromosomes**. The introduction of a normally active gene into a chromosome at a position adjacent to a transcriptionally inactive, condensed heterochromatin domain leads to significant repression of the transcription process. This is believed to occur through spreading of the heterochromatin structure into the normally active gene (see **Heterochromatin**). Thus, the expression of a gene depends on its chromosomal position, and hence the term position effect. This phenomenon provides a useful screen for genes and their products that suppress or enhance the repression of transcription because of position effects on a suitable **reporter gene**, often one influencing *Drosophila* eye coloration. Many of the genes modifying the position effect have been characterized and encode structural components of chromatin itself or are capable of modifying the organization of chromatin through enzymatic mechanisms.

The general properties of “modifier” genes that enhance or suppress position effect variegation provide some insight into how particular proteins might work to change chromosomal structure. The **phenotypes** of modifier mutations strongly depend on the number of gene copies within the cell. This has led to the proposal that the modifier proteins that assemble heterochromatin act through simple mass action. The more modifier protein in the nucleus, the more repressive heterochromatin is assembled. It is probable that certain modifier proteins cooperate to assemble multimeric complexes that alter the structure of entire chromosomal domains (1). The distribution of Polycomb along the chromosomal domain including the **bithorax complex**, which encodes three **homeotic genes**, supports this hypothesis. Orlando and Paro (2) found that the Polycomb protein is associated with transcriptionally inactive chromatin over more than 200 kbp of DNA, whereas it is absent from regions of genetic activity. The boundaries of Polycomb-associated chromatin and Polycomb-free chromatin are very distinct, which implies that regulatory elements exist that initiate and terminate Polycomb binding (2).

An important aspect of the functional role of the Polycomb proteins is that they are not involved in establishing the expression state of a particular gene but maintain the repressed state through **DNA replication** and chromosomal duplication. Because they act through large multicomponent complexes, it is possible that they might subdivide through a replicative event, thereby maintaining a repressive chromatin structure. Specialized mechanisms exist within the chromosome to enable the removal of Polycomb from chromatin and the resetting of chromatin to a transcriptionally competent state (3). The molecular machines involved in this process include the Brahma protein, which is structurally related to a component of the yeast general activator complex SNF2/SWI2. Evidence for competition between transcriptional activators and repressors, such as Polycomb, comes from experiments in which Polycomb-mediated transcriptional silencing is overcome by the expression of high levels of the transcriptional activator Gal4. **Immunostaining** of **polytene chromosomes** revealed that the Polycomb complex is removed from the vicinity of the Gal4 binding site without

the need for DNA replication and **cell division**.

Gottschling and colleagues (4) successfully established that position effect variegation also occurs in the chromosomes of *S. cerevisiae*. When a gene is located near a [telomere](#), its transcriptional activity is reduced (4). Transcriptional repression at yeast telomeres is caused by the assembly of a distinct chromatin structure that initiates at the telomere. The sequence of the DNA is important, because internal tracts of telomeric DNA also act as silencers of transcription in *S. cerevisiae*. Mutations in the amino-terminal tails of histones H3 and H4 relieve transcriptional silencing. Histone H4 at the telomeres is hypoacetylated, and mutations in the acetyltransferases that **acetylate** the amino terminus also relieve silencing. DNA **methyltransferases** expressed in yeast have more restricted access to telomeric chromatin than to most of the chromosome. This implies that telomeric chromatin either is more compacted or is sequestered away from freely diffusible **trans-acting** factors in the nucleoplasm. Several other similarities also exist between yeast and *Drosophila* heterochromatin. In both cases, genes within heterochromatin are normally maintained in a stably repressed state but occasionally escape repression (5).

The efficiency of transcriptional repression decreases when a gene is placed further from the *S. cerevisiae* telomere (over a 10- to 20-kbp range). This result supports the idea that a repressive chromatin structure originates at the telomere and spreads along the chromosome. The expression of protein SIR3 influences the extent of silencing at the telomeres, suggesting that it is essential for assembling repressive chromatin (6). SIR3 and the histone tails are also involved in sequestering yeast chromosomal telomeres at the periphery of the nucleus. Thus telomeric silencing could be related to the sequestration of this portion of the chromosome within a transcriptionally incompetent compartment of the nucleus adjacent to the nuclear envelope (7).

Many of the modifiers of the position effect in yeast chromosomes are shared between the telomeres of the chromosomes and the silent **mating-type** cassettes in yeast. Mutations in four silent information regulator genes, SIR1, 2, 3, and 4, cause derepression of the silent mating-type cassettes. Mutations in three of these, SIR2, 3, and 4, influence telomeric repression. The assembly of specialized chromosomal domains at the silent mating-type cassettes is reflected in the more limited access of [restriction enzymes](#) to DNA at these sites than at bulk chromatin. The silencers at the mating-type cassettes are required to reestablish the transcriptionally repressed state after each round of replication. This indicates that repressive chromatin structures do not self-template the reassembly of a repressive structure. It is possible that both the telomeres and the mating-type cassettes are sequestered in specialized nuclear structures or compartments in which the transcriptional machinery does not function efficiently.

Recent results also implicate heterochromatin-mediated silencing mechanisms in the [nucleolus](#), the site of ribosomal RNA transcription. SIR2, SIR4, H2A, and H2B influence the transcription of genes requiring RNA polymerase II that are integrated into ribosomal DNA. Remarkably, silencing also controls aging in yeast. Interference with the targeting of silencing complexes by expression of mutant forms of SIR4 allows yeast to live longer (8). Nucleolar organization is implicated in this phenomenon, because mutant SIR4 proteins localize to the nucleolus. In the fission yeast *Schizosaccharomyces pombe*, the determinants of gene silencing and heterochromatin assembly at the silent mating-type cassettes are shared with the centromeres (9).

The existence of these diverse sites of heterochromatin-mediated silencing has led to the idea that competition exists to limit components and that the telomere might act as a molecular sink to form a reservoir of silencing factors. Alternatively, individual targeting factors, such as SIR1 at the mating-type loci, might modulate the stability of heterochromatin assembly.

Position effects in mammalian chromosomes have been a recurrent problem for **transgenic** research, because highly variable levels of transcriptional activity follow from the random introduction of reporter genes into the genome (10, 11). These effects are relieved by introducing a **locus control region** (LCR) that exerts a dominant transcriptional activation function over a chromatin domain (10)

to 100 kbp). The mechanism of this activation function remains to be determined. Communication between LCRs, **enhancers** and **promoters**, however, either directly or through modifications of chromatin structural components, are favored hypotheses. Recent evidence suggests that both the locus control regions and enhancers act in cis to suppress position-effect variegation actively. In this regard, locus control regions function as operationally defined “powerful” enhancers. The coexistence of heterochromatin domains that transmit repressive effects and the definition of the extensive long-range activation function of LCRs emphasize the necessary compartmentalization of the chromosome into discrete functional units. These discrete functional units are prevented from influencing each other in a natural chromosomal context because of, in part, the existence of special chromosomal regions that prevent the transmission of chromatin structural features associated with the boundaries of repressive or active domains. These specialized chromosomal regions are known as insulators (see [Domain, Chromosomal](#)).

The original evidence for an insulator function within the chromosome came from genetic experiments in *Drosophila*. Each boundary of the 87A7 **heat-shock** locus is defined by a pair of nuclease **hypersensitive** sites bordering a 250- to 300-bp segment of DNA. These specialized chromatin structures (scs) are located at the junctions between the decondensed chromatin of the transcriptionally active 87A7 heat-shock locus and adjacent condensed chromatin. The scs elements have three functional properties: (1) they establish a domain of independent genetic activity at many distinct chromosomal positions; (2) scs elements are necessary at each edge of the domain; and (3) independently the elements are neither inhibitory nor stimulatory to transcriptional activity within the domain ([12](#)). Subsequent work found that introducing an scs element between an enhancer and a promoter blocks communication between the two regulatory DNA sequences. Thus, the scs elements prevent both the transmission of repressive effects on transcription from proximity to heterochromatin and the transmission of stimulatory effects on transcription from an enhancer. How this insulation is achieved is unknown. Moreover, the nature of the nucleoprotein complex assembled on the scs elements has yet to be defined. Fortunately, similar phenomena have been described that are associated with a well-defined nucleoprotein complex between the *Drosophila* suppressor of hairy-wing protein [su(Hw)] and the gypsy retrotransposon ([13](#)).

Insertion of a gypsy element as far as 10 to 30 kbp away from a promoter causes a mutant phenotype. The mutant phenotype requires the su(Hw) protein to interact with the inserted gypsy element at a 350-bp region containing 12 copies of a 10-bp sequence separated by AT-rich sequences. The su(Hw) protein has a molecular weight of 100 kDa, and its sequence includes several motifs characteristic of eukaryotic transcription factors, including 12 **zinc fingers**, a **leucine zipper** and two acidic **domains**. The complex of the su(Hw) protein with a gypsy element has many of the properties of an insulator element. The complex blocks enhancer activity when placed between an enhancer and a promoter, and when the complex is placed at the boundaries of a gene-containing fragment, the gene is protected from the repressive effects of heterochromatin on transcription ([14](#)).

Although in certain circumstances the su(Hw) protein does not independently stimulate or repress transcription of a reporter gene, the su(Hw) protein occasionally functions as a transcriptional activator. This suggests that the function of an insulator may be conferred on sequences by [DNA-binding proteins](#) that under other circumstances might have more conventional roles in the transcription process. Although it is clear that the su(Hw) protein does not bind to scs elements, it seems likely that these elements form large nucleoprotein complexes with a similar composition.

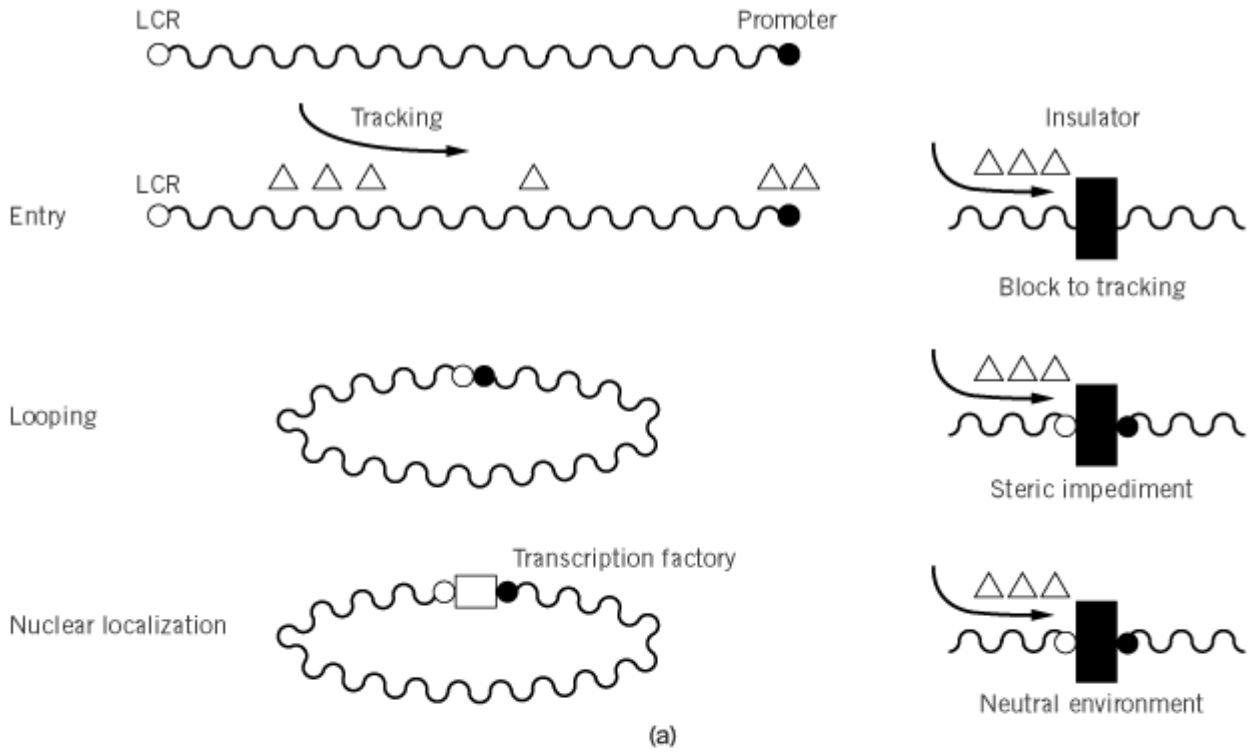
How might insulators function? Any models must explain how characteristics of both repressive or active chromatin are restricted to particular chromosomal domains. Several models have been suggested to explain the activities of LCRs and enhancers (Fig. [1](#)). These elements might function as entry points for transcription factors, RNA polymerase, or other components of the transcription machinery, which then might track along the DNA until reaching the promoter. Alternatively, the LCR or enhancer complex might associate with the promoter complex by stable looping of the intervening DNA or chromatin, forming a complex that increases the efficiency with which RNA polymerase is recruited and used. Another possibility is that the LCR or enhancer complex might



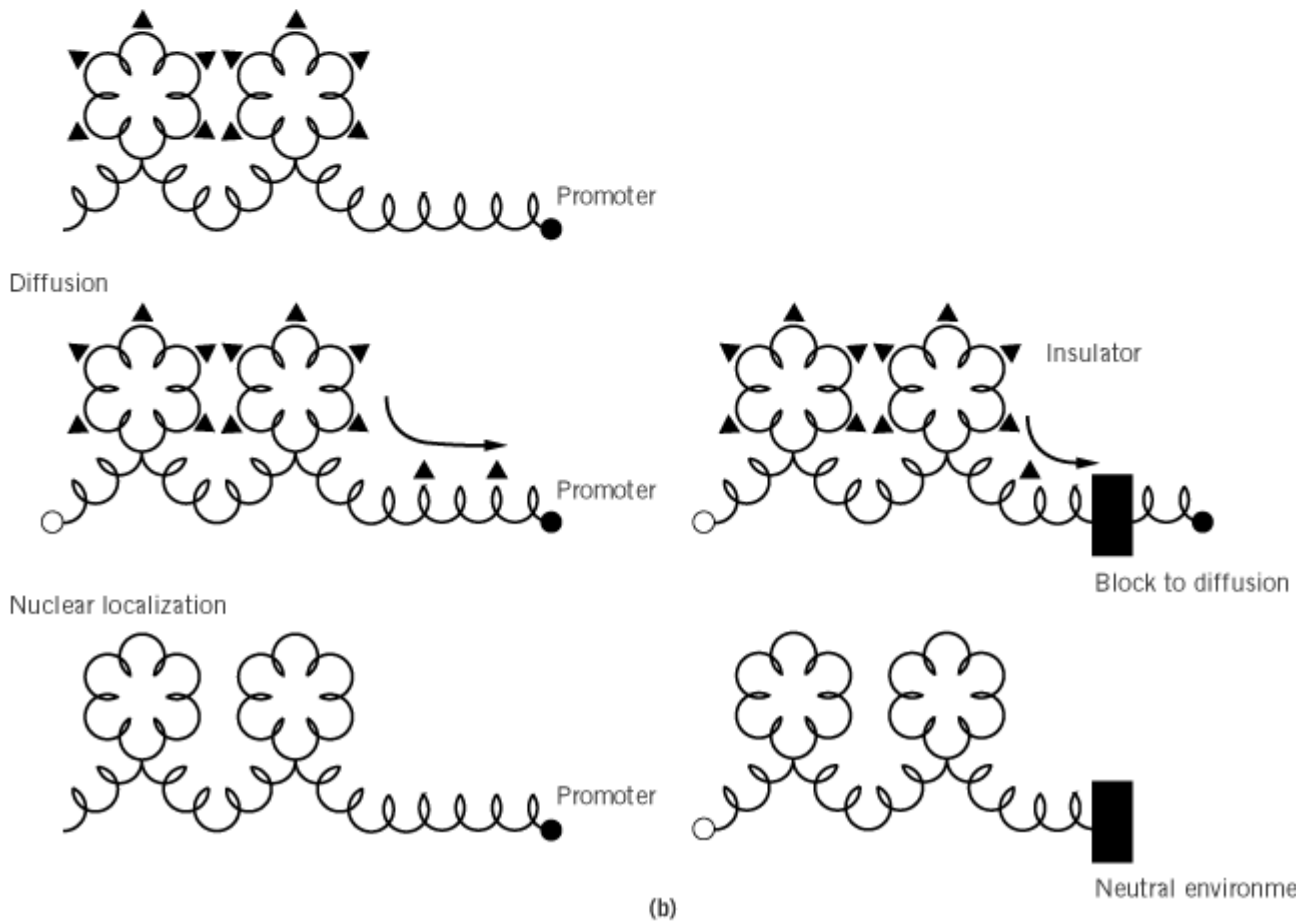
cause the gene to assemble into a chromatin structure capable of being transcribed through its association with nuclear compartments (or **organelles**) that act as transcription factories or associate with proteins that modify repressive chromatin structure by disrupting nucleosomes. Similar models can explain the repressive effects of heterochromatin. Repressive chromatin proteins, such as the *Drosophila* HP1 protein or histone deacetylases, may undergo local diffusion to adjacent DNA sequences, enlarging the heterochromatin domain. Alternatively, heterochromatin may be sequestered in a transcriptionally incompetent region of the nucleus.

**Figure 1.** Models to explain the actions of insulators on the activities of locus control regions/enhancers and heterochromatin. (a) Enhancers could communicate with promoters by acting as entry sites for processive enzymes that track along DNA (entry), they could make direct contact with the promoter to enhance function by looping out intervening DNA (looping), or they could direct the chromatin domain containing the promoter to a compartment of the nucleus competent for transcription (nuclear localization). Insulators could block each of these functions as indicated. (b) Proximity to heterochromatin leads to position effect variegation on adjacent genes. This could be caused by simple diffusion of repressive heterochromatin proteins to otherwise active genes in euchromatin (diffusion), or it could be caused by inappropriate nuclear localization in an environment that is not competent for transcription (nuclear localization). Insulators could block each of these functions as indicated.

Insulating activation



Insulating repression



In considering these models, it is important to recognize that the eukaryotic nucleus is a highly organized structure in which DNA is compacted into nucleosomes and the chromatin fiber by its association with histone proteins. Although it is possible that insulators prevent protein tracking or diffusion between active and repressive domains, it is difficult to envisage how this might occur in the nucleus—where DNA segments separated linearly by many kilobases can be juxtaposed by the folding of the DNA helical axis in three dimensions—without invoking some specific attachment of inactive chromatin domains, insulators, and active chromatin domains to a nuclear framework. Similar attachments might be required to prevent the juxtaposition, of LCRs, enhancers, and promoters as a result of DNA looping. Perhaps the most economical suggestion is that insulators are nucleoprotein complexes that associate neither with regions nor structures in the nucleus where “transcription factories” load on to DNA nor with regions or structures from which the transcriptional machinery is excluded (15, 16). Instead, the insulators might associate with distinct “neutral” nuclear structures. The “neutral” nuclear structures would tether promoter elements where the transmissible activating effects of enhancers or silencing effects of heterochromatin could not occur. This absence of transmissible effects could be accounted for by the exclusion of particular transcriptional coactivators or corepressor from the “neutral” nuclear structures.

### Bibliography

1. A. Moehrle and R. Paro (1994) *Dev. Genet.* **15**, 478–484.
2. V. Orlando and R. Paro (1993) *Cell* **75**, 1187–1198.
3. J. W. Tamkun et al. (1992) *Cell* **68**, 561–572.
4. D. E. Gottschling, O. M. Aparicio, B. L. Billington, and V. A. Zakian (1990) *Cell* **63**, 751–762.
5. S. Henikoff (1990) *Trends Genet.* **6**, 422–426.
6. H. Renauld et al. (1993) *Genes Dev.* **7**, 1133–1145.
7. W. W. Franke (1974) *Int. Rev. Cytol. Suppl.* **4**, 71–236.
8. B. K. Kennedy et al. (1997) *Cell* **89**, 381–391.
9. R. C. Allshire, J. P. Javerzat, N. J. Redhead, and G. Cranston (1994) *Cell* **76**, 157–169.
10. F. Grosveld, G. B. van Assendelft, D. R. Greaves, and G. Kollias (1987) *Cell* **51**, 975–985.
11. S. Fiering et al. (1995) *Genes Dev.* **9**, 2203–2213.
12. R. Kellum and P. Schedl (1991) *Cell* **64**, 941–950.
13. V. G. Corces and P. K. Geyer (1991) *Trends Genet.* **7**, 69–73.
14. R. R. Roseman, V. Pirrotta, and P. K. Geyer (1993) *EMBO J.* **12**, 435–442.
15. D. A. Jackson, A. B. Hassan, R. J. Errington, and P. R. Cook (1993) *EMBO J.* **12**, 1059–1065.
16. F. Palladino et al. (1993) *Cell* **75**, 543–555.

### Suggestion for Further Reading

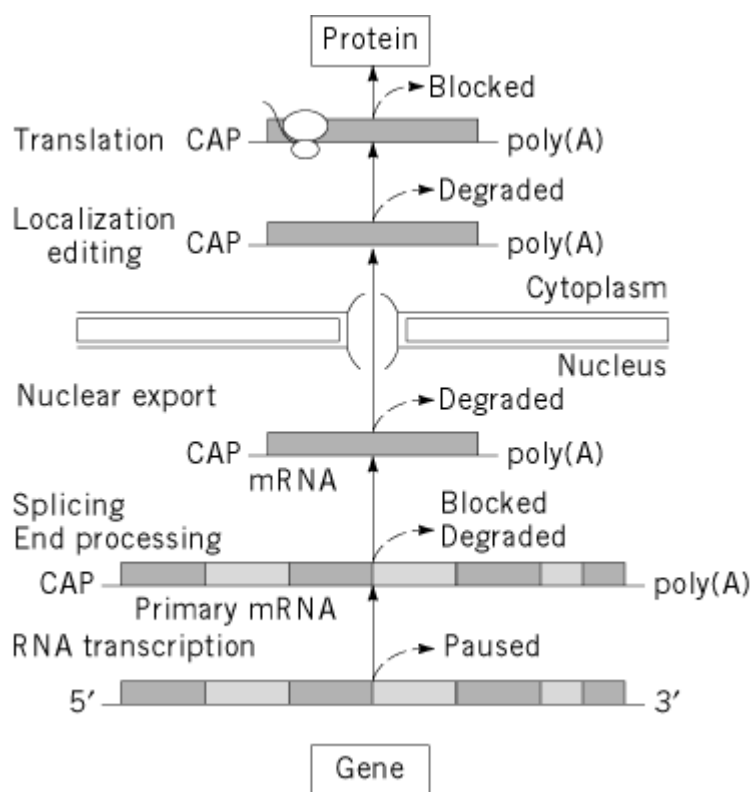
17. A. Wolffe (1998) *Chromatin: Structure and Function*, 3rd ed., Academic Press, London.

### Post-Transcriptional Regulation

The synthesis of RNA in [transcription](#) is a selective process. Eukaryotic cells contain three types of **RNA polymerase** that transcribe DNA into RNA, whereas bacteria contain only one. The **enzyme** RNA polymerase II converts about 1% of eukaryotic cell DNA into [messenger RNA](#) (mRNA), the

template that codes for [translation](#) of new [proteins](#) in **protein biosynthesis**. Cells adopt many strategies to decide how and when a given **gene** should be expressed. Within eukaryotic cells there are numerous post-transcriptional steps available for regulating protein production, commencing with transcription of the primary mRNA in the [nucleus](#) and culminating with its eventual translation into protein in the cytoplasm (see Fig. 1). Newly synthesized precursor mRNAs, or primary mRNA transcripts, are usually quite large (10–20 kb), due to the presence of **intron** sequences, which are subsequently removed by [RNA splicing](#). The 5' and 3' ends of the primary transcript require modification in the nucleus, and only then is the mature mRNA transported to the cytoplasm, where it is directed to the desired location and translated into protein. Therefore, several steps are involved in the conversion of a primary mRNA transcript into a [polypeptide chain](#), and each level provides a locus for cellular fine-tuning (see Fig. 1).

**Figure 1.** The cellular pathway of protein synthesis: from gene transcription to mRNA translation. The major processes in this pathway are shown on the left side of the figure. Some of the potential effects of regulating these events are indicated by broken arrows. The gene-coding exons (dark gray) and introns (light gray) are as shown.



## 1. In the Nucleus

The process of protein biosynthesis comprises several discrete steps, some of which are likely to be integrated. The nucleus is a dynamic cellular structure, and the transcription of genes into mRNA occurs at specific sites in the nucleus, during which the growing mRNA chain becomes complexed with many different proteins.

### 1.1. Transcriptional Attenuation

Some genes in bacteria and eukaryotes are blocked at the level of transcriptional **attenuation**, where the transcribing polymerase pauses at a specific site on the DNA template. This event can regulate the rate of expression of certain genes, such as the [c-myc oncogene](#) in differentiated leukemia cells (1).

## 1.2. Capping

Primary mRNA transcripts (but not other RNAs) are modified at the 5' end by the addition of a 7-methylguanosine [cap](#). This “capping” reaction occurs almost immediately after the RNA chain is initiated, and the modification is required for subsequent events, which include the splicing, nuclear export, stability, and translation of each mRNA. The cap structure forms a stable complex with nuclear proteins that mediate some of these nuclear processes (2). Once the mature mRNA reaches the cytoplasm, the mRNA cap is bound by a different set of proteins that help facilitate translation.

## 1.3. 3'-End Processing and Polyadenylation

The primary mRNA transcripts of eukaryotes undergo more complex and regulatable processing at the 3' end than in bacteria, and they contain two types of cleavage signal (AAUAAA and a less well-conserved GU-rich element nearby). This cleavage/[polyadenylation](#) sequence is located in the 3'-untranslated region (UTR), the noncoding end of the mRNA that follows the translation **stop codon**. Once the newly formed primary transcript is cleaved, a poly(A) polymerase enzyme adds a polyadenylate tail of about 100–200 A residues to the 3' terminus of the RNA (3). Specific binding proteins associate with the poly(A) tail and are capable of enhancing the stability, translation, and nuclear export of the mature mRNA (4). The degree of polyadenylation can be regulated in the nucleus by splicing factors, such as U1A snRNP (3). Many genes contain multiple cleavage/polyadenylation sites, giving rise to different length transcripts in certain types of cell or tissue.

The combination of 5' cap and 3' poly-A tail identify the mature mRNA as a complete coding unit, and these end modifications are required for later translation of the mRNA template in the cytoplasm. The **histone** genes represent the only known exception. Histone mRNAs are synthesized without introns and are not polyadenylated; they contain a 3'-end stem-loop structure that mediates an alternate processing pathway for this mRNA (5).

## 1.4. Splicing

In higher eukaryotic cells, the growing primary transcript forms large complexes with different proteins, many of which include the abundant heterogeneous nuclear **ribonucleoprotein** particle (hnRNP) proteins. At the junctions between intron and exon, a [spliceosome](#) (RNA splicing complex) develops on the pre-mRNA and is responsible for removal of intron sequences from the immature transcript (6). The spliceosome machinery comprises many proteins, often complexed with small stable RNAs (6). Completion of the splicing process is required for continued processing of the primary transcript, and eventual export of the mature mRNA to the cytoplasm.

## 1.5. Alternative RNA Splicing

The process of pre-mRNA splicing involves not only distinguishing the junctions between exon-coding sequences and the intervening introns but also excising unwanted introns and religating the correct exon termini. The detailed orchestration of this procedure also provides a means for altering the gene product, essentially by religating different exons together. This is known as [alternative splicing](#), and it provides a potent mechanism for generating not only different mRNAs but also different polypeptide chains, from a single gene. Some genes contain ambiguous splice sites, and as a result generate several alternative gene products in most cell types and tissues. The regulation of alternative splicing is observed more frequently during the tissue-specific expression of genes, and during [development](#) of higher eukaryotes (7). Alternative splicing can be effected by specific exon or intron sequences, examples of which occur in the genes for human [growth hormone](#), cardiac troponin T, [fibronectin](#), and the nerve-cell-specific form of the *c-src* **oncogene**. The sex-determination genes of *Drosophila melanogaster* (sex-lethal, transformer and doublesex) currently represent the best characterized regulatory system of alternative splicing. In this example, both the pre-mRNA target sequences and the regulatory **RNA-binding** factors that block or activate splice site selection are well defined (6, 7).

## 1.6. Pre-mRNA Stability

Most RNA synthesized in the nucleus is degraded, and inappropriate splicing or end modification may sentence a primary mRNA to the same fate (see [RNA Degradation In Vitro](#)). The RNA-degrading enzymes are **ribonucleases**, and they operate in both bacterial and eukaryotic cells. The expression of many genes is, in fact, determined by regulated degradation of the primary nuclear transcript. This type of regulation can sometimes correlate with disease. For example, the gene encoding a cancer-associated proteinase, [urokinase](#), is efficiently transcribed in both malignant and benign cancer cells derived from rat mammary tumors, but the **proteinase** is not well expressed in the nonmalignant cells, due largely to nuclear degradation of the primary transcript (8). Certain examples of **nonsense codon**-mediated mRNA decay, signaled by the presence of a premature stop codon, also occur in the nucleus. Although progress in this area has been slow, it seems likely that regulated nuclear RNA degradation, like the constitutive breakdown of excised intronic RNA, is somehow linked to RNA splicing and modulated by RNA-binding proteins, such as the hnRNPs (9).

### 1.7. RNA Export

Once spliced and processed, the mature mRNA transcript travels to the cytoplasm through the [nuclear pore complex](#). The mRNA molecules do not exit the nucleus unescorted; they are chaperoned by proteins, some of which contain nuclear export signals (see [Nuclear Import, Export](#)). The mRNA translocates through the pores as part of a large ribonucleoprotein complex, of which some nuclear proteins are later replaced by cytoplasmic factors, to facilitate further movement of the mRNA or its translation into protein. The most common proteins to accompany exported mRNAs include the cap-binding complex proteins CBP20 and CBP80 and the hnRNPA1 shuttling protein (10). The cap-binding complex proteins also mediate U **small nuclear RNA** (snRNA) export in mammals and yeast. Different types of RNA, such as [transfer RNA](#) and **ribosomal RNA**, lack a cap structure, and their export requires interaction with specific carrier proteins. TFIIA and **La** proteins are implicated in export of 5 S rRNA. The tRNA export receptor (termed exportin-t) was recently identified (11).

There is growing evidence that different types of mRNAs interact with very specific regulatory proteins, which either impede or enhance the rate of nuclear export. **Viruses** provide the best studied examples of regulated mRNA export. For example, unspliced intron-containing **HIV-1** pre-mRNA is effectively exported, while cellular host pre-mRNAs remain anchored in the nucleus by the spliceosome. The HIV-1 pre-mRNA contains specific signals that redirect its nuclear location and target binding of the HIV Rev protein, which facilitates rapid nuclear export of the viral pre-mRNA. It is anticipated that several cellular mRNA-binding carrier proteins, containing export signals analogous to that of the Rev protein, will soon be identified. One potential candidate is HuR, a nuclear-cytoplasmic shuttling protein that binds to AU-rich mRNA instability elements, which may possibly enhance both the nuclear transport and stability of specific mRNAs (12). There is, in addition, evidence for vectorial transport of RNAs from the nucleus. The transcripts of the *Drosophila* **pair-rule gene** family (eg, even-skipped, hairy, and runt) exit from the apical side of the nucleus, and this directionality is dependent on sequences in the 3'-untranslated region of the RNAs (13).

## 2. In the Cytoplasm

### 2.1. mRNA Localization

In the cytoplasm, eukaryotic mRNAs interact with [ribosomes](#) that continuously translate the mRNA sequence into specific polypeptide chains. Genes that encode secreted proteins contain a unique sequence that, when translated at the amino terminus of the growing polypeptide chain, directs the translation complex to the [endoplasmic reticulum](#), a **polyribosome**-coated organelle that surrounds the nuclear envelope. Some processed mRNAs contain sequences at the 3' end that signal movement of the RNA, prior to translation, to defined locations in the cytoplasm. In this way, simply by relocating the mRNA, the cell is able to direct efficiently the site-specific accumulation of certain proteins.

Some mRNAs appear to be routed throughout the cell, not by passive [diffusion](#), but through

association with proteins that interact with the cell's [cytoskeleton](#). There is evidence in *D. melanogaster* oocytes that the mRNAs encoding the developmental proteins [bicoid](#) and [oskar](#) are transported in association with Staufen protein along microtubules to the anterior and posterior poles of the oocyte (13). Again, microtubule association is inferred for movement to the vegetal pole of Vg1 mRNA in *Xenopus laevis* oocytes. While the muscle [b-actin](#) mRNA appears to travel along microtubules in nerve cells, it associates only with actin filaments in fibroblast cells (13). Interestingly, the localization of b-actin is responsive to [growth factors](#), and b-actin mRNA in chicken fibroblasts moves rapidly to the leading edge of the cell, where it is translated into actin filaments that influence cell motility. The proteins that bind to the actin and Vg1 mRNAs contain potential nuclear export signals, and they possibly chaperone the transcripts from the site of transcription/splicing in the nucleus, to that of translation in the cytoplasm. Other examples of localized mRNAs have been identified, and current study is focused on understanding the molecular mechanisms responsible for cytoplasmic mRNA movement.

## 2.2. mRNA Translation

The translation process itself can be influenced, most often by mRNA-binding proteins. The mRNA signal sequences that define the AUG [initiation codon](#) differ between bacteria and eukaryotic cells, although in all cases the mRNA is scanned by ribosomes in a 5'–3' direction, and then translated into protein. In eukaryotes, translation initiation is dependent on the 5'-cap structure, and the translation process is subject to several levels of cellular control. mRNAs are translated in the cytoplasm by free or membrane-bound **polysomes**. The 5' cap is recognized by the **initiation factor** eIF-4E, which recruits other proteins to form a functional cap-binding complex. The cap-binding complex, in turn, targets binding of the small 40 S ribosomal subunit (as part of a 43 S complex with additional factors). The resulting pre-initiation complex then scans along the mRNA, and following recognition of the initiation codon and binding of the 60 S ribosomal subunit, the mRNA open reading frame is translated (14). The cell may control translation by modifying some of the general initiation factors, or by shifting the balance of these factors in the cytoplasm.

Translation initiation of specific mRNAs can be regulated by RNA stem-loop structures near the 5' cap. There are a few examples both in eukaryotes and bacteria, but the best studied case is the post-transcriptional regulation of vertebrate iron homeostasis. This is mediated by specific RNA stem loops near the 5' cap of **ferritin** mRNAs and some other transcripts. When iron levels are low, iron regulatory proteins 1 and 2 bind the stem-loops and inhibit translation, presumably by interfering sterically with formation of the 43 S pre-initiation complex (14). High levels of intracellular iron relieve the translational block. Translation control signals can also occur at the 3' end of mRNA. While all mRNAs undergo constitutive polyadenylation in the nucleus, the poly(A) tails of certain mRNAs are modified in the cytoplasm during early development in vertebrates, amphibians, and insects. Specific enzymes lengthen or shorten the poly(A) tail in the cytoplasm; a longer poly(A) tail (150–200 nucleotides) generally enhances translation of maternal messages during oocyte maturation and fertilization, whereas shorter poly(A) tails of less than 50 nucleotides block translation (14).

mRNA translation control signals have also been identified between the 5' cap and 3' poly(A) tail. Such elements were observed in the 3' UTR of the 15-lipoxygenase mRNA and the mouse protamine 2 gene; both RNA signals bind to proteins that regulate translation. Other mRNAs contain similar regulatory elements within the actual protein coding sequence. An interesting example of this type is [thymidylate synthase](#), which also represents the best defined case of translational autoregulation. Thymidylate synthase is an enzyme involved in DNA precursor biosynthesis, and it doubles in function as an RNA regulator, capable of binding and regulating translation of its own mRNA, in addition to mRNAs encoded by other genes (15). These examples are likely to represent a form of cross-talk between the 5'-bound pre-initiation complex and the distant site of regulation, but the mechanisms involved are not yet defined.

Finally, additional methods of translation control include signal-mediated pausing or stalling of the ribosomes, and ribosomal [frameshifting](#), in which ribosome translocation is shifted one base out of

frame, thus producing a protein of different sequence. There are only a few examples of pausing and frameshifting, and the latter event is used primarily by eukaryotic and bacterial viruses.

### 2.3. mRNA Stability

Bacterial cells divide much faster than eukaryotic cells, and bacterial mRNAs are degraded at correspondingly faster rates. Most bacterial mRNAs are stable only for a matter of minutes, while some eukaryote mRNAs are stable for many hours. The mRNAs encoding **stress-response** or growth-responsive proteins, such as [transcription factors](#) and [growth factors](#), are often degraded with minutes. Cellular stresses such as viral infection, growth-factor signaling, hypoxia, and other environmental changes can decrease, or extend, the rate of mRNA degradation.

In eukaryotes, the rate of cytoplasmic mRNA decay is regulated by mRNA-binding proteins. In addition to the poly(A) tail, control elements have been identified in the mRNA coding sequence and in the 3' UTR. The interaction of specific proteins can protect the mRNA against attack from ribonucleases. Deadenylation, or shortening of the poly(A) tail, is an initial step in the decay of many mRNAs. The formation of a complex between the poly(A)-binding protein and the poly(A) tail stabilizes mammalian mRNAs, but it displays the opposite effect in yeast ([16](#)).

The 3'-untranslated region of many unstable growth factor and **oncogene** mRNAs contain instability signals, sequences rich in adenosine and uridine. These AU-rich elements fall into distinct classes, bind a variety of cellular proteins, and vary in their ability to modulate mRNA stability *in vitro*. Some of the AU-rich element binding proteins, such as the HuR protein, shuttle between nucleus and cytoplasm, and thereby they might associate with mRNAs soon after synthesis in the nucleus, acting as mRNA chaperone, as well as guardian against degradation. In other instances, the RNA stability signal is a hairpin or stem-loop structure, such as that occurring at the 3' end of histone mRNA, or the cluster of five [iron-response elements](#) (IREs) in the 3' UTR of [transferrin](#) receptor mRNA. When iron levels are low, the iron regulatory proteins bind to the IRE stem loops and protect transferrin receptor mRNA from degradation. This model for regulated mRNA stability is critical for maintaining iron homeostasis in the cell, and many other examples of such control are likely to occur.

There are several candidate ribonucleases specific for mRNA, but none has yet proved to play a major role in cytoplasmic mRNA degradation. General degradation of mRNAs occurs in a 3'–5' direction and often is preceded by deadenylation. Some mRNAs, like transferrin receptor, are instead cleaved at specific sites by endonucleases, thereby generating an unprotected 3' end susceptible to rapid exonuclease activity ([16](#)). Bound regulatory proteins are likely to modulate degradation by either specific endonucleases or general exonucleases.

### 2.4. mRNA Editing

Finally, various examples of [RNA Editing](#) (ie, specific nucleotide sequence alterations) have been identified in several organisms. One of the more intricate and extensive types of editing involves insertion and/or deletion of uridine residues in the mitochondrial RNAs of the trypanosome parasites. This type of editing is directed by a small class of RNA, termed **guide RNAs**, and results in a range of RNAs of varied sequence and reading frame. Additional types of editing have been observed in plants and higher eukaryotes. Several examples of single base alterations mediated by specific enzymes have been reported, leading to gene-specific changes in mRNA reading frame, or creation of new [stop codons](#). As in the case of alternative splicing, a cell can use RNA editing to diversify the end products derived from a single gene.

### Bibliography

1. D. L. Bentley and M. Groudine (1986) *Nature* **321**, 702–706.
2. E. Izaurralde, J. Lewis, C. McGuigan, M. Jankowska, E. Darzynkiewicz, and I. W. Mattaj (1994) *Cell* **78**, 657–668.



3. E. Wahle and W. Keller (1996) *Trends Biochem. Sci.* **21**, 247–250.
4. Y. Huang and G. G. Carmichael (1996) *Mol. Cell Biol.* **16**, 1534–1542.
5. W. F. Marzluff (1992) *Gene Expression* **2**, 93–97.
6. M. J. Moore, C. C. Query, and P. A. Sharp (1993) in *The RNA World*, R. F. Gesteland and J. F. Atkins, eds., Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 303–357.
7. M. McKeown (1992) *Annu. Rev. Cell Biol.* **8**, 133–155.
8. B. R. Henderson, W. P. Tansey, S. M. Phillips, I. A. Ramshaw, and R. F. Kefford (1992) *Cancer Res.* **52**, 2489–2496.
9. G. Dreyfuss, M. J. Matunis, S. Pinol-Roma, and C. G. Burd (1993) *Annu. Rev. Biochem.* **62**, 289–321.
10. E. Izaurralde and I. W. Mattaj (1995) *Cell* **81**, 153–159.
11. U. Kutay, G. Lipowsky, E. Izaurralde, F. R. Bischoff, P. Schwarzmaier, E. Hartmann, and D. Gorlich (1998) *Mol. Cell* **1**, 359–369.
12. X. C. Fan and J. A. Steitz (1998) *EMBO J.* **17**, 3448–3460.
13. D. St. Johnston (1995) *Cell* **81**, 161–170.
14. M. W. Hentze (1995) *Curr. Opin. Cell Biol.* **7**, 393–398.
15. E. Chu and C. J. Allegra (1994) *BioEssays* **18**, 191–198.
16. J. Ross (1996) *Trends Gen.* **12**, 171–175.

### Suggestions for Further Reading

17. J. E. Darnell (1982) Variety in the level of gene control in eukaryotes, *Nature* **297**, 365–371. (An earlier treatise on the predicted variability of post-transcriptional control mechanisms. This was a landmark review at the time, and remains an easy to read introduction to the field.)
18. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing, New York, pp. 223–240, 365–378, 453–468. (An excellent general text, both informative and well illustrated.)
19. J. E. G. McCarthy and H. Kollmus (1995) Cytoplasmic mRNA-protein interactions in eukaryotic gene expression, *Trends Biochem. Sci.* **20**, 191–197. (Summarizes many of the important RNA–protein interactions that mediate different aspects of post-transcriptional control.)

### Post-Translational Modifications

All naturally occurring [proteins](#) and [polypeptide chains](#) are primarily made up of the 20 [amino acids](#), 21 including [selenocysteine](#), dictated by the [genetic code](#). To achieve the myriad of detailed functions required to support life in its varied forms, including the essential processes of translocation, regulation, activation, and turnover, a large number of additional modifications are required. These changes are introduced into both the polypeptide [backbone](#) and the majority of the amino acid [side chains](#). For the most part, these occur either during the formation of the polypeptide chain during [translation](#) or at some point following chain termination. Thus, these alterations are labeled as either co- or post-translational. In some cases, however, reactions may occur at both stages, or at a time when that distinction is not clear. There are also a limited number of amino acid modifications that can occur after the formation of the charged [transfer RNA](#) species, but before

incorporation of the amino acid during translation, and it has been suggested that these be referred to as “pre-translational” to distinguish them from the co-/post-translational changes. The extent to which co-/post-translational modifications occur must await the complete characterization of all the proteins specified by the [genome](#) of any organism. Even excluding the actual **proteolysis** associated with **protein degradation**, it can be safely concluded that they will be very widespread, affecting nearly all proteins.

## 1. Function and Reversibility of Protein Modifications

Protein synthesis is predominantly a cytoplasmic process. A very small amount of synthesis occurs in organelles such as the **mitochondrion**, which reflects its endosymbiotic origin from prokaryotes, and hence can actually be thought of as being similar to cytosolic synthesis in these organisms. This is also manifested in many of the distinguishing characteristics of this process in mitochondria. Proteins must function in other cellular locations as well, however, and in these cases they must be translocated in an efficient fashion. Inability of the cell to place a protein correctly in its proper environment, despite the fact that it may otherwise be functional, is generally equivalent to a mutation or deletion that results in loss of function. Proteolytic events direct these translocation processes, as well as several additional changes that occur as a result of the transfer. In the case of both the [endoplasmic reticulum](#) (ER) and the mitochondrion, the translocation is accompanied by the transfer of the protein molecule across active [membranes](#) in an energy-dependent fashion, requiring the protein molecules to be in an unfolded state. [Molecular chaperones](#), proteins that assist in the folding/unfolding process, are essential to complete these transfers, although they themselves do not introduce any modifications directly. Transport of proteins out of the cytoplasm into other organelles, **vesicles**, or eventually to the extracellular milieu, is a vectorial process, requiring essentially irreversible covalent changes (1).

In addition to the modifications that can occur on proteins during and after translocation, numerous covalent changes also occur on proteins that remain in the cytoplasm, or that are eventually transferred to the nucleoplasm. This group of reactions is designed primarily to induce the formation of active structures or to modulate or regulate these activities, and they involve almost exclusively the modification of amino acid side chains. In contrast, proteins that function in the compartments of the endoplasmic reticulum, [Golgi apparatus](#) and *trans*-Golgi networks, including the extracellular space, are frequently subject to additional backbone processing to generate active molecules.

The distinction between reversible and irreversible modification is of major significance, because these two forms of modification are used to achieve distinct cellular goals. Covalent modifications that are designed to drive a process in a single direction, such as the transfer of a protein from the cytoplasm to the extracellular environment, are almost always accompanied by irreversible change. Similarly, covalent modifications that are meant to contribute to the destruction of the cell and its components are usually of a similar nature. In addition, modifications that induce constitutive activity also tend to be permanent modifications, and generally require protein turnover and new synthesis to regulate their levels of activity. In contrast, changes that fine-tune the activity of enzymes, but are seldom actually involved in the catalytic process, are often reversible in nature. These tend to be associated with rapid, short-range responses, while the former are related more to long-range effects. Since most chemical modifications of proteins are formed through enzymatic [catalysis](#), it follows that modification of a reversible type will have specific enzymes that catalyze both the addition and removal of the modifying group.

It should be emphasized that the reversibility of a chemical modification must be defined in a biological, not chemical, context. Thus, for example, the introduction of [disulfide bonds](#) into proteins is fundamentally irreversible under most physiological conditions, even though their cleavage by reduction or oxidation can be readily achieved, and often is, in the laboratory. *N*-Terminal acetylation, as a cotranslational event, is also basically an irreversible change because there are no deacetylase enzymes for removing this modification on peptides larger than two or three amino acids, despite the fact that the amide bond is readily hydrolyzed under a number of conditions.

Clearly, the irreversible nature of many structural modifications results from the lack of appropriate enzymes to catalyze the opposite reaction, rather than from an unfavorable thermodynamic situation.

It is also important to distinguish the controlled or regulated modifications that occur as the result of cellular function from those that happen in an uncontrolled manner. These can be both enzymatic and nonenzymatic in nature. Intracellular random proteolysis is an excellent example. Although it is unclear to what extent such reactions occur, the recent identification of a number of intracellular [proteases](#) (eg, [caspases](#)) suggests that this may be a more extensive phenomenon than has been previously recognized. Even inactive forms of these enzymes can have some latent activity, albeit considerably less than the mature enzyme. It must be presumed that there are viable scavenger pathways for removing the fragments that are generated by these cleavages. Nonenzymatic modifications also result from the generation of various species, such as free radicals, that can arise from a variety of environmental conditions. These reagents can induce a host of chemical modifications, some of which are deleterious, particularly in biological macromolecules. This is particularly true of proteins and nucleic acids, and certainly a significant amount of cellular and genetic damage, which may be a major contributor to many pathologies such as cancer and aging, results from these modifications. Although there are appropriate mechanisms for removing or reversing inappropriately generated derivatives of both proteins and nucleic acids, the principal means for removing altered protein molecules is degradation and resynthesis.

## 2. Non-Enzymatic Modifications

Although apparently nonenzymatic modifications of proteins may be relatively common, it is often difficult to verify that any given modification is truly nonenzymatic under physiological conditions. For example, cyclization of amino-terminal glutamine residues to pyrrolidone carboxylic acid is well known to occur readily in solution under mild conditions. Thus, it might be logically concluded that this modification would occur and be found to a significant extent in all proteins with this sequence. There is, however, substantial evidence that there is a cyclase enzyme for carrying out this modification, at least with some substrates.

Many of the nonenzymatic modifications involve cleavage of the polypeptide backbone, including transpeptidation (see text below) and the modification of the **α-amino group** by reducing compounds, most commonly, sugars. Other modifications that occur nonenzymatically include the deamidation of [asparagine](#) and [glutamine](#) residues; the introduction of side-chain **crosslinks** to stabilize large multicomponent complexes, such as **blood clots**; the oxidation of certain amino acid side chains; and the conversion of [serine](#) and [cysteine](#) residues, usually in substituted form, to dehydroalanine, which, in turn, can readily react with [lysine](#) and [histidine](#) residues nonenzymatically, to form a new amino acid that can serve as an inter- or intramolecular crosslink.

The physiological relevance of these changes is for the most part difficult to determine. For example, it has been suggested that the deamidation of asparagine and glutamine in proteins, which varies in rate according to the surrounding amino acid sequence and the local conformation, may act as a biological clock, controlling turnover in at least some protein molecules. This hypothesis remains unsubstantiated. Some of the changes that occur, such as the formation of glycosylated derivatives of the α-amino group of long-term stable proteins, can be useful markers for monitoring pathological events, for example, the levels of the sugar-modified hemoglobin A<sub>1c</sub> are useful for monitoring diabetes. As noted above, random modifications by highly reactive species that may effectively inactivate the function of a molecule cannot be easily characterized, but they clearly contribute to some pathologies.

## 3. Enzymatic Modifications

### 3.1. Peptide Bond Hydrolysis

The cleavage of peptide bonds, formed during protein biosynthesis, is a major form of post-translational modification and can occur throughout the lifetime of a protein molecule. In functional

terms, such reactions can be grouped into three categories: (1) stabilization/translocation, (2) activation, and (3) degradation. The enzymes that catalyze all of these modifications act by removing single amino acids from either termini (exopeptidases) or by cleaving internal peptide bonds (endopeptidases). They can generally be sorted into several large families that are characterized by the reaction mechanism utilized (2).

Stabilization–translocation modifications occur primarily at the amino terminus, usually during protein biosynthesis; that is, they are co-translational. Virtually all protein synthesis is initiated with [methionine](#) or *N*-formyl-methionine in prokaryotes. The formyl group is added to the methionine as a “pre” modification. The formyl group is first removed by a deformylase in prokaryotes, and the methionine can be removed subsequently by a highly specific aminopeptidase. The specificity of this enzyme is well conserved and is defined primarily by the residue adjacent to the methionine (3). The 7 smallest amino acid side chains (as defined by the radius of gyration) are substrates, while the remaining 13 are generally not. This specificity matches well to that of the ***N*-end rule**, a mechanism that defines the rate of degradation of proteins and depends on recognition of amino-terminal residues (4). Thus, the retention of methionine by the proteins with the 13 largest amino acids protects them from premature degradation. In most organisms, it appears that proteins with the smaller amino acids at the *N*-terminus represent the largest class of proteins, as well as the largest amount of proteins on a mass basis. Thus, the majority of proteins are degraded, presumably to release methionine to allow for more protein synthesis initiation, as well as for other purposes.

As noted above, translocational reactions also involve protein hydrolysis, but usually at an internal site. Proteins to be transferred into the ER or mitochondrion have an *N*-terminal **signal sequence** that is first removed by specific enzymes appropriately located on the lumen side of the membrane to be transversed and subsequently degraded. The efficient removal of this entity prevents the protein from returning to the cytosolic compartment. The sequences defining the cleavage sites are well defined in both cases (5, 6).

The peptide bond cleavages that produce active proteins from a prohormone or proenzyme are of a more diverse nature. For the most part, these reactions occur in the ER continuum. However, other molecules can also be activated with entirely different functions, such as the structural protein [collagen](#). The processing of these precursors is well defined in many cases and involves a variety of proteinases (2, 7). Activation reactions occur primarily as internal cleavages, but they often involve more than one hydrolyzable bond. In contrast, this activation mechanism is not common for proteins that remain in the cytoplasm or localize to the nucleoplasm. When activation of proteins in these latter compartments is required, it is usually achieved by reversible modification of one or more side chains.

The degradation and turnover of proteins is considerably less well understood than either of the two processes described above. Significant amounts of protein turnover are associated with lysosomes and the hydrolytic enzymes found in this organelle, following both [endocytosis](#) and intracellular vesicle fusion. Another major entity involved in degradation is the [proteasome](#), a very complex structure in the cytoplasm and ER. Additional molecules associate with this core structure to provide it with the specificity for recognizing entities marked for degradation. This usually involves the formation of poly-[ubiquitin](#) “tags” that have been added to proteins ready for turnover. There would appear to be a variety of signals that generate these modifications, one of which is the *N*-terminus (*N*-end rule) noted above. However, the nature of the other signals is not well understood. It is noteworthy that ubiquitinylation of proteins for purposes other than signaling degradation is widespread. A number of cytosolic proteinases that may also be involved in protein turnover have been identified. Of particular note are the calpains—enzymes that are activated by increased levels of calcium—and the caspases, some of which are associated with intracellular hormone activation and appear to be widely involved in instituting [cell death](#), or [apoptosis](#).

A specific set of reactions involving the cleavage of peptide bonds, which can be viewed as a subset of the activation group, are found in proteins that undergo cleavages (either enzymatically or

nonenzymatically) leading to an eventual reorientation of internal protein sequences through [protein splicing](#) events that resemble those seen with [messenger RNA](#). Concanavalin A is a case in point, although other examples, not involving cyclization, have also been defined. The concanavalin A reaction involves multiple backbone cleavages, presumably catalyzed by a specific proteinase, although this is still an area of controversy, and the resynthesis of peptide bonds in an autocatalytic fashion. Such transpeptidation reactions are highly unfavorable in an aqueous environment and presumably occur in a “compartment” in which water has been substantially eliminated, analogous to the nonaqueous conditions used to carry out these reactions in the laboratory.

### 3.2. Amino and Carboxyl Termini

In addition to the removal of the methionine or *N*-formyl methionine *N*-terminal residue and the nonenzymatic formation of glycosylated derivatives, a number of other *N*-terminal modifications occur, mostly as a result of post-translational modifications. The most prevalent of these is acetylation catalyzed by a series of *N*<sup>a</sup>-acetyltransferases (3). The proteins with the four smallest amino acids (glycine, alanine, serine, and threonine), exposed after the initiator methionine has been removed, and those retaining methionine followed by aspartate, glutamate, or asparagine residues, represent the largest set of substrates for this modification. Other acylations can also occur, however, including by [fatty acids](#) ranging up to **myristoyl** moieties. Unsaturated fatty acids and pyruvyl groups (ketoacyl) as additional substituents have also been reported. These entities are all added post-translationally to the same residues that are usually acetylated. It is not known how this distinction is accomplished at a cellular level.

These various substitutions fulfill a variety of functional needs. *N*<sup>a</sup>-Acetylation probably prevents the loss of reducing sugars by reaction with  $\alpha$ -amino groups, which may be important for the energy metabolism of the cell. It has also been suggested that this modification stabilizes these proteins, presumably from attack by other [aminopeptidases](#), although there is little evidence to support this theory. The appearance of this modification during evolution with the development of eukaryotes is more in keeping with the former view. The longer fatty acid derivatives are usually related to association with membranes. Such modifications allow a more transient interaction than do transmembrane segments within a polypeptide chain, but still allow sufficient contact to permit other modifications, reactions, and molecular interactions to take place. Although the occurrence of these modifications is limited, they are generally of marked significance.

Modifications of the  $\alpha$ -carboxyl group include the formation of both amide and methyl esters, and are usually found on small bioactive peptides that have been formed from larger precursor molecules. *C*-Terminal amide groups result from the unusual cleavage of a *C*-terminal glycine residue between the  $\alpha$ -amino nitrogen and the  $\alpha$ -carbon, to produce the amidated derivative of the adjacent residue.  $\alpha$ -Carboxyl groups can also be activated with ATP by the formation of mixed anhydrides, and this mechanism can be used to add individual amino acids or, in the case of ubiquitin, to transfer the entire protein onto the side chain of lysine residues to form a polyubiquitin derivative. (Amino acids can also be added to the  $\alpha$ -amino group using charged tRNA.) The  $\alpha$ -carboxyl group is also the site of the formation of *N*-ethanolamineglycanphosphoinositides.

### 3.3. Side-Chain Modifications

Only five of the amino acids (glycine, alanine, valine, isoleucine, leucine) do not appear to have significant side-chain modifications in proteins. All the other 15 residues contain moieties that can be modified, in some cases extensively, to produce 150–200 different derivatives. These are briefly summarized, by side chain, in Table 1. For a more detailed treatment of these individual modifications, the reader is referred to the chapter by Krishna and Wold (see Suggestions for further Reading).

### 3.4. Oxidation

One of the most common side-chain modifications that occurs in proteins is the oxidation of cysteine to cystine to produce disulfide bonds. All proteins that are extruded into the lumen of the ER are exposed to the machinery for inserting these crosslinks, and nearly all such cysteines are so

converted. The cystine residues formed are found as both intra- and intermolecular bridges. The mechanism for forming these bonds clearly involves [thiol–disulfide exchange](#) between the cysteine [thiol groups](#) of the protein and oxidized molecules containing disulfide bonds that have been previously formed (see [Hydrophilic](#)). For the most part, these protein disulfide bonds are considered to impart greater rigidity and stability to proteins that exist and function in the extracellular milieu. Disulfides do occur in a few intracellular proteins, when the residues involved generally cycle between the reduced and oxidized states. Higher levels of oxidation of cysteine or cystine can also occur, producing progressively the sulfenic, sulfinic, and sulfonic acids, but these have not been detected in natural proteins, except under more extreme conditions or in laboratory experiments. However, they may form as products of environmental insults.

Methionine residues can also be oxidized to the corresponding sulfoxide, as an example of a common environmental “insult” that can lead to inactivation. The existence of a reductase to convert it back to methionine supports the view that this is the principal situation producing this alteration. As with cysteine, there is a higher oxidation state (sulfone), but this is seldom found except under severe oxidizing conditions. The sulfoxide modification renders the methionine residue more polar, and hence more [hydrophilic](#), which has been suggested to be important in molluscan ligament proteins, where this modification has been found.

### 3.5. Glycosylation

Carbohydrate molecules are added to protein side chains in both the ER/Golgi compartment and the cytoplasm. In the former, two types of residues serve as attachment sites: asparagine (*N*-**glycosylation**) and serine/threonine (*O*-**glycosylation**). The *N*-linked class is quite large and very complex. In brief, a core oligosaccharide, attached to the lipid carrier dolichol, is donated to the protein acceptor in the ER, probably primarily as a cotranslational event. This structure can then be further processed by both excision and addition reactions. The array of different carbohydrate structures so produced depends on the cell and the organ, reflecting the complement of enzymes present. The attachment sites occur only in —Asn–X–Ser/Thr sequences, but there are certainly also conformational requirements and restrictions that further limit the residues that are modified. Nonetheless, extracellular proteins, in particular the extracellular **domains** of transmembrane proteins, are heavily modified. *O*-Linked sugars are also added in the ER continuum, but they are rarer in occurrence. Four main classes are defined by the monosaccharide (mannose, *N*-acetylgalactosamine, galactose, and fucose) attached to the serine hydroxyl group, and they are variously distributed in different organisms.

There are many functions associated with such a varied group of modifications. The heavy glycosylation of the transmembrane proteins may be important in aiding their flow through the ER continuum and contribute to the buffy coat of eukaryotic cells. These moieties are also apparently important for regulating the surface lifetime of these proteins and may contribute, in individual cases, to their function. The glycosyl groups are also clearly important in maintaining the circulating levels of serum proteins and can be involved in targeting their turnover.

*O*-Linked *N*-acetylglucosamine is also found on certain intracellular proteins as a monomeric addition. In contrast to the extracellular modifications, this entity is readily removed, and it has been suggested that they may act in a similar fashion as phosphoryl groups in reversibly regulating diverse functions.

### 3.6. Phosphorylation

The addition of phosphoryl groups to the side chains of serine, threonine, and tyrosine, and to a much lesser extent several other amino acids (see [Table 1](#)), is certainly the most prevalent cytoplasmic modification aside from *N*-terminal alterations. The reactions are catalyzed by an extensive family of **kinases** using nucleoside triphosphates, usually ATP, as phosphate donors. For the most part, the protein kinases are specific for either serine/threonine or tyrosine residues, since the former represent esters and the latter, mixed anhydrides. A similar situation exists for the protein **phosphatases** that remove these groups. These modifications control the activity of key enzymes in

most metabolic pathways and are the primary basis for [signal transduction](#) pathways for [hormones](#) and [growth factors](#). The role of this modification for the other amino acids, such as histidine, lysine and arginine, is less well defined but is probably not as important to metabolic regulation or signal transduction.

**Table 1. Major Post-translational Modifications of Amino Acid Side Chains in Proteins**

|  |  |
|--|--|
| Arginine                               | Lysine                                 |
| <i>N<sup>W</sup></i> -ADP-ribosylation | <i>N<sup>ε</sup></i> -Acetylation      |
| <i>N<sup>W</sup></i> -Methylation      | <i>α</i> -Hydroxylation                |
| <i>N<sup>W</sup></i> -Phosphorylation  | <i>N<sup>ε</sup></i> -Methylation      |
| Asparagine                             | <i>N<sup>ε</sup></i> -Ubiquitinylation |
| <i>N</i> -Glycosylation                | <i>N<sup>ε</sup></i> -Phosphorylation  |
| Deamidation ( <i>N</i> -terminal)      | <i>N<sup>ε</sup></i> -Retinylation     |
| <i>N</i> -Methylation                  | <i>N'</i> -Lipoylation                 |
| <i>β</i> -Hydroxylation                | <i>N<sup>ε</sup></i> -Biotinylation    |
| Aspartate                              | Crosslinking to (multiple entities)    |
| <i>β</i> -Carboxylation                |  |
| <i>erythro-β</i> -Hydroxylation        | Methionine                             |
| <i>O</i> -Phosphorylation              | Oxidation                              |
| Methylation                            | Phenylalanine                          |
| Cysteine                               | Glycosylation (via oxidation)          |
| Oxidation (to cystine)                 | Proline                                |
| Thioacylation                          | Hydroxylation                          |
| Dehydration (to dehydroalanine)        | Serine                                 |
| Thioether formation                    | <i>O</i> -Phosphorylation              |
| Thioester formation                    | <i>O</i> -Acetylation                  |
| Thioglycerol formation                 | <i>O</i> -Glycosylation                |
| ADP ribosylation                       | <i>O</i> -Acylation                    |
| Glutamate                              | Selenocysteine formation               |
| <i>β</i> -Carboxylation                | Threonine                              |
| Methylation                            | <i>O</i> -Phosphorylation              |
| <i>O</i> -ADP ribosylation             | <i>O</i> -Acylation                    |
| Glutamine                              | <i>O</i> -Glycosylation                |
| <i>N</i> -Methylation                  | <i>O</i> -Methylation                  |
| Transamidation (crosslinking)          | Tryptophan                             |
| Histidine                              | Interindolic crosslinking              |
| <i>N</i> -Phosphorylation              | Glycosylation                          |
| <i>N<sup>P</sup></i> -Methylation      | Bromination                            |
| Iodination                             | Tyrosine                               |
| Dipthamide formation                   | Iodination (halogenation)              |
|  | Hydroxylation                          |

*O*-Phosphorylation  
*O*-Sulfation  
*O*-Adenylation  
(esterification)  
Interphenolic crosslinking

---

### 3.7. Acetylation

In addition to the acetylation of the  $\alpha$ -amino group, acetyl derivatives of both lysine and serine residues have been reported, although the latter is not common. As with the *N*-terminus, the acetyl group is derived from acetyl-CoA. The role of these modifications is largely unknown.

### 3.8. Acylation

In addition to acetylation, there are a large number of acyl-group additions to the side chains of serine, threonine, cysteine, and lysine residues. These are generally formed by various fatty acid or other **hydrophobic** moieties, such as isoprenoid compounds derived from the synthesis of cholesterol. As with acetylation, the carboxyl group of the donors is usually activated by CoA esterification prior to transfer. The function of these derivatives is varied but often is associated with membrane interactions, as is the case for the fatty acylation of the  $\alpha$ -amino group.

### 3.9. Methylation

Methylation is a widespread modification and occurs on the basic, acidic, and hydroxyl-containing amino acids. Mono-, di-, and tri- (when possible) derivatives of the basic residues are known. The methyl esters of aspartate and glutamate are also common, while the methyl ethers of serine and threonine have not been conclusively demonstrated and have been inferred only from the acid stability of some protein-bound methyl groups. The formation of methyl aspartate actually uses either D-aspartate or L-isoaspartate as acceptor and may be important in tagging these residues for further modification or turnover. The functional significance of the methylation of the basic residues is largely unknown (see [Hydroxylation \(Lysine, Proline\)](#)).

### 3.10. Carboxylation

Additional carboxyl groups can be added to aspartate and glutamate in a vitamin K–dependent manner. This modification adds a strong calcium binding site to these residues and is commonly found in ribosomal proteins (Asp) and in blood clotting and calcification reactions (Glu).

### 3.11. Hydroxylation

Lysine and proline residues, as well as aspartate and asparagine, can be hydroxylated, usually to produce sites suitable for further additions. With lysine and proline, these modifications are very important in establishing or stabilizing crosslinks in various structural proteins (see [Mass Spectrometry](#)). The hydroxylation of aspartate and asparagine appears to be functionally related to the  $\gamma$ -carboxylation of glutamate, as it occurs in many of the same proteins.

## 4. Analysis

The technology for determining the type, number, and sites of post-translational modifications in proteins has improved in precision and accuracy, to the point that direct identification on even small amounts of polypeptides is now possible. Advances in [mass spectrometry](#) and **chromatographic** separation techniques have not only increased the experimental capabilities but also raised the criteria for establishing these structural modifications. The success of these determinations depends, however, on the quantity and purity of the protein sample and on the instrumentation available (8).

The principal tools used to detect post-translational modifications are incorporation of



**radiolabeling**, [Mass Spectrometry](#), and [mass spectrometry](#). These approaches vary in their ability to detect, verify, or quantify alterations in protein structure, and they are often used in combination. Additional experimental tools can aid in both the detection and enrichment of modified proteins or peptides derived from proteolytic digests of complex mixtures, such as **affinity columns** (eg, antibody-conjugated, nickel-chelate) or **antibodies** specific for defined protein sequences or the modification (eg, antibodies to phosphotyrosine).

#### 4.1. Radioactivity

Radioactive precursors are frequently used to identify modifications such as acylation or phosphorylation with intact cells, subcellular fractions, or purified proteins and enzymes. After incubation with a specific precursor molecule, the cellular proteins are typically separated by one- or two-dimensional [gel electrophoresis](#), and then analyzed by both [autoradiography](#) and protein staining, so that the incorporation of the precursor molecule into a particular band or spot can be identified. Knowledge of the migration behavior of a protein in the electrophoretic system helps connect the protein with the modification. Use of a specific antibody to **immunoprecipitate** or immunolocalize the protein of interest provides further evidence.

For modifications such as phosphorylation, which are most commonly found on serine, threonine, or tyrosine residues, the protein(s) is (are) usually hydrolyzed and subjected to amino acid analysis by two-dimensional thin-layer electrophoresis and chromatography to determine which residues are modified. If radioactive methods are also used for the identification of the site(s) of modification, the labeled protein can be digested with proteinases followed by high-performance liquid chromatographic (**HPLC**) separation or another chromatographic method. Direct identification by Edman sequencing, and counting a fraction of the product of each cycle, can help identify the position of the modification. Some radioactive products are not stable or otherwise amenable to the chemistries used, and they often yield ambiguous results. They are also rarely used with more precise localization by mass spectrometry analysis because of the resulting contamination of the instrument. In some cases, the identity of the labeled peptide has been deduced by its properties and by the comigration of synthetic compounds used to “prove” the site of the modification.

#### 4.2. Edman Degradation

Edman degradation techniques have been productively used to identify many modified amino acids and their positions within a protein or peptide sequence. Edman sequencers employ controlled reaction and reagent delivery systems, coupled to automated HPLC separation of the PTH amino acids. The Edman reaction is specific for peptide bonds, and the elution positions of many natural and modified PTH–amino acids are well documented (9). While these chromatographic characteristics have been used as proof of modification in the past, a mass measurement of the modified peptide is now an important confirmation. Detection of an unusual peak in a PTH–amino acid chromatogram can lead to an unexpected observation of a protein modification, which can then be followed through by other methods. Since the sequencers now available can analyze proteins in the low picomole/high femtomole range, this remains a valuable analytical tool for determining protein modifications. It is important to note that not all modified amino acids are stable under the conditions of Edman chemistry, and some modifications may elute outside the defined chromatographic separations, or they may be missed because of poor resolution. In addition, if specific efforts are not expended to purify modified peptides, detection of a modification by Edman sequencing can easily be overlooked because of the greater abundance of the unmodified amino acid.

#### 4.3. Mass Spectrometry

Innovations in mass spectrometry, particularly in ionization techniques and in coupling to separation devices such as HPLC, have substantially advanced protein analyses. With careful technique and/or instrument modifications, many of the newer instruments are able to reach the subpicomole range. Analysis in the 25–200-fmol range is possible, but not common. Even routine analysis of small amounts of protein may, however, require repetition of experiments and experimental redesign, either because of the nature of the protein or because of contamination with other natural proteins or with experimentally derived materials. Both Edman degradation and mass spectrometry, but

particularly the latter technology, are very susceptible to problems developed during sample preparation. Nonetheless, given the constraints and the potential, it is now possible to undertake identification of proteins and protein modifications in bands from one-dimensional electrophoretic gels or single spots from [two-dimensional gel electrophoresis](#).

A careful measurement of the mass of the intact natural protein can be very informative about the extent of its processing and modification. Tables now exist to help in preliminary identification of the types of modifications possible (10). In a relatively simple situation, it is possible to obtain data consistent with the presence of a polypeptide chain with a mass equal to that predicted by its gene sequence, plus the *N*-acetylated form and/or that lacking the initiator methionine residue (if these modifications occur), and such measurements, with an analysis of the differences between predicted and expected masses, represent a useful first step in determining the type of modification. Some modifications, such as formation of *N*-terminal pyrrolidone carboxylic acid from glutamine, result in the difference of a single Dalton, and these are difficult to detect directly and unambiguously without combined analysis by Edman degradation or mass spectrometry analysis of the blocked peptide. Importantly, the measurement of a single mass of a phosphorylated polypeptide, as an example, could represent a population of polypeptides phosphorylated to the same extent, but at different sites. It is well known that different sites may be phosphorylated or dephosphorylated at different stages of signal transduction. Thus, this analysis is only the first step in determining the full range of post-translational modifications in a given protein

Proteins in solution are most amenable to direct mass measurement, although methods are being developed for analysis of intact proteins from unstained [polyacrylamide](#) gels or [blotting matrices](#). As yet, these measurements are more problematic and not as accurate. Methods have been developed to measure accurately the mass of immunoprecipitated proteins, providing an improved signal by enriching the target protein. [Membrane proteins](#) are still more of a challenge, but advances have been made in analyzing both intact proteins and hydrophobic peptides in unusual solvent mixtures or in the presence of certain nonionic [[Error: anchor with id emb1184-anc-0054 not found in linking document](#)].

#### 4.4. Applications

To determine the location of a modification within a polypeptide sequence by mass spectrometry, digests are prepared using special preparations of proteolytic enzymes. Differences between observed and predicted masses usually locate the desired modified peptides. Depending on the instrumentation available, combining offline HPLC with mass measurement, or online HPLC coupled to mass spectrometry, can improve the “coverage” of the protein sequence and the identification of all predicted and modified peptides. Sequence analysis, to verify both the sequence of the peptide and the position of the modification within the sequence, can be accomplished by mass spectrometry analysis, which can often be usefully supplemented by Edman sequencing, not only to provide an orthogonal analytical approach but also in cases where fragmentation patterns may be challenging to interpret.

Techniques have been developed to take advantage of the properties of an individual group to aid in its identification. This is particularly true for phosphate groups, which can be monitored not only for added mass, or for fragility of the bond, but also by performing precursor ion scans. Chelating resins or antibodies, such as to phosphotyrosine, are tools that can be very helpful in enriching the desired species in the peptide mixture. Glycosyl groups are also notoriously difficult to analyze. In these cases, it is often useful to couple experiments with use of purified glycosidases to aid in verifying the presence of the carbohydrate, or to determine its location. Structural characterization of the carbohydrate moiety itself is often the most difficult part of the determination.

The strategy for identifying the partners in disulfide linkages, although relatively unchanged from the first such analyses, has been substantially augmented by mass spectrometry, which has added speed, precision, and sensitivity. Proteins are digested under conditions that minimize disulfide exchange, in the absence of thiol groups or at low pH. Mass measurements are made before and after

reduction. The differences in mass and the appearance of peptides with new masses identify the peptides involved. When more than one disulfide is present in a protein, and the peptide pairs are similar in mass and cannot be determined by mass measurement alone, it is necessary to purify the disulfide peptides first. The larger the protein, the more likely this is to be the case. Difficulties also arise when half-cystine residues are adjacent or near each other in sequence, since few proteolytic enzymes will cleave at these sites. Similar protocols can be used to analyze other types of protein crosslinks.

## 5. Conclusion

Post-translational modifications provide essential extensions of the structures produced by the 20 genetically coded amino acids (including the backbone itself). These alterations govern the protein's translocation, function, and turnover. For the most part, they cannot be accurately predicted from the sequence alone and must be determined by direct measurements. These analyses are complicated by the fact that many modifications, particularly those involved in intracellular functions, represent only a small percentage of the actual population of the mature molecule, indicative of the transient events that they control. Others are unstable and may be partially or completely lost during experimental manipulations. Finally, it is also important to avoid preconceived notions in analyzing for post-translational modifications, since new derivatives are still being discovered.

## Bibliography

1. G. M. Fuller and D. S. Shields (1998) *Molecular Basis of Medical Cell Biology*, Appleton & Lange, Stamford, CT.
2. J. S. Bond and A. J. Barrett, eds. (1993) *Proteolysis and Protein Turnover*, Portland Press, London.
3. R. A. Bradshaw, W. W. Brickey, and K. W. Walker (1998) *Trends Biochem. Sci.* **23**, 263–267.
4. A. Varshavsky (1996) *Proc. Natl. Acad. Sci. USA* **93**, 12142–12149.
5. W. Neupert (1997) *Annu. Rev. Biochem.* **66**, 863–917.
6. T. A. Rapoport, B. Jungnickel, and U. Kutay (1996) *Annu. Rev. Biochem.* **65**, 271–303.
7. A. I. Smith, ed. (1995) *Peptidases and Neuropeptide Processing*, Academic Press, San Diego.
8. J. Coligan, B. Dunn, H. Ploegh, D. Speicher, and P. Wingfield (1998) *Current Protocols in Protein Science*, Wiley, New York.
9. M. W. Crankshaw and G. A. Grant (1993) *Identification of Modified PTH-Amino Acids in Protein Sequence Analysis*, ABRF, Bethesda, MD.
10. (site currently unavailable).

## Suggestions for Further Reading

11. R. G. Krishna and F. Wold (1998) "Posttranslational modifications", in *Proteins: Analysis and Design*, R. H. Angeletti, ed., Academic Press, Orlando, FL, pp. 121–206. (A detailed description of individual modifications.)
12. B. L. Martin (1996) "Post- and co-translational modifications of proteins, enzyme catalyzed, in" *and Molecular Medicine*, R. A. Meyers, ed., VCH, Weinheim, Germany, Vol. **5**, pp. 49–56.
13. R. L. Lundblad (1995) *Techniques in Protein Modification*, CRC Press, Boca Raton, FL.
14. S. M. Arfin and R. A. Bradshaw (1996) "Intracellular protein degradation, mechanisms of regulated, in" *and Molecular Medicine*, R. A. Meyers, ed., VCH, Weinheim, Germany, Vol. **3**, pp. 346–354.

## Potential Functions (Force Fields)

**Computer simulation** of macromolecules is based on a mathematical description of the dependence of their energy on their structure. Such a dependence is called the *potential surface*. Molecular potential surfaces can be evaluated in principle by using quantum mechanical approaches, but such approaches are too expensive for effective modeling of large molecules. Alternatively, one can use the fact that macromolecules are assembled by the same type of bonds that connect the atoms in small molecules. Thus one can describe large molecules as a collection of small molecular fragments, where the overall potential surface is expressed as a sum of contributions from bonded atoms and interactions between nonbonded atoms. Such a representation is usually done by sets of analytical functions that present an approximation to the true potential surface and are called *potential functions* or *force fields*. The functional forms and parameters of molecular force fields are taken from studies of small molecules with the implicit assumption that these functions are transferable from small to large molecules.

Molecule potential functions are usually given in the form

$$U(\mathbf{s}) = U_{b,\theta}(\mathbf{b}, \theta) + U_{\phi}(\phi) + U_{nb}(\mathbf{r}) \quad (1)$$

where  $\mathbf{s}$  is the vector of internal coordinates composed of  $\mathbf{b}$ ,  $\mathbf{q}$ ,  $\mathbf{f}$ , and  $\mathbf{r}$ , which are, respectively, the vectors of bond lengths, bond angles, torsional angles, and nonbonded distances. The first two terms define a very deep potential well; and because the molecule stays in most cases inside this well (except in extreme cases of bond dissociation), it is reasonable to approximate this part of the potential surface by its quadratic expansion, which is given by

$$U_{b,\theta}(\mathbf{b}, \theta) = \frac{1}{2} \sum_i K_{b,i} (b_i - b_{0,i})^2 + \frac{1}{2} \sum_i K_{\theta,i} (\theta_i - \theta_{0,i})^2 + \text{cross terms} \quad (2)$$

where  $U$  is usually given in kcal/mol,  $b$  in Å and  $q$  in radians. The torsional potential  $U(\mathbf{f})$  is a periodic function, which can be described by the leading terms in the Fourier expansion of the potential

$$U_{\phi}(\phi) = \frac{1}{2} \sum_i K_{\phi,i} (1 - \cos n\phi_i) \quad (3)$$

The nonbonded potential,  $U_{nb}$ , can be described by an atom–atom interaction potential of the form

$$U_{nb} = \sum_{ij} A_{ij} r_{ij}^{12} - B_{ij} r_{ij}^6 + C q_i q_j / r_{ij} + U_{ind}(\mathbf{r}) \quad (4)$$

where  $r_{ij}$  is the distance between the indicated atoms, the  $q_i$  are the residual atomic charges, and  $U_{ind}$  is the many-body inductive effect of the electronic polarizabilities. The constant  $C$  is 332 kcal/mol Å when the  $q$ 's are given in atomic units. The introduction of potential functions opened the way for the use of computers in conformational analysis (eg, Refs. (1-3)). The earliest use of potential functions in modeling proteins was reported by Levitt and Lifson (4).

The accuracy of the given set of potential functions depends, of course, on the specific set of parameters (the  $K_b$ ,  $K_q$ ,  $A$ ,  $B$ , etc.). These parameters can be optimized by using them to calculate different independent molecular properties (eg, energies, structures, and vibrations) and then fitting the calculated properties to the corresponding observed properties by a systematic change of the

potential parameters in a least-squares procedure. In general, the parameters  $b_0$  and  $q_0$  are sensitive to structural information and the parameters  $K_b$  and  $K_q$  are sensitive to molecular vibrations, while the  $K_f$  are determined by information about torsional barriers (which can also be obtained effectively by quantum mechanical calculations). The parameters  $A$  and  $B$  of the nonbonded potential are obtained by fitting properties of molecular crystals (3) (eg, sublimation energies, crystal structures, and lattice vibrations). Nevertheless, all properties depend in one way or another on all parameters; thus it is important to perform the fitting procedure simultaneously for different properties. Such a simultaneous fitting is the basis of the so-called consistent force field (CFF) approach (3) (for related methods see, for example, Refs. (2) and (5)).

The advent of powerful quantum mechanical approaches and the increase in availability of computer time allows one to start refining potential functions by fitting them to *ab initio* calculations (eg, Ref. 6). Such approaches should, however, involve careful attention to experimental facts. For example, charge distributions obtained from *ab initio* calculations with small basis sets are very unreliable, and, in general, studies that were based on empirical fitting have been much more reliable than those based on *ab initio* calculations until the late 1980s.

Simulation of macromolecules must reflect solvent effects, which are not just small perturbations but major contributors to the overall energetics and force. In fact, modeling of macromolecules in a vacuum is quite irrelevant to the behavior of such molecules in solution and to proteins. Here, one can use all-atom solvent models (5, 7, 8) or simplified solvent models (9, 10) as a part of the overall potential function. One can also use implicit solvent models (11, 12). The selection of proper representation of the solvent should involve special attention to the proper boundary conditions, because it is impossible to include an infinite system in the simulation. Periodic boundary conditions (7) and spherical boundary conditions (8, 9, 13) involve advantages and limitations, and improper treatment of the surface region in spherical systems might lead to incorrect polarization (8, 14).

In modeling complex processes such as **protein folding** or [protein–protein interactions](#), it is frequently advantageous to use simplified protein models (15, 16) where the potential of groups of atoms or of an entire amino acid residue is represented by a single interaction center.

## Bibliography

1. J. B. Hendrickson (1961) *J. Am. Chem. Soc.* **83**, 5437.
2. U. Burkert and N. L. Allinger (1982) *Molecular Mechanics*, American Chemical Society, Washington, D.C.
3. S. Lifson and A. Warshel (1968) *J. Chem. Phys.* **49**, 5116.
4. M. Levitt and S. Lifson (1969) *J. Mol. Biol.* **46**, 269–279.
5. P. Kollman (1993) *Chem. Rev.* **93**, 2395.
6. C. E. Dykstra (1993) *Chemical Reviews* **93**, 2339–2353.
7. M. P. Allen and D. J. Tildesley (1987) *Computer Simulation of Liquids*, Oxford University Press, Oxford, U.K.
8. G. King and A. Warshel (1989) *J. Chem. Phys.* **91**, 3647.
9. A. Warshel (1979) *J. Phys. Chem.* **83**, 1640.
10. D. Bratko, L. Blum, and A. Luzar (1985) *J. Chem. Phys.* **83**, 6367–6370.
11. C. S. Still, A. Tempczyk, R. C. Hawley, and T. Hendrickson (1990) *J. Am. Chem. Soc.* **112**, 6127–6129.
12. S. Miertus, E. Scrocco, and J. Tomasi (1981) *J. Chem. Phys.* **55**, 117.
13. C. L. Brooks III and M. Karplus (1983) *J. Chem. Phys.* **79**, 6312.
14. A. C. Belch and M. Berkowitz (1985) *Chem. Phys. Lett.* **113**, 278.
15. M. Levitt and A. Warshel (1975) *Nature* **253**, 694.

16. J. N. Onuchic, P. G. Wolynes, Z. Luthey-Schulten, and N. D. Socci (1995) Proc. Natl. Acad. Sci. USA **92**, 3626–3630.

## POU Domain

The POU “domain” is a conserved sequence motif defining a family of eukaryotic [transcription factors](#) that contains two structurally independent **DNA-binding** domains connected by a flexible linker of variable length and sequence. Both **domains** make sequence-specific contacts with DNA and together bind with high affinity to a conserved octameric binding site, ATGCAAAT. The N-terminal domain, POUS, is a **four-helix bundle** similar in structure to the bacterial **helix-turn-helix** domain, while the distal domain, POUH, is structurally homologous to the **homeobox** domain.

### Suggestion for Further Reading

W. Herr and M. A. Cleary (1991) The POU domain: versatility in transcriptional regulation by a flexible two-in-one DNA-binding domain. *Genes Dev.* **9**, 1679–1693.

## Prebiotic Evolution

The origin of life on Earth comprised a long series of steps: from the synthesis of small molecules within the primordial atmosphere or near hydrothermal vents, through the formation of biomonomers and biopolymers, culminating in the emergence of a self-replicating, autonomous organism. This philosophical outlook, if not the intimate details, began with the Russian biochemist Alexander Oparin and the British biologist J. B. S. Haldane, who, in the 1920s independently proposed a sequential model for the origin of life (1). Although the process of Darwinian selection may have modulated the populations of genetic macromolecules once the stage of an RNA (or “pre-RNA”) world developed, the term “prebiotic evolution” is used here to describe the presumed earlier era of synthesis and degradation that preceded self-replication. A common theme is that the ingredients for life were generated by the flow of energy (sunlight, lightning, or thermal radiation) through the primordial hydrosphere so that the putative mechanisms for the origin of life should be compatible with the conditions that would have prevailed in the early atmosphere and oceans.

The time frame for the emergence of life must be constrained by the physical and biological history of the Earth, but firm dates are difficult to establish. The age of our planet is approximately 4.5 billion years, and life cannot be more ancient unless it came from an extraterrestrial environment. The latter possibility should not be dismissed and may gain credence if the existence of past (or current) life on other planets is confirmed (2). Regardless of the source, organisms would not have survived until the Earth cooled and meteoritic bombardment subsided sufficiently such that the oceans remained in a liquid state. The conditions necessary for sustainable life may not have persisted until about four billion years ago. However, evidence of extant organisms appears in fossilized stromatolites in Western Australia from 3.5 billion years ago, and possibly in apatite inclusions from rocks in Greenland dated at almost 3.9 billion years, suggesting that the appearance of life occurred quite rapidly on a geologic time scale once the conditions were favorable (3).

A historic demonstration of the feasibility of prebiotic simulations was performed by Stanley Miller (4) during the fall of 1952 in the laboratory of Harold Urey at the University of Chicago. Based on Urey's cold-accretion theory for the origin of the planets, Miller subjected a gaseous mixture of methane, ammonia, hydrogen, and water to an electrical discharge, analogous to the effect of lightning in the atmosphere of the young Earth. Chromatographic analysis revealed the presence of three biological amino acids (glycine, alanine, and aspartic acid) along with other products. Further work by Miller and by many other research groups extended the suite of presumed prebiotic amino acids to include glutamic acid, leucine, isoleucine, serine, and threonine. All of the chiral amino acids were obtained as a racemic mixture of left- and right-handed forms, as expected from the achiral starting materials.

The apparent success of these early experiments, which depended on the accessibility and sensitivity of assays specific for amino acids, heightened interest in the nascent discipline of origins-of-life studies. However, prescient objections to the reducing atmosphere of the Miller-Urey simulations were raised first by Philip Abelson, a geochemist at the Carnegie Institution, who argued that the hydrogen-rich gases would have been rapidly replaced by a secondary atmosphere in which carbon was present as carbon dioxide or carbon monoxide, while nitrogen was probably present as molecular dinitrogen (5). By the early 1980s, a growing body of computational and experimental evidence slowly led to a revised view of the dominant atmosphere during the era before life began. Unfortunately, as shown by Miller and others, these non-reducing gas mixtures ( $\text{CO}_2/\text{N}_2/\text{H}_2\text{O}$  or  $\text{CO}/\text{N}_2/\text{H}_2\text{O}$ ) give dramatically lower yields and less variety in amino acids produced by electric discharge (6). If gas-phase syntheses of the Miller-Urey design were important in the origin of life, they must either have proceeded during a brief period when the Earth was very rich in hydrogen, or there may have been another unidentified source of reducing equivalents that maintained a source of hydrogen over a longer stretch of geologic time. Alternatively, the formation of amino acids and other organic molecules may have been favored near submarine hydrothermal vents, where reducing equivalents would have been present in the extruded gases; the technical challenges inherent in high pressures and temperatures have necessarily restricted the number of such simulations, but compounds as complex as pyruvic acid have been detected using formic acid as the carbon source (7). Nevertheless, the existence of both thermodynamic and kinetic barriers to the reduction of  $\text{CO}_2$  (the preferred starting point for any prebiotic synthesis) raises many questions about the availability of the organic precursors in hydrothermal models (8).

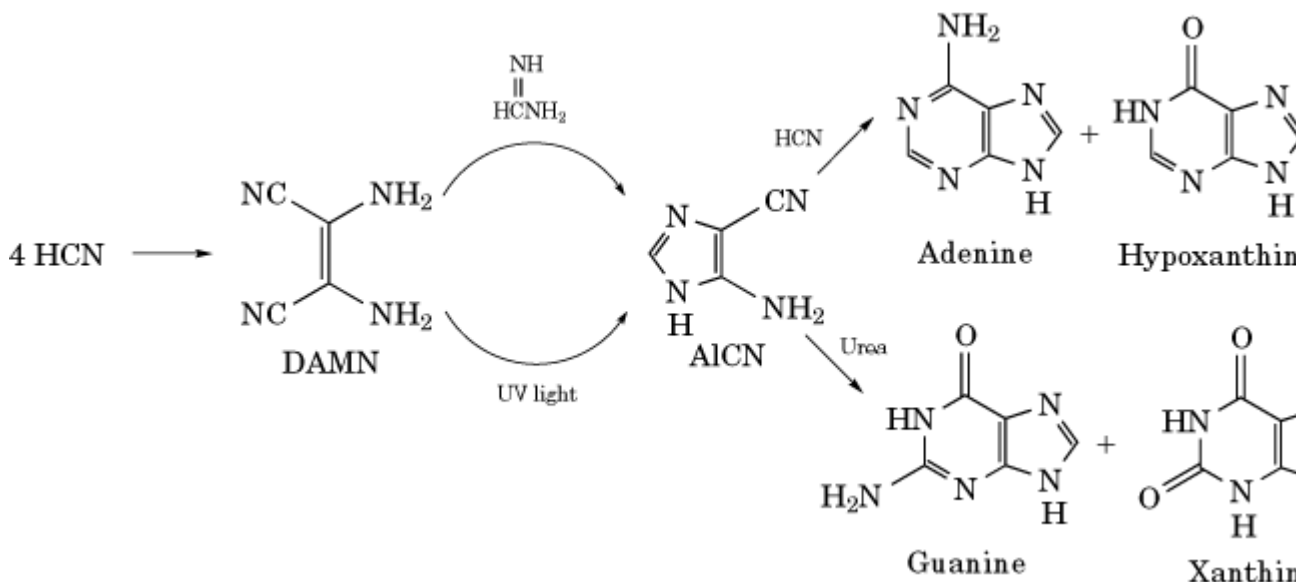
A very different source of the vital ingredients may exist beyond the Earth, in the form of comets, interplanetary dust, and asteroidal debris (9). These materials are believed to contain carbon compounds, some of which can survive passage through the atmosphere when such objects approach our planet. An intriguing possibility is that their cargo could already be enriched in the left-handed form of amino acids that are required to make modern proteins (10). While much remains to be understood about the chemical processing of these extraterrestrial bodies, there is great interest in their possible contribution to the primordial soup.

## 1. The Cyanide Paradigm

John Oró at the University of Houston in 1960 made a startling observation: adenine, a constituent of RNA as well as the nicotinamide and flavin cofactors, was formed after acid hydrolysis of ammonium cyanide solutions, which were deemed prebiotic starting materials based on the prior detection of HCN in the Miller-Urey electric discharge reactions (11). Although the conditions involved high concentrations of both cyanide and ammonia, James Ferris and his coworkers (12) at Rensselaer Polytechnic Institute obtained a yield of 0.03% after acid hydrolysis of a six-month reaction at room temperature of 0.1 M HCN, adjusted to pH 9.2 with ammonia and kept in the dark; adenine was also detected after hydrolysis at pH 8.5 (12). Ferris, working with Leslie Orgel (13) at the Salk Institute, demonstrated that adenine could be synthesized photochemically using near-UV light, via an intermediate known as 4-amino-5-cyanoimidazole (AICN, Fig. 1). This imidazole

derivative can be prepared through a dark reaction by combining diaminomaleonitrile (DAMN) with formamidine (14). Urea, widely regarded as a prebiotic compound formed (among other ways) by the hydrolysis of HCN oligomers, reacts with AICN to give another important purine, guanine (15).

**Figure 1.** Synthesis of purines from HCN, showing two possible routes to the intermediate 4-amino-5-cyanoimidazole (AICN) from the HCN tetramer, DAMN.

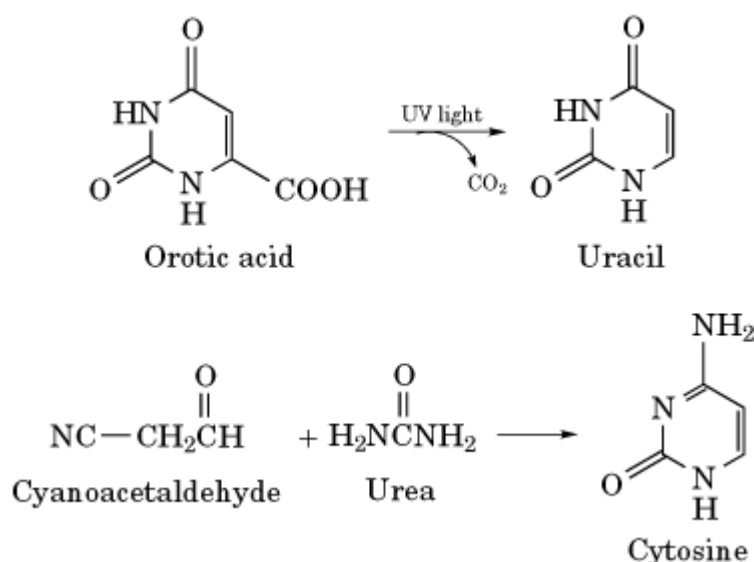


The adenine that is released upon alkaline hydrolysis of the HCN oligomers is probably formed by a separate pathway (not shown in Fig. 1) that does not require light. The initial products in the oligomerization are now well-characterized, proceeding through a stepwise self-addition to give the tetramer, DAMN, which is a precursor to the imidazole derivatives (Fig. 1). The higher oligomers constitute a heterogeneous mixture that is poorly defined, but their composition probably does not correspond to any known biopolymers (12).

HCN oligomerization also provides a source of the RNA component, uracil, which was first identified in such a mixture by Alan Schwartz and Andries Voet (16) at the University of Nijmegen. Ferris and Joshi (17) showed the presence of orotic acid after alkaline hydrolysis of HCN oligomers and further demonstrated that this pyrimidine undergoes a photochemical decarboxylation (Fig. 2) to yield uracil. Cytosine is more difficult to synthesize, but Miller and Robertson (18) have published an efficient pathway based on urea and cyanoacetaldehyde as starting materials; cyanoacetaldehyde is formed by hydrolysis of cyanoacetylene (a product of methane and dinitrogen atmospheres subjected to an electrical discharge, which thus requires a source of reduced carbon). Hydrolysis of cytosine yields uracil so that these syntheses provide an indirect route to this pyrimidine.

**Figure 2.** Possible steps in the prebiotic synthesis of pyrimidines.





Finally, HCN also provides a route to amino acids, notably glycine, alanine (alpha and beta isomers) and aspartic acid (12). Thus, an especially attractive feature of cyanide as a precursor to biomolecules is that it yields constituents of proteins and nucleic acids by a common reaction system. Moreover, Leslie Orgel and his colleagues (14) proposed an innovative mechanism for concentrating HCN by eutectic freezing: at  $-23.4^{\circ}\text{C}$ , a solution of 74.5 mole percent HCN (25 M) is obtained. Despite the low temperature, oligomerization is quite rapid under such conditions, which might have occurred in glacial climates.

An unresolved dilemma is how HCN itself might have been made on the early Earth (19). As Miller and others have shown, cyanide is readily formed when reducing gas mixtures of methane and ammonia are subjected to an electrical discharge, and this species was proposed as an important intermediate in the Miller-Urey syntheses of amino acids. By the same token, HCN is not formed from more oxidized atmospheres of  $\text{CO}_2/\text{N}_2/\text{H}_2\text{O}$  or  $\text{CO}/\text{N}_2/\text{H}_2\text{O}$ . One possibility proposed by Fujio Egami (20) at the Mitsubishi-Kasei Institute of Life Sciences is that HCN resulted from the reaction of formaldehyde with hydroxylamine. Formaldehyde is a plausible photoproduct of oxidized atmospheres ( $\text{CO}_2$  and  $\text{H}_2\text{O}$ ), but the source of  $\text{NH}_2\text{OH}$  remains to be elucidated (21). If the synthesis of HCN can ultimately be reconciled with a non-reducing environment on the primordial Earth, this achievement would greatly strengthen the paradigm that cyanide played an important role in the synthesis of amino acids, purines, and pyrimidines.

## 2. Activation Processes

Modern biochemistry exploits the unique properties of phosphate as an activating group to drive unfavorable reactions, and the earliest biosynthetic processes may have used phosphate in an analogous fashion. A striking illustration of the activating power of this group is found in glycolaldehyde phosphate, which Albert Eschenmoser (22) of the Federal Technical Institute (ETH) in Zurich has proposed as a precursor of the RNA sugar, ribose. Under strongly alkaline conditions, glycolaldehyde phosphate reacts with formaldehyde to give ribose 2,4-diphosphate as the major product, with glyceraldehyde 2-phosphate serving the role of intermediate; by contrast, formaldehyde alone gives a complex mixture with only a trace of ribose. More recently, Gustaf Arrhenius and his colleagues (23) at the Scripps Institution of Oceanography have shown that certain mineral catalysts (known as layered double hydroxides) can effect the formation of ribose 2,4-diphosphate from glyceraldehyde 2-phosphate and glycolaldehyde phosphate at near-neutral pH. While the subsequent steps from this derivative to nucleotides have yet to be elucidated, the work of Eschenmoser represents a significant advance in stereoselective synthesis.

A formidable challenge in prebiotic research has been the demonstration of possible condensation mechanisms by which the elements of water ( $H_2O$ ) could be removed in order to couple two molecules. One approach has involved the heating of the ingredients in the dry state to drive off the water. This method yields nucleosides in 2% to 10% yield when purines (adenine or guanine) are heated in the presence of ribose, but the reaction fails when the pyrimidines cytosine or uracil are substituted (24). Repeated wet-dry cycles, which simulate the periodic flooding and evaporation of a tidal lagoon, effect the self-condensation of amino acids to form modest amounts of oligopeptides up to the pentamer in a process that is catalyzed by a common clay mineral called kaolinite (25). While such a reaction system is relatively inefficient, with total yields near 1% after 2 months, it represents a more natural model than studies carried out with transient condensing agents as the activating species (26).

Phosphorylation of nucleosides and other substrates represents another step that has met with limited success. Organic condensing agents such as cyanamide and cyanoacetylene are effective but not under the low micromolar concentrations of phosphate that likely prevailed in the primordial ocean (26). Inorganic polyphosphates including the trimer, which has been detected in volcanic fumaroles, can convert adenosine into a mixture of mono-, di-, and triphosphates, but the high temperatures associated with volcanic environments render them unsuitable for complex organic synthesis (27). A structurally related compound, cyclotriphosphate or “trimetaphosphate,” has been employed in the preparation of glycolaldehyde phosphate and other phosphate esters (28); however, trimetaphosphate concentrations would be limited by the geothermal formation of longer polymers from which it is derived. Mixtures of urea and inorganic phosphate give high yields of phosphorylated products under mild heating, but this approach requires high concentrations of urea that may not have been available on the early Earth (29). Further work must be done to elucidate plausible prebiotic pathways to phosphate esters.

In summary, a significant body of literature (see Suggestions for Further Reading) has emerged that demonstrates the feasibility of prebiotic syntheses of specific compounds under particular conditions. The efficacy of cyanide as a precursor to purines, amino acids, and, to a lesser extent, pyrimidines suggests that HCN likely had a role in the formation of the first biomolecules. However, much work remains to reconcile such pathways with geochemical conditions that might have prevailed on the early Earth and to elucidate how a genetic macromolecule might have formed from a dilute primordial soup. Despite the many problems in nucleotide assembly, the postulate of a so-called RNA world has provided remarkable insights into the interrelated roles of replication and catalysis in the origins of life.

## Bibliography

1. H. Kamminga (1988) *Biosphere* **18**, 1–11.
2. D. S. McKay, E. K. Gibson, Jr., K. L. Thomas-Keprta, H. Vali, C. S. Romanek, S. J. Clemett, X. D. F. Chillier, C. R. Maechling, and R. N. Zare (1996) *Science* **273**, 924–930.
3. J. M. Hayes (1996) *Nature* **384**, 21–22.
4. S. L. Miller (1953) *Science* **117**, 528–529.
5. P. H. Abelson (1966) *Proc. Natl. Acad. Sci. U.S.A.* **55**, 1365–1372.
6. G. Schlesinger and S. L. Miller (1983) *J. Mol. Evol.* **19**, 376–382.
7. G. D. Cody, N. Z. Boctor, T. R. Filley, R. M. Hazen, J. H. Scott, A. Sharma, and H. S. Hatten, Jr. (2000) *Science* **289**, 1337–1340.
8. M. A. A. Schoonen, Y. Xu, and J. Bebie (1999) *Biosphere* **29**, 5–32.
9. C. Chyba and C. Sagan (1992) *Nature* **355**, 125–132.
10. J. R. Cronin and S. Pizzarello (1997) *Science* **275**, 951–955.
11. J. Oró (1960) *Biochem. Biophys. Res. Commun.* **2**, 407–412.
12. J. P. Ferris, P. C. Joshi, E. H. Edelson, and J. G. Lawless (1978) *J. Mol. Evol.* **11**, 293–311.

13. J. P. Ferris and L. E. Orgel (1966) *J. Amer. Chem. Soc.* **88**, 1074.
14. R. A. Sanchez, J. P. Ferris, and L. E. Orgel (1967) *J. Mol. Biol.* **30**, 223–253.
15. R. A. Sanchez, J. P. Ferris and L. E. Orgel (1968) *J. Mol. Biol.* **38**, 121–128.
16. A. B. Voet and A. W. Schwartz (1982) *Origins Life* **12**, 45–49.
17. J. P. Ferris and P. C. Joshi (1979) *J. Org. Chem.* **44**, 2133–2137.
18. M. P. Robertson and S. L. Miller (1995) *Nature* **375**, 772–774.
19. J. F. Kasting (1993) *Science* **259**, 920–926.
20. Kamaluddin, H. Yanagawa, and F. Egami (1979) *J. Biochem.* **85**, 1503–1507.
21. J. P. Pinto, C. R. Gladstone, and Y. L. Yung (1980) *Science* **210**, 183–185.
22. A. Eschenmoser and E. Loewenthal (1992) *Chem. Soc. Rev.* **21**, 1–16.
23. R. Krishnamurthy, S. Pitsch, and G. Arrhenius (1999) *Biosphere* **29**, 139–152.
24. W. D. Fuller, R. A. Sanchez, and L. E. Orgel (1972) *J. Mol. Evol.* **1**, 249–257.
25. N. Lahav, D. White, and S. Chang (1978) *Science* **201**, 67–69.
26. J. Hulshof and C. Ponnampereuma (1976) *Origins Life Evol. Biosphere* **7**, 197–224.
27. Y. Yamagata, H. Watanabe, M. Saitoh, and T. Namba (1991) *Nature* **352**, 516–519.
28. R. Krishnamurthy, G. Arrhenius, and A. Eschenmoser (1999) *Origins Life Evol. Biosphere* **29**, 333–354.
29. R. Reimann and G. Zubay (1999) *Origins Life Evol. Biosphere* **29**, 229–247.

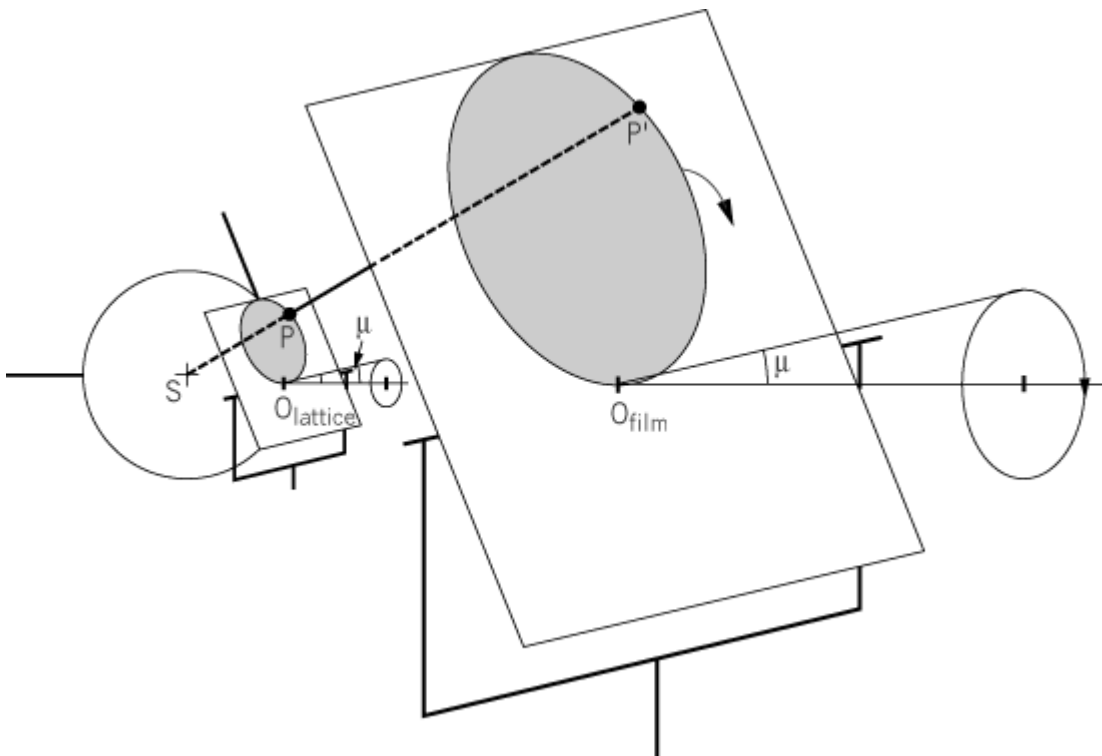
### Suggestions for Further Reading

30. A. Brack, ed. (1998) *The Molecular Origins of Life: Assembling Pieces of the Puzzle*, Cambridge University Press, New York.
31. J. P. Ferris and W. J. Hagan Jr. (1984) HCN and chemical evolution: the possible role of cyano compounds in prebiotic synthesis. *Tetrahedron* **40**, 1093–1120.
32. S. F. Mason (1991) *Chemical Evolution: Origin of the Elements, Molecules, and Living Systems*, Oxford University Press, New York.

### Precession Photograph

Precession cameras were widely used in protein [X-ray crystallography](#) in earlier days. Their advantage is that they give an undistorted image of the [reciprocal space](#) of lattice planes (Fig. 1). Moreover, the symmetry in the crystal and its [space group](#) can easily be determined and the crystal quality checked. The precession camera is not suitable for efficient three-dimensional data collection because it records only one lattice plane in each exposure. Therefore, it is no longer a popular instrument in protein X-ray crystallography.

Figure 1.



### Suggestion for Further Reading

J. Drenth (1995) *Principles of Protein X-ray Crystallography* Springer, New York, pp. 38–46.

## Precipitation

The solubility of a macromolecule is a complex function of a number of factors, such as pH, temperature, the chemical and physical nature of the macromolecular surface, and the nature of the precipitant (or solubilizer) used. The purposes of precipitating [proteins](#) out of solution may be varied.

1. During the purification, preparation of a particular protein out of a tissue extract or cell homogenate, the approach can be that of partial precipitation, or fractionation by precipitation using mild nondenaturing precipitants. This uses the solubility properties of the various components of the mixture to remove those that are less soluble than the desired product, or to separate the product in precipitated form from more soluble impurities. Once the more soluble and less soluble components of the original mixture have been removed, the desired protein can be purified further by cutting sharper fractions, based on solubility in a given agent, say,  $\text{Na}_2\text{SO}_4$ , until only a single component is left that precipitates very sharply at a constant very narrow concentration range of the precipitant. This represents pure protein by the criterion on solubility (1).
2. Proteins may be precipitated out of solution in the form of crystals for the purpose of [X-ray crystallography](#) structural studies. This involves very special techniques (see [Crystallization](#)).
3. At times, when it is necessary to obtain a protein-free solution for the purpose of analysis for other

components, in which proteins may interfere, proteins are precipitated in a denatured, generally coagulated state. This is done by the addition of general strong precipitants, with *trichloroacetic acid* being an additive of choice.

Precipitation of proteins during purification or crystallization is usually accomplished by the gradual addition at high concentration (0.3 to 10M) of small molecules that interact weakly, and gently, with the proteins. All these are known to be preferentially excluded from the protein [accessible surface](#) at the environmental conditions (such as temperature) used (2) (see [Preferential Hydration](#)). Many of them are structure stabilizers (see [Stabilization And Destabilization By Co-Solvents](#)). The most common precipitants are [sulfate salts](#) (ammonium sulfate, sodium sulfate, magnesium sulfate); ammonium sulfate is the most commonly used precipitant because of its high solubility in water and strong salting out properties (see [Salting In, Salting Out](#)).

Some organic solvents can also be used as precipitants (3). These include ethanol, acetone, and butanol. They must be used, however, at low concentration (about 5 to 10%) since they tend to induce protein **denaturation**. Two organic compounds, *2-methyl-2,4-pentanediol* (MPD) and *polyethylene glycol* (PEG), have been found to be excellent protein precipitants and good crystallizing agents. They may be used at a high concentration (eg, 50% MPD), but their use is limited to room temperature at which they are preferentially excluded from proteins. At higher temperatures, however, both become good protein denaturing agents.

#### Bibliography

1. J. H. Northrop, M. Kunitz, and R. M. Herriott (1948) *Crystalline Enzymes*, Columbia Univ. Press, New York.
2. T. Arakawa and S. N. Timasheff (1985) *Meth. Enzymol.* **114**, 49–77.
3. A. McPherson (1985) *Meth. Enzymol.* **114**, 112–120.

#### Suggestion for Further Reading

4. E. J. Cohn and J. T. Edsall (1943) *Proteins, Amino Acids and Peptides*, Reinhold, New York.

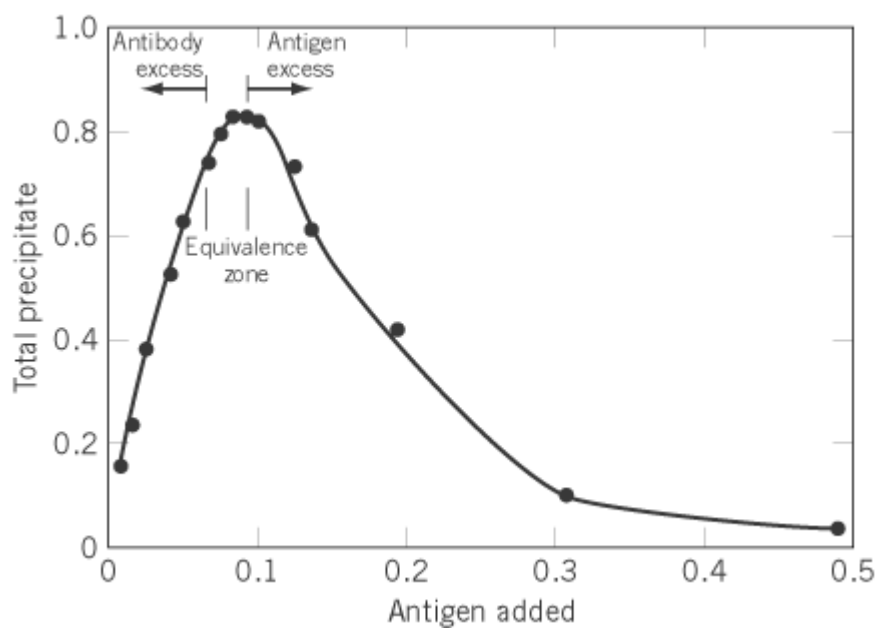
### Precipitin Reaction

The precipitin reaction is a phase-change phenomenon observed in **antibody-antigen reactions**. Kraus (1), its discoverer, demonstrated formation of a precipitate when spent bacterial culture broth was mixed with a specific antiserum. Marrack and Smith (2) showed that recovery of precipitate was largely independent of solution conditions or the presence of nonspecific proteins, so the precipitin reaction was not a simple phase change. Marrack (3) proposed correctly that the precipitate was an extended lattice in which multivalent [antibody](#) molecules form links with two or more [antigen](#) molecules, which, in turn, can form multiple links with other antibody molecules.

Since lattice formation requires multiple intermolecular links per monomer, the precipitin reaction is very sensitive to the ratio of antibody and antigen. A precipitin curve, shown in Figure 1, refers to the rise, peak, and fall in amount of precipitate recovered as a fixed amount of serum is titrated with antigen. The maximum amount of precipitate forms when antigen and antibody are present in similar molar amounts (the equivalence zone). At a low antigen:antibody ratio, few antibody molecules will bind more than one antigen molecule, hence the multiple intermolecular links necessary for a lattice again do not form (antibody excess zone). At a high antigen:antibody ratio, all antibody combining

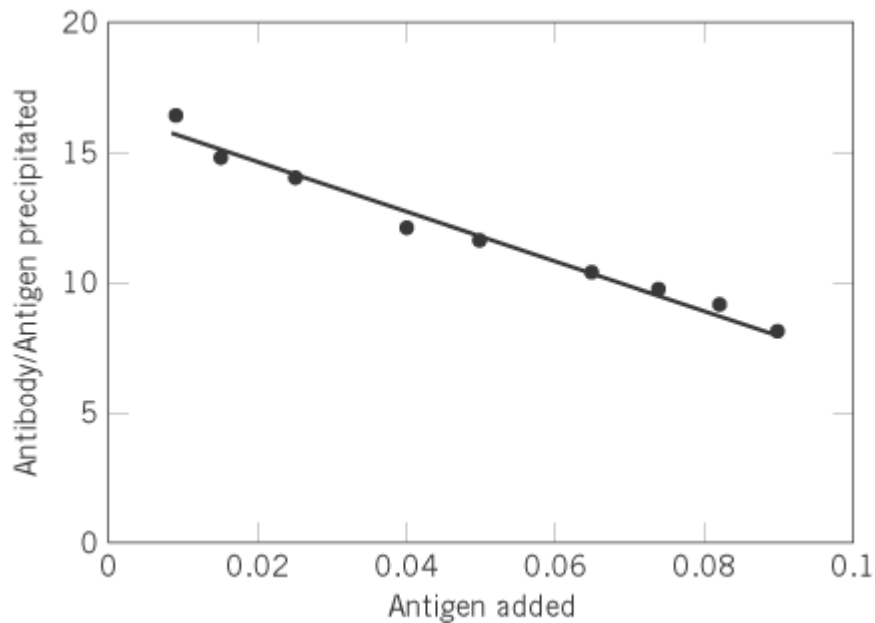
sites will be saturated with antigen, and few antigen molecules will be shared between two antibodies (antigen excess zone). The zones are defined empirically as those sets of reactions whose supernatants give additional precipitate upon adding a small increment of antibody (in the antigen excess zone) or antigen (in the antibody excess zone). In the equivalence zone, little antibody or antigen remains in the supernatant, and no additional precipitate can be induced. With some antigens and antisera, complete antibody precipitation occurs over an extended range of antigen additions, while excess antigen is still observed in the supernatant. In such cases, the term “inhibition zone” refers to the region of the antigen excess zone in which the absolute recovery of precipitate declines.

**Figure 1.** Precipitin curve. An antibody is titrated with antigen, and the resulting precipitates are isolated and quantitated. (Data are from Ref. 5.)



Heidelberger and Kendall (4, 5) analyzed the composition of precipitates and discerned a general relation that formed the basis of a quantitative assay of antigen or antibody preparations of unknown concentration. In this method, a precipitin curve is obtained, and the experimental points in the equivalence and antibody excess zones are identified as described above. (Points in the antigen excess zone are not used.) These authors found for many antibody–antigen systems that a straight line results if the ratio of antibody to antigen in the precipitate is plotted against the amount of antigen added, illustrated in Figure 2. (A frequent simplifying assumption is that all antigen added is precipitated.) The linearity of the relation allowed easy interpolation; hence the antigen content of an unknown could be determined by comparing the antibody:antigen ratio in the precipitate from the unknown to the results of precipitin reactions from standards of known concentration.

**Figure 2.** Precipitin curve analyzed by the method of Heidelberger and Kendall (4). Data from the equivalence and antibody excess zones in Figure 1 are replotted as the ratio of antibody to antigen found in each precipitate versus the quantity of antigen added to initiate the precipitation.



Although the precipitin curve technique illustrated in the figures is no longer used for analysis, the precipitin reaction itself has been incorporated into many immunological techniques (see [Immuno-electrophoresis](#)). One of the most visually elegant applications is the Ouchterlony double-diffusion assay (see [Immunoassays](#)), in which one can infer properties such as molecular weight and cross-reactivity of an antigen or antiserum preparation from the shape of a precipitin line on an agar plate (6). A turbidimetric assay based on the precipitin reaction is widely used in clinical laboratories. In this type of immunoassay, antibody that is free in solution reacts with microscopic antigen-coated latex beads (or free antigen reacts with antibody-coated beads) to form large cross-linked aggregates that are easily and sensitively detected by [light scattering](#) (7).

#### Bibliography

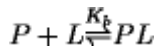
1. R. Kraus (1897) *Wien. Klin. Wochenschr.* **10**, 736–738.
2. J. Marrack and F. C. Smith (1931) *Br. J. Exp. Pathol.* **12**, 30–35.
3. J. R. Marrack (1934) *The Chemistry of Antigens and Antibodies. MRC Special Report*, 1934, London: HMSO.
4. M. Heidelberger and F. E. Kendall (1935) *J. Exp. Med.* **61**, 563–591.
5. M. Heidelberger and F. E. Kendall (1935) *J. Exp. Med.* **62**, 697–720.
6. Ö. Ouchterlony (1948) *Acta Pathol. Microbiol. Scand.* **25**, 186–191.
7. J. M. Singer and C. M. Plotz (1956) *Am. J. Med.* **21**, 888–892.

#### Suggestions for Further Reading

8. E. D. Day (1966) "Precipitation reactions", in *Foundations of Immunochemistry*, Baltimore, Williams & Wilkins, pp. 111–148.
9. M. Heidelberger (1939) Quantitative absolute methods in the study of antigen–antibody reactions. *Bacteriol. Rev.* **3**, 49–95.
10. E. A. Kabat (1961) "Precipitin reaction", in *Experimental Immunochemistry*, Thomas, Springfield, Ill., pp. 22–96.

## Preferential Binding

The equilibrium [binding](#) of ligands to proteins (nucleic acids) is classically measured by [equilibrium dialysis](#), or related techniques, such as gel **gel filtration**. The binding can be expressed by a simple mass action equation



where  $P$  is protein and  $L$  the ligand, with a binding equilibrium constant

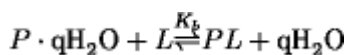
$$K_b = \frac{[PL]}{[P][L]}$$

The binding isotherm is expressed by

$$\bar{\nu} = \frac{nK_b[L]}{1 + K_b[L]}$$

where  $n$  is the extent of binding (moles ligand per mole protein) at a concentration of free ligand  $[L]$  and  $n$  the number of sites.

In the case of extremely weakly interacting ligands, which is the case with agents that stabilize (eg, sugar, glycerol) or destabilize ([urea](#), **guanidinium chloride**) the folded structures of protein or nucleic acids, they are required at high concentrations ( $\cong 1 M$ ). Consequently, they are more properly considered to be cosolvents, and displacement of water from the binding sites must be taken into account explicitly ([1](#)). The binding process is depicted in Figure [1](#), and the equilibrium equation that expresses the water–ligand exchanges at a site is



The binding equilibrium constant becomes an exchange constant

$$K_b = K_{ex} = \frac{[PL][H_2O]}{[P \cdot H_2O][L]}$$

For the simple case where one ligand molecule displaces one water molecule, the binding isotherm becomes

$$\bar{\nu} = \frac{n(K_{ex} - 1/m_w)m_L}{1 + K_{ex}m_L}$$

where  $m_1$  and  $m_3$  are the molal (moles per 1000 g  $H_2O$ ) concentrations of water and ligand (cosolvent), respectively. Since  $M_1 = 55.56$  moles  $H_2O$  per 1000 g  $H_2O$ ,  $(1/m_1) = 0.018m^{-1}$ . This means that the actual experimental result of a binding measurement by, say, equilibrium dialysis can yield *negative*, as well as positive, values of the extent of binding,  $n$ . When  $K_{ex} > 0.018m^{-1}$ , the measured binding is positive. When  $K_{ex} < 0.018m^{-1}$ , the value of  $n$  is negative. Operationally, the magnitude of preferential binding (in fact, the extent of any binding, weak or strong) is obtained from the measurements of the ligand concentration inside and outside the dialysis bag (in terms of



equilibrium dialysis):

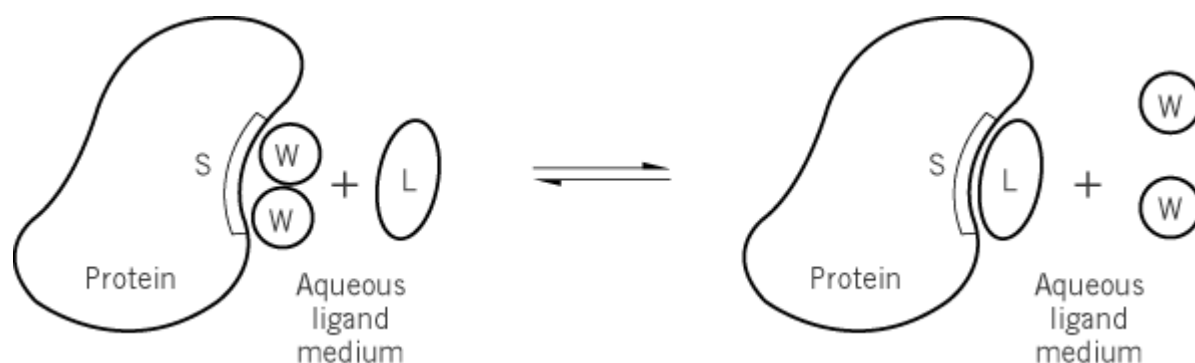
$$\bar{v} \left( \frac{\text{Moles ligand}}{\text{Mole protein}} \right) = \frac{(\text{Conc. ligand inside bag} - \text{Conc. ligand outside bag})}{\text{Conc. protein}}$$

It is clear, therefore, that the magnitude and *sign* of  $n$  are simply the expression of the difference between two measured concentrations. In more rigorous notation, this can be stated as

$$\bar{v} = \left( \frac{\partial m_L}{\partial m_{pr}} \right)_{T, \mu_w, \mu_L}$$

ie, variation of the ligand concentration induced by the presence of protein at a given temperature  $T$  and at dialysis equilibrium (this is indicated in the notation by the chemical potentials of water  $\mu_w$  and ligand  $\mu_L$  being identical inside and outside the bag). In this way, it has been found that many substances which stabilize protein structure give negative values of binding. For example, equilibrium dialysis of **ribonuclease A** (RNase A) in 1  $M$  sucrose at pH 7 yields the following result:  $n = -7.6$  moles sugar bound per mole RNase A. In other words, there is a deficiency of sugar molecules on the protein surface relative to the sugar concentration in the bulk solvent. This is referred to as **preferential exclusion**. Since an insufficiency of sugar exists, there must be an excess of water. This excess of water is referred to as [preferential hydration](#).

**Figure 1.** Schematic representation of the replacement of water W molecules by a ligand L molecule at a binding site S on a protein. [Reprinted with permission from S. N. Timasheff (1992) *Biochemistry* **31**, 9857–9864. Copyright 1997 American Chemical Society.]



Total binding is the number of molecules of a solvent component (water or cosolvent) that form contacts with the protein (or nucleic acid) surface. Binding measured by the usual techniques is *preferential binding*. It is related to total binding by

$$\bar{v} = \left( \frac{\partial m_L}{\partial m_{pr}} \right)_{T, \mu_w, \mu_L} = B_L - \frac{m_L}{m_w} B_w \quad (1)$$

where  $B_L$  is the total binding of the ligand to the protein and  $B_w$  is that of water. In other words,  $B_L$  and  $B_w$  are the numbers of molecules of ligand and water that at any moment occupy sites on the surface of the protein molecule (are in contact with the protein surface). It is clear, therefore, that equilibrium dialysis result do not necessarily give the total number of ligand molecules that occupy sites on (are bound to) a protein molecule. Whether this is true or not depends on the magnitude of the second term on the right-hand side of Eq. (1). If binding of the ligand is strong (as is the case

with [enzyme](#) substrates, cofactors, various effectors, etc.), it is measured at a low concentration of the ligand, typically  $<10^{-5}M$ . This renders the second term of Eq. (1) negligibly small, with the consequence that equilibrium dialysis does yield the actual number of ligand molecules bound to the protein. This permits use of the classical binding equations, as expressed by various standard plots, such as the **Scatchard**, double reciprocal or **Hill plots** [(2)]. In the case of weakly interacting cosolvents that must be used at high concentration ( $>1 M$ , eg, 8  $M$  urea), the second term becomes significant and  $n \neq B_L$ , so that equilibrium dialysis does not give total binding, ie, the true number of ligand molecules in contact with (occupying sites on) the protein.

#### Bibliography

1. J. A. Schellman (1990) *Biophys. Chem.* **37**, 121–140.
2. J. Wyman and S. J. Gill (1990) *Binding and Linkage*, University Science Books, Mill Valley, CA, Chap. "2–3".

#### Suggestions for Further Reading

3. S. N. Timasheff (1993) *Ann. Rev. Biophys. Biomol. Struct.* **22**, 67–97.
4. S. N. Timasheff (1995) In *Protein-Solvent Interactions* (R. B. Gregory, ed.), Marcel Dekker, New York, Chap. "11".

### Preferential Hydration

Preferential hydration is the excess of water in the immediate domain of a protein relative to the water concentration in the bulk solvent (eg, in 3  $M$  glycerol). It is related to the [preferential binding](#) of a ligand by a simple reciprocal relation (see [Binding](#)):

$$\left(\frac{\partial m_w}{\partial m_{pr}}\right)_{T, \mu_w, \mu_L} = -\frac{m_w}{m_L} \left(\frac{\partial m_L}{\partial m_{pr}}\right)_{T, \mu_w, \mu_L} \quad (1)$$

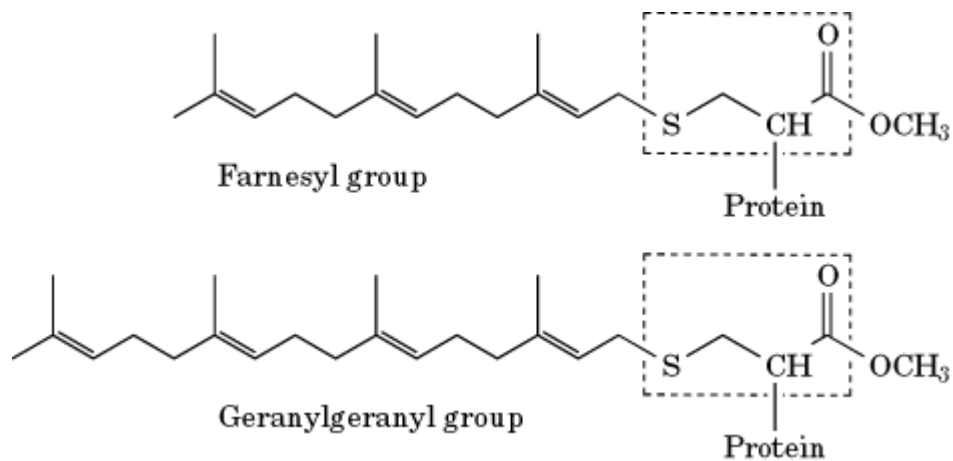
The term on the left-hand side is the excess of water molecules on the surface of the protein relative to the concentration of water in the bulk mixed solvent. As we see, this expression tells us that there is an excess of water on the protein surface if there is a deficiency of ligand, and vice versa. Understanding this and similar relations (see [Binding](#)) is of major importance in the stabilization and destabilization of protein structure by the addition of cosolvents, such as [urea](#) or sucrose. Preferential hydration must not be confused with total [hydration](#).

### Prenylation

Prenylation or isoprenylation is a [post-translational modification](#) process in which [cysteine](#) residues close to the C-terminal regions of some eukaryotic [proteins](#) are biosynthetically modified with an isoprenoid lipid: the 15-carbon farnesyl group or the 20-carbon geranylgeranyl group (see Fig. 1 and Table 1). Prenylation provides some proteins with a hydrophobic [membrane anchor](#), and is important

for their correct localization within the cell. Prenylation is one of several processes that attach lipid membrane anchors to proteins (see [Membrane Anchors](#)).

**Figure 1.** Modification of C-terminal cysteine residues by prenyl groups. The C-terminal cysteine residue of the protein is outlined by the dotted line. The thiol group is thioether-linked to either a farnesyl or a geranylgeranyl group, and the exposed carboxyl group is methylated.



**Table 1. Examples of Prenylated Proteins**

|   |
|---|
| <i>Farnesylated</i>                           |
| Ras proteins                                  |
| Transducin $\gamma$ subunit                   |
| Rhodopsin kinase                              |
| Nuclear lamins A and B                        |
| Fungal mating pheromones <sup>a</sup>         |
| <i>Geranylgeranylated</i>                     |
| $\beta$ subunits of heterotrimeric G-proteins |
| Ras-related G-proteins (Rho/Rac/Rap/Ral/Rab)  |

<sup>a</sup> Peptides.

Isoprenoids are branched unsaturated hydrocarbons that are synthesized in eukaryotic cells from [acetyl Coenzyme A](#) by the first part of the metabolic pathway that is used to synthesize cholesterol and other sterols. Attachment of isoprenoids to proteins is a post-translational process with four main steps: 1) recognition of the C-terminal sequence by one of three distinct prenyltransferases (1); 2) prenylation of a cysteine residue(s) located at or close to the C-terminus using farnesylpyrophosphate or geranylgeranylpyrophosphate as the substrate; 3) proteolysis of the C-terminal residues exposes the carboxyl group on the prenylated cysteine; and 4) the isoprenylated cysteine is recognized by a methyltransferase, which methylates the carboxyl group using S-adenosyl methionine as the methyl donor. Steps 1) to 3) take place in the cytosol, whereas step 4)

occurs on the cytoplasmic surface of the endoplasmic reticulum or the plasma membrane. Thus efficient methylation requires prior isoprenylation to localize the protein at the membrane surface. The thioether linkage between the cysteine and the prenyl group is chemically very stable and probably not subject to metabolic turnover. However, the carboxylic ester linkage to the methyl group is relatively labile, and may be removed after attachment. These steps differ substantially between proteins, depending on the sequence motif at the C-terminus:

1. *Cys-a-a-X*. If X is **serine**, **methionine**, or [glutamine](#), it is recognized by farnesyl transferase, and the cysteine residue will be farnesylated. If X is leucine, it is recognized by geranylgeranyltransferase I, and the cysteine residue will be geranylgeranylated. The identity of the “a” residues (usually aliphatic) is less important, but can influence whether isoprenylation takes place or not. Farnesyl transferase and geranylgeranyltransferase I are both heterodimers; they have identical a subunits, whereas the a subunits have only 30% identity. Farnesylation can also occur at the C-terminus of a variety of fungal mating pheromone peptides, and in yeast the same enzyme is used for farnesylating both proteins and peptides. Although farnesyl groups have relatively low affinity for membranes themselves, they can enhance the membrane association due to other lipid groups. Farnesyl groups, because of their small size, may also play an important role in protein–protein interactions by binding directly to specific sites on other proteins ([2](#), [3](#)).
2. *Cys-Cys*, *Cys-X-Cys* or *Cys-Cys-X-X*. These double cysteine motifs are restricted to the Rab subgroup of Ras-related small **G-proteins**. The Rab protein first forms a complex with Rab escort protein (REP). The Rab–REP complex is then recognized by geranylgeranyltransferase II. After prenylation, REP remains bound to Rab until it is delivered to the membrane. REP probably has a dual role: recognition of Rab and masking the two geranylgeranyl groups until they can be inserted into the appropriate membrane. Both cysteines are geranylgeranylated, and consequently proteolysis cannot occur. The C-terminus is not methylated in those Rab proteins ending with the sequence Cys-Cys ([4](#)).

Many of the prenylated proteins are involved in [signal transduction](#) or **vesicle** traffic, and the prenyl group, by facilitating rapid and reversible binding to membranes, plays an essential role in these functions ([5](#), [6](#)). The membrane affinity of the prenylated proteins can be influenced by four different mechanisms (for a general discussion of factors which can affect membrane affinity of lipid anchored proteins, see [Membrane Anchors](#)):

1. The attachment of a palmitate residue (see [Palmitoylation](#)) to a cysteine close to the C-terminus reinforces the binding (eg, as in H- or N-Ras). Palmitoylation only occurs in membranes, however, so prenylation is required for it to take place ([7](#)).
2. The presence of basic residues close to the C-terminus will result in electrostatic attraction to the negatively charged bilayer surface (as in K-Ras) and increase membrane affinity ([8](#)).
3. Methylation converts the C-terminal residue from a negatively charged, hydrophilic group to an uncharged, hydrophobic group and increases membrane affinity approximately 10-fold ([5](#), [6](#)). The increase in affinity is due to the hydrophobicity of the methyl group, rather than a reduction in electrostatic repulsion, because methylation gives comparable increases in binding to uncharged membranes. Methylation can have a profound influence on the cellular distribution of farnesylated proteins, because the farnesyl group is too short to provide an effective anchor by itself. Turnover of the methyl group has also been observed, and it is possible that repeated cycles of methylation and demethylation are used to regulate protein function.
4. The membrane affinity will be reduced by soluble carrier proteins, which are able to bind to the isoprenyl group(s) and mask them from the aqueous environment. This mechanism is important for the repeated releasing and recycling of Rab proteins during membrane vesicular traffic processes ([9](#), [10](#)).

1. P. J. Casey and M. C. Seabra (1996) *J. Biol. Chem.* **271**, 5289–5292.
2. M. Sinensky (2000). *Biochim. Biophys. Acta* **1484**, 93–106.
3. R. M. Barton and H. J. Worman (1999). *J. Biol. Chem* **274**, 30008–30018.
4. J. B. Pereira-Leal, A. N. Hume, and M. C. Seabra (2001). *FEBS Lett.* **498**, 197–200.
5. S. Shahinian and J. R. Silvius (1995). *Biochemistry* **34**, 3813–3822.
6. J. R. Silvius and F. l'Heureux (1994). *Biochemistry* **33**, 3014–3022.
7. L. Liu, T. Dudler, and M. H. Gelb (1996). *J. Biol. Chem.* **271**, 23269–23276.
8. K. Cadwallader, H. Paterson, S. G. MacDonald, and J. F. Hancock (1994). *Mol. Cell. Biol.* **14**, 4722–4730.
9. S. R. Pfeffer, A. B. Dirac-Svejstrup, and T. Soldati (1995). *J. Biol. Chem.* **270**, 17057–17059.
10. G. J. Quellhorst Jr, C. M. Allen, and M. Wessling-Resnick (2001). *J. Biol. Chem*, in press.

### Suggestions for Further Reading

11. R. S. Bhatnagar and J. I. Gordon (1997) Understanding covalent modification of proteins by lipid: Where cell biology and biophysics mingle. *Trends Cell Biol.* **7**, 14–20.
12. S. Clarke (1992) Protein isoprenylation and methylation at carboxyl-terminal cysteine residues. *Annu. Rev. Biochem.* **61**, 355–386.
13. F. L. Zhang and P. J. Casey (1996) Protein prenylation: Molecular mechanisms and functional consequences. *Annu. Rev. Biochem.* **65**, 241–269.
14. M. H. Gelb (1997) Protein prenylation, et cetera: Signal transduction in two dimensions. *Science* **275**, 1750–1751.
15. P. J. Casey and J. E. Buss (1995) Lipid Modifications of Proteins, *Meth. Enzymol.* **250**.

### Pre-Protein, Pre-Pro-Protein

Many [proteins](#) are synthesized as precursors in the form of pre-proteins or pre-pro-proteins that carry pre-sequences and or pro-sequences (see [Pro-Sequence](#)). Pre-sequences usually function as **signal peptides** for **protein targeting**, while pro-sequences play a crucial role in the folding of pro-proteins, whose examples are increasing.

Most proteins destined for [protein secretion](#) are synthesized in the cytoplasm as protein precursors (pre-proteins) containing an additional *N*-terminal sequence called the signal peptide, which plays a crucial role in protein recognition of the cell secretory machinery and in initiation of protein membrane translocation (1). The signal peptides vary between 18 and 35 amino acids and have no extensive sequence homology, but share similarities in their amino acid properties. The canonical signal peptide contains three characteristic regions: an *N*-terminal domain positively charged with one or two basic residues, a hydrophobic core, and a polar *C*-terminal domain containing the signal peptide cleavage site (2). Introduction of negatively charged and **hydrophobic** amino acids into the *N*-terminal basic region reduces the efficiency of secretion (3), which probably interact with anionic phospholipids of the [membrane](#) (4). Substitutions of hydrophilic amino acids for hydrophobic core residues abolish protein membrane translocation, suggesting that this region is involved in the interaction with components of secretory machinery.

Upon translocation, the signal peptide is cleaved by a membrane bound [signal peptidase](#) (5), and the pre-protein precursors are converted into mature proteins. Signal peptide cleavage sites are rather

regular and described by the “-3, -1 rule” (6) or “A-X-B” model (7) originally proposed by statistical evaluation of primary structures of known signal peptides. This model predicts the presence of small neutral residues at positions -3 and -1, whose structural regularity may be necessary for recognition of the signal peptide cleavage site by signal peptidase, although not essential for protein translocation itself (8). Consistently, amino acid substitutions at these positions of precursors prevent the processing.

Pre-pro-proteins have additional amino acid stretches (pro-peptides or pro-sequences) located between the signal peptide and the mature part of the protein (see [Pro-Protein](#)).

### Bibliography

1. J. W. Izard and D. A. Kendall (1994) *Mol. Microbiol.* **13**, 765–773.
2. G. von Heijne (1990) *J. Membr. Biol.* **115**, 195–201.
3. G. P. Vlasuk, S. Inouye, H. Ito, K. Itakura, and M. Inouye (1983) *J. Biol. Chem.* **258**, 7141–7148.
4. M. A. Nesmeyanova et al. (1997) *FEBS Lett.* **403**, 203–207.
5. R. E. Dalbey and G. von Heijne (1992) *Trends Biochem. Sci.* **17**, 474–478.
6. G. von Heijne (1983) *Eur. J. Biochem.* **133**, 17–21.
7. D. Perlman and H. O. Halvorson (1983) *J. Mol. Biol.* **167**, 391–409.
8. J. D. Fikes, G. A. Barkocy-Gallagher, D. G. Klapper, P. J. Bassford, Jr. (1990) *J. Biol. Chem.* **265**, 3417–3423.

### Suggestions for Further Reading

9. J. W. Izard and D. A. Kendall (1994) *Mol. Microbiol.* **13**, 765–773.
10. J. Eder and A. R. Fersht (1995) *Mol. Microbiol.* **16**, 609–614.
11. A. L. Karamyshev, Z. N. Karamysheva, A. V. Kajava, V. N. Ksenzenko, and M. A. Nesmeyanova (1998) *J. Mol. Biol.* **277**, 859–870.

## Pre-Replicative Complexes

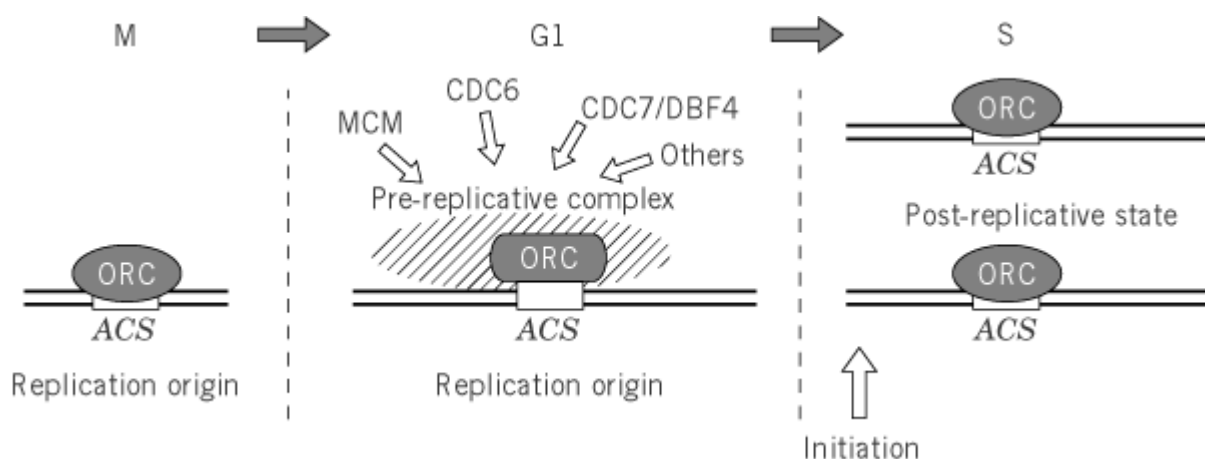
Prior to [initiation of DNA replication](#), the [replication origins](#) must be prepared by forming DNA–protein complexes at the sites during the pre-replication period. The concept of a pre-replication stage arose from reconstitution studies of replication reactions *in vitro* using the replication origins of *Escherichia coli* (oriC), **bacteriophage** lambda and [SV40 virus](#) (1, 2). Dissection of these replication reactions revealed that they are composed of multiple steps and in the first step the initiator protein assembles at the origin DNA and activates it in an ATP-dependent manner. Formation of this protein–DNA complex induces a change in the structure of the DNA around the origin and directs it to accept several proteins necessary for DNA synthesis: [DNA helicases](#), **primases**, and **DNA polymerases**. The initial active-complex state is relatively stable and retained for a certain period if the later DNA synthesis proteins are absent. Consequently, this initial complex can be isolated and is called a *pre-initiation* or *pre-priming* complex.

In eukaryotic [chromosomes](#), preparation for triggering of the replication origins involves a long period and complicated procedures, since DNA synthesis from a replication origin has to occur within a narrow window of time during one [cell cycle](#). Furthermore, the initiation site has to be distinguished clearly throughout the cell cycle from other regular chromosomal regions where stable

[chromatin](#) structures are formed. Actually, the predicted licensing mechanism that restricts the firing of each replication origin to once a cell cycle (see [Licensing Factor](#)) proposes a marking mechanism of potential origin sites at the end of mitosis to initiate DNA synthesis in the next **S phase**. Therefore, eukaryotic cells seem to start the preparation much earlier than the actual initiation time by establishing the competency for the initiation in parallel with the progression of cell-cycle events during the **G1 phase**. The molecular features of DNA–protein complexes at replication origins during the pre-replication period have been elucidated only in the yeast *Saccharomyces cerevisiae* (for a review, see (3)). A protein complex that binds specifically to the yeast replication origin ([autonomously replicating sequence](#), ARS), the ARS conserved sequence (ACS) has been isolated and called the [origin recognition complex](#) (ORC). It binds to the origin throughout the cell cycle in an ATP-dependent manner and has an [ATPase](#) activity that is regulated by the ACS DNA. Thus, it is thought that ORC functions as a landmark for replication origins in the huge chromosomal structure and also has some active role in activating the replication origins. ORC is necessary to initiate DNA synthesis from the origins but is not sufficient for activation. Additional factors, MCM proteins (licensing factor), CDC6, and CDC7/DBF4 kinase (activators), interact with ORC and will be directly involved in the activation step.

The DNA–protein complex at the ARS was characterized by *in vivo* DNase I [footprinting](#), which observed both an ORC-specific protection pattern and a second protection covering a wider region. The second protection pattern appeared at the end of mitosis, persisted throughout the G1 phase, and disappeared right after entry into the S phase. The wider nature of the protection pattern seemed to correspond to the formation of a larger protein complex on the ACS-bound ORC in the pre-replication stage. Further analysis demonstrated that the factors interacting with ORC, as indicated above, are all required for the formation and maintenance of the wider protection in the G1 phase and that the same factors are essential for the initiation of DNA replication from the origin (4). Therefore, it is most likely that these ORC-interacting factors form a large complex with ORC at the ACS, and the assumed complex built up on ORC early in the G1 phase is considered the pre-replicative complex competent to initiate DNA replication. During progression from the G1 to S phase, some highly programmed signaling mechanism transduces the final signal to each origin to initiate at the right time (Fig. 1). Then the pre-replicative complex will be disrupted and shifted to the post-replicative complex, simultaneously with the initiation of the origin.

**Figure 1.** Assembly of a pre-replicative complex at the yeast replication origin in G1 phase.



Knowledge about the actual assembly of the pre-replicative complexes in nuclei is limited. It was previously reported that replicating DNA and several replication proteins colocalize in nuclei and form [replication foci](#) in the S phase. It might be expected that pre-replicative complexes will

assemble at the same sites prior to the formation of replication foci. However, no components in the pre-replicative complexes reported in yeast have been demonstrated to be involved in such foci prior to the initiation of the S phase. Exceptionally, a eukaryotic single-stranded DNA-binding protein, RPA, exhibits punctuated assembly in nuclei prior to DNA replication in the *Xenopus* egg extract reaction (5). RPA is an essential component for [SV40 virus](#) DNA replication *in vitro*, and it facilitates the formation of unwound DNA at the SV40 replication origin in the pre-initiation stage. These RPA foci were formed at the end of mitosis, colocalized with some early replication foci upon the initiation of the S phase, and were linked with replication-associated unwinding of the DNA. Thus, they are called *pre-replication foci* (or *pre-replication centers*), which serve as the precursors for replication foci and will be one reflection of the formation of the pre-replicative complex or its closely related protein assembly, although there is no direct evidence to connect them.

## Bibliography

1. J. F. X. Diffley and B. Stillman (1990) *Trends Genet.* **6**, 427–432.
2. T. A. Baker and A. Kornberg (1992) In *DNA Replication*, 2nd ed. W. H. Freeman, New York, pp. 521–545.
3. S. P. Bell (1995) *Curr. Opin. Gen. Dev.* **5**, 162–167.
4. J. F. Diffley et al. (1994) *Cell* **78**, 303–316.
5. Y. Adachi and U. K. Laemmli (1994) *EMBO J.* **13**, 4153–4164.

## Pribnow Box

The Pribnow box is a conserved promoter element that is recognized by the main form of *E. coli* RNA polymerase holoenzyme (see [TATA Box](#), Bacterial promoters). It is also called the -10 region to indicate its distance from the start site of transcription. This element was independently recognized as a region of sequence conservation in the mid-1970s by Drs. David Pribnow and Heinz Schaller when they visually inspected the first available sequences. Computer analysis of several hundred promoter sequences has established the sequence TATAAT (on the nontemplate strand) as the consensus sequence of the Pribnow box. Many eukaryotic and archaeal promoters have a sequence similar to that of the Pribnow box at approximately -30, with consensus sequence TATAAA.

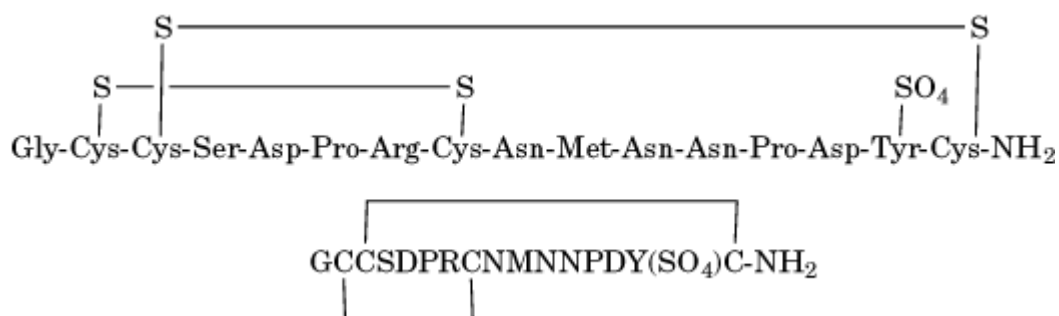
## Primary Structure

Protein structure is classified in a hierarchical manner into [primary structure](#), **secondary structure**, [tertiary structure](#), and [quaternary structure](#). The primary structure of a protein is the sequence, or order, of [amino acid](#) residues in the [polypeptide chain](#) and can be represented in a variety of ways (two examples are given in Fig. 1). By convention, the sequence is numbered from the [N-terminus](#) to the [C-terminus](#) of the polypeptide chain. The primary structure also includes information about the number of polypeptide chains in a protein and any covalent modifications, such as [disulfide bond](#) formation, **phosphorylation**, [sulfation](#), or **glycosylation** (see [N-Glycosylation](#) and [O-Glycosylation](#)). In addition, nonprotein groups, or **prosthetic groups**, such as heme, metal ions, or pigments, form



part of the primary structure of a protein. The significance of a protein's primary structure is that it determines both the three-dimensional structure and the function of that protein.

**Figure 1.** Alternative representations of the primary structure of the 16 residue peptide  $\alpha$ -conotoxin EpI (1). The upper representation gives the amino acid sequence using the three-letter code, the lower representation uses the one-letter code for different amino acid residues. Both representations show the presence of two disulfide bonds, [sulfation](#) of the [tyrosine](#) residue, and amidation of the C-terminus (see [Post-Translational Modifications](#)).



Determination of a protein's primary structure can be done either directly using a combination of biochemical and chemical techniques (see [Protein Sequencing](#)) or, more often, indirectly by identifying the nucleotide sequence of the corresponding **gene** or complementary DNA. If two different proteins have similar primary structures, they are said to be **homologous** and are likely to have similar tertiary structures and functions. [Sequence databases](#) such as SWISS-PROT (1) hold information about the primary structures of many thousands of proteins and can be searched to identify homologous proteins.

*In vivo*, the primary structure of each protein is genetically encoded. [Mutations](#) in the [genetic code](#), or errors in reading it, can lead to changes in the protein primary structure. Such variations from the normal or **wild-type** sequence can result in changes to both the structure and function of a protein and can affect the viability of the cell and/or the organism. On the other hand, changes may be engineered into the primary structure of a protein *in vitro* by [site-directed mutagenesis](#) to investigate the role of specific residues.

Theoretically, it is possible to predict the tertiary structure of a protein directly from knowledge of its primary structure (see [Protein Structure Prediction](#)). In practice, however, this so-called **protein folding** problem is not yet solved, except in those instances where protein structures can be predicted by **homology modeling**.

[See also [Protein Structure](#).]

#### Bibliography

1. A. Bairoch and R. Apweiler (1988) *Nucleic Acids Res.* **26**, 38–42.
2. M. Loughnan et al. (1998) *J. Biol. Chem.* in press.

#### Suggestions for Further Reading

3. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
4. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, New York.
5. N. J. Darby and T. E. Creighton (1993) *Protein Structure*, IRL Press, Oxford, U.K.

## Primase

[Initiation Of DNA Replication](#) by **DNA polymerases** requires **primers** that provide a 3'-OH terminus on the [template](#) DNA. In some particular systems, long RNA **transcripts**, DNA fragments, or [proteins](#) covalently linked with nucleotides function as primers for DNA synthesis, but in general this function is performed by short RNA fragments newly synthesized on the template. During [DNA replication](#), the two antiparallel DNA strands in double-stranded DNA have two different modes of synthesis, in which one strand is synthesized continuously (**leading strand**) and another discontinuously (**lagging strand**; see also [Okazaki Fragments](#)). Since the latter mechanism requires frequent priming in certain intervals on a DNA strand, an [enzyme](#) that is able to produce short RNA primers repeatedly is an essential component of most DNA replication systems. Indeed, all cells have an essential specialized enzyme, called the *primase*, that synthesizes short RNA primers for DNA replication. Originally, this enzyme was discovered through studies on the initiation of replication of *E. coli* single-stranded DNA **bacteriophages**. Among more than 10 replication proteins needed for the replication of these phages, the DnaG protein was identified as a factor essential for the initiation stage. Subsequent reconstitution of the reaction with purified proteins demonstrated that the role of this protein is the production of short RNA segments on a single-stranded DNA, which is subsequently used for DNA synthesis by DNA polymerase (1).

DnaG primase is inert on its own, and additional proteins cooperate in the priming process. Reconstitution experiments revealed that a mobile protein complex called the *primosome* is assembled and migrates along the template (2). One major member of primosome is the DnaB [DNA helicase](#), which interacts functionally with primase and activates the enzyme. A *primosome assembly site* (PAS) was isolated from the *E. coli* fX174 phage genome as a characteristic stem-loop structure. An [ATPase](#) named PriA recognizes the PAS and forms a complex with two other primosome assembly proteins (PriB, C). Subsequently, DnaB helicase and DnaG primase are recruited to the complex. In the case of replication origins, complexes of initiator proteins also assemble the DnaB-DnaG primosome, as reported for *E. coli* oriC or the [lambda phage](#) origin (for a review, see (3)). The assembled primosome migrates on the template strand in a 5' to 3' direction, hydrolyzing ATP and synthesizing many primers. The direction of migration is opposite to that of DNA synthesis, which explains the mechanism of priming for the discontinuous, lagging strand. During DNA synthesis, the primosome will be a part of the [replication fork](#) complex and will manage the synthesis of primer RNA by the interval corresponding to the Okazaki fragment length. A similar tight link of primases and DNA helicases, and their involvement in the replication fork complexes, has also been reported in bacteriophage T4 and T7 replication (4, 5).

Eukaryotic primases have slightly different features from those of prokaryotic systems. They were detected as a component of DNA polymerase  $\alpha$ , which is an essential DNA polymerase for chromosomal replication. This polymerase has a large subunit of about 180 kDa, a middle subunit of about 70 kDa, and two small subunits of 60 to 50 kDa (6). The primase activity is detected in the subcomplex of two small subunits, and maintenance of the complex is necessary for this activity. The primase catalytic activity is harbored in the 50-kDa subunit, and the 60-kDa one functions as its accessory subunit. Studies of the DNA elongation process using an *in vitro* SV40 replication reaction revealed that eukaryotes have multiple DNA polymerases and the DNA polymerase  $\alpha$  complex is specialized for the synthesis of Okazaki fragments (7). Thus, the association of primase with DNA polymerase  $\alpha$  would be an adapted feature of eukaryotes to synthesize Okazaki fragments efficiently by specializing one of many DNA polymerases for the priming process, although the mechanism of switching from the priming reaction to DNA synthesis in one protein complex is still unclear.

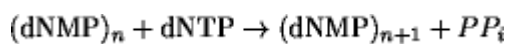
The nature of the sequence specificity of the primer RNA to initiate primer synthesis, and its ultimate length, have been studied with several *E. coli* replicons. The results indicate that the preferred sequence for DnaG priming is GTC, and most of the primers have pppApG at their 5' end, which initiated from base pairing with the middle T. RNA–DNA junctions are observed at fixed sites from the primer start, and the primer length is predominantly  $11 \pm 1$  nucleotides. Priming by T4 or T7 phage primers also has a sequence preference, but the sequences are slightly different from those of DnaG. The nature of the preferred sequence for eukaryotic primase has not been elucidated well, but the length of the primer is 12 to 14 nucleotides for *Drosophila* primase and 8 to 12 nucleotides for mouse or yeast primases (1).

## Bibliography

1. T. A. Baker and A. Kornberg (1992) *DNA Replication*, 2nd ed., W. H. Freeman, pp. 275–298.
2. K.-I. Arai and A. Kornberg (1981) *J. Biol. Chem.* **256**, 5260–5266.
3. K. J. Marians (1992) *Ann. Rev. Biochem.* **61**, 673–719.
4. M. Venkatesan et al. (1982) *J. Biol. Chem.* **257**, 12426–12434.
5. J. Bernstein and C. C. Richardson (1989) *J. Biol. Chem.* **264**, 13066–13073.
6. W. C. Copeland and T. S.-F. Wang (1993) *J. Biol. Chem.* **268**, 26179–26189.
7. T. Tsurimoto et al. (1990) *Nature* **346**, 534–539.

## Primer

**DNA polymerase** plays a primary role in [DNA replication](#), catalyzing the synthesis of a DNA chain complementary to the [template](#) DNA strand as follows:



The reaction absolutely requires three components: (1) a deoxynucleosidetriphosphate (dNTP) as a substrate, (2) single-stranded DNA as a template, and (3) a short RNA or DNA as a *primer* (1). No known DNA polymerase is able to initiate DNA synthesis without a primer. The primer must be annealed to the template by base-pairing, and its 3'-terminus must possess a free 3'-OH group. Chain growth is exclusively  $5' \rightarrow 3'$  since the polymerization mechanism is a nucleophilic attack by the 3'-OH group of the primer on the  $\alpha$ -phosphate of the incoming dNTP. The products are a new primer, longer by one nucleotide, and one inorganic pyrophosphate. A 2',3'-dideoxynTP can be incorporated by some DNA polymerases, but it blocks chain elongation because it lacks a free 3'-OH group. This chain termination property is utilized in some DNA sequencing methods (2) (see **Sanger method** and **Sequencing, DNA**).

Why does DNA polymerase require the primer for chain elongation? It has been thought that chain growth by the covalent extension of primer makes it possible for DNA polymerase to copy DNA very accurately. All the replicative DNA polymerases possess proofreading capacity, which allows for removal of nucleotides incorrectly inserted by the polymerase and increases the fidelity of DNA synthesis. The proofreading process involves melting the primer terminus to translocate it from the polymerase active site to the  $3' \rightarrow 5'$  exonuclease active site (3). Since the distance between both active sites is about four nucleotides in length, the elongating DNA in the molecule of DNA polymerase must be longer than four residues. If not, the short piece of DNA would not be set back properly on the template sequence. It is also known that the 3'-end of RNA annealed to template

DNA is not subjected to proofreading by DNA polymerase. Thus, the DNA polymerase can start DNA synthesis from primer RNA without an abortive proofreading process occurring. It should be noted that the primer RNA, which is made with a much higher error rate than DNA, is removed upon completion of [Okazaki Fragments](#).

A common primer for the DNA replication of **chromosomal** DNA is a short RNA transcript synthesized *de novo* by a [primase](#) in both prokaryotic and eukaryotic cells (pp. 279–297 (1)). Numerous **bacteriophage** and **plasmids**, as well as some animal **viruses**, possess their own primase for DNA replication. This is probably due to the **replicon**-specific replicative apparatus. In prokaryote cells, primase is usually a single polypeptide chain, and its priming function requires other factors, forming a multiprotein complex called the *primosome* on single-stranded template DNA (4). A transient and functional interaction between primase and replicative DNA helicase is important for the priming reaction. In eukaryotes, however, primase is associated with DNA polymerase activity, as an integral part of DNA polymerase  $\alpha$  (5). Since no eukaryotic DNA helicase has been assigned as the replicative helicase, the interaction between primase and helicase is not known. In contrast, a tight complex of primase and helicase, and even a single polypeptide chain carrying both primase and helicase activities, has been found to be a priming enzyme for bacteriophages (6).

Primases can start a new chain when copying a duplex DNA but, like DNA polymerases, they require single-stranded DNA as the template (7). Primases absolutely require an rNTP (ATP in most cases) to initiate the primer, but they can utilize dNTPs in place of rNTPs. Under physiological concentrations, rNTPs are preferentially used during primer synthesis. In prokaryotes, primer synthesis is initiated at preferred sequences on the template: 3'-GTC for *E. coli* primase (DnaG), 3'-CTG(G/T) for T7 phage primase (gp4), and 3'-TTG for T4 phage primase (gp61). The sequence preference of the initiation of primer synthesis by eukaryotic primases is uncertain. All primases found in various organisms extend the primer chain in the 5' → 3' direction to a particular length: 10 to 12 nucleotides by DnaG protein, 4 nucleotides by T7 gp4, 5 nucleotides by T4 gp61, and 8 to 12 nucleotides by eukaryotic primases (8-11).

There are several ways to make primer other than the priming process by primase. It has been shown that **RNA polymerase** can act as a priming enzyme, producing a special kind of RNA primer for the initiation of replication of some bacteriophage and plasmid DNA (12, 13). In the DNA replication of particular types of bacteriophage (such as *Bacillus subtilis* phage  $\phi$ 29) and animal viruses (such as [Adenovirus](#)), protein priming promotes the initiation step, in which the 3'-OH group of an [amino acid](#) residue within the termination protein is recognized as a primer end by DNA polymerase (14). The end of a preexisting DNA strand can also serve as the primer. The 3'-OH terminus generated in nicks or gaps of DNA is a common primer for DNA synthesis to complete [DNA repair](#) or the [recombination](#) processes.

## Bibliography

1. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York, pp. 106–109.
2. F. Sanger, S. Nicklen, and A. R. Coulson (1977) Proc. Natl. Acad. Sci. USA **74**, 5463–5467.
3. D. L. Ollis, P. Brick, R. Hamlin, N. G. Xuong, and T. A. Steitz (1985) Nature **313**, 762–766.
4. K. Arai, R. Low, J. Kobori, J. Shlomai, and A. Kornberg (1981) J. Biol. Chem. **256**, 5273–5280.
5. R. C. Conaway and I. R. Lehman (1982) Proc. Natl. Acad. Sci. USA **79**, 2523–2527.
6. J. A. Bernstein and C. C. Richardson (1989) J. Biol. Chem. **264**, 13066–13073.
7. L. Rowen and A. Kornberg (1978) J. Biol. Chem. **253**, 758–764.
8. T. Kitani, K. Yoda, T. Ogawa, and T. Okazaki (1985) J. Mol. Biol. **184**, 45–52.
9. S. Tabor and C. C. Richardson (1981) Proc. Natl. Acad. Sci. USA **78**, 205–209.
10. Y. Kurosawa and T. Okazaki (1979) J. Mol. Biol. **135**, 841–861.

11. B. Y. Tseng and C. N. Ahlem (1983) *J. Biol. Chem.* **258**, 9845–9849.
12. J. M. Kaguni and A. Kornberg (1982) *J. Biol. Chem.* **257**, 5437–5443.
13. Y. Sakakibara and J. Tomizawa (1974) *Proc. Natl. Acad. Sci. USA* **71**, 1403–1407.
14. J. M. Hermoso, E. Mendez, F. Soriano, and M. Salas (1985) *Nucl. Acid Res.* **13**, 7715–7728.

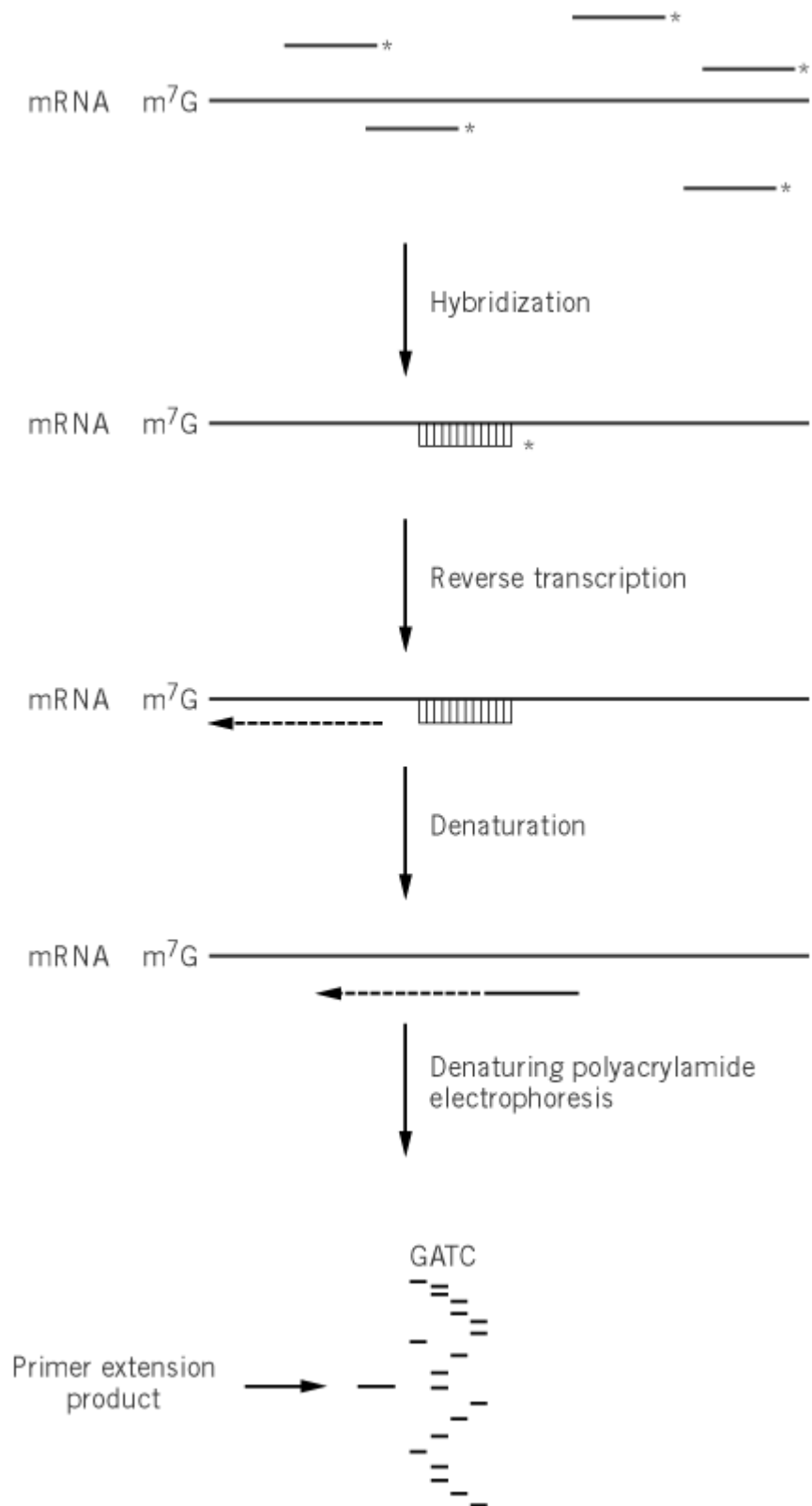
### Suggestions for Further Reading

15. M. L. DePamphilis, ed. (1996) *DNA Replication in Eukaryotic Cells*, Cold Spring Harbor Laboratory Press, New York.
16. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.

### Primer Extension

Primer extension is a technique that can be used for a variety of purposes, such as mapping the 5' ends of [messenger RNA](#) and the ribonuclease cleavage sites in RNA. A DNA oligonucleotide complementary to the RNA is synthesized and radiolabeled at the 5' end with  $g^{32}P$  ATP and **polynucleotide kinase** (Fig. 1). The primer is then annealed to the RNA, followed by the addition of nucleoside triphosphates and **reverse transcriptase** (RT). The reverse transcriptase synthesizes the complementary strand until it falls off the [template](#) RNA at its 5' end. The primer-extended product is then purified and subjected to [electrophoresis](#) on a denaturing [polyacrylamide](#) gel to determine its length. In the case of mRNA, a DNA **sequencing** reaction is carried out in parallel on the genomic sequences encoding the mRNA, which permits the site at which the reverse transcriptase stopped, that is, the site of initiation of [transcription](#), to be mapped to the exact nucleotide (1, 2).

**Figure 1.** The primer extension reaction. A small DNA oligonucleotide, labeled at the 5' end (\*) with a 5' radioactive phosphate, is annealed to the RNA. The primer is extended using reverse transcriptase and nucleoside triphosphate. The extended products are resolved by denaturing polyacrylamide gel electrophoresis (PAGE) and visualized by [autoradiography](#). The exact size of the site can be judged by reference to a sequencing ladder generated from a corresponding fragment of the gene. For example, see reference (2), where this technique was used to map cleavage sites in ribosomal RNA precursors.



### Bibliography

1. T. Maniatis, E. F. Fritsch, and J. Sambrook (1989) *Molecular Cloning: A Laboratory Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 7.79–7.81.
2. C. Allmang, Y. Henry, J. P. Morrissey, H. Wood, E. Petfalski, and D. Tollervey (1996) *RNA* **2**, 63–73.



## Prion

The term “prion” was coined by Prusiner (1) as a shortened form of “proteinaceous infective agent”, to describe the agent he thought responsible for the transmissibility of certain spongiform encephalopathies such as scrapie, BSE, or “Mad Cow” disease and Creutzfeldt Jakob disease (CJD) (see Scrapie). What makes the prion concept so interesting and controversial is that it breaks with current orthodoxy that nucleic acids are necessary for the process of disease infection. Prusiner has long maintained that the evidence of scrapie infection ruled out involvement of nucleic acids, leaving a self-replicating proteinaceous agent as the only other possibility, however heterodox that might seem. The threat that these diseases present to the human population, their insidious nature involving very long incubation periods, the uncertain routes of infection, and their particularly unpleasant symptoms has heightened general public concern in their study and the parallel search for therapies. For his persistence in exploring the prion hypothesis and demonstrating its apparent correctness, Prusiner was awarded a Nobel Prize in 1997.

The scrapie agent was first discovered in the brains of the Syrian hamster and identified by antibody-labelling with the prion protein in both infected and uninfected brains. This suggested that the agent was a single protein that exists in a different isoform from that in healthy tissue. Subsequent studies have shown that the prion protein is responsible for both toxicity and infectivity of scrapie. The prion protein appears to occur in two distinct forms: PrP<sup>C</sup>, the normal cellular form and PrP<sup>Sc</sup>, the infective scrapie form. The two forms can be readily distinguished by their different sensitivities to proteases: PrP<sup>C</sup> is a 33–35 kDa protein rapidly degraded by protease K, while PrP<sup>Sc</sup> is more stable to this protease forming a 27–30 kDa proteinase-resistant core, known as PrP27–30. In the presence of detergents, PrP<sup>C</sup> remains soluble, while PrP27–30 polymerizes into amyloid fibrils (2, 3) (see Amyloid). As the PrP is a single copy gene encoded on a single exon, these different forms do not derive from differential splicing. The difference between the PrP<sup>C</sup> and PrP<sup>Sc</sup> forms appears to reside in the three-dimensional, rather than the covalent, structures of the proteins. Spectroscopic examination of PrP<sup>C</sup> and PrP<sup>Sc</sup>, using Fourier transform infra-red spectroscopy (FTIR) and circular dichroism, shows that PrP<sup>C</sup> contains around 42%  $\alpha$ -helix but only about 3%  $\beta$ -sheet (4), while in PrP<sup>Sc</sup> the helix content declines to 20% to 30%, and the  $\beta$ -sheet rises to about 43% and even further to 47% to 54% in the PrP27-30 proteolytic product (5-7). The high  $\alpha$ -helical content of PrP<sup>C</sup> form has been confirmed by an NMR determination of its largest folding domain, residues 121–231 (8), which shows a globular structure with three  $\alpha$ -helices and a small two-stranded  $\beta$ -sheet. Solid-state NMR of the adjacent peptide, residues 109–122, suggests it has a  $\beta$ -structure under some conditions and an  $\alpha$ -helix under others (9). This supports the proposal that the change associated with infectivity has the character of an  $\alpha$ - to  $\beta$ -structural conversion. It is interesting to note that these studies have not uncovered any unusual structural feature in the prion protein that might explain the structural conversion.

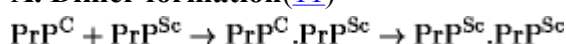
In Prusiner’s view, infectivity of the prion occurs because PrP<sup>Sc</sup> molecules can convert normal PrP<sup>C</sup> molecules into the PrP<sup>Sc</sup> form, which, in turn, may convert further cellular prions into the scrapie form, as in some kind of chain reaction. How this process occurs *in vivo* is not absolutely certain, but it has been found possible in the presence of guanidinium to slightly destabilize the proteins and generate a proteinase-resistant form of the prion protein from a simple mixture of PrP<sup>C</sup> and PrP<sup>Sc</sup> *in vitro* (10). One possible mechanism for the conversion of PrP<sup>C</sup> to the PrP<sup>Sc</sup> form by PrP<sup>Sc</sup> is through the formation of a mixed dimer (11) (see Scheme A in Table 1). Alternatively, a hypothesis of nucleated, or seeded, polymerization has been proposed (13) in which the structural conversion involves the recruitment of PrP<sup>C</sup> molecules into a growing PrP<sup>Sc</sup> polymer (see Scheme B in Table 1). Evidence is increasing that the PrP<sup>Sc</sup> form is not dimeric but at least tetrameric or oligomeric (14,



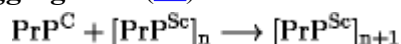
15). As PrP<sup>27–30</sup> is known to form amyloid and amyloid formation from  $\alpha$ -helical proteins is known to involve  $\alpha$ - to  $\beta$ - structural conversion (see [Amyloid](#)), an amyloid fibril may function as the prion-converting agent. This view is given support by sedimentation and ultrafiltration studies, which suggest that converting activity is associated with heterogeneous particles that are several times larger than the solubilized PrP molecule (16). On the other hand, a number of observations (17) do not correlate infectivity with the presence of visible or high molecular weight prion amyloid. However, it is not certain that the involvement of small amyloid fragments or templates that might act as initiators of prion amyloid can be entirely ruled out.

**Table 1. Proposed Mechanisms of Prion Structural Conversion**

**A. Dimer formation(11)**



**B. Aggregation (13)**



Other studies have suggested alternative mechanisms for the structural conversion of PrP<sup>C</sup> and have implicated another molecule, designated “protein X,” in the process of formation of PrP<sup>Sc</sup> from PrP<sup>C</sup> (18). It is proposed that protein X may bind to PrP<sup>C</sup> in such a way that subsequent binding of PrP<sup>Sc</sup> to produce a ternary complex, and conversion of the cellular to the scrapie form, is facilitated (19). Another proposal (20) suggests that the single disulfide bond in the prion may be destabilized and reforms as an intermolecular bond linking PrP<sup>Sc</sup> units in an amyloid fibril. This model has the additional advantage of adding yet greater stability to the already very stable amyloid fibril, thereby providing an explanation of the extreme thermal and chemical stability of the scrapie agent.

About 10% of the incidence of CJD is familial and, thus, inherited. PrP is encoded by a single gene on chromosome 20 in man, consisting of a single exon. Some 20 different mutations in the human prion gene segregate with dominantly inherited disease, and five of these have been genetically linked to human familial prion diseases (Table 2). Plotting these mutations on the known structure of the 121–231 domain of the PrP<sup>C</sup> molecule (8) shows a tendency to cluster within, or close to,  $\alpha$ -helices, possibly stabilizing their conversion to  $\beta$ -structures (10). Alternatively, the placement of many variant residues in the hydrophobic core of the prion domain may result in a large-scale unfolding and refolding of the prion protein, as indicated experimentally in the  $\alpha$ - to  $\beta$ -conversion in lysozyme (21), whose proportions of  $\alpha$ -helix and  $\beta$ -sheet broadly match those of PrP<sup>C</sup> (see [Amyloid Precursor Protein](#)). This is supported by the proposal that molten globule prion intermediates (22) may be involved in the  $\alpha$ - to  $\beta$ -conversion PrP<sup>C</sup> to PrP<sup>Sc</sup> conversion, as they are in lysozyme.

**Table 2. Some Mutations in the Prion Protein in Human Prion Diseases**

**Creutzfeldt-Jakob disease Gerstmann-Sträussler-Scheinker disease**

A117V

P102L

D178N if V129\*

F198S

E200K

*Fatal Familial Insomnia*

D178N if M129\*

---

\* polymorphic site

It has also proved possible to establish transgenic animal models of the disease (1823), which have been used to test a number of important avenues of investigation. One of these is the most puzzling aspect of prion diseases: the question of strains (24, 25). The scrapie diseases, like many other infectious diseases, exist in a number of sub-forms or phenotypes, which implies that differences in the prion protein molecule are capable of being passed on. Because no nucleic acids are present in the scrapie diseases to replicate amino acid variation, how can the representation of different disease phenotypes occur in the protein-only prion system? This question needs to be considered against a clear experimental background that scrapie strains can be transmitted. In a crucial experiment (25), mice expressing a chimeric mouse-human PrP gene were inoculated with prions from two different subtypes of CJD and from Fatal Familial Insomnia (FFI). The different mutations in the PrP molecules from CJD and FFI (Table 2) result in slight differences in chain conformation, which are expressed as different fragmentation patterns when treated with proteinase K; the CJD proteins yield a 21 kDa proteinase-resistant fragment, whereas in the FFI protein it is 19 kDa. Inoculation of the mice with FFI brain extracts induced the formation of the 19 kDa PrP<sup>Sc</sup> fragment, whereas inoculation with CJD extracts induced formation of the 21 kDa fragment. The only reasonable conclusion is that this experiment demonstrates that the probably subtle structural differences in variant PrP<sup>Sc</sup> molecules can be “imprinted” on the same set of PrP<sup>C</sup> molecules in their conversion to the PrP<sup>Sc</sup> form, leading to prions carrying specific information of their strain of origin. As the resultant PrP<sup>Sc</sup> molecules are covalently identical, the differences representing the strains must reside in differences in the folding of the polypeptide chains or possibly in their pattern of glycosylation. One possible model is analogous to crystallization; crystals can replicate themselves in the presence of monomers by the process of seeding or nucleation. Different PrP<sup>Sc</sup> molecules could give rise to structurally slightly different amyloid fibrils (fibril formation can be considered as one-dimensional crystallization); fragments of these, in turn, could act as nuclei recruiting into a growing fibril only those prions of the same kind, thereby replicating themselves. Experimental evidence seems now to support this mode (26, 27).

The study of the prion has challenged the two long-held fundamental dogmas. First, that the amino acid sequence is the sole determinant of the native protein fold (28) and, second, that protein structural information can only be transmitted through the mediation of DNA or RNA. It has also indicated how protein-protein information transmission might occur. All of this is in the context of identifying protein misfolding as a hitherto unsuspected cause of certain widespread inherited and infectious diseases.

#### Bibliography

1. S. B. Prusiner (1982) *Science* **216**, 136–144.
2. M. P. McKinley, A. Taraboulos, and L. Kenaga (1991) *Lab. Invest.* **65**, 622–630.
3. S. B. Prusiner et al (1983) *Cell* **35**, 349–358.
4. K.-M. Pan, M. A. Baldwin, J. T. Nguyen, M. Gasset, A. Serban, A. Serban, D. Groth, I. Mehlhorn, Z. Huang, R. J. Fletterick, F. E. Cohen, and S. B. Prusiner (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 10962–10966.
5. B. Caughey, A. Dong, K. S. Bhat, D. Ernst, S. F. Hayes, and W. S. Caughey (1991)

Biochemistry **30**, 7672–7680.

6. M. Gasset M. A. Baldwin, R. J. Fletterick, and S. B. Prusiner (1993) Proc. Natl. Acad. Sci. U.S.A. **90**, 1–5.
7. J. T. Nguyen et al (1995) J. Mol. Biol. **252**, 412–422.
8. R. Riek, S. Horemann, G. Wider, M. Billeter, R. Glockshuber, and K. Wurthlich (1996) Nature **382**, 180–182.
9. J. Heller et al (1996) Prot. Sci. **5**, 1655–1661.
10. F. E. Cohen, K-M. Pan, Z. Huang, M. Baldwin, R. J. Fletterick, and S. B. Prusiner (1994) Science **264**, 530–531.
11. P. M. Harrison, P. Bamborough, V. Daggett, S. B. Prusiner, and F. E. Cohen (1997) Curr. Opin. Struct. Biol. **7**, 53–59.
12. S. A. Priola, B. Caughey, K. Wehrly, and B. Chesebro (1995) J. Biol. Chem. **270**, 3229–3305.
13. J. T. Jarrett and P. T. Lansbury (1993) Cell **73**, 1055–1058.
14. K. Jansen et al (2001) Biol. Chem. **382**, 683–691.
15. M. Morillas, D. L. Vanik, and W. K. Surewicz (2001) Biochemistry **40**, 6982–6987.
16. B. Caughey, D. A. Kocisko, G. J. Raymond, and P. T. Lansbury (1995) Chem. Biol. **2**, 807–817.
17. H. Wille, M. A. Baldwin, F. E. Cohen, S. J. DeArmond, and S. B. Prusiner (1996) In *The Nature and Origin of Amyloid Fibrils* (G. R. Bock and J. A. Goode, eds.), John Wiley & Sons, Chichester, England, pp. 181–199.
18. G. C. Telling, et al (1995) Cell **83**, 79–85.
19. K. Kaneko, L. Zulianello, M. Scott, C. M. Cooper, A. C. Wallace, T. L. James, F. E. Cohen, and S. B. Prusiner (1997) Proc. Natl. Acad. Sci. U.S.A. **94**, 10069–10074.
20. E. Welker, W. J. Wedemeyer, and H. A. Scheraga (2001) Proc. Natl. Acad. Sci. **98**, 4334–4336.
21. D. R. Booth et al (1997) Nature **385**, 787–793.
22. J. Safer, P. Roller, D. Gajdusek, and C. Gibbs (1994) Biochemistry **33**, 8375–8383.
23. M. Fischer et al EMBO J. (1992) **15**, 1255–1264.
24. R. A. Bessen et al (1995) Nature **375**, 698–700.
25. G. C. Telling et al (1996) Science **274**, 2079–2082.
26. T. R. Serio et al (2000) Science **289**, 1317–1321.
27. P. Chien and J. S. Weissman (2001) Nature **410**, 223–227.
28. D. L. Minor, Jr., and P. S. Kim (1996) Nature **380**, 730–734.

### **Suggestions for Further Reading**

29. S. B. Prusiner (2001) The Shattuck lecture — Neurodegenerative diseases and prions. New Engl. J. Med. **344**, 1516–1526.
30. S. B. Prusiner (1998) Prions. Proc. Natl. Acad. Sci. USA **95**, 13363–13383.
31. B. Caughey and B. Chesebro (1997) Prions: Protein and the transmissible spongiform encephalopathies. Trends Cell Biol. **7**, 56–62.
32. A. R. Clarke, G. S. Jackson, and J. Collinge (2001) The molecular biology of prion propagation. Philos. Trans. R. Soc. Lond. Biol. Sci. **356**, 185–195.

### **Probe Hybridization**

The ability of complementary nucleic acids to form hybrid duplexes has been exploited in a large number of powerful techniques for identification and manipulation of the genetic information stored in DNA and used by the cell via RNA. Most of these **hybridization** techniques rely on the use of a labeled probe nucleic acid to hybridize with a target nucleic acid. After removal of any unreacted probe, detection of the remaining labeled probe identifies and quantifies the hybrid duplex and, consequently, the regions of complementarity between the probe and the target.

If the target is one strand of a duplex, **denaturation** must precede hybridization. Usually, denaturation is accomplished by heat or by exposure to alkaline conditions. Often the number of targets is quite low, perhaps only a few copies. There are a number of amplification techniques that will produce large numbers of copies of the target and thus increase the amount of hybrid duplex produced with a concomitant increase in the observed signal. Solid-phase hybridization requires an addition preparative step, namely, immobilization of the target.

The length of a hybridization probe can vary from about 15 bases to several hundred or more, with the appropriate choice depending on the specific experiment. Long probes provide greater specificity. They can incorporate many detectable groups per probe, thereby increasing the signal intensity. Long probes typically are generated by random-primed labeling, nick translation, or **PCR** (1), and require more work to prepare than do short chemically synthesized probes.

Short probes (15–30 bases) can be synthesized in large quantities. Incorporation of nonisotopic [reporter groups](#) is easily accomplished. Short probes bind less tightly to the target than do long probes. In many circumstances weaker binding is advantageous. The hybridization experiment can then be sensitive to single defects of complementarity, including mismatches, deletions, bulges, and small internal loops.

Detection and quantification of the hybrid duplex is provided via the labeled probe. Labeling, whether isotopic or nonisotopic, is an essential feature of the probe. All the early work in hybridization was performed with isotopic labels. More recently, sensitive **chemiluminescent** and bioluminescent detection schemes have come into common use. The nonisotopic techniques have advantages in speed, cost, and safety, and are approaching isotopic techniques in sensitivity.

Immobilization of the target on a surface is the most powerful method for avoidance of the competitive equilibrium problem. Because the target and its original partner strand are immobilized, the original duplex cannot reform. The target is immobilized on a [nitrocellulose](#) or nylon filter, although there are many other solid-phase formats. In most techniques, the probe is delivered to the target; however, in [DNA chip](#) technology the target is delivered to a microarray of immobilized oligonucleotides.

[Blotting](#) techniques provide a means for quantifying and identifying specific target nucleic acids. In dot blots or slot blots, the target is deposited on a small area of a **blotting matrix**. A series of targets, each deposited in a discrete spot on the filter, can be hybridized simultaneously to a single probe. In **Southern** and **Northern blots** and related techniques, target nucleic acids are distributed into discrete bands on a gel by [electrophoresis](#). The target nucleic acids are then transferred to a filter membrane, on which a conventional hybridization procedure is performed. The extent of hybridization is detected and quantified for each band on the original gel.

When the spatial distribution of target is sought, [in situ hybridization](#) techniques are employed. Tissue, cell, or [chromosome](#) samples are prepared and fixed. The nucleic acids are denatured. The hybridization reaction is performed by placing a probe-containing solution directly onto the sample. Visualization involves one of several [microscopy](#) techniques.

Bibliography

1. C. Kessler (1992) in *Nonisotopic DNA Probe Techniques*, L. Kricka, ed., Academic Press, San Diego. pp. 29–92.

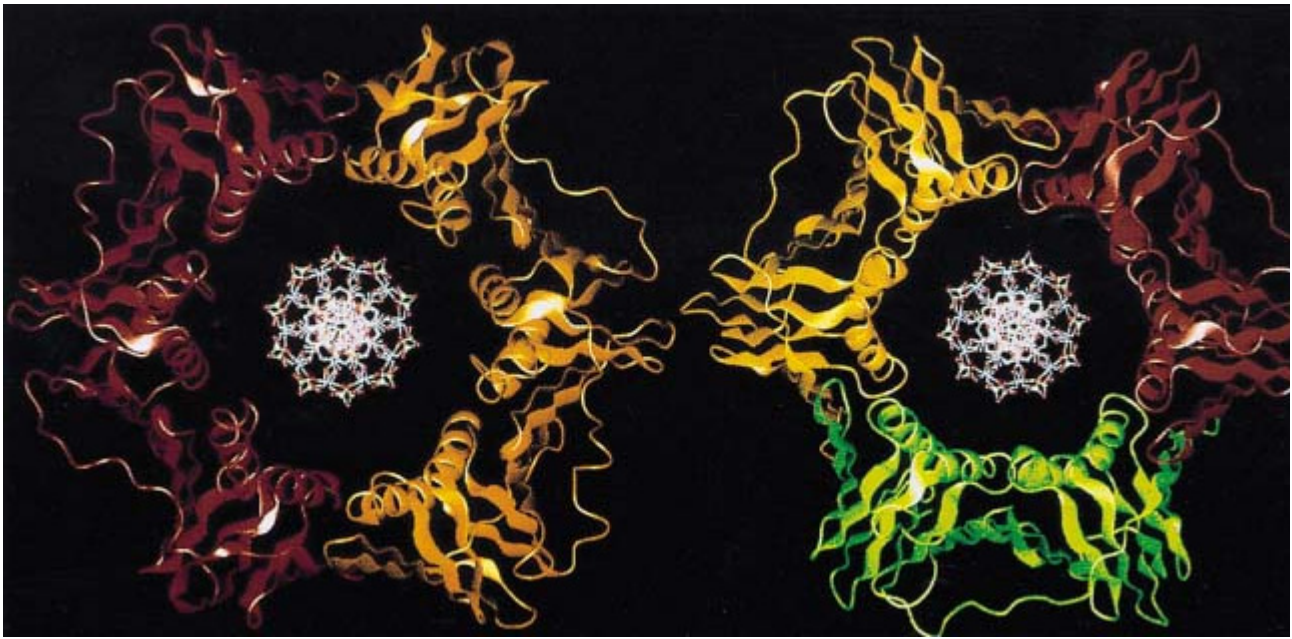
### Suggestions for Further Reading

2. P. Tijssen (1993) *Hybridization with Nucleic Acid Probes: Part 1, Theory and Nucleic Acid Preparation; Part 2, Probe Labeling and Hybridization Techniques*, Elsevier, Amsterdam.
3. L. Kricka, ed. (1992) *Nonisotopic DNA Probe Techniques*, Academic Press, San Diego.
4. D. C. Darling and P. M. Brickell (1994) *Nucleic Acid Blotting: The Basics*, IRL Press, Oxford.
5. G. H. Keller and M. M. Manek (1993) *DNA Probes*, 2nd ed., Stockton Press, New York.

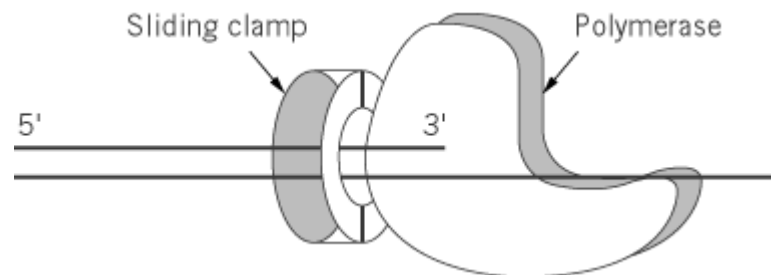
### DNA Polymerase Sliding Clamps

Chromosomal replicases are multiprotein assemblies composed of a DNA polymerase and several accessory proteins and are characterized by remarkable speed and processivity. The term “processivity” refers to the fact that the polymerase is capable of polymerizing thousand of nucleotides without dissociating from the DNA. In several well-characterized systems (bacteria, archaea, eukarya, and bacteriophage-T4), the replicases have been shown to be similar in function, structure, and overall organization (see Clamp Loaders of DNA Polymerases), and in all these systems the remarkable processivity is achieved, in part, by the DNA polymerase accessory proteins. The polymerase accessory proteins can be divided into two groups: a complex of proteins known as the “clamp loader,” and a processivity factor called the “sliding clamp” (1). The sliding clamp is a ring-shaped protein (Fig. 1) that encircles DNA and anchors the polymerase to the DNA template (Fig. 2). The sliding clamp cannot assemble itself around DNA but is loaded onto DNA by the clamp loader (see Clamp Loaders of DNA Polymerases). The arrangement of the polymerase and sliding clamp allows for rapid and processive DNA synthesis (Fig. 2). The sliding clamps have also been shown to play a role in DNA repair and gene transcription and to interact with cellular components involved in progression and control of the cell cycle.

**Figure 1.** Computer generated images of *E. coli*,  $\beta$  subunit, and yeast PCNA. Ribbon representation of polypeptide backbones of trimer of yeast PCNA (**right**) and a dimer of the *E. coli*  $\beta$  subunit (**left**). Strands of  $\beta$  sheets are shown as fl ribbons, and a helices are shown as spirals. The subunits within each ring are distinguished by different colors. A hypothetical model of DNA is placed in the center of each ring. Adapted from (8).



**Figure 2.** A model for the action of a sliding clamp as a processivity factor. The sliding clamp encircles the DNA and binds to the polymerase, thus tethering it to the DNA for processive DNA synthesis.



## 1. Structure and Biochemical Properties of the Sliding Clamp

The sliding clamps of bacteria, archaea, eukaryotes and phage-T4 are similar in both structure and function (2-4). In *Escherichia coli*, the sliding clamp is the  $\beta$  subunit; in archaea and eukaryotes, it is the product of the proliferating cell nuclear antigen (PCNA) gene, and in phage-T4, it is the product of gene-45 (gp45). The  $\beta$  subunit is active as a homo-dimer, whereas PCNA and gp45 are active as homo-trimers (Fig. 1) (3, 5). The molecular mass of the dimer of the  $\beta$  subunit is 81 kDa (40.5 kDa as monomer), the PCNA trimer is 87 kDa (27 kDa as monomer), and gp45 trimer is 74 kDa (24.5 kDa as monomer) (6). Thus, the molecular mass of these oligomers are similar, and it would be shown that the  $\beta$  subunit, PCNA, and gp45 have similar structures.

### 1.1. The Three-Dimensional Structure

The processivity factors are ring-shaped proteins that encircle DNA. The overall structure of these sliding clamps are very similar; the PCNA and the  $\beta$  subunit rings are superimposable (Fig. 1) (5, 7-12). Each ring has similar dimensions and a central cavity large enough to accommodate a double-stranded DNA (dsDNA) molecule (Fig. 1) ((7-12); reviewed in: (3-5)). Although PCNA and gp45 are trimers and the  $\beta$  protein is a dimer, they each have sixfold symmetry. The symmetry derives from the two (PCNA and gp45) or three (the  $\beta$  subunit) globular domains, each within each monomer that have the same chain fold. Although the domains in each monomer do not have

significant sequence similarity, they are superimposable in three dimensions (7). Each domain contains two  $\alpha$ -helices that line the central cavity and are perpendicular to the DNA. The  $\alpha$ -helices are supported by a continuous layer of  $\beta$ -sheet structures all around the outside that also form the intermolecular boundaries (7-12).

The sliding clamps are quite acidic proteins and would be repelled by DNA. The charge distribution on the ring, however, is asymmetric. The outer surface has strong negative electrostatic potential, which may prevent the clamp from nonspecific interactions with DNA. On the other hand, the inside of the central cavity, where the DNA is located, has a net positive electrostatic potential and, thus, may strengthen the interaction with DNA upon proper assembly of the ring around by the clamp loader (see Clamp Loaders of DNA Polymerases).

Although the three prototypic clamps have a similar structure, they differ in their relative stability, in solution and on DNA. The  $\beta$  subunit forms a stable dimer in solution and is also stable on DNA while PCNA is not as stable. On the other hand, gp45 is much less stable in solution and on DNA (6).

## 2. The Biochemical Properties

The first indication for the toroidal shape of the sliding clamps came from the study of the  $\beta$  subunit of the *E. coli* replicase. The  $\beta$  subunit was shown to bind tightly to circular DNA but dissociate readily upon linearization of the plasmid by sliding off the ends (6, 13). This dissociation of  $\beta$  from linear DNA could be prevented if the DNA ends were blocked (6, 13). Because the affinity of  $\beta$  to DNA depended on the geometry of the DNA molecule, these observations suggested a toroidal shaped structure for the sliding clamp and a topological mode of binding.

Evidence for the topological binding of PCNA to DNA came from the observation that PCNA alone can support processive replication by the polymerase in the absence of clamp loader (see Clamp loaders of DNA Polymerases) only when the DNA is linear and has a double-stranded end (14). These results demonstrated that PCNA could thread itself onto the end of the dsDNA and slide along the duplex until it reaches the 3' terminus where it interacts with the polymerase to initiate processive DNA synthesis (Fig. 2). Photocrosslinking experiments also demonstrated the sliding property of PCNA. PCNA can be crosslinked to DNA after its assembly around circular DNA. This crosslinking did not occur if the DNA was linearized, suggesting that the PCNA had dissociated from the DNA via sliding over the ends (15).

Cryoelectron microscopy was used to observe gp45, which appeared to encircle DNA; the results also indicated that it might slide along the duplex (16). The ability of gp45 to slide along DNA was confirmed using photochemical cross-linking analysis. gp45 could be cross-linked to circular, but not linear, DNA, suggesting that the protein could slide over linear DNA (15).

## 3. Functions of the Processivity Factors

### 3.1. Processivity Factors and DNA Replication

To date, the best understood function for the sliding clamps is their essential role in chromosomal DNA replication. The processivity factors endow its cognate polymerase with the ability to polymerize thousands of nucleotides without dissociating from the DNA template. The  $\beta$  subunit is the processivity factor of DNA polymerase (pol) III holoenzyme; PCNA is a part of the pol  $\delta$  holoenzyme, and gp45 binds the product of gene-43 (phage-T4 polymerase).

In bacteria the processivity factor is essential for DNA replication. *dnaN* is the gene that encodes the *E. coli*  $\beta$  subunit (17). *E. coli* cells harboring a temperature-sensitive mutation in the *dnaN* gene stop DNA synthesis when shifted to the restrictive temperature (18). This observation, together with early biochemical analysis (19), indicated a role for the  $\beta$  subunit in DNA replication. Further studies demonstrated that the  $\beta$  subunit is responsible for the processivity of both DNA polIII and polIII

(Table 1). The polymerase of the *E. coli* polIII holoenzyme incorporates approximately 20 nucleotides/s and is processive for only 11 nucleotides before dissociating from the DNA template. Upon binding of the  $\beta$  subunit, polIII becomes a very rapid (750 nucleotides/s) and highly processive (>50,000 nucleotides) enzyme (20). Upon the assembly of the  $\beta$  subunit around the DNA template by the accessory complex (see Clamp Loaders of DNA Polymerases), the clamp loader can be removed, leaving the clamp around the DNA. Then, the  $\beta$  ring can associate with the polymerase (Fig 2) to initiate rapid and processive DNA synthesis (21). Thus, at least in prokaryotes, the formation of a processive holoenzyme can be seen as a two step process. In the first step, a  $\beta$  clamp is assembled around the DNA primer by the clamp loader to form the “preinitiation complex.” In the second step, the polymerase assembles with the preinitiation complex to form the “initiation complex” capable of rapid and processive DNA synthesis (see Clamp Loaders of DNA Polymerases). Therefore, the mechanism by which a  $\beta$  ring confers processivity to polymerase is by direct contact with it, tethering the polymerase to the DNA (Fig. 2). The polymerase would then pull  $\beta$  along during polymerization.

**Table 1. Partial List of Sliding Clamp Interacting Proteins**

| <b>Clamp</b>       | <b>Interacts With:</b>            | <b>Role of Interaction</b>          |
|--------------------|-----------------------------------|-------------------------------------|
| <b>PCNA</b>        | DNA polymerase $\delta$           | Processive DNA synthesis            |
|                    | DNA polymerase $\epsilon$         | Processive DNA synthesis            |
|                    | RFC                               | Loading/unloading of PCNA           |
|                    | Fen-1                             | Stimulation of Fen-1                |
|                    | DNA ligase I                      | Okazaki fragment maturation         |
|                    | D-type cyclins                    | Inhibits DNA replication            |
|                    | p21                               | Inhibits processive DNA replication |
|                    | Gadd45                            | DNA repair                          |
|                    | MSH                               | DNA repair                          |
|                    | <b><math>\beta</math> subunit</b> | DNA polymerase II                   |
| DNA polymerase III |                                   | Processive DNA synthesis            |
| g-complex          |                                   | Loading/unloading of $\beta$        |
| DNA ligase         |                                   | Okazaki fragment maturation         |
| <b>gp45</b>        |                                   | MutS                                |
|                    | gp43                              | Processive DNA synthesis            |



*E. coli* RNA polymerase Activation of transcription

---

In eukaryotes, the early indication of a role for PCNA in DNA synthesis came from the pattern of expression during the cell cycle. PCNA is a nuclear protein with elevated expression during the S-phase (DNA synthesis) of the cell cycle (22, 23) whereas in quiescent and senescent cells there are very low levels of the protein.

Several *in vivo* and *in vitro* studies determined the important role played by PCNA in DNA replication. Expression of anti-sense RNA in exponentially growing cells caused a suppression of DNA replication and cell-cycle progression (24). Deletion of the PCNA gene in both the budding yeast (25) and the fission yeast (26) demonstrated that PCNA is an essential protein that is important for DNA replication. *In vitro* replication assays have been employed to establish the role of PCNA as the processivity factor of Pold (Table 1), the replicase of eukaryotic cells. PCNA enables pold to replicate a template with long single-stranded stretches (27) and is required for Simian virus 40 (SV-40) *in vitro* replication (28), while antibodies generated against PCNA inhibit pold-dependent *in vitro* replication (29). PCNA also plays an important role in coordinating leading and lagging strand synthesis (30).

PCNA also stimulates the activity of pol<sup>ε</sup> (Table 1) (31). The precise role of pol<sup>ε</sup> in chromosomal DNA synthesis, however, is not yet clear; it may play a role in Okazaki fragment maturation. Pol<sup>ε</sup> was shown to be important for DNA repair (23, 32), and its stimulation by PCNA might be important in this process.

The role of PCNA in chromosomal DNA synthesis might not be limited to its effect on the polymerase. PCNA was shown to stimulate the activity of Fen-1 and DNA ligase I (Table 1), proteins that are important for Okazaki fragment maturation on the lagging strand (33). Thus, not only does the processivity factor play an important role in DNA synthesis but also in the final steps leading to the formation of a mature duplex DNA. PCNA may also interact with other members of the replication machinery, although these interactions have not yet been elucidated.

Early studies using temperature-sensitive mutants of gp45 (and gp44/62, its clamp loader) demonstrated its importance for DNA replication (34). Further studies using purified protein in *in vitro* replication assays have been used to establish the role of gp45 as the processivity factor for phage-T4 polymerase (gp43) (35).

### 3.2. Processivity Factors and DNA Repair

Both  $\beta$  and PCNA are also involved in DNA repair, at least in part, via the stimulation of their respective polymerases (32, 36). As discussed above, the  $\beta$  subunit forms a dimeric ring that functions as the processivity factor of polII and PolIII, both of which are involved in DNA repair processes (36, 37). It was also shown that the protein interacts with mutS, a protein needed for mismatch repair (36-38). It was also demonstrated that exposure of *E. coli* cells to UV-irradiation induces a shorter form of  $\beta$ , called  $\beta^*$  (39).  $\beta^*$ , comprised of the C-terminal 2/3 of normal  $\beta$ , was shown to behave as a trimer, presumably forming a trimeric ring, and was found to stimulate DNA synthesis by polIII. Thus,  $\beta^*$  might serve as an alternative clamp in DNA repair processes. In eukaryotes, PCNA is the processivity factor of pol<sup>ε</sup>, an enzyme that has been implicated in DNA repair (23, 32). PCNA becomes localized in the cell nucleus after UV-irradiation of cells not in S-phase (40); PCNA transcription is stimulated upon UV-irradiation (41), and several mutant forms of

PCNA have been shown to be defective in DNA repair processes (42). These *in vivo* observations are consistent with *in vitro* repair assays that demonstrated a role for PCNA in DNA excision-repair (43, 44 reviewed in: 32, 36).

PCNA may also be involved in DNA repair in mechanisms that do not involve the polymerase (23). Upon DNA damage, there is an elevation in the level of several genes. Two of these genes, p21 and Gadd45, have been shown to bind PCNA (Table 1) (33). The role of the interactions between PCNA and these proteins in DNA repair is not yet fully understood. It was shown, however, that p21-PCNA interaction inhibits DNA replication, thus linking cell-cycle control processes to the DNA replication machinery (23).

### 3.3. Transcription and Other Processes

Several well characterized examples suggest a role for viral-encoded sliding clamps in RNA transcription. It was demonstrated that gp45 serves as a transcriptional activator of phage-T4 late genes (45, 46). Two other viral-encoded proteins are also involved in this mechanism: the gene-33 product (gp33) is a transcriptional coactivator, and the gene-55 product (gp55) recognizes the late promoter. The gp45 protein, in association with gp33 and gp55, binds to the *E. coli* RNA polymerase and directs it to the promoters of the phage genes expressed in late stages of infection (45-47) (Table 1). The clamp loader is also involved in this process as it was demonstrated that gp45 has to be assembled around DNA in order to activate transcription (46).

At least one eukaryotic virus has been shown to utilize a sliding clamp for the regulation of gene transcription. A PCNA homologue from an insect virus has been shown to be important for the expression of several viral genes. This PCNA, however, is not essential for viral DNA replication (48).

These two examples imply that cellular-encoded sliding clamps (b and PCNA) may have a similar function and may play a role in gene transcription. Upon completion of DNA replication, some rings may be left around the duplex DNA (49). These rings perhaps then bind to RNA polymerase, directly or indirectly, and regulate gene expression.

PCNA may also play a role in chromatin remodeling. Studies conducted in starfish and *Drosophila* have indicated a function for PCNA in the assembly of chromatin (50, 51).

## 4. Concluding Remarks

The DNA polymerase sliding clamp encircles DNA and endows the polymerase with its high processivity by tethering the polymerase to the DNA. In the last several years, however, it has become apparent that the sliding clamps play diverse roles in nucleic acid metabolism (e.g., DNA repair and transcription). The proteins were also shown to be involved in cell-cycle progression via their interaction with diverse cellular proteins (23). The list of proteins that interact with the clamp is not complete. Future studies will elucidate the scope of processes in which the sliding clamps are involved.

The first proteins found to have a ring-shaped structure and encircle DNA were the DNA polymerase sliding clamps. Since that structure was demonstrated for the sliding clamps, other ring-shaped protein complexes that encircle DNA have been shown to play a role in DNA metabolism (52, 53). It will be interesting to know how general this structure is, how many other proteins have a similar structure, and how many other cellular processes utilize these proteins.

## Bibliography

“DNA Polymerase Sliding Clamps” in , Vol. 3, pp. 1954–1958, by Zvi Kelman, University of Maryland, Biotechnology Institute, Rockville, MD and Lori M. Kelman, Montgomery College, Germantown, MD; “Helix-Turn-Helix Motif” in (online), posting date: January 15, 2002, by by Zvi

Kelman, University of Maryland, Biotechnology Institute, Rockville, MD and Lori M. Kelman, Montgomery College, Germantown, MD.

1. C.-C. Huang, J. E. Hearst, and B. M. Alberts (1981) *J. Biol. Chem.* **256**, 4087–4094.
2. B. Stillman (1994) *Cell* **78**, 725–728.
3. M. M. Hingorani and M. O'Donnell (2000) *Curr. Biol.* **10**, R25–R29.
4. M. M. Hingorani and M. O'Donnell (2000) *Curr. Organic Chem.* **4**, 887–913.
5. Z. Kelman and M. O'Donnell (1995) *Nucleic Acids Res.* **23**, 3613–3620.
6. N. Yao, J. Turner, Z. Kelman, P. T. Stukenberg, F. Dean, D. Shechter, Z.-Q. Pan, J. Hurwitz, and M. O'Donnell (1996) *Genes to Cell* **1**, 101–113.
7. X. -P. Kong, R. Onrust, M. O'Donnell, and J. Kuriyan (1992) *Cell* **69** 425–37.
8. T. S. R. Krishna, X.-P. Kong, P. M. Burgers, and J. Kuriyan (1994) *Cell* **79** 1233–1243.
9. J. M. Gulbis, Z. Kelman, J. Hurwitz, M. O'Donnell, and J. Kuriyan (1996) *Cell* **87**, 297–306.
10. Y. Shamoo and T. A. Steitz (1999) *Cell* **99**, 155–156.
11. I. Moarefi, D. Jeruzalmi, J. Turner, M. O'Donnell, and J. Kuriyan (2000) *J. Mol. Biol.* **296**, 1215–1223.
12. S. Matsumiya, Y. Ishino, and K. Morikawa (2001) *Prot. Sci.* **10**, 17–23.
13. P. T. Stukenberg, P. S. Studwell-Vaughan, and M. O'Donnell (1991) *J. Biol. Chem.* **266**, 11328–11334.
14. P. M. Burgers and B. L. Yoder (1993) *J. Biol. Chem.* **268**, 19923–19926.
15. R. L. Tinker, G. A. Kassavetis, and E. P. Geiduschek (1994) *EMBO J.* **13**, 5330–5337.
16. E. P. Gogol, M. C. Young, W. L. Kubasek, T. C. Jarvis, and P. H. von Hippel (1992) *J. Mol. Biol.* **224**, 395–412.
17. P. M. J. Burgers, A. Kornberg, and Y. Sakakibara (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5391–5395.
18. Y. Sakakibara and T. Mizukami (1980) *Mol. Gen. Genet.* **178**, 541–553.
19. S. H. Wickner (1978) *Ann. Rev. Biochem.* **47**, 1163–1191.
20. M. Mok and K. J. Mariani (1987) *J. Biol. Chem.* **262**, 16644–16654.
21. Z. Kelman and M. O'Donnell (1995) *Ann. Rev. Biochem.* **64**, 171–200.
22. R. Bravo and J. E. Celis (1980) *J. Cell. Biol.* **84**, 795–802.
23. Z. Kelman (1997) *Oncogene* **14**, 629–640.
24. D. Jaskulski, J. K. DeRiel, W. E. Mercer, B. Calabretta, and R. Baserga (1988) *Science* **240**, 297–304.
25. G. A. Bauer and P. M. J. Burgers (1990) *Nucleic Acids Res.* **18**, 261–265.
26. N. H. Waseem, K. Labib, P. Nurse, and D. P. Lane (1992) *EMBO J.* **11**, 5111–5120.
27. C.-K. Tan, C. Castillo, A. G. So, and K. M. Downey (1986) *J. Biol. Chem.* **261**, 12310–12316.
28. G. Prelich, M. Kostura, D. R. Marshak, M. B. Mathews, and B. Stillman (1987) *Nature* **326**, 471–475.
29. C.-K. Tan, K. Sullivan, X. Li, E. M. Tan, K. M. Downey, and A. G. So (1987) *Nucleic Acids Res.* **15**, 9299–9308.
30. G. Prelich and B. Stillman (1988) *Cell* **53**, 117–126.
31. S.-H. Lee, Z.-Q. Pan, A. D. Kwong, P. M. J. Burgers, and J. Hurwitz (1991) *J. Biol. Chem.* **266**, 22707–22717.
32. R. D. Wood (1996) *Ann. Rev. Biochem.* **65**, 135–167.
33. Z. Kelman and J. Hurwitz (1998) *Trends Biochem. Sci.* **23**, 236–238.
34. R. H. Epstein, A. Bolle, C. M. Steinberg, E. Kellenberger, E. Boy de la Tour, R. Chevalley, R. S. Edgar, M. Susman, G. H. Denhardt, and A. Lielausis (1963) *Cold Spring Harbor Symp. Quant. Biol.* **28**, 375–394.

35. N. G. Nossal and B. M. Alberts (1983) *Bacteriophage T4*. (C. K. Mathews, E. M. Kutter, G. Mosig, and P. B. Berget eds.). American Society for Microbiology, Washington DC, pp. 71–81.
36. A. Sancar (1996) *Ann. Rev. Biochem.* **65**, 43–81.
37. Z. Livneh, O. Cohen-Fix, R. Skaltier, and T. Elizur (1993) *CRC Crit. Rev. Biochem. Mol. Biol.* **28**, 465–513.
38. F. J. Lopez de Saro and M. O'Donnell (2001) *Proc. Natl. Acad. Sci.* **98**, 8376–8380.
39. R. Skaliter, T. Paz-Elizur, and Z. Livneh (1996) *J. Biol. Chem.* **271**, 2478–2481.
40. J. E. Celis and P. Madsen (1986) *FEBS Lett.* **209**, 277–283.
41. X.-R. Zeng, Y. Jiang, S.-J. Zhang, H. Hao, and M. Y. W. T. Lee (1994) *J. Biol. Chem.* **269**, 13748–13751.
42. R. Ayyagari, K. J. Impellizzeri, B. L. Yoder, S. L. Gary, and P. M. J. Burgers (1995) *Mol. Cell. Biol.* **15**, 4420–4429.
43. A. F. Nichols and A. Sancar (1992) *Nucleic Acid Res.* **20**, 2441–2446.
44. M. K. K. Shivji, M. K. Kenny, and R. D. Wood (1992) *Cell* **69**, 367–374.
45. D. R. Herendeen, G. A. Kassavetis, J. Barry, B. M. Alberts, and E. P. Geiduschek (1989) *Science* **245**, 952–58.
46. E. P. Geiduschek (1995) *Semi. Virol.* **6**, 25–33.
47. E. N. Brody, G. A. Kassavetis, M. Ouhammouch, G. M. Sanders, R. L. Tinker, and E. P. Geiduschek (1995) *FEMS Microb. Lett.* **128**, 1–8.
48. D. R. O'Reilly, A.M. Crawford, and L. K. Miller (1989) *Nature* **337**, 606.
49. A. Yuzhakov, J. Turner, and M. O'Donnell (1996) *Cell* **86**, 877–886.
50. A. Nomura (1994) *J. Cell Sci.* **107**, 3291–3300.
51. D. S. Henderson, S. S. Banga, T. A. Grigliatti, and J. B. Boyd (1994) *EMBO J.* **13**, 1450–1459.
52. M. M. Hingorani and M. O'Donnell (2000) *Nat. Rev. Mol. Cell Biol.* **1**, 22–30.
53. M. M. Hingorani and M. O'Donnell (1998) *Curr. Biol.* **8**, R83–R86.

### **Suggestions for Further Reading**

54. M. M. Hingorani and M. O'Donnell (1998) Toroidal proteins: running rings around DNA. *Curr. Biol.* **8**, R83–R86.
55. M. M. Hingorani and M. O'Donnell (2000) Sliding clamps: A (tail)ored fit. *Curr. Biol.* **10**, R25–R29.
56. Z. Kelman (1997) PCNA: Structure, functions and interactions. *Oncogene* **14**, 629–640.
57. Z. Kelman and J. Hurwitz (1998) Protein-PCNA interactions: A DNA-scanning mechanism? *Trends Biochem. Sci.* **23**, 236–238.
58. Z. Kelman and N. O'Donnell (1995) Structural and functional similarities of prokaryotic and eukaryotic sliding clamps. *Nucleic Acids Res.* **23**, 3613–3620.

### **Prochiral**

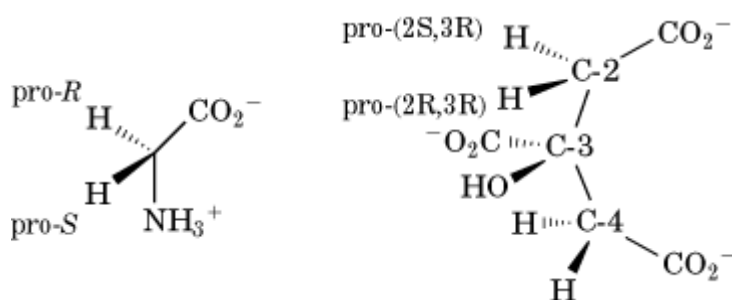
Atoms, substituents, or whole molecules can be described as prochiral. When describing an atom, prochiral indicates that the substitution of that atom with a heavier isotope will alter the **chirality** of the molecule. The phrase prochiral center is employed to describe an atom with two identical substituents such that the substitution of either one with a different isotope would make that atom a

chiral center (see **Chiral**). If the isotopic substitution of the substituents would generate **enantiomers**, the groups are enantiotopic, whereas if the substitution would produce **diastereomers**, the groups are diastereotopic (1, 2). The concept of enantiotopic and diastereotopic groups is important in biochemistry and molecular biology because these groups interact differently with chiral molecules. The classic example of the difficulties presented by prochiral chemistry was interpreting the isotopic labeling of citrate during its biosynthesis (3).

As an example of prochirality, the  $\alpha$ -carbon in [glycine](#) is a prochiral center because the substitution of either of the two identical substituents (for glycine, the H<sup>a</sup> atoms) with a heavier isotope would generate a chiral molecule. The two H<sup>a</sup> atoms are enantiotopic because their replacement would generate different enantiomers of glycine. To differentiate the enantiotopic groups, they are preceded with the prefix “pro-” preceding the specification of the configuration generated. In chemical structures, the prochirality of an atom can be identified by the subscript *R* or *S* (1), which indicates the [configuration](#) of the chiral center as determined by the Cahn, Ingold, and Prelog rules (4) (see [Configuration](#)). In the peptide glycyl-L-alanine, the glycine H<sup>a</sup> atoms are diastereotopic because the substitution of either glycine H would generate diastereomers.

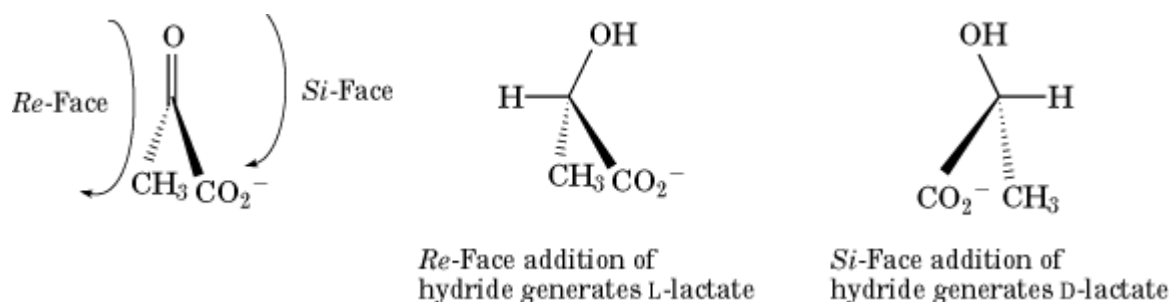
In citrate, the carboxy methyl groups are enantiotopic and C-3 is a prochiral center, because the substitution of any of the atoms in either carboxymethyl group will generate chirality at C-3 (Fig. 1). The upper carboxymethyl group in Figure 1 is the pro-*R* group because isotopic substitution anywhere in the functional group would generate (3*R*)-citrate. C2 and C4 are also prochiral centers, because the substitution of the methylene hydrogen atoms would generate a chiral center at the substituted carbon. Both hydrogen atoms on C-2 of citrate are diastereotopic, because the substitution of deuterium at either of these two positions would generate diastereomers. They are diastereomers because two new chiral centers (C2 and C3) would be generated, and the chirality at C3 is the same for both hydrogens, whereas the chirality at C2 would differ.

**Figure 1.** The stereochemistry of glycine and citrate. The enantiotopic hydrogen atoms of glycine are labeled (left). The forward H is pro-*S* because substitution with <sup>2</sup>H at this H would generate (*S*)-[2-<sup>2</sup>H]glycine, as suggested by the counterclockwise arc connecting the atoms in order of priority. In citrate (right), the hydrogen atoms at C2 are diastereotopic; C2 and C3 are prochiral centers because the substitution of <sup>2</sup>H at either of these hydrogen atoms would generate the labeled diastereomers.



Planar double-bonded carbon and nitrogen atoms have the potential to be prochiral if they have three different substituents, because the two faces of the molecule are nonequivalent (2). The faces of the molecule are differentiated by the designation of *Re* or *Si* based on the Cahn, Ingold, and Prelog priority (4) of the three substituents. As shown in Figure 2, the two faces of pyruvate are enantiotopic because addition to opposite faces would generate enantiomers. If the addition to opposite faces would generate diastereomers, the faces are diastereotopic.

**Figure 2.** The stereochemistry of pyruvate. The two faces of pyruvate are different. When viewed from the left, the circle connecting the three different substituents in their Cahn/Ingold/Prelog priority order describes a clockwise motion, whereas when viewed from the right, the circle describes a counterclockwise motion. The faces are enantiotopic because the addition of an identical group to each face would generate enantiomers.



### Bibliography

1. D. Arigoni, and E. L. Eliel (1969) In *Stereochemistry* (E. L. Eliel and N. L. Allinger, eds.), Vol. 4, Wiley-Interscience, New York, pp. 127–243.
2. B. Testa (1982) In *Stereochemistry* (C. Tamm, ed.), New Comprehensive Biochemistry, Vol. 3, Elsevier, Amsterdam, pp. 1–48.
3. A. G. Ogston (1978) *Nature* **276**, 676.
4. R. S. Cahn, C. K. Ingold, and V. Prelog (1966) *Angew. Chem. Int. Ed.* **5**, 385–415.

### Suggestions for Further Reading

5. K. R. Hanson (1975) Reactions at prochiral centers. Interdependence in the estimation of enzyme stereospecificity toward prochiral centers and the configurational purity of labeled substrates. *J. Biol. Chem.* **250**, 8309–8314.
6. K. R. Hanson and I. A. Rose (1975) Interpretations of enzyme reaction stereospecificity. *Acc. Chem. Res.* **8**, 1–10.
7. H. Hirschman, and K. R. Hanson (1971) Elements of stereoisomerism and prostereoisomerism. *J. Org. Chem.* **36**, 3293–3306.
8. V. Prelog and G. Helmchen (1972) Pseudoasymmetry in organic chemistry. *Helv. Chim. Acta* **55**, 2581–2598.

### Product Inhibition

It would be expected that the addition of a reaction product will inhibit the rate of the forward reaction catalyzed by an [enzyme](#). As a substrate for the reaction in the reverse direction, it must give rise to the same enzyme forms as that from which it was released in the forward reaction. In practice, however, it may be difficult to demonstrate the inhibition if the inhibition constant for the product is relatively high.

Product inhibition occurs either by partly reversing the reaction or by competing with a substrate for a particular form of the enzyme. The inhibition patterns produced by a product can be very

informative in the case of enzyme reactions that produce multiple products. Rules for the qualitative prediction of product inhibition patterns have been developed by Cleland (1, 2); they can be summarized as:

1. A product inhibitor affects the *slope* of a double-reciprocal plot (see [Lineweaver–Burk Plot](#)) when
  1. the product inhibitor and the variable substrate combine reversibly with the same enzyme form, or
  2. the product inhibitor and the variable substrate combine reversibly with different enzyme forms that are connected by a series of reversible steps along which interaction can occur. The release of product at zero concentration or the addition of substrate at a saturating concentration are considered to be irreversible steps that break reversible connections.
2. A product inhibitor affects the intercept of a double-reciprocal plot when it combines reversibly with an enzyme form other than the one with which the variable substrate combines and saturation with the variable substrate does not overcome the inhibition.

Application of these rules permits analysis of product inhibition patterns that make it possible to distinguish between different **kinetic mechanisms**. Product inhibition patterns for Bi–Bi reaction mechanisms, with two substrates and two products, are illustrated in Table 1.

**Table 1. Product Inhibition Patterns for Bi–Bi Reaction Mechanisms**

| Mechanism <sup>a</sup>             | Product Inhibitor | Variable Substrate |                 |
|------------------------------------|-------------------|--------------------|-----------------|
|                                    |                   | A                  | B               |
| Ordered                            | P                 | NC <sup>b</sup>    | NC <sup>c</sup> |
| Equilibrium, ordered               | Q                 | C                  | NC              |
|                                    | P                 | NC                 | NC              |
| RER <sup>d</sup> plus dead-end EBQ | Q                 | C                  | C               |
|                                    | P                 | C                  | C               |
| RER plus dead-end EBQ and EAP      | Q                 | C                  | NC              |
|                                    | P                 | NC                 | C               |
| Ping-Pong                          | Q                 | C                  | NC              |
|                                    | P                 | NC                 | C               |
|                                    | Q                 | C                  | NC              |

<sup>a</sup> Mechanisms are illustrated under **Kinetic mechanisms**.

<sup>b</sup> C, linear competitive; NC, linear noncompetitive; UC, linear uncompetitive.

<sup>c</sup> Inhibition becomes UC if B is present at a saturating concentration.

<sup>d</sup> RER, rapid-equilibrium, random.

Product inhibition patterns can also be predicted by using the complete rate equation for a particular mechanism, setting to zero the concentration of the product that is not to be added, and then rearranging the equation in double-reciprocal form, with each of the two substrates as the variable one (2, 3). The resulting equations are also of importance with respect to the quantitative analysis of product inhibition data. It is only when the quantitative relationship between true and apparent constants is in hand that it becomes possible to determine true values for the kinetic constants associated with a particular mechanism.

### Bibliography

1. W. W. Cleland (1963) *Biochim. Biophys. Acta* **67**, 188–196.
2. W. W. Cleland (1963) *Biochim. Biophys. Acta* **67**, 104–137.
3. W. W. Cleland (1986) in *Investigations of Rates and Mechanisms*, vol. **6** (C. F. Bernasconi, ed.), John Wiley & Sons, New York, pp. 791–870.

### Proenzyme

Proenzymes are [enzymes](#) that are synthesized in cells as inactive precursors that must undergo modification before they are catalytically active. They often function as [proteinases](#), when they are also known as **zymogens**, and they are produced in this latent form because of their destructive potential. An additional protective measure that applies to some of the activated proenzymes is their inactivation through the formation of tight complexes with specific proteins (see [Proteinase Inhibitors](#)). The activation of proenzymes often occurs by **proteolysis**, a process in which peptide bonds are irreversibly cleaved and peptides may be liberated from the molecule (1). Proteolysis may occur autocatalytically or by the action of another proteinase or as a result of the combination of the two activities. Most of the proenzymes are proteinases responsible for digestion, but [thrombin](#) plays an important role in [blood clotting](#). A list of proenzymes and the reactions involved in their conversion to active enzyme forms are given in Table 1. (See top of next page)

**Table 1. Activation of Proenzymes (see Proenzyme article, previous page)**

| Proenzyme           | Active form                      | Chemistry of activation   |
|---------------------|----------------------------------|---|
| Chymotrypsinogen    | <b>Chymotrypsin</b>              | Cleavage of Arg-Ile bond by <b>trypsin</b> and autocatalysis                                    |
| Pepsinogen          | <b>Pepsin</b>                    | Cleavage of Leu-Ile bond, removal of residues from N-terminus at lower pH, and by autocatalysis |
| Procarboxypeptidase | <a href="#">Carboxypeptidase</a> | Removal of peptide from N-terminus by trypsin   |
| Proelastase         | <a href="#">Elastase</a>         | Removal of residues from N-terminus by trypsin  |
| Prothrombin         |                                  | Removal of peptide from N-terminus  |



|             |                                      |  |
|-------------|--------------------------------------|--|
|             | <a href="#">Thrombin<sup>a</sup></a> | by cleavage of Arg-Thr bond by serine endopeptidase and then cleavage of Arg-Ile bond by same enzyme |
| Trypsinogen | <b>Trypsin<sup>b</sup></b>           | Cleavage of Lys-Ile bond by <a href="#">enterokinase</a> and autocatalysis <sup>a</sup>              |

---

<sup>a</sup> Inactivated by interaction with antithrombin III.

<sup>b</sup> Inactivated by interaction with [trypsin inhibitors](#).

## Bibliography

1. R. M. Stroud, A. A. Kossiakoff, and J. L. Chambers (1977) *Ann. Rev. Biophys. Bioeng.* **6**, 177–193.

## Profilin

Profilin (mw 12 000–15 000) is an abundant actin monomer-binding protein in eukaryotic cells (see [Actin-Binding Proteins](#)) ((1); review). Particularly high concentrations of profilin are found in lymphoid cells and brain cells. In lymphoid tissues the concentration of unpolymerized actin is about 100  $\mu$ M and in platelets it may be as high as 200  $\mu$ M. In both cases about 50% of the unpolymerized actin can be accounted for by profilin:b-actin and profilin:g-actin isoforms. The remaining unpolymerized actin appears to be sequestered by b-thymosins (see below). In tissue cultured cells, profilin is concentrated in regions of high motile activity (2), where it is thought to bring actin monomers to growing (+)-ends of filaments exposed through transmembrane signalling events. It binds to a number of cell cortex-associated proteins participating in the control of actin filament formation, eg the VASP family of proteins (3), the WASP-binding protein WIP (4), the formins (5-7) and the Arp2 protein (8).

*Sacharomyces cerevisiae* has one essential profilin gene (9), *Dictyostelium discoideum* has two (10), and plants contain multiple isoforms of profilin (10). *Acanthamoeba castellani* also contains two isoforms of profilin (11), as do mammalian cells (12).

Genetic studies suggest that profilin plays a key role in the establishment of actin-dependent cell polarity in *Saccharomyces cerevisiae*, in the formation of the fruiting body in *Dictyostelium discoideum*, and during early development in *Drosophila* (10, 13). In transgenic mice, homozygous profilin “knock-outs” are lethal emphasizing the physiological importance of profilin in mammalian cells (14).

The structure and ligand-binding by different isoforms of profilins have been elucidated by mutagenesis (15), crystallography (16-19) and NMR (20). All isoforms of profilin have similar folds even though the amino acid sequence is highly divergent. In *Arabidopsis* profilin, the residues involved in the actin interaction have diverged, but the general character of involved residues are conserved. In the profilin:b-actin crystal, profilin makes two contacts with actin. The largest contact involves 2260  $\text{Å}^2$  of buried surface area, and is the contact between profilin and actin in solution. It engages subdomains 1 and 3 of actin at the (+)-end of the actin monomer (16). This leaves the (-)-

end of the monomer free to bind to the fast growing (+)-end of an actin filament in the polymerization process. The less extensive profilin:actin interaction site (1187 Å<sup>2</sup> b.s.a.) seen in the crystal is formed by the N-terminal helix of profilin binding to subdomain 4 of a different actin monomer (16). So far, this contact has not been observed in solution.

*In vitro*, profilin inhibits the polymerization of actin with varying efficiency depending on the ionic conditions, and on whether the homologous nonmuscle actin or heterologous muscle  $\alpha$ -actin is used (21-23). In the presence of physiological concentrations of Mg<sup>2+</sup> ions, profilin is a relatively poor actin sequesterer. Apparently, the conformation of the (–)-end is not disturbed by profilin, and the profilin:actin complex can add onto the (+)-end of filaments. Thus both actin and the profilin:actin complex can support growth of actin filaments at the fast growing (+)-end (22, 24). The characteristics of the association/dissociation of actin monomers at the (–)-end of actin filaments are such that profilin effectively prevents the addition of actin at that end.

Beta thymosins, which inhibit growth at both ends and occur at high concentrations in cells, are thought to be the primary actin-sequestering proteins, and profilin, with its higher affinity for actin, serves to deliver actin to the barbed end (24). In the presence of proteins that cap the (+)-end of actin filaments, profilin is an efficient actin sequestering protein, since it effectively prevents polymerization at the (–)-end under all conditions (25, 26). This implies that the control of actin polymerization in cells may depend primarily on regulating the association/dissociation of proteins that cap the (+)-end of preexisting filaments, or by regulating the nucleation of new filaments.

Profilin increases the rate of **nucleotide** exchange on actin more than 1 000-fold (27, 28), and stabilizes the actin in the nucleotide-free state. Since actin with ATP bound polymerizes faster than ADP-actin, profilin might promote repolymerization of actin. ADP monomers coming off of filaments during depolymerization, by rapidly converting them to ATP-actin monomers.

In a search for **proline** hydroxylase, Tanaka and Shibata found that affinity columns of poly(L-proline) (PLP), in addition to the hydroxylase, also bound profilin and actin. This allowed the development of a simple procedure to obtain a mixture of profilin and the two isoforms of profilin:nonmuscle actin (29, 30), from which the profilin:b-actin and profilin:g-actin isoforms can be isolated (31). The poly(L-proline) binding site on profilin has been identified. It is comprised of highly invariant hydrophobic and aromatic amino acids involving the N- and the C-terminal helices of profilin (15). Mutagenesis, modelling and direct structure determination has shown a bonding pattern in the profilin:polyproline interaction that is analogous to that of SH3 domains in signal transduction proteins (15, 18, 19). The VASP family of proteins, the WIP protein, and the formins all contain proline rich sequences with capacity to bind profilin. Profilin has also been found to bind to Arp2/Arp3 complex which is a multiprotein complex implicated in the mechanism of actin polymerization (see below).

*In vitro*, profilin alone or in complex with nonmuscle actins binds the phospholipid PtdIns 4,5-bisphosphate in micellar form, as well as in vesicles together with other phospholipids. Profilin:b/g-actin complexes also bind to PtdIns 4,5-bisphosphate *in vitro*, something which results in dissociation of the complexes with release of polymerization-competent actin. This would seem to indicate that the appearance of PtdIns 4,5-bisphosphate in the plasma membrane could recruit profilin:actin for polymerization in the advancing cell edge. However, that this is an oversimplification is indicated by the observations that the hydrolysis of PtdIns 4,5-bisphosphate by phospholipase Cg-1 is inhibited when profilin is associated with the PtdIns 4,5-bisphosphate substrate, and that phosphorylation of the enzyme with an activated tyrosine kinase receptor overrides this inhibition. This points at profilin also playing a role in the control of the metabolism of the polyphosphoinositides. Furthermore, profilin binds to the regulatory subunit of PtdIns 4,5-bisphosphate kinase and stimulates the formation of PtdIns 3,4,5-trisphosphate ((1)-for refs).

These observations point at profilin being an important factor also in the control of the release of

inositol trisphosphate and thus of the generation of  $\text{Ca}^{2+}$  pulses in the cell. The primary targets for  $\text{Ca}^{2+}$  regulation in cells is the actomyosin system with global changes in structure and activity as the result. It is now known that many of the microfilament-associated proteins interact with components formed in the phosphatidylinositol-cycle as a result of receptor-mediated activation of the cell, and that these interactions modulate the activity of these proteins vis-à-vis actin. Although the exact physiological roles of these interactions remain to be elucidated, it suggests that the phosphatidylinositol-cycle is directly involved in controlling the microfilament-based motility cycle.

## Bibliography

1. K. Schlüter, B. M. Jockusch, and M. Rothkegel (1997) *Biochim. Biophys. Acta* **1359**, 97–109.
2. F. Buss, C. Temm-Grove, S. Henning, and B. M. Jockusch (1992) *Cell Motil. Cytoskeleton* **22**, 51–61.
3. M. Reinhard, M. Halbrügge, U. Scheer, C. Wiegand, B. M. Jockusch, and U. Walter. (1992) *EMBO. J.* **11**, 2063–2070.
4. N. Ramesh, I. M. Anton, J. J. Hartwig, and R. S. Geha (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14671–14676.
5. M. Evangelista, K. Blundell, M. S. Longtine, C. J. Chow, N. Adames, J. R. Pringle, M. Peter, and C. Boone (1997) *Science* **276**, 118–122.
6. N. Watanabe, P. Madaule, T. Reid, T. Ishizaki, G. Watanabe, A. Kakizuka, Y. Saito, K. Nakao, B. M. Jockusch, and S. Narumiya (1997) *EMBO J.* **16**, 3044–3056.
7. J. Petersen, O. Nielsen, R. Egel, and I. M. Hagan (1998) *J. Cell Sci.* **111**, 867–876.
8. R. D. Mullins, J. F. Kelleher, J. Xu, and T. D. Pollard (1998) *Mol. Biol. Cell* **9**, 841–852.
9. V. Magdolen, U. Oechsner, G. Muller, and W. Bandlow (1988) *Mol. Cell Biol.* **8**, 5108–5115.
10. M. Haugwitz, A. A. Noegel, J. Karakesisoglou, and M. Schleicher (1994) *Cell* **79**, 303–314.
11. S. Huang, J. M. McDowell, M. J. Weise, and R. B. Meagher (1996) *Plant Physiol.* **111**, 115–126.
12. D. A. Kaiser, M. Sato, R. F. Ebert, and T. D. Pollard (1986) *J. Cell Biol.* **102**, 221–226.
13. A. Lambrechts, J. van Damme, M. Goethals, J. Vandekerckhove, and C. Ampe (1995) *Eur. J. Biochem.* **230**, 281–286.
14. T. Fujiwara, K. Tanaka, A. Mino, M. Kikyo, K. Takahashi, K. Shimizu, and Y. Takai (1998) *Mol. Biol. Cell* **9**, 1221–1233.
15. W. Witke, A. H. Sharpe, and D. J. Kwiatkowski (1993) *Mol. Biol. Cell* **4** 149a.
16. C. Björkegren, M. Rozycki, C. E. Schutt, U. Lindberg, and R. Karlsson (1993) *FEBS Lett.* **333**, 123–126.
17. C. E. Schutt, J. C. Myslik, M. D. Rozycki, N. C. Goonesekere, and U. Lindberg (1993) *Nature* **365**, 810–816.
18. E. S. Cedergren-Zeppezauer, N. C. W. Goonesekere, M. D. Rozycki, J. C. Myslik, Z. Dauter, U. Lindberg, and C. E. Schutt (1994) *J. Mol. Biol.* **240**, 459–475.
19. K. S. Thorn, H. E. Christensen, R. Shigeta, D. Huddler, L. Shalaby, U. Lindberg, N. H. Chua, and C. E. Schutt (1997) *Structure* **5**, 19–32.
20. N. M. Mahoney, P. A. Janmey, and S. C. Almo (1997) *Nat. Struct. Biol.* **4**, 953–960.
21. W. J. Metzler, A. J. Bell, E. Ernst, T. B. Lavoie, and L. Mueller (1994) *J. Biol. Chem.* **269**, 4620–4625.
22. H. Larsson and U. Lindberg (1988) *Biochim. Biophys. Acta* **953**, 95–105.
23. E. Korenbaum, P. Nordberg, C. Björkegren-Sjögren, C. E. Schutt, U. Lindberg, and R. Karlsson (1998) *Biochemistry* **37**, 9274–9283.
24. V. K. Vinson, E. M. De La Cruz, H. N. Higgs, and T. D. Pollard (1998) *Biochemistry* **37**, 10871–10880.

25. D. Pantaloni and M.-F. Carlier (1993) *Cell* **75**, 1007–1014.
26. M.-F. Carlier and D. Pantaloni (1993) *Sem. Cell Biol.* **5**, 183–191.
27. F. Markey, H. Larsson, K. Weber, and U. Lindberg (1982) *Biochim. Biophys. Acta* **704**, 43–51.
28. S. C. Mockrin and E. D. Korn (1980) *Biochemistry* **19**, 5359–5362.
29. P. J. Goldschmidt-Clermont, L. M. Machesky, S. K. Doberstein, and T. D. Pollard (1991) *J. Cell Biol.* **113**, 1081–1089.
30. M. Tanaka and H. Shibata (1985) *Eur. J. Biochem.* **151**, 291–297.
31. U. Lindberg, C. E. Schutt, E. Hellsten, A.-C. Tjäder, and T. Hult (1988) *Biochim. Biophys. Acta* **967**, 391–400.
32. M. Segura and U. Lindberg (1984) *J. Biol. Chem.* **259**, 3949–3954.

## Progenote

The current literature uses the term progenote in two different ways: 1) it signifies an organizational level in [evolution](#) when prokaryotic organization preceded cells; or 2) it is used to denote the last common ancestor of all extant life. In some scenarios that describe early cellular evolution, it was assumed that the last common ancestor was at a preprokaryotic level of organization; however, subsequent analyses of the molecular evolution of different cellular components suggest that the last common ancestor was a prokaryote. Based on this realization, the term progenote should be more properly used to denote a hypothetical preprokaryotic stage in cellular evolution, distinct from the last common ancestor.

In 1977 Woese and Fox ([1](#)) defined progenote as a hypothetical stage in the evolution of cells where typical prokaryotic cellular organization preceded organisms:

Eucaryotes did arise from procaryotes, but only in the sense that the procaryotic is an organizational, not a phylogenetic distinction. In analogous fashion procaryotes arose from simpler entities. The latter are properly called progenotes, because they are still in the process of evolving the relationship between genotype and phenotype.

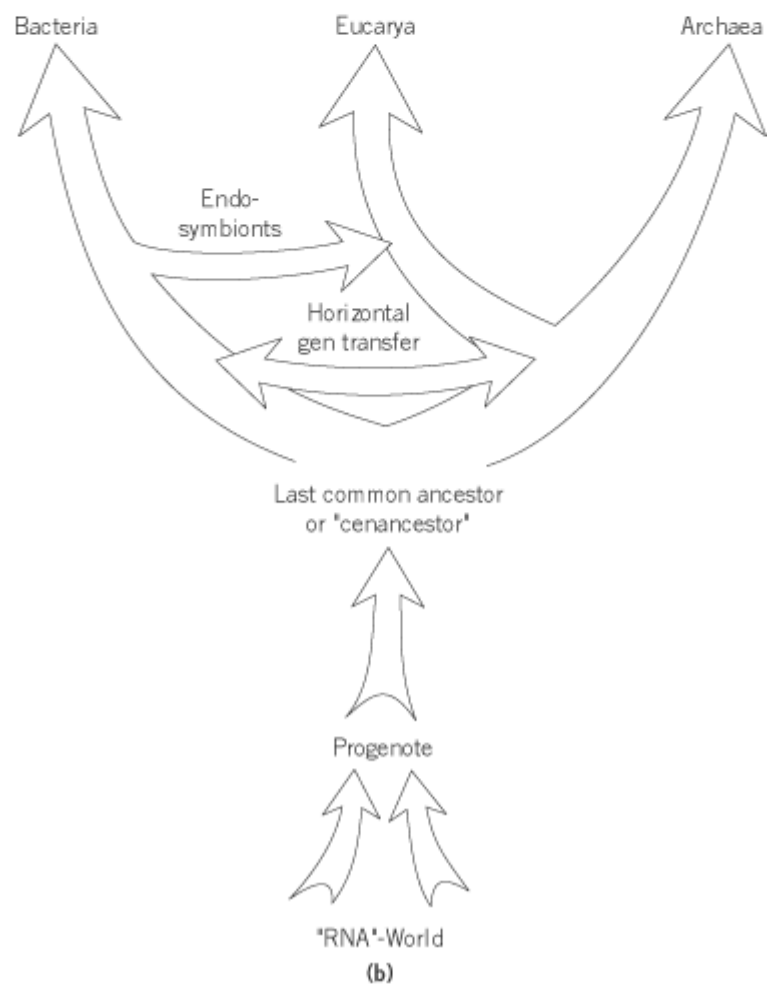
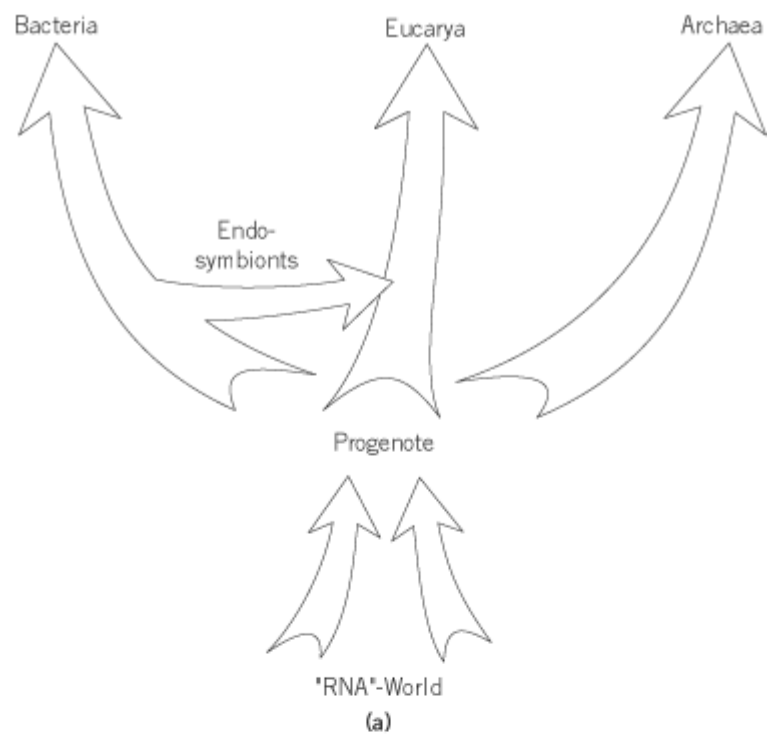
Woese and Fox intended to define an organizational level simpler than, and preceding, the prokaryotic level. At the progenotic level, genes and encoded [proteins](#) were smaller, and the accuracy of [transcription](#) and [translation](#) was lower than at the prokaryotic level. As a result, sequence evolution occurred more rapidly.

At the prokaryotic level, organisms contain a [genome](#) that encodes a multitude of biochemical and structural functions. Among the genome-encoded functions at the prokaryotic level are genome replication, translation of genomic information into functioning molecules, and formation of a semipermeable barrier between the organism and its environment. Prokaryotic organization is so complex that the likelihood is close to zero that a prokaryote spontaneously assembled from activated [nucleotide](#) and [amino acid](#) precursors in a primordial soup (see [Prebiotic Evolution](#)). A solution to this problem is to assume intermediate steps that successively evolve into more complex structures. An intermediate step usually assumed in the origin of life is self-replicating RNA-like molecules (see [RNA World](#)). However, because of the limited accuracy of these early replicators, the size and information contained in these self-replicating molecules is limited ([2](#)).

How a self-propagating network of biochemical reactions that maintains a boundary with the environment, that is, an autopoietic network ([3](#)), might have evolved from simple self-replicating

molecules is a major puzzle in the evolution of life. Most scenarios that describe the evolution of cellular life include progenote-like organisms that are intermediate between the RNA world and the first prokaryotes [e.g., ((4); see Fig. 1b). However, the progenote concept that postulates an organism without strict coupling between **genotype** and **phenotype**, partially negates the major conceptual breakthrough associated with the RNA world, that Darwinian [natural selection](#) already acted on simple self-replicating RNA molecules. Several questions remain open regarding the evolution of progenotes: How does natural selection act on an organism without strict coupling between genotype and phenotype? Can progenotes evolve by natural selection without a phenotype encoded by their genes? Is it feasible that selection took place only (or mainly) at the level of individual molecules, not at the organismal level? Are alternative scenarios reasonable that assume tight coupling between genotype and phenotype at the level of self-replicating molecules that maintain this coupling throughout the transition stages from the RNA world to the prokaryotic level (2)

**Figure 1.** Comparison of two scenarios of the early evolution of cellular life: **(a)** the scenario envisioned when the term progenote was first defined (5, 6, 8); **(b)** the view that the last common ancestor was a prokaryote and that the progenote represents an earlier stage in the evolution of life (4, 10, 13).



At the same time that the term progenote was introduced to describe a pre-prokaryotic stage of cellular organization, an alternative meaning of the term progenote originated. Woese and Fox (1) argued that the last common ancestor of bacteria and the eukaryotic nucleocytoplasmic component was a progenote. Later, this argument was extended to also include a third line of descent, the archaea or archaebacteria ((5)–(7)). Woese and collaborators suggested that the three domains or Urkingdoms might have evolved independently from the progenote, that the optimization of the transcriptional and translational machinery occurred in parallel in the three lines of descent, and that the last common ancestor might be equally related to each of the three domains ((5)–(8)) (see Fig. 1a). Unfortunately, the initial definitions were ambiguous and potentially contradictory. Describing the evolution of bacteria and the eukaryotic nucleocytoplasmic component, Woese and Fox (1) summarize,

The two lines of descent, nevertheless shared a common ancestor, that was far simpler than the procaryote. This primitive entity is called a progenote, to denote the possibility that it had not yet completed evolving the link between geno and phenotype.

Taken literally, this paragraph labels the ancestor of bacteria and eukaryotic nucleocytoplasm a progenote and justifies the choice of name by stating that the last common ancestor might have been at a preprokaryotic level of organization. As a result, the term progenote is often used in the sense of progenitor to denote the last common ancestor of archaea, bacteria, and eukaryotes, not in the intended sense as a contrast to genote or eugenote, that is, organisms that have a “precise, accurate link between genotype and phenotype” (7).

Determining the properties of the last common ancestor is an important matter. [Horizontal Gene Transfer](#) and the fusion of formerly independent lineages have turned the tree of life into a net of life (9). Characters found in all three cellular lineages might have been present in the last common ancestor, but these shared characters also might have evolved much later in one of the lineages and spread into the other two domains by horizontal transfer (10–12). The congruence of many molecular phylogenies supports the scenario depicted in Figure 1b, and suggests that the last common ancestor had DNA and RNA polymerases, complex ribosomes made of both rRNA and proteins, and membranes already used for chemiosmotic coupling (13). However, if molecular phylogenies reflect horizontal gene transfer more than shared ancestry, recurring patterns might result from frequent gene transfers between some organisms and make the nature of the last common ancestor hard to discern. The use of ancient duplicated genes allows resolving the deep tripartite division of life into two successive bifurcations (14, 15). The current majority consensus considers that the archaea is a sister group of the eukaryotes ((4, 13, 16); see Fig. 1b). According to this view (4, 13), the last common ancestor was a prokaryote that had a DNA genome, elaborate transcriptional and translational machinery, and strong coupling between geno- and phenotype. Although DNA replication and transcription appear to have been further optimized independently in the three domains (17), the last common ancestor appears to have been a prokaryote, not a progenote. Therefore, to avoid confusion, the last common ancestor of all extant life should be denoted the universal ancestor (8) or cenancestor (18), and the term progenote should be reserved to denote a hypothetical preprokaryotic stage in cellular evolution (1).

## Bibliography

1. C. R. Woese and G. E. Fox (1977) *J. Mol. Evol.* **10**, 1–6.
2. M. Eigen and P. Schuster (1978) *Naturwissenschaften* **65**, 341–369.
3. H. Maturana and F. Varela (1980) *Autopoiesis and Cognition*, D. Reidel, Dordrecht, Holland.
4. W. F. Doolittle (1996) In *Evolution of Microbial Life* (D. McL. Roberts, P. Sharp, G. Alderson, and M. Collins, eds.), Society for General Microbiology Symposium 54, University Press, Cambridge, U. K., pp. 1–21.
5. C. R. Woese, L.J. Magrum, and G.E. Fox (1978) *J. Mol. Evol.* **11**, 245–251.
6. C. R. Woese and G. E. Fox (1977) *Proc. Natl. Acad. Sci. U.S.A.* **74**, 5088–5090.

7. C. R. Woese (1987) *Microbiol. Rev.* **51**, 221–271.
8. G. E. Fox, E. Stackebrandt, R. B. Hespell, J. Gibson, J. Maniloff, T. A. Dyer, R. S. Wolfe, W. E. Balch, R. S. Tanner, L. J. Magrum, L. B. Zablen, R. Blakemore, R. Gupta, L. Bonen, B. J. Lewis, D. A. Stahl, K. R. Luehrsen, K. N. Chen, and C. R. Woese (1980) *Science* **209**, 457–463.
9. E. Hilario and J. P. Gogarten (1993) *BioSystems* **31**, 111–119.
10. J. P. Gogarten (1995) *Trends Ecol. Evol.* **10**, 147–151.
11. C. R. Woese (2000) *Proc. Natl. Acad. Sci. U.S.A.* **97**, 8392–8396.
12. W. F. Doolittle (1999) *Science* **284**, 2124–2128.
13. J. P. Gogarten, E. Hilario, and L. Olendzenski (1996) In *Evolution of Microbial Life* (D. McL. Roberts, P. Sharp, G. Alderson, and M. Collins, eds.), Society for General Microbiology Symposium 54, University Press, Cambridge, U. K., pp. 267–292.
14. J. P. Gogarten, H. Kibak, P. Dittrich, L. Taiz, E. J. Bowman, B. J. Bowman, M. F. Manolson, R. J. Poole, T. Date, T. Oshima, J. Konishi, K. Denda, and M. Yoshida (1989) *Proc. Natl. Acad. Sci. U.S.A.* **86**, 6661–6665.
15. N. Iwabe, K-I. Kuma, M. Hasegawa, S. Osawa, and T. Miyata (1989) *Proc. Natl. Acad. Sci. U.S.A.* **86**, 9355–9359.
16. C. R. Woese, O. Kandler, and M. L. Wheelis (1990) *Proc. Natl. Acad. Sci. U.S.A.* **87**, 4576–4579.
17. D. R. Edgell and W. F. Doolittle (1987) *Cell* **89**, 995–998.
18. W. M. Fitch and K. Upper (1987) *Cold Spring Harbor Symp. Quant. Biol.* **52**, 759–767.

### Suggestions for Further Reading

19. D. McL. Roberts, P. Sharp, G. Alderson and M. Collins, eds. (1996) *Evolution of Microbial Life*, Society for General Microbiology Symposium 54, University Press, Cambridge, U. K.
20. H. J. Morowitz (1992) *Beginnings of Cellular Life*. Yale University Press, New Haven, CT.
21. C. R. Woese (1987) *Microbiol. Rev.* **51**, 221–271.
22. W. F. Doolittle (2000) *Sci. Am.* **282**, 90–95.
23. D. W. Deamer and G. R. Fleischaker, eds. (1994) *Origins of Life*, Jones and Bartlett, London.

### Programmed Cell Death

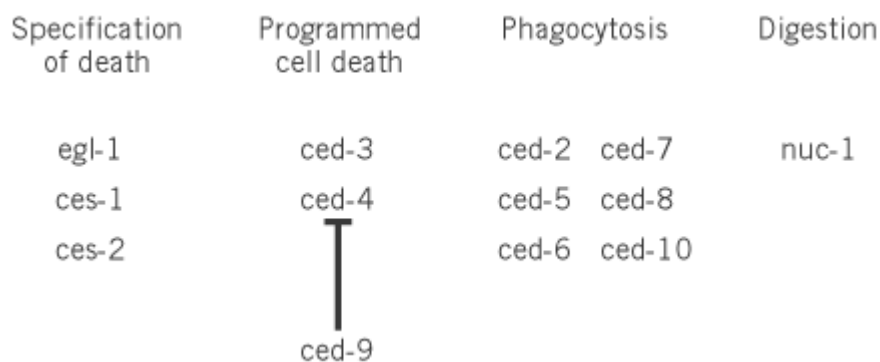
*Programmed cell death* (PCD) is the term used to describe an active form of [cell death](#) that occurs during invertebrate and vertebrate [development](#). PCD shares all the attributes of [apoptosis](#), but the term apoptosis is generally restricted to the description of active mammalian cell death, whereas PCD describes all other active cell deaths. PCD pathways have been characterized in several invertebrate systems ([1](#)), but none so clearly as that occurring in the [nematode](#) worm *Caenorhabditis elegans*.

A definitive plan of **gene expression** during cell development and death is seen in *C. elegans* ([2-4](#)). The cell lineage development pattern within this animal has been extensively investigated, which established that exactly 131 cells of the 1090 within the developing nematode die. Some of the genes involved in these deaths are shown in [Fig. 1](#). Of the genes involved in nematode PCD, three have been investigated extensively. **Cell death gene 3** (*ced-3*) and *ced-4* are directly responsible for the death of specific cells during development. Loss-of-function mutants of *ced-3* or *ced-4* result in an



animal where cells that normally die during development do not, but instead take up the same differentiated form as their sister cells (2). Ced-3 and Ced-4 gene products act within the cells that die and not as killer signals produced by neighboring or surrounding cells (5). Ced-3 and Ced-4 killing is suppressed by another gene product, Ced-9. Although *ced-9* gain-of-function mutants exhibit no cell death, cell death does occur in *ced-9* loss-of-function mutants. Moreover, this death occurs not only in the 131 cells that normally die, but also in other cells within the nematode (6). These observations suggest that *ced-9* is a general genetic suppressor of cell death within these animals.

**Figure 1.** The genetic control and pathway of programmed cell death in *C. elegans*. The ‘killing’ of the cell is carried out by the products of two genes, *ced-3* and *ced-4*, whose activities are regulated by *ced-9*. Many of the genes that act within this pathway have been conserved throughout the evolution of multicellular organisms, demonstrating the importance of regulated cell death (2, 7, 9)



The notion that cell death is essential for the development of the organism is verified by the conservation of death genes throughout multicellular [evolution](#). For example, mammalian homologues of two of the nematode *ced* genes have been identified. Gene *ced-9* is a homologue of the mammalian anti-apoptotic proto-oncogene *bcl-2* (7), the **B-cell** leukemia/lymphoma gene (8) that suppresses apoptotic cell death, and a host of other *bcl-2*-like genes that define a gene family involved in cell death regulation (9, 10). **Transgenic** expression of *bcl-2* in *ced-9* loss-of-function nematode mutants suppresses cell death, demonstrating the conservation of function between these two homologues (11). This conservation of function suggests that genes with functions similar to those of *ced-3* and *ced-4* should exist in mammalian cells. This hypothesis was borne out by the discovery of a *ced-3* mammalian homologue, the *interleukin-1b-converting enzyme* (ICE) (12). ICE is a [thiol proteinase](#) that activates pro-interleukin-1b by **proteolytic** cleavage to its active form during an inflammatory response (13). Expression of ICE in fibroblasts results in apoptosis, as does expression of Ced-3, demonstrating the conservation of function between these two homologues (12). There are currently 14 members of the Ced-3/ICE gene family in mammals (14-16) (more recently termed [caspases](#) for cysteine **asp**artate proteinases (17)) and two in *Drosophila*. A mammalian homologue of *ced-4* has been identified, and its function is slowly being deciphered (18-21). Ced-4 binds both Ced-9 and Ced-3, resulting in a protein complex that somehow acts to regulate PCD (see [Apoptosis](#)). Although the exact role of Ced-4 in this complex is unclear, Ced-4 induces death in both mammalian and yeast cells, suggesting some intrinsic death capacity (19, 21).

Other genes within the *C. elegans* PCD pathway are involved in the [phagocytosis](#) of dead cells. Mutant worms that lack one or more of these genes do not phagocytose dead cells; instead, the cells are left within the cellular tissues with no apparent deleterious effects (22), unlike unphagocytosed cells in mammalian tissues. Once again, some of these genes are also conserved in higher organisms, demonstrating the conservation of the PCD pathway as a whole throughout evolution (23).

The fruit fly *Drosophila* also has conserved PCD pathways. *Drosophila* embryos exhibit PCD at approximately 7 hours after egg laying (stage 11 embryos), with cell death becoming widespread throughout the embryo at stages 12 and 13. More specific regions of cell death are also seen during dorsal closure (stage 14) and advanced head involution (stage 15) (24). PCD in *Drosophila* climaxes with prominent cell death throughout the central nervous system as the ventral nerve cord condenses. As in *C. elegans*, the genes required for PCD in *Drosophila* are being sought and characterized. One such gene, *reaper* (*rpr*), is required for virtually all PCD in the developing embryo. *rpr*-deficient embryos have many extra cells and fail to hatch, demonstrating the requirement for *rpr*-regulated cell death during development (25). These embryos are also resistant to X-irradiation-induced death at low doses. That cell death occurs at higher radiation doses suggests that the death pathway is intact but less sensitive in the absence of *rpr* (26). Two other genes, Grim (27) and Head Involution Defective (*hid*) (28), are also required for some of the PCD occurring during development. Acting downstream of these genes are the Ced-3-like caspase homologues. Two of these genes have so far been identified, *Drosophila* cell death proteinase-1 (DCP-1) (29) and *Drosophila* ICE (DrICE) (30). Like their mammalian counterparts, these caspase homologues cleave specific protein substrates during PCD. However, unlike the hierarchical arrangement of the caspases in mammalian cells, DrICE appears to be able to autoactivate and carry out PCD in *Drosophila* cells without requiring DCP-1 activity [Fraser, personal communication].

Other insect models have been successfully used to document genes regulating PCD. Following metamorphosis, the intersegmental muscles of the tobacco hawk moth, *Manduca sexta*, are no longer required and undergo degeneration, a process triggered by a reduction in the levels of an ecdysone hormone (31-33). This type of PCD involves the activation of several genes, plus the process of attachment of multiple ubiquitin molecules to other proteins, which is probably required for tagging them for protein degradation. Additionally, several genes are downregulated during the commitment of these cells to PCD, and over 30 new polypeptides are translated during muscle degeneration (34).

Examples of PCD are found in all multicellular organisms so far studied, and many of the genes that regulate these deaths have been conserved throughout evolution. This makes the study of invertebrate cell death invaluable in identifying how PCD is regulated in more complex organisms such as mammals.

## Bibliography

1. M. D. Jacobson, M. Weil, & M. C. Raff (1997) Programmed cell death in animal development. *Cell* **88**, 347–354.
2. H. M. Ellis, and H. R. Horovitz (1986) Genetic control of programmed cell death in the nematode *C. elegans*. *cell* **44**, 817–829.
3. J. E. Sulston, E. Schierenberg, J. G. White, and J. N. Thomson (1983) The embryonic cell lineage of the nematode *Caenorhabditis elegans*. *Dev. Biol.* **100**, 64–119.
4. E. M. Hedgecock, J. E. Sulston, and J. N. Thomson (1983) Mutations affecting programmed cell deaths in the nematode *Caenorhabditis elegans*. *Science* **220**, 1277–1279.
5. J. Y. Yuan, and H. R. Horvitz (1990) The *Caenorhabditis elegans* genes *ced-3* and *ced-4* act cell autonomously to cause programmed cell death *Dev. Biol.* **138**, 33–41.
6. M. O. Hengartner, R. E. Ellis, and H. R. Horvitz (1992) *Caenorhabditis elegans* gene *ced-9* protects cells from programmed cell death. *Nature* **356**, 494–499.
7. M. Hengartner, and H. Horvitz (1994) *C. elegans* cell survival gene *ced-9* encodes a functional homolog of the mammalian proto-oncogene *bcl-2*. *Cell* **76**, 665–676.
8. Y. Tsujimoto, L. R. Finger, J. Yunis, P. C. Nowell, and C. M. Croce (1984) Cloning of the chromosome breakpoint of neoplastic B cells with the t(14;18) chromosome translocation. *Science* **226**, 1097–1099.
9. G. Williams, and C. Smith (1993) Molecular regulation of apoptosis—genetic controls on cell-death. *Cell* **74**, 777–779.

10. J. C. Reed (1997) Double identity for proteins of the Bcl-2 family. *Nature* **387**, 773–776.
11. D. L. Vaux, I. L. Weissman, and S. K. Kim (1992) Prevention of programmed cell-death in *Caenorhabditis elegans* by human bcl-2 *Science* **258**, 1955–1957.
12. J. Y. Yuan, S. Shaham, S. Ledoux, S., H. M. Ellis, and H. R. Horvitz (1993) The *C. elegans* cell-death gene *ced-3* encodes a protein similar to mammalian interleukin-1 beta converting enzyme. *Cell* **75**, 641–652.
13. N. A. Thornberry, H. G. Bull, J. R. Calaycay, K. T. Chapman, A. D. Howard, M. J. Kostura, D. K. Miller, S. M. Molineaux, J. R. Weidner, J. Aunins, and et al. (1992) A novel heterodimeric cysteine protease is required for interleukin-1 beta processing in monocytes. *Nature* **356**, 768–774.
14. M. Whyte (1996) ICE/Ced-3 proteases in apoptosis. *Trends Cell Biol.* **6**, 245–248.
15. A. Takahashi, and W. Earnshaw (1996) Ice-related proteases in apoptosis. *Curr. Opin. Gen. & Dev.* **6**, 50–55.
16. Y. Lazebnik, A. Takahashi, G. Poirier, S. H. Kaufmann, and & W. Earnshaw (1995) Characterization of the execution phase of apoptosis *in vitro* using extracts from condemned-phase cells. *J. Cell Sci.* **19**, 41–49.
17. E. Alnemri, D. Livingston, D. Nicholson, G. Salvesan, N. Thornberry, W. Wong, and J. Yuan (1996) Human ICE/CED-3 protease nomenclature. *Cell* **87**, 171.
18. D. Wu, H. D. Wallen, and G. Nunez (1997) Interaction and regulation of subcellular localization of CED-4 by CED-9. *Science* **275**, 1126–1129 (see comments).
19. A. M. Chinnaiyan, K. O'Rourke, B. R. Lane, and V. M. Dixit (1997) Interaction of CED-4 with CED-3 and CED-9: a molecular framework for cell death. *Science* **275**, 1122–1126 (see comments).
20. M. S. Spector, S. Desnoyers, D. J. Hoepfner, and M. O. Hengartner (1997) Interaction between the *C. elegans* cell-death regulators CED-9 and CED-4 *Nature* **385**, 653–656.
21. C. James, S. Gschmeissner, A. Fraser, and G. Evan (1997) Ced-4 induces chromatin condensation in *S. pombe* and is inhibited by direct physical association with Ced-9. *Curr. Biol.* **7**, 246–252.
22. R. E. Ellis, D. M. Jacobson, and H. R. Horvitz (1991) Genes required for the engulfment of cell corpses during programmed cell death in *Caenorhabditis elegans* *Genetics* **129**, 79–94.
23. M. F. Luciani, and G. Chimini (1996) The ATP binding cassette transporter ABC1, is required for the engulfment of corpses generated by apoptotic cell death. *Embo. J.* **15**, 226–235.
24. J. M. Abrams, K. White, L. I. Fessler, and H. Steller (1993) Programmed cell death during *Drosophila* embryogenesis. *Development* **117**, 29–43
25. K. White, M. E. Grether, J. M. Abrams, L. Young, K. Farrell, and H. Steller (1994) Genetic control of programmed cell death in *Drosophila*. *Science* **264**, 677–683 (see comments).
26. H. Steller, J. M. Abrams, M. E. Grether, and K. White (1994) Programmed cell death in *Drosophila*. *Phil. Trans. R Soc. Lond. B Biol. Sci.* **345**, 247–250.
27. P. Chen, W. Nordstrom, B. Gish, and J. M. Abrams (1996) grim, a novel cell death gene in *Drosophila*. *Genes Dev.* **10**, 1773–1782.
28. M. E. Grether, J. M. Abrams, J. Agapite, K. White, and H. Steller (1995) The head involution defective gene of *Drosophila melanogaster* functions in programmed cell death. *Gene Dev.* **9**, 1694–1708.
29. Z. Song, K. McCall, and K.H. Steller (1997) DCP-1, a *Drosophila* cell death protease essential for development. *Science* **275**, 536–539.
30. A. Fraser, and G. Evan (1997) Identification of a *Drosophila melanogaster* ICE/CED-3-related protease, drICE. *EMBO. J.* **16**, 2805–2813.
31. M. Schwartz, J. W. Truman (1982) Peptide and steroid regulation of muscle degeneration in an insect. *Science* **215**, 1420.

32. L. M. Schwartz, J. W. Truman (1983) Hormonal control of rates of metamorphic development in the tobacco hornworm *Manduca sexta*. *Dev. Biol.* **99**, 103.
33. L. M. Schwartz, L. Kosz, L. and B. K. Kay (1990) Gene activation is required for developmental programmed cell death. *Proc. Natl. Acad. Sci. USA* **87**, 6594.
34. A. G. Wadewitz, and R. A. Lockshin (1988) Programmed cell death: Dying cells synthesize a co-ordinated, unique set of proteins in two different episodes of cell death. *FEBS Lett.* **241**, 19–23.

### Suggestions for Further Reading

35. M. O. Hengartner (1994) Programmed cell death. A rich harvest. *Curr. Biol.* **4**, 950–952.
36. R. E. Ellis, J. Y. Yuan, and H. R. Horvitz (1991) Mechanisms and functions of cell death. *Ann. Rev. Cell Biol.* **7**, 663–698.
37. H. Steller, J. M. Abrams, M. E. Grether, and K. White (1994) Programmed cell death in *Drosophila*. *Phil. Trans. R Soc. Lond. B Biol. Sci.* **345**, 247–250.

## Prokaryotic Genetics

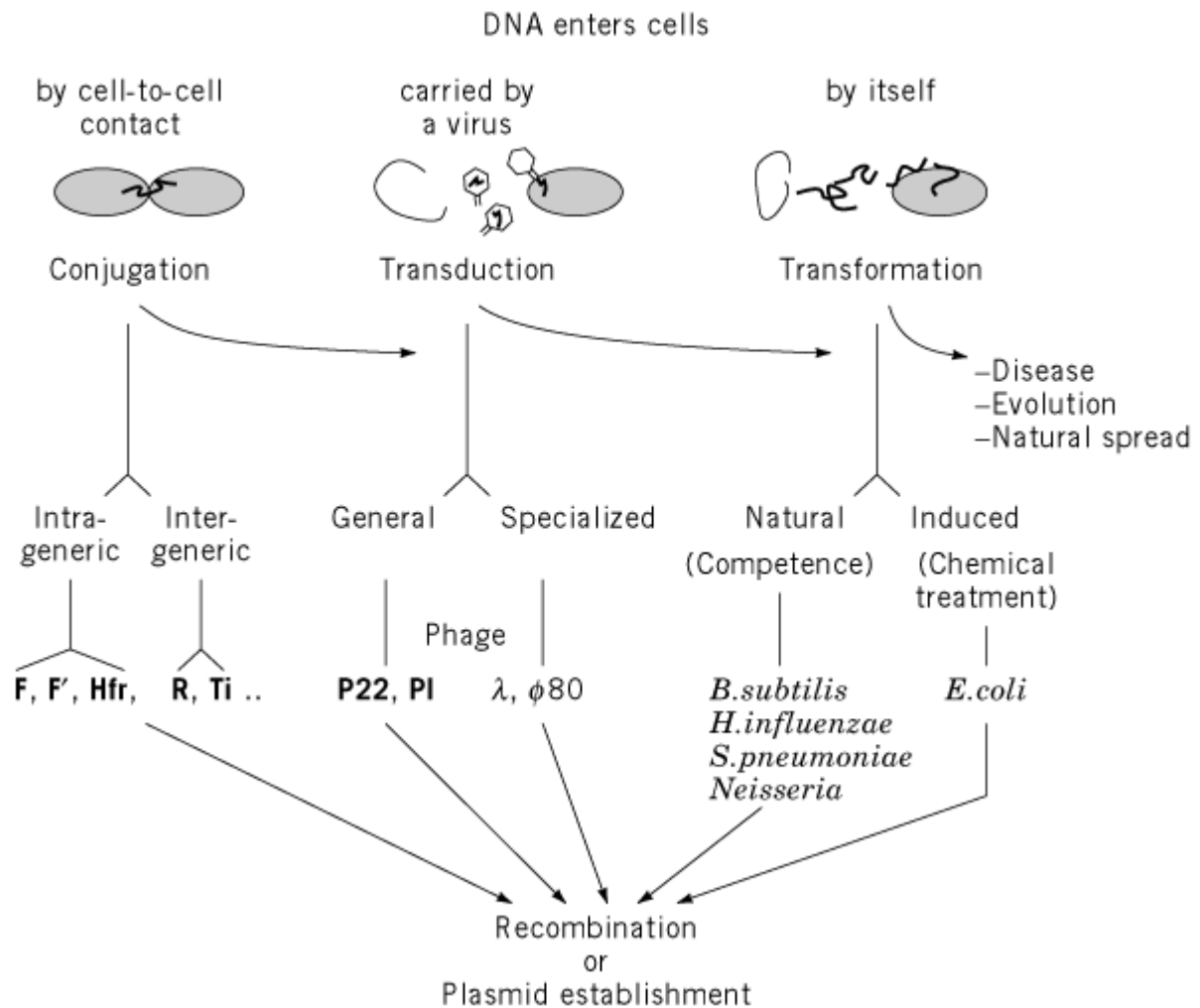
Prokaryotic genetics took biology's center stage for about 30 years, transforming the classic genetic study of abstract factors of inheritance into the assertively molecular style of genetics that we know today. The impact of prokaryotic genetics on biology has been succinctly described by Thomas Brock:

Almost every area of modern biology under active research study today owes a major debt to bacterial genetics. Molecular biology, immunology, cancer research, medical virology, epidemiology, genetics of higher organisms, evolution, taxonomy, cell biology, and developmental biology all depend on concepts that arose first from studies in bacterial genetics.

To this must now be added whole-genome sequencing, scarcely underway in 1990, when the above words were written, but now usurping the function of genetic crosses in the mapping and definition of genes.

The purpose of this section is to highlight some of the key phenomena on which prokaryotic genetics is based. In practice, genetics is based on exploiting natural modes of gene transfer from one cell to another (Fig. 1). Individual entries are devoted to each of the known methods of transfer—[transformation](#), **transduction**, and **conjugation**—as well as to certain devices that that played a central role in putting gene transfer to work: **replica plating**, development of [competence](#) for transformation, **merodiploidy**, **Hfr's and F-primes**, and the **F-plasmid**. While the heyday of this traditional prokaryotic genetics is over, an appreciation of its essentials, as represented by the topics presented here, remains important as genetic investigations spread to the wider reaches of the microbial world.

**Figure 1.** The primary methods of transfer of genetic information in prokaryotes. Aspects for which there are individual entries are in bold.



## Proliferating Cell Nuclear Antigen (PCNA)

Proliferating cell nuclear antigen (PCNA) was originally identified as an [antigen](#) of the autoimmune disease systemic lupus erythematosus (1). This protein localizes in **nuclei** in a **cell-cycle**-dependent manner and was also called *cyclin*, due to its periodic appearance in **S-phase** cells. The name of PCNA is now used for this protein so as to avoid confusion with the [cyclin](#) family of proteins that bind to cyclin-dependent protein **kinases** (cdk). Staining of growing cells with an [antibody](#) against PCNA showed a punctuated pattern in S-phase nuclei that colocalized with bromodeoxyuridine foci, or [replication foci](#) (2). This observation suggested that this protein was closely involved in [DNA replication](#). This speculation was confirmed by the second discovery of PCNA as the essential replication factor in **SV40 virus** replication *in vitro* (3). Concomitantly, PCNA was identified as an accessory protein of **DNA polymerase d** (pol d) (4).

Without PCNA, pol d has limited activity and synthesizes only several tens of nucleotides into DNA, but the addition of PCNA greatly stimulates the activity and makes the polymerase very processive, to synthesize long DNA strands. These characteristics implied that PCNA might function as a **eukaryotic** processivity factor as the b-subunit (b-clamp) of *Escherichia coli* DNA polymerase III.

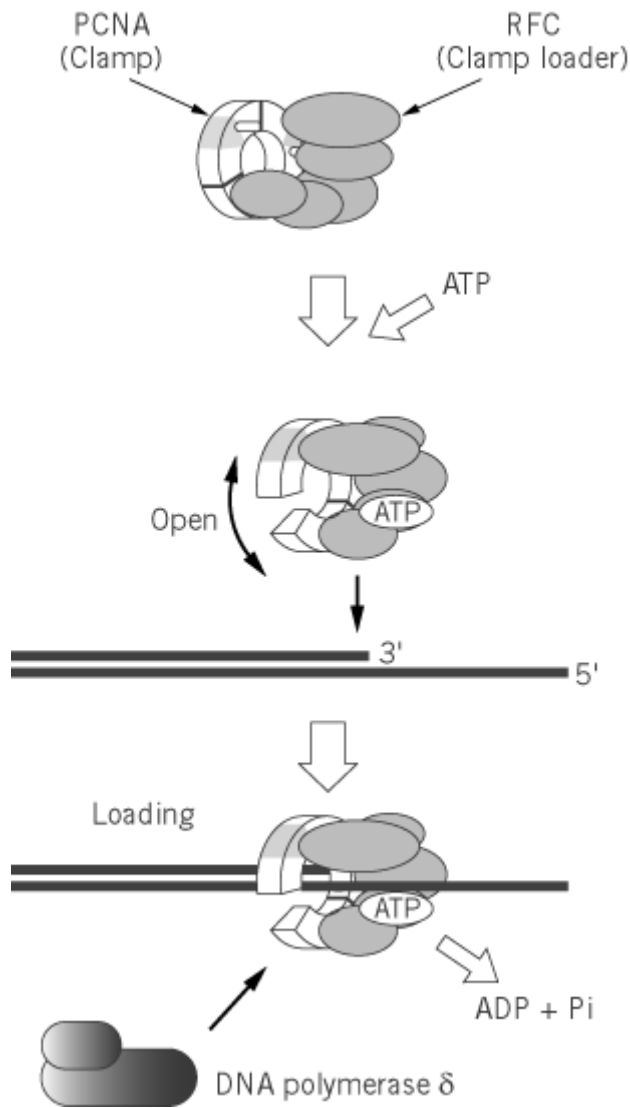
The omission of PCNA from the complete SV40 replication reaction results in an accumulation of short [Okazaki Fragments](#), rather than the full-length products, indicating that PCNA is essential to coordinate the elongation of [leading and lagging strands](#) (5). Another newly identified replication factor, RFC (replication factor C), was also required at the same elongation stage and to coordinate two-strand synthesis, suggesting that RFC, PCNA, and pol d function at the same DNA elongation stage of eukaryotic DNA replication (6).

Comparable studies of the elongation process of *E. coli*, **bacteriophage**, and eukaryotic systems have revealed that there is a striking functional conservation in their replication apparatus. Factors called DNA **sliding clamps** were commonly identified in all of them, such as PCNA in eukaryotes, DNA polymerase III (pol III) b-subunit in *E. coli* and gene45 protein in T4 phage. The name DNA sliding clamp came from the experimental observation that once the clamp molecule binds on circular DNA, it stably associates with the DNA, but dissociates quickly upon linearization of the DNA (7). This property suggests the protein has a ring structure that binds the DNA topologically and can slide freely along the DNA. Comparison of the predicted **secondary structures** of clamp molecules demonstrated striking conserved alignments, although there are almost no similarities in their [primary structures](#). The monomer sizes of b-clamp, PCNA, and gene45 are 41, 29, and 25 kDa, respectively, but their native molecular weights are similar. This is because PCNA and gene45 have homotrimer structures, whereas, in contrast, the b-clamp has a homodimer structure.

The structures determined by [X-ray crystallography](#) of the b-clamp and PCNA demonstrated that they have very similar ring structures composed of monomers joined head to tail (8). The monomers of b-clamp and PCNA have three and two internal repeating structures, respectively, so both rings are composed of six repeating domains that generate a hexagonal shape. Each domain unit has a stable triangular structure comprising nine antiparallel [beta strands](#) facing outward and two **alpha-helices** facing the center of the ring. The inside surfaces are strongly positively charged, and the hole has a size sufficient to accommodate double-stranded DNA. Thus, this class of proteins has a new sequence-nonspecific, topological type of DNA binding, which permits these proteins to slide along a duplex DNA stably but freely in either direction.

A crucial process for PCNA function is its loading onto DNA prior to its association with pol d. However, loading of the ring molecule onto DNA without the free ends will be topologically difficult. Another accessory protein, RFC, has the activity of loading PCNA at the 3'-end of a [primer](#) on circular DNA in the presence of ATP. During this process, ATP-bound RFC interacts with PCNA and may induce the opening of its ring to pass template DNA through it. Since further ATP hydrolysis by RFC is necessary to load pol d on the DNA, the second structural change of the RFC-PCNA complex will occur by ATP hydrolysis and induce the association of pol d with PCNA (Fig. 1) (see also [Clamp Loaders, Processivity Complex](#)).

**Figure 1.** Loading of PCNA onto a template DNA and formation of pol d complex.



Besides its function as a replication factor, PCNA also plays crucial roles in [DNA repair](#) processes. Studies with *in vitro* nucleotide [excision repair](#) and [mismatch repair](#) reactions indicated that PCNA functions as a component of the DNA polymerase complex for repair DNA synthesis. In addition, if a cell is irradiated by UV light, PCNA accumulates in nuclei, where the DNA repair reaction occurs, although the cell is not in **S phase** (9). PCNA also interacts with various proteins involved in **cell-cycle** regulation, including CDK/cyclin, p21(Cip1/Waf1), and GADD45 (for a review, see (10)). The binding of p21 to PCNA resulted in inhibition of the elongation reaction by pol d, suggesting a mechanism to regulate the functions of the replication apparatus by these cell-cycle regulators. The importance of these secondary function of PCNA is still ambiguous, but suggest a novel mechanism to directly connect cell-cycle control and the DNA replication apparatus.

#### Bibliography

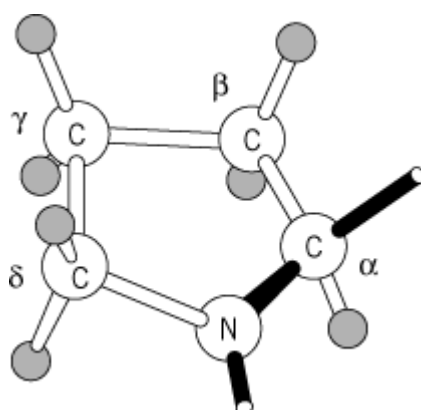
1. K. Miyachi K. et al. (1978) J. Immun. **121**, 2228–2234.
2. R. Bravo and H. Macdonald-Bravo (1987), J. Cell Biol. **105**, 1549–1554.
3. G. Prelich et al. (1987) Nature **326**, 471–475.
4. G. Prelich et al. (1987) Nature **326**, 517–520.
5. G. Prelich and B. Stillman (1988) Cell **53**, 117–126.
6. T. Tsurimoto and B. Stillman (1989) EMBO J. **8**, 3883–3889.

7. P. T. Stukenberg et al. (1991) *J. Biol. Chem.* **266**, 11328–11334.
8. J. Kuriyan and M. O'Donnell (1993) *J. Mol. Biol.* **234**, 915–925.
9. R. L. Gregory et al. (1996) *Curr. Biol.* **6**, 189–199.
10. Z. Kelman (1997) *Oncogene* **14**, 629–640.

## Proline (Pro, P)

The imino acid proline is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to four **codons**—CCU, CCC, CCA, and CCG—and represents approximately 5.1% of the residues of the proteins that have been characterized. The prolyl residue incorporated has a mass of 97.12 Da, a **van der Waals volume** of  $90 \text{ \AA}^3$ , and an [accessible surface area](#) of  $143 \text{ \AA}^2$ . Pro residues are changed infrequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [alanine](#) and [serine](#) residues.

The side chain of Pro is unique in that it is bonded covalently to the nitrogen atom of the peptide backbone, which is indicated by the solid bonds:



Therefore, the peptide backbone at Pro residues has no amide hydrogen for use as a donor in **hydrogen bonding** or in resonance stabilization of the [peptide bond](#) of which it is part. The cyclic five-membered ring is invariably puckered, with the C<sup>α</sup>, C<sup>γ</sup>, C<sup>δ</sup>, and N atoms approximately coplanar, but with the C<sup>β</sup> atom displaced about  $0.5 \text{ \AA}$  from the plane. This ring imposes rigid constraints on rotation about the N–C<sup>α</sup> bond of the backbone; the value of the torsion angle  $\phi$  is restricted to approximately  $-60^\circ$  (see [Ramachandran Plot](#)). The peptide bond preceding Pro residues also has an intrinsic tendency to adopt the [cis configuration](#) approximately 10% of the time in model peptides (1). This configuration is also found frequently at certain Pro peptide bonds in folded [protein structures](#). The synthetic polypeptide Poly(Pro) adopts one of two regular conformations, depending on the conditions, known as Poly(Pro)I and II. Form I has all *cis* peptide bonds, whereas form II has all *trans*. The values of  $\phi$  and  $\psi$  are very similar for both, but form I is a right-handed helix with 3.3 residues per turn, whereas form II is a left-handed helix with 3.0 residues per turn (2). The two forms can be interconverted by changing the solvent; the interconversion occurs by a “zipper” mechanism in which the intrinsically slow *cis–trans* isomerization of each peptide bond starts at one end and progresses sequentially along the polypeptide chain. The *I* → *II* interconversion starts at the amino end of the chain, whereas the reverse interconversion starts at the other end (3).



The side-chain atoms of Pro residues are [nonpolar](#) and chemically inert, but only about 18% of Pro residues are fully buried in [protein structures](#). Instead, they often occur at the protein surface, because they are frequently used in reverse [turns](#). In contrast, the geometry of the Pro residue is incompatible with the **a-helix** and **b-sheet** types of **secondary structure**, and the Pro residue destabilizes an a-helix by about 3 kcal/mol in model peptides. Nevertheless, Pro residues fit well at the *N*-terminus of an a-helix, and single Pro residues can be accommodated in long a-helices by distorting the helical geometry locally. Pro residues have special roles in the triple helix of **collagens**, where they are often modified by [hydroxylation](#) at the  $\alpha$  carbon when in the sequence–Xaa–Pro–Gly– or at the  $\beta$  carbon when in the sequence –Gly–Pro–.

#### Bibliography

1. G. N. Ramachandran and A. K. Mitra (1976). *J. Mol. Biol.* **107**, 85–92.
2. G. N. Ramachandran et al. (1966) *Biochim. Biophys. Acta* **112**, 168–170.
3. L. N. Lin and J. F. Brandts (1980) *Biochemistry* **19**, 3055–3059.

#### Suggestion for Further Reading

4. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

### Promiscuous Dna

Interactions between the [nucleus](#) and **organelles** of **cytoplasm** have a major role in eukaryotic cellular metabolism. **ATPases**, some [cytochromes](#), for example, the yeast cytochrome *c* oxidase complex, are assembled posttranslationally in [mitochondria](#) and **chloroplasts** from some polypeptide subunits synthesized in the nucleus, others in the organelle. The genetic materials of the nucleus, mitochondrion, and chloroplast are usually maintained discretely in each organelle without duplication elsewhere. Exceptions are known as *promiscuous DNA*, defined as a nucleotide sequence that occurs in more than one of these three membrane-bound organelles. This should be distinguished from the DNA of promiscuous **plasmids**, which can be transferred to a wide range of prokaryotic cells.

Promiscuous DNA was first reported upon the discovery that the genomes of mitochondria and chloroplasts of **maize** have a 12-kb DNA sequence in common. The most attractive conclusion is that the chloroplast DNA became incorporated in the mitochondrial genome after [gene duplication](#) and [transposition](#) (1).

A similar phenomenon was found in **Saccharomyces cerevisiae**. A piece of DNA is common to the mitochondrial and nuclear genomes. There is good evidence that the “jumping” occurred in the direction nucleus → mitochondrion (2).

A third example comes from the sea urchin, *Strongylocentrotus purpuratus*, whose sections of mitochondrial DNA, coding for cytochrome oxidase subunit 1 and the 3' 16 S RNA of [ribosomes](#), are found in the nuclear genome. This is probably a result of a **chromosomal** transposition about 25 million years ago, followed by chromosomal rearrangements and single nucleotide substitutions (3).

#### Bibliography

1. D. B. Stern and D. M. Lonsdale (1982) *Nature* **299**, 698–702.
2. F. Farrelly and R. A. Butow (1983) *Nature* **301**, 296–301.
3. H. T. Jacobs et al. (1983) *J. Mol. Biol.* **165**, 609–632.

## Pronase

The K-1 strain of *Streptomyces griseus* synthesizes and secretes a mixture of [proteinases](#) that for many years has been commercially available under the proprietary name of “Pronase.” Initially offered as a purified proteinase of broad specificity and widely adopted for this purpose, it became clear that pronase is composed of a variety of exo- and **endopeptidases** ([1](#)), thus accounting for its remarkable ability to catalyze the hydrolysis of many types of [peptide bonds](#). Among the different constituents are more than a dozen separate enzymes having the specificities of **trypsin**, **chymotrypsin**, [elastase](#), [subtilisin](#), [carboxypeptidase](#), and [aminopeptidase](#). Many of these have been isolated in pure form, starting with the crude mixture, and subsequently characterized both structurally and kinetically. The commercial product is still referred to as “an unusually nonspecific proteinase,” although it is now well known to be a mixture of proteinases. It is free of **nuclease** activity and hence is convenient for removing proteins in the course of procedures for the isolation of **nucleic acids**.

### Bibliography

1. M. W. Awad, Jr. et al. (1976) In *Proteolysis and Physiological Regulation* (D. W. Ribbons and K. Brew, eds.), Academic Press, New York, pp. 77–91.

## Pro-Protein

Many [proteins](#) destined for [protein secretion](#) undergo a fast maturation process during translocation through membranes by cleavage of [signal peptides](#). Some secretory proteins, however, go through an intermediate stage of the relatively long-lived intracellular form termed pro-protein (See [Pre-Protein](#), [Pre-Pro-Protein](#) and [Pro-Sequence](#)). Some serum proteins (eg, [Albumin](#)) or [hormones](#) (eg, [insulin](#) and [glucagon](#)) are synthesized as longer inactive precursors, from which specific polypeptides (or pro-sequences) are cleaved to generate mature active molecules ([1](#), [2](#)). In general, proteolytic conversion of the pro-protein to the mature form involves endoproteolytic cleavage at the site after an amino acid recognition sequence, such as Arg–Arg or Lys–Arg. Additional sequences can be cleaved at the *N*-terminus (eg, in the case of proalbumin) or at both ends of the pro-protein (eg, proglucagon). In proinsulin, the extra sequence, named the C peptide, is located internally in the pro-protein.

### Bibliography

1. L. C. Lopez, M. L. Frazier, C. J. Su, A. Kumar, and G. F. Saunders (1983) *Proc. Natl. Acad. Sci. USA* **80**, 5485–5489.
2. D. F. Steiner, S. P. Smeekens, S. Ohagi, and S. J. Chan (1992) *J. Biol. Chem.* **267**, 23435–23438.

## Pro-Sequence

A pro-sequence is part of a protein and is required for protein folding or its maintenance in an inactive form (see [Pre-Protein](#), [Pre-Pro-Protein](#) and [Pro-Protein](#)). Pro-sequences are removed from exported protein after their translocation across a [membrane](#). Several of these pro-peptides are composed of long polypeptide chains. For example, the propeptide in pro- $\alpha$ -lytic proteinase has 174 amino acid residues; there are 91 in pro-carboxypeptidase Y, and 77 in pro-[subtilisin](#). Long pro-peptides are typical for [proteinases](#). The pro-sequence is usually attached to the *N*-terminus of the proteinase (although pro-sequences are sometimes also found to be *C*-terminal extensions, or a combination of the two). There is evidence that some pro-sequences are necessary for proper folding of the proteins. The concept of pro-sequence-assisted protein folding was originally demonstrated for [serine proteinases](#). The best-studied examples include the proteinases subtilisin BPN', which is secreted from *Bacillus amyloliquefaciens*,  $\alpha$ -lytic protease from *Lysobacter enzymogenes*, and the vacuolar carboxypeptidase Y from *Saccharomyces cerevisiae*. The pro-sequence in all three cases is important for mediating proper folding of the corresponding proteinase domain *in vivo* and *in vitro*. In the absence of pro-sequences, these enzymes fail to form an active proteinase conformation. After pro-proteins are folded properly, their pro-sequences are cleaved off in an intramolecular (via autocatalysis) or intermolecular (via a second proteinase) reaction. Pro-sequences of  $\alpha$ -lytic protease and subtilisin are known to accelerate proper folding of the proteinase domain by lowering a high-energy barrier on the folding pathway; that is, the pro-sequence directly catalyzes the folding reaction. Therefore, the pro-sequence can be classified as an intra-**molecular chaperone**.

### Suggestion for Further Reading

J. Eder and A. R. Fersht (1995) *Mol. Microbiol.* **16**, 609–614.

## Protamine

The classic protamines are small basic proteins, originally defined in teleosts, that are rich in [arginine](#) residues, but deficient in [lysine](#) and [cysteine](#). Their length is about 30 amino acid residues, and clusters of arginine residues are distributed throughout the sequence. Protamine–DNA complexes often represent the final state of [chromatin](#) compaction in fish [sperm](#). The organization of sperm chromatin containing protamines does not include **nucleosomes** (1). Protamines fold up into  $\alpha$ -helical structures when they bind to double-stranded **nucleic acid**.

Mammalian protamines are more complex than those in fish. They are arginine-rich and contain cysteine residues, but remain deficient in lysine. Three distinct regions exist within the approximately 50-residue mammalian protamine sequence, an N-terminal segment of 15 residues that always begins with the same four amino acids, a central region of 25 residues that is arginine-rich, with 3 to 4 clusters of arginine clusters separated by neutral amino acids, and a variable C-terminus. The cysteine residues are distributed throughout the protein at relatively conserved positions. Protamine genes are generally specific to the male germ line in their expression (2).

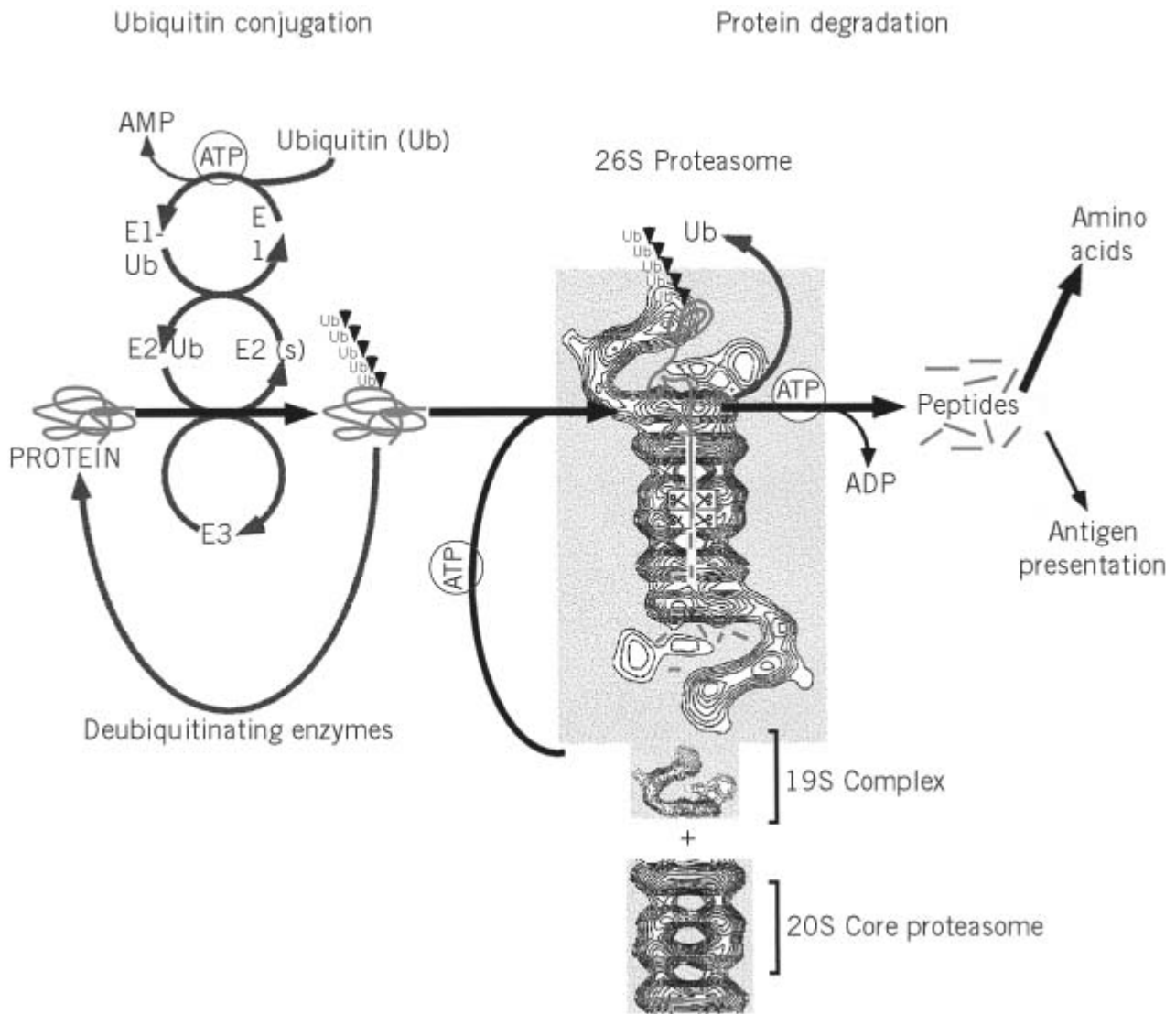
## Bibliography

1. D. Poccia (1986) *Int. Rev. Cytol.* **105**, 1–65.
2. K. C. Kleene (1996) *Mol. Reprod. Dev.* **43**, 268–281.

## Proteasome

Individual [proteins](#) in eukaryotes are marked for **protein degradation** by being covalently linked to [ubiquitin](#) (Ub). They are then digested **proteolytically** into small **peptides** by the 26 S proteasome particle (Fig. 1). The key proteolytic component of this large complex is the 600-kDa 20 S particle. It is a major cell constituent, comprising up to 1% of the cellular protein. It has a barrel-shaped **structure** composed of four stacked rings, each containing seven subunits surrounding a central cavity. The two inner b rings enclose a large central chamber containing the sites where proteins are degraded, and the two outer a rings surround a small opening through which protein substrates enter. Because of the opening's small diameter, polypeptides must first be unfolded to enter it.

**Figure 1.** Mechanism of ubiquitin-conjugated protein degradation by the proteasome. Key: E1, ubiquitin-activating protein; E2, ubiquitin carrier protein; E3, ubiquitin-protein ligase; Ub, ubiquitin.



Unlike traditional [proteinases](#), the proteasome does not simply cleave a protein and release the fragments; on the contrary, a protein substrate, once bound, is cut multiple times into small peptides ranging from 3–25 residues in length before the next protein molecule is attacked. This highly processive mechanism ensures that partially digested proteins do not accumulate within cells. Peptides released by the proteasome are rapidly hydrolyzed to [amino acids](#) in the cytosol, although some peptides are used in **major histocompatibility class (MHC) I antigen presentation** .

The presence of a powerful proteolytic enzyme system within the cell requires the evolution of mechanisms to safeguard against the nonspecific digestion of essential proteins. The first safeguard is the requirement for conjugation with ubiquitin before degradation of the protein. Second, the [active sites](#) of the proteasome are geographically isolated within the central chamber and away from the cytosolic milieu. Third, only unfolded proteins can enter the narrow opening in the rings on either end of the 20S proteasome.

On either end of the 20S proteasome, there is a large 19S (700 kDa) regulatory particle. This structure provides specificity to proteolysis by binding ubiquitinated substrates and catalyzing the entry of the polypeptide into the 20S particle. Interestingly, the 19S regulatory complex also contains at least six different **ATPases**. It is likely that the consumption of ATP by the 19S complex enables it to unfold protein substrates, to inject them into the 20S proteasome, and to activate the proteolytic activity of the particle. Thus, the marking of substrates by Ub conjugation, the organization of the

26S complex, and the energy requirement for proteolysis all appear to have evolved to provide a remarkable degree of selectivity and regulation to the degradative process.

The mechanism by which peptide bonds are cleaved in the proteasome is also unique. Proteasomes do not fit into the standard classification of proteolytic enzymes according to their active sites (eg, [serine proteinases](#), [thiol proteinases](#), [carboxyl proteinases](#), or [metalloproteinases](#)). The sequences of their  $\beta$  subunits are not **homologous** to those of known proteolytic enzymes, and the pattern of sensitivity to various inhibitors differs from that of any known proteinase family. The recent [X-ray crystallography](#) studies and [site-directed mutagenesis](#) of different amino acid residues in the proteasome have uncovered a new type of proteolytic mechanism. The active-site nucleophile of the proteasome is the hydroxyl group of a [threonine](#) residue at the amino terminus of the  $\beta$  subunit. This novel mechanism has permitted selective inhibitors of the proteasome to be synthesized.

The mammalian proteasome contains three types of active sites that differ in their specificity for different types of bonds. One site cleaves preferentially after basic amino acid residues, the others after large **hydrophobic** residues, small neutral residues, or acidic residues. These activities, which function together to catalyze the complete digestion of proteins, are associated with different  $\beta$  subunits. Intriguing adaptations exist in proteasome subunit composition that can aid in the generation of diverse antigenic peptides. For example, in disease states, the cytokine **g-interferon**, which enhances antigen presentation in most cells, induces the expression of three novel  $\beta$  subunits that are incorporated into the 20S proteasome in place of the  $\beta$  subunits normally present. These components alter the way in which peptide bonds are cleaved to favor the production of the kinds of peptides that are used preferentially for MHC class I antigen presentation.

#### 1. 20 S Proteasomes of Prokaryotes

Many [Gram-Positive Bacteria](#), like archaeobacteria, contain a 20S proteasome, resembling in size and shape that in eukaryotes but composed of only one type of a subunit and one type of  $\beta$  subunit. These 20 S proteasomes appear *in vivo* to function in protein breakdown in association with ATPase complexes called PAN, or Proteasome-Activating nucleotidase. This hexameric-ring structure is a member of the AAA family of ATPases and is homologous to the six ATPases found in the 19S regulator of the eukaryotic proteasome. Addition of ATP and PAN to the 20S proteasome allows the ATP-dependent degradation of proteins but does not influence peptide hydrolysis. Many features suggest that this PAN complex was the evolutionary precursor of the 19S component of the eukaryotic 26S proteasome that degrades ubiquitinated proteins in an ATP-dependent fashion.

#### Suggestions for Further Reading

W. Baumeister, J. Walz, F. Zuhl, and E. Seemuller (1998) The proteasome: Paradigm of a self-compartmentalizing protease. *Cell* **92**, 367–380.

A. Ciechanover (1994) The ubiquitin-proteasome proteolytic pathway. *Cell* **79**, 13–21.

W. E. Mitch and A. L. Goldberg (1996) Mechanisms of muscle wasting. The role of the ubiquitin-proteasome pathway. *New Engl. J. Med.* **335**, 1897–1905.

O. Coux, K. Tanaka, and A. L. Goldberg (1996) Structure and functions of the 20S and 26S proteasomes. *Ann. Rev. Biochem.* **65**, 801–847.

A. L. Goldberg and K. L. Rock (1992) Proteolysis, proteasomes and antigen presentation. *Nature* **357**, 375–379.

M. Groll, L. Ditzel, J. Lowe, D. Stock, M. Bochtler, H. D. Bartunik, and R. Huber (1997) Structure of 20S proteasome from yeast at 2.4-Ångstrom resolution. *Nature* **386**, 463–471.

E. Seemuller, A. Lupas, D. Stock, J. Lowe, R. Huber, and W. Baumeister (1995) Proteasome from *Thermoplasma acidophilum*—a threonine protease. *Science* **268**, 579–582.

## Protein A

Protein A is a cell-surface [protein](#) isolated from *Staphylococcus aureus* that binds mammalian **immunoglobulins**. The reactivity of an extract of *S. aureus* cells with all human sera was initially thought to reflect an innate human anti-staphylococcal immune specificity (1). Later studies showed that the binding reaction could be competed with Fc fragments of human IgG, rather than Fab fragments (see [Immunoglobulin Structure](#)). This finding pointed to a non-specific interaction of Protein A with the heavy-chain constant region of IgG, present in Fc, rather than a normal antibody–antigen recognition reaction mediated through the variable region that is present in Fab (2). A distinct immunoglobulin binding mode was subsequently discovered in which Protein A bound to antibody variable domains (3, 4).

Native protein A consists of a single [polypeptide chain](#) with five homologous ~60-residue **domains**, arranged from the *N* terminus in the order EDABC (5-7). An additional region following the C domain is used for covalent linkage to the peptidoglycan cell wall of *S. aureus* (8-11). The E domain has less homology to A, B, C, and D than the other four domains have with each other, and it was not identified in the initial protein sequence determination (7). Soluble protein A is released as an intact fragment by digestion of the peptidoglycan with the staphylolytic endopeptidase lysostaphin (12). Each domain is capable of binding Fc regions (6, 13). Structural studies of domain B and of Protein A showed that domain B is a three-helix bundle that binds near the junction of the immunoglobulin C<sub>H2</sub> and C<sub>H3</sub> domains and makes contacts with both immunoglobulin domains (14, 15).

The role of Protein A in staphylococcal infection remains speculative. The primary importance of Protein A in molecular biology is as an extremely versatile reagent for immunological studies. Protein A crosslinked to Sepharose or other supports is used as an affinity adsorbent for immunoglobulins (16). The serum of virtually all mammalian species has immunoglobulin classes that will bind Protein A, and other classes that do not, allowing rapid fractionation (17, 18). Within a species, the heavy-chain **isotypes** that bind protein A often do so with differential sensitivity to pH and may be separated quantitatively by pH step or gradient elution (19, 20).

A second use of Protein A is as an immunosorbent. To prepare this immunosorbent, *S. aureus* cells expressing surface Protein A are fixed with formalin, then heat-killed. Following this treatment, the Protein A remains attached to the cell wall and functional in Fc binding. Treated cells added to an antibody solution rapidly adsorb antibodies and immune complexes, and they retain these adsorbed molecules while the cells are repeatedly washed and pelleted by [centrifugation](#). This stability makes the fixed cells ideal for capture of antibody and antibody–antigen complexes in immunoassays (21). Little nonspecific adsorption of other assay components to the cells is observed; hence uncomplexed antigen is free in solution. Fixed staphylococci have also been used as an immunosorbent for isolating surface antigens from animal cells or cell lines (22). Cells expressing the antigen to be isolated are first **radiolabeled** biosynthetically or chemically. Cells are then lysed, usually concomitant with [detergent](#) solubilization of the antigen. A specific antiserum or [monoclonal antibody](#) is added, and time is allowed for antigen binding to occur. The Protein A immunosorbent is then added to trap antibody–antigen complexes, from which antigen can be isolated by denaturing gel [electrophoresis](#). The advantage of this technique over using a second antibody for immunoprecipitation of antibody–antigen complexes is that in the latter case precipitate formation is slow, seldom quantitative, and highly dependent on the ratio of the first and second antibodies.

Recombinant fusion proteins incorporating one or more domains from Protein A or one of the other receptors take advantage of their high-affinity interaction with immunoglobulins and have become

widely used in biotechnology (23). Such proteins are easily detected with immunological reagents and can be purified in a single step by [affinity chromatography](#) on immobilized IgG media.

Protein A is the most extensively used and studied **Gram-positive** bacterial immunoglobulin receptor, but it is not the only one. Other proteins with similar functions have been identified (24), and they often have complementary specificities. Protein G from staphylococci binds to all human IgG **isotypes**, interacting with both the Fc region and the heavy-chain C<sub>H1</sub> constant domain (25, 26). Protein G overlaps the protein A binding site on Fc, yet it has a very different three-dimensional structure. Protein L from *Peptococcus magnus* is structurally homologous to protein G (27), but it binds to conserved regions in a subset of kappa light-chain variable domains (28, 29).

## Bibliography

1. K. Jensen (1958) *Acta Pathol. Microbiol. Scand.* **44**, 421–428.
2. A. Forsgren and J. Sjöquist (1966) *J. Immunol.* **97**, 822–827.
3. M. Inganäs, S. G. O. Johansson, and H. H. Bennich (1980) *Scand. J. Immunol.* **12**, 23–31.
4. M. Inganäs (1981) *Scand. J. Immunol.* **13**, 343–352.
5. S. Löfdahl, B. Guss, M. Uhlén, L. Philipson, and M. Lindberg (1983) *Proc. Natl. Acad. Sci. USA* **80**, 697–701.
6. J. Sjödahl (1977) *Eur. J. Biochem.* **73**, 343–351.
7. J. Sjödahl (1977) *Eur. J. Biochem.* **78**, 471–490.
8. M. Uhlén, B. Guss, B. Nilsson, S. Gatenbeck, L. Philipson, and M. Lindberg (1984) *J. Biol. Chem.* **259**, 1695–1702.
9. O. Schneewind, P. Model, and V. A. Fischetti (1992) *Cell* **70**, 267–281.
10. O. Schneewind, A. Fowler, and K. F. Faull (1995) *Science* **268**, 103–106.
11. H. Ton-That, K. F. Faull, and O. Schneewind (1997) *J. Biol. Chem.* **272**, 22285–22292.
12. J. Sjöquist, B. Meloun, and H. Hjelm (1972) *Eur. J. Biochem.* **29**, 572–578.
13. T. Moks, L. Abrahmsén, B. Nilsson, U. Hellman, J. Sjöquist, and M. Uhlén (1986) *Eur. J. Biochem.* **156**, 637–643.
14. J. Deisenhofer (1981) *Biochemistry* **20**, 2361–2370.
15. H. Torigoe, I. Shimada, A. Saito, M. Sato, and Y. Arata (1990) *Biochemistry* **29**, 8787–8793.
16. H. Hjelm, K. Hjelm, and J. Sjöquist (1972) *FEBS Lett.* **28**, 73–76.
17. J. Goudswaard, J. A. van der Donk, A. Noordzij, R. H. van Dam, and J.-P. Vaerman (1978) *Scand. J. Immunol.* **8**, 21–28.
18. D. D. Richman, P. H. Cleveland, M. N. Oxman, and K. M. Johnson (1982) *J. Immunol.* **128**, 2300–2305.
19. P. L. Ey, S. J. Prowse, and C. R. Jenkin (1978) *Immunochemistry* **15**, 429–436.
20. R. C. Duhamel, P. H. Schur, K. Brendel, and E. Meezan (1979) *J. Immunol. Meth.* **31**, 211–217.
21. S. Jonsson and G. Kronvall (1974) *Eur. J. Immunol.* **4**, 29–33.
22. S. W. Kessler (1975) *J. Immunol.* **115**, 1617–1624.
23. B. Nilsson, L. Abrahmsén, and M. Uhlén (1985) *EMBO J.* **4**, 1075–1080.
24. G. Kronvall (1973) *J. Immunol.* **111**, 1401–1406.
25. L. Björck and G. Kronvall (1984) *J. Immunol.* **133**, 969–974.
26. J. P. Derrick and D. B. Wigley (1992) *Nature* **359**, 752–754.
27. M. Wikström, T. Drakenburg, S. Forsén, U. Sjöbring, and L. Björck (1994) *Biochemistry* **33**, 14011–14017.
28. E. B. Myhre and M. Erntall (1985) *Mol. Immunol.* **22**, 879–885.
29. L. Björck (1988) *J. Immunol.* **140**, 1194–1197.



### Suggestions for Further Reading

30. J. W. Goding (1978) Use of staphylococcal Protein A as an immunological reagent. *J. Immunol. Meth.* **20**, 241–253.
31. S. Ståhl, P.-Å. Nygren, A. Sjölander, and M. Uhlén (1993) Engineered bacterial receptors in immunology. *Curr. Opin. Immunol.* **5**, 272–277.

### Protein Blots (Western Blots)

SDS–Polyacrylamide gel electrophoresis (SDS–PAGE ) is one of the most versatile and important separation methods for analyzing [proteins](#). Hundreds of polypeptides can be resolved on a single gel by this process, and **two-dimensional electrophoresis** expands the possibility of separation to thousands and more. These techniques separate the components on the basis of their physical properties, not on their activity or ability to interact with a given probe, eg, [immunoglobulin](#), [lectin](#) or **ligand**. Western blotting (also called protein blotting ) extends [gel electrophoresis](#) of proteins to provide such information ([1](#), [2](#)).

The production of a protein blot involves transferring the protein pattern from the gel onto a blotting matrix (see [Blotting](#)). This is commonly achieved by electroelution, subjecting the slab gel to a second electric field perpendicular to the surface of the gel, so that the polypeptides are eluted onto a blotting matrix, such as a [nitrocellulose](#) membrane filter. Numerous factors affect the quality of the transfer, such as the molecular mass of the peptides, the type of gel, the composition of the transfer buffer, the type of matrix, and the transfer conditions (e.g., voltage/current and time) ([1-4](#)). The transfer apparatus also affects the quality of the blot. Generally, systems that produce uniform or regulated gradient electric fields should be employed ([5](#), [6](#)).

Before a blot is probed, it is stained for protein to reveal the pattern of all transferred proteins. Stains, such as [Ponceau S](#) or Amido black are useful. Then the blot must be quenched, i.e., all available unoccupied surface of the matrix must be blocked with an inert blocking agent . Commonly, bovine [serum albumin](#), gelatin, or milk ([casein](#)) have been used, and nonionic [detergents](#), such as Tween-20 or Triton X-100 are often included in the quenching, washing, and incubation solutions. Quenching minimizes nonspecific adsorption of the probe to the blot and reduces the irrelevant background to a minimum. Then the quenched blot is incubated with a probe to identify any protein with which it associates specifically (see [Blot Overlays](#)). For example, when **antibodies**, [lectins](#), **ligands**, or even cells are used as probes, immunoblots , lectin blots, ligand blots, and cell blots, respectively are carried out. After probing, the blot undergoes extensive washes to remove nonbound probe.

Identifying the complex between the probe and the protein bound to the surface of the blot achieved is by [autoradiography](#) when the probe is **radiolabeled**. Alternatively, a second probe can be used to identify the complex. For example, an enzyme conjugate of goat antimouse IgG might be used to reveal a murine [monoclonal antibody](#) bound to its corresponding antigen. [Avidin](#) or [streptavidin](#) enzyme conjugates are often required when biotinylated probes are used (see [Avidin-Biotin System](#)). Nonradioactive detection systems routinely employ the typical **ELISA** reagents that produce precipitating colored products. However, **luminescent** and **chemiluminescent** procedures are also sensitive and easy to use ([7](#), [8](#)).

In addition to detecting bimolecular interactions, protein blots have gained importance as an essential step in peptide sequence analysis. In such cases, a protein is transferred a gel by using a blotting

matrix that withstands the reactions of the [Edman Degradation](#). Polyvinyl difluoride (PVDF) membrane filter is the blotting matrix most commonly used for this application. ([9](#), [10](#)).

## Bibliography

1. J. M. Gershoni (1988) *Methods Biochem. Anal.* **33**, 1–58.
2. H. Towbin and J. Gordon (1984) *J. Immunol. Methods* **72**, 313–340.
3. D. E. Garfin and G. Bers (1989) In *Protein Blotting; methodology, research and diagnostic applications* (B. A. Baldo and E. R. Tovey, eds.), Karger, Basel, pp. 5–42.
4. A. De Maio (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 11–32.
5. M. Bittner, P. Kupferer, and C. F. Morris (1980) *Anal. Biochem.* **102**, 459–471.
6. J. M. Gershoni, F. E. Davis, and G. E. Palade (1985) *Anal. Biochem.* **144**, 32–40.
7. R. E. Geiger (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press Oxford, UK, pp. 133–138.
8. I. Durrant and S. Fowler (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 141–152.
9. M. A. Mansfield (1994) In *Protein Blotting: A Practical Approach* (B. S. Dunbar, ed.), IRL Press, Oxford, UK, pp. 33–52.
10. C. Eckerskorn and F. Lottspeich (1993) *Electrophoresis* **14**, 831–838.

## DNA Degradation *In Vivo*

Low levels of DNA degradation occur continuously in all cells from bacteria to man in conjunction with DNA repair processes (see [DNA Repair](#)). Probably 90% of this degradation can be attributed to excision repair during which only a few nucleotides are released from each damaged site by the actions of repair endo- and exonucleases and recycled in the overall process. High levels of DNA degradation occur within a population of cells in individual cells that need elimination from that population. This need occurs either because their DNA has become too heavily damaged to be successfully repaired, because the cells have been infected with bacteriophages or viruses, or, during development and tissue remodeling in multicellular organisms, because these cells have received signals instructing them to commit suicide. In higher eukaryotes, this cell suicide is known as apoptosis (see [Apoptosis](#)). In this case, the chromatin DNA is fragmented and packaged with other cellular materials into apoptotic bodies that are engulfed by adjacent healthy cells in the tissue and digested in their lysosomes.

Except for a small body of work on *Escherichia coli*, little is known about cell suicide in bacteria. When DNA bacteriophages that have not adapted to growth in *E. coli* are cleaved by host restriction nucleases, the resulting ds-DNA fragments are degraded to small oligo-nucleotides by a nuclease with  $Mg^{++}$ - and ATP-dependent endo- and exonuclease activities encoded by the *recB*, *recC*, and *recD* genes. This nuclease acts normally in recombination and in recombinational ds-break repair, a minor DNA repair pathway ([1](#)). Some bacteriophages that have adapted to *E. coli* contain genes that encode specific inhibitors of the *recBCD* nuclease. This nuclease is also responsible for degrading genomic DNA damaged beyond repair by UV light or ionizing radiation ([1](#)), but the details of the process are not known. The enzyme has homologs in many other species of bacteria.

Most chromatin DNA fragmentation in eukaryotic cells follow an ordered path that yields first

300 kbp double strand (ds) fragments, then 50 kbp ds-fragments, and, last, a range of small ds-fragments that are multiples of mononucleosome-sized DNA, 180 to 200 bp (2). This fragmentation pattern reflects the different levels of packing of the DNA in the chromatin. Six 50 kbp loops form 300 kbp rosettes that are stacked together. Cleavage between rosettes releases the 300 kbp fragments, and cleavage at the bases of the loops releases 50 kbp linear ds-DNA that is readily cleaved between nucleosomes to small fragments that have 3'-OH and 5'-P termini. When electrophoresed in agarose, the latter appear as a characteristic "ladder" of DNA. All stages of the fragmentation are  $Mg^{++}$ -dependent, but the subsequent cleavage to small fragments is also activated by  $Ca^{++}$ . This may indicate that there is more than one nuclease involved in chromatin DNA cleavage (2). In a few cases, the earlier stages of fragmentation are sufficient for apoptosis to proceed to completion.

Random DNA degradation occurs during cell necrosis as a result of very heavy damage to the cell or tissue injury and is accompanied by lysosome disruption and cytolysis. DNase II, an acid DNase normally sequestered in the lysosomes, will generate a ladder of ds-DNA fragments when incubated *in vitro* with isolated nuclei at acidic pH. However, only a small decrease in intracellular pH (0.3 pH U) occurs during apoptosis (3), not enough to result in appreciable activation of DNase II even if it had a nuclear location. DNase II is not metal ion-dependent and makes ds-breaks with 3'-P and 5'-OH groups, termini not found on the apoptotic ds-DNA fragments. The corresponding acid DNase of the flatworm, *Caenorhabditis elegans*, plays a role in apoptosis, namely in helping to digest the DNA in the apoptotic bodies engulfed by healthy cells. A loss of function of this nuclease caused by mutation of the *nuc 1* gene leads to accumulation of the apoptotic bodies in the lysosomes (4).

At least three nucleases have been identified as responsible for DNA degradation during apoptosis in higher eukaryotes: endo-exonuclease (EE), caspase-activated DNase (CAD), and the mitochondrial nuclease (mitNuc). These are present in cells either in inactive forms (EE and CAD) that are activated by proteases or sequestered (mitNuc). EE was first purified from *Neurospora crassa* (5). It is a  $Mg^{++}$ -dependent enzyme with homologs in other fungi and yeast, where it acts in recombination and in recombinational ds-break repair (5). It may be the eukaryotic counterpart of the bacterial recBCD nuclease. It has endonuclease activity with both DNA and RNA and exonuclease activity with DNA. EE is the major degradative nuclease in nuclei of human leukemia cells, present entirely in inactive form, mainly bound to the nuclear matrix (6). A nearly equal amount of inactive EE is also bound to the membranes of the endoplasmic reticulum. Proteolysis of EE has been detected in response to different apoptotic agents and yielded polypeptides identical in sizes to the various  $Ca^{++}$ -,  $Mg^{++}$ -endonucleases isolated previously from apoptotic cells by others (6). EE is activated *in vitro* by treatment with the serine protease trypsin or with caspase-3 (7). A mammalian EE, synergistically activated by  $Ca^{++}$ , has been isolated from monkey CV-1 cells (8). Genetic evidence is needed to further delineate the precise role of EE in apoptosis.

CAD is present in the cytosol and in nuclei of mammalian cells (9) in the form of a single inactive polypeptide called DNA fragmentation factor (DFF). Activation of CAD takes place when DFF is cleaved by one of the specialized and highly regulated proteases called caspases. This occurs *in vitro* with caspase-3 or caspase-7. CAD is a  $Mg^{++}$ -dependent specific DNase that cleaves the DNA of isolated nuclei to nucleosome-sized fragments. However, CAD is not present in some lower organisms such as yeast or the flatworm (10). In addition, DNA degradation occurs during apoptosis in cells from DFF-knockout mice (10). This indicates that CAD is not essential and that other nucleases also play a role in the process.

In cells of DFF-knockout mice, tumor necrosis factor causes high levels of nuclear DNA degradation. In this case, it is due to the release of the mitNuc from the inter-membrane space of the mitochondria (10). Yeast mitNuc was previously shown to be an endo-exonuclease and is closely related to the mammalian mitNucs (5). The release of mouse mitNuc is mediated by the activation of caspase-8. This protease cleaves a protein called Bid to a truncated form (tBid) that induces the release of several apoptotic factors including mitNuc (10). However, mitNuc is not released from

mitochondria in HL-60 cells during apoptosis induced by two other agents, the chemotherapeutic drug etoposide (VP-16) or the anti-rheumatic drug hydroxychloroquine ([11](#)). Thus, it is unlikely that mitNuc is involved in the DNA degradation caused by these apoptotic agents.

### Bibliography

1. G. R. Smith (1988) *Microbiol. Rev.* **52**, 1–28.
2. P. R. Walker, S. Pandey, and M. Sikorska (1995) *Cell Death Differ.* **2**, 93–100.
3. M. A. Barry and A. Eastman (1993) *Arch. Biochem. Biophys.* **300**, 440–450.
4. J. Hevelone and P. S. Hartman (1988) *Biochem. Genet.* **26**, 447–460.
5. M. J. Fraser and R. L. Low (1993) In *Nucleases*, 2 ed. (R. J. Roberts, S. M. Linn, and S. Lloyd, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
6. M. J. Fraser et al (1996) *J. Cell Sci.* **109**, 2343–2360.
7. X. W. Meng et al (2000) *Apoptosis* **5**, 243–254.
8. C. Couture and T. Y.-K. Chow (1992) *Nucleic Acids Res.* **20**, 1379–1385.
9. M. Enari et al (1998) *Nature* **391**, 43–50.
10. Y. Li et al (2001) *Nature* **412**, 95–99.
11. X. W. Meng, et al (1998) *Apoptosis* **3**, 395–406.

### Suggestions for Further Reading

12. M. J. Fraser (1996) *Endo-exonucleases: Actions in the life and death of cells*. R. G. Landes Co., Austin, Texas.

### Protein Detection

Many methods for detecting or measuring [proteins](#) in solution or on a [blotting](#) membrane have been developed. The Kjeldahl method detects all of the nitrogen atoms of proteins as ammonia after heating them with sulfuric acid. It is used only with very complex samples. The Lowry (see later), **biuret**, and [ninhydrin](#) assays are more widely used colorimetric methods. The UV method utilizes the **absorbance** of [tryptophan](#) and [tyrosine](#) residues at 280 nm and is a very convenient method for determining protein content ([1](#)). Detecting or determining proteins by the binding of dyes (see [Coomassie Brilliant Blue](#) and [Ponceau S](#)) and after reaction with a dye that becomes fluorescent (see [Fluorescamine](#)) are sensitive and convenient methods. Many of the methods mentioned previously can be employed to detect proteins on a membrane or in a polyacrylamide gel. The most sensitive detection method is [silver staining](#), which detects 2 to 5 ng of protein/band. For most of these assays, the responses of different proteins vary quantitatively, depending upon their amino acid compositions, which makes quantification difficult.

To determine or detect a particular protein, its biological function can be utilized (see [Overlay Assay: Enzyme Zymography of Plasminogen Activators and Inhibitors](#)). **Immunoblotting** is one of the most widely used methods of this type.

#### 1. Lowry Assay

The Lowry method ([2](#)) is a sensitive, colorimetric, protein determination method that measures 1 to 20 µg of protein. It is a combination of the phenol method and the biuret method. The Folin and

Ciocalteu phenol reagent contains sodium molybdate(VI), sodium tungstate(VI), and phosphoric acid, and it reacts with phenols. Phenol reagent is added to the protein sample after the **biuret** reaction, and the absorbance at 750 nm is measured.

### Bibliography

1. W.H. Van Esand J.H. Wisse (1963) *Anal. Biochem.* **6**, 135–143.
2. J.A. Smith (1995) In *Current Protocols in Molecular Biology*, John Wiley, New York, pp. 10.1.1.–3.

### Suggestions for Further Reading

3. T.M. DeSilva (1995) "Protein detection in gels using fixation", In *Current Protocols in Protein Science* (J.E. Coligan et al. eds.), Wiley, pp. 10.5.1–12.)
4. J.A. Ursitti (1995) "Protein detection in gels without fixation", In *Current Protocols in Protein Science* (J.E. Coligan et al. eds.), Wiley, pp. 10.6.1–8.
5. S. Harper and D. W. Speicher (1995) "Detection of proteins on blot membranes", In *Current Protocols in Protein Science* (J.E. Coligan et al. eds.), Wiley, pp. 10.8.1–7.

## Protein Disulfide Isomerase

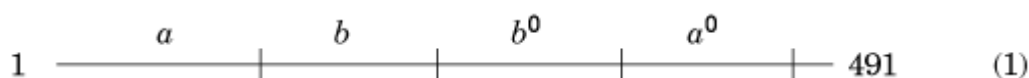
*Protein disulfide isomerase* (PDI) is the major catalyst in the [endoplasmic reticulum](#) of **eukaryotes** of **protein folding** that involves [disulfide bond](#) formation. It catalyzes the formation of protein disulfide bonds, their reduction, and their rearrangement, using domains that are **homologous** to [thioredoxin](#). Its catalytic activity results primarily from the somewhat unusual [thiol–disulfide exchange](#) properties of its thioredoxin-like active sites, plus a tendency to bind unfolded and misfolded polypeptide chains. In [Gram-Negative Bacteria](#), such as *Escherichia coli*, its multiple activities are catalyzed by several simpler proteins.

### 1. Eukaryotic PDI

PDI was the first catalyst of protein folding to be discovered, being detected independently by two groups in 1963 ([1](#), [2](#)) as part of the classic studies of C. B. Anfinsen and colleagues on the refolding of reduced **ribonuclease A** that accompanies disulfide bond formation. It was apparent that the oxidation of reduced ribonuclease was much too slow to account for its rate of biosynthesis *in vivo*, so a catalyst was sought; it was named PDI when it appeared that its prime function was to catalyze protein disulfide rearrangements, or isomerization. The importance of PDI in catalyzing disulfide formation and breakage was overlooked because those studies were using air oxidation to generate protein disulfide bonds. When a more physiological procedure was used, using reduced and oxidized [glutathione](#), GSH and GSSG, respectively, it became apparent that PDI catalyzes every step in the disulfide-linked folding of a protein ([3](#)). Glutathione is the primary sulfur reagent present within the endoplasmic reticulum, occurring in the millimolar concentration range. There the ratio of GSSG to GSH is between 1:1 to 1:3 ([4](#)), more oxidizing than the cytosol, where the ratio is closer to 1:100 (see [Disulfide Bonds](#)).

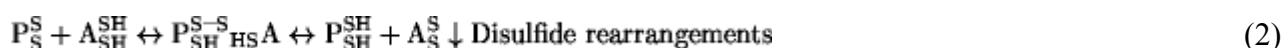
PDI occurs at high concentrations in the ER, approximately  $10^3\text{M}$ , where it comprises about 12% by weight of the total protein. At these concentrations, it is believed to be a dimer. It is a [calcium-binding protein](#), as are many proteins of the ER, although this has no known effect on its functional properties ([5](#)). The PDI polypeptide chain is also used as the b-subunit of prolyl-4-hydroxylase ([6](#))

and as one subunit of the heterodimeric microsomal triglyceride transfer protein (7); in both cases, it appears to play primarily a structural role. The primary structure of the mammalian PDI polypeptide chain (8) suggested that it is comprised of four domains, or modules, plus an acidic C-terminus that is believed to be involved in binding  $\text{Ca}^{2+}$  ions, and this has been confirmed by dissection studies (9):



The  $a$  and  $a'$  domains are clearly **homologous** to thioredoxin, with the typical pair of cysteine residues at the active site, and they have the same fold (10); surprisingly, the  $b$  domain also has the thioredoxin fold, but with the active-site cysteine residues deleted (11); the  $b'$  domain is homologous to  $b$  and therefore probably also has the thioredoxin fold.

Disulfide bond formation, breakage, and rearrangement in proteins are believed to occur by PDI reacting with their disulfide or thiol groups ( $\text{P}_S^S$  or  $\text{P}_{SH}^{SH}$ , respectively) using the pair of cysteine residues at the thioredoxin-like active sites of the PDI  $a$  and  $a'$  domains,  $\text{A}_{SH}^{SH}$  and  $\text{A}_S^S$ , through an intermediate mixed disulfide between the two:



The activities of the individual  $a$  and  $a'$  domains account quantitatively for the activity of PDI in inserting disulfide bonds into protein by the mechanism of Eq. 2 (12). The interconversion between the dithiol and disulfide forms of PDI might be accomplished by its reaction with the high concentrations of GSSG and GSH in the endoplasmic reticulum, or it could be enzyme-catalyzed (13). Protein disulfide isomerizations would occur if a different protein cysteine residue reacted with the mixed disulfide bond between protein and PDI to generate a new disulfide bond. Nevertheless, the individual  $a$  and  $a'$  domains have incomplete disulfide isomerization activity, which requires the presence of the  $b$  and  $b'$  domains.

The disulfide bonds at the PDI active sites are relatively unstable, so they are readily transferred to a substrate protein that can readily form disulfide bonds (12). The instability of the PDI disulfide bonds is an indirect effect that arises primarily because the accessible active site thiol groups of the  $a$  and  $a'$  domains have very low  $pK_a$  values; that for the  $a$  domain is 4.5 (14). Relative to a normal Cys residue thiol group, this thiolate anion is stabilized by 5.7 kcal/mol; most of this is believed to be due to interaction with the partial positive charges at the N-terminus of the  $\alpha$ -helix where it is located, with 1.1 kcal/mol being due to electrostatic interaction with a nearby His residue (14). The thiolate anion is the reactive form of a thiol group, so its stabilization in PDI is undoubtedly at least part of the reason why this thiol group is so reactive. Stabilization of the thiolate anion also stabilizes the reduced dithiol form of the protein over the disulfide form, so it indirectly destabilizes the disulfide bond.

In addition to the intrinsic reactivities of the thiol and disulfide groups at the active sites of PDI, weak binding interactions between the catalyst and a substrate protein are also believed to contribute to catalysis. The details of these phenomena are best understood in the analogous catalysts from *E. coli*.

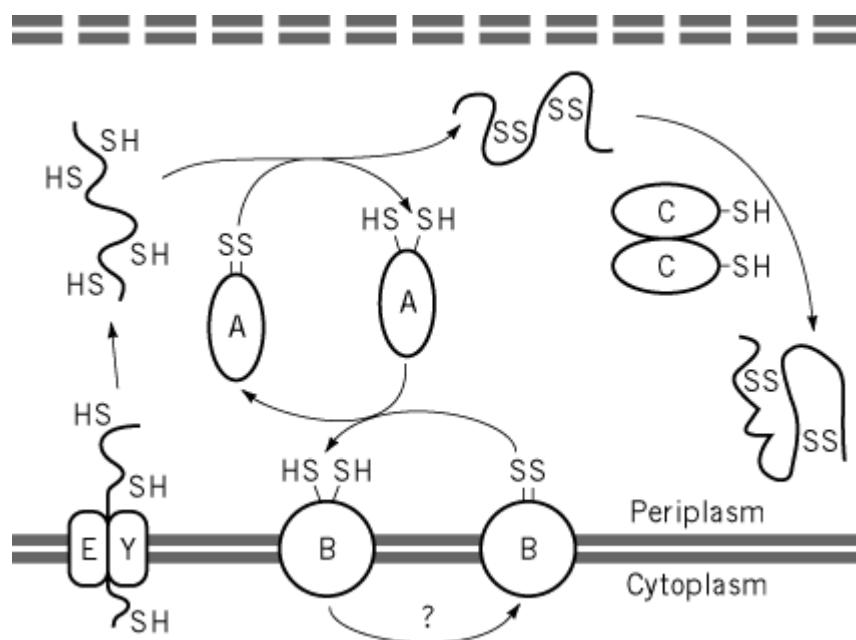
## 2. Gram-negative bacteria

Gram-negative bacteria have a **periplasmic** compartment between the inner and outer membranes,

where **secreted proteins** fold and incorporate disulfide bonds. Genetic analysis has identified a number of proteins that are required for this process (15-17).

Disulfide bonds are believed to be generated in newly secreted proteins by reaction with the unstable disulfide bond of DsbA, as in Eq. 2 (Fig. 1). This protein has an active site and a fold like that of thioredoxin, but with another domain inserted into the polypeptide chain (18). DsbA has a disulfide bond that is one of the least stable known (19), due to its reduced form being stabilized by **electrostatic interactions** with the thiolate anion form of the accessible cysteine thiol group, which has a greatly reduced  $pK_a$  of 3.5 (20). This disulfide bond is less stable than when the protein is unfolded, so it actually destabilizes the folded conformation of DsbA (19). Upon transfer to a newly synthesized protein (21), the DsbA disulfide bond is regenerated by the plasma membrane-bound DsbB protein, and may be coupled to the redox potential of the cytosol (17) (Fig. 1).

**Figure 1.** Present model for folding and disulfide formation in a newly secreted protein in the periplasm of Gram-negative bacteria. At lower left, a newly synthesized protein with four cysteine residues is being secreted through the plasma membrane using the *Sec* system. This protein incorporates two disulfide bonds by reacting with the disulfide form of two molecules of protein DsbA (labeled A). The disulfide form of DsbA is regenerated by the membrane protein DsbB. These disulfide bonds are isomerized by the dimeric DsbC to generate the native protein.



DsbA has very little protein disulfide isomerase activity, but this function appears to be catalyzed by a separate protein, DsbC. About half the polypeptide chain of this homodimeric protein is likely to have a folded conformation like thioredoxin (22). Its properties are very similar to those of DsbA, and why it is a better isomerase is not clear. It is likely that binding interactions with the substrate are crucial, and they are known to occur in both DsbA and DsbC (23).

Further genes involved in disulfide bond formation in the *E. coli* periplasm are being discovered and await characterization (24). In addition to the thioredoxin motif in all of these proteins, the additional domains are likely to be present to restrict the interactions between these various proteins and for their interaction with substrate proteins.

## Bibliography

1. R. F. Goldberg, C. J. Epstein, and C. B. Anfinsen (1963) *J. Biol. Chem.* **238**, 628–635.

2. P. Venetianer and F. B. Straub (1963) *Biochim. Biophys. Acta* **67**, 166–168.
3. T. E. Creighton, D. A. Hillson, and R. F. Freedman (1980) *J. Mol. Biol.* **142**, 43–62.
4. C. Hwang, A. J. Sinskey, and H. F. Lodish (1992) *Science* **257**, 1496–1502.
5. A. Zapun, T. E. Creighton, P. J. E. Rowling, and R. B. Freedman (1992) *Proteins Struct. Funct. Genet.* **14**, 10–15.
6. T. Pihlajaniemi, T. Helaakoski, K. Tasanen, R. Myllylä, M.-L. Huhtala, J. Koivu, and K. I. Kivirikko (1987) *EMBO J.* **6**, 643–649.
7. J. R. Wetterau, K. A. Combs, S. N. Spinner, and B. J. Joiner (1990) *J. Biol. Chem.* **265**, 9800–9807.
8. J. C. Edman, L. Ellis, R. W. Blacher, R. A. Roth, and W. J. Rutter (1985) *Nature* **317**, 267–270.
9. N. J. Darby, J. Kemmink, and T. E. Creighton (1996) *Biochemistry* **35**, 10517–10528.
10. J. Kemmink, N. J. Darby, K. Dijkstra, M. Nilges, and T. E. Creighton (1996) *Biochemistry* **35**, 7684–7691.
11. J. Kemmink, N. J. Darby, K. Dijkstra, M. Nilges, and T. E. Creighton (1997) *Curr. Biol.* **7**, 239–245.
12. N. J. Darby and T. E. Creighton (1995) *Biochemistry* **34**, 16770–16780.
13. J. Lundstrom and A. Holmgren (1990) *J. Biol. Chem.* **265**, 9114–9120.
14. T. Kortemme, N. J. Darby, and T. E. Creighton (1996) *Biochemistry* **35**, 14503–14511.
15. J. C. A. Bardwell, K. McGovern, and J. Beckwith (1991) *Cell* **67**, 581–590.
16. D. Missiakas, C. Georgopoulos, and S. Raina (1994) *EMBO J.* **13**, 2103–2020.
17. A. Rietsch, D. Belin, N. Martin, and J. Beckwith (1996) *Proc. Natl. Acad. Sci. USA* **93**, 13048–13053.
18. J. L. Martin, J. C. A. Bardwell, and J. Kuriyan (1993) *Nature* **365**, 464–468.
19. A. Zapun, J. C. A. Bardwell, and T. E. Creighton (1993) *Biochemistry* **32**, 5083–5092.
20. J. W. Nelson and T. E. Creighton (1994) *Biochemistry* **33**, 5974–5983.
21. N. J. Darby and T. E. Creighton (1995) *Biochemistry* **34**, 3576–3587.
22. D. Frishman (1996) *Biochem. Biophys. Res. Comm.* **219**, 686–689.
23. N. J. Darby, S. Raina, and T. E. Creighton (1998) *Biochemistry* **37**, 783–791.
24. D. Missiakas and S. Raina (1997) *Ann. Rev. Microbiol.* **51**, 179–202.

### **Suggestions for Further Reading**

25. R. B. Freedman (1989) Protein disulfide isomerase: Multiple roles in the modification of nascent secretory proteins. *Cell* **57**, 1069–1072.
26. R. Noiva and W. J. Lennarz (1992) Protein disulfide isomerase. A multifunctional protein resident in the lumen of the endoplasmic reticulum. *J. Biol. Chem.* **267**, 3553–3556.
27. N. J. Bulleid (1993) Protein-disulfide-isomerase: Role in biosynthesis of secretory proteins. *Adv. Protein Chem.* **44**, 125–150.
28. J. L. Martin (1995) Thioredoxin—a fold for all reasons. *Structure* **3**, 245–250.

### **Protein Engineering**

Protein engineering is a multidisciplinary technology for the design and production of new [proteins](#)



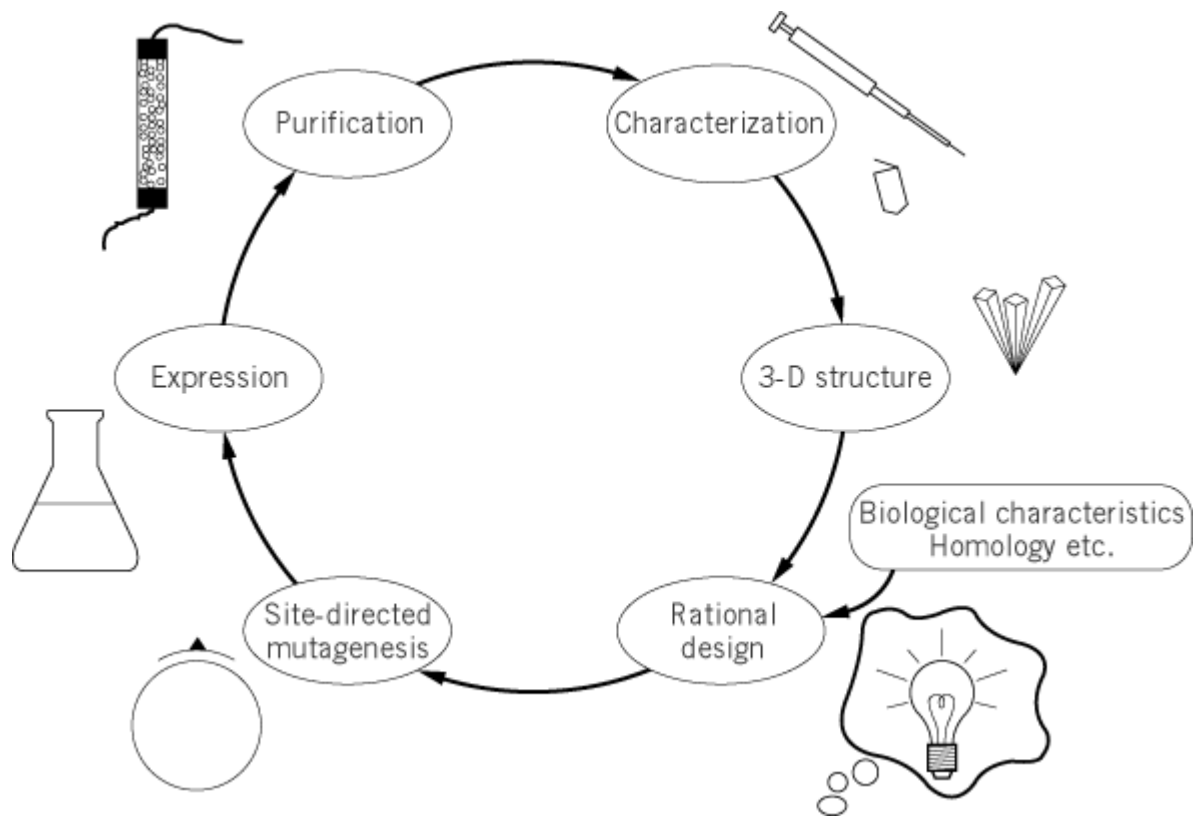
with predetermined properties. There are two basic motivations for pursuing protein design. The first is academic, aimed at increasing our knowledge about the fundamental properties of natural proteins, such as their folding pathways, stability (see [Protein Stability](#)) and catalytic properties (see [Enzymes](#)). Studies of the first type are often referred to as structure-function studies because they aim at establishing the relationship between how local parts of a protein's structure (individual amino acid residues or specific segments) contribute to its net functional behavior. The second is a practical motivation to be able to construct “tailor-made” proteins with improved properties to fulfill specific tasks in industry, medicine, agriculture, and so forth. Because the conditions under which industrially important proteins are used often differ from those of their natural environment, stability is a necessary consideration for most engineered proteins. An increase in catalytic activity is another goal of great interest. The ultimate ambition of the protein engineer is to be able to create a protein with any desired structure and function.

Three strategies for design dominate the field. The oldest strategy, “classical” protein engineering, is based on the reengineering of already existing proteins, whereas the second, *de novo* design, uses nature's rules to design entirely new proteins. These two are rational approaches by which the scientist tries to predict which sequence is needed to obtain a certain structure and function of interest. The third category of methods encompasses strategies mimicking natural [evolution](#), making these approaches more “irrational.”

## 1. Classical Protein Engineering

Classical protein engineering follows a cyclic diagram (Fig. 1), which is reiterative and can, in principle, be repeated infinitely. It is desirable to have considerable information about both the structural and functional aspects of the system to be studied. Structural information derives from **X-ray crystallographic** or NMR (nuclear magnetic resonance) studies, two techniques that complement each other. Functional information can derive from biochemical analysis, such as enzymatic assays, [affinity labeling](#), chemical modification, and [crosslinking](#) and includes knowledge about any sites of [post-translational modifications](#). This structural and functional knowledge of the protein forms the basis for deciding which amino acid substitutions should be introduced by [site-directed mutagenesis](#) in order to solve the present problem.

**Figure 1.** The classical protein engineering cycle uses the structure of the protein of interest as the starting point for designing mutants to explore a given problem. The mutations are introduced into the protein's gene using site-directed mutagenesis, and the mutant protein is expressed in a suitable expression system. Thereafter, the resulting protein is purified and characterized with respect to alterations in its function. The functionality profile is interpreted in structural terms and, ultimately, the structure of the mutant protein is determined in order to obtain a complete understanding of its functional behavior.



Gradually, computer tools are being developed to assist the protein engineer. These range from **homology modeling** to predict the structure of a related protein and the expected impact of a [mutation](#) on its structure or function to programs that will dress a structure of interest with the optimized sequence. There are two fundamentally different ways of computer modeling: a subjective, intuitive, and interactive approach and an objective, analytic approach based on energetic criteria and other mathematically formulated biophysical concepts.

The design of point mutants should focus on which amino acid residue to mutate and what to replace it with. A guideline for the identification of important residues is found in alignments of the protein of interest with related proteins (1) (see [Sequence Analysis](#)). Residues that are directly involved in a common function, such as [catalysis](#), will be among the most conserved; and residues that are important for structural reasons will also show a high degree of conservation. Consequently, mutations in residues that are required for structure formation or stability can also have dramatic effects on activity, leading to misinterpretations of the role of these residues. Conversely, the nonconserved residues are likely to be involved in functions unique to each protein, such as substrate binding and specificity, or these residues may have resulted from [neutral mutations](#) and may be dispensable for function.

Two sources of information can be used in attempts to determine which new amino acid should be introduced. Natural evolution has used *in vivo* mutations to engineer stable protein structures for specific functions. A similar overall folding pattern is maintained for functionally related molecules, irrespective of the accumulation of natural amino acid substitutions. Considerable variation can be tolerated in external loops, while the protein cores are maintained as well-packed regions with topologically equivalent side chains in each of the species. Evolutionary studies suggest that the interior of a protein can tolerate a volume change of no more than one methyl group without a change in the polar/[nonpolar](#) nature of the mutated residue. In contrast, residues near the surface can be permitted drastic substitutions, providing there is enough exposed surface for the side chain to reach the aqueous solvent (see [Accessible Surface](#)). Overall, however, the surface should appear

hydrophilic whereas the core is stabilized by **hydrophobic interactions** (2). Another source of information comes from other protein engineering studies, but the rules for “safe” mutations deduced from these two types of studies are more or less the same.

Two possible kinds of alterations can be introduced into a protein (3): (1) substitution of one amino acid by an isosteric residue of different function (eg, Glu by Gln), or (2) replacement of one amino acid by another of the same nature but different structure (eg, Glu by Asp). The first kind of substitution attempts to probe function while keeping structure constant, whereas the second kind enables the study of how function depends on structure. It is generally advisable to refrain from the mutation and introduction of **glycine** and **proline residues**, because these residues often have a great impact on the structure of the protein.

One should be very cautious when designing mutants because a badly designed mutant will produce more questions than answers. Mutations are easy to make, but the results of the subsequent functional characterization can turn out to be laborious, or even impossible, to evaluate in a sensible manner.

Apart from point mutations, the introduction of insertions or deletions (**indels**) can turn out to be a valuable tool for the engineering of proteins (4). Indels often give rise to structural alterations that could not have been obtained using substitutions alone. Nature seems to have used indels as a way of optimizing the position of side chains for a specific biological function. So far, protein engineers have used indels mainly for the mapping of protein structure and function, but the introduction of indels should provide a new tool for engineering changes in protein activity. The consequences of this type of mutation are very difficult to predict, making a random approach combined with a screening/selection strategy the most feasible procedure (see below). Another mutagenesis strategy that is also being used for the identification of functionally important regions is the “charged-to-alanine scanning mutagenesis.” This technique involves the progressive replacement of each of the charged amino acids in the protein, ie, Asp, Glu, Arg, Lys, and His, with Ala. This primary screen should be followed by the creation of additional mutations to assess the contributions of the neighboring polar and neutral amino acids in the relevant regions (5).

The mutations used in classical protein engineering are most often introduced into the corresponding **gene** using **site-directed mutagenesis**, and an extensive number of methods are available for introducing mutations *in vitro* (6). **PCR** has also proved to be a powerful technique for the generation of both site-specific and random mutations.

The engineered gene must be expressed using an **expression system**. Most often, a simple and fast expression system will be the first choice for the exploration of the effects of a series of point mutations. The use of a homologous expression system can be an advantage because the mutant proteins have a better chance of being processed correctly in their native environment, leading to easier folding and more stable expression. It is beneficial to choose a generally applicable expression strategy, such as a **fusion gene** approach, to avoid the need for individual optimization for each mutant. Often, a few milligrams of the mutant protein will provide enough material for the subsequent functional characterization. If a more thorough structural analysis is desired, larger amounts will be needed, and the generation of microheterogeneity in terms of folding or addition of post-translational modifications will be a concern. Finally, if the mutant turns out to be of practical/commercial value, more work needs to be put into the development of the most stable, efficient, and safe expression system, taking into account the intended use of the engineered protein.

After expression, the mutant protein must be purified. Once again, a fusion gene approach will minimize the need for time-consuming optimizations and, often, the protein of interest can be obtained in pure form in one **chromatography** step. One further advantage is that the mutant protein is easy to separate from the **wild-type** protein, which may be a problem if a homologous expression strategy is applied. For some applications, the fusion partner needs to be cleaved off in a subsequent step but, frequently, it will not affect the function of the mutant seriously, enabling the functional

characterization to be carried out directly on the fusion protein. The wild-type protein should be expressed and purified in exactly the same manner as the mutant protein, to serve as a reference in the functional characterization.

The purified protein can be characterized in various ways to determine if it meets one's requirements. First of all, one should make certain that the protein actually contains the mutation. This can be done by [mass spectrometry](#) analysis or peptide sequencing. It is most often assumed, however, that the mutant gene remains unaltered in the expression host. The use of a **recombination**-deficient host increases the likelihood of this. The next task is to ensure that the overall structure of the mutant protein is preserved. In some cases, the solubility and stable expression of a mutant protein may be taken as evidence of its structural integrity, given that mutations disturbing the overall structure of a protein often result in **denaturation**, insolubility, and **proteolytic** degradation. The most exact and detailed methods for obtaining structural information are, of course, X-ray crystallography and NMR, but both methods are time-consuming, and faster methods are available that usually will supply sufficient information about the preservation of structure, such as measurement of the thermostability of the mutant protein using, for example, **calorimetry** or the application of various kinds of [spectroscopy](#) (see [Circular Dichroism](#); [Vibrational Spectroscopy](#)). If no global or long-range structural distortions are apparent, one can proceed to the functional assays that will investigate the role of the substituted amino acid in terms of catalytic activity, involvement in **ligand binding**, and so forth. The ultimate characterization of a mutant will be the determination of its three-dimensional structure, which allows a complete evaluation of the functional characterization in structural terms. Furthermore, it is the starting point for a new round of improvements in the protein engineering cycle.

Traditional site-directed mutagenesis provides the possibility of substituting one amino acid residue with each of the other 19 natural amino acids. This is a fairly limited selection, and the range narrows when one considers that some of these substitutions will obviously give rise to unplanned structural changes. Methods are available, however, that permit the insertion of unnatural amino acids into a protein. The principle is that a purified and chemically mischarged [suppressor tRNA](#) is added to an *in vitro* transcription/translation system containing the gene of interest with the cognate **nonsense codon** at the position to be mutated (7, 8).

As an example, the engineering of increased stability is often a requirement for proteins with industrial applications. One of the most obvious sources of information is a comparison of sequences and structures of proteins from a **thermophilic** origin with those of their **mesophilic** counterparts (9). A number of "rules" for achieving thermostability can be set up on this basis: (1) certain residues tend to be circumvented (eg, Cys, Asn, and Gln); (2) thermostable proteins often have compact structures, and internal cavities are avoided; (3) the sizes of loops connecting **secondary structure** elements are reduced; and (4) interfaces are designed to be hydrophobic. Most often, several mutations need to be introduced into a thermolabile protein in order to increase its stability, which is a rather tedious task. Protein minimization (10), in which proteins are decreased in size via a combination of design, selection, and screening, is another approach to increasing the stability of a protein.

The engineering of proteins using site-directed mutagenesis can be rather complicated. The simplest approach for building new proteins has been transferring entire **domains** from one protein to another (11). Functional sites can also be transferred from one protein to another, with conservation of structural integrity and gain in function. The maintenance of structure, however, requires that the three-dimensional structure of the protein used as a scaffold be consulted.

## 2. Strategies for Design of Structure *de Novo*

The problem to be solved when designing a novel protein structure is to create an amino acid sequence that will fold into a stable and unique three-dimensional structure. A polypeptide chain can, in principle, take up an astronomically large number of conformations. Therefore, the **entropic** cost

of fixing the chain in a unique conformation is rather high. The task is to design a fold with a number of favorable interactions whose formations are associated with a gain in **free energy** that exceeds the decrease in conformational entropy.

The design process often occurs in a series of steps representing the hierarchy of forces required for stabilizing [tertiary structures](#), starting with hydrophobic forces and adding more specific interactions as required to obtain a unique functional protein (12). The simplest strategy of *de novo* design is the modular approach, in which fragments forming units of **secondary structure** are assembled into oligomeric structures. The driving force for this process is most often the burial of hydrophobic side chains in the center of the assembly. The modules are easy to synthesize, but the practical applicability of the approach is limited. The modular approach has been most successful for the assembly of **alpha-helices** into bundles, because  $\alpha$ -helices are stabilized by intramolecular [hydrogen bonds](#) that make them rather stable as single units. In contrast, **beta-strands** need to be linked together by interstrand hydrogen bonds to form stable assemblies; a  $\beta$ -strand is not stable individually. There are some general rules for the design of secondary structure building blocks for the construction of protein bodies. The modules need to be amphiphilic, and the presence of hydrophobic and [hydrophilic](#) residues needs to follow the regular pattern of the secondary structure elements in question. The assembly of modules can be facilitated by decreasing the entropic freedom of the system by linking its individual elements covalently. Several kinds of linkers, such as peptide loops, ligand-binding sites, covalent [crosslinking](#), or synthetic templates have been used successfully.

The next step in *de novo* design is to create single polypeptide chains that fold into a defined three-dimensional structure without the need for any linkers. The structure of interest needs to be stabilized by the introduction of appropriate interactions. Yet one also needs to design against unwanted, alternative structures by destabilizing such structures. Interactions stabilizing a protein structure are of many types, eg, [van der Waals interactions](#), [hydrogen bonds](#), [salt bridges](#), and **hydrophobic interactions**. Furthermore, the amino acids have certain propensities for existing in different kinds of secondary structure elements. Nevertheless, two rules that govern protein design in nature seem sufficient when proteins are designed *de novo*: (1) soluble proteins fold to maximize the burial of hydrophobic residues and exposure of hydrophilic residues, and (2) proteins are composed of building blocks of secondary structure elements, which are stabilized by a repeated hydrogen bonding pattern.

The ultimate ambition in *de novo* design is, of course, to design proteins that not only take up a predetermined structure but are also equipped with a specified function. The *de novo* design of function has led to proteins and peptides with specified properties for binding and catalysis. Synthetic [membrane proteins](#), **ion channels**, and new polypeptide materials are other examples of *de novo* design of function.

### 3. “Irrational Design” of Proteins via Directed Evolution

Rational design of proteins is hampered by how little we understand how precisely the secondary and tertiary structures of a protein are related to the amino acid sequence and how structure and function are related; but we are now learning how to mimic nature's way of designing new proteins. In nature, proteins are built in an “irrational” manner by the force of [evolution](#) driving the process by [mutation](#) and [natural selection](#). A library of the protein of interest containing randomly introduced mutations is subjected to screening or selection to find variants with new biological functions (examples can be found in references (13-15)). Directed evolution (also known as *in vitro* evolution) has an advantage over rational design in that it bypasses the need for understanding structure–function relations. This approach will also benefit from structural knowledge, however, as only a very small fraction of all possible protein sequences are available to analysis because of the practical limits of libraries. A library containing  $10^{12}$  different random peptide motifs enables the display of only ten residues in all possible combinations of the 20 amino acids. It is therefore very important to

be able to target the sequence randomization to critical areas. The terms *screening* and *selection* are often confused in the literature, even though they are fundamentally different techniques. Screening conditions apply if all members of a library are present when one chooses clones for further analysis. Screening methods include simple visual inspection and calorimetric assays (16). Under selection conditions, in contrast, only those clones of potential interest appear, which aids the sampling of much larger libraries. Selection techniques include genetic selection, **phage display**, the peptide-on-plasmid-technique, **polysome display**, and the **two-hybrid system** (17). The first and last methods represent *in vivo* techniques, whereas the selections of the others take place *in vitro*.

The generation of **catalytic antibodies** is a special way of mimicking nature, and the process can be carried out both *in vitro*, as described above, or *in vivo*, taking advantage of the immune system of the immunized animal, which possesses tremendous powers of randomization. A catalytic antibody is an antibody that has been raised against a **transition state analog** (18). The resulting antibody will therefore have a much stronger affinity for the transition state than for the substrate or product, which is how an enzyme provides a more favorable energetic route from substrate to product. The difficulty lies in choosing the best analogue of the transition state because the transition state itself cannot be synthesized.

Protein engineering is a multidisciplinary technology involving computer technologists, X-ray crystallographers, molecular biologists, fermentation experts, biochemists, and other specialists. Protein engineering will face many challenges before the ultimate goal of designing a protein purposely endowed with a specific structure and function is reached. One of the most serious challenges is probably to induce these scientists representing different disciplines to cooperate effectively.

#### Bibliography

1. J. U. Bowie, J. F. Reidhaar-Olson, W. A. Lim, and R. T. Sauer (1990) *Science* **247**, 1306–1310.
2. D. Bordo and P. Argos (1991) *J. Mol. Biol.* **217**, 721–729.
3. J. R. Knowles (1987) *Science* **236**, 1252–1258.
4. D. Shortle and J. Sodek (1995) *Curr. Opin. Biotech.* **6**, 387–393.
5. C. S. Gibbs and M. J. Zoller (1991) *Methods* **3**, 165–173.
6. M. K. Trower (ed.) (1996) *In Vitro Mutagenesis Protocols*, *Methods in Molecular Biology* (J. M. Walker, series ed.), vol. **57**, Humana Press, Totowa, N.J.
7. M. Ibba (1995) *Biotech. Gen. Eng. Rev.* **13**, 197–216.
8. D. Mendel, V. W. Cornish, and P. G. Schultz (1995) *Annu. Rev. Biophys. Biomol. Struct.* **24**, 435–462.
9. R. J. M. Russell and G. L. Taylor (1995) *Curr. Opin. Biotech.* **6**, 370–374.
10. B. C. Cunningham and J. A. Wells (1997) *Curr. Opin. Struct. Biol.* **7**, 457–462.
11. L. Bülow (1990) *Biochem. Soc. Symp.* **57**, 123–133.
12. J. W. Bryson et al. (1995) *Science* **270**, 935–941.
13. W. P. C. Stemmer (1994) *Nature* **370**, 389–391.
14. A. Cramer, E. A. Whitehorn, E. Tate, and W. P. Stemmer (1996) *Nature Biotechnol.* **14**, 315–318.
15. J. C. Moore and F. H. Arnold (1996) *Nature Biotechnol.* **14**, 458–466.
16. H. Zhao and F. H. Arnold (1997) *Curr. Opin. Struct. Biol.* **7**, 480–485.
17. J. W. Smith and E. Ruoslahti (1997) *Biotechnology and Genetic Engineering Reviews* **14**, 51–65.
18. D. B. Smithrud and S. J. Benkovic (1997) *Curr. Opin. Biotech.* **8**, 459–466.

#### Suggestions for Further Reading

19. P. R. Carey (ed.) (1996) *Protein Engineering and Design*, Academic Press, San Diego, California. Discusses design strategies, computer modeling, heterologous expression, in vitro mutagenesis, proteins crystallography, and spectroscopic and calorimetric methods for characterization. A few examples of applications are included.
20. D. L. Oxender and C. F. Fox (eds.) (1987) *Protein Engineering, Tutorials in Molecular and Cell Biology*, Alan R. Liss, Inc., New York. A classic with many good examples.
21. A. R. Rees, M. J. E. Sternberg, and R. Wetzel (eds.) (1993) *Protein Engineering: A Practical Approach*, IRL Press, Oxford, U.K. Covers all aspects of classical protein engineering and touches the engineering of antibody combining sites as well as phage display.
22. P. Wrede and G. Schneider (1994) *Concepts in Protein Engineering and Design: An Introduction*, Walter De Gruyter. Addresses mainly novices in the field of protein engineering.
23. In the journal *Current Opinion in Biotechnology*, one issue per volume (usually issue 4) is dedicated to reviewing the latest advances in the field of protein engineering, enabling the reader to become fully updated.

## Protein Evolution

The of can be studied at the levels of **nucleotide sequences, primary structures, secondary structures**, and . A nucleotide substitution within a **gene** coding for a protein can be classified as either synonymous substitution or nonsynonymous substitution (1). A synonymous substitution is the so-called silent substitution that never causes amino acid changes. On the other hand, nonsynonymous substitutions alter the amino acid sequence. The ratio ( $f$ ) of the number of nonsynonymous substitutions ( $K_A$ ) over that of synonymous substitutions ( $K_S$ ) can be used as a measure of the functional constraints on a protein. It has constantly been found that the rate of amino acid substitution in the functionally more important parts of the protein is slower than in the less important parts. With the assumption that  $K_S$  is a more direct reflection of the mutation rate, because there are no changes of amino acids (although this assumption is not always valid),  $f$  is an indicator of the degree of functional constraints for the whole, as well as parts of a protein.

For amino acid sequences, the rate of amino acid substitution is estimated as a measure of an evolutionary rate of change of proteins (for more details, see ). Regarding secondary structures, most functional protein **domains** are composed of a certain combination of **alpha-helix, beta-sheet** structures, , and irregular parts. One of the strong driving forces of protein evolution should be , in which those secondary structures may play an important role. and have elucidated a large number of three-dimensional . Moreover, methods for the prediction of tertiary structures of proteins from the amino acid sequences have improved tremendously (2). This provides a unique opportunity for us to make throughout comparisons of tertiary structures among different proteins, possibly leading to the elucidation of the ancient evolution of proteins and genes.

### Bibliography

“Protein Evolution” in , Vol. 3, p. 1990, by T. Gojobori; “Protein Evolution” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. M. Nei and T. Gojobori, (1986) *Mol. Biol. Evol.* **3**, 418–426.
2. J. U. Bowie, R. Luthy, and D. Eisenberg, (1991) *Science* **253**, 164–170.

## Protein Folding *In Vitro*

Proteins are biologically active only after they have adopted their native, three-dimensional folded **conformations** (see [Protein Structure](#)). Yet the genetic information used in **protein biosynthesis** specifies only the linear sequence of amino acid residues, the [primary structure](#). Natural proteins can often be **denatured**, or unfolded, and then **renatured**, or refolded, to the original conformation. When they cannot, it is usually because the unfolded protein has precipitated, aggregated, or been subjected to covalent modification. Therefore, the information for the **secondary structure** and the [tertiary structure](#) resides in the primary structure. Consequently, it should be possible to predict the three-dimensional structure of a protein from just its primary structure, if the process of protein folding were understood. Moreover, many proteins produced for pharmaceutical or industrial uses are generated initially in insoluble, unfolded, and inactive forms in [inclusion bodies](#), and they must be folded before they can be used.

Protein folding is simply a conformational change, an isomerization (unless [disulfide bonds](#) are formed), but its unusual aspect is the enormous number of conformations that an unfolded protein can adopt. If each amino acid residue adopts an average of  $j$  conformations, a polypeptide chain with  $(N+1)$  residues ( $N$  peptide units between them to define the conformation) could adopt up to  $j^N$  different conformations. The value of  $j$  is believed to be approximately 8, so a relatively small polypeptide chain of 100 amino acid residues should be able to sample some  $10^{89}$  different conformations. If the rate constant for an unfolded conformation to change is  $k_v$ , the average time to sample all of these conformations is given by

$$\tau = (Nk_v)^{-1} j^N \quad (1)$$

Unfolded conformations cannot change more rapidly than  $10^{13}$  times per second (and probably do so some  $10^9$  to  $10^{10}$  times per second), so it would require, on average, more than  $10^{66}$  years to sample  $10^{89}$  conformations. Even if there were only two conformations possible per residue, there would still be  $10^{30}$  conformations for  $N = 100$ , and  $10^7$  years would be required for random searching. Nevertheless, many proteins refold *in vitro* within seconds or minutes, some within a millisecond. Clearly, protein folding does not occur by a random searching of all possible conformations to find the unique native conformation, and there are likely to be pathways of folding.

A further complication is that unfolded proteins under denaturing conditions (eg, 8 M [urea](#) or 6 M **guanidinium chloride** (GdmCl)) approximate [random coils](#), so that each of the  $10^{15}$  to  $10^{18}$  molecules in a typical experimental sample would be expected to have a different conformation at each instant (and a slightly different one some  $10^{-10}$  seconds later). Therefore, each molecule is initiating folding from a different starting point.

### 1. Refolding of Small, Single-Domain Proteins

The native conformational states of proteins may often be unfolded reversibly by adding **denaturants**, increasing or decreasing the temperature, varying the pH, applying high pressures, or cleaving disulfide bonds (see [Denaturation, Protein](#)). At equilibrium, the unfolding transitions of single-domain proteins are usually two-state, and only the fully folded, native (N) and unfolded (U) states are populated. In this case, unfolding of the native conformation is cooperative, and partly folded molecules are unstable relative to the U or N states under all conditions. The most common

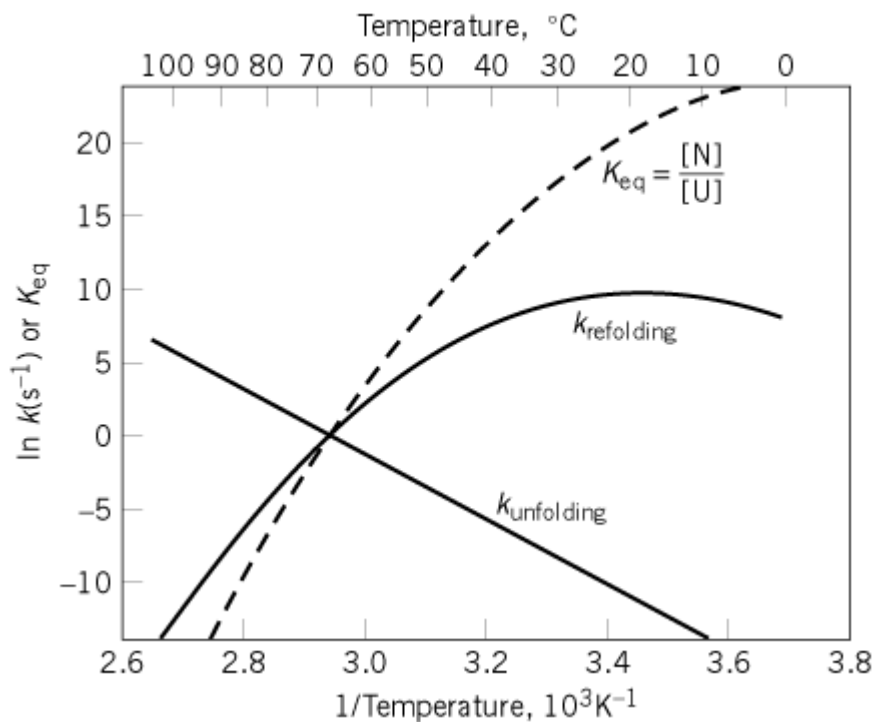


exception is the [molten globule](#) conformation, which predominates under intermediate denaturing conditions with some proteins.

Experimental *in vitro* studies of protein refolding kinetics generally start with the fully unfolded protein under unfolding conditions and, to initiate refolding, change the conditions abruptly to favor the folded state. Then changes in the conformational properties of the protein molecule population are monitored as a function of time. Studies of a number of small, model proteins have given the following general picture of how proteins fold (although there are always exceptions, and there is no generally accepted view of the protein folding mechanism.)

Generally, the fully folded state N begins to appear immediately without a detectable lag period and at a rate depending on both the identity of the protein and the refolding conditions. The rate is independent of the unfolding conditions used. The more physiological the refolding conditions, the greater the rate of folding (Fig. 1). Unfolding is observed by the reverse process of abruptly changing conditions from folding to those favoring unfolding, for example, by rapidly adding a denaturant or by changing the pH or temperature. The more denaturing the conditions, the more rapid the rate of unfolding.

**Figure 1.** Typical temperature dependence of the rates and equilibria of protein folding transitions not involving intrinsically slow isomerizations. The natural logarithms of the rate constants for unfolding and refolding are plotted as a function of  $(\text{temperature})^{-1}$  in an Arrhenius plot. A similar plot of the equilibrium constant  $K_{eq}$  between the folding (N) and unfolded (U) states is a van't Hoff plot. The curvature of the van't Hoff plot results from the greater apparent **heat capacity** of U than N. The linear Arrhenius plot for the rate of unfolding indicates that the folding transition state has the same heat capacity as N. The greater heat capacity of U is reflected entirely in the curvature of the Arrhenius plot for the rate of refolding because  $\ln K_{eq} = \ln k_{refolding} - \ln k_{unfolding}$ . The data used to construct this diagram are for hen egg-white [lysozyme](#) at pH 3 (2, 3), extrapolated to the absence of GdmCl. Although  $k_{refolding} = k_{unfolding}$  at  $K_{eq} = 1$ , it is a coincidence that the rate constants also have the value  $1\text{s}^{-1}$  at this temperature, so that all three curves intersect at a common point.

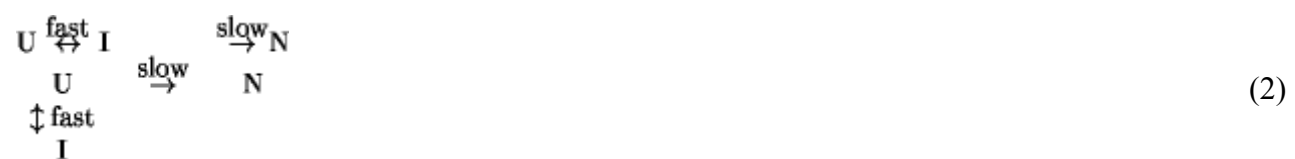


In both unfolding and refolding, all of the molecules have the same probability of undergoing the

transition, and a single rate constant is observed for the population, unless the molecules differ covalently or conformationally in a way that is only slowly interconverted. The most apparent instances of the latter are when there are both *cis* and *trans* isomers of peptide bond residues. (see [Cis/Trans Isomerization](#)). The *trans* form is intrinsically more stable, and the *cis* form is usually not present substantially, except when adjacent to [proline](#) residues. This *cis* form can occur in a folded protein, which usually has such a peptide bond *cis* in all of the molecules. In the unfolded state, however, there is an equilibrium mixture of both isomers. Consequently, some unfolded molecules have one or more incorrect isomers. Such isomerization is intrinsically slow, and the folding of these molecules is slowed or prevented by the incorrect isomer. The isomerization can be rate-limiting for their folding. The following discussion is limited to those small model proteins that show no such intrinsically slow process in folding.

The ratio of the rate constants for unfolding and refolding generally agrees with the equilibrium constant for the transition, so microscopic reversibility applies. There appears to be a classical [transition state](#) for the unfolding/refolding transition. That all of the many conformationally diverse molecules of a population of unfolded protein refold with the same rate constant, independent of how the protein was unfolded, indicates that all of the molecules must equilibrate rapidly and reversibly before undergoing the same **rate-limiting step**.

Upon transferring the unfolded protein to refolding conditions, some proteins adopt partly folded or molten globule conformations very rapidly, more rapidly than the native conformation appears. Such partly unfolded species are often considered responsible for the rapidity of folding, and much effort has gone into characterizing them. This is difficult because they are populated only transiently and are converted to N. These intermediate species, however, are usually in rapid equilibrium with the unfolded state, and it is generally not possible to distinguish between the two possible kinetic models, in which they are either on- or off-pathway intermediates I:



Many proteins do not adopt such partly folded species but remain unfolded until converting to N in an apparently all-or-nothing transition. These proteins also refold more rapidly than those that adopt partly folded conformations, although they are also the smaller proteins. Therefore, the presence of stable, partly folded intermediates is not necessary for rapid protein folding.

Such partly folded species are generally not detected as intermediates during unfolding, which is almost always an all-or-nothing transition, even with proteins that adopt partly folded intermediates in refolding. Some exceptions have, however, been reported ([1](#)).

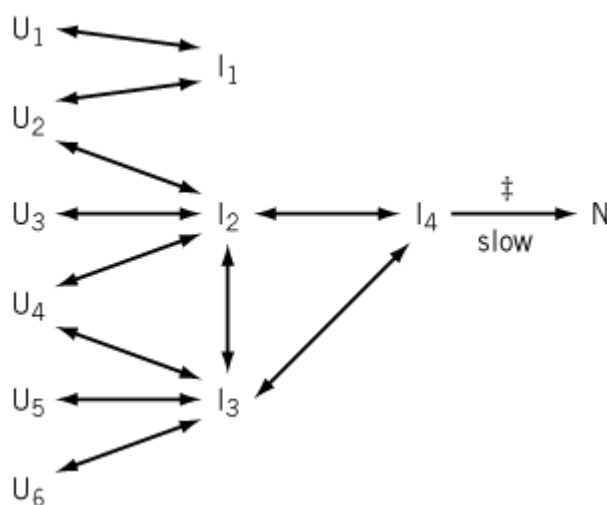
The [transition state](#) for the folding transition is characterized indirectly by measuring the rate of unfolding or refolding as a function of changing the conditions (Fig. [1](#)) or by altering the covalent structure of the protein. Plots of the logarithm of the rate constant versus the denaturant concentration are generally closely linear, which suggests that the nature of the transition state and the pathway remains constant. The particular data of Figure [1](#) indicate that the transition state in that case has the same **heat capacity** as the native state but somewhat lower **enthalpy** (estimated by the slopes of the Arrhenius plots). Similar observations are found with other proteins, usually by varying the denaturant concentration rather than the temperature. Studies using [protein engineering](#) to alter the structure of the protein systematically indicate that many, if not all, of the stabilizing interactions of the native state have been disrupted in the transition state. The transition state is close to the final conformation structurally but lacks the cooperativity that gives net stability to the fully folded conformation.

**Ligands** that bind tightly to the native protein generally do not increase the refolding rate of small proteins. Instead, they usually decrease the unfolding rate, and this is why they increase the stability of the folded conformation. These observations indicate that the transition state does not bind the ligand tightly.

Folded proteins are often cleaved by [proteinases](#) at specific sites on their surfaces, and the folded conformation is maintained. Such cleaved proteins can be unfolded and dissociated into two or three polypeptide fragments. Very often these fragments recombine and regenerate the folded conformation.

The general kinetic scheme suggested by the experimental observations for folding a single protein domain is illustrated in Figure 2.

**Figure 2.** A general kinetic model for protein folding indicated by the experimental observations, in the absence of intrinsically slow conformational isomerizations.  $U_i$  are various unfolded molecules with different conformations at the start of folding,  $I_i$  are partially folded molecules, and  $N$  is the fully folded protein. All kinetic steps indicated by arrows are rapid, except for that labeled “slow.” “ $\ddagger$ ” indicates the occurrence of the overall transition state. All steps are reversible, except for that indicated with a single-headed arrow, which occurs only in the indicated direction under conditions strongly favoring the folded state. As indicated, all of the unfolded molecules rapidly equilibrate under refolding conditions with a few partly folded species, which are also in rapid equilibrium. All of the molecules pass through a common slow step, which involves going through a transition state that is a distorted form of the native-like conformation. The intermediates  $I_i$  might be stable or unstable and therefore populated transiently or not, but all intermediates that occur after the rate-determining step are very unstable relative to  $N$ .



## 2. Kinetic Determination of Folding

If proteins cannot fold randomly and nonrandom folding pathways are crucial, the resulting folded state may not be the most stable conformation possible, but could be instead the form most kinetically accessible. If a kinetic pathway of folding is so vital, it should be possible to block folding by interfering with that pathway, and a protein might fold normally, solely for kinetic reasons, to a metastable state that is not the most stable thermodynamically. Examples are known, but only relatively few.

A number of bacterial [proteinases](#), such as [subtilisin](#) and  **$\alpha$ -lytic protease**, are synthesized as inactive precursors that have amino-terminal prosegments which are subsequently removed proteolytically to generate the active, native proteinase. These proproteinases unfold and refold *in vitro*, but the mature

forms do not refold. They refold only when the pro segment is added. The negligible rate of refolding of the mature protein, when the native protein is very stable, indicates a kinetic block to folding that is alleviated by the pro segment.

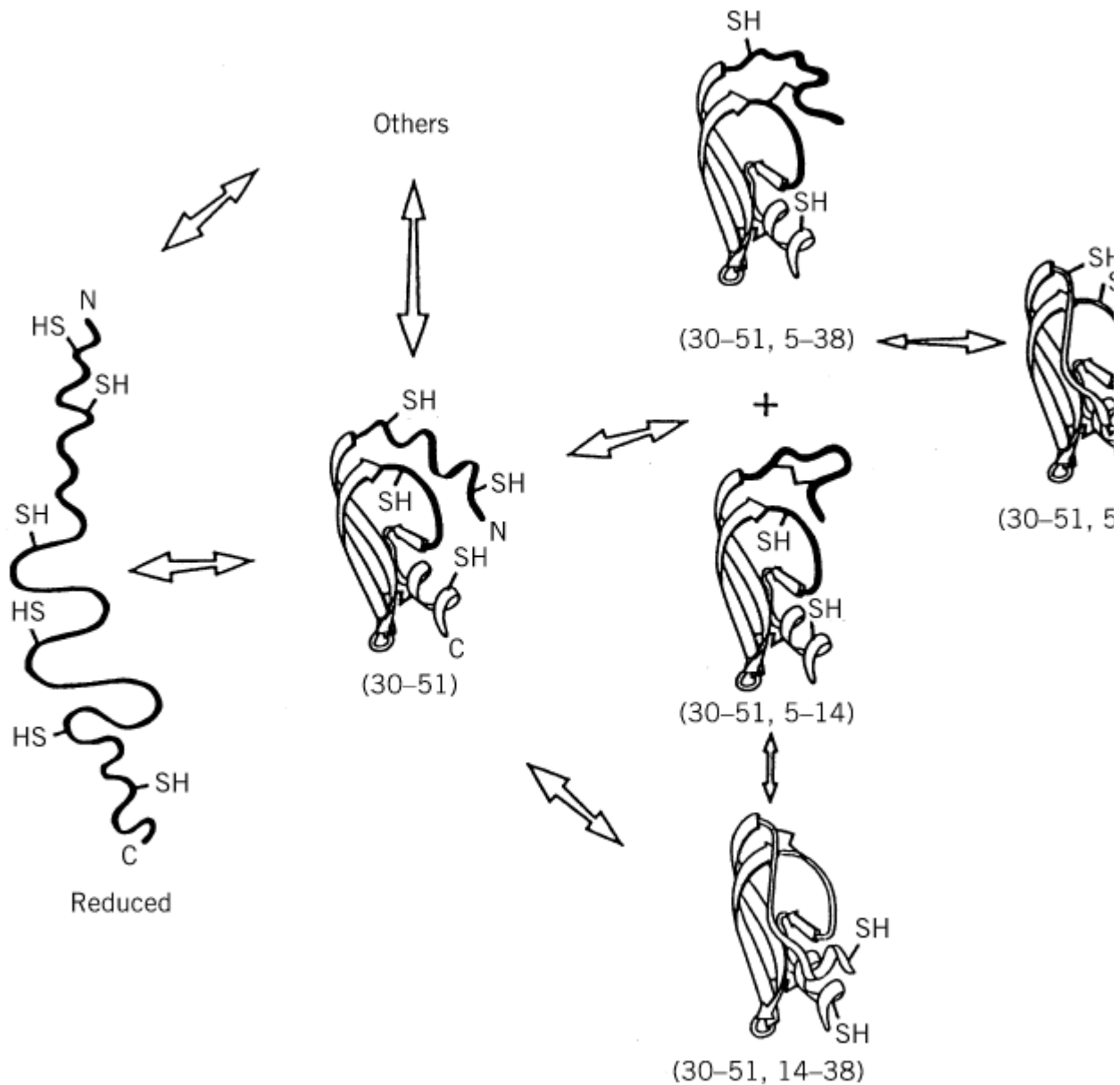
Within the [serpin](#) proteinase inhibitor family, plasminogen activator inhibitor –1 is synthesized in a form that is active as an inhibitor but relatively unstable. The active form slowly converts to a more stable form but is inactive as an inhibitor and known as the latent form. If this latent form is unfolded and then refolded, the active metastable form is regenerated before again undergoing the slow conversion to the inactive but stable form. Therefore, folding does not produce the most stable folded conformation directly but only through a metastable intermediate. The difference between the two forms involves a large change in a **b-sheet** of the protein.

### 2.1. Folding Coupled to Disulfide Formation

Many proteins that contain [disulfide bonds](#) between the [thiol groups](#) of **cysteine residues** become unfolded if these bonds are reduced, even in the absence of denaturants. The reduced protein remains unfolded even under physiological conditions. If the disulfides are permitted to form again, the protein can regenerate the native disulfide bonds and conformation. Then folding is coupled to disulfide formation. The great advantages of this are that disulfide bond formation and breakage can be controlled experimentally (using **thiol-disulfide exchange** with a disulfide reagent) and that any disulfide bonds in a protein can be trapped in a form that can be stable indefinitely. Cysteine residues can also be replaced or their thiol groups blocked irreversibly to decrease the number of disulfide possibilities, and the effect on the folding process can be used to dissect the pathway. Therefore, all the disulfide intermediates can be identified and characterized, and their roles in the folding process can be determined, usually unambiguously.

The best characterized disulfide folding pathway is that of [BPTI](#) (Fig. 3). Reduced BPTI is a very unfolded polypeptide chain, approximating a random coil, but with weak local interactions between residues close in the primary structure. Consequently, the initial formation of disulfide bonds is approximately random (after correcting for any differences in thiol group reactivity). The one-disulfide intermediates are not random, however, because that with the Cys30–Cys51 disulfide bond adopts a stable, partly folded conformation: this restricts which disulfide bonds can be formed subsequently. There is a kinetic block in forming either the 30–51 or 5–55 disulfide bonds, if the other is already present, which would generate the native-like (30–51,5–55) intermediate. Instead, this intermediate is normally formed most readily by intramolecular disulfide rearrangements of two intermediates with nonnative second disulfide bonds (Fig. 3).

**Figure 3.** The productive disulfide folding pathway of BPTI. R is the fully reduced protein. Intermediates are indicated by letters. The major disulfide intermediates are depicted, and their conformations determined by NMR analysis are indicated by the backbone. The backbone is unfolded or very flexible. The relative rates of the intramolecular step in forming each disulfide bond are indicated by the arrowhead. The wider the arrowhead, the greater the rate in that direction. The brackets indicate that the one-disulfide intermediates are in equilibrium. The “+” between two species indicates that they have the same kinetic roles. The nonproductive quasi-native-like intermediate is also shown.



The disulfide pathways elucidated to varying extents with a few other small proteins show similar properties but have variations. Initial disulfide formation in the reduced protein is approximately random until a stable folded conformation is adopted, which can either favor or disfavor formation of further disulfide bonds. Some proteins, such as BPTI (Fig. 3), adopt partly folded conformations. Some, such as  $\alpha$ -lactalbumin, adopt the molten globule conformation, and others remain unfolded until the entire folded conformation appears, such as ribonuclease A. In the absence of a folded conformation, further disulfide bonds are formed more slowly and are less stable. There is a kinetic block in forming a disulfide bond if that bond will become buried in a resulting stable folded conformation. Such kinetic blocks are caused by the high energy barrier that is most apparent in the reverse direction, upon reducing a buried disulfide bond. This kinetic barrier is usually overcome most readily in disulfide formation by intramolecular protein disulfide rearrangements in place of the intermolecular process of direct protein disulfide formation. The two processes involve, however, the same energy barrier and the same conformational transitions.

When the protein has all but one or two of the native disulfide bonds, it can adopt the stable native

conformation. For example, the very same native conformation is observed with BPTI in which any one of the three native disulfide bonds is missing. These quasi-native species indicate that disulfide bonds merely stabilize the native conformation and do not determine it. The stability of the quasi-native conformation varies, depending on which disulfide bond is missing. The quasi-native conformation inhibits formation of native disulfide bonds that will be buried but favors formation of those on the surface.

### 3. Multiple Domains

Large proteins are divided into **domains**, each of which is usually comparable to a small single-domain protein, although they vary in the extent to which the domains interact with and stabilize each other. When such interactions are not very great, a single domain excised from such a protein usually maintains the same conformation and refolds to it after being unfolded. The same principles are expected to apply to the folding of a single domain within a multidomain protein, and it is believed that multidomain proteins fold modularly by their individual domains folding and then associating. The final step, in which the domains interact, is often rate-limiting. Very often, the isolated domain folds more rapidly than when in the intact protein, suggesting that the unfolded domains interact and interfere with each other during folding.

The folded domains and subunits present during refolding of these complex proteins often recognize and bind their specific ligands. Consequently, the presence of such a ligand in this instance increases the rate of refolding and assembly.

Multidomain proteins also appear to be especially susceptible to aggregation during refolding. This is thought to be caused by specific complementary interactions between the domains of different molecules, similar to those that should occur between the domains of the same polypeptide chain.

#### 3.1. Multiple Subunits

Proteins with [quaternary structure](#), consisting of multiple polypeptide chains, are often dissociated and unfolded under denaturing conditions. Upon transfer to folding conditions, they can refold and reassemble to regenerate the original tertiary and quaternary structures. The folding of the individual domains and their reassembly is often observed as separate transitions. Assembly can be followed by covalent **cross-linking** at different times. Whether folding or assembly is rate-limiting depends on the protein concentration, because assembly is more rapid at higher protein concentrations. Kinetic results with a number of proteins are consistent with individual subunits folding and then reassembling, but the assembly process also involves some changes in the tertiary structure of the individual subunits. The large size and complexity of multisubunit proteins makes it difficult to elucidate the details. It is usually further complicated by the tendency of the protein to aggregate during refolding, especially at high protein concentrations. The interactions responsible for the aggregation are specific, because other proteins do not usually have an effect, unless very closely related.

In some dimeric proteins, the two polypeptide chains are highly entwined in their native conformations, and it is impossible for the same conformation to be stable with a single chain in the absence of the other subunit. Nevertheless, such proteins refold and reassemble. This process is probably similar to that in which proteolytic fragments of a single polypeptide chain reassociate and refold (see previous discussion). The two processes of folding and assembly must be coordinated and occur in a single, concerted step.

#### 3.2. Differences with Folding *in Vivo*

[Protein folding \*in vivo\*](#), during or after biosynthesis of the polypeptide chain, differs most fundamentally from refolding *in vitro* in that domains can fold sequentially as the polypeptide chain is synthesized. It also involves the concerted action of **chaperones** and catalysts of *cis-trans* isomerization (eg, [cyclophilin](#)) and disulfide bond formation, breakage, and rearrangement (eg,

[protein disulfide isomerase](#)).

## Bibliography

1. T. Kiefhaber, A. M. Labhardt, and R. L. Baldwin (1995) *Nature* **375**, 513–515.
2. S.-I. Segawa and M. Sugihara (1984) *Biopolymers* **23**, 2473–2488.
3. W. Pfeil and P. L. Privalov (1976) *Biophys. Chem.* **4**, 41.

## Suggestions for Further Reading

4. T. E. Creighton (1990) Protein folding, *Biochem. J.* **270**, 1–16.
5. T. E. Creighton, ed. (1992) *Protein Folding*, W. H. Freeman, New York.
6. R. H. Pain, ed. (1994) *Mechanisms of Protein Folding*, IRL Press, Oxford.
7. A. R. Fersht (1997) Nucleation mechanisms in protein folding, *Curr. Opin. Struct. Biol.* **7**, 3–9.

## Protein Folding *In Vivo*

Protein folding is the process by which the linear information contained in the [amino acid](#) sequence of a [polypeptide chain](#) gives rise to the unique three-dimensional [conformation](#) of the functional [protein structure](#). How folding is achieved with high efficiency constitutes one of the basic problems in biology. The discovery that unfolded polypeptides can refold spontaneously *in vitro* (1) (see [Protein Folding In Vitro](#)) suggested that the folding (acquisition of [tertiary structure](#)) and assembly (formation of [quaternary structure](#)) of newly-synthesized polypeptides *in vivo* also occur spontaneously, without the involvement of further components. It is now clear, however, that in many cases efficient folding in the cell depends on a machinery of preexisting proteins, the **molecular chaperones**, whose primary task, often in an energy-dependent manner, is to prevent protein misfolding and aggregation (2, 3). Thus *in vivo* the principle of [self-assembly](#) of proteins is replaced by a process of assisted self-assembly. Another important difference between protein folding *in vitro* and *in vivo* relates to the fact that in the cell folding occurs in the context of [translation](#) during **protein biosynthesis**. As a result, folding can initiate cotranslationally, and this mechanism is important for the successful folding of large modular polypeptides with multiple **domains**.

### 1. Protein Aggregation

Many proteins refold via compact globular intermediates that contain varying amounts of **secondary structure** but lack the stable tertiary structure that defines the native state (see [Molten Globule](#)). These intermediates expose to solvent some **hydrophobic** amino acid residues (which mostly will be buried upon correct folding) and via these residues tend to associate with one another to form aggregates. Aggregate formation is usually an irreversible off-pathway step. *In vitro*, the extent of aggregation can often be controlled by lowering the protein concentration and the temperature of the refolding mixture and by adjusting other physical parameters of the reaction, such as pH and ionic strength. In contrast, cellular conditions dictate that, without the intervention of molecular chaperones, aggregation would outcompete correct folding, at least for a significant fraction of newly synthesized polypeptides. The concentration of unfolded polypeptides emerging from [ribosomes](#) (ie, “nascent” polypeptide chains) in the cytosol is very high; it reaches  $\sim 35 \mu\text{M}$  ( $\sim 1\text{mg/ml}$  for a 30-kDa chain) in *Escherichia coli*, assuming a uniform distribution of ribosomes. The local concentration of nascent chains is significantly greater, however, because of macromolecular crowding (see [Excluded Volume](#)) and because translating ribosomes are organized in

**polyribosomes.** *Macromolecular crowding* refers to the fact that a large fraction of the cellular volume is occupied by proteins and other **macromolecules** at a total concentration of ~300 g/L and is therefore not available to other macromolecules (4, 5). Crowding is predicted to result in an increase by several orders of magnitude in **association constants** for unfolded polypeptides over those in dilute solution.

The risk of newly synthesized polypeptides aggregating is enhanced by the inability of nascent chains to fold into stable tertiary structures, at least during the early phase of translation. Stable folding requires the presence of a complete protein domain (usually ~100–300 amino acid residues in length) that can fold independently (see [Protein Structure](#)). As the C-terminal ~30 amino acid residues of a translating polypeptide are tethered to the ribosome, and are thus topologically restricted, nascent chains remain unfolded until an entire domain has emerged from the ribosome. Aggregation of these nascent chains is thought to be prevented by the co-translational binding of molecular chaperones, including members of the **heat-shock** proteins 70 (Hsp70) and Hsp40 (DnaJ) families (3).

## 2. Molecular Chaperones

Molecular chaperones were originally defined as proteins that mediate the correct assembly of other proteins, but are not themselves components of the final functional structures (2). Chaperones occur ubiquitously, and many of them are classified as **stress response** proteins, although their functions are essential under normal growth conditions (3, 6-9). Most chaperones function by stabilizing an otherwise unstable conformer of another protein—and by controlled binding and release, may facilitate its correct fate *in vivo*, whether this is folding, oligomeric assembly, transport to a particular subcellular compartment, or disposal by degradation (9). Molecular chaperones do not contribute steric information for correct folding, but rather prevent incorrect interactions between (and perhaps also within) nonnative polypeptides, thus typically increasing the yield, but not the rate, of folding reactions. These properties distinguish the chaperones from so-called folding catalysts, [protein disulfide isomerases](#), and [peptidylprolyl cis–trans isomerases](#), which accelerate intrinsically slow steps in the folding of some proteins, namely, the rearrangement of disulfide bonds in secretory proteins (see [Protein Secretion](#)) and the *cis–trans* isomerization of peptide bonds preceding [proline](#) residues, respectively (10, 11).

Molecular chaperones fall into several structurally unrelated families of proteins, including the members of the Hsp70, Hsp40 (DnaJ), and Hsp90 families, as well as the **chaperonins** (Hsp60) and the so-called small heat-shock proteins. Most of these proteins are soluble, but **membrane-bound** chaperones, such as **calnexin** in the [endoplasmic reticulum](#) (ER), exist as well. A table summarizing the main classes of molecular chaperones is found in the entry **Molecular chaperone**. Chaperones may be expressed constitutively or on exposure of cells to stresses, such as high temperature. Those chaperones with a broad spectrum of protein substrates generally suppress protein aggregation by recognizing hydrophobic amino acid residues or [accessible surfaces](#) that are exposed by [unfolded proteins](#) or incompletely folded polypeptide chains. The need for an increase in cellular chaperone capacity at elevated growth temperatures (or upon exposure of cells to other stresses) is explained by the tendency of preexisting proteins to unfold under these conditions. Proper folding is usually achieved by the controlled dissociation of the complexes between unfolded polypeptides and chaperones, which often occurs in an ATP- and cofactor-dependent mechanism.

Of the chaperones with a well-documented role in assisting the folding of newly synthesized polypeptides (*de novo* folding), the members of the Hsp70 and chaperonin classes have been studied most extensively. They represent two basic paradigms of ATP-dependent chaperone action.

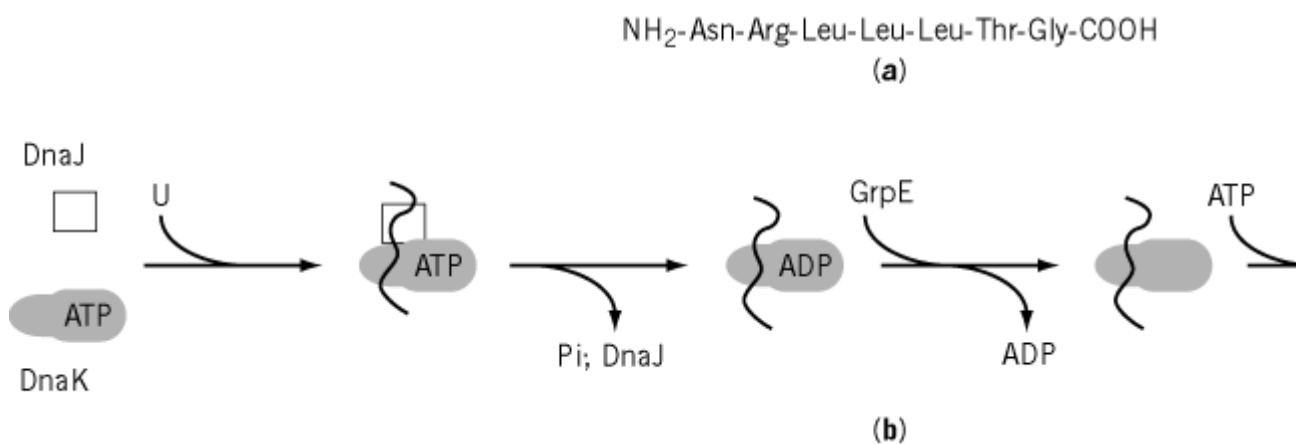
## 3. Hsp70s

The Hsp70s are a family of highly conserved **ATPases** of relative molecular mass ~70,000 found in **prokaryotes** (see [DnaK/DnaJ Proteins](#)) as well as in eukaryotes, where they occur in the **cytosol**,



mitochondria, **chloroplasts**, and the ER (12) (see [BiP \(Hsp70\)](#)). Protection of nascent chains from aggregation is thought to be the main role of Hsp70 in de novo protein folding (13-15). In addition, Hsp70s have other important functions in protein metabolism under both stress- and nonstress conditions, including functions in membrane translocation of proteins and the degradation of misfolded proteins. This versatility of the Hsp70s results from their basic function, which is to bind and release hydrophobic peptide segments that are generally exposed in unfolded polypeptides. These peptide segments are on average seven amino acid residues in length and must be enriched in hydrophobic amino acids, such as [leucine](#) (Fig. 1a) (16-18). These segments are bound in an extended conformation within a binding cleft of the ~18 – kDa C-terminal domain of Hsp70, whose high resolution [X-ray crystallography](#) structure is known (19). Peptide binding and release by the C-terminal domain is regulated via ATP-dependent conformational changes in an N-terminal ATPase domain of ~45 kDa, which has the three-dimensional structure of [actin](#) (20). Peptide binds to the ATP-state of Hsp70 in which the peptide binding cleft in the C-terminal domain is open (Fig. 1b). In the ADP state, the peptide cleft is closed, preventing the peptide from dissociating. Efficient peptide binding depends on essential cofactors of the Hsp40 (DnaJ) family. *E. coli* DnaJ and several of its eukaryotic homologues are chaperones in their own right, as well as activators of the Hsp70 ATPase (21-23). DnaJ can present an unfolded polypeptide to ATP-bound Hsp70 (DnaK in *E. coli*). DnaJ-catalyzed hydrolysis of ATP results in stable peptide binding. In *E. coli* (and in mitochondria) an additional cofactor, GrpE, catalyzes the exchange of Hsp70-bound ADP for ATP and thereby facilitates polypeptide release (21, 24). It is noteworthy that the Hsp70 reaction cycle resembles the regulation of certain [GTP-binding proteins](#) (3, 23) (Fig. 1b).

**Figure 1.** The Hsp70 chaperone system. (a) Example of a peptide with high affinity for Hsp70. This peptide was cocrystallized with the peptide binding domain of DnaK (19). Peptide segments with similar properties occur on average at a distance of 50–100 amino acids. (b) Mechanism of the *E. coli* Hsp70 system DnaK (Hsp70), DnaJ, and GrpE. Note that *E. coli* DnaJ binds to the unfolded polypeptide, but this may not be the case for other Hsp40 homologs, such as mammalian Hsp40. Whenever U is released, provided all structural elements necessary for folding are available (eg, on release of a nascent chain from the ribosome), the polypeptide will fold. of interaction with the Hsp70 system.



Hsp70 functions essentially as a [buffer](#) for unfolded and incompletely folded polypeptides, thus reducing the concentration of aggregation-sensitive folding intermediates. The net result of Hsp70 action is the binding and release of the polypeptide chain in an unfolded conformation. On release, the unfolded polypeptide may fold spontaneously to the native state, provided all structural elements necessary for folding are available, or it may be transferred to another chaperone (see text below) or rebind to Hsp70. On the basis of these properties, the Hsp70 chaperone system is ideally suited to protect nascent polypeptide chains against aggregation and to assist in their folding. As long as a polypeptide chain (or a domain) is not yet completely synthesized, release of nascent polypeptide from Hsp70 will not result in folding but in rebinding to Hsp70 as the chain continues to expose

hydrophobic residues (see Fig. 1). Indeed, a large fraction of nascent polypeptide chains interact with Hsp70 *in vivo* (15), but whether this interaction is generally required for folding remains to be demonstrated.

#### 4. Chaperonins

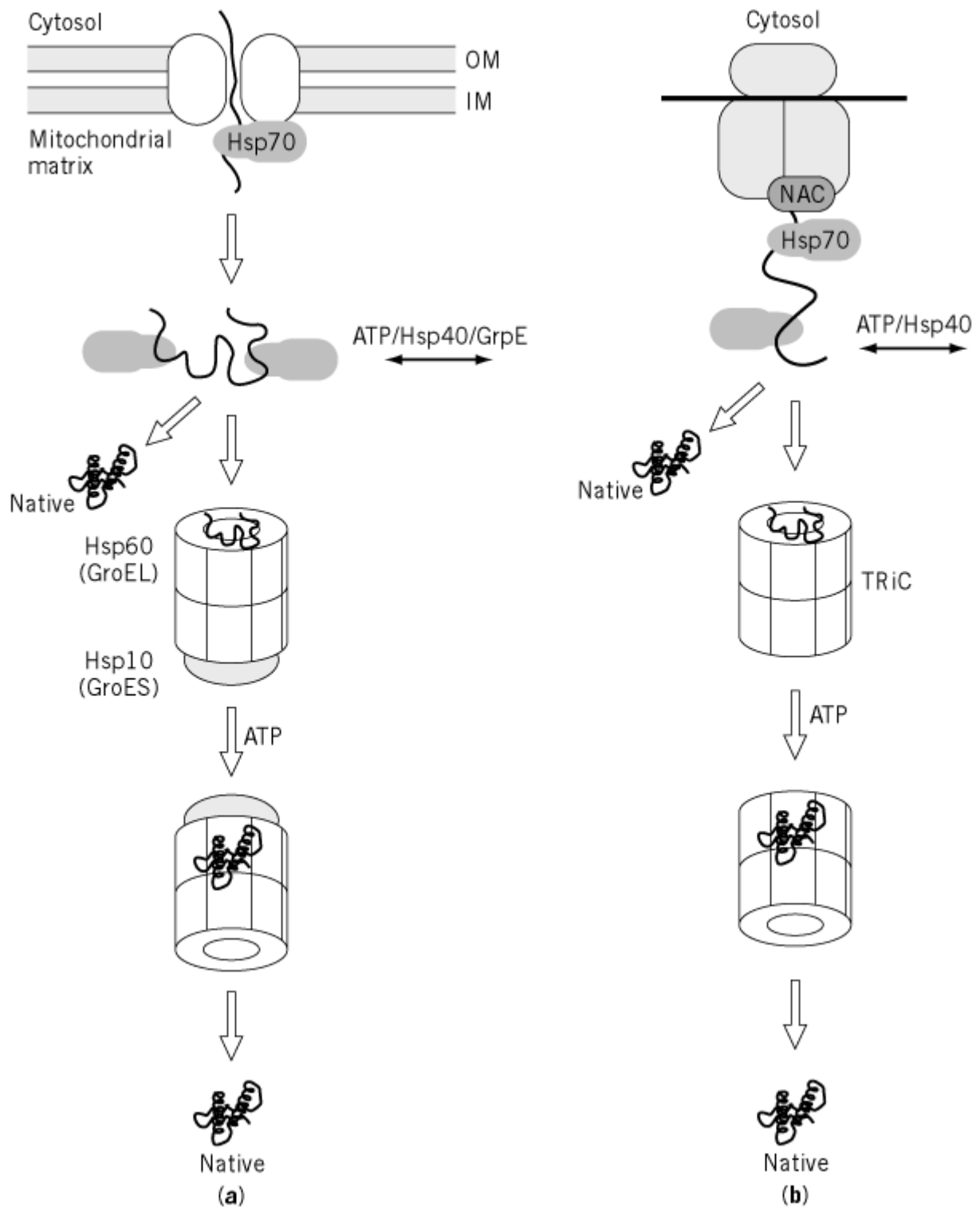
In contrast to the Hsp70s, the chaperonins (also classified as Hsp60s) are large cylindrical protein complexes consisting of two rings of ~60 – kDa subunits that are stacked back to back. In the case of eubacterial chaperonins, such as GroEL of *E. coli*, there are seven identical subunits per ring, whereas the chaperonins of **archaeobacteria** and the eukaryotic cytosol are heterooligomeric and may contain up to eight different subunits (3, 25). The salient structural feature of these ~800 – kDa complexes is a central cavity in which a single polypeptide molecule can fold and avoid aggregation with other unfolded polypeptides. Unfolded polypeptide chain binds to hydrophobic patches at the inner wall of the chaperonin cavity. In the case of GroEL, binding of the ring-shaped cofactor GroES to the opening of the GroEL cylinder then displaces the polypeptide from these binding sites into an enclosed folding cage (26-29). The ATPase activity of the chaperonin regulates the closing and opening of the cage. Binding to chaperonin may also unfold polypeptides that have been kinetically trapped in misfolded states, thereby giving them another chance to fold upon release into the cavity. (See also [Chaperonin](#).)

Chaperonins are essential for cell viability under all growth conditions (30-32), because they are required for the folding of a subset of newly synthesized polypeptides. In *E. coli*, GroEL interacts with approximately 15% of total newly synthesized cytosolic proteins at growth temperatures of 30–37°C, and with up to 30% or more under heat stress at 42°C (33). In addition to chaperonin-dependent proteins that interact more or less quantitatively with GroEL, other proteins transit GroEL with only a few percent of the total population of molecules. These latter proteins do not depend on GroEL for folding, at least under nonstress conditions. GroEL acts predominantly post-translationally, and most of its substrates fall into the size range of 10 to 55 kDa, the upper size limit of the folding cage (33). According to *in vitro* studies, GroEL-dependent proteins have relatively slow folding rates and therefore are highly aggregation-sensitive. When expressed in cells containing insufficient or dysfunctional chaperonin, these proteins aggregate or are degraded proteolytically. The mechanisms of folding for polypeptides larger than ~55 kDa will be discussed below.

#### 5. Chaperone Pathways in Folding

The Hsp70 and chaperonin systems can act sequentially. This sequential action has been demonstrated for protein folding in the cytosolic compartment and in mitochondria and chloroplasts (3). Considering the origin of mitochondria from endosymbiotic bacteria, the matrix of mitochondria—specifically, the space surrounded by the inner mitochondrial membrane—is evolutionarily related to the bacterial cytosol. This compartment contains homologues of all the major bacterial chaperones, DnaK (Hsp70), DnaJ, GrpE, GroEL, and GroES. Although mitochondria still synthesize a small number of proteins in the matrix, most mitochondrial proteins are imported from the cytosol (34). A prerequisite for effective membrane [translocation](#) is that these polypeptides be maintained in an unfolded state by cytosolic chaperones (see [Translocation](#)). On import, they interact first with mitochondrial Hsp70 that is bound to the inner surface of the inner membrane and then with soluble Hsp70 in the matrix (Fig. 2a). Membrane-bound Hsp70 participates actively in the translocation process itself, whereas soluble Hsp70 assists in folding. A subset of imported proteins has to be transferred from Hsp70 to the mitochondrial GroEL homologue, Hsp60, for successful folding (35, 36). *In vitro* reconstitution experiments using the *E. coli* chaperones established the mechanistic significance of this chaperone relay (22). The Hsp70 system prevents aggregation of unfolded polypeptides, preserving their folding competence, whereas the chaperonin mediates folding to the native state. This pathway is not necessarily unidirectional; polypeptides that cannot fold on the chaperonin may be transferred back to Hsp70 and may eventually be degraded (3).

**Figure 2.** Sequential action of the Hsp70 and chaperonin systems in folding. **(a)** Protein folding in mitochondria. OM and IM, outer and inner mitochondrial membranes, respectively. Hsp60, mitochondrial GroEL; Hsp10, mitochondrial GroES. Mitochondrial Hsp70 interacts first with the incoming polypeptide because of its ability to recognize extended peptide motifs and its specific association with the translocation machinery of the inner mitochondrial membrane. Release of newly translocated polypeptide from Hsp70 is dependent on ATP and the action of mitochondrial Hsp40 (DnaJ) and GrpE (see Fig. 1). The Hsp70 and Hsp60 (chaperonin) systems are functionally distinct in that only Hsp60 can release the substrate protein in a fully folded state. Small, rapidly folding proteins will either not interact with Hsp60 or will do so only with low efficiency. **(b)** Folding of newly synthesized polypeptides in the eukaryotic cytosol. TRiC, the cytosolic chaperonin. Cytosolic Hsp70 interacts with the nascent polypeptide because of its ability to recognize extended peptide motifs. Binding of NAC very close to the peptidyl-transferase center may precede that of Hsp70 for most cytosolic proteins. Release of newly translated polypeptide from Hsp70 is dependent on ATP and the action of Hsp40 (DnaJ). This step is probably GrpE-independent in eukaryotes. Most proteins fold upon release from Hsp70, but a subset of proteins needs assistance by the chaperonin for folding. Although folding of these proteins occurs in the central cavity of TRiC (as in the case of GroEL), TRiC is independent of a GroES cofactor. The function of GroES to form a lid on the opening of the chaperonin cylinder appears to be integrated into the structure of the TRiC subunits (50).



There is increasing evidence from *in vivo* studies that the Hsp70/chaperonin pathway plays an important role in the prokaryotic and eukaryotic cytosol, both for the folding of newly synthesized polypeptides and during the refolding of stress-denatured polypeptides (3, 37-41). Hsp70 interacts with a wide array of nascent chains in eukaryotes (15). While the majority of proteins probably do not have to transit further chaperones to complete folding, a subset of proteins, presumably those that fold slowly, must be transferred from Hsp70 to chaperonin to reach their native state (39, 40) (Fig. 2b). The major substrate proteins in eukaryotes that follow such a pathway are the cytoskeletal proteins [actin](#) and [tubulin](#), which are critically dependent on the cytosolic chaperonin (CCT), or

TRiC) for folding (32, 40, 42, 43). The situation is similar in bacteria in that only a fraction of all newly formed polypeptide chains bind to chaperonin (33). For example, folding of certain bacterial forms of **ribulose biphosphate carboxylase** expressed in *E. coli* requires interaction with DnaK (Hsp70) and with GroEL, whereby the two systems must act sequentially and not in parallel (41).

Although the basic mechanistic principles of chaperone action are now well understood, the complexity of chaperone-assisted folding pathways *in vivo* is only beginning to be appreciated. For example, in *E. coli* cytosol a significant fraction of nascent chains interact with **trigger factor**, a 48-kDa chaperone that has an affinity for ribosomes and possesses both chaperone and peptidylprolyl *cis*–*trans* isomerase activities (44). In eukaryotes many nascent chains interact first with nascent-chain-associated complex (NAC) before they form a complex with Hsp70 (45) (Fig. 2b). NAC is a heterodimer of 33-kDa and 21-kDa subunits. NAC binds also to ribosomes and prevents their association with the membrane of the endoplasmic reticulum, except when a secretory precursor polypeptide is synthesized. For cytosolic proteins, NAC binding may help to recruit Hsp70.

A multiplicity of chaperone components cooperate in mediating proper folding, **disulfide bond** formation, and **glycosylation** of secretory proteins in the lumen of the endoplasmic reticulum (ER) (see **Protein Secretion**). Correct folding and assembly is a prerequisite for the packaging of proteins into **vesicles** that travel from the ER via the **Golgi apparatus** to the cell surface, and the ER lumen is effectively a highly concentrated solution of chaperones and protein folding catalysts. These include, among others, the Hsp70 homologue **BiP**, various Hsp40s, the Hsp90 homologue Grp94, **protein disulfide isomerase**, **calnexin**, and **calreticulin**. The latter two chaperones recognize certain carbohydrate modifications that are typical of incompletely folded polypeptides and retain these polypeptides in the ER until folding is completed (46). It is noteworthy that the ER does not contain a chaperonin homologue.

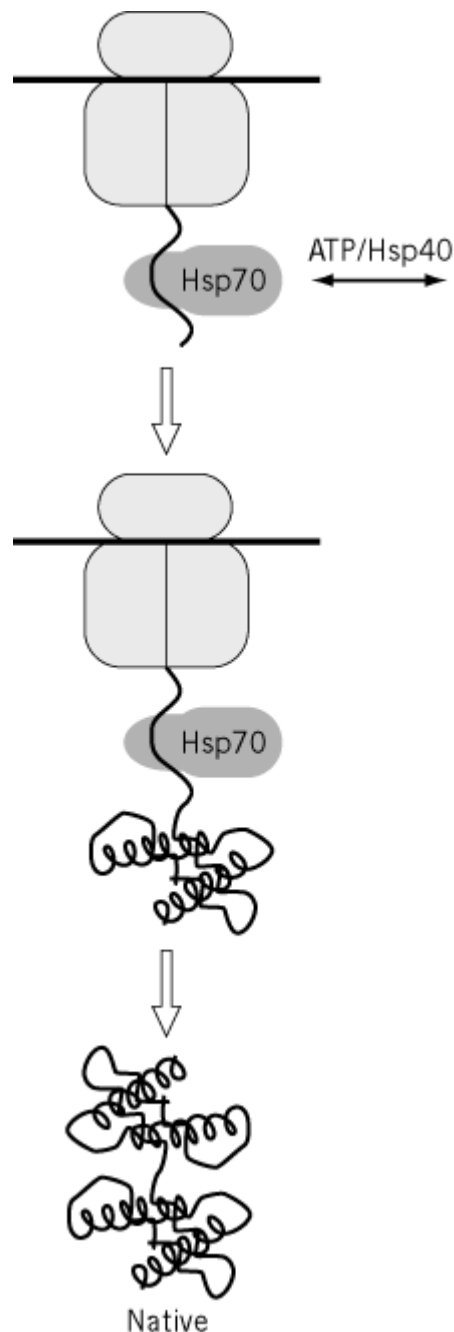
## 6. Cotranslational Folding of Multidomain Proteins

Analysis of the size distribution of proteins in several completely sequenced genomes indicates that eukaryotes have a proportionally larger number of modular polypeptides consisting of multiple protein domains, than do bacteria (47, 48). For example, in the yeast *Saccharomyces cerevisiae* the average protein has a length of 496 amino acids (~55 kDa), and ~38% of all yeast proteins are larger than 55 kDa, including ~1450 soluble proteins (48). In contrast, the average length of an *E. coli* protein is only 317 residues (~35 kDa), and only 13% of all *E. coli* proteins exceed 55 kDa, the size cutoff of the GroEL/GroES folding cage. The size distribution of protein domains (the “folding units”) is uniform across all three kingdoms of life (bacteria, archaea, eukarya) in the range of 100–300 residues. Thus a genome encoding proportionally longer polypeptides must encode more and/or longer multidomain polypeptides. Since these proteins frequently do not refold efficiently *in vitro*, it would be expected that their folding *in vivo* is particularly chaperonin-dependent. However, the volume capacity of the central cavity of the eukaryotic cytosolic chaperonin, TRiC, is probably not much greater than that of GroEL. Moreover, TRiC is of low abundance in many eukaryotic cells and is thought to interact with only a restricted subset of polypeptides (see text above). Neither do other abundant chaperone proteins, including the Hsp90 system (49), play a general role in *de novo* protein folding. How do large modular polypeptides manage to fold efficiently *in vivo*?

In modular polypeptides, domains are often joined by flexible linker segments. Such proteins are able to fold cotranslationally as their domains emerge sequentially from the ribosome (3, 47). This mechanism allows a high efficiency of folding by reducing unproductive intramolecular interactions between concurrently folding domains. Such interactions may occur during the collapse of the unfolded chain into a disorganized globule and may explain the tendency of multidomain proteins to misfold *in vitro* (see **Protein Folding In Vitro**). Mechanistically, cotranslational folding reduces the problem of folding a large polypeptide to the folding of its independent modules or domains, those structural units most able to fold spontaneously. Sequential domain folding probably relies predominantly on the protection of nascent chains by Hsp70 until a complete domain has been synthesized and emerged from the ribosome (Fig. 3). ATP-dependent release of Hsp70 may then

allow domain folding. The ATPase activity of Hsp70 in the eukaryotic cytosol (~1ATP hydrolyzed per molecule per minute) seems to be adjusted to the speed of translation (about two to three amino acids per second) such that Hsp70 would bind and release the nascent chain once during the synthesis of a polypeptide domain of average length. Transfer from Hsp70 to the chaperonin may be necessary only for proteins that are unable to fold into stable structures during translation. This transfer is the case for some multidomain proteins in which the domains are constructed of discontinuous sequence segments of the polypeptide chain. Here a continuous chain forms part of a domain, then leaves the compact region to form part or all of another domain, after which it returns to complete the previous domain. Actin, one of the main substrate proteins of the eukaryotic cytosolic chaperonin, is composed of discontinuous domains and forms stable tertiary structure only post-translationally. Similarly, proteins whose domains are structurally unstable in isolation (and are ultimately stabilized by interactions with other domains or subunits) may also require sequestration in the chaperonin folding cage for post-translational folding. Thus a combination of chaperonin-independent (cotranslational) and chaperonin-dependent (post-translational) mechanisms is operative in the eukaryotic cytosol (48).

**Figure 3.** Protein folding pathways in the eukaryotic cytosol. Model for the cotranslational folding of a multidomain protein in the eukaryotic cytosol. Folding is assisted by the Hsp70 system independently of the chaperonin. Folding of a completed domain occurs as Hsp70 dissociates from the nascent chain in an ATP-dependent manner.



Evidence has been presented that the bacterial translation–folding machinery has a reduced capacity to support the cotranslational and sequential folding of multidomain proteins, when domain folding is slow compared to the rapid speed of bacterial translation (15–20 amino acids per second) (47). It has been proposed that bacterial proteins are generally selected for efficient post-translational folding. Inefficiency of cotranslational folding in bacteria would help to explain not only why many eukaryotic multidomain proteins misfold upon bacterial expression, but also perhaps why the bacterial protein complement is structurally less complex than that of eukaryotic cells. Modular polypeptides are believed to have evolved by random [gene fusion](#) events. Thus, if cotranslational folding in bacteria were even partially constrained, the evolution of multidomain proteins by domain shuffling would be less frequent than in organisms more generally able to support sequential domain folding. Future research will have to explore to what extent differences in folding mechanism between bacterial and eukaryotic cells may be responsible for the explosive evolution of modular polypeptides in eukaryotes.

## Bibliography

1. C. B. Anfinsen (1973) *Science* **181**, 223–230.
2. R. J. Ellis (1987) *Nature* **328**, 378–379.
3. F. U. Hartl (1996) *Nature* **381**, 571–580.
4. R. J. Ellis and F. U. Hartl (1996) *FASEB J.* **10**, 20–26.
5. S. B. Zimmermann and A. P. Minton (1993) *Annu. Rev. Biophys. Biomol. Struct.* **22**, 27–65.
6. R. J. Ellis and S. van der Vies (1991) *Annu. Rev. Biochem.* **60**, 321–347.
7. M.-J. Gething and J. Sambrook (1992) *Nature* **355**, 33–45.
8. W. J. Welch (1991) *Curr. Opin. Cell. Biol.* **3**, 1033–1038.
9. J. P. Hendrick and F. U. Hartl (1993) *Annu. Rev. Biochem.* **62**, 349–384.
10. R. B. Freedman (1995) *Curr. Opin. Struct. Biol.* **5**, 85–91.
11. F. X. Schmid (1993) *Annu. Rev. Biophys. Biomol. Struct.* **22**, 123–143.
12. E. A. Craig, B. D. Gambill, and R. J. Nelson (1993) *Microbiol. Rev.* **57**, 402–414.
13. R. P. Beckmann, L. A. Mizzen, and W. J. Welch (1990) *Science* **248**, 850–854.
14. R. J. Nelson, T. Ziegelhoffer, C. Nicolet, M. Werner-Washburne, and E. A. Craig (1993) *Cell* **71**, 97–105.
15. D. K. Eggers, W. J. Welch, and W. J. Hansen (1997) *Mol. Biol. Cell* **8**, 1559–1573.
16. G. C. Flynn, J. Pohl, M. T. Flocco, and J. E. Rothman (1991) *Nature* **353**, 726–730.
17. S. Blond-Elguindi et al. (1993) *Cell* **75**, 717–728.
18. S. Rüdiger, L. Germeroth, J. Schneider-Mergener, and B. Bukau (1997) *EMBO J.* **16**, 1501–1507.
19. X. T. Zhu et al. (1996) *Science* **272**, 1606–1614.
20. K. M. Flaherty, C. DeLuca-Flaherty, and D. B. McKay (1990) *Nature* **346**, 623–628.
21. K. Liberek, J. Marszalek, D. Ang, C. Georgopoulos, and M. Zylicz (1991) *Proc. Natl. Acad. Sci. USA* **88**, 2874–2878.
22. T. Langer et al. (1992) *Nature* **356**, 683–689.
23. A. Szabo et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 10345–10349.
24. C. J. Harrison, M. Hayer-Hart, M. Di Liberto, F. U. Hartl, and J. Kuriyan (1997) *Science* **276**, 431–435.
25. K. Braig et al. (1994) *Nature* **371**, 578–586.
26. J. Martin, M. Mayhew, T. Langer, and F. U. Hartl (1993) *Nature* **366**, 228–233.
27. M. Mayhew et al. (1996) *Nature* **379**, 420–426.
28. J. S. Weissman, H. S. Rye, W. A. Fenton, J. M. Beechem, and A. L. Horwich (1996) *Cell* **84**, 481–490.
29. Z. Xu, A. L. Horwich, and P. B. Sigler (1997) *Nature* **388**, 741–749.
30. O. Fayet, T. Ziegelhoffer, and C. Georgopoulos (1989) *J. Bacteriol.* **171**, 1379–1385.
31. A. L. Horwich, K. B. Low, W. A. Fenton, I. N. Hirshfield, and K. Furtak (1993) *Cell* **74**, 909–917.
32. H. Kubota, G. Hynes, and K. Willison (1995) *Eur. J. Biochem.* **230**, 3–16.
33. K. L. Ewalt, J. P. Hendrick, W. A. Houry, and F. U. Hartl (1997) *Cell* **90**, 491–500.
34. W. Neupert (1997) *Annu. Rev. Biochem.* **66**, 863–917.
35. J. Ostermann, A. L. Horwich, W. Neupert, and F.-U. Hartl (1989) *Nature* **341**, 125–130.
36. S. Rospert et al. (1996) *EMBO J.* **15**, 764–774.
37. A. Gragerov et al. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 10341–10344.
38. G. A. Gaitanaris, A. Vysokanov, S.-Z. Hung, M. Gottesman, and A. Gragerov (1994) *Mol. Microbiol.* **14**, 861–869.



39. J. Frydman, E. Nimmesgern, K. Ohtsuka, and F. U. Hartl (1994) *Nature* **370**, 111–117.
40. J. Frydman and F. U. Hartl (1996) *Science* **272**, 1497–1502.
41. S. K. Checa and A. M. Viale (1997) *Eur. J. Biochem.* **248**, 848–855.
42. Y. Gao, J. O. Thomas, R. L. Chow, G. H. Lee, and N. J. Cowan (1992) *Cell* **69**, 1043–1050.
43. M. B. Yaffe et al. (1992) *Nature* **358**, 245–248.
44. T. Hesterkamp, S. Hauser, H. Lütcke, and B. Bukau (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4437–4441.
45. B. Wiedmann, H. Sakai, T. A. Davis, and M. Wiedmann (1994) *Nature* **370**, 434–440.
46. A. Helenius, E. S. Trombetta, D. N. Herbert, and J. S. Simons (1997) *Trends Cell Biol.* **7**, 193–200.
47. B. Netzer and F. U. Hartl (1997) *Nature* **388**, 343–349.
48. W. J. Netzer and F. U. Hartl (1998) *Trend Biochem. Sci.* **23**, 68–73.
49. D. F. Nathan, M. H. Vos, and S. Lindquist (1997) *Proc. Natl. Acad. Sci. USA* **94**, 12949–12956.
50. M. Klumpp, W. Baumeister, and L.-O. Essen (1997) *Cell* **91**, 263–270.

### Suggestions for Further Reading

51. R. I. Morimoto, A. Tissieres, and C. Georgopoulos, eds. (1994) *The Biology of Heat Shock Proteins and Molecular Chaperones*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–593. (Contains a collection of articles by leading authors on the function of various heat shock proteins and molecular chaperones.)
52. R. J. Ellis, ed. (1996) *The Chaperonins*, Academic Press, San Diego, CA, pp. 1–323. (Contains 10 chapters devoted solely to the structure, function, and regulation of the chaperonins.)
53. M. J. Gething, ed. (1997) *Guidebook to Molecular Chaperones and Protein Folding Catalysts*, Oxford Univ. Press, Oxford, U.K. (Contains detailed information about the components of the cellular machinery of protein folding.)

## Protein Kinase A (Cyclic Amp-Dependent Protein Kinase)

The effects of [cyclic AMP](#) (cAMP) as a [second messenger](#) are mediated by the activation of cAMP-dependent protein kinase, also known as protein kinase A (PKA). PKA exists as a tetramer, with two catalytic (C) subunits, and two regulatory (R) subunits. The binding of cAMP to the R subunits results in their dissociation from C, releasing the catalytic subunit from the inhibitory effect of R, and activating it for phosphorylation of acceptor proteins. Because the binding of cAMP to R is **cooperative**, small changes in cAMP levels can produce a large increase in protein phosphorylation in cells ([1](#)).

The physiological substrates for PKA are many and varied. This enzyme participates in the regulation of secretion, metabolism, [transcription](#), and almost every regulatory process known in [signal transduction](#). The importance of PKA as a physiological effector has been studied pharmacologically, with peptide and organic inhibitors. There have also been a number of transgenic mice made lacking **alleles** containing different components of the PKA family of proteins.

### Bibliography

1. S. S. Taylor (1989) *J. Biol. Chem.* **264**, 8443–8446.

## Protein Kinase C

**Hormone**-stimulated hydrolysis of phospholipids by [phospholipases](#) gives rise to two potential [second messengers](#), calcium and diacylglycerol. Both of these molecules are directed to the activation of another family of protein kinases, protein kinase C (PKC). PKC is a family of enzymes, presently known to comprise up to nine members, with different patterns of expression and sensitivity to phospholipids, diacylglycerol, and calcium. In resting cells, these enzymes are largely cytoplasmic, but they translocate to the plasma membrane upon the generation of diacylglycerol by [phospholipase C](#).

PKC is a multidomain protein. The family members contain a conserved catalytic domain, along with additional calcium-, phospholipid-, and substrate-binding sites that vary considerably. In general, PKC  $\alpha$  and  $\gamma$  are completely dependent on calcium, while  $\beta$  is less calcium-sensitive. Other forms of the enzyme are calcium-insensitive (1).

The activation of PKC in different cells results in a large number of biological effects. PKC can **phosphorylate** a variety of substrates, regulating processes in metabolism, differentiation, and growth. Its substrates include [transcription factors](#), other enzymes, receptors, and **ion channels**, to name a few. Moreover, PKC is involved in several examples of cross-talk and feedback regulation in [signal transduction](#) (1).

### Bibliography

1. Y. Nishizuka (1992) Science **258**, 607–614.

## Protein Motif

A protein motif, sometimes referred to as a [supersecondary structure](#), is the term used to describe certain common combinations of **secondary structure** elements that are observed frequently in [protein structures](#). Protein motifs provide a useful way of categorizing distinctive and recurring components of protein structure. In the hierarchy of protein structure classification, the protein motif falls between that of secondary structure and [tertiary structure](#). Classification of these preferred patterns of packing secondary structure elements is somewhat arbitrary, and there is sometimes overlap between the structural terms of protein motifs and protein **domains**. However, unlike domains, protein motifs do not necessarily fold independently of the rest of the polypeptide chain, and they may not always have an intrinsic function.

Proteins having structurally similar motifs can share significant sequence identity in the equivalent regions of the polypeptide chain, even if there is very low sequence identity overall. The identification of regions of conserved sequence within proteins means that novel protein motifs can be predicted using sequence information only, in the absence of any three-dimensional structural data. The [Kringle domain](#) is an example of such a sequence motif, in that it is easily identified from

protein sequences based on three characteristic [disulfide bonds](#). Other protein motifs are purely structural and have little or no sequence conservation; examples include the **b-meander**, the **Greek key**, and the **jelly roll**.

[See also [Protein Structure](#).]

#### Suggestions for Further Reading

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

E. G. Hutchinson and J. M. Thornton (1996) PROMOTIF—A program to identify and analyze structural motifs in proteins. *Protein Sci.*, **5**, 212–220.

## Protein Secretion

Most of the [proteins](#) made by a cell are retained intracellularly; only a specialized subset is secreted outside the cell. Proteins secreted by cells have several major functions, including signaling to other cells, formation of an insoluble [extracellular matrix](#) surrounding the cell, and degradation of extracellular material. Molecular biology has helped explain how newly synthesized proteins are selected for secretion and transportation across the **hydrophobic** barrier of the plasma [membrane](#) (see **Protein biosynthesis**).

The mechanism for transporting newly synthesized proteins is highly conserved from bacteria to mammals. A key difference, however, is that bacteria translocate the proteins directly across the plasma membrane to the outside world, whereas eukaryotic cells translocate them into a specialized intracellular organelle, the [endoplasmic reticulum](#). Newly synthesized proteins are conveyed from the endoplasmic reticulum to the cell surface via a series of carrier **vesicles**. It is therefore useful to consider prokaryotic and eukaryotic protein secretion separately.

### 1. Protein Export from Bacteria

Bacteria are extremely simple from a cell biological point of view. **Gram-positive** bacteria have only a single membrane, the plasma membrane, separating their cytoplasm from the outside world. Gram-negative bacteria have a second outer membrane, separated from the inner membrane by periplasmic space. In this case, the newly synthesized protein that is to be exported must cross the outer membrane in addition to the inner plasma membrane.

The basic elements of translocation across the membrane appear to be universal ([1](#)). A small pore or hole through which unfolded proteins can pass traverses the membrane. Energy is used to push unfolded protein through the hole. **Molecular chaperone** proteins exist on the cytoplasmic side of the pore to maintain a protein in an unfolded state, which allows it to be threaded through the pore. The chaperones distinguish proteins that are to be exported from cytoplasmic proteins because the exported proteins have **signal sequences**, regions of contiguous amino acids that are usually at the amino terminus of the protein. Once the protein has crossed the membrane, additional chaperones can facilitate correct **protein folding**.

Much of the research on protein export by bacteria has utilized *Escherichia coli*, a gram-negative bacterium. In *E. coli* there are at least two different ATP-driven transport systems, one that uses the

*secA* system (see [Sec Mutants/Proteins](#)) to pump proteins into the periplasmic space (2), and a second that uses a member of the ABC (ATP-binding cassette) family of proteins to export proteins across both *E. coli* membranes simultaneously (3).

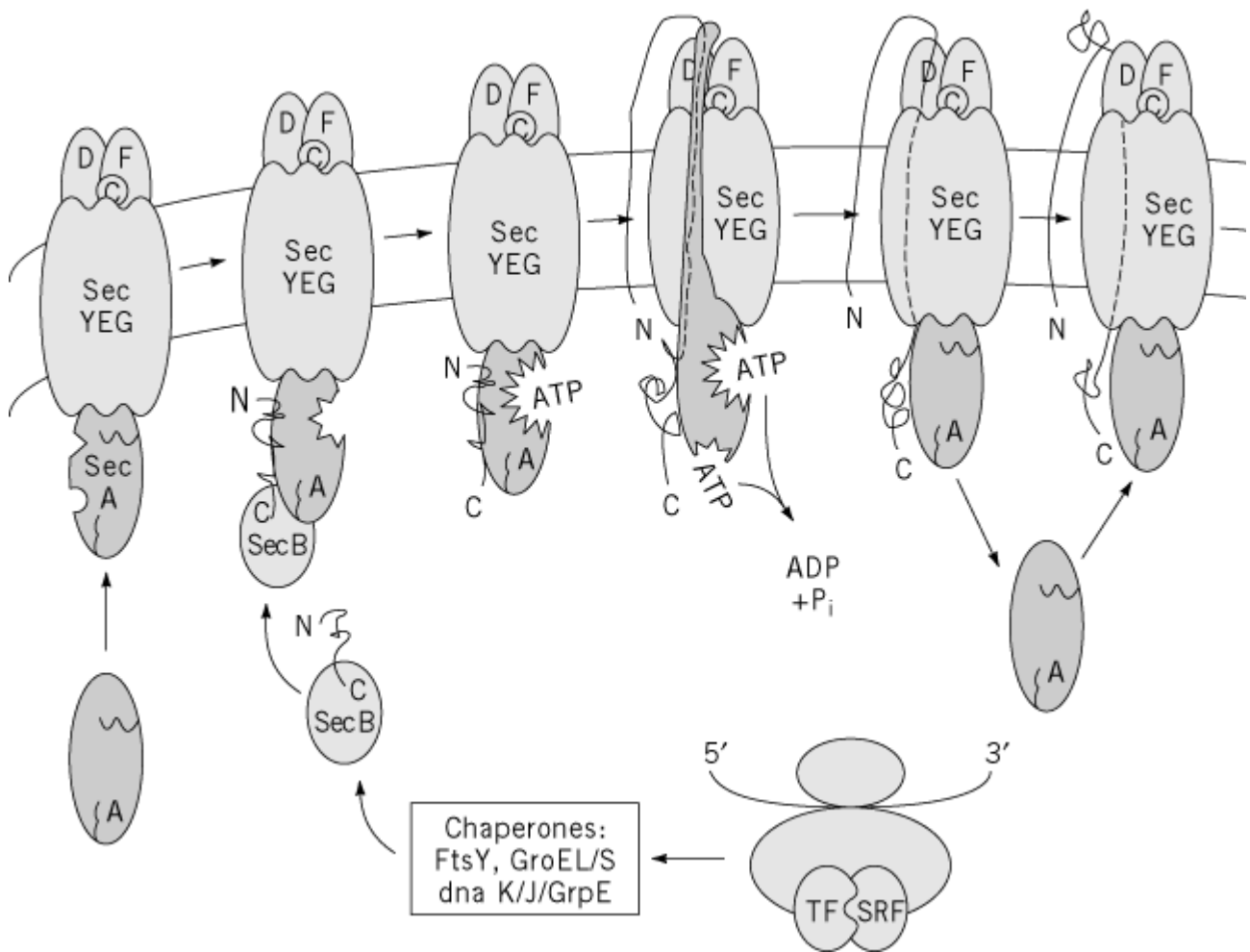
### 1.1. The *secA* System of Protein Transport

When *E. coli* makes a protein that is destined for export, it is made as a preprotein containing a signal sequence (4). In most proteins, these signal sequences are a stretch of contiguous amino acids at the amino terminus of a protein, with a positively charged amino terminus, at least six sequential hydrophobic amino acids, and a cleavage site for the signal peptidase that removes the signal after translocation. The presence of the signal sequence in a newly synthesized protein causes a protein chaperone to bind to it, retarding its folding. There appear to be at least two major chaperone systems in *E. coli*. One involves the tetrameric cytosolic protein *secB*, which binds to newly synthesized proteins that contain a signal sequence. It binds throughout the length of the newly synthesized [polypeptide chain](#), keeping it from adopting its native conformation (2). A different class of chaperones binds to the nascent chain as the secreted protein is being synthesized on the [ribosome](#). An example of such a chaperone is a ribonucleoprotein, the **signal recognition protein** (SRP). In *E. coli* SRP consists of a 4.5 S RNA molecule and a single protein, Ffh, which is a [GTPase](#); that is, it hydrolyzes guanosine triphosphate (5).

The proteins destined for export, complexed with chaperones that keep them from folding, interact with receptors on the bacterial plasma membrane to initiate translocation across the membrane. The *secB* receptor is *secA*, and the SRP receptor is another GTPase, FtsY (5). The signal sequence has therefore two functions: (1) attracting a chaperone and (2) helping to target the complex to plasma membrane receptors.

The actual translocation event is through a transmembrane pore. In *E. coli* the pore is made up of the *sec* YEG complex, consisting of three major membrane proteins (*secY*, *secE*, and *secG*) with several accessory proteins. The *sec* YEG complex forms the pore through which unfolded proteins pass on their way out of the cell. Translocation involves pushing the unfolded peptide through the pore by a mechanism that resembles a sewing machine needle or a piston (6). A peripheral membrane protein (*secA*) takes a 25-residue loop, formed by the amino-terminal signal peptide and the contiguous amino acids, and pushes the loop through the membrane. The amino-terminal amino acids remain cytoplasmic. To do this, the *secA* protein goes through a remarkable, ATP-driven, conformational change, which causes an arm of the *secA* protein to be inserted transiently into and partially across the membrane. After the pushing step, ATP is hydrolyzed and *secA* is returned to a cytoplasmic location, ready for a second translocation pumping cycle. The signal peptide is cleaved off by a “signal peptidase” located in the periplasmic space. Further translocation in 25-residue steps can be driven by additional cycles of protein activity. After the initial step, however, further translocation is facilitated by a transmembrane [proton motive force](#) (Fig. 1).

**Figure 1.** The SecA system of protein transport. Chaperones bind to a nascent protein in order to keep it an unfolded state to receptor SecA at the membrane. SecA translocates the protein across the membrane through the protein pore (SecY/E/G conformational changes). [Figure 1B from F. Duong, J. Eichler, A. Price, M. R. Leonard, and W. Wickner (1997) *Cell* 91



Bacteria may push protein across the translocation channel, rather than pull it, because the source of energy, ATP, is exclusively cytoplasmic. Once in the periplasmic space, secreted proteins form [disulfide bonds](#) and fold into **proteinase**-resistant conformations. Such folding is necessary before the proteins are transported through holes, poorly described, in the outer membrane.

### 1.2. ABC Protein-Mediated Export

Some proteins destined for export are recognized in the cytoplasm of *E. coli* by an entirely different class of signal sequences. In these cases, the sequence is at the **C-terminus** and is not cleaved off by a signal peptidase after transport. Proteins with this type of signal sequence are closely related to each other and include such proteins as [toxins](#), proteinases, and [lipases](#). The signal sequence is recognized by a different class of plasma membrane protein, ATP-driven protein translocators of the ATP-binding cassette (ABC) family ([3](#), [7](#)). ABC proteins have two cytoplasmic ATP-binding **domains** and two hydrophobic domains, with six transmembrane sequences. They can be either a single polypeptide or be made up of several polypeptides. A complex of a bacterial ABC export with two accessory proteins allows it to export a cytoplasmic protein across both the inner and outer bacterial membrane at the same time. Assembly of the transporting complex is triggered by substrate binding ([8](#)).

### 2. Protein Export from Eukaryotes

Unlike bacteria, eukaryotic cells are packed with intracellular membranes. A considerable fraction of those membranes is involved in protein export. Thus protein export in eukaryotes involves translocation across membranes but also transport of the newly synthesized protein through

intracellular compartments (9). Protein export in eukaryotes is thus a much more complex process than in bacteria. It becomes necessary to know the intracellular pathway taken by a newly synthesized protein before it reaches the eukaryotic cell surface (see also **Protein targeting**).

### 2.1. Intracellular Trafficking of Exported Protein

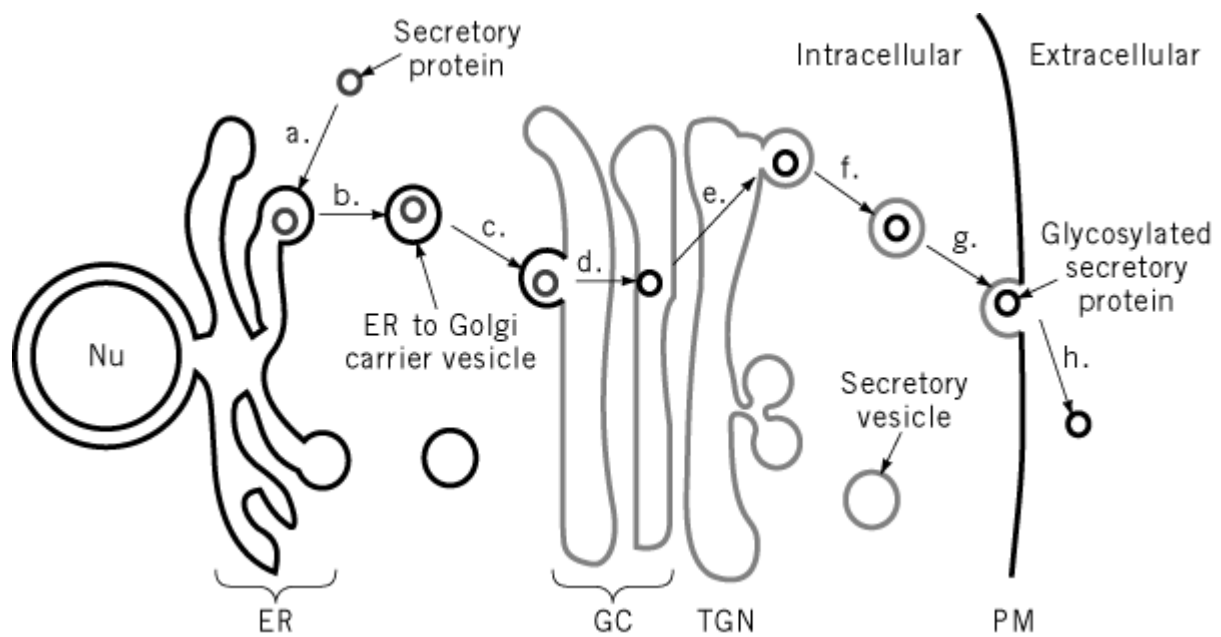
For the most part, eukaryotes do not secrete proteins directly across the plasma membrane. Instead, the newly synthesized proteins are translocated from the cytoplasm into an intracellular organelle, the **endoplasmic reticulum** (ER). The ER is a series of membranous cisternae and tubules that spread throughout the cytoplasm and is continuous with nuclear membrane. As we shall see, the mechanisms for translocating proteins across ER membranes and bacterial membranes both involve a transmembrane pore. Why have an ER? A common speculation on the evolutionary advantage of the ER is that it provides a controlled milieu in which exported proteins can fold and oligomerize, without being exposed to the rigors of the extracellular world. In this regard, the lumen of the ER resembles the **periplasmic** space in gram-negative bacteria.

A protein, correctly folded and oligomerized in the ER, must be transported to the outside of the cell. The initial step in transport from the ER is the formation of small membrane vesicles, usually called carrier vesicles, that contain the exported proteins as cargo (10, 11). Formation of carrier vesicles occurs at specialized ER regions called transitional zones, and it requires a cytoplasmic coat, to impart curvature to the vesicle membrane and to pinch it off from the donor ER (12).

After a carrier vesicle is formed, it must recognize and fuse with its target. Recognition and fusion (see **Exocytosis**) involve proteins on the vesicle (v-SNARE) and on the target membrane (t-SNARE) (13). When a carrier vesicle leaves the ER, it does not go directly to the plasma membrane but instead fuses with an organelle, the **Golgi complex**, which is a mandatory way station on the secretory pathway of eukaryotes. The Golgi complex is a stack of membranous cisternae, similar morphologically to a stack of pancakes. The carrier vesicles enter the *cis* end of a Golgi stack and exit from the *trans* side. A protein to be exported goes through a series of glycosylation steps in which six- or nine-carbon sugars are added to or removed from oligosaccharide chains attached to **serine**, **threonine** (**O-glycosylation**), or **asparagine** residues (**N-glycosylation**). The added sugars can often protect the exported protein from rapid **proteolytic** degradation after it is secreted into the extracellular world. A unique class of oligosaccharides is added to newly synthesized lysosomal enzymes that allows them to be diverted out of the secretory pathway to primary **lysosomes**.

A second budding event takes place from the *trans* region of the Golgi complex. A second coating mechanism causes the formation of a **secretory vesicle** containing glycosylated proteins for export. The post-Golgi secretory vesicles move to the plasma membrane, either by diffusion or by transport along microtubules. When the secretory vesicles reach the plasma membrane, their membranes and the plasma membrane fuse to form one continuous bilayer, a process known as **exocytosis**. The fusion step results in the release of the exported protein into the extracellular medium (Fig. 2).

**Figure 2.** The eukaryotic secretory pathway. Secretory proteins enter the membrane trafficking system by (a) translocation from the cytoplasm into the endoplasmic reticulum (ER) (b) In the ER, proteins are packaged into ER to Golgi complex (GC) carrier vesicles that (c) deliver proteins to the GC. As proteins progress through the GC they become (d) glycosylated. Once the (e) reach the *trans*-Golgi network (TGN), they (f) are sorted into secretory vesicles in which proteins are (g) transported to the plasma membrane (PM). At the cell surface secretory proteins are (h) released into the extracellular environment.



This pathway of protein export or secretion is the one followed in most organisms, including yeast, protozoa, and mammalian cells such as muscle cells and fibroblasts. It is commonly referred to as the constitutive pathway, to distinguish it from a specialized pathway found in specialized cells, such as endocrine and exocrine cells, and in cells of the **hematopoietic** system, such as neutrophils or [cytotoxic T lymphocytes](#). In such cells, newly synthesized proteins are diverted out of the normal biosynthetic pathway to be stored in secretory granules (14). Fusion of the secretory granule with the plasma membrane is usually triggered by an external stimulus. Because export of this class of secreted proteins is controlled by a stimulus, the triggered release of protein from a storage pathway has been called regulated release.

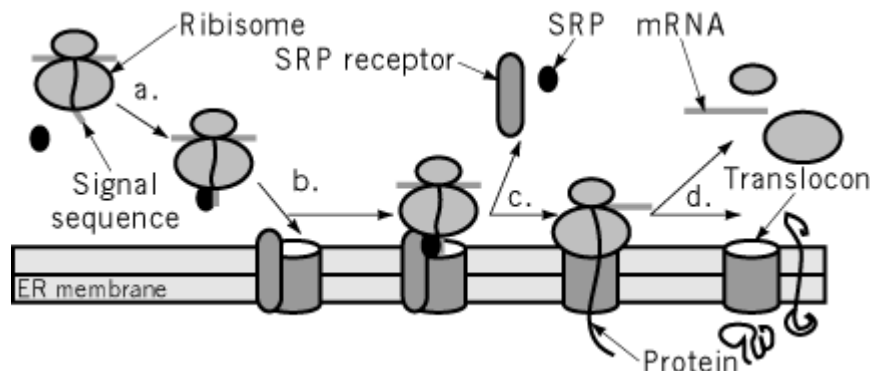
The protein export pathway of eukaryotes is highly conserved. Protozoa, yeast, and mammalian cells use essentially the same mechanisms to translocate proteins into an ER, to glycosylate and sort them in the Golgi, and to export them across the plasma membrane by exocytosis.

## 2.2. Translocation Across the ER Membrane

Translocation of newly synthesized proteins across the ER membrane shows many similarities to translocation across the plasma membrane protein of bacteria (1, 15, 16). Proteins are prevented from folding in the cytoplasm. They are fed across the plasma membrane through a translocon, a proteinaceous pore, which has three subunits very similar to the bacterial proteins made by the *secY*, *E*, and *G* genes. By [electron microscopy](#), these pores are rings about 8 to 10 nm in diameter, with a central pore of 2 nm, sufficient to allow the passage of an extended, hydrated peptide of 1.1 nm in diameter. These pores can now be recognized (17). In yeast, proteins traverse pores in the ER by two different types of translocation mechanisms. One is an ATP-driven process that translocates proteins whose synthesis is complete. The other couples translation to the translocation process. In this transport mode, the ribosome is attached to the proteinaceous transport pore, the translocon, and feeds the nascent chain through the pore as it is being synthesized. Mammalian cells only have the co-translation mode of translocation. When translocation is co-translational, the nascent chain is recognized in the cytoplasm by a [signal recognition particle](#), which stops further protein synthesis until the complex of ribosome, nascent chain, and signal recognition particle reaches the endoplasmic reticulum (Fig. 3).

**Figure 3.** Translocation across the ER membrane. (a) Signal recognition particle (SRP) binds to a signal sequence on a

nascent protein and halts the protein's synthesis by the ribosome. (b) SRP brings the protein and ribosome to the translocon on the ER membrane by binding to its receptor. (c) Protein synthesis restarts when the ribosome is correctly positioned and the protein has been threaded into the translocon. (d) Proteins can either be translocated completely into the lumen of the ER to become secretory proteins or translocated partially through the membrane to become integral membrane proteins.



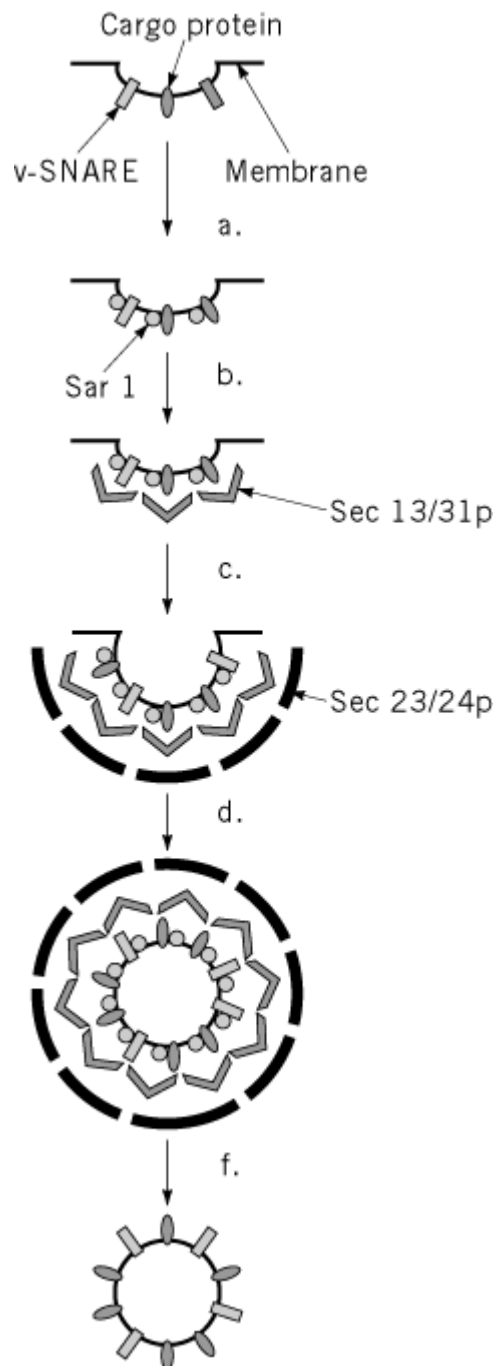
As the nascent chain enters the cytoplasm, it is glycosylated on suitable asparagine residues and forms disulfide bonds. Folding is completed by association with accessory proteins such as [BiP](#), [calnexin](#), [peptidyl prolyl \*cis/trans\* isomerases](#), and [protein disulfide isomerase](#) (18). Proteins also oligomerize in the lumen of the ER. For example, [protocollagens](#) trimerize with an extended [coiled-coil](#) configuration, accompanied by cross linking between hydroxyproline and hydroxylysine residues (see [Hydroxylation \(Lysine, Proline\)](#)). If cross linking of collagen protofibrils is inhibited, the improperly folded collagen cannot leave the ER and is quickly degraded. A similar observation is made with the secreted protein  [\$\alpha\$ 1-antitrypsin](#). In the disease,  $\alpha$ -antitrypsin deficiency, a mutation in the protein can prevent its correct folding. These two examples demonstrate that there is a quality-control system that surveys the proteins leaving the ER and prevents the export of improperly folded proteins. Recent data suggest that improperly folded ER proteins are exported from the ER back to the cytoplasm, where they are degraded by [proteasomes](#) (19).

### 2.3. Traffic of Secretory Proteins from the ER

To take correctly folded and oligomerized proteins from the ER, a vesicle forms in the transitional elements and includes proteins to be exported, but excludes resident proteins of the ER lumen, such as [BiP](#) (20). The coat that causes the vesicle to form is now known as COPII. Yeast COPII contains four subunits, *sec31p*, *sec13p*, *sec23p*, and *sec24p* (see [sec mutants](#)). Assembly of a COPII coat requires a small GTPase, *Sar1*, and a [guanine nucleotide exchange factor](#), *Sec12p*, in the ER membrane (12) (Fig. 4). The coated vesicle leaving the Golgi carries with it a complement of v-SNARE molecules (see [Exocytosis](#)) to allow it to fuse with the *cis*-Golgi network. In yeast, these are *Sec22p*, *Bos1p*, and *Bet1p*. Resident proteins such as [BiP](#) may be excluded from the lumen of the coated vesicle because they are oligomerized into complexes that are too big to enter the small vesicle. To some extent, exported proteins are those that lack a retention signal and so are not retained in the ER. Export of secreted proteins would then be by default, because they lack information to go anywhere else. There is evidence, however, that positive sorting occurs (21) (Fig. 5). In yeast, the secreted protein invertase is recognized by a membrane-bound ER protein (*Emp24p*) that is required for its transport to the Golgi (22). Furthermore, cargo proteins are concentrated as they leave the ER (23, 24). Since most soluble resident proteins in the ER lumen are not glycosylated, an attractive hypothesis is that exported proteins are recognized by a [lectin](#), which concentrates them in budding vesicles. A protein, [ERGIC-53](#), recycles between the ER and the Golgi and is a lectin with the capacity to bind the mannose residues found on newly synthesized secretory proteins (25). Proteins such as [ERGIC-53](#) might bind secreted proteins and actively carry them to the Golgi complex, in the same way that the [mannose phosphate receptor](#) carries newly formed lysosomal enzymes to the prelysosomal compartment.

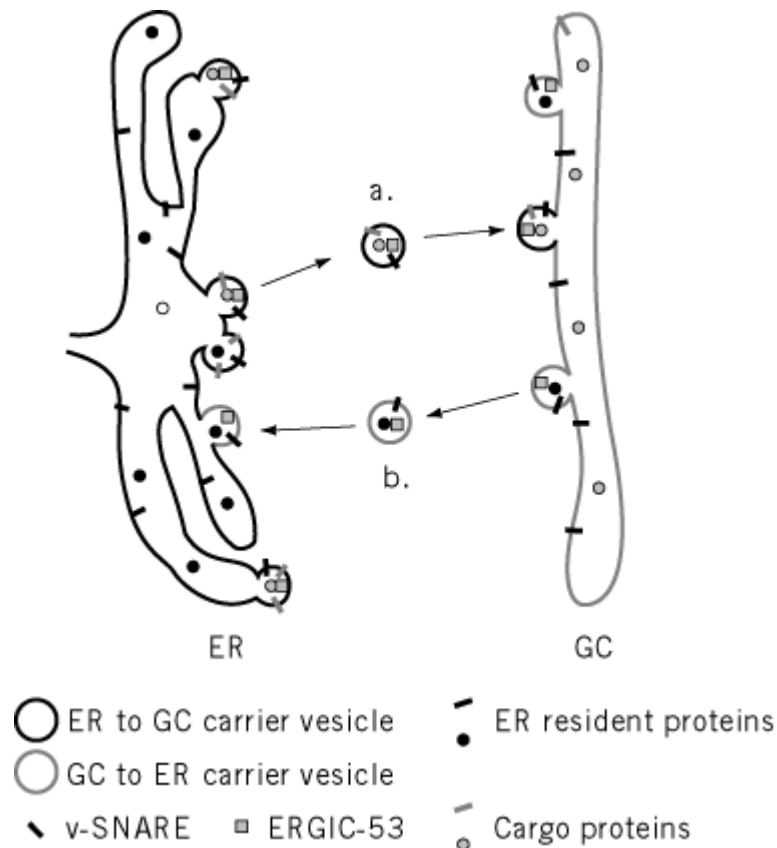


**Figure 4.** Formation of ER to GC carrier vesicles. ER to GC carrier vesicles are formed with the COPII coat (Sar1, Sec 13/31p, and Sec 23/24p). Components are thought to bind sequentially with (a) the GTPase Sar1 (in its GTP state), binding to v-SNARE and cargo molecules. (b) Sec 13/31p binds to Sar1 and then (c) Sec 23/24p binds to the whole complex. By unknown mechanisms the COPII coat (d) causes a vesicle to bud off from the ER membrane, and (e) uncoating leads to the (f) mature carrier vesicle, which can then go on to fuse with the Golgi complex. Formation of other types of vesicles is likely to involve similar steps using different coating proteins.



**Figure 5.** Sorting of proteins between the ER and GC. (a) Cargo proteins destined for secretion and v-SNAREs needed for carrier vesicle fusion to the GC are sorted away from the ER resident proteins into COPII coated vesicles. After delivery to the GC, the cargo proteins continue through the secretory pathway. (b) However, v-SNAREs needed for the ER to GC trafficking step and ER resident proteins that have been mistakenly transported to the GC can be

retrieved from the secretory pathway through retrograde transport back to the ER by COPI-coated carrier vesicle. ERGIC-53, a likely carrier protein, cycles between the ER and GC and may assist in sorting the secretory proteins from proteins that are retained in the ER.



Proteins such as ERGIC-53 and the v-SNARE proteins can recycle back to the ER from the Golgi. The vesicle for Golgi to ER flow is made up of coatamers, referred to as COPI coatamers. Coatamers consist of seven subunits. The small GTPase active during coatamer coat formation is an **ADP ribosylation** factor (ARF1).

Fusion of the ER-derived carrier vesicle with the Golgi membranes requires the v-SNAREs *Sec22p*, *Bos1p*, and *Bet1p* in the vesicle membrane, and the cognate *cis*-Golgi t-SNARE, the *sed5* protein. In addition, a second class of small GTPases are required that are members of the *rab* family (26). The *rab* GTPases are *ras*-like proteins that are specific to a certain organelle or a certain trafficking step. In yeast, mutations in the *rab* protein *ypt1* inhibit fusion with the Golgi and cause the accumulation of carrier vesicles. The *ypt1* protein appears to play a role in v-SNARE/t-SNARE interactions (see [Exocytosis](#)) involved in *cis*-Golgi targeting. Two of the v-SNAREs in yeast, *Bos1p* and *Sec22p*, normally form a complex. The complex does not form if *ypt1* is absent, nor is a v-SNARE/t-SNARE complex formed, suggesting that *ypt1* plays a role in activating the v-SNAREs, allowing a v-SNARE/t-SNARE complex to form. It is also possible, however, that *ypt1* is also activating the t-SNARE complex on the target membrane, by removing its chaperone, *Sly1p*. This chaperone is a member of the *Sec1p* family of syntaxin-family Chaperones (see [Exocytosis](#)).

#### 2.4. Modification of Secreted Proteins in the Golgi

As secreted proteins pass through the Golgi complex, they can be modified by glycosylation, [sulfation](#), or [proteolysis](#). Three major types of glycosylation occur, N-linked, O-linked, and glycosaminoglycan addition. In mammalian cells, the mannose-rich, asparagine-linked oligosaccharides that had been added in the ER are first trimmed in the Golgi complex and then rebuilt. Mannoses are first removed and then replaced by other sugars. On a stump of the original

mannose tree, a multichain oligosaccharide is generated by the orderly addition of the sugars *N*-acetylglucosamine, galactose, and sialic acid (*N*-acetylneuraminic acid). The previous form of glycosylation is referred to as **N-glycosylation**, since the sugars are added to an asparagine residue. If a sugar chain is added to a serine or threonine residue, it is referred to as **O-glycosylation**. The first sugar added is usually *N*-acetylgalactosamine. The oligosaccharide chains are branched but are usually shorter and more variable than N-linked oligosaccharides.

The third type of glycosylation results in the formation of proteoglycans. Proteoglycans are formed by the addition of xylose to serine residues as the protein moves from the ER to the Golgi (27), followed by the addition of a long, highly charged, unbranched oligosaccharide called a glycosaminoglycan. Glycosaminoglycans are polymers of a disaccharide unit, one of which is usually uronic acid. The other member of the oligosaccharide pair varies, depending on the type of glycosaminoglycan. Part of the strong negative charge on the proteoglycans is due to the addition of sulfate groups to the glycosaminoglycans. Sulfate can also be added directly to **tyrosine** residues on secreted proteins. Sulfation of proteins only takes place in the late Golgi. Radioactive sulfate is thus commonly used to label exported proteins selectively and to identify a set of proteins at a late Golgi step of transport. Glycosaminoglycan chain synthesis can be initiated in the late Golgi by growing cells in the presence of a membrane-permeable derivative of xylose (28). This technique has also been used to identify membrane traffic from the late Golgi to the surface of the cell.

Several proteins, membrane associated and secreted, are subject to proteolytic processing as they pass through the Golgi complex. The proteolysis is highly specific, occurring usually to the carboxyl side of a four-residue consensus sequence: Arg–X–(Lys/Arg)–Arg (29). The Golgi enzyme that carries out this proteolysis is a membrane-attached form of the **subtilisin** family, called furin. Secreted proteins cleaved by furin include the hepatocyte and **nerve growth factors** and **pro-albumin**.

Exported proteins can take one of two pathways to leave the Golgi complex, the constitutive and the regulated [Fig. 4 in (14)]. In the regulated pathway, newly synthesized proteins are stored in specialized secretory granules until cells receive a signal to secrete. Storage does not occur in the constitutive pathway; nascent proteins go directly from the Golgi to the cell surface. Little is known about sorting into constitutive carrier vesicles, the coat required to generate a carrier vesicle, or the v-SNAREs and t-SNAREs required for fusion with the cell membrane. Much more is known about the regulated secretory pathway, which packages newly synthesized proteins into secretory granules (qv). Similarly, study of the regulated secretory pathways has revealed a great deal about the final step, fusion with the plasma membrane, exocytosis.

## 2.5. Nonclassical Protein Secretion

The *secA*-mediated pathway of bacterial protein export uses a translocon or transmembrane pore that is closely similar to the ER translocon in eukaryotes. There is also a eukaryote equivalent of the ABC transporter system in bacteria. The yeast protein, Ste6p, is an ABC protein that exports a-factor, a small yeast mating-type peptide, from the cells. There is no evidence, however, that eukaryotic ABC proteins can translocate full-length proteins across the plasma membrane.

Yeast can export proteins without using the classical, ER to Golgi pathway. Mutations in the classical pathway that block invertase secretion have no effect on this pathway, called the nonclassical pathway. The gene for a transmembrane protein, *nce2*, that is important for nonclassical secretion has been **cloned** and sequenced, but its role in protein export remains unclear (30).

Mammalian cells export proteins, such as **interleukin-1b**, basic **fibroblast growth factor**, and **thioredoxin**, that do not have a classical signal sequence (31). In addition, drugs such as **monensin** and **brefeldin A** that normally block classical secretion are without effect on these proteins. Their secretion has been linked to the **heat-shock** response. Little is currently known about the molecular mechanisms that are involved.

## Bibliography

1. G. Schatz and B. Dobberstein (1996) *Science* **271**, 1519–1526.
2. W. Wickner and M. R. Leonard (1996) *J. Biol. Chem.* **271**, 29514–29516.
3. R. Binet, S. Letoffe, J. M. Ghigo, P. Delepelaire, and C. Wandersman, (1997) *Gene* **192**, 7–11.
4. F. Duong, J. Eichler, A. Price, M. R. Leonard, and W. Wickner (1997) *Cell* **91**, 567–573.
5. J. D. Miller, H. D. Bernstein, and P. Walter (1994) *Nature* **367**, 657–659.
6. A. Economou, J. A. Pogliano, J. Beckwith, D. B. Oliver, and W. Wickner (1995) *Cell* **83**, 1171–1181.
7. F. Duong, A. Lazdunski, and M. Murgier (1996) *Mol. Microbiol.* **21**, 459–470.
8. R. Binet, S. Letoffe, J. M. Ghigo, P. Delepelaire, and C. Wandersman (1997) *Gene* **192**, 7–11.
9. G. Palade (1975) *Science* **189**, 347–358.
10. J. E. Rothman (1994) *Nature* **372**, 55–63.
11. M. J. Kuehn and R. Schekman (1997) *Curr. Opin. Cell Biol.* **9**, 477–483.
12. R. Schekman and L. Orci (1996) *Science* **271**, 1526–1533.
13. J. E. Rothman (1996) *Prot. Sci.* **5**, 185–194.
14. T. L. Burgess and R. Kelly (1987) *Annu. Rev. Cell Biol.* **3**, 243–293.
15. A. K. Corsi and R. Schekman (1996) *J. Biol. Chem.* **271**, 30299–302302.
16. P. Walter and A. E. Johnson (1994) *Annu. Rev. Cell Biol.* **10**, 87–119.
17. D. Hanein, K. E. Matlack, B. Jungnickel, K. Plath, K. U. Kalies, K. R. Miller, T. A. Rapoport, and C. W. Akey (1996) *Cell* **87**, 721–732.
18. U. Tatu and A. Helenius (1997) *J. Cell Biol.* **136**, 555–565.
19. R. R. Kopito (1997) *Cell* **88**, 427–430.
20. R. D. Teasdale and M. R. Jackson (1996) *Annu. Rev. Cell Develop. Biol.* **12**, 27–54.
21. J. E. Rothman and F. T. Wieland (1996) *Science* **272**, 227–234.
22. F. Schimmoller, B. Singer-Kruger, S. Schroder, U. Kruger, C. Barlowe, and H. Riezman (1995) *EMBO J.* **14**, 1329–1339.
23. W. E. Balch, J. M. McCaffery, H. Plutner, and M. G. Farquhar (1994) *Cell* **76**, 841–852.
24. J. F. Presley, N. B. Cole, T. A. Schroer, K. Hirschberg, K. J. Zaal, and J. Lippincott-Schwartz (1997) *Nature* **389**, 81–85.
25. C. Itin, A. C. Roche, M. Monsigny, and H. P. Hauri (1996) *Mol. Biol. Cell* **7**, 483–493.
26. S. R. Pfeffer (1996) *Annu. Rev. Cell Develop. Biol.* **12**, 441–461.
27. B. M. Vertel, L. M. Walters, N. Flay, A. E. Kearns, and N. B. Schwartz (1993) *J. Biol. Chem.* **268**, 11105–11112.
28. R. A. Chavez, S. G. Miller, and H. P. Moore (1996) *J. Cell Biol.* **133**, 1177–1191.
29. J. B. Denault and R. Leduc (1996) *FEBS Lett.* **379**, 113–116.
30. A. E. Cleves and R. B. Kelly (1996) *Curr. Biol.* **6**, 276–278.
31. A. E. Cleves (1997) *Curr. Biol.* **7**, R318–320.

## Protein Sequencing

The introduction of protein sequence determination started the era of macromolecular

characterization. It established the homogeneity of [proteins](#), permitted analytical automation, initiated the path to increasing speed and sensitivity, and was instrumental also in the interpretation of three-dimensional [protein structure](#) analysis and, later, DNA [sequence analysis](#). Both the speed and sensitivity of chemical protein sequence determination have increased enormously. Initially, each analytical step required days or more in time and millimoles in amounts, but now only minutes and picomoles or less are necessary. At the same time, the major chemical method, the [Edman Degradation](#), has stayed more or less unchanged. It now constitutes one of the oldest molecular methods that is still much in use. Recently, however, even this method has been challenged by other approaches, [mass spectrometry](#) in particular, but further chemical sequencers are also still being launched. Probably, both chemical and mass spectrometric protein sequence analysis will remain for a considerable time. However, protein sequence determination is not just a question of the degradative analysis; it also involves sample preparation, correlation with DNA sequences, including [genome](#) sequencing projects, the screening of [databases](#), [NMR](#) and [X-ray crystallography](#), as well as [epitope](#) analysis and [protein structure predictions](#), all of which constitute different aspects of integrated, modern molecular biology. In short, sequence determination is part of the era of the [proteome](#) that is now emerging at the protein identification stage of the genomic sequencing breakthrough. The various stages of protein analysis using chemical degradation and mass spectrometry are presented here, together with the current stage and perspectives of future methodology.

## 1. Development of sequence determination

Sequence determination of proteins has been possible for about 50 years ([1-5](#)). During the second half of the twentieth century, development has been rapid, bringing structural analysis from something requiring the effort of a whole department, a small protein, large amounts, several years of work, and still some uncertainty, to something requiring just a single person, any size of protein, minute amounts, a limited time, and great reliability. This progress has been paralleled by methodological and instrumental developments in essentially nine steps (Table [1](#)).

**Table 1. Important Steps in the Development of Methods and Instruments in Protein Sequence Analysis**

---

|   |
|---|
| 1. The first analysis and recognition of defined structures (1940s to 1950s)                        |
| 2. The Edman reaction (~1950)   |
| 3. Manual methods: “Dansyl-Edman” and later follow-ups (DABITC) (from ~1960 on)                     |
| 4. Automation: the protein sequenator or sequencer (1967 on)  |
| 5. Development of HPLC (~1977)  |
| 6. Polybrene, membranes, valves, conversion, on-line HPLC: 2nd generation sequencers (end of 1970s) |
| 7. Mass spectrometry (1970s)  |
| 8. Third-generation sequencers, including C-terminal degradation (end of 1990s)                     |
| 9. Electrospray, nanospray, present-generation “Protein mass spectrometers” (end of 1990s)          |

---

Each one of these steps has meant a corresponding leap in analytical speed, reproducibility,

sensitivity, and ease of analysis. The successive development is clearly traceable in the scientific literature, and in this case especially in the Proceedings volumes of the series of the special methods conferences, MPSA (for Methods in Protein Sequence Analysis, 1975–1992, and Methods in Protein Structure Analysis from 1994 to the present) that originated from still earlier meetings that roughly coincided with the introduction of automatic sequencers (6, 7). MPSA continued for a long time as the leading methodological meeting, but other conferences also covering the subject now compete in progress reports, in particular the Protein Society (which started in 1986) and its journal, *Protein Science* (started in 1992), and the ABRF (Association of Biomolecular Resource Facilities, from 1988; *The Journal of Biomolecular Techniques*, electronic since 1997). ABRF and MPSA still concentrate on the methodological developments, with ABRF being the larger source. Details of the MPSA meetings and their proceedings are given in Table 2. In retrospect, most of the major advances can be followed in these proceedings.

**Table 2. MPSA Proceedings**

| MPSA<br>Number, | Proceedings Volume |  |
|-----------------|--------------------|--|
|                 | Year               | Publisher (editor), Year of Publication                      |
| I,              | 1975               | Pierce (Laursen), 1975                                       |
| II,             | 1977               | Elsevier (Previero and Coletti-Previero), 1977               |
| III,            | 1979               | Elsevier (Birr), 1980  |
| IV,             | 1981               | Humana (Elzinga), 1982                                       |
| V,              | 1984               | (Not Published) (Walker)                                     |
| VI,             | 1986               | Humana (Walsh), 1987   |
| VII,            | 1988               | Springer (Wittmann-Liebold), 1989                            |
| VIII,           | 1990               | Birkhäuser (Jörnvall et al.), 1991                           |
| IX,             | 1992               | Plenum (Imahori and Sakiyama), 1993                          |
| X,              | 1994               | Plenum (Atassi and Appella), 1995                            |
| XI,             | 1996               | J. Prot. Chem. 16 (van der Rest and Vandekerckhove),<br>1997 |
| XII,            | 1998               | To be held (Kyriakidis and Choli-Papadopoulou)               |

## 2. Methodology and instrumentation

### 2.1. Initial Analysis of a Protein Sequence and Definition of the Concept

During the end of the 19th century and early parts of the 20th century, proteins were gradually purified, giving insight into the fact that they constitute defined molecules of exact structure. During the same time, the constituent [amino acids](#) were characterized, culminating with threonine in the mid-1930s (8). By the use of many methods, including partial acid hydrolysis, and reaction with a labeling reagent that is stable during hydrolysis and specific for the protein [amino group](#), fluorodinitrobenzene (FDNB) (9), the first protein [primary structure](#), that of [insulin](#), was determined in the late 1940s and early 1950s (1, 2), for which Frederick Sanger received the 1958 Nobel Prize in Chemistry. With this sequence determination, proteins were recognized as defined molecules, marking the start of the modern era of structural analysis.

## 2.2. Edman Degradation and Protein Sequencers

The Edman degradation and its use in automatic protein sequencers is described in the entry [Edman Degradation](#).

## 2.3. Mass Spectrometry

The nonvolatility of peptides was the prime limitation to the use of [mass spectrometry](#) (MS) with proteins. Many advances changed this recently. One especially important step was the introduction of triple quadrupole instruments ([10](#)), essentially tandem mass spectrometers (MS/MS), where a first quadrupole, the first “mass spectrometer” or MS<sub>1</sub>, performs peptide mass selection; a second quadrupole acts as a collision gas cell where peptide fragmentation can occur; and the third quadrupole, the second mass spectrometer, or MS<sub>2</sub>, analyzes the peptide fragments generated in the collision cell by their mass, thus allowing peptide sequence analysis by collision-induced dissociation or collision-activated dissociation mass spectrometry ([11](#)).

Another crucial step in the MS sequencing approach was the introduction of novel ion-producing techniques, to make MS accessible to proteins in general (see [Mass Spectrometry](#)). Desorption methods of ionization such as fast atom bombardment ([12](#)), plasma desorption ([13](#)), and [matrix-assisted laser desorption ionization](#) (MALDI) ([14](#)), made it possible to transfer peptides and proteins into the gas phase, making them available to mass spectrometry. MALDI, coupled to time-of-flight (TOF) instruments ([14](#)), made mass spectrometers easy to handle and brought the technique within the economic means of many protein laboratories. Hence, MALDI-TOF instruments are now used in many laboratories, but mostly for peptide mass measurements, for peptide identification and screening, rather than sequence analysis. With further developments, however, such as the introduction of reflectrons ([15](#)), delayed extraction ([16](#)), and post-source decay ([17](#)), these instruments can also give some sequence information via peptide cleavages. However, more important for sequence analysis, and the real breakthrough via MS/MS instruments, was the introduction of [electrospray ionization](#) (ESI) ([18](#)). This opened the way to sample introduction via solvents, and hence to on-line applications, with the MS analysis immediately following a [chromatography](#) separation. Similarly, further mass analysis systems, including ion traps ([19](#)), made additional progress in some MS/MS techniques. Finally, full-scale computerization and on-line databank screenings allowed further speed and interpretations ([20](#)). It is possible that MS may become the method that eventually replaces the Edman method.

## 2.4. Present-Generation Chemical Sequencers

Recent chemical developments of automatic protein sequencers have involved primarily further miniaturization. The most novel type of chemical sequencers now have phenylthiohydantoin (PTH) detection columns in the sub-millimeter diameter range, UV detectors with volumes in the nanoliter range, application possibilities in the sub-picomole range, and overall sensitivities at that level. An important complement has also been developed on the side of sample preparation, with instruments now [blotting](#) proteins and peptides in nanoliter volumes from chromatography separations directly onto [blotting matrices](#) for subsequent sequencer analysis ([21](#)).

## 2.5. C-Terminal Sequencing

Other chemical advancements have involved development of sequence determination from the C-terminus. The principle has long been known, using isothiocyanate degradation ([22](#)). The conditions were harsh, however, using strong acids, acetic anhydride, and repeated activation of the C-terminal [carboxyl group](#) in each step. Consequently, the yields were low, side-reactions plentiful, and erroneous peptide bond cleavages frequent. In addition, the secondary amine of [proline](#) residues could not react at all, and yields for carboxyl and hydroxyl residues were especially poor. Hence, degradative methods from the C-terminus have not been used much, and commercial C-terminal sequencers did not exist, in spite of all the progress with the N-terminal (see [Edman Degradation](#)). Instead, C-terminal analysis long relied on hydrazinolysis (now outdated) ([23](#)), reductive methods to get the corresponding alcohol (also outdated) ([24](#)), and, in particular, enzymatic approaches, utilizing [carboxypeptidases](#) ([25](#)). Initially, the available carboxypeptidases had too strict substrate specificities

to be useful for protein analysis in general. Lately, however, carboxypeptidases have appeared with wide specificities, and hence good applicability to proteins in general (26-28). In conjunction with mass spectrometric analysis in MALDI-TOF instruments, this has opened a new route to small-scale C-terminal analysis in the “ladder sequencing mode” (see below); (29, 30).

Progress has also been substantial in chemical aspects of C-terminal degradations of proteins. Several adjustments allowed the method to work for a few cycles in most cases (31). Recently, further modifications allow, in one step, simultaneous cleavage of the C-terminal residue and activation of the next, by the use of an additional chemical modification involving S-alkylation of the thiohydantoin. This has now made it possible to degrade several proteins in reasonable yield, and in several cases to follow the C-terminal sequence for up to 10 cycles (32). This new approach has just become available in commercial instruments. Although C-terminal degradation is still less sensitive than N-terminal degradation and less reliable, does not reach as far, has difficulties with some residues (carboxylic and hydroxylic), and is still not feasible at all with proline residues, C-terminal analysis is now becoming practical. Just a few cycles of C-terminal sequence information are sometimes sufficient, especially for identification of the correct recombinant proteins and corresponding gene constructs; for this purpose, there are now both C- and N-terminal instruments.

In conclusion, present-day chemical sequencers have reached the sub-picomole range for N-terminal analyses, and degradations from both ends on a larger scale. The chemistry for C-terminal degradations has started to develop. For N-terminal degradations, miniaturization has progressed substantially, with sensitivity increased about  $10^8$ -fold and speed about  $10^2$ -fold in ~50 years, using the same basic chemistry of Edman degradation.

## 2.6. Present-Generation “Protein Mass Spectrometers”

Substantial progress has also been made in mass spectrometry. Electrospray ionization has been miniaturized. The use of small capillaries in the “microspray” mode (33) or “nanospray” mode (34) transfers ions into the gas phase more efficiently and increases the sensitivity, as do the infusion of small volumes at low flow rates (10 nL/min), and signal averaging (33, 34).

Similarly, introduction of ion traps has allowed the more efficient collection of ions, increasing the sensitivity. Furthermore, use of TOF mass analyzers as the second MS of MS/MS instruments (35) has increased both the speed and the accuracy (mass accuracy of 0.1 Da and sensitivity at the attomole scale) to new levels in instruments just released. Such instruments for sequence analysis, coordinated with MALDI-TOF instruments for peptide mass determination and screening, are easy to run and require less work than ordinary mass spectrometers. Together, they bring MS to a stage beyond that of chemical protein sequencers, and will perhaps one day overtake the whole sequencing market at a future step. For the moment, TOF combination instruments on the mass spectrometry side, plus the current chemical sequencers, have brought protein sequence determination to a new stage of perfection, speed, and sensitivity.

## 3. Perspectives and further methodology

Three levels of future protein sequencing may be predicted: “conventional approaches, tissue characterizations, and further correlations.

### 3.1. “Conventional” Approaches, with Protein Purification and a Column End-Step

For this approach, MS/MS mass spectrometers are now in routine use at the femtomole scale, and miniaturized modern chemical sequencers work in the sub- or low picomole range. Regarding [mass spectrometry](#), the nanospray approach to electrospray ionization has made it routine to analyze proteins in solution. Similar approaches with chemical sequencers and with sample preparation using [capillary zone electrophoresis](#) (36) or micro-HPLC (18) mean that sample preparation will continue to make progress along with both MS and chemical sequencer instrumentation. This combination makes it possible to analyze virtually any protein prepared by column chromatography. This development is expected to continue, with further miniaturization, plus on-line shortcuts. Regarding



sensitivity, however, we may just have passed half or even more of the major leaps in this “conventional” mode of analysis! Half a century of progress has seen a sensitivity increase of about  $10^8$ , down to about  $\sim 10^{-14}$  moles, which is about half-way down through the magnitudes to one atom (the inverse of Avogadro's number,  $0.16 \times 10^{-23}$ ). In other words, we might expect to reach the ultimate one-molecule sensitivity level before or just after another similarly large leap. Also, in some special approaches of protein detection and analysis, using fluorescence correlation spectroscopy, science can already approach the one-molecule level of analysis (37). Although much still remains to make it all routine, to increase speed, and to lower the cost, it is possible that something like half of the major leaps have already been seen in the “conventional” approach of sequence analysis via column preparations of individually purified proteins.

### 3.2. Tissue Characterizations Using Two-Dimensional Gel Separations for Sample Preparation

This is one of two other types of analysis that offer great possibilities. With proper care, [two-dimensional gel electrophoresis](#) can now separate almost thousands of proteins into distinct positions (38). The patterns obtained can be stored, compared, and analyzed for both relative and absolute changes. In this manner, all major proteins of any tissue can be rapidly screened. Comparisons between patterns from normal or diseased organs allow direct detection of tumor markers, signal proteins, and other special forms. Similarly, gene expression in different tissues, and at different ages, allows conclusions to be drawn about developmental and differentiation patterns. In all cases, the corresponding protein spots, including all those that can be detected by ordinary or silver-staining methods (see [Silver Stain](#)), can be recovered and identified by sequence analysis.

The present procedure involves recovery of the corresponding gel piece, proper washing to allow subsequent analysis, drying, addition and penetration of a **proteinase** (usually **trypsin**), digestion of the protein in the gel (39), and recovery of the fragments, most conveniently via subsequent mass spectrometry. The MS can be performed simply by accurate mass determination of all the peptide fragments obtained from the action of the proteinase; this is usually sufficient to identify a known protein from a sequence database. If this tryptic fragment mass analysis is not sufficient for identification, the MS/MS instrument can simply be set instead to analyze directly for sequence (via use of the collision cell) of all those fragments that are not identified by just their masses. Notably, this entire analysis is theoretically possible from just one spot recovered from a two-dimensional gel. Several such approaches are listed in the latest MPSA proceedings volumes (eg (40, 41)), and are within reach of most protein analysis centers (42). It is to be expected that all major proteins in most tissues will soon have been identified as to their nature and function.

### 3.3. Further Correlations of Databank Technology, Protein Analyses, and Separation Modes

The great variety of [post-translational modifications](#) can also be identified through protein sequencing methods. Two-dimensional gel separations are excellent also for separation of multiple forms of the same protein, independent of whether the multiplicity is derived from additions/deletions (size differences, noticeable in one of the two directions) or from changes in charge (the other direction). It should also be noted that mass spectrometry can be used to detect also noncovalent associations, simply by using proper energy levels in the ionization steps of ion production.

## 4. Summary

In short, much of the proteome era to come will use the separation, analysis, and computer methods now available to characterize all cellular proteins functionally through sequence analysis. Because of the large equipment involved, much of this era may be expected to occur in “protein analysis centers” rather than in individual laboratories or groups. Still, continuous miniaturization and small-scale approaches like the “ladder” modes mentioned above (28, 29) may also keep individual laboratories and groups genuinely in-scale and fully contributing.

## Bibliography

1. F. Sanger and H. Tuppy (1951) *Biochem. J.* **49**, 463–490.
2. F. Sanger and E. O. P. Thompson (1953) *Biochem. J.* **53**, 353–374.
3. P. Edman (1950) *Acta Chem. Scand.* **4**, 283–293.
4. P. Edman (1953) *Acta Chem. Scand.* **7**, 700–701.
5. P. Edman (1956) *Nature* **177**, 667–668.
6. P. Edman and G. Begg (1967) *Eur. J. Biochem.* **1**, 80–91
7. R. A. Laursen (1971) *Eur. J. Biochem.* **20**, 89–102.
8. R. H. McCoy, C. E. Meyer, and W. C. Rose (1935) *J. Biol. Chem.* **112**, 283–302.
9. F. Sanger (1945) *Biochem. J.* **39**, 507–515.
10. R. A. Yost and C. G. Enke (1978) *J. Am. Chem. Soc.* **100**, 2274–2275.
11. F. W. McLafferty, P. F. Bente, III, R. Kornfeld, S.-C. Tsai, and I. Howe (1973) *J. Am. Chem. Soc.* **95**, 2120–2129.
12. M. Barber, R. S. Bordoli, R. D. Sedgwick, and A. N. Tyler (1981) *J. Chem. Soc. Chem. Commun.* 325–327.
13. R. D. Macfarlane and D. F. Torgerson (1976) *Science* **191**, 920–925.
14. M. Karas and F. Hillenkamp (1988) *Anal. Chem.* **60**, 1299–2301.
15. B. A. Mamyryin, V. J. Karataev, D. V. Shmikk, and V. A. Zagulin (1973) *Sov. Phys. JETP.* **37**, 45–48.
16. R. S. Brown and J. J. Lennon (1995) *Anal. Chem.* **67**, 1998–2003.
17. B. Spengler, D. Kirsch, R. Kaufmann, and E. Jaeger (1992) *Rapid Commun. Mass Spectrom.* **6**, 105–108.
18. J. B. Fenn, M. Mann, C. K. Meng, S. F. Wong, and C. M. Whitehouse (1989) *Science* **246**, 64–70.
19. R. G. Cooks and R. E. Kaiser (1990) *Accounts Chem. Res.* **23**, 213–219.
20. J. R. Yates (1998) *J. Mass Spectrom.* **33**, 1–19.
21. K.-L. Hsi, M. L. Kochersperger, W. E. Werner, H. Ly, S. Sandell, and P. M. Yuan (1995) *Protein Science* **4**, Suppl. 2, 150, 540-M.
22. P. Schlack and W. Kumpf (1926) *Z. Physiol. Chem.* **154**, 125–170.
23. S. Akabori, K. Ohno, and K. Narita (1952) *Bull. Chem. Soc. Japan*, **25**, 214–218.
24. C. Fromageot, M. Jutisz, D. Meyer, and L. Penasse (1950) *Biochim. Biophys. Acta* **6**, 283–289.
25. J. I. Harris and C. H. Li (1955) *J. Biol. Chem.* **213**, 499–507.
26. R. Hayashi (1977) *Meth. Enzymol.* **47**, 84–93.
27. H. Tschesche (1977) *Meth. Enzymol.* **47**, 73–84.
28. R. Walter, W. H. Simmons, and T. Yoshimoto (1980) *Mol. Cell. Biochem.* **30**, 111–127.
29. D. H. Patterson, G. E. Tarr, F. E. Regnier, and S. A. Martin (1995) *Anal. Chem.* **67**, 3971–3978.
30. V. Bonetto, A.-C. Bergman, H. Jörnvall, and R. Sillard (1997) *Anal. Chem.* **69**, 1315–1319.
31. J. M. Bailey, F. Nikfarjam, N. R. Shenoy, and J. E. Shively (1992) *Protein Science* **1**, 1622–1633.
32. V. L. Boyd, M. Bozzini, G. Zon, R. L. Noble, and R. J. Mattaliano (1992) *Anal. Biochem.* **206**, 344–352.
33. M. R. Emmett and R. M. Caprioli (1994) *J. Am. Soc. Mass Spectrom.* **5**, 605–613.
34. M. Wilm, A. Shevchenko, T. Houthaeve, S. Breit, L. Schweigerer, T. Fotsis, and M. Mann (1996) *Nature* **379**, 466–469.
35. H. R. Morris, T. Paxton, A. Dell, J. Langhorne, M. Berg, R. S. Bordoli, J. Hoyes, and R. H. Bateman (1996) *Rapid Commun. Mass Spectrom.* **10**, 889–896.

36. A.-C. Bergman and T. Bergman (1996) *FEBS Lett.* **397**, 45–49.
37. L. Edman, Ü. Mets, and R. Rigler (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6710–6715.
38. H. Leffers, K. Dejgaard, B. Honore, P. Madsen, M. S. Nielsen, and J. E. Celis (1996) *Electrophoresis* **17**, 1713–1719.
39. U. Hellman, C. Wernstedt, J. Gónez, and C.-H. Heldin (1995) *Anal. Biochem.* **224**, 451–455.
40. R. Aebersold, L. N. Amankwa, H. Nika, D. T. Chow, E. J. Bures, H. D. Morrison, D. Hess, M. Affolter, and J. D. Watts (1995) In *Methods in Protein Structure Analysis* (M. Z. Atassi and E. Appella, eds.) Plenum, pp. 3–14.
41. K. Gevaert, H. De Mol, J.-L. Verschelde, J. Van Damme, S. De Boeck, and J. Vandekerckhove (1997) *J. Prot. Chem.* **16**, 335–342.
42. A.-C. Bergman, C. Linder, K. Sakaguchi, M. Sten-Linder, A. A. Alaiya, B. Franzén, M. C. Shoshan, T. Bergman, B. Wiman, G. Auer, E. Appella, H. Jörnvall, and S. Linder (1997) *FEBS Lett.* **417**, 17–20.

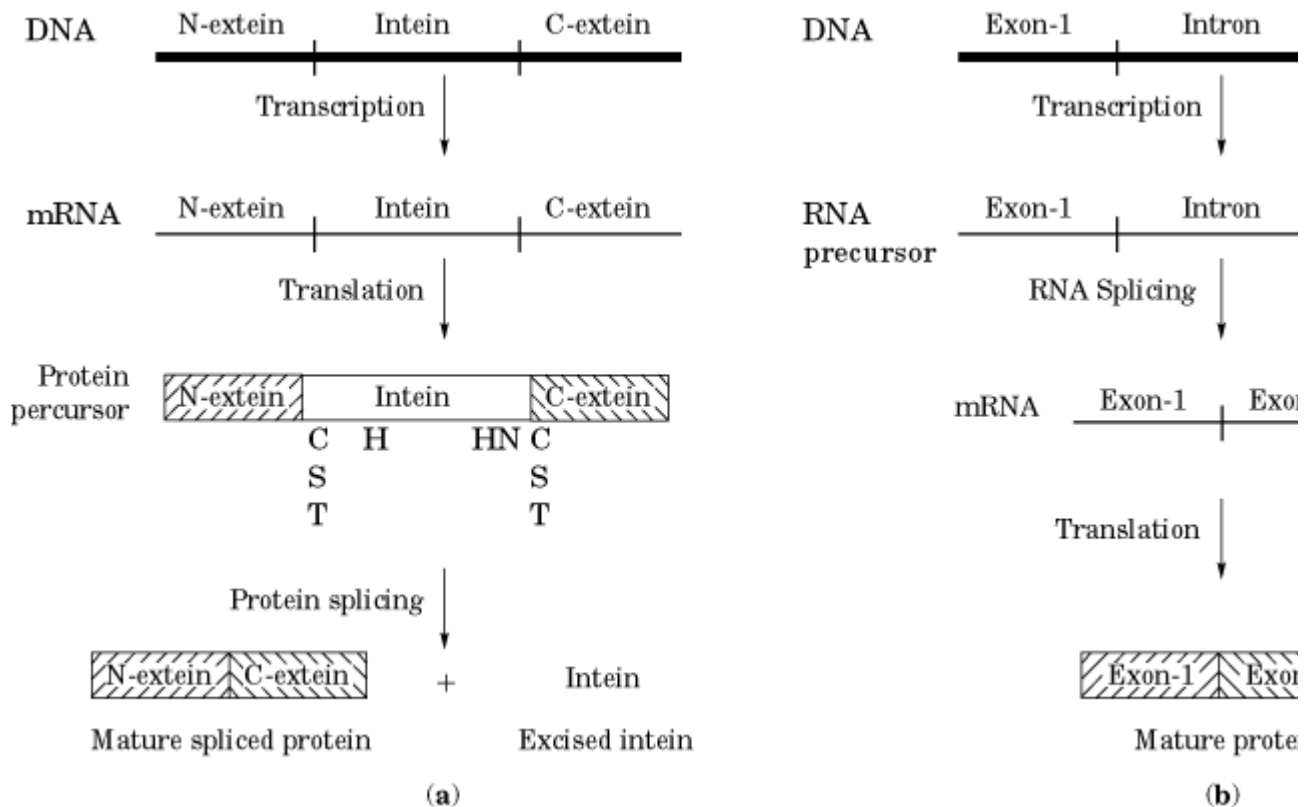
### Suggestions for Further Reading

43. *Methods in Protein Structure Analysis*. Proceedings from the conference series with the same name, usually published biannually. Latest issue (*J. Prot. Chem.* **16**, 1997) from the XI Conference, at Annecy 1996.
44. *Techniques in Protein Chemistry*. Proceedings from the annual meetings of the Protein Society, usually published yearly by Academic Press.
45. *ABRF News*. Publication of the Association of Biomolecular Resource Facilities.

## Protein Splicing

Two types of **intervening sequences** within **genes**, **introns** ( [RNA splicing](#) elements) and **inteins** (protein splicing elements) are not present in the mature [protein](#) (1). Introns are spliced out of a precursor RNA before translation of the [messenger RNA](#), whereas inteins are posttranslationally excised from precursor proteins by a process termed protein splicing (Fig. 1). The discovery of protein splicing in 1990 (2, 3) fundamentally changed our view of genetic information flow from DNA to RNA to protein. Protein splicing results in the expression of two proteins from a single gene and thereby is an exception to the rule of one gene, one protein. Protein splicing is defined as the synthesis of a precursor protein followed by the excision of an internal intein domain and the simultaneous ligation of the external protein domains (termed exteins ) to form the excised intein and the mature protein product of the extein gene (1). The ligation of the flanking extein domains to form a mature protein differentiates protein splicing from other types of **posttranslational** processing of precursor proteins or [polyproteins](#) (4). Protein splicing requires proteolytic cleavage of the precursor protein at both splice junctions and formation of a native [peptide bond](#) between the exteins. This complex pathway is mediated by a highly coordinated series of simple chemical reactions (four nucleophilic attacks). The elements that direct these chemical reactions are contained within the intein plus the first C-extein residue, which are capable of splicing if cloned in foreign protein contexts. However, splicing in foreign proteins is generally not as efficient as splicing in the native context, suggesting that the exteins may play a role in the splicing reaction by helping to align the splice junctions or by enabling proper folding of the intein catalytic core.

**Figure 1.** Comparison of protein (a) and RNA (b) splicing. Genes that contain protein splicing elements (inteins) are transcribed normally. After synthesizing a precursor protein, however, the intein is precisely excised, and the flanking exons are ligated together to form two stable protein products. The five conserved residues that actively participate in the splicing are shown below the protein precursor (in one-letter code). They are a Ser (S), Thr (T), or Cys (C) immediately following both splice sites, a His (H) about 100 residues from the intein N-terminus, and the dipeptide His-Asn (HN) at the intein C-terminus (immediately adjacent to the downstream splice junction). RNA splicing differs from protein splicing in that the precursor RNA is spliced before forming the mature mRNA, which is then translated into a single mature protein product.



The intein, often a bifunctional protein, is also a site-specific **endonuclease** capable of cleaving DNA at the intein insertion site in genes that lack the intein (also termed the intein “homing site”) (1, 5, 6). This type of endonuclease is called a homing endonuclease (7). Homing endonucleases were originally described as **open reading frames** in mobile self-splicing introns in which homing endonuclease activity is required for intron mobility via the double-strand, break-repair pathway (8, 9). Many inteins are members of the dodecapeptide class of homing endonucleases and are thus potential genetic **mobile elements** (7-10). Inteins and self-splicing introns also have similar phylogenetic distributions in **Archaeobacteria**, **eubacteria**, and lower **eukaryotes**. Inteins and self-splicing introns have not been found in higher eukaryotes.

In early papers, inteins were called spacers, protein introns, protozymes, or intervening protein sequences (IVPS), among other names (2-4, 6, 11-15)). In 1994, however, the majority of workers in the protein splicing field established a standard nomenclature for protein splicing elements consisting of the terms intein and extein, as defined previously (1). The name of each intein (gene or protein) consists of a three-letter genus and species designation, followed by the name of the extein in which the intein resides. If multiple inteins are present in the same protein, they are given a numerical suffix. For example, the second intein in the *Thermococcus litoralis* DNA polymerase gene is called *Tli* pol intein-2. If endonuclease activity has been demonstrated, the intein is also given an endonuclease designation following the restriction enzyme nomenclature convention with the addition of the prefix “PI-” (eg, PI-*Tli*I).

Protein precursors that contain inteins are synthesized just like any other protein (Fig. 1) (see

**Protein biosynthesis**). The gene is first **transcribed** into mRNA, and then the entire mRNA that contains the intein sequence is **translated** to produce a precursor protein. The RNA is not spliced. Silent substitutions at the extein/intein junction that change the RNA sequence, but not the protein sequence, do not affect protein splicing as they would RNA splicing. Then, the precursor protein is spliced to yield the free intein and the ligated exteins. Splicing of many inteins, including the *Tli* pol inteins, *Psp* pol intein-1, *Sce* VMA intein, and *Mtu* recA intein, does not depend on host cell factors because these inteins splice when expressed in many different heterologous systems, such as *Escherichia coli*, eukaryotic *in vitro* translation systems, or **baculovirus** (4, 6, 12, 16). However, the *Mle* recA intein does not splice when expressed in *E. coli* (13), suggesting that either some inteins may require host factors for splicing or that some expression systems do not provide a suitable environment for properly folding or splicing inteins.

Irrespective of the expression system, splicing is so rapid in the native protein that the precursor is rarely observed in cell extracts after staining of **SDS-PAGE** or **Western blotting**. The inability to demonstrate splicing of a purified precursor protein was a major obstacle to proving that inteins splice as proteins and not as RNA. A major breakthrough in understanding protein splicing was the establishment of the MIP *in vitro* protein splicing system, consisting of the *Psp* pol intein-1 cloned in-frame with the *E. coli* **maltose-binding protein** (MBP) at its N-terminus and a fragment of *Dirofilaria immitis* paramyosin at its C-terminus (16-18). Because MIP contains the maltose-binding protein, it can be rapidly purified by **affinity chromatography** over amylose resin. Splicing of MIP is blocked by inducing the fusion protein at low temperatures (12° to 20°C) and purifying it at pH 8.5 or above (MIP splicing is slowed at pH 8.5 and totally inhibited at pH 10). Then, splicing of the purified precursor protein is induced by raising the reaction temperature to 30° to 65°C and lowering the pH to 6 to 7.5. Splicing of the *Psp* pol intein-1 in its native DNA polymerase context is not temperature-dependent in *E. coli*, even though it was derived from an extreme thermophile (*Pyrococcus spp.*, strain GB-D) that normally grows at 95°C. This suggests that the foreign protein environment in MIP is responsible for the observed temperature-dependent splicing of this fusion protein.

The MIP *in vitro* splicing reaction was followed by examining the size of the proteins present on **Coomassie Brilliant Blue**-stained SDS-PAGE or by probing Western blots with **antisera** specific for each of the three domains. When purified MIP precursor is incubated *in vitro* over time, the full-length fusion protein gradually disappears, and a concomitant increase occurs in a larger, slowly migrating protein, instead of the expected splicing products. This larger protein still contains all three MIP domains. The expected spliced products do not appear until the slowly migrating protein begins to disappear. These data suggested that the slowly migrating protein is a splicing intermediate. Formation of the slowly migrating intermediate is reversible. At pH 10, the slowly migrating species reverts back to the full-length precursor. N-terminal amino acid sequencing indicates that the slowly migrating protein is a branched protein that has two N-termini comprising the MBP N-terminus (the N-extein) and the intein N-terminus. The N-extein (MBP) is linked to the remainder of the precursor (IP) by an alkali-labile bond. The formation of stable branched proteins is very rare in biology.

### 0.1. Identification of Inteins and Conserved Intein Motifs

Because inteins are in-frame insertions that are absent from homologous extein genes, intein junctions were initially determined by comparison with sequenced extein homologues. In many cases, pairwise comparisons with extein homologues or insertion of the intein in a conserved extein motif clearly indicated the intein boundaries. These predictions were confirmed when the N-termini of the free *Sce* VMA, *Tli* pol-2, and *Psp* pol-1 inteins were sequenced by **Edman degradation** (4, 6, 16). Characterization of many inteins indicates that they are moderately sized proteins (300 to 550 residues) with two known exceptions, *Ppu* dnaB intein (150 residues) and *Mxe* gyrA intein (198 residues) (19). However, these small inteins have yet to be analyzed for splicing or endonuclease activity.

Analysis of 36 intein sequences revealed four conserved splice junction residues (see the precursor protein in Fig. 1) and eight conserved intein motifs (Blocks A to H) that are considered important for

splicing or endonuclease activity (19, 20). Known inteins begin with either Ser or Cys, although it should also be possible for Thr to occur at this position. These amino acids contain side-chain hydroxyl (Ser or Thr) or **thiol** (Cys) groups. Inteins end in Asn, and 33 of 36 sequenced inteins have a His preceding the C-terminal Asn. C-exteins begin with Ser, Thr, or Cys. Thus, the amino acid at the downstream side of both splice junctions is limited to Ser, Thr, or Cys. As discussed later, these residues participate in the chemical reactions required in protein splicing. The four conserved amino acids at the two splice junctions reside in larger conserved intein terminal motifs, Blocks A and G. Another conserved intein motif, called Block B, is usually approximately 100 residues from the intein N-terminus and contains a second conserved His, which, it is thought, is involved in N-terminal splice junction cleavage. Two of the intein motifs are also the homing endonuclease dodecapeptide motifs (Blocks C and E) that catalyze DNA cleavage. It has been suggested by many researchers that the region between and including the dodecapeptide motifs functions as the homing endonuclease and that the regions flanking these motifs direct splicing. Further study, especially structural analysis, is required to support this hypothesis.

## 0.2. Inteins, Endonuclease Activity, and Mobility

The dodecapeptide motifs in the dodecapeptide class of homing endonucleases are also found in many inteins. However, these motifs are absent or mutated in many inteins, suggesting that not all inteins are active homing endonucleases. Endonuclease activity is not required for protein splicing (21). Endonuclease activity has been demonstrated in inteins from yeast [PI-*Sce* I, (5)], *T. litoralis* [PI-*Tli* I and II, (6)] and *Pyrococcus spp.*, strain GB-D (PI-*Psp* I). Homing endonuclease activity in self-splicing introns is a prerequisite for intron mobility. *Sce* VMA intein mobility has been demonstrated during **meiosis** in yeast (10). Homing endonucleases have large (15 to 40 nucleotide) recognition sequences that correspond to the intein or intron homing site and are absent in the indigenous genome when the intein or intron is present (7). The homing site is composed of the intein or intron sequences that flank the insertion site in genes that lack the intein or intron. When an intein or intron is present in the gene, the homing site is disrupted, and the endonuclease cannot cleave the DNA. Therefore, homing endonucleases can only cut DNA of alleles that do not contain the intein or intron. Mutational analysis indicates that homing endonucleases will still cleave target DNA containing one or more mutations. In fact, very few specific positions within the recognition sequence are strictly required for DNA cleavage. This lack of complete sequence specificity is advantageous for enzymes that initiate mobility because individual, species, or genus variation would otherwise limit cleavage activity and mobility.

Inteins that encode PI-*Tli* II and PI-*Psp* I are inserted in the same position in their respective DNA polymerase genes. These enzymes are **isoschizomers** (different endonucleases that have the same recognition site). Inteins in the same site in extein homologues are intein **alleles** (19). Intein alleles can be transmitted horizontally as mobile genetic elements and are more closely related to each other than to nonallelic inteins. Not every individual in a species necessarily contains an intein identified in one member of that species (13, 22).

Although the mechanism of intein mobility has not been determined, it is considered analogous to intron mobility mediated by the dodecapeptide class of homing endonucleases (8, 9). Intein mobility is initiated when an extein allele that lacks the intein enters the cell during sexual **reproduction, conjugation, transduction, bacteriophage infection, or plasmid transfer**. The intein endonuclease activity cleaves the extein gene at the intein homing site, making a double-strand break. This double-strand break initiates very efficient **gene conversion** by the double-strand, break-repair pathway, and most extein genes acquire the intein.

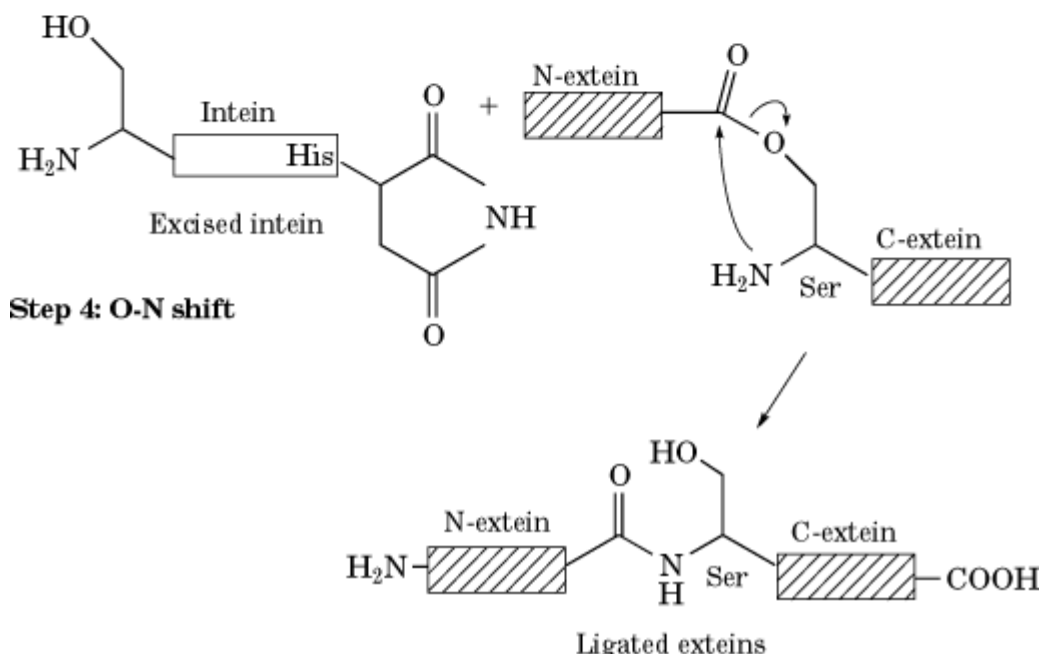
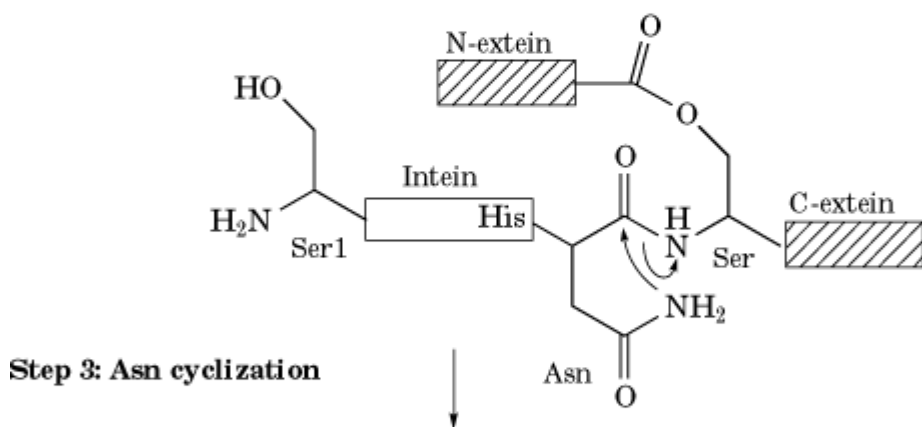
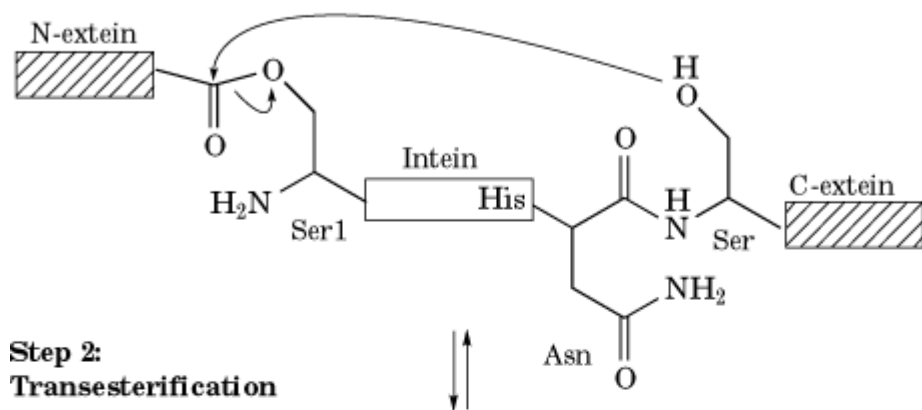
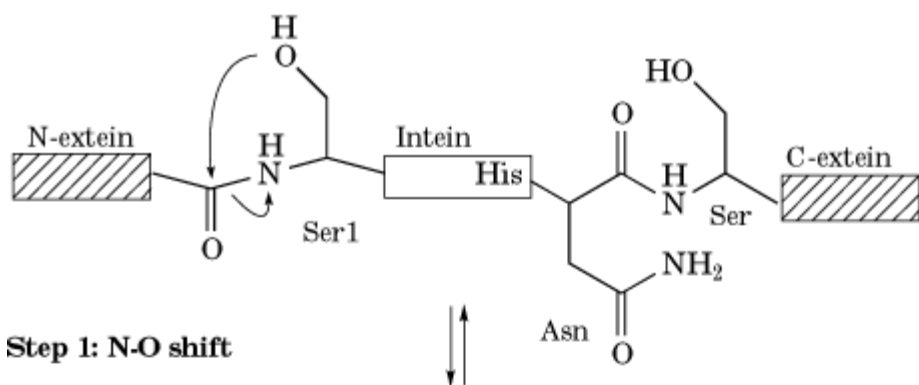
## 0.3. The Mechanism of Protein Splicing

The mechanism of protein splicing in the *Psp* pol intein-1 was determined by using a combination of genetic, chemical, and biophysical techniques (16-18, 23, 24). Analysis of the *Sce* VMA intein confirmed this splicing pathway (25). The protein splicing pathway requires four nucleophilic attacks to complete the splicing reaction. The chemically important amino acids are the intein N-terminal residue, the intein C-terminal residue, and the downstream extein N-terminal residue (see

precursor protein in Fig. 1). The intein catalytic core is formed by bringing the conserved amino acids at the two splice junctions into proximity by properly folding the intein. Although the protein splicing mechanism was initially defined in the *Psp* pol intein-1 system with serine residues following both splice junctions, the mechanism applies equally well to inteins that have either Thr or Cys in these positions. The *Sce* VMA intein has cysteine after each splice junction. Ser, Thr, or Cys have similar side-chain groups (hydroxyl or thiol) and chemical properties. However, because of the differences in their nucleophilic properties and  $pK_a$ , there are slight differences in the reactions involving each of these residues.

The mechanism of protein splicing is presented in Fig. 2. The first step in the protein splicing pathway is a chemical rearrangement involving the intein N-terminal residue. This reaction is called an acyl or N-O shift (or rearrangement). If the residue is a Cys, it is called an acyl or N-S shift. The side-chain hydroxyl (or thiol) attacks the carbonyl (C=O) of the adjacent upstream amino acid. This nucleophilic attack results in cleavage of the peptide bond at the N-terminal splice junction and shifting of the N-extein onto the side chain of the intein N-terminus. Thus, the N-extein is connected to the remainder of the precursor by an ester (or thioester) bond. Acyl shifts involving Ser, Thr, or Cys are reversible, and the equilibrium strongly favors the peptide bond, not the ester form. However, the environment created by the intein catalytic core enhances formation of the ester. The second step in the splicing pathway is a transesterification reaction. The side-chain hydroxyl (or thiol) of the Ser (or Thr/Cys) at the downstream splice junction attacks the carbonyl at the upstream splice junction, resulting in cleavage of the upstream splice site and transfer of the N-extein to the side chain of the Ser (or Thr/Cys) at the downstream splice junction, forming the previously described slowly migrating branched protein intermediate. The third step in the splicing pathway is resolution of the branched intermediate by cyclization of the intein C-terminal Asn residue to form a succinimide ring, which results in cleavage of the downstream splice junction. This occurs when the side-chain amide nitrogen of Asn attacks the Asn main-chain carbonyl (C=O). The succinimide ring at the intein C-terminus is susceptible to hydrolysis to form Asn or isoasparagine (aspartic acid amide) (see also [Deamidation](#)). After removal of the intein, the N-extein is still linked to the side chain of the Ser (or Thr/Cys) at the downstream splice junction. This structure is identical to the ester form of Ser (or Thr/Cys) after an acyl shift. Once the intein is removed, this ester is no longer stabilized and, in Step 4, the Ser (or Thr/Cys) spontaneously undergoes an O-N (or S-N) acyl shift to form a native peptide bond between the exteins.

**Figure 2.** The chemical mechanism of protein splicing. The intein is depicted with a Ser residue following both splice junctions. It should be understood that the mechanism is essentially unchanged when a Thr or Cys is at either or both of these positions. Protein splicing involves four nucleophilic attacks. The first step is an N-O acyl shift at the intein N-terminus, followed by a transesterification reaction that results in forming the branched intermediate. The intein is cleaved from the branched intermediate by cyclization of the intein C-terminal Asn. The native peptide bond is formed between the two exteins after a spontaneous O-N shift. See text for further details. The residues involved in protonation and deprotonation reactions are not shown.





Mutagenetic data indicate that the penultimate His residue of the intein is required for branch resolution (succinimide formation) and thus may be involved in deprotonating the intein C-terminal Asn, directly or via a charge relay system. A few inteins lack this intein penultimate His (19, 26). Other residues in these inteins must assist the cyclization of the intein Asn. However, these residues do not necessarily have to be adjacent to the Asn in primary amino acid sequence, but they need only be correctly positioned in the three-dimensional structure of the enzyme. It has been proposed that the conserved His in Block B assists in the acyl shift at the N-terminal splice junction (20). The identification of all the residues involved in these nucleophilic attacks awaits further mutagenetic experiments and the determination of the crystal structure of a splicing precursor.

The splicing pathway is finely tuned and highly coordinated. Even conservative substitutions of the splice junction Ser, Thr, or Cys residues result in reduced or blocked splicing. When the splicing reaction is perturbed by conservative mutation or expression of the intein in unacceptable foreign protein contexts, the most commonly observed phenomenon is cleavage at one or both splice junctions in the absence of extein ligation. Extein ligation is the most sensitive step in the protein splicing pathway. C-terminal cleavage occurs when Asn cyclization predominates or precedes branch formation. N-terminal cleavage occurs when the upstream ester bond is hydrolyzed or attacked by other nucleophiles before branch formation. Intein N-terminal cleavage and autocleavage of several proteins, including glycosylasparaginase (27) and the hedgehog protein (28), proceed essentially by the same pathway and have an acyl shift resulting in ester (thioester) formation at a Ser, Thr, or Cys, followed by cleavage via nucleophilic attack on the resultant ester (or thioester).

#### 0.4. Protein Splicing and Protein Engineering

Understanding the protein splicing pathway opens up avenues for [protein engineering](#). For example, mutating the C-terminal splice junction Asn and Ser (or Thr/Cys) results in inteins that fail to splice but instead cleave the N-terminal splice junction. Likewise, mutating the intein N-terminus (Ser, Thr, or Cys) results in inteins that cleave only at their C-terminal splice junctions. A controllable N-terminal cleavage system for protein purification has been developed using the *Sce* VMA intein (25, 29). Furthermore, inteins may provide a new method for controlling expression of active proteins. The target protein is inactivated when an intein is present, but it can be activated by splicing the intein. One can imagine numerous methods for controlling splicing. For example, cold-sensitive mutants of the thermostable *Psp* pol intein-1 splice at temperatures over 30°C, but not below 30°C. Temperature-sensitive mutants of inteins from **mesophiles** splice at low temperatures but not higher. Removable blocking agents incorporated into essential junction residues prevent splicing until they are removed (30). Controllable splicing elements provide a means of expressing toxic proteins or making controllable knockouts. Finally, modified inteins provide a means of introducing electrophilic esters and thioesters into proteins, which then can be exploited to tag proteins with nucleophilic reagents (**fluorescent** labels, radioactive Cys residue, peptides) at their C-termini or as substrates for protein semisynthesis (31).

#### Bibliography

1. F.B. Perler, E.O. Davis, G.E. Dean, F.S. Gimble, W.E. Jack, N. Neff, C.J. Noren, J. Thorner, and M. Belfort (1994) *Nucleic Acids Res.* **22**(7), 1125–1127.
2. R. Hirata, Y. Ohsumi, A. Nakano, H. Kawasaki, K. Suzuki, and Y. Anraku (1990) *J. Biol. Chem.* **265**, 6726–6733.
3. P.M. Kane, C.T. Yamashiro, D.F. Wolczyk, N. Neff, M. Goebel, and T.H. Stevens (1990) *Science* **250**, 651–657.
4. A.A. Cooper, Y. Chen, M.A. Lindorfer, and T.H. Stevens (1993) *EMBO J.* **12**(6), 2575–2583.
5. M. Bremer, F.S. Gimble, J. Thorner, and C. Smith (1992) *Nucleic Acids Res.* **20**, 5484.
6. F.B. Perler, D.G. Comb, W.E. Jack, L.S. Moran, B. Qiang, R.B. Kucera, J. Benner, B.E. Slatko,

- D.O. Nwankwo, S.K. Hempstead, C.K.S. Carlow, and H. Jannasch (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5577–5581.
7. J.E. Mueller, M. Bryk, N. Loizos, and M. Belfort (1994) In *Nucleases* (S.M. Linn, R.S. Lloyd, and R.J. Roberts, eds.), Cold Spring Harbor Press, Cold Spring Harbor, New York, pp. 111–143.
  8. A.M. Lambowitz and M. Belfort (1993) *Annu. Rev. Biochem.* **62**, 587–622.
  9. M. Belfort and P.S. Perlman (1995) *J. Biol. Chem.* **270**, 30237–30240.
  10. F.S. Gimble and J. Thorner (1992) *Nature* **357**, 301–306.
  11. E.O. Davis, S.G. Sedgwick, and M.J. Colston (1991) *J. Bacteriol.* **173**(18), 5653–5662.
  12. E.O. Davis, P.J. Jenner, P.C. Brooks, M.J. Colston, and S.G. Sedgwick (1992) *Cell* **71**, 201–210.
  13. E.O. Davis, J.S. Thangaraj, P.C. Brooks, and M.J. Colston (1994) *EMBO J.* **13**(3), 699–703.
  14. H.H. Gu, J. Xu, M. Gallagher, and G.E. Dean (1993) *J. Biol. Chem.* **268**, 7372–7381.
  15. Y. Anraku and R. Hirata (1994) *J. Biochem.* **115**, 175–178.
  16. M. Xu, M.W. Southworth, F.B. Mersha, L.J. Hornstra, and F.B. Perler (1993) *Cell* **75**, 1371–1377.
  17. M. Xu and F.B. Perler (1996) *EMBO J.* **15**, 5146–5153.
  18. M. Xu, D.G. Comb, H. Paulus, C.J. Noren, Y. Shao, and F.B. Perler (1994) *EMBO J.* **13**, 5517–5522.
  19. F.B. Perler, G.J. Olsen, and E. Adam *Nucleic Acids Res.*, in press.
  20. S. Pietrokovski (1994) *Protein Sci.* **3735**(124), 2340–2350.
  21. R.A. Hodges, F.B. Perler, C.J. Noren, and W.E. Jack (1992) *Nucleic Acids Res.* **20**, 6153–6157.
  22. H. Fsihi, V. Vincent, and S.T. Cole (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3410–3415.
  23. Y. Shao, M.Q. Xu, and H. Paulus (1995) *Biochemistry* **34**(34), 10844–10850.
  24. Y. Shao, M.-Q. Xu, and H. Paulus (1996) *Biochemistry* **35**, 3810–3815.
  25. S. Chong, Y. Shao, H. Paulus, J. Benner, F.B. Perler, and M. Xu (1996) *J. Biol. Chem.* **271**, 22159–22168.
  26. C.J. Bult, O. White, G.J. Olsen, L. Zhou, R.D. Fleischmann, G.G. Sutton, J.A. Blake, L.M. FitzGerald, R.A. Clayton, J.D. Gocayne, A.R. Kerlavage, B.A. Dougherty, J. Tomb, M.D. Adams, C.I. Reich, R. Overbeek, E.F. Kirkness, K.G. Weinstock, J.M. Merrick, A. Glodek, J.L. Scott, N.S.M. Geoghagen, J.F. Weidman, J.L. Fuhrmann, D. Nguyen, T.R. Utterback, J.M. Kelley, J.D. Peterson, W. Paul, S.W. Sadow, C. Hanna, M.D. Cotton, K.M. Roberts, M.A. Hurst, B.P. Kaine, K.H. Borodovsky, C.M. Fraser, H.O. Smith, C.R. Woese, and J.C. Venter (1996) *Science* **273**, 1058–1073.
  27. C. Guan, T. Cui, V. Rao, W. Liao, J. Benner, C.L. Lin, and D. Comb, (1996) *J. Biol. Chem.* **271**(3), 1732–1737.
  28. J.A. Porter, S.C. Ekker, W.J. Park, D.P. von Kessler, K.E. Young, C.H. Chen, Y. Ma, A.S. Woods, R.J. Cotter, E.V. Koonin, and P.A. Beachy (1996) *Cell* **86**(1), 21–34.
  29. M.-Q. Xu (1997) *The NEB Transcript* **8**(2), 1–5.
  30. S.N. Cook, W.E. Jack, X. Xiong, L.E. Danley, J.A. Ellman, P.G. Schultz, and C.J. Noren (1995) *Angew. Chem. Int. Ed. Engl.*, **34**, 1629–1630.
  31. P.E. Dawson, T.W. Muir, I. Clark-Lewis, and S.B. Kent (1994) *Science* **266** (5186), 776–779.

### **Suggestions for Further Reading**

32. M. Belfort, M.E. Reaban, T. Coetzee, and J.Z. Dalggaard (1995). Prokaryotic introns and inteins: A panoply of form and function, *J. Bacteriol.* **177**, 3897–3903.
33. J. Clyman (1995). Some microbes have splicing proteins, *ASM News* **61**, 344–347.
34. E.O. Davis and P.J. Jenner (1995). Protein splicing—the lengths some proteins will go to.

Antonie Van Leeuwenhoek **67**, 131–137.

35. D.A. Hickey (1994). Protein introns: Optional or essential? *Trends Genet.* **10**, 147–149.

36. D.A. Shub and H. Goodrich-Blair (1992). Protein introns: A new home for endonucleases, *Cell* **71**, 183–186.

## For proteins from **Protein Stability**

Protein stability involves the ability of the native, folded [protein structure](#) to withstand the disruptive, denaturing influence of the external environment, its resistance to unfolding (see [Proteins](#)). This depends on the type of disruptive influence (eg, extremes of temperature, pressure, pH, or concentration of **denaturants**) and on the criteria for judging the nativeness of the folded conformation. The criteria might be a specific functional property of a given protein (eg, [catalysis](#), **ligand binding**, etc.), a structural property (eg, light **absorption**, [circular dichroism](#), [NMR](#) spectra, [hydrodynamic volume](#)), or a thermodynamic property (eg, **heat capacity**, **enthalpy**). However, when protein denaturation (unfolding) is a highly **cooperative** process (see [Proteins](#)), all of these changes occur simultaneously with variation of the denaturing conditions. In that case, the stability of a protein against any external factor can be specified by the value of the factor that induces 50% of the maximal change in the protein properties. These semiquantitative measures of protein tolerance to various factors have practical importance for specifying the range of external conditions in which that protein preserves its native structure and functions. However, it is not straightforward to compare tolerances against different denaturing factors. They must be correlated. For example, denaturation by temperature depends on the solvent conditions, particularly pH, and denaturation by pH depends on the temperature at which the experiment is carried out. Consequently, these half-denaturation values are not very informative in themselves for understanding the physical and energetic basis of the native protein structure.

With regard to energetics, it is much more useful to specify protein stability by the amount of work required to disrupt its native structure by any means, that is, to specify it as the Gibbs **free energy** difference between the unfolded (denatured) and the native states of the protein. This Gibbs free energy difference is a multidimensional function that depends on all the parameters which specify the external conditions. It can be determined through thermodynamic analysis of the processes of protein denaturation by varying one of the parameters and keeping the others fixed. It can be determined quantitatively, however, only if the unfolding process is reversible, complete, and highly **cooperative**, close to a two-state transition. Therefore, it can be used only for small globular proteins for which denaturation can be approximated by a two-state transition.

For large proteins, denaturation is usually not a two-state transition but proceeds in discrete steps corresponding to the unfolding of the individual structural **domains** present. Therefore, one has to specify thermodynamically the stability of each of the domains individually, plus any interactions between them. The latter might be stabilizing (positive cooperation) or destabilizing (negative cooperation) ([1-3](#)). Estimating of the mutual influence between domains requires determining of the stability of each domain with and without its neighboring domains or reversing of the sequence of domain unfolding by changing the environmental conditions.

Estimation of the Gibbs free energy difference between the native and denatured (unfolded) states is usually based on measuring the equilibrium constant  $K_{eq}$  between these states in the unfolding transition zone where both are populated, using any observable parameter sensitive to the state of the

protein,  $Y(x)$ , that depends on the variable parameter  $x$  that specifies the environmental condition, (eg, temperature, pressure, concentration of denaturant, etc.). For a two-state monomolecular transition

$$\Delta G(x) = -RT \ln K_{\text{eq}} = -RT \ln \left\{ \frac{Y(x) - Y_n}{Y_d - Y(x)} \right\}$$

where  $Y_n$  and  $Y_d$  are values of the parameter that characterizes the pure native and the pure denatured states, respectively (see [Proteins](#)). However, the main interest is not the value of  $\Delta G(x)$  in the denaturation transition zone, but the value under some standard conditions. Extrapolation of this function to the chosen standard condition is the main problem in measuring of protein stability. When protein stability is probed by varying the concentration of a denaturant (eg, [urea](#) or **guanidinium** chloride, GDMCl), the values of the  $\Delta G$  function determined in the presence of rather high concentrations of denaturant have to be extrapolated to zero concentration (4). This is usually done by assuming that this function is linear, although this has no clear basis (see [Proteins](#)). When pressure is used to probe the stability of a protein, the values of  $\Delta G$  obtained at the very high pressure at which the protein denatures have to be extrapolated to normal pressure. Again this is usually done by assuming that this functional dependence is linear (5) because not much is known about protein volume depending on pressure, which determines this function. Only in thermal denaturation and denaturation by pH can extrapolation to standard conditions be done reliably because the functional dependence of  $\Delta G$  values on temperature and pH can be determined experimentally.

In thermal denaturation by a monomolecular, two-state transition, we can determine **calorimetrically** the **enthalpy** and the **entropy** of the transition over the entire temperature range in which the difference between the partial **heat capacities** of protein in the unfolded (denatured) and native states,  $DC_p(T) = DC_p(T)^U - DC_p(T)^N$ , is known (see Eqs. (10) and (11) in [Proteins](#)). Using these functions, we can determine the Gibbs free energy difference between the unfolded (denatured) and folded (native) states of protein at any temperature  $T$  in the temperature range considered (6-8):

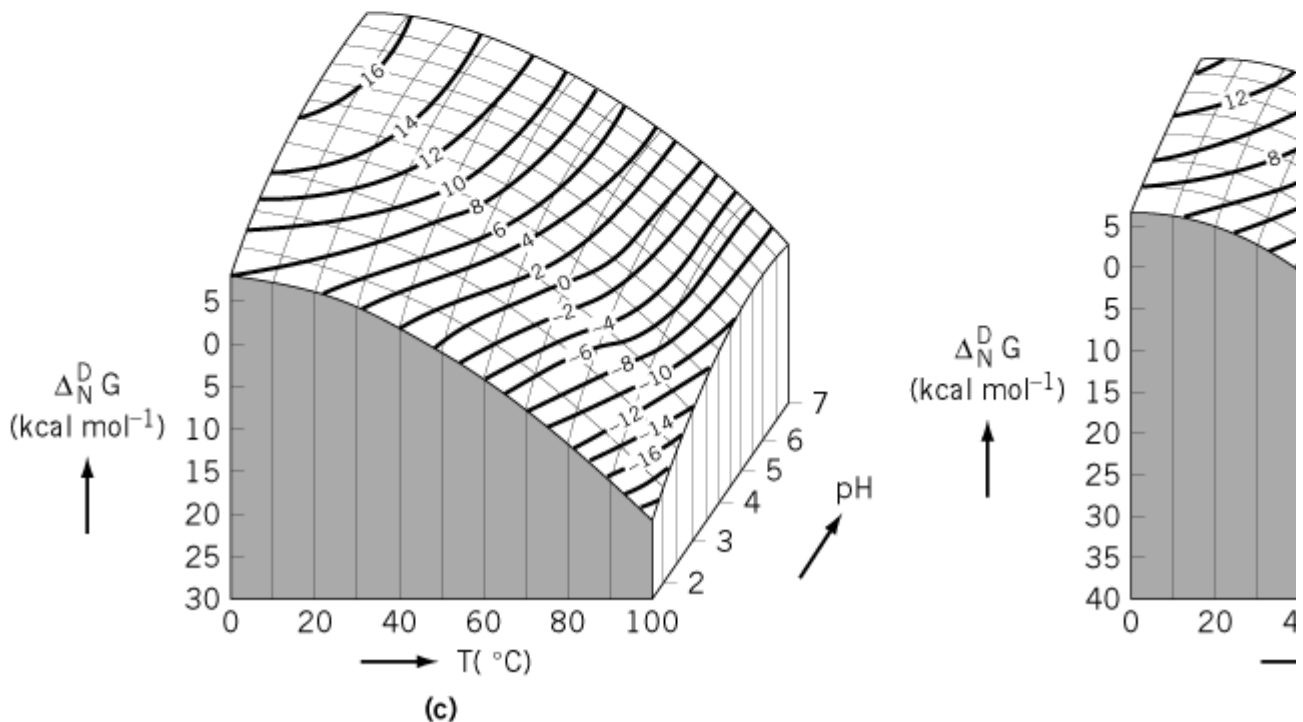
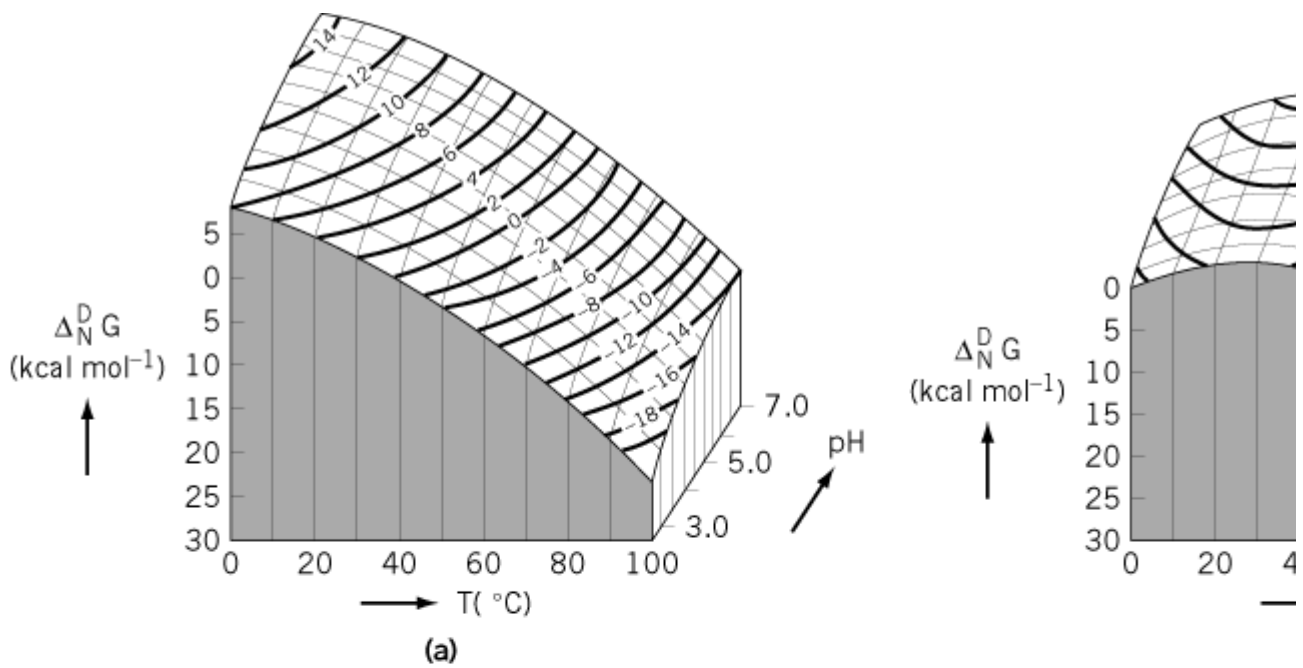
$$\begin{aligned} \Delta G(T)_{\text{pH}} &= \Delta H(T)_{\text{pH}} - T \Delta S(T)_{\text{pH}} = \Delta H(T_t)_{\text{pH}} \frac{T_t - T}{T_t} \\ &+ \int_{T_t}^T \Delta C_p(T) dT - T \int_{T_t}^T \frac{\Delta C_p(T)}{T} dT \end{aligned}$$

In of pH denaturation, potentiometric titration of a protein at some fixed temperature  $T$  can be used to determine the difference between the Gibbs free energies of the unfolded and folded states as a function of pH (9):

$$\Delta G(\text{pH})_T = 2.3RT \int_{\text{pH}_t}^{\text{pH}} \Delta \nu(\text{pH}) d\text{pH}$$

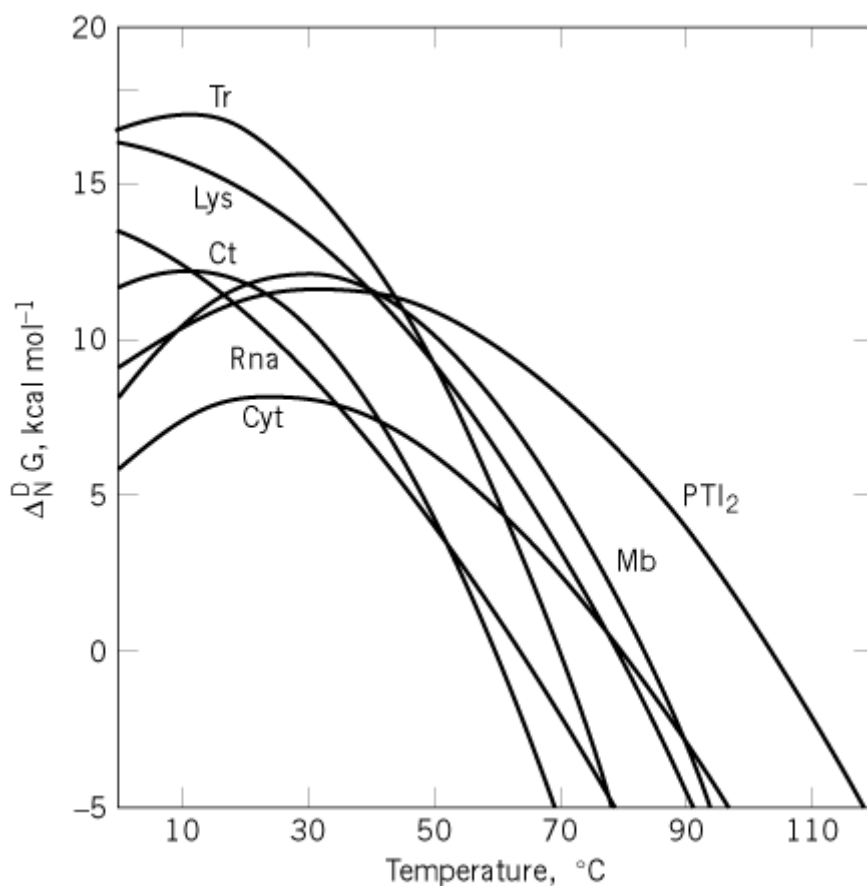
Combining both types of information (10), one can construct  $\Delta G(T, \text{pH})$  functions for the protein (Fig. 1).

**Figure 1.** The Gibbs free energy difference between the unfolded and folded states of (a) **ribonuclease A**, (b) met-**myoglobin**, describing protein stability as a function of temperature and pH. Reproduced from Ref. 10.



The Gibbs free energy difference between the unfolded and the native states of a protein usually decreases in magnitude at acidic and alkaline pH values. It decreases at high temperatures because of the increased dissipative thermal motion, and it also decreases at low temperatures because of [hydration](#) effects (11-14). Because  $(dDG/dT)_{T_{\max}} = -DS(T_{\max}) = 0$  at  $T_{\max}$ , the maximum stability of a protein is achieved at the temperature at which the entropy of unfolding is zero, and therefore the structure is stabilized enthalpically. For most **mesophilic** proteins this maximum takes place at temperatures close to physiological or lower (Fig. 2).

**Figure 2.** Temperature dependence of the molar unfolding Gibbs free energies of various single-domain globular proteins. Tr: bovine **trypsinogen**; Lys, hen **lysozyme**; Ct, bovine **chymotrypsinogen**; Rna, **ribonuclease A**; Cyt, **cytochrome c**; Mb, met-**myoglobin**; PTI<sub>2</sub>, **BPTI**, assumed to be dimeric. Reproduced from Ref. 11.



For proteins from **thermophiles**, which function at much higher temperatures (15), one can expect either a shift of the maximum of the Gibbs free energy of protein stabilization to the higher temperatures or a flattening of the Gibbs free energy function, that is spreading of its positive values over a broader temperature range, including those at which thermophiles function. The first requires increasing the enthalpic interactions that stabilize the native state. This is achieved by improved packing of the **nonpolar** groups in the protein interior or by a decrease of the configurational entropy of protein unfolding, which is achieved by **disulfide** cross-linking of the polypeptide chain. Flattening of the Gibbs free energy function requires decreasing the heat capacity increment of protein unfolding. This can be achieved by increased hydration of polar and charged groups upon protein unfolding, at the expense of nonpolar groups. Comparison of the structures of homologous proteins from mesophiles and thermophiles does not reveal significant differences in their packing densities (16-20). Because of the steep dependence of **van der Waals** interactions on distance, however, even the smallest change in packing density, beyond what can be resolved by **X-ray crystallography** could result in a significant increase of the van der Waals contribution to the stabilization of protein structure. Indirect evidence of a higher packing density of thermophiles is their lower rates of **hydrogen exchange** (17). On the other hand, comparative analysis of structures shows that proteins from thermophiles have more internal **salt bridges** (charged pairs) and **hydrogen bonds** than their mesophilic homologues (21-24). If so, hydration of the charged and **polar** groups upon protein unfolding should decrease the heat capacity increment of unfolding and thus flatten the Gibbs energy function that describes the temperature dependence of protein stability. Unfortunately we do not know much about the real Gibbs free energy function of proteins from thermophiles because these proteins unfold at too high temperatures, where the covalent structure becomes

modified. Consequently, the unfolding of these proteins is usually irreversible.

Determining of the contribution of the various factors to the stabilization of native protein structure is complicated by the fact that many factors contribute with opposite signs and with different dependence on temperature. The details of this balance and the contributions of various factors to the energetics of protein structure are still some of the most disputable subjects in protein science. The lack of experimental information on the magnitude of various contributions to protein energetics has led to a wide variety of opinions. During the past several decades, attention has swung between complete concentration on hydrogen bonding (25) and complete neglect of its role in stabilizing of compact native states. According to a current widely accepted point of view, the main driving force in polypeptide chain folding is the **hydrophobic** interaction, which is responsible for the low solubility of nonpolar liquids in water (26). However, the protein interior can hardly be regarded as a liquid phase, although if one considers it a crystal-like phase, it is more appropriate to describe its energetics by van der Waals interactions between tightly packed groups and by the hydration effects upon exposure of these groups to water (27). Estimates of these effects come from the data obtained by studying the transfer of model compounds from the gaseous phase into water. More recent analysis has led to the conclusion that the main contributors to the stabilizing of compact native protein structures are, first, hydrogen bonding and secondly, van der Waals interactions, and the hydration effects on nonpolar groups. These stabilizing interactions are opposed by the hydration of aromatic and polar groups and by the dissipative force of thermal motion, which is proportional to the increase in configurational entropy upon unfolding,  $T\Delta S^{\text{conf}}$  (14). However, if the energy of hydrogen bonding between polar groups is combined with the hydration effects of these groups and the energy of van der Waals interactions between nonpolar groups is combined with the hydration effect on these groups, the overall contributions stabilizing the protein structure of polar groups (ie, hydrogen bonding) and of nonpolar groups (ie, hydrophobic effect) are almost of the same order of magnitude, and the stabilizing effects of both of them decrease with decreasing temperature (14). Therefore, both of these factors are responsible for the decrease of protein stability at low temperatures (cold denaturation ; see [Proteins](#)). As the charged groups are commonly exposed to the aqueous solvent and are surrounded by counterions, the [electrostatic interactions](#) between them cannot be of major importance for protein stabilization (26). However, internal salt bridges might play a role in protein stabilization because hydration of the individual groups upon unfolding should decrease the heat capacity increment of unfolding and thus flatten the Gibbs energy function. This might be the reason that more salt bridges are found in extremely thermostable proteins from thermophilic microorganisms (21-24). An essential factor stabilizing the native state of protein is also **disulfide** cross-linking, which significantly decreases the configurational entropy gain upon protein unfolding (28).

The maximum stability under conditions close to physiological does not differ much for different mesophilic proteins and is on the order of 30 to 60 kJ/mol for small globular proteins and for the cooperative structural units (**domains**) of multidomain proteins (14, 29, 30). Small proteins consist of about 100 amino acid residues, so it appears that each amino acid residue contributes on average about 0.5 kJ/mol to stabilizing the native state. This value is five times smaller than the energy of thermal motion at room temperature,  $RT = 2.5\text{kJ/mol}$ . This raises two questions: (1) why the native structure is stable at room temperature, and even higher; and (2) why protein stability is so low.

The formal answer to the first question is obvious. A native protein structure is stable only because all of its elements cooperate and it can unfold only cooperatively, that is its stability is determined by the combined contributions of all of the residues. Not much is known, however, about the mechanism of cooperativity. It is supposed only that specific tight packing of groups in the protein interior plays a significant role. As for the low contribution of each residue to the stabilizing the native structure and the low overall stability of this structure, perhaps they are needed because otherwise the polypeptide chain could not search for the proper native conformation during **protein folding** but would collapse into a compact aggregate, the [molten globule](#) (31).

## Bibliography

1. P.L. Privalov and S.A. Potekhin (1986) *Methods Enzymol.* **131**, 4–51.
2. S.A. Potekhin and P.L. Privalov (1982) *J. Mol. Biol.* **159**, 519–535.
3. J.F. Brandts et al. (1989) *Biochemistry*, **28**, 8588–8596.
4. C.N. Pace (1990) *Trends Biochem. Sci.* **15**, 14–17.
5. S.A. Hawley (1971) *Biochemistry* **10**, 3257–3264.
6. P.L. Privalov and N.N. Khechinashvili (1974) *J. Mol. Biol.* **86**, 665–684.
7. P.L. Privalov and G.I. Makhatadze (1990) *J. Mol. Biol.* **213**, 385–391.
8. E. Freire (1995) *Methods Enzymol.* **259**, 144–243.
9. W. Pfeil and P.L. Privalov (1976) *Biophys. Chem.* **4**, 23–32.
10. W. Pfeil and P.L. Privalov (1976) *In Biochemical Thermodynamics* (M. N. Jones ed), Elsevier, Amsterdam, The Netherlands, pp. 75–115.
11. P.L. Privalov (1979) *Adv. Protein Chem.* **33**, 167–241.
12. W. Pfeil (1986) *In Thermodynamic Data for Biochemistry and Biotechnology* (H.-J. Hinz, ed.), Springer Verlag, Berlin, pp. 349–376.
13. G.I. Makhatadze and P.L. Privalov (1994) *Biophys. Chem.* **51**, 291–309.
14. G.I. Makhatadze and P.L. Privalov (1995) *Adv. Protein Chem.* **47**, 307–425.
15. R. Jaenicke (1981) *Ann. Rev. Biophys. Bioeng.* **10**, 1–67.
16. R. Jaenicke (1996) *Adv. Protein Chem.* **48**, 181–267.
17. R. Jaenicke (1996) *FASEB J.* **10**, 84–92.
18. G. Auerbach, U. Jacob, M. Grattinger, H. Schurig, and R. Jaenicke (1997). *J. Biol. Chem.* **378**, 327–329.
19. R.J.M. Russell, D.W. Hough, M.J. Danson, and G.L. Taylor (1994) *Structure* **2**, 1157–1167.
20. D.C. Rees and M.W.W. Adams (1995) *Structure* **3**, 251–254.
21. J.-H. Lim et al. (1997) *J. Mol. Biol.* **270**, 259–274.
22. S. Macedo-Ribeiro et al. (1996) *Structure* **4**, 1291–1301.
23. W. Pfeil et al. (1997) *J. Mol. Biol.* **272**, 591–596.
24. A. Goldman (1995) *Structure* **3**, 1277–1279.
25. M.I. Pauling (1960) *The Nature of the Chemical Bond and Structure of Molecules and Crystals: an introduction to modern structural chemistry*, 3rd ed., Cornell University Press, Ithaca, NY.
26. K.A. Dill (1990) *Biochemistry* **29**, 7133–7155.
27. P.L. Privalov and S.J. Gill (1988) *Adv. Protein Chem.* **39**, 191–234.
28. C.N. Pace and G.R. Grimsley (1988) *Biochemistry* **27**, 3242–3246.
29. C.N. Pace (1990) *Trends Biotech.* **8**, 93–98.
30. A.D. Robinson and K.P. Murphy (1997) *Chem. Rev.* **97**, 1251–1267.
31. P.L. Privalov (1982) *Adv. Protein Chem.* **35**, 1–104.

## Suggestions for Further Reading

32. T.E. Creighton (1991) Stability of folded conformations, *Curr. Opinion Struct. Biol.* **1**, 5–16.
33. K.A. Dill and D. Stitger (1995) Modeling protein stability as heteropolymer collapse, *Adv. Protein Chem.* **46**, 59–103.
34. R. Jaenicke (1996) Structure and stability of hyperthermostable proteins, *Adv. Protein Chem.* **48**, 181–269.
35. W. Kauzman (1959) Some factors in the interpretation of protein denaturation, *Adv. Protein Chem.* **14**, 1–63.
36. G.I. Makhatadze and P.L. Privalov (1995) Energetics of protein structure, *Adv. Protein Chem.*



47, 307–425.

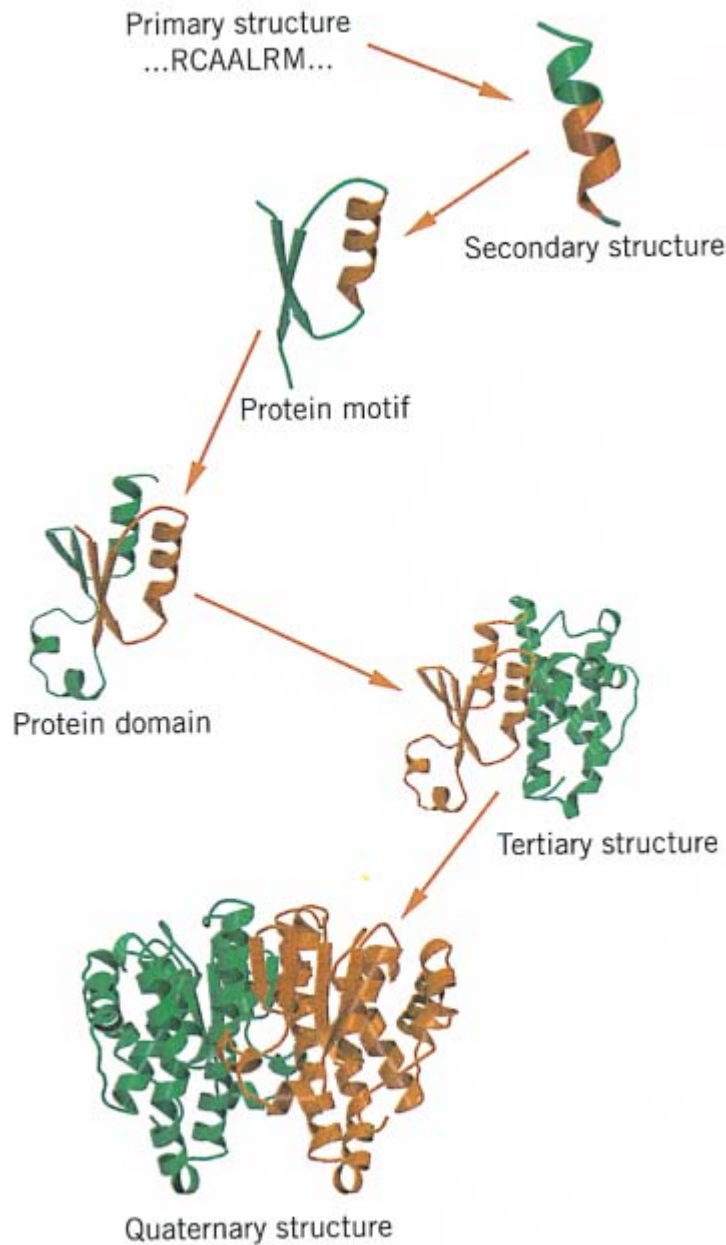
37. C.N. Pace, et al. (1996) Forces contributing to the conformational stability of proteins, *FASEB J.* **10**, 75–83.
38. W. Pfeil (1986) "Unfolding of proteins", In *Thermodynamic Data for Biochemistry and Biology* (H.-J. Hinz, ed). Springer-Verlag, Berlin, pp. 349–376.
39. P.L. Privalov (1989) Thermodynamic problems of protein structure, *Ann. Rev. Biophys. Chem.* **18**, 47–69.
40. P.L. Privalov (1992) "Physical basis of the stability of folded conformations of proteins", In *Protein Folding* (T.E. Creighton, ed.), Freeman, New York.

## Protein Structure

The structures of [proteins](#) are the key to understanding both their function and [evolution](#). Determination and analysis of the stable folded three-dimensional structures of proteins also provides the foundation for [protein structure prediction](#) from amino acid sequences, design of novel proteins with desirable catalytic functions, and design of new medicines with specific and selective therapeutic activities.

Proteins are [biopolymers](#) formed from [amino acids](#) linked together by [peptide bonds](#) into a [polypeptide chain](#). They can consist of many thousands of atoms, and their structures are extremely complex. Protein structure has therefore been classified into different levels, beginning at the simplest level of [primary structure](#) on through **secondary structure**, **protein motifs**, **domains**, [tertiary structure](#), and [quaternary structure](#) (Fig. 1).

**Figure 1.** Different levels of protein structure, from primary to quaternary, are shown using the structure of human p-class [glutathione S-transferase](#) as an example (6). In each succeeding level of protein structure, that portion corresponding to the preceding level is shown in orange. Helices are depicted as coils, and b-strands are depicted as arrows. The secondary structure example shown is an a-helix, and the protein motif is a bab motif (where b is a b-strand, a is an a-helix, and the motif produces a parallel b-sheet). The highlighted protein domain is a [thioredoxin](#) domain, and the green domain is an all-helical domain. The active enzyme is a homodimer, as shown in the quaternary structure. This figure was generated using Molscrip (7) and Raster3D (8, 9).



The primary structure is the sequence, or order, of amino acid residues in the polypeptide chain. By convention, the sequence is numbered from the [N-terminus](#) to the [C-terminus](#). The primary structure also includes information about the number of polypeptide chains in a protein and covalent modifications such as [disulfide bond](#) formation, **phosphorylation**, [sulfation](#), [N-glycosylation](#), and [O-glycosylation](#).

The secondary structure of a protein refers to the **backbone conformation** of the polypeptide chain. There are several different types of secondary structure elements found in proteins including the [a-helix](#), [b-sheet](#), and loops and [turns](#) that reverse the direction of the polypeptide chain. In protein structures, linked secondary structure elements have preferred modes of packing together. These preferred packing arrangements of secondary structure units are called [protein motifs](#) or supersecondary structure. Examples of different types of protein motifs include the [Greek key motif](#) and the [jelly roll motif](#).

The next level of protein structure is the protein domain, although there is some overlap between the structural levels of protein motif and protein domain. The protein domain (also referred to as a *protein fold* or *module*) usually consists of 50 to 150 amino acid residues and is considered to be the building block of protein structure, function, and evolution. In general, it is defined as that part of a protein's structure that can fold independently of the remainder of the protein. Individual domains within a protein structure usually have distinct functions, including catalytic, regulatory, binding, recognition or oligomerization roles. Examples of protein domains include the [EGF motif](#) and the [Kringle domain](#).

The polypeptide chain of a protein folds into one or more domains; the tertiary structure describes the three-dimensional arrangement of these protein domains. Proteins can also incorporate one or more polypeptide chains, called *subunits*. The three-dimensional arrangement of these individual subunits, which may or may not have equivalent sequences and tertiary structures, is described by the quaternary structure of the protein.

Certain rules govern the structures adopted by proteins. For example, the side chains of most globular proteins are distributed in a nonrandom manner. Most **hydrophobic** residues are located in the inner core of the structure, and charged side chains are generally found on the surface. In addition, over 90% of amino acid residues are involved in secondary structure elements such as  $\alpha$ -helices,  $\beta$ -strands, and turns. These rules can be used to assess the quality of experimentally determined structures and to assist and validate protein structure prediction. Three-dimensional structures of proteins can be determined experimentally by protein [X-ray crystallography](#), **nuclear magnetic resonance (NMR)**, or [cryoelectron microscopy](#). Coordinates that describe the atomic structures of all published protein structures are deposited with the Brookhaven Protein Data Bank (PDB) ([1](#)) (see [Structure Databases](#)). The PDB is an extremely useful resource for searching, analyzing, and retrieving protein structure information. Databases that define different structural types of protein domains have also been developed ([2-5](#)).

### Bibliography

1. F. C. Bernstein et al., *J. Mol. Biol.* **112**, 535–542 and <http://www.pdb.bnl.gov/>
2. A. G. Murzin, S. E. Brenner, T. Hubbard, and C. Chothia (1995) *J. Mol. Biol.* **247**, 536–40 and <http://scop.mrc-lmb.cam.ac.uk/scop>
3. C. A. Orengo et al. (1997) *Structure* **5**, 1093–1108 and <http://www.biochem.ucl.ac.uk/bsm/cath>
4. A. S. Siddiqui and G. J. Barton (1995) *Protein Sci.* **4**, 872–884 and (site currently unavailable)
5. R. Sowdhamini, S. D. Rufino, and T. L. Blundell (1996) *Folding and Design* **1**, 209–220 and <http://www-cryst.bioc.cam.ac.uk/~ddbbase/>
6. A. J. Oakley et al. (1997) *Biochemistry* **36**, 576.
7. P. J. Kraulis (1991) *J. Appl. Cryst.* **24**, 946–950.
8. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
9. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

### Suggestions for Further Reading

10. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
11. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.
12. M. Perutz (1992) *Protein Structure: New Approaches to Disease and Therapy*, W. H. Freeman, New York.

## Protein Structure Prediction

The goal of protein structure prediction is to infer the unique three-dimensional [protein structure](#) of any **protein** simply from its [amino acid](#) sequence, or [primary structure](#). There are sufficient data available of known protein structures and sequences to test any prediction method, but the problem has not been solved in spite of many approaches attempted over the past thirty years.

The difficulty arises from two sources. The first is the extremely large number of conformations that a polypeptide chain may adopt, often called Levinthal's paradox ([1](#)). The ideal way to accomplish structural prediction would be to simulate the entire **protein folding** process on the computer. This is possible, in principle, because no unknown factors are involved in the *in vitro* folding process. Calculations to minimize the **free energy** of the conformation should be applicable, taking into consideration the various kinds of interactions within a protein molecule and between the protein and the surrounding medium, [water](#). In fact, this type of computer simulation is effective for small-sized [polypeptide chains](#), such as a [peptide](#) of 10 amino acid residues, but not for a protein of 100 residues. The difference between the calculations for a 10-residue peptide and a 100-residue protein is not a factor of 10 in magnitude, because it increases exponentially. The magnitudes of the two situations may be  $2^{10}$  (about  $10^3$ ) versus  $2^{100}$  (about  $10^{30}$ ), based on the very conservative assumption that every residue, on average, has two different choices of conformation. This vast increase in computational time makes the prediction impractical.

The second difficulty arises from the intrinsic nature of protein folding. Folding is a **cooperative** transition between the **denatured** and the folded states. Folding and unfolding are so cooperative that individual protein molecules remain in one of the two states at any instant without adopting intermediate conformations to any substantial extent. In forming the cooperatively folded structure, each part of the protein molecule is influenced by the entire structure and vice versa. Because of this complexity, there is no simple logical one-to-one correspondence between the sequence and the structure. A segment of polypeptide with a certain sequence may adopt an **alpha-helical** conformation in one protein, but a [beta-strand](#) in another. Therefore, the sequence/structure relationship is only probabilistic, rather than deterministic, making prediction difficult.

Several lines of approach to protein structural prediction have been developed thus far. The folding simulation approach has encountered the problem of the size effect described previously. This approach is important, however, because its success would also explain the mechanism of protein folding, which is one of the greatest enigmas of basic molecular biology. Another type of approach uses rules relating sequence and structure that are derived empirically from the known protein structures. After use for several decades, this approach is successful with **secondary structure prediction**, but it has not been able to address [tertiary structures](#). Another, more recent approach is known as [threading protein sequences](#). This approach is quite different in logic from the others. It assumes that the folded conformation is known and present in the structure [databases](#) and simply examines whether the new sequence is likely to fit any of these known structures. This threading method is an extension of the methods that search for [homology](#) between protein sequences (see [Sequence Analysis](#)) and therefore is more practical than the other methods (see [Homological Modeling](#)). In fact, several proteins of unknown structure have been successfully predicted by the threading method, which is an epoch-making event in the field of protein structure prediction. On the other hand, there is an inherent limitation to the method. It is impossible to predict any truly novel structures because it assumes that the structure is already known. Also, it would not be able to reveal the mechanism of protein folding because the logic behind the algorithm skips the entire process of folding and assumes the final structure in advance.

## 1. Ab initio Predictions

All of the approaches of structural prediction that do not require a known protein structure are lumped together as ***ab initio* methods** (from the Latin “from the beginning”, or from the sequence). Considering the limitations of the threading procedure, it is still worth pursuing the development of *ab initio* methods, which would be applicable to any protein, including those of novel structure. Although this is a very difficult task, success of *ab initio* methods would also provide understanding of the mechanism of protein folding because any successful algorithm would have to incorporate the key factors of protein folding in it.

One of the earliest *ab initio* approaches to simulate the folding of a protein was that of Levitt (2), using [BPTI](#). To reduce the computation time required even for this very small protein (58 residues), the geometry of the polypeptide chain was simplified, as were the atomic force fields employed to estimate the energy. Energy minimization drove the conformation of BPTI from an initial open structure to a folded, globular structure with the lowest energy. The final structure had some resemblance to the known native structure of BPTI. There was a root mean square (rms) deviation of 6 Å between the computed and known structures. Another line of approach is to attempt to pack the secondary structure elements,  $\alpha$ -helices and  $\beta$ -sheets, into a globule using empirical rules (3) or posing general constraints (4). In these cases, it is assumed that the locations of all secondary structural elements along the chain are known in advance. If the algorithm works well, the combination with secondary structural prediction could provide a complete prediction scheme starting simply from the amino acid sequence.

More recently, a series of folding simulations on a lattice was intensively studied by Skolnick et al. (5). Lattice models require that each residue occupy only a point intersection on the lattice and adjacent residues in the sequence must be on adjacent lattice points. They have the advantage of dramatically reducing the number of possible conformations, actually digitizing them into countable objects. Lattice models mimic quite well the overall shape of a protein, but they are poor at mimicking the local structure properly. It is not difficult to imagine the difficulty of representing an  $\alpha$ -helix on a cubic lattice. Using lattices more sophisticated and better suited for realistic protein models, however, the interactions between the various residues occupying neighboring lattice points can be defined to approximate the energy of the conformation. A method that simulates the entire folding process of large proteins, including one with a [TIM barrel](#) fold, has been developed (5, 6). This achievement was not claimed as a prediction because for some interactions specific to the native protein being studied were built in to the computations to guide the folding process in the right direction. If those few guiding potentials can be replaced by general potentials, such results could be called predictions.

After the advent of the threading approach, it should be possible to incorporate the fundamental idea of threading into an *ab initio* method. The difference between the two is not so great. If the backbone conformation were no longer fixed but were allowed to vary, the procedure would be converted from threading to *ab initio*. The trial by Jones (7) is one such example, in which partial, [supersecondary structures](#) are first threaded by a query sequence and then the threaded structures are combined into one as the starting conformation for refinement to a final, novel structure. A combination of the threading algorithm and the *ab initio* energy calculation seems to be a promising direction to be pursued because the efficiency of [threading protein sequences](#) has been verified. When *ab initio* predictions open the door to success, a new era of protein science will begin.

### Bibliography

1. C. Levinthal (1968) *J. Chim. Phys.* **65**, 44–45.
2. M. Levitt (1976) *J. Mol. Biol.* **104**, 59–107.
3. O. B. Ptitsyn and A. A. Rashin (1975) *Biophys. Chem.* **3**, 1–20.
4. F. E. Cohen, T. J. Richmond, and F. M. Richards (1979) *J. Mol. Biol.* **132**, 275–288.

5. J. Skolnick and A. Kolinski (1990) *Science* **250**, 1121–1125.
6. A. Godzik, J. Skolnick, and A. Kolinski (1992) *Proc. Natl. Acad. Sci. USA* **89**, 2629–2633.
7. D. T. Jones (1997) *Proteins: Struct. Function Genet. (Suppl. 1)*, 185–191.

## Protein Targeting, Intracellular

A central feature of the cell is the spatial separation of various biochemical reactions, due to the existence of separate areas containing specific sets of proteins. However, the protein composition of those areas may change due to import or export of proteins in response to extracellular or intracellular signals (eg, the binding of a soluble [kinase](#) to a membrane-bound receptor during [signal transduction](#)); some proteins may be directed to **proteolysis** [eg, targeting of cytoplasmic proteins for in lysosomes degradation ([1](#))]; and new proteins must be translocated from their site of synthesis into the correct area in order to maintain its protein composition. Therefore, cells need a controlled exchange of proteins between the different areas. This entry will concentrate on one aspect, the intracellular protein transport between cellular compartments, emphasizing the fate of newly synthesized proteins that are exported from the cytoplasm.

Cellular compartmentalization is achieved by the separation of [hydrophilic](#) areas with the help of **hydrophobic** membranes. Typical compartments are (a) the cytoplasm, (b) the extracellular space, (c) organelles, like the nucleus or the peroxisome, or (d) parts of organelles, like the different stacks of the [Golgi apparatus](#). The most simple form of compartmentalization is found in many prokaryotes, where the only [membrane](#) present is the plasma membrane, which constitutes the borderline between the interior and the exterior of the cell. In eukaryotes, this compartmentalization is more complex, due to the existence of the various membrane-surrounded cellular organelles. Moreover, the surface of several eukaryotic cell types may also be subdivided. In epithelial cells, for example, the plasma membrane and the adjoining extracellular space of the apical and basolateral surfaces differ considerably in their protein composition.

Because **protein biosynthesis** is restricted to only very few compartments—the cytoplasm of prokaryotes and eukaryotes, the matrix of [mitochondria](#), and the stroma of plastids (see [Chloroplast](#))—all proteins that perform their function in other compartments must be transported into them. In addition, proteins may get relocalized from one compartment to another one in order to regulate their function—for example, the shuttling of [transcription factors](#) between nucleoplasm and cytoplasm, or the internalization of receptor proteins from the plasma membrane into [endosomes](#). Proteins synthesized in the cytoplasm are targeted to all other compartments of the cell or leave the cell completely ([protein secretion](#)). Often the proteins do not enter their final compartment directly from the cytoplasm but have to cross other compartments to reach their destination (see text below). Therefore, the transport of proteins synthesized in the cytoplasm may be divided into the following major steps: the [targeting](#) of the protein from the cytoplasm to the correct translocation sites of an adjoining compartment; the direct [transfer](#) of the protein from the cytoplasm into this compartment; and, if desired, the [sorting](#) to the final compartment.

The signals determining the intracellular transport pathway of proteins are encoded in their [primary structures](#). Proteins that reach their final compartment *in transit* have usually two types of transport signals: One encodes for the translocation site that is used to leave the cytoplasm, and a second one triggers further sorting to the final compartment. Moreover, [membrane proteins](#) need additional signals triggering their integration into the lipid bilayer ([2](#)) and determining their orientation within the membrane ([3](#), [4](#)). The same region of one protein may bear two signals. For example, several

membrane proteins have signal-anchor sequences that function as a signal sequence and as a membrane anchor (5).

Signals may be short defined clusters of amino acids, like the peroxisomal targeting signal 1 (PTS1) (6). Other signals may cover longer regions of the protein, containing more variable sequences that are characterized by more general features—for example, the [signal peptide](#) involved in transport to the [endoplasmic reticulum](#) (ER), the leader peptide of prokaryotic secretory proteins, the mitochondrial targeting sequences, or the plastid targeting sequences (7, 8). But a modification of a protein may also serve as a signal. For example, lysosomal proteins bear special sugar modifications as sorting signals (9), and some plasma membrane proteins are directed to the apical plasma membrane due to the presence of a post-translational lipid modification, the [GPI anchor](#) (10). In the end, even here the primary structure of the protein is decisive, because it determines if this modification may occur (11, 12).

## 1. Targeting

The synthesis of proteins exported from the cytoplasm starts on free [ribosomes](#), regardless of whether the transport occurs post-translationally or co-translationally. Therefore, the first step during the export of a newly synthesized protein from the cytoplasm is directing it to the correct translocation site. Although this is always triggered by specific protein-coded signals, there may be signals within the [messenger RNA](#) of exported proteins, leading to concentration of the mRNA near the appropriate translocation sites and facilitating the targeting process (13).

With the exception of cyanobacteria, prokaryotic cells have only the plasma membrane in contact with the cytoplasm. Nevertheless, these cells may have several protein export pathways that use different types of translocation sites and targeting signals (14-18) (see Table 1). Although eukaryotic cells contain many compartments adjoining the cytoplasm, the bulk of protein export occurs across membranes of only a few organelles—the ER membrane, the outer membrane of mitochondria (or plastids), the [nuclear envelope](#), and perhaps the peroxisomal membrane (see Table 1). Secretory proteins, plasma membrane proteins, and proteins destined to compartments of organelles forming the so-called secretory pathway (Golgi apparatus, endosomes, transport and storage vesicles or lysosomes/vacuoles) are usually first translocated into the ER. However, a few secretory proteins and some proteins of the yeast vacuole are transported directly across the plasma membrane (19) or the vacuolar membrane (20), respectively.

**Table 1. Pathways of the Intracellular Protein Transport of Newly Synthesized Proteins<sup>a</sup>**

| Synthesizing Compartment | Exit Membrane   | Type of Translocation Site   | Targets  |
|--------------------------|-----------------|--|--|
| Prokaryotic cytoplasm    | Plasma membrane | SecY-complex   | Inner membrane space, plasma membrane, outermembrane extracellular space |
|                          |                 | ABC-type transporter (Type I secretion) (eg, secretion of hemolysin) | Inner membrane space, extracellular space                                |

|                      |                                     |   |   |
|----------------------|-------------------------------------|---|---|
|                      |                                     | Contact-dependent secretion system (Type III secretion), (eg, transport of YpkA Ser/Thr kinase) | Extracellular space (plasma membrane and cytoplasm of eukaryotic host)  |
|                      |                                     | Unknown (Sec-independent secretion of TMAO reductase)   | Inner membrane, inner membrane space, extracellular space   |
| Eukaryotic cytoplasm | Thylakoid membrane of cyanobacteria | SecY complex  | Lumen and membrane  |
|                      | Nuclear membrane                    | Nuclear pore  | Compartments of the nucleus   |
|                      | ER membrane                         | Sec61 complex   | Extracellular space, plasma membrane, compartments of the secretory pathway (Golgi apparatus, lysosomes/vacuoles, endosomes, transport vesicle, storage vesicles) |
|                      | Outer mitochondrial membrane        | TOM complex and TIM complex   | Outer membrane, inner membrane, inner membranespace and matrix of mitochondria  |
|                      | Outer plastid membrane              | TOC complex and TIC complex   | Outer envelope, inner envelope, inner membranespace, stroma, and thylakoids of plastids   |
|                      | Peroxisomal membrane?               | Unknown   | Peroxisome  |
|                      | Plasma membrane                     | ABC-type transporter (eg, secretion of a-factor)  | Extracellular space   |
|                      |                                     | Unknown (eg, secretion of basic fibroblast growthfactor)  | Extracellular space   |



|                      |                    |   |                        |
|----------------------|--------------------|---|------------------------|
|                      | Unknown            | Unknown (eg, aminopeptidase I transport)                    | Vacuolar lumen (yeast) |
| Mitochondrial matrix | Inner membrane     | Unknown   | Inner membrane         |
| Stroma of plastids   | Inner envelope     | Unknown   | Inner envelope         |
|                      | Thylakoid membrane | SecY  | Thylakoids             |
|                      |                    | Unknown (eg, pH-dependent import of 23-kDa protein of PSII) | Thylakoids             |

<sup>a</sup> ABC, ATP-binding cassette; TOM, translocase-of-the-outer-membrane; TIM, translocase-of-the-inner-membrane; TOC, translocase of the outer chloroplast envelope; TIC, translocase of the inner chloroplast envelope; PSII, photosystem II. For some pathways, typical substrates are included. The table does not cover pathways exclusively involved in the insertion of membrane anchors.

In some pathways, correct targeting requires recognition of the protein-coded signal by a specific cytoplasmic receptor. One example is the co-translational transport of proteins across the ER membrane. In this case, targeting starts as soon as the signal sequence that is responsible for the targeting of the nascent chain has left the ribosome. Briefly, the signal sequence is recognized by its cytoplasmic receptor, the [signal recognition particle](#) or SRP, which also binds to the ribosome ([21](#), [22](#)). The SRP-ribosome-nascent chain complex is then transported to the translocation sites in the ER. Here the interaction between SRP and its membrane-bound receptor triggers release of the signal sequence ([23](#)), which is then recognized by its membrane receptor, a step which is essential for the following transport across the membrane ([24](#)). Another example is the post-translational import of nuclear proteins. Specific cytoplasmic factors have been discovered that bind the targeting signals of distinct proteins and carry them through the nuclear pores, finally releasing those proteins into the nucleoplasm. Examples are the importin complex, which recognizes the “classical” basic nuclear localization sequence ([25](#)), and transportin, which interacts with nuclear import signals of the M9-type ([26](#)). Some mitochondrial proteins are also imported using a specific targeting factor: the mitochondrial import stimulating factor (MSF), which recognizes the respective targeting signal ([27](#)) and triggers the binding of the substrate to the Tom70/Tom37 receptors at the translocation sites of the outer membrane ([28](#)). For most of the targeting pathways, however, specific cytoplasmic factors are still missing. It may be that in these cases the signals are recognized only directly at the membranes by their corresponding receptors at the translocation sites [eg, Sec61 at the ER ([24](#)), the Tom20/Tom22 ([29](#)) at the mitochondria, or the various peroxisomal targeting sequence receptors at the peroxisomes ([30](#))]. This recognition might be supported by general [molecular chaperones](#), like cytosolic Hsp70 ([31-33](#)), keeping the proteins in a loosely folded state and helping to present the targeting signal to the environment (see [BiP \(Hsp70\)](#)).

## 2. Transport at the Translocation Site

In the co-translational mode of translocation, the nascent chain is crossing the membrane in an extended conformation due to the tight binding of the ribosome to the ER translocation site ([34](#)). But in many post-translational pathways also, translocation of proteins across the membrane requires the substrate to adopt a more or less unfolded conformation. This is observed during mitochondrial

import (35), transport across the ER (36), SecA-dependent transport across bacterial membranes (37), and the import into plastids (38) (39). It should be noted that some proteins may be transported across the plastid envelope in a nonextended conformation (40). As far as is presently known, these translocation sites consist of membrane–protein complexes (translocases) that often form channel-like structures (39) (41–44). These translocases allow the transport of hydrophilic proteins across the hydrophobic membrane bilayer and may also trigger the integration of hydrophobic membrane anchors of membrane-spanning proteins into the bilayer (45, 46). For this purpose, their structures must be very dynamic. For example, in the case of the ER translocation site, it was found that the channel opens both perpendicularly to the membrane, to allow the transfer between cytoplasm and lumen (47), and laterally, to allow the transfer of membrane anchors from the channel into the bilayer (48). Moreover, although these channels allow the signal-dependent transport of relatively large hydrophilic molecules, they are tight enough to prevent uncontrolled ion fluxes between the compartments. In mitochondria and plastids, the translocation sites may even traverse both the outer and inner membranes of the organelle, allowing a direct delivery of cytoplasmic proteins to the inner membrane and the matrix or stroma of the organelle. In mitochondria, these two membrane-spanning structures are formed due to the dynamic interaction of the translocase-of-the-outer-membrane (TOM) complex with the translocase-of-the-inner-membrane (TIM) complex (49, 50) (see Table 1). Nevertheless, the TOM complex may also function independently from the TIM complex during protein transport to the outer membrane or the inner membrane space (51). A similar situation may exist within the two membranes of the chloroplast envelope (52). During or shortly after translocation, the proteins are folded. In the ER, in mitochondria, in plastids, and in the periplasmic space of certain bacteria, various chaperones exist that assist the folding of proteins that traversed the membrane in an extended conformation and. Some of these chaperons also are involved in the directionality of the post-translational transport across the membrane. Examples are Kar2p, which is involved in protein import into the ER lumen (36, 53) (see [Disulfide Bonds](#)), and mitochondrial Hsp70, which is essential for the import into the mitochondrial matrix (54). Folding is accompanied by **post-translational** modifications, like the removal of amino-terminal targeting signals (55, 56), the formation of [disulfide bonds](#), or [N-glycosylation](#). In the case of the transport into the ER, the modifying enzymes may be part of the translocation site (OST, signal peptidase).

Correct functioning of the folding and modification machinery is in many cases a prerequisite for further transport of proteins. In mitochondria and plastids, the chaperones present in the matrix and the stroma, respectively, have to guarantee that proteins that must cross further membranes (see below) are kept in a translocation-competent form (57). Incorrect folding/modification of proteins leaving the ER may lead to their retention in the ER, often followed by their **protein degradation** (58). Correct modification in the ER may also be essential to create the signals for further sorting (eg, GPI anchor, or mannose 6-phosphate modification).

In various pathways, proteins leave the cytoplasm in a more or less completely folded state. The translocation site used for the transport of proteins into the nucleus is also a hydrophilic proteinaceous channel, the **nuclear pore**. However, it is far larger than the channels mentioned above, and it allows the unregulated exchange of molecules of up to 60 kDa between nucleoplasm and cytoplasm (59). For other export pathways that probably also transport folded proteins (the transport into peroxisomes, the direct transport of proteins into yeast vacuoles, and the **Sec**-independent transport of bacterial proteins across the plasma membrane) the translocation site is unknown. Delivery of proteins to the yeast vacuole may not rely on the transport across a hydrophobic bilayer through channel-like structures, but on vesicle-mediated processes similar to those found during autophagy (61). In the case of import into peroxisomes, very complex pathways are presently being discussed, including ER-dependent steps and vesicle-mediated processes (62).

### 3. Sorting to the Final Compartment

Sorting of proteins from the ER to their final compartment employs a system of vesicular budding and fusion, without any further membrane-crossing steps (63). This vesicular flow mediates a specific and controlled transport of substances between the different organelles, including the plasma

membrane. Proteins that end up in the compartments of other organelles of the secretory pathway, as well as proteins that stay in the ER, need a sorting signal. Secretory proteins and plasma membrane proteins may reach their destination by default (64), but many of them rely on specific sorting signals too (eg, proteins of the apical plasma membrane or secretory proteins that are subject of regulated secretion) (65). The signals may code for a true retention in the target compartment (eg, signals involved in the so-called kin-recognition of the membrane proteins of the Golgi apparatus), or they may trigger the packaging of proteins into specific transport vesicles. Sorting may result in forward transport along the secretory pathway, or it may lead to a retrieval of proteins that left their “final” compartment accidentally. For example, a carboxy-terminal located KDEL sequence (-Lys-Asp-Glu-Leu) guides soluble proteins back from the *cis*-Golgi to the ER (66). There is still not much known about how the packaging into the transport vesicles occurs. Well-studied examples are membrane proteins containing signals in the cytoplasmic portion of the protein—for example, the KKXX-motif (Lys-Lys-X-X) of ER-proteins (67) or the di-leucine-containing motif of endosomal proteins (68). These signals may be recognized by components of the cytoplasmic coat of transport vesicles (69, 70). Since the assembly of the coat is essential for the budding of these vesicles, this interaction triggers the inclusion of the membrane proteins into these vesicles (see [Membrane Anchors](#)). Soluble proteins or membrane proteins that have the sorting signal located in the lumen of the organelle may be recognized by membrane receptors [eg, the KDEL receptor (71) or the **mannose 6-phosphate receptors** (72)]. These receptors in turn contain signals that trigger their sorting, together with the bound ligand, into the correct transport vesicle (73). In some cases, the membrane anchors of proteins define the final compartment [eg, some membrane proteins of the Golgi apparatus (74, 75) (see [GPI Anchor](#))], probably due to their tendency to partition preferentially into a specific lipid environment. Similar mechanisms are discussed for the sorting of membrane proteins with [GPI anchors](#).

Some mitochondrial and plastid proteins that left the translocation sites and reached the matrix or stroma, respectively, are also subject to sorting events, such as proteins destined for the chloroplast thylakoids, the inner envelope membrane of chloroplasts (76), or the mitochondrial inner membrane (57) (conservative sorting). In contrast to what is known about the secretory pathway, this additional sorting includes a second step of crossing a membrane. Both mitochondria and chloroplasts are descendants of formerly free-living bacteria according to the endosymbiont hypothesis. Therefore, matrix and stroma are homologous to the bacterial cytoplasm, and the inner membranes and thylakoid membrane are homologous to the plasma membrane, respectively. Signals triggering the export of proteins from the mitochondrial matrix or the plastid stroma are similar to the bacterial leader sequences that trigger the secretion of prokaryotic proteins (77, 78). They are removed by a **proteinase** that is **homologous** to the bacterial leader peptidase (79, 80). Even more striking, translocation of proteins from the stroma into the thylakoids occurs via translocation sites, containing membrane proteins homologous to those engaged in the protein export from the bacterial cytoplasm (81, 82) (see Table 1).

## Bibliography

1. S. A. Hayes and J. F. Dice (1996) *J. Cell Biol.* **132**, 255–258.
2. J. E. Rothman and J. Lenard (1977) *Science* **195**, 743–747.
3. G. von Heijne and Y. Gavel (1988) *Eur. J. Biochem.* **174**, 671–678.
4. M. Spiess (1995) *FEBS Lett.* **369**, 76–79.
5. J. Lipp and B. Dobberstein (1986) *Cell* **106**, 1813–1820.
6. S. Subramani (1993) *Annu. Rev. Cell Biol.* **9**, 445–478.
7. G. von Heijne (1987) *Biochim. Biophys. Acta* **947**, 307–333.
8. M. G. Carlos, S. Brunak, and G. von Heijne (1997) *Curr. Opin. Struct. Biol.* **7**, 394–398.
9. S. Kornfeld and I. Mellman (1989) *Annu. Rev. Cell Biol.* **5**, 483–525.
10. S. K. Powel et al. (1991) *Nature* **353**, 76–77.
11. T. J. Baranski, P. L. Faust, and S. Kornfeld (1990) *Cell* **63**, 281–291.

12. I. W. Caras (1991) *J. Biol. Chem.* **113**, 77–85.
13. T. Lithgow, J. M. Cuezva, and P. A. Silver (1997) *Trends Biochem. Sci.* **22**, 110–113.
14. A. G. Pugsley (1993) *Microbiol. Rev.* **57**, 50–108.
15. W. Wickner (1994) *Science* **266**, 1197–1198.
16. V. Koronakis and C. Hughes (1993) *Semin. Cell Biol.* **4**, 7–15.
17. C. J. Hueck (1998) *Microbiol. Molecular Biol. Rev.* **62**, 379–389.
18. C. L. Santini et al. (1998) *EMBO J.* **17**, 101–112.
19. A. Müsch et al. (1990) *Trends Biochem. Sci.* **15**, 86–88.
20. D. J. Klionsky, R. Cueva, and D. S. Yaver (1992) *J. Cell Biol.* **119**, 287–299.
21. P. Walter, I. Ibrahim, and G. Blobel (1981) *J. Cell Biol.* **91**, 545–550.
22. T. Kurzchalia et al. (1986) *Nature* **320**, 634–636.
23. R. Gilmore, P. Walter, and G. Blobel (1982) *J. Cell Biol.* **95**, 470–477.
24. B. Jungnickel and T. A. Rapoport (1995) *Cell* **82**, 261–270.
25. D. Görlich et al. (1994) *Cell* **79**, 767–778.
26. V. W. Pollard et al. (1996) *Cell* **86**, 985–994.
27. T. Komiyama et al. (1994) *J. Biol. Chem.* **269**, 30893–30897.
28. T. Komiyama et al. (1997) *EMBO J.* **16**, 4267–4275.
29. A. Mayer et al. (1995) *EMBO J.* **14**, 4204–4211.
30. P. Rehling, M. Albertini, and W. H. Kuhnau (1996) *Ann. N. Y. Acad. Sci.* **804** 34–46.
31. W. J. Chirico, M. G. Waters, and G. Blobel (1988) *Nature* **332**, 805–810.
32. T. Komiyama, M. Sakaguchi, and K. Mihara (1996) *EMBO J.* **15**, 399–407.
33. P. A. Walton et al. (1994) *J. Cell Biol.* **125**, 1037–1046.
34. P. Whitley, I. Nilsson, and G. von Heijne (1996) *J. Biol. Chem.* **271**, 6241–6244.
35. M. Eilers and G. Schatz (1986) *Nature* **322**, 228–232.
36. S. Sanders et al. (1992) *Cell* **69**, 353–365.
37. P. J. Schatz and J. Beckwith (1990) *Annu. Rev. Genet.* **24**, 215–248.
38. D. Walker et al. (1996) *J. Biol. Chem.* **271**, 4082–4085.
39. S. C. Hinnah et al. (1997) *EMBO J.* **16**, 7351–7360.
40. S. A. Clark and S. M. Theg (1997) *Mol. Biol. Cell* **8**, 923–934.
41. A. S. Gaikward and M. G. Cumsky (1994) *J. Biol. Chem.* **269**, 6437–6443.
42. F. M. Valette et al. (1994) *J. Biol. Chem.* **269**, 13367–13374.
43. D. Hanein et al. (1996) *Cell* **87** 721–732.
44. R. Beckmann et al. (1997) *Science* **278**, 2123–2126.
45. D. Görlich and T. A. Rapoport, *Cell* **75**, 615–630.
46. D. G. Millar and G. C. Shore (1996) *J. Biol. Chem.* **271**, 25823–25829.
47. K. S. Crowley et al. (1994) *Cell* **78**, 461–467.
48. B. Martoglio et al. (1995) *Cell* **81**, 207–214.
49. J. Berthold et al. (1995) *Cell* **81** 1085–1093.
50. M. Horst et al. (1995) *EMBO J.* **14**, 2293–2297.
51. C. Sirrenberg et al. (1997) *J. Biol. Chem.* **272**, 29963–29966.
52. S. V. Scott and S. M. Theg (1996) *J. Cell Biol.* **132**, 63–75.
53. S. Panzner et al. (1995) *Cell* **81**, 561–570.
54. H. Schneider et al. (1994) *Nature* **371**, 768–774.
55. G. Blobel and B. Dobberstein (1975) *J. Cell Biol.* **67**, 852–862.

56. G. Hawlitschek et al. (1988) *Cell* **53**, 795–806.
57. E. E. Rojo, R. Stuart, and W. Neupert (1995) *EMBO J.* **14**, 3445–3451.
58. J. S. Bonifacio and R. D. Klausner (1991) *Curr. Opin. Cell Biol.* **3**, 592.
59. R. Peters (1986) *Biochim. Biophys. Acta* **864**, 305–359.
60. Q. Yang, M. P. Rout, and C. W. Akey (1998) *Mol. Cell* **1**, 223–234.
61. S. V. Scott et al. (1997) *J. Cell Biol.* **138**, 37–44.
62. W. H. Kuhnau and R. Erdmann (1998) *Curr. Biol.* **8**, R299–R302.
63. J. E. Rothman and F. T. Wieland (1996) *Science* **272**, 227–234.
64. S. Pfeffer and J. E. Rothman (1987) *Annu. Rev. Biochem.* **56**, 829–852.
65. K. Matter and I. Mellmann (1994) *Curr. Opin. Cell Biol.* **6**, 545–554.
66. S. Munro and H. R. B. Pelham (1987) *Cell* **48**, 899.
67. M. Jackson, T. Nilsson, and P. A. Peterson (1990) *EMBO J.* **9**, 3153–3163.
68. I. S. Trowbridge, J. F. Collawn, and C. R. Hopkins (1993) *Annu. Rev. Cell Biol.* **9**, 129–161.
69. P. Cosson and F. Letourneur (1994) *Science* **263**, 1629–1631.
70. J. Dietrich et al. (1994) *EMBO J.* **13**, 2156–2166.
71. M. J. Lewein, D. J. Sweet, and H. R. B. Pelham (1990) *Cell* **61**, 1359.
72. S. Kornfeld (1992) *Annu. Rev. Biochem.* **61**, 307–330.
73. S. Ogata and M. Fukuda (1994) *J. Biol. Chem.* **269**, 5210–5217.
74. M. S. Bretscher and S. Munro (1993) *Science* **261**, 1280–1281.
75. S. Munro (1995) *EMBO J.* **14**, 4695–4704.
76. J. Lubeck, L. Heins, and J. Soll (1997) *J. Cell Biol.* **137**, 1279–1286.
77. G. von Heijne, J. Steppuhn, and R. G. Herrmann (1989) *Eur. J. Biochem.* **180**, 535–545.
78. A. M. Chaddock et al. (1995) *EMBO J.* **14**, 2715–2722.
79. J. Nunnari, T. D. Fox, and P. Walter (1993) *Science* **262**, 1997–2004.
80. B. K. Chaak et al. (1998) *J. Biol. Chem.* **273**, 689–692.
81. L. M. Roy and A. Barkan (1998) *J. Cell Biol.* **141**, 385–395.
82. A. M. Settles et al. (1997) *Science* **278**, 1467–1470.

### **Suggestions for Further Reading**

83. R. Lill, F. E. Nargang, and W. Neupert (1996) *Curr. Opin. Cell Biology* **8**, 505–512.
84. T. A. Rapoport, B. Jungnickel, and U. Kutay (1996) *Annu. Rev. Biochem.* **65**, 271–303.
85. S. V. Scott, and D. Klionsky (1998) *Curr. Opin. Cell Biol.* **10**, 523–529.
86. P. Walter and A. E. Johnson (1994) *Annu. Rev. Cell Biol.* **10**, 87–119.

### **Protein-DNA Recognition**

The binding of a protein to a specific DNA sequence is largely dependent on two types of interaction. The principal basis for sequence selectivity is direct contact between the polypeptide chain and the exposed edges of the base pairs, primarily in the major groove of B-DNA. These contacts may involve either [hydrogen bonds](#) or [van der Waals interactions](#). Small molecules, such as [water](#) molecules, which are tightly and rigidly bound to a protein and are thus integral components of

the **macromolecular** structure, may also participate in these interactions and so provide binding specificity to the protein by proxy. These direct interactions are supplemented by the sequence-dependent bendability or deformability of DNA which limits the energetically favorable conformations of a particular binding site, and thereby, imposes additional sequence-dependent constraints on the binding affinity. Relative to direct contacts to the base pairs, this is a second-order effect which modulates the affinity for the whole binding site.

The binding energy available from direct interactions with the base pairs, although significant, is not in general sufficient by itself to allow the formation of a stably bound complex for binding sites of average length (6–15 bp). The required additional binding energy may be provided by direct electrostatic interactions between basic [amino acid](#) residues and the negatively charged sugar-phosphate backbone. The spatial constraints imposed by this type of interaction may also serve to constrain the configuration of the DNA when bound to protein. It is the difference between the binding energies for the sequence-dependent and sequence-independent components of the interaction that is the measure of the sequence selectivity of a DNA-binding protein.

In many DNA-protein complexes, the DNA is substantially bent. This is achieved by a variety of mechanisms, notably the induction of bends by spatial constraint on a rigid protein surface (the histone octamer, FIS), the insertion of hydrophobic residues between adjacent bases in the DNA duplex (HMG domain proteins, TBP, lac repressor) and charge neutralization on one face of the double helix (histone octamer, CAP).

Typically a DNA sequence recognition motif recognizes a sequence of 3–4 bp. This is insufficient to allow highly selective discrimination among all sequences in a genome. In practice the effective site size for recognition can be increased by use of a larger protein assembly. This assembly is often a stable dimer (as in helix-turn-helix and bZip proteins). However, in other cases cooperative interactions between proteins bound at contiguous or distant DNA sites are required for the formation of an assembly. One example of such interactions is the cooperative binding of  $\lambda$  C<sub>I</sub> repressor dimers to the leftward and rightward operators in  $\lambda$  DNA, each of which contains three C<sub>I</sub> binding sites. Because stable binding to any two of these sites is dependent on interaction between separate dimers, occupation is sensitive to small changes in concentration of repressor molecules over a certain range. Interactions can also occur between distant binding sites on the DNA such that simultaneous occupation generates a loop of intervening DNA. Loop formation of this type need only require a single stable protein assembly (as in the case of the tetramer of the *lac* repressor containing four helix-loop-helix motifs) or may require cooperative interactions between proteins bound at the separate sites (as have been postulated for interactions between eukaryotic enhancer and promoter elements).

## **Protein–Protein Interactions**

Noncovalent protein-protein interactions are ubiquitous in cells and living organisms. They are the basis for building every biological structure, from the nanometer scale of the molecule to the micrometer scale of the whole cell, and above. In the [nucleus](#), the [nucleosome](#) that holds eukaryotic **DNA** in a compact form, and the machinery that **replicates** and **transcribes** it, are noncovalent protein-protein complexes. In the cell [membrane](#), the **channels**, **pumps** and energy-producing systems are assemblies of proteins held together by noncovalent forces. So are the coats of **viruses** and the fibers of muscles. These particular assemblies are permanent and made of **polypeptide chains** that, under normal conditions, are found only as complexes. Such permanent interactions often occur concomitantly with folding of the individual subunits. Other protein–protein interactions

are readily reversible and require preformed partners to meet and recognize each other. The [quaternary structure](#) of [oligomeric proteins](#) is an example of a permanent interaction, whereas [antibody](#) recognition of a native protein [antigen](#) is an example of folding and association as two separate steps. The noncovalent forces involved are the same in the two cases, and the many intermediate situations of all kinds make the distinction between permanent and nonpermanent assembly somewhat arbitrary. Nevertheless, oligomeric proteins are dealt with in a separate entry, and here we concentrate on the pairwise and reversible recognition of preformed proteins in solution, a simple example of a protein–protein interaction. In this particular case, the crucial concepts of **affinity** and specificity can be given a structural and a **thermodynamic** basis, and we may trust that the rules that apply here will also apply to more elaborate multi-component assemblies.

## 1. Detecting Protein–Protein Interactions

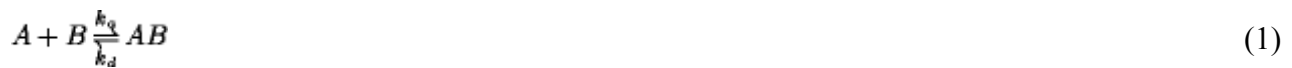
It is often possible to detect the interaction of two proteins initially by genetic methods. [Suppressor mutations](#) are indicators of a protein–protein interaction. If a suppressor mutation occurs in a gene other than that containing the original defect, it may be that the two gene products form a functional complex. The first mutation destabilizes the interaction; the second compensates and restores it. Dominant-negative [mutations](#) are another indicator; in heterozygous cells in which normal and mutant polypeptide chains are expressed, the two **alleles** can interact—if the mutant chains form an inactive assembly with the wild type chains, this will lead to a negative **phenotype**. On the other hand, there are alternative explanations for both suppressor and dominant-negative mutations, and neither interpretation is general. In contrast, the very popular yeast [two-hybrid system](#) was specifically designed for detecting *in vivo* partner proteins. The test relies on the observation that transcriptional activator proteins in **eukaryotes** often consist of two independent folding units, a DNA-binding domain and a transcription-activation domain, that must be part of the same molecular assembly but need not be covalently bonded. To detect an interaction between proteins A and B, a genetic construction is prepared in which A is linked to the DNA-binding domain and B, to the activation domain. If A and B form a stable complex in the cell, the two transcriptional activator domains will be brought together, and transcriptional activation of a [reporter gene](#) will be detected. In a commonly used construction, the A gene fused to the sequence coding for the DNA-binding domain of the *Escherichia coli* [LexA repressor](#). Then, a gene library is used to fuse a large number of unknown B proteins to the activation domain of the yeast Gal4 activator. Yeast colonies that express the reporter gene under LexA regulation contain hybrid genes for proteins that interact with A *in vivo*, and they can be isolated and characterized further. The method requires the two hybrid genes to be expressed, produce proteins that fold and enter the [nucleus](#), and interact in a way that brings the DNA-binding domains and activation domains into the appropriate proximity; thus it has a high incidence of false negatives. Nevertheless, it is so simple and powerful that it is possible in principle to search for pairwise interactions between all proteins encoded in the *Saccharomyces cerevisiae* genome.

In contrast to the genetic approach, biochemical methods require the purification of at least one of the partners of the putative complex. Once purified, partner A can be chemically coupled to activated [agarose](#) or another matrix support. The derivatized matrix is then used to isolate interacting proteins from cell extracts by [affinity chromatography](#). Alternatively, antibodies prepared against the purified protein, preferably [monoclonal antibodies](#), can be used to precipitate it, together with any protein it binds to. Such coimmunoprecipitation has many successes to its credit, for example, isolation of the **retinoblastoma** protein from precipitates obtained with antibodies against a protein from [adenovirus](#). In this case, false positives are common, and an interaction that is detected may involve several partners, rather than being direct. Further characterization requires the identification of all the major components in the coprecipitate, and eventually their purification to homogeneity, opening the way to a detailed structural and functional analysis.

## 2. Affinity and Stoichiometry of Protein–Protein Complexes

Two parameters characterize protein–protein interaction at equilibrium: the stoichiometry and the

affinity of the partners for each other. The one-to-one association reaction of two components A and B may be written as follows:



where  $k_a$  is the rate constant for the bimolecular association reaction and  $k_d$  is that for the monomolecular dissociation. Their ratio is equal to either  $K_a$ , the **association constant**, or  $K_d$ , the **dissociation constant**. This is related to the concentrations of A, B, and AB at thermodynamic equilibrium by the law of mass action:

$$\frac{[A][B]}{[AB]} = \frac{1}{K_a} = K_d = \frac{k_d}{k_a} \quad (2)$$

where  $K_d$  has the dimension of concentration. Its value determines the stability of the complex AB at equilibrium; therefore, it is a measure of the affinity of the two components for each other. If  $x = [A]$ , the equilibrium concentration of free A, and  $y = [AB]$ , the concentration of bound A, Equation (2) can be written

$$\frac{y}{x} = \frac{[B]}{K_d}$$

Assuming B to be in excess over A, a high affinity implies that most of A is in the bound state. In other words, the  $y/x$  ratio is large, and therefore  $K_d$  must be small relative to [B]. Affinity is essentially a relative concept, not an intrinsic property, as is often assumed. For a [hormone receptor](#), a ligand with a  $K_d$  of 1  $\mu M$  can be considered to have low affinity and unlikely to be physiologically relevant, because hormone concentrations are usually much less than micromolar. Yet, this is approximately the  $K_d$  value for the dissociation of human [hemoglobin](#) into ab dimers when oxygen is bound, which is usually considered a stable complex. In red blood cells, the hemoglobin concentration is well above millimolar, and oxyhemoglobin remains a tetramer, whereas it would readily dissociate in a dilute solution.

Determination of the dissociation constant and of the stoichiometry of reaction is an important step in the study of protein–protein interaction. It relies on measuring the equilibrium concentrations of the complex and its components. This generally requires having pure proteins and a means to detect their association at concentrations of the order of the dissociation constant. When  $K_d$  is micromolar or greater, a wide range of physical and chemical methods are available, among which [fluorescence quenching](#), [sedimentation equilibrium centrifugation](#), and [microcalorimetry](#) may be cited. The data may be obtained by titrating a fixed concentration of component B with increasing concentrations of component A. They are analyzed by drawing the binding isotherm, which relates the concentration  $y$  of bound A to that of the free species  $x$ :

$$y = [B_0] \frac{x}{x + K_d} \quad (3)$$

Here,  $[B_0]$  is the total concentration of B, or more appropriately, the total concentration of binding sites on B. Equation 3 is derived from the law of mass action (Eq. 2) by noting that  $[B_0] = [B] + [AB]$ . It can be rearranged to yield a linear [Scatchard Plot](#) in which  $y/x$  is plotted against  $y$ :

$$\frac{y}{x} = \frac{1}{K_d}([B_0] - y) \quad (4)$$



The Scatchard plot has slope  $-1/K_d$  and intersects the horizontal axis at  $y = [B_0]$ , from which the stoichiometry may be derived.

When  $K_d$  is much less than micromolar, more sensitive detection methods are needed. Binding isotherms may be obtained by **radiolabeling** one of the components, or by indirect titration with antibodies. However, these methods are also inadequate when  $K_d$  is in the subnanomolar range. For such very tight complexes, kinetic measurements are a useful alternative to equilibrium methods; the two rate constants in Equation (2) are measured directly, and the dissociation constant is derived as their ratio. This particular relationship between rate and equilibrium constants is valid only if the reaction mechanism is that of Equation (1). More elaborate mechanisms are possible, including mono-molecular steps that precede or follow the bimolecular association step. These may in principle be detected in kinetic experiments by exploring the concentration dependence of the rate of the reaction and analyzing its departure from second-order kinetics.

Assuming Equation (1) to be a valid approximation, the two rate constants can be measured in the same or in separate experiments. Mixing experiments yield  $k_a$ : B is added in excess to A, and pseudo-first-order association kinetics are measured. They should follow a single exponential curve with a concentration-dependent relaxation rate  $1/\tau_a$ :

$$\frac{1}{\tau_a} = k_a[B] + k_d$$

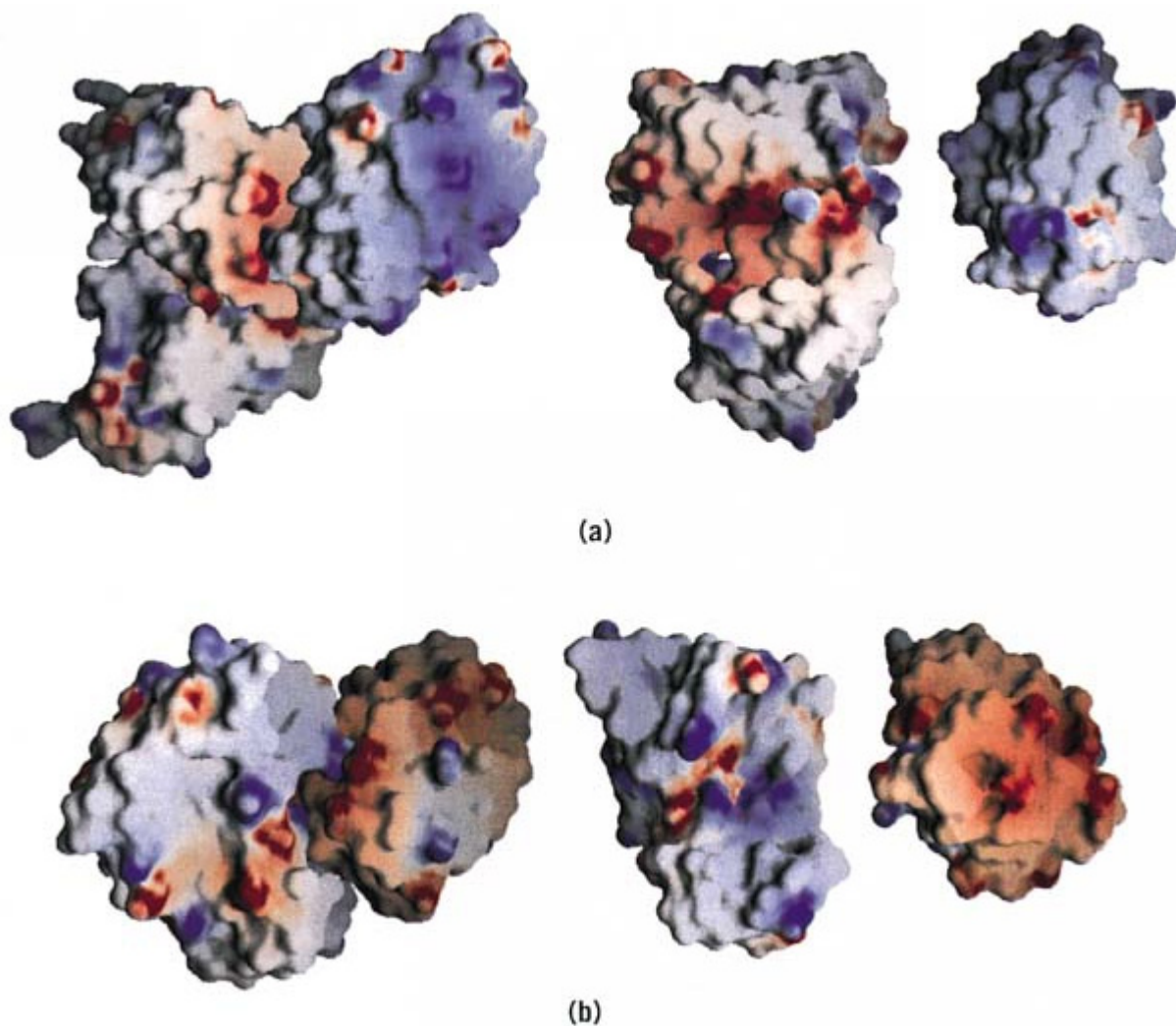
A plot of  $1/\tau_a$  versus  $[B]$  should be linear and yield  $k_a$  as the slope and  $k_d$  by extrapolation to  $[B] = 0$ , although such an extrapolation is often unreliable. Most observed  $k_a$  values for association of two protein molecules are in the range  $10^5 - 10^9 M^{-1} s^{-1}$ . This implies that  $\tau_a$  is less than a second when  $[B]$  is greater than micromolar, so the experiment must be performed in a rapid-mixing device, such as a stopped-flow apparatus. To determine  $k_d$  properly, the dissociation kinetics could in principle be followed by diluting the complex below  $K_d$ , but this is not practical even when the required sensitivity can be achieved, because measuring subnanomolar protein concentrations is not reliable. Thus, the experiment is usually performed by chasing a radioactive ligand bound to the complex with an excess of cold ligand. The rate of release of the radioactivity yields  $k_d$ . Values observed in protein-protein complexes cover a wide range:  $10^{-6} - 10^3 s^{-1}$ , which implies that dissociation is completed in milliseconds for the weaker complexes, whereas it takes days or months for the tighter ones.

Deriving the value of  $K_d$  from kinetics instead of equilibrium studies can be convenient in other circumstances. **Surface plasmon resonance** measurements are well suited to kinetic studies of association-dissociation reactions. This technique uses a very sensitive optical method to detect changes of the refractive index at the surface of a gold chip placed in a continuous flow of buffer. First, purified partner A is chemically immobilized on the surface of the chip. Then partner B is added to the buffer, and its binding to the immobilized component is detected by recording the optical signal due to the change in refractive index of the chip. Partner B needs not to be pure so long as no other component of the solution binds to A, but in all cases estimates of the stoichiometry must take into account possible artifacts of the immobilization procedure. After binding is completed, pure buffer is used, and the dissociation of partner B from the complex on the chip yields a signal of opposite sign. As a result of the continuous flow of material or buffer, the system cannot reach true equilibrium, only a steady state. Still, provided the buffer flow rate is low compared to the reaction rate, the steady state should be sufficiently close to equilibrium for the measurement to yield approximately correct kinetic and equilibrium parameters.

### 3. The structural basis for recognition

Many noncovalent protein–protein complexes have been analyzed by [X-ray crystallography](#) and some by high-resolution [NMR](#). Data sets of atomic coordinates, available through the Brookhaven Protein Data Bank (PDB), provide a structural basis for understanding protein association [1](#), [2](#). Other than oligomeric proteins, the best-represented examples of protein–protein interactions in the PDB are the complexes between antibodies and protein antigens and between [proteinases](#) and protein [proteinase inhibitors](#). There are also examples of [enzymes](#) other than proteinases that form complexes with protein substrates or inhibitors. Figure [1a](#) shows hen egg [lysozyme](#) bound to the Fab fragment of monoclonal antibody HyHEL5, a typical example of an antibody–antigen complex ([3](#)); in Figure [1b](#), the bacterial ribonuclease **barnase** is bound to the protein inhibitor barstar ([4](#)). The dissociation constants are of the order of  $10^{-10}M$  for lysozyme-HyHEL5 and  $10^{-13}M$  for barnase-barstar. Both X-ray crystallography structures illustrate the tight and highly specific association of two proteins that fold independently of each other and, for lysozyme and barnase, at least, perform their biological functions in the unbound state.

**Figure 1.** Examples of antibody–antigen and enzyme–inhibitor complexes. On the left, the complexes are shown as they are in the X-ray structure; on the right, the components are separated by opening the complex like a book and exposing the two protein surfaces that form the interface. The surface is colored according to its electric charge—red is negative; blue, positive; white, neutral. **(a)** Hen egg lysozyme complexed to the Fv part of monoclonal antibody HyHEL5 (file 3hfl in the Protein Data Bank). The combining site of the antibody (middle panel) forms a cleft lined with negative charges, in which part of the positively charged lysozyme surface (right panel) fits ([3](#)). **(b)** Complex of the bacterial ribonuclease barnase with its natural inhibitor barstar (PDB file 1brs). The inhibitor (right panel) is negatively charged and covers the active site of the enzyme (middle panel), which carries several positive charges ([4](#)). (Figure made with GRASP, K. Sharp and B. Honig, Columbia University, New York). See color insert.



The individual structures of lysozyme, barnase, and barstar are known. They differ from those seen in the complexes by only small changes, typically no greater than 1 Å ( $1 \text{ \AA} = 10^{-10} \text{ m}$ ) in amplitude, in the polypeptide main-chain fold and in the conformations of a number of amino acid [side chains](#). Thus, the partner molecules of the complexes shown in Figure 1 associate to a good approximation as rigid bodies. Their interface involves two complementary surfaces that are largely preformed. Quasi rigid-body association of this type is frequent in protein–protein recognition, but is definitely not general. The PDB shows examples where significant or major conformational changes accompany association. A case in point is the complex formed between **trypsinogen** and bovine pancreatic trypsin inhibitor ([BPTI](#)). Its X-ray structure is essentially identical to that of the trypsin–BPTI complex, a very high affinity complex in which the two components undergo little conformation change. Yet free trypsinogen is not like free trypsin; it contains disordered loops of polypeptide chain that constitute its activation domain. These loops are fully ordered in the mature enzyme, trypsin, and in the trypsinogen–BPTI complex (5). Thus, BPTI binding induces a disorder-to-order transition in trypsinogen that converts it to the mature enzyme conformation, and this change is reflected in its lower affinity than trypsin for BPTI.

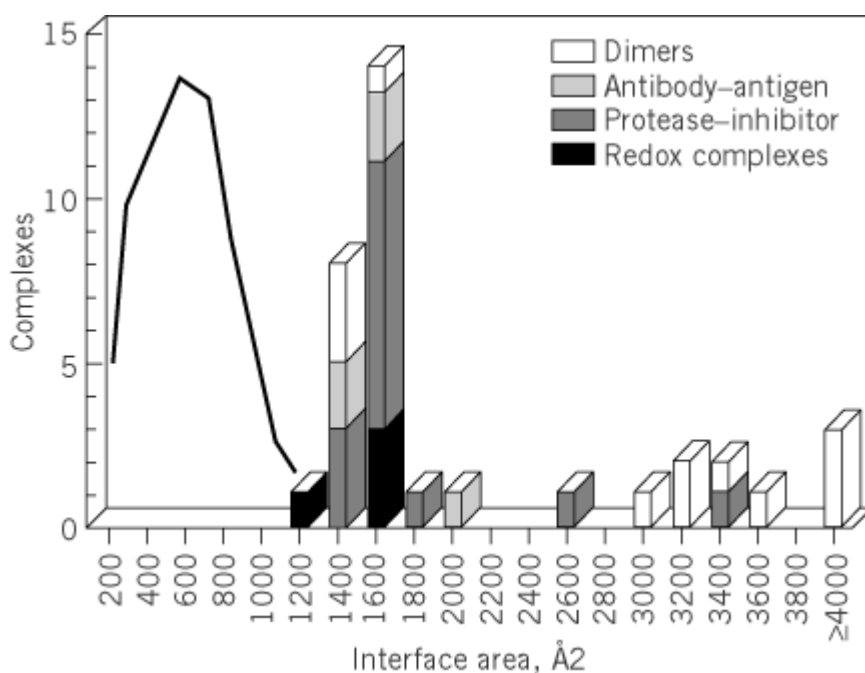
In addition to specific complexes of the kind just discussed, the PDB contains many examples of weaker, nonspecific protein–protein interactions: those that build protein crystals. Their physicochemical basis is the same as that for antibody–antigen or enzyme–inhibitor association, but crystal contacts are not subjected to [natural selection](#) by [evolution](#) or by the immune system. From a structural point of view, the obvious difference is that crystal contacts are less extensive than in

specific complexes. On the basis of their atomic coordinates, the extent of the contact between molecules A and B can be measured from their interface area:

$$B = A_A + A_B - A_{AB} \quad (5)$$

where  $A_{AB}$  is the solvent-accessible surface area of the complex,  $A_A$  and  $A_B$  are those of the dissociated components, and  $B$  represents the area of the protein surface that becomes buried at their interface when the two proteins associate. Its value is proportional to the number of pairs of atoms from the two components that are in contact in the complex. Figure 2 (6) is a histogram of interface areas observed in the PDB. For crystal packing contacts, the distribution is broad and centered at  $B \approx 500 \text{ \AA}^2$ . For protein-protein complexes, there is a well-defined peak centered at  $B \approx 1500 \text{ \AA}^2$ . No interface area is less than  $1000 \text{ \AA}^2$ , yet there are interfaces burying  $2000 \text{ \AA}^2$  and more. The most plausible interpretation of the histogram is that any stable, specific protein-protein association requires removing approximately  $1500 \text{ \AA}^2$  from contact with water,  $750 \text{ \AA}^2$  per partner molecule, whereas the random contacts illustrated by crystal packing bury much less. A  $1500\text{-\AA}^2$  interface typically involves 15–20 residues from each partner, with widely varying contributions of individual residues. A single [arginine](#) or [tryptophan](#) residue can lose more than  $150 \text{ \AA}^2$  of accessible surface area. Arginine is almost never buried inside protein monomers, yet buried arginine residues are not uncommon at protein-protein interfaces, and they are major contributors to both the interface area and the polar interactions. Arginine and the [nonpolar](#) residue [leucine](#) are actually the two most abundant residues in protein-protein interfaces, a finding that indicates that polar interactions (illustrated by arginine) and the [hydrophobic effect](#) (by leucine) are nearly equally important in protein-protein interactions, as they are in [protein stability](#).

**Figure 2.** Histogram of interface areas in protein-protein complexes. The data are for 5 antigen-antibody, 15 proteinase-inhibitor and 3 redox-protein complexes, one of which is the cytochrome c-cytochrome c peroxidase complex. The interface areas observed in 12 homodimeric proteins are included, to show that most homodimers have larger interfaces than do protein-protein complexes. The thin line is a similar histogram for the 1320 pairwise interface areas of crystal contacts in a sample of monomeric protein crystals. These small interfaces illustrate weak and more nonspecific protein-protein interactions (6).



A few protein–protein complexes have small interface areas near  $1200 \text{ \AA}^2$ ; **cytochrome c**–cytochrome c [peroxidase](#) is an example. This is a [redox enzyme](#)–substrate complex that must be stable transiently for the electron transfer reaction to occur; the affinity cannot be very high in such a case. Yet beyond that qualitative statement, there is no simple correlation between the value of  $B$  and the affinity, for the major peak at  $B \approx 1500 \text{ \AA}^2$  includes complexes with  $K_d$  values ranging from  $10^{-7}$  to  $10^{-13} \text{ M}$ . On the other hand, there is a relation between interface area, affinity, and the occurrence of conformation changes. The lysozyme–HyHEL5, barnase–barstar, and trypsin–BPTI complexes, all of which associate as quasi-rigid bodies, have interface areas near  $1500 \text{ \AA}^2$ . The trypsinogen–BPTI complex is similar, yet its dissociation constant is six orders of magnitude greater than that for trypsin. This can be safely attributed to the conformation change in trypsinogen. Components of the complexes that have very large interfaces often undergo major conformation changes. An example is [hirudin](#), from the medicinal leech, binding to the blood proteinase [thrombin](#). The complex has an interface area of  $3300 \text{ \AA}^2$ , half of which is due to the *C*-terminal tail of hirudin filling a fibrinogen-binding site on the proteinase that is distinct from its [active site](#) (7). This tail is disordered in free hirudin. Like the trypsinogen–BPTI association mentioned above, association of hirudin–thrombin causes a disorder-to-order transition in hirudin. From that point of view, it resembles the assembly of subunits in oligomeric proteins, where the subunits generally form large interfaces and subunit folding is tightly coupled with oligomer assembly.

Interfaces in the lysozyme–HyHEL5, barnase–barstar, and trypsin–BPTI complexes comprise a nonpolar component, mostly from amino acid [side chains](#), and a polar component, originating from both the main chain and the side chains. The nonpolar component contributes to stabilizing the complexes through the hydrophobic effect, the polar component through [hydrogen bonds](#). The relative proportion of the nonpolar and polar components of the interface area is about the same in the three examples, and it is also very close to that on the average protein surface, which is 55% nonpolar and 45% polar. In other complexes, there are indications that nonpolar patches of protein surface serve as a binding site. Still, the polar fraction of the interface area is always significant and in the range 30–50%. With very few exceptions, the many polar groups that are buried form hydrogen bonds. The number of hydrogen bonds per unit buried surface area is similar in all complexes, about one per  $150 \text{ \AA}^2$ . Thus, the complexes forming the major peak at  $B \approx 1500 \text{ \AA}^2$  have  $10 \pm 4$  hydrogen bonds bridging their components. In addition to direct hydrogen bonds between protein groups, one finds [water](#) molecules bridging polar groups of the two proteins, usually at the periphery of the contact region. In the lysozyme–HyHEL5, barnase–barstar and trypsin–BPTI examples, a majority of both the direct and indirect hydrogen bonds are [salt bridges](#) or charged bonds; that is, hydrogen bonds where either the donor group, the acceptor group, or both in a salt bridge, is an anion or a cation. These features are observed in many other protein–protein complexes (8).

#### 4. Thermodynamics of association

Thermodynamic state functions for the dissociation reaction are related to the equilibrium constant  $K_d$  and its temperature derivative (Table 1). The equilibrium constant yields the standard-state **free enthalpy** of dissociation,  $DG_d$ , whereas the temperature derivative gives the **enthalpy** of dissociation,  $DH_d$ , by the van't Hoff law. Note that the standard-state concentration of  $1 \text{ M}$  is a convention and that the values of the free energy and **entropy** changes depend on it. The temperature dependence of the affinity yields  $DH_d$  from the slope of a van't Hoff plot (plot of  $\ln K_d$  vs  $1/T$ ). The van't Hoff plot need not be linear, because  $DH_d$  can be temperature-dependent. This dependence is expressed in the **heat capacity** of dissociation at constant pressure,  $DC_d$ . In recent years, the development of isothermal titration calorimetry (ITC) has made it possible to measure the enthalpy change directly. Sensitive microcalorimeters measure the heat of mixing of the two components,

which is equal to  $-DH_d$  per mole of complex. Provided the reaction mechanism is that of Equation (1), the bimolecular association of two free subunits, the calorimetric enthalpy value should be the same as the one derived by applying the van't Hoff law to the temperature dependence of  $K_d$ .

**Table 1. Thermodynamical State Functions**

|  |  |
|--|--|
| Standard-state <sup>a</sup> free-energy change | $DG_d = DH_d - TDS_d$<br>$= -RT \ln \frac{K_d}{c_0}$                 |
| Enthalpy change (van't Hoff)                   | $\Delta H_d = R \frac{d(\ln K_d)}{d(1/T)}$                           |
| Entropy change                                 | $\Delta S_d = -\frac{d(\Delta G_d)}{dT}$                             |
| Heat-capacity change                           | $\Delta C_d = \frac{d(\Delta H_d)}{dT} = T \frac{d(\Delta S_d)}{dT}$ |
| Gas constant                                   | $R \approx 2 \text{ cal}\cdot\text{mol}^{-1}\cdot\text{K}^{-1b}$     |

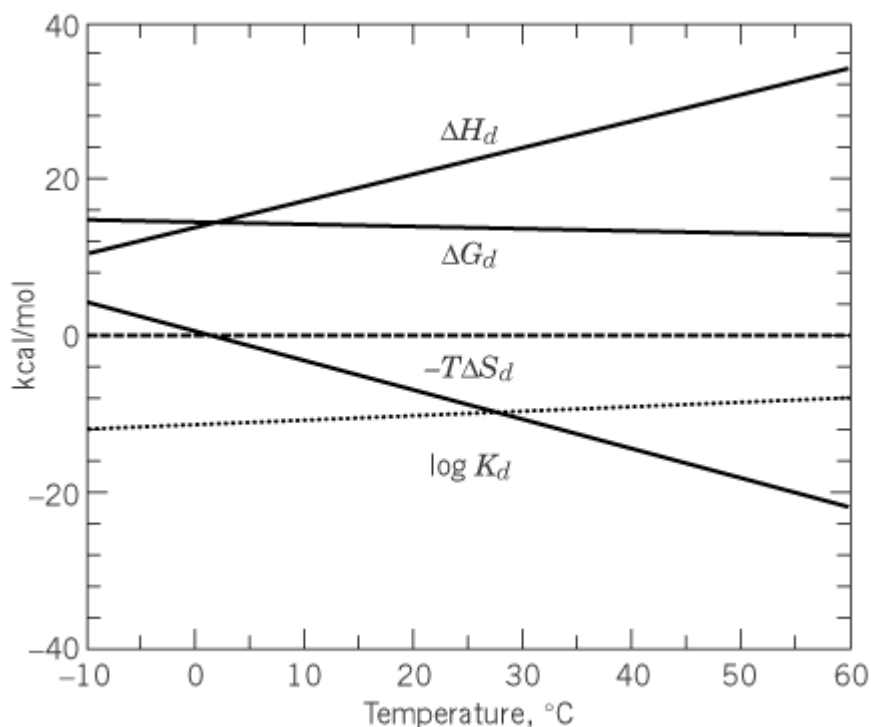
<sup>a</sup> Standard state: pressure  $p_0 = 1 \text{ bar}$ , concentration  $c_0 = 1 \text{ M}$  (1 mol/L).

<sup>b</sup> 1 cal = 4.19 J.

Biologically relevant protein–protein interactions may have extremely different stabilities. The value of  $K_d$  ranges from at least millimolar to less than picomolar concentrations, and therefore the value of  $DG_d$  varies from 4 to 17 kcal/mol. The relative contributions of  $DH_d$ , the enthalpic contribution to  $DG_d$  and of  $-TDS_d$ , the entropic contribution, depend greatly on the system. This also depends strongly on the temperature, for calorimetric data indicate that  $DC_d$  is usually positive and large.

This is illustrated by Figure 3 (9) for the lysozyme-HyHEL5 monoclonal antibody complex. At 25°C,  $K_d$  is less than 1 nM, a value rather typical of such complexes. ITC measurements demonstrate  $DH_d$  to be large and positive at 25°C, sufficient to explain the stability of the complex. Thus, lysozyme-HyHEL5 association is enthalpy-driven: a favorable enthalpy change overcomes the destabilizing effect of the positive entropy change. On the other hand, calorimetric analysis of the lysozyme–HyHEL5 complex also yields a large positive value of the heat capacity change,  $DC_d$ . As a consequence, both  $DH_d$  and  $DS_d$  become smaller at lower temperatures. These variations compensate each other in  $DG_d$ , which undergoes a much smaller change;  $DS_d$  even changes sign near 0°C, so that at lower temperatures, lysozyme–HyHEL5 association would be entropy-driven. The change of sign does not imply a different mechanism; it is simply a consequence of the temperature dependence of the state functions. Protein–protein complexes with negative (destabilizing) dissociation enthalpies and negative (stabilizing) dissociation entropies at 25°C also exist. They do not differ in any fundamental way from enthalpy-driven complexes, such as lysozyme-HyHEL5. The large positive  $DC_d$ , which leads to a strong temperature dependence of the enthalpy and entropy, is attributable largely to the hydrophobic effect. When a complex dissociates, many nonpolar chemical groups that were buried at the interface become **hydrated**. Nonpolar group hydration has a large positive heat capacity and is the dominant term in  $DC_d$ . Other components of  $DH_d$  may be greater in absolute value than nonpolar group hydration, yet they are much less temperature-sensitive and therefore contribute less to the heat-capacity change. For example, polar groups at the interface also become hydrated on dissociation, but polar-group hydration has very little heat capacity, even though its enthalpy is large.

**Figure 3.** Thermodynamics of an antibody–antigen complex, that of lysozyme and antibody HyHEL5 (9). The dependences on temperature of the  $\Delta H_d$  and  $\Delta C_d$  state functions were derived from calorimetric measurements, whereas  $\Delta G_d$  and  $-T\Delta S_d$  were measured by the temperature dependence of the dissociation constant,  $K_d$ . (1 cal = 4.19 J).



## 5. Site-Directed Mutagenesis of Protein–Protein Interactions

In many protein-protein complexes, only 15–20 amino acid residues of each partner contribute directly to the interface, and if the structure is known, it is feasible to substitute each one by [site-directed mutagenesis](#) and to test how this affects the stability of the complex. Such measurements have been made on several systems, including antigen–antibody complexes, the human [growth hormone](#)–growth hormone **receptor** complex, and the barnase–barstar complex (10-12). The effect of a mutation on the stability of the complex can be evaluated as the ratio of the dissociation constants or, equivalently, as  $\Delta\Delta G_d$ , the change in free energy of dissociation from the wild type to the mutant:

$$\Delta\Delta G_d = \Delta G_d^{\text{wt}} - \Delta G_d^{\text{mut}} = RT \ln \frac{K_d^{\text{mut}}}{K_d^{\text{wt}}} \quad (6)$$

The accompanying change in the enthalpy and the entropy of dissociation can be measured by microcalorimetry. Substitution of an amino acid by [alanine](#) effectively deletes a side chain and all its interactions, hopefully without introducing major perturbations in the structure. It is often the best choice when testing side-chain contributions. Deleting a polar side chain that makes hydrogen bonds across the interface can have a dramatic effect, the  $K_d$  of the complex increasing by up to five orders of magnitude, equivalent to  $\Delta\Delta G_d = 7$  kcal/mol. This is more free energy than expected for breaking a hydrogen bond in aqueous solution, but the hydrogen bond partner of the mutated residue is left buried in an unfavorable environment at the interface, unless water molecules replace the missing

side chain. Thus, when deleting in partner A an [arginine](#) side chain that forms a salt bridge with an [aspartic acid](#) residue of partner B, the observed  $DDG_d$  value may represent the cost of dehydrating the carboxylate rather than the salt bridge energy. The latter may be recovered by also mutating the aspartate residue, measuring dissociation constants for all four complexes AB,  $A^{mut}B$ ,  $AB^{mut}$ , and  $A^{mut}B^{mut}$ , and calculating:

$$\Delta G_d^{inter} = \Delta\Delta G_d^{AmutB} + \Delta\Delta G_d^{ABmut} - \Delta\Delta G_d^{AmutBmut} \quad (7)$$

Through Equation (6), this relates to the ratio of the dissociation constants:

$$\Delta G_d^{inter} = RT \ln \frac{(K_d^{AmutB})(K_d^{ABmut})}{(K_d^{AB})(K_d^{AmutBmut})}$$

If there is no interaction between the two mutated residues, the effect of the two substitutions should be additive and  $DG_d^{inter} \approx 0$ . If they form a salt bridge, its energy is lost in all three mutant complexes. Equation (7) indicates that it contributes directly to  $DG_d^{inter}$ , whereas the cost of dehydrating the two side chains cancels out. Mutagenesis data collected in several systems reveal that a single pairwise interaction may account for as much as 6 kcal/mol of the free energy of dissociation. Residue pairs that form salt bridges and charged hydrogen bonds yield the largest values; pairs making neutral hydrogen bonds or nonpolar interactions rarely exceed 3 kcal/mol. It should be kept in mind that polar interactions also include a **van der Waals** component and that the mutagenesis method cannot assess the contribution of the main chain. In complexes of known three-dimensional structure, peptide groups are involved in at least half of the hydrogen bonds at protein–protein interfaces; side chain-to main chain bonds are very common, and main chain–main chain bonds also occur.

## 6. Complementarity and specificity

Complementarity is an essential feature of both stability and specificity. Figure 1 illustrates how it is expressed in both the shape of the surface and its chemical composition. In the lysozyme–HyHEL5 and barnase–barstar complexes, the interacting surfaces are bumpy, yet fairly flat in their general appearance. There is a global fit that exists prior to their association, supplemented with a knob-into-hole type of fit, which may result largely from side-chain adjustments. Two extreme models of shape complementarity are (1) the lock-and-key fit of two rigid pre-formed shapes, and (2) the hand-into-glove fit, where one partner adapts to the shape of the other. Whereas neither model is a satisfactory description of reality, the two systems of Figure 1 are closer to the first than the second. The overall shape complementarity is not only apparent visually but is also sufficient for a computer search to be able to select the right mode of association. This can be performed by **docking** lysozyme onto the combining site of the antibody, or barstar onto barnase. The search usually succeeds when docking is performed using protein structures from a dissociated complex or the free protein structures, when their conformations are not optimized for association (13). The chemical and electrostatic complementarities of the interacting surfaces are equally obvious in Figure 1. Barnase is positively charged around its active site, which binds **RNA** substrates, and barstar has a very negatively charged surface to interact with it. In lysozyme–HyHEL5, **electrostatic** effects are more subtle, but in both complexes buried polar groups on one surface find hydrogen-bonding partners on the other.

When interface residues are changed by mutation, the complementarity is perturbed. Deleting a large side chain creates a hole that water may or may not fill, depending on the size of the hole; removing (or adding) a polar group leaves a hydrogen bond donor or acceptor unsatisfied. This is achieved experimentally by [site-directed mutagenesis](#) or by using natural evolutionary variants, such as bird lysozymes that are closely related to hen lysozyme and bind to the same monoclonal antibody, but with different affinities. Specificity is often evaluated from the effect of such single-residue



substitutions on affinity. Because a mutation on partner A can, in principle, be compensated by modifying the surface of partner B, the double mutant complex  $A^{\text{mut}}B^{\text{mut}}$  may be more stable than either single mutant species  $A^{\text{mut}}B$  and  $AB^{\text{mut}}$ , demonstrating that a point mutation can alter specificity. However, point mutants are not appropriate tests of the specificity required for function *in vivo* or in immunochemical experiments for antibodies. In these cases, specificity is the capacity to identify a given molecular species among many others, not to discriminate between closely related molecules such as two bird lysozymes. When an immunological test is performed, the lysozyme [epitope](#) competes with others for interaction with the antibody-combining site. The competition is not with a mutant version of the epitope, but with features on the molecular surface of all the other chemical species present in the sample. How many show sufficient complementarity with the combining site to compete with the cognate epitope? The answer to this question obviously depends on the composition of the solution, and it determines the response to the test. To visualize nonspecific protein–protein interactions associated with random contacts in such mixtures, we may assume that they create transient interfaces comparable to crystal contacts. These interfaces are small, and very few of them have surface areas that can compare with the cognate interface. Because the number of polar interactions increases with the area of the interface, this remark also applies to hydrogen bonds and salt bridges. Thus, the presence on the surfaces of two proteins of regions covering a sufficient area (on the order of  $750 \text{ \AA}^2$  each), and having complementary shapes and chemical compositions, appears to be both necessary and sufficient for specific recognition (1).

## 7. Conclusion

A combination of structural, biochemical, and thermodynamic studies has increased our understanding of protein–protein interactions in recent years. Site-directed mutagenesis, guided by a known three-dimensional structure, has proved to be a powerful tool in elucidating the roles of individual amino acid residues and of pairwise interactions. In systems that assemble as approximately rigid bodies, general rules have been established. They may help in designing protein surfaces that show sufficient complementarity to form stable complexes. Much less can be said at present of those systems where major conformation changes occur or where, as in trypsinogen–BPTI and thrombin–hirudin, association is coupled with folding of the components. In these systems, the bimolecular step becomes part of the folding process, underlining the close relationship between the two and recalling that protein–protein interaction is as essential to structure as it is to function.

## Bibliography

1. J. Janin and C. Chothia (1990) The structure of protein–protein recognition sites, *J. Biol. Chem.* **265**, 16027–16030.
2. S. Jones and J. M. Thornton (1996) Principles of protein–protein interaction, *Proc. Natl. Acad. Sci. USA* **93**, 13–20.
3. S. Sheriff, E. W. Silvertown, E. A. Padlan, G. H. Cohen, S. J. Smith-Gill, B. C. Finzel, and D. R. Davies (1987). Three-dimensional structure of an antibody–antigen complex, *Proc. Natl. Acad. Sci. USA* **84**, 8075–8079
4. V. Guillet, A. Laphorn, R. W. Hartley, and Y. Mauguén (1993) Recognition between a bacterial ribonuclease, barnase, and its natural inhibitor, barstar, *Structure* **1**, 165–177.
5. M. Marquart, J. Walter, J. Deisenhofer, W. Bode, and R. Huber (1983) The geometry of the reactive site and of the peptide groups in trypsin, trypsinogen and its complexes with inhibitors. *Acta Crystallogr.* **B39**, 480–492.
6. J. Janin and F. Rodier (1995) Protein–protein interaction at crystal contacts, *Proteins* **23**, 580–587.
7. T. J. Rydel, A. Tulinsky, W. Bode, and R. Huber (1991) The refined structure of the hirudin–thrombin complex, *J. Mol. Biol.* **221**, 583–601.
8. D. Xu, C. J. Tsai, and R. Nussinov (1997) Hydrogen bonds and salt bridges across protein–protein interfaces, *Protein Eng.* **10**, 999–1012.

9. K. A. Hibbits, D. S. Gill, and R. C. Willson (1994) Isothermal titration calorimetric study of the association of hen egg lysozyme and the anti-lysozyme antibody Hy-HEL5. *Biochemistry* **33**, 3584–3590.
10. D. R. Davies and G. H. Cohen (1996) Interactions of protein antigens with antibodies, *Proc. Natl. Acad. Sci. USA* **93**, 7–12.
11. J. A. Wells (1996) Binding in the growth hormone receptor complex, *Proc. Natl. Acad. Sci. USA* **93**, 1–6
12. G. Schreiber and A. R. Fersht (1995) Energetics of protein–protein interactions: analysis of the barnase–barstar interface by single mutations and double mutant cycles, *J. Mol. Biol.* **248**, 478–486.
13. J. Janin (1996) Protein–protein recognition, *Prog. Biophys. Mol. Biol.* **64**, 145–166.

### Suggestions for Further Reading

14. C. Bränden and J. Tooze (1991) *Introduction to Protein Structure*, Garland Publishing, New York.
15. E. M. Phizicky and S. Fields (1995) Protein–protein interactions: methods for detection and analysis, *Microbiol. Rev.* **59**, 94–123.
16. J. Janin (1996) Protein–protein recognition, *Prog. Biophys. Mol. Biol.* **64**, 145–166.
17. S. Jones and J. M. Thornton (1995) *Prog. Biophys. Mol. Biol.* **63**, 131–165.

## Proteinase Inhibitors

Inhibitors of **proteolytic enzymes**, *proteinases*, encompass a wide variety of molecules, ranging from protein inhibitors (see [Proteinase Inhibitors](#), [Proteins](#)) to naturally occurring, low-molecular-weight inhibitors; and from broad-specificity, synthetic inhibitors to highly specific, low-molecular-weight inhibitors. Among the latter (Table 1), there are group-specific reagents that are effective against most of the members of each of the four major types of proteinase (**carboxyl**, **thiol**, [serine](#), and [metalloproteinase](#)), and these can be reversible or irreversible. Thus, [pepstatin](#) reversibly inhibits [pepsin](#), renin, [cathepsin D](#), and other carboxyl proteinases. Isolated from the culture filtrate of **Streptomyces** species, it is a [transition state analogue](#) that forms a tight binding complex with the enzyme. Highly specific inhibitors have been synthesized that are directed against renin, the human immunodeficiency virus ([HIV](#)) proteinase, and numerous other pathogen-related carboxyl proteinases. **Iodoacetate** and *p*-mercuribenzoate are irreversible inhibitors of thiol proteinases because they are general alkylating agents for [thiol groups](#). Metals such as Hg<sup>2+</sup> or Ag<sup>+</sup> are relatively nonspecific reagents for thiol groups and also inhibit thiol proteinases. Another more specific, irreversible inhibitor of this class of proteinases is E-64 [*trans*-epoxysuccinyl-L-leucylamido (4-guanidino) butane] isolated from culture filtrates of *Aspergillus japonicus*. While it covalently modifies thiol groups in thiol proteinases, it does not modify them in other enzymes.

**Table 1. Proteinase Inhibitors** <sup>a</sup>

| Type | Proteinases Inhibited |
|------|-----------------------|
|------|-----------------------|

|                          |              |                                  |
|--------------------------|--------------|----------------------------------|
| Amastatin                | Peptidic     | Metallo (aminopeptidases)        |
| Antipain                 | Peptidic     | Serine/thiol (trypsin-like)      |
| $\alpha_1$ -Antitrypsin  | Protein      | Serine                           |
| Aprotinin                | Protein      | Serine                           |
| Bestatin                 | Peptidic     | Metallo (aminopeptidases)        |
| Cystatin                 | Protein      | Thiol                            |
| Chymostatin              | Peptidic     | Serine/thiol (chymotrypsin-like) |
| 3,4-Dichloroisocoumarin  | Heterocyclic | Serine                           |
| E-64                     | Epoxide      | Thiol                            |
| Elastatinal              | Peptidic     | Serine (elastase-like)           |
| EDTA                     | Chelator     | Metallo                          |
| Leupeptin                | Peptidic     | Serine/thiol (trypsin-like)      |
| $\beta_2$ -Macroglobulin | Protein      | Many                             |
| Pepstatin A              | Peptidic     | Carboxyl (some)                  |
| PMSF                     | Sulfonyl     | Serine                           |
| Phosphoramidon           | Peptidic     | Metallo (some)                   |
| TLCK                     | Alkylating   | Serine (trypsin-like)            |
| TPCK                     | Alkylating   | Serine (chymotrypsin-like)       |
| Trypsin inhibitors       | Proteins     | Serine                           |

---

<sup>a</sup> This is a representative list of the many proteinase inhibitors that are known. The ones listed are active against particular classes of proteinases, although in several cases inhibitory activity is limited. There is a very large number of proteinase inhibitors that have been designed or found to be specific for just one proteinase such as angiotensin converting enzyme inhibitors or AIDS proteinase inhibitors. A listing of these is beyond the scope of this table. Those interested in the structures and/or mechanism of action of any of the above inhibitors are referred to specific listings.

There are many irreversible inhibitors of serine proteinases, some of which are described elsewhere in this volume, such as [DIFP](#) (diisopropylfluorophosphate), [PMSF](#) (phenylmethylsulfonyl fluoride), [TLCK](#), and [TPCK](#). A slowly reversible inhibitor of a wide range of serine proteinases is 3,4-*dichloroisocoumarin*. Some reversible inhibitors that inhibit serine as well as thiol proteinases are [leupeptin](#), *antipain*, and *chymostatin*, all of which are aldehyde derivatives of amino acids.

Metalloproteinases are inhibited by reagents that remove the active site metal ion, most notably [EDTA](#) (ethylenediaminetetraacetic acid) and 1,10-phenanthroline (orthophenanthroline). A naturally occurring fungal product, *phosphoramidon*, is a potent but reversible inhibitor of [thermolysin](#) and similar bacterial metalloendoproteinases. It has served as a model for the synthesis of many more-specific metalloproteinase inhibitors.

There are naturally occurring inhibitors, such as *amastatin* and *bestatin*, that are specific for [aminopeptidases](#). A particularly effective inhibitor of [carboxypeptidase A](#) is L-benzylsuccinate. Carboxypeptidase Y, which is useful in protein analysis, is inactivated by serine proteinase inhibitors and by thiol reagents.

#### Suggestion for Further Reading

A. J. Barrett and G. Salvesen (eds.) (1986) *Proteinase Inhibitors*, Elsevier, Amsterdam.

## Proteinase Inhibitors, Protein

Proteinase inhibitors, protein are proteins that form tight complexes with proteolytic enzymes. Such complexes are commonly totally devoid of enzymatic activity but some, especially those with macroglobulins [see [Macroglobulins](#)], retain activity against small but not large substrates. In view of the very large number of biological functions of proteinases, the spatial and temporal control of their activity is essential. Nature appears to employ two major mechanisms for this control. Many proteinases are synthesized as zymogens, where the removal of a propeptide is required for the expression of maximal activity. Once activated, proteinases are controlled by protein inhibitors. These are extremely abundant, as each organism contains many kinds of protein proteinase inhibitors, some of them at very high molar concentrations. The two mechanisms of control of proteolysis often, but not always, overlap. In many, but not all, zymogens, the propeptide is a protein inhibitor covalently tethered to the enzyme.

Although overall function of protein proteinase inhibitors—the control of unwanted proteolysis—is obvious for most inhibitors, their target enzymes are not. The presumed biological purpose of an inhibitor is to inhibit its target enzymes. *In vitro*, most inhibitors are capable of inhibiting a much greater number of enzymes, called the cognate enzymes. There seems to be little likelihood of the ovomucoid from the egg white of laughing kookaburra coming in contact with bovine chymotrypsin outside a biochemical laboratory. In some cases, the target enzymes seem obvious. The worm *Ascaris lumbricoides* contains large amounts of inhibitors of digestive enzymes—trypsin, chymotrypsin, elastase, and pepsin. As the worm spends a large fraction of its lifetime in the digestive tract of its mammalian host, it seems intuitive to presume that the host's digestive enzymes are the targets. Similarly, the saliva of a large number of blood-sucking animals contains inhibitors for thrombin or for enzymes preceding thrombin in the blood-clotting cascade. It seems natural to assume that the targets are the blood-clotting enzymes of the prey. Humans with genetic deficiency of a serpin—a<sub>1</sub> proteinase inhibitor are often subject to emphysema. Emphysema is a consequence of excessive proteolysis of the lungs by leukocytic enzymes. As these enzymes are inhibited by a<sub>1</sub> proteinase inhibitor, they are among the target enzymes for this inhibitor. In several other cases, genetic diseases provided a clue for the target.

For most protein proteinase inhibitors, however, the targets are unknown. One reason for this is that the targets are not limited to cognate, endogenous enzymes of the organism. As already shown, inhibitors are often involved in interorganismic interaction. In particular, defensive function is often postulated for inhibitors, and the list of potential predators—vertebrates, insects, and bacteria—is huge. Historically, the first protein proteinase inhibitors known were trypsin inhibitors. A custom developed of screening numerous plant, animal, and bacterial tissues with bovine cationic trypsin. A very large number of proteins that inhibit bovine trypsin were found and frequently named as tissue X trypsin inhibitors without regard to the likelihood that any trypsin in general or bovine trypsin in particular was the target. The a<sub>1</sub> proteinase inhibitor discussed above was first isolated as a bovine trypsin inhibitor and named a<sub>1</sub> antitrypsin. The name diverted many researchers into the detailed study of bovine trypsin—a<sub>1</sub> proteinase inhibitor interaction, even though later it was shown that trypsin is not among the targets and is not an especially efficient cognate enzyme. Similar histories occurred for many other inhibitors, and it must be expected that many of the substances called trypsin inhibitors do not have trypsin as their primary target.

Some protein proteinase inhibitors appear to have been recruited to perform functions other than inhibition of proteinases. Such is the case with the serpin, ovalbumin, and with several members of

the Kunitz bovine pancreatic trypsin inhibitor family, which function as dendrotoxins and as  $\text{Ca}^{++}$  channel blockers. In such proteins, the inhibitory activity has either disappeared, as in ovalbumin, or remain in vestigial form. An alternate possibility is that inhibitory activity of some proteins is accidental. The inhibitory activity arose without any biological pressure to inhibit. The exquisitely good fit between enzymes and inhibitors in many enzyme-inhibitor complexes appears to belie such an accidental origin.

Protein proteinase inhibitors are a huge and highly diverse class of proteins. However, a classification based on the mechanistic class of proteinases they inhibit proved moderately successful. They can be conveniently divided into the following classes. [See separate entries for each class.]

[Macroglobulins](#)

[Serine proteinase inhibitors](#)

Cysteine proteinase inhibitors

Aspartyl proteinase inhibitors

[Metalloproteinase inhibitors.](#)

Macroglobulins trap proteinases belonging to all four mechanistic classes. The remaining four classes inhibit only some enzymes belonging to its mechanistic class. This classification lasted for 20 years. Some exceptions have been found. For example, the viral serpin CrmA, a serine proteinase inhibitor, inhibits the interleukin-converting enzyme, ICE, which is a cysteine proteinase.

In some proteinase inhibitors–proteinase associations, it is impossible or difficult to define an equilibrium constant for the interaction  $K_a$  or its inverse, the inhibition constant  $K_I$ . Such is the case for complexes with human  $\alpha_2$  macroglobulin, where the presence of covalent isopeptide bonds between the trapped enzyme and the inhibitor leads to effective irreversibility. Similarly, in many enzyme–serpin interactions [see [Serpins](#)], the main route of dissociation is to the enzyme and the reactive site hydrolyzed inhibitor. In contrast to standard-mechanism protein inhibitors of serine proteinases, such a reactive site hydrolyzed inhibitor is no longer effective. Therefore, instead of inappropriate  $K_a$  or  $K_I$  values, such systems are characterized by the second-order association rate constant  $k_{\text{on}}$ . The reported  $k_{\text{on}}$  values range from  $10^3\text{M}^{-1}\text{s}^{-1}$  to somewhat above  $10^7\text{M}^{-1}\text{s}^{-1}$  at optimal pH value. The same range is also reported for inhibitors for which enzyme-inhibitor equilibrium is established. The typical  $k_{\text{on}}$  value is  $10^6\text{M}^{-1}\text{s}^{-1}$ .

Most inhibitors are characterized by enzyme inhibitor association equilibrium constants  $K_a$  or by their inverses, the inhibition constants  $K_I$ . The reported values for  $K_I$  range from  $10^{-3}$  to  $10^{-14}\text{M}$ ; these range limits are not a function of the system but rather of the ability of available measurement techniques. It is quite likely that both much weaker and much stronger enzyme-inhibitor pairs will be found as the range of measuring techniques widens. The association equilibrium constant  $K_a$  is given by

$$K_a = \frac{1}{K_I} = \frac{k_{\text{on}}}{k_{\text{off}}} \quad (1)$$

where  $k_{\text{off}}$  is the first-order dissociation rate constant for the enzyme inhibitor complex. As the  $k_{\text{on}}$  value is typically  $10^6\text{M}^{-1}\text{s}^{-1}$  and varies relatively little, it is seen that, for  $K_I$  greater than  $10^{-7}$  to  $10^{-8}\text{M}$ ,  $k_{\text{off}}$  is quite large and the dissociation half-life is fast on the laboratory time scale. At  $K_I$ ,  $10^{-8}\text{M}$  half-life is 70 s, and this rises to  $7 \times 10^7\text{s}$  (2 yr) for  $K_I$  of  $10^{-14}\text{M}$ . These very slow dissociation half-lives tend to mislead many investigators who do not expect them. A great majority of protein

inhibitors of proteinases are strictly competitive as binding of the inhibitor to the active site completely prevents the binding of substrate to the active site. When concentrated substrate is added to an enzyme inhibitor complex, however, it competes for the enzyme so slowly that the phenomenon is often missed. Thus, many protein inhibitors are incorrectly labeled *noncompetitive*.

Two techniques of measuring  $K_a$  ( $1/K_I$ ) are widely employed. In one,  $k_{on}$  and  $k_{off}$  are both determined, and the equilibrium constant is derived from the ratio. In the other, appropriate known concentrations of enzyme and inhibitor are preincubated together until they achieve equilibrium. Then the concentration of the components is determined. Both methods generally employ the determination of free enzyme by detection of the hydrolysis of highly sensitive synthetic substrates. Such substrates are well worked out for many proteinases.

Because we are still discussing all the protein proteinase inhibitors together, it is relatively hard to make strong and sweeping generalizations about their behavior. Some ( $\alpha_2$  macroglobulins, serpins) are moderately large and rather unstable proteins. Others are very small (marinostatins have only a dozen amino acid residues, whereas squash family inhibitors have about 30 residues). Some are exceptionally stable to heat, denaturing acids, and proteolysis. Some of these—such as bovine pancreatic trypsin inhibitor (Kunitz), Bowman-Birk family inhibitors, and avian ovomucoids—are isolated by denaturing most proteins in the sample so that only the native inhibitor remains. It is quite common to encounter single polypeptide chains that consist of more than one inhibitory domain. The phenomenon is well exemplified by the Kazal family of protein inhibitors of serine proteinases. Pancreatic secretory trypsin inhibitors and acrosin inhibitors in sperm consist of but a single domain. Submandibular inhibitors in carnivore mammals have two tandem domains on the same polypeptide chain. Avian ovomucoids have three tandem domains, and alligator ovomucoids have four. Ovoinhibitor, a protein present both in avian egg white and in avian blood, has seven such domains. There are 10 Kazal domains in addition to many other domains in agrin. Similar multiple tandem domains are present in the bovine pancreatic trypsin inhibitor (Kunitz). In that family, many such inhibitors are called bikunins, trikunins, and so forth. Most inhibitors in the plant inhibitor Bowman-Birk family consist of two homology regions each with its own reactive site. In the potato II family, multidomain and tandem domains are common.

Although serine proteinase inhibitors are the most studied, the presence of multiple tandem domains on the same polypeptide chain is not limited to them. High molecular weight kininogens, histidine-rich glycoproteins and fetuins consist of several tandem repeats of cystatin—a widely known cysteine proteinase inhibitor. An 85-kDa inhibitor with eight cysteine proteinase-inhibiting domains from potatoes has been described.

Multiple domains, each of which can combine with an enzyme, sometimes arise from strong and specific noncovalent association. All the members of *Streptomyces* subtilisin inhibitor family, SSI, are homodimers of this type. They form ternary complexes with enzymes. Similarly, a serine proteinase inhibitor from *E. coli*—ecotin—is a homodimer. A Kazal inhibitor from red sea turtle egg whites, testudin, is a disulfide bridged heterodimer of two Kazal inhibitory domains. Finally, human  $\alpha_2$  macroglobulin is a noncovalent homodimer of two disulfide bridged homodimers. The biological reasons for multiple inhibitory domains in the same molecule vary. Ovoinhibitor is present in the blood of birds. Its seven domains and glycosylation make it just large enough to avoid renal filtration. The SSI dimer strongly contributes to the stabilization of this otherwise fragile protein. The multidomain ornamental tobacco potato II family inhibitor appears to be an in vivo precursor of separate inhibitory domains.

#### Suggestions for Further Reading

W. Bode and R. Huber (1992) Natural protein proteinase inhibitors and their interaction with proteinases. *Eur. J. Biochem.* **204**, 433–451.

M. Laskowski Jr. and I. Kato (1980) Protein inhibitors of proteinases. *Ann. Rev. Biochem.* **49**,

## Proteinase K

*Proteinase K* is a [serine proteinase](#) of the [subtilisin](#) family obtained from *Tritirachium album*. It catalyzes the hydrolysis of [peptide bonds](#) involving the carbonyl group of **hydrophobic** aliphatic and aromatic amino acid residues (1). It has quite broad specificity and is capable of degrading native proteins (2). As a consequence, it is often added to cell or nuclear extracts to inactivate **nucleases** during the isolation of DNA or RNA. It can also be used to release **nucleic acids** from protein–nucleic acid complexes and to modify proteins and **glycoproteins** on cell surfaces. It can also be used to examine the topological orientation of membrane-associated proteins (3). It is stable in solution at 4°C, pH 8.0 in the presence of 1 mM calcium for several months. Serine proteinase inhibitors such as **DIPF** or [PMSF](#) will inactivate it, but [TLCK](#) and [TPCK](#) do not.

### Bibliography

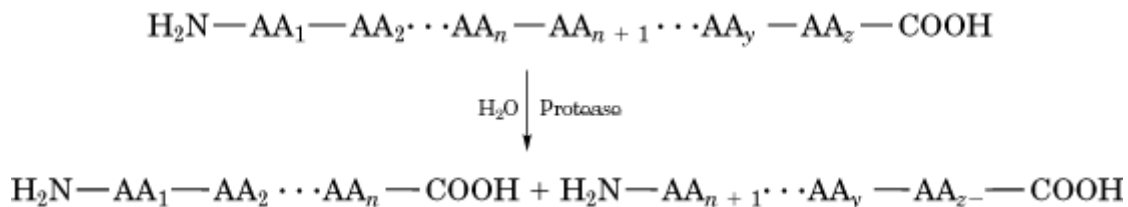
1. J. K. Dattagupta et al. (1975) *J. Mol. Biol.* **97**, 267–271.
2. H. G. Leberherz et al. (1986) *Biochem. J.* **233**, 51–56.
3. J. M. Pratt (1993) In *Proteolytic Enzymes: a Practical Approach* (R. J. Beynon and J. S. Bond, eds.), IRL Press, Oxford, U.K., pp. 181–191.

## Proteinases

Proteinases are [enzymes](#) whose substrates are [proteins](#) and whose catalytic mechanism involves the hydrolysis of one or more of the [peptide bonds](#) that constitute the **polypeptide** backbone of the protein molecule (Fig. 1). Proteins are the functional molecules of cells. Some are the catalysts that facilitate essentially every chemical reaction that occurs in biology; others are inhibitors, activators, or modulators of these reactions; still others have a structural function; and some participate in various processes of cellular communication. Indeed every conceivable biological function relies directly or indirectly on one or more proteins. Some of these functions are required for very short periods of time (minutes to hours), whereas others are more persistent (days to weeks). The cell has developed elaborate schemes based on gene regulation (see **Gene expression**) to ensure that proteins are synthesized and are available when needed. Similarly, they have developed mechanisms to eliminate proteins when they are no longer needed and to recycle their constituent amino acids. These mechanisms involve hydrolysis of the peptide bonds that link the amino acids into polymeric chains, a process (known as proteolysis) that is catalyzed by enzymes known as proteinases (also referred to as proteases, proteolytic enzymes, peptide hydrolases, or [peptidases](#)). As one might imagine, the realm of proteinases is about as extensive as that of proteins in general; and because of the catastrophic havoc they could create if allowed to run rampant, they must be subject to tight control. Consequently, an intricate system of checks and balances has evolved to limit the location, duration, and specific target of proteinase activity. The Enzyme Commission of the International Union of Biochemistry and Molecular Biology has classified proteinases according to their

numerical system as E.C. 3.4.

**Figure 1.** Schematic depiction of the (endo)proteinase-catalyzed hydrolysis of a single peptide bond in a protein [here represented as a linear polymeric chain of  $z$  amino acids (AA)]. For simplicity, the amino ( $\text{H}_2\text{N}$ -) and carboxyl (-COOH) ends of the protein and its cleavage products are indicated in the nonionized state.



Four general classes of proteinases have been defined based on the particular hydrolytic mechanism that is used for peptide bond cleavage.

1. The [carboxyl proteinases](#) (E.C. 3.4.23) employ a pair of [carboxyl groups](#), the side chains of two [aspartic acid](#) residues, to activate a [water](#) molecule so that it can lyse a peptide bond in a protein substrate.
2. The [serine proteinases](#) (E.C. 3.4.21) are characterized by a unique [catalytic triad](#) of amino acid side chains: an aspartyl carboxyl group, a histidyl imidazole group, and a seryl hydroxyl group.
3. The third class of proteinases has a cysteinyl [thiol group](#) that is required for activity, and hence these enzymes are called [thiol proteinases](#) (E.C. 3.4.22).
4. The fourth class of proteinases is unique in that it has a metal ion, zinc, as part of the active site (see [Metalloproteinases](#)) (E.C. 3.4.24).

Proteinases are also classified on the basis of where they act on a protein chain. Those that remove amino acids sequentially from the end of the chain are called *exoproteinases*. If they remove amino acids from the amino terminus, they are [aminopeptidases](#) (E.C. 3.4.11); and if they act at the carboxyl terminus, they are [carboxypeptidases](#) (E.C. 3.4.16–18). Proteinases that cleave bonds in the middle of the protein chain are *endoproteinases*.

Proteinases often act collectively, as with the digestive proteinases of the pancreas, the compartmentalized proteinases of the **lysosome**, the cascading proteinases of [blood clotting](#), or the sequentially acting proteinases of the renin–angiotensin–**kinin** system. A rather specialized group of proteinases constitutes what is known as the [proteasome](#), a macromolecular assembly of numerous different proteinases that degrades cytosolic proteins that have been specifically marked for that purpose by covalent attachment of a small protein called [ubiquitin](#) (see **Protein degradation**). The proteasome also plays a role in processing antigens as part of the immune response and also functions critically in the control of the cell cycle.

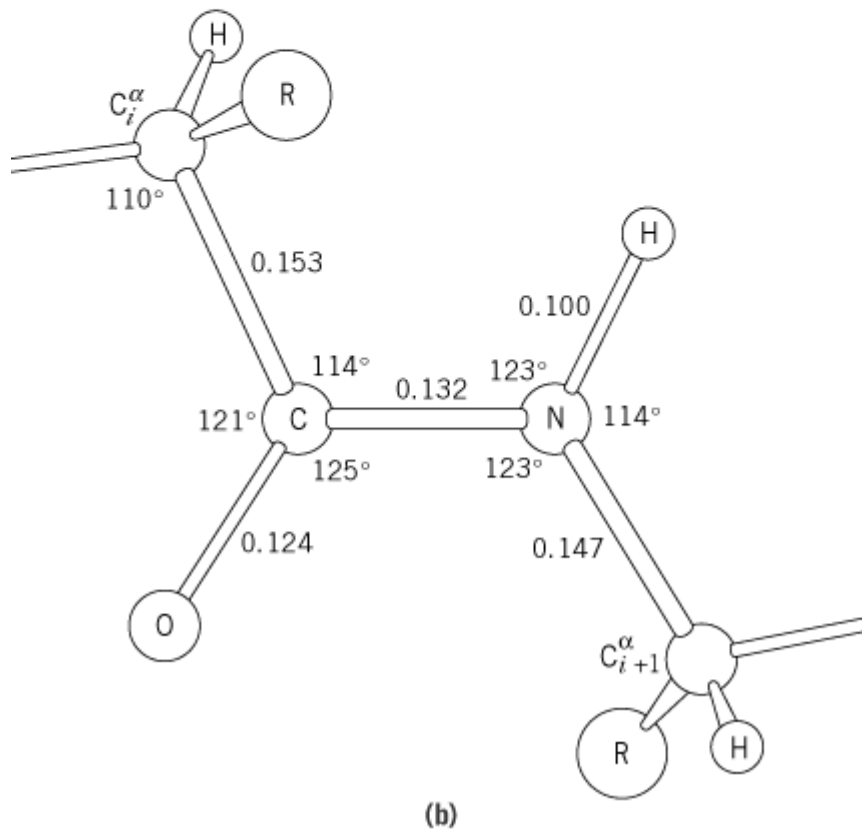
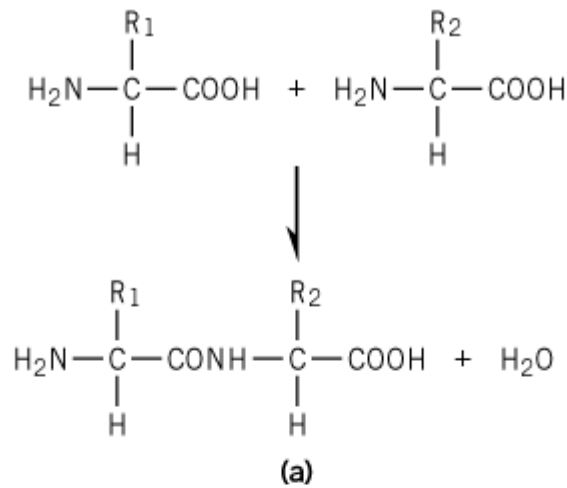
## Proteins

Proteins are **macromolecules** that support both the structural and functional aspects of life, along with other molecules such as nucleic acids, polysaccharides, [lipids](#), vitamins, etc. The name derives



from the Greek “protos,” meaning “first,” as suggested by Berzelius in 1838. Proteins are most vividly recognized for their functions; eg, [enzymes](#) as **catalysts**, muscle and **flagella** as motor machines, [growth factors](#) and [hormones](#) as regulatory elements. In contrast, bones, tendons, nails, and hair are notable manifestations of structural proteins. Proteins are linear copolymers of 20 rather limited kinds of L-**α-amino acids** (plus [selenocysteine](#)), which are linked through condensation polymerization, with the loss of one H<sub>2</sub>O unit from the **α-carboxyl group** of one amino acid and the **α-amino group** of another (see Fig. 1). The polypeptide **backbone** of a protein is thus a long repeat of amino acid residues in the form of —NH—CH(R<sub>i</sub>)—CO—, R<sub>i</sub> being one of 20 kinds of [side chains](#). The group of atoms between two amino acid residues, —CO—NH—, contains the [peptide bond](#). Hence, proteins are collectively [polypeptide chains](#) of biological origins. The N—C peptide bond has a partial double-bond character, through resonance with the carbonyl double bond, that prevents its free rotation. Thus, the geometry of a peptide bond is virtually fixed in a planar conformation, usually the *trans* isomer (see [Cis/Trans Isomerization](#)). This is a very important characteristic when one considers the conformational possibilities of proteins (1). One end of a linear polypeptide bearing H<sub>2</sub>N is called the amino- or **N-terminus**, and the other with —COOH is called the carboxyl- or **C-terminus** (2).

**Figure 1.** (a) Structure of L-α-amino acids and polymerization mechanism by dehydration between the amino and carboxyl groups. (b) Structure of the peptide bond.



Each polypeptide chain is therefore distinguished by the number of different kinds of amino acid residues it contains (its amino acid composition) and the specific order of such residues, aligned from the N- to C-terminus (its amino acid sequence, or the [primary structure](#)). The chemical nature of the side chains and their distribution along the main chain determine the **conformational** state of a polypeptide chain under any particular conditions. The side chains of the 20 amino acids are therefore classified according to their most notable chemical properties, such as (1) ionizable ([hydrophilic](#)), (2) [polar](#) but not ionizable (hydrophilic), (3) aliphatic (**hydrophobic**), (4) aromatic (hydrophobic), (5) crosslinker, and (6) others, as in Table 1. (More details are found under [Amino acids](#) and entries on the individual residues.)

**Table 1. Classification of the 20 Natural Amino acids That Constitute Proteins**

---

| Ionizable              | Asp, Glu, Lys, Arg, His, Tyr, Cys       |
|------------------------|---|
| Polar but nonionizable | Ser, Thr, Asn, Gln                      |
| Aliphatic-containing   | Ala, Val, Leu, Ile, Pro (ring), Met (S) |
| Aromatic               | Phe, Tyr, Trp                           |
| Crosslinker            | Cys                                     |
| Others                 | Gly                                     |

---

Unlike synthetic **polymers**, the arrangement of amino acid residues in the primary structure of a protein is not random, but is precisely determined by the genetic information stored in the chromosomal DNA or RNA. In humans, there are more than 100,000 proteins with various activities. Although a human has the capacity to synthesize the majority of the amino acids from carbohydrates and other nutrients, proteins are formed in most part from recycled amino acids, both from the daily intake of food stuff and from the degradation products of proteins in the body. So long as they recover their original chemical structures, recycled amino acids from any protein may be used for the resynthesis of any other protein by linking them in different orders according to the genetic instruction. One of the most unifying views of biochemistry is that the fundamental metabolism is common among all the organisms, eg, the energy-yielding metabolic cycle called the glycolysis pathway is one of the bases of life, for *Escherichia coli* as well as humans. Enzymes participating in the corresponding steps in the glycolysis pathway in various organisms have virtually the same functions and the same name, but when looked at more closely, their amino acid sequences differ from species to species, thus defining each enzyme as species-specific. The same muscle protein we take as beef must, therefore, be hydrolyzed (digested) to individual amino acids in the digestive organs and transported via the bloodstream to the liver and other tissues to be repolymerized into muscle and other proteins, all of which are specific to the human body. Hydrolysis of proteins into amino acids is effected by a class of proteins known as **proteolytic** enzymes, or **proteinases**, such as **pepsin** in the stomach, and **trypsin** and **chymotrypsin** in the small intestine. Proteins are also constantly synthesized and degraded in the cell for the regulation of cellular activities and for recycling amino acids. Intracellular protein degradation is carried out in acidic vesicles called lysosomes. They contain a group of enzymes called **cathepsins** for the hydrolysis of proteins, along with many other enzymes for the degradation of nucleic acids, polysaccharides, glycolipids, etc. Both in the stomach and lysosomes, an acidic condition is effective to loosen the three-dimensional structure of proteins and to expose their peptide bonds for hydrolysis (2).

Proteins are synthesized inside living cells on a supra-molecular machinery called the **ribosome**. After receiving genetic information prescribing an amino acid sequence from the nucleus in the form of **messenger RNA** (mRNA), ribosomes interpret the base sequence of mRNA into the expected amino acid sequence with the help of **transfer RNAs** (tRNA) and enzymes called **aminoacyl tRNA synthetases**. Once each tRNA carrying a specific **anticodon** is conjugated on the enzyme with a corresponding amino acid, it will bind to the mRNA on ribosomes. As tRNAs are aligned along the mRNA, amino acids at the other end of the tRNA are polymerized into a nascent polypeptide chain.

The polypeptide chains formed in this way are not biologically active without going through the processes called **post-translational modification** and **protein folding**. These are, respectively, the covalent modifications and the formation of a specific three-dimensional structure that are needed for their full biological activities. Polypeptide chains with unrelated primary structures usually form different three-dimensional structures having different activities. Thus, one **gene** dictates one particular amino acid sequence and a corresponding activity through the unique three-dimensional structure of the protein. During the folding process, **disulfide bonds** between pairs of **cysteine**

**residues** may be generated. Post-translational modifications include, among others, proteolytic cleavage of peptide bonds, [O-glycosylation](#) at specific [serine](#) and [threonine](#) residues, and **N-glycosylation** at [asparagine](#), [hydroxylation](#) of [proline](#) and/or [lysine](#) residues, **methylation** of Lys residues, and covalent linkage of fatty acids (see [Membrane Anchors](#)).

**Prosthetic groups** can be added covalently or noncovalently. Post-translational modifications cannot yet be predicted always from the sequences of genes, so the chemical analysis of proteins must be conducted carefully to reveal every possible modification not directly prescribed by structural genes. After going through folding and post-translational modifications, a protein assumes its “native” conformation, complete with its biological activity. The three-dimensional structure of the native conformation is called the [tertiary structure](#), which may be dissected into several **secondary structures** known as the [alpha-helix](#), [beta-sheets](#), [turns](#), and various irregular structures. Secondary structures are generally those favored by the local sequence of residues, whereas the tertiary structure is built up by the long-range interactions between residues distant in the primary structure.

### 1. Molecular Weight

The molecular weights of protein range from relatively small values, such as 5,807 for [insulin](#) (51 amino acid residues) or 6,380 for rubredoxin, up to many millions. If we limit our scope to a single polypeptide chain, rather than including multichain complexes, the largest thus far known is the muscle titin (connectin), which acts as a structural support of myofibrils in muscle, with  $M_w = 3.5 \times 10^6$ . Many polypeptide chains smaller than insulin are found as enzyme inhibitors, hormones or growth factors, [toxins](#), etc. and are usually referred to as **oligopeptides** or polypeptides, rather than proteins.

### 2. Subunit or Quaternary Structure

Protein molecules are often composed of multiple polypeptide chains, termed subunits, within the arrangement called the [quaternary structure](#). Subunits may be identical (homo) or different (hetero) from each other, and they are held together primarily by noncovalent forces, such as hydrophobic and/or [electrostatic interactions](#) and [hydrogen bonds](#); sometimes, disulfide bonds are also found either within or between subunits. Two (dimer), four (tetramer), and even-numbered subunit structures predominate, whereas three (trimer) or odd-numbered ones are less common (see [Quaternary Structure](#)).

### 3. Origin and Evolution of Proteins

The **evolutionary** origin of proteins is not yet known, but it is certain that life on the Earth as we know it is not possible without proteins. DNA carries the genetic information, but proteins implement the machinery of life. Nevertheless, RNA molecules may have initially played the roles of proteins in an early stage of evolution, known as the [RNA world](#). Yet proteins are more efficient than RNA molecules, so they gradually replaced them, except in ribosomes, tRNA, and ribozymes. The most important events in protein evolution since then are considered to be the enlargement of protein structures through [gene duplication](#) phenomena, where the structural genes for proteins were repeatedly duplicated, producing larger and larger proteins. During biological evolution, it was imperative that at least one [active site](#) in an enzyme must retain its activity to keep the life process going. Gene duplication is a convenient strategy to satisfy that condition, while allowing the rest of the multiple active sites free to change their functions from generation to generation. Many mutations can, therefore, be accumulated in a protein that went through gene duplication events without totally abolishing the biological activity of the resulting protein, and that protein can be reborn one day as a remodeled protein with a newly acquired function. Today, it is possible to trace the **phylogenetic** relations between proteins based mainly on a statistical comparison of their primary structures and to estimate the structure of their common ancestor and how many years have

elapsed since they branched off. Such [divergent evolution](#) is the most frequent, but occasionally there are pairs of proteins with similar activities, but without an obvious sequence or structural similarity except in the active site regions. These cases are referred to as [convergent evolution](#), meaning that the two proteins have come a long way during their evolution from different ancestors, to express similar activities from necessity.

#### 4. Methods in Protein Chemistry

Proteins are biological macromolecules, so they are usually discovered and purified on the basis of a specified activity. Occasionally, a protein is purified without knowledge of its activity or, most recently, amino acid sequences of many putative proteins are deduced solely from the **nucleotide sequences** of their genes, using the [genetic code](#). In such cases, researchers are forced to speculate on the biological activities of unknown proteins by analogy.

##### 4.1. Purification

Purification of proteins from various biological sources has been a major concern of protein biochemistry since the early days of its history. Enzymes are purified using (1) **chromatographic methods (size exclusion, affinity, and ion exchange)**, (2) [electrophoresis](#), (3) differential and [density gradient centrifugation](#), (4) **salting out** methods using, eg, ammonium **sulfate** as precipitant, (5) macromolecular phase separation using, eg, polyethylene glycol, etc. In purification procedures, the activity of the targeted sample is used to follow the direction of fractionation at each step. Because a finite amount of actual protein must be recovered, it has been extremely difficult to purify proteins that exist in very small amounts. The recent development of molecular biology [cloning](#) techniques has, however, changed protein purification quite dramatically by enabling the researchers to produce rare proteins of, eg, human origin easily and in abundance from *E. coli* using [recombinant DNA](#) technology, provided the degree of post-translational modifications is not extensive. For example, glycosylation cannot be achieved unless eukaryotic cells are used instead of *E. coli* for the expression of recombinant DNA (3) (see [Protein Engineering](#)).

##### 4.2. Chemical Structure

The chemical structure of a protein is studied, first, from the quantitative analysis of its amino acid composition, by [amino acid analysis](#). Classically, proteins are hydrolyzed into their constituent amino acids by heating them under evacuated conditions at 110°C in 6 M HCl for 24 hours. To protect [tryptophan](#) from oxidative degradation, a small amount of thioglycol is added. Separation and quantification of each amino acid are completed with an automated ion-exchange chromatography system known as an automatic amino acid analyzer. Cysteine residues in both the thiol and disulfide forms are oxidized to cysteic acid and cannot be differentiated, and [glutamine](#) and [asparagine](#) are converted to [glutamic acid](#) and [aspartic acid](#), respectively, and cannot be distinguished from their acid counterparts. Despite such ambiguities, a quantitative amino acid analysis must be done before the next step of the determination of amino acid sequences.

The primary structure, or the amino acid sequence, of a protein is usually elucidated by a chemical procedure known as the [Edman Degradation](#), which uses PITH (phenylisothiohydantoin) as the starting reagent to modify the N-terminal amino acid residue. After release and identification of the N-terminal residue as a PTH (phenylthiohydantoin) amino acid, the process is repeated on the shortened polypeptide chain, to identify the next amino acid residue, and so forth. This is still the most efficient and reliable way to determine the primary structure (see [Protein Sequencing](#)), provided the amino terminal is not blocked, eg, by acetylation or formylation.

##### 4.3. Mass Spectrometry

Proteins are not normally volatile, but [mass spectrometry](#) has become a powerful tool in the sequence determination of peptides and of the oligo-saccharide structures of **glycoproteins**. The extraordinary sensitivity of modern mass spectrometry enables researchers to distinguish a difference of one atomic mass unit from a total of 10,000 units.

#### 4.4. X-ray Crystallography

The three-dimensional structures of proteins are usually determined by [X-ray crystallography](#). The first protein to be crystallized was [albumin](#), the first enzyme was urease, and the first structure to be determined was that of [myoglobin](#). The procedure of structure determination is, in brief, as follows. First, the protein is crystallized out of aqueous solution by, eg, adjusting the pH closer to the [isoelectric point](#) (pI) of the protein, where its solubility is the lowest, or by changing the ionic strength and/or the activity of [water](#) by adding polymeric substances such as polyethylene glycol. Once a single crystal is obtained, it is usually irradiated with a monochromatic X-ray beam, and as many so-called “reflections” as possible are collected. This defines the crystal parameters. Each reflection is the Bragg reflection from a particular crystal plane with three integer indexes; when enough of them are collected, and their phases estimated (see [Phase Problem](#)), the reflections can be recombined to reconstruct inside a unit cell of the crystal the distribution of electron density and, therefore, that of the atoms. A unit cell contains at least one molecule of the protein, so the result will give the 3D distribution of all the protein atoms, except for the hydrogens, and consequently the structure of a protein molecule.

**Neutron diffraction** can also be used to determine the structures of proteins, but it is usually more consuming of time and money and is used primarily for special cases where the location of the hydrogen atoms is important. [Electron crystallography](#) can be used to determine the structure of membrane proteins in two-dimensional crystalline arrays. **Single-particle reconstruction** of individual protein molecules with different orientations from [electron microscopy](#) has also been used to determine the structures of noncrystalline proteins, especially those of tubular or fibrous bacteriophages and of muscle origins.

#### 4.5. Nuclear Magnetic Resonance (NMR)

[NMR](#) has increasingly come to be used to determine the structure of small proteins in solution. It has the advantage over X-ray crystallography that proteins do not have to be crystallized. However, the accuracy of the resulting structure is somewhat less than those obtained from the X-ray diffraction method, and only smallish proteins can be studied; the upper limit in molecular weight is of the order of 50,000 at the present time.

#### 4.6. Chemical Modification

The reaction of proteins with chemical reagents has been the classical method to study the contribution of particular side chains to their activities, or to change their physical and chemical properties. For example, polyethylene glycol groups can be linked to [lysine](#) side chains of proteins, to alter its solubility properties, particularly in organic solvents, change the immunological properties, in some cases to increase the stability, and so forth. The  $\epsilon$ - [amino group](#) of Lys residues is often targeted as the site of chemical modification, using reagents such as [anhydrides](#) and succinimide, whereas the [thiol groups](#) of cysteine residues react readily with many reagents, such as **iodoacetamide**. Chemical modification has, however, been largely superseded of late by [site-directed mutagenesis](#) and [protein engineering](#).

### 5. Other Methods Frequently used in Protein Chemistry

#### 5.1. Circular Dichroism (CD)

CD is conveniently used to estimate the **secondary structure** content of a protein, by comparing the CD spectrum of a protein with those of the [alpha-helix](#), [beta-sheet](#), and [turns](#). It is based on the chiral interaction of circularly polarized light with optically active molecules. A helical structure itself is chiral and interacts with left- and right-polarized light in different ways, as do nonhelical structures. A  $\beta$ -sheet, being another type of helix, is also chiral and shows a different spectrum from that of the  $\alpha$ -helix. Consequently, the CD spectrum of a protein can be used to estimate the percentage of its residues in helical,  $\beta$ -sheet, and turn secondary structures.

#### 5.2. Electron Microscopy (EM)

Large protein molecules can be observed by **EM**, either after staining them with a chemical

containing heavy atoms such as tungsten or uranium, or unstained by [cryoelectron microscopy](#). The staining method called “negative staining” is frequently used to visualize proteins “negatively” as empty regions after staining with the heavy-metal reagent. Three-dimensional reconstruction of protein images taken from different angles is increasingly used for the structural determination of noncrystallizable proteins. The arrangement of subunits within an oligomeric protein molecule is often deduced from electron micrographs of negatively stained protein molecules. Observation is made in a high vacuum, which is far from the native condition for biomolecules, and the stain can perturb the structure. The recent inventions of cryoelectron microscopy and [atomic force microscopy](#) (AFM) have greatly increased the possibility to visualize individual protein molecules under near-native conditions.

### 5.3. Chromatography

Most classical types of chromatography are now available as **HPLC**, for improved performance when working with small samples. Prepacked columns with uniform-sized gel particles having a small size and high mechanical stability guarantee a superior performance over manually operated systems. Ion-exchange, gel permeation, and affinity chromatography are most frequently used for the purification of proteins. [Reversed phase chromatography](#) is particularly suitable for separating oligopeptides.

### 5.4. Polyacrylamide Gel Electrophoresis (PAGE)

PAGE is a mainstay in protein chemistry, with a variety of modifications suited to different purposes. Among them, [disc electrophoresis](#) and [SDS-PAGE](#) under reducing conditions are the two most frequently used variations, with great resolving power for the former, and as a convenient method for the estimation of the molecular weight for the latter. Electrophoretic mobility measurements made at several pH values provide the isoelectric points of proteins, as does [isoelectric focusing](#). [Capillary zone electrophoresis](#) is a new development for all types of electrophoresis for analytical purposes, using very small amounts of sample.

### 5.5. Centrifugation

It was by ultracentrifugation that proteins were first shown to be homogeneous macromolecules, each of which had a definite molecular weight. Today, ultracentrifuges are widely used for preparative purposes to separate proteins and other biological macromolecules on the basis of their [sedimentation coefficients](#). The analytical ultracentrifuge is specifically built for the determination of the molecular weights and [hydrodynamic volumes](#) of proteins by running a **sedimentation equilibrium** and **sedimentation velocity** experiments. The sedimentation equilibrium method was one of a few last resorts for determining the absolute molecular weight, until the development of [mass spectrometry](#) techniques. Other methods of determination of protein molecular weights include, in the order of specialization, [light scattering](#), X-ray scattering, and neutron scattering.

### 5.6. UV and Visible Spectroscopy

The **absorbance spectroscopy** of proteins detects absorption bands in the visible and UV regions. The characteristic absorption of many proteins around 280 nm is due to the presence of Trp and Tyr residues. It is used to determine the concentration of a protein in solution, if its extinction coefficient is known. Proteins with a low content of Trp and/or Tyr (eg, [collagen](#)) have a very low absorption at 280 nm. Proteins with chromophores as prosthetic groups have visible colors. The red color of [hemoglobin](#) or [cytochrome](#) solutions is a good example. Some of the copper-containing proteins have an intense blue color. Flavin-containing proteins are often yellow in solution.

The backbone structure of a protein gives characteristic infrared and Raman spectra, due to its peptide groups (see [Vibrational Spectroscopy](#)). The absorption bands called amide I and II appear at specific wavelengths for  $\alpha$ -helix and  $\beta$ -sheet structures. The relative intensities of these bands can be used to estimate the secondary structure of a protein.

### 5.7. Fluorescence Spectroscopy

**Fluorescence** is a sensitive and versatile technique for the study of various conformational aspects of

proteins. When irradiated with UV light, many proteins give a fluorescence emission spectrum in the wavelength region of 300 to 400 nm. The wavelength of maximum emission ( $\lambda_{em}$ ) is highly variable and depends on the physical environment of Trp and Tyr residues, which are solely responsible for the protein fluorescence. The emission from Tyr residues has  $\lambda_{em}$  around 304 nm when no Trp residues are present, but if the protein has both Trp and Tyr, the emission spectrum is usually that of Trp alone, because of the resonance **energy transfer** from Tyr to Trp. Fluorescence spectroscopy is a very sensitive tool to monitor any structural change of protein molecules.

### 5.8. Rapid Kinetics

Methods such as stopped-flow and temperature-jump methods have been applied extensively to the study of the rates of enzyme catalysis or structural changes in proteins. Many valuable [kinetic](#) parameters have been obtained, and important intermediate species in such reactions have been proposed and identified. The time-scales of typical structural changes of proteins, and therefore of many biochemical reactions, including enzyme catalysis, range from microseconds to milliseconds.

## 6. Physical Properties of Proteins

### 6.1. Stability

Protein molecules, unlike synthetic polymers, normally have specific three-dimensional [protein structures](#) with unique thermodynamic and kinetic stabilities. The native three-dimensional structure of a protein is thermodynamically stabilized by the **cooperative** effects of [hydrogen bonds](#), **hydrophobic**, **electrostatic**, and [van der Waals interactions](#), plus in some cases covalent [disulfide bonds](#) (4, 5). This folded conformation can be destabilized by extremes of temperature, pH, and pressure and by **denaturants** (see [Protein Stability](#)). It can be unfolded reversibly, when the unfolded protein refolds upon changing the conditions, or irreversibly, when it is prevented from refolding, usually by covalent modification or by aggregation. Unfolding can be monitored by changes in either the physical or functional properties of the protein, as **denaturation** usually inactivates the protein.

### 6.2. Folding

**Protein folding** is the process of forming a specific three-dimensional structure of the protein from its newly synthesized state or after being unfolded. After C. B. Anfinsen showed that denatured and reduced **ribonuclease A** can be refolded into the correct native structure with full specific activity (6), the *in vitro* refolding of many proteins has been studied, and many have been shown to occur spontaneously. Subsequently, it has been recognized that the resulting “Anfinsen's dogma” applies primarily to small proteins without complicated **domain** or quaternary structures. The pathway of folding has been investigated extensively with a hope to isolate and characterize the intermediate state(s) that are important in determining the rate and direction of folding. Many unexpected and exciting details in the folding process of proteins have been steadily unraveled, but the molecular mechanism of protein folding still remains one of the most outstanding unresolved problems in science (6, 7). Recent developments on the role of proteins called **molecular chaperones**, which seem to help other proteins fold to their native conformations, are giving insight into how proteins gain their functional structures *in vivo*. (see [Protein Folding In Vivo](#)).

## 7. Classification of Proteins

Proteins are traditionally classified as either *globular* or [fibrous proteins](#). The individual domains or subunits of globular proteins are approximately spherical (see [Protein Structure](#)), whereas fibrous proteins usually have elongated structures. The classification is useful because of the lack of simple and visual words to describe complex structures of proteins. When approximated as ellipsoids of revolution, typical enzymes have an axial ratio of the long and short semi-axes of no more than five and thus are called “globular,” whereas proteins like silk fibroin, **collagens** in bones and tendons, [keratins](#), and muscle myosin are examples of fibrous proteins with greater axial ratios.



Another classical way of classifying proteins is based on their solubility properties in aqueous salt solutions. Historically, the fractionation of blood proteins was one of the most important contributions of protein chemistry to practical clinical applications. It was based on their solubilities in solutions with different concentrations of ammonium **sulfate**, at different pH's, or with different ionic strengths. **Globulins** were thus defined as those proteins that precipitated in 50% saturated ammonium sulfate and were insoluble in pure water. Proteins that were not precipitated in 50% saturated ammonium sulfate and were soluble in pure water were called **albumins**. Such names still persist, but today proteins are named after their functions as much as possible.

Yet another way of classifying proteins is based on whether a protein is functionally complete by itself or needs extra nonamino acid molecules other than polypeptide chain(s). Those that function alone are called *simple* proteins, and those that need other molecules are *conjugated* proteins. The polypeptide part of a conjugated protein is termed the **apoprotein**, and the other part is called the **prosthetic group**. The fully active complex is the **holoprotein**. Examples of conjugated proteins are (1) **glycoproteins** and mucoproteins (mucoids), where the former contain smaller amounts of carbohydrates than the latter; (2) nucleoproteins; (3) **lipoproteins**; (4) proteins, mostly enzymes, with small molecules, such as **pyridoxal phosphate** (PLP) and related compounds, **biotin**, flavins, NAD and NADP, folate, various metal ions in **metalloproteins**, retinal, vitamin **B<sub>12</sub>**, thiamine, **fatty acids**, etc.

Proteins are also classified according to their functions or by the major sites of their presence or synthesis. Thus, functional groupings are recognized, such as (1) **enzymes** (bio-catalysts); (2) enzyme inhibitors like **proteinase inhibitors**; (3) **hormones**; (4) **channel** proteins; (5) transport proteins; (6) structural proteins; (7) **DNA-binding proteins**; (8) **antibodies** (immunoglobulins); (9) **growth factors**; (10) lens proteins; (11) **ribosomal** proteins; (12) storage proteins, etc. According to their origins, proteins are often grouped as (1) blood (or plasma or serum) proteins; (2) egg proteins (egg yolk and egg white proteins); (3) milk proteins; (4) **membrane proteins**; (5) plant seed proteins, etc. Historically, different groups of proteins have been studied by researchers in different fields; eg, blood proteins by medical biochemists, milk and egg proteins by agricultural chemists, with specific purposes peculiar to their respective fields.

#### Bibliography

1. G. E. Schulz and R. H. Shimer (1979) *Principles of Protein Structure*, Springer-Verlag, New York.
2. A. White, P. Handler, E. Smith, R. L. Hill, and I. R. Lehman (1978) *Principles of Biochemistry*, McGraw-Hill Kogakusha, Tokyo.
3. K. J. Wilson (1990) In *Proteins: Form and Function* (R. L. Bradshaw and M. Purton, eds.), Elsevier, Cambridge, pp. 15–19.
4. C. N. Pace (1990) In *Proteins: Form and Function* (R. L. Bradshaw and M. Purton, eds.), Elsevier, Cambridge, pp. 117–123.
5. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.
6. C. B. Anfinsen, E. Harber, M. Sela, and F. H. White (1961) Proc. Natl. Acad. Sci. USA **47**, 1309–1314.
7. S. Lapanje (1978) *Physicochemical Aspects of Protein Denaturation*, Wiley, New York.
8. K. A. Dill and H. S. Chan (1997) From Levinthal to pathways to funnels. Nature Struc. Biol. **4**, 10–19.

## Proteolysis, Proteolytic

*Proteolysis*, or *proteolytic cleavage*, is the enzymatically catalyzed hydrolysis of [peptide bonds](#) in proteins [from *proteo*, the combining form for protein (from the Greek, meaning of primary importance) and lysis (from the Greek meaning dissolution)]. A large number of [enzymes](#) can catalyze peptide bond hydrolysis (see [Proteinases](#)). Some enzymes only act on one particular protein, others act on almost all proteins. Some may hydrolyze a single peptide bond in a protein (a process known as limited proteolysis, most commonly seen in [zymogen](#) activation), others may act only at selective peptide bonds in proteins (eg, those involving one particular type of residue, such as [arginine](#)), and yet others may have very broad specificity and cleave many types of peptide bonds.

Proteolysis is essential for an enormous range of biological phenomena ranging from digestion of dietary protein in the gastrointestinal tract, to [blood clotting](#), [hormone](#) production, and blood pressure regulation. It also has numerous commercial applications. Proteolysis can be inhibited by naturally occurring anti-proteinases, which are either proteins, peptide analogues, or nonpeptide organic or inorganic molecules (see [Proteinase Inhibitors](#)). Synthetic proteinase inhibitors are being developed as important drugs by the pharmaceutical industry.

## Proteome

The word “proteome” was created at the 1994 Conference on Genome and Protein Maps in Siena, Italy ([1](#)) and appeared in the literature for the first time in 1995 ([2-5](#)). Since then, the word has been widely adopted and within two years appeared in more than fifty publications ([6-8](#)).

Proteome means PROTEin complement expressed by a genOME. The proteome is all the expressed [genes](#) or [proteins](#) of a [genome](#). A cell type may display numerous subproteomes under different growing conditions, nutrition status, health, or disease. “Proteomics” is the use of quantitative measurements of the level of a protein or gene expression to characterize biological processes and to decipher the mechanism and control of **gene expression**. Proteome research or proteomics is the best path between genome and function studies.

The word proteome is not synonymous with [two-dimensional electrophoresis](#) (2D-PAGE), but 2D-PAGE is currently the most powerful method to separate thousands of proteins in parallel and to display many of them. This parallel process is necessary, as discussed below, because of the tremendous complexity of nature and the proteome.

### 1. Nature's Complexity

In the past, we were taught that one gene expresses one protein or [enzyme](#), but this statement is no longer strictly tenable. First, RNA **splicing** can produce several species of [messenger RNA](#) (mRNA), and the mRNA can be edited and modified before its translation. So, a single gene may give rise to several mRNAs. Second, most, if not all, **eukaryotic** proteins are modified covalently during or after [translation](#) (see [Post-Translational Modifications](#)); for example, more than 30% of eukaryote proteins are **phosphorylated** ([9](#)). Therefore, a single gene sequence might provide one or several [polypeptide chains](#), each of which may then be modified in various ways, leading to many versions of the final protein.

Based on a study of more than a hundred proteins identified by 2D-PAGE, one gene in *Escherichia*

*coli* produces on average 1.3 different polypeptides, in *Saccharomyces cerevisiae* 3 variant proteins, and in humans 3 cellular proteins and more than 20 proteins secreted into fluid (5). If the human genome contains 100,000 genes, which seems a reasonable estimate, between 300,000 and 500,000 protein variants might be expressed at different times; they will represent the full human proteome.

Interestingly, pigeons have 80 pairs of chromosomes; we have only 23. Frog and salamander cells contain between 1.5 and 20 times the amount of nucleic acid we have (see [C-Value](#)). Nature might have increased the complexity at the protein level of humans, in part, to keep the size of the genome within bounds while preventing a growing number of errors during replication.

Full genomes of several organisms have been sequenced recently. However, the life of these organisms is understood only at the level of the information content. Such information dictates how the amino acid residues are assembled to form the [polypeptide chain](#) (see [Translation](#)), the precursor of the mature co-and/or post-translationally modified proteins. We need to characterize most of these proteins, the proteome, before we can understand the structure and/or function of complex biological systems.

Each cell or tissue has different subproteomes, depending on their differentiation state, nutrition status, health, or disease. Proteins from many cell or tissue subproteomes can be detected in body fluids. It is indeed well known that dying cells release their contents into the lymph and blood. As an example, myocyte cytosolic proteins such as creatine phosphokinase or troponin are released into the circulation in myocardial infarction. One therefore could anticipate that most soluble human proteins from any cell or tissue could be found sometimes in the bloodstream, even if only in minute amounts. At least 100,000 different proteins could probably be found in plasma at certain times.

It is then quite obvious that the study of cell, tissue, or even body fluid subproteomes requires massively parallel and rapid protein analysis and identification techniques. The selected methodologies should be very sensitive over a large dynamic range. The speed of all processes should be maximized by automation and size reduction.

## 2. Parallel Concentration and Separation of Protein Samples

Most biological samples are highly complex, and the range of protein concentrations encountered is tremendous, exceeding a factor of  $10^6$  for cell subproteomes and  $10^{12}$  for body fluids. Consequently, concentration and separation steps are necessary today for most proteome analyses. Filtration, [centrifugation](#), or recycling [electrophoresis](#) provides means to concentrate proteins that are present only in minute amounts. Liquid [chromatography](#) and/or electrophoretic techniques allow the powerful separation of proteome content, especially when they are linked to [mass spectrometry](#).

Currently, high-resolution 2D-PAGE techniques are still the most powerful for concentrating and separating in parallel thousands of proteins (10-15). The two dimensions are, first, an [isoelectric focusing](#) separation, where the proteins are separated along a pH gradient by charge or isoelectric point, and, second, an [SDS-PAGE](#) separation done orthogonally to the first one, where the proteins are separated by size. For the first dimension, the pH gradient can be generated by carrier ampholytes (10, 11, 16) or immobilized pH gradients (17-20) (see **Isoelectric focusing**). Immobilized pH gradients allow high loading capacity (21, 22), by reswelling precast dry IPG strips with sample solution (13, 23-25), and give highly reproducible separations even between laboratories, since the gradient is covalently attached to the gel matrix (26). Alkaline immobilized gradients have been designed that extend from pH 6 to 12 (27, 28). The second dimension SDS PAGE separation may be replaced by mass spectrometers, which give very precise protein mass measurements (29-34).

After an SDS-PAGE second dimension, the protein can be **blotted** or eluted from the gel for further chemical analysis, such as [Edman Degradation](#) amino acid microsequencing (34, 35). Several gel, or

membrane, punching machines have been designed to transfer gel or membrane spots automatically to robotized chemical stations and mass spectrometers. The parallel 2D-PAGE process then becomes sequential. The blotted membrane can also be scanned by matrix-assisted laser desorption ionization (MALDI) mass spectrometry (33). The mass spectrometry techniques can identify proteins, as well as simply detect them (30-32, 36-44).

### 3. Protein Detection and Identification

Numerous protocols for staining proteins have been designed (see [Coomassie Brilliant Blue](#), [Gel Electrophoresis](#), [Ponceau S](#)). The most sensitive until recently has been **silver staining** (45-56), although fluorescent dyes now provide similar sensitivity. If further protein identification and characterization are desired, however, no staining methods that modify the protein should be used (57, 58). The limit of detection varies from protein to protein, but at least 0.1 ng/mm<sup>2</sup> is generally required. Most of the colloidal stains provide a linear correlation over several factors of 10 between the protein concentration and the staining intensity. The slope of this correlation varies from protein to protein, however, with the staining properties being influenced by the amino acid composition, and even the amino acid sequence, of the polypeptide chain (59).

Proteins can also be detected if they are **radiolabeled**. <sup>35</sup>S-methionine labeling is most frequently used. In some cases, proteins are <sup>14</sup>C-labeled sequentially with each amino acid to obtain in parallel the amino acid composition of each spot (60, 61). This provides a massively parallel way to detect, and also to identify, many 2D-PAGE spots. Derivatization with Fmoc (9-fluorenylmethyl chloroformate) after acidic hydrolysis provides information about the amino acid composition at high detection sensitivity (pmol), but such analyses are sequential, not parallel (62, 63). The amino acid composition can be used to identify a protein and is successful on average 50% of the time when the genome is known. It is useful also for cross-species comparison when a genome is unknown (64). The identification of a few residues of amino acid sequence from the N-terminal, C-terminal, or even internal regions of the polypeptide chain is very powerful and discriminating: this has been called *sequence tagging* (65-68).

For small genomes, such as that of *E. coli*, a 3-residue sequence tag from the C-terminus, or 4-residues from the N-terminus, are usually sufficient for unambiguous protein identification. Once the full human genome is known, a 6-residue tag from either end will probably suffice. In addition, comparisons between the predicted and measured isoelectric points (pI) and the mass further refine the protein identification process. Indeed, the pI of a polypeptide chain that has not been post-translationally modified can be predicted, and it corresponds with very high accuracy to the value measured in immobilized pH gradients under denaturing conditions (69, 70). Prediction of the mass by SDS-PAGE is rather unreliable, but developments in mass spectrometry fortunately circumvent this problem.

Peptide **fingerprinting** and partial sequencing are certainly the most attractive tools to identify a protein rapidly when the corresponding genome is known (64, 71-73). In addition, [peptide mapping](#) allows the detection and identification of numerous post-translational modifications. Proteins are cleaved either chemically (by cyanogen bromide, CNBr) or enzymatically (**trypsin**, Lys C, etc.) in the gel, in the eluted spot solution, or on the membrane. The masses of the peptide fragments produced are then compared to the theoretical masses calculated from the known proteins of that species. When the masses of four or five fragments correspond to expected masses, the protein is identified with confidence. Any differences between the remaining unmatched peptide masses and the expected ones are often the result of one or more post-translational modifications, which then can be deduced from the difference. Table 1 lists several cyber tools available on the World Wide Web in 1997 for protein identification and characterization.

**Table 1. Several Sequence Analysis Tools Available on ExPASy Server <sup>a, b</sup>**

---

| <b>Protein identification</b>       |   |
|-------------------------------------|---|
| Swiss-Shop ( <i>L</i> )             | A sequence alerting system for SWISS-PROT that allows one to obtain automatically (by e-mail) new sequence entries relevant to one's field(s) of interest.                                  |
| AACompIdent ( <i>L</i> )            | Identify a protein by its amino acid composition.   |
| AACompSim ( <i>L</i> )              | Compare the amino acid composition of a SWISS-PROT entry with all other entries.  |
| MultiIdent ( <i>L</i> )             | Identify proteins with pI, MW, amino acid composition, sequence tag, and peptide mass fingerprinting data.  |
| TagIdent ( <i>L</i> )               | Get the SWISS-PROT proteins closest to a given pI and Mw and identify proteins with a sequence tag.   |
| PeptideMass ( <i>L</i> )            | Calculate masses of peptides and their post-translational modifications for a SWISS-PROT or TREMBL entry or user sequence.  |
| Compute pI/Mw ( <i>L</i> )          | Compute the theoretical pI and Mw from a SWISS-PROT or TREMBL entry or for a user sequence.   |
| ProteinProspector                   | A variety of tools from UCSF (MS-Fit, MS-Tag, MS-Digest, etc) for mining sequence databases in conjunction with mass proteometry experiments (see also the mirror in UCL-Ludwig in the UK). |
| PROWL                               | Protein information Retrieval On-line WWW Lab from the MS groups at Rockefeller and New York Universities.  |
| PeptideSearch                       | Peptide mass fingerprint tool from EMBL Heidelberg.   |
| MOWSE                               | From Daresbury Laboratory.  |
| <b>Pattern and profile searches</b> |   |
| ScanProsite ( <i>L</i> )            | Scan a sequence against PROSITE or a pattern against SWISS-PROT and TREMBL.   |
| ProfileScan                         | Scan a sequence against the profile entries in PROSITE.   |
| Pfam                                | Scan a sequence against the PFAM HMM protein families db. A HMM search at Washington University/Sanger Centre (Hinxton, UK).  |
| FPAT                                | Regular expression searches in protein databases.   |
| PRATT                               | Interactively generates conserved patterns from a series of unaligned proteins.   |
| <b>Primary structure analysis</b>   |   |
| ProtParam ( <i>L</i> )              | Physico-chemical parameters of a protein sequence (composition, extinction coefficient, etc).   |
| ProtScale ( <i>L</i> )              | Amino acid scale representation (hydrophobicity, other conformational parameters, etc).   |
| SAPS                                | Statistical analysis of protein sequences at ISREC (also available at EBI).   |

|                                       |  |
|---------------------------------------|--|
| PSORT                                 | Prediction of protein sorting signals and localizationsites.                             |
| SignalP                               | Prediction of signal peptide cleavage sites.   |
| NetOGlyc                              | Prediction of type O-glycosylation sites in mammalian proteins.                          |
| NetPicoRNA                            | Prediction of proteinase cleavage sites in picornavial proteins.                         |
| Coils                                 | Prediction of coiled coil regions in proteins (Lupas's method).                          |
| Paircoil                              | Prediction of coiled coil regions in proteins (Berger's method).                         |
| REPRO                                 | Recognition of protein sequence repeats at EMBL.   |
| Protein Colourer                      | Tool for coloring your amino acid sequence.  |
| RandSeq ( <i>L</i> )                  | Random protein sequence generator.   |
| <b>Secondary structure prediction</b> |  |
| AntheProt                             | Institute of Biology and Chemistry of Proteins (IBCP)/Lyon.                              |
| BCM PSSP                              | Baylor College of Medicine.  |
| DSC                                   | Discrimination of protein secondary structure at BMM (ICRF/London).                      |
| GOR                                   | Garnier, Osgoodthorpe, and Robson (GOR) secondary structure prediction method (at SBDS). |
| nnPredict                             | University of California at San Francisco (UCSF).  |
| PredictProtein                        | PHDsec, PHDacc, PHDhtm, PHDtopology, PHDthreader, MaxHom, EvalSec from EMBL.             |
| PREDATOR                              | Protein secondary structure prediction from single sequence at EMBL (Argos' group).      |
| PSA                                   | BioMolecular Engineering Research Center (BMERC)/Boston.                                 |
| SSPRED                                | Protein secondary structure prediction from aligned sequences at EMBL (Argos' group).    |
| <u>Tertiary structure</u>             |  |
| Swiss-Model ( <i>L</i> )              | An automated knowledge-based protein modeling server.                                    |
| Swiss-PdbViewer ( <i>L</i> )          | A program to analyze and superimpose protein three-dimensional structures.               |

---

<sup>a</sup> The tools marked by (*L*) are local to the ExPASy server. The remaining tools are pointers to other servers.

<sup>b</sup> Last modified October 21, 1997 by ELG and simplified by the author.

#### 4. Protein Characterization, Modeling, and Function Prediction

First, any covalent modifications of the protein need to be identified. Many dozens of different types of modification have been reported, and many influence the protein's charge, [hydrophobicity](#), conformation, and **half-life** (see [Post-Translational Modifications](#)). Numerous methods have been used to identify specific post-translational modifications: (1) blotting with [lectins](#) to detect specific

sugars (74, 75) (2) detecting phosphoamino acids by immunoblotting or  $^{32}\text{P}$  radiolabeling, (3) Edman degradation sequencing, and (4) mass spectrometry.

Then, characterizing the protein and modeling its three-dimensional structure can be used to predict its function. The shape of the protein and how it interacts with other molecules provide functional information. Often, a protein structure can be predicted directly from its amino acid sequence when the structures of **homologous** proteins belonging to the same family have already been determined (see **Homology modeling**) (76, 77).

From separating the proteins to their identification and characterization, or to the final modeling and prediction of function, the extensive use of [databases](#) and computers is required. An up-to-date repository of genome and proteome information is absolutely essential in proteomic research, and a set of related cyber tools greatly enhances the potential to make significant discoveries.

## 5. Proteome Databases and Related Cyber Tools

Annotated protein databases (78-82) serve as a platform to link any virtual or real experiment with genome and biomedical information. Virtual protein chemistry tools can be applied to the protein database to predict results (Table 1). Comparison between experimental and computer data permit the identification and characterization of numerous proteins, as mentioned above.

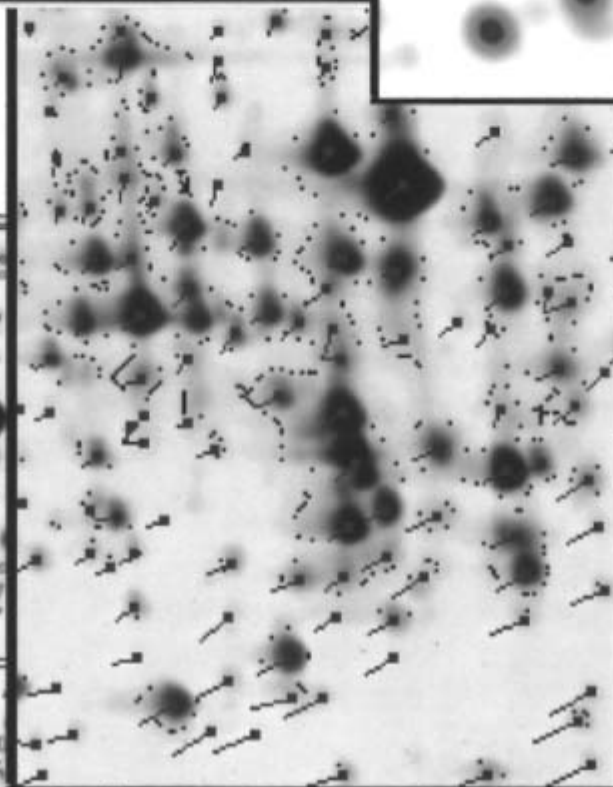
2D-PAGE image analysis provides qualitative and quantitative information on changes in protein or gene expression and can display [epigenetic](#) regulation (Fig. 1). Principle component analysis (83, 84), heuristic clustering tools (85), neuronal networks (86-90), and wavelet decomposition algorithms (91, 92) cluster gel images in meaningful groups and highlight the most significant proteome alteration or modification in a wide variety of experiments, in toxicology, health, or disease states. Protein databases establish a link between 2D-PAGE image databases (93) and genome and biomedical information (site currently unavailable) (Fig. 2). Many other databases should be linked to the protein structure database, such as those for protein motifs, three-dimensional structures, mass spectra data, etc (see [Databases](#)). Protein physicochemical parameters should be used to display and link the result of all experiments in proteome research such as the precise protein size, charge, composition, amino acid sequence, etc. Artificial two-dimensional protein maps can be reconstructed to display the expected results of any proteome analysis.

**Figure 1.** Image analysis software (MELANIE) showing the analysis of a 2D-PAGE gel from the detection, quantification, and identification over the Internet.

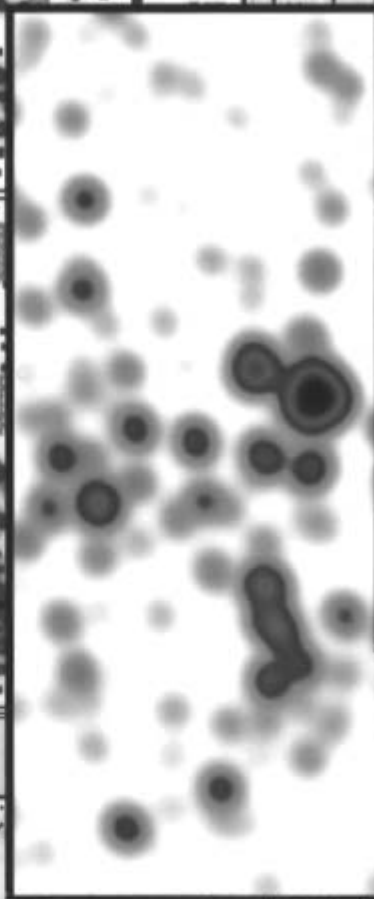
Analyze Database Help  
 of Features...  
 of DNA Features...  
 Gels

McLennan - [McLView]  
 File Edit Select Show Stack Image Process Analyze Database Help

94-0005.mcl "Signal"  
 94-0005.mcl "Signal"  
 94-0005.mcl "Signal"  
 94-0005.mcl "Signal"



| FeatureID | VOL   | %VOL  | AREA  | %A    |
|-----------|-------|-------|-------|-------|
| 344       | 0.009 | 0.003 | 0.276 | 0.121 |
| 345       | 0.009 | 0.003 | 0.245 | 0.137 |
| 346       | 0.045 | 0.016 | 0.613 | 0.353 |
| 347       | 0.327 | 0.118 | 1.531 | 1.099 |
| 348       | 0.208 | 0.075 | 1.011 | 1.019 |
| 349       | 0.141 | 0.051 | 0.888 | 0.737 |
| 350       | 2.899 | 1.047 | 6.860 | 1.832 |
| 351       | 0.401 | 0.145 | 1.282 | 1.428 |



AC

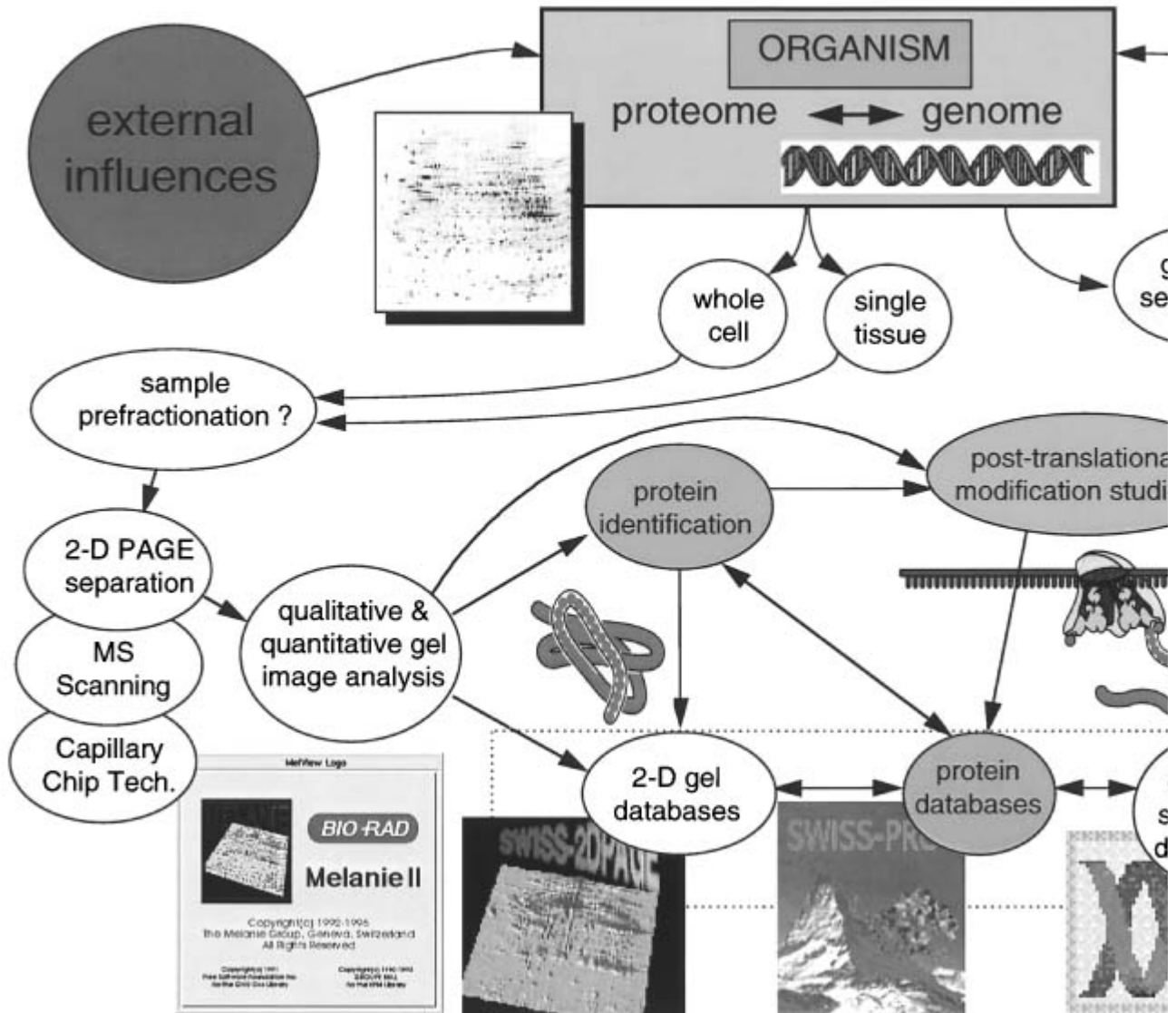
VOL 3

VOL 1.2

Gels-analyzed/94-0002 a  
 Gels-analyzed/94-0005 b  
 Gels-analyzed/ECOLI c  
 Gels-analyzed/Synthetic d



**Figure 2.** Diagram of the relationships between proteome and genome research, highlighting the essential roles of genom WWW servers.



## 6. Present and Future Frontiers

The behavior of proteins is difficult to predict (Table 2). Their building blocks are far more complex than those of nucleic acids. Their solubility is often poor in water and even sometimes in the presence of **chaotropic** agents and strong **detergents**. Many proteins are unstable or quite reactive and have a short half-life. Their concentrations in a very complex biological fluid can vary tremendously. There is no necessary correlation between mRNA level and protein concentration (94). Indeed, **expressed sequence tag** (EST) databases are derived from mRNA abundance and not from the amount of protein. Therefore, we cannot predict protein abundance from genome or mRNA studies. Today, separation techniques such as 2D-PAGE or **capillary zone electrophoresis** and mass spectrometry show only the tip of the iceberg as far as the complexity of organisms goes, but they will undoubtedly progress. Even when the technological problems are solved, real physiological experiments will be required in many cases to study further epigenetic regulation in the entire

organism and to measure the effects of external influences ([95-97](#)). At that point, proteome research will bring a new dimension to combinatorial chemistry and genome studies and will hasten new therapeutic developments.

**Table 2. Differences between Proteins and Nucleic Acids**

| Chemical Family                                | Nucleic Acids                                     | Proteins/Peptides  |
|--|---|--|
| Belongs to                                     | Genomes   | Proteomes  |
| Level  | Information                                       | Product  |
| Number of building blocks                      | 4+  | ≥20  |
| Solubility in water                            | High  | Often very low   |
| Natural or artificial blockage when sequencing | No, but sometimes problems with GC rich stretches | Yes, very often  |
| Prediction of behavior                         | Easy  | Difficult  |
| Number of specific cleavage enzymes            | Very many (>300)                                  | Few (<12)  |
| Possible propagation or amplification          | Easy  | No (except <a href="#">prions</a> or enzymes that catalyze their formation from pro-enzyme?) |
| Number of residues sequenced/day (1997)        | >10,000   | <50  |

### Bibliography

1. M. R. Wilkins, C. Pasquali, R. D. Appel, K. Ou, O. Golaz, J. C. Sanchez, J. X. Yan, A. A. Gooley, G. Hughes, I. Humphery-Smith, K. L. Williams, and D. F. Hochstrasser (1996) *Bio/Technology* **14**, 61–65.
2. M. R. Wilkins, J. C. Sanchez, A. A. Gooley, R. D. Appel, I. Humphery-Smith, D. F. Hochstrasser, and K. L. Williams (1995) *Biotechnol. Genet. Eng. Rev.* **13**, 19–50.
3. V. C. Wasinger, S. J. Cordwell, A. Cerpa-Poljak, J. X. Yan, A. A. Gooley, M. R. Wilkins, M. W. Duncan, R. Harris, K. L. Williams, and I. Humphery-Smith (1995) *Electrophoresis* **16**, 1090–1094.
4. A. Shevchenko, O. N. Jensen, A. V. Podtelejnikov, F. Sagliocco, M. Wilm, O. Vorm, P. Mortensen, A. Shevchenko, H. Boucherie, and M. Mann (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14440–14445.
5. M. R. Wilkins, J. C. Sanchez, K. L. Williams, and D. F. Hochstrasser (1996) *Electrophoresis* **17**, 830–838.
6. K. L. Williams and D. F. Hochstrasser (1997) In *Proteome Research: New Frontiers in Functional Genomics* (M. R. Wilkins, K. L. Williams, R. D. Appel, and D. F. Hochstrasser, eds.), Springer-Verlag, Berlin-Heidelberg-New York, pp. 1–12.
7. D. Swinbanks (1995) *Nature* **378**, 653.
8. P. Kahn (1995) *Science* **270**, 369–370.

9. E. H. Fisher (1997) In *Proteome Research: New Frontiers in Functional Genomics* (M. R. Wilkins, K. L. Williams, R. D. Appel, and D. F. Hochstrasser, eds.), Springer-Verlag, Berlin-Heidelberg-New York.
10. J. Klose (1975) *Humangenetik* **26**, 231–243.
11. P. H. O'Farrell (1975) *J. Biol. Chem.* **250**, 4007–4021.
12. D. F. Hochstrasser, M. G. Harrington, A. C. Hochstrasser, M. J. Miller, and C. R. Merrill (1988) *Anal. Biochem.* **173**, 424–435.
13. B. Bjellqvist, J. C. Sanchez, C. Pasquali, F. Ravier, N. Paquet, S. Frutiger, G. J. Hughes, and D. Hochstrasser (1993) *Electrophoresis* **14**, 1375–1378.
14. B. Bjellqvist, C. Pasquali, F. Ravier, J. C. Sanchez, and D. Hochstrasser (1993) *Electrophoresis* **14**, 1357–1365.
15. A. Gorg, G. Boguth, C. Obermaier, A. Posch, and W. Weiss (1995) *Electrophoresis* **16**, 1079–1086.
16. J. Klose and U. Kobalz (1995) *Electrophoresis* **16**, 1034–1059.
17. B. Bjellqvist, K. Ek, P. G. Righetti, E. Gianazza, A. Gorg, R. Westermeier, and W. Postel (1982) *J. Biochem. Biophys. Meth.* **6**, 317–339.
18. A. Gorg, W. Postel, S. Gunther, J. Weser, J. R. Strahler, S. M. Hanash, L. Somerlot, and R. Kuick (1988) *Electrophoresis* **9**, 37–46.
19. A. Gorg, W. Postel, and S. Gunther (1988) *Electrophoresis* **9**, 531–546.
20. A. Gorg (1993) *Biochem. Soc. Trans.* **21**, 130–132.
21. S. M. Hanash and J. R. Strahler (1989) *Nature* **337**, 485–486.
22. S. M. Hanash, J. R. Strahler, J. V. Neel, N. Hailat, R. Melhem, D. Keim, X. X. Zhu, D. Wagner, D. A. Gage, and J. T. Watson (1991) *Proc. Natl. Acad. Sci. USA* **88**, 5709–5713.
23. T. Rabilloud, C. Valette, and J. J. Lawrence (1994) *Electrophoresis* **15**, 1552–1558.
24. T. Rabilloud, M. Vincon, and J. Garin (1995) *Electrophoresis* **16**, 1414–1422.
25. J. C. Sanchez, V. Rouge, M. Pisteur, F. Ravier, L. Tonella, M. Moosmayer, M. R. Wilkins, and D. F. Hochstrasser (1997) *Electrophoresis* **18**, 324–327.
26. J. M. Corbett, M. J. Dunn, A. Posch, and A. Gorg (1994) *Electrophoresis* **15**, 1205–1211.
27. P. G. Righetti, A. Bossi, A. Gorg, C. Obermaier, and G. Boguth (1996) *J. Biochem. Biophys. Meth.* **31**, 81–91.
28. A. Gorg, C. Obermaier, G. Boguth, A. Csordas, J. J. Diaz, and J. J. Madjar (1997) *Electrophoresis* **18**, 328–337.
29. J. A. Castro, C. Koster, and C. Wilkins (1992) *Rapid Commun. Mass Spectrom.* **6**, 239–241.
30. H. H. Rasmussen, E. Mortz, M. Mann, P. Roepstorff, and J. E. Celis (1994) *Electrophoresis* **15**, 406–416.
31. E. Mortz, O. Vorm, M. Mann, and P. Roepstorff (1994) *Biol. Mass Spectrom.* **23**, 249–261.
32. S. D. Patterson (1995) *Electrophoresis* **16**, 1104–1114.
33. C. Eckerskorn, K. Strupat, D. Schleuder, D. Hochstrasser, J. C. Sanchez, F. Lottspeich, and F. Hillenkamp (1997) *Anal. Chem.* **69**, 2888–2892.
34. S. L. Cohen, and B. T. Chait (1997) *Anal. Biochem.* **247**, 257–267.
35. N. LeGendre and P. Matsudaira (1988) *Biotechniques* **6**, 154–159.
36. J. R. D. Yates, S. Speicher, P. R. Griffin, and T. Hunkapiller (1993) *Anal. Biochem.* **214**, 397–408.
37. M. Mann and M. Wilm (1994) *Anal. Chem.* **66**, 4390–4399.
38. S. D. Patterson, D. Thomas, and R. A. Bradshaw (1996) *Electrophoresis* **17**, 877–891.
39. A. L. McCormack, D. M. Schieltz, B. Goode, S. Yang, G. Barnes, D. Drubin, and J. R. R. Yates (1997) *Anal. Chem.* **69**, 767–776.

40. J. R. R. Yates, A. L. McCormack, D. Schieltz, E. Carmack, and A. Link (1997) *J. Protein Chem.* **16**, 495–497.
41. A. Shevchenko, M. Wilm, and M. Mann (1997) *J. Protein Chem.* **16**, 481–490.
42. A. Shevchenko, I. Chernushevich, W. Ens, K. G. Standing, B. Thomson, M. Wilm, and M. Mann (1997) *Rapid Commun. Mass Spectrom.* **11**, 1015–1024.
43. A. R. Dongre, J. K. Eng, and J. R. R. Yates (1997) *Trends Biotechnol.* **15**, 418–425.
44. G. Li, M. Waltham, N. L. Anderson, E. Unsworth, A. Treston, and J. N. Weinstein (1997) *Electrophoresis* **18**, 391–402.
45. R. C. D. Switzer, C. R. Merrill, and S. Shifrin (1979) *Anal. Biochem.* **98**, 231–237.
46. C. R. Merrill, R. C. Switzer, and M. L. Van Keuren (1979) *Proc. Natl. Acad. Sci. USA* **76**, 4335–4339.
47. C. R. Merrill, D. Goldman, S. A. Sedman, and M. H. Ebert (1981) *Science* **211**, 1437–1438.
48. C. R. Merrill, M. L. Dunau, and D. Goldman (1981) *Anal. Biochem.* **110**, 201–207.
49. C. R. Merrill, D. Goldman, and M. L. Van Keuren (1983) *Methods Enzymol.* **96**, 230–239.
50. C. R. Merrill, D. Goldman, and M. L. Van Keuren (1984) *Methods Enzymol.* **104**, 441–447.
51. C. R. Merrill and M. E. Pratt (1986) *Anal. Biochem.* **156**, 96–110.
52. T. Rabilloud, G. Carpentier, and P. Tarroux (1988) *Electrophoresis* **9**, 288–291.
53. T. Rabilloud (1990) *Electrophoresis* **11**, 785–794.
54. T. Rabilloud (1992) *Electrophoresis* **13**, 429–439.
55. T. Rabilloud, V. Brodard, G. Peltre, P. G. Righetti, and C. Ettori (1992) *Electrophoresis* **13**, 264–266.
56. T. Rabilloud, L. Vuillard, C. Gilly, and J. J. Lawrence (1994) *Cell. Mol. Biol. (Noisy-le-grand)* **40**, 57–75.
57. A. Shevchenko, M. Wilm, O. Vorm, and M. Mann (1996) *Anal. Chem.* **68**, 850–858.
58. M. Wilm, A. Shevchenko, T. Houthaeve, S. Breit, L. Schweigerer, T. Fotsis, and M. Mann (1996) *Nature* **379**, 466–469.
59. C. R. Merrill, M. E. Bisher, M. Harrington, and A. C. Steven (1988) *Proc. Natl. Acad. Sci. USA* **85**, 453–457.
60. K. E. Latham, J. I. Garrels, and D. Solter (1993) *Methods Enzymol.* **225**, 473–489.
61. J. I. Garrels, B. Futcher, R. Kobayashi, G. I. Latter, B. Schwender, T. Volpe, J. R. Warner, and C. S. McLaughlin (1994) *Electrophoresis* **15**, 1466–1486.
62. J. X. Yan, M. R. Wilkins, K. Ou, A. A. Gooley, K. L. Williams, J. C. Sanchez, O. Golaz, C. Pasquali, and D. F. Hochstrasser (1996) *J. Chromatogr. A* **736**, 291–302.
63. O. Golaz, M. R. Wilkins, J. C. Sanchez, R. D. Appel, D. F. Hochstrasser, and K. L. Williams (1996) *Electrophoresis* **17**, 573–579.
64. M. R. Wilkins and K. L. Williams (1997) *J. Theor. Biol.* **186**, 7–15.
65. M. R. Wilkins, E. Gasteiger, J. C. Sanchez, R. D. Appel, and D. F. Hochstrasser (1996) *Curr. Biol.* **6**, 1543–1544.
66. E. Mortz, P. B. O'Connor, P. Roepstorff, N. L. Kelleher, T. D. Wood, F. W. McLafferty, and M. Mann (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8264–8267.
67. A. A. Gooley, K. Ou, J. Russell, M. R. Wilkins, J. C. Sanchez, D. F. Hochstrasser, and K. L. Williams (1997) *Electrophoresis* **18**, 1068–1072.
68. M. R. Wilkins, K. Ou, R. D. Appel, J. C. Sanchez, J. X. Yan, O. Golaz, V. Farnsworth, P. Cartier, D. F. Hochstrasser, K. L. Williams, and A. A. Gooley (1996) *Biochem. Biophys. Res. Commun.* **221**, 609–613.
69. B. Bjellqvist, G. J. Hughes, C. Pasquali, N. Paquet, F. Ravier, J. C. Sanchez, S. Frutiger, and D. Hochstrasser (1993) *Electrophoresis* **14**, 1023–1031.
70. B. Bjellqvist, B. Basse, E. Olsen, and J. E. Celis (1994) *Electrophoresis* **15**, 529–539.

71. C. O'Brien (1996) *Mol. Med. Today* **2**, 316.
72. M. J. Wise, T. G. Littlejohn, and I. Humphery-Smith (1997) *Electrophoresis* **18**, 1399–1409.
73. M. R. Wilkins, I. Lindskog, E. Gasteiger, A. Bairoch, J. C. Sanchez, D. F. Hochstrasser, and R. D. Appel (1997) *Electrophoresis* **18**, 403–408.
74. P. Gravel, O. Golaz, C. Walzer, D. F. Hochstrasser, H. Turler, and L. P. Balant (1994) *Anal. Biochem.* **221**, 66–71.
75. O. Golaz, P. Gravel, C. Walzer, H. Turler, L. Balant, and D. F. Hochstrasser (1995) *Electrophoresis* **16**, 1187–1189.
76. M. C. Peitsch, M. R. Wilkins, L. Tonella, J. C. Sanchez, R. D. Appel, and D. F. Hochstrasser (1997) *Electrophoresis* **18**, 498–501.
77. M. C. Peitsch (1997) *Large-Scale Comparative Protein Modelling*, Springer-Verlag, Berlin-Heidelberg-New York, pp. 177–186.
78. A. Bairoch (1996) *Nucleic Acids Res.* **24**, 221–222.
79. A. Bairoch and R. Apweiler (1997) *Nucleic Acids Res.* **25**, 31–36.
80. A. Bairoch, P. Bucher, and K. Hofmann (1997) *Nucleic Acids Res.* **25**, 217–221.
81. R. Apweiler, A. Gateau, S. Contrino, M. J. Martin, V. Junker, C. O'Donovan, F. Lang, N. Mitaritonna, S. Kappus, and A. Bairoch (1997) *Ismb* **5**, 33–43.
82. A. Bairoch and R. Apweiler (1997) *J. Mol. Med.* **75**, 312–316.
83. T. Pun, D. F. Hochstrasser, R. D. Appel, M. Funk, V. Villars-Augsburger, and C. Pellegrini (1988) *Appl. Theor. Electrophor.* **1**, 3–9.
84. C. Roch, T. Pun, D. F. Hochstrasser, and C. Pellegrini (1989) *Comput. Med. Imaging Graph* **13**, 383–391.
85. R. Appel, D. Hochstrasser, C. Roch, M. Funk, A. F. Muller, and C. Pellegrini (1988) *Electrophoresis* **9**, 136–142.
86. F. H. Grus and C. W. Zimmermann (1997) *Electrophoresis* **18**, 1120–1125.
87. C. Kesmir, I. Sondergaard, and K. Jensen (1995) *Electrophoresis* **16**, 927–933.
88. J. N. Weinstein, T. Myers, J. Buolamwini, K. Raghavan, W. van Osdol, J. Licht, V. N. Viswanadhan, K. W. Kohn, L. V. Rubinstein, A. D. Koutsoukos, and et al. (1994) *Stem Cells (Dayt)* **12**, 13–22.
89. M. A. Kratzer, B. Ivandic, and A. Fateh-Moghadam (1992) *J. Clin. Pathol.* **45**, 612–615.
90. I. Sondergaard, B. N. Krath, and M. Hagerup (1992) *Electrophoresis* **13**, 411–415.
91. A. Zahnd, J. D. Tissot, and D. F. Hochstrasser (1993) *Appl. Theor. Electrophor.* **3**, 321–328.
92. A. Zahnd, M. Funk, P. Vaudaux, D. Lew, J. R. Scherrer, and D. F. Hochstrasser (1994) *Appl. Theor. Electrophor.* **4**, 19–24.
93. R. D. Appel, A. Bairoch, J. C. Sanchez, J. R. Vargas, O. Golaz, C. Pasquali, and D. F. Hochstrasser (1996) *Electrophoresis* **17**, 540–546.
94. L. Anderson and J. Seilhamer (1997) *Electrophoresis* **18**, 533–537.
95. R. Strohman (1994) *Biotechnology (N Y)* **12**, 156–164.
96. R. C. Strohman (1995) *Integr. Physiol. Behav. Sci.* **30**, 273–282.
97. R. C. Strohman (1997) *Nat. Biotechnol.* **15**, 194–200.
98. A. A. Gooley and N. H. Packer (1997) *Proteome Research: New Frontiers in Functional Genomics* (M. R. Wilkins, K. L. Williams, R. D. Appel, and D. F. Hochstrasser, eds.), Springer-Verlag, Berlin-Heidelberg-New York, pp. 65–91.

### Suggestions for Further Reading

99. Overall view of proteome research M. R. Wilkins, K. L. Williams, R. D. Appel, and D. F. Hochstrasser (eds.) (1997) *Proteome Research: New Frontiers in Functional Genomics*, Springer-Verlag, Berlin-Heidelberg-New York, pp. 1–243. *Sample preparation*

100. T. Rabilloud, C. Adessi, A. Giraudel, and J. Lunardi (1997) Improvement of the solubilization of proteins in two-dimensional electrophoresis with immobilized pH gradients, *Electrophoresis* **18**, 307–316. *First dimension separation*
101. A. Gorg, G. Boguth, C. Obermaier, A. Posch, and W. Weiss (1995) Two-dimensional polyacrylamide gel electrophoresis with immobilized pH gradients in the first dimension (IPG-Dalt): The state of the art and the controversy of vertical versus horizontal systems, *Electrophoresis* **16**, 1079–1086. *Second dimension separation*
102. D. F. Hochstrasser, M. G. Harrington, A. C. Hochstrasser, M. J. Miller, and C. R. Merrill (1988) Methods for increasing the resolution of two-dimensional protein electrophoresis, *Anal. Biochem.* **173**, 424–435. *Protein silver staining*
103. T. Rabilloud, L. Vuillard, C. Gilly, and J. J. Lawrence (1994) Silver-staining of proteins in polyacrylamide gels: A general overview, *Cell Mol. Biol. (Noisy-le-grand)* **40**, 57–75. *Mass-spectrometry approaches*
104. S. D. Patterson and R. Aebersold (1995) Mass spectrometric approaches for the identification of gel-separated proteins, *Electrophoresis* **16**, 1791–1814. *Computer analysis*
105. M. R. Wilkins, D. F. Hochstrasser, J. C. Sanchez, A. Bairoch, and R. D. Appel (1996) Integrating two-dimensional gel databases using the Melanie II software, *Trends Biochem. Sci.* **21**, 496–497. *Clinical applications*
106. D. F. Hochstrasser and J. D. Tissot (1993) "Clinical applications of 2-D PAGE." In *Advances in Electrophoresis* (A. Chrambach, M. J. Dunn, and B. J. Radola, eds.), VCH, pp. 270–375. *Short selection of Web sites to surf*
107. (site currently unavailable)
108. <http://biobase.dk/cgi-bin/celis>
109. <http://www-lecb.ncifcrf.gov/2dwgDB/>
110. <http://www.proteome.com/YPDhome.html>
111. <http://www.harefield.nthames.nhs.uk/nhli/protein>

## Proton Gradient

The term “proton gradient” is a somewhat inappropriate name for the transmembrane electrochemical proton potential that is involved in [chemiosmotic coupling](#). This term obscures the fact that the gradient is actually composed of two terms: the difference between the proton activity (concentration) or pH, and the membrane potential. See [Chemiosmotic Coupling](#) for further details.

## Proton Motive Force

The *proton motive force*, abbreviated  $Dp$ , is the electrochemical proton potential across a [membrane](#) expressed in terms of volts, which is the energy available to drive ATP formation, metabolite transport, or motility. At 25°C, the value of  $Dp$  in millivolts is approximately equal to 59 (DpH) + DY, where DY is the membrane potential in millivolts. See [Chemiosmotic Coupling](#) for

further details.

## Protoplast Fusion

The cells of bacteria, fungi, and plants may be converted to protoplasts (or *spheroplasts*) by the enzymatic degradation of their cell walls. This opens up a range of new possibilities for the investigation of cells. Protoplasts may be fused with other protoplasts from strains of the same or different species or even with animal cells (1) or liposomes (2). The motivations for carrying out such protoplast fusions are diverse and include genetic investigations of species that are not amenable to normal genetic study, the improvement of strains of organisms of agronomic or biotechnological importance, the study of nucleo-cytoplasmic interactions and the introduction of foreign materials into cells (see [Transfection](#)).

The spontaneous fusion of protoplasts is usually quite rare, and the event must be induced in some way. The fusion technique itself involves two processes: 1) bringing the two protoplasts into close membrane contact; and 2) limited, localized disruption of adjacent membranes, permitting the formation of cytoplasmic continuities between the cells. The presence of residual cell wall material on the surface of spheroplasts may interfere with the formation of close contacts between the cytoplasmic membranes, but the fusion of spheroplasts is quite often possible although it occurs at a lower frequency than protoplast fusion (3, 4).

Two basic techniques have been employed for inducing fusion. The first, and currently most widely used, technique involves chemical induction. Polyethylene glycol (PEG) causes a non-specific aggregation of protoplasts and also causes protoplasts to shrink by withdrawing water. This serves to bring cytoplasmic membranes into close contact and, of itself, is adequate to bring about a small number of protoplast fusions. For efficient protoplast fusion, PEG treatment must be accompanied by calcium-ion treatment that results in local disturbances in the membrane, leading to fusion. The optimal molecular mass and concentration of PEG should be investigated as a preliminary step to any previously untried fusion. For plant and fungal protoplasts, it is often desirable to use re-crystallized PEG, as degradation products or other impurities may interfere with fusion yields (5, 6). Degradation products formed during the autoclaving of PEG are similarly toxic and, when good yields are critical, it is necessary to sterilize by filtration. Calcium chloride is the usual source of calcium ions, but improved fusion yields have been achieved for *Candida albicans* and *Saccharomyces cerevisiae* (7) when calcium acetate or propionate is used.

The alternative to the chemical induction of protoplast fusion is *electrofusion*. This was first achieved by Senda et al (8), who used a DC pulse to fuse pairs of plant protoplasts micro-manipulated into contact. The process has been refined and applied to bacterial (9), fungal (10), and plant protoplasts (11). Contact is established between protoplasts by exposing them to an alternating, weakly inhomogeneous field that leads to polarization of the cells and their migration to regions of higher field density in a process known as *dielectrophoresis*. Because they form dipoles, the protoplasts line up on meeting, parallel to the field lines, in what have been termed “pearl chains.” The AC field also causes the lateral diffusion of membrane proteins and formation of protein-free regions at points of close contact (12). The fusion of adjacent cells is induced by a DC-field pulse of short duration and high intensity. This leads to electrical breakdown at the area of membrane contact causing pore formation and cytoplasmic continuity. The technique has a number of advantages over chemical fusion. It permits controlled observation of the fusion, limits the possibility for multiple cell fusion, and results in fusion products that have not been subjected to prolonged toxic effects of the fusogenic agent. The major disadvantage is that relatively small numbers of protoplasts may be processed compared with chemical fusion. It is possible to combine the PEG-induced aggregation of

protoplasts with electrical-field-induced fusion.

Whichever technique is chosen for protoplast fusion, it is necessary to have some means of selecting the fused cells (fusants) and a suitable regimen for the regeneration of a cell wall and reversion to dividing cells. Complementary auxotrophic mutations in the parent protoplasts are most commonly used to select fusants that should grow on a minimal medium. Other selection techniques may be based on (1) the complementation of resistances to antibiotics or other inhibitors (2), in the case of yeasts, the restoration of respiratory competence in a respiration-deficient mutant and (3), in plants, the complementation of an albino mutation. The use of mutant parent strains may not always be desirable, particularly in fusions aimed at producing new strains with agronomically or biotechnologically improved characteristics, where the presence of such mutations, even in a complemented situation, may be deleterious. In such situations, it may be advantageous to fuse protoplasts that have been differentially labelled, eg, with contrasting fluorescence signals.

Fusant cells with hybrid fluorescence may be selected either visually or using a *fluorescence-activated cell sorter* (FACS) (see [Flow Cytometry](#)). It is not always necessary that partners in a fusion should be individually viable. Protoplasts from cells killed by chemical, irradiation, or heat treatments may be fused with viable protoplasts. In some situations, this may improve the chances of obtaining hybrids; in others, it can be used to bias the genetic input of the parents. The heat killing of one partner in a bacterial fusion may inactivate restriction systems that could interfere with the success of the fusion. Regeneration and reversion conditions for hybrids may need to be investigated, as it is not necessary that hybrid protoplasts will behave similarly to the parents. It is not always possible to regenerate hybrids under direct conditions of selection. When bacterial protoplasts are fused, a wider variety of hybrid types can be isolated if they are initially grown on rich medium. Plant hybrids may need to be grown alongside nurse cells that provide growth factors for the protoplasts. The general requirements for the regeneration of protoplasts are discussed in [Protoplasts](#).

Natural genetic systems in bacteria have three characteristics (13): 1) Genetic transfer is unidirectional; 2) DNA alone, and no cytoplasm, is transferred; and 3) It is an unusual event for the full genome to be transferred. Protoplast fusion seemed to offer great potential for studying genetic recombination between whole genomes and interactions between the nucleus and cytoplasm. Much work has been done with the *Bacillus* species, as these are easy to convert to protoplasts. Attempts at genetic mapping by recombination analysis in protoplast fusion hybrids, however, have not always been very successful. When protoplasts of complementary auxotrophic strains of *B. subtilis* were fused and fusants selected on a minimal medium, complemented diploids were extremely hard to find (14). When fusants were regenerated on a complete medium and well-isolated colonies were examined, practically no prototrophs were found. A significant number of fusants that grew on minimal medium plus the nutritional requirements of either of the parents, but not on a minimal medium alone, were discovered. These were referred to as *noncomplementary diploids* (NCDs). The physical presence of chromosomes from each parent in NCDs was demonstrated, and it was postulated that only one of the two chromosomes was being expressed. The silence of the chromosome is not the result of methylation but rather of a total shutdown of [transcription](#) (15). The genomes of the NCDs can undergo recombination, and this is the probable origin of the small numbers of prototrophic fusants obtained. Gabor and Hotchkiss (16) examined the recombination patterns and found that recombination was frequent and multiple in some NCDs and reciprocal recombinants were frequently uncovered, often together within a single regenerated colony. Recombination occurred in all chromosomal regions, but the frequencies of recombination were considerably greater in the chromosomal replication origin and termination regions. The latter was attributed to structural characteristics at the origin and termination regions associated with membrane attachments. The genetic manipulation of organisms when the requisite genes are known and accessible is best performed by [transformation](#). Protoplast fusion for genetic manipulation is most useful when target genes are not known or a polygenic trait is to be transferred (17). It is also of value for the manipulation of organellar genetic systems.

The use of protoplast fusion for genetic manipulation has perhaps its greatest success, and potential



with plants and has involved intraspecific, interspecific, intergeneric, and intertribal fusions. Intraspecific fusions are generally performed to effect the pooling of desirable characteristics from varieties within a species. Attempts at interspecies (and higher-order) hybrid production have aimed to increase the gene pool available for the improvement of a crop or ornamental plant. Particular objectives have been the acquisition of disease and pest resistance, characteristics often found in wild species but lost from cultivated species or the breeding of new types of ornamental plants. With respect to crop plants, it is generally important that hybrids are fertile and they maintain the advantageous characteristics of the cultivated parent. Somatic hybrids often display low fertility, in particular for pollen, but this may be improved and desirable crop plant genes maintained by using the hybrid as the female parent in a back-cross to the cultivated species. The last few years have seen a rapid growth in the range of species for which protoplast formation, fusion, and regeneration are possible. Notably, techniques have been developed for woody species and members of the Poaceae (grasses), many of which are of economic importance and had previously proved recalcitrant, particularly in protoplast regeneration.

Whatever the outcome of a fusion process with respect to the contributing nuclei, cytoplasmic union is an inevitable outcome. When this occurs without nuclear fusion or the exchange of genetic information between the nuclei, the product is a *cybrid* (cytoplasmic hybrid). Mixing of the cytoplasm may be a slow process, but the time taken for the regeneration of protoplasts ensures fairly thorough cytoplasmic mixing before cell division occurs. In the yeasts *Saccharomyces cerevisiae* (18) and *Kluyveromyces lactis* (19), cybrid formation is the most common outcome of the fusion of two protoplasts. When nuclear fusion does not take place, loss of one of the two nuclei is likely to occur during the early divisions of the regenerating protoplasts. Fusion of a respiratory deficient, mitochondrial DNA mutant strain of *Candida utilis* with a respiratory competent *S. cerevisiae* gave rise to fusants having *C. utilis* nuclear and mitochondrial genomes, the latter transformed to respiratory competence (20). In addition, respiratory competence has been restored in a respiratory deficient strain of *S. cerevisiae* by fusion of protoplasts with mini-protoplast-containing functional mitochondria. This demonstrates the possibility of selectively introducing organelles into specific cells via protoplast fusion. In fusions involving protoplasts of respiratory competent strains, the mitochondrial genomes may not be retained equally in the hybrids. In a fusion of *K. lactis* and *K. fragilis*, it was demonstrated that the mitochondrial genome of the latter was always preferentially retained. Protoplast fusion may also be used in fungi to demonstrate the mitochondrial DNA location of a mutation (21). The transfer of dsRNA viruses between *Aspergillus* species has been achieved by protoplast fusion, although the strains involved were mating incompatible, demonstrating the potential of this technique to overcome mating barriers (22).

After fusion, the nuclei exist in a common cytoplasm for some time before [karyogamy](#), or the loss of one nucleus occurs. Stable [heterokaryons](#) are characteristic of a wide range of filamentous fungi but not of yeasts or plant cells. However, heterokaryosis can be maintained by the continued application of selection in fusants when karyogamy does not take place. In *Saccharomyces cerevisiae*, the incidence of karyogamy can be greatly increased if the nuclei of the parents are synchronized by treatment with either *a*- or *a*-mating pheromones before protoplast isolation (18, 23).

The fusion of haploid strains of *Saccharomyces cerevisiae* usually produces stable **diploids**. *Kluyveromyces lactis*, however, gives stable diploids only if haploids of similar mating type are fused; *ala* hybrids produced by protoplast fusion **sporulate** spontaneously. In some cases with filamentous fungi, protoplast fusion produces stable, vigorously growing heterokaryons. In a similar manner to a fungal parasexual cycle, karyogamy may occur, producing heterozygous diploids and recombinants on segregation. This is the case for *Aspergillus nidulans* (24), *Penicillium chrysogenum* (25), *Aspergillus oryzae* (26), and *Aspergillus niger* (27). In the case of *Paecilomyces fumosoroseus* (28) or *Aspergillus chrysogenum* (29), heterokaryons grow very slowly, and heterozygous diploids are rarely found. Recombinant haploids can usually be obtained, probably resulting from karyogamy followed by early haploidization. In a number of instances, heterokaryons stabilized by the apparent random loss of chromosomes from one parent and gave rise to a genome containing the chromosomes from one parent together with a few chromosomes from the other.

The products of interspecific fungal protoplast fusion are less predictable. As a general rule, interspecific heterokaryons formed from filamentous fungi grow badly. Sporulation gives spores characteristic of either parent, often with a predominance of one of them (30). When the selection of yeast protoplast fusants is based on auxotrophic [complementation](#), a common outcome is that these contain the full genome of one of the parents plus a small number (often one) of the chromosomes of the other parent (31-33). Aneuploids and diploids have been formed as a result of interspecific fusions of yeast protoplasts, but the fusants display the characteristics associated with the dominant contributor of DNA. It is not clear whether the retention of a single chromosome of one partner in a hybrid results from: 1) karyogamy followed by chromosome loss, or 2) from chromosome transfer between nuclei in a transitory heterokaryon.

As for plant protoplasts, the most significant potential for protoplast fusion is for the transfer of polygenic traits between species. When a characteristic to be transferred is determined by a small number of known genes, a direct molecular genetic approach to transfer is the most useful. When protoplast fusion is used, the fusant progeny are likely to have acquired an uncontrolled number of unwanted genes, as well as those being sought. Nevertheless, there have been many attempts to produce fusion hybrids with biotechnological value. Many of these have sought to construct strains in which the ability to metabolize novel carbon sources has been transferred from a poorly fermenting species to the efficient ethanol producer *Saccharomyces cerevisiae*. These include the fermentation of lactose, from *Kluyveromyces lactis* (34) and *Kluyveromyces fragilis* (35); of xylose, from *Pachysolen tannophilus* (36); and of cellulose hydrolysates, from *Zygosaccharomyces fermentati* (37). These studies have met with mixed success. Typically, fermentation levels in the hybrids are better than those of the parent strain that metabolized the substrate but still somewhat inferior to *S. cerevisiae* in terms of general fermentation efficiency. The stability of the hybrids is often an additional problem. Hybrids produced as a result of the fusion of protoplasts of the yeast *Pichia stipitis* stabilized during culturing as a result of ploidy reduction. Attempts have been made to improve the fermentation characteristics of yeasts by raising the [ploidy](#) level using protoplast fusion. In the case of xylose fermentation, raising ploidy levels increased rates of ethanol formation with *P. tannophilus* (38) and *Candida shehetae* (39) but not with *Pichia stipitis* (40). For filamentous fungi, improvements in citric acid synthesis by *Aspergillus niger* (27) and cephalosporin production by *Aspergillus chrysogenum* (29) have been achieved by the fusion of different strains. In many other recorded cases, however, protoplast fusion did not give rise to improved performance (30, 41). There is clearly some potential for the improvement of biotechnological performance using protoplast fusion with fungi, but the approach is so empirical that more rational strategies are preferable wherever possible.

## Bibliography

1. Q. F. Ahkong, J. I. Howell, J. A. Lucy, F. Safwat, M. R. Davey, and E. C. Cocking (1975) *Nature* **255**, 66–67.
2. J. F. Makins (1983) In *Protoplasts 1983* (L. Potrykus, C. T. Harms, A. Hinnen, R. Hatter, P. J. King, and R. D. Shillito, eds.), Berkhiiuser Verlag, Basel, Switzerland, pp. 197–207.
3. P. Van Solingen and J. Van der Plaats (1977) *J. Bacteriol.* **130**, 946–947.
4. K. Kavanagh and P. A. Whittaker (1990) *Biotechnol. Appl. Biochem.* **12**, 57–62.
5. P. K. Chand, M. R. Davey, J. B. Power, and E. C. Cocking (1988) *J. Plant Physiol.* **133**, 480–485.
6. K. Kavanagh, M. Ghannoum, I. Mansour, and P. A. Whittaker (1990) *Biotechnol. Tech.* **4**, 281–284.
7. K. Kavanagh, M. Walsh, and P. A. Whittaker (1991) *FEMS Microbiol. Lett.* **81**, 283–286.
8. M. Senda, J. Takeda, S. Abe, and T. Nakamura (1979) *Plant Cell Physiol.* **20**, 1441–1443.
9. H. J. Ruthe and J. Adler (1985) *Biochim. Biophys. Acta* **819**, 105–113.
10. R. Schnetter, U. Zimmermann, and C. C. Emeis (1984) *FEMS Microbiol. Lett.* **24**, 81–85.

11. U. Zimmermann and P. Scheurich (1981) *Planta* **151**, 26–32.
12. B. Hahn-Hagerdahl, K. Hosono, A. Zachrisson, and C. H. Born-man (1986) *Physiol. Plant* **67**, 359–364.
13. K. Fodor, K. Rostas, and L. Alfoldi (1980) In *Advances in Proto-plast Research* (L. Ferenczy, L. Farkas, and G. Lazare, eds.), Pergamon, Oxford, UK, pp. 19–28.
14. P. Schaeffer, B. Cami, and R. D. Hotchkiss (1976) *Proc. Natl. Acad. Sci. U.S.A.* **73**, 2151–2155.
15. N. Guillen, C. Sanchez-Rivas, and L. Hirschbein (1983) *Mol. Gen. Genet.* **191**, 81–85.
16. M. H. Gabor and R. D. Hotchkiss (1982) In *Genetic Exchange: A Celebration and a New Generation* (U. N. Streips, S. H. Goodgal, W. R. Guild, and G. A. Wilson, eds.), Marcel Dekker, New York, pp. 283–292.
17. S. Waara and K. Glimelius (1995) *Euphytica* **85**, 217–233.
18. B. P. G. Curran and B. L. A. Carter (1986) *Curr. Genet.* **10**, 943–945.
19. V. C. Bugeja and P. A. Whittaker (1986) *Microbiol. Lett.* **31**, 69–74.
20. M. S. Richard, M. R. G. Van Broock, and L. I. C. Figueroa (1987) *Curr. Microbiol.* **16**, 109–112.
21. A. J. Morgan, J. Heritage, and P. A. Whittaker (1978) *Microbiol. Lett.* **4**, 103–107.
22. A. D. van Diepingen, A. J. Deets, and R. F. Hoekstra (1998). *Fungal Genet. Biol.* **25**, 171–180.
23. V. C. Bugeja and P. A. Whittaker (1988) *FEMS Microbiol. Lett.* **51**, 101–104.
24. L. Ferenczy (1976) In *Cell Genetics in Higher Plants* (D. Dudits, G. L. Farkas, and P. Maliga, eds.), Akademiai Kaido, Budapest, Hungary, pp. 171–182.
25. J. Anne (1977) *Agricultura (Louvaine)* **25**, 1–117.
26. K. Uchida (1980) In *Molecular Breeding and Genetics of Applied Microorganisms* (K. Sakaguchi and M. Okanishi, eds.), Academic Press, Tokyo, pp. 103–105.
27. J. L. Azevedo and R. Bonatelli (1982) In *Overproduction of Microbial Products* (V. Krumphanzl, B. Sikyta, and Z. Vanek, eds.), Academic Press, London, pp. 439–450.
28. G. Riba (1978) *Entomophaga* **23**, 417–421.
29. P. P. Hamlyn and C. Ball (1979) In *Genetics of Industrial Microorganisms* (O. K. Sebek and A. I. Laskin, eds.), American Society for Microbiology, Washington, DC, pp. 185–191.
30. J. Anne (1983) In *Protoplasts 1983* (I. Potrykus, C. T. Harms, A. Hinnen, R. Hatter, P. J. King, and R. D. Shillito, eds.), Berkhiiuser Verlag, Basel, Switzerland, pp. 167–178.
31. J. Skala, Y. Luta, and Z. Kotylak (1988) *Curr. Genet.* **13**, 101–104.
32. H. Kobori, Y. Takata, and M. Osumi (1991) *J. Ferment. Bioeng.* **72**, 439–444.
33. T. Yamazaki and H. Nomura (1991) *J. Ferment. Bioeng.* **77**, 202–204.
34. M. Taya, H. Honda, and T. Kobayashi (1984) *Agric. Biol. Chem.* **48**, 2239–2243.
35. F. Fahranak, T. Seki, D. D. R. Ryu, and D. Ogrydziak (1986) *Appl. Environ. Microbiol.* **51**, 362–367.
36. H. Heluane, J. F. T. Spencer, D. Spencer, L. De Figueroa, and D. A. S. Callieri (1993) *Appl. Microbiol. Biotechnol.* **40**, 98–100.
37. A. Pina, I. L. Calderon, and T. Benitez (1986) *Appl. Environ. Microbiol.* **51**, 995–1003.
38. R. Maleszka, A. P. James, and H. Schneider (1983) *J. Gen. Microbiol.* **129**, 2495–2500.
39. E. Johannsen, L. Eagle, and G. Bredenham (1985) *Curr. Genet.* **9**, 315–319.
40. A. S. Gupthar (1987) *Curr. Genet.* **12**, 605–610.
41. C. Ball (1982) In *Overproduction of Microbial Products* (V. Krumphanzl, B. Sikyta, and Z. Vanek, eds.), Academic Press, London, pp. 515–534.

## Protoplasts

The cells of most **bacteria**, plants and **fungi** are surrounded by a more or less rigid cell wall. This is a multifunctional structure separated from the protoplasmic contents of the cell by the plasma membrane. Cell walls provide protection for cells against mechanical damage and allow cells to survive in a medium of lower osmotic potential than that of its protoplasm. They may also provide a locus for enzymes, ligands, and receptors, and, in some cases, they represent a barrier to the penetration of exogenous materials. These cell walls may be removed from most bacterial, plant, and fungal cells, leaving protoplasm surrounded by a plasma or cytoplasmic membrane. Such cells are referred to as *protoplasts*.

Protoplasts are osmotically sensitive, as in a suspension medium hypotonic to that of its protoplasm, water enters the cells by osmosis. This inflow is prevented in walled cells by the rigidity of the cell wall. It has been calculated (1) that the hydrostatic pressure of the cytoplasm of bacterial cells suspended in a dilute aqueous medium may be as much as 30 atm. In intact cells, this is contained by the inward pressure of the wall. In contrast, protoplasts transferred to a hypotonic medium swell and burst. It is consequently necessary for them to be suspended in an iso-osmotic or hypertonic medium for the protoplasts to survive. In iso-osmotic conditions, protoplasts assume a spherical shape. It is often difficult to be certain that a cell wall has been totally removed because cells become osmotically sensitive and adopt a spherical shape before the cell wall is completely removed, so neither of these can be regarded as an adequate indicator of protoplast formation. Cells from which the cell wall has been incompletely removed are referred to as spheroplasts. When it is not certain whether protoplasts or spheroplasts are being used, these are often referred to as osmotically fragile cells.

### 1. Bacteria

Bacterial protoplasts or spheroplasts can be prepared by the treatment of cells with the enzyme [lysozyme](#) (murein hydrolyase) that digests the peptidoglycan murein component of their cell walls. This was first achieved using **Gram-positive** *Bacillus* cells (2, 3). In the case of Gram-negative bacteria (4), the situation is significantly more complicated, as the typical wall (or envelope) structure comprises a relatively mobile inner or cytoplasmic membrane, a thin layer of murein and an outer membrane rich in lipopolysaccharide and stabilized by divalent cations. Both lipopolysaccharide and murein contribute to the rigidity of the envelope. The outer membrane prevents access of lysozyme to the murein, so it is usual to include [EDTA](#) (ethylene diamine tetraacetic acid) in the isolation medium to remove divalent cations. This destabilizes the outer membrane, causing the loss of significant quantities of lipopolysaccharide, and permits access for lysozyme to the murein layer. This technique invariably leaves significant quantities of outer membrane in place, so the usual products are spheroplasts rather than protoplasts.

Regeneration of bacterial protoplasts or spheroplasts to normal cells is variable in its success. Removal of the induction regime (lysozyme/EDTA) may permit regeneration of the cell wall if a suitable medium is provided. Success may depend on the growth phase of the cells used initially or the temperature at which the cells were grown and the regeneration is attempted (5). The nature of the osmotic stabilizer used to prevent the lysis of bacterial protoplasts can be critical, especially if the regeneration and reversion of protoplasts are required. Those most commonly used are inorganic salts (KCl, NaCl, MgSO<sub>4</sub>) or nonmetabolizable sugars or sugar derivatives (sucrose, mannitol, sorbitol), or combinations of these.

Most, and possibly all, species of bacteria are capable of converting to *L-forms*. These, named after the Lister Institute where they were first discovered (6), resemble normal protoplasts or spheroplasts in being spherical in shape and osmotically sensitive but differ in their replication capabilities.

**Electron microscopic** (EM) examination of L-form cells shows that they can be either devoid of any obvious cell wall material (protoplast-type L-forms), or they may contain some cell wall material (spheroplast-type L-forms). Large-scale conversion of *Bacillus subtilis* and other species of bacilli was first achieved by the removal of cell walls with lysozyme and plating the resultant protoplasts (spheroplasts) onto hypertonic soft agar medium (7). L-forms are most effectively induced by treatment with inhibitors of cell wall biosynthesis, especially **penicillin** and other **b-lactam** antibiotics. It is not always clear what heritable change leads to the L-form transition. Although there have been reports of changes in DNA sequence associated with the transition (8, 9), the ways in which L-forms are induced and the abilities of most L-forms to revert fairly readily suggest that, in most cases, the transformation is likely to be epigenetic in character. It has been suggested that the transition from an unstable to a stable L-form in *Proteus mirabilis* may be the result of a genetic change (4).

L-forms may be stable or unstable. The latter readily revert to the normal bacterial growth pattern and must be maintained in the presence of penicillin; the former are more difficult to revert, and, in some cases, reversion has not been possible. They grow by a budding process and shed envelope components as vesicular fragments. Murein is present in normal quantities and fully cross-linked but lacks its normal organization. Stable L-forms are usually devoid of an outer membrane and murein layer (10, 11). The cytoplasmic membrane of stable L-forms of *Proteus mirabilis* has been shown (4) to be modified from its normal composition, containing significant quantities of lipopolysaccharide and with a shift toward shorter-length fatty acids. L-form bacteria exhibit several properties that distinguish them from normal cells, such as resistance to **bacteriophage** in protoplast but not in spheroplast types, resistance to **antibiotics**, whose primary mechanism is to inhibit cell wall assembly, and increased sensitivity to antibiotics that act intracellularly. Early work suggested that L-forms of pathogenic bacteria were not pathogenic unless they showed a significant tendency to revert spontaneously. More recently (12), L-forms of some pathogens have been shown to retain pathogenicity. Some other pathogens show a tendency to adopt the L-form when patients are treated with penicillin.

## 2. Filamentous Fungi and Yeasts

Protoplasts may be formed from both filamentous fungi and yeasts. The usual method for making fungal protoplasts is to treat a culture with an enzyme cocktail that digests the major polysaccharide components of the cell wall:  $\beta$ -glucan, mannan, and chitin. Digestion of the cell wall was first achieved for *Saccharomyces cerevisiae* by Giaja (13) using the gastric juices of the snail *Helix pomatia* as the digesting enzyme mixture, although no protoplasts were formed in the absence of an osmotic stabilizer. These were obtained by Eddy and Williamson (14) also using snail gut juice (*Suc d'Helix pomatia*). Emerson and Emerson (15) produced the first protoplasts in filamentous fungi. There are now several commercially available sources of snail gut juice ( $\beta$ -glucuronidase, glucosylase, helicase), and these are known to be complex mixes of glucanases, mannanases, lipases, and proteinases. Subsequently, other enzyme preparations have become available—some from bacteria, principally *Bacillus* and *Streptomyces* spp., others from fungi. The most widely used fungal enzyme preparations are Novozyme 234 (from *Trichoderma harzianum*) and Fungelase (from *Trichoderma viride*), which contain chitinase as well as the other enzymes mentioned.

Enzymic release of protoplasts from yeasts follows one of two pathways. Either the cell wall is gradually eroded until a protoplast (or spheroplast) remains or, when a suitably sized portion of cell membrane is exposed, the protoplast may be extruded through what is left of the cell wall. In yeasts this often occurs through a bud scar. In the case of filamentous fungi, protoplasts are often smaller than individual fungal cells, particularly in the case of coenocytic fungi, and are variable in size and nuclear content. The growing tips of hyphae are generally more susceptible to enzymic digestion, so protoplasts tend to be released from these regions. In some species, however, the release of

protoplasts may occur at any position on the hyphal surface. Because of the method of release, the osmotically fragile products of filamentous fungi are usually protoplasts rather than spheroplasts. The rate at which fungal protoplasts are formed may be limited by the accessibility of digesting enzymes to the inner layers of the cell walls. In particular, the access of glucanases to the inner glucan component may be restricted by the outer mannoprotein layer. The inclusion of disulfide-reducing agents (e.g., *b*-mercaptoethanol, dithiothreitol) may open up the mannoprotein matrix and enhance the rate of protoplast release. This is most important for stationary phase cells, from which it is often quite difficult to release protoplasts. For some species, however, the release of protoplasts is not enhanced by the presence of reducing compounds, and their presence may be deleterious insofar as they may interfere with the reversion of the protoplasts.

Sorbitol or mannitol are mostly preferred as osmotic stabilizers for yeasts and inorganic salts for filamentous fungi. Optimal stabilizers and concentrations of these differ from species to species, from strain to strain, and with growth conditions. For any new species or strain, it is usually necessary to establish optimal conditions for protoplast formation and maintenance. The minimum condition is that the osmolality of the suspending medium be at least isotonic with that of the cytoplasm. Hypotonicity results in lysis; hypertonicity may be tolerated up to a point but can result in a reduction of the rate of protoplast release (16, 17). There is a cell wall-less mutant strain of *Neurospora crassa* (18) that grows on the surface of osmotically stabilized agar as a cluster of multinucleate protoplasts. This mutant, to an extent similar to bacterial L-forms, has been used as a source of protoplasts whose cell membranes have not been exposed to lytic enzymes.

Protoplasts from some fungal species (eg, *Schizosaccharomyces pombe*, *Nadsonia elongata*) can regenerate a cell wall in liquid suspension (19). In many species, however, regeneration requires embedding protoplasts in a solid or semisolid medium (20, 21). This is necessary to prevent the diffusion and loss of cell wall materials secreted from the protoplast. In some species, it has proved difficult to regenerate dividing cells from true protoplasts but relatively easy to regenerate from spheroplasts (22). The first polymer to appear on the surface of a regenerating protoplast of the yeast *Candida albicans* appears to be chitin, which is then followed by *b* (1, 3) and *b* (1, 6) glucan. Once these polymers have been deposited on the protoplast surface, mannans and mannoproteins are secreted and, if restrained by the presence of a physical barrier, aggregate to form a nascent cell wall. In the case of yeasts, the initial regenerated cells are spherical because the cell wall is assembled on the surface of a spherical protoplast. The typical shape of the yeast cell is reinstated only in cells arising during subsequent budding divisions. Protoplasts of filamentous fungi may also regenerate a spherical cell that subsequently produces either a germ tube or pseudohyphal structure before the establishment of true hyphae (23).

Once the cell wall has assembled around the protoplast, the process of reversion to the normal cellular state can take place. While in the protoplast state, nuclear division (karyokinesis) can proceed, but cytoplasmic division (cytokinesis) is impossible. Consequently, many protoplasts may be multi-nucleate by the time the cell wall has regenerated. During the first few generations, the “regenerated” cell reverts to the “normal” state, that is, the number of nuclei per cell is eventually reduced to one, the cell wall changes from one having a high amount of chitin to one with a relatively low amount and cellular metabolism re-commences.

### 3. Plants

Plant protoplasts have been prepared from plant types ranging from algae to flowering plants and from a variety of tissues. Many protoplasts are totipotent and, given suitable conditions, will regenerate cell walls, initiate and sustain cell division, and regenerate complete plants. The source of plant material for protoplast production is very important, particularly if good regeneration is required. Suspension cultures have been widely used, early exponential phase cultures giving the best yields (24). Sequential subculturing, however, leads to diminished protoplast yields, as well as to the accumulation of mutations and loss of totipotency (25). If regeneration is required, protoplasts isolated from leaf mesophyll cells of living plants grown under strictly controlled environmental

conditions are usually most consistent. Shoot culture under axenic conditions is finding increasing usage as a source of mesophyll and shoot apical protoplasts (26, 27). This is particularly useful for woody plants. Other sources of protoplasts include roots, petals, cotyledons, somatic embryos, and pollen.

The first-ever protoplast isolation was completed by Klercker (28), who cut onion scale epidermis using a sharp knife under 1 M sucrose. A small number of protoplasts were released from cells whose cell walls had been cut without damaging the protoplasts. Enzymic isolation of plant protoplasts in quantity was later achieved by Cocking (29). In intact tissues, plant cell walls are surrounded by a layer of pectin that must be removed before a tissue may be dis-aggregated. The major component of plant cell walls is cellulose; a minor component, particularly in aging cells, is hemicellulose. Various enzyme products are available for the removal of plant cell walls. When purified enzymes are used, a mixture of pectinase, cellulase, and hemicellulase may be required, but cruder preparations often contain sufficient amounts of all the enzymes necessary for cell wall removal. An empirical approach may be necessary in any new investigation, as preparations that give good protoplast yields in one system may be damaging to the protoplasts in another. Fowke and Cutler (30) give detailed protocols for the isolation of protoplasts from soybean suspension cultures, pea leaves, somatic embryos, and the filamentous green alga *Ulothrix*. Protoplast yields may be enhanced by adding, in addition to the necessary enzymes and osmotic stabilizer (usually a sugar alcohol): 1) Tween 80 to encourage protoplast release (31), 2) potassium dextran sulfate (32) or bovine serum albumin (33) to improve protoplast stability, or 3) antibiotics to prevent the growth of bacteria present from the source material (34). Protoplasts may be purified from residual plant material by flotation on either sucrose or Ficoll gradients (30). Healthy protoplasts appear spherical, and cytoplasmic streaming is apparent. Viable protoplasts, on treatment with fluorescein diacetate, produce **fluorescence** inside the plasma membrane.

The regeneration of plant protoplasts was first reported by Takebe et al. (35) for tobacco and has now been described for well over 300 species from at least 50 families. The process is very complex but can be considered to take place in three stages: 1) preparation for cell division involves the synthesis of a cell wall and the onset of [DNA replication](#) and **organelle** assembly; 2) cell division results in the establishment of a callus (undifferentiated cell mass); and 3) morphogenesis comprises differentiation to produce a somatic embryo. Establishment of cell division of protoplasts has been achieved by embedding in agar medium or overlaying a liquid protoplast suspension onto an agar surface. The accumulation of toxic compounds in the region of the regenerating cells may interfere with their recovery and, for this reason, it may be useful to include filter paper at the liquid agar interface or activated charcoal in the agar (34). Other problems stem from the presence of toxic materials in some agars or the temperature shock experienced by protoplasts on suspension in molten agar. Success has been achieved in many cases by the use of agarose or alginate as low-temperature gelling agents.

The nutritional requirements for protoplast culture are generally similar to those required for plant cell suspension cultures and are complex. Plant protoplasts are remarkably variable in their nutritional demands, and slight variations in the quality of the protoplast source or of medium constituents have given rise to difficulties in reproducing results between laboratories. The nature of the nitrogen source and the levels of growth regulators provided are particularly critical.

#### 4. Applications

Probably the most important use of protoplasts, and the one that gave the impetus to the rapid developments in protoplast technology during the last 20 years, is genetic manipulation of species using [protoplast fusion](#). However, protoplasts have also found many other uses. One of the most important has been the isolation of cell organelles undamaged by the mechanical fragmentation required for walled cells. This approach has been used for the preparation of intact [nuclei](#), [mitochondria](#), [chloroplasts](#), and vacuoles, as well as for cytoplasmic and tonoplast membranes. Protoplasts have also been employed for the study of cytoplasmic membrane permeability and

transport properties, for the investigation of the biochemistry of cell wall synthesis, and to assist in the introduction of macromolecules into cells.

## Bibliography

1. R. E. Marquis and T. R. Corner (1976) In *Microbial and Plant Protoplasts* (J. F. Peberdy, A. H. Rose, H. J. Rogers, and E. C. Cocking, eds.), Academic Press, London, UK, pp. 1–22.
2. J. Tomcsik and S. Guex-Holzer (1952) *Schweiz. Z. Allg. Path. Bakt.* **15**, 517–525.
3. C. Weibull (1953) *J. Bacteriol.* **66**, 688–695.
4. H. H. Martin (1983) In *Protoplasts 1983* (L. Potrykus, C. T. Harms, A. Hinnen, R. Hiitter, P. J. King, and R. D. Shillito, eds.), Berkhliuser Verlag, Basel, Switzerland, pp. 213–225.
5. R. H. Baltz and S. Matsushima (1981) *J. Gen. Microbiol.* **127**, 137–146.
6. E. Klieneberger (1935) *J. Pathol. Bacteriol.* **40**, 93–105.
7. O. E. Landman and S. Halle (1963) *J. Mol. Biol.* **7**, 721–738.
8. B. H. Hoyer and J. R. King (1969) *J. Bacteriol.* **97**, 1516–1517.
9. P. B. Wyrick, M. McConnell, and H. J. Rogers (1973) *Nature* **244**, 505–507.
10. C. Weibull (1965) *J. Bacteriol.* **90**, 1467–1480.
11. J. Gumpert, E. Schuhmann, and U. Taubeneck (1971) *Z. Allg. Mikrobiol.* **11**, 19–33.
12. W. N. Pachas (1986) In *The Bacterial L-Forms* (S. Madoff, ed.), Marcel Dekker, New York, NY pp. 287–318.
13. J. Giaja (1914) *C. R. Soc. Biol.* **77**, 2–4.
14. A. A. Eddy and D. H. Williamson (1957) *Nature* **179**, 1252–1253.
15. S. Emerson and M. R. Emerson (1958) *Proc. Natl. Acad. Sci. USA* **44**, 669–671.
16. H. Eyssen (1977) *Agricultura (Louvain)* **25**, 21–44.
17. K. Kavanagh and P. A. Whittaker (1991) *Biomed. Lett.* **5**, 313–316.
18. S. Emerson (1963) *Genetica* **34**, 62–182.
19. O. Necas, A. Svoboda, and M. Havelkova (1968) *Folia Biol. Prague* **14**, 80–85.
20. O. Necas (1961) *Nature* **184**, 1664–1665.
21. A. Svoboda (1966) *Exp. Cell Res.* **44**, 640–642.
22. K. Kavanagh and P. A. Whittaker (1990) *Biotechnol. Appl. Biochem.* **12**, 57–62.
23. J. F. Peberdy (1979) *Ann. Rev. Microbiol.* **33**, 21–39.
24. H. Uchimiya and T. Murashige (1974) *Plant Physiol.* **54**, 936–944.
25. M. Bayliss (1980) *Intl. Rev. Cytol. Suppl.* **11A**, 113–144.
26. H. Binding, R. Nehls, R. Kock, J. Finger, and G. Mordhorst (1981) *Z. Pflanzenphysiol.* **101**, 119–130.
27. M. A. L. Smith and B. H. McCown (1982) *Plant Sci. Lett.* **28**, 149–156.
28. J. A. Klercker (1882) *Ofvers Vetensk. Akad. Forh. Stockholm* **49**, 463–471.
29. E. C. Cocking (1960) *Nature* **187**, 962–963.
30. L. C. Fowke and A. J. Cutler (1994) In *Plant Cell Biology: A Practical Approach* (N. Harris and K. J. Oparka, eds.), IRL Press, Oxford, UK, pp. 177–197.
31. H. Lang and H. W. Kohlenbach (1982) *Planta Med.* **46**, 78–81.
32. K. C. Sink and R. P. Niedz (1982) In *Plant Tissue Culture 1982* (A. Fujiwara, ed.), Jpn. Assoc. Plant Tissue Culture, Tokyo, Japan, pp. 583–584.
33. K. L. Lai and L. F. Liu (1982) *Japan. J. Crop Sci.* **51**, 70–74.
34. M. R. Davey (1983) In *Protoplasts 1983* (I. Potrykus, C. T. Harms, A. Hinnen, R. Hutter, P. J. King, and R. D. Shillito, eds.), Berkhliuser Verlag, Basel, Switzerland, pp. 19–22.
35. I. Takebe, G. Labib, and G. Melchers (1971) *Naturwissenschaften* **58**, 318–320.



## Suggestions for Further Reading

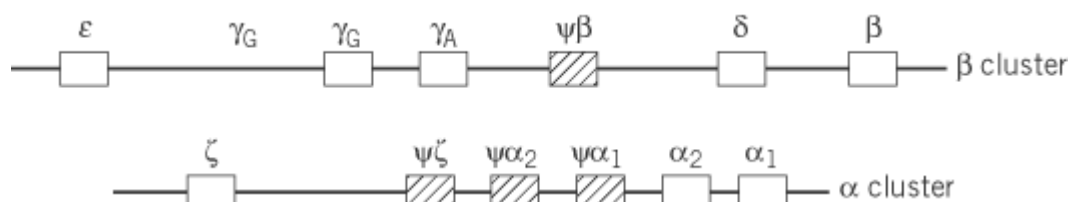
36. K. Kavanagh and P. A. Whittaker (1996) Application of protoplast fusion to the non-conventional yeast. *Enzyme Microbiol. Technol.* **18**, 45–51.
37. P. Matsushima and R. H. Baltz (1986) Protoplast fusion. In *Manual of Industrial Microbiology and Biotechnology* (A. L. Demain and N. A. Solomon, eds.), American Society for Microbiology, Washington, DC, pp. 170–183.
38. J. B. Power and J. V. Chapman (1985) Isolation, culture and genetic manipulation of plant protoplasts. In *Plant Cell Culture: A Practical Approach* (R. A. Dixon, ed.), IRL Press, Oxford, UK, pp. 37–66.

## Pseudogenes

Pseudogenes are functionless **gene** variants that are present as a result of ancient historical accident. Two general types of origin of pseudogene are postulated: (1) duplication within tandemly repetitive gene clusters, most easily envisaged as due to **unequal crossing over**; and (2) **reverse transcription** of **messenger RNA** followed by reinsertion of the **DNA** into the **chromosomes**. In both cases the sequence of the newly generated gene copy can decay by mutation without adverse effects on fitness. The new tandem repeat of (1) will be functionally redundant, and the reinserted reverse transcript of (2) will generally be without function from the start because it will be separated from its **promoter** sequence.

Pseudogenes that apparently arose as additional, redundant, tandem repeats were first identified in the two globin **gene clusters** of mammals (Fig. 1). **Globin**-encoding genes are expressed at different stages of development to provide the polypeptide chains for the range of globin tetramers:  $\zeta_2 e_2$ ,  $\zeta_2 g_2$ ,  $a_2 e_2$ , and  $a_2 g_2$  in the embryo;  $a_2 e_2$  in the fetus; and  $a_2 d_2$  and  $a_2 b_2$  in the adult. In mammals generally, the a and a-like genes are present in one cluster and the b and b-like genes in another on different chromosomes. Figure 1 shows the arrangement in humans. Pseudogenes are present in both human gene clusters and are designated by the prefix “y.” Almost exactly the same arrangement is found in chimpanzee, gorilla, orangutan, and baboon, so the pseudogenes must have originated more than 100 million years ago. During the intervening time they have presumably been free to decay by mutation. Their present inactivity is due to multiple mutations, some of which would have been inactivating individually (1).

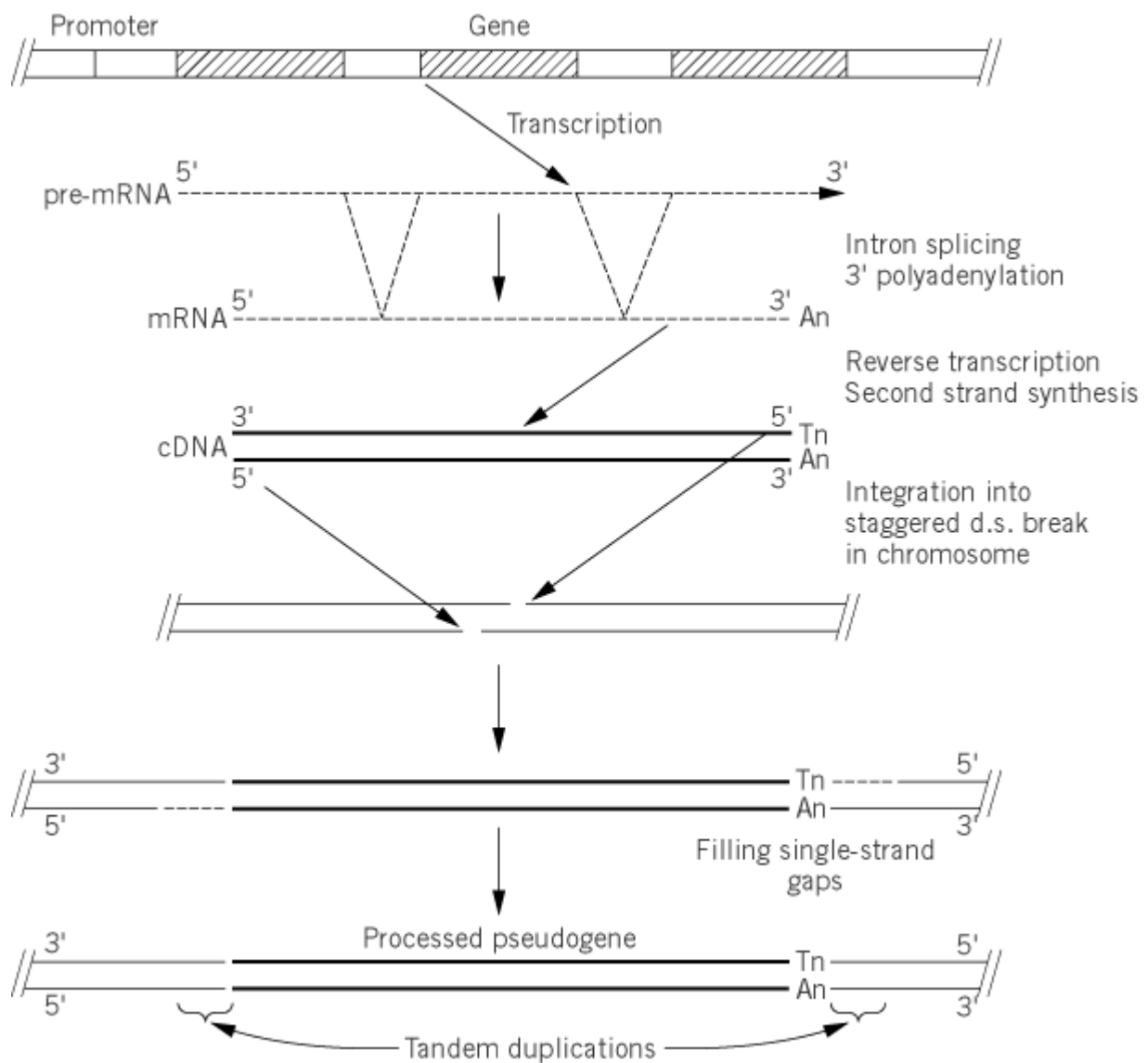
**Figure 1.** The arrangement of human globin genes. Pseudogenes are crosshatched and given the prefix y. They have undergone many disabling mutations, but retain recognizable similarity to their functional progenitors, including intron sequences.



More remotely related mammals, for example, rabbit and mouse, have differently arranged a and b clusters, but these also include pseudogenes. This suggests that any tandemly repetitive arrangement of genes should be subject to the same generation of redundancy and decay. Indeed, several examples in support of this prediction are known, for example, in [immunoglobulin](#) and [histocompatibility](#) gene clusters and the repetitive 5 S **ribosomal** RNA genes.

Pseudogenes originating through reverse transcription of RNA differ most obviously from those arising from genomic DNA duplication in that they have lost any **introns** that may have been present in their progenitors. For this reason, they are often called processed pseudogenes. They may also show their mRNA origin through the presence of sequences of T residues at the 5' end of their ancestrally transcribed sequence and an opposite sequence of A's at the 3'-end of the complementary strand, which is presumably a relic of the [poly-A](#) sequence at the 3'-terminus of mRNA. Processed pseudogenes may also be flanked by short [direct repeats](#), a common feature of sequences that have been inserted into the genome in new positions (see [Transposon](#), [Retrotransposons](#)) (Fig. 2).

**Figure 2.** The presumed mode of origin of a processed pseudogene by reverse transcription of mRNA and insertion of double-stranded cDNA into a double-strand break at some apparently arbitrary chromosomal locus. The flanking tandem repeats, which are present often but not always, are usually 10 to 20 base pairs long. They occur when the double-strand break is staggered, leaving single-strand gaps to be filled in by repair synthesis. The poly-A/T tracts ( $A_n/T_n$ ) are on the order of 10 bp.



Because their sequences arise solely from RNA transcripts, processed pseudogenes lack any nontranscribed promoter sequences. Thus, if they are derived from messenger RNA, the product of **RNA polymerase II**, which depends on upstream promoters, they should not generally be capable of further transcription. Processed pseudogenes in this category should all be due to independent events. However, analysis of a total of 14 pseudogenes derived from the human gene that encodes argininosuccinate synthetase (2, 3) has shown that certain pairs have closely similar patterns of mutation, suggesting further replication following the initial reverse transcription. It remains true that processed pseudogene derivatives of protein-encoding genes, though collectively fairly numerous in mammalian genomes and a source of confusion in hunts for real genes, do not **amplify** individually to high copy number.

In contrast, the relatively small RNA molecules transcribed by RNA polymerase III can develop superabundant pseudogene clones. Examples are the transfer RNA-related elements in several mammalian genomes and, most famously, the **Alu elements** of humans, thought to be derived originally from the gene for the RNA component of the 7SL signal particle. This capacity for indefinite proliferation is attributed to the fact that polymerase III promoters are downstream of the transcription start point and hence within the transcribed region, so that reverse transcripts are capable of further rounds of transcription.

## Bibliography

1. P. Jagadeeswaran et al. (1982) Cold Spring Harbor Symp. Quant. Biol. **47**, 1081–1082.
2. S. O. Freytag, H.-G. O. Bock, A. L. Beaudet, and W. E. O'Brien (1984) J. Biol. Chem. **259**, 3160–3166.
3. H. Nomiya et al. (1986) J. Mol. Biol. **192**, 221–233.

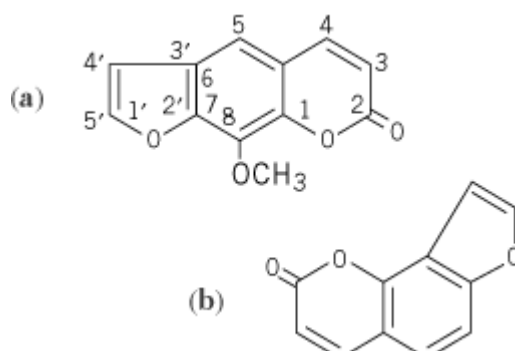
## Psoralen

The psoralen family of compounds was originally discovered as natural products with the ability to cause intense photosensitization reactions (exacerbated sunburn) in human skin. Under controlled conditions, they became a major tool in the dermatologic armory for the treatment of several skin disorders. Simultaneously, it was found that they could also be used to unravel the details of molecular contacts between DNA and proteins, as well as between DNA strands. Among the wide variety of compounds that can be photoactivated, the photochemistry of psoralen is probably the most studied and the best understood.

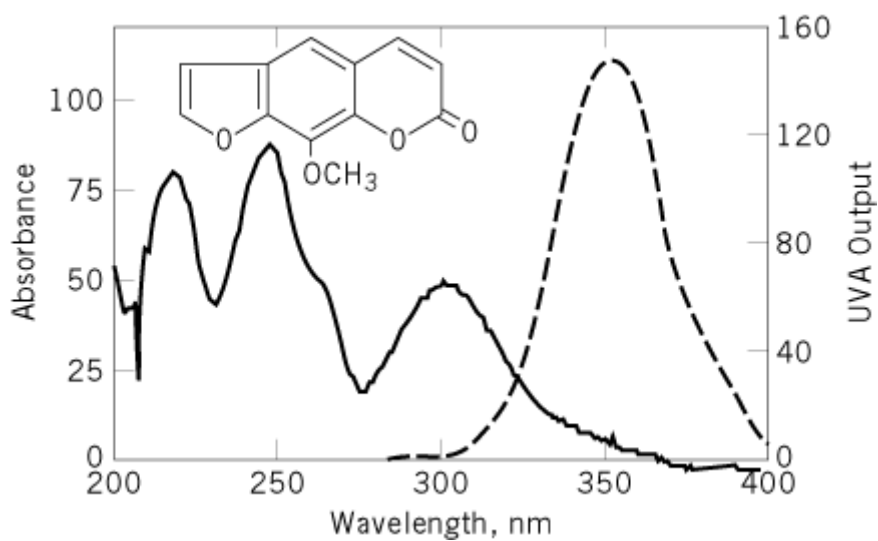
### 1. Properties

Psoralens are members of the furocoumarin family. They are tricyclic structures (Fig. 1) with an extended aromatic system that gives rise to strong ultraviolet **absorption** bands near 220, 250 and 300 nm (Fig. 2) (1). Also shown in Figure 2 is the output of the ultraviolet lamps typically used to activate psoralen. The spectroscopic properties of representative psoralens are listed in Table 1. After absorption of a photon by a molecule in the ground state, an electron can be promoted to an excited singlet state. It may then return to the ground state by the emission of a photon or by the release of energy in the form of heat (radiationless collisional deactivation) or light (**fluorescence** or phosphorescence). Excited state forms of the molecule can react in different ways: photo-addition, photo-dimerization, and photo-oxidation of nearby moieties (nucleic acids, proteins, or [membranes](#)). Alternatively, energy may be transferred to molecular oxygen, leading to highly reactive singlet oxygen ( $^1\text{O}_2$ ) (2). Although many molecules may absorb UV light, photochemistry rarely occurs. Furthermore, it has been shown that psoralen photochemistry in a solution can differ drastically from that within the confines of a cell.

**Figure 1.** The chemical structures of psoralen (a) and angelicin (b).



**Figure 2.** The UV absorption spectrum (solid line) for 8-methoxypsoralen (8-MOP; inset) and the typical output of UVA lamps used to photoactivate psoralens (dashed line).



**Table 1. Spectroscopic Properties of Psoralens(1)**

| Compound                 | Absorbance            |  | Fluorescence          |               |
|--------------------------|-----------------------|--|-----------------------|---------------|
|                          | $\lambda_{\max}$ (nm) | Extinction Coefficient ( $M^{-1}cm^{-1}$ ) | $\lambda_{\max}$ (nm) | Quantum Yield |
| Psoralen                 | 295                   | 10,600                                     | 450                   | 0.019         |
| 8-methoxy                | 303                   | 11,500                                     | 495                   | 0.013         |
| 5-methoxy                | 308                   | 14,500                                     | -                     | -             |
| 3-carbethoxy             | 320                   | 11,500                                     | 490                   | 0.025         |
| 4,5',8-Trimethylpsoralen | 293                   | 7,950                                      | 460                   |               |
| 4'-aminomethyl           | 298                   | 10,000                                     | 450                   |               |
| Angelicin                | 300                   | 9,350                                      |                       |               |
| 3-methyl                 | 300                   | 10,900                                     |                       |               |
| 5-methyl                 | 305                   | 11,410                                     |                       |               |
| 4,5'-dimethyl            | 298                   | 9,350                                      |                       |               |
| 4,6,4'-trimethyl         | 298                   | 8,650                                      |                       |               |

## 2. Photoadduct Formation

## 2.1. Nucleic Acids

The best understood psoralen photochemistry is that which occurs when psoralen is intercalated between DNA base pairs. A 2 + 2 photocycloaddition product is formed between pyrimidine (primarily thymidine at 5'-TpA sites) when the intercalation complex absorbs UV radiation (3). Furan side monoadducts can absorb additional long-wavelength ultraviolet radiation (UVA, 320–400 nm) and, if located at a proper site, undergo a second photoaddition reaction, thereby forming an interstrand **crosslink**. The psoralen isomers (angelicins) are angular psoralens (Fig. 1) that have been studied because of their ability to form only mono-addition photoadducts on one DNA strand (4). Due to its angular structure, the angelicin 4',5'-monoadduct cannot form crosslinks.

## 2.2. Proteins, Lipids

Although DNA psoralen photochemistry has attracted the most attention, psoralens also react with proteins as well as other cell constituents (5). In studies of subcellular fractions of rat epidermis after treatment with 8-methoxypsoralen (8-MOP) and UVA light, it was found that 17% of the 8-MOP was bound to DNA, whereas a substantial amount was bound to proteins (57%) and lipids (26%) (6). Proteins photomodification is known to alter the ability of endopeptidases to recognize the usual sites of incision (7). Much greater doses of 8-MOP and UVA are required to affect proteins than doses necessary to damage DNA. Even low levels of protein modifications, however, could participate in the concerted cellular events that lead to clinical effects. Recently, the first psoralen–amino acid photoadduct was described (8).

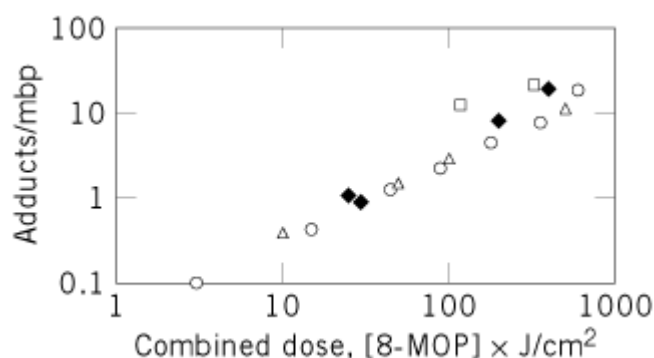
The ability of 8-MOP and trimethylpsoralen (TMP) to react in an oxygen-independent way with unsaturated **fatty acids** was first shown by Kittler et al. (9). Specht et al. (10, 11) and Caffieri et al. (12) independently isolated the lipid-psoralen adducts. **NMR** studies indicated that these adducts resulted from cyclo-additions of the 3,4-double bond of the psoralen to the central double bond of the fatty acid. Recently, Dall'Acqua et al. (13) suggested that these adducts have structures similar to **diacylglycerol** (DAG) and proposed that they could play a role in cell **signal transduction** events. In more recent studies, it has been observed that some photoactivated psoralens can crosslink DNA and proteins (14).

## 3. Impact on Cells

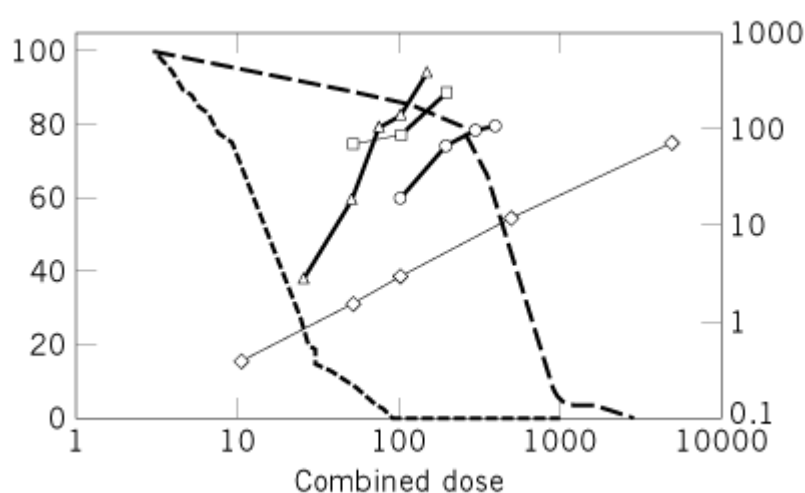
Over the nearly three decades that psoralen photochemotherapy has been employed by dermatologists, its efficacy has been explained on the basis of cytotoxicity, immune modulation, and/or **apoptosis**, somewhat depending on what has been in vogue. The initial observation of psoralen crosslinking led to the potent antiproliferative/cytotoxic paradigm that persisted until the evidence arose for gene induction and immune modulation.

The formation of psoralen–DNA adducts has been studied extensively in different cell types under *in vitro* conditions. [<sup>3</sup>H]8-MOP has been used to quantify the number and type of adducts formed, using **HPLC** and radioactivity incorporation analysis (15). In Figure 3, the number of photoadducts induced in a variety of cells, and treated *in vitro* is shown as a function of the combined dose of psoralen and light (8-MOP concentration in ng/ml multiplied by UVA dose in J/cm<sup>2</sup>). Adduct formation is directly proportional to the dose ( $r = 1.0$ ). The effect of 8-MOP and UVA in response to mitogenic stimulation by phytohemagglutinin (PHA) and exclusion of trypan blue has been studied in human lymphocytes. The former response is abrogated at relatively low doses of 8-MOP and UVA, with a complete block at a combined dose greater than ~3050ng J/ml cm<sup>2</sup> (Fig. 4). Much greater doses are required to affect membrane integrity; hence, a distinctly different curve is seen for trypan blue exclusion (16). The antiproliferative effect of 8-MOP and UVA on murine keratinocytes was studied *in vitro* by Tokura et al. (17), who correlated the number of adducts with inhibition of [<sup>3</sup>H]thymidine incorporation after PHA stimulation. A minimum number of 0.9 adducts per million bases was required to inhibit DNA synthesis. Moreover, it was shown that the photoadduct removal rate in keratinocytes was dependent on the extracellular calcium concentration, being reduced in low-calcium media (17).

**Figure 3.** Dose dependence for 8-MOP photoadduct formation in cells treated with 8-MOP and UVA light. Adduct numbers per megabase pair (mbp) are plotted vs combined dose of 8-MOP (ng/ml) and UVA light ( $J/cm^2$ ).  $\Delta$ , human lymphocytes;  $\circ$ , murine keratinocytes;  $\square$ , bovine aorta smooth muscle cells;  $f$ , Jurkat cells.



**Figure 4.** Dose dependence for cellular responses to psoralen and light. linear fit of the photoadduct data from Figure 3. The biological effects of these adducts are most significant at the lower modest doses (1 to 100), over which the response to PHA stimulation (□) as measured by tritiated thymidine incorporation) is completely abrogated by a combined dose of 100. The impact on cell viability (as measured by trypan blue exclusion (heavy dash) measured three days after treatment) only occurs at much higher doses, reflecting the lower efficiency of 8-MOP photochemistry with proteins and lipids. Superimposed on these results are data (solid lines) illustrating the induction of class I MHC molecules on murine lymphoma cells, either total ( $f$ ), (48), or surface-bound ( $\circ$ ) ((49),(50)). Finally, the induction of apoptosis in human lymphocytes is shown ( $\Delta$ ) (22). Taken together, these data illustrate that any selected property is likely to have its own characteristic dose dependence. Furthermore, these different effects are likely to reflect the combined effects of 8-MOP/UVA at different cellular sites involving a range of biomolecules. The range of these dose-dependent effects may be represented by response to mitogen (perhaps best reflecting the impact of DNA photoadducts) and cell viability (reflecting the effects of photomodification of lipids and proteins). Regarding the latter, it is interesting to note that the impact on viability takes several days to develop. If viability is measured immediately following the 8-MOP/UVA treatment, there appears to be little difference between cell populations treated with different doses of 8-MOP and UVA.



Extensive DNA repair studies have been carried out in bacterial systems as well as in mammalian cells. Several repair mechanisms have been described, including excision repair (the most common), postreplication repair, and photoreactivation (18). (See [Photolyase/Photoreactivation](#).) It has often

been assumed that crosslinks would not be repaired, but an excision-recombination mechanism has been proposed to account for crosslink repair in bacteria and bacteriophage (19, 20).

Repair mechanisms in mammalian cells have not been described as completely, but it has been shown in various human and murine cells that DNA-psoralen adducts induced by 8-MOP and UVA can be removed (see Table 1). Freshly isolated human lymphocytes do not repair 8-MOP photoadducts when treated with 100 ng/ml 8-MOP and 1 J/cm<sup>2</sup> of light, but do remove photoadducts after treatment with only 10 ng/ml 8-MOP. In PHA-stimulated lymphocytes, relatively more adducts are formed, but more of these are repaired in 24 hours (16). Repair in murine keratinocytes seems to be more efficient than in lymphocytes, which may be due to the increased metabolic activity of keratinocytes.

#### 4. Apoptosis

Although modest doses of 8-MOP and UVA are antiproliferative, higher doses can be toxic to cells. Thus, depending on the dose employed, the resulting effects can range from cytostatic to apoptotic to necrotic. Marks et al. (21) demonstrated the induction of apoptosis in human lymphocytes by using 300 ng/ml 8-MOP and 10 J/cm<sup>2</sup> UVA light (21). This dose is much greater, however, than the therapeutic dose used in patients (100–200 ng/ml 8-MOP and 1–2 J/cm<sup>2</sup> UVA). More recently, Vowels et al. (22) described the induction of apoptosis with more therapeutically relevant doses of 8-MOP and UVA (see Fig. 3).

#### 5. PUVA Therapy for Psoriasis and Skin Cancer

If apoptosis fails, mutations may occur, which can lead to skin cancer. Thus, the therapeutic use of psoralens and UVA (PUVA) is not without drawbacks. PUVA therapy for psoriasis increases the incidence of squamous cell carcinoma (23). This is thought to be due to the formation of mutagenic psoralen photoadducts in the keratinocytes that are repaired incorrectly or not at all. Based on mutation frequencies in human fibroblasts treated with physiologic doses of 8-MOP and UVA, Burger and Simons (24) calculated the number of induced mutations in human *hprt* genes per phototherapeutic session ( $1.2 \times 10^{-5}$ ) and per 30 years of maintenance therapy ( $1.3 \times 10^{-2}$  mutants per cell). This is two orders of magnitude greater than the spontaneous mutation rate. Petersheim et al. observed chromosomal aberrations and sister chromatid exchanges in treated lymphocytes immediately after photopheresis (25). Although the damage was repaired in 72 hours, no molecular information was presented to demonstrate the accuracy of the repair processes. The implications of these data are not clear, but close monitoring for mutations occurring in treated cells would appear to be prudent.

It has been thought that crosslinks would primarily lead to the observed mutations, because the mutations occur at crosslinkable sites (5'TpA). Recent evidence indicates that this is probably not the case. Chiou et al. (26) studied the induction of mutations in the *hprt* gene in human fibroblasts. Using 8-MOP plus low doses of UVA (0.0060 J/cm<sup>2</sup>) to induce primarily monoadducts, mutations at 5'-TpA sites were observed. Fewer *hprt* mutations per 10<sup>6</sup> clonable cells were detected when a protocol favoring crosslinks (second irradiation of treated cells in the absence of any free psoralen) was employed. Because the numbers of adducts were not reported, drawing firm conclusions is difficult, but monoadducts might be at least as mutagenic as crosslinks.

Gunther et al. (27) studied the mutations induced by UVA and 8-MOP (producing both crosslinks and monoadducts) and those induced by UVA and 5-methylangelicin (5-MA; producing only monoadducts), related to adduct distribution. Mutagenesis was examined in a mouse fibroblast cell line carrying a recoverable, chromosomally integrated lambda phage shuttle vector, using the *supF* gene as a mutation reporter gene. Both psoralens generated predominantly T:A to A:T transition mutations with some T:A to G:C transversions. Most of the mutations occurred at TpA and ApT



sites, both of which are conducive to crosslink formation. However, 5-MA induces only monoadducts. In addition, it was shown by HPLC analysis that, under the 8-MOP/UVA conditions used, only 20% crosslinks and 80% monoadducts, were formed, strongly indicating that monoadducts, as well as crosslinks, are critical premutagenic lesions. Earlier, different studies concluded that crosslinks were more mutagenic than monoadducts, showing one mutation hot spot that correlated with crosslink formation (28, 29). These studies were performed by irradiating isolated vector DNA containing the *supF* gene with very high doses of 8-MOP and UVA. Gunther et al. (27) irradiated cells containing DNA with a chromosomally integrated *supF* transgene. Thus, the different results of these two kinds of studies could be attributed to the extra- and intracellular PUVA treatments used, which would be expected to lead to different adduct yields and distributions and, hence, to different mutation frequencies and spectra.

## 6. p53 Mutations in Murine vs Human Skin Cancer

Mice exclusively exposed to a regimen of psoralen and UVA radiation developed skin cancers over a 42-week period (30). Analysis of p53 mutations in the resulting tumors showed evidence of psoralen photochemistry at the expected 5'-TpA hot spots. In humans, however, the picture is significantly different. In two independent studies, a paucity of PUVA-type mutations at 5'TpA sites were found, and UVB-type mutations at dipyrimidine sites were more prevalent (31, 32). These data suggest that the incidence of skin cancers in PUVA patients may be traceable to some degree of immunosuppression in the PUVA-treated skin. In fact, these cancers may arise from previous solar damage. It may be possible to dissect the relative contribution of previous actinic damage and therapy-induced damage by analysis of p53 mutations in these patients.

Although it is common to think that DNA damage is mutagenic and cytotoxic, it has also been shown that it can lead to gene induction. In fact, many studies have illustrated how DNA damage can induce immune-modulating gene expression. An example of this was demonstrated by Bernstein et al. (33) in a transgenic murine system. Transgenic mice containing a chromosomally integrated reporter gene construct with the promoter for human elastin linked to the gene for chloramphenicol acetyl transferase (CAT) were used to demonstrate the dose-dependent induction of CAT activity (33). Events such as these link DNA damage to gene induction and can explain changes in cytokine production in phototreated cells (34, 35)

## 7. Therapeutic Applications

Psoralen photochemotherapy has been widely used for the treatment of the hyperproliferative skin disease psoriasis (36). 8-MOP is administered orally or topically and, after an appropriate time, the affected skin is exposed to UVA light. Clearing can be achieved after 8–12 weeks, with three treatments each week. The PUVA therapy can be gradually decreased, so that only maintenance therapy at monthly intervals is required. For vitiligo, a similar therapy can be followed, leading to repigmentation of the affected skin after 200–300 treatments (37)

It has long been thought that 8-MOP DNA photoadducts inhibit cell division and lead to the beneficial effects of PUVA therapy. However, immunosuppression due to PUVA also appears to play an important role (38). Although PUVA therapy has a high efficacy, it also exhibits some undesirable side effects, such as skin phototoxicity and increased risk of squamous cell cancer (see text above), both of which are thought to be due to the ability of psoralens to undergo photoreactions with DNA.

## 8. Other Applications

A completely different application of psoralens and UVA light is the sterilization of blood products, mainly platelet concentrates (38, 39). It has been shown that 8-MOP and 4'-aminomethyl-4,5',8-trimethylpsoralen (AMT) are able to reduce the titers of enveloped viruses by factors of  $10^5$ – $10^6$

while preserving platelet function. AMT has the advantage of being water soluble, and it is not harmful to the DNA-lacking platelets in the presence of scavengers of reactive oxygen species (40). AMT exerts some mutagenic effects in the dark, however, so it must be removed before transfusion. Other psoralens with undisclosed structures have reported to be very efficient for platelet concentrate sterilization (41), but the lack of published data makes evaluation of these claims very difficult.

## 9. A Tool for the Study of Macromolecular Structures

In studies employing adventitious interactions of proteins and nucleic acids, the ability of psoralen to be photoactivated *in situ* and to lock in structures has been used to identify molecular contact points between the interacting molecules (42). In other studies, psoralens have been chemically linked to oligonucleotides that recognize specific DNA sequences (including double helical structures) to deliver psoralen adducts either to specific nucleic acid base sequences or to proteins in contact with oligonucleotides (43). The internal structure of 16 S ribosomal RNA (rRNA) in the small ribosomal subunit has been studied using a site-directed crosslinking approach (44, 45). Psoralen derivatives were transferred into specific positions in the 16 S rRNA aided by complementary deoxyoligonucleotides, which were completely removed after the transfer reaction. Another photoreactive group, azidophenacyl, was then attached to the psoralen moiety before reconstitution of the 16 S rRNA into the 30 S subunit. Re-irradiation of the derivatized subunits with UV radiation produced intramolecular crosslinks in the 16 S rRNA, as evidenced by gel electrophoresis under denaturing conditions. Crosslinking points in the 16 S rRNA were determined after preparative separation of the products by gel electrophoresis, followed by reverse transcriptase analysis.

In another study by the same group, this technique was refined by pre-forming monoadducts in an oligonucleotide (46). Thus, monoadducts were targeted to specific nucleotides in the pre-mRNA, leading to the subsequent formation of intramolecular RNA-RNA crosslinks after another round of photoactivation. This decreased the number of crosslinked products in comparison to nonspecific psoralen crosslinking and confirmed the locations of previously determined free psoralen crosslinks in the human precursor mRNA. New crosslinks consistent with an alternative secondary structure were also observed, as were a small number of crosslinks representative of higher-order interactions. The use of psoralen isomers such as angelicins, and/or activation with other UV and visible wavelengths, could lead to more selective photochemical reactions that have the potential to reveal other binding modes and sites. The recent elucidation of a psoralen amino acid photoadduct is noteworthy in this context (see text above).

**Table 2. Repair of 8-MOP Photoadducts in Different Cell Types<sup>a</sup>**

| Cell Type                    | 8-MOP (ng/ml)/UVA<br>(J/cm <sup>2</sup> ) | Adducts/mbp | % repaired<br>(24 h) |
|------------------------------|---|-------------|----------------------|
| Human lymphocytes<br>(15)    |   |             |                      |
| Resting                      | 10/1                                      | 0.40        | 25                   |
| PHA stimulated               | 10/1                                      | 0.80        | 52                   |
| Murine keratinocytes<br>(17) | 15/1                                      | 0.48        | 54                   |
| Murine fibroblasts<br>(27)   | 1080/0.1                                  | 4.4         | ND                   |
|                              | 1080+60'400nm light                       | 2.9         | 66                   |

|                         |  |      |    |
|-------------------------|--|------|----|
| Murine lymphoma<br>(46) | 100/1                                    | 4.2  | 54 |
| Bovine SMC (47)         | 1000+12J/cm <sup>2</sup> 419 nm<br>light | 13.5 | 25 |

---

<sup>a</sup> Key: mbp, megabase pair; ND, not done; PHA, phytohemagglutinin; SMC, smooth muscle cells.

## 10. Historical Acknowledgments

Five figures stand out for their contributions to the field of psoralen photochemistry, photobiology, and photomedicine. A half century ago, Aaron Lerner and Tom Fitzpatrick pioneered the application of psoralen photochemotherapy for vitiligo. In the 1970s, this was expanded by John Parrish to include therapy for psoriasis. It was in this same era that John Hearst developed techniques that led to a fuller development of psoralen photochemistry. In the 1980s, largely through the efforts of Margaret Kripke, we came to a better understanding of the immune-modulating effects of photoactivated psoralen.

## Bibliography

1. F. P. Gasparro (1994) In *Extracorporeal Photochemotherapy: clinical aspects and the molecular basis for efficacy* (F. P. Gasparro, ed.), R. G. Landes Press, Georgetown, Texas (Medical Intelligence Unit), pp. 13–36.
2. N. J. Turro (1978) *Modern Molecular Photochemistry*, Benjamin Cummings, Menlo Park, 579–611.
3. D. Kanne, K. Straub, H. Rapoport and J. E. Hearst (1982) *Biochemistry* **21**, 861–871.
4. F. Bordin, F. Dall'Acqua, and A. Guiotto (1991) *Pharmac Ther* **52**, 331–363.
5. W. R. Midden (1988) In *Psoralen-DNA Photobiology*, vol. **II** (F. P. Gasparro, ed.), C. R. C. Press Inc., Boca Raton, Florida, 1–50.
6. G. M. J. Beijersbergen van Henegouwen et al. (1989) *J Photochem Photobiol B: Biol* **3**, 631–635.
7. I. M. Schmitt, S. Chimenti, and F. P. Gasparro (1995) *J Photochem Photobiol B: Biol* **27**, 101–107.
8. S. Sastry (1997) *Photochem Photobiol* **65**, 937–944.
9. L. Kittler and G. Lober (1984) *Stud Biophys* **101**, 69–72.
10. K. G. Specht et al. (1987) *Photochem Photobiol* **45 S**, 51.
11. K. G. Specht, L. Kittler, and W. R. Midden (1988) *Photochem Photobiol* **47**, 537–541.
12. S. Caffieri, G. Tamborrino, and F. Dall'Acqua (1987) *Med Biol Environ* **15**, 11–14.
13. F. Dall'Acqua et al. (1994) *Photochem Photobiol* **59**, 55S–56S.
14. F. Bordin et al. (1994) *J Photochem Photobiol B: Biol* **26**, 197–201.
15. G. Olack, P. Gattolin, and F. Gasparro (1993) *Photochem Photobiol* **57**, 941–949.
16. F. P. Gasparro et al. (1991) In *Photobiology—the Science and its Applications* (E. Riklis, ed.), Plenum Press, New York, 951–962.
17. Y. Tokura, R. L. Edelson, and F. P. Gasparro (1991) *J Invest Dermatol* **96**, 942–949.
18. C. A. Smith (1988) In *Psoralen-DNA Photobiology*, vol. **2** (F. P. Gasparro, ed.), C. R. C. Press Inc., Boca Raton, Florida, 87–116.
19. B. Van Houten et al. (1986) *Proc Natl Acad Sci USA* **83**, 8077–8081.
20. K. Dye and S. I. Ahmad (1995) *J Gen Virol* **76**, 723–726.

21. D. I. Marks and R. M. Fox (1991) *Biochem Cell Biol* **69**, 754–760.
22. B. R. Vowels, E. K. Yoo, and F. P. Gasparro (1996) *Photochem Photobiol* **63**, 572–576.
23. R. S. Stern et al. (1994) *Cancer* **73**, 2759–2764.
24. P. M. Burger and J. W. I. M. Simons (1979) *Mutat Res* **63**, 371–380.
25. U. M. Petersheim et al. (1991) *Arch Dermatol Res* **283**, 81–85.
26. C. C. Chiou and J. L. Yang (1995) *Carcinogenesis* **16**, 1357–1362.
27. E. J. Gunther et al. (1995) *Cancer Res* **55**, 1283–1288.
28. A. Bredberg and N. Nachmansson (1987) *Carcinogenesis* **8**, 1923–1927.
29. E. Sage and A. Bredberg (1991) *Mutat Res* **263**, 217–222.
30. A. J. Nataraj, H. S. Black, and H. N. Ananthswamy (1996) *Proc Natl Acad Sci USA* **93**, 7961–7965.
31. A. J. Nataraj, P. Wolf, L. Cerroni, and H. N. Ananthswamy (1997) *J Invest Dermatol* **109**, 238–243.
32. X. M. Wang et al. (1997) *Photochem Photobiol* **66**, 1294–1299.
33. E. F. Bernstein et al. (1996) *Photochem Photobiol* **64**, 369–374.
34. D. P. Fivenson, B. J. Nickoloff, and G. M. Saed (1994). *J Invest Dermatol* **102**, 585A.
35. G. M. Saed and D. P. Fivenson (1994) *Biochem Biophys Res Commun* **203**, 935–942.
36. J. A. Parrish et al. (1974) *N Engl J Med* **291**, 1207–1211.
37. H. Honigsmann et al. (1987) In *Dermatology in General Medicine* (T. B. Fitzpatrick, A. Z. Eisen, and K. Wolf, eds.), McGraw-Hill, New York, pp. 1728–1754.
38. M. L. Kripke (1984) *Natl Cancer Inst Monogr* **66**, 247–251.
39. E. Ben-Hur et al. (1996) *Transf Med Rev* **10**, 15–22.
40. H. Margolis-Nunno et al. (1994) *Transfusion* **34**, 802–810.
41. S. Wollowitz et al. (1995) *Photochem Photobiol* **61**, 90S.
42. P. L. Wollenzien, D. C. Youvan, and J. E. Hearst (1977) *Proc Natl Acad Sci USA* **75**, 1642–1646.
43. P. A. Havre et al. (1993) *Proc Natl Acad Sci USA* **90**, 7879–7883.
44. D. Mundus and P. Wollenzien (1997) *Nucleic Acids Symp Ser* **36**, 171–174.
45. J. Teare and P. Wollenzien (1990) *Nucleic Acids Res* **18**, 855–864.
46. A. C. E. Moor and F. P. Gasparro (1996) *Clinics in Dermatology* **14**, 353–365.
47. B. E. Sumpio, G. Li, L. I. Deckelbaum, and F. P. Gasparro (1994) *Circ Res* **75**, 208–213.
48. A. Felli et al. (1995) *J Invest Dermatol* **104**, 647A.
49. A. C. E. Moor et al. (1995) *J Photochem Photobiol B: Biol* **29**, 193–198.
50. I. M. Schmitt et al. (1995) *Tissue Antigens* **46**, 45–49.

### **Suggestions for Further Reading**

51. *Psoralen–DNA Photobiology*, vol. **I** and **II** (1988) (F. P. Gasparro, ed.), C. R. C. Uniscience Volume C. R. C. Press Boca Raton, Florida.
52. *Extracorporeal Photochemotherapy* (1994) (F. P. Gasparro, ed.), R. G. Landes, Georgetown, Texas (Medical Intelligence Unit).
53. *Photoimmunology* (1995) (J. Krutmann and C. A. Elmetts, eds.), Blackwell Science, Oxford, U. K.
54. *The Science of Photobiology* (1989) (K. C. Smith, ed.), Plenum Press, New York.
55. *CRC Handbook of Organic Photochemistry and Photobiology* (1994) (W. M. Horapool and P.-S. Song, eds.), CRC Press, Boca Raton, Florida,

## Puff, Chromosomal

Puffs represent sites in [polytene chromosomes](#) that are active in gene [transcription](#). The morphology of puffs is such that they contain chromosomal material that has become dramatically decondensed. The process of generating a puff has been likened to untwisting the strands of a rope (1). It is important to note, however, that transcription also occurs at nonpuffed bands in polytene chromosomes (2). Therefore it is probable that puffing represents a specialized adaptation of chromosomal structure or is a consequence of very high levels of transcriptional activity. Evidence for transcription at sites of puffing derives from the **immunofluorescent** staining of **RNA polymerase II** and the incorporation of **radioactive** uridine into RNA at puff sites, as revealed by [autoradiography](#) (3).

The appearance of puffs in polytene chromosomes is developmentally regulated. Analysis of puffing patterns has been very informative for studies of **gene expression**. In some organisms, like the midge *Chironomus*, puffs become enormous, thereby providing very valuable cytological insights into transcription and RNA processing (see [Balbiani Ring](#)). In *Drosophila melanogaster*, the study of puffing induced by the hormone [ecdysone](#) has been especially valuable. Ecdysone induces several puffs in salivary gland nuclei during the late third larval instar stage of development. These fall into two classes: (1) the early puffs, which appear within minutes of hormone addition and increase in size over a 1- to 4-hour period before diminishing; and (2) the late puffs, which appear after 3 hours, reach their maximal activity after 5 to 7 hours, and then regress. Inhibition of protein synthesis using drugs, such as [cycloheximide](#), prevents the appearance of the late puffs but not the early puffs. This result indicates a requirement for protein synthesis to generate of late puffs, most probably using [messenger RNA](#) derived from the genes in the early puffs (4). Some of the early puffs fail to diminish if protein synthesis is inhibited, suggesting that autoregulatory circuits exist in which the early puffs are turned off by their own gene products. The early puffs encode several regulatory [DNA-binding proteins](#) that carry out both gene activation and repression (5). Interestingly, chromosomal duplications or deletions of sites of early puffing indicate that more copies lead to greater and more rapid activation of the late puffs, whereas fewer copies lead to a reduced level of late puffing activity (6). Genes within the puffs, activated by ecdysone, encode proteins, such as dopa decarboxylase, an enzyme involved in cuticle formation and pigment synthesis in hypodermal cells (7), and the glue proteins required to attach the pupa to its substrate (8).

### Bibliography

1. W. Beerman (1964) *J. Exp. Zool.* **157**, 49–62.
2. J. J. Bonner and M. L. Pardue (1977) *Cell* **12**, 227–234.
3. M. Jamrich, A. L. Greenleaf, and E. K. F. Bautz (1977) *Proc. Natl. Acad. Sci. USA* **74**, 2079–2083.
4. M. Ashburner, C. Chihara, P. Meltzer, and G. Richards (1973) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 655–662.
5. L. D. Urness and C. S. Thummel (1995) *EMBO J.* **14**, 6239–6246.
6. V. K. Walker and M. Ashburner (1981) *Cell* **26**, 269–277.
7. G. P. Kraminsky et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4175–4179.
8. S. K. Beckendorf and F. C. Kafatos (1976) *Cell* **9**, 365–373.

### Suggestion for Further Reading

9. D. dePomerai (1990) *From Gene to Animal*, 2nd ed., Cambridge University Press, Cambridge,

UK.

## Pulse-Chase Experiments

[Radioisotopes](#) provide one of the most sensitive methods of tracing the chemical pathway of a moiety in a biological system, especially in a complex milieu like that of an entire cell. Because a radioactive atom differs from its nonradioactive counterpart only in the weight of its atomic nucleus, the cell treats any compound containing a radiolabel virtually the same as an unlabeled one. The [kinetics](#) of a process can be followed by performing a pulse-chase experiment using radiolabeled material. The labeled material is added for only a brief period of time (the pulse). This is followed by adding a large excess of the unlabeled material (the chase), which effectively stops utilization of the labeled material. Then the fate of the radioactive material that entered the biological process during the brief pulse is followed as a function of time.

If the radiolabeled material is metabolized, the pathway followed by the radioactive atoms can be determined by fractionating the system as a function of time and determining into which other molecules they are transformed. Metabolic pathways have been determined in this way, as have the pathways of assembly of complex structures, such as **viruses**. If the radiolabeled material is not metabolized but transported within an organ or cell, its location can be determined by [autoradiography](#). Such methods have assisted in elucidating the pathway taken by proteins secreted to the cell exterior (see [Topogenesis](#)) after their synthesis in the [endoplasmic reticulum](#). Pulse-chase experiments using radiolabeled proteins also demonstrate that most of them are degraded and then replaced by new synthesis. (see **Protein degradation**)

### Suggestion for Further Reading

B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland, New York.

## Pulsed Field Gel Electrophoresis

During [gel electrophoresis](#), long elastic polymers, such as **DNA** chains, and even large polystyrene spheres, become so entangled during their passage through the gel as ultimately to stop migrating. A change in the direction of the electric field can liberate these species from their entanglement. Net migration in one direction can be accomplished by alternating the field direction in pulses. The periodic change in the direction of the field also provides a mechanism for separation based on size; the longer chains are more sluggish in following the periodic changes of field direction than the shorter ones. Pulsed field gel electrophoresis (PFGE) of large DNA is based on the two mechanisms of liberation from entanglement and size separation ([1](#)).

Strictly, this is not a gel electrophoretic separation method because the separation is not directly related to the concentration of the gel; the gel acts as an anticonvection medium and, more importantly, as an array of obstacles entangling, stretching, and releasing the elastic, reptating DNA

chain (2).

A pulsed power supply can also be used to advantage in the electrophoresis of [proteins](#) to provide improved resolving capacity. This presumably functions by suppressing the band spreading that would occur as a result of effective heat dissipation during the off-power phase in cases of inadequate joule heat dissipation capacity of the apparatus.

### Bibliography

1. R. M. Gemmill (1991) *Adv. Electrophoresis* **4**, 1–48.
2. B. Birren and E. Lai (1990) *Methods* **1**, 129–214.

### Suggestions for Further Reading

3. C. L. Smith (ed.) (1993) Paper Symposium: Changing directions in electrophoresis. *Electrophoresis* **14**, 249–370.

## ***Pumilio* Gene**

The *pumilio* is a **maternal effect** gene required for development of the fruit fly *Drosophila melanogaster*. *pumilio* is required during oogenesis for the development and maintenance of the germ-line [stem cell](#) lineage and in the embryo for patterning of the abdomen. *pumilio*, in concert with **nanos**, represses [translation](#) of maternally provided *hunchback* [messenger RNA](#) in the *Drosophila* embryo. Translational regulation is mediated by sequences in the 3' untranslated region of *hunchback*, and Pumilio binds specifically to these regulatory sequences. Pumilio is a member of a conserved family of **RNA-binding** proteins with a novel RNA-binding motif. Ongoing studies seek to understand the interaction of the Pumilio protein with its RNA target.

### 0.1. Protein and RNA Structure

The *pumilio* gene spans over 160 kbp. The primary transcript contains **12 introns** (one over 120 kbp) and produces two mature mRNAs (6.8 kb and 6.7 kb) that encode identical proteins, but use different first exons. The *pumilio* mRNA is detected at high levels in ovaries, 0 to 8 hour embryos, and adult females (1, 2). Low levels of the transcript are also detected in 8- to 24-h embryos and pupae. The Pumilio protein is a 1533-amino-acid-residue polypeptide with a predicted molecular weight of 160 kDa. Analysis of Pumilio protein reveals eight imperfect repeats of 36 residues each (1). Recently it has been shown that these repeats constitute the *pumilio* RNA-binding domain (3, 4). This RNA-binding motif (known as the *pumilio* homology domain or PUM-HD) is conserved from yeast to humans (3, 5).

### 1. *pumilio* in Oogenesis

A *Drosophila melanogaster* ovary consists of approximately 16 ovarioles each representing an independent egg assembly line (for review see Ref. 6). Every ovariole is a linear array of developing egg chambers, beginning at the anterior with the germarium, where the progeny of the germ-line and somatic stem cells are organized into egg chambers. At the tip of the germarium, the germ-line stem cell produces a new stem cell and a cystoblast by asymmetric division. The cystoblast undergoes four rounds of division with incomplete cytokinesis, resulting in 16 cells connected by cytoplasmic bridges; one of these cells becomes the oocyte while the other 15 become nurse cells. Somatic derived follicle cells surround the egg chamber. As the egg chambers mature, the nurse cells synthesize and deposit factors into the oocyte that are required for the development of the oocyte and

early embryo. Late in oogenesis, the nurse cells contract and dump the contents of their cytoplasm into the oocyte.

Pumilio protein is expressed at high levels in the germ line stem cells and at lower levels in the dividing cystoblasts and the soma (7, 8). Females *trans*-heterozygous for strong *pumilio* alleles lay a few eggs, but rapidly become sterile due to a loss of germ-line stem cells from the ovary (7, 8). Analysis of *pumilio* mutant ovaries demonstrates that germ cells often fail to be incorporated into the ovarioles, suggesting that *pumilio*, like *nanos*, is required for germ-line cells to populate the ovary. Ovarioles with germ-line cells lack germ-line stem cells, and germline cells differentiate directly into egg chambers. This suggests that *pumilio* is required for establishment of the germ-line stem cell fate. Additionally, *pumilio* mutant ovaries from older flies often contain tumorous growths resulting from overproliferation of somatic cells, suggesting a role for *pumilio* in the soma as well.

## 2. *pumilio* in Embryonic Patterning

In the late oocyte and early embryo, *pumilio* mRNA and protein are distributed throughout the embryo (2). Pumilio activity is required for the **translational repression** of maternal *hunchback* RNA in the posterior half of the embryo. Hunchback is a zinc-finger [transcription factor](#) that controls [transcription](#) of the abdominal gap genes (9). In addition to Pumilio, repression of maternal *hunchback* translation is mediated by sequences within the *hunchback* 3'-untranslated region (known as *nanos* [response elements](#) or NREs) and the activity of the posterior determinant *nanos* (10). Pumilio protein binds to the *hunchback* NREs (3, 4, 11). Mutations that disrupt Pumilio binding *in vitro* have a strong effect on the regulation of *hunchback* *in vivo* (4). Together Pumilio and Nanos act through the *hunchback* NREs to direct removal of the maternal *hunchback* [poly \(A\)](#) tail in the posterior of the embryo, suggesting that deadenylation mediates translational repression (12). Ectopic expression of Nanos in the eye **imaginal** disc results in disruption of ommatidial development (4). Loss of Pumilio function suppresses the activity of ectopic Nanos, suggesting that *nanos* and *pumilio* are cooperating to regulate some gene product required for eye development. Additionally, ectopic *nanos* can repress the expression of reporter genes bearing the *hunchback* NREs and an [internal ribosome entry site](#), suggesting that *nanos* and *pumilio* regulate translation independent of the **5'-cap** (4).

## 3. The Pumilio Homology RNA-Binding Domain

The *pumilio* homology domain has been identified in proteins or [expressed sequence tags](#) from human, [mouse](#) (*M. musculus*), rat (*R. norvegicus*), [nematodes](#) (*C. elegans*), yeast (*S. cerevisia* and *S. pombe*), plants (*A. thaliana*), rice (*O. sativa*), and corn (*Z. mays*) (3, 5). The high degree of sequence conservation among these very diverse species suggests that the PUM-HD is an ancient protein motif. Furthermore, a human PUM-HD can bind to the *hunchback* NREs, albeit with slightly altered sequence specificity (3). *fem-3* is a gene required to determine the sexual fate of *C. elegans* hermaphrodite germ line (13, 14). The regulation of *fem-3* is very analogous to that seen for *Drosophila hunchback*. *fem-3* contains a small *cis*-element within its 3'-untranslated region (known as the point-mutation element) that, when mutated, causes aberrant translation of *fem-3* (15, 16). Recently, FBF-1 and FBF-2 were cloned based on their ability to bind to the wild type, but not mutant, point-mutation element. FBF-1 and FBF-2 each contain a PUM-HD, suggesting that these proteins share not only a conserved protein motif with *pumilio*, but also a conserved function (5). Other PUM-HD family members have recently been shown to have a role in aging and recovery from stress, but their mechanism of action remains unclear (17).

## Bibliography

1. D. D. Barker, C. Wang, J. Moore, L. K. Dickinson, and R. Lehmann (1992) *Genes Dev.* **6**, 2312–2326.
2. P. M. Macdonald (1992) *Development* **114**, 221–232.
3. P. D. Zamore, J. R. Williamson, and R. Lehmann (1997) *RNA* **3**, 1421–1433.



4. R. P. Wharton, J. Sonoda, T. Lee, M. Patterson, and Y. Murata (1998) *Mol. Cell* **1**, 863–872.
5. B. Zhang, M. Gallegos, A. Puoti, E. Durkin, S. Fields, J. Kimble, and M. P. Wickens (1997) *Nature* **390**, 477–484.
6. A. C. Spradling (1993) In *The Development of Drosophila melanogaster*, Vol. **1** (M. B. Arias and A. M. Bates, eds.), Cold Spring Harbor Press, Plainview, New York, pp. 1–70.
7. A. Forbes and R. Lehmann (1998) *Development* **125**, 679–690.
8. H. Lin and A. C. Spradling (1997) *Development* **124**, 2463–2476.
9. M. Hulskamp, C. Pfeifle, and D. Tautz (1990) *Nature* **346**, 577–580.
10. R. P. Wharton and G. Struhl (1991) *Cell* **67**, 955–967.
11. Y. Murata and R. P. Wharton (1995) *Cell* **80**, 747–756.
12. C. Wreden, A. C. Verrotti, J. A. Schisa, M. E. Lieberfarb, and S. Strickland (1997) *Development* **124**, 3015–3023.
13. M. K. Barton, T. B. Schedl, and J. Kimble (1987) *Genetics* **115**, 107–119.
14. T. A. Rosenquist and J. Kimble (1988) *Genes Dev.* **2**, 606–616.
15. J. Ahringer and J. Kimble (1991) *Nature* **349**, 346–348.
16. J. Ahringer, T. A. Rosenquist, D. N. Lawson, and J. Kimble (1992) *EMBO J.* **11**, 2303–2310.
17. Y. Kikuchi, Y. Oka, M. Kobayashi, Y. Uesono, A. Toh-e, and A. Kikuchi (1994) *Mol. Gen. Genet.* **245**, 107–116.

### Suggestions for Further Reading

18. A. Forbes and R. Lehmann (1998) *nanos* and *pumilio* have critical roles in the development and function of *Drosophila* germline stem cells. *Development* **125**, 679–690.
19. D. R. Gallie (1998) A tale of two termini: a functional interaction between the termini of an mRNA is a prerequisite for efficient translation initiation. *Gene* **216**, 1–11.
20. P. E. Kuwabara (1998) Gametogenesis: keeping the male element under control. *Curr. Biol.* **8**, R278–R281.
21. R. P. Wharton, J. Sonoda, T. Lee, M. Patterson, and Y. Murata (1998) The Pumilio RNA-binding domain is also a translational regulator. *Mol. Cell* **1**, 863–872.
22. P. D. Zamore, J. R. Williamson, and R. Lehmann (1997) The Pumilio protein binds RNA through a conserved domain that defines a new class of RNA-binding proteins. *RNA* **3**, 1421–1433.

### Pure Line

A population of identical **homozygotes** constitutes a pure line. Pure lines “breed true”; that is, they maintain their phenotype indefinitely if genetic crosses are restricted to members of the pure line (see [True Breeding](#)).

The concept of pure line was introduced by Johannsen ([1](#)). Pure lines are obtained through repeated self-fertilization of heterozygous individuals or consanguineous crosses in heterozygous populations. Under [Mendelian inheritance](#), the offspring of a single self-fertilization will be heterozygous for one-half as many genes as the parent. Repeated self-breeding approaches full homozygosity asymptotically. Several generations of self-fertilization are needed to obtain an acceptable pure line (97% of the genes that were originally heterozygous would have turned homozygous after five

generations). A population of an autogamous organism is a mixture of pure lines; and, if started from a single organism, the whole population is a single pure line.

Repeated crosses between close relatives also approach full homozygosity, but more slowly than self-fertilization. Thus, about 20 generations of crosses between a brother and a sister from each generation may be carried out to produce a pure line of mammals.

Populations of allogamic organisms, in which crosses occur at random or even with preference for distantly related partners (outbreeding), maintain a high level of heterozygosity. In these populations, rare recessive alleles that are deleterious when homozygous have little effect on the fitness of the population, because most of them are maintained in heterozygotes. The total number of such alleles may be high, because they may occur at many genes. The overall deleterious effect on the population (the *genetic load* of the population) may be considerable. Thus, human populations suffer on average a genetic load that is equivalent to that caused by the presence in each individual of four recessive lethal alleles (or a larger number of less disabling alleles). Inbreeding in allogamic organisms bring the deleterious recessive alleles to homozygosity; the immediate consequence is an increase in the frequency of defective offspring, or, in another words, an increase in the genetic load of the population. This phenomenon is called *inbreeding depression* or *inbreeding degeneration*. As inbreeding continues, the deleterious alleles are selected out and eventually disappear. The original heterozygous populations are often more fit than the resulting pure lines because they profit from [heterosis](#) and balanced polymorphisms; the main advantage of pure lines is the quick production of many individuals with the same well-adapted genotype, while the allogamy continuously generates new genotypes.

#### Bibliography

1. W. Johannsen (1903) *Über Erbllichkeit in Populationen und reinen Linien*, Gustav Fischer, Jena.

#### Suggestion for Further Reading

2. R. Frankham (1995) Conservation genetics. *Annu. Rev. Genet.* **29**, 305–327. (This article discusses inbreeding depression in natural populations.)

## Purine Ribonucleotide Metabolism

### 1. Biosynthesis of the Purine Ring

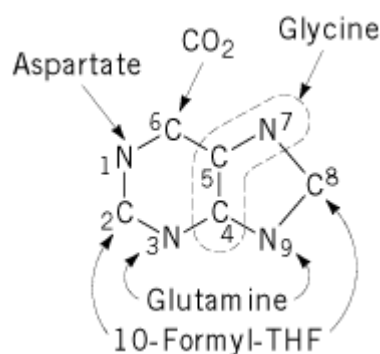
The biosynthetic pathway for the purine ring was described in the 1950s, in a classic series of experiments carried out in the laboratories of John Buchanan and G. Robert Greenberg. Birds excrete most of their nitrogen in the form of uric acid, an oxidized purine. Therefore, by feeding **radiolabeled** precursors to pigeons and then chemically degrading the uric acid crystallized from their droppings, it was possible to identify precursors to each position in the ring, and this led to identification of the reactions and isolation of the [enzymes](#) involved. This pathway is called the *de novo* pathway because it involves synthesis of the purine ring from low-molecular-weight precursors. A separate set of pathways is referred to as salvage pathways because they involve reutilization of preformed purine ring-containing compounds, usually nucleosides or nucleobases released by nucleic acid degradation (see [Salvage Pathways To Nucleotide Biosynthesis](#)).

### 2. *De Novo* Biosynthesis of Purine Nucleotides

#### 2.1. Biosynthesis of Inosinic Acid

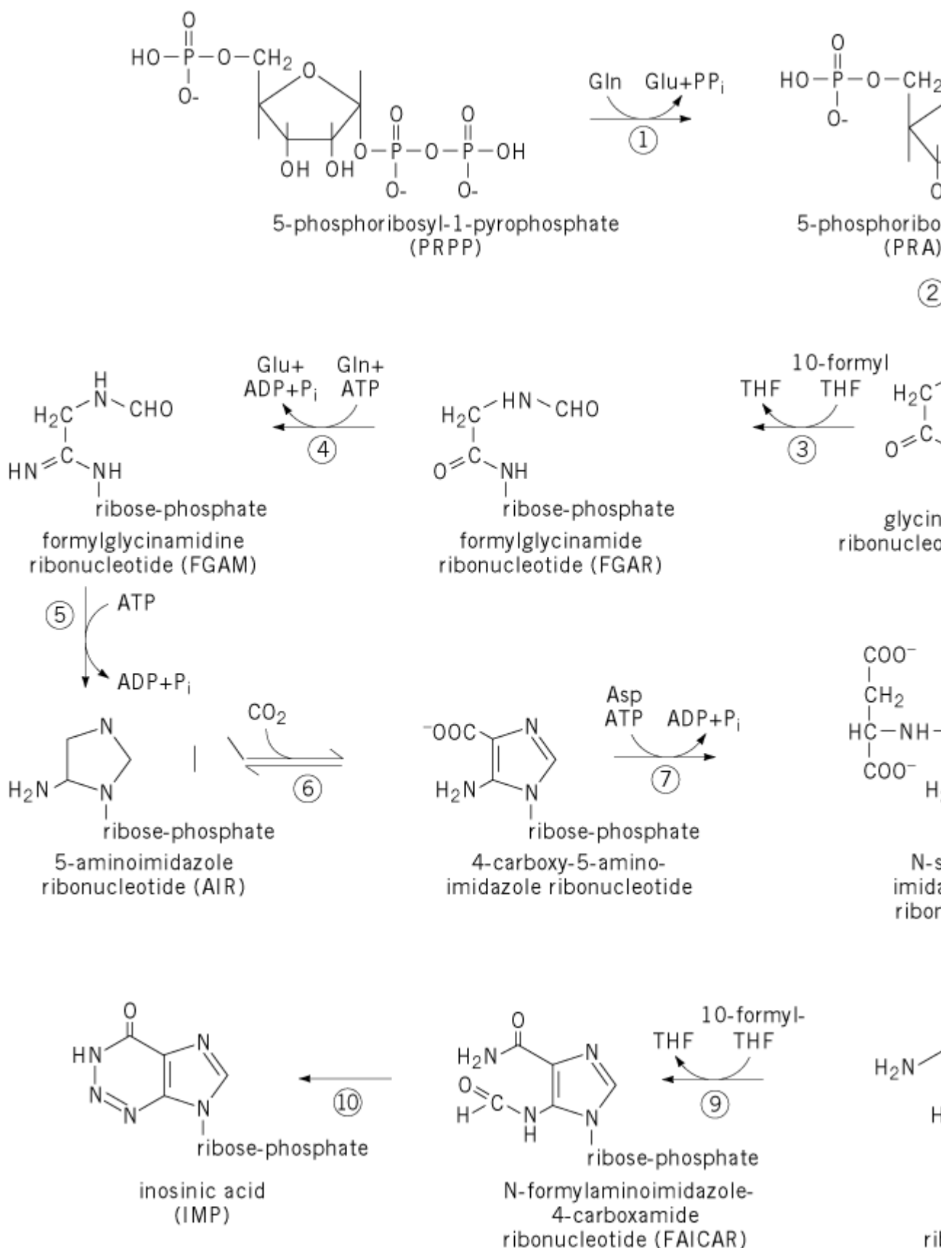
The early radiolabeling studies identified the precursors of the purine ring (Fig. 1) as glycine, glutamine amide nitrogen, CO<sub>2</sub>, aspartate amino nitrogen, along with the “one-carbon pool,” which included formate or the b-carbon of serine but which is now known to come directly from the formyl group of N<sup>10</sup>-formyltetrahydrofolate (10-formyl-THF). The pathway is sensitive to inhibition by folate antagonists, such as methotrexate, and glutamine antagonists, such as azaserine, and this sensitivity largely explains the chemotherapeutic efficacy of these classes of antimetabolites. Aside from a number of parasitic protozoans, which lack the capacity for purine ring synthesis and depend on salvage pathways, the *de novo* purine synthetic pathway is virtually identical among all organisms examined to date.

**Figure 1.** Metabolic sources of the atoms of the purine ring, as determined by administration of labeled precursors.



That pathway, summarized in Figure 2, involves the stepwise assembly of the purine ring at the nucleotide level, on a ribose 5-phosphate backbone. In the first committed reaction (step 1), the glutamine amide group displaces the pyrophosphate at C-1 of 5-phosphoribosyl-1-pyrophosphate (PRPP), to give 5-phosphoribosylamine, the simplest possible nucleotide (its base is ammonia). The enzyme, PRPP amidotransferase, is the principal point for control of the overall pathway. There follows the ATP-dependent incorporation of glycine (step 2), the first of two formyltransferase reactions involving N<sup>10</sup>-formyltetrahydrofolate (step 3), a second glutamine amidotransferase reaction (step 4), and a ring closure (step 5), yielding the imidazole portion of the bicyclic purine ring. The next reaction (step 6) is a CO<sub>2</sub>-fixation reaction unusual in that it does not require [biotin](#) as a cofactor. Next (step 7), an ATP-dependent reaction links the aspartate amino nitrogen to the carboxyl group formed in the previous reaction. There follows an a,b-elimination reaction (step 8) in which the carbon skeleton of aspartate is released as fumarate, with its nitrogen becoming part of a carboxamide group. A second formyltransferase (step 9) creates a product that undergoes an intramolecular condensation (step 10), yielding inosinic acid (IMP), the first completed purine nucleotide. IMP has hypoxanthine as its purine base.

**Figure 2.** The *de novo* biosynthetic pathway to inosinic acid. Enzyme names: 1, PRPP amidotransferase; 2, GAR synthetase; 3, FGAR amidotransferase; 4, FGAM cyclase; 5, AIR carboxylase; 6, SAICAR synthetase; 7, SAICAR lyase; 8, SAICAR transaminase; 9, AICAR transaminase; 10, AICAR cyclase.

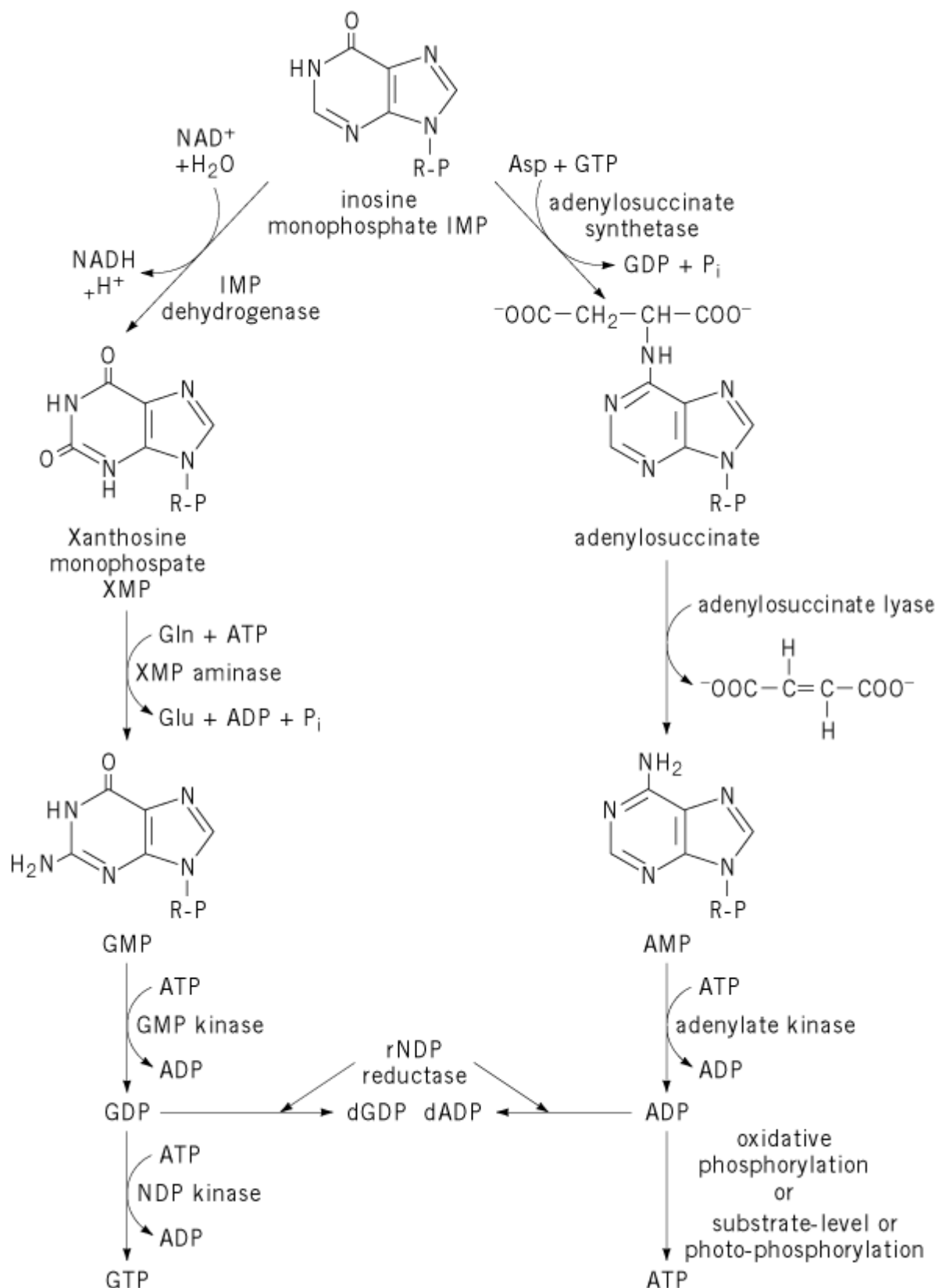


## 2.2. Conversion of IMP to Adenine and Guanine Nucleotides

IMP represents a branch point in the synthesis of adenine and guanine nucleotides, as shown in Figure 3. En route to GMP, IMP dehydrogenase oxidizes the hypoxanthine base of IMP to the

xanthine base of xanthosine monophosphate (XMP). An ATP-dependent glutamine amidotransferase converts XMP to GMP. En route to AMP, IMP undergoes a reaction sequence involving aspartate, which is very similar to reactions 7 and 8 of the IMP synthetic pathway, except that GTP, not ATP, is the energy cofactor. The release of fumarate is the last step in AMP synthesis.

**Figure 3.** Pathways from inosinic acid to ATP and GTP. R-P is a ribose 5-phosphate moiety. NDP kinase is nucleoside diphosphate kinase. A few organisms reduce ribonucleotides to deoxyribonucleotides at the triphosphate level rather than, as shown here, the diphosphate level.



AMP and GMP are converted to the respective 5'-diphosphates by specific nucleotide kinases. The resultant ADP and GDP are converted to the triphosphates by [nucleoside diphosphate kinase](#), an active and nonspecific enzyme that transfers the γ-(outermost) phosphate of any nucleoside 5'-triphosphate to any nucleoside 5'-diphosphate. Because the reaction has an equilibrium constant close to 1, the direction of the reaction is determined primarily by concentrations of substrates and products. Because ATP is by far the most abundant nucleotide in most aerobic cells, the principal

function of nucleoside diphosphate kinase is to catalyze the ATP-dependent conversion of nucleoside diphosphates to triphosphates, using ATP that was produced by oxidative phosphorylation.

In most organisms, the ribonucleoside diphosphates (purine and pyrimidine) serve as precursors for biosynthesis of deoxyribonucleotides (see [Ribonucleotide Reductases](#) and [Deoxyribonucleotide Biosynthesis And Degradation](#)). In some organisms, however, those precursors are the respective triphosphates.

### 2.3. Multifunctional Enzymes and Enzyme Complexes

Vertebrate cells contain a number of the various enzyme activities of IMP biosynthesis in the form of multifunctional enzymes. This was first suspected when **cloned** vertebrate **cDNAs** encoding purine synthetic enzymes were found to complement multiple genetic defects in purine synthesis after [transformation](#) into *Escherichia coli* (1). By this means, it was found that a single enzyme catalyzes the second, third, and fifth reactions shown in Figure 2. Similar evidence indicated that the sixth and seventh reactions are catalyzed by a bifunctional enzyme. Moreover, the two transformylase enzymes (reactions 3 and 9) in some animals constitute part of a tightly associated multienzyme complex that also contains several activities of tetrahydrofolate metabolism and single-carbon mobilization. The metabolic rationale for all these enzyme associations has not been established, but it may well involve the cell's attempt to utilize scarce or unstable intermediates more efficiently by facilitating their transfer from active site to active site within the same reaction sequence.

### 2.4. Regulation

Control of purine nucleotide synthesis involves both allosteric and genetic regulation. In most cells, PRPP synthetase, which synthesizes the first intermediate in IMP synthesis, is inhibited by AMP, ADP, and GDP, whereas PRPP amidotransferase (reaction 1), the primary control point for the overall reaction (2), is inhibited allosterically by AMP, ADP, GMP, and GDP. In *E. coli*, biosynthesis of the enzymes of IMP synthesis is inhibited by a **repressor** encoded by the *purR* gene. This protein binds either hypoxanthine or guanine, and the resultant protein-purine base complex binds to DNA sites upstream from **promoters** for several purine (and pyrimidine) biosynthetic enzymes. The crystal structure of the PurR repressor (3) shows it to be closely related to the well-known [Lac repressor](#), which controls the lactose utilization operon by similar mechanisms.

Conversion of IMP to AMP and GMP is also regulated allosterically. GMP inhibits IMP dehydrogenase, the enzyme that converts IMP to GMP, whereas AMP controls its own formation by inhibiting the addition of aspartate to AMP to form adenylosuccinate (Fig. 3). Also, it may have regulatory significance that ATP is involved in the conversion of IMP to GMP, whereas GTP is required for one of the reactions leading from IMP to AMP.

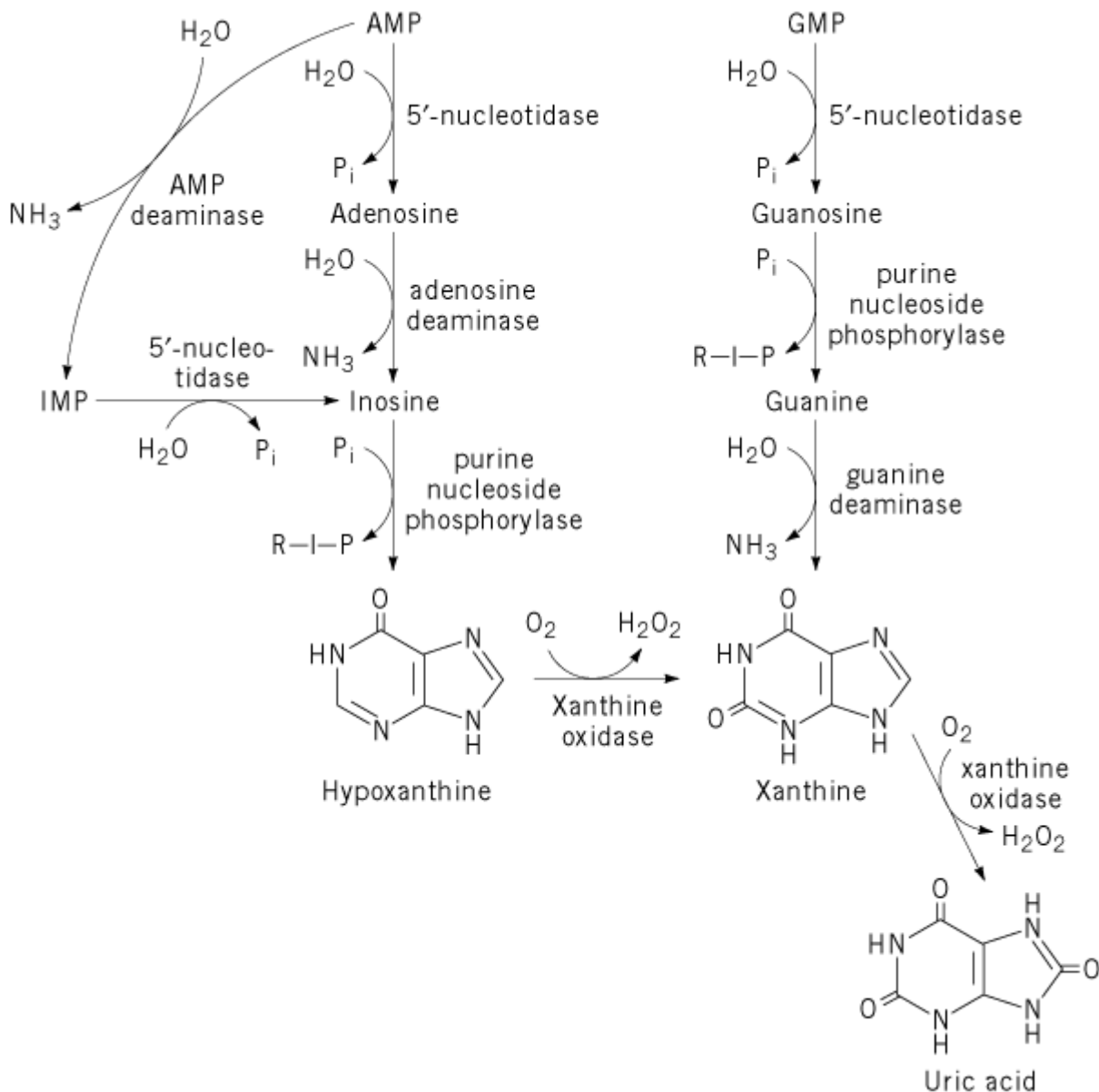
## 3. Catabolism of purine nucleotides

Nucleotides released by enzymatic digestion of nucleic acids are rather efficiently reutilized for nucleic acid biosynthesis in most cells. However, pathways of nucleotide degradation are significant, as shown by sometimes unexpected and severe consequences of genetic deficiencies in humans of particular enzymes of purine degradation, as dealt with in the next section.

In primates, the end product of purine nucleotide catabolism is uric acid, which is excreted as such. Pathways leading to uric acid vary considerably in different tissues and cells. Most of the reactions involved are shown in Figure 4. Note, for example, that AMP degradation can begin either with deamination, to yield IMP, or with hydrolysis, to yield adenosine. In mammals, the deamination pathway is particularly active in muscle tissue. Both pathways lead to the nucleoside inosine, which is cleaved by inorganic phosphate and purine nucleoside phosphorylase, yielding ribose 1-

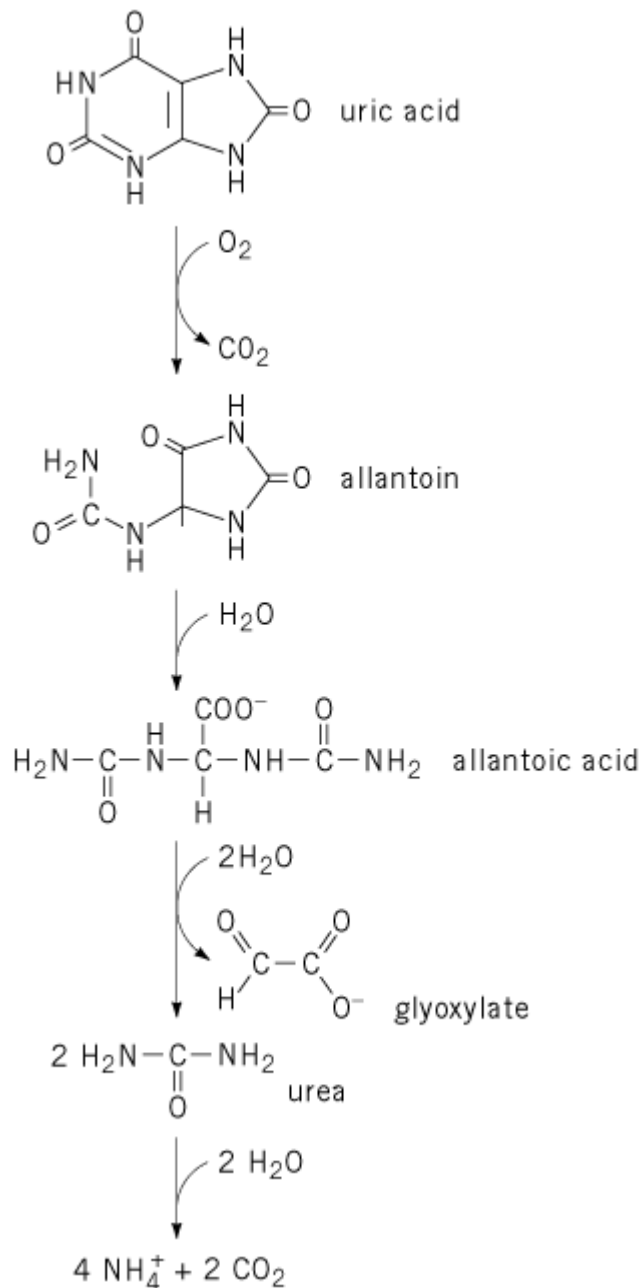
phosphate and hypoxanthine. Hypoxanthine is oxidized by the versatile molybdenum- and iron-containing enzyme, xanthine oxidase, to xanthine, which is also produced by guanine nucleotide catabolism. Xanthine is also acted on by xanthine oxidase, to give uric acid. As noted, the process ends here in primates. However, most animals further degrade uric acid to allantoin and then to allantoic acid. Some fishes excrete allantoic acid, but most aquatic animals further catabolize allantoic acid to urea and, in the case of marine invertebrates, to ammonia. These latter pathways are summarized in Figure 5.

**Figure 4.** Pathways of purine nucleotide catabolism to uric acid. R-1-P is ribose 1-phosphate.



**Figure 5.** Metabolic degradation of uric acid.





#### 4. Clinical abnormalities of purine metabolism

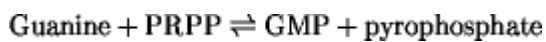
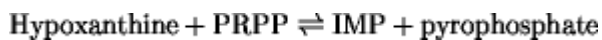
Six classes of metabolic disorders involving purines have been described (4-9). Three of these conditions—adenine phosphoribosyltransferase deficiency (4), adenylate deaminase deficiency (5), and xanthine oxidase deficiency (6)—are quite rare and will not be discussed here. The three conditions that will be described either are relatively common or their study has revealed important metabolic principles and relationships, or both.

##### 4.1. Hyperuricemia and Gout

Gout refers to a family of diseases in which prolonged elevation of uric acid levels in tissues and blood leads to its crystallization in the joints, causing intermittent attacks of an acute inflammatory arthritis (7). Prolonged hyperuricemia does not always lead to gout attacks; factors leading to the precipitation of urate salts are not thoroughly understood.

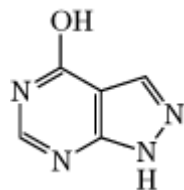
Normally, about two-thirds of the uric acid produced in purine catabolism is excreted through the kidneys; the remainder is further broken down by intestinal bacteria. Renal malfunction can lead to elevation of blood uric acid levels, which is a cause of gout. Gout also results from abnormally high purine nucleotide synthesis, which leads to degradation of the excess. Three specific biochemical changes are known to cause this condition. First, hyperactivity of PRPP synthetase elevates intracellular concentrations of PRPP, the substrate for PRPP amidotransferase, the first and rate-controlling step in the *de novo* pathway (reaction 1, Fig. 2), and this increases flux through the whole pathway.

A second form of gout results from partial or complete deficiency of a purine salvage enzyme, hypoxanthine-guanine phosphoribosyltransferase (HGPRT). This enzyme catalyzes the following reactions:



It is not yet clear why this deficiency accelerates purine synthesis. One model states that in the absence of these salvage pathways, GMP levels decline which, in turn, modulates the feedback inhibition of PRPP amidotransferase by this nucleotide. Alternatively, it has been proposed that decreased flux through this salvage pathway causes PRPP to accumulate and this, in turn, accelerates flux through PRPP amidotransferase by substrate level control.

The third enzymatic deficiency leading to gout is a deficiency of glucose-6-phosphatase (glycogen storage disease type I). The relationship between this abnormality and uric acid accumulation is still obscure. Whatever the basis for hyperuricemia, the most effective drug for treating the condition is the xanthine oxidase inhibitor allopurinol.



Inhibition of xanthine oxidase causes hypoxanthine and xanthine—both of which are both more soluble and less toxic than uric acid—to accumulate, thereby preventing uric acid precipitation.

#### 4.2. Lesch–Nyhan Syndrome

As noted, a partial HGPRT deficiency leads to gout. A total deficiency of this enzyme has far more serious consequences. Lesch–Nyhan syndrome (8) involves not only severe hypericemia and gout but, in addition, the nervous system develops abnormally, leading to spasticity and behavioral problems, including aggressive behavior toward others and self-mutilation. Because the gene for HGPRT lies on the [X-chromosome](#), the Lesch–Nyhan syndrome is **sex-linked**, having been observed only in males. Although the gouty arthritis of Lesch–Nyhan syndrome usually responds well to allopurinol, there is no known cure for the developmental and neurological abnormalities, and afflicted individuals rarely live beyond age 20. Moreover, although it is clear that all symptoms of this condition arise from the HGPRT deficiency, the specific relationship between the enzyme deficiency and the neuropathology is not understood.

#### 4.3. Immunodeficiencies Caused by Purine Abnormalities (9)

In 1972, a patient with an inheritable combined immunodeficiency involving both **B** and **T cells** was found to be deficient in the enzyme adenosine deaminase (ADA), which converts adenosine to inosine in purine nucleotide catabolism (see Fig. 4). This surprising finding was followed by the

discovery of other families in which the same enzyme deficiency was associated with severe combined immunodeficiency and, by now, several hundred such families have been described. In 1975, a second form of immunodeficiency was found to result from a different purine abnormality, namely, a deficiency of purine nucleoside phosphorylase, an enzyme whose primary function is to degrade guanosine to guanine in purine catabolism. Often, purine nucleoside phosphorylase deficiency involves only defective T-cell function in the immune system.

What is the relationship between deficiencies in two obscure purine catabolic enzymes and defective immune function? A clue came when it was found that erythrocytes of ADA-deficient patients contained high levels of deoxyadenosine and dATP. Since human erythrocytes are without nuclei, the presence of a DNA precursor seemed gratuitous. Subsequently, it was found that dATP accumulates in many tissues. It arises because ADA acts on deoxyadenosine as well as adenosine. Lymphoid tissues have very high activities of purine salvage enzymes, which reutilize products released in nucleic acid breakdown of cells undergoing [apoptosis](#). Accumulation of deoxyadenosine, when its catabolism is blocked, leads to its conversion to dATP in these tissues, and it accumulates in both red and white blood cells. dATP is a potent allosteric inhibitor of [ribonucleotide reductase](#), and its accumulation can block the white cell proliferative response that results from immunochallenge by inhibiting an essential step in replication of DNA—the synthesis of its precursors. There are also indications that excessive accumulation of dATP leads to an ATP deficiency in some cells. By contrast, in PNP-defective cells, the deoxyribonucleotide that accumulates is primarily dGTP, which is a less potent inhibitor of ribonucleotide reductase. This may explain the somewhat milder immune dysfunction associated with PNP deficiency. However, a completely different mechanism for the toxic effect has recently come to light from *in vitro* studies showing that dATP, in combination with cytochrome c, triggers a chain of protease activation steps leading ultimately to apoptosis ([10](#)).

Adenosine deaminase deficiency is the first condition to be treated by gene therapy.

#### Bibliography

1. H. Zalkin and J. Dixon (1992) *Prog. Nucleic Acid Res. Mol. Biol.* **42**, 259–287.
2. T. Yamaoka, M. Kondo, S. Honda, H. Iwahana, M. Moritani, S. Ii, K. Yoshimoto, and M. Itakura (1997) *J. Biol. Chem.* **272**, 17719–17725.
3. M. A. Schumacher, K. Y. Choi, H. Zalkin, and R. G. Brennan (1994) *Science* **266**, 763–770.
4. H. A. Simmonds, A. S. Sahota, and K. J. Van Acker (1995) In *The Metabolic and Molecular Basis of Inherited Disease*, 7th ed. (C. R. Scriver, A. L. Beaudet, W. S. Sly, and D. Valle, eds.) McGraw-Hill, New York, pp. 1707–1724.
5. R. L. Sabina and E. W. Holmes (1995) In *Metabolic and Molecular Basis of Inherited Disease*, pp. 1769–1780.
6. H. A. Simmonds, S. Reiter, and T. Nishino (1995) *Metabolic and Molecular Basis of Inherited Disease*, pp. 1781–1798.
7. M. A. Becker and B. J. Roessler (1995) *Metabolic and Molecular Basis of Inherited Disease*, pp. 1655–1678.
8. B. J. F. Rossiter and C. T. Caskey (1995) *Metabolic and Molecular Basis of Inherited Disease*, pp. 1679–1706.
9. M. S. Hershey and B. S. Mitchell (1995) *Metabolic and Molecular Basis of Inherited Disease*, pp. 1725–1768.
10. P. Li, D. Nijhawan, I. Budihardjo, S. M. Srinivasula, M. Ahmad, E. S. Alnemri, and X. Wang (1997) *Cell* **91**, 479–489.

#### Suggestion for Further Reading

11. R. L. Blakley and S. J. Benkovic (1984) *Folates and Pterins*, Vol. I, Academic Press, New

York. Various chapters treat the biochemistry of folate coenzymes in detail, as well as their involvement in purine metabolism.

## Puromycin

### 1. Introduction

Puromycin, first described in 1952, is a secondary metabolite of *Streptomyces alboniger* that blocks **protein biosynthesis**. Puromycin is a structural analogue of the 3' end of aminoacyl-[transfer RNA](#), but differs from tRNA insofar as the aminoacyl residue is linked to the ribose via an amide bond rather than an ester bond. Puromycin, like aminoacyl-tRNA, binds to the A site of the [ribosome](#) peptidyl-transferase center. When the A site is occupied by puromycin, [peptidyl-transferase](#) links the peptide residues of the peptidyl-tRNA in the ribosomal P site covalently to puromycin. Since the amide bond cannot be cleaved by the ribosome, no further peptidyl transfer takes place, and the peptidyl-puromycin complex falls off the ribosome.

Puromycin has facilitated the elucidation of ribosome structure and function, particularly the peptidyl-transferase reaction. The puromycin reaction provided the first demonstration that ribosomes can catalyze the formation of [peptide bonds](#) in the absence of cytoplasmic factors and without an added energy source, as well as the basis for distinguishing the A site from the P site (peptidyl-tRNA bound to the P site reacts with puromycin, whereas peptidyl-tRNA bound to the A site does not). The fragment reaction served to localize the peptidyl-transferase center on the 50S subunit of the bacterial ribosome.

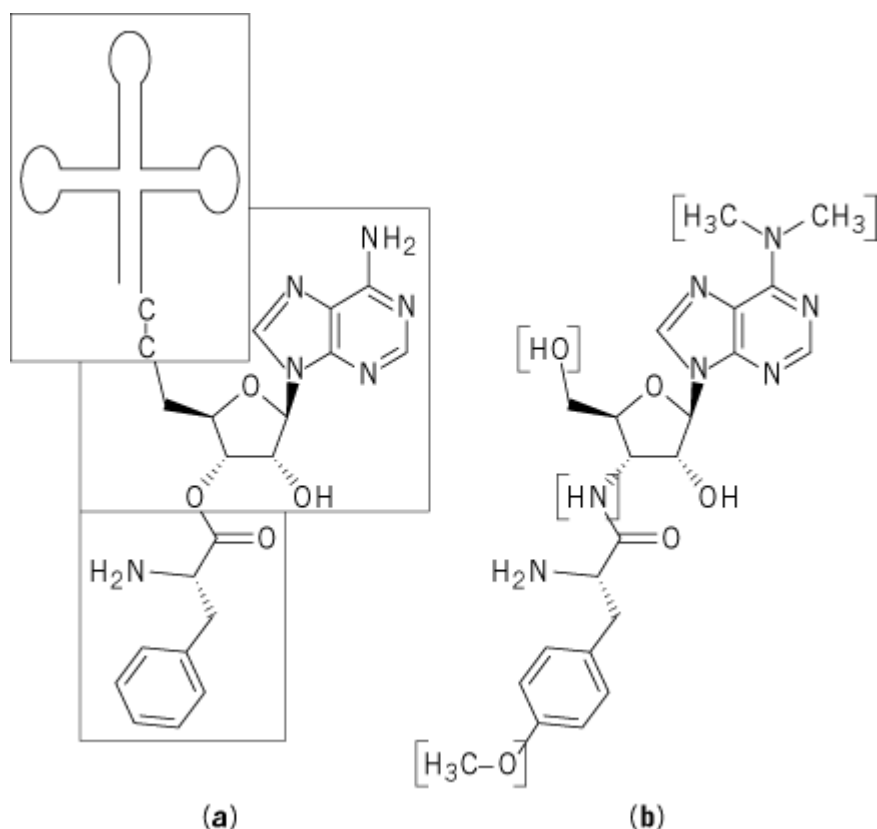
### 2. Antibacterial Spectrum

In light of the structural resemblance of puromycin to aminoacyl-tRNA, it is not surprising that the spectrum of activity of this antibiotic is broad ([1](#)). The spectrum includes not only eubacteria and archaeobacteria but also eukaryotes; consequently, puromycin is too toxic to be used clinically.

### 3. Structure–Activity Relationships

Puromycin essentially mimics aminoacyl-tRNA, in particular, phenylalanyl-tRNA (Phe-tRNA). The chemical structure of puromycin, 3'-(*α*-amino-*p*-methoxyhydrocinnamamido)-3'-deoxy-*N,N*-dimethyladenosine, is shown in [Figure 1](#). In order for puromycin to retain its activity, both the nucleoside and amino acid moieties are required, and the amino acid moieties must be in the L configuration and have an unsubstituted amino group (i.e., the same requirement as that for amino acids incorporated into protein) ([2](#)). Derivatives with aromatic side chains are more active than aliphatic ones ([3](#)), which suggests that the aromatic aminoacyl residues may enhance the ability of puromycin to bind to the ribosomes.

**Figure 1.** Comparison of the chemical structure of puromycin and Phe-tRNA. (a) Puromycin mimics the amino acid part of Phe-tRNA (bottom box) and its 3'-terminal A (middle box). The tRNA body itself (top box) is absent in the puromycin molecule. (b) The differences between the two molecules are bracketed.

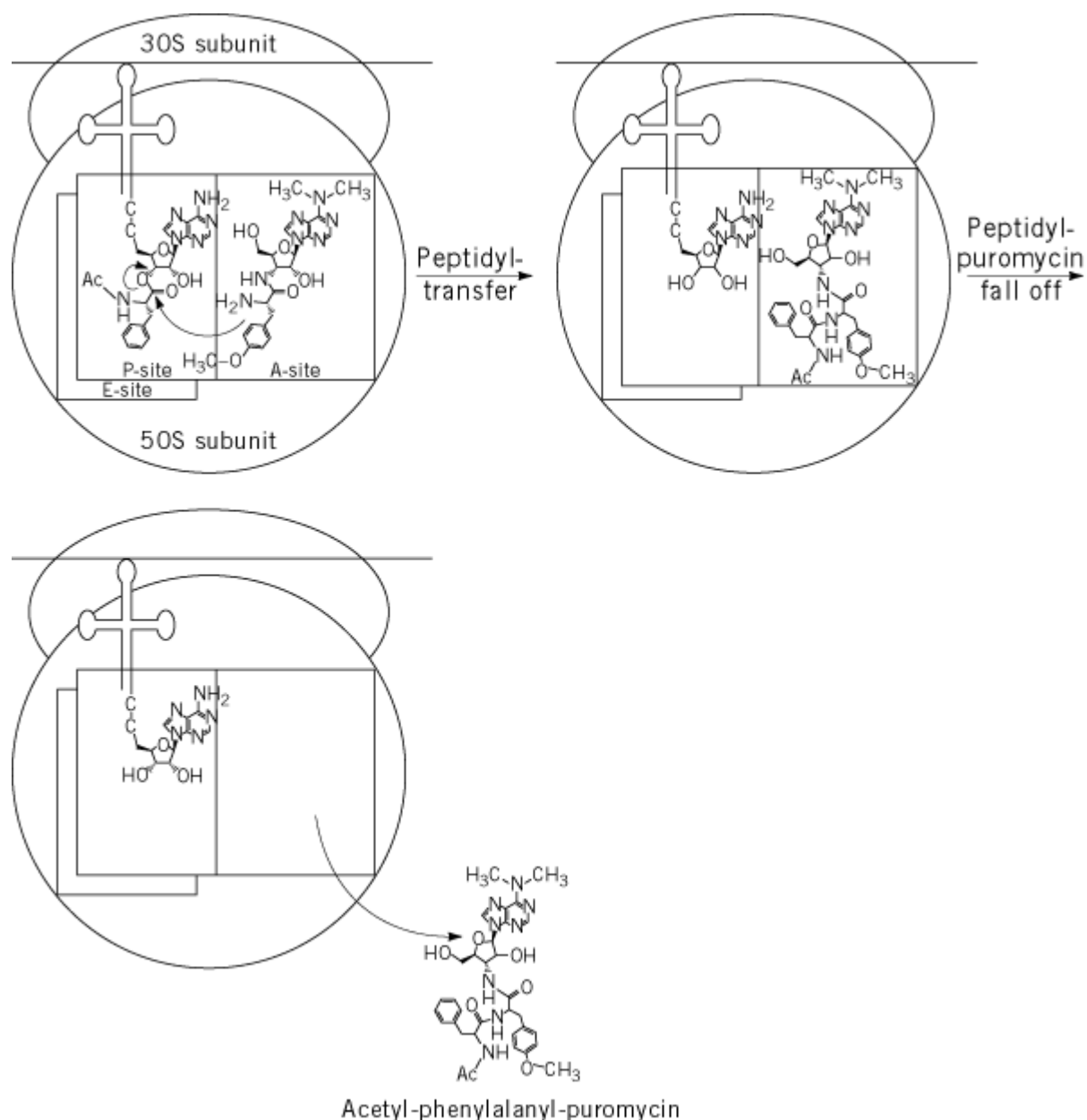


All tRNA molecules that are active in translation contain a conserved 3'-terminal sequence ending in an adenosine (A; see [Transfer RNA](#)). Similarly, puromycin retains the final A, although in a slightly modified form. Puromycin analogues having other nucleoside residues have little (I,C) or no (U,G) activity ([4](#)). The nucleoside and amino acid moieties of puromycin are linked by an amide bridge, rather than by an ester bridge as occurs in aminoacyl-tRNA. Since the amide bond is more stable than the ester bond, it cannot be cleaved by the ribosome.

#### 4. Mechanism of Action

Puromycin blocks protein synthesis as follows (Fig. [2](#)): first, puromycin binds to the A site of the peptidyl-transferase center of the ribosome in the pre-translocational state, where the P site is occupied by peptidyl-tRNA. Peptidyl-transferase then covalently links the antibiotic to the peptide residues of the peptidyl-tRNA in the P site. Because the ribosome cannot cleave the amide bond of puromycin and because peptidyl-puromycin, which lacks the whole tRNA body, has a lower affinity for the ribosome than peptidyl-tRNA, peptidyl-puromycin falls off the ribosome. High puromycin concentrations are needed to inhibit translation completely, because (i) the binding of puromycin to the ribosome is weak ( $K_{\text{diss}} = 3 \times 10^{-4} M$ ) ([5](#)), (ii) one ribosome can transfer several puromycin molecules to peptidyl-puromycin, and (iii) once peptidyl-puromycin has fallen off the ribosome, it does not bind again and has no further antibacterial activity. The degradation of **polysomes** observed with puromycin *in vivo* and *in vitro* is caused by the release of peptidyl-puromycin from the ribosome, followed by ribosome subunit dissociation.

**Figure 2.** Puromycin reaction with P-site bound peptidyl-tRNA, eg, Ac-Phe-tRNA as the peptidyl-tRNA. See text for details.



The puromycin reaction (6) is a classic *in vitro* experiment that has been used to explore ribosome function and to investigate the intrinsic peptidyl-transferase catalytic activity of the ribosome. This reaction, in which ribosomes are combined with puromycin and peptidyl-tRNA or its analogues (eg, Ac-Phe-tRNA), provided the first demonstration that ribosomes could catalyze the formation of peptide bonds in the absence of cytoplasmic factors and without an added energy source, such as GTP.

In the fragment reaction (7), 50S ribosomal subunits are combined with puromycin and a fragment of the tRNA molecule (CACCA-Leu-Ac) in the presence of 33% ethanol. Under these conditions, peptidyl-puromycin is still formed, demonstrating that the entire tRNA molecule is not necessary for peptidyl transfer. The fragment reaction also demonstrated for the first time that the entire ribosome is not required for peptidyl transfer and that the peptidyl-transferase center is located on the 50S subunit. Additionally, puromycin is used to investigate the inhibition mechanism of different translation inhibitors [eg, oxazolidinones (17) and azithromycin (18)].

The classic definitions of the tRNA binding sites on the ribosome are based on the ability of

peptidyl-tRNA to react with puromycin: peptidyl-tRNA bound to the P site reacts with puromycin, whereas peptidyl-tRNA bound to the A site does not. The third tRNA binding site, the E site, is specific for deacylated tRNA, which does not undergo the puromycin reaction. Consequently, puromycin cannot be used to characterize this tRNA-binding site.

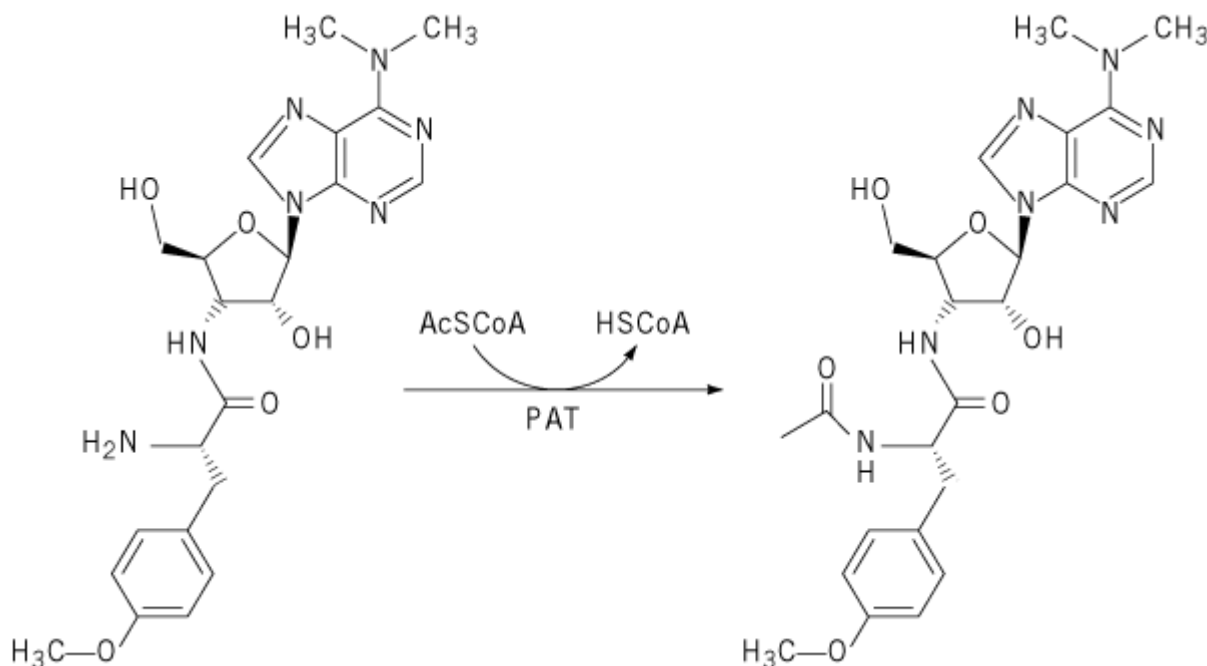
## 5. Ribosomal Binding Site

Our knowledge of the topography of the peptidyl-transferase center of the ribosome is largely derived from experiments with puromycin. In *E. coli* [crosslinking](#) experiments, a large number of proteins on the 50 ribosomal subunit, including L15, L18, L23, and L29, were labeled by puromycin and p-azidopuromycin ([8-10](#)). Most of these proteins are located in or next to the peptidyl-transferase center of the ribosome, and several can be **affinity-labeled** with other antibiotics that interact with this center (**erythromycin**, lincosamides, streptogramins, [chloramphenicol](#)). In contrast to the other peptidyl-transferase inhibitors, puromycin crosslinks proteins of both the large ribosomal subunit and the small subunit (eg, S7, S14, S18) ([8-10](#)). This indicates that the peptidyl-transferase center of the ribosome is located on the 50S subunit next to the 30S/50S subunit interface. P-azidopuromycin also cross-links nucleotides G2502 and G2504 of the central loop of 23 S rRNA ([11](#)) (see [Peptidyl Transferase](#), [Chloramphenicol](#)); thus, this loop may be involved, directly or indirectly, in the peptidyl-transferase reaction ([12](#)).

## 6. Resistance

Puromycin can be inactivated by puromycin-*N*-acetyltransferase (PAT), which is found in the producing organism, *Streptomyces alboniger* ([13](#)). The resulting *N*-acetyl puromycin (Fig. [3](#)) cannot bind to ribosomes and fails to inhibit translation. Puromycin-*N*-acetyltransferase is commonly used as a dominant selectable marker for genetic studies of mammalian cell lines ([14](#)); the gene encoding this enzyme is used as a [reporter gene](#) in gene expression experiments ([15](#)).

**Figure 3.** Inactivation of puromycin by puromycin-*N*-acetyltransferase (PAT).



In other organisms, efflux pumps have been identified that export puromycin from the bacterial cell.

The genes encoding such pumps (eg, *acrAB* and *envCD* in *E. coli* and *bmr* and *blt* in *Bacillus subtilis*) confer a multiple-[drug resistance](#) phenotype (16).

So far, no resistance mechanism has been discovered that operates by ribosomal modification. This is not surprising, in light of the similarity between puromycin and tRNA: ribosomes that could not bind puromycin would probably be unable to bind tRNA and therefore be unable to synthesize protein.

### Bibliography

1. R. Amils et al. (1990) in *The Ribosome: Structure, Function, and Evolution*, W. E. Hill et al., eds., American Society for Microbiology, Washington, DC, p. 645.
2. R. J. Harris et al. (1971) *Biochim. Biophys. Acta* **240**, 244.
3. J. P. Waller et al. (1966) *Biochim. Biophys. Acta* **119**, 566.
4. I. Rychlick et al. (1969) *J. Mol. Biol.* **43**, 13.
5. R. Fernandez-Munoz and D. Vazquez (1973) *Mol. Biol. Rep.* **1**, 27.
6. R. R. Traut and R. E. Monro (1964) *J. Mol. Biol.* **10**, 63.
7. R. E. Monro (1971) *Meth. Enzymol.* **20**, 472.
8. B. S. Cooperman et al. (1990) in *The Ribosome: Structure, Function and Evolution*, W. E. Hill et al. eds., American Society for Microbiology, Washington, DC, p. 491.
9. F. Krassnigg et al. (1978) *Eur. J. Biochem.* **87**, 439.
10. B. Wittmann-Liebold et al. (1995) *Biochem. Cell Biol.* **73**, 1187.
11. G. Steiner et al. (1988) *EMBO J.* **7**, 3949.
12. I. Nitta et al. (1998) *Science* **281**, 666.
13. M. Sugiyama et al. (1985) *J. Gen. Microbiol.* **131**, 1999.
14. M. Taniguchi et al. (1998) *Nucleic Acids Res.* **26**, 679.
15. C. Milke et al. (1995) *TIG* **11**, 258.
16. H. Nikaido (1994) *Science* **264**, 382.
17. Patel U, Yan YP, Hobbs J, Kaczmarczyk J, Slee AM, Pompiano DL, Kurilla, MG, Bobkova EV (2001) *J. Biol. Chem.* **276** 37199–37205.
18. Dinos GP, Michelinaki M, Kalpaxis DL (2001) *Mol. Pharmacol.* **59**, 1441–1445.

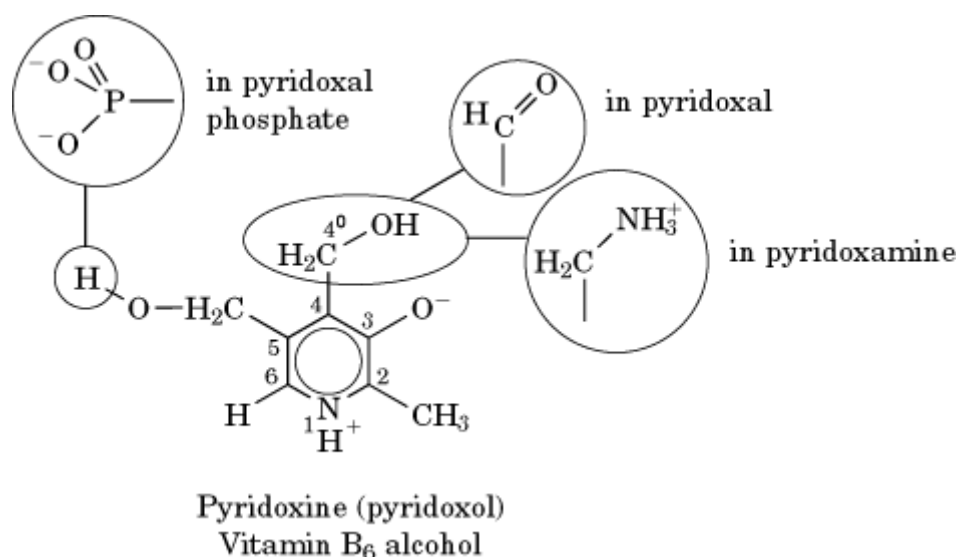
### Suggestion for Further Reading

19. R. E. Monro (1971) Ribosomal peptidyltransferase: the fragment reaction, *Meth. Enzymol.* **20**, 472–481.

## Pyridoxal Phosphate

Pyridoxal 5'-phosphate (PLP), a phosphate ester of the aldehyde form of vitamin B<sub>6</sub>, is an essential **coenzyme** for numerous [Enzymes](#) (Table 1), many of which function in all living cells. PLP is a rather simple derivative of 3-hydroxypyridine. Its structure and that of related forms of vitamin B<sub>6</sub> are as follows: (Structure 1) Within the human body, dietary vitamin B<sub>6</sub> is converted into the coenzyme forms PLP and pyridoxamine 5'-phosphate (PMP). For all these compounds, the dipolar ionic ring structure shown predominates in solution and is also found in most enzymes.





**Table 1. Groups of Enzymes Utilizing Pyridoxal or Pyridoxamine Phosphate as a Coenzyme<sup>a</sup>**

*Removal of alpha hydrogen as H<sup>+</sup>*

Aminotransferases

Aspartate aminotransferase ([14](#), [15](#))

D-Amino acid aminotransferase ([16](#))

Branched chain aminotransferase ([17](#))

Gamma-aminobutyrate aminotransferase ([18](#))

Serine:pyruvate aminotransferase ([19](#))

Alanine racemase ([20](#))

Aminocyclopropane carboxylate synthase ([21](#))

2-Amino-3-oxobutyrates-CoA ligase ([22](#)) AKB Synthase

*Beta elimination and replacement reactions*

L- and D-Serine dehydratases (deaminases) ([23](#), [24](#))

Tyrosine phenol-lyase ([25](#))

Alliinase ([26](#))

Cystathionine b-lyase (cystathionase) ([27](#))

O-Acetylserine sulfhydrylase (cysteine synthase) ([28](#))

Tryptophan synthase ([29-32](#))

*Removal of alpha carboxylate as CO<sub>2</sub>* ([33](#))

Diaminopimelate decarboxylase ([34](#))

Glycine decarboxylase (requires lipoyl group) ([35](#))

Glutamate decarboxylase ([36](#), [37](#))

DOPA decarboxylase ([38](#))

Dialkylglycine decarboxylase ([39](#))

*Removal of side chain by aldol cleavage*

Serine hydroxymethyltransferase ([40](#))

### Reactions of ketimine intermediates

Threonine synthase (3)

### Other enzymes

Lysine 2,3-aminomutase (41)

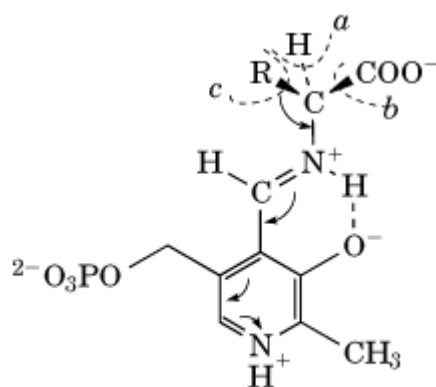
Glycogen phosphorylase (42)

Pyridoxamine phosphate (PMP) in synthesis of 3,6-dideoxy-L-arabino-hexose (43)

---

<sup>a</sup> Names of a few representative enzymes are also listed.

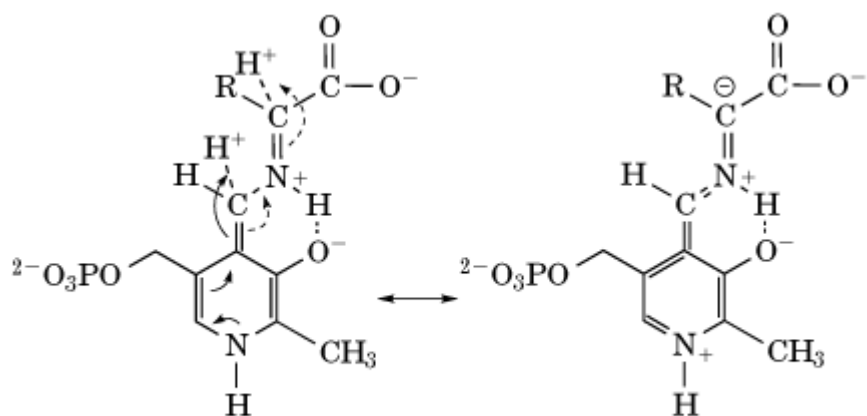
Within the active sites of most PLP-dependent enzymes, the carbonyl group of the coenzyme forms a **Schiff Base** with an amino group of a specific **lysine** side chain. This Schiff base is commonly called the *internal aldimine*. When an **amino acid** substrate binds into the active site, a series of rapid reactions occur, many of which have been characterized spectroscopically and kinetically (1-6). The **amino group** of the substrate is deprotonated and then adds to the carbon atom of the C=N Schiff base group of the internal aldimine, forming an intermediate geminal diamine. After another proton transfer, elimination of the lysyl side-chain amino group leads to a stable **hydrogen-bonded** Schiff base with the substrate, the external aldimine: (Structure 2) In this structure, the protonated ring nitrogen provides a powerful electron-accepting center that can assist in breaking any one of the



Schiff base

three bonds marked *a*, *b*, and *c* on Structure (2), as indicated by the curved arrows (7). Which type of reaction occurs is determined by the structure of the enzyme **active site** and the way in which the amino acid part of the Schiff base is held. The bond to be cleaved must be nearly perpendicular to the plane of the coenzyme ring and to the conjugated p-electron system of the Schiff base. (8) Side-chain cleavage (*c*) occurs only for  $\beta$ -hydroxyamino acids and related compounds.

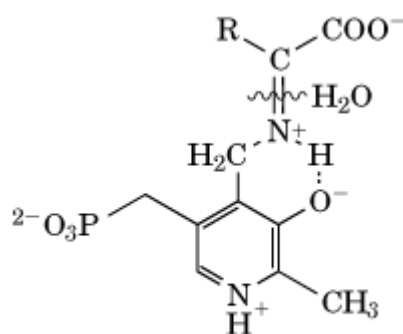
The product of the bond cleavage has been called the *quinonoid* or *carbanionic* intermediate, and the product of removal of  $H^+$  by a catalytic base



Quinonoid-carbanionic structure

is illustrated with cleavage of bond a: (Structure 3) The two resonance structures shown illustrate the quinonoid and carbanionic attributes. The carbanionic structure is stabilized by the adjacent positively charged hydrogen-bonded proton. Another resonance form has the negative charge on the 4' carbon rather than the a carbon. It is possible that the proton on the Schiff base nitrogen may be shifted onto the phenolate oxygen of the ring in this intermediate (9). There are still uncertainties about the details of the reaction mechanism.

The quinonoid-carbanionic form shown in Structure (3) can undergo reprotonation at either Ca or C4'. Protonation at Ca, but without stereospecificity, leads to racemization of the amino acid (see [Stereoisomers](#)) and is observed for bacterial alanine and glutamic acid racemases. Transaminases (aminotransferases) catalyze protonation at either Ca (with retention of the configuration) or C4'. The initial product in the latter case is another Schiff base called the *ketimine*: (Structure 4) Its hydrolysis produces pyridoxamine phosphate (PMP), the product of a half-reaction of a transaminase.



Ketimine

Removal of the  $\alpha$ -hydrogen is also a first step in a number of other reactions, such as synthesis in plants of the ethylene precursor aminocyclohexane carboxylate and numerous reactions involving elimination of the b or g substituent. For example, the bacterial tryptophan indole-lyase releases indole from [tryptophan](#), the remainder of the molecule being converted to pyruvate and ammonia. The g-elimination reactions are more complex and pass through a ketimine step before the elimination.

Removal of the  $\alpha$ -carboxylate as carbon dioxide is characteristic of a large number of decarboxylases, including glutamate and dihydroxyphenylalanine (DOPA) decarboxylases of the brain. Cleavage of the side chains is of more limited occurrence but accomplishes another series of important metabolic reactions. Best known of these is cleavage of [serine](#) by serine hydroxymethylase to give [glycine](#) and formaldehyde. The latter does not leave the active site but is captured by

tetrahydrofolate as methylenetetrahydrofolate.

There are many medical and technological aspects to the study of PLP. Amino acid racemases are targets for antibacterial drugs, and several enzymes are targets for herbicides and pesticides. Enzyme deficiency diseases are known. Poor removal of homocysteine leads to homocysteinemia, which is linked to atherosclerosis (10, 11). PLP has some blocking effect on the virus HIV (12). PLP is widely used as a labeling reagent for amino groups in proteins by reduction of Schiff bases formed with borohydride (13), which may be **radiolabeled** with deuterium or tritium.

## Bibliography

1. C. M. Metzler, J. Mitra, D. E. Metzler, M. W. Makinen, C. C. Hyde, P. Rogers, and A. Arnone (1988) *J. Mol. Biol.* **203**, 197–220.
2. P. A. Clark, J. N. Jansonius, and E. L. Mehler (1993) *J. Am. Chem. Soc.* **115**, 1894–1902.
3. B. Laber, K.-P. Gerbling, C. Harde, K.-H. Neff, E. Nordhoff, and H.-D. Pohlenz (1994) *Biochemistry* **33**, 3413–3423.
4. C.-C. Hwang, E. U. Woehl, D. E. Minter, M. F. Dunn, and P. F. Cook (1996) *Biochemistry* **35**, 6358–6365.
5. H. Hayashi and H. Kagamiyama (1995) *Biochemistry* **34**, 9413–9423.
6. A. Graf von Stosch (1996) *Biochemistry* **35**, 15260–15268.
7. D. E. Metzler, M. Ikawa, and E. E. Snell (1954) *J. Am. Chem. Soc.* **76**, 648–652.
8. H. C. Dunathan (1971) *Adv. Enzymology* **35**, 79–134.
9. C. M. Metzler, A. G. Harris, and D. E. Metzler (1988) *Biochemistry* **27**, 4923–4933.
10. R. Clarke, L. Daly, K. Robinson, E. Naughten, S. Cahalane, B. Fowler, and I. Graham (1991) *New Engl. J. Med.* **324**, 1149–1155.
11. M. Watanabe, J. Osada, Y. Aratani, K. Kluckman, R. Reddick, M. R. Malinow, and N. Maeda (1995) *Proc. Natl. Acad. Sci. USA* **92**, 1585–1589.
12. J. M. Salhany and L. M. Schopfer (1993) *J. Biol. Chem.* **268**, 7643–7645.
13. B. Perez-Ramirez, A. Iriarte, and M. Martinez-Carrion (1994) *J. Protein Chem.* **13**, 67–76.
14. J. M. Goldberg and J. F. Kirsch (1996) *Biochemistry* **35**, 5280–5291.
15. V. Malashkevich, B. Strokopytov, V. Borisov, Z. Dauter, K. Wilson, and Y. Torchinsky (1995) *J. Mol. Biol.* **247**, 111–124.
16. S. Sugio, G. A. Petsko, J. M. Manning, K. Soda, and D. Ringe (1995) *Biochemistry* **34**, 9661–9669.
17. T. R. Hall, R. Wallin, G. D. Reinhart, and S. M. Hutson (1993) *J. Biol. Chem.* **268**, 3092–3098.
18. M. D. Toney, S. Pascarella, and D. De Biase (1995) *Protein Sci.* **4**, 2366–2374.
19. C. Uchida, T. Funai, T. Oda, K. Ohbayashi, and A. Ichiyama (1994) *J. Biol. Chem.* **269**, 8849–8856.
20. K. Hoffmann, E. Schneider-Scherzer, H. Kleinkauf, and R. Zocher (1994) *J. Biol. Chem.* **269**, 12710–12714.
21. E. Hohenester, M. F. White, J. F. Kirsch, and J. N. Jansonius (1994) *J. Mol. Biol.* **243**, 947–949.
22. H. Tong and L. Davis (1995) *Biochemistry* **34**, 3362–3367.
23. H. Ogawa, T. Gomi, K. Konishi, T. Date, H. Nakashima, K. Nose, Y. Matsuda, C. Peraino, H. C. Pitot, and M. Fujioka (1989) *J. Biol. Chem.* **264**, 15818–15823.
24. K. D. Schnackerz, J. H. Ehrlich, W. Giesemann, and T. A. Reed (1979) *Biochemistry* **18**, 3557–3563.
25. A. A. Antson, T. V. Demidkina, P. Gollnick, Z. Dauter, R. L. Von Tersch, J. Long, S. N. Berezchnoy, R. S. Phillips, E. H. Harutyunyan, and K. S. Wilson (1993) *Biochemistry* **32**, 4195–4206.
26. E. Block, J. Z. Gillies, C. W. Gillies, A. A. Bazzi, D. Putnam, L. K. Revelle, D. Wang, and X.

- Zhang (1996) *J. Am. Chem. Soc.* **118**, 7492–7501.
27. T. Clausen, R. Huber, B. Laber, H.-D. Pohlenz, and A. Messerschmidt (1996) *J. Mol. Biol.* **262**, 202–224.
  28. V. D. Rege, N. M. Kredich, C.-H. Tai, W. E. Karsten, K. D. Schnackerz, and P. F. Cook (1996) *Biochemistry* **35**, 13485–13493.
  29. C. C. Hyde, S. A. Ahmed, E. A. Padlan, E. W. Miles, and D. R. Davies (1988) *J. Biol. Chem.* **263**, 17857–17871.
  30. A. Peracchi, S. Bettati, A. Mozzarelli, G. L. Rossi, E. W. Miles, and M. F. Dunn (1996) *Biochemistry* **35**, 1872–1880.
  31. M. C. Yee, V. Horn, and C. Yanofsky (1996) *J. Biol. Chem.* **271**, 14754–14763.
  32. P. Pan, E. Woehl, and M. F. Dunn (1997) *Trends Biochem. Sci.* **22**, 22–27.
  33. E. Sandmeier, T. I. Hale, and P. Christen (1994) *Eur. J. Biochem.* **221**, 997–1002.
  34. J. G. Kelland, L. D. Arnold, M. M. Palcic, M. A. Pickard, and J. C. Vederas (1986) *J. Biol. Chem.* **261**, 13216–13223.
  35. A. Kume, H. Koyata, T. Sakakibara, Y. Ishiguro, S. Kure, and K. Hiraga (1991) *J. Biol. Chem.* **266**, 3323–3329.
  36. M. G. Erlander, N. J. K. Tillakaratne, S. Feldblum, N. Patel, and A. J. Tobin (1991) *Neuron* **7**, 91–100.
  37. G. Baum, Y. Chen, T. Arazi, H. Takatsuji, and H. Fromm (1993) *J. Biol. Chem.* **268**, 19610–19617.
  38. V. N. Malashkevich, P. Filippini, U. Sauder, P. Dominici, J. N. Jansonius, and C. B. Voltattorni (1992) *J. Mol. Biol.* **224**, 1167–1170.
  39. M. D. Toney, E. Hohenester, S. W. Cowan, and J. N. Jansonius (1993) *Science* **261**, 756–759.
  40. P. Stover, M. Zamora, K. Shostak, M. Gautam-Basak, and V. Schirch (1992) *J. Biol. Chem.* **267**, 17679–17687.
  41. W. Wu, K. W. Lieder, G. H. Reed, and P. A. Frey (1995) *Biochemistry* **34**, 10532–10537.
  42. N. G. Oikonomakos, S. E. Zographos, K. E. Tsitsanou, L. N. Johnson, and K. R. Acharya (1996) *Protein Sci.* **5**, 2416–2428.
  43. X. He, J. S. Thorson, and H.-w Liu (1996) *Biochemistry* **35**, 4721–4731.

### Suggestions for Further Reading

44. D. Dolphin, R. Poulson, and O. Avramovic, eds. (1986) *Vitamin B<sub>6</sub> Pyridoxal Phosphate*, Wiley, New York.
45. J. N. Jansonius and M. G. Vincent (1987) "Structural basis for catalysis by aspartate aminotransferase." In *Biological Macromolecules and Assemblies*, Vol. **3**: Active Sites of Enzymes (F. A. Jurnak and A. McPherson, eds), Wiley, New York, pp. 187–285.
46. T. Fukui, H. Kagamiyama, K. Soda, and H. Wada, eds. (1991) *Enzymes Dependent on Pyridoxal Phosphate and Other Carbonyl Compounds as Cofactors*, Pergamon Press, Oxford, U.K.
47. K. Dakshinamurti, ed. (1990) *Vitamin B<sub>6</sub>*, Vol. **585**, Annals of the New York Academy of Sciences, New York.
48. F. Hayashi, H. Wada, T. Yoshimura, N. Esaki, and K. Soda (1990) Recent topics in pyridoxal 5 $\phi$ -phosphate enzyme studies. *Ann. Rev. Biochem.* **59**, 87–110.
49. E. Sandmeier, T. I. Hale, and P. Christen (1994) Multiple evolutionary origin of pyridoxal-5 $\phi$ -phosphate-dependent amino acid decarboxylases. *Eur. J. Biochem.* **221**, 997–1002.

## Pyrophosphatase

Inorganic pyrophosphatase (PPase) was first discovered in animal tissues in 1928 (1). It hydrolyzes pyrophosphate ( $PP_i$ ) to two molecules of inorganic phosphate ( $P_i$ ). Soluble cytoplasmic PPase (s-PPase) is believed to be both ubiquitous and essential (2-4), playing a central role in cell metabolism by hydrolyzing the truly staggering quantity of  $PP_i$  produced as a by-product of the biosyntheses of nucleic acids, proteins, lipids, and polysaccharides. Klemme (5) calculated that bacterial  $PP_i$  levels would rise to 3 M in an hour in the absence of PPase, and several kilograms of  $PP_i$  are generated in adult humans per day. The required high flux of  $PP_i$  hydrolysis results from the high specific activity of s-PPase [turnover numbers at 25°C and neutral pH of  $\sim 100500\text{ s}^{-1}$  (6, 7),  $10^{10}$  faster than the uncatalyzed rate (8)] and its relative abundance, typically 0.1–0.5% of cell protein by weight.

PPase catalysis of  $PP_i$  hydrolysis drives biopolymer synthesis (9), but the high rate of  $PP_i$  production results in steady-state cellular  $PP_i$  concentrations that in some cells are well above equilibrium levels. Thus, at conditions close to physiological, the concentration of  $PP_i$  in equilibrium with 10 mM  $P_i$  is about 0.1 mM (10, 11), whereas in bacterial, plant, and yeast cells  $PP_i$  levels of the order of millimolar or greater are observed (5, 12, 13). Even in mammals, where the  $PP_i$  concentrations are generally much lower, millimolar levels have been found in blood platelets (14), starved hepatocytes (15), and hepatocytes metabolizing acetate in the presence of  $Ca^{2+}$  (16).

The steady-state levels of  $PP_i$  can be important for metabolic regulation. For example, there is a linear relationship between the  $PP_i$  concentration and the error rate in DNA polymerization (17) (see [DNA Replication](#)), and a model for  $PP_i$  regulation of the termination of [transcription](#) in *Escherichia coli* has been presented (18). A disorder of  $PP_i$  metabolism is suspected as a cause of calcium pyrophosphate dihydrate crystal deposition disease, a major arthropathy in the elderly that has also been linked to osteoarthritis (19, 20). Although as yet no linkage has been established between these diseases and a PPase abnormality, s-PPase has been shown to be effective in dissolving calcium pyrophosphate dihydrate crystals (21).

In addition to s-PPase, plants and certain bacteria have a membrane-bound PPase, structurally unrelated to s-PPase, which works as a reversible **proton pump**. In these organisms, the  $PP_i$  level is much higher, and  $PP_i$  is used as a phosphoryl and energy donor instead of, or in parallel with, ATP (22, 23). **Mitochondrial** PPase is a third kind of PPase, which, though structurally and functionally similar to s-PPase (3), contains noncatalytic subunits that anchor it to the inner mitochondrial membrane.

s-PPases have been studied from a large number of organisms (6, 22). Here we focus on the two best understood s-PPases, those of *Saccharomyces cerevisiae* (Y-PPase) and *E. coli* (E-PPase). Both of these enzymes have been **cloned** and expressed, and the structural and functional properties of wild-type enzyme, as well as many variants, have been characterized (7). Comprehensive recent reviews of E- and Y-PPases may be found in Refs. 6, 24, and 25.

### 1. Substrate Specificity

Isolated PPase requires added divalent metal ion for activity. Only four of many divalent ions tested

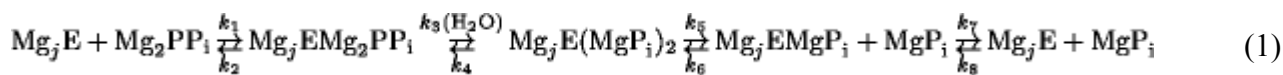
have been shown to confer PPase with appreciable ( $\geq 5\%$  of maximal)  $PP_i$  hydrolysis activity, with relative activities in the order  $Mg^{2+} > Zn^{2+} > Co^{2+} \approx Mn^{2+}$  (26). E-PPase shows similar specificity toward divalent metal ions (27). For Y-PPase, the lower effectiveness of  $Zn^{2+}$  and  $Co^{2+}$  as cofactors for  $PP_i$  hydrolysis than  $Mg^{2+}$  is due mainly to the slower rates of product  $P_i$  release from enzyme in the presence of each of these ions.

With  $Mg^{2+}$  as cofactor, the reaction is almost totally specific for  $PP_i$  as substrate; in contrast, with  $Zn^{2+}$  and/or  $Mn^{2+}$ , other pyrophosphates, such as ATP and serial pyrophosphate, are hydrolyzed at rates within a factor of 10 of that of  $PP_i$  itself (28-30).  $CaPP_i$  is not a substrate, but is a potent competitive inhibitor of PPase, with a  $K_i$  considerably lower than the  $K_m$  ( **Michaelis constant**) of  $Mg PP_i$  under corresponding conditions (31). In the presence of  $Mg^{2+}$ , high specificity is also evident for the back reaction from  $P_i$  to synthesize enzyme-bound  $PP_i$ . Thus, enzyme catalysis of thiophosphate, phosphoramidate, or fluorophosphate hydrolysis, which would be analogous to  $P_i-H_2O$  oxygen exchange (Eq. (1)), proceeds very slowly if at all. Fluoride shows very potent inhibition of Y-PPase (but not E-PPase), by formation of a stable, isolatable complex:  $PPase \cdot (Mg^{2+})_2 \cdot PP_i \cdot F^-$  (32).

Paralleling its substrate specificity, PPase displays a high binding specificity for its natural ligands,  $PP_i$  and  $P_i$ . The two  $P_i$  sites per PPase subunit differ in affinity. Both the higher-affinity (P1) and the lower-affinity (P2) sites bind  $P_i$  more tightly as a function of added metal ion (33). Some changes are tolerated for the atom bridging the phosphoryl groups. In the presence of  $Mg^{2+}$ , both  $O_3PCHOHPO_3$  and  $O_3PNHPO_3$  are reasonably good **competitive inhibitors**, although  $O_3PCH_2PO_3$  is not. With  $Mn^{2+}$ , the specificity is somewhat relaxed:  $PPP_i$  has a  $K_m$  similar to that of  $PP_i$ , and  $O_3PCHOHPO_3$  binds quite tightly. The binding of analogues of  $P_i$  has been examined only in the presence of  $Mn^{2+}$ . None of the analogues tested, methyl phosphonate, phosphoramidate, or thiophosphate, bind comparably to  $P_i$  in site P1, but methyl phosphonate binds comparably to  $P_i$  in site P2. Thus, site P1 appears the primary source of binding specificity (33).

## 2. Kinetic Mechanism

Y-PPase and E-PPase share a common mechanism for catalysis of  $PP_i$  hydrolysis,  $P_i-H_2O$  oxygen exchange, and  $PP_i:Pi$  exchange (34, 35).  $PP_i$  hydrolysis proceeds via single-step direct phosphoryl transfer to water, without formation of a phosphorylated enzyme intermediate (Eq. (1)):



$$j = 1, 2 : K_1 = k_1/k_2; K_3 = k_3/k_4; K_5 = k_5/k_6; K_7 = k_7/k_8$$

The outstanding features of Equation 1 are as follows: (a) All three forward rate constants,  $k_3$ ,  $k_5$ , and  $k_7$ , are partially rate-determining for  $PP_i$  hydrolysis; (b) synthesis of enzyme-bound  $PP_i$  from product  $P_i$  proceeds quite readily, but release of bound  $PP_i$  (step 2) is very slow; (c) the equilibrium

constant for hydrolysis of enzyme-bound  $PP_i$ ,  $K_3$ , is only 4–5, in contrast to much higher value for the hydrolysis of  $PP_i$  in solution; and (d) the first  $P_i$  released (step 5) is the electrophilic  $P_i$ , containing an oxygen atom from the nucleophilic water (34). The value of  $k_{cat}$  goes through a maximum as a function of pH (at pH 8 for E-PPase, 7 for Y-PPase), implying that both acidic and basic groups are important for activity (7, 36). The value of  $k_{cat}/K_m$  also displays a pH optimum near neutrality for each enzyme.

### 3. Structural Studies on PPases

Y-PPase is a homodimer (286 amino acid residues per subunit), typical of eukaryotic PPases (37), and E-PPase is typical of prokaryotic PPases in being a homo-hexamer (175 residues each) (38). Although the enzymes have only about 27% sequence identity (6), they have a very similar [protein structure](#), consisting of eight **beta-strands** and two **alpha-helices** (39) and belonging to the [nucleotide-binding motif](#). The E-PPase subunit may be thought of as embedded in the Y-PPase subunit, and the common parts of the two molecules are best described as a distorted, highly twisted five-stranded b-barrel with four excursions, the latter creating a very large [active site](#). About 14 polar active-site residues, highly conserved in s-PPases, are significant for catalysis (7), and 12 of them are in the excursions. In addition to E-PPase and Y-PPase, more than a dozen soluble PPases have been **cloned** and sequenced (6, 24).

High-resolution structures have been determined for the  $Mg_{1.5}$  (39) and  $Mg_{2.5}$  (40) complexes of E-PPase and for the  $Mn_2$  and  $Mn_4(P_i)_2$  complexes of Y-PPase (25, 41), where the numbers refer to the stoichiometries of ligand per monomer. These ligands are all bound at the active site, with the exception of the half  $Mg^{2+}$  site in E-PPase. The latter reflects  $Mg^{2+}$  binding at the trimer:trimer interface, which stabilizes the enzyme but otherwise has little effect on enzyme activity (42). In these structures, the E-PPase active site residues are virtually superposable on those of Y-PPase, demonstrating the highly conserved nature of the PPase active site. The 2.2-Å resolution structure of the product  $Mn_4(P_i)_2$  complex of Y-PPase (25) has been the most useful for considerations of mechanism (see text below).

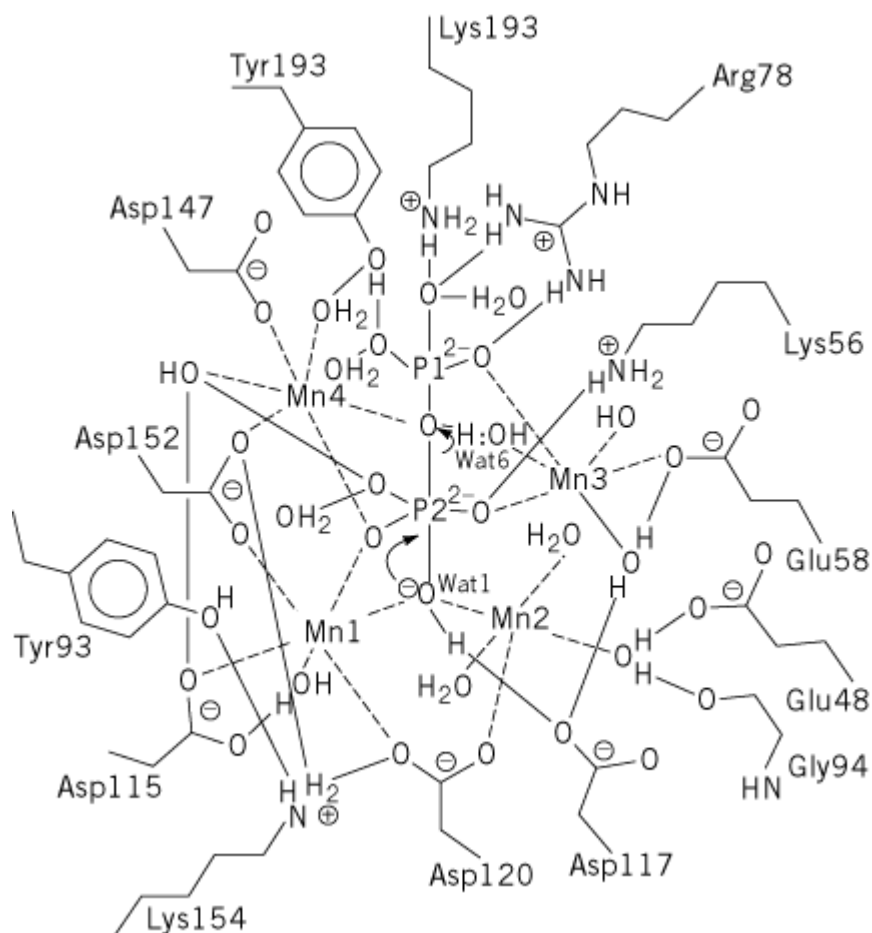
### 4. Chemical Mechanism

PPase is presently unique among phosphoryl transfer enzymes in having four divalent metal ions at the active site, although it is probable that not more than three of these are required for hydrolysis (43). The structure of the  $Mn_4(P_i)_2$  complex of Y-PPase, along with functional studies of wild-type and variant PPases, allows formulation of a detailed proposal for the catalytic mechanism (Fig. 1). In the proposed mechanism, PPase catalyzes  $PP_i$  hydrolysis by lowering the  $pK_a$  of the leaving group in site P1 (studies of phosphoryl transfer reactions in solution show that the rates of such reactions increase by 10- to 20-fold for each unit lowering of the leaving group  $pK_a$ ), by forming an incipient hydroxide ion that functions as a stronger nucleophile than water, and by shielding the charge on the electrophilic phosphorus in site P2, thus permitting attack by the hydroxide anion. According to this mechanism, the leaving phosphoryl group is activated by coordination to Arg78, Tyr192, Lys193, Met3, Met4, and Water6, the  $pK_a$  of the nucleophilic water is lowered by coordination to Met1, Met2, and Asp117, and the charge on the electrophilic phosphoryl group is neutralized through coordination to Lys56, Tyr93, and Met1 to Met4. The extraordinary array of Lewis acids (four divalent metal ions; Lys, Arg, and Tyr side chains, bound waters) within the active site activate  $PP_i$  hydrolysis through coordination to almost all of the available electron pairs in the oxygens of bound  $PP_i$ . The essential base implicated by both steady-state and pre-steady-state pH-dependent studies is assigned to the presumptive nucleophilic water, Wat1, which has an abnormally low  $pK_a$  (~ 6) by virtue of its coordination to Met1, Met2, and Asp117. These same studies implicate the essential acid



in a step following  $PP_i$  hydrolysis, which is believed to be protonation of the phosphate in site P1 prior to its dissociation from enzyme (44).

**Figure 1.** Mechanism and active-site structure of Y-PPase as deduced from the crystal structure of the  $Mn_4(P_i)_2$ -Y-PPase complex. From Ref. 25, with permission.



Interestingly, the proposed PPase mechanism has much more in common with those proposed for inorganic phosphatase and **DNA polymerases** and exonucleases than with that proposed for **ATPase**, which catalyzes a chemically more analogous reaction (25). Studies of the crystal structures and functional properties of critical active-site variants, such as Arg78Lys and Asp117Glu, that test various aspects of the proposal in Figure 1, are underway.

#### Bibliography

1. H. D. Kay (1928) *Biochem. J.* **22**, 1446–1448.
2. J. Chen, A. Brevet, M. Fromant, F. Leveque, J.-M. Schmitter, S. Blanquet, and P. Plateau (1990) *J. Bacteriol.* **172**, 5686–5689.
3. M. Lundin, H. Baltscheffsky, and H. Ronne (1991) *J. Biol. Chem.* **266**, 12168–12172.
4. U. Sonnewald (1992) *Plant J.* **2**, 571–581.
5. J. Klemme (1976) *Z. Naturforsch. Teil C.* **31**, 544–550.
6. A. A. Baykov, B. S. Cooperman, A. Goldman, and R. Lahti (1998) In *Inorganic Polyphosphates* (H. C. Schröder, ed.), Springer-Verlag, Berlin.

7. P. Pohjanjoki, R. Lahti, G. Belgurov, A. A. Baykov, A. Goldman, and B. S. Cooperman (1998) *Biochemistry* **37**, 1754–1761.
8. B. S. Cooperman (1982) *Methods Enzymol.* **87**, 526–548.
9. A. Kornberg (1962) In *Horizons in biochemistry* (M. Kasha and D. Pullman, eds.), Academic Press, New York, pp. 251–264.
10. H. Flodgaard and P. Fleron (1974) *J. Biol. Chem.* **249**, 3465–3474.
11. L. de Meis (1984) *J. Biol. Chem.* **259**, 6090–6097.
12. C. C. Black Jr., D.-P. Xu, S. S. Sung, L. Mustardy, N. Paz, and P. P. Kormanik (1987) In *Phosphate Metabolism and Cellular Regulation in Microorganisms*, (A. Torriani-Gorini, F. G. Rothman, S. Silver, A. Wright, and E. Yagil, eds.), American Society of Microbiology, Washington, D.C., pp. 264–268.
13. S. A. Ermakova, S. E. Mansurova, T. S. Kalebina, E. S. Lobakova, I. O. Selyach, and I. S. Kulaev (1981) *Arch. Microbiol.* **128**, 394–397.
14. M. Fukami, C. Dangelmaier, J. Bauer, and H. Homsen (1980) *Biochem. J.* **192**, 99–105.
15. A. Davidson and A. Halestrap (1988) *Biochem. J.* **254**, 379–384.
16. R. L. Veech and W. L. Gitomer (1988) *Adv. Enzyme Regul.* **27**, 313–343.
17. P. Herbomel and J. Ninio (1980) *C. R. Acad. Sci., Ser. D* **291**, 881–884.
18. R. B. Kent and S. K. Guterman (1982) *Proc. Natl. Acad. Sci. USA* **79**, 3992–3996.
19. L. M. Ryan, R. L. Wortmann, B. Karas, M. P. Lynch, and D. J. McCarty, (1986) *J. Clin. Invest.* **7**, 1689–1693.
20. C. T. Baldwin et al. (1995) *Am. J. Hum. Genet.* **56**, 692–697.
21. Y. Xu, T. Cruz, P.-T. Cheng, and K. P. H. Pritzker (1991) *J. Rheumatol.* **18**, 66–71.
22. M. Baltscheffsky and H. Baltscheffsky (1992) In *Molecular Mechanisms in Bioenergetics* (L. Ernster, ed.), Elsevier, Amsterdam, pp. 331–348.
23. P. A. Rea and R. J. Poole (1993) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **44**, 157–180.
24. B. S. Cooperman, A. A. Baykov, and R. Lahti (1992) *TIBS* **17**, 262–266.
25. P. Heikinheimo, J. Lehtonen, A. A. Baykov, R. Lahti, B. S. Cooperman, and A. Goldman (1996) *Structure* **4**, 1491–1508.
26. K. M. Welsh, A. Jacobyansky, B. Springs, and B. S. Cooperman, (1983) *Biochemistry* **22**, 2243–2248.
27. J. Josse and S. C. K. Wong (1971) In *The Enzymes*, 3rd ed, Vol **4** (P. D. Boyered, ed.), Academic Press, New York, pp. 499–527.
28. W. E. Höhne and P. Heitmann (1974) *Acta Biol. Med. Ger.* **33**, 1–14.
29. M. J. Schlesinger and M. J. Coon (1960) *Biochim. Biophys. Acta.*, **41**, 30–36.
30. S. M. Avaeva, S. N. Kara-Murza, and M. M. Botvinik (1967) *Biokhimiya* **32**, 205–209.
31. O. A. Moe and L. G. Butler (1972) *J. Biol. Chem.* **247**, 7315–7319.
32. A. A. Baykov, A. A. Artjerkov, and S. M. Avaeva, (1976) *Biochim. Biophys. Acta* **429**, 982–992.
33. B. S. Cooperman, A. Panackal, B. Springs, and D. J. Hamm (1981) *Biochemistry* **20**, 6051–6060.
34. B. Springs, K. M. Welsh, and B. S. Cooperman (1981) *Biochemistry* **20**, 6384–6391.
35. A. A. Baykov, A. S. Shestakov, V. N. Kasho, A. V. Vener, and A. H. Ivanov (1990) *Eur. J. Biochem.* **194**, 879–887.
36. T. Salminen, J. Käpylä, P. Heikinheimo, A. Goldman, J. Heinonen, A. A. Baykov, B. S. Cooperman, and R. Lahti (1995) *Biochemistry* **34**, 782–791.
37. L. Kolakowski, M. Schlösser, and B. S. Cooperman (1988) *Nucleic Acids Res.* **16**, 10441–10452.

38. R. Lahti, T. Pitkaranta, E. Valve, I. Ilta, E. Kukko-Kalske, and J. Heinonen (1988) *J. Bacteriol.* **170**, 5901–5907.
39. J. Kankare, T. Salminen, R. Lahti, B. S. Cooperman, A. A. Baykov, A. Goldman (1996) *Biochemistry* **35**, 4670–4677.
40. E. H. Harutyunyan, V. Y. Oganessyan, N. N. Oganessyan, S. M. Avaeva, T. I. Nazarova, N. N. Vorobyeva, S. A. Kurilova, R. Huber, and T. Mather (1997) *Biochemistry* **36**, 7754–7760.
41. E. H. Harutyunyan, I. P. Kuranova, B. K. Vainshtein, W. E. Höhne, V. S. Lamzin, Z. Dauter, A. V. Teplyakov, and K. S. Wilson (1996) *Eur. J. Biochem.* **239**, 220–228.
42. I. S. Efimova, A. Salminen, P. Pohjanjoki, J. Lapinniemi, N. N. Magretova, B. S. Cooperman, A. Goldman, R. Lahti, and A. A. Baykov (1999) *J. Biol. Chem.* **274**, 3294–3299.
43. A. A. Baykov, T. Hyytiä, S. E. Volk, V. N. Kasho, A. V. Vener, A. Goldman, R. Lahti, and B. S. Cooperman (1996) *Biochemistry* **35**, 4655–4661.
44. P. Halonen and B. S. Cooperman, in preparation.

### Suggestions for Further Reading

45.

### Mechanisms of Phosphoryl Transfer

46. J. R. Knowles (1980) Enzyme-catalyzed phosphoryl transfer reactions. *Annu. Rev. Biochem.* **49**, 877–919.
47. K. A. Maegley, S. J. Admirael, and D. Herschlag (1996) Ras-catalyzed hydrolysis of GTP: a new perspective from model studies. *Proc. Natl. Acad. Sci. USA* **93**, 8160–8166.
48. E. E. Kim and H. W. Wyckoff (1991) Reaction mechanism of alkaline phosphatase based on crystal structures. *J. Mol. Biol.* **218**, 449–464.
49. L. S. Beese and T. A. Steitz (1991) Structural basis for the 3 $\phi$ -5 $\phi$  exonuclease activity of *E. coli* DNA polymerase I: a two metal ion mechanism. *EMBO J.* **10**, 25–33.

### PP<sub>i</sub> Metabolism

50. H. G. Wood and J. E. Clark (1988) Biological aspects of inorganic polyphosphates. *Annu. Rev. Biochem.* **57**, 235–260.
51. I. S. Kulaev and V. M. Vagabov (1983) Polyphosphate metabolism in micro-organisms. *Adv. in Microbiol. and Physiol.* **24**, 83–171.

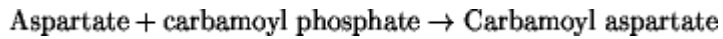
## Pyrimidine Ribonucleotide Metabolism

### 1. Biosynthesis of the Pyrimidine Ring

The biosynthesis of pyrimidine **nucleotides** is comparable to that of **purine nucleotides** in several important respects. Both rings are synthesized by *de novo* pathways that begin with amino acids and their metabolites as precursors. The *de novo* pathway to pyrimidines is universal among organisms studied thus far, as is the purine pathway. Like purines, pyrimidine bases and nucleosides can be used in “**salvage**” pathways that involve reutilization of components released in nucleic acid degradation. Salvage capacities vary considerably among organisms and among different cell types in a single organism.

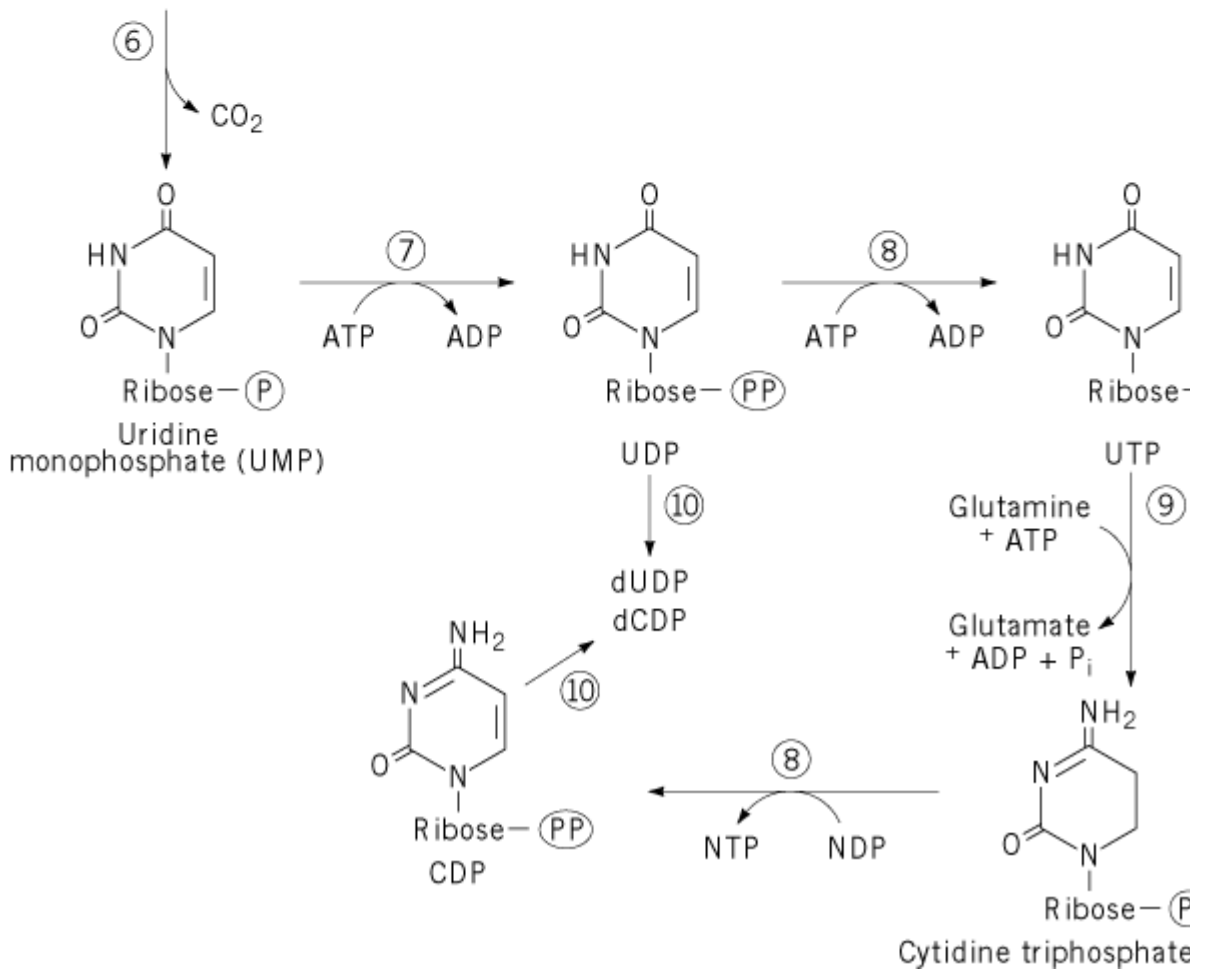
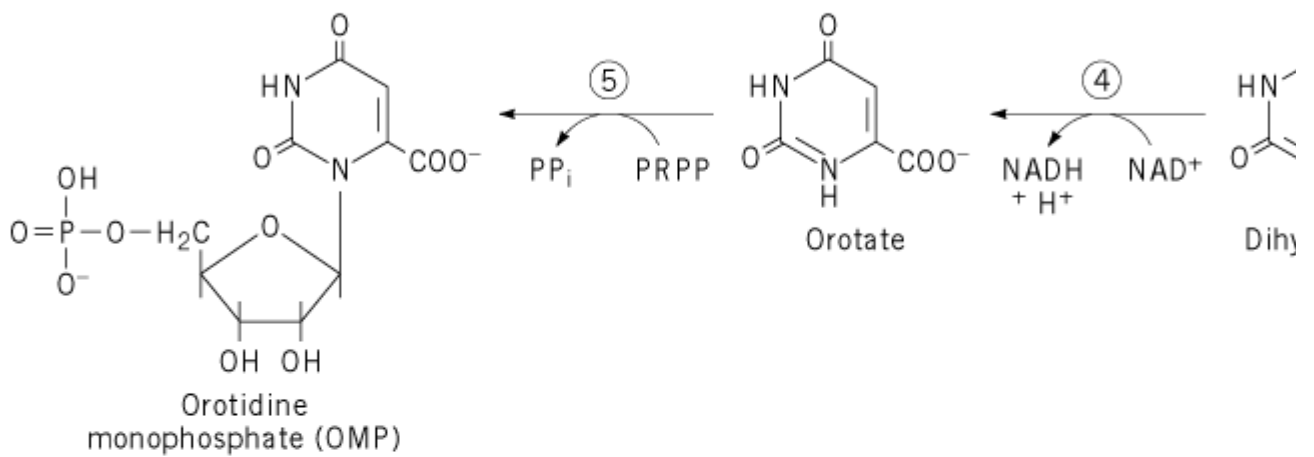
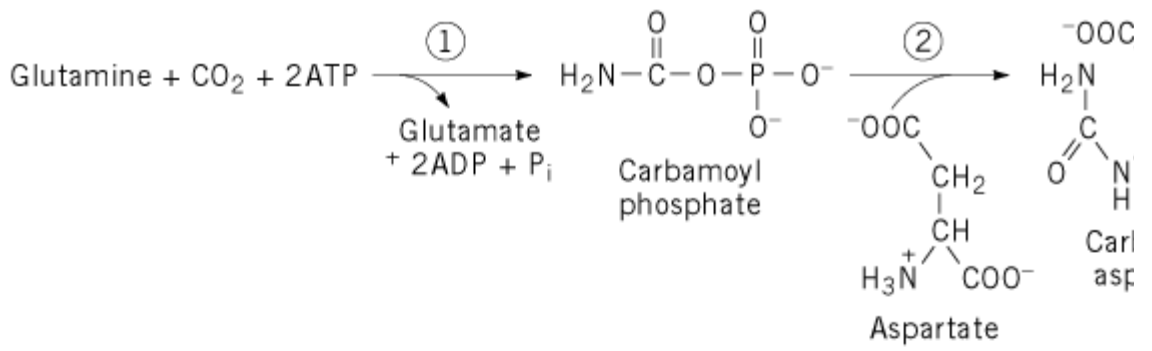
There are also important differences between purines and pyrimidines in the nature of the *de novo* synthetic pathway. The purine ring is assembled at the nucleotide level whereas, in pyrimidine synthesis, the sugar and phosphate are incorporated near the end of the pathway. Also, purine synthesis is branched, with guanine and adenine nucleotides formed by separate pathways from a common intermediate, inosinic acid. In contrast, pyrimidine ribonucleotide synthesis is unbranched, with uridine nucleotides serving as precursors to cytidine nucleotides.

In bacteria, the principal regulated step, and the committed reaction for the pathway,



is catalyzed by [aspartate transcarbamoylase](#) (ATCase), the [enzyme](#) that has perhaps told us more than any other enzyme about **allosteric** mechanisms in **enzyme regulation**. The entire pathway is shown in Figure [1](#). Control of ATCase involves allosteric inhibition by a pyrimidine end product, CTP, and “feedforward” activation by a purine nucleotide, ATP. The product, carbamoyl aspartate, contains all six of the carbon and nitrogen atoms that will eventually appear in the pyrimidine ring.

**Figure 1.** *De novo* biosynthetic pathway to pyrimidine nucleotides. Enzyme names are as follows: 1, carbamoyl phosphate synthetase; 2, aspartate transcarbamoylase; 3, dihydroorotase; 4, dihydroorotate dehydrogenase; 5, orotate phosphoribosyltransferase; 6, orotidylate decarboxylase; 7, UMP kinase; 8, nucleoside diphosphate kinase; 9, CTP synthetase; 10, ribonucleoside diphosphate reductase. PRPP is 5-phosphoribosyl-1-pyrophosphate. The circle drawn around the letter P denotes organic phosphate, or inorganic triphosphate, as indicated. For the conversion of CTP to CDP by nucleoside diphosphate kinase, the preferred phosphate donor, has not been identified.



Carbamoyl phosphate synthetase (CPS), which provides one of the substrates for aspartate transcarbamoylase, is often considered the initial reaction in pyrimidine synthesis. Carbamoyl phosphate is a precursor to both pyrimidines and arginine. Bacterial cells contain a single carbamoyl phosphate synthetase, so that the reaction is not committed to pyrimidine synthesis, and control is exercised at the next reaction, catalyzed by ATCase. **Eukaryotic** cells, however, contain two carbamoyl phosphate synthetases in different compartments—one in the cytosol, specialized for pyrimidine synthesis, and one in [mitochondria](#), used for the synthesis of arginine (and urea in those organisms possessing a urea cycle for amino acid degradation). Thus, one usually thinks of carbamoyl phosphate synthetase in eukaryotic cells as catalyzing the first step in pyrimidine synthesis (reaction 1 in Fig. 1), with ATCase as reaction 2.

In all cells, carbamoyl aspartate cyclizes (reaction 3) to give a ring compound, dihydroorotate, and this undergoes dehydrogenation (reaction 4) to give the first pyrimidine—orotate, which is 6-carboxyuracil. A phosphoribosyltransferase reaction with 5-phosphoribosyl-1-pyrophosphate (PRPP) gives a pyrimidine nucleotide, orotidine-5'-phosphate (reaction 5). Decarboxylation (reaction 6) yields uridine monophosphate, which undergoes two successive phosphorylations (reactions 7 and 8) to give uridine triphosphate, UTP. UTP is the substrate for an amidotransferase enzyme, CTP synthetase, which transfers the amide nitrogen of glutamine in the ATP-dependent conversion of UTP to CTP (reaction 9).

In most organisms, the synthesis of deoxyribonucleotides occurs at the nucleoside diphosphate level (reactions 10), with reduction of the ribose sugar to deoxyribose *in situ* (see [Ribonucleotide Reductases](#)). UDP is formed as an intermediate in UTP synthesis, whereas CDP is probably synthesized in a [nucleoside diphosphate kinase](#)—catalyzed reaction in which CTP serves as a phosphate donor. Because nucleoside diphosphate kinase has very low specificity for both phosphate donor and phosphate acceptor, this represents the most straightforward route for conversion of CTP to CDP. There is a report of a myokinase -type activity in mammalian cells that reversibly transfers phosphate from CTP to CMP, yielding two CDPs (1). If this enzyme plays a major role in CDP synthesis, much of the CMP probably comes from catabolism of cytidine-containing phospholipids, CDP-choline, or CDP-diacylglycerol. CMP is not an intermediate in pyrimidine synthesis *de novo*.

## 2. Regulation of Pyrimidine Synthesis

As noted earlier, bacterial ATCase presents a particularly well-understood example of allosteric regulation. Control of pyrimidine ring synthesis occurs also at the level of gene [transcription](#), with an interesting **attenuation** process (2). Classic attenuation control involves [operons](#) of amino acid biosynthesis, where the rate of [translation](#) of a [messenger RNA](#) determines the folding pattern of a nascent mRNA which, in turn, determines whether the conformation adopted by that mRNA will lead to premature transcriptional termination. In the case of pyrimidine synthesis, the intracellular concentration of UTP inversely controls the transcription of the genes *pyrB1* and *pyrE*, which encode subunits of ATCase and orotate phosphoribosyltransferase, respectively; at high concentrations of UTP, the transcription rates are low. Low UTP levels limit the rate of transcriptional chain elongation because the  $K_m$  of **RNA polymerase** for this nucleotide is high. Under these conditions, [ribosomes](#) translating the nascent mRNA catch up to RNA polymerase, and the interaction interferes with formation of an RNA hairpin loop that acts as a transcription terminator. By contrast, high UTP levels allow RNA polymerase to move fast enough to prevent interference with formation of the transcription terminator.

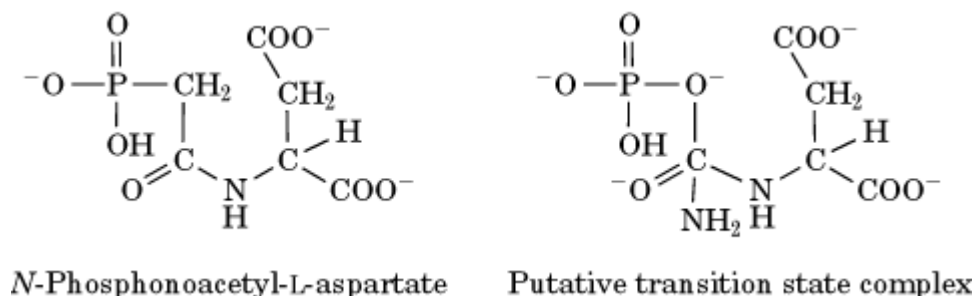
Another important control mechanism involves allosteric regulation of CTP synthetase. This enzyme is activated by GTP and inhibited by its product, CTP (3). The physiological importance of this regulation is seen in the dramatic effects in mammalian cells of mutations that abolish the

feedback inhibition by CTP. Cells bearing these mutations have grossly elevated pools of CTP and dCTP and abnormally small pools of uridine and thymidine nucleotides. The pool imbalance creates a **mutator** phenotype, and the dTTP deficiency creates a thymidine auxotrophy. Also, a recent report states that phosphorylation of CTP synthetase by protein [kinase C](#) decreases the enzyme's sensitivity to CTP inhibition (4).

### 3. Multifunctional Enzymes in Eukaryotic Pyrimidine Synthesis

In animal cells, the six enzymes of UMP biosynthesis are encoded by just three genes because of the existence of two multifunctional enzymes. This finding came to light with the synthesis of *N*-phosphonoacetyl-L-aspartate (PALA), which was produced as a [bisubstrate analogue](#) in the ATCase reaction and a potential inhibitor of pyrimidine synthesis (Fig. 2). PALA does inhibit pyrimidine synthesis. However, when cell lines resistant to PALA were developed, they were found to have greatly elevated levels not only of ATCase, the target enzyme, but also of carbamoyl phosphate synthetase and dihydroorotase. Analysis showed that the three activities are carried on one 220-kDa [polypeptide chain](#) and that three or six chains associate to form one enzyme molecule, called by the acronym the CAD protein (5). Through [sequence analysis](#) and limited **proteolysis**, it has been established that each chain of this multifunctional protein contains an N-terminal carbamoyl phosphate synthetase **domain**, a C-terminal ATCase domain, and an internal dihydroorotase domain. From the names of these three domains comes the term CAD. The ATCase domain displays significant [homology](#) to the catalytic subunit of bacterial ATCase. There is no regulatory subunit counterpart in the eukaryotic enzyme, and control of CAD activities and of the overall pathway is usually exerted at the level of carbamoyl phosphate synthetase.

**Figure 2.** Structures of *N*-phosphonoacetyl-L-aspartic acid and the presumed transition state in the aspartate transcarbamoylase reaction.



The fifth and sixth activities in the pathway—*orotate* phosphoribosyltransferase and *OMP* decarboxylase—also form a bifunctional enzyme, called UMP synthase. In mammals, these two activities are carried on a single polypeptide chain of about 52 kDa. The enzyme exists as either a monomer or a homodimer, with binding of **ligands**, particularly *orotidine* monophosphate, stimulating dimerization (6). Through sequence comparison with organisms containing separate forms of these enzymes, it has been established that the N-terminal part of the UMP synthase polypeptide is the *orotidine* phosphoribosyltransferase domain; the decarboxylase activity is carried in the C-terminal portion.

What is the biological significance of a protein design that carries either two or three activities in a single polypeptide chain? One likely answer is that it allows for coordinate control of enzymes involved in a single pathway. Another plausible answer is that it permits metabolic channeling—the facilitated transfer of a scarce or unstable metabolite from one active site to the next, so as to drive the overall pathway more efficiently. However, kinetic analysis has not established that either the CAD protein or UMP synthase displays significant channeling activity *in vitro*. Moreover, the

subcellular localization of pathway enzymes makes channeling *in vivo* seem improbable, at least so far as the overall pathway is concerned. The CAD protein, catalyzing reactions 1, 2, and 3, is located in the cytosol. The product of reaction 3, dihydroorotate, must then move to the outer surface of the inner mitochondrial membrane, where the fourth enzyme, dihydroorotate dehydrogenase, is located. The product, orotate, then moves back to the cytosol for the final two reactions catalyzed by UMP synthase. The necessary migration of an intermediate from one cell compartment to another clearly precludes channeling from steps 3 to 5.

#### 4. Pyrimidine Nucleotide Degradation

Pathways of pyrimidine catabolism are shown in Figure 2. Nucleic acid breakdown usually generates nucleoside 5'-monophosphates, which are cleaved to the respective nucleosides by pyrimidine 5'-nucleotidase. Cytidine is converted to uridine by deamination (the same enzyme also converts deoxycytidine to deoxyuridine). Uridine phosphorylase cleaves uridine, giving uracil plus ribose-1-phosphate. Uracil is reduced to dihydrouracil, which then undergoes a hydrolytic ring opening to give  $\beta$ -ureidopropionate. A further hydrolysis gives  $\beta$ -alanine (a precursor to coenzyme A) plus  $\text{CO}_2$  and  $\text{NH}_3$ . A parallel pathway, using mostly the same enzymes, acts on thymine nucleotides. The final product is  $\beta$ -aminoisobutyrate.

#### 5. Disorders of Pyrimidine Metabolism

In contrast to purine metabolic disorders, hereditary deficiencies of pyrimidine metabolic enzymes are extremely rare (6). The most common such deficiency, hereditary orotic aciduria, has been described in just 15 cases. This condition arises from a deficiency of UMP synthase, the bifunctional enzyme that catalyzes reactions 5 and 6 in *de novo* pyrimidine synthesis. The condition involves excessive excretion of orotate, which usually crystallizes in the urine. A megaloblastic anemia is seen, presumably resulting from insufficient pyrimidine nucleotides to support the nucleic acid synthesis needed for blood cell division and maturation. In addition, some form of mental retardation is seen in orotic aciduria. In all cases tested thus far, the patient responds well to treatment with uridine, which can replete the pyrimidine nucleotide pools through salvage pathways.

Deficiencies of pyrimidine 5'-nucleotidase and dihydropyrimidine dehydrogenase have been described (reactions 1 and 4, respectively, in Fig. 3). The former condition involves a hemolytic anemia; also, the erythrocytes contain abnormally high levels of UTP and CTP. In the latter condition, high levels of uracil and thymine are seen in both blood and urine. Finally, two cases have been described involving deficiency of dihydropyrimidase (reaction 5, Fig. 3), the enzyme that converts dihydrouracil to  $\beta$ -ureidopropionate. The condition involves excessive excretion of dihydrouracil. Because of the rarity of these conditions, few metabolic or clinical details are available.

#### Bibliography

1. P. Chiba and J. G. Cory (1987) *Cancer Biochem. Biophys.* **9**, 353–358.
2. K. F. Jensen, F. Bonekamp, and P. Poulsen (1986) *Trends Biochem. Sci.* **11**, 362–365.
3. B. Aronow and B. Ullman (1987) *J. Biol. Chem.* **262**, 5106–5112.
4. W-L. Yang, M. E. C. Bruno, and G. M. Carman (1996) *J. Biol. Chem.* **271**, 11113–11119.
5. P. F. Coleman, D. P. Suttle, and G. R. Stark (1977) *J. Biol. Chem.* **252**, 6379–6385.
6. D. R. Webster, D. M. O. Becroft, and D. P. Suttle (1995) In *The Metabolic and Molecular Basis of Inherited Disease*, 7th ed. (C. R. Scriver, A. L. Beaudet, W. S. Sly, and D. Valle, eds.) McGraw-Hill, New York, pp. 1799–1837.

#### Suggestions for Further Reading

7. Pyrimidine metabolism is described in all standard biochemistry textbooks, for example, C. K.



Mathews and K. E. van Holde (1996) *Biochemistry*, 2nd Ed., Benjamin/Cummings, Menlo Park, CA, pp. 796–798. In addition, reference 6 above includes a complete presentation of pyrimidine nucleotide synthesis and degradation, as well as a description of clinical disorders and an exhaustive (353-entry) reference list.

## Qb Replicase

The RNA-dependent **RNA polymerase** of **bacteriophage Qb** was first isolated in 1965 by Haruna and Spiegelman, who named it Qb replicase (1). Today, this enzyme still represents the prototype of RNA **virus**-replicating polymerases, and it has remained unique in that it allows the specific replication and amplification of an infectious viral RNA by a complex of soluble, stable, and defined protein components *in vitro*. Early studies had established that the viral RNA **genome** (the plus strand) was replicated by primer-independent end-to-end RNA synthesis in the 5' to 3' direction. A free, single-stranded complementary strand (minus strand) is synthesized as an intermediary product and serves as a **template** for the synthesis of single-stranded plus strands, whereas double-stranded RNA is devoid of any template activity (2). Investigations with Qb replicase have resulted over the years in a surprising number of novel concepts of more general relevance, such as: (1) specific template recognition by a polymerase (1); (2) recruitment of host proteins as subunits of a viral polymerase (3, 4); (3) role of a viral polymerase in the temporal control of **protein biosynthesis** (2); (4) an internal RNA binding site as an **enhancer** of synthesis initiation (5); (5) activation of an RNA template by a host factor acting as an RNA **molecular chaperone** (6); (6) secondary and tertiary structure as determinants of template recognition (7, 8); (7) role of secondary structure in elongation of RNA synthesis (9, 10); (8) sequence **evolution** *in vitro* (11) and the quasispecies concept (12); (9) *de novo* synthesis of replicatable RNA (13); (10) **site-directed mutagenesis** (14); (11) use of small RNA templates as **vectors** for foreign sequences (15); and (12) development of such systems for diagnostic purposes (16). Most of these aspects will be discussed here.

### 1. Protein components of the *in vitro* Qb RNA replication system

Qb replicase is isolated from Qb-infected *Escherichia coli* cells as an **enzyme** consisting of four different subunits, only one of which (subunit b, 65.5 kDa) is phage-specific. The other subunits are recruited from the host cell and were identified as the **ribosomal** protein S1 (subunit a, 61.2 kDa) and the protein synthesis **elongation factors** EF-Tu and EF-Ts (subunits g, 43.2 kDa, and d, 30.3 kDa). In addition to this holoenzyme form of replicase, a smaller core enzyme species, consisting of only subunits b, g, and d, was found to be equally active for all templates except the Qb plus strand. Of the three core subunits, b (phage-specific) represents the catalytic subunit, and EF-Tu/Ts (subunits g and d) appear to be important for the assembly and stabilization of the active enzyme complex. The role of protein S1 (subunit a) as an **RNA-binding protein** consists of mediating the recognition of the Qb plus-strand RNA as a template (see below). The template activity of the Qb plus strand, in contrast to the other known template RNAs, depends further on the presence of an additional host-derived RNA binding protein called the “host factor” (17), or Hfq protein (6×11kDa). Hfq protein was recently found to function in the cell as an “RNA chaperone” required to activate translation for the induction of the stationary phase and other functions (18).

Replicases of other RNA phages make use of the same host-derived subunits (19, 20), but their host factor requirements may be different, as shown for phage f2 (21) and its close relative MS2 (22).

### 2. Template Specificity

For an RNA to be replicated and amplified by replicase, both the RNA itself and its complementary strand must be efficient templates. This is the case for Qb RNA and the Qb minus strand in the presence of the holoenzyme and host factor. The core enzyme without the host factor is sufficient to copy the minus-strand RNA into a plus strand, but cannot by itself use the latter as an efficient template. However, the core enzyme is efficient for replicating and amplifying a family of small RNAs, collectively called 6 S RNA, that arise *in vitro* (and possibly also *in vivo*) in the absence of added template RNA, either by uninstructed *de novo* synthesis or through elongation and recombination of small RNA fragment contaminants (see below). RNAs of closely related phages (eg, SP) have partial activity (23), but more distant phage RNAs (eg, MS2), like other viral and cellular RNAs, are inactive. Synthetic templates like poly(C) and mixed polyribonucleotides rich in C function as templates for the synthesis of a complementary strand, which however remains paired with the template and is itself not active as a template.

### 3. Recognition of the Qb plus-strand template

Replicase holoenzyme forms tight complexes with the Qb plus strand by binding at two internal RNA regions called the S- and M-sites (24), mapping at positions 1247 to 1346 and 2545 to 2867 of the Qb genome (4217 nucleotides), respectively. **Electron microscopic** examination of such complexes revealed the formation of RNA loops between these sites, showing that the two interactions occurred simultaneously on one and the same strand (25). Similar looped complexes were observed by binding of S1 protein alone, which suggests that these internal interactions are mediated by the  $\alpha$  subunit (26). Binding of replicase to the S-site, located at the start of the coat protein cistron, is not necessary for template recognition but is thought to represent a **translational repression** mechanism by which replicase prevents the attachment of ribosomes during the time it uses the plus strand as a template (2). This stands in contrast to the M-site, in which a branched secondary structure element (nucleotides 2696 to 2754) has an enhancer-like function essential for high template activity (5). Binding of the 3'-end, ie, the site of initiation of synthesis, is mediated by the host factor, which in addition binds to two internal sites adjacent to the M- and S-sites (26, 27). Qb phages adapted to growth in host factor-less *E. coli* strains contained mutations altering the 3'-terminal folding of the RNA and releasing the sequence at the immediate 3'-end from base-pairing (6). This is in agreement with the concept of the host factor as an RNA chaperone that modifies the secondary and tertiary structures of the 3'-end so as to make it accessible to replicase.

In summary, the recognition of the Qb plus strand by replicase occurs through contacts at three (or possibly even more) different RNA sites that map far apart on the genome but must be held in close vicinity by the tertiary structure of the RNA (8). Interestingly, it is the host proteins S1 and Hfq that are the mediators of these interactions, whereas the catalytic  $\beta$  subunit is very probably responsible for the direct contacts at the 3'-end of the template that allow synthesis to start. Initiation occurs with GTP, which is complementary to the template's essential 3'-terminal  $\frac{1}{4}$  CCC motif; the additional A residue found at the 3'-end of most RNA chains is not essential.

### 4. The Qb minus strand and other RNAs as templates

The interactions of replicase with the Qb minus strand and the replicating 6 S RNAs are expected to be different from those with the plus strand, because the former RNAs are efficient templates for core replicase in the absence of proteins S1 and Hfq. Deletion analysis of the Qb minus strand identified two structural features required for template activity (7). One consists of a sequence that folds into an imperfectly base-paired stem-loop structure located near the 5'-terminus of the minus strand. A highly **homologous** sequence is found in MDV-1 RNA (a 6 S RNA) and was characterized as an element essential for recognition of this template (28). The other essential element is formed by two short complementary sequences forming a helical stem by long-range base-pairing near the 3'-end of the minus strand. On the basis of stability calculations, the 3'-terminal  $\frac{1}{4}$  CCC(A) sequence, which represents another essential structure, appears in the minus strand to exist in an unpaired

conformation, in agreement with the fact that this template does not require a host factor.

Of the many short-chain RNAs characterized as templates, their only common features (29) consisted of a single stranded  $\frac{1}{4}$  CCC(A) 3'-end, as well as sequences rich in secondary structure, as discussed below.

Recently, the SELEX technique was used to generate two families of short RNA ligands that bound strongly to replicase holoenzyme and were active templates (30). One contained a pseudoknot structure rich in A and C residues in the unpaired loops and was found to crosslink to S1 protein. The other contained a polypyrimidine tract and crosslinked to EF-Tu (the g subunit). At present, it appears difficult to integrate these findings into a common picture with those from the Qb plus and minus strand.

#### 5. Secondary structure required for the release of single-stranded product RNA

Short-range intrastrand base-pairing in product and template RNAs appears to be necessary for the release of the product RNA in a single-stranded form, presumably by competitively inhibiting the formation of long product-template duplexes. In fact, longer tracts of template sequence unable to form local stem-loop structures resulted in the formation of duplexes between product and template strands that could not be replicated further (9). Moreover, RNA phage mutants carrying insertions of such tracts *in vivo* quickly evolved to revertants that had eliminated the inserts (10).

#### 6. Evolution *in vitro* and *in vivo*, error rates, recombination, and the question of *de novo* synthesis

Spiegelman and coworkers recognized very early that their RNA amplification system provided a unique opportunity to study evolution *in vitro* by serial transfer experiments (11). Studies focusing on this aspect were instrumental in the development of the “quasi-species” concept by M. Eigen (12). Estimates of the phage Qb mutation rate (about  $10^{-4}$ ; (31)) and of the sequence heterogeneity of the Qb genome in phage populations determined by Weissmann and coworkers (32) were in full agreement with this notion.

The occurrence of RNA [recombination](#) in RNA phage replication was first experimentally demonstrated *in vivo* (33), but had been suspected before, because sequences homologous to cellular RNAs were found in several 6 S RNA species (34). Later *in vitro* studies confirmed the process, which was thought to take place via copy-choice (35) or by a novel type of mechanism (36).

The claim that rigorously purified replicase under specific conditions can be observed to synthesize small replicating RNAs *de novo* without template instruction (13) encountered opposition because of the difficulty of excluding contamination by traces of 6 S RNA in any laboratory engaged in *in vitro* work with replicase (37). The doubts have stimulated, however, a large amount of careful work, making today a convincing case for *de novo* synthesis, such as the demonstration that different incubation mixtures of identical ingredients result in entirely unrelated RNA sequences arising after widely fluctuating latency periods (29). The stochastically arising RNA species and their propagation and evolution to high-efficiency replication could be followed quantitatively by fluorescence techniques in a sophisticated capillary incubation apparatus (13).

#### 7. Practical applications of replicase

Not surprisingly, the unique properties of replicase as a nucleic acid amplification system have raised interest in developing practical uses. Most efforts were directed toward using replicating RNAs as vectors for the amplification of foreign sequences (15). Success was achieved mostly with relatively short inserts, because of the high tendency of the system toward mutation and recombination and the requirement of a strong secondary structure, as outlined above. Nevertheless, a sensitive diagnostic assay using short inserts of ribosomal RNA sequences from several pathogenic

agents was developed and marketed as an alternative to PCR (16).

## 8. Outlook

Despite the described attempts to get “something useful” out of Qb replicase, it appears likely that the main fascination with this enzyme will remain within the realm of basic science. The most intriguing questions may concern the precise mechanics of the interplay of the macromolecular components that make the system behave in the way it does. No doubt, knowledge of the three-dimensional structures of these components would be a great step forward toward this understanding. Efforts toward determining a structure of replicase by **X-ray crystallography** have not been successful thus far, but should be continued and extended to substructures like the core enzyme, S1 protein, and host factor. As the crystallography of RNA is still in its infancy, large RNA structures like the Qb plus and minus strand may not be accessible very soon, but interesting information could come from specific and well-studied small RNAs like MDV-1, especially if co-crystals with replicase and its subassemblies could be analyzed. In the meantime, the recent rapid advances in computer modeling of RNA structures (38) will continue and, together with **phylogenetic** comparisons and chemical or biochemical probing techniques, should give us improved models of template RNAs and their complexes with replicase. Similarly, we can expect that functional studies with site-directed RNA mutants, and especially the elegant evolutionary approach in which deficient mutants are allowed to evolve back to high viability (8), will continue to advance our understanding of this exemplary RNA-protein recognition process.

## Bibliography

1. I. Haruna and S. Spiegelman (1965) *Proc. Natl. Acad. Sci. USA* **54**, 579–587.
2. C. Weissmann (1974) *FEBS Lett.* **40**, S10–S18.
3. T. Blumenthal, T. A. Landers, and K. Weber (1972) *Proc. Natl. Acad. Sci. USA* **69**, 1313–1317.
4. A. J. Wahba et al. (1974) *J. Biol. Chem.* **249**, 3314–3316.
5. D. Schuppli, G. Miranda, S. Qiu, and H. Weber (1998) *J. Mol. Biol.* **283**, 585–593.
6. D. Schuppli et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 10239–10242.
7. D. Schuppli, I. Barrera, and H. Weber (1994) *J. Mol. Biol.* **243**, 811–815.
8. J. Klovins, V. Berzins, and J. van Duin (1998) *RNA* **4**, 948–957.
9. V. D. Axelrod, E. Brown, C. Priano, and D. R. Mills (1991) *Virology* **184**, 595–608.
10. R. C. L. Olsthoorn and J. van Duin (1996) *J. Virol.* **70**, 729–736.
11. S. Spiegelman (1971) *Quart. Rev. Biophys.* **4**, 213–253.
12. M. Eigen (1993) *Scienti. American* **269** (July), 32–39.
13. C. K. Biebricher, M. Eigen, and J. S. McCaskill (1993) *J. Mol. Biol.* **231**, 175–179.
14. R. A. Flavell, D. L. O. Sabo, E. F. Bandle, and C. Weissmann (1974) *J. Mol. Biol.* **89**, 255–272.
15. Y. Wu, D. Y. Zhang, and F. R. Kramer (1992) *Proc. Natl. Acad. Sci. USA* **89**, 11769–11773.
16. B. B. Stone et al. (1996) *Mol. Cell. Probes* **10**, 359–370.
17. M. T. Franze de Fernandez, L. Eoyang, and J. T. August (1968) *Nature* **219**, 588–590.
18. A. Muffler, D. Fischer, and R. Hengge-Aronis (1996) *Genes Dev.* **10**, 1143–1151.
19. N. V. Fedoroff and N. D. Zinder (1971) *Proc. Natl. Acad. Sci. USA* **68**, 1838–1843.
20. H. Mori, Y. Fukami, and I. Haruna (1978) *J. Biochem.* **84**, 681–686.
21. N. V. Fedoroff and N. D. Zinder (1973) *Nature New Biol.* **241**, 105–108.
22. S. Qiu et al. (1997) *Virology* **227**, 211–214.
23. T. Yonesaki, K. Furuse, I. Haruna, and I. Watanabe (1982) *Virology* **116**, 379–381.
24. F. Meyer, H. Weber, and C. Weissmann (1981) *J. Mol. Biol.* **153**, 631–660.
25. H. J. Vollenweider, Th. Koller, H. Weber, and C. Weissmann (1976) *J. Mol. Biol.* **101**, 367–377.

26. G. Miranda et al. (1997) *J. Mol. Biol.* **267**, 1089–1103.
27. I. Barrera, D. Schuppli, J. M. Sogo, and H. Weber (1993) *J. Mol. Biol.* **232**, 512–521.
28. T. Nishihara, D. R. Mills, and F. R. Kramer (1983) *J. Biochem.* **93**, 669–674.
29. C. K. Biebricher and R. Luce (1993) *Biochemistry* **32**, 4848–4854.
30. D. Brown and L. Gold (1996) *Proc. Natl. Acad. Sci. USA* **93**, 11558–11562.
31. E. Batschelet, E. Domingo, and C. Weissmann (1976) *Gene* **1**, 27–32.
32. E. Domingo, D. Sabo, T. Taniguchi, and C. Weissmann (1978) *Cell* **13**, 735–744.
33. K. Palasingham and P. N. Shaklee (1992) *J. Virol.* **66**, 2435–2442.
34. A. V. Munishkin et al. (1991) *J. Mol. Biol.* **221**, 463–472.
35. C. K. Biebricher and R. Luce (1992) *EMBO J.* **11**, 5129–5135.
36. A. B. Chetverin, H. V. Chetverina, A. A. Demidenko, and V. I. Ugarov (1997) *Cell* **88**, 503–513.
37. A. B. Chetverin, H. V. Chetverina, and A. V. Munishkin (1991) *J. Mol. Biol.* **222**, 3–9.
38. M. Zuker and A. B. Jacobson (1998) *RNA* **4**, 669–679.

### Suggestions for Further Reading

39. T. Blumenthal (1982). Q replicase. *The Enzymes* **15**, 267–279.
40. T. Blumenthal and G. G. Carmichael (1979). RNA replication: Function and structure of Q replicase. *Ann. Rev. Biochem.* **48**, 525–548.
41. R. I. Kamen (1975) "Structure and function of the Q replicase". In *RNA Phages* (N. D. Zinder, ed.), Cold Spring Harbor Lab, Cold Spring Harbor, NY, pp. 203–234.
42. J. van Duin (1988) "Single-stranded RNA bacteriophages". In *The Bacteriophages* (R. Calendar, ed.), Plenum Press, New York, pp. 117–167.

## Quaternary Structure

The *quaternary structure* describes the way two or more **polypeptide chains** are arranged within the same **protein** molecule. The term, suggested by J. D. Bernal, was used first by J. Kendrew in 1959 [(1)] to further develop the already classical hierarchy of K. Linderström-Lang. If (1) **primary structure** refers to the amino acid sequence, (2) **secondary structure** is **alpha-helices** or **beta-sheets**, and (3) **tertiary structure** is the fold of the polypeptide chain, quaternary structure refers to the assembly of several polypeptide chains, as in four-chain **hemoglobin**, or of many chains in **virus** capsids or muscle fibers. In 1959, Kendrew had just solved the crystal structure of the single-chain **myoglobin**, which illustrated secondary structure as the  $\alpha$ -helix and produced the first image of a tertiary structure. Soon afterward, M. Perutz obtained the crystal structure of hemoglobin, making quaternary structure directly visible in atomic detail. We now have many more examples in proteins ranging from **enzymes** to **transcription factors** and **membrane** components.

Proteins with more than one polypeptide chain are *oligomeric* or *multimeric*, and their chains are also called *subunits*. Oligomeric proteins may be built of identical subunits, when they are called *homo-oligomers* and their quaternary structure can be coded in the formula  $a_n$  to indicate that they are made of  $n$  copies of an  $a$  subunit. Other oligomeric proteins are referred to as *hetero-oligomers*. The quaternary structure formula  $a_n b_m$  describes a hetero-oligomer that contains two types of subunits  $a$  and  $b$  in  $m$  and  $n$  copies, respectively. **Gtpases** involved in cellular **signal transduction** are abg hetero-

trimers. Mammalian hemoglobins may be written as either  $a_2b_2$  or  $(ab)_2$ . **Immunoglobulins G** also have this formula, but unlike the a and b chains of hemoglobin, the **heavy** and **light chains** of an immunoglobulin are covalently linked by [disulfide bonds](#).  $a_2$  dimers and  $a_4$  or  $a_2b_2$  tetramers are the most common formulae in globular oligomeric proteins. More complex quaternary structures do exist, however, especially in membrane proteins. Thus, the formula  $a_2bgd$  describes the quaternary structure of the nicotinic [acetylcholine receptor](#) of the neuromuscular junction, which comprises four different chains with unequal stoichiometries.

In order to write a formula such as  $a_2b_2$ , we need to know the number and type of the constituent polypeptide chains. This cannot be deduced from just the knowledge of the **gene** sequences, although genetic evidence may point to the presence of more than one type of chain. Determination of the formula usually requires expressing and purifying the protein. After purification, [SDS-PAGE](#) or [mass spectrometry](#) are commonly used to separate the individual polypeptide chains and estimate their apparent molecular weights. If a homogeneous protein yields more multiple species with different molecular weights, it is probably a hetero-oligomer, but quantitative peptide **sequencing** will be needed to determine the relative amounts of each chain type and to eliminate the possibility that they derive from **proteolysis** or [post-translational modification](#) of a single gene product.

Whether or not the constituent polypeptide chains are of one or several types, their actual number in the protein is ascertained by measuring the molecular weight of the native protein complex. The most commonly used methods are gel filtration and **sedimentation equilibrium** ultracentrifugation. However, neither is as easy and sensitive as SDS-PAGE, and both are much less accurate than mass spectrometry, two techniques that yield the molecular weights of the components, but not of the assembly. Moreover, the results from gel filtration and ultracentrifugation may be misinterpreted if the assembly is not fully stable, so that partial dissociation or aggregation takes place under conditions where the measurement is made. Chemical [crosslinking](#) followed by SDS-PAGE is an alternative to molecular weight determination. Covalent crosslinking an  $a_n$  protein can, in principle, yield  $n$  bands on the gel having apparent molecular weights corresponding to species  $a$ ,  $a_2$ ,  $1/4$ ,  $a_n$ . However, not all the constituent chains may form crosslinks readily, and the conditions where a complete ladder of  $n$  bands can actually be observed may be difficult to obtain. When the quaternary structure does not dissociate spontaneously during a separation, the number of subunits can be determined by hybridizing two distinguishable forms of the subunit, say,  $A$  and  $B$ , and then separating the various oligomers on the basis of this difference. For a tetramer, this should yield five species:  $A_4$ ,  $A_3B_1$ ,  $A_2B_2$ ,  $A_1B_3$ , and  $B_4$ . In practice, molecular weight determination is more difficult and less reliable for a native protein than for a **denatured** protein. Indeed, we lack that basic information for the vast majority of the proteins known only as putative gene products or bands on a gel. In addition, quaternary structures are sometimes revised, and tabulated values should be verified.

Beyond counting polypeptide chains, the quaternary structure is part of the three-dimensional structure as a whole. [NMR](#) and [X-ray crystallography](#) studies can determine the symmetry relationships between subunits, the location of their interfaces, and the nature of their atomic contacts, together with the fold of the polypeptide chains. In crystals, contacts between subunits coexist with intermolecular, crystal packing contacts. The subunit contacts can usually be recognized just on the basis of their greater extent, but ambiguities sometimes arise. They must be resolved by other methods, such as the [site-directed mutagenesis](#) of the proposed subunit contacts. To characterize the contacts by NMR, special isotope labeling techniques have been developed. In a dimer, one subunit is labeled with  $^{13}\text{C}$ , the other with  $^{15}\text{N}$ , and isotope-edited **NOE** signals are specifically detected between CH hydrogen atoms of the first subunit and NH hydrogen atoms of the other. However, the large sizes of oligomeric proteins generally make the NMR resonance lines broad, so the spectra are poorly resolved and difficult to assign. Indeed, neither NMR nor crystallography can dispense with the determination of the molecular weight of native protein.

The subunit contacts that make up the quaternary structure are generally noncovalent and form interesting examples of specific [protein–protein interactions](#). They almost always make a major contribution to the stability of the native protein and often its function also. There are many oligomeric proteins in which a **ligand-binding** site exists at the interface between subunits. This is the case of the **antigen-binding** sites of immunoglobulins and of the DNA-binding sites in the great majority of [DNA-binding proteins](#). Enzyme [active sites](#) at subunit interfaces are not uncommon, and **allosteric** mechanisms of regulation largely rely on properties of the quaternary structure. These aspects are discussed under [Oligomeric proteins](#), together with the evolutionary aspects of quaternary structure.

#### Bibliography

1. J. Kendrew (1959) Structure and function in myoglobin and other proteins. *Fed. Proc.* **18**, 740–751.

#### Suggestion for Further Reading

2. C. Bränden and J. Tooze (1991) *Introduction to Protein Structure*, Garland Publishing, New York.

### R-Factor (Crystallographic)

The final result in [X-ray crystallography](#) is presenting a molecular model of the structure being determined. If this were an ideal model and if the packing of the models in the crystal structure were accurately known, the calculated the amplitude of the [structure factor](#)  $F_{\text{calc}}$  would be equal to the observed amplitude  $F_{\text{obs}}$  for each reflection. In practice this ideal situation does not exist, and the difference between the two values indicates model's inadequacies. This difference, averaged over all of the measured reflections, is expressed as the crystallographic R-factor or reliability or residual index  $R$ :

$$R = \frac{\sum_{hkl} |F_{\text{obs}}| - k|F_{\text{calc}}|}{\sum_{hkl} |F_{\text{obs}}|}$$

or alternatively

$$R = \frac{\sum_{hkl} |F_{\text{obs}}| - k|F_{\text{calc}}|}{\sum_{hkl} |F_{\text{obs}}|} \times 100\%$$

where  $k$  is a factor scaling the  $F_{\text{calc}}$  to the  $F_{\text{obs}}$  's.

If the atoms were to be placed entirely at random in a non-centrosymmetrical structure like a protein,  $R$  would have the value 0.59. For correct structures refined to, for example, 2.0 Å resolution, the  $R$ -factor is of the order of 0.20. If it is much higher, the structure is not likely to be correct.

A better estimate for the reliability of a protein structure is the free  $R$ -factor ([1](#)). It is calculated for a random selection of 5 to 10% of the observed reflections. These reflections are not used in the refinement, so refining and testing are completely independent.

## Bibliography

1. A. T. Brünger (1993) *Acta Crystallogr.* **D49**, 24–36.

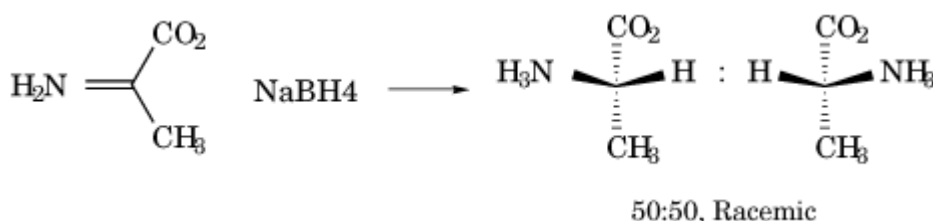
## Suggestion for Further Reading

2. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Racemic And Racemization

A mixture containing equal amounts of both **enantiomers** of a compound is a *racemic mixture*, which may also be called a *racemate*. The original compound identified as containing equal parts of two enantiomers was racemic acid as isolated from grapes (1). This compound is now known as ( $\pm$ ) tartaric acid. When a **chiral** compound is formed by the reaction of two achiral reactants, the product must be a racemic mixture (2). An example of such a reaction is the reduction of the imine of pyruvate by sodium borohydride as shown in Figure 1. The product alanine is chiral, but both enantiomers must be generated in identical amounts.

**Figure 1.** Generating the chiral products L- and D-alanine from achiral reactants. The two products are produced in equal amounts, to give a racemic mixture of DL-alanine.



The process of separating a racemic mixture into its two enantiomeric components is called *racemic resolution*. Any process that catalyzes the interconversion of enantiomers, ie, a *racemization*, will necessarily result in a racemic mixture being formed. This must be so because the free energy of formation of the two enantiomers, in the absence of any other chiral compound, must be identical. During the chemical process of [peptide synthesis](#), racemization of the activated amino acid derivatives has been a major problem. [Enzymes](#), such as proline racemase, will racemize a solution of either the L- or D-enantiomer by catalyzing the interconversion of the two enantiomers, generating a racemic mixture of the substrate.

## Bibliography

1. L. Pasteur (1853) *Pharmacol. J.* **XIII**, 111.
2. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York, Chap. "4".

## Suggestion for Further Reading

3. March, J. (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, pp. 102–107.



## Radiation Hybrid

The fusion of **somatic cells** from different species to create somatic cell hybrids has provided an exceptionally useful tool for mapping of genes to particular [chromosomes](#) because the chromosomes of one species are always progressively lost in such interspecific somatic cell hybrids (1). Using hybrid cell lines, it is also possible to use recognizable chromosome breakpoints to localize genetic loci more accurately. This process is greatly facilitated by inducing chromosome breaks with radiation (2). This technique uses a somatic cell hybrid between a rodent and a human cell that contains a single human chromosome, all of the rodent chromosomes, and a selectable marker that confers resistance to the antibiotic [neomycin](#). The chromosomes in this hybrid cell line are fragmented by irradiation such that cell death normally occurs. However, the irradiated hybrid cell line is rescued by fusion with another rodent cell line that is neomycin-sensitive. Growth in neomycin allows selection to operate in favor of the resulting radiation hybrid cell lines, rather than the neomycin-sensitive unfused cell line. Fragments of the human chromosomes are retained and incorporated into the rodent cell chromosomes. The different radiation hybrid cell lines generated in this way can be examined for the presence or absence of molecular markers from the human chromosome.

The production of radiation hybrids has provided a useful methodology for mapping how close genes are in the human chromosome. Two genes that are close together remain together following radiation, whereas those that are far apart are separated. By estimating the frequency of separation, it is possible to build up a physical map of the chromosome providing gene order and distance (3).

### Bibliography

1. R. L. Stallings and M. J. Siciliano (1981) *Somat. Cell Genet.* **7**, 683–697.
2. S. J. Goss and H. Harris (1977) *J. Cell Sci.* **25**, 17–28.
3. M. A. Walter and P. N. Goodfellow (1993) *Trends Genet.* **9**, 352–356.

### Suggestion for Further Reading

4. M. S. Clark and W. J. Wall (1996) *Chromosomes. The Complex Code*, Chapman and Hall, London.

## Radioactivity

Radioactivity is the property of certain elements to undergo spontaneous transformation of their atomic nuclei, the release of energy, and the formation of new elements (decay products). Radioactive nuclei are unstable and seek a more stable configuration by releasing energetic particles or photons of energy. These particles may include alpha particles, beta particles, and positrons. Photons include gamma rays, X rays, and neutrinos, which are discrete quanta of energy without mass or charge. These emissions impart energy to matter by creating tracks of ionized molecules. The emission of charged particles or gamma rays may be measured with radiation detectors. This property of radioactive materials makes them useful for a great number of practical applications in

the physical and biomedical sciences.

Every chemical element has one or more radioactive isotopes, and the total number of known radioactive and stable isotopes is more than 1500. Radioactive isotopes (radioisotopes or radionuclides) of a given element differ in the number of neutrons in the nucleus and, hence, in total atomic mass. A radioactive label (radiolabel) is a radionuclide or radioisotope in a chemical compound that replaces a stable isotope of the same element. It is used to mark the compound for detection by instruments that measure radioactivity. Radioactive labeling is useful for tracking the uptake, retention, metabolism, or clearance of chemical compounds, or for investigating metabolic pathways, [enzyme](#) kinetics, or chemical reactions. Radioactive labels are sometimes called radioactive tracers.

The radioactivity of a particular nuclide is determined by the configuration of its atomic nucleus and is independent of the chemical and physical state of the radioisotope and its environment (temperature and pressure). Radioactivity takes many different forms. The process of radioactive disintegration results in changes to the nucleus in atomic number and in its number of nucleons (protons plus neutrons), as indicated in Table 1.

**Table 1. Changes in Atomic Number and the Number of Nucleons by Radioactive Transformation<sup>a</sup>**

| Radiation Emitted                  | Change in     |                    |
|------------------------------------|---------------|--------------------|
|                                    | Atomic Number | Number of Nucleons |
| Alpha particle                     | -2            | -4                 |
| Beta (minus) particle              | +1            | 0                  |
| Beta (plus) particle (or positron) | -1            | 0                  |
| Gamma or X ray                     | 0             | 0                  |

<sup>a</sup> From F. H. Attix (1986) *Introduction to Radiological Physics and Radiation Dosimetry*, John Wiley & Sons, New York.

An example of radioactive transformation is the decay of phosphorus-32 to sulfur-32 by emission of a beta (minus) particle (or electron),  ${}_{-1}^0e$ :



Beta decay and positron emission also are accompanied by emission of energy in the form of a nonionizing neutrino. Decay schemes for radionuclides may be complex.

The total energy of the gamma rays and charged particles emitted during radioactive decay is equivalent to the net decrease in the rest mass of the disintegrating atom as it changes from parent to decay product. Its energy, momentum, and electronic charge are conserved. The emitted energy is either kinetic energy of particles in motion or quantum energy of photons, each of which degrades into heat. The ionization of matter through which radiation passes may result in direct or indirect chemical changes and radiation damage.

## 1. Half-Life

The transformation kinetics of radioactive decay are unique to each radioisotope, and each has its own characteristic, constant decay rate. The time required for any given radioisotope to decay in amount to one-half of its original amount is a measure of the rate at which radioactive transformation takes place. The physical half-life ( $t_{1/2}$ ) of a radioactive atom may range from fractions of a second to billions of years and is unique to each radionuclide. Naturally existing radionuclides have long physical half-lives or are created by the decay of radioisotopes with long half-lives.

If  $N$  is the number of identical radioactive atoms and  $\lambda$  is the radioactive decay constant ( $s^{-1}$ ; reciprocal seconds), then  $1/\lambda$  is the activity, and the rate of change in  $N$  at any time  $t$  is equal to the activity

$$\frac{-dN}{dt} = \lambda N \quad (2)$$

Integrating from  $t = 0$  (when  $N = N_0$ ), we obtain

$$\int_{N_0}^N \frac{dN}{N} = - \int_0^t \lambda dt \quad (3)$$

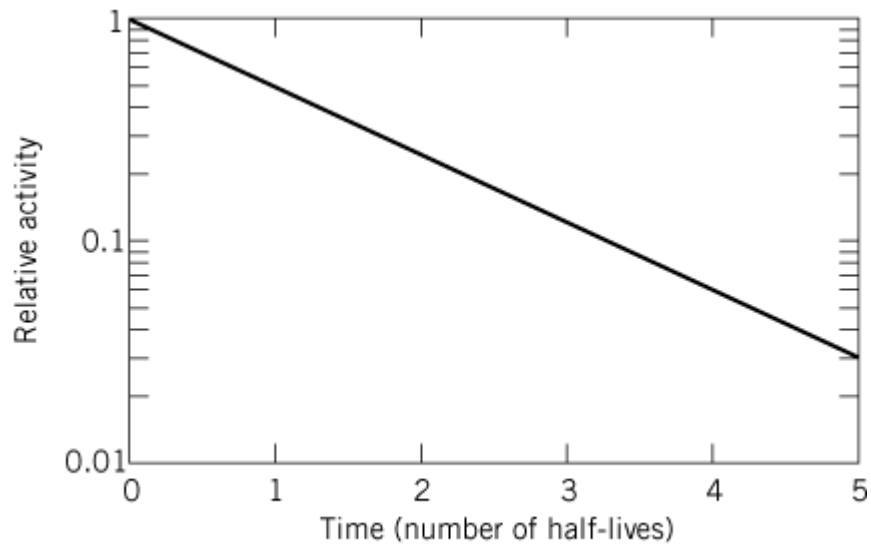
$$\therefore \ln N - \ln N_0 = -(\lambda t - 0) \quad (4)$$

$$\ln \frac{N}{N_0} = -\lambda t \quad (5)$$

$$\frac{N}{N_0} = e^{-\lambda t} \text{ or } N = N_0 e^{-\lambda t} \quad (6)$$

Equation 6 describes the observed law of radioactive decay and is a formula that is useful for determining the amount of a radioactive material remaining after a period of time  $t$  (Fig. 1).

**Figure 1.** Semilog plot of the exponential decay of a radioactive material in terms of the number of half-lives that have transpired.



The unit of activity was originally the curie (Ci), the number of disintegrations per second taking place in 1 g of  $^{226}\text{Ra}$ , where

$$1 \text{ Ci} = 3.7 \times 10^{10} \text{ disintegrations per second} \quad (7)$$

In recent years, the curie has been replaced with the special S.I. unit becquerel (Bq), which is one disintegration per second. However, the units of curie, millicurie, and microcurie are still commonly used.

$$1 \text{ Ci} = 3.7 \times 10^{10} \text{ Bq} \quad (8)$$

$$1 \text{ mCi} = 3.7 \times 10^7 \text{ Bq} \quad (9)$$

$$1 \mu\text{Ci} = 3.7 \times 10^4 \text{ Bq} \quad (10)$$

## 2. Mean Life

The mean life of a radioactive material is the time required for an original amount  $N_0$  to decay to  $1/e$  of the original amount. Thus

$$\frac{N}{N_0} = \frac{1}{e} = 0.3679 = e^{-\lambda\tau} \quad (11)$$

where  $\tau$  is the mean life (s) of the material:

$$\ln e^{-1} = -1 = -\lambda\tau \quad (12)$$

$$\tau(\text{mean life}) = \frac{1}{\lambda} = 1.443T_{1/2} \quad (13)$$

### 3. Specific Activity

The specific activity (Bq/g) is the activity (Bq) of a radioactive material per unit mass (g) or volume. The mass or volume may refer to the element itself or to the medium in which the radioactive material is contained. For example, the specific activity of a carbon-14-labeled compound is the radioactivity (Bq) of the carbon-14 divided by the total mass of all the compound molecules in a given volume.

### 4. Serial Transformation

If the decay products of a radioactive material are themselves radioactive, a decay chain is said to exist. The ingrowth of the first decay product is dependent on the rate of decay of the parent, and so forth through each daughter-product decay, until a stable isotope finally ends the chain. Three natural radioactive series exhibit long decay chains of successive members: the thorium-232 chain, with 12 members, concludes with stable lead-208; the uranium-238 series, with 19 members, concludes with stable lead-206, and the uranium-235 series, with 14 members, concludes with stable lead-207. A fourth series starting with plutonium-241, with 15 members and concluding with bismuth-209, has been created artificially. Each radioactive isotope of a decay chain emits alpha or beta particles and possibly also gamma rays.

An example of serial transformation is given by the decay of krypton-90, a fission product of uranium-235. Each step involves beta (minus) decay:



The rate of formation of a daughter radionuclide, or the number of daughter atoms present ( $N_d$ ) from its parent ( $N_p$ ), is equal to the rate of formation of the parent; however, each will have different decay rates. If the half-life of the parent (p) is much longer than that of the daughter (d) and only the parent is present at time zero ( $t = 0$ ), the formation of the daughter product is described by secular equilibrium:

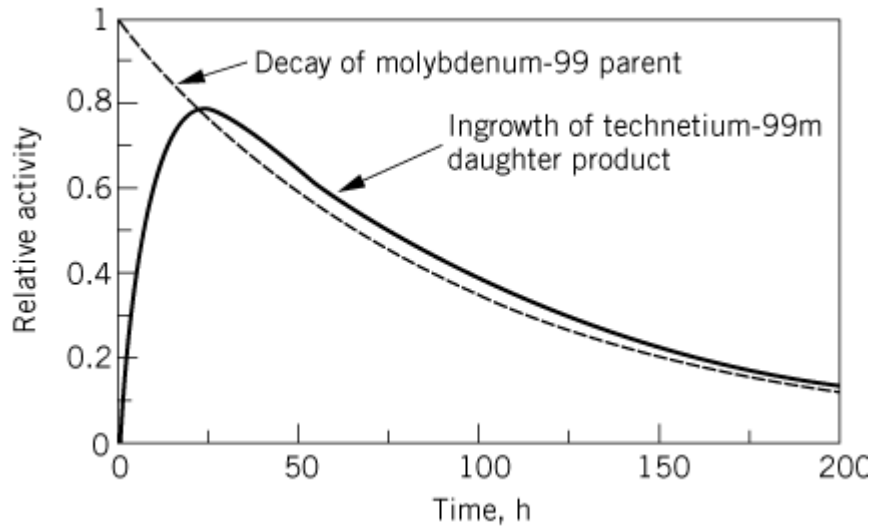
$$(N_d) = \frac{\lambda_p N_p}{\lambda_d} (1 - e^{-\lambda_d t}) \quad (15)$$

At equilibrium,  $\lambda_p N_p = \lambda_d N_d$ , and the activities of each radioactive member are equal ( $A_p = A_d$ ). A state of transient equilibrium exists when the half-life of the parent is somewhat greater than that of the daughter and both half-lives are relatively short. An example is the ingrowth of technetium-99m (the most widely used radioisotope in diagnostic nuclear medicine) from parent molybdenum-99 (Fig. 2). The number of daughter atoms present at any time  $t$  is given by

$$(N_d) = \frac{\lambda_p N_p}{\lambda_d - \lambda_p} (e^{-\lambda_p t} - e^{-\lambda_d t}) \quad (16)$$

When transient equilibrium occurs, the daughter radioisotope undergoes radioactive transformation at the same rate as it is produced, decreasing in amount with time according to the decay rate of its parent, as shown in Figure 2.

**Figure 2.** Radioactive decay of molybdenum-99 and ingrowth of technetium-99 m with time (an example of *transient equilibrium*). The amount of radioactivity from each isotope is plotted as a function of time.



## 5. Alpha Particles

An alpha particle is equivalent to a helium nucleus (two protons plus two neutrons). Alpha decay usually occurs in nuclei of heavy radioactive atoms. An example of alpha decay is the transformation of radium-226 to radon-222 with a half-life of 1600 years:



Alpha particles are essentially monoenergetic. During formation of radon-222, two electrons are ejected from the outermost electron shell. The +2 charged helium ( $\text{He}^{2+}$ ) nucleus is ejected along a straight path, giving up energy and producing ion pairs along the way. When it comes to rest, it captures two electrons from its environment and becomes a stable, neutral helium atom. Most of the 4.78 MeV is given up to the absorbing medium as kinetic energy. During 95% of the time there is emission of a 4.70 MeV alpha particle and 0.08 MeV as radon-222 recoil energy. In the remaining 5% of the time there is emission of a 4.6-MeV alpha particle, and 0.18 MeV is given off as gamma rays.

The range of the alpha particle in a unit-density medium is only a few micrometers. The thickness of skin or a sheet of paper is sufficient to prevent most alpha-particle penetration. The range,  $R$ , of an alpha particle in air at  $0^\circ\text{C}$  and 760 mmHg pressure may be estimated from

$$R(\text{cm}) = 0.56E_\alpha(\text{MeV}) \text{ for } E < 4 \text{ MeV} \quad (18)$$

$$R(\text{cm}) = 1.24E_\alpha(\text{MeV}) - 2.62 \text{ for } 4 < E < 8 \text{ MeV} \quad (19)$$

The range of alpha particles in media other than air may be estimated from

$$R_m(\text{mg}/\text{cm}^2) = 0.56A^{0.33}R_{\text{air}} \quad (20)$$

where  $A$  is the atomic number of the absorbing medium and  $R_{\text{air}}$  is the range of the alpha particle in

air (cm).

Alpha-particle tracks may give rise to low energy, secondary electron tracks, called *delta rays*, which radiate outward to distances of tens of nanometers from the primary particle track.

## 6. Beta Particles

Beta particles are electrons that are ejected from the nucleus of an unstable, beta-emitting atom. Beta particles carry a single negative charge and small mass (only about 1/1800th) that of a proton or neutron). Beta emission appears as the change in a nucleus of one neutron into a proton, and occurs among radionuclides with greater numbers of neutrons than protons in the nucleus. An example was given in equation 1 above.

Beta particles are not monoenergetic, but rather are emitted with a continuous energy distribution, ranging from near zero to the theoretical maximum energy. To comply with the law of conservation of energy, each beta particle is accompanied by emission of a neutrino, whose energy makes up the difference between the theoretical maximum energy of the beta particle and its observed kinetic energy. Gamma rays may accompany beta emission in order to reach the ground energy state of the daughter product. For example, potassium-42 decays 82% of the time by a maximum energy of 3.55 MeV to calcium-42, and 18% of the time by a maximum energy of 2.04 MeV, together with a gamma ray of 1.53 MeV.

When beta particles slow to rest, they transfer a negative charge to the absorber. Their range is a few millimeters in unit-density tissue. The beta energy range may be measured by adding successively thicker absorbers until a count rate cannot be detected. An absorber that stops one-half the beta particles is about one-eighth the range of beta particles. An estimate of the range of beta particles may be obtained from

$$R(\text{mg}/\text{cm}^2) = 412E^{1.265-0.0954 \ln E} \quad (21)$$

for electron energies between 0.01 and 2.5 (MeV).

## 7. Positrons

Positrons are positively charged electrons that are emitted from atomic nuclei where the neutron:proton ratio is low and sufficient energy is not available for alpha-particle decay. Positron emission represents the transformation within the nucleus of a proton into a neutron. In other ways, the emission of positrons is similar to that of beta (minus) particles, which have similar mass and range in tissue. When a positron comes to rest, it quickly combines with an electron, and the two particles annihilate and give off two gamma-ray photons, whose energies are equal to the mass equivalent of the positron plus the electron (two photons of 0.511 MeV). An example of positron decay is the transformation of sodium-22 to neon-22:



## 8. Electron Capture

Some radionuclides decay by the process of electron capture, in which a *K* orbital electron is captured into the nucleus, uniting with a proton or hydrogen nucleus and changing it into a neutron:



When an atom decays by electron capture, a characteristic X ray (photon) is emitted as an electron from an outer orbit falls into the energy level of the captured electron.

## 9. Gamma Rays

Gamma rays are monoenergetic photons or quanta of energy with discrete frequency that are emitted from the nucleus during radioactive decay to remove excess energy. Like X rays, gamma rays are highly penetrating electromagnetic photons of energy. Gamma rays usually accompany beta decay and always accompany positron decay. Gamma rays are attenuated by matter, and the efficiency of the shielding increases with atomic number.

## 10. Internal Conversion

Radioactive decay by internal conversion takes place when an unstable nucleus of a gamma-emitting nucleus gives off excess excitation energy by imparting energy to an orbital *K*- or *L*-shell electron, ejecting it from the atom. Characteristic X rays are emitted as outer-shell orbital electrons collapse inward to fill vacant energy levels produced by ejected electrons. If the characteristic X rays are absorbed by an inner orbital electron, internal conversion may take place, ejecting the electron (called an Auger electron).

## 11. Radiation Absorbed Dose

Radiation absorbed dose is the energy deposited by radiation per unit mass of the absorbing medium. The radiation absorbed dose in units of gray (Gy) to a unit-density medium containing an alpha-emitter is

$$D(\text{Gy}) = 1.61 \times 10^{-10} E_{\alpha}(\text{MeV}) \frac{N_{\alpha}}{g} \quad (24)$$

where  $E_{\alpha}$  is the average alpha-particle energy,  $N_{\alpha}$  is the number of alpha particles emitted in the medium, and  $g$  is the mass of the medium.

The radiation absorbed dose rate ( $D^{\circ}$ , in grays per second) from a beta-emitting radioisotope under conditions of charged-particle equilibrium can be estimated from

$$D^{\circ}(\text{Gy/s}) = 1.61 \times 10^{-10} E_{\beta}(\text{MeV}) \frac{N_{\beta}}{g} \quad (25)$$

where  $E_{\beta}$  is the average beta-particle energy per disintegration and  $N_{\beta}$  is the number of beta-particle disintegrations taking place per second, all per gram of medium. If the mass of the medium is small, some of the beta energy may escape, and conditions for charged-particle equilibrium may not be met.

### Suggestions for Further Reading

H. Cember. (1983) *Introduction to Health Physics*, 2nd ed., Pergamon Press, New York.

R. D. Evans (1955) *The Atomic Nucleus*, McGraw-Hill, New York.

Y. Wang, ed. (1969) *Handbook of Radioactive Nuclides*, Chemical Rubber Company, Cleveland, Ohio.



## Radioimmunoassay

Radioimmunoassay (RIA) is a technique for determining the content of a particular molecule in a sample, based on displacement of a radioactive [antigen](#) from an [antibody](#) or **antiserum**. Berson and Yalow and their colleagues developed the RIA method, and first applied it to the assay of human and animal [insulins](#) in plasma (1, 2). RIA has subsequently been used most extensively in endocrinology, as the technique is sufficiently sensitive for detection of vanishingly small quantities of [hormones](#), but it is by no means restricted to hormones and has also been applied to all classes of biomolecule. Many modifications of the original technique are in common use. Here we discuss the classical RIA and several of the most common variations.

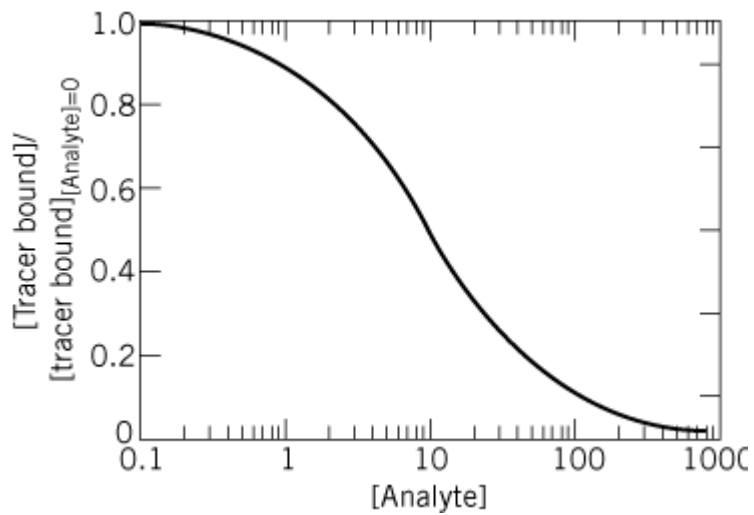
### 1. Classic RIA

Materials necessary for a classic RIA are:

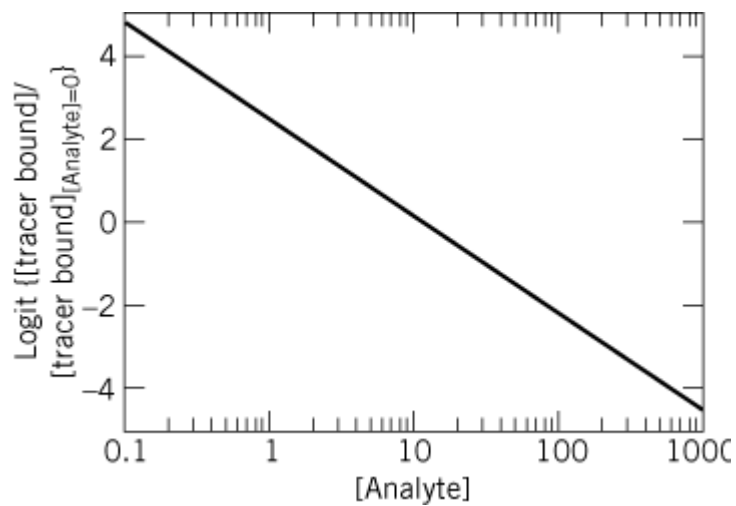
- The unknown, that is, a sample of a chemically known antigen whose concentration is unknown and to be determined.
- Standards for calibrating the assay. These must be chemically identical to the unknown and of known concentration.
- A polyclonal or [monoclonal antibody](#) preparation that binds the unknown specifically.
- A **radiolabeled** ligand, termed a *tracer*, that binds to the antibody competitively with the unknown. Most often the tracer is chemically identical to the unknown, except for the radioactive tag, which usually is, or contains, an [iodine isotope](#). Identity is not required, however, as long as binding of the tracer and the unknown to the antibody are mutually exclusive.
- A means for rapidly separating bound and free tracer. Separation is often accomplished using an antiglobulin antibody to **immunoprecipitate** the antibody–tracer complex, or an immobilized affinity adsorbent that reacts with the antibody constant region, such as [protein A](#).

The RIA is initiated by adding the antibody to mixtures of a constant amount of tracer and increasing amounts of either the calibration standard or the unknown. Ideally, the tracer will be present in negligible molar quantity, relative to the **ligand-binding** capacity of the antibody, so that saturating half of the antibody sites with standard or unknown will displace half of the tracer. The necessary controls are (1) antibody added to tracer with a large excess of the calibration standard, to measure background radioactivity under conditions of complete tracer displacement, and (2) antibody added to tracer alone, to measure the radioactivity bound when no displacement occurs. After a suitable time for equilibration of the mixtures, bound and free tracer are separated, and one or both are measured. A dose–response curve, based on the calibration standards, can be plotted directly as the fraction of bound tracer retained or displaced as a function of the amount of standard added. This fraction is relative to the quantity of tracer bound in the absence of added unlabeled standard, not to the total amount of tracer added to the reaction. An example of such a direct plot is shown in [Figure 1](#).

**Figure 1.** Dose–response curve typical of a classic RIA. Addition of increasing amounts of ligand progressively displaces the radioactive tracer from the antibody in the assay, generating an S-shaped curve. The abscissa in the figure is a concentration scale in units of the amount of unlabeled antigen added. The ordinate is a fractional scale, where 1.0 corresponds to the tracer bound to antibody in the absence of added antigen.



**Figure 2.** Logit transformation of RIA data from Figure 1. The logit transform converts the S-shaped dose–response curve into a linear form. The advantage of the logit plot is that a straight line is more readily analyzed, and outlying points are more easily identified, than for the direct plot.



A preferable treatment of data is to construct a logit–log plot (3). The logit transformation for a fraction  $f$  is defined as

$$\frac{\text{logit } f}{1 - f} = \ln f \quad (1)$$

For RIA data, this corresponds to

$$\text{logit } f = \ln \frac{\text{fraction of tracer retained}}{\text{fraction of tracer displaced}} \quad (2)$$

A plot of  $\text{logit } f$  versus  $\log_{10}$  (standard) is often linear over much of its range; hence it is more easily analyzed than the S-shaped direct plot.

On the basis of the dose–response behavior of the standards, the amount of tracer displaced by the

unknown can be related graphically or algebraically to the quantity of antigen that must have been present in the unknown. Software for statistical analysis of RIA data is also widely available.

## 2. Solid-Phase RIA

This method is similar to the classical RIA, but with a great practical simplification: as the first step of the assay, the antibody is adsorbed to the surface of plastic assay tubes (4). (Other media such as beads and disks can also be used.) Adsorption is accomplished simply by adding a dilute antibody solution to the tube: a considerable amount of the antibody sticks noncovalently, but essentially irreversibly, to the plastic. The remaining components of the assay mixture—tracer, standards, and antigen—are initially present in solution. On their addition to the tubes, the tracer and antigen bind to the antibody, in ratios determined by their relative concentrations. Unbound assay components are subsequently removed, simply by decanting and washing. The radioactivity retained in the tubes is measured, and the data are analyzed as for the classic RIA.

## 3. Immunoradiometric Assay

The classic RIA is a competition assay, in which excess sample or antigen competes with tracer in binding to a limiting amount of antibody. An immunoradiometric assay (IRMA) is a noncompetitive design in which the antigen itself is the limiting quantity; hence IRMA is in principle orders of magnitude more sensitive than the classical assay (5). An experimentally important aspect of IRMA is that the antibody is the labeled component of the assay, rather than the tracer. Radiolabeling antibody by **iodination** is technically straightforward and often accomplished more readily than synthesis of labeled tracer (6-9). In the original IRMA method, antigen and labeled antibody are incubated together, and a complex is allowed to form. The mixture is next exposed to an immunosorbent, usually consisting of immobilized antigen, that will complex free antibody through the combining site. The radiolabeled antibody binds, unless it is already binding the antigen. Bivalency of IgG introduces a complication, since one IgG molecule has the ability to bind both antigen and immunosorbent. Nevertheless, measuring the radioactivity in either the solid or solution phase, and comparison of test samples with the standard curve measured using known amounts of antigen, gives an accurate measure of the quantity of antigen in the unknown.

## 4. Sandwich RIA

The sandwich RIA technique (10, 11), also known as a “two-site IRMA,” is a solid-phase method that uses two different antibodies, one radiolabeled, to detect a limiting amount of antigen. Both antibodies must be specific for the antigen, but they must bind at independent sites. As a first step, an unlabeled “capture” antibody specific for the antigen is immobilized. (An alternative to using immobilized reagents is to immunoprecipitate the antigen or capture antibody.) Antigen that is unlabeled and free in solution is added and allowed to bind to the capture antibody. Unbound components are washed away. The labeled antibody used for detection is added, and it also binds to the now-immobilized antigen. The unbound detection antibody is washed away, and the radioactivity remaining is measured. The antigen allows a “sandwich” structure to form, and the amount of labeled antibody incorporated into the sandwich is a direct measure of the amount of antigen in the sample.

This assay requires that the two antibody molecules bind to the antigen simultaneously, so the sandwich assay is most appropriate for analysis of [protein](#) or other large, **macromolecular** antigens. For the sandwich structure to form, the two antibodies must also bind to nonoverlapping **epitopes** on the antigen, so the use of [monoclonal antibodies](#) for capture, or for both capture and detection, is advantageous.

## Bibliography

1. S. A. Berson, R. S. Yalow, A. Bauman, M. A. Rothschild, and K. Newerly (1956) J. Clin.

- Invest. **35**, 170–190.
2. R. S. Yalow and S. A. Berson (1960) *J. Clin. Invest.* **39**, 1157–1175.
  3. D. Rodbard, P. L. Rayford, J. A. Cooper, and G. T. Ross (1968) *J. Clin. Endocrin. Metab.* **28**, 1412–1418.
  4. K. Catt and G. W. Tregear (1967) *Science* **158**, 1570–1572.
  5. L. E. M. Miles and C. N. Hales (1968) *Nature* **219**, 186–189.
  6. W. M. Hunter and F. C. Greenwood (1962) *Nature* **194**, 495–496.
  7. J. I. Thorell and B. G. Johansson (1971) *Biochim. Biophys. Acta* **251**, 363–369.
  8. P. J. Fraker and J. C. Speck (1978) *Biochem. Biophys. Res. Commun.* **80**, 849–857.
  9. A. E. Bolton and W. M. Hunter (1973) *Biochem. J.* **133**, 529–539.
  10. E. Haberman (1970) *Z. Klin. Chem. Klin. Biochem.* **8**, 51–55.
  11. L. Wide (1969) *Acta Endocrinol. Supp.* **142**, 207–221.

### Suggestions for Further Reading

12. A. E. Bolton and W. M. Hunter (1986) in *Handbook of Experimental Immunology*, D. M. Weir, ed., Blackwell, London, pp. 26.1–26.56. (Comprehensive review of RIA methodology.)
13. R. Edwards (1997) "Radiolabelled immunoassays, in" *Principles and Practice of Immunoassay*, C. P. Price and D. J. Newman, eds., Macmillan, London, pp. 325–348.
14. D. Rodbard and G. R. Frazier (1975) Statistical analysis of radioligand assay data, *Meth. Enzymol.* **37**, 1–22. (Comprehensive discussion of RIA models, graphic display, and data analysis.)
15. R. S. Yalow (1978) Radioimmunoassay: a probe for the fine structure of biologic systems, *Science* **200**, 1236–1245. (Nobel Prize lecture by one of the inventors of RIA.)

## Radioisotopes

Isotopes are atomic species of the same atomic number (belonging to the same element) that have different mass numbers. The number of elements in the periodic table is about 110, and each one has more than one isotope; the total number of known isotopes is more than 1500. Each isotope of a given element has the same number of protons in its atomic nucleus, but differs in the number of neutrons in its nucleus. Isotopes of an element cannot be distinguished chemically because they have the same electronic structure and undergo the same chemical reactions.

Although some isotopes are stable, the nuclear configurations of radioisotopes (or radionuclides) are unstable, and they spontaneously undergo a radioactive transformation (or decay) to a more stable energy state (see [Radioactivity](#)). The half-life of each radioisotope is the time required for exactly one-half of the atoms to undergo radioactive transformation. Radioisotopes may decay to either stable or other radioactive species. Decay from one radioisotope to another is called a decay series.

Radioisotopes occur in small amounts in nature as the result of the decay of long-lived primordial materials (such as uranium-238). Atmospheric reactions with solar particles also produce radioactive species. Approximately 50 radionuclides occur naturally in the atmosphere, ocean, or the earth's crust; these include carbon-14, potassium-40, radon-222, radium-226, and uranium-238.

Radioisotopes can be produced artificially by nuclear high energy reactions that combine atomic

nuclei. The first human-made or artificial radioisotopes were made by Frederic and Irene Joliot-Curie in 1933, who irradiated a thin aluminum foil with alpha particles and observed tracks in a cloud chamber that diminished in intensity with a half-life of about 3 min, due to phosphorus-30 beta-plus decay. When they replaced aluminum with a boron foil, they found new activity, with a half-life of 14 min, due to nitrogen-13 beta-plus decay (1). In 1934, Lawrence produced small amounts of new radioisotopes at the Berkeley cyclotron using deuteron bombardment reactions on stable-element targets. Fermi produced heavier radioisotopes of the same element by neutron bombardment (or activation), also in 1934 (1). Hevesey conducted the first biological studies with radioisotope tracers. These developments made it possible to discover, produce, and test a large number of scientifically significant radioisotopes during the decade that followed. Radioisotope production continues today. In general, neutron-rich radioisotopes are produced in nuclear reactors, whereas neutron-lean radioisotopes are produced in charged-particle accelerators. A carrier-free radioisotope is one that is not produced or mixed with any other isotope of the same element. The specific activity of a radioisotope preparation is the radioactivity (bequerels or curies) exhibited per unit mass or volume of the radioactive material.

Radioisotopes can be detected easily and identified using radiation detection instruments or photographic film (see [Autoradiography](#) and [Fluorography](#)). Therefore, they have numerous practical applications in the physical, chemical, and biomedical sciences. Among the most important applications in biomedical research are those that involve the tagging of a radioisotope to a biomolecule to permit tracking of the molecule in reaction processes and metabolism. Animal tissues containing radioisotopes are analyzed by nuclear radiation-detection techniques such as liquid scintillation counting, gamma spectroscopy, alpha spectrometry, and neutron activation analysis—depending on the relevant radiation emission.

In studies of life processes, the most important radioisotopes are those of **hydrogen, carbon, sulfur,** and **phosphorous**, because these elements are present in practically all cellular components essential to maintaining life. Some of the more common radionuclides used in biomedical research are given in Table 1, together with their mass number, physical half-life, beta-particle yield and energies, and associated gamma-ray energies.

**Table 1. Common Radioisotopes Used in Biomedical Research, with Principal Radioactive Emissions, Yields, and Energies (2)**

| Element    | Mass (amu) <sup>a</sup> | Physical Half-Life | Beta Particle Yield | Average Beta Energy (MeV) <sup>b</sup> | Gamma-Ray Yield | Gamma Energy (MeV) |
|------------|-------------------------|--------------------|---------------------|--|-----------------|--------------------|
| Hydrogen   | 3                       | 12.35 years        | 1.0                 | 0.00568                                | —               | —                  |
| Carbon     | 11                      | 20.38 min          | 0.998               | 0.386                                  | 2.0             | 0.511              |
| Carbon     | 14                      | 5730 years         | 1.0                 | 0.0495                                 | —               | —                  |
| Phosphorus | 32                      | 14.29 days         | 1.0                 | 0.695                                  | —               | —                  |
| Sulfur     | 35                      | 87.44 days         | 1.0                 | 0.076                                  | —               | —                  |
| Calcium    | 45                      | 163 days           | 1.0                 | 0.0771                                 | —               | —                  |
| Iodine     | 125                     | 59.6 days          | —                   | —                                      | 0.0667          | 0.0355             |
| Iodine     | 131                     | 8.021 days         | 1.0                 | 0.182                                  | 0.0606          | 0.284              |
|            |                         |                    |                     |  | 0.812           | 0.364              |

<sup>a</sup> Atomic mass units.

<sup>b</sup> Million electronvolts.

In molecular biology, the most important application of radioisotopes is the radioactive labeling of nucleic acids and proteins. Radioactively labeled cells, such as **organelles** and **chromosomes**, can be imaged on high speed X-ray film in a process called **autoradiography**.

### Bibliography

1. M. Brucer (1990) *A Chronology of Nuclear Medicine*, Heritage Publications, Inc., St. Louis, Mo.
2. National Council on Radiation Protection and Measurements (1985) *A Handbook of Radioactivity Measurements Procedures*, 2nd ed., National Council on Radiation Protection and Measurements, Bethesda, Md.

### Suggestions for Further Reading

3. G. de Hevesey (1961) *Adventures in Radioisotope Research, the Collected Works of George de Hevesy* (2 vols.), Pergamon Press, New York.
4. R. D. Evans (1955) *The Atomic Nucleus*, McGraw-Hill, New York.
5. J. C. Harbert, W. C. Eckelman, and R. D. Neumann (1996) *Nuclear Medicine: Diagnosis and Therapy*, Thieme Medical Publishers, Inc., New York.
6. E. J. Hall (1994) *Radiobiology for the Radiologist*, J. B. Lippincott Co., Philadelphia.
7. Y. Wang, ed. (1969) *Handbook of Radioactive Nuclides*, Chemical Rubber Company, Cleveland, Ohio.




## Radius Of Gyration

The radius of gyration,  $R_g$ , is a parameter that can be determined for biological molecules in solution using **small-angle scattering** techniques with X-rays, neutrons, or light. The  $R_g$  for a particle is defined as the root-mean-square distance of all elemental scattering volumes from their center of mass weighted by their scattering densities, and it is a simple but useful measure of the overall shape of the particle.  $R_g$  values for homogeneous geometrical shapes can be calculated analytically using the relationships in Table 1(1). These relationships can be very useful in interpreting experimentally determined  $R_g$  values. Table 2 shows how  $R_g$  values vary for a single protein of a specific molecular weight and partial specific volume, assuming increasingly asymmetric shapes from a perfect sphere to a cylindrical rod. Also shown in Table 1 are the  $R_g$  values for a six-subunit assembly assuming various arrangements of the subunits. It is readily seen that  $R_g$  gives a simple measure of the asymmetry in a particle. The mathematical formalisms used for determining  $R_g$  values from neutron and X-ray scattering data are somewhat different to those used for **light scattering**. We describe here the conventional equations used in the interpretation of neutron and X-ray scattering data. The formalisms used in **light scattering** are given in its own topic.

**Table 1. Formula for Calculating  $R_g$  Values of Simple Geometric Shapes**

|  |   |
|--|---|
| Sphere radius $R$                                | $R_g^2 = \frac{3}{5} R^2$                                 |
| Hollow-sphere radii $R_1$ and $R_2$              | $R_g^2 = \frac{3}{5} \frac{R_2^5 - R_1^5}{R_2^3 - R_1^3}$ |
| Ellipsoid, semi-axes $a, b, c$                   | $R_g^2 = \frac{a^2 + b^2 + c^2}{5}$                       |
| Prism, edge lengths $A, B, C$                    | $R_g^2 = \frac{A^2 + B^2 + C^2}{12}$                      |
| Elliptic cylinder, semi-axes $a, b$ , height $h$ | $R_g^2 = \frac{a^2 + b^2}{4} + \frac{h^2}{12}$            |
| Hollow cylinder, height $h$ , radii $R_1, R_2$   | $R_g^2 = \frac{R_1^2 + R_2^2}{2} + \frac{h^2}{12}$        |

**Table 2.  $R_g$  Values for Various Shaped Objects**

|   | $R_g^a$ (Å) | Maximum Linear Dimension (Å) |
|---|-------------|------------------------------|
| Sphere, radius 10 Å, volume 4189 Å <sup>3</sup>   | 7.75        | 20                           |
| Prolate ellipsoid, semi-axes 6.325, 6.325, 25 Å, volume 4189 Å <sup>3</sup>                     | 11.87       | 50                           |
| Cylindrical rod, cross-sectional radius 3.6515 Å, length 100 Å, volume 4189 Å <sup>3</sup>      | 28.98       | 100.27                       |
| Different arrangements of six identical spheres, radii 10 Å, total volume 25,134 Å <sup>3</sup> | 35          | 120                          |
|              | 20.66       | ~70                          |
|              | 21.44       | 60                           |
|              |             |                              |

<sup>a</sup>  $R_g$  values calculated using formulae in Table 1 for the single domain structures, each of which has

the same volume.  $R_g$  values for the multiple domain structures were calculated using the relationship  $R_g^2 = R_{gm}^2 + \frac{1}{N} \sum_i l_i^2$ , where  $R_{gm}$  is the radius of gyration of each identical monomer,  $l_i$  are the distances of each monomer from the center of mass of the entire assembly,  $N$  is the number of monomers, and the sum is over all monomers (6). Note that the lower rosette arrangement has a larger radius of gyration than the rectangular array above it, even though it has a smaller maximum linear dimension. This is because the rosette arrangement has a “hole” in the middle.

## 1. Determination of $R_g$ from Small-Angle Scattering Data

The [small-angle scattering](#) of X-rays or neutrons yields the [scattering intensity distribution](#),  $I(Q)$ , which is related to the vector length or pair distribution function by a Fourier transformation:

$$I(Q) = 4\pi \int P(r) \frac{\sin Qr}{Qr} dr \quad (1)$$

$$P(r) = \frac{1}{2\pi^2} \int I(Q) Q \cdot r \sin(Q \cdot r) dQ \quad (2)$$

$Q$  is the amplitude of the scattering vector, which is equal to  $4\pi(\sin\theta)/\lambda$  ( $\lambda$  is wavelength of the scattered radiation, and  $\theta$  is half the scattering angle).  $P(r)$  is the probable frequency distribution of all possible vector lengths,  $r$ , between scattering centers (or small volume elements) within a particle, weighted by the product of the scattering densities at the respective centers.  $R_g$  can be calculated as the second moment of the pair distribution function  $P(r)$ :

$$R_g^2 = \frac{\int P(r)r^2 dr}{2 \int P(r) dr} \quad (3)$$

Alternatively, one can determine  $R_g$  from scattering data using the Guinier approximation. In 1939 Guinier (2) showed that for sufficiently small  $Q$ -values:

$$I(Q) = I(0)e^{-\frac{Q^2 R_g^2}{3}} \quad (4)$$

A plot of  $\log [I(Q)]$  versus  $Q^2$  thus yields  $R_g$  and the zero angle scatter,  $I(0)$ , from the slope and intercept, respectively.  $[I(Q)]_{Q \rightarrow 0}$  is directly proportional to the square of the molecular weight,  $M$ , of the scattering particle. For spherical particles, the Guinier approximation holds for  $Q < 1.3/R_g$ , and as the scattering particle becomes increasingly asymmetric, the approximation breaks down at even lower  $Q$  values.

When one dimension of a particle is much greater than the other two (eg, a rod), Guinier (2) also showed that one can approximate the scattering for certain small-  $Q$  values as

$$QI(Q) = I_c(0)e^{-\frac{Q^2 R_c^2}{2}} \quad (5)$$

where  $R_c$  is the radius of gyration of cross section  $[QI_c(Q)]_{Q \rightarrow 0}$  is directly proportional to the mass per unit length,  $M_L$ , of the scattering particle. For a rod of radius  $R$ ,  $R_c = R\sqrt{2}$ . The range of  $Q$  for which equation (5) is valid depends on the shape of the cross section and the aspect ratio of the rod



(3). For infinite cylindrical rods, equation (4) is valid for  $Q$  values  $<0.8/R_c$ . For finite rods there will be a “rollover” near the origin. Explicitly for rods of radius  $R$  and aspect ratio  $A = L/2R$ , the scattering intensity will decrease for values of  $Q < 5/2AR$  that is, the higher the aspect ratio, the lower the value of  $Q$  where the rollover begins. For a rod-shaped object,  $R_g$  and  $R_c$  are related by

$$R_g^2 - R_c^2 = \frac{L^2}{2} \quad (6)$$

For particles with two dimensions much greater than the third (eg, a disk), the small  $Q$  scattering can be approximated as

$$Q^2 I(Q) = I_t(0)e^{-Q^2 R_t^2} \quad (7)$$

where  $R_t$  is the radius of gyration of thickness. For a disk of thickness  $T$ ,  $R_t = T/\sqrt{12}$ . Equations (5) and (7) break down for particles with low axial ratios, and so they must be used with caution.

### 1.1. Changes in $R_g$ Values Give Insights into Biological Function

In a series of scattering experiments on the dumbbell-shaped calcium-binding protein calmodulin, it was shown that upon binding a wide variety of amphipathic helices, calmodulin undergoes a dramatic conformational collapse involving its two globular lobes coming into close contact. The collapse is facilitated by a flexible helix linking calmodulin's two globular domains, and it is characterized by an approximately 25% reduction in  $R_g$  (reviewed in Ref. 4). This conformational flexibility has proven key to understanding how calmodulin binds and activates a wide variety of target enzymes. Another example of the utility of determining  $R_g$  values is the study by Mangel et al. (5) in which they characterized an extremely large ligand induced conformational change in native human Glu-plasminogen which was shown to have an  $R_g$  value of 39 Å. Upon occupation of a weak lysine-binding site that regulates the activation of plasminogen, the protein's shape irreversibly changes to give an  $R_g$  value of 56 Å. This change in shape is achieved without any accompanying change in secondary structure, and the data have been interpreted in terms of a rearrangement of the five domains of the protein structure from a compact, closed structure that is relatively rigid, to an open flexible structure in which the individual domains no longer interact with each other and hence are more accessible. The increased flexibility in the open form is postulated to account for the observation that this form is more readily activated, because it would facilitate urokinase binding. The open form is also proposed to be key to plasminogen's role in fibrinolysis.

### Bibliography

1. P. Mittelbach (1964) *Acta Phys. Austriaca* **19**, 53–102.
2. A. Guinier (1939) *Ann. Phys. (Paris)* **12**, 161.
3. R. P. Hjelm (1985) *J. Appl. Crystallogr.* **18**, 452–460.
4. J. Trehwella (1994) In *Structural Biology: State of the Art*, Vol. I, Proceedings of the Eighth Convention, State University of New York, Albany, 1993 (R. H. Sarma and M. H. Sarma, eds.), Adenine Press, New York, pp. 43–57.
5. W. F. Mangel, B. Lin, and V. Ramakrishnan (1990) *Science* **248**, 69–73.
6. P. Zipper and H. Durchschlag (1980) *Monatsh. Chem.* **111**, 1367–1390.

### Suggestions for Further Reading

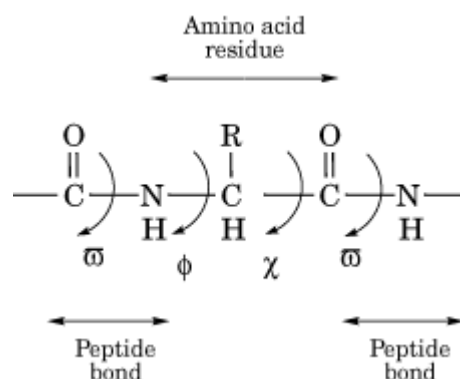
7. C. R. Cantor and P. R. Schimmel (1980) In *Biophysical Chemistry, Techniques for the Study of Biological Structure and Function*, W. H. Freeman, San Francisco, pp. 812–814.

- O. Glatter and O. Kratky (1982) *Small-Angle X-Ray Scattering*, Academic Press, New York.
- A. Guinier and G. Fournet (1955) *Small-Angle Scattering of X-Rays*, Wiley, New York.

## Ramachandran Plot

The Ramachandran plot is a two-dimensional graph of the phi ( $\phi$ ) and psi ( $\psi$ ) [backbone](#) angles for each [amino acid](#) residue of a [protein](#); it is a simple method of assessing the quality of a [protein structure](#). It is named after G. N. Ramachandran, who first determined the values of  $\phi$  and  $\psi$  that are theoretically permitted ([1](#)). The symbols  $\phi$  and  $\psi$  are used to represent the dihedral angles of the backbone N–C $\alpha$  and C $\alpha$ –C bonds, respectively (Fig. [1](#)). Omega ( $\omega$ ) is the angle of the [peptide bond](#), C–N, but this bond has partial double bond character, so  $\omega$  can only adopt values close to 180° (called *trans*) or 0° (*cis*). Of these two, the *trans* peptide  $\omega$  conformation is observed in most cases, unless the following residue is [proline](#) (see [Cis/Trans Isomerization](#)).

**Figure 1.** Section of a polypeptide chain showing two peptide bonds (–CONH–) and one amino acid residue (–NH–C $\alpha$ –(R)(H)–CO–). The functional group R varies, depending on the type of amino acid. The dihedral angles  $\omega$ ,  $\phi$  and  $\psi$  of the peptide backbone are labeled.

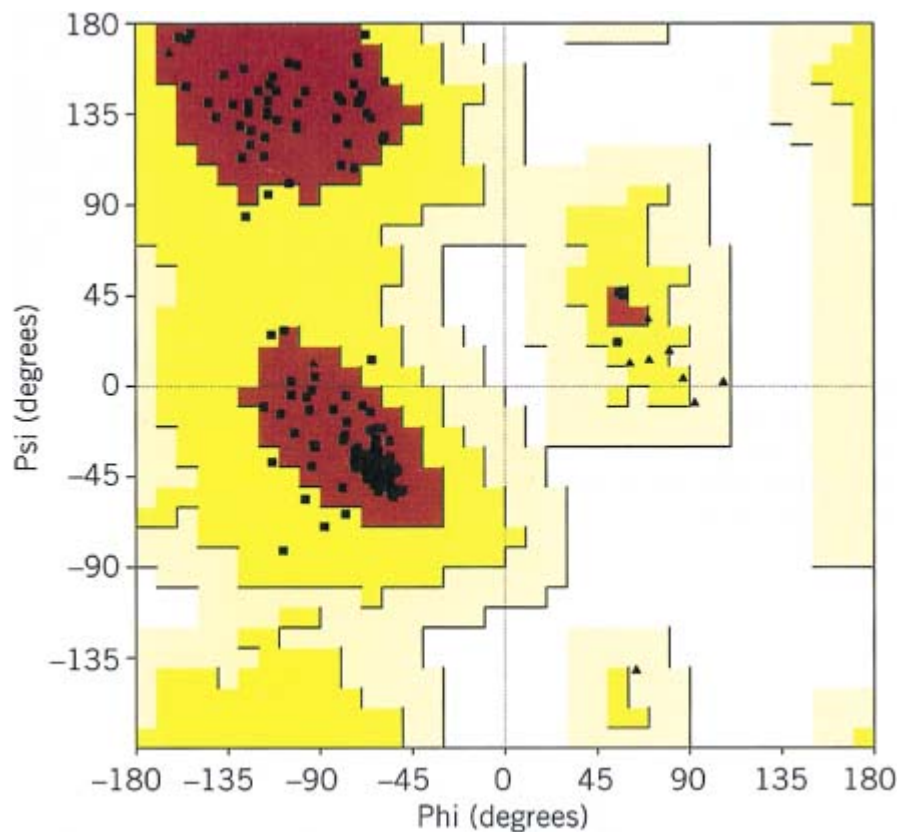


The bonds represented by  $\phi$  and  $\psi$  are single bonds, and in principle all bond rotations are feasible. Only certain combinations of the backbone angles  $\phi/\psi$  are possible in [polypeptide chains](#), however, due to steric restraints imposed by the amino acid [side chains](#). These are defined as the most favorable or “allowed” regions of the Ramachandran plot, and they include the **conformations** of regular **secondary structure** such as **a-helix** or **b-strand**. In general, the  $\phi$  angle only adopts negative values. Of the 20 natural amino acids, 18 have only slightly differing allowed values of  $\phi$  and  $\psi$ . The [glycine](#) residue, with only a hydrogen for its side chain, has a high degree of backbone conformational flexibility and can adopt  $\phi/\psi$  angles that are forbidden for other amino acid residues (positive  $\phi$  angles, for example). In contrast, the cyclic nature of the [proline](#) residue restricts its value of  $\phi$  to approximately –60°.

The residues in natural folded proteins generally adopt fully allowed values of  $\phi$  and  $\psi$ ; exceptions usually have functional importance. Generally, over 90% of nonglycine amino acid residues in a protein structure have backbone conformations corresponding to the a, b, or turn types of secondary

structure. This can be visualized using plots of the observed  $\phi/\psi$  combinations for each residue of a protein structure (Fig. 2). The Ramachandran plot is a very useful means of assessing the quality of a protein structure and of identifying residues with unusual backbone conformations.

**Figure 2.** Ramachandran plot for a typical protein structure. The  $\phi/\psi$  values of each nonglycine residue are plotted as squares, and those for glycine residues are shown as triangles. The allowed or most favored regions of the Ramachandran plot are shown in red. The top left section corresponds to  $\beta$ -strand secondary structure, the middle left section corresponds to  $\alpha$ -helix and the small middle right section corresponds to left-handed helix. Over 90% of the nonglycine residues fall within the most favored (red) regions of the Ramachandran plot. Additional allowed regions are shown in varying shades of yellow. This figure was generated using Procheck (2). See color insert.



#### Bibliography

1. C. Ramakrishnan and G. N. Ramachandran (1965) *Biophys J.* **5**, 909–933.
2. R. A. Laskowski, M. W. MacArthur, D. S. Moss, and J. M. Thornton (1993) *J. Appl. Crystallogr.* **26**, 283–291.

#### Suggestions for Further Reading

3. G. J. Kleywegt and T. A. Jones (1996) Phi/Psi-chology: Ramachandran revisited. *Structure* **4**, 1395–1400. (Update on the use of  $\phi/\psi$  plots for protein structure analysis.)
4. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
5. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, and New York.
6. N. J. Darby and T. E. Creighton (1993) *Protein Structure*, IRL Press, Oxford, U.K.

## Random Coil

Random coil describes the multiconformational state of **peptides** and unfolded [proteins](#) and [polypeptide chains](#). In addition, it is sometimes used to refer to the unstructured **conformations** adopted by the [N-terminus](#), [C-terminus](#), and certain loop regions of stable folded [proteins](#). A polypeptide chain in a random coil state adopts multiple interconverting conformations in which the average conformation of each **amino acid residue** is independent of the conformations of all residues other than those immediately adjacent in the [primary structure](#). Proteins usually adopt the random coil state under strongly **denaturing** conditions, such as high concentrations of [guanidinium salts](#) or [urea](#).

## Random X-Inactivation

It is biologically necessary to make the expression of genes from the two **X-chromosomes** in female cells comparable to that from the single X-chromosome in male cells. To achieve this goal one of the two X-chromosomes in female cells is inactivated (see [X-Chromosome Inactivation](#)). The two X-chromosomes from the female cell derive from the male gamete (paternal,  $X^P$ ) and from the female gamete (maternal,  $X^M$ ). Both the  $X^M$  and  $X^P$  chromosomes are functional in the very early female mammalian [embryo](#). Once cytodifferentiation starts in the embryonic and extraembryonic tissues of the embryo, however, X-chromosome inactivation occurs (1). In the mouse embryo, the trophoctoderm and primitive endoderm preferentially inactivate the paternally contributed X-chromosome  $X^P$ . In all other cell lineages that eventually give rise to the adult animal, there is a random inactivation of  $X^P$  and  $X^M$  (2). This random X-inactivation occurs in all eutherian (placental) mammals, making females functional mosaics of heterozygous X-linked genes. In marsupials, in contrast, the paternally-derived X-chromosome is always inactivated in the female animal (3).

### Bibliography

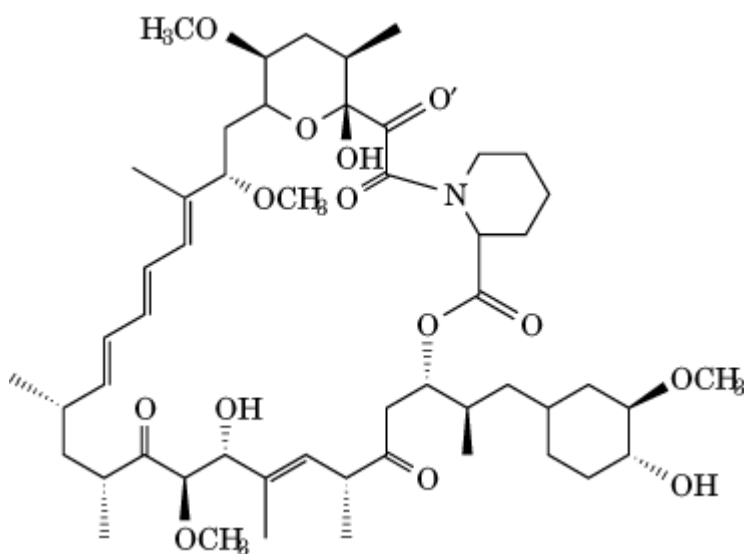
1. S. G. Grant and V. M. Chapman (1988) *Ann. Rev. Genet.* **22**, 199–217.
2. W. I. Frels and V. M. Chapman (1980) *J. Embryol. Exp. Morph.* **56**, 179–189.
3. G. B. Sharman (1971) *Nature* **230**, 231–232.

## Rapamycin

Rapamycin (also known as Sirolimus, Rapamune) is a 31-membered macrocyclic lactam **antibiotic** of the **FK506**-type (Fig. 1), isolated from the soil actinomycete *Streptomyces hygroscopicus* in 1975 (1). The imidic  $\alpha$ -ketoacyl-pipecolinyl bond is subject to [cis/trans isomerization](#), leading to 38 % *cis* isomer in aqueous solution (2). It was initially evaluated as a potent antifungal agent, but later it was

found that it exerts antiproliferative action by disrupting the signaling pathway of [lymphokines](#) (IL-2, IL-4, IL-6) which promote [T cell](#), B cell, and mast cells. Some nonimmune cell lines are also blocked in proliferation, although rapamycin does not exhibit the properties of a general inhibitor of cell proliferation. The mode of action in these specific effects on the immune response is different from that of the immunosuppressants [cyclosporin A](#) and FK506, in that cells are arrested in the late **G<sub>1</sub> phase**, just prior to the S phase. Moreover, rapamycin suppresses a much wider spectrum of cellular signals. Supporting the belief that it has a pathway distinct from that of other immunosuppressants, only rapamycin inhibits the chronic organ rejection process characterized by graft arteriosclerosis, thus prolonging graft survival. In addition to inhibiting the lymphokine response, the production of antibodies and lymphokines through Ca<sup>2+</sup>-independent pathway is impaired (3).

**Figure 1.** The structure of rapamycin.



Typically, these biological effects are obtained at rapamycin levels (picomolar to micromolar) below those at which it is toxic. Rapamycin does not always act synergistically with cyclosporin A and FK506, as in the inhibition of T cell and B cell proliferation, but it often affects identical intracellular signaling events (4). Sometimes, FK506 and rapamycin reciprocally antagonize the parameters tested, suggesting a common cellular receptor for both compounds. Indeed, at doses in the picomolar range of concentration, both drugs bind reversibly to FK506-binding proteins (FKBP). Among mammalian FKBP, the cytosolic FKBP12 is responsible for mediating rapamycin effects (5). Binding always occurs with a concomitant inhibition of the [peptidyl prolyl cis/trans isomerase](#) (PPIase) activity of the FKBP, but this inhibition fails to account for the immunosuppressive effects. Instead, several lines of evidence show an indirect positive action of the rapamycin/FKBP12 complex, because rapamycin treatment induces the inhibition of several cellular serine/threonine [kinases](#), including mammalian p70 S6 kinase.

The upstream target, which may represent the protein that physically interacts with the drug/receptor complex, is the 289-kDa RAFT1 protein in mammals (also known as FRAP or RAPT1) (6-8). Similar to its **homologues** TOR1 and TOR2 in yeast, this protein contains a C-terminal **domain** with amino acid homology to several [phosphatidylinositol](#) 4- and 3- lipid kinases. It can be postulated that a putative protein/lipid kinase activity of RAFT1/TOR takes part in the **phosphorylation** cascade that governs the function of p70 s6 kinase in controlling **protein biosynthesis**. The effector region of the rapamycin/FKBP12 complex consists of a solvent-exposed part of rapamycin along with residues

from FKBP12 in the loops starting with residues 40 and 80. This region is the unique feature of the complex for targeting the RAFT1/TOR family (9).

## Bibliography

1. C. Vezina, A Kudelski, and S. N. Sehgal (1975) *J. Antibiot.* **28**, 721–732.
2. T. Zarnt, K. Lang, H. Burtscher, and G. Fischer (1995) *Biochem. J.* **305**, 159–164.
3. F. J. Dumont and Q. X. Su, (1995) *Life Sci.* **58**, 373–395.
4. S. Wera et al. (1995) *Endocrine Res.* **21**, 623–633.
5. D. A. Fruman (1995) *Eur. J. Immunol.* **25**, 563–571.
6. M. I. Chiu, H. Katz, and V. Berlin (1994) *Proc. Natl Acad. Sci. USA* **91**, 12574–12578.
7. E. J. Brown et al. (1994) *Nature* **369**, 756–758.
8. Y. Q. Chen et al. (1994) *Biochem. Biophys. Res. Commun.* **203**, 1–7.
9. J. A. Kallen, R. Sedrani, and S. Cottens (1996) *J. Amer. Chem. Soc.* **118**, 5857–5861.

## Suggestions for Further Reading

10. S. N. Sehgal and C. C. Bansbach (1993) Rapamycin—In vitro profile of a new immunosuppressive macrolide, *Ann. NY Acad. Sci.* **685**, 58–66.
11. F. J. Dumont and Q. Su (1996) Mechanism of action of the immunosuppressant rapamycin, *Life Sci.* **58**, 373–395; most comprehensive data collection.
12. S. N. Sehgal (1995) Rapamune (Sirolimus, rapamycin): An overview and the mechanism of action, *Therapeutic Drugs Monitoring* **17**, 660–665.

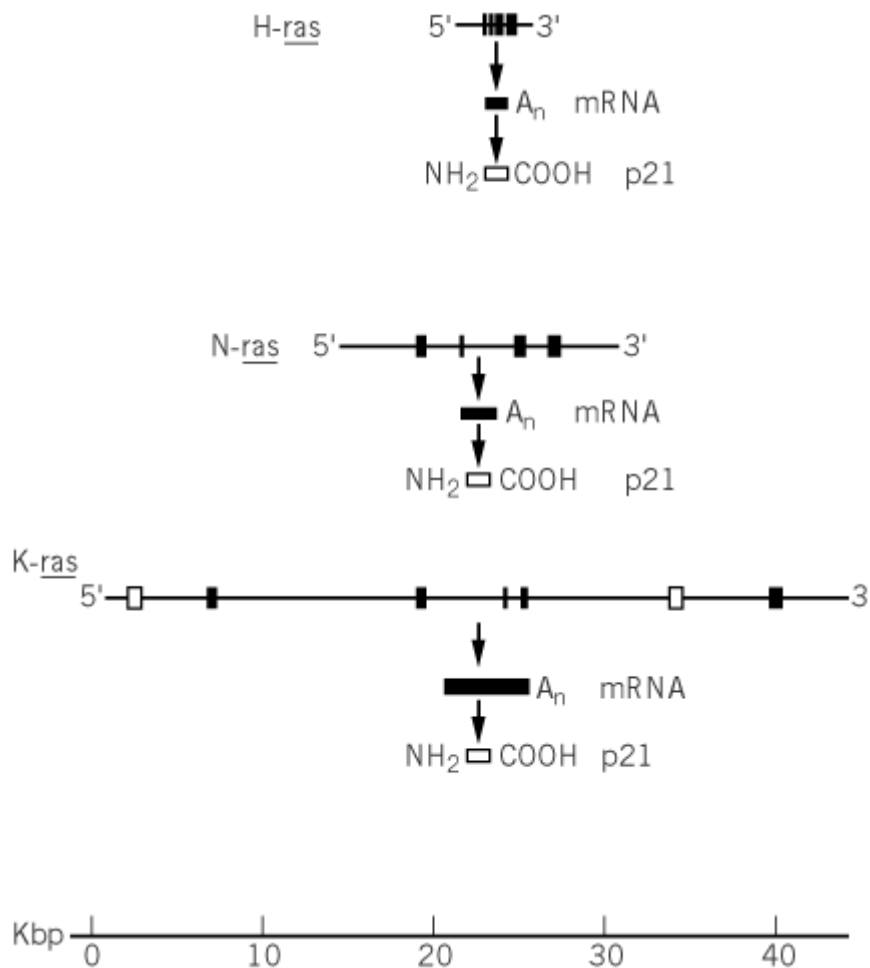
## Ras Oncogene

The *ras* gene is an important gene that has been transduced by three rodent [retroviruses](#) (Harvey sarcoma virus, Balb MSV, and Kirsten sarcoma virus). It has been studied extensively because of its role in **growth-factor-associated signal-transduction** pathways (1). The viral *ras* genes code for 21-kDa proteins, also known as p21, that belong to the G-protein superfamily. G-proteins bind to GTP with high affinity and hydrolyze it to GDP (see [GTP-Binding Proteins](#); [Gtpases](#); [Heterotrimeric G Proteins](#)). The GTP-bound Ras protein is biochemically active, whereas the GDP-bound form is inactive. More importantly, the proportion of GTP-bound p21 Ras is tightly regulated in normal cells by feedback mechanisms, and it normally represents less than 5% of the total Ras protein. Comparison of *v-* and *c-ras* revealed that the viral oncogenes contain two point mutations, one in codon 12 and a second in codon 59, both of which impair the intrinsic GTPase activity of the mutant proteins and render them resistant to negative regulation (1, 2). As a result, the *v-Ras* proteins remain in a constitutively activated (or GTP-bound) state, which contributes to their oncogenic activity.

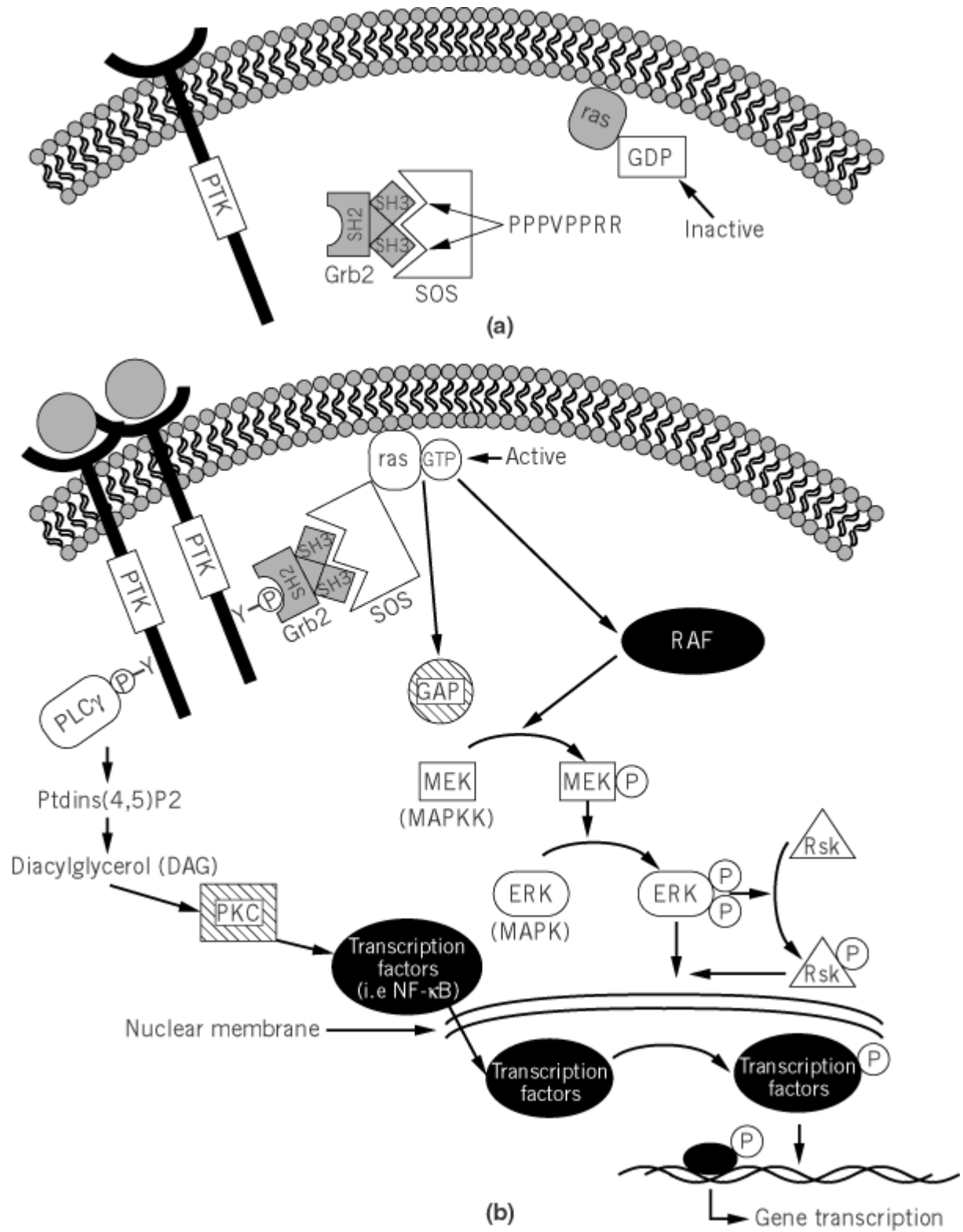
The human genome contains three closely related *ras* genes, all of which code for 21-kDa proteins, exhibit a high degree of [homology](#), and are known as H-*ras*, K-*ras*, and N-*ras*. (Fig. 1) (1). Of these, the H-*ras* gene is most similar to the *v-ras* genes derived from the Harvey sarcoma virus and BALB-MSV, whereas the K-*ras* gene is most similar to the oncogene encoded by the Kirsten sarcoma virus. The N-*ras* gene was originally isolated from a neuroblastoma following transfection of the tumor DNA into NIH/3T3 cells and analysis of the sequences responsible for inducing transformed foci. Approximately 20% of human tumors contain activated *ras* genes. A detailed analysis of the tumor-derived oncogenic *ras* genes showed that almost all of these genes contain

mutations either in the 12th or the 61st **codon** that result in constitutively activating the protein products. The ras proteins are membrane-bound [Gtpases](#) and play a critical role in transmitting growth signals transmitted by growth-factor receptors following their interaction with their ligands. These interactions lead to an increase in the cellular concentration of GTP-Ras, and the activated Ras in turn activates Raf and the Map kinase pathway (Fig. 2).

**Figure 1.** Comparison of the genomic structures of the *ras* proto-oncogenes. All of these genes encode four exons. *K-ras* encodes a fifth exon which is alternatively spliced during transcription. Although all of the encoded proteins are identical in size, the length of their transcripts and genetic loci differ considerably.



**Figure 2.** Schematic representation of receptor tyrosine kinase (RTK) signaling. (a) In a resting cell, RTKs exist as nondimerized, nonphosphorylated transmembrane proteins. The Grb-2 and Sos adaptor proteins are bound through the two Grb2 SH3 domains and the proline-rich regions of Sos. In addition, ras remains in an inactive, GDP-bound state. (b) Once the RTK has bound to its cognate ligand, it dimerizes and undergoes autophosphorylation on multiple tyrosine residues. Now, the phosphorylated receptor can bind to the Grb2:Sos adaptor complex. Because the function of adaptor proteins is to link multiple signal-transducing proteins within a single pathway, the Grb:Sos complex triggers the formation of GTP-bound, activated ras. Subsequently, activated ras initiates the ordered, sequential phosphorylation of RAF, MEK, ERK, and Rsk which is essential for cell growth and differentiation. The phosphorylation of several classes of transcription factors enables them to bind to DNA and induce gene transcription.



## Bibliography

1. M. Barbacid (1987) *Annu. Rev. Biochem.* **56**, 779–827.
2. G. Bollag and F. McCormick (1991) *Annu. Rev. Cell Biol.* **7**, 601–632.



## Reading Frame

**Nucleotide sequences** determine amino acid sequences by nucleotide triplet words named **codons**. Codons are arrayed on [messenger RNA](#) in a nonoverlapping, side-by-side manner. A nucleotide sequence can be read as codons in three different ways (or reading frames). During initiation of translation, the ribosome selects one of the reading frames by starting from the proper [initiation codon](#) and **translates** its reading frame until the stop codon. Maintenance of the reading frame is a highly accurate process, with an estimated error rate of  $5 \times 10^{-5}$  per codon. The reading frame maintenance relies primarily on three mechanisms: (i) unambiguous codon–[anticodon](#) triplet pairing between transfer RNA and mRNA; (ii) sufficient stability of codon–anticodon pairing throughout each elongation cycle; (iii) correct entry of aminoacyl-tRNA to a codon adjacent to that paired with peptidyl-tRNA. Theoretically any segment of nucleotide sequence can specify three different amino acid sequences. Some mRNAs, particularly those of viruses, actually contain overlapping protein coding sequences on different reading frames.

A change of reading frame during translation is called translational [frameshifting](#). Some translational frameshifting is genetically programmed and is required for expression of protein products. Similar reading frame changes can occur at transcriptional or post-transcriptional levels. [Frameshift mutations](#) are insertions or deletions of one or two nucleotides within the protein coding region, which cause reading frame changes. By frameshift mutations, the amino acid sequence corresponding to the 3' to the mutation site is completely disturbed, and the polypeptide product is often truncated. Some mutagens, such as [acridine dyes](#), induce frameshift mutations. Earlier genetic analysis of frameshift mutations and their intergenic [suppressor mutations](#) had led to the conclusion, even before the assignment of the meaning of the codon, that the [genetic code](#) consists of triplets.

A protein-coding region flanked by the [initiation codon](#) and the [stop codon](#) is, of course, devoid of any termination codons in the same reading frame. In contrast, termination codons tend to appear frequently in the other two reading frames of the protein coding region or regions that do not specify proteins. A reading frame following the initiation codon without interruption by any premature termination codons is called an “open reading frame.” A long open reading frame is likely to be a protein coding region. Hence, searching open reading frames serves as a useful way to predict protein coding regions from nucleotide sequences not yet characterized.

## Readthrough

During protein biosynthesis, readthrough in [translation](#) is a nonstandard decoding of the [stop codon](#) to an [amino acid](#), allowing the [ribosome](#) to continue translation beyond the termination codon (1). This phenomenon is also called [nonsense suppression](#). Readthrough is caused by either mutations in the translation machinery or by specific signals on the [messenger RNA](#). Genetically programmed readthrough in response to the mRNA signals is included in the recoding that represents reprogrammed genetic decoding [see [Frameshifting](#) (2)]. In a broader sense, translation beyond the stop codon by any mechanism, including frameshifting, is sometimes referred as readthrough. The term readthrough is also used for [transcription](#) beyond the transcript terminator.

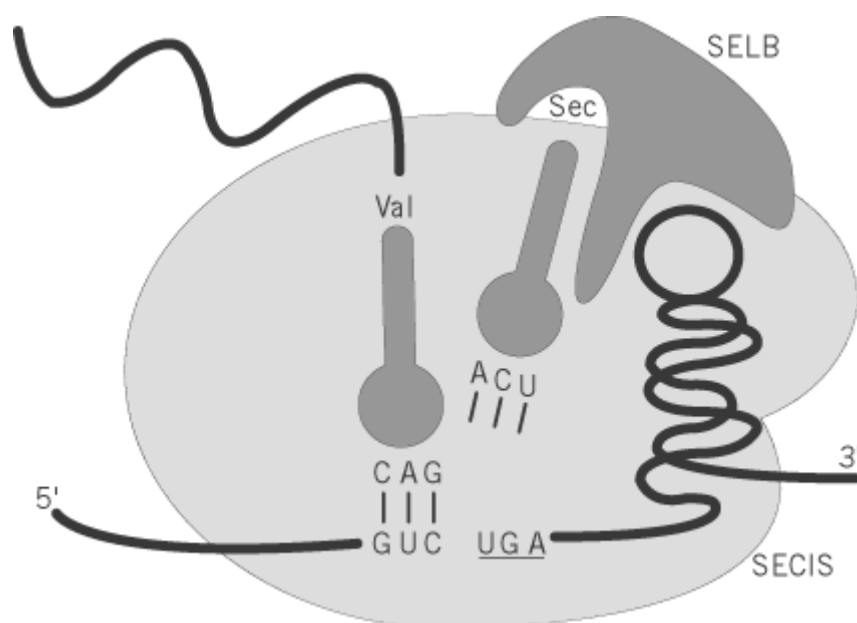
Nonsense suppressor tRNAs are mutant [transfer RNAs](#) that decode one or more termination codons,

usually due to a point mutation in the [anticodon](#). Some of the naturally occurring tRNAs show similar readthrough activity at the stop codon and are referred to as natural suppressor tRNAs. Certain mutations of ribosomes or [release factors](#) exhibit various nonsense suppressor activities in the absence of a suppressor tRNA. These mutations enhance the frequency of stop-codon misreading with a natural tRNA by affecting the normal termination process.

The programmed readthrough is further divided into two classes. In one class, canonical amino acids are inserted for the termination codons. At least two cellular genes utilize this type of readthrough. The gene for *E. coli* [pili](#) assembly is expressed by readthrough of a UAG codon. *Drosophila kelch* mRNA uses UGA readthrough. A number of virus genes have termination codons to be readthrough. A C-type [retrovirus](#), murine leukemia virus (MuLV), has a UAG codon that is decoded as a [glutamine](#) at residue an efficiency of 5% to express the *gag-pol* [polyprotein](#). An animal virus, sindbis, reads a terminating UGA as tryptophan at 10% efficiency. Various RNA plant viruses and bacteriophages also use UGA or UAG readthrough. These phenomena are mRNA- and codon-specific, involving specific signals on the mRNAs. In the case of MuLV, the UAG codon is accompanied by a downstream pseudoknot (see [RNA Structure](#)) structure that stimulates the readthrough.

The other category of programmed readthrough is insertion of a [selenocysteine](#) residue at certain UGA codons to synthesize selenoproteins (3). Selenoproteins exist in many organisms of all the three superkingdoms—archaea, eubacteria, and eukarya—except for some yeast species. In higher eukaryotes, nearly 10 selenoproteins have been identified. The UGA codons in the mRNA are decoded by a specific selenocysteinyl-tRNA encoded by a chromosomal gene. Selenocysteine tRNAs have a unique secondary structure. The tRNA is first aminoacylated with serine, which is subsequently converted to selenocysteine on the tRNA by specific enzymes. Stem loop structures of the mRNA, termed selenocysteine insertion sequence (SECIS), serves as a signal for the insertion. The SECIS is located just 3' to the UGA codon in bacterial selenoprotein mRNA, while it is on the 3' untranslated region in animal selenoprotein mRNAs. In bacteria, a special [elongation factor](#), SelB, instead of EF-Tu, delivers selenocysteyl-tRNA to the A site of the ribosome (Fig. 1). The precise mechanism of selenocysteine insertion in animals is not known.

**Figure 1.** Schematic representation of decoding of the UGA codon to selenocysteine (Sec) on the ribosome using the specific tRNA, elongation factor SELB, and the SECIS signal.



## Bibliography

1. E. J. Murgola (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society of Microbiology Press, Washington, D.C., pp. 491–509.
2. J. F. Atkins and R. F. Gesteland (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society of Microbiology Press, Washington, D.C., pp. 471–490.
3. C. Baron and A. Böck (1995) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society of Microbiology Press, Washington, D.C., pp. 529–544.

## Suggestion for Further Reading

4. R. F. Gesteland and J. F. Atkins (1996) *Annu. Rev. Biochem.* **65**, 741–768.

## Reassociation, Nucleic Acids

Single-stranded nucleic acid polymers that possess complementary base sequences can combine to form duplexes when mixed under appropriate solution conditions. For complementary short oligonucleotides, the reassociation process is fast and complete. For complementary polymeric nucleic acids, the annealing process is more complex, slower, and not necessarily complete. Monitoring of the rate of reassociation of the duplex is a classical method for evaluating the relative complexity of genomic DNA. The duplex DNA is fragmented by exposure to shear gradients or [restriction enzymes](#). Fragmentation is followed by thermal **denaturation** to produce single-stranded fragments. The reassociation process is monitored as a function of time. From these data, plots are constructed of the fractional extent of reassociation as a function of  $C_0t$ , where  $C_0$  is the initial nucleotide concentration and  $t$  is time. The value of  $C_0t$  at which the fractional extent of reassociation is 0.5, is a measure of the complexity of the genome (see [C0t Curve](#)).

Double-stranded DNA exhibits hypochromicity; the absorbance of the duplex is less than the sum of the absorbance of the constituent single strands. Changes in absorbance on reassociation to form duplex DNA can be used to estimate the extent of reassociation at any point in time.

Under conditions of low sodium phosphate concentration, both single-stranded and double-stranded DNA will bind to **hydroxyapatite**. At moderate sodium phosphate concentrations, 0.1–0.2 *M*, double-stranded DNA is retained on the column, but single-stranded DNA is eluted. Double-stranded DNA is eluted at 0.5 *M* sodium phosphate concentration. Because double-stranded and single-stranded DNA can be discriminated by their retention on the hydroxyapatite column, the time course for the reassociation process can be monitored.

## Suggestion for Further Reading

- V. A. Bloomfield, D. M. Crothers, and I. Tinoco (1974) *Physical Chemistry of Nucleic Acids*, Harper & Row, New York.

## Receptors, Hormonal

Hormone receptors act as the essential link between circulating [hormones](#) and molecular responses within target cells. Cell growth and function is regulated by a variety of extracellular signals, including not only hormones (in addition to the classical hormones, several so-called vitamins are now recognized to be hormones) but also [growth factors](#) and cytokines. Because there are very close similarities between the receptors for all three of these categories of effector, it seems sensible to include all three within this article. For example, the receptor for granulocyte macrophage [colony stimulating factor](#) (GM-CSF), a class I cytokine receptor, is very similar in both structure and function to the receptors for [growth hormone](#) and prolactin (1). These are all now grouped as members of the class I cytokine receptor superfamily.

This classification of GM-CSF receptor along with growth hormone (GH) receptor immediately implies that there are a range of receptor families. Firstly, although the majority of hormones, growth factors, and cytokines have their receptors in the plasma membrane of target cells, molecules such as steroids, [thyroid hormones](#), and vitamins A and D are all membrane-soluble, so they can cross the cell membrane down a concentration gradient. For this reason, their receptors are located within the cell. Thus, the first classification is by site: Receptors are either (a) within the plasma membrane or (b) within the soluble part of the cell.

The receptors found in the soluble part of the cell are all members of the same superfamily and are called the *steroid receptor superfamily*, which will be discussed in more detail later. The remaining receptors, perhaps best called transmembrane receptors, can be subclassified in a variety of ways. All have an extracellular **domain** (usually quite heavily glycosylated), a membrane-spanning domain, and an intracellular domain. The simplest is to classify them according to the nature of the membrane-spanning domain, but adding in some component of their cell signaling mechanism. This method of classification gives rise to four classes of transmembrane receptor. Firstly, there is the class of receptor in which there is *a single membrane-spanning chain* and the internal domain has *no intrinsic tyrosine kinase activity*. Examples of this are the growth hormone receptor, prolactin receptor, GM-CSF receptor, **interleukin-3** receptor, and so on. Secondly, there is the *single membrane-spanning chain that does have intrinsic tyrosine kinase activity*, and this is characterized by the *erbB* family (eg, **EGF receptor**, *erbB-2*, etc.), the [insulin](#) receptor, platelet-derived growth factor (PDGF) receptor, and so on. Thirdly, there is the *seven-helix membrane-spanning chain coupled to G-protein*. This represents the majority of hormone receptors and neurotransmitters—for example, receptors for the gonadotrophs (leutinizing hormone, follicle stimulating hormone), [glucagon](#), or adrenaline. For completion, there is the group of several (usually around 12) membrane-spanning subunits that permit the *gated* passage between extracellular and intracellular compartments for ion exchange.

It is not the purpose of this entry to consider *cell signaling*, and readers interested in the details of activation by receptors of **G-proteins**, [adenylate cyclase](#), **ras**, *raf*, mitogen-associated protein (**MAP**) **kinases**, **tyrosine kinases** [including Janus kinases (**Jaks**)], **signal transducers** and activators of [transcription](#) (Stats) should consult the appropriate entry. However, hormone receptors would be of little use if the binding of ligand to the receptor did not result in activation of one or more of these pathways. In the case of the growth hormone, prolactin, and GM-CSF family of receptors, there is clear evidence of overlap between the Jak-Stat and Ras/Raf/MAP kinase pathways (2).

Hormone receptors must be biologically selective. This means that any ligand that activates a specific receptor must be associated with the physiological responses associated with that receptor. For example, *estradiol-17 $\beta$*  is an 18-carbon steroid that induces estrogenic responses through the [estrogen receptor](#). *Estrogen-17 $\alpha$*  has virtually no estrogenic activity. In studies on the binding of

ligand to receptor, estradiol-17 $\beta$  binds to estrogen receptor with a dissociation constant of about  $10^{-10}$  M, whereas the 17 $\alpha$  isomer has a 100-fold lower affinity for the receptor, indicating that the hydroxyl group on carbon-17 plays an important role in the binding of ligand to receptor. However, diethylstilbestrol (DES), which is not a steroid and, on paper, bears little resemblance to the fused ring system of the steroids, has estrogenic activity at an equivalent dose to estradiol. Not surprisingly, in competition assays, estradiol and DES turn out to have very similar binding affinities for the estrogen receptor.

The class 1 cytokine receptor superfamily has already been mentioned, and it gives another interesting example of the importance of the nature of *ligand receptor interaction*. The three cytokines GM-CSF, IL-3, and IL-5 all have similar three-dimensional structures, in that they each contain a region with four  $\alpha$ -helices. The receptors for all three ligands are all heterodimers made up of  $\alpha$ - and  $\beta$ -subunits, but the  $\beta$ -subunit is common to all three. On examining ligand–receptor interaction (3), it was found that the  $\alpha$ -helix nearest the *N*-terminal end of the molecule interacts specifically with the common  $\beta$ -subunit of the receptor, whilst the  $\alpha$ -helix nearest the *C*-terminal end gives the biological specificity because it will only interact with the  $\alpha$ -subunit of the receptor that is specific for that ligand. The common  $\beta$ -subunit is the one that is involved in the cell signaling processes.

This use of a common subunit for cell signaling is not restricted to the class 1 cytokine receptor superfamily. For example, the receptors for luteinizing hormone (LH), follicle stimulating hormone, b-human chorionic gonadotrophin, and thyroid stimulating hormone are again all heterodimers with the  $\beta$ -subunit again being common—this time to act as the activator of the G-protein signaling system.

One of the interesting features of hormone receptors is that of selectivity in terms of target cell type. For example, a fat cell (adipocyte) will have receptors for LH, glucagon, and so on. If the cell is fully activated by LH, then that activity level of the target enzyme, hormone-sensitive [lipase](#), will be maximal ie, addition of glucagon will give no further activity. However, if the dose of LH used is below saturating, then addition of an appropriate dose of glucagon would push lipase activity up to the maximum. If the fat cells are replaced by Sertoli cells, however, there are no glucagon receptors and so glucagon will have no effect on such cells. Sertoli cells do have LH receptors but do not contain hormone-sensitive lipase. Instead, they are rich in the enzymes required for metabolism of cholesterol, and binding of LH to its receptors on Sertoli cells will result in an increase in synthesis and secretion of testosterone. Thus a common receptor can, nevertheless, give rise to different physiological end-points in different target cells.

The steroid receptor superfamily, which represents the receptors for those molecules which can diffuse across the plasma membrane, has been an important target for studies on regulation of gene expression. It has been recognized for many years that the effects induced by the majority of plasma membrane receptors are short-term (can be seen within seconds of ligand binding to receptor and are complete within a few hours, at most). Responses induced through the steroid receptor superfamily are generally long-term. The reason for this is that the hormones that act through transmembrane receptors are regulating the activity of preexisting enzymes, usually by increasing or decreasing the level of **phosphorylation** of specific [serine](#), [threonine](#) or [tyrosine](#) residues. The steroid receptor superfamily acts through modulation of transcription of specific genes, so that it is the amount of total protein (enzyme) that is changed, rather than the activity of a preexisting protein.

Historically, long before receptor activation of gene transcription was proven, it was recognized that steroids induce synthesis of fresh enzyme. For example, [glucocorticoid](#) induction of tyrosine aminotransferase synthesis in hepatocytes has been known for many years. All members of the steroid hormone receptor superfamily contain similar domains (4). In particular, they all contain a DNA-binding domain that is remarkably well conserved among different members of the superfamily (receptors for estrogen, androgen, progesterone, glucocorticoid, mineralocorticoid,

thyroid hormones, vitamin A derivatives, and vitamin D derivatives). They all contain a ligand-binding domain, and most (not glucocorticoid receptor) contain a nuclear localization sequence (see [Nuclear Import, Export](#)). This means that, as each receptor comes off its **polysome**, it is transported into the nucleus, and the empty receptor is located in the soluble part of the nucleus of target cells. Another common feature is that both the C- and N-termini of the molecule contain transactivation regions that are essential for the final activation of gene expression through interaction with various [transcription factors](#) (6).

Empty steroid receptor is found in association with various **heat-shock** proteins. In simplistic terms, these proteins are responsible for transporting the receptor to the nucleus and for preventing its binding to the DNA prior to the arrival of the hormone. For example, sites for binding a dimer of heat shock protein-90 (hsp90) are recognized in both the **DNA-binding** domain and the ligand-binding domain. Once ligand arrives and binds to the ligand-binding site, the receptor molecule undergoes some degree of **allosteric** change such that the hsp-90 molecule disengages and exposes not only the DNA-binding domain, but also a dimerization site. Thus two molecules of receptor, each with hormone attached, come together to form a dimer. This dimer now has high affinity for the specific sequence of nucleotides found in the appropriate [hormone response element](#), which is normally (but not always) upstream of the structural gene(s) known to be activated by that hormone within the particular target cell. In order to get full activation of the gene, the dimer must have both transactivation sites functional, and all appropriate transcription factors must be in place.

Functional hormone receptors are essential to life. Many endocrine diseases are now recognized to occur because of changes in the activity of different hormone (growth factor or cytokine) receptors. For example, cases of diabetes can be due to loss of, or fall in, activity of the insulin receptor. Assays of receptor function are now becoming very important in various areas of medical diagnosis and treatment. A variety of diseases of the digestive tract may be ascribed to over- or underactivity of the epidermal growth factor receptor; and controlled use of selective tyrosine kinase inhibitors that block the function of the EGF receptor, but do not significantly effect other tyrosine kinase-mediated responses, may be very useful in treating such clinical problems.

Treatment for breast cancer is now closely dependent on the detection of estrogen receptor in the primary disease. Patients whose tumors do not contain functional estrogen receptor will not respond to endocrine therapies and thus can be put immediately onto another type of therapy. Previously, it was often the case that all patients received initial additive endocrine therapy and were only switched to alternative therapies once they had failed to show any response to the endocrine therapy. Now, they should be getting a more appropriate therapy at an earlier stage of the disease, at which time it has more chance of success.

Because of the increasing clinical importance of hormone receptor assays, there is much more interest in the reliability of such assays. A number of external quality control assays have been set up. Biochemical assay of the estrogen receptor is very well established, and the results from very large studies are available to give a good guide to expected results (7). As many labs begin to switch from biochemical to immunohistochemical (IHC) assays, it is important to establish the same external QA for these. One such QA scheme has been set up for the IHC of estrogen receptors in breast cancer biopsies and has shown that remarkably good agreement can be reached by quite large numbers of participating labs (8).

## 1. Summary

Hormone receptors mediate the physiological response of specific target cells to external signals. Depending on the natures of both the ligand and the receptor, the response can be short- or long-term. The chemical nature of the ligand determines whether the receptor will be membrane-bound or present within the internal soluble components of the target cell. A target cell for any specific external effector is defined by the presence of that receptor on or within the cell.

Plasma membrane receptors may be subclassified according to the structure of the transmembrane region and the nature of the subsequent cell signaling mechanism. These receptors act, principally, by altering the activity of preexisting proteins through increasing or decreasing the amount of phosphorylation on specific serine, threonine, or tyrosine residues. Soluble receptors are normally found within the nucleus and, when activated by ligand, become dimerized in such a way that they acquire high affinity for specific nucleotide sequences (hormone response elements) associated with the appropriate structural genes.

Much current medical research is directed toward modulating the activities of these different types of receptor through either directly blocking their action or interfering with the consequences of their activation (selective kinase inhibitors, phosphatases, etc.).

### Bibliography

1. C. Bole-Feysot, V. Goffin, M. Edery, N. Binart, and P. A. Kelly (1998) *Endocrine Reviews* **19**, 225–268.
2. J. N. Ihle (1996) *Bioessays* **18**, 95–98.
3. A. B. Shanafelt, A. Miyajima, T. Kitamura, and R. Kastelein (1991) *EMBO J.* **10**, 4105–4109.
4. D. J. Mangelsdorf, C. Thummel, M. Beato, P. Herrlich, G. Schutz, K. Umesono, B. Blumberg, P. Kastner, M. Mark, P. Chambon, and R. E. Evans (1995) *Cell* **83**, 835–839.
5. R. White and M. G. Parker (1998) *Endocrine-Related Cancer* **5**, 1–14.
6. W. B. Pratt and D. O. Toft (1997) *Endocrine Reviews* **18**, 306–360.
7. R. Leake (1997) *Endocrine-Related Cancer* **4**, 289–296.
8. S. Romain, C. Laine Bidron, P. M. Martin, and H. Magdelenat (1995) *Eur. J. Cancer* **31A**, 411–417.
9. C. Jarvis, H. Patrick, L. Hopwood, and A. Dobson (1998) *J. Cell. Path.* **2**, 81–82.

### Suggestions for Further Reading

10. E. E. Baulieu and P. A. Kelly (1990) *Hormones*, Chapman and Hall, London.
11. C. Brook and N. Marshall (1996) *Essential Endocrinology*, Blackwell Science, Oxford.

## Receptors Linked To Tyrosine Kinases

Many cells are regulated by ligands that bind to cell-surface receptors with a single **transmembrane** domain, but without intrinsic [enzyme](#) activities in their cytoplasmic domains. While there are numerous examples of such receptors, these systems are best characterized by the **cytokine receptor superfamily**. Signaling by these receptors depends upon their interaction with cytoplasmic **tyrosine kinases**, which utilize a series of [protein–protein interactions](#) similar to those employed by the **receptor tyrosine kinases**.

### 1. Interactions of Cytokine Receptors

[Interferons](#), many of the **interleukins**, **growth hormone**, and numerous other molecules bind to the cytokine family of receptors (1). These receptors are characterized by (i) a single **transmembrane** domain with large extracellular regions and (ii) a shorter intracellular domain that does not contain any tyrosine kinase activity. In most cases, these receptors initiate signaling pathways by interacting directly with other signaling proteins, usually tyrosine kinases that are predominately cytoplasmic. In

fact, the receptors characteristically have discrete binding domains for these kinases, as well as a number of tyrosine-phosphorylation sites for interaction with **SH2**-containing proteins.

Although the activation of *src* family kinases by cytokines has been extensively studied, these enzymes are not activated by most cytokines. In contrast, cytokine binding to its receptor does lead to the activation of the JAK family of tyrosine kinases. They can couple cytokine receptors to a variety of downstream signaling pathways, but are specific for activation of a unique family of [transcription factors](#), known as STATs (see [JAK/STAT Signaling](#)).

#### Bibliography

1. J. F. Bazan (1990) Proc. Natl. Acad. Sci. USA **87**, 6934–6938.

### Recessive Lethal Mutations

Lethal mutations are [mutations](#) that inactivate essential genes. Mutants with dominant lethal mutations exist only if the lethality is conditional (see [Conditional Lethal Mutations](#)). In contrast, recessive lethal mutations, even if nonconditional, are propagated in a **heterozygote** because the essential function is supplied by the wild-type **allele** on the other [chromosome](#). The presence of the recessive lethal is confirmed upon mating by the absence of a class of progeny. Recessive lethal mutations are also propagated in **haploids** if they have more than one copy of the relevant gene. For example, an *Escherichia coli* cell carrying an extra copy of a chromosomal gene on a **plasmid**, a **bacteriophage**, or an **episome** is a partial diploid (called a [merodiploid](#)). A common use for such a strain is to study the regulation of an essential gene by inserting a [reporter gene](#) into one of the copies.

### Reciprocal Space

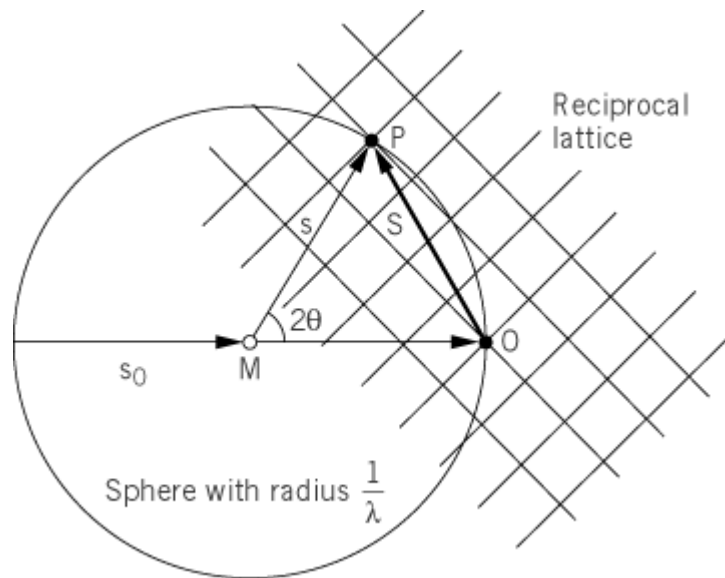
X-ray crystallography uses the concept of reciprocal space, an imaginary space that contains the reciprocal lattice. This lattice is also a nonreal lattice related to the crystal lattice. It is designed as an extremely useful tool to construct the directions of scattering by the crystal in combination with the Ewald sphere (Fig. [1](#)). The reciprocal lattice has a unit cell whose vector axes are **a\***, **b\*** and **c\***, and the collection of these unit cells forms reciprocal space:

1. **a\*** is perpendicular to the **b** and **c** axes of the real unit cell and has a length  $\frac{1}{d(100)}$ , where  $d(100)$  is the spacing between the planes in the set of planes (100).
2. **b\*** is perpendicular to the **a** and **c** axes of the real unit cell and has a length  $\frac{1}{d(010)}$ , where  $d(010)$  is the spacing between the planes in the set of planes (010).
3. **c\*** is perpendicular to the **a** and **b** axes of the real unit cell and has a length  $\frac{1}{d(001)}$ , where  $d(001)$  is the spacing between the planes in the set of planes (001).

Each reciprocal lattice point corresponds to a set of planes ( $hkl$ ) in the real lattice. Scattering of the incident X-ray beam by the crystal occurs if a reciprocal lattice point, for example, point P in Fig. [1](#), passes through the surface of the Ewald sphere.



**Figure 1.** The Ewald sphere as a tool to construct the direction of the scattered beam. The sphere has radius  $1/\lambda$ . The origin of the reciprocal lattice is at O.  $s_0$  indicates the direction of the incident beam and  $s$  of the scattered beam.



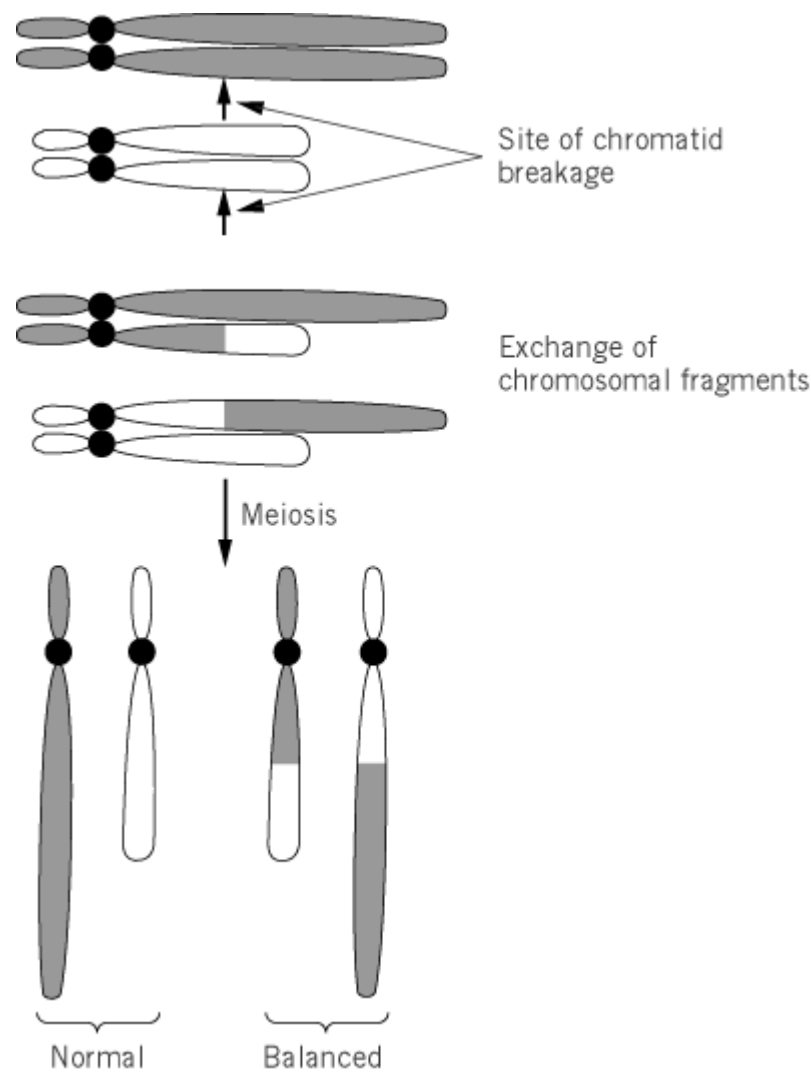
Suggestion for Further Reading

J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Reciprocal Translocation

Reciprocal **translocations** occur when breaks in two nonhomologous chromosomes are followed by the reciprocal exchange of chromosome fragments and repair of the broken ends (Fig. 1). Such reciprocal translocations involve no loss of chromosomal material. During **meiosis**, complete pairing of homologous chromosomal regions is accomplished by association of all four chromosomes, the two involved in the reciprocal translocation and the two in their original form. If **chiasmata** are formed at the end of each chromosome, a quadrivalent is formed (see **Multivalents**). The formation of fewer chiasmata may generate trivalents and univalents. The resulting problems in meiosis lead to abnormal **gametes** and cause **developmental** abnormalities. Many such defective fetuses spontaneously abort, but the frequency of human births in which balanced translocations have occurred is about 1 in 250.

**Figure 1.** Chromatid exchanges between nonhomologous chromosomes leading to balanced chromosomal translocations.



### Suggestion for Further Reading

M. S. Clark and W. J. Wall (1996) *Chromosomes. The Complex Code*, Chapman and Hall, London, UK.

### Recombinant DNA

Recombinant DNA is any **DNA** molecule that has been manipulated by *in vitro* procedures to create a novel sequence. The recombinant molecules produced may have a modified base sequence or contain sequences from two or more different genes or organisms. They are usually introduced into an organism to create a novel **protein** or **RNA** molecule that alters the properties of the organism.

Major recombinant DNA techniques include [site-directed mutagenesis \(1\)](#), [cloning](#), and polymerase chain reaction (**PCR**). Key enzymes required for recombinant DNA work are type-II [restriction enzymes](#), [DNA Ligase](#), **DNA polymerase**, **reverse transcriptase** and DNA **phosphatase**. Some of the uses of recombinant DNA techniques are to clone **genes**, to **knock out** a gene or genes in an

organism (2), to create **antisense** RNAs in an organism so as to interfere with **gene expression** (3), and to add a gene or genes to an organism to engineer a new or modified metabolic pathway (4).

### Bibliography

1. S. N. Ho, H. D. Hunt, R. M. Horton, J. K. Pullen, and L. R. Pease (1989) *Gene* **77**, 51–59.
2. C. G. Porter and M. J. Dallman (1997) *Transplantation* **64**, 1227–1235.
3. C. Robinson-Benion and J. T. Holt (1995) *Methods Enzymol.* **254**, 363–375.
4. F. A. Skraly, B.-L. Lytle, and D. C. Cameron (1998) *Appl. Environ. Microbiol.* **64**, 98–105.

### Suggestion for Further Reading

5. D. Watson, M. Gilman, J. Witkowski, and M. Zoller (1992) *Recombinant DNA*, 2nd edition, W. H. Freeman, New York.

## Recombinant Proteins

A recombinant protein is a [protein](#) encoded by a gene—on [recombinant DNA](#)—that has been **cloned** in a system that supports [transcription](#) of the gene into [messenger RNA](#) and [translation](#) of the resulting mRNA into protein (see [Expression Systems](#)). Recombinant DNA technology (also known as *gene cloning* or *molecular cloning*) comprises a number of techniques leading to the transfer of genetic information (deoxyribonucleic acid, or DNA) from one organism to another. The term *recombinant DNA* originates from the idea that insertion of DNA into a cloning vector is a special form of genetic [recombination](#). Genetic recombination is the formation of a new **genotype** by rearrangement of genes, resulting in gene combinations different from those present in the parents.

Several protocols exist for the creation of recombinant DNA. However, a typical cloning experiment could be carried out as outlined below:

1. The DNA from the donor organism is isolated, cleaved with [restriction enzymes](#), and **ligated** to a cloning **vector** that has been digested with the same enzymes as the donor DNA, causing the formation of a new, recombinant DNA molecule (a so-called DNA construct).
2. The DNA construct is transferred (**transformed**) into a host cell, where it can exist independently and be maintained as an extrachromosomal entity.
3. The host cells that take up the recombinant DNA can be identified and isolated from those that do not because of the presence of a selection marker on the cloning vector conferring, for example, [antibiotic resistance](#) to the host.
4. The DNA construct can be further manipulated to ensure that the cloned gene is expressed by the host, producing the recombinant protein product.

More recently, the polymerase chain reaction (**PCR**) has proved a valuable tool for the amplification of specific DNA fragments. The technique allows the amplification of tailor-made fragments equipped with restriction sites that facilitate subsequent cloning.

How can a gene encoding a specific protein of interest be identified so that it can be expressed as a recombinant protein? First, a [DNA library](#) representing the entire [genome](#) (or mRNA) of the donor organism is created by fragmentation of the DNA, followed by the cloning of each of the resulting fragments into a cloning **vector**. The resulting recombinant DNA molecules are transformed into a host. Now a short stretch of DNA with a sequence corresponding to part of the gene of interest, a probe, is created. This probe is used to “fish out” the library clones containing the matching DNA by

**hybridization** techniques, where the probe is allowed to find a complementary DNA sequence with which it **base-pairs**. The clone encoding the protein of interest can also be identified directly by the presence of that particular protein. For this strategy, an [expression library](#) is needed in which the cloned genes can be expressed, giving rise to the corresponding recombinant protein. The protein can now be detected, eg, by immunological techniques or by its activity. If the gene of interest is of **eukaryotic** origin, the expression library needs to contain **cDNA**, which is free of **introns** (the noncoding parts of a gene).

The design of an effective expression system is the key to obtaining a recombinant protein with the properties of interest. Basically, an expression system consists of the expression vector and a host for delivering the enzymatic machinery necessary for transcription and translation of the gene. The choice of host often depends on the intended application of the recombinant protein, as well as the origin of that protein. The expression vector is a complex tool, with a number of elements to adjust for optimal output. It allows the regulation of when, where, and in which form the protein should appear. It is furthermore possible to equip the recombinant protein with different types of “handles” by gene fusion. The protein of interest is then expressed as a **fusion protein**. The ideal fusion partner endows the recombinant protein with helpful features that facilitate its purification, detection, and folding (see [Fusion Gene](#), [Fusion Protein](#)).

Recombinant proteins have become of immense scientific and commercial value since the first cloning experiment, giving rise to a hybrid molecule that would express the foreign DNA within it, was successfully performed in 1973 (1). In 1978, the successful expression of human [insulin](#) in *Escherichia coli* was reported and, in 1982, the drug became commercially available (2). The immediate advantage of the technology was a more reliable source of the protein than the natural one. Furthermore, the protein could be produced in a more reproducible manner at a lower cost. The recombinant drug had an advantage over the commercially available porcine insulin in that it did not provoke any allergic reaction in the diabetics using it. This is an example of a first-generation therapeutic, a naturally occurring protein produced by recombinant DNA technology. The second generation of drugs are those being designed by use of [protein engineering](#). These proteins are redesigned versions of naturally existing proteins that give them improved properties in terms of stability, activity, and specificity. As an example, a second-generation tissue [plasminogen](#) activator (tPA) with a lower clearance rate was engineered by **site-specific mutagenic** replacement of three amino acid residues and introduction of a new N-glycosylation site. The new drug retained its catalytic properties (3). The third generation of therapeutic proteins will be produced not in a fermentor but by the patients themselves, thus starting the era of gene therapy.

Recombinant proteins have found use not only in health care but also for various commercial purposes in industry—not least in the production of foods and beverages. Also, the detergent industry makes profound use of recombinant proteins. Recombinant [lipases](#), [proteinases](#), and cellulases are important ingredients in washing powders. Many of these have been engineered to be able to withstand the harsh environment existing in a washing machine.

Recombinant proteins are obtained by using recombinant DNA technology. Only after the protein has passed the battery of tests needed before its release is approved can it finally be used commercially. The recombinant protein technology provides a protein source that, in almost all cases, exceeds the quantity and quality provided by the natural source.

#### Bibliography

1. S. N. Cohen, A. C. Y. Chang, H. W. Boyer, and R. B. Helling (1973) Proc. Natl. Acad. Sci. USA **70**, 3240–3244.
2. I. S. Johnson (1983) Science **219**, 632–637.
3. B. A. Keyt et al. (1994) Proc. Natl. Acad. Sci. USA **91**, 3670–3674.

#### Suggestions for Further Reading

4. Bernard R. Glick and Jack J. Pasternak (1998) *Molecular Biotechnology: Principles and Applications of Recombinant DNA*, American Society for Microbiology, Washington, D.C. A textbook for courses in biotechnology, covering both the underlying scientific principles and the wide-ranging industrial, agricultural, pharmaceutical, and biomedical applications of recombinant DNA technology. Strongly recommended.
5. Rocky S. Tuan (ed.) (1997) *Recombinant Gene Expression Protocols and Recombinant Protein Protocols: Detection and Isolation*, Methods in Molecular Biology, vols. **62** and **63**, Humana Press, Totowa, New Jersey. Provides balanced presentations of both background information and practical procedures.

## Recombinase

A major advance was made in the understanding of [gene rearrangement](#) of [immunoglobulin](#) (Ig) genes with the discovery [by Baltimore and coworkers (1)] of recombinases, encoded by the RAG1 (for recombinase activating genes) and RAG2. These genes are in opposite orientation on a 10-kbp piece of genomic DNA and were identified by screening a [cDNA library](#) derived by differential hybridization for its ability to promote an artificial substrate for V–J recombination after transfection into a fibroblast cell line. It was also realized that both RAG1 and RAG2 were necessary to ensure [recombination](#). This was also proven *in vivo* after generating of RAG1<sup>/</sup> or RAG2<sup>/</sup> targeted mice, because Ig gene rearrangement was abolished in both mutants, leading to a major severe primary immune deficiency, with [B-cell](#) differentiation blocked at the very early proB stage. Incidentally, [T-cell](#) development is similarly blocked, indicating that the recombinase is also required for [T-cell receptor](#) gene rearrangement.

It was not clear initially whether the RAG genes directly encoded the recombinases or whether they produced an activator of the enzymes; more recently, **gel retardation** experiments with appropriate synthetic substrates revealed that both RAG proteins clearly bind to the heptamer (CACAGTG) and nonamer (ACAAAAACC) recombination signal sequences (RSSs). Binding of recombinase to these sites requires that they be made accessible, as shown by the fact that, prior to gene rearrangement, low levels of germline transcripts of [V genes](#) may be identified. Also, in fibroblasts transfected with the synthetic substrate for recombination and the RAG cDNAs, endogenous Ig genes do not recombine, because this cell line does not contain the necessary enzymatic equipment that ensures DNA accessibility at the Ig locus in specialized B cells.

As shown by McBlane et al. (2), RAG proteins are necessary and sufficient to recognize and to cleave an RSS. The cleavage occurs in two steps. First, a nick is made at the 5' end of the signal heptamer, leaving a 5' phosphoryl group on the signal part and a 3' hydroxyl on the coding end. Second, the 3' OH is joined to the opposite phosphate on the other strand, forming a hairpin structure on the coding end and a blunt signal end. The two steps are strictly RAG-dependent. Evidently, cleavage and hairpin formation must take place on the two gene segments that are going to connect to each other (eg, a D to a J) and come into close contact. Both hairpin-like structures are then cleaved (this occurs at a random position), and partial modification of the coding end takes place, first by removal of some nucleotides and then by adding new nucleotides by terminal deoxynucleotidyl transferase (TdT), thus introducing **N-diversity**. The final step is the joining of the two coding ends, which requires several enzymes and factors, of which a major contributor is a DNA-dependent protein [kinase](#) (DNA-PK). Binding to DNA is devoted to the Ku subunits of DNA-PK, whereas the catalytic subunit (DNA-PKcs) is presumably activated after the complex is bound to

DNA. Simultaneously, the signal joints are also formed, making a piece of circular DNA that can be identified in preB cells. Mutations in the Ku DNA-binding subunits and/or DNA-PKcs of the catalytic subunit will result in severe blockage of lymphocyte differentiation, such as the *scid* (severe combined immune deficiency) mutation in the mouse, in which a regulatory subunit is inactivated, leading to the impossibility of making coding joints.

### Bibliography

1. M. A. Oettinger, D. G. Schatz, C. Gorka, and D. Baltimore (1990) RAG-1 and RAG-2, adjacent genes that synergistically activate V(D)J recombination. *Science* **248**, 1517–1523.
2. J. F. McBlane et al. (1995) Cleavage at a V(D)J recombination signal requires only RAG1 and RAG2 proteins and occurs in two steps. *Cell* **83**, 387–395.

### Suggestion for Further Reading

3. M. Gellert (1997) Recent advances in understanding V(D)J recombination. *Adv. Immunol.*, **64**, 39–64.

## Recombination

Genetic recombination encompasses a variety of processes that produce new linkage arrangements in the two interacting parental [chromosomes](#), or segments of a single chromosome. Several types of recombination can be distinguished according to the parental substrates. In **homologous**, or generalized, [recombination](#), the substrates share extensive **nucleotide sequence** identity, and recombination can occur anywhere within the **homologous** region. In site-specific recombination, the substrates share only limited nucleotide sequence homology (about 10 to 50 base pairs), and recombination occurs at a unique point within these sequences. In [transposition](#), a segment of DNA in one substrate is bounded by special sequences. This segment, or a copy of it, is inserted into random or nearly random sites in the second substrate. [Illegitimate recombination](#) includes events that do not fall into these categories. These events often occur between substrates that share only a few homologous base pairs and may reflect aberrant [DNA replication](#).

Recombination can also be classified according to the types of products produced. In homologous recombination, reciprocal recombination (crossing over) produces the two genetically reciprocal recombinant types in the same event. For example, *AB* and ++ parents produce *A+* and *+B* recombinants. Nonreciprocal recombination (gene conversion) produces one, but not the other, recombinant type. For example, *AB* and ++ parents produce an *A+* recombinant but not a *+B*. In site-specific recombination, if the two special sequences are in the opposite (indirect) orientation on one chromosome, recombination between the sites inverts the intervening segment. If the sites are in the same (direct) orientation, recombination deletes (or excises) the intervening segment. If the excised segment is circular, reversal of this reaction integrates one substrate into other. In transposition, the [transposable element](#) between the special sites may replicate and leave the transposable element in the substrate and introduce a copy of the element into the second substrate. Alternatively, the element may be excised from the first substrate and inserted into the second (cut-and-paste mechanism). In this case, the first substrate may be left broken, or the ends may be joined or otherwise repaired.

Recombination between DNA substrates is most thoroughly characterized, but recombination between RNA substrates, such as in RNA **viruses**, also occurs. RNA recombination may result from a copy-choice mechanism in which a replication complex “jumps” from one template to another at a

homologous point.

### Suggestions for Further Reading

D. E. Berg and M. M. Howe (eds.) (1989) *Mobile DNA*, American Society for Microbiology, Washington, D.C.

D. G. Catcheside (1977) *The Genetics of Recombination*, University Park Press, Baltimore, MD.

R. Kucherlapati and G. R. Smith (eds.) (1988) *Genetic Recombination*, American Society for Microbiology, Washington, D.C.

D. F. R. Leach (1996) *Genetic Recombination*, Blackwell Science, Oxford, England. This clearly written and illustrated book is an excellent introduction to all of the types of recombination discussed previously.

K. B. Low (ed.) (1988) *The Recombination of Genetic Material*, Academic Press, San Diego.

F. W. Stahl (1979) *Genetic Recombination: Thinking About it in Phage and Fungi*, Freeman, San Francisco, CA. This book gives mathematical treatments of recombination not readily found elsewhere.

H. L. K. Whitehouse (1982) *Genetic Recombination: Understanding the Mechanisms*, Wiley, New York.

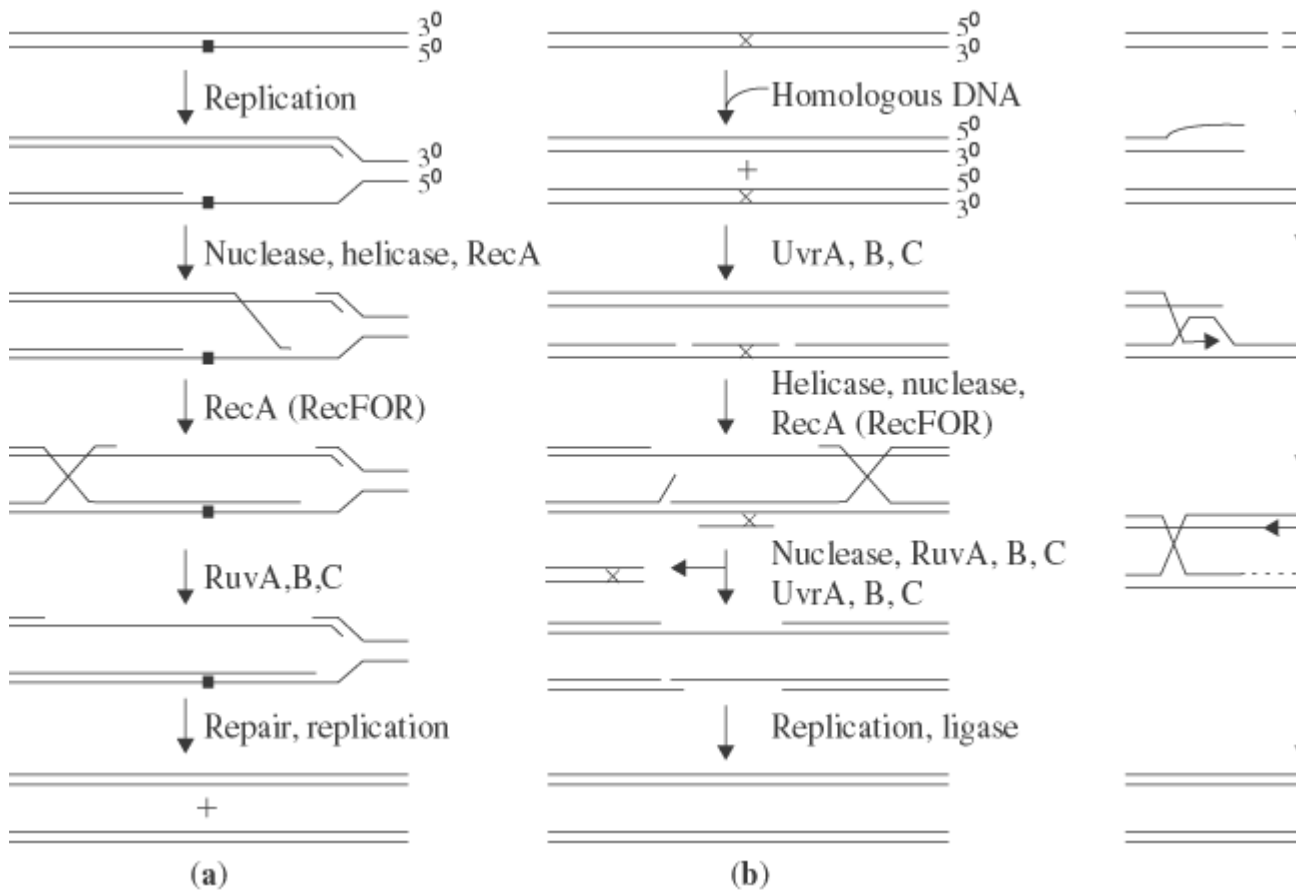
J. H. Wilson (ed.) (1985) *Genetic Recombination*, Benjamin Cummings, Menlo Park, CA. This volume is a collection of articles on mitotic, meiotic, and bacterial recombination from The Annual Reviews series from 1975–1983.

E. Wimmer, C. U. Hellen, and X. Cao (1993) Genetics of poliovirus, *Annu. Rev. Genet.* **27**, 353–436. This review contains a section on RNA recombination.

## Recombinational Repair

Most [DNA repair](#) mechanisms rely on the redundant information inherent to the structure of the DNA double helix to fix the damage: The base that is mismatched, modified, adducted, or fragmented is removed to generate a single-stranded gap, which is then filled in by **DNA polymerase**, using the complementary strand as template, and ligated. Lesions that affect both strands of the duplex cannot be repaired by this general reaction mechanism. Such lesions can be rectified only by retrieving the lost information from a homologous duplex in **diploid** organisms, or from the sister chromatid in [haploid](#) organisms such as *Escherichia coli*, which under exponential growth conditions is partially diploid for a major fraction of its [chromosome](#). Three types of lesion affect both strands and are repaired by recombinational repair: post-replication gaps across from a lesion, interstrand cross-links, and double-strand breaks (Fig. [1](#)).

**Figure 1.** Mechanisms of recombinational repair. **(a)** In post-replication repair a lesion in one strand leads to gap formation. The gap is invaded by the complementary strand from the sister duplex. Following further processing by nucleases and DNA junction is formed which is resolved by RuvABC resolvase. The remaining lesion in the duplex is then removed by excision repair. **(b)** ABC excinuclease makes dual incisions in one strand, and the cross-linked oligomer is displaced by the recombination reaction. This generates a Holliday structure and a “dangling” oligomer cross-linked to the duplex. This structure is recognized as a mobile DNA structure and is released by dual incisions. The Holliday structure is resolved and the gaps resulting from recombination are filled in by polymerases and ligated. **(c)** In double-strand break the RecBCD helicase/nuclease unwinds the duplex from both ends and generates a structure which can be processed by the RecA strand transfer activity. Further action of RecBCD and perhaps other nucleases generates a double-Holliday structure which is resolved by resolvases.



## 1. Post-replication Repair

When the *E. coli* replication machinery encounters certain nucleotide adducts, such as pyrimidine dimers, it stops replicating and reinitiates about 1000 base pairs beyond the adduct, generating a single-stranded gap that contains a damaged nucleotide (1). At the same time, a normal duplex is produced from the complementary strand. Thus, replication of a damaged duplex gives rise to one duplex with two normal strands and one partial duplex with a lesion in one strand and a gap in the other. The duplex with the defect is repaired by a process that involves both [recombination](#) and [excision repair](#). The **RecA** protein forms a helical filament at the post-replication gap and promotes homologous pairing with the intact sister duplex. This is followed by reciprocal strand exchange, so that the gap is “transferred” from the damaged duplex to the undamaged duplex, concomitant with the formation of a **Holliday intermediate** (2, 3). The latter is resolved by a **resolvase** encoded by the *ruvABC* or *rusA* genes. Filling in the gap by **DNA polymerase**, using the intact strand as template, yields two uninterrupted duplexes, one of which still contains a damaged base which can now be eliminated by a conventional excision repair reaction (4).

At present, there is no evidence that such a single-strand gap across damage is generated by the replication machinery of mammalian cells (5). Hence, post-replication repair has a different meaning in these cells, namely, the elimination of base lesions from DNA following replication of the damaged strand. In mammalian cells, the damaged strand can be converted into a duplex by “translesion synthesis” or by “template switching.” In translesion synthesis, the replication machinery simply synthesizes across the damaged base (6, 7), frequently by inserting the wrong base. In template switching, the [replication fork](#) stops at the lesion site on the damaged strand, but continues DNA synthesis on the undamaged strand. Then the newly synthesized strand is used as template for the strand of the sister duplex that had been blocked by the lesion (8). Once the synthesis (error-free) past the lesion is accomplished, the nascent strand complementary to the



damaged strand switches back to its parental strand. The end result of translesion replication and template switching is, again, the creation of two duplexes with no discontinuities. The lesion that remains following replication is eventually removed by excision repair. Even if the lesion is not removed by excision repair, however, the post-replication mechanisms outlined above can be repeated through many rounds of replication and consequently aid cell survival by ensuring the inheritance of uninterrupted duplexes to the daughter cells. In contrast to the bacterial system, however, the eukaryotic post-replication repair phenomenon remains ill-defined. For example, because of multiple **origins of replication** in eukaryotes, small post-replication gaps (7) may be generated by utilizing adjacent replication origins, and the resulting gaps may be processed as in prokaryotes.

## 2. Interstrand Cross-link Repair

Many carcinogenic or chemotherapeutic chemicals, including formaldehyde, nitrous acid, nitrogen and sulfur mustards, mitomycin C, cisplatin, and [psoralen](#) plus light, cross-link the two strands of the duplex covalently. Such cross-links prevent separation of the two strands and hence constitute absolute blocks for [transcription](#) and replication. Because of this effect on strand separation, cross-links constitute a more challenging lesion to cellular repair machinery than the post-replication gaps. Neither strand of the cross-linked DNA can be replicated to generate an intact duplex and a doubly damaged duplex that could recombine to produce two contiguous duplexes. Instead, the cross-link must be processed by the excision repair machinery. The mechanism of interstrand cross-link repair is relatively well-understood in *E. coli* (9). First, the (A)BC excinuclease incises the ninth phosphodiester bond 5' of the cross-linked base, as well as the third phosphodiester bond 3', in one strand only. Second, the resulting "excised" oligomer is displaced by RecA, which may form a filament at the gap and promote strand invasion of the homologous duplex. The recombination reaction generates a gapped duplex from the initially undamaged DNA and a triple-strand intermediate, in which the cross-linked oligomer (12-mer) is flipped out of the duplex. This latter structure is recognized as a form of monoadduct by the (A)BC excinuclease, which cuts out the dangling oligomer and produces a 12-nucleotide gap in one strand. The end-product of dual incision–recombination–dual incision reactions is to generate two duplexes, each with about a 12-nucleotide gap in one strand. These gaps are filled and ligated to complete repair.

The mechanism of interstrand cross-link repair in eukaryotes is not known. In yeast, both excision repair and homologous recombination systems are required to eliminate cross-links *in vivo*, suggesting a mechanism analogous to the prokaryotic excision repair system (10, 11). In mammalian cells, in contrast, mutants defective in the basal subunits of excision nuclease are not extremely sensitive to cross-linking agents, in contrast to mutants defective in homologous recombination, which are 30- to 90-fold more sensitive to mitomycin C than are wild-type cells (12, 13). It is possible that a homologous recombination mechanism, in which strand transfer is initiated from nicks produced by replication origins near the cross-links, leads to eventual elimination of cross-links. It has been found *in vitro*, however, that the excision-repair system acts very efficiently on psoralen cross-links in a very unusual manner. The enzyme system releases a 22- to 28-nucleotide oligomer from one strand on the immediate 5' side of the cross-linked base, by making dual incisions on the 5'-side of the cross-linked base (14). The resulting gap adjacent to the cross-link may act as a recombinogenic signal, leading to eventual removal of the cross-link by a mechanism rather like that of prokaryotes.

## 3. Double-Strand Break Repair

The basic mechanism of double-strand break repair in *E. coli* appears to be quite similar to the mechanism of cross-link repair. Again, RecA aided by RecBCD helicase/nuclease promotes strand invasion from the site of the double-strand break, using the homologous duplex as a target. The proposed reaction pathway generates two Holliday junctions, which are resolved by the RuvABC or *rusA* resolvase systems (15, 16). In mammalian cells, the double-strand break repair does not appear

to employ homologous recombination. Instead, the DNA-dependent protein kinase (DNA-PK) complex appears to play a crucial role in restoring the continuity of the broken duplex. The complex consists of a Ku70/Ku80 heterodimer, which has affinity for DNA ends, binds to the site of a double-strand break and recruits the DNA-PK catalytic subunit (DNA-PK<sub>cs</sub>), a 470-kDa protein with **serine/threonine kinase** activity to the site of the double-strand break.

The details of double-strand break repair in mammalian cells are not known. It is well-established, however, that joining of the coding sequence (CJ formation) during V(D)J recombination of the immune system (see [Immunoglobulin Biosynthesis](#)) and joining double-strand breaks caused by ionizing radiation employ the same enzyme system. As a consequence, mutants defective in Ku70, Ku80, or DNA-PK<sub>cs</sub> are defective in V(D)J recombination and are extremely sensitive to ionizing radiation.

### Bibliography

1. W. D. Rupp and P. Howard-Flanders (1968) *J. Mol. Biol.* **31**, 291–304.
2. W. D. Rupp, C. E. I. Wilde, D. L. Reno, and P. Howard-Flanders (1971) *J. Mol. Biol.* **61**, 25–44.
3. Z. Livneh and I. R. Lehman (1982) *Proc. Natl. Acad. Sci. USA* **79**, 3171–3175.
4. S. C. West, E. Cassuto, and P. Howard-Flanders (1981) *Nature* **294**, 659–662.
5. A. R. Lehman (1972) *Eur. J. Biochem.* **31**, 438–445.
6. R. Meneghini and P. Hanawalt (1976) *Biochim. Biophys. Acta* **425**, 428–437.
7. J. M. Clarkson and R. R. Hewitt (1976) *Biophys. J.* **16**, 1155–1164.
8. N. P. Higgins, K. Kato, and B. Strauss (1976) *J. Mol. Biol.* **101**, 417–425.
9. S. Cheng, B. Van Houten, H. Gamper, A. Sancar, and J. E. Hearst (1988) *J. Biol. Chem.* **263**, 11451–11460.
10. W. J. Jachymczyk, R. C. von Borstel, M. R. Mowat, and P. J. Hastings (1981) *Mol. Gen. Genet.* **182**, 196–205.
11. N. Magana-Schwencke, A. J. Henriques, R. Chanet, and E. Moustacchi (1982) *Proc. Natl. Acad. Sci. USA* **79**, 1722–1726.
12. C. A. Hoy, L. H. Thompson, C. L. Mooney, and E. P. Salazar (1985) *Cancer Res.* **45**, 1737–1743.
13. L. H. Thompson (1996) *Mutat. Res.* **363**, 77–88.
14. T. Bessho, D. Mu, and A. Sancar (1997) *Molec. Cell. Biol.* **17**, 6822–6830.
15. R. A. Holliday (1964) *Genet. Res.* **5**, 282–304.
16. J. W. Szostak, T. L. Orr-Weaver, R. J. Rothstein, and F. W. Stahl (1983) *Cell* **33**, 25–35.

### Suggestions for Further Reading

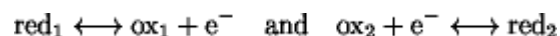
17. W. K. Kaufmann (1989) Pathways of human cell post-replication repair. *Carcinogenesis* **10**, 1–11.
18. H. Naegali (1994) Roadblocks and detours during DNA replication: mechanisms of mutagenesis in mammalian cells. *Bioessays* **16**, 557–564.
19. A. I. Roca and M. M. Cox (1997) RecA protein: structure, function, and role in recombinational DNA repair. *Prog. Nucleic Acids Res. Mol. Biol.* **56**, 129–223.
20. M. C. Whitby, G. J. Sharples, and R. G. Lloyd (1995) The RuvAB and RecG proteins of *Escherichia coli*. *Nucleic Acids Mol. Biol.* **9**, 66–83.
21. M. M. Cox (1997) Recombinational crossroads: eukaryotic enzymes and the limits of bacterial precedents. *Proc. Natl. Acad. Sci. USA* **94**, 11764–11766.

## Redox Potential

### 1. Redox Potential

#### 1.1. Redox reactions and redox couples

The molecule donating an electron is the reductant and becomes oxidized in the reaction, whereas the molecule accepting the electron is the oxidant and becomes reduced. In fact each redox exchange  $\text{red}_1 + \text{ox}_2 \longleftrightarrow \text{ox}_1 + \text{red}_2$  is composed of two reactions:



The reductant and its corresponding oxidized form constitute a redox couple, and two redox couples participate in the redox reaction. Each redox couple is characterized by a standard oxido-reduction potential ( $E_o'$ ), that quantitatively defines its tendency to lose an electron. In a redox reaction, the redox couple having more affinity for the electron will be the oxidant, whereas the couple having more tendency to donate electrons will be the reductant. The same redox couple may thus be reductant or oxidant depending on the redox potential of the other couple in the reaction. This is the way the electron transfer chains like the mitochondrial or the anaerobic bacterial chain and the photosynthesis in chloroplasts work: the reductant of the (*i*) step is the oxidant of the step (*i-1*).

#### 1.2. Electrode potential

The redox potential can be measured electrochemically as the electromotive force generated by connecting a half-cell of the redox sample with a reference half-cell. The reference half-cell consists of a Platinum electrode (immersed in a 1 M  $H^+$  solution and saturated with  $H_2$  gas at 1 atm); the sample half-cell consists of an electrode immersed in a solution of the redox couple under standard conditions (1 M of the electron donor and acceptor, 25°C and pH 7). The electrodes are connected to a voltmeter. A salt bridge (KCl solution) assures the electrical continuity between the half-cells. The electrons flow in the direction determined by the redox potential of the sample relatively to the reference, which is the Voltage observed when the experiment starts (standard concentration of 1 M) since the redox potential of the Hydrogen electrode is arbitrarily set at 0 Volt. The redox potentials ( $E_o'$ ) of relevant redox couples are thus expressed relatively to the  $H_2/2H^+$  couple.

For a redox reaction ( $\text{ox} + e^- \rightarrow \text{red}$ ), the Nerst equation relates the concentration of the redox species with the redox potential:

$$E = E^o + \frac{RT}{nF} \ln \frac{[\text{ox}]}{[\text{red}]}$$

where  $E$  is the redox potential,  $E^o$  is the redox potential for components in their standard state at pH 0,  $R$  is the gas constant ( $8.3144 \text{ J K}^{-1} \text{ mol}$ ),  $F$  is the faraday ( $9.65 \times 10^4 \text{ J V}^{-1} \text{ mol}$ ),  $n$  is the number of electrons transferred,  $T$  is the absolute temperature. In biology, the equation is generally written as

$$E_h = E_{m,x} + 0.06/n \log \frac{[\text{ox}]}{[\text{red}]}$$

where  $E_h$  is the redox potential referred to the standard hydrogen electrode,  $E_{m,x}$  is the mid-point redox potential (when  $[\text{ox}] = [\text{red}] \neq \text{standard state}$ ) at a defined pH (of  $x$ ) and  $2.303 RT/F = 0.06$  for 29.4°C. The value of  $E_{m,x}$  for a defined redox couple depends on the relative stability of the oxidized

and reduced state: a more negative the  $E_m$  results from a more stable oxidized form and the couple is a stronger electron donor. Conversely, any factor stabilising the reduced form makes the couple a better electron acceptor having a more positive redox potential.

### 1.3. Redox potential, free energy and equilibrium constant

The free energy change associated to a redox reaction ( $\text{red}_1 + \text{ox}_2 \leftrightarrow \text{ox}_1 + \text{red}_2$ ) is related to the difference in redox potential of the reactants by the formula

$$\Delta G^{0'} = -nF\Delta E^{0'}$$

in which  $n$  is the number of electrons transferred,  $F$  is the energy change as one mole of electrons falls through a potential of one volt ( $23.06 \text{ kcal V}^{-1}\text{mol}^{-1}$ ),  $\Delta E^{0'}$  is the difference of redox potential of the reactants in volts,  $\Delta G^{0'}$  is the free energy change per mole and is expressed in kilocalories. Since we know that at constant temperature and pressure a reversible reaction goes on till an equilibrium position defined by a concentration of reactants related by

$$[\text{ox}_1]/[\text{red}_1] \times [\text{red}_2]/[\text{ox}_2] = \text{Keq}$$

in which  $\text{Keq}$  is the equilibrium constant of the reaction at a certain constant temperature (K), and that  $\text{Keq}$  is related to the  $\Delta G^{0'}$  free energy change by

$$\Delta G^{0'} = -RT \ln \text{Keq}$$

where  $R$  is the gas constant ( $1.987 \text{ cal mol}^{-1} \text{ K}^{-1}$ ) and  $T$  is the absolute temperature (K).

We can also say that

$$-nF\Delta E^{0'} = -RT \ln \text{Keq} \quad \text{or} \quad \Delta E^{0'} = RT/nF \times \ln \text{Keq}$$

The last relation allows calculating the equilibrium concentrations of redox couples by knowing their redox potentials, or conversely the redox potential of a solution when the relative concentrations of the redox couples are known.

### 1.4. How to measure the redox potential

The oxidation of a reductant ( $\text{red}_1$ ) by an oxidant ( $\text{ox}_2$ ), at fixed temperature and pH, can be followed potentiometrically by measuring the variation of the electromotive force of a solution of the reductant along its titration with the oxidant. In practice with a Pt electrode connected with a reference one (i.e. calomel electrode), the redox potential of the  $\text{red}_1$  solution is recorded after each addition of  $\text{ox}_2$ . The redox potential increases as the  $\log [\text{ox}_1]/[\text{red}_1]$  increases and is equal to  $E_m$  when the concentration  $[\text{ox}_1] = [\text{red}_1] = 50\%$ . When  $[\text{ox}_1]$  is close to 100% the redox potential varies rapidly and the titration comes to the end: the point of equivalence has been reached, i.e. one equivalent of  $\text{ox}_2$  has been added to  $\text{red}_1$ . In biological systems, the redox proteins generally carry coloured groups having a definite absorption in the visible range of wavelengths, which allows a precise determination of the relative concentration of the oxidised and reduced species.

## 2. Redox Enzymes

Cell life depends on the availability of energy readily usable and storable for the many different specialised functions. This energy is extracted from nutrients by submitting them to a series of

chemical transformations associated to a negative variation of free energy. This transformation can be coupled to another chemical reaction with a positive variation of free energy, but smaller than the previous one. This second reaction is used to store energy. It consists generally of the formation of energy-rich bond between ADP and phosphate (Pi), leading to an ATP molecule. The covalent bond energy is available upon need through dissociation in ADP + Pi. Many of the physiological reactions in a cell, fundamental both in the metabolic pathways of degradation and in the storage of energy *via* production of ATP, involve the electron transfer from one molecule, the donor (or reducing agent), to another, the acceptor (or oxidant), and are thus called redox reactions. Many molecules can act as oxidising and reducing agents in non-biological redox reactions, whereas in the metabolic pathway of living organisms few molecules have been kept by the evolution as redox agents for many different substrates.

### 2.1. Where and how the electron transfers take place

In prokaryotes, the enzymatic systems metabolising nutrients associated to the complex synthesising ATP are localised in the cytoplasmic membrane. In eukaryotes, the respiratory chain and the oxidative phosphorylation are compartmented in mitochondria. In aerobic organisms, the biological energy comes from oxidation of substrates like glucose to the final state of CO<sub>2</sub> and H<sub>2</sub>O and terminal electron acceptor is O<sub>2</sub>. In order to allow the liberation of discrete amount of free energy and its optimal conversion to ATP or reducing effectors, like NADH or FADH<sub>2</sub>, the oxidation of glucose take place by steps. Three main metabolic pathways are involved in the whole process, i.e. glycolysis, tricarboxylic acid cycle and respiratory chain. Under anaerobic conditions, glucose is converted to pyruvate and NADH by glycolysis, and electron acceptors other than O<sub>2</sub> are used for the oxidation of NADH. Depending on the environmental conditions, some bacteria have developed respiratory chains making use of inorganic substrate like nitrate or sulphate instead of pyruvate; moreover, the final acceptor may be nitrate, Mn<sup>4+</sup> ions, Fe<sup>3+</sup> ions, sulphate or CO<sub>2</sub> instead of oxygen.

In the simpler case of Fe in the heme group or in Iron-Sulfur clusters and Cu in copper proteins, one electron is taken in going from the oxidised to reduced state. Conversely, the redox proteins carrying pyridine nucleotides, flavins or quinones as prosthetic groups take two electrons for total reduction. The two electrons can be transferred in a single step or in two successive steps. In the former case, the hydride (H<sup>-</sup>) transfer implies the formation of a substrate carbonium ion (C<sup>+</sup>), or alternatively the proton (H<sup>+</sup>) abstraction mechanism involves a substrate carbanion transition state (C<sup>-</sup>) followed by the transfer of 2 electrons. In the case of 2 successive electron transfers, flavins and quinones may also take one electron only and get to a stable semireduced state (semiquinone radical). Moreover, the gain of one or more electrons may be associated with the ionisation of a group in the protein and result in gaining one or more proton. In this last case protons take part in the reaction and the concentration of H<sup>+</sup> in the solution may influence the redox equilibrium.

Looking at the redox potential of the biological cofactors in the Table I and II, it appears evident that the couple NAD<sup>+</sup>/NADH is the one having the more negative potential and is thus the one with stronger tendency to transfer electrons and to become oxidised by other cofactors. A thermodynamically favoured sequence would be composed by NADH (-320 mV), flavin (-200 mV), ubiquinone (+100 mV), cyt c (+300 mV), cyt a<sub>3</sub> (+550 mV) and O<sub>2</sub>(+820 mV), which is the actual electron transfer pathway found in animal, plant and prokaryotic cells. Copper proteins having a potential in the range 200 to 500 mV, are sometimes present just before oxygen.

**Table I. Oxidoreduction potentials of some biologically relevant compounds and free redox groups**

| Redox Couple   | $E_{m,7}$ (mV) |
|--|----------------|
| $\text{NAD(P)}^+ + \text{H}^+ + 2\text{e}^- \rightarrow \text{NAD(P)H}$                                | -320           |
| $\text{S} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{H}_2\text{S}$                                  | -230           |
| $\text{Riboflavin} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Riboflavin-H}_2$                      | -200           |
| $\text{FMN} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{FMNH}_2$                                     | -219           |
| $\text{Pyruvate} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{lactate}$                               | -190           |
| $\text{Protoheme IX (Fe}^{3+}) + \text{e}^- \rightarrow (\text{Fe}^{2+})$                              | -115           |
| $2\text{H}^+ + 2\text{e}^- \rightarrow \text{H}_2$ *taken as reference                                 | 0.0 *          |
| $\text{Dehydroascorbate} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Ascorbate}$                     | +60            |
| $\text{Ubiquinone} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{Ubiquinone-H}_2$                      | +100           |
| $\text{MetBleu} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{MetBleu}$                                | +110           |
| $\text{Cu}^{2+} + 2\text{e}^- \rightarrow \text{Cu}^+$   | +150           |
| $\text{K}_3\text{Fe}^{3+}(\text{CN})_6 + \text{e}^- \rightarrow \text{K}_3\text{Fe}^{2+}(\text{CN})_6$ | +360           |
| $\text{NO}_3^- + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{NO}_2^- + \text{H}_2\text{O}$             | +480           |
| $\text{SO}_4^{2-} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{SO}_3^{2-} + \text{H}_2\text{O}$       | +480           |
| $\text{Fe}^{3+} + \text{e}^- \rightarrow \text{Fe}^{2+}$   | +770           |
| $1/2 \text{O}_2 + 2\text{H}^+ \rightarrow \text{H}_2\text{O}$  | +820           |

**Table II. Oxidation-Reduction Potentials of Some Types of Redox Proteins**

| Redox Couple   | $E_{m,7}$ (mV) |
|--|----------------|
| Fe in [Fe-S] cluster   | -570 -300      |
| Flavodoxins  | -500 -400      |
| Flavoproteines ox / red  | -400 +150      |
| Cytochrome b $\text{Fe}^{3+}/\text{Fe}^{2+}$   | +80 0          |
| $\text{PQQ} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{PQQH}_2$                     | +90            |
| $\text{TTQ} + 2\text{H}^+ + 2\text{e}^- \rightarrow \text{TTQH}_2$                     | +100           |
| Ubiquinone ox / red  | +100           |
| $\text{Cyt a Fe}^{3+} + \text{e}^- \rightarrow \text{Cyt a Fe}^{2+}$                   | +290           |
| Type I Copper proteins   | +350 +240      |
| Cytochromes c $\text{Fe}^{3+}/\text{Fe}^{2+}$  | +400 +100      |
| Ferredoxines ox / red  | +420           |
| $\text{Cyt a}_3 \text{Fe}^{3+} + \text{e}^- \rightarrow \text{Cyt a}_3 \text{Fe}^{2+}$ | +550           |
| Copper proteins $\text{Cu}^{2+} + 2\text{e}^- \rightarrow \text{Cu}^+$                 | +600 +300      |

## 2.2. Role of the apoprotein

The apoprotein is involved in at least four important issues: it stabilises the prosthetic group, it modulates the redox potential by binding the redox centre, it insures that the sequence of electron transfer is well followed, and that only specific partners of a given pathway will interact in an

ordered sequence.

The specific role of the apoprotein in the modulation of the redox potential has been the subject of extensive studies and is still matter of debates. In the field of hemoproteins, several properties of the protein have been proposed to influence the redox potential, like the polarity of the environment, the accessibility of the redox group to the solvent, the strength of the axial ligands or the charges on propionates. The role of electrostatic energy in this context has come out strongly in the recent years from the realisation that the determination of the redox properties cannot be done by considering a protein as a non-polar sphere, and that permanent dipoles must be taken into account (1). Another factor playing a role is the presence of charged and ionisable groups. A classical electrostatic analysis can reproduce the midpoints potentials of 4 hemes in *Rhodospseudomonas viridis* photosynthetic reaction centre (2). Moreover, a simple electrostatic hypothesis involving the positive charges on ferric hemes has been proposed recently (3). It allows in particular to explain both the difference in redox potential and the negative cooperativity in heme-heme redox interaction found in the artificially synthesised multi-heme 62 residues di- $\alpha$ -helical peptides.

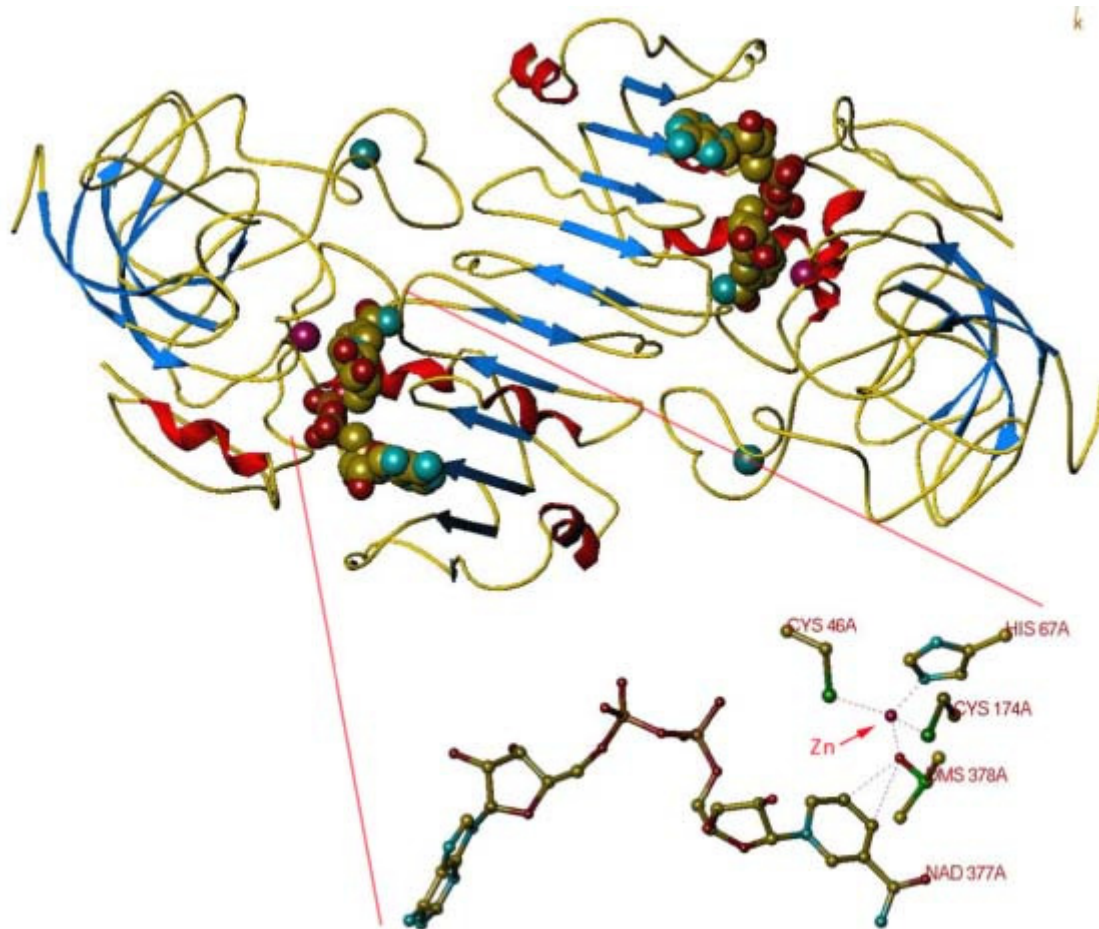
### 2.3. Groups transferring electrons

#### 2.3.1. Nicotinamide Coenzyme (4, 5)

The redox active centre of NAD (Figure 1) and the functionally equivalent NADP is the nicotinamide ring. These are the coenzymes of most dehydrogenases, which are generally believed to react by hydride transfer. Two electrons reduction and addition of one hydrogen produce a dihydropyridine cycle. Depending on the chemical reaction carried out, different classes of dehydrogenases have been defined (Glyceraldehyde-3-phosphate, Malate, Alcohol (Figure 1) [entry "ADH"], Glutamate, Isocitrate, Lactate [entry ""] Aldehyde dehydrogenase, Steroid reductase, Dihydrofolate reductase). In all of them the coenzyme participate to the reaction as a substrate: the reduced coenzyme dissociate from the apoprotein at the end of each catalytic cycle, thus the apoprotein does not significantly influence the redox potential of the coenzyme. In contrast to the flavin group, a fundamental character of this coenzyme is the stability of the reduced dihydropyridine towards the spontaneous oxidation in air. In spite of their different aminoacid sequence, the coenzyme binding site of all the NAD-dehydrogenases is very similar: a super-secondary structure composed of a parallel six-stranded  $\beta$  sheet flanked by a helices on both side of the sheet (Rossmann' fold). (4, 5)

**Figure 1.** The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows,  $\alpha$ -helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the "balls-and-sticks" representation, with the same color coding as above. The PDB site where the structures are available is "<http://www.pdb.bnl.gov/>".

Three dimensional structure representation of horse liver alcohol dehydrogenase dimer (PDB entry 6ADH). The NAD captures the reducing hydride from the Zn bound alcoholate, which, after loosing a proton becomes an aldehyde. The catalytic Zn is depicted as a magenta sphere, the structural Zn as a green sphere. Zn ligands are indicated in the close-up view.



### 2.3.2. Flavin Coenzyme (5)

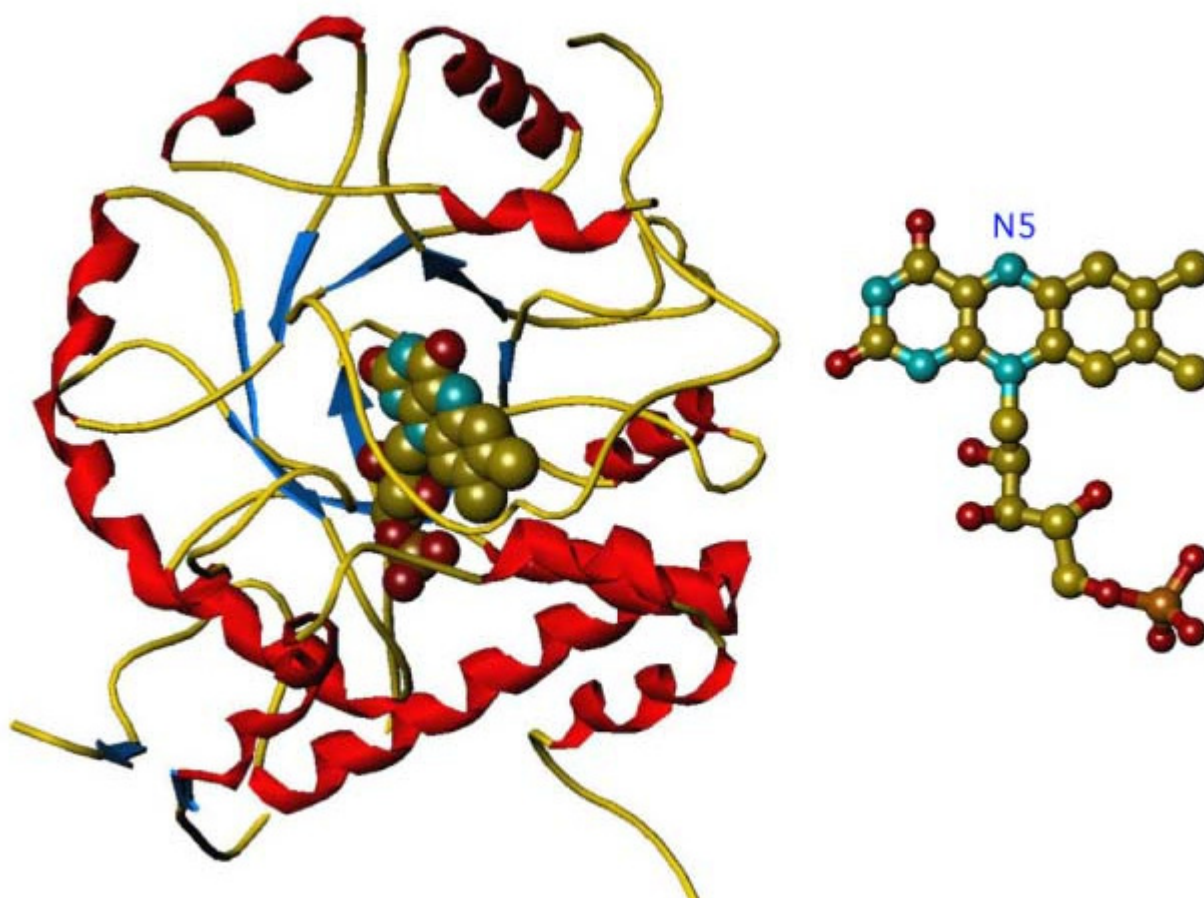
As in the case of nicotinamide coenzymes, FMN (Figure 2) and FAD are functionally equivalent. The electron-active part of the coenzyme is the highly conjugated isoalloxazine ring. The  $2e^-$  reduction results in the formation of 1,5-dihydroflavin, whereas the  $1-e^-$  reduction gives rise to a radical species, flavin semiquinone, which has different charge depending on the enzyme system. Unlike nicotinamide coenzyme, flavin does not freely dissociate from the apoprotein and has to be re-oxidised at the end of each catalytic cycle. The reoxidation of 1,5-dihydroflavin by  $O_2$  in solution is a relatively fast process ( $t_{1/2} < 1$  second). The strong binding of flavin prosthetic group make it possible to modulate the redox potential of the flavin by differently binding the 3 redox species of flavin to the apoprotein. The selective stabilisation of the oxidised form makes the redox potential of the coenzyme more negative than that in aqueous solution ( $-208$  mV); by stabilisation of the semiquinone the 1-electron potentials become more separated and the enzyme works essentially as a 1-electron transferring protein. The redox potential of flavoenzymes ranges between  $-465$  and  $+149$  mV; electron donors and acceptors compatible with these values determine the localisation of the flavoenzyme at specific places in electron transfer chains. Flavoenzymes work as oxidases [entry “Oxidoreductase”] (D-amino acid oxidase, AMP-sulphate reductase, Monoamine oxidase, Glycolate oxidase (Figure 2), Lactate oxidase, Hydroxyacid oxidase) or dehydrogenases (D-lactate dehydrogenase, Succinate dehydrogenase, p-OH-benzoate hydrogenase, Flavocytochrome  $b_2$ ) depending whether  $O_2$  is a good electron acceptor or other electron transfer proteins like quinone, cytochromes or iron-sulfur cluster are the acceptor; moreover, different types are known among oxidases (Oxidases, Oxidases-decarboxylases, Monooxygenases, Dioxygenases), which differ for the type of product(s) ( $H_2O$ ,  $H_2O_2$  etc.). Metalloflavoenzymes contain transition-metal ions like Fe or Mo and can work as oxidases or dehydrogenases. Flavodoxins are a class apart having a particularly



negative redox potential ( $-400$  to  $-500$  mV), functioning as  $1e^-$ -transfer proteins and having the dimethyl-benzenoid cycle of the isoalloxazine exposed to solvent. Most of the dehydrogenases transfer to  $1e^-$ -acceptor, which implies that the  $1e^-$ -reduced flavinsemiquinone should be a stable radical, at least under specific conditions. Actually, the peculiar role of flavoenzymes in biological oxidation is to be the relay between  $2e^-$  and  $1e^-$  electron transfer. (5)

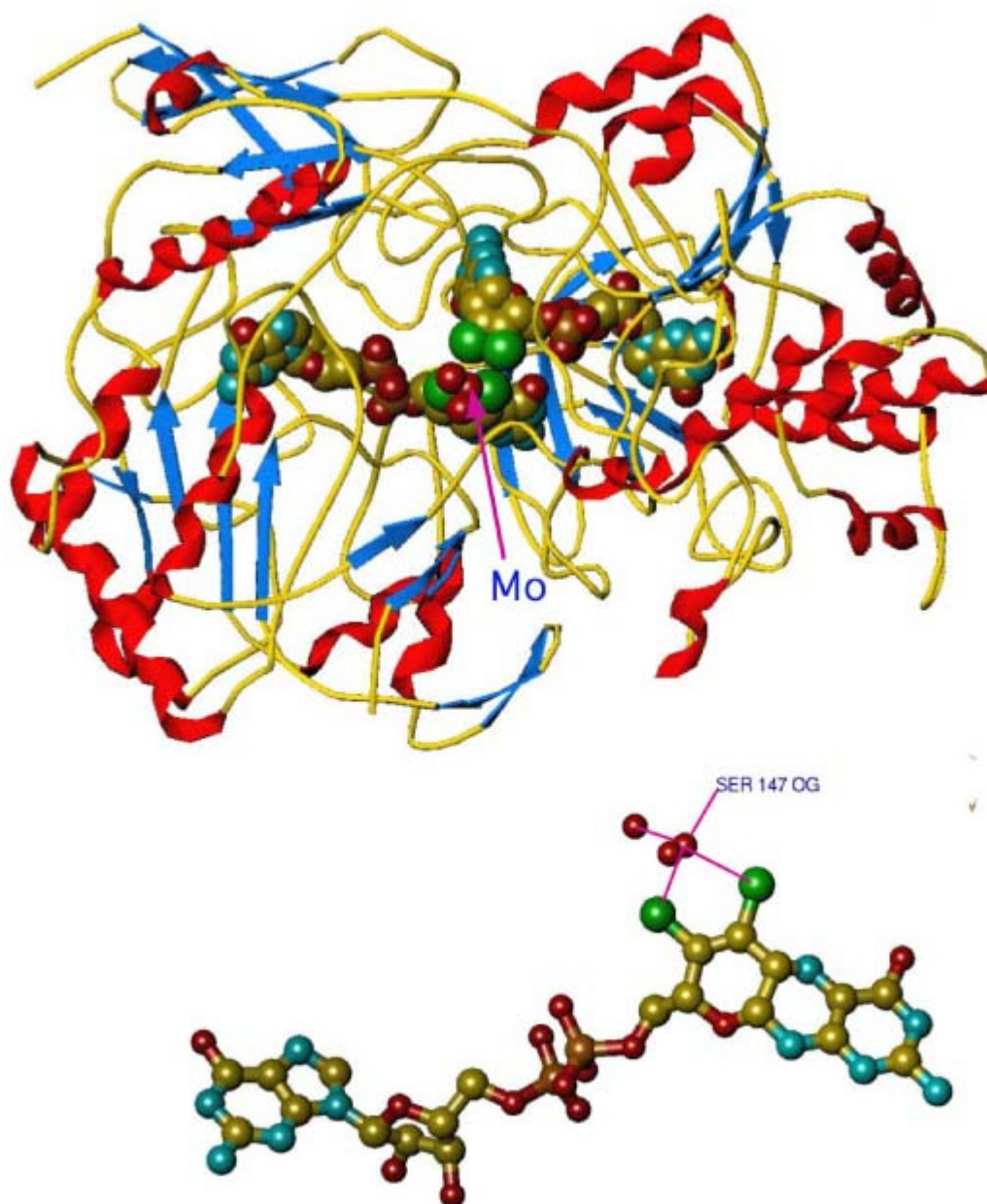
**Figure 2.** Three dimensional structure representation of spinach glycolate oxidase (PDB entry 1GOX). The folding motif is a  $\beta_8\alpha_8$  barrel, analogous to that of triose phosphate isomerase (“TIM barrel”). The FMN co-factor is located at the C-terminus extremity of the internal  $\beta$ -strands of the barrel. The nitrogen N5 is the electron acceptor.

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is “<http://www.pdb.bnl.gov/>”.



**Figure 3.** Three dimensional structure representation of the periplasmic DMSO reductase from *Rhodobacter capsulatus* (PDB entry 1DMS). The enzyme catalyzes the reduction of highly oxidized substrates like dimethyl sulfoxide to dimethyl sulfide. At the molybdenum redox center, two single electrons are transferred to the substrate dimethyl sulfoxide, generating dimethyl sulfide and water. The bis(molybdopterin guanine dinucleotide, MGD) has the molybdenum ion bound to the cis-dithiolene group of only one MGD molecule. Three additional ligands are bound to the Mo (see close-up view).

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is “<http://www.pdb.bnl.gov/>”.



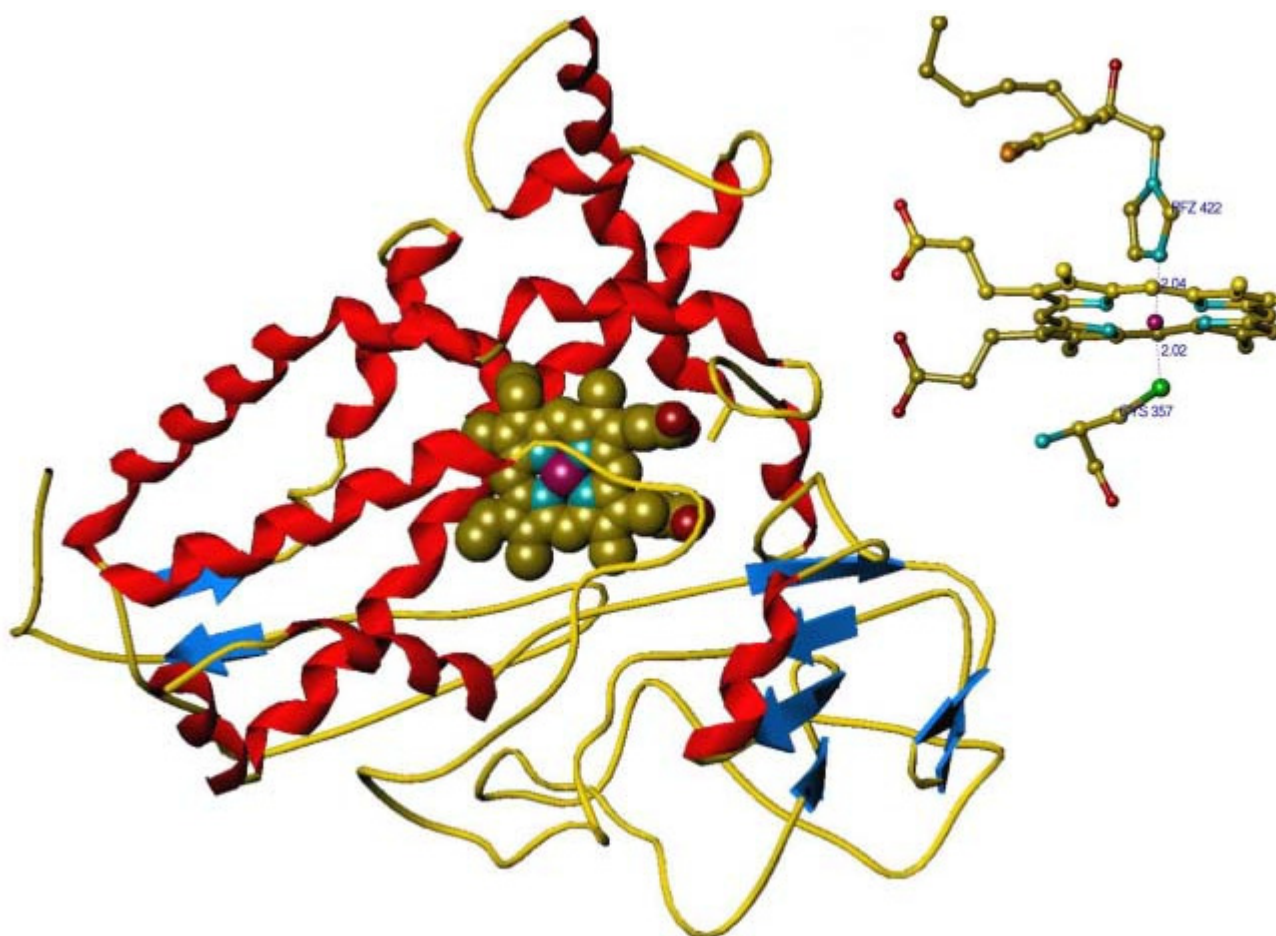
### 2.3.3. Metalloproteins (6)

[entries “Metalloproteins” “Metal Requiring Enzymes” “Iron Binding Proteins”]. Metals like molybdenum or tungsten are found associated to other redox centre in redox enzymes involved in the metabolism of nitrogen, sulfur and carbon. Generally the metal incorporates into the protein as Mo-cofactor (Mo-co) containing a Mo or W atom coordinated to a pterin cofactor, the structure of which has been determined (Figure 3). Two classes of these enzymes are known catalysing the transfer of an oxygen atom (DMSO reductase from *Rhodobacter* (Figure 3) and *E. coli*, mitochondrial Sulphite oxidase) or the hydroxylation of aldehydes or aromatic compounds (Xanthine oxidase and *Desulfovibrio gigas* Aldehyde oxidoreductases which also contain a flavin coenzyme, and *Pyrococcus furiosus* Aldehyde ferredoxin reductase). These reactions are associated with the Mo atom cycling between the oxidised state Mo(VI) to the reduced state Mo(IV). In DMSO reductase, the substrate takes  $2e^-$  from the Mo-co(IV) producing Mo(VI). In sulphite reductase, the  $2e^-$  reduction

of Mo-co from (VI) to (IV), is followed by a  $1e^-$  oxidation by a  $1e^-$  electron acceptor, which generate an intermediate Mo(V) state. It often occurs that the acceptor redox centre is contained in the same enzyme molecule as the Mo-co. (6)

**Figure 4.** Three dimensional structure representation of the P<sub>450</sub> cytochrome from *Pseudomonas putida* (PDB entry 1PHB). The *b* heme iron has a cysteine and the imidazolyl inhibitor as axial ligands. The C5 of camphor, a substrate, is oxidized in the alcoholic product.

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows,  $\alpha$ -helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is “<http://www.pdb.bnl.gov/>”.



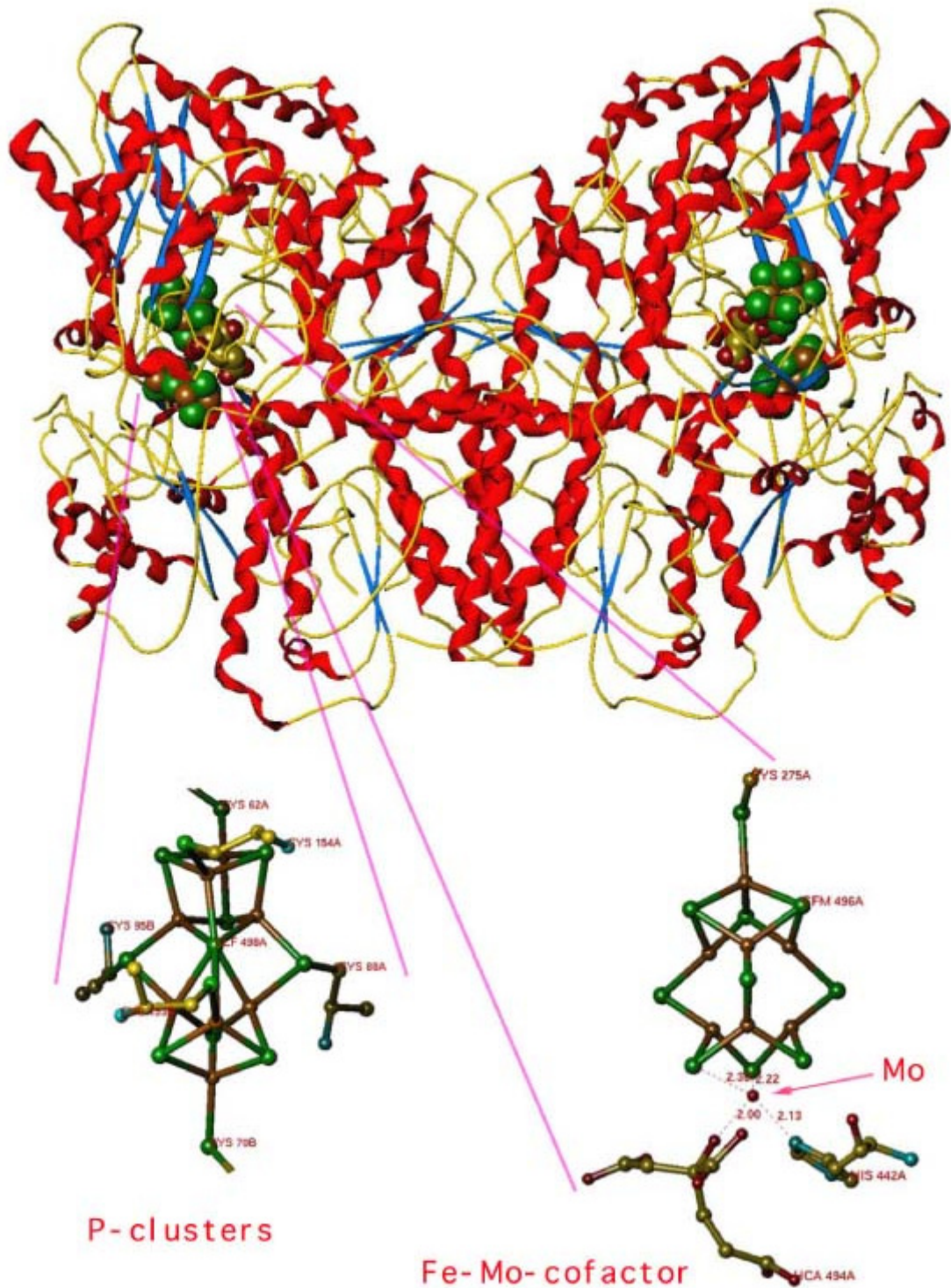
#### 2.3.4. Hemoproteins (5)

[entries “Iron Binding Proteins” “Hemoglobin” “Cytochromes” “O<sub>2</sub> Binding Proteins”. The prosthetic group generally bound to hemoproteins is protoporphyrin IX (Figure 4), in which the nitrogen atoms of the 4 pyrroles rings are the 4 equatorial ligands of the iron; other ligands from the apoprotein or not fill the 2 axial positions. The tetrapyrrole macrocycle has different peripheral side chains in cytochrome *b*, *a* and *d*; the *c* heme has the same peripheral side chains as the *b* heme but is covalently bound to the apoprotein by cysteine residues. The different environment in the apoprotein is strongly modulating the reduction potential of the 4 types of hemoproteins. Hemoproteins have different biological roles: binding and transport of O<sub>2</sub> (Hemoglobin, Myoglobin), electron transfer (Cytochromes), reduction of peroxides (Peroxidases, Catalases) and terminal component in

hydroxylation or desaturation (P<sub>450</sub> (Figure 4) [entry “Cytochrome P450”], Cytochrome b<sub>5</sub>). (5).

**Figure 5.** Three dimensional structure representation of the nitrogenase from *Azotobacter vinelandii* (PDB entry 3MIN). Nitrogenases perform the ATP-dependent reduction of N<sub>2</sub> to ammonia. Mo-Fe nitrogenase contains very unusual Fe clusters. The P-cluster contain two 4Fe-4 S clusters bridged by two cysteine thiol ligands and liganded by 6 cysteines SH groups. Among the six, two cysteines bridge the two clusters (bottom, left). The catalytic Fe-Mo cofactor (bottom right) has 4Fe-3 S and 1Mo-3Fe-3 S clusters bridged by 3 non protein ligands (not shown in the figure). The Mo is liganded by 3 sulfurs, an imidazole from His 442 and 2 oxygens of homocitrate (bottom right). The two clusters may be coupled in the sequence of N<sub>2</sub> to ammonia reduction.

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is [“http://www.pdb.bnl.gov/”](http://www.pdb.bnl.gov/).



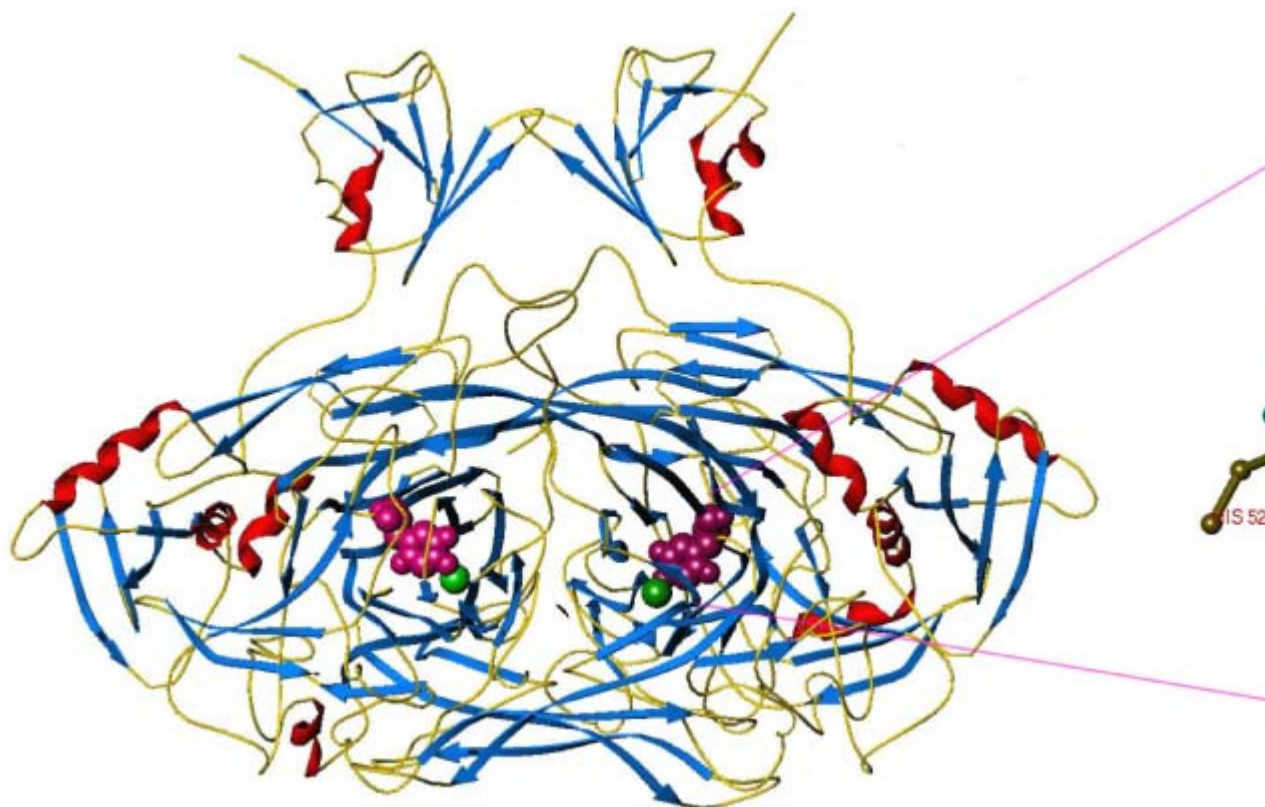
### 2.3.5. Iron-Sulfur cluster proteins (7)

[entries “Iron Sulfur Proteins” “Ferredoxin”]. Different classes of iron-sulfur containing proteins are defined on the basis of their iron content. The structure around the iron atom is believed to influence the redox potential (+770 mV of free  $\text{Fe}^{3+}/\text{Fe}^{2+}$ ). For example, Rubredoxin of anaerobic bacteria have 1 iron liganded by 4 sulfur atoms of cysteines, with a distorted tetrahedron geometry and its redox potential range between + 20 and – 60 mV. Ferredoxins contain 2 iron atoms and 2 inorganic sulfur atoms, 4 cysteine sulfurs are the other ligands of the irons and have a negative redox potential –240 to –460 mV. Bacterial ferredoxins have  $[4\text{Fe}_4\text{S}]$  clusters with a distorted cubic array, 4 cysteine sulfurs are the other ligands of the irons in each cluster. Their redox potential is between 0

and 645 mV. Rubredoxins and ferredoxins act only as electron transfer proteins, however, many redox enzymes contains iron-sulfur clusters in addition to other redox centres: among them we mention monooxygenases and dioxygenases (Methane, 4-methoxybenzoate monooxygenase, Benzene, Toluene, Naphtalene, Phtalate, Benzoate dioxygenase), Molybdopterin enzymes (Xanthine, Aldehyde oxidase, Purine hydroxylase, Carbon monoxide oxidoreductase), Siroheme enzymes (Sulphite, NADH-Nitrite reductase), Fe and Ni-hydrogenase, mixed-metal containing Nitrogenase (Figure 5), Iron-Sulfur non-redox enzymes like Aconitase, Fumarase and Dihydrohydroxyacid dehydratases. (7)

**Figure 6.** Three dimensional structure representation of the copper amine oxidase from *E.coli* (PDB entry 1OAC). Copper amine oxidases are quinonenzymes that catalyze the oxidative deamination of primary amines to aldehydes, with reduction of  $O_2$  to  $H_2O_2$ . The dimer contains two single copper ions and a covalently bound cofactor, TPQ (see text; figure left: TPQ magenta, copper green). The coordination of the copper ion consist of the TPQ cofactor and three histidines (close-up view: copper magenta).

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site were the structures are available is “<http://www.pdb.bnl.gov/>”.



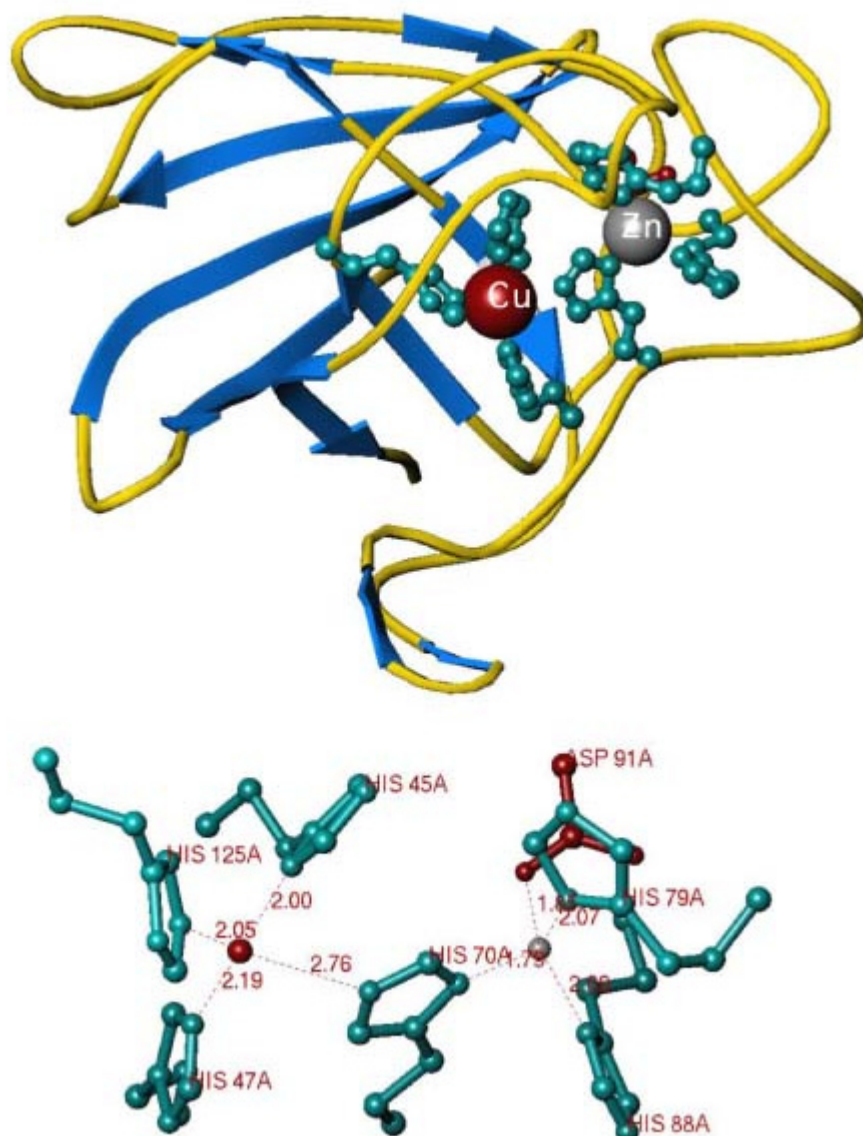
### 2.3.6. Quinones (8, 9)

It is now established that post-translational modifications of aromatic aminoacid side chains produce redox cofactors of quinoproteins. Tyrosine appears to be more frequently involved in this mechanism (Figure 6). The redox behaviour of quino-dehydrogenases is very similar to that of dehydrogenases using flavins or nicotinamide coenzyme, and based on the 3 redox forms experimentally found, hydride transfer as well as consecutive  $1e^-$  transfer steps seem possible. The group most frequently found in the quinoproteins known to date, is PQQ (2,7,9-tricarboxy-1H-pyrrolo[2,3-f]quinoline-4,5-

dione or methoxatin), that contains a pyrrole ring fused to a quinoline ring with an o-quinone in it. Due to the high redox potential of the PQQ/PQQH<sub>2</sub> couple (+ 90 mV), this group is able to oxidise most reduced flavin and nicotinamide coenzymes. Methanol DH, Galactose oxidase, Dopamine-b-hydroxylase and lipoxygenase have PQQ as redox cofactor. Another group of quinonoproteins contains TPQ (2,4,5-trihydroxy-phenylalanine-quinone) or TTQ (4-(2'-tryptophyl)-tryptophan-6,7-dione). Copper-containing amine oxidase carries a TPQ cofactor (Methylamine oxidase, Phenylethylamine oxidase, Diamine oxidase, and Amine oxidase). In copper amine oxidase (Figure 6), the copper atom was proposed to play a key role in the tyrosine to TPQ conversion reaction (9). Methylamine DH from *Thiobacillus versutus* and *Paracoccus denitrificans* contains TTQ as redox cofactor. The redox potential of the couple TTQ/TTQH<sub>2</sub> is high (+100 mV). (8, 9)

**Figure 7.** Three dimensional structure representation of the super oxide dismutase from *Photobacterium leiognathi* (PDB entry 1YAI). The eukaryotic super oxide dismutases are Cu, Zn enzymes catalyzing the dismutation of O<sub>2</sub><sup>-</sup>

The 3-dimensional structures are depicted using the following convention: β-strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is “<http://www.pdb.bnl.gov/>”.



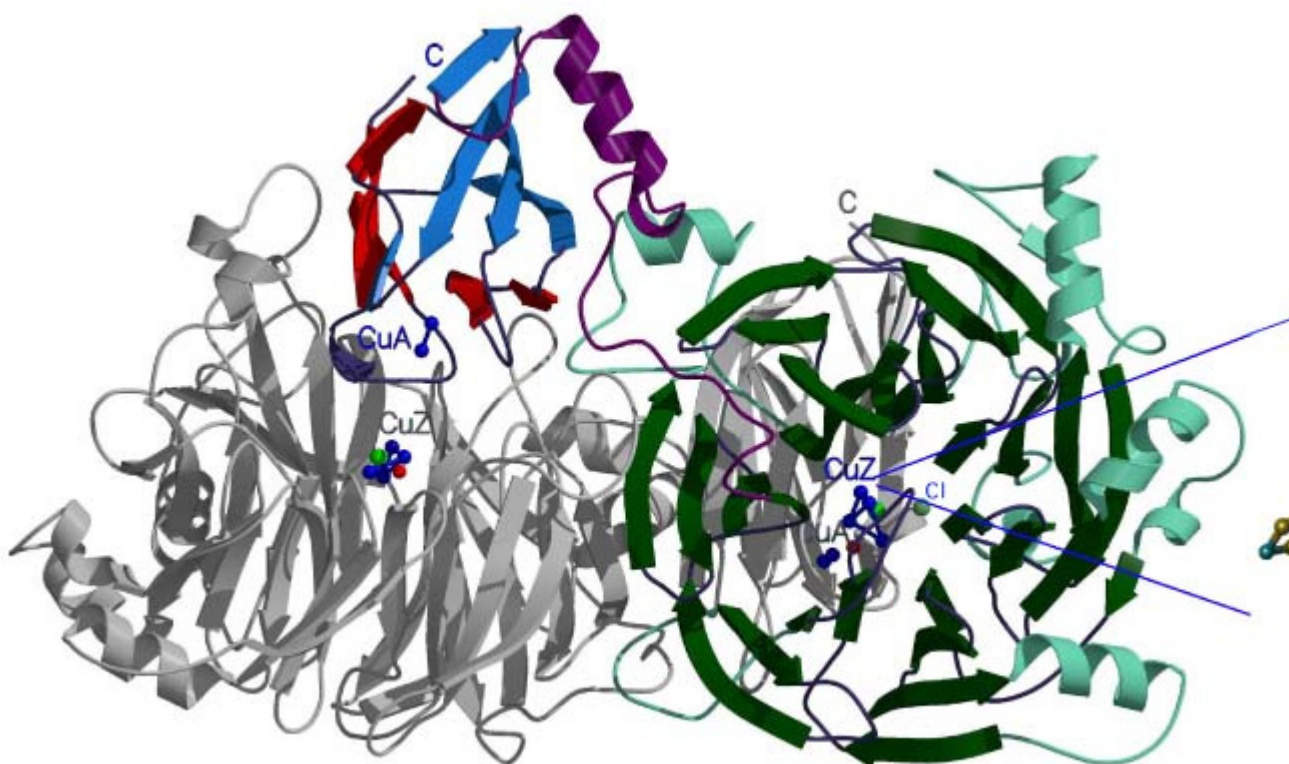
### 2.3.7. Copper containing enzymes (10-13)

Copper-containing redox proteins are characterised by very different coordination geometry of the copper atom. We can distinguish type I and type II copper proteins by their extinction coefficient at 600 nm, which is very intense for type I ( $0.7$  to  $10 \text{ mM}^{-1} \text{ cm}^{-1}$ ), and type III by its strong absorbance at 330 nm. Type I proteins are generally involved in electron transfer reactions and the redox change of the copper is associated to the bleaching of the blue absorption (Azurin, Plastocyanin, Pseudoazurin, Cucumber basic blue protein, Amicyanin). Two histidines, one cysteine and one methionine are coordinating the copper atom in a distorted tetrahedral arrangement. Their redox potential ranges between  $+240 \text{ mV}$  and  $+350 \text{ mV}$ . Type II proteins are involved in chemical reactivity and catalysis is not accompanied by redox or absorption change of the copper. The latter proteins reduce  $\text{O}_2$  either to  $\text{H}_2\text{O}$  (Super oxide dismutase, dopamine b hydroxylase) or to  $\text{H}_2\text{O}_2$  (Amine oxidase, Galactose oxidase). Super oxide dismutase (Figure 7) carries out the dismutation of two super oxide anions  $\text{O}_2^-$  in  $\text{O}_2$  and  $\text{H}_2\text{O}_2$ . Hemocyanine and Tyrosinase similarly coordinate copper of type III, with trigonal bipyramidal geometry. Multiple types of copper can also be present on the same protein (Ascorbate oxidase (type I, II, III), Ceruloplasmin (type I, II), Laccase, Cytochrome oxidase, Nitrite reductase (type I and II), Nitrous oxide reductase) (10). Nitrous Oxide Reductase (Figure 8), recently solved, contains in addition to a CuA binuclear centre, similar to that of cytochrome oxidase (11), a catalytic centre CuZ of a totally novel type: four copper atoms, held in place by seven histidines, are arranged in a distorted tetrahedron; an inorganic Sulfur atom, at the centre of the tetrahedron bridges the four copper atoms (12, 13).

**Figure 8.** Three dimensional structure representation of the nitrous oxide reductase from *P.nautica* (PDB entry 1QNI). Nitrous oxide reductase catalyses the last step of dissimilatory denitrification in anaerobic bacteria, namely the reduction of  $\text{N}_2\text{O}$  to  $\text{N}_2$ . The enzyme is a dimer of two identical subunit, each containing a CuA centre with 2 copper atoms, at which electrons enter into the system, and a catalytic CuZ centre, composed of four copper atoms (close-up view: CuZ with the seven coordinating histidines; copper mauve, sulfur green, oxygen red). One monomer is represented in grey. Secondary structure elements of the second monomer are dark and light green in the propeller domain, carrying the CuZ, red and blue in the cupredoxin domain, carrying the CuA centre, and purple in the linker region. In the dimer, the functional unit is formed by the CuA centre of one monomer and the CuZ of the other monomer.

The 3-dimensional structures are depicted using the following convention:  $\beta$ -strands are arrows, alpha-helices are spirals and turn or coil segments are ribbon. In the 3 dimensional structures, except otherwise mentioned, coenzymes or prosthetic groups are depicted as CPK models, with the color convention for atoms: oxygen red, carbon yellow, nitrogen blue, sulfur green, and phosphate orange. In the close-up figures, except otherwise mentioned, coenzymes or prosthetic groups are depicted in the “balls-and-sticks” representation, with the same color coding as above. The PDB site where the structures are available is “<http://www.pdb.bnl.gov/>”.





## References

“Redox Potential” in , Vol. 4, pp. 2112–2177, by Mariella Tegoni and Christian Cambillau, CNRS, AFMB UMR 6098, 31 Chemin Joseph Aiguier, 13402 Marseille, France; “*ara* Operon” in (online), posting date: January 15, 2002, by Mariella Tegoni and Christian Cambillau, CNRS, AFMB UMR 6098, 31 Chemin Joseph Aiguier, 13402 Marseille, France.

1. A. Warshel & A. Papazyan (1998) *Current Opinion in Structural Biology*, **8**, 211–217.
2. M.R. Gunner & B. Honig (1991) *Proc. Natl. Acad. Sci.*, **88**, 9151–9155.
3. D.E. Robertson, R.S. Farid, C.C. Moser, J.L. Urbauer, S.E. Mullholland, R. Pidikiti, J.D. Lear, A.J. Wand, W.F. DeGrado and P.L. Dutton (1994) *Nature*, **368**, 425–432.
4. M.G. Rossmann, A. Liljas, C. Bränden and L.J. Banaszak (1975) *The Enzymes*, vol **XI**, P.D. Boyer ed, 61–396, Academic Press.
5. C. Walsh (1979) *Enzymatic Reaction Mechanism*, Section III, W.H. Freeman and Co., San Francisco.
6. C. Kisker, H. Schindelin and D.C. Rees (1997) *Ann. Rev; Biochem.* **66**, 233–267.
7. R. Cammack (1992) *Adv. in Inorg. Chem.* **38**, 281–322.
8. J.A. Duine (1991) *Eur. J. Biochem.* **200**, 271–284.
9. M. Fontecave, H. Eklund (1995) *Structure*, **3**, 1127–1129.
10. E.T. Adman (1991) *Adv. in Prot. Chem.* **42**, 145–197.
11. S. Iwata, C. Ostermeier, B. Ludwig & H. Michel (1995) *Nature*, **376**, 660–669.
12. K. Brown, M. Tegoni, M. Prudencio, A.S. Pereira, S. Besson, J.J.G. Moura, I. Moura & C. Cambillau (2000) *Nature Strucural Biol.* **7**, 191–195.
13. K. Brown, K. Djinovic-Carugo, T. Haltia, I. Cabrito, M. Saraste, J.J.G. Moura, I. Moura, M. Tegoni, & C. Cambillau (2000) *J. of Biol. Chem.* **275**, 41133–41136.

## Suggestions for Further Reading

14. J. Garet Morris (1968) *A Biologist's Physical Chemistry, Chapter on Oxidation and Reduction*,

Edward Arnold Ltd., London.

15. C. Walsh (1979) *Enzymatic Reaction Mechanism*, Section III, W.H. Freeman and Co., San Francisco.
16. G.W. Pettigrew & G.R. Moore (1987) *Cytochromes c Biological Aspects*, Chapter "1", Springer-Verlag.
17. G.W. Pettigrew & G.R. Moore (1990) *Cytochromes c Evolutionary, Structural and Physicochemical Aspects*, Chapter "7", Springer-Verlag.
18. D.G. Nicholls & S.J. Ferguson (1992) *Bioenergetics 2*, Academic Press.
19. C. Bränden & J. Tooze (1991) *Introduction to Protein Structure*, Garland Publishing.

## Relaxation Spectrometry, Relaxation Time

Kinetic analysis often requires measuring very rapid reactions (see [Kinetics](#)), for which relaxation spectrometry is very useful. The temperature or pressure of a reaction mixture can be altered very rapidly, within 1  $\mu\text{s}$ , which perturbs an equilibrium between reactants and products by a small amount. Then the approach to the new equilibrium can be followed on the microsecond timescale, permitting the measurement of kinetic rate constants as great as  $10^6 \text{ s}^{-1}$ . The equilibrium studied can be a normal chemical reaction or it could be simply the association of two molecules to form a complex.

The two major advantages of relaxation spectrometry are that it operates on very short timescales and that the perturbation of the equilibrium is small in magnitude. For example, in a temperature-jump apparatus, the temperature of the solution is changed by about  $5^\circ$  in 1  $\mu\text{s}$ , which alters the equilibrium of a reaction by a small amount, depending on the enthalpy change of the reaction. The shift to the new equilibrium occurs and can be followed spectrophotometrically. Specific absorbance or fluorescent probes in the macromolecule can be used, or changes in protonation can be followed using pH-sensitive dyes. The advantage gained by perturbing the equilibrium by only a small amount simplifies the kinetic analysis. Each step in a reaction is observed as a single, first-order, kinetic reaction. The reciprocal of such a rate constant is known as the relaxation time.

Because relaxation spectrometry involves approaching equilibrium, both forward and reverse steps of each reaction are involved. As with other kinetic analysis involving approaching equilibrium (see [Kinetics](#)), the rate constant measured for each step is the sum of the forward and reverse rate constants for that step.

Another relaxation technique, which does not require perturbation of the system, is **nuclear magnetic resonance** (NMR). It also can be used to measure the rate constants for reactions and for the interaction between protein and ligand.

### Suggestions for Further Reading

- C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part III, *The Behavior of Biological Macromolecules*, W. H. Freeman, San Francisco, pp. 887–925.

## Release Factor

Termination of translation during protein biosynthesis takes place on [ribosomes](#) in response to a [stop codon](#), rather than a sense **codon**, in the “decoding” site (A site). Translation termination requires two classes of polypeptide release factors (RFs): a class I factor, codon-specific RF (RF-1 and RF-2 in prokaryotes; eRF-1 in eukaryotes), and a class II, nonspecific RF (RF-3 in prokaryotes; eRF-3 in eukaryotes) that binds guanine nucleotides and stimulates class-1 RF activity (Table 1). The fact that the two release factors from bacteria exhibit codon specificity suggests that they must interact directly with the codon. However, the underlying mechanism for translation termination represents a long-standing coding problem of considerable interest, because it entails protein-RNA recognition instead of the well-understood codon–anticodon pairing during the **messenger RNA**-transfer RNA interaction (1, 2).

**Table 1. Release Factors**

| Class | Factor  |   |
|-------|---|---|
|       | Prokaryotes   | Eukaryotes  |
| I     | RF-1 (UAG/UAA specific)                                   | eRF-1 (reads all three stop codons—omnipotent)                              |
| II    | RF-2 (UGA/UAA specific)                                   |   |
|       | RF-3 (stimulates RF-1 / 2 activity; recycle factor;GTPase | eRF-3 (stimulate eRF-1 activity; GTPase; prion-likeactivity of yeast eRF-3) |

### 1. Release Factor Genes

The genes encoding bacterial RF-1 (*prfA*) and RF-2 (*prfB*) have been identified. Mutations in these genes often cause misreading of stop signals, increased [frameshifting](#), or temperature-sensitive growth of the cells. The decrease in the cellular level and efficiency of RFs, therefore, leads to increased translational [readthrough](#) as the result of an abnormally long pausing of ribosomes at the stop signal. Unlike RF-1 and RF-2, RF-3 has received little attention since its initial characterization in the 1970s; its biological significance in protein synthesis has been a long-standing puzzle. After two decades of silence, the gene for RF-3, *prfC*, was discovered simultaneously by two groups (3, 4). Mutations in *prfC* caused misreading of all three stop codons, and an excess of RF-3 stimulates the formation of ribosomal termination complexes and increases RF-1 or RF-2 activity. The existence of a protein with RF activity in eukaryotes was demonstrated some 20 years ago in rabbit reticulocytes. After two decades of investigation, a eukaryotic protein family with the properties of an RF, designated eRF-1, was discovered (5). It includes a Sup45 protein of *Saccharomyces cerevisiae* that is involved in omnipotent suppression of three **nonsense codons** during translation. Eukaryotic counterparts of bacterial RF-3 were identified recently and are referred to as eRF-3 (6). The eRF-3 family includes a Sup35 protein of *S. cerevisiae* that carries G-domain motifs and is involved in

omnipotent suppression of nonsense codons.

## 2. Release Factor tRNA Mimicry

Upon accumulation of class I RF sequences from different organisms, conservation of protein motifs has been detected in prokaryotic and eukaryotic RFs, as well as in the C-terminal portion of elongation factor EF-G. EF-G is a translocase protein that forwards peptidyl tRNA from the A site to the P site on the ribosome (see [Elongation Factors \(EFs\)](#)). Because this C-terminal part of EF-G, domain IV, mimics the shape of the anticodon helix of tRNA (7), a provocative model has been proposed in which class I RFs mimic a tRNA shape for binding to the ribosomal A site and also mimic an anticodon for pairing to the stop codon (“RF-tRNA mimicry” hypothesis) (8). This view is supported by protein-RNA [cross-linking](#) and [footprinting](#) data that provide evidence for close contact between the stop codon and the release factor in the decoding site of the ribosome.

## 3. Class II Release Factors

Analogous to the initiation and elongation steps of translation, the termination step involves hydrolysis of GTP to GDP by RF-3 or eRF-3. The basic function of the G domain, like other [Gtpases](#), is to switch the protein conformation between two alternative states, a GTP-bound ON state and a GDP-bound OFF state. The model of RF-tRNA mimicry predicts that class II GTP/GDP-binding proteins, RF-3 and eRF-3, may be an EF-Tu-like vehicle protein to bring class I proteins to the A site of the ribosome or an EF-G-like translocase protein.

## 4. Prion-like Activity of Yeast eRF-3

The yeast eRF-3, Sup35, is a non-Mendelian **prion-like** element called [*PSI*<sup>+</sup>] that was uncovered some 30 years ago as a modifier of tRNA-mediated nonsense suppression in *S. cerevisiae* (9). Sup35 has several N-terminal tandem peptide repeats (PQGGYQQYN) similar to those in other prion proteins of mammals. Overexpression of the wild-type *SUP35* gene results in the *de novo* appearance of the corresponding [*PSI*<sup>+</sup>] phenotype, revealing that the [*PSI*<sup>+</sup>] factor is a self-modified protein analogous to mammalian prions. Moreover, **molecular chaperone** proteins, such as Hsp104, can cure cells of prions without affecting their viability (10). Therefore, Sup35 is likely to assume two functionally distinct conformations that differentially influence the efficiency of translation termination. In normal yeast strains, most Sup35 protein is soluble and functions in termination in a complex with its partner protein Sup45 (eRF-1). In [*PSI*<sup>+</sup>] strains, most Sup35 is insoluble, and this insolubility is inherited from generation to generation. It is tempting to speculate that this N-terminal “switching” domain of Sup35 is involved in some specific or global control such as [cell cycle](#) regulation. The biological significance of the prion-like property of Sup35 remains to be determined.

## Bibliography

1. W. P. Tate and C. M. Brown (1992) *Biochemistry* **31**, 2443–2450.
2. Y. Nakamura, K. Ito, and L. A. Isaksson (1996) *Cell* **87**, 147–150.
3. O. Mikuni, K. Ito, J. Moffat, K. Matsumura, K. McCaughan, T. Nobukuni, W. Tate, and Y. Nakamura (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5798–5802.
4. G. Grentzmann, D. Brechemier-Baey, V. Heurgue, L. Mora, and R. H. Buckingham (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5848–5852.
5. L. Frolova et al. (1994) *Nature* **372**, 701–703.
6. G. Zhouravleva et al. (1995) *EMBO J.* **14**, 4065–4072.
7. P. Nissen et al. (1995) *Science* **270**, 1464–1472.
8. K. Ito, K. Ebihara, M. Uno and Y. Nakamura (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5443–5448.
9. I. Stansfield and M. F. Tuite (1994) *Curr. Genet.* **25**, 385–395.

10. S. Lindquist (1997) *Cell* **89**, 495–498.

### Suggestions for Further Reading

11. Y. Nakamura and K. Ito (1998) *Genes Cells* **3**, 265–278.

12. R. H. Buckingham, G. Grentzmann, and L. Kisselev (1997) *Mol. Microbiol.* **24**, 449–456.

13. K. Ito, M. Uno, and Y. Nakamura (1998) *Proc. Natl. Acad. Sci. USA* **95**, 8165–8169.

## Renaturation

The term renaturation is used to describe a number of different processes in which the specific native three-dimensional structures of biological macromolecules are refolded *in vitro*. It can pertain to [proteins](#) or to **nucleic acids**. Both such types of molecules can be denatured, a process in which they lose their native structures (see [Denaturation, Nucleic Acids](#); [Denaturation, Protein](#)). Denaturation is generally accomplished by adding **denaturants** or by exposure to extremes of temperature, pH or pressure. The products are usually very unfolded, [random coil](#)-like individual polypeptide or polynucleotide chains, although less disordered chains, such as [molten globule](#) polypeptides, can result with only mildly denaturing conditions. Such unfolded chains often contain the information for their folded three-dimensional conformation in their linear sequence of amino acid residues or nucleotides and can regenerate it by a [self-assembly](#) process when the conditions are made less denaturing and more physiological. Renaturation involves refolding of the chains, which can involve only changes in conformation of the macromolecule. It can also include the addition of **prosthetic groups**, **coenzymes**, and other ligands, plus the introduction of [disulfide bonds](#) between **cysteine** residues in proteins.

With a protein, renaturation is described by [protein folding in vitro](#), starting with the unfolded, denatured state. The structure regenerated is the **secondary structure** and the [tertiary structure](#) of a single-**domain** protein that remains monomeric. [Quaternary structure](#) is also regenerated if the protein consisted of multiple domains or subunits.

Renaturation of a nucleic acid involves simply regenerating its folded conformation from a single polynucleotide chain, as with a [transfer RNA](#). The term is also used for the regeneration of the specific **double helix** from complementary individual single strands of DNA or RNA (see [Hybridization, Nucleic Acids](#)).

If renaturation does not occur, it is usually because the unfolded chains have precipitated or aggregated. Molecular **chaperones** prevent such aggregation with proteins by sequestering the unfolded molecule and permitting it to refold in the absence of aggregation or other nonproductive interactions (see [Protein Folding In Vivo](#)). Alternatively, renaturation might not occur spontaneously because the original molecule, which self-assembled spontaneously after its biosynthesis, has subsequently been modified covalently. Chemical modifications of proteins and nucleic acids prevent their self-assembly. Some covalent modifications are used *in vivo* to ensure that macromolecules do not renature spontaneously. For example, a peptide [hormone](#) like [insulin](#) is synthesized as a precursor, proinsulin, which folds and is then processed proteolytically to mature insulin. The single-chain proinsulin spontaneously folds and forms three disulfide bonds, and it can be renatured readily. Processed, two-chain insulin, on the other hand, does not readily self-assemble if it is unfolded and the three disulfide bonds are reduced. The principle of self-assembly can be expected to apply only to the original covalent forms of proteins and nucleic acids. Subsequent

covalent modifications can prevent renaturation, although not necessarily.

### Suggestions for Further Reading

R. Rudolph, G. Böhm, H. Lilie, and R. Jaenicke (1997) "Folding proteins, In" *Protein Function, A Practical Approach* (T. E. Creighton, ed.) IRL Press, Oxford, pp. 57–99.

## Reovirus

Reoviruses are icosahedrally symmetric, nonenveloped **viruses** that are 85 nm in diameter, whose morphology has been clarified by [cryoelectron microscopy](#) and [single-particle reconstruction](#). The name reovirus is an acronym for respiratory and enteric orphan virus, which means that these viruses can be isolated from human respiratory and enteric tracts, but are not associated with serious human diseases. However, reovirus causes a variety of diseases in experimental animals, mice in particular, and this virus has been used extensively as a tool for studying viral pathogenesis in a number of organs, such as the brain, heart, intestines, and liver.

The [genomes](#) of reoviruses are characterized by 10 discrete segments of double-stranded RNA (dsRNA) which can be well-separated by polyacrylamide [gel electrophoresis](#): three large (L1–L3), three medium (M1–M3), and four small (S1–S4). The reovirus genome is enclosed by eight structural proteins arranged in two concentric capsids. All of the **genes** from a type 3 strain encompassing 23,548 base pairs in length have been **cloned** and sequenced. All 10 segments contain short noncoding regions; the 5'- and 3'-untranslated regions range between 12 and 32 bases and between 35 and 83 bases, respectively. Sequences at the 5' and 3' ends of the 10 plus-strands are highly conserved, and the terminal sequences 5' GCUC---- and ----UCAUC 3' are found in all strains, while segment-specific 5'- and 3'-subterminal sequences are also present. The untranslated regions are likely to include sequences for [messenger RNA](#) packaging, recognition by the viral **RNA polymerase** for initiating plus- and minus-strand synthesis, and determining translational efficiency. The 5' end of each mRNA has a eukaryotic [-cap](#) structure, <sup>m7</sup>GpppGmpC, but the 3' end has no [poly \(A\)](#) tract. Each of the segments is monocistronic, except for the S1 gene, which encodes s1 and s1s derived from alternative [reading frames](#).

Replication of reovirus occurs entirely in the cytoplasm. Sialic acid residues on **glycoproteins** of the cell surface are the receptors for the outer capsid protein s1, which forms a trimer with head-and-tail morphology projecting from each of the vertices of the virion. Additional host cell factors also appear to ensure reovirus infection. Concerning the enhancement of infection efficiency conferred by [epidermal growth factor](#) receptor and v-erbB, the virus has been suggested to make opportunistic use of an already-activated [signal transduction](#) pathway. After virus entry by [endocytosis](#), virions are subject to proteolytic processing of s3 and m1/m1C, and the resultant particle is called the infectious subviral particle (ISVP). ISVP then converts to the core, transcriptionally activated particle. Viral core is a multienzyme complex consisting of l3, l2, l1, m2, and s2 proteins. l3 serves as the RNA-dependent **RNA polymerase**, and l2 provides the guanylyltransferase and methyltransferase activities for the capping process. [RNA helicase](#) and nucleotide triphosphatase activities reside on l1 and m2.

[Transcription](#) occurs in a conserved and asymmetric fashion in the core, and the mRNA is released from the core through the viral channel surrounded by the l2 pentamer. Each mRNA is transcribed at a rate proportional to the reciprocal of its segment length, providing viral proteins, and also acts as a template to make double-stranded RNA. Translation of the individual reovirus mRNAs is regulated,

and the translation frequency of mRNAs differs as much as 100-fold. Uncomplexed s3 with m1 down-regulates dsRNA-activated protein kinase, an [interferon](#)-induced enzyme that inhibits initiation of translation by **phosphorylating** eIF-2a, by sequestering its dsRNA cofactor. Later on in infection, newly made capsid proteins assemble with plus-strand single-stranded RNA into capsid. Nonstructural proteins with ssRNA binding activity, such as mNS and sNS, may play a role in the assortment of the mRNA molecules. However, the mechanism of how one complete set of 10 mRNAs is correctly packaged into each new core remains unknown. The synthesis of dsRNA (the RNA replication process) takes place in the assembled particles, in which the replicase catalyzes a single round of minus-strand synthesis.

When cells are coinfecting with two different reovirus strains, a high percentage of the progeny contain novel assortments of gene segments. This process is known as reassortment. New virus strains can arise readily, and evolution occurs by such means. The reassortment is most efficient when the coinfecting viruses are closely related, and sequence divergence in the noncoding regions greatly affects the efficiency of this phenomenon. There is virtually a random distribution of most segments, but some parental **alleles** (such as L1-L2 and L1-M1) exhibit preferential cosegregation in the reassortants. Certain constellations of segments were isolated repeatedly in *in vivo* experiments.

Reassortment has played a major role in identifying (mapping) the viral genes responsible for several biological phenotypes, especially in determining the molecular basis of pathogenesis. For example, the following associations are found: The M2 gene encoding **myristoylated** protein m1 is a major determinant of neurovirulence; the M1 and L2 genes are correlated with myocarditis; and S1 is associated with pneumonia and growth in intestinal tissue. It is noteworthy that the genetic background of the recipient virus, onto which the genes of a donor virus are imposed, can alter the expression of the donor genes. In addition, reovirus reassortants appear commonly to contain mutations in genes that improve their fitness for independent replication.

To date, the genetic engineering of reovirus genes has been restricted to studies in which mutant gene products are analyzed for activity *in vitro* or in transfected cells. Self-assembled core-like particles have also been generated in cells infected with recombinant [vaccinia virus](#) expressing several reovirus gene segments, and this system is useful for studies on morphogenesis and [protein-protein interactions](#). Ultimately, an important development would be a system for introducing synthetically modified or foreign nucleic acid into viable virus particles, because this would revolutionize various aspects of reovirus studies; this is hampered, however, by the mysterious assortment mechanism of the segmented genome.

#### Suggestions for Further Reading

M. L. Nibert, L. A. Schiff, and B. N. Fields (1996) "Reoviruses and Their Replication". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 1557–1596.

### Repeated DNA Sequence Interspersion

A large fraction of the **DNA** in the [genomes](#) of **eukaryotes** is made up of families of similar sequences, with members scattered widely throughout the DNA, called interspersed repeated sequences. The members of such a family are related to each other in DNA sequence but are not precise copies, and some may be quite **divergent**. In typical higher eukaryotes, there are many such families, with the number of copies ranging from less than 10 to almost a million. The typical length for short interspersed repeats ([Sines](#)) is 300 nucleotides. These sequence families are distinct from

**microsatellite** and [minisatellite DNA](#), which are principally clusters of repeated copies arrayed in tandem and repeated fairly accurately. There are also families of longer interspersed repeats ([Lines](#)), many of which are [transposable elements](#). It is very likely that almost all of the short and long interspersed repeats are the result of insertions of DNA sequences. Much is known about the transposable elements in *Drosophila melanogaster* ([1](#)). No clear function is known for most interspersed sequences, but a few of them have been identified with roles in the regulation of **gene expression** ([2, 3](#)).

The most extreme case of interspersion is the set of about half a million copies of the 283-bp [Alu sequence](#) repeats of the human genome, which is considered to be a [retrotransposon](#) ([4, 5](#)). RNA transcripts are derived from a few genes transcribed by **RNA polymerase III**, including the 3' [poly-A](#) tail. Other parts of the genome supply the **reverse transcriptase**, and the RNA transcripts are converted into DNA, which are inserted into the genome by an unknown mechanism. At least half a million copies have been inserted, with an average spacing of only about 4 kbp. They occur in all regions of the chromosomes, except that they (*1*) are correlated with [G bands](#), (*2*) are principally excluded from **exon** sequences coding for protein, and (*3*) are more frequent in the 3' untranslated regions and **introns**. The pattern is nearly random, except that copies are clustered near some genes. The Alu family can be considered the prototypical interspersed repeat, but other families have entirely distinct mechanisms of replication and insertion ([6](#)). The mammalian-wide interspersed repeat (MIR) is an example ([7](#)), since it does not include the 3' poly-A sequence. This sequence is widespread throughout mammalian genomes and has a conserved inner core sequence. Members can be widely divergent within the genome of an individual, but some sequences show little or no divergence between distant mammals. The conservation of the MIR sequence suggests that it has a function that offers an advantage to the host, or perhaps just to the MIR itself.

#### Bibliography

1. M. G. Kidwell and D. Lisch (1997) Proc. Natl. Acad. Sci. USA **94**, 7704–7711.
2. R. J. Britten (1996) Proc. Natl. Acad. Sci. USA **93**, 9374–9377.
3. R. J. Britten (1997) Gene **205**, 177–182.
4. C. W. Schmid (1996) Prog. Nucleic Acid Res. **53**, 283–319.
5. R. J. Britten (1996) in *Human Genome Evolution*, (M. Jackson, T. Strachan, and G. Dover, eds.), BIOOS, Oxford, pp. 211–228.
6. A. F. A. Smit (1996) Curr. Opin. Genet. Develop. **6**, 743–748.
7. A. F. A. Smit and A. D. Riggs (1995) Nucleic Acids Res. **23**, 98–102.

#### Repertoire

Because the [immune response](#) is specific, a basic question is to understand what are the structural and genetic bases that allow the immune system to make the necessary huge number of different recognition molecules—that is, [immunoglobulins](#) (Igs) and [T-cell receptors](#) (TCRs) that are required to interact specifically with all the potential [antigens](#) of the living world. This is what immunologists call the *repertoire problem*.

Is it possible to give an estimate of the number of potential [epitopes](#) that may eventually be encountered by an individual? Several theoretical approaches have been proposed. One is based on the average size of the [antibody](#) combining site, which may be taken as  $600 \text{ \AA}^2$ , although large



variations do exist from one antibody to another (see [Immunoglobulin Structure](#)). Given this size, one may question how many different chemical organic structures could theoretically fill one antibody combining site. Calculations have been made, and the number is astronomic, about  $10^{17}$  different structures. Another possible approach would be to start from an estimate of the number of gene products that can be generated by the living world. To stick to a minimal value, one may limit oneself to animal species (which already excludes viruses, bacteria, and plants!). The evaluation of the number of such species is somewhat greater than  $10^6$ . Taking an average number of  $10^5$  genes/species leads to  $10^{11}$  different protein molecules. One now has to take into account the **allelic** variants and the number of individual epitopes that can be defined at any protein surface. This easily comes to  $10^{14}$ . This order of magnitude is reasonably consistent with the  $10^{17}$  previously given, because we excluded a large part of the living world in the later estimate. Both values are simply enormous, so the repertoire problem is a major issue. Because **T-cell** epitopes are linear [peptides](#), with a variable length [9 for those presented by the [major histocompatibility complex](#) (MHC) class I molecules and a little more for those bound to class II], a calculation may also be made of the potential number of different peptides that might be generated. Again, the number would be within somewhat similar ranges, perhaps slightly lower.

What is the situation on the immune system side? Both Igs and TCRs are generated from mosaic genes that must rearrange before becoming functional. Because rearrangements are made on a combinatorial basis, one may give an estimate of a minimum number of discrete Igs and TCRs that can be generated from this mechanisms. Taking into account [gene rearrangements](#) and random association of chains (H and L for Ig, ab or gd for TCRs), one gets values around  $10^7$  in both cases. This, however, does not take into account **N diversity**, which brings an additional increment ( $10^3$  or more), nor does it take into account [somatic hypermutations](#), which are, however, strictly restricted to Ig molecules. One must also take into consideration the fact, that despite its specificity, there is some degeneracy in immune recognition. It is therefore very clear that the immune system may potentially face recognition of any epitope. It certainly is quite remarkable that all this diversity is generated by fewer than 500 genes—that is, less than 1% of the entire genome (as expressed from the total number of gene segments, with the actual space actually occupied by the Ig and TCR loci being much lower).

Repertoire diversity is the result of a number of cumulative causes, especially combinatorial mechanisms:

- Degeneracy of recognition
- Combinatorial gene rearrangements of H and L chains of Ig or of a/b and g/d for TCR, including D gene fusion, different [reading frame](#) use, and so on
- Random association in polypeptide chain pairing (HL, ab, gd)
- [Junctional diversity](#)
- N diversity
- Somatic hypermutations (restricted to Ig chains)

All theoretical combinations are not, however, necessarily used, leading to the distinction between the potential and actual repertoires. For example, some gene segments (**V genes**, **D genes**, or **J genes**) are used more frequently than others. This may result from pure mechanistic reasons, such as subtle structure in RSS (recombination signal sequences) regions, varying the extent to which they are triggered by [recombinases](#). It seems also clear that not all possible H–L polypeptide chain pairs are encountered and that preferential associations take place. It also has long been observed that the expressed repertoire could vary in life, with some genes preferentially expressed at birth and no longer expressed in the adult. This may of course reflect selective mechanisms that possibly involve [idiotype](#) interactions implying a degree of connectivity that may vary with age. Ultimately, the expressed repertoire of any individual evidently reflects its own history and the various and numerous contacts it may have encountered with the outside world.

See also entries [Antigen](#), [Antibody](#), [Immunoglobulin](#), [V Genes](#), **D gene**, and **J gene**.

### Suggestions for Further Reading

F. Melchers (1997) Control of the sizes and contents of precursor B cell repertoires in bone marrow. *Ciba Found. Symp.* **204**, 172–182.

K. Rajewsky (1996) Clonal selection and learning in the antibody system. *Nature* **381**, 751–758.

A. Nobrega, M. Haury, A. Grandien, E. Malanchere, A. Sundblad, and A. Coutinho (1993) Global analysis of antibody repertoires. II. Evidence for specificity, self-selection and the immunological “homunculus” of antibodies in normal serum. *Eur. J. Immunol.* **23**, 2851–2859.

## Repetition Frequency

*Repetition frequency* is the frequency with which particular DNA sequences are repeated tandemly within a [genome](#) (see [Repetitive DNA](#)). The basis for the phenomenon is not understood. To account for the observations, a model of sequence-dependent, **unequal crossing-over** and [gene amplification](#) (*slippage replication*) has been simulated ([1](#)). It was deduced that DNA, whose sequence is not maintained through [natural selection](#), will exhibit repetitive patterns over a wide range of rates of genetic [recombination](#) resulting from the interaction of unequal crossing-over and slippage replication, both processes that depend on sequence similarity. At high crossing-over frequencies, the nucleotide patterns generated by the simulations were simple and highly regular, with short, nearly identical sequences repeated in tandem. Decreasing recombination rates increase the tendency to longer and more complex repeat units. Periodicities have been observed down to very low recombination rates (one or more orders of magnitude lower than the [mutation](#) rate). At such low rates, most of the sequences contain repeats that have an extensive substructure and a high degree of heterogeneity among themselves. Other high-order structures are superimposed on a tandem array. These structures have been compared to various real structures of tandemly repeated DNA known from eukaryotes ([1](#)).

Notwithstanding the interest in this type of simulation and in the many observations on several long fragments of DNA, it is clear that only the total sequence and characterization of **exons** and **introns** in several **eukaryotes** will give us information about the periodicity of interspersed sequences, if any.

### Bibliography

1. W. Stephan (1989) *Mol. Biol. Evol.* **6**, 198–212.

## Repetitive DNA

*Repetitive DNA*, in which multiple copies of the same or very similar sequences are in the [genomes](#) of **eukaryotes**, was first detected from the kinetics of reassociation of fragmented, **denatured** genomic DNA ([1](#)) (see [C0t Curve](#)). Two types of repetitive DNA were detected: ([1](#)) the fast-

renaturing component consists of rather short sequences tandemly repeated thousands of times, such as those found at the [telomeres](#) and the *variable number tandem repeats* (VNTR;). The slower component consists of longer sequences ([Sines](#), including [Alu sequences](#) in primates, and [Lines](#), such as Kpn in primates and MIF in rodents), which are repeated several hundred to several hundred thousand times.

Higher eukaryotic genomes often consist of more than 50% repetitive DNA. There are some extreme variations among vertebrates. For example, *Fugu rubripes*, the puffer fish, which has apparently all or almost the same number of structural **genes** as other vertebrates, has a genome much smaller than that of primates because of considerably smaller introns and probably a lesser number of repetitive sequences.

## Bibliography

1. J. R. Wu et al. (1977) Proc. Natl. Acad. Sci. USA **74**, 4382–4386.

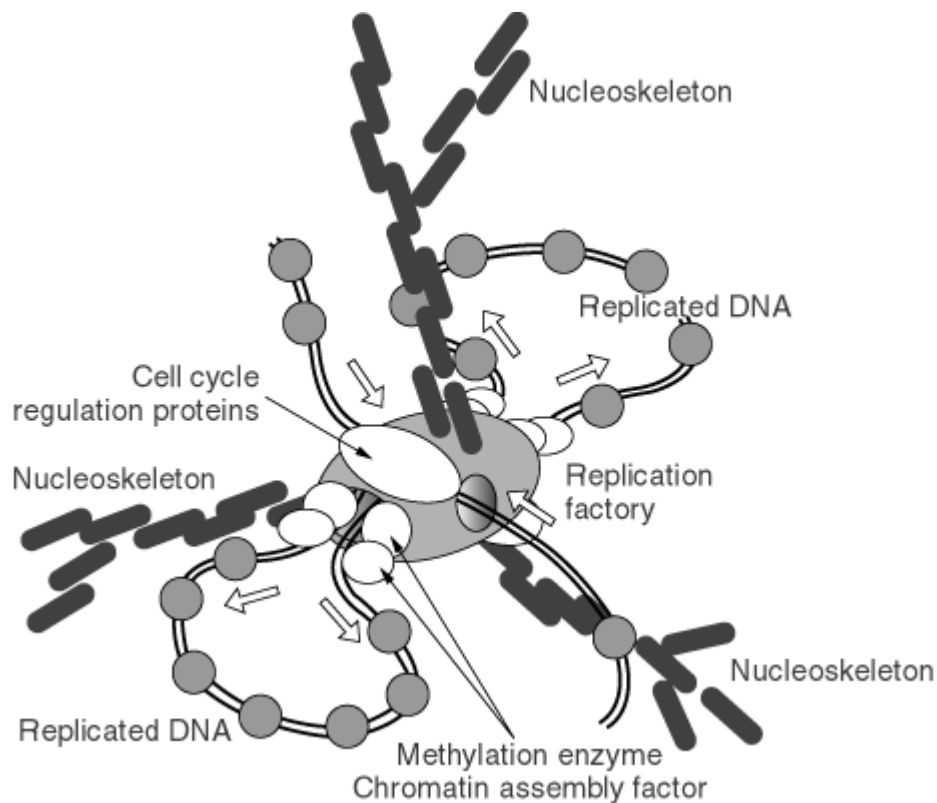
## Replication Foci (Factories, Centers)

During **S phase** of the [cell cycle](#), multiple [replication origins](#) of eukaryotic [chromosomes](#) initiate [DNA replication](#), and numerous [replication forks](#) travel along the DNA for long distances. The positions in mammalian cell [nuclei](#) where DNA synthesis occurs were analyzed by the pulse-labeling of newly synthesized DNA with bromodeoxyuridine (see [5-Bromouracil](#)) and visualization of the sites by fluorescence [microscopy](#). The results showed that the replication sites do not distribute diffusely throughout nuclei but form discrete foci called *replication foci* (or *replication centers*). In rat fibroblasts, more than a hundred small foci appear in early S phase, and the foci subsequently become fewer and larger through several waves of appearance (1). This profile is thought to result from a transition of replication activation from early to late origins during S phase. Similar foci were observed when [sperm](#) with their membranes removed replicated in *Xenopus* [egg](#) extracts (2). These foci are thought to correspond to structures where replicating DNA and replication machineries assemble in an organized manner. Each focus is thought to be composed of several tens of replicating forks from synchronously initiated origins that cluster in a chromosome domain subjected to a common type of **gene expression**. In general, transcriptionally active gene loci are replicated early in S phase, whereas inactive loci or repetitive DNA loci are replicated later.

One replication protein, [proliferating cell nuclear antigen](#) (**PCNA**), was originally identified as an autoimmune [antigen](#) that specifically localized in S-phase nuclei. The immunostaining showed a characteristic punctuated pattern in nuclei and colocalized with replication foci (3). Thus, replication foci are the sites of replication protein assembly at specific nuclear compartments. Many **replication proteins**, such as **DNA polymerase**  $\alpha$ , PCNA (4), RPA ([single-stranded DNA binding protein](#); (5), and [DNA Ligase I](#) (6), are associated with replication foci. Assembly of RPA in nuclei also occurs prior to the start of DNA synthesis and is poised for initiation. This means that the formation of replication foci has been initiated at the pre-replication stage at particular sites in nuclei, and S-phase entry signals may change the complexes to replication foci by the addition or activation of a subset of replication proteins. Replication foci include several additional components, such as **cyclin-dependent protein kinases** (5), cytosine [methyltransferase](#) (7), and [chromatin](#) assembly factor 1 (8). Therefore, replication foci are not merely assemblies of the eukaryotic replisome, but functionally organized apparatuses that respond to progression of the cell cycle, modification of the replicated DNA, and formation of chromatin structure, in addition to DNA replication.

One question about replication foci is how the protein assembly is spatially organized in the nucleus. The structure of replication sites in nuclei was visualized by [electron microscopy](#) using HeLa cells encapsulated in [agarose](#) microbeads, followed by labeling with **biotin-dUTP** of the permeabilized cells (4). Sites of DNA synthesis could be seen in specific dense structures attached to a diffuse nucleoskeleton, appearing at the end of G1 phase, increasing in size and decreasing in number according to the progression of the S phase, similar to the replication foci observed by light microscopy. From their morphology, the dense structures are called *ovoid bodies*, which are distinguishable from other nuclear bodies. The ovoids are strung along the nucleoskeleton and associate with PCNA, in addition to replicating DNA. After short labeling of cells, all replication sites are associated with ovoids, but longer labeling caused the sites to spread into adjacent chromatin. These observations strongly suggest that the ovoid bodies are replication protein assemblies fixed on the nucleoskeleton (Fig. 1) that function as a factory where replication occurs as the template moves through it. In this respect, the protein assemblies at ovoids are called *replication factories*, which are identical structurally to replication foci.

**Figure 1.** Schematic drawing of a replication factory. One factory unit has a set of replication apparatus, cell-cycle regulation proteins, methylation enzymes, and chromatin assembly factors. An actual factory is composed of more than 10 units and many replicating DNAs. These complexes are fixed on the nucleoskeleton.



## Bibliography

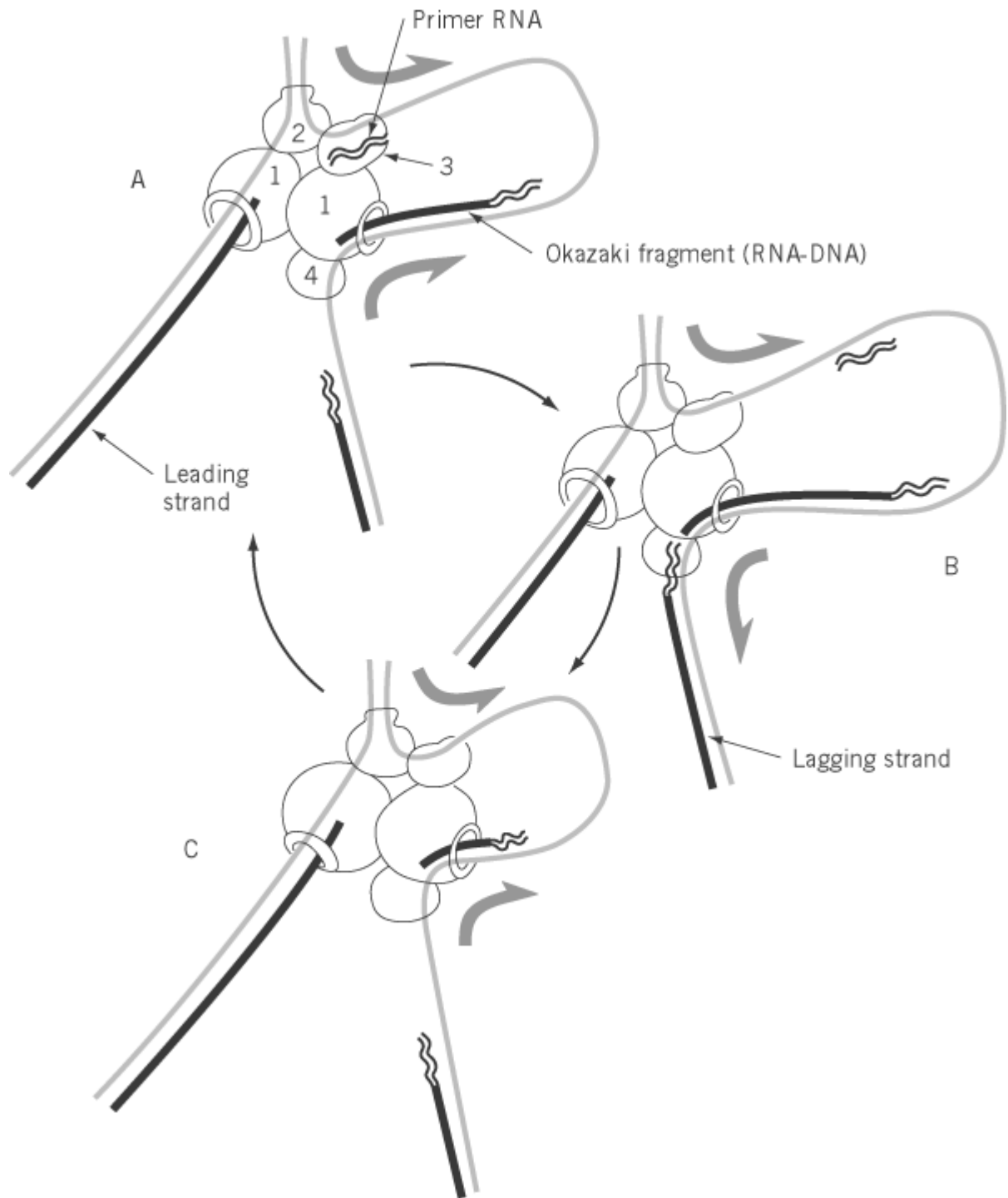
1. H. Nakamura et al. (1986) *Exp. Cell Res.* **165**, 291–297.
2. J. J. Blow and R. A. Laskey (1986) *Cell* **47**, 577–587.
3. R. Bravo and H. Macdonald-Bravo (1987) *J. Cell Biol.* **105**, 1549–1554.
4. P. Hozak et al. (1993) *Cell* **73**, 361–373.
5. M. C. Cardoso et al. (1993) *Cell* **74**, 979–992.
6. A. Montecucco et al. (1995) *EMBO J.* **14**, 5379–5386.

7. R. L. Adams (1995) *BioEssays* **17**, 139–145.
8. T. Krude (1995) *Exp. Cell Res.* **220**, 304–311.

## Replication Fork (Y-Fork Intermediate)

Elongation of [DNA replication](#) proceeds by unwinding the DNA double helix and synthesizing complementary strands on each separated single-stranded template. This generates a fork-like structure at the growing point of replicating chromosomes. The fork structure can be detected as Y-fork intermediates by [electron microscopy \(EM\)](#) and by [two-dimensional gel electrophoresis](#) of replication intermediates (1). The actual structure *in vivo* is thought not to be a simple fork, but to have a complicated configuration to coordinate the asymmetrical synthesis of [leading and lagging strands](#) in opposite directions. In *Escherichia coli*, several models for the replication fork have been proposed, the most accepted being the concurrent synthesis of two strands by an asymmetric complex of DNA polymerase III holoenzyme proposed by Kornberg (2). In this model, [DNA helicase](#) (DnaB)-[primase](#) complex (the primosome) located at the fork precedes the replication to produce two single-stranded loops, one of which contains RNA primer synthesized by the primosome. **DNA polymerase III** holoenzyme located behind the primosome forms an asymmetric dimer structure containing two different catalytic sites; one polymerizes the leading strand progressively and the other synthesizes the lagging strand discontinuously from RNA primer. Concurrent replication of two single-stranded loops by the single holoenzyme results in a gradual decrease in the size of the leading-strand loop, but causes a dynamic change in the lagging-strand loop, first increasing in size and then suddenly shrinking back to a simple fork structure. This generates a cyclic and dynamic change in structure at the growing fork that is comparable to the movement of a trombone musical instrument (Fig. 1). In **eukaryotes**, the fork structure seems to be more complicated than that of **prokaryotes**. Studies on **SV40 virus** DNA replication *in vitro* have shown that at least two DNA polymerases, polymerase d/e and polymerase  $\alpha$ -primase complex, are involved separately in the synthesis of leading strand and lagging strand, respectively (see [DNA Replication](#)) (3). In addition, the dissociation of any [nucleosome](#) structure precedes the fork movement, and reconstruction of the nucleosome occurs concomitantly with the synthesis of new DNA strands (4). However, essentially the same structure as that of the *E. coli* fork was proposed where a heterogeneous complex between DNA polymerases d and a replaces DNA polymerase III holoenzyme. The principle of the fork structure model is that the primosome, polymerases, and other proteins involved in the synthesis of leading and lagging strands form a supramolecular complex, the replisome, through which the chromosome moves, making a dynamic and cyclic change of the Y-fork structure. Assembly of multiple replisomes on the **nuclear** membrane or matrix is thought to form a large structure like a replication focus that is visible by cytological methods (5).

**Figure 1.** Dynamic structure of replication fork. A model for the cooperative synthesis of leading and lagging strands by a replication machinery composed of (1) dimeric form of DNA polymerase III, (2) [DNA helicase](#), (3) [primase](#), and (4) DNA polymerase I and [DNA Ligase](#). The structure of the replication fork changes in a cyclic manner, from A to B to C to A. The partially single-stranded loop moves through the machinery back and forth (shown by shaded arrows) as the elongation (A and C) and completion of **Okazaki fragment** synthesis occur (B). Primer RNA is synthesized only at certain intervals.



### Bibliography

1. B. J. Brewer and W. L. Fangman (1987) *Cell* **51**, 463–471.
2. H. Maki, S. Maki, and A. Kornberg (1988) *J. Biol. Chem.* **263**, 6570–6578.
3. B. Stillman (1994) *Cell* **78**, 725–728.
4. B. Stillman (1986) *Cell* **45**, 555–565.
5. M. Nakamura, T. Morita, and C. Sato (1994) *Exp. Cell Res.* **165**, 291–297.

## Replication Origin

A replication origin is a site at which a **chromosome** or **minichromosome** initiates the **DNA replication** process. Kornberg and colleagues explored in some considerable detail the molecular mechanisms controlling the highly regulated initiation of chromosomal replication in *Escherichia coli*. This event normally occurs once per cell generation at a single site selected from the entire *E. coli* **genome** ( $4 \times 10^6$ bp). This unique chromosomal origin (oriC, 245 bp long) is recognized by a sequence-specific DNA-binding protein, DnaA, that associates with four noncontiguous 9-bp repeats. Because the DnaA protein functions specifically to determine the site at which replication initiates, it can be called a specificity factor. The protein interacts with individual 9-bp repeats, but only when four are placed in the correct positions relative to each other is a large complex of DnaA protein and the oriC DNA sequence formed. This large nucleoprotein complex is recognized by the other proteins required for replication.

The DnaA protein interacts with DNA and also associates through direct protein-protein contacts with other DnaA molecules. This results in a cooperative association of 20 to 30 DnaA molecules with oriC. DNA is wrapped around the complex of the DnaA protein. The DnaA protein–DNA complex formed at oriC facilitates a specific opening of the **double-helix** in an adjacent AT-rich DNA sequence. The DnaC protein mediates the association of the DnaB helicase with this single-stranded AT-rich sequence, unwinding a substantial segment of the double helix, at which replicative enzymes, such as DNA **primase** and **DNA polymerase**, begin to act. Interestingly, the protein most analogous to a histone in *E. coli*, the HU protein, facilitates the formation of a functional nucleoprotein complex at oriC, as do topoisomerases. The role of these proteins is probably to facilitate the correct topological arrangement of DNA for the subsequent replicative events.

Several general rules follow from this analysis. The initiation of replication requires exceptional levels of control, as do similar events in **eukaryotic** cells. Although DNA-binding proteins interact with specific sites with high affinity (with **dissociation constants**,  $K_D = 10^{-9}$ – $10^{-13}$ M), they also bind DNA nonspecifically ( $K_D = 10^{-6}$ – $10^{-3}$ M). Thus the basis of the exceptional precision of replicative events in *E. coli* is unlikely to follow from the binding of a single protein to a single DNA site. Instead, the cooperative interaction of a particular protein (the DnaA protein) with multiple sites over a 200- to 300-bp region of DNA is required for initiation. The precise organization of DNA into these complexes is necessary for other proteins to recognize the origin of replication. This precise organization also requires the mediation of proteins that alter the DNA conformation (HU or IHF) and remove topological constraints to the folding of DNA (topoisomerases).

The precision of regulated events in eukaryotic cells is determined by following the prokaryotic paradigm and by utilizing multiple, high-affinity, sequence-specific **DNA-binding proteins** that recognize multiple related sequence elements and by masking many of the nonspecific sites by folding them into **chromatin** (1). The first examples of multiple binding sites within DNA for particular sequence recognition proteins that regulate replication or transcription in eukaryotic cells were observed in viral systems. The **SV40** genome has three binding sites for the **T Antigen** at the origin of replication (ORI), which is a nucleosome-free region in the **minichromosome**. **Scanning transmission electron microscopy** reveals that trimers and tetramers of T antigen bind at each of these sites. Like DnaA, the T antigen also causes significant changes in local DNA structure at the origin. Other viral genomes have a similar requirement that multiple binding sites be occupied by a particular virally encoded protein to initiate replication. The nuclear antigen (EBNA-1) protein of **Epstein–Barr virus** has six specific binding sites in the viral replicative origin region. Thus eukaryotic viruses are likely to regulate replication comparably to *E. coli*.

In *Saccharomyces cerevisiae* chromosomes, a multiprotein complex has been identified that interacts with over 150 bp of DNA at several functionally defined origins of chromosomal replication (or [autonomously replicating sequences](#), [ARS]) (2). This origin recognition complex (ORC) requires that ATP binds specifically to DNA, a property shared with SV40 T antigen and the *E. coli* DnaA protein. Although only a short consensus sequence is necessary for ORC binding, the final DNA sequence organized by the ORC is much more extensive and resembles the *in vivo* **footprint** at a yeast origin of replication.

The utilization of origins of replication within the chromosomes of metazoans is much more complex than that of *S. cerevisiae* (3). Although it is possible to define limited segments of DNA as small as 500 bp where bidirectional replication is initiated, no consensus sequences that independently direct replicative initiation have been defined. Extended chromosomal regions as large as several kilobase pairs in length can facilitate replicative initiation. A model was proposed to explain these results reflecting the Jesuit maxim that “many are called, but few are chosen.” This suggests that, although there are many detectable sites at which replication might begin, the assembly of DNA into nucleosomal arrays, followed by the subsequent assembly of higher order chromatin structures, represses many potential origins and potentiates the activity of others. This might also contribute to the selective utilization of ARS elements as origins of replication in yeast. This model does not offer any explanation as to why and how a replication-competent nucleoprotein complex is assembled at an origin (4).

Arguments that nuclear organization influences the initiation of DNA replication derive from *in vitro* experiments in *Xenopus* egg extracts and *in vivo* experiments with the chromosomes of *Drosophila* and *S. cerevisiae*. Chromatin assembly and nuclear structures are both necessary for replication in *Xenopus* egg extracts. Additional evidence that chromatin structure influences origin utilization comes from the suppression of origin utilization in regions of *S. cerevisiae* and *Drosophila* chromosomes that contain heterochromatin. Nuclear scaffold or [matrix attachment regions](#) have also been proposed as facilitating aspects of the replicative process (see [Domain, Chromosomal](#)). Chromosomal replication within the nucleus clearly occurs within morphologically defined factories attached to a nuclear scaffold (5). Finally, the nuclear envelope has a regulatory role in determining the activity of the “licensing” process that enables replicative initiation.

#### Bibliography

1. S. - Y. Lin and A. D. Riggs (1975) *Cell* **4**, 107–111.
2. S. P. Bell and B. Stillman (1992) *Nature* **357**, 128–135.
3. D. M. Gilbert et al. (1993) *Cold Spring Harbor Symp. Quant. Biol.* **58**, 475–485.
4. M. L. DePamphilis (1993) *Ann. Rev. Biochem.* **62**, 29–59.
5. P. Hozak, A. B. Hassan, D. A. Jackson, and P. R. Cook (1993) *Cell* **73**, 361–373.

#### Suggestion for Further Reading

6. A. Kornberg and T. A. Baker (1991) *DNA replication*, 2 ed., W. H. Freeman, New York.

#### Replicative Form

The circular **genomes** of small, single-stranded DNA **bacteriophages**, such as M13 and FX174, are converted to a double-stranded form, called the **replicative form (RF)**, upon infection of a host



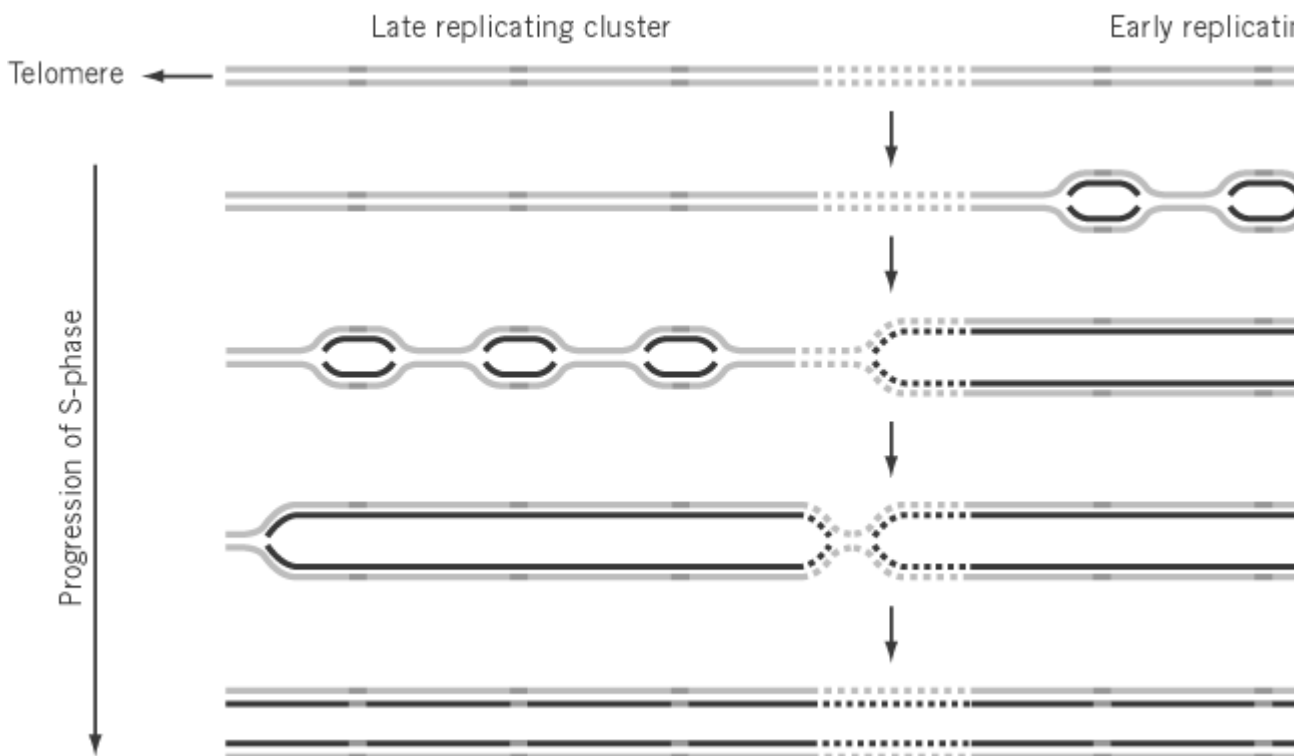
bacterium (see [Single-Stranded DNA Replication](#)). Synthesis of the complementary strand of the original single strand is primed using a RNA [primer](#) made at a specific site on the viral strand and elongated by the host **DNA polymerase III** holoenzyme from the primer terminus. The initial product of complementary strand synthesis is a duplex circle DNA with a small gap, called RFII. It consists of an intact viral DNA and a nearly full-length linear complementary DNA plus the RNA primer. The RNA primer is then removed, and the resultant gap filled by DNA polymerase I. Subsequently, the nick is sealed by [DNA Ligase](#), and the duplex DNA is converted to the **supercoiled** form (RFI) by DNA gyrase (see [DNA Topology](#)). The RFI DNA serves as a [template](#) for multiplication of the viral genome and for [transcription](#) of [messenger RNA](#) for the production of phage proteins.

## Replicon

[DNA replication](#) of a [genome](#) is a self-regulatory process, in which each genome or portion of a genome is replicated independently of each other. Such a unit of replication is called a *replicon*. The concept of a replicon was proposed by Jacob et al. in 1963 (1) as a hypothetical entity that consisted of a replicator, a signal sequence on the **chromosome**, and an initiator, a positive regulator that recognizes its own replicator to initiate replication. The hypothesis was proposed to explain the replication and segregation of chromosome and **plasmid** independently of each other in *Escherichia coli*. In other words, both a *cis*-regulatory element, the replicator, and a *trans*-regulatory element, the initiator, are specific for each replicon. **Prokaryote** and **virus** genomes are found to be composed of a single replicon, each of which contains its own replicator and initiator, as was prophesied by the replicon hypothesis. Both replicators and initiators are variable in structure in plasmid, **bacteriophage**, and viral genomes. In contrast, the initiator, DnaA protein, and its recognition sequence, DnaA-box, are conserved in **eubacteria**. In all cases, the replicator is identical to the [replication origin](#), as replication initiates from within the replicator.

In contrast to prokaryotic genomes, those of **eukaryotes** are composed of multiple replicons of variable size, 10 to 330 kbp (2). A cluster of replicons seems to be replicated simultaneously, but different clusters at varying times during **S phase**, depending on their location along the chromosome. In general, clusters near the **centromere** tend to replicate early, and those near the [telomere](#) replicate late in S phase (Fig. 1). No general picture of a replicon as to replicator and initiator in eukaryotes has been devised, since the fine-structure analysis of chromosomal origins is limited to a single species of eukaryotes, *Saccharomyces cerevisiae* (3). About 20 [autonomously replicating sequences](#) (**ARS**) have been isolated from *S. cerevisiae* chromosomes, some of which are found to function as origins for the chromosome. Every ARS contains an 11-bp ARS consensus sequence (ACS) that is recognized and bound by a protein complex consisting of six subunits, the [origin recognition complex](#) (ORC). The ACS and ORC correspond to the DnaA-box and the DnaA of bacterial replicons. However, the ORC itself is not an active regulator but requires other regulators for the activation of origins. Since the origins of the yeast genome are initiated at different times in S-phase, and vary in the efficiency of their use in the [cell cycle](#), they seem to be regulated specifically by factors that act on regulatory elements of the ORC-ACS complex common to the origins. Replicons of chromosomes containing multiple origins are difficult to define, because termination sites between two adjacent origins are not clear, and in many cases inefficient origins are passively replicated by the **replication fork** coming from neighboring active origins.

**Figure 1.** A scheme of replication of multiple replicons of eukaryotic chromosomes.



Replicons of higher eukaryotes are ambiguous, because replication initiates randomly from a region ranging from 0.5 to 55 kbp (4). Such a region is called an “initiation zone” and contains no specific signal sequences corresponding to the yeast ACS, although six subunits of ORC were found to be widely conserved from yeast to human cells. The size of a chromosome (more than 10 kb) and the structure of its [chromatin](#) seem to be more important than sequence in determining the replicator function of replicons in higher eukaryotes.

#### Bibliography

1. F. Jacob, S. Brenner, and F. Cuzin (1963) *Cold Spring Harbor Quant. Biol.* **28**, 329–348.
2. R. Hand (1978) *Cell* **15**, 317–325.
3. C. S. Newlon (1996) in *DNA Replication in Eukaryotic Cells* (M. L. Depamphilis, ed.), Cold Spring Harbor Laboratory Press, New York, p. 45–86.
4. M. L. DePamphilis (1996) in *DNA Replication in Eukaryotic Cells* (M. L. Depamphilis, ed.), Cold Spring Harbor Laboratory Press, New York, p. 45–86.

#### Reporter Genes

Reporter genes are used throughout molecular and cell biology, from areas such as gene [cloning](#), gene expression, and functional studies, to use in high throughput drug screens and the development of biosensors. Reporter genes provide visual and quantitative data on gene expression. The primary requirements of reporter genes are that their products be readily detectable, with low background activities from endogenous gene products, and that high levels of ectopic expression of the reporter genes themselves have no adverse effects in the cells in which they are expressed (1). For detection,

reporters fall into three main categories: first, enzymes that convert chromogenic or fluorogenic substrates that can be observed in cells *in situ*, or by spectroscopy in cell or tissue lysates; second, enzymes that catalyze the transfer of radiolabeled groups, which can be detected by liquid scintillation counting or [autoradiography](#); third, bioluminescent proteins, or enzymes that catalyze bioluminescent reactions and that have been cloned and developed as some of the most sensitive reporter genes.

A summary of the reporter genes described here is given in Table 1. See also the following individual entries: [Alkaline Phosphatase](#); [Beta-Glucuronidase](#).

**Table 1. Summary of the Common Reporter Gene Systems**

| Reporter                                    | Origin                        | Substrate  | Major Applications  |
|---|-------------------------------|--|---|
| Chloramphenicol acetyl transferase (CAT)    | <i>E. coli</i> , transposon 9 | [ <sup>3</sup> H]- or [ <sup>14</sup> C] chloramphenicol | <i>In vitro</i> gene expression analysis                                    |
|   |                               | + <i>n</i> -Butyryl CoA                                  | Identification of <i>cis</i> - or <i>trans</i> -acting regulatory sequences |
| b-Galactosidase                             | <i>E. coli lac</i> operon     | Xgal   | <i>In situ</i> analysis of gene expression in transgenic animals            |
|   |                               | FDG  | Labeling living cells for sorting by flow cytometry                         |
|   |                               | CMFDG  | Normalizing transfection efficiencies <i>in vitro</i>                       |
| b-Glucuronidase (GUS)                       | <i>E. coli gus</i> operon     | X-glu  | <i>In vivo</i> gene expression in transgenic plants and yeast               |
|   |                               | MUGlcU   | Analysis of protein targeting and transport                                 |
|   |                               | CFDG-GlcU  | Analysis of plant–pathogen, plant–symbiont interactions                     |
| Human placental alkaline phosphatase (hpAP) | Human placental gene          | BCIP/NBT   | <i>In vitro</i> gene expression studies                                     |
|   |                               | MUP  | <i>In vivo</i> expression studies in transgenic animals                     |
|   |                               | HNPP/Fast Red TR   | Dual labeling with a second reporter  |
|   |                               | ELF-97   | Enzyme-linked tag   |

|                                    |   |  |   |
|------------------------------------|---|--|---|
| Luciferase                         | Firefly                                   | Luciferin + ATP<br>+ O <sub>2</sub> + Mg <sup>++</sup> | of molecular probes<br><i>In vitro</i> readout of<br>gene transcription /<br>high throughput drug<br>screens  |
|                                    | Sea pansy                                 | Coelentrarin + O <sub>2</sub>                          | Dual luciferase<br>assays normalize<br>variable transfection<br>efficiencies  |
| Green fluorescent<br>protein (GFP) | Jellyfish<br><i>Aequorea<br/>victoria</i> | None   | Noninvasive gene<br>expression analysis<br>in livingcells and<br>organisms<br>Developmental<br>studies over time<br>Protein tagging to<br>study dynamics of<br>intracellular events |
| b-lactamase                        | <i>E. coli</i><br>plasmid                 | PADAC /<br>Nitrocephin                                 | <i>In vitro</i> gene<br>expression analysis /<br>high-throughput drug<br>screens  |
|                                    |   | CCF2   | Kinetic analysis of<br>gene expression in<br>livingcells  |
|                                    |   | CCF2/AM  | Cell sorting by flow<br>cytometry   |

---

## 1. Chloramphenicol Acetyltransferase

[Chloramphenicol acetyltransferase](#) (E.C. 2.3.1.28; CAT) was one of the first reporter genes to be developed and, although less versatile than some of the more recently developed reporter genes, is still widely used (2, 3). The enzyme confers resistance to [chloramphenicol](#) on bacteria and is derived from [transposon 9](#) of *Escherichia coli*. It is a trimeric protein with a subunit size of 25 kDa. Significantly, there is no endogenous mammalian counterpart to CAT, so background activities are negligible; furthermore, the protein is relatively stable in mammalian cells, and the [messenger RNA](#) transcript has a short half-life.

CAT catalyzes the transfer of acetyl groups from acetyl-coenzyme A to chloramphenicol:



When radiolabeled substrate (eg, [<sup>14</sup>C]chloramphenicol) is added to a cell lysate, the amount of the acetylated chloramphenicol produced can be measured and is directly proportional to the amount of enzyme present; significantly, this relationship remains linear over three orders of magnitude. The acetylated and nonacetylated forms of chloramphenicol have different solubilities in organic solvents and may be separated by differential extraction or by using **thin-layer chromatography** (2).

Variations on this basic assay include the use of [<sup>3</sup>H]chloramphenicol and the use of *n*-butyryl CoA

instead of acetyl CoA, because *n*-butyrylated chloramphenicol partitions more efficiently into the organic phase upon extraction than the acetylated form. The amount of radiolabeled *n*-butyrylated chloramphenicol extracted may then be determined by liquid scintillation counting (4).

Although CAT assays are very sensitive, CAT is not as versatile as other reporter gene systems described here. It is used primarily to quantify gene expression in transient expression studies—for example, to assay gene **enhancer** or **promoter** activities. At present, CAT activity needs to be assayed in cell extracts, which restricts its use to *in vitro* studies. Nevertheless, CAT has been used in a number of transgenic studies in which expression at a cellular resolution was not required. Despite the lack of an *in situ* assay for CAT activity, the presence of CAT protein in a cell can be detected by immunocytochemistry (5). The availability of good **antibodies** against CAT also allows quantitative analysis of the amount of CAT protein produced by using **enzyme-linked immunosorbent assays** (ELISAs), rather than measuring CAT activity (6).

## 2. b-Galactosidase

**Beta-galactosidase** (E.C. 3.2.1.23) is encoded by the *lacZ* gene from the **lac operon** of *E. coli* and is one of the reporter genes used most widely to study gene expression in transgenic animals. It has been used extensively in all of the major model organisms, including *Caenorhabditis elegans*, *Drosophila*, *Xenopus*, and the **mouse**. Importantly, high levels of ubiquitous *lacZ* expression can be obtained with no apparent adverse effects on development or cell biology (7). An advantage of *lacZ* is its ability to form functional fusion genes. *LacZ* is often used as a fusion gene with the selectable markers for **neomycin** or hygromycin resistance (*neo<sup>r</sup>*, *hyg<sup>r</sup>*). These fusion genes, bgeo and bhyg, allow direct selection of cells that express the reporter gene and have formed the basis of “gene-trap” vectors that are used to identify, mutate, and clone novel expressed genes simultaneously, principally in mouse embryonic **stem cells** (8, 9). The construction of fusion genes also enables the subcellular localization of the reporter to be studied. For example, in cell labeling studies, b-galactosidase is often used carrying a nuclear localization signal (10), which provides a more intense stain within cell bodies, and an *N*-terminal fusion with the neurofilament protein Tau effectively labels axonal projections for studies of neural development (11). Exceptionally, fusion genes carrying an *N*-terminal cleavable **signal sequence** do not show b-galactosidase activity (12). b-Galactosidase is not secreted, but it becomes sequestered in the lumen of the **endoplasmic reticulum** (ER). However, b-galactosidase activity can be restored by adding a transmembrane domain to the fusion gene; in this case, the b-galactosidase becomes anchored within the membrane of the ER, and the b-galactosidase domain of the protein is held in the cytosol. The restoration of b-galactosidase activity of fusion genes carrying secretory and membrane-encoded signal sequences by inclusion of a transmembrane domain at the *N*-terminus of the CD4 gene has formed the basis of an elegant gene-trap strategy to identify novel genes for membrane-bound and secreted proteins (12).

The enzyme comprises a homotetramer, with a subunit size of 116 kDa, and it catalyzes the hydrolysis of various b-galactosides, such as lactose. *LacZ* is most widely used in transgenic experiments, where its expression can be detected *in situ*, in fixed whole tissue preparations, using histochemical substrates such as X-gal (5-bromo-4-chloro-3-indoly-galactopyranoside). Hydrolysis of X-gal in the presence of a ferricyanide/ferrocyanide catalyst releases an indolyl that is oxidized to an indoxyl that forms an indigo blue precipitate (13). This reaction has proved particularly valuable for analyzing gene expression patterns when *lacZ* is expressed from specific promoter sequences, or integrated into endogenous genes by homologous recombination. The stain is clearly visible, does not spread to adjacent nonexpressing cells, and is stable to further histological processing. Although the b-galactosidase reporter gene is derived from *E. coli*, many eukaryotic cells also contain some b-galactosidase activity. For the most part, background staining can be minimized by fixation and by staining at pH 7.5 to 8, because the endogenous b-galactosidases are largely lysosomal and have an optimum pH of about 3.5.

In addition to X-gal histochemistry, a number of other substrates have been developed that have

increased the versatility of b-galactosidase as a reporter gene (13, 14). Hydrolysis of the chromogenic substrate 2-nitrophenyl b-D-galactopyranoside (ONPG) allows quantitative determination of b-galactosidase activity in cell extracts, by measuring the **absorbance** at 420 nm. Quantitative measurement of b-galactosidase is less sensitive than that of other reporter gene activities routinely analyzed in cell extracts, such as luciferase or CAT, but b-galactosidase has a valuable use in normalizing transfection efficiencies when used in combination with the more sensitive reporters.

The most sensitive substrates for b-galactosidase are the **fluorescent** substrates, such as fluorescein di-b-D-galactopyranoside (FDG), which is hydrolyzed to release fluorescein and can be detected using a luminometer or can be observed using standard fluorescein isothiocyanate fluorescence microscopy (13). Although fluorescent substrates increase the sensitivity of detection, higher background levels of activity are also detected due to the endogenous b-galactosidase activity of many cells.

Fluorogenic substrates may also be used to study living cells. For this purpose, derivatives of FDG (C<sub>12</sub>FDG and CMFDG, from Molecular Probes Inc.) have been developed that improve the loading of cells with the substrate and reduce the leakage of released fluorescein from the cells. A major use of this technology has been to sort expressing from nonexpressing cells in a population, using [flow cytometry](#) (15). Sorted cells may then be used for a wide variety of applications. Despite the problem of background fluorescence, living cells may also be observed *in situ* using fluorescent and [confocal microscopy](#). This technology allows cells to be observed in real time and to be identified *in situ* for a subsequent manipulation, such as specific cell injection or performing physiological experiments through patch clamping (16).

### 3. Luciferase

A number of genes for bioluminescent proteins have been cloned and used in reporter systems (17). Amongst the first such gene, and one of the most widely used, is the **luciferase** gene, *luc*, from the firefly *Photinus pyralis*. Luciferase (E.C. 1.13.12.5) is a monooxygenase, a monomer of 61 kDa, that catalyzes the oxidation of its substrate beetle luciferin, in the presence of Mg<sup>2+</sup>, oxygen, and ATP, with the emission of green to yellow light (550 to 570 nm), which can be measured using a luminometer. The luminescent reaction is extremely sensitive; a flash of light is observed 0.3 s after the addition of ATP to the enzyme and substrate mixture, and this flash of light then decays over about 15 s. The light output is directly proportional to the luciferase concentration. Importantly, the enzyme activity is very closely coupled to protein synthesis, and its activity remains linear over a remarkable 10<sup>8</sup> concentration range of the enzyme. Because the light emission is so fast, it was necessary initially to use a specialized luminometer in which the ATP is injected into the sample at the time of recording. However, current protocols include the addition of coenzyme A to the substrate mix, and this reduces the decay of light emission, extending it to several minutes (18). This allows multiple reactions to be set up and assayed using less sophisticated equipment.

In addition to firefly luciferase, luciferase genes have been cloned from bacteria (*lux*) and from the sea pansy coelenterate *Renilla reniformis* (*Rluc*). Despite the common feature of light emission, the luciferases are quite distinct; they do not share common evolutionary origins and have different chemistries and substrate specificities. *Renilla* luciferase is a 31-kDa protein and catalyses the oxidation of coelenterazine, with the emission of blue light at 480 nm. Unlike the firefly luciferase, *Renilla* luciferase gives a low level of nonenzymatic autoluminescence, which reduces its sensitivity by about 10-fold. Seldom used alone, *Renilla* luciferase is often used in conjunction with the firefly luciferase in dual labeling studies. The *Renilla* enzyme is particularly useful as an internal control to normalize for variations in [transfection](#) efficiencies between samples. In such assays, the firefly and *Renilla* enzymes may be assayed from the same cell extracts by the sequential addition of the luciferin and coelenterazine substrates, determining the ratio of luminescence from the two enzymes (19).

The major use of luciferase is to measure gene activity in transient expression assays. For this purpose, luciferase assays are more rapid, more sensitive (two to three times), and cheaper than CAT assays. For these reasons, this reporter gene system has been developed for the requirements of high-throughput screens of gene expression, in which assays can be performed in 96-well luminometry plates. Despite the high sensitivity of detection of luciferase activity cell extracts, efficient protocols have not been developed to monitor expression in intact cells or in tissue sections.

#### 4. b-Lactamase

The most recent reporter gene to be developed for use in animal cells is TEM **b-lactamase** (penicillin amido-b-lactamhydrolase; E.C. 3.5.2.6), the product of the *E. coli* **ampicillin**-resistance gene *Amp<sup>r</sup>*, which functions to cleave **penicillins** and cephalosporins (20). The gene encodes a 27-kDa monomeric soluble enzyme with a high catalytic efficiency. Initial studies using the chromogenic substrates 3-(2,4-dinitrostyryl)-(6R,7R)-7-(2-thienylacetamido)-ceph-3-em-4-carboxylic acid (PADAC) and nitrocephin showed that b-lactamase is a very sensitive reporter, with as few as 10<sup>7</sup> to 10<sup>9</sup> molecules being detectable, which compares well with the limits of detection of other reporters, such as CAT, luciferase, and human placental [alkaline phosphatase](#) (hpAP). Furthermore, as for CAT and hpAP, the steady-state levels of enzyme activity are representative of the steady-state mRNA levels. The enzyme may be secreted, and its activity assayed in the culture medium, or it may be assayed in cell extracts. In either case, the reaction simply requires addition of substrate and measurement of the change in absorbance, measured at 570 nm for PADAC or 470 nm for nitrocephin.

In addition to available chromogenic substrates, Tsien and colleagues have developed fluorogenic substrates for b-lactamase that allow gene activity in single living mammalian cells to be monitored in real time, with very high sensitivity (21). In the system developed, the assay is based on fluorescence resonance **energy transfer** in which the b-lactam ring of a cephalosporin substrate is linked to two fluorophores. Excitation (at 409 nm) of the first, shorter-wavelength, donor fluorophore (coumarin) leads to excitation and green fluorescence emission (520 nm) from the second, longer-wavelength, acceptor fluorophore (fluorescein). Cleavage of the substrate by b-lactamase then separates the two fluorophores, so that excitation of coumarin leads to fluorescence emission of blue light (447 nm) only. Thus the shift in the fluorescence emission wavelength from 520 nm (green) to 447 nm (blue) can be measured, and the ratio of the intensities of the two wavelengths indicates the extent of cleavage, or the level of reporter gene activity.

To study b-lactamase in living mammalian cells, an esterified derivative of the substrate CCF2 (CCF2/AM) was made that could cross cell membranes readily to load cells, without necessitating additional cell permeabilization or shock. The ester groups are then removed in the cell by endogenous nonspecific esterases, to trap the active CCF2 within the cell. One shortcoming of the substrate is that it is free to diffuse throughout the cell; consequently it cannot be used to study subcellular events in the way that **green fluorescent protein** can.

Work to develop the b-lactamase/CCF2 substrate system used cotransfection of b-lactamase under the control of a tandem trimer of nuclear factor of activated T cells (NF-AT) binding sites, together with the M1 muscarinic receptor gene. Stimulation of cells with the muscarinic agonist carbachol then led to a 100-fold increase in b-lactamase expression per cell, with a corresponding switch from green to blue fluorescence. Importantly, b-lactamase activity showed a good dose response to the agonist and revealed kinetic changes in gene transcription at both the population and the single-cell levels. This study also showed b-lactamase to be an excellent reporter of transcriptional readout in drug screens for receptor agonists and antagonists.

Finally, expressing and nonexpressing cells can be sorted easily by flow cytometry. Because the substrate and detection system cause minimal damage to living cells, cells may be selected for

continued culture based on their levels of reporter gene expression. Such cell selection by flow cytometry may circumvent the long and tedious process of picking and analyzing transfected cell clones selected by co-transfection with a drug-resistance gene.

Although b-lactamase is a relatively new addition to the list of genetic reporters, it is likely to become increasingly important in gene expression analysis, because the system combines the sensitivity of enzymatic cleavage of the substrate with the ability to observe, to sort, and, potentially, to manipulate living cells.

## 5. Reporter Gene Vectors

While significant advances have been made in reporter gene technologies, generating new reporter variants and developing novel substrates, advances have been made similarly in developing gene expression vectors that improve the transcriptional and translational efficacy of the reporter genes. A large number of different reporter gene expression vectors are currently available from a range of companies. In many cases, this has involved customizing the genes for the different contexts in which they are used. Many of the reporters are derived from prokaryotic genes or eukaryotic cDNAs. To improve translational efficiency, most reporter genes are engineered to contain an artificial eukaryotic translation initiation sequence (Kozak sequence) (22) and a [polyadenylation](#) signal. In many cases a small **intron** is included in the transcript, because full RNA processing has been found to contribute to message stability and transport. Reporter vectors also include a variety of genetic elements that increase their versatility. To reduce the need for cotransfection, reporter genes may be expressed as part of a bicistronic message, in which translation of the reporter gene may be initiated downstream of a test gene, at a viral internal ribosome entry site (23). Alternatively a gene of interest and the reporter may be expressed in opposite directions from a bidirectional promoter (24).

Traditionally reporter genes have been used to study gene regulation by cloning *cis-acting* test enhancer/promoter sequences into a multiple cloning site upstream of the reporter gene. In addition, reporter cassettes are becoming available in which the reporter genes are placed under the control of a variety of specific **transcription factor** binding sites or **response elements**, which then report on the induction of *trans-acting* events. Such reporter-gene cassettes can give transcriptional readouts in response to a variety of external stimuli, such as receptor binding, [protein–protein interactions](#), and **virus infections**, which induce specific [signal transduction](#) cascades and transcription factor binding. Examples include the NF-TA–b-lactamase reporter described above (21), and reporters under the control of **cyclic AMP** response elements (CRE) (25) or **retinoic acid** response elements (RAREs) (26).

In summary, a range of reporter genes are in common use, some of which are better suited to some applications than others. Many parameters need to be considered in selecting a reporter for a specific experiment, namely message and protein turnover, suitability for *in vitro* or *in vivo* work, assay sensitivity, assay detection, substrate availability, and cost. The field of reporter gene technology is a rapidly advancing field, and novel reporters and substrates will continue to be developed to meet the changing demands of biologists.

## 6. Acknowledgments

NDA is supported by an Advanced Fellowship from the Biotechnology and Biological Sciences Research Council, United Kingdom. I also thank Bridget Snook for discussion and help with the manuscript.

## Bibliography

1. J. Alam and J. L. Cook (1990) *Anal. Biochem.* **188**, 245–253.
2. C. M. Gorman, L. F. Moffat, and B. H. Howard (1982) *Mol. Cell Biol.* **2**, 1044–1051.
3. B. R. Cullen et al. (1987) *Methods Enzymol.* **152**, 684–704.



4. B. Seed and J.-Y. Sheen (1988) *Gene* **67**, 271–277.
5. K. Hanaoka et al. (1991) *Differentiation* **48**, 183–189.
6. A. E. Reifel-Miller et al. (1996) *Biotechniques* **21**, 1033–1036.
7. B. P. Zambrowicz et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3789–3794.
8. G. Freidrich and P. Soriano (1991) *Genes Dev.* **5**, 1513–1523.
9. D. Natarajan and C. A. Boulter (1995) *Nucleic Acids Res.* **23**, 4003–4004.
10. C. Bonnerot et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 6795–6799.
11. C. A. Callahan and J. B. Thomas (1994) *Proc. Natl. Acad. Sci. USA* **91**, 5972–5976.
12. W. C. Skarnes, J. E. Moss, S. M. Hurlley, and R. S. P. Beddington (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6592–6596.
13. R. Sullivan and C. W. Lo (1997) *Methods Mol. Biol.* **63**, 229–246.
14. C. E. Olesen, J. J. Fortin, J. C. Voyta, and I. Browstein (1997) *Methods Mol. Biol.* **63**, 61–70.
15. G. P. Nolan, S. Fiering, J. F. Nicolas, and L. A. Herzenberg (1988) *Proc. Natl. Acad. Sci.* **85**, 2603–2607.
16. N. J. Wright and Y. Zhong (1995) *J. Neurosci.* **15**, 1025–1034.
17. I. Bronstein et al. (1994) *Anal. Biochem.* **219**, 169–181.
18. K. V. Wood (1991) in *Bioluminescence and Chemiluminescence: Current Status*, (A. Szalay, L. J. Krika, and P. Stanley, eds.), Wiley, New York.
19. W. W. Lorenz et al. (1993) in *Bioluminescence and Chemiluminescence: Status Report* (A. Szalay, L. J. Krika, and P. Stanley, eds.), Wiley, New York.
20. J. T. Moore, S. T. Davis, and I. K. Dev (1997) *Anal. Biochem.* **247**, 203–209.
21. G. Zlokarnik et al. (1998) *Science* **279**, 84–88.
22. M. Kozak (1986) *Cell* **44**, 283–292.
23. P. Mountford et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 4303–4307.
24. U. Baron et al. (1995) *Nucleic Acids Res.* **17**, 3605–3606.
25. M. J. Castanon and W. Spevak (1994) *Biochem. Biophys. Res. Comm.* **198**, 626–631.
26. H. C. Tsou et al. (1994) *Exp. Cell Res.* **214**, 27–34.

### Suggestions for Further Reading

27. D. Groskreutz and E. Schenborn (1996) "Reporter Systems." In *Recombinant Proteins: Detection and Isolation* (R. Tuan, ed.), Humana Press, Clifton, NJ.
28. S. R. Gallagher, ed. (1992) *GUS Protocols*, Academic Press, New York.
29. R. P. Haugland (1996) *Handbook of Fluorescent Probes and Research Chemicals*, 6th ed., Molecular Probes Inc., Eugene, Oregon.

### Reporter Groups

Biochemical reactions are most easily followed by routine spectroscopic methods, such as **absorbance spectroscopy** and [fluorescence spectroscopy](#). Spectroscopic methods are fast and sensitive, they allow the on-line detection of changes, and they are nondestructive. The chromophore that changes its absorbance or fluorescence during the reaction of interest (which can be an [enzyme](#) activity assay, **ligand binding**, protein association, or a conformational change) is called a reporter

group.

In the simplest case, the reporter group is an intrinsic constituent of one of the reactants, such as an aromatic amino acid residue of a protein, or a firmly bound **prosthetic group**, such as a heme group. Absorbing or fluorescing **coenzymes** or substrates such as **NADH**, [pyridoxal phosphate](#), or **FAD** are also very useful reporter groups. If natural substrates do not contain chromophores, substrates that are modified with chromophoric groups can be employed. Widely used are nitrophenyl esters or amides as chromogenic substrates for [proteinases](#), esterases, or **phosphatases**, or chromogenic sugar derivatives as substrates for glycosidases. For **nucleotide-binding** proteins, a variety of fluorescing derivatives of adenosine nucleotides are available.

Proteins can be labeled by the covalent attachment of chromophoric reporter groups. Often fluorescing groups, such as dansyl, pyrene, rhodamine, or fluorescein, are attached. These groups are usually linked to [lysine](#) or [cysteine](#) residues. Accessible cysteine residues can be introduced at specified positions of a protein by [site-directed mutagenesis](#). It is important to check that labeling is site-specific and that the activity of the protein is not changed by the covalently-linked reporter group. Detailed lists of chromophores and fluorophores that are useful as reporter groups for biochemical reactions are found in (1-3).

### Bibliography

1. M. R. Eftink (1991) in *Methods of Biochemical Analysis* (C. H. Suelter, ed.), Vol. **35**, Wiley, New York, pp. 127–205.
2. C. R. Bagshaw and D. A. Harris (1987) In *Spectrophotometry and Spectrofluorimetry: A Practical Approach* (D. A. Harris and C. L. Bashford, eds.), IRL Press, Oxford, UK, pp. 91–114.
3. A. N. Glazer, R. J. DeLange and D. S. Sigman (1975) "Chemical Modification of Proteins". *Selected Methods and Analytical Procedures*, North-Holland, Amsterdam.

### Repressors

Encoded within the DNA of every living organism are the **genes** that determine the structures of molecules that constitute the cells. At any given time during its life, the cells that comprise the organism only need the products of a subset of all the genes carried by that cell. The production of unneeded gene products is a waste of precious resources. Moreover, in multicellular organism, these unneeded gene products can disrupt the orderly progression of the organism's growth and development or can cause disease. Hence, the production or expression of the genes within a cell is tightly regulated. Early in the study of gene regulatory mechanisms, Jacob and Monod (1) advanced the idea that gene regulation could occur by preventing or repressing the expression of genes. According to their idea, molecules called *repressors* prevented the cellular machinery from reading and synthesizing the gene products in a process known as [transcription](#). The [operon](#) hypothesis of Jacob and Monod was so persuasive that it was initially thought that all gene regulators would be repressors and function by negatively regulating transcription. Subsequent studies revealed that gene expression is also positively regulated by transcriptional activators. In contrast to negatively regulated, genes whose expression is positively regulated are not expressed in the absence of the regulatory molecule. That is, the gene products are not made unless a regulatory molecule stimulates the cellular machinery to transcribe the gene.

Gene expression can be regulated not only at the level of transcription, but also at the subsequent step of [translation](#) of gene products into [proteins](#) and at the level of regulation of protein and/or

activity. This entry will focus exclusively on the negative regulation of gene transcription. Both transcription and translation are complex multistep processes, however, so there may be parallels in the overall strategy of negative regulation in these two processes. Hence, as the study of the control of translation matures, a variety of mechanisms by which negative regulators function in this process are likely to be found.

An understanding of how repressors of gene transcription function requires an introduction to the process of transcription. Although transcription initiation in eukaryotes is a more complex process than in prokaryotes, the fundamental steps appear to be similar in the two classes of organisms. Thus, for the purposes of simplicity, the description of the steps leading up to the assembly of the cellular machinery capable of juggernaut-like transcription will focus on prokaryotes as a model. Key differences between prokaryotic and eukaryotic transcription initiation will be pointed out in this context.

Initiation of transcription in prokaryotes is comprised of a series of ordered steps. At the outset, the **RNA polymerase** enzyme that is responsible for synthesis of the transcribed product recognizes and binds to a specific sequence in DNA called the **promoter**. Most bacterial promoters contain two conserved regions, namely, the  $-35$  and  $-10$ . Specific binding promoter by prokaryotic RNA polymerases is governed by the **sigma factor** of RNA polymerase (2, 3). Generally, the strength of a promoter is related directly to how tightly the  $\sigma$  subunit binds the  $-35$  and  $-10$  sequences. Most promoters with either a near-consensus  $-10$  or  $-35$  sequence are transcribed with moderate efficiency. The  $\sigma$  subunit is only needed for promoter-specific transcription initiation; in the absence of  $\sigma$ , the RNA polymerase can initiate transcription at the ends of DNA or at nicks or gaps in a nonspecific manner. Moreover, the RNA polymerase-associated  $\sigma$  subunit is not needed for the elongation process and is released from the enzyme during elongation of initiated transcripts. In a step subsequent to specific binding of RNA polymerase to the promoter forming a “closed complex,” several lines of evidence suggest that the enzyme undergoes an isomerization step signaling a more intimate interaction between the enzyme and DNA. In this complex, the so-called intermediate complex, the DNA substrate has not yet unwound, but the enzyme–DNA complex is more resistant to dissociation than is the closed complex (4).

The responsibility for RNA synthesis in eukaryotes is divided between three RNA polymerase enzymes. In all three, specific initiation by the RNA polymerase requires the assembly of a “preinitiation complex” (PIC) (5). The process of forming the PIC will be illustrated by considering assembly of the complex needed to direct transcription of **messenger RNAs** by RNA polymerase II. Unlike prokaryotic RNA polymerase, the promoter-specificity-determining proteins of the PIC assemble at the promoter prior to binding of the RNA polymerase. These proteins bind in an ordered fashion on the promoter DNA, in a complex that is nucleated by the **transcription factor** TFIID, which contains the DNA recognition subunit TATA binding protein (TBP) (for review, see Ref. 6) (see **TATA Box**). This subunit binds in a sequence-specific manner to the TATAAAA sequence located between 20 and 45 bp upstream of the transcriptional start site. Association of the RNA polymerase with this part of the PIC nucleates the binding of several other factors, leading to the formation of an RNA polymerase–DNA complex that is tightly engaged with the DNA, but has not yet unwound the DNA. This complex has been likened to the intermediate complex seen in the transcriptional initiation reaction of prokaryotic RNA polymerases. The protein–DNA and protein–protein contacts that govern the formation of the PIC are relatively weak. Hence, in eukaryotes, efficient gene expression often requires the activity of transcription activator proteins. These proteins often function to recruit members of the PIC to the promoter and/or stabilize their interaction with DNA or each other.

Subsequent to the tight engagement between RNA polymerase and promoter DNA in both prokaryotic and eukaryotic RNA polymerase–promoter DNA complexes, a second isomerization occurs that exposes the bases of the template strand of the DNA. This DNA-unwinding step leads to the formation of the open promoter complex. In this configuration, the RNA polymerase is poised to begin synthesizing RNA rapidly from the nucleotide triphosphate precursors. Within prokaryotic

RNA polymerase, formation of the open promoter complex occurs without the input of energy or the assistance of additional proteins. In eukaryotic RNA polymerases, strand separation requires the hydrolysis of ATP and may require the action of the TFIIF subunit. The ATP is used both to catalyze the **helicase** activity of the PIC and to **phosphorylate** the C-terminal repeated sequences contained on the largest subunit of the RNA polymerase enzyme. Both of these actions appear to be required to release the RNA polymerase from the PIC.

The final step of transcript initiation is the progression of the RNA polymerase enzyme beyond the first few bases of the transcript. At certain promoters, prokaryotic RNA polymerase enters an idling mode wherein it repetitively synthesizes abortive transcripts between 3 and 12 bases in length. Transition from this abortively initiated complex to the elongating RNA polymerase enzyme is marked by the loss of the  $\sigma$  subunit. Although eukaryotic RNA polymerases may initiate abortively, the process does not appear to be a major step along the initiation pathway. Transition to the elongating form of the eukaryotic enzyme may require additional cofactors that promote or inhibit juggernaut synthesis by the enzyme.

In both prokaryotes and eukaryotes, repressors of transcription initiation could, in principle, interfere with any step along the initiation pathway. Moreover, they could act directly on the enzyme itself or function in an indirect manner. The remainder of this entry will focus on specific repressors that act at particular steps along the initiation pathway. The mechanism of each repressor's action will be discussed to the extent that it is known.

### 1. Repressors of RNA Polymerase Binding at the Promoter

The classical model of repressor function put forth within the operon theory of Jacob and Monod (1) proposes that repressors block access of RNA polymerase to the promoter by occluding the RNA polymerase binding site. Although several prokaryotic repressors use this mechanism, it is not as prevalent as previously thought. The most clear-cut example of repression by promoter occlusion is the **lambda cI repressor** blockage of RNA polymerase binding at the  $P_R$  promoter (7). The binding sites for  $\lambda$  cI repressor surround the conserved  $-10$  and  $-35$  sequences in the promoter. Examination of the RNA polymerase and  $\lambda$  cI repressor binding sites reveals the occupancy of the repressor sites by protein blocks access of RNA polymerase to the conserved sequences of the  $P_R$  promoter. This observation does not mean that RNA polymerase cannot bind at in the vicinity of the  $P_R$  promoter sequence. In fact, studies with homologous systems show the RNA polymerase and another cI repressor can simultaneously occupy the DNA at the promoter (J. Xu and G. B. Koudelka, unpublished results). This observation suggests that alternative binding sites for RNA polymerase may exist at  $P_R$ . Hence  $\lambda$  cI repressor-mediated repression at  $P_R$  may occur not by inhibiting RNA polymerase binding at  $P_R$ , but instead by preventing it from forming a complex at the promoter that would be able to initiate transcription.

In eukaryotes, the transcriptionally competent RNA polymerase assembles at the promoter only after formation of a PIC; there are no repressors that specifically prevent access of RNA polymerase to the promoter DNA. However, several mechanisms exist that prevent the assembly of the PIC or association of RNA polymerase with the complex in either a global or promoter-specific manner.

The assembly of the PIC can be inhibited by the formation of **chromatin**. **Nucleosomes** are the primary component of chromatin. DNA associates with nucleosomes in a relatively non-sequence-specific manner by wrapping around this large multisubunit protein. DNA binding by nucleosomes obstructs access of the proteins that form the PIC to the promoter by burying them within the nucleoprotein complex (8). Adding the **histone** H1 to nucleosome–DNA complexes links them together, stabilizing and compacting the nucleoprotein structure, enhancing the binding inhibitory effect (9). Repression of transcription by chromatin can occur globally, resulting in the inactivation of entire regions of the **genome**, or locally, inhibiting expression of only one gene. Local inhibition

appears to result from the presence of nucleosomal positioning sequences near the promoter. These sequences bind with high affinity to the nucleosome, resulting in precise placement of the nucleosome over the assembly point for the pre-initiation complex (10). Global transcriptional repression by chromatin (also known as **silencing**) probably results from the folding of large blocks of nucleosomal DNA into an especially compacted form that is refractory to the formation of active transcription complexes (11). The compacted form of the DNA is very stable, presumably as a result of the attachment of these silenced region to the proteins of the nuclear superstructure (nuclear matrix). The most extensive example of silencing occurs in female mammals, where one of the two [X chromosomes](#) is not used as a transcription template in each cell. [X chromosome inactivation](#) is extraordinarily stable: Once an X chromosome has been inactivated, it is maintained in this state in all subsequent cell divisions, despite repeated replications of the DNA (12).

As alluded to above, most eukaryotic promoters require a constellation of transcription activators in order to be efficiently transcribed. Hence, complete repression of a promoter could be achieved either by interfering with the action of each of a number of activator proteins (see text below) or by acting at a step common to all initiation reactions, the binding of components of the PIC to the promoter DNA. The former mechanism provides a method by which expression levels can be modulated over a wide range. This latter mechanism is analogous to the classical ideas of repressor function and are distinguished from repressors of activated transcription by their ability to repress basal transcription. An example of this type of repressor is the Dr1 protein (13). Proteins that belong to this class of repressors must function early in the PIC assembly pathway, because the more complete complexes are refractory to the repressing effects of these proteins.

## 2. Repressors That Prevent RNA Polymerase Open Complex Formation

In bacteria, the DNA-bound RNA polymerases progress through at least two isomerization steps prior to strand separation to form the open complex. Although repressors could act at either of these two steps, the transient nature of the intermediate complexes has prohibited detailed studies on the rates of formation and of breakdown of this species. This prohibition prevents a complete understanding of the mechanisms of action of repressors that act on these isomerization steps.

One of the best understood and mechanistically interesting repressors of RNA polymerase isomerization is the MerR protein. This protein influences the transcription rate from mercury-dependent promoters by binding within the promoter. MerR binds as a homodimer to a region of dyad symmetry within the merOP region (14, 15). The mer operon contains two overlapping and divergently opposed promoters, designated  $P_R$  and  $P_T$ . MerR is responsible for controlling transcription from both of these promoters. Binding of nonliganded MerR to this region represses transcription from both the structural gene promoter,  $P_T$ , and the divergently oriented merR promoter,  $P_R$ . Upon binding of either liganded or nonliganded MerR protein, the  $P_R$  promoter is repressed, apparently through steric interference of the RNA polymerase, for the MerR binding site overlaps the +1 transcription start site (16, 17). At  $P_T$ , however, the binding of nonliganded MerR to its site between the -35 and -10, partially overlapping the -35 hexamer, results in the recruitment of S70 RNA polymerase holoenzyme but maintains repression (14). Insight into why the merR-RNA polymerase-promoter is transcriptionally inactive can be ascertained by considering what merR does to activate the  $P_T$  promoter. Binding of the allosteric inducer Hg(II) to the MerR-DNA complex results in a structural rearrangement that underwinds the DNA between the -35 and -10 and relaxes the MerR-induced DNA bend (18, 19). This results in [dihedral angles](#) between the -35 and -10 elements that would be similar to those found in a promoter with an optimal 17-bp spacer between these two elements. Additionally, the relaxation of the MerR-induced DNA bend would then place both the -35 and -10 elements in contact with the RNA polymerase, facilitating transcription. So it appears that repression by merR occurs by stabilizing a conformation of the DNA in the RNA polymerase-promoter that is incompatible with open complex formation.

Although the assembly of PICs on eukaryotic promoters occurs in the absence of additional components, efficient transcriptional initiation at most promoters requires the presence of transcriptional activator proteins bound to DNA at sites nearby or at a distance. Hence progression from the relatively inactive PIC to actively transcribing RNA polymerase can be repressed by interfering with the action of these activator proteins.

There are three classes of mechanism by which repressors may interfere with the action of a transcriptional activator. In the first two classes, the repressor prevents DNA binding by the activator. In one case, binding of the repressor to a site that either completely or partially overlaps the binding site of the activator protein competitively inhibits DNA binding by the activator. In another case, the repressor may bind to an activator and prevent its DNA binding. Often repressors of this type interact with activator proteins that have multiple subunits. Usually, the repressor shares homology with one subunit of the activator protein, but is unable to contribute a critical component to DNA recognition. In this type of inhibition, the repressor inhibits DNA binding by functioning as a negative dominant mutant subunit of the activator protein. The third type of repressor “quenches” activator function by preventing transmission of the signal from the activation surface of the DNA-bound activator protein to the promoter.

Repressors that function by competing directly with activator DNA binding are not common in either prokaryotes or eukaryotes. An example of this type of repressor is the *Krüppel* (*Kr*) protein. In *Drosophila*, *Kr* binding sites overlap the binding sites of several activator proteins. One of the better-characterized examples is in the *eve* promoter, where the *Kr* sites overlap the binding sites of the [bicoid](#) protein. *Krüppel* inhibition of bicoid function alters the pattern formation during embryogenesis in flies (20). By contrast to the paucity of repressors that function by competitive DNA binding, there are many repressors that compete for association with one of the activator subunits. Regulation of the entire class of basic helix–loop–helix proteins occurs ubiquitously by this mechanism in eukaryotic developmental programs. For example, during development in several organisms, the ability of myoD to bind DNA and activate genes required for muscle development (21) is completely inhibited by the interaction of this protein with the Id (22) or **Notch** (23) proteins. These repressor proteins contain a region that is identical to the oligomerization domain of myoD, but lack the basic region used by myoD to make contacts with DNA. Therefore an Id-myoD heterodimer is inactive in transcriptional activation because it is incapable of binding DNA.

Repression by quenching of an activator surface requires that both the activator and repressor occupy the DNA at the same time. This phenomenon has been observed in a wide variety of systems and often involves an interaction between two proteins that, on their own, function as transcription activators. For example, YY1 activates transcription of a number of different promoters (24). In the *c-fos* promoter, however, binding of YY1 blocks the **cyclic AMP**-dependent induction of this gene. YY1 induces a bend in the promoter DNA that prevents the interaction of the protein bound at the cAMP [response element](#) with the PIC (25). Repression by quenching does not always require that the repressor protein be a DNA-binding protein. In this type of repression, a bound protein may recruit a “[corepressor](#)” that interferes with the interaction of the activator protein with a coactivator. In one well-studied case, the activation of a set of promoters by DNA-bound Mad-Max heterodimer is repressed by the binding of a member of the mSin3A/B complex. This protein in turn recruits a histone deacetylase that deacetylates the histones at the promoter, which in turn inhibits transcription initiation (26, 27).

### 3. Repressors that Prevent RNA Polymerase Escape from the Promoter

The observation that some repressors function after synthesis of the first phosphodiester bond appears to be antithetical to the conservatism of cellular metabolic regulation. The wisdom of this mode of transcriptional regulation appears, however, upon considering the need of an organism to respond extraordinarily rapidly to a change in environmental conditions. According to this idea, the RNA polymerase negatively regulated post-initiation is “primed” and ready to go in response to the induction signal. Not surprisingly then, two proteins that control the carbon source utilization of *E.*

*coli*, the [Lac repressor](#) and Gal repressors (see [Gal Operon](#)), are thought to exert their repressive effects, at least in part, by this mechanism ([28](#), [29](#)).

One of the best-studied examples repressors of promoter escape is the repression of the early A2c promoter of f29 phage by the p4 repressor. In this system, p4 and RNA polymerase bind cooperatively at the promoter, forming a complex that incorporates nucleoside triphosphates into abortive RNA transcripts. However, the RNA polymerase does not proceed to make full-length transcripts until the level of p4 falls below that needed to occupy its binding site ([30](#)).

Until now, repressors that regulate promoter escape have only been identified in prokaryotic organisms. On the one hand, their absence in eukaryotes may represent a fundamental difference in the time scale in which the two types of cell must respond to the environment to survive. On the other hand, these repressors may exist in eukaryotic organisms, and it may just be a matter of time before we discover yet another example of the commonality of transcriptional regulation between prokaryotes and eukaryotes.

### Bibliography

1. F. Jacob and J. Monod (1961) *J. Mol. Biol.* **3**, 318–356.
2. T. Gardella, H. Moyle, and M. M. Susskind (1989) *J. Mol. Biol.* **206**, 579–590.
3. P. Zuber, J. Healy, H. L. Carter, 3d, S. Cutting, C. P. Moran Jr., and R. Losick (1989) *J. Mol. Biol.* **206**, 605–614.
4. P. L. DeHaseh, M. L. Zupancic, and M. T. Record Jr. (1998) *J. Bacteriol.* **180**, 3019–3025.
5. T. Matsui, J. Segall, P. A. Weil, and R. G. Roeder (1980) *J. Biol. Chem.* **255**, 11992–11996.
6. L. Zawel, and D. Reinberg (1993) *Prog. Nucleic Acid Res. Mol. Biol.* **44**, 68–108.
7. D. K. Hawley, A. D. Johnson, and W. R. McClure (1985) *J. Biol. Chem.* **260**, 8618–8626.
8. A. N. Imbalzano, H. Kwon, M. R. Green, and R. E. Kingston (1994) *Nature* **370**, 481–485.
9. R. T. Kamakaka, M. Bulger, and J. T. Kadonaga (1993) *Genes Dev.* **7**, 1779–1795.
10. G. Felsenfeld (1992) *Nature* **355**, 219–224.
11. D. H. Rivier and J. Rine (1992) *Curr. Opin. Genet. Dev.* **2**, 286–292.
12. T. Goto and M. Monk (1998) *Microbiol. Mol. Biol. Rev.* **62**, 362–378.
13. J. A. Inostroza, F. H. Mermelstein, I. Ha, W. S. Lane, and D. Reinberg (1992) *Cell* **70**, 477–489.
14. A. Heltzel, I. W. Lee, P. A. Totis, and A. O. Summers (1990) *Biochemistry* **29**, 9572–9584.
15. T. V. O'Halloran, B. Frantz, M. K. Shin, D. M. Ralston, and J. G. Wright (1989) *Cell* **56**, 119–129.
16. P. A. Lund, S. J. Ford, and N. L. Brown (1986) *J. Gen. Microbiol.* **132**, 465–480.
17. N. N. Ni'Bhriain, S. Silver, and T. J. Foster (1983) *J. Bacteriol.* **155**, 690–703.
18. A. Z. Ansari, M. L. Chael, and T. V. O'Halloran (1992) *Nature* **355**, 87–89.
19. A. Z. Ansari, J. E. Bradner, and T. V. O'Halloran (1995) *Nature* **374**, 371–375.
20. D. Stanojevic, S. Small, and M. Levine (1991) *Science* **254**, 1385–1387.
21. H. Weintraub, R. W. Davis, S. Tapscott, M. Thayer, M. Krause, R. Benezra, T. K. Blackwell, D. Turner, R. Rupp, and S. Hollenberg (1991) *Science* **251**, 761–766.
22. R. Fairman, R. K. Beran-Steed, S. J. Anthony-Cahill, J. D. Lear, W. F. Stafford, W. F. DeGrado, P. A. Benfield, and S. L. Brenner (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10429–10433.
23. R. Kopan, J. S. Nye, and H. Weintraub (1994) *Development* **120**, 2385–2396.
24. E. Seto, Y. Shi, and T. Shenk (1991) *Nature* **354**, 241–245.
25. S. Natesan and M. Z. Gilman (1993) *Genes Dev.* **7**, 2497–2509.
26. D. E. Ayer, C. D. Laherty, Q. A. Lawrence, A. P. Armstrong, and R. N. Eisenman (1996) *Mol. Cell Biol.* **16**, 5772–5781.

27. P. J. Hurlin, C. Queva, P. J. Koskinen, E. Steingrimsson, D. E. Ayer, N. G. Copeland, N. A. Jenkins, and R. N. Eisenman (1996) *EMBO J.* **15**, 2030.
28. H. E. Choy and S. Adhya (1993) *Proc. Natl. Acad. Sci. USA* **90**, 472–476.
29. J. Lee and A. Goldfarb (1991) *Cell* **66**, 793–798.
30. F. Rojo, M. Mencia, M. Monsalve, and M. Salas (1998) *Prog. Nucleic Acid. Res. Mol. Biol.* **60**, 29–46.

### Suggestions for Further Reading

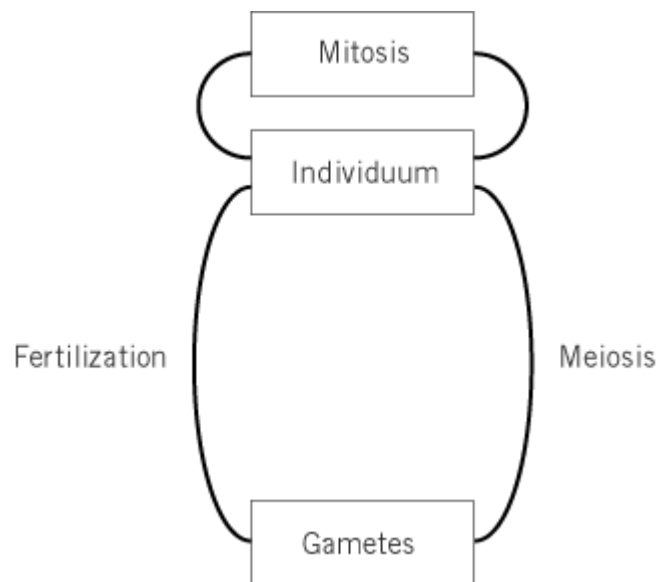
31. M. T. Record, et al. (1996) *E. coli and Salmonella: Cellular and Molecular Biology*, American Society for Microbiology, Washington, D.C., p. 792. (A comprehensive review of transcription initiation pathways.)
32. A. Hochschild and S. L. Dove (1998) *Cell* **92**, 597–600. (A state-of-the-art review of activation and repression.)
33. A. D. Johnson (1995) *Cell* **81**, 655–658. (The price of repression.)

## Reproduction, Animal

Reproduction is the process of genetic transfer to the next generation. In animals, reproduction involves two [sexes](#), **meiosis** and multicellularity ([1](#)). Thus, reproduction in animals is different from reproduction in **prokaryotes**, such as **bacteria**. A further characteristic of animal reproduction is that animals are mobile and therefore able to spread their own **genes** in a much more efficient way than, for example, **plants** can do. Another characteristic that is related to this subject is the necessity for sexual reproduction. Unlike plants or bacteria, animals have no vegetative reproduction pathways. The rare **parthenogenic** events that occur in some insects like bees and sometimes in frogs are very special exceptions and cannot be considered to be the rule. Also, apart from some fishes (which can be [polyploid](#)) or certain stages after reproductive events in some insects, most animals tend to be **diploid** individuals. During gametogenesis, the diploid genotype has to be reduced in order to provide **germ cells** that are [haploid](#); these events, therefore, involve both **mitotic** and meiotic divisions. The joining of the maternal and paternal nuclei during [karyogamy](#) results in the creation of a diploid individual (Fig. [1](#)).

**Figure 1.** Diagrammed relationship between the individual, gametes, meiosis, mitosis, and fertilization. As a rule, mammals are diploid, whereas their gametes are reduced during meiosis to a haploid stage. After a male and a female gamete join during fertilization, a diploid individual is generated. Cells other than the germ cells of the individual divide mitotically.





In animals, there is a strict division between **somatic** and **germ line cells**. Even after just the first stages of cell division following [fertilization](#), germ cells are determined, and they migrate to the primordial urogenital tract, where they develop according to their genetic constitution into either spermatogonia or oogonia. Somatic cells are always diploid, apart from the **G2** and **M phase** of the [cell cycle](#), in which the [genome](#) is temporarily **tetraploid**. There are no certain stages in animals and in the animal reproductive cycle as there are in plants during which the whole individual undergoes a haploid or a diploid phase. Although there is a strict separation of somatic and germ cells, neither cell type acts completely independently of the other. The **phenotype**, ie, the morphology and the ethology of animals, is most certainly determined by actions that are regulated by the germ cells. For instance, puberty causes animals to change their phenotype. In frogs, maturation can easily be seen in the metamorphosis of tadpoles to mature frogs. Most of these developmental processes are, in terms of animal reproduction, determined by the action of [hormones](#). For this reason, it is not surprising that a number of hormone **receptors** are present, not only in organs that are directly related to reproductive events but also in organs and somatic tissues in which no direct link with reproductive activity is presupposed. This is, for example, the case with the [estrogen receptor](#), which is expressed in a number of tissues.

Animal reproduction is not restricted to the events of the production of germ cells, reduction to a haploid stage, fertilization, and early development of the embryo. Regarding certain aspects of reproduction, some [homology](#) even to reproduction in plants can probably be found. A very important and unique aspect of reproduction in animals, however, is the event of giving birth. Interestingly, there are vast differences among species. The predominant one can be seen when comparing animals in which extracorporal fertilization is the rule (as in spawning frogs) with more sophisticated species, like chordates. In extracorporal fertilization, only the release of the mature gametes into a liquid medium (in most cases, water) has to be regulated. From then on, gamete interaction is directed by [protein-protein interactions](#), as genetic control by the [nucleus](#) has not been observed. **Chemotactic** attraction then directs the motile gamete, which is normally that of the male, to the **oocyte**. In species that give birth, not only the gametogenesis, behavioral, and ethological aspects have to be considered to achieve reproductive success, but the hormonal machinery to induce labor, milk secretion, and changes in behavior are also important. This is of sufficient interest to elucidate the molecular basis for the development of these metabolic events.

Animal reproduction also comprises the events that have to take place after birth of the offspring, known as brood care. Although brood care is an important element in most of the primates, in other species the intensity of the brood care varies. In species with extracorporal fertilization, the intensity

of brood care can be extremely low, as can be seen in frogs or in sea urchins. In sticklebacks, the brood care involves at least the construction of a nest, whereas, in birds, brood care also entails feeding the offspring. An interesting exception to this is the cuckoo, which has developed a mechanism that allows it to abuse the brood care of other singing birds. In higher species, brood care also involves the phase of imprinting and teaching social behavior. This is extremely important in terms of animal reproduction because it can be shown that artificially promoted faulty imprinting can lead to animals that have lost the drive to mate. The brood care does not occur in plants or in bacteria and is therefore a characteristic of reproduction in animals, especially in higher mammals with a sophisticated social structure.

## Bibliography

1. Ruvinsky A (1997) *Acta Biotheor* **45**, 127–141.

## Response Element

A response element is a short DNA sequence, usually between 5 and 20 base pairs in length, that is found in all **genes** that show a particular pattern of **gene expression**. Moreover, it can confer that pattern of expression on another gene when artificially linked to it.

Different genes have a wide variety of different expression patterns, with particular genes being expressed in specific tissues of the body (see [Development](#)), while others are activated in response to particular [hormones](#) or other inducers. Such regulation is primarily produced at the level of [transcription](#) by controlling which specific genes are transcribed into RNA in a particular tissue or in response to a particular stimulus; all the other stages of gene expression, such as [RNA splicing](#), transport to the cytoplasm, and [translation](#) into protein follow more or less automatically. This conclusion led to a detailed study of the DNA sequences of individual genes in the hope of identifying elements that might confer the response to a particular stimulus or produce a particular tissue-specific pattern of transcription.

Such studies at the DNA sequence level led to the identification of a number of short response elements that were found in the regulatory regions of all genes that responded to a particular stimulus, but not in genes that did not respond. A selection of such response elements is listed in Table 1(1), and several of these are discussed in individual entries (see [Glucocorticoid Response Element](#), [Hormone Response Elements](#), [Metal Response Element](#), [Serum Response Element](#), and [Sterol Response Element](#)). The purpose of this article is to provide a more general overview of the nature of response elements and the manner in which they function.

**Table 1. Sequences that Confer Response to a Particular Stimulus**

| Consensus Sequences | Response to | Protein Factor           | Gene Containing Sequences |
|---------------------|-------------|--------------------------|---------------------------|
| CTNGAATNTT<br>CTAGA | Heat        | Heat-shock transcription | hsp70, hsp83, hsp27, etc. |

|                         |                                     |                                 |  |
|-------------------------|-------------------------------------|---------------------------------|--|
| T/G T/A<br>CGTCA        | Cyclic AMP                          | factor<br>CREB/ATF              | Somatostatin<br>fibronectin, a-<br>gonadotrophin c-fos,<br>hsp70         |
| TGAGTCAG                | Phorbol esters                      | AP1                             | Metallothionein IIA<br>a <sub>1</sub> -antitrypsin<br>collagenase        |
| CC(A/T) <sub>6</sub> GG | Growth factor in<br>serum           | Serum response<br>factor        | c-fos, <i>Xenopus</i> g-actin  |
| RGRACNNN                | Glucocorticoid                      | Glucocorticoid<br>receptors     | Metallothionein IIA<br>tryptophan<br>oxygenase,<br>uteroglobin, lysozyme |
| TGTYCY<br>RGGTCANNN     | Estrogen                            | Estrogen<br>receptor            | Ovalbumin,<br>conalbumin,<br>vitellogenin                                |
| TGACCY<br>RGGTCAT       | Thyroid<br>hormone<br>retinoic acid | Thyroid<br>hormone<br>receptors | Growth hormone<br>myosin heavy chain                                     |
| GACCY<br>TGCGCCCGCC     | Heavy metals                        | Mep-1                           | Metallothionein genes  |
| AGTTTCNN                | Interferon-a                        | Stat-1 Stat-2                   | Oligo A synthetase<br>guanylate-binding<br>protein                       |
| TTTCNC/T<br>TTNCNNNAA   | Interferon-g                        | Stat-1                          | Guanylate-binding<br>protein, Fc g receptor <sup>a</sup>                 |

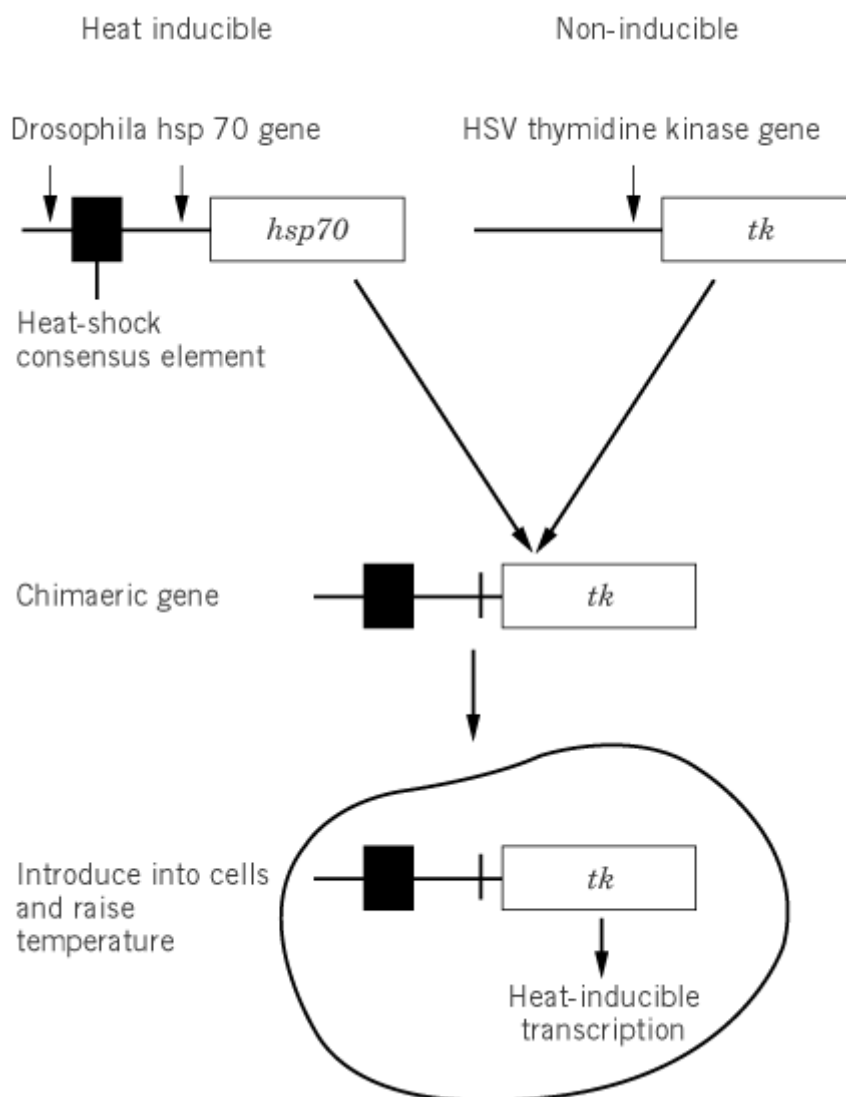
<sup>a</sup> N, any base; R, purine; Y, pyrimidine.

Although the presence of a specific short DNA sequence in all genes that respond to a particular stimulus strongly suggests that this element confers the response to that stimulus, it is necessary to prove that this is the case. In order to do this, the putative response element must be transferred to a gene that is not normally responsive to the stimulus and consequently shown to respond to it.

One of the earliest examples where this was achieved concerns the heat shock element (HSE), which is found in the regulatory regions of genes whose expression is increased in response to elevated temperature. Thus, for example, the gene encoding the 70-kDa heat shock protein (hsp70) contains an HSE within the gene **promoter** that is located approximately 90 bases upstream of the transcriptional start site. Pelham (2) took this HSE and linked it to the non-heat-inducible **thymidine kinase** (tk) gene of the eukaryotic **herpes** simplex virus. He then introduced this hybrid gene into mammalian cells and observed that the hybrid gene could be activated by elevated temperature, leading to increased thymidine kinase production, even though the tk gene was not normally heat inducible (Fig. 1). Hence the HSE constitutes a bona fide response element that can confer the response to heat upon another gene that is not normally heat-inducible. Such a direct demonstration

of the role of a particular element in producing a particular response has now been provided by similar methods for the other sequences that are listed in Table 1. Thus, for example, linking a glucocorticoid response element to another gene renders that gene responsive to glucocorticoid even though that gene normally does not respond in this manner, and so on.

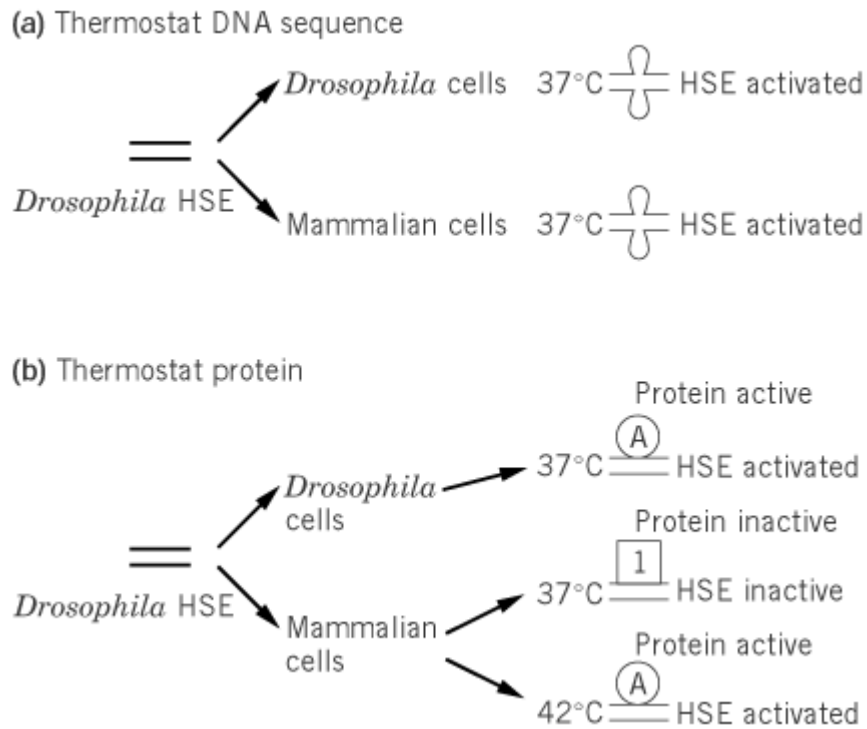
**Figure 1.** Demonstration that the heat shock element mediates heat inducibility. Transfer of this sequence to the thymidine kinase gene, which is not normally inducible by heat, renders this gene heat-inducibile.



The above example concerning the HSE also illustrates another aspect of response elements. Thus, the HSE used by Pelham was taken from the hsp70 gene of the fruit fly *Drosophila melanogaster*, but the hybrid gene was introduced into mammalian cells. The successful functioning of the fly gene in mammalian cells not only indicates that this process is evolutionarily conserved but permits a further conclusion about the way in which the effect operates. Thus in the cold-blooded *Drosophila*, 37°C would represent a thermally stressful temperature at which the heat shock response would normally be active. Thus if the HSE acts as a temperature sensing thermostat (Fig. 2a), it would be set to go off at 37°C rather than at 42°C, which is the thermally stressful temperature for mammalian cells. In fact, however, the hybrid gene in this experiment was inactive at 37°C in the mammalian cells and was induced only at 42°C. Hence this response element sequence does not act itself as a thermostat but rather must act by being recognized by a cellular protein that is activated only at an

increased temperature characteristic of the mammalian cell heat shock response (Fig. 2b).

**Figure 2.** Predicted effects of placing the *Drosophila* heat shock element (HSE) in a mammalian cell if the element acts as a thermostat detecting increased temperature directly (a) or if it acts by binding a protein that is activated by raised temperature (b). Note that only possibility (b) can account for the observation that the *Drosophila* HSE activates transcription in mammalian cells only at the mammalian heat shock temperature of 42°C and not at the *Drosophila* heat shock temperature of 37°C.



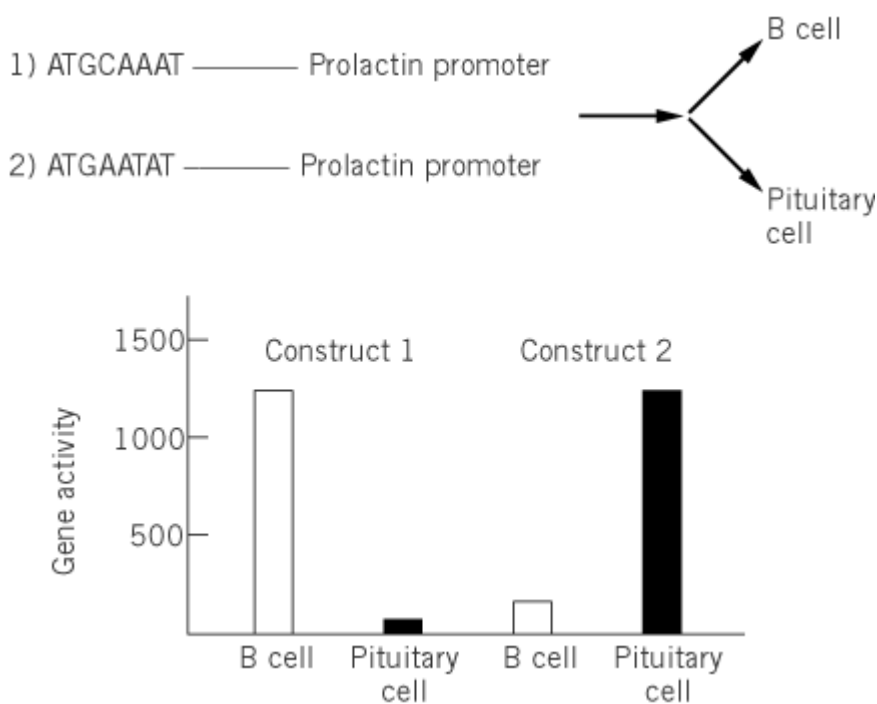
This experiment thus indicates that response elements act by binding specific proteins known as [transcription factors](#). These transcription factors are specifically synthesized or activated post-translationally in response to a specific stimulus and then activate gene transcription via their specific response element. Hence the specificity provided by different response elements producing different patterns of gene activation arises because they bind different transcription factors, which become active in response to these stimuli.

The nature of the transcription factors that bind to each of the response elements listed is shown in Table 1. In the case of the HSE, it achieves its effect by binding a protein known as the heat shock factor (HSF). In response to elevated temperature, HSF undergoes a transition from an inactive monomeric form to a trimeric form that can bind to the HSE and activate transcription of the heat shock genes (3).

Although response elements were first identified in genes that respond to specific stimuli, similar DNA sequence elements in the promoters of tissue-specific genes also play a critical role in producing their tissue-specific pattern of gene expression, by binding factors that are present in an active form only in a specific tissue. Thus, for example, the promoter of the [immunoglobulin](#) heavy- and light-chain genes contain a response element known as the octamer motif (ATGCAAAT) (4) that is important in their specific expression in B lymphocytes and not in other cell types. The octamer motif can confer B-cell-specific expression on an unrelated promoter (5), indicating its critical role in this process. Interestingly, the related sequence ATGAATAA/T is found in the genes for prolactin and for [growth hormone](#), which are expressed specifically in the anterior pituitary gland

(6). If this short sequence is inserted upstream of a promoter, the gene is expressed only in pituitary cells, indicating its critical importance in producing pituitary-specific gene expression. In contrast, if this short sequence is replaced with the octamer motif, which differs by only 2 bases, it will direct expression of the same promoter only in B cells (7) (Fig. 3).

**Figure 3.** Linkage of the octamer motif ATGCAAAT (1) and the related motif ATGAATAT (2) to the prolactin promoter and introduction into B cells and pituitary cells (a). Note that the octamer-containing construct 1 directs a high level of activity in B cells but not in pituitary cells, whereas construct 2 directs a high level of gene activity in pituitary cells but not in B cells (b).



Thus, specific short DNA sequences can produce either tissue-specific patterns of transcription or result in a gene responding to a particular stimulus. Such sequences play a central role in the regulation of gene transcription acting by binding specific regulatory transcription factors.

#### Bibliography

1. E. H. Davidson, H. T. Jacobs, and R. J. Britten, (1983) *Nature* **301**, 468–470.
2. H. R. B. Pelham (1982) *Cell* **30**, 517–528.
3. R. Morimoto (1993) *Science* **259**, 1409–1410.
4. J. O. Mason, G. T. Williams, and M. S. Neuberger (1985) *Cell* **41**, 479–487.
5. T. Wirth, L. Staudt, and D. Baltimore (1987) *Nature* **329**, 176–178.
6. C. R. Nelson, V. R. Albert, H. P. Elsholtz, L. E. W. Lu, and M. G. Rosenfeld (1988) *Science* **239**, 1400–1405.
7. H. P. Elsholtz, V. R. Albert, M. N. Treacy, and M. G. Rosenfeld (1990) *Genes Dev.* **4**, 43–51.

#### Suggestions for Further Reading

8. D. S. Latchman (1998). *Gene Regulation: A Eukaryotic Perspective*, 3rd ed., Stanley Thornes, Cheltenham (especially Chapter 7).
9. D. S. Latchman (1998). *Eukaryotic Transcription Factors*, 3rd ed. Academic Press, London, San Diego (especially Chapter 1).

## Restriction Enzymes

*Restriction enzyme* refers here to the bacterial type II restriction **endonucleases**—one type of a general class of enzymes that function in restriction–modification systems (see [Restriction–Modification Systems](#) for other types of restriction enzymes). Type II restriction enzymes are ubiquitous and simple in composition. They are used extensively in [recombinant DNA](#) work and, in fact, made the technology possible. The restriction enzymes are  $Mg^{2+}$ -requiring, homodimeric proteins that cleave at defined phosphodiester bonds within or adjacent to specific sequences of 4–8 base pairs (the “restriction site”) of double-stranded DNA. The phosphodiester bond cleavage of each DNA strand produces DNA fragments with 5' phosphate and 3' hydroxyl termini. The scission of both strands results in either staggered-end or blunt-end double-strand DNA fragments (see [Staggered Cut](#)). Depending on the endonuclease used, the termini of the duplex DNA fragments may be completely duplex (blunt) or have protruding 5' or 3' single strands (otherwise known as “staggered”, “sticky”, or cohesive ends). Almost all type II restriction enzymes cut at DNA sequences that are **palindromic**, that is, have a dyad axis of symmetry, so that their complementary strand has the same sequence. Examples are the target sites for *EcoRI* (GAATTC) and *FinII* (CCGG). Some restriction enzymes classified as type II are monomeric enzymes and cut at defined distances outside specific asymmetric recognition sequences.

Approximately 2600 different restriction enzymes are known, with over 230 cleavage specificities (1). Bacteria synthesize minute amounts of the enzymes to serve their biological function (see [Restriction–Modification Systems](#)), but the overexpression of cloned restriction enzyme genes substantially facilitates their preparation and purification. The cognate modification methyltransferase partner of the restriction enzyme prevents endonucleolytic cleavage by methylating specific adenosine or cytosine nucleotides in the recognition sequence (see [Methyltransferase, DNA](#)), resulting in a sequence that is refractory to endonuclease hydrolysis.

Under optimized reaction conditions and presumably *in vivo*, restriction enzymes are highly selective for their canonical restriction sites. Changes in restriction enzyme buffer conditions, however, such as pH, ionic strength, divalent metal ion substitutions ( $Ca^{2+}$  or  $Mn^{2+}$  for  $Mg^{2+}$ ), or the addition of organic solvents, can induce inaccurate cutting ([star activity](#)). Different restriction enzymes that recognize and cut at the same position in identical recognition sites are referred to as [isoschizomers](#). For example, the *EcoRI* and *RsrI* restriction enzymes both recognize and cut between the G and A of the canonical sequence, GAATTC. Other restriction enzymes, termed *neoschizomers*, recognize identical DNA sequences but cut at different positions within the sequence. For instance, *Asp718I* (GGTACC) and *KpnI* (GGTACC) cleave between the two Gs and between the two Cs, respectively.

### 1. Protein Structures

Little amino acid similarity (15–25% identity) exists among the more than 60 known sequences of restriction enzymes (2), indicative of a diversity of target sites and of mechanisms of recognition and catalysis. However, [X-ray crystallography](#) structures of type II restriction enzymes that have little or no amino acid sequence [homology](#) (*EcoRI*, *BamHI*, *EcoRV*, *FokI* (a type II enzyme), *PvuII* and *Cfr10I*), reveal a common structural motif consisting of a five-stranded **b-sheet** with flanking **α-helices** containing the [active site](#) (3, 4).

Which bonds in the DNA are cleaved is determined by the spatial positioning of the restriction

enzyme around the DNA double helix. Cleavable phosphodiester bonds generating 5' or 3' overhangs are positioned in the major or minor groove, respectively, whereas cleavable bonds generating blunt cuts are on opposite sides of the helix and near the minor groove (5). Accordingly, restriction enzymes that generate 5' overhangs, such as *EcoRI* and *BamHI*, dimerize to align critical catalytic residues adjacent spatially to cleavage sites separated by 4 base pairs, with the enzyme facing the major groove. Blunt-end cutters such as *EcoRV* and *PvuII* dimerize so that their active sites are located on opposite sides of the DNA helix, with the enzymes facing the minor groove. A restriction enzyme generating a 3' overhang has yet to be crystallized with DNA. The orientation with respect to the grooves of the DNA may also dictate which motifs are used in recognition. *EcoRI* and *BamHI* use a **four-helix bundle** to make base contacts in the major groove, whereas *EcoRV* and *PvuII* use loops or strands in the major and minor grooves (6).

Interestingly, the monomeric type II enzyme, *FokI*, which recognizes a pentameric, asymmetric sequence, and cuts a short distance away to leave a 5' overhang, has a bilobed structure with the catalytic site in one **domain** and the target recognition site in the other. Also, *FokI* has an active site that is similar to *BamHI*, uses  $\alpha$ -helices for DNA recognition, and has been modeled to cleave in the major groove (4).

Crystal structures of complexes of the enzymes with their target DNA reveal that some restriction enzymes may utilize DNA distortion and bending in their recognition or catalysis. Cocrystals of DNA with *EcoRI* and *EcoRV* have the DNA bent and unwound, whereas cocrystals of *BamHI* and *PvuII* contain regular, B-form DNA. Crystal structures will continue to prove valuable for understanding the recognition and catalytic properties of restriction enzymes.

## 2. DNA Recognition

For the type II restriction enzymes, the dimer first binds DNA nonspecifically and then transverses the DNA by random linear [diffusion](#), apparently pausing at sites resembling the true restriction site, but continuing until the correct recognition sequence is reached. Linear diffusion, rather than three-dimensional diffusion with random collision, speeds the location of a restriction site by reducing the search to one dimension (7). Restriction enzymes such as *EcoRI* and *BamHI* bind much more tightly to their recognition sites than to nonspecific DNA, whereas others such as *EcoRV* and *TaqI* bind to nonspecific or target DNA with equal affinities. These latter restriction enzymes achieve specific sequence recognition and cleavage on binding  $Mg^{2+}$  (8). The differential target recognition properties of restriction enzymes indicate that specificity discrimination may occur at the level of both the initial enzyme–DNA complex and the correct [transition state](#) during catalysis. Specific DNA binding by restriction endonucleases is dependent on a network of [electrostatic interactions](#) and [hydrogen bonds](#) involving both side-chain and main-chain atoms of the enzyme, plus [water](#) molecules, with bases and the sugar–phosphate backbone of the DNA. **Hydrophobic** interactions also contribute importantly to DNA binding. All the interactions, along with deformation of the DNA, promote the precise geometry required for successful catalysis. Suboptimal reaction conditions (see [Star Activity](#)) or methylation of DNA hinders or alters the structural geometry needed for recognition specificity and catalysis by restriction enzymes.

## 3. Catalysis

Depending on the restriction enzyme and the reaction conditions, cleavage of both strands of the recognition sequence occurs in either a sequential or concerted manner (9, 10). Scission of the duplex DNA recognition sequence probably occurs by direct nucleophilic attack of water on the phosphorus atom of the phosphodiester bond, rather than through a covalent enzyme intermediate, since there is inversion of configuration at the scissile phosphate bond (9). The mechanism by which the water is activated to form the nucleophilic OH group has yet to be determined. Several possible ligands in the active site could contribute to removing a proton from water, the most obvious candidate is the catalytically essential  $Mg^{2+}$  ion. Acidic catalytic site residues also may serve as



general bases (11) to abstract a proton: Asp90 in *EcoRV*, Glu111 in *EcoRI*, or Glu113 in *Bam HI* (6, 11, 12). It has also been suggested that the negatively charged oxygen atom of the neighboring phosphodiester group may serve in proton abstraction (13). Whereas some restriction enzymes, such as *EcoRI*, probably use one  $Mg^{2+}$  ion for catalysis, *EcoRV* apparently needs two; one  $Mg^{2+}$  is proposed to generate a catalytically competent complex and the other to activate water for hydrolysis (12, 14). Cleavage by restriction enzymes is independent of the **supercoiling** of the DNA but dependent on the length and sequence context of flanking sequences surrounding the restriction site (10). With turnover numbers in the  $min^{-1}$  (reciprocal minute) range, restriction enzymes are catalytically slow enzymes, and the product release is the rate-limiting step (9).

#### 4. Applications

Numerous restriction enzymes with a large variety of sequence specificities are commercially available, and they provide an indispensable tool for **cloning** and recombinant DNA work. The discovery of restriction enzymes allowed the manipulation of large DNA, which has led to a greater understanding of its structure and information content. Restriction enzymes also provide valuable opportunities for studying DNA–protein interactions and enzyme **kinetics**. Several applications for the manipulation and analysis of DNA using restriction enzymes are discussed in the **Restriction Fragment**, **Staggered Cut**, and **Restriction Map** articles.

#### Bibliography

1. G. G. Wilson (1996) personal communication.
2. A. Jeltsch, M. Kroger, and A. Pingould (1995) *Gene* **160**, 7–16.
3. D. Bozic, S. Grazulis, V. Siksnys, and R. Huber (1996) *J. Mol. Biol.* **255**, 176–186.
4. D. A. Wah, J. A. Hirsh, L. F. Dorner, I. Schildkraut, and A. K. Aggarwal (1997) *Nature* **388**, 97–100.
5. J. Anderson (1993) *Curr. Opin. Struct. Biol.* **3**, 24–30.
6. M. Newman, T. Strzelecka, L. F. Dorner, I. Schildkraut, and A. K. Aggarwal (1995) *Science* **269**, 656–663.
7. A. Jeltsch, J. Alves, H. Wolfes, G. Maass, and A. Pingoud (1994) *Biochemistry* **33**, 10215–10219.
8. J. Heitman (1993) in *Genetic Engineering*, Vol. **15**, J. K. Setlow, ed., Plenum Press, New York, pp. 57–107.
9. A. Pingould, J. Alves, and R. Geiger (1993) in *Methods in Molecular Biology*, Vol. **16**, *Enzymes of Molecular Biology*, M. M. Burrell, ed., Humana Press, Inc., Totowa, N. J. p. 167.
10. R. J. Roberts and S. E. Halford. (1993) in *Nucleases*, 2nd ed., S. M. Linn, R. S. Lloyd, and R. J. Roberts, eds., Cold Spring Harbor Laboratory Press, New York, pp. 35–88.
11. A. Jeltsch, J. Alves, G. Maass, and A. Pingoud (1992) *FEBS Lett.* **304**, 4–8.
12. D. Kostrewa, and F. K. Winkler (1995) *Biochemistry* **34**, 683–696.
13. A. Jeltch, M. Pleckaityte, U. Selent, H. Wolfes, V. Siksnys, and A. Pingould (1995) *Gene* **157**, 157–162.
14. G. S. Baldwin, I. B. Vipond, and S. E. Halford (1995) *Biochemistry* **34**, 705–714.

#### Suggestions for Further Reading

15. H.-Y. Eun (1996) "Restriction endonucleases and modification methylases, in" *Enzymology Prime for Recombinant DNA*, Academic Press, San Diego, pp. 233–306.
16. A. Pingould, J. Alves, and R. Geiger (1993) "Restriction enzymes, in" *Methods in Molecular Biology*, Vol. **16**, *Enzymes of Molecular Biology*, M. M. Burrell, ed., Humana Press, Inc., Totowa, N. J. p. 167.

## Restriction Fragment

Restriction fragments are discrete pieces of double-stranded DNA having unique termini that result from the cleavage of larger DNA by restriction endonucleases, enzymes that recognize specific sequences (usually 4–8 bp in length, in duplex DNA and catalyze the cleavage of both strands within or adjacent to the sequence to generate **cohesive ends** or **blunt ends** (see [Restriction Enzymes](#) and [Staggered Cut](#)). The discovery of restriction endonucleases allowed the manageable analysis and manipulation of an otherwise huge, experimentally inaccessible macromolecule. The use of restriction endonucleases also led to the development of engineered DNA **vectors** (**plasmids**, **cosmids**, **bacteriophages**, and *yeast artificial chromosomes* (YACs)), which are self-replicating DNAs that are used as carriers to clone and characterize restriction fragments. Important features of these vectors include selectable phenotypic markers (such as [antibiotic resistance](#) genes), unique restriction sites grouped in *multiple cloning sites*, where restriction fragments can be introduced, **promoters** for the expression of the inserted DNA, and mechanisms for [DNA replication](#). Vector or chromosomal DNA can be digested into restriction fragments with one or more restriction endonucleases and analyzed for size and fragment distribution using [gel electrophoresis](#) (see [Restriction Map](#)). Restriction-fragment size is dependent on the type of restriction endonuclease used and the particular sequence. For instance, assuming a random distribution bases in DNA with a 50% (G + C) composition, a given tetranucleotide sequence will occur randomly approximately every 256 base pairs, whereas a hexanucleotide sequence occurs approximately every 4096 base pairs. Large restriction fragments can be generated by (1) rare-cutting enzymes, like the octanucleotide (GC<sup>^</sup>GGCCGC, where “<sup>^</sup>” indicates the hydrolysis site in the DNA recognition sequence) cutter R•NotI (where “R•” designates restriction endonuclease), (2) partial R•endonuclease digestion resulting from short digestion duration or dilute enzyme concentrations, or (3) limited methylation of R•endonuclease sites with the cognate [methyltransferase](#) enzyme. The properly sized spectrum of generated restriction fragments can be covalently **ligated** into [cloning](#) vectors to form a collection of chimeric DNA molecules that can replicate as individual **clones** once introduced individually into a host organism, such as *Escherichia coli*, by transformation, [transfection](#), or transduction. Such a population of clone-containing fragments encompassing all of the original DNA is called a “**gene library**.” The genetic information encoded on an individual restriction fragment or library of fragments is easily amplified and manipulated.

### 1. Methods for Analyzing Restriction Fragments

Several methods are used to analyze restriction fragments separately or in a collection of clones. In [Southern blotting](#), restriction fragments separated on an [agarose](#) gel by electrophoresis are processed by transferring the DNA from the gel to a [nitrocellulose](#) filter, which provides a stable support for **hybridization** with single-strand labeled probes containing known sequences of interest. If a restriction fragment contains the complementary sequence, the probe will hybridize, and subsequent analysis of the filter for the label will identify the fragment. It can then be isolated and further characterized by sequencing. Another related method, called colony blotting, involves hybridizing DNA probes to colonies containing recombinant DNA that have been fixed to nitrocellulose filters and lysed. Analysis of the filter can then identify clones carrying sequences of interest.

**Chromosome walking** uses a [genomic library](#) of clones to determine consecutive DNA fragments for mapping or locating a particular region in the genome. In this technique, a cloned DNA fragment is isolated, and the ends of the insert fragment are subcloned and used as probes to find clones containing the same DNA sequences, specifically, DNA adjacent to that generating the probes. Variations in DNA sequences between individual organisms are revealed by digesting DNA with various restriction endonucleases to reveal *restriction fragment length polymorphisms* (RFLP).

RFLPs result from changes in the fragmentation pattern of a restriction digest due to DNA sequence differences, which can be as small as single-base-pair substitutions. RFLP mapping provides a “genetic fingerprint” of the DNA of organisms. DNA fragments containing regulatory sequences, such as, **promoters** or **operators**, can be analyzed by placing them in vectors containing adjacent reporter genes that encode a protein with a measurable activity. Expression of the reporter gene is thereby placed under the control of the regulatory sequence. For instance, an antibiotic resistance gene can be used to report promoter activity. Only clones carrying the expressed antibiotic resistance gene will grow on media containing the antibiotic. Another approach uses enzymatic or functional assays of cell lysates of clones to identify those carrying a DNA fragment containing a gene that encodes the protein of interest.

The generation of restriction fragments by restriction endonucleases has contributed greatly to biotechnology. Fragmentation and sequencing of DNA allows characterization of genes, as well as probing for similar genes in different organisms using hybridization techniques. Sequencing of restriction fragments also reveals the **open reading frames** (ORF) of new genes, and their spatial organization with respect to other genes. RFLP mapping is important in forensic and medical analysis. Restriction fragments provide a starting point for biochemists and molecular biologists to **mutagenize** genes to determine how resultant changes are reflected in protein structure or function. Specially designed cloning vectors containing powerful promoters can overexpress cloned genes of interest to increase protein production. Molecular cloning of restriction fragments is an invaluable tool that has yielded a wealth of information with widespread applications.

## **Restriction Fragment Length Polymorphism (RFLP)**

RFLP is a simple method of detecting variability in the sequences of **DNA** within the [genomes](#) of a population of individuals (ie, **polymorphism**) using [mutations](#) at the sites of action of [restriction enzymes](#) (1, 2). Restriction enzymes are **endonucleases** that cut DNA at specific, very short base sequences, producing [restriction fragments](#) whose average size reflects the frequency with which the specific restriction sequence is found along the DNA. Point mutations may cause restriction sites to disappear, may cause the restriction fragment affected to increase in size by addition of the adjacent one, or may cause new restriction sites to appear, splitting one restriction fragment into two smaller fragments. By separating the resulting restriction fragments by [gel electrophoresis](#) in agar, on the basis of their lengths, the patterns of DNA bands can be compared in different individuals. This technique is limited by its ability to detect only mutation events occurring within enzyme restriction sites, which are sequences of only 4 to 6 nucleotides separated by hundreds or thousands of bases. This limitation can be overcome to some extent by using several restriction enzymes; however, in a study of vertebrate mitochondrial DNA using 20 restriction enzymes that are specific for sequences of six bases, only about 3% of the genome will be sampled (3).

Consequently, RFLP is not a very powerful method for analyzing polymorphisms. If applied to a random selection of DNA, it will reveal, approximately, a level of polymorphism comparable to that found with protein **isozymes**. An improvement is the use in **Southern blotting** of specific probes revealing fragments belonging to a restricted region whose polymorphism is under study. In the case of a species whose genome is practically unknown, but for which a large number of probes is available, randomly cloned genome fragments may generate RFLP patterns of some interest. It is best, however, to isolate the section of DNA of interest, expand it by the polymerase chain reaction (**PCR**), subject it to the subsequent action of several restriction enzymes, and analyze the pattern of bands obtained by electrophoresis. Given the very large number of restriction enzymes that are available now, it will often be possible to design simple RFLP tests to reveal the alternatives.

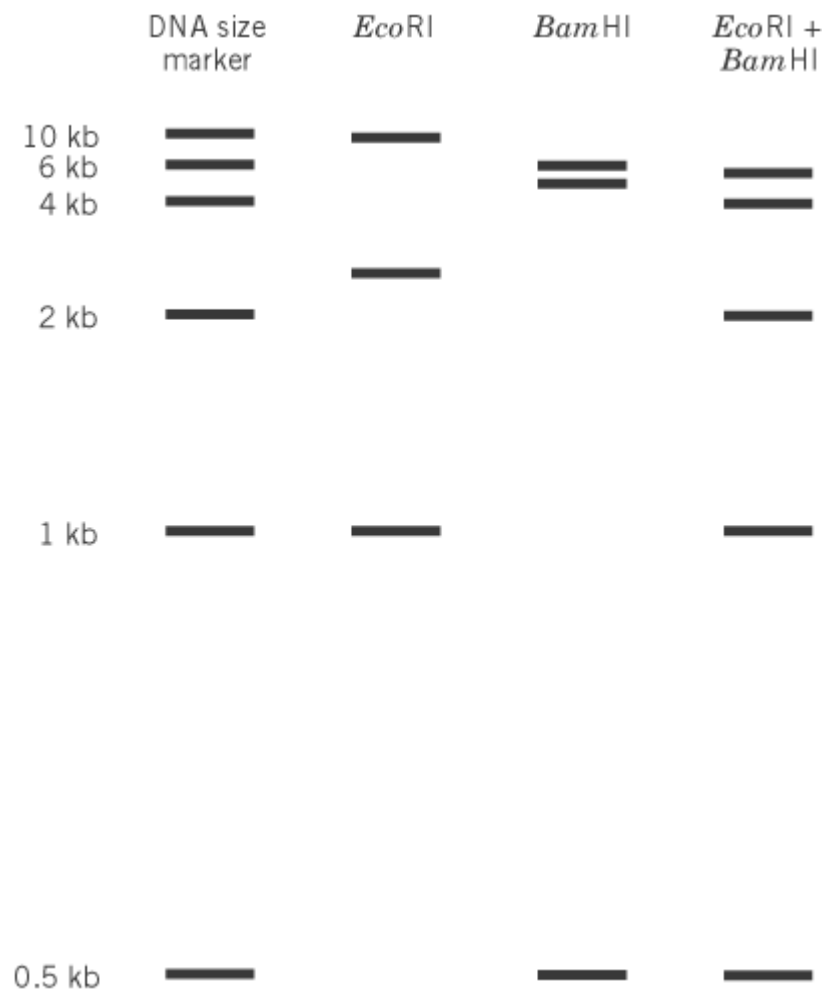
## Bibliography

1. W. M. Brown (1980) *Proc. Natl. Acad. Sci. USA* **77**, 3605–3609.
2. T. E. Dowling, C. Moritz, J. D. Palmer, and L. H. Rieseberg (1996) in *Molecular Systematics* (D. M. Hillis, C. Moritz, and B. K. Mable, eds.), Sinauer Associates, Sunderland, Massachusetts, pp. 249–320.
3. E. Bermingham, G. Seutin, and R. E. Ricklefs (1996) in *Molecular Genetic Approaches in Conservation* (T. B. Smith and R. K. Wayne, eds.), Oxford University Press, NY, pp. 104–124.

## Restriction Map

Restriction mapping is a technique for determining the linear order of restriction sites (see [Restriction Enzymes](#)) along a DNA molecule. The [restriction fragments](#) that are generated by digestion of DNA with different restriction enzymes, alone or in combination, and separated by [gel electrophoresis](#), are used to deduce a restriction map. The resulting map gives the ordered arrangement of the restriction fragments. The choice of restriction enzyme, whether a rare or frequent cutter, depends upon the size of the DNA to be mapped and the frequency of its sites. For instance, rare-cutting restriction enzymes were used for initial mapping of large [genomes](#) like those of *Escherichia coli* and *Homo sapiens*, containing  $4.2 \times 10^6$  base pairs and  $3.3 \times 10^9$ , base pairs respectively. The resulting large restriction fragments isolated were further mapped by digestion with more frequently cutting enzymes. The fragments are separated by electrophoresis according to the inverse of their molecular weights, specifically, their lengths. Fragment size is determined by comparing the sample with a collection of standards. To illustrate, a linear 13-kb piece of DNA is mapped using two enzymes, *EcoRI* and *BamHI*. A diagram of the electrophoretic separation of the digests in an [agarose](#) gel is given in [Figure 1](#). The DNA was digested with *EcoRI*, *BamHI*, and both enzymes. The sizes of fragments produced by single-enzyme digests indicate the number of restriction sites for each enzyme and the distances between them, while the fragments from the double digest indicate the relative distances between the *BamHI* and *EcoRI* sites ([Fig. 2](#)). *EcoRI* digestion produced three restriction fragments, 1, 2.5, and 9.5 kb in size. *BamHI* digestion also produced three fragments, 2, 5, and 6 in size. These findings indicate that each enzyme has two restriction sites on the linear DNA. A linear DNA will give  $(n+1)$  fragments, where  $n$  is the number of restriction sites. The double digest reveals five restriction fragments—0.5, 1, 2, 4, and 5.5 kb in size—indicating four cleavages. The correct, consecutive ordering of fragments must be determined by reconciling each of the single-enzyme digests with the fragments arising from the double digest. All possible fragment orders for both enzymes are listed in [Figure 2](#), as well as the order that gives the restriction map that agrees with the double-digest fragment sizes. In the example given, all fragment sizes from the two single digests and from the double digest sum to 13 kb each, the size of the original DNA. When determining a restriction map, one must remember that some fragments may be larger than the resolution power of the gel or that digestion may generate different fragments of the same size, which will migrate together. In such cases, alternative enzymes must be employed. Another approach uses the timed partial digestion of uniquely end-labeled, linearized DNA, which generates fragments ranging from the full-length DNA to the shortest fragment containing the label. The label identifies one end of the molecule. A restriction map can then be determined by a comparative analysis of the collection of fragments produced.

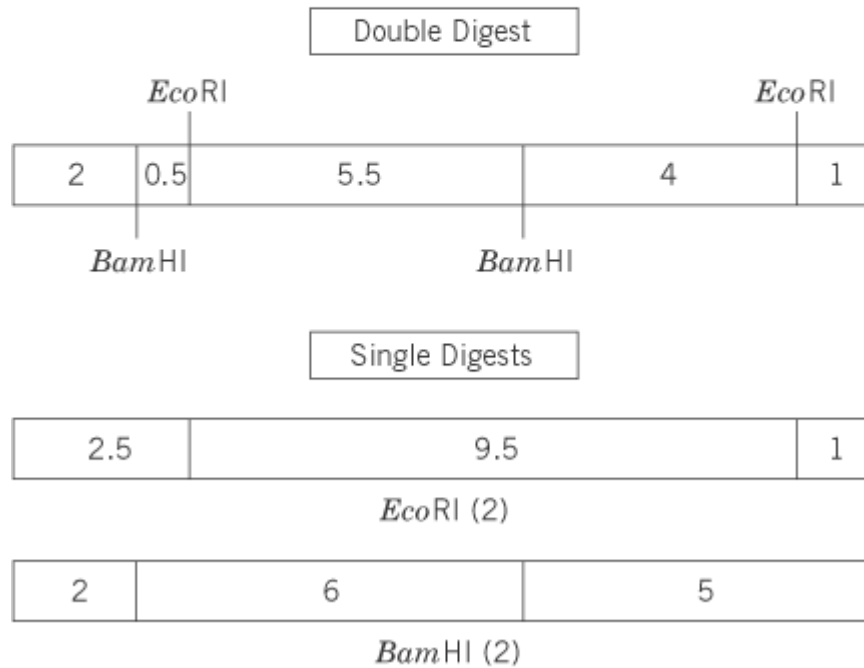
**Figure 1.** Idealized pattern of DNA fragments after digestion of a linear 13-kb DNA with *EcoRI* and *BamHI*, separately or combined, separation by electrophoresis, followed by visualization.



**Figure 2.** Restriction map derived from the data in Figure 1: (a) all possible combinations of linear restriction fragment orders from the individual *Eco*RI or *Bam*HI digests; (b) diagram of the final restriction maps showing the correct, unique order of restriction fragments from a double digest with both *Eco*RI and *Bam*HI of the 13 kb linearized DNA. Fragmentation patterns (2) of *Eco*RI and (2) of *Bam*HI combine to give five fragments produced by the double-digest.

- |   |   |
|---|---|
| Possible <i>Eco</i> RI<br>fragment orders | Possible <i>Bam</i> HI<br>fragment orders |
| 1. 1 kb --- 2.5 kb --- 9.5 kb             | 1. 2 kb --- 5 kb --- 6 kb                 |
| 2. 2.5 kb --- 9.5 kb --- 1 kb             | 2. 2 kb --- 6 kb --- 5 kb                 |
| 3. 9.5 kb --- 1 kb --- 2.5 kb             | 3. 5 kb --- 2 kb --- 6 kb                 |

(a)



(b)

Restriction digestion breaks DNA into fragments of a manageable size for analysis (see [Restriction Fragment](#)). Restriction maps are thus helpful in determining gene order or chromosomal organization. Generation of restriction maps from different organisms reveals similarities or differences in their genomic arrangements. One important application of restriction mapping is genetic fingerprinting: the typing of individual DNA molecules based on differences in restriction map patterns. The differences in map patterns are reflected in *restriction fragment length polymorphisms* (RFLPs), which result from changes in restriction sites or distances between restriction sites due to local sequence variations. In some cases, RFLPs can be used as genetic markers for diagnosis of diseases, and can aid in isolating the gene or identifying DNA sequence aberrations causing the disease. Applications for analysis and manipulation of restriction fragments are discussed under the topic [restriction fragment](#).

### Restriction–Modification Systems

Restriction–modification (R-M) systems in bacteria consist of pairs of [enzymes](#) that presumably provide protection from exogenous foreign DNA, for example, from infecting **bacteriophage**. This phenomenon was observed almost 50 years ago, when some strains of bacteria, those with an R-M

system, were observed to be less susceptible to bacteriophage infection and propagation than others. Their infection by bacteriophage was “restricted”; titers of bacteriophage produced were one to five orders of magnitude lower than in other strains, which lack an R-M system (1). Those bacteriophages that did propagate on the restricted hosts, however, were subsequently “immune” to the restriction and were competent to grow normally on those host strains.

The R-M systems act to prevent or to “restrict” the entrance of foreign DNA into the host bacterial cell by fragmenting it enzymatically with the restriction endonuclease, or [restriction enzyme](#). The foreign DNA is recognized because it is not methylated, or “modified,” whereas the DNA of the host is, due to the action of a host [methyltransferase](#) (see [Methylation, DNA](#)). Over 25% of 10,000 bacterial strains surveyed have been found to contain at least one R-M modification system; some strains contain as many as six (1, 2). R-M systems are dispensable, as bacteria without them are viable and display no deficiencies other than greater sensitivity to bacteriophage infection. The genes for R-M systems can be located on the [chromosome](#), on a **plasmid**, or on a prophage, and those for each pair of endonuclease and methyltransferase are usually found in proximity. More than 160 R-M systems have been cloned, and most of their genes have been sequenced (3).

## 1. R-M Enzymes

The R-M phenomenon is explained by the action of the two enzymatic activities making up each R-M system. One activity is a restriction endonuclease, or restriction enzyme, which cleaves DNA, whereas the other is a methyltransferase, or methylase, which methylates the host DNA. These two activities can occur on two separate proteins or on an enzyme complex displaying both activities. The restriction enzyme recognizes a specific DNA sequence (usually 4–8 bp long) and catalyzes phosphodiester bond cleavage at sites within or outside of the sequence to produce a double-strand cut in the duplex DNA. The enzyme thus fragments the DNA, which makes it susceptible to further degradation by general exonucleases. The DNA recognition sequence, or “restriction site,” can be unique or degenerate (representing more than one unique sequence), continuous or interrupted by unrecognized base pairs, **palindromic** or asymmetric (see [Staggered Cut](#)). The partner modification methyltransferase activity recognizes the same sequence as its cognate restriction enzyme and catalyzes the transfer of a methyl group from the cofactor [S-adenosyl-L-methionine](#) (AdoMet) to specific recipient bases in the sequence; this methylation prevents cleavage by the partner endonuclease. Sequences methylated on either one or both strands are resistant to endonuclease cleavage. The methyltransferases fall into two major mechanistic classes: one group methylates the 5-carbon of cytosine to form [5-methylcytosine](#) (5mC) and the other the exocyclic amino group of either cytosine to form *N*-4-methylcytosine (N4mC) or of adenine to produce *N*-6-methyladenine (N6mA) (see [Methylation, DNA](#)).

## 2. R-M System Nomenclature

R-M systems are named with an italicized three letter abbreviation that is derived from the name of organism in which that R-M system resides. The first italicized letter (uppercase for protein and lowercase for gene designations) comes from the genus, the lowercase second and third letters from the species. Any strain designation follows the italicized abbreviation, and Roman numerals are used to designate different R-M proteins within the same organism. For example, *EcoRI*, *EcoRII*, and *EcoRV* are different systems from *Escherichia coli* strain R. For the enzymes, the letters R or M, followed by a dot (•) are used before the name: R•*EcoRI* or M•*EcoRI* are the cognate enzyme pair, endonuclease and methylase, respectively, from *E. coli* strain R. For the genes of R-M systems, the name designation is completely italicized and followed with the letter (R) or (M) to designate the activity, eg, *ecoRIM* and *ecoRIR*, for the methylase and endonuclease genes, respectively (2, 4).

## 3. R-M Gene Expression and Organization

Expression of the genes for R-M systems is assumed to be tightly regulated. The host genome carrying an R-M system must be continuously protected from its restriction endonuclease even

during physiologically difficult times, such as intermittent starvation with concomitant decreases in AdoMet concentration. The genes for R-M systems are spatially linked and tandemly organized. The adjacent genes can be arranged (1) in parallel, with the 5' end of one gene following the 3' end of the other; (2) convergently, with the 3' ends of the genes in proximity; or (3) divergently, with the 5' ends of the genes near one another. Comparisons of the sequences of R-M enzymes have provided information about their evolutionary relationships (5), as well as the means of their classification. Three different types of R-M systems (I, II, and III) are recognized, and different assemblies of proteins exist for each system (Table 1).

**Table 1. Restriction Modification Systems**

| System                     | Type   |   |   |   |
|----------------------------|--|---|---|---|
|                            | I  | II  | IIs   | III   |
| Enzymes                    | 3 subunits:<br>R, M, S (R <sub>2</sub> M <sub>2</sub> S)                           | Independent<br>subunits:<br>R and M<br>proteins (R <sub>2</sub><br>and M) | Independent<br>subunits:<br>R and M<br>proteins (R<br>and M)    | 2 subunits:<br>R, M (RM)  |
| Cofactors                  | Restriction:<br>Mg <sup>2+</sup> and ATP;<br>AdoMet is<br>allosteric effector      | Restriction:<br>Mg <sup>2+</sup><br><br>Modification:<br>AdoMet           | Restriction:<br>Mg <sup>2+</sup><br><br>Modification:<br>AdoMet | Restriction:<br>Mg <sup>2+</sup> , ATP is<br>allosteric<br>effector;<br>stimulated by<br>AdoMet<br><br>Modification:<br>AdoMet,<br>stimulated by<br>ATP |
| Recognition<br>site        | Asymmetric <sup>a</sup> , 3–<br>5 half-sites <sup>c</sup> with<br>spacer region    | Palindromic <sup>b</sup><br>4–8 bp  | Asymmetric,<br>4–7 bp   | Asymmetric, 5–<br>6 bp  |
| Cleavage site              | Variable<br>distances from<br>recognition site,<br><br>ATP-driven<br>translocation | Within<br>recognition<br>sequence   | Under 20 bp<br><br>3' of<br>recognition<br>site                 | 25–30 bp<br><br>3' of<br>recognition site   |
| Modification<br>properties | Opposite strands<br>of each half-site<br><br>(M <sub>2</sub> S)                    | Within<br>recognition<br>sequence   | Methylation<br>of one or both<br>strands                        | One strand of<br>each<br>recognition site   |



|             | Prefer hemimethylated sites over nonmethylated sites                   |                                    |  |                                       |
|-------------|--|------------------------------------|--|---------------------------------------|
| Sequence    | Conferred by S subunit R and M mutually exclusive                      | Intrinsic to enzyme                | Intrinsic to enzyme  | Conferred by M subunit                |
| Specificity |  | Separate                           | Separate   | Simultaneous if all cofactors present |
| Enzymatic   |  |                                    |  |                                       |
| Activities  |  |                                    |  |                                       |
| Example     | <i>EcoAI</i> :<br>GAG(N <sub>7</sub> )GTCA<br>CTC(N <sub>7</sub> )CACT | <i>EcoRI</i> :<br>GAATTC<br>CTTAAG | <i>FokI</i> :<br>GGATG(N <sub>9</sub> )<br>CCTAC(N <sub>13</sub> ) | <i>HinfIII</i> :<br>CGAAT<br>GCTTA    |

<sup>a</sup> Recognition sites not having the same sequence in both strands of DNA.

<sup>b</sup> Recognition sites having the same sequence in both strands of DNA.

<sup>c</sup> A site containing only half of the necessary DNA sequence for complete enzyme recognition.

### 3.1. Type I R-M Systems

Fewer than 20 type I R-M systems have been found, but they are present in *E. coli*, *Salmonella typhimurium*, and *C. freundii*. Type I systems contain both the restriction and modification functions in the same multisubunit enzyme. The enzyme is at least a pentameric complex (R<sub>2</sub>M<sub>2</sub>S) consisting of three nonidentical subunits, referred to as R (for restriction), S (for DNA sequence specificity), and M (for methylation). The enzymes require Mg<sup>2+</sup>, AdoMet, and ATP for activity, or as **allosteric** effectors (6). If the restriction sequence is fully methylated, ATP hydrolysis drives the dissociation of the enzyme from DNA. If the restriction site is hemimethylated, the enzyme methylates the other strand and dissociates. If unmethylated, the R subunit cleaves at random DNA sites after ATP-driven translocation of the enzyme to variable distances (up to 10<sup>3</sup> bp) from the recognition sequence (6, 7). The asymmetric recognition sequence comprises two half-sites, each 3–5 bp long, separated by a nonspecific sequence of 6–8 bp. In all type I enzymes, the methyltransferase subunit forms specific N6mA residues at each half-site, preferring hemimethylated DNA substrates over unmodified ones (8). The S subunit provides DNA sequence recognition. The type I systems are grouped into three different families: IA, IB, and IC, which are differentiated by their genetic location, immunological cross-reactivity, and gene sequences. The genes of the IA and IB families are located on the chromosome and have the gene order *hsdR*, *hsdM*, and *hsdS* (*hsd* indicates *host specificity determinant*). The genes of the IC family are encoded on a plasmid in the order *hsdM*, *hsdS*, and *hsdR*. All families of the type I R-M systems have two adjacent **transcriptional** units, with the *hsdM* and *hsdS* genes transcribed from one **promoter** and *hsdR* from the other. Although type I enzymes (such as *EcoB* and *EcoK*) were the first restriction endonucleases to be discovered, type II enzymes have proved more numerous, simpler in composition, and more useful in practical applications.

### 3.2. Type II R-M Systems

More than 2600 Type II systems, with greater than 230 different specificities exist (3). Type II systems have discrete restriction and modification enzymes. Each enzyme of a pair recognizes the same DNA sequence. Most type II enzymes recognize palindromic, duplex DNA sequences, such as

GAATTC, whose complementary strand has the same 5'–3' sequence. The restriction endonucleases are homodimeric, require  $Mg^{2+}$ , and cleave phosphodiester bonds within or immediately adjacent to (type IIs; see below) the recognition sequence to leave a staggered or blunt double-strand cut (see [Staggered Cut](#) and [Restriction Enzymes](#)). The methyltransferases are monomeric and require the cofactor AdoMet. Methylation takes place on both strands of the DNA duplex within the recognition sequence, rendering the sequence refractory to cleavage. Interestingly, little amino acid sequence similarity exists between partner endonucleases and the methyltransferases, suggesting that the enzymes evolved independently ([5](#), [9](#)) and reflecting the fundamentally different chemical reactions they catalyze and the mechanisms they use for DNA recognition. Those type II R-M systems in which the two genes are aligned consecutively are believed to operate as single transcriptional units, with the order of the genes being unimportant. Those with the genes arranged divergently or convergently may be subject to independent transcriptional control. The genes for a type II R-M system can be located on the chromosome or on a plasmid. Some type II systems also contain an **open reading frame** encoding a “controller” or “C” protein, with sequence similarity to some [DNA-binding proteins](#). These controller proteins probably regulate the expression of the genes. For example, disruption of the *BamHIC* gene in the BamHI R-M system leads to an increase in modification and a decrease in restriction cleavage ([2](#)).

### 3.3. Type IIs Enzymes

A subset of the type II R-M systems, termed type IIs, share the property of having independent restriction and modification enzymes, but they recognize uninterrupted and asymmetrical sequences (4–7 bp long) and cleave DNA 3' to the recognition sequence, up to 20 bp away, leaving a staggered double-strand cut (see [Staggered Cut](#)). These endonucleases act as monomers and require  $Mg^{2+}$  for activation and cleavage. Type IIs target sequences are asymmetric, and the methyltransferases act as monomers to methylate one strand at a time within the recognition sequence. In some type IIs systems, methylation of both strands is performed by a pair of methyltransferases, which may (*eg*, *HgaI*) or may not, (*eg*, *Alw26I*) be of the same class (5mC, N4mC, or N6mA) (see [Methyltransferase, DNA](#)). In another system, *FokI*, methylation is accomplished by a fused, bifunctional enzyme ([2](#)). There are over 80 characterized type IIs enzymes, representing over 35 specificities ([6](#)). Their sequence specificities are similar to those of type III enzymes.

### 3.4. Type III R-M Systems

Only a few type III R-M systems are known, with four different specificities identified ([6](#)). A single bifunctional enzyme catalyzes both the endonuclease and the methyltransferase restriction activities. The enzymes are composed of two nonidentical subunits: the M subunit (encoded by the *mod* gene) and the R subunit (encoded by the *res* gene). The R subunit must be complexed with the M subunit for restriction activity, because the M subunit provides the sequence specificity for the enzyme. The two enzymatic activities compete for the uninterrupted, asymmetric DNA recognition sequence, which is usually 5–6 bp long. Restriction activity requires that two copies of the recognition site be proximal on the DNA and in opposite orientations to one another. Such an array of sites may also be viewed as one symmetrical sequence separated by a spacer region of undefined length ([10](#)). Type III restriction activity is absent when either a single site or two recognition sites in the same orientation are present. Cleavage takes place 25–30 bp away to the 3' side of the DNA recognition sequence. If one or both of the DNA strands is (are) methylated, no cleavage occurs. The M subunit can act independently as a methyltransferase, requiring AdoMet and methylating only one strand of the duplex recognition sequence at a time, which is sufficient to inhibit the restriction reaction. Methylation is independent of the number and orientation of the restriction sites, suggesting that the enzyme reacts with single sites only ([10](#)).

### 3.5. Other R and M Systems

Other variations of restriction and modification systems exist. For example, endonucleases that recognize and cleave only methylated DNA have been described. Three such methylation-dependent restriction systems have been isolated from the K12 strain of *E. coli*: (1) *Mrr* (methyladenine recognition and restriction) cleaves DNA with sequences containing N6mA and 5mC bases; (2)

*McrA* (modified cytosine restriction) cleaves sequences containing 5mC; and (3) *McrBC* cleaves DNA containing 5mC, N4mC, or 5-hydroxymethylcytosine (11). No corresponding methyltransferases are known for these systems. These systems are thought to restrict foreign methylated DNA, and their presence must be recognized when trying to clone methylated DNA. Among the type II enzymes, R-*DpnI* cleaves only at methylated GmATC sites, whereas *DpnII* cleaves only at unmethylated GATC sites. These and other enzymes with similar properties can be used to assess the methylation state of DNA.

Restriction-independent methylases have also been found, such as the adenine-specific Dam and DNA-cytosine (Dcm) methyltransferases. The Dam enzyme methylates the adenine in the sequence GATC and has various functions in methyl-directed [mismatch repair](#), [DNA replication](#), and **gene regulation** (12). Dcm functions in very short patch [DNA repair](#), which serves to repair deaminated cytosine bases using the G-containing template of the complementary strand of the DNA as a guide (2, 5). Finally, **intron**-encoded endonucleases that catalyze methylation-independent restriction have been described. These enzymes recognize large specific DNA sequences (18bp) and cleave within the recognition sequence or up to 20 bp away. They are believed to function in site-specific intron [transposition](#) and share some characteristics with type I restriction endonucleases (6).

#### 4. Antirestriction Systems

Bacteriophage have evolved diverse antirestriction mechanisms to defeat the R-M systems of bacteria. Such evasive measures include (1) production of phage-encoded proteins that inhibit host R-M enzymes or destroy R-M cofactors, (2) stimulation of the host modification function, (3) phage self-modification of DNA (using modified bases in their genomes), and (4) evolutionary elimination of restriction sites from the phage genome. As examples, phage T3 contains the gene for the enzyme AdoMet hydrolase, which destroys AdoMet; the Ral protein produced by phage  $\lambda$  inhibits restriction and stimulates the methylation function of type IA enzymes (11). The T-even phages carry glycosylated hydroxymethylcytosine bases in an effort to evade restriction (5). On the other hand, a restriction enzyme encoded on bacterial plasmid RtsI, *PvuRtsI1*, specifically cleaves only phage DNA sequences containing hydroxymethylcytosine residues (13). Evolution has clearly promoted the development of both defensive and offensive mechanisms in the competition between bacteria and bacteriophage.

#### 5. Concluding Remarks

R-M systems probably provide an “immune system” for bacteria that shares characteristics with that of eukaryotes (see [Immune Response](#)). For example, foreign, “non-self-”DNA is distinguished from “self-”DNA (11). A provocative, alternative hypothesis for the existence of R-M systems suggests that they evolved because of the “selfishness” of their genes. The sociobiologically derived “selfish gene” theory dictates that under certain circumstances natural selection accommodates proliferation of genes that are potentially deleterious to the host carrying them (14). In support of this view, two groups have shown that a daughter bacterial cell not receiving the plasmid encoding the R-M system of the parent is subject to chromosomal DNA cleavage by residual restriction endonuclease remaining after cell division if the residual methylase activity cannot protect all the restriction sites efficiently (15, 16). Thus, selection for progeny carrying the R-M system exists in the absence of their “immunity” role. Other rationalizations for the existence of R-M systems include roles for them in DNA [recombination](#) or repair, regulation of gene expression, and acquisition of foreign genes, but these roles are less certain (2).

The discovery of R-M enzymes has revolutionized several areas of research and biotechnology (see [Cloning](#)). Isolation and characterization of R-M genes and purification of the proteins have led to greater understanding of protein–DNA and enzyme-cofactor interactions. For example, a novel example was found in the [X-ray crystallography](#) structures of two type II 5mC methyltransferases with their target DNA: an extrahelical cytosine base is flipped out of the DNA into pockets in the enzymes (see [5-Methylcytosine](#)). The Type II restriction enzymes provide a vast array of sequence-

specific DNA cleavage tools that allow manipulation and manageable analysis of an otherwise formidable macromolecule. Several applications of restriction endonucleases and methyltransferases are discussed in the entries [Restriction Fragment](#), [Restriction Map](#), and [Staggered Cut](#). R-M systems provide ideal models for studying basic genetic and biochemical mechanisms, as well as serving as practical tools for expanding other areas of biological research.

### Bibliography

1. R. J. Roberts and S. E. Halford. (1993) in *Nucleases*, 2nd ed. S. M. Linn, R. S. Lloyd, and R. J. Roberts, eds., Cold Spring Harbor Laboratory Press, New York, pp. 35–88.
2. G. G. Wilson and N. E. Murray (1991) *Annu. Rev. Genet.* **25**, 585–627.
3. G. G. Wilson (1996) personal communication.
4. W. Szybalski, R. M. Blumenthal, J. E. Brooks, S. Hattman, and E. A. Raleigh (1988) *Gene* **74**, 279–280.
5. T. A. Bickle and D. H. Kruger (1993) *Microbiol. Rev.* **57**, 434–450.
6. H.-Y. Eun (1996) in *Enzymology Primer for Recombinant DNA*, Academic Press, San Diego, pp. 233–306.
7. D. T. F. Dryden, S. S. Sturrock, and M. Winter (1995) *Nat. Struct. Biol.* **2**, 632–635.
8. I. A. Taylor, K. G. Davis, D. Watts, and G. G. Kneale (1994) *EMBO J.* **13**, 5772.
9. R. Korona, B. Korona, and B. R. Levin (1993) *J. Gen. Microbiol.* **139**, 1283–1290.
10. S. Saha and D. N. Rao (1995) *J. Mol. Biol.* **247**, 1–9.
11. J. Heitman (1993) in *Genetic Engineering*, Vol. **15**, J. K. Setlow, ed., Plenum Press, New York. pp. 57–107.
12. V. U. Nwosu (1992) *Biochem. J.* **283**, 745–750.
13. L. Janosi, H. Yonemitsu, H. Hong, and A. Kaji (1994) *J. Mol. Biol.* **242**, 45–61.
14. J. J. Bull, I. J. Molineux, and J. H. Werren (1992) *Science* **256**, 65.
15. T. Naito, K. Kusano, and I. Kobayashi (1995) *Science* **267**, 897–899.
16. S. Kulakauskas, A. Lubys, and S. D. Ehrlich (1995) *J. Bacteriol.* **177**, 3451–3454.

### Suggestions for Further Reading

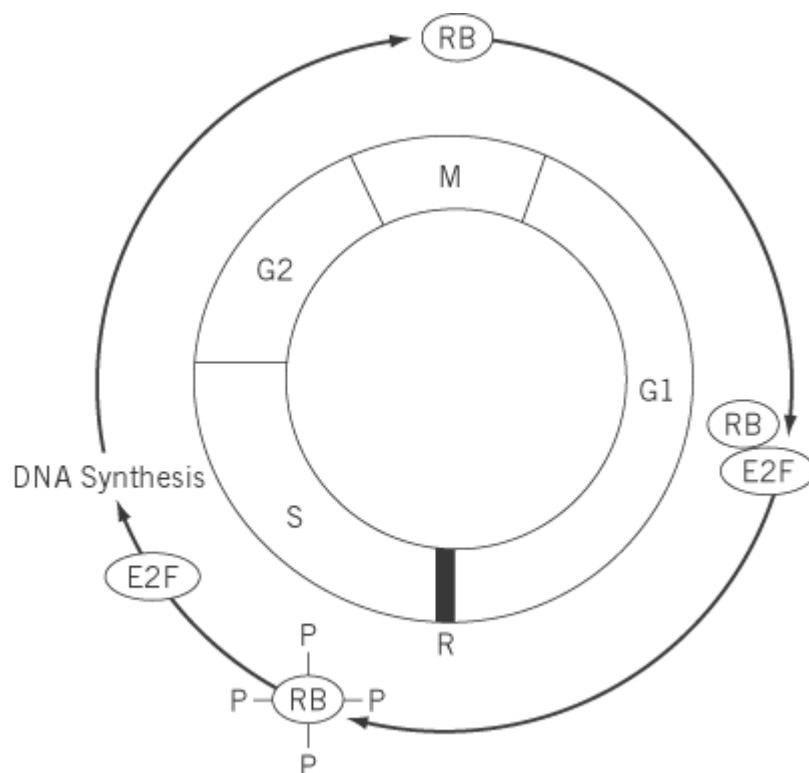
17. G. G. Wilson and N. E. Murray (1991) Restriction and modification systems. *Annu. Rev. Genet.* **25**, 585–627.
18. R. J. Roberts and S. E. Halford. (1993) "Type II restriction endonucleases", in *Nucleases*, 2nd ed., S. M. Linn, R. S. Lloyd, and R. J. Roberts, eds., Cold Spring Harbor Laboratory Press, New York, pp. 35–88.
19. H.-Y. Eun (1996) "Restriction endonucleases and modification methylases", in *Enzymology Primer for Recombinant DNA*, Academic Press, San Diego, pp. 233–306.
20. J. Heitman (1993) "On the origins, structures and functions of restriction-modification systems, in" *Genetic Engineering*, Vol. **15**, J. K. Setlow, ed., Plenum Press, New York. pp. 57–107.

### Retinoblastoma Gene

The retinoblastoma (Rb) gene was the first [tumor suppressor gene](#) identified by genetic analysis of a **chromosomal** deletion observed in retinoblastoma patients. Familial retinoblastoma is a rare childhood disease that occurs in 1 out of 20,000 patients. In a majority of the cases, it occurs as a

sporadic form of cancer, where the tumor is formed only in one eye, but in others the patients develop tumors in both eyes. The Rb gene is located on chromosome 13q14.2 and is often deleted in tumors of the retinoblastoma. The Rb gene is also deleted or mutated in a number of human cancers, including lung, breast, pancreas, bladder, prostate, and osteocarcinomas. The protein encoded by the retinoblastoma gene is a 105-kDa polypeptide chain, also known as p105RB, which can be **phosphorylated**. The phosphorylation state of p105RB is tightly regulated during the [cell cycle](#), is maximal during S-phase and minimal soon after mitosis. Stimulation of quiescent [T cells](#) leads to hyperphosphorylation of p105RB, whereas differentiation of myeloid cells is associated with very low levels of phosphorylated p105RB. Interestingly, only the hypophosphorylated form of RB has tumor suppressor activity. A simple model predicts that p105RB is phosphorylated in the late G<sub>1</sub>/S-phase and dephosphorylated in the late M-phase, and this has been observed (Fig. 1). In its dephosphorylated state, the Rb gene product binds to a group of [transcription factors](#) known as the E2F class of proteins, which are required for [DNA synthesis](#) (1, 2). [Growth Factors](#) stimulate the phosphorylation of the p105RB protein by activating a group of kinases known as CDC kinases, which carry out the phosphorylation. The p34CDC2 kinase is a candidate kinase (1). Phosphorylation in turn reduces the affinity of p105Rb for E2F, allowing these transcription factors to bind to DNA and to activate transcription of genes required for DNA synthesis. Interestingly, growth inhibitors, such as [transforming growth factor b](#) (TGF-b), which inhibits cell proliferation, prevent phosphorylation of p105 RB, even in the G<sub>1</sub>/S-phase. These observations suggest that the phosphorylation of p105RB is a critical regulatory event in cell proliferation.

**Figure 1.** The role of Rb in cell cycle regulation. The Rb protein is in an unphosphorylated state in the G<sub>1</sub>-phase of the cell cycle. In this state, it associates with E2F and blocks its ability to transcribe genes required for DNA synthesis. Phosphorylation of Rb makes it unable to bind to E2F, which in turn results in the transcription of genes that are required for DNA synthesis. Then the cell cycle progresses from the G<sub>1</sub>- to the S-phase.



p105RB binds to the [adenovirus](#) E1A protein, and this association is essential for adenoviral transformation. Hence, the transforming DNA viruses induce transformation by sequestering the

hypophosphorylated form of RB105. The biochemical mechanism of tumor suppression by p105RB is at the level of the expression of the proto-oncogenes *c-myc* and *c-fos* because RB expression suppresses the expression of both genes.

Recent studies have shown that the mammalian gene codes for two other Rb-related genes, known as p107 and p130, which is also known as Rb-2. Like p105RB, these two proteins act as cell-cycle regulators in cells where they are expressed and thus function as tumor suppressor genes.

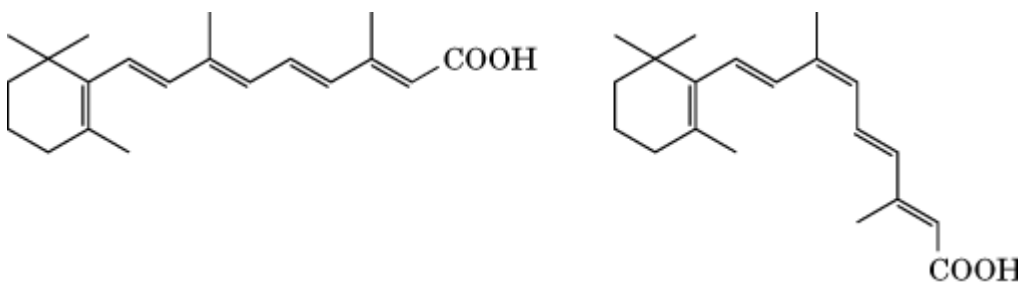
### Bibliography

1. X. Graña and E. P. Reddy (1996) *Oncogene* **11**, 211–220.
2. R. A. Weinberg (1991) *Science* **254**, 1138–1145.

## Retinoic Acids

Retinoic acids (RAs) are natural intracellular oxidation products of vitamin A (retinol) that mediate the biological effects of vitamin A (1). The most active compound is the all-*trans*-retinoic acid (t-RA), [2E, 4E, 6E, 8E]-3,7,-dimethyl-9-(2,6,6-trimethylcyclohex-1-en-1-yl) nona-2,4,6,8-tetraenoic acid (retinoic acid, or vitamin A acid) (Fig. 1). A stereoisomer of t-RA, 9-*cis*-retinoic acid (9-*cis*-RA) is also biologically active. The RAs play a major role in vertebrate pre- and postnatal development, are modulators of cellular proliferation and differentiation, and are effectors of morphogenic changes (2-4).

**Figure 1.** Structure of all-*trans*-retinoic acid and 9-*cis*-retinoic acid.



The biological effects of RAs result from their action on gene [transcription](#) and are mediated by specific **receptors** belonging to the superfamily of nuclear receptors. Like [steroid hormone receptors](#), RA receptors represent ligand-inducible [transcription factors](#). Two such families of receptors have been described, the RAR and the RXR (5-7). Three genes that comprise the RAR family (a, b, and g) have been identified. Each of these subtypes can be further divided into different protein **isoforms**, which arise either from different **promoter** usage or from [alternative splicing](#). Although the various receptor isoforms can substitute for one another, they are not truly equivalent functionally. The RXR family also consists of three subtypes a, b, and g and is **homologous**, but with primary structures and ligand specificity that differ substantially from those the RAR family. The highest degree of similarity between RAR and RXR lies in the DNA-binding domain (63% identical), and the rest of the protein molecules show about 30% identity. RARs and the thyroid

hormone receptors belong to the same receptor subclass, due to their high degree of similarity and to an identical **P element** in the first **zinc finger** of the DNA-binding domain (see [Steroid Hormone Receptors](#)). The physiological ligand for RARa, RARb, and RARg is t-RA. Originally, RXR was considered to be an **orphan receptor**, but 9-*cis*-RA was subsequently demonstrated to be its ligand. RAR binds to RA response elements (direct repeats of AGGTCA, separated by five base pairs) as a homodimer, and as a heterodimer it binds to RXR. RXR can act as a homodimer, binding to RXR response elements (direct repeats of AGGTCA, separated by one base pair), but also as a heterodimer with RAR or with the receptors for [thyroid hormones](#), vitamin D, and peroxisome proliferator (8). The heterodimerization with these receptors increases the apparent binding of each receptor to their cognate response elements and increases the regulatory potential of RAs. A physical interaction of RAR with AP1 proteins has been demonstrated, accounting for the repressive action of RAs on several genes. RAR can be **phosphorylated** by protein kinase A *in vitro* and *in vivo*. In [transfection](#) experiments, it was shown that phosphorylation of Ser369 of RAR increases its capacity to activate a [reporter gene](#) carrying a RAR response element.

Retinoids are important regulators of epidermal cell homeostasis, skin formation, and maintenance and have been used to treat skin disorders (cystic acne, psoriasis, epithelial cancers). *In vitro*, RAs inhibit differentiation of keratinocytes, inhibit the proliferation and squamous differentiation of head and neck cancers, and inhibit growth of estrogen receptor-positive human breast carcinoma cells in culture. RAs are required for eye and limb morphogenesis and for spermatogenesis. They induce remission in promyelocyte leukemias (PMLs). Translocation of the RARa gene and its fusion in-frame with the PML gene is a key event in the genesis of acute promyelocytic leukemia (APL). The PML gene regulates **hematopoiesis** and controls cell growth and tumorigenesis. Its function is essential for the tumor growth-suppressive activity of RA and for the ability of RA to induce terminal myeloid differentiation of precursor cells by transactivation of specific genes regulating the [cell cycle](#) (eg, gene p2/WAF1/CIP1). Disruption of the PML gene by fusion with RARa and translocation induces the pathway to APL (9). Administration of excess RA to mice embryos can lead to teratogenic effects and limb malformations (10).

### Bibliography

1. B. F. Tate, A. A. Levin, and J. F. Grippo (1994) Trends Endocrinol. Metab. **5**, 189–194.
2. M. B. Sporn and A. B. Roberts (1983) Cancer Res. **43**, 3034–3040.
3. C. Thaller and G. Eichele (1987) Nature **327**, 625–628.
4. P. Kastner, M. Mark, and P. Chambon (1995) Cell **83**, 809–869.
5. A. Zelent et al. (1989) Nature **339**, 714–717.
6. J.-H. Xiao et al (1995) J. Biol. Chem. **270**, 3001–3011.
7. D. J. Mangelsdorf and R. M. Evans (1995) Cell **83**, 841–850.
8. X. K. Zhang and M. Pfahl (1993) Trends Endocrinol. Metab. **4**, 156–162.
9. Z. G. Wang et al. (1998) Science **279**, 1547–1551.
10. J. G. Rutledge et al. (1994) Proc Natl Acad Sci **91**, 5436–5440.

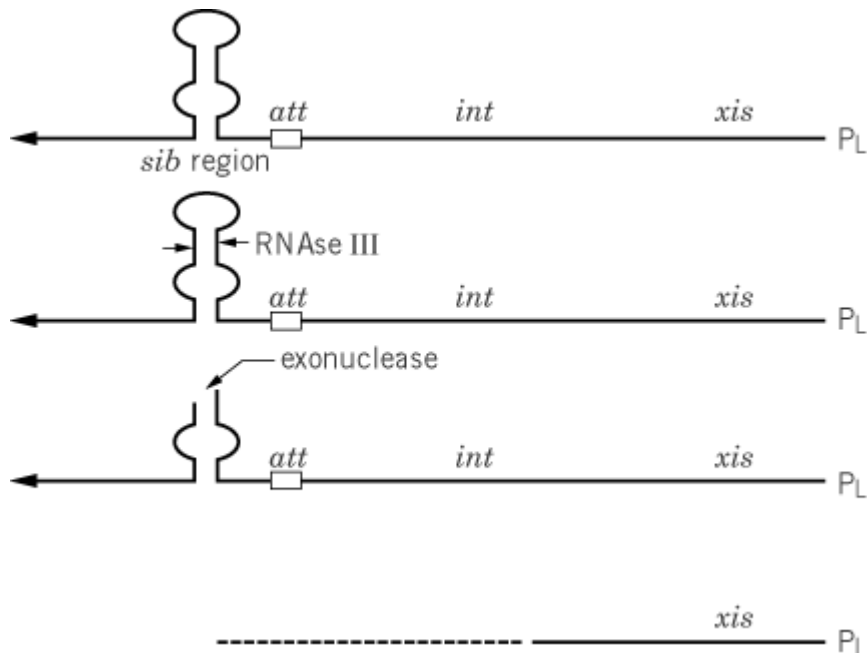
### Suggestions for Further Reading

11. N. B. Sporn and A. B. Roberts, and D. S. Grodman, eds. (1984) *The Retinoids*, Academic Press, Orlando, FL.
12. A. L. Means and L. J. Gudas (1995) The roles of retinoids in vertebrate development. Annu. Rev. Biochem. **64**, 201–233.
13. P. Chambon (1994) The retinoid signaling pathway: molecular and genetic analysis. Semin. Cell Biol. **5**, 115–125.

## Retroregulation

Retroregulation of gene expression was discovered in the *int* gene expression of [lambda phage \(1\)](#). The *int* gene encodes the [integrase](#) that catalyzes integration (and excision) of  $\lambda$  DNA to (or from) the *Escherichia coli* chromosome. It is transcribed from the  $P_L$  promoter of  $\lambda$  in the lytic cycle. The  $P_L$ -derived *int* transcript is, however, very labile, and the level of integrase expression is low. Point mutations in the region 260 bp downstream of the *int* gene are known to increase the expression of *int*. Because these mutations are clustered in the *b2* region of  $\lambda$ , they are named *sib* for something in *b2*. The *sib* region does not encode a protein, but its RNA transcript can form an extended stem-loop structure that provides a substrate structure for a double-stranded RNA-dependent ribonuclease, RNase III. Therefore, the wild-type *sib* RNA is cleaved by RNase III and rapidly degraded by 3'-to-5' exonuclease, such as [polynucleotide phosphorylase](#) (PNPase). The *sib* mutations destroy the RNase III recognition structure and stabilize the *int* mRNA, resulting in the increased synthesis of integrase; hence, this is referred to as retroregulation. In this case, retroregulation prevents  $\lambda$  phage from integrating into the bacterial genome in the lytic cycle. Although little integrase protein is synthesized from the  $P_L$  promoter after phage infection, much integrase is expressed after induction of  $\lambda$  prophage, where the *sib* sequence is no longer placed downstream of *int* because  $\lambda$  phage attachment site, *attP*, truncates these two regions.

**Figure 1.** Retroregulation of  $\lambda$  phage *int* gene. The *int* mRNA transcribed from the  $P_L$  promoter is unstable because of dsRNA cleavage with RNase III and subsequent degradation by PNPase. The *int* transcript is stabilized by *sib* mutations which destroy the RNase III recognition structure.



The *int* gene has its own **promoter** upstream of the [initiation codon](#). The *int* promoter is dependent on the  $\lambda$  transcription activator protein, *cII*. When *cII* activates transcription, the *int* transcript terminates at a r-independent terminator located within the *sib* region. Hence, the *int* transcript does



not form a RNase III-recognition site and is stable. On the other hand, the  $P_L$  transcript reads through this terminator because of the **antitermination** mechanism.

The 1.1 and 1.2 genes of **T7 phage** provide another example of retroregulation. They are encoded in the T7 early transcript, and the distal site (R5) can regulate these upstream genes by the RNase III-dependent mechanism similar to the *int* gene.

## Bibliography

1. H. Echols and G. Guarneros (1983) In *Lambda II* (R. Hendrix, J. Roberts, F. Stahl, and R. Weisberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 75–92.

## Retrotransposons

Retrotransposons are a large group of [transposable elements](#) that includes [retroviruses](#) and other elements that resemble retroviruses, but have no extracellular phase (1). Well-studied retroviruses include human immunodeficiency virus ([HIV](#)) and Moloney murine leukemia virus (MoMuV). Well-studied retroviral-like elements include the [Ty elements](#) of yeast and the gypsy and copia elements of *Drosophila*.

The [transposition](#) substrate for retrotransposons is a DNA copy of the element generated by **reverse transcription** of an RNA copy of the “proviral” form of the element integrated into the genomic DNA. At the ends of this DNA copy are the special sequences, again arranged as [inverted terminal repeats](#), that will be acted upon by the [integrase](#). The chemical steps of transposition—that is, the DNA breakage and joining reactions that process this reverse transcription product and insert it into target DNA—are fundamentally identical to the reactions that mediate the translocation of other DNA-only transposable elements in bacteria, *Drosophila*, and [Caenorhabditis elegans](#) (2). Moreover, the retrotransposon integrases are related in structure to the [transposases](#) of DNA-only elements, including those from bacteria. Indeed, the [active site](#) of the bacterial transposase MuA (3), is virtually superimposable with the catalytic region of the HIV integrase (4), although there is little primary structure [homology](#) between these proteins.

A retroviral provirus found in genomic DNA has at its ends two directly repeated segments several hundred base pairs long called [long terminal repeats](#) (LTRS); at the termini of these segments are short sequences that are recognized by the integrase and mark the sites of [recombination](#) (Fig. 1). (Only the sequences at the outside tips of the LTRs are recognized as recombination sequences; recombination does not occur at the two internal sites.) The interior of the element encodes several [polyproteins](#) (gag, pol, and env) that are cleaved by element-encoded [proteinases](#) to yield their several final products. Gag includes MA, a matrix-associated protein, CA, the major structural component of the viral capsid, and NC, a nucleocapsid component that interacts with the viral genome. Pol includes reverse transcriptase, [ribonuclease H](#), and integrase; the order of these domains does vary among different elements. Env includes several capsid proteins involved in receptor recognition and viral entry. Retrotransposons have a structure similar to retroviruses, but lack an extracellular phase. Indeed, many retrotransposons lack or contain a defective version of env, which is key for the extracellular phase.

**Figure 1.** Retroviruses and retrotransposons. Retroviruses and retrotransposon generate a mobile DNA copy of the

element by using an element-encoded reverse transcriptase and ribonuclease H to convert an RNA copy of the element into a double-stranded DNA copy. Exposed 3'OH ends on this DNA segment, resulting either from reverse transcription or from processing by the element-encoded integrase, attack the target DNA at staggered positions (see Fig. 2 in [Transposable Elements](#)). These elements encode multiple polyproteins, gag, pol, and env: gag includes capsid components, pol includes reverse transcriptase, ribonuclease H, and integrase, and env includes functions for entry into new cells. In some, retrotransposons, the Ty1-copia class (Ty1 being a yeast element and copia a *Drosophila* element), no env gene is recognizable. In the Ty3-gypsy class (Ty3 being a yeast element and gypsy a *Drosophila* element), a defective version of env can sometimes be observed. These elements and other retrotransposons are very widespread.



Transposition begins by the transcription of the provirus by the host **RNA polymerase** to yield a viral RNA copy; some of these molecules serve as [messenger RNA](#) for synthesis of the viral proteins. Two copies of the RNA are packaged into the virion, which includes reverse transcriptase and integrase. Viruses bud from the cell and then bind to and enter another cell; retrotransposons move only intracellularly. While in the virion, the RNA is converted to double-stranded DNA. This resulting DNA, with assembled viral proteins, enters into the [nucleus](#), where integration occurs. As with other transposons, insertions of these elements can disrupt genes or bring host genes under the control of viral [enhancers](#).

The substrate for transposition is the DNA copy of the element. Recombination often initiates with the cleavage by integrase of several nucleotides from the 3' ends of the DNA, to expose the actual element termini, although in some elements this trimming step is not necessary. Thus, as with DNA-only elements, exposing the 3'OH ends of the transposon is a key step. The integrase then promotes the attack of these 3'OH ends on staggered positions in the target DNA; because of this stagger, the newly inserted transposon is flanked by short gaps that will be repaired by the host [DNA repair](#) machinery, resulting in target site duplications of a characteristic length for each element.

### Bibliography

1. J. D. Boeke and J. P. Stoye (1997) In *Retroviruses* (H. Varmus, S. Hughes, and J. Coffin, ed.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 343–435.
2. K. Mizuuchi (1992) *J. Biol. Chem.* **267**, 21273–21276.
3. P. Rice and K. Mizuuchi (1995) *Cell* **82**, 209–220.
4. F. Dyda, A. B. Hickman, T. M. Jenkins, A. Engleman, R. Craigie, and D. R. Davies (1994) *Science* **266**, 1981–1986.

## Retroviruses

Retroviruses comprise a class of single-stranded RNA **viruses** whose genetic information is propagated through a double-stranded DNA intermediate integrated into the host [chromosome](#). First documented in 1904 in the form of a nontumor-forming agent, equine infectious anemia virus and retroviruses of avian, mammalian, reptilian and piscine origin have now been documented (Table 1). They are associated with a variety of neurological disorders, wasting diseases, immunodeficiencies, and malignancies. Clinically significant retroviruses include type 1 and type 2 *human immunodeficiency virus* (**HIV-1** , HIV-2), the causative agents of *acquired immunodeficiency syndrome* (AIDS), and type 1 and 2 human *T-cell leukemia virus* (HTLV-1 , HTLV-2), which manifest themselves in adult T-cell leukemia (ATL), tropical spastic paraparesis (TSP) and HTLV-associated myelopathy (HAM). Despite the increasing public health problem posed by HIV and HTLV infection, retroviruses have also made a fundamental contribution to biotechnology and molecular medicine. Following the seminal discovery in 1970 that **reverse transcriptase** (RT) catalyzes the flow of genetic information from RNA into DNA (1, 2), this enzyme for **cdNA** synthesis has found widespread use in producing medically important [recombinant proteins](#), such as **cytokines**, [growth factors](#), and [hormones](#). Furthermore, since 1989, molecularly engineered retroviral vectors are increasingly used as delivery vehicles for gene therapy. A recent and ambitious application of retroviral vectors involves the proposal to introduce attenuated strains of HIV as live viral vaccines to control the spread of HIV infection and AIDS.

**Table 1. Retroviral Genera**

| Retrovirus<br>Abbreviation | Full Name                    | Genus            |
|----------------------------|------------------------------|------------------|
| MMTV                       | Mouse mammary tumor virus    | B-type           |
| ASLV                       | Avian sarcoma-leukosis virus | Avian C-type     |
| ALV                        | Avian leukosis virus         | Avian C-type     |
| RSV                        | Rous sarcoma virus           | Avian C-type     |
| FeLV                       | Feline leukemia virus        | Mammalian C-type |
| GALV                       | Gibbon ape leukemia virus    | Mammalian C-type |
| MLV                        | Murine leukemia virus        | Mammalian C-type |
| REV                        | Reticuloendotheliosis virus  | Mammalian C-type |
| SNV                        | Spleen necrosis virus        | Mammalian C-type |
| MPMV                       | Mason–Pfizer monkey virus    | D-type           |
| SMRV                       | Squirrel monkey retrovirus   | D-type           |
| BLV                        | Bovine leukemia virus        | HTLV-BLV         |
| HTLV                       | Human T-cell leukemia virus  | HTLV-BLV         |

|      |                                      |              |
|------|--------------------------------------|--------------|
| BIV  | Bovine immunodeficiency virus        | Lentivirus   |
| CAEV | Caprine arthritis-encephalitis virus | Lentivirus   |
| EIAV | Equine infectious anemia virus       | Lentivirus   |
| FIV  | Feline immunodeficiency virus        | Lentivirus   |
| HIV  | Human immunodeficiency virus         | Lentivirus   |
| SIV  | Simian immunodeficiency virus        | Lentivirus   |
| HSRV | Human spuma retrovirus               | Spumavirus   |
| SFV  | Simian foamy virus                   | Spumavirus   |
| CSRV | Corn snake retrovirus                |              |
| SnRV | Snakehead fish retrovirus            |              |
| VRV  | Viper retrovirus                     | Unclassified |
| WDSV | Walleye dermal sarcoma virus         |              |

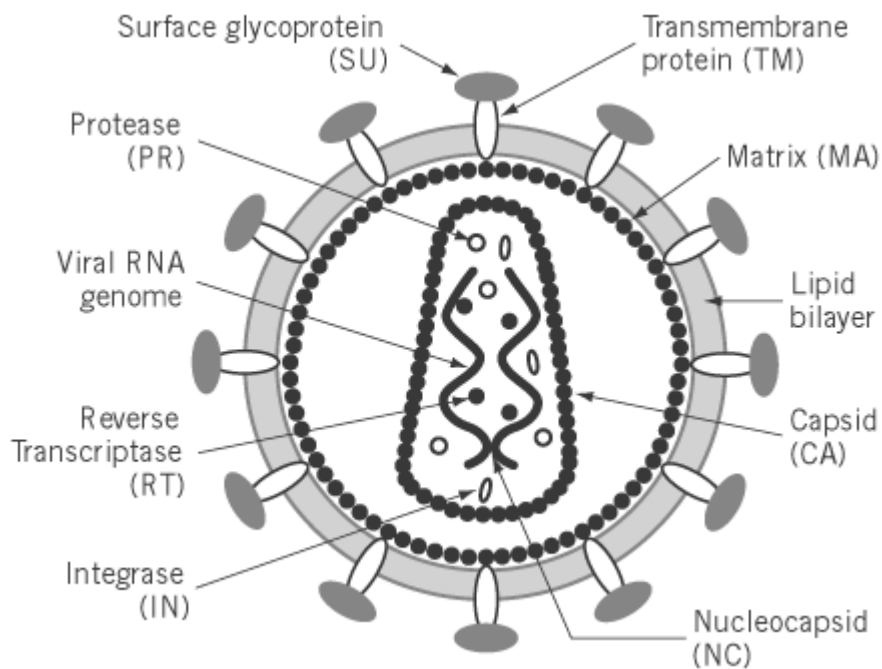
---

## 1. General features of retroviruses

### 1.1. Structure, Classification, and Genomic Organization

Despite significant differences in genetic complexity, host range, and interaction with their host, retroviruses share a common virion structure, a schematic version of which is illustrated in Fig. 1. Virions are typically 80 to 100 nm in diameter, contain an electron-dense core, and derive their lipid bilayer envelope from the host during budding. The virion surface is decorated with the envelope (*env*) gene products, comprised of a large surface **glycoprotein** (SU) (3) connected through [disulfide bond](#) linkages to a smaller **transmembrane** protein (TM). Matrix protein (MA), the amino-terminal product of the *gag* precursor [polyprotein](#), is located between the virion envelope and capsid, and it possibly interacts with the TM subunit. Consistent with the notion of membrane association, the amino termini of MA proteins are often modified by the addition of a fatty acid, most commonly a **myristoyl group**. A second *gag* product, *capsid protein* (CA), forms the major internal structure of the virion, the core shell. Although the core shell assumes a highly ordered structure, the manner in which this is accomplished through packing of CA monomers remains to be established. An electron-dense virion core harbors two identical copies of the **positive-stranded** RNA [genome](#), **hybridized** to a host-derived **tRNA** replication primer at the *primer binding site* (PBS) and likely coated over its entire length with the nucleocapsid protein (NC). The *pol*-coded enzymes protease (PR), RT, and [integrase](#) (IN) constitute additional core-associated proteins. Immediately following virion entry, RT and IN are required to synthesize and integrate, respectively, the double-stranded DNA provirus into the host genome, and PR serves to process multi-domain *gag* and *gag-pol* precursor polypeptides during virion assembly at later stages in the retroviral life cycle. In addition to these “generic” core-associated proteins, virion cores of several lentiviruses contain a fourth *pol*-derived enzyme, *deoxyuridine triphosphatase* (dUTPase), whereas those of primate lentiviruses harbor the accessory virus-coded proteins Vpr and Vpx.

**Figure 1.** Schematic representation of a typical retrovirus virion, illustrating the disposition of structural and enzymatic proteins. Within the virion core, the dimeric RNA genome is most likely completely coated with the NC product of the *gag* gene. The two-letter abbreviation for retroviral proteins is indicated in parentheses.



Before nucleotide sequence data defining genomic organization, [electron microscopy](#) was used to distinguish retroviruses according to morphological features, leading to their designation as type *A*, *B*, *C* or *D*.

1. *A* particles are strictly intracellular nucleocapsid structures (containing a dimeric RNA genome, *gag* and *gag-pol* precursors), examples of which are the intracisternal *A* particles (IAP) of several rodent species and the immature form of mouse mammary tumor virus (MMTV). Fully formed immature cores of the *B*-type, *D*-type and spumaretroviruses are included in this group.
2. *B* particles are defined as the extracellularly enveloped form of MMTV. Immediately following budding, *B* particles mature to form condensed acentric nucleocapsid, which contrasts with the hollow appearance of the immature intracellular form. A prominent feature of *B* particles is decoration of their surface with the SU product of the *env* gene.
3. *C* particles comprise the simple avian and mammalian viruses, examples of which include [Rous sarcoma virus \(RSV\)](#) and *Moloney murine leukemia virus* (MOMLV). Such particles can be visualized as crescent-shaped patches on the cell membrane during budding, that is, fully formed intracellular structures are not generally detected. An electron-dense core morphology is also assumed by *C* particles immediately after budding, reflecting **proteolytic** processing of structural precursor polyproteins by the virally coded PR. Unlike their *B*-type counterparts, surface projections are barely visible on *C*-type viruses.
4. *D* particles are exemplified by *Mason–Pfizer monkey virus* (MPMV). These share many features of *B*-type viruses (including an intracellular nucleocapsid) but differ primarily in that their surface projections are less prominent.

More recently, seven genera of retroviruses have been defined, based on combining of nucleotide sequence homology and genomic complexity. Examples of these are provided in Table [1](#). Although *B*-type and *D*-type viruses remain unchanged, now those of *C*-type have been subdivided into avian and mammalian retroviruses. Both classes of *C*-type viruses contain “simple” genomes, that is, they encode only structural and enzymatic proteins from their *gag*, *pol*, and *env* **open reading frames**. Distinct from these are several groups of “complex” lentiviruses that encode as many as six accessory proteins on their genomes, in addition to *gag*, *pol*, and *env* products:

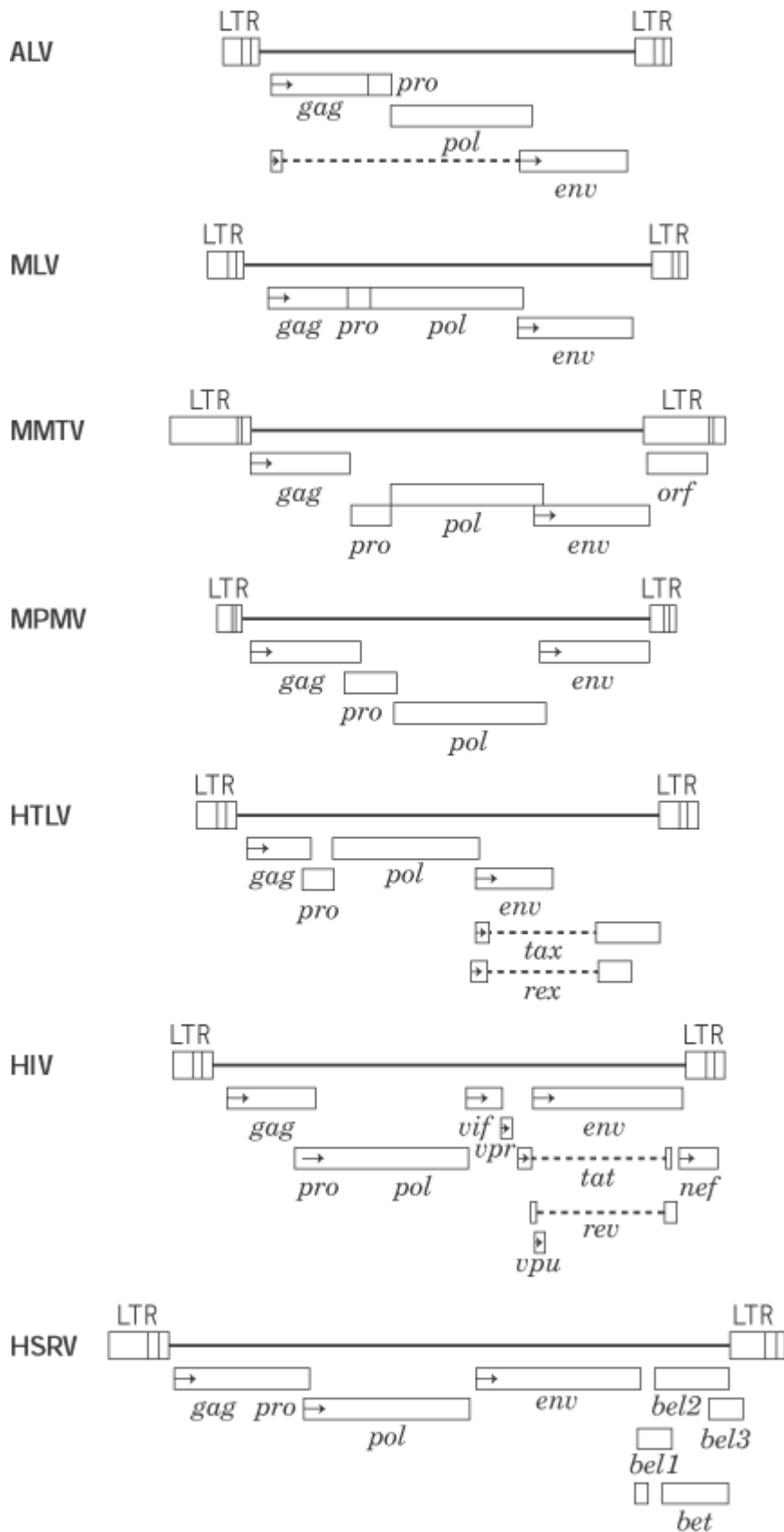
1. *Lentiviruses* are a diverse group of mammalian retroviruses responsible for neurological and

immunological diseases but not directly implicated in malignancies.

2. *HTLV-BLV* genus includes a small number of complex retroviruses associated with B- and T- cell lymphomas, in addition to certain neurological disorders.
3. *Spumaviruses* (alternatively designated foamy viruses because infected cells display highly vacuolated or “foamy” syncytia) comprise a class of complex retroviruses of simian, human, and feline origin with which no disease has currently been associated. A hallmark of spumaviruses is the length of their genomes (>11,000 nucleotides). Additional features of their life cycle that distinguish spumaviruses as a unique class of retroviruses include (i) [alternative splicing](#) mechanisms for production of an independent **mRNA** encoding enzymatic proteins of the *pol* gene; (ii), the presence of significant amounts of double-stranded proviral DNA in the virion, and (iii) an unusual site of transport for the preintegrative nucleoprotein complex.

The organizations of several well-characterized retroviral genomes are indicated in Fig. 2. All proviral genomes are characterized by terminal noncoding regions ([long-terminal repeats](#) or LTR) that contain **cis-acting** functions necessary for replication, between which lie genes for proteins that have structural (*gag* and *env* open reading frames) and enzymatic roles (*pol* open reading frame). The same figure also highlights considerable differences in genomic complexity among different retrovirus genera. The simple *-gag-pol-env-* organization of ALV and MLV can be contrasted with HIV-1, which contains at least six accessory proteins derived via multiple RNA splicing events, the best characterized of which are the **transactivators** Tat and Rev. Through a complex series of interactions with highly ordered structures assumed by the viral RNA genome (the TAR loop with Tat and *rev*-responsive element or RRE with Rev), these proteins play important roles in controlling [transcription](#) of the viral genome and cellular localization of the RNA transcripts, respectively. Potential roles for additional HIV-1 proteins have included steps in virion maturation (Vif), nuclear localization of the viral preintegrative complex (Vpr), enhancement of virion release (Vpu) and CD4 receptor downregulation (NEF). Although controversy still surrounds the biological role(s) of these factors, several studies have shown that viruses within which Nef has been deleted or prematurely terminated are attenuated (i.e., less virulent), suggesting an important role for this protein in pathogenicity.

**Figure 2.** Coding regions of simple and complex retrovirus genomes. The provirus is presented in each case and is bracketed by LTR. Each box under this represents an open reading frame, within which horizontal arrows denote the initiation codon. Dotted lines designate coding regions joined by splicing events. The retroviral protease (*pro*) can exist as a component of the *gag* or *pol* open reading frames or in an independent reading frame.



Counterparts to HIV Tat and Rev proteins have been detected in most, if not all, lentiviruses. In the HTLV-BLV class the Tat counterpart, Tax, fulfills a similar role of transcriptional activation but through an alternative mechanism. Rather than interacting with a control element located in the RNA

transcript, the target of Tax action is a segment of the retroviral promoter that contains three 21-bp repeat elements. Current evidence suggests that Tax does not interact directly with these motifs but rather with a member of the [cyclic AMP](#) response (CREB) family, which shares the same recognition element on the LTR promoter, thereby stimulating host transcriptional machinery. In contrast to Tax, HTLV Rex functions similar to the lentiviral proteins, in this case interacting with its RxRE target sequence on nascent viral RNA to control the levels of singly and multiply spliced transcripts in the cytoplasm. Of the several *bel* genes (between *env* and LTR) described for spumaviruses, the best characterized is the *belI* product, which functions analogously to HTLV Tax.

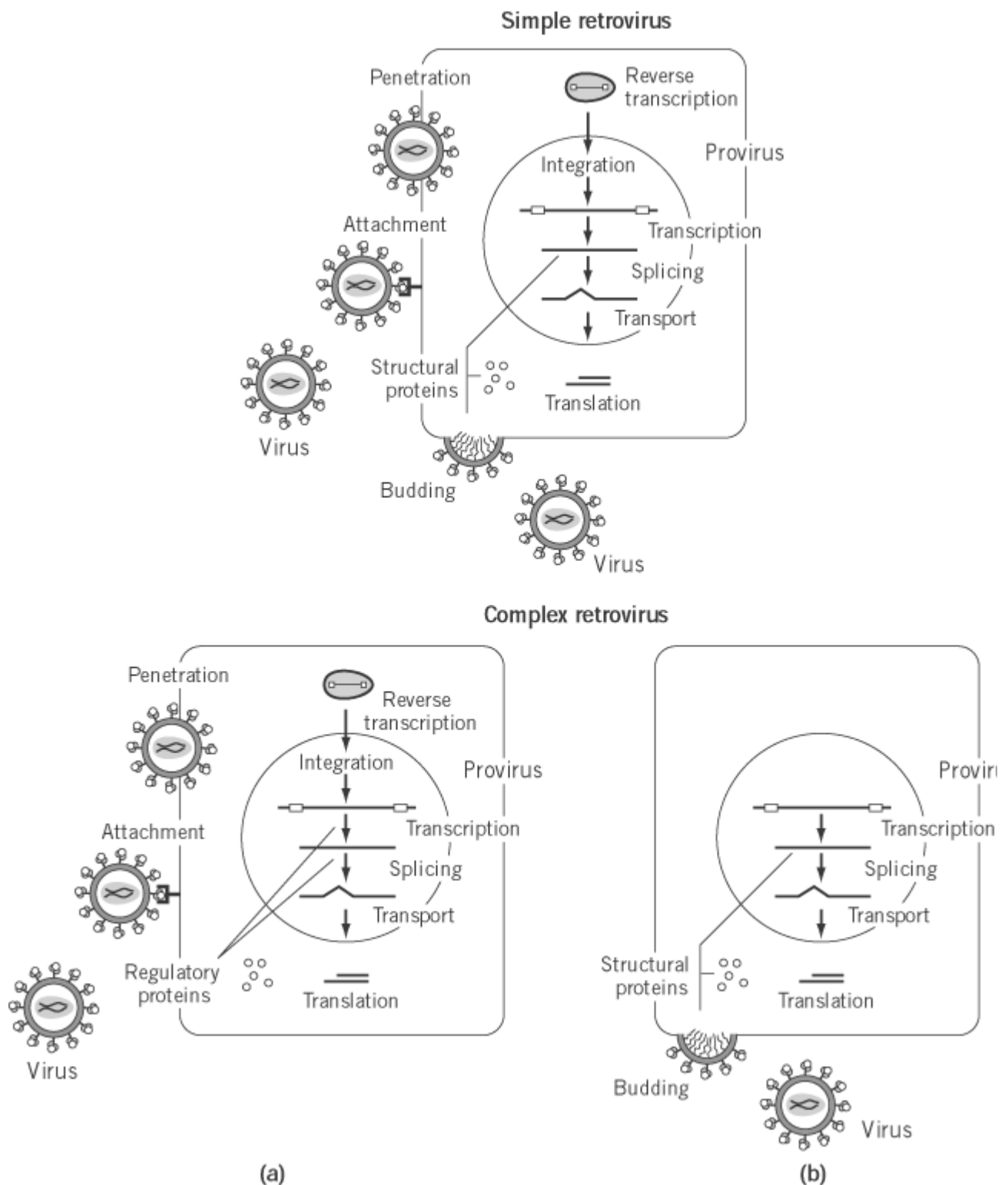
## 2. The retroviral life cycle

The life cycles of simple and complex retroviruses are indicated in Fig. 3. Although both follow a similar pathway leading to establishment of the integrated provirus, they differ in subsequent events leading to maturation, where complex retroviruses have established control mechanisms for maintaining virion RNA within the nucleus. Salient features of the life cycle include

- attachment to specific **receptor(s)** on the surface of the target cell
- penetrating and uncoating to release virion cores into the cytoplasm
- converting the single-stranded RNA genome into double-stranded proviral DNA
- transporting the “preintegrative” nucleoprotein intermediate into the [nucleus](#)
- integrating proviral DNA into the host genome
- synthesizing viral RNA from the integrated provirus
- transporting viral RNA to the **cytoplasm**
- synthesizing polyprotein precursors that encode structural and enzymatic proteins
- assembling nucleic acid and protein components at the cell membrane and budding
- maturing infectious virions

**Figure 3.** Life cycles of simple and complex retroviruses. For simple retroviruses, the life cycle can be envisioned as two discrete sets of events, (1) establishment of the integrated provirus and (2) expression of the provirus, followed by virion maturation. Although complex retroviruses establish the provirus similarly, then regulatory proteins (Tat, Rev, Rex) are synthesized from multiply spliced RNA (**a**) which aid in transporting unspliced and singly-spliced messages to the cytoplasm. Thereafter, (**b**) events follow the same pattern as simple retroviruses.





## 2.1. Attachment

The interaction of virions with a cell-surface receptor is mediated by the highly glycosylated SU component of the retroviral *env* gene. One of the most extensively characterized of these is the HIV receptor CD4, which functions in the host cell to aid in signaling the interaction of T-helper cells with **antigen-presenting** cells. Receptors closely resembling transport systems have been identified for both MLV (Rec-1 and Ram-1) and GALV (Glv-1). Tva, the subgroup A ALSV receptor, is structurally distinct and bears some resemblance to the **low-density lipoprotein** receptor. Candidate receptors for BLV and FIV have been identified but require further characterization. Although the

prevailing dogma envisages a single class of receptor-mediated attachment, the observation that murine model systems bearing the human CD4 receptor support attachment of HIV, but not the subsequent events, means that a second receptor may be involved. This notion has recently been proven experimentally by the identification of chemokine receptors CXCR-4 (fusin) and CCR5 as important coreceptors in HIV infection (4). Although CXCR-4 mediates infection of CD4<sup>+</sup> cells by T-cell line adapted HIV-1 strains, CCR5 allows entry of primary/macrophage tropic strains. Other macrophage-tropic and dual-tropic strains exploit the additional chemokine receptors CCR3 and CCR2b.

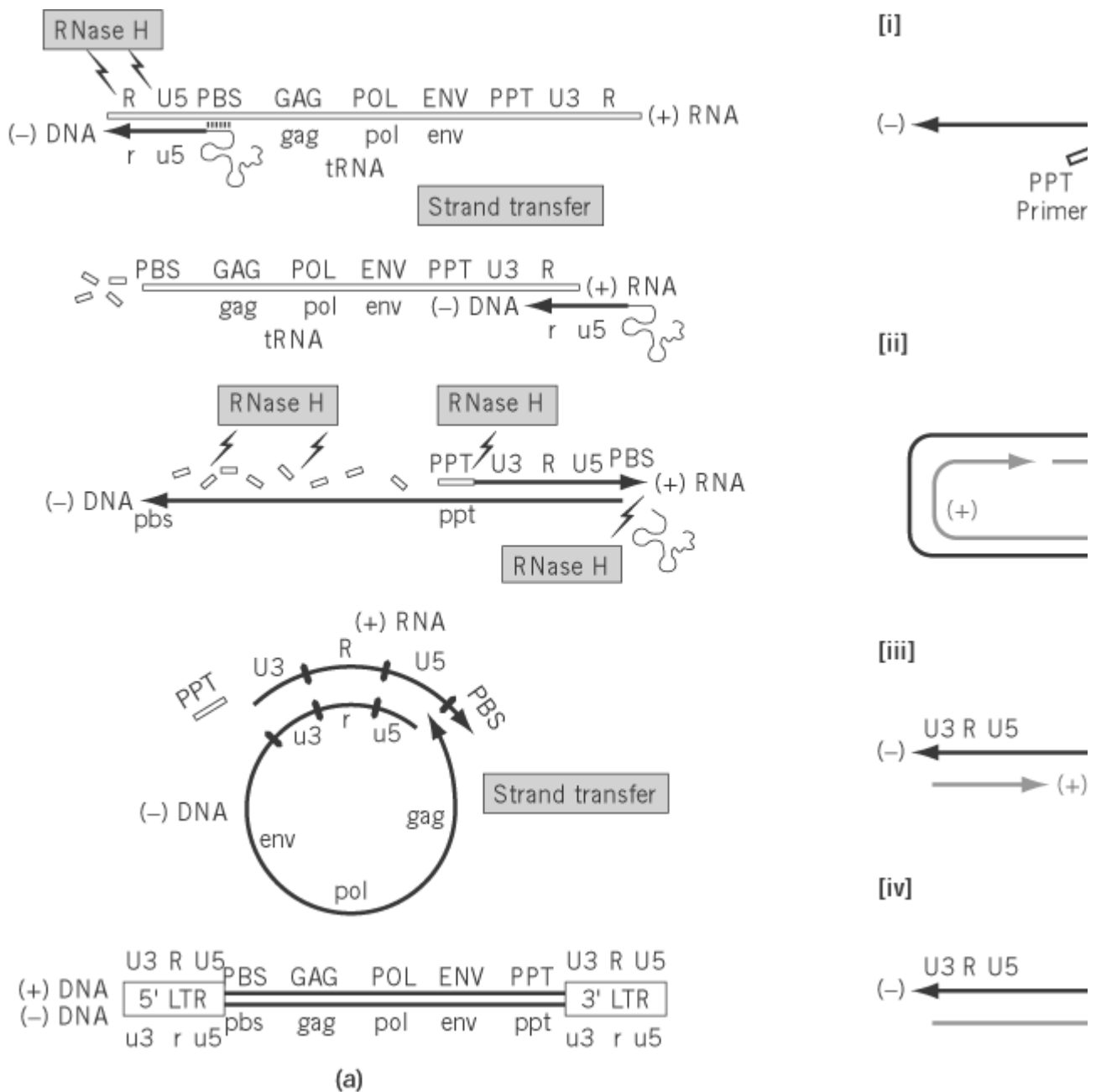
## 2.2. Fusion and Uncoating

Penetration and release of the virion nucleoprotein core occurs by one of two distinct mechanisms: (i) fusion between the viral envelope with the cell membrane or (ii) receptor-mediated [endocytosis](#) accompanied by fusion with the **endosomal** membrane. Such events in retroviruses are pH-independent and are mediated by **hydrophobic** amino acids at the amino terminus of the envelope transmembrane (TM) component. Studies with HIV provide compelling evidence that a cellular peptidyl-prolyl isomerase, [cyclophilin A](#) (CyPA), participates in events between membrane fusion and initiation of reverse transcription (5). CyPA is incorporated into virions late in the infection cycle through direct interaction with the CA component of *gag* precursor polyprotein, suggesting that it may participate directly in uncoating. The virion structure presented in Fig. 1 predicts that MA provides an outer shell and would not be associated with core particles after uncoating. However, this is not the case with HIV, where a small proportion of MA molecules remain core-associated, possibly to assist in transporting the preintegrative nucleoprotein complex into the nucleus.

## 2.3. Proviral DNA Synthesis

Proviral DNA synthesis occurs in virion cores whose structure is relaxed to allow access to deoxynucleoside triphosphates. Through the elegant series of events depicted in Fig. 4, positively-stranded viral RNA is converted into double-stranded proviral DNA by a single virus-coded enzyme, RT. However, because signals that determine initiation and termination of transcription by the host **RNA polymerase** reside within the integrated provirus, these by default are absent from the RNA genome of an invading virus. Therefore, a mechanism must be established for regenerating them in the course of the reverse transcription cycle.

**Figure 4. (a)** Model for synthesizing proviral DNA (closed lines) from the single-stranded RNA genome (open lines). First, the RNA genome that contains the coding information for protein synthesis is designated (+) RNA, from which a complementary (-) DNA strand is synthesized. Subsequent use of the (-) DNA template generates the (+) DNA strand. Between these events, the (+) RNA genome is degraded. The LTR elements are only partially represented at either terminus of the RNA genome and are regenerated following two discontinuous (+) strand syntheses in lentiviruses and termination of 3' PPT-initiated synthesis at the center of the RNA genome. In Step (1), (+) strand synthesis initiates from both the 3' and 5' ends of the RNA genome. PPT-initiated synthesis terminates within the tRNA template, and the tRNA is released. Step (2), second-strand transfer, involves the transfer of the 3' end of the (+) strand to the 5' end of the (-) strand, creating a double-stranded DNA intermediate with a gap. Step (3), bidirectional strand displacement synthesis via elongation of the (-) strand and the (+) strand, fills the gap. The (+) strand synthesis is terminated after CTS-controlled strand displacement at the center of the retroviral genome.



First, **minus-strand** DNA synthesis is initiated from a host-derived tRNA that is hybridized via its 3'-terminal 18 nucleotides to the primer binding site (PBS) immediately downstream from U5. Examples of tRNA used to initiate replication include tRNA<sup>Trp</sup> in RSV, tRNA<sup>Pro</sup> in MLV, and tRNA<sup>Lys,3</sup> in HIV. As U5 and R sequences are copied, RT-associated [ribonuclease H](#) (RNase H) activity degrades RNA of the RNA–DNA replication intermediate, thereby making R sequences of minus DNA accessible for hybridization to their complement at the 3' end of the RNA genome. Through a process designated “*strand transfer*” or “*switching*”, nascent minus-strand DNA is relocated to the 3' end of the same (intramolecular) or a different RNA genome (intermolecular). Continued RNA-dependent DNA synthesis creates a first copy of the LTR-U3-R-U5- hallmark on minus-strand DNA, and RNase H activity continues to degrade RNA of the replicative intermediate. The exception to this is a purine-rich RNA sequence located near the 3' end of the RNA genome, designated the *polypurine tract* (PPT), which resists hydrolysis to serve as a primer for plus-strand DNA synthesis. Plus-strand, DNA-dependent DNA synthesis ensues, using newly synthesized minus-strand DNA as a template, up to and including the first 18 nucleotides of the tRNA replication

primer (i.e., this is used as both [primer](#) and [template](#) during replication). At this step, the replication machinery transiently pauses at a methylated base of the tRNA, allowing RNase H-mediated removal of the intact minus-strand primer. As a consequence, plus-strand PBS sequences are made available for hybridization to their minus-strand counterparts to complete a second, intramolecular strand switch (concomitant with which the PPT primer is removed). Subsequent bidirectional DNA-dependent DNA synthesis generates a double-stranded proviral DNA flanked by complete copies of the LTR.

Although this model holds true for many retroviruses, several lentiviruses and yeast [retrotransposons](#) harbor a central plus-strand discontinuity, that reflects use of a second primer in the middle of the genome (the central or cPPT) (6). Although the advantage of an additional PPT is unclear, its alteration on the HIV genome has severe consequences for replication. The reverse transcription cycle that accounts for both 3'- and cPPT-initiated plus-strand synthesis is presented in Figure 4 (b). Documentation of a plus-strand discontinuity at the center of proviral DNA imposes an additional constraint on these retroviruses, namely, to provide a mechanism for terminating plus-strand synthesis immediately downstream from the cPPT. Preliminary evidence for this has been provided in HIV, where localized distortion of double-helical DNA geometry (“bending”), causes pausing and dissociation of the replication machinery. The DNA sequence and structure implicated in such events has been designated the *central termination sequence* (CTS). Finally, it should be pointed out that both copies of RNA in the virion core are used to synthesize a single proviral DNA, which indicates that relocation of nascent DNA between genome templates occurs frequently during replication. Because the retroviral polymerase is inherently more error-prone than host **DNA polymerases**, such recombination has the potential to enhance the rate of [evolution](#), in addition to allowing RT to bypass lesions in the RNA genome.

#### 2.4. Transport of Proviral DNA

Following DNA synthesis, the nucleoprotein complex composed of linear DNA, the core protein CA, RT, and possibly NC, is transported through the nuclear membrane (see **Nuclear import, export**). In some cases (eg, MLV), breakdown of the nuclear membrane at cell division is necessary to support these events, evidenced by the observation that inhibiting the onset of **mitosis** delays proviral DNA integration. An alternative strategy is employed by lentiviruses, which do not require the infected cell to enter mitosis. Nonenzymatic HIV-1 proteins, such as MA-derived peptides and Vpr, have been detected as constituents of the nucleoprotein complex, from which it has been postulated that nuclear localization signals of these accessory proteins assume a supportive role in transport (7).

Within the nucleus of infected cells, two variants of circular proviral DNA have been detected that contain either one or two copies of the LTR. Originally proposed as integration intermediates, it has been unequivocally demonstrated that these represent the products of aberrant reverse transcription that have no biological function.

#### 2.5. Integration

Integration is a mechanism whereby proviral DNA is stably inserted into the host genome colinearly with the way it is synthesized by RT. This process is catalyzed by the integrase (IN) component of the *pol* open reading frame and is independent of host factors and external energy sources (8). The immediate LTR termini contain short, imperfect, repeat sequences (**att sites**) that as shown by mutagenesis studies, constitute essential *cis*-acting sequences recognized by IN. In the absence of an exhaustive analysis, there is insufficient evidence to suggest preferential sites of integration in the host genome. During integration (the catalytic mechanism for which is described later), proviral DNA sequences are characterized by loss of two nucleotides at either end, whereas target sequences immediately adjacent to the integrated proviruses contain a short (4 to 6 bp) duplication, whose size is characteristic of a particular retrovirus. Examples of this are a 5-bp duplication in HIV and a 6-bp duplication in ALSV.

#### 2.6. Viral RNA Synthesis

In essence, integration completes the first half of the retroviral life cycle, after which viral DNA is preserved in the **germ line** and transcribed from the LTR **promoter** by host-coded **RNA polymerase II**. Duplication of terminal sequences in the provirus implies that transcriptional signals of equal strength reside in both the 5'- and 3'-LTR. In fact, only the former is productively used, suggesting that regulatory signals downstream of the 5'-LTR are involved. The full-length genomic transcript initiates at the upstream U3/R junction with a "**5'-capped**" guanosine nucleoside and is **polyadenylated** immediately after the downstream R/U5 junction. In addition to LTR-derived RNA, spumaviruses initiate RNA synthesis from a site upstream of the *bel-1* gene, and additional transcriptional initiation sites have also been identified for MMTV. A later section deals with host-coded factors active on the LTR promoter and the role of the transactivators in augmenting gene expression.

## 2.7. Splicing and Transport of Viral RNA

In all retroviruses, the intact transcript is both the RNA genome and messenger RNA for synthesizing *gag* and *gag-pol* precursor polyproteins, whereas a singly spliced subgenomic transcript encodes the *env* precursor. Different mechanisms have evolved to maintain a balance between unspliced and spliced mRNA.

Simple retroviruses exploit a *cis*-acting sequence, designated the *negative regulator of splicing* (NRS), which interacts with components of the cellular machinery to reduce splicing efficiency. In contrast, to maintain this balance, complex retroviruses have evolved more elaborate systems, involving the accessory proteins Rex (HTLV) and Rev (HIV, Fig. 2), which are derived via multiple splicing shortly after integration (9). These accessory proteins interact with signals located within either the viral *env* mRNA (Rev) or the R/U3 region of 3' LTR RNA (RxRE), which otherwise provide a barrier to the export of transcripts harboring such sequences that are characterized by complex secondary structures and designated Rex- or Rev-responsive element (RxRE and RRE, respectively). Additional interactions of this RNA-protein complex with components of the cellular transport machinery facilitate transfer of intact and singly spliced mRNA to the cytoplasm. As outlined in Fig. 3, low-level synthesis of transactivator proteins (Tat, Rev, Rex), shortly after infection or after transcription is activated, has the consequence that unspliced and singly spliced RNA are maintained in the nucleus. However, a threshold level of Rev or Rex is eventually achieved, so that intact and singly spliced RNA are efficiently exported to the cytoplasm and made available to the translational machinery for events leading to virion assembly. Although such elaborate mechanisms are absent in simple retroviruses, RNA containing introns (ie, unspliced) must likewise be transported to the cytoplasm, implying an alternative mechanism. Preliminary data from *D*-type MPMV and RSV suggest that this is facilitated by the interaction of a viral RNA sequence, designated the *constitutive transport element* (CTE), with host-coded factors.

## 2.8. Protein Synthesis

Spliced mRNA provide the template for synthesizing TM and SU components of the viral envelope (in this case on rough **endoplasmic reticulum-associated polyribosomes**), and accessory proteins and **oncogenes**, whereas intact, unspliced RNA fulfills two roles. A portion of these become genomes, and another directs synthesis of virion structural (*gag*) and enzymatic proteins (*pol*). The mechanism by which *gag* and *pol* components are synthesized is a good example of genetic "streamlining" in retroviruses, where different mechanisms have evolved to allow high-level synthesis of structural proteins and low-level synthesis of virion enzymes from the same mRNA (10). MLV achieves this via *transitional suppression*, a mechanism more often associated with **prokaryotic** organisms. A UAG **termination codon** is located between the *gag* and *pol* open reading frames, and to the larger extent terminates synthesis of the *gag* precursor. In about 5% of cases, however, the host translational machinery inserts a **glutamine** residue at this position, thereby overriding termination to synthesize a *gag-pol* precursor polyprotein. Because the *gag* and *pol* genes in most other retroviruses are not in the same reading frame, an alternative method must be employed. This involves ribosomal **frameshifting**, where slippage of the translating ribosome allows reading the base before (-1 frameshift) or after (+1 frameshift) together with the *gag* stop codon,

thereby creating the triplet for an amino acid. A single frameshift event is required to cotranslate the *gag* and *pol* genes in viruses, such as ALV and HIV, whereas two frameshift events are necessary when the PR open reading frame is out of frame with both of these (eg, MMTV, HTLV, and MPMV). In a frequently recurring theme, spumaviruses, the exception to these models of *pol* gene synthesis, have evolved an independently spliced mRNA for the *pol* gene that has its own translational [initiation codon](#) (11).

## 2.9. Assembly and Budding

Following their synthesis, *gag*, *pol*, and *env* gene products, together with a dimeric RNA genome, assemble in an ordered fashion at the cytoplasmic face of the cell membrane, gradually becoming enveloped in forming an immature virion. Membrane targeting of *gag* and *gag-pol* precursors is aided by N-terminal fatty acid modifications of MA that locate this protein to its ultimate fate as the inner shell of the virion. Lack of myristylation on *gag* precursors of ALSV and some lentiviruses suggests participation of additional, as yet unidentified, components in membrane targeting. An interaction involving a *cis*-acting element located at the 5' end of viral RNA (the *encapsidation* or Y sequence) and *gag*-coded NC is critical for packaging the genome. Based on the general nucleic acid binding properties of NC and the amount present in virions, it is most likely that this protein interacts with and “coats” the RNA genome. CA, located between MA and NC, is appropriately positioned to provide the virion core shell after proteolytic maturation. *Gag-pol* precursors are most likely oriented at the membrane similarly. They often serve an additional role of sequestering (via the *pol*-coded RT) the appropriate host tRNA isoacceptor required to prime reverse transcription events after subsequent infection. Alternatively, complementarity between viral PBS sequences and the 3' terminus of the cognate replication primer may achieve the same goal. Only low-level DNA synthesis activity has been reported for RT embedded within the ~170kDa *gag-pol* precursor, which may reflect a control mechanism preventing premature reverse transcription before virion budding. Spumaviruses are the exception to this. Virions contain significant amounts of proviral DNA, which may lie in the observation that the *pol* gene is expressed independently of *gag*.

## 2.10. Virion Maturation

Immature cores that contain the appropriate protein and nucleic acid components gradually bud from the cell, together with a portion of the plasma membrane concentrated with the Env gene products. Concurrent with or after budding, the virally-coded protease is activated and is responsible for proteolytic maturation of *gag* and *pol* products. Maturation is accompanied by significant morphological changes. The most notable of these is formation of an electron-dense virion core from “doughnut”-shaped precursor polyproteins that previously concentrated at the inner surface of the budding particle. At this stage, virions are considered capable of initiating a new cycle of infection by interacting with cell-surface receptors.

## Bibliography

1. H. Temin and S. Mizutani (1970) *Nature* **226**, 1211–1213.
2. D. Baltimore (1970) *Nature* **226**, 1209–1211.
3. J. Leis, D. Baltimore, J. M. Bishop, J. Coffin, E. Fleissner, S. P. Goff, S. Oroszlan, H. Robinson, A. M. Skalka, H. M. Temin, and V. Vogt (1988) *J. Virol.* **62**, 1808–1809.
4. Y. Feng, C. C. Broder, P. E. Kennedy, and E. A. Berger (1996) *Science* **272**, 872–877.
5. J. Luban, K. A. Bossolt, E. K. Franke, G. V. Kalpana, and S. P. Goff (1993) *Cell*, **73**, 1067–1078.
6. P. Charneau, M. Alizon, and F. Clavel (1992) *J. Virol.* **66**, 2814–2820.
7. A. G. Burkinskaya, A. Ghorpade, N. K. Heinzinger, T. E. Smithgall, R. E. Lerwis, and M. Stevenson (1996) *Proc. Natl. Acad. Sci. USA* **93**, 367–371.
8. M. D. Andrade and A. M. Skalka (1996) *J. Biol. Chem.* **271**, 19633–19636.
9. J. Karn, M. G. Gait, M. J. Churcher, D. A. Mann, I. Mikaelian, and C. Pritchard (1994) In *RNA-Protein Interactions* (K. Nagai and I. W. Mattaj eds.), Oxford University Press, New York, pp. 192–220.

10. T. Jacks (1990) In *Retroviruses. Strategies of Replication* (R. Swanstrom and P. K. Vogt, eds.) Springer-Verlag, New York, pp. 93–124.
11. J. Enssle, I. Jordan, B. Mauer, and A. Rethwilm (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4137–4141.

### Suggestions for Further Reading

12. B. R. Cullen (1993) *Human Retroviruses*, Oxford University Press, New York.
13. J. A. Levy (1994) *The Retroviridae*, Vol. 1–4, Plenum, New York.
14. B. N. Fields, D. M. Knipe, and P. M. Howley (1996) In *Fundamental Virology*, 3rd ed., Lippincott-Raven, Philadelphia, New York, pp. 763–916.
15. J. M. Coffin, S. H. Hughes, and H. Varmus (1997) *Retroviruses*, Cold Spring Laboratory Press, Cold Spring Harbor, New York.

## Reverse Translation

**Translation** is the biosynthesis of a **protein** from a **messenger RNA** template on **ribosomes**. Reverse translation is not a biological process. Instead, it is inferring DNA sequence from the amino acid sequence of a protein. Reverse translation is often employed to make a **hybridization** probe or a **PCR** primer used to **clone** the **gene** encoding the protein of interest (see **Cloning**) (1-3).

In most cases, the amino acid sequence of a protein is directly inferred from the DNA or mRNA sequence coding for the protein because each nucleic acid triplet (**codon**) specifies either a single amino acid or a termination signal (see **Genetic Code**). The converse, determining the DNA or mRNA sequence coding for a specific amino acid, is more complex because the genetic code is “degenerate” (see **Degeneracy of the Genetic Code**). In the nuclear DNA of **eukaryotes**, 61 codons specify 20 amino acids, so many amino acids are coded by more than one codon. This means that reverse translation of a protein does not produce a single nucleotide sequence. Instead, it results in a population of different sequences that, if translated, would all code for the same amino acid sequence. To identify the actual genomic sequence that codes for the protein *in vivo*, it is necessary to clone and sequence the gene for the protein. The first step in cloning is to synthesize a mixture of oligonucleotides (oligos) that corresponds to all of the potential coding sequences determined by reverse translation. This pool of oligos is used as a “degenerate” (mixed) hybridization probe to isolate the corresponding DNA or cDNA clone from a library. Alternatively, reverse translation is used to design two sets of “degenerate” PCR primers to amplify the gene from genomic DNA.

When using reverse translation to design degenerate oligos for gene cloning, several factors must be taken into account. A 14-base oligo is sufficiently long to identify the gene of interest specifically, but the five-residue stretch of protein that is reverse translated to produce this oligo must be chosen carefully. The more protein sequence that is known, the easier it is to find an appropriate amino acid stretch to reverse translate. Because **serine**, **leucine**, and **arginine** are each coded by six different codons, these residues should be avoided. Protein sequences containing **tryptophan** and **methionine** residues are preferred, because they are each coded by only one triplet codon. Fewer different oligo sequences are needed to cover all reverse translation possibilities if less “degenerate” amino acids are chosen, and a less complex set of oligos makes a more efficient probe or PCR primer. Different organisms preferentially use particular codons to specify amino acids, and this **codon usage bias** should also be taken into account when designing oligos by reverse translation. Additionally, computer programs are available to help design synthetic genes and degenerate probes and primers

by using reverse translation (4).

Once many of the large-scale genome sequencing projects are complete (eg, the human, *Caenorhabditis elegans*, and *Arabidopsis* sequencing projects), and with the ever increasing number of [expressed sequence tags](#) (ESTs) available, reverse translation most frequently will be used to isolate genes from organisms where little sequence data is available. For well-studied organisms, database searches (where a computer program compares known protein sequences and translated nucleotide sequences to look for similarities) will replace the need to reverse translate and clone to determine the nucleotide sequence.

#### Bibliography

1. S. C. Vendeland, M. A. Beilstein, J. Y. Yeh, W. Ream, and P.D. Whanger (1995) *Proc. Natl. Acad. Sci. USA* **92**, 8749–8753.
2. N. Dekeyzer, Y. Engelborghs, and G. Volckaert (1994) *Protein Eng.* **7**, 125–130.
3. A. Landa, A. Rojo-Dominguez, L. Jimenez, and D. A. Fernandez-Velasco (1997) *Eur. J. Biochem.* **247**, 348–355.
4. T. Tamura, S. R. Holbrook, and S. H. Kim (1991) *Biotechniques* **10**, 782–784.

### Reversed-Phase Chromatography

In general, any **chromatographic** separation, in which the mobile phase is more polar than the stationary phase and polar solutes move with this mobile phase rather than remain with the less polar stationary phase, is called reversed-phase chromatography. More specifically, reversed-phase chromatography (particularly in **HPLC**) refers to those procedures in which the stationary phase consists of silica beads that are completely covered with covalently bound hydrophobic *n*-alkyl chains (eg, octyl, C8, or octadecyl, C18). In reversed-phase chromatography, the adsorption of a solute to the reversed-phase matrix is generally driven by **hydrophobic** interactions, and the solutes migrate in decreasing order of net charge, extent of ionization, and **hydrogen-bonding** capabilities.

Because of its excellent resolution, experimental ease, high recoveries and excellent reproducibility, reversed-phase HPLC now plays a critical role in the biochemical/biomedical studies (see reviews 1, section V of Ref. 2, pp. 273–421, 3 and 4). Detailed practice is beyond the scope of this volume, and interested readers are directed to the reviews cited previously and to a number of excellent monographs (5, 6).

#### Bibliography

1. M. T. W. Hearn and M. I. Aguilar (1988) in *Modern Physical Methods in Biochemistry, Part B* (A. Neuberger and L. L. M. Van Deenen, eds.), Elsevier Science Amsterdam, pp. 107–142.
2. C. T. Mant and R. S. Hodges (eds.) (1991) *High-Performance Liquid Chromatography of Peptides and Proteins: Separation, Analysis, and Conformation*, CRC Press, Boca Raton.
3. M. T. W. Hearn (1989) in *Protein Purification: Principles, High Resolution Methods, and Applications* (J.-C. Janson, and L. Rydén, eds.), VCH, New York, pp. 175–206.
4. R. D. Shah and C. A. Maryanoff (1997) in *HPLC: Practical and Industrial Applications* (J. Swadesh, ed.), CRC Press, Boca Raton, pp. 111–169.
5. A. M. Krstulovic and P. R. Brown (1982) *Reversed-Phase High -Performance Liquid Chromatography: Theory, Practice, and Biomedical Applications*, Wiley, New York.



6. G. Szepesi (1992) *How to Use Reversed-Phase HPLC*, VCH, New York.

## Reversion, Revertant

A reversion is a [mutation](#) that restores the wild-type **phenotype** to a [mutant](#). A revertant is an individual who has a reversion mutation. Reversion assays are used to investigate the mechanism by which a [mutagen](#) acts. Reversional assays are powerful genetic techniques for studying mutations because revertants are usually easy to detect. Reversion assays rely on restoring a wild-type phenotype, and this gain of function can usually be selected for, as opposed to the loss of function that must usually be screened for (see **Selection**).

The most straightforward revertants are **genotypic** in which the original wild-type **gene** sequence is restored. The most widely used reversion assay was developed by Bruce Ames. Using various strains of *Salmonella typhimurium* with known mutations in the [his operon](#), the frequency of reversion to histidine **prototrophy** is measured after exposure to [mutagens](#). Because the mutations that cause the various examples of histidine **auxotrophy** are known ([base-pair substitutions](#) or [frameshift mutations](#)), the nature of the mutation that caused the reversion is known.

Phenotypic reversion also occurs by mutation at sites different from those of the original mutation. The second mutation can occur in the same gene, in which case it is a *second-site revertant*. In a [protein](#) gene product, the second amino acid change in the protein apparently corrects the first amino acid change. Compensating nucleotide changes are also known in [transfer RNA](#) molecules. If the second mutation occurs in a different gene, the phenomenon is known as **suppression**. Such mutations give valuable information about functional [protein–protein interactions](#).

## RGS Proteins

RGS proteins, **Regulators of G protein Signaling**, are a family of **membrane-associated multi-domain** proteins that modulate the activities of the [heterotrimeric G proteins](#). Although the first known RGS protein, Sst2p in yeast, was discovered in 1987 (1), the existence of the family and its mechanism of action was not known until 1996 (2). At least 30 mammalian RGS proteins have been identified at the level of [complementary DNA](#) sequences, about a dozen are known in *Caenorhabditis elegans*, and family members have been identified in several fungi, multiple animal phyla, and higher plants. RGS proteins in animals act on members of the G<sub>i</sub> and G<sub>q</sub> families of heterotrimeric [GTP-binding proteins](#), with varying degrees of specificity within this large group. Their specificity in fungi and plants is unknown, but it may be limited practically by the small number of Ga subunits in these kingdoms.

The RGS proteins that have been studied biochemically have been found to be **GTPase-activating** proteins (GAPs) for their G-protein targets (3). They bind the activated, GTP-bound form of the Ga subunit and increase the rate at which it hydrolyzes bound GTP. Stimulation of this deactivation rate can be as much as 100-fold. RGS proteins are thus intrinsically inhibitory to G-protein signaling, and overexpression of RGS proteins in animal cells or yeast does indeed inhibit signaling. The

physiological roles of RGS proteins are less well established. In yeast, Sst2p is a feedback inhibitor of the G-protein-mediated **mating** factor signaling pathway. Sst2p [transcription](#) is induced by this pathway, and Sst2p rescues yeast from cell-cycle arrest induced by mating factor. Genetic studies of the flb protein in *Aspergillus* and the Egl10 protein in *C. elegans* also suggest that they act by damping a G-protein signal. RGS proteins may not be physiologically inhibitory in all cases, however; they can also lower basal or background signaling in the absence of stimulation by receptors, and they are probably vital for accelerating the return of signals to background levels upon termination of stimulation by agonists. How RGS proteins are themselves regulated, other than by transcription of their [messenger RNA](#), is unknown, although the GAP activities of several RGS proteins were recently shown to be inhibited when their Ga target is linked to a **palmitoyl group**. RGS proteins appear to be expressed at low levels, in some cases well below those of their presumed Ga targets, and limited association with subpopulations of G proteins may add to the specificities of their effects.

RGS proteins can be recognized from their [primary structures](#) because of a highly conserved domain of about 115 amino acid residues in length. This “RGS box” forms a **four-helix bundle** that binds end-on to the switch II region of the G-protein  $\alpha$  subunit (4). Sequences of RGS proteins diverge markedly on either side of the RGS box, and the functions of these *N*- and *C*-terminal domains are unknown. RGS that have been truncated to leave only the RGS box retain GAP activity, but have diminished affinities for Ga substrates, so the regions outside the box are thought to increase their affinity for G proteins and, perhaps, to enhance selectivity among them. The *N*-terminal regions of several RGS proteins may also contribute to the protein's ability to bind to membranes.

#### Bibliography

1. C. Dietzel and J. Kurjan (1987) *Molec. Cell. Biol.* **7**, 4169–4177.
2. M. R. Koelle and H. R. Horvitz (1996) *Cell* **84**, 115–125.
3. D. M. Berman, T. M. Wilkie, and A. G. Gilman (1996) *Cell* **86**, 445–452.
4. J. J. G. Tesmer, D. M. Berman, A. G. Gilman, and S. R. Sprang (1997) *Cell* **89**, 251–261.

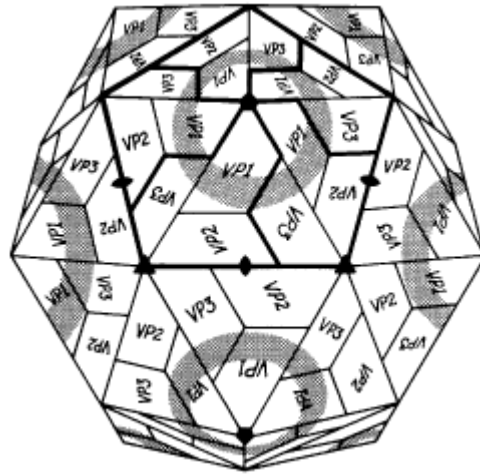
#### Suggestions for Further Reading

5. D. M. Berman and A. G. Gilman (1998) Mammalian RGS proteins: barbarians at the gate. *J. Biol. Chem.* **273**, 1269–1272.
6. M. R. Koelle (1997) A new family of G-protein regulators—the RGS proteins. *Curr. Opin. Cell Biol.* **9**, 143–147.

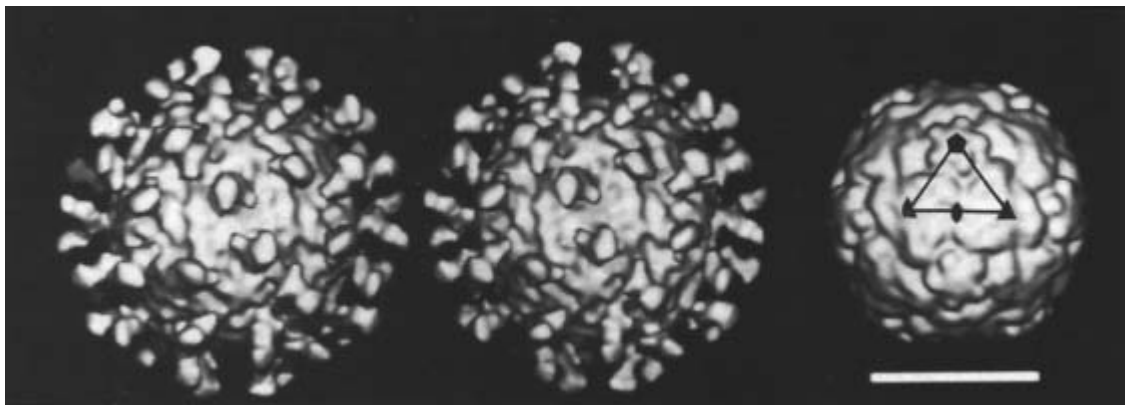
## Rhinovirus

Rhinovirus is a major cause of the common cold, and it is a **virus** member of *Picornaviridae* family (see [Poliovirus](#)). Members of this family are small, icosahedral, nonenveloped viruses. Their structures have been determined by [X-ray crystallography](#) (Fig. 1 and 2). A small shell with a diameter of 30 nm consists of 60 copies each of proteins VP1, VP2, VP3, and VP4, which encapsulate one copy of single-stranded RNA [genome](#) of positive polarity. Like other positive-strand RNA viruses, the genome of rhinovirus is infectious and functions as [messenger RNA](#) after entry into host cell cytoplasm. [Translation](#) of the viral mRNA occurs by entry of a [ribosome](#) in an internal nucleotide sequence within the 5'-untranslated region of the RNA. A unique long open reading frame encodes a single large [polyprotein](#) that is the precursor for all the virus-specific proteins. The optimal temperature for growth of rhinovirus is 33°C in cultured cells.

**Figure 1.** Schematic diagram of a human rhinovirus particle. The icosahedral symmetry, subunit organization, and canyon (shaded) are shown. Thick lines encircle the five protomers of VP1-VP3. The fourth viral protein, VP4, is inside the capsid.(From Ref. [1](#), with permission.).



**Figure 2.** Structures of human rhinovirus. The two pictures on the left give a stereoview of a single-particle reconstruction by cryoelectron microscopy of human rhinovirus 16 in a complex with D1D2, viewed along an icosahedral two-fold axis in approximately the same orientation as in Figure [1](#). D1D2 is a polypeptide chain consisting of domains 1 and 2 of ICAM-1. The right view is a shaded-surface view of rhinovirus 14, computed from the known atomic structure and truncated to 20-Å resolution. The triangular outline of one icosahedral asymmetric unit corresponding to that in Figure [1](#) is indicated. Bar = 200 Å.(From Ref. [1](#), with permission.)



More than 100 different serotypes of human rhinovirus are known. With one exception (rhinovirus type 87), human rhinoviruses use one of two types of receptors to gain entry into the cell. The majority of human rhinovirus, referred to as the major group, bind to the cell-surface molecule known as intercellular adhesion molecule 1 (ICAM-1), a member of the [immunoglobulin](#) superfamily. It has been demonstrated that 60 attachment sites for ICAM-1, called the canyon, exist per virion (Fig. [2](#), see top of next page). Interaction of the receptor with the virion particle results in an alteration of the conformation of the virion structure, which is considered to be important for establishment of the virus infection. The minor group (serotypes 1A, 1B, 2, 29, 30, 31, 44, 47, 49, 62) enter host cells via the low-density lipoprotein receptor (LDLR) and the  $\alpha_2$  [macroglobulin](#) receptor/LDLR-related protein, which is also a member of the LDLR family.

Because of the numerous serotypes of human rhinovirus, it has not been possible to develop effective vaccines against the common cold. Thus, most effort has focused on the development of effective antiviral reagents. Elucidation of the three-dimensional structure of rhinoviruses made it possible to study the interaction of drugs with the virus capsid proteins. As a result, many compounds are described to bind to the **hydrophobic** pocket at the floor of canyon on the virion particle, so that the receptor cannot bind the virus. Although topical administration of canyon inhibitors failed to provide a favorable effect, these inhibitors have been used effectively to understand the molecular basis of virus-receptor interaction.

## Bibliography

1. N. H. Olson et al. (1993) Proc. Natl. Acad. Sci. USA **90**, 507–511.

## Suggestions for Further Reading

2. R. R. Rueckert (1996) "*Picornaviridae: The Viruses and Their Replication*". in *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 609–654.
3. G. Stanway (1994) "Rhinoviruses". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1253–1259.

## Rho GTPases

Ras homology (RhO) proteins constitute one of the families of cytosolic 20- to 30-kDa **GTPases** **homologous** to p21<sup>Ras</sup>. They are involved in regulating the actin-**cytoskeleton** organization, gene [transcription](#), membrane trafficking, [development](#), and cell proliferation. The three key members were originally described on the basis of their effects on cultured 3T3 fibroblasts:

- RhoA was implicated in forming stress fibers, bundles of [actin](#) filaments that contain myosin and  $\alpha$ -actinin.
- Rac1 induces a meshwork of actin filaments associated with lamellipodia and membrane ruffles.
- Cdc42Hs produces actin-rich protrusions called filopodia.

More esoteric members of this continually expanding family include the human proteins RhoB, RhoC, RhoD, RhoE, RhoG, Rac2, Rac3, TC10, Rnd1, Rnd2, and several homologues found in yeast and invertebrates, such as *Drosophila melanogaster* and *Caenorhabditis elegans*. Because of space limitations, this article focuses on the principal three mammalian Rho GTPases.

### 1. Overview

Like all [GTP-binding proteins](#), Rho [Gtpases](#) cycle between the biologically active, GTP-bound form, and the inactive GDP-bound form. The GDP-bound form is found predominantly in the cytosol, tightly complexed with Rho guanine nucleotide dissociation inhibitor (GDI). The active state of the GTPase has a lifetime limited by the rate at which it hydrolyzes bound GTP to GDP and inorganic phosphate. Because of the very low intrinsic catalytic rates of the GTPases, their biological activity is down-regulated by GTPase activating proteins, or GAP s. Upon GAP-assisted GTP hydrolysis, the GTPase reverts to the GDP-bound ground state and is solubilized in the cytosol by RhoGDI, completing the cycle.

The biological cycle of a Rho GTPase is initiated in response to extracellular signals with the exchange of GDP for GTP, a process that requires dissociation of the RhoGDI–GTPase complex. A

family of [guanine nucleotide exchange factors](#) (GEFs) is responsible for stimulating the exchange of GDP for GTP. In the active, GTP-bound state, the Rho GTPases are associated with the cell membrane. This interaction, it is thought, depends primarily on the geranylgeranyl [membrane anchor](#) groups on the C-terminal [cysteine](#) residues. In their active form, the Rho GTPases bind to and stimulate specific downstream effectors, often protein [kinases](#), thus activating a network of signaling pathways that span the entire cellular space from the focal adhesion complexes to the [nucleus](#).

Signaling pathways that involve the principal Rho GTPases are intertwined. In 3T3 fibroblasts, Cdc42Hs activates Rac, which in turn activates Rho. All Rho proteins function downstream of a Ras-mediated cascade. It is quite possible, however, that the hierarchical patterns differ, depending on the cell type and the physiological phenomena involved.

## 2. Molecular structure

At present, our knowledge of the molecular structures of Rho proteins is derived from **X-ray crystallographic** studies of Rac1•GppNHp (1), RhoA•GDP (2) and RhoAVal14•GTPγS (3), all in their unprenylated forms. In addition, Cdc42Hs has been studied in both GDP and GTP-analog bound states by multidimensional [NMR](#) (4). Like Ras, Rho [protein structures](#) consist of a single **domain** with a six-stranded [beta-sheet](#) surrounded by **alpha-helices** (Fig. 1). The following description applies specifically to RhoA, which is used as a representative example. The b-sheet contains two antiparallel (B2 and B3) and five parallel [beta-strands](#) (B3, B1, B4-B6), five a-helices (A1, A3, A3', A4, A5), and three  $3_{10}$ -helices (H1–H3). When compared to Ras, Rho proteins contain three insertions and one deletion. The 13-residue insertion (Asp 124 to Gln 136 in RhoA) between strand B5 and helix A4 is the most distinctive feature that clearly differentiates Rho from other Ras-related GTPases. The amino acid sequence in this insert varies considerably among different Rho proteins, suggesting functional implications. Interestingly, the conformation of this helix does not depend on the type of nucleotide bound to the protein.

Conformational differences between RhoA•GDP and RhoA•GTPγS are restricted to the switch I and switch II regions (the nomenclature follows that used for Ras). Dramatic changes in switch I from the GDP- to the GTP-bound state are exemplified by displacements of 5.4 Å and 6.4 Å (1 Å = 10<sup>-10</sup>m) at Pro36 and Phe39, respectively (Fig. 1b). The side chains of Pro36 and Tyr34 become oriented toward the nucleotide, so that the ring of Tyr34 stacks on Pro36. Three **hydrophobic** residues, Val35, Val38, and Phe39, become solvent exposed, suggesting their involvement in target binding following GTP-induced activation. An important difference between the GDP and GTPγS-bound RhoA relates to the mode of Mg<sup>2+</sup> ion binding. In the former, the Mg<sup>2+</sup> ion is coordinated by three [water](#) molecules, a g-phosphate oxygen atom, and two protein ligands, including the main-chain carbonyl of Thr37, whereas in the GTPγS structure the coordination is similar to p21<sup>Ras</sup> in its active form and involves the side chain of Thr37.

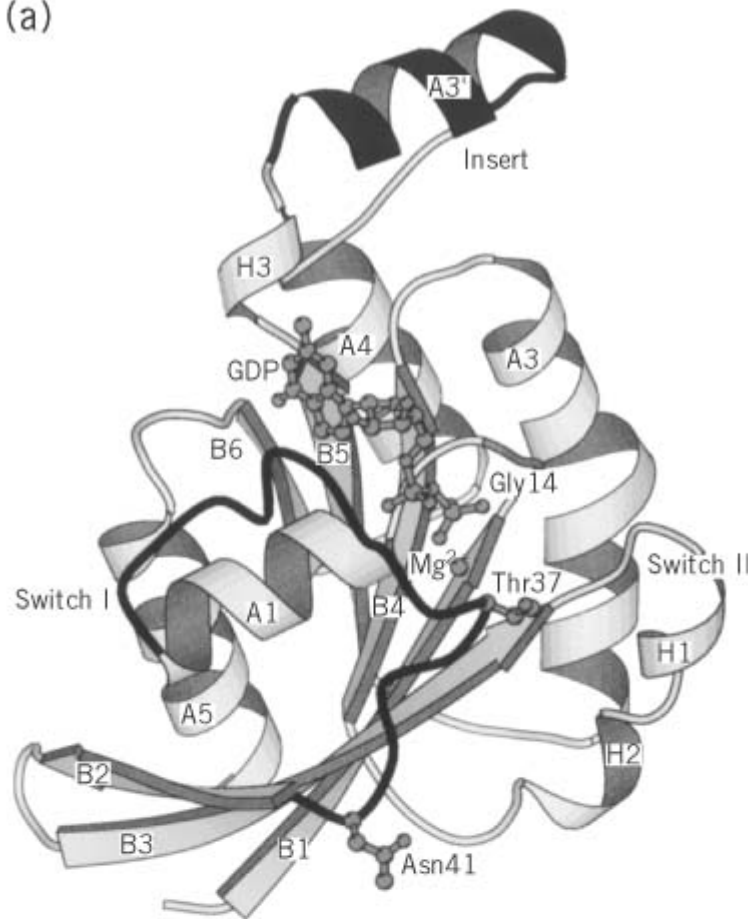
Switch II contains a segment of the polypeptide chain between strands B3 and B4, including Gln63, which is essential for the catalytic activity of the protein, both intrinsic and GAP-stimulated. In RhoA•GDP, residues 63 to 65 are disordered, whereas in the complex with GTPγS the segment is well defined by two  $3_{10}$  helices, H1(residues 64 to 66) and H2(residues 70 to 72), separated by a short loop. A similar conformation is found in Rac•GppNHp, suggesting that it is physiologically relevant. The sequence in this region is relatively well conserved among Rho GTPases and is probably less involved in target selection.

Both GTP and GDP bind to Rho GTPases with nanomolar **dissociation constants**. The nucleotides interact with the phosphate-binding [P loop](#), typical of all small GTPases, and contain a **consensus sequence** Gly-X-Gly-X-X-Gly-Lys-Thr/Ser. Mutation of the second Gly in this motif (Gly14) to Val renders Rho active by inhibiting its catalytic properties, both intrinsic and GAP-induced. As inferred from the RhoA•GTPγS complex, GTP makes 21 direct and nine water-mediated [hydrogen bonds](#) to

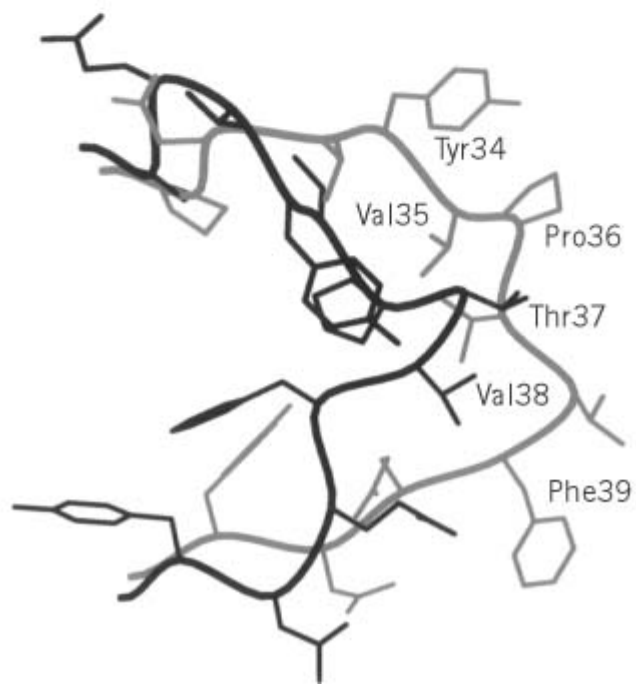
the protein that involve six residues of the P loop, four residues in switch I, three residues in switch II, and four residues involved in base recognition. This coordination of the nucleotide rationalizes the structural changes in the switch I and switch II regions. Figure 1 shows the representative structural features of the RhoA GTPase.

**Figure 1.** Structural features of RhoA GTPase. **(a)** The structure of RhoA in an inactive conformation with bound GDP and  $Mg^{2+}$ . The secondary structural elements are labeled, and the switch I and switch II regions are highlighted. Also indicated are the locations of the insert helix and residues Gly14, Thr37, and Asn41. **(b)** The conformational changes in switch I upon activation of RhoA. Alignment of the inactive and active structures of RhoA reveals dramatic changes in the switch I region. The GDP-bound conformation is shown as dark line from the same perspective as in A, and the aligned GTP<sub>S</sub>-bound conformation is shown as the lighter line. PDB codes 1FTN and 1A2B.

(a)



(b)



All Rho GTPases undergo complex **posttranslational modifications**. First, Type I geranylgeranyl (20-carbon) transferase attaches the geranylgeranyl group covalently through a thioether linkage to the cysteine residue in the C-terminal Cys-Ala-Ala-X motif, that typically contains Leu in the “X” position. Following prenylation, the three terminal residues are cleaved off by a specific [proteinase](#) and the terminal prenylated Cys is carboxymethylated (5). Unmodified Rho GTPases (eg, those expressed in *Escherichia coli*) fail to associate with membranes and do not interact with GDI (see GDI later), but in most cases bind and activate their downstream effectors in complexes with slowly hydrolyzable GTP analogs, such as GTPγS or GppNHp. Although prenylation is a necessary factor for translocating Rho to membranes, it may not be sufficient. Emerging evidence suggests that specific proteins in caveolae—unique microdomains of membranes involved in [signal transduction](#)—may be involved in Rho binding (6). The structural aspects of these interactions are not known.

### 3. Regulatory proteins

Like all GTPases, Rho proteins are subject to tight regulation. This is accomplished in a cell by accessory proteins that act downstream and upstream of the Rho•GTP complex. A brief review follows of what is known about these regulatory proteins.

#### 3.1. Guanine Nucleotide Exchange Factors (GEFs)

GEFs activate Rho GTPases by a >10 fold stimulation of GDP dissociation and concomitant exchange of GDP for GTP. Almost twenty proteins that have Rho GEF function have been identified to date in mammalian [genomes](#). Tight binding is observed for GEFs and nucleotide-depleted GTPases. GEFs cannot act on the GDP-bound cytosolic pool of Rho proteins, which are in complex with GDI, and it is believed that additional accessory molecules are involved. The ezrin, radixin, moesin system (ERM) has been recently implicated in the dissociation of the GTPase-RhoGDI complex (7, 8).

The mechanism of  $Mg^{2+}$  coordination in RhoA•GDP and mutagenetic studies are consistent with a model in which GEFs act on Thr37 (RhoA numbering) of Rho to destabilize the bound  $Mg^{2+}$  ion before GDP dissociation (9). The catalytic activity of GEFs resides in the Dbl homology (DH) domain, typically linked to a PH ( **pleckstrin homology**) domain, responsible for targeting GEF (10). However, several proteins that have DH domains do not show any GEF activity and therefore may recruit Rho to specific locations. Finally, other domains, such as src homology 3 (SH3), and the diacylglycerol-binding “zinc butterfly” motif are also found in GEFs, suggesting additional functions (11).

Some GEFs are quite selective with respect to the GTPase that they activate, and others have broad specificities. For example, Lbc and Lsc are distinctly Rho-specific with respect to both nucleotide exchange and binding, whereas Dbl binds to all three Rho proteins and stimulates GDP for GTP exchange in Cdc42Hs and Rho (12).

Many Rho GEFs, including the prototypical Dbl, are **oncoproteins**, although their oncogenic potential may stem primarily from functions other than their GEF activity.

Relatively little is known about the upstream activation events that stimulate the action of GEFs on Rho. **Tyrosine kinase** and [G-protein coupled receptors](#), it is believed, are involved in the activation step (13, 14). The recent discovery of p115RhoGEF, a GTPase activating protein for  $G_{\alpha_{12}}$  and  $G_{\alpha_{13}}$  provides evidence of a direct link between heterotrimeric G-proteins and Rho GTPases (15, 16).

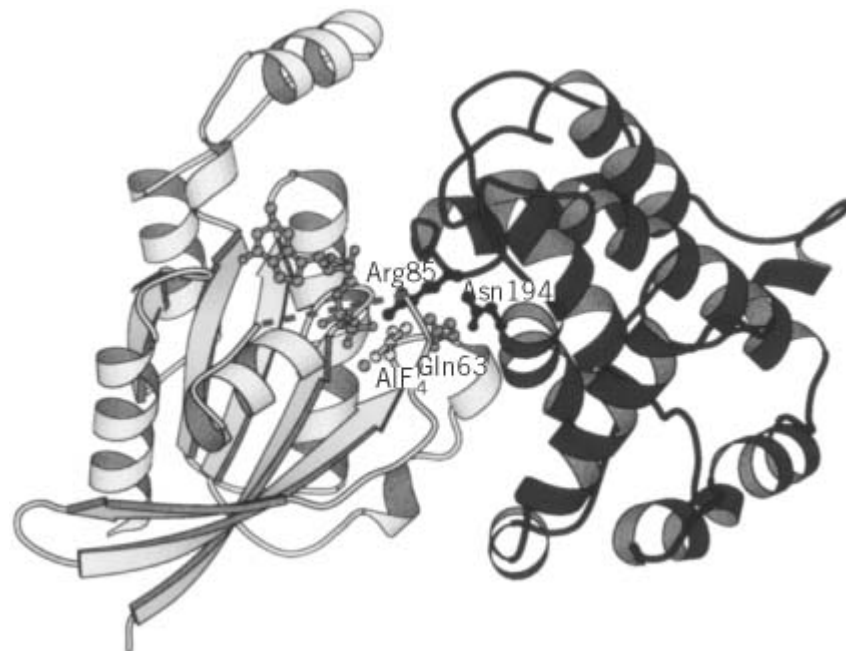
#### 3.2. GTPase-Activating Proteins (GAPs)

RhoGAPs enhance the rate of hydrolysis of GTP by Rho GTPases by a factor of up to  $10^5$ . They are multidomain, often large proteins, but their GTPase activity resides in one particular domain, known as the BH-domain or just the GAP domain (17). The breakpoint cluster region homology (BH)



domains are approximately 200 residues long and constitute a **homologous** family. The BH-domain of the **phosphoinositide 3-kinase p85a** subunit is unique in that it shows no enhancement of hydrolysis although it binds Cdc42Hs and Rac (18). BH domains enhance the catalytic activity of Rho GTPases. GAPs bind strongly to Rho GTPases only in the presence of  $\text{AlF}_4$ , when a [transition-state analog](#) is produced by the complex of  $\text{AlF}_4$  and RhoGDP [this phenomenon was first observed in RasGAP, which is required to downregulate Ras, also a GTPase (19)]. The structure of the BH domain of p50RhoGAP, a ubiquitous GAP that activates all Rho proteins, has been determined by X-ray crystallography on its own and in complexes with Cdc42Hs•GppNHP and RhoA•GDP, the latter in the presence of  $\text{AlF}_4$  (20-22). The tertiary structure of this molecule includes nine helices, four of which form the core bundle (Fig. 2). Though originally not recognized, there is a structural relationship between the RhoGAPs and the RasGAPs, implying possible homology (23). Like RasGAP, the Rho GAPs operate by the “arginine-finger” mechanism. When bound to the transition-state Rho•GTP, GAP donates an Arg residue (Arg85 in the case of p50RhoGAP) directly to the active site. This arginine is functionally equivalent to Arg178 of the [heterotrimeric G protein](#)  $G_{i\alpha 1}$ . In addition to their direct participation in catalysis, GAPs also stabilize the switch regions in the GTPases.

**Figure 2.** The structure of p50RhoGAP complexed with RhoA GDP in the presence of  $\text{AlF}_4$ . The transition-state complex is mimicked by the binding of p50RhoGAP (dark area on right side) to RhoA with bound GDP (area near  $\text{AlF}_4$ ) and  $\text{AlF}_4$  (light areas on left side). PDB code 1TX4. Arg85 of p50RhoGAP is highlighted, as is Gln63 from RhoA.



The crystal structures also allow rationalizing the observed effects of the Gly14→Val mutation of Rho. Gly14 is within the interface region, and any side chain in this position interferes with GAP binding. Accordingly, the mutant resists the GAP-induced GTP-hydrolysis and continues to activate its targets.

The structures of the p50RhoGAP complexes show clear interactions with both switch I and switch II residues of Rho. Specifically, Gln63 in RhoA is clearly resolved in the complex that contains  $\text{AlF}_4$ , suggesting that GAP assists in the stabilizing the catalytically competent conformation that

orients the hydrolytic water molecule during its nucleophilic attack on the  $\gamma$ -phosphate. This stabilization can be seen as an **allosteric** contribution to catalysis.

RhoGAPs downregulate the Rho GTPases by accelerating their GTP hydrolysis and concomitant return of the protein to its inert, GDP-bound state, and also by binding to those **epitopes** on the GTPase that are thought to be involved in the interactions with effectors. It has been postulated that some GAPs may serve as specific Rho effectors, and studies of *n*-chimaerin substantiate this proposal (24). The GAPs are targeted to the appropriate location through their SH2 and SH3 domains, their PH domains, and proline-rich segments that may interact with SH3 domains on other proteins. Some GAPs contain domains harboring catalytic activities, for example, kinases. Others contain GTPase binding domains, GEF domains, or novel domains of unknown function. Rho-specific BH-domains have also been identified in the heavy chains of class IX myosins (25).

The specificity of GAPs toward different Rho family GTPases varies significantly. Most GAPs are active on at least two GTPases. For example, MgcRacGAP acts on Rac and Cdc42Hs (26), whereas Gaf activates Cdc42Hs and RhoA, but not Rac (27). Specific activities are typically assayed *in vitro* on either isolated BH-domains or GST-fused BH-domains. It is not clear to what extent the activity of GAPs is modulated *in vivo* by targeting through other domains.

### 3.3. GDP Dissociation Inhibitor (RhoGDI)

RhoGDI extracts Rho GTPases from membranes and arrests the cycling of these proteins at the GDP-bound point (8, 28-30). Both prenylation and **proteolytic** truncation of the C-terminal three residues from the Cys-Ala-Ala-X motif in Rho are essential for RhoGDI binding, but carboxymethylation is not. Although GDIs have been described for other GTPases, notably Rab, the RhoGDI is a unique protein that has three closely homologous isoforms identified to date: RhoGDI $\alpha$ , GDI/D4 (67% sequence identity)(31) and GDI $\gamma$  (50% sequence identity) (32). Although RhoGDI $\alpha$  interacts with all three key Rho proteins, the specificities of the other isoforms have not been determined. The ubiquitously expressed RhoGDI $\alpha$  is made up of a single polypeptide chain of 204 residues. There are two domains, the N-terminal domain (residues 1 to 64) which is disordered in solution in the absence of a GTPase, and an ordered domain, which begins approximately at residue 65 and has a  $\beta$ -sandwich fold (33, 34). The N-terminal 22 residues do not play any functional role, whereas the fragment 23 to 64 is essential for the guanine nucleotide exchange inhibition. Binding to the GTPase, with a dissociation constant in the high nanomolar range, is attributed largely to the  $\beta$ -sandwich domain (a truncated 42 to 204 protein binds to Cdc42Hs only 100-fold less tightly). A narrow hydrophobic cleft in the C-terminal domain is expected to bind isoprenes. The structure of the RhoGTPase-RhoGDI complex is currently not known, but mutagenetic results indicate the insertion helix characteristic of Rho proteins (see above) is involved in the GDP-dissociation inhibition by GDI (35).

## 4. Downstream effectors

To date, a number of potential targets of Rho proteins have been identified, primarily by the yeast **two-hybrid system** (36). In many cases, especially those that involve enzymes, the functional significance of these interactions has been substantiated by activation assays. The following are among the targets thought to be physiologically most relevant.

### 4.1. ROK $\alpha$ /ROK $\beta$

These proteins, more than 1300-residues long, coprecipitate with RhoA in its active form and associate with the latter physically and functionally *in vivo* and *in vitro* (37, 38). The two proteins differ only with respect to their N-termini. ROK $\alpha$  is nine residues longer. The catalytic N-terminal domain was identified on the basis of amino acid sequence homology to the Ser/Thr **protein kinase A** family. The catalytic activity of ROK kinases is enhanced by binding of active RhoA. Downstream of the kinase domain is a putative **coiled-coil** segment of 600 residues, followed by a PH-domain and a cysteine-rich **zinc-finger** motif. The kinase activity of ROK induces the formation of stress fibers in the absence of the active RhoA, strongly suggesting that the reorganization of the cytoskeleton by

RhoA is mediated through ROK. Among the putative targets of ROK is the myosin-binding subunit (MBS) of myosin light-chain phosphatase (39, 40).

#### 4.2. Protein Kinase N (PKN)

PKN is a 120-kDa **serine/threonine kinase** that interacts with RhoA in a GTP-dependent fashion (41, 42). The kinase domain in this protein is homologous to **protein kinase C**. PKN is probably involved in the reorganization of cytoskeletal elements, a function mediated through its direct interaction with  $\alpha$ -actinin (43).

#### 4.3. p21 Associated Kinases (PAK)s

These Ser/Thr kinases associate with both Rac and Cdc42Hs in a GTP-dependent manner. At least three **isoforms** (PAK1/a, PAK2/g, and PAK3/b) are found in mammalian tissues, albeit they have different tissue distributions (44-46). PAK3 is released from the GTPase upon activation, and the other two remain bound. The protein kinase domains located in the C-terminal part are highly conserved and contain a distinguishing Gly-Thr-Pro-Tyr/Phe-Trp-Met-Ala-Pro-Glu (GTPY/FWMAPE) signature motif, a possible site of regulatory **phosphorylation**. The N-terminal fragment contains a Cdc42Hs/Rac interactive binding (CRIB) domain, also known as the **p21-binding domain** (PBD) or **GTPase-binding domain** (GBD) domain, which mediates the interaction of PAK with the GTPase. CRIB-domains have been found in other proteins, and they may effectively help in identifying Rho targets (47, 48). An N-terminal Pro-rich fragment, it is believed, interacts with SH3-containing proteins. Although PH-domains occur in yeast PAK homologues, no such PAKs have been found to date in mammalian genomes. At present, downstream signaling pathways mediated by PAKs are not well understood. Involvement in cytoskeletal rearrangement is very likely.

#### 4.4. Plenty of SH3s (POSH)

This ubiquitously expressed 93-kDa protein interacts with Rac (but not Rho and Cdc42Hs) in a GTP-dependent manner (49). It is implicated specifically in gene transcription regulation, that is, Rac-mediated JNK activation. In Cos-1 cells, POSH induces JNK activation and NF- $\kappa$ B nuclear translocation. The protein contains a potential zinc-finger structure at the N-terminus and four SH3 domains. The two C-terminal SH3 domains are sufficient to induce NF- $\kappa$ B translocation, but not JNK activation. Overexpression of POSH is toxic to cells and results in **apoptosis**.

#### 4.5. p67<sup>phox</sup>

This protein is a Rac target and a component of the NADPH oxidase complex. It is essential for superoxide formation, along with Rac1, cytochrome b, and p47<sup>phox</sup> (50). The Rac1 binding site is cryptic and has been mapped to the N-terminal fragment, which can inhibit other Rac-mediated pathways (51).

#### 4.6. Other Targets

Other targets include raphilin (42) and rhotekin (42), both of which are homologues of PKN devoid of kinase activity and targets of RhoA; citron (52) (target of Rho and Rac); phospholipase D (RhoA, possibly Rac) (53); p140mDia (RhoA) (54); p140Sra-1 (Rac1) (55); phosphoinositide 3-kinase (Rac, Cdc42Hs) (56); IQGAP (Cdc42Hs) (57); and several others.

The Rho-family targets contain a variety of GTPase-binding motifs. Similarly, Rho GTPases utilize several surface epitopes in their interactions with the effectors. Lack of structural information about the complexes precludes a more detailed analysis.

### 5. The mechanism of *Clostridium* toxins

Several Rho proteins are targets for potent **toxins** produced by various species of *Clostridium*. Toxins A and B are from *C. difficile* monoglucosylate Thr37 (58, 59). The same residue is modified by the *C. novyi*  $\alpha$ -toxin, which transfers an *N*-acetyl-glucosamine moiety onto it (60). These toxins target all three major Rho proteins. The C3 toxin from *C. botulinum* and the staphylococcal

epidermal-cell differentiation inhibitor catalyze specific **ADP-ribosylation** of RhoA on Asn41 ([61-63](#)). Because Asn41 is close to the critical region of switch I, ADP-ribosylation may cause steric hindrance to prevent forming productive Rho-effector complexes.

### 5.1. Physiological and Pathophysiological Phenomena Regulated by Rho

Extensive studies have implicated Rho GTPases in a variety of regulatory mechanisms, specifically those that control cytoskeletal organization and gene transcription. In addition, numerous targets have been identified, typically kinases that become activated when bound to the GTPase. Given the wide ranging effects of Rho on cells, it is often difficult to identify specific physiological effects that the proteins exert on the tissues or systems. Three examples of established physiological functions of Rho are given here.

### 5.2. Calcium Sensitization of Smooth-Muscle

The regulatory effect of  $\text{Ca}^{2+}$  on smooth-muscle contraction is amplified by the  $\text{Ca}^{2+}$  sensitization effect, which involves inhibition of myosin phosphatase ([64](#)). RhoA and its target ROKb mediate the process. Abnormal smooth-muscle contractility is a factor in hypertension. The pyridine derivative Y-27632 inhibits ROKb kinase activity, and it consequently inhibits smooth-muscle contraction. Administered to hypertensive rats, the compound drastically reduces their blood pressure, suggesting a new approach to hypertension chemotherapy ([65](#)).

### 5.3. Development

The ability of Rho to reorganize actin filaments is vital for morphogenetic pathways during development. This has been demonstrated primarily in *Drosophila*, where Rac and Cdc42 homologues are expressed ubiquitously in the nervous system and mesoderm and play a key role in neural and muscle differentiation ([66](#)). Strong evidence supports the role of Rac and RhoA in neuronal development in humans. Interestingly, X-linked mental retardation has been linked to disruption of RhoA-mediated developmental pathways ([67](#), [68](#)).

### 5.4. Cell Transformation

Although Rho proteins are not generally regarded as potent oncoproteins, they are essential for p21<sup>Ras</sup>-induced **neoplastic transformation** ([69](#)). It has been shown that RhoA suppresses p21<sup>Waf1/Cip1</sup>, an inhibitor of the cell cycle that, in RhoA's absence, would interfere with Ras driving the cells into the S-phase of the **cell cycle** ([70](#)).

## 6. Summary

The studies of Rho-mediated signaling pathways constitute one of the most active contemporary areas of biomedical research. This article provides only a superficial and very selective overview. Interested readers are urged to refer to current literature for the status of the field.

## Bibliography

1. M. Hirshberg, R. Stockley, G. Dodson, and R. Webb. (1997) *Nat. Struct. Biol.* **4**, 147–151.
2. Y. Wei et al. (1997) *Nat. Struct. Biol.* **4**, 699–702.
3. K. Ihara et al. (1998) *J. Biol. Chem.* **273**, 9656–9666.
4. J. L. Feltham et al. (1997) *Biochemistry* **36**, 8755–8766.
5. P. Adamson, C. J. Marshall, A. Hall, and P. A. Tilbrook (1992) *J. Biol. Chem.* **267**, 20033–20038.
6. D. Gingras et al. (1998) *Biochem. Biophys. Res. Commun.* **247**, 888–893.
7. K. Takahashi et al. (1997) *J. Biol. Chem.* **272**, 23371–23375.
8. T. Sasaki and Y. Takai (1998) *Biochem. Biophys. Res. Commun.* **245**, 641–645.
9. R. Li and Y. Zheng (1997) *J. Biol. Chem.* **272**, 4671–4679.
10. M. F. Olson et al. (1997) *Oncogene* **15**, 2827–2831.

11. R. A. Cerione and Y. Zheng (1996) *Curr. Opin. Cell Biol.* **8**, 216–222.
12. J. A. Glaven et al. (1996) *J. Biol. Chem.* **271**, 27374–27381.
13. S. Farah et al. (1998) *J. Biol. Chem.* **273**, 4740–4746.
14. C. D. Nobes, P. Hawkins, L. Stephens, and A. Hall (1995) *J. Cell Sci.* **108**, 225–233.
15. M. J. Hart et al. (1998) *Science* **280**, 2112–2114.
16. T. Kozasa et al. (1998) *Science* **280**, 2109–2111.
17. N. Lamarche and A. Hall (1994) *Trends Genet.* **10**, 436–440.
18. A. Musacchio, L. C. Cantley, and S. C. Harrison (1996) *Proc. Natl. Acad. Sci. USA* **93**, 14373–14378.
19. R. Mittal, M. R. Ahmadian, R. S. Goody, and A. Wittinghofer (1996) *Science* **273**, 115–117.
20. T. Barrett et al. (1997) *Nature* **385**, 458–461.
21. K. Rittinger et al. (1997) *Nature* **388**, 693–697.
22. K. Rittinger et al. (1997) *Nature* **389**, 758–762.
23. K. Rittinger, W. R. Taylor, S. J. Smerdon, and S. J. Gamblin (1998) *Nature* **392**, 448–449.
24. R. Kozma, S. Ahmed, A. Best, and L. Lim (1996) *Mol. Cell. Biol.* **16**, 5069–5080.
25. P. Post, G. Bokoch, and M. Mooseker (1998) *J. Cell Sci.* **111**, 941–950.
26. A. Toure et al. (1998) *J. Biol. Chem.* **273**, 6019–6023.
27. J. D. Hildebrand, J. M. Taylor, and J. T. Parsons (1996) *Mol. Cell. Biol.* **16**, 3169–3178.
28. T. Ueda et al. (1990) *J. Biol. Chem.* **265**, 9373–9380.
29. Y. Fukumoto et al. (1990) *Oncogene* **5**, 1321–1328.
30. G. M. Bokoch, B. P. Bohl, and T. H. Chuang (1994) *J. Biol. Chem.* **269**, 31674–31679.
31. J. M. Lelias et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 1479–1483.
32. C. N. Adra et al. (1997) *Proc. Natl. Acad. Sci. USA* **94**, 4279–4284.
33. Y. Q. Gosser et al. (1997) *Nature* **387**, 814–819.
34. N. H. Keep et al. (1997) *Structure* **5**, 623–633.
35. W. J. Wu, D. A. Leonard, R. A. Cerione, and D. Manor (1997) *J. Biol. Chem.* **272**, 26153–26158.
36. P. Aspenstrom and M. F. Olson (1995) *Methods Enzymol.* **256**, 228–241.
37. T. Leung, E. Manser, L. Tan, and L. Lim (1995) *J. Biol. Chem.* **270**, 29051–29054.
38. T. Ishizaki et al. (1996) *EMBO J.* **15**, 1885–1893.
39. K. Kimura et al. (1996) *Science* **273**, 245–248.
40. T. Matsui et al. (1996) *EMBO J.* **15**, 2208–2216.
41. M. Amano et al. (1996) *Science* **271**, 648–650.
42. G. Watanabe et al. (1996) *Science* **271**, 645–648.
43. H. Mukai et al. (1997) *J. Biol. Chem.* **272**, 4740–4746.
44. G. A. Martin, G. Bollag, F. McCormick, and A. Abo (1995) *EMBO J.* **14**, 1970–1978.
45. E. Manser et al. (1995) *J. Biol. Chem.* **270**, 25070–25078.
46. S. Bagrodia et al. (1995) *J. Biol. Chem.* **270**, 22731–22737.
47. P. D. Burbelo, D. Drechsel, and A. Hall (1995) *J. Biol. Chem.* **270**, 29071–29074.
48. M. G. Rudolph et al. (1998) *J. Biol. Chem.* **273**, 18067–18076.
49. N. Tapon, K. Nagata, N. Lamarche, and A. Hall (1998) *EMBO J.* **17**, 1395–1404.
50. D. Diekmann et al. (1994) *Science* **265**, 531–533.
51. S. Ahmed et al. (1998) *J. Biol. Chem.* **273**, 15693–15701.
52. P. Madaule et al. (1995) *FEBS Lett.* **377**, 243–248.
53. C. D. Bae, D. S. Min, I. N. Fleming, and J. H. Exton (1998) *J. Biol. Chem.* **273**, 11596–11604.

54. N. Watanabe et al. (1997) EMBO J. **16**, 3044–3056.
55. K. Kobayashi et al. (1998) J. Biol. Chem. **273**, 291–295.
56. K. F. Tolias, L. C. Cantley, and C. L. Carpenter (1995) J. Biol. Chem. **270**, 17656–17659.
57. J. W. Erickson, R. A. Cerione, and M. J. Hart (1997) J. Biol. Chem. **272**, 24443–24447.
58. I. Just et al. (1995) Nature **375**, 500–503.
59. I. Just et al. (1995) J. Biol. Chem. **270**, 13932–13936.
60. J. Selzer et al. (1996) J. Biol. Chem. **271**, 25173–25177.
61. E. J. Rubin, D. M. Gill, P. Boquet, and M. R. Popoff (1988) Mol. Cell. Biol. **8**, 418–426.
62. M. Yamamoto et al. (1993) Oncogene **8**, 1449–1455.
63. M. Sugai, C. H. Chen, and H. C. Wu (1992) Proc. Natl. Acad. Sci. USA **89**, 8903–8907.
64. M. C. Gong et al. (1996) Proc. Natl. Acad. Sci. USA **93**, 1340–1345.
65. M. Uehata et al. (1997) Nature **389**, 990–994.
66. L. Luo, Y. J. Liao, L. Y. Jan, and Y. N. Jan (1994) Genes Dev. **8**, 1787–1802.
67. P. Billuart et al. (1998) Nature **392**, 923–926.
68. L. Luo et al. (1996) Nature **379**, 837–840.
69. R. G. Qiu, A. Abo, F. McCormick, and M. Symons (1997) Mol. Cell. Biol. **17**, 3449–3458.
70. M. F. Olson, H. F. Paterson, and C. J. Marshall (1998) Nature **394**, 295–299.

### Suggestions for Further Reading

71. V. Aelst and C. D'Souza-Schorey (1997) Rho GTPases and signaling networks, Genes Dev. **11**, 2295–2322. A superb overview of current knowledge with an exceptionally exhaustive bibliography.
72. A. Hall (1998) Rho GTPases and the actin cytoskeleton, Science **279**, 509–514. An outstanding review by one of the pioneers and experts in the field.
73. M. Geyer and A. Wittinghofer (1997) GEFs, GAPs, GDIs and effectors: Taking a closer (3D) look at the regulation of Ras-related GTP-binding proteins, Curr. Opin. Struct. Biol. **7**, 786–792. A brief review leading to a good bibliography.
74. S. J. Gamblin and S. J. Smerdon (1998) GTPase-activating proteins and their complexes. Curr. Opin. Struct. Biol. **8**, 195–201. A good overview with a comprehensive bibliography.
75. K. Aktories (1997) Rho proteins: Targets for bacterial toxins, Trends Microbiol. **5**, 282–288.

## Rhodopsin

Rhodopsin is the key player in the biology of vision; it is a [membrane protein](#) with seven [a-helices](#) that functions as the photoreceptor in rod cells of the retina in the vertebrate and invertebrate eye ([1](#)). Many eukaryotic G-protein-coupled receptors (see [Heterotrimeric G Proteins](#)), among them the  $\beta$ -adrenergic receptor, are **homologous** to rhodopsin. In these receptors, the ligands bind to the pocket occupied by the **prosthetic group** retinal in rhodopsin. Bovine rhodopsin consists of 348 residues and resides in the disc membranes of the outer segments of rod cells. The photoactive center of rhodopsin comprises retinal, or vitamin A aldehyde, bound through a protonated [Schiff Base](#) to a [lysine](#) residue in the C-terminal [a-helix](#) of the protein. Upon absorption of a photon, the bound chromophore isomerizes from the 11-*cis* form to all-*trans* retinal. The [cis/trans isomerization](#) is followed by the sequential appearance of spectroscopic intermediates (known as the rhodopsin

photocycle). Finally, the isomerized chromophore dissociates from the protein, followed by regeneration of the dark form of rhodopsin with 11-*cis* retinal.

An intermediate called metarhodopsin II, in which the retinal is in the all-*trans* configuration and the Schiff base unprotonated, is formed within milliseconds during the photocycle. This is linked to structural changes on the cytoplasmic face of rhodopsin, where it interacts with [transducin](#), a heterotrimeric G protein. This structural change involves a rigid-body movement of  $\alpha$ -helix F (2) and leads to the dissociation of the  $\alpha$ -subunit from transducin. The  $\alpha$ -subunit of transducin activates a **phosphodiesterase**, which hydrolyses [cyclic GMP](#). The decrease in the cyclic GMP concentration closes cation channels in the plasma membrane and causes hyperpolarization of the membrane of the rod cells. Cycles of hyperpolarization in light and depolarization in dark generate the nerve impulses required for vision. The lifetime of metarhodopsin II is of the order of 1 min. Because of the longevity of this state, the activity of metarhodopsin II is also regulated by rhodopsin kinase (see **Phosphorylation**) and by the binding of an inhibitory protein, called *arrestin*, to the phosphorylated site near the C-terminus of rhodopsin on the cytosolic side of the membrane.

A structural model of rhodopsin at an intermediate resolution has been constructed using [electron microscopy](#) of two-dimensional crystals (3, 4). The model shows that the arrangement of the seven transmembrane helices in rhodopsin is distinct from that of another light-transducing, retinal-containing protein, [bacteriorhodopsin](#). Bacteriorhodopsin from the archeon *Halobacterium salinarium* (formerly *H. halobium*), with its high-resolution structure and wealth of mutagenetic, spectroscopic, and biophysical data, is the most thoroughly characterized membrane protein. It contains 248 residues, has seven transmembrane  $\alpha$ -helices, and carries the same prosthetic group as rhodopsin, namely, a retinal bound to a lysine residue via a protonated Schiff base. In both rhodopsin and bacteriorhodopsin, the absorption of a light quantum brings about a photocycle with associated rotational and tilting movements of some transmembrane  $\alpha$ -helices. Despite these similarities, there is no sequence homology between bacteriorhodopsin and rhodopsin.

Bacteriorhodopsin is a light-driven **proton pump**. As with rhodopsin, absorption of a photon causes isomerization of the bound retinal in bacteriorhodopsin, but in this case from all-*trans* to 13-*cis*. This primary photochemical event is followed by a characteristic multistep photocycle, in which some steps are linked to slight movements of some  $\alpha$ -helices and proton translocation across the cell membrane (5). There are two aqueous channels, one from each side of the membrane to the Schiff base in the middle. Two [aspartic acid](#) residues, one on the cytoplasmic side and the other on the extracellular side of the Schiff base, are located in the narrow channels and participate directly in proton transfer across the membrane. During the photocycle, the *cis-trans* isomerization of the retinal (which remains bound to the protein) leads to the movement between inward and outward channels of the Schiff base, which can be reversibly protonated. This movement shuttles the protons across the membrane dielectric and leads to the active transport.

In *Halobacterium*, two homologues of bacteriorhodopsin are involved in **membrane transport** of anions (halorhodopsin) and in phototaxis (sensory rhodopsin) (5). The acidic channel residues of bacteriorhodopsin are substituted with other amino acids, such as [threonine](#) in halorhodopsin, which contributes toward the maintenance of osmotic balance in an environment with high salinity by **active transport** of chloride and nitrate. The sensory rhodopsin is involved in the control of light-induced **flagellar** motion. The exact mechanism of its function is not known. It is obvious, however, that the same rhodopsin fold can be employed in different biological functions through a small number of changes in the internal channels.

## Bibliography

1. H. G. Khorana (1992) Rhodopsin, photoreceptor of the rod cell. *J. Biol. Chem.* **267**, 1–4.
2. D. L. Farrens, C. Altenbach, K. Yang, W. L. Hubbell, and H. G. Khorana (1996) Requirement of rigid-body motion of transmembrane helices for light activation of rhodopsin. *Science* **274**, 768–770.

3. G. F. X. Schertler, C. Villa, C., and R. Henderson (1993) Projection structure of rhodopsin. *Nature* **362**, 770–772.
4. V. M. Unger, P. A. Hargrave, J. M. Baldwin, and G. F. X. Schertler (1997) Arrangement of rhodopsin transmembrane alpha-helices. *Nature* **389**, 203–206.
5. D. Oesterhelt (1998) Structure and mechanism of the family of retinal proteins from halophilic archaea. *Curr. Opin. Struct. Biol.* **8**, 489–500.

## Ri Plasmid

The root-inducing, or Ri, plasmid is the cause of hairy root formation on dicotyledenous **plants** that are infected by *Agrobacterium rhizogenes*. Similarly to the [Ti plasmid](#), a defined region of **DNA** of the Ri plasmid, the transferred or T-DNA (see **T-complex**), is transferred from the **bacteria** to the infected plant cell and integrated into the plant [genome](#).

Ri plasmids have been less studied than the Ti plasmids. Like the Ti plasmid, Ri plasmids can be characterized by their **opine** specificity, agropine or mannopine. There are similarities in sequence between Ti and Ri plasmids, most notably in the T-DNA and the virulence (*vir*) regions ([1](#)). Transfer of the T-DNA from Ri plasmids to the plant cell most probably uses the same mechanism as T-DNA transfer from Ti plasmids ([2](#)). **Vectors** for transferring foreign DNA to plant cells based on Ri plasmids have been developed for use in [plant genetic engineering](#) ([3](#)).

## Bibliography

1. G. A. Huffman, F. F. White, M. Gordon, and G. Nester (1984) *J. Bact.* **157**, 269–276.
2. P. Zambryski, J. Tempe, and J. Schell (1989) *Cell* **56**, 193–201.
3. J. Stougaard, D. Abildsten, and K. A. Marcker (1987) *Mol. Gen. Genet.* **207**, 251–255.

## Suggestion for Further Reading

4. G. Kahl and J. Schell (1982) *Molecular Biology of Plant Tumors* Academic Press, London.

## Ribonuclease H

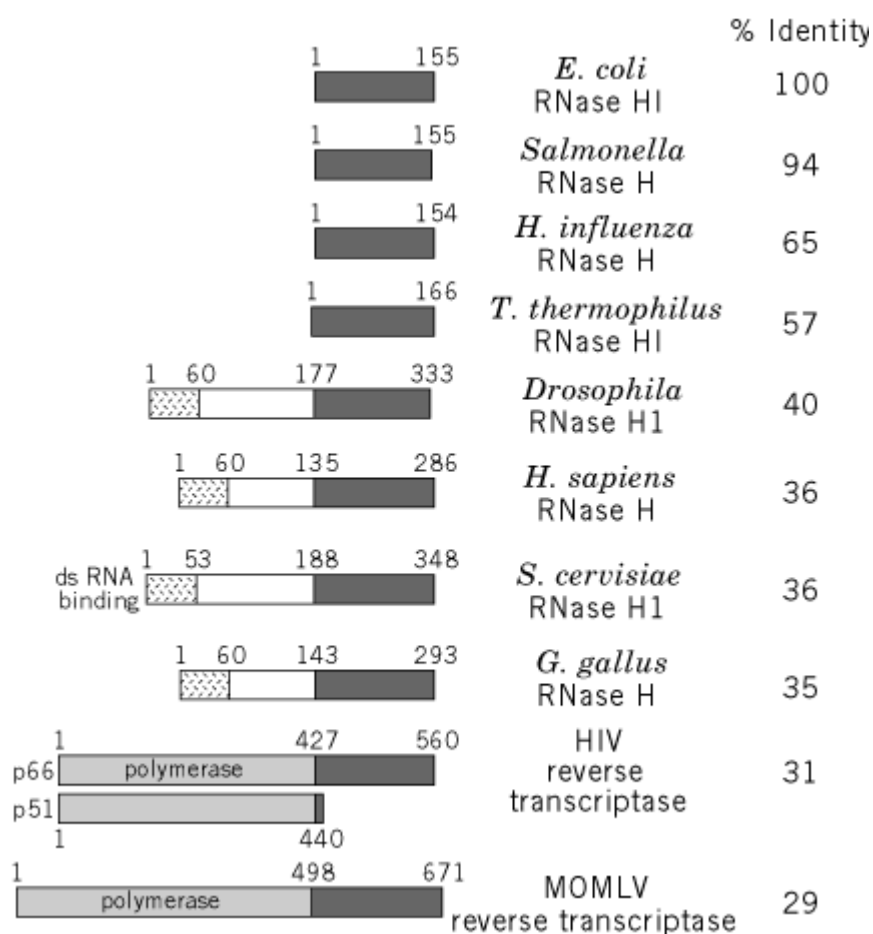
Ribonuclease H (RNase H) is an [enzyme](#) that cleaves RNA from RNA•DNA hybrids in a divalent cation ( $Mg^{2+}$  or  $Mn^{2+}$ )-dependent manner. RNase H (H for “hybrid”) can act in both an endonucleolytic and an exonucleolytic manner, yielding 3′ hydroxyl and 5′ phosphate groups as products (unlike the prototypical ribonuclease, **ribonuclease A**, which leaves 3′ phosphates and 5′ hydroxyls). The 3′ hydroxyl groups left by RNase H are then competent substrates for both **DNA polymerase** and [DNA Ligase](#). RNase H activity is important for retroviral reverse transcription by [retroviruses](#), removal of RNA primers in [DNA replication](#), and the action of [antisense oligonucleotides](#). RNase H is commonly used in molecular biology applications where cleavage of RNA in a heteroduplex with DNA is required, such as second-strand synthesis of [complementary](#)



[DNA](#) (cDNA) or in the selective targeting of RNA by the hybridization of a sequence-specific oligonucleotide.

There are several classes of RNase H proteins, but the best-characterized is a large family of structurally-similar enzymes that exist in prokaryotes (including proteins whose genes have been cloned (see [Cloning](#)) from *Escherichia coli*, *Salmonella*, and *Thermus thermophilus*) and in eukaryotes (including yeast, trypanosomes, *Drosophila*, chickens, mice, and humans). This class also includes an essential domain of retroviral reverse transcriptases (RT) (Fig. 1). These RNases H are small (~150 residues) and have a mixed a/b [protein structure](#). The prokaryotic members of this family contain only the RNase H domain, whereas eukaryotic RNases H contain an additional N-terminal domain that is known to bind double-stranded RNA. In retroviral reverse transcriptases, the RNase H domain is covalently linked to the N-terminal polymerase domain.

**Figure 1.** Primary structures of various ribonucleases H. The RNase H domain is shaded dark; other domains are indicated. The percent identity of each RNase H domain to the sequence of *E. coli* RNase HI is indicated.



## 1. Roles of RNase H

RNase H plays diverse roles in biology, some of which are not yet fully understood. The requirements for RNase H as part of RT in the replication of retroviruses, such as human immunodeficiency virus ([HIV](#)) and Moloney murine leukemia virus (MoMLV), are well-understood. Whereas MoMLV RT is a monomeric, multi-domain protein, HIV RT is a heterodimer composed of two polypeptide chains: p66 and p51 (Fig. 1). The 66-kDa subunit contains both the polymerase and RNase H domains, but the 51-kDa subunit carries only the polymerase domain and lacks the RNase

H domain, which was removed by HIV protease. In all reverse transcriptases, RNase H functions to degrade the genomic RNA following its use as a template for DNA synthesis by the polymerase domain. This occurs in a roughly sequence-independent, exonucleolytic manner, although more specific patterns of endonucleolytic cleavage are required for removal of the initial [transfer RNA](#) primer and the creation of a second-strand primer (see [Reverse Translation](#)). A number of mutations dramatically alter the exonucleolytic function, but not the endonucleolytic function ([1-4](#)) suggesting some independence in the two modes of RNase H cleavage.

The biological importance of RNase H in bacteria and yeast is less clear. These proteins are not essential for viability, so it has been difficult to pinpoint their specific roles. Two RNases H exist in *E. coli*: RNase HI and HII. RNase HI prevents initiation of replication at origins other than *oriC* ([5-7](#)). Any false origins generated by RNA polymerase are eliminated by RNase H, so therefore RNase HI serves as a specificity factor for replication initiation. Furthermore, mutations in *dnaA*, the gene for a protein needed for initiation of DNA replication at the normal origin, *oriC*, are suppressed by the loss of RNase HI activity. RNase HI is also required for viability in cells that lack functional RecBCD proteins ([8](#)). Interestingly, RecBCD is a **nuclease** involved in [DNA repair](#), but whether RNase H is needed in some capacity for DNA recombination and repair is not yet certain. *E. coli* RNase HI also contributes to the required removal of RNA primers used in the leading and lagging strands of DNA replication. The vast majority of this activity, however, is contributed by the 5' to 3' exonuclease domain of DNA polymerase I. Finally, RNase HI is necessary for the replication of ColE1-based **plasmids**. Following production of an RNA primer by RNA polymerase on ColE1, RNase HI processes the RNA, making it competent for elongation by DNA polymerase. As mentioned earlier, a second enzyme, *E. coli* RNase HII, has also been identified ([9](#)). This protein has very weak activity and no sequence homology with RNase HI, and its cellular significance is yet unknown.

The budding yeast *Saccharomyces cerevisiae* contains an RNase H1 enzyme with a C-terminal domain similar to that of *E. coli* RNase HI and an N-terminal region that contains a 50-residue RNA-binding motif that is also seen in [Cauliflower Mosaic Virus](#) P6 protein. This region binds double-stranded RNA in the absence of the RNase H domain ([10](#)) (Fig. [1](#)). Interestingly, increasing  $Mg^{2+}$  levels decrease the enzyme's affinity for double-stranded RNA while increasing its activity on RNA•DNA hybrids. It is unknown, however, if or how this relates to the control of the RNase H portion, which cleaves RNA only within heteroduplexes with DNA.

In mammalian cells, two different classes of RNases H have been identified. Type I RNases H are larger (68 to 90 kDa) and can utilize either  $Mg^{2+}$  or  $Mn^{2+}$  as cofactors. Type II RNases H are only ~40kDa and are inhibited by the presence of  $Mn^{2+}$ . Type-I RNase H from calf thymus acts *in vitro* as an endonuclease to cleave RNA replication primers one nucleotide from the RNA/DNA junction ([11](#)). Gene sequences encoding type-II RNases H in mice and humans have been deposited into the GenBank database (see [Sequence Databases](#)). They are very similar to other RNases H, including those from bacteria, lower eukaryotes, and retroviruses, and they also contain an N-terminal domain similar to the yeast duplex RNA-binding region. Given their similarity to prokaryotic RNases HI, these mammalian RNases H are probably better termed RNase HI despite the original "Type II" designation applied to these activities prior to the cloning of their genes.

RNase H activity is also required for the action of [antisense oligonucleotides](#). Antisense technology works to eliminate specific populations of mRNA *in vivo*. Complementary DNA oligonucleotides are introduced into the cells of interest and hybridize to the target mRNA; this creates a RNA•DNA hybrid suitable for RNase H hydrolysis. RNase H cleavage of the hybrid then specifically destroys the mRNA target.

## 2. Structure/function of RNase H

The three-dimensional structure of RNases H from *E. coli*, *Thermus thermophilus*, and HIV have

been determined by [X-ray crystallography](#) (12-15). Each has a similar fold containing [alpha-helices](#) flanking a five-stranded beta-sheet (see [Beta-Pleated Sheet](#)). RNase H belongs to a superfamily of structurally similar nucleic acid-modifying proteins termed the “polynucleotidyl transferases” (16, 17). In addition to RNase H, members of this superfamily include retroviral [integrases](#), Mu [transposase](#), and *E. coli* RuvC. These proteins all require metal ion cofactors, leave 3' hydroxyl and 5' phosphate groups, and probably share a similar catalytic mechanism.

A combination of biochemical and structural experiments has elucidated the importance of several specific regions to the catalytic function and structural stability of RNase H. Five highly conserved residues (three aspartate, a glutamate, and a histidine residue) map to similar metal-binding active sites in each protein. The “basic helix/loop” region seen in both *E. coli* and *T. thermophilus* RNases H contains a large number of positively-charged residues that contribute to the enzyme's affinity for its nucleic acid substrate (18). In a manner that is not yet understood, RNase H is able to discriminate RNA in a heteroduplex with DNA from duplex RNA and from single-stranded RNA. Regions of the protein that interact with the substrate have been mapped using [NMR](#) analysis (19). Nonetheless, a high-resolution description of RNase H binding nucleic acid would help resolve questions regarding its specificity in catalysis.

Despite the importance of the basic helix/loop to substrate recognition in *E. coli* RNase HI, this region is absent from the HIV RNase H domain. The HIV RNase H domain is inactive when expressed independently of the polymerase domain. Grafting the *E. coli* RNase H basic helix/loop onto the HIV domain partially reactivates it (20, 21), suggesting that the inactivity of the isolated RNase H domain is due, at least in part, to lack of substrate affinity. Furthermore, the C-terminal helix and loop of the HIV protein is known to be dynamically disordered in solution (22). [Site-directed mutagenesis](#) of this region in *E. coli* has demonstrated its importance for the activity of the protein (23-25), suggesting that the destabilized C-terminal region also contributes to the inactivity of the isolated HIV RNase H domain.

A single  $Mg^{2+}$  ion binds the active site of *E. coli* RNase HI (12, 26), whereas two  $Mn^{2+}$  ions were identified in the HIV RNase H crystal structure (15). Based upon these results, several different catalytic mechanisms have been proposed (15, 19, 23). It is still uncertain whether one or two metal ions are required for catalysis, but recent experiments support the idea that one is needed to activate the enzyme and that a second metal ion can inhibit this activity (27). A large number of mutations that inactivate the enzyme with  $Mg^{2+}$  do not eliminate its  $Mn^{2+}$ -dependent activity (20, 21, 25, 28), suggesting potential differences in the way metal ions activate RNase H. Interestingly, while  $Mg^{2+}$  has traditionally been considered to be the more physiologically relevant cation, due to its higher cellular concentrations (~1-5mM), *E. coli* RNase HI requires a much lower concentration of  $Mn^{2+}$  (~2  $\mu$ M) to achieve maximum activity (27, 28).

RNase H also serves as an excellent model for studies on **protein folding** and [protein stability](#) (29).

## Bibliography

1. B. M. Wohrl, S. Volkmann, and K. Moelling (1991) *J. Mol. Biol.* **220**, 801–818.
2. S. Volkmann, B. M. Wohrl, M. Tisdale, and K. Moelling (1993) *J. Biol. Chem.* **268**, 2674–2683.
3. N. M. Cirino, C. E. Cameron, J. S. Smith, J. W. Rausch, M. J. Roth, S. J. Benkovic, and S. F. LeGrice (1995) *Biochemistry* **34**, 9936–9943.
4. M. Ghosh, K. J. Howard, C. E. Cameron, S. J. Benkovic, S. H. Hughes, and S. F. Le Grice (1995) *J. Biol. Chem.* **270**, 7068–7076.
5. B. de Massy, O. Fayet, and T. Kogoma (1984) *J. Mol. Biol.* **178**, 227–236.
6. G. Lindahl and T. Lindahl (1984) *Mol. Gen. Genet.* **196**, 283–289.
7. T. Ogawa, G. G. Pickett, T. Kogoma, and A. Kornberg (1984) *Proc. Natl. Acad. Sci. USA* **81**,

1040–1044.

8. M. Itaya and R. J. Crouch (1991) *Mol. Gen. Genet.* **227**, 424–432.
9. M. Itaya (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8587–8591.
10. S. M. Cerritelli and R. J. Crouch (1995) *RNA* **1**, 246–259.
11. R. A. Bambara, R. S. Murante, and L. A. Henricksem (1997) *J. Biol. Chem.* **272**, 4647–4650.
12. W. Yang, W. A. Hendrickson, R. J. Crouch, and Y. Satow (1990) *Science* **249**, 1398–1405.
13. K. Katayanagi, M. Miyagawa, M. Matsushima, M. Ishikawa, S. Kanaya, M. Ikehara, T. Matsuzaki, and K. Morikawa (1990) *Nature* **347**, 306–309.
14. K. Ishikawa, M. Okumura, K. Katayanagi, S. Kimura, S. Kanaya, H. Nakamura, and K. Morikawa (1993) *J. Mol. Biol.* **230**, 529–542.
15. J. Davies, Z. Hostomska, Z. Hostomsky, S. R. Jordan, and D. A. Matthews (1991) *Science* **252**, 88–95.
16. W. Yang and T. Steitz (1995) *Structure* **3**, 131–134.
17. P. Rice, R. Craigie, and D. R. Davies (1996) *Curr. Opin. Struct. Biol.* **6**, 76–83.
18. S. Kanaya, N. C. Katsuda, and M. Ikehara (1991) *J. Biol. Chem.* **266**, 11621–11627.
19. H. Nakamura et al. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 11535–11539.
20. S. J. Stahl, J. D. Kaufman, T. S. Vikić, R. J. Crouch, and P. T. Wingfield (1994) *Protein Eng* **7**, 1103–1108.
21. J. L. Keck and S. Marqusee (1995) *Proc. Natl. Acad. Sci. USA* **92**, 2740–2744.
22. R. Powers, G. M. Clore, A. Bax, D. S. Garrett, S. J. Stahl, P. T. Wingfield, and A. M. Gronenborn (1991) *J. Mol. Biol.* **221**, 1081–1090.
23. Y. Oda, M. Yoshida, and S. Kanaya (1993) *J. Biol. Chem.* **268**, 88–92.
24. M. Haruki, E. Noguchi, C. Nakai, Y. Y. Liu, M. Oobatake, M. Itaya, and S. Kanaya (1994) *Eur. J. Biochem.* **220**, 623–631.
25. E. R. Goedken, T. M. Raschke, and S. Marqusee (1997) *Biochemistry* **36**, 7256–7263.
26. C. B. Black and J. A. Cowan (1994) *Inorg. Chem.* **33**, 5805–5808.
27. J. L. Keck, E. R. Goedken, and S. Marqusee (1998) *J. Biol. Chem.* **273**, 34128–34133.
28. J. L. Keck and S. Marqusee (1996) *J. Biol. Chem.* **271**, 19883–19887.
29. A. K. Chamberlain and S. Marqusee (1997) *Structure* **5**, 859–863.

### **Suggestions for Further Reading**

30. R. J. Crouch (1990) Ribonuclease H: from discovery to 3D structure. *New Biology* **2**, 771–777.
31. Z. Hostomsky, Z. Hostomska, and D. A. Matthews (1993), "Ribonucleases H" in *Nucleases*, 2nd Edition, Cold Spring Harbor Laboratory Press, Plainview, N.Y., pp. 341–376.
32. J. M. Whitcomb and S. H. Hughes (1992), Retroviral reverse transcription and integration: Progress and problems, *Ann. Rev. Cell Biology* **8**, 275–306.

### **Ribonuclease P**

[Transfer RNA](#) (tRNA) molecules are synthesized as long RNA transcripts that must be trimmed down to the appropriate size and then processed further (see [tRNA Biosynthesis](#)). Removal of the 5'-leader sequence is catalyzed by the enzyme ribonuclease P (RNaseP). Hydrolysis of the

phosphodiester bond yields a mature tRNA carrying a 5'-phosphate and a leader sequence carrying a 3'-hydroxyl. Early genetic (1, 2) and biochemical experiments (3) established the *Escherichia coli* enzyme as a [ribonucleoprotein](#), consisting of both RNA and protein. In 1983, Guerrier-Takada et al. (4) made the remarkable discovery that RNA from the *E. coli* and *Bacillus subtilis* RNaseP, in the absence of their protein subunits, could catalyze the removal of 5' leaders from tRNA precursors *in vitro*. The elements necessary for binding and catalysis were RNA-specific, because the RNA alone had the same affinity for pre-tRNA as the holoenzyme (5). Such catalytically active RNA molecules are now known as [ribozymes](#). The RNA from bacterial RNaseP is generally able to process tRNA precursors in an *in vitro*, protein-free reaction, but that from **eukaryotes** and **archaeobacteria** have not yet been found to be catalytically active in the absence of their protein partners. RNA subunits have been identified in the RNaseP from [mitochondria](#) (6), but not from [chloroplasts](#) (7), leaving open the possibility that 5'-tRNA processing in chloroplasts is fundamentally different.

Insight into the structure of the bacterial RNaseP RNA has been possible because a large database supports convincing alignments of primary sequence and covariation analyses necessary to identify homologous nucleotides on the basis of location in the secondary structure (8). The continuously growing database of bacterial RNA sequences (8) strengthens the validity of the inferred secondary structure and identifies possible elements of tertiary structure. In spite of the differences in primary sequence and the occurrence of discrete helical elements at distinct sites, a core containing similar sequences and secondary structure has been identified. Indeed, a 211-nucleotide RNA missing all the evolutionarily variable structures found in the smallest, naturally occurring, catalytically active 276-nucleotide RNA from *Mycoplasma fermentans* (9) retains catalytic activity and defines a phylogenetic minimal RNaseP RNA. This and other such "minimal RNA" require increased ionic strength and display lower thermal stability, indicating that the variant helix elements contribute to the stability of essential structural domains that contain the active site of the ribozyme.

No single nucleotide is essential to catalysis (10), suggesting that, although specific nucleotides are important to promote interactions between structural domains, the elements responsible for binding and catalysis are distributed through the sequence. Recent additions to the archaeobacteria database reveal that these RNA are similar in primary and secondary structure to those from **eubacteria**, in spite of their great evolutionary distance (11). Although secondary structure models of eukaryotic RNA (12) produced to date do not have the rigor of the bacterial and archaeobacteria models, it is clear that all RNaseP RNA share short elements of primary sequence similarity and are predicted to share some common structural features (13). This supports the conclusion that they are all related and, despite some lacking *in vitro* RNA activity, serve as the catalytic subunit of RNaseP. Two three-dimensional structures of the *E. coli* RNaseP RNA, on the basis of chemical probing and [crosslinking](#), in conjunction with secondary structure models, have been proposed (14, 15) but the challenges of arriving at a definitive structure and identifying the [active site](#) remain.

The Cs in the CCA 3' end of tRNA interact with the well-conserved G292 and G293 through canonical base pairing in the RNA from *E. coli* and other bacteria (16). Many crosslinking, [footprinting](#), and binding experiments (17) have identified nucleotides that interact with either the tRNA substrate and/or the protein subunit. While they are found in separate regions of the secondary structure model, the tertiary models bring them together in the conserved core. Whether the bases identified in these studies function directly in catalysis or are required for structure, protein binding, or metal binding is currently under investigation. Despite ease of assay, a refined secondary structure and the rapidity that *in vitro* reactions using RNA alone bring to mutagenic and chemical probing experiments, a mechanism has yet to be established.

All RNaseP RNAs are made as precursors. In *E. coli*, ribonuclease E is responsible for processing the 3' end and, although it is not known what processes the 5' end, ribonucleases E, P, and III have been eliminated (18). In the yeast nucleus, Rpr1r is made as a precursor with an 84-nucleotide 5' leader and a 3' trailer of about 33 nucleotides (19). In yeast mitochondria, Rpm1r is cotranscribed

with two flanking tRNA molecules. Although separation of the tRNA from the transcript does not depend on Rpm2p, the protein subunit of the mitochondrial enzyme, Rpm2p, has a critical function in further RNA processing steps necessary to make mature Rpm1r (20).

Given that the RNA is the catalytic subunit of RNaseP, what is the function of the RNaseP protein (s)? The *in vivo* role of the protein is not yet known, but it is essential. *In vitro* it stimulates the rate of cleavage and of product release, and it affects substrate specificity (21, 22). One way that the protein might affect the reaction in these ways is by counteracting electrostatic repulsions between RNA phosphates in the substrate and the RNA subunit of the enzyme (4, 23). Or, it may stabilize the structure of the RNaseP RNA. The fact that the presence of the protein can overcome many of the mutations that affect the activity of the RNA by itself supports the idea that it serves as a stabilizer for the RNaseP RNA structure (24).

The RNaseP proteins identified to date are structurally diverse. The 14-kDa (kilodalton) eubacterial proteins have a core of highly conserved basic and **hydrophobic** residues and semiconserved and conserved aromatic residues (25). Stoichiometry of the RNA and protein subunits is 1:1 (25). Only two eukaryotic RNaseP proteins, both yeast proteins, Pop1p (26) and Rpm2p (27), had been identified by the end of 1996. Rpm2p appears to be the major protein subunit of the **mitochondrial** RNaseP, whereas Pop1p appears to be one of several proteins important to RNaseP activity in the **nucleus**. Mutations in *POPI* cause defects in the RNA processing enzymes MRP and nuclear RNaseP (28). While *S. cerevisiae* MRP and nuclear RNaseP compositions appear complex, it is interesting that purification of nuclear RNaseP from *S. pombe* reveals a single protein the size of Pop1p (29). Rpm2p copurifies with mitochondrial RNaseP activity (30), and biochemical, immunological, and genetic evidence demonstrates that Rpm2p is a subunit of mitochondrial RNaseP (30). Rpm2p and Pop1p are about the same size, but approximately 100 kDa larger than the bacterial proteins. They do not show any obvious sequence similarity with each other or with the bacterial proteins.

## Bibliography

1. P. Schedl and P. Primakoff (1973) Proc. Natl. Acad. Sci. USA **70**, 2091.
2. H. Sakano, S. Yamada, T. Ikemura, Y. Shimura, and H. Ozeki (1974) Nucleic Acids Res. **1**, 355.
3. B. C. Stark, R. Kole, E. J. Bowman, and S. Altman (1978) Proc. Natl. Acad. Sci. USA **75**, 3717–3721.
4. C. Guerrier-Takada, K. Gardiner, T. Marsh, N. Pace, and S. Altman (1983) Cell **35**, 849–857.
5. C. Reich, G. J. Olsen, B. Pace, and N. R. Pace (1988) Science **239**, 178–181.
6. C. A. Wise and N. C. Martin (1991) J. Biol. Chem. **266**, 19154.
7. M. J. Wang, N. W. Davis, and P. Gegenheimer (1988) EMBO J. **7**, 1567–1574.
8. J. W. Brown, J. M. Nolan, E. S. Haas, M. A. Rubio, F. Major, and N. R. Pace (1996) Proc. Natl. Acad. Sci. USA **93**, 3001–3006.
9. R. W. Siegel, A. B. Banta, E. S. Haas, J. W. Brown, and N. R. Pace (1996) RNA **2**, 452–462.
10. D. S. Waugh and N. R. Pace (1993) FASEB J. **7**, 188–195.
11. E. S. Haas, D. W. Armbruster, B. M. Vucson, C. J. Daniels, and J. W. Brown (1996) Nucleic Acids Res. **24**, 1252–1259.
12. A. J. Tranguch and D. R. Engleke (1993) J. Biol. Chem. **268**, 14045–14055.
13. J. L. Chen and N. R. Pace (1997) RNA **3**, 557–560.
14. M. E. Harris, J. M. Nolan, A. Malhotra, J. W. Brown, S. C. Harvey, and N. R. Pace (1994) EMBO J. **13**, 3953–3963.
15. E. Westhof and S. Altman (1994) Proc. Natl. Acad. Sci. USA **91**, 5133–5137.
16. S. G. Svard, U. Kagardt, and L. A. Kirsebom (1996) RNA **2**, 463–472.
17. N. R. Pace and J. W. Brown (1995) J. Bacteriol. **177**, 1919–1928.

18. U. Lundberg and S. Altman (1995) *RNA* **1**, 327–334.
19. J.-Y. Lee, C. E. Rohlman, L. A. Molony, and D. R. Engelke (1991) *P. Mol. Cell. Biol.* **11**, 721–730.
20. V. Stribinskis, G.-J. Gao, P. Sulo, Y.-L. Dang, and N. C. Martin (1996) *Mol. Cell Biol.* **16**, 3429–3436.
21. J. W. Brown and N. R. Pace (1992) *Nucleic Acids Res.* **20**, 1451–1456.
22. S. Altman, L. Kirsebom, and S. Talbot (1995) in *tRNA: Structure, Biosynthesis, and Function*, D. Söll and U. L. RajBhandary, eds., American Society for Microbiology Press, Washington, D.C., Vol.6, pp. 67–78.
23. K. J. Gardiner, T. L. Marsh, and N. R. Pace (1985) *J. Biol. Chem.* **260**, 5415–5419.
24. N. Lumelsky and S. Altman (1988) *J. Mol. Biol.* **202**, 443–454.
25. S. Talbot and S. Altman (1994) *Biochemistry* **33**, 1406–1411.
26. Z. Lygerou, P. Mitchell, E. Petfalski, B. Seraphin, and D. Tollervey (1994) *Genes Devel.* **8**, 1423–1433.
27. M. J. Morales, Y. L. Dang, Y. C. Lou, P. Sulo, and N. C. Martin (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9875–9879.
28. Z. Lygerov, P. Mitchell, E. Petfalski, B. Seraphin, and D. Tollervey (1994) *Genes. Dev.* **12**, 1423–1433.
29. S. Zimmerly, D. Drainas, L. A. Sylvers, and D. Soll (1993) *Eur. J. Biochem.* **217**, 501–507.
30. Y. L. Dang and N. C. Martin, *J. Biol. Chem.* **268**, 19791–19796.

## Ribonucleoprotein (RNP)

Ribonucleoprotein (RNP) is widely defined as a biologically important molecular assembly composed of both **RNA** and [protein](#). Such RNP complexes are key elements in many steps throughout the expression of DNAs genetic information into functional proteins. For example, protection of the ends of [chromosomes](#) against [DNA degradation](#) is mediated by the enzyme [telomerase](#), a protein-RNA complex that requires RNA for its activity. The replication of RNA viruses involves interactions between RNA-dependent **RNA polymerase** or **reverse transcriptase** and viral RNAs. The machinery for gene [transcription](#) contains RNA polymerase (I, II, or III) in action on a DNA template, with associated nascent RNA. Further processing of the nascent RNAs (eg, [transfer RNA](#) modification; **ribosomal RNA** modification, processing, and storage; processing of RNA polymerase II transcripts including [RNA splicing](#)) also occur while the RNAs are associated with their respective specific proteins in the form of RNP complexes. **Protein biosynthesis** takes place within a distinct class of RNP complexes—the [ribosomes](#). Finally, specific proteins associate with [messenger RNA](#) to control and perform its degradation. The term “RNP” was originally coined for the complex that assembles on the precursor messenger RNA (pre-mRNA) in the nucleus of eukaryotic cells. Later it was realized that these complexes carry out the processing of nuclear pre-mRNA. This chapter focuses on this class of RNP complexes.

Most eukaryotic transcripts that code for proteins are transcribed in the nucleus by RNA polymerase II to yield pre-mRNA molecules. These primary transcripts undergo several obligatory processing events within the nucleus to produce mature messenger RNAs (mRNA) that can then exit from the nucleus into the cytoplasm to be translated into proteins. These nuclear processing activities include 5'-end capping, pre-mRNA splicing, 3'-end processing, [RNA Editing](#), and, finally, RNA transport

from the nucleus to the cytoplasm. During transcription, the nascent pre-mRNA is packaged with proteins and additional components to form RNP particles, in which the transcript is processed and remains packaged in the nucleus until it is exported to the cytoplasm. The packaging of the pre-mRNA in nuclear RNP particles is proposed to play a dual role. On one hand, it protects the pre-mRNA from untimely degradation. On the other hand, it provides the machinery that enables accurate processing at defined and precise locations within the pre-mRNA sequence, and it facilitates the processing of the pre-mRNA in a temporally and spatially regulated fashion. It is therefore expected that the native nuclear RNP particle will indeed harbor all the components that are required for the above-mentioned processing activities.

The term RNP was originally given, in the early 1960s, to the complexes in which nuclear transcripts of RNA polymerase II were found packaged with proteins (1). In early studies, these transcripts were collectively termed heterogeneous nuclear RNA (hnRNA) due to their diverse length. Accordingly, the RNP particles that package hnRNA were termed heterogeneous nuclear RNP (hnRNP) particles (2). At the time this terminology was coined, some of the pre-mRNA processing events had not yet been discovered. Today, a more accurate definition of RNA polymerase II transcripts is pre-mRNA, and the respective RNP particles are thus defined as pre-mRNP particles.

RNP particles that package mRNA in the cytoplasm are termed mRNP particles. The protein compositions of nuclear pre-mRNP and cytoplasmic mRNP particles are thought to be different, although recent studies indicate that some of the nuclear proteins accompany the mRNA into the cytoplasm (3). Although the mRNP particles have not yet been fully characterized, they likely play a role in the translation and stability of the mRNA, and they may also have an effect on the localization of the mRNA in the cytoplasm.

The extensive research on pre-mRNA splicing *in vitro* led to the characterization of a group of uridine-rich, small, nuclear RNAs (U snRNAs) as essential pre-mRNA splicing factors (4, 5). For this function, the U snRNAs are packaged with specific sets of proteins to form complexes that are called U small nuclear RNP (U snRNP) complexes (6). Purified snRNP particles have [sedimentation coefficients](#) of 12–25 S in density gradients. The studies of splicing *in vitro* led to the definition of yet another RNP complex—the 60 S [spliceosome](#), which is the RNP complex that assembles on one **intron** flanked by two exons and is capable of removing that intron. The spliceosome is composed of spliceosomal U snRNPs and non-snRNP proteins (7-9).

Early biochemical and electron microscopy studies dealt mainly with the general population of pre-mRNP complexes. With the development in the early 1980s of **genetic engineering** techniques, the study of specific RNP particles was made possible through the use of probes for specific pre-mRNAs. This led to the identification of RNP particles in which specific pre-mRNAs, as well as most **polyadenylated** nuclear RNAs, are packaged. These particles sediment in sucrose density gradients as 200 S complexes and are thus termed large nuclear RNP (lnRNP) particles (10). In the following sections, early studies of RNP particles will be discussed, with emphasis on mammalian nuclear RNP particles. lnRNP particles will be described in detail, and their relationship to other RNP complexes implicated in pre-mRNA splicing will be discussed.

## 1. Early RNP models

In situ [electron microscopy](#) studies of [chromatin](#), using sections of fixed and stained nuclei or cells, revealed already in the 1950s that nascent RNA transcripts are packaged in RNP particles (11-13). Early attempts in the 1960s and 1970s to isolate the nuclear RNP complexes of <sup>3</sup>H-uridine-labeled nascent RNA polymerase II transcripts used mild digestion by exogenous or endogenous ribonucleases (RNases). Fractionation of the resultant nuclear supernatants in sucrose gradients gave RNP complexes that sedimented in the 30 S–40 S region of the gradients. These RNP complexes were composed of RNA of variable length (500–1000 bases), and their major protein component was a set of six proteins of 35–45 kDa, designated A1, A2, B1, B2, C1, and C2 (1, 14-16). The short



length of the RNA suggested that the 40 S RNP particles were the products of degradation of the native pre-mRNP particles by endogenous **nucleases**. These observations raised the intriguing question of whether the 40 S particles represented a structural element of the native complexes or they were the outcome of an uncontrolled degradation. Attempts to inhibit the endogenous ribonucleases with specific inhibitors gave a heterogeneous population of RNP complexes, sedimenting at 40 S–250 S in sucrose gradients, which were later termed hnRNP particles. The protein composition of hnRNP particles was similar to that of the 40 S particles, but the hnRNP RNA was longer (2, 17). Electron micrographs of such preparations were interpreted as displaying a repeating unit (18). This led to a “beads-on-a-string” model for the hnRNP particle, in which, by analogy to **chromatin**, the 40 S particle was analogous to the **nucleosome**, and the six hnRNP core proteins were analogous to the **histone** core proteins. Although this conclusion seemed logical at the time, it was not corroborated by subsequent systematic visualization of spread preparations of chromatin (19). Furthermore, because these studies dealt with the general, uridine-labeled population of nuclear RNA, whether or not the sizes of specific pre-mRNAs were preserved during the RNP isolation was not known. It was therefore impossible to determine unambiguously whether the 40 S–250 S hnRNP particles represented the native pre-mRNP particles or degradation intermediates thereof.

## 2. Protein Composition of hnRNP Particles

Determining the protein composition of RNP particles has been hampered by difficulties in isolating intact particles and by the uncertainty in the degree of preservation of the pre-mRNA during isolation. A methodology that circumvents these problems involves identifying nuclear RNP proteins that are directly associated with pre-mRNAs by UV-induced **crosslinking** in situ and isolation of polyadenylated RNA under denaturing conditions. The proteins that are crosslinked to the RNA can thus be identified. Using this method, several proteins were identified as associated with pre-mRNA, including the 30–45-kDa proteins previously identified as A1, A2, B1, B2, C1, and C2, as well as proteins of 120, 68, and 53 kDa (3, 20). **Monoclonal antibodies** were raised against these proteins and were used for indirect **immunoprecipitation** of proteins complexed in RNP particles. About 20 RNP proteins were identified and were designated alphabetically from hnRNP protein A to hnRNP protein U (21). The hnRNP proteins are very abundant. They include not only the previously identified A, B, and C proteins but also additional abundant proteins. Most hnRNP proteins bind RNA, as was shown by their UV-induced crosslinking and their association with single-stranded DNA. Their association with RNA is attributed to the presence of an RNA-binding domain (RBD) in most of the hnRNP proteins, or to the presence of additional RNA binding motifs such as the KH motif and an RGG box (Arg-Gly-Gly repeats interspersed with several aromatic residues) (22) (see **RNA-Binding Proteins**). The hnRNP proteins are very diverse owing to the presence of **isoforms** generated by **alternative splicing** and also owing to **post-translational modifications**, such as **phosphorylation** of **serine** and **threonine** residues (eg, the A/B, C, and U proteins) and **methylation** of **arginine** residues (eg, A1 and A2 and potential sites in U and K proteins). The composition of human hnRNP proteins (21) is described in the text below.

### 2.1. A/B Proteins

These proteins have apparent molecular weights of 34–40 kDa and isoelectric points (pI) of 8.4 to 9.0. A1, A2, and B1 each have two RBD motifs and a glycine-rich domain at the **C-terminus**. A2 and B1 are identical, except for an additional 12 residues at the **N-terminus** of B1, and are probably derived from the same gene by alternative splicing. The RBD of A1 is rather **homologous** to that of A2 and B1, whereas the glycine-rich domain is more diverse. Several variants of A1 have been identified, some resulting from alternative splicing and some owing to post-translational modifications. The A/B proteins are phosphorylated *in vivo*, and the A proteins are methylated on arginine residues at the C-terminal glycine-rich domain.

### 2.2. C1/C2 Proteins

These proteins of apparent molecular weight of 41 and 43 kDa (pI = 5.9) are identical, except for 13 additional residues present in C2 in the middle of the protein; they are probably derived from the

same gene by alternative splicing. The amino terminus contains an RBD domain, whereas the C-terminal domain contains a negatively charged region. The C proteins undergo phosphorylation *in vivo*.

### 2.3. D, E, G, F/H Proteins

These proteins are less well characterized. Their respective molecular weights are: 44–48, 36–43, 43, 53, and 56 kDa, and their respective pIs are 7 to 7.8, 7.3, 9.5, and 6.1 to 7.1. RBD domains were found in the E, G, and F/H proteins.

### 2.4. I Protein

The hnRNP I protein has an apparent molecular weight of 59 kDa and a pI of 8.5. This protein is identical to PTB, a protein that binds to the polypyrimidine tract at the 3'-end of introns and is proposed to play a role in RNA splicing. The I protein has four noncanonical RBD domains. Analysis of cDNA clones revealed several isoforms that are probably derived from alternative splicing. On nuclease digestion, the I protein is more readily released from RNP complexes than other RNP proteins.

### 2.5. K/J Proteins

These proteins have apparent molecular weights of 66 and 64 kDa, and pIs of 6.1 and 6.4, respectively. The two proteins are immunologically related and lack the RBD motif but have a KH motif thought to play a role in RNA binding.

### 2.6. L Protein

This protein has an apparent molecular weight of 64–68 kDa and a pI of 7.4–7.7. It has four noncanonical RBD domains and is similar in structure to protein I.

### 2.7. M, N, P, Q, R, S, and T Proteins

These proteins are less well characterized. Their respective molecular weights are 68, 70, 72, 76–77, 82, 105, and 113 kDa, and their respective pIs are 7.8–8.2, 8.7–8.9, 9.0, 8.3, 8.0, 8.8, and 8.4.

### 2.8. U Protein

This protein has an apparent molecular weight of 120 kDa and a pI of 6.6–7.2. It is an abundant protein that undergoes phosphorylation. It does not contain an RBD motif; instead, it has an RGG motif that binds RNA.

Studies of the protein composition of pre-mRNP particles from cells of other organisms, including vertebrates, avians, and amphibians, revealed proteins of similar composition to the 20-hnRNP A–U proteins identified in humans. In addition, invertebrate RNP proteins, especially those of the fruit fly *Drosophila melanogaster*, revealed domain structures similar to those of the human RBD motif (21).

## 3. Alternative RNP Models

An important issue that should be considered in formulating a model for the native nuclear pre-mRNP particle is the discovery in the late 1970s that pre-mRNA is composed of exon and intron sequences and undergoes a series of splicing reactions. In these reactions, the introns are removed and the exons are ligated to produce the mature mRNA; this is an obligatory step to allow the export of the mRNA to the cytoplasm to code for proteins. The discovery of RNA splicing required a conceptual change in the early RNP model for two main reasons. First, the model should account for the new functional role of the pre-mRNPs; second, it should accommodate the components that are essential for splicing—ie, the U snRNPs and the other protein splicing factors (7-9). Taking this new information into account, revised models were put forward. One approach implies that the overall assembly of a pre-mRNP particle occurs by the assembly of individual 60 S spliceosomes on each intron of the pre-mRNA. According to another approach, the overall structure is built in a modular fashion from 40 S RNP particles assembled on exons and 60 S spliceosomes assembled on introns (23). Yet another model describes an RNP fibril composed of hnRNP proteins and pre-mRNA, in

which spliceosomal components are associated with intron junctions (21).

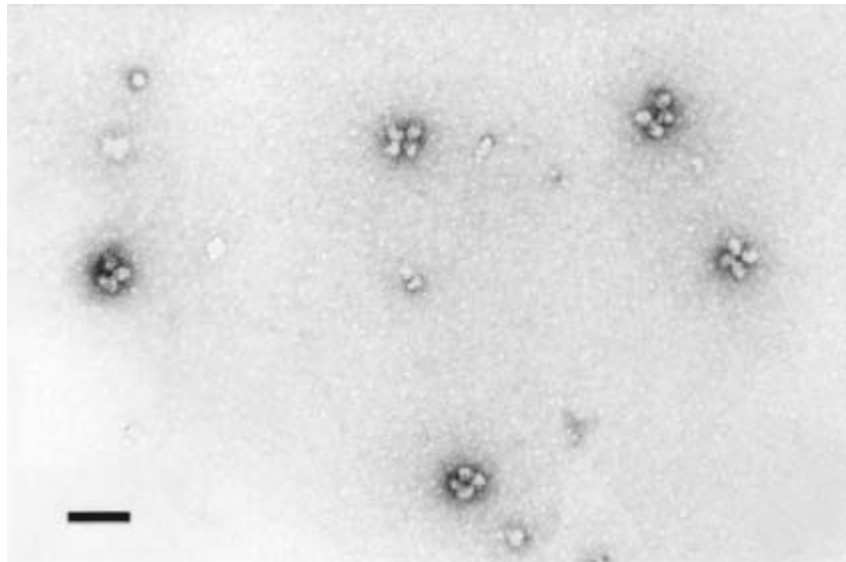
The beads-on-a-string concept that underlies some of the models of the pre-mRNP particle predicts that mild ribonuclease digestion of the RNP should yield a set of RNA fragments whose sizes approximate an integer multiple of a basic unit length (ie, an RNA ladder). This is analogous to the beads-on-a-string model of chromatin, in which a ladder of DNA fragments that are multiples of about 200 bp is generated on mild deoxyribonuclease digestion. This ladder represents the repeating unit of chromatin—the nucleosome—which is composed of a repeating unit of DNA complexed with a repeating unit of the histone core proteins (24). In contrast to this prediction, attempts to generate a ladder of RNA fragments by nuclease digestion of RNP particles did not succeed, and the 40 S RNP particles thus generated did not have a fixed length of RNA (25). Thus, a new pre-mRNP model was required—one that conforms to the present knowledge on pre-mRNA splicing and to the proposed intron and exon compositional differences.

A different, unprejudiced approach to the identification of the native pre-mRNP complex was taken with the development of a methodology for isolating native nuclear RNP particles from mammalian cell nuclei, using preservation of the entire length of specific pre-mRNAs as a criterion for nativeness (26). The protocol utilizes several potent RNase inhibitors during all steps of the preparation. It involves preparation of clean nuclei and their lysis by mild microsonication. This is followed by precipitation of the chromatin in the presence of tRNA (to avoid nonspecific adhesion of RNP proteins) and fractionation of the nuclear supernatant in a sucrose gradient. The advantage of following the fate of a specific transcript in developing the protocol for isolation of native RNP is that two criteria that had not been applied before—namely, whether the specific transcript was released into the nuclear supernatant as RNP particles in a nearly quantitative yield and in an intact, undegraded form—were now used to judge the quality of the protocol. The resulting RNP particles, and the model that was proposed based on their biochemical and structural characterization are discussed below.

#### 4. RNP Particles of Mammalian Cells—lnRNP Particles

The protocol described above was first tested by following the release of RNA encoding the multifunctional protein carbamoyl phosphate synthetase, aspartate transcarbamoylase, dihydroorotase (abbreviated CAD) from the nuclei of Syrian hamster cells into RNP particles. More than 85% of the nuclear CAD mRNA was packaged in RNP particles that sedimented at the 200 S region in a sucrose gradient (26). These 200 S particles contained not only 25-kb CAD pre-mRNA molecules but also their splicing intermediates and mature 7.9-kb CAD RNA (28). Subsequent experiments showed that other mammalian specific RNA transcripts that differ greatly in size (eg, the 36 kb pre-mRNA and the 1.6 kb mRNA of the [dihydrofolate reductase](#) (DHFR) gene, and the 3.5 kb pre-mRNA and the 1.8-kb mRNA of the [b-actin](#) gene), as well as practically all polyadenylated nuclear RNAs, were each found assembled in 200 S RNP particles (10, 27, 28). These results were corroborated by visualization by electron microscopy, which revealed large compact particles, 50 nm in diameter, composed of substructures (Fig. 1) (27). The large RNP particles were therefore termed large nuclear RNP (lnRNP) particles. The unique feature of the lnRNP particles is that different transcripts, independent of their length and the number of introns they contain (eg, 37 introns for CAD pre-mRNA (29) and 5 for b-actin (30)) are all found packaged in particles of the same size and hydrodynamic properties. These surprising results could not be explained by the beads-on-a-string RNP model or by related models, because these models predicted that the size of an RNP particle should be proportional to the length and/or number of introns of the RNA it packages.

**Figure 1.** Visualization of lnRNP particles. Electron micrographs of negatively stained particles from the 200 S peak of a HeLa cell nuclear supernatant fractionated in a sucrose density gradient. Bar = 100 nm.



To determine the composition of lnRNP particles, monoclonal antibodies directed against lnRNP components (31) as well as antibodies against U snRNP components were utilized in **Western blot** and **immunoprecipitation** analyses. These experiments revealed that the lnRNP particles contain as integral components all spliceosomal U snRNPs (U1, U2, and U4/U6.U5 snRNPs) (32, 33), hnRNP proteins, and all known essential protein splicing factors, including the SR proteins, U2AF, and PTB (33-36) (see Table 1).

**Table 1. Components of Large Nuclear RNP (lnRNP) Particles<sup>a</sup>**

| snRNAs                                     | Probes                              | Reference |
|--|-------------------------------------|-----------|
| U1   | cDNA, anti-sense RNA                | (32)      |
| U2   | cDNA, anti-sense RNA                | (32)      |
| U4   | cDNA, anti-sense RNA                | (33)      |
| U5   | cDNA, anti-sense RNA                | (33)      |
| U6   | cDNA, anti-sense RNA                | (32)      |
| Proteins                                   | Antibodies                          | Reference |
| hnRNP core proteins A1, A2, B1, B2, C1, C2 | SLE autoantibodies, MAb iD2         | (10)      |
| snRNP proteins                             | SLE autoantibodies, anti-200 S MAbs | (32, 31)  |
| 56 kD antigen                              | Myositis autoantibodies             | (64)      |
| SR proteins                                | MAb104                              | (36)      |
| U2AF                                       | Anti-U2AF <sup>65</sup>             | (36)      |
| PTB  | Anti-PTB                            | (36)      |
| 783  | Anti-200 S MAb 783                  |           |

|       |                       |      |
|-------|-----------------------|------|
| 88 kD | Anti-200 S MAb 53 / 4 | (34) |
| 35 kD | Anti-200 S MAb 15 / 7 | (31) |
| 32 kD | Anti-200 S MAb 84 / 3 | (31) |
| 45 kD | Anti-200 S MAb 36     | (31) |
| 45 kD | Anti-200 S MAb 85     | (31) |

---

<sup>a</sup> Cells: Syrian hamster; HeLa.

#### 4.1. Stability of InRNP Particles

Mild RNase digestion of InRNP particles resulted in their conversion to 30 S particles, apparently without going through intermediate stages that give rise to an RNA ladder (10, 26). A similar observation had been made in 1979 by Wahrman and Augenlicht, who reported that limited nuclease digestion of hnRNP in nuclei revealed some discrete bands of RNA that appeared transiently, early in the course of digestion. But no multiples of a monomeric size RNP were evident (25).

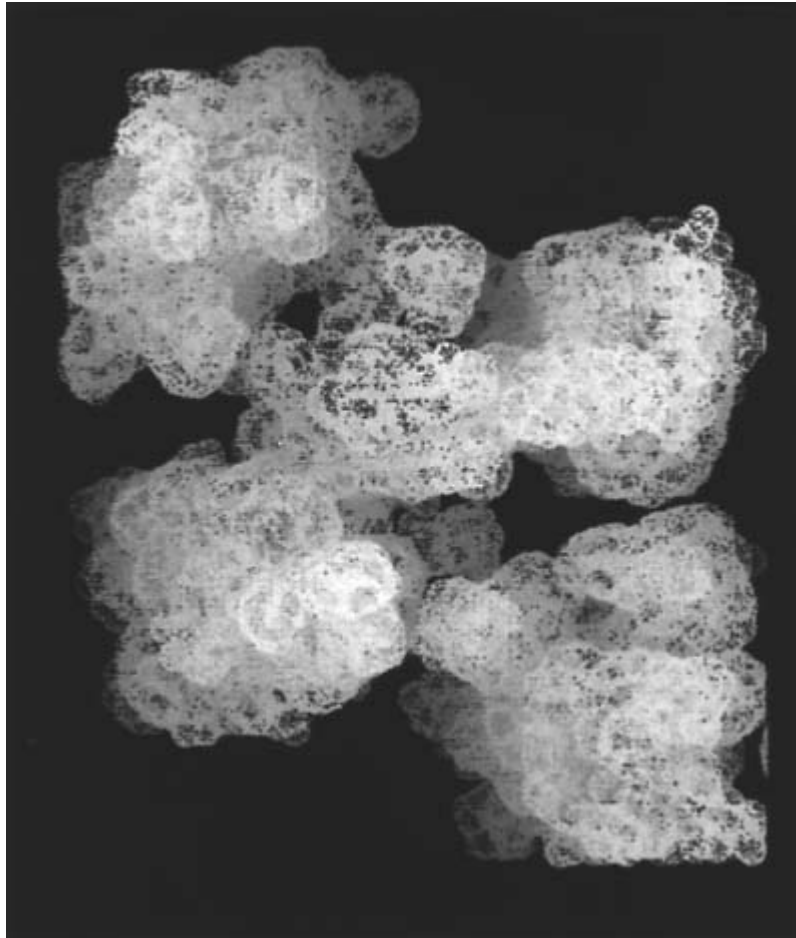
The 200 S InRNP particles require  $Mg^{2+}$  cations for their integrity, as they are more resistant to RNase digestion when incubated in the presence of magnesium cations than when incubated in the presence of EDTA. The addition of EDTA results in conversion of the 200 S InRNP particles into RNP particles that sediment at the 70 S region of the sucrose gradient. This shift is not due to pre-mRNA degradation but is caused by dissociation of some InRNP components (eg, U1 snRNP), leaving the pre-mRNA packaged with only a subset of the components of the 200 S InRNP. Importantly, the dissociation into 70 S particles is reversible, as indistinguishable 200 S InRNP particles are obtained after back addition of  $Mg^{2+}$  cations to the EDTA-dissociated particles (33). The importance of  $Mg^{2+}$  ions to the preservation of the InRNP particles is in accord with the known effect of  $Mg^{2+}$  ions on RNA folding (37) and with the proposed role of magnesium cations in RNA splicing (38). The sensitivity of InRNP particles to environmental conditions is consistent with the known labile nature of RNP particles, and it may explain some of the difficulties in earlier attempts to isolate intact RNP particles.

#### 4.2. Structure of InRNP Particles

Three-dimensional image reconstruction of isolated InRNP particles by automated electron tomography was performed at an optimal resolution of 2 nm (see [Single Particle Reconstruction](#)). The reconstructions revealed a quadrangular compact structure whose dimensions are 50×50×35nm (Fig. 2). The model is composed of four major subunits of similar dimensions that are connected to each other. An additional domain is sometimes observed toward the center of the particle. The four major subunits of the reconstructed models have an elongated globular shape pointing toward the center of the particle, where each subunit is connected to its two neighboring subunits at an internal radius close to the center of the particle (39, 40). Although an assessment of the interactions involved in these contacts cannot yet be made, it is plausible that the RNA transcript (mRNA) is participating in these inter-subunit interactions, although [protein-protein interactions](#) also can contribute. The automatic electron tomography studies were reinforced by mass measurements of the InRNP particles using [scanning transmission electron microscopy](#) (STEM). These studies confirmed the similarity of the individual InRNP particles and by measuring the mass of hundreds of particles revealed that the mass of the tetrameric InRNP particle is 21±1.6mDa. Furthermore, the apparent similarity in dimensions of the four major subunits was also confirmed by the STEM mass measurements, showing that the InRNP particles are composed primarily of four similar subunits, each having a mass of 4.8±0.5mDa (41). Interestingly, a rough estimate of the mass of the 60 S spliceosome, whose composition is similar to that of the InRNP particle, gives an almost identical value of 4.9 mDa. These observations have raised the possibility that the monomeric subunit of the native InRNP particle is equivalent to the in vitro spliceosome. The model that emerges from these studies is of an InRNP particle that is composed mainly of four similar subunits, presumably 60 S

spliceosomes. Thus, the lnRNP particle can be viewed as a supra-spliceosomal complex.

**Figure 2.** A three-dimensional model of an lnRNP particle. The model was reconstructed by automated electron tomography from a tilt-series of 71 images (39).



## 5. RNP Particles in Other Systems

The [lampbrush chromosome](#) loops, seen especially in amphibian oocytes, are a very suitable system to study the assembly of pre-mRNP particles. Studies on lampbrush chromosome loops of frogs and salamanders using Miller spreads (42) revealed the extended DNA template that forms the axis of the lateral loops that are transcribed simultaneously by several RNA polymerase II molecules. The nascent transcripts that are being packaged during transcription in RNP particles make the strands vary from thin to thick, depending on the length of the transcript (18, 43). Most lampbrush chromosome loops were shown to contain several hnRNP proteins, U snRNAs, snRNP proteins, and non-snRNP protein splicing factors, including the SR proteins family. Further, for all factors the signal extended throughout the length of the transcription unit. The conclusion of this study was that the nascent transcripts on lampbrush chromosome loops are assembled into a unitary hnRNP/snRNP particle that is maintained throughout the transcription process (44). These observations were also supported by in situ studies of perichromatin fibrils, which are proposed to represent extranucleolar RNP (45). These findings are consistent with studies on the composition of the lnRNP particles from mammalian cells, in which U snRNPs, non-snRNP and hnRNP proteins were found to be associated with the pre-mRNA in the lnRNP particles.

Another system that attracted attention is the **puffs** of [polytene chromosomes](#). These interphase

chromosomes replicate many times without separation of the replication products and exhibit a clearly distinguishable structure with a series of bands and interbands. The puffs in these chromosomes depict the extreme loosening of a band that accompanies DNA transcription activation in the band. One example of such studies are the puffs of the giant polytene chromosomes in the dipteran *Chironomus tentans* salivary glands (11, 18). Three large puffs, designated the [Balbiani ring](#) (BR), are sites of active transcription. The 35- to 40-kbp BR genes encode secretory polypeptides that are major components of salivary gland cells. In early electron microscopy studies of in situ sections, the very large, nascent BR pre-mRNA was observed to be assembled into 50-nm RNP particles. Furthermore, BR particles were seen as RNP globules of about 50 nm diameter in the nuclear sap and as elongated particles in the **nuclear pores** during translocation from the nucleus to the cytoplasm (46). The translocation of the BR RNP through the nuclear pore was further studied by electron microscopy methods, including electron tomography (47).

Interesting information concerning the assembly of BR RNP can be obtained from in situ studies, as the BR gene is transcribed simultaneously by several RNA polymerases. Thus, the proximal part of the gene represents initial stages of RNP assembly on short transcripts, whereas the distal part represents RNP particles close to maturation. In the proximal part of the active gene, the RNP fiber increases gradually in length and later coils into a thicker fiber. In the distal region, no further increase in the length of the RNP fiber is observed. Instead, the growing RNP fibers are packed into globular particles that increase only slightly in size along the gene and reach a constant diameter of 50 nm (48). The BR RNP particles were reported to contain hnRNP-like proteins and the Sm and snRNP proteins (47). Electron microscopy tomographic studies of individual BR RNP particles at 8-nm resolution revealed compact particles of 50 nm in diameter in which a thick RNP ribbon is bent into a horseshoe-shaped structure (49, 50). The use of osmium amine-B, a nucleic acid-specific stain, revealed particulate substructures (51, 52), which may be related to the subunit structure observed for the InRNP particles.

## 6. Coupling of Transcription to the Assembly of pre-mRNP

Numerous studies have established that nascent pre-mRNA transcripts are assembled during transcription with hnRNP proteins, U snRNP particles, and non-snRNP proteins, to form pre-mRNP complexes. The cotranscriptional assembly model is further supported by recent studies that demonstrate that several splicing factors are associated with the C-terminal domain repeat of RNA polymerase II large subunit (53, 54). Furthermore, the 5'-end capping enzymes, which are responsible for the capping activity of the 5'-end of pre-mRNAs, associate with the C-terminal domain, but only with its phosphorylated form (55, 56). This observation fits nicely with findings that 5'-capping occurs at early stages of transcription, ie, transcripts about 20 nucleotides long are already capped (57). Also, association of the 5'-capping enzymes only with the phosphorylated C-terminal domain, which is the stage when the polymerase becomes processive, ensures that only transcription units destined to be fully transcribed will be capped at their 5'-end. A yet unresolved question is why transcripts of RNA polymerase II assemble into pre-mRNP particles, whereas transcripts of RNA polymerase I and III do not, even though RNA-binding domains are highly abundant in most pre-mRNP proteins. Because this distinction must occur early in transcription, when pre-mRNP assembly initiates, it may be attributable to the association of the 5'-capping enzymes with the C-terminal domain. This association is unique to RNA polymerase II transcripts, as the other RNA polymerases do not have a C-terminal domain. Thus, capping the transcripts of RNA polymerase II not only protects them from degradation but also “marks” them as pre-mRNAs to be assembled as pre-mRNP complexes. Recent observations that splicing factors and factors involved in 3'-end processing also associate with the C-terminal domain (58-60) support this view.

## 7. Naturally Assembled pre-mRNPs are Large

The InRNP particle is a very large molecular assembly that sediments at the 200 S region in sucrose gradients and has a diameter of 50 nm (10, 27, 29). An interesting note is that similarly large RNP

particles have been observed in other systems using different methodologies. For example, 150–300 S complexes were observed by Wassarman and Steitz (61) on centrifugation of [HeLa Cells](#) nuclear extracts active in splicing. These complexes contained all the spliceosomal U snRNAs. Remarkably, the U2-U6 snRNA base-pairing interaction, which is typical of active spliceosomes, was found only within the large particles. These complexes probably represent particles assembled on endogenous HeLa transcripts. Another example are the BR RNP particles that were studied both in situ and in isolation and were shown to have a compact structure of about 50 nm in diameter (46, 49, 51). Finally, electron microscopy of spread preparations of chromatin, prepared under conditions that more closely resemble the physiologic ones than the standard conditions used for Miller spreads, revealed compact structures of about 60 nm in diameter (19). It therefore seems likely that the very large particle of 50 nm in diameter represents the native pre-mRNP particles in a broad range of living organisms.

## 8. A Model for the Native pre-mRNA Processing Machinery

Among the nuclear RNP complexes isolated thus far, the InRNP particle is unique in the sense that it contains intact pre-mRNA and all known components of the splicing machinery, and it can be viewed as a supraspliceosome complex. A speculative model can thus be formulated, assuming that the supraspliceosome is composed of four spliceosomes, each being capable of processing a single intron. In that case, it is quite simple to imagine how a pre-mRNA consisting of four introns is packaged in such a structure. The RNA would be displayed on the surface of the supraspliceosome in such a way that each intron would be looped around one subunit, thereby bringing the 5' and the 3' splice sites into juxtaposition. Pre-mRNAs with three or less introns would be assembled on the supraspliceosome occupying the respective number of subunits and leaving the remaining ones unoccupied. Notably, even RNA molecules that do not contain introns, such as the intronless nuclear histone H4 mRNA, are packaged in 200 S InRNP particles (10). The packaging of intronless pre-mRNAs in a supraspliceosome should not be surprising, because the assembly of RNP complexes in living cells initiates on an exon sequence as soon as the nascent pre-mRNA is synthesized (18). This is also suggested by studies of histone transcription units of lampbrush chromosomes of frog oocytes, in which spliceosomal U snRNPs, hnRNP proteins, and non-snRNP splicing factors were observed to be assembled on histone transcripts (44). The packaging of pre-mRNAs having more than four introns requires that the pre-mRNA be processed in steps. The prediction here would be that each step constitutes splicing of four introns at a time. Once a splicing event occurred, the ligated exons would translocate to make room for the association of another pair of exons that are part of the next set of four exon pairs to be spliced.

The concept of a supraspliceosome complies with the necessity to splice simultaneously more than one intron in a multi-intron message and fits well with current notions that the multiple spliceosomes that form on multi-intron pre-mRNAs must communicate to identify splice sites correctly and ensure the ligation of exons in the correct order (62). This issue becomes even more important in the case of regulated splicing that requires communication between alternative 5' and 3' splice sites that do not belong to the same intron. In this case, the supraspliceosome might act as a single mold on which alternative splice sites can be checked and approved for splicing. Finally, the organization of the splicing machinery in a multifunctional structure can also account for the fact that the introns of a multi-intron pre-mRNA are not removed by a linear progression of the splicing machinery from one end of the transcript to the other (eg, the APRT pre-mRNA, in which the third intron is removed first (63)), as would have been required for a processive mono-functional particle. The supraspliceosome, in contrast, allows the removal of introns in a nonlinear fashion while maintaining processivity throughout the entire splicing event. At the present stage, the InRNP model may help us understand the structure–function relationship of the naturally assembled spliceosomes, and it should especially help us design experiments to establish the structure and function of these RNP particles.

## Bibliography

1. O. P. Samarina, A. A. Krichevskaya and G. P. Georgiev (1966) *Nature* **210**, 1319–1322.



2. O. P. Samarina and A. A. Krichevskaya (1981) in *The Cell Nucleus* (H. Busch, ed.), Academic Press, New York, pp. 1–48.
3. G. Dreyfuss (1986) *Ann. Rev. Cell Biol.* **2**, 459–498.
4. J. A. Steitz, D. L. Black, V. Gerke, K. A. Parker, A. Krämer, D. Frendewey and W. Keller (1988) in *Structure and Function of Major and Minor Small Nuclear Ribonucleoprotein Particles* (M. L. Birnstiel, ed.), Springer-Verlag, Heidelberg, pp. 115–154.
5. H. D. Madhani and C. Guthrie (1994) *Ann. Rev. Genet.* **28**, 1–26.
6. R. Lührmann, B. Kastner and M. Bach (1990) *Biochim. Biophys. Acta* **1087**, 265–292.
7. A. Krämer (1995) in *Pre-mRNA Processing* (A. I. Lamond, ed.), R. G. Landes Company, Austin, Texas, pp. 35–64.
8. A. Krämer (1996) *Ann. Rev. Biochem.* **65**, 367–409.
9. J. M. Moore, C. C. Query and P. A. Sharp (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 303–358.
10. R. Sperling and J. Sperling (1990) in *The Eukaryotic Nucleus, Molecular Biochemistry and Macromolecular Assemblies* (P. R. Strauss and S. H. Wilson, eds.), Telford Press, Caldwell, New Jersey, pp. 453–476.
11. W. Beermann and G. F. Bahr (1954) *Exp. Cell Res.* **6**, 195–201.
12. J. G. Gall (1956) *J. Biophys. Biochem. Cytol. Suppl.* **2**, 393–396.
13. H. Swift (1959) *Brookhaven Symp. Biol.* **12**, 134–151.
14. T. Pederson (1974) *J. Mol. Biol.* **83**, 163–183.
15. A. L. Beyer, M. E. Christensen, B. W. Walker and W. M. LeSturgeon (1977) *Cell* **11**, 127–138.
16. T. E. Martin, P. Billings, J. Pullman, B. Stevens and A. Kinniburgh (1977) *Cold Spring Harb. Symp. Quant. Biol.* **42**, 899–909.
17. M. Jacob, G. Devilliers, J. P. Fuchs, H. Gallinaro, R. Gattoni, C. Judes and J. Stevenin (1981) in *The Cell Nucleus* (H. Busch, ed.), Academic Press, New York, pp. 194–259.
18. J. Sommerville (1981) in *The Cell Nucleus* (H. Busch, ed.), Academic Press, New York, pp. 1–55.
19. A. L. Beyer and Y. N. Osheim (1990) in *The Eukaryotic Nucleus, Molecular Biochemistry and Macromolecular Assemblies* (P. R. Strauss and S. H. Wilson, eds.), Telford Press, Caldwell, New Jersey, pp. 431–451.
20. S. Piñol-Roma, Y. D. Choi, M. J. Matunis and G. Dreyfuss (1988) *Genes Dev.* **2**, 215–227.
21. G. Dreyfuss, M. J. Matunis, S. Piñol-Roma and C. G. Burd (1993) *Annu. Rev. Biochem.* **62**, 289–321.
22. C. G. Burd and G. Dreyfuss (1994) *Science* **265**, 615–621.
23. G. Dreyfuss, M. S. Swanson and S. Piñol-Roma (1990) in *The Eukaryotic Nucleus, Molecular Biochemistry and Macromolecular Assemblies* (P. R. Strauss and S. H. Wilson, eds.), Telford Press, Caldwell, New Jersey, pp. 501–517.
24. R. D. Kornberg (1977) *Annu. Rev. Biochem.* **46**, 931–954.
25. M. Z. Wahrman and L. H. Augenlicht (1979) *Biochem. Biophys. Res. Commun.* **87**, 395–402.
26. R. Sperling, J. Sperling, A. D. Levine, P. Spann, G. R. Stark and R. D. Kornberg (1985) *Mol. Cell Biol.* **5**, 569–575.
27. P. Spann, M. Feinerman, J. Sperling and R. Sperling (1989) *Proc. Natl. Acad. Sci. USA* **86**, 466–470.
28. E. Miriami, J. Sperling and R. Sperling (1994) *Nucleic Acids Res.* **22**, 3084–3091.
29. R. A. Padgett, G. M. Wahl and G. R. Stark (1982) *Mol. Cell Biol.* **2**, 293–301.
30. J. Leavitt, P. Gunning, P. Porreca, S. Y. Ng, C. S. Lin and L. Kedes (1984) *Mol. Cell Biol.* **4**,

1961–1969.

31. D. Offen, P. Spann, R. Sperling and J. Sperling (1987) *Mol. Biol. Rep.* **12**, 183–184.
32. R. Sperling, P. Spann, D. Offen and J. Sperling (1986) *Proc. Natl. Acad. Sci. USA* **83**, 6721–6725.
33. E. Miriami, M. Angenitzki, R. Sperling and J. Sperling (1995) *J. Mol. Biol.* **246**, 254–263.
34. G. Ast, D. Goldblatt, D. Offen, J. Sperling and R. Sperling (1991) *EMBO J.* **10**, 425–432.
35. R. Sperling and J. Sperling (1998) In *RNP Particles, Splicing and Autoimmune Diseases* (J. Schenkel, ed.), Springer, Heidelberg, pp. 29–47.
36. S. Yitzhaki, E. Miriami, J. Sperling and R. Sperling (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8830–8835.
37. T. Pan, D. M. Long and O. C. Uhlenbeck (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 271–302.
38. T. A. Steitz and J. A. Steitz (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6498–6502.
39. R. Sperling, A. J. Koster, C. Melamed-Bessudo, A. Rubinstein, M. Angenitzki, Z. Berkovitch-Yellin and J. Sperling (1997) *J. Mol. Biol.* **267**, 570–583.
40. O. Medalia, A. J. Koster, A. Tocilj, M. Angenitzki, J. Sperling, Y. Z. Berkovitch and R. Sperling (1997) *J. Struct. Biol.* **120**, 228–236.
41. S. Müller, B. Wolpensinger, M. Angenitzki, A. Engel, J. Sperling and R. Sperling (1998) *J. Mol. Biol.* **283**, 383–394.
42. O. L. J. Miller and B. A. Hamkalo (1972) *Annu. Rev. Cytol.* **33**, 1–25.
43. J. G. Gall (1992) in *Advances in Developmental Biochemistry* (P. M. Wassarman, ed.), JAI Press, Greenwich, Connecticut, pp. 1–29.
44. Z. Wu, C. Murphy, H. G. Callan and J. G. Gall (1991) *J. Cell Biol.* **113**, 465–483.
45. S. Fakan (1994) *Trends Cell Biol.* **4**, 86–90.
46. B. J. Stevens and H. Swift (1966) *J. Cell Biol.* **31**, 55–77.
47. B. Daneholt (1997) *Eur. J. Cell Biol.* **74**, 407–416.
48. K. Andersson, B. Björkroth and B. Daneholt (1980) *Exp. Cell Res.* **130**, 313–326.
49. U. Skoglund, K. Andersson, B. Strandberg and B. Daneholt (1986) *Nature* **319**, 560–564.
50. A. Lönnroth, K. Alexciev, H. Mehlin, T. Wurtz, U. Skoglund and B. Daneholt (1992) *Exp. Cell Res.* **199**, 292–296.
51. A. L. Olins, D. E. Olins, H. A. Levy, M. B. Shah and D. P. Bazett-Jones (1993) *Chromosoma* **102**, 137–144.
52. A. L. Olins, D. E. Olins, V. Olman, H. A. Levy and D. P. Bazett-Jones (1994) *Chromosoma* **103**, 302–310.
53. E. J. Steinmetz (1997) *Cell* **89**, 491–494.
54. K. M. Neugebauer and M. B. Roth (1997) *Genes Dev.* **11**, 3279–3285.
55. S. McCracken, N. Fong, E. Rosonina, K. Yankulov, G. Brothers, D. Siderovski, A. Hessel, S. Foster, S. Shuman and D. L. Bentley (1997) *Genes Dev.* **11**, 3306–3318.
56. E. J. Cho, T. Takagi, C. R. Moore and S. Buratowski (1997) *Genes Dev.* **11**, 3319–3326.
57. M. Salditt-Georgieff, M. Harpold, S. Chen-Kiang and J. Darnell (1980) *Cell* **19**, 69–78.
58. M. J. Mortillaro, B. J. Blencowe, X. Wei, H. Nakayasu, L. Du, S. L. Warren, P. A. Sharp and R. Berezney (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8253–8257.
59. A. Yuryev, M. Patturajan, Y. Litingtung, R. V. Joshi, C. Gentile, M. Gebara and J. L. Corden (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6975–6980.
60. S. McCracken, N. Fong, K. Yankulov, S. Ballantyne, G. Pan, J. Greenblatt, S. D. Patterson, M. Wickens and D. L. Bentley (1997) *Nature* **385**, 357–361.

61. D. A. Wassarman and J. A. Steitz (1993) Proc. Natl. Acad. Sci. USA **90**, 7139–7143.
62. B. L. Robberson, G. J. Cote and S. M. Berget (1990) Mol. Cell. Biol. **10**, 84–94.
63. O. Kessler, Y. Jiang and L. A. Chasin (1993) Mol. Cell. Biol. **13**, 6211–6222.
64. H. Arad-Dann, D. A. Isenberg, Y. Shoenfeld, D. Offen, J. Sperling and R. Sperling (1987) J. Immunol. **138**, 2463–2468.

## Ribonucleotide Reductases

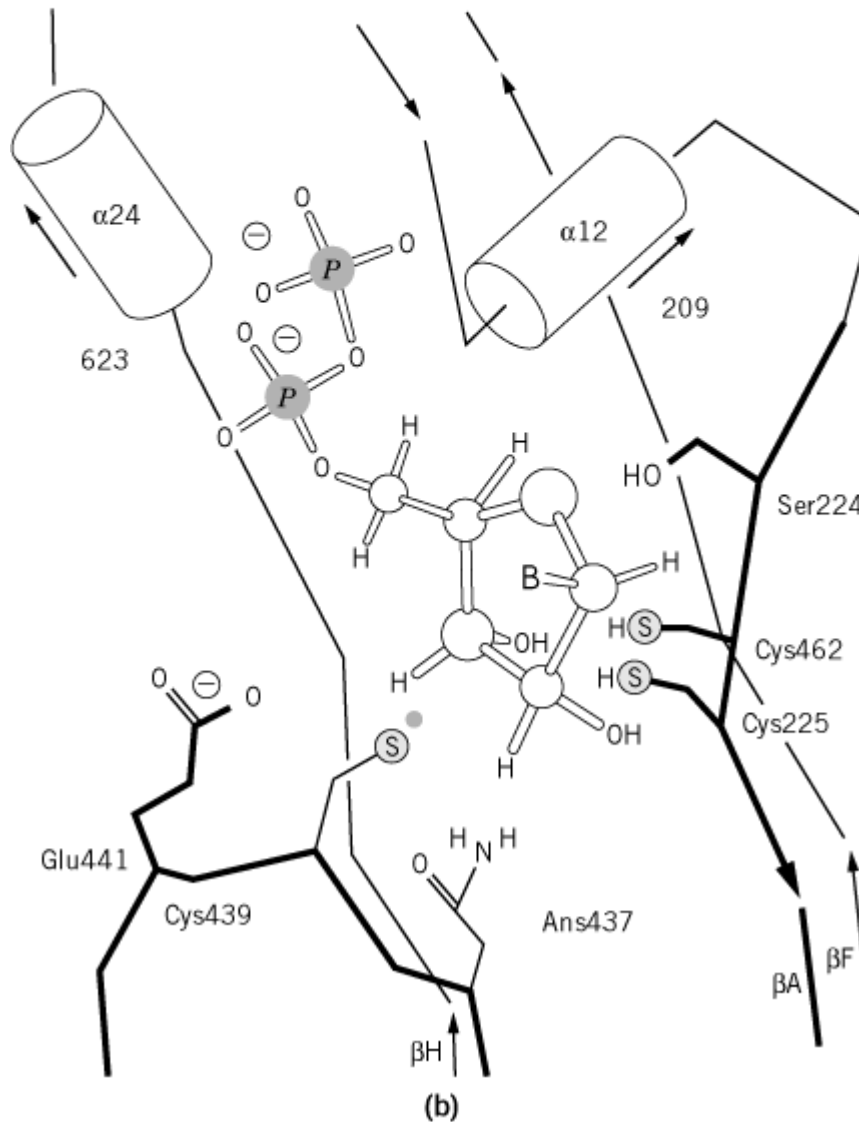
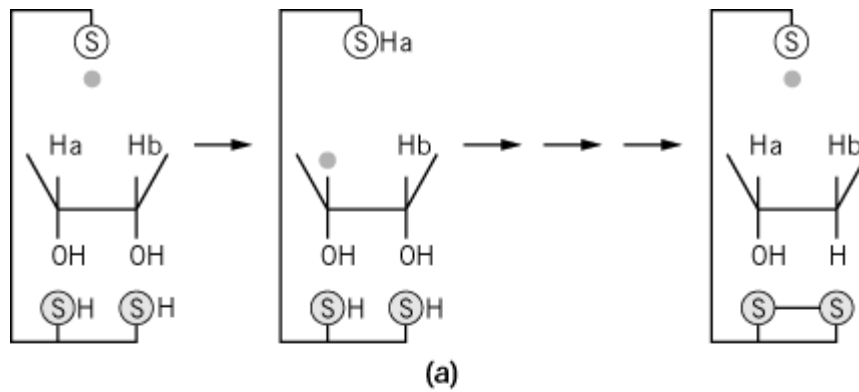
In all living cells, a ribonucleotide reductase (RNR) provides the deoxyribonucleoside triphosphates (dNTPs) for [DNA replication](#) and **repair** (1). A single [enzyme](#) replaces the OH group at C2' of the ribose moiety of the four common ribonucleotides with hydrogen. The reaction occurs with either ribonucleoside di- or triphosphates as substrates. Different organisms employ widely different forms of three different classes of the enzyme and, in some cases, contain more than one enzyme. *Escherichia coli* has three RNRs, two aerobic class I enzymes, and one anaerobic class III. All classes of RNR share two exceptional features:

1. The polypeptide chain of the active enzyme harbors a free radical amino acid residue that participates in the catalytic process; the mechanism for radical generation sets the classes apart
2. The specificity toward the four ribonucleotides is tightly controlled by **allosteric** effects that are remarkably similar for the three classes.

### 1. Reaction Mechanism

All ribonucleotide reduction proceeds via controlled free-radical-based chemistry, in which a free radical amino acid residue of an RNR generates a substrate radical by abstracting a hydrogen atom from C3' of the substrate, to facilitate the leaving of the OH group on the vicinal C2' (2) (Fig. 1a). A thiol group of a cysteine residue performs this function (3) (Fig. 1). Two additional redox-active cysteine residues then provide the reducing equivalents for the subsequent reduction at C2'. This general mechanism has strong experimental support for class I and II enzymes. The three-dimensional structure of the catalytic site of the *E. coli* class I enzyme beautifully fits this mechanism (4) (Fig. 1b). For class III RNR, the evidence for a similar mechanism is indirect.

**Figure 1.** (a) Free radical mechanism for the reduction of a ribonucleotide. A protein radical (thiyl for class I and II) abstracts the hydrogen from C3' of the ribotide, thereby transferring the radical function to the substrate. Subsequently, C2' of the activated substrate is reduced by redox-active cysteine residues, and the abstracted hydrogen is restored to C3', with regeneration of the protein radical. (After Ref. 2, simplified.) (b) Three-dimensional structure of the large protein of *E. coli* with a substrate modeled into the catalytic site. The suggested thiyl of Cys439 is located on the b-face of the ribose ring in close vicinity to C3', whereas the redox-active Cys462 and Cys225 are strategically located on the opposite a face, ready for reduction of C2'. (Reproduced from Ref. 1 with permission of TIBS. Courtesy of U. Uhlin.)



## 2. Three Classes of RNR

### 2.1. Structural Aspects

RNR can be grouped into three separate classes on the basis of structural and functional differences (5) (Table 1). Each class has a distinct quaternary structure: Class I and III are  $a_2b_2$  heterotetramers (for class I,  $a_2$  is called R1 and  $b_2$  is called R2), while class II has an  $a$  (or  $a_2$ ) homomeric structure. Within each class, the amino acid sequences are related, whereas **homologies** were originally not

found between the classes. However, the sequences of some class II RNR from primitive bacteria were recently found to contain stretches of homology also with class I and III enzymes.

**Table 1. Properties of RNRs Belonging to Different Classes**

|              | Class Ia  | Class Ib  | Class II                 | Class III                     |
|--------------|-----------|-----------|--------------------------|-------------------------------|
| Respiration  | Aerobic   | Aerobic   | Aerob/anaerobic          | Anaerobic                     |
| Structure    | $a_2b_2$  | $a_2b_2$  | $a(a_2)$                 | $a_2b_2$                      |
| Metal center | Fe–O–Fe   | Fe–O–Fe   | –Co–                     | 4Fe–4 S                       |
| Radical      | Tyr. .cys | Tyr. .cys | AdB <sub>12</sub> . .cys | AdoMet. .gly<br>AdoMet. .cys? |
| Reductant    | Redoxins  | Redoxins  | Redoxins                 | Formate                       |

The large a-polypeptide of all three classes has a molecular mass of between 80 and 100 kDa, with binding sites for substrates and for all allosteric effectors (see text below). It is the business end of RNR, where the reaction outlined in Figure 1 takes place. The small b-polypeptide, close to 40 kDa for class I and 20 kDa for class III, functions in radical generation.

The tertiary structures of both the large  $a_2$  (4) and the small  $b_2$  proteins (6) of the *E. coli* class I RNR and of the  $b_2$  protein of the mouse class I RNR (7) have been determined to high resolution. A plausible structure was suggested for the  $a_2b_2$  holoenzymes. The b-polypeptide contains a m-oxygen-linked diferric cluster with a stable tyrosyl radical in close vicinity. The exceptional stability of this radical (weeks in solution at cold room temperature) is explained by its extreme **hydrophobic** environment. The tyrosyl harnesses the free radical reactivity of the enzyme between turnover. During catalysis, when substrate is bound to the a-subunit, a coupled electron and proton transfer process between the tyrosyl radical of b and Cys439 of a (*E. coli* numbering) generates the thiyl radical required for abstraction of hydrogen from C3' of the substrate. This transfer occurs over a long distance (in *E. coli* estimated to be 25 to 30 Å), probably via a specific pathway involving bonded amino acid residues from both a and b, several of which have been identified from [site-directed mutagenesis](#) and considerations of the three-dimensional structure (8, 9).

Class II RNRs contain no stable free radical. Adenosylcobalamin replaces the function of b and forms a transient free radical during catalysis by homolytic cleavage of the carbon–cobalt bond. This radical, in turn, generates the thiyl radical on the cysteine residue involved in substrate activation (3).

Class III enzymes are constructed in a still different manner. Their  $b_2$ -protein harbors a 4Fe–4 S cluster bridging the two polypeptide chains (10) (see [Ferredoxins](#)), but no free radical amino acid residue. Instead, a glycy radical is located at the carboxy-terminal end of the a protein. Under anaerobic conditions the radical is very stable, but it is lost on admission of oxygen. The  $b_2$ -protein generates this radical (see text below).

Class I can be divided further into two subgroups (Ia and Ib) for which the two aerobic RNR of *E. coli* are the prototypes (Table 1). Class Ib RNR lack the first 50 amino-terminal residues, which in Ia RNR form the first two a-helices (11). One of the two allosteric sites of class I RNR is located in this

region, and it is lacking in class Ib enzymes.

## 2.2. Occurrence in Nature

Class I RNRs occur only in aerobic organisms, because the generation of their tyrosyl radical requires oxygen. **Eukaryotes** contain subclass Ia, whereas **microorganisms** usually contain subclass Ib (1). *Enterobacteriaceae* and some related bacteria are an exception and use Ia as the physiologically active enzyme. Several *Enterobacteriaceae* in addition contain a potentially active Ib form whose physiological function is unclear, because Ib does not complement conditional mutants of Ia unless several copies of Ib DNA are introduced into the cell. The *E. coli* Ia enzyme was the first discovered RNR, and extensive studies of its function have spearheaded our understanding of ribonucleotide reduction.

Members of class II occur among almost all groups of bacteria (1) (Table 2). They function both with and without oxygen and are common in both aerobic and anaerobic bacteria. The latter include several types of **archaeobacteria**. Nature's choice between class II and I is not obvious. Closely related bacteria may have either class. Class II has not been found in eukaryotes, with *Euglena gracilis* being an exception.

**Table 2. Distribution of Various Classes of RNRs in Nature**

|                           | Class |    |     |     |
|---------------------------|-------|----|-----|-----|
|                           | Ia    | Ib | II  | III |
| Eukaryotes                | +     |    | (+) |     |
| Eubacteria                |       |    |     |     |
| Thermotoga                |       |    | +   |     |
| Green nonsulfur bacteria  |       |    | +   |     |
| Deinococcus and relatives |       |    | +   |     |
| Gram-positive (high GC)   |       |    | +   | +   |
| Gram-positive (low GC)    | +     | +  | +   | +   |
| Purple bacteria           |       |    |     |     |
| a                         |       |    | +   |     |
| b                         |       |    | +   |     |
| g                         | +     | +  | +   | +   |
| Cyanobacteria             |       |    | +   |     |
| Green sulfur bacteria     |       |    | +   |     |
| Archebacteria             |       |    |     |     |
| Methanogens               |       |    |     | +   |
| Halophils                 |       |    | +   |     |
| Extreme thermopiles       |       |    | +   |     |

Class III enzymes (12) were discovered only recently, and their full potential is probably not yet appreciated. So far, all our knowledge rests on studies of the *E. coli* and **bacteriophage T4** enzymes,

and caution should be exercised in extrapolation to other members of this class. The existence of other active class III RNR has been deduced from DNA sequences found in various facultative aerobes and strict anaerobes (Table 2). Interestingly, *Methanobacteria* use class III and not class II enzymes, in spite of the abundance of cobalamines in these microorganisms. The protein sequences deduced from genes are all homologous to that of the anaerobic *E. coli* RNR, including a sequence motif around the radical-carrying glycyl residue at the carboxy-terminus.

### 3. Generation of the Stable Free Radical

For class I enzymes, the tyrosyl radical can be generated *in vitro* in the iron-free form of the  $b_2$  protein (the apoprotein) by treatment with ferrous iron and oxygen (13). Ferrous iron is bound in the appropriate site and oxidized to ferric iron, with generation of the tyrosyl radical. Intermediate steps probably involve the initial formation of a diferric peroxide, followed by several unstable intermediates, whose detailed structure is under discussion (13).

In cells, iron exists mainly as ferric iron and must be reduced before it can generate the radical. Reduction is catalyzed by a ubiquitous flavin reductase, ferric reductase (14). This enzyme also reduces a diferric center inside a form of  $b_2$  (met-protein) that has lost the tyrosyl radical (eg, by treatment with the radical scavenger hydroxyurea) but maintains its diferric center. In the presence of oxygen, flavin reductase then regenerates the radical and restores activity to the enzyme.

The generation of the glycyl radical of class III enzymes requires **S-adenosylmethionine** (AdoMet). An  $a_2b_2$  holoenzyme that lacks the glycyl radical binds AdoMet to its  $b_2$  moiety, but only after its  $4Fe-4S$  center has been reduced (15). Reduction can be made either chemically or by an enzyme system consisting of NADPH, flavodoxin, and flavodoxin reductase. Bound AdoMet is cleaved reductively to methionine + 5'-deoxyadenosine, driven by the oxidation of the iron-sulfur cluster. Simultaneously, the glycyl radical is formed on a  $a_2$ . The 5'-deoxyadenosyl radical is a probable intermediate in the reaction, as in the adenosylcobalamin-driven radical generation of class II enzymes. Two other enzymes, pyruvate formate lyase (16) and lysine 2,3-amino mutase (17), use AdoMet in a similar fashion for radical generation.

Similar to the stable tyrosyl radical of class I enzymes, the glycyl radical of class III may not activate the substrate directly but generate a second, transient protein-derived radical for abstraction of the hydrogen on C3'.

### 4. Electron Donors

Figure 1a shows that after one catalytic turnover the two redox-active cysteine residues of the enzyme have formed a **disulfide bond** that requires reduction before the next cycle. For class I and II enzymes, NADPH ultimately provides the required electrons. Several different small redoxins (**thioredoxins** and **glutaredoxins**) act as intermediates between NADPH and the disulfide on the enzyme. The reduced forms of these small proteins first interact with and reduce two cysteine residues at the carboxy-terminus (cysteine 754 and 759 of the *E. coli* class Ia RNR), which, in turn, by intramolecular transthiolation, reduce the disulfide at the active site (18).

Formate (19), rather than NADPH, is the reductant with class III RNR. Tritium from [ $^3H$ ] COOH is not incorporated into the product deoxyribotide but is instead incorporated into water, demonstrating formation of an intermediate carrying exchangeable protons. This tentatively suggests that redox-active protein-bound cysteine residues also function with class III RNR.

### 5. Inhibitors of RNR

The requirement of RNR for DNA synthesis gives inhibitors of the enzyme a potential medical

interest. Hydroxyurea (20), a scavenger of the tyrosyl radical, has been used clinically for some time. The compound has also been employed to distinguish between class I and II RNR and to study the function of the tyrosyl radical of class I enzymes. Experiments with substrate analogs that carry substitutions at C2' and that act as mechanism-based inhibitors were of paramount importance for the understanding of the radical chemistry catalyzed by RNR (2). Some such inhibitors have found their way into the clinic. Other potentially very interesting drugs are peptidomimetics that prevent binding of  $a_2$  to  $b_2$ . Such compounds can specifically inhibit the **herpes simplex virus** RNR without affecting the mammalian enzyme (21).

## 6. Allosteric Regulation

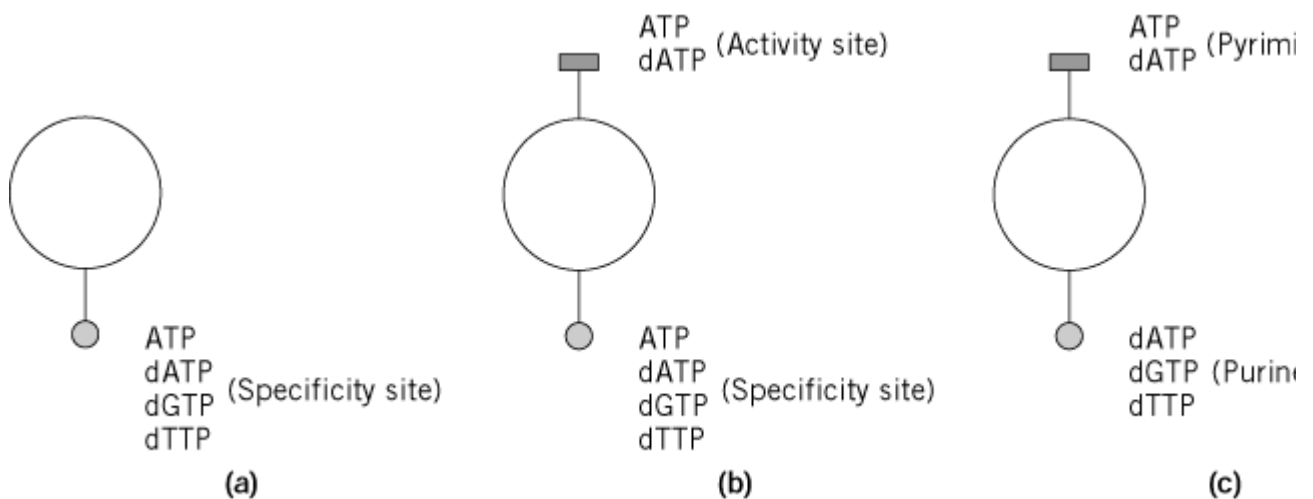
All three classes regulate their substrate specificity toward the four common ribonucleotides by binding nucleoside triphosphates to a specific allosteric site of the protein (the specificity site) (5). Class I and some class II RNRs distinguish between nucleoside *diphosphates* as substrates and nucleoside *triphosphates* as effectors. This distinction provides them with a more sophisticated regulation.

A specificity site binds either ATP, dATP, dTTP, or dGTP. Each binding event transforms the RNR into a form that reduces a specific substrate: with class I and II enzymes, ATP and dATP favor reduction of pyrimidine substrates, dTTP favors guanine nucleotides, and dGTP favors adenine nucleotides. The degree of specificity varies between different enzymes; it is very high for the mammalian enzyme. The effects are remarkably similar for the three classes in spite of the large differences in primary structure. The molecular mechanism is not known, but a direct interaction between effector and substrate nucleotides can be excluded, since the two binding sites are too far apart. Instead, it appears likely that subtle conformational changes induced by effector binding transform the substrate binding site and adapt it to a given substrate.

Class Ia RNRs contain a second allosteric site that binds only ATP and dATP (5). This site is named the activity site, because binding of ATP activates the enzyme, whereas dATP turns it off. The anaerobic *E. coli* class III enzyme is physiologically regulated in the same way as the Ia enzyme, but by a slightly different mechanism. It also has two separate allosteric sites, but both determine its substrate specificity. Because dATP binding to either site inhibits the enzyme, the final result is the same as for class Ia. Figure 2 shows models for the allosteric regulation of the various RNR.

**Figure 2.** Models for allosteric effector binding. (a) Large protein of a class Ib or class II RNR. A single site binds all effectors for the regulation of substrate specificity; (b) Large protein of a class Ia RNR. A specificity site binds the effectors that regulate substrate specificity of the RNR, and a separate activity site regulates enzyme activity by binding either ATP (activating) or dATP (inhibiting). (c) Large protein of a class III RNR. One site (the pyrimidine site) binds either ATP (activation of CTP and reduction) or dATP (inhibitory). The second site (the purine site) binds either dGTP (activating ATP reduction), dTTP (CTP reduction), or dATP (inhibitory).





## 7. Evolution(1)

The diversity of RNR may at first sight suggest that each of the three classes has evolved independently and that existing similarities are caused by **convergent** evolution. **Divergent** evolution is, however, favored by the similarity in allosteric effects and also by the ubiquitous and similar involvement of cysteine residues in the catalytic mechanism (1, 5). The properties of class III RNRs make them a favorite candidate for being the closest relative of a hypothetical “ur” reductase. The use of formate and of an iron–sulfur cluster fits a primitive enzyme. AdoMet can be considered a forerunner of adenosylcobalamin, the corresponding radical generator of class II RNR. Also, the fact that class III RNRs are strictly anaerobic enzymes speaks in their favor. Ribonucleotide reduction should have preceded [photosynthesis](#) during [evolution](#) and therefore was originally an anaerobic process. With the appearance of oxygen, other mechanisms had to come into play, leading to the emergence of class II and eventually class I RNR.

### Bibliography

1. P. Reichard (1997) *Trends Biochem. Sci.* **22**, 81–86.
2. J. Stubbe (1990) *Adv. Enzymol.* **63**, 349–417.
3. S. Licht, G. J. Gerfen, and J. Stubbe (1996) *Science* **271**, 477–481.
4. U. Uhlin and H. Eklund (1994) *Nature* **370**, 533–539.
5. P. Reichard (1993) *Science* **260**, 1773–1777.
6. P. Nordlund, B-M. Sjöberg, and H. Eklund (1990) *Nature* **345**, 593–598.
7. B. Kauppi et al. (1996) *J. Mol. Biol.* **262**, 706–720.
8. U. Rova et al. (1995) *Biochemistry* **34**, 4267–4275.
9. B-M. Sjöberg (1995) *Nucleic Acids Mol. Biol.* **9**, 192–221.
10. S. Ollagnier et al. (1996) *J. Biol. Chem.* **271**, 9410–9416.
11. A. Jordan et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 12892–12896.
12. P. Reichard (1993) *J. Biol. Chem.* **268**, 8383–8385.
13. A. Gräslund and M. Sahlin (1996) *Annu. Rev. Biophys. Biomol. Struct.* **25**, 259–286.
14. G. Spyrou et al. (1991) *J. Bacteriol.* **173**, 3673–3679.
15. S. Ollagnier et al. (1997) *J. Biol. Chem.* **272**, 24216–24223.
16. J. Knappe et al. (1993) *Biochem. Soc. Trans.* **21**, 731–734.
17. P. A. Frey (1993) *FASEB J.* **7**, 662–670.
18. A. Åberg et al. (1988) *J. Biol. Chem.* **264**, 12249–12252.

19. E. Mulliez et al. (1995) Proc. Natl. Acad. Sci. USA **92**, 8759–8762.
20. B. J. Kennedy (1992) Semin. Oncol. **19**, 21–26.
21. M. Liuzzi et al. (1994) Nature **372**, 695–698.

## Ribophorins

Ribophorins I and II are transmembrane **glycoproteins** characteristic of the [endoplasmic reticulum](#) (ER) membrane of all cells. They are restricted to the “rough” subcompartment of the ER, and they derive their name from their association with the [translation](#) machinery of **protein biosynthesis**. This association has been demonstrated in several experimental ways: (i) ribophorins are co-enriched with membrane-bound **polysomes** through multiple purification steps ([1](#)), with a molar ratio of approximately one of each ribophorin per [ribosome](#); (ii) ribophorins can be chemically **cross-linked** to ribosomes; (iii) their distribution in the plane of the ER membrane coincides with the ribosomes; and (iv) **antibodies** to the cytosolically exposed **epitopes** of ribophorins inhibit the co-translational translocation of a secretory protein into microsomes ([2](#)).

The tight association of the ribophorins with ribosomes initially led to the suggestion that they serve as ribosome receptors in the rough ER. It is now known, however, that the ribophorins are not involved in the translocation *per se*, but rather in the process of Asn-linked **N-glycosylation**. Ribophorin I and II were identified as two subunits of the mammalian oligosaccharyltransferase (OST) complex, which also contained a 48-kDa subunit. The transmembrane segment of ribophorin I is homologous to a recently proposed dolichol recognition **consensus sequence** ([3](#)), suggesting that ribophorin I could be involved in glycolipid binding and delivery.

The *Saccharomyces cerevisiae* OST complex consists of six nonidentical subunits. The a subunit is encoded by the Ost1 gene, which has significant sequence identity to mammalian ribophorin I. The d subunit is encoded by the Swp1 gene, which shares homology with the C-terminal half of mammalian ribophorin II. Expression of the Ost1 protein is essential for vegetative growth of yeast, and conditional Ost1 mutants have a pleiotropic effect of underglycosylation of soluble and membrane-bound glycoproteins ([4](#)). Microsomal membranes from Ost1 mutant yeast are markedly deficient in the *in vitro* transfer of high-mannose oligosaccharide from an exogenous lipid-linked donor to a glycosylation-site acceptor tripeptide, Asn–X–Thr/Ser ([4](#)).

Given these data, the previous experiments showing association of ribophorins with the translational machinery may indicate the existence of a supramolecular complex in the ER membrane, which performs docking of ribosomes, translocation of nascent chains across the membrane, Asn-linked glycosylation, and retrograde translocation of misfolded proteins out of the ER.

The human and rodent genes encoding ribophorins are single-copy genes. The human ribophorin I gene is located on chromosome 3q21, and the mouse gene is located on chromosome 6 ([5](#)). The human ribophorin II gene is on chromosome 20q12–q13.1 ([6](#)), while the murine gene is on chromosome 2, close to the *Src* proto-**oncogene** ([7](#)). Human ribophorin I and II are polypeptide chains of 68.5 and 69.3 kDa, respectively. [Primary structure](#) analysis, coupled with [proteinase](#) protection studies of intact microsomes, indicate that both ribophorins are synthesized as precursors having cleavable amino-terminal **signal sequences** of 22 to 23 amino acid residues ([8](#)). Ribophorin I and II are largely lumenally disposed: Their N-terminus is in the lumen, they span the membrane once with a 23-residue putative **transmembrane** domain, and they have 150- and 70-amino-acid-residue cytoplasmic domains, respectively.

Rat ribophorin I is glycosylated on one out of three potential Asn sites, and the carbohydrate remains in the high-mannose form (5), as expected from a resident ER protein. Detailed analysis of the glycosylation changes indicates that, unlike other ER proteins, the ribophorins always reside in the ER and do not cycle significantly to the Golgi apparatus (9, 10).

### Bibliography

1. G. Kreibich, M. Czako-Graham, R. Grebenau, W. Mok, E. Rodriguez-Boulan, and D. D. Sabatini (1978) *J. Supramol. Struct.* **8**, 279–302.
2. Y. H. Yu, D. D. Sabatini, and G. Kreibich (1990) *J. Cell Biol.* **111**, 1335–1342.
3. D. J. Kelleher, G. Kreibich, and R. Gilmore (1992) *Cell* **69**, 55–65.
4. S. Silberstein, P. G. Collins, D. J. Kelleher, P. J. Rapiejko, and R. Gilmore (1995) *J. Cell Biol.* **128**, 525–536.
5. A. Behal, K. Prakash, P. D'Eustachio, M. Adesnik, D. D. Sabatini, and G. Kreibich (1990) *J. Biol. Chem.* **265**, 8252–8258.
6. C. Loffler, V. V. Rao, and I. Hansmann (1991) *Hum. Genet.* **87**, 221–222.
7. G. Pirozzi, Z. M. Zhou, P. D'Eustachio, D. D. Sabatini, and G. Kreibich (1991) *Biochem. Biophys. Res. Commun.* **176**, 1482–1486.
8. C. Crimando, M. Hortsch, H. Gausepohl, D. I. Meyer (1987) *EMBO J.* **6**, 75–82.
9. M. G. Rosenfeld, E. E. Marcantonio, J. Hakimi, V. M. Ort, P. H. Atkinson, D. Sabatini, and G. Kreibich (1984) *J. Cell Biol.* **99**, 1076–1082.
10. N. E. Ivessa, C. De Lemos-Chiarandini, Y. S. Tsao, A. Takatsuki, M. Adesnik, D. D. Sabatini, and G. Kreibich (1992) *J. Cell Biol.* **117**, 949–958.

### Suggestions for Further Reading

11. G. Kreibich, M. Czako-Graham, R. Grebenau, W. Mok, E. Rodriguez-Boulan, and D. D. Sabatini (1978) Characterization of the ribosomal binding site in rat liver rough microsomes: ribophorins I and II, two integral membrane proteins related to ribosome binding. *J. Supramol. Struct.* **8**, 279–302.
12. S. Silberstein, P. G. Collins, D. J. Kelleher, P. J. Rapiejko, and R. Gilmore (1995) The alpha subunit of the *Saccharomyces cerevisiae* oligosaccharyltransferase complex is essential for vegetative growth of yeast and is homologous to mammalian ribophorin I. *J. Cell Biol.* **128**, 525–536.
13. C. Crimando, M. Hortsch, H. Gausepohl, and D. I. Meyer (1987) Human ribophorins I and II: the primary structure and membrane topology of two highly conserved rough endoplasmic reticulum-specific glycoproteins. *EMBO J.* **6**, 75–82.

### Ribosomes

Over thirty years ago, the central dogma of molecular biology stated the way genetic information flows: **DNA** is transcribed into [messenger RNA](#) (mRNA) that, in turn, is translated into [proteins](#) (1). The *ribosome*, discovered in the mid 1950s, is the universal cellular organelle facilitating the [translation](#) step, by catalyzing the sequential polymerization of amino acids according to the blueprint encoded in the mRNA. Although the ribosome catalyzes a rather simple chemical reaction, the formation of [peptide bonds](#), the process of **protein biosynthesis** is highly complicated and

sophisticated, and it depends on a large range of recognitions and interactions. The ribosome guarantees its fidelity and high efficiency by providing highly specific sites with various affinities for a large variety of incoming and outgoing molecules involved in protein synthesis. In **prokaryotes**, during exponential cell growth, the ribosomes may account for up to 50% of the dry cell mass and are distributed in the **cytoplasm**. **Eukaryotic** ribosomes are found also in [mitochondria](#) and **chloroplasts**. Ribosomes in all organisms are giant [ribonucleoprotein](#) (RNP) particles consisting of two subunits of unequal size, known as the large and small subunits. Two-thirds of the ribosomal mass is ribosomal RNA (rRNA), the rest is composed of 50 to 82 different ribosomal proteins (r-proteins), depending on the ribosomal source (Table 1). The names of the r-proteins are composed of L or S (depending on whether the protein is from the large or small subunit), and a running number, according to the position of this protein on **two-dimensional electrophoresis** gels.

**Table 1. Ribosomal Components**

| Ribosome Source                           | Avg. Sedimentation Coefficient(range) | rRNA Chains (Large Sub/SmallSub) | Approx. No. r-Proteins |
|---|---------------------------------------|----------------------------------|------------------------|
| Bacterial                                 | 70 S                                  | 5 S, 23 S/16 S                   | 50–60                  |
| Mitochondria (mammals)                    | 55 S                                  | 16 S/12 S                        | 80–90                  |
| Chloroplasts                              | 70 S                                  | 23 S, 5 S, 4.5 S/16 S            | 50–60                  |
| Mitochondria (fungi, protozoans, mammals) | 70 S (55 S–80 S)                      | 21–24 S/15–17 S                  | 65–90                  |
| Archaeobacteria                           | 70 S                                  | 5 S, 23 S/16 S                   | 65–75                  |
| Plant mitochondria                        | 75 S                                  | 5 S, 26 S/16 S                   | 70                     |
| Eukaryotes (cytoplasm)                    | 80 S                                  | 5 S, 5.8 S, 26–28 S/17–18 S      | 70–90                  |

Owing to the fundamental significance of ribosomes, they have been the target of numerous biochemical, biophysical, and genetic studies (see Appendix 1). These resulted in the elucidation of the gross structure of the ribosome, as well as the approximate locations of several functional sites. Due to recent significant technical advances, current ribosome research is characterized by major conceptual revisions resolving previous ambiguities and introducing substantial spatial rearrangements (2-5).

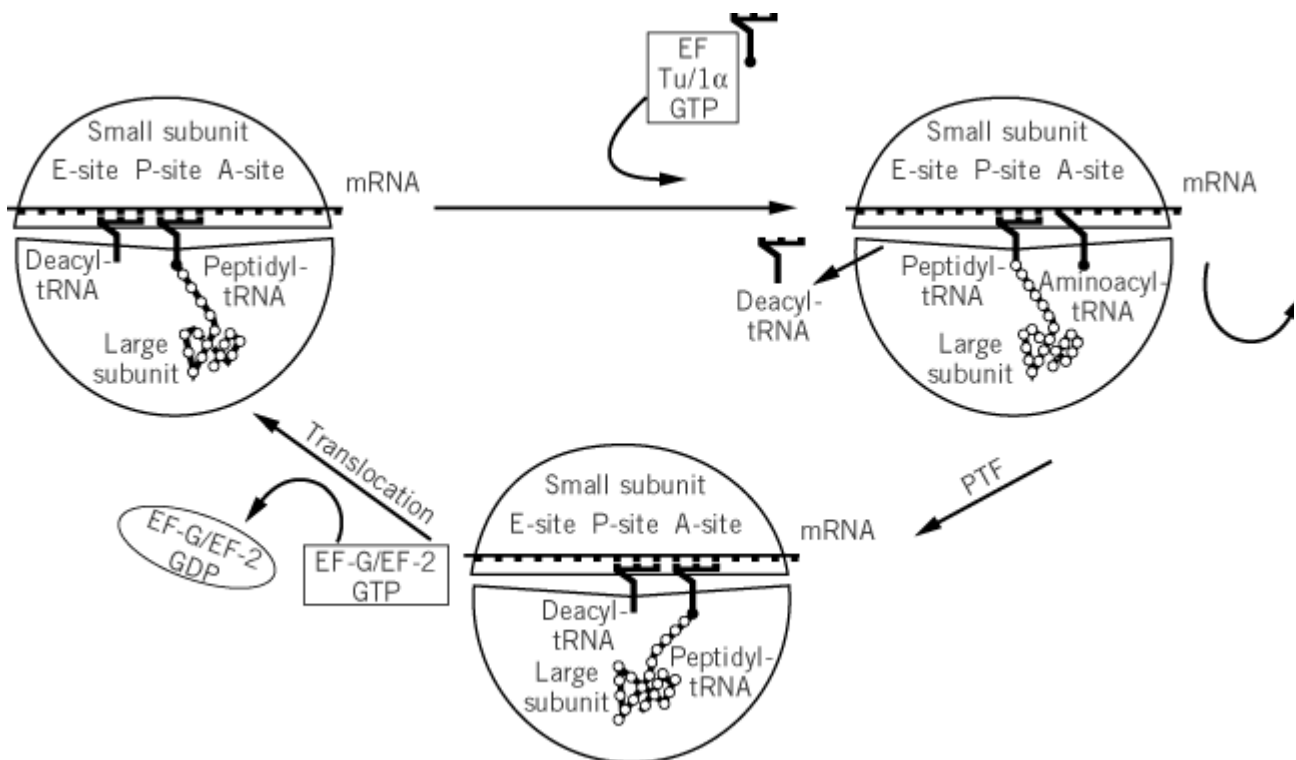
### 1. The Biosynthetic Process: An Overview

Protein biosynthesis can be divided into three functional steps: initiation, elongation, and termination. The initiation step requires the formation of the initiation complex, which is made of the small ribosomal subunit, [initiation factors](#) (IF), energy-rich compounds (GTP), initiator formylmethionine-**transfer RNA** (fMet-tRNA), and the mRNA molecule, with the [initiation codon](#) (usually AUG) in a favorable context. When the large subunit binds to the initiation complex, the elongation cycle can start.

The core of translation is the elongation cycle (reviewed in Ref. 6). In this step, one amino acid at a time is integrated into the growing nascent chain. The amino acids are brought to the ribosome in an activated state, bound to their corresponding tRNA through a high-energy phosphodiester bond. To ensure incorporation of the correct amino acid, the [anticodon](#) of the aminoacyl-tRNA must match the codon in the translated mRNA (see [Aminoacyl tRNA Synthetases](#)). In the early 1960s, Watson (7) and Lipmann (8) suggested a model for the elongation cycle with two tRNA binding sites on the ribosome: the A-site (acceptor site for aminoacyl-tRNA) and P-site (peptidyl-tRNA site). In the beginning of the 1980s, the existence of a third site, E (exit)-site, that binds only deacylated tRNA was proposed (9).

In the first round of the elongation cycle, only one tRNA is bound to the ribosome (the initiator tRNA, which is Met-tRNA in eukaryotes and fMet-tRNA in prokaryotes). During all other elongation rounds, two tRNA molecules are bound to the ribosome. Elongation starts with the f/Met-tRNA at the P-site. Subsequently, the decoding stage takes place, and an aminoacyl-tRNA carrying the amino acid coded by the next mRNA triplet is selected from the tRNA pool and delivered to the A-site by the [elongation factor](#) (EF) Tu in prokaryotes and EF-1-a in eukaryotes, which is a [GTP-binding protein](#), in the form of a ternary complex, aminoacyl-tRNA/EF/GTP. Once the right aminoacyl-tRNA is bound to the A-site, the EF leaves the ribosome as EF-GDP, and the ribosome carries out its intrinsic enzymatic task, the formation of the peptide bond. After peptide bond formation (the pre-translocational stage), the growing peptide is bound to the tRNA (as peptidyl-tRNA) at the A-site. At this point, a second elongation factor (EF-G in bacteria, EF-2 in eukaryotes, both **G proteins**) binds to the pre-translocational ribosome and catalyzes the translocation of the peptidyl-tRNA from the A-site to the P-site. At the same time, the deacylated tRNA from the P-site moves to the E-site (as illustrated in Fig. 1). In this way, the ribosome moves by one codon and reaches the post-translocational state (peptidyl-tRNA at the P-site, deacylated tRNA at the E-site). Now it is ready for the next round of elongation. Once the aminoacyl-tRNA corresponding to the next codon binds to the ribosome at the A-site, the deacylated tRNA bound at the E-site leaves the ribosome and the elongation cycle repeats itself, incorporating a new amino acid to the nascent peptide with each round of elongation. It has recently been suggested that during translocation the ribosome moves like a rigid frame along the mRNA between the three binding sites (A, P, and E), carrying two tRNA molecules so that the microtopography of the sites adjacent to the decoding region does not change (10).

**Figure 1.** Schematic diagram of the process of the elongation cycle.



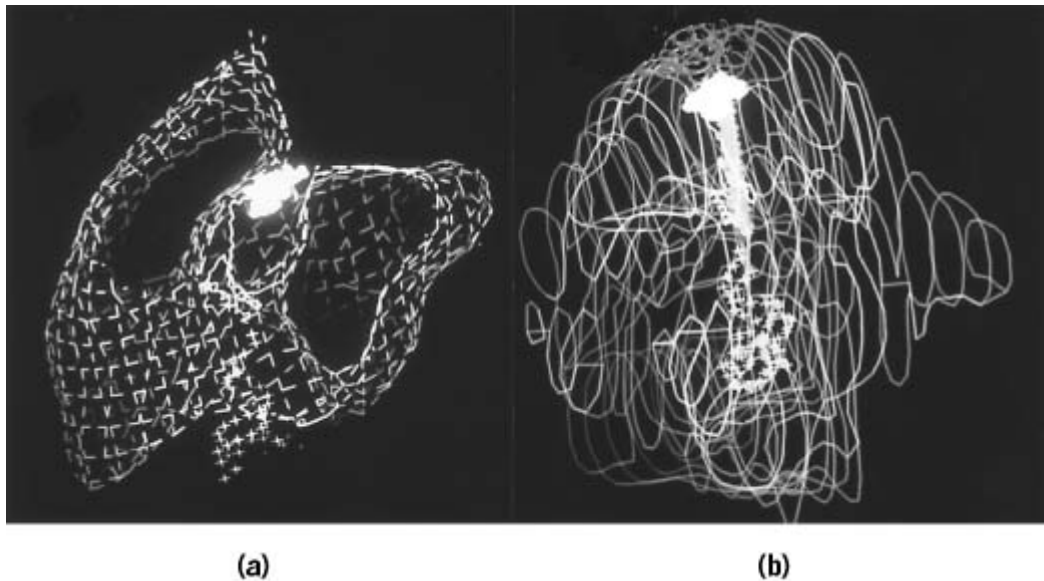
The elongation process ends when a [stop codon](#) is present at the A-site. The mechanism of the stop codon recognition by the termination factors and the subsequent release of the nascent protein from the ribosome is still not completely understood. In prokaryotes, it requires two [release factors](#) and one that stimulates them; in eukaryotes, a **homologous** family has been identified. It is known that the interaction of the release factors with the peptidyl transferase center catalyzes the addition of a [water](#) molecule instead of an amino acid to the peptidyl tRNA. This reaction frees the carboxyl end of the growing polypeptide chain from its attachment to the tRNA and promotes the release of the mature protein from the ribosome.

Additional compounds involved in this process are the [signal recognition particle](#) (SRP) in eukaryotes [or the corresponding complex in prokaryotes ([11](#))], which facilitates the targeted translocation of the nascent peptides; **chaperonins**, responsible for the correct folding of the nascent peptide; hydroxylase, acetylase, and aminopeptidase, a series of enzymes catalyzing amino acid modifications that can act cotranslationally; and regulators of ribosomal function under normal or stressful conditions, heat or starvation shock (eg, the *relA* gene product).

The various enzymatic activities associated with the process of protein biosynthesis take place in an internal ribosomal gap at the interface between the two subunits ([12](#), [13](#)). This assignment is based on extensive biochemical investigations that showed that the ribosome masks most of the components participating in the biosynthetic process: a stretch of about 30 nucleotides of the mRNA, the aminoacylated tRNA molecules ([14](#)), and a significant part of the nascent polypeptide chain ([15-17](#)). Spatial considerations allow the placement of three molecules of tRNA in this void, as well as the factors participating in the elongation cycle ([17-22](#)). Two of the tRNA molecules can be positioned so that their anticodons are close to the presumed rRNA-rich path of the bound mRNA, on the surface of the small subunit, and their CCA-termini pointing so that the growing peptide chain may extend into a tunnel spanning the large subunit. This tunnel originates at the subunit interface and terminates on the opposite side, at a location compatible with the exit site of the nascent polypeptide chain identified by [immunoelectron microscopy](#), and is thus suggested to be the path for the nascent proteins. Indeed, *in vitro* cotranslational folding was shown for the synthesis of full-length proteins (eg, rhodanese and ricin) while being bound to the ribosome ([23](#)). A feasible

description of these assignments is given in Figure 2.

**Figure 2.** Structural models of bacterial ribosomes. (a) A computer graphics display of the outer contour of the 70 S ribosome (13) (in blue). The image was reconstructed from negatively stained crystalline arrays of 70 S ribosomes from *B. stearothermophilus* at 47 Å resolution (17, 18). A model-built tRNA molecule was placed in the gap at the subunit interface so that its CCA 3' end points into the tunnel within the 50 S subunit. Following an *in vitro* cotranslational experiment (23), the main chain of the MS2 coat protein was placed along the tunnel in a partially unfolded conformation, maintaining the native **beta-strands** and the native (crystallographically determined) conformation of the segment of residues 1 to 47. The C-terminus was placed in the vicinity of the proposed peptidyl transferase center and the N-terminus at the exit domain of the tunnel. The region 1 to 47 includes all atoms and is shown as a space-filling structure. (b) A slice of 50 Å thickness of the 50 S subunit (within the 70 S shown on the left), as reconstructed from electron diffraction (11, 12) with the same components.



The synthesis of a complete protein of average size takes 20 to 60 seconds. Nevertheless, multiple initiations of translation would ensure a faster and efficient translation of mRNA molecules. Accordingly, *polysomes* (or **polyribosomes**) can be formed along a single molecule of mRNA from several ribosomes spaced as close as 25 codons apart.

## 2. Proposed Functional Relevance of Ribosomal Components

### 2.1. The Ribosomal Components

The natural tendency of the ribosomes to disintegrate led originally to the assumption that they are nonspecific aggregates. Later on, the contrary was established, and the composition of ribosomal particles from various sources was determined (24). By far the best biochemically characterized bacterial ribosome is that of *Escherichia coli*. It contains about quarter of a million atoms, and has a molecular weight of about 2.3 million daltons and **sediments** with a coefficient of 70 S. About two-thirds of the mass of the ribosome is composed of three chains of rRNA (with a total of about 4500 nucleotides). One copy of each of its 58 different proteins is present in the ribosome, with the exception of the tetrameric L12. The ribosomes of **archaeobacteria** and eukaryotes are somewhat larger (eg, eukaryotic ribosomes migrate with a sedimentation coefficient of 80 S), reflecting the higher complexity of the eukaryotic cellular environment.

### 2.2. Ribosomal RNA Chains

Although it has been suggested that the original, primitive ribosome may have been composed solely of rRNA (see [RNA World](#)), until recently it was assumed that the catalytic activities of the ribosome

are carried out mainly by the r-proteins and that the rRNA molecules have a more passive role in providing the scaffold for the ribosome and in binding the mRNA. Recently, it became clear that the ribosomal functions are no longer due solely to r-proteins, and the ribosomal RNA has been proven to play an active part in the ribosomal functions (25). The prominent catalytic activities of the rRNA are:

1. The [peptidyl transferase](#) activity (PTF) resides in the large subunit. In some organisms, it was found that a model assay of this activity is relatively resistant to proteolysis (26). The minimal set of components indispensable for peptidyl transferase activity was identified as a stretch of approximately 100 nucleotides, at the central loop of domain V of the 23 S RNA, plus a small fraction of the r-proteins (about 5 to 10% of the total mass). The aminoacyl end of the bound tRNA has been shown to make contacts with this rRNA loop, and mutations in this region abolish the PTF reaction completely or make the ribosomes resistant to antibiotics that target the PTF activity.
2. The [GTPase](#) center is associated with a highly conserved stretch of 23 S rRNA that binds the moderately conserved protein L11. Whereas the rRNA is essential for this activity, cells lacking protein L11 are viable.
3. The first step in protein biosynthesis, the formation of the preinitiation complex, depends in prokaryotes on the base-pairing interaction of the small ribosomal subunit with a region of the mRNA, called the [Shine–Dalgarno sequence](#), located 3 to 10 bases 5' of the initiation codon (usually AUG). Moreover, the rRNA from the small subunit provides the decoding site, where the mRNA and anticodon loop of the tRNA interact.

Crystallographic studies led to the determination of the structure of two rRNA domains, both part of the 5 S rRNA: a synthetic ribonucleic oligomer of 12 base pairs, imitating Helix A; and stretches of 29 nucleotides containing the sarcin/ricin loop as well as 62 nuclease-resistance nucleotides (helices I and IV, loop E) and a part of it, a dodecamer containing the minimum 11 base pairs required for binding the protein L25. The last, composed of loop E, exhibits an irregular geometry of cross-strand purine stacks.

### 2.3. The r-Proteins

Because of the complexity of the ribosome and its functions, it is still not possible to assign a function to each ribosomal protein. It is assumed that the r-proteins play an important role in the proper folding of the rRNA and enabling it to function efficiently. As mentioned above, some of the r-proteins are essential for the enzymatic properties, such as the PTF (eg, L2), whereas others are involved in the binding of tRNA (S7 and S8). In addition, there is evidence for the bifunctionality of a number of r-proteins (in prokaryotes), which function as regulators of translation by their ability to bind to the polycistronic mRNA coding for them. In some cases, they bind to similar structural motifs in the rRNA and mRNA, suggesting that these proteins interact in a similar fashion with mRNA and rRNA (eg, the interactions of L1 and S8, which possess a stable structure even in isolation); in others, the structural motifs of the rRNA and mRNA target sites do not resemble each other (eg, S15, which exhibits significant flexibility in isolation).

Complete sequences of the r-proteins from *E. coli* and many other bacteria and eukaryotes have been determined (24, 26). Some ribosomal proteins undergo [post-translational modifications](#). Among others are acylation (eg, L17/L12 in *E. coli*) and phosphorylation (eg, the eukaryotic P1, P2, and P0). Several structural motifs have been suggested for the r-proteins. Among them are clusters of basic/acidic residues; amino acid sequence repeats in shared elements; **zinc finger** domains; basic regions; [leucine zipper](#) motifs; and carboxyl extensions of **ubiquitin**-like proteins.

For over two decades, the production of crystals of r-proteins useful for [X-ray crystallography](#) was extremely poor. This, together with the observation that some ribosomal proteins lose their *in situ* conformation upon isolation, led to the assumption that the conformations of almost all r-proteins are dictated by their *in situ* supporting environment. Furthermore, no correlation has been found between



the crystallizability of individual r-proteins and the degree of their evolutionary conservation, their localization within the ribosomes, or their involvement in primary contacts with rRNA. An appropriate example of a highly conserved r-protein that is intimately bound to r-RNA in a major functional center (GTPase activity), and undergoes significant conformational changes upon isolation from the ribosome, is protein L11. Interestingly, L11 regains its natural fold and can be reconstituted into core particles lacking it, even when a large chemical moiety with a molecular weight approaching a third of its own is bound to it (17).

The recent increasing sophistication in instrumentation, the implementation of powerful genetic techniques, and the use of ribosomes from **thermophilic** bacteria resulted in major progress in the structure determination of isolated r-proteins. Though it remains to be seen whether the structures of the isolated ribosomal components bear resemblance to their *in situ* conformation, it has been suggested that components possessing an intrinsic characteristic fold may crystallize, provided they are not damaged during their preparation. The structures that have been determined (fully or partly) by either X-ray crystallography or solution heteronuclear NMR are: S4, S5, S6, S7, S8, S17, L1, L6, L9, L7/I12 (the C terminal fragment), L14, L21, L22, L25, L30 (27, 28). Most share the split b-a-b fold, called the “common” motif and abbreviated as RRM (RNA recognition motif) (see [RNA-Binding Proteins](#)). A few proteins show different folds, called “unique” and “multiple” in which the interacting regions are built mainly of loops. One protein (S8) exhibits versatility in its RNA contact sites, since one interacts with rRNA and the other is involved in binding tRNA. Two novel RNA-binding domains have been recently detected. One can be aligned with **homeodomains** (DNA-binding proteins) consisting predominantly of **alpha-helices** connected by a turn and exhibiting structural flexibility, acquiring stability upon RNA binding (S15 and in the C-terminal of L11). The other is a b-ribbon arm, similar to that found in DNA-bending proteins, observed in S7. This protein acts as the main regulatory element for one of the r-protein operons and is crucial for tRNA binding and assembly of the small subunit.

#### 2.4. Assembly and Reconstitution

In prokaryotes, the assembly of ribosomes occurs in the cytoplasm and is coupled to the [transcription](#) of rRNA molecules. Thus, r-proteins bind to rRNA while being synthesized. In prokaryotes, the *in vivo* ribosome assembly requires between 2 and 3 min. In eukaryotes, the situation is different and more complicated. The r-proteins are synthesized in the cytoplasm and then imported into the [nucleolus](#) (a substructure of the [nucleus](#)), where ribosomal assembly takes place. Once the ribosomes are assembled, they are exported back into the cytoplasm. Thus, ribosome assembly and transport take between 30 min (for the small subunit) and 1 hour (for the larger one).

In prokaryotic ribosomes, both subunits can be separated into their components and then reconstituted *in vitro* to fully active articles, even after partial or total unfolding has occurred. Interestingly, the reconstitution process is performed under nonphysiological conditions and takes considerably longer than the *in vivo* assembly (90 min vs 3 min). The ability of the ribosomes to reconstitute *in vitro* shows that the information required to obtain the active [quaternary structure](#) of the ribosome resides within the ribosomal components. Originally, it was assumed that the r-proteins governed the assembly process and the rRNA chains undergo significant conformational changes throughout the assembly process. However, it has been suggested recently that rRNA may influence the conformations of some r-proteins (eg, S15, L11), as conformational changes in these r-proteins were induced by their interactions with rRNA (29).

Reconstitution experiments led to the construction of ribosome assembly maps, showing the sequential binding of the different r-proteins to the rRNA molecules during assembly. It was found that two ribosomal proteins initiate the assembly of each ribosomal subunit. These proteins are defined as structural inducers, since they bind directly to the respective RNA without cooperativity during the onset of assembly and are assumed to induce the creation of *in situ* microenvironments that serve as assembly nuclei within the ribosome. It has been shown that, in general, the *in vitro* assembly patterns imitate the *in vivo* formation of ribosomes, but recent studies indicate the involvement of nonribosomal cell products in the *in vivo* ribosome assembly. These include the

chaperonin DnaK (30) and [RNA Helicases](#) belonging to the **DEAD box** protein family.

## 2.5. Evolution vs Universality

Ribosomes from the three kingdoms (eubacteria, eukaryotes, and archaeobacteria) vary considerably in their size and the number of their components. However, they also exhibit a high degree of conservation with respect to their architecture and some r-proteins and rRNA regions (31, 32). As the ribosome is basically an RNA enzyme (a **ribozyme**), it can be regarded as ancient on the evolution scale. The existence of r-proteins conserved throughout all three **phylogenetic** kingdoms implies that they appeared in the first stages of evolution. With the progression from primitive to higher organisms, the number of the r-proteins increases, indicating their participation in the more complex functions of the eukaryotic ribosomes.

The ribosomes from archaeobacteria are thought to be vestiges of transition stages in the evolution from prokaryotes to eukaryotes. Thus, it has been frequently remarked that the amino acid sequences of archaeobacterial r-proteins are closer to their eukaryotic homologues than their bacterial counterparts, whereas the organization of their genes mimics that of eubacteria. Furthermore, in several *in situ* substructures, the level of conservation is so high that ribosomal proteins from one source can be fully exchanged by their homologues from other ribosomes, even when they belong to two different kingdoms or function under totally different conditions. It is not surprising that protein L11 from *E. coli* binds well to cores of *Bacillus stearothermophilus* lacking it because of the high homology of these ribosomes, but the universality of the internal substructure of L1 and of a segment of the 23 S rRNA is rather unpredictable. Despite the evolutionary distance between the eubacteria and archaeobacteria, chimeric complexes of L1 were reconstituted between halophilic components and their corresponding mates from *E. coli*.

## 2.6. Inhibitors of Ribosomal Function

Many antibiotics act directly on ribosomes, and most of them interact in one way or another with rRNA. Their binding sites were determined in different ways, including by rRNA mutations leading to drug resistance. The antibiotic functions of the PTF drugs involve interference with the reaction itself ([puromycin](#) and [chloramphenicol](#)) or with the movements required to perform the elongation cycle. **Tetracycline** is the classical antibiotic that inhibits the binding of the aminoacyl-tRNA. The aminoglycosides, [streptomycins](#) and the family of the gentamycins, [kanamycins](#), and [neomycins](#) block A-site occupation and stimulate misreading, resulting in the incorporation of the wrong amino acids and leading to the production of nonfunctional proteins. **Erythromycin** and lincosamide probably stimulate the dissociation of peptidyl-tRNA by blocking the entrance to the tunnel that conveys the exiting nascent peptide. Thiostrepton and spectinomycin prevent translocation, ie, they inhibit the conformational change between the pre- and post-translocational states.

Among the most potent inhibitors of protein synthesis are naturally occurring peptide [toxins](#) acting mainly as **nucleases** and known as ribosomal-inactivating proteins (RIP). This group includes a-sarcin, **ricin**, arbin, mitogillin, restrictocin, shiga toxin, and vero toxin. Most hydrolyze a single phosphodiester bond in the large subunit rRNA, whereas [colicin](#) E3 acts on the 16 S RNA (Table 1) and pokeweed antiviral protein (PAP) interacts with the EF binding site. a-Sarcin and ricin bind to a loop in the large subunit rRNA that includes a universally conserved dodecamer sequence, which appears even in ribosomes that are not sensitive to these toxins (eg, those from *Haloarcula marismortui*). This domain is involved in aminoacyl-tRNA and elongation factor binding, as well as GTPase activity.

## 3. Structural Information

### 3.1. Approximate Shapes and Positions

A large number of traditional, as well as specifically designed structural methods (see Appendix 1), have been employed for shedding light on the structural organization of the ribosome and elucidating its [quaternary structure](#). The relative positions of the centers of mass of the *E. coli* r-proteins have been determined, and the spatial *in situ* proximities between several ribosomal components could be

approximated (2-5). In parallel, attempts to determine accurately the secondary structure of rRNA, pinpoint tertiary structure elements, and assign to them functional relevance, resulted in proposals for fairly detailed topological models for the organization of the rRNA (32). Computational and **phylogenetic** analyses of possible tertiary interactions in the rRNA, combining results of chemical, physical, and functional experiments with energy minimization, were also carried out.

For over three decades, a wide variety of [electron microscopy](#) techniques have been the methods of choice for viewing ribosomes at various levels of detail. Since the late 1980s, along with the refinement of sophisticated and powerful microscopical and image reconstruction techniques, these studies enabled feasible assignments of functional relevance to a few structural features. The first reconstructed three-dimensional models were obtained from periodically ordered monolayers (ordered arrays) of eukaryotic and prokaryotic ribosomes and their large subunits that occurred naturally (12) or were grown *in vitro* (13). Despite their low resolution, these models revealed several key features, associated mainly with internal vacant spaces, cavities, gaps, tunnels, and partially filled hollows, that had not been detected earlier. Consequently, the ribosome, which was traditionally conceptualized as a compact network, was shown to be rather spongy.

### 3.2. Higher Resolution Studies

The accumulated knowledge about ribosomes has not yet revealed the molecular mechanism of protein biosynthesis. To extend the limits of our understanding of this process, an accurate and reliable model of the ribosome is required. Such a model should be obtained by X-ray crystallography. Two approaches have been taken. One focuses on isolated ribosomal components (r-proteins and rRNA described above); the second aims at the elucidation of the structure of entire ribosomal particles. Being ribonucleoprotein complexes that are notoriously flexible, unstable, and prepared routinely as conformationally mixed populations, ribosomes provide an extremely complicated system for crystallographic studies. On the other hand, the natural periodic organization of ordered helical or two-dimensional arrays of ribosomes has been observed in eukaryotic cells (eg, lizard, chicken, amoebae, and human) exposed to stressful conditions, such as suboptimal temperatures, wrong diet, or lack of oxygen, and it has been hypothesized that these periodically ordered forms are the physiological mechanism for temporary storage of ribosomes, aimed at preserving their integrity and activity for the expected better future. Thus, we have had the crystallization attempts aimed at extending the natural tendency to form periodic arrays into the growth of well-ordered three-dimensional crystals. Perhaps the most striking and unexpected achievement of the last decade is the growth of usable crystals of ribosomal particles from **halophilic** or thermophilic bacteria (17). As a strong correlation was found between the activity of the ribosomes and the quality of their crystals, it has been suggested that these ribosomes are more stable than those from eubacteria and retain their integrity and activity during the isolation and crystallization processes. Furthermore, far beyond the initial expectations, one of these crystal forms (of the 50 S subunit from *H. marismortui*) diffracts to almost atomic resolution, 2.7 Å (34, 35), the other (of the small subunit from *T. thermophilus*) to 3.0 Å (36). As initial phases could be derived at medium resolution (31), the way to structural analysis of the ribosome has been paved.

#### APPENDIX 1: Summary of methods used to study ribosomal structure and/or function

- Tests for ribosomal activity: RNA-binding; Poly-U-dependent Poly-Phe synthesis; natural mRNA *in vitro* translation; the use of synthetic tRNA and analogs
- Total ribosomal reconstitution (from isolated r-proteins and r-RNA) for ribosome assembly studies
- Selective removal of r-proteins under mild chemical conditions (salt or organic solvents) followed by partial *in vitro* reconstitution (ribosomal cores+split proteins)
- Heterologous ribosomal reconstitution (using rRNA with r-proteins from a different species and vice versa)
- Reconstitution in the absence of single components, to determine the effect of single r-protein omission in either assembly or ribosomal function
- Chemical [footprinting](#), hydroxyl radical cleavage, and protection studies of ribosome–ligand

- interactions (eg, tRNA, mRNA, and antibiotics)
- Chemical modification probing of selected ribosomal moieties
  - Specific cleavages by **nucleases**: a-sarcin, ricin, ribonuclease H, etc.
  - rRNA-r-protein, rRNA-rRNA, mRNA-rRNA, and subunit [crosslinking](#) by mild UV irradiation or chemical agents to map contact topography or functional sites
  - Photo-activated reagents in [affinity labeling](#) of ribosomal ligands or components
  - Fluorescence **energy transfer** to follow the binding of ribosomal ligands (eg, tRNA) or the dynamics of the growth of the nascent polypeptide and its folding
  - Single-site mutations (in rRNA) leading to antibiotic resistance or ribosome inactivation
  - Probing exposed single-stranded rRNA regions with complementary DNA oligonucleotides
  - Modeling of rRNA, based on [distance geometry](#), energy minimization, or phylogenetic conservation
  - Sequence determination and secondary-structure predictions of r-proteins and rRNA
  - Comparisons between species and determination of homologous conserved regions
  - [Light scattering](#) and **hydrodynamic** measurements for size/shape estimation
  - Triangulation: Reconstitution or binding of one or a few protonated components (r-proteins, tRNAs) into deuterated ribosomes. The locations of the centers of mass are approximated by **neutron scattering**, including proton-spin [contrast variation](#).
  - [Electron microscopy](#) (EM), dark field and tunneling, negative and positive staining
  - Electron microscopy coupled with three-dimensional image reconstruction using (1) optical diffraction of tilt series of ordered two-dimensional arrays (monolayers) or (2) single particles embedded in vitrified ice and viewed with normal- or low-dose electrons, followed by image processing and angular reconstitution of isolated particles
  - [Immunolectron microscopy](#): gross r-protein mapping through electron microscope localization of **antibodies** bound to the r-proteins that were the [immunogen](#)
  - **Neutron diffraction** coupled with contrast variation, investigating the gross localization of the ribosomal components, benefiting from the different scattering properties of protein and RNA
  - Heteronuclear three-dimensional [NMR](#) spectroscopy of isolated small r-proteins or fragments of the larger ones
  - X-ray crystallography of ribosomal particles, as well as crystallizable ribosomal components

## Bibliography

1. F. Crick (1970) *Nature* **227**, 561–563.
2. B. Hardesty and G. Kramer, eds. (1986) *Structure, Function and Genetics of Ribosomes*, Springer-Verlag, New York.
3. E. W. Hill, A. Dahlbert, R. A. Garrett, P. B. Moore, D. Schlessinger and J. R. Warner eds. (1990) *The Ribosomes: Structure, Function and Evolution*, American Society for Microbiology, Washington, DC.
4. K. Nierhaus F. Franceschi, A. P. Subraminian, V. A. Erdmann, and B. Wittmann-Liebold, eds. (1993) *The Translation Apparatus*, New York, Plenum Press.
5. A. T. Matheson, J. Davies, W. E. Hill, and P. P. Dennis, eds. (1995) *Frontiers in Translation*, J. Biochem. Cell Biol. **73**.
6. J. Czworkowski and P. B. Moore (1996) *Prog. Nucl. Acid Res. Mol. Biol.* **54**, 293–332.
7. J. D. Watson (1963) *Science* **140**, 17–26.
8. F. Lipmann (1963) *Prog. Nucl. Acid Res.* **1**, 135–262.
9. H. J. Rheinberger, H. Sternbach, and K. N. Nierhaus (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5310–5314.

10. K. N. Nierhaus (1996) *Nature* **379**, 491–492.
11. J. D. Miller, H. D. Bernstein, and P. Walter (1994) *Nature* **367**, 657–659.
12. R. A. Milligan and P. N. T. Unwin (1986) *Nature* **319**, 693–696.
13. A. Yonath, K. R. Leonard, and H. G. Wittmann (1987) *Science* **236**, 813–816.
14. H. F. Noller (1991) *Ann. Rev Biochem* **60**, 191–227.
15. L. I. Malkin and A. Rich (1967) *J. Mol. Biol.* **26**, 329–334.
16. G. Blobel and D. D. Sabatini (1970) *J. Cell. Biol.* **45**, 130–145.
17. Z. Berkovitch-Yellin, W. S. Bennett, and A. Yonath (1992) *CRC Crit. Rev. Biochem. Mol. Biol.* **27**, 403–444.
18. H. A. S. Hansen et al. (1990) *Biochim. Biophys. Acta* **1050**, 1–5.
19. H. Stark et al. (1995) *Structure*, **3**, 815–821.
20. H. Stark et al. (1997) *Cell* **88**, 19–28.
21. J. Frank et al. (1995) *Nature*, **376**, 3, 441–444.
22. J. Frank (1997) *Curr. Opin. Struct. Biol.* **7**, 266–272.
23. M. Eisenstein et al. (1994) in *Biophysical Methods in Molecular Biology* (G. Pifat, ed.), Balaban Press, Rehovot, Israel, p. 213.
24. H. G. Wittmann (1983) *Ann. Rev. Biochem.*, **52**, 35–47.
25. H. F. Noller, V. Hoffarth, and L. Zimniak (1992) *Science* **256**, 1416–1419.
26. I. G. Wool, Y. L. Chan, and A. Glück (1996) in *Translational Control*, Cold Spring Harbor Lab. Press, New York, pp. 685–732.
27. A. Liljas and M. Garber (1995) *Curr. Opin. Struct. Biol.* **5**, 721–727.
28. V. Ramakrishnan and S. W. White (1998) *TIBS* **23**, 208–212.
29. A. Yonath and F. Franceschi (1997) *Nature Struct. Biol.* **4**, 3–5.
30. J. H. Alix and M. F. Guerin (1993) *Proc. Natl. Acad. Sci. USA* **90**, 9725–9729.
31. P. Dube et al. (1998) *Structure* **6**, 389–399.
32. A. Yonath and F. Franceschi (1998) *Structure* **6**, 678–684.
33. R. Brimacombe (1995) *Eur. J. Biochem.* **230**, 365–383.
34. K. von Böhlen et al. (1991) *J. Mol. Biol.* **222**, 11–15.
35. A. Yonath et al. (1998) *ACTA Cryst.* **A54**, 945–955.
36. A. Harms et al. (1999) *Structure*, in press.

### **Suggestions for Further Reading**

37. B. Hardesty and G. Kramer, eds. (1986)
38. E. W. Hill, A. Dahlbert, R. A. Garrett, P. B. Moore, D. Schlesinger and J. R. Warner eds. (1990)
39. A. T. Matheson, J. Davies, W. E. Hill, and P. P. Dennis, eds. (1995)
40. K. Nierhaus F. Franceschi, A. P. Subraminia, V. A. Erdmann, and B. Wittmann-Liebold, eds. (1993)

### **Ribozyme/Catalytic RNA**

A ribozyme is defined as an **RNA** molecule capable of **catalyzing** a chemical reaction in the absence

of [proteins](#). Ribozymes can promote their reactions with rate enhancements as large as  $10^{11}$  and specificities of better than 1 site in 1000. Catalytic RNA can fold into a compact three-dimensional structure that defines and efficiently activates a single phosphate linkage for a chemical reaction. All of the known biological ribozymes catalyze transesterification reactions of phosphate diesters and participate in RNA maturation or processing. Although the majority of catalytic RNAs in biology perform autocatalytic functions, all of them can be converted into true enzymes capable of multiple turnover catalysis by altering the covalent connectivity of the RNA, so that the “substrate” can be added in *trans* to the active site. The discovery of catalytic RNA has led to the speculation that ribozymes might have contributed to self-replication and **protein biosynthesis** in early evolution, an idea termed the [RNA World](#) hypothesis (1).

Catalytic RNAs are found in a broad variety of biological situations. Four small catalytic RNAs were identified in **satellite RNA** from plants and animals. These include the hammerhead, the hairpin, the hepatitis delta virus (HDV), and the Varkud satellite (VS) ribozymes. All of these RNAs perform the same cleavage reaction, but they differ substantially in their **secondary** and tertiary structures. In each case, the satellite RNA is proposed to replicate via a complementary RNA intermediate by a **rolling circle** mechanism. The ribozymes participate in the replication by self-cleaving the tandem satellite repeats into monomer units.

A second and more common category of catalytic RNA comprises the self-splicing **introns** (2). These RNAs provide the active site necessary to complete their own RNA **splicing** from an RNA transcript, with concomitant ligation of the flanking exons. Self-splicing introns are divided into two classes (group I and group II) based on differences in their conserved structures and reaction mechanisms. In addition to splicing, both classes of introns can catalyze reverse splicing reactions, which allows them to act as [mobile elements](#) for the horizontal transfer of genetic information.

A third example of a ribozyme is [ribonuclease P](#) (RNase P), which participates in transfer RNA (tRNA) processing (3). It acts in *trans* to catalyze the removal of nucleotides from the 5'-end of the pre-tRNA. The cellular form of this enzyme consists of both an RNA and a protein component (M1 RNA and C5 protein in *Escherichia coli*). Both are biologically essential; however, the RNA subunit alone is sufficient to carry out the tRNA processing reaction, albeit at a slower rate than the holoenzyme.

## 1. Reaction Mechanisms

Each reaction mechanism of these ribozymes involves in-line attack by an activated nucleophile at a phosphodiester linkage. This results in inversion of the stereochemical [configuration](#) at the phosphate group. The primary difference among the various catalytic RNA reaction mechanisms is the identity of the nucleophile and the leaving group (4). All four of the small catalytic RNAs utilize an internal nucleophile, the 2'-OH located immediately adjacent to the scissile phosphate, to displace the 5'-O of the phosphodiester. This one-step reaction mechanism produces RNAs with 2'-3' cyclic phosphate and 5'-OH termini. Identical products result from alkaline hydrolysis of RNA, although the ribozyme reaction is specific to a single phosphodiester linkage.

The group I introns perform two consecutive transesterification reactions that employ external nucleophiles (5). In the first step of splicing, the 3'-OH of an exogenous guanosine **cofactor** attacks the phosphate linkage at the 5'-exon-intron boundary. This reaction covalently adds the guanosine onto the 5'-end of the intron and releases the 5'-exon. The second step of splicing is essentially the reverse of the first step but utilizes the 3'-OH of the guanosine at the intron-3'-exon boundary as the leaving group.

The group II introns employ a reaction mechanism that is a hybrid between these two classes (2). It utilizes an internal nucleophile, the 2'-OH of a bulged A nucleotide but the hydroxyl group is not

adjacent to the scissile phosphate. The phosphate is some distance away, at the 5'-exon-intron boundary. As with the group I intron, the 3'-OH of the 5'-exon is the leaving group for the first step of group II splicing and the nucleophile for the second step. The intron is released as a **lariat structure**, with the 5'-end of the intron covalently attached to the 2'-OH of the bulged A nucleotide. This is the same reaction mechanism employed in the more complex process of messenger RNA precursor splicing that is catalyzed by the [spliceosome](#).

Unlike the reaction mechanisms of the other naturally occurring catalytic RNAs that primarily use either the 2' or 3'-OH for chemical reactions, RNaseP utilizes water as its nucleophile (3). The hydrolysis reaction displaces the 5'-leader sequence to produce a tRNA with a terminal 5'-phosphate group. Because it both employs an external nucleophile and acts on an external substrate, RNaseP is the only natural RNA known to act as a true multiple turnover [enzyme](#) *in vivo*.

## 2. Metals in Catalysis

All of the naturally occurring catalytic RNAs (with the possible exception of the hairpin ribozyme) require divalent metal cations for activity (6, 7). Although the metal specificity varies substantially between ribozymes, it can usually be met by  $Mg^{2+}$  or  $Mn^{2+}$ . These metals play both a structural and a catalytic role. A dramatic example of divalent metals in RNA folding is the P5abc subdomain of the *Tetrahymena* group I intron. Three  $Mg^{2+}$  ions coordinate to disparate phosphates within the subdomain (8). This allows the RNA to fold inside out, ie, the phosphates point into the structure and the nucleotide bases point out to solvent. Some structural metals in RNA can often be substituted with polycations such as **spermidine** or cobalt hexamine, which emphasizes the importance of charge neutralization for RNA folding.

Other metals included in the ribozyme fold are required for catalysis. These metals activate the nucleophile or stabilize the leaving group within the transesterification reaction. A well-characterized example is the group I intron, where a different metal is required for each of these functions (9, 10). The metal-binding sites are highly selective and cannot be substituted with a generic polycation. A ribozyme [active site](#) with two metal ions is analogous to what occurs within the **DNA polymerase** enzyme and might be a general feature of catalytic RNAs (11).

## 3. Catalytic RNA Structure

To understand RNA catalysis it is necessary to understand the three-dimensional (3D) arrangement of the nucleotides, metal ions, water molecules, and co-factors that comprise the RNA structure. The very concept of a catalytic RNA with a substrate binding site requires that the structure be much more complex than the single-stranded description given in early textbooks. Several methods, including [chemical modification](#), **phylogenetic** comparison, [X-ray crystallography](#), and [NMR spectroscopy](#), have shown that RNA can adopt a complex globular structure with many deviations from the canonical **Watson-Crick base pairs** and the traditional **double helix**. Although a surprisingly small amount of high-resolution structural information is available on catalytic RNAs, a set of structural motifs is gradually being identified that includes several surprises. For example, RNA can form a pseudo base pair between two consecutive A nucleotides in a single-stranded region (termed an *A-platform*) that serves as a platform for tertiary helix-stacking interactions (12). There are also examples of **pseudoknots**, G · U [wobble pairing](#) receptors, GAAA **tetraloops**, and U-turns (13-16). The **consensus sequences** of these structural motifs are seen repeatedly in the [phylogenetic databases](#), which suggests that they have been used as building blocks to create a variety of RNA structures.

Although a large domain of the *Tetrahymena* group I intron has also been reported, the 3D structure of only one complete catalytic RNA is currently available. The structure of the hammerhead ribozyme was determined with an all-deoxy substrate and with a substrate containing a single 2'-O-methyl substitution in place of the 2'-OH nucleophile (14, 17). Both structures are very similar. A

third structure was determined with an all-ribose substrate in the absence of divalent metal ions (18). Metal soaking and freeze-trapping experiments have suggested that only small conformational changes are necessary for the ribozyme to move from the ground state to the chemical [transition state](#).

#### 4. Other potential ribozymes

In addition to the ribozymes described above, other essential cellular processes might also be catalyzed by RNA. The [ribosome](#) and the [spliceosome](#) have essential and highly conserved RNA components. It is a matter of intense interest and investigation whether the active sites of these molecular machines are composed of RNA or protein. For example, a ribosome treated extensively with [detergent](#) and [proteinase K](#), conditions that **denature** and degrade the protein components of the ribosome but leave the RNA intact, could still catalyze rudimentary [peptide bond](#) formation (19). Although this was highly suggestive of an RNA active site in the ribosome, at least three proteins remained intact after the [proteinase](#) treatment, and all activity was lost if these proteins were removed. It remains uncertain whether the RNA or the few remaining proteins comprise the active site, but the intriguing possibility remains that the ribosome might also be a ribozyme.

#### 5. New unnatural ribozymes

Several new RNA catalysts have been identified from large sequence pools by SELEX, a combinatorial amplification technique that exerts evolutionary pressure on an RNA population to select and optimize for a particular catalytic function. The range of reactions catalyzed by selected ribozymes in vitro is now well beyond that known for natural RNAs (20). Some of the notable activities identified by this approach include an RNA-based **RNA polymerase** that utilizes nucleotide triphosphates and releases pyrophosphate as the leaving group, a **polynucleotide kinase** that can transfer the  $\gamma$ -phosphate from ATP onto the 5' or internal 2'-OH of an RNA substrate, and RNAs that utilize  $\text{Ca}^{2+}$  or  $\text{Pb}^{2+}$  instead of  $\text{Mg}^{2+}$  for catalysis. RNAs have also been identified that catalyze chemical reactions other than transesterification at phosphate. An aminoacyl-RNA synthetase was identified that can transfer an aminoacyl group from Phe-AMP onto the 3'-end of an RNA substrate (see Aminoacyl tRNA synthetases). In vitro selection was also used to identify ribozymes that catalyze formation of carbon-nitrogen and carbon-sulfur bonds. In one example, in which a [transition state analogue](#) was used as the bait for selection, an RNA was isolated that could catalyze the isomerization of a bridged biphenyl linkage. A variety of catalysts such as these would be necessary in a prebiotic RNA world where RNA would act as both the information storage molecule and the catalyst for reactions. This Darwinian approach to in vitro evolution has also demonstrated that nucleic acid-based catalysts are not restricted to RNA. **DNA** catalysts have been identified that promote reactions with rate enhancements comparable to the small, naturally occurring ribozymes (21). Although the diversity of reactions remains rather limited, there is ample reason to believe that the chemical reactions catalyzed by RNA will continue to expand.

#### Bibliography

1. G. F. Joyce and L. E. Orgel (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 1–25.
2. T. R. Cech (1986) *Cell* **44**, 207–210.
3. S. Altman, L. Kirsebom, and S. Talbot (1993) *FASEB J.* **7**, 7–14.
4. T. R. Cech and B. L. Bass (1986) *Ann. Rev. Biochem.* **55**, 599–629.
5. T. R. Cech (1990) *Ann. Rev. Biochem.* **59**, 543–568.
6. T. Pan, D. M. Long, and O. C. Uhlenbeck (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 271–303.
7. A. M. Pyle (1993) *Science* **261**, 709–714.
8. J. H. Cate, R. L. Hanna, and J. A. Doudna (1997) *Nature Struct. Biol.* **4**, 553–558.



9. J. A. Piccirilli, J. S. Vyle, M. H. Caruthers, and T. R. Cech (1993) *Nature* **362**, 85–88.
10. L. B. Weinstein, B. C. N. M. Jones, R. Cosstick, and T. R. Cech (1997) *Nature* **388**, 805–808.
11. T. A. Steitz and J. A. Steitz (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 6498–6502.
12. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, A. A. Szewczak, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1696–1699.
13. H. A. Heus and A. Pardi (1991) *Science* **253**, 191–194.
14. H. W. Pley, K. M. Flaherty, and D. B. McKay (1994) *Nature* **372**, 68–74.
15. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1678–1685.
16. S. A. Strobel, L. Ortoleva-Connelly, S. P. Ryder, J. H. Cate, and E. Moncoeur (1998) *Nature Struct. Biol.*, in press.
17. W. G. Scott, J. T. Finch, and A. Klug (1995) *Cell* **81**, 991–1002.
18. W. G. Scott, J. B. Murray, J. R. P. Arnold, B. L. Stoddard, and A. Klug (1996) *Science* **274**, 2065–2069.
19. H. F. Noller, V. Hoffarth, and L. Zimniak (1992) *Science* **256**, 1416–1419.
20. T. Pan (1997) *Cur. Op. Chem. Biol.* **1**, 17–25.
21. R. R. Breaker (1997) *Cur. Op. Struc. Bio.* **1**, 26–31.

### Suggestions for Further Reading

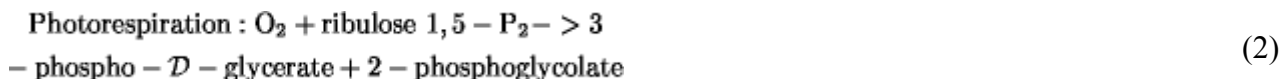
22. *The RNA World* (1993) (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Press, Cold Spring Harbor, N Y. Includes several excellent reviews of RNA biology and catalysis. Of particular interest is the article by T. R. Cech (“Structure and mechanism of the large catalytic RNAs: Group I and Group II introns and ribonuclease P,” pp. 239–269).

## Ribulose Bisphosphate Carboxylase/Oxygenase (Rubisco)

Ribulose bisphosphate carboxylase/oxygenase (Rubisco) is probably the best-studied plant [enzyme](#), partly because of its great abundance (constituting 30–50% total protein in leaves) but also because of its importance in the metabolism of all cells. By catalyzing carbon dioxide fixation during [photosynthesis](#), this enzyme is responsible for virtually all of the reduced carbon found in living organisms. Rubisco couples the inorganic and organic carbon pools on earth and is the most significant route for linking these pools together by the synthesis of carbohydrate. It is, however, a poor **catalyst**, having both a low affinity for carbon dioxide and a small [turnover number](#) ( $3s^{-1}$ ). **Autotrophic** organisms must devote a major part of their synthetic capacity to produce sufficient enzyme to sustain life. The high concentrations of Rubisco and its [messenger RNA](#) have provided many opportunities to study the molecular biology of plants. This includes gene organization and **gene expression** in both **nuclei** and [chloroplasts](#), **protein targeting** to chloroplasts, **molecular chaperones**, enzyme assembly, and many other aspects of plant cell biology.

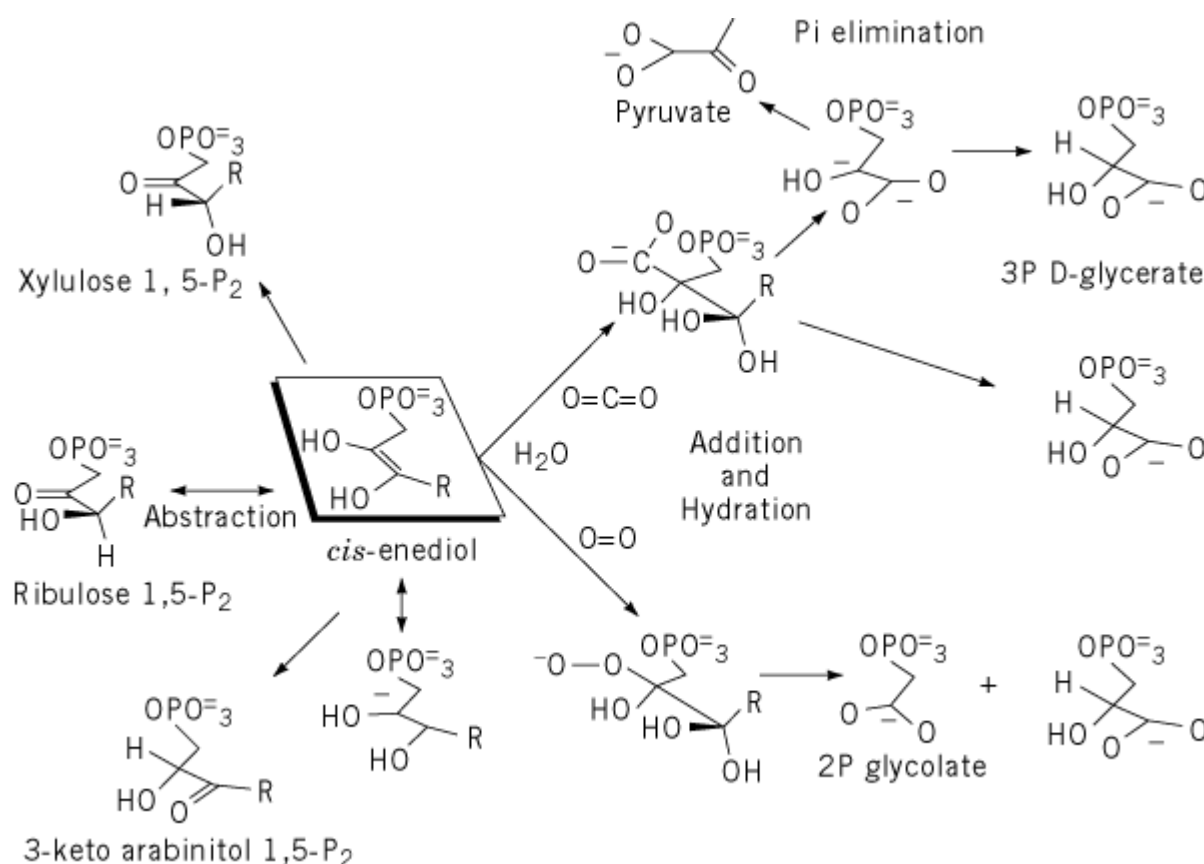
### 1. Reactions Catalyzed by Rubisco

Rubisco catalyzes the first reaction in the pathways of both photosynthesis and photorespiration ([1](#)), which are respectively its carboxylase and oxygenase activities:



Photosynthesis and photorespiration are interlocking metabolic cycles, with Rubisco determining the relative rates of carbon flow through the two pathways (Fig. 1). Rubisco requires both  $\text{CO}_2$  and  $\text{Mg}^{2+}$  as obligatory **cofactors**, resulting in carbamylation of an **active-site** lysine residue by  $\text{CO}_2$  and coordination of  $\text{Mg}^{2+}$  (1, 2).

**Figure 1.** The reactions catalyzed by Rubisco. The first intermediate of catalysis is the C2,C3 *cis*-enediol form of ribulose-bisphosphate (ribulose 1,5- $\text{P}_2$ ) after abstraction of the C3 proton. The enediol can partition a number of ways, the majority into the products of carboxylation (upper reactions) or oxygenation (lower reactions). However, a number of misprotonated isomers of ribulose 1,5- $\text{P}_2$ , for example, xylulose-bisphosphate, have been detected with the wild-type enzyme that are produced in quantity by mutations of specific amino acid residues involved in proton transfer. Phosphate elimination of the carbanion forms of intermediates are also produced by some mutants. R,  $\text{—CHOH—CH}_2\text{OPO}_3^-$ ; 3P D-glycerate, 3-phospho glycerate; 2P glycolate, 2-phosphoglycolate. (Gutteridge and Gatenby, 1995; the material is copyrighted by the American Society of Plant Physiologists and is reprinted with permission.)



Five discrete partial reactions have been described for the carboxylation reaction (Eq. 1) of ribulose bisphosphate (3). The initial formation of a C2,C3-enediol intermediate of the ribulose bisphosphate substrate is followed by its carboxylation, where the enediol reacts with  $\text{CO}_2$  at the C2 position. The

resulting six-carbon intermediate is hydrolytically cleaved to two molecules of 3-phosphoglycerate, resulting in an overall gain of carbon during photosynthesis.

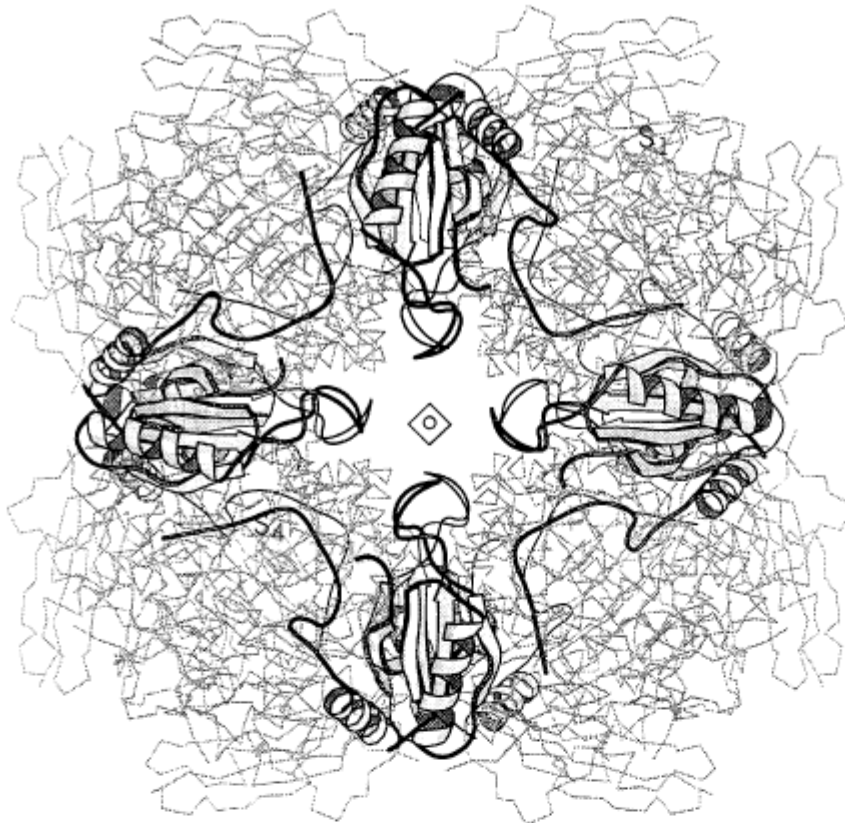
In the oxygenase reaction (Eq. 2) of photorespiration, however, the enediol intermediate reacts with molecular oxygen to form a hydroperoxy derivative that breaks down to 2-phosphoglycolate and 3-phosphoglycerate. Photorespiration consumes ribulose biphosphate, which results in no net gain of carbon, and also requires the consumption of energy to recycle the lost carbon. Partitioning of ribulose biphosphate between the two reactions of photosynthesis and photorespiration can vary significantly between different photosynthetic organisms (4), which has stimulated interest in improving the carboxylation efficiency of crop plants.

Rubisco activity is modulated in plants by an inhibitor and an activator. The inhibitor 2'-carboxy arabinitol 1-phosphate (2CAIP) accumulates in some plants during darkness and binds to the active site of Rubisco (5, 6). 2CAIP is degraded by a specific phosphatase, which presumably allows Rubisco to function during photosynthesis in the light. Rubisco can be severely inhibited by a range of sugar biphosphates, including substrate analogues. The enzyme *Rubisco activase* has the ability to relieve the inhibition caused by sugar biphosphates (7), possibly by interacting with Rubisco and altering the affinity of the enzyme for biphosphates.

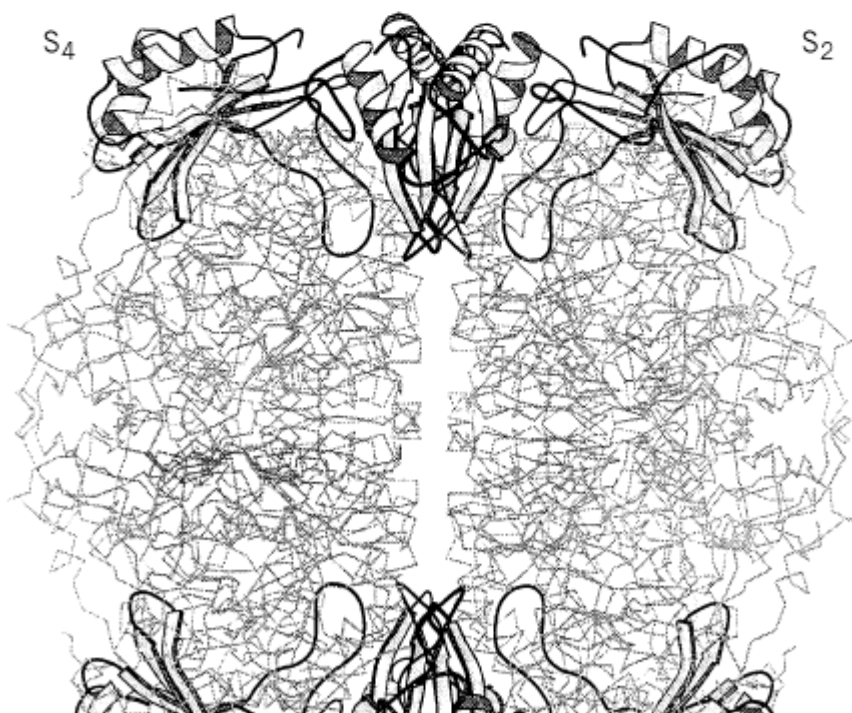
## 2. Structure of Rubisco

Two distinct forms of Rubisco found in different species can be differentiated by their [quaternary structures](#). Form I Rubisco is the most common structure found in photosynthetic bacteria and plants (Fig. 2a). It has a hexadecameric structure ( $L_8S_8$ ) with an equivalent number of both large (L) (50–55-kDa) and small (S) (12–18-kDa) subunits. Form II Rubisco has a much simpler structure, containing L subunits alone. It is less efficient than the form I enzyme and is present in some photosynthetic bacteria that are sensitive to oxygen.

**Figure 2.** The organization of the subunits of  $L_8S_8$  Rubisco. In (A) the display is from above the fourfold axis of symmetry of the enzyme. The four L-subunit dimers are shown only schematically by a light trace of the positions of the Ca atoms of the backbone, whereas the structural elements of the S subunits have been displayed with Molscript. The four S subunits shown reside at the top of the molecule situated between each L-subunit dimer, with a loop extending into, but not obscuring, the prominent central channel of the  $L_8$  core. In (B), the other four S subunits are found to occupy the same position at the bottom of the core. (This figure was kindly provided by Steven Gutteridge and is reprinted with permission from JAI Press Inc.).



S<sub>3</sub>  
Spinach rubisco  
(a)



Structural analysis of the  $L_2$  Rubisco dimer from *Rhodospirillum rubrum* reveals that amino acid residues essential for the function of each of the two active sites are located on both L subunits (8). Although the quaternary structure of  $L_8S_8$  Rubisco from plants is clearly more complex than the *R. rubrum* homodimer, X-ray crystallographic studies reveal that the  $L_8$  core of  $L_8S_8$  molecules can be built up from four *R. rubrum*-like dimers arranged around a fourfold axis (9). The eight S-subunits occupy large crevices between the ends of the  $L_2$ -like dimers at the top and bottom of the  $L_8$  core, making extensive interactions with L subunits to stabilize the  $L_8$  core. The S subunits do not contribute directly to the active site, but they do influence substrate affinities and turnover through contacts to elements that form the site (10).

### 3. Gene Organization and Expression

In most **prokaryotes**, the organization of the Rubisco S-subunit (*rbcS*) and L-subunit (*rbcL*) genes is relatively simple, with both genes usually closely linked and **transcribed** together. To the 5' side of the initiator methionine residue of many *rbc* genes is a sequence similar to a **Shine–Dalgarno** site for **ribosome** binding, and sufficiently conserved to allow correct translation when **cloned** into *Escherichia coli* (11, 12).

Plant cells compartmentalize Rubisco in chloroplasts, but the genetic information is shared between chloroplast and nucleus. The *rbcL* genes are present in chloroplast DNA (13), and their transcription and **translation** in plastids uses sequences that are similar to those found in prokaryotes (14, 15), to the extent of allowing direct expression when transferred to *E. coli* (16). The S-subunit genes are located on nuclear **chromosomes** and have a more complex structural arrangement. The *rbcS* genes contain **introns** and are present as small **multigene families** that are often closely linked (17). Light-induced expression is mediated by both phytochrome and blue light photoreceptors (18), and positive and negative regulatory sequences are located in **cis-acting** transcriptional control regions. The *rbcS* **promoters** also appear to contain nuclear **matrix attachment regions** (MARs) (19), which may be important for their expression. The highest level of *rbcS* mRNA is found in leaves, but it is also found in the photosynthetic tissues in stems, petals and pods. S subunits are synthesized on free cytoplasmic **polyribosomes** as precursor molecules with an *N*-terminal **transit peptide** (20). The S-subunit precursors are imported post-translationally into chloroplasts in a process requiring ATP, and the transit peptide is removed (21, 22).

### 4. Rubisco Assembly

Following import into chloroplasts and removal of the transit peptide, mature S subunits are assembled with chloroplast-synthesized L subunits to give the active  $L_8S_8$  Rubisco holoenzyme (21, 22). This assembly process requires the assistance of another chloroplast protein (23) now known as **chaperonin** 60 (cpn60) (24, 25). In fact, studies on the assembly of Rubisco in chloroplasts and bacteria (23, 26, 27) led to the discovery of the molecular chaperone cpn60 and its role in the correct folding of Rubisco and many other proteins (24, 25, 28). Productive folding of Rubisco requires  $Mg^{2+}$ , ATP hydrolysis, and a smaller cochaperonin molecule (29). Cpn60-mediated folding of Rubisco in bacteria uses a cochaperonin oligomer with 10-kDa subunits (24), but the folding of Rubisco in chloroplasts seems to involve a co-chaperonin oligomer with 21-kDa subunits (30). The reason for this larger cochaperonin in chloroplasts and the mechanistic details of the Rubisco assembly process in plants are currently under investigation.

### Bibliography

1. H. M. Miziorko and G. H. Lorimer (1983) *Annu. Rev. Biochem.* **52**, 507–535.
2. G. H. Lorimer and H. M. Miziorko (1980) *Biochemistry* **19**, 5321–5328.
3. T. J. Andrews and G. H. Lorimer (1987) in *The Biochemistry of Plants*, Vol. **10**, M. D. Hatch and N. K. Boardman, eds., Academic Press, San Diego, pp. 131–218.
4. D. B. Jordon and W. H. Ogren (1981) *Nature* **291**, 513–515.
5. S. Gutteridge et al. (1986) *Nature* **324**, 274–276.
6. J. A. Berry et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 734–738.
7. A. R. Portis (1992) *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **43**, 415–437
8. G. Schneider, Y. Lindqvist, C.-I. Branden, and G. H. Lorimer (1986) *EMBO J.* **5**, 3409–3415.
9. M. S. Chapman et al. (1988) *Science* **241**, 71–74.
10. G. Schneider et al. (1990) *EMBO J.* **9**, 2045–2050.
11. A. A. Gatenby, S. M. van der Vies, and D. Bradley (1985) *Nature* **314**, 617–620.
12. S. M. van der Vies, D. Bradley, and A. A. Gatenby (1986) *EMBO J.* **5**, 2439–2444.
13. L. McIntosh, C. Poulsen and L. Bogorad (1980) *Nature* **288**, 556–560.
14. D. Bradley and A. A. Gatenby (1985) *EMBO J.* **4**, 3641–3648.
15. W. Gruissem and G. Zurawski (1985) *EMBO J.* **4**, 3375–3383.
16. A. A. Gatenby, J. A. Castleton, and M. W. Saul (1981) *Nature* **291**, 117–121.
17. C. Dean et al. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 4964–4968.
18. R. Fluhr and N.-H. Chua (1986) *Proc. Natl. Acad. Sci. USA* **83**, 2358–2362.
19. I. Meier, B. Groning, S. Michalowski, and S. Spiker (1997) *FEBS Lett.* **415**, 91–95.
20. B. Dobberstein, G. Blobel, and N.-H. Chua (1977) *Proc. Natl. Acad. Sci. USA* **74**, 1082–1085.
21. N.-H. Chua and G. W. Schmidt (1978) *Proc. Natl. Acad. Sci. USA* **75**, 6110–6114.
22. S. M. Smith and R. J. Ellis (1979) *Nature* **278**, 662–664.
23. R. Barraclough and R. J. Ellis (1980) *Biochim. Biophys. Acta* **608**, 19–31.
24. P. Goloubinoff, A. A. Gatenby, and G. H. Lorimer (1989) *Nature* **337**, 44–47.
25. T. H. Lubben, G. K. Donaldson, P. V. Viitanen, and A. A. Gatenby (1989) *Plant Cell* **1**, 1223–1230.
26. M. V. Bloom, P. Milos and H. Roy (1983) *Proc. Natl. Acad. Sci. USA* **80**, 1013–1017.
27. A. A. Gatenby, T. H. Lubben, P. Ahlquist, and K. Keegstra (1988) *EMBO J.* **7**, 1307–1314.
28. S. M. Hemmingsen et al. (1988) *Nature* **333**, 330–334.
29. P. Goloubinoff, J. T. Christeller, A. A. Gatenby, and G. H. Lorimer (1989) *Nature* **342**, 884–889.
30. F. Baneyx et al. (1995) *J. Biol. Chem.* **270**, 10695–10702.

### **Suggestions for Further Reading**

31. S. Gutteridge and A. A. Gatenby (1995) Rubisco synthesis, assembly, mechanism, and regulation. *Plant Cell* **7**, 809–819.
32. S. Gutteridge and T. Lundqvist (1994) Structural elements involved in the assembly and mechanism of action of rubisco. *Adv. Mol. Cell Biol.* **10**, 287–335.
33. F. C. Hartman and M. R. Harpel (1993) Chemical and genetic probes of the active site of D-ribulose-1,5-bisphosphate carboxylase/oxygenase: A retrospective based on the three-dimensional structure. *Adv. Enzymol.* **67**, 1–75.
34. G. Schneider, Y. Lindqvist, and C.-I. Branden (1992) Rubisco: structure and mechanism. *Annu. Rev. Biophys. Biomol. Struct.* **21**, 119–143.
35. R. J. Spreitzer (1993) Genetic dissection of rubisco structure and function. *Annu. Rev. Plant Physiol. Plant Mol. Biol.* **44**, 411–434.

## Ricin And Ribosome-Inactivating Plant Toxins

Ricin is the prototype of a large family of toxic **glycoproteins** expressed by several **plant** tissues, particularly by seed endosperm cells (1). These plant **toxins** are synthesized as pre-pro-toxins, with a **leader sequence** of about 35 residues, which are then processed proteolytically to A-B type toxins. The mature form of isotype D of ricin (from *Ricinus communis* seeds) has an A chain of 267 residues **disulfide bond**-linked to a B chain of 262 residues (2).

The ricin B chain binds to a large number of different glycolipid and glycoprotein cell surface **receptors** containing a terminal galactose residue. Ricin is then internalized by **endocytosis** inside **clathrin**-coated and uncoated **vesicles** that merge in **endosomal** compartments, from where it can follow different intracellular trafficking pathways (3). Some molecules are recycled back to the cell surface, whereas others are degraded inside **lysosomes**, and still other ricin molecules are transported retrogradely through the **Golgi** stacks to the **endoplasmic reticulum** (3). There is evidence that ricin has to reach the Golgi apparatus before the A chain can translocate into the cytosol. Other A-B plant toxins, such as modeccin, abrin, and viscumin, may follow different intracellular routes, depending on their cell surface receptors (1).

These toxins arrest cell protein biosynthesis by inactivating the 60 S subunit of **ribosomes** (hence, they are called RIPs, ribosome-inactivating proteins) (1). Once in the cytosol, their A chains display a hydrolytic activity specific for the N-glycosidic bond of Adenine4324 of 28S ribosomal RNA (4), which is at the center of a ribosome area involved in binding of aminoacyl-**transfer RNA**, involving **elongation factor** 1, and in GTP hydrolysis and translocation, catalyzed by elongation factor 2 (5). The same enzymatic activity is associated with the A chain of Shiga and Shigalike toxins produced by *Shiga spp.* and *E. coli spp.*

Many plants also produce single-chain RIPs consisting of solely the A chain (1). These RIPs are a thousand times less toxic than the A-B type toxins, because of their very inefficient entry into cells. They are increasingly used to assemble immunotoxins, whose cell binding and intracellular routing is determined by the antibody part (see **Antibody-Conjugated Toxins**). Also, **fungi** of the genus *Aspergillus* produce a family of protein toxins of potential usefulness in the preparation of immunotoxins. They are single polypeptide chains of about 150 amino acid residues with two intramolecular **disulfide bonds** endowed with a **ribonuclease** activity specific for the phosphodiester bond on the 3' side of Guanine4325 of 28 S ribosomal RNA (6). They and ricin recognize the same essential GAGA sequence, which is present in a conserved loop of the 28 S ribosomal RNA molecule (5).

### Bibliography

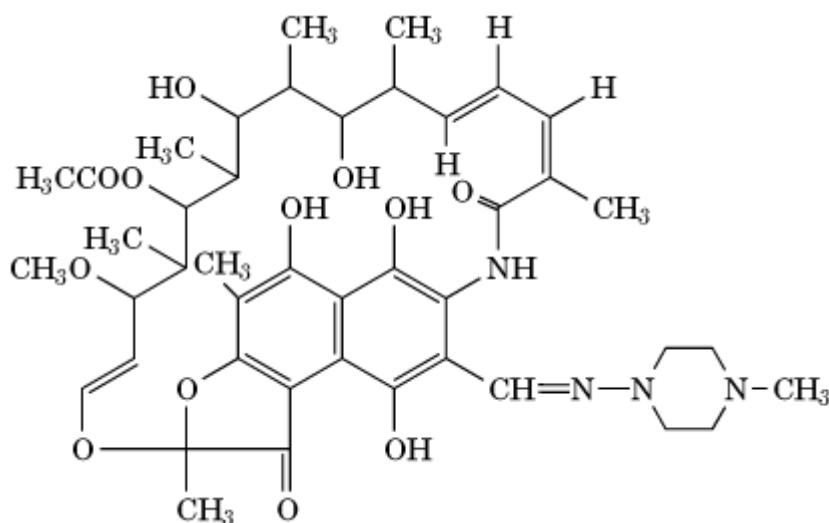
1. L. Barbieri, M. G. Battelli, and F. Stirpe (1993) Biochim. Biophys. Acta **1154**, 237–282.
2. J. M. Lord, L. M. Roberts, and J. D. Robertus (1994) FASEB J. **8**, 201–208.
3. K. Sandvig and B. van Deurs (1994) FEBS Lett. **346** 99–102.
4. Y. Endo and K. Tsurugi (1987) J. Biol. Chem. **262**, 8128–8230.
5. I. G. Wool, A. Gluck, and Y. Endo (1992) Trends Biochem. Sci. **17**, 266–269.
6. B. Fleischer (1995) Rev. Med. Microbiol. **6**, 49–57.

## Rifampicin

Rifampicin (RIF), also known as rifampin (1971), is a semisynthetic antibacterial agent derived from rifamycin B, which was isolated from *Nocardia mediterranei* in 1957 (1960). RIF has wide bactericidal activity against **Gram-positive** and **Gram-Negative Bacteria**, both intracellular and extracellular organisms. Since it was introduced in 1968, it has had wide clinical use, especially as an effective antituberculosis drug; even now, it is a very important agent for the treatment of many infections. Recently, it was reported that RIF activates the human **glucocorticoid** receptor and acts as an immunodepressant (1998).

RIF is a semisynthetic derivative of rifamycin B, a fermentation product of the mold *Nocardia mediterranei*, which was reclassified in 1969 as *Streptomyces mediterranei* (1969). RIF is 3-(4-methyl-piperazinyl-iminomethyl)-rifamycin SV ( $C_{43}H_{58}N_4O_{12}$ , molecular weight 822.95) (1996) (Fig. 1), which belongs to a group of macrocyclic compounds (ansamycin) (1972). RIF is used as a first-line drug in antitubercular therapy, along with isoniazid (INH), pyrazinamid, and **streptomycin** (or ethambutol), as recommended by the World Health Organization (WHO) (1991).

**Figure 1.** The structure of rifampicin.



RIF exhibits bactericidal activity by inhibiting RNA synthesis (**transcription**) and, consequently, the initiation of a new round of **DNA replication** (198319841977). It inhibits transcription by binding to prokaryotic DNA-dependent **RNA polymerase**, which consists of four subunits (a, b, b', d, encoded by the *rpoA*, *rpoB*, *rpoC*, and *rpoD* genes, respectively). RIF binds only to the holoenzyme of RNA polymerase ( $\alpha_2\beta\beta'$  and d) and not to the individual b subunit of RNA polymerase, even though mutation in the *rpoB* gene encoding the b subunit causes RIF resistance. RIF inhibits the initiation of transcription by interfering with the first translocation of RNA polymerase on template DNA and/or with the dissociation of the b subunit from template DNA (19761979), but it does not affect chain elongation. It also binds to viral DNA-dependent RNA polymerase and **reverse transcriptase**. In contrast, eukaryotic nuclear RNA polymerase is less sensitive to RIF. Therefore, the toxicity of RIF



is highly selective.

Resistance to RIF is conferred by drug inactivation, the alteration of either its target sites, or the changes in [membrane permeability \(199519941994\)](#). Drug inactivation occurs by phosphorylation, glucosylation, and ribosylation in *Nocardia* and rapidly growing mycobacteria ([1995](#)). The alterations of target sites occur by missense, deletion, and insertion [mutations](#) in a restricted central portion of the *rpoB* gene that encodes the  $\beta$  subunit of RNA polymerase ([1988](#)), which in *Escherichia coli* consists of 1342 amino acid residues ([1981](#)). Jin and Gross ([1988](#)) classified the mutated region into three clusters (I, II, and III). Clusters I and II comprise amino acid residues 507–533 and 563–572, respectively, whereas cluster III possesses one affected site at amino acid 687. In addition, Lisitsyn et al ([1984](#)) found one mutation site at residue 146. Resistance mutations occur most frequently in cluster I. Mutations in clusters I and II confer especially high-level RIF resistance, whereas mutations in cluster III or at position 146 confer only low-level RIF resistance. Alterations in membrane permeability were reported in *Neisseria meningitidis* ([1996](#)), *Pseudomonas aeruginosa* ([1982](#)), and mycobacteria ([1977](#)). Recently, increased efflux of the drug from the cell was found in resistant *Pseudomonas fluorescens*, which was mediated by a **plasmid** ([1998](#)). It was reported that RIF-resistant strains of *Francisella tularensis* have reduced virulence ([1994](#)).

Because of the emergence of multiple **drug resistant** *M. tuberculosis* strains, the mechanism of RIF resistance has been studied on the basis of data from *E. coli* ([19931995](#)). The  $\beta$  subunit of RNA polymerase from *M. tuberculosis* consists of 1178 amino acids ([1994](#)). The region corresponding to that mutated in the *rpoB* gene of *E. coli* was sequenced. Nearly 95% of RIF-resistant *M. tuberculosis* strains have mutations in this region of the *rpoB* gene. Of these, more than 90% have missense mutations in the region corresponding to cluster I of the *E. coli rpoB* gene; mutations at residues 526 and 531 were found most often. Even in cluster I, hot spots occur frequently in resistant mutations, and the level of RIF resistance depends on which and where the mutation occurred, as had been found in *E. coli* ([19961996](#)).

RIF has been used also for the treatment of infection by *M. avium-intracellulare* complex and *M. leprae*. The in vitro RIF susceptibility varies significantly among *M. avium-intracellulare* complex serovars. Most *M. avium-intracellulare* complex organisms are resistant in vitro, even though only one of 31 resistant strains possessed a missense mutation in the region equivalent to cluster I of the *rpoB* gene of *E. coli* ([1994](#)). RIF resistance of *M. avium-intracellulare* complex organisms is assumed to be conferred by other mechanisms, presumably by decreased uptake or increased efflux of the drug. RIF-resistant strains of *M. leprae* have mutations in cluster I of the *rpo B* gene, although mutations occur most frequently at residue 531, instead of at residue 526 ([1993](#)).

Rifabutin, a spiro-piperidyl derivative of rifamycin S ([1987](#)), has been used for some RIF-resistant *M. tuberculosis* and *M. avium-intracellulare* complex infections in AIDS patients. Mutations at positions 513, 526, or 531 in the *rpo B* gene had high-level cross resistance to RIF and rifabutin, whereas mutations at residue 511, 516, or 522 do not confer rifabutin resistance.

RIF shows wide bactericidal activity against Gram-positive and Gram-negative bacteria and both intracellular and extracellular microorganisms ([1996](#)). It is able to enter phagocytes and kill living intracellular organisms ([1983](#)). RIF is bactericidal for wild-type *M. tuberculosis* and affects intracellular, slowly replicating bacilli in caseous lesions, as well as actively replicating tubercle bacilli in open pulmonary cavities. However, it shows bacteriostatic activity against enterococci. It has a broad antibacterial spectrum and is also effective against *P. aeruginosa*, *M. tuberculosis*, *M. leprae*, *M. kansasii*. It shows variable activity against *M. avium-intracellulare* complex and *Legionella* species. It is also highly active against *Staphylococcus aureus*, *N. meningitidis*, and *Haemophilus influenzae*. Its antibacterial activity is high; for example, *S. aureus*, *M. tuberculosis*, and *E. coli* are inhibited at concentrations of 0.002, 0.2, and 20 mg/ml, respectively.

RIF is well absorbed from the gastrointestinal tract and is well distributed to almost all body tissues and fluids. RIF is activated in the liver and excreted in the bile; therefore, the dosage should be

adjusted for patients with severe hepatic dysfunction. About 30% to 40% of the administered drug is excreted in the urine.

Adverse effects include hepatotoxicity, gastrointestinal and hypersensitivity reactions (drug fever, skin rashes, and eosinophilia), and a red-orange discoloration of urine, tears, saliva, sweat, and other body fluids (1996). RIF also induces increased hepatic metabolism of a wide variety of other drugs. RIF-induced hepatitis occurs more frequently during concurrent isoniazid therapy for tuberculosis. The drug antagonizes the effect of oral contraceptives and diminishes the anticoagulant activity of warfarin. Recently, it was reported that RIF activates the human glucocorticoid receptor, a nuclear [steroid hormone receptor](#), and acts as an immunodepressant (1998).

## Bibliography

- W. Wehrli and M. Staehelin (1971) *Bacteriol. Rev.* **35**, 290–309.
- P. Sensi, A. M. Greco and R. Ballotta (1960) *Antibiot. Ann.* 262–270.
- C. Calleja, J. M. Pascussi, J. C. Mani, P. Maurel and M. J. Vilarem (1998) *Nature Med.* **4**, 92–96.
- J. E. Thiemann, G. Zucco and G. Pelizza (1969) *Arch. Microbiol.* **67**, 147–155.
- S. Budavari (1996) *The Merck Index: An Encyclopedia of Chemicals, Drugs, and Biologicals*, Merck Research Laboratories, Whitehouse Station, N.J., pp. 1413–1414.
- K. L. Rinehart (1972) *Accounts Chem. Res.* **5**, 57.
- WHO (1991) *Guidelines for the treatment in adults and children in national tuberculosis programmes*, WHO/Tub, 161.
- W. Wehrli (1983) *Rev. Infect. Dis.* **5**(suppl.3), S407–S411.
- T. Atlung (1984) *Mol. Gen. Genet.* **197**, 125–128.
- M. M. Bagdasarian, M. Izakowska and M. Bagdasarian (1977) *J. Bacteriol.* **130**, 577–582.
- F. Y. Wu, L. R. Yarbrough and C. W. Wu (1976) *Biochemistry* **15**, 3254–3258.
- A. L. Sonenshein and H. B. Alexander (1979) *J. Mol. Biol.* **127**, 55–72.
- E. R. Dabbs, K. Yazawa, Y. Mikami, M. Miyaji, N. Morisaki, S. Iwasaki, and K. Furihata (1995) *Antimicrob. Agents Chemother.* **39**, 1007–1009.
- B. G. Spratt (1994) *Science* **264**, 388–393.
- H. Nikaido (1994) *Science* **264**, 382–388.
- D. J. Jin and C. A. Gross (1988) *J. Mol. Biol.* **202**, 45–58.
- Y. A. Ovchinnikov, G. S. Monastyrskaya, V. V. Gubanov, S. O. Guryev, O. Yu Chertov, N. N. Modyanov, V. A. Grinkevich, I. A. Makarova, T. V. Marchenko, I. N. Polovnikova, V. M. Lipkin and E. D. Sverdlov (1981) *Eur. J. Biochem.* **116**, 621–629.
- N. A. Lisitsyn, E. D. Sverdlov, E. P. Moiseyeva, O. N. Danilevskaya and V. G. Nikiforov (1984) *Mol. Gen. Genet.* **196**, 173–174.
- F. J. Abadi, P. E. Carter, P. Cash and T. H. Pennington (1996) *Antimicrob. Agents Chemother.* **40**, 646–651.
- B. L. Angus, A. M. Carey, D. A. Caron, A. M. Kropinski and R. E. Hancock (1982) *Antimicrob. Agents Chemother.* **21**, 299–309.
- J. Hui, N. Gordon and R. Kajioka (1977) *Antimicrobiol. Agents Chemother.* **11**, 773–779.
- S. Chandrasekaran and D. Lalithakumari (1998) *J. Med. Microbiol.* **47**, 197–200.
- N. Bhatnagar, E. Getachew, S. Straley, J. Williams, M. Meltzer and A. Fortier (1994) *J. Infect. Dis.* **170**, 841–847.
- A. Telenti, P. Imboden, F. Marchesi, D. Lowrie, S. Cole, M. J. Colston, L. Matter, K. Schopfer and T. Bodmer (1993) *Lancet* **341**, 647–650.
- J. M. Musser (1995) *Clin. Microbiol. Rev.* **8**, 496–514.
- L. P. Miller, J. T. Crawford and T. M. Shinnick (1994) *Antimicrob. Agents Chemother.* **38**, 805–811.

- H. Taniguchi, H. Aramaki, Y. Nikaido, Y. Mizuguchi, M. Nakamura, T. Koga and S. Yoshida (1996) *FEMS Microbiol. Lett.* **144**, 103–108.
- H. Ohno, H. Koga, S. Kohno, T. Tashiro and K. Hara (1996) *Antimicrob. Agents Chemother.* **40**, 1053–1056.
- C. Guerrero, L. Stockman, F. Marchesi, T. Bodmer, G. D. Roberts and A. Telenti (1994) *J. Antimicrob. Chemother.* **33**, 661–663.
- N. Honore and S. T. Cole (1993) *Antimicrob. Agents Chemother.* **37**, 414–418.
- R. J. O'Brien, M. A. Lyle and D. E. Snider Jr. (1987) *Rev. Infect. Dis.* **9**, 519–530.
- G. L. Mandel and W. A. Petri Jr. (1996) In *Goodman & Gilman's The Pharmacological Basis of Therapeutics* (J. G. Hardman and L. E. Limbird, eds.), McGraw-Hill, New York, pp. 1159–1161.
- G. L. Mandell (1983) *Rev. Infect. Dis.* **5** (Suppl. 3), S463–S467.

## RNA Blots (Northern Blots)

Northern blots are similar to [Southern blots](#), but the polynucleotide being **blotted** and immobilized is **RNA**, rather than DNA (see [Blotting](#) and [Southern Blots \(DNA Blots\)](#)). In this technique, for example, samples of total cellular RNA are resolved by [electrophoresis](#) in denaturing [agarose](#) gels containing **formaldehyde**, they are subsequently transferred to a **blotting matrix**, ultimately to be probed with **radioactive** DNA to identify RNA molecules complementary to the DNA probe ([1-3](#)). The major obstacle in developing this procedure was that RNA did not bind efficiently to [nitrocellulose](#) matrices under conditions suitable for DNA blots. In 1977, however, chemically reactive blotting matrices [such as diazobenzylxymethyl (DBM) filters] were used to bind the blotted RNA covalently ([4](#)), thus illustrating the feasibility of “Northern blotting,” the term coined as a play on words in contrast to Southern's DNA transfers ([5](#)).

Since the original protocols for RNA blotting were introduced, conditions for using of nitrocellulose have been developed ([6](#)). Moreover, a common alternative matrix is a nylon membrane ([3](#)). In principle, the analyses and applications of RNA blots do not differ from the methods used for Southern blots, except for the additional requirement that caution should be applied in handling RNA samples because this nucleic acid is markedly more sensitive to degradation than DNA.

## Bibliography

1. J. Meinkoth and G. Wahl (1984) *Anal. Biochem.* **138**, 267–284.
2. (1988) *Nucleic Acid Hybridization: A Practical Approach* (B. D. Hames and S. J. Higgins, eds.), IRL Press, Oxford, UK.
3. L. G. Davis, M. D. Dibner, and J. F. Battey (1986) *Basic Methods in Molecular Biology*, Elsevier, New York, pp. 143–146.
4. J. C. Alwine, D. J. Kemp, and G. R. Stark (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5350–5354.
5. J. C. Alwine, D. J. Kemp, B. A. Parker, J. Reiser, J. Renart, G. R. Stark, and G. M. Wahl (1979) *Methods Enzymol.* **68**, 220–242.
6. P. S. Thomas (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5201–5205.

## RNA Cap

The RNA cap refers to the [cap](#) modification that is characteristic to **eukaryotic** RNAs **transcribed** by **RNA polymerase II**. It consists of an inverted 7-methyl-guanosine nucleotide linked via a 5'-5'-triphosphate bridge to the first base of the [messenger RNA](#) precursor. The cap plays an important role in the metabolism of RNA polymerase II transcripts. It provides stability against 5'-3' **exonucleases** ([1](#)), is required for efficient pre-mRNA **splicing** ([2](#), [3](#)), [polyadenylation](#) ([4](#), [5](#)), and contributes to the export of U **small nuclear RNA** (snRNA) from the [nucleus](#) to the **cytoplasm** ([6](#), [7](#)). It also increases the rate of mRNA export to a lesser extent, approximately twofold ([6](#)). In the cytoplasm, the cap increases the efficiency of mRNA [translation](#) ([8](#)).

### Bibliography

1. Y. Furuichi, A. LaFiandra, and A. J. Shatkin (1977) *Nature* **266**, 235–239.
2. A. Krainer, T. Maniatis, B. Ruskin, and M. Green (1984) *Cell* **36**, 993–1005.
3. M. Konarska, R. Padgett, and P. Sharp (1984) *Cell* **38**, 731–736.
4. C. Cooke and J. C. Alwine (1996) *Mol. Cell. Biol.* **16**, 2579–2584.
5. S. M. Flaherty, P. Fortes, E. Izaurralde, I. W. Mattaj, and G. M. Gilmartin (1997) *Proc. Nat. Acad. Sci. USA* **94**, 11893–11898.
6. A. Jarmolowski, W. C. Boelens, E. Izaurralde, and I. W. Mattaj (1994) *J. Cell. Biol.* **124**, 627–635.
7. J. Hamm and I. W. Mattaj (1990) *Cell* **63**, 109–118.
8. A. Shatkin (1985) *Cell* **40**, 223–224.

### Suggestion for Further Reading

9. J. D. Lewis and E. Izaurralde (1997) The role of the cap structure in RNA processing and nuclear export. *Eur. J. Biochem.* **247**, 461–469.

## RNA Degradation *In Vitro*

The ultimate fate of all cellular **RNA** molecules is their destruction. The final products are mononucleotides, which can be used in new **RNA biosynthesis**. RNA degradation *in vivo* involves site-specific cleavages performed by a diverse ensemble of cellular **ribonucleases** (RNase). These [enzymes](#) include both **endonucleases** and **exonucleases** that recognize specific RNA sequences and structural features and can act in a coordinated manner to establish a specific decay pathway. RNA degradation can also regulate **gene expression** by controlling the levels of [messenger RNA](#) (mRNA). For example, an increased rate of mRNA degradation lowers the steady-state level of that mRNA and, consequently, the amount of [protein](#) expressed from it. Cells can also control gene expression quickly by coupling rapid mRNA degradation with cessation of its biosynthesis by [transcription](#). The physical (or chemical) **half-life** of an RNA is distinct from its functional half-life. The physical half-life is a measure of the rate of chemical degradation of an RNA species, whereas the functional half-life is a measure of the rate at which it becomes inactive and perhaps remains chemically stable. Functional inactivation (eg, repression of an mRNA to prevent its [translation](#)) can be caused by a change in RNA structure or by binding of a protein factor. Often, mRNA functional inactivation

precedes physical inactivation. There are significant differences in the mechanisms and pathways of RNA decay in **prokaryotic** and **eukaryotic** cells because of their different mRNA structural features; the coupling of bacterial transcription and translation; and the subcellular compartmentalization of RNA synthesis and utilization in the eukaryotic cell. However, there is growing evidence for the functional conservation of a number of enzymes involved in RNA decay.

## 1. Bacterial mRNA Decay

Bacterial mRNA are generally short-lived species that have physical half-lives typically ranging from 1 to 3 minutes (1, 2). Degradation of an mRNA can occur even before its synthesis is complete, when full-length transcripts are rarely observed. However, some mRNA have half-lives greater than the cell division time (which for *Escherichia coli* is 20 minutes at 37°C). The lifetime of an mRNA depends on the accessibility of its potential cleavage sites and on the amount and location of the **nucleases**. The substrate for degradation is not naked mRNA, but one that is covered with translating **ribosomes**. Translocating ribosomes can occlude cleavage sites, so the frequency of translation can control the rate of decay. A major mRNA decay pathway begins with a site-specific cleavage, often near the 5' end. This cleavage can be **rate-limiting** and is usually carried out by RNase E or in several cases by RNase III (2, 3). Cleavage blocks further translation initiation, which allows internal cleavage sites to appear after passage of the last translocating ribosome. This pathway fragments mRNA with an overall 5' to 3' directionality (1, 2). Further degradation of the mRNA fragments is accomplished by the combined action of **poly(A)** polymerase and one or more exoribonucleases. Poly(A) polymerase adds a poly(A) tail of variable length to the 3' ends of the RNA fragments, creating substrates for the 3' → 5' exoribonucleases RNase II or **polynucleotide phosphorylase**. The former enzyme is a **phosphodiesterase** that yields ribonucleoside 5'-monophosphates, whereas the latter uses phosphate as a cosubstrate to produce ribonucleoside 5'-diphosphates (3). It is not known if these exonucleases completely digest the RNA fragments, but a third 3' → 5' exonuclease, oligoribonuclease, degrades short RNA to mononucleotides (3).

The major mRNA decay pathway is mediated by a multisubunit complex called the RNA degradosome (4, 5). RNase E and polynucleotide phosphorylase are integral components of this complex. Other components include enolase, the heat-shock **chaperonin**, **DnaK**, and an **RNA helicase**. The RNA helicase may melt RNA secondary structure, providing more efficient degradation. Other enzymatic activities occur within this complex but are only partially characterized. RNA degradosome activity may be regulated by the physiological state of the cell. There are mRNA decay pathways, however, that are independent of the RNA degradosome (6). Stable RNA in cells includes **ribosomal** RNA and **transfer RNA**, which are normally not degraded. However, stable RNA degradation can be triggered under certain conditions, such as carbon starvation (3), and provide precursors for new RNA synthesis in the absence of *de novo* nucleotide synthesis. The nucleases involved in this process are not known.

## 2. Eukaryotic RNA Decay

The physical half-lives of eukaryotic mRNA range from minutes to hours. Two key structural features of eukaryotic mRNA—the 7-methylguanosine **cap** and the poly(A) tail at the 3' end—are key participants in translation and also in degradation. The compartmentalized structure of the eukaryotic cell also plays a key role in RNA degradation. Thus, the nuclear membrane separates the synthesis and maturation of mRNA in the **nucleus** from its translation and degradation in the **cytoplasm**. RNA decay in **mitochondria** and **chloroplasts** also represent distinct pathways. A diverse, yet interrelated set of decay pathways has been revealed by genetic and biochemical studies in **yeast** (7, 8). A major pathway for mRNA decay is initiated by the shortening of the poly(A) tail, which is accomplished by a poly(A)-specific, 3' → 5' exoribonuclease. The rate of poly(A) tail shortening depends on specific sequences in the mRNA 3' **untranslated region** (3'-UTR), which are recognized by specific proteins. A decrease in poly(A) tail length below a certain size triggers the hydrolytic

removal of the Cap structure. The yeast decapping activity, Dcp1, releases  $m^7GDP$  and creates an RNA chain with a 5' terminal monophosphate (5'-pRNA) (9). A second, nuclear-localized decapping activity may also exist. The 5'-pRNA is a substrate for 5' → 3' exoribonuclease action. In yeast this is accomplished by XRN1 that creates 5'-mononucleotides (7, 8) [in contrast, a 5' → 3' exoribonuclease activity has not been detected in any bacterial cell (1-3)]. It is assumed that this step is rapid (and therefore not rate-limiting) because intermediates of this step have not been detected. A second deadenylation-dependent pathway involves shortening of the poly(A) tail, followed by 3' → 5' exonucleolytic digestion of the remainder of the mRNA (7, 8). In this pathway, the 5' cap is not obligatorily removed.

Several deadenylation-independent mRNA decay pathways have been characterized. In one pathway, decay is initiated by decapping, followed by 5' → 3' exonucleolytic digestion. This pathway may participate in a “surveillance” mechanism, in which the presence of a premature **termination codon**—produced by [mutation](#) or by a failure to **splice** or **edit**—in conjunction with downstream sequences provides a signal for mRNA decapping and 5' → 3' decay (7, 8). The process of mRNA surveillance is conserved evolutionarily, reflecting its importance in ensuring fidelity of gene expression. A third deadenylation-independent decay pathway involves one or more endonucleolytic cleavages upstream of the poly(A) tail. These events functionally inactivate the mRNA and permit further degradation by a pathway that is not yet well defined. This latter pathway may be specific to mRNA that encodes functionally related proteins, and it may depend on protein recognition of an AU-rich sequence element in the 3'-UTR (7, 8). A fourth pathway that occurs in animal cells involves the activation of a latent ribonuclease. **Interferon** or **viral** infection of animal cells stimulates the production of 2'-5' oligo(A) that activates RNase L, which endonucleolytically cleaves cellular and viral RNA. This inactivates protein biosynthesis, thereby limiting viral replication (7). Whether RNase L participates in normal cell RNA metabolism is unclear. Finally, [histone](#) mRNA, which are non-polyadenylated species, are degraded by a 3' → 5' exonuclease (10).

Most studies on eukaryotic RNA decay have used yeast as the model organism. Attempts to relate fungal RNA decay pathways to those in animal and plant cells have been hampered by the difficulty in performing comparable biochemical genetic analyses with these cells. Nonetheless, the basic steps in RNA decay may be conserved because activities have been characterized in animal and plant cells that shorten the poly(A) tail, remove the 5' cap, or exonucleolytically digest RNA with 5'→3' directionality (10). There is also growing evidence that the Cap and Poly(A) tails interact functionally during translation. A model relating eukaryotic mRNA translation and mRNA decay proposes that when translation initiation is disrupted by a disengagement of the mRNA 5'↔3'-end interaction, the mRNA is shunted into one of the decay pathways (7, 8).

## Bibliography

1. J. G. Belasco and G. Brawerman (1993) *Control of Messenger RNA Stability*, Academic Press, New York.
2. D. P. Nierlich and G. J. Murakawa (1996) *Prog. Nucleic Acids Res. Mol. Biol.* **52**, 153–216.
3. A. W. Nicholson (1997) in *Ribonucleases: Structures and Functions* (J. F. Riordan and G. D'Alessio, eds.), Academic Press, New York, pp. 1–49.
4. A. J. Carpousis, G. Van Houwe, C. Ehretsmann, and H. M. Krisch (1994) *Cell* **76**, 889–900.
5. B. Py, H. Causton, E. A. Mudd, and C. F. Higgins (1994) *Mol. Microbiol.* **14**, 717–729.
6. W-M. Woo and S. Lin-Chao (1997) *J. Biol. Chem.* **272**, 15516–15520.
7. A. Jacobson and S. W. Peltz (1996) *Ann. Rev. Biochem.* **65**, 693–739.
8. C. Beelman and R. Parker (1995) *Cell* **81**, 179–183.
9. T. E. LaGrandeur and R. Parker (1996) *Biochimie* **78**, 1049–1055.
10. J. Ross (1996) *Trends Genet.* **12**, 171–175.

## RNA Editing

RNA editing can best be defined as changing the **nucleotide sequence** of an **RNA** so that of a mature RNA differs from that encoded by the **genomic** sequence. In **eukaryotes**, RNA editing is widespread, occurring in organisms as diverse as **yeast** and humans. Many different classes of RNA, [transfer RNA](#), [ribosome](#) RNA, and [messenger RNA](#), are edited to varying extents, and in many cases the changes introduced by the editing process have been shown to be functionally important. In recent years, much has been learned about both the types of modifications that occur and the factors involved in the modification reaction (for extensive reviews, see refs. [1-3](#)).

There are two broad classes of RNA editing. The first is restricted to site-specific base modifications that give rise to either nucleotide **transitions** or **transversions**. The second is the less well understood insertion/deletion editing that occurs primarily in kinetoplastid protozoan parasites.

### 1. Site-Specific Base Modification Editing

This type of editing includes nucleotide transitions or transversions in all classes of RNA. Many different RNAs are edited in a site-specific manner, and the number of RNAs found to be edited is growing. Several independent mechanisms direct the editing machinery to the site to be edited, depending on the example. The first is a “mooring” sequence in the RNA that is specifically recognized by the editing complex, and editing is directed to this site; this occurs in apolipoprotein B mRNA. The mRNA for ApoB is deaminated at a specific cytosine to a uracil, by a deamination reaction that converts a CAA codon (coding for glutamine) to a UAA [stop codon](#). In this case, editing produces two **isoforms** of the ApoB protein.

The second type of modifications are restricted to double-stranded regions of RNA (see text below). In general most of the edited sites are not edited to completion. One of the best understood examples of this class of RNA editing is in the case of the ionotropic glutamate gated **ion channel** receptors (gluR) in mammals, which involves deamination of adenosine to inosine at specific sites in the pre-mRNA. The functional receptor is encoded by a **gene family** of five receptor subunits—gluR-B, gluR-C, gluR-D, gluR-5, and gluR-6—which can form either homo- or heteropentameric transmembrane channels. Editing has been detected in all five of the glutamate receptor subunits. However, it is the editing events in exons 11 of gluR-B, gluR-5, and gluR-6 that is best double-stranded RNA, which results from base pairing between sequences from exons and introns. This means that editing has to take place before the introns have been removed by **splicing** of the RNA; by implication, it must occur during [transcription](#). The Q/R editing site in exon 11 of gluR-B is found to be edited at frequencies approaching 100%. A CAG codon coding for glutamine is converted to a CIG codon for arginine by deamination of the adenosine residue. The functional consequence of editing at this site is a glutamate receptor that has altered permeability to calcium ions (see [Calcium Signaling](#)). In **transgenic** mice in which the Q/R site of one copy of gluR-B made editing incompetent, the mice developed seizures and died ([5](#)). This demonstrates that the editing event is essential for normal brain function. An enzyme that preferentially catalyzes this reaction has been characterized and cloned: ADAR2 (adenosine deaminase that acts on RNA) belongs to a family of such enzymes that show differing specificities for other editing sites. *In vitro*, ADAR2 efficiently edits the Q/R site, in the absence of cofactors, by deamination of adenosine to inosine ([6](#)). It remains to be determined whether ADAR2 also catalyzes this reaction *in vivo*.

### 2. Insertion/Deletion Editing

This type of editing is found in the [mitochondria](#) of kinetoplastid protozoan parasites. It involves the

precise insertion or deletion of uridylates in the pre-mRNA to generate a functional mRNA. Unlike site-specific base modification editing, which occurs at only a few specific sites, the edited kinetoplastid RNAs are very highly modified. The first case of editing was described by Benne and co-workers (7), who found a precise insertion of four uracil nucleotides in the Cox II mRNA from trypanosome mitochondria that was not encoded by the genomic DNA. Since then many examples of such editing have been described in kinetoplastids. At the present time, the factors involved in this editing are not well characterized, although a mechanistic description of the editing process is becoming apparent.

The sites in the pre-mRNA to be edited are defined by small RNAs that are complementary to edited RNA sequences. These are commonly referred to as [guide RNA](#). The guide RNAs have three domains: (1) the 5' region, which is complementary to the substrate pre-mRNA and acts as an "anchor"; (2) the central domain, which contains the information necessary to insert and/or delete uridylates in the pre-mRNA to make the edited sequence, and is normally around 30–40 nucleotides in length; and (3) the 3' end of the guide RNA, which is characterized by a poly U tail. The editing reaction takes place in a large ribonucleoprotein complex that has been termed the *editosome*. The site to be edited is cleaved immediately 5' to the base-paired guide RNA; Us are either inserted or deleted, and then the RNA halves are religated. This reaction is repeated until the region covered by the guide RNA is completely edited.

### Bibliography

1. H. C. Smith and M. P. Sowden (1996) *Trends Genet.* **12**, 418–424.
2. H. C. Smith, J. M. Gott, and M. R. Hanson (1997) *RNA* **3**, 1105–1123.
3. M. A. Oconnell (1997) *Curr. Biol.* **7**, R437–R439.
4. P. H. Seeburg (1996) *J. Neurochem.* **66**, 1–5.
5. R. Brusa, F. Zimmerman, D. Koh, D. Feldmeyer, P. Gass, P. H. Seeburg, and R. Sprengle (1995) *Science* **270**, 1677–1680.
6. M. A. Oconnell, A. Gerber, and W. Keller (1997) *J. Biol. Chem.* **272**, 473–478.
7. R. Benne, J. Vandenburg, J. P. J. Brakenhoff, P. Sloof, J. H. Vanboom, and M. C. Tromp (1986) *Cell* **46**, 819–826.

### Suggestions for Further Reading

8. B. L. Bass (1997) RNA editing and hypermutation by adenosine deamination, *Trends Biochem. Sci.* **22**, 157–162.
9. R. Schoepfer, H. Moneyer, B. Sommer, W. Wisden, R. Sprengel, T. Kuner, H. Lomeli, A. Herb, M. Kohler, N. Burnashev, W. Gunther, P. Ruppersberg, and P. Seeburg (1994) Molecular-biology of glutamate receptors, *Prog. Neurobiol.* **42**, 353–357.
10. M. L. Kable, S. Heidmann, and K. D. Stuart (1997) RNA editing: getting U into RNA. *Trends Biochem. Sci.* **22**, 162–166.
11. H. Vanderspek, D. Speijer, G. J. Arts, J. Vandenburg, H. Vansteeg, P. Sloof, and R. Benne (1990) RNA editing in transcripts of the mitochondrial genes of the insect trypanosome *Crithidia fasciculata* *EMBO J.* **9**, 257–262.

### RNA Helicases



RNA helicases form a class of [proteins](#) that **catalyze** the unwinding of double-stranded RNA driven by the hydrolysis of a nucleoside triphosphate (NTP). Therefore, RNA helicases are associated with **nucleic acid**-stimulated nucleoside triphosphatase (NTPase) activity. In most cases, the NTPase activity displays a specificity for ATP or dATP, but some of the helicases use any one of the four NTPs or dNTPs equally well ([1](#)).

RNA helicases are distinct from *RNA unwindases*, a class of enzymes that unwind RNA by converting adenine into inosine residues. This modification creates local I-U mismatches and thereby lowers the melting temperature of RNA duplexes in a reaction that does not require  $Mg^{2+}$  and a nucleoside triphosphate ([2](#)).

Most of the RNA helicases known thus far belong to class II of the helicase [superfamily](#) ([3](#)). All class II helicases possess a core domain that contains seven conserved amino acid motifs. A bipartite NTP-binding site, consisting of the residues Gly-Cys-Gly-Lys-Thr (or GCGKT; site A) and Phe-Ile-Leu-Asp-Asp (or FILDD; site B), is found in motif I. a DEXD/H signature (Asp-Glu-X-Asp/His) shows up in motif II, and a putative RNA-binding site XRXGRXXR (or X-Arg-X-Gly-Arg-X-X-Arg) has been assigned to motif IV of some RNA helicases, such as the eukaryotic translation [initiation factor](#) (eIF-) 4A ([4](#)) and the plum pox potyvirus CI helicase ([5](#)). The DEXD/H signature is found in many RNA helicases, with eIF-4A being the prototype of the so-called **DEAD box** family, whereas the motif DEXH is found only within a distinct subfamily of RNA helicases ([6](#)). Numerous open reading frames from both prokaryotes and eukaryotes contain DEXD/H motifs and on this basis have been assigned helicase functions. The existence of these motifs in a protein sequence is suggestive of a helicase function, but as long as NTP-dependent unwinding has not been demonstrated experimentally, such an assignment should be regarded as preliminary and tentative.

RNA helicases unwind two complementary RNA strands to yield single-stranded products. The time-course of unwinding can be followed by subjecting the reaction products to [gel electrophoresis](#), where the double-stranded RNA substrate migrates more slowly than the single-stranded products. In principle, there are two ways to unwind RNA: the helicase might move 3'–5' with respect to the single-stranded RNA to which it has bound, or it might move in the opposite direction (ie, 5'–3'). Correspondingly, the directionality of unwinding can be determined by offering single-stranded RNA substrates with duplexes at either end. Only eIF-4A in conjunction with eIF-4B does not display a preferred directionality on unwinding of the secondary structure at the 5'-end of mRNA ([7](#)).

Some of the canonical RNA helicases also unwind double-stranded DNA (see [DNA Helicase](#)) once they have bound to a single-stranded RNA entry site. For example, the *Escherichia coli* **Rho** transcription [termination factor](#) induces the release of an RNA transcript from its template DNA, when the factor has migrated to the RNA-intrinsic rho-dependent terminator sequence ([8](#)). Some DEXD/H-containing RNA helicases even unwind DNA/DNA duplexes, yielding a class of enzymes that unwind both RNA and DNA. Members of this type of helicase are the nonstructural protein 3 of hepatitis C virus ([9](#)), the nucleic acid-dependent phosphohydrolase II of vaccinia virus ([10](#)), and the “maleless” (MLE) protein of *Drosophila* ([11](#)). Also, the **SV40 virus** large [T Antigen](#), which is not a member of the DEXD/H-family of helicases, unwinds both RNA and DNA ([12](#)).

It is thought that RNA helicases first bind to a single-stranded template and then undergo a conformational change that is induced by NTP binding and hydrolysis. The NTP-driven conformational change is an example of an **energy transduction** process that transduces chemical energy into mechanical work that is used for translocation of the helicase and disruption of the RNA duplex. The [protein structure](#) of the NS3 RNA helicase of hepatitis C virus might give an example of the coupling between NTP hydrolysis and RNA unwinding ([13](#)). In addition to the **RNA-binding** site of the catalytic core, many RNA helicases of the DEXD/H type contain auxiliary binding sites for nucleic acids that contribute to the specificity and/or affinity of substrate binding but may not be subjected to conformational changes driven by NTP hydrolysis ([14](#)). The presence of auxiliary

binding sites complicates currently discussed unwinding models, where nucleotide binding and hydrolysis coordinate alternating DNA binding affinities during translocation and unwinding (15). In the case of the Rho helicase, the auxiliary RNA binding sites are thought to contribute to a “tethered tracking” mechanism, in which Rho forms tight primary binding interactions with the recognition region of the RNA and remains bound there, while transient secondary RNA binding interactions coupled to ATP hydrolysis serve to scan along the RNA to contact the DNA/RNA duplex (16).

Functional aspects of RNA helicases have been addressed mainly by studying **conditional mutants** of the respective **genes**, both in yeast and in bacteria. In that way, it has been shown that RNA helicases are involved in numerous aspects of RNA metabolism, such as [RNA splicing](#), assembly of [ribosomes](#), and initiation of [translation](#) (17). A manifold of RNA helicases exist in yeast, with more or less defined functions. A group of RNA helicases, known as the pre-mRNA processing (PRP) factors prp2, prp5, prp16, prp28, and prp43, participate in consecutive steps of pre-mRNA splicing. PRP proteins are either tightly associated or transiently interacting with the [spliceosome](#) (18). Among these, prp2, prp16, prp22, and prp43 contain a DEAH signature in motif II, whereas prp5 and prp28 display a DEAD signature in the same motif. The PRP proteins facilitate interactions between [small nuclear RNPs](#) and pre-mRNA, and they are involved in splice site recognition, spliceosome activation, mRNA release, and spliceosome disassembly. The well-known ATP dependence of spliceosome-catalyzed mRNA processing is thought to involve RNA helicase functions that may also contribute to the accuracy of the splicing reaction. [Ribosome](#) biogenesis is another process that requires RNA helicases, with possible functions in tRNA transcription (19), rRNA maturation (20), and ribosome assembly (21). RNA helicases may assist endo- and exonucleolytic pre-rRNA processing steps, they may regulate interactions between small nucleolar RNA (snoRNA) and pre-rRNA, and they may rearrange rRNA-protein interactions during ribosome assembly.

Some RNA helicases of the DEXD box family promote the **nuclear export** of RNA. The RNA helicase Dbp5 from yeast (22) has been shown to be a **nuclear pore**-associated protein that is essential for mRNA export (23). Moreover, **homologues** of Dbp5 helicase seem to be present in most eukaryotic organisms (23). Also, the human homologue of the Drosophila MLE helicase is thought to be involved in the nuclear export of unspliced type D RNA of [retroviruses](#) (24).

The apparent homology between many yeast and mammalian RNA helicases can be exploited for rescuing functional yeast mutants by [complementation](#) with the mammalian counterpart. This might determine the suspected physiologic role of the mammalian homologue, once the function of the yeast enzyme has been defined. Due to the lack of suitable genetic approaches, however, the physiologic function of many mammalian RNA helicases is still unknown. Nevertheless, RNA helicases of metazoan origin are thought to be involved in embryogenesis, differentiation, and spermatogenesis, as well as in cellular growth and division (17). The MLE helicase from Drosophila is a well-characterized RNA helicase that is involved in the hypertranscription of the single X chromosome of the male gender. MLE-promoted hypertranscription equalizes the amount of transcripts obtained from the single X chromosome of males with that obtained from both X chromosomes of females (25). MLE is highly conserved between [nematodes](#), dipterans, and mammals, but it is missing from yeast. A less complex form of MLE has been found in the plant [Arabidopsis](#) by comparisons with protein [sequence databases](#), although this enzyme has different N- and C-termini compared to the conserved nematode to mammalian counterparts (26). Therefore, computer-assisted comparisons of the sequences of various RNA helicases may give further hints to the structures and functions of distinct members of this important class of enzymes.

## Bibliography

1. S. W. Matson and K. A. Kaiser-Rogers (1990) *Annu. Rev. Biochem.* **59**, 289–329.
2. K. Nishikura (1992) *Ann. N. Y. Acad. Sci.* **28**, 240–250.
3. A. E. Gorbalenya, E. V. Koonin, A. P. Donchenko, and V. M. Blinov (1989) *Nucleic Acids Res.* **17**, 4713–4730.

4. A. Pause, N. Methot, and N. Sonenberg (1993) *Mol. Cell. Biol.* **13**, 6789–6798.
5. A. Fernandez, S. Lain, and J. A. Garcia (1995) *Nucleic Acids Res.* **23**, 1327–1332.
6. D. A. Wassarman and J. A. Steitz (1991) *Nature* **349**, 463–464.
7. F. Rozen, I. Edery, K. Meerovitch, T. E. Dever, W. C. Merrick, and N. Sonenberg (1990) *Mol. Cell. Biol.* **10**, 1134–1144.
8. C. A. Brennan, A. J. Dombroski, and T. Platt (1987) *Cell* **48**, 945–952.
9. C. L. Tai, W. K. Chi, D. S. Chen, and L. H. Hwang (1996) *J. Virol.* **70**, 8477–8484.
10. C. D. Bayliss and G. L. Smith (1996) *J. Virol.* **70**, 794–800.
11. C.-G. Lee, K. A. Chang, M. I. Kuroda, and J. Hurwitz (1997) *EMBO J.* **16**, 2671–2681.
12. M. Scheffner, R. Knippers, and H. Stahl (1989) *Cell* **57**, 955–963.
13. J. L. Kim, K. A. Morgenstern, J. P. Griffith, M. D. Dwyer, J. A. Thomson, M. A. Murcko, C. Lin, and P. R. Caron (1998) *Structure* **6**, 89–100.
14. S. Zhang and F. Grosse (1997) *J. Biol. Chem.* **272**, 11487–11494.
15. T. M. Lohman, K. Thorn, and R. D. Vale (1998) *Cell* **93**, 9–12.
16. E. J. Steinmetz and T. Platt (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1401–1405.
17. S. R. Schmid and P. Linder (1992) *Mol. Microbiol.* **6**, 283–291.
18. J. P. Staley and C. Guthrie (1998) *Cell* **92**, 315–326.
19. R. Gururajan, L. Mathews, F. J. Longo, and D. L. Weeks (1994) *Proc. Natl. Acad. Sci. USA* **15**, 2056–2060.
20. P. L. Weaver, C. Sun, and T. H. Chang (1997) *Mol. Cell. Biol.* **17**, 1354–1365.
21. D. Kressler, J. de la Cruz, M. Rojo, and P. Linder (1998) *Mol. Cell. Biol.* **18**, 1855–1865.
22. S. S. I. Tseng, P. L. Weaver, Y. Liu, M. Hitomi, A. M. Tartakoff, and T. H. Chang (1998) *EMBO J.* **17**, 2651–2662.
23. C. A. Snay-Hodge, H. V. Colot, A. L. Goldstein, and C. N. Cole (1998) *EMBO J.* **17**, 2663–2676.
24. H. Tang, G. M. Gaietta, W. H. Fischer, M. H. Ellisman, and F. F. Wong-Staal (1997) *Science* **276**, 1412–1415.
25. M. I. Kuroda, M. J. Kernan, R. Kreber, B. Ganetzky, and B. S. Baker (1991) *Cell* **66**, 935–947.
26. W. Wei, D. Twell, and K. Lindsey (1997) *Plant J.* **11**, 1307–1314.

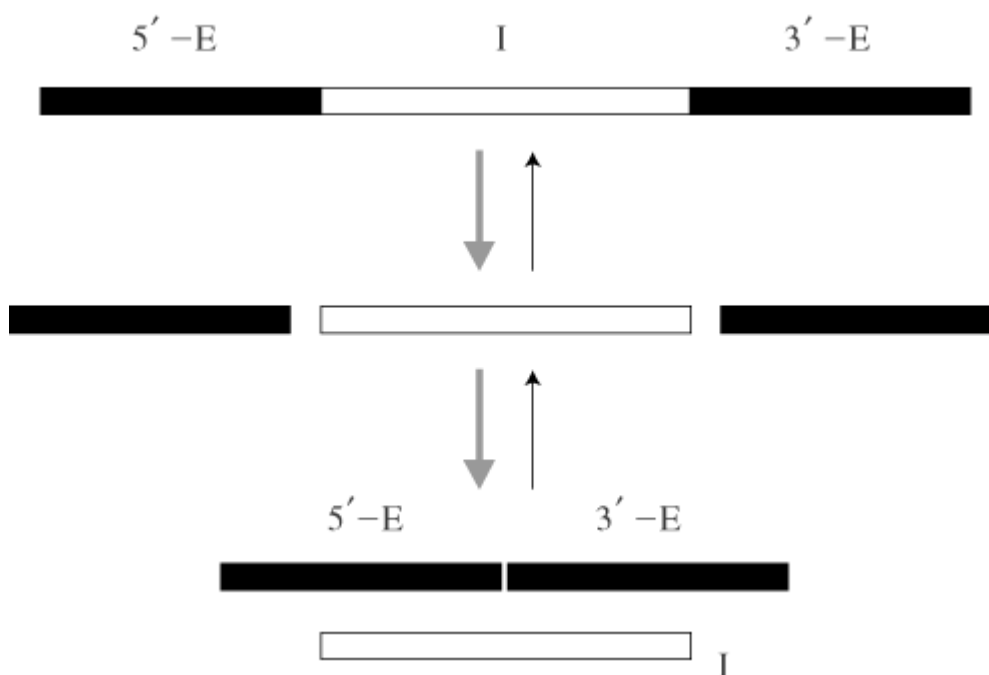
### **Suggestions for Further Reading**

27. J. S. Anderson and R. Parker (1997) RNA turnover: the helicase story unwinds. *Curr. Biol.* **6**, 780–782.
28. L. E. Bird, H. S. Subramanya, and D. B. Wigley (1998) Helicases: a unifying structural theme? *Curr. Opin. Struct. Biol.* **8**, 14–18.
29. G. Kadare and A. L. Haenni (1997) Virus-encoded RNA helicases. *J. Virol.* **71**, 2583–2590.
30. E. V. Koonin (1992) A new group of putative RNA helicases. *Trends Biochem. Sci.* **17**, 495–497.
31. A. Kramer (1996) The structure and function of proteins involved in mammalian pre-mRNA splicing. *Annu. Rev. Biochem.* **65**, 367–409.
32. K. J. Marians (1997) Helicase structures: a new twist on DNA unwinding. *Structure* **5**, 1129–1134.

### **RNA Ligases**

RNA ligase is an [enzyme](#) that can catalyze the ligation of RNA molecules via phosphodiester bonds. Several species of RNA ligases are currently known from all three kingdoms of organisms: bacteria, archaeobacteria, and eukaryotes. The majority of them are involved in [RNA splicing](#), splicing [transfer RNA](#). Another important function of RNA ligation is in [RNA Editing](#). The basic reaction steps related to RNA ligation are shown schematically in Figure 1, which indicates that every RNA ligation follows one or more cleavages of the precursor RNA. Therefore, we can categorize ligation reactions into three types, depending upon whether they are directed by an intron, exon (see [Introns](#), [Exons](#)), or [guide RNA](#) (gRNA) and also depending upon what connects the two RNA segments that are to be ligated together (Table 1). Intron-directed means that most of the information necessary for determining the cleavage sites and the two RNA molecules to be ligated is conveyed by an intron moiety, the portion between the two segments to be ligated, which will be removed. This is very reasonable in splicing of [messenger RNA](#), because each RNA molecule to be ligated is an exon, which should retain its own information other than that for specifying the cleavage sites. This ligation reaction is usually carried out by protein catalysts. Conversely, exon-directed indicates that the exon parts, which are to be joined, are responsible for those requirements. In this case, an elaborate construct must have been exploited in order to reconcile the exon-specific information with the information specifying the cleavage sites (see [tRNA Ligase](#)). On the other hand, gRNA-directed means that the ligation reaction requires an additional RNA element, the guide RNA. Currently, the origin of gRNA is far more mysterious than that of introns. The functions of some RNA ligases are not known (Table 1).

**Figure 1.** The basic steps involved in RNA ligation. All of the phenomena, splicing (mRNA and tRNA) and editing (insertion and deletion), can be reduced to the forward and reverse steps shown in this figure, so long as the details are not considered; only insertional editing involves the reverse steps. An RNA molecule is divided into three parts that are designated (from left to right) as 5'-E (exon), I (intron), and 3'-E. The overall reaction considered here consists of two cleavages and one ligation (or one cleavage and two ligations, in the case of a single event of insertional editing). It involves the recognizing of the sites for cleavage and the joining of segments 5'-E and 3'-E.



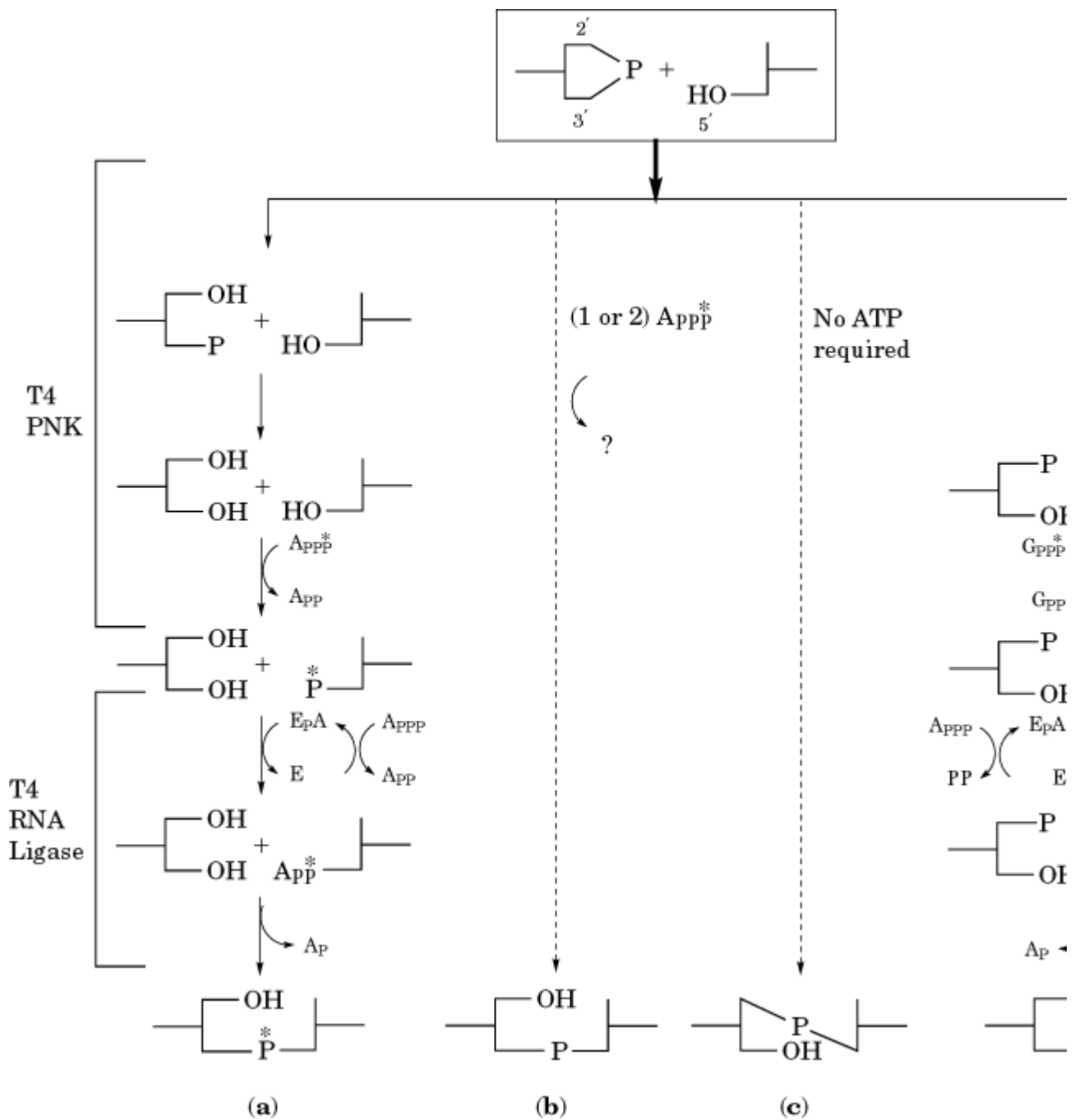
**Table 1. Categories of RNA Ligation**

| Name of phenomena or key enzymes                | Function | Pairing information <sup>a</sup> | Reactions                                   |
|---|----------|----------------------------------|---|
| Self-splicing                                   |          | Intron-directed                  | Transesterification (Non-protein catalysis) |
| Spliceosome                                     | Splicing |                                  |   |
| tRNA Splicing                                   |          | Exon-directed                    | Cleavage and ligation                       |
| Insertional/Deletional editing                  | Editing  | gRNA-directed                    | Cleavage, insertion/deletion, and ligation  |
| T4 RNA ligase<br><i>E. coli</i> 2'-5'RNA ligase | Unknown  | Unknown                          | (Unknown)+ligation                          |

<sup>a</sup> Information necessary for determining the two fragments to be ligated and bringing them neighbored.

The reaction mechanisms of RNA ligases are summarized in Figure 2. In general, the RNA-joining reaction begins with two kinds of RNA molecules, a 2',3'-cyclic phosphate acceptor and a 5'-hydroxyl donor. Eukaryotic tRNA ligase has a whole set of activities for this mechanism: (i) opening of the 2',3'-cyclic phosphate to a 2'-phosphomonoester (phosphodiesterase activity), (ii) addition of a phosphate group to the 5'-hydroxyl of the donor (kinase activity), (iii) transfer of an adenylyl group from an adenylylated enzyme, which is produced in advance by another kinase activity of the same enzyme, to the 5'-phosphate of the donor (adenylylation activity), and (iv) attack of the acceptor 3'-hydroxyl on the activated donor phosphoanhydride, to form a 3',5'-phosphodiester bond (bond formation activity). The mature tRNA is then generated by the action of a distinct enzyme, 2'-phosphotransferase, which removes the 2'-phosphate from the product. T4 RNA ligase accomplishes the same reaction with steps iii and iv, although the donor for T4 RNA ligase must have a 2'-hydroxyl structure (not 2'-phosphate). Therefore, tRNA ligase has the combined functions of T4 **polynucleotide kinase** and T4 RNA ligase, as indicated in Figure 2.

**Figure 2.** Reaction mechanisms of RNA ligases. (a) T4 RNA ligase; (b) animal pathway RNA ligase (ligases that resemble exclusively in the cells of vertebrates); (c) 2',5' RNA ligase; (d) tRNA ligase. The 3' end of the left-hand RNA segment ribose indicated; the 5' end of the right-hand RNA segment has the 5' ribose position indicated. A<sub>PPP</sub>:ATP; G<sub>PPP</sub>:GTP; P bridging phosphate in the product (bottom) is not the starting phosphate incorporated in the 2',3'-cyclic phosphodiester bond. The source of the ATP is indicated similarly.



RNA editing enzymes in kinetoplasts of *Trypanosoma brucei* and *Leishmania tarantolae* are also assumed to have a reaction mechanism similar to that of T4 RNA ligase, although a transesterification mechanism that does not require a ligase for joining of the RNAs cannot be completely ruled out (1). *Escherichia coli* 2',5' RNA ligase and animal pathway ligase have been reported, but their functions and detailed reaction mechanisms are not well understood (Fig. 2). Recent [genome](#) sequencing analyses made it clear that the archaeon *Methanococcus jannaschii* and the **gram-positive bacterium** *Bacillus stearothermophilus* each have a homologue of the *E. coli* 2',5' RNA ligase, but the bacterium *Haemophilus influenzae* does not. Interestingly, the archaeon does not contain homologues of T4 RNA ligase or tRNA ligase, although it does retain the homologue of eukaryotic tRNA splicing endonuclease, supporting the suggestion that 2',5' RNA ligase may be a relevant ligase for tRNA splicing in those organisms. The eukaryotic and archaeal endonucleases are very similar in substrate recognition and protein structures (2, 3).

The common reaction steps of enzyme adenylation (or guanylation) and nucleotidyl transfer are shared among RNA ligases (T4 RNA ligase and tRNA ligases), DNA ligases (those of **T4 phage**, *Saccharomyces cerevisiae*, and *Homo sapiens*), and mRNA capping enzymes. In all three, a [lysine](#) residue at the [active site](#) of the enzymes, which is highly conserved, is transiently adenylylated (or guanylylated in the case of mRNA capping) using ATP (or GTP), and then the nucleotide is transferred to the 5'-phosphate group (5'-pyrophosphate in the case of mRNA capping) of a donor molecule. The activated diphosphoanhydride form facilitates ligation of both DNA and RNA. In the case of mRNA capping, the same reaction occurs without a further reaction, generating G(5') pppN<sub>1</sub>N<sub>2</sub><sup>1/4</sup>. Being similar in function, structure, and reaction mechanism, these enzymes are believed to have had a common progenitor.

Ligation might be performed by the reverse of the reactions of endonucleases like **ribonuclease A**, **RNase T<sub>1</sub>**, and RNase PH. These ribonucleases usually cleave RNA to produce a transient 2',3'-cyclic phosphate and a 5'-hydroxyl at the newly exposed ends, so the substrate structures are the same as those used in authentic RNA ligation (although it would not require that the reaction mechanism be the same in detail). There is, however, no evidence that such a reverse ribonuclease reaction occurs under physiological conditions.

**Ribozymes** (RNA enzymes) capable of RNA ligation have been produced using the technique of *in vitro* selection (4, 5). Currently, two types of ligation reactions are used, mimicking the reactions of (i) RNA polymerase [attack of the 3'-hydroxyl on the 5'-triphosphate (4)] or (ii) RNA ligase [attack of the 3'-hydroxyl on the 5'-5'- diphosphate (5)]. These studies may lead to elucidation of the mechanism of **self-splicing** and development of highly useful ligation techniques.

## 1. T4 RNA Ligase

T4 RNA ligase is encoded by gene 63 of T4 phage, has a molecular weight of 43.5 kDa, and is composed of 375 amino acid residues of known [primary structure](#) (6). Originally, T4 RNA ligase was discovered as an enzyme that intramolecularly ligates the terminal 3'-hydroxyl of a polyriboadenylate with its own 5'-phosphate, generating a circular RNA. Then, this enzyme was found to have a more general function, to ligate the 3'-hydroxyl of a single-stranded RNA/DNA with the 5'-phosphate of the same or another single-stranded RNA/DNA, although the ligation of RNA molecules is 200-fold more efficient than is that of DNA. The biological role of T4 RNA ligase is not known, though it has a distinct function, carried out by a different part of the enzyme, in **T4 phage** tail-fiber attachment as part of phage morphogenesis. An intriguing hypothesis is based on the observation that one strain of *E. coli* (CTr5x) that contains an [anticodon](#) nuclease (ie, a nuclease targeted at the anticodon site of the host tRNA<sup>Lys</sup>) is resistant to T4 mutants that carry a defect in either RNA ligase or **polynucleotide kinase** (PNK) activity. The interpretation of this is that these T4 mutants cannot proliferate in *E. coli* cells that have an anticodon nuclease activity induced by T4 stp protein (part of the host **restriction/modification** system) due to their inability to repair the resulting breakage of the tRNA<sup>Lys</sup> which is now essential for the T4 growth. To restore the tRNA, the following serial reactions are required: 2',3'-cyclic phosphodiesterase, 3'-phosphatase, and the 5'-kinase activities of T4 PNK and the ligation activity of T4 RNA ligase (see Fig. 2). This hypothesis, which considers T4 RNA ligase to be a kind of tRNA ligase, also implies that T4 RNA ligase may have arisen from an ancestral tRNA splicing enzyme.

The enzymatic properties of T4 RNA ligase have been investigated thoroughly, and the reaction mechanism is shown in Figure 2. For ligation, first, the **ε-amino group** of Lys99 of T4 RNA ligase is adenylylated, using one ATP molecule. Second, a donor RNA/DNA properly located on the enzyme accepts the adenylyl group on its 5'-terminal phosphate, generating a 5',5'-phosphoanhydride bond. Finally, the 3'-hydroxyl of an acceptor attacks nucleophilically the activated phosphoanhydride of the adenylylated donor, forming a 3',5'-phosphodiester and releasing

AMP. Asp101 of this enzyme is known to be important in the second and the final steps of ligation. Short-circuiting of the first and the second steps by use of a previously 5'-adenylylated donor is possible and enhances the overall reactions, indicating those steps to be rate-limiting.

Because T4 RNA ligase is very permissive in its substrate recognition, various techniques that use this enzyme have been developed: (i) enzymic ligation of oligoribo/oligodeoxyribonucleotides (7), (ii) 3'- or 5'-end labeling of RNA/DNA with **biotin**, **fluorescent dyes**, or **radioisotopes** (8), (iii) introduction of nicotinamide nucleotides into RNA (9), (iv) end-closing of DNA/RNA (10), and so forth.

### Bibliography

1. J. D. Alfonzo, O. Thiemann, and L. Simpson (1997) *Nucleic Acids Res.* **25**(19), 3751–3759.
2. H. Li, C. R. Trotta, and J. Abelson (1998) *Science* **280**(10), 279–284.
3. S. Fabbri et al. (1998) *Science* **280**(10), 284–286.
4. D. P. Bartel and J. W. Szostak (1993) *Science* **261**, 1411–1418.
5. A. J. Hager and J. W. Szostak (1997) *Chem. Biol.* **4**(8), 607–617.
6. K. Rand and M. J. Gait (1984) *EMBO J.* **3**, 397–402.
7. X.-H. Zhang and V. L. Chiang (1996) *Nucleic Acids Res.* **24**(5), 990–991.
8. Y. Kinoshita, K. Nishigaki, and Y. Husimi (1997) *Nucleic Acids Res.* **25**(18), 3747–3748.
9. K. Harada and L. E. Orgel (1993) *Proc. Natl. Acad. Sci. USA* **90**, 1576–1579.
10. D. Beaudry and J-P. Perreault (1995) *Nucleic Acids Res.* **23**(15), 3064–3066.

### Suggestions for Further Reading

11. O. Uhlenbeck and R. I. Gumport (1982) "T4 RNA Ligase". in *The Enzymes* (P. Boyer, ed.), Academic Press, New York, pp. 31–60.
12. M. L. Kable, S. Heidmann, and K. D. Stuart (1997) RNA editing: getting U into RNA. *Trends Biochem. Sci.* **22**(5), 162–166.
13. E. A. Arn and J. Abelson (1998) "RNA Ligases: Function, Mechanism, and Sequence Conservation". In *RNA Structure and Function* (R. W. Simons and M. Grunberg-Manago, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 695–726.

## RNA Polymerases, DNA-Dependent

DNA-dependent RNA polymerases (RNAPs) are essential components of **gene** expression in organisms that use DNA to encode their genetic information. These **enzymes** catalyze the formation of phosphodiester bonds between ribonucleotides to form a single-strand RNA polymer using DNA as a template (see **Transcription**). DNA-dependent RNAPs thus produce an RNA copy of a gene that may, after further modification, function directly, in the case of **transfer RNA**, **ribosomal RNA**, or **small nuclear RNA** (snRNA), or undergo further decoding by the protein **translation** machinery, in the case of **messenger RNA**. DNA-dependent RNAPs can be divided into two distinct **homology** classes, multi- and single-subunit RNAPs. The multisubunit RNAPs are found throughout the living world; they function as the main cellular RNAPs in eubacteria, archaea, and eukaryotes. Some virally-encoded RNAPs are also members of the multisubunit RNAP group. The single-subunit RNAPs (ssRNAPs) were first discovered in **bacteriophages T7** and T3, but related members of this group are also found in eukaryotic **mitochondria** and **chloroplasts**.

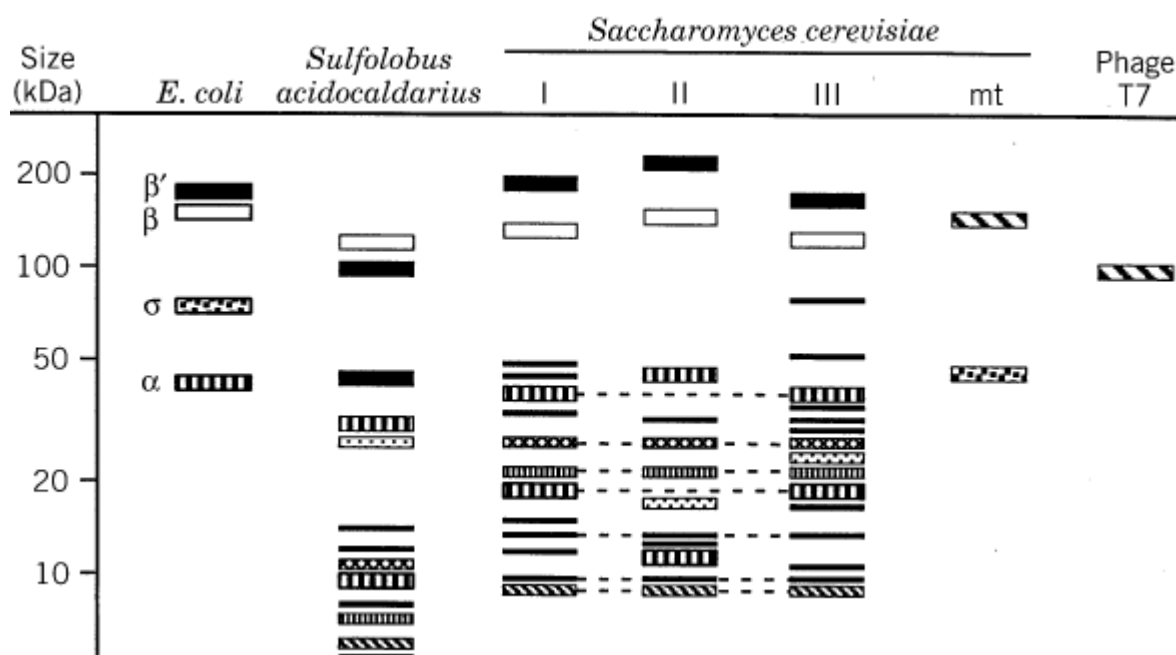


The evolutionary relationship between the two classes of DNA-dependent RNAPs is not clear, because there is no obvious sequence similarity between the single-subunit RNAP and the subunits of the multicomponent cellular RNAP. As described below, however, the two classes of RNAPs have several mechanistic similarities. Both classes also undergo a similar transcription cycle involving (a) DNA **promoter** recognition, (b) opening of the promoter DNA, (c) initiation of transcription, (d) a conformational change releasing the RNAP-promoter interaction and triggering a processive mode of transcript elongation, and (e) termination and transcript release.

## 1. Multisubunit RNAPs

The multisubunit RNAPs are of ancient origin. As shown in Figure 1, the cellular RNAPs from the three domains of life are clearly related, suggesting that the last common ancestor of the three lineages contained a multisubunit RNAP (1, 2). Additional multisubunit RNAPs function within chloroplasts (3) and in certain cytoplasmic DNA **viruses** (4, 5). The multisubunit chloroplast RNAP is very similar to the eubacterial enzyme, consistent with the endosymbiotic origin of the organelle, whereas the viral RNAPs are most closely related to eukaryotic nuclear RNAPs.

**Figure 1.** Relationships between the subunits of multisubunit and of selected bacterial, archaeal, eukaryotic nuclear and mitochondrial, and bacteriophage DNA-dependent RNA polymerases. This figure represents a schematic **SDS-PAGE** separating the subunits of the indicated RNAPs according to their masses. Amino acid sequence homologies are indicated by similar shading of the thicker bars used to denote some of the subunits. Dotted lines connect subunits of the eukaryotic nuclear RNAPs encoded by the same genes. In the following list, the parenthesis include the mass in kDa and in some cases the name of the encoding gene. *Escherichia coli* RNAP is composed of  $\beta'$  (160, *rpoC*),  $\beta$  (155, *rpoB*),  $\alpha$  (40, *rpoA*), and one of several  $\sigma$  factors; the major vegetative factor (70, *rpoD*) is shown here. The archaeal RNAP from *S. acidocaldarius* is composed of B (122), A' (101), A'' (44), D (30), E (27), G (13.8), F (12), H (11.8), L (10), I (9.7), K (9.7), N (7.5), and M (5.5) subunits. The *S. cerevisiae* nuclear RNAP I is composed of (190, *RPA190*), (135, *RPA135*), (49, *RPA49*), (43, *RPA43*), (40, *RPC40*), (34, *RPA34*), (27, *RPB5*), (23, *RPB6*), (19, *RPC19*), (14, *RPA14*), (14, *RPB8*), (12, *RPA12*), (10, *RPC10*), and (10, *RPB10*) subunits. RNAP II is composed of (220, *RPB1*), (150, *RPB2*), (45, *RPB3*), (32, *RPB4*), (27, *RPB5*), (23, *RPB6*), (17, *RPB7*), (14, *RPB8*), (13, *RPB9*), (13, *RPB11*), (10, *RPC10*), and (10, *RPB10*) subunits. RNAP III is composed of (160, *RPC160*), (128, *RPC128*), (82, *RPC82*), (53, *RPC53*), (40, *RPC40*), (37), (34, *RPC34*), (31, *RPC31*), (27, *RPB5*), (25, *RPC25*), (23, *RPB6*), (19, *RPC19*), (17), (14, *RPB8*), (11), (10, *RPC10*), and (10, *RPB10*) subunits. The *S. cerevisiae* mitochondrial RNAP (mt) is composed of (145, *RPO41*) and (43, *MTF1*) subunits. The T7 RNAP is a single polypeptide chain of 100 kDa. The data used to create this figure came from Refs. 2, 11, 24, 38, and 63.



The multisubunit RNAPs include a core enzyme that is fully capable of RNA chain elongation and termination. Most of these core enzymes can initiate transcription *in vitro* from single-stranded DNA or, in some cases, from model copolymer templates like poly dAT. However, the core RNAPs are incapable of initiating transcription selectively at a specific promoter sequence on a natural DNA template. The multisubunit RNAPs therefore rely on associated factors for promoter recognition and, in some cases, for the initial steps of transcription. The required initiation factors in eubacteria and eukaryotes differ greatly. In eubacteria, specific [sigma factors](#) associate with the core RNAP and confer promoter-specific transcription initiation properties to the otherwise nonspecific enzyme. While not similar at the level of [primary structure](#) or three-dimensional [protein structure](#) to sigma factors, [TATA box](#) binding protein (TBP) plays an analogous role in both archaea and eukaryotes. TBP binds the TATA box sequence in promoter DNA and, with the aid of several other [transcription factors](#), recruits the RNAP to promoter regions of the DNA.

As shown in [Figure 1](#), the multisubunit RNAPs of all three kingdoms are clearly related, but the eukaryotic and archaeal RNAPs are more closely related to each other, in terms of both subunit composition and similarity of the individual subunits, than they are to the eubacterial RNAP. The eubacterial core RNAP is composed of only three subunits,  $b'$ ,  $b$ , and  $a$ . The eukaryotic and archaeal RNAPs include homologues of the eubacterial core polypeptides, but they also have many additional subunits. The eubacterial RNAP therefore represents a minimal multisubunit core RNAP that functions similarly to the larger RNAPs in the basic mechanisms of initiation, elongation, and termination of transcription. Indeed, much of what we know about the multisubunit RNAPs is derived from studies of *Escherichia coli* RNAP.

### 1.1. Eubacterial RNAP

Early studies of the bacterial RNAP revealed that it is a multisubunit enzyme found in two forms. The core RNAP consists of three subunits  $a$ ,  $b'$ , and  $b$ , with a stoichiometry of 2:1:1 ( $a_2b'b$ ), while the holoenzyme also includes a  $\sigma$  sigma factor (see [Fig. 1](#)). Core RNAP functions in elongation and termination, while the holoenzyme is responsible for initiating transcription. Although differences occur between the promoters and initiation factors of multisubunit RNAPs, the basic steps of transcription are similar for all RNAPs ([6](#), [7](#)) (see also [Transcription](#) and [Sigma Factors](#)). They include, but are not limited to, the following six steps.

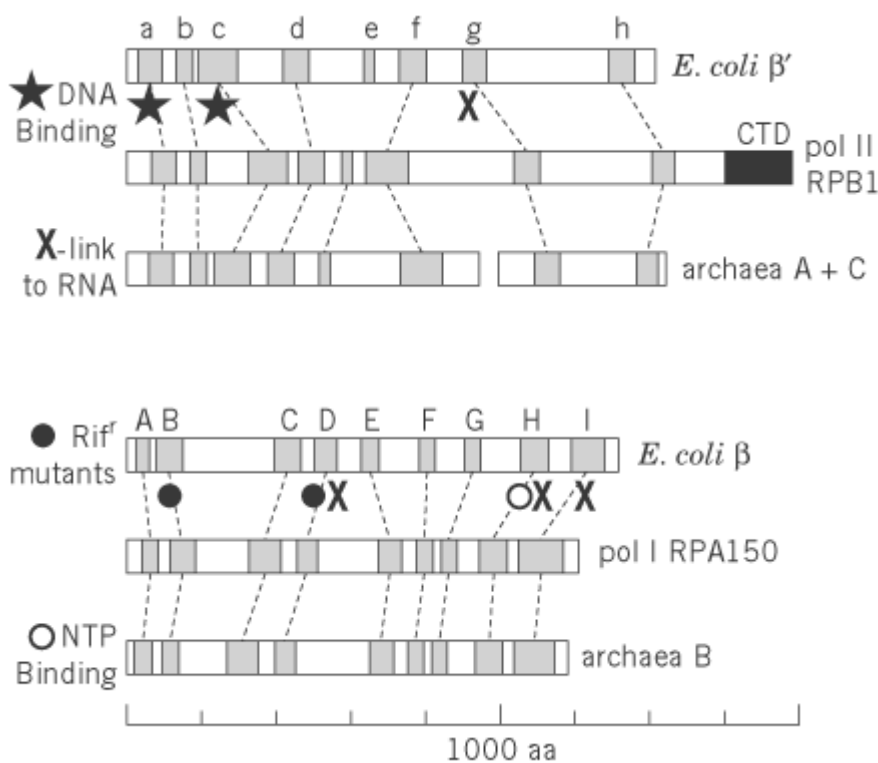
1. The holoenzyme binds to the promoter DNA.
2. The promoter is isomerized into an open complex in which the DNA strands are separated.
3. Transcription is initiated in a poorly processive process where many abortive short transcripts are released while the RNAP maintains contact with the promoter.
4. The RNAP undergoes a significant change called *promoter escape*, which occurs after the transcript reaches a length of about 8 to 12 nucleotides. Escape is marked by the release of the  $\sigma$  factor and by a transition to a processive form of transcription.
5. The transcript is elongated.
6. The transcript is released in a process called termination.

[Footprinting](#) of *E. coli* RNAP/promoter DNA complexes revealed that the RNAP initially protects a large region of DNA spanning from  $-55$  to  $+20$  relative to the start site of transcription ([7](#), [8](#)). During the isomerization and early initiation steps, the upstream border of the footprint remains fixed, while the DNA is opened to form a 18-bp “bubble” within the protected region. The downstream border of the footprint also remains fixed until the RNAP escapes from the promoter. At that point, the downstream border of the footprint jumps forward, and the footprint decreases in size coincident with the release of the sigma factor. The open DNA bubble is maintained through contacts with the RNAP and through the formation of an RNA/DNA hybrid 8 to 12 bp long. Destabilization of the RNA/DNA hybrid is associated with termination of transcription.

The structures of the *E. coli* core and holoenzyme have been analyzed by [X-ray crystallography](#) methods (9, 10). Although the reported resolution is currently low, these initial studies reveal a fundamental difference in the structures of the two forms of the RNAP. Both structures contain a channel of a size that could accommodate binding of the template DNA. The channel of the core RNAP appears to be closed, while the channel of the holoenzyme is open and potentially accessible to DNA, consistent with its role in initiation. Furthermore, the closed form of the channel in the core RNAP leads to a plausible mechanism for the processivity of the core enzyme. Closing of the open channel when  $\sigma$  dissociates may “clamp” the RNAP into a more stable elongating form.

The two largest subunits of the *E. coli* RNAP ( $\beta'$  and  $\beta$ ) form the catalytic core of the enzyme. These subunits have been shown to function in DNA binding and transcription elongation and termination (8).  $\beta'$  and  $\beta$  are homologous to the largest subunits of all multisubunit RNAPs (see Fig. 1). As shown in Figure 2,  $\beta'$ ,  $\beta$ , and their homologues contain eight (labeled **a–h**) and nine (labeled **A–I**) collinear segments, respectively, that have strong amino acid sequence similarities (11). Mutational analyses support the idea that these conserved segments are critical for the function of RNAP; many mutations affecting elongation and termination map to these universally conserved regions (8).

**Figure 2.** Primary structure relationships of the largest subunits of the bacterial, eukaryotic, and archaeal RNA polymerases. Similarities between the *E. coli*  $\beta'$  subunit and its eukaryotic and archaeal orthologues are shaded and labeled **a–h**. Similarities between the *E. coli*  $\beta$  subunit and its orthologues are labeled **A–I**. pol I and pol II are the yeast RNAP I and RNAP II, respectively. Regions of known function for the *E. coli* subunits are indicated by symbols below the shaded regions. The data used to create this figure came from Refs. 11, 13, and (14).



$\beta'$  contains a Cys<sub>4</sub>- zinc finger-like motif near the *N*-terminus within conserved region **a** (Fig. 2) that is thought to function in **DNA-binding**. This motif is conserved in the corresponding subunits of all multisubunit RNAPs. Additionally, within conserved region **c** there is a small region of similarity between  $\beta'$  (and its homologues) and a DNA-binding region of **DNA polymerase I**. In confirmation

of the role of this region in DNA binding, mutations there destabilize the binding of the RNAP to promoter DNA during initiation (12). Finally, the 3' end of the nascent RNA transcript is capable of [crosslinking](#) to conserved region **g**, strongly suggesting that this region of b' may form part of the active site of the RNAP (8).

The b subunit is responsible for binding the substrate nucleoside triphosphates (Fig. 2) (13). The b subunit also contains three regions (**D**, **H**, **I**) that are in close proximity to the 5' end of the nascent transcript, with mutational studies placing **H** near the [active site](#) (8). Additionally, analysis of **rifampicin**-resistant RNAP variants has identified mutations in several regions of the b subunit that result in a lower affinity for the drug (Fig. 2) (14). Rifampicin blocks the path of nascent RNA emerging from the catalytic center, so that only short transcripts are made; therefore, the rifampicin binding site is thought to be near the catalytic center. This is supported by crosslinking studies showing that the rifampicin binding site is within 15 Å of the initiating nucleotide and within 2 Å of the -2 and -3 positions of the template DNA (14). Recent work using photocrosslinking agents is resulting in more detailed information about the location of the catalytic center and the path of RNA and DNA in the elongating complex (7).

In both eubacteria and archaea, the genes coding for the b and b'-like subunits are organized in an [operon](#), with the b-like subunit gene always preceding that for the b'-like subunit. The operon organization, in addition to providing coordinate regulation, may allow the subunits to assemble cotranslationally, which might predict that the b C-terminus is in close proximity to the b' N-terminus in the RNAP. This idea is supported by the finding that the two genes are naturally fused in the bacterium *Helicobacter pylori* (15), and experimental fusion of the two *E. coli* genes results in a catalytically active protein that is functional *in vivo* (16). In addition, crosslinking studies have confirmed that these domains are close in space in the folded protein (17).

The a subunit functions in three main processes (8). First, it plays a vital role in the assembly of the core RNAP. The core assembles in an ordered fashion  $2a > a_2 > a_2b > a_2bb'$ . Truncated proteins lacking the C-terminal third of a are competent for core assembly, showing that the N-terminal domain (a-NTD) mediates assembly. Additional studies show that discrete regions of a are required for each step of RNAP assembly. These regions of a are conserved in subunits of the other multisubunit RNAPs, suggesting that there may be similarities in the assembly of all of these enzymes. Second, using a DNA-binding function in the C-terminal domain (a-CTD), a binds to DNA sequences that are found upstream of some strong promoters and augments the activity of RNAP at those sites. Finally, a is required for transcriptional activation by a number of activators; holoenzyme containing only the a-NTD is functional for basal transcription, but cannot be activated at certain promoters. The importance of the a-CTD in activation was confirmed by experiments showing that the CRP activator makes direct contacts with the a-CTD (18).

The s subunit (see [Sigma Factors](#)) is responsible for promoter-specific DNA binding and specific initiation (19, 20). Prokaryotic organisms contain multiple sigma factors that function at various stages of growth or under different environmental conditions. Normally, the primary s factor ( $\sigma^{70}$  in *E. coli*) is used to transcribe the **housekeeping genes** of the cell. Alternative s factors coordinate expression of specific sets of genes that are required for more specialized functions. Each s factor recognizes and binds to different subsets of promoters to direct transcription of specific genes. In addition to its role in initiation of basal transcription, s is also required for transcriptional activation at certain promoters.

There are many ancillary factors that associate with the core RNAP during transcription to alter its properties. They include factors such as NusA, which increases the efficiency of pausing and termination, Gre factors, which stimulate elongation by rescuing paused RNAP complexes, and Rho protein, which functions in termination of transcription at specific sites.

## 1.2. Eukaryotic RNAPs

Eukaryotes contain three nuclear RNAPs. The three enzymes differ in their subnuclear location, sensitivity to  $\alpha$ -amanitin, and subunit composition (see Fig. 1), as well as in the types of RNA that they transcribe. RNAP I is located in the [nucleolus](#), where its sole function is the synthesis of precursor ribosomal RNA from the multiple rRNA genes. Although only a single type of promoter is utilized by RNAP I, its products may represent up to four-fifths of the RNA being synthesized in a rapidly growing cell (21). The nucleoplasmic RNAP II synthesizes mRNA and some snRNA genes and is therefore responsible for transcribing many different genes with widely divergent promoters. RNAP III, also nucleoplasmic, is responsible for the synthesis of 5 S rRNA, tRNAs, and the remaining snRNAs.

Much of the study of eukaryotic RNAPs has focused on the mechanisms of transcription initiation. The subunit structures are shown in Figure 1. Although they are more complex than the eubacterial enzyme, all represent core RNAPs that are unable to recognize promoter DNA in the absence of additional factors. Although, as described below, each RNAP has its own complement of initiation factors, they have one factor in common: Each eukaryotic nuclear RNAP requires the TATA-binding protein (TBP) for transcription initiation (22). Indeed, all eukaryotic promoters can be divided into two classes, depending on whether or not they contain a TATA box. In TATA-containing promoters, found in some but not all RNAP II- and RNAP III-dependent genes, TBP and its associated factors play an important role in binding the promoter and nucleating the association of additional factors required for recruitment of the RNAP. In TATA-less promoters, including RNAP I promoters and some RNAP II and III promoters, additional factors are required to bind first to the promoter and then to recruit TBP.

#### 1.2.1. Subunit Composition of Eukaryotic RNAPs

The eukaryotic nuclear RNAPs are built on the same basic plan as the eubacterial enzyme (Fig. 1). Each enzyme contains two large subunits that are homologous to b and b' of the eubacterial RNAP and two polypeptides with amino acid sequences similar to the a subunit. The largest subunits contain several regions of amino acid sequence and functional similarity to b and b' (Fig. 2). These subunits, like their bacterial counterparts, bind DNA and contact the nascent RNA (23). The second-largest subunit has also been implicated in substrate binding. As indicated in Fig. 1 by dotted lines, five of the smaller RNAP subunits are encoded by single genes and are common to all three enzymes, forming a shared core, while two subunits are shared by RNAP I and RNAP III. Additionally, each enzyme contains unique subunits. Despite the large number of components, most of the subunits are essential in *Saccharomyces cerevisiae* (24). The overall shape and structure of RNAP I and RNAP II from *S. cerevisiae* have been investigated by [electron microscopy](#) of two-dimensional crystals (25, 26), demonstrating similar overall shapes and containing the thumb-like projection and channel as described above for the core RNAP from *E. coli*.

#### 1.2.2. Features of RNAP I

To initiate transcription from rRNA promoters, RNAP I requires two DNA-binding transcription factors, UBF (upstream binding factor) and SL1 (TIF-IB) (21, 27). UBF is an **HMG box** protein that binds DNA as a dimer (see [DNA-Binding Proteins](#)). SL1 contains TBP and three additional factors unique to rRNA synthesis. SL1 and UBF form a complex on the core rRNA promoter that is remarkably stable and can catalyze multiple rounds of transcription initiation. Both proteins can also bind an upstream control element (UCE) that functions to enhance the function of the core promoter. RNAP I is recruited to the initiation complex through interactions with SL1. The activity of RNAP I is very sensitive to the growth state of the cell, and additional factors are involved in modulating the level of rRNA synthesis.

#### 1.2.3. Features of RNAP II

Because of its pivotal role in the regulation and expression of eukaryotic proteins, the mechanism of initiation by the eukaryotic RNAP II has been the focus of intense research. Early studies identified several factors required for selective initiation by RNAP II including TFIIB, TFIID, TFIIE, TFIIIF and TFIIH (28-30). TFIID is the RNAP II-specific, multisubunit form of TBP (22, 31). TFIIH is a multisubunit factor that also plays a critical role in nucleotide [excision repair](#) (32). Although TFIIB

is a single polypeptide chain, TFIIE and TFIIIF are also multisubunit factors, and the minimal initiation-competent form of RNAP II contains more than 30 proteins!

Several protein factors are known to associate with, and modify the activity of, the elongating form of RNAP II. TFIIIF, also a component of the initiation complex, increases the rate of RNA chain elongation (33). Like the *E. coli* Gre factors, TFIIS stimulates paused RNAP complexes that are arrested by transcriptional blocks. Both the Gre factors and TFIIS function by stimulating the inherent ability of RNAPs to reverse the polymerization reaction and to remove incorrectly incorporated bases (7). Two recently described RNAP II elongation factors, elongin (34) and P-TEFb (35), appear to play important roles in the regulated expression of subsets of mammalian genes. In particular, P-TEFb is critically important for high level of expression of genes of the **HIV virus** (35, 36).

The largest subunit of RNAP II contains a unique C-terminal domain (CTD, see Fig. 2) consisting of repeats of a heptad sequence, Tyr-Ser-Pro-Thr-Ser-Pro-Ser, not found in any other RNAP (37-39). The number of repeats varies among organisms, with yeast RNAP II containing 26 repeats, while mouse RNAP II has 52. Deletion analysis of the CTD reveals that a minimum number of repeats is required for viability in *S. cerevisiae*. This domain is the site of complex protein associations that have broad implications in the regulation of RNAP II and mRNA synthesis and processing. Extensive **phosphorylation** of the CTD tail has been observed shortly after the transition between initiation and elongation, implicating the CTD in the conversion of RNAP II into an elongation-competent form. Although several [kinases](#) have been identified that can phosphorylate the CTD, it is of particular interest that the initiation factor TFIIF and the elongation factor P-TEFb possess CTD kinase activity (32, 35). These factors appear to play a role in the transition between initiation and elongation.

The unphosphorylated form of the CTD apparently nucleates the formation of a large complex of proteins referred to as the *mediator*. A high-molecular-weight form of RNAP II, including the mediator and several of the TFII initiation factors, has been dubbed the “holoenzyme” (40). Mediator proteins are required for transcriptional activation through contact with transcriptional activator proteins. The CTD also binds a number of additional proteins involved in the subsequent processing of pre-mRNA transcripts. These proteins include cleavage and [polyadenylation](#) factors, as well as proteins that could associate the [spliceosome](#) with the transcription machinery (41).

#### 1.2.4. Features of RNAP III

Although RNAP III is the most complex of the eukaryotic nuclear RNAPs, with 17 subunits (Fig. 1), it still requires three additional basal transcription factors, TFIIA, TFIIB, and TFIIC, to initiate transcription selectively (27, 42). TFIIC from yeast contains six subunits, while TFIIB, the RNAP III TBP-containing factor, is composed of at least three subunits. TFIIB also contains a homologue of the RNAP II transcription factor TFIIB, known as BRF (for TFIIB related factor). TFIIB is the central transcription factor for RNAP III and is required to recruit the polymerase to the promoter. TFIIC serves as an assembly factor for TFIIB. It binds to elements that are internal to the gene and directs positioning of TFIIB at RNAP III promoters. TFIIA, a single polypeptide chain containing nine zinc-finger domains that bind DNA or RNA, is a gene-specific transcription factor for the 5 S rRNA genes. TFIIA binds to sequences within the 5 S rRNA gene and functions as an adaptor for TFIIC. Most genes transcribed by RNAP III, including 5 S rRNA and tRNA genes, utilize promoter sequences internal to the gene coding region. However, some RNAP III transcripts depend on upstream promoters, in some cases including a TATA box (22).

#### 1.3. Archaeal RNAP

The single RNAP from *Sulfolobus acidocaldarius* is composed of 13 subunits, with an overall structure more like the eukaryotic nuclear enzymes than the eubacterial RNAP ([Fig. 1 (2)]. The b' homologue of the archaeal RNAP is split into two subunits, A' and A'': A' corresponds to the first two-thirds of b', while A'' is homologous to the C-terminal one-third. The b' homologue is split in all known archaea, as well as in cyanobacteria and some chloroplasts (43, 44). The b homologue is also

divided in halophilic and methanobacteria (45-47). Like the eukaryotic nuclear RNAPs, two of the subunits of the archaeal RNAP have some similarity to the  $\alpha$  subunit of bacterial RNAP (Fig. 1). The *S. acidocaldarius* RNAP includes five homologues of the eukaryotic nuclear RNAP subunits, in addition to the  $\alpha$ ,  $\beta$ , and  $\beta'$  homologues. Three of these homologues are similar to the shared core subunits of the nuclear RNAPs, whereas four additional subunits do not resemble proteins from either eubacteria or eukaryotes.

Most archaeal gene promoters have an A-box motif (TTTA[T/A]A) centered around position  $-27$  that plays an important role in selection of the start site. The location and composition of the A-box closely resembles the TATA box found in many RNAP II promoters. In addition, the archaeal RNAP requires homologues of TBP and TFIIB for initiation (48), and it utilizes a homologue of the eukaryotic transcription elongation factor TFIIS (2). Overall, the archaeal transcription apparatus is much more similar to that found in the eukaryotic nucleus than to that in eubacteria.

## 2. Single-Subunit RNAPs

While the multisubunit RNAPs show a clear evolutionary history, the origin of single-subunit RNAPs is uncertain (49). **Phylogenetic** analyses group the ssRNAPs into three nonoverlapping and well-defined clusters (phage, nuclear-encoded, and mitochondrial plasmid encoded). The phage ssRNAPs include four enzymes encoded by bacteriophage T7 and T3, *Salmonella* phage SP6, and *Klebsiella* phage K11. The nuclear-encoded class was limited for a number of years to a single example, the yeast mitochondrial RNAP. Recently, genes for several putative mitochondrial RNAPs have been cloned and sequenced, including those from other fungi (50), plants (51, 52), and mammals (53). The ssRNAPs of the third group are found within the mitochondria of certain plants and fungi, where they are encoded by linear plasmids. It is possible that the plasmid-borne genes represent an early step in the acquisition of the ssRNAP from a phage, prior to transfer to the nucleus (reviewed in (49)).

One additional RNAP has recently been discovered that may represent a fourth group of ssRNAPs. In *Arabidopsis thaliana*, a second chloroplast RNAP has been identified, in addition to the chloroplast-encoded eubacterial-like enzyme (3). This nuclear-encoded ssRNAP appears to transcribe a subset of chloroplast genes. Phylogenetic comparisons show that the chloroplast ssRNAP is closely related to the mitochondrial ssRNAPs, and it seems to have arisen from a duplication of the nuclear-encoded mitochondrial ssRNAP gene (51).

### 2.1. Phage RNAPs

The phage RNAPs are the best characterized group of the ssRNAPs, with the phage T7 RNAP serving as a model for the entire class (54). The T7 RNAP is a 99-kDa enzyme. Like the multisubunit RNAPs, T7 RNAP is capable of *de novo* RNA synthesis, RNA chain elongation, and transcription termination at appropriate signals. Despite its small size, however, the phage RNAP is itself a holoenzyme with the ability to recognize and to bind promoter DNA without the addition of  $\sigma$ - or TBP-like factors. Promoters for the four phage RNAPs are similar and consist of two functional domains as defined by the analysis of point mutations (55). The region from  $-17$  to  $-6$  of the promoter (relative to the start site of transcription) constitutes the binding domain, while the initiation domain covers the region from  $-5$  to  $+6$ .

Its structure indicates that T7 RNAP is related to DNA polymerases, showing the most similarity to the Klenow fragment of DNA polymerase I (pol I) (56, 57). Although not strikingly similar at the amino acid sequence level, both T7 RNAP and pol I resemble a cupped right hand containing three subdomains: the “fingers,” “palm,” and “thumb.” There is conservation of sequence within four functional motifs that are found in many DNA-directed RNAPs (54, 58). Two of the motifs, designated “A” and “C,” contain **aspartate** residues that chelate metal ions within the catalytic pocket of the polymerases. The “B” motif is located within the fingers domain and contributes to substrate binding. These motifs are required for catalytic function and cluster within a small part of the folded structure of the RNAP. The fourth motif, designated T/DxxGR, is found in all DNA-

directed polymerases, including the largest subunits of the multimeric RNAPs, and is involved in template strand contacts. Mutations within the 'B' motif allow T7 RNAP to use deoxynucleotide triphosphate substrates, converting the RNAP to an enzyme more like a DNA polymerase (59, 60).

Footprinting studies at various points of transcription have revealed several features about the T7 RNAP (61). Before initiation, the T7 RNAP footprint protects ~20 bases. During the initial stages of transcription, the downstream border of the footprint extends, while the upstream border remains fixed, with the result that the footprint increases in length, reaching a maximum length of ~30 bases. The static nature of the upstream border suggests that the RNAP maintains contact with the promoter during this early abortive stage of transcription. After the transition to processive transcription, the upstream border of the footprint moves downstream, until the footprint returns to the preinitiation length. The transition from abortive to processive transcription is thought to depend on an RNA/RNAP interaction and a conformational change in the structure of the RNAP (62). On the basis of the structural data, it is easy to envision a conformational change in which the thumb or finger domain closes around the DNA, thus preventing its premature release from the RNAP. Similarly, hairpin structures that terminate transcription may disrupt the RNA/RNAP interaction, reversing the conformational change and leading to a release of the RNAP, as is also envisioned for the multisubunit RNAPs (7).

## 2.2. Mitochondrial RNAP

The yeast mitochondrial core RNAP is encoded by a nuclear gene (*RPO41*) with sequence similarity to T3 and T7 phage RNAPs (Fig. 1). The similarity extends along the entire length of phage RNAPs and includes the amino acid residues that make up the A, B, C, and T/DxxGR motifs mentioned above (63). The similarity to phage RNAPs is surprising, because it was thought that the mitochondrial transcription machinery would reflect the eubacterial origin of the mitochondria (3). Until recently, the possibility existed that the yeast mitochondrial RNAP was an evolutionary anomaly, not representative of mitochondrial RNAPs of other eukaryotes. However, the recent identification of related ssRNAPs in other fungi, plants, and mammals has confirmed the yeast enzyme as a useful model for the class. Homologues have also been identified in algae and protozoa, including some of closest relatives of amitochondriate eukaryotes (64), demonstrating that the T7-like RNAP was recruited at an early stage in the evolution of the mitochondria.

Despite homology to T7 RNAP, yeast mitochondrial RNAP requires a second factor for promoter recognition and initiation of transcription (Fig. 1) (63). In yeast, the factor is also encoded by a nuclear gene (*MTF1*). This promoter recognition factor functions like sigma factors during the initial stages of transcription, interacting with the core RNAP in solution to provide specific promoter recognition. The core RNAP/factor complex recognizes a short nonanucleotide promoter sequence (ATATAAGTA) that resembles the -10 TATAAT region of the eubacterial sigma 70 promoter. Like sigma, the mitochondrial factor is required only for initiation, and it is released shortly after. Whether or not a separate s-like initiation factor is required for all mitochondrial RNAPs is still unclear, although the reported composition of the *Xenopus* mitochondrial RNAP includes a 40-kDa dissociable initiation factor (65). The mitochondrial RNAPs, therefore may represent a hybrid between the initiation factor-requiring multisubunit enzymes and the self-sufficient ssRNAPs.

## 3. Concluding Remarks

This brief description of DNA-dependent RNA polymerases omits many important details about mechanism and regulation of these critical enzymes. In addition, many interesting enzymes, including several of bacteriophage and viral origin, have not been described. The factors regulating the activity of the RNA polymerases represent a very large and rapidly growing family of proteins, and the reader is encouraged to explore the references below and additional sources to learn more about these complex enzymes and their role in transcribing the genetic information in DNA into RNA.

## Bibliography



1. R. I. Scheinman et al. (1995) *Science* **270**, 283–286.
2. R. Langer, J. Hain, P. Thuriaux, and W. Zillig (1995) *Proc. Natl. Acad. Sci. USA* **92**, 5768–5772.
3. M. Gray and F. Lang (1998) *Trends Microbiol.* **6**, 1–3.
4. B. Moss (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway eds.), Raven Press, New York, pp. 185–205.
5. K.-C. Sonntag and G. Darai (1996) *Virus Genes* **11**, 271–284.
6. W. R. McClure (1985) *Annu. Rev. Biochem.* **54**, 171–204.
7. P. H. von Hippel (1998) *Science* **281**, 660–665.
8. C. A. Gross, C. L. Chan, and M. A. Lonetto (1996) *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **351**, 475–482.
9. A. Polyakov, E. Severinova, and S. A. Darst (1995) *Cell* **83**, 365–373.
10. S. A. Darst, E. W. Kubalik, and R. D. Kornberg (1989) *Nature* **340**, 730–732.
11. P. Thuriaux and A. Sentenac (1992) in *The Molecular and Cellular Biology of the Yeast Saccharomyces: Gene Expression* (E. W. Jones, J. R. Pringle, and J. R. Broach, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–48.
12. L. M. Heisler et al. (1996) *J. Biol. Chem.* **271**, 14572–14583.
13. R. A. Young (1991) *Annu. Rev. Biochem.* **60**, 689–715.
14. D. J. Jin and Y. N. Zhou (1996) *Methods Enzymol.* **273**, 300–319.
15. N. Zakharova, P. S. Hoffman, D. E. Berg, and K. Severinov (1998) *J. Biol. Chem.* **273**, 19371–19374.
16. K. Severinov, R. Mooney, S. A. Darst, and R. Landick (1997) *J. Biol. Chem.* **272**, 24137–24140.
17. E. Nudler et al. (1998) *Science* **281**, 424–428.
18. Y. Chen, Y. W. Ebright, and R. H. Ebright (1994) *Science* **265**, 90–92.
19. C. Gross, M. Lonetto, and R. Losick (1992) in *Transcriptional Regulation* (K. Yamamoto and S. McKnight, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 129–176.
20. J. D. Helmann (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.), Raven Press, New York, pp. 1–17.
21. M. R. Paule (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.) Raven Press, New York, pp. 83–106.
22. N. Hernandez (1993) *Genes Dev.* **7**, 1291–1308.
23. M. Sawadogo and A. Sentenac (1990) *Annu. Rev. Biochem.* **59**, 711–754.
24. J. Archambault and J. Friesen (1993) *Microbiol. Rev.* **57**, 703–724.
25. S. A. Darst, A. M. Edwards, E. W. Kubalek, and R. D. Kornberg (1991) *Cell* **66**, 121–128.
26. P. Schultz et al. (1993) *EMBO J.* **12**, 2601–2607.
27. E. P. Geiduschek and G. A. Kassavetis (1995) *Curr. Opin. Cell Biol.* **7**, 344–351.
28. R. C. Conaway and J. W. Conaway (1993) *Annu. Rev. Biochem.* **62**, 161–190.
29. L. Zawel and D. Reinberg (1995) *Annu. Rev. Biochem.* **64**, 533–561.
30. R. G. Roeder (1996) *TIBS* **21**, 327–334.
31. W. P. Tansey and W. Herr (1997) *Cell* **88**, 729–732.
32. J. Svejstrup, P. Vichi, and J. Egly (1996) *TIBS* **21**, 346–350.
33. T. Aso, J. Conaway and R. Conaway (1995) *FASEB J.* **9**, 1419–1427.
34. T. Aso, W. S. Lane, J. W. Conaway, and R. C. Conaway (1995) *Science* **269**, 1439–1440.
35. Y. Zhu, T. Pe'ery, J. Peng, Y. Ramanathan, N. Marshal, et al. (1997) *Genes Dev.* **11**, 2622–2632.

36. H. Mancebo, G. Lee, J. Flygare, J. Tomassini, P. Luu, et al. (1997) *Genes Dev.* **11**, 2633–2644.
37. A. Emili and C. J. Ingles (1995) *Curr. Biol.* **5**, 204–209.
38. N. A. Woychik and R. A. Young (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.) Raven Press, New York, pp. 227–242.
39. M. Dahmus (1996) *J. Biol. Chem.* **271**(32), 19009–19012.
40. A. J. Koleske and R. A. Young (1995) *TIBS* **3**, 113–116.
41. E. Steinmetz (1997) *Cell* **89**, 491–494.
42. G. A. Kassavetis, C. Bardeleben, B. Bartholomew, et al. (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.) Raven Press, New York, pp. 107–126.
43. W. W. Xie, K. Jager, and M. Potts (1989) *J. Bacteriol.* **171**, 1967–1973.
44. K. J. Bergsland and R. Haselkorn (1991) *J. Bacteriol.* **173**, 3446–3455.
45. G. Puhler, H. Leffers, F. Gropp, P. Palm, H. P. Klenk, et al. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 4569–4573.
46. H. Leffers et al. (1989) *J. Mol. Biol.* **206**, 1–17.
47. B. Berghofer, L. Krockel, C. Kortner, M. Truss, J. Schallenberg, et al. (1988) *Nucleic Acids Res.* **16**, 8113–8128.
48. S. A. Qureshi, S. D. Bell, and S. P. Jackson (1997) *EMBO J.* **16**, 2927–2936.
49. N. Cermakian, T. M. Ikeda, P. Miramontes, B. F. Lang, M. W. Gray, et al. (1997) *J. Mol. Evol.* **45**, 671–681.
50. B. Chen, A. R. Kubelik, S. Mohr, and C. A. Breitenberger (1996) *J. Biol. Chem.* **271**, 6537–6544.
51. B. Hedtke, T. Borner, and A. Weihe (1997) *Science* **277**, 809–811.
52. A. Weihe, B. Hedtke, and T. Borner (1997) *Nucleic Acids Res.* **25**, 2319–2325.
53. V. Tiranti, A. Savoia, F. Forti, M. D'Apollito, M. Centra, et al. (1997) *Hum. Mol. Genet.* **6**(4), 615–625.
54. R. Sousa (1996) *TIBS* **21**, 186–190.
55. T. Li, H. H. Ho, M. Maslak, C. Schick, and C. T. Martin (1996) *Biochemistry* **35**, 3722–3727.
56. R. Sousa, Y. Chung, J. Rose, and B. Wang (1993) *Nature* **364**, 593–599.
57. D. Jeruzalmi and T. A. Steitz (1998) *EMBO J.* **17**, 4101–4113.
58. O. Poch, I. Sauvaget, M. Delarue, and N. Tordo (1989) *EMBO J.* **8**(12), 3867–3874.
59. R. Sousa and R. Padilla (1995) *EMBO J.* **14**, No. 18, 4609–4621.
60. D. A. Kostyuk, S. M. Dragan, D. L. Lyakhov, V. O. Rechinsky, V. L. Tunitskaya, et al. (1995) *FEBS Lett.* **369**, 165–168.
61. R. Ikeda and C. Richardson (1986) *Proc. Natl. Acad. Sci. USA* **83**, 3614–3618.
62. C. T. Martin, D. K. Muller, and J. E. Coleman (1988) *Biochemistry* **27**, 3966–3974.
63. S. Jang and J. Jaehning (1994) in *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.), Raven Press, New York, NY, pp. 171–184.
64. N. Cermakian, T. M. Ikeda, R. Cedergren, and M. W. Gray (1996) *Nucleic Acids Res.* **24**, 648–654.
65. I. Antoshechkin and D. F. Bogenhagen (1995) *Mol. Cell. Biol.* **15**, 7032–7042.

### Suggestions for Further Reading

66. R. C. Conaway and J. W. Conaway, eds. (1994) *Transcription: Mechanisms and Regulation* Raven Press, New York. (This book contains several excellent reviews on eukaryotic nuclear, viral and mitochondrial RNA polymerases and regulatory factors.)
67. J. Archambault and J. Friesen (1993) *Genetics of eukaryotic RNA polymerases I, II, and III.*

Microbiol. Rev. **57**: 703–724.

68. C. A. Gross, C. L. Chan, and M. A. Lonetto (1996) A structure/function analysis of Escherichia coli RNA polymerase. Philos. Trans. R. Soc. Lond. B. Biol. Sci. **35**, 475–482.
69. R. Langer, J. Hain, P. Thuriaux, and W. Zillig (1995) Transcription in Archaea: similarity to that in Eucarya. Proc. Natl. Acad. Sci. USA **92**: 5768–5772.
70. P. Thuriaux and A. Sentenac (1992) "Yeast Nuclear RNA Polymerases". In: *The Molecular and Cellular Biology of the Yeast Saccharomyces: Gene Expression* (E. W. Jones, J. R. Pringle and J. R. Broach, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–48.

## RNA Polymerases, RNA-Dependent

In all organisms the major flow of genetic information is from **DNA** to **RNA** using DNA-dependent **RNA polymerases**, with the DNA itself being replicated by DNA-dependent **DNA polymerases**. Some of the DNA synthesis, primarily of specific viruses ([retroviruses](#)) and retroelements, is from RNA by reverse transcriptase. However, there are many **viruses** of both prokaryotes and eukaryotes that have RNA **genomes** and that replicate by the synthesis of RNA from the genome. This replication is performed by RNA-dependent RNA polymerases (EC 2.7.7.48) (RdRp) (also termed "*RNA replicases*") and usually involves other additional [enzyme](#) activities such as [RNA Helicases](#) and capping and **methylation** enzymes; these are frequently on the same protein molecule as the polymerase itself. There is also some evidence for RdRp activity in uninfected cells, and there was initially considerable controversy as to whether RNA viruses were replicated by host or viral enzyme (s) (1)

RdRp can be considered to be an ancient enzymatic activity. It is widely suggested that, at an early stage in the evolution of life, there was the "[RNA World](#)" (2) with the initial synthetic and processing activities being performed by catalytic RNAs, or **ribozymes**. In this scene, the early proteins would have been better enzymes than their RNA counterparts and thus would eventually dominate (2). It is likely that RdRp activity was among these early proteins.

### 1. Uninfected Cells

RdRp activity has been found in uninfected cells from both animals and plants. Mouse erythroleukemia cells contain a cytoplasmic RNA synthesis activity which is resistant to [actinomycin D](#), leads to the formation of (–) strand globin [messenger RNA](#) and can use poly(A)-oligo(U) as a template-primer combination (3, 4). This activity required  $Mg^{2+}$  but was inhibited by  $Mn^{2+}$ . The best-characterized plant RdRp is that from tomato leaves (5, 6). This 128-kDa polypeptide chain sedimented at 6.6 S in its native form, which indicated that it was probably a monomer. Its optimum activity was at pH 7.8. It required a divalent cation,  $Mg^{2+}$ , which is much more effective than  $Mn^{2+}$ . The enzyme activity could transcribe RNA from both RNA and DNA templates, most of the transcripts being very short. Transcription could take place either with or without a primer; in the latter case priming was at, or close, to the 3' end of the template.

As noted above, there was controversy about the role of host RdRp in viral RNA replication, which was exacerbated by the fact that the host RdRp activity increases upon virus infection. However, it is now generally assumed that the host RdRp is not involved, at least significantly, in viral replication. Its function is unknown, but the short lengths of RNA that it transcribes might be involved in some control mechanism, such as post-transcriptional suppression of **gene expression** (7).

## 2. Viral RdRps

Viral RdRps have three remarkable features: (i) they amplify their RNA template many-fold during a very short infection period; (ii) they specifically replicate the viral genomic RNA in the presence of a great excess of host RNAs; (iii) they copy the entire template RNA, and also express subgenomic RNAs, in most cases without using endogenous primers. In all viral enzyme complexes, the catalytic polymerase activity is virus-encoded. The detailed mechanisms by which viral RNA replicases operate are unknown, but there would appear to be differences between them, as is exemplified by the current understanding of different viruses.

### 2.1. Bacteriophage Q $\beta$

The RdRp holoenzyme of coliphage Q $\beta$  is capable of *in vitro* synthesis of (–)-strand RNA from a (+)-strand template; it consists of four subunits, three of which are host-derived proteins [**ribosomal** protein S1 (subunit a) and translation **elongation factors** EF-Tu(g) and EF-TS(d)], and the fourth(b) is virus-encoded and is the catalytic subunit. Details of the function and structure of the basic holoenzyme are reviewed by Blumenthal and Carmichael (8) and by Ishihama and Barbier (9) and a working model has been proposed (10). Recent evidence indicates that a host factor is also required for efficient *in vivo* RNA replication (11). As with other RdRps, that of Q $\beta$  is template-specific, and the template recognition involves complex interactions between the holoenzyme and the RNA. The tertiary structure of the RNA is recognized as being important in bringing together the sites with which the proteins of the holoenzyme interact (12). Recent studies suggest that template recognition and specificity involve both the carboxyl-terminal region of the virus-encoded subunit and the host-derived S1 protein (12, 13) (see also **b Replicase**).

### 2.2. Poliovirus

The core enzyme of poliovirus RdRp (the 3D gene product) is processed by the adjacent gene product (3C) from the **polyprotein** translated from the viral genome, and it forms a membrane-associated complex with the primer **VPg protein** (which covalently attaches to the 5' end of the viral RNA) and protein 3C. Although there is not a clear picture of the detailed mechanism by which poliovirus, and other picornaviruses, replicate their RNA, there are various observations that have to be incorporated in any model (14). These include: (i) the need for the replication complexes to be membrane-associated for initiation and elongation of minus- and plus-strand synthesis; (ii) the VPg is covalently linked to all newly synthesized RNA; (iii) the primer for the processive RdRp is unknown although there are various hypotheses; (iv) the precursor 3CD protein binds to the 5'-terminal cloverleaf structure of the plus strand in the presence of cellular and/or viral factors; (v) various cellular proteins have been implicated in RNA synthesis, including a report that human protein Sam68 interacts with poliovirus 3D protein (15).

### 2.3. Plant (+)-Strand RNA Viruses

Understanding of the functioning of the RdRps of plant (+)-strand RNA viruses is limited because of the difficulties in isolating actively replicating systems from the membrane-bound replication complexes in infected cells. The replication complexes of several viruses have been shown to contain host-encoded proteins, as well as the predicted virus-encoded products. The replication complex of **tobacco mosaic virus** (TMV) associates with a protein related to GCD10 protein, which is the RNA-binding subunit of yeast eIF-3 (16) and that of brome mosaic virus associates with the 41-kDa subunit of wheat eIF-3 (17); the replication complex of cucumber mosaic virus is associated with a 50-kDa host protein (18). As the TMV and brome mosaic virus replicase complexes are associated with different eIF-3 subunits, it is thought likely that these host proteins play different roles in their respective virus replications (16).

The replication of RNA by RdRp is considered to have several stages covering initiation to elongation (19). There is recent evidence pointing towards the involvement of the **transfer RNA**-like structures at the 3' end of virion RNAs of several plant viruses in the initiation and priming of at least (–)-strand synthesis (19, 20).

## 2.4. Double-Stranded RNA-Viruses and (–)-Strand RNA Viruses

As the genomes of double-stranded or (–)-strand RNA viruses cannot be translated directly, the virus particles carry the RdRp to initiate synthesis of mRNA on entry into the infected cell.

The RdRp of [influenza virus](#), which has a (–)-strand segmented genome, has essentially two functions, the synthesis of mRNA from the genome segments and the replication of the virion RNA via the complementary RNA. Influenza mRNA synthesis is primed at the 3' end of the virion RNA by 10- to 13-nucleotide capped RNA fragments that are “cap-snatched” from the 5' ends of host heterogeneous nuclear RNA (**hnRNA**). The exact details of priming of the complementary RNA are not yet fully understood.

## 3. Groups of RdRps

Kamer and Argos ([21](#)) recognized various sequence motifs characteristic of RdRps. The most conserved of these is the central Gly-Asp-Asp (GDD) triplet flanked by 5-residue segments that are mainly **hydrophobic** amino acids. This is taken as being suggestive of a b-hairpin structure comprising two **hydrogen-bonded** antiparallel **beta-strands** connected by a short exposed loop containing the GDD amino acids. Mutation of the GDD box can abolish the replicase function (reviewed in ref. [22](#)).

Alignment of the RdRp sequences has been extended with the recognition of eight motifs ([23](#), [24](#)). This led to the classification of the RdRps into three supergroups (Table [1](#)), with a number of lineages within each supergroups. The supergroups extend across viruses that infect animals, plants and bacteria and bring together viruses whose genome structures and expression have properties in common. However, these groupings have not been fully supported by other analyses ([25-27](#)) and care has to be taken in deriving evolutionary relationships from these analyses.

**Table 1. Grouping of RdRps<sup>a</sup>**

| Supergroup 1                              | Supergroup 2                | Supergroup 3        |
|---|-----------------------------|---------------------|
| Picornaviruses (a) <sup>a</sup>           | Leviviruses (b)             | Togaviruses (a)     |
| Arteriviruses (a)                         | Flaviviruses (a)            | Caliciviruses (a)   |
| Astroviruses (a)                          | Pestiviruses (a)            | Bromoviruses (p)    |
| Caliciviruses (a)                         | Dianthoviruses (p)          | Capilloviruses (p)  |
| Coronaviruses (a)                         | Enamovirus RNA2 (p)         | Carlaviruses (p)    |
| Nodaviruses (i)                           | Luteoviruses subgroup I (p) | Closteroviruses (p) |
| Comoviruses (p)                           | Machlomoviruses (p)         | Furoviruses (p)     |
| Enamovirus RNA1 (p)                       | Necroviruses (p)            | Hordeiviruses (p)   |
| Luteoviruses subgroup II (p)              | Tombusviruses (p)           | Idaeoviruses (p)    |
| Potyviruses (p)                           | Umbraviruses (p)            | Potexviruses (p)    |
| Sequiviruses (p)                          |                             | Tobamoviruses (p)   |
| Sobemoviruses (p)                         |                             | Tobraviruses (p)    |
| Barnaviruses (f)                          |                             | Trichoviruses (p)   |
|   |                             | Tymoviruses (p)     |
| Properties in common for most of subgroup |                             |                     |

|                        |                              |                              |
|------------------------|------------------------------|------------------------------|
| 1 genome segment       | 1 to several genome segments | 1 to several genome segments |
| Polyprotein expression | Individual gene translation  | Individual gene translation  |
| VPg                    | Capped RNA                   | Capped RNA                   |

---

<sup>a</sup> Ref. [24](#).

<sup>b</sup> Hosts of viruses: (a) = higher animal; (i) = insect; (p) = plant; (f) = fungus; (b) = bacteria.

#### 4. Structure

The three-dimensional structure of the RdRp of poliovirus has been determined by [X-ray crystallography](#) to 0.26 nm (2.6 Å) resolution ([28](#)). The overall shape of this polymerase resembles those of other polymerases, being likened to a right hand. The palm domain that contains the catalytic core structure is very similar to that of other polymerases, but the structures of the “fingers” and the “thumb” subdomains differ from those of other polymerases. Extensive regions of interactions between neighboring molecules were observed in the crystals, which suggests that an unusual higher order structure might be important in polymerase function.

#### 5. Mechanism Action

As noted above, it has been suggested that the functioning of RdRps has several stages ([19](#)). By analogy with the functioning of DNA-dependent DNA polymerases, these were proposed as being: (i) template binding, (ii) promoter localization, (iii) melting the template to give a transcriptionally open complex, (iv) nucleotide substrate binding, (v) formation of the first phosphodiester bond, (vi) promoter clearance, and (vii) progressive elongation. Relatively little is known about any of these functional stages for RdRps, and it is likely that there will be differences between the various virus systems for at least some of the stages. For example, there is evidence for priming of some of the plant (+)-strand viruses at the tRNA-like 3' terminal structures, whereas other plant viral RNAs do not have these structures. Also, there appear to be two different priming systems for influenza mRNA and complementary RNA synthesis. There also has to be priming for the formation of the (–) strand, most likely at the 3' end of the (+)-strand, and then for the synthesis of the (+) strand, again most likely at the 3' end of the (–)-strand. The latter is the complement of the 5' end of the (+)-strand and, for most viruses, the 5' and 3' genome sequences bear little resemblance.

This lack of detailed information is mainly due to the difficulties of studying these enzymes systems, especially those of eukaryotes that are associated with membranes. Furthermore, as well as the RdRp itself, the replication complex also contains several other virus-encoded and host-encoded activities. Some of the other virus-encoded activities, such as helicases (probably involved in the melting of the template stage), capping, and methylation enzymes are reviewed by Buck ([22](#)).

RdRps lack proofreading activity, and thus, there is a high rate of error in the synthesis of the new RNA strand. It is estimated that 1 in  $10^3$  to  $10^4$  nucleotides is misincorporated, which gives a high rate of genome [mutation](#) ([29](#)). This, coupled with the high rate of replication, can lead to rapid evolution of RNA viruses.

#### Bibliography

1. H. Fraenkel-Conrat (1986) *Crit Rev. Plant Sci.* **4**, 213–226.
2. W. Gilbert (1986) *Nature* **319**, 618.
3. V. Volloch (1986) *Proc. Natl. Acad. Sci. USA* **83**, 1208–1212.

4. V. Volloch, B. Schweitzer, and S. Rits (1987) *J. Cell Biol.* **105**, 137–143.
5. W. Schiebel, B. Haas, S. Marinkovi, A. Klanner, and H. L. Sanger (1993) *J. Biol. Chem.* **268**, 11851–11857.
6. W. Schiebel, B. Haas, S. Marinkovi, A. Klanner, and H. L. Sanger (1993) *J. Biol. Chem.* **268**, 11858–11867.
7. W. G. Dougherty and T. D. Parks (1995) *Curr. Opin Cell Biol.* **7**, 399–405.
8. T. Blumenthal and G. G. Carmichael (1979) *Ann. Rev. Biochem.* **48**, 525–548.
9. A. Ishihama and P. Barbier (1994) *Arch. Virol.* **134**, 235–258.
10. D. Brown and L. Gold (1996) *Proc. Natl. Acad. Sci. USA* **93**, 11558–11562.
11. S. Qiu, D. Schuppli, H.-C. T. Tsui, M. E. Winkler, and H. Weber (1997) *Virology* **227**, 211–214.
12. G. Miranda, D. Schuppli, I. Barrera, C. Hausherr, J. M. Sogo, and H. Weber (1997) *J. Molec. Biol.* **267**, 1089–1103.
13. Y. Inokuchi and M. Kajitani (1997) *J. Biol. Chem.* **272**, 15339–15345.
14. E. Wimmer, C. U. T. Helen, and X. Cao (1993) *Annu. Rev. Genet.* **27**, 353–436.
15. A. E. McBride, A. Schlegel, and K. Kirkegaard (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2296–2301.
16. T. A. M. Osman and K. W. Buck (1997) *J. Virol.* **71**, 6075–6082.
17. R. Quadt, C. C. Kao, K. S. Browning, R. P. Hershberger, and P. Ahlquist (1993) *Proc. Natl. Acad. Sci. USA* **90**, 1498–1502.
18. R. J. Hayes and K. W. Buck (1990) *Cell* **63**, 363–368.
19. J.-H. Sun and C. C. Kao (1997) *Virology* **233**, 63–73.
20. R. N. Singh and T. W. Dreher (1997) *Virology* **233**, 430–439.
21. G. Kamer and P. Argos (1984) *Nucl. Acids Res.* **12**, 7269–7282.
22. K. Buck (1997) *Adv. Virus Res.* **47**, 159–251.
23. E. V. Koonin (1991) *J. Gen. Virol.* **72**, 2197–2206.
24. E. V. Koonin and V. V. Dolja (1993) *Crit Rev. Biochem. Mol. Biol.* **28**, 375–430.
25. R. Goldbach and P. De Haan (1994) in (S. S. Morse, ed.), *The evolutionary Biology of Viruses* Raven Press, New York, pp. 105–119.
26. J. Bruenn (1991) *Nucl. Acids Res.* **19**, 217–225.
27. P. M. DeA. Zanotto, M. J. Gibbs, E. A. Gould, and E. C. Holmes (1996) *J. Virol.* **70**, 6083–6096.
28. J. L. Hansen, A. M. Long and S. C. Schultz (1997) *Structure* **5**, 1109–1122.
29. J. Holland, K. Spindler, F. Horodyski, E. Grabau, S. Nichol, and S. VanderPol (1982). *Science* **215**, 1577–1585.

### **Suggestions for Further Reading**

30. A. Ishihama and P. Barbier (1994) Molecular Anatomy of Viral-Directed RNA Polymerases; *Arch. Virol.* **134**, 235–258.
31. K. Buck (1996) Comparison of the Replication of Positive-Stranded RNA Viruses of Plants and Animals, *Adv. Virus Res.* **47**, 159–251.

### **RNA Sequencing**

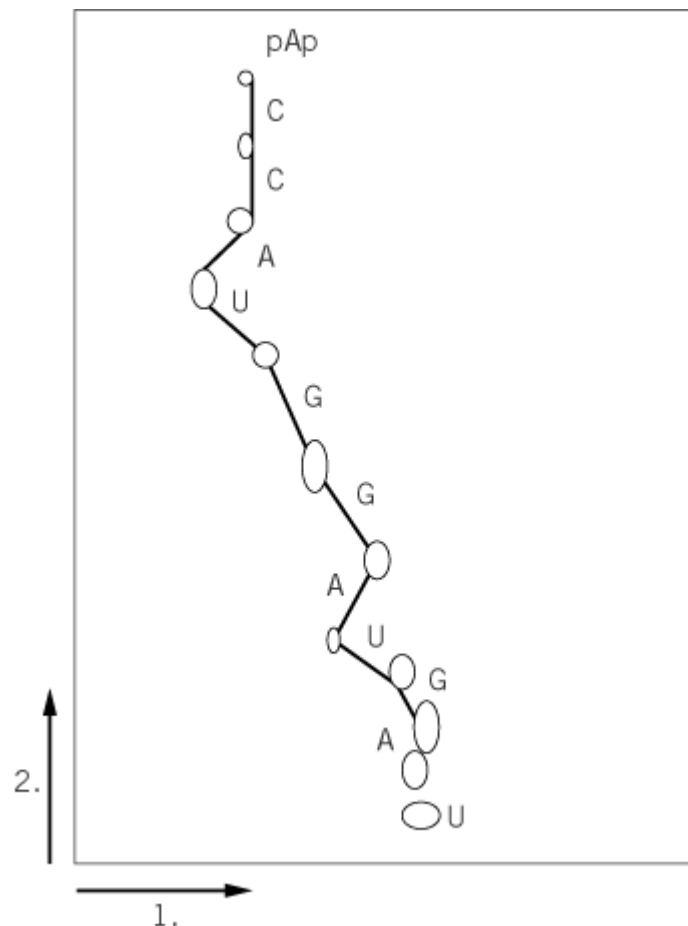
The majority of **RNA** sequences have been obtained by sequencing of either **cloned DNA** for **RNA genes** or **complementary DNA** prepared by **transcription** of RNA with **reverse transcriptase**. Such sequences of RNA genes do not, however, provide information about post-transcriptional modification of RNA. In addition to **RNA splicing**, cleavage, addition of **cap**, **polyadenylation**, and **RNA Editing**, the primary transcripts of eukaryotic and prokaryotic RNAs are often processed by specific modification of the nucleotide residues. Almost 100 base and ribose modifications of the standard four ribonucleotides are known (1), and they are located primarily in the functionally important regions of RNAs. Consequently, rapid methods for RNA sequencing, similar to those for **DNA sequencing**, have been developed, in addition to special chemical, enzymatic, and spectrophotometric techniques to identify modified nucleotides in RNA.

## 1. Sequencing Oligoribonucleotides

Methods to sequence RNA were developed almost two decades before rapid DNA sequencing technology became available. The early sequencing of **transfer RNA** used column **chromatography** to separate oligoribonucleotides. Guanine-specific **ribonuclease T1** (RNase T1) from *Aspergillus oryzae* and pyrimidine-specific **ribonuclease A** (RNase A) from bovine pancreas were used to achieve alternative cleavages of RNA to oligoribonucleotides. Spectroscopic characterization, determination of nucleoside composition, and stepwise degradation provided partial sequences from which the sequence of the original RNA was derived. Two-dimensional separation by high-voltage **electrophoresis** and chromatography of oligoribonucleotides obtained by ribonuclease cleavage, so-called fingerprinting (2)—considerably increased the resolution and speed of RNA sequencing. Two oligoribonucleotides that differ by one nucleotide in length have a characteristic mobility shift in the two-dimensional fingerprint depending on the nucleotide by which they differ. These shifts depend mainly on the  $pK_a$  values of the particular nucleotide residue. After partial cleavage of end-labeled oligoribonucleotides, followed by two-dimensional separation and chromatography, a sequence of 8–15 nucleotides could be read in one experiment (Fig. 1) (3).

**Figure 1.** Two-dimensional separation of radioactively end-labeled oligoribonucleotides obtained by limited cleavage by an endonuclease. First dimension: electrophoresis; second dimension: chromatography. The shift in mobility on removal of each 3'-terminal nucleotide provides its identity (3).





## 2. Rapid Sequencing by Electrophoresis of End-labeled RNA on Polyacrylamide Gels

The methods developed for rapid sequencing of DNA have been modified for RNA sequencing. 5'-End-labeling is usually achieved by dephosphorylation of RNA with [alkaline phosphatase](#), followed by phosphorylation with phage T4 **polynucleotide kinase** and [ $g\text{-}^{32}\text{P}$ ]ATP (4). Alternatively, the 3'-end of an oligoribonucleotide can be phosphorylated by cytidine-3',5' [ $^{32}\text{P}$ ] diphosphate ( $[^{32}\text{P}]\text{pCp}$ ) and T4 [RNA ligase](#) in the presence of ATP (5). Limited base-specific cleavage of end-labeled RNA can be achieved either enzymatically (6, 7) or chemically (8) (Table 1). The oligoribonucleotides generated are separated by electrophoresis, and the bands with the label are detected by [autoradiography](#). For comparison, a nonspecifically cleaved sample is usually prepared by hydrolysis of an end-labeled RNA.

**Table 1. Cleavage Reactions for RNA Sequencing**

(a) Enzymatic cleavage (6, 7)

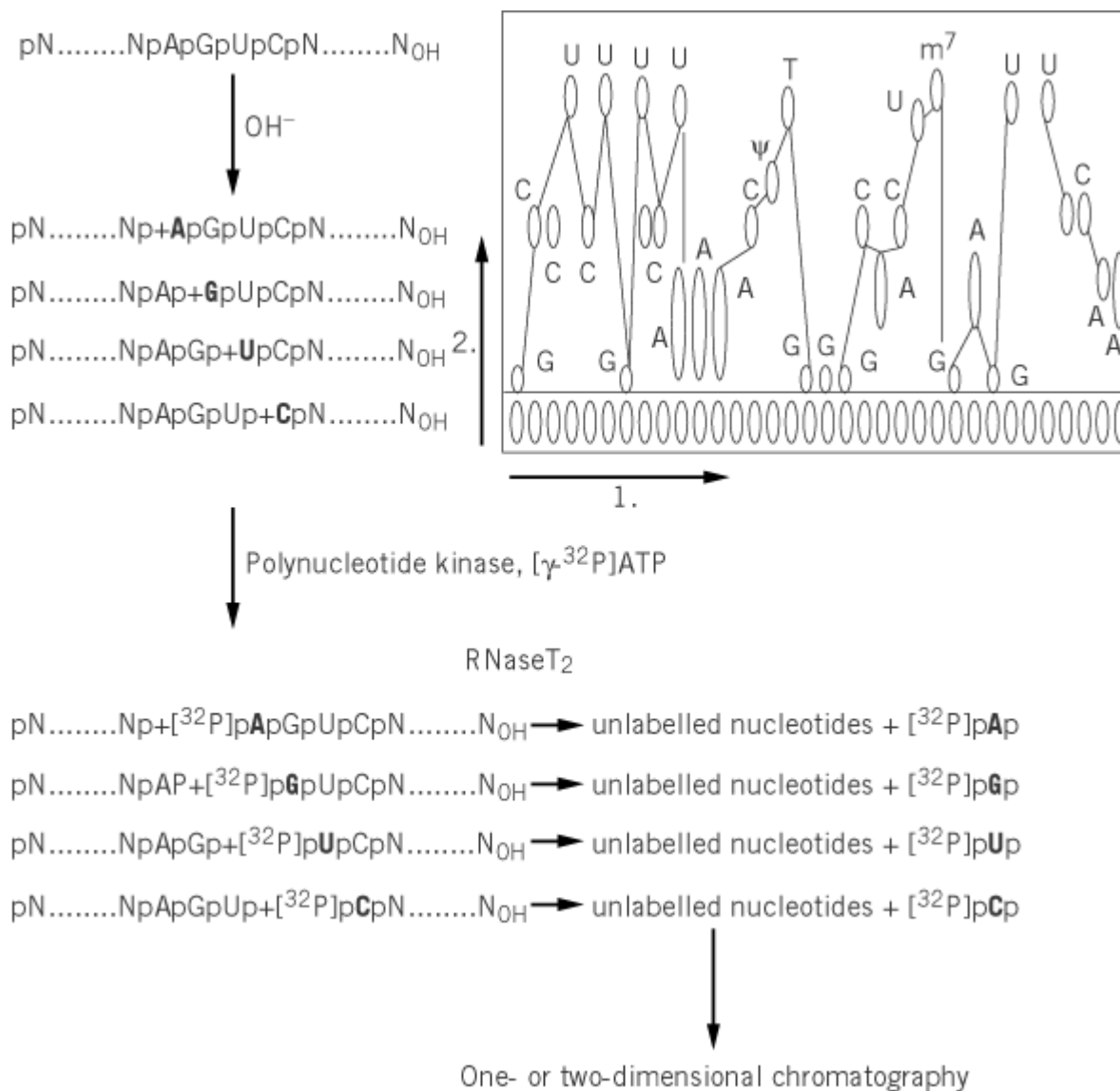
| Ribonuclease (RNase)                                  | Specificity                        | Product                                      |
|---|------------------------------------|--|
| RNase T <sub>1</sub> from <i>Aspergillus oryzae</i>   | $\frac{1}{4}\text{GpN}\frac{1}{4}$ | $\frac{1}{4}\text{Gp} + \text{N}\frac{1}{4}$ |
| RNase U <sub>2</sub> from <i>Ustilago sphaerogena</i> | $\frac{1}{4}\text{ApN}\frac{1}{4}$ | $\frac{1}{4}\text{Ap} + \text{N}\frac{1}{4}$ |
| RNase Phy M from <i>Physarum polycephalum</i>         | $\frac{1}{4}\text{ApN}\frac{1}{4}$ | $\frac{1}{4}\text{Ap} + \text{N}\frac{1}{4}$ |
|   | $\frac{1}{4}\text{UpN}\frac{1}{4}$ | $\frac{1}{4}\text{Up} + \text{N}\frac{1}{4}$ |

|                                    |                             |                                |
|------------------------------------|-----------------------------|--------------------------------|
| RNase A from bovine pancreas       | $\frac{1}{4}UpN\frac{1}{4}$ | $\frac{1}{4}Up + N\frac{1}{4}$ |
|                                    | $\frac{1}{4}CpN\frac{1}{4}$ | $\frac{1}{4}Cp + N\frac{1}{4}$ |
| RNase from <i>Bacillus cereus</i>  | $\frac{1}{4}UpN\frac{1}{4}$ | $\frac{1}{4}Up + N\frac{1}{4}$ |
|                                    | $\frac{1}{4}CpN\frac{1}{4}$ | $\frac{1}{4}Cp + N\frac{1}{4}$ |
| (b) Chemical cleavage (8)          |                             |                                |
| Reagent                            | Cleavage                    | Specificity                    |
| Dimethylsulfate, NaBH <sub>4</sub> | Aniline (pH 4.5)            | G                              |
| Diethylpyrocarbonate               | Aniline (pH 4.5)            | A > G                          |
| Hydrazine                          | Aniline (pH 4.5)            | U > C                          |
| Hydrazine in 3M NaCl               | Aniline (pH 4.5)            | C > U                          |

---

A serious disadvantage of the rapid sequencing methods is their inability to detect modified nucleotides in RNA. Methods to detect them were developed that combine polyacrylamide [gel electrophoresis](#) (PAGE) separation of oligoribonucleotides according to their length, with chromatographic identification of the radioactively labeled nucleotide present at the cleavage site ([9](#), [10](#)). These methods were especially successful in determining the sequences of tRNAs. A mixture of oligoribonucleotides is prepared by nonspecific cleavage in aqueous solution in the presence of **formamide** at elevated temperature. The 5'-hydroxyl groups of the oligoribonucleotides formed by this cleavage are then labeled by T4 RNA kinase and [g-<sup>32</sup>P]ATP. After separation by gel electrophoresis, bands of the different oligoribonucleotides are blotted onto chromatographic plates and digested *in situ* to mononucleotides by RNase T2. For each band, the terminal <sup>32</sup>P-labeled nucleotide is identified separately by one- or two-dimensional chromatography (Fig. [2](#)).

**Figure 2.** Sequencing of RNA by the post-labeling method. Limited alkaline hydrolysis leads to mixture of oligoribonucleotides, which are end-labeled and then separated by polyacrylamide gel electrophoresis (PAGE; first dimension) and by chromatography (second dimension). After hydrolysis of each fragment by ribonuclease T2, the identity of the labeled nucleotide at the terminus is identified by chromatography ([9](#)).



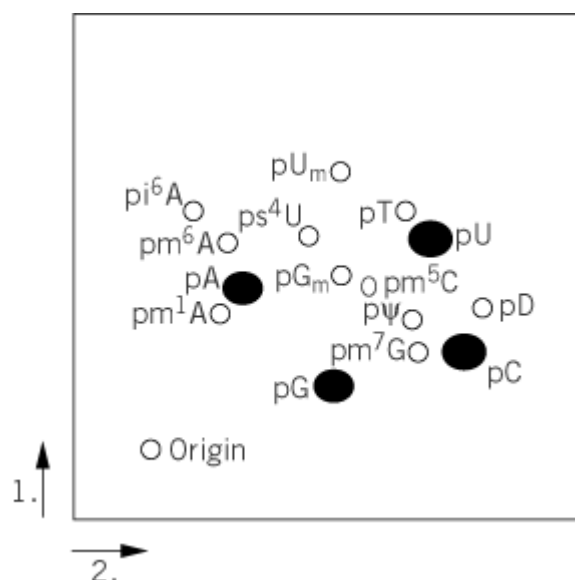
Sequencing of RNA by [primer extension](#) using AMV reverse transcriptase (11) is very similar to the common enzymatic **Sanger method** of DNA sequencing with chain-terminating inhibitors. A synthetic oligodeoxyribonucleotide complementary to part of the RNA to be sequenced is used as a primer for DNA synthesis in the presence of radioactively labeled nucleoside-5'-triphosphates and one of the four nucleoside-5'-triphosphate **dideoxynucleotide** analogues. The sequence is then read from the ladder produced by PAGE. More common, however, is transcription of the RNA to DNA, cloning, and determination of the sequence by standard DNA-sequencing technology. Reverse transcription is now used primarily for identification of positions occupied by bulky modified nucleotides, where the reverse transcriptase frequently terminates DNA synthesis. The corresponding positions are then detected as enhanced bands on the PAGE sequencing ladder. The [primer extension](#) method was successfully used to identify regions of RNA susceptible to chemical reagents (12) or to identify modified nucleotides in **ribosomal** RNA (13).

### 3. Identification of Modified Nucleotides in RNA

Rapid sequencing of RNA may provide hints to the presence of modified nucleotides in RNA, because the bands of oligonucleotides terminating there may be irregular, missing, or enhanced on the sequencing gels. This is due to the different reactivities of modified nucleotides to chemical reagents and ribonucleases. However, analysis by **thin-layer chromatography** (4) or **HPLC** (14)

and comparison with synthetic standards still is the most common way of detecting modified nucleotides in RNA (Fig. 3).

**Figure 3.** Identification of modified nucleoside-5'-phosphates by two-dimensional chromatography (10). RNA is digested by nuclease  $P_1$  and the nucleotides separated by first dimension: electrophoresis; second dimension: chromatography.  $i^6A$  denotes  $N^6$ -isopentenyladenosine;  $U_m$ , 2'-*O*-methyluridine;  $m^6A$ ,  $N^6$ -methyladenosine;  $s^4U$ , 4-thiouridine;  $m^1A$ , 1-methyladenosine;  $G_m$ , 2'-*O*-methylguanosine;  $m^5C$ , 5-methylcytidine,  $Y$ , pseudouridine;  $D$ , 5,6-dihydrouridine;  $m^7G$ , 7-methylguanosine;  $T$ , 5-methyluridine.



Considerable progress has been achieved in analysis of modified nucleotides by coupling liquid chromatography and electrospray ionization [mass spectrometry](#). This method was used successfully to analyze modified nucleosides in tRNA and ribosomal RNA (15, 16). In view of the rapid progress in instrumentation, the limits of this method have not yet been defined.

### Bibliography

1. J. A. McCloskey and P. F. Crain (1998) *Nucleic Acids Research* **26**, 198–200.
2. A. D. Branch, B. J. Benefeld and H. D. Robertson (1989) *Methods Enzymol* **180**, 130–154.
3. A. Diamond and B. Dudock (1983) *Methods Enzymology* **100**, 431–453.
4. Y. Kuchino, S. Nishimura (1989) *Methods Enzymology* **180**, 154–163.
5. T. E. England and O. C. Uhlenbeck (1978) *Nature* **275**, 560–561.
6. A. Simoncsits, G. G. Brownlee, R. S. Brown, J. R. Rubin and H. Guilley (1977) *Nature* **269**, 833–836.
7. H. Donis-Keller, A. M. Maxam and W. Gilbert (1977) *Nucleic Acids Research* **4**, 2527–2538.
8. D. H. Peattie (1979) *Proc. Natl. Acad. Sci. USA* **76**, 1760–1764.
9. J. Stanley and S. Vasilenko (1978) *Nature* **274**, 87–89.
10. Y. Kuchino, N. Hanyu and S. Nishimura (1987) *Methods Enzymology* **155**, 379–396.
11. C. S. Hahn, E. G. Straunss and J. H. Strauss (1989) *Methods Enzymology* **180**, 121–130.
12. S. Stern, D. Moazed, H. F. Noller (1988) *Methods Enzymology* **164**, 481–489.
13. A. Bakin and J. Ofengand (1993) *Biochemistry* **32**, 9754–9762.
14. C. W. Gehrke, K. C. Kuo, R. A. McCune, K. O. Gerhardt and P. F. Agris (1982) *J. Chromatogr.*

230, 297–308.

15. J. A. Kowalak, S. C. Pomerantz, P. F. Crain and J. A. McCloskey (1993) *Nucleic Acids Res.* **24**, 4577–4585.
16. J. A. Kowalak, E. Bruenger and J. A. McCloskey (1995) *Biol. Chem.* **270**, 17758–17764.

### Suggestions for Further Reading

17. S. M. Weissman, ed. (1983) *Methods of DNA and RNA Sequencing* (1983), Praeger Publishers, New York. (Useful compilation of the most common methods for RNA sequencing.)
18. *Modification and Editing of RNA* (1998) (H. Grosjean and R. Benne, eds.) ASM Press, Washington D.C. (Collection of articles dealing with the major problem in RNA sequencing, ie, modified nucleotides.)

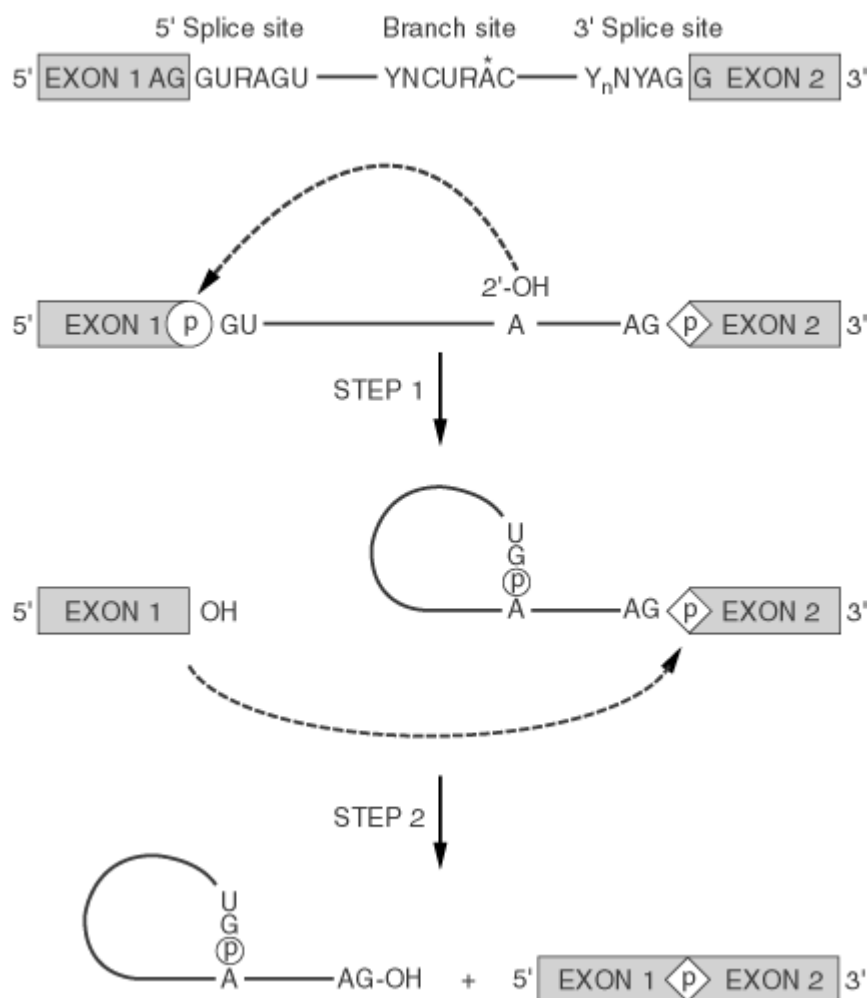
## RNA Splicing

Splicing of [messenger RNA](#) precursors is a critical step in the long chain of events required for the expression of most eukaryotic **genes**. The majority of these genes are transcribed as a large precursor mRNA (pre-mRNA) in which the coding regions are interrupted by intervening sequences (**introns**). In higher eukaryotes, nuclear pre-mRNAs typically contain multiple introns (in some cases greater than 50) with variable lengths extending up to 200,000 nucleotides. The conversion of pre-mRNA to functional mRNA requires the precise excision of these introns and the subsequent ligation of the coding sequences (exons), a process called *splicing*. Even the smallest of errors in the splicing reaction can have deleterious effects (for example, shifting the [reading frame](#) of the mRNA) that result in the production of altered, nonfunctional protein. The splicing of some pre-mRNAs is regulated during development or in a tissue-specific manner, such that **alternatively spliced** mRNAs, which are either nonfunctional or code for different protein isoforms, are generated from the same pre-mRNA. Alternative splicing thus represents an additional level of genetic regulation that enhances the genetic capacity of eukaryotes.

Much has been learned about the mechanism of nuclear pre-mRNA splicing through biochemical studies in both mammalian and yeast systems, as well as via genetic approaches in yeast. In particular, the development of an *in vitro* splicing assay has played a key role in identifying components involved in the splicing reaction and elucidating the splicing pathway. These studies have demonstrated that the basic chemical mechanism of splicing is conserved between higher and lower eukaryotes. The removal of an intron and subsequent ligation of the flanking exons is a two-step process that involves two temporally distinct transesterification (ie, phosphate transfer) reactions (see Fig. 1). The phosphodiester bonds of the pre-mRNA that are to be cleaved and then ligated are precisely defined by **consensus sequences** located around the 5' and 3' [splice sites](#). The first step of splicing is initiated by a nucleophilic attack on the phosphate at the 5' splice site by the 2'-OH group of an adenosine (designated the branch point), which is normally located 18 to 40 nucleotides upstream of the 3' splice site. This results in cleavage of the 3',5' phosphodiester bond at the 5' splice site and the concomitant formation of an unusual 2',5' phosphodiester bond between the first nucleotide of the intron and the branch site adenosine. This reaction produces two splicing intermediates, the 5' exon and a circular tailed molecule (ie, a so-called lariat) containing the intron and 3' exon. In the second step, the 3' OH of the 5' exon, which is released by the first transesterification reaction, carries out a nucleophilic attack on the phosphate at the 3' splice site. This results in excision of the intron in the form of a lariat and the simultaneous ligation of both

exons via a 3',5' phosphodiester bond. Despite the fact that pre-mRNA splicing requires the hydrolysis of ATP, the phosphates present in the products of the splicing reaction (ie, the spliced mRNA and excised intron) are derived from the pre-mRNA substrate and not from ATP (1-3). Splicing thus essentially involves the exchange of one pre-mRNA substituent for another on the phosphodiester bond at the 3' splice site.

**Figure 1.** Schematic representation of the two-step splicing pathway of nuclear pre-mRNA introns. Boxes and solid lines represent exon and intron sequences, respectively. The consensus sequences found at the mammalian 5' and 3' splice sites and branch site of U2-dependent introns are indicated, where N = any base, Y = pyrimidine, and R = purine. The branch site adenosine is marked with an asterisk, and the polypyrimidine tract is indicated by Y<sub>n</sub>. The nucleophilic attacks on the splice sites by the 2' OH of the branch site adenosine (step 1) and of the 3' OH of the cleaved 5' exon (step 2) are depicted by dashed arrows. The phosphate groups at the 5' and 3' splice sites, which are conserved in the splicing products, are also indicated.



Similar reaction intermediates and products are also observed during the splicing of group II introns, which are found in pre-mRNA molecules synthesized in plant and fungal organelles (4) (see [Self-Splicing Introns](#)). This has led to the hypothesis that group II and nuclear pre-mRNA introns may be evolutionarily related and that both types of splicing may involve similar catalytic mechanisms. Group II introns, in contrast to nuclear pre-mRNA introns, can be spliced autocatalytically *in vitro* in the absence of proteins or other factors. The self-splicing nature of these introns lies in their ability to fold into an elaborate, highly conserved intramolecular structure that favorably aligns the 5' and 3' splice sites and the branch site for the two cleavage/ligation reactions of splicing. Nuclear pre-

mRNA introns, on the other hand, possess only short conserved [cis-acting](#) sequence elements that are confined to the 5' and 3' splice sites, the polypyrimidine tract (present only in higher eukaryotes), and the branch site (see [Splice Sites](#)). The folding of nuclear pre-mRNA introns into a catalytically favorable conformation thus requires the presence of a large number of *trans*-acting factors. These include the [small nuclear RNP-\(snRNPs\)](#), evolutionarily highly conserved ribonucleoprotein (RNA-protein) complexes, and non-snRNP proteins. These factors interact stepwise with the pre-mRNA to form the [spliceosome](#), a large ribonucleoprotein complex that catalyzes both steps of splicing. Two distinct spliceosomes have to date been identified. The major or U2-dependent spliceosome is found in all eukaryotes and is responsible for the excision of so-called U2-dependent introns, which comprise the vast majority of nuclear pre-mRNA introns. The recently identified minor or U12-dependent spliceosome catalyzes the removal of the less abundant U12-dependent introns, which at present have been identified in only a subset of eukaryotes ([5-7](#)).

Despite recent advances in the splicing field, our understanding of the complex process of nuclear pre-mRNA splicing is far from complete. Currently little is known about the three-dimensional structure of the spliceosome and the nature of the spliceosomal [active sites](#) (i.e., whether the **catalyst** is RNA and/or protein) remains a matter of intense debate and investigation. One of two structurally distinct active sites appears to be responsible for each catalytic step of splicing ([8, 9](#)). During spliceosome assembly, the RNA components of the spliceosomal snRNPs (ie, the UsnRNAs) base-pair in a dynamic fashion with the short conserved regions of the pre-mRNA around the 5' and 3' splice sites, as well as with one another. This leads to the formation of a dynamic network of RNA–RNA interactions (see [Spliceosome](#)) that apparently provides the structural framework necessary for the catalysis of pre-mRNA splicing (reviewed in Refs. [10](#) and [11](#)). In fact, structural elements that closely resemble highly structured functional domains characteristic of self-splicing group II introns have been identified in the spliceosomal RNA network, supporting the long-standing hypothesis that nuclear pre-mRNA splicing is largely, if not exclusively, catalyzed by RNA ([12](#)). In addition to the snRNAs, however, proteins also play important roles in pre-mRNA splicing, particularly during the assembly of the spliceosome. Over 70 spliceosomal proteins have been identified thus far, and much effort is still currently invested in their structural and functional characterization. Proteins are involved in a number of [protein–protein interactions](#) and protein–RNA interactions that are required for the formation of a catalytically active spliceosome (see [Spliceosome](#)). In addition, several spliceosomal proteins possessing enzymatic activity (eg, RNA unwindase or protein isomerase activity), which are thought to catalyze conformational changes in the spliceosome during the splicing process, have also been identified (reviewed in Refs. [13](#) and [14](#)). More recently, proteins have even been proposed to contribute directly to the spliceosome's active sites. At present, a clearer understanding of how the splicing machinery functions awaits more detailed information about the nature and dynamics of the myriad of protein–protein and protein–RNA interactions formed within the spliceosome.

## Bibliography

1. R. A. Padgett, M. M. Konarska, P. J. Grabowski, S. F. Hardy, and P. A. Sharp (1984) *Science* **225**, 898–903.
2. M. M. Konarska, P. J. Grabowski, R. A. Padgett, and P. A. Sharp (1985) *Nature* **313**, 552–557.
3. R.-J. Lin, A. J. Newman, S.-C. Cheng, and J. Abelson (1985) *J. Biol. Chem.* **260**, 14780–14792.
4. A. Jacquier (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.) IRL Press, Oxford U.K., pp. 1–36.
5. S. L. Hall and R. A. Padgett (1996) *Science* **271**, 1716–1718.
6. W.-Y Tarn and J. A. Steitz (1996) *Cell* **84**, 801–811.
7. P. A. Sharp and C. B. Burge (1997) *Cell* **91**, 875–879.
8. M. J. Moore and P. A. Sharp (1993) *Nature* **365**, 364–368.
9. E. J. Sontheimer, S. Sun, and J. A. Piccirilli (1997) *Nature* **388**, 801–805.
10. H. D. Madhani and C. Guthrie (1994) *Annu. Rev. Genet.* **28**, 1–26.

11. T. W. Nilsen (1998) In *RNA Structure and Function* (R. W. Simons and M. Grunberg-Manago, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 279–307.
12. H. D. Madhani and C. Guthrie (1992) *Cell* **71**, 803–817.
13. C. L. Will and R. Lührmann (1997) *Curr. Opin. Cell Biol.* **9**, 320–328.
14. J. P. Staley and C. Guthrie (1998) *Cell* **92**, 315–326.

### Suggestions for Further Reading

15. M. J. Moore, C. C. Query, and P. A. Sharp (1993) "Splicing of Precursors to mRNA by the Spliceosome". In *The RNA World* (R. F. Gesteland and J. F. Atkins eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 303–358.
16. A. Krämer (1995) "The Biochemistry of Pre-mRNA Splicing". In *Pre-mRNA Processing* (A. I. Lamond, ed.), RG Landes Company, Austin, TX, pp. 35–64.

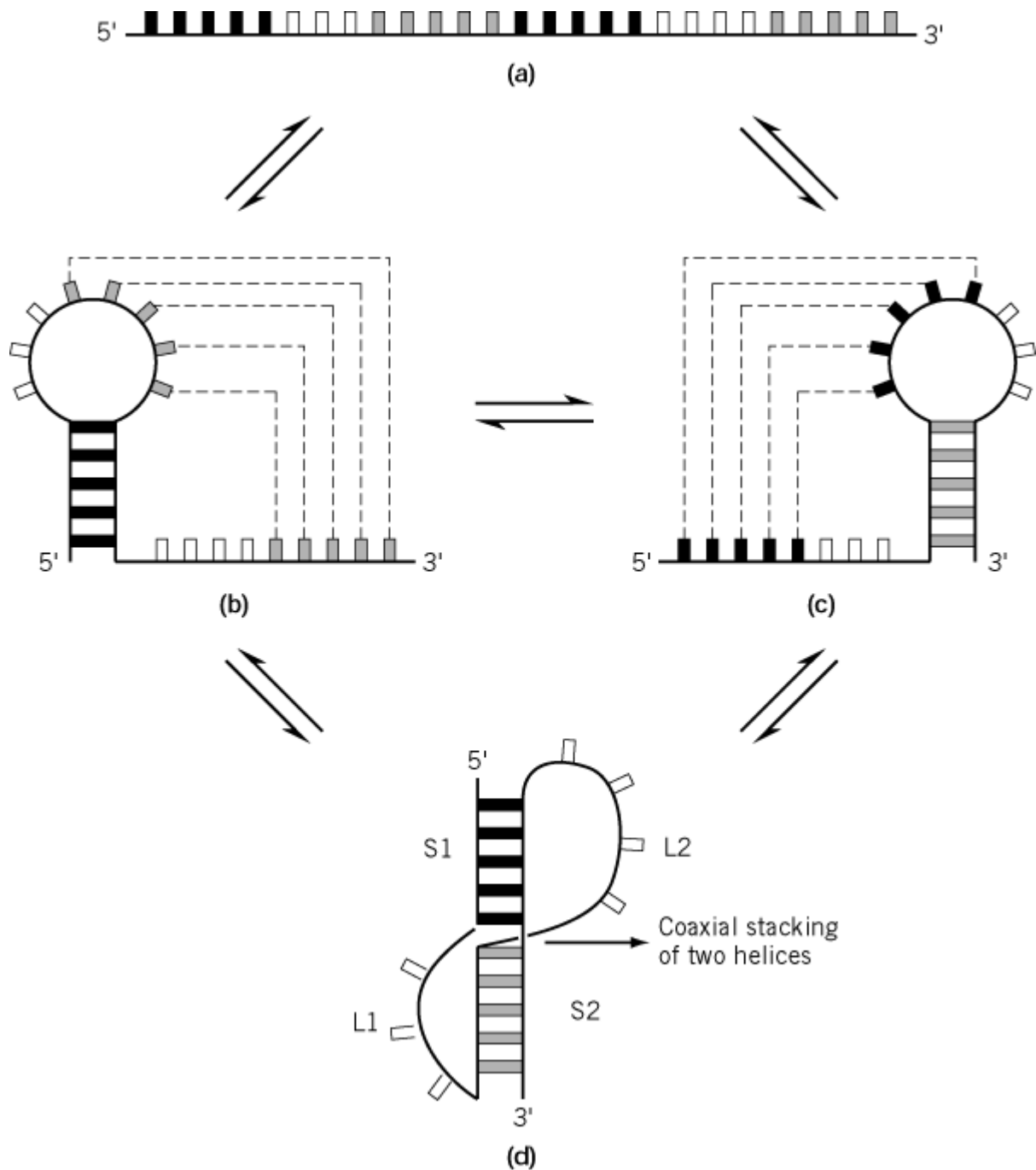
## RNA Structure

The diversity of function of RNA molecules in living organisms ranges from the **enzyme**-like activity of **ribozymes** to storage of genetic information in RNA **viruses**. RNA molecules adopt diverse structures in response to these functional requirements. RNA has a covalent structure very similar to that of **DNA**, the only differences being the change from a 2'-deoxyribose sugar in DNA to a ribose sugar in RNA and from a methyl group in thymidine to a hydrogen atom in uracil. However, their functional differences lead to correspondingly different structures. The requirement for storage of genetic information imposes the double-helical structure on most DNA molecules, whereas RNA molecules adopt an array of structures to rival [protein structures](#) in their complexity. This is not the result of an intrinsic limitation of DNA stereochemistry, but rather the result of different functional requirements.

It is convenient to describe RNA structure in hierarchical terms (Fig. 1), comparable to those used in describing protein structure: **primary**, **secondary**, **tertiary**, and **quaternary** structures. The primary structure refers to the sequence of an RNA molecule. Unlike proteins, which in most cases function only when properly folded, many RNAs function as unstructured, single-stranded species. For example, [messenger RNA](#) must be unfolded for the genetic message to be **translated**, and stable RNA secondary structures inhibit **protein biosynthesis**.

**Figure 1.** Hierarchy of RNA folding. **(a)** The primary structure corresponds to the RNA sequence. **(b, and c)** The secondary structure (in this case, two stem-loops or hairpins) form by Watson–Crick base pairing between complementary nucleotides. **(d)** Tertiary interactions (coaxial stacking of double helices) lead to the final three-dimensional structure (in this case, a pseudoknot).



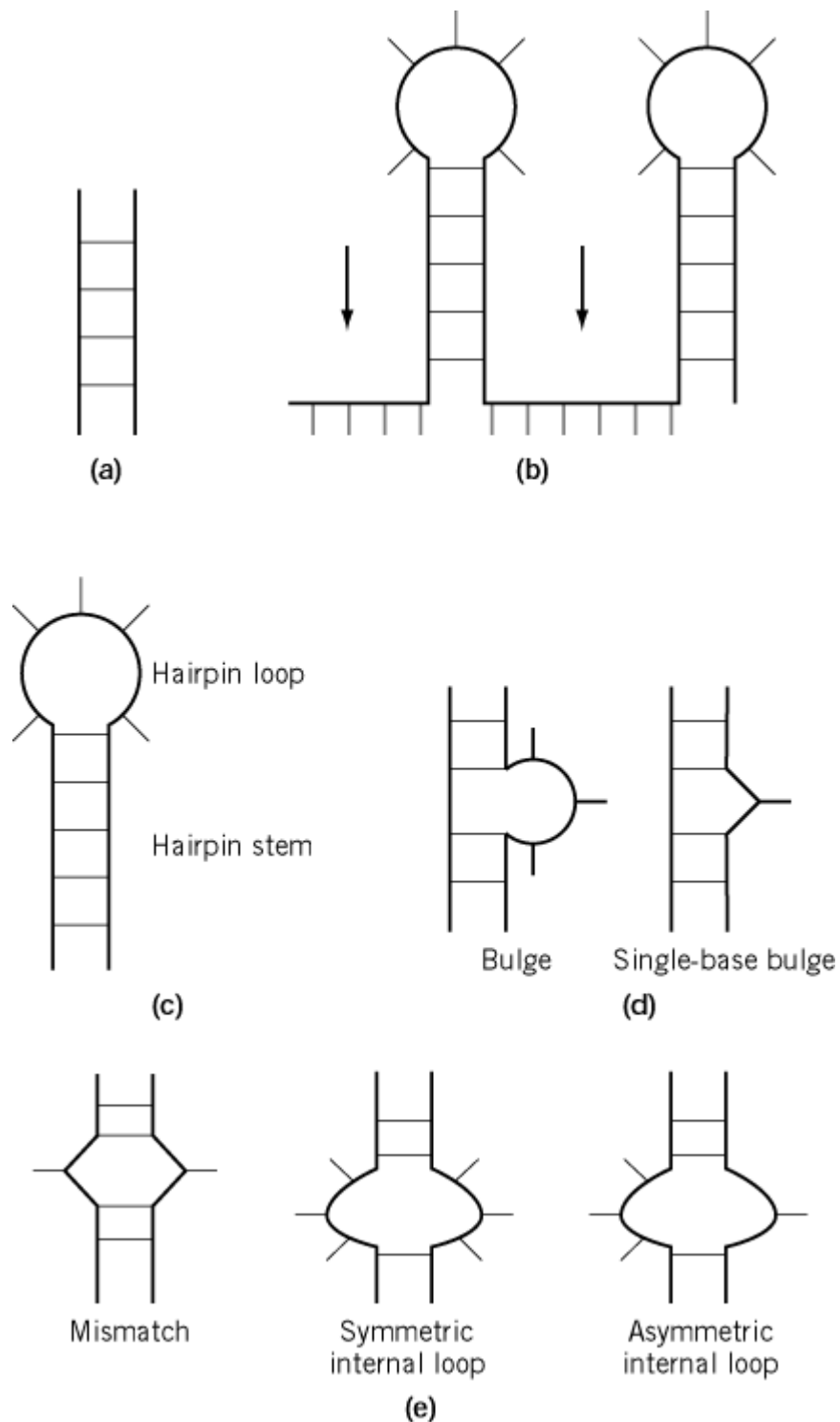


## 1. Secondary Structure of RNA

Secondary structure in DNA and RNA is dominated by Watson–Crick base pairing, leading to the formation of double-helical structures of varying length. Isolated base pairs are not thermodynamically stable, but formation of several consecutive base pairs readily occurs, resulting in a variety of possible arrangements (Fig. 2). These so-called secondary structure motifs represent the building blocks through which the most complex RNA three-dimensional structures are constructed. There is a fundamental difference between RNA and protein structures. Protein secondary structure is generally only marginally stable in the absence of stabilizing tertiary structure interactions, whereas RNA secondary structure is often stable on its own. Thus, formation of the secondary structure dominates the process of RNA folding, and RNA tertiary structure forms through relatively weak interactions between pre-formed secondary structure motives. Because of this property, RNA secondary structure can often be predicted successfully from **thermodynamics**

(1) or [phylogeny](#) (2).

**Figure 2.** RNA secondary structure motifs. (a) Duplexes; (b) Single-stranded regions; (c) hairpins; (d) bulges; (e) internal loops.



RNA secondary structure is dominated by the formation of double helices stabilized by Watson–Crick base pairs between complementary stretches. Unlike DNA, these helices are relatively short, seldom more than 8 to 10 base pairs in length, and are interrupted by single-stranded nucleotides

forming loop elements: hairpins, bulges, and internal loops (Fig. 2). These, together with the helical junction that is formed when more than two double helices come together, are the secondary structure motives and the building blocks upon which the most complex RNA structures are built.

In many RNAs, more than half of all nucleotides are incorporated into double-stranded helices. Duplex regions can form through very long-range interactions, and these interactions are crucial to determine and stabilize the overall fold of an RNA molecule. For example, opposite strands in double helices within ribosomal RNA can be separated by as many as 2000 nucleotides. In RNA, G·U pairs are almost as common as the canonical G·C or A·U base pairs, and they introduce slight distortions in double-helical structure that are recognized by proteins and other RNAs (3, 4). RNA (and DNA) double helices have an antiparallel right-handed helical conformation. RNA double helices adopt the A-form structure, which differs significantly from the canonical B-form adopted by DNA double helices (see [DNA Structure](#)). RNA duplex structures are not uniform, although their variability is less than that in DNA. These variations depend on the sequence and structural context of the helix relative to the global three-dimensional structure. A-form double helices differ from B-DNA in the conformation of the sugar and the displacement of the bases from the helical axis. These local differences lead to very diverse shapes, with profound consequences for recognition by proteins and other ligands (5).

The most common element of RNA secondary structure is the hairpin (or stem-loop) (Fig. 2). For example, bacterial 16S ribosomal RNAs contains approximately 30 phylogenetically conserved hairpins (2). A hairpin forms when the phosphodiester backbone folds back on itself to form a double-helical tract (called the *stem*), leaving unpaired nucleotides to form a single-stranded region, called the *loop*. Hairpin loops represent the most extensively studied RNA structure motif (apart from double helices). Small hairpin loops contain a high degree of structure, whereas longer loops (containing more than 7 to 8 unpaired nucleotides, such as those present in [transfer RNA](#)) are generally more poorly structured and thermodynamically less stable. Most loops in ribosomal RNA are small, 4 to 9 nucleotides in length, perhaps reflecting these thermodynamic preferences.

### 1.1. Tetraloops

Hairpins containing four nucleotides (tetraloops) are unusually common in cellular RNAs, and their sequences cluster within three exceptionally common families: UNCG, GNRA, and CUUG (where N is any of four nucleotides and R is a purine). For example, 70% of all tetraloops in 16S-like ribosomal RNAs from all organisms belong to the UNCG and GNRA families. The exceptionally common tetraloop structures have high thermodynamic stabilities, UNCG being the most stable. Capping a double-helical tract with a UNCG tetraloop is thermodynamically equivalent to extending the double-helical stem by two base pairs. Despite the differences in sequence, there are extensive structural similarities between these three families. In each of them, the first and last unpaired nucleotides form non-Watson–Crick base pairs to close the loop, reducing the number of unpaired nucleotides to only two. Longer loops, such as those found in tRNA, preserve extensive base stacking interactions that presumably stabilize the loop structure, but are generally characterized by significant conformational flexibility.

### 1.2. Bulges and Internal Loops

Bulges and internal loops form when two double-helical tracts are separated on either one (bulge) or both strands (internal loops) by one or more unpaired nucleotides. Internal loops containing equal numbers of bases on each strand are symmetric, whereas they are asymmetric when the number of bases are different. For example, single base mismatches are symmetric internal loops of two nucleotides. The presence of an internal loop or bulge reduces the thermodynamic stability, when compared to a perfect double helix, but unpaired nucleotides are more readily accessible to protein or nucleic acid ligands, which often recognize such sites. Non-Watson–Crick base pairs readily form within internal loops, while unpaired nucleotides within a bulge may stack within the helix or be bulged outside. The presence of an internal loop or bulge can induce bending in an RNA molecule; the extent of bending depends on the RNA sequence within the loop and can change upon ligand binding (6, 7). Thus, these motifs are ideal sites for conformational switches, where ligand binding

can result in long-range conformational changes.

## 2. Tertiary Structure of RNA

Interactions between two or more secondary structure elements give rise to RNA tertiary structure and define the overall folding of RNA molecules. In essentially all cases investigated thus far, RNA secondary structure elements maintain their three-dimensional structure even when extracted from very complex tertiary structures. Noncanonical base pairs, unpaired bases, and the backbone functional groups (the negatively charged phosphate groups and the unique 2'-hydroxyl group of RNA) are very important for tertiary interactions. Unpaired bases can twist or flip out of an helical patch to define unique surfaces for recognition by other RNAs during formation of the tertiary structure. The RNA secondary structure helps in orienting key residues into appropriate positions for tertiary interactions to occur. For example, the geometry of the four-way junction in tRNA, combined with a conserved length of the double-helical regions, help in positioning T-loop and D-loop nucleotides in close proximity to facilitate the tertiary loop-loop base pairs that define the L shape of all tRNAs (8).

Tertiary interactions often consist of base stacking and [hydrogen bond](#) interactions. Helical stacking between the terminal base pairs of two helices enables the building of the molecule into an extended helix (coaxial stacking) (Fig. 1). A common example of a structural module built upon base-stacking interactions is provided by the adenosine platform motif. Consecutive adenine bases within an internal loop can form a pseudo-A·A non-Watson-Crick base pair, to create a platform capable of mediating long-range tertiary interactions (9). Base triplets form when a preformed base pair becomes involved in another set of hydrogen bonds with a third base (8); this can occur on either the minor groove or major groove side of an RNA double helix. Hydrogen bonds between bases and backbones are observed when double helices are packed together in compact structures (10). Finally, divalent metal ions, especially hydrated magnesium ions, are often used to screen the negatively charged phosphate groups along the helical backbone, in order to build a compactly folded structure with close packing of negatively charged phosphates (11).

Unpaired nucleotides embedded within a secondary structure motif often form tertiary structures by interacting with unpaired nucleotides from another secondary structure module. These interactions can involve any secondary structure motif: hairpin loops, internal loops, and bulge loops. The nature of the tertiary contacts can be intercalation, base triplet formation, and Watson-Crick base pairing between complementary loop sequences. The loop-loop interaction module is best illustrated by the tertiary interaction of tRNA (8). The folding of tRNA is initiated by the formation of four hairpin loops via Watson-Crick base pairing, resulting in the organization of a four-way helical junction. Tertiary interactions, by helical stacking, further organize the tRNA into two extended helical domains. This arrangement positions all unpaired residues in the loop regions close in space and facilitates the formation of the tertiary base pairing and base triplets that lock the tRNA into the L-shaped three-dimensional structure.

### 2.1. Pseudoknots

Pseudoknots form when complementary primary sequences of a hairpin or internal loop and a single-stranded region interact with each other by Watson-Crick base pairing (12). When a pseudoknot forms between a hairpin loop and a complementary single-stranded region, then formation of two alternative hairpin structures can occur (Fig. 1). The formation of a pseudoknot creates an extended helical region through helical stacking of the hairpin double-helical stem and the newly formed loop-loop interaction helix. Although the pseudoknot is only marginally more stable than the two hairpins, tertiary interactions (such as base triplets) between unpaired nucleotides in the bridging loops and between base pairs within the extended helix can increase the stability of this structure.

### 2.2. Magnesium Ions

Double-helical regions must be packed together to build compact RNA structures, as found for example in the catalytic core of catalytic RNA ribozymes. This process is opposed by the strong

electrostatic repulsion between the negatively charged phosphate groups. RNA molecules overcome this repulsion through the direct or indirect coordination of divalent metal ions. In group I [self-splicing introns](#), a magnesium-organized ion core containing five magnesium ions constructs an exterior surface that facilitates the close packing of noncanonical loops into the minor groove of a double helix. Within the interior of this ion core, base-stacking and hydrogen-bonding interactions between nucleotides form specific metal ion binding sites ([10](#), [11](#)).

### 3. Quaternary Structure of RNA

There are relatively few well-characterized examples of the association of RNA molecules to form supramolecular quaternary structures, but these are relatively important. For example, during pre-mRNA splicing, messenger RNAs associate with five major [ribonucleoprotein \(RNP\)](#) particles called [small nuclear RNPs](#); such snRNPs interact with each other and with mRNAs. These interactions and their dynamic disruption and formation by means of RNA–RNA quaternary interactions are essential for [RNA splicing](#) to occur ([13](#)).

In most examples characterized thus far, the quaternary association of RNA molecules occurs by conventional Watson–Crick base pairing. For example, small regulatory RNAs with longer complementary sequences within RNA molecules ([antisense RNAs](#)) form intermolecular duplexes during the control of gene expression in both prokaryotes and eukaryotes ([14](#)). Similarly, [guide RNA](#) recognizes complementary sequences to identify sites where mRNAs are edited post-transcriptionally ([15](#)) (see [RNA Editing](#)). Although more complex RNA quaternary structures do not rely exclusively on Watson–Crick pairing, base pairs are still very important. So-called “kissing-hairpins” form between self-complementary loop nucleotides in two stem-loop structures ([16](#), [17](#)). These structures provide protein-recognition sites during the regulation of prokaryotic plasmid copy number, and possibly during the dimerization of the [HIV](#) genome. The best-characterized example of supramolecular association of RNA molecules by means of non-Watson–Crick interactions is provided by so-called “G-quartet structures” ([18](#)). These structures form readily *in vitro* for RNA and DNA sequences containing stretches of guanidines or uracils ([18](#)), but it is not clear whether these structures occur at all *in vivo*.

### 4. Summary

In conclusion, formation of Watson–Crick base pairs provide a simple structural code of RNA secondary structure formation, and helical pairing can be predicted very successfully by phylogeny and thermodynamics. As demonstrated for the first time by tRNA, tertiary interactions between the secondary structural elements fold RNA molecules into their three-dimensional structures. Helical regions generally define the RNA secondary structure. Tertiary contacts between secondary structure regions of the RNA, often involving nucleotides in single-stranded regions, fold the RNA into its three-dimensional structure. RNA quaternary structures are less well understood, yet relatively important in gene expression and its regulation.

### Bibliography

1. J. A. Jaeger, D. H. Turner, and M. Zuker (1989) Proc. Natl. Acad. Sci. USA **86**, 7706–7710.
2. R. R. Gutell, B. Weiser, C. R. Woese, and H. F. Noller (1985) Prog. Nucleic Acid Res. Mol. Biol. **32**, 155–216.
3. J. A. Doudna, B. P. Cormack, and J. W. Szostak (1989) Proc. Natl. Acad. Sci. USA **86**, 7402–7406.
4. K. Gabriel, J. Schneider, and W. H. McClain (1996) Science **271**, 195–197.
5. T. A. Steitz (1990) Q. Rev. Biophys. **23**, 205–280.
6. M. Zacharias and P. J. Hagerman (1995) Proc. Natl. Acad. Sci. USA, **92**, 6052–6056.
7. F.-H. T. Allain, C. C. Gubser, P. W. A. Howe, K. Nagai, D. Neuhaus, and G. Varani (1996) Nature **380**, 646–650.

8. J. E. Ladner, A. Jack, J. D. Robertus, R. S. Brown, D. Rhodes, B. F. C. Clark, and A. Klug (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4414–4418.
9. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, A. A. Szewczak, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1696–1699.
10. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1678–1685.
11. J. H. Cate, R. L. Hanna, and J. A. Doudna (1997) *Nature Struct. Biol.* **4**, 553–558.
12. D. E. Draper (1990) *Curr. Opin. Cell Biol.* **2**, 1099–1103.
13. J. A. Wise (1993) *Science* **1993**, **262**, 1978–1978.
14. M. Wickens and K. Takayama (1994) *Nature* **367**, 17–18.
15. J. Scott (1995) *Cell* **81**, 833–836.
16. K. Y. Chang and I. Tinoco Jr. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8705–8709.
17. J. P. Marino, R. S. J. Gregorian, G. Csankovszki, and D. M. Crothers (1995) *Science* **268**, 1448–1454.
18. C. Cheong and P. B. Moore (1992) *Biochemistry*, **31**, 8406–8414.

### Suggestion for Further Reading

19. I. Tinoco Jr. and J. R. Wyatt (1993) "RNA Structure". In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.

## RNA Structure Prediction

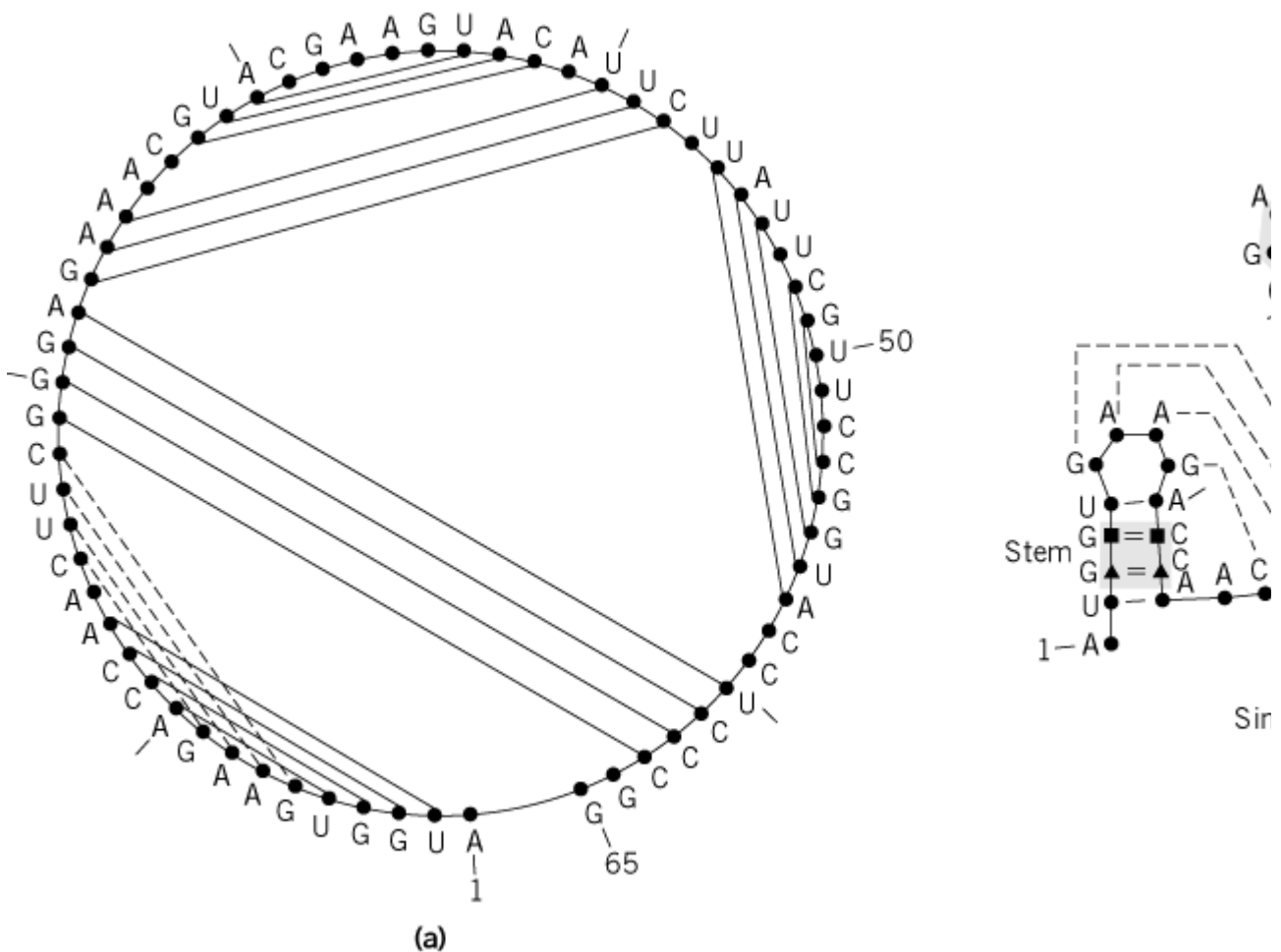
RNA molecules play many roles in the cell, and their biological functions are definitely not restricted to the relatively simple reading and decoding of the sequence of bases along the polynucleotide backbone that is exemplified by a simple view of [messenger RNA](#). On the contrary, many RNA molecules exert their biological activities through the chemical architectures they form, as **polypeptide chains** do (see [Protein Structure](#)). In aqueous solutions containing simple cations like sodium or magnesium ions, single-stranded RNA molecules, depending on their base sequence, either remain mostly unstructured (or **randomly coiled**) or they fold back on themselves to form complex three-dimensional (3-D) structures. The driving force for the folding is the stacking between the bases, which minimizes their exposure to [water](#), with the specific molecular recognition between RNA segments occurring overwhelmingly by Watson–Crick base-pairing mediated by [hydrogen bonds](#). Some RNA molecules are able to self-assemble into complex objects because they also contain additional tertiary contacts between segments of the polynucleotide chain (eg, the self-splicing group I **introns**). RNA molecules can also be observed in states containing only the **secondary-structure** double helices, without tertiary interactions, and they often need protein cofactors for folding into their biologically active conformations *in vivo*.

### 1. Partitioning Between Secondary and Tertiary Structures

The term “secondary structure” includes all segments that can build double-stranded helices by any combination of the isosteric Watson–Crick pairings, but this has some ambiguity. A secondary structure can be broken down into recurrent separable elementary motifs, such as the helical regions (stem structures, pseudoknots) and linking nonhelical elements (hairpins and internal loops, bulges and multiple junctions) (Fig. [1](#)). In secondary structure, a pseudoknot is a specific RNA motif that

results from standard Watson–Crick pairing between a single-stranded stretch that is located between two paired strands and a second, distal single-stranded region. The single-stranded regions may belong to a hairpin loop, an internal loop, or a 3' (or 5') dangling end, but at least one of them must occur between strands that form a double helix. When both single-stranded regions belong to a hairpin loop, they are said to form a *loop-loop* motif, which is formally equivalent to a pseudoknot. On the other hand, the two-dimensional (2-D) structure reduces the secondary structure to the set of Watson–Crick base pairs that form a planar graph (ie, without crossing edges) when the sequence of bases is arranged along a circle and the base pairs are connected by edges. Consequently, pseudoknots should be considered as belonging to the secondary structure, perhaps even the tertiary structure. This distinction is important, since the most efficient and most frequently used algorithms for predicting the two-dimensional structure do not take into account pseudoknots. At the next level of organization, that of the active tertiary structure, a three-dimensional architectural motif is a simple but recurrent arrangement containing a few secondary-structure elements that interact with a specific geometry and topology. The combination of such substructures leads to compact domains, which often fold autonomously and independently of the rest of the RNA architecture. A variety of observations support a view of RNA folding in which the three-dimensional architecture results from the cooperative compaction of separate, preformed, and stable substructures, which might undergo only minor and local rearrangement during the process. In summary, the secondary structure implies a local folding, whereas the tertiary structure refers to the global three-dimensional architecture of the RNA. The introduction of modular units, hierarchically organized and folded, circumvents most of the numerical nightmares inherent in the purely mathematical prediction of RNA structure.

**Figure 1.** (a) Secondary structure with the sequence arranged within a circle graph. A pseudoknot element (represented as a crossing edge) makes the graph nonplanar. (b) The more conventional representation shows basic structural elements. The free energy of the structure is the sum of the contributions of all the elements. According to Nussinov's first proposition (8), each structural element is a  $k$ -loop ended by a closing base pair (black triangles in the figure). In Figure 1b, five different  $k$ -loops are represented and labeled: a 0-loop (no interior pair); two stacked base pairs, an internal loop, and a bulge form 1-loops (one interior pair); and junction pairs is shown in the example). The black circles have no particular meaning.



Energetically, the secondary structure is the main component of RNA architecture, while tertiary structure contributes only slightly to the Gibbs **free energy** stability of the native state. Therefore, determination of the secondary structure is an essential step in the study of the structure-function relationships of an RNA molecule. The state in which an RNA molecule exists can be monitored by UV [absorption spectroscopy](#) as a function of temperature, RNA concentration, and the concentrations of ions. A molecule with definite secondary and tertiary structures will normally display two melting peaks: a first sharp low-temperature transition corresponding to the melting of the tertiary structure, and a second broad high-temperature transition to the melting of the secondary structure (1). In a complex and compact three-dimensional structure, the sugar-phosphate chain folds several times on itself and, by necessity, negatively charged phosphate groups come into close contact. In order to relieve the resulting electrostatic repulsions, positively charged cations are necessary. Biologically, the most prevalent and efficient cations are magnesium ions and polyamines. Magnesium cations are “hard” ions that interact favorably with the “hard” negatively charged oxygen atoms of the phosphate groups. Polyamines always carry positively charged [amino groups](#), and their flexibility and small size allow them to snuggle in helical grooves and in-between helical sugar-phosphate backbones. By monitoring the UV absorbance at low and high magnesium concentrations, one can distinguish the tertiary melting peak (it moves to still lower temperatures with decreasing concentrations of magnesium ions) from the secondary peak (it remains rather invariant). Native [gel electrophoresis](#) at various temperatures also distinguishes between folded and unfolded RNA molecules. Most important, electrophoretic methods are useful to ascertain the presence (and to evaluate the yields) of dimeric or higher oligomeric RNA species. The measurement of UV absorbance at various RNA concentrations allows one to calculate the concentrations of dimers and monomers and to measure the **melting temperature** of both dimers and monomers. Oligomerization may also depend on the concentrations of cations or polyamines.



## 2. Prediction of Secondary Structure

Computer programs exist that automatically produce a set of possible two-dimensional structures for a given RNA sequence. All are based on thermodynamic considerations, and the subsequent optimization of a set of criteria approximating the total free energy. In most cases, however, the selection of one secondary structure from the several that are usually produced by these programs and are indistinguishable in energy within the approximations used requires additional chemical or biological information stemming from experimental probing (biological approach) or comparisons of related sequences (**phylogenetic** or comparative approach). When a set of **homologous** sequences is available, which have common ancestry and function, one can search for a consensus core of secondary-structure elements, common to all the sequences, which should include a consensus three-dimensional architecture with the given function. In practice, one aims to organize the sequences so that the Watson–Crick paired regions align vertically, by searching for base covariations, ie, regions demonstrating compensatory base changes (eg, an A-U pair changing into a C-G pair) horizontally in all sequences. The more compensatory base change events there are in the sequences, the more firmly the secondary structure will be established. The efficiency of the comparative approach stems from the fact that molecular three-dimensional architectures evolve much more slowly than sequences. However, both the thermodynamic and phylogenetic methods are fraught with problems related to statistical relevance. With only four bases to choose among, purely coincidental compensatory base changes (or covariations between positions) are bound to occur. In the thermodynamic approach, they are mathematically resolved on the basis of the given set of thermodynamic parameters. Phylogenetically, the level of ambiguity can be reduced with additional sequences presenting additional covariations. For the establishment of a two-dimensional structure, and especially of a secondary structure, one should ideally employ all three of the approaches described above, with the caveat that the experimental data be gathered under conditions that favor a stable and functional structure. In practice, a comparison of foldings based solely on thermodynamics with the structures of 16S and 16S-like **ribosomal** RNA derived by the comparative approach shows that the quality of the thermodynamic predictions is variable, with percentages of correctly predicted base pairs ranging between 10 and 90% (2, 3).

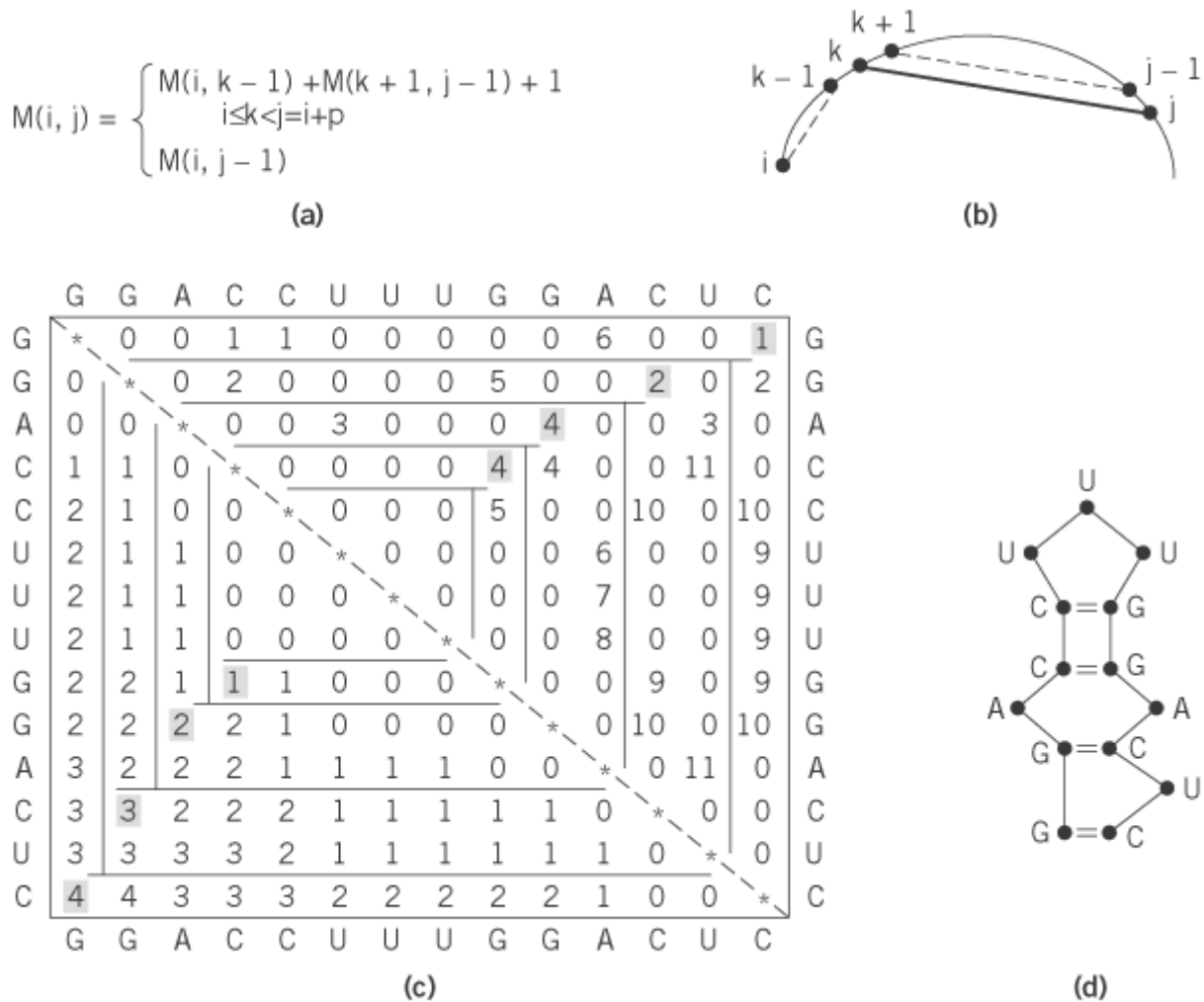
### 2.1. The Thermodynamic Approach

The thermodynamic approach is based on the empirical observation that the two-dimensional structure of an RNA molecule can be decomposed into elementary motifs that are identifiable and recurrent. It originates in the work of Tinoco and co-workers (4), who first generated from UV absorbance melting studies a set of thermodynamic parameters for the stability of structures formed by short oligonucleotides. Structural stability is measured by the decrease in free energy accompanying the transition from the unfolded or **denatured** state to the native state. The total decrease in free energy is equal to the sum of the independent contributions from each elementary motif present in the structure (the Tinoco–Uhlenbeck approximation), if we assume that tertiary interactions are weaker than secondary interactions and the sum of free energies of secondary elements is a reasonable approximation of the total free energy. The stability of a given base pair depends only on its immediate neighbors (nearest-neighbor approximation), considering that the stability is essentially related to base stacking and hydrogen bonding. The base pairs considered are the standard Watson–Crick pairs, G-C and A-U, as well as the often functionally important [wobble pair](#) G-U (5, 6).

When a single sequence is available, computational methods rely on the assumption that the native secondary structure is based on the two-dimensional structure with the lowest energy, or at least belongs to one of the suboptimal predicted two-dimensional structures (7). Stereochemistry requires that at least three nucleotides separate two paired strands. As mentioned above, an additional constraint is the lack of prediction of pseudoknots. Thus, because efficient algorithms are usually based on the decomposition into substructures (which is mathematically not possible in the presence of pseudoknots), one solution is to account for pseudoknots in a second step, either visually or by algorithmic fudging.

Efficient computer programs use  $O(N^3)$  recursive algorithms based on dynamic programming principles to fold sequences containing up to one thousand bases. A dynamic programming algorithm solves the problem by combining solutions corresponding to subproblems. It solves every subproblem just once and then saves the answers in a table, thereby avoiding the work of recomputing the answer every time the subproblem is encountered. The dynamic programming approach of RNA folding was first proposed by Nussinov and Jacobson (8) on the basis of the decomposition of a structure into base-paired structural elements (Fig. 2). An elegant decomposition into elementary structures proposed later by Zuker (Fig. 1) has made it possible to consider more complete thermodynamic data, at the price of increased computational time and storage requirements, depending on the way in which energy is assigned to loops (7). Both the free energies of the base pairs and additional experimental data are encoded in the energy function, resulting in the fast computation of the optimal two-dimensional structure.

**Figure 2.** Illustration of the principles underlying the dynamic programming approach. (a) The simplest recurrence proposed by Nussinov (8) for maximizing the number of base pairs. The same principle is applied by programs developed by Zuker (7), which include a more complex recurrence including loop consideration. (b) A schematic view of the principle of recurrence. For each increasing  $i, j$  subsection, the variable  $k$  is allowed to assume each position from  $i$  to  $j-1$  to test the ability of base  $k$  to pair with base  $j$ . (c) At each point, the total number of base pairs in the section  $i, j$  is computed. For each subsection, the maximum number of pairs that can be formed is saved in  $M(i, j)$ , and the value of  $k$  that yields this number is saved in  $M(j, i)$ . If  $j$  cannot pair with any  $k$  in the subsection,  $M(i, j) = M(i, j-1)$ . The maximum number of base pairs that can be formed in the folding of the example is given by reading  $M(14, 1)$ . (d) The secondary structure is obtained by reading in the upper half matrix partners that give the maximum number of pairs in the considered subsection.

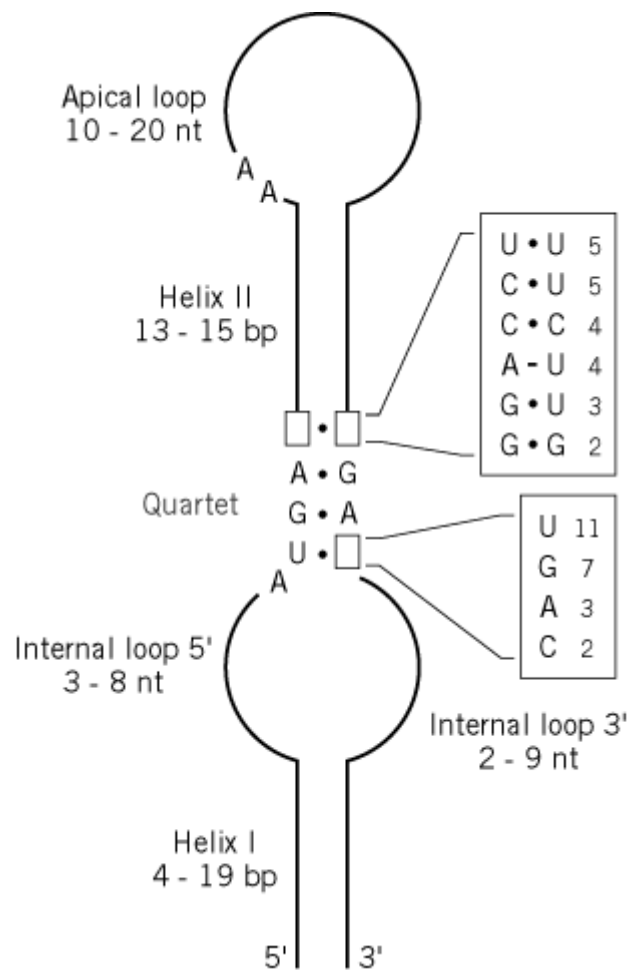


Nevertheless, the predictions depend strongly on the thermodynamic parameters used, which argues either for a lack of reliability of the optimization approaches or for incompleteness and imprecision in the energetic parameters. Although a better knowledge of the underlying thermodynamic parameters would certainly lead to better predictions, it is worth remembering the starting assumptions and limitations of the mathematical modeling. First, the Watson–Crick and wobble base pairs are not the only base–base interactions occurring between two RNA strands: Various A-A or A-G pairs, Hoogsteen [base pairs](#), or pyrimidine-pyrimidine pairings are frequently present in internal loops eg, loop E of eukaryotic 5 S rRNA ([9](#)) or the SECIS element in **selenocysteine**-coding mRNA ([10](#)) (see Figs. [3](#) and [4](#)) or hairpin loops (eg, the thymine loop of [transfer RNA](#) ([11](#)) or the -GNRA-tetraloops ([12](#))). The three-way junction of the hammerhead **ribozyme** is organized as well around a core of tandem sheared A-G pairs and a non-Watson–Crick A-U pair ([13](#)). With pseudoknots, one of the two constituting helices will, in a favorable situation, be predicted on the basis of its stability, and visual inspection of the remaining single-stranded regions normally points to the second helix. However, regions that should form noncanonical base pairs are not considered as such by the prediction programs, but as potential Watson–Crick base-pairing regions. Thus, depending on the relative weights of the various helices, such noncanonical regions will be involved in standard helices (or in “internal loops” penalized by positive free energy increments) leading to incorrect local (and sometimes global) predictions. The manual or automated analysis of energy [dot plots](#) ([14-16](#)) offers an alternative to test new folding models and to guide the RNA sequence folding according to available information. In short, those prediction programs should be primarily used for obtaining a solid framework on which further refinements can be applied.

**Figure 3.** Sequence alignments of mRNAs coding for various selenocysteine containing proteins (Gpx = glutathione peroxidase). The paired segments, denoted Helix I/I' and Helix II/II' corresponding to the 5' and 3' strands, are underlined. The invariant loops are highlighted (black on grey). At the base of helix II, the quartet made of four non-Watson–Crick pairs that constitute the Insertion Sequence (SECIS element) is boxed and highlighted (white on black). Gaps have been introduced to maximize structure elements and sequence similarities. Adapted from ([10](#)) with permission.

|                 |      | Helix I | Int.loop 5' |   | Helix II    |                 | Apical loop | Helix                        |
|-----------------|------|---------|-------------|---|-------------|-----------------|-------------|------------------------------|
|                 |      |         |             |   | Quartet     |                 |             |                              |
| Bovine cell.GPx | 699  | - UUUU  | --- GCCC    | A | <b>UGAA</b> | GGUGUUCCCUCU    | AA          | ACCUACG-----UGGAGGAAU        |
| Human cell.GPx  | 668  | - UUUC  | --- AUCU    | A | <b>UGAG</b> | GGUGUUUCCUCU    | AA          | ACCUACGA-----GGGAGGAAC       |
| Mouse cell.GPx  | 911  | - UCUU  | ---- CC     | A | <b>UGAU</b> | GGUGUUUCCUCU    | AA          | AUUUGCA-----CGGAGAAAC        |
| Rat cell.GPx    | 1011 | - UUUU  | ---- CC     | A | <b>UGAC</b> | GGUGUUUCCUCU    | AA          | AUUUACA-----UGGAGAAAC        |
| Bovine pl.GPx   | 1309 | - CAUC  | ---- GCC    | A | <b>UGAA</b> | GGAGGGGC- CCG   | AA          | G- CCCGCGUGGG-----CGG- GCCUC |
| Human pl.GPx    | 1342 | - GUCU  | ---- UC     | A | <b>UGAG</b> | GGAGGGGC- CC-   | AA          | AGCCCUUGUGGGC-----GGA- CCUC  |
| Mouse pl.GPx    | 1215 | - UGUC  | ---- UCU    | A | <b>UGAA</b> | GGAGGGGC- CCG   | AA          | G- CCCUUGUGGG-----CGG- GCCUC |
| Rat pl.GPx      | 1199 | - UGUC  | ---- UCC    | A | <b>UGAA</b> | GGAGGGGC- CCG   | AA          | G- CCCUUGUGGG-----CGG- GCCUC |
| Human GPx - GI  | 831  | - UCAC  | ---- AGA    | A | <b>UGAU</b> | GGCACC UU- CCU  | AA          | ACCCUCA-----UGG- GUGGU       |
| Rabbit GPx      | 653  | - UUUC  | --- AUCC    | A | <b>UGAG</b> | GGCGUUCCCCCG    | AA          | AACAAA-----UGGAGGAAC         |
| Pig heart PHGPx | 2712 | - GCAC  | ---- CC     | A | <b>UGAC</b> | AGU- CUGC- CUA  | AA          | AACCAGCC- CUGGU- GGG- GCAG-  |
| Rat PHGPX       | 720  | - GCAC  | ---- UC     | A | <b>UGAC</b> | GGU- CUGC- CUG  | AA          | AA- CAGCCC CUGGU- GGG- GCAG- |
| S.mansonii GPx  | 539  | - UUCU  | - CGCUAU    | A | <b>UGAC</b> | GAU GGCAAUCUC   | AA          | AUGUUCA-----UUGGUUGCC-       |
| Selenoprotein W | 367  | - CCGC  | ---- UUC    | A | <b>UGAC</b> | AGGAAGGA- CUG   | AA          | AUGUCUCAAAGACCUUG-UGGUCUUUC  |
| Hum.sel.P 1st   | 1445 | - UGCU  | ---- UUA    | A | <b>UGAG</b> | AAUAG- AAACGU   | AA          | ACUAUGACCUAG-----GGGUUUUCU   |
| Rat sel.P 1st   | 1467 | - UUAC  | --- AUUG    | A | <b>UGAG</b> | AACAG- AAACAU   | AA          | ACUAUGACCUAG-----GGGUUU- CU  |
| Hum.sel.P 2nd   | 1881 | - AUAG  | - - -UCA    | A | <b>UGAU</b> | GGUUUAAUAGGU    | AA          | ACCAA-----CCCUA- UAA         |
| Rat sel.P 2nd   | 1846 | - AUAA  | - - -UCA    | A | <b>UGAC</b> | GGUUUAAUAGAG    | AA          | ACUGAG-----UCCUA- UGA        |
| Rat DI type 1   | 1528 | - AUUU  | --- GUUU    | A | <b>UGAU</b> | GGUC- ACAGUGU   | AA          | AGUUCACAC-----AGCUGUGAC      |
| Hum DI type1    | 1732 | - AUUU  | --- GUUU    | A | <b>UGAU</b> | GGCC- ACAGCCU   | AA          | AGUACACAC-----GGCUGUGAC      |
| Rat DI type3    | 1684 | - CUGC  | ---- UG     | A | <b>UGAC</b> | GAACC- GCCUCU   | AA          | CUGGGCUUGACCAC-----GGGUCGGCU |
| Rat DI type3    | 1602 | - CCCC  | ACUGCUG     | A | <b>UGAC</b> | GAACU AU- CU CU | AA          | CUGGUCUUGACCAC-----GAGCUAGUU |
| Xenop. DI type3 | 1300 | - UGUU  | - -UGCAA    | A | <b>UGAC</b> | GACCGAUU- UUG   | AA          | AUGGUCUCACGGCCAA- AAACUCGUG  |

**Figure 4.** Schematic view of the secondary structure of the SECIS element corresponding to the sequence alignment of I tandem sheared G-A pairs [involving hydrogen bonds between N2(G) and N7(A), as well as between N3(G) and N6(A)] between pyrimidines. Adapted from [10](#) with permission.



The comparative approach is based on the assumption that the function and folding architecture have and, consequently, that a consensus secondary structure should be derivable by comparing RNA sequences of choice when a set of homologous sequences is available for RNA with the same biological function arranged in groups and subgroups (ideally of similar size), either according to the phylogenetic classification. The overall robustness of the approach increases with the diversity of the sequences and the number of them, whereas the accuracy of each prediction depends on the number of covariation events in each alignment. The first step in an alignment consists of establishing the paired regions along each sequence, which should be arranged with their lengths of the paired regions juxtapose vertically (see the example of Fig. 3). In a second step, the covarying bases can be highlighted by a vertical alignment with the inclusion of blanks or gaps in a fashion similar to multiple sequence alignments (see Figs. 3 and 4).

Finding the conserved core of the secondary structure is a difficult task that brings together, within a framework of aligning sequences and assessing covariation. The first task is usually done by hand, with the help of computer-aided similarity search (17, 18). The second task consists of searching for nucleotide interactions by measuring covariation at RNA positions in the alignment (19, 20). It is most effective when an alignment is available but, at times, requires readjusting the alignment. Computationally expensive approaches that attempt to find a conserved core by iteratively converging a cycle, combining the alignment search and secondary-structure assessment, have also been used. Usually, one sequence is aligned to a hypothetical secondary-structure model (21).

### 3. Prediction of Tertiary Structure

Although secondary structure is not yet sufficient to predict function, a proper understanding of the functional role of a macromolecule requires knowledge of its precise molecular organization in space. In the absence of experimental data, molecular modeling attempts to construct and propose a three-dimensional architecture for a mixture of theoretical and experimental data. Hence, prediction methods range from the most mathematically rigorous, solely on computer algorithms, to the most pragmatic and operational ones, in which insights come from

experiment. Modeling is best considered a heuristic tool that should help in the rationalization of experimental data and most important, should suggest new relations between the various components of the modeled structure. A three-dimensional model, mutagenesis of a macromolecule will, by necessity, be somewhat random, but can be informative. At best, mutagenesis experiments performed under such conditions will confirm the secondary structure of the molecule, since there is no tertiary model able to organize the data at a higher level. Such experiments are useful for bootstrapping a three-dimensional structure, which will serve as a framework for organizing existing mutagenesis experiments.

Construction of the tertiary structure of an RNA molecule always assumes and starts from a given set of secondary contacts. Tertiary contacts can be gained through chemical modifications (the importance of specific atomic positions and their protections cannot be explained solely by the secondary structure; see **Footprinting nucleic acids**), but they do not yield directly the partners if we assume a single conformer in solution), and most efficiently by careful analysis. In the case of the prediction of the secondary structure, the approaches can be divided into those relying on algorithms for automatic folding prediction and those relying on previously accumulated knowledge. One should weigh the advantages of mathematical objectivity and automation while, on the other hand, being aware of biased human decisions that include the weighing and integration of highly variable, diverse, and sometimes conflicting data. The [distance geometry](#) method (22) falls in the first category. Problems do occur when applying this method to structures with chiralities and for avoiding knots in the structures.

Another method (23) exploits a pseudoatom approach, with either one pseudoatom per helix or one per nucleotide. Appropriate potential functions have been developed. The use of spherical pseudoatoms leads to a loss of detail, and, most important, all fine interactions that control RNA folding are ignored. A third approach is a satisfaction algorithm, which searches conformational space so that, for a given set of input constraints (distances, angles, etc.), all possible models are produced (24). The manual approach involves the extensive use of X-ray crystallography and NMR structures. With developments in the production and purification of RNA, the number of structures published at an increasing pace. The structures can be used to extract the structure of a fragment of interest, which is then assembled manually on a computer graphics screen, using interactive modeling procedures, and restrained least-squares minimization, [molecular mechanics](#), or [molecular dynamics](#) programs (25). These methods imply some human judgments that ultimately depend on the available [database](#), as well as the stereochemical knowledge of the modular. However, the human mind can quickly grasp three-dimensional relations and take into account diverse experimental data. The solvent-accessible surface of the final model can be used to validate the structure against experimental reactivities of specific positions to chemical reagents (26).

#### 4. Perspectives

Several biologically important RNA families have been identified and their secondary and tertiary structures described above. In some cases of low conservation of bases along the sequence, the known sequence and the secondary (or even the tertiary) structure can be viewed as a biological signal to search for. To increase and refine the knowledge of a given signal or to identify genomic sequences as specific RNA families, tRNA is certainly the best example, since it is now possible to scan new genomic databases with pre-identified tRNA sequences (27). In the future, it will be necessary to scan genomic sequences rapidly to identify other RNA families. Some programs already offer a dedicated language to specify and search for such sequence-structure motifs. At the increasing pace at which three-dimensional RNA structures are produced, such developments should be useful. The secondary and tertiary structures of RNA sequences identified as functionally important in genomes.

#### Bibliography

1. P. Brion and E. Westhof (1997) *Ann. Rev. Biophys. Biomol. Struct.* **26**, 113–137.
2. D. A. M. Konings and R. R. Gutell (1995) *RNA* **1**, 559–574.
3. D. H. Mathews, T. C. Andre, J. Kim, D. H. Turner, and M. Zuker (1998) in *Molecular Modelling of Nucleic Acids* (eds. J. SantaLucia Jr., et al.), ACS Symposium Series 682, American Chemical Society, Washington, DC.
4. P. N. Borer, B. Dengler, I. Tinoco, and O. C. Uhlenbeck (1974) *J. Mol. Biol.* **86**, 843–853.
5. D. H. Turner and N. Sugimoto (1988) *Ann. Rev. Biophys. Chem.* **17**, 167–192.

6. M. J. Serra and D. H. Turner (1995) *Methods Enzymol.* **259**, 242–261.
7. M. Zuker (1989) *Science* **244**, 48–52.
8. R. Nussinov and A. B. Jacobson (1980) *Proc. Natl. Acad. Sci. USA* **77**, 6309–6313.
9. C. C. Correll, B. Freeborn, P. B. Moore, and T. A. Steitz (1997) *Cell* **91**, 705–712.
10. R. Walczak, E. Westhof, P. Carbon, and A. Krol (1996) *RNA* **2**, 367–379.
11. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer-Verlag, Berlin.
12. H. W. Pley, K. M. Flaherty, and D. B. McKay (1994) *Nature* **372**, 111–113.
13. H. W. Pley, K. M. Flaherty, and D. B. McKay (1994) *Nature* **372**, 68–74.
14. J. S. McCaskill (1990) *Biopolymers* **29**, 1105–1119.
15. M. Zuker and A. B. Jacobson (1995) *Nucl. Acids Res.* **23**, 2791–2798.
16. C. Gaspin and E. Westhof (1995) *J. Mol. Biol.* **254**, 163–174.
17. F. Corpet (1988) *Nucl. Acids Res.* **16**, 10881–10890.
18. J. D. Thompson, T. J. Gibson, F. Plewniak, F. Jeanmougin, and D. G. Higgins (1997) *Nucl. Acids Res.* **25**, 4909–4916.
19. D. K. Y. Chiu and T. Kolodziejczak (1991) *CABIOS* **7**, 347–352.
20. F. Michel and E. Westhof (1990) *J. Mol. Biol.* **216**, 585–610.
21. S. Eddy and R. Durbin (1994) *Nucl. Acids Res.* **22**, 2079–2088.
22. J. M. Hubbard and J. E. Hearst (1991) *Biochemistry* **30**, 5458–5465.
23. A. Malhotra, R. K. Z. Tan, and S. Harvey (1990) *Proc. Natl. Acad. Sci. USA* **87**, 1950–1955.
24. F. Major, D. Gautheret, and R. Cedergren (1993) *Proc. Natl. Acad. Sci. USA* **90**, 9408–9412.
25. E. Westhof, B. Masquida, and L. Jaeger (1996) *Folding & Design* **1**, R78–R88.
26. E. Westhof, P. Romby, P. J. Romaniuk, J. P. Ebel, C. Ehresmann, and B. Ehresmann (1989) *J. Mol. Biol.* **216**, 585–610.
27. B. Billoud, M. Kontic, and A. Viari (1996) *Nucl. Acids Res.* **24**, 1395–1403.

### Suggestions for Further Reading

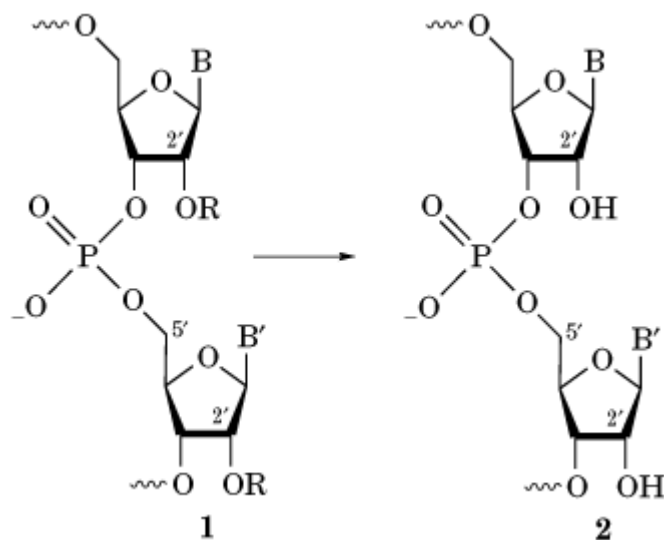
28. N. B. Leontis and J. SantaLucia Jr., eds. (1998) *Molecular Modelling of Nucleic Acids*, ACS Symposium Series 700, American Chemical Society, Washington, DC. The most recent and complete overview of molecular modeling algorithms of RNA as well as DNA.
29. F. Michel and M. Costa (1998) "Inferring RNA structure by phylogenetic and genetic analyses". (R. W. Symons and M. Grunberg-Manago, eds.), Cold Spring Harbor Lab, New York, Monograph 10. An artful description of the subtleties of comparative sequence analysis for RNA structure prediction.
30. P. Schuster, P. Stadler, and A. Renner (1997) *RNA structures and folding: From conventional to predictions*. *Curr. Opin. Struct. Biol.* **7**, 229–235. A quick overview of the recent advances in the field.
31. M. S. Waterman (1995) *Introduction to Computational Biology: Maps, Sequences, and Genome Analysis*. A remarkable book for the mathematically inclined reader.
32. E. Westhof, P. Auffinger, and C. Gaspin (1996) "DNA and RNA structure prediction". In *DNA and RNA: A Practical Approach* (M. J. Bishop and C. J. Rawlings, eds.), pp. 255–278, Oxford University Press. Includes with references to e-mail and web addresses.
33. M. Zuker (1989) "The use of dynamic programming algorithm in RNA secondary structure prediction". In *DNA Sequences* (M. S. Waterman ed.), pp. 159–184, CRC Press, Boca Raton, FL. A very useful presentation.

Until a few years ago, much more progress had been made in the chemical synthesis of oligo- and polydeoxyribonucleotides (DNA sequences) (see [DNA Synthesis](#)) than in oligo- and polyribonucleotides (RNA sequences). There are two main reasons for this. First, until fairly recently, there had been a very much greater demand for DNA than for RNA sequences in biological research. Secondly, due to the presence of the 2'-hydroxy functions, the chemical synthesis of RNA requires the use of an additional protecting group and is therefore inherently more complicated than the chemical synthesis of DNA. Nevertheless, in the past decade or so, significant progress has been made in synthesizing RNA sequences, in both solution and on a solid support. As in [DNA synthesis](#), the three main factors to be taken into account in synthesizing RNA sequences are (1) the choice of suitable protecting groups, (2) the development of phosphorylation procedures that are suitable for introducing the internucleotide linkages, and (3) the RNA sequences themselves.

### 1. Protecting the 2'-Hydroxy Functions

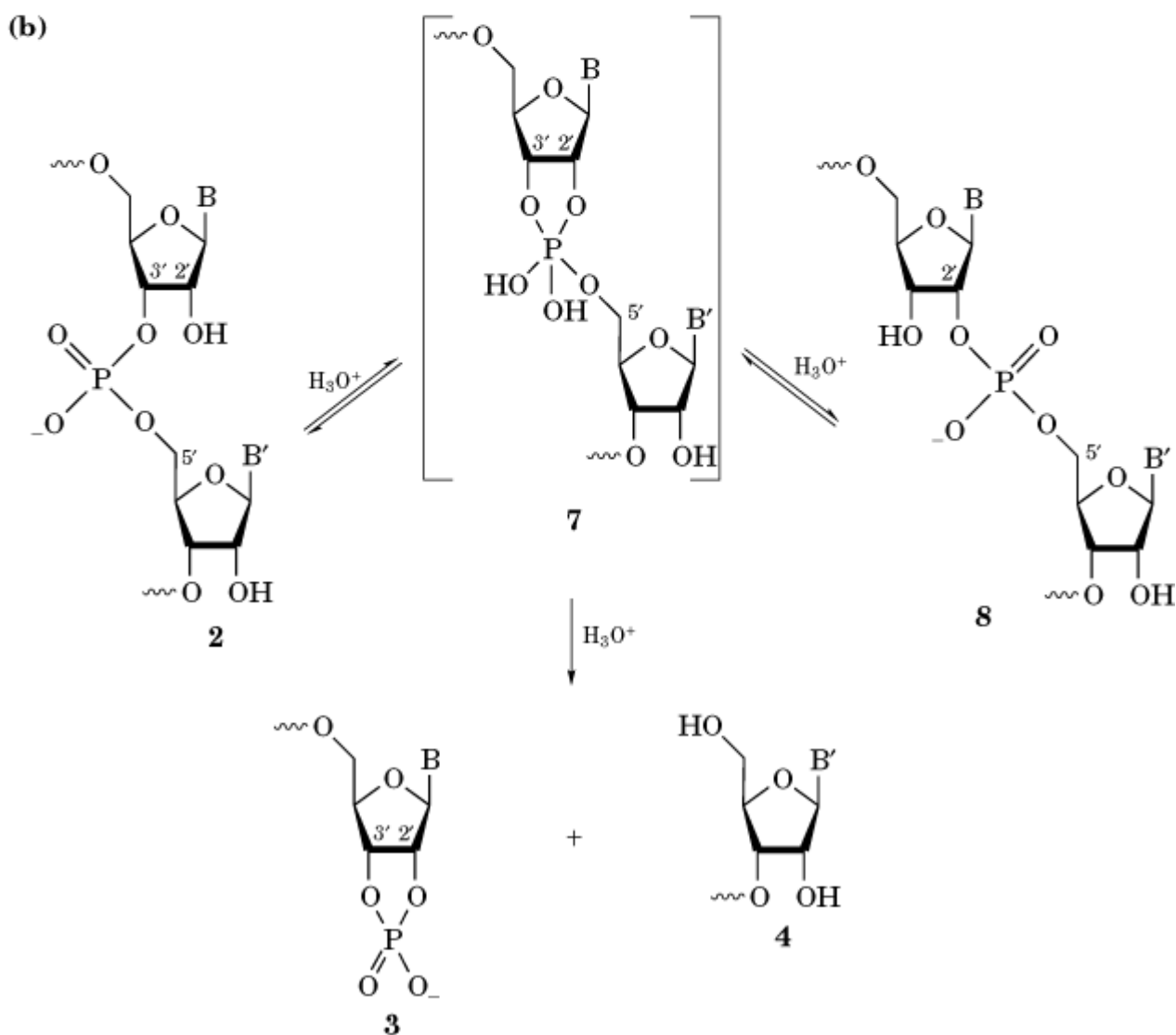
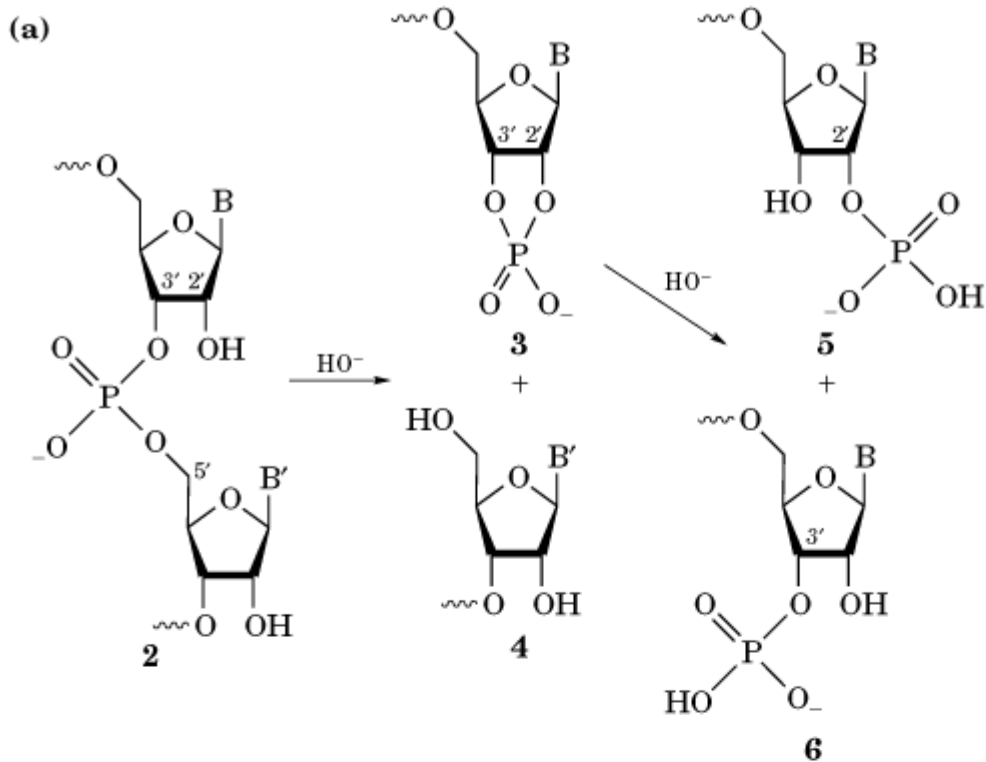
Almost certainly the most critical decision that has to be taken in the chemical synthesis of RNA sequences is the choice (1) of the protecting group (R, as in 1) for the 2'-hydroxy functions. This protecting group must remain intact until the final unblocking step (Fig. 1) at the end of the synthesis. It must then be possible to remove it under very mild conditions that do not promote the attack of the released 2'-hydroxy functions (as in 2) on the vicinal (3' → 5')-internucleotide phosphodiester linkages, thereby leading to their cleavage or migration (Fig. 2).

**Figure 1.** Final unblocking step.



**Figure 2.** (a) Cleavage of the internucleotide linkage under conditions of basic hydrolysis (b) Cleavage and migration of the internucleotide linkage under conditions of acidic hydrolysis.

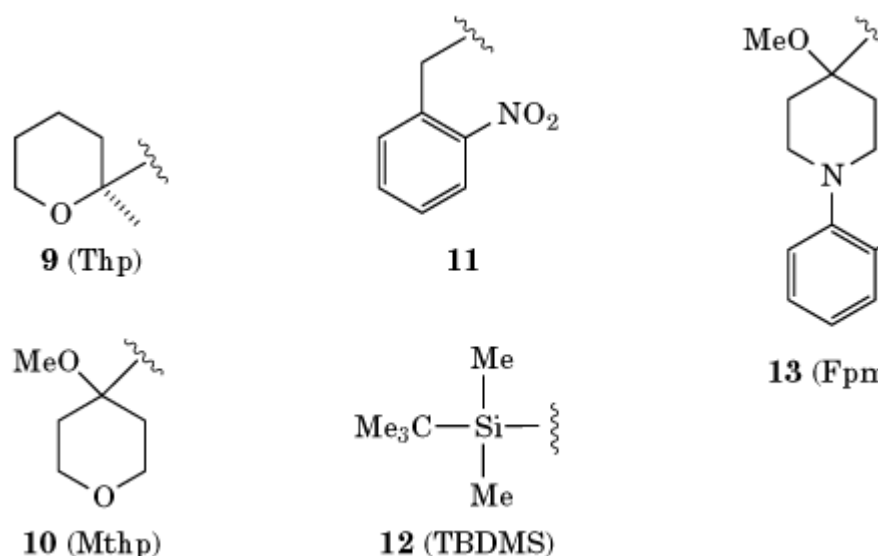




A number of the groups that have been used to protect the 2'-hydroxy functions are removed hydrolytically either under basic or acidic conditions. If the 2'-protecting groups are removed under basic conditions, cleavage but not migration of the internucleotide linkages can occur (2). This is illustrated in Figure 2a. Under basic conditions, the initially formed 2',3'-cyclic phosphate 3 undergoes further hydrolysis to give a mixture of the corresponding 2'- and 3'-phosphates (5 and 6, respectively). If the 2'-protecting groups are removed under acidic conditions, both cleavage and migration of the internucleotide linkages can occur (2). This is illustrated in Figure 2b. In the context of RNA synthesis, cleavage of the internucleotide linkages is clearly undesirable, in that it inevitably leads to diminished yields. However, phosphoryl migration, which gives rise to material with one or more (2' → 5')-internucleotide linkages (as in 8), is a much more serious matter. Even in the case of relatively low molecular weight oligoribonucleotides, it is virtually impossible to remove isomeric contaminants containing unnatural (2' → 5')-internucleotide linkages, and thereby obtain pure RNA sequences. Therefore, if acid-labile groups are used to protect the 2'-hydroxy functions in the synthesis of RNA sequences, such protecting groups must be removed under conditions of acidic hydrolysis that are mild enough to avoid the occurrence of phosphoryl migration.

Some of the groups that have been used most widely and indeed most successfully to protect the 2'-hydroxy functions in the chemical synthesis of RNA sequences are illustrated in Figure 3. Several of these protecting groups [ie, the tetrahydropyran-2-yl (Thp, 9), 4-methoxytetrahydropyran-4-yl (Mthp, 10), and 1-(2-fluorophenyl)-4-methoxypiperidin-4-yl (Fpmp, 13) groups] are acetal systems that undergo hydrolysis under mild conditions (especially in the case of the Fpmp 13 group) of acidic hydrolysis (13). These protecting groups have the considerable advantage that they are completely stable under the ammonolytic conditions that are needed to remove virtually all of the other protecting groups from fully assembled RNA sequences. The 2-nitrobenzyl protecting group 11 (4) is removed photolytically (at wavelengths > 280 nm) and is, of course, also stable to ammonolysis. The *tert*-butyldimethylsilyl (TBDMS 12) protecting group (5), which has been used very widely in the automated solid-phase synthesis of RNA sequences, has the advantage that it can be removed in the final unblocking step under nonacidic conditions. However, there are at least two distinct disadvantages associated with the use of the TBDMS protecting group. First, it is partially removed under the standard ammonolytic unblocking conditions (concentrated aqueous ammonia, 55°C). Second, as the TBDMS group readily undergoes base-catalyzed migration (6) from the 2'- to the 3'-hydroxy function of a ribonucleoside (and vice versa), particular care has to be taken in the preparation of the appropriate monomeric building blocks. Otherwise the synthetic RNA sequences will inevitably be contaminated with material containing (2'→5')-internucleotide linkages.

**Figure 3.** Some protecting groups for the 2'-hydroxy functions in RNA synthesis.

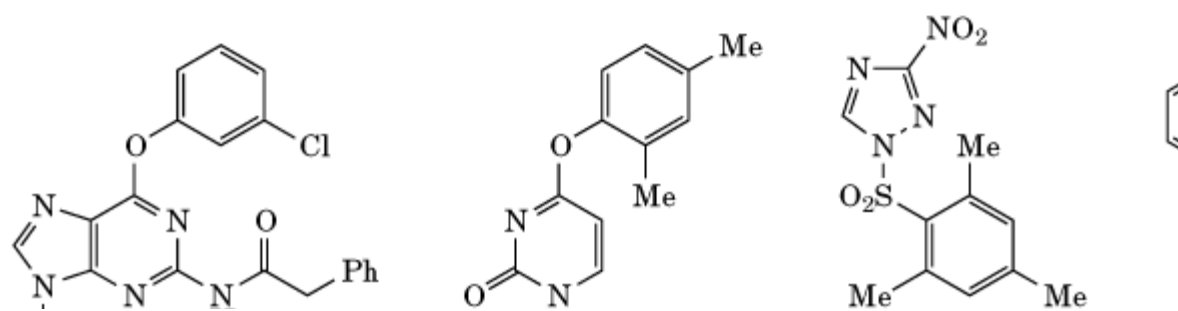
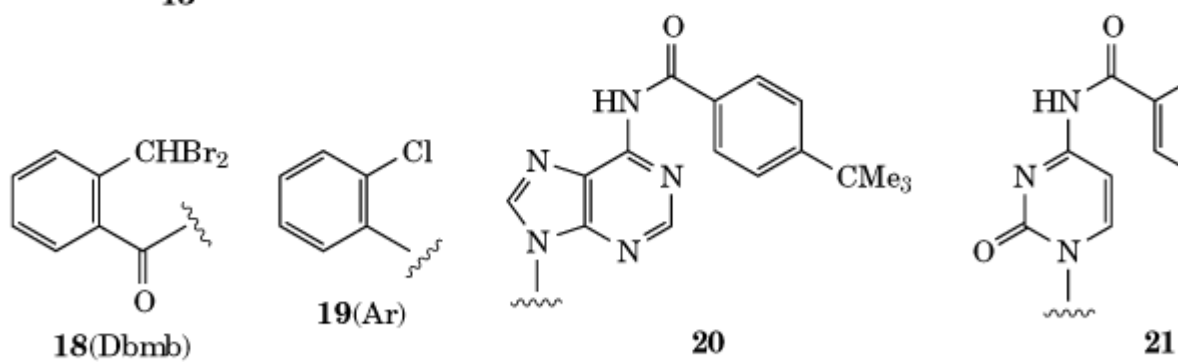
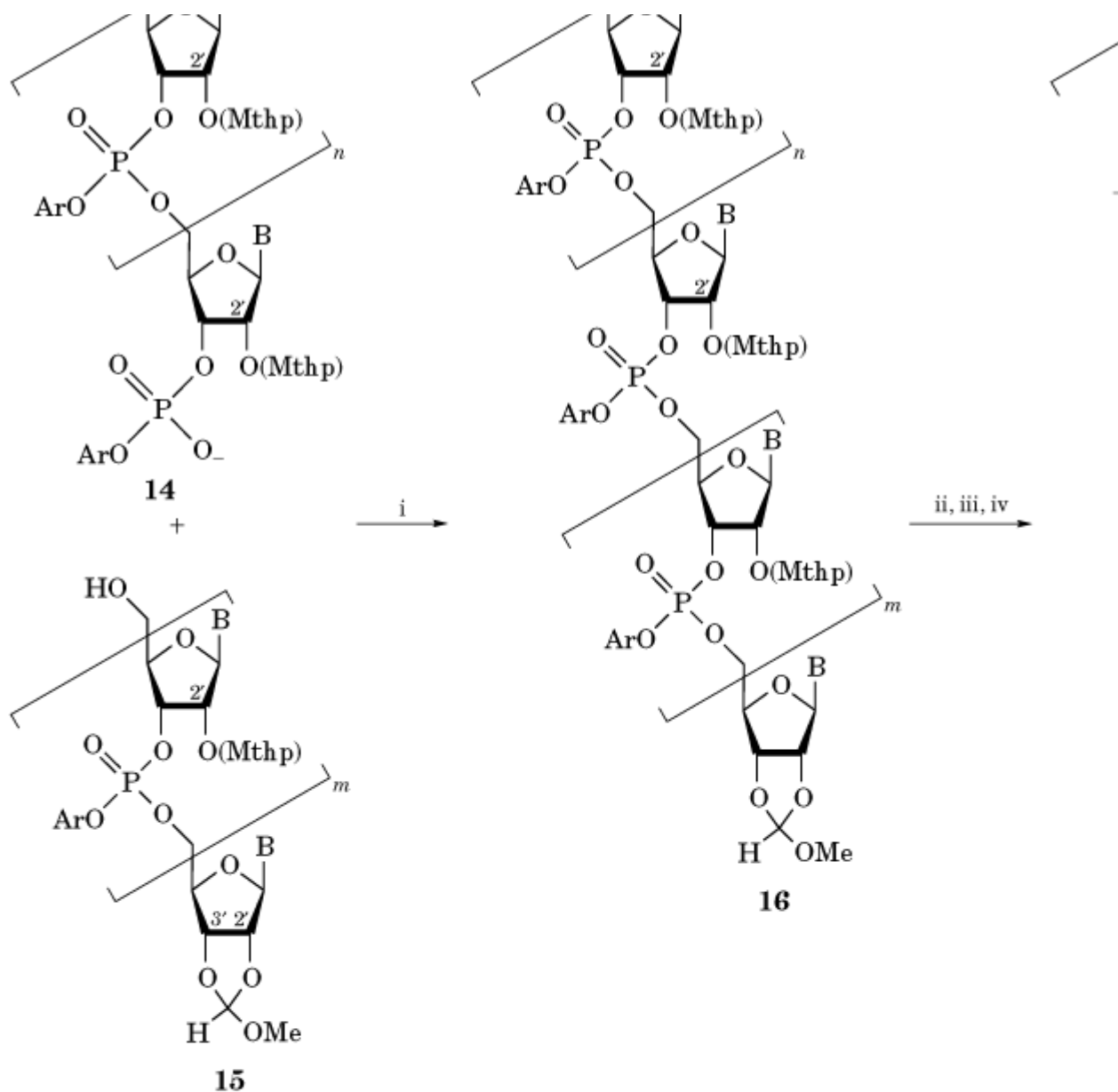


Apart from the particular problem associated with 2' protection, the principles governing the chemistry are the same as those governing the synthesis of DNA sequences. Similar protecting groups are used for the internucleotide linkages, and the same three principal phosphorylation methods (ie, the phosphotriester, phosphonate approaches) that are used in the [DNA synthesis](#) are also used in the synthesis of RNA sequences. In the case of [DNA synthesis](#), of these three phosphorylation methods only the phosphotriester approach ([1](#)) is used in the synthesis of RNA sequences, and the phosphoramidite approach ([8-11](#)) again appears to be the method of choice for the synthesis of RNA sequences.

## 2. Solution-Phase RNA Synthesis

Before a really satisfactory procedure for automated solid-phase synthesis had been developed, relatively long RNA sequences were prepared by the phosphotriester approach in solution. As an example of this approach used successfully ([7](#)) in the preparation of the 3'-terminal decamer, nonadecamer, and heptatriacontamer of tRNA<sup>Ala</sup> is illustrated in outline in Figure 4. The achiral Mthp group **10** was used to protect the 2'-hydroxyl group, and the methoxymethylene group was used to protect the 2',3'-terminal vicinal diol system (as in **15**). The 5'-terminal hydroxyl group (as in **14** and **16**) was protected with an acyl group [such as the 2-(dibromomethyl)benzoyl group], which is removable under what were effectively very mild basic conditions. Internucleotide linkages were prepared by the phosphotriester approach (Ar, **19**), and the adenine, cytosine, guanine, and uracil residues were protected as in **20**, **21**, **22**, and **23**, respectively. Uracil residues were protected on O-6 and O-4, respectively, with aryl groups to prevent their occurrence in subsequent coupling steps. Coupling (Fig. 4, step i) was effected with 1-(mesitylene-2-sulfonyl)-3-nitro-1,2,4-triazole in pyridine solution. After the fully protected RNA sequences had been assembled, unblocking was effected by an unblocking step (step ii), the 2-chlorophenyl protecting groups were removed from the internucleotide linkages, and the O-aryl groups were removed from the guanine and uracil residues by treatment with (*E*)-2-nitrobenzyltrimethylammonium tetramethyl-guanidine **26**. All of the acyl protecting groups were then removed by treatment with sodium hydroxide (step iii). Finally, the Mthp **10** and the 2',3'-terminal methoxymethylene protecting groups were removed by treatment with acidic hydrolysis. Solution-phase synthesis is labor-intensive and is likely to compete with automated solid-phase synthesis (see below) only if relatively large quantities of RNA sequences are required, say, for chemotherapeutic applications. The methoxymethylene group is more robust than the Mthp **10** group, and is also removable under milder conditions of acidic hydrolysis. In the case of solution-phase oligonucleotide synthesis it would be advisable to replace the Mthp by the Fpmp or another protecting group.

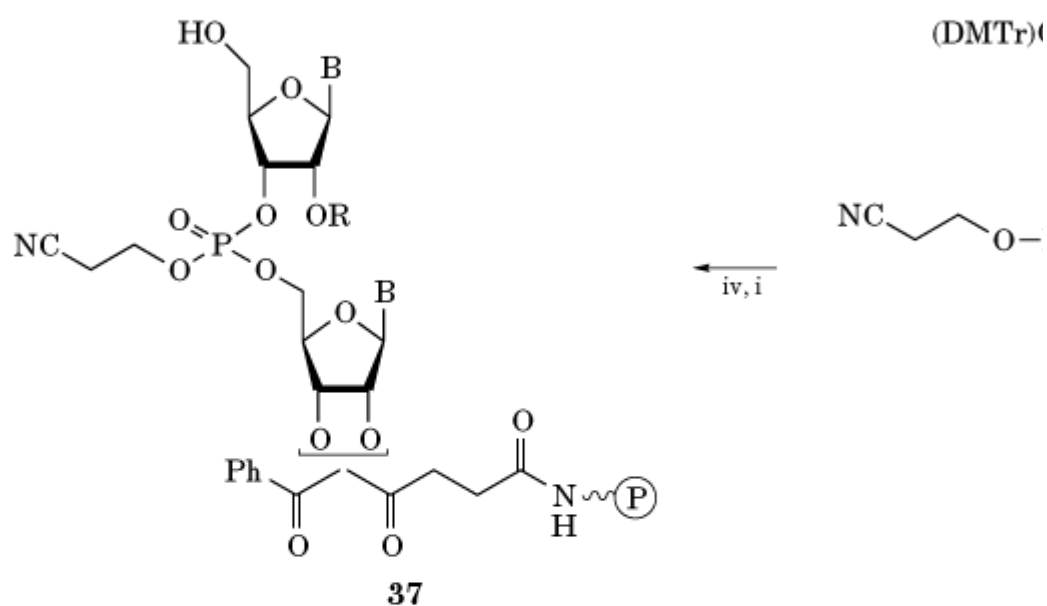
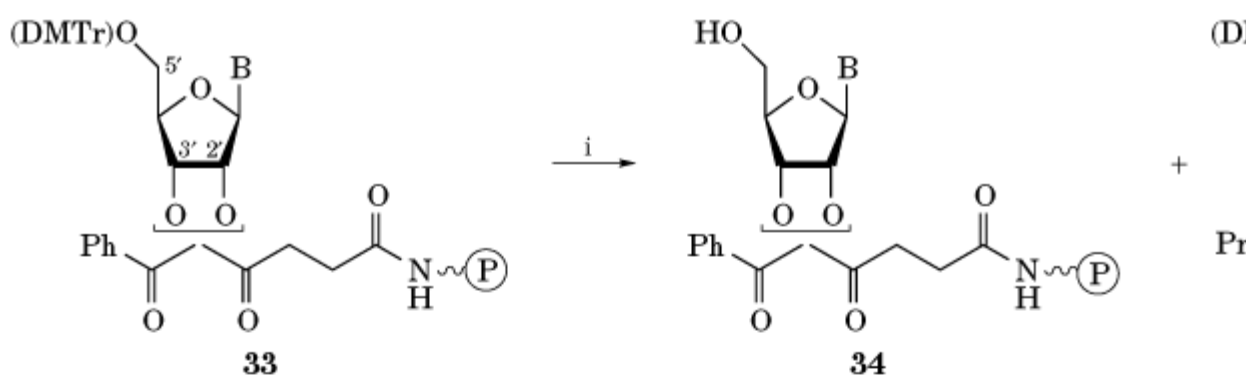
**Figure 4.** See text for description of steps i–iv.



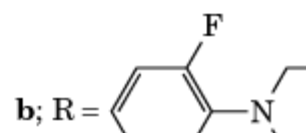
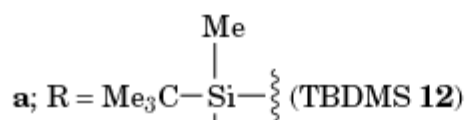
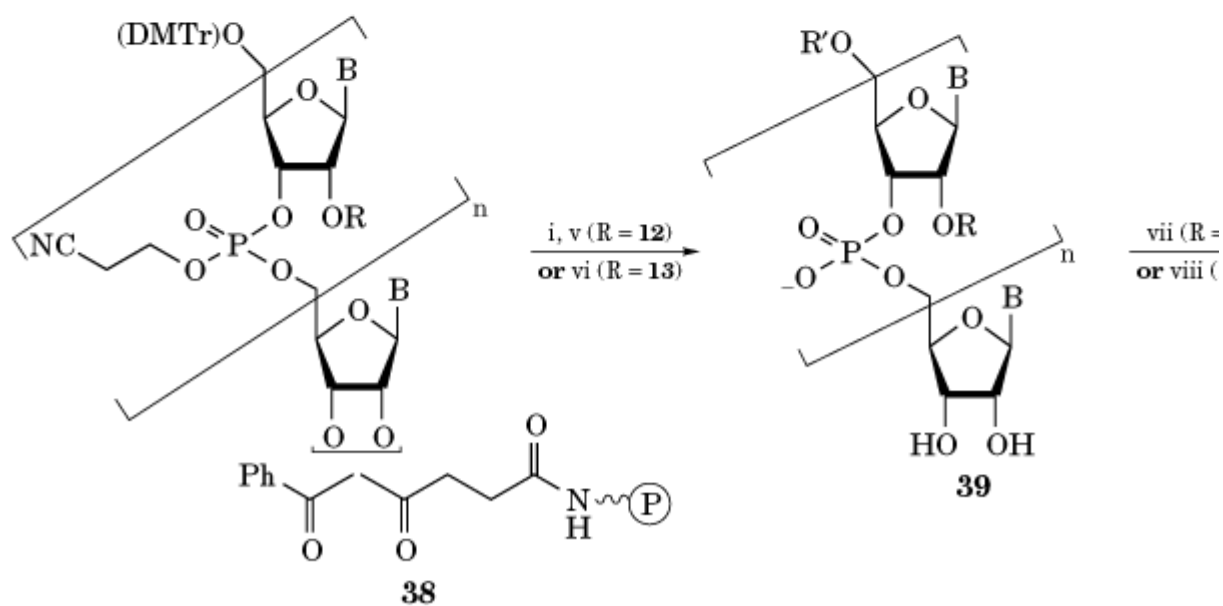
### 3. Solid-Phase RNA Synthesis

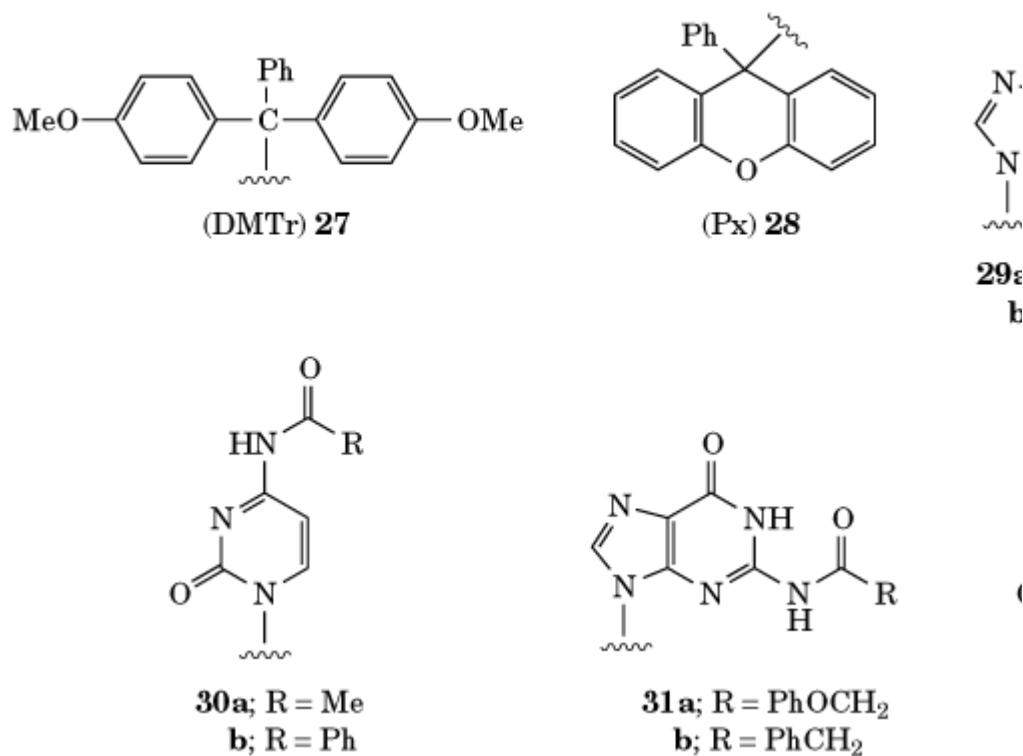
The advantages of automated solid-phase synthesis in the preparation of relatively small (ie, milligrams) are fully discussed in the early [DNA Synthesis](#). The same advantages apply to the automated solid-phase RNA synthesis. Indeed, apart from the final unblocking procedures, the protocols involved in solid-phase DNA and RNA synthesis are closely similar. Thus, the first step in the synthetic cycle of solid-phase RNA synthesis (8-11) is the attachment of a protecting group (Fig. 5a, step i). It is essential that this step should proceed both rapidly and quantitatively. The 4,4-dimethyl-2-tritylphenylmethyl (DMTr, **27**) and 9-phenylxanthen-9-yl (Px, **28**) ([11](#)), which are particularly acid-labile, are the 5'-protecting groups of choice. Such "trityl" groups have an additional advantage in that they can be estimated spectrophotometrically, and the efficiency of each coupling step can thereby be monitored. In the synthesis of DNA sequences, the adenine, cytosine, and guanine base residues are protected with *N*-acyl groups (as in **32**), and the uracil moieties are left unprotected (as in **32**). The base moieties of modified ribonucleosides are protected in an appropriate manner. Again, as in the solid-phase synthesis of DNA sequences, 2-cyanoethyl phosphoramidite blocks **35** are generally used, and the 3'-terminal ribonucleoside residue (as in **33**) is usually attached to the solid support [P, usually controlled-pore glass (CPG) or polystyrene]. The other secondary amine function is conveniently protected with an acyl (e. g., benzoyl) group, and it does not matter which one is actually linked to the solid support.

**Figure 5. Reagents:** i, 3%  $\text{Cl}_3\text{C}\cdot\text{CO}_2\text{H}$ ,  $\text{CH}_2\text{Cl}_2$ ; ii, 1*H*-tetrazole, MeCN; iii,  $\text{Ac}_2\text{O}$ , 2,6-lutidine, 1-methylimidazole, THF; iv,  $\text{NH}_3$ , EtOH– $\text{H}_2\text{O}$ ; v, conc. aq.  $\text{NH}_3$ , 55°C; vi,  $\text{Et}_3\text{N}\cdot 3\text{HF}$ ; viii, 0.5 *M* aq. NaOAc buffer (pH 3.25), 30°C. (P) = solid support



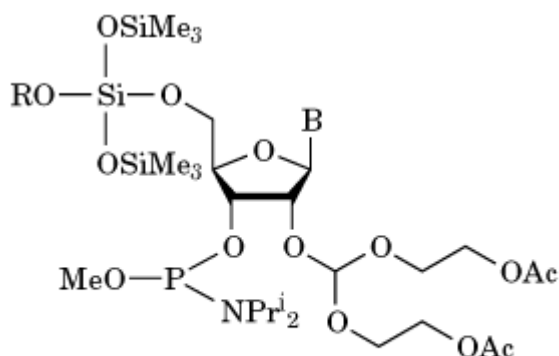
(b)





In solid-phase RNA synthesis, it is essential that the 2'-protecting groups should be completely stable during the "detritylation" steps at the beginning of each synthetic cycle (Fig. 5a, step i); it also should be completely stable under the ammonolysis conditions that lead to the detachment of the base from the solid support, the removal of the *N*-acyl protecting groups from the base moieties, and the removal of the phosphoramidite from the internucleotide linkages in the penultimate unblocking step (Fig. 5b, step v or vi). At present, the most suitable groups for the protection of the 2'-hydroxy functions in solid-phase RNA synthesis are phosphoramidite building blocks **35a** and **35b** are available commercially. In the case of TBDMS-protected phosphoramidites, adenine, cytosine, and guanine residues are now usually protected (9) with particularly base-sensitive groups [29a, 30a, and 31a, respectively]; in the case of the Fmp-protected phosphoramidites, adenine residues are protected (11) with more robust acyl groups (as in 29b, 30b, and 31b, respectively). The reasons for the consideration of the unblocking steps involved (Fig. 5b). The TBDMS-protected RNA sequences **38a** are normally treated with aqueous alcoholic ammonia (9) (step v) under mild conditions to remove the TBDMS-protecting groups. The resulting partially protected intermediates **39a**, R' = H, are then treated directly with triethylamine trihydrofluoride (step vii) (8, 9) or with tetra-*n*-butylammonium fluoride solution to remove all of the TBDMS-protecting groups and give the fully unblocked products **40**. The partially protected material **39b**, R' = DMTr, without any loss of the Fmp-protecting groups. This 2'-protecting group is completely stable to base-catalyzed hydrolysis and to any contaminant ribonucleases. A considerable advantage is that the latter material may be regarded as stabilized RNA and can readily be purified (ie, freed from the support) and then fully unblocked under mild conditions of acidic hydrolysis (step viii). If the final unblocking step is carried out under carefully controlled conditions (pH 3.25, 30°C), is not unnecessarily prolonged, cleavage and migration of the bases can be avoided. Relatively high molecular weight RNA sequences have been prepared successfully from **35a** and **35b**. The coupling rates observed are somewhat slower than for the corresponding DNA phosphoramidites.

efficiencies appear to be slightly lower. It is not yet clear which of the two 2'-protecting groups **12** a



**41**; R = cyclo-octyl or cyclododecyl

Finally, the increasing demand for synthetic RNA sequences will undoubtedly stimulate further research improvements in the methodology of both solid-phase and solution-phase synthesis. Indeed, it has rephase synthesis based on protected ribonucleoside phosphoramidites of general structure **41** leads to now forms the basis of a custom synthesis service.

#### Bibliography

1. C. B. Reese (2000) In *Current Protocols in Nucleic Acid Chemistry* (S. L. Beaucage, D. E. Berg eds.), Wiley, New York, pp. 2.2.1–2.2.24.
2. P. Järvinen, M. Oivanen, and H. Lönnberg (1991) *J. Org. chem.* **56**, 5396–5401.
3. C. B. Reese, M. V. Rao, H. T. Serafinowska, E. A. Thompson, and P. S. Yu (1991) *Nucleosides*
4. E. Ohtsuka, S. Tanaka, and M. Ikehara (1978) *J. Am. Chem. Soc.* **100**, 8210–8213.
5. K. K. Ogilvie, K. L. Sadana, E. A. Thompson, M. A. Quillian, and J. B. Westmore (1974) *Tetrahedron*
6. S. S. Jones and C. B. Reese (1979) *J. Chem. Soc., Perkin Trans.* **1**, 2762–2764.
7. J. M. Brown, C. Christodoulou, A. S. Modak, C. B. Reese, and H. T. Serafinowska (1989) *J. Chem. Soc., Perkin Trans.* **1**, 1767.
8. F. Wincott, A. DiRenzo, C. Shaffer, S. Grimm, D. Tracz, C. Workman, D. Sweedler, C. Gonzalez (1995) *Nucleic Acids Res.* **23**, 2677–2684.
9. B. Sproat, F. Colonna, B. Mullah, D. Tsou, A. Andrus, A. Hampel, and R. Vinayak (1995) *Nucleic Acids Res.* **23**, 2685–2694.
10. M. V. Rao and K. McFarlane (1995) *Nucleosides Nucleotides* **14**, 911–915.
11. M. V. Rao, C. B. Reese, V. Schehlmann, and P. S. Yu (1993) *J. Chem. Soc., Perkin Trans.* **1**, 43–52.
12. S. A. Scaringe, F. E. Wincott, and M. H. Caruthers (1998) *J. Am. Chem. Soc.* **120**, 11820–11821.

#### Suggestions for Further Reading

13. E. Ohtsuka and S. Iwai (1987) In *Synthesis and Applications of DNA and RNA* (S. A. Narang, ed) pp. 115–136.
14. C. B. Reese (1989) In *Nucleic Acids and Molecular Biology*, Vol. **3** (F. Eckstein and D. M. J. Li eds.), pp. 164–181.
15. S. L. Beaucage and M. H. Caruthers (1996) In *Bioorganic Chemistry: Nucleic Acids* (S. M. Hecht ed.), Wiley, New York, pp. 36–74.



## RNA World

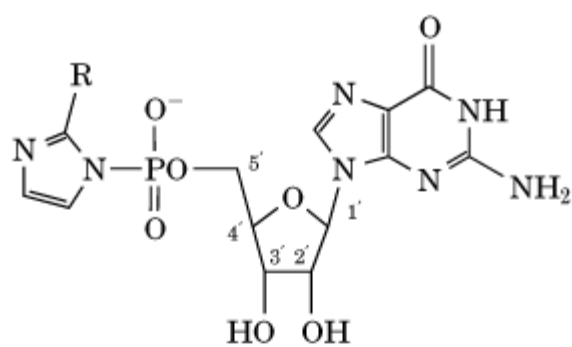
The “RNA World” is a term coined in 1986 by Walter Gilbert to describe the hypothesis that the origin of life proceeded through a stage that did not require DNA (1). The evidence from modern biochemistry that RNA is the more ancient nucleic acid obtains from many different inferences, including: 1) the 2'-deoxyribose nucleotides are synthesized enzymatically from the corresponding ribose precursors in modern biochemistry, 2) the DNA base thymine is formed by thymidylate synthase from the corresponding uracil nucleotide, and 3) RNA oligonucleotides function as primers in DNA biosynthesis. The molecular biology of DNA points repeatedly to the requirement for pre-existing RNA components, suggesting that the evolutionary origin of DNA depended on an ancient world in which RNA was the library of genetic information storage. In this view, DNA arose as a more stable vector by virtue of its superior resistance to self-cleavage, associated with the 2'-hydroxyl group, and the ability of thymine to circumvent the deleterious thermal and photochemical addition reactions typical of uracil.

A major boost to the supposed primacy of RNA was the discovery by Thomas Cech and Sidney Altman in the early 1980s that RNA could also serve a catalytic function as a so-called “ribozyme.” Beginning with the self-splicing of introns, the types of reactions promoted by ribozymes has expanded with the help of *in vitro* selection techniques to include various ligation processes that mimic the action of RNA polymerases (2). An important step toward the emergence of a true “replicase” has come with the creation of a ribozyme that catalyzes elongation of a primer sequence by 11 to 14 nucleotides, with fidelities from 92.0% (for the incorporation of adenosine) to 99.96% (for guanosine) in a mixture of nucleoside triphosphates (3). Not only can ribozymes promote the formation of peptide bonds (4), but the recently unveiled X-ray structure of the ribosome suggests that an adenine base within an RNA chain (instead of an amino acid in a protein) plays the primary catalytic role in the active site: thus, the ribosome is itself a ribozyme (5).

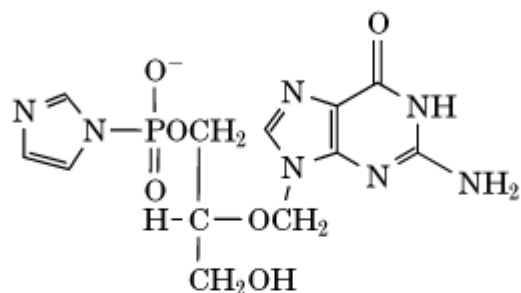
### 1. Oligomerization of Phosphorimidazolides

The origin of the RNA world has posed a formidable challenge. Self-replication must have preceded the evolution of ribozymes, and, thus, Leslie Orgel and his coworkers (6) have extensively studied the non-enzymatic template-directed oligomerization of activated nucleotides (eg, Ia in Fig. 1). Modern biochemistry employs nucleoside triphosphates, in which pyrophosphate serves as the leaving group, but these compounds react too slowly in laboratory simulations. A seminal discovery in 1980 was the zinc-catalyzed formation of oligoguanylates up to 40 nucleotides in length using a polycytidylate template; moreover, the internucleotide linkage in the presence of this metal ion was predominantly the “natural” 3'-5' bond (7). Activated precursors containing 2-methylimidazole (Ib in Fig. 1) gave even better results without the need for added metal ions; subsequent work in Orgel's group established the fidelity of template copying at better than 99% for the Watson-Crick base when present in a mixture with the other three nucleotide phosphorimidazolides (8). While the precursors were chosen for their chemical efficacy rather than their prebiotic relevance, the results established that template-directed replication of a polynucleotide did not require the presence of a complex protein polymerase.

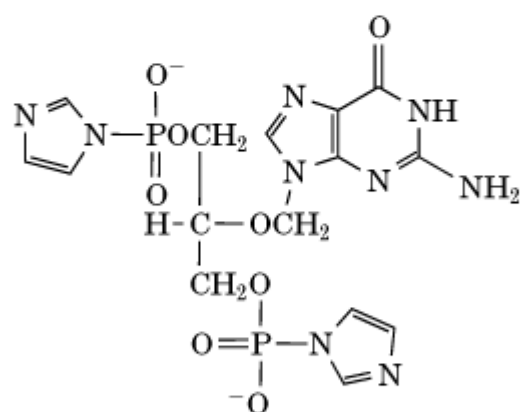
**Figure 1.** Activated guanine nucleosides and analogs: (a) guanosine 5'-phosphorimidazole (R = H) and guanosine 5'-phospho-2-methylimidazole (R = CH<sub>3</sub>); (b) 9-[(1-hydroxy-3-phosphorimidazole-2-propoxy)methyl]guanine; (c) 9-[(1,3-diphosphorimidazole-2-propoxy)methyl]guanine; (d) PNA guanine dimer.



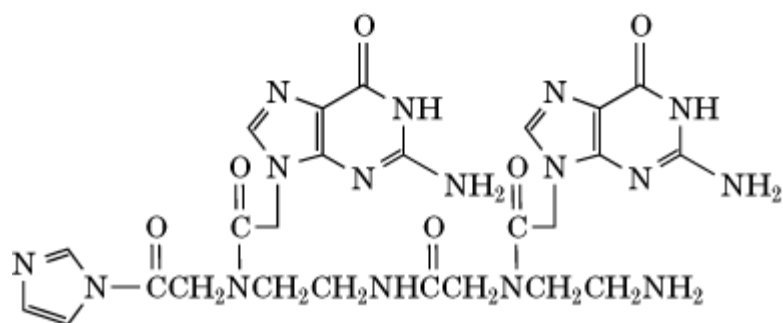
(a)



(b)



(c)



(d)

Despite the aesthetic appeal of this model system, Orgel and his associates have highlighted serious

limits to the templates and activated nucleotides that can be used. 2-Methylphosphorimidazolides prepared from a mixture of D- and L-guanosine do not oligomerize efficiently on a polycytidylate template due to chain termination by the “wrong” stereoisomer (9). Polyuridylates cause problems because of their tendency to form triple helices with adenine derivatives, while polyguanylates assemble into tetraplexes that inhibit replication (10). Mixed templates of guanosine and cytosine work as long as the amount of cytosine is high, but research on the deoxyribonucleotide series has shown that thymidine (and, by implication, uridine) effectively blocks the reaction (11). Because the goal of any self-replicating RNA model requires that the product also serve as a template, the problems of thymine (and uracil) arise whenever adenine is present in the strand being replicated. If these results with DNA analogs apply to their RNA counterparts, then the implications for the RNA world could be: 1) a smaller set of bases was employed early in the origins of replication (with severe limits on the amount of information storage); 2) a different suite of bases was used (with the attendant difficulties for the requisite transition); or 3) a natural catalyst (as yet unidentified) was necessary to achieve the conformation for thymine-containing templates to be replicated. Current research on the second and third alternatives may provide a solution to this dilemma.

Clay minerals on the early Earth could have played a role in the origin of the first oligonucleotide templates. In a remarkable series of papers by James P. Ferris and his colleagues (12), montmorillonite has been shown to promote the oligomerization of nucleoside phosphorimidazolid monomers: oligonucleotides up to 11 U long with a predominance of 3'-5' linkages have been obtained from the activated adenosine derivative. The chain length can be extended further (up to the 50-mer) by repeated feeding of the oligoadenylates with additional adenosine phosphorimidazolid (13). The pyrophosphate-linked diadenosine (AppA) appears to be an intermediate in the reaction, and its presence as a starting material augments the regioselectivity of the oligoadenylate product. Montmorillonite also promotes the self-condensation of uridine phosphorimidazoles (14) and cytidine phosphorimidazolides (15), but the linkages are mainly of the non-biological (2' to 5') type. An intriguing observation is that a mineral-bound template of exclusively 2',5'-oligocytidylate is a catalyst for the reaction of activated guanosine nucleotides (15). These results demonstrate that naturally occurring minerals cannot only accelerate such reactions but, in some cases, control the orientation of the monomers such that the products contain the biochemically important 3'-5' internucleotide bonds.

## 2. Pre-RNA Alternative Worlds

A major question in the study of the origins of life is what might have preceded the RNA world, and much attention has centered on simpler monomers that possess fewer stereocenters (or even none). The glycerol derivative (II) shown in Figure 1 is attractive because glycerol is a product of prebiotic simulations, and the proposed nucleotide analog has only one stereogenic carbon, but the activated phosphate undergoes a facile intramolecular cyclization that prevents template-directed polymerization (16). A related structural analog, the 1,3-diphosphoimidazolid (III), does react catalytically on a polycytidylate template to give pyrophosphate-linked oligomers up to the 20-mer (17).

Other models for the pre-RNA world have sought to avoid the problems associated with phosphate (low reactivity and dilute natural concentrations) through peptide-linked nucleic acids (PNAs). Although the specific monomer units employed by Nielsen and his colleagues (18) are not asserted to be prebiotic compounds, he and Orgel have demonstrated that the guanosine PNA dimer (IV) undergoes oligomerization up to the decamer on a decadeoxycytidine template. Interestingly, the reaction of the phosphoimidazolid (but not 2-methylimidazolid) of guanosine is promoted by a complementary cytosine-containing analog, thus showing that information transfer between peptide nucleic acid and RNA (in either direction) is feasible. However, studies of autocatalysis in oligopeptide ligation test the assumption that chains of amino acids by themselves cannot replicate (19, 20).

In conclusion, RNA oligonucleotides up to 50 base units long can be synthesized on clay mineral surfaces, and this size lies within the lower regime in which catalytic activity becomes feasible. Self-replication of polynucleotides is possible without the need for enzymes, but only cytidine-rich templates can be efficiently copied. The challenge is to find pathways to this RNA world, either with the discovery of new catalysts or through a genetic takeover from a nucleic acid analog that is truly “prebiotic.”

### Bibliography

1. W. Gilbert (1986) *Nature* **319**, 618.
2. M. C. Wright and G. F. Joyce (1997) *Science* **276**, 614–617.
3. W. K. Johnston, P. J. Unrau, M. S. Lawrence, M. E. Glasner, and D. P. Bartel (2001) *Science* **292**, 1319–1325.
4. B. Zhang and T. R. Cech (1997) *Nature* **390**, 96–100.
5. P. Nissen, J. Hansen, N. Ban, P. B. Moore, and T. A. Steitz (2000) *Science* **289**, 920–930.
6. L. E. Orgel (1995) *Acc. Chem. Res.* **28**, 109–118.
7. R. Lohrmann, P. K. Bridson, and L. E. Orgel (1980) *Science* **208**, 1464–1465.
8. T. Inoue and L. E. Orgel (1982) *J. Mol. Biol.* **162**, 204–217.
9. G. F. Joyce, G. M. Visser, C. A. A. van Boeckel, J. H. Van Bloom, L. E. Orgel, and J. van Wastrenen (1984) *Nature* **310**, 602–604.
10. G. F. Joyce and L. E. Orgel (1986) *J. Mol. Biol.* **188**, 433–441.
11. A. R. Hill Jr., L. E. Orgel, and T. Wu (1993) *Biosphere* **23**, 285–290.
12. J. P. Ferris and G. Ertem (1992) *Science* **257**, 1387–1388.
13. J. P. Ferris, A. R. Hill Jr., R. Liu, and L. E. Orgel (1996) *Nature* **381**, 59–61.
14. P. Z. Ding, K. Kawamura, and J. P. Ferris (1996) *Biosphere* **26**, 151–171.
15. G. Ertem and J. P. Ferris (1997) *J. Am. Chem. Soc.* **119**, 7197–7210.
16. A. R. Hill Jr., L. D. Nord, L. E. Orgel, and R. K. Robins (1988) *J. Mol. Evol.* **28**, 170–171.
17. J. Visscher and A. W. Schwartz (1988) *J. Mol. Evol.* **28**, 3–6.
18. C. Böhler, P. E. Nielsen, and L. E. Orgel (1995) *Nature* **376**, 578–581.
19. D. H. Lee, J. R. Granja, J. A. Martinez, K. Severin, and M. R. Ghadiri (1996) *Nature* **382**, 525–527.
20. S. Yao, I. Ghosh, R. Zutshi, and J. Chmielewski (1998) *Nature* **396**, 447–450.

### Suggestions for Further Reading

21. D. P. Bartel and P. J. Unrau (1999) Constructing an RNA World. *Trends Genet* **15**, M9–M13.
22. J. P. Ferris (1993) Catalysis and Prebiotic RNA Synthesis. *Origins Life Evol. Biosphere* **23**, 307–15.
23. R. F. Gesteland, T. R. Cech, and J. F. Atkins, eds. (1999) *The RNA World*, Cold Spring Harbor Press, Plainview, New York.

### RNA-Binding Proteins

**RNA** is structurally more versatile than **DNA** and has more diverse biological roles. It is an informational molecule that serves as genetic material and the mediator of genetic information from

DNA to **protein**, it also plays a vital structural role in many [ribonucleoprotein](#) particles and, as is now recognized, can itself have catalytic activity. Whereas in the hypothetical [RNA world](#), RNA would have acted unaided by other macromolecules, RNA almost invariably functions in modern organisms associated with RNA-binding proteins, a wide variety of which have now been identified. Although many RNA-binding proteins are uncharacterized structurally and do not fall into known categories, **primary sequence** analysis has led to the identification of a number of recurring RNA-binding **motifs** in functionally diverse RNA-binding proteins (1, 2). Interestingly in most cases, these are structurally and therefore probably evolutionarily distinct from **DNA-binding** motifs. RNA molecules are classified into various types based on their function or localization, such as [transfer RNA](#) (tRNA), [messenger RNA](#) (mRNA), **ribosomal RNA** (rRNA), viral RNA (vRNA), **small nuclear RNA** (snRNA), and small cytoplasmic RNA (scRNA). Table 1 gives some examples of RNA-binding proteins associated with each form, and Table 2 summarizes with examples the known RNA-binding motifs. Understanding how RNA-binding proteins specifically interact with their target RNA to form functional complexes is a key problem in structural biology, that has ramifications throughout molecular biology from [transcription](#) to **protein biosynthesis** and from replication of **viruses** to [development](#) of **embryos** and [genetic disease](#) .

**Table 1. Selected RNA-Binding Proteins, their RNA Targets, and their Functions**

| RNA                           | RNA binding proteins                        | Function   |
|-------------------------------|---|--|
| tRNA                          | tRNA (m <sup>5</sup> U54)methyl-transferase | General tRNA modification                          |
|                               | tRNA-guanine transglycosylase               | Specific tRNA modification                         |
|                               | aminoacyl-tRNA synthetase                   | tRNA aminoacylation                                |
|                               | met-tRNA <sup>fmet</sup> formyltransferase  | Prokaryote translation initiation                  |
|                               | IF3   | Prokaryote translation initiation                  |
| M1 RNA                        | EF-Tu                                       | tRNA transport to ribosome                         |
|                               | C5 protein                                  | Subunit of <i>E. coli</i> RNase-P endoribonuclease |
| mRNA/hnRNA                    | RNA polymerase II                           | Transcription                                      |
|                               | GreA  | Transcription termination                          |
|                               | hnRNP proteins                              | Pre-mRNA binding                                   |
|                               | Sex-lethal                                  | Pre-mRNA splicing                                  |
|                               | mRNA capping enzyme                         | 5' capping enzyme                                  |
|                               | <i>E. coli</i> RNase III                    | RNA processing                                     |
|                               | ds RNA adenosine deaminase                  | mRNA editing                                       |
|                               | CBP20, CPB80                                | Nuclear cap-binding protein subunits               |
|                               | eIF-4F                                      | Cytoplasmic cap binding protein                    |
|                               | Poly-A-polymerase, CStF                     | 3' poly-adenylation                                |
| Poly-A-binding protein (PABP) | mRNA stability                              |  |
| aCP-1, aCP-2                  | a-Globin mRNA stability                     |  |
| Staufen                       | mRNA localization                           |  |

|            |                                      |  |
|------------|--------------------------------------|--|
|            | Endo-, exoribonucleases              | mRNA degradation                           |
|            | IRE-binding protein, <i>E. coli</i>  | Translational regulation                   |
|            | Threonyl-tRNA synthetase, T4 regA    | Translational regulation                   |
| rRNA       | Thymidylate synthase                 | Translational regulation                   |
|            | RNA polymerase III                   | rRNA transcription                         |
|            | Nucleolin                            | rRNA processing                            |
|            | Nucleolar 2'-O-methyltransferase     | rRNA modification                          |
|            | TFIIIA                               | 5s rRNA storage                            |
| snRNA      | Sx, Lx                               | Small and large subunit ribosomal proteins |
|            | B, B', D1, etc.                      | Core snRNP proteins (Sm-proteins)s         |
| vRNA       | U1A, U2B', etc.                      | Specific SnRNP proteins                    |
|            | Human RNA helicase A                 | RNA helicases                              |
|            | RNA-dependent RNA polymerase         | Replications                               |
|            | Viral nucleocapsid proteins          | Assembly, protections                      |
|            | MS2 coat protein                     | Viral capsid, translational repressors     |
| RNA ligase | dsRNA-dependent protein kinase       | Translation regulations                    |
|            | Reverse transcriptase                | Copying RNA to DNAs                        |
|            | DNA ligase                           | Ligating RNAs                              |
|            | HIV-Tat                              | Transcriptional activators                 |
| scRNA      |                                      |  |
| SRP RNA    | Signal recognition particle proteins | Targeting of secretory proteins            |
| gRNA       | TUTase                               | mRNA editings                              |

**Table 2. Various RNA-Binding Domains and Sequence Motifs<sup>a, b</sup>**

| Sequence Motif          | Representative Proteins and Reference Describing Structure (in parentheses) | Target RNA                   |
|-------------------------|---|------------------------------|
| RNP domain              | U1A spliceosomal protein (7) (9,10)   | U1 snRNA<br>U1A protein mRNA |
| (RNA recognition motif) | U1 70k spliceosomal protein (1,6)   | U1 snRNA                     |
|                         | hnRNP protein C (1)   | mRNA precursor mRNA          |

|                           |   |                              |
|---------------------------|---|------------------------------|
|                           | hnRNP protein A1 <sup>a</sup>               | precursor                    |
|                           | Poly(A)-binding protein (1)                 | mRNA poly(A) tail            |
| dsRNA-binding motif       | <i>E. coli</i> RNase III (12)               | RNA transcript               |
|                           | <i>Drosophila</i> Staufen (11)              | Maternal bicoid mRNA         |
| KH domain                 | hnRNP protein K (1)                         | mRNA precursor               |
|                           | Fragile X protein (13)                      | Unknown                      |
|                           | Vigilin (14)                                | tRNA?                        |
| Zn-finger                 | TFIIIA (18)                                 | 5 S rRNA                     |
| S1 domain                 | Polynucleotide phosphorylase (15)           |                              |
|                           | Ribosomal protein S1                        | mRNA                         |
|                           | Initiation factor 1                         |                              |
| Sm domain                 | Spliceosomal core proteins (2)              | snRNAs                       |
| OB domain                 | Class IIb aminoacyl tRNA synthetase (19–21) | tRNA                         |
| RGG box                   | hnRNP proteins (1)                          | mRNA precursor               |
| <i>Alu</i> binding module | SRP 14 / 9 heterodimer                      | <i>Alu</i> domain of SRP RNA |

<sup>a</sup> Y. Shamoo, U. Krueger, L. M. Rice, K. R. Williams, and T. A. Steitz (1997) *Nature Struct. Biol.* **4**, 215–222; R.-M. Xu, L. Jokhan, X. Cheng, A. Mayeda, and A. R. Krainer (1997) *Structure* **5**, 559–570.

<sup>b</sup> M. Görlach, M. Friedrichs, G. Dreyfuss, and L. Mueller (1992) *Biochemistry* **31**, 9254.

Some of the important questions concerning the nature of protein-RNA recognition can be illustrated by considering tRNAs (see [Transfer RNA](#)). They function as adaptor molecules that specify [amino acid](#) residues corresponding to particular triplet **codons**. There are 20 [aminoacyl-tRNA synthetases](#), one for each [amino acid](#), and each of these enzymes must specifically recognize and aminoacylate *only* its cognate tRNA with its cognate amino acid ([3](#), [4](#)). In contrast, the prokaryotic [elongation factor](#) Tu (EF-Tu), which introduces aminoacylated tRNA (aa-tRNA) into the A-site of the [ribosome](#), binds to all aa-tRNA (except initiator-Met tRNA and **selenocysteinyl**-tRNA) and hence must recognize features common to all elongator aa-tRNAs ([5](#)). These are two extreme examples of tRNA recognition mechanisms that are essential in maintaining the fidelity and efficiency of protein biosynthesis. The various tRNA are transcribed as precursor RNA molecules and must be processed correctly by specific nucleases and a surprisingly large number of different base-modifying enzymes to form mature molecules. These modifying enzymes also have diverse specificities. Some recognize almost all tRNA [eg, tRNA (m<sup>5</sup>U54)methyltransferase]. Others recognize a subset of tRNA (eg, tRNA-guanine transglycosylase), even a unique tRNA species.

A wide variety of RNA-binding proteins are involved in producing, processing, transporting, translating, and degrading mRNA, the most abundant single-stranded RNA in cells. mRNA in eukaryotes is **capped** at the 5' end and **polyadenylated** at the 3' end, and **introns** must be correctly excised by the **splicing** machinery. Mature mRNAs are transported from the [nucleus](#) to the

**cytoplasm** for [translation](#). **Gene expression** is regulated at the transcriptional level and also at the translational level by a variety of mechanisms mediated by mRNA-binding proteins. The controlled **degradation** of mRNA is also an important process, involving a complex set of endo- and exonucleases. All of these processes must interact precisely with a large number of specific or nonspecific RNA-binding proteins and enzymes.

RNA is an essential structural and functional component of many ribonucleoproteins (RNP). Assembly of these RNP requires specific recognition of RNA structure and/or base identity by RNP protein subunits. The ribosome itself contains rRNA molecules, which are matured by splicing and modifying enzymes in the [nucleolus](#) in association with small nucleolar RNP (snoRNP). The rRNA may actually play the key catalytic role in peptidyl transfer. Pre-mRNA splicing is a complicated regulated process, involving a number of small nuclear RNPs containing snRNA that form a dynamic complex known as the [spliceosome](#). The substrate of the spliceosome is not naked RNA, but pre-mRNA (also known as heterogeneous nuclear RNA, **hnRNA**) complexed with various proteins that form particles known as hnRNP. The mammalian [signal recognition particle](#) (SRP) is an scRNP consisting of a 300-nucleotide 7 S RNA associated with six protein subunits. It binds to [signal peptides](#) at the N-terminal end of **secreted proteins** as they emerge from the ribosome and directs the nascent-chain/ribosome complex to receptors on the [endoplasmic reticulum](#) membrane in eukaryotic cells. Systems homologous to SRPs are found in all living cells. [Ribonuclease P](#), a tRNA-processing endonuclease, is an RNP in which the RNA component plays the catalytic role. [Telomerase](#) is a complex enzyme consisting of template RNA and protein subunits that adds multiple repeats of a particular sequence onto the ends of chromosomal DNA. Many viruses contain RNA as their genetic element, and virally-coded RNA-binding proteins play essential roles in the very diverse modes of replication and assembly of viruses.

How do RNA-binding proteins recognize their RNA-binding sites? RNA is distinguished from DNA most readily by the presence of the ribose 2'-OH and because the extended RNA double helix is predominantly in the **A-form** and has helical parameters significantly different from DNA, which is normally in the **B-form**. These general features permit nonspecific double-stranded RNA-binding proteins to recognize their correct nucleic acid substrate (eg, dsRNA-dependent protein [kinase](#)). On the other hand, specific RNA-binding proteins have more difficulty recognizing undistorted A-form RNA, which is characterized by a deep, narrow **major groove** and a shallower **minor groove**. The bases in the major groove are not readily accessible to protein side chains for specific recognition by **hydrogen bonding** except near the beginning or ends of helices. By contrast, the readily accessible minor groove contains less information for base discrimination. This problem is in reality bypassed by the fact that most cellular RNA is single-stranded, which forms irregular structures comprising short helices resulting from **Watson-Crick pairing** of short complementary stretches interspersed with hairpins, internal loops, bulges, or **pseudoknots**. Many RNA-binding proteins specifically recognize irregular RNA structures whose bases are more exposed for hydrogen bonding or stacking interactions with protein side chains. Indeed the irregularities may be specifically induced or enhanced by the protein/RNA interaction. Some RNA molecules fold further into complex tertiary structures (eg, tRNA) whose unique three-dimensional shape may be sufficient to define specific backbone interactions with an RNA-binding protein (eg, seryl-tRNA synthetase recognizing tRNA<sup>Ser</sup>).

## 1. Domains and Motifs Found in RNA-Binding Proteins

The majority of the known RNA-binding proteins have modular structures that contain an RNA-binding **domain** combined with other auxiliary domains (1, 2). Four RNA-binding sequence motifs have been found in RNA-binding proteins from diverse species: therefore, it is considered that they arose early in [evolution](#). These correspond to (1) the RNP domain, (2) the KH domain, (3) the dsRNA-binding domain, and (4) the S1 domain. [X-ray crystallography](#) has shown or [NMR](#) that they contain an a/b fold similar to those found in some ribosomal protein subunits, and it has been suggested that these RNA-binding motifs may have evolved from ribosomal proteins. This might



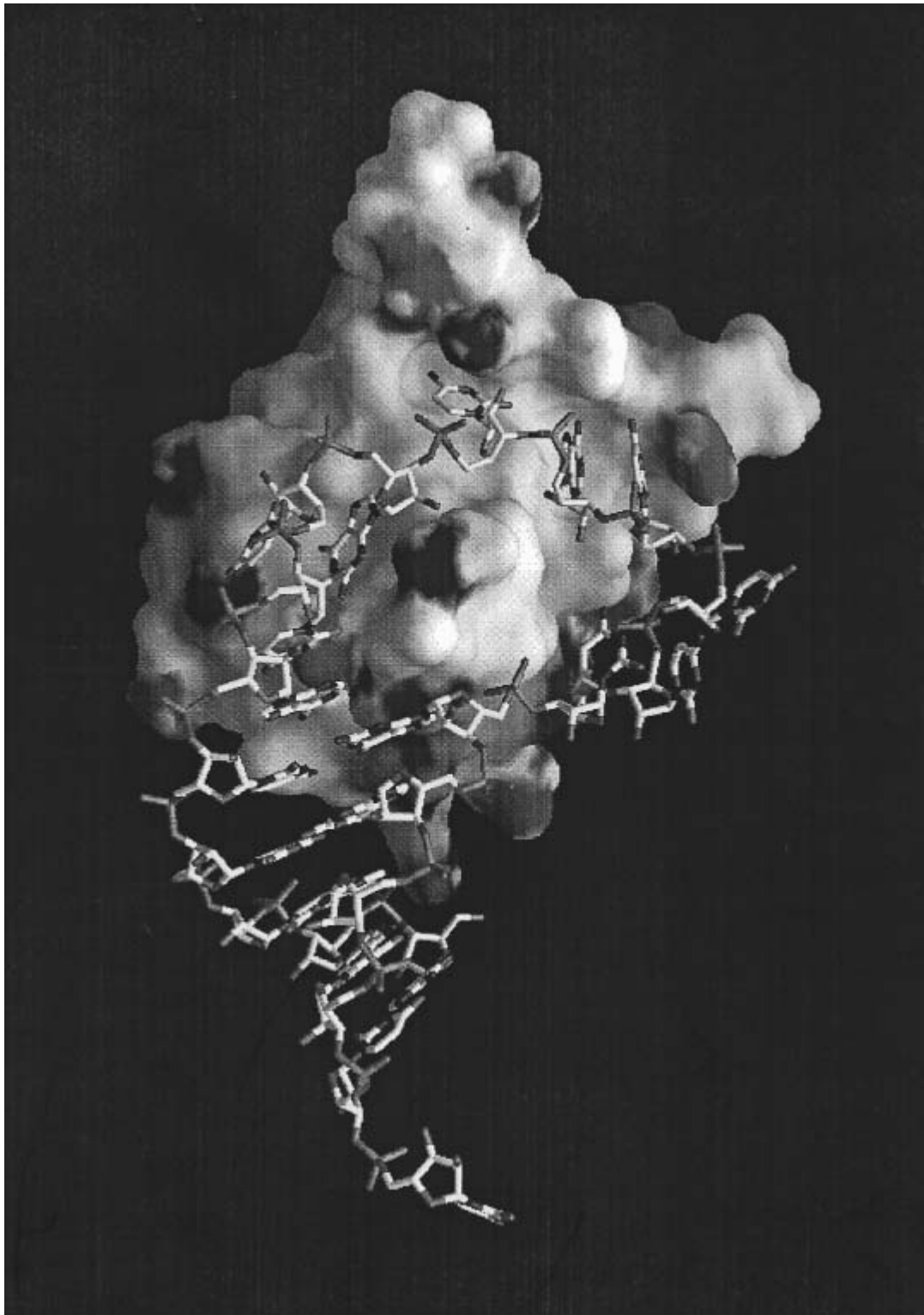
help to explain why RNA-binding domains are generally distinct from DNA-binding domains. Notable exceptions are the **zinc finger**, OB fold, and **homeodomain**. It is an interesting observation that many known RNA-binding proteins are all  $\beta$  or  $\alpha/\beta$  proteins that contain an exposed  **$\beta$ -sheet**. The crystalline structure of the complex between U1A protein and an RNA hairpin discussed later suggests why a  $\beta$ -sheet is a good RNA-binding surface.

### 1.1. RNP Domain

The RNP domain, also known as RNA recognition motif (RRM) or RNP-consensus sequence (RNP-CS) type RNA-binding domain, is found in more than 200 distinct RNA-binding proteins from diverse species (1, 2, 6). It is also the second most common protein sequence motif in the entire genome of the nematode *Caenorhabditis elegans*. These facts suggest that this module appeared early in evolution. The RNP domain is about 80 amino acid residues long and contains two highly conserved short sequence motifs called RNP octamer (RNP1) and RNP hexamer (RNP2). Some proteins contain only a single RNP domain, but others contain multiple copies. For example, the protein that binds to the polyadenylate (**poly A**) tail of mRNA in eukaryotes contains four tandem copies of the RNP domain.

The three-dimensional structure of the RNP domain first determined for U1A spliceosomal protein, indicated that it consists of a four-stranded antiparallel  $\beta$ -sheet flanked on one side by two  $\alpha$ -**helices**. The RNP1 and RNP2 motifs are located in the two middle  $\beta$ -strands of the sheet, and the side chains of the three highly conserved aromatic residues within RNP1 and RNP2 are projected onto the surface of the  $\beta$ -sheet (in U1A one of these residues is replaced by a glutamine) (7). U1A protein is a protein component of U1 snRNP, a large RNA/protein complex involved in pre-mRNA splicing, and binds to an RNA hairpin that contains 10 nucleotides in the loop. The crystal structure of a complex between U1A protein and an RNA hairpin representing its binding site has been determined at 1.9 Å resolution (8). The RNA hairpin loop binds to the surface of the  $\beta$ -sheet as an open structure, and the first seven nucleotides of the loop are fitted into a groove formed on the surface of the  $\beta$ -sheet (Fig. 1). The polypeptide loop between the  $\beta_2$  and  $\beta_3$  strands protrudes through the RNA loop. The bases of the seven nucleotides are splayed, and the protein-RNA contacts are made almost exclusively by RNA bases. These seven bases stack onto an adjacent base, a protein side chain, or both and also form an intricate hydrogen-bond network with protein side chains and the amide and carbonyl groups of the protein main chain. U1A protein also binds to the 3'-untranslated region of its own pre-mRNA and prevents its polyadenylation by directly interacting with poly-A polymerase. The binding site contains two internal loops, each with the AUUG(C/U)AC heptamer found in U1 snRNA hairpin II (6, 9, 10).

**Figure 1.** Crystal structure of a complex between U1A spliceosomal protein and its hairpin RNA binding site in U1 snRNA (8). The contact surface of the protein molecule is shown, whereas the RNA is a skeletal model. The AUUGCAC sequence in the ten-nucleotide loop fits into the groove on the surface of the  $\beta$ -sheet and binds tightly through stacking and hydrogen bond interactions with the protein.



## 1.2. Double-Stranded RNA-Binding Domain (dsRBD)

The dsRBD is a short sequence motif found in multiple copies in RNA-binding proteins from diverse origins including *Escherichia coli*, *Drosophila*, *Xenopus*, and mammals. In the *Drosophila* embryo, Staufen protein binds to maternal bicoid mRNA and plays an important role in establishing the anterior-posterior polarity through mRNA localization ([11](#)). *E. coli* ribonuclease III is an important enzyme involved in processing ribosomal and transfer RNA. The dsRBD is also found in **adenosine deaminase**, which is a key enzyme in [RNA Editing](#), and in the double-stranded RNA-dependent protein kinase, which plays an important role in viral gene expression. The NMR structures of

dsRBD from RNase III and Staufen protein show that the module contains a three-stranded antiparallel  $\beta$ -sheet that has two  $\alpha$ -helices packed on one side (11, 12). The dsRBD shows strong similarities to the N-terminal domain of *Bacillus stearothermophilus* ribosomal protein S5 at the levels of both amino acid sequence and three-dimensional structure (11). This strongly suggests that dsRBD may have evolved from a ribosomal protein. It is not yet known how this domain binds RNA, but [site-directed mutagenesis](#) experiments suggest that loops between  $\beta$ 1 and  $\beta$ 2 and between  $\beta$ 3 and  $\alpha$ 2 may be involved in RNA-binding. The major groove of dsRNA is narrow and deep, and the bases are not accessible by protein side chains for sequence-specific interactions. It is likely that each dsRBD binds to the minor groove or to the phosphate backbone in the major groove. Biochemical experiments suggest that a set of dsRBD, together with some auxiliary domains, may recognize a structural feature of folded RNA.

### 1.3. KH Domain

Heterogeneous nuclear ribonucleoprotein (hnRNP) K is one of the nuclear proteins that bind to precursors of mRNA (1). There are three copies of a sequence motif within hnRNP K, and homologous sequences have been found in many proteins from diverse species, from bacteria to man, that interact with RNA. Human Fragile X (FMR1) gene contains two copies of the KH domain. A mutation in one of these is responsible for severe hereditary mental retardation (1, 13). The natural RNA target of the FMR1 gene product has not been identified, but the mutation that causes severe hereditary mental retardation abolishes the binding of poly(U) to the FMR1 protein. This suggests that the ability of FMR1 to bind to RNA is essential for its *in vivo* function. Two proteins that contain three copies of the KH domain are associated with the 3'-untranslated region of  **$\alpha$ -globin** mRNA. These proteins affect the *in vivo* half-life of the  $\alpha$ -globin mRNA. Ribosomal protein S3 from *E. coli* and other bacteria contains the KH domain, suggesting that the KH domains found in higher organisms have evolved from a ribosomal protein (1). The NMR structure of the KH domain from human vigilin has an  $\alpha/\beta$  structure that contains a three-stranded  $\beta$ -sheet with three  $\alpha$ -helices (14).

### 1.4. S1 Domain

Ribosomal protein S1 from *Escherichia coli* contains six copies of a short sequence motif of approximately 80 amino acid residues. S1 protein is directly involved in RNA binding because it can be **cross-linked** to mRNA in the translational initiation complex. The S1 domain has been found in proteins from bacteria and eukaryotes, including [polynucleotide phosphorylase](#), [initiation factor 1](#) (IF1), NusA, **ribonucleases** II and E from *E. coli*, yeast [RNA helicase](#) PRP22, eukaryotic initiation factor eIF2a, and eIF2a kinase inhibitor. An NMR study of the S1 domain of polynucleotide phosphorylase shows that this domain consists of a five-stranded antiparallel  $\beta$ -barrel (15).

### 1.5. MS2 Bacteriophage Coat Protein

One of the best studied examples of an RNA-binding protein is the MS2 **bacteriophage** coat protein. The MS2 phage contains a genomic RNA 3569 nucleotides long, packaged into an icosahedral protein shell that consists of 180 copies of the coat protein. A capsid protein dimer binds to an RNA hairpin formed near the ribosomal-binding site of the virally encoded **replicase** gene, thereby inhibiting translation of the replicase mRNA. The interaction also triggers assembly of the coat protein and packaging of the genomic RNA. A synthetic RNA hairpin that represents the binding site was soaked into a crystal of the empty phage particle, and its structure was determined to 2.8 Å resolution (16). The RNA, which contains a **tetraloop** and a bulged adenosine, binds to the surface of the continuous  $\beta$ -sheet across the dyad axis. The side chain of a [tyrosine](#) residue stacks onto a cytidine in the tetraloop, and the bulged adenosine and an adenosine in the tetraloop bind to equivalent sets of residues from each subunit.

### 1.6. Zinc-Containing RNA-Binding Proteins

[Transcription factor](#) IIIA is a transcriptional activator of the 5 S rRNA gene, but it also binds to its gene product 5 S rRNA, and functions as a storage or transport protein (17). It contains nine **zinc fingers**, each of which folds into a domain that contains two  $\beta$ -strands and one  $\alpha$ -helix stabilized by the coordination of two [histidine](#) and two [cysteine](#) residues to a zinc ion. Crystallographic analyses of zinc finger proteins in complex with dsDNA show that the  $\alpha$ -helix fits into the major groove of the

DNA and forms many sequence-specific contacts with the DNA bases (18). It is believed that 5 S rRNA folds into a branched double helix that contains many bulges which may widen the major groove of RNA to permit entry of the recognition helices. The HIV nucleocapsid protein is another example of a zinc-containing, RNA-binding protein. Others are *E. coli* alanyl-tRNA synthetase and tRNA guanine transglycosylase (28).

### 1.7. Other RNA-Binding Modules

The bases of the anticodon loops of tRNA<sup>Asp</sup> and tRNA<sup>Lys</sup> are splayed over the surface of a b-barrel domain of their cognate tRNA synthetases and are recognized in a sequence specific manner (see later) (19, 20). Many structural homologues of this b-barrel domain, which is known as the *OB-fold*, have been found in proteins that bind oligonucleotides (both ssRNA and ssDNA) or oligosaccharides (21).

The Sm proteins that form the core of spliceosomal small nuclear ribonucleoprotein particles (snRNP) contain a conserved sequence motif. It is predicted that this domain has an a/b fold, but its structure is yet to be determined (2). Some hnRNP proteins contain multiple repeats of an Arg-Gly-Gly (RGG) sequence that are believed involved in RNA binding (1).

The crystal structure of the SRP 14 / 9 heterodimer, which binds to the *Alu* domain of the mammalian signal recognition particle RNA, has been determined. SRP9 and SRP14 are structurally homologous and contain the same a-b-b-b-a fold, related to but distinct from the dsRNA-binding module. The heterodimer has pseudo two-fold symmetry and is saddle-like, comprising a strongly curved six-stranded amphipathic b-sheet. The four helices are packed on the convex side, and the exposed concave surface is lined with positively charged residues.

Both HIV tat and rev proteins bind the TAR and RRE elements, respectively, to their target RNA sequences, by arginine-rich peptides. The solution structure of a 14-residue arginine-rich peptide from HIV tat complexed with HIV TAR has been determined by NMR (26). The peptide forms a b-hairpin that interacts in the RNA major groove.

## 2. tRNA-Binding Proteins

### 2.1. Aminoacyl-tRNA Synthetases

The fidelity of protein synthesis depends to a large extent on the extreme specificity with which aminoacyl-tRNA synthetases charge their cognate tRNA with their cognate amino acid. In *E. coli*, there are at least 46 different tRNA molecules that have anticodons which correspond to the various amino acids. The seryl-tRNA synthetase, for example, has to charge the six serine isoacceptors (including one for selenocysteine) selectively and ignore the others. Because tRNA superficially have similar secondary and tertiary structures, what is the molecular basis for the specific recognition between aminoacyl-tRNA synthetases and tRNA? In many cases, extensive biochemical studies have revealed the so-called tRNA identity elements (see [Aminoacyl tRNA Synthetases](#)).

A more detailed picture of specific tRNA recognition and catalysis by aminoacyl-tRNA synthetases is emerging from crystallographic studies of aminoacyl tRNA synthetases complexed with various combinations of their three substrates: ATP, cognate amino acid, and cognate tRNA, plus the activated amino acid intermediate, the aminoacyl-adenylate. Now, crystal structures are known of 14 of the 20 aminoacyl tRNA synthetases, five of which are in complexes with cognate tRNA (3, 4). The three systems, for which the most extensive structural data exists on protein-tRNA recognition, are the class I glutamyl system and the class II aspartyl and seryl systems. These show strikingly different modes of specific synthetase-tRNA interaction but share certain general features, including a fairly large synthetase-tRNA interactive interface characterized by (1) nonspecific backbone contacts, often involving basic residues, both of which increase binding affinity and aid correct positioning and orienting the tRNA; (2) discriminatory base-specific interactions restricted to a few regions, principally the anticodon and the acceptor stem. The second general feature is mutual induced fit by which protein-RNA contacts are made as a result of conformational changes in either

or both macromolecules. This includes ordering of protein loops and reorienting and stabilizing domains (eg, SerRS), base-pair breaking in the acceptor stem, and 3'-end distortion (eg, tRNA<sup>Gln</sup>) and destacking of bases in the anticodon loop (eg, tRNA<sup>Gln</sup>, tRNA<sup>Asp</sup>).

## 2.2. Glutamyl-tRNA Synthetase (GlnRS)

GlnRS is a monomeric class I synthetase whose specificity for tRNA<sup>Gln</sup> is largely determined by interactions with identity elements in the tRNA acceptor stem and anticodon stem-loop, both of which have severe distortions from the structure found in uncomplexed tRNA. In the *E. coli* complex, the tRNA anticodon stem is extended from five to seven base pairs by two extra non-Watson–Crick base pairs (23). The three anticodon bases (CUG) are splayed to fit into three separate recognition pockets formed at the interface between the distal two b-barrel domains of the protein. In the active site of the synthetase, the tRNA is oriented so that specific interactions can be made within the acceptor stem's minor groove to identity determinants in base pairs 2 and 3. On the other hand, the tRNA 3'-end reaches the catalytic center only by forming an unusual hairpin turn. This conformation requires breaking the first U1-A72 base pair and is stabilized, in part, by a hydrogen bond between the discriminator base G73 and the phosphate of A72.

## 2.3. Aspartyl-tRNA Synthetase (AspRS)

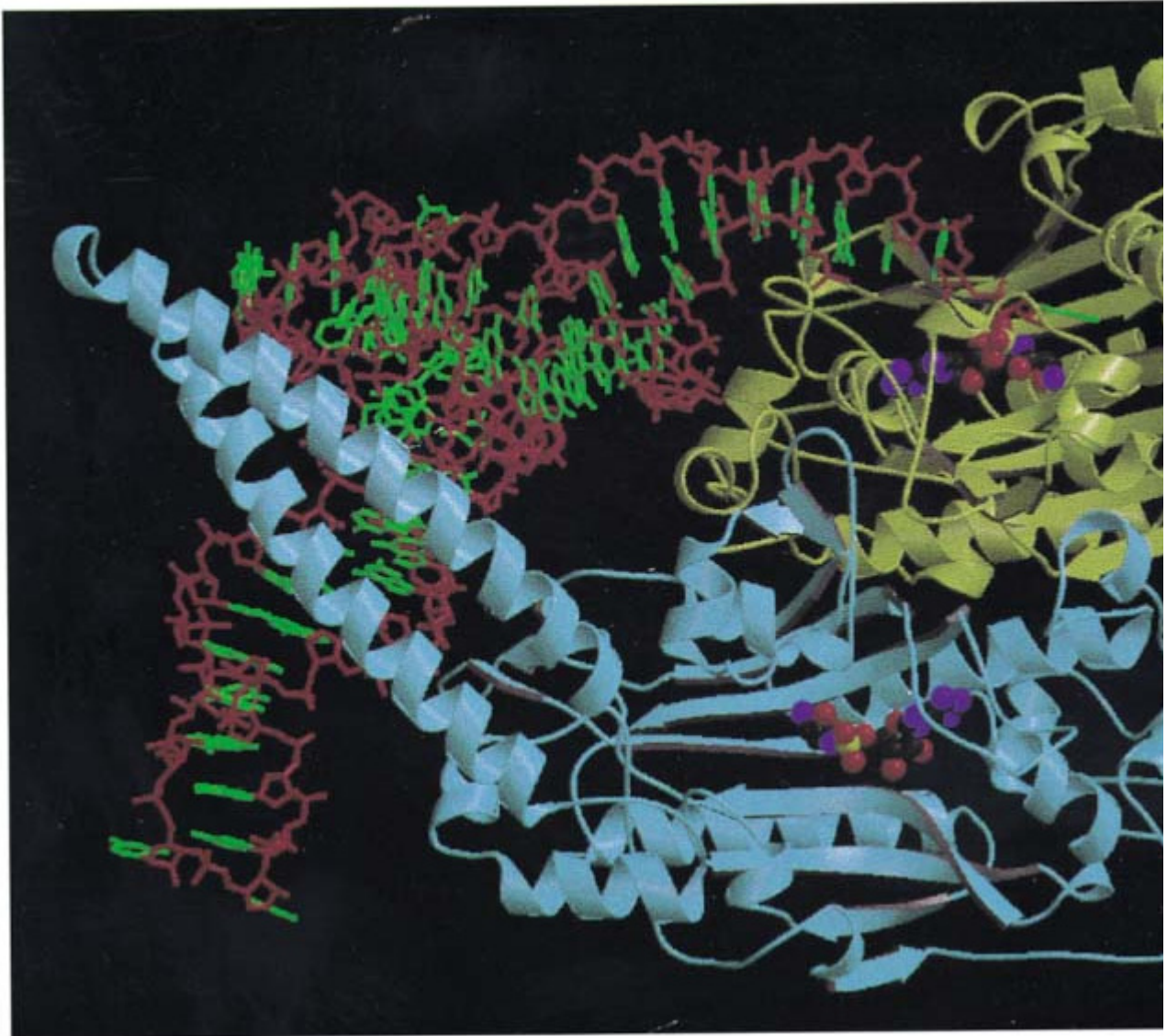
AspRS is a dimeric class IIb synthetase that binds two tRNAs symmetrically, although each tRNA interacts predominantly with only one subunit. In contrast to class I, class II synthetases interact with the acceptor stem of their cognate tRNA from the major groove side (27) by the so-called motif 2 loop, which in yeast AspRS makes base-specific interactions with the discriminator base G73 and base pair U1-A72 (19). The major groove side recognition also means that the single-stranded 3'-end of the tRNA enters the synthetase active site without significant distortion from its normal helical path, again in strong contrast to class I. Anticodon recognition by AspRS is performed by an N-terminal, five-stranded, b-barrel domain (OB fold) (21). The normal compact structure of the free tRNA anticodon loop undergoes a large conformational change, and the five anticodon loop bases are exposed to the exterior. The three anticodon bases (GUC) lie across the b-sheet surface and are recognized by specific hydrogen bonding interactions (19). Recent crystallographic results on another closely related class IIb synthetase complexed with its cognate tRNA, *T. thermophilus* lysyl-tRNA synthetase, show a very similar interaction between the anticodon of tRNA<sup>Lys</sup> (CUU) and the N-terminal, b-barrel domain of LysRS (20). Whereas the central U35 interacts identically in the two systems (by stacking with a conserved phenylalanine residue and hydrogen bonding with conserved glutamine and arginine residues), the specificity for base 36 and to a lesser extent for base 34 is idiosyncratic for each synthetase.

## 2.4. Seryl-tRNA Synthetase (SerRS)

SerRS is also a dimeric class II synthetase, but its mode of tRNA recognition differs significantly from that of AspRS (Fig. 2). SerRS is characterized by a unique 100-residue, N-terminal domain that is folded into a 60 Å long, solvent exposed and flexible, antiparallel coiled coil, known as the helical arm (24, 25). The variety of serine anticodons from two distinct groups means that the serine isoacceptors share no common anticodon base. Thus the anticodon is not an identity element of tRNA<sup>Ser</sup>, whereas the long variable arm is, a feature shared only by tRNA<sup>Leu</sup> and tRNA<sup>Tyr</sup>. The main features of the recognition of tRNA<sup>Ser</sup> by *T. thermophilus* SerRS are (1) the tRNA binds across the two subunits of the dimer; (2) upon tRNA binding the helical arm of the synthetase is stabilized in a new orientation and binds between the TYC loop and the long, variable-arm of the tRNA; (3) contacts with the tRNA long, variable-arm backbone extend until the sixth base pair, explaining the need for a minimum length of the tRNA variable arm; (4) the synthetase makes many backbone contacts, but few base-specific interactions, and is principally recognizing the unique shape of tRNA<sup>Ser</sup>. The unique shape of tRNA<sup>Ser</sup> is largely determined by bases 20A and 20B inserted into the D-loop, both of which play novel roles in tertiary interactions in the tRNA core. In particular, the base of G20B is stacked against the first base pair of the long variable arm and thus defines the spatial orientation of the latter; (5) the anticodon stem loop is not in contact with the synthetase; (6)

the motif 2 loop of SerRS (longer than that of AspRS) makes base contacts down to the fourth base pair within the acceptor stem's major groove, but they are, however, only weakly discriminatory (25).

**Figure 2.** Overall view of the ternary complex of seryl-tRNA synthetase, tRNA<sup>Ser</sup>(GGA), and a nonhydrolyzable analog monomer 1 is in yellow, and monomer 2 is in blue. Only the backbone and its secondary structure are shown (arrows are tRNA backbone is in red, and the bases are in green). The tRNA is viewed looking down the anticodon stem, which is not in the figure. The long variable arm of the tRNA crosses the helical arm of monomer 2 perpendicularly and emerges at the Ser-AMS molecule is represented by spheres. See color insert.



## 2.5. Elongation Factor EF-Tu

EF-Tu is a **G-protein** that, in the activated GTP-bound state, binds all aminoacylated elongator tRNA with higher affinity than uncharged tRNA and delivers them to the ribosome. The crystalline structure of the ternary complex of yeast tRNA<sup>Phe</sup>, *T. aquaticus* EF-Tu and GDPNP has recently been determined and reveals how activated EF-Tu contacts relatively limited and conserved regions of aa-tRNA (5). The only parts of the aa-tRNA in contact with EF-Tu are the aminoacylated CCA end, the 5'-phosphate, and the T-stem helix. The long stem of the L-shaped tRNA projects away

from the protein, so that the entire complex is extremely elongated. The CCA-3' end binds in a cleft between domains 1 and 2 of EF-Tu, and conserved residues from domain 2 interact with the base and phosphate of A76. Main-chain hydrogen bonds from EF-Tu are made to the ester group linking the carboxyl group of the amino acid to the 3'-OH of the ribose and to the free amino group of the amino acid. The amino acid itself projects into a pocket formed between domains 1 and 2. The tRNA 5'-phosphate is bound by conserved basic residues from helix B of domain 1 and from domain 2. One side of the T-stem helix is packed against a depression between domains 1 and 3. The contacts are made by nonconserved residues from domain 3.

It has been pointed out that the EF-Tu.GTP.tRNA<sup>Phe</sup> ternary complex extraordinarily similar in shape to that of the elongation factor EF-G.GDP binary complex. This mimicry suggests that the two complexes might bind to the same ribosomal state.

### Bibliography

1. C. G. Burd and G. Dreyfuss (1996) *Science* **265**, 615–621.
2. K. Nagai (1994) *Curr. Opin. Struct. Biol.* **6**, 53–61.
3. J. G. Arnez and D. Moras (1994) *RNA-protein interactions* (ed, K. Nagai and I. W. Mattaj, eds.), IRL Press, Oxford, pp. 52–81.
4. S. Cusack (1995) *Nature Struct. Biol.* **2**, 824–831.
5. P. Nissen, M. Kjeldgaard, S. Thirup, G. Polekhina, L. Reshetnikova, B. Clark, and J. Nyborg. (1995) *Science* **270**, 1464–1472.
6. K. Nagai, C. Oubridge, N. Ito, J. Avis, and P. Evans (1995) *Trends Biochem. Sci.* **20**, 235–240.
7. K. Nagai, C. Oubridge, T.-H. Jessen, J. Li, and P. R. Evans. (1990) *Nature* **348**, 515–520.
8. C. Oubridge, N. Ito, P. R. Evans, C. H. Teo, and K. Nagai (1994) *Nature* **372**, 432–438.
9. F.-H. T. Allain, C. C. Gubser, P. W. A. Howe, K. Nagai, D. Neuhaus, and G. Varani (1996) *Nature*, **380**, 646–650.
10. C. W. G. van Gelder, S. I. Gunderson, E. J. R. Jansen, W. C. Boelens, M. Polycarpou-Schwartz, I. W. Mattaj, and W. J. van Venrooij (1993) *EMBO J.* **12**, 5191–5200.
11. M. Bycroft, S. Grünert, A. G. Murzin, M. Proctor, and D. St. Johnston (1995) *EMBO J.* **14**, 3563–3571.
12. A. Kharrat, M. J. Macias, T. J. Gibson, M. Nilgens, and A. Pastore (1995) *EMBO J.* **14**, 3572–3584.
13. H. Siomi, M. Choi, M. C. Siomi, R. L. Nussbaum, and G. Dreyfuss. (1994) *Cell* **77**, 33–39.
14. G. Musco, G. Stier, C. Joseph, M. A. C. Morelli, M. Nilges, T. J. Gibson, and A. Pastore (1996) *Cell* **85**, 237–245.
15. M. Bycroft, T. J. P. Hubbard, M. Proctor, S. M. V. Freund, and A. G. Murzin (1997) *Cell* **88**, 1–20.
16. K. Valegård, J. B. Murray, P. G. Stockley, N. J. Stonehouse, and L. Liljas (1994) *Nature* **371**, 623–626.
17. J. Miller, A. D. McLachlan, and A. Klug (1985) *EMBO J.* **4**, 1609–1614.
18. N. P. Pavletich and C. O. Pabo (1991) *Science* **252**, 809–817.
19. J. Cavarelli, B. Rees, M. Ruff, J. C. Thierry, and D. Moras (1993) *Nature* **362**, 181–184.
20. S. Cusack, A. Yaremchuk, and M. Tukalo (1996) *EMBO J.* **15**, 6321–6334.
21. A. G. Murzin (1993) *EMBO J.* **12**, 861–867.
22. M. A. Rould, J. J. Perona and T. A. Steitz (1991) *Nature* **352**, 213–218.
23. V. Biou, A. Yaremchuck, M. Tukalo, and S. Cusack (1994) *Science* **263**, 1404–1410.
24. S. Cusack, A. Yaremchuk, and M. Tukalo (1996) *EMBO J.* **15**, 2834–2842.
25. J. D. Puglisi, L. Chen, S. Blanchard, and A. D. Frankel (1995). *Science* **270**, 1200–1203.

26. M. Ruff et al., (1991) *Science* **252**, 1682–1689.  
 27. C. Romier, K. Reuter, D. Suck, and R. Ficner (1996) *EMBO J.* **15**, 2850–2857.

### Suggestions for Further Reading

28. K. Nagai and I. W. Mattaj (1994) *RNA-Protein Interactions*, IRL Press, Oxford.  
 29. R. Gestland and J. Atkins (1993) *The RNA world*, Cold Spring Harbor Press, Cold Spring Harbor, New York.  
 30. D. E. Draper (1995). Protein-RNA recognition, *Annu. Rev. Biochem.* **64**, 593–620.  
 31. D. Söll and U. L. RajBhandary (1995). *tRNA, Structure, Biosynthesis and Function*, ASM Press, Washington, D.C.

## ROESY Spectrum

A conventional [nuclear Overhauser effect](#) (NOE) experiment in nuclear magnetic resonance ([NMR](#)) spectroscopy involves detection of changes in the intensities of NMR signals associated with a particular spin of the molecule being examined, when other spins are perturbed in some way. The changes in signal intensity are the result of altered populations of nuclear spin energy levels associated with both sets of spins and reflect their physical proximity. There may be magnetic interactions between spins while they are in coherent states and these may lead to changes of signal intensity. This second kind of Overhauser effect is detected in a “rotating frame” NOE (or ROE) experiment. The ROE may be determined in a one-dimensional format or be incorporated into multidimensional experiments. A two-dimensional (2D) ROE experiment is often referred to as ROESY (rotating frame Overhauser effect spectroscopy). In such a 2D experiment, the appearance of a cross peak at the chemical shift coordinates (A, B) requires that there be a ROE between the spins characterized by these chemical shifts ( $s_A, s_B$ ). The power of experiments that produce ROEs or NOEs is that changes in signal strength are strongly dependent on the distance between interacting spins. Thus, the observation of a ROE or NOE provides a constraint that can be used to define the [tertiary structure](#) of the macromolecule being studied.

For the standard NOE experiment wherein the detected signal intensity changes arise from alterations of spin energy level populations, the effect is given by

$$f_I\{S\} = I_p/I_0 - 1 \quad (1)$$

where  $f_I\{S\}$  indicates the NOE on the signal from spin I when there is a perturbation of level populations associated with spin S:  $I_0$  is the normal intensity of the signal for spin I observed from a sample at thermal equilibrium before the analyzing RF pulse, and  $I_p$  is the intensity of the same signal when there has been a perturbation of spins S before the analyzing pulse. The value of  $f_I\{S\}$  depends on the gyromagnetic ratios of spins I and S, how these spins move in the sample, the strength of the magnetic field used for the NMR experiment, and the details of how the energy level populations associated with spin S are perturbed during the course of an experiment. Values of  $f_I\{S\}$  ranging from 0.5 to -1.0 are possible when both I and S are protons. It should be noted that, given the definition of  $f_I\{S\}$ , some set of experimental conditions may exist for which  $f_I\{S\}$  becomes zero. Under these conditions, there is little or no change in signal intensity, even though spins I and S



might be very close to each other. These conditions are typically experienced when the mass of the molecule under study is in the range of 500–2000 D.

In contrast, the ROE on signal intensities is always positive for near-neighbor nucleus-nucleus interactions and thus remains detectable under all experimental conditions (1-3). This very significant advantage is countered, however, by the need to correct the experimental data for various artifacts and by the possibility that coherence transfers may produce effects that are unrelated to the ROE. Analysis of ROE data in conjunction with conventional NOE data may provide important insights into existing chemical exchange processes (4). Although most ROESY experiments with biological macromolecules involve interactions between hydrogen  $^1\text{H}$  atoms, useful heteronuclear ROESY experiments also are feasible (5). The elements of the ROESY experiment can be built into experiments that produce three-dimensional or higher NMR spectra. (See also [Nuclear Overhauser Effect \(NOE\)](#), [NOESY Spectrum](#).)

#### Bibliography

1. A. A. Bothner-By, R. L. Stephens, J. M. Lee, C. D. Warren, and R. W. Jeanloz (1984) *J. Am. Chem. Soc.* **106**, 811–813.
2. A. Bax and D. G. Davis (1985) *J. Magn. Reson.* **63**, 207–213.
3. D. Neuhaus and M. P. Williamson (1989) *The Nuclear Overhauser Effect in Structural and Conformational Analysis*, V.C.H., New York, Chapter "9", pp. 307–350.
4. J. Fejzo, W. M. Westler, S. Macura and J. L. Markley (1991) *J. Magn. Reson.* **92**, 20–29.
5. G. W. Kellogg and B. I. Schweitzer (1993) *J. Biomol. NMR* **3**, 577–595.

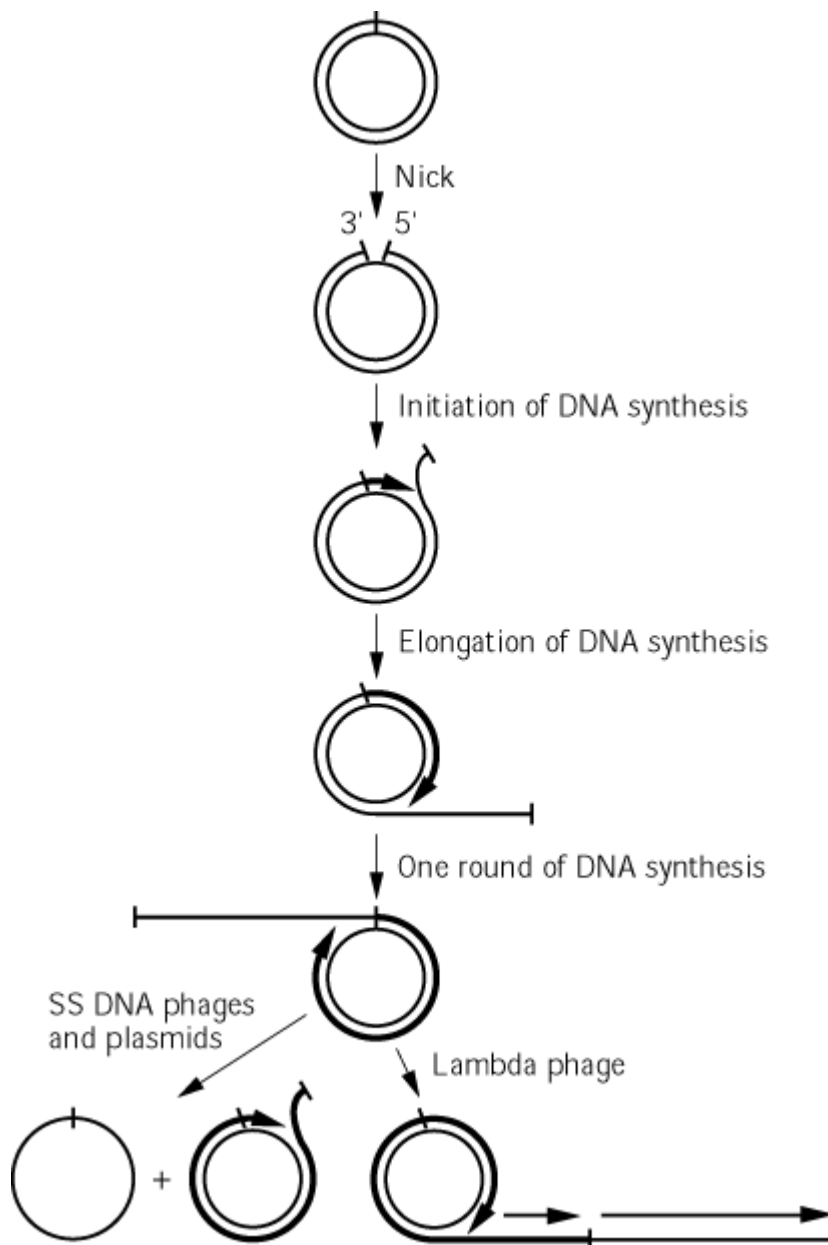
#### Suggestions for Further Reading

6. C. J. Bauer, T. A. Frenkiel, and A. N. Lane (1990) *J. Magn. Reson.* **87**, 144–152.
7. T. E. Bull (1992) *Prog. NMR Spectroscopy* **24**, 377–410.
8. D. Canet (1996) *NMR Concepts and Methods*, Wiley, New York, Chapter "4", pp. 139–196.
9. J. Cavanagh, W. J. Fairbrother, A. G. Palmer III, and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic Press, San Diego.

### Rolling Circle DNA Replication

In [DNA replication](#), the **DNA polymerase** cannot initiate the synthesis of a new DNA strand and must rely on a priming device. In general, an RNA [primer](#) is synthesized at or near a [replication origin](#) to start synthesis of the **leading strand**. However, a DNA primer terminus can be generated by a **nuclease**-generated nick at a specific place in some circular duplex DNA, and replication will then proceed unidirectionally, as shown in Figure 1. This mode of replication is called *rolling circle replication* and is found for replication of the [replicative form \(RF\)](#) form of **bacteriophage** single-stranded genomes of [Gram-negative bacteria](#) and of the multicopy **plasmids** of Gram-positive bacteria (see [Single-Stranded DNA Replication](#)). Rolling circle replication is also observed in the late stage of the replication of the [lambda phage](#) genome and in the process of the **conjugative** transfer of bacterial plasmids.

**Figure 1.** Scheme for rolling circle replication (see text for details).



DNA synthesis initiates using the free 3'-OH end at the nick as a primer, and a [replication fork](#) proceeds around the template. In the process, the newly synthesized strand displaces the old strand from the [template](#). In the case of replication of the RF form of single-stranded phage **genomes** and of plasmids of Gram-positive bacteria, the displaced old strand is cleaved off after one round of replication and is converted into the circular, double-stranded form. In contrast, in phage lambda replication, the replication fork proceeds a number of revolutions around the template without cleavage of the displaced strand, and the displaced strand becomes double-stranded as it is peeled off. The linear [concatemer](#) thus created is cleaved into one unit length and packaged into the phage particles. In the conjugation process of plasmids, the displaced strand is transferred into the new cell.

The initial nick of rolling circle replication is introduced by an **endonuclease** specific for each system. The similarity of the amino acid sequences of the initiator endonucleases and of the proteins (known as relaxases) involved in the initiation and termination of conjugative DNA strand transfer clearly indicate that they are evolved from a common ancestor (1, 2).

Bibliography

1. T. V. Ilyina and E. V. Koonin (1992) Nucl. Acids Res. **20**, 3279–3285.
2. D. R. Byrd and S. W. Matson (1997) Mol. Microbiol. **25**, 1011–1022.

## Rossmann Fold

In the 1970s, Rossmann and colleagues ([1](#), [2](#)) identified, analyzed, and characterized the similarity between the [protein structures](#) of [lactate dehydrogenase](#) and malate dehydrogenase, two [enzymes](#) that bind the dinucleotide NAD. Both protein structures incorporate two b–a–b–a–b [protein motifs](#), and each half of the bound dinucleotide binds to one of these units. The b–a–b–a–b structural motif, which forms a parallel **b-sheet** surrounded by **a-helices**, is referred to as either a Rossmann fold or a **mononucleotide-binding motif**.

[See also [Nucleotide-Binding Motif](#).]

### Bibliography

1. S. T. Rao and M. G. Rossmann (1973) J. Mol. Biol. **76**, 241–256.
2. M. G. Rossmann, D. Moras, and K. W. Olsen (1974) Nature **250**, 194–199.

## R<sub>0</sub>T Curve

### 1. DNA–RNA Hybridization

During purification of [messenger RNA](#), its quantity and the extent of purification can be determined by **hybridization** techniques. This also permits determining the size of the mRNA independently by [gel electrophoresis](#).

Hybridization involves **annealing** a large excess of RNA with double-stranded DNA previously **radiolabeled** by **nick translation** or during its synthesis. Such labeling techniques result in a high specific radioactivity, so the necessary excess of RNA can be obtained easily. Heating double-stranded DNA denatures it to its single-stranded components and slow cooling of these single strands permit [renaturation](#). This technique is used in **C<sub>0</sub>t curves** to characterize DNA samples. Hall and Spiegelman ([1](#)) demonstrated that the same experiment can be performed with heat-denatured DNA and the corresponding mRNA, and slow cooling produces DNA–mRNA hybrids that were observed in cesium chloride [density gradient centrifugation](#). When mRNA was heated with genetically unrelated DNA, no hybrid molecules were formed.

#### 1.1. Kinetics of Association

DNA–RNA hybridization is a second-order reaction (see [Kinetics](#)). If  $t$  is the time,  $D$  the concentration of single-stranded DNA in moles of nucleotides/liter,  $R$  the concentration of excess mRNA in the same units, and  $k$  the rate constant for the reaction, then the equation describing the reaction is the following:

$$dD/dt = -kDR$$

The value of  $k$  is expressed in units of  $M^{-1}s^{-1}$ .

DNA self-reassociation is avoided by using a low concentration with RNA in excess. Then the RNA hybridizes much faster than the DNA, and its concentration remains practically constant throughout the reaction. The reaction follows pseudo-first order kinetics, and the equation above simplifies to the following:

$$\frac{dD}{dt} = -kR_0$$

$R_0$  is the mRNA concentration at  $t = 0$ .

By substituting  $D = D_0$  at  $t = 0$ , the integrated equation becomes

$$\ln \frac{D_0}{D} = kR_0t$$

When 50% of the DNA (cDNA) strands have been converted into DNA-RNA hybrids, at  $t = t_{1/2D}$  becomes  $D_0/2$ , and  $R_0t$  becomes  $R_0t_{1/2}$ . This yields

$$\ln \frac{D_0}{D_0/2} = kR_0t_{1/2}$$

from which the value of  $k$  can be derived by

$$k = \frac{\ln 2}{R_0t_{1/2}}$$

Therefore the velocity of the association reaction is inversely proportional to the value of  $R_0t_{1/2}$ , which itself is directly related to the size of a pure RNA or DNA molecule.

## 1.2. Determining mRNA size and evaluating a given mRNA in an RNA mixture

Knowing the size  $a$  of a given mRNA, the unknown size  $x$  of another mRNA is determined by the following equation:

$$\frac{\text{Size of mRNA } a}{\text{Size of mRNA } x} = \frac{R_0t_{1/2}(\text{mRNA } a)}{R_0t_{1/2}(\text{mRNA } x)}$$

For ovalbumin mRNA, which is 2000 nucleotides long, the  $R_0t_{1/2}$  value is  $4 \times 10^{-3} M^{-1}s^{-1}$ . For prolactin mRNA (900 nucleotides long) the  $R_0t_{1/2}$  value is accordingly smaller:  $1.8 \times 10^{-3} m^{-1}s^{-1}$  (2).

The values of  $R_0t_{1/2}$  allow determining the fraction of a certain mRNA in a mixture by comparing hybridization of the corresponding cDNA with the corresponding "pure" mRNA and with the same mRNA in a mRNA mixture. Because the size of the desired mRNA is constant, its reassociation rate depends only on its concentration in the mixture, according to the equation

$$\frac{R_0t_{1/2}(\text{"pure" mRNA})}{R_0t_{1/2}(\text{unknown sample})} = \frac{\% \text{mRNA (unknown sample)}}{\% \text{mRNA (pure mRNA)}}$$

The course of hybridization is monitored by determining the number of DNA–RNA hybrids at various reaction times. This is done in two ways (2):

1. by applying the reaction mixture to a [chromatography](#) column of **hydroxyapatite** in 0.05 M **phosphate buffer** at pH 6.8 and at a temperature of 60°C. The single-stranded DNA elutes with 0.14 M phosphate buffer, whereas the double-stranded DNA–RNA molecules elutes with 0.45 M phosphate buffer.
2. by digestion with [S1 nuclease](#), which digests single-stranded DNA, whereas RNA–DNA hybrids remain unaffected and are precipitated with trichloroacetic acid.

## 2. Determining the Number of Different Sequences in mRNA Transcribed from a cDNA Library

To judge whether a **cDNA** library is representative, ie, whether the number of different clones carrying a particular cDNA insert reflects the relative fraction of the corresponding mRNA molecules in the cell from which the library has been prepared, the number of different mRNA sequences in the mRNA population is determined. This is done by hybridizing highly labeled cDNA with an excess of mRNA. The amount of hybridized DNA at any time  $t$  is plotted as a function of  $R_0$ . The kinetic curve yields the number of kinetic classes of mRNA molecules and their  $R_0 t_{1/2}$  values. The values are corrected to obtain the values that would have been obtained if each of the classes had been considered on its own. The values are also corrected for the nonspecific poly-dA sequences in the cDNA from the mRNA [poly A](#) tails, which do not contribute to hybridization. Then the corrected  $R_0 t_{1/2}$  values are used to determine the number of different sequences in each kinetic class by comparing them with the values of mRNA of known complexity ((3),(4)).

### References

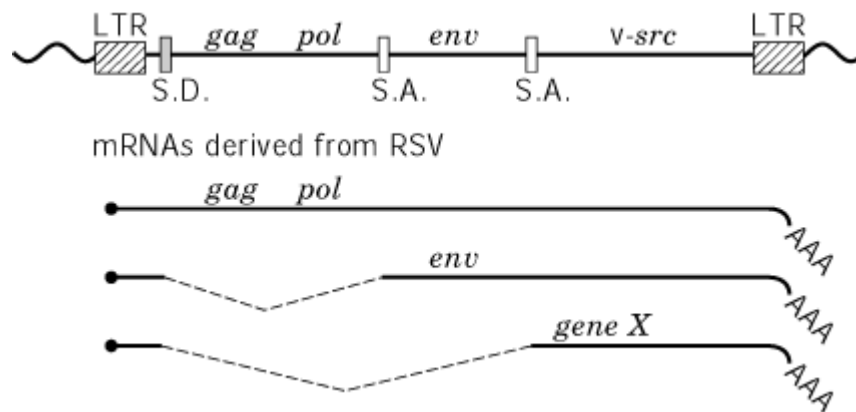
1. B. D. Hall and S. Spiegelman (1961) Proc. Natl. Acad. Sci. USA **47**, 134–146.
2. R. A. Maurer (1980) J. Biol. Chem. **255**, 854–859.
3. J. O. Bishop et al. (1974), Nature **250**, 199–204.
4. S. J. Flint (1980) In J. H. Setlow and A. Hollaender, eds., *Genetic Engineering*, Plenum Press, New York, Vol. **2**, pp. 47–82.

## Rous Sarcoma Virus (RSV)

Rous sarcoma virus (RSV) was discovered in the extract of chicken sarcoma by Dr. P. Rous in 1911 as the first infectious reagent that causes tumors, and it is now known to be the prototype oncogenic **retrovirus** (*Oncovirinae*). It is widely accepted that ALV (avian leukosis virus), a replication-competent retrovirus carrying the **genes** *gag* (internal structural proteins), *pol* (**reverse transcriptase** and DNA **endonuclease**) (see [HIV \(Human Immunodeficiency Virus\)](#) and [Hepatitis B Virus](#)), and *env* (envelope glycoproteins), had once picked up a gene present in the chicken chromosome, now designated as cellular *src* (*c-src*), during its replication cycle as a rare event. RSV therefore carries an extra gene designated viral *src* (*v-src*) (because it is derived from *c-src*) at the extreme 3' region of the genome, and it can be expressed from a subgenomic [messenger RNA](#) produced by [RNA splicing](#) (Fig. 1). By changing from *c-src* upon its integration into the virus genome and by

accumulation of mutations during the subsequent propagation of virus, this virus gene had acquired the ability to cause [neoplastic transformation](#) in chickens. High-level expression of *c-src*, however, does not cause such transformation. The products of both the *v-src* and *c-src* genes are now known to be **tyrosine-kinases**, with molecular weights of 60 kDa, and  $p60^{v-src}$  has much higher kinase activity than does  $p60^{c-src}$ .

**Figure 1.** Proviral DNA structure of Rous sarcoma virus. The line at the top indicates the chromosomal DNA containing proviral DNA. SD, splicing donor site; and SA, splicing acceptor site. Three transcripts are produced from the LTR promoter, as shown below.



This prototype virus, however, is exceptional among RNA tumor viruses in that it is replication-competent. All the other oncogenic retroviruses acquired proto-**oncogenes** at the expense of a part of their genomes that is essential for virus replication, and they therefore require the association of other replication-competent virus such as ALV (often called [helper virus](#)) for their propagation. Because of the advantage of this exceptional genome structure, RSV has been historically used in laboratory for the analysis of neoplastic transformation in cell culture; virus titers, which were originally assayed by sarcoma formation in the injected chicken, are now determined much more easily by formation of a focus (a group of transformed cells) in infected chicken embryo fibroblasts (CEFs). RSV-infected CEF assume refractile morphology and acquire the ability to grow in soft agar (anchorage-independent growth). These phenotypes unique to cellular transformation are maintained even after cell divisions, because the virus genome is also replicated as a proviral DNA that is integrated into the host [chromosome](#).

One of the breakthroughs in oncogene research was the isolation of unique **temperature-sensitive mutants** of RSV. These mutant RSVs were replication-competent independent of the temperature, but the transformation phenotypes of the infected cells were temperature-dependent. At 37°C, infected cells assume fully transformed morphology, while cells become flatter and assume very similar morphology to that of uninfected CEF at 40°C to 41°C. These observations led to the discovery of **oncogenes** and revealed the following basic ideas about them:

1. RSV encodes a gene responsible for the initiation of cellular transformation (the oncogene now known to be *v-src*, described above).
2. The oncogene is not necessary for the viral replication.
3. Continuous function of the oncogene is necessary for the maintenance of cellular transformation. (This observation excluded another hypothesis, the “hit and run theory,” in which the virus infection process is supposed to cause irreversible phenotypic changes to the host cells.)

It is also noteworthy that RSV-based avian retrovirus vectors are now widely used as a powerful tool for basic research. By transfecting plasmid DNA carrying a proviral DNA in which the v-src sequence was substituted with an exogenous gene, replication-competent virus vectors expressing the exogenous gene can be recovered from the culture fluid of the transfectants and used to introduce foreign genes into several avian primary cultures, embryos, or adults quite efficiently.

#### Suggestion for Further Reading

J. M. Coffin, S. H. Hughes, and H. E. Varmus (1977) *Retroviruses*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.

### S-Adenosyl-L-Methionine

S-Adenosyl-L-methionine (AdoMet), one of nature's most interesting biological objects, is a sulfonium compound, first described by Cantoni (1), that is synthesized in every cell from L-methionine and adenosine triphosphate (ATP) by *methionine adenosyl transferase* (MAT).



MAT catalyzes a most unusual reaction where ATP serves as the donor of its adenosine moiety, while its triphosphate chain is cleaved nonrandomly by the intrinsic triphosphatase activity of MAT itself, with the  $\alpha$  and the  $\beta$  phosphates yielding inorganic pyrophosphate ( $\text{PP}_i$ ) and the  $\gamma$  phosphate yielding orthophosphate,  $\text{P}_i$  (2). It is noteworthy that the triphosphatase activity of MAT is stimulated by AdoMet (3). MAT is found in every cell, and the amino acid sequences of enzymes isolated from species that are separated by more than a billion years in evolutionary divergence (see [Divergent Evolution](#)), such as *Escherichia coli* (4) and humans (5, 6), exhibit an extraordinary degree of [homology](#).

Elucidation of the structure of AdoMet as a sulfonium compound and clarification of the mechanism of its biosynthesis revealed that, in the reaction catalyzed by MAT, the chemical energy of the pyrophosphate bonds in ATP is utilized for the biosynthesis of a new class of energy-rich compounds. At physiological pH, the free energy of hydrolysis of a sulfonium compound is equivalent to that of the pyrophosphate bond in ATP (see [Adenylate Charge; Free Energy Relationships](#)), and the reaction is essentially irreversible, because it is accompanied by the release of a proton (7).

AdoMet was first identified as the source of methyl groups for biological transmethylation reactions (8), but it was anticipated that the energetic equivalence of the bonds linking the three ligands to the sulfur atom would enable AdoMet to function biologically as a source of either alkyl or adenosyl groups, thereby endowing AdoMet with a unique biochemical versatility (9). Experimentally, Tabor et al. (10, 11) made the important discovery that S-adenosylpropylamine, a sulfonium compound resulting from the enzymatic decarboxylation of AdoMet (12), can serve as a donor of propylamine in the synthesis of **spermidine** from putrescine, and of **spermine** from spermidine (11). AdoMet decarboxylase is particularly interesting from the mechanistic point of view, since it has a covalently bound pyruvate that is required for activity (13). AdoMet decarboxylase and the

propylaminotransferases are widely distributed in nature having been found in bacteria, yeast (14), and vertebrates. The synthesis of polyamines is of great biological importance although, quantitatively, it is less significant than biological transmethylation.

The biochemical availability of the groups linked to the sulfonium center of AdoMet was demonstrated by Nishimura et al. (15), who showed that AdoMet can serve as a donor of the aminobutyryl group to S-RNA, and more recently by Knappe (16) and by Reichardt and collaborators, who demonstrated that AdoMet can serve as an adenosine donor (17, 18). Finally Adams and Yang (19, 20) discovered in fruit ripening that AdoMet is cleaved to thiomethyl-adenosine and 1-aminocyclopropane-1-carboxylic acid. This reaction is of great biological significance, in that 1-aminocyclopropane-carboxylic acid is further metabolized to [ethylene](#), a plant [hormone](#) that initiates fruit ripening and regulates many aspects of plant growth and development. Thiomethyl-adenosine is cleaved phosphorolytically to adenine and thiomethyl ribose 1-phosphate, in the first step of a salvage pathway that results in the regeneration of methionine (21).

In these reactions, AdoMet functions as a donor of the three ligands attached to the sulfonium atom. In addition, it has been shown that its 6 amino group is utilized in the synthesis of 7,8-diaminopelargonic acid, an intermediate in the synthesis of [biotin](#) (22). The mechanism of this reaction is entirely unknown, but it demonstrates the biochemical versatility of AdoMet.

In addition, AdoMet plays a role as a positive or as a negative allosteric effector (see [Allostery](#)) in the regulation of the key reactions that affect the metabolism of *homocysteine*. This amino acid can either be remethylated to methionine by 5-methyl-tetrahydrofolate (see [Aminopterin](#), [Methotrexate](#), [Trimethoprim](#), and [Folic Acid](#)) or conjugated with serine to yield cystathionine. As an inhibitor of methylene-tetrahydrofolate reductase (23, 24), AdoMet regulates the availability of methyl-tetrahydrofolate, the key intermediate needed for the methylation of homocysteine in the *de novo* synthesis of methionine, whereas as an allosteric effector of cystathionine synthase (25) AdoMet favors commitment of homocysteine to the transsulfuration pathway.

Quantitatively, AdoMet's role of serving as a methyl donor predominates. With the exception of the *de novo* synthesis of methionine, AdoMet is the sole source of the methyl groups utilized by a myriad of substrate-specific methyltransferases for the synthesis of a great variety of compounds, such as creatine, phosphatidylcholine, yeast sterols, plant alkaloids, fungal and bacterial antibiotics, ubiquinone, triterpenes, lignins, and methyl chloride. AdoMet is also the source of the methyl groups required for the [post-translational modification](#) of proteins and nucleic acids.

Inversion of configuration at the stereospecific sulfonium center is a constant feature of the enzymatic methyl group transfer reactions (26-29). This indicates that methyltransferases, and by extrapolation alkyltransferases in general, operate by a mechanism involving direct transfer of the methyl, or alkyl, group from the sulfonium atom to the acceptor substrate, thus precluding a methylated intermediate or a ping-pong mechanism (see [Enzymes](#)).

S-Adenosyl-L-homocysteine (AdoHcy) is a product (30) and consequently a competitive inhibitor (31, 32), of all the reactions in which AdoMet participates as a methyl donor. AdoHcy also regulates fruit ripening by inhibiting the formation of ethylene (33) from AdoMet. In eukaryotes and most prokaryotes, the only pathway for the further metabolism of AdoHcy is catalyzed by *adenosylhomocysteine hydrolase* (AdoHcyase), an enzyme that, as first shown by de la Haba and Cantoni (34), catalyzes the reversible hydrolysis of AdoHcy to homocysteine and adenosine. The equilibrium of this reaction lies far in the direction of synthesis; physiologically, however, the reaction proceeds in the hydrolytic direction because homocysteine and adenosine are efficiently removed by *methionine synthase* and by adenosine deaminase, respectively. It was shown by Eloranta (35) that in most tissues the activity of AdoHcyase is 100 to 1000 times greater than the activity of MAT and that consequently the intracellular concentration of AdoHcy is 20 to 100 times smaller than that of AdoMet.



The mechanisms regulating the utilization of AdoMet by many different competing reactions are not well understood. Many methyltransferases have been purified and studied in detail. The kinetics of most are relatively simple, and the physiological activity of these enzymes appears to be related directly to the availability of methyl acceptor substrates, to their affinity ( $K_m$ ) for AdoMet, and to their inhibition constant ( $K_i$ ) for AdoHcy. Comparison of these parameters in a number of enzymes revealed that the kinetic features of different methyltransferases can be very different, and they appear to vary independently. This led to the suggestion that the intracellular ratio of AdoMet/AdoHcy, also known as the “*methylation ratio*”, might play a key role in the regulation of AdoMet utilization (36).

The validity of this hypothesis can be explored by modulating experimentally the intracellular AdoMet/AdoHcy ratio. This ratio will decrease whenever inhibition of the enzymes responsible for the removal of adenosine and/or homocysteine shifts the equilibrium of AdoHcyase towards AdoHcy synthesis, or when removal of AdoHcy generated by methyltransfer reactions is prevented by inhibition of AdoHcyase. It is also possible to take advantage of the broad specificity of AdoHcyase for adenosine by supplying adenosine analogs, like 3-deazaadenosine (DZA), that, being adenosine deaminase-resistant, can be converted to 3DZAHcy (37). Modulation of the AdoMet/AdoHcy ratio or the generation of analogues containing groups other than adenosine in response to the administration of various AdoHcyase inhibitors and/or substrates, indicates that in different tissues the physiological responses may be ascribed to the inhibition of specific methyltransferases (38). These results therefore appear to support the hypothesis that the utilization of AdoMet by various competing reactions is regulated by the intracellular AdoMet/AdoHcy ratio. It should be noted, however, that interpretation of these experiments is complex because the accumulation of AdoHcy and/or of other analogues may also result in an increase in the intracellular level of AdoMet, due to feedback inhibition of AdoMet-dependent methyl or alkyl transferases.

Unexpectedly and most importantly, AdoMet has been found to be a safe and effective agent in the treatment of certain forms of clinical depression (39, 40). While it has been clearly established that administration of AdoMet results in a significant increase in the concentration of AdoMet in the cerebrospinal fluid (CSF) (41), it has not been determined whether the therapeutic effects of AdoMet are related to its ability to function as a methyl donor. In the absence of an animal model for depressive disorders, it has unfortunately been possible only to formulate hypotheses to account for the mode of action of AdoMet in such illnesses.

## Bibliography

1. G. L. Cantoni (1953) *J. Biol. Chem.* **204**, 403–416.
2. G. L. Cantoni and J. Durell (1957) *J. Biol. Chem.* **225**, 1033–1048. (Abstract)
3. S. H. Mudd (1962) *J. Biol. Chem.* **237**, PC1372–PC1375.
4. G. D. Markham, J. DeParasis, and J. Gatmaitan (1984) *J. Biol. Chem.* **259**, 14505–14507.
5. S. Horikawa and K. Tsukada (1991) *Biochem. Int.* **25**, 81–90.
6. L. Alvarez, F. Corrales, A. Martin-Duce, and J. M. Mato (1993) *Biochem. J.* **293**, 481–486.
7. G. L. Cantoni (1960) in *Onium compounds. Handbook of Comparative Biochemistry* (M. Florkin and H. Mason, eds.), Academic Press, New York, p. 181.
8. G. L. Cantoni and P. J. Vignos, Jr. (1954) *J. Biol. Chem.* **209**, 647–659.
9. G. L. Cantoni (1952) in *A Symposium on Phosphorus Metabolism* (W. D. McElroy and B. Glass, eds.), Vol.2, p. 129, Johns Hopkins Univ. Press, Baltimore MD.
10. H. Tabor, C. M. Rosenthal, and C. W. Tabor (1958) *J. Biol. Chem.* **233**, 907–914.
11. C. W. Tabor and H. Tabor (1984) *Ann. Rev. Biochem.* **53**, 251–282.
12. C. Kutzbach and E. L. R. Stokstadt (1971) *Biochim. Biophys. Acta* **250**, 459–477.
13. G. D. Markham, C. W. Tabor, and H. Tabor (1982) *J. Biol. Chem.* **257**, 12063–12068.
14. K. K. Kashiwagi, S. K. Taneja, T. Y. Liu, C. W. Tabor, and H. Tabor (1990) *J. Biol. Chem.* **265**,

22321–22328.

15. S. Nishimura, Y. Taya, Y. Kuchino, and Z. Ohashi (1974) *Biochem. Biophys. Res. Commun.* **57**, 702–708.
16. J. Knappe, F. A. Neugebauer, H. P. Blaschkowski, and M. Ganzler (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1332–1335.
17. R. Eliasson, M. Fontecave, H. Jornvall, M. Krook, E. Pontis, and P. Reichard (1990): *Proc. Natl. Acad. Sci. USA* **87**, 3314–3318.
18. J. Harder, R. Eliasson, E. Pontis, M. D. Ballinger, and P. Reichard (1992) *J. Biol. Chem.* **267**, 25548–25552.
19. D. O. Adams and S. F. Yang (1979) *Proc. Natl. Acad. Sci. USA* **76**, 170–174.
20. Y. Yu, D. O. Adams, and S. F. Yang (1979) *Arch. Biochem. Biophys.* **198**, 280–286.
21. P. S. Backlund, Jr., C. P. Chang, and R. A. Smith (1982) *J. Biol. Chem.* **257**, 4196–4202.
22. G. L. Stoner and M. A. Eisenberg (1975) *J. Biol. Chem.* **250**, 4037–4043.
23. D. A. Jencks and R. G. Matthews (1987) *J. Biol. Chem.* **262**, 2485–2493.
24. K. Hashinaka and M. Yamada (1992) *Arch. Biochem. Biophys.* **293**, 40–45.
25. J. D. Finkelstein, W. E. Kyle, J. J. Martin, and A.-M. Pick (1975) *Biochem. Biophys. Res. Commun.* **66**, 81–87.
26. J. Luthy, J. Retey, and D. Arigoni (1969) *Nature* **221**, 1213–1215.
27. H. G. Floss, L. Mascaro, M. Tsai, and R. W. Woodard (1979) in *Transmethylation*. Elsevier/North Holland, New York, p. 135.
28. J. W. Cornforth, S. A. Reichard, P. Talalay, H. L. Carrell, and J. P. Glusker (1987) *Lancet* **99**, 7292–7300.
29. G. de la Haba et al. (1959) *J. Am. Chem. Soc.* **81**, 3975–3980.
30. G. L. Cantoni and E. Scarano (1954) *J. Am. Chem. Soc.* **76**, 4744–4745.
31. K. D. Gibson, J. D. Wilson, and S. Udenfriend (1961) *J. Biol. Chem.* **236**, 673–679.
32. R. T. Borchardt (1977) in *The Biochemistry of Adenosylmethionine* (F. Salvatore, E. Borek, V. Zappai, H. G. Williams-Ashman, and F. Schlenk, eds.), Columbia Univ. Press, New York, p. 151.
33. G. Miura and P. K. Chiang (1985) *Anal. Biochem.* **147**, 217–221.
34. G. de la Haba and G. L. Cantoni (1959) *J. Biol. Chem.* **234**, 603–608.
35. T. O. Eloranta (1977) *Biochem. J.* **166**, 521–529.
36. G. L. Cantoni and P. K. Chiang (1980) In: *Natural Sulfur Compounds* (D. Cavallini, G. E. Gaull, and V. Zappia, eds.), Plenum Press, New York, p. 67.
37. H. H. Richards, P. K. Chiang, and G. L. Cantoni (1978) *J. Biol. Chem.* **253**, 14476–14480.
38. P. K. Chiang and G. L. Cantoni (1979) *Biochem. Pharmacol.* **28**, 1897–1902.
39. G. M. Bressa (1994) *Acta Neurol. Scand. Supplement* 154, 7–14.
40. G. L. Cantoni, S. H. Mudd, and V. Andreoli (1989) *Trends Neurosci.* **12**, 319–324.
41. T. Bottiglieri and K. Hyland (1994) *Acta Neurol. Scand. Supplement* 154, 19–26.

## S1 Nuclease

S1 nuclease is a 32-kDa **nuclease** that is specific for hydrolyzing of single-stranded **RNA** or **DNA**

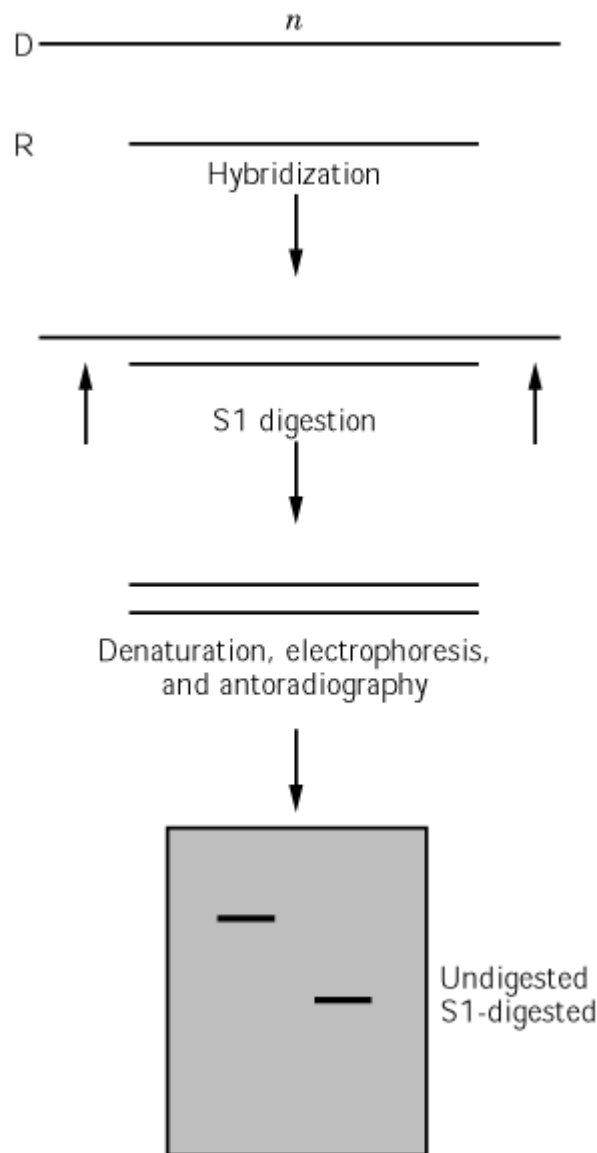
molecules into 5'-mononucleotides. It is a zinc-requiring [metalloprotein](#) that is inactivated by chelating agents, such as EDTA and citrate, and by phosphate concentrations as low as 10 mM. Its enzyme activity is optimal at pH 4.0 to 4.3, halved at pH 4.9, and negligible above pH 6.0. The protein is thermostable (1) and resistant to several **denaturing** agents, such as [urea](#), **SDS**, and **formamide** (2, 3). Digestion of single-stranded DNA is five times more efficient than that of single-stranded RNA, and 75,000 times more efficient than digestion of double-stranded DNA. The low level of cleavage of double-stranded DNA is minimized by increasing the ionic strength to about 0.2 M salt (2, 4-6), whereas it is increased by negative **supercoiling** of the double-helical DNA, UV irradiation, and depurination (3, 6-9).

Because S1 nuclease does not degrade double-stranded DNA or DNA–RNA hybrids, it is used widely to remove single-stranded regions from such duplexes. It is extremely useful for (1) measuring the extent of specific **hybridization** between single strands of DNA and/or RNA; (2) probing the existence of duplex DNA regions; (3) removing **cohesive, sticky ends** generated by [restriction enzymes](#) from single-stranded protruding DNA; (4) increasing the specificity of nucleic acid hybridization (10); (5) localizing [intron, exon](#) boundaries; (6) isolating of duplex regions in single-stranded viral [genomes](#) (11, 12); (7) probing strand breaks in duplex DNA molecules (2, 9, 13); (8) cleaving double-stranded regions that have lower duplex stability (14, 15); (9) localizing [inverted repeat](#) sequences (14, 15); (10) introducing deletion mutations at **D loops** in duplex DNA; and (11) mapping the genomic regions involved in interactions with DNA-binding proteins (16). Its most frequent use in mapping DNA sequences that encode RNA is described further here, because it exemplifies the potential and limits of the technique.

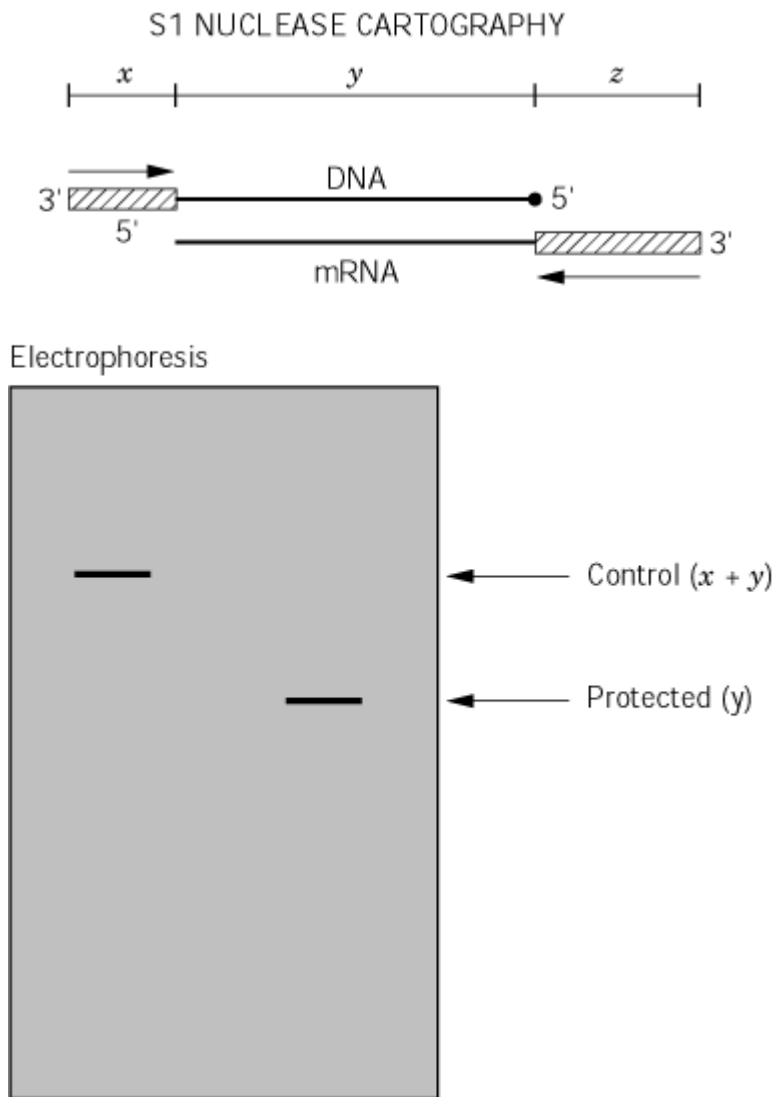
## 1. Mapping DNA sequences that encode RNA

A method for mapping RNA molecules onto the DNA templates from which they were **transcribed** was first described by Berk and Sharp (17, 18). The rationale is illustrated in Figure 1. An RNA transcript R is complementary to  $n$  nucleotides of the coding strand of a DNA, D. Hybridization of the denatured DNA and RNA leads to a hybrid duplex molecule that has single-stranded overhangs. These are removed by digestion with S1 nuclease. The size of the resulting DNA fragment, that is, the value of  $n$ , can be determined by subjecting it to [electrophoresis](#) on a [polyacrylamide](#) gel like those used in [DNA sequencing](#). S1 nuclease does not remove the [poly A](#) track or the **-cap** of eukaryotic [messenger RNA](#) when they are base-paired, but the 5'-cap structure, it has been thought, sterically hinders S1 digestion of single-stranded mRNA at the phosphodiester bond adjacent to the first nucleotide (19). When mapping is done with an end-labeled probe, one can determine the polarity and the map position of the RNA on the corresponding DNA sequences (19). This strategy is very useful for mapping the 5'-end of mRNA transcripts, thereby identifying promoter sequences (Fig. 2), and mapping the exons in spliced eukaryotic mRNA (Fig. 3).

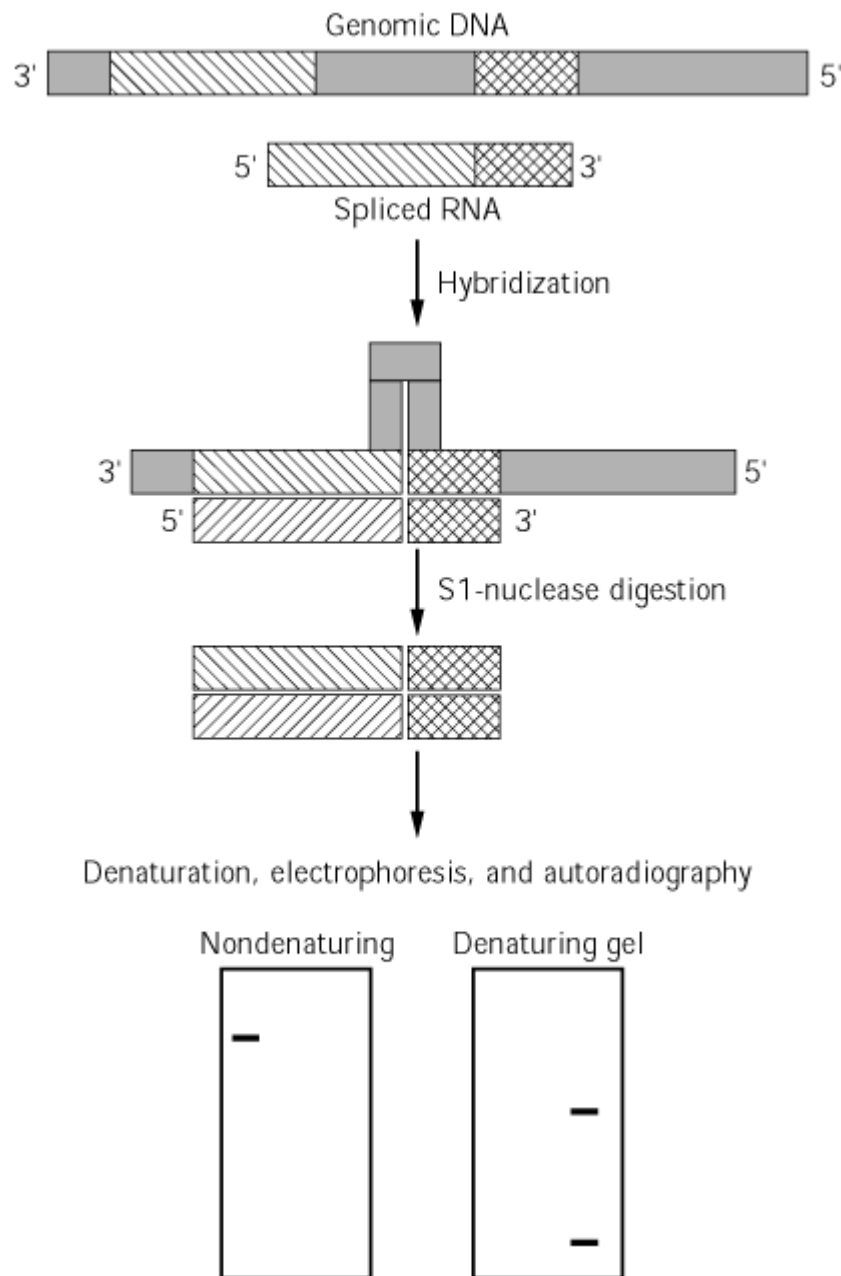
**Figure 1.** Rationale for S1-nuclease mapping



**Figure 2.** Mapping of 5'-proximal promoter sequences. Promoter sequences localized upstream of the 5' terminus of eukaryotic spliced RNAs can be mapped after hybridization of RNA preparations with a labeled DNA template (corresponding to the coding strand) that contains the putative promoter sequences (as suggested by DNase hypersensitivity, sequence analysis, etc.). In this example, a DNA fragment XYZ encodes an RNA complementary to sequence YZ. The promoter sequences are localized in X. The 5'-end of the RNA delineates the X-Y junction. The first step consists of preparing a single-stranded, 5'-labeled DNA probe (XY) that contains the 5'-proximal promoter sequences and does not cover the 3-end of the RNA molecule. This conserves the 5'-labeled end (the dark sphere) after S1 digestion. Then, this DNA strand (top) is hybridized to the mRNA (bottom). The DNA and RNA overhangs are digested by S1 nuclease, and hybridized regions are protected. The S1-digested hybrid is denatured before to analysis by denaturing gel electrophoresis. Only the labeled DNA fragments are visualized by autoradiography. Left lane: control undigested labeled probe (XY); right lane: S1-digested sample in which only the protected Y region of labeled probe is protected. The length of the protected DNA fragment (Y) is smaller than that of the initial probe (XY). The difference in size (X) identifies the starting point (5'-end) of the RNA molecule.



**Figure 3.** Mapping exon-intron junctions. In this example, a spliced mRNA composed of two exons is hybridized with the complementary, labeled DNA template (coding strand). Following hybridization of complementary sequences, a loop of single-stranded DNA (corresponding to the intron) is generated. Upon incubation of the partial DNA–RNA duplex with S1 nuclease, the free 3'- and 5'-proximal DNA overhangs and the loop are digested. The size of the resulting DNA–RNA hybrid, which contains only exon sequences, is determined by gel electrophoresis in a non-denaturing buffer. Establishing the individual sizes of the two exons requires gel electrophoresis under denaturing conditions to dissociate the hybrid strands. Only the labeled DNA fragments are observed by autoradiography.



## Labeling the DNA fragment

Determining the size of the S1-resistant DNA fragment relies on measuring a difference in size between the initial DNA probe and the protected segment. The DNA is usually detected specifically by [autoradiography](#), which requires that it be **radiolabeled** and that the label be within the protected segment and survive the S1 nuclease digestion. Cloned DNA fragments are labeled uniformly or specifically at either the 5'- or 3'-termini. Labeling at the 5'-terminus is usually done with **polynucleotide kinase** in the presence of  $\gamma$ - $^{32}\text{P}$ -ATP, whereas 3'-labeling is achieved with  $\alpha$ - $^{32}\text{P}$ -nucleoside triphosphates and either the **Klenow fragment** of *Escherichia coli* **DNA polymerase I**, T4 polynucleotide kinase, or nucleotidyl terminal transferase.

## 2. Preparing single-strand templates

The previous labeling methods label both strands, so it is necessary to purify the labeled coding strand or to use conditions that favor hybridization between the RNA and DNA strands (see later). Generally, the best results are obtained with separated strands. After denaturation, the two strands of a DNA fragment may be separated by electrophoresis in polyacrylamide or [agarose](#) gels if they

have different base compositions. It is not possible to predict how each strand will migrate, however, so the S1 mapping must be done with each strand. Hybridization of polyUG to the denatured strands may accentuate differences in their electrophoretic mobilities (20, 21). The double-stranded DNA fragment is digested asymmetrically from its two ends, so that the two single strands have different mobilities after denaturation (see Fig. 3). Alternatively, single-stranded DNA may be isolated by cloning in filamentous **bacteriophage** or plasmid vectors based upon such phage (see [DNA Sequencing](#)).

### 2.1. Use of Double-Stranded DNA Fragments

The rate of hybridization of RNA molecules to single strands of DNA is similar to that of DNA single-strand renaturation (22, 23). Therefore, the rate of hybridization is maximal in the presence of a high salt concentration (1 M NaCl) and a temperature corresponding to the **melting temperature** ( $T_m$ ) of 25 to 30°C. In practice, hybridization is often done in the presence of formamide to decrease the  $T_m$  and to hybridize at a lower temperature (24, 25). The melting temperature of duplex DNA molecules can be calculated (26) from the base composition, the ionic strength, and the formamide concentration by using the equation

$$T_m(^{\circ}\text{C}) = 81.5 + 0.5(\%G + C) + 16.6 \log[\text{Na}^+] - 0.6(\%\text{formamide}) \quad (1)$$

At high concentrations of formamide (70 to 80%), the  $T_m$  of a DNA-DNA duplex is 5 to 10°C lower than the  $T_m$  of the corresponding RNA-DNA hybrid (26). Under such conditions, a hybridization temperature between these two  $T_m$  values considerably favors association of the RNA with its coding DNA strand. Most RNA-DNA hybridizations are done at temperatures between 40 and 60°C with a buffer that provides constant pH and ionic strength, such as 0.04 M NaCl and 1 mM EDTA. To establish the optimal conditions for hybridization, the  $T_m$  of the DNA duplex used should be determined experimentally. Then, hybridization with the RNA is done at a temperature 2 to 4°C higher. The best results are usually obtained when excess DNA is used. Under such conditions, RNA-DNA hybridization is nearly complete in 3 hours with adenovirus 2 DNA (35 kbp) at a concentration of 10 µg/mL (27). Extrapolating, it can be predicted that nearly complete hybridization of a DNA that has  $n$  kilobases will be obtained under similar conditions in 3 hours at a DNA concentration of  $(n/3.5)$  mg/mL.

In practice, the labeled, double-stranded DNA and the RNA molecule are mixed and reextracted together once. After resuspending the pellet in a “double-strand hybridization buffer” of 50 mM PIPES (pH 6.4), 1 mM EDTA, 0.4 M NaCl, 80% formamide, the mixture is incubated at 65°C to denature the nucleic acids. Then, it is transferred to the appropriate temperature for hybridization.

### 2.2. Use of Single-Stranded DNA Fragments

The probe and the target single-stranded DNA are resuspended in a “single-strand hybridization buffer” of 0.3 M PIPES (pH 6.4), 30 mM EDTA, 2.5 M NaCl, and incubated as before.

### Bibliography

1. T. Ando (1966) *Biochim. Biophys. Acta* **114**, 158–168.
2. V. M. Vogt (1973) *Eur. J. Biochem.* **33**, 192–200.
3. H. Hofstetter, A. Schambock, Van den Berg, and C. Weissmann (1976) *Biochim. Biophys. Acta* **454**, 587–591.
4. T. E. Shenk, R. P. W. Rhodes, and P. Berg (1975) *Proc. Natl. Acad. Sci. USA* **72**, 989–993.
5. V. M. Vogt (1980) in *Methods in Enzymology*, **65**, (L. Grossman and K. Moldave, eds.),

Academic Press, New York, pp. 248–255.

6. P. Beard, J. F. Morrow, and P. Berg (1973) *J. Virol.* **12**, 1303–1313.
7. G. N. Godson (1973) *Biochim. Biophys. Acta* **308**, 59–67.
8. M. Mechali, A. M. Recondo, and M. Girard (1973) *Biochem. Biophys. Res. Commun.* **15**, 1306–1320.
9. K. Shishido and T. Ando (1975) *Biochim. Biophys. Acta* **390**, 125–132.
10. B. J. Bratina, M. Viebahn, and T. M. Schmidt (1996) *Methods Mol. Cell. Biol.* **6**.
11. K. Shishido and Y. Ikeda (1970) *J. Biochem.* **67**, 759–765.
12. K. Shishido and Y. Ikeda (1971) *Biochem. Biophys. Res. Commun.* **44**, 482–489.
13. J. E. Germond, V. M. Vogt, and B. Hirt (1974) *Eur. J. Biochem.* **16**, 591–600.
14. D. M. Lilley (1980) *Proc. Natl. Acad. Sci. USA* **77**, 6468–6472.
15. N. Panayotatos and R. Wells (1981) *Nature* **289**, 466–470.
16. R. Meyer, J. Grassberg, J. Scott, and A. Kornberg (1980) *J. Biol. Chem.* **255**, 2897–2901.
17. A. J. Berk and P. A. Sharp (1977) *Cell* **12**, 721–732.
18. A. J. Berk and P. A. Sharp (1978) *Proc. Natl. Acad. Sci. USA* **75**, 1274–1278.
19. R. F. Weaver and C. Weissman (1979) *Nucleic Acids Res.* **7**, 1175–1193.
20. G. S. Hayward (1972) *Virology* **49**, 342–344.
21. R. W. Goldbach, R. F. Evans, and P. Borst (1978) *Nucleic Acids Res.* **5**, 2743–2754.
22. J. G. Wetmur and N. Davidson (1968) *J. Mol. Biol.* **31**, 349–370.
23. R. J. Britten and D. E. Kohne (1968) *Science* **161**, 529–540.
24. Bonner, K. Kugn, and Beklow (1967) *Biochemistry* **6**, 3650–3653.
25. B. L. McConaughy, C. D. Laird, and B. I. McCarthy (1969) *Biochemistry* **8**, 3289–3295.
26. J. Casey and N. Davidson (1978) *Nucleic Acids Res.* **5**, 1539–1552.
27. P. A. Sharp, A. J. Berk, and S. M. Berget (1980) in *Methods in Enzymology*, **65**, (L. Grossman and K. Moldave, eds.), Academic Press, New York, pp. 750–768.

## Salt Bridge (Salt Linkage, Ionic Bond)

A salt bridge is a pair of ionized groups with opposite charge that are in very close proximity in a macromolecule and contact each other at their **van der Waals radii**. The ion pair is stabilized by attractive [electrostatic interaction](#) and may even have some [hydrogen bond](#) character. Intra- and intermolecular salt bridges are frequently observed in proteins (1) and at protein–protein and protein–nucleic acid interfaces (2).

### Bibliography

1. D. J. Barlow and J. M. Thornton (1983) *J. Mol. Biol.* **168**, 867–885.
2. C. O. Pabo and R. T. Sauer (1992) *Ann. Rev. Biochem.* **61**, 1053–1095.



## Saltatory DNA Replication

Particles of **SV40 virus** are produced in **permissive** simian cells. Infection of nonpermissive rodent cells by SV40 results in integration of the viral [genome](#) into the host [chromosome](#). The transformed rodent cells have no free SV40 genomes, but SV40 [genomes](#) are produced when these nonpermissive cells are fused with permissive simian cells. The chromosomes from the monkey cells provide replicative functions not present in the rodent cells, making SV40 [DNA replication](#) possible in the [heterokaryons](#). Repeated rounds of [initiation of DNA replication](#) from the integrated SV40 [replication origins](#) result in polytenization of that region (see [Polytene Chromosome](#)), which may liberate a linear DNA duplex containing the integrated SV40 genome. The liberated DNA will then **recombine** with itself to form a circular viral genome. In cell lines with tandem duplications of the integrated SV40 DNA, the efficiency of infectious virus formation is increased because of the terminal duplication of the liberated linear DNA. In the absence of tandem duplication, cyclization by nonhomologous recombination may form mutants with the deletion of viral DNA or insertion of host DNA ([1](#)).

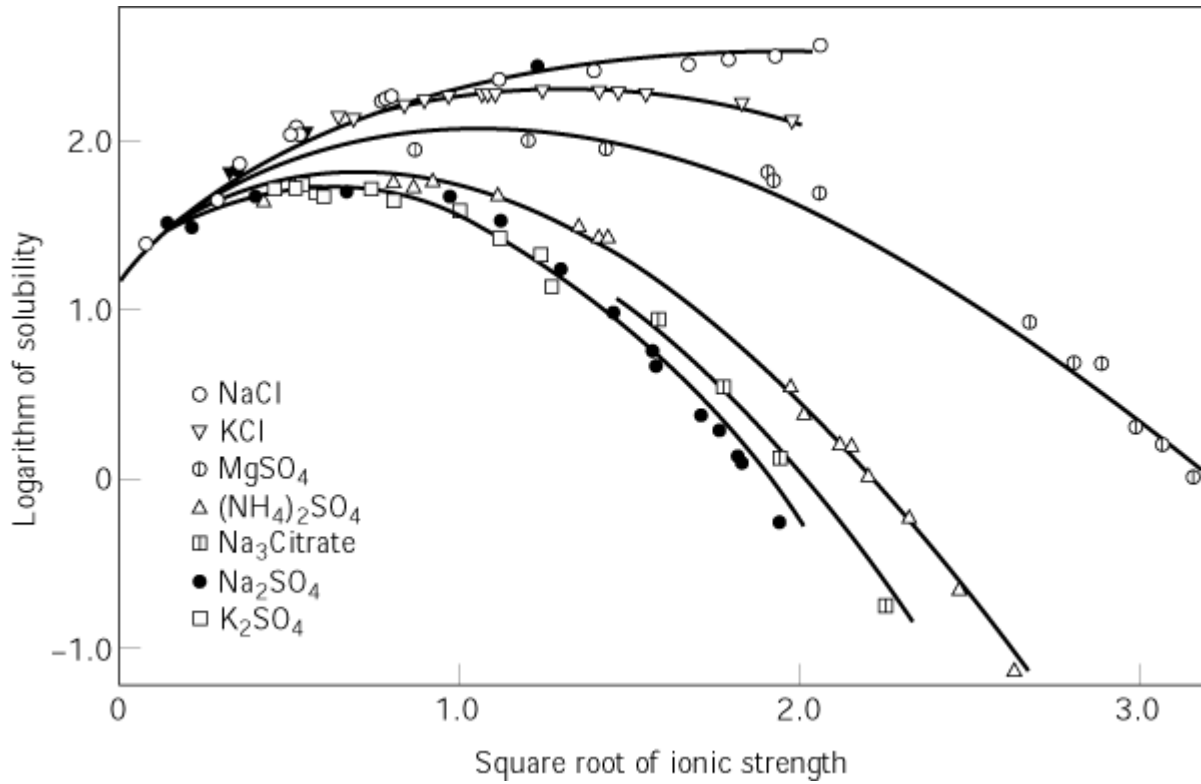
### Bibliography

1. A. Kornberg and T. A. Baker (1992) In *DNA Replication*, 2nd ed., W. H. Freeman, New York, pp. 690–699.

## Salting In, Salting Out

The effect of salts on [protein](#) solubility is described by the salting in, salting out effect ([1](#)): The dependence of the solubility of a protein on salt concentration is one in which, for any salt, at low salt concentration the solubility of a protein increases; then, after the protein passes a point of maximal solubility, the further addition of salt reduces its solubility (see [Precipitation](#)). Typical patterns are shown in Figure [1](#) ([2](#)).

**Figure 1.** The solubility of carboxyhemoglobin in various electrolytes at 25°C. (Reprinted from [2](#)) with permission of The American Society of Biological Chemists.)



This complex dependence has been decomposed into two contributions (3). At low salt concentration, proteins are salted in; in other words, their solubility increases. For any protein, this is independent of the nature of the salt and is simply the consequence of the fact that proteins carry positively and negatively charged groups on their [accessible surface](#). These are surrounded by an atmosphere of the ions of the salt, which increases the solubility according to the simple *Debye–Hückel equation*. Thus, if  $s_s$  and  $s_w$  are the solubilities of the protein in salt solution and pure water, respectively, the salting in is described by the following equation:

$$\log(s_s/s_w) \simeq Z^2 \ell I^{\frac{1}{2}} \quad (1)$$

where  $Z$  is the net charge of the protein,  $I$  its ionic strength, and  $\ell$  a combination of constants.

At high salt concentration (approximately  $> 0.5 M$ ), salting in becomes saturated, and the effect of salting out takes over, ie, the protein solubility decreases, following the law that the logarithm of the solubility ( $S$ , usually given in units of grams of protein per  $1000 \text{ cm}^3$  of solvent) is linear in salt concentration expressed as ionic strength:

$$\log S = \beta - K_s I \quad (2)$$

where  $b$  and  $K_s$  are empirical constants. The parameter  $K_s$ , called the salting out constant, is characteristic of the particular salt. The order of salting out effectiveness is close to invariant for most proteins and follows the [Hofmeister series](#) (2). The magnitude of  $K_s$  is defined by the preferential interaction of the particular salt with the protein, as

$$K_s = (1/2.303RT)(\partial\mu_{pr}/\partial m_s)_{T,P,m_{pr}} \quad (3)$$

where  $m_s$  is the molal concentration of the salt (3). Salts that induce [preferential hydration](#) (*negative binding* in [equilibrium dialysis](#); see [Binding](#)) reduce the solubility of proteins. Conversely, a solvent system that shows [preferential binding](#) increases the solubility of proteins.

### Bibliography

1. E. J. Cohn and J. D. Ferry (1943) in *Proteins, Amino Acids and Peptides* (E. J. Cohn and J. T. Edsall, eds.), Van Nostrand Reinhold, New York, pp. 586–622.
2. A. A. Green (1932) *J. Biol. Chem.* **95**, 47–66.
3. T. Arakawa and S. N. Timasheff (1985) *Meth. Enzymol.* **114**, 49–77.

## Salvage Pathways To Nucleotide Biosynthesis

Most **nucleotide** biosynthesis in most cells occurs via the nearly ubiquitous *de novo* synthetic pathways, starting from [amino acids](#) and their derivatives (see [Purine Ribonucleotide Metabolism](#) and [Pyrimidine Ribonucleotide Metabolism](#)). However, most cells possess capabilities for taking up nucleosides and nucleobases and converting them to nucleotides. Because these processes involve reutilization of previously synthesized purine and pyrimidine rings, they are called *salvage pathways*. These pathways are much shorter and simpler than the *de novo* pathways. On the other hand, there is much more variability from organism to organism, and from tissue to tissue in the same organism, in salvage synthetic capabilities. The metabolic importance of salvage pathways has come to light, first, along with the realization of the serious consequences of hereditary deficiencies in certain salvage enzymes and, second, with the many ways in which salvage pathways are being exploited to create or enhance the effectiveness of chemotherapy against a variety of diseases.

### 1. Transport of Nucleosides and Nucleobases

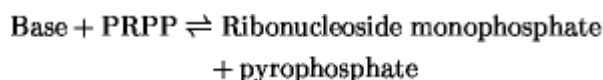
Substrates for salvage pathways can come from intracellular nucleic acid degradation (see [DNA Degradation In Vivo](#)). For most cells, however, salvage precursors are derived from the extracellular environment. An extreme example is protozoan parasites, which derive all their nucleic acid precursors from salvage substrates in the blood of infected organisms. In fact, these parasites have evolved to the extent that they completely lack *de novo* pathways and are totally dependent on salvage (1, 2); development of specific salvage inhibitors now constitutes one of the most active research areas of biochemical parasitology.

Nucleosides are hydrophilic molecules and do not readily diffuse through membranes; to an extent, the same is true for nucleobases. Thus, most cells contain specific **transporter** systems, which take up these molecules by **facilitated diffusion** (3, 4). In recent years, however, concentrative transport systems have been described, especially for nucleosides. These systems depend on sodium and involve cotransport of the nucleoside with  $\text{Na}^+$ . Most animal cells contain a broad-specificity nucleoside transporter that functions by facilitated diffusion, the active transport systems being limited primarily to specialized cells, for example, those that take up and use adenosine as a physiological regulator. The nucleobase transport systems are less well characterized, and the

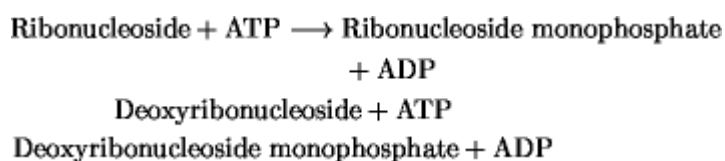
question of whether cells contain multiple base transporters or a single broad-specificity transporter has not yet been resolved.

## 2. Salvage Pathways

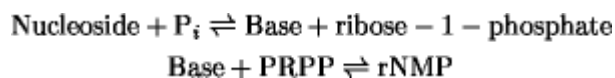
Nucleobases are anabolized by phosphoribosyltransferases, which use 5-phosphoribosyl-1-pyrophosphate (PRPP),



while nucleosides are anabolized primarily by [kinases](#).



Nucleosides can also be salvaged by a nucleoside phosphorylase followed by a phosphoribosyltransferase.



In principle, this route could be used for nucleotide synthesis, starting with a base, given that nucleoside phosphorylases are readily reversible. However, these enzymes generally act within cells in the direction of nucleoside degradation.

Salvage enzymes vary considerably in tissue distribution and specificity. Mammalian cells salvage purines primarily at the free base level through the action of two phosphoribosyltransferases: hypoxanthine guanine phosphoribosyltransferase (HGPRT), which converts hypoxanthine and guanine to inosinic acid (IMP) and GMP, respectively; and adenine phosphoribosyltransferase (APRT), which specifically converts adenine to AMP. HGPRT is the enzyme missing in cells of males suffering from Lesch–Nyhan syndrome (see [Purine Ribonucleotide Metabolism](#)). Pyrimidines, by contrast, are more likely to be salvaged at the nucleoside level. In the *de novo* pathway, orotic acid is converted to orotidylate by a phosphoribosyltransferase, but comparable pyrimidine salvage enzymes are of extremely limited distribution.

Protozoa that rely on nucleotide salvage pathways, particularly for purines, bear a distinctive set of enzymes— nucleoside hydrolases. These organisms take up purine nucleosides from the blood of infected animals, for example, and then hydrolyze the nucleoside to ribose plus the base; the base is then salvaged by a phosphoribosyltransferase (5). Because the nucleoside hydrolases are quite distinct from known proteins, they are being explored as targets for drugs that would specifically inhibit the growth of an infecting parasite by blocking this salvage route, which is essential for the parasite but not the host.

As discussed more thoroughly below, nucleoside and base analogues are widely used as antiviral, antimicrobial, and anticancer drugs. Because of the interest in developing or improving the

effectiveness of drugs that act as DNA synthesis inhibitors, there has been particular emphasis on the deoxyribonucleoside kinases (6). Four such enzymes, with overlapping specificities, are found in human cells. Properties and subcellular distributions of these enzymes are summarized in Table 1. It is of particular interest that deoxycytidine kinase (dCK) is also the principal enzyme for utilization of purine deoxyribonucleosides (the  $K_m$  values are much higher for the purines than for deoxycytidine; hence the name of the enzyme).

**Table 1. Properties of the Human Deoxyribonucleoside Kinases**

| Enzyme Cycle Regulation     | Natural Substrates | Intracellular Distribution | Cell Cycle Regulation |
|-----------------------------|--------------------|----------------------------|-----------------------|
| Deoxycytidine kinase (dCK)  | dCyd, dAdo, dGuo   | Cytosol                    | Constitutive          |
| Thymidine kinase 1 (TK1)    | dThd, dUrd         | Cytosol                    | S-phase-specific      |
| Thymidine kinase 2 (TK2)    | dThd, dUrd, dCyd   | Mitochondria               | Constitutive          |
| Deoxyguanosine kinase (dGK) | dGuo, dAdo, dIno   | Mitochondria <sup>o</sup>  | Constitutive          |

Source: Arnér and Eriksson (6)

A key to understanding some clinical disorders is knowledge of the tissue distribution of salvage enzymes. For example, the immunodeficiency state that results from adenosine deaminase deficiency can be related to this factor (see [Purine Ribonucleotide Metabolism](#)). Most or all tissues in an affected individual accumulate adenosine and deoxyadenosine as a result of the genetic block to purine catabolism. Nucleoside salvage enzymes are particularly active in the reticuloendothelial system, where [apoptosis](#) leads to degradation and reutilization of cell components, including nucleic acids. Thus, blood cells accumulate excessive amounts of ATP, and in particular, dATP, which interferes with white cell proliferation as part of the [immune response](#), through its inhibition of ribonucleotide reductase.

A salvage enzyme of considerable interest is the family of thymidine kinases specified by [herpes viruses](#). As exemplified by the herpes simplex thymidine kinase, these enzymes have two distinctive properties. First, they are bifunctional enzymes, with thymidylate kinase activity as well. Thus, the enzyme is evidently designed to catalyze two sequential reactions. Second, the nucleoside kinase activity has much broader substrate specificity than most cellular nucleoside kinases. As discussed below, antiviral chemotherapeutic strategies exploit this distinctive property.

### 3. Salvage Enzymes as Selectable Markers

Because nearly all cells possess *de novo* nucleotide synthetic capabilities, the salvage enzymes are usually not required for cell viability. In addition, a large number of nucleobase and nucleoside

antimetabolites are available, most of them inhibitors of specific enzymes (see [Nucleotides, Nucleosides, And Nucleobases](#)). These factors create favorable conditions for the use of salvage enzymes as selectable genetic markers, ie, genetic characteristics that promote the selective survival or growth of desired cell types. For example, the mammalian enzyme encoding HGPRT is widely used as a selectable marker in genetic analysis. 6-Thioguanine is metabolized by HGPRT to give the thiol analogue of inosinic acid, which is toxic. Cells lacking HGPRT can be selected for because they grow in 6-thioguanine-containing medium, while other cells are killed. Hence, one can estimate [mutation](#) rates by culturing wild-type cells in the presence of thioguanine and enumerating the cells that grow. Another advantage of this gene for genetic analysis is that it is carried on the [X-chromosome](#). Thus, in male-derived cell lines, only one mutation, rather than two independent events is required to give the resistant phenotype.

A comparable selectable marker is the gene encoding thymidine kinase. 5-Bromodeoxyuridine is anabolized similarly to thymidine;



However, incorporation of bromodeoxyuridine into DNA is a lethal event. Thus, thymidine kinase deficiency leads to a BrdUrd-resistant phenotype. Because most large DNA viruses encode a thymidine kinase, this is a useful system for manipulating viral genes, for example, the use of [vaccinia virus](#) as a vector in the generation of multivalent vaccines (7).

On the other hand, one can select for the presence of active salvage enzymes. The best example is the use of “HAT medium” in **somatic cell** genetic analysis and in preparing [monoclonal antibodies](#). These techniques involve fusing cells of different origins and culturing in HAT medium to select for those cells that have undergone fusion. HAT is an acronym for the medium's constituents,— hypoxanthine, aminopterin, and thymidine. Aminopterin inhibits [dihydrofolate reductase](#), blocking the synthesis of tetrahydrofolate needed for *de novo* synthesis of purine nucleotides and thymidine nucleotides. Thus, cells can grow in HAT medium only if they express active thymidine kinase and HGPRT, for salvage synthesis of thymidine and purine nucleotides, respectively. In monoclonal antibody production, one of the cell lines to be fused lacks thymidine kinase, and the other lacks HGPRT. Thus, only cells resulting from a fusion event have functional copies of both enzymes and can grow.

A final example relates to the use of a selectable marker to force expression of a nonselected marker. For [cloning](#) into [expression systems](#) in mammalian cells, one often incorporates into the cloning vector the *Escherichia coli xpt* gene, which encodes a distinctive phosphoribosyltransferase that acts on xanthine and guanine. During and after transformation of cells with the recombinant DNA, the cells are cultured in the presence of mycophenolic acid, which blocks *de novo* guanine nucleotide synthesis by inhibiting IMP dehydrogenase, the enzyme that converts IMP to XMP (which would then be converted to GMP). Thus, the only cells that can grow are those that have taken up and expressed the *xpt* gene, which bypasses this metabolic block. Being carried on the same vector, the gene of interest is also cloned and/or expressed, even though its expression was not directly selected for.

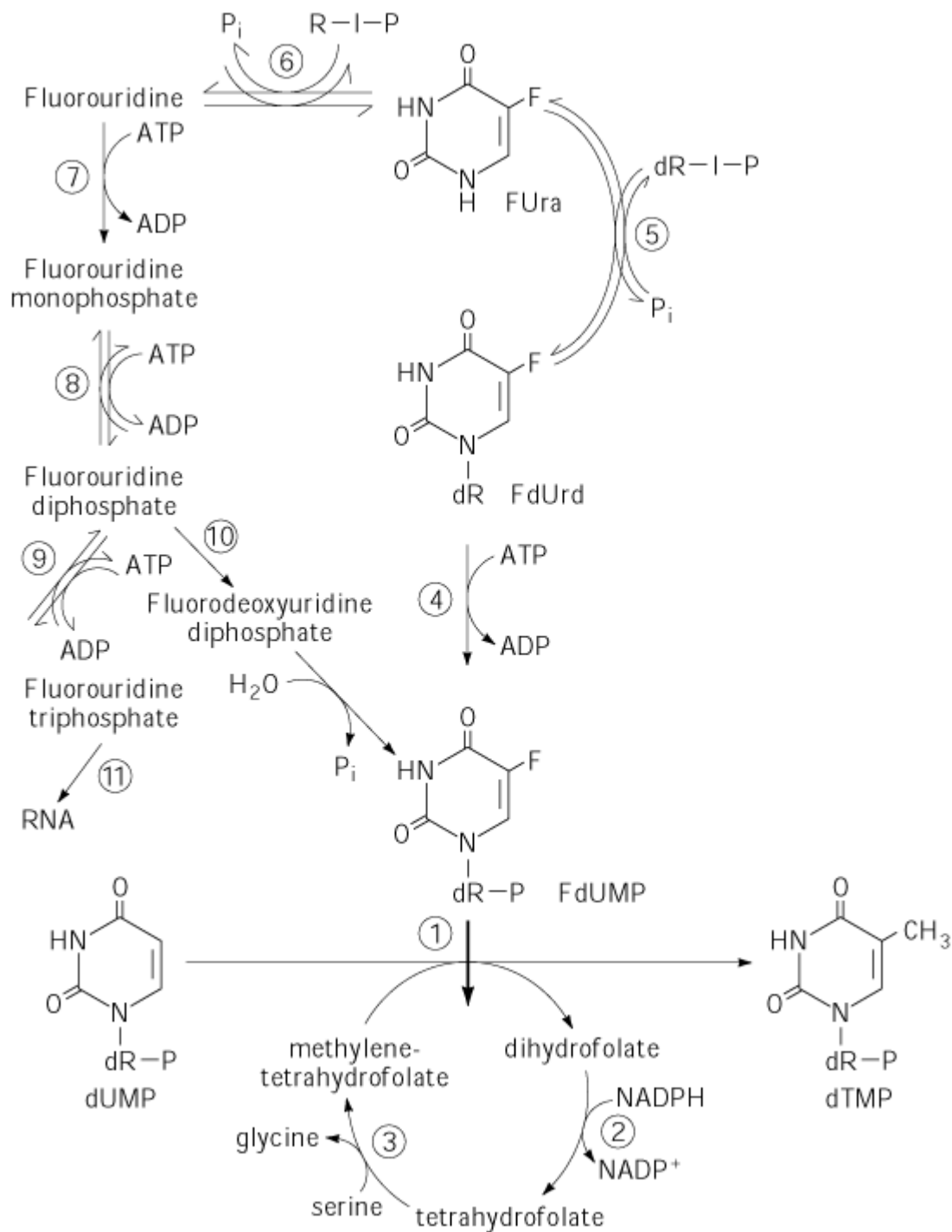
#### 4. Salvage Enzymes and Molecular Pharmacology

A large number of antiviral, antibacterial, antiparasitic, and anticancer drugs act by inhibiting or otherwise manipulating pathways in nucleotide metabolism and thereby interfering with nucleic acid synthesis. Several examples are cited in [Nucleotides, nucleosides, and nucleobases](#). In most cases, the active species is a nucleotide. Because nucleotide molecules are charged, and because specific transport systems for nucleotides don't exist, these molecules penetrate membranes poorly, if at all.

Therefore, in order to generate a therapeutically effective nucleotide intracellularly, its precursor must be administered extracellularly as a nucleobase or nucleoside analogue. Effective use of such drugs demands extensive understanding of salvage pathways in the target cell, subcellular distribution of the enzymes involved, degradative enzymes that might compete with the pathway leading to the desired nucleotide analogue, levels of the target enzyme, and cell-cycle regulation of all the enzymes involved (8).

As an example of the factors involved, consider the fluorinated pyrimidines, 5-fluorouracil (FUra) and 5-fluorodeoxyuridine (FdUrd), analogues used for four decades in treating various cancers. It was established in 1958 (9) that these drugs are converted *in vivo* to 5-fluorodeoxyuridine monophosphate (FdUMP), an analogue of deoxyuridine monophosphate, the substrate for [thymidylate synthase](#) (Fig. 1; see [Deoxyribonucleotide Biosynthesis And Degradation](#)), and that FdUMP is a potent inhibitor of thymidylate synthase and, hence of [DNA replication](#). Figure 1 shows also the metabolic pathways that both activate these analogues and divert them from their desired endpoint (10). From the figure, one can see that coadministration with FdUrd of a thymidine phosphorylase inhibitor should increase the effectiveness of the analogue by blocking its catabolism. Note that there are multiple routes for activation of FUra; note also that FUra can enter pools of RNA precursors which, in principle, could limit its selectivity by diminishing the specificity of its effect against DNA synthesis. There is evidence, however, that, in some tumors, the effectiveness of FUra actually depends in part on its incorporation into RNA, where it stimulates translational miscoding.

**Figure 1.** Metabolism of fluorinated pyrimidines by salvage and other enzymes. dR-P is deoxyribose 5'-phosphate, dR-1-P deoxyribose-1-phosphate, and R-1-P ribose-1-phosphate. Enzymes are 1, thymidylate synthase; 2, dihydrofolate reductase; 3, serine transhydroxymethylase; 4, thymidine kinase; 5, thymidine phosphorylase; 6, uridine phosphorylase; 7, uridine kinase; 8, uridylate kinase; 9, nucleoside diphosphate kinase; 10, ribonucleoside diphosphate reductase; 11, RNA polymerase.



## Bibliography

1. R. L. Berens, E. C. Krug, and J. J. Marr (1995) in *Biochemistry of Parasitic Organisms and its Molecular Foundation* (J. J. Marr and M. Muller, eds.), Academic Press, London, pp. 89–117.
2. D. J. Hammond and W. E. Gutteridge (1996) *Mol. Biochem. Parasitol.* **13**, 243–261.
3. P. G. W. Plagemann, R. M. Wohlhueter, and C. Woffendin (1988) *Biochim. Biophys. Acta* **947**, 405–443.
4. D. A. Griffith and S. M. Jarvis (1996) *Biochim. Biophys. Acta* **1286**, 153–181.
5. D. W. Parkin (1996) *J. Biol. Chem.* **271**, 21713–21719.
6. E. S. J. Arnér and S. Eriksson (1995) *Pharmacol. Ther.* **67**, 155–186.



7. D. W. Grosenbach and D. E. Hruby (1995) in *Methods in Molecular Genetics: Molecular Virology Techniques, Part B* (K. W. Adolph, ed.), **7B**, 45–64.
8. R. I. Christopherson and S. D. Lyons (1990) *Medicinal Res. Rev.* **10**, 505–549.
9. S. S. Cohen, J. G. Flaks, H. D. Barner, M. R. Loeb, and J. Lichtenstein (1958) *Proc. Natl. Acad. Sci. USA* **44**, 1004–1012.
10. B. Ardalan and R. Glazer (1981) *Cancer Treatment Reviews* **8**, 157–167.

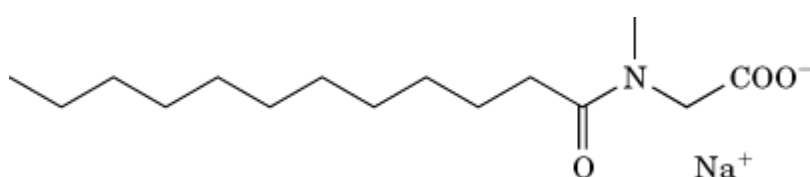
### Suggestions for Further Reading

11. Reference 8 above is a thorough review of purine and pyrimidine analogues as chemotherapeutic agents.
12. G. B. Elion (1989) *Science* **244**, 41–47; describes the history of the development of 6-thioguanine as an antileukemic drug, allopurinol as a gout treatment, and acyclovir as an antiviral agent.
13. H. Mitsuya, ed. (1997) *Anti-HIV Nucleosides: Past, Present, and Future*, R. G. Landes Co., Georgetown, TX. An informative reference, with five review articles discussing strategies for developing more effective nucleoside analogues for treating HIV infections.

### Sarkosyl

Sodium *N*-lauroylsarcosine, or Sarkosyl, is a [detergent](#) that is structurally related to **SDS**, except that Sarkosyl has an additional polar and rigid [peptide bond](#) linkage within the **hydrophobic** backbone (Fig. 1). Such rigidity added to the N-terminus of the hydrophobic chain could lead to a decrease in its ability to be inserted freely into hydrophobic [membrane](#) bilayers and proteins. This could explain why Sarkosyl is milder than SDS in its ability to **denature** and disrupt membrane and protein structures. Because of a similar difference in side-chain structure, the bile acid detergent CHAPS is rendered milder than the strongly denaturing sodium cholate, which inactivates integral membrane proteins, such as the serotonin 1A receptor (1). The relatively mild nature of Sarkosyl has been exploited at least in a few cases. Important examples include the purification of *Escherichia coli* **RNA polymerase** sigma factors by solubilizing [inclusion bodies](#) (2) and the solubilization and purification of the [scrapie](#) protein, the [prion](#) PrP<sup>Sc</sup> (3, 4). Although comparison of the structural features of Sarkosyl with those of other detergents suggests that it might be effective in solubilizing functionally active membrane proteins, the major use of Sarkosyl has been to isolate **DNA** and **RNA**. The principal reason for this is that Sarkosyl is highly efficient in dissociating [nucleosomes](#) and [ribosomes](#) (5). The detergent also denatures **nucleic acids** and inhibits certain [enzymes](#), such as **deoxyribonucleases** and **ribonucleases**, which could degrade these molecules.

**Figure 1.** The structure of Sarkosyl.



## Bibliography

1. P. Banerjee, J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
2. R. R. Burgess (1996) *Methods Enzymol.* **273**, 145–149.
3. K. U. Grathwohl, M. Horiuchi, N. Ishiguro, and M. Shinagowa (1996) *Archiv. Virol.* **141**, 1863–1874.
4. T. Muramoto, M. Scott, F. E. Cohen, and S. B. Prusiner (1996) *Proc. Natl. Acad. Sci. USA* **93**, 15457–15462.
5. C. H. Chang and D. S. Luse (1997) *J. Biol. Chem.* **272**, 23427–23434.

## Satellite DNA

The regions adjacent to the [centromeres](#) in [chromosomes](#) of higher **eukaryotes** are composed of very long blocks of highly [repetitive DNA](#), in which simple sequences are repeated a thousand times or more. Often, these repeated sequences have a base composition very different from that of the remainder of the DNA, so they can be separated from the bulk DNA by cesium chloride [density gradient centrifugation](#) of slightly fragmented DNA. For this reason, such DNA is referred to as *satellite DNA*. The localization of satellite DNA near the centromeres allows localization of the centromeres by *in situ hybridization* and [autoradiography](#) in squashed **metaphase** cells, using **radiolabeled** RNA copies of the satellite sequence (2).

Different satellites can occur in the same species. For example, **Drosophila virilis** possesses three centromeric satellites (1) whose sequences are highly related. None of these satellite sequences is a pure repeat of the basic consensus sequence. The reason behind the evolution of several sequences is not clear.

Usually, satellite DNA is not transcribed into [messenger RNA](#), as it lacks **promoter** sites where RNA chains can be initiated. **DNA polymerase** must recognize the satellite sequences, however, because the DNA from centromeres is replicated at the same rate as the rest of the chromosomal DNA. It is assumed that the satellite sequences bind particular proteins essential for the centromere's function of holding the daughter [chromatids](#) together during metaphase, such as playing a role in the assembly of the [kinetochores](#) responsible for their attachment to the mitotic spindle or guiding the topoisomerase that eventually separates the sister chromatids (see [DNA Topology](#)). Thus satellite sequences are examples of structural rather than genetic regions of DNA.

## Bibliography

1. J. G. Gall and D. D. Atherton (1974) *J. Mol. Biol.* **85**, 633–634.
2. M. L. Pardue and J. G. Gall (1970) *Science* **168**, 1356–1358.

## Satellites

Satellites are subviral agents that depend for their replication on co-infection of a host cell with a [helper virus](#) (HV). The **nucleotide sequences** of the nucleic acids of satellites are substantially different from those of their HV or their hosts, although some satellites may share with their HV short sequences, often at the termini. Thus, satellites are different from defective interfering particles or RNAs, which are wholly derived from the genome of the HV. One chimeric replicon with sequences derived from both a satellite and a HV has been described (1). The relationship between helper virus and satellite is usually specific in that only one HV species, or a group of related species, can sustain the replication of a particular satellite species; on the other hand, a particular HV may sustain the replication of different, unrelated satellites. There is no correlation between the taxonomic relationships of satellites and those of their HVs. Thus, satellitism may have evolved independently several times. The presence of the satellite may modulate the level of accumulation and the pathogenicity of the HV. Because of these properties, satellites can be considered as molecular parasites of their HVs. Satellites are not needed for the normal multiplication of HVs, but there are satellite-like RNA molecules that may be necessary for completion of the HV's life cycle under natural conditions. Examples are the satellite-like single-stranded RNAs associated with groundnut rosette umbravirus or with beet necrotic yellow vein benyvirus, which are needed for transmission of the HV by their respective aphid or fungal vector.

Satellites are a heterogeneous group, rather than a single taxonomic unit. There are satellites that encode [proteins](#) that are expressed *in vivo*, and others that are noncoding. When the encoded protein is a structural protein that encapsidates the satellite nucleic acid, the satellite is called a *satellite virus*. Satellite viruses are found as distinct nucleoprotein components in preparations of the HV. The first satellite virus was described in 1962 associated with tobacco necrosis virus (TNV) (2). Satellites that do not encode their structural proteins are called *satellite nucleic acids*. The first satellite nucleic acid was described associated with tobacco ringspot nepovirus in 1969 (3). The different groups of satellite viruses and nucleic acids, and their abbreviations, are listed in Table 1.

**Table 1. Satellites and Their Types, with Indication of Size of Nucleic Acid and Accession Number of Typical Sequence Variants<sup>a</sup>**

| Satellite   | Size (nt) | Accession |
|---|-----------|-----------|
| <b>Viruses</b>  |           |           |
| Subgroup I  |           |           |
| Tobacco necrosis satellite virus (STNV)                 | 1239      | J02399    |
| Panicum mosaic satellite virus (SPMV)                   | 826       | M17182    |
| St. Augustine decline satellite virus (SSADV)           | 824       | L10083    |
| Tobacco mosaic satellite virus (STMV)                   | 1058      | M24782    |
| Maize white line mosaic satellite virus (SMWLMV)        | 1168      | M55012    |
| Subgroup II   |           |           |
| Chronic bee-paralysis associated satellite virus (CPVA) | ~ 1100    | NR        |
| <b>Nucleic acids</b>                                    |           |           |
| ssRNA satellites  |           |           |
| <i>Large satRNAs with messenger properties</i>          |           |           |

|  |                  |          |
|--|------------------|----------|
| Arabidopsis mosaic nepovirus (ArMV) large sat-RNA        | 1104             | D00664   |
| Chicory yellow mottle nepovirus (CYMV) large sat-RNA     | 1165             | D00686   |
| Grapevine bulgarian latent nepovirus (GBLV) sat-RNA      | ~ 1500           | NR       |
| Grapevine fanleaf nepovirus (GFLV) sat-RNA               | 1114             | D00442   |
| Myrobalan latent ringspot nepovirus (MLRV) sat-RNA       | ~ 1400           | NR       |
| Strawberry latent ringspot nepovirus (SLRV) sat-RNA      | 1118             | X69826   |
| Tomato black ring nepovirus (TBRV) sat-RNA               | 1374             | X05689   |
| Bamboo mosaic potyvirus (BaMV) sat-RNA                   | 836              | L22762   |
| <i>Small linear satellite RNAs</i>                       |                  |          |
| Turnip crinkle carmovirus (TCV) sat-RNA                  | 230              | X12749   |
|  | 355 <sup>b</sup> | X12750   |
| Cucumber mosaic cucumovirus (CMV) sat-RNA                | 335              | M18872   |
| Peanut stunt cucumovirus (PSV) PSV sat-RNA               | 393              | Z98198   |
| Pea enation mosaic enamovirus (PEMV) sat-RNA             | 717              | U03564   |
| Tobacco necrosis necrovirus (TNV) small sat-RNA          | 620              | NR       |
| Chicory yellow mottle nepovirus (CYMV) small sat-RNA     | 457              | D00721   |
| Artichoke mottled crinkle tombusvirus (AMCV) sat-RNA     | ~700             | NR       |
| Cymbidium ringspot tombusvirus (CyRSV) sat-RNA           | 619              | D0072    |
| Tomato bushy stunt tombusvirus (TBSV) sat-RNA            | 822              | Af022788 |
| <i>Small circular satellite RNAs</i>                     |                  |          |
| Barley yellow dwarf luteovirus (BYDV) sat-RNA            | 322              | M63666   |
| Arabidopsis mosaic nepovirus (ArMV) small sat-RNA        | 300              | M21212   |
| Tobacco ringspot nepovirus (TRSV) sat-RNA                | 359              | M14879   |
| Lucerne transient streak sobemovirus (LTSV) sat-RNA      | 324              | X01984   |
| Rice yellow mottle sobemovirus (RYMV) sat-RNA            | 212              | AF039909 |
| Solanum nodiflorum mottle sobemovirus (SNMV) sat-RNA     | 377              | J02386   |
| Subterranean clover mottle sobemovirus (SCMV) sat-RNA    | 332              | M33000   |
| Velvet tobacco mottle sobemovirus (VTMoV) sat-RNA        | 366              | J02439   |
| dsRNA satellites   |                  |          |
| M satellite of <i>Saccharomyces cerevisiae</i> virus L-A | 1801             | U78817   |
| Satellite of <i>Trichomonas vaginalis</i> virus T1       | 497              | U15991   |
| ssDNA satellites   |                  |          |
| Tomato leaf curl geminivirus (TLCV) sat-DNA              | 682              | U74627   |

<sup>a</sup> nt, nucleotides; ss, single-stranded; ds, double-stranded; NR, not reported.

<sup>b</sup> A chimeric replicon with sequences derived from both satellite and helper virus.

## 1. Satellite viruses

All known satellite viruses have a single-stranded RNA of 800 to 1200 nucleotides. In addition to the open reading frame (ORF) encoding their coat protein, some satellite viruses contain further ORFs; it remains unclear if the encoded products have any role *in vivo* (4). The structure of their isometric,

17-nm diameter particles, built of 60 protein subunits, differs from that of their HV. The particle structures of STNV, STMV and SPMV (see Table 1) have been determined at high resolution by [X-ray crystallography](#) (5). In spite of the different structures of the TNV virus particle and of the coat protein of its satellite STNV, both are able to bind specifically to zoospores of the vector fungus *Oplidium brassicae*.

Best-characterized are the satellites of the tobacco necrosis satellite virus (STNV) subgroup. The 1000- to 1200-nucleotide RNA of STNV has no methylated cap structure or genome-linked protein (Vpg protein) at its 5'-end. Unlike most plant virus RNAs, it has a phosphorylated 5'-terminus. Interference with the accumulation of the TNV has been described for STNV. The particles of STNV may also contain a noncoding satellite RNA of about 620 nucleotides, which depends on TNV for its replication and on STNV for its encapsidation. This is a good example of the complexity of the dependence relationships in satellitism.

Satellite viruses are also found associated with chronic bee-paralysis virus (CPV) (6). Its RNA consists of three species, about 1.1 kb, which can be encapsidated in 17-nm isometric particles built of coat proteins encoded by the satellite or by CPV. The satellite interferes with CPV replication.

## 2. Satellite nucleic acids

Satellite nucleic acids are single-stranded RNA, single-stranded DNA, or double-stranded RNA satellites associated to viruses of the same type of nucleic acid. Most of the characterized satellites are associated with plant viruses, and most are single-stranded RNA.

### 2.1. Single-Stranded RNA Satellites

These are classified into three subgroups.

(i) Subgroup I includes *satellites with messenger properties*. The RNA molecules are between 0.8 and 1.5 kb in size. These satellites encode nonstructural proteins that are expressed *in vivo*. The most studied satellites of this subgroup are those associated with nepoviruses. Large single-strand RNA satellites of nepoviruses have a 3'-terminal [Poly A](#) sequence and a 5'-terminal Vpg protein that, in the analyzed instances, is indistinguishable from that on the HV genome and, therefore, is encoded by it. The encoded nonstructural proteins of different satellite RNAs share some conserved domains, and in some cases have been shown to be needed for replication of the satellite RNA. It may be speculated that these proteins are needed to adapt the virus replication complex to the satellite RNA. For most large satellite RNAs of nepoviruses, an effect on the accumulation or pathogenicity of the HV has not been shown. This may depend on the experimental system: the large satellite RNA of ArMV was shown to modulate the symptoms of the HV, depending on the species of host plant. Another satellite RNA in this subgroup has been found associated with bamboo mosaic potexvirus. The encoded nonstructural protein is not essential for satellite replication (7). Interestingly, this protein shares significant sequence similarity with the structural protein of the satellite virus of PMV, suggesting an evolutionary relationship between the satellite virus and the satellite RNA.

(ii) The second subgroup of single-stranded RNA satellites are *small, linear RNAs*, of less than 800 nucleotides and with no circular forms. Although they may contain potential ORFs, and *in vitro* [translation](#) products have been described for some variants of the satellite RNA of CMV, evidence is strong that these molecules are noncoding RNAs. Thus, their biological functions must depend on the RNA structure, which would mediate the direct interaction of the satellite RNA with components of the HV and/or the host plant. Experimental analyses of secondary structure have been done for the satellite RNAs of CMV and TCV, and there is evidence that several have highly-paired secondary structures. This may account for the high survivability of satellite RNAs *in vivo* and *in vitro*, as well as for their high infectivity. Most satellite RNAs in this subgroup modify the symptoms induced by the HV. The most frequent modification is symptom attenuation, but modulation to more severe symptoms (eg intense chlorosis or necrosis) also occurs. Symptom modulation depends on the triple interaction between the strain of satellite RNA, strain of HV, and species or genotype of the host

plant. For the satellite RNA of CMV, there has been an extensive analysis of symptom determinants in the satellite RNA. For the satellite RNAs of CMV and TCV, the HV and host determinants have been identified as well. The presence of the satellite RNA often results in depression of the accumulation of the HV, but the relationship between this and the attenuation of symptoms is unclear: depression of HV accumulation also occurs when the interaction between satellite RNA and HV results in symptom exacerbation. Also, for some combinations of satellite RNA and HV, symptom attenuation is not accompanied by a reduction in the accumulation of the HV. The satellite RNA may also interfere with movement of the HV within the infected plant, as shown for the satellite RNA of TCV (8).

(iii) The third subgroup of single-stranded RNA satellites are *small circular RNAs*. These satellite RNAs actually occur as both circular and linear molecules. The HV may encapsidate circular forms (for sobemoviruses) or linear forms (for nepoviruses and luteoviruses). In both cases, replication is through a rolling circle mechanism (see [Rolling Circle DNA Replication](#)). Small circular single-stranded RNA satellites interfere with the accumulation and symptoms induced by the HV. These satellites have been described under the heading [Virusoids](#).

The single-stranded RNA satellites may replicate by different mechanisms than their HV. Thus, the replication machinery of the HV must be adapted to the replication of the satellite RNA. Satellite-encoded factors (for the large, encoding satellites) or unidentified host factors may play a role in modifying the virus replication complex. It is to be noted that for those systems that have been analyzed in detail, the efficiency of satellite RNA replication depends on the HV, as well as on the host plant.

Because of the attenuation of the symptoms of the HV, single-stranded RNA satellites have been proposed as efficient agents for the biocontrol of HV-induced plant diseases, both in classical cross-protection programs and by engineering transgenic plants. The high genetic variability of satellite RNAs under experimental and natural conditions, and the possibility of their evolution from attenuating to highly pathogenic types, should be considered when evaluating the potential risks of such control programs.

## 2.2. Double-Stranded RNA Satellites

Other nucleic acid satellites are the double-stranded RNA satellites associated with viruses in the family *Totiviridae*. Of these, the best studied are the M satellites, which are responsible for the killer phenotype and for immunity to the killer toxin associated with the L-A virus of *Saccharomyces cerevisiae*. The multiplication of the M satellite depends on both the helper virus L-A and on host chromosomal genes.

## 2.3. Single-Stranded DNA Satellites

A 682-nucleotide, circular, single-stranded DNA satellite has been found associated with tomato leaf curl geminivirus (TLCV). The satellite DNA contains no ORF, is encapsidated in TLCV particles, and is totally dependent for its replication on the replication-associated protein of TLCV. It shows no sequence similarity with the HV DNA, except for a motif that is universally conserved in geminiviruses and involved in single-stranded DNA replication, and for the binding site for the replication-associated protein. At odds with most RNA satellites, this satellite DNA can be supported by distantly related geminiviruses (9).

Satellites have been considered interesting models to approach the analysis of fundamental aspects of virology, such as virus replication, encapsidation, pathogenesis, recombination, and variability. Also, the often complex dependence relationships between satellites and HV pose interesting questions about [evolution](#). Because of this, much research has focused on satellites. Still, major questions remain largely unanswered: (i) how satellites are replicated by the replication machinery of the HV, (ii) how they affect the accumulation (ie, replication and/or movement) of the HV, (iii) how they modify symptoms of the HV, and (iv) what is their origin, evolution, and fate in populations of their HV.

## Bibliography

1. A. E. Simon and S. H. Howell (1986) *EMBO J.* **5**, 3423–3428.
2. B. Kassanis (1962) *J. Gen. Microbiol.* **27**, 477–488.
3. I. R. Schneider (1969) *Science* **166**, 1627–1629.
4. G. Routh, J. A. Dodds, L. Fitzmaurice, and T. E. Mirkov (1995) *Virology* **212**, 121–127.
5. N. Ban, S. B. Larson, and A. McPherson (1995) *Virology* **214**, 571–583.
6. H. A. Overton, K. W. Buck, L. Bailey, and B. V. Ball (1982) *J. Gen. Virol.* **63**, 171–179.
7. N. S. Lin, Y. S. Lee, B. Y. Lin, C. W. Lee, and Y. H. Hsu (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3138–3142.
8. Q. Z. Kong, J. L. Wang, and A. E. Simon (1997) *Plant Cell* **9**, 2051–2063.
9. I. B. Dry, L. R. Krake, J. E. Rigden, and M. A. Rezaian (1997) *Proc. Natl. Acad. Sci. USA* **94**, 7088–7093.

## Suggestions for Further Reading

10. C. W. Collmer and S. H. Howell (1992) Role of satellite RNA in the expression of symptoms caused by plant viruses. *Annu. Rev. Phytopathol.* **30**, 419–442.
11. R. I. B. Francki (1985) Plant virus satellites. *Ann. Rev. Microbiol.* **39**, 151–174.
12. C. Fritsch, M. A. Mayo, and O. Hemmer (1993) Properties of satellite RNA of nepoviruses. *Biochimie* **75**, 561–567.
13. F. Garcia-Arenal and P. Palukaitis (1999) Structure and function relationships of satellite RNAs of cucumber mosaic virus. In P. K. Vogt (Ed.), *Satellites and defective viral RNAs*. *Curr. Topics Microbiol. Immunol.* **239**, 37–63.
14. M. J. Roossinck, D. Sleat, and P. Palukaitis (1992) Satellite RNAs of plant viruses: structures and biological effects. *Microbiol. Rev.* **56**, 265–279.
15. R. B. Wickner (1996) Double-stranded RNA viruses of *Saccharomyces cerevisiae*. *Microbiol. Rev.* **60**, 250–265.

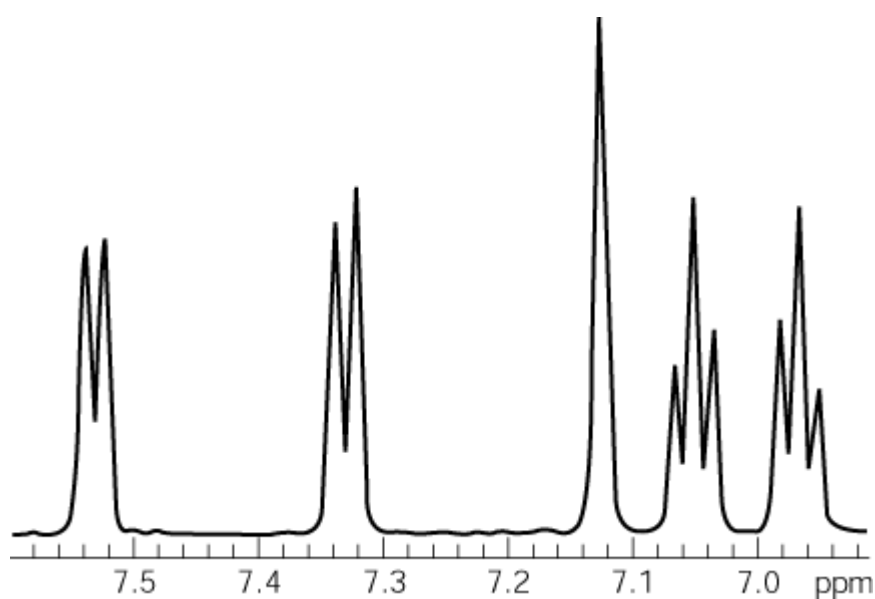
## Scalar Coupling

Nuclear magnetic resonance (NMR) spectra of liquid or gaseous samples of a molecule are usually more complex than would be expected based on the number of chemically distinct nuclei in the molecule. The additional complexity (fine structure) in the spectra arises from spin-spin coupling, also called J-coupling, spin-spin splitting, or scalar coupling. The exact resonance frequencies of NMR transitions associated with a particular nucleus depend on the magnetic field at that nucleus, which is the sum of many components. The largest part is provided by the magnetic field of the NMR instrument, as modified by the chemical shift effects of the electrons of the subject molecule and the electrons of surrounding molecules. (See [Chemical Shift](#).) Smaller contributions to the total magnetic field are made by other nuclei in the vicinity of the subject nucleus. Surrounding nuclei are usually “spin 1/2” nuclei and, according to quantum mechanics, these nuclei can have only a “spin-up” or “spin-down” state. Either of these states will potentially contribute to the magnetic field experienced at a neighboring nucleus, with the result that the NMR signals associated with that nucleus may be a collection of signals (a multiplet). Each line of the multiplet represents a particular combination of neighboring nuclear-spin orientations. Analysis of multiplet structure can provide important details of a molecule, because scalar coupling effects are dependent on the number and

relative orientations of the coupled nuclei. Scalar coupling is a necessary condition for coherence transfer experiments such as correlation spectroscopy (**COSY**) and total correlation spectroscopy (**TOCSY**), and it is fundamental to **isotope-filtering** NMR experiments.

Consider the one-dimensional proton NMR spectrum of *N*-acetyltryptophan shown in Figure 1. The proton (hydrogen atom) attached to C-4 of the indole ring of this molecule has a shielding parameter that corresponds to 7.53 ppm. The NMR spectrum of proton H-4 is more complicated than would be expected because H-4 feels the orientation of adjacent proton H-5. Two possible orientations exist for H-5, so H-4 experiences two different magnetic fields, and these are represented in the proton NMR spectrum as two signals at slightly different frequencies.

**Figure 1.** A portion of the proton NMR spectrum of *N*-acetyltryptophan corresponding to the protons attached to the indole ring, obtained at 500 MHz. (The full spectrum is given in Figure 3 of [NMR](#).) The acetylated amino acid was dissolved in  $d_6$ -DMSO.

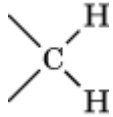
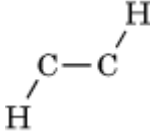
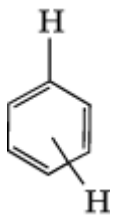
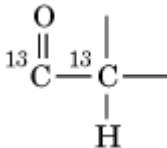
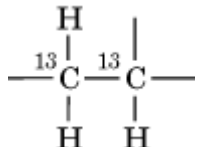
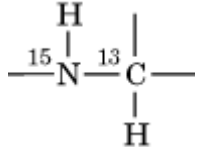


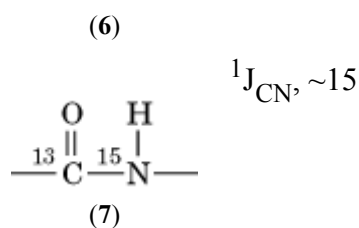
The size of the scalar coupling effect is measured in terms of a parameter called the scalar coupling constant, or spin-spin coupling constant, usually symbolized  $J$  and always measured in Hz. Subscripts may follow the  $J$  that denote the identities of the spins whose interaction is represented by the coupling constant. A superscript may precede the  $J$  that indicates the number of chemical bonds that separate the interacting nuclei. Experiments show that  $J$  can be either positive or negative. Not all NMR experiments are sensitive to the sign of  $J$ , however, and often only the magnitude of  $J$  is discussed. In the case of proton H-4 of *N*-acetyltryptophan, the spin-spin coupling interaction between H-4 and H-5, represented by the symbol  ${}^2J_{\text{H}_\beta\text{H}_\beta'}$ , has a value for the coupling constant of about +7Hz.

The value of a spin-coupling constant depends on several factors, including the gyromagnetic ratios of the interacting nuclei and the nature of the chemical bonding between them. The larger the product of the gyromagnetic ratios of the coupling partners, the larger the  $J_{\text{AB}}$ . A tendency exists for coupling constants to decrease with the number of intervening bonds, that is,  ${}^1J_{\text{AB}} > {}^2J_{\text{AB}} > {}^3J_{\text{AB}}$ .  ${}^3J_{\text{AB}}$  has an important dependence on the [dihedral angle](#) between the first and third covalent bonds, which can be the source of useful stereochemical information. Table 1 indicates the ranges observed for various spin coupling constants between spin 1/2 nuclei that make up biological macromolecules.



Table 1. Typical Values for Coupling Constants, ( $^nJ_{AB}$ )

| Chemical System  | J, Hz  |
|--|--|
|  <p>(1)</p>   | $^2J_{HH}$ , -15 to -10  |
|  <p>(2)</p>   | $^3J_{HH}$ , 2 to 12 (depends on dihedral angle)   |
|  <p>(3)</p>  | $^3J_{HH}$ , 7 to 9 ( <i>ortho</i> ) $^4J_{HH}$ , 1.5 to 2.5 ( <i>meta</i> ) $^5J_{HH}$ , 0 to 1 ( <i>para</i> ) |
| $^{13}C-H$   | $^1J_{CH}$ , 125 to 250 (depends on carbon hybridization)  |
| $^{13}C-C-H$   | $^2J_{CH}$ , -5 to 5   |
|  <p>(4)</p> | $^1J_{CC}$ , ~55   |
|  <p>(5)</p> | $^1J_{CC}$ , 30 to 40  |
| $^{15}N-H$   | $^1J_{NH}$ , 80 to 100   |
|             | $^1J_{CN}$ , ~11   |



A detailed understanding of the spin-spin coupling phenomenon requires a quantum-mechanical treatment. An important result of such an analysis is that J-coupling between nuclei that have the same chemical shift is not observable experimentally in most of the commonly performed NMR experiments with liquids. Thus, the spin coupling constant between two protons of a methyl group, which is a two-bond proton-proton interaction ( ${}^2J_{\text{HH}}$ ), is not observable in a standard one-dimensional proton NMR spectrum.

An observed spin-spin coupling constant of  ${}^nJ_{\text{AB}}$  is potentially made up of two types of contributions. Interactions between the A and B nuclei and the electrons that bind A and B to each other and to the remainder of the molecule provide the means whereby the orientation of one nuclear spin is sensed by the other. There are several possible mechanisms by which nuclei-electron interactions can take place; in all cases, the fine structure introduced by them, as reflected by J values, is independent of the laboratory magnetic field used for the NMR experiment ( $B_0$ ) and, for all practical purposes, independent of the orientation of the molecule being studied in this magnetic field. A direct through-space interaction between spin 1/2 nuclei is also possible. This nuclear dipole–nuclear dipole interaction is strongly dependent on internuclear distance and on orientation of the molecule relative to  $B_0$ . Importantly, in samples where sufficient molecular motion exists so that all possible orientations relative to the magnetic field are rapidly sampled, this contribution to spin-spin coupling averages to zero, leaving only the electron-mediated contributions. These are the conditions that normally obtain in nonviscous solutions at magnetic fields less than ~12 T and normal laboratory temperatures. In organized systems (such as liquid crystals) or in very high magnetic fields, sufficient ordering may exist so that the dipole-dipole contribution to J is not completely averaged away, and an appreciable structure-dependent contribution to the observed coupling constant through this mechanism will be present.

Often, the multiplet structure produced in an NMR spectrum by J-coupling is undesirable because it provides no useful information. A variety of experimental techniques, generally known as decoupling, can be used to remove the effects of spin-spin coupling from an NMR spectrum. Under decoupling conditions, the spectrum is simplified, perhaps reduced to single peaks at the positions corresponding to the chemical shifts (shielding parameters) of the nuclei being examined. Bear in mind, however, that decoupling methods do not remove spin-spin couplings. These are fundamental aspects of the electronic and geometrical properties of the sample that cannot be removed by external intervention. Only the effects of spin coupling—multiplicity in NMR spectra—are removed by a decoupling procedure. The  ${}^{13}\text{C}$  spectrum shown in Figure 4 of [NMR](#) was obtained with proton decoupling. Because proton- ${}^{13}\text{C}$  spin coupling effects are absent from the spectrum, the spectral line positions depend only on the  ${}^{13}\text{C}$  chemical shifts.

#### Suggestions for Further Reading

R. J. Abraham, J. Fisher and P. Loftus (1988) *Introduction to NMR Spectroscopy*, Wiley, New York.

- F. A. Bovey (1988) *Nuclear Magnetic Resonance Spectroscopy*, Academic, San Diego.
- R. K. Harris (1983) *Nuclear Magnetic Resonance Spectroscopy: a physicochemical view*, Pitman, Marshfield, Mass.
- S. W. Homans (1992) *A Dictionary of Concepts in NMR*, Clarendon, Oxford.
- R. Kitamaru (1990) *Nuclear Magnetic Resonance: principles and theory*, Elsevier, New York.
- R. S. Macomber (1998) *A Complete Introduction to Modern NMR Spectroscopy*, Wiley, New York.
- C. H. Yoder and C. D. Schaeffer, Jr. (1987) *Introduction to Multinuclear NMR*, Benjamin/Cummings, Menlo Park.

## Scanning Electron Microscopy

The scanning electron microscope (SEM) has been used for three decades to produce direct, 3-D topographic images of biological samples. Until recently, the resolution in SEM has been too low to directly observe most [macromolecules](#). Advances in specimen preparation, which are sufficient to preserve fine structural details at the molecular level, and instrumentation have extended the possible spatial resolution from 5–10 nm to 2–4 nm. In the SEM, the image is constructed by using a fine electron beam that is raster scanned over the surface of a specimen. Secondary electrons from the specimen are emitted through excitation caused by the primary electrons in the beam striking the specimen. The intensity of the secondary electrons is detected and used to modulate the brightness of another electron beam synchronously scanned over the surface of a cathode ray tube (1). Even though its resolution does not approach that routinely attained in the TEM, the high-resolution SEM (HRSEM) has the advantage of high contrast and high signal-to-noise. However, it must be noted that the SEM provides information only about the surface, whereas the TEM will provide information about the internal structure as well. To increase HRSEM signal yield, enhance contrast, and reduce charging artifacts, coating the specimen with a thin metal layer is performed. However, it is possible to view the fractured surface of a frozen hydrated specimen directly in the SEM without coating the surface with metal (2). For high-resolution work, the metal coat must be thin, with a grain size small enough to reduce decoration artifacts while still producing a high-resolution secondary electron signal.

The preparation of biological specimens can be broadly divided into two methods, chemical and cryo. Though easier to perform, chemical fixation is slower than cryofixation and consideration of pH and osmotic strength must be made. Furthermore, the more widespread adoption of high-pressure freezing to inhibit ice crystal formation has made cryofixation more attractive. One advance that has permitted cryo-HRSEM to characterize biological specimens at near molecular resolution is the employment of an in-lens field emission gun as the source of the electron beam (3). The beam diameter from such a gun is less than 1 nm. However, practical resolution has not achieved this dimension because of the nature of the metal coat and beam damage during observation in the HRSEM. Furthermore, the resolution is limited by the penetration of the electron beam into the specimen, since secondary electrons are emitted from a volume around the site at which the probe meets the specimen. HRSEM can be performed at either low voltage with bulk specimens, or high voltage with thin specimens. The low-voltage electron beam penetrates less deeply into bulk specimens, resulting in the high-contrast signal being generated at only a few nanometers radius from the point of beam impact. Hence, high topographic spatial resolution can be achieved (1). On the other hand, the beam diameter can be made smaller at high voltage, therefore representing a greater potential for high-resolution scanning of thin specimens (4).

HRSEM is being used more often to probe finer details of biological specimens than in the past. Cryo-HRSEM has been used to visualize surface features on [reovirus](#) particles (5). It was found that because SEM imaged only the surface of each particle, obscuring of structural detail by the internal structure or by the opposite surface was avoided, an important consideration with large structures such as viruses. Further results with cryo-HRSEM provided evidence that topographical views of [actin](#) filaments, [microtubules](#), and [clathrin](#) cages could show individual subunits, thus permitting identification by molecular morphology (6). The success in visualizing the 3-D surface architecture of these cytoskeletal elements was attributed to controlled freeze-drying, thin metal coating, cryo-transfer and cryo-observation. In order for SEM visualization to be accomplished, the ice surrounding frozen samples must be removed, either by freeze-drying (ice sublimation in a vacuum) or freeze-substitution (ice dissolved by a chemical solvent), followed by critical point drying. It was found for cytoskeletal structures that freeze-drying provided a superior topography. HRSEM of [mitochondria](#) showed a greater variation of [cristae](#) morphology than previously recognized and also elucidated how serial-section TEM produced unreliable results because of limited resolution (7).

### Bibliography

1. D. C. Joy and J. B. Pawley (1992) *Ultramicroscopy* **47**, 80–100.
2. J. B. Pawley, P. Walther, and S.-J. Shih (1990) *J. Microscopy* **161**, 327–335.
3. T. Nagatani et al. (1987) *Scanning Microscopy* **1**, 901–909.
4. R. Herman and M. Muller (1991) *Scanning Microscopy* **5**, 653–664.
5. V. E. Centronze et al. (1995) *J. Struct. Biol.* **115**, 215–225.
6. Y. Chen et al. (1995) *J. Microscopy* **179**, 67–78.
7. P. J. Lea et al. (1994) *Microsc. Res. Tech.* **27**, 269–277.

### Scanning Hypothesis

Most cellular [messenger RNAs](#) initiate [translation](#) via a scanning mechanism in eukaryotes. According to the scanning hypothesis (1), eukaryotic [ribosomes](#) bind to the 5' end of the mRNA and subsequently migrate along the mRNA until they encounter an AUG codon embedded in an appropriate context, whereupon translation initiation occurs. One of the first steps in initiation of translation of most mRNAs in mammalian cells is recognition of the [-Cap](#) structure and the interaction with it of eIF-4F. At some point in this process, the small ribosomal subunit with associated factors binds to the mRNA and scans to the initiation codon. The scanning model postulates that the 40S ribosomal subunit stops at the first AUG if that codon occurs in a favorable context. But if the first AUG codon occurs in a suboptimal context—for example, in the absence of the critical purine at position –3 or a G at position <+4—some 40S subunits will bypass the first AUG and initiate instead at a downstream site. Two independently initiated proteins may thus be produced from one mRNA by context-dependent leaky scanning. Sometimes ribosomal leaky scanning is used for protein expression. For example, translation of the [hepatitis B virus](#) *pol* gene from the viral pre-genome RNA involves such a mechanism (2).

In addition, internal [start codons](#) in eukaryotes can be reached by a reinitiation mechanism, originally thought to be the exclusive trait of eubacterial ribosomes. Reinitiation at a downstream codon may be possible when the 5-proximal AUG triplet is followed shortly by a terminator codon and fails to form the proper [initiation complex](#) (3).

A similar mechanism of ribosomal scanning is used for translational reinitiation in eubacteria (4). According to this scanning model, the terminated but not released ribosome reaches neighboring initiation codons by lateral diffusion along the mRNA. This ribosomal scanning-like movement is bidirectional, has a radius of action of more than 40 nucleotides in the model system used, and activates the first encountered restart site. The ribosomal reach in the upstream direction is less than in the downstream one, probably due to dislodging by elongating ribosomes. The proposed model has parallels with the scanning mechanism postulated for eukaryotic translational initiation and reinitiation. There are no data to indicate that the two kinds of scanning are basically different. A quantitative difference between prokaryotic and eukaryotic ribosomal scannings is the distance that can be screened. In eukaryotes, this is far greater, in both the 5'-Cap-directed and in termination-dependent states.

### Bibliography

1. M. Kozak (1983) *Microbiol. Rev.* **47**, 1–45.
2. C. G. Lin and S. J. Lo (1992) *Virology* **188**, 342–352.
3. A. G. Hinnebusch (1996) in *Translational Control* (J. W. B. Hershey, M. B. Mathews and N. Sonenberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 199–244.
4. M. R. Adhin and J. van Duin (1990) *J. Mol. Biol.* **213**, 811–818.

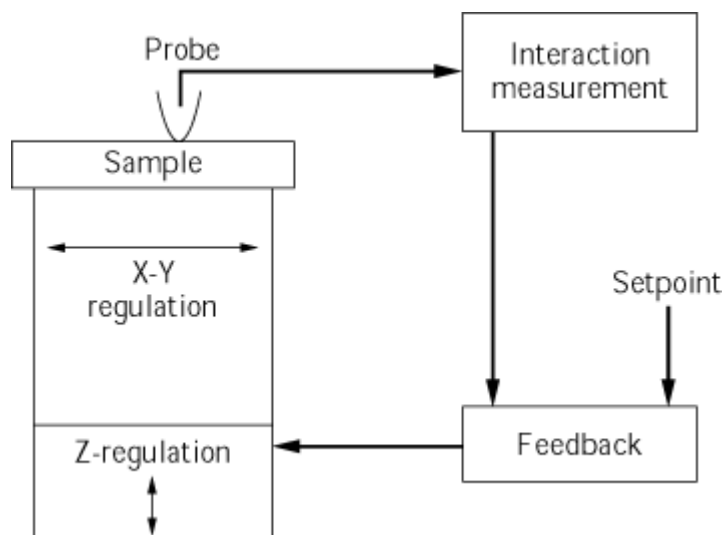
### Suggestion for Further Reading

5. M. Kozak (1995) *Proc. Natl. Acad. Sci. USA* **92**, 2662–2666.

## Scanning Probe Techniques

Scanning probe microscopes provide a unique window to the microworld of cells, subcellular structures, and biological [macromolecules](#). A common feature of all scanning probe microscopes is that the image of a specimen surface is obtained by using a sharp probe scanning in the near field. The specimen surface is scanned while maintaining a very close spacing between the sample and the probe (Fig. 1). Such scanning produces a very short depth of focus on the molecular scale. As a result, the physical features and physicochemical properties of only the top layers of rigidly bound molecules are imaged. The lateral (x and y) resolution in an scanning probe microscope image of inorganic molecules is typically 20 Å and, with the best instruments and the right sample, can be better than a few angstroms. The height (z) resolution is typically less than 1 Å. The scanning probe instruments described here include [scanning tunneling microscopy](#), [atomic force microscopy](#) (AFM), and [magnetic force microscopy](#). Of these, AFM has the greatest potential for biological studies because it can image the three-dimensional (3-D) surface structures of a wide array of native biological specimens in a liquid environment and thus provides an avenue for observing biochemical and physiological processes at molecular resolution in real time. In addition to imaging molecular structures, AFM is also an exquisite sensor of molecular forces, such as [hydrogen bonds](#) and [van der Waals](#) and [electrostatic interactions](#). In addition, by applying a controlled imaging force, surfaces can be altered (nanodissection), created (nanomachining), or moved (nanotransplant) in real time. In its very short period of existence, AFM has imaged biological structures ranging in size from whole cells to macromolecular complexes, such as chromosomal **DNA**, [proteins](#) in extended 2-D arrays, individual multisubunit proteins, various other biological macromolecules, and a variety of [membranes](#). Dynamic processes, such as polymerization or dissolution of macromolecular complexes and crystal growth, and physicochemical properties, such as elasticity, viscosity, and various chemical forces in biological specimens, have been studied.

**Figure 1.** Schematic illustration of the basic operating principles of scanning probe microscopes, in particular the atomic force microscope (AFM) and scanning tunneling microscope (STM). The interaction between a probe tip and the specimen surface is measured continuously while raster scanning in the X and Y directions. With a feedback system, the height (Z position) of the sample is adjusted to keep the tip-specimen interaction force constant. For STM, the interaction force is measured by measuring the tunneling current between the tip and the specimen surface. In AFM, the extent of interaction is measured in terms of the forces of interaction between the specimen surface molecule and the probe, which is usually a tip attached to a cantilevered spring.



Suggestion for further reading

H. E. Wickramasinghe (1992) "Scanned probes old and new", In *Scanned Probe Microscopy* (H. K. Wickramasinghe, ed.), American Institute of Physics, New York, pp. 9–22.

## Scanning Transmission Electron Microscopy

The scanning transmission electron microscope (STEM) combines the imaging mode of the [scanning electron microscope](#) with the specimen characteristics of the conventional transmission electron microscope (CTEM). In the STEM, the electron source, usually a small field emission tip, is demagnified by electron lenses to produce a very small electron probe which is scanned across the specimen in a raster pattern. Unlike the SEM, the STEM is used with thin specimens, and electrons that pass through the specimen are detected and either displayed on a cathode ray tube whose scan is synchronized with the scan of the electron probe, or the signal is digitized and stored in computer memory. CTEMs can be adapted to become STEMs, but the resolution of these hybrid microscopes is not as good as either a dedicated CTEM or a dedicated STEM. A dedicated STEM can produce a probe of 0.2–0.3 nm and hence has resolution roughly equivalent to a CTEM.

A STEM can be used to image the same types of specimens as CTEMs, [macromolecules](#), viruses, thin sections, and so on, but have little advantage over the CTEM at the cost of increased complexity, purchase price, and maintenance of dedicated STEMs. The STEM, however, has an

established niche in molecular biological applications of [microscopy](#) through accurate and precise mass measurements of individual macromolecules and complexes (1, 2). The principle of mass determination relies on quantifying a specimen's differential capacity to scatter electrons elastically. Quantitative image processing is performed on STEM dark-field images of unstained specimens. Unstained particles produce a dark-field signal that is proportional to their local mass density which means that not only the total mass of an individual molecule can be determined, but its resolved domains can also be determined by integrating over appropriate regions of the image (3).

The STEM attains relatively high spatial resolution of different mass elements by the use of a field emission source of electrons. Electrons that pass through the thin specimens used in STEM are collected by a detector placed below the specimen. Since elastically scattered electrons used for mass measurements are scattered over a wide angular range, they can be selectively detected with an annular detector placed below the specimen. Unscattered and inelastically scattered electrons pass through a hole in the center of the annular detector and can be recorded separately by a second circular detector placed below. STEM overcomes limitations in mass analyses, such as impurities, low amount of material, large size of the macromolecule or complex, and conformational variability and polymorphism found in more traditional methods (eg, [chromatography](#), sedimentation, small-angle X-ray, or **neutron scattering**). Because of the STEM's high spatial resolution for mass measurements, discrete subpopulations can be distinguished within a mixed population and treated separately for statistical measurements. STEM also has the advantage of invariant mass measurements regardless of conformation or viewing geometry. Incisive mass specification can be made with as few as 100 molecules with an accuracy of 1–2% (4).

In addition to determining the total mass of a molecule, it is often feasible to measure the mass distribution within it; this is referred to as *mass-mapping*. Distinct features resolvable in images can be classified and computer averaged, thus increasing the signal for more precise mass-mapping with extended resolution. For example, molecular edges can be localized to within 0.5 nm and the radial distribution of density within complexes can be accurately determined (5). The versatility of mass-mapping can be appreciated from the determination of masses ranging from single protein subunits to large viruses (2). An instructive example of the potential of STEM mass analysis was made with isolated gap junction membranes (see Gap Junction and Membranes) to determine the protein composition of the hexameric unit (6).

### Bibliography

1. J. S. Wall and J. Hainfeld (1986) *Ann. Rev. Biophys. Chem.* **15**, 355–376.
2. D. Thomas et al. (1994) *Biol. Cell* **80**, 181–192.
3. D. Walzthony et al. (1984) *EMBO J.* **3**, 2621–2626.
4. M. S. Hamilton et al. (1989) *J. Ultrastr. Mol. Struct. Res.* **102**, 221–228.
5. A. C. Steven et al. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6363–6367.
6. G. Sosinsky (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9210–9214.

## Scanning Tunneling Microscopy

### 1. Principle of Operation

The scanning tunneling microscope (STM), invented by the Nobel prize winners Binnig and Rohrer, is a [scanning probe technique](#) that provides information at atomic resolution about structures that are electrically conducting. STM is based on the quantum mechanical principle that electrons travel

across the classically forbidden energy potential zone (i.e., tunnel) from one surface to another without physical contact if there is overlap in their probabilities of finding their positions, as determined by their electron density functions at the two surfaces. In such a situation, applying a small potential difference between the corresponding surfaces causes electrons to tunnel through the potential barrier that exists between the surfaces of the sample and the probe, which gives rise to the tunneling current. The tunneling current density  $i$  through a planar potential barrier at small voltages  $V$  is inversely proportional to the exponential power of distance and can be approximated by

$$i = Ax \exp[-(\phi)^{1/2} \times s] \quad (1)$$

In this expression,  $f$  is proportional to the average potential barrier height of the two surfaces, and  $s$  is the tunnel distance in angstroms. The constant  $A$  contains several important features, such as the applied bias voltage  $V$ . For a greater applied bias voltage  $V$ , the effective barrier height becomes a function of  $V$ . For a nonplanar barrier, like the barrier between a conducting surface and an STM tip, the tunneling current cannot be described so easily because, in addition in these circumstances, the current is extremely sensitive to the tip-sample separation and is manifested by an exponential dependence on the tunnel distance  $s$ . Generally, the possibility of finding any tunneling current beyond a few angstroms from a weakly conducting surface, such as a biological material, and for a small bias potential, is negligible.

## 2. Modes of operation

The most commonly used STM operating mode is “constant current topography” imaging. In this mode of imaging, the tunneling current is kept at a constant value by feedback electronics, and the change in the scanner height necessary to maintain constant tunneling current is registered as topographic information (see Fig. 1 of [Scanning Probe Techniques](#)). Such height change is registered at each lateral x-y coordinate as the tip scans over the surface. This mode is most suitable to measuring the topography of nonflat surfaces. Because of the finite reaction time of the feedback, however, the scan speed is limited. Another way of operating the STM is the “constant height mode.” For this mode, the height of the probe is fixed, and the tunneling current is measured at each x-y coordinate during the scanning. This mode allows much faster scanning of atomically flat surfaces because the feedback does not have to respond to surface features. This is important because fast scanning allows dynamic studies. It also minimizes data acquisition time, and image distortion caused by thermal drift and creep in the piezo scanner.

STM images give pure topographic information about surfaces only when the surface has uniform electronic properties because it actually measures the electronic structure of the sample and probe via barrier heights. The resolution is limited primarily by the size, shape, and mechanical stability of the tip and is on the order of 1 Å laterally, and 0.1 Å vertically.

STM can be operated in a more **spectroscopic** mode, measuring current-voltage characteristics. For each lateral position of the probe, the feedback is temporarily disabled after which the tunneling current is measured while ramping the bias voltage. The variation in the tunneling current at specific bias voltages can be visualized by constructing real-space images of surface electronic states. Spatially resolved current-voltage curves have been used to examine the electronic structure of semiconductor and metallic surfaces. The energy resolution of the spectroscopic imaging mode is influenced by several factors. Some of these factors are the thermal limit, which depends on the width of the edge of the Fermi distribution of the electrons in the electrode surfaces; an energy-dependent lifetime broadening of electronic states; and the uncertainty principle, which plays a role because the lateral resolution and tip radius are both of the order of angstroms.

## 3. Mechanical design



The mechanical design of a STM system can be modified for specific experimental needs. The most general criteria that must be satisfied by a STM setup are (1) an x, y, and z scanning range that allows for a scanning window sufficiently large to find and identify features of interest; (2) a resolution of 0.1 Å laterally and 0.01 Å in the z-direction; and (3) rigid construction of the tip-sample junction part with high resonance frequency to minimize vibrational disturbances. The scanners use mostly polycrystalline ceramic materials, like lead zirconate titanate and barium titanate.

The operating environment of STM studies depends on the nature of the contamination. Most inert samples are imaged in air, but for more accurate and reliable studies on clean surfaces, a STM situated in a vacuum chamber may be necessary. The mechanical construction in this case is more complicated because of the need for vibration isolation in the vacuum chamber, for in situ cleaning, and for possible use of Auger electron spectroscopy in the same chamber which is low-energy electron diffraction. Another approach that imposes harsher requirements on a piezo scanner is STM imaging in cryogenic conditions because the low temperature reduces the piezoelectric response of the ceramics.

#### 4. Probe preparation

The physical shape, size, and chemical identity of the STM probe employed influences the resolution obtained in a STM scan and also the electronic structure. Three important properties make a “good” probe for STM. First, a large blunt macrostructure ensures high mechanical resonance frequency, leading to low hysteresis and an increase in the data acquisition rate. Second, the microstructure of the tip at the very apex needs a single site of atomic approach to the sample. If this is not met, multiple tip imaging occurs because electrons tunnel through several points on the tip apex. Third, the purity of the tip is important because any impurities may lead to barriers through which electrons do not tunnel. Most STM tips are made of metal wires, such as tungsten, platinum-iridium, or gold. Platinum-iridium tips are used most because platinum is inert to oxidation, produces a pure tip, and iridium adds stiffness to the probe. They are cut and sharpened in different ways, from simply using a wire cutter to ion milling or electrochemical etching.

#### 5. Applications

The first STM work was done on conducting or semiconducting surfaces. A precise image of the atomic arrangements of those surfaces was obtained, mostly with the STM running in ultrahigh vacuum, primarily to avoid contamination from gases, mainly oxygen, and organic materials present under ambient conditions. Although it is difficult to image thicker biological samples with STM, small molecules like aromatic molecules, fatty acids, hydrocarbons, and liquid crystal molecules are routinely imaged. Energy levels between the various electron orbitals can be precisely calculated in some cases, and the use of a substrate like graphite can shift these orbitals in the adsorbed liquid crystalline molecules. Numerous similar applications of STM are described elsewhere.

Although attempts have been made to obtain structural information about larger biological macromolecular specimens, which are nonconducting, the mechanisms of image generation and their interpretation are unclear. Some of the early imaging of large molecules, such as **DNA**, with STM were later proven to be artifacts, when it was found that features that look very similar to a double-helical DNA strand could be obtained from a clean sample of graphite. More recently, a new way of imaging DNA and possibly other macromolecules with STM was described. The technique imaged DNA and [tobacco mosaic virus](#) on mica surfaces under a relative humidity of 65% and found clear images with a tunneling current in the subpicoampere range. The electrical conductivity is explained by a layer of [water](#) on the surface that has conductivity five times higher than bulk water.

Suggestions for Further Reading

G. Binnig and H. Rohrer (1987) Scanning tunneling microscopy - from birth to adolescence, Rev. Mod. Phys. **59**, 615–625.

P. K. Hansma, V. B. Elings, C. E. Bracker, and O. Marti (1988) Scanning tunneling microscopy and atomic force microscopy - application to biology and technology, Science **242**, 209–216.

## Scatchard Plot

The Scatchard plot is a graphical method of analyzing equilibrium **ligand binding** data (1). It is used to determine the number of ligand-binding sites on a **receptor**, whether these sites show **cooperative** interactions, whether more than one class of site exists, and the respective affinities of each site. The experimental parameters used for a Scatchard plot are the free ligand concentration [L] and the average number of ligand molecules bound to a receptor,  $\bar{n}$ , at a particular ligand concentration at equilibrium. The formal definition of  $\bar{n}$  is

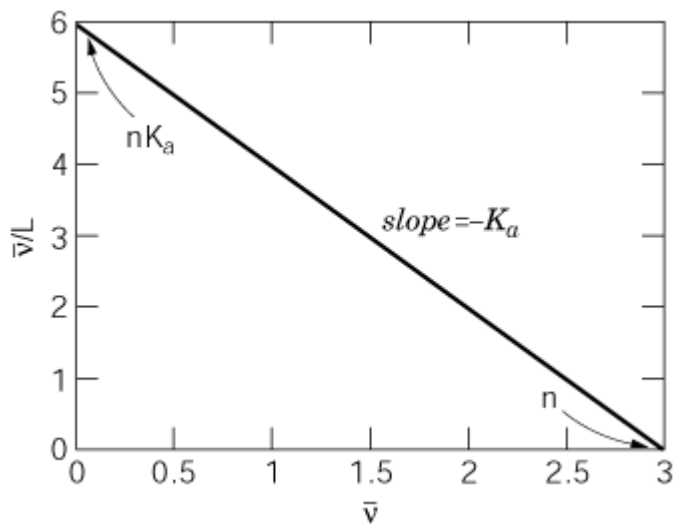
$$\bar{n} = \frac{\text{moles bound ligand}}{\text{moles receptor}} = \frac{\sum_{i=1}^n i[R_i]}{\sum_{i=0}^n [R_i]} \quad (1)$$

where  $n$  is the maximum number of ligands that can bind and  $[R_i]$  is the concentration of the receptor species with  $i$  ligands bound. In the simplest case, of a receptor with a single class of noninteracting sites, receptor–ligand binding follows the Scatchard equation:

$$\frac{\bar{n}}{[L]} = nK_a + \bar{n}K_a \quad (2)$$

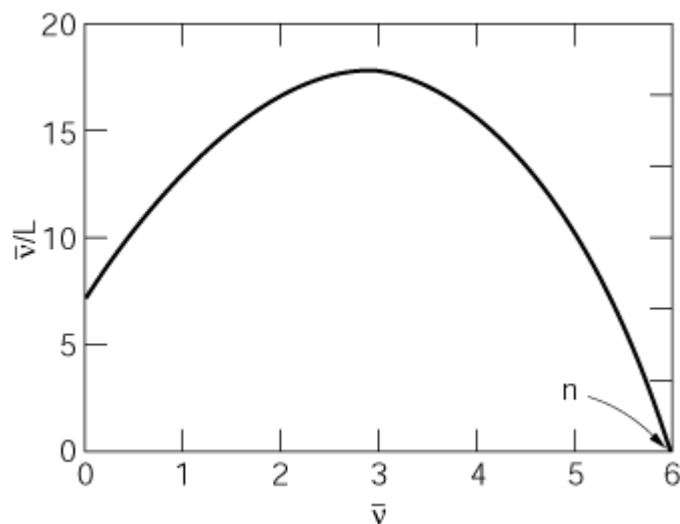
A linear plot results, as shown in Figure 1. The intercept on the abscissa equals the total number of binding sites in the receptor, and the slope equals  $-K_a$ , where  $K_a$  is the intrinsic equilibrium **association constant**.

**Figure 1.** Scatchard plot for a receptor with three equivalent, noninteracting ligand-binding sites.



Nonlinear Scatchard plots are diagnostic of several more complex types of ligand–receptor interaction. A plot that is concave downward indicates positive **homotropic** cooperativity between binding sites on an **allosteric** receptor. An example is shown in Figure 2. The abscissa intercept again equals the maximum number of binding sites. The curvature of the plot corresponds to the degree of cooperativity in the system studied, but it is difficult to obtain an experimentally meaningful parameter from the curvature by inspection. Nevertheless, the linearity of Scatchard plots is very sensitive to cooperative interactions, and the downward-concave form of the plot in Figure 2 is considered diagnostic of allosteric proteins.

**Figure 2.** Scatchard plot for a receptor that exhibits cooperative ligand binding. The example illustrates a molecule that follows the Monod–Wyman–Changeux **concerted model** of allostery, with six binding sites and 10-fold preferential binding to the R state.

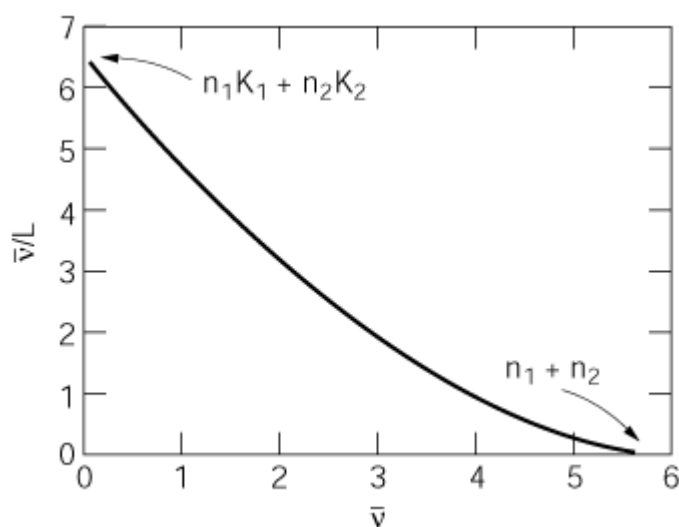


A second type of nonlinear Scatchard plot arises when multiple classes of binding sites with distinct affinities are present on a receptor. For receptors with multiple ( $m$ ) classes of noninteracting sites, the Scatchard equation becomes a sum:

$$\frac{\bar{v}}{[L]} = \sum_{j=1}^m \frac{n_j K_j}{1 + [L]K_j} \quad (3)$$

where  $n_j$  and  $K_j$  are the number of sites and association constant of the  $j$ th class. The corresponding Scatchard plot (Fig. 3) is multiphasic, superficially resembling a linear combination of simpler plots. The abscissa intercept equals the total number of binding sites on the receptor, but further interpretation of multiphasic Scatchard plots is difficult.

**Figure 3.** Scatchard plot for a receptor with two classes of ligand binding sites, three high affinity, and three low affinity sites. The sites do not interact, and the affinities differ by a factor of 10.



Several limitations of Scatchard analysis have been recognized. One is that  $n$  and  $n/[L]$  both contain experimental errors; hence calculating the slope and intercept of a linear Scatchard plot by simple least squares is not statistically valid. A better approach for statistical analysis is to fit  $n$  as a function of  $[L]$  directly to the equation for a binding model (2). In addition, nonlinear Scatchard plots are notoriously difficult to interpret. Scatchard analysis amplifies the prominence of binding data obtained at low ligand concentrations. These data can distort extrapolations to the total number of binding sites that will be filled at a saturating ligand concentration (3). Another source of overinterpretation stems from the fact that binding behavior consistent with multiple classes of sites is also consistent with other binding models (4). In the case of multiphasic plots like that in Figure 3, discerning the true molecular properties of a receptor requires independent experiments, in addition to Scatchard analysis.

#### Bibliography

1. G. Scatchard (1949) *Ann. NY Acad. Sci.* **51**, 660.
2. J. E. Fletcher and A. A. Spector (1968) *Comput. Biomed. Res.* **2**, 164–175.
3. I. M. Klotz (1982) *Science* **217**, 1247–1249.
4. I. M. Klotz and D. L. Hunston (1984) *J. Biol. Chem.* **259**, 10060–10062.

#### Suggestion for Further Reading

5. I. M. Klotz (1985) Ligand-receptor interactions: facts and fantasies, *Quart. Rev. Biophys.* **18**, 227–259. (Comprehensive review detailing pitfalls in formulation of ligand binding models and

analysis of binding data.)

## Scattering Intensity Distribution

The scattering of light, X-rays, or neutrons by atoms in matter provides information on the structure and dynamics of the material being sampled. The scattering experiments depend upon measuring the scattering intensity distribution, generally as a function of scattering angle. Presented here is a brief description of the theory for the elastic, coherent scattering of X-rays and neutrons. Coherence refers to the fact that the amplitudes of the scattered waves are additive and hence can produce interference. Elastic scattering means that there is no energy difference between the incident and scattered waves. When X-rays or neutrons are scattered coherently and elastically by atoms, the scattered waves interfere in a manner related to the spatial distribution of the atoms in the sample. This interference provides the foundation for [small-angle scattering](#) and [crystallography](#) applications. The basic principles of [light scattering](#) are the same as those for X-rays, but the wavelengths are much longer and different formalisms are used. [Light scattering](#) is therefore discussed as a separate topic in this series.

X-rays and neutrons can be considered as plane waves with wavelengths,  $\lambda$ , and can be represented as

$$\Psi(z) = e^{ik_i z} \quad (1)$$

where  $k_i = 2\pi/\lambda$  is the amplitude of the incident wave vector which is taken to be in the direction  $z$ . A wave scattered by an atom at a position in space designated by the vector  $\mathbf{r}$  will be a spherical wave of the form

$$\Psi(\mathbf{r}) = -\frac{A}{r} e^{ik_j \cdot \mathbf{r}} \quad (2)$$

where  $\mathbf{k}_j$  is the scattered wave vector, and  $A$  is the scattering amplitude for the atom. X-rays are electromagnetic radiation and are scattered via interactions with the electrons in a sample. Thus, for X-rays,  $A$  is proportional to the number of electrons in the atom and is generally designated as  $f$ . Because X-ray wavelengths are of the same order as the dimensions of the electron clouds of the scattering atoms, there is an angular dependence to [X-ray scattering](#) amplitudes. Neutrons are neutral particles and are scattered by atomic nuclei. Neutron scattering amplitudes are generally given in units of length and are designated as  $b$ . Because nuclei have dimensions much smaller than the wavelengths of neutrons used in scattering experiments,  $b$  values have no dependence on scattering angle. The total coherent scattering from a molecule made up of atoms is the summation of the coherent scattering over all atom pairs  $(i,j)$ :

$$\sum_i \sum_j A_i A_j e^{-i\mathbf{Q} \cdot \mathbf{r}_{ij}} \quad (3)$$

where  $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$ .  $\mathbf{Q}$  is the scattering vector equal to the difference between the incident and scattered wave vectors ( $\mathbf{k}_f - \mathbf{k}_i$ ); its amplitude is  $4p(\sin q)/\lambda$ , where  $2q$  is the scattering angle (Fig. 1).  $\mathbf{Q}$  is also the momentum transfer vector. In scattering experiments, one can describe molecules as continuous density distributions rather than sums of discrete point atom scatterers. Scattering densities are calculated as the sum of the scattering amplitudes of atoms within a finite volume element divided by its volume,  $\sum A_i/V$ . The total coherent, elastic scattering,  $I(\mathbf{Q})$ , from a molecule in a vacuum then can be written as

$$I(\mathbf{Q}) = \left| \int \rho(\mathbf{r}) e^{-i\mathbf{Q}\cdot\mathbf{r}} d^3\mathbf{r} \right|^2 \quad (4)$$

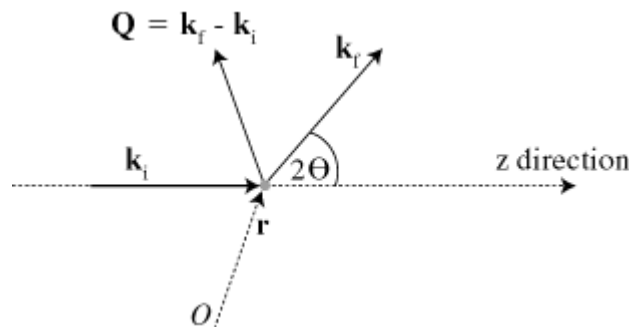
where  $\rho(\mathbf{r})$  is the scattering density distribution within the molecule, and the integration is over the volume of the molecule. For a randomly oriented molecule in a solvent with scattering density  $\rho_s$ , equation (4) becomes

$$I(Q) = \left\langle \left| \int \Delta\rho(\mathbf{r}) e^{-i\mathbf{Q}\cdot\mathbf{r}} d\mathbf{r} \right|^2 \right\rangle \quad (5)$$

$$= \int \int \Delta\rho(\mathbf{r}_1) \Delta\rho(\mathbf{r}_2) \frac{\sin Q|\mathbf{r}_1 - \mathbf{r}_2|}{Q|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (6)$$

The brackets  $\langle \dots \rangle$  indicate averaging over all orientations of the particle, and  $\Delta\rho(\mathbf{r}) = \rho(\mathbf{r}) - \rho_s$  is the scattering density difference or “contrast” between the molecule and the solvent (1, 2). Equation (6) is known as the Debye equation 3 and is the basic equation for scattering. The total scattering intensity distribution for a solution of monodisperse molecules will be directly proportional to both the number density (concentration) and the square of the molecular weight of the particles. Analysis of the scattering intensity distribution at small angles ([small-angle scattering](#)) yields structural parameters such as the [radius of gyration](#),  $R_g$ , molecular weight,  $M$ , and the vector length distribution function,  $P(r)$ , for the scattering particle. [Contrast variation](#) techniques involve the deliberate manipulation of contrast in order to increase the information content of the scattering experiment. For polydisperse solutions, or a solution of nonidentical particles, the time- and ensemble-averaged scattering intensity distribution function is measured; hence average structural parameters are determined.

**Figure 1.** Geometrical representation of the scattering vector  $\mathbf{Q}$ .  $\mathbf{k}_i$  and  $\mathbf{k}_f$  are the incident and scattered wave vectors, respectively, of a wave that interacts with an atom at a point  $\mathbf{r}$  from an arbitrary origin  $O$ .



For scattering from samples ordered in one, two, or three dimensions, the scattering from an individual molecule is convoluted with the repeating lattice structure to yield a diffraction pattern that has discrete intensity maxima. These diffraction patterns can give information on the repeat distances in the lattices, as well as provide higher resolution structural information on the ordered molecules. For well-ordered three-dimensional crystals, complex diffraction patterns are obtained that can be indexed according to the crystal lattice indices  $h, k, l$ . The [unit cell](#) scattering is described mathematically in terms of the square of the [structure factor](#)  $F_{hkl}$ , which is the ratio of the radiation scattering by any real sample to a point scatterer at the origin:

$$I(Q_{hkl}) = F_{hkl}^2 = \left| \int_j A(Q)_j e^{-2\pi i(hx_j + ky_j + lz_j)} \right|^2 \quad (7)$$

where  $x, y$ , and  $z$  are the coordinates of each atom in the crystallized molecules in one unit cell, and the summation is over all atoms. Diffraction data from three-dimensional crystals of biological polymers (proteins, polynucleotides) can be used to solve the structure of the crystallized polymer at high resolution ([crystallography](#)) if one knows both the phases and amplitudes of the structure factors  $F_{hkl}$ . The structure factor amplitudes are readily calculated as the square root of the measured intensities of the diffraction peaks. These amplitudes are combined with experimentally and/or theoretically deduced model phases in order to calculate (by Fourier transformation) an electron density distribution function that is usually interpreted using the known contiguous sequence of chemical groups in the polymer.

#### Bibliography

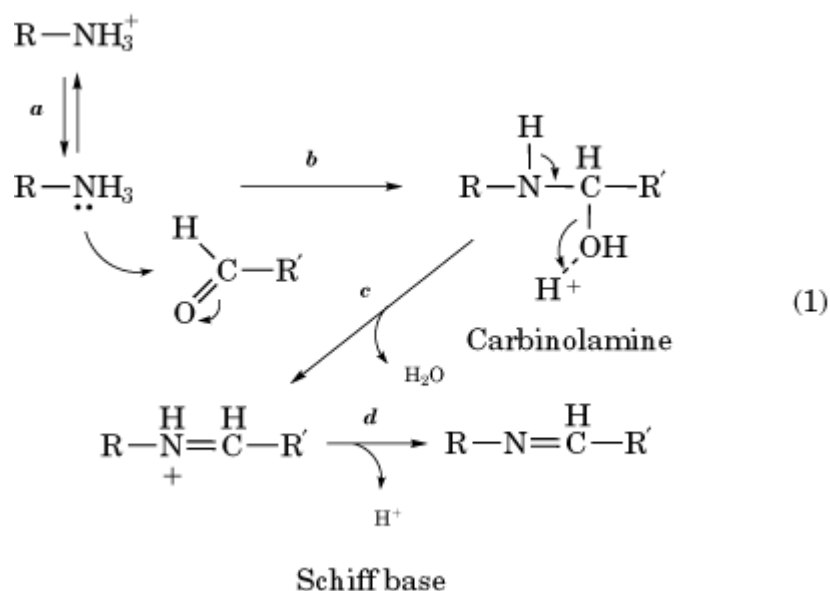
1. W. L. Bragg and M. F. Perutz (1952) Proc. R. Soc. **A213**, 425.
2. K. Ibel and H. B. Sturhmann (1975) J. Mol. Biol. **93**, 255–265.
3. P. Debye (1915) Ann Physik. (Leipzig) **46**, 809–823.

#### Suggestions for Further Reading

4. G. E. Bacon (1975) *Neutron Diffraction*, Clarendon Press, Oxford, U.K.
5. T. L. Blundell and L. N. Johnson (1976) *Molecular Biology. An International Series of Monographs and Textbooks: Protein Crystallography*, Academic Press, London.
6. C. R. Cantor and P. R. Schimmel (1980) In *Biophysical Chemistry, Part II: Techniques for the Study of Biological Structure and Function*, W. H. Freeman, San Francisco, pp. 687–846.
7. L. A. Feigen and D. I. Svergun (1987) *Structure Analysis by Small-Angle X-Ray and Neutron Scattering*, Plenum Press, New York.
8. O. Glatter and O. Kratky (1982) *Small-Angle X-Ray Scattering*, Academic Press, New York.

## Schiff Base

Schiff bases (also known as Schiff's bases, imines, or azomethines) are unsaturated compounds formed by the condensation of [amino groups](#) with aldehydes or ketones (1-4): (Eq. 1)



The amine will be reactive in its un-ionized form and may require deprotonation (step *a*). The free amino group adds to the carbonyl group (step *b*) to form an intermediate known as a carbinolamine. It loses water (step *c*) to form an *N*-protonated Schiff base, which may lose a proton (step *d*). All the reaction steps reach equilibrium, often quite rapidly. Because the  $pK_a$  of the protonated Schiff base is usually low whereas that of a protonated primary amino group is high, the equilibria will be pH-dependent. Consider formation of a Schiff base by reaction of an unprotonated amine with a carbonyl compound to give an unprotonated Schiff base, with formation constant  $K_f$  at high pH. If the  $pK_a$  value of the amine is 10 and that of the Schiff base is 5, it is easy to show that the formation constant for reaction of the protonated amine to form protonated Schiff base at low pH, eg, below pH 3, will be only  $10^{-5}K_f$ . Consequently, the Schiff base will not be formed to a significant extent at low pH. If, however, the  $pK_a$  of the amine is unusually low, or if the Schiff base proton is held by a [hydrogen bond](#) not present in the free amine, a Schiff base may be quite stable at neutral pH. This is the case for Schiff bases of [pyridoxal phosphate](#).

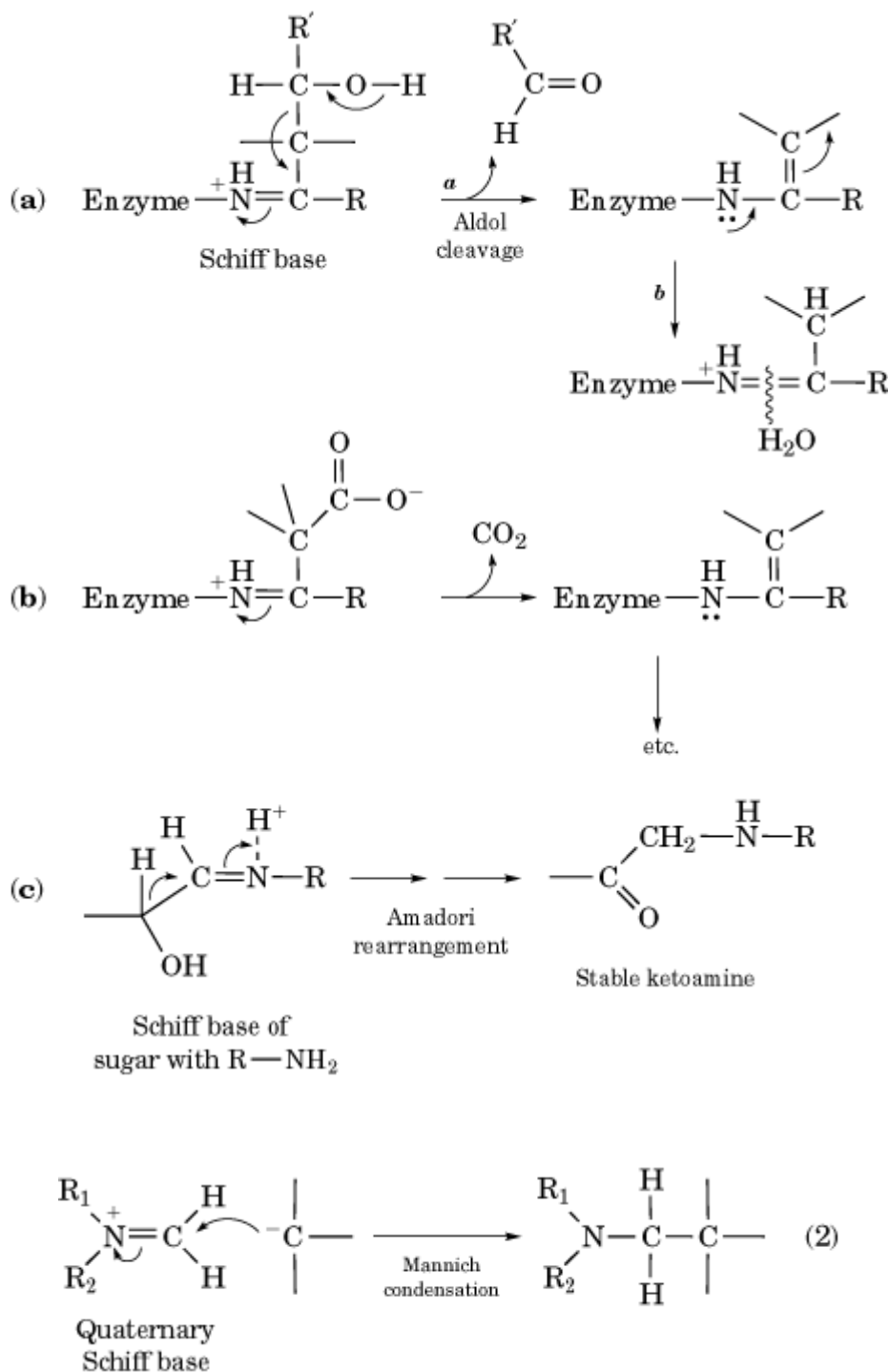
Schiff bases are chemically reactive, and various adducts can be formed. Especially popular is a reaction with a borohydride ion, which reduces the double bond of the Schiff base, converting the labile imine linkage to a stable secondary amine. This reaction has proved useful in linking pyridoxal phosphate to amino groups of [enzymes](#) or substrates, in attaching aldehydes to amino groups of receptors, and so forth (5, 6).

Schiff bases participate in many biological processes. In a common class of aldolases, protonated Schiff bases are formed with ketone substrates such as fructose 1,6-bisphosphate. The protonated imine serves as an electron-accepting group that facilitates carbon-carbon bond cleavage by the enzyme (Fig. 1a). (7, 8) A similar function is found for bacterial acetoacetate decarboxylase (Fig.



**1b** (9). The Amadori rearrangement, prominent in carbohydrate chemistry, including nonenzymatic glycation and the Maillard browning reaction, also proceeds through a Schiff base mechanism (Fig. **1c**). (10, 11) The Mannich condensation useful in organic synthesis (3, 12) is an addition reaction of a quaternary Schiff base, often derived from formaldehyde, with a carbanionic center (**Eq. 2**).

**Figure 1.** Participation of Schiff bases in three biologically important reactions: (a) type I aldolase; (b) beta-oxoacid decarboxylase, eg. acetoacetate decarboxylase; (c) Amadori rearrangement.



## Bibliography

1. V. Malatesta and M. Cocivera (1978) *J. Org. Chem.* **43**, 1737–1742.

2. R. L. Reeves (1966) in *The Chemistry of the Carbonyl Group* (S. Patai, ed.), Interscience, London, pp. 567–619.
3. S. Patai, ed. (1970) *The Chemistry of the Carbon-Nitrogen Double Bond*, Interscience, London.
4. G. Tennant (1979) in *Comprehensive Organic Chemistry* (I. O. Sutherland, ed.), Pergamon, Oxford, U.K., Vol. 2, pp. 385–590.
5. R. Edger and D. B. Rifkin (1977) *Biochim. Biophys. Acta* **470**, 70–83.
6. B. Perez-Ramirez, A. Iriarte, and M. Martinez-Carrion (1994) *J. Protein Chem.* **13**, 67–76.
7. A. J. Morris and D. R. Tolan (1994) *Biochemistry* **33**, 12291–12297.
8. J. Jia, U. Schörken, Y. Lindqvist, G. A. Sprenger, and G. Schneider (1997) *Protein Sci.* **6**, 119–124.
9. L. A. Highbarger, J. A. Gerlt, and G. L. Kenyon (1996) *Biochemistry* **35**, 41–46.
10. N. Mori and J. M. Manning (1986) *Anal. Biochem.* **152**, 396–401.
11. J. W. Baynes and V. M. Monnier, eds. (1989) *The Maillard Reaction in Aging, Diabetes and Nutrition*, Alan Liss, New York.
12. N. Risch and M. Arend (1995) in *Methods of Organic Chemistry*, Vol. E21b (G. Helmchen, R. W. Hoffman, J. Mulzer, and E. Schaumann, eds.), Georg Thieme Verlag, Stuttgart, Germany, pp. 1908–2202.

### Suggestions for Further Reading

13. D. L. Leussing (1986) "Model reactions". In *Vitamin B<sub>6</sub> Pyridoxal Phosphate Chemical, Biochemical and Medical Aspects*, Vol. 1, Part A (D. Dolphin, R. Poulson, and O. Avramovic, eds.), Wiley, New York, pp. 69–115.
14. R. G. Kallen, T. Korpela, A. E. Martell, Y. Matsushima, C. M. Metzler, D. E. Metzler, Y. V. Morozov, I. M. Ralston, F. A. Savin, Y. M. Torchinsky, and H. Ueno (1985) "Chemical and spectroscopic properties of pyridoxal and pyridoxamine phosphates". In *Transaminases* (P. Christen and D. E. Metzler, eds.), Wiley, New York, pp. 37–108.

### Scrapie

Scrapie is the archetype of the transmissible spongiform encephalopathies (TSE), a closely related group of rare and fatal neurodegenerative diseases affecting humans and other mammals. They are of particular scientific interest because the transmissible infectious agent appears to be a single type of [protein](#) molecule, the [prion](#) protein, rather than an organism containing **nucleic acid** ([1](#), [2](#)). The diseases are also unique in that they are both hereditary and infectious. Their importance has been further highlighted in the past 10 years by the recognition that they can be transmitted across the species barrier, from cattle to humans, in contaminated foodstuffs. The characteristic signs of the disease include dementia, progressive loss of locomotion coordination (ataxia), and wasting, leading inevitably to death. Postmortem histological examination of the brain ([3](#), [4](#)) shows a spongiform appearance with [amyloid](#) plaques surrounding the voids. The incubation times of the diseases tend to be long; in mice it is about 1 year, in sheep about 5 years, and in humans it may be 10–50 years, depending on the disease strain. The spongiform encephalopathies can be transmitted between animals by a number of different routes. The most effective, and useful in laboratory studies, is by intracerebral inoculation of tissue homogenate. Other routes include feeding ([5](#)) and the inadvertent intramuscular injection ([6](#)) of infected material during medical treatment and corneal grafting ([7](#)). In the United Kingdom, scrapie may have been transferred from sheep to cattle (as bovine spongiform

encephalopathy (BSE)) on a large scale, and thence to the human population (as new-variant Creutzfeldt–Jakob disease) (8) through infected beef products, so far on a small scale.

The spongiform encephalopathies include the following:

**Scrapie:** a disease of sheep and goats whose altered behavior prior to death includes the obsessive scraping of the fleece, hence the Scottish dialect name of “scrapie.” It has been known for at least 250 years and is endemic in all countries except Australia and New Zealand. Experimental sheep-to-sheep transfer of scrapie by inoculation was first carried out 60 years ago, and it is thought that natural transmission occurs through pasture contaminated with placental tissue carrying the agent, but vertical transmission, ewe-to-lamb, is disputed. There is no evidence of direct transmission from sheep to humans. Scrapie can be bred out of flocks, indicating it has a genetic component. Similar diseases include transmissible mink encephalopathy in mink and chronic wasting disease in muledeer and elk. Spongiform encephalopathies have been reported in the United Kingdom in domestic cats and several types of deer, possibly from feeding with infected sheep or cattle remains.

**Bovine spongiform encephalopathy (BSE), or “mad cow” disease:** a disease of cattle first recognized in 1986 originating in, and largely confined to, the United Kingdom. The United Kingdom has reported 168,000 cases of BSE in the period 1986–1997 in an epidemic that peaked in 1992–1993 and is now diminishing rapidly. Other European countries have reported no more than a few hundred cases, often of cattle imported from the United Kingdom, but no cases of BSE have been found in the United States, which has nearly 100-fold more cattle than the United Kingdom. For commercial or political reasons, BSE cases are likely to be reported conservatively in every country. Epidemiological evidence suggests that BSE originated from the feeding (5) of inadequately rendered sheep remains from scrapie-containing flocks, but a spontaneous origin in cattle herds has not been entirely ruled out. The practice of using meat and bone meal in cattle feeds had been common for several decades, but reductions in the severity of the rendering conditions permitted by the U.K. government in the late 1970s and early 1980s apparently allowed the infectious scrapie agent to escape destruction and to pass into cattle. It is also known that BSE is transmissible from cow-to-calf (vertical transmission) at a very low rate.

## 1. Human Spongiform Encephalopathies

There are five recognized diseases in the human population associated with TSEs.

**Creutzfeldt–Jakob Disease (CJD):** a rare and fatal neurodegenerative disease in humans first described in Germany in the early 1920s independently by H. G. Creutzfeldt and A. Jakob. It normally affects people above 60 years of age and is presented as a rapidly progressive dementia, with a characteristic electroencephalogram signature. CJD may have one of three origins: hereditary transmission as an autosomal, dominant trait; acquisition by ingestion; or acquisition by an iatrogenic route, such as surgery, corneal transplant, or [growth hormone](#) injection. The majority of cases are sporadic (85%), 10–15% are familial, and the remainder are iatrogenic. The incidence of CJD is given as about 1 per million per year, with little geographic variation. It is possible, however, that CJD is rather more common than this because of its misdiagnosis as Alzheimer's disease (9) (see [Amyloid Precursor Protein](#)), with which it shares a remarkable number of etiological and pathological similarities (10).

**New-variant CJD, or v-CJD:** appears to be distinct from normal CJD (8) and thus far limited to the United Kingdom. Cases have appeared from 1996 at a rate of about 10–20 per year. The new variant affects an age group below 40 years, it progresses more slowly than classic CJD, and the brain plaques are very similar to those of BSE in cattle brains. It is now thought to be a human form of BSE resulting from the ingestion of infectious cattle tissue. In view of the major epidemic of BSE in the U.K. cattle herds, and the prolonged incubation period, a parallel increase in cases of v-CJD at a later stage in the human population in the United Kingdom cannot be ruled out.

**Kuru:** most probably a sub-form of CJD, discovered in Papua, New Guinea, but now extinct. It is likely that its transmission was associated with ritual tribal practices of eating the brains of ancestors.

**Gerstmann–Straussler–Scheinker disease:** a very rare inherited midlife disease in humans manifest by ataxia.

**Fatal familial insomnia (FFI):** another very rare human midlife inherited disease in which dementia follows difficulty in sleeping.

It is evident that all these diseases are closely related and may represent essentially the same disease in different hosts. The nature of the transmissible agent has been very controversial. The initial assumption that it was a slow-acting **virus** has been eroded, as no virion could be isolated, no **antibody** reaction observed, and, most compelling of all, by the ability of the infectivity to survive irradiation, heat, and enzyme treatments that should destroy nucleic acids. In its place, the prion hypothesis was proposed by Prusiner in 1982 (2) that posited a proteinaceous infectious particle as the scrapie agent (see **Prion**). The essence of the prion hypothesis is that the prion protein (PrP) exists in two states, a normal, cellular form (PrP<sup>C</sup>), and a pathologic scrapie form (PrP<sup>Sc</sup>). The normal and scrapie forms of the prion differ not covalently but in the three-dimensional fold of their polypeptide chains (11). Spectroscopic examination has demonstrated that the PrP<sup>C</sup> molecule is a predominantly **alpha-helical** protein, whereas PrP<sup>Sc</sup> molecule has much more **beta-sheet** structure. Moreover, PrP<sup>C</sup> is sensitive to **proteolysis**, but PrP<sup>Sc</sup> has a large proteinase resistant core, which spontaneously forms insoluble **amyloid** fibrils. Scrapie diseases result from the ability of the scrapie form of the prion (PrP<sup>Sc</sup>) to induce the structurally quite different cellular form (PrP<sup>C</sup>) to take up the structure of the PrP<sup>Sc</sup> form in a chain reaction (12). In this way the scrapie form of the prion protein can effectively replicate itself, and also act as an infectious agent, capable of the induction of more scrapie prions in other hosts.

That a single protein molecule could be responsible for a group of inherited and communicable diseases was a radical idea naturally attracting much skepticism. However, a great deal of experimental verification of the prion hypothesis has been obtained (12), and it has recently received the imprimatur of the Nobel committee. It seems highly probable that the prion protein is indeed the source of the infectious scrapie agent, although the viral hypothesis has still not been entirely abandoned. One of the aspects that encourages retention of the viral origin of the scrapie diseases is the issue of “strains” (3, 4) (see **Prion**). It is known that scrapie prions appear as multiple strains that differ somewhat in their disease expression, for example, Gerstmann–Straussler–Scheinker disease, FFI, and different subtypes of CJD (see list above) are expressed by different point mutations in the prion gene. When these different prion strains are inoculated into mouse brain containing only the single type of mouse prion, the resultant scrapie prions exhibit the characteristics of the inoculated strain. In viral diseases, these different disease signatures would be encoded on the viral DNA/RNA, but in the prion hypothesis they are thought to be encrypted on scrapie prion proteins in the form of slightly different three-dimensional folds of their polypeptide chains. These different folds are assumed to be capable of being “imprinted” on the same cellular prion proteins in their conversion to the scrapie form (14), thereby expressing the particular strain of the infection. A similar line of argument has been used to explain why prions are able to cross some species barriers, but not others (15). The experimental evidence suggests that infection crosses species barriers more easily when the host prion proteins have a closer amino acid **homology** with the infecting prions. For example, **transgenic** mice carrying both mouse and hamster prion genes express more mouse prions when inoculated with mouse prions, and more hamster prions when inoculated with hamster prions. This suggests that prions preferentially recruit homologous prion proteins. Similarly, the spread of scrapie from sheep to cows may be facilitated because the sequences of sheep and cattle PrP differ at only seven positions. It is probable, however, that it is the locations as well as the number of differences that is crucial: how else to explain that scrapie can pass from sheep to cattle, and then to humans, while scrapie cannot apparently pass directly from sheep to humans?

The prion hypothesis has taken our understanding of the scrapie diseases a considerable way forward, and in doing so has added some radically new concepts about [protein structure](#) to molecular biology. The weight of evidence is quite strongly in its favor, but nevertheless some important aspects of the role of prions in molecular mechanisms of the scrapie diseases remain to be uncovered. For example:

1. What is the precise nature of the infectious scrapie agent?
2. In what molecular complex does the structural conversion of the prion proteins take place, and how does the “imprinting” of prion strains occur within it?
3. How are the neurotoxic effects of scrapie prions expressed?
4. How does the agent pass from the gut to the brain?
5. What is its relation to Alzheimer's disease?

More radical ideas will be needed to answer these questions—or to overturn the prion hypothesis.

### Bibliography

1. J. S. Griffith (1967) *Nature* **215**, 1043–1044.
2. S. B. Prusiner (1982) *Science* **216**, 136–144.
3. M. E. Bruce, A. G. Dickinson, and H. Fraser (1976) *Neuropathol. Appl. Neurobiol.* **2**, 471–478.
4. S. J. DeArmond et al. (1985) *Cell* **41**, 221–235.
5. R. Anderson, C. Donnelly, N. Ferguson, and M. Woolhouse (1996) *Nature* **382**, 779–788.
6. C. R. Buchanan, M. A. Preece, and R. D. Milner (1991) *Br. Med. J.* **302**, 824–828.
7. P. Duffy, J. Wolf, G. Collins, A. G. DeVoe, B. Streeten, and D. Cowan (1974) *N. Engl. J. Med.* **290**, 692–693.
8. R. G. Will, J. W. Ironside, M. Zeidler, S. N. Cousins, K. Estibeiro, and A. Alperovich (1996) *Lancet* **347**, 921–925.
9. P. J. Harris and G. W. Roberts (1991) *Br. J. Psych.* **158**, 457–470.
10. S. J. DeArmond (1993) *Curr. Opin. Neurol.* **6**, 872–881.
11. M. Gasset, J. M. A. Baldwin, R. J. Fletterick, and S. B. Prusiner (1993) *Proc. Natl. Acad. Sci. USA* **90**, 1–5.
12. F. E. Cohen, K.-M. Pan, Z. Huang, M. Baldwin, R. J. Fletterick, and S. B. Prusiner (1994) *Science* **264**, 530–531.
13. R. A. Bessen, D. A. Kocisko, G. J. Raymond, S. Nandan, P. T. Lansbury, and B. Caughey, (1995) *Nature* **375**, 698–700.
14. G. C. Telling et al. (1996) *Science* **274**, 2079–2082.
15. S. B. Prusiner (1995) *Sci. Am.* **272**, 48–57

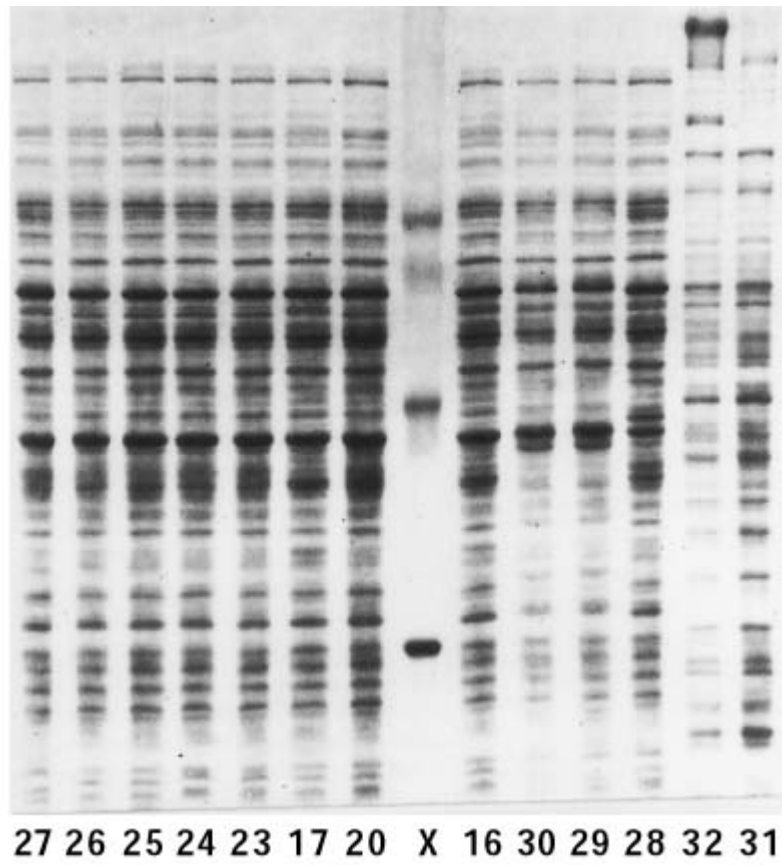
### Suggestions for Further Reading

16. S. B. Prusiner, J. Collinge, J. Powell, and B. Anderton (1992) *Prion Diseases of Humans and Animals*, Ellis Horwood, New York.
17. S. B. Prusiner (1996) Molecular biology and pathogenesis of prion diseases, *Trends Biochem. Sci.* **21**, 482–487.

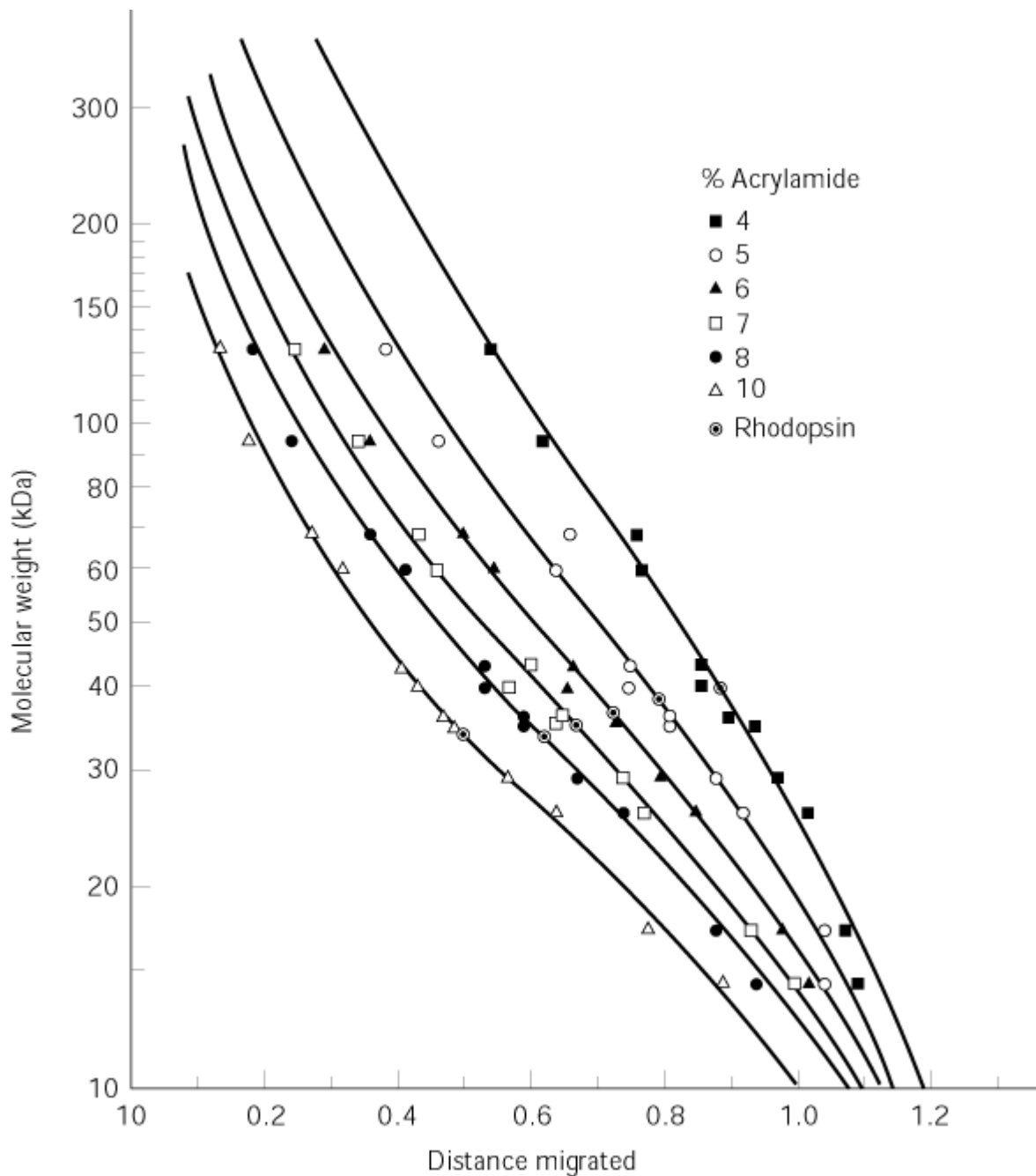
[Electrophoresis](#) of fully **denatured** and dissociated [polypeptide chains](#) derived from [proteins](#) by binding of, and saturation with, the negatively charged detergent **SDS**, usually designated SDS-PAGE, is by far the most popular mode of [gel electrophoresis](#) in [polyacrylamide](#) gels (usually abbreviated PAGE). There are two reasons why SDS-PAGE is preferred: (1) it provides the resolving capacity and **molecular sieving** of PAGE, especially when conducted as [disc electrophoresis](#) (Fig. 1); and (2) it can give an estimate of the molecular weight of each polypeptide chain on the basis of electrophoresis at a single gel concentration. Molecular weights are often measured for protein using this technique alone, but it must be remembered that the results are valid only to the extent that the following simplifying conditions apply:

1. All SDS-denatured proteins tend to have the same net negative charge density: The binding of SDS to the protein usually masks the intrinsic charge of the protein and because of the bound SDS, gives it a standard negative charge density, that is largely independent of pH.
2. All SDS-denatured proteins tend to have the same conformation: The binding of SDS, when coupled with reduction of any [disulfide bonds](#) and dissociation into individual polypeptide chains at high temperatures, produces a conformation like a [random coil](#) for all SDS-polypeptides; therefore, mobility differences due to specific protein conformations are abolished.
3. The mobility is inversely proportional to the size of the polypeptide chain: Because the surface charge densities of proteins are equalized [see (1)], the [Ferguson plots](#) of all SDS-polypeptides should intersect at zero gel concentration; in this case, the ratios of the mobilities for all proteins at various gel concentrations are the same.
4. The mobility is inversely proportional to the logarithm of the molecular weight: A plot of log (molecular weight) versus migration distance is linear except at very large and very small molecular weights. Such a linear plot provides a simple means of deriving the molecular weight of an unknown protein (such as rhodopsin in the example of Fig. 2) by comparison to a set of standard proteins.

**Figure 1.** SDS-PAGE patterns of polypeptides derived from total cell extracts of various strains (numbered) of *Achromobacter*. Lane X contains the protein standards (from top to bottom) ovotransferrin (77 kDa), albumin (66 kDa), ovalbumin (43 kDa), carbonic anhydrase (30 kDa), and myoglobin (17 kDa). The polyacrylamide gel was of 10% (w/v) acrylamide, 2.7% Bisacrylamide, in a discontinuous Tris-chloride-glycinate buffer system [Fig. 2, M. Costas (1990) *Electrophoresis* **11**, 382].



**Figure 2.** Standard curve of  $\log(\text{molecular weight})$  versus relative distance migrated in SDS-PAGE. Twelve standard proteins of known size were separated by SDS-PAGE in polyacrylamide gels of six different concentrations. The protein of molecular weight to be determined was rhodopsin. [Fig. 2 of R. N. Frank and D. Rodbard (1975) *Arch. Biochem. Biophys.* **171**, 1–13].



The four simplifying conditions of SDS-PAGE given above are, however, not entirely or always valid, since:

1. The mobilities are usually constant between pH 7 and 11 but not at lower pH values.
2. The derivatization of proteins with SDS produces polypeptides of uniform random-coiled conformations only when the reaction conditions are sufficiently severe. In particular, the reaction time needs to be sufficient, which in many cases requires prolonged boiling with SDS. Once conformational uniformity and an approximately uniform net charge are attained (see condition 3 above), the mobilities and separations are more reproducible than those of native proteins.
3. Equalization of the surface charge densities of polypeptides by reaction with SDS is incomplete for proteins with extreme net charges or for those containing glyco-, lipo-, or nucleo-moieties. Such incomplete saturation of the protein surface with SDS accounts for the failure of the Ferguson plots of many SDS-polypeptides to intersect at or near the ordinate (1).



4. When the charge densities among proteins are not fully equalized by derivatization with SDS, the ratio of mobilities between size standards and unknowns becomes dependent on gel concentration. Therefore, the molecular weight obtained from a comparison with the mobilities of standards will vary with gel concentration; in that case, accurate molecular weights of SDS-polypeptides, like those of native proteins, can be derived only from the slope of a Ferguson plot (2).
5. Even if surface charge and conformational equalization among proteins are established for particular proteins, the sigmoidal nature of the plot of log(molecular weight) versus migration distance (Fig. 2) needs to be recognized (2). Thus, the approximation of linearity in the central segment of the sigmoidal curve needs to be restricted to a relatively narrow range of migration distances.

There are other common problems with SDS-PAGE:

1. Its SDS concentrations are greater than the [critical micelle concentration](#) (cmc), so that the SDS is present as micelles. Micellar SDS bound to the [tracking dye](#) unstacks at high gel concentrations in disc electrophoresis; so a tracking dye that marks the moving boundary front at low gel concentrations fails to do so at high concentrations (1).
2. Polyacrylamide is oxidative, and so disulfide bonds may be re-formed after proteins enter the gel; also, oxidation of agents used to reduce any disulfide bonds in the original sample may introduce artifactual new components in the gel pattern.
3. There is no theoretical understanding of SDS-PAGE in discontinuous buffer systems that takes account of moving boundaries set up by monomeric or micellar SDS and its complexes with buffer constituents.

For the foregoing reasons, it is always advisable to establish the accuracy of any important molecular weight value by Ferguson plot analysis.

#### Bibliography

1. M. Wyckoff, D. Rodbard, and A. Chrambach (1977) *Anal. Biochem.* **78**, 459–482.
2. D. Rodbard (1976) in *Methods of Protein Separation*, Vol. 2 (N. Catsimpoalas, ed.) Plenum Press, New York, pp. 145–180.

#### Suggestions for Further Reading

3. G. S. Makowski and M. L. Ramsby (1997) "Protein molecular weight determination by sodium dodecyl sulfate polyacrylamide gel electrophoresis". In *Protein Structure: A practical approach*, 2nd ed. (T. E. Creighton, ed.), IRL Press, Oxford, U.K. pp. 1–27.

#### Sec Mutants/Proteins

A compelling approach to any biological process is to generate [mutants](#) that affect it. One such process is [protein secretion](#). Almost all cells in almost all organisms are capable of secreting proteins. Thus it is possible to study secretion in two of the organisms most amenable to genetic analysis, *Escherichia coli* and *Saccharomyces cerevisiae*. Mutants in bacteria and yeast with defective secretion have been called sec mutants.

Proteins are secreted from bacteria by translocating them in an unfolded state directly from the cytoplasm across the plasma [membrane](#) into extracellular space. The *sec* mutations in *E. coli*, therefore, affect the unfolding process (*secB*, Table 1), the translocase or pore in the plasma membrane through which the nascent protein passes (*secY*, *secE*, *secG*, etc.), and the motor that drives the protein through the pore (*secA*).

Yeast, like all eukaryotic cells, translocate newly synthesized proteins across the membrane of the [endoplasmic reticulum](#) (ER). This process closely resembles translocation across the plasma membrane in bacteria. The protein must be kept unfolded by a **molecular chaperone** (Table 1) and is translocated through a pore whose basic constituents are **homologous** to the *secY*, *secE*, and *secG* proteins in *E. coli* (Table 2). There is no eukaryotic equivalent of the ATP-driven (adenosine triphosphates) translocating machinery, encoded in *E. coli* by *secA*.

**Table 1. Mutations Affecting Chaperones**

| Name           | Organism       | Location    | Notes  |
|----------------|----------------|-------------|--|
| <i>SecB</i>    | <i>E. coli</i> | Cytoplasm   | Delivers protein to <i>secB</i>  |
| <i>DNAJK</i>   | <i>E. coli</i> | Cytoplasm   | ?  |
| <i>GroE</i>    | <i>E. coli</i> | Cytoplasm   | Affects pre- <b>b-lactamase</b> . GroEL and GroES together form a complex multisubunit structure ( <a href="#">chaperonin</a> )          |
| <i>Ffh</i>     | <i>E. coli</i> | Cytoplasm   | Part of bacterial ribonucleoprotein <a href="#">signal recognition particle</a> (SRP) complex; interacts with 4-5S RNA                   |
| <i>sec65</i>   | yeast          | Cytoplasm   | Homologue of 19-kDa protein of mammalian SRP. Required for assembly of integral membrane proteins  |
| <i>SRP54</i>   | Yeast          | Cytoplasm   | Homologue of 54-kDa protein of mammalian SRP, also a GTPase  |
| <i>Ssa 1-4</i> | Yeast          | Cytoplasm   | Members of yeast <b>heat-shock</b> protein family required for ATP-dependent posttranslational translocation                             |
| <i>Kar2p</i>   | Yeast          | Lumen of ER | Analogous to mammalian <a href="#">BiP</a> , required for translocation into ER  |
| <i>PDI1</i>    | Yeast          | Lumen of ER | A <a href="#">protein disulfide isomerase</a> that catalyzes <a href="#">disulfide bond</a> formation as the nascent chain enters the ER |
| <i>Fkb2</i>    | Yeast          | Lumen of ER | <b>Peptidyl proline</b> <i>cis-trans</i> <a href="#">isomerase</a>   |

**Table 2. Mutations Affecting Membrane Translocation**

| Name | Organism | Location | Notes |
|------|----------|----------|-------|
|------|----------|----------|-------|

|                              |                |                    |  |
|------------------------------|----------------|--------------------|--|
| <i>secY</i> ,<br>E, G        | <i>E. coli</i> | Plasma<br>membrane | Forms transmembrane channel  |
| <i>secA</i>                  | <i>E. coli</i> | Cytoplasm          | ATP-driven shuttle, pushing protein through translocon. Recognizes newly synthesized protein, complexed with <i>secB</i>                                 |
| <i>secD</i> ,<br><i>secF</i> | <i>E. coli</i> | Plasma<br>membrane | Accessory proteins for bacterial translocon  |
| <i>sec61</i> ,<br>62, 63     | Yeast          | ER                 | Forms transmembrane channel  |
| <i>SRP</i><br>101            | Yeast          | ER<br>membrane     | Homologue of mammalian 61p, a, b, g<br>A GTPase that acts as receptor for a signal recognition particle (SRP); homologue of mammalian $\alpha$ -subunits |
| <i>ft-sY</i>                 | <i>E. coli</i> | ?                  | Bacterial homologue of SRP receptor  |

In eukaryotes, protein secretion also requires membrane traffic from the endoplasmic reticulum to the **Golgi complex**, and from there to the plasma membrane. Each of the two membrane trafficking steps involves segregating cargo proteins from endogenous proteins, forming a carrier vesicle, and fusing the carrier vesicle with the appropriate target. The same categories of proteins are involved in each step: namely, sortases, coating molecules, and SNARE complexes to allow fusion. There are also conserved [Gtpases](#) (guanosine triphosphate hydrolases) that regulate the coating and targeting steps. Finally, the movement of carrier vesicles sometimes requires [cytoskeleton](#) elements, and so secretion can be perturbed by mutations that affect the cytoskeleton.

The steps at which *sec* mutants act are described in the entries [protein secretion](#), [exocytosis](#), and [secretory vesicles/granules](#). Indeed, much of our understanding of these topics has come from analysis of mutants defective in secretion. For convenience, mutations are listed here by their function. It is doubtful that all mutations that affect secretion have been identified, but the overall patterns seem to be clear.

In several cases, yeast mutations that affect secretion were not isolated through a direct screen for a secretion deficiency. Instead, they were selected as mutations that acted as **suppressors** of a deficiency, or as giving synthetic lethality when combined with a weak secretion-deficiency allele. In other cases, they were identified by homology to a mammalian protein of known function. Finally, some were products of other selections but were found also to have a “*sec*” phenotype. Thus not all of the mutants listed in Tables [1](#) and [2](#) have the conventional designation of *sec* followed by a number or letter.

## 1. Chaperones

**Molecular chaperones** are proteins that help keep newly synthesized proteins unfolded until they pass through the transmembrane channel and then help fold them correctly after translocation. Representative examples are given in Table [1](#).

## 2. Membrane Translocation

The newly synthesized proteins cross membranes via a proteinaceous channel. Proteins that recognize the [signal peptides](#) help present the newly synthesized protein to the channel. Mutations

are found in channels and in the receptors that recognize signal peptides.

### 3. Budding from the endoplasmic reticulum

In eukaryotes, newly synthesized proteins leave the ER in a coated vesicle. The COPII coatomers are clearly involved in the ER to Golgi pathway. Since the *ret1* mutation blocks retrieval of dilysine-containing ER proteins from the Golgi, COPIs are involved in intra-Golgi or Golgi to ER traffic. Mutants in COPI may only affect ER to Golgi traffic indirectly. Cargo can be selected from ER-resident proteins by specific “sortases” such as *SHR3* or *emp24*. (See Table 3.)

**Table 3. Mutations in ER to Golgi Traffic**

| Name               | Organism | Location    | Notes  |
|--------------------|----------|-------------|--|
| <i>SHR3</i>        | Yeast    | ER membrane | A sortase required specifically for amino-acid <b>permease</b> export from the ER                                    |
| <i>emp24</i>       | yeast    | ER membrane | Required for invertase secretion. Believed to be a carrier molecule  |
| <i>sec23/24</i>    | Yeast    | Cytosol     | Form a heterodimer; <i>sec23p</i> has GAP activity. With <i>sec 13/31</i> , form a COPII coat on ER-budding vesicles |
| <i>sec13/31</i>    | Yeast    | Cytosol     | Form heterodimers. With <i>sec23/24p</i> , form a COPII coat on ER-budding vesicles                                  |
| <i>saur1</i>       | yeast    | Cytosol     | A yeast GTPase, homologous to ARF1, which mediates coating and uncoating   |
| <i>sec12</i>       | yeast    | ER membrane | Sar1 exchange factor   |
| <i>RET 1/sec26</i> | yeast    | ER membrane | Homolog of the a-, b-, b'-, and g-subunits of <i>sec27</i> , <i>sec21</i> mammalian coatomer (COPI)                  |
| <i>ARF1/2</i>      | yeast    | ER membrane | Yeast homologues of the small mammalian GTPase, ARF1   |

### 4. Fusion of ER-Derived Carrier Vesicles with Golgi Membranes

ER-derived vesicles fuse with the *cis*-Golgi network. Recognition appears to be mediated by specific SNAREs. There is also a parallel recognition pathway involving another small GTPase, *ypt1* or *rab1*, which, like all *rab* proteins, has an isoprene addition to its C-terminal tail. (See Table 4.)

**Table 4. Mutations Affecting Fusion with Golgi Membranes**

| Name | Organism | Location | Notes |
|------|----------|----------|-------|
|------|----------|----------|-------|

|              |       |                    |  |
|--------------|-------|--------------------|--|
| <i>Bct1</i>  | Yeast | ER carrier vesicle | Similar in structure to a v-SNARE  |
| <i>Bos1</i>  | Yeast | ER carrier vesicle | Suppressor of <i>bet1</i> mutations. Required for fusion competence. Homologue of synaptobrevin/VAMP       |
| <i>sec22</i> | Yeast | ER carrier vesicle | v-SNARE  |
| <i>ypt1</i>  | Yeast | ER carrier vesicle | A small ras-like GTPase, homologue of mammalianrab1  |
| <i>bet2</i>  | Yeast | Cytoplasm          | A geranyl geranyl transferase for <i>ypt1</i>  |
| <i>sec18</i> | Yeast | Cytoplasm          | An ATPase required for all fusion reactions; homologue of NSF ( <i>N</i> -ethylmaleimide sensitive factor) |
| <i>sec17</i> | Yeast | Cytoplasm          | Peripheral membrane protein that binds <i>sec18p</i> to membranes; homologue of a-SNAP                     |
| <i>end2</i>  | Yeast | Golgi membranes    | Receptor for proteins carrying HDEL sequence; homologue of mammalian KDEL receptor                         |
| <i>sed5</i>  | Yeast | Post-ER structures | A putative t-SNARE; homologue of mammalian syntaxin  |
| <i>sly1</i>  | Yeast | Golgi membranes    | Homologue of <i>sec1</i> , which interacts with syntaxin at the cell surface                               |

## 5. Golgi to Plasma Membrane Traffic

The basic requirements for transfer between Golgi and plasma membrane are identical to those that are involved in ER to Golgi traffic, since the steps of vesicle formation and selective vesicle fusion are common to both. Sometimes the same proteins—for example, *sec17p* (a-SNAP), *sec18p* (NSF) (NEM-sensitive factor), and ARF1 (ADP-ribosylation factor)—are common to both steps. Table 5 shows mutations that affect exocytoses.

**Table 5. Mutations that affect Exocytosis**

| Name          | Organism | Location           | Notes  |
|---------------|----------|--------------------|--|
| <i>Snc1,2</i> | Yeast    | Secretory vesicles | Homologues of VAMP/synaptobrevin. Redundant gene products  |
| <i>sec7</i>   | Yeast    | Golgi              | Peripheral membrane protein believed to coat Golgi-derived carrier vesicles  |
| <i>kex2</i>   | Yeast    | Golgi membrane     | <a href="#">Serine proteinase</a> that cleaves protein precursors in lumen of Golgi that have Lys-Arg or Arg-Arg sequences. Homologue of mammalian furin |
| <i>kex1</i>   | Yeast    | Golgi              | Removes the two basic residues after   |

|              |       |                          |   |
|--------------|-------|--------------------------|---|
|              |       | membranes                | <i>kex2</i> has made its proteolytic cleavage   |
| <i>sec1</i>  | Yeast | Plasma membrane          | Homologue of <i>nsec1</i> , which binds tightly to syntaxin. Mutants cause secretory vesicle accumulation |
| <i>sec9</i>  | Yeast | Plasma membrane          | Homologue of SNAP-25  |
| <i>sec8,</i> | Yeast | Cytoplasm                | Forms a large 19S complex required for <i>sec15</i> , etc. exocytosis, in addition to SNAREs              |
| <i>sec4</i>  | Yeast | Secretory vesicles       | A small GTPase of the ras family  |
| <i>DSS4</i>  | Yeast |                          | A <b>guanine nucleotide exchange protein</b> acting on <i>sec4</i>  |
| <i>act1</i>  | Yeast | Cytoplasmic cytoskeleton | Mutants cause accumulation of 100-nm secretory vesicles   |
| <i>MYO-2</i> | Yeast | Cytoplasmic cytoskeleton | A novel yeast myosin of the class V-type; causes accumulation of 100-nm vesicles                          |

---

## Second Messengers

In most cases of hormonal [signal transduction](#), the  $\alpha$  subunits of [heterotrimeric G proteins](#) interact directly with one or more effector proteins in the plasma [membrane](#), to induce the release of soluble second messengers that in turn amplify the hormonal signal into the cell. While the number of effectors continues to grow, four distinct types have been studied in detail: [adenylate cyclases](#), [phospholipases](#), [ion channels](#), and phosphodiesterases. All of these enzymes are known to be transmembrane proteins, or proteins that closely associate with the plasma membrane, and they interact with receptor systems to generate second messengers, especially [cyclic AMP](#), diacylglycerol, [inositol phosphates](#), and [calcium](#).

### 1. Effectors of Signal Transduction

#### 1.1. Adenylyl Cyclase

The first well-characterized effector system in signal transduction centered around the discovery of cyclic AMP as a second messenger by Sutherland et al. It is now known that adenylyl cyclases are a large multigene family, with a conserved basic structure. Thus far, nine basic forms of adenylyl cyclase have been identified, with different sensitivities to calcium, G proteins, and regulation by **phosphorylation** (1). Like many ion channels and other membrane transporters, the adenylyl cyclases are characterized by 12  [\$\alpha\$ -helices](#): One intracellular loop occurs between the sixth and seventh transmembrane domain (C1), and another (C2) occurs at the C-terminal tail. The amino terminus of the protein is also thought to be intracellular. While the  [\$\alpha\$ -helices](#) diverge significantly between the family members, the intracellular loops are conserved. The C1 and C2 regions are thought to contain the catalytic activity of the enzyme and to be responsible for much of its regulation.

In addition to regulation by both the  $\alpha$  and  $\beta\gamma$  subunits of [heterotrimeric G proteins](#), some forms of adenylate cyclase are also subject to regulation by the diterpene forskolin, by calcium/[calmodulin](#), and by **phosphorylation** catalyzed by several [kinases](#). It has been speculated that the  $\alpha$  subunit interacts directly with the C1 and/or C2 domains of cyclase, although it remains unclear how the inhibitory actions of G proteins are mediated. Calmodulin is thought to interact with the amino-terminal region of C1 in adenylate cyclase I, whereas phosphorylation by CAM kinase in adenylate cyclase III is thought to occur in the C-terminal tail (1) (see [Calcium Signaling](#)).

### 1.2. Phospholipases

Another second messenger system that plays a key role in the regulation of cellular function involves the generation of small molecules derived from phospholipid metabolism. The hormone-sensitive enzymes that catalyze these reactions are [phospholipases](#). The phosphatidylinositol cycle in particular (see [Inositol Lipids and Phosphates](#)) results in the generation of two potential second messengers, the membrane-associated diacylglycerol, which can activate [protein kinase C](#), and inositol trisphosphate, which can lead to the mobilization of intracellular calcium levels.

### 1.3. Ion Channels

There are many examples of **ion channels** that are regulated by heterotrimeric G-protein  $G_{\alpha}$  subunits in a manner analogous to that described for adenylate cyclase and [phospholipase C](#) (2). Numerous studies have focused on voltage-gated  $Na^{+}$ ,  $K^{+}$ , and  $Ca^{2+}$  channels in different tissues, particularly endocrine and neuronal cells. In general, receptors linked to slow synapses in the central nervous system are coupled to ion channels through G proteins. These changes in ion permeability are slower than those induced by ligand-gated channels (see text below) and are generally longer-lasting. In many cases, ion channels are regulated by  $G_i$  subunits.

The voltage-sensitive ion channels that are regulated by G-protein interactions have been the subject of extensive investigation. While there has been significant progress in elucidating the subunit structures of these channels and their regulation by various [toxins](#), the molecular details of G-protein interaction remains cloudy. Probably the best-studied system has been the regulation of  $K^{+}$  channels by the muscarinic [acetylcholine receptor](#). In this case, the binding of acetylcholine to the receptor activates a  $G_i$  protein, leading to the dissociation of  $\alpha_i$  from  $\beta\gamma$ . It is then thought that  $\alpha_i$  can directly activate the  $K^{+}$  channel, although  $\beta\gamma$  subunits can interact with the channel as well (2, 3).

### 1.4. Phosphodiesterases

Cyclic nucleotide phosphodiesterases (PDEs) are a large superfamily of [enzymes](#) that play an important role in the termination of [cyclic AMP](#) (cAMP) and [cyclic GMP](#) (cGMP) signaling. Some forms of PDE specifically recognize one cyclic nucleotide, while others are less specific. The best-studied of these enzymes are the cGMP-specific forms of PDE. These enzymes are regulated directly by G proteins, such as [rhodopsin](#) in rod outer segments. In this system, the activation of rhodopsin by light leads to activation of the [heterotrimeric G protein](#) known as [transducin](#). In its GTP-bound, activated state, transducin binds directly to cGMP PDE, inducing its activation through the release of the  $\beta\gamma$  subunit of PDE (4). For the other isoforms of PDE, the modes of regulation are less well established, although some may involve feedback **phosphorylation**.

## 2. Targets

Many of the effects of second messengers on cells are mediated by the activation of protein **tyrosine kinases** and **serine-threonine kinases**. Indeed, protein **phosphorylation** is the major currency by which signals are transmitted in cells. While there have been over 1000 protein kinases identified to date, only a subset are thought to be activated by soluble second messengers. The second messenger-activated kinases can be categorized based on the upstream processes that lead to their regulation. They are [protein kinase A](#) (or cyclic AMP-dependent protein kinase), [protein kinase C](#), and CAM kinase, which is a critical target in [calcium signaling](#).

## Bibliography

1. R. Taussig and A. G. Gilman (1995) *J. Biol. Chem.* **270**, 1–4.
2. A. M. Brown and L. Birnbaumer (1990) *Annu. Rev. Physiol.* **52**, 197–213.
3. B. Hille (1992) *Neuron* **9**, 187–195.
4. M. D. Houslay and G. Milligan (1997) *Trends Biochem. Sci.* **22**, 217–224.

## Secondary Structural Prediction Of Proteins

The **secondary structure** of [proteins](#) is the two-dimensional arrangement of the **polypeptide** backbone, defined primarily by a network of [hydrogen bonds](#). The two most common and most regular secondary structures are the [alpha-helix](#) and the [beta-sheet](#). The beta-sheet is made up of [b-strands](#) arranged in either parallel or antiparallel fashion. The goal of secondary structural prediction is a relatively modest aspect of **protein structural prediction** to infer the location of a-helices and b-strands along the polypeptide chain without necessarily determining how the latter associate into a b-sheet. Therefore, the prediction problem is formally no longer two-dimensional, but it is addressed to the correspondence between two linear strings, the amino acid sequence and the linear sequence of elements of secondary structure in the sequence. The secondary structure is often denoted by letters: A, a-helical; B, b-stranded; C, coil or other conformations. Any predictive algorithm produces a string of A's, B's, or C's along the sequence being analyzed, and then this outcome must be compared with the real secondary structure, when known. The success rate of the prediction is measured simply as the percent of residues that were predicted correctly, counted on a residue-by-residue basis. The random level of the three-state predictive accuracy might be expected to be close to 33%, but it is actually about 40%, because of an uneven distribution of the three states in known structures: on average, the coil is the most abundant, and the b-strand the rarest.

A wide variety of predictive methods have been developed to date, but the basic logic of most of them is like that developed by Chou and Fasman (1). It is assumed that individual amino acid residues have their own intrinsic propensities to form a-helix, b-strand, or coil conformations. For example, every [alanine](#) residue has a certain propensity for A, B, and C conformations, and every [valine](#) has different propensities, as do the other amino acid residues. These propensities are assumed to be additive along a sequence, to give a propensity for any segment of residues to adopt each conformation. The propensities for individual residues are estimated from the known protein structures. This type of logic can be extended to a “directional” propensity (2), which incorporates the influence of near-neighbor residues, or it can be extended to pairs (3) or triplets (4) of residues. [Neural networks and genetic algorithms](#) operate similarly, predict the secondary structural state of the central residue in each window of about 10 residues, and take into account the influences of all the residues within the window on the conformational state of the central residue (5). The parameters used are estimated by the computer “learning” from the known protein structures using the neural network algorithm. The simple application of neural networks, assuming only the input and output layers, is almost identical in logic to the directional propensity method (6), but it becomes more complicated upon adding “hidden layers” inserted between the input and output layers. The “hidden Markov” method (HMM) is another technique that was developed in informatics theory and applied to protein secondary structural prediction (7).

The number of parameters required for the simple method of Chou and Fasman (1) is only 60 (20 amino acids × 3 states each), about 1000 for the “directional” propensity (20 amino acids × 17 positions of neighboring residues × 3 states) (2), and more than 10,000 for the method using residue-



pair propensities ( $20 \times 20$  amino acid pairs  $\times$  10 positions  $\times$  3 states) (3). Only 29 protein structures with a total of about 6,000 residues were known at the time Chou and Fasman developed their method (8). This was sufficient to derive the 60 parameters for their method but was dubious for the second method and very insufficient for the third. As the size of the structural [databases](#) increased, much more data became available, making it possible to derive reliable parameters for more complex methods, and the predictive accuracy has also increased. Whereas the initial predictive accuracy for the Chou–Fasman technique was a significant 55% or less (9, 10), that figure has now reached 70% or more (11, 12). For example, a neural network composed of five independent networks can use a number of homologous sequences aligned with each other (11). Such a treatment is logical because proteins that have homologous sequences adopt basically the same fold. So using a number of them increases the information available (see [Homological Modeling](#)). The predictive methods just described use as input only the local sequence information, such as the neighboring residues within a narrow window. Recently, algorithms to incorporate global features of the input sequence have also been developed (13, 14). Thus, secondary structure predictions are still under development, and the structural database is growing by 40% each year.

How much further can the predictive accuracy increase? It is almost certain that 100% accuracy is impossible. One of the reasons is simply technical. Assignment of secondary structure in real protein structures involves some ambiguity because the local structure can be irregular. The computer program DSSP (15) is widely used to assign the ranges of  $\alpha$ -helix and  $\beta$ -strands along a polypeptide chain automatically from the atomic coordinates of the structure determined experimentally. X-ray crystallographers, on the other hand, often assign the secondary structure by manual inspection and record their own assignments in the Protein Data Bank (PDB) file (16). It would not be correct to consider one assignment more correct than the other. The DSSP assignments sometimes appear too strict to the eye, although they are the more objective, so long as the atomic coordinates of the model are correct. When these two types of assignments are compared with each other in 200 well-resolved and unrelated structures, the percentage of identical three-state assignments is about 90%. This figure may indicate the upper limit to the predictive accuracy.

## Bibliography

1. P. Y. Chou and G. D. Fasman (1974) *Biochemistry* **13**, 211–221; 222–245.
2. J. Garnier, D. J. Osguthorpe, and B. Robson (1978) *J. Mol. Biol.* **120**, 97–120.
3. J.-F. Gibrat, J. Garnier, and B. Robson (1987) *J. Mol. Biol.* **198**, 425–443.
4. K. Nagano (1977) *J. Mol. Biol.* **109**, 251–274.
5. N. Qian and T. J. Sejnowski (1988) *J. Mol. Biol.* **202**, 865–884.
6. K. Nishikawa and T. Noguchi (1991) *Methods Enzymol.* **202**, 31–44.
7. K. Asai, S. Hayamizu, and K. Handa (1993) *CABIOS* **9**, 141–146.
8. P. Y. Chou and G. D. Fasman (1978) *Ann. Rev. Biochem.* **47**, 251–276.
9. K. Nishikawa (1983) *Biochim. Biophys. Acta* **748**, 285–299.
10. W. Kabsch and C. Sander (1983) *FEBS Lett.* **155**, 179–182.
11. B. Rost and C. Sander (1993) *J. Mol. Biol.* **232**, 584–599.
12. A. A. Salamov and V. V. Solovyev (1995) *J. Mol. Biol.* **247**, 11–15.
13. D. Frishman and P. Argos (1996) *Protein Eng.* **9**, 133–142.
14. M. Ito, Y. Matsuo, and K. Nishikawa (1997) *CABIOS* **13**, 415–423.
15. W. Kabsch and C. Sander (1983) *Biopolymers* **22**, 2577–2637.
16. F. C. Bernstein, T. F. Koetzle, G. J. B. Williams, E. F. Meyer, M. D. Brice, J. R. Rodgers, O. Kennard, T. Shimanouchi, and M. Tasumi (1977) *J. Mol. Biol.* **112**, 535–542.

## Secondary Structure, Protein

[Protein structure](#) is classified in a hierarchical manner into [primary structure](#), secondary structure, [tertiary structure](#), and [quaternary structure](#). The secondary structure of a protein refers to the **backbone conformation** of the [polypeptide chain](#). There are several different types of secondary structure elements found in proteins, including the regular  [\$\alpha\$ -helix](#) and  [\$\beta\$ -sheet](#) as well as the nonregular [omega loops](#) and [turns](#) that reverse the direction of the polypeptide chain. Regular secondary structure is defined by repeating backbone [dihedral angles](#) and [hydrogen bonds](#) of the peptide backbone. Nonregular secondary structure does not have repeating backbone angles, and the backbone hydrogen-bonding pattern is irregular.

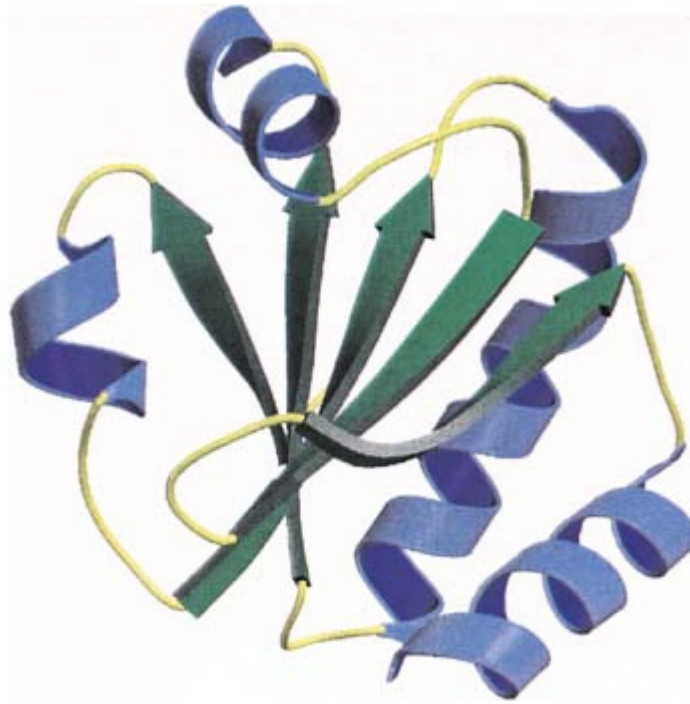
The peptide backbone dihedral angles around the  $N-C_{\alpha}$  and  $C_{\alpha}-C$  bonds are referred to as phi ( $\phi$ ) and psi ( $\psi$ ), respectively. Only certain combinations of  $\phi/\psi$  angles are allowed due to steric restraints imposed by the amino acid side chain;  $\phi$  is almost always negative (see [Ramachandran Plot](#)). The **glycine residue**, with only a hydrogen for its [side chain](#), has a high degree of backbone conformational flexibility and can adopt  $\phi/\psi$  angles that are forbidden for other amino acid residues (for example, positive  $\phi$  angles). Generally, over 90% of nonglycine amino acid residues in a protein structure have backbone conformations close to those of the  $\alpha$  or  $\beta$  types of secondary structure. This can be visualized in the Ramachandran plot, which shows the allowed combinations of  $\phi/\psi$  angles as discrete regions or areas and is used to plot observed  $\phi/\psi$  combinations for each residue in a protein structure.

Different types of secondary structure have different  $\phi/\psi$  combinations and different [hydrogen bond](#) patterns. For example, the  $\alpha$ -helix has  $\phi$  and  $\psi$  angles of close to  $-60^{\circ}$  and  $-50^{\circ}$ , respectively, and is stabilized by ( $i+4$ ) hydrogen bonds—that is, hydrogen bonds formed between the backbone carbonyl oxygen of one residue ( $i$ ) and the backbone nitrogen of another located four residues further on in the sequence ( $i+4$ ). In contrast, the  $\beta$ -strand structure is characterized by an almost fully extended polypeptide backbone conformation ( $\phi$  and  $\psi$  angles of close to  $-120^{\circ}$  and  $140^{\circ}$ , respectively) that forms hydrogen bonds with backbone atoms of adjacent  $\beta$ -strands to form  $\beta$ -sheets.

In the absence of a complete three-dimensional structure of a protein, its secondary structure may be predicted or determined experimentally. **Secondary structure prediction** is based on the observation that individual amino acid residues have different tendencies to form  $\alpha$ -helices,  $\beta$ -sheets, or turns. The sequence of amino acid residues in a protein, its [primary structure](#), is therefore used to predict its secondary structure. [Circular dichroism](#) measurements of a protein can be used to estimate the relative proportion of each type of secondary structure in a protein, and [nuclear magnetic resonance \(NMR\)](#) can be used to assign the secondary structure type for each amino acid residue.

The three-dimensional structures of proteins consist of many thousands of atoms and are very complex. Often, their structures are represented simply by showing the backbone trace of the protein, with the secondary structure components identified (Fig. 1).

**Figure 1.** Representation of the backbone structure of the oxidoreductase protein [thioredoxin](#) (1) showing the secondary structure elements.  $\alpha$ -Helices are shown as blue coils, and  $\beta$ -strands are shown as green arrows. Loops and turns are in yellow. This figure was generated using Molscript (2) and Raster3D (3, 4). See color insert.



[See also [Protein Structure](#), [Alpha-Helix \(310-Helix and Pi-Helix\)](#), [Beta-Sheet](#), and [Turns](#).]

#### Bibliography

1. S. K. Katti, D. M. Le Master, and H. Eklund (1990) *J. Mol. Biol.* **212**, 167–184.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

5. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
6. T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, New York.
7. J. S. Richardson (1981) The anatomy and taxonomy of protein structure. *Adv. Protein Chem.* **34**, 167–339.
8. W. Kabsch and C. Sander (1983) Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637.
9. A. V. Efimov (1993) Patterns of loop regions in proteins. *Curr. Opin. Struct. Biol.* **3**, 379–384.

#### Secretion Vector

A secretion vector is a specialized expression **vector** (see [Expression Systems](#)) that directs the production of [protein](#) for locations other than the cytoplasm. Transport out of the cytosol can be achieved by fusing the gene of interest with a nucleotide sequence encoding the [signal peptide](#) of a

readily **secreted protein**. The use of a secretion strategy has many potential advantages, including high yields, simplified purification, improved stability; escape from cytoplasmic proteins, especially [proteinases](#), prevention of toxicity; avoidance of an N-terminal methionine residue extension (mainly a problem in **prokaryotic** expression systems); and localization in a more oxidizing environment where [disulfide-bond](#) formation may proceed and the protein may fold correctly, alleviating formation of [inclusion bodies](#) (which is mainly a problem that occurs on overexpression of a protein of **eukaryotic** origin in a prokaryotic host).

Secretion vectors can be designed for a broad spectrum of hosts, including both **Gram-negative** and Gram-positive bacteria, hosts ranging from primitive eukaryotes to higher eukaryotes, including mammals.

## 1. Gram-negative Bacteria

In Gram-negative bacteria, such as *Escherichia coli*, proteins can be targeted to the **periplasm** (the compartment between the cytoplasmic and the outer membrane), the outer membrane, or the culture medium. By definition, export takes place in the first two cases, whereas only in the latter case does true secretion occur. Proteins directed to the periplasm can be specifically released by osmotic shock treatment, providing an enrichment of the protein. Three pathways exist for protein targeting in Gram-negative bacteria: the type I pathway (also known as the *single-step pathway*, the *signal sequence-independent pathway*, or the *hemolysin pathway*); the type II pathway (also known as the *two-step pathway* or the *general secretion pathway*, or GSP); and the type III secretion pathway.

### 1.1. Type I Pathway

In the type I pathway, the proteins are translocated across both the cytoplasmic and the outer membrane in a concerted manner. The best-known protein secreted in this way is hemolysin. The signals targeting the proteins to the type I pathway are located at the carboxyl end of the molecule, but there is disagreement on the **consensus sequence** of this “secretion signal.” The genes essential for type I secretion in *E. coli* include *hlyB*, *hlyD*, and *tolC*. The first two are not found in laboratory strains of *E. coli* but can be introduced into them. The key element in this pathway, the carboxy-terminal domain of the hemolysin protein, can be fused to a “passenger” protein. The resulting **fusion protein** will, in many cases, be efficiently secreted into the extracellular medium. If the native form of the protein of interest is required, a specific [proteinase](#) cleavage site should be introduced between the “passenger” and the “carrier.” A drawback of this pathway is the unpredictability of what makes a good passenger (1).

### 1.2. Type II Pathway

Most secretion vectors use the type II pathway, in which the protein to be secreted is equipped with an amino-terminal **signal sequence** (also known as a *leader sequence* or *presequence*). This sequence marks the protein for transport through the cytoplasmic membrane, which is accomplished by the **Sec protein** complex in *E. coli* and presumably by **homologous** pathways in other species. These signal sequences share characteristic features with the signal sequences directing proteins to the extracellular medium in Gram-positive bacteria or the lumen of the [endoplasmic reticulum](#) (ER) of eukaryotes: the N-terminal region is rich in positively charged amino acid residues, the core domain is hydrophobic, and the C-terminal region contains the signal recognized by a [signal peptidase](#), which removes the signal peptide during the translocation process. The signal sequences of secreted *E. coli* proteins, such as [OmpA](#), [maltose-binding protein](#), [alkaline phosphatase](#), and  $\beta$ -lactamase, have been used with success for the export of heterologous proteins. As another example, the **promoter** and secretion signals of staphylococcal [protein A](#) (SpA) have been shown to be functional also in the Gram-negative *E. coli*. Protein A fusions are secreted to the periplasm of the bacteria and, in some cases, to the culture medium; from there, they can be readily purified by immunoglobulin G [affinity chromatography](#) (see [Fusion Gene, Fusion Protein](#)).

Proteins to be secreted through the outer membrane reside temporarily in the periplasm, where folding and assembly are carried out, before transport through the outer membrane can take place

(2). In principle, the fusion of a signal sequence to the protein of interest should direct the protein into the periplasmic space. Not all proteins are compatible with the type II pathway, however. In some examples, chimeric heterologous proteins are incorrectly folded, which blocks their transport into the periplasm. The overexpression of various **molecular chaperones**, eg, **DnaK/DnaJ** or various components of the Sec system, will, in some cases, relieve this problem. In other cases, however, more specific folding catalysts are needed. If the folding process in the periplasm is inefficient, the accumulated, misfolded proteins may induce the **heat-shock** response, leading to increased proteolysis.

In general, the success of a secretion strategy using a bacterial host is most likely with proteins that are naturally secreted bacterial exoproteins. More problems are often encountered when the protein of interest is of eukaryotic origin, and the magnitude of the problems seems to correlate with the size of the protein in question. Certain large eukaryotic proteins require the addition of [post-translational modifications](#) in order to be secreted properly. Also, the assistance of accessory factors not present in the bacteria may be needed.

### 1.3. Type III Pathway

The type III pathway is present in certain pathogenic bacteria and has been found to secrete proteins involved in virulence toward animals and **plants**. The signal sequence involved in this pathway is positioned at the N-terminus of the proteins. It does not have the characteristics of a classical signal peptide, however, and is not cleaved off during the translocation.

### 1.4. Leakage through the Outer Membrane

Another mode of secretion is semispecific secretion, in which the expression of a recombinant protein in the periplasm causes the outer membrane to destabilize and become permeable, allowing the protein to diffuse from the cell in a partly specific manner. The production level using this approach can be rather high, but the bacteria will often be killed by the destruction of their outer membrane. A similar result is obtained by the use of so-called leaky mutants with a defective outer membrane, from which the periplasmic proteins are slowly released (3). Another approach to inducing limited leakage of the outer membrane utilizes the coexpression of the [bacteriocin](#) release gene (4). Protoplast type L-forms of various bacteria, including *E. coli*, lack the outer cell wall and the periplasmic space, which permits direct secretion of recombinant proteins into the extracellular medium (5).

## 2. Gram-positive Bacteria

The outer membrane of Gram-negative bacteria presents a strong barrier to the secretion of proteins. Gram-positive bacteria do not have this outer membrane to prevent exported proteins from being released from the cell. Gram-positive *Bacillus* have become very popular hosts for expression, mainly as a result of their potential to secrete proteins directly into the culture medium. Signal sequences of Gram-negative (type II pathway) and Gram-positive microorganisms share the same structural features. One of the serious disadvantages of using *B. subtilis* as an expression host is the secretion of high levels of proteinases. Proteinase-deficient strains that grow normally have been constructed, but particularly susceptible proteins may still be difficult to produce. Another *Bacillus* species, *B. brevis*, has an advantage over *B. subtilis* in that its level of extracellular proteinase activity is very low (6). Mammalian proteins are often toxic to *B. brevis*, and it is therefore necessary to screen a large number of transformants for particular clones that produce the mammalian protein in large amounts. Furthermore, the productivity is highly dependent on growth conditions (eg, medium composition and temperature), and these have to be optimized individually for each new protein produced.

## 3. Eukaryotic Secretion Systems

The eukaryotic secretory pathway is similar to the type II pathway of Gram-negative bacteria. The emergence of the signal peptide from the [ribosome](#) synthesizing the protein to be secreted causes its

binding to the [signal-recognition particle](#) (SRP) and cessation of translation. The resulting complex is targeted to the rough [endoplasmic reticulum](#) (ER) via an interaction with the SRP receptor. Translation of the protein is then resumed, and the polypeptide chain is translocated across the ER membrane into the lumen of the ER, where the signal peptide is processed by the ER membrane-bound signal peptidase. Now maturation of the protein starts, including its folding and addition of post-translational modifications. This includes quality control to ensure that the protein fold is compatible with the secretion machinery. If this is found not to be the case, the protein is degraded, and this stage is where the major loss of heterologous proteins takes place. As the next step, the protein is transported in **vesicles** to the [Golgi apparatus](#), where further maturation occurs. Finally, the protein is transported in vesicles to the cell wall and secreted into the extracellular medium. Heterologous proteins have a propensity for folding incorrectly in the ER, either because the levels of factors required for proper folding/processing are too low to handle the increased throughput or because the [enzymes](#) required to perform the authentic post-translational modifications are lacking. As an example, [disulfide-bond](#) formation catalyzed by [protein disulfide isomerase](#) (PDI) takes place in the ER lumen.

To ensure efficient secretion, it is usually advisable to choose the leader sequence of a homologous or closely related, naturally well-secreted protein. However, examples exist in which a signal sequence of a more distantly related species turned out to be the most efficient, such as the human versus honeybee signal sequence in the baculovirus/insect cell system (7). The precise natures of signal sequences are not fully understood. It is not possible to graft such signals onto a particular gene and to guarantee efficient export or secretion of the resulting product. Proteins that are not normally secreted may contain sequences that are incompatible with the export apparatus, and the attachment of a signal sequence to such a protein may kill the host cell.

As mentioned earlier, a secretion or export strategy may solve a protein degradation problem but, in some cases, it may also worsen the problem. Overloading the secretion machinery may block secretion and cause precursor molecules to accumulate, leading to the induction of a **stress response**; this often involves the production of several proteinases.

A special variant of secretion allows the display of peptides and proteins on the surface of a cell (8). This demands that the corresponding gene be fused to, or inserted into, the gene encoding an outer membrane protein or a cell surface structure. Such cells may be useful as live vaccines, as whole cell adsorbents, in biocatalysis or bioremediation, or in drug hunting. Also, the study of protein–ligand interactions may be facilitated. In eukaryotes, intracellular targeting to **organelles**, such as [mitochondria](#) or the [nucleus](#), is another option that permits the study of the effects of a recombinant protein within a specific subcellular location (9).

Even though several examples of the successful application of a secretion strategy exist, a more thorough understanding of secretion mechanisms is required before predictable manipulations of secretion systems can be made to secrete native recombinant proteins that are not naturally targeted to the extracellular compartment.

## Bibliography

1. M. A. Blight and I. B. Holland (1994) *Trends Biotechnol.* **12**, 450–455.
2. A. P. Pugsley (1993) *Microbiol. Rev.* **57**, 50–108.
3. R. E. Hancock (1984) *Annu. Rev. Microbiol.* **38**, 237–264.
4. S. W. Altmann et al. (1995) *Protein Expr. Purif.* **6**, 722–726.
5. J. Gumpert et al. (1996) *J. Basic Microbiol.* **36**, 89–98.
6. S. Udaka and H. Yamagata (1993) *Meth. Enzymol.* **217**, 23–33.
7. B. S. Mroczkowski et al. (1994) *J. Biol. Chem.* **269**, 13522–13528.
8. G. Georgiou et al. (1996) *Nature Biotechnol.* **15**, 29–34.
9. L. Persic et al. (1997) *Gene* **187**, 1–8.

### Suggestions for Further Reading

10. R. J. Gouka, P. J. Punt, and C. A. M. J. J. van den Hondel (1997) Efficient production of secreted proteins by *Aspergillus*: progress, limitations and prospects. *Appl. Microbiol. Biotechnol.* **47**, 1–11, Concentrates on *Aspergillus* but gives a valuable overview of the secretion process in eukaryotes in general and points out putative problems and the corresponding solutions.
11. M. Sandkvist and M. Bagdasarian (1996) Secretion of recombinant proteins by Gram-negative bacteria. *Curr. Opin. Biotechnol.* **7**, 505–511. Review of the secretion apparatus of Gram-negative bacteria. Fine reference source.
12. J. A. Stader and T. J. Silhavy (1990) Engineering *Escherichia coli* to secrete heterologous gene products. *Methods Enzymol.* **185**, 166–187. Useful overview of the secretory pathways of *E. coli* and how to engineer the organism to secrete recombinant proteins.
13. S.-L. Wong (1995) Advances in the use of *Bacillus subtilis* for the expression and secretion of heterologous proteins. *Curr. Opin. Biotech.* **6**, 517–52. Reviews the use of *Bacillus subtilis* as a secretion host. Provides many references.

### Secretory Vesicles/Granules

[Enzymes](#), [growth factors](#), [extracellular matrix](#) proteins, and signaling molecules are all secreted by cells by fusion of a secretory vesicle with the plasma membrane, releasing the vesicular contents (see [Exocytosis](#)). All cells have constitutive secretory vesicles, which carry newly synthesized proteins directly from the Golgi complex to the cell surface (see [Protein Secretion](#)). Dedicated secretory cells, such as neuronal, endocrine, and exocrine cells, divert classes of secretory proteins out of the constitutive pathway into a specialized class of secretory vesicles, the secretory granules (1), which are stored in the cytoplasm until the cell receives an appropriate stimulatory signal. The term *granule* is a historical misnomer, deriving from the observations of early morphologists, who saw the granular content in their [electron microscopy](#), before the limiting membrane was seen. The protein concentration of secretory granules can be so high, about 0.1 g/ml, that the proteins condense to give secretory granules a solid core that is dense in electron micrographs. A commonly-used alternative name is therefore dense core secretory granule or vesicle. Granulocytes and platelets also have dense core secretory granules that resemble, but are not identical to, those of neurons and endocrine and exocrine cells. Dense core secretory granules are all examples of secretory vesicles derived from the biosynthetic pathway, whose main function is to carry newly synthesized proteins to the cell surface (2). A fourth major group of secretory vesicles, the synaptic vesicles, are generated by [endocytosis](#) from the cell surface (3, 4). Because they do not form at the **Golgi complex**, synaptic vesicles cannot contain newly-synthesized proteins. Instead they secrete small molecules, such as acetylcholine, glutamate, glycine, catecholamines, and g-amino-butyric acid, which they take up directly from the cytoplasm using specialized membrane **transporters**.

Regulated secretory vesicles accumulate in the cytoplasm because their exocytosis is normally inhibited. The accumulation of such vesicles in the cytoplasm is a defining morphological feature of endocrine and exocrine cells, granulocytes, and neurons. An appropriate extracellular signal can remove the inhibition, leading to a massive release of the contents stored in the cytoplasmic vesicles. It is the capacity of such cells to trigger release from a stored pool that gave rise to the term *regulated* secretory vesicles, to contrast them with constitutive secretory vesicles, whose exocytosis occurs in the absence of an extracellular stimulus.

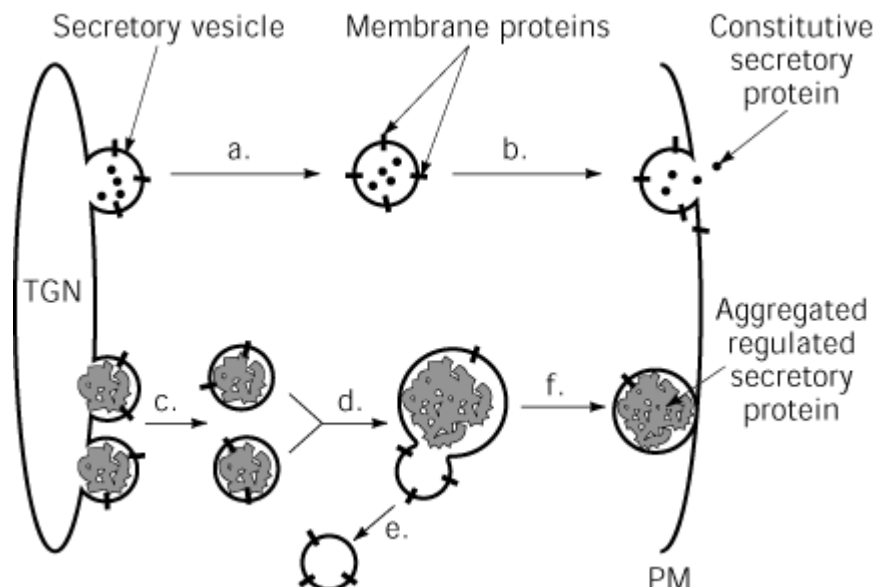
## 1. Secretory Granules from Endocrine Cells

### 1.1. Formation

In endocrine cells, proteins in the *trans* region of the Golgi complex, the *trans*-Golgi Complex, are segregated into two export pathways. Whereas some exit by the conventional constitutive route, others are segregated into the specialized regulatory pathway. Proteins that are segregated into a regulated pathway have a feature absent from constitutively secreted proteins. Thus, if a chimera is made of a regulated pathway protein and a constitutive pathway one, the chimera is sorted into the regulated pathway (5). The feature that allows regulated pathway proteins to be sorted away from other secreted proteins may not be the usual **signal sequence**, a short stretch of contiguous amino acids residues, but it appears to be a region that encourages the formation of large protein complexes in the milieu of the *trans*-Golgi complex (6). An exception to this might be the [hormone](#) pro-opiomelanocortin, which has been reported to have a signal peptide recognized by a receptor (7). Condensation of proteins that will become the contents of secretory granules can often be seen in the lumen of *trans*-Golgi complex membranes, prior to secretory granule formation.

To form an immature secretory granule, a portion of the *trans*-Golgi network pinches off, trapping an aggregate of the secretory granule proteins. Several of these immature granules may fuse to give a precursor that contains all of the content of a secretory granule, but an excessive amount of membrane. To mature, excess membrane is removed from the secretory granule by a series of vesiculation steps (Fig. 1). The immature secretory granule retains some of the features of the Trans Golgi Network, since the excess membrane includes proteins that are destined to go to prelysosomal compartments (8). Maturation of the secretory granule requires coating molecules, [clathrin](#), and the heterotetrameric adaptor, AP1, recruited to the membranes by the small [GTPase](#), adenosyl ribosylating factor-1 (ARF-1) (9, 10), as in the case of other coating events in membrane traffic (see [Protein Secretion](#)).

**Figure 1.** Formation of constitutive and regulated secretory vesicles. (a–c) Both constitutive and regulated secretory vesicles bud from the *trans*-Golgi complex (TGN). (b) Constitutive secretory vesicles go directly to the cell surface, where they immediately fuse with the plasma membrane (PM) and release their protein contents. (c) Immature regulated secretory vesicles are formed by budding aggregates of regulated secretory proteins from the TGN. (d) Immature secretory vesicles fuse to form larger granules, (e) from which excess membrane and proteins can be retrieved and trafficked elsewhere in the cell. (f) The now mature secretory granule goes to the PM, where its exocytosis is blocked until the cell receives the proper signal to release its specialized proteins.





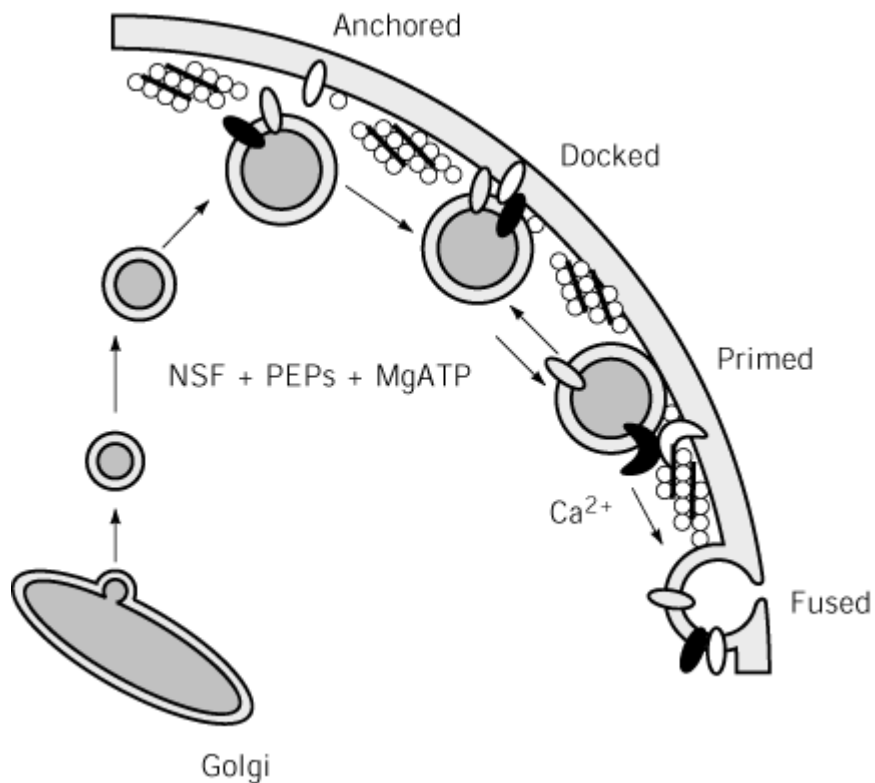
The core of the secretory granule contains proteins specific to the class of endocrine cells, (eg, hormones such as [insulin](#) or [growth hormone](#)). Most secretory granules also contain two other classes of protein, independent of their hormone content. One of these is the chromogranins: highly charged, and often sulfated, proteins that facilitate the condensation of hormones into a proteinaceous core (6). The second is a family of furin-related proteinases (see [Protein Secretion](#)) that act to convert prohormones into their biologically-active forms (11). These prohormone convertases act on prohormones after they have been sorted into secretory granules. Prohormone convertases and chromogranins also have sorting domains that target them to the regulated pathway rather than the constitutive pathways. The exact nature of the sorting domain remains obscure. Unlike retention signals for the [endoplasmic reticulum](#), or internalization signals from the plasma membrane, no amino acid sequence motif has yet been identified for the sorting of convertases.

## 1.2. Docking and Fusion of Secretory Granules

To release their protein contents, secretory granules must first make contact with the plasma membrane, a process called docking. Then the secretory granule membrane fuses with the plasma membrane, a process called [exocytosis](#) (qv). Before docking occurs, the secretory granule migrates from its site of formation, the *trans*-Golgi network, to the plasma membrane. Movement is along [microtubules](#), probably mediated by a “plus-end” directed [motor protein](#) of the [kinesin](#) family (1). Disruptors of microtubule structure, such as [colchicine](#) or [vinblastine](#), inhibit the regulated secretion of hormones. The next step in reaching the cell surface requires the secretory granules to pass through the layer of cortical [cytoskeleton](#) comprised of [actin](#), which underlies the plasma membrane. This process, mobilization, requires that the cortical cytoskeleton be temporarily removed. Actin-depolymerizing enzymes, such as gelsolin, have been linked to the mobilization process (12). Finally, the secretory granules can associate with the plasma membrane, ready for exocytosis. The molecular links between secretory granules and the plasma membrane are not known, but they are sufficiently strong that they can survive homogenization.

Some of the steps involved in the fusion reaction are now known. Two of the first steps require ATP. Before the secretory granule is competent to fuse, the [phosphatidylinositol](#) phospholipids in the cytoplasmic leaflet of its phospholipid bilayer must be correctly phosphorylated (see [Inositol Lipids and Phosphates](#)). This requires a phosphatidylinositol 5' [kinase](#) to be present in the cytosol. The reaction is facilitated by a phosphatidylinositol-transfer protein present in the cytoplasm (13, 14) (Fig. 2).

**Figure 2.** Docking and fusion of secretory granules. Secretory granules travel from the Golgi, where they are formed, to the plasma membrane, where they become anchored once they pass through the layer of cortical actin cytoskeleton. Interaction between proteins (probably SNAREs) on the opposing membranes allows the vesicles to dock there. Vesicles become primed for exocytosis by two ATP (adenosine triphosphate)-dependent steps, dissociation of the SNAREs by NSF (NEM-sensitive factor) and remodeling of the phospholipids by the PEPs (priming of exocytosis proteins). Influx of Ca<sup>2+</sup> from nearby Ca<sup>2+</sup> channels triggers fusion of the vesicles with the plasma membrane and exocytosis of the vesicle contents. [Redrawn from Fig. 4e in A. Banerjee, V. A. Barry, B. R. GasGupta, and T. F. J. Martin (1996) *J. Biol. Chem.* **271**, 20223–20226.]



In addition to correctly phosphorylated phospholipids, the secretory granule must have a functional v-SNARE (15) (see [Exocytosis](#)). Treatment of secretory cells with tetanus [toxin](#), which is a v-SNARE-specific [endopeptidase](#), blocks exocytosis (16). The secretory-granule v-SNARE forms a complex with two t-SNAREs, syntaxin and SNAP-25. This complex must be dissociated by the [ATPase](#), *N*-ethylmaleimide sensitive factor, prior to the fusion event (17).

The third requirement for exocytosis is an increase in cytoplasmic calcium concentration (see [Calcium Signaling](#)). Hormone release is triggered when an external signal triggers the opening of a **calcium channel**, raising the internal calcium level from its usual 0.1  $\mu\text{M}$  to about 5  $\mu\text{M}$  (18). The elevated calcium triggers exocytosis by interacting with a membrane-bound protein on the secretory granule, presumably [synaptotagmin](#) (19).

### 1.3. Exocrine Secretory Granules

Exocrine cells are epithelial cells dedicated to protein secretion from their apical surfaces. They include cells of the exocrine pancreas and the intestine, which secrete digestive enzymes, mammary gland cells secreting [casein](#), and mucin- and saliva-secreting cells. These cells exhibit polarized secretion (that is, secretion through a specialized region of the plasma membrane). In addition, exocrine cells are characterized by the very large size of their secretory granules, often called [zymogen](#) granules. Their size is such that they can be seen readily in the light microscope. Early morphologists noted the disappearance of zymogen granules on stimulation of exocrine cells and the appearance of their content extracellularly. Before the phenomenon of exocytosis was discovered, it was assumed that the extracellular signal caused the “granules” to dissolve in the cytoplasm, allowing the proteins to appear extracellularly by unknown transport mechanisms. This view of exocrine secretion remained in vogue for a remarkably long time, even after exocytosis was discovered.

What is known of zymogen granule formation and exocytosis shows strong parallels with what is found in endocrine cells. The contents of the zymogen granules begin to condense in the *trans*-Golgi network. Sorting signals for targeting proteins into zymogen granules must exist (20), and

aggregation is likely to participate (21). The large exocrine secretory granules grow by the fusion of many smaller immature granules. The granules also have a member of the synaptobrevin family of v-SNAREs that is required for exocytosis (22). Targeting to apical membranes appears to be due to selective expression of the appropriate t-SNARE (syntaxin) on the apical surface (23).

In addition to having polarized secretion, secretion exclusively via the apical membrane into the lumen of the gland or the gut, exocrine cells differ from endocrine cells in showing piggy-back exocytosis. This means that once an exocrine secretory granule has fused with the apical plasma membrane, a second secretory granule can fuse with the membrane of the first secretory granule, and so on. As a result, massive exocytosis leads to a highly convoluted apical plasma membrane.

#### 1.4. Secretory Granules in Cells of Hematopoietic Lineage

Cells of the hematopoietic lineage, [stem cells](#), have a form of regulated exocytosis that differs in some ways from what is found in endocrine and exocrine cells. First, the secretory vesicle resembles a **lysosome** more than a secretory granule. In [cytotoxic T lymphocytes](#), it contains [cathepsins](#) and membrane markers of lysosomes. It can contain small vesicles and so resembles a multivesicular body, which is a lysosomal precursor. Sorting of granzymes to the storage vesicles involves the **mannose phosphate receptor**, which normally sorts lysosomal enzymes (24, 25). Secretion in cytotoxic T lymphocytes is polarized in an unusual way. The Golgi complex and microtubule-organizing center rotate around the nucleus until they face the target cell. Finally, exocytosis from mast cells, neutrophils, and cytotoxic T lymphocytes is readily stimulated by GTPγS. It seems likely that protein secretion from hematopoietic cells should be seen as a modification of the conventional lysosomal pathway and is thus only superficially similar to regulated secretion from endocrine cells.

#### 1.5. Constitutive Secretory Vesicles

The constitutive secretory vesicle was initially defined on the pathway for the externalization of membrane proteins in regulated secretory cells. Exported proteins that lack a sorting domain directing them to the regulated pathway also take the constitutive pathway. In cells such as fibroblasts or **antibody**-secreting lymphocytes, there is no major post-Golgi pool of secreted proteins. Proteins arrive at and are released from the cell surface a few minutes after leaving the Golgi complex, instead of being stored for as long as several days in secretory granules, as is common for the regulated pathways. Cells with no post-Golgi storage capacity are usually considered to have only constitutive secretion.

Because the constitutive secretory vesicles are short-lived, it has been difficult to purify enough of them to characterize them biochemically.

Formation of constitutive secretory vesicles in the *trans*-Golgi complex requires a coating mechanism. The early observation that their formation is blocked by **brefeldin A** (26) is consistent with the recruitment of coats by the small [GTPase](#), ADP-ribosylation factor (27). One potential brefeldin A-sensitive coating molecule is p200 (28).

#### 1.6. Synaptic Vesicles

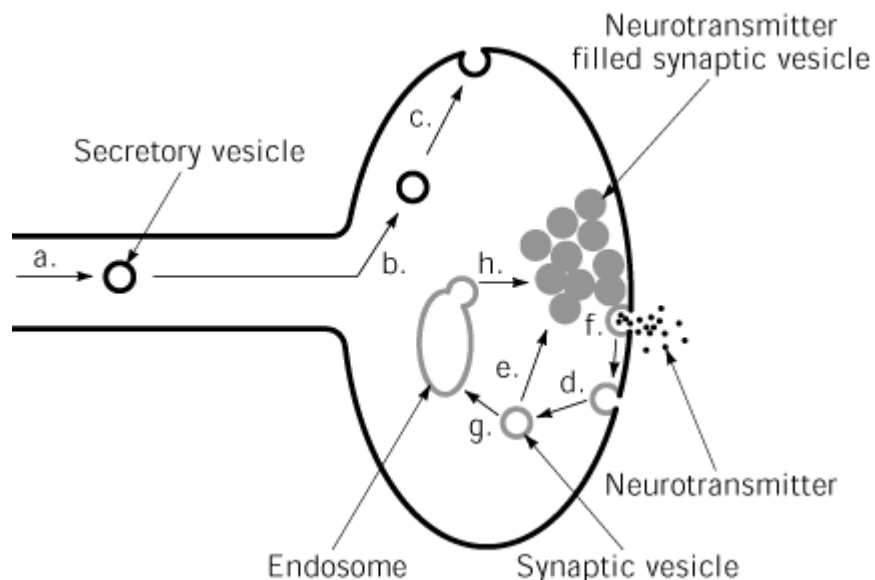
Two cell types dedicated to intercellular signaling are endocrine cells and neurons. Both cell types store their chemical signals in secretory vesicles and release them by exocytosis on stimulation. The crucial difference between the two systems is speed. While an endocrine response can take minutes, communication between neurons takes a millisecond. The properties of the synaptic vesicle make most sense when they are seen as ultra-high-speed secretory vesicles.

#### 1.7. Formation of Synaptic Vesicles

Formation of synaptic vesicles is one of the fastest membrane trafficking events known. Each stimulus causes the release of about 1% of the synaptic vesicles stored in the nerve terminal. Since stimuli can arrive at the terminals at a rate of up to 50 per second, the vesicle pool would be depleted in seconds if there were no efficient way to replenish it (29). Since it can take days for membranes to arrive at the nerve terminal by axonal transport from the cell body, synaptic vesicles

need to be made locally at the nerve terminal. To do this, the membrane proteins of the synaptic vesicle are retrieved after neurotransmitter release and recycled quickly back into fresh synaptic vesicles. The biosynthetic machinery is sufficiently fast to replenish the complete store of nerve terminal synaptic vesicles in 10 to 20 seconds (29). Nerve terminals adapt the process of endocytosis present in all cells. For example, **dynamin**, a large GTPase that polymerizes rapidly into helical structures, is required for the formation of clathrin-coated vesicles in all mammalian cells (30), while AP2 is a heterotetrameric adaptor protein that recruits clathrin coats to plasma membrane proteins awaiting internalization. Mutations in dynamin or in AP2 prevent synaptic vesicle formation (31). Clathrin-coated vesicles are abundant in nerve terminals and contain almost exclusively synaptic vesicle membrane proteins. The synaptic vesicles may form directly from the plasma membrane or from a short-lived endosomal intermediate in the nerve terminal (Fig. 3).

**Figure 3.** Formation of synaptic vesicles. Synaptic vesicle proteins are synthesized in the cell body and (a) transported through the secretory pathway and down the axon of the neuron to the (b) nerve terminal, where they are (c) fused with the plasma membrane. In the nerve terminal, synaptic vesicle proteins undergo multiple rounds of exo/endocytosis. Synaptic vesicles are (d) formed by endocytosis, (e) filled with neurotransmitter, and (f) fused with the plasma membrane to release neurotransmitter in response to stimuli. They may form directly (d–f) from the plasma membrane or (d, g, h) through an endosomal intermediate.



When they have formed, they are quickly filled with neurotransmitter to concentrations as high as 0.4 M. To concentrate the neurotransmitters to this extent, synaptic vesicles have a transporter with 12 transmembrane-spanning regions that makes use of a transmembrane proton electrochemical potential, generated by a **proton pump** in the synaptic vesicle plasma membrane. Several drugs of clinical significance, such as amphetamines and reserpine, interfere with the uptake of neurotransmitter by the transporter (32).

### 1.8. Docking and Fusion of Synaptic Vesicles

The newly formed and filled synaptic vesicle quickly joins the pool of synaptic vesicles in the nerve terminal. These vesicles are not free to diffuse in the nerve terminal but are restricted to a region of the active zone, the site of synaptic vesicle exocytosis. The **actin-binding protein** synapsin is believed to restrict the free movement of synaptic vesicles (33). Of the vesicles in the vicinity of the active site, about 10% are docked at specific docking sites on the nerve-terminal plasma membrane (34). Although protein-containing dense core secretory granules (also called dense core synaptic vesicles) can be found in nerve terminals, they cannot dock at the synaptic vesicle release sites. Docked synaptic vesicles are separated from the plasma membrane by less than the width of a membrane.

Release of neurotransmitter by exocytosis requires the synaptic vesicle v-SNARE, called synaptobrevin or VAMP, and two t-SNAREs on the plasma membrane, syntaxin and synaptosome-associated protein-25 (SNAP-25) (see [Exocytosis](#)). *Drosophila* mutants defective in synaptobrevin or syntaxin do not have normal synaptic transmission ([35](#), [36](#)). Furthermore, it is now known that the blocks in neurotransmitter release caused by bacterial [neurotoxins](#) are due to their endoproteolytic activity on the three SNARE molecules ([37](#)) (see [Exocytosis](#)). Although neurotransmission is blocked by the bacterial neurotoxins, the vesicles can still dock at the plasma membrane, showing that the SNAREs are not essential for docking ([38](#)). The t-SNAREs syntaxin and SNAP-25 are not localized to the active sites and so cannot by themselves specify the site of docking of synaptic vesicles ([38](#)).

The small rab3 GTPases are associated with the resting synaptic vesicle and dissociate from it when the synaptic vesicle fuses with the nerve-terminal plasma membrane ([39](#)). **Knockout** mice defective in rab3A have no apparent defect in docked vesicles, but they are unable to mobilize synaptic vesicles as efficiently from the reserve pool to the docked pool ([33](#)).

### 1.9. Regulation of Synaptic Vesicle Exocytosis

Synaptic vesicles docked at active zones fuse with the plasma membrane within 100  $\mu$ s after an electrical stimulus reaches the nerve terminal. The voltage change caused by the electrical signal opens voltage-sensitive calcium channels in the nerve terminal, and the calcium that enters triggers synaptic vesicle fusion. An important finding was that the cytoplasmic calcium concentration has to reach 50  $\mu$ M before fusion is triggered ([34](#)). Since such calcium concentrations are found only in close proximity to calcium channels, the docking/fusion site and the calcium channels are physically linked to each other. Indeed, it is the calcium channels that define the location of the active zones of vesicle fusion. After a round of exocytosis, the nerve terminal needs to be reset to be capable of a second exocytotic event in about 10 ms. In particular, the calcium concentration at the release site must quickly return to near normal (0.1  $\mu$ M). This is accomplished by rapid diffusion of calcium away from the release site at the mouth of the calcium channel ([40](#)).

The synaptic vesicle uses the [calcium-binding protein](#) synaptotagmin as a sensor to detect changes in cytoplasmic calcium. Mouse knockout mutants lacking synaptotagmin lose almost all calcium-dependent exocytosis from nerve terminals ([41](#)). The fusion machinery is normal in such mice, because they can still release neurotransmitter in the absence of calcium, when stimulated by hyperosmolarity in the bathing fluid or by the toxin of the black widow spider,  $\alpha$ -latrotoxin. Thus synaptotagmin is a regulator of the fusion machinery.

### Bibliography

1. T. L. Burgess and R. B. Kelly (1987) *Annu. Rev. Cell Biol.* **3**, 243–293.
2. S. Urbe, S. A. Tooze, and F. A. Barr (1997) *Biochim. Biophys. Acta* **1358**, 6–22.
3. P. De Camilli and K. Takei (1996) *Neuron* **16**, 481–486.
4. O. Cremona and P. De Camilli (1997) *Curr. Opin. Neurobiol.* **7**, 323–330.
5. H. H. Moore and R. B. Kelly (1986) *Nature* **321**, 443–446.
6. S. Natori and W. B. Huttner (1996) *Proc. Natl Acad. Sci. USA* **93**, 4431–4436.
7. D. R. Cool, E. Normant, F. Shen, H. C. Chen, L. Pannell, Y. Zhang, and Y. P. Loh (1997) *Cell* **88**, 73–83.
8. R. Kuliawat, J. Klumperman, T. Ludwig and P. Arvan (1997) *J. Cell Biol.* **137**, 595–608.
9. A. S. Dittie, N. Hajibagheri, and S. A. Tooze (1996) *J. Cell Biol.* **132**, 523–536.
10. Y. G. Chen and D. Shields (1996) *J. Biol. Chem.* **271**, 5297–5300.
11. Y. Rouille, S. J. Duguay, K. Lund, M. Furuta, Q. Gong, G. Lipkind, A. A. Oliva, Jr., S. J. Chan, and D. F. Steiner (1995) *Front. Neuroendocrinol.* **16**, 322–361.
12. L. Zhang, M. G. Marcu, K. Nau-Staudt, and J. M. Trifaro (1996) *Neuron* **17** 287–296.

13. J. C. Hay and T. F. Martin (1993) *Nature* **366**, 572–575.
14. J. C. Hay, P. L. Fiset, G. H. Jenkins, K. Fukami, T. Takenawa, R. A. Anderson, and T. F. Martin (1995) *Nature* **374**, 173–177.
15. J. E. Rothman (1996) *Protein Science*, **5**, 185–194.
16. R. Regazzi, K. Sadoul, P. Meda, R. B. Kelly, P. A. Halban, and C. B. Wollheim (1996) *EMBO J.* **15**, 6951–6959.
17. A. Banerjee, V. A. Barry, B. R. DasGupta, and T. F. J. Martin (1996) *J. Biol. Chem.* **271**, 20223–20226.
18. R. H. Chow, J. Klingauf, C. Heinemann, R. S. Zucker, and E. Neher (1996) *Neuron* **16**, 369–376.
19. L. A. Elferink, M. R. Peterson, and R. H. Scheller (1993) *Cell* **72**, 153–159.
20. V. Colomer, K. Lal, T. C. Hoops and M. J. Rindler (1994) *EMBO J.* **13**, 3711–3719.
21. F. A. Leblond, G. Viau, J. Laine, and D. Lebel (1993) *Biochem. J.* **291**, 289–296.
22. J. E. Braun, B. A. Fritz, S. M. Wong, and A. W. Lowe (1994) *J. Biol. Chem.* **269**, 5328–5335.
23. H. Y. Gaisano, M. Ghai, P. N. Malkus, L. Sheu, A. Bouquillon, M. K. M. K. Bennett, and W. S. Trimble (1996) *Mol. Biol. Cell* **7**, 2019–2027.
24. G. M. Griffiths (1997) *Sem. Immunol.* **9**, 109–115.
25. G. M. Griffiths and S. Isaacs (1993) *J. Cell Biol.* **120**, 885–896.
26. G. J. Strous, P. van Kerkhof, G van Meer, S. Rijnboutt, and W. Stoorvogel (1993) *J. Biol. Chem.* **268**, 2341–2347.
27. J. P. Simon, I. E. Ivanov, M. Adesnik and D. D. Sabatini (1996) *J. Cell Biol.* **135**, 355–370.
28. N. Narula and J. L. Stow (1995) *Proc. Nat. Acad. Sci. USA* **92**, 2874–2878.
29. T. A. Ryan, S. J. Smith and H. Reuter (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5567–5571.
30. H. Damke, T. Baba, A. M. van der Blik, and S. L. Schmid SL. (1995) *J. Cell Biol.* **131**, 69–80.
31. M. Gonzalez-Gaitan and H. Jackle 1997 *Cell* **88**, 767–776.
32. R. H. Edwards (1993) *Ann. Neurol.* **34**, 638–645.
33. T. W. Rosahl, D. Spillane, M. Missler, J. Herz, D. K. Selig, J. R. Wolff, R. E. Hammer, R. C. Malenka, and T. C. Sudhof (1995) *Nature* **375**, 488–493.
34. R. S. Zucker (1996) *Neuron* **17**, 1049–1055.
35. K. Broadie, A. Prokop, H. J. Bellen, C. J. O'Kane, K. L. Schulze, and S. T. Sweeney (1995) *Neuron* **15**, 663–673.
36. K. L. Schulze, K. Broadie, M. S. Perin, and H. J. Bellen (1995) *Cell* **80**, 311–320.
37. C. Montecucco and G. Schiavo (1995) *Quart. Rev. Biophys.* **28**, 423–472.
38. K. Broadie, A. Prokop, H. J. Bellen, C. J. O'Kane, K. L. Schulze, and S. T. Sweeney (1995) *Neuron* **15**, 663–673.
39. G. Fischer von Mollard, B. Stahl, A. Khokhlatchev, T. C. Sudhof, and R. Jahn (1994) *J. Biol. Chem.* **269**, 10971–10974.
40. R. Llinas, R. M. Sugimori, and R. B. Silver (1995) *Neuropharmacology* **34**, 1443–1451.
41. M. Geppert, Y. Goda, R. E. Hammer, C. Li, T. W. Rosahl, C. F. Stevens, and T. C. Sudhof (1994) *Cell* **79**, 717–727.

## Sedimentation Coefficient (*S*-Value)

During [centrifugation](#), the rate of movement, or the velocity of sedimentation, of a macromolecule through a solution under the influence of a centrifugal field ( $w^2s$ ) can be characterized by its sedimentation coefficient, or *s*-value. As a macromolecule sediments, it moves at the velocity

$$v = dr_b/dt = r\omega^2s \quad (1)$$

where  $r$  is the radial position of the boundary at time  $t$  and  $w$  is the angular velocity of the centrifuge rotor in radians/sec. Specifically, the *s*-value is defined as the velocity of movement divided by the centrifugal field strength  $r w^2$ . Then the *s*-value is given by

$$s = v/r\omega^2 \quad (2)$$

The *s*-value can be calculated graphically or by computer (see [Sedimentation Velocity Centrifugation](#)). When characterizing macromolecular sedimentation, it is important to measure the value of *s* from different experiments at several different macromolecular concentrations to evaluate the effects of intermolecular interactions ([1](#), [2](#)). In the simplest case of nonideality, where the sedimentation of one molecule interferes with another, the *s*-value decreases with increasing macromolecular concentration.

The *s*-value and the [diffusion](#) coefficient  $D$  are sufficient to determine the **molecular weight** ( $M_{s,D}$ , see [Sedimentation Velocity Centrifugation](#)) of a molecule from the Svedberg equation :

$$M_{s,D} = \frac{s}{D} \frac{RT}{(1 - \bar{v}\rho)} \quad (3)$$

where  $R$  is the gas constant,  $T$  the absolute temperature,  $r$  the density of the solvent, and  $v$  the [partial specific volume](#) of the macromolecule (see [Sedimentation Velocity Centrifugation](#)). Finally, the *s*-value is a function of the shape of the macromolecule and of its [frictional coefficient](#) ([1](#), [2](#)), and it depends on the solvent viscosity. The *s*-value is also sensitive to temperature because the viscosity of [water](#) is strongly temperature-dependent. Thus, macromolecular *s*-values (and also diffusion coefficients) are routinely corrected to water as the solvent and at a standard temperature (20°C), designated  $s_{20,w}$  and  $D_{20,w}$ , respectively.

The *s*-value has the dimensions of time and it is usually expressed in [Svedberg units](#)  $S$  ( $1 S = 1 \times 10^{-13} s$ ).

### Bibliography

1. K. E. Van Holde (1971) *Physical Biochemistry*, Prentice–Hall, Inc., Englewood Cliffs, NJ, pp. 70–121.
2. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part II: Techniques for the Study

## Sedimentation Equilibrium Centrifugation

Sedimentation equilibrium experiments measure the radial distribution of molecules at chemical and hydrodynamic equilibrium under a constant centrifugal force (see [Centrifugation](#)). This equilibrium distribution is determined by the opposing tendencies of the molecules to sediment and to **diffuse** (see [Sedimentation Velocity Centrifugation](#)). Consequently, the technique permits direct measurement of the **molecular weight** of the molecule in solution and of  $M_w$  and  $M_z$ , the weight- and z-average molecular weight distributions (see [Analytical Ultracentrifugation](#)). It is also one of the few direct methods, other than titration **calorimetry** and [equilibrium dialysis](#), for accurately measuring the binding constants between macromolecules.

The behavior of an ideal solute at sedimentation equilibrium (Fig. 1) and infinite dilution obeys an exponential (Eq. (1)) relationship and its linear transform (Eq. (2)):

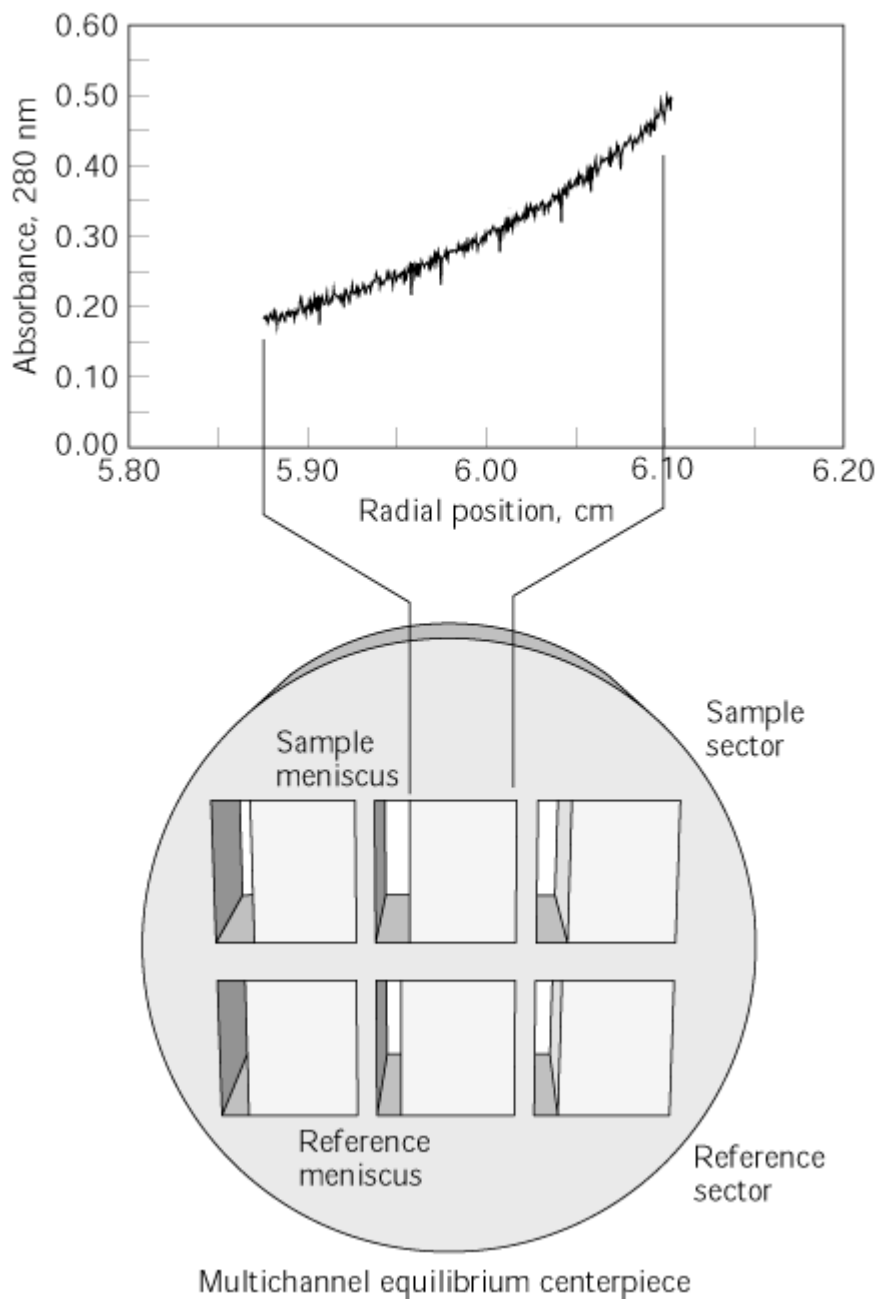
$$C_r = C_m e^{\left\{ \frac{M_w \omega^2 (1 - \bar{v} \rho) (r_i^2 - r_m^2)}{2RT} \right\}} \text{ or } C_r = C_m e^{\{M_w C_0 (r_i^2 - r_m^2)\}} \quad (1)$$

$$\ln C_r = \frac{M_w \omega^2 (1 - \bar{v} \rho) (r_i^2 - r_m^2)}{2RT} + \ln C_m \quad (2)$$

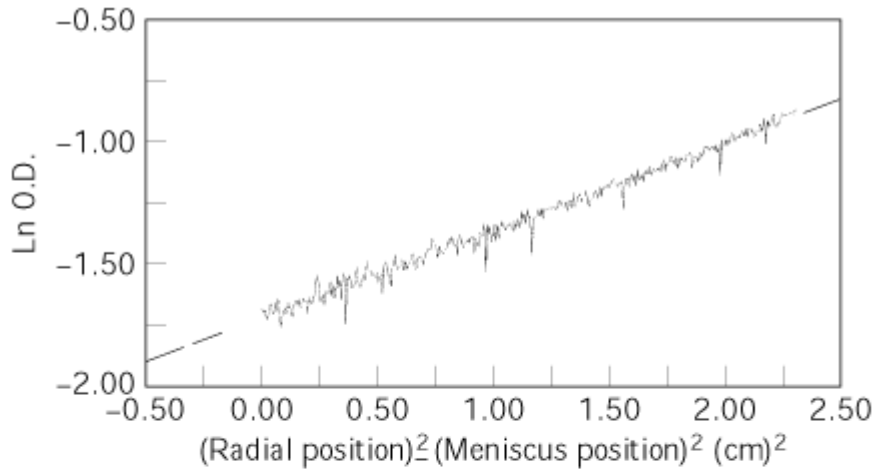
where  $C_r$  and  $C_m$  are, respectively, the concentrations of the solute at radius  $r_i$  and at the meniscus  $r_m$ ,  $M_w$  is the molecular weight of the solute and  $\bar{v}$  its [partial specific volume](#),  $\omega$  is the angular velocity of the rotor,  $\rho$  is the density of the solvent,  $R$  is the gas constant,  $T$  the absolute temperature, and the  $C_0$  term is the combined set of equation constants [ $\omega^2(1 - \bar{v} \rho)/2RT$ ]. In the linear transform of Eq. 2 (Fig. 2), the apparent solute weight-average molecular weight is defined by the slope of the plot of  $\ln C_r$  versus  $(r_i^2 - r_m^2)$ . Deviations from linearity are typically caused by solute paucidispersity, polydispersity, and/or nonideality (see [Analytical Ultracentrifugation](#)). For most dilute samples of otherwise homogeneous proteins, polydispersity can arise from self-association of solute monomers into larger oligomeric forms or from the dissociation to monomers from a normally [oligomeric protein](#).

**Figure 1.** Sedimentation equilibrium profile scanned by absorption optics. The raw data are plotted as sample absorbance (sample solution minus reference solution) versus radial distance from the center of rotor rotation. The figure inset depicts the sample and reference channels of the middle of a three-channel centrifuge centerpiece. The sample and reference solution menisci are indicated.

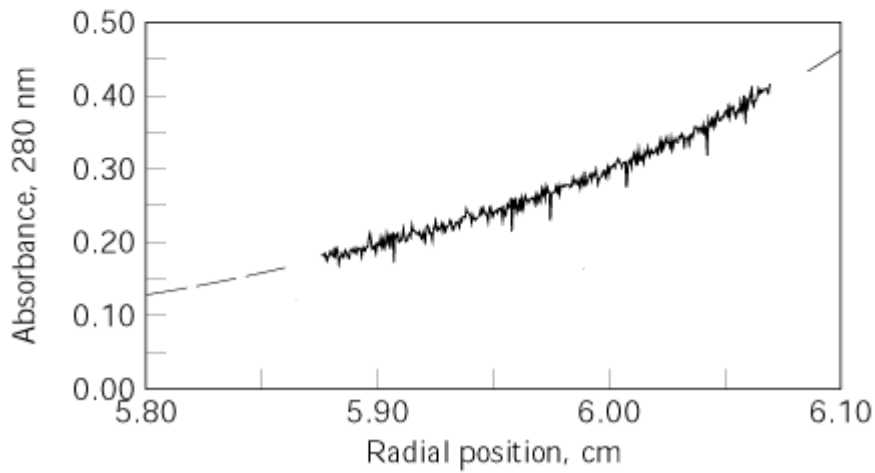




**Figure 2.** Analysis of sedimentation equilibrium data. In the upper panel, the data for a single sedimenting monodisperse protein are plotted according to Eq. (2) and fitted to a straight line (solid line). The slope of the line is directly proportional to the solute weight-average molecular weight  $M_w$ . In the lower panel, the same data are plotted according to Eq. (1), fitting them to a single exponential (solid line).



(a)



(b)

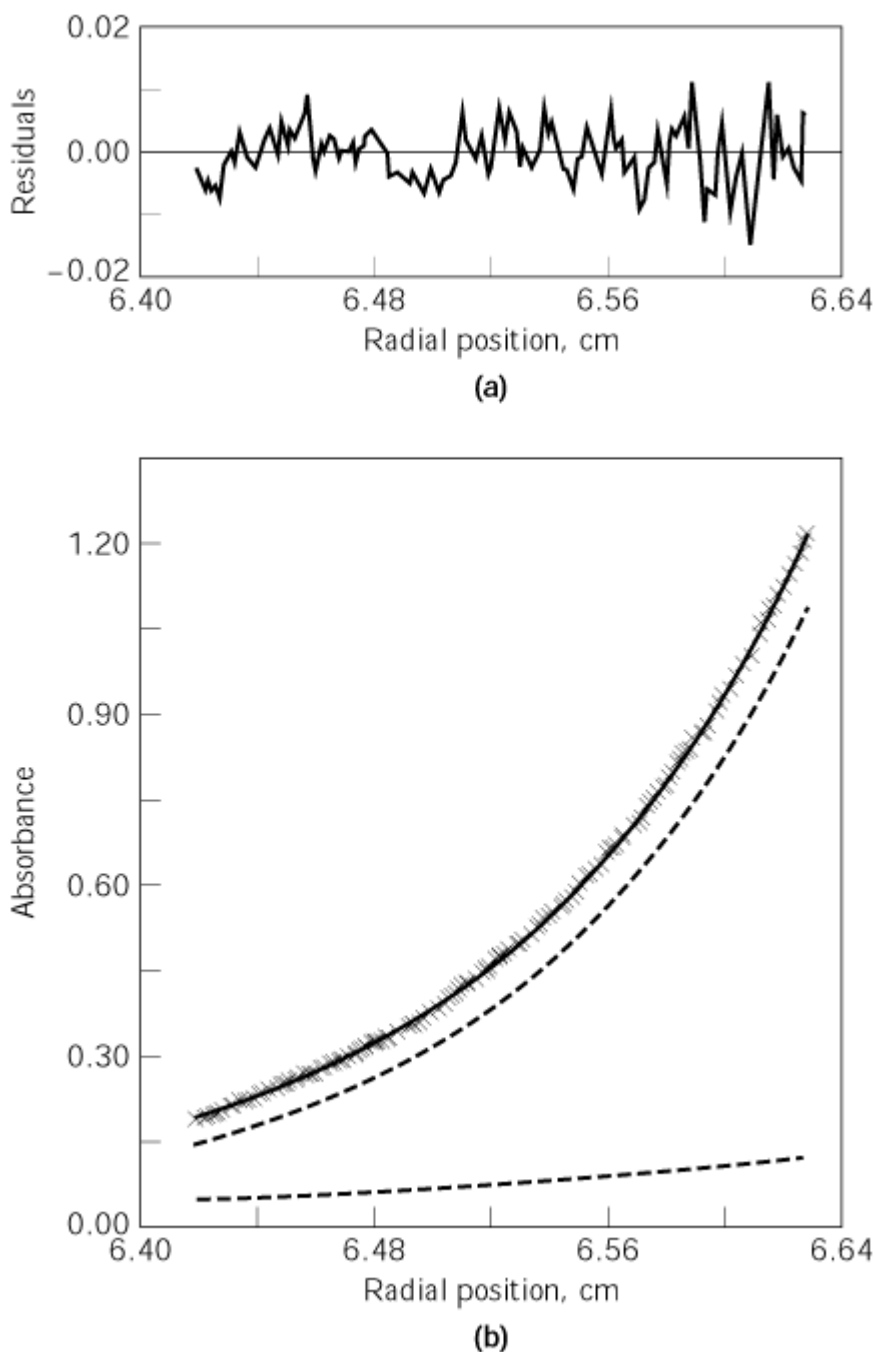
The distribution of each species in a polydisperse solution during centrifugation is described by an equation of the form of Eq. (1). Rather than dealing with logarithmic transformations of data in Eq. (2), the preferred analysis involves using best-fit exponential sums, as in Eq. (3) (1, 2). For example, at equilibrium, the distribution of total protein concentration between a monomer and a dimer as a function of radial position is described as follows:

$$P_{r^2} = A_{11}e(C_0 \cdot \overline{M}_w \cdot r^2) + A_{12}e(2 \cdot C_0 \cdot \overline{M}_w \cdot r^2) \quad (3)$$

where  $r^2$  is the square of the radial position relative to the sample meniscus and ( $P_{r^2}$ ) is the absorbance at that radial position. The monomer molecular weight ( $M_w$ ) for the two-term fit, along with the preexponential constants, are allowed to change during data fitting. The fitting routine iteratively adjusts the values of these coefficients to minimize the difference between the observed data and values calculated from the analytical expression. The number of exponential terms employed in the fitting routines is determined by analysis of both goodness-of-fit parameters (usually the chi-square value) and visual examination of the residual differences between the fitted sums and the actual data set. Such a fit is illustrated in Fig. 3. Thus, this analysis permits evaluating the molecular weights and the radial distributions of the components in solution. The radial distributions of the individual species present can be used to calculate binding constants directly for

the association of the solution components (3). It can be thought of as combining mass transport (under the influence of the centrifugal field) with [equilibrium dialysis](#) (3).

**Figure 3.** Fit of two exponential terms for a monomer and a dimer to sedimentation equilibrium data for the association of the subunits of **cytomegalovirus proteinase** (Ref. 4, with permission). Lower panel: raw data (x), fitted exponential sums of monomer (---, lower shallow dashed curve) and dimer (---, upper steep dashed curve). The thick black steep curve passing through the raw data points (x) is the calculated sum of monomer and dimer at each radial position. Upper panel: residual differences between the fitted sum of the two exponential terms and the actual data.



Unlike indirect techniques that correlate changes in solute spectroscopic properties (such as **absorbance**, **fluorescence** or [circular dichroism](#) ) upon binding, which are often assumed to represent two-state events (free versus bound), sedimentation equilibrium analysis requires no such assumptions. In a carefully performed experiment, both completely bound and free states can exist,

literally side-by-side, in the same centrifuge cell. Along with such techniques as isothermal titration **calorimetry**, where the actual thermochemical energetics of ligand binding are measured, sedimentation equilibrium analysis is one of the few direct techniques capable of such measurements.

### Bibliography

1. S. J. Edelstein, M. J. Rehmar, J. S. Olson, and Q. H. Gibson (1970) *J. Biol. Chem.* **245**, 4372–4381.
2. R. H. Crepeau, C. P. Hensley, and S. J. Edelstein (1974) *Biochemistry* **13**, 4860–4865.
3. I. Z. Steinberg and H. K. Schachman (1966) *Biochemistry* **12**, 3728–3747.
4. S. W. Snyder, R. P. Edalji, F. G. Lindh, K. A. Walter, L. Solomon, S. Pratt, K. Steffy and T. F. Holzman (1996) *J. Protein Chem.* **15**, 763–774.

### Suggestion for Further Reading

5. H. Fugita (1962) *Mathematical Theory of Sedimentation Analysis*, Academic Press, New York. The most elegant, definitive, and detailed monograph written to date on the mathematics behind transport processes as applied to sedimentation velocity and equilibrium centrifugation analyses. Unfortunately, it is out of print and very difficult to find.

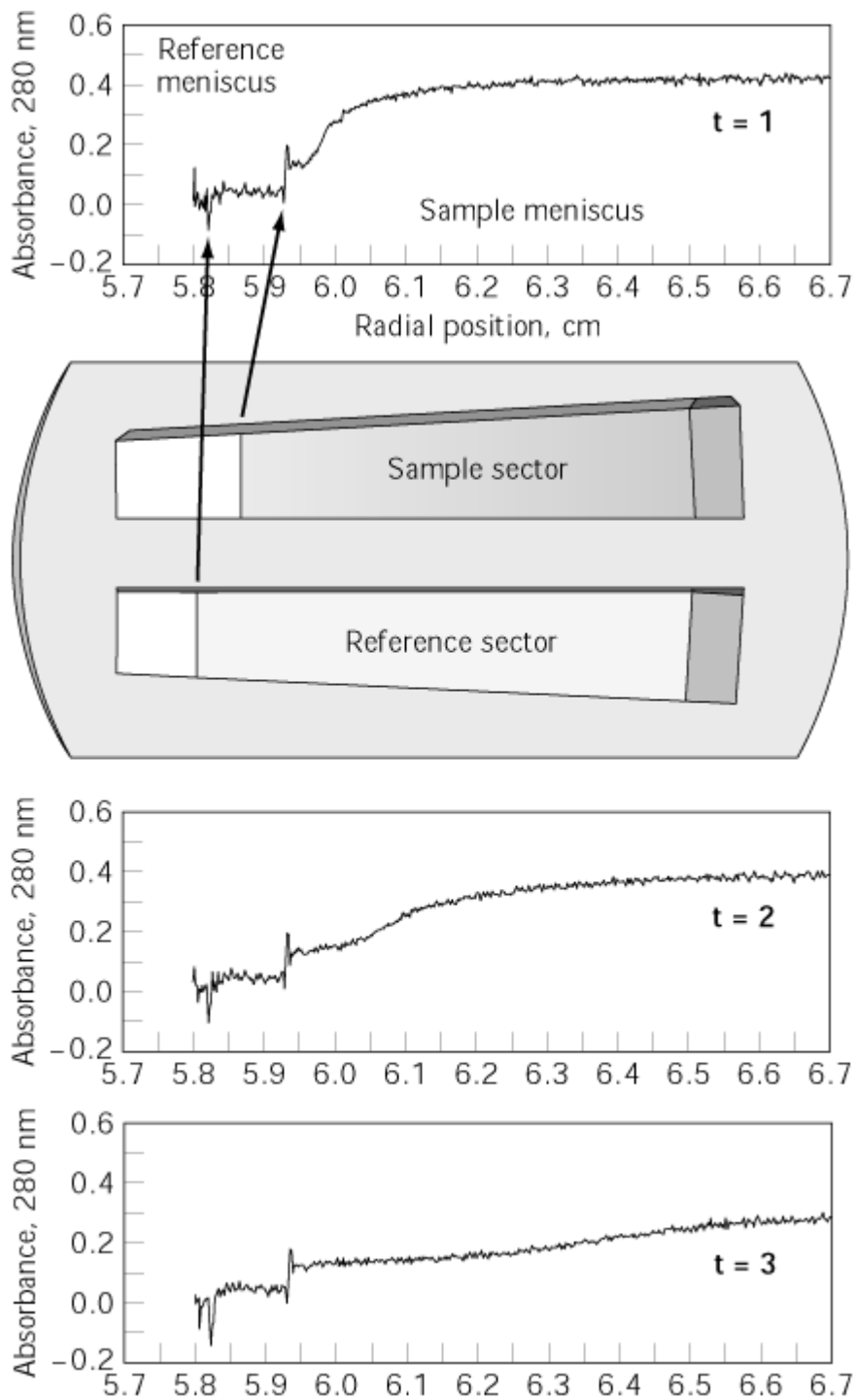
## Sedimentation Velocity Centrifugation

Sedimentation velocity centrifugation experiments measure the rate of movement of molecules through solution under centrifugal force and permit determining the [sedimentation coefficient](#) of a single macromolecule or the individual time-average distributions in a solution containing multiple sedimenting components. In addition to [dynamic light](#) scattering analyses, sedimentation velocity experiments are also useful for estimating the [diffusion](#) coefficients and **frictional coefficients** of the macromolecules and their size in terms of the radii of equivalent spherical molecules. From a single, short, sedimentation velocity run, it is possible to measure both the sedimentation and diffusion coefficients and to combine them to calculate  $M_{s,D}$ , the sedimentation velocity-derived **molecular weight**. Finally, such information can be used with results from other physical measurements to evaluate the shape and [hydration](#) of a macromolecule in solution ([1](#), [2](#)).

### 1. Simple Graphical Analytical Methods

A sedimentation velocity experiment is illustrated in [Figure 1](#). A sample of macromolecule is placed in one side of a centrifuge cell with only the buffer in the other side. It is centrifuged at a rate sufficient to sediment the macromolecule, depleting it at the meniscus and accumulating it at the bottom of the cell. At varying times, the sample cell is scanned for its absorbance at a wavelength where the macromolecule absorbs, giving the profile of the macromolecular concentration versus radial position. A series of profiles at different times represents the moving boundary of molecules sedimenting from the meniscus through the solution under the influence of the centrifugal field.

**Figure 1.** Sedimentation velocity experiment profiles at three different times ( $t = 1$  to 3) in a centrifugation run scanned by absorption optics. The sample and reference solution menisci are indicated. The figure inset depicts the sample and reference channels of the centrifuge centerpiece.



The basic principle of sedimentation velocity is straightforward. The velocity of movement of a sedimenting boundary is defined by

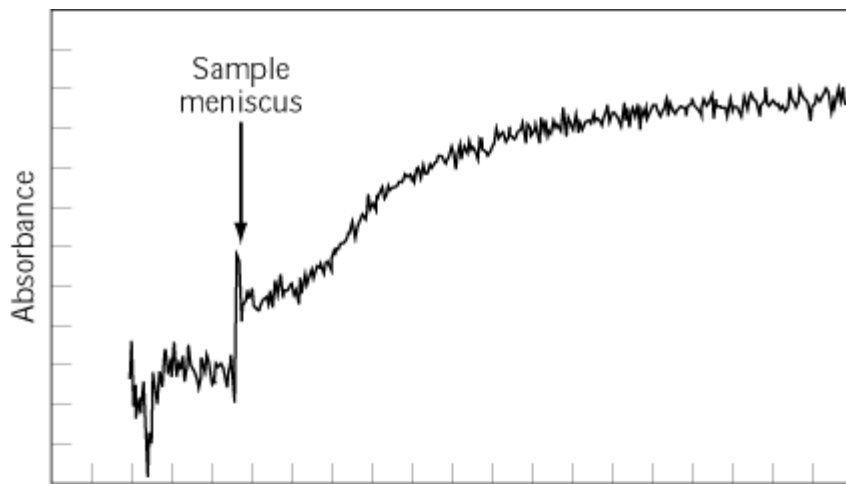
$$v = dr_b/dt = r\omega^2 s \tag{1}$$

or in its integrated form

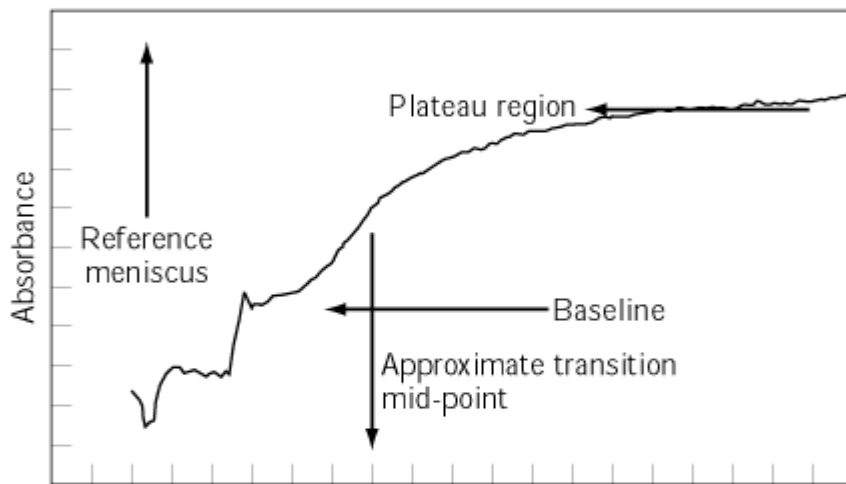
$$\ln \left( \frac{r_{b(t)}}{r_{b(t_0)}} \right) = \omega^2 s (t - t_0) \quad (2)$$

where  $r_{b(t_0)}$  is the position of the sedimenting boundary at time zero,  $r_{b(t)}$  is the position of the sedimenting boundary at the time of measurement,  $\omega$  is the angular velocity of the rotor in radians/sec, and  $s$  is the [sedimentation coefficient](#). The sedimentation coefficient can be determined from Eq. (2) by a plot of  $\ln [r_{b(t)}/r_{b(t_0)}]$  versus  $(t-t_0)$ . There are typically three graphical/computer methods to choose the position of the sedimenting boundary (Figure 2). The first is simply to estimate or calculate the  $r_{b(t)}$  value at the midpoint of the sedimenting boundary based on the (height) distance between the plateau region and the absorbance baseline of the sedimenting sample. The second is to use a first-derivative transformation of the data and to choose the  $r_{b(t)}$  value at the peak position of the first derivative of each scan. The third and most accurate method for single-scan analysis using Eq. (2) is calculating the second-moment  $r_{b(t)}$  boundary position (3).

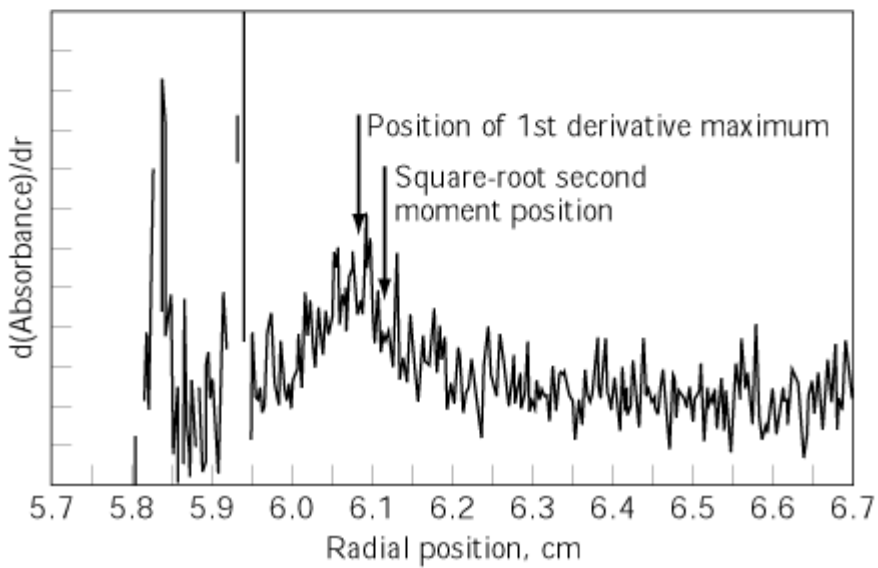
**Figure 2.** Three different ways of analyzing sedimentation velocity data to identify the position of the sedimenting boundary. Raw data scans (upper panel) or smoothed data scans (middle panel) are evaluated at various times for the approximate values of  $r_{b(t)}$  at (1) the transition midpoint (middle panel), (2) first-derivative maximum (lower panel), or (3) second-moment boundary position (lower panel) and plotted according to Eq. (2).



(a)



(b)



(c)

## 2. Computer-Based Analysis Methods

Until recently, graphical linear-fit analyses employing sedimentation velocity Eq. (2) have, been the most widely applied methods for analyzing sedimentation velocity data. Readily availability modern software for computerized graphical analysis has led, however, to the easy application of sophisticated hydrodynamic analytical methods to sedimentation velocity data to account for both the sedimentation and the diffusion of the macromolecule. In this way, it is possible to obtain values for both the sedimentation coefficient  $s$  and the translational diffusion coefficient  $D$ . For example, now two powerful analytical methods can be routinely applied to velocity data with only a modest investment of time and mathematical expertise.

### 2.1. Single-Scan Fitting to the Lamm Equation

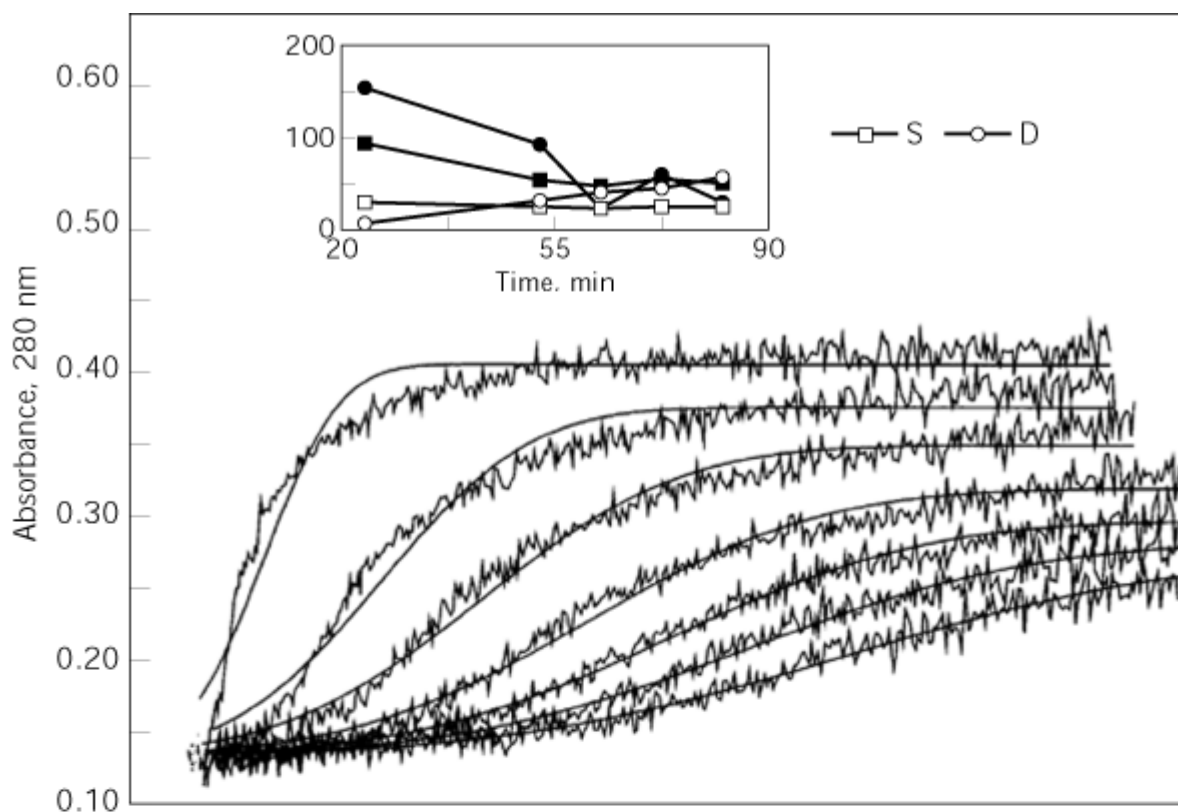
One method now in common use is based on forms of Faxén's two-component solution to the *Lamm equation* (4). The useful expressions are as follows:

$$c = [(c_0 e^{-\tau})/2] \left( 1 - \Phi \left\{ \frac{1 - (xe^{-\tau})^{1/2}}{[\varepsilon(1 - e^{-\tau})]^{1/2}} \right\} \right) \quad (3)$$

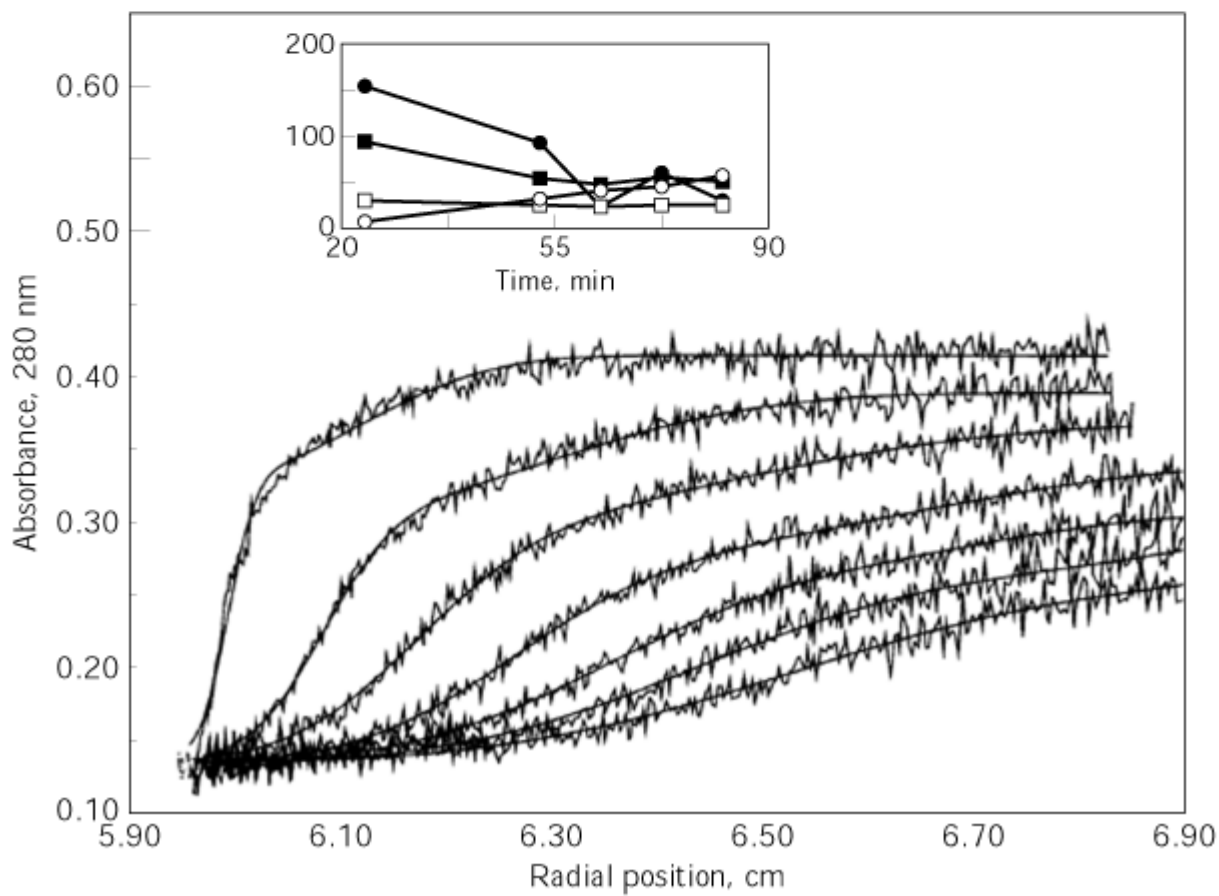
where  $x = (r/r_0)^2$ ,  $t = 2sw^2t$ ,  $\varepsilon = 2D/s(wr_0)^2$ ,  $F$  is the error function of the term enclosed in brackets, and  $s$  and  $D$  are, respectively, the [sedimentation coefficient](#) and the [diffusion](#) coefficient. Under suitable conditions (4), Eq. (3) can approximate the concentration profile of the sedimenting solute boundary closely (i.e.,  $c$  versus  $r$ , as illustrated in Fig. 3). Recently, Behlke and Ristau (5) have described a very useful extension of this method, and Demeler et al. (6) demonstrated its value in modeling the sedimentation behavior of  $n$ -component samples. A short review of recent advances in this approach is provided by Laue (7).

**Figure 3.** Sedimentation velocity analysis of a 40-residue, b-amyloid peptide (Ref. 9 with permission). The absorbance profiles of the sedimenting peptide are given by the original curves after seven different times of centrifugation. In (a), the solid smooth lines indicate a single-component fit of the Lamm Eq. (3) to the data. The inset shows the values of  $s$  and  $D$  obtained from the scans at different times. In (b), the smooth solid lines represent a two-component fit (with two sedimenting species) of the Lamm Eq. (3) to the sample data. The inset represents the  $s$  and  $D$  values for the two different components at the various scan times. The sedimentation coefficients are in units of [Svedbergs](#), and the diffusion coefficients are in units of Ficks.





(a)

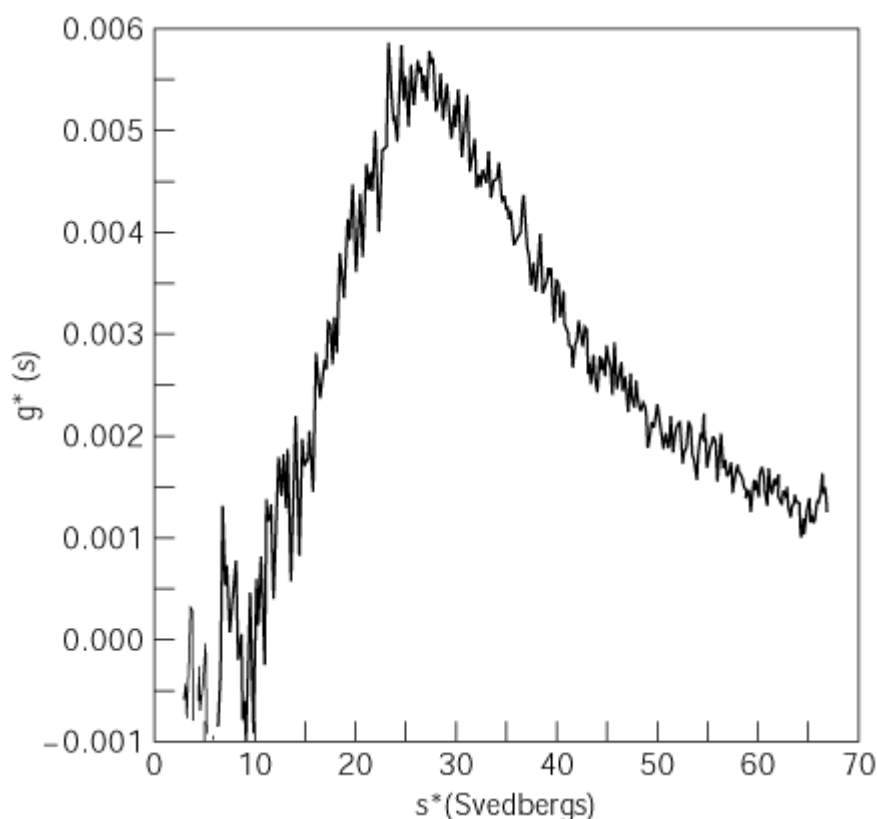


(b)

## 2.2. Sedimentation Coefficient Distribution Analysis

A second modern approach, termed  $g^*(s)$  analysis, was developed and has recently been reviewed by Stafford (8). Rather than analyzing a single scan, it uses a set of scans in a centrifugal run to compute the distribution of sedimentation coefficients of the various components of the sedimenting solute. This approach is extremely useful in analyzing otherwise poorly defined multicomponent samples. Because of this, it has found wide application in studying the size distributions of “real-world” samples of aggregated peptides (Fig. 4) and in defining the oligomeric state of recombinant proteins (9).

**Figure 4.** The apparent distribution of sedimentation coefficients for the data of Fig. 3 (Ref. 9 with permission). Whereas the data of Figure 3 show two sedimenting components of about 28S and 64S, the  $g^*(s)$  analysis shows that the sample of b-peptide is better represented as a distribution of species ranging from about 10S to 70S. The amyloid b-peptide and its aggregated forms may play a role in familial Alzheimer's disease.



## 3. The Sphere-Equivalent Radius

As indicated above, fitting to the Lamm equation yields the values of both  $s$  and  $D$ . One way of converting these sedimentation/diffusion coefficient values into more familiar terms is to employ the relationship between the diffusion coefficient and the **frictional coefficient** of the sedimenting particle (see Table 2 of [Centrifugation](#)). Although the details of the relationships between the frictional coefficient and molecule shape are beyond the scope of this discussion (see references [2](#), [10](#), [11](#) and [Diffusion](#) for more thorough treatments), a simple approximation is possible and is often very useful as a molecular “yardstick.” The diffusion coefficient of a sphere is directly proportional to temperature and inversely proportional to the Stokes’ law frictional coefficient. Thus, knowing the viscosity of the solvent, it is possible to compute the radius of the equivalent sphere  $r_{\text{sphere}}$  of the solute particle under study. This is typically done using the diffusion coefficient obtained from fitting to the Lamm equation. Alternatively,  $D$  can be estimated using the simple ratio of area/height

of the first derivative of the sedimenting boundary (Fig. 1), as described by Fugita (4):

$$(A/H)^2 = D(2\pi/s\omega^2)[\exp(2s\omega^2t) - 1] \quad (4)$$

Equation (2) provides the sedimentation coefficient  $s$ , while Eq. (4) gives the value of  $D$ . Using the relationship between the frictional coefficient and the diffusion coefficient and combining all of the various constants, it can be shown that  $r_{sphere}$  in nanometers is equal to  $(7.323 \times 10^{-9}) \times (T \text{ in kelvin} / D \text{ in cm}^2/\text{sec})$  for a particle in a solvent with the viscosity of water. Thus a diffusion coefficient of 1 Fick, or  $1.0 \times 10^{-7} \text{ cm}^2/\text{sec}$ , represents an  $r_{sphere}$  of 21.5 nm at 20°C in water. Using this relationship with some typical proteins, it is possible to estimate their diameters from just their measured diffusion coefficients. For example, hen egg-white [lysozyme](#) has a diffusion coefficient of  $11.2 \times 10^{-7} \text{ cm}^2/\text{sec}$  (10) and an  $r_{sphere}$  of ~1.9 nm at 20°C in water, and bovine [serum albumin](#) has a diffusion coefficient of  $6.97 \times 10^{-7} \text{ cm}^2/\text{sec}$  (10) and an  $r_{sphere}$  of ~3.1 nm at 25°C in water. For comparison, the **X-ray crystallographic** structure of hen lysozyme has rough dimensions of 3.0×3.0×4.5 nm.

#### 4. $M_{s,D}$ , The Molecular Weight Derived from Sedimentation Velocity Data

The molecular weights of solutes are measured most accurately by [sedimentation equilibrium centrifugation](#), but the sedimentation and diffusion coefficients obtained from sedimentation velocity experiments can be combined to give an approximate molecular weight (2, 10, 11) from the expression

$$M_{s,D} = \frac{s}{D} \frac{RT}{(1 - \bar{v}\rho)} \quad (5)$$

where  $R$  is the gas constant,  $T$  the absolute temperature,  $r$  the density of the solvent, and  $v$  the [partial specific volume](#) of the macromolecule. It is most useful to calculate the  $M_{s,D}$  values during the sedimentation velocity phase of an analytical ultracentrifugation run profile. This is especially true for runs using “long-column” centerpieces (Fig. 1) and extended equilibrium run times (>48 hours). However, it is also sometimes useful for data collected with multichannel equilibrium centerpieces, even though the time to equilibrium is shorter because of the reduced column height in the centerpieces. Because the value of  $M_{s,D}$  can be calculated quickly (within a few minutes), it is often useful in choosing the final speed for the sedimentation equilibrium phase of a centrifugation run after the high-speed, velocity phase of the run (see [Analytical Ultracentrifugation](#)).

#### Bibliography

1. C. Tanford (1961) *Physical Chemistry of Macromolecules*, Wiley, New York, pp. 317–456.
2. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part II: Techniques for the Study of Biological Structure and Function, W. H. Freeman, San Francisco, CA.
3. R. J. Goldberg (1953) *J. Chem. Phys.* **57**, 194–202.
4. H. Fugita (1962) *Mathematical Theory of Sedimentation Analysis*, Academic Press, New York.
5. J. Behlke and O. Ristau (1997) *Biophys. J.* **72**, 428–434.
6. B. Demeler, H. Saber, and J. C. Hansen (1997) *Biophys. J.* **72**, 397–407.
7. T. M. Laue (1997) *Biophys. J.* **72**, 395–396.

8. W. F. Stafford (1994) "Acquisition and interpretation of data for biological and synthetic polymer systems" In *Modern Analytical Ultracentrifugation* (T. M. Schuster and T. M. Laue, eds.), Birkhäuser Boston., pp. 119–137.
9. S. W. Snyder, U. S. Lador, W. S. Wade, G. T. Wang, L. W. Barrett, E. D. Matayoshi, H. J. Huffaker, G. A. Krafft, and T. F. Holzman, (1994) *Biophys. J.* **67**, 1216–1228.
10. K. E. Van Holde (1971) *Physical Biochemistry*, Prentice–Hall, Englewood Cliffs, NJ, pp. 70–121.
11. D. Eisenberg and D. Crothers (1979) *Physical Chemistry with Applications to the Life Sciences*, Benjamin–Cummings, P Menlo Park, CA, pp 701–745.

## Selection, Genetic

Selection is the process by which organisms with certain traits gain reproductive advantage over organisms of the same population without those traits. Provided that the traits are heritable, selection increases the prevalence of the selected traits in the population, ie, the population becomes adapted to the selective conditions. The source of heritable variation is [mutation](#), so selection is the differential **reproduction** of mutant individuals. Selection has somewhat different connotations depending on whether it is used in a genetic sense or in an **evolutionary** sense.

In genetic usage, selection almost always means allowing only individuals with a given trait, usually mutants, to reproduce. Thus, selection allows the researcher to isolate an extremely rare mutant so long as the mutant has a selectable **phenotype**. A common example is [antibiotic resistance](#) in **bacteria**, which is caused by a mutation or the acquisition of a **plasmid** or [transposable element](#) carrying drug resistance. In the presence of the antibiotic, the resistant bacteria survive and give rise to **clones** of resistant progeny, whereas sensitive bacteria are killed. Likewise, with the proper selection, rare cells or organisms with other desirable characteristics can be harvested from a large population. Typical genetic selections are for amino acid **prototrophies**, for the use of specific carbon sources or for suppression of a [conditional lethal mutation](#). The unwanted individuals do not have to die, but only fail to reproduce.

Selection is distinguished from two other genetic manipulations, *enrichment* and *screening*. Enrichment is also a selective process, but unwanted individuals do not fail to reproduce (and in this way it is more like natural selection, see later). Typically, several rounds of enrichment for growth on a poorly used substrate, for example, are needed to obtain a reasonably pure culture of mutant cells better able to use the substrate. Screens, on the other hand, are nonselective but allow the researcher to identify the desired cell or organism. For example, a typical step in cloning a gene is to insert it into a plasmid so that another gene, for which an easy assay exists, is disrupted. Genes are commonly cloned into the *lacZ* gene, which encodes **beta-galactosidase**. Cells with active b-galactosidase turn blue on medium containing X-gal (5-bromo-4-chloro-3-indolyl-b-D-galactopyranoside), so the progeny of cells receiving the disrupted *lacZ* gene containing the cloned gene are white and can be isolated for further analysis (see [Operons](#)).

Selections are often devised to find cells that have gained a trait, whereas screens must usually be used to identify cells that have lost a trait. But there are notable exceptions. For example, mutations that inactivate nonessential genes are selected with toxic analogues of the normal enzyme's substrate. And part of the intellectual challenge of genetics is to devise selection procedures to find desired genes or proteins. Recently, the yeast [two-hybrid system](#) for identifying interacting proteins has been used to select for inactivating mutations by making a successful interaction toxic to the cell ([1](#)).

In **evolutionary** usage, selection is the process by which adaptive changes take place in populations over time. [Natural selection](#) increases the mean fitness of a population by enriching it for more fit individuals and decreasing the prevalence of less fit individuals, thus producing changes in the frequencies of the genetic alleles associated with variation in fitness. There are two ways in which fitness is measured: (1) the number of offspring an individual produces; and (2) the change in the frequency of a certain allele with time. Because selection operates on the whole organism, the first of these has more evolutionary meaning. However, an organism's phenotype is the totality of its traits, each of which may or may not contribute to its fitness. Selection operates to favor particular traits encoded by particular genetic alleles. While this is what is usually meant by evolution, it is not necessarily true that all of the traits of a population are adaptive. Neutral genetic alleles can become prevalent in a population by chance, a process known as [genetic drift](#).

Although fitness can be measured absolutely, a more meaningful measure is relative fitness, i.e., the fitness of individuals with a particular **genotype** relative to individuals with another genotype. The strength of the selective pressure on individuals with a given genotype, the selection coefficient, is one minus their relative fitness. Note that the selective pressure can be positive or negative, but because the fittest genotype is usually given the value of 1, the selection coefficient takes values of 0 to 1. Negative selection eliminates variants from the population, whereas positive selection favors them. Fitness can be defined only for a given set of conditions. Today's fittest genotype may be at a disadvantage if conditions change.

#### Bibliography

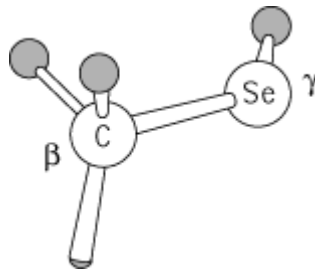
1. M. Vidal, R. K. Brachmann, A. Fattaey, E. Harlow, and J. D. Boeke (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10315–10320.

#### Suggestions for Further Reading

2. R. A. LaRossa (1996) "Mutant selections linking physiology, inhibitors, and genotypes" In *Escherichia coli and Salmonella; cellular and molecular biology*, 2nd ed. (F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology Press, Washington, DC, pp. 2527–2587.
3. A. Adams, D. E. Gottschling, C. Kaiser, and T. Sterns (1998) *Methods in Yeast Genetics: A Laboratory Course Manual*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
4. A. Griffiths, J. H. Miller, D. Suzuki, R. Lewontin, and W. Gelbart (1996) *An Introduction to Genetic Analysis*, 6th ed., W. H. Freeman, New York.
5. M. Kimura (1983) *The Neutral Theory of Molecular Evolution*, Cambridge University Press, Cambridge, UK.
6. J. Maynard Smith. (1998) *Evolutionary Genetics*, 2nd ed., Oxford University Press, Oxford, UK.

#### Selenocysteine

The side chain of the [amino acid](#) *selenocysteine* differs from [cysteine](#) only in having a selenium atom in place of the sulfur of Cys:



This amino acid residue is generally abbreviated as Sec, and its selenol group is an essential component in the [active sites](#) of a few important [enzymes](#) in both prokaryotes and eukaryotes, such as formate dehydrogenase and glutathione peroxidase. The Sec residue occurs at a specific position in every polypeptide chain of these enzymes. It is inserted during **protein biosynthesis** using the [genetic code](#) of the [messenger RNA](#) and can be considered the 21st amino acid of proteins.

The gene sequences for proteins with Sec residues indicate that the **codon** specifying the selenocysteine residue is UGA, one of the normal termination codons. In the appropriate context, however, this codon is **translated** using a special [transfer RNA](#) specific for selenocysteine. In both **bacteria** and **eukaryotes**, this tRNA is charged initially with serine (see **AMINOACYL TRNA SYNTHETASES**), which is then converted enzymatically to selenocysteine. This conversion is catalyzed by a specific enzyme, using the intermediate selenophosphate, which is produced from selenide and ATP by selenophosphate synthetase, also generating orthophosphate and AMP. Although bacteria produce free selenocysteine, the amino acid is apparently not used in protein biosynthesis.

Whether a UGA codon is translated as Sec or as a termination signal is dependent on other sequences in the mRNA. In prokaryotes, the 40-base sequence immediately downstream from the UGA must be present, and this sequence is believed to adopt a **stem-loop structure**. Analogous structures are believed to be important in eukaryotes, but in this case the sequence responsible is in the 3' untranslated region of the mRNA. In bacteria, a special translation [elongation factor](#) is required; it acts in place of EF-Tu, and transports the selenocysteyl-tRNA<sup>Sec</sup> to the [ribosome](#).

A selenol group normally ionizes with a  $pK_a$  of 5.2, so it would be fully ionized at physiological pH; in contrast, cysteine [thiol groups](#) have considerably higher  $pK_a$  values and are only slightly ionized at pH 7. Consequently, the selenol and thiol groups have different chemical properties, and the enzymes containing selenocysteine residues are much less active if these residues are replaced by cysteine.

#### Suggestions for Further Reading

T. C. Stadtman (1996) Selenocysteine, *Annu. Rev. Biochem.* **65**, 83–100.

S. C. Low and M. J. Barry (1996) Knowing when not to stop: selenocysteine incorporation in eukaryotes, *Trends Biochem. Sci.* **21**, 203–207.

### Selenomethionine

Selenomethionine (Se-methionine) differs from the normal amino acid [methionine](#) only in that it has a selenium atom in place of the normal sulfur atom. Selenomethionine can be incorporated biosynthetically into proteins in place of methionine. In practice, this is most facile using

microorganisms that cannot synthesize methionine and depend on its presence in the growth medium. Selenomethionine is simply added instead.

The presence of selenomethionine residues is very useful in protein [X-ray crystallography](#). The change of S (16 electrons) into Se (34 electrons) is an ideal [isomorphous replacement](#) and can be applied to solve the [phase problem](#) in protein X-ray crystallography. This requires X-ray data for the native protein and for its Se-methionine derivative. However, even more important is the application of this derivative in the multiple wavelength anomalous dispersion technique (see **MAD**) by taking advantage of the relatively strong [anomalous dispersion](#) of the Se atom near its K-absorption edge (0.98 Å) (1).

## Bibliography

1. W. A. Hendrickson, J. R. Horton and D. M. LeMaster (1990) *EMBO J.* **9**, 1665–1672.

## Suggestion for Further Reading

2. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Self-Assembly

### 1. Definition

Self-assembly is the formation of a well-defined, complex structure under equilibrium conditions spontaneously from noncovalent interactions between two or more molecular components. Self-assembly encompasses the self-association of [proteins](#) and **nucleic acids**, such as the formation of different forms of a double-stranded **DNA** double helix from two complementary polydeoxynucleotide strands (see [Reassociation, Nucleic Acids](#)), the association of multiple copies of the same [polypeptide chain](#) (see [Oligomeric Proteins, Quaternary Structure](#)), the formation of complexes comprised of more than one type of polypeptide chain (ie, *heteromultimeric* proteins), and the reconstitution of an infectious **viral** particle, comprised of one or more proteins, plus nucleic acid.

Self-assembly can also be viewed as part of a more general topic of *supramolecular* assembly or self-organization, which encompasses the formation of a wide variety of structures with masses greater than 1 MDa and sizes greater than 10 nm. These include pathogens such as viruses; large **multienzyme complexes** involved in biosynthetic pathways (eg, *pyruvate dehydrogenase*) and in **protein degradation** (eg, the [proteasome](#)); multiprotein complexes that bind to nucleic acids and are active in [DNA replication](#), **RNA biosynthesis** and **protein biosynthesis** (eg, [nucleosomes](#), [polymerases](#), [spliceosomes](#), [ribosomes](#)), transmembrane protein complexes (eg, [acetylcholine receptor](#), respiratory enzyme complexes, [photosynthetic reaction centers](#)); and the protein complexes encasing transport **vesicles** (**clathrin-coated vesicles** and circulating **lipoproteins**).

The foregoing definition of self-assembly implies that all the necessary information is encoded within the covalent structures of the components and hence requires complete reversibility of the assembly process. Reversibility is demonstrated *in vitro* by the reassembly of the completely dissociated structure upon return to the initial solution conditions. Dissociation into individual molecules can be effected in many different ways: (a) increased hydrostatic pressure, (b) addition or removal of specific **ligands**, such as **coenzymes**, substrates, **allosteric** effectors or other **prosthetic** groups, (c) decrease in concentration, (d) alteration of pH, temperature, or ionic strength, and (e) the addition of dissociating and denaturing **chaotropic** agents such as [urea](#), [guanidinium salts](#), and

[detergents](#) (neutral, cationic, or anionic). At present, overexpression in *Escherichia coli* of the individual wild-type or mutagenized subunits makes it possible to elucidate the pathways of self-assembly and the roles in the self-assembly process of various protein **domains** within each subunit.

Although self-assembly of a limited number of supramolecular structures has been demonstrated to occur *in vitro*, it is now known that the assembly of many large, complex supramolecular structures *in vivo* requires various cellular factors and generally involves several steps, some of which are irreversible. Self-organization of intracellular spatial structures encompasses self-assembly and almost certainly occurs under conditions that are distant from the equilibrium conditions under which the latter is studied *in vitro*.

## 2. Historical Perspective

Historically, the concept of self-assembly arose primarily from the work on the reconstitution of **tobacco mosaic virus (TMV)** by H. Fraenkel-Conrat and R. Williams. In 1955 they were able to reassemble its isolated RNA and protein components into a native-like, infectious rod-like particle, consisting of exactly 2130 identical protein subunits (each of 158 amino acid residues) stacked in a helical array around a single-stranded RNA molecule comprising 6390 nucleotides. A subsequent key development was the demonstration in 1961 by C. Anfinsen and colleagues that the enzyme **ribonuclease A** (a single [polypeptide chain](#) of 124 residues) could regain its four [disulfide bonds](#) and three-dimensional structure, plus the resulting enzymatic activity, from the completely unfolded, reduced, and inactive form. This was recognized to indicate that all the necessary information for the reversible folding of a polypeptide chain into a well-defined three-dimensional structure resided within its sequence of amino acids (see **Protein folding**). Subsequently, in 1968, M. Nomura and collaborators accomplished the complete reassembly of the 30S ribosome from a mixture of 30 proteins and the 16S rRNA. This and numerous other studies of multimeric proteins and enzymes provided a firm molecular basis for the concept of self-assembly. The observations of intracellular arrays of linear polymers of cytoplasmic proteins, [microtubules](#), [intermediate filaments](#), and [actin filaments](#) in the early 1970s provided the impetus for investigations of the structure and assembly of supramolecular spatial structures responsible for the maintenance of cell morphology and the related processes of cell movement, division, and [differentiation](#).

## 3. Intermolecular Forces Responsible for Self-Assembly

Although the covalent structure of each component encodes all the information necessary for its three-dimensional structure and its interactions with the other components in the self-assembly process, the intermolecular forces involved in both processes are all noncovalent, being energetically weaker than covalent bonds by more than an order of magnitude (<10kcal/mol vs ~100kcal/mol). These intermolecular forces represent both specific, geometrically directed interactions, like [hydrogen bonds](#) and [salt bridges](#) (which may be regarded as hydrogen-bonded ion pairs), and relatively nonspecific **electrostatic** and [van der Waals interactions](#), including the [hydrophobic effect](#). The resulting energetic stabilization depends on the amount of [accessible surface](#) area buried in the contact between the interacting components and, most critically, on the complementarity in the chemical nature and shapes of the interacting surfaces.

## 4. Role of Subassemblies and Assembly Pathways

Simultaneous encounter and assembly of three or more independent molecules is most unlikely on kinetic grounds, and self-assembly usually occurs by the stepwise interactions of pairs of molecules or subassemblies. In the case of a relatively simple, hypothetical assembly system consisting of only three components A, B, and C, there are three possible pathways for the formation of the ABC assembly via formation of the subassemblies AB, AC, or BC. One of the three pathways may be the only pathway that is used because it is kinetically and energetically most favorable, perhaps through a nucleation event involving the formation of one of the initial binary complexes. Pathways and subassemblies can play an important role in determining the accuracy of the final assembly structure



(1). For example, forming a homooligomer of 1000 subunits via stepwise addition of individual subunits with a probability of being correct of 0.999 for each step (ie, only a 0.1% error rate) is only 0.368 ( $0.999^{1000}$ ). However, if two subassemblies of 500 subunits are formed first and then associate to the final oligomer, the probability of correct assembly rises to 0.605. The two optimum assembly pathways would be from 25 subassemblies, each of 40 subunits, or from 40 subassemblies, each of 25 subunits (1). Furthermore, a 50% probability of assembling the correct 1000-mer structure using the optimum pathways requires that each step in the assembly occurs correctly with a probability of at least 0.99. This value is decreased (ie, there is an increased probability of obtaining the correct structure), with an increase in the number of different subassemblies—for example, to 0.98 with self-association to 10-mers, followed by self-assembly of 10-mers into 100-mers, followed in turn by self-assembly of 100-mers to the correct 1000-mer. Thus, successful assembly of a multisubunit structure can be favored by maximizing the number of subunits per subassembly and the total number of subassemblies. The existence of subassemblies can also increase the speed with which a multisubunit structure can be assembled. If each subunit or subassembly addition in the foregoing example requires 1 s, then the assembly of the 1000-mer by stepwise addition would take 999 s to complete; from 40 subassemblies of 25 subunits each, only 63 s (39+24) would be necessary (1).

In view of these simple considerations, it is reasonable that assembly of large spatial structures occurs via well-defined pathways and involves the formation of specific subassemblies. A. Klug and collaborators worked out the details of TMV assembly from the coat protein and the single-stranded RNA. The initiation event was found to be the interaction with a hairpin loop of the RNA of a self-assembled two-layer disc of 34 protein subunits. This was followed by an elongation phase of the self-assembly process proceeding in both directions through further addition of discs of protein subunits (2).

## 5. Experimental Methods of Investigating Self-Assembled Structures

The detailed structures of small and large proteins, of unassembled viral capsid protein oligomers, and of many capsids themselves have been determined to high resolution by [X-ray crystallography](#). The development of [cryoelectron microscopy](#) and three-dimensional image reconstruction (see [Single Particle Reconstruction](#)) techniques capable of achieving a resolution of 1 nm and better has provided a very powerful method for the investigation of supramolecular assemblies. Cryoelectron microscopy can be used with large particles that are unstable, short-lived or present in quantities too small for crystallization. The two methods complement each other, in that X-ray crystallography can provide high-resolution structures of individual components, and cryoelectron microscopy provides lower resolution structures of the large complexes consisting of these components (3). Furthermore, [immunolectron microscopy](#) can localize [epitopes](#) within protein subunits participating in a supramolecular structure, providing a method to identify specific protein subunits in the structure. With capsids, immunolectron microscopy can be used to assist in the mapping of the movement of protein subunits during capsid maturation.

## 6. Advantages of Assembled Structures

Self-assembly of multisubunit protein complexes provides a number of structural and functional advantages over equally large, single-chain proteins (4, 5).

1. Coding efficiency: Less genetic information is required to build a small protein subunit, with a concomitant reduction in biosynthetic error.

2. Quality control via association and dissociation of subunits: Errors in subunit synthesis are minimized, because mutated subunits can be excluded from the assembly process.

3. Enhancement of function through the association/dissociation equilibria of the subunits: assembly of monomers to the active complex at locations distant from the site of monomer synthesis and the

ability to either transmit or respond to allosteric regulatory signals via intersubunit communication (allostery).

4. The self-assembly and dissociation of multisubunit structures can be driven by small alterations in external conditions (in pH, ion concentration, etc.), thereby providing the possibility of added refinement in their functions.

5. Multimerization can provide specificity for interaction of the complex with another partner in supramolecular structures: many proteins involved in the regulation of gene [transcription](#) form homo- and heterodimers by more than one mechanism, thereby enabling them to bind to the substrate DNA with exquisite specificity.

#### 7. Self-Assembly *In Vivo* is More Complex than Self-Assembly *In Vitro*

Since Anfinsen's seminal experiments, the duration of refolding of some proteins *in vitro* was found to be reduced in the presence of the catalysts [peptidyl prolyl cis–trans isomerase](#) and [protein disulfide isomerase](#). Furthermore, the inability of [insulin](#) and some other protein [hormones](#) to refold following complete unfolding was shown by D. Steiner in 1974 to be due to **proteolytic** cleavage subsequent to the folding of a precursor form. Although substantial information had been accumulated by 1987 on the self-assembly of many homo- and heteromultimeric protein complexes (6), as well as of viruses, the first intimation of the involvement of additional factors in self-assembly was the finding that the plant enzyme, **ribulobisphosphate carboxylase**, required **molecular chaperone** proteins for proper folding and assembly of its two different subunits, L and S, to the final L<sub>8</sub>S<sub>8</sub> form. Although the molecular chaperones primarily inhibit misfolding and misassembly and do not direct assembly per se, it has become apparent that the narrow definition of self-assembly provided at the outset has very limited applicability to the self-organization of the many large supramolecular assemblies *in vivo*. It needs to be emphasized that the self-organization of large cellular structures and their maintenance requires a dynamic interplay between assembly and regulatory processes; not only are their pathways multistep processes, but they also occur under nonequilibrium conditions.

#### 8. Control of Self-Assembly

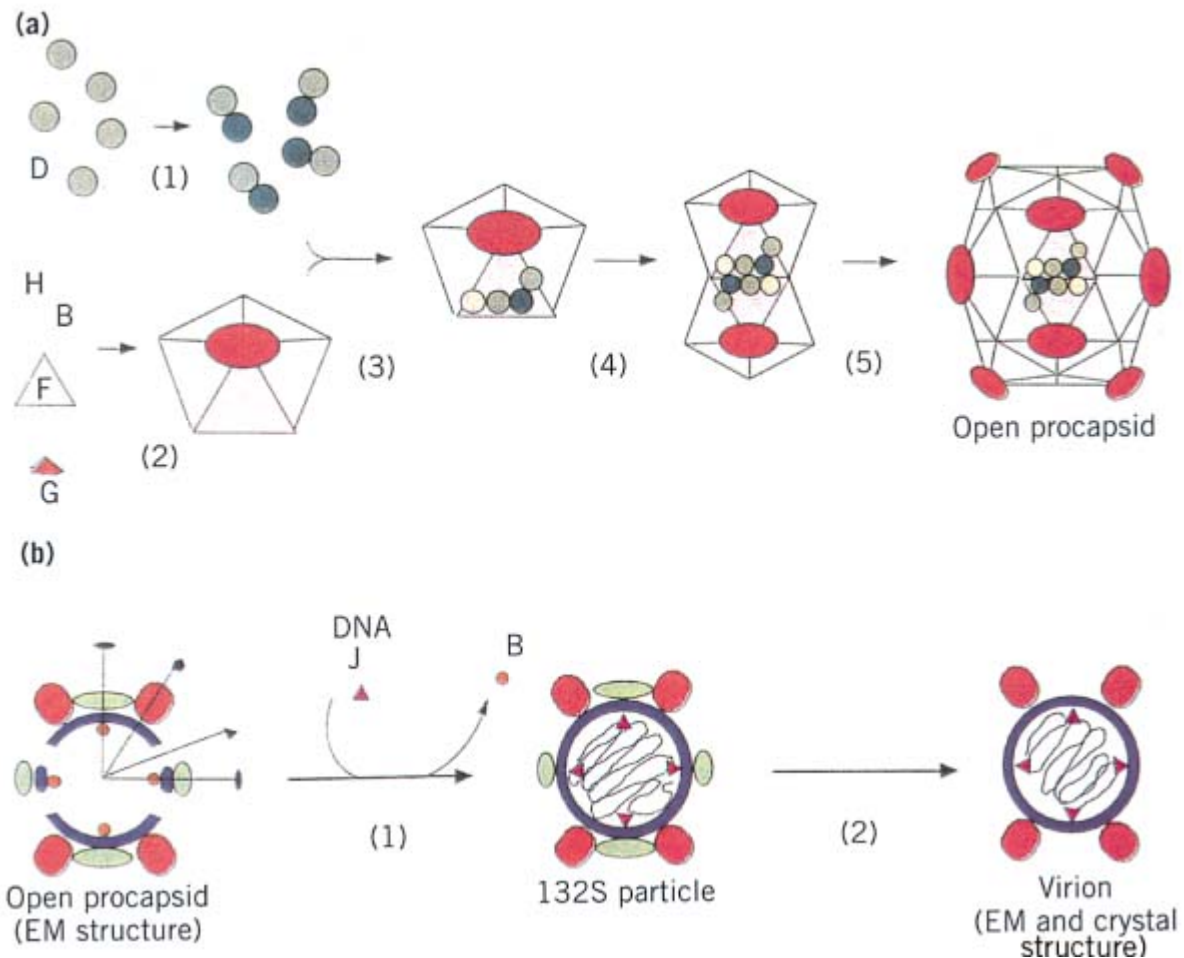
Self-assembly *in vitro* and the assembly of supramolecular structures *in vivo* proceed via specific assembly pathways, not by random association of various components, and generally require a nucleation event, which could be the formation of a dimer or higher complex. The first step results in a conformational alteration, which creates the sites for subsequent addition of protein subunits, leading to a protein-determined cooperative assembly pathway, also called *conformational switching* (7). In the case of TMV, the nucleation event triggers a conformational change in the bound disc of 34 protein subunits that provides binding sites for the stepwise addition of further subunits of the coat protein. It is not yet clear which of the observed differences between the structures of polymerized coat protein with and without the single-stranded RNA lead to the creation of the preferential binding sites.

In the case of TMV and other simple viruses, the length of the RNA appears to determine the length of the tubular coat protein particles formed: The ratio was found to be constant at about 20 bases of RNA per 1 nm length of assembled structure. This control mechanism can be regarded as a case of template-assisted assembly, which occurs widely. Another possible example may be the roles as protein rulers played by the giant muscle proteins **titin** and **nebulin**, both multidomain proteins ~1 μm long, in the assembly of **myosin** and [actin](#) filaments, respectively. It is also possible for the template to be removed prior to completion of the assembly. Well-known examples of the second possibility occur widely in **herpesviruses** and tailed **bacteriophages**, where scaffolding proteins direct the assembly of capsid coat proteins, by forming a three-dimensional scaffold onto which the shell is assembled with correct **icosahedral symmetry**. This is followed by the exit of the

scaffolding polypeptides and their subsequent reutilization or their removal via proteolysis. Because the scaffolding proteins in these cases are not part of the final assembly, they can be considered to act as chaperones.

Conformational alterations resulting at one or more steps of the assembly process appear to be a widely applicable regulatory mechanism. One example is conditional assembly, wherein the protein complex provides a signal for subsequent processing; an example is the homo- and heterodimeric forms of cystine-knot [growth factors](#), which nucleate the assembly of different multisubunit receptor complexes. Another is the major cooperative conformational changes that occur during procapsid assembly and maturation of double-stranded DNA bacteriophages. Covalent alterations, predominantly **proteolysis**, have been found to occur during the formation of the infectious virions of many viruses, particularly the more complex ones. Proteolysis makes the assembly irreversible, it may lead to conformational alterations required for subsequent association or increased stability, and it could be necessary for the removal of scaffolding proteins ([7](#)). [Figure 1](#) illustrates the complex pathway of the icosahedral *E. coli* bacteriophage FX174 capsid assembly based on recent X-ray crystallographic and cryoelectron microscopic work ([8](#)).

**Figure 1.** The assembly pathway of bacteriophage FX174. (**a**) Assembly of the open procapsid. Two scaffolding proteins, B and D, are required for the formation of the 108S, 36-nm-diameter procapsid from pentameric intermediates containing proteins F, G, and H. (1) The scaffolding D protein (green) forms asymmetric dimers (light and dark green); (2) the main capsid protein F forms pentameric precursors, complexed with proteins B, G, and H; (3) the D dimers bind to the pentamers in at least two distinct steps, resulting in an altered conformation of one of the D subunits (very light green); (4) dimerization of the pentameric precursors occurs through the binding of D dimers and requires protein B; (5) generation of the complete icosahedral open procapsid. (**b**) Maturation of the open procapsid presented in cross section, showing the scaffolding protein B (orange); one fivefold, one threefold and two twofold symmetry axes are indicated. (1) Packaging of DNA and J protein, accompanied by removal of B protein, leads to 132S particles; (2) removal of the external scaffolding protein D from the 132S particle generates the mature infectious virion. EM, electron microscopy; crystal, X-ray crystallography. (Adapted from Ref. [8](#), with the permission of Dr. M. G. Rossman and the journal *Nature*.) See color insert.



In addition to the use of templates, regulation of self-assembly can be achieved by the binding of additional proteins to the self-assembling structure to stop the process. For example, the lengths of actin filaments *in vivo* are precisely controlled: 0.06  $\mu\text{m}$  in the erythrocyte membrane skeleton and 1.1  $\mu\text{m}$  in the sarcomeres of skeletal muscle. Control of actin oligomerization is achieved by capping—that is, the binding of different proteins or protein complexes to the two distinct ends of each actin filament.

## 9. Self-Organization of Lipid Bilayers and Membrane Proteins

Phospholipids, triacylglycerols, and cholesterol, which constitute the membranes of cells and intracellular organelles, are **amphipathic** molecules consisting of a **hydrophilic** group attached to a larger **hydrophobic** moiety, the latter being either a long fatty-acid chain or a fused nonaromatic ring system in the case of cholesterol. They have the ability to self-organize spontaneously *in vitro* into structures in which their **polar** groups form the interface with water, excluding their hydrophobic moieties from the aqueous environment: spherical micelles, **vesicles**, and **lipid bilayers**. It should be pointed out that lipid micelles and bilayers are not precisely defined structures, and their sizes can be variable. All cell membranes and intracellular organelle membranes are lipid bilayers. Some of the large, complex viruses are enclosed by a lipid bilayer that is acquired from the plasma membrane of the infected host cell. The self-assembly of **lipoprotein** particles and of membrane-bound structures, such as **receptors**, represents processes in which the constituent nascent polypeptide chains interact cotranslationally with the membrane. Many membrane proteins will insert spontaneously and correctly into artificial lipid bilayers and liposomes, yet still require chaperones and insertional machinery *in vivo*.

## 10. Nanostructures and Supramolecular Chemistry

The constant drive toward increasing miniaturization of electronic components has led to a general and intense interest in nanostructures, which are assemblies of bonded atoms that have dimensions in the range of 1 to 100 nm. In addition to physical micromethods for forming nanostructures, nanochemistry encompasses the synthesis and characterization of ultralarge molecules and molecular assemblies that can be achieved through covalent polymerization, self-organization, controlled formation of covalent bonds, and molecular self-assembly (9). The latter two topics represent supramolecular chemistry, a discipline involving the self-assembly into ordered structures of synthetic inorganic and organic molecules and synthetic **peptides** and **polypeptides**. The origin of this rapidly evolving area lies in the seminal work carried out by Lehn, Cram, and Pedersen, followed by the discovery of fullerene, a self-assembled structure of 60 carbon atoms, by Smalley, Curl, and Kroto.

Efforts have been made by organic chemists to obtain large, supramolecular complexes via self-assembly of synthetic organic monomers, using either templates or directed interactions. Several growing areas include the design of artificial and mimetic proteins using self-assembly on templates, as well as metal-ion-assisted assembly, self-assembly of designed peptides into biopolymers having desirable material properties, and the use of self-assembled monolayers for fundamental studies of protein adsorption and cell adhesion.

### Bibliography

1. P. B. Berget (1985) in *Virus Structure and Assembly* (S. Casjens ed.), Jones and Bartlett, Boston, pp. 150–168.
2. A. Klug (1983) *Angew. Chem. Int. Ed. Engl.* **22**, 565–582.
3. A. C. Steven, B. L. Trus, F. P. Booy, N. Cheng, A. Zlotnick, J. R. Caston, and J. F. Conway (1997) *FASEB J.* **11**, 733–742.
4. I. M. Klotz, D. W. Darnall, and N. R. Langerman (1975) In *The Proteins*, Vol. **3** (H. Neurath and R. L. Hill, eds.), Academic Press, New York, pp. 294–411.
5. A. D. S. Goodsell and A. J. Olson (1993) *Trends Biochem. Sci.* **18**, 65–68.
6. R. Jaenicke (1987) *Prog. Biophys. Mol. Biol.* **49**, 117–237.
7. S. Casjens (1997) in *Structural Biology of Viruses* (W. Chiu, R. M. Burnett and R. L. Garcea, eds.), Oxford University Press, Oxford, U.K., pp. 3–37.
8. T. Dokland, R. McKenna, L. L. Ilag, B. R. Bowman, N. L. Incardona, B. A. Fane, and M. G. Rossman (1997) *Nature* **389**, 308–313.
9. G. M. Whitesides, J. P. Mathias, and C. T. Seto (1991) *Science* **254**, 1312–1319.

### Suggestions for Further Reading

10. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
11. T. E. Creighton (1993) *Proteins, Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.
12. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts, and J. D. Watson (1994) *The Molecular Biology of the Cell*, 3rd ed., Garland, New York.
13. J. Frank (1996) *Three-Dimensional Electron Microscopy of Macromolecular Assemblies*, Academic Press, San Diego, CA.
14. W. Chiu, R. M. Burnett, and R. L. Garcea, eds. (1997) *Structural Biology of Viruses*, Oxford University Press, Oxford, U.K.

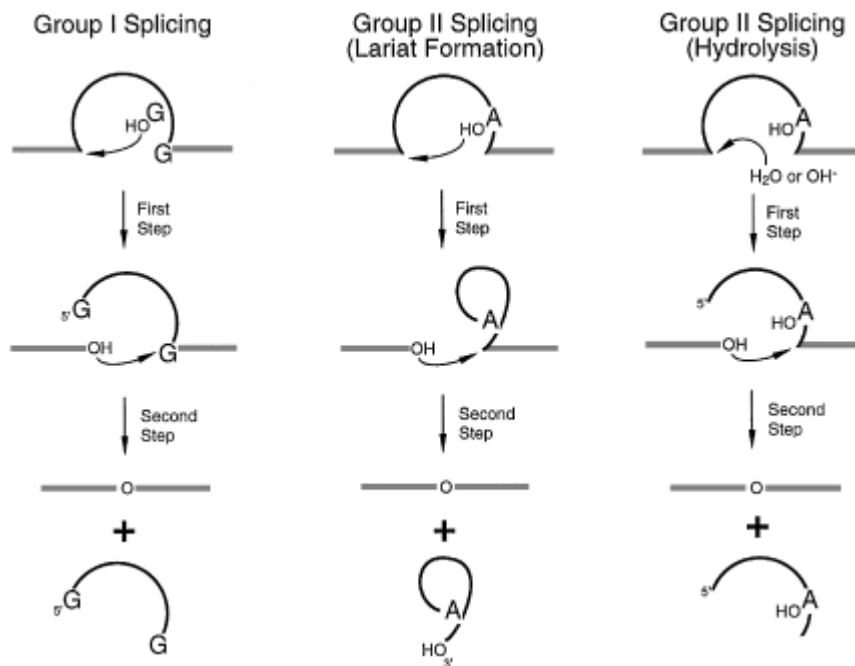
## Self-Splicing Introns

Group I and group II introns are two distinct classes of self-splicing **RNA** molecules that contain their own **active site** for **intron** removal and exon ligation (see [RNA Splicing](#)). These two classes of **ribozymes** are distinguished by their mechanisms of splicing and by their unique structures ([1](#)).

### 1. Reaction Mechanisms

Both classes of self-splicing introns perform two consecutive transesterification reactions in the process of exon ligation. The first step of splicing in a group I intron involves nucleophilic attack at the 5'-**splice site** by the 3'-OH of an exogenous, bound guanosine cofactor (Fig. [1](#)) (reviewed in Ref. [2](#)). This reaction adds the guanosine onto the 5'-end of the intron and releases the 5'-exon. The second step is analogous to the reverse of the first step. Following a conformational rearrangement of the active site in which the exogenous guanosine is replaced by the G nucleotide at the 3'-terminus of the intron, the 5'-exon attacks the 3'-exon boundary. This releases the intron and ligates the 5' and 3'-exons. Both the first and second steps of the splicing reaction are fully reversible because no net energy is consumed. The result of these reactions is that the flanking exons are ligated and the intron is released as a linear molecule with an uncoded G at the 5'-end.

**Figure 1.** Two-step reaction mechanisms for the self-splicing group I and group II classes of introns.



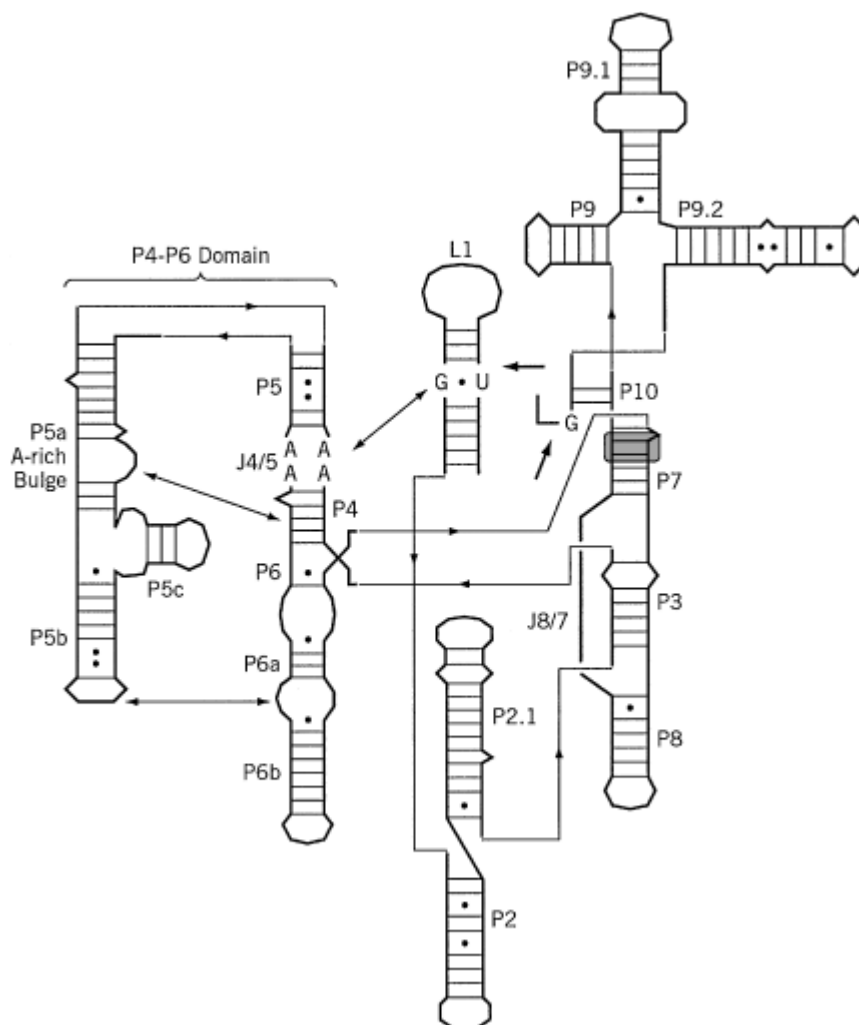
Unlike the group I introns, a group II intron utilizes an internal nucleophile for the first step of splicing (reviewed in Ref. [3](#)) (Fig. [1](#)). The 2'-OH of a highly conserved bulged A nucleotide located within domain 6 of the intron attacks the 5'-splice site. This results in release of the 5'-exon and formation of a **lariat structure** whose 5'-end of the intron is covalently attached to the 2'-OH of the bulged A. Alternatively, the 5'-exon can be released by hydrolysis in which [water](#) is the nucleophile. The second step of splicing involves attack by the 5'-exon on the 3'-splice site. This results in ligation of the flanking exons and release of the lariat or linear intron. The reaction catalyzed by the

group II intron follows the same mechanism as that employed in the more complex process of [messenger RNA precursor splicing](#) catalyzed by the [spliceosome](#) (4). For this reason, it has been proposed that group II intron splicing is an evolutionary precursor to pre-mRNA splicing.

## 2. Structures

Group I and group II introns catalyze their self-splicing reactions by folding into a distinct [tertiary structures](#) comprised of many conserved secondary structural elements (5, 6). The hallmarks of a group I intron include a common secondary structure of 10 paired segments (termed P1–P10) (Fig. 2) and several single-stranded joiner (or J) segments between the double helices (5). The “catalytic core” of the intron is made up of about 120 nucleotides that are highly conserved among group I introns isolated from a broad diversity of biological sources. The conserved nucleotides are clustered within paired regions P4, P6, P3, and P7 and the joiner regions J3/4, J4/5, J6/7, and J8/7. There is also a universal requirement for a G · U **wobble pair** at the 5'-exon boundary. The intron includes a guanosine binding site located in the major groove of the P7 helix, an internal [guide RNA](#) sequence that base pairs to the 5'-exon to form the P1 helix, a catalytic cleft within joiner regions J4/5 and J8/7 that orients the P1 helix into the active site for nucleophilic attack, and a terminal G nucleotide that defines the 3'-splice site. Although these elements are quite distant within the intron's linear sequence, they converge spatially within the tertiary structure to catalyze the two splicing steps.

**Figure 2.** Secondary structure of the *Tetrahymena* group I intron showing the conserved helices P1–P10, the 5' and 3' splice sites (large arrows), and the guanosine binding site (gray box). The double-headed arrows indicate a few sites of tertiary interaction between different regions of the intron.

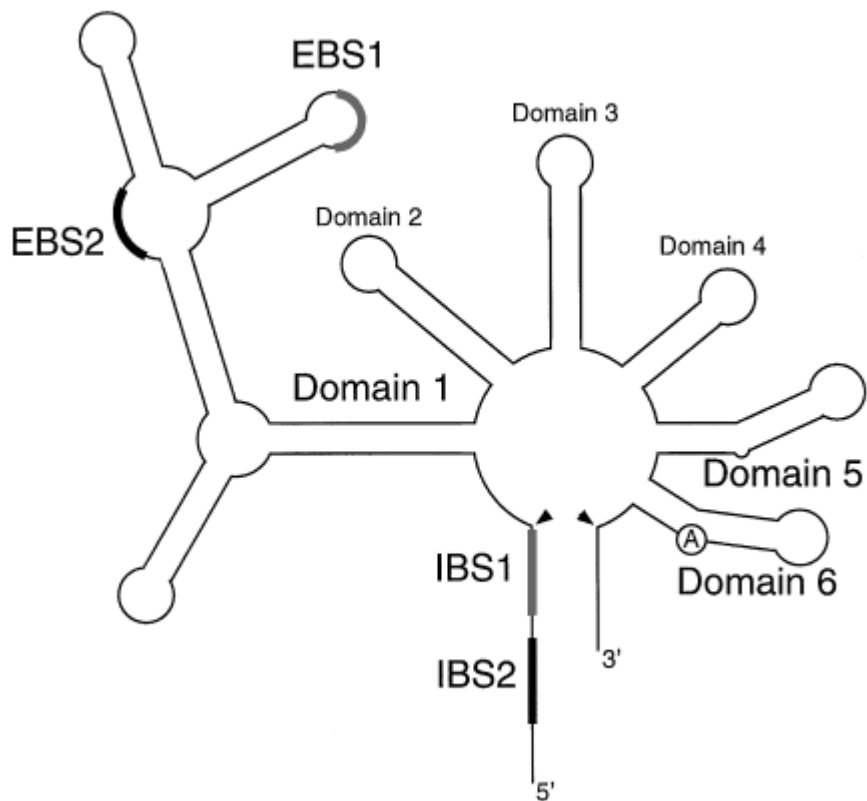


An excellent example of this structural convergence was observed within the P4–P6 domain of the *Tetrahymena* group I intron (7). The crystal structure of this independently folding domain of the intron was recently reported (8). This 160-nucleotide fragment is the largest RNA structure currently known and provides numerous insights into the molecular basis of RNA folding (see [RNA Structure](#)). The structure is formed from two long helical segments (helices P6, P4, P5 and helices P5b, P5a) that pack side by side. A bend of about 150° between the helices allows the two helical elements to form two extensive arrays of tertiary interactions. These contacts form between the minor groove of the P4 helix and the A-rich bulge of P5a and between the GAAA tetraloop of P5b and a tetraloop receptor in P6a/6b. Several new motifs of RNA structure and for metal binding were identified within this folding domain, including the “A-platform” and the “ribose zipper” (8, 9). When viewed from the side, the P4–P6 helix is essentially one RNA helix thick. It comprises about half of the intron active site but does not include the P1 helix or the guanosine binding site. Models have been proposed for the complete structures of four different group I introns.

Although there are substantially fewer conserved nucleotides, a group II intron also has a distinct structure (6) (Fig. 3). The intron is organized into six double-helical domains (D1–D6) that originate from a central wheel. Each of the domains has a particular function in the activity of the intron, although D1, D5, and D6 form the intron core. Domain 1 is an extended multihelical element that includes two looped regions (termed EBS1 and EBS2) complementary to two sites in the 5'-exon (termed IBS1 and IBS2). Interdomain interactions also allow domain I to serve as the scaffold upon which the intron active site is built. In contrast to domain 1, domain 5 is a relatively small (34-nucleotide) helical element, yet a large percentage of the phylogenetically conserved nucleotides is concentrated in this domain. It constitutes the catalytic center of the group II intron. D5 catalyzes 5'-splice site hydrolysis even when the domain is added *trans* to the rest of the intron (10). The chemical groups responsible for this catalytic enhancement map into the major groove of the D5 helix in the nucleotides surrounding a conserved G · U wobble pair (11). Domain 6 includes the branch-point adenosine whose 2'-OH nucleophilically attacks the 5'-splice site.

**Figure 3.** Schematic secondary structure of the group II intron. The bulged A, which acts as the nucleophile in the first step of splicing, is circled in Domain 6. The complementary sequences (EBS and IBS) between Domain I of the intron and the 5'-exon are shown as shaded lines. Arrows represent the splice sites.





### 3. Biological Occurrence

Group I introns are found within **genes** for mRNA, **ribosomal RNA**, and [transfer RNA](#) (1). They are extremely widespread across phylogeny and have been found in mitochondrial, chloroplast, and nuclear **genomes** of diverse eukaryotes, although they have not yet been observed in vertebrates. The discovery of a group I intron in **T4 bacteriophage** was the first example of RNA splicing observed in a prokaryote. Group I introns have also been found in eubacterial genomes. There are currently more than 450 examples of group I introns in the genomic [databases](#). All known group II introns are located within eukaryotic organelles, including plant and fungal mitochondrial DNA and the majority of introns in plant chloroplasts (1).

### 4. Metals in Folding and Catalysis

Both group I and group II introns are metalloenzymes, which require divalent metal cations for activity (12). Although the metal specificity varies substantially among group I introns, it is usually satisfied by  $Mg^{2+}$  or  $Mn^{2+}$ . Group II splicing is substantially less efficient and requires nonphysiological concentrations of monovalent and divalent cations (as high as 2.0 M KCl and 100 mM  $MgCl_2$ ), which play both structural and catalytic roles. A dramatic example of divalent metals in RNA folding is in the P5abc subdomain of the *Tetrahymena* group I intron. Three  $Mg^{2+}$  ions coordinate to separate phosphate groups within the subdomain (13). This allows the RNA to fold inside-out, that is, the phosphates point into the structure and the nucleotide bases point out to the solvent. Some structural metals in RNA can often be substituted with polycations, such as **spermidine** or cobalt hexamine, which emphasizes the importance of charge neutralization in RNA folding.

Biochemical evidence has implicated two metal ions in the chemical [transition state](#) of group I intron splicing (14, 15). One of these metals activates the nucleophile, and the second stabilizes the leaving group during the transesterification reaction. These metal-binding sites are highly selective and

cannot be substituted with a generic polycation. A two-metal active site has also been proposed for the group II intron reaction mechanism, although the evidence for this mechanism is not as complete (11). A ribozyme active site that has two metal ions is analogous to those seen in protein [polymerases](#) that catalyze the transesterification reactions of replication and [transcription](#) (16).

## 5. Accessory Protein Factors

Although some group I and group II introns undergo efficient self-splicing *in vitro*, several have no *in vitro* splicing activity. This is true of the majority of the group II introns that have been isolated. Splicing of these introns is likely to be promoted by accessory protein factors that assist the RNA in forming the appropriate active structure. For example, the group I intron in *Neurospora* pre-rRNA has no self-splicing activity unless complexed with the CYT-18 accessory protein (17). Even the yeast mitochondrial group II introns that have *in vitro* splicing activity require nuclear genes to splice efficiently *in vivo* (18).

## 6. Multiple-Turnover Catalysis

Both the group I and group II introns can be converted into ribozymes capable of multiple-turnover catalysis. The group I intron was converted to a true enzyme by eliminating the 3'-splice site and breaking the connection of the 5'-exon to the rest of the intron. An oligonucleotide analog of the 5'-exon was added to the ribozyme as a substrate (19) where it bound the internal guide sequence within the intron and was cleaved in a reaction equivalent to the first step of splicing. In this form, the ribozyme is not covalently altered during the reaction, and the RNA can cleave multiple substrates. Under multiple-turnover conditions where the substrate is in excess of the ribozyme concentration, the rate-limiting step is not the chemical reaction but release of the 5'-exon (20). A multiple-turnover form of the group II intron was created by deleting the 3'-exon and domains 4 and 5 (10). Upon adding independent domain 5, 5'-splice site hydrolysis of the truncated RNA was catalyzed. Analysis of this construct demonstrated that the chemical step is slower than the association or dissociation rates for substrate binding (21).

## Bibliography

1. T. R. Cech (1993) in *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, N.Y., pp. 239–269.
2. T. R. Cech and D. Herschlag (1996). In *Nucleic Acids and Molecular Biology* (F. Eckstein and D. M. J. Lilley, eds.), Springer, New York, pp. 1–17.
3. A. M. Pyle (1996) in *Nucleic Acids and Molecular Biology* (F. Eckstein and D. M. J. Lilley, eds.), Springer, New York, pp. 75–107.
4. T. R. Cech (1986) *Cell* **44**, 207–210.
5. F. Michel and E. Westhof (1990) *J. Mol. Biol.* **216**, 585–610.
6. F. Michel, K. Umenson, and H. Ozeki (1989) *Gene* **82**, 5–30.
7. F. L. Murphy and T. R. Cech (1993) *Biochemistry* **32**, 5291–5300.
8. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1678–1685.
9. J. H. Cate, A. R. Gooding, E. Podell, K. Zhou, B. L. Golden, A. A. Szewczak, C. E. Kundrot, T. R. Cech, and J. A. Doudna (1996) *Science* **273**, 1696–1699.
10. K. A. Jarrell, R. C. Dietrich, and P. S. Perlman (1988) *Mol. Cell. Biol.* **12**, 1950–1958.
11. B. B. Konforti, D. L. Abramovitz, C. M. Duarte, A. Karpeisky, L. Beigelman, and A. M. Pyle (1998) *Mol. Cell* **1**, 433–441.
12. A. M. Pyle (1993) *Science* **261**, 709–714.
13. J. H. Cate and J. A. Doudna (1996) *Structure* **4**, 1221–1229.
14. J. A. Piccirilli, J. S. Vyle, M. H. Caruthers, and T. R. Cech (1997) *Nature* **362**, 85–88.

15. L. B. Weinstein, B. C. N. M. Jones, R. Cosstick, and T. R. Cech (1997) *Nature* **388**, 805–808.
16. T. A. Steitz and J. A. Steitz (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6498–6502.
17. G. Mohr, M. G. Caprara, Q. Guo, and A. M. Lambowitz (1994) *Nature* **370**, 147–150.
18. G. Wiesenberger, M. Waldherr, and R. J. Schweyen (1992) *J. Mol. Biol.* **267**, 6963–6969.
19. A. J. Zaugg and T. R. Cech (1986) *Science* **231**, 470–474.
20. D. Herschlag and T. R. Cech (1990) *Biochemistry* **29**, 10159–10171.
21. A. M. Pyle and J. B. Green (1994) *Biochemistry* **33**, 2716–2725.

### Suggestion for Further Reading

22. An excellent description of the initial discovery of the self-splicing group I intron is by T. R. Cech (1990) *Biosci. Rep.* **10**, 239–261.

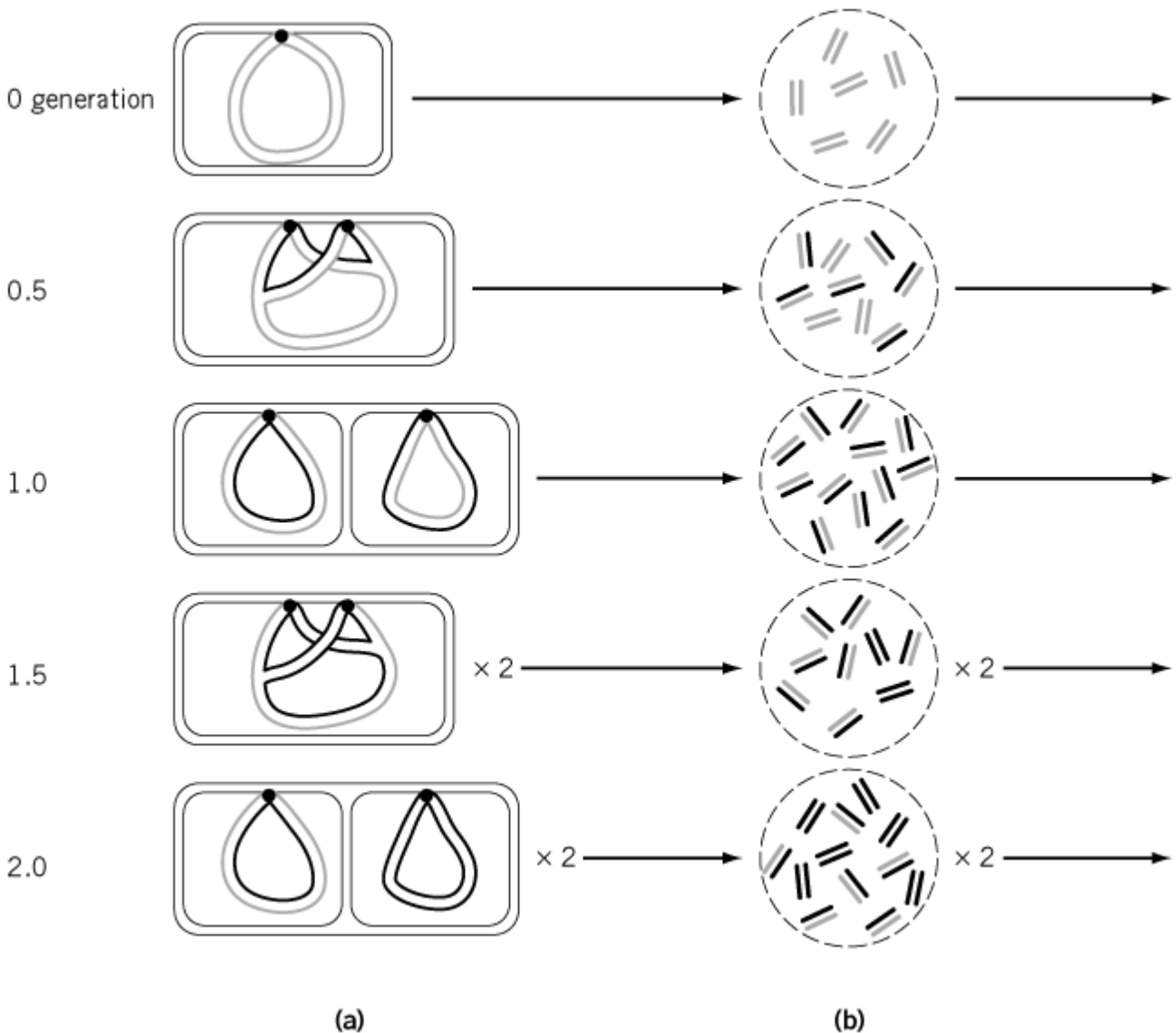
## Semi-Conservative DNA Replication

The double-helical model of **DNA** proposed by Watson and Crick in 1953 predicted how the molecule could self-duplicate, by assembling deoxyribonucleotides according to the complementary Watson–Crick base-pair rule of A to T and G to C on the two separated strands used as [templates](#). In other words, the double helix predicts a *semi-conservative* mode of [DNA replication](#) in which each of the original two strands is conserved in the two-daughter double-stranded DNA. Semi-conservative replication was demonstrated in an ingenious experiment by Meselson and Stahl in 1958 using CsCl [density gradient centrifugation](#) to separate DNA with different densities by labeling with a heavy isotope ([1](#)). The principle of the method is that each large molecule, like a nucleic acid or protein, attains its specific buoyant density in a solution of high density like 6 M of CsCl. When such a solution is subjected to ultracentrifugal force, it generates a concentration gradient along the gravity axis, so large molecules in the solution move to the concentration of the solute corresponding to their buoyant densities and there form sharp bands.

*Escherichia coli* was grown in  $(\text{NH}_3)_2\text{SO}_4$  labeled with heavy nitrogen  $^{15}\text{N}$  for several generations, to make high-density DNA labeled with  $^{15}\text{N}$  in both strands,  $^{15}\text{N} : ^{15}\text{N}$ . The culture was then shifted to a medium containing a rich nitrogen source of normal density,  $^{14}\text{N}$ . According to the semi-conservative mode of replication, newly synthesized DNA should contain a hybrid density, one strand heavy and one strand light,  $^{15}\text{N} : ^{14}\text{N}$ . All the DNA molecules are expected to have hybrid density after one generation growth in light nitrogen, and a third molecule of  $^{14}\text{N} : ^{14}\text{N}$  would appear in the following generations. As expected, the three types of molecules varying in densities were identified, and changes of molecular species from one to the other during the growth of cells after the shift from the heavy to the light medium were exactly as expected from the model (Fig. [1](#)). The Meselson–Stahl experiment not only proved the semi-conservative replication of DNA, but it also demonstrated directly that DNA is composed of two strands, as the molecule of hybrid density produced two types of single-stranded molecules differing in density upon **denaturation** by heat or alkali. The semi-conservative replication found in *E. coli* is ubiquitous in all [genomes](#) that are composed of double helical DNA, from **plasmids** to humans.

**Figure 1.** Demonstration of semi-conservative DNA replication. (a) Cells grown in  $^{15}\text{N}$  medium and containing chromo:

were further incubated in  $^{14}\text{N}$  medium for 2.0 generations. Chromosomes after the various generations indicated are drawn strand synthesized in the  $^{14}\text{N}$  medium. **(b)** DNA fragments isolated from each cell population are illustrated schematically chromosome is divided in six fragments (in reality into more than 1000 fragments). **(c)** The DNA preparation in **(b)** is used to obtain three discrete bands of DNA with different buoyant densities. The amount of DNA in each band is determined graphically. HH corresponds to heavy (grey)/heavy (grey), HL to heavy (grey)/light (black) and LL to light (black)/light (black)



## Bibliography

1. M. Meselson and F. Stahl (1958) Proc. Natl. Acad. Sci. USA **44**, 671–682.

## Seminal Ribonuclease

Bovine seminal ribonuclease (BS-RNase), or simply “seminal ribonuclease”, as the enzyme is present only in the seminal plasma of cattle and water buffalo bulls, is an eccentric [enzyme](#) and

**protein** in many ways, not the least of which is its absence in all other mammals investigated thus far. It is the only dimeric RNase of the entire **ribonuclease A superfamily** of more than 100 homologous members (see [Homology](#)) (1); it is present in seminal plasma at extremely high concentrations (over 1.5 mg/ml); it is endowed with surprising special biological functions distinct from, but dependent on, its catalytic action, including a cytotoxic action selective for tumor cells, an immunosuppressive action, and an aspermatogenic action. The structural and biological properties of the enzyme have been extensively reviewed (2). BS-RNase can be isolated as a soluble protein from bull seminal plasma, or extracted from bovine seminal vesicles, where it is produced (3). It has also been prepared as a [recombinant protein](#) in *Escherichia coli* (4) and in eukaryotic cells (5).

## 1. Structural Properties of Seminal RNase

The subunit of BS-RNase shares with RNase A from bovine pancreas, the archetype of the vertebrate RNase A superfamily, more than 80% of its amino acid residues. BS-RNase differs in that its two subunits are linked non-covalently, as well as by two [disulfide bonds](#), which bridge Cys31 and Cys32 of one subunit with Cys32 and Cys31, respectively, of the partner subunit. Seminal RNase is a homodimer; however, Asn67 of the sequence is readily deamidated (see [Deamidation](#)) in one or in both subunits, so that BS-RNase as isolated is a mixture of three components: one with intact Asn67 in both subunits, one with Asn67 deamidated in one subunit, the third with Asn67 deamidated in both subunits (6, 7).

The [tertiary structure](#) of BS-RNase subunit is the classical RNase A fold: a kidney-shaped structure, comprising a three-stranded curved [beta-sheet](#), in which the N-terminal **alpha-helix** (helix-I) is inserted, linked to helix-II by a linker peptide segment. A third  $\alpha$ -helix runs across the subunit. [X-ray crystallography](#) analysis of BS-RNase (8) has shown that each subunit is not autonomous, in that its tertiary structure is constructed with the N-terminal  $\alpha$ -helix from one subunit, and the main protein body from the other. Biochemical studies (9) have shown that the intertwined structure (termed MXM) is not the only [quaternary structure](#) available to covalently linked seminal RNase chains. The protein also can adopt another quaternary conformation T (termed M = M) in which each subunit folds onto itself, without interchange of parts between subunits. The two quaternary forms are in equilibrium with each other in solution, with the intertwined structure as the more stable form, present in a ratio of 2 to 1 with the other form (9).

Although the enzyme is isolated as a covalent dimer, and its special biological actions depend on its dimeric structure, selective reduction of its intersubunit disulfide bonds, followed by protection of the exposed [thiol groups](#) and mild treatment with [urea](#), dissociates the dimer into stable, enzymatically active monomers (10). When the urea treatment is omitted, only the M = M form (with no interchange of parts between subunits) dissociates (9). Upon refolding, after complete **denaturation** and reduction, BS-RNase chains fold first into superactive monomers (11), which later associate into M = M dimers; these eventually transform partly into, and equilibrate with, the MXM form of dimers (9). Thus BS-RNase unfolds in a 3-state transition, in that a stable, folded monomer is an equilibrium intermediate between the associated dimer and the unfolded monomer (see [Protein Folding In Vitro](#)). This finding may have an evolutionary significance, as refolding may recapitulate the evolutionary history of a protein (12). It has been hypothesized that an ancestor of the **gene** encoding the uniquely dimeric RNase encoded a stable monomer, which after destabilizing genetic alterations could then find a lower free energy minimum, and thus a new level of stability, in the association with an identical monomer. The key genetic alterations probably occurred within the stretches of sequence encoding the linker peptide connecting helix-I and -II and the helix-II, which compose the intersubunit interface. In particular, these would include substitutions providing the evolved protein chain with a Pro19 residue in the linker peptide, and Leu28, Cys31 and Cys32 residues in helix-II. These residues have been demonstrated in a [protein engineering](#) experiment to be essential for the transformation of a monomeric RNase such as RNase A into a stable dimeric and biologically active RNase (13, 14).

## 2. Biological Properties of Seminal RNase

As an enzyme, BS-RNase is unique among the vertebrate RNase A superfamily in that it is endowed with **allosteric** properties (15, 16). In the second, hydrolytic step of the ribonuclease reaction, the substrate regulates the enzyme activity allosterically, with **negative cooperativity** at low concentrations, and positive **cooperativity** at high concentrations. In contrast, the first catalytic step, the transphosphorolysis of the phosphodiester bonds, follows typical **Michaelis-Menten kinetics**. On the other hand, this finding can hardly be reconciled with the extracellular nature of the enzyme, and with its high concentration in seminal plasma.

Of greater interest biologically is the ability of the enzyme to degrade effectively both single- and double-stranded RNA, under conditions in which most homologous extracellular RNases only cleave single-stranded substrates (17). This activity is due not only to the dimeric structure of the enzyme, and to its high positive net charge (18, 19), but also, and more specifically, to the presence at key positions of the seminal RNase chain of residues Gly38 and Gly111, which in the RNases with low or no activity on double-stranded RNA are replaced by acidic residues (20).

Also of interest is the resistance of seminal RNase to the cytosolic **ribonuclease inhibitor**, which is clearly due to its dimeric structure (21). The docking of a sensitive monomeric RNase, such as RNase A (22), involves the very protein surface that becomes a subunit interface in dimeric BS-RNase. Thus when dimeric BS-RNase dissociates, as presumably happens in the more reducing environment of the cytosol, the protein becomes sensitive to the action of the cytosolic RNase inhibitor (23).

As anticipated above, the interest in seminal RNase is also directed towards its surprising and diverse biological functions. It has been proposed (24) that they are mere reflections in laboratory assay systems of the actual physiological action of the enzyme, which is yet to be conclusively deciphered. On the other hand, the high concentrations of the enzyme in seminal plasma, and the disappearance of immunosuppression in bull semen deprived of the enzyme, have led to the hypothesis (25) that the immunosuppressive activity of the enzyme is the basis for its real physiological action, which consists in the immune protection of the non-self sperm cells in the female host.

Given the general interest in the antitumor action of seminal RNase since its early discovery (26), many efforts have been directed towards the elucidation of the structural and functional determinants of this activity and of its mechanism of action. The dimeric structure of the protein and the integrity of its catalytic function have been found to be essential, since neither monomeric nor inactivated derivatives of the protein display antitumor activity (27), as is the interchange of N-terminal  $\alpha$ -helices between subunits (28, 29). This is also apparent from the results of the **site-directed mutagenesis** experiments cited above, by which monomeric RNase A, devoid of antitumor action, was transformed into a dimeric protein with antitumor activity (13, 14). The antitumor activity is correlated with the interchange of parts between subunits, and both depend on the simultaneous presence in the subunit chains of a Pro19 residue at the linker segment connecting the exchanged helix-I to helix-II, of the Cys31 and Cys32 residues responsible for the intersubunit disulfide bonds, and of Leu28 residue, which participates in a **hydrophobic** interaction at the intersubunit interface.

Some of the steps of the targeting pathway of the protein have also been defined: the protein is first secreted and concentrated on the **extracellular matrix**, where it is bound specifically; then it is internalized by non-receptor mediated **endocytosis**; through unknown pathway(s), it is eventually delivered to the cytosol, where it specifically degrades **ribosomal RNA**, thus blocking **protein biosynthesis**, and leading to the death of the cell.

Besides the immunosuppressive activity, the other “special,” ie, noncatalytic, bioactions of seminal RNase are its aspermatogenic, antitumor, antiviral, and embryotoxic activities (2). Perhaps the most interesting of them is the antitumor activity, discovered by Josef Matousek and coworkers (26) almost 30 years ago. Thus, it is not surprising that many efforts have been directed towards the elucidation of the structural and functional determinants of this activity and its mechanism of action

(30). The main determinants have been proposed to be the integrity of the catalytic action and the dimeric structure of the protein (27). The former is clearly related to the ability to degrade rRNA, which the seminal enzyme exerts in malignant, but not in normal, cells (31). As for the significance of the dimeric structure in the protein antitumor action, this has been related to its ability to evade as a dimer the neutralizing effect of the cytosolic RNase inhibitor (see above). The main data leading to this conclusion are: 1) The higher cytotoxic activity of the MXM form (29) compared with that of the M=M form; the latter, upon reduction of the intersubunit disulfides by the cytosolic reducing environment, dissociates into monomers, whereas MXM remains a dimer, albeit non-covalent, hence, resistant to the cytosolic inhibitor (32); 2) Monomeric non-cytotoxic RNases can be made cytotoxic by rendering them resistant to the inhibitor (33, 34).

As for the mechanism of anti-tumor action of seminal RNase, a more complete picture is emerging. The protein first accumulates at the extracellular matrix, where it binds with high affinity to specific, saturable, reversible binding sites (31), and then it is internalized by nonreceptor mediated endocytosis and packaged in endosomes; it reaches the *Trans*-Golgi Network (TGN), and, eventually, the cytosol (A. Bracale et al, unpublished results), where it degrades rRNA. This engenders a block of protein synthesis and apoptotic death (35). The observation that TGN and cytosol are reached by the protein only in tumor cells justifies the selectivity of the anti-tumor action of seminal RNase, but it does not shed light on the reason why this toxic pathway is followed by the enzyme only in tumor cells. Clearly, the resistance to the cytosolic inhibitor does not explain all seminal RNase cytotoxic potential and its selectivity.

#### Bibliography

1. J. J. Beintema, H. J. Breukelman, A. Carsana, and A. Furia (1997) In *Ribonucleases: Structures and Functions* (G. D'Alessio and J. F. Riordan, eds.), Academic Press, San Diego, pp. 245–269.
2. G. D'Alessio, A. Di Donato, L. Mazzarella, and R. Piccoli (1997) In *Ribonucleases. Structures and Functions* (G. D'Alessio and J. F. Riordan, eds.), Academic Press, San Diego, pp. 383–423.
3. M. Tamburrini, R. Piccoli, R. De Prisco, A. Di Donato, and G. D'Alessio (1986) *Ital. J. Biochem.* **35**, 22–32.
4. M. de Nigris, N. Russo, R. Piccoli, G. D'Alessio, and A. Di Donato (1993) *Biochem. Biophys. Res. Commun.* **193**, 155–160.
5. N. Russo, M. De Nigris, A. Di Donato, and G. D'Alessio (1993) *FEBS Lett.* **318**, 242–244.
6. A. Di Donato, P. Galletti, and G. D'Alessio (1986) *Biochemistry* **25**, 8361–8368.
7. A. Di Donato, M. A. Ciardiello, M. de Nigris, R. Piccoli, L. Mazzarella, and G. D'Alessio (1993) *J. Biol. Chem.* **268**, 4745–4751.
8. L. Mazzarella, S. Capasso, D. Demasi, G. Di Lorenzo, C. A. Mattia, and A. Zagari (1993) *Acta Cryst. D* **49**, 389–402.
9. R. Piccoli, M. Tamburrini, G. Piccialli, A. Di Donato, A. Parente, and G. D'Alessio (1992) *Proc. Natl. Acad. Sci. USA* **89**, 1870–1874.
10. G. D'Alessio, M. C. Malorni, and A. Parente (1975) *Biochemistry* **14**, 1116–1122.
11. A. Parente and G. D'Alessio (1985) *Eur. J. Biochem.* **149**, 381–387.
12. G. D'Alessio (1995) *Nature Struct. Biol.* **2**, 11–13.
13. A. Di Donato, V. Cafaro, and G. D'Alessio (1994) *J. Biol. Chem.* **269**, 17394–17396.
14. A. Di Donato, V. Cafaro, I. Romeo, and G. D'Alessio (1995) *Protein Sci.* **4**, 1470–1477.
15. R. Piccoli, A. Di Donato, and G. D'Alessio (1988) *Biochem. J.* **253**, 329–336.
16. A. Di Donato, R. Piccoli, and G. D'Alessio (1987) *Biochem. J.* **241**, 435–440.
17. A. Floridi and M. Libonati (1969) *Eur. J. Biochem.* **8**, 81–97.
18. M. Libonati, S. Sorrentino, R. Galli, R. La Montagna, and A. Di Donato (1975) *Biochim. Biophys. Acta* **407**, 292–298.
19. S. Sorrentino, M. Lavitrano, R. De Prisco, and M. Libonati (1985) *Biochim. Biophys. Acta* **827**,

135–139.

20. J. C. Opitz, M. I. C., M. Haugg, K. Trautwein-Fritz, S. A. Raillard, T. M. Jermann, and S. A. Benner (1998) *Biochemistry* **37**, 4023–4033.
21. B. S. Murthy and R. Sirdeshmukh (1992) *Biochem. J.* **281**, 343–348.
22. B. Kobe and J. Deisenhofer (1996) *J. Mol. Biol.* **264**, 1028–43.
23. B. S. Murthy, C. De Lorenzo, R. Piccoli, G. D'Alessio, and R. Sirdeshmukh (1996) *Biochemistry* **35**, 3880–3885.
24. G. D'Alessio (1993) *Trends Cell Biol.* **3**, 106–109.
25. M. Tamburrini, G. Scala, C. Verde, M. R. Ruocco, A. Parente, S. Venuta, and G. D'Alessio (1990) *Eur. J. Biochem.* **190**, 145–148.
26. J. Matousek (1973) *Experientia* **29**, 858–859.
27. S. Vescia, D. Tramontano, G. Augusti Tocco, and G. D'Alessio (1980) *Cancer Res.* **40**, 3740–3744.
28. J. S. Kim, J. Soucek, J. Matousek, and R. T. Raines (1995) *J. Biol. Chem.* **270**, 31097–31102.
29. V. Cafaro, C. De Lorenzo, R. Piccoli, A. Bracale, M. R. Mastronicola, A. Di Donato, and G. D'Alessio (1995) *FEBS Lett.* **359**, 31–34.
30. R. J. Youle and G. D'Alessio (1997) In *Ribonucleases: Structures and Functions* (G. D'Alessio and J. F. Riordan, eds.), Academic Press, San Diego, CA.
31. M. R. Mastronicola, R. Piccoli, and G. D'Alessio (1995) *Eur. J. Biochem.* **230**, 242–249.
32. B. S. Murthy, C. De Lorenzo, R. Piccoli, G. D'Alessio, and R. Sirdeshmukh (1996) *Biochemistry* **35**, 3380–3385.
33. P. A. Leland, L.W. Schultz, B.-M. Kim, and R. T. Raines (1999) *Proc. Natl. Acad. Sci. U.S.A.* **95**, 10407–10412.
34. N. Russo, A. Antignani, and G. D'Alessio (2000) *Biochemistry* **39**, 3585–3591.
35. J. Cinatl Jr., J. Cinatl, R. Kotchetkov, J. U. Vogel, B.G. Woodcock, J. Matousek, P. Pouckova, and B. Kornhuber (1999) *Int. J. Oncol.* **15**, 1001–9.

### Suggestions for Further Reading

36. G. DiAlessio and J. F. Riordan (eds.) (1997) *Ribonucleases. Structures and Functions*, Academic Press, San Diego, CA, USA.
37. R. J. Youle and G. D'Alessio (1997) "Antitumor Ribonucleases", In *Ribonucleases. Structures and Functions*, (G. D'Alessio and J. F. Riordan, eds.), Academic Press, San Diego, CA, USA.

### Semliki Forest Virus

Semliki Forest virus (SFV) is a member of the **virus** genus *Alphavirus* in the family *Togaviridae*. SFV was originally isolated in Western Uganda in 1942 from mosquitoes, *Aedes abnormalis*, and since then it has also been found in African and Eurasian countries. In nature, SFV circulates between birds and mosquitoes, while experimentally SFV has a wide host range of mammalian, avian, and insect cells. In rare cases, SFV causes rash, fever, or arthritis in humans. SFV and Sindbis virus, the prototype alphavirus, are most extensively studied molecular biologically among the 26 members in the genus.

The genome is a positive-strand RNA molecule with 11,422 nucleotides, excluding the [cap](#) and the



3' [poly\(A\)](#) tail. The virion is 70 nm in diameter and consists of a nucleocapsid 40 nm in diameter, which contains the genomic RNA encapsidated with 240 copies of capsid (C) protein, and a surrounding lipid bilayer [membrane](#) that anchors 80 spikes, each composed of three heterotrimeric units (240 in total) formed by **glycoproteins** E1, E2, and E3. Both the nucleocapsid and envelope spikes are arrayed in a  $T = 4$  symmetry. The **sedimentation coefficient** of the virion is 280S (in water at 20°C), and the buoyant density in sucrose is 1.22 g/mL.

The **genes** for four nonstructural proteins required for viral RNA replication and transcription—nsP1, nsP2, nsP3, and nsP4, are positioned in that order in the 5' two-thirds of the genomic RNA and are **translated** initially as [polyprotein](#) P1234. P1234 is cleaved into intermediate and four mature proteins by a papain-like [thiol proteinase](#), located in the C-terminal half of nsP2. Genetic and biochemical studies, as well as amino acid sequence analysis, indicate that nsP1 has capping enzyme activity, nsP2 has a nucleoside triphosphate-binding [RNA helicase](#) activity in the N-terminal half and a proteinase activity in the C-terminal half, and nsP4 is the viral **RNA polymerase**. The functions of nsP3 have not been elucidated, but it is a component of the replication complex and is **phosphorylated** *in vivo*.

The genes for the structural proteins C, E3, E2, 6K, and E1 are positioned in that order from the 5' terminus and are expressed from a subgenomic RNA that is collinear with the 3' one-third of the genome and is transcribed from an internal **promoter** on the full-length complementary strand. The structural proteins are translated as a polyprotein and processed into mature products by capsid autoproteinase at the C—E3 bond, by cellular [signal peptidase](#) at the E2–6K and 6K–E1 bonds, and by a cellular proteinase at the E3–E2 bond.

The virions enter the host cells by receptor-mediated [endocytosis](#), followed by fusion of the viral envelope with a host endosomal membrane, triggered by low pH. After release of the nucleocapsid into the cytoplasm, disassembly of the nucleocapsid appears to be triggered by its binding to [ribosomes](#). Subsequent RNA replication and [transcription](#) occurs entirely in the cytoplasm. The regulatory mechanism of plus-strand and minus-strand RNA synthesis was elucidated from studies with Sindbis virus. Four nonstructural proteins are translated as polyprotein P1234, in which **proteolytic** cleavage at the nsP3–nsP4 bond occurs rapidly in *cis* (intramolecularly), and the resulting P123 and nsP4 form a replication complex for full-length minus-strand RNA synthesis. Polyprotein P123 is unstable due to its intrinsic proteinase activity, and the cleaved products together with nsP4 form a stable replication complex for plus-strand genomic and subgenomic RNA synthesis. At late stages of infection, when the nsP2 proteinase concentration has increased, the nsP2–nsP3 bond of nascent polypeptide chain is cleaved, preventing formation of the P123–nsP4 complex and shutting off minus-strand RNA synthesis. Plus-strand genomic and subgenomic RNA synthesis continues, because this replication complex is composed of fully cleaved, metabolically stable products. On the other hand, structural proteins are translated from subgenomic RNA, which is produced in a large excess over the genomic RNA. Capsid proteins and a genomic RNA molecule assemble to form a nucleocapsid. The assembled nucleocapsid interacts with the cell plasma membrane, which is occupied by the viral glycoproteins. Mature virus particles are released from the cell membrane by budding.

Full-length [complementary DNA](#) of the SFV genome was **cloned** immediately downstream of a bacteriophage **RNA polymerase** promoter in a plasmid **vector**. Transcription of the cDNA insert with a bacteriophage RNA polymerase generates infectious RNA *in vitro*, which leads to production of infectious virions upon [transfection](#) of host cells. This “infectious” cDNA clone has proved to be a powerful tool for molecular genetic analysis of SFV replication and pathogenesis. The full-length clones have been also successfully used for development of transient [expression systems](#), so that genes of interest can be expressed in the cytoplasm of the transfected or infected cells.

#### Suggestions for Further Reading

R. E. Johnston and C. J. Peters (1996) "Alphaviruses". In *Fields Virology*, 3rd ed. (B. N. Fields et

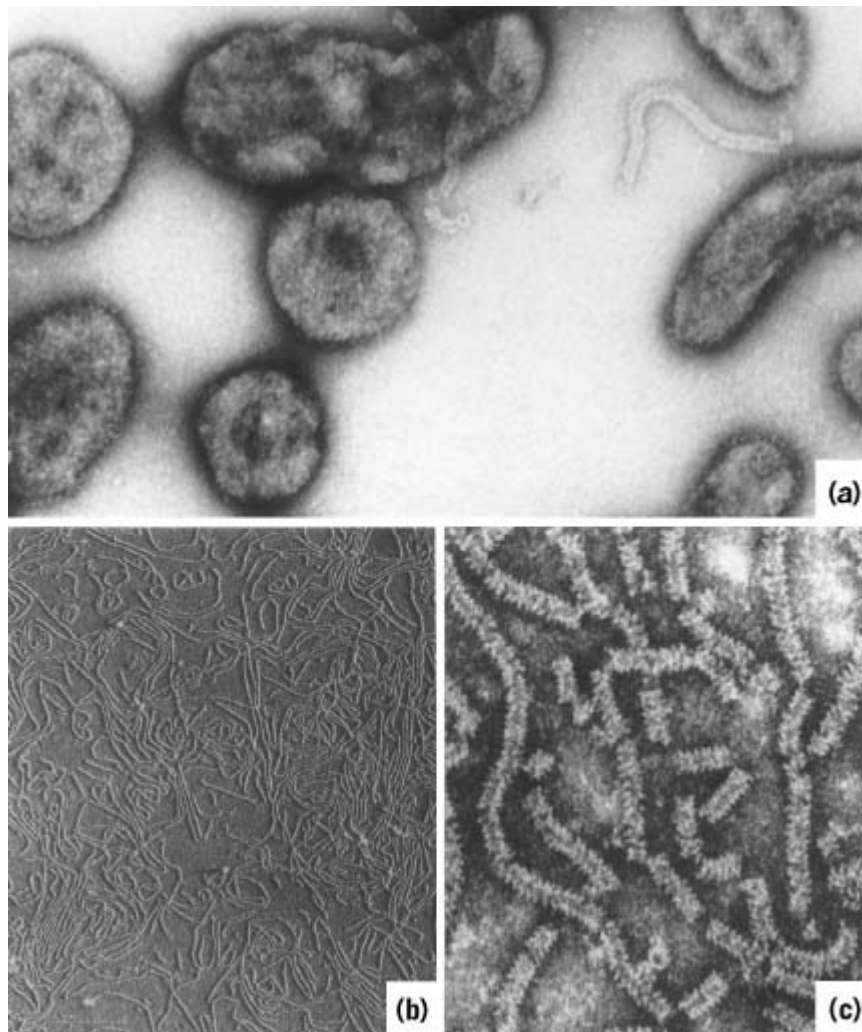
al., eds.), Lippincott-Raven, Philadelphia, pp. 843–898.

M. J. Schlesinger (1994) "Sindbis and Semliki Forest viruses". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1330–1333.

## Sendai Virus

Sendai virus (murine parainfluenza virus 1) belongs to the genus *Respirovirus* in the subfamily of *Paramyxovirinae* of the family *Paramyxoviridae*. The **virus** particles are 150 nm or more in diameter, and they are pleomorphic but usually spherical in shape (Fig. 1). The **genome** is a single molecule of linear, single-stranded, negative-sense RNA of 15,384 nucleotides. The genome RNA is tightly associated with the nucleocapsid (N) subunit proteins and **RNA polymerase**, which consists of a phosphoprotein (P) and a large (L) protein, forming a helical ribonucleoprotein complex (RNP) or nucleocapsid (Fig. 1). The virion consists of the RNP surrounded by a lipid envelope, which is derived from the host cell plasma **membrane** and contains two virus-specific **glycoproteins**, the hemagglutinin-neuraminidase (HN) and the fusion (F) protein. These glycoproteins are present as homooligomers, forming spike-like projections of 8 nm in length. HN binds to the receptor sialic acid residues on the host cell surface, whereas the F protein facilitates virus entry by mediating fusion of the viral envelope with the host cell plasma membrane. F is derived from the biologically inactive precursor F<sub>0</sub> glycoprotein through proteolytic processing by a host cell endoprotease. There is a matrix (M) protein layer between the envelope and RNP. The M protein is important for assembly and stabilization of the virion structure.

**Figure 1.** Sendai virus particles. Sendai virus virions were negatively stained with uranyl acetate (a). RNPs were revealed by shadowing (b) and negative staining (c). Magnifications: (a) ×130,000, (b) ×16,000, (c) ×220,000. (Courtesy of Dr. Y. Hosaka.)



The genes encoding the individual proteins are organized in the viral genome in the order 3'-leader-N-P-M-F-HN-L-trailer-5'. The short terminal leader and trailer regions contain the **promoters** for replication. The RNP, but not the naked RNA, can serve as template for both [transcription](#) and replication. The **RNA polymerase** enters the 3' terminal region and generates each mRNA successively by a stop–start mechanism controlled by the specific *cis*-acting signals present at each gene boundary. Therefore, the viral gene expression is generally monocistronic, although expression of the P gene is a notable exception. The P [messenger RNA](#) directs the synthesis of not only P protein, but also a nested set of accessory proteins, C', C, Y1, and Y2, collectively referred to as C proteins, by using multiple [initiation codons](#) in the +1 [reading frame](#) relative to that of P. The P gene further generates an mRNA encoding another accessory protein, V, by cotranscriptional RNA editing featuring the insertion of a single guanine residue at a specific site. The P and V proteins thus have the same amino-terminus but differ in their carboxyl-terminal halves. After translation of the mRNAs and accumulation of viral proteins, the same RNA polymerase copies the same RNP template, but somehow ignores all junction stop–start signals and the editing site, to generate antigenomic RNP with a full-length positive strand RNA. This RNA serves as the template to generate genomic RNP. All replication steps take place in the cytoplasm, and maturation occurs by budding from the plasma membrane.

Sendai virus was first isolated in 1952 in Sendai, Japan, by inoculating mice with the autopsy lung tissue from a newborn child with pneumonia. In 1955 it was initially named hemagglutinating virus of Japan (HVJ) by the Japanese Society for Virology. The question soon arose, however, as to whether it was the causative agent of a human disease, because the virus was found to be prevalent

worldwide in mice and other rodents and to cause enzootic or epizootic infections in these animals. Indeed, no direct evidence supporting the involvement of Sendai virus in any human disease is yet been available. Instead, the Sendai virus system in the mouse has been widely used as a model for studying acute respiratory infections. Recently, a technology has been established that allows infectious virus to be generated solely from cDNA, permitting the genetic engineering of Sendai virus. This technology is not only solving previously unsettled issues, including the roles of the Sendai virus accessory proteins, but also opening the possibility to use Sendai virus as a novel expression vector.

### Suggestions for Further Reading

- N. Ishida and M. Homma (1978) Sendai virus. *Adv. Virus. Res.* **23**, 349–384.
- Y. Nagai (1993) Protease-dependent virus tropism and pathogenicity. *Trends. Microbiol.* **1**, 81–87.
- A. Kato et al. (1997) The paramyxovirus, Sendai virus, V protein encodes a luxury function required for viral pathogenesis. *EMBO J.* **16**, 578–587.
- Y. Nagai (1999) Paramyxovirus replication and pathogenesis. Reverse genetics transforms understanding. *Rev. Med. Virol.* **9**, 83–99.
- Y. Yonemitsu, C. Kitson, S. Ferrari, R. Farley, U. Griesenbach, D. Judd, R. Steel, P. Scheid, J. Zhu, P. K. Jeffery, A. Kato, M. K. Hasan, Y. Nagai, M. Fukumura, M. Hasegawa, D. M. Geddes, and E. W. F. W. Alton (2000) Efficient gene transfer to airway epithelium using recombinant Sendai virus. *Nat. Biotechnol.* **18**, 970–973.

## Senescence

Tissue homeostasis is regulated not only through the elimination of excess cells by [apoptosis](#), but also by the inability of most cells to proliferate infinitely; at some stage, their further growth is arrested. Thus, cell populations are consistently replenished by a small pool of proliferating [stem cells](#), with the daughter cells they produce eventually arresting, limiting the number of cells within the population. The maintenance of cell survival in the absence of division is termed *senescence*; a senescent cell is essentially a cell that is genetically dead. For long-lived animals, senescence is a sensible option because reducing the number of cells with the capacity to divide also reduces their capacity to both mutate and form tumors (see [Cell Death](#)). Senescence is tightly regulated, with clear differences existing between proliferating and senescent cells in the proteins expressed ([1](#)). Some of the proteins involved in regulating senescence also regulate apoptosis, but the link between these two pathways, if indeed there is one, is unclear.

Primary cells in culture will undergo a finite number of passages before terminally arresting ([2](#), [3](#)) (see [Cell growth](#)). Once in this state, cells cannot be triggered to reenter the [cell cycle](#), suggesting that each cell must have a clocklike mechanism that counts each division and arrests this cell after the critical number of doublings ([3-5](#)). This clock like mechanism is now believed to occur in the [telomeres](#) of [chromosomes](#) ([6-9](#)). Telomeres are, however, not the only mechanism through which senescent cells are arrested. Senescent cells have alterations in the expression of proteins that regulate the cell cycle. In particular, senescent cells express low levels of [cyclins](#) and cyclin-dependent kinases (cdk), have no phosphorylated (active) [retinoblastoma](#) (Rb) protein and express elevated levels of the cell-cycle inhibitor p21. Senescent cells also express particular senescent genes, as shown by the [complementation](#) analysis of immortalized cells ([10](#)). **Cell fusions** between transformed cell lines can result in the hybrid cell either remaining immortal or becoming senescent. If the fused cells remain immortalized, they must contain the same genetic mutation allowing

immortalization, whereas the senescent hybrid cells must have different immortalization mutations and are therefore assigned to different complementation groups. Several putative genes that regulate senescence have been identified using this method, but their function is not yet known (1).

The regulation of cellular senescence has been described in two stages, mortality stage 1 (M1) and mortality stage 2 (M2), which are controlled by two different mechanisms. This model was derived to explain the behavior of cells transfected with inducible SV40 large [T-antigen](#). Proteins from a variety of tumor **viruses**, including T-antigen, E1A and E1B from Adenovirus type 5, and E6/E7 from human papilloma virus, extend the life span of fibroblasts in culture. Cells expressing these proteins have effectively bypassed the first stage of senescence, because these proteins sequester both [p53](#) and Rb, two proteins thought to be critical for arresting cells in M1 ([11](#), [12](#)). SV40 large-T antigen-transfected cells do not reach senescence; instead, they enter *crisis*. During crisis, cells exhibit both mitosis and apoptosis and, because death outweighs proliferation, the population is eventually lost. Only rarely does an immortalized cell emerge from crisis.

Both p53 and Rb have well-described roles in cell-cycle control ([13-17](#)) (see [Cell Cycle](#)). Activation of p53 results in either G<sub>1</sub> arrest, preventing cell replication in the presence of DNA damage ([18](#), [19](#)), or apoptosis ([20](#), [21](#)). Several functions have been ascribed to p53, including the control of [transcription](#) and **replication**, either of which may contribute to growth suppression. The induction of cell-cycle arrest by p53 is likely to be facilitated by its induced expression of p21<sup>Waf-1/Cip-1</sup> ([22-24](#)), which efficiently inhibits the activities of the cyclin-dependent kinases ([25](#)), one of which, cdk-2, phosphorylates Rb, inhibiting its interaction with the [transcription factor](#) E2F ([25-27](#)). Thus, the removal of p53 and Rb eliminates two interlocked cell-cycle inhibition steps, allowing uncontrolled cellular progression; furthermore, mutations within these cells will be undetected due to the absence of functional p53. Not only will this cooperation suppress the M1 phase of senescence, it will also prevent apoptosis, suggesting that the regulation of M1 and apoptosis overlap.

Cells that evade M1 progress instead to the M2 crisis. Only a cell that gains a mutation in one of the M2-regulating genes is able to overcome crisis and emerge as an immortalized cell. M2 is thought to be regulated by telomere shortening. The hypothesis proposed is that telomere shortening triggers the M1 phase and cells are arrested by p53 and Rb. However, in large T-antigen-expressing cells, M1 arrest is suppressed. Hence, the telomeres continue to shorten, leading to the chromosomal damage that triggers M2 and crisis ([12](#)). Cells that evade M2 are thought to enable the reexpression of telomerase, an outcome that is also found in many tumor cells. Thus, an understanding of the mechanistic pathways that facilitate senescence is critically important if the ways in which tumor cells become immortalized are to be fully understood.

## Bibliography

1. C. A. Afshari, and J. C. Barrett (1996) "Molecular genetics of in vivo cellular senescence". In *Cellular Aging and Cell Death* (N. J. Holbrook, G. R. Martin, and R. A. Lockshin, eds.), Wiley-Liss, New York, pp. 109–122.
2. L. Hayflick (1984) Intracellular determinants of cell aging. *Mech. Ag. Dev.* **28**, 177–185.
3. L. Hayflick (1976) The cell biology of human aging. *N. Engl. J. Med.* **295**, 1302–1308.
4. R. T. Dell'Orco, J. G. Mertens, and P. F. Kruse, P. F. Jr. (1974) Doubling potential, calendar time, and donor age of human diploid cells in culture. *Exp. Cell Res.* **84**, 363–366.
5. T. W. Roberts, and J. R. Smith (1980) The proliferative potential of chick embryo fibroblasts: Population doublings vs. time in culture. *Cell Biol. Int. Rep.* **4**, 1057–1063.
6. C. B. Harley, A. B. Futcher, and C. W. Greider (1990) Telomeres shorten during ageing of human fibroblasts. *Nature* **345**, 458–460.
7. G. B. Morin (1991) Recognition of a chromosome truncation site associated with alpha-thalassaemia by human telomerase. *Nature* **353**, 454–456.
8. C. W. Greider, J. C. Harley (1996) "Telomeres and telomerase in cell senescence and

- immortalisation". In *Cellular Aging and Cell Death* (N. J. Holbrook, G. R. Martin, R. A. Lockshin, eds.), Wiley-Liss, New York, pp. 123–138.
9. C. M. Counter, H. W. Hirte, S. Bacchetti, and C. B. Harley (1994) Telomerase activity in human ovarian carcinoma. *Proc. Natl. Acad. Sci. USA* **91**, 2900–2904 (see comments).
  10. O. M. Pereira Smith, and J. R. Smith (1988) Genetic analysis of indefinite division in human cells: identification of four complementation groups. *Proc. Natl. Acad. Sci. USA* **85**, 6042–6046.
  11. J. W. Shay, O. M. Pereira Smith, and W. E. Wright (1991) A role for both RB and p53 in the regulation of human cellular senescence. *Exp. Cell Res.* **196**, 33–39.
  12. W. E. Wright and J. W. Shay (1996) "Mechanisms of escaping senescence in human diploid cells". In *Cellular Aging and Cell Death* (N. J. Holbrook, G. R. Martin and R. A. Lockshin, eds.), Wiley-Liss, New York, pp. 153–166.
  13. K. Buchkovich, L. A. Duffy, and E. Harlow (1989) The retinoblastoma protein is phosphorylated during specific phases of the cell cycle. *Cell* **58**, 1097–1105.
  14. R. E. Hollingsworth, C. E. Hensey, and W.-H. Lee (1993) Retinoblastoma protein and the cell cycle. *Curr. Opin. Gen. Devel.* **3**, 55–62.
  15. P.-L. Chen, P. Scully, J.-Y. Shew, J. Y. J. Wang and W.-H. Lee (1989) Phosphorylation of the retinoblastoma gene product is modulated during the cell cycle and cellular differentiation. *Cell* **58**, 1193–1198.
  16. H. J. Xu, S. X. Hu, and W. F. Benedict (1991) Lack of nuclear RB protein staining in G0/middle G1 cells: Correlation to changes in total RB protein level. *Oncogene* **6**, 1139–1146.
  17. W. El-Deiry, T. Tokino, V. Velculescu, D. Levy, R. Parsons, J. Trent, D. Lin, W. Mercer, K. Kinzler, and B. Vogelstein (1993) *WAF1*, a potential mediator of p53 tumor suppression. *Cell* **76**, 817–825.
  18. D. P. Lane (1992) Cancer. p53, guardian of the genome. *Nature* **358**, 15–16 (news; comment).
  19. Y. Yin, M. Tainsky, F. Bischoff, L. Strong, and G. Wahl (1992) Wild-type p53 restores cell cycle and inhibits gene amplification in cells with mutant p53 alleles. *Cell* **70**, 937–938.
  20. R. E. Yonish, D. Grunwald, S. Wilder, A. Kimchi, E. May, J. J. Lawrence, P. May and M. Oren (1993) p53-mediated cell death: Relationship to cell cycle control. *Mol. Cell Biol.* **13**, 1415–1423.
  21. D. Lane (1993) A death in the life of p53. *Nature* **362**, 786 (commentary).
  22. Y. Xiong, G. J. Hannon, H. Zhang, D. Casso, R. Kobayashi, and D. Beach (1993) p21 is a universal inhibitor of cyclin kinases. *Nature* **366**, 701–704 (see comments).
  23. J. Harper, G. Adami, G. Wei, N. Keyomarsi, K.S. Elledge (1993) The p21 Cdk-interacting protein Cip1 is a potent inhibitor of G1 cyclin-dependent kinases. *Cell* **76**, 805–816.
  24. W. S. el Deiry, J. W. Harper, P. M. O'Connor, V. E. Velculescu, C. E. Canman, J. Jackman, J. A. Pietenpol, M. Burrell, D. E. Hill, and Y. Wang (1994) WAF1/CIP1 is induced in p53-mediated G1 arrest and apoptosis. *Cancer Res* **54**, 1169–1174.
  25. J. Pines (1993) Cyclins and their associated cyclin-dependent kinases in the human cell cycle. *Biochem. Soc. Trans.* **21**, 921–925.
  26. J. A. Lees, K. J. Buchkovich, D. R. Marshak, C. W. Anderson, and E. Harlow (1991) The retinoblastoma protein is phosphorylated on multiple sites by human cdc2. *Embo. J.* **10**, 4279–864.
  27. B. T. Lin, S. Gruenwald, A. O. Morla, W. H. Lee, and J. Y. Wang (1991) Retinoblastoma cancer suppressor gene product is a substrate of the cell cycle regulator cdc2 kinase. *Embo. J.* **10**, 857–864.

### Suggestion for Further Reading

28. N. J. Holbrook, G. R. Martin, and R. A. Lockshin (1996) "Cellular aging and cell death". In *Modern Cell Biology*, Vol. **319** (J. B. Harford, ed.), Wiley-Liss, New York.



## Sephadex, Sepharose, Sephacryl

Sephadex, Sepharose, and Sephacryl are trade names for three types of **gel-filtration** matrices made by the Pharmacia Fine Chemicals Company (1). Sephadex is a bead-formed gel prepared by cross-linking dextran with epichlorohydrin. The gel is extremely **hydrophilic** due to the large number of hydroxyl groups, and it therefore swells readily in water and electrolyte solutions. The G-types of Sephadex differ in their degree of cross-linking and hence in their degree of swelling and their fractionation range. Sephadex also swells in dimethylsulfoxide and formamide. Sephadex is insoluble in other nonaqueous solvents (unless it is chemically degraded). It is stable in water, salt solutions, organic solvents, and alkaline and weakly acidic solutions. Sephadex does not melt, and it may be sterilized wet at neutral pH or dry by autoclaving for 30 minutes at 120°C without affecting its chromatographic properties.

Sepharose is a bead-formed gel prepared from **agarose** by a purification process that removes the charged polysaccharides and gives a gel with only a very small number of residual, charged groups (see **Electroendosmosis**). Sepharose is stable in water and salt solutions over the pH range 4 to 9 and in the absence of oxidizing agents. It melts by heating above 40°C, and the bead structure may be irreversibly damaged by freezing. Cross-linked Sepharose, as prepared by reaction with 2,3-dibromopropanol under strongly alkaline conditions, is called Sepharose CL. It has greatly increased thermal and chemical stability. For example, Sepharose CL is used in aqueous media in the pH range 3 to 14.

Sephacryl is prepared by covalently cross-linking allyl dextran with *N,N'*-methylene bisacrylamide (see **Polyacrylamide**) and gives a rigid gel with a carefully controlled range of pore sizes. It is usually used in aqueous solutions. However, its gel structure allows replacing the water with organic solvents that have a much smaller effect on pore size than in the case of Sephadex. Sephacryl is insoluble in all solvents unless it is chemically degraded. It may be used in the range pH 2 to 11, with eluents containing **detergents** (eg, **SDS**), and with structure-breaking media, such as 6 M **guanidinium** chloride and 8 M **urea**. Sephacryl does not melt and may be autoclaved repeatedly at 120°C and neutral pH without significantly affecting its chromatographic properties.

### Bibliography

1. Pharmacia (1985) *Gel Filtration: Theory and Practice*, Rahms i Lund, Uppsala, Sweden.

## Sequence Analysis

Sequence analysis is devoted to the extraction of information contained in nucleic acids and proteins, and is one of the most important areas of **bioinformatics**. Many of the sequence-analysis methods are aimed at the prediction of functional and structural features in biological sequences. Given the availability of complete **genomes** for several organisms, and thus of their entire gene pools, prediction of individual protein structure, function, and interactions can now be complemented by methods that reconstruct entire biochemical pathways and other cellular processes.

Nevertheless, computational approaches are only beginning to contribute to the closure of the gap between genotypic and phenotypic information (1). As yet, too little is known about the signals in



genomic *DNA*, so elementary processes such as gene regulation cannot currently be predicted from sequence alone. Computational methods do better on the prediction of protein-coding genes, although this is a complicated process: Coding sequences comprise only 3% of the human genome. Therefore numerous, often weak, signals have to be combined (**promoters**, [RNA splicing](#) junction and branch sites, [polyadenylation](#) signals, etc.) in order to recognize exons and **introns** and to assemble the putative exons into complete genes (2). [Database](#) searches using sequence similarity are amongst the most powerful tools in this process, because ever more **homologues** of newly sequenced genes are found to be already stored in publicly available databases. However, because gene prediction methods currently only have an accuracy in the range 60% to 80% (3, 4), wrongly assembled proteins are accumulating in databases. Falsely annotated coding sequences, together with direct sequencing errors, frequently cause major problems in analysis of protein sequence data (5, 6).

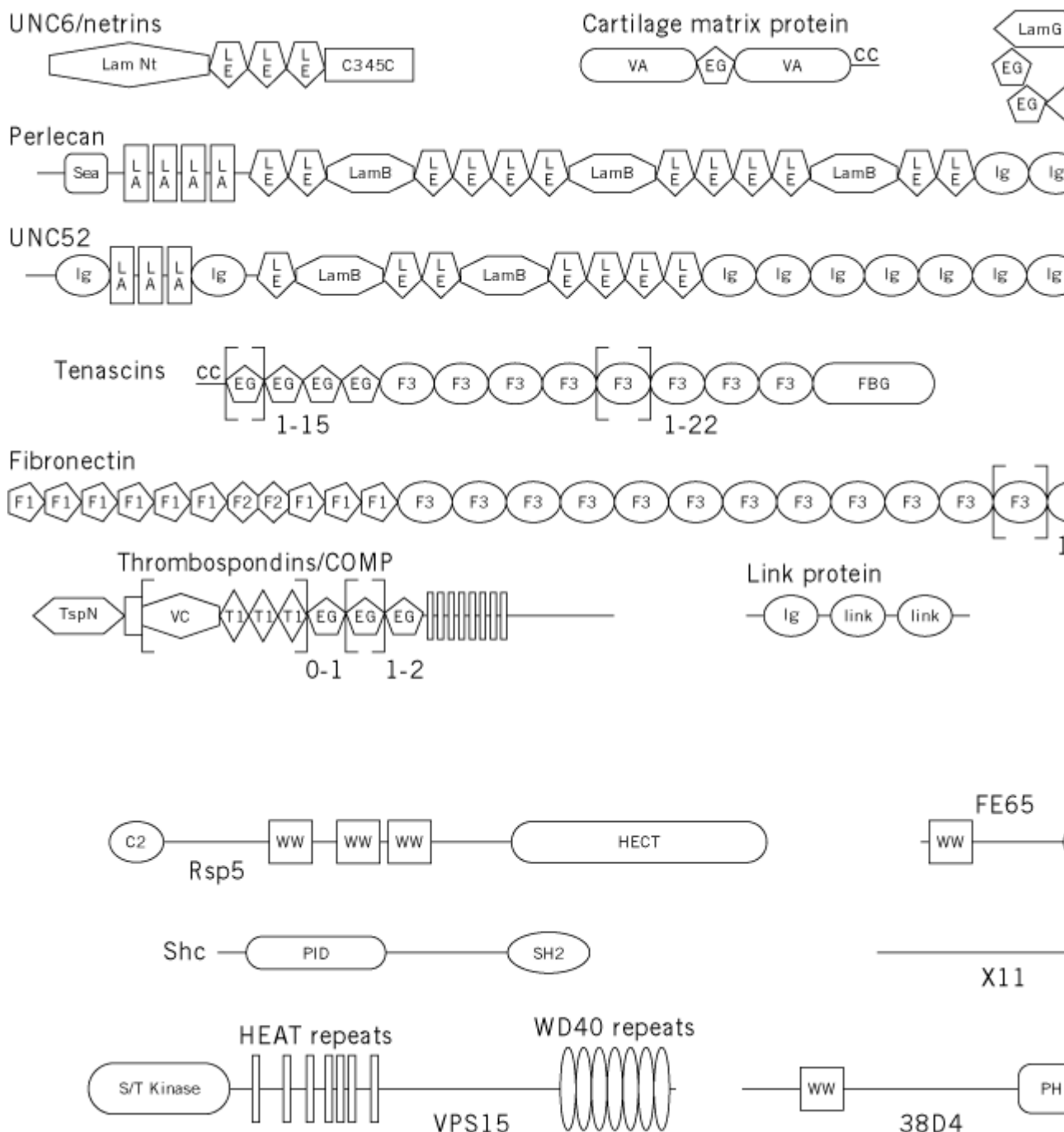
When a gene is identified and translated into a protein using the [genetic code](#), a variety of fast methods can be applied to identify its intrinsic features (4). Examples are the prediction of **transmembrane** regions and, more specifically, of **signal peptides**. Typically, these methods **thread sequences** of interest against precomputed weighting matrices (but more sophisticated algorithms may use trained **neural networks**, hidden Markov models, etc.). Other intrinsic features of proteins that should be characterized are nonglobular regions, internal repeats, **secondary structure**, and **posttranslational modifications**. Nonglobular regions of proteins include **coiled coil** or other fibrous elements, transmembrane regions, short internal repeats and segments of [random coil](#). These regions usually show a strong compositional bias in favor of certain amino acid residues (7). Nonglobular regions should usually be filtered out before a database search is carried out, because most algorithms are optimized for the identification of globular regions (8).

Sequence database searches are often the best way to obtain features indicative of individual protein function. If a homologue with a known [tertiary structure](#) can be identified, structural information can also be transferred from the hit to the query (see **Homology modeling**). Database searches are best done in an iterative procedure: If the first scan revealed some significantly similar proteins in the database, reciprocal searches with the clear homologues should be done and all proteins identified this way should be aligned together in order to utilize family information (9). Sequence conservation within an alignment is a valuable source of information about **divergence** of sequences and their [evolution](#), as well as of functional and structural features. More often than not, several homologues are found within one species; that is, the query sequence belongs to a [multigene family](#). Based on a multiple alignment (and given a certain expectation for the grouping of the sequences according to our current understanding of species evolution), [phylogenetic](#) trees can be built that help to discriminate **orthologous** from [paralogous genes](#)—that is, the equivalent genes in different species, as distinct from homologous but nonequivalent genes belonging to the same multigene family (10). Moreover, the conservation within aligned sequences can be used to perform much more sensitive database searches in order to identify distant relatives. For this purpose, the entire alignment can be used as query, using the full family information, although it is often useful to restrict the further searches to particularly conserved local regions, [protein motifs](#) (9, 11). A highly conserved motif often hints at a functional site, usually overlaid on a conserved structural unit. Motifs are often sufficient to predict the presence of a functional activity, so sequences can also be scanned against motif databases that contain signatures of well-characterized protein families (12).

A complication in the elucidation of individual protein function and structure is the modular architecture of many proteins; that is, they contain several functionally and structurally independent building blocks, called autonomous **domains**. Ideally, definition of domains should form the basis of all analyses of genes and proteins, but it is often hard to deduce the domain structure if only the sequence is known. A domain is best described in structural terms, because physical parameters can be applied (ie., ratio of diameter versus surface area) that provide clear thresholds for domain definitions (9, 11). The average size of proteins in databases is more than 300 amino acid residues, yet globular domains only rarely exceed 200 residues in length; hence the majority of proteins contain more than one domain. If proteins are known to harbor a certain domain, or such a domain has been identified in a database search, it is advisable to perform separate searches with the

remaining parts of the proteins. Homologous domains can be found in a diverse set of proteins with unrelated overall functions, as exemplified in Figure 1. Therefore, genetic mechanisms must have existed (and presumably still do) that allowed their horizontal spread within genomes. For animals, [exon shuffling](#) appears to be one such mechanism, especially with regard to extracellular proteins (13), but, because it requires introns, this cannot be the driving force for the fast evolution of certain modular bacterial proteins (especially in soil bacteria), many of which show **horizontal transmission** between genomes. Both plasmid transfer and efficient [recombination](#) mechanisms are likely to contribute to the horizontal transfer of domains between bacterial genomes.

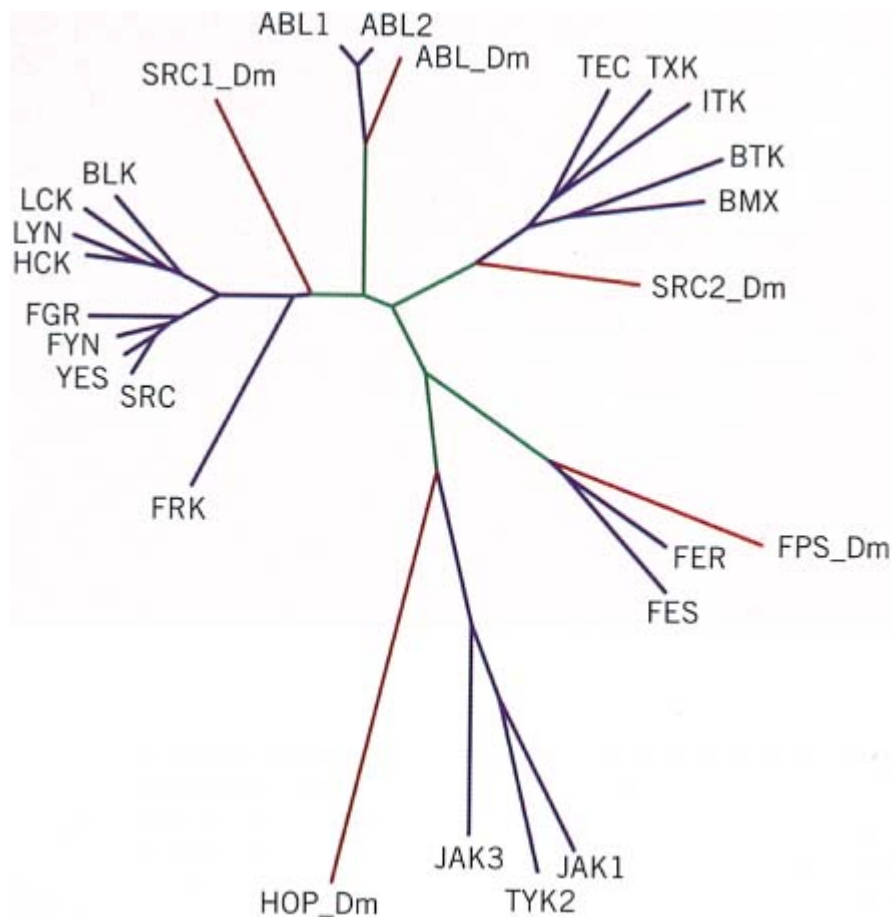
**Figure 1.** A selection of multidomain proteins. Most eukaryotic extracellular proteins are partly or wholly assembled from more than 50 widespread cytoplasmic signalling domains known to date (31). In the nucleus, there is an extensive chromatin-associated domains. These sets of domains are together estimated to occur in more than 20,000 human proteins.



The domain phenomenon also has to be considered when comparing whole genomes, with the aim of identifying orthologues and thence reconstructing metabolic pathways. This has become feasible now that several complete bacterial genomes are available, while, especially for *Escherichia coli*, a vast amount of biochemical knowledge has also been generated. The comparative analysis of entire genomes is rather new, though progressing quickly (14-16). Orthologues from different species are collected, and newly incoming bacterial genomes can be scanned against those collections for the purpose of annotation. Databases of known pathways are becoming linked to sequence information and apparent gaps in, or the absence of, certain pathways have then to be interpreted. Secondary metabolism in a bacterium can then be confined to a restricted number of “leftover genes” that are not present in other bacteria.

In addition to following the evolution of genes and proteins, it is now becoming possible to analyze the evolution of entire genomes. One important result to have emerged from extensive sequence analyses is that the vertebrate genome is essentially octaploid (a fourfold genome amplification, because somatic cells are already diploid). This would be accounted for by two consecutive genome duplications in the vertebrate common ancestral lineage. It is routinely found that a given gene in the model organism *Drosophila* has up to four counterparts in human or mouse, called tetralogues, because, although they are paralogues, they are all orthologous with respect to *Drosophila* (17). In some cases, individual tetralogues may have been lost by deletion, as is documented within the *Hox* gene clusters. Figure 2 shows the situation for a set of human and *Drosophila* tyrosine kinases. As a consequence of **polyploidy**, tetralogous genes originally had an essentially total functional overlap. Although their functions may have diverged subsequently (especially for extracellular proteins), there are many instances of vertebrate functional redundancy, most clearly indicated by the effect of mouse gene knockouts with notoriously small phenotypic effects—for example, for *Src*, MyoD, RXR, and so on. Functional redundancy enormously complicates analyses of, for example, signaling pathways and gene regulation: It becomes essential to obtain the full set of redundant genes and proteins. Until such time as a complete vertebrate genome sequence is determined, the [expressed sequence tag](#) (EST) databases [large collections of randomly sequenced [complementary DNAs](#) (18)] are a vital resource for identifying tetralogues.

**Figure 2.** Neighbor-joining tree prepared from an alignment of human and *Drosophila* tyrosine kinases, showing multiple human kinases for each related *Drosophila* kinase. Branches leading to human kinases are colored blue; to *Drosophila* kinases, red; with green for branches that precede the split of the human–*Drosophila* lineages. The least divergent human groups, such as the ABL and SRC groups, are likely to have the most conserved and overlapping functions. The more diverged BTK group has incompletely overlapping functions. *Drosophila* orthologues were not yet known for the SRC and LYN tetralogies. Local chromosomal gene duplications may account for some of the multiple human genes; thus the BTK group exceeds the tetralogy with five known members, but BTK and BMX are close together on the X chromosome, implying a relatively recent duplication. See color insert.



Vertebrate tetraploidy is only one of very many polyploidies that have occurred in eukaryotic evolution. Polyploidy is probably one of the main routes to the establishment of multigene families (although not the only one) (19). For example, the full sequence for the genome of Baker's yeast revealed traces of an ancient genome duplication, estimated to have occurred 100 million years ago (20). Many other medium- to large-scale genetic events will also be revealed as more genomes become completely sequenced. Genetic mechanisms leading to inversions, slippage, plasmid operon mobility, and so on, will in future be better detected and quantified.

As the amount of sequence data increases, analyses rely more and more on sophisticated [databases](#) and knowledge retrieval systems. Therefore, adequate computer hardware and network facilities are becoming essential even to experimental biologists. Amongst the academic disciplines, biologists are indisputably the single greatest beneficiaries of modern computer networking developments such as the World Wide Web (WWW), which has allowed legions of web servers to become accessible to any biologist with a personal computer and telephone line. Because of continual development in sequence analysis tools, many programs and web servers in common use at the time of writing are likely to be replaced in the near to medium term. Therefore only those that the authors regard as having longer-term utility, or particular current importance, will be mentioned by name in this section.

## 1. Searching Sequence Databases

As well as the sequences themselves, the databases contain annotation that can provide very useful information about the sequences—and increasingly, cross-links to other forms of databases (or other biological resources)—which may not be based on primary sequence information—for example, tertiary structure databases or medical databases. Therefore tools have been developed for two vital

forms of database retrieval; firstly, by sequence similarity and, secondly, by keyword search through the text annotation. In practice, sequence-based searches are likely to induce information extraction, and vice versa, so the tools are highly complementary.

### 1.1. Search by Sequence Similarity

Underlying similarity searches is the concept of *homology*, which implies that sequences have diverged from common ancestral genes and genomes; hence traces of the evolutionary history may still be seen in the present sequences. Typically, both functional and structural constraints apply to the sequences, so that conservation within a set of related sequences mirrors those constraints, and current database search methods tend to be optimized to identify these signals.

The sensitivity of the search methods is inversely related to their computational speeds, which are primarily determined by the algorithm that is used for [aligning sequences](#). Searches for identical sequence matches are useful for checking whether a newly isolated cDNA or protein is already known from other experimental approaches. For (nearly) identical sequences, very fast methods can be used to find database matches. Searches for more divergent sequence matches are undertaken in the hope of finding related genes or proteins that can shed light on the possible functions of the query sequence. The fast search programs of the BLAST (8, 21) and FASTA (22) series, which first find ungapped matches and then use more sensitive routines in matching regions, may be accessed via the web but are fast enough to be installed and run on workstations. More sensitive searches using single sequences or multiple alignments [as profiles or hidden Markov models (HMMs)] for the input query, as well as dynamic programming algorithms such as Smith–Waterman (23) for the search, are mostly used in central facilities because of their larger CPU requirements. Fortunately, the most sensitive alignment methods are increasingly being ported to fast dedicated machines that can be accessed by web servers.

### 1.2. Statistical Analysis of DB Matches

Many programs provide statistics for the significance of matches in the search output (24, 25). These can be helpful for assessing whether a match is true or false and are essential for automatic large-scale analyses, in order to be able to apply thresholds and cutoffs. Common indicators of statistical significance are  $E$  and  $p$  values ( $E$ , expected ratio of false positives;  $p$ , probability of a chance match; they can be converted into each other using  $p = 1 - e^{-E}$ ). While they are a useful indication for the significance of a hit (8, 21), it is important to be aware that they are based on assumptions that very often do not hold. For example, they are affected by choice of amino acid mutation matrix and assume an unbiased amino acid composition in the database and that both query and database proteins are devoid of low-complexity regions. Also, the values are dependent on the size of the database so that, as data accumulate, the increased noise level causes increased  $p$  and  $E$  values for a given hit over time. On the other hand, it should be noted that there is a compensatory effect filling the finite space of real sequences: For any given protein, the chance of a database containing a sufficiently close and detectable homologue steadily increases.

HMM-based searches are able to provide intrinsic probabilities that a given match is related to the query. This is expected to become very useful (26). Nevertheless, whether these values are more robust than the current  $E$  and  $p$  values will depend on the approximations in the HMM. For example, an HMM that does not have a good model for reduced sequence complexity will be just as error prone as any other statistic.

### 1.3. Retrieval of Information from Databases

Modern sequence analysis packages contain integrated retrieval systems that try to add biologically relevant information to the hit detected by database sequence similarity searches. Biological context often helps in the evaluation of weak similarities. Species information and functional information, including cellular localization and tissue expression patterns, are examples of the useful features that add value to the sequence data itself.

Molecular biological databases originated as stores for DNA and protein sequences and

macromolecular structures. However, increasingly diversified databases reflecting different functional criteria are appearing and will continue to proliferate. Databases for single organisms, for mutations causing human inherited diseases, for biochemical pathways, and for developmental pathways are obvious topics. The most practically useful databases must include extensive cross-indexing to other classes of databases. In this respect the Swiss-Prot protein sequence database, with links to many other databases, provides a pioneering example of the standards that need to be met (27).

Efficient retrieval of database entries requires fast text matching and keyword indexing strategies. These biological needs are providing a stimulus for computational methods development in database retrieval, as exemplified by the Sequence Retrieval System (SRS), which combines efficient indexing, a flexible parsing language, and extensive database cross-links, to provide a powerful general database query tool (28). SRS, installed at many web sites, and the ENTREZ and PubMed servers at the NCBI (29) together provide universal network access to many different biological databases.

#### 1.4. World Wide Web

WWW servers have been set up for most forms of sequence analysis. They offer advantages to the biological community in that central sites can offer services to users everywhere. Web servers may be at places where analysis programs are being actively developed and thus offer the latest, most sensitive methodology and/or they may be sites that are well equipped with powerful computers, so that computationally demanding jobs can be offered. However, the WWW also has some disadvantages. There is an increase in system complexity: web interaction requires two computers, a network, and compatible web software all to be functioning. To run stably, web sites have a high maintenance load. They may disappear when the person responsible leaves an institute. Stable servers are most likely to be found at the larger bioinformatics institutes, such as the NCBI in the United States, the EBI and EXPASY in Europe, and the many sites providing facilities at the national level. Some tasks are poorly suited to web servers, which must turn around jobs rapidly to allow heavy access. For example, web servers offering multiple-sequence alignment tend to be highly restrictive, because the calculation time progresses arithmetically with respect to the number of sequences and, worse still, increases proportionately to the product of the sequence lengths. A web site cannot afford to be blocked by giant alignment tasks, so these have to be done on a local machine. Such disadvantages mean that it is likely that local computational resources will continue to play an important role in sequence analysis, and users will still need to install and run analysis programs directly.

#### Bibliography

1. P. Bork, C. Ouzounis, and C. Sander (1994) *Curr. Opin. Struct. Biol.* **4**, 393–403.
2. R. Guigo (1997) *J. Mol. Med.* **75**, 289–293.
3. M. Burset and R. Guigo (1996) *Genomics* **34**, 353–367.
4. C. Burge and S. Karlin (1997) *J. Mol. Biol.* **268**, 78–94.
5. P. Bork and A. Bairoch (1996) *Trends Genet.* **12**, 425–427.
6. E. Birney, J. D. Thompson, and T. J. Gibson (1996) *Nucleic Acids Res.* **24**, 2730–2739.
7. J. C. Wootton and S. Federhen. (1996) *Methods Enzymol.* **266**, 554–573.
8. S. F. Altschul, M. S. Boguski, W. Gish, and J. C. Wootton (1994) *Nature Genet.* **6**, 119–129.
9. P. Bork and T. J. Gibson (1996) *Methods Enzymol.* **266**, 162–184.
10. W. M. Fitch (1970) *Syst. Zool.* **19**, 99–106.
11. P. Bork and E. V. Koonin (1996) *Curr. Opin. Struct. Biol.* **6**, 366–376.
12. A. Bairoch and P. Bucher (1994) *Nucleic Acids Res.* **22**, 3583–3589.
13. L. Patthy (1994) *Curr. Opin. Struct. Biol.* **4**, 383–392.
14. R. L. Tatusov, A. R. Mushegian, P. Bork, N. P. Brown, M. Borodovsky, W. S. Hayes, K. E. Rudd, E. V. Koonin, et al. (1996) *Curr. Biol.* **6**, 279–291.

15. A. R. Mushegian and E. V. Koonin (1996) Proc. Natl. Acad. Sci. USA **93**, 10268–10273.
16. E. V. Koonin, A. R. Mushegian, M. Y. Galperin, and D. R. Walker (1997) Mol. Microbiol. **25**, 619–637.
17. J. Spring (1997) FEBS Lett. **400**, 2–8.
18. M. D. Adams et al. (1991) Science **252**, 1651–1656.
19. S. Ohno (1970) *Evolution by Gene Duplication*, George Allen and Unwin, London.
20. K. H. Wolfe and D. C. Shields (1997) Nature, **387**, 708–713.
21. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman (1990) J. Mol. Biol. **215**, 403–410.
22. W. R. Pearson (1995) Protein Sci. **4**, 1145–1160.
23. T. F. Smith and M. S. Waterman (1981) Adv. Appl. Math. **2**, 482–489.
24. J. F. Collins and A. F. W. Coulson (1990) Methods Enzymol. **183**, 474–487.
25. M. S. Waterman and M. Vingron (1994) Proc. Natl. Acad. Sci. USA **91**, 4625–4628.
26. S. Eddy (1996) Curr. Opin. Struct. Biol. **6**, 361–365.
27. A. Bairoch and R. Apweiler (1997) Nucleic Acids Res. **25**, 31–36.
28. T. Etzold, A. Ulyanov, and P. Argos (1996) Methods Enzymol. **266**, 114–128.
29. G. D. Schuler, J. A. Epstein, H. Ohkawa, and J. A. Kans (1996) Methods Enzymol. **266**, 141–161.
30. P. Bork, K. A. Downing, B. Kieffer, and I. D. Campbell (1996) Q. Rev. Biophys. **29**, 119–167.
31. P. Bork, J. Schultz, and C. P. Ponting (1997) Trends Biochem. Sci. **22**, Supplement TIBS CO4.

### Suggestions for Further Reading

32. M. J. Bishop and C. J. Rawlings (1987) *Nucleic Acid and Protein Sequence Analysis: A Practical Approach*, IRL Press, Oxford, UK.
33. R. F. Doolittle (1990) *Molecular Evolution: Computer Analysis of Protein and Nucleic Acid Sequences, Methods in Enzymology*, Vol. **183**, Academic Press, San Diego, CA.
34. R. F. Doolittle (1996) *Computer Methods for Macromolecular Sequence Analysis, Methods in Enzymology*, Vol. **266**, Academic Press, San Diego, CA.
35. A. M. Griffin and H. G. Griffin (1994) *Methods in Molecular Biology: Computer Analysis of Sequence Data*, Humana Press, Totowa, NJ.

### Sequence Codes

**Nucleotide sequences** do not only encode the [primary structures](#) of [proteins](#). Numerous sequence-directed processes and sequence-dependent structures and interactions also exist. The corresponding sequence instructions are “read” from the DNA or RNA molecule, each in its own way, via one or another specific molecular interaction or a whole network of such interactions. In the case of the triplet [genetic code](#), the reading device is the [ribosome](#), in a complex combination with [transfer RNAs](#) and various protein factors (see **Protein biosynthesis**). In the case of [gene splicing](#), the sequence signals are recognized by the [spliceosome](#), which, in addition to the nuclear RNA transcript, involves various small nuclear RNAs and proteins (see [RNA Splicing](#)). Numerous, relatively simple sequence-specific DNA-protein and RNA-protein interactions also exist, in which the respective sequences are read by a single protein (see [DNA-Binding Proteins](#) and [RNA-Binding](#)

[Proteins](#)).

The existence of codes other than the classic [translation](#) triplet code is already suggested by the degeneracy of the triplet code. Eight [amino acids](#) are actually encoded by doublets, with the third bases of the codons being arbitrary, ie, fully degenerate. Seven amino acids are encoded by triplets, with the third base being either C or U, and five amino acids have degenerate A or G in the third positions of their triplets. Three of the amino acids (Leu, Ser, and Arg) are encoded by two sets of triplets each, with two different doublets in the first and second positions. Only two residues, Met and Trp, are each encoded by a single, nondegenerate triplet. Such freedom in the choice of codons allows significant changes in the nucleotide sequence without changing the encoded protein sequence. This, in principle, makes it possible to utilize the interchangeable bases of the mRNA sequence for purposes, other than protein coding, eg, for some additional, different sequence codes (1).

Despite their traditional classification as noncoding, sequences other than those for protein coding carry some important information as well (2). Among numerous examples are cases of sequence conservation. Comparison of human and mouse [T-cell receptor](#) alpha locus sequences—about 80 kbp each, both largely noncoding—shows greater than 70% similarity (3). High conservation is also observed in the 5' and 3' untranslated regions, as well as in the **introns** (4). Such conservation of the noncoding sequences suggests that they actually do encode some message(s).

Each sequence-specific interaction, either complex or simple, is represented by a specific sequence pattern. Visual or computer recognition of a given pattern in the sequence allows location of corresponding functional sites (regions) in the long DNA or RNA molecules. The pattern therefore serves as a code for recognition of the sites in the sequences. The code can be defined as a sequence pattern that corresponds to one or another biological or biomolecular function. The code can be general, even universal, of broad use in various species, or specific, eg, species or gene specific. For example, the triplet code is general, but the sequence pattern responsible for recognition of *Escherichia coli* DNA by [Lac repressor](#) represents a specific code.

Among known general sequence codes other than the triplet code are [transcription](#) signals in **promoters**, such as TATAAA-box in eukaryotes and TATAAT- and TTGACA-boxes in bacteria (see [TATA Box](#)), coding for the initiation of transcription. Another widely known sequence code is the gene-splicing code, the GT-AG rule (see [RNA Splicing](#)) and some sequence preferences around the intron-exon junctions. Eukaryotic sequences around these junctions resemble to varying degrees the **consensus sequences** agGUaag and y<sub>10-15</sub>ncAG, where the capitalized nucleotides are most conserved, n is any nucleotide, and y is a pyrimidine. The splicing code is likely to contain some other as-yet undetected sequence features, as the current computer recognition algorithms based on the above patterns are not as precise as the natural gene-splicing process.

A complex set of sequence rules describes details of the shape of DNA that are important for DNA-protein interactions and DNA folding in the cell. The DNA structure is not monotonously uniform. It is modulated by sequence-dependent local deviations from the standard geometry, which, for example, may accumulate to a net DNA curvature. The first manifestation of the intrinsic DNA curvature was observed in the predicted 10.5 base sequence periodicity. Local deflections of the DNA axis at certain dinucleotide steps, primarily at AA(TT) dinucleotides, repeated at the DNA helical pitch distance, result in an accumulating unidirectional deflection of the DNA axis, making it curved. The first experimental indication of the existence of DNA curvature was the anomalous **electrophoretic** mobilities of certain DNA fragments (5). More direct experimental evidence was provided by [electron microscopy](#) of the curved DNA (6) and by making small circles of DNA with designed periodic sequences (7). The sequence rules to describe the DNA curvature or the generally nonplanar sequence-dependent **writhe** of DNA are not yet fully developed. The simplest description is based on the original wedge model. The geometry of every base-pair step in this model is described by three angles: wedge roll, wedge tilt, and twist. The full set of these angles for all



dinucleotide steps has been estimated from the available experimental data. Following the sequence and deflecting the DNA axis at every step according to the wedge and twist angles from the table of the dinucleotide “codons,” one can calculate the predicted path of the DNA axis for any given sequence.

The [chromatin](#) code is a set of rules directing sequence-specific positioning of the [nucleosomes](#). Although nucleosome positioning may also involve several other factors (8), the sequence plays a substantial if not a major role. Sequence-dependent deformational anisotropy (bendability) of DNA appears to be an underlying principle of the nucleosome sequence specificity. The anisotropy can be ascribed to certain sequence elements periodically distributed along the nucleosome DNA, following the DNA helical repeat. In particular, these are dinucleotides AA and TT, as in the case of DNA curvature, however, with different phase relationships between the dinucleotides (9). Other dinucleotides, such as CC(GG) (10), TA, and perhaps some additional sequence elements (11), also contribute to the nucleosome sequence pattern. Further elucidation of the chromatin sequence code is important for understanding the role of the nucleosomes in gene expression (12).

A unique property of the nucleotide sequences is the superposition of the codes they carry—that is, a given base in a given position along the sequence may be involved simultaneously in several messages of the same or a different nature as was theoretically predicted in 1971 (1). For example, the coding sequences of the genome of hepatitis B virus massively overlap, so that the same regions of the sequence simultaneously code for two different proteins with different amino-acid sequences (13). The promoters of 5S rRNA genes are located within the genes themselves (14). The sequences at the ends of the exons are involved not only in the protein coding but also in the gene splicing sequence patterns (15). Many other cases of such overlapping are known; these are only immediate simple examples of a general phenomenon of superposition of many different sequence codes. Such overlapping is possible owing to degeneracy of the codes, well exemplified by the triplet code, so that alternative base combinations may be used simultaneously to satisfy several superimposed messages. Of course, an informational limit for such superposition, should exist, in which the degeneracy becomes insufficient to accommodate all messages without loss of quality. One current theory, by Zuckerkandl, on the nature of the intervening sequences is based on the notion of spatial separation of two otherwise conflicting sequence messages—protein-coding pattern and sequence determinants of the chromatin structure (2). The intervening sequences are proposed to have been introduced to take the load of the chromatin code, whereas the exons would remain primarily responsible for the protein-coding function.

## Bibliography

1. T. Schaaap (1971) *J. Theor. Biol.* **32**, 293–298.
2. E. Zuckerkandl (1986) *J. Mol. Evol.* **24**, 12–27.
3. B. F. Koop and L. Hood (1994) *Nature Genet.* **7**, 48–53.
4. L. Duret, F. Dorkeld and C. Gautier (1993) *Nucl. Acids Res.* **21**, 2315–2322.
5. J. C. Marini et al. (1982) *Proc. Natl. Acad. Sci. USA* **79**, 7664–7668.
6. J. M. Griffith et al. (1986) *Cell* **46**, 717–724.
7. L. Ulanovsky et al. (1986) *Proc. Natl. Acad. Sci. USA* **83**, 862–866.
8. R. Kornberg and Y. Lorch (1995) *Curr. Opin. Cell Biol.* **7**, 371–375.
9. I. Ioshikhes et al. (1996) *J. Mol. Biol.* **262**, 129–139.
10. A. Bolshoy (1995) *Nature Struct. Biol.* **2**, 446–448.
11. P. T. Lowary and J. Widom (1998) *J. Mol. Biol.* **276**, 19–42.
12. R. E. Kingston, C. A. Bunker and A. N. Imbalzano (1996) *Genes Dev.* **10**, 905–920.
13. R. H. Miller et al. (1989) *Hepatology* **9**, 322–327.
14. S. Sakonju, D. F. Bogenhagen and D. D. Brown (1980) *Cell* **19**, 13–25.
15. M. S. Gelfand (1992) *J. Mol. Evol.* **35**, 239–252.

### Suggestions for Further Reading

16. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.
17. E. N. Trifonov (1996), Interfering contexts of regulatory sequence elements. *CABIOS* **12**, 423–429. Review of various codes.

### Sequence Databases

A sequence database is a repository of all publicly available **nucleotide sequences** or amino acid sequences that are experimentally determined by [DNA sequencing](#) and [protein sequencing](#) methods. The primary nucleotide sequence databases, which comprise the International Nucleotide Data Banks, are GenBank by the U.S. National Center for Biotechnology Information, EMBL by the European Bioinformatics Institute, and DDBJ (DNA Data Bank of Japan) by the Japanese National Institute of Genetics. These three databases have virtually the same content, for they daily exchange data received from individual researchers and large sequencing centers. The primary protein sequence databases are SWISS-PROT by the University of Geneva and Protein Information Resource (PIR) by the National Biomedical Research Foundation in Georgetown University, and its collaborators in Germany and Japan. The contents of these two are not the same; SWISS-PROT is a well-curated database, whereas PIR contains a [superfamily](#) classification. Because the majority of protein sequences are now determined from DNA sequences, the protein databases also rely on direct submissions to the Nucleotide Data Banks. A third database Protein Research Foundation (PRF) in Japan also covers peptide and protein sequences reported in the literature. The sequence databases can be viewed as a collection of entries, each of which is uniquely identified by its entry name or accession number. The content of an entry generally consists of definition, source, references, and other text information, a table describing the biological features of the sequence data, and the information about the sequence itself. Because of the expansion of the sequencing of entire [genomes](#) and of [complementary DNA](#) (cDNA), most of the databases are updated daily and released over the Internet. The WWW address is shown in Table 1 for each of the organizations that produce the primary sequence databases.

**Table 1. WWW Addresses for the Sequence Databases**

| Database   | Address  |
|------------|--|
| GenBank    | <a href="http://www.ncbi.nlm.nih.gov">www.ncbi.nlm.nih.gov</a>       |
| EMBL       | <a href="http://www.ebi.ac.uk">www.ebi.ac.uk</a>                     |
| DDBJ       | <a href="http://www.ddbj.nig.ac.jp">www.ddbj.nig.ac.jp</a>           |
| SWISS-PROT | <a href="http://expasy.hcuge.ch">expasy.hcuge.ch</a>                 |
| PIR        | <a href="http://www-nbrf.georgetown.edu">www-nbrf.georgetown.edu</a> |
| PRF        | <a href="http://www.prf.or.jp">www.prf.or.jp</a>                     |

Amino acid sequence determination became possible when Frederick Sanger sequenced [insulin](#) in 1953 (1). In the mid-1960s, Margaret O. Dayhoff initiated an effort to collect protein sequence data to study molecular [evolution](#). Her collection was not fully computerized but was published in the series *Atlas of Protein Sequence and Structure* between 1968 and 1978 (2). The PIR database was established in 1984 as a descendant of her collection. In 1977, Sanger published the first complete genome sequence of a tiny virus, bacteriophage  $\phi$ X174 (3). This was the beginning of the era of DNA sequencing. The sequence databases of GenBank at the Los Alamos National Laboratory and EMBL at the European Molecular Biology Laboratory were officially launched in 1982, and DDBJ joined in 1984. Over the years, these databases have undergone transitions to cope effectively with ever increasing amounts of sequence data. Originally, it was the databases' responsibility to find, enter, and annotate published sequences, but now it is the authors' responsibility to annotate and submit the data, and to obtain the accession number because many journals require data submission as a precondition of publication. The main use of the sequence databases is [homology](#) searches to retrieve related sequences (see [Bioinformatics](#)). Because there are multiple databases for the same type of data and because there are identical sequences in the same database, attempts are being made to derive computationally a nonredundant set of sequences from multiple sources. The sequence databases also form a backbone resource from which other value-added [databases](#) are developed.

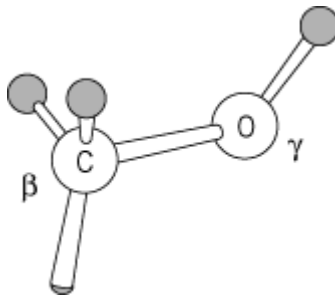
#### Bibliography

1. F. Sanger and H. Tuppy (1951) *J. Biochem.* **49**, 463–481 and 481–490; F. Sanger and E. O. P. Thompson (1953) *J. Biochem.* **53**, 353–366 and 366–374.
2. M. O. Dayhoff (1978) *Atlas of Protein Sequence and Structure*, Vol. **5**, Supplement 3, National Biomedical Research Foundation, Washington, D.C.
3. F. Sanger et al. (1977) *Nature* **265**, 687–695.

#### Serine (Ser, S)

The [amino acid](#) serine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to six **codons**—UCU, UCC, UCA, UCG, AGU, and AGC. This is unique, in that the first four codons are not interchangeable with the last two by single base substitution mutations, so there are two types of serine codons. Ser represents approximately 6.9% of the residues of the proteins that have been characterized. The seryl residue incorporated has a mass of 87.08 Da, a **van der Waals volume** of  $73 \text{ \AA}^3$ , and an [accessible surface](#) area of  $122 \text{ \AA}^2$ . Ser residues are changed frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [alanine](#), [threonine](#), [glycine](#), and [asparagine](#) residues.

The side chain of Ser is dominated by its hydroxyl group:



This hydroxyl group is normally no more reactive chemically than is that of ethanol, so there are few chemical reactions that can modify Ser residues in a protein specifically and readily. The only useful general reaction is acetylation with acetyl chloride in aqueous trifluoroacetic acid. When in the appropriate environment, however, and with a potentially reactive group held in the correct position, the Ser hydroxyl group can function as a potent nucleophile, as occurs in the [serine proteinases](#) .

The hydroxyl group is very polar and [hydrophilic](#) and can function as either a donor or acceptor in [hydrogen bonds](#). Consequently, Ser residues can be fully buried in protein structures, and about 22% are, especially if their hydroxyl group is paired in a hydrogen bond. The Ser hydroxyl group is situated sterically to interact with polar groups of the polypeptide backbone, which affects its conformation and reactivity. For example, peptide bonds adjacent to Ser residues are especially susceptible to acid hydrolysis. Ser residues also moderately favor the **alpha-helical** conformation in model peptides, and they are particularly adept in serving as the terminal “capping” residue at the end of the  $\alpha$ -helix, when the hydroxyl group of the side chain participates in hydrogen bonding with the backbone. In folded [protein structures](#), Ser residues occur most frequently at reverse [turns](#), again because the hydroxyl group can hydrogen bond with the backbone.

In native proteins, Ser residues are frequently **post-translationally modified** by the addition of oligosaccharides (see [O-Glycosylation](#)) and phosphate groups.

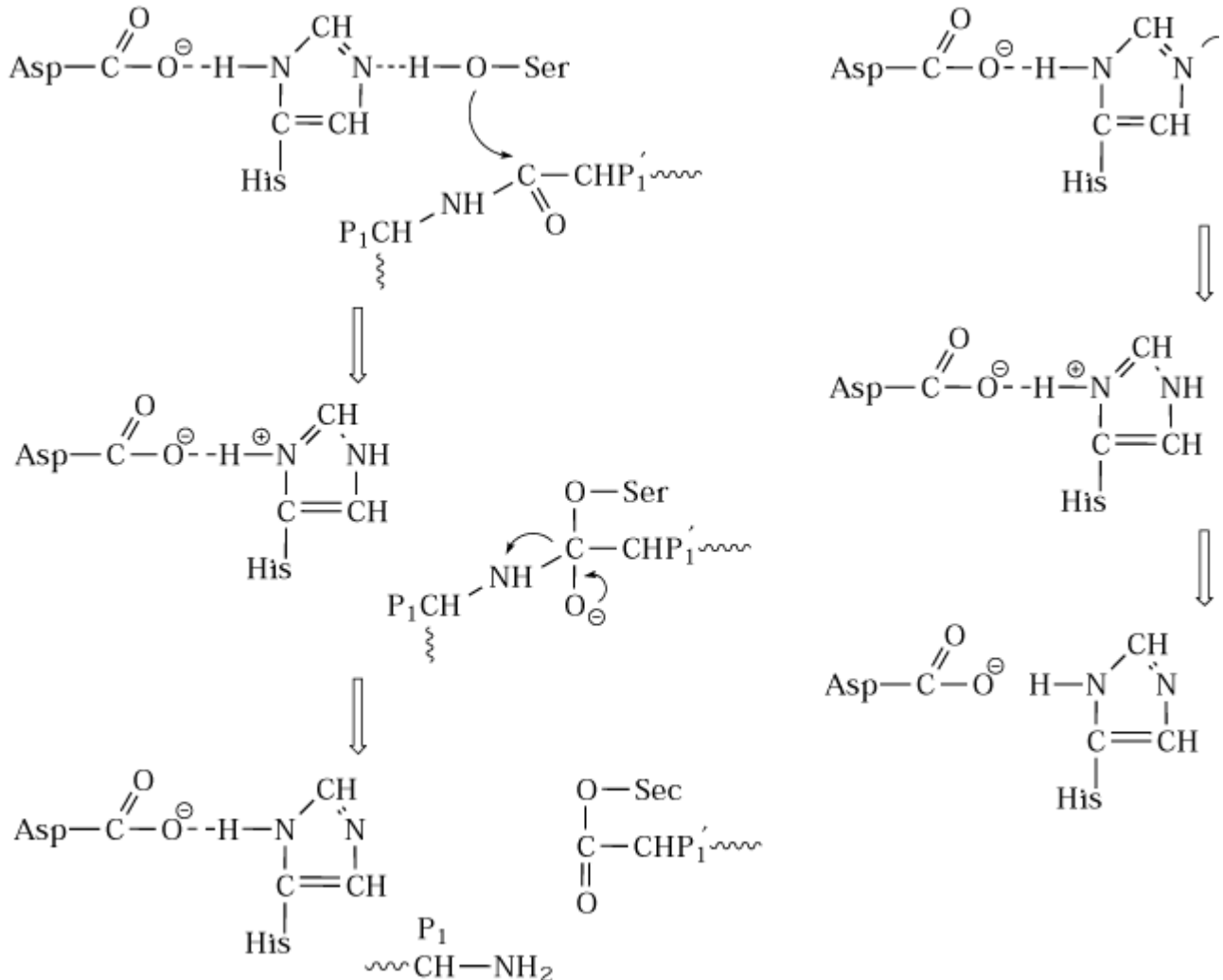
#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Serine Proteinase

The term *serine proteinase* defines a class of proteolytic [enzyme](#) (one that catalyzes the hydrolysis of [peptide bonds](#) in proteins) in which the hydroxyl group of a [serine](#) residue participates in the catalytic process (E.C. 3.4.21) (1). The three-dimensional structure of the serine proteinase (2) brings together what is called a “[catalytic triad](#),” wherein a carboxyl group of an [aspartic acid](#) residue is aligned with the imidazole side chain of a [histidine](#) residue, which in turn interacts with, and thereby activates, the serine hydroxyl group. When a substrate (protein or peptide) binds to the [active site](#) of the serine proteinase (the region of the enzyme where catalysis occurs), the serine hydroxyl group interacts with the carbonyl group of the peptide bond that is to be hydrolyzed and forms an acylenzyme (Fig. 1). Subsequent hydrolysis of this acylenzyme involves the imidazole side chain again and a water molecule, and this completes the proteolytic event. Members of this family are also referred to as *alkaline proteinases* because they tend to act optimally at a pH close to 8.

**Figure 1.** Schematic representation of the catalytic reaction of serine proteinases. On the left, the active-site serine residue and histidine components of the catalytic triad, attacks the carbonyl carbon in the peptide bond of the substrate. In a subsequent step, the particular peptide attacked is determined by the ability of the enzyme to recognize  $P_1$  and  $P_1'$ , which are the particular amino acid side chains. The resulting tetrahedral intermediate is converted to an acyl-enzyme intermediate, with formation of the carbonyl carbon of the acyl intermediate is attacked by an activated water molecule to form another tetrahedral intermediate, with regeneration of the starting enzyme.



There are two evolutionary families of serine proteinases. One includes the notable examples of **trypsin**, **chymotrypsin**, and **elastase**, digestive enzymes secreted by the pancreas. These three proteins have rather similar amino acid sequences, indicative of **divergent evolution** from a common ancestor, and almost identical three-dimensional structures. They act on a wide range of proteins, but each has a unique specificity: Trypsin cleaves bonds where a **lysine** or **arginine** residue contributes the scissile carbonyl group; chymotrypsin requires an aromatic or bulky **hydrophobic** side chain; and elastase acts only adjacent to small, uncharged side chains.

The second class of serine proteinases includes bacterial proteins related to **subtilisin**. They have completely different amino acid sequences from the trypsin family and are not **homologous** to them. Nevertheless, they have the exact same “catalytic triad” arrangement. This is a clear indication of what is known as **convergent evolution**.

Serine proteinases are involved in **blood clotting**, where activation and localization of function

require exquisitely fine tuning. Indeed, the presence in plasma of the [serpin](#) class of **serine proteinase inhibitors** (3) to prevent clotting throughout the vasculature is essential for hemodynamic homeostasis. On the other hand, another serine proteinase, tissue plasminogen activator, can be useful in treating heart attack and stroke victims because it catalyzes the formation of **plasmin**, yet another serine proteinase that dissolves unwanted blood clots.

In addition to protein inhibitors, serine proteinases are characteristically inhibited by a variety of reagents that chemically either block the reactive serine hydroxyl group or modify the histidine imidazole group. The classic serine reagent is diisopropylfluorophosphate (**DIPF**), an early nerve gas whose action originally defined serine proteinases. Phenylmethane sulfonyl fluoride (see [PMSE \(Phenylmethylsulfonyl Fluoride\)](#)) is easier to work with and has become widely accepted as a serine proteinase inhibitor. Chloromethyl ketones (see [TLCK \(N-P-Tosyl-Lysine Chloromethyl Ketone\)](#) or [TPCK \(N-P-Tosyl-Phenylalanine Chloromethyl Ketone\)](#)) have been designed to be inhibitors of trypsin and trypsin-like serine proteinases or chymotrypsin and chymotrypsin-like serine proteinases, respectively. More recently, 3,4-dichloroisocoumarin has been proposed as a very potent catalytic mechanism-based inhibitor of all serine proteinases (4).

#### Bibliography

1. R. M. Stroud (1974) *Sci. Am.* **231**, 24–88.
2. D. M. Blow (1976) *Acc. Chem. Res.* **9**, 145–152.
3. R. Carrell and J. Travis (1985) *Trends Biochem. Sci.* **10**, 20–24.
4. J. W. Harper, K. Hemmi, and J. C. Powers (1985) *Biochemistry* **24**, 1831–1841.

#### Serine Proteinase Inhibitors, Protein

Serine proteinase inhibitors, protein are arguably the most intensely studied of protein enzyme inhibitors. Serine proteinase inhibitors are proteins that inhibit only enzymes that are serine proteinases [see [Proteinase Inhibitors, Protein](#)]. There are a few exceptions. Inhibitors of this type are not all homologous. Indeed, a large number of families of protein inhibitors of serine proteinases were described (Table 1). It should be noted that the division into animal, plant, and microbial inhibitors is far from absolute. As an example, a leech protein [see [Eglin C](#)] is a member of the potato I family. Similarly, there are plant serpins. Nonetheless, the division seems convenient, indicating as it does the source of most known members of the family. Another caveat should be borne in mind. Families of proteins are structurally related groups. It does not follow that all the members of the family are proteinase inhibitors. For example, ovalbumin, while a serpin, is not a proteinase inhibitor. Neither are several other serpins. Some members of the bovine pancreatic trypsin inhibitor (Kunitz) family, such as dendrotoxins and bungarotoxins, act as  $K^+$  and  $Ca^{2+}$  channel blockers. It seems unlikely that the dendrotoxins and the  $Ca^{++}$  blockers are efficient inhibitors of proteinases. It seems more likely that inhibitors were recruited to serve as toxins, although the opposite scenario can also be contemplated. Some authors discuss the possibility that in those families, such as the cereal family, where there are only a few inhibitors, the inhibitory properties may have arisen accidentally.

**Table 1. Standard-Mechanism Canonical Inhibitors (in Bold)**

---

## Protein Inhibitors of Serine Proteinases

| Animal   | Plant  | Microorganism                                       |
|--|--|---|
| <b>BPTI (Kunitz) family<sup>a</sup></b>                                  | <b>Bowman–Birk family<sup>b</sup></b>                          | <b>SSI family<sup>c</sup></b>                       |
| <b>PSTI (Kazal) family<sup>d</sup></b><br><a href="#">Ascaris family</a> | <b>STI (Kunitz) family<sup>e</sup></b><br><b>Squash family</b> | <b>Ecotin family</b><br><b>Marinostatin family</b>  |
| <b>Chelonian in family<sup>f</sup></b>                                   | <b>Potato I family</b>   |   |
| <b>Antistasin family</b>   | <b>Potato II family</b>  |   |
| <b>Silkworm family</b>   | <b>Cereal family</b>   |   |
| <b>Grasshopper family</b>  | <b>Arrowhead family</b>  |   |
|  | <b>Rapeseed family</b>   |   |
| Hirudin family <sup>g</sup>  |  | Subtilin family propeptides                         |
| Serpin family  |  |   |
| TAP family <sup>h</sup>  |  |   |
|  |  | <i>Streptomyces</i> <sup>i</sup> family propeptides |

<sup>a</sup> Bovine pancreatic trypsin inhibitor (Kunitz) [see [BPTI \(Bovine Pancreatic Trypsin Inhibitor\)](#)].

<sup>b</sup> [See Bowman–Birk inhibitor.]

<sup>c</sup> *Streptomyces* subtilisin inhibitor family.

<sup>d</sup> Pancreatic secretory trypsin inhibitor (Kazal).

<sup>e</sup> [See Soybean trypsin inhibitor.]

<sup>f</sup> The most studied member of this family is secretory leukocyte proteinase inhibitor (SLPI).

<sup>g</sup> [See Hirudin.]

<sup>h</sup> Tick anticoagulant peptide.

<sup>i</sup> Propeptides of a lytic proteinase and *Streptomyces griseus* proteinases *A* and *B* as well as glutamic acid-specific proteinase.

The mechanism of association of serine proteinases with their protein inhibitors has been a matter of great interest. This mechanism, which is well understood for a large subset of those inhibitors called standard-mechanism inhibitors and listed in Table 1 in bold, is given in equation (1).

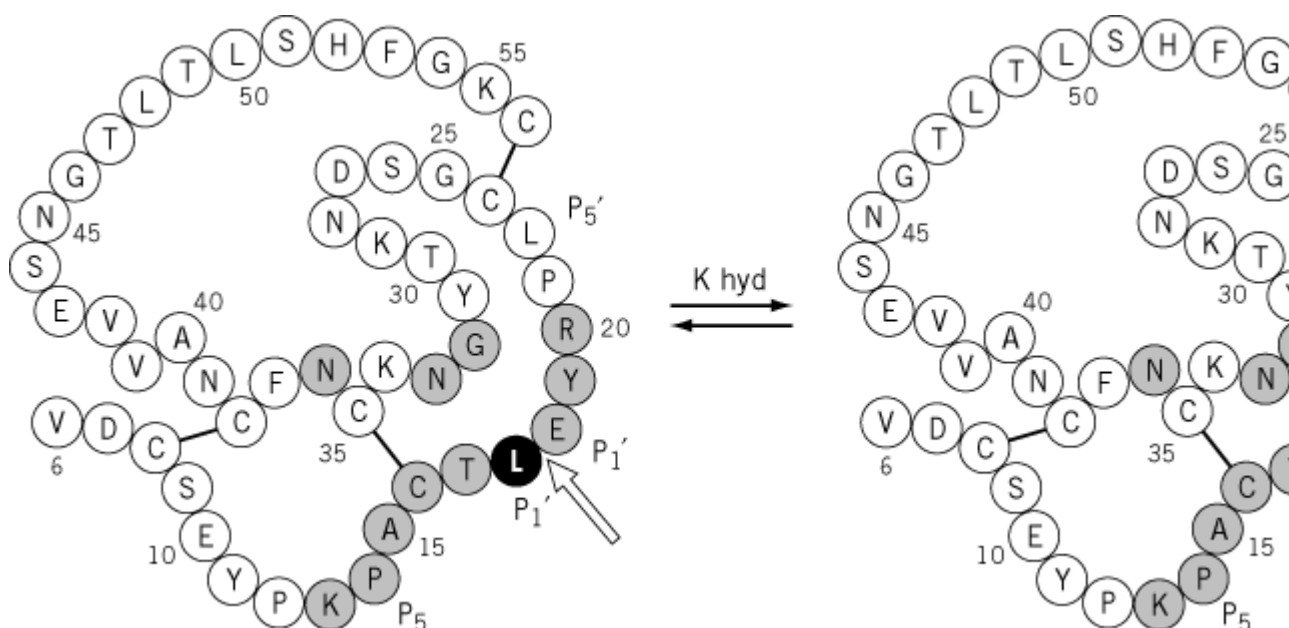


where  $E$  is the enzyme (serine proteinase) and  $C$  the stable proteinase inhibitor complex.  $I$ , the virgin inhibitor, has all its peptide bonds intact, whereas  $I^*$ , the modified inhibitor, has one specific peptide bond hydrolyzed (Fig. 1). With less than a stoichiometric amount of enzyme added,  $I$  and  $I^*$  attain equilibrium, given by  $K_{\text{hyd}}$ .

$$K_{\text{hyd}} = \frac{(I^*)}{(I)} \quad (2)$$

Equation (1) asserts that the complex  $C$  is the same substance whether made by mixing  $E+I$  or  $E+I^*$ . In addition, it turns out that, for the great majority of standard-mechanism inhibitors,  $I$  and  $I^*$  are equally effective thermodynamically. This is because, for most inhibitors, the equilibrium constant for reactive site peptide bond hydrolysis [eq. (2)] is very close to unity at neutral pH.

**Figure 1.** The amino acid sequence of virgin (left) and of modified (right) turkey ovomucoid third domain, a standard-mechanism inhibitor of serine proteinases. In most such inhibitors, the equilibrium constant  $K_{\text{hyd}}$  is the near unity. The residues are numbered sequentially from Val<sup>6</sup> and also (Schechter and Berger notation) from the reactive site peptide by an arrow in the virgin inhibitor. Leu<sup>18</sup> (in black) is the  $P_1$  residue and Glu<sup>19</sup> the  $P_1'$ . The shaded residues make contact with enzyme in complexes. The bars between C residues indicate disulfide bridges. The sequence from Lys<sup>14</sup> ( $P_6$ ) to Arg<sup>21</sup> ( $P_1'$ ) is a contiguous contact loop; a portion of it from  $P_4$  to  $P_3'$  is the canonical region.

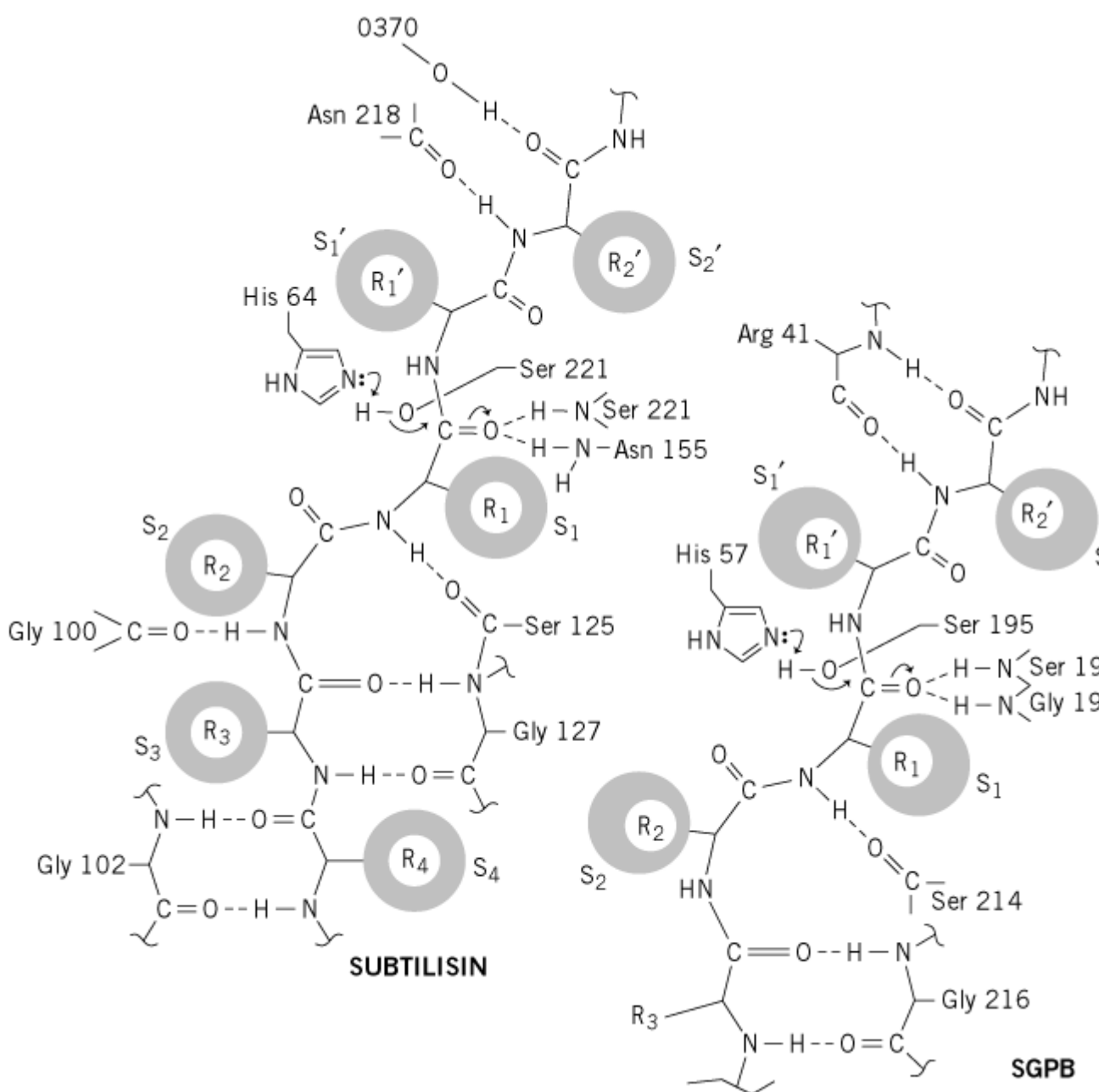


The detailed nature of  $C$  was solved by many X-ray crystallographic structure determinations. In all these complexes, the reactive site peptide bond is intact. The oxygen atom OG1 of the catalytic serine of the enzyme makes a very close approach to the C atom of the  $P_1$  residue of the inhibitor (for the Schechter and Berger notation, see Fig. 1). The distance is shorter than for van der Waals closest approach but longer than C—O covalent single bond. This close approach suggests a tetrahedral intermediate. The peptide bond remains planar, however, and the C atom is trigonal, not tetrahedral. The O atom of the  $P_1$  residue inserts into the oxyanion hole (Fig. 2), and the N atom



donates a hydrogen bond to the enzyme. About a dozen residues of the inhibitor make close contacts with the enzyme. They consist of a (sequentially) contiguous set surrounding the reactive site and of a few discontinuous residues. The latter set seems very different in different families and occasionally even in subfamilies. A portion of the contiguous set comprising the  $P_4$  to  $P_3'$  residues is called the canonical set. The canonical residues are said to have closely similar  $\phi$  angles in free inhibitors and in complexes with all cognate enzymes but, most importantly, in all inhibitor families (see Table 2). It appears that all canonical inhibitors exhibit the standard mechanism and vice versa.

**Figure 2.** Diagrammatic representation of the interaction between the active sites of serine proteinases and standard-mechanism canonical protein inhibitors: (left) for the subtilisin class of enzymes and (right) for the chymotrypsin class of enzymes exemplified by *Streptomyces griseus proteinase B*, SGPB. In both cases, an imperfect, antiparallel beta structure is formed between enzyme and inhibitor. In the chymotrypsin group, one strand of the enzyme interacts with the combin loop of the inhibitor. In subtilisin, there are additional interactions at  $P_4$  and  $P_2$  with another strand of the enzyme. (From C. A. McPhalen and M. N. G. James (1988) *Biochemistry* 27, 6582–6598.)



**Table 2. The Ramachandran Angles  $\phi$  and  $\psi$  of a Kazal Family and of a Potato I Family Inhibitor Free and in Complex with Bovine Chymotrypsin in the Canonical Portion of the Combining Loops**

|                                       | $P_4$  |        | $P_3$  |        | $P_2$  |        | $P_1$  |        | $P_1'$ |        | $P_2'$ |        | $P_3'$ |        | PDB file          |
|---------------------------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|-------------------|
|                                       | $\phi$ | $\psi$ | $\phi$ | $\psi$ | $\phi$ | $\psi$ | $\phi$ | $\psi$ | $\phi$ | $\psi$ | $\phi$ | $\psi$ | $\phi$ | $\psi$ |                   |
| OMSVP3 <sup>a</sup><br>free           | -163   | 149    | -155   | 131    | -174   | 87     | -96    | 9      | -139   | 58     | -99    | 93     | -69    | 130    | 2 ovo             |
| OMTKY3 <sup>b</sup> -<br>chymotrypsin | -136   | 129    | -150   | 131    | -160   | 68     | -32    | 107    | -159   | 74     | -107   | 113    | -76    | 142    | 1 cho             |
| Eglin c free                          | -71    | 153    | -163   | 132    | -150   | 74     | -82    | 12     | -159   | 99     | -128   | 111    | -123   | 129    | Note <sup>c</sup> |
| Eglin c-<br>chymotrypsin              | -81    | 156    | -166   | 144    | -153   | 79     | -97    | 40     | -166   | 90     | -121   | 131    | -129   | 129    | 1 acb             |

<sup>a</sup> Silver pheasant ovomucoid third domain (member of the Kazal family).

<sup>b</sup> Turkey ovomucoid third domain (member of the Kazal family). The only difference between the two inhibitors is at  $P_1$ . Met in OMSVP3 and Leu in OMTKY3. Eglin c is a member of the potato I family.

<sup>c</sup> K. Hipler, J. P. Priestle, J. Rahuel, and G. Grütter (1992) X-ray crystal structure of the serine proteinase inhibitor eglin c at 1.95 Å resolution. *FEBS* **309**, 139-145.

Serine proteinases exhibit many different specificities. Words such as *trypsinlike*, *chymotrypsinlike*, and *elastase-like* have become commonplace descriptors. The protein inhibitors combine with all these enzymes, but the strength of association differs greatly. In some cases,  $K_a$  is so small ( $K_I$  so great) that the inhibitor is said to be ineffective. In other cases, very strong inhibition is observed. Aside from discriminating between different substrate specificities of enzymes, inhibitors from some families, eg, BPTI (Kunitz), appear to inhibit enzymes only of the chymotrypsin group of serine proteinases. Others inhibit members of both the chymotrypsin group and the subtilisin group, eg, PSTI (Kazal), ovomucoids, and so forth. It is highly likely that the chymotrypsin/subtilisin discrimination is a consequence of the family's scaffolding. In contrast, the better-studied substrate specificity is caused by the nature of the inhibitor's residues that make contact with the enzyme (shaded in Fig. 1). The simplest of these is  $P_1$ , the primary specificity residue. This residue, on complex formation, becomes embedded in the primary specificity  $S_1$  pocket of the enzyme. When  $P_1$  is Lys and Arg, strong inhibition of trypsin is anticipated and frequently found.  $P_1$  Tyr, Trp, Phe, Leu, and Met (in that order) are optimal for chymotrypsin. Ala and Leu for porcine pancreatic elastase and Cys, Leu, and Met for *Streptomyces griseus* proteinases A and B. Glutamic acid-specific proteinase from *Streptomyces griseus*, Glu SGP, is strongly inhibited by avian ovomucoid third domains with  $P_1$  Glu and Asp. The strength of interaction and specificity of inhibitors does not depend solely on the nature of  $P_1$ . The other contact residues (Fig. 1), in aggregate, contribute much more to the strength and specificity of association than  $P_1$  does. The effects of contact residue changes in various positions in the inhibitor on the free energy of association are nearly additive. The

conclusions about the effect of contact residues on inhibition emerged both from the study of natural inhibitors, where hypervariability of contact residues is observed, and from intentional changing of these residues by recombinant DNA technology or by chemical and enzymatic modification. In the recombinant approach, the design of protein inhibitors to become strong or specific inhibitors of a new target enzyme became a paradigm of protein engineering technology. Recently, phage display has often been employed to aid in rapid screening of the set of contact residues. Probably even more striking is the manipulation of the contact residues by nature in the process of evolution. In most proteins, both structurally and functionally important residues are conserved among closely related species. This is not the case for several families of protein inhibitors of serine proteinases. Here, the structurally important residues are still conserved, but the functional enzyme-inhibitor residues are often the most variable in the inhibitor molecule. The term *hypervariability* of the reactive site region was adopted from the antibody literature, where hypervariability describes the variation between numerous antibodies in the same individual. Among the inhibitors, the term is used to describe the variation of one inhibitor among closely related species. Positive Darwinian selection was invoked as the driving force maintaining hypervariability among inhibitors from related species. The suggestion is controversial.

#### Suggestions for Further Reading

W. Boder and R. Huber (1992) Natural protein proteinase inhibitors and their interaction with proteinases. *Eur. J. Biochem.* **204**, 433–451.

T. E. Creighton and N. J. Darby (1989) Functional evolutionary divergence of proteolytic enzymes and their inhibitors. *Trends Biochem. Sci.* **14**, 319–324.

M. Laskowski Jr. et al. (1988) Positive Darwinian selection in evolution of protein inhibitors of serine proteinases. *Cold Spring Harbor Symp. Quant. Biol.* **52**, 545–553.

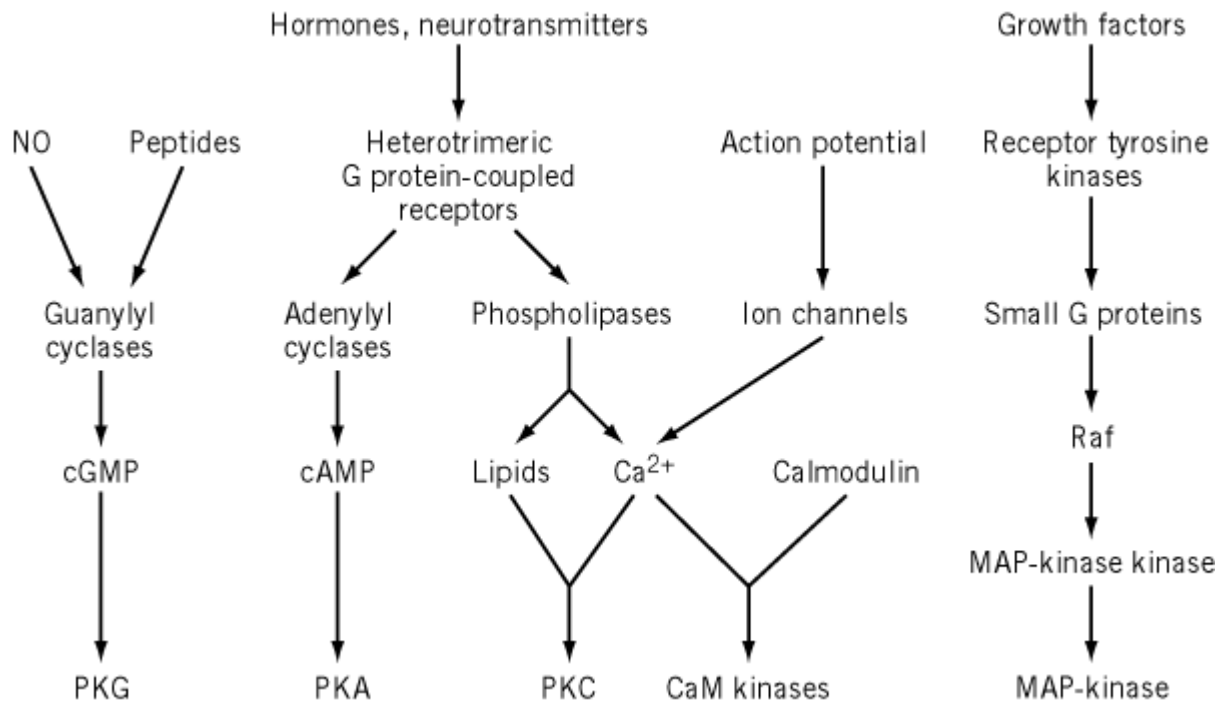
M. Laskowski Jr. and I. Kato (1980) Protein inhibitors of proteinases. *Ann. Rev. Biochem.* **49**, 593–626.

K. Ozawa and M. Laskowski Jr. (1966) The reactive site of trypsin inhibitors. *J. Biol. Chem.* **241**, 3955.

## Serine/Threonine Kinases and Phosphatases

[Serine](#) and [threonine](#) are the amino acid residues in [proteins](#) that are most commonly phosphorylated. Their phosphorylation is used as a regulatory device in many [signal transduction](#) pathways of prokaryotic and eukaryotic cells. The presence of a phosphate group instead of an hydroxyl group on a serine or a threonine side chain modifies dramatically its size and charge. Phosphorylation can have important consequences on the properties of the protein in the vicinity of the phosphorylated residue or can lead to global rearrangements of the protein. Phosphorylation of proteins on serine or threonine residues is the result of a dynamic equilibrium between the activities of the [kinases](#) and **phosphatases** acting on these residues. Almost all serine/threonine protein kinases (Ser/Thr kinases) have a conserved catalytic **domain** that is a member of a very large gene family, which contains also tyrosine and dual-specificity protein kinases ([1](#)). Phosphoserine/phosphothreonine protein phosphatases (Ser/Thr phosphatases) are [metal-requiring enzymes](#) that belong to two highly conserved gene families ([2-5](#)). Both Ser/Thr kinases and phosphatases are tightly regulated within cells by intracellular and extracellular signals. Major pathways activating Ser/Thr kinases in response to extracellular stimuli are illustrated schematically in [Fig. 1](#). In many instances, their precise localization is also tightly controlled by specific targeting sequences or by interaction with anchoring proteins or targeting subunits ([6-8](#)).

**Figure 1.** Examples of regulation of Ser/Thr kinases by extracellular signals in higher eukaryotes. Various signals such as hormones, neurotransmitters, nitrous oxide (NO), growth factors, or action potential can activate signal transduction pathways, leading to the stimulation of protein kinases, including cGMP-dependent protein kinase (PKG), cAMP-dependent protein kinase (PKA), protein kinase C (PKC), Ca<sup>2+</sup>/calmodulin-dependent protein kinases (CaM kinases), and mitogen-activated protein kinases (MAP kinases). These kinases phosphorylate a number of protein substrates including enzymes, transcription factors, cytoskeletal proteins, receptors, and ion channels, the phosphorylation of which is responsible for the biological effects of the extracellular signals. The pathways are highly simplified, and the numerous cross-talks that take place between them are ignored.



## 1. Ser/Thr Kinases

### 1.1. Structure and Regulation

Most Ser/Thr kinases have in common a conserved catalytic domain, whose three-dimensional [protein structure](#) is well-characterized by [X-ray crystallography](#) (9, 10). Schematically, this catalytic domain comprises a small lobe, mostly of [b-sheet](#), a large lobe mostly of [a-helix](#), and a cleft between the two, in which catalysis takes place. The small lobe is involved in ATP binding, and the large lobe is involved in substrate binding. For the kinase to be active, a critical loop of the large lobe, termed the *activation loop* or *T loop*, must be positioned correctly. This requires in many cases the phosphorylation of the loop on one or several residues. In some cases (eg, cAMP-dependent protein kinases, or [protein kinase A](#)) this phosphorylation is constitutive. In others (eg, [MAP kinases](#)) phosphorylation is a major regulatory step carried out by an activating protein kinase. In addition to this conserved catalytic domain, Ser/Thr kinases comprise other domains that are highly variable from one kinase to the other. These domains play a role in targeting and in oligomerization, as well as in regulation. In several cases, the regulatory domain contains a substrate or a pseudosubstrate sequence that obstructs the catalytic site when the Ser/Thr kinase is in an inactive state. Activation, usually by binding of a regulatory molecule, removes this intramolecular inhibition, allowing access of substrates to the active site. Activation can also occur, usually under nonphysiological conditions, following limited **proteolysis**, which releases a free, constitutively active catalytic domain. There are a very large number of Ser/Thr kinases, which have been classified on the basis of the sequence of their catalytic domain (Table 1) (1). Usually, this classification is in good agreement with the mode

of regulation of the Ser/Thr kinases, and it will be used here for an overview of this family of enzymes.

**Table 1. Protein Ser/Thr Kinases<sup>a</sup>**

---

**1. AGC group**

*Cyclic nucleotide-regulated*

cAMP-kinase (PKA)

cGMP-kinase (PKG)

*Protein kinases C (PKC)*

Classical PKC (cPKC)

New PKC (nPKC)

Atypical PKC (aPKC)

*PKN*

*RAC/Akt/PKB*

*G-protein-coupled receptor kinases (GRKs)*

*p70-S6kinase*

*p90-S6kinase (rsk)*

**2. CaMK group**

*Calcium/calmodulin-regulated protein kinases (CaMK) dedicated to one substrate*

Phosphorylase kinase

Myosin light chain kinase (MLCK)

eEF2 kinase (CaM-kinase III)

CaM-kinase kinase

*Multifunctional CaM kinases (many substrates)*

CaMK I

CaMK II

CaMK IV or Gr

*AMP-regulated protein kinase*

**3. GMGC group**

*Cyclin-dependent kinases*

*Mitogen activated kinases and related kinases*

MAPkinases (ERK 1, ERK 2)

JNK

p38

*Glycogen synthase kinase3 (GSK3)*

*Protein kinase CK2 (casein kinase II)*

**4. Other protein Ser/Thr kinases and dual specificity protein kinases**

*Polo-like*

*MEK 1/2 [eg, MAP-kinase kinase (MAPKK)]*

*MEKK*

*Raf kinases*

*Pak/Ste 20*  
*p160 ROCK*  
*Wee 1 mik 1*  
*Protein kinase CK1 (casein kinase I)*  
*Heme-regulated kinase*  
*Double-stranded RNA-regulated*  
*Serine kinase receptors: Transmembrane receptors for TGF, activin, etc.*

---

<sup>a</sup> Classification based on Ref. [1](#).

### 1.2. The AGC Group

This group comprises cAMP-dependent protein kinases (PKA), cGMP-dependent protein kinases (protein kinase G, PKG), [protein kinase C](#) (PKC), and related enzymes. PKA was the first protein kinase identified to be regulated by a [second messenger](#), [cyclic AMP](#) (cAMP). PKA is the main target for cAMP in eukaryotic cells and is probably the most extensively studied protein kinase. The catalytic (C) and regulatory (R) subunits are coded by separate **genes** and form a heterotetramer (R<sub>2</sub>C<sub>2</sub>) ([11](#), [12](#)). Upon binding of cAMP (two molecules of cAMP per catalytic subunit), the two catalytic subunits are released, and they phosphorylate neighboring proteins or diffuse to the [nucleus](#), where they phosphorylate proteins involved in gene [transcription](#). PKG is activated in a similar manner, but the catalytic and regulatory domains are located on the same polypeptide chain. PKCs form a subfamily of Ser/Thr kinases characterized by the optional presence of **zinc-finger motifs** (C1 domain) and a Ca<sup>2+</sup>-phospholipid-binding domain (C2 domain) ([13-15](#)). The C1 domain provides a binding and activation site for diacylglycerol and for its nonmetabolizable structural analogues, the carcinogenic **phorbol esters** derived from croton oil. The C2 domain allows a Ca<sup>2+</sup>-dependent interaction of the enzyme with membrane phospholipids. Whereas the subgroup of “classical” PKCs have both C1 and C2 domains, the “novel” PKCs lack the C2 domain, and the “atypical” PKCs lack the C2 domain and one of the two zinc fingers of the C1 domain. Other kinases in this group include PKB (also known as Akt), which is recruited to the membrane by interaction of its **pleckstrin homology** domain with phosphatidylinositol-2,3,4-P<sub>3</sub> (see [Inositol Lipids and Phosphates](#)) and is phosphorylated by an activating kinase ([16](#), [17](#)). As such, it mediates some of the actions of **phosphatidylinositol-3-kinase**. Another interesting group of enzymes are the [G-protein-coupled receptor](#) kinases (GRK), which phosphorylate [rhodopsin](#) or seven-transmembrane-domain receptors when they are in the ligand-bound form, controlling their desensitization ([18](#), [19](#)).

### 1.3. The CaMK Group

This group comprises Ca<sup>2+</sup>/**calmodulin**-activated enzymes (or CaM kinases). Two types of CaM kinases have been characterized: those that have a narrow substrate specificity (eg, phosphorylase kinase, myosin light-chain kinase, eEF2 kinase) and those that have a broad substrate specificity (CaM kinases I, II, and IV) ([20](#)). In addition to activation by the binding of the Ca<sup>2+</sup>/calmodulin complex, CaM kinases I and IV are activated by a different Ca<sup>2+</sup>/calmodulin-dependent enzyme, constituting a Ca<sup>2+</sup>/calmodulin-activated protein kinase cascade ([21](#)). CaM kinase II stands apart in this group, because of its oligomeric structure (8 to 12 subunits) and biochemical properties, which endow it with an interesting capability to integrate Ca<sup>2+</sup> transients ([22](#)). The AMP-regulated protein kinase is structurally related to CaM kinases, but is CaM-independent. It plays an important role in metabolic controls ([23](#)).

### 1.4. The CMGC Group

The cyclin-dependent protein kinases (CDKs) are critical regulators of the [cell cycle](#) that are activated by a combination of factors, including phosphorylation of the activation loop,

dephosphorylation of amino-terminal residues, and interaction with proteins whose levels vary during the cell cycle, the [cyclins](#) (24-27). The MAP kinases (mitogen-activated protein kinases) form a group of protein kinases that contain essentially only the catalytic domain (28, 29). They are activated by phosphorylation of the activation loop on both a threonine and a tyrosine residue by a specific group of dual-specificity protein kinases, the MAP-kinase kinases. These MAP-kinase kinases are themselves activated by phosphorylation of their activation loop, by a third group of activating kinases, in a characteristic kinase cascade that has been highly conserved during evolution. Some MAP kinases (also called ERKs, extracellular signal-regulated kinases) are activated in response to growth factor stimulation. Others, such as Jun-N-terminal kinase (JNK), or p38, are implicated in cascades activated by cytotoxins and cellular stress (eg, osmotic shock, [tumor necrosis factor](#)) and are sometimes known as SAP kinases (stress-activated protein kinases) (see Ref. 28 for a discussion of the nomenclature of MAP kinases). The same group also includes glycogen-synthase kinase-3 (GSK3), which is not only involved in the control of glycogen synthesis, as its name implies, but also in differentiation, cell survival, and other functions (30). Protein kinase CK2 is a ubiquitous and highly conserved enzyme, existing as a dimer of catalytic and noncatalytic subunits. Although it has numerous substrates, its regulation and function are still poorly understood.

### 1.5. Miscellaneous Ser/Thr Kinases and Dual-Specificity Protein Kinases

Several protein kinases (eg, Raf, p160-ROCK, PAK) interact with the GTP-bound form of small G proteins (such as Ras, Rho, or Rac, respectively). This interaction allows their translocation to the plasma membrane, which is often completed by phosphorylation by activating kinases. Raf activation triggers the ERK pathway by phosphorylating MEK, the dual specificity kinase that activates ERK (31). Related enzymes are involved in the activation of other dual-specificity MAP-kinase-kinases that phosphorylate and activate SAP-kinases JNK or p38. Protein kinases activated by heme or by double-stranded RNA play a role in the control of **protein biosynthesis** and antiviral defenses, respectively. Protein kinase CK1 is a highly conserved monomeric enzyme, whose regulation and function are poorly understood. Some Ser/Thr kinases have a **transmembrane** segment and are activated by binding of a polypeptide chain to their extracellular domain, in a manner similar to the activation of growth-factor [tyrosine kinase receptors](#) (32). These Ser/Thr kinases include receptors for [transforming growth factor b](#), activin, and inhibin.

### 1.6. Ser/Thr Kinase Inhibitors

A number of Ser/Thr kinase inhibitors have been discovered and/or designed. Some are peptide pseudosubstrates, usually based on the sequence of the inhibitory domain of the kinase itself or of an associated protein. The best-known example is the PKA inhibitor (termed PKI or Walsh inhibitor) naturally expressed as a small regulatory protein that may be involved in the nucleocytoplasmic traffic of PKA (33). The other group of Ser/Thr kinase inhibitors is comprised of nonpeptide compounds that interact either with the protein kinase catalytic domain or with the regulatory domain, preventing its activation (eg, binding sites for cAMP in PKA or for diacylglycerol in PKC). Such inhibitors are powerful tools for studying phosphorylation pathways in intact cells and may provide interesting therapeutic leads (28, 34, 35).

## 2. Ser/Thr Phosphatases

On the basis of sequence analysis, Ser/Thr phosphatases belong to two gene families, but they are all metalloenzymes and have a similar structural organization (5, 36). In many cases, they are found as catalytic subunits that interact with other subunits responsible for their targeting and regulation (2, 3, 36, 37). Ser/Thr phosphatases are also classified according to their functional properties, which were determined before the cloning of their [complementary DNAs](#) (see Refs. 38 and 39 and Table 2).

**Table 2. Protein Phospho-serine and Phospho-threonine Phosphatases**

---

| Classification Based on Sequence                       | Classification Based on Function  |
|--|---|
| <b>PPP gene family</b>                                 | <b>Protein phosphatases 1</b><br><i>Inhibited by phospho-inhibitor 1 and inhibitor 2</i>                                    |
| PPP1C  | Protein phosphatases 1 (PP1)  |
|  | <b>Protein phosphatases 2</b><br><i>Resistant to phospho-inhibitor 1 and inhibitor 2</i>                                    |
| PPP2C  | Protein phosphatase 2A (PP2A)<br><i>Insensitive to divalent cations</i>   |
| PPP3C  | Protein phosphatase 2B (PP2B) = calcineurin<br><i>Activated by <math>Ca^{2+}</math> and <math>Ca^{2+}</math>/calmodulin</i> |
| <b>PPM gene family</b>                                 | Protein phosphatase 2C (PP2C)<br><i>Activated by <math>Mg^{2+}</math></i>   |
| <b>Other phosphatases from the PPP family</b>          |   |
| Related to PP1: PPY, Ppz, Ppq                          |   |
| Related to PP2A: PPX (PP4), PP6 (Sit4, Ppe)            |   |
| PP5 (contains a tetratricopeptide repeat) <sup>a</sup> |   |

<sup>a</sup> Source: References [39](#), [43](#), and [47](#).

### 2.1. PPP Group

This gene family includes several of the best-characterized Ser/Thr phosphatases. The PP1 subfamily is comprised of a catalytic subunit that interacts with a number of other proteins including targeting subunits (eg, for glycogen particles, myofibrils, and dendritic spines) and inhibitor proteins ([2](#), [3](#), [40](#)). The activity of some of the inhibitor proteins acting on these phosphatases can be regulated by phosphorylation. For example, inhibitor-1 and DARPP-32, two related proteins, are potent inhibitors of PP1 only when they are phosphorylated by PKA, providing a mechanism by which extracellular signals that increase cAMP levels can inhibit a protein phosphatase ([41](#)). Enzymes of the PP2A subfamily are heterotrimers comprised of a catalytic subunit and two regulatory subunits ([42](#)). Like PP1, PP2A is mostly regulated by interaction with other proteins. A number of naturally occurring toxins including okadaic acid, calyculin, microcystin LR, and tautomycin produced by various microorganisms (dinoflagellates, cyanobacteria, etc.) are potent inhibitors of PP1 and PP2A ([34](#)). These compounds, which are usually membrane permeant, can be environmental hazards, but they are also extremely useful tools to study the role of PP1 and PP2A in biological processes. The CDNAs for several novel Ser/Thr phosphatases related to PP1 or PP2A have been **cloned** recently ([43](#)) (Table [2](#)). Although PP2B is structurally related to PP1 and PP2A, it is far less sensitive to the toxins listed above. PP2B is expressed ubiquitously in eukaryotic cells, but it is particularly abundant in brain; it was independently isolated as a major brain  $Ca^{2+}$ /calmodulin binding phosphatase and called *calcineurin* ([44](#)). PP2B/calcineurin is a heterodimer, comprised of a catalytic subunit (A) and a  $Ca^{2+}$ -binding subunit related to calmodulin (B). Calcineurin is activated by binding of  $Ca^{2+}$ /calmodulin to A and of  $Ca^{2+}$  to B; hence it is a phosphatase that responds directly to  $Ca^{2+}$  transients. Calcineurin is inhibited by the complex between the immunosuppressant drugs



[cyclosporin A](#) and [FK506](#), in combination with their intracellular receptor proteins, the [immunophilins](#) ([cyclophilin](#) and FK-binding protein, respectively) (45). These drugs are used to prevent rejection following organ grafting and act by blocking the  $\text{Ca}^{2+}$ -dependent, calcineurin-mediated activation of **interleukin-2** transcription in lymphocytes.

## 2.2. The PPM Gene Family

This group of phosphatases is less well-characterized than the PPPs. It corresponds to the PP2C enzymes, which are activated *in vitro* by high concentrations of  $\text{Mg}^{2+}$  (39). In cells, they appear to be monomeric and little is known about their regulation. PP2C has a very widespread expression among eukaryotes, and it may have relatives in prokaryotes (see Ref. 46). The lack of a specific inhibitor, however, has hindered study of the functional role of this enzyme.

## Bibliography

1. S. K. Hanks and T. Hunter (1995) *FASEB J.* **9**, 576–596.
2. S. Shenolikar and A. C. Nairn (1991) *Adv. Cyclic Nucleotide Protein Phosphorylation Res.* **23**, 1–121.
3. S. Wera and B. A. Hemmings (1995) *Biochem. J.* **311**, 17–29.
4. J. E. Villafranca, C. R. Kissinger, and H. E. Parge (1996) *Curr. Opin. Biotechnol.* **7**, 397–402.
5. D. Barford (1996) *Trends Biochem. Sci.* **21**, 407–412.
6. M. J. Hubbard and P. Cohen (1993) *Trends Biochem. Sci.* **18**, 172–177.
7. M. C. Faux and J. D. Scott (1996) *Trends Biochem. Sci.* **21**, 312–315.
8. T. Pawson and J. D. Scott (1997) *Science* **278**, 2075–2080.
9. S. S. Taylor, D. R. Knighton, J. Zheng, L. F. Ten Eyck, and J. M. Sowadski (1992) *Annu. Rev. Cell Biol.* **8**, 429–462.
10. L. N. Johnson, M. E. M. Noble, and D. J. Owen (1996) *Cell* **85**, 149–158.
11. S. S. Taylor, J. A. Buechler, and W. Yonemoto (1990) *Annu. Rev. Biochem.* **59**, 971–1005.
12. S. S. Taylor, D. R. Knighton, J. Zheng, J. M. Sowadski, C. S. Gibbs, and M. J. Zoller (1993) *Trends Biochem. Sci.* **18**, 84–89.
13. Y. Asaoka, S. Nakamura, K. Yoshida, and Y. Nishizuka (1992) *Trends Biochem. Sci.* **17**, 414–417.
14. C. Tanaka and Y. Nishizuka (1994) *Annu. Rev. Neurosci.* **17**, 551–567.
15. A. C. Newton (1995) *J. Biol. Chem.* **270**, 28495–28498.
16. B. M. Marte and J. Downward (1997) *Trends Biochem. Sci.* **22**, 355–358.
17. J. Downward (1998) *Curr. Opin. Cell Biol.* **10**, 262–267.
18. K. Palczewski and J. L. Benovic (1991) *Trends Biochem. Sci.* **16**, 387–391.
19. J. G. Krupnick and J. L. Benovic (1998) *Annu. Rev. Pharmacol. Toxicol.* **38**, 289–319.
20. P. I. Hanson and H. Schulman (1992) *Annu. Rev. Biochem.* **61**, 559–601.
21. T. R. Soderling (1996) *Biochim. Biophys. Acta Protein Struct. Mol. Enzymol.* **1297**, 131–138.
22. A. P. Braun and H. Schulman (1995) *Annu. Rev. Physiol.* **57**, 417–445.
23. D. G. Hardie and D. Carling (1997) *Eur. J. Biochem.* **246**, 259–273.
24. M. Dorée and S. Galas (1994) *FASEB J.* **8**, 1114–1121.
25. J. Pines (1995) *Biochem. J.* **308**, 697–711.
26. E. Lees (1995) *Curr. Opin. Cell Biol.* **7**, 773–780.
27. D. O. Morgan (1997) *Annu. Rev. Cell Dev. Biol.* **13**, 261–291.
28. P. Cohen (1997) *Trends Cell Biol.* **7**, 353–361.
29. T. S. Lewis, P. S. Shapiro, and N. G. Ahn (1998) *Adv. Cancer Res.* **74**, 49–139.
30. J. R. Woodgett, S. E. Plyte, B. J. Pulverer, J. A. Mitchell, and K. Hughes (1993) *Biochem. Soc.*

Trans. **21**, 905–907.

31. R. Marais and C. J. Marshall (1996) *Cancer Surv.* **27**, 101–125.
32. J. Massagué and F. Weis-Garcia (1996) *Cancer Surv.* **27**, 41–64.
33. D. A. Fantozzi, S. S. Taylor, P. W. Howard, R. A. Maurer, J. R. Feramisco, and J. L. Meinkoth (1992) *J. Biol. Chem.* **267**, 16824–16828.
34. C. MacKintosh and R. W. MacKintosh (1994) *Trends Biochem. Sci.* **19**, 444–448.
35. L. Meijer (1996) *Trends Cell Biol.* **6**, 393–397.
36. S. Shenolikar (1994) *Annu. Rev. Cell Biol.* **10**, 55–86.
37. M. C. Mumby and G. Walter (1993) *Physiol. Rev.* **73**, 673–699.
38. T. S. Ingebritsen and P. Cohen (1983) *Science* **221**, 331–338.
39. P. Cohen (1989) *Annu. Rev. Biochem.* **58**, 453–508.
40. M. Bollen and W. Stalmans (1992) *Crit. Rev. Biochem. Mol. Biol.* **27**, 227–281.
41. P. Greengard, A. C. Nairn, J. A. Girault, et al. (1998) *Brain Res. Rev.* **26**, 274–284.
42. R. E. Mayer-Jaekel and B. A. Hemmings (1994) *Trends Cell Biol.* **4**, 287–291.
43. P. T. W. Cohen (1997) *Trends Biochem. Sci.* **22**, 245–251.
44. C. B. Klee, G. F. Draetta, and M. J. Hubbard (1988) *Adv. Enzymol.* **61**, 149–200.
45. J. Liu, J. D. Farmer, W. S. Lane, J. Friedman, I. Weissman, and S. L. Schreiber (1991) *Cell* **66**, 807–815.
46. P. Bork, N. P. Brown, H. Hegyi, and J. Schultz (1996) *Protein Sci.* **5**, 1421–1425.
47. P. E. Visconti, G. D. Moore, J. L. Bailey, et al. (1995) *Development* **121**, 1139–1150.

## Serotype

Microorganisms express at their surface a large array of molecules, including polysaccharides, that are frequently the target of **antibodies** produced by the host [immune response](#). This is the case of extracellular pathogens, like staphylococci or streptococci, that may express one or the other of a large collection of polysaccharides, differing from each other by a limited number of monosaccharide units. These small differences are sufficient to be recognized by different antibodies. The individual antigenic specificities so defined are given the name of *serotypes*. Serology of bacteria has been extensively developed for many years, so that microbiologists have a long list of identified serotypes that are of importance for diagnosis, prognosis, and treatment of a number of infectious diseases. For example, one single bacterial species such as *Salmonella enteritidis* may have 2000 distinct serotypes; *Streptococcus pneumoniae*, which causes bacterial pneumonia in humans, has nearly 100. Each strain, with its characteristic serotype, behaves as a discrete [antigen](#), which provides a way for the bacteria to escape the immune system by [antigenic variation](#). Whenever mutants occur in a pathogen, it may acquire an unexpected virulence and spread dangerously over very large populations of people or animals.

See also entries **Alloantibody**, [Antigen](#), [Immunogen](#), and [Antisera](#).

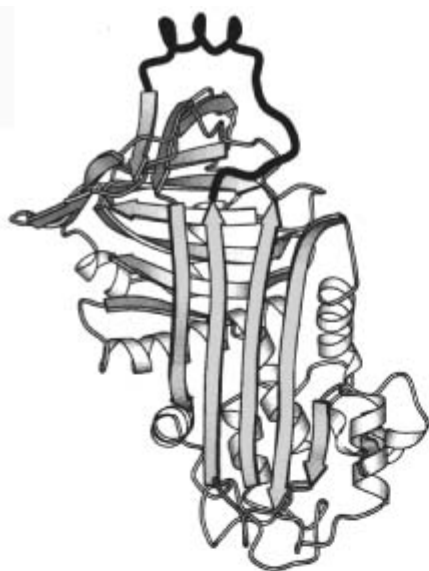
### Suggestion for Further Reading

- B. B. Finlay and S. Falkow (1989) *Microbiol. rev.* **53**, 210–230.

## Serpins

*Serpins* is an acronymic name given to a family of **serine proteinase inhibitors** that share a complex, but well-conserved, [tertiary structure](#) (1-3). Members of the family are diversely present in eukaryotes, plants, and **viruses**, and are evident in everyday life, from the foam protein of beer—the barley Z proteinase inhibitor—to the white of the breakfast egg—the noninhibitory serpin [ovalbumin](#). Most notably, though, the serpins are the principal proteinase inhibitors in human plasma; [antithrombin](#) controls the proteolytic [blood clotting](#) cascade; C<sub>1</sub>-inhibitor controls **complement** activation; the [plasminogen](#) activator inhibitors, PAI-1 and PAI-2, control fibrinolysis; and  $\alpha_1$ - **antitrypsin**, also known as  $\alpha_1$  proteinase inhibitor, modulates connective tissue restructuring. Altogether, the inhibitory serpins make up approximately 10% in molar terms, of the proteins in human plasma. Also present in plasma, although in smaller concentrations, are other serpins that have lost their inhibitory activity but have taken on other functions vital to life; examples are the vasopressor peptide source angiotensinogen, and the thyroxine- and **corticosteroid**-binding [globulins](#). The reason for the evolutionary success of the serpins is their possession, uniquely among the many families of serine proteinase inhibitors, of a mobile reactive-site loop (Fig. 1). It is the ability of this loop to change its conformation profoundly (4) that enables the serpins to bind to their target [proteinases](#) as a virtually irreversible complex (5-7). The possession of alternative folding topologies also explains the variety of other functions that have evolved among members of the family (8). It provides the evolutionary advantage of the potential to modulate inhibitory activity, as occurs in the spontaneous refolding of the fibrinolysis inhibitor PAI-1 (9) or on the conformational activation of antithrombin by heparin (10). This flexibility of fold carries with it a price, as serpins, because of their complicated mechanism of action, are especially vulnerable to [mutations](#) affecting their conformational mobility (11). In humans, such mutations in serpins are a significant cause of disease, including familial emphysema, cirrhosis, thrombosis, and allergic hypersensitivity.

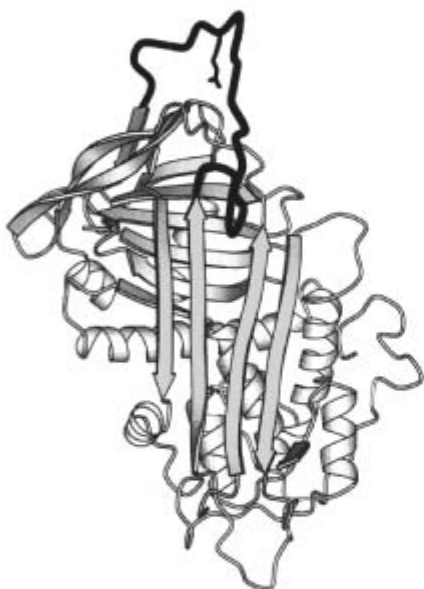
**Figure 1.** Mobility of the reactive loop of serpins (black) and its ability to refold as the central strand of the A b-sheet of the molecule (gray shaded bands). (a) The noninhibitory serpin ovalbumin has a relatively fixed helical loop; (b)  $\alpha_1$ -antitrypsin has a loop in the optimal inhibitory—canonical—conformation; (c) in antithrombin the loop is partly inserted into the A-sheet, with the reactive center arginine in an inaccessible internal orientation; (d) the latent inactive conformation of antithrombin has a completely inserted intact loop; (e) cleavage of  $\alpha_1$ -antitrypsin (shown in b) at its reactive center gives immediate insertion of the amino-terminal section of the loop; (f) peptides homologous to the reactive loop sequence can insert into the A-sheet displacing and disordering the reactive loop. (Crystallographic structures, with the F-helix removed for clarity, as detailed in Ref. 8.)



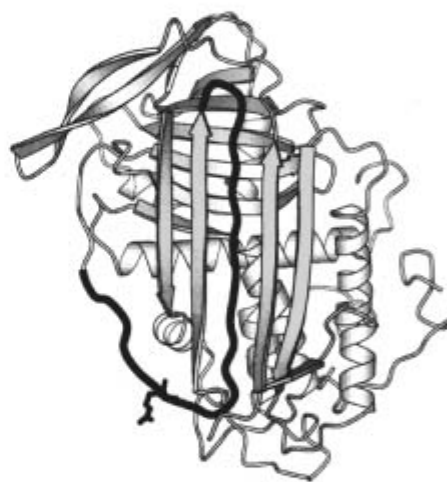
(a)



(b)



(c)



(d)



(e)



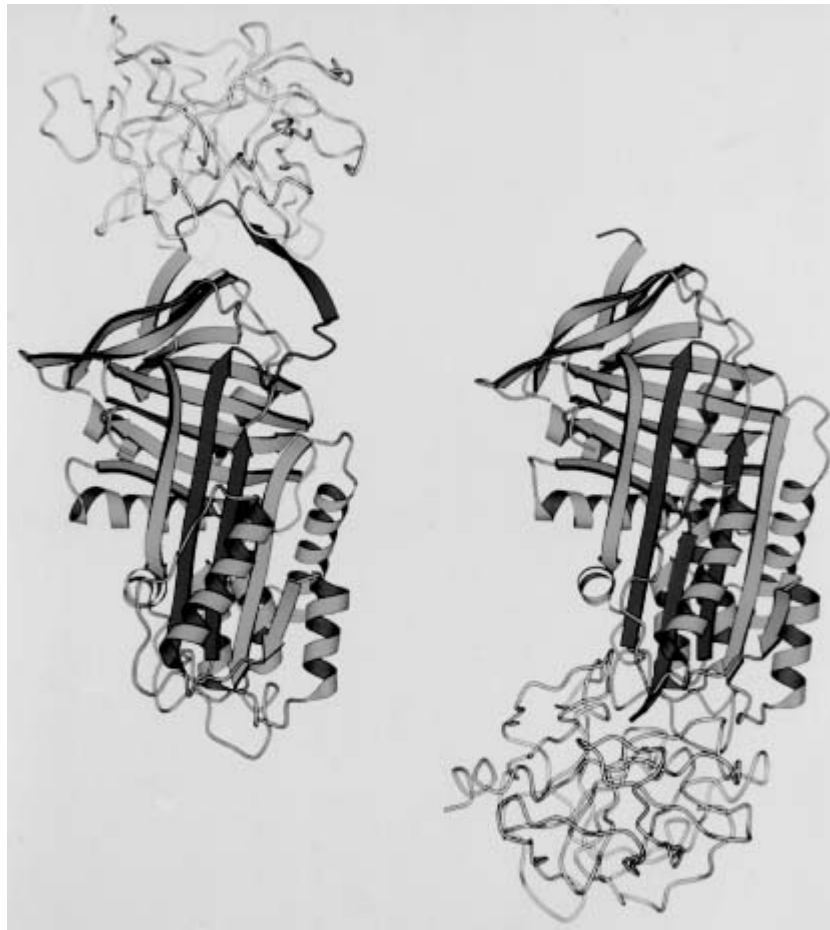
(f)

## 1. Structure and Function

There are now more than 300 known members of the serpin family, all of which share the same core [protein structure](#) of some 400 amino acid residues folded in nine [alpha-helices](#) and three [beta-sheets](#). Considerable variations can occur at the *N*-terminus of the molecule, as with the glycosaminoglycan-binding extensions in antithrombin and another anticoagulant serpin, heparin cofactor II, and the long glycosylated extension in C<sub>1</sub>-inhibitor. Similarly, functional **domains** have been added to the *C*-terminus, as in the plasma fibrinolysis inhibitor a<sub>1</sub>-antiplasmin. The central core of the molecule, however, is remarkably conserved, with almost superimposable backbone outlines from one serpin to another (3). The dominant feature of the molecule is the A β-sheet, which changes from a five-stranded conformation in the active inhibitory state to a six-stranded sheet in the inactive, complexed form of the molecule. This remarkable change, illustrated in Figure 1, results from the movement of the reactive site loop of the molecule, hinged at the end of [beta-strand](#) 5, from an externally exposed position to insertion as the central, strand 4, of the sheet. This shift in conformation is accompanied by an equally remarkable increase in molecular stability, typically from a **T<sub>m</sub> (melting temperature)** of 56°C in the inhibitory form to 100°C or more in the inserted, six-stranded, form. This stressed to relaxed, *S* → *R* change (12), was first shown to occur on **proteolytic** cleavage of the reactive-site loop, a process that is now realized to model the shift accompanying complex formation with the cognate proteinase. The *S* → *R* change in stability also occurs when the intact loop is incorporated into the A-sheet, as with the spontaneous transition of PAI-1 to the six-stranded latent conformation, or the inducible transition to the latent form in antithrombin, a<sub>1</sub>-antitrypsin, and a<sub>1</sub>-antichymotrypsin (13).

In the active inhibitory form of the serpins, the peptide loop containing the reactive center is exposed at the pole of the molecule. This peptide loop presents the reactive-center peptide bond—denoted the P<sub>1</sub>–P'<sub>1</sub> bond—as an ideal substrate. The cognate proteinase cleaves this bond to form an acyl-enzyme linkage that is stabilized by the rapid insertion of the amino-terminal portion of the reactive loop P<sub>1</sub>–P<sub>14</sub> into the A-sheet (5). This movement of the loop takes with it the linked proteinase to the opposite pole of the molecule (Fig. 2) thereby distorting and rendering inaccessible the acyl-enzyme bond, hence giving a stable, virtually irreversible complex of the two. Although no [X-ray crystallography](#) structure of the complex has yet been identified, there is support for this mechanism from numerous observations with both natural and **site-directed mutations** (14). Mutations in the hinges of the reactive loop slow the rate of loop insertion, with a proportionate shift of the interaction of the serpin with the proteinase, changing from inhibitor to substrate.

**Figure 2.** Mechanism of inhibition. A model based on the structure of intact (15) and cleaved (4) a<sub>1</sub>-antitrypsin, showing on the left how an outlined serine proteinase precisely docks with the canonically conformed reactive loop of a<sub>1</sub>-antitrypsin. On the right, cleavage with formation of an acyl-enzyme linkage is believed to result in a 70-Å shift of the proteinase driven by the insertion of the cleaved loop as the fourth strand of the A-sheet. (Figure prepared by J. P. Abrahams.)



Although there is uncertainty as to the final locking conformation of the complex, there is now a firm structural basis for the initial docking of inhibitor and proteinase. This had been a puzzle as in comparison to other families of serine proteinase inhibitors, the conformation of the reactive-site loop of the serpins is mobile and variable, from the totally exposed and well-ordered  $\alpha$ -helix of ovalbumin to the partly obscured and A-sheet-inserted loop of antithrombin (Fig. 1). The other nonserpins families of serine proteinase inhibitors are all much smaller molecules, which have diverse, but tightly fixed, conformations of their exposed reactive loops. The surprising finding is that, although these other families of serine proteinase inhibitors have evolved separate and unrelated structures, they all share the same canonical conformation that makes a complementary fit with the active-site cleft of the serine proteinases. The structure of wild-type  $\alpha_1$ -antitrypsin (15) has confirmed that the reactive loop of the serpins can also adopt this same precise conformation, and in the case of  $\alpha_1$ -antitrypsin, but not other serpins, the loop is held in this canonical loop conformation by side-chain [hydrogen bonds](#) to the body of the molecule. Thus the general mechanism of complex formation can be divided into two stages (Fig. 2). The first is docking, with the approach of the proteinase to the serpin inducing a shift in the inhibitory loop to the canonical conformation, followed by cleavage at the  $P_1$ - $P'_1$  [peptide bond](#), to give an acyl-enzyme linkage. The second change, immediate to the first, is to the locking conformation, with complete insertion of the  $P_1$ - $P_{14}$  loop into the A-sheet and transport of the proteinase some 60–70 Å to the distal pole of the molecule.

## 2. Modulation by Glycosaminoglycans (Heparin)

Antithrombin is the principal inhibitor of coagulation of human plasma, but its activity is relatively limited, until it binds to the heparan sidechains that line the microvasculature. Heparin, an animal

extract, is the therapeutic equivalent of the heparans; both are glycosaminoglycans containing a core pentasaccharide fragment that binds to the D and A helices of antithrombin. Crystallographic structures (10) of antithrombin and its complex with heparin reveal an elegant mechanism of activation and illustrate the way in which the mobile structure of the serpins has been adapted by evolution to modulate their function.

The reactive loop of antithrombin differs from that of other serpins in being partially inserted into the A-sheet, with an inward orientation of the reactive-center arginine residue, rendering it inaccessible to proteinases. The molecule as a whole is in a constrained conformation that is released on the binding of heparin, with a closure of the A-sheet, a shift of the D helix, and an elongation of both the D and A helices. This shift of the molecule from its constrained state results in the expulsion from the A-sheet of the reactive loop, which is then freed to take up the inhibitory conformation, with exposure of the reactive-center arginine. This conformational change fully activates antithrombin as an inhibitor of factor Xa, but in addition, the longer-chain heparins provide a bridging template for an even greater acceleration of the inhibition of thrombin.

A related but intriguingly different mechanism also takes place with the activation of heparin cofactor II. This binds another glycosaminoglycan, dermatan sulfate, to give a similar conformational activation, but with a supplementary contribution from the *N*-terminal extension of the inhibitor (16). The interaction of heparin with antithrombin also involves linkage and movement of the unstructured *N*-terminus, but this is more markedly so with heparin cofactor II. The much longer *N*-terminus of heparin cofactor II contains a sequence **homologous** to that of the thrombin inhibitor from leech, [hirudin](#) (17). Hirudin binds directly to an anion-binding site on thrombin, sterically blocking the [active-site](#) cleft of thrombin. There is good evidence (16) that, on activation with dermatan sulfate, the amino-terminus of heparin cofactor II moves to give a hirudin-like linkage and an analogous blockage of the active site of thrombin.

### 3. Molecular Pathology and Conformational Disease

More than 100 different mutations have now been identified in the plasma serpins as being responsible for associated diseases (11). In many cases, these affect specific functional domains, such as the heparin-binding site in antithrombin or the reactive site itself. By far the largest group of dysfunctional mutations, however, are those that affect the mobile mechanism of A-sheet opening and reactive-loop hinging and insertion. Such mutations can cause a spontaneous *S* → *R* transition to give the formation of the inactive latent form, or, if the change occurs at higher concentrations of the serpins, to allow the opening of the A-sheet with intermolecular linkage between the reactive loop of one molecule and the A-sheet of the next. These structurally well-defined mechanisms of monomeric transitions and of loop-sheet polymerization provide the best-defined examples of what is now being recognized as a new clinical entity—the conformational diseases (18).

The pathological transition of the active inhibitor to the latent form is most apparent with the mutants of antithrombin where the occurrence of the transition results in thrombosis. The most frequent example of such changes in the serpins, however, is observed with the common genetic variants of plasma  $\alpha_1$ -antitrypsin (19); 1 in 10 people of northern European descent are carriers of either the severe Z-deficiency variant or the milder S variant. Homozygotes for the Z variant are likely to develop the degenerative lung disease emphysema in early adult life, and in childhood they are at risk of developing a fatal liver cirrhosis. Homozygotes for the Z mutation accumulate  $\alpha_1$ -antitrypsin in the [endoplasmic reticulum](#) of their hepatocytes. The blockage of processing results in a plasma deficiency of  $\alpha_1$ -antitrypsin and consequently a failure to protect the lung tissue against the [elastase](#) released by inflammatory white cells. The mutation in Z antitrypsin results in the replacement of conserved Glu342 (residue P17) by a lysine. This site is at the base of the hinge of the reactive-center loop, and the effect of the mutation is to destabilize the A-sheet and hence favor the insertion of the loop from another Z antitrypsin molecule, to give loop-to-sheet polymerization (20). This results in the accumulation of intracellular tangles of fibrils at the site of synthesis in the liver, with

consequent cell death and eventually the development of liver cirrhosis.

#### 4. Cell Biology and the Serpins

There is now increasing interest in a large subclass of the serpins that have intracellular or pericellular functions (21). A consistent, though not invariable, finding in this subclass is the presence on the reactive loop of residues that readily oxidizes: [cysteine](#) or [methionine](#) residues. The likely reason is to confine the inhibitory activity of these serpins to the reducing environment of the cell and its immediate pericellular radius. The proximity to intracellular proteolytic processes makes these serpins candidate regulators of cell growth, differentiation, and death. The intra- and pericellular serpins can be divided into two main categories: those that resemble the chicken egg protein ovalbumin and those encoded by poxviruses (22).

In recent years an increasing number of ovalbumin-type serpins have been described that share the properties of lacking a conventional secretion [signal peptide](#) for [protein secretion](#), having a truncated C-terminus and possessing a similar overall gene structure. Several of these proteins have been described in human tissues, and, although in most cases their intracellular function is uncertain, some important clues are beginning to emerge. A good example is PI9 (a granzyme B inhibitor) that is expressed in the cytosol of activated [T cells](#) and natural killer (NK) cells and is thought to protect these killer cells from their own pro- **apoptotic** proteinase, granzyme B (23, 24). Other intracellular serpins may protect cells from [thiol proteinases](#) of [lysosomes](#) (25).

The idea that intracellular serpins could regulate cell growth was strengthened by the identification of maspin, a serpin that is expressed in mammary epithelial cells and downregulated on malignant differentiation. When transformed cells are induced to express maspin, their migratory capacity is markedly reduced (26). The mechanism for this effect is unclear, as this serpin is found both intra- and extracellularly. Furthermore, it is uncertain whether maspin is capable of acting as a proteinase inhibitor. Other intracellular serpins include plasminogen activator inhibitor 2, squamous cell carcinoma antigens 1 and 2, monocyte–neutrophil elastase inhibitor, and bomapin (present in haemopoietic cells) (22).

The other group of intracellular serpins is expressed by poxviruses. Orthopoxviruses produce three serpins: SPI-1, SPI-2 (*crmA*), and SPI-3. Studies of viruses deficient in one or other serpins suggest that they have different activities. SPI-2 (*crmA*) inhibits the [caspase](#) interleukin 1- $\beta$ -converting enzyme, thereby preventing production of mature interleukin 1 by the host cell and dampening the immune response to infection (27). When overexpressed in mammalian cells, SPI-2 (*crmA*) inhibits Fas- and [tumor necrosis factor](#)-mediated apoptosis, but it is uncertain whether this effect occurs physiologically (28). SPI-1 inhibits virus-induced apoptosis, and SPI-3 surprisingly seems to modulate cell–cell fusion during **virus infection** (29, 30).

#### Bibliography

1. R. W. Carrell and D. R. Boswell (1986) in *Proteinase Inhibitors*, A. Barrett and G. Salvesen, eds., Elsevier Biomedical Press, Amsterdam, pp. 403–419.
2. J. Potempa, E. Korzus, and J. Travis (1994) *J. Biol. Chem.* **269**, 15957–15960.
3. R. Huber and R. W. Carrell (1989) *Biochemistry* **28**, 8951–8966.
4. H. Loebermann, R. Tokuoka, J. Deisenhofer, and R. Huber (1984) *J. Mol. Biol.* **177**, 531–556.
5. H. T. Wright and J. N. Scarsdale (1995) *Proteins* **22**, 210–225.
6. R. Engh, R. Huber, W. Bode, and A. Schulze (1995) *Trends Biotechnol.* **13**, 503–510.
7. M. Wilczynska, M. Fa, P.-I. Ohlsson, and T. Ny (1995) *J. Biol. Chem.* **270**, 29652–29655.
8. J. Whisstock, R. Skinner, and A. M. Lesk (1998) *Trends Biochem. Sci.* **23**, 63–67.
9. J. Mottonen, A. Strand, J. Symersky, R. M. Sweet, D. E. Danley, K. F. Geoghegan, R. D. Gerard, and E. J. Goldsmith (1992) *Nature* **355**, 270–273.



10. L. Jin, J. P. Abrahams, R. Skinner, M. Petitou, R. N. Pike, and R. W. Carrell (1997) *Proc. Natl. Acad. Sci. USA* **94**, 14683–14688.
11. P. E. Stein and R. W. Carrell (1995) *Nature Struct. Biol.* **2**, 96–113.
12. R. W. Carrell and M. C. Owen (1985) *Nature* **317**, 730–732.
13. W. Chang and D. Lomas (1998) *J. Biol. Chem.* **273**, 3695–3701.
14. R. W. Carrell and P. E. Stein (1996) *Biol. Chem. Hoppe-Seyler* **377**, 1–17.
15. P. R. Elliott, D. A. Lomas, R. W. Carrell, and J. P. Abrahams (1996) *Nature Struct. Biol.* **3**, 676–681.
16. D. M. Tollefsen (1997) in *Chemistry and Biology of Serpins*, F. C. Church et al., eds., Plenum Press, New York, pp. 35–44.
17. M. T. Stubbs and W. Bode (1995) *Trends Biochem. Sci.* **20**, 23–28.
18. R. W. Carrell and D. A. Lomas (1997) *Lancet* **350**, 134–138.
19. C.-B. Laurell and S. Eriksson (1963) *Scand. J. Clin. Lab. Invest.* **15**, 132–140.
20. D. A. Lomas, D. L. Evans, J. T. Finch, and R. W. Carrell (1992) *Nature* **357**, 605–607.
21. F. L. Scott, P. B. Coughlin, C. Bird, L. Cerruti, J. A. Hayman, and P. Bird (1996) *J. Biol. Chem.* **271**, 1605–1612.
22. E. Remold-O'Donnell (1993) *FEBS Lett.* **315**, 105–108.
23. C. A. Sprecher, K. A. Morgenstern, S. Mathewes, J. R. Dahlen, S. K. Schrader, D. C. Foster, and W. Kisiel (1995) *J. Biol. Chem.* **270**, 29854–29861.
24. J. R. Sun, C. H. Bird, V. Sutton, L. McDonald, P. B. Coughlin, T. A. Dejong, J. A. Trapani, and P. I. Bird (1996) *J. Biol. Chem.* **271**, 27802–27809.
25. C. Schick, P. A. Pemberton, G. P. Shi, Y. Kamachi, S. Cataltepe, A. J. Bartuski, E. R. Gornstein, D. Bromme, H. A. Chapman, and G. A. Silverman (1998) *Biochemistry* **37**, 5258–5266.
26. Z. Zou, A. Anisowicz, M. J. C. Hendrix, A. Thor, M. Neveu, S. Sheng, K. Rafidi, E. Seftor, and R. Sager (1994) *Science* **263**, 526–529.
27. C. A. Ray, R. A. Black, S. R. Kronheim, T. A. Greenstreet, P. R. Sleath, G. S. Salvesen, and D. J. Pickup (1992) *Cell* **69**, 597–604.
28. M. Tewari and V. M. Dixit (1995) *J. Biol. Chem.* **270**, 3255–3260.
29. M. A. Brooks, A. N. Ali, P. C. Turner, and R. W. Moyer (1995) *J. Virol.* **69**, 7688–7698.
30. P. C. Turner and R. W. Moyer (1992) *J. Virol.* **66**, 2076–2085.

### **Suggestions for Further Reading**

31. F. C. Church, ed. (1997) *Chemistry and Biology of the Serpins*, Plenum Press, New York.
32. P. G. W. Gettins, P. A. Patston, and S. T. Olson (1996) *Serpins: Structure, Function and Biology*, R. G. Landes Company, Georgetown, TX.

## **Serum Albumin**

Serum proteins are operationally classified into a, b, and g fractions, plus [albumin](#), according to their relative mobilities upon [electrophoresis](#) at pH 8.6. With an isoelectric point of 4.9, serum albumin is the most acidic protein in serum, except for a small amount of pre-albumin, and the only albumin. Human serum albumin (HSA) and bovine serum albumin (BSA) are two of the most studied

proteins, because of their abundance, and the most popular standard proteins in many protein assay systems. Some of the important physical constants of bovine serum albumin are given in Table 1(1).

**Table 1. Physical Parameters of Serum Albumin**

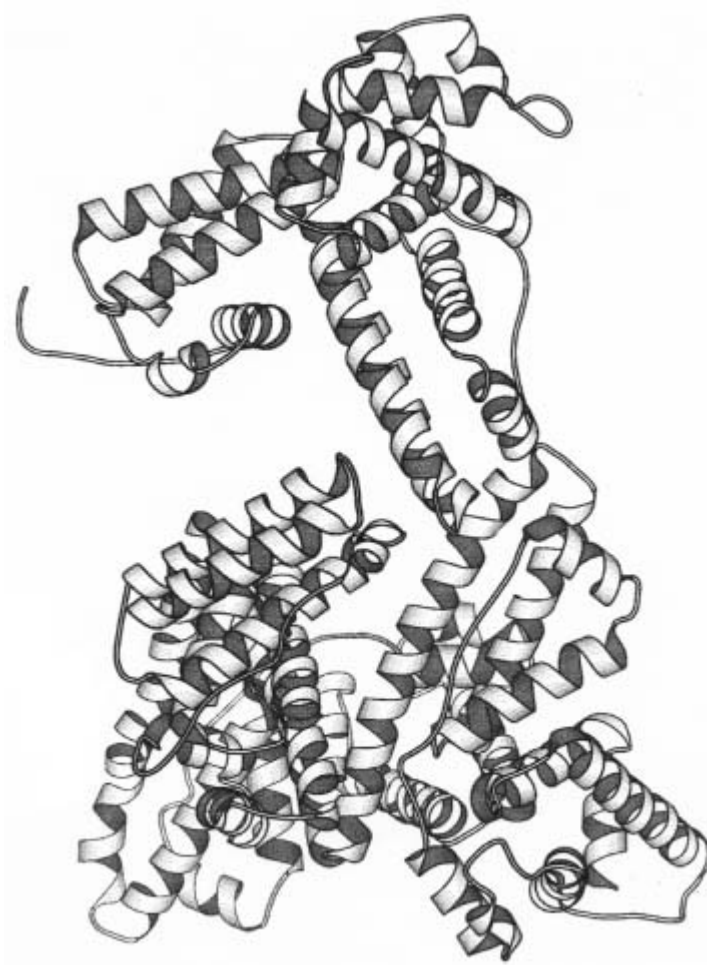
| Property  | Bovine | Human  | Rat    |
|---|--------|--------|--------|
| Molecular weight  |        |        |        |
| From composition  | 66,267 | 66,439 | 65,871 |
| From physical data  | 66,700 | 69,000 | 64,300 |
| Sedimentation coefficient (S)                                 |        |        |        |
| Monomer   | 4.5    | 4.6    | 4.2    |
| Dimer   | 6.7    | 6.7    |        |
| Diffusion coefficient ( $10^{-11} \text{m}^2 \text{s}^{-1}$ ) | 5.9    | 6.1    |        |
| Partial specific volume ( $\text{cm}^3/\text{g}$ )            | 0.733  | 0.733  |        |
| Intrinsic viscosity ( $\text{cm}^3/\text{g}$ )                | 0.041  | 0.042  |        |
| Refractive index increment (578 nm, $\times 10^3$ )           | 1.90   | 1.89   |        |
| Optical absorbance (279 nm, 1 mg/ml)                          | 0.667  | 0.531  | 0.59   |
| Reproduced from (1) with permission.                          |        |        |        |

Albumin is the most abundant protein in the serum, occurring at a concentration of about 35 to 45 g/L and comprising almost 50 to 60% of the total plasma proteins. Because of this and its relatively small molecular weight, it accounts for about 80% of the colloid osmosis of the blood. Albumin is produced in the liver as pro-albumin at the rate of approximately 120 to 200 mg/kg body weight daily for an adult male and 120 to 150 mg/kg for a female. In the first several years after birth, the production level is highest at 180 to 300 mg/kg weight. Albumin is one of the most soluble proteins and can be dissolved up to a concentration 30% (w/v) in water. To prepare serum albumin, the globulins of serum are precipitated at 50% saturation of ammonium **sulfate**, the remaining solution is adjusted to pH 4.4, and the precipitate is collected and dialyzed against water. Most such preparations contain various degrees of micro-heterogeneity: (1) mercaptoalbumin with one thiol group and nonmercaptoalbumin, with the thiol group blocked covalently, (2) different combinations of disulfide bonds due to disulfide interchange reactions, (3) varying amounts of bound fatty acids and other ligands that are mostly insoluble in plasma in the absence of albumin.

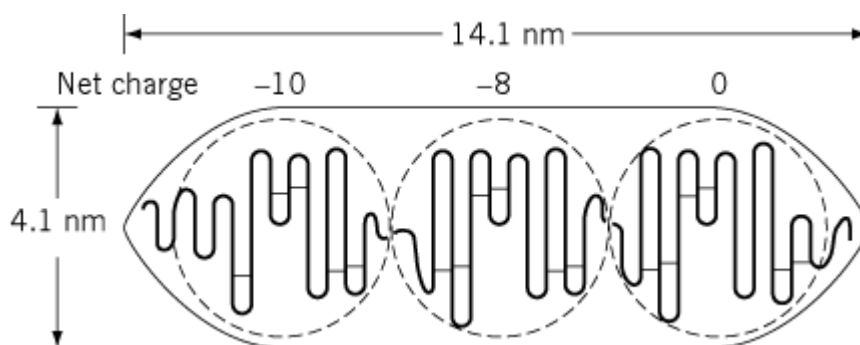
Human serum albumin is made of 585 amino acid residues, and its three-dimensional [protein structure](#) is known (2, 4) (Fig. 1). It is usually a monomeric protein, with a molecular weight of 66,300, although some dimer is almost always present. It has 17 [disulfide bonds](#) and a single cysteine [thiol group](#). The amino acid sequence shows the presence of three regions with high sequence similarities, residues 1 to 191, 192 to 384, 385 to 585, and these correspond to three **domains** in the native conformation (Fig. 2). The gross structure of albumin may be approximated by a prolate of revolution with the major and minor axes being 4.1 and 14.1 nm, respectively. About 50% of the amino acid residues are in the **alpha-helical** conformation.

**Figure 1.** The three-dimensional structure of human serum albumin. It is composed of three homologous regions.

Based on PDB data.



**Figure 2.** A model of serum albumin, emphasizing the tandem repeat of three semiglobular units that are related to each other in primary and tertiary structures and in function. Reproduced from (3) (Fig. 5, p. 179) with permission.



Serum albumin is relatively stable against **denaturants** such as [urea](#) or **guanidinium** chloride, but at pH 4 to 4.5 it changes from the native N state to a structurally loose and electrophoretically fast-moving species known as the F state. The transition from N to F accompanies an increase of about 40 titratable [carboxyl groups](#) that are apparently masked in the N state. The N state is stable in alkaline solution up to pH 10 to 11 (3).

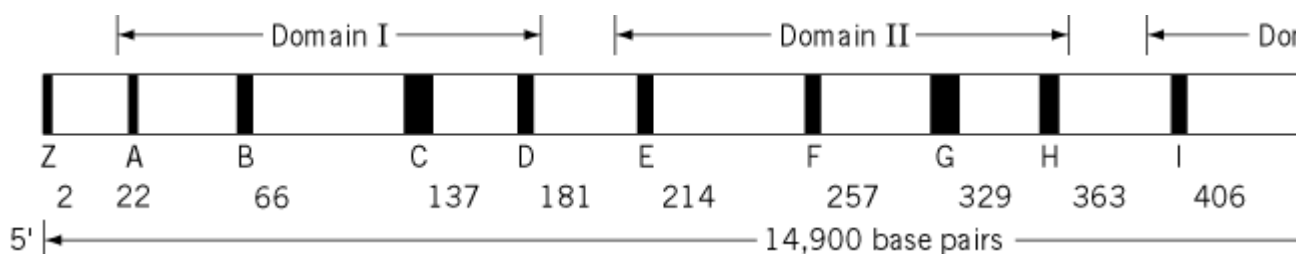
Serum albumin functions as a carrier of insoluble long-chain [fatty acids](#) (mainly stearic, palmitic, and oleic acids) having two very strong and five strong binding sites, with association constants of  $1.1 \times 10^8$  and  $4.0 \times 10^6 M^{-1}$ , respectively. Other ligands known to bind to albumin include bilirubin, chloride ion, testosterone, salicylic acid, phenoxymethyl- penicillin, and synthetic dyes. The function of albumin may be replaced by other blood proteins, mainly those in globulin fractions, since it has been shown that some patients of analubinaemia who completely lack serum albumin have led almost normal lives.

Only one copy of the albumin **gene** has been detected in the human and rat [genomes](#). The gene for a-fetoprotein is in close proximity. The albumin gene consists of 15 exons separated by 14 **introns** and is 16,000 base pairs long, even though the total coding region for albumin requires only 1830 base pairs. The protein is initially produced as pro-albumin, which has an extra hexapeptide, Arg—Gly—Val—Phe—Arg—Arg, at its [N-terminus](#). It is synthesized in the liver and secreted into the blood after cleavage of the hexapeptide. Substitution of the penultimate Arg residue prevents cleavage of the hexapeptide from pro-albumin, but its secretion into the blood is not hampered. The average half-life of albumin in blood is about 19 days.

The structure of the albumin **gene** is given in Fig. 3. In the case of humans, the gene is located at the proximal end of the long arm of chromosome 14, in bands q11-22 (5), where the genetically linked a-fetoprotein and vitamin D-binding protein genes are also located. In the case of mouse, the albumin and a-fetoprotein genes are in tandem on chromosome 5, and the gene for a-fetoprotein lies 13.5 kbp downstream of the albumin gene, showing a close genetic linkage. The albumin gene has 15 exons and 14 **introns**. Expression of the albumin gene is codominant, with both genes present in the diploid genome showing copenetrance. **Polymorphism** of human albumin is rare.

Bisalbuminemia, where both the normal albumin A and a variant are present in the plasma, occurs with a frequency of 0.0003 to 0.001. In a variant called albumin B, Gly570 is substituted by Lys, and in albumin Mexico-2, Asp550 by Gly. Proalbumin with unexcised hexapeptide has also been found in the circulation, together with normal albumin. Variations in the amino acid sequence in the hexapeptide pro-region are also known. Familial dysalbuminemic hyperthyroninemia, in which there is an abnormally tight binding of thyroxin by albumin, is caused by a single amino acid replacement. In more severe cases, analbuminemia, patients show a complete lack of circulating albumin. By using a strain of rats that have a similar lack of albumin, it has been shown that the animals have an intact albumin gene, but the processing of the **transcriptional** products to [messenger RNA](#) (mRNA) is prevented (see [RNA Splicing](#)). Analbuminaemia humans and rats may seem healthy under normal conditions, but their health can be threatened under stress or disease conditions (3).

**Figure 3.** The gene structure of rat serum albumin. The solid bars labeled with capital letters are exons, and the 14 intron correspond to the terminal amino acid residue in each exon. Reproduced from (3) (Fig. 9, p. 209) with permission.



The counterpart of serum albumin in the fetal serum is a glycoprotein known as a-fetoprotein, which has limited sequence similarity with albumin. Phylogenetically, albumin is found in the serum of amphibians and higher vertebrates, and less frequently in bony fish (3).

## Bibliography

1. G. E. Schulz and R. H. Shirmer (1979) *Principles of Protein Structure*, Springer-Verlag, New York.
2. M. Perutz (1992) *Protein Structure: New Approach to Disease and Therapy*, W. H. Freeman and Co., New York.
3. T. Peters, Jr. (1985) Serum albumin. *Adv. Prot. Chem.* **37**, 161–245.
4. D. C. Carter and X. M. He (1990) Structure of human serum albumin, *Sci.* **249**, 302–303.
5. P. P. Minghetti, D. E. Ruffner, W. J. Kuang, O. E. Dennison, J. W. Hawkins, W. G. Beattie, and A. Dugaiczky (1986) *J. Biol. Chem.* **261**, 6747–6757.

## Suggestion for Further Reading

6. B. Blombäck and L. A. Hanson (1979) *Plasma Proteins*, John Wiley & Sons, New York.

## Serum Dependence

A large number of serum-free media for [tissue culture](#) of [cell lines](#) have been developed ([1](#), [2](#)), and many of these are now commercially available, but the bulk of cell culture on the laboratory scale is still carried out in the presence of serum. Serum provides the cells with [growth factors](#), attachment factors, hormones, nutrients, minerals, and [trypsin inhibitor](#) activity. It also buffers the medium against pH fluctuations and traces of toxins. On the other hand, it is a poorly defined natural product that carries the risk of contamination with **viruses** and [mycoplasma](#), is variable in quality and constitution, is sometimes in limited supply, and adds complexity to the isolation of cell products from the medium. Over the past 20 years, several recipes have been generated for the replacement of serum, usually with defined supplements containing selenium, [insulin](#), [transferrin](#), steroids, and lipids, and often with additional minerals and intermediary metabolites. The main advantages of these serum-free recipes are definition, reproducibility, and selectivity for the cell type that may be propagated. Serum dependence is also a feature that is most strongly expressed in normal cells, as one of the recognized features of [Neoplastic Transformation](#) is a reduction in the requirement for serum, reflecting the increase in autonomous growth control typical of transformed cells.

### 1. Constitution of Serum

Data on serum analysis are traditionally derived from clinical laboratories and depend on conventional analyses and assays based on the major constituents, and those that show fluctuations that are of interest in diagnosis of disease or nutritional deficiency. It is now possible, however, to obtain data from serum suppliers, with greater emphasis on those constituents of greater significance to cell culture, and whose determination usually requires [immunoassays](#) or functional bioassays.

The major constituents of serum are listed in Table [1](#). There are undoubtedly other constituents that may be of importance, but these are the main ones that have been identified. Those that have been shown to be required in serum-free medium, in addition to the regular constituents found in serum-containing medium, are given in boldface type. The major roles of these constituents are listed in Table [2](#).

**Table 1. Constituents<sup>a</sup> of Serum**

|   | <i>Concentration Range<br/>of Each<sup>b</sup></i> |                     |
|---|--|---------------------|
| <a href="#">Albumin</a>   | 20–50 mg/mL  |                     |
| Amino acids   | 0.01–1.0 $\mu$ M                                   |                     |
| Calcium   | 4–7 mM   |                     |
| Chlorides   | 100 $\mu$ M  |                     |
| Cholesterol   | 10 $\mu$ M   |                     |
| <a href="#">Fatty acids</a>   | 0.1–1.0 $\mu$ M                                    |                     |
| Fetuin  |  | In FB only          |
| <a href="#">Fibronectin</a>   | 1–10 $\mu$ g/mL                                    |                     |
| Globulins   | 1–15 mg/mL   |                     |
| Glucose   | 0.6–1.2 mg/ml                                      |                     |
| <a href="#">Growth Factors</a> : EGF, PDGF, IGF,FGF,<br>IL-1, IL-6, Insulin                           | 1–100 ng/mL  |                     |
| <b>Hydrocortisone</b>   | 10–200 nM  |                     |
| Hexosamine  | 0.0–1.2 mg/mL                                      | Highest in<br>human |
| <b>Iron</b>   | 1–10 $\mu$ M                                       |                     |
| Lactic acid   | 0.5–2.0 mg/mL                                      | Highest in<br>FB    |
| <b>Linoleic acid</b>  | 0.01–0.1 $\mu$ M                                   |                     |
| Phospholipids   | 0.7–3.0 mg/mL                                      |                     |
| <b>Polyamines</b> putrescine,spermidine   | 0.1–1.0 $\mu$ M                                    |                     |
| Potassium   | 5–15 mM  |                     |
| <a href="#">Proteinase inhibitors</a> :a <sub>1</sub> -antitrypsin, a <sub>2</sub> -<br>macroglobulin | 0.5–2.5 mg/mL                                      |                     |
| Protein, total  | 40–80 mg/mL  |                     |
| Pyruvic acid  | 2–10 $\mu$ g/mL                                    |                     |
| <b>Selenium</b>   | 0.01 $\mu$ M                                       |                     |
| Sodium  | 135–155 mM   |                     |
| Total lipids  | 2–10 mg/mL   |                     |
| <a href="#">Transferrin</a>   | 2–4 mg/mL  |                     |
| <b>Tri-iodotyrosine</b>   |  |                     |
| Urea  | 170–300 $\mu$ g/mL                                 |                     |
| Vitamins  | 10 ng–10 $\mu$ g/mL                                |                     |
| Vitamin A   | 10–100 ng/mL                                       |                     |
| <b>Zinc</b>   | 0.1–1.0 $\mu$ M                                    |                     |

<sup>a</sup> Constituents in bold type are those that have been used to supplement serum-free media.

<sup>b</sup> Concentration range is very approximate and only intended to convey the order of magnitude.

**Table 2. Major Functions of Serum in Cell Culture**

---

| <b>Function</b>  | <b>Constituent</b>   |
|--|--|
| Mitogenic activity   | Growth factors, platelet-derived growth factor(PDGF) and insulin-like growth factors (IGF-1,IGF-2) |
| Attachment factors   | Fibronectin, fetuin  |
| Antagonize trypsin used in subculture                                  | $\alpha_1$ -Antitrypsin, $\alpha_2$ -macroglobulin   |
| Bind toxins  | Albumin  |
| Detoxify free radicals   | Selenium   |
| Stimulate nutrient uptake, eg, glucose and amino acids                 | Insulin  |
| Bind and transport iron  | Transferrin  |
| Provide intermediary metabolites                                       | Pyruvate, a-ketoglutarate, adenosine, etc  |
| Viscosity to buffer cells against mechanical damageduring manipulation | Albumin  |

---

Many of the growth factors in serum are derived from platelets; hence serum from naturally clotted blood is usually superior to serum separated from the cellular constituents by centrifugation. The growth factors released by platelets include [platelet-derived growth factor](#) (PDGF) and [transforming growth factor b](#) (TGFb), which have effects that are cell-type-specific. PDGF is a mitogen for mesodermally derived cells, such as fibroblasts, and glial cells, such as astrocytes; TGFb is cytostatic for many epithelial cells. Hence, the presence of serum is likely to favor fibroblastic cells, rather than epithelium, in cultures from normal tissues, and this has added weight to the argument in favor of using serum-free selective media to grow normal epithelial cells. Transformed cells may be exceptional, as they often produce a spectrum of autocrine growth factors, giving them a degree of autonomous control over mitogenesis. In addition, they may also produce TGFa, which reverses the cytostatic effect of TGFb and gives it a mitogenic and transformation-like effect. The inhibitory effect of TGFb can also be reduced by the presence of a feeder layer of 3T3 primitive embryonic mouse mesodermal cells, which appears to degrade or block activation of TGFb in cocultures with epithelium. Such feeder layers are frequently used to generate cultures of epithelial cells, such as epidermal keratinocytes (3), in serum-containing medium.

Growth of cells in the absence of serum frequently requires modification of the substrate to enable cell attachment and, for normal cells at least, the cell spreading that is necessary for cell proliferation. This modification can be a simple chemical treatment with, for example, poly-D-lysine, which neutralizes the negatively charged plastic and produces a slight positive charge. The incorporation of **spermidine** or putrescine into the medium (4) may play a similar role. Other protocols require [extracellular matrix](#) constituents, such as [collagen](#), [fibronectin](#), or [laminin](#) (5), which will interact with specific [integrin](#) receptors on the cell surface and promote adhesion to the substrate. In many cases, the cells may be capable of generating these matrix constituents themselves, as is suggested by the improved survival of cells reseeded into the flask from which they were released by **trypsin** treatment, or on to a substrate conditioned by other cells, such as vascular endothelium (6). It has been suggested that serum contains many such attachment factors that condition the substrate, modify its charge, and allow attachment of extracellular matrix molecules important in cell adhesion. Fibronectin is known to be present in serum, although in a modified form, and may participate in this.

## 2. Effect of Transformation

Reduced serum dependence is frequently used as a criterion for the identification of transformed cells. While normal cells will block at a restriction point in the G<sub>1</sub> phase of the [cell cycle](#) in the absence of serum and will re-enter cycle when serum is restored (7, 8), transformed cells will tend to progress into S phase regardless of the serum concentration, although they may arrest at later points in the cycle because of nutritional deficiencies. This distinction is created by the reduced dependence of transformed cells on exogenous growth factors (see [Neoplastic Transformation](#)) as a result of the expression of autonomous control mechanisms, such as autocrine growth factor production, for example, PDGF b-subunit in gliomas (9), or modifications in [signal transduction](#) leading to unregulated, permanently active induction of mitogenesis, such as *ras* mutations in colon carcinoma (10) (see [Oncogenes, Oncoproteins](#)).

Transformed cells are also able to proliferate without the degree of cell spreading required with normal cells, even to the extent of proliferating in suspension. Hence, there is reduced dependence on attachment factors present in serum for cells to initiate proliferation. Transformed cells are not devoid of [cell adhesion molecules](#) (CAMs), but they frequently have modifications affecting the selectivity of attachment and, possibly more important, alterations to the extracellular domain, which modify the interactions with the [actin](#) cytoskeleton, which, in turn, may make redundant cell spreading and the formation of adhesion plaques.

## 3. Serum-Free Media

Although serum supplementation is still widely used, there are several problems associated with the use of serum, most of which derive from its variability and undefined nature. Although a fairly precise analysis can be performed of the major constituents of serum, it is the minor constituents, those present in small amounts, but with high specific activity, that are difficult to detect and most likely to be variable. Some of these are listed in Table 1 and include hormones and growth factors. Although batch testing minimizes this problem, each batch has a limited shelf life, and it is unlikely that the next batch will be identical. Hence, although the nutritional content and physical properties of serum can be reproduced fairly accurately, the signaling content cannot, giving rise to considerable physiological variation from batch to batch. With the dramatic increase in the use of cell culture by the pharmaceutical industry, acute worldwide shortages have been avoided only by the development of low serum or serum-free formulations by drug companies, many of which, unfortunately, are not available in the public domain. In order to meet the requirements of Good Laboratory Practice (The United Kingdom Compliance Programme: Department of Health, London, 1989) and Good Manufacturing Practice (Medicines Control Agency, "Rules and Guidance for Pharmaceutical Manufacturers and Distributors 1997," The Stationery Office Ltd., London, 1977), pharmaceutical companies have also been obliged to address the problems of viral contamination, which most research laboratories still choose to ignore. As the elimination of viruses from serum has proved to be extremely difficult, the obvious alternative is serum-free media, which have the additional advantages of minimizing the risk of [mycoplasma](#) infections and reducing the contamination of cell products with serum proteins, thereby facilitating downstream processing.

One of the major advantages of serum-free media has been the development of media that are selective for individual cell types. The range is now quite extensive, and many are available commercially from companies such as Bio Whittaker, Sigma, Invitrogen (Gibco), and ICN. Some examples are given in Table 3; for a more extensive list see Barnes et al. (11), Freshney, (12), Davis, (13) and Jayme and Gruber (2).

**Table 3. Serum-Free Selective Media**



| Medium        | Cell Type                     | Reference  |
|---------------|-------------------------------|--|
| MCDB 105      | Human fibroblasts             | McKeehan, et al. (1977) <i>InVitro</i> <b>13</b> , 399;<br>Bettger et al.(1981) <i>Proc. Natl. Acad. Sci. USA</i> <b>78</b> 5588 |
| MCDB 110      | WI38, MRC5,<br>IMR90          |  |
| MCDB 202      | Chicken fibroblasts           | McKeehan et al. (1977), <i>InVitro</i> <b>13</b> , 399   |
| MCDB 402      | Mouse 3T3 cells               | Shipley and Ham, (1981) <i>InVitro</i> <b>17</b> , 656   |
| MCDB 411      | Mouse neuroblastoma C1300     | Agy, et al. (1981) <i>In Vitro</i> <b>17</b> , 671   |
| MCDB 153      | Human keratinocytes           | Boyce and Ham, (1983) <i>J. Invest.Dermatol.</i> <b>81</b> , 33s   |
| MCDB 170      | Mammary epithelium            | Hammond et al. (1984) <i>Proc. Natl.Acad. Sci. USA</i> <b>81</b> , 5435  |
| Iscove's      | Hemopoietic cells             | Iscove and Melchers (1978) <i>J. Exp.Med.</i> <b>147</b> , 923   |
| LHC           | Bronchial epithelium          | Lechner and LaVeck (1985) <i>J. Tissue Cult. Meth.</i> <b>9</b> , 43   |
| HITES         | Small cell lung cancer        | Carney et al. (1981) <i>Proc. Natl.Acad. Sci. USA</i> <b>78</b> , 3185   |
| WAJC 404      | Prostatic epithelium          | Chaproniere and McKeehan (1986) <i>Cancer Res.</i> <b>46</b> , 819   |
| DMEM:F12, 1:1 | Various, with supplementation | Sato (1979) in <i>Methods inEnzymology</i> , W. B. Jakoby and I. H. Pastan, eds., AcademicPress, New York, pp. 94–109            |

Unfortunately, a move to serum-free medium is not without its difficulties. Individual recipes may be required for each cell type maintained, and the time required to develop media for specific cell types is quite significant, possibly several years. The problem is diminishing with an increase in the supply of serum-free media from commercial sources, but it is still substantially more expensive than serum-containing media and may be beyond the reach of laboratories on modest budgets. There are also problems with subculture, as serum is normally responsible for the inhibition of trypsin that is used to liberate cells when they are reseeded. Trypsin damage can be reduced by using purified trypsin at a reduced temperature ([14](#)) and/or by incorporating a [trypsin inhibitor](#) in the medium at reseeded.

The major additions to medium to replace serum are [insulin](#), selenium, and iron-saturated [transferrin](#), and these are to be found in most serum-free formulations. In addition, hydrocortisone, **cholera toxin**, or isoprenaline are often added to increase proliferation of epithelial cells. Hydrocortisone probably acts by increasing attachment via the induction of proteoglycans, which may also activate **cytokines** and growth factors ([15](#), [16](#)), cholera toxin and isoprenaline increase the intracellular concentration of [cyclic AMP](#), which is mitogenic in some epithelial cells. Lipid precursors, such as cholesterol and linoleic acid, ethanolamine, and phosphoethanolamine, high density **lipoproteins** (HDL), and crude lipid preparations, such as soya bean lipid, are often included and may contribute to the biosynthesis of [membranes](#) or [second messengers](#). Thiol compounds such as [b-mercaptoethanol](#) are often used to inhibit oxidative stress induction from free radicals. In addition to

insulin and hydrocortisone, other growth factors and hormones include [epidermal growth factor](#) (EGF), [fibroblast growth factor](#) (FGF), **insulin-like growth factor-1** (IGF-1, somatomedin C) and -2 (IGF-2, multiplication stimulating activity (MSA)) ([17](#)), PDGF, follicle-stimulating hormone (FSH), prolactin, tri-iodotyrosine, estradiol, and prostaglandins such as PGF<sub>2a</sub>. In addition to selenium and iron, other trace elements include copper, zinc, manganese, molybdenum, tin, nickel, vanadium, and silicon. Finally, in the absence of the detoxifying effect of serum proteins, there is a requirement for highly purified reagents and water.

## Bibliography

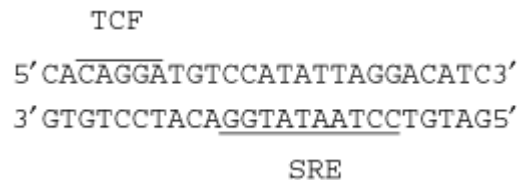
1. R. Maurer (1992) in *Animal Cell Culture, a Practical Approach*, 2nd ed., R. I. Freshney, ed., IRL Press at Oxford Univ. Press, Oxford, pp. 15–46.
2. D. W. Jayme and D. F. Gruber (1994) in *Cell Biology, a Laboratory Handbook*, J. E. Celis, Academic Press, New York, pp. 18–24.
3. J. G. Rheinwald and H. Green (1975) *Cell* **66**, 331–344.
4. W. J. Bettger, S. T. Boyce, B. J. Walthall, and R. G. Ham. (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5588–5592.
5. K. A. Elliget, and J. F. Lechner (1992) in *Culture of Epithelial Cells*, R. I. Freshney, ed., Wiley-Liss, New York, pp. 181–196.
6. D. Gospodarowicz, D. Delgado, and I. Vlodavsky (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4094–4098.
7. H. Chang, and R. Baserga (1977) *J. Cell. Physiol.* **92**, 333–343.
8. A. Lindgren, B. Westermark, and J. Ponten (1975) *Exp, Cell Res.* **95**, 311–319.
9. F. S. Vassbotn, A. Ostman, N. Langeland, H. Holmsen, B. Westermark, C.-H. Heldin, and M. Nister (1994) *J. Cell. Physiol.* **158**, 381–389.
10. R. A. Weinberg, ed. (1989) *Oncogenes and the Molecular Origins of Cancer*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
11. D. W. Barnes, D. A. Sirbasku, G. H. Sato, eds. (1984) *Cell Culture Methods for Molecular and Cell Biology*. 4 vols., Alan R. Liss, New York.
12. R. I. Freshney (2000) *Culture of Animal Cells, a Manual of Basic Technique*, Wiley-Liss, New York, pp. 106–109.
13. J. M. Davis, ed. (1994) *Basic Cell Culture*, IRL Press at Oxford Univer. Press, Oxford, pp. 58–91.
14. W. L. McKeehan (1977) *Cell Biol. Intl. Rep.* **1**, 335–343.
15. N. Yevdokimova, and R. I. Freshney (1997) *Bri J. Cancer*, **76**, 261–289.
16. M. Klagsbrun, and A. Baird (1991) *Cell* **67**, 229–231.
17. N. C. Dulak, and H. M. Temin (1973) *J. Cell. Physiol.* **81**, 153–160.

## Serum Response Element

The serum response element (SRE) was originally identified within the **promoter** of the c-fos **proto-oncogene** and is responsible for the activation of this **gene** following addition of serum to cultured cells ([1](#)). Subsequently an SRE was also identified in a number of other genes whose [transcription](#) is enhanced by treatment of cells with serum or [growth factors](#) ([2](#), [3](#)). The SRE has a **consensus sequence** of the form CC A/T<sub>6</sub> GG, in which two C residues are followed by a run of six residues

that can be either A or T and then followed by a further two G residues. The sequence of the SRE within the c-Fos promoter is illustrated in Figure 1.

**Figure 1.** Structure of the c-Fos SRE and the adjacent binding site for TCF. The SRE is underlined and the TCF site is overlined.



Subsequent studies demonstrated that the SRE acts by binding a ubiquitous 67-kDa protein known as the serum response factor (SRF). Although the binding of SRF to the SRE is essential for the response to growth factors in serum, it is not sufficient for this to occur. Thus, further studies identified a 62-kDa protein known as ternary complex factor (TCF), which could only bind to the SRE in partnership with SRF and was required for SRE function (4). Subsequent studies identified TCF as a member of the *ets* family of DNA-binding proteins (5). Although a binding site for this factor is located adjacent to the SRE in the c-Fos promoter (Fig. 1), TCF cannot bind to this site unless it has also formed a protein-protein interaction with SRF bound at the SRE. Hence binding of TCF to its binding site is dependent on the prior association of SRF with the SRE.

Interestingly, TCF is a substrate for phosphorylation by a number of different kinases that are activated following addition of serum or growth factors to cultured cells. Such phosphorylation enhances both the ability of TCF to associate with SRF and its ability to activate transcription (6). Hence the addition of serum or growth factors to cells results in kinase activation, which in turn stimulates the activity of TCF and hence causes gene activation by the SRF/TCF complex bound to the SRE and its adjacent sequences.

#### Bibliography

1. R. Treisman (1986) *Cell* **46**, 567–574.
2. R. Treisman (1990) *Semin. Cancer Biol.* **1**, 47–58.
3. R. Treisman (1987) *EMBO J.* **6**, 2711–2717.
4. P. E. Shaw, H. Schroter, and A. Nordheim (1989) *Cell* **56**, 563–572.
5. A. D. Sharrocks, A. L. Brown, Y. Ling, and P. R. Yates (1997) *Int. J. Biochem. Cell Biol.* **29**, 1371–1387.
6. H. Gille, A. D. Sharrocks, and P. E. Shaw (1992) *Nature* **358**, 414–417.

#### Suggestions for Further Reading

7. R. Triesman (1992) The serum response element. *Trends Biochem. Sci.* **17**, 423–426.
8. C. S. Hill and R. Triesman (1995) Transcriptional regulation by extra cellular signals: Mechanisms and specificity. *Cell* **80**, 199–211.
9. R. Triesman (1995) Journal to the surface of the cell: Fos regulation and the SRE. *EMBO J.* **14**, 4905–4913.

## Sex

Biological [evolution](#) is driven by the [natural selection](#) of the fittest individuals that pass their own **genes** on to the next generation. At the very early stages of the biological evolution, **reproductive** success was enhanced by the invention of a polarity within a species. In higher **eukaryotes**, this has led to the two forms of beings within one species, the [male](#) and [female](#) sexes. The evolutionary advantage of sexual reproduction is not obvious at first sight because the investment is twice as high as with asexual reproduction (1). However, the reproductive success of this strategy in most eutherians is superior to asexual reproduction. The most important advantage of sexual reproduction is the enlarged basis for natural selection. In systems without sexual reproduction and [recombination](#), as in **prokaryotes**, genetic adaptation to a changing environment is achieved by a high **mutation** rate and a short generation interval. The phenomenon of genetic recombination during gametogenesis and fertilization, however, introduces a variability that would otherwise necessarily depend solely on the mutation rate of the genome (2). Even more important in asexual reproduction, the risk of losing beneficial mutations from the population is enhanced by deleterious mutations (1, 3). Therefore, the invention of genetic recombination and sex seems to be one of the basic requirements for the evolution of individuals with longer generation times; consequently, it is responsible for the diversity of higher species.

Sexual reproduction enables a selection that occurs not only on a single path but independently in the two sexes. Further, selection may be different in the two sexes. This has led to distinctly different sex-specific strategies between males and females (see [Male](#); [Female](#)). Regarding a common male behavior, it is a prevalent strategy of males to attempt to spread their genetic information widely in the population. Because all males do so, they suffer from severe competition for females. Especially in higher mammals, this has led to population structures in which some males are dominant enough to transfer their genes on a big scale while the other males remain inferior with no reproductive success at all. This competition has influenced the male **phenotype**. This becomes apparent in the development of male germ cells that are specialized to spread genetic information.

In females, the sex selection is different. The chance to reproduce is limited by the number of **oocytes** that reach the stage of ovulation. The reproductive success in females does not depend so much on competition within females, and so the risk of achieving no reproductive success is considerably lower than in males. On the other hand, the total reproductive success of females in a life span cannot be as high as reproductively successful males because the female's investment in its offspring is generally greater and resources are limited. In many species, females therefore exhibit a phenotype that is adapted not for producing high numbers of gametes but more so to the preservation and the successful raising of the offspring. Even in animals with extracorporal fertilization, the investment of the female is still greater than in the male. This is because the female produces the egg and gene products necessary for the initial growth and differentiation of the new individual. In mammals, the high investment of the female becomes apparent in the long period of maternal care during pregnancy.

### 1. Sex Determination

Sex can be determined on a **chromosomal** level using cytogenetic techniques. During the **meiosis** of spermatogenesis, when the reduction of the **diploid** to the [haploid](#) genome occurs, sex chromosomes are divided; as a consequence, haploid spermatozoa carry either an **X-** or a [Y-Chromosome](#). After fertilization, when the two haploid genomes join, either two X-chromosomes or one X- and one Y-chromosome come together to form either a female or a male, respectively. It is hitherto still usual to attribute the heterogametic sex to the male. **Crossing-over** between the X- and Y-chromosomes is very rare, and the sex-determining factor is almost always located on the Y-chromosome in mammals. In birds, on the other hand, the heterogametic sex is the female. An arrangement of ZW

chromosomal constitution determines female, while ZZ determines male.

At the molecular level, the sex determining factor is a major gene (SRY) that acts as a [transcription factor](#). The SRY gene belongs to a conserved gene family (4). This gene is necessary in mammals for the formation of a male **phenotype**. The conservation among species appears to be greatest in the **HMG** box (high mobility group domain), which encodes a DNA-binding domain and is closely related to the SOX gene family (5). Mutations within this domain go together with sex reversal in XY females (6). Interestingly, almost no sequence similarity is observed in the other regions of the gene or protein, indicating either that these sequences have evolved in species differently or that they are not involved in sex determination (7).

The sex is not determined in all species by the presence of a Y-chromosome carrying the SRY gene. Sex determination can depend on the presence or absence of genes or on their heterozygous or homozygous state. In a number of insects, sex determination is based on the number and the homozygous or heterozygous status of **alleles**. In honeybees, homozygosity of certain alleles leads to male individuals (the drones), which are produced **parthenogenetically**. In most animals, the genetic sex is the basis for the physiological development of the two sexes, as it induces the hormonal and, subsequently, the gonadal sex, as well as the behavioral sex.

Hormonal sex can be determined by measuring the levels of the sex-specific hormones, eg, the amount of testosterone. The levels of hormones have a considerable effect on the phenotypic appearance. During embryogenesis, a cascade of gene products, especially testosterone and MIS (Mullerian-inhibiting hormone), which are expressed by the embryonic testis, influence the undifferentiated sex glands in such a way that they develop a functionally active testis (8). The seminiferous tubules of the developing testis are responsible for the expression of a hormone that subsequently affects the development of the final reproductive tract. Mullerian ducts undergo a regression on MIS expression and form Wolfian ducts. In mammals without the testis-determining factor, the primordial urogenital tract develops to Mullerian ducts to give a female phenotype.

Hormones influence not only the gonadal sex during embryogenesis but also the behavioral sex of adults. High levels of testosterone increase aggression; also, behavioral changes dependent on the function of the [estrogen receptor](#) have been reported (9). In most cases, genetic, chromosomal, gonadal, and behavioral sex go together, although aberrations, such as the [hermaphrodite](#), are also reported.

#### Bibliography

1. R. J. Redfield (1994) *Nature* **369**, 145–147.
2. R. S. Howard, and C. M. Lively (1994) *Nature* **367**, 554–557.
3. J. R. Peck (1994) *Genetics* **137**, 597–606.
4. D. C. Page, R. Moshier, E. M. Simpson, E. M. Fisher, G. Mardon, J. Pollack, B. McGillivray, A. de la Chapelle, and L. G. Brown (1987) *Cell* **51**, 1091–1104.
5. J. W. Foster, and J. A. Graves (1994) *Proc. Natl. Acad. Sci. USA* **91**, 1927–1931.
6. L. S. Whitfield, R. Lovell-Badge, and P. N. Goodfellow (1993) *Nature* **364**, 713–715.
7. P. K. Tucker, and B. L. Lundrigan (1993) *Nature* **364**, 715–717.
8. W. H. Shen, C. C. Moore, Y. Ikeda, K. L. Parker, and H. A. Ingraham (1994) *Cell* **77**, 651–661.
9. S. Ogawa, D. B. Lubahn, K. S. Korach, and D. W. Pfaff (1997) *Proc. Natl. Acad. Sci. USA* **94**, 1476–1481.

#### Suggestion for Further Reading

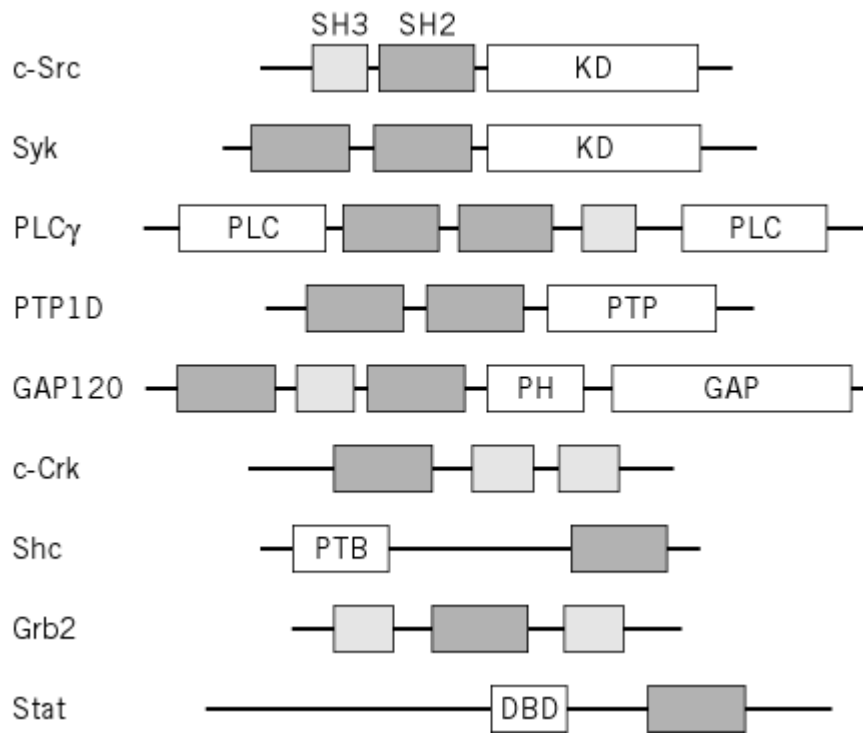
10. A. Ruvinsky (1997) Sex, meiosis and multicellularity. *Acta Biotheor.* **45**, 127–141.

## SH2, SH3 Domains

**Src homology** (SH) domains 2 and 3, SH2 and SH3, are protein **domains** that bind to phosphotyrosine-containing sequence motifs and proline-rich motifs, respectively. The SH2 domain was originally identified as a conserved amino acid sequence in the noncatalytic region of **tyrosine kinases** encoded by viral **oncogenes** such as *v-src*, *v-yes*, and *v-fps* (1) (see [SRC Genes](#)). The most conserved sequence, however, was that of the catalytic region of these kinases, which was referred to as the Src homology 1 (SH1) domain. The SH3 domain is also present in the noncatalytic region of many tyrosine kinases, but it was originally found as a sequence similarity between phosphatidylinositol-specific [phospholipase C](#) and a novel viral oncogene product, Crk (2). Both SH2 and SH3 domains are now known to be present in a wide variety of [signal transduction](#) molecules and have been shown to mediate [protein-protein interactions](#) (3).

Three-dimensional structures of several SH2 and SH3 domains have been elucidated. The SH2 domains consist of two [a-helices](#), with a core of antiparallel [b-sheet](#) in between. The binding of phosphotyrosine-containing peptide to SH2 is bifurcated. The primary interaction is between the phosphotyrosine and basic residues of the SH2 structure. The second interaction, which confers the ligand specificity of each SH2 domain, is provided by the recognition of amino acid residues immediately carboxyl terminal of the phosphotyrosine. The SH3 domains contain a five-stranded antiparallel b-sheet that forms two **hydrophobic** pockets. SH3 ligands have the minimal **consensus sequence** X-Pro-Pro-X-Pro, where X tends to be an aliphatic residue. Each X-Pro pair binds to one of the two hydrophobic pockets of SH3 in either orientation (ie, amino to carboxy terminal or carboxy to amino terminal). Some SH2- and SH3-containing proteins are composed of only SH2 and SH3 domains. These proteins are believed to function as phosphotyrosine-dependent adaptor proteins. Examples of SH2 and SH3 proteins are shown in Figure 1.

**Figure 1.** Schematic structures of SH2 and SH3 proteins. Illustrated are five enzymes (c-Src and Syk tyrosine kinases: PLC $\gamma$ , phospholipase C $\gamma$ ; PTP1D, protein-tyrosine phosphatase 1D; and GAP120), three adapter proteins (c-Crk, Shc, and Grb2), and one [transcription factor](#) (Stat). SH2 and SH3 domains are shown by hatched and gray bars, respectively. KD represents the kinase domains of c-Src and Syk tyrosine kinases. PLC and PTP are the catalytic domains of phospholipase C $\gamma$  and protein-tyrosine phosphatase 1D, respectively. GAP is the domain of GAP120 that is responsible for Ras [GTPase](#) activation; PH, the [pleckstrin-homology domain](#) that binds to phosphoinositides; PTB, the phosphotyrosine-binding domain that is structurally distinct from the SH2 domain; DBD, DNA-binding domain.



### Bibliography

1. T. Pawson (1988) Non-catalytic domains of cytoplasmic protein-tyrosine kinases: regulatory elements in signal transduction. *Oncogene* **3**, 491–495.
2. B. J. Mayer, M. Hamaguchi, and H. Hanafusa (1988) A novel viral oncogene with structural similarity to phospholipase C. *Nature* **332**, 272–275.
3. T. Pawson (1995) Protein modules and signaling networks. *Nature* **373**, 573–580.

### Suggestion for Further Reading

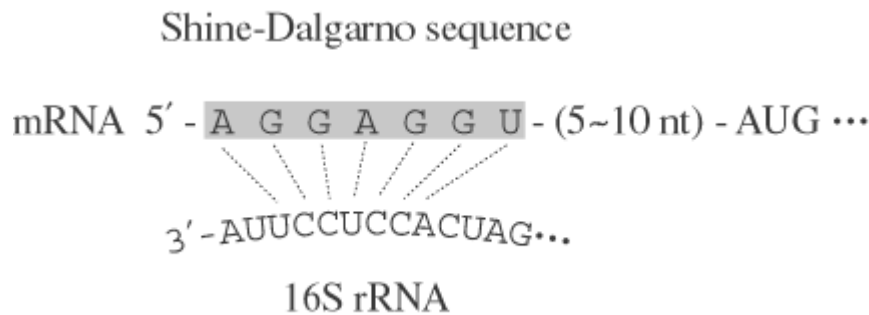
4. D. Lin and T. Pawson (1997) Protein modules in signal transduction. *Trend. Cell Biol.* **7**, pull-out centerfold.

### Shine–Dalgarno Sequence

To initiate protein biosynthesis in bacteria, a bacterial 30S **ribosomal** subunit binds to a specific site on [messenger RNA](#) that is located 5 to 10 bp upstream of the [initiation codon](#). This ribosome-binding sequence is called the Shine–Dalgarno sequence (or S–D sequence), which is named after its discoverers, J. Shine and L. Dalgarno (1). As shown in Fig. 1, it consists of seven or fewer nucleotides that are complementary to the 3' terminus of bacterial 16S rRNA. The diversity of the number of the complementary bases in the S–D region (3 to 9 bases) or of the distance separating it from the initiator codon (5 to 9 bases) accounts, in part, for differences in mRNA translatability. Much evidence supports the idea that the S–D sequence binds to the 16 S rRNA and controls the frequency of translation initiation. For example, oligonucleotides complementary to the 3' terminus of 16S rRNA inhibit translation, and base substitutions in the S–D sequence that reduce the

complementarity to 16S rRNA severely restrict translation, while compensatory changes in the 16S rRNA sequence restore the translatability. It is known in several cases that gene expression is regulated at the translational level by sequestering or exposing the S–D sequence with a regulatory protein or mRNA structure. Although most of the bacterial genes use the S–D sequence to mark the initiation site, there are several exceptions that do not use a S–D sequence. In these cases, alternate sequences that are complementary to other sites of the 16S rRNA are located immediately upstream or downstream of the initiation codon. The latter signal is referred to as the downstream box (2).

**Figure 1.** Base-pairing between the Shine–Dalgarno sequence of mRNA and the 16S rRNA 3' tail.



#### Bibliography

1. J. Shine and L. Dalgarno (1974) Proc. Natl. Acad. Sci. USA **71**, 1342–1346.
2. K. Ito, K. Kawakami, and Y. Nakamura (1993) Proc. Natl. Acad. Sci. USA **90**, 302–306.

#### Short-Chain Dehydrogenases/Reductases SDRS

The SDRs are a large family of dehydrogenase and reductase [enzymes](#) that includes insect [alcohol dehydrogenase](#) and **hormone**-converting and regulatory enzymes in humans and all living forms. They frequently have subunits with ~250 residues and a sequence Tyr-X-X-X-Lys at the active site. For examples, references and further details, see [Alcohol Dehydrogenase \(ADH\)](#).

#### Shotgun Experiments

Shotgun experiments are based on the theory of “randomly divide up into manageable fragments and conquer.” Shotgun experiments use a [genomic library](#) or a [cDNA library](#) as a source of individual clones that, as far as possible, represent the entire [genome](#) or the entire [messenger RNA](#) complement of an organism. It is important for any subsequent experiment that the library be constructed to maximize the probability that any one DNA or mRNA sequence in the starting material is



represented at least once among clones in the library. In practice, this ideal is hard to achieve, because for various reasons certain DNA sequences are hard to clone and are omitted from the library. These reasons include nonuniform distribution of [restriction enzyme](#) sites and the presence of [repetitive DNA](#) and [poison sequences](#). Therefore shotgun experiments are experiments that use as a starting material a library cloned in so as to be inclusive and unbiased. Different types of shotgun experiments are described below.

## 1. Shotgun cloning

Shotgun cloning, one type of shotgun experiment, is an important technique for analyzing whole genomes or any DNA molecule too large to be analyzed directly. The DNA is randomly fragmented into smaller pieces by methods, such as mechanical shearing or partial digestion with a restriction enzyme, and then cloned into a suitable **vector** (see [Cloning](#)). Mechanical shearing is the most random and thus the most unbiased method of fragmenting DNA because the DNA sequence has no effect on breakage sites. However, sheared DNA cannot be cloned directly. The sheared ends must be filled in and linkers must be added before the DNA is cloned into a vector (see [Fill-In Reaction](#) and **Linker fragment**). The set of recombinant vectors that together contain the shotgun-cloned DNA are termed a **library**. Such libraries can be used in a number of different types of shotgun experiments, such as shotgun sequencing, or isolating individual genes or functional sequences, such as **origins of replication** or **promoters** (see [Cloning](#)) (1-3). An example of a typical shotgun cloning experiment is the following: *Bacillus subtilis* DNA was shotgun cloned into the **plasmid** pBR322, the *B. subtilis* library was mobilized into a cold-sensitive mutant of *Escherichia coli*, and a fragment of *B. subtilis* DNA that suppressed the mutant *E. coli* **phenotype** was identified (4).

## 2. Shotgun Sequencing

The theory behind shotgun **sequencing** is that by sequencing large numbers of random library inserts, the original shotgun cloned sequence is pieced together by using the overlaps between the different clones to order them (5). Shotgun sequencing is used to sequence smaller fragments of DNA, such as **yeast artificial chromosome** (YAC) inserts or **cosmid** inserts, to identify genes known to be contained in the fragment of DNA, or it can be used to sequence entire genomes. Many of the procedures involved in shotgun sequencing have been automated, including isolating of the DNA from random clones, carrying out the sequencing reactions, and reading the sequence. A number of entire bacterial genomes have recently been sequenced by shotgun methods, including *Archaeoglobus fulgidus*, the first sulfur-metabolizing organism to have its genome sequenced (6). In the case of *A. fulgidus*, several different libraries were used in the sequencing strategy, including a small-sized insert library and a medium-sized insert library, each created by randomly shearing genomic DNA, and a large-sized insert [lambda phage](#) library created by partially digesting genomic DNA with a [restriction enzyme](#). Because some parts of the genome are more difficult to clone than others, using different libraries created by different methods that select for differently sized DNA molecules is useful to ensure that all sequences are represented in at least one library. The final genomic sequence of *A. fulgidus* was based on randomly determining 29,642 sequences from the three libraries. The average sequence was determined 6.8 times. Even with such sequencing redundancy, other methods, such as combinatorial **PCR** or primer walking, were needed to fill in gaps in the sequences (6).

A technique termed “ordered shotgun sequencing” has been described for yeast artificial chromosomes (YACs) that reduces the redundancy of shotgun sequencing. A YAC is subcloned into plasmids, and the ends of the plasmid inserts are sequenced and analyzed for overlaps. Plasmid inserts with minimal overlaps are sequenced fully in the first round and then the next set of plasmids to be sequenced is also chosen for minimal overlap with the ends of the previously sequenced plasmids (7). A similar method of ordered shotgun sequencing has been proposed as the strategy for sequencing the human genome. Human bacterial artificial chromosome libraries (BACs) would be generated, BAC ends would be sequenced, and BACs with minimal overlap would be sequenced fully by shotgun sequencing (8).

### 3. Other “shotgun experiments”

As mentioned, use of the word “shotgun” to describe an experiment implies that a library is randomly screened to obtain a gene or genes of interest or that random library clones are analyzed (eg, by sequencing) to provide a representative survey of the entire library. Such general usage of the word “shotgun” results in a wide variety of different “shotgun experiments” in the literature. In one unusual shotgun experiment, “shotgun **antisense**” was used to create mutants and identify genes in *Dictyostelium* by transforming random antisense cDNA clones into *Dictyostelium*. Clones conferring interesting **phenotypes** were selected for further analysis (9). In another example, a novel V-ATPase was isolated by “shotgun screening” an [expression library](#) with an antiserum against the V-ATPase (10). In a different type of experiment, “shotgun analysis” of cDNA populations was used to identify tissue-specific genes (11).

#### Bibliography

1. L. Axelsson, A. Holck, S. E. Birkeland, T. Aukrust, and H. Blom (1993) *Appl. Environ. Microbiol.* **59**, 2868–2875.
2. W. J. Meijer, G. Venema, and S. Bron (1995) *Nucleic Acids Res.* **23**, 612–619.
3. R. N. Waterhouse, D. J. Silcock, H. L. White, H. K. Buhariwalla, and L. A. Glover (1993) *Mol. Ecol.* **2**, 285–293.
4. M. G. Craven, D. J. Henner, D. Alessi, A. T. Schauer, K. A. Ost, M. P. Deutscher, and D. I. Friedman (1992) *J. Bacteriol.* **174**, 4727–4735.
5. J. C. Roach (1995) *Genome Res.* **5**, 464–473.
6. H.-P. Klenk et al. (1997) *Nature* **390**, 364–370.
7. E. Y. Chen, D. Schlessinger, and J. Kere (1993) *Genomics* **17**, 651–656.
8. J. C. Venter, H. O. Smith, and L. Hood (1996) *Nature* **381**, 364–366.
9. T. P. Spann, D. A. Brock, D. F. Lindsey, S. A. Wood, and R. H. Gomer (1996) *Proc. Natl. Acad. Sci. USA* **93**, 5003–5007.
10. R. Graf, A. Lepier, W. R. Harvey, and H. Wiczorek (1994) *J. Biol. Chem.* **269**, 3767–3744.
11. K. Kato (1992) *Trends Neurosci.* **15**, 319–323.

#### Suggestion for Further Reading

12. A. Favello, L. Hillier, and R. K. Wilson (1995) Genomic DNA sequencing methods. *Methods Cell Biol.* **48**, 551–569. Discussion of shotgun sequencing strategy.

#### Sickle Cell Disease

Sickle cell disease results from homozygosity for a glutamic acid to valine mutation in the b-globin polypeptide chain of hemoglobin (Hb), producing sickle cell Hb (Hb S). Red cells with only Hb S form sickle shapes under low oxygen pressure. These can block circulation in capillaries and also cause problems in the spleen. There is no very effective treatment. Most persons who are homozygous for Hb S die young. The frequency of Hb S is high in African populations and in some groups in India and Arab countries. Heterozygotes appear to have an advantage with respect to falciparum malaria, thus countering the low survival of homozygotes. The Hb S mutation has arisen independently five times.

## 1. Clinical Picture and Genetics

Sickle cell disease (sickle cell anemia) was first described in 1910 by J.B. Herrick. It is the major inherited disorder in most populations of African descent. It also occurs in the Mediterranean region and in some Arabic and Indian groups. Its most characteristic feature is the distorted shape of deoxygenated erythrocytes. They acquire large projections, suggesting the shape of a sickle to early observers. These sickled cells can clog capillaries, leading to even lower oxygenation and complete blocks in circulation. Both oxygenated and deoxygenated red cells in a person with sickle cell disease have increased adherence to vascular endothelium, a property that increases the likelihood of exposure to low oxygen pressure in capillaries.

The blocked capillaries, which can occur anywhere in the body, produce a variety of consequences, depending on the organs involved. Heart damage, brain damage, and damage to the intestinal tract are common occurrences. Renal damage results in inability to concentrate urine. The interruption of blood supply can be very painful, and sickle cell disease is characterized by painful crises as one or another tissue is affected. In many young children, the collection of sickled cells in the spleen causes loss of splenic function and leads to septicemia, a common cause of death. The anemia associated with loss of red cells and decreased production of hemoglobin causes increase in bone marrow and results in bone deformity. A "tower skull" is often seen in patients. Eventually, most patients succumb, usually in the teenage years.

The genetic basis of sickle cell anemia was established in 1949 James V. Neel. Familial clustering had been noted earlier, but the blood of many persons in these same families would sickle without other evidence of disease. Eventually it was recognized that persons whose red cells sickle in vivo are homozygous for a gene that is benign in the heterozygous state. Red cells of heterozygotes sickle in the more stringent anaerobic conditions used in the laboratory but not in vivo. If the allelic variations at this genetic locus are represented by *A* and *S*, the three possible genotypes are *AA* (homozygous wildtype), *AS* (heterozygous carrier), and *SS* (homozygous affected). Heterozygotes are often described as having *sickle cell trait*, as opposed to the sickle cell disease/anemia of homozygotes. There are no well-documented problems associated with heterozygosity for Hb S.

## 2. Molecular Genetics

In 1949, L. Pauling and coworkers compared the electrophoretic mobility of hemoglobin (Hb) from persons with and without sickle cell disease and found that the sickle cell Hb migrates more slowly toward the positive electrode under mild alkaline conditions, indicating a more positive net charge. Hb from heterozygotes showed two populations of molecules, one of which matched normal Hb (Hb A) and the other of which matched sickle cell Hb (Hb S). Thus there is a simple correspondence between the three genotypes and the three electrophoretic patterns.

Vertebrate Hb's are tetrameric, consisting of two each of two related but distinct polypeptide chains. In normal adult human Hb, the chains are labeled *a* and *b*, and the molecular formula is written  $a_2b_2$ . The two types of chain are products of separate genes, the *a*-globin gene (*HBA*) being on human chromosome 16 at 16pter-p13.3 and the *b*-globin gene (*HBB*) on chromosome 11 at 11p15.5. In 1956, Vernon Ingram demonstrated that the only difference between Hb's A and S is replacement of one glutamic acid by valine. Subsequently it was shown that the replacement is in amino acid residue six of the *b* chain. The N-terminal sequence for the *b* chain of Hb A is Val-His-Leu-Thr-Pro-Glu-Glu-Lys-, and that for Hb S is Val-His-Leu-Thr-Pro-Val-Glu-Lys. The Hb S allele can therefore be designated *HBB*, *GLU6VAL*.

These N-terminal segment of the *b* chain that includes this amino acid substitution lies on the exterior of the folded globin. The replacement of glutamic acid by valine creates a new nonpolar site. This allows deoxygenated Hb S in high concentrations to aggregate into fiberlike structures that form

a gel and that align to distort the cell shape and make the cells subject to rupture. In heterozygotes, in which a majority of the Hb is Hb A, sickling does not occur under physiological conditions but can be induced if the oxygen pressure is lowered sufficiently. Even in persons who are homozygous for Hb S, a small increase in fetal hemoglobin ( $\alpha_2\gamma_2$ ) can interfere with sickling and ameliorate the symptoms of sickle cell disease.

### 3. Population Genetics

In view of the fact that most persons with sickle cell disease do not survive to reproduce, one would expect the frequency of the Hb S allele to be very low. Yet, in West Africa, the Hb S allele frequency can be as high as 15%. The expected frequency of homozygotes is 2.25% and of heterozygotes is 25.5%. The proportion of Hb S alleles in homozygotes and therefore subject to strong negative selection would be 15%. The frequency of this allele in the population should decrease rapidly over time. However, it does not.

The explanation appears to be the selective advantage of heterozygotes because of their greater resistance to falciparum malaria as compared to the normal Hb AA population. Whatever the initial allele frequencies, eventually the number of Hb S alleles lost through sickle cell disease will equal the excess number of Hb A alleles lost because of malaria. This *balanced polymorphism* involving falciparum malaria and sickle cell disease in Africa also occurs with falciparum malaria and other hemoglobin disorders, including Hb C and  $\beta$ -thalassemia in West Africa and Hb E, Hb Constant Spring, and  $\alpha$ -thalassemia in Southeast Asia. Approximately 8% of African-Americans are heterozygous for Hb S and some 6 per thousand are affected with sickle-cell disease. Because falciparum malaria does not occur in the U.S., the frequency of the Hb S allele should diminish, although very slowly.

### 4. Origins of the Hb S Mutations

Did the single nucleotide substitution that constitutes the Hb S mutation occur once and spread because of selection, or did it occur multiple times? This question can now be answered because of the availability of a number of polymorphic DNA markers that lie in the vicinity of the *HBB* locus. Because these markers are close to each other, recombination during meiosis is very rare, and a particular combination of "alleles" at these marker loci is transmitted as a unit, called a *haplotype*. There are several hundred possible combinations of marker variations, but only a portion have been observed worldwide. Any one population will have a number of haplotypes, some of which will be shared with other populations, depending on the prehistoric origins of the populations and the migrations and hybridizations that have occurred.

When a mutation occurs, it necessarily occurs within a specific haplotype combination. The nucleotide substitution in codon 6 of the  $\beta$ -globin gene, if it arose only once, should occur in association with a single haplotype or in related haplotypes that could have arisen by rare recombination events from the original haplotype. On the other hand, if the mutation occurred more than once, it should be found in distinctive haplotype backgrounds.

Analysis of Hb S/haplotype combinations indicates that the mutation must have occurred at least five times. Four are found predominantly in Africa and are designated Bantu, Benin, Senegal, and Cameroon. A fifth is found in the populations of India and Saudi Arabia in which the sickle cell gene occurs. The Hb S allele in Sicily and other Mediterranean areas occurs with the Benin haplotype, which is otherwise very rare there. This suggests an African origin for the Mediterranean alleles.

There is some difference in severity of sickle cell disease also associated with haplotypes. Persons homozygous for the Arabian haplotype are least severely affected and those with the Bantu haplotype are the most severely affected. Because the Hb S mutation does not differ among the haplotypes, it is likely that regulatory elements in the  $\beta$ -globin complex vary, perhaps resulting in

more or less fetal hemoglobin that interferes with sickling.

## 5. Diagnosis and Treatment

The diagnosis of homozygosity for Hb *S* is readily accomplished by gel electrophoresis of red cell lysates from blood of newborns or adults. Examination of Hb in fetal red cells cannot be used routinely for diagnosis because of the low production of  $\beta$ -globin during this period. Direct examination of DNA has made prenatal diagnosis possible, however. Most earlier procedures were based on the fact that the DNA sequence for codons 5 to 7 is 5'-CCTGAGGAG-3' in the case of Hb *A* and 5'-CCTGTGGAG-3' for Hb *S*. The restriction enzyme MstII cleaves at the sequence CCTNAGG, which is present in Hb *A* but not in Hb *S*. Thus the restriction fragment patterns will differ for the two alleles. The difference can be observed using a Southern blot with an appropriate  $\beta$ -globin probe. Several variations on polymerase chain reaction (PCR) amplification of the altered DNA segment have also been developed. These have the advantage of requiring minute quantities of DNA, such as single cells from an in vitro-fertilized 8-cell embryo or from rare fetal cells in the maternal circulation.

Effective treatment of sickle cell disease has yet to be developed. Current approaches have been directed toward increasing the level of fetal Hb, which interferes with sickling. Hydroxyurea increases fetal Hb and has some beneficial effect on sickling, but the levels required and the uncertain long term effects have been problems. Nevertheless, it is the treatment of choice at present. Bone marrow replacement should be effective but generates its own major problems.

### Additional Reading

Bunn H.F. and Forget B.G., *Hemoglobin: Molecular, Genetic and Clinical Aspects*. Saunders, Philadelphia, 1986. pp. 690

Steinberg M.H., Management of sickle cell disease, *N. Eng. J. Med.* **340**, 1021–1030 (1999).

*This entry in On-line Mendelian Inheritance in Man provides summaries of studies of sickle cell disease and links to important sites and literature . (site currently unavailable).*

*This entry in On-line Mendelian Inheritance in Man provides information and links for the HBB (-globin) locus as well as a list of alleles. Allelic variant.0243 refers specifically to Hb S. (site currently unavailable).*

## Side Chain

Side chain is the general term to describe the variable part of a [polymer](#), the functional group(s) attached to the constant or regularly repeating [backbone](#). Different kinds of polymers have different kinds of side chains. For homopolymers, the side chains are all chemically equivalent (for example, [polyacrylamide](#) has amide side chains). Copolymers have two or more types of side chains and can be very complex. Biological examples of copolymers include the **nucleic acids** and [proteins](#). In the nucleic acid **DNA**, which has a sugar phosphate backbone, there are four possible side chains corresponding to the purine and pyrimidine bases **cytosine**, **adenine**, **thymidine**, and **guanine**. Proteins are even more complex, with 20 possible [amino acids](#), each with a different side chain on the [polypeptide chain](#) backbone. The 20 different protein side chains are chemically very diverse, including [polar](#), [nonpolar](#), aromatic, sulfur-containing, positively charged, and negatively charged functional groups. Further complexity is possible for protein side chains through [post-translational modifications](#), such as [phosphorylation](#), [N-glycosylation](#), [O-glycosylation](#), [sulfation](#), and [disulfide](#)

[bond](#) formation, that modify the side-chain chemistry. The order or sequence of amino acids, and therefore of side chains, in a polypeptide chain influences both the structure and function of the protein. Furthermore, the large variety of available amino acids leads to an immense number of possible side-chain combinations. This diversity in amino acid building blocks is the basis for the enormous range of [protein structures](#) and functions.

[See also [Backbone](#) and [Polymer](#).]

#### Suggestions for Further Reading

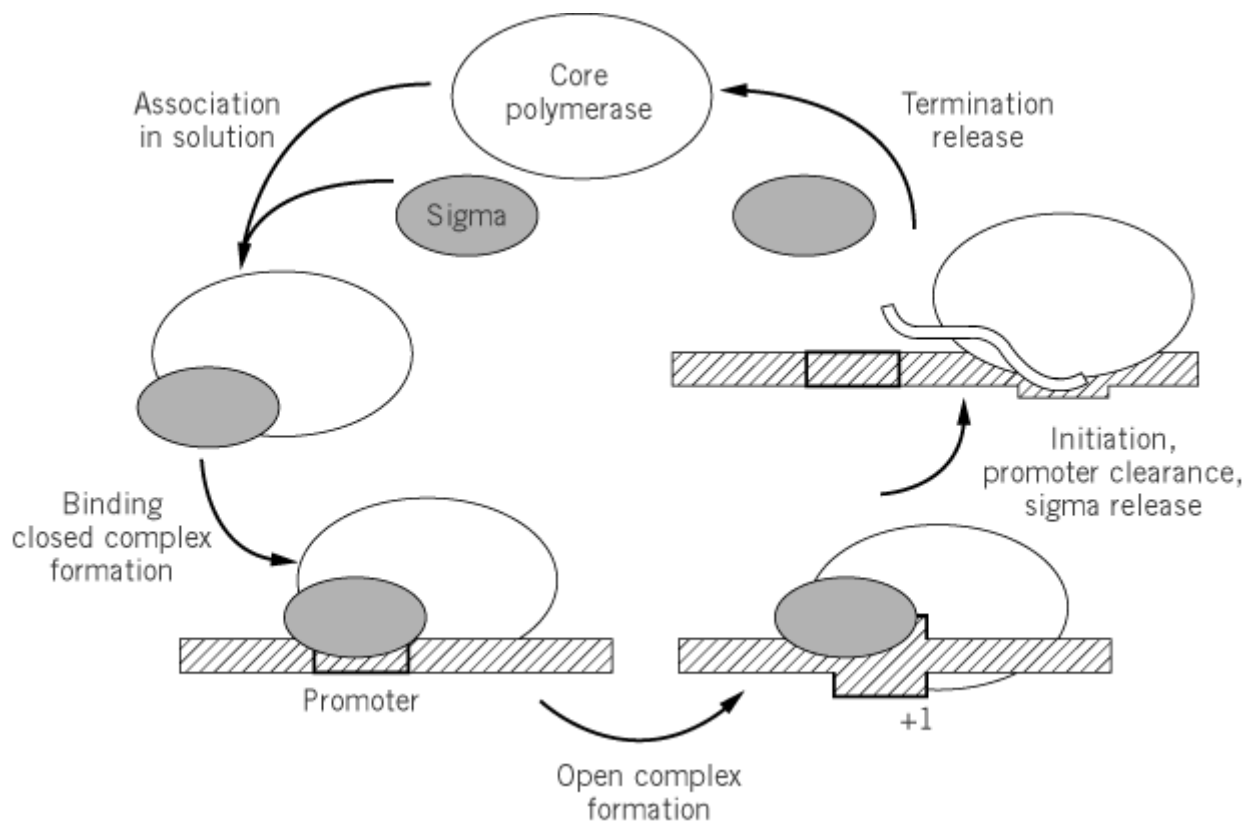
P. Munk (1989) *Introduction to Macromolecular Science*, Wiley-Interscience, New York.

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

## Sigma Factors

Sigma factors are subunits of the DNA-dependent **RNA polymerase** (RNAP) holoenzyme from eubacteria that are required for initiation of [transcription](#) at specific **promoters** ([1](#), [2](#)). Sigma factors function in a number of roles during transcription initiation: (i) They bind to the core RNAP; (ii) inhibit nonspecific transcription; (iii) recognize specific DNA sequences at the promoter; (iv) contribute to promoter melting; (v) interact with certain transcriptional activators; and (vi) play a critical role in promoter clearance (Fig. [1](#)). Sigma factors therefore bind to the nonspecific core RNAP and convert the enzyme into a form that is specific for initiation. Consistent with s; factors playing a role only during initiation, the factor is released from RNAP during the transition from initiation to elongation.

**Figure 1.** The cycle of interaction between s factor and core RNA polymerase showing the steps of transcription initiation.



Sigma factors were originally identified as a dissociable subunit of the *Escherichia coli* RNAP required for promoter recognition (1). The first  $\sigma$  factor to be identified was the major vegetative  $\sigma^{70}$  factor required for recognition of most *E. coli* promoters. Transcriptional studies in *Bacillus subtilis* and its bacteriophage led to the identification of multiple alternative  $\sigma$  factors that alter promoter recognition, thereby being capable of changing the patterns of transcription. It is now clear that all eubacteria have multiple sigma factors that control genes for specific functions. These functions include: **sporulation**; response to **heat shock**; entry into and maintenance of stationary phase; expression of **flagella** genes; and control of nitrogen metabolism. Some  $\sigma$  factors are specific for classes of bacteria, including the factors involved in sporulation and flagella synthesis, whereas other  $\sigma$  factors, such as those that control nitrogen metabolism, heat shock, and the stationary phase response, are found throughout the bacterial world.

## 1. Overview of Sigma Function

As shown in Figure 1,  $\sigma$  factor-dependent initiation of transcription can be described as a cyclic process. First, the core RNAP associates with a  $\sigma$  factor in the absence of DNA. This association alters the conformation of both the RNAP and the  $\sigma$  factor, allowing the holoenzyme to recognize and bind specific promoter DNA sequences. Eubacterial promoters are comprised of sequences upstream of the start site of transcription that are recognized by  $\sigma$  factors (see below for specific examples). The RNAP initially binds to these recognition sequences, to form a closed complex with the double-stranded DNA of the promoter. After initial binding, the DNA strands are separated, forming an open promoter complex, which is competent for initiation of transcription. The early stages of transcription are only slightly processive, and the enzyme dissociates readily. Many short RNA transcripts are produced while the RNAP maintains contacts with the promoter. When the nascent transcript reaches a critical length, the RNAP undergoes an important transformation, called *promoter clearance*. This process is marked by dissociation of the  $\sigma$  factor from the transcription complex and the concomitant release of contacts with the promoter. At this point, the core RNAP enters an elongation mode of transcription, while the released  $\sigma$  factor is free to interact with a new

core RNAP. After termination of transcription, the core RNAP is released from the DNA and can reassociate with a sigma factor to initiate a new round of transcription.

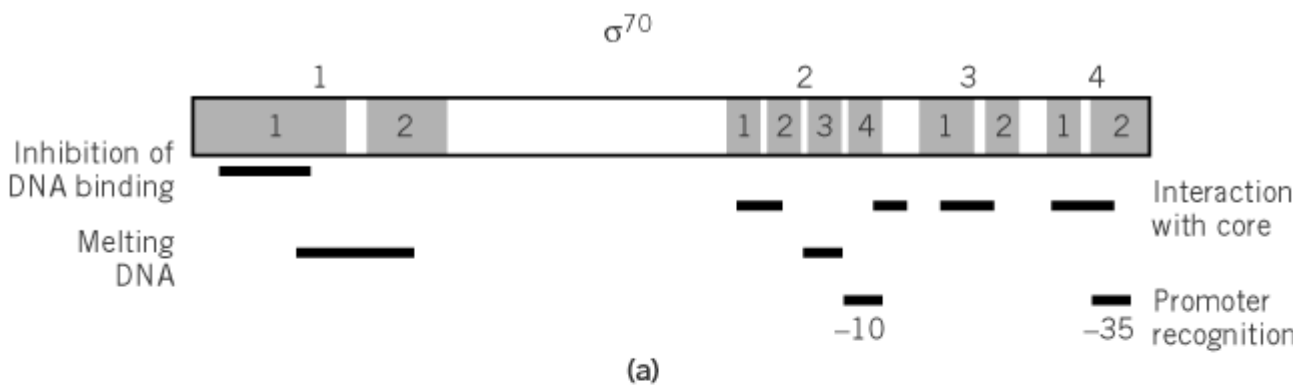
## 2. Classes of Sigma Factors

### 2.1. The $\sigma^{70}$ Family

Most  $\sigma$  factors are similar at the amino acid sequence level to the *E. coli*  $\sigma^{70}$  factor identified originally and constitute a large family of related proteins. The  $\sigma^{70}$  family has been divided into three groups (1). Group 1 includes the primary  $\sigma$  factors. These factors direct the transcription of most of the genes required for exponential growth, and they are therefore essential for cell survival. Group 2  $\sigma$  factors are not essential for growth, but are quite similar at the amino acid sequence level to Group 1  $\sigma$  factors. This group includes the  $\sigma^{38}$  ( $\sigma^S$ ) stationary phase  $\sigma$  factor found in *E. coli* and many other bacteria. Group 3 includes the alternative  $\sigma$  factors. They are required for transcription of certain regulons and have diverged considerably (usually with less than 30% sequence identity) from the Group 1  $\sigma$  factors. Group 3  $\sigma$  factors can be further divided into five closely related clusters. Three of the five clusters include sporulation-specific  $\sigma$  factors, while the remaining two clusters consist of flagellar and heat-shock  $\sigma$  factors. Alignments of the  $\sigma^{70}$  family have identified four regions of high conservation (1-4), see Fig. 2A). These regions have been further subdivided as functional information has been obtained. The primary or Group 1  $\sigma$  factors contain all four of these regions, whereas alternative  $\sigma$  factors may contain only a subset. Regions 2 and 4, which are implicated in DNA binding, are conserved in many of the  $\sigma$  factors, whereas regions 1 and 3 are often absent from the alternative and phage-encoded  $\sigma$  factors. The functions of the conserved regions in the  $\sigma^{70}$  class of factors are described below.

**Figure 2.** Features of selected  $\sigma$  factors. (a) *Escherichia coli*  $\sigma^{70}$ . Conserved regions (numbered 1–4) are labeled above the bar depicting the primary structure of  $\sigma^{70}$ . Conserved subregions are shown as shaded boxes. Regions of known function are shown as lines below the  $\sigma^{70}$  sequence. The data used to derive this figure came from references 1, 2, 13, 20, 22, 23, and 27. (b) Summary of selected sigma factors from *E. coli*, *B. subtilis*, phage T4, and yeast mitochondria (24, 36, 43). Consensus promoter recognition sequences of the  $\sigma$  factors are shown. Nucleotide spacing between the recognition elements is indicated where applicable.





|                                | Factor (kDA)    | Gene           | Consensus promoter |                | Function         |
|--------------------------------|-----------------|----------------|--------------------|----------------|------------------|
| <i>E. coli</i> vegetative      | $\sigma^{70}$   | <i>rpoD</i>    | TTGACA             | 16-18 TATAAT   | General nitrogen |
|                                | $\sigma^{54}$   | <i>rpoN</i>    | CTGGNA             | 6 TTGCA        | Heat shock       |
|                                | $\sigma^{32}$   | <i>rpoH</i>    | CCCTTGAA           | 13-15 CCCGATNT | Flagellar        |
|                                | $\sigma^{28}$   | <i>fliA</i>    | CTAAA              | 15 GCCGATAA    |                  |
| <i>B. subtilis</i> sporulation | $\sigma^H(30)$  | <i>spoOH</i>   | RWAGGAXXT          | 14 HGAAT       | Stationary       |
|                                | $\sigma^F(?)$   | <i>spollAC</i> | GCATR              | 15 GGHRARHTX   | Prespore         |
|                                | $\sigma^E(29)$  | <i>spollGB</i> | ZHATAXX            | 14 CATACAHT    | Mother cell      |
| Phage T4                       | gp55            | TATAAATA       |                    |                |                  |
| Yeast mitochondria             | Mtf1p (43) MTF1 | ATATAAGTA      |                    |                |                  |

(b)

## 2.2. The $s^{54}$ Family

A second class of  $s$  factors is similar to *E. coli*  $s^{54}$  (3, 4). These  $s$  factors have very little amino acid sequence similarity to the  $s^{70}$  family.  $s^{54}$  recognizes a unique class of promoter elements characterized by the presence of unique sequences at nucleotides 2 12 and 2 24 relative to the transcription start site (Fig. 2). Unlike the  $s^{70}$  class of factors,  $s^{54}$ -containing holoenzyme is incapable of forming an open complex on its own. Activation of the RNAP and open complex formation requires an activator protein bound upstream of the promoter. When the activator protein is **phosphorylated** in response to an environmental signal, it triggers open complex formation, in a process that requires ATP hydrolysis. After promoter clearance,  $s^{54}$  remains associated with the promoter DNA (5). Unlike  $s^{70}$ ,  $s^{54}$  can bind some promoters in the absence of core RNAP. The  $s^{54}$  family of  $s$  factors therefore undergoes a unique cycle of transcription distinct from that of the  $s^{70}$  family.

$s^{54}$  can be divided into three regions, designated I to III, respectively: the *N*-terminal region; a poorly conserved acidic central region; and a long *C*-terminal region. Region I contains only the *N*-terminal 25 to 50 amino acid residues. This region is glutamine-rich and is required for activation. Certain mutations in this region permit transcription in the absence of activator proteins (6). Region II is required for binding to core RNAP (7). And finally, Region III is a promoter-binding domain (8). A **helix–turn–helix** motif in this domain is required for recognition of the 2 12 region of the promoter (9).

### 3. Functional Domains of Sigma Factors

#### 3.1. Interactions with DNA

One of the main functions of  $\sigma$  factors is to bind DNA, directing initiation of transcription to the promoter. For a number of years, it was not certain whether  $\sigma$  factors bound DNA or altered the structure of the core RNAP to allow the enzyme to bind DNA. The discovery of multiple  $\sigma$  factors with unique recognition sites supported the idea that the different  $\sigma$  factors directly recognized the promoter DNA. Current data suggest that the majority, if not all, of contacts with the promoter in its 210 region are made by the  $\sigma$  subunit (10). A critical development for  $\sigma$ -specific DNA binding came from genetic analyses of promoter mutants with altered recognition sites (1, 2). These studies identified compensatory alterations in  $\sigma$  factor genes that suppress promoter sequence alterations, confirming that  $\sigma$  factors make direct contacts with the DNA. Suppressor mutation studies revealed that two regions of  $\sigma$  are critical for binding DNA (see Fig. 2a). Region 4.2 recognizes the 235 element of the  $\sigma^{70}$  promoter. This region is predicted to form a helix–turn–helix motif similar to that found in many DNA-binding proteins. Two residues in the second  $\alpha$ -helix of this motif affect promoter binding in *E. coli*  $\sigma^{70}$  and *B. subtilis*  $\sigma^A$ . Furthermore, mutations that switch the promoter specificities of *B. subtilis*  $\sigma^B$ ,  $\sigma^F$ , and  $\sigma^G$  also map to this region. The second DNA-binding region, 2.4, recognizes the 210 element of the promoter. Genetic studies have identified amino acid residues in this region that make specific contacts with bases in the 210 region of the promoter. Recent structural studies have confirmed earlier speculation that this region forms a DNA-recognition  $\alpha$ -helix (11).

Intact  $\sigma^{70}$  does not bind to promoter DNA, but direct recognition can be observed after removal of an *N*-terminal inhibitory region (12). Subregion 1.1 (see Fig. 2a) prevents the  $\sigma$  factor from binding to DNA in the absence of core RNAP by direct interactions with the *C*-terminal region of  $\sigma$  that includes subregion 4.2 (13). Binding of  $\sigma$  factor to core RNAP must alter the inhibitory effect of this domain. In support of this idea, recent evidence suggests that the core RNAP makes direct contact with region 1.1 (see “Sigma-Core Interactions,” below). Inhibition of DNA binding by the *N*-terminal regions of sigma factors may be a common mechanism of regulation. Although they have no obvious amino acid sequence similarity to region 1.1 of  $\sigma^{70}$ , the *N*-terminal regions of *E. coli*  $\sigma^{32}$  and *B. subtilis*  $\sigma^K$  also alter the ability of these factors to bind DNA (13).

Recent work has revealed that base-specific recognition of the 210 region of the promoter is directed through interactions with the nontemplate strand (10, 14). *E. coli* RNAP holoenzyme can bind single-strand oligonucleotides containing the consensus 210 region of the nontemplate strand, but not of the template-strand. However, whether the holoenzyme first recognizes double-stranded DNA and then later converts the binding to single-strand interaction is not known. Alternatively, single-strand binding may be the only base-specific interaction, after melting of the DNA by holoenzyme. In either case, specific binding to the nontemplate strand would leave the template strand accessible to the catalytic center of the RNAP.

#### 3.2. Sigma-Core Interactions

The complex interactions between the sigma subunit and the core RNAP play a critical role in the regulation of the transcriptional cycle. In the previous section, we reviewed the requirement for core RNAP to unmask the sigma DNA-binding activity. Additional regulatory interactions have been observed (2). For example, binding of  $\sigma^{70}$  decreases the affinity of core RNAP for nonpromoter DNA by at least  $10^4$ -fold. Furthermore, the affinity of  $\sigma$  for core is relatively high (**dissociation constant** of about 2 nM). However, the affinity decreases more than  $10^3$ -fold after the transition to the elongating form of core RNAP. Association with  $\sigma$  factor also greatly alters the structure and conformation of the core RNAP. This conclusion is supported by low-resolution structure studies of the core and holoenzyme (15, 16). Both the core and holoenzyme contain a channel presumed to bind the template DNA, but the channel of the core RNAP is closed. Association with  $\sigma$  factor

appears to open the channel of the holoenzyme, making it accessible to DNA for promoter binding. The transition to the closed form of the core RNAP after release of  $\sigma$  factor is thought to contribute to the processive nature of the elongation complex, serving to “clamp” the RNAP to the DNA.

Despite the importance of the  $\sigma$ -core interaction, little is known about the specific amino acid residues that mediate this association. Recent work has indicated that several of the conserved regions of  $\sigma$  factors are required for this interaction (see Fig. 2a). For a number of years, region 2.1 has been implicated in core binding. Deletion analysis of  $\sigma^{70}$  identified a segment of 30 residues that may contribute to interaction with core RNAP (17). This segment forms an  $\alpha$ -helix (18) and contains most of subregion 2.1. Mutational studies have reinforced the importance of region 2.1 in core interaction and have identified additional regions of  $\sigma$  that are required for interaction with the core RNAP. Many of these studies involve single amino acid substitutions in  $\sigma$  factors that significantly reduce the interaction with the core RNAP. [Site-directed mutagenesis](#) of specific regions of sigma factors has confirmed the importance of amino acid residues in subregion 2.1 and also identified subregion 2.2 as important for core interaction (19-21). Additional core-binding mutants have been identified by the analysis of defective  $\sigma^{32}$  alleles (22). Mutations in subregion 4.2 and C-terminal to 2.4, in a region known as the RpoH-box, were found to be defective for core interaction. Several other mutants in subregions 3.1 and 4.1 were identified that had core-binding defects that were apparent only when in competition with  $\sigma^{70}$  for core binding.

An additional study [footprinting proteins](#) has confirmed that the interface between  $\sigma$  and core RNAP is not simple. Footprints of  $\sigma^{70}$  alone and bound to core RNAP showed that regions 1.1, 2, 3, and 4.2 are all protected by interaction with core (23). While these studies show that multiple regions of the  $\sigma$  subunit contact the core RNAP, there may be some differences between  $\sigma$  factors in their core binding sites, because the  $\sigma^{32}$  RpoH box is not conserved in other  $\sigma$  factors. However, subregion 2 probably plays a pivotal role in core binding, because this region has been implicated in core binding in many  $\sigma$  factors. Interestingly, similarities to  $\sigma^{70}$  region 2 have been identified in the core-binding region of  $\sigma^{54}$ , the otherwise unrelated  $\sigma$  factor. A small region in  $\sigma^{54}$  is nearly identical to residues found in subregion 2.1 of  $\sigma^{70}$ . Mutations in this conserved patch of residues demonstrate that this region is required for the  $\sigma^{54}$ -core interaction (7). This work suggests that the two types of  $\sigma$  factors may have a common core binding site, despite their very different properties and lack of overall amino acid sequence similarity.

#### 4. DNA Melting By Sigma Factors

Early predictions of how  $\sigma$  factor would bind to promoter DNA were based on comparisons to [single-stranded DNA binding proteins](#), which use aromatic amino acid residues to interact with exposed bases and basic residues for charge neutralization and stabilization of the bound complex. Sigma factor subregion 2.3 contains many aromatic and basic amino acids and was therefore predicted to be involved in interactions with single-stranded DNA at the promoter (24). This prediction was consistent with photocrosslinking experiments that showed contact between  $\sigma^{70}$  and the 2 3 position of the nontemplate strand (25). Evidence for single-stranded DNA binding was provided by mutagenesis of specific residues in region 2.3 in a primary  $\sigma$  factor ( $\sigma^A$  of *B. subtilis*) (26). Many of these mutations impair DNA melting by the RNAP holoenzyme. However, the melting defect can be overcome by raising the reaction temperature or by using a **supercoiled** template. Several of the mutant  $\sigma$  factors exhibit a *trans*-dominant lethal phenotype *in vivo* consistent with their defect in DNA melting (26).

A recent investigation of  $\sigma^{70}$  region 1, already implicated in autoinhibition of DNA binding and interactions with core RNAP, has identified an additional possible role in DNA melting (27). Deletion of region 1.1 resulted in reduced rates of open complex formation and formation of the initial transcribing complex. Deletion of both regions 1.1 and 1.2 resulted in a more severe

phenotype. The mutant holoenzyme complex arrests after initial binding to the promoter. These results suggest that region 1, as well as subregion 2.3, plays a critical role in the early stages of open complex formation.

## 5. Sigma Factors of *Escherichia Coli*

In eubacterial cells there are several different sigma factors that bind to the core RNAP to direct transcription of distinct subsets of genes. Therefore, in addition to the intrinsic strength of each promoter DNA sequence, the degree to which each of approximately 4000 genes is transcribed is in part determined by the properties of  $\sigma$  factors. In *E. coli*, seven  $\sigma$  factors ( $\sigma^{70}$ ,  $\sigma^{32}$ ,  $\sigma^{28}$ ,  $\sigma^F$ ,  $\sigma^S$ ,  $\sigma^{54}$ , and FecI) have been identified (28, 29). Each *E. coli* promoter is recognized by one or more of the different  $\sigma$  factors. A summary of some features of several *E. coli*  $\sigma$  factors is presented in Figure 2b. Most genes expressed in exponentially growing cultures are transcribed by the  $\sigma^{70}$  holoenzyme. This  $\sigma$  factor binds a bipartite promoter with recognition sequences centered at nucleotides 210 and 235 relative to the start site of transcription. Promoters recognized by alternative  $\sigma$  factors have distinct recognition sequences, but with similar spacing between the recognition sites.

Genes responsive to heat shock are transcribed by  $\sigma^{32}$  ( $\sigma^H$ ) and  $\sigma^{24}$  ( $\sigma^E$ ) (24, 30).  $\sigma^{32}$  is the main heat-shock  $\sigma$  factor and is expressed at low levels during steady-state growth. Expression and stability of  $\sigma^{32}$  is dramatically increased upon a shift to high temperature. The other heat-shock  $\sigma$  factor,  $\sigma^{24}$ , is responsive to extracytoplasmic stress signals (31).  $\sigma^{54}$  is responsible for transcription of genes in response to a number of signals, including availability of nitrogen (3, 4). As described above,  $\sigma^{54}$  requires input from additional binding proteins to activate transcription in response to extracellular signals.  $\sigma^{38}$  ( $\sigma^S$ ) is required for survival of starvation in stationary phase. Some genes required for stationary phase can only be transcribed by the  $\sigma^{38}$  holoenzyme; other genes that are expressed in exponential growth or stationary phase can be recognized by either  $\sigma^{70}$  or  $\sigma^{38}$  (32). It has been proposed that  $\sigma^{38}$  is also important for expression of genes in stationary phase cells, based in part on the observation that the intracellular level of  $\sigma^{38}$  increases with decreases in cell growth rate (28).  $\sigma^{28}$  is required for synthesis of **flagella** and **chemotaxis** genes.

The most recent sigma factor identified in *E. coli* is FecI, which responds to external stimuli to regulate a small set of genes required for ferric citrate transport (29). Although homologues of these sigma factors are found in most bacteria, including *B. subtilis* as described below, there are some interesting exceptions. For example, the complete sequence of *Treponema pallidum* (33), an obligate human parasite and the causative agent of syphilis, demonstrated that this organism lacks homologues of  $\sigma^{38}$  and  $\sigma^{32}$ . The parasitic lifestyle of this organism apparently does not require response to starvation or heat stress.

## 6. Sigma Factors of *Bacillus Subtilis*

The **gram-positive bacterium** *B. subtilis* has more  $\sigma$  factors than the **gram-negative** *E. coli*, reflecting the central role these factors play in the complex sporulation pathway of *B. subtilis* (24, 34-36). In *B. subtilis*, there are 11 known sigma factors; five are required for sporulation, and the other six are used only during vegetative growth. The main vegetative sigma,  $\sigma^{43}$  ( $\sigma^A$ ), although it is missing 245 amino acid residues between conserved regions 1 and 2 found in  $\sigma^{70}$ , is 50% identical to *E. coli*  $\sigma^{70}$  and recognizes similar promoters. Like *E. coli*, *B. subtilis* contains alternative  $\sigma$  factors that regulate **stress response**, stationary-phase growth, and flagellar biosynthesis (34). In the case of the flagellar genes, the *B. subtilis*  $\sigma$  factor can substitute for *E. coli*  $\sigma^{28}$  to support flagellar expression.

Furthermore, *B. subtilis* SigL is homologous to  $\sigma^{54}$  and similar in function (37). Finally, *B. subtilis* sigX can complement an *E. coli* fecI mutant, although sigX does not appear to regulate ferric citrate utilization in *B. subtilis* (38). Therefore, in at least several cases, the alternative  $\sigma$  factors of *E. coli*

and *B. subtilis* are analogous.

In Figure 2B, the consensus promoter elements are shown for three of the five *B. subtilis*  $\sigma$  factors that are required for sporulation (36).  $\sigma^H$  is the only sporulation  $\sigma$  that is also active in vegetative cells. In addition to regulating genes required for the initiation of sporulation,  $\sigma^H$  regulates stationary-phase genes and genes needed for sporulation competence. The four remaining sporulation  $\sigma$  factors are expressed only during the sporulation program. Two  $\sigma$  factors ( $\sigma^F$  and  $\sigma^G$ ) function within the prespore, while the other two ( $\sigma^E$  and  $\sigma^K$ ) function within the mother cell. In both cell types, one  $\sigma$  factor regulates early spore development, while the other becomes active later during sporulation

## 7. Anti-Sigma Factors and Competition between Sigma Factors

Studies of transcription in *E. coli* show that there are approximately 2000 core RNAP molecules per cell in exponential cultures (39). At any given time, roughly two-thirds of the core RNAPs are engaged in transcriptional elongation, leaving 600 or 700 RNAP molecules free for interaction with a  $\sigma$  factor. On the other hand, there are about 700 molecules of  $\sigma^{70}$  in the cell, in addition to lesser amounts of the alternative  $\sigma$  factors. Therefore, there is likely to be competition between the primary and alternative  $\sigma$  factors for interaction with core RNAP. Direct competition between  $\sigma$  factors for access to RNAP is also involved in the case of **T4 phage** replication, where transcription of late genes is controlled by the phage-encoded  $\sigma$  factor gp55. The fact that gp55 competes very poorly with  $\sigma^{70}$  for core interaction led to a proposal that a  $\sigma$  antagonist might participate in the T4 developmental process (40). Subsequent investigations identified a 10-kDa T4 phage-encoded protein that binds to  $\sigma^{70}$  and inhibits its activity. This protein, named AsiA, belongs to a now fairly large and growing group of proteins known as *anti- $\sigma$  factors*.

The anti- $\sigma$  factors are an unrelated collection of proteins that bind  $\sigma$  factors to regulate their activity, often in unique ways. For example, after the formation of flagella, the anti- $\sigma$  factor FlgM of *S. typhimurium* binds to  $\sigma^F$  and inhibits its further activity (41). While the flagellum is being synthesized, however, the FlgM protein is actively transported through the immature flagellar tubes and out of the cell. This provides a unique mechanism for regulating the complex developmental program of flagellar formation. Anti- $\sigma$  factors have now been identified for a number of  $\sigma$  factors including: inhibitors of *B. subtilis* sporulation factors  $\sigma^F$  and  $\sigma^B$ ; *E. coli* heat shock  $\sigma$  factors  $\sigma^{32}$  and  $\sigma^{24}$ ; and a possible inhibitor of *E. coli*  $\sigma^{70}$  in stationary-phase cultures (2, 39). Gene regulation in eubacterial cells is therefore modulated by both the presence of different types of  $\sigma$  factors and by their interactions with anti- $\sigma$  factors.

## 8. Sigma Factor Relatives in Phage and Eukaryotes

Sigma-like factors have been identified in the bacteriophage T4 and SPO1 (1). These factors interact with the host core RNAP to direct transcription of middle and late genes in the phage genome. The bacteriophage SPO1  $\sigma$  factors, gp34 and gp28, have regions somewhat similar to regions 2, 3, and 4 of the cellular sigma factors. The T4-encoded sigma, gp55, has some weak amino acid sequence similarity to region 2 of the  $\sigma^{70}$  family. It is interesting that, consistent with the lack of a similarity to conserved region 4, which is required for recognition of the 2 35 element of the bacterial promoter, gp55 recognizes promoters with only a 2 10 consensus (see Fig. 2b).

Two factors with weak but significant sequence similarity to  $\sigma$  factors have been identified in eukaryotic cells. One factor, RAP30, is a component of the required RNA polymerase II initiation factor TFIIF. Like  $\sigma$  factors, RAP30 binds stably to its RNAP (see also text below) and promotes binding of the RNAP to the promoter, while reducing the affinity of the enzyme for nonpromoter DNA (42). A second  $\sigma$ -like factor, Mtf1p, is required for initiation of mitochondrial transcription in

yeast (see Fig. 2b) (43). A number of similarities between Mtf1p and  $\sigma$  factors have been noted, including release of the factor shortly after initiation of transcription. However, Mtf1p binds to an RNAP that is homologous to the T7 phage single-subunit RNAP, rather than to a multisubunit RNAP.

Core-binding studies of the phage and eukaryotic factors provide additional evidence for their similarity to  $\sigma$  factors. The core-binding sites of RAP30 and gp55 have been mapped to regions similar to  $\sigma$  regions 2.1 and 2.2, which are required for core binding in  $\sigma$  factors. Remarkably, human RAP30 also has the ability to interact with the *E. coli* core RNAP (44). The interaction of RAP30 with core RNAP can be disrupted by the addition of  $\sigma^{70}$ , suggesting that the two factors bind to the same region of the RNAP. Like  $\sigma$  factors, multiple regions of Mtf1p are required for binding to the core (45). The Mtf1p core-binding regions include regions similar to subregions 2.1, 2.2, and 3 of the eubacterial  $\sigma$  factors, which, as described above, have also been implicated in  $\sigma$ -factor-core RNAP interactions.

## 9. Concluding Remarks

Although  $\sigma$  factors are used throughout the eubacterial world to direct core RNAP to certain promoters, in both archaea and in the eukaryotic nucleus the unrelated **TATA box**-binding protein (TBP) appears to be the critical determinant for identification of promoters. However, recent work has revealed that multiple forms of TBP exist, and additional factors modulate the activity of the major TBP at different promoters in ways that are specific for tissue and developmental stage (46, 47). At some level, therefore, the regulation of gene expression by  $\sigma$  factors described in this brief review is used in all forms of life. In addition, the cyclic interactions of initiation factors with a nonspecific core RNA polymerase are also found in all types of cells (see [RNA Polymerases, DNA-Dependent](#)). The study of  $\sigma$  factors therefore continues to provide important lessons for transcriptional regulation of all genes.

## Bibliography

1. M. Lonetto, M. Gribskov, and C. A. Gross (1992) *J. Bacteriol.* **174**, 3843–3849.
2. J. D. Helmann (1994) In *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.) Raven Press, New York, pp. 1–17.
3. M. J. Merrick (1993) *Mol. Microbiol.* **10**, 903–909.
4. M. R. Atkinson and A. J. Ninfa (1994) In *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.) Raven Press, New York, pp. 323–342.
5. Y. Tintut, J. T. Wang, and J. D. Gralla (1995) *Genes Dev.* **9**, 2305–2313.
6. J. T. Wang, A. Syed, M. Hsieh, and J. D. Gralla (1995) *Science* **270**, 992–994.
7. Y. Tintut and J. D. Gralla (1995) *J. Bacteriol.* **177**, 5818–5825.
8. C. Wong, Y. Tintut, and J. D. Gralla (1994) *J. Mol. Biol.* **236**, 81–90.
9. M. Merrick and S. Chambers (1992) *J. Bacteriol.* **174**, 7221–7226.
10. M. T. Marr and J. W. Roberts (1998) *Science* **276**, 1258–1260.
11. A. Malhotra, E. Severinova, and S. A. Darst (1996) *Cell* **87**, 127–136.
12. A. J. Dombroski et al. (1992) *Cell* **70**, 501–512.
13. A. Dombroski, W. Walter, and C. Gross (1993) *Genes Dev.* **7**, 2446–2455.
14. C. W. Roberts and J. W. Roberts (1996) *Cell* **86**, 495–501.
15. A. Polyakov, E. Severinova, and S. A. Darst (1995) *Cell* **83**, 365–373.
16. S. A. Darst, E. W. Kubalik, and R. D. Kornberg (1989) *Nature* **340**, 730–732.
17. S. A. Lesley and R. R. Burgess (1989) *Biochemistry* **28**, 7728–7734.
18. A. Malhotra, E. Severinova, and S. Darst (1996) *Cell* **87**, 127–136.
19. M. F. Shuler et al. (1995) *J. Bacteriol.* **177**, 3687–3694.

20. D. Joo, N. Ng, and R. Calendar (1997) *Proc. Natl. Acad. Sci. USA* **94**, 4907–4912.
21. J. P. Leonetti, K. Wong, and E. P. Geiduschek (1998) *EMBO J.* **17**, 1467–1475.
22. D. Joo, A. Nolte, R. Claendar, Y. Zhou, and D. Jin (1998) *J. Bacteriol.* **180**, 1095–1102.
23. H. Nagai and N. Shimamoto (1997) *Genes Cells* **2**, 725–734.
24. J. D. Helmann and M. J. Chamberlin (1988) *Annu. Rev. Biochem.* **57**, 839–872.
25. R. B. Simpson (1979) *Cell* **18**, 277–285.
26. Y. L. Juang and J. D. Helmann (1994) *J. Mol. Biol.* **235**, 1470–1488.
27. C. Wilson and A. Dombroski (1997) *J. Mol. Biol.* **267**, 60–74.
28. A. Ishihama (1993) *J. Bacteriol.* **175**, 2483–2489.
29. A. Angerer, S. Enz, M. Ochs, and V. Braun (1995) *Mol. Microbiol.* **18**, 163–174.
30. C. Gross, M. Lonetto, and R. Losick (1992) In *Transcriptional Regulation* (K. Yamamoto and S. McKnight, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 129–176.
31. A. De Las Penas, L. Connolly, and C. A. Gross (1997) *J. Bacteriol.* **179**, 6862–6864.
32. K. Tanaka et al. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 8303.
33. C. M. Fraser, S. J. Norris, G. M. Weinstock, O. White, G. G. Sutton, et al. (1998) *Science* **281**, 375–388.
34. W. G. Haldenwang (1995) *Microbiol. Rev.* **59**, 1–30.
35. R. Losick and P. Stragier (1992) *Nature* **355**, 601–604.
36. J. Errington, A. Feucht, P. J. Lewis, M. Lord, T. Magnin, et al. (1996) *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **351**, 537–542.
37. M. Debarbouille, I. Martin-Verstaete, F. Kunst, and G. Rapoport (1991) *Proc. Natl. Acad. Sci. USA* **88**, 9092–9096.
38. S. Brutsche and V. Braun (1997) *Mol. Gen. Genet.* **256**, 416–425.
39. M. Jishage and A. Ishihama (1998) *Proc. Natl. Acad. Sci. USA* **95**, 4953–4958.
40. S. Malik, K. Zalenskaya, and A. Goldfarb (1987) *Nucleic Acids Res.* **15**, 8521–8530.
41. K. L. Brown and K. T. Hughes (1995) *Mol. Microbiol.* **16**, 397–404.
42. J. W. Conaway and R. C. Conaway (1991) *J. Biol. Chem.* **266**, 17721–17724.
43. S. Jang and J. Jaehning (1994) In *Transcription: Mechanisms and Regulation* (R. C. Conaway and J. W. Conaway, eds.), Raven Press, New York, pp. 171–184.
44. S. McCracken and J. Greenblatt (1991) *Science* **253**, 900–902.
45. P. F. Cliften et al. (1997) *Genes Dev.* **11**, 2897–2909.
46. M. Chang and J. A. Jaehning (1997) *Nucleic Acids Res.* **25**, 4861–4865.
47. S. Buratowski (1997) *Cell* **91**, 13–15.

### Suggestions for Further Reading

48. M. R. Atkinson and A. J. Ninfa (1994) "Mechanism and Regulation of Transcription from Bacterial  $s^{54}$ -Dependent Promoters". In *Transcription: Mechanisms and Regulation*, (R. C. Conaway and J. W. Conaway, eds.), Raven Press, New York, pp. 323–342.
49. J. D. Helmann (1994) "Bacterial Sigma Factors". In *Transcription: Mechanisms and Regulation*, (R. C. Conaway and J. W. Conaway, eds.), Raven Press, New York, pp. 1–17.
50. R. Landick and J. W. Roberts (1996) The shrewd grasp of RNA polymerase. *Science* **273**, 202–203.
51. M. Lonetto, M. Gribskov, and C. A. Gross (1992) The sigma 70 family: sequence conservation and evolutionary relationships. *J. Bacteriol.* **174**, 3843–3849.
52. A. Malhotra, E. Severinova, and S. Darst (1996) Crystal structure of a  $s^{70}$  subunit fragment from *E. coli* RNA polymerase. *Cell* **87**, 127–136.

## Signal Peptidase

Signal peptidases are the [enzymes](#) that cleave [signal peptides](#) from pre-proteins during their translocation across [membranes](#) (see [Protein Targeting, Intracellular](#)). The type I signal peptidases that are present in eubacteria are simple monomeric proteins, but the [signal peptidases](#) of the [endoplasmic reticulum](#) and [mitochondria](#) exist as multisubunit complexes. Whereas detailed analyses of signal peptides have been ongoing for more than 20 years, the structure and function of the individual subunits of these protein complexes has become clearer only through more recent studies.

### 1. Endoplasmic Reticulum Signal Peptidase

Unlike eubacterial representatives of the type I signal peptidase family, purified ER signal peptidase is a multisubunit protein complex, thus leading to its alternative name signal peptidase complex (SPC) (1). The SPC purified from the ER membrane of canine pancreas contains five membrane-bound subunits (2). Naming of these subunits has been based on their size as measured by [gel electrophoresis](#) (see [SDS-PAGE](#)). The subunits are SPC25, SPC22/23, SPC21, SPC18, and SPC12 (Table 1). SPC22/23 migrates on SDS-PAGE gels as two distinct species exhibiting molecular masses of 22 and 23 kDa. This heterogeneity is due to SPC22/23 being the only **glycoprotein** subunit of the complex. SPC18 and SPC21 are highly homologous to each other (Fig. 1) and bear limited homology to eubacterial leader peptidase: The homology is confined to five distinct regions, denoted A through E (Fig. 2). Because two of these five regions contain probable [active-site](#) residues in the eubacterial type I enzymes (indicated by an asterisk), the mammalian SPC may contain two catalytic subunits. The remaining SPC subunits show no apparent [homology](#) to type I signal peptidases.

**Figure 1.** Alignment of the amino acid sequences of SPC21, SPC18, and Sec11p. Sequence identity is indicated by using capital letters.



|        |     |     |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |   |   |   |   |   |   |   |   |    |   |     |     |   |     |     |
|--------|-----|-----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|---|---|---|---|---|---|---|----|---|-----|-----|---|-----|-----|
| SPC21  | mvr | gav | gt | gl | pa | Sg | LD | if | gD | lR | kM | NK | RQ | LY | QV | LN | Fa | MI | 40 |    |    |    |    |    |    |    |   |   |   |   |   |   |   |   |    |   |     |     |   |     |     |
| SPC18  |     |     |    |    |    |    | m  | l  | S- | LD | f  | l  | d  | D  | v  | R  | r  | M  | NK | RQ | LY | QV | LN | Fg | MI | 28 |   |   |   |   |   |   |   |   |    |   |     |     |   |     |     |
| Sec11p |     |     |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    | M  | N  | l  | R  | f  | e  | l  | q  | k | l | L | N | v | c | f | l | 16 |   |     |     |   |     |     |
| SPC21  | V   | S   | S  | A  | L  | M  | I  | W  | K  | G  | L  | i  | V  | l  | T  | G  | S  | E  | S  | P  | I  | V  | V  | L  | S  | G  | S | M | E | P | A | F | H | R | G  | D | L   | L   | F | 80  |     |
| SPC18  | V   | S   | S  | A  | L  | M  | I  | W  | K  | G  | L  | m  | V  | i  | T  | G  | S  | E  | S  | P  | I  | V  | V  | L  | S  | G  | S | M | E | P | A | F | H | R | G  | D | L   | L   | F | 68  |     |
| Sec11p | f   | a   | S  | A  | y  | M  | f  | W  | q  | G  | L  | a  | i  | a  | T  | n  | S  | a  | S  | P  | I  | V  | V  | L  | S  | G  | S | M | E | P | A | F | q | R | G  | D | i   | L   | F | 56  |     |
| SPC21  | L   | T   | n  | f  | R  | E  | D  | P  | I  | R  | a  | G  | E  | I  | V  | V  | F  | k  | V  | E  | G  | R  | d  | I  | P  | I  | V | H | R | V | i | K | v | H | E  | K | d   | N   | G | d   | 120 |
| SPC18  | L   | T   | N  | r  | v  | E  | D  | P  | I  | R  | V  | G  | E  | I  | V  | V  | F  | r  | i  | E  | G  | R  | e  | I  | P  | I  | V | H | R | V | L | K | i | H | E  | K | q   | N   | G | h   | 108 |
| Sec11p | L   | w   | N  | -  | R  | n  | t  | f  | n  | q  | V  | G  | d  | v  | V  | y  | e  | V  | E  | G  | k  | q  | I  | P  | I  | V  | H | R | V | L | r | q | H | n | n  | h | a   | d   | k | 95  |     |
| SPC21  | I   | K   | F  | L  | T  | K  | G  | D  | N  | N  | e  | V  | D  | D  | R  | G  | L  | Y  | K  | e  | G  | Q  | n  | W  | L  | E  | K | - | K | D | V | V | G | R | A  | R | G   | F   | L | P   | 159 |
| SPC18  | I   | K   | F  | L  | T  | K  | G  | D  | N  | N  | A  | V  | D  | D  | R  | G  | L  | Y  | K  | q  | G  | Q  | h  | W  | L  | E  | K | - | K | D | V | V | G | R | A  | R | G   | F   | v | P   | 147 |
| Sec11p | q   | f   | l  | L  | T  | K  | G  | D  | N  | N  | A  | g  | n  | D  | i  | s  | L  | Y  | a  | n  | k  | i  | y  | L  | n  | K  | s | K | e | i | V | G | t | v | k  | G | y   | F   | P | 135 |     |
| SPC21  | Y   | v   | G  | m  | V  | T  | I  | i  | M  | N  | D  | Y  | P  | K  | F  | K  | Y  | A  | L  | L  | a  | v  | m  | G  | a  | y  | V | L | L | k | R | E | S |   |    |   |     | 192 |   |     |     |
| SPC18  | Y   | i   | G  | i  | V  | T  | I  | i  | M  | N  | D  | Y  | P  | K  | F  | K  | Y  | A  | v  | L  | f  | l  | L  | G  | L  | f  | V | L | v | h | R | E |   |   |    |   | 179 |     |   |     |     |
| Sec11p | q   | l   | G  | y  | i  | T  | i  | w  | i  | s  | e  | n  | k  | y  | a  | K  | f  | A  | L  | L  | G  | m  | L  | G  | L  | s  | r | L | L | g | g | E |   |   |    |   | 167 |     |   |     |     |

**Figure 2.** Regions of homology in the type I signal peptidase family. Regions A through E show the strongest homology among proteins in the type I signal peptidase family. Sequence identity is indicated using capital letters. Amino acids thought to make up the catalytic serine/lysine dyad of eubacterial signal peptidases are indicated using an asterisk. The sequences displayed are from *E. coli* (Lep), the canine ER (SPC21 and SPC18), the ER of the yeast *S. cerevisiae* (Sec11p), and the mitochondrial inner membrane of the yeast *S. cerevisiae* (Imp1p and Imp2p).

|          |     |                    |
|----------|-----|--------------------|
| <b>A</b> |     |                    |
| Lep      | 72  | IVliv              |
| SPC21    | 60  | IVVVL              |
| SPC18    | 49  | IVVVL              |
| Sec11p   | 37  | IVVVL              |
| Imp1p    | 22  | flhii              |
| Imp2p    | 23  | vllti              |
| <br>     |     |                    |
| <b>B</b> |     |                    |
|          |     | *                  |
| Lep      | 88  | SG-SMmPTLl         |
| SPC21    | 65  | SG-SMEPafH         |
| SPC18    | 54  | SG-SMEPafH         |
| Sec11p   | 42  | SG-SMEPafq         |
| Imp1p    | 37  | rGeSMlPTLs         |
| Imp2p    | 38  | kGtSMqPTLn         |
| <br>     |     |                    |
| <b>C</b> |     |                    |
| Lep      | 127 | R-GD-Ivvf          |
| SPC21    | 74  | R-GD-lLFL          |
| SPC18    | 63  | R-GD-lLFL          |
| Sec11p   | 51  | R-GD-ILFL          |
| Imp1p    | 65  | kmGDcIvaL          |
| Imp2p    | 71  | lsrDdIilf          |
| <br>     |     |                    |
| <b>D</b> |     |                    |
|          |     | *                  |
| Lep      | 137 | EdPkldyiKRavGLPGDK |
| SPC21    | 99  | EGrdIPIVhRVikvh-eK |
| SPC18    | 88  | EGreIPIVhRVLkih-eK |
| Sec11p   | 75  | EGkqIPIVhRVLrqhnnh |
| Imp1p    | 76  | tdPnhrIcKRVtGmPGDl |
| Imp2p    | 83  | tnPrkvycKRVkGLPfdt |
| <br>     |     |                    |
| <b>E</b> |     |                    |
| Lep      | 272 | GDNrdNSaDSR        |
| SPC21    | 126 | GDN--NeVDdR        |
| SPC18    | 115 | GDN--NaVDdR        |
| Sec11p   | 102 | GDNnag-nDis        |
| Imp1p    | 130 | GDNlshSlDSR        |
| Imp2p    | 123 | GDNyfhSiDSn        |

**Table 1. Homology (% Identity and % Similarity) Between ER Signal Peptidase Subunits from Canine and Yeast Cells**

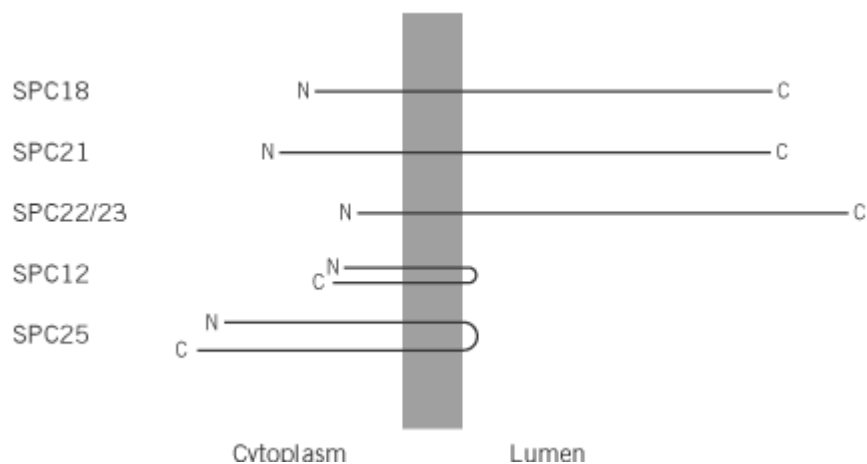
---

| Canine | Yeast | Identity (%) | Similarity (%) |
|--------|-------|--------------|----------------|
|--------|-------|--------------|----------------|

|            |        |    |    |
|------------|--------|----|----|
| SPC18      | Sec11p | 48 | 75 |
| SPC21      | Sec11p | 45 | 75 |
| SPC12      | Spc1p  | 24 | 63 |
| SPC25      | Spc2p  | 20 | 55 |
| SPC22 / 23 | Spc3p  | 22 | 55 |

Membrane topology studies of the SPC subunits reveal two distinct groups: single-spanning proteins orienting their C-termini to the ER lumen and proteins containing two membrane-spanning segments (3, 4). Proteins in the first group, SPC22/23, SPC21, and SPC18, possess much larger luminal domains than cytoplasmic domains (Fig. 3). SPC25 and SPC12, which represent the second topological group, are oriented predominately toward the cytoplasmic side of the membrane. Because the catalytic site of the SPC is located on the luminal side of the membrane, it is thought that SPC25 and SPC12 do not contribute significantly to the catalytic site. This hypothesis is supported by the fact that a functional signal peptidase containing only two subunits has been purified from the ER of hen oviduct (5). These two subunits are homologous to the SPC22/23 and SPC18 subunits of the canine SPC.

**Figure 3.** Membrane topologies of the subunits of ER signal peptidase from canine cells. The deduced topologies of the SPC18, SPC21, SPC22 / 23, SPC12, and SPC25 subunits are represented. The diagram depicts the relative sizes of the cytoplasmic and luminal domains of each subunit but is not drawn to scale.



ors to proteins resident within the ER, transported from the ER to another compartment of the secretory pathway, secreted from the cell, and resident within the **lysosome**. The function of the SPC subunits has been elucidated in the yeast *Saccharomyces cerevisiae*. The first of the signal peptidase genes to be cloned from yeast was *SEC11*, which encodes the 19-kDa subunit (6). This gene has adopted the *SEC* (for secretion; see [Sec Mutants/Proteins, Proteins](#)) nomenclature because *SEC11* was identified in a selection for mutations inhibiting protein secretion. That *sec11* mutations inhibit secretion is due to the fact that defects in cleaving signal peptides can result in abnormal maturation of the cargo and, consequently, retention of immature precursors in the ER lumen (see [Unfolded Protein Response](#)). The genes encoding the other yeast subunits are named by the more conventional SPC nomenclature and are numbered, starting with *SPC1*, according to the order in which the genes have been isolated (see text below). Sequencing of the *SEC11* gene demonstrates that the corresponding protein, Sec11p, is highly related (~45% identity and ~75% similarity) to two subunits of the canine enzyme, SPC18 and SPC21 (Fig. 1) (7). This high degree

of homology is greater than that found for the other SPC subunits (Table 1). In addition, gene-disruption analyses have shown that Sec11p is essential for cell growth, and *sec11* [temperature-sensitive mutations](#) inhibit cleavage of signal peptides in yeast cells incubated at the nonpermissive temperature (6). Consistent with these observations, Sec11p appears to contain the five distinct regions of homology typified by proteins in the type I signal peptidase family (Fig. 2), suggesting that Sec11p contains all or part of the catalytic site.

A high-copy **suppressor** of a *sec11* mutation has been isolated and found to encode the SPC12 homologue Spc1p (8). Overexpression of Spc1p permits growth of *sec11* temperature-sensitive mutant cells at their nonpermissive temperature. Spc1p shares 24% identity and 63% similarity with SPC12 at the primary structure level (Table 1). Sequence analysis of the cloned *SPC1* gene predicts a molecular weight for the encoded protein of 11 kDa and a topology like its mammalian homologue. Despite its presence in the partially purified SPC and homology to mammalian SPC12, Spc1p is not required for cell growth or signal peptidase activity in yeast cells. It is plausible, however, that another polypeptide can replace Spc1p functionally, although analysis of the complete genome of *S. cerevisiae* has revealed no other SPC12 homologue. The only phenotype found to be associated with the *spc1* null mutant is that expression of high levels of a membrane protein fragmented by signal peptidase (see text below) causes a partial accumulation of uncleaved precursors in mutant cells. This finding suggests that Spc1p may be important for efficient signal peptidase activity.

Spc2p (21 kDa) has been identified by comparing open reading frames in the yeast genome to the sequence of canine SPC25 (9). Spc2p shares 20% identity and 55% similarity to SPC25, and it copurifies with Sec11p, supporting the idea that Spc2p is the SPC25 homologue (Table 1). The *SPC2* gene has been disrupted and, as with the *SPC1* gene, found to be nonessential for cell growth and signal peptidase activity under normal laboratory conditions. Furthermore, disruption of both *SPC1* and *SPC2* in the same cell is nonlethal. To date, the only phenotype associated with the *spc2* null mutant is that cell viability and signal peptidase activity are diminished at high temperatures. A specific function for Spc2p is suggested by the fact that the mammalian homologue, SPC25, can be cross-linked to the b-subunit of the translocon (Sec61p complex; see [Endoplasmic Reticulum](#)), suggesting that Spc2p may form a bridge between signal peptidase and the ER translocation machinery and thereby ensure co-translocational signal peptide cleavage (10).

The *SPC3* gene has been isolated using a suppressor approach similar to the analysis that yielded *SPC1* (11). The *SPC3* gene product, Spc3p, is 22% identical and 55% similar to SPC22/23 of the mammalian SPC (Table 1). The calculated molecular mass of the Spc3p polypeptide chain is 21 kDa, but it migrates on [SDS-PAGE](#) gels as a 25-kDa polypeptide. The difference between its calculated and apparent molecular masses is due to the presence of two *N*-linked glycans on the Spc3p polypeptide (see [N-Glycosylation](#)) (12). That Spc3p is a **glycoprotein** is consistent with the fact that it is homologous to SPC22/23. Furthermore, Spc3p is present in the partially purified SPC from yeast cells. Disruption of the *SPC3* gene is lethal to yeast, and *spc3* [temperature-sensitive mutations](#) inhibit cleavage of signal peptides at the nonpermissive temperature (11). The data thus indicate that two subunits, Sec11p and Spc3p, are required for ER signal peptidase activity in yeast cells.

The subunit compositions of the ER signal peptidases that have been purified from mammalian cells and partially purified from yeast cells are summarized in Table 1. The mammalian SPC contains five subunits, but the yeast SPC has four (12, 13). The major difference between these enzymes is the presence of two Sec11p-type subunits in the mammalian SPC (SPC18 and SPC21), accounting for the extra subunit. The presence of two Sec11p homologues appears to be unique to higher eukaryotic systems, because analysis of the complete sequence of the *S. cerevisiae* genome reveals only one Sec11p-type protein. The SPC purified from hen oviduct contains two subunits that correspond to the two essential subunits of the yeast SPC, namely, Sec11p and Spc3p (5). Why the purified avian SPC contains only two subunits is probably explained by differences in the

purification procedures used. Based therefore on genetic data in the yeast system and biochemical evidence from the avian system, Sec11p and Spc3p seem to form a two-subunit core enzyme that is responsible for the signal peptide cleavage reaction. The core enzyme from the ER is the only type I signal peptidase known to contain two polypeptide chains unrelated to each other. As discussed below, the mitochondrial type I signal peptidase also contains two subunits, but these subunits are homologous to each other.

### 1.2. Protein Fragmentation by ER Signal Peptidase

It appears that ER signal peptidase functions in the fragmentation of several abnormal membrane proteins. The internal site(s) cleaved are usually placed after a **transmembrane** segment that is oriented so that the C-terminus of the transmembrane segment is placed on the luminal side of the ER membrane (14, 15). The types of proteins that are fragmented in this manner include a number of chimeric membrane proteins, a mutated membrane protein (the invariant chain of the **major histocompatibility** antigen), and an unassembled membrane protein subunit (the H2 subunit of the human asialoglycoprotein receptor). The presence of a small uncharged amino acid, particularly at the 21 position, seems to be important for cleavage at these internal sites, similar to the “rules” governing cleavage of N-terminal signal peptides. In addition, because the membrane proteins recognized by signal peptidase are abnormal, cleavage sites that are probably cryptic within proteins undergoing normal maturation may be exposed in an abnormal protein.

The subunit requirements for the fragmentation reaction are the same as for signal peptide cleavage. That is, both Sec11p and Spc3p are required for fragmentation, whereas Spc1p and Spc2p are not. The kinetics of the two processes, however, differ dramatically when monitored *in vivo*. Signal peptides can be cleaved while the protein cargo is still being synthesized on a ribosome (see [Translation](#)), meaning that signal peptide cleavage is often a co-translational process. On the other hand, protein fragmentation by signal peptidase occurs several minutes after the fully synthesized polypeptide chain enters the ER lumen (15). As the fragmentation process proceeds rather slowly, other maturation events, such as glycosylation, can precede fragmentation. This raises the possibility that recognition of a fragmentation site may be affected negatively by modifications present at or near a cleavage site or by folding events that bury a potential cleavage site within the structure of the protein.

## 2. Mitochondrial Type I Signal Peptidase

The mitochondrial type I signal peptidase cleaves the second signal located after the matrix targeting signal of proteins disposed to the mitochondrial inner membrane. Mitochondrial type I signal peptidase (named IMP, for inner membrane proteinase) is bound to the inner membrane of mitochondria and is oriented so that its catalytic site is located within the intermembrane space, which lies between the inner and outer mitochondrial membranes. The IMP contains two subunits, Imp1p and Imp2p, both of which are catalytic and bear the five distinct homology regions characteristic of type I signal peptidases (Fig. 2) (16, 17). Imp1p and Imp2p differ, however, in that Imp1p cleaves signal peptides containing the unconventional amino acid asparagine at the 21 position, whereas Imp2p cleaves conventional signal peptides containing a small uncharged residue at this position.

## 3. Catalytic Mechanism of ER and Mitochondrial Type I Signal Peptidase

The eubacterial type I signal peptidases appear to utilize a serine/lysine dyad for catalysis (1). Because serine and lysine residues in Imp1p and Imp2p align with the catalytic serine and lysine residues of their eubacterial counterparts, both subunits of the mitochondrial enzyme may contain a serine/lysine catalytic dyad (Fig. 2). In contrast, despite the functional relationship between ER and eubacterial type I signal peptidases, the Sec11p-type subunits contain a histidine residue that aligns to the catalytic lysine residue of leader peptidase (Fig. 2). From this, the following models are plausible but have not yet been tested: (i) ER signal peptidase utilizes a histidine residue for

catalysis, meaning that ER signal peptidase belongs to a subfamily that is functionally related to the type I enzymes. (ii) ER signal peptidase utilizes a serine/lysine dyad, and therefore current sequence alignments have failed to match catalytic residues in Sec11p to those present in the eubacterial enzymes. (iii) One or more catalytic residues are contained in Spc3p, the other essential subunit of ER signal peptidase.

### Bibliography

1. R. E. Dalbey et al. (1997) The chemistry and enzymology of the type I signal peptidases. *Protein Sci.* **6**, 1129–1138.
2. E. A. Evans et al. (1986) Purification of microsomal signal peptidase as a complex. *Proc. Natl. Acad. Sci. USA* **83**, 581–585.
3. G. S. Shelness et al. (1993) Membrane topology and biogenesis of eukaryotic signal peptidase. *J. Biol. Chem.* **268**, 5201–5208.
4. K. U. Kalies and E. Hartmann (1996) Membrane topology of the 12- and 25-kDa subunits of the mammalian signal peptidase complex. *J. Biol. Chem.* **271**, 3925–3929.
5. R. K. Baker and M. O. Lively (1987) Purification and characterization of chicken oviduct microsomal signal peptidase. *Biochemistry* **26**, 8561–8567.
6. P. C. Bohni et al. (1988) *SEC11* is required for signal peptide processing and yeast cell growth. *J. Cell Biol.* **106**, 1035–1042.
7. G. S. Shelness and G. Blobel (1990) Two subunits of the canine signal peptidase complex are homologous to yeast SEC11 protein. *J. Biol. Chem.* **265**, 9512–9519.
8. H. Fang et al. (1996) The homologue of mammalian SPC12 is important for efficient signal peptidase activity in *Saccharomyces cerevisiae*. *J. Biol. Chem.* **271**, 16460–16465.
9. C. Mullins et al. (1996) Structurally related Spc1p and Spc2p of yeast signal peptidase complex are functionally distinct. *J. Biol. Chem.* **271**, 29094–29099.
10. K.-U. Kalies et al. (1998) The b subunit of the Sec61 complex facilitates cotranslational protein transport and interacts with the signal peptidase during translocation. *J. Cell Biol.* **141**, 887–894.
11. H. Fang et al. (1997) In addition to *SEC11*, a newly identified gene, *SPC3*, is essential for signal peptidase activity in the yeast endoplasmic reticulum. *J. Biol. Chem.* **272**, 13152–13158.
12. H.-A. Meyer and E. Hartmann (1997) The yeast SPC22 / 23 homolog Spc3p is essential for signal peptidase activity. *J. Biol. Chem.* **272**, 13159–13164.
13. J. T. YaDeau et al. (1991) Yeast signal peptidase contains a glycoprotein and the *Sec11* gene product. *Proc. Natl. Acad. Sci. USA* **88**, 517–521.
14. M. Yuk and H. Lodish (1993) Two pathways for degradation of the H2 subunit of the asialoglycoprotein receptor in the endoplasmic reticulum. *J. Cell Biol.* **123**, 1735–1749.
15. C. Mullins et al. (1995) A mutation affecting signal peptidase inhibits degradation of an abnormal membrane protein in *Saccharomyces cerevisiae*. *J. Biol. Chem.* **270**, 17139–17147.
16. M. Behrens et al. (1991) Mitochondrial inner membrane protease 1 of *Saccharomyces cerevisiae* shows sequence similarity to the *Escherichia coli* leader peptidase. *Mol. Gen. Genet.* **228**, 167–176.
17. J. Nunnari et al. (1993) A mitochondrial protease with two catalytic subunits of nonoverlapping specificities. *Science* **262**, 1997–2004.

### Suggestion for Further Reading

18. G. von Heijne (1994) *Signal Peptidases*, R.G. Landes Company, Austin, TX.

## Signal Peptide

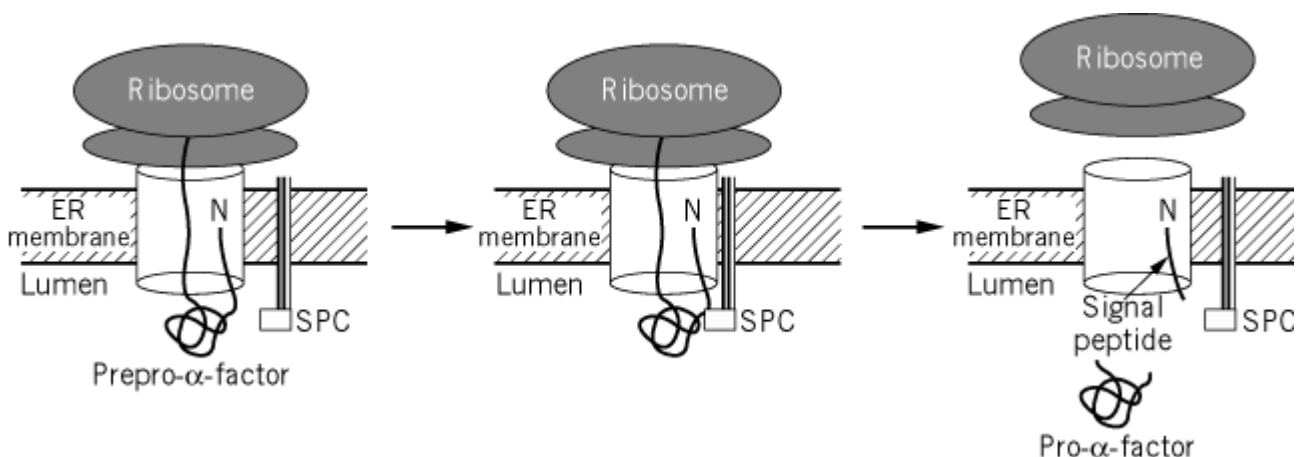
Signal peptides are [cis-acting](#) sequences that effect the transport of [polypeptide chains](#) from the cytoplasm to a membrane-bound cellular compartment (see [Protein Targeting, Intracellular](#)). There are various targeting signals present in eukaryotic cells, and they exhibit much diversity at the [primary structure](#) level, accounting for the specificity required in routing proteins to different organelles. Many signal peptides are cleaved after targeting is achieved, including the type found in precursors targeted to the [endoplasmic reticulum](#) (ER) membrane and the **mitochondrial** inner membrane. This type of signal peptide is cleaved by [enzymes](#) resident within these two organelles that are related functionally to the type I [signal peptidases](#) present in eubacteria. Unlike their eubacterial counterparts, however, the ER and mitochondrial signal peptidases exist as multisubunit complexes. It has also been shown that defects in the signal peptide cleavage event are linked to some human diseases.

### 1. Structures of Signal Peptides

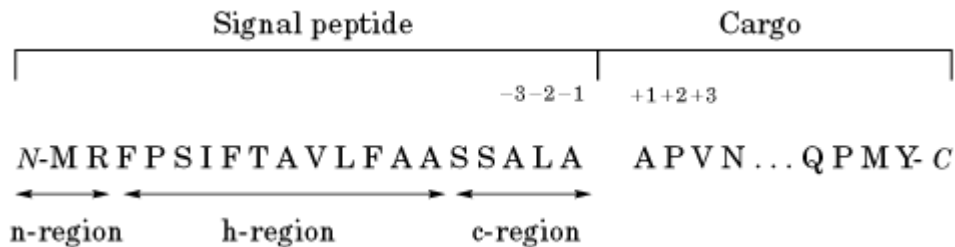
#### 1.1. Endoplasmic Reticulum

Signal peptides of protein precursors targeted to the ER membrane are cleaved by ER signal peptidase. As depicted by the example shown in [Figure 1](#), targeting to and translocation across the ER membrane of the yeast pheromone precursor prepro- $\alpha$ -factor results in cleavage of its signal peptide, thereby liberating the cargo pro- $\alpha$ -factor. The signal peptide of this precursor, like many others targeted to the ER, is placed at the *N*-terminus and has a typical length of 19 amino acid residues, although signals containing 40 or more residues have also been observed ([1](#)). Consistent with this variation in length, there is only modest conservation in the primary structures of various signal peptides, leading to the idea that only common structural features may be important to their function and recognition. Indeed, a general motif is apparent that consists of three regions: (i) a positively charged stretch of amino acids (the n-region), (ii) a central **hydrophobic** core (the h-region), followed by (iii) a polar amino acid stretch that includes the cleavage site (the c-region) ([Fig. 2](#)).

**Figure 1.** The signal peptide cleavage event. The diagram depicts the precursor prepro- $\alpha$ -factor in the process of being cleaved by ER signal peptidase, also termed signal peptidase complex (SPC). Shown here is the cotranslational cleavage of a signal peptide, although some precursors, including prepro- $\alpha$ -factor, can also be cleaved post-translationally.



**Figure 2.** Sequence of a typical signal peptide. The amino acid sequence of the signal peptide from prepro-a-factor and the partial amino acid sequence of pro-a-factor (cargo) are depicted. The positions of the n-region, h-region, and c-region of this typical signal peptide are indicated by arrows. Representative amino acids in the signal peptide (starting with 21) and the cargo (starting with 11) are indicated.



The length of the n-region can vary considerably, from only 1 residue to almost 20. The net charge of the n-region can also vary, but it is often found to be about 12, including the charge of the *N*-terminal [amino group](#). The h-region has an average length of 12 residues and is hydrophobic. [Mutagenesis](#) studies have, in fact, revealed the importance of an uninterrupted stretch of uncharged amino acids within the h-region (2). The c-region exhibits less variation, the average length being 5 to 6 residues (1). At the *C*-terminus of the c-region is the cleavage site, which marks the *C*-terminal end of the signal peptide. Within the c-region, there is a requirement for small uncharged amino acids at the 21 and 23 positions from the cleavage site. Evidence also exists that one or more amino acid residues located within the *N*-terminal part of the cargo is important for determining the site of cleavage by signal peptidase (3).

### 1.2. Anchor Sequences of Membrane Proteins.

Signal peptides that function in targeting many [membrane proteins](#) to the ER membrane are often found internally within the polypeptide chain, and they are not cleaved after targeting in normal situations (4). These internal signal peptides have a dual role, in that they also help to anchor the polypeptide chain within the [lipid](#) bilayer, thus explaining the alias signal/anchor sequences. Signal/anchor sequences span the membrane in an *N*(cytoplasm)–*C*(lumen) orientation. Their structure is like that of cleavable signals, except that they have an extended hydrophobic region containing 18 or more amino acid residues that is necessary to span and become anchored to the lipid bilayer, and they may lack a signal-peptide cleavage site.

### 1.3. Mitochondria Inner Membrane

Signal peptides like those found in proteins targeted to the ER are also present in proteins of the mitochondrial inner membrane. Unlike ER targeting signals, however, these inner-membrane signals are not found at the *N*-terminus of the polypeptide chain. Instead, the *N*-terminal position is occupied by a matrix targeting signal that directs all or part of the polypeptide chain to the mitochondrial matrix, where cleavage by a [metalloproteinase](#) occurs. Upon cleavage of the matrix signal, the inner-membrane signal is then exposed at the *N*-terminus, where it targets the polypeptide chain to the mitochondrial inner membrane. The structure of the inner membrane signal is like its ER counterpart: It contains a hydrophobic h-region flanked on either side by polar (sometimes charged) sequences. Why the mitochondrial inner-membrane signal does not target proteins to the ER may be due to a dominance of the mitochondrial matrix signal over the second signal. One important difference between mitochondrial inner-membrane signals and other signal peptides cleaved by type I signal peptidases is that some inner-membrane signals contain the unconventional amino acid [asparagine](#) at the –1 position.

## 2. Signal Peptides in Disease



Proteolytic removal of *N*-terminal signal peptides is often required for proper maturation of the protein cargo. It therefore seems plausible that mutations inhibiting the cleavage event in humans might be responsible for some inherited diseases. There are at least two diseases that appear to result from the defective cleavage of specific signal peptides. A mutation causing substitution of the –1 alanine with threonine of the pre-pro-vasopressin signal peptide leads to inefficient cleavage by ER signal peptidase and is a possible cause for familial central diabetes insipidus (5). A mutation that results in a glycine to arginine substitution at the –3 position of human coagulation factor X inhibits cleavage by ER signal peptidase and is linked to a severe bleeding diathesis (6). A third inherited disease also appears to result from a signal peptide mutation. A change of cysteine to arginine at the –8 position of the pre-pro-parathyroid hormone alters the hydrophobic nature of the h-region (7). Because both targeting to the ER membrane and cleavage are inhibited by this mutation, it is unclear whether only one or both of these defects is responsible for the disease phenotype.

## Bibliography

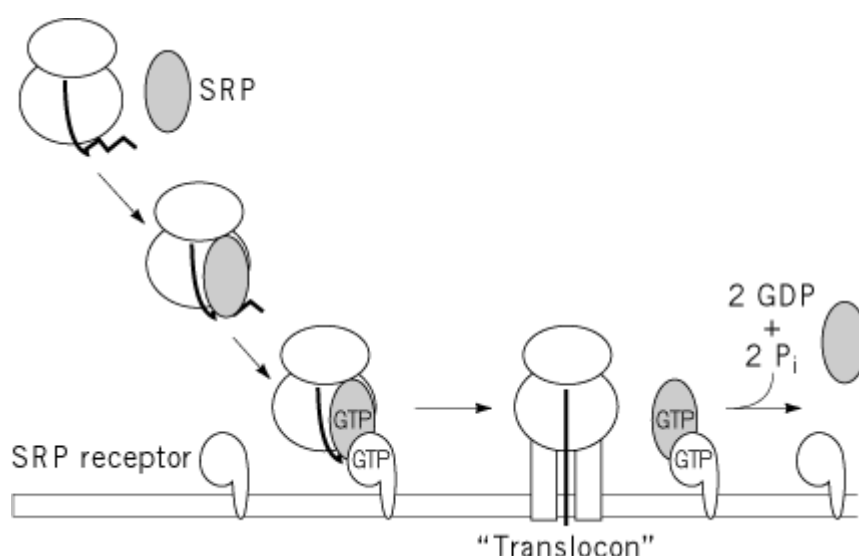
1. G. von Heijne (1985) Signal sequences the limits of variation. *J. Mol. Biol.* **184**, 99–105.
2. D. S. Allison and E. T. Young (1989) Mutations in the signal sequence of prepro-a-factor inhibit both translocation into the endoplasmic reticulum and processing by signal peptidase in yeast cells. *Mol. Cell. Biol.* **9**, 4977–4985.
3. S. F. Nothwehr et al. (1990) Residues flanking the COOH-terminal C-region of a model eukaryotic signal peptide influence the site of its cleavage by signal peptidase and the extent of coupling of its co-translational translocation and proteolytic processing *in vitro*. *J. Biol. Chem.* **265**, 21797–21803.
4. M. Friedlander and G. Blobel (1985) Bovine opsin has more than one signal sequence. *Nature* **318**, 338–343.
5. M. Ito et al. (1993) Possible involvement of inefficient cleavage of preprovasopressin by signal peptidase as a cause for familial central diabetes insipidus. *J. Clin. Invest.* **91**, 2565–2571.
6. M. Racchi et al. (1993) Human coagulation factor X deficiency caused by a mutant signal peptide that blocks cleavage by signal peptidase but not targeting and translocation to the endoplasmic reticulum. *J. Biol. Chem.* **268**, 5735–5740.
7. A. C. Karaplis et al. (1995) Inefficient membrane targeting, translocation, and proteolytic processing by signal peptidase of a mutant preproparathyroid hormone protein. *J. Biol. Chem.* **270**, 1629–1635.

## Signal Recognition Particle

The signal recognition particle (SRP) is a small cytoplasmic [ribonucleoprotein](#) that plays a crucial role in intracellular protein sorting (see [Protein Targeting, Intracellular](#)). In mammalian cells, essentially all proteins that are secreted or that are routed through the secretory and endocytic pathways depend on SRP to initiate their journey. All of these proteins contain either **N-terminal signal peptide sequences** or **transmembrane** regions that earmark them for translocation across, or integration into, the membrane of the [endoplasmic reticulum](#) (ER). SRP recognizes these sequences as soon as they emerge from [ribosomes](#) during [translation](#) (1, 2) (Fig. 1). The binding of SRP causes a transient inhibition of further polypeptide elongation, or “translation arrest” (3). Subsequently, SRP guides or “targets” the ribosome-bound nascent polypeptide chains to transport sites in the ER (4). There, the release of the nascent chains and their insertion into protein transport channels (“translocons”) is facilitated by an interaction between SRP and the heterodimeric SRP receptor (SR)

(5-8). Once the nascent chain is inserted into the translocon, translation resumes, and the growing polypeptide chain is fed gradually into the lumen of the ER. Finally, SRP dissociates from the membrane, and the entire cycle is repeated. Although SRP initiates the protein sorting process by acting as an adaptor between the translation and transport machineries, it plays no direct role in transport per se. An SRP-based co-translational targeting mechanism presumably evolved in part to ensure that translocation begins before [polypeptide chains](#) are large enough to fold tightly and to become translocation-incompetent as a result.

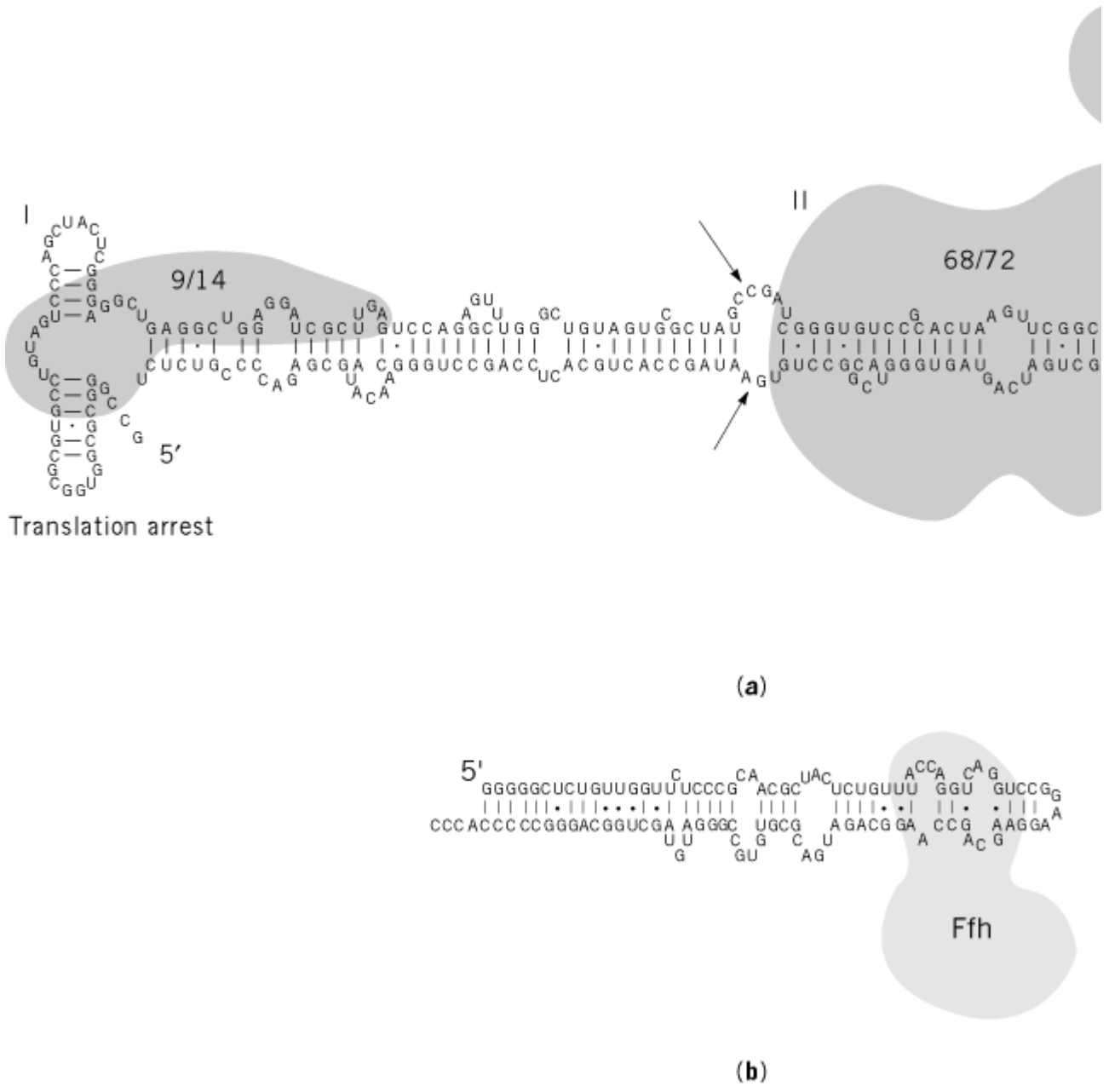
**Figure 1.** The SRP targeting cycle. As described in the text, SRP binds to targeting sequences as they are synthesized and then inhibits further polypeptide elongation (“translation arrest”). SRP-ribosome-nascent chain complexes then travel to the ER, where an interaction between SRP and the SRP receptor catalyzes the release of the nascent chain and its insertion into the translocon. Translation of the polypeptide chain resumes at that point. Finally, SRP dissociates from the membrane, and the cycle is repeated. GTP binds to both SRP54 and the SRP receptor  $\alpha$ -subunit when SRP arrives at the ER membrane and facilitates a regulated release of the nascent chain. GTP hydrolysis at the end of the targeting cycle promotes dissociation of SRP from the ER membrane.



SRP was first identified in mammalian cells and purified on the basis of its ability to support the translocation of secreted proteins into small ER vesicles called “**microsomes.**” The particle is composed of six polypeptide chains of 9, 14, 19, 54, 68, and 72 kDa and a single, highly base-paired RNA of approximately 300 nucleotides (7SL RNA) (9, 10). The [sedimentation coefficient](#) of SRP in sucrose gradients (11 S) is much smaller than would be expected for a 300-kDa spherical particle (see [Hydrodynamic Volume](#)), suggesting that it has a rod-like shape. Examination of SRP by [electron microscopy](#) confirmed that it has an extended structure (11). The shape of SRP suggests that it may need to bridge a considerable distance in order to perform its various functions. Consistent with this notion, the translation arrest and signal recognition activities appear to reside at opposite ends of the particle (see text below and Fig. 2).

**Figure 2.** The structure of SRP. (a) Mammalian SRP has an elongated, rod-like structure. Treatment of the intact particle in cleavage of the 7SL RNA (arrows) and produces two discrete domains (“*Alu* domain” and “S domain”). The *Alu* domain mediates translation arrest. The S domain contains all of the other protein subunits. SRP54 binds to targeting signals and receptor. The region of SRP RNA contacted by each subunit was determined by chemical and enzymatic footprinting and facilitate comparison of SRP RNA molecules from different species, the RNA has also been divided into four structural domains encompassing the 5' and 3' ends, II is the main stem, and III and IV are the stem-loop structures that diverge from domain I. Bacteria and Mycoplasma is a much smaller particle that contains a ~100 nucleotide RNA (4.5 S RNA) and a homologue that contains the protein binding site is very closely related in sequence and structure to the equivalent region in domain

from *Annu. Rev. Cell Biol.* **10**, Copyright 1994.)



Significant progress has been made in obtaining a detailed structural and functional map of SRP. The particle can be divided into two discrete domains by limited **nuclease** digestion (Fig. 2) (12). One domain, which contains both the 5' and 3' ends of the RNA and a heterodimer of the 9- and 14-kDa subunits (SRP9 / 14), is required for SRP to mediate translation arrest. This domain is called the “*Alu* domain” based on the homology of the RNA sequence with the *Alu* family of [repetitive DNA](#) elements (see [Alu Sequences](#)). [X-ray crystallographic](#) analysis has shown that SRP9 and SRP14 are structurally related and that together they form an **RNA-binding** surface comprised of a concave [beta-sheet](#) lined with positively charged amino acids (13). The second domain of SRP is called the “S domain” and contains the remainder of the polypeptide chains and RNA. The 54-kDa subunit (SRP54) binds to signal sequences and to SR (14-16). This protein can be separated into two segments by limited **proteolytic** digestion, one which contains a [GTPase](#) activity (“NG domain”), and one which contains an usually large number of [methionine](#) residues (“M domain”) (17). The

function of the other three subunits is less well understood, and their amino acid sequences have proved uninformative. The 19-kDa subunit (SRP19) is required for efficient binding of SRP54 to 7SL RNA (18). The 68-kDa and 72-kDa subunits bind to the RNA as a heterodimer (SRP68 / 72), although only SRP68 clearly contacts the RNA directly (19). Based on a variety of data, SRP68 / 72 has been proposed to play a role in the interaction of SRP with SR or, alternatively, to influence the affinity of SRP54 for GDP (20, 21). Although the protein subunits bind to 7SL RNA with high affinity (22), SRP can be disassembled easily, and the purified components can be mixed together to reconstitute a completely functional particle (18).

SRP54 binds to signal sequences and transmembrane segments with a high degree of specificity. Many studies indicate that SRP54 recognizes the relatively long stretches of **hydrophobic** amino acids found in these targeting sequences. For example, SRP binds to the signal sequence of preprolactin efficiently, but the binding is abolished when the **leucine** residues are replaced by the hydrophilic amino acid analogue **b-hydroxy leucine** (1). Although signal sequence binding has been mapped to the M domain of SRP54 (17, 23, 24), the mechanism by which signal sequences are identified correctly on the basis of hydrophobicity in spite of considerable variation in amino acid sequence has not been fully elucidated. The crystal structure of the M domain reveals a hydrophobic groove formed by several **alpha-helices** that is likely to be the signal sequence-binding pocket (25). The groove is lined with methionine residues, which differ from other large hydrophobic amino acids in that their side chains can adopt several different energetically favorable conformations. It has been proposed that the “flexible” methionine residues impart a great deal of plasticity to the signal sequence-binding pocket and thereby allow SRP54 to bind to hydrophobic sequences that vary greatly in amino acid composition (26).

The binding and release of signal sequences by SRP is a highly regulated process that involves the participation of several GTPases. In addition to SRP54, both subunits of SR have GTPase activity (16, 27). The observation that the GTPase domains of SRP54 and the SR a-subunit (SRa) are highly **homologous** has led to the proposal that these two domains participate in a homotypic binding reaction (26). Prior to contact between the SR and the complex SRP-ribosome-nascent chain, the GTP binding sites of both SRP54 and SRa appear to be empty (17, 28). A concerted binding of GTP by SRP54 and SRa subsequently stabilizes the interaction between SRP and SR (Fig. 1), and a tight SRP-SR complex is formed in the presence of nonhydrolyzable GTP analogues (28, 29). The binding of GTP is coupled to release of the signal sequence and its insertion into the translocon. Finally, hydrolysis of GTP by both proteins facilitates dissociation of SRP from the membrane (29). Because SRP54 and SRa have a very low affinity for GDP (17, 30), they may return to a nucleotide-free state spontaneously. The SRP54-SRa “concerted switch” presumably serves to prevent premature release of the signal sequence and to ensure proper coupling between protein targeting and protein translocation. The SR b-subunit (SRb) GTPase may also play a role in coordinating the two processes, but this possibility remains to be explored.

## 1. The Evolution of SRP Structure and Function

Analysis of a large number of **genomes** has shown that SRP and SR are widely distributed in all three kingdoms of life. Indeed, it is likely that all cells contain SRP. During the course of **evolution**, however, the particle appears to have undergone several significant size reductions (21). Higher eukaryotes and archaea have SRP RNAs that are similar in size and structure, suggesting that the “primordial” 7SL RNA was approximately 300 nucleotides in length. **Gram-Positive Bacteria**, such as *Bacillus subtilis*, however, have SRP RNAs that lack a segment called domain III (Fig. 2). Many **gram-negative** organisms and **Mycoplasma** have an even smaller SRP RNA of about 100 nucleotides (“4.5 S RNA”) that corresponds to domain IV of 7SL RNA, but the remainder of the RNA seems to be missing altogether (Fig. 2). Despite considerable sequence divergence among the different SRP RNAs, all are predicted to be highly base-paired and all contain a highly conserved sequence motif in domain IV that mediates the attachment of SRP54 (31-33). The only exception is the unusually large 600-nucleotide scR1 RNA of *Saccharomyces cerevisiae* SRP, which cannot be aligned with other SRP RNAs using **secondary structure prediction** programs.

Like SRP RNA, SRP54 is also very highly conserved; homologues have been clearly identified in all organisms that have been examined. The *E. coli* SRP54 homologue (designated “Ffh” for fifty-four homologue) has been shown to form a complex with 4.5 S RNA *in vivo* (34). Purified *S. cerevisiae* SRP has subunits that are homologous to SRP19, SRP68, SRP72, and SRP14, a 7-kDa protein that may correspond to SRP9, and an additional 21-kDa protein not found in higher eukaryotes (35). Unlike SRP54, these proteins are not highly homologous to their mammalian counterparts. Significant evolutionary drift may explain why genes that encode these subunits have not been identified in the genomes of microorganisms in which their existence is predicted by the presence of a large SRP RNA.

SRa is also ubiquitous in nature, but SRb has been identified only in eukaryotes. Like its eukaryotic counterpart, the *E. coli* homologue of SRa (“FtsY”) forms a tight complex with Ffh in the presence of the GMP analogue GMP-PNP (36), suggesting that interaction among components in the SRP targeting pathway is also conserved.

Despite the changes in SRP structure during evolution, recent studies indicate that the targeting function has been highly conserved. Bacterial SRP facilitates the transport of proteins across the cytoplasmic membrane, which is evolutionarily related to the eukaryotic ER and contains a homologous translocation machinery. Although it might seem surprising that an SRP that has lost most of its subunits would perform a similar function, a particle consisting of only mammalian SRP54 and *E. coli* 4.5 S RNA promotes protein translocation into microsomes under certain experimental conditions (37). Yeast and bacteria differ from mammalian cells, however, in that many proteins can use SRP-independent targeting pathways to gain access to the [protein secretion](#) pathway. Depletion of SRP in *S. cerevisiae* has a profound effect on the translocation of some proteins, but only a small effect on others (38). In that organism, SRP-dependence correlates well with the hydrophobicity of the signal sequence (39). In bacteria, even fewer proteins require SRP to traverse the cytoplasmic membrane. Inhibition of SRP function in *E. coli* blocks the insertion of many integral [membrane proteins](#) into the cytoplasmic membrane, but has essentially no effect on protein secretion (40, 41). SRP-independent targeting can occur post-translationally, and its prevalence in rapidly growing organisms suggests that they favor efficient export over the high degree of control afforded by SRP. **Molecular chaperones** such as hsp70 and the *E. coli* SecB protein (see [Sec Mutants/Proteins](#)) play an important role in SRP-independent targeting (42, 43), probably because they help to maintain proteins in a loosely folded, translocation-competent conformation.

## Bibliography

1. P. Walter, I. Ibrahimi, and G. Blobel (1981) *J. Cell Biol.* **91**, 545–550.
2. M. Friedlander and G. Blobel (1985) *Nature* **318**, 338–343.
3. P. Walter and G. Blobel (1981) *J. Cell Biol.* **91**, 557–561.
4. P. Walter and G. Blobel (1981) *J. Cell Biol.* **91**, 551–556.
5. R. Gilmore, G. Blobel, and P. Walter (1982) *J. Cell Biol.* **95**, 463–469.
6. R. Gilmore, P. Walter, and G. Blobel (1982) *J. Cell Biol.* **95**, 470–477.
7. D. I. Meyer, E. Krause, and B. Dobberstein (1982) *Nature* **297**, 647–650.
8. S. Tajima, L. Lauffer, V. L. Rath, and P. Walter (1986) *J. Cell Biol.* **103**, 1167–1178.
9. P. Walter and G. Blobel (1980) *Proc. Natl. Acad. Sci. USA* **77**, 7112–7116.
10. P. Walter and G. Blobel (1982) *Nature* **99**, 691–698.
11. D. W. Andrews, P. Walter, and F. P. Ottensmeyer (1987) *EMBO J.* **6**, 3471–3477.
12. V. Siegel and P. Walter (1986) *Nature* **320**, 81–84.
13. D. E. A. Birse, U. Kapp, K. Strub, S. Cusack, and A. Aberg (1997) *EMBO J.* **16**, 3757–3766.
14. U. C. Krieg, P. Walter, and A. E. Johnson (1986) *Proc. Natl. Acad. USA* **83**, 8604–8608.

15. T. V. Kurzchalia, M. Wiedmann, A. S. Girschovich, E. S. Bochkareva, H. Bielka, and T. A. Rapoport (1986) *Nature* **320**, 634–636.
16. J. D. Miller, H. Wilhelm, L. Gierasch, R. Gilmore, and P. Walter (1993) *Nature* **366**, 351–354.
17. D. Zopf, H. D. Bernstein, A. E. Johnson, and P. Walter (1990) *EMBO J.* **9**, 4511–4517.
18. P. Walter and G. Blobel (1983) *Cell* **34**, 525–533.
19. H. Lütcke, S. Prehn, A. J. Ashford, M. Remus, R. Frank, and B. Dobberstein (1993) *J. Cell Biol.* **121**, 977–985.
20. V. Siegel and P. Walter (1988) *Cell* **52**, 39–49.
21. S. Althoff, D. Selinger, and J. A. Wise (1994) *Nucleic Acids Res.* **22**, 1933–1947.
22. F. Janiak, P. Walter, and A. E. Johnson (1992) *Biochemistry* **31**, 5830–5840.
23. H. Lütcke, S. High, K. Römisch, A. J. Ashford, and B. Dobberstein (1992) *EMBO J.* **11**, 1543–1551.
24. D. Zopf, H. D. Bernstein, and P. Walter (1993) *J. Cell Biol.* **120**, 1113–1121.
25. R. J. Keenan, D. M. Freymann, P. Walter, and R. M. Stroud (1998) *Cell* **94**, 181–191.
26. H. D. Bernstein, M. A. Poritz, K. Strub, P. J. Hoben, S. Brenner, and P. Walter (1989) *Nature* **340**, 482–486.
27. T. Connolly and R. Gilmore (1989) *Cell* **57**, 599–610.
28. P. J. Rapiejko and R. Gilmore (1997) *Cell* **89**, 703–713.
29. T. Connolly, P. J. Rapiejko, and R. Gilmore (1991) *Science* **252**, 1171–1173.
30. J. D. Miller, S. Tajima, L. Lauffer, and P. Walter (1995) *J. Cell Biol.* **128**, 273–282.
31. J. C. R. Struck, H. Y. Toschka, T. Specht, and V. A. Erdmann (1988) *Nucleic Acids Res.* **16**, 7740.
32. M. A. Poritz, K. Strub, and P. Walter (1988) *Cell* **55**, 4–6.
33. H. Wood, J. Luirink, and D. Tollervey (1992) *Nucleic Acids Res.* **20**, 5919–5925.
34. M. A. Poritz, H. D. Bernstein, K. Strub, D. Zopf, H. Wilhelm, and P. Walter (1990) *Science* **250**, 1111–1117.
35. J. D. Brown, B. C. Hann, K. F. Medzihradzky, M. Niwa, A. L. Burlingame, and P. Walter (1994) *EMBO J.* **13**, 4390–4400.
36. J. D. Miller, H. D. Bernstein, and P. Walter (1994) *Nature* **357**, 657–659.
37. S. Hauser, G. Bacher, B. Dobberstein, and H. Lütcke (1995) *EMBO J.* **14**, 5485–5493.
38. B. C. Hann and P. Walter (1991) *Cell* **67**, 131–144.
39. D. T. Ng, J. D. Brown, and P. Walter (1996) *J. Cell Biol.* **134**, 269–278.
40. J. W. de Gier, P. Mansournia, Q. A. Valent, G. J. Phillips, J. Luirink, and G. von Heijne (1996) *FEBS Lett.* **399**, 307–309.
41. N. D. Ulbrandt, J. A. Newitt, and H. D. Bernstein (1997) *Cell* **88**, 187–196.
42. C. A. Kumamoto and J. Beckwith (1985) *J. Bacteriol.* **163**, 267–274.
43. W. J. Chirico, M. G. Waters, and G. Blobel (1988) *Nature* **332**, 805–810.
44. K. Strub, J. Moss, and P. Walter (1991) *Mol. Cell. Biol.* **11**, 3949–3959.
45. V. Siegel and P. Walter (1988) *Proc. Natl. Acad. Sci. USA* **85**, 1801–1805.

### **Suggestions for Further Reading**

46. P. Walter and A. E. Johnson (1994) Signal sequence recognition and protein targeting to the endoplasmic reticulum membrane. *Annu. Rev. Cell Biol.* **10**, 87–119. A comprehensive review of the literature on SRP structure and function that contains an extensive bibliography.
47. J. S. Millman and D. W. Andrews (1997) Switching the model: a concerted mechanism for GTPases in protein targeting. *Cell* **89**, 673–676. A detailed discussion of the GTPase cycles of SRP54 and SR .

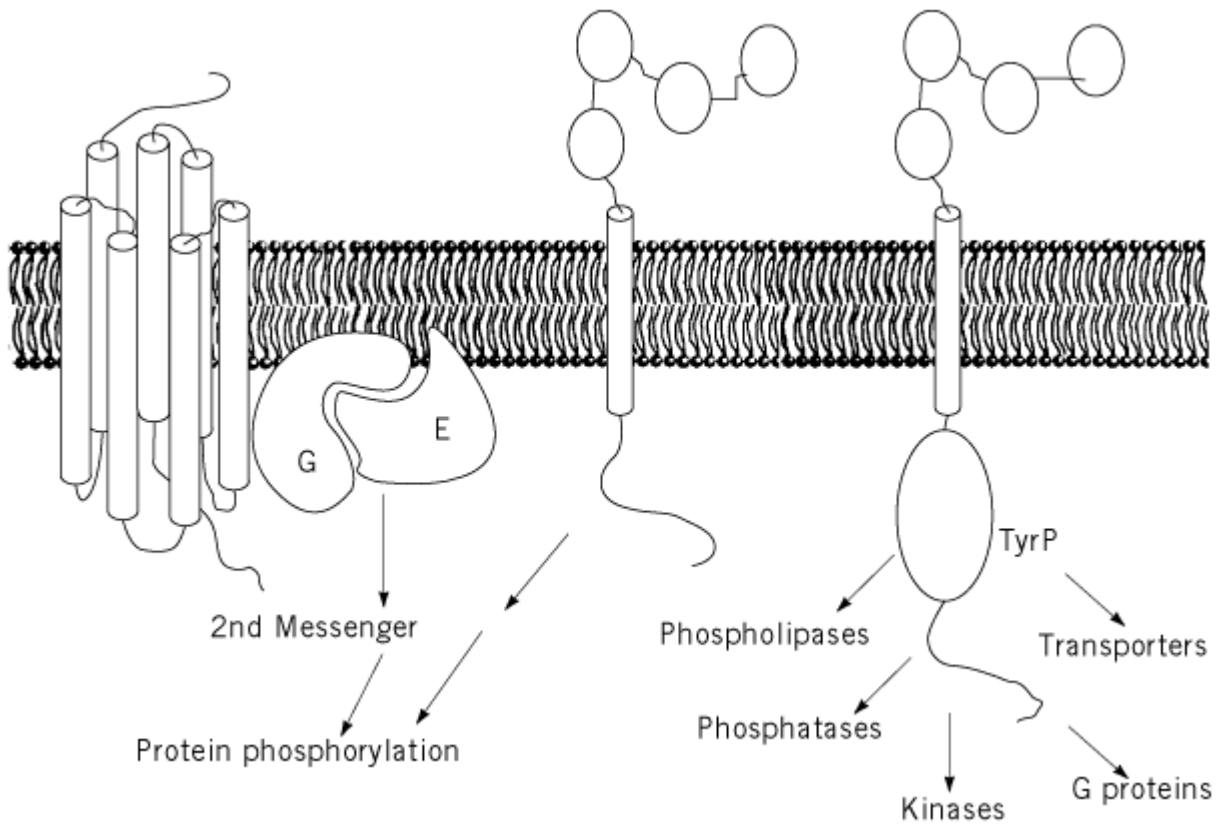
48. J. W. De Gier, Q. A. Valent, G. von Heijne, and J. Luirink (1997) The *E. coli* SRP: preferences of a targeting factor. *FEBS Lett.* **408**, 1–4. A review of studies on the function of bacterial SRP.
49. N. Larsen, T. Samuelsson, and C. Zwieb (1998) The signal recognition particle database (SRPDB). *Nucleic Acids Res.* **26**, 177–178. A description of the annotated database that contains a compendium of SRP RNA and protein sequences, sequence alignments, and links to crystal structures. The internet address has changed to <http://psyche.uthct.edu/dbs/SRPDB/SRPDB.html>

## Signal Transduction

Signal transduction encompasses a wide variety of biochemical processes that mediate the response of a cell to extracellular stimuli. These processes have been most extensively characterized in endocrine and neuronal cells, whose primary function is cellular communication. On the other hand, signal transduction governs the regulation of metabolism, growth, differentiation, and death in virtually all cells. In simple terms, signal transduction from the plasma [membrane](#) can be organized into three discrete components: receptors, effectors, and targets. Cell-surface receptors are bifunctional proteins that traverse the plasma membrane, interacting with extracellular molecules and transmitting information to the cytoplasmic portion of the membrane. This information is conveyed by an intermediary system, referred to as an *effector*, that generates different types and combinations of intracellular messengers. [Second messenger](#) systems generate further specificity in the response and provide an amplification of the extracellular signal. Finally, these messages are received by their intracellular targets, which often are [enzymes](#) that may undergo covalent modifications, altering their activity. In reality, signal transduction systems are complex networks of convergent and divergent pathways, with cross-talk mechanisms and feedback loops, allowing different [hormones](#) to elicit their characteristic effects on cells.

The precise relay systems used by receptors can be predicted to some extent by their structure and distribution. In general, there are four types of receptors that generate signals (Fig. 1): (i) [G-protein-coupled receptors](#) (GPCRs); these receptors are known to interact with a family of GTP-binding transducing molecules, which in turn regulate the activities of a variety of membrane proteins. (ii) **Ion channels**; these receptors allow the flow into the cell of specific ions, and they are directly regulated by ligand binding. (iii) Receptors with intracellular enzyme activities; such receptors bind ligands on the extracellular surface, resulting in the activation of an intracellular enzyme activity, such as a protein [kinase](#). These receptors typically have a single **transmembrane** region, along with large extracellular and intracellular **domains**; many of these proteins are receptors for [growth factors](#). (iv) Receptors lacking enzyme activities; these receptors also contain a single transmembrane domain, typically bind to **cytokines**, and interact specifically with intracellular proteins that modulate a variety of processes.

**Figure 1.** Structures of receptors involved in signal transduction.



Dramatic advances have been made in our understanding of the molecular mechanisms involved in signal transduction. It appears that many of the important components in signaling pathways have been identified, and studies on their structure and function are likely to produce a fairly detailed picture of what role they play in governing cell growth, differentiation, and survival. However, we still have much to learn regarding the molecular mechanisms that ensure specificity in signal transduction, especially regarding the temporal and spatial relationships of different pathways to each other, as well as the intrinsic limitations on signaling, and ways in which pathways cross-communicate within cells.

The details of signal transduction are described in [G-protein-coupled receptors](#); [Second messengers](#); [Heterotrimeric G proteins](#); [Phosphorylation](#); [Phospholipases](#); [Protein kinases A and C](#); [Calcium signaling](#); [Tyrosine kinase receptors](#); [Ion channel receptors](#); [MAP kinases](#); [Phosphatidylinositol kinases](#); [Receptors linked to tyrosine kinases](#); [JAK/STAT signaling](#).

## Silencer, Gene

A silencer is a [cis-acting](#) regulatory sequence that lowers the rate of initiation of [transcription](#) of eukaryotic genes. Silencer elements have been identified from a variety of organisms. The first and best-characterized silencer was identified near the promoters for mating-type genes at the HML and HMR loci in the yeast *Saccharomyces cerevisiae* (1). The yeast-mating-type genes, *a* and *a'*, at HML and HMR are transcriptionally silent and become active only when they replace the mating-type allele present at the MAT locus during the process known as mating-type switching. The silenced mating-type genes require E regulatory sites known as HMLE and HMRE, respectively. Molecular



analyses have identified a DNA element at the E regulatory site, called the silencer, capable of repressing transcription of a target **promoter** in *cis*. This element comprises binding sites for regulatory proteins. Interactions of the silencer element with these proteins are necessary and sufficient for forming the repressive [chromatin](#) structure (2, 3). This element was named a silencer because of its similarity to transcriptional [enhancers](#), despite their opposite effects on transcription: (i) it functions in either orientation; (ii) it exerts a silencing effect relatively independently of its position with respect to its target promoter; (iii) it is capable of repressing promoters other than its normal target. Subsequently, however, the term silencer has been also used to describe negatively acting regulatory elements that do not conform to the criteria used in the original definition (4).

### Bibliography

1. A. H. Brand et al. (1985) *Cell* **41**, 41–48.
2. S. Loo and J. Rine (1995) *Annu. Rev. Cell Dev. Biol.* **11**, 519–548.
3. S. G. Holmes, M. Braunstein, and J. R. Broach (1996) in V. E. A. Russo, R. A. Martienssen, and A. D. Riggs, eds. *Epigenetic Mechanisms of Gene Regulation*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 467–487.
4. S. Ogbourne and T. M. Antalis (1998) *Biochem. J.* **331**, 1–14.

### Suggestions for Further Reading

5. V. E. A. Russo, R. A. Martienssen, and A. D. Riggs, eds. (1996) *Epigenetic Mechanisms of Gene Regulation*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
6. S. C. R. Elgin (1995) *Chromatin Structure and Gene Expression*, IRL Press, Oxford.

### Silent Mutation

A silent mutation is a [mutation](#) that changes a **codon** into another codon that specifies the same [amino acid](#) because of the [degeneracy of the genetic code](#). By definition, a silent mutation can only be a base substitution mutation. If a silent mutation comes to fixation in a population or subpopulation, it is called a [synonymous substitution](#).

### Silver Stain

Silver stains are used for the visualization of [proteins](#) and nucleic acids on [polyacrylamide](#) gels and [blotting matrices](#). They have achieved widespread application because of their great sensitivity (1, 2). Most silver stains can detect proteins and nucleic acids in the nanogram range. This level of detection surpasses by 10- to 100-fold those that can be achieved with the commonly used nonfluorescent organic stains, such as the [Coomassie brilliant blue](#) stains for proteins. Silver stains are also employed in applications that previously employed **radiolabeling** because they can achieve comparable levels of detection without the environmental and economic concerns associated with the use of [radioisotopes](#) (1). The ability of silver stains to detect proteins and nucleic acids depends on a differential oxidation-reduction potential that is contingent on the presence of proteins or nucleic acids in the gels or membranes. This differential oxidation-reduction potential in the gel or

membrane permits selective reduction of ionic silver to the metallic state, resulting in the visualization or detection of proteins or nucleic acids (3). For example, if a site containing a protein or a nucleic acid has a higher reducing potential than the surrounding gel or matrix, the protein will be positively stained. Conversely, if the protein-occupied site has a lower reducing potential than the surrounding gel or matrix, the protein will appear to be negatively stained. It is possible to alter these relative oxidation-reduction potentials by changing the chemistry of the staining procedure so that proteins separated on polyacrylamide gels stain either negatively or positively (4).

## 1. Silver Stain Protocols

There are three basic silver-staining protocols (3): (1) the diamine or ammoniacal stains, (2) the nondiamine chemical reduction stains, and (3) the silver stains based on photo-development.

The diamine or ammoniacal stains maintain silver ion concentrations at a low level by the formation of silver diamine complexes with ammonium hydroxide. Image development is initiated by acidification of the ammoniacal silver solution, usually with citric acid in the presence of formaldehyde. The citric acid lowers the concentration of free ammonium ions. This results in the liberation of silver ions for reduction by formaldehyde to metallic silver, modulated by the presence of nucleic acid or protein. The diamine or ammoniacal silver stains have proved particularly useful for the staining of proteins separated on polyacrylamide gels that are thicker than 1 mm.

Nondiamine chemical reduction stains use silver nitrate to provide silver ions for reaction with protein and nucleic acid sites under acidic conditions. Image development is initiated by placing the gel in an alkaline solution, utilizing sodium carbonate and/or hydroxide and other bases to maintain the alkaline pH, while the silver ions are selectively reduced to metallic silver with formaldehyde. The formic acid produced by the oxidation of formaldehyde is buffered by the sodium carbonate. This type of stain is relatively rapid and simple to perform. The nondiamine stains generally work best with gels 1 mm or less in thickness.

Photodevelopment silver stains depend on the energy from photons of light to reduce ionic silver to metallic. The presence of proteins or nucleic acids in the gel or on the membrane enhance the photoreduction of ionic to metallic silver in gels impregnated with silver chloride. Photodevelopment stains are very rapid and simple. They can permit the visualization of protein patterns within 10 minutes after an [electrophoresis](#) separation. However, current photodevelopment protocols lack the sensitivity of the other silver-staining methods.

It is also possible to combine these protocols. For example, by combining photodevelopment with a nondiamine stain, a protocol was developed to detect proteins and nucleic acids on membranes. This combined protocol relies on the formation of silver halide precipitates in gels or on membranes, and it employs light to initiate the formation of silver nucleation centers on the silver halide crystals, followed by the use of a chemical development to deposit additional silver on the silver nucleation centers. This protocol provides a silver stain that is rapid and sensitive. It can detect proteins and nucleic acids in the nanogram range on thin membranes in under 15 minutes (5).

## 2. Quantitative Applications of Silver Stains

In silver staining, the accumulation of metallic silver generally displays a reproducible but limited linear relationship between the density of silver in the developed image and the concentration of protein or nucleic acid on the gel or membrane ((2),45). For most proteins, the linear portion of this relationship extends over a range of concentrations from 0.02 to 2 ng/mm<sup>2</sup>. Above 2 ng/mm<sup>2</sup>, saturation generally occurs. Saturation can be recognized by the occurrence of bands or spots with centers that are less intensely stained than the regions near the edges. As the staining-response curve for each protein species may vary, quantitative comparisons between different gels should be limited to bands or spots of the same protein on each gel. For example, the [actin](#) spot on one gel should be

compared only with an actin spot on another gel and not with a [transferrin](#) spot. Furthermore, gels to be compared should be run under similar conditions (eg, percentage acrylamide, stacking gel specifications), as migration distance may affect band or spot compression which, in turn, may influence the staining reactions.

The specific silver-staining curve displayed by each protein species is not limited to silver detection methods; similar unique responses have also been demonstrated with Coomassie blue and with the commonly used Lowry protein assay. These protein-specific staining curves arise because of the unique collection of reactive staining groups contained in each type of protein. By studying [amino acid](#) homopolymers, individual amino acids, and [peptides](#) of known sequence, it has been determined that the principal reactive groups are the free [amino groups](#) and the sulfur atoms (see [Cysteine \(Cys, C\)](#) and [Methionine \(Met, M\)](#) residues) contained by the proteins ([5](#), [6](#)). The importance of the basic amino acids has been further substantiated by studies of the relationship between the specific amino acid contents of proteins and their ability to stain with silver. Studies of silver-staining reactions of nucleic acids and their precursors indicate that the **purine** bases may be important. However, because individual nucleoside and nucleotide derivatives stain poorly, silver stains may require cooperative effects between several functional groups on the same molecule to form complexes with the silver ions prior to their reduction to metallic silver ([5](#)). Despite complexities due to the variations in the staining response and to the difficulties in controlling staining end points, it has been possible to compare proteins found in one gel with those in another gel quantitatively, using computerized image analysis and densitometry. The identification of proteins that occur in the same amounts in different samples has proved helpful because these proteins can serve as normalization standards ([7](#)). The identities of silver-stained proteins may also be determined directly from the stained gels by the use of **mass spectroscopy** ([8](#)).

### 3. Historical Perspective

Silver staining was introduced in 1979 as a general protein detection method for proteins separated by [polyacrylamide](#) gel electrophoresis (PAGE) ([1](#), [2](#)). However, the observation that silver nitrate, the main ingredient of silver stains, can blacken organic substances was reported as early as the twelfth century by one of the great alchemists, Albert the Great (Albertus Magnus), the eldest son of Count von Bollstädt. Although there were numerous experiments between the twelfth and nineteenth centuries concerning the ability of silver nitrate to stain organic materials, ranging from artistic applications to the development of “invisible inks,” the first modern science application is credited to Krause. He used it, in 1844, to stain fresh tissue for microscopic analysis ([9](#)). Soon thereafter, Golgi and Cajal revolutionized our understanding of the central nervous system by their analysis of sections of the human brain, stained with silver ([10](#), [11](#)). Whether the use of silver to enhance histological imaging was initiated by the development of silver-based photography is not known, but there was clearly a cross-fertilization of ideas. These early efforts with silver staining occurred just as the revolution of photographic imaging was beginning. Although Scheele recognized, as early as 1777, that the blackening of silver chloride crystals by light was caused by the formation of metallic silver, it wasn't until 1839 that William Fox Talbot utilized this ability of light to reduce ionic silver to metallic as the basis of the photographic process ([12](#)). By 1862, experiments to enhance the photographic process resulted in the introduction of organic reducing agents to help with the development of images that were initiated by exposure of silver halides to light. Cajal participated in these early experiments to enhance photographic techniques, and some of the protocols were clearly adapted to his development of histological stains. One of the main compounds that Cajal used for developing his photographic images and later for his histological stains, formaldehyde, is currently used in many of the silver-stain protocols for the visualization of proteins and nucleic acids separated by gel electrophoresis.

### Bibliography

1. C. R. Merrill, R. C. Switzer III, and M. L. Van Keuren (1979) Proc. Natl. Acad. Sci. USA **76**, 4335–4339.

2. R. C. Switzer III, C. R. Merrill, and S. Shifrin (1979) *Anal. Biochem.* **98**, 231–237.
3. C. R. Merrill (1986) *Acta Histochem. Cytochem.* **19**, 655–667.
4. C. R. Merrill and D. Goldman (1984) In *Two-Dimensional Gel Electrophoresis of Proteins* (J. E. Celis and R. Bravo, eds.), Academic Press, New York, pp. 93–108.
5. C. R. Merrill and M. Pratt (1986) *Anal. Biochem.* **156**, 96–110.
6. B. L. Nielsen and L. R. Brown, (1984) *Anal. Biochem.* **144**, 311–315.
7. C. R. Merrill, G. J. Creed, J. Joy, and A. D. Olson (1993) *Appl. Theoret. Electropho.* **3**, 329–333.
8. M. Wilm, A. Shevchenko, T. Houthaeve, S. Breit, L. Schweigerer, T. Fotsis, and M. Mann (1996) *Nature* **379**, 466–469.
9. C. Krause (1844) In *Handwörterbuch der Physiologie* (R. Wagner, ed.), Vieweg and Sohn, Braunschweig, Germany, pp. 108–186.
10. C. Golgi (1873) *Gazz. Med. Ital. Lombardia* **33**, 244–246.
11. S. R. Cajal (1903) *Trab. Lab. Invest. Biol. Univ. Madrid* **2**, 129–222.
12. B. Newhall (1983) *Latent Image, the Discovery of Photography*, Univ. of New Mexico Press, Albuquerque, NM, pp. 58–59.

### Suggestion for Further Reading

13. J. M. Eder (1945) *History of Photography*, trans. E. Eptean, Columbia University Press, New York, pp. 22–24. This book provides considerable historic information on the use of silver in image formation, particularly the early history. Although it is concerned primarily with photography, many of the principles apply to the silver staining of proteins and nucleic acids on gels.

### Simulated Annealing

Simulated annealing is a general procedure for finding a solution to an optimization problem ([1](#), [2](#)). Nuclear magnetic resonance (NMR) experiments can be used to produce constraints on the possible **secondary structures** and [tertiary structures](#) of a biological [macromolecule](#). The experimental observations usually include [nuclear Overhauser effects](#) (NOEs), which provide information about the distances between pairs of nuclei of the molecule, and three-bond coupling constants, which can provide estimates of [dihedral angles](#). Constraints may also be produced by experiments that define the orientation of [chemical shift](#) tensors or dipolar interactions within the molecular framework. Determination of the three-dimensional structure consists of finding a conformation that is consistent with all constraints. This is a kind of optimization problem; simulated annealing thus provides a method for finding tertiary structures that are consistent with known constraints.

The use of simulated annealing for determination of tertiary structure is based on the type of empirical molecular force field that is usually used in molecular modeling. Such a force field is a collection of mathematical expressions that attempt to account for the energies associated with the deformation of chemical bonds, nonbonded [van der Waals interactions](#) between atoms of a molecule, and [electrostatic interactions](#) between charged groups of the structure. Molecular modeling programs usually contain a facility for conformational energy minimization. This feature uses an algorithm to adjust the positions of the atoms of a molecule so that the total energy, computed using the force field of the program, is at a minimum.

The simplest approach to finding a molecular conformation consistent with a set of NMR constraints is to include a term in the molecular force field that adds nothing to the total energy of the molecule when a given constraint is satisfied, but adds a large, unfavorable energy contribution when that constraint is violated. With such terms present, the process of energy minimization will tend to drive the molecule toward a conformation in which all constraints are satisfied. The defect in this procedure is that, although all energy minimization algorithms are capable of finding a minimum energy, no way exists to guarantee that it will be the conformation of lowest energy. Significant energy maxima (barriers) may separate local energy minima from the global minimum, and these must be overcome before the conformation that satisfies all constraints represented in the force field is identified. Simulated annealing is useful at this stage of optimization problem.

If the mathematical expressions in a force field can be differentiated, it is possible to calculate the force experienced by any atom of the molecule for any specific conformation. By employing Newton's second law ( $F = ma$ ), one can calculate the acceleration of each atom and, therefore, the change of position of an atom with time. Using a molecular force field in conjunction with the solution of a family of equations that embody the second law is the basis for [molecular dynamics](#) simulations.

The method of simulated annealing is based conceptually on the behavior of materials as they undergo the transition from the liquid state to the solid. If a sample starts at high temperature and is cooled rapidly, a polycrystalline or amorphous material may be obtained. This disordered sample represents a material that is not in its lowest free energy state. If the liquid is cooled slowly (annealed), a single crystal can often be formed. In the crystal, molecules are ordered in a way that represents the conformation of lowest possible energy, that is, the global minimum energy.

A direct proportionality exists between the kinetic energy of a system of moving particles and the temperature of the system. To find the global energy minimum for a system that includes NMR constraints, a molecular dynamics simulation is run in which the system is subjected to a schedule of temperature changes. The system is held at a high temperature for a period of time and then cooled slowly. It is hoped that during the high temperature phase, there will be sufficient motions of the atoms of the structure that barriers separating local minima from the global minimum are overcome. The cooling phase of the calculation potentially "anneals" the system into the conformation representing the global energy minimum—presumably the conformation in which experimentally derived constraints are best satisfied. In practice, several cycles of heating and cooling (annealing) are used, and the best schedule for these must often be worked out by trial-and-error. (See also [Distance Geometry](#).)

#### Bibliography

1. S. Kirkpatrick, C. D. Gelatt, Jr. and M. P. Vecchi (1983) *Science* **220**, 671–680.
2. W. H. Press, B. P. Flannery, S. A. Teukolsky and W. T. Vetterling (1988) *Numerical Recipes in C*, Cambridge University Press, Cambridge, Chapter "10".

#### Suggestions for Further Reading

3. J. Cavanagh, W. J. Fairbrother, A. G. Palmer III and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic Press, San Diego.
4. P. Güntert (1997) In *Protein NMR Techniques* (D. G. Reid, ed.), Humana, Totowa, New Jersey, pp. 157–194.
5. A. R. Leach (1996) *Molecular Modeling: principles and applications*, Longman, Singapore.

## Sines

Distinctive families of both short and long [interspersed DNA elements](#), *SINEs*, and [Lines](#), respectively, are repeated within [genomes](#), in addition to the abundant primate family of short repeats called [Alu sequences](#). SINEs are typically shorter than 500 bp and occur on the order of hundreds to thousands of times within each genome. Each family of SINEs comprises a set of segments that are nonidentical but closely similar in length and sequence. The members of any given family are dispersed throughout the genome. SINEs are detected by [probe hybridization](#) using **radiolabeled** probes. They are found between **genes**, inside [introns](#) and within [satellite DNA](#), but never in regions coding for protein. SINEs are the major component of dispersed [repetitive DNA](#) in all mammalian genomes. Most SINEs contain an intragenic **RNA polymerase III promoter** that initiates [transcription](#) at the 5' end of the repeated DNA sequence and which, it has been proposed, facilitates the [transposition](#) and amplification of these sequences by an RNA-intermediate mechanism.

The term SINEs is restricted to short **retrotransposons**, which have internal RNA polymerase III promoter sites in a region derived from a structural RNA (usually a [transfer RNA](#)), that are amplified via [cDNA](#) intermediates and then enter the host genome. There is considerable evidence that these DNA [transposable elements](#) result from self-primed **reverse transcription** of their RNA transcripts with subsequent DNA integration. In contrast to [retroviruses](#) and retrotransposons, SINEs do not encode the [enzymes](#), such as reverse transcriptases, required for their amplification. Thus it is presumed that they borrow these enzymes from other sources. In the human genome, however, two families of [mobile elements](#) have been identified which have the sequence characteristics of transposons that move directly from DNA to DNA, rather than requiring the reverse transcription of an RNA intermediate. One type of element comprises a coding region for a [transposase](#) flanked by short terminal repeats of 31 or 32 bp (1).

Alu sequences are a type of SINE but, in contrast to Alu sequences, the other SINEs are not confined to primates. SINEs have been characterized in the genomes of many organisms, including mammals and other vertebrates (*Xenopus borealis*), salmonid fishes, in many **plants**, in the [nematode](#) *Caenorhabditis elegans*, in several filamentous **fungi** of the genus *Podospora*, and in [Dictyostelium discoideum](#). A SINE of 470 bp, occurring in 100 copies per haploid genome, isolated and characterized in *Magnaporthe grisea*, the rice blast fungus, as an insertion element within an inverted repeat transposon, shows the typical features of a mammalian SINE. At its 5' end, a secondary structural analysis reveals a tRNA-related region that could fold into a tRNA-like cloverleaf structure (2). DANA, the first SINE isolated from *Danio rerio*, the [zebrafish](#), is unique in its substructure of distinct cassettes. In contrast to classical SINE elements, it appears to have been assembled by insertions of short sequences into a progenitor tRNA-derived element. Once associated with each other, these subunits are amplified as a new transposable element with such success that DANA-related sequences form 10% of the modern zebrafish genome (3). Three different kinds of SINEs were isolated and characterized in *Octopus vulgaris*. Two of them seem to have been derived from tRNA<sup>Arg</sup>, and the third is from a tRNA that could not be identified because of sequence divergence from the original tRNA sequence (4).

A family of SINEs described in the hamster genome constitutes approximately 0.3 to 0.5% of the genome. The repeats are about 300 bp long and are highly divergent (differing by up to 30% from the **consensus sequence**). In contrast to the usual SINEs, the repeats are not flanked by short direct repeats and lack sequences corresponding to the RNA polymerase III promoter (5). Other SINEs had already been found highly homologous to tRNA genes, which suggested that many SINE mammalian families are amplified tRNA [pseudogenes](#) (6).

## Bibliography

1. G. T. Morgan (1995) *J. Mol. Biol.* **254**, 1–5.
2. P. Kachroo, S. A. Leong, and B. B. Chattoo (1995) *Proc. Natl. Acad. Sci. USA* **92**, 11125–11129.
3. Z. Izsvak et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1077–1081.
4. K. Ohshima and N. Okada (1994) *J. Mol. Biol.* **243**, 25–37.
5. C. Miles and M. Meuth (1989) *Nucleic Acids Res.* **17**, 7221–7228.
6. G. R. Daniels and P. L. Deininger (1985) *Nature* **317**, 619–622.

## Single Particle Reconstruction

Many biological [macromolecules](#) are very difficult to crystallize and thus are not amenable for electron or X-ray crystallography. Fortunately, techniques have been developed over the past two decades for image averaging and three-dimensional reconstruction of individual macromolecules with low or no symmetry. Presently, the attainable resolution for reconstructing these molecules in 3-D is 1.5–4 nm. Radiation damage of biological material remains the principal problem, since low-dose imaging requires averaging of thousands of individual images in order to boost the specimen signal to a level sufficient to allow construction of an atomic model. In theory, low-dose images provide enough information about the position and orientation of individual macromolecules or assemblies to provide atomic resolution if the molecular weight is greater than  $\sim 10^5$  dal and more than 10,000 particles are averaged (1). However, the quality of the best images as yet does not match theoretical calculations, which means that the molecular weight limit and the number of averaged particles required for atomic resolution may need to be an order of magnitude higher. The limitation in resolution comes about because of inaccuracies in determining the five orientation parameters needed to align low-dose images that are inherently noisy. Larger molecules or assemblies can be more accurately aligned simply because more orientational information is produced by the imaging process. The poor image contrast at high resolution of ice-embedded specimens needs to be overcome before atomic resolution can be challenged. This low contrast is thought to be caused by blurring of the image due to charging of the specimen by the electron beam and to beam-induced specimen movement (2).

Computer analysis of single-particle images was developed initially from methods to align individual images of macromolecules based upon cross-correlation functions versus a common reference. For 2-D images, three alignment parameters must be determined, a rotational orientation angle and two translational vectors. These can be determined from a rotational cross-correlation function (performed on the modulus of the Fourier transform of the images, since this function is translationally invariant) and from a subsequent translational cross-correlation function, respectively (3). Typically a second alignment procedure is performed using the averaged structure from the first alignment as a reference in the second. A second major development in single-particle methods was the use of correspondence analysis (a type of multivariate statistical analysis) to determine the most significant factors by which a collection of images varies (4). Typically images of identical molecules will cluster into groups that differ in their orientations, and if the orientation of each group can be determined, a 3-D reconstruction can be calculated using established computer algorithms such as the filtered backprojection. The strategy for collecting different tilted views of a specimen varies, but one common approach is the “random conical tilt” procedure (5). In this example, two electron micrographs are recorded of an area containing dozens or hundreds of molecules. The first micrograph is a low-dose exposure with the specimen grid tilted at some predetermined angle; the second exposure is a higher-dose exposure of the untilted grid. The higher-dose image is used to

identify individual particles and the orientations using cross-correlation alignment and correspondence analysis. From this information, the same particles can be identified in the low-dose, high-tilt exposure, and their orientations can be calculated knowing the orientations of the untilted views, the known tilt angle, and the orientation of the tilt axis. The 3-D structure of the molecule can then be calculated from the various images of the molecules identified in the low-dose image.

Although performed at relatively low resolution, single-particle approaches have provided significant insight into the operation of large macromolecular complexes. The structure and function of [ribosomes](#) have been explored (6-8), which has also produced the sites of association of the 30S and 50S subunits (9) and a direct visualization of tRNAs bound to the 70S ribosome (10, 11). Other large complexes studied in 3-D are the NADH–dehydrogenase complex (see Dehydrogenase (12)), the calcium release channel/ryanodine receptor (13, 14), the [nuclear pore complex](#) (15), and [hemoglobins](#) (16, 17). Not all structural characterizations need to be performed in 3-D as evidenced by the ligand and subunit locations found in 2-D in photosystem I and II (18, 19).

### Bibliography

1. R. Henderson (1995) *Quarterly Rev. Biophys.* **28**, 171–193.
2. R. Henderson (1992) *Ultramicroscopy* **46**, 1–18.
3. M. van Heel, M. Schatz, and E. Orlova (1992) *Ultramicroscopy* **46**, 307–316.
4. J. Frank (1990) *Quarterly Rev. Biophys.* **23**, 281–329.
5. M. Radermacher (1988) *J. Elec. Microsc. Tech.* **9**, 359–394.
6. P. A. Penczek, R. A. Grassucci, and J. Frank (1994) *Ultramicroscopy* **53**, 251–270.
7. H. Stark et al. (1995) *Structure* **3**, 815–821.
8. J. Frank (1997) *Curr. Op. Struct. Biol.* **7**, 266–272.
9. K. R. Lata et al. (1996) *J. Mol. Biol.* **262**, 43–52.
10. R. K. Agrawal et al. (1996) *Science* **271**, 1000–1002.
11. H. Stark et al. (1997) *Cell* **88**, 19–28.
12. V. Guenebaut et al. (1997) *J. Mol. Biol.* **265**, 409–418.
13. M. Radermacher et al. (1994) *J. Cell Biol.* **127**, 411–423.
14. E. V. Orlova et al. (1996) *Nature Struct. Biol.* **3**, 547–552.
15. J. E. Hinsha, B. O. Carragher, and R. A. Milligan (1992) *Cell* **69**, 1133–1141.
16. F. de Haas et al. (1996) *Proteins* **26**, 241–256.
17. F. de Haas et al. (1996) *J. Mol. Biol.* **264**, 111–120.
18. E. J. Boekema et al. (1995) *Proc. Natl. Acad. Sci.* **92**, 175–179.
19. C. Lelong et al. (1996) *EMBO J.* **15**, 2160–2168.

### Single-Stranded DNA Binding Protein

During several **DNA** metabolic reactions, double-stranded DNA is converted to single-stranded DNA temporarily by the actions of [DNA helicase](#) or **exonucleases**. The exposure of a single-stranded DNA region is necessary for [enzymes](#) or the complementary DNA strand to interact with the region, but this also produces regions that form unfavorable **secondary structures** and are susceptible to attack by **nucleases**. To avoid such problems and make the enzyme reactions efficient, all organisms have single-stranded DNA binding proteins (SSBs) that preferentially bind only single-



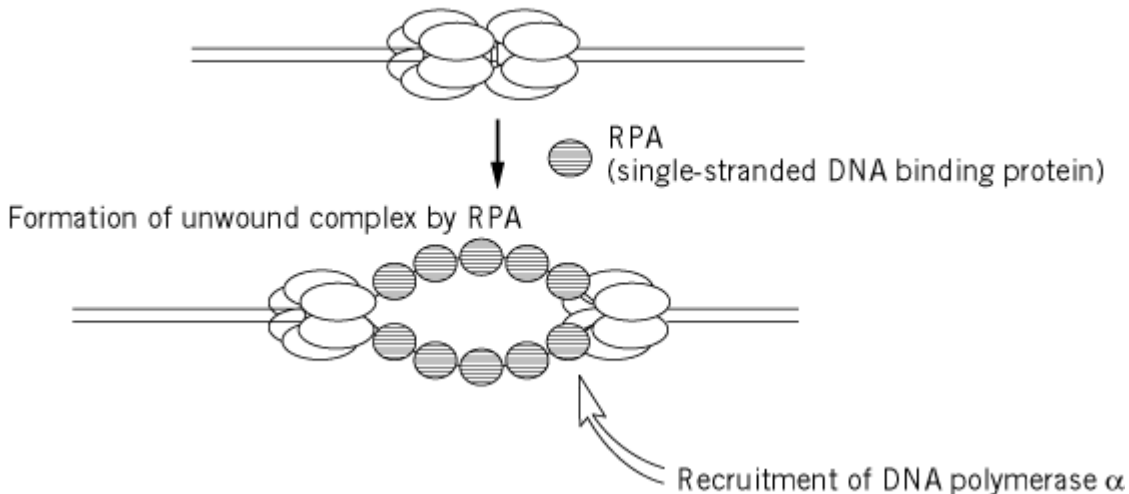
stranded DNA. Due to their involvement in vast tasks, they appear to be multifunctional and one of the major players in various cellular processes.

The first discovery of SSB was **bacteriophage** T4 gene32 protein. [Affinity chromatography](#) of T4 phage-infected cell lysate on single-stranded DNA cellulose successfully identified a protein that bound strongly and preferentially to single-stranded DNA (1). Genetic studies had shown that the gene product was essential for [recombination](#) and [DNA replication](#), indicating the importance of this type of **DNA-binding protein** in these processes. Following this discovery, many SSBs were identified from various systems, eg, from *Escherichia coli*, bacteriophages, and **viruses**. Although there are almost no similarities in their protein structures, they commonly bind to single-stranded DNA with more than  $10^3$  times greater affinity than to duplex DNA and have a crucial role in DNA metabolism. In general, their binding to single-stranded DNA is **cooperative**; this means that each molecule prefers to bind adjacent to another molecule along the DNA strand. As a result, many molecules of SSB coat DNA and make it an active target for many enzymes (2). This aspect of SSBs sometimes induces the **denaturation** of duplex DNA under physiological conditions, so that they were initially called unwinding or **helix-destabilizing proteins**. This function is distinguishable from DNA helicase action, since the denaturation of duplex DNA by SSB does not require ATP hydrolysis. SSB interacts with the backbone of single-stranded DNA rather than the bases. Thus, the binding is almost nonspecific for the DNA sequence, and the bound DNA is extended. The extended conformation facilitates base pairing of complementary strands and increases the rate of **DNA polymerase** migration; this function explains the requirement of SSB for replication and recombination.

SSBs have been sought from eukaryotic cells, and several single-stranded DNA-specific binding proteins were identified, but none of them were obviously demonstrated to be involved in DNA metabolism. The essential cellular SSB was first identified from human cell extracts as a component essential for *in vitro* [SV40 virus](#) DNA replication and named RPA (replication protein A, also called RFA or HSSB; for a review, see (3)). During the replication reaction, RPA interacts with the viral-coded initiator protein, large [T Antigen](#), and in collaboration with its helicase activity forms a [replication origin](#) unwound complex (Fig. 1). Another interaction between RPA and DNA polymerase  $\alpha$ /primase facilitates the recruitment of this polymerase on the unwound DNA and initiates DNA synthesis. In the following DNA elongation step, RPA stimulates DNA polymerases through either specific or nonspecific interactions (4). Therefore, RPA has critical roles throughout this viral DNA replication and will have the same function during chromosomal DNA replication. Related to this point, studies by immunostaining with anti-RPA **antibodies** have revealed the tight interaction of RPA with the replication apparatus. RPA is localized in the [nucleus](#), exists in the [replication foci](#) prior to the initiation of DNA synthesis, and remains there during **S phase** (5).

**Figure 1.** Unwinding of the SV40 virus replication origin by SV40 T-antigen and RPA.

### ATP dependent binding of Tag helicase to SV40 replication origin



Unlike prokaryotic SSBs, which are monomers or oligomers of identical subunits, human RPA has three heterogeneous subunits of 70, 32, and 14 kDa. This oligomeric structure is conserved among eukaryotic SSBs, so homologous SSBs in heterotrimeric complexes have been identified in all tested eukaryotes from yeast to human (3). Since eukaryotes have several specific cellular processes, such as the mitotic [cell cycle](#), DNA damage signaling, [chromatin](#) formation, and [transcription](#) activation, this complex structure of RPA might reflect its extra roles in these processes, in addition to DNA metabolism. For example, RPA interacts directly with nucleotide [excision repair](#) proteins XPA, XPG, and XPF (xeroderma pigmentosa group A, G, and F proteins) and plays crucial roles in damage recognition and cleavage (6). Another type of RPA-interacting protein includes transcription activators such as GAL4, VP16 (7), and tumor suppressor [p53](#) protein (8), suggesting the involvement of RPA in transcription regulation. Indeed, RPA has been isolated as a [transcription factor](#) for several genes in yeast (9). In addition to [protein-protein interactions](#), the 32k-Da subunit of RPA is **phosphorylated** in a cell-cycle-dependent manner and by DNA damage. An up-shift of the mobility of the 32-kDa subunit in [gel electrophoresis](#) was observed upon the phosphorylation that takes place efficiently in the presence of single-stranded DNA (10). Although the phosphorylation of RPA has been studied intensively, no direct evidence to connect it with the modulation of RPA functions could be obtained. However, if we consider this specific phosphorylation, the localization of RPA in nuclei, and its interaction with several important factors for various cellular processes, we see that eukaryotic single-stranded DNA-binding proteins are not merely DNA metabolic proteins, but will also play a role in coordinating DNA metabolism with several cellular processes.

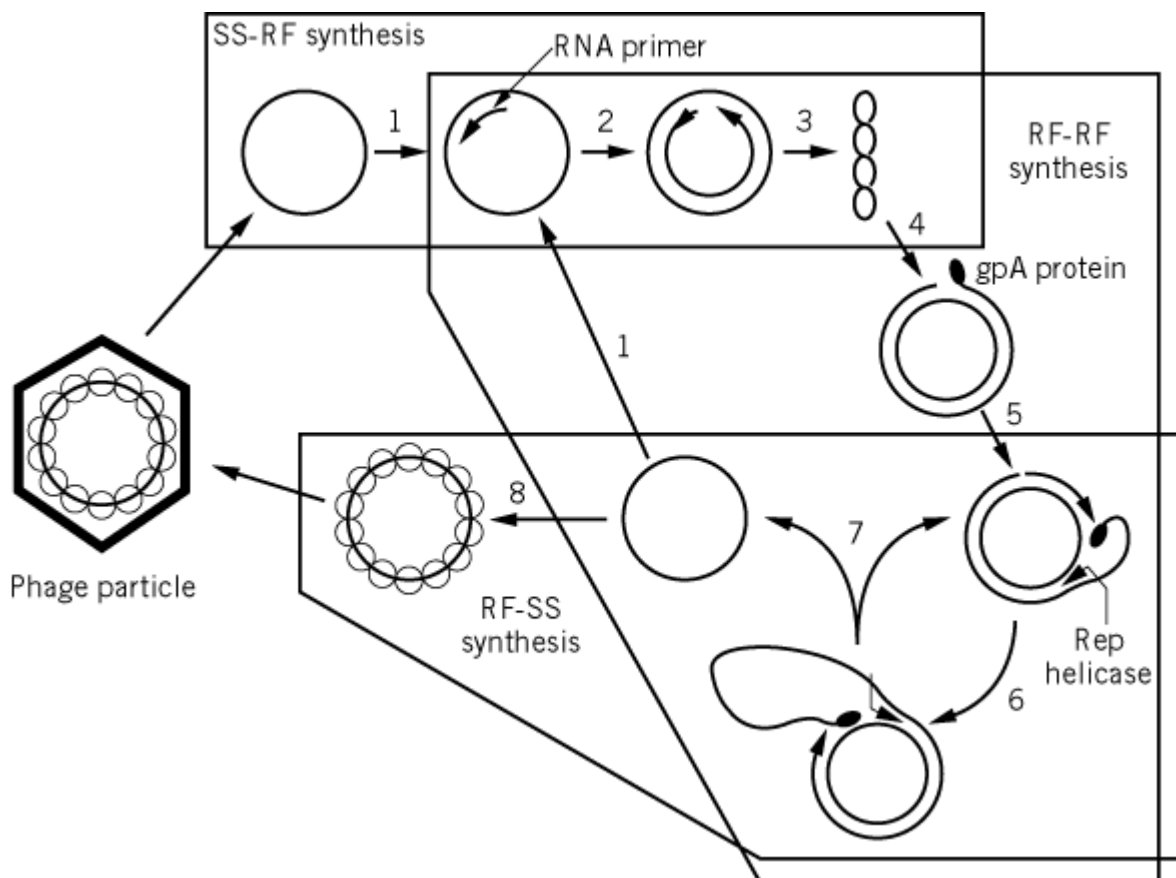
### Bibliography

1. B. Alberts and G. Herrick (1971) *Meth. Enz.* **21**, 198–217.
2. T. A. Baker and A. Kornberg (1992) *DNA Replication*, 2nd Ed. New York, W. H. Freeman, pp. 325–339.
3. M. S. Wold (1997) *Ann. Rev. Biochem.* **66**, 61–92.
4. M. K. Kenny et al. (1990) *J. Biol. Chem.* **265**, 7693–7700.
5. Y. Adachi and U. K. Laemmli (1994) *EMBO J.* **13**, 4153–4164.
6. L. Li et al. (1995) *Mol. Cell Biol.* **15**, 5396–5402.
7. Z. He et al. (1993) *Cell* **73**, 1223–1232.
8. A. Dutta et al. (1993) *Nature* **365**, 79–82.
9. R. M. Luche et al. (1993) *Mol. Cell Biol.* **13**, 5749–5761.
10. F. Fang and J. W. Newport *J. Cell Sci.* **106**, 983–994.

## Single-Stranded DNA Replication

Certain small **bacteriophages** contain a circular, single-stranded **DNA** (ssDNA) as **genome**. Replication of the single-stranded DNA has been characterized extensively for two groups of *Escherichia coli* phage, filamentous phages represented by M13 and spherical phages represented by fX174 (1). Their replication can be divided into three steps (Fig.1);(1) conversion of the ssDNA genome to a double-stranded form, called the **replicative form** (RF); (2) multiplication of RF DNA by **rolling circle replication**; and (3) generation of ssDNA genome for packaging into phage from the RF DNA. Thus, replication of the ssDNA genome depends on the synthesis by DNA polymerase of a complementary strand in a sophisticated manner.

**Figure 1.** Scheme for single-stranded DNA replication of *E. coli* phage fX174. Three stages of the replication, SS → RF, RF → RF and RF → SS synthesis, are indicated by boxes (see text for details).



### 1. Conversion of ssDNA into Double-Stranded RF Form

Immediately after infection by the ssDNA genome (the plus strand), synthesis of the complementary strand (the minus strand) begins using host **enzymes**. As all the phage mRNA is synthesized using

the minus strand as [template](#) strand, no phage [messenger RNA](#) is synthesized at this stage. Synthesis of the minus strand is primed using an RNA [primer](#) made at a specific site on the plus strand, called minus *ori* or single-strand *ori* (sso), (step 1 in Fig. 1). Then the host elongation proteins synthesize the minus strand from the primer terminus. Three different modes of priming are known: by **RNA polymerase** (2), by DnaG primase (3), and by a seven-component *primosome* (4). A general feature of minus *ori* is an extensive region of **secondary structure**, which probably stabilizes the recognition sequences for the primer synthesis against melting by [single-stranded DNA binding protein](#) (SSB). After completion of the minus-strand synthesis, the initial product is a duplex circular DNA (RFII) consisting of an intact viral DNA and a nearly full-length linear complementary DNA plus the RNA primer (step 2 in Fig. 1). Then DNA polymerase removes the RNA primer and fills the resultant gap. Subsequently, the nick is sealed by [DNA Ligase](#), and the duplex DNA is converted to the superhelical form (RFI) by the action of DNA gyrase (step 3 in Fig. 1) (see [DNA Topology](#)). The RFI DNA serves as a template for RNA [transcription](#), and the production of phage proteins begins.

## 2. Rolling Circle Replication of RF DNA

Multiplication of RF DNA is initiated by introducing a nick by an initiator **nuclease** encoded in the phage genome at a particular site in the plus strand (plus *ori* or double-strand *ori*, dso) (step 4 in Fig. 1). The free 3'-OH end at the nick site serves as a primer and is elongated by DNA polymerase III holoenzyme, as the host Rep [DNA helicase](#) peels off the existing plus strand (step 5 in Fig. 1). When synthesis of the new plus strand reaches the plus *ori* (step 6 in Fig. 1), the nick is introduced again at the junction of the old and newly synthesized plus strand, to form linear ssDNA of one unit length (step 7 in Fig. 1). The linear form is then converted into the circular form, a further reaction catalyzed by the initiator nuclease. Thus, one round of the rolling circular replication produces one new plus strand, and it is used as a template to form another RF DNA.

The molecular mechanism of action of the initiator nuclease of ssDNA phage has been best characterized for that of fX174, gpA (product of gene A) (5). The gpA nuclease recognizes and binds to a specific sequence in the plus *ori* and introduces a nick in a nearby sequence. Then the gpA protein covalently attaches to the 5'-P end at the nick by a phospho bond to a **tyrosine residue**, and the energy liberated by the cleavage of the DNA strand is stored in the protein–DNA complex (see step 4 in Fig. 1). GpA has affinity for the host Rep DNA helicase and recruits it to the nicked site. Furthermore, the interaction locks the gpA protein to the template throughout the replication (see steps 5 and 6 in Fig. 1). When synthesis of the new plus strand reaches the regenerated plus *ori*, the gpA nuclease introduces a second cleavage in the plus strand, and the 5'-P group of the old strand attached to the gpA protein is transferred to the newly created 3'-OH end, to produce circular ssDNA. At the same time, the new 5'-P end is transferred to the Tyr residue of the gpA protein, and the next cycle of the rolling circle replication begins (see step 7 in Fig. 1). The fX174 circle is synthesized in about 10 sec, and 20 or more ssDNA circles are released from a single rolling circle intermediate. It should be noted that **supercoiling** of the template is required for the first cleavage, whereas the second cleavage occurs on the relaxed template. An explanation for this difference may be that the supercoiling is required for specific binding of the gpA protein to the template and that, since the gpA is already bound to the DNA, recognition of the nicking site is sufficient for the second cleavage.

Covalent linkage between the initiator nuclease and DNA strand is not observed for filamentous phage, and the mechanism by which the energy liberated by the cleavage is stored is not clear in this case. Their rolling circle replication is limited to a single round.

## 3. Generation of ssDNA Genome for Packaging

Accumulation of supercoiled RF DNA provides the many copies of the template for [transcription](#) of **genes** encoding phage structural proteins and phage proteins necessary for the assembly of phage particle. When a certain amount of such proteins has accumulated, synthesis of the minus strand (SS

to RF synthesis) is inhibited, and the accumulated single-stranded plus (viral) DNA is packaged into the phage particles (step 8 in Fig. 1). This switch from RF synthesis to phage assembly depends on accumulation of a specific phage gene product. For example, gp5 protein of filamentous phage blocks the minus strand synthesis by coating the displaced viral strand, to form the [nucleoprotein](#) precursors for packaging.

#### 4. Similar Replication Mode of Plasmids of Gram-Positive Bacteria

Single-stranded DNA phage have not been reported for [Gram-positive](#) bacteria. However, many small plasmids of Gram-positive bacteria follow in essential detail the pattern and strategy used by ssDNA phage of Gram-negative bacteria for multiplication of RF DNA (6). These plasmids encode an initiator protein that introduces a site-specific nick and produce a single-stranded circle as replication intermediate. The ssDNA is converted to the double-stranded form by the host enzymes, as is the case for the SS to RF conversion of ssDNA phages. The similarity in the amino acid sequences of the initiator endonucleases of the two groups and of their recognition sequences clearly indicates that they have evolved from a common ancestor (7).

#### Bibliography

1. A. M. Campbell (1996) *In Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 2325–2338.
2. K. Geider, E. Beck, and H. Schaller (1978) *Proc. Natl. Acad. Sci. USA* **75**, 645–649.
3. J. P. Bouche, L. Browen, and A. Kornberg (1978) *J. Biol. Chem.* **253**, 765–769.
4. J. Schlomai, L. Polder, K. Arai, and A. Kornberg (1981) *J. Biol. Chem.* **256**, 5233–5238.
5. R. Hanai and J. Wang (1993) *J. Biol. Chem.* **268**, 23830–23836.
6. D. R. Helinski, A. E. Toukdarian, and R. Novick (1996) *In Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 2295–2324.
7. T. V. Ilyina and E. V. Koonin (1992) *Nucl. Acids Res.* **20**, 3279–3285.

#### Suggestion for Further Reading

8. A. Kornberg and T. A. Baker (1992) *DNA Replication*, 2nd ed., W. H. Freeman, New York.

### Site-Directed Mutagenesis

Site-directed mutagenesis is a method used to alter the **nucleotide sequence** of cloned **DNA** at a pre-defined position or *site*. This method is one form of *in vitro* [mutagenesis](#), because it involves a series of biochemical steps using purified reagents and is performed on **cloned** DNA in a small test tube. Site-directed mutagenesis can be used to change only a single base pair, to change a few base pairs, or to create more extensive sequence changes, such as deletions and insertions. The ability to create site-directed mutations in cloned DNA largely became possible with the advent of methods to synthesize oligonucleotides, short segments of DNA, to create the desired DNA changes (see [DNA Synthesis](#)). Although methods using chemical modification of DNA can create site-directed mutations in cloned DNA (see [Mutagenesis](#)), site-directed mutagenesis is most easily and precisely accomplished using an oligonucleotide; hence, the method is often termed *oligonucleotide-directed mutagenesis*. This entry describes the basic principles of oligonucleotide-based site-directed mutagenesis and highlights a number of procedures currently in use.

## 1. Site-Directed Mutagenesis is a Powerful Tool for Studying Nucleic Acids and Proteins

The ability to alter cloned DNA at will provided molecular biologists with a powerful tool with which to study the role of cloned DNA in cellular processes. Prior to the development of site-directed mutagenesis, biologists' only source of DNA variants was from rare spontaneous mutations that occurred randomly *in vivo* or from imprecise methods for *in vitro* chemical mutagenesis of cloned DNA. In contrast, oligonucleotide-based site-directed mutagenesis changes the sequence of a DNA fragment precisely as designed by the molecular biologist. Cloned genes can be studied in a biological system by comparing the function of the normal gene to variants containing site-directed mutations.

This revolutionary technique is not limited to the study of **genes**. Because the flow of genetic information is normally DNA → RNA → protein (according to the central dogma), site-directed mutagenesis of a cloned DNA fragment can also be used to study **RNA** and **proteins**. For example, site-directed mutagenesis has been used to uncover nucleotides in messenger RNA precursors that direct the **RNA splicing** of **introns** (1), to characterize amino-acid residues critical for **protein–protein interactions** (2), and to probe the **catalytic** mechanism of **enzymes** (3). There are many additional examples where site-directed mutagenesis has helped to solve basic problems of biology, and it has also been applied in the biotechnology industry to generate novel therapeutics.

Since the first report of oligonucleotide-based site-directed mutagenesis (4) by M. Smith and colleagues in 1978, site-directed mutagenesis has developed into a number of related methodologies that are now part of every molecular biologist's toolbox. Kits containing the reagents and protocols needed to conduct site-directed mutagenesis are commercially available and make the technique relatively straightforward to perform. In 1993, the method was recognized as being truly revolutionary by the award of the Nobel Prize for Chemistry to Smith for his “fundamental contributions to the establishment of oligonucleotide-based, site-directed mutagenesis and its development for protein studies” (5).

## 2. Two General Methods for Oligonucleotide-based Site-directed Mutagenesis

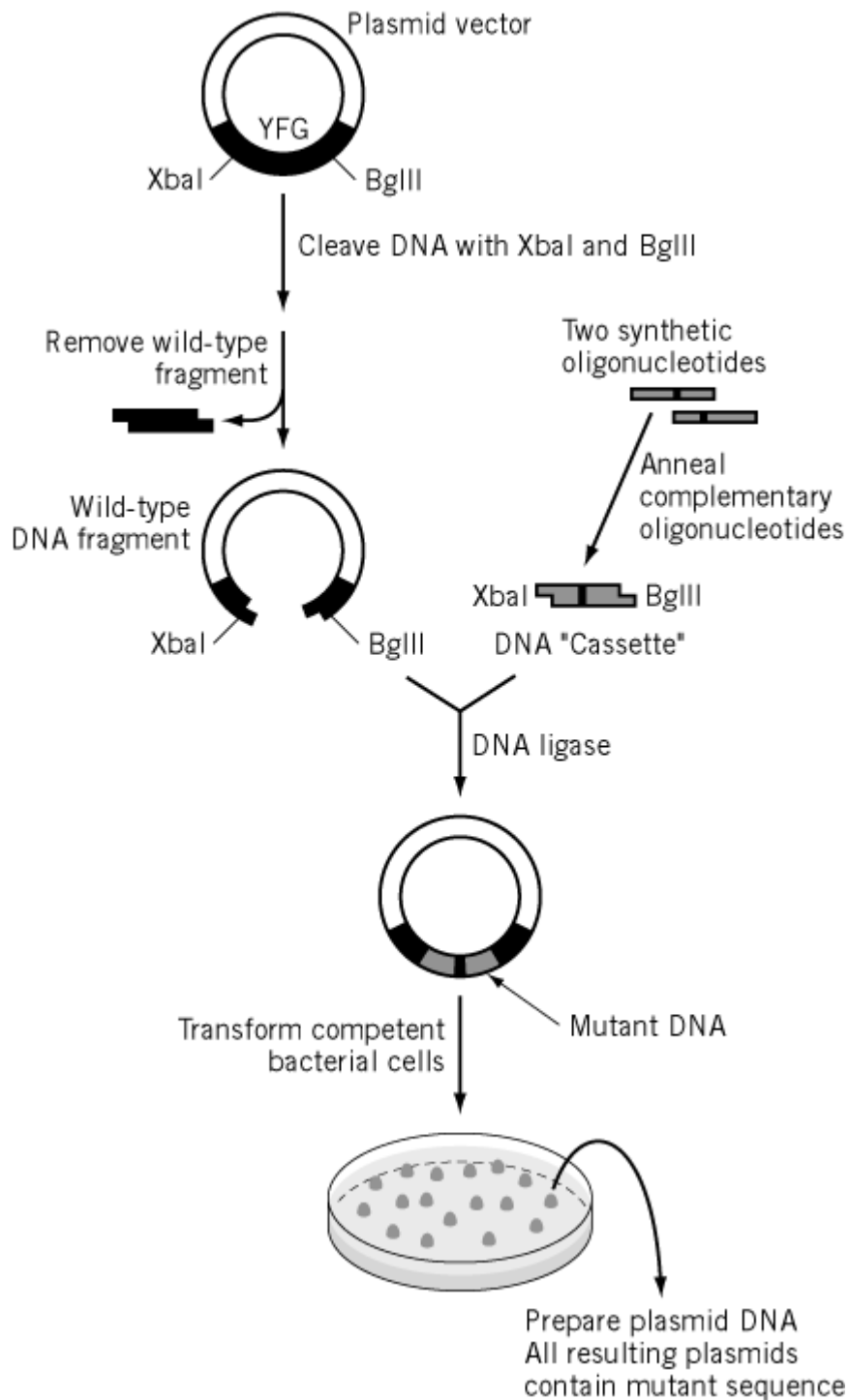
Current procedures for oligonucleotide-based site-directed mutagenesis can be grouped into two categories: (1) cassette mutagenesis and (2) enzymatic extension of a mutagenic oligonucleotide annealed to a DNA template. Cassette mutagenesis is, in practice, the simpler of the two approaches. As described in the next section, cassette mutagenesis is a derivative of total gene synthesis (6) in which a gene is constructed by ligating a series of synthetic oligonucleotides.

Site-directed mutagenesis by enzymatic extension of mutagenic oligonucleotides stemmed from research in the 1960s and 1970s on DNA polymerization *in vitro*, the genetics of bacteriophage  $\phi$ X174, and the preparation of synthetic oligonucleotides. Researchers were studying the ability of the enzyme **DNA polymerase I** from *Escherichia coli* to synthesize DNA in a test tube using an oligonucleotide primer hybridized to a single-stranded DNA template from  $\phi$ X174 **bacteriophage** (7). These experiments demonstrated that a single-stranded DNA template could be converted to a double-stranded molecule *in vitro*. Importantly, the DNA polymerase needed a short oligonucleotide primer to initiate DNA synthesis. Additional experiments on  $\phi$ X174 studied the ability of **restriction fragments** from wild-type  $\phi$ X174 annealed to a mutant form of  $\phi$ X174 to create wild-type phage by a process called *marker rescue* (8). These experiments, and research on oligonucleotide synthesis, led to the first site-directed mutagenesis experiment, where an oligonucleotide with only 12 nucleotides was used to change a single G to A in the  $\phi$ X174 DNA. Since the initial experiments on  $\phi$ X174 (4, 9), site-directed mutagenesis protocols have been developed for DNA fragments cloned into a variety of **vectors**. These methods can be subdivided further, depending on the particular DNA polymerase used. Early methods used DNA polymerases derived from *E. coli* or bacteriophage T4 or T7. Recent procedures have been developed for DNA polymerases from thermostable bacteria and the polymerase chain reaction (**PCR**). The details of these methodologies are discussed below.

## 2.1. Synthetic DNA Cassettes

In cassette mutagenesis (see Fig. 1), a restriction fragment from the cloned DNA of interest is replaced by another restriction fragment containing the desired nucleotide changes (10-12). The process begins by digesting the cloned DNA with one or two [restriction enzymes](#) to release a DNA fragment containing the (wild-type) sequence to be changed. The new (mutant) sequence is constructed from two synthetic, complementary single-stranded oligonucleotides, which are first mixed together in a buffered solution so they hybridize, thereby forming a duplex DNA molecule, but with [cohesive, sticky ends](#) to hybridize with restriction sites of the vector. To position the synthetic fragment back into the vector, the ends of the synthetic DNA must be the same as that in the wild-type fragment. The remainder of the cloned DNA in the vector is mixed with the synthetic fragment, and the two are ligated into place by the action of the enzyme [DNA Ligase](#) from bacteriophage T4. Following the joining of the two DNA molecules, DNA from the ligation reaction is introduced into competent *E. coli* cells, and recombinant molecules containing the desired mutation are selected.

**Figure 1.** Cassette mutagenesis. A plasmid containing a clone of your favorite gene (YFG; black segment) is cleaved with two restriction enzymes, eg. XbaI and BglII, each of which has only one restriction site in the entire plasmid. The reaction mixture is separated by [agarose](#) gel electrophoresis, and the larger fragment is purified from the gel. Two single-stranded oligonucleotides are synthesized by automated DNA synthesis. The sequences of the oligonucleotides are complementary to each other and differ from the wild-type sequence at only a single position (black stripe) containing the desired changes. The oligonucleotides are mixed in a solution that promotes hybridization of the two strands by virtue of their complementary sequences. The ends of the duplex fragment are single-stranded, cohesive, sticky ends that join with XbaI and BglII sites. The DNA cassette is mixed with the isolated fragment, and the two molecules are covalently joined by the action of T4 DNA ligase. The ligated DNA is transformed into *E. coli*, and drug-resistant colonies are selected. Plasmid DNA is prepared from individual bacterial colonies. Since the two linear fragments themselves cannot transform *E. coli*, all colonies contain plasmids with the mutant sequence.



Cassette mutagenesis generally results in nearly 100% of the resulting clones having the desired new (mutant) sequence. The difficulty with this method is that unique, conveniently spaced restriction endonuclease recognition sites to excise a small cassette in the desired place are often not present. In this case, one can turn to enzymatic extension methods.

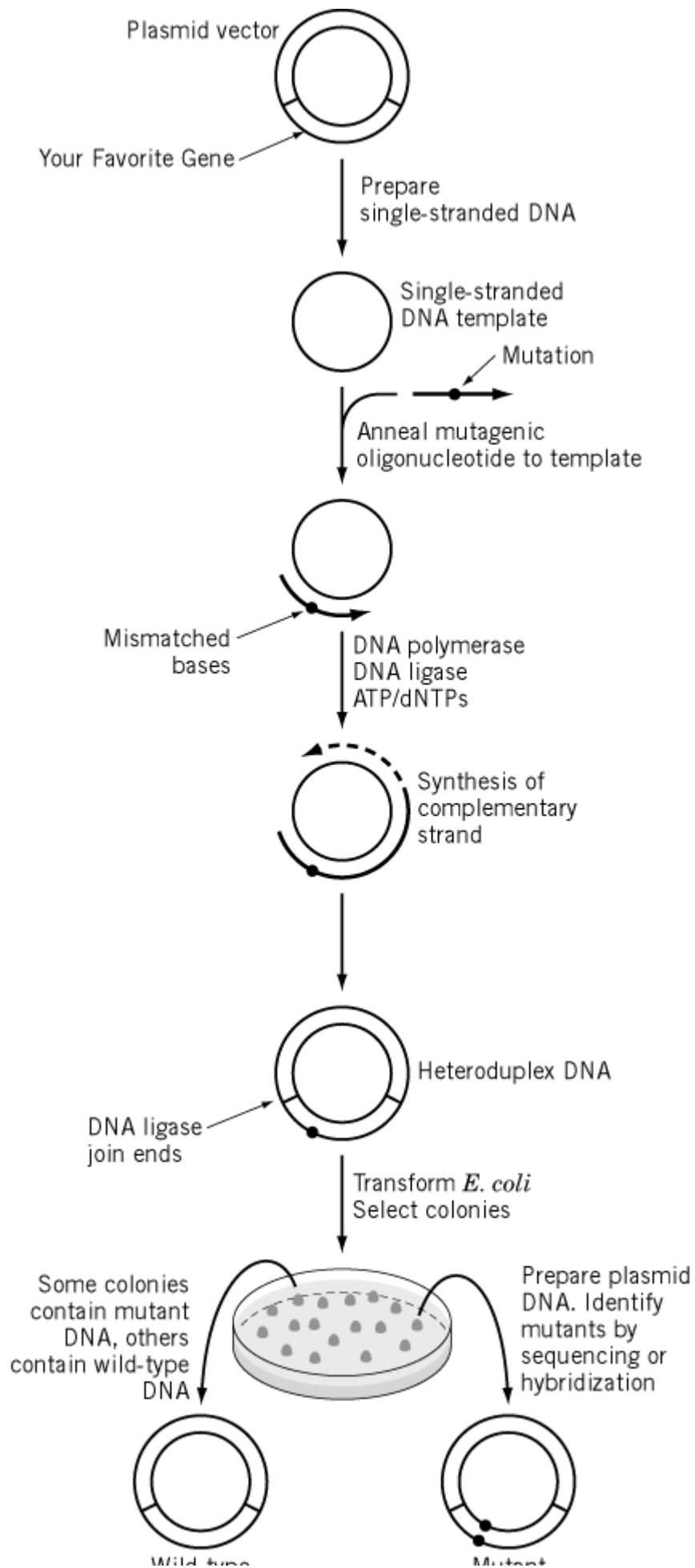
## 2.2. Enzymatic Extension of a Mutagenic Oligonucleotide

The principle and end result of site-directed mutagenesis by oligonucleotide extension are similar to those of cassette mutagenesis: An oligonucleotide encoding the desired new DNA sequence is inserted into a cloned DNA fragment. In the basic method (Fig. 2), a short oligonucleotide, typically 20 to 40 nucleotides in length, is hybridized to its complementary sequence in a circular, single-



stranded, wild-type DNA template. The sequence of the oligonucleotide is designed so that the new (mutant) sequence is in the middle of the oligonucleotide and flanked by wild-type sequences on the ends. Upon annealing the oligonucleotide to the template DNA, the wild-type sequences within the oligonucleotide are perfectly complementary to the template DNA, but the mutant sequences are mismatched. The wild-type sequences need to be long enough so the oligonucleotide will form a stable duplex with the template DNA. Next, a DNA polymerase and nucleotide precursors (dNTPs) are added to synthesize the remainder of the complementary strand using the mutagenic oligonucleotide as a primer for DNA synthesis. In the past, derivatives of *E. coli* DNA polymerase were used in site-directed mutagenesis. Recently, DNA polymerases from bacteriophage T4 or T7 have replaced the *E. coli* enzyme. These DNA polymerases are more processive than the *E. coli* enzyme (see [DNA Polymerase Sliding Clamps](#)) and do not displace or proofread the mutagenic oligonucleotide, thereby improving the frequency of the mutant plasmids obtained.

**Figure 2.** Site-directed mutagenesis by enzymatic extension of a mutagenic oligonucleotide. A DNA fragment encoding YFG is cloned into a plasmid vector. Single-stranded DNA is prepared by one of several methods. A synthetic oligonucleotide is annealed to the DNA template. The sequence is perfectly complementary to the wild-type sequence of the template, except for the mismatched nucleotides (black dot) that contain the mutant DNA sequence. The remainder of the complementary strand is synthesized by a DNA polymerase using the mutagenic oligonucleotide as a primer. Upon completing the circle, the ends of the newly synthesized DNA strand are joined by DNA ligase, forming a heteroduplex of two covalently closed DNA molecules. One strand contains the mutant sequence, and the other contains the wild-type sequence. The heteroduplex DNA is introduced into competent *E. coli* cells. Drug-resistant colonies that contain plasmid DNA are selected. The two strands separate by DNA replication in the bacterial cells. Plasmid DNA is isolated from the bacteria and sequenced to identify clones with the desired mutation.



Once synthesis of the complementary DNA strand has proceeded completely around the circular template, the ends of the newly synthesized strand are joined enzymatically by T4 DNA ligase. In this way, a single-stranded DNA molecule is converted into a double-stranded molecule, where all sequences will be wild-type except the region containing the mismatched nucleotides. The resulting double-stranded *heteroduplex* DNA molecule, composed of one wild-type strand and one mutant strand, is introduced into competent *E. coli* cells, where the two strands segregate by [DNA replication](#). Theoretically, each colony should contain both wild-type and mutant plasmids; in practice, however, only one sequence is often obtained from a colony due to [mismatch repair](#) prior to replication of the two strands. Generally, the frequency of mutant plasmids is lower than that of the wild type. This is most likely due to incomplete DNA replication *in vitro*. A number of selection or screening methods have been developed to enrich for plasmids with the desired mutation (see below).

### 3. Modifications and Improvements to Site-directed Mutagenesis Methods

The first oligonucleotide-directed mutagenesis experiments were performed using DNA from the single-stranded bacteriophage  $\phi$ X174. Because  $\phi$ X174 is not a convenient cloning vector, mutagenesis of cloned DNA fragments was performed using plasmids such as **pBR322** ([13](#)) or vectors derived from the filamentous [M13 phage](#) ([14-17](#)) and fd ([18](#)). The phage vectors offered an easy way to prepare a single-stranded DNA template. However, protocols were also developed for oligonucleotide-directed mutagenesis in double-stranded plasmid vectors ([19, 20](#)). In addition, **phagemid** vectors ([21, 22](#)) were developed that normally are double-stranded but can be induced to produce single-stranded phage-like DNA. These vectors facilitated template preparation for site-directed mutagenesis.

The early methods for oligonucleotide-directed mutagenesis resulted in a low frequency of plasmids with the desired mutation. Often, the fraction of mutants was less than 1 out of 100 plasmids. As discussed above, this was probably due to a variety of technical and biological problems. A number of methods were introduced to make it easier to find a plasmid carrying the desired mutant sequence. For some experiments, mutant molecules can be identified if the mutation happens to create or destroy a restriction endonuclease site. In a more general method, the double-stranded heteroduplex molecules were purified by centrifugation ([17, 23](#)). A simpler procedure was to make the mutagenic oligonucleotide **radioactive** and to use it as a hybridization probe ([24](#)). Under certain stringent conditions, the probe would hybridize with only mutant molecules and not with wild-type molecules. Other high-frequency mutagenesis methods used two oligonucleotide primers ([25, 26](#)); one was the mutagenesis oligonucleotide, and the other was positioned upstream. The second primer helped to complete synthesis of the complementary strand.

A powerful adaptation of the original procedure ([27](#)), developed by Kunkel in 1985, yielded mutants with efficiencies greater than 50%. First, the DNA vector was grown in a *dut*<sup>-</sup>, *ung*<sup>-</sup> strain of *E. coli* that causes uracil to be incorporated into the DNA template. Single-stranded template was prepared and annealed with the mutagenic oligonucleotide; the complementary strand was synthesized using standard deoxyribonucleotide triphosphates. In this way, the mutant strand contains thymine and the wild-type strand contains uracil. The heteroduplex molecule is introduced into an *ung*<sup>+</sup> *E. coli* strain, where the wild-type strand is destroyed by uracil deglycosylase (the product of the *ung* gene) cleavage of sites containing uracil.

Several other methods are available that increase the efficiency of mutagenesis. One method incorporates thiophosphoryl nucleotides into the newly synthesized strand, allowing the wild-type strand to be converted into the mutant sequence *in vitro* ([28](#)). Yet another method uses a second primer that alters a restriction site or corrects a mutation in a drug-resistance marker ([29, 30](#)). In this

way, plasmids containing the desired mutation can be enriched by drug selection or restriction enzyme digestion. By most of these improved methods, the fraction of mutant molecules obtained was often greater than 50%. Once these highly efficient methods were established, mutants were identified simply by randomly picking one or two plasmids and subjecting them to [DNA sequencing](#). Finally, the ability to perform multiple mutagenesis experiments in parallel has been addressed by the development of **solid-phase** mutagenesis methods ([31](#)). By this approach, the creation of site-directed mutants can be automated.

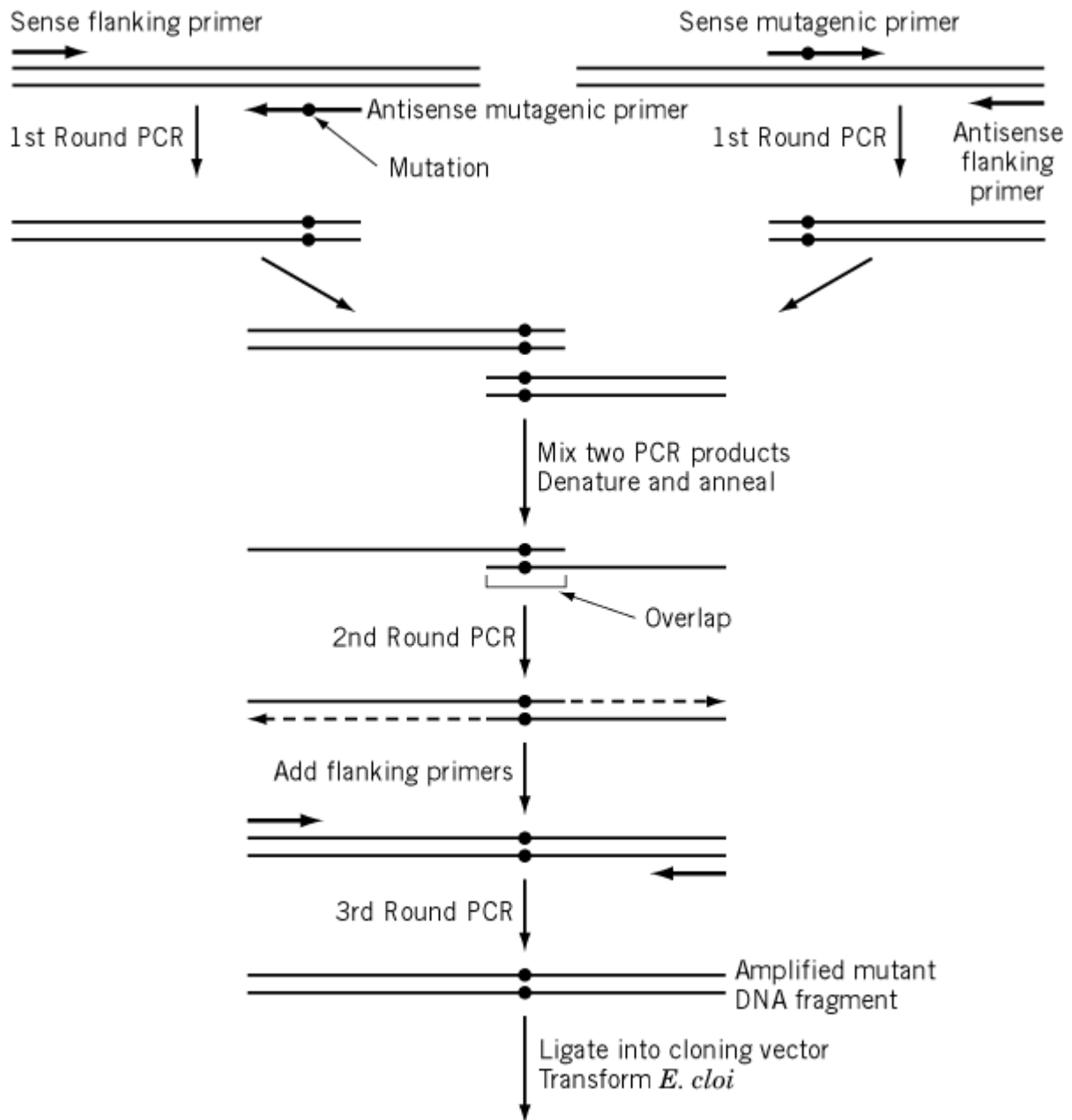
#### 4. Oligonucleotide-Directed Mutagenesis using the Polymerase Chain Reaction

New approaches to molecular biology problems are constantly being developed; they are usually simpler, cheaper, and more efficient. Soon after the polymerase chain reaction (PCR) was reported ([32](#)), it was realized that it could be adapted to create site-directed mutants (reviewed in ([33](#))). The PCR procedure was developed to amplify a segment of DNA from a minute amount of a DNA sample. The basic method uses two oligonucleotides as primers for DNA synthesis from a DNA template (see **PCR**). The advantage of using PCR for site-directed mutagenesis stems from its simplicity and speed. The disadvantage of PCR-based methods is the possibility of undesired mutations due to the error-prone nature of certain thermostable DNA polymerases. Today, there are dozens of variations for PCR-based site-directed mutagenesis. Several representative approaches are described here.

##### 4.1. Overlap-extension Method

The *overlap-extension* method requires four oligonucleotide primers and three separate amplification reactions ([34](#), [35](#)). Two complementary mutagenic primers induce the mutation into the desired sequence of DNA, and two flanking primers amplify the mutant fragment and facilitate cloning of the PCR fragment into a suitable vector. By this approach (Fig. [3](#)), a variety of mutations can be created, such as single base-pair changes, deletions, and insertions. First, two separate PCR reactions are set up in parallel. One reaction has the “sense” mutant primer and an “anti-sense” flanking primer 3' to the mutation site. The other reaction contains the anti-sense mutant primer and the sense flanking primer 5' to the mutation site. The two amplified fragments contain mutations at the 5' or 3' terminus, respectively. In the second round of amplification, the two fragments from the first round of PCR are purified, then used as templates for amplification using only the flanking primers. After amplification, the mutation is contained within the target DNA segment, which is cloned into appropriate vectors for DNA sequencing and subsequent functional studies.

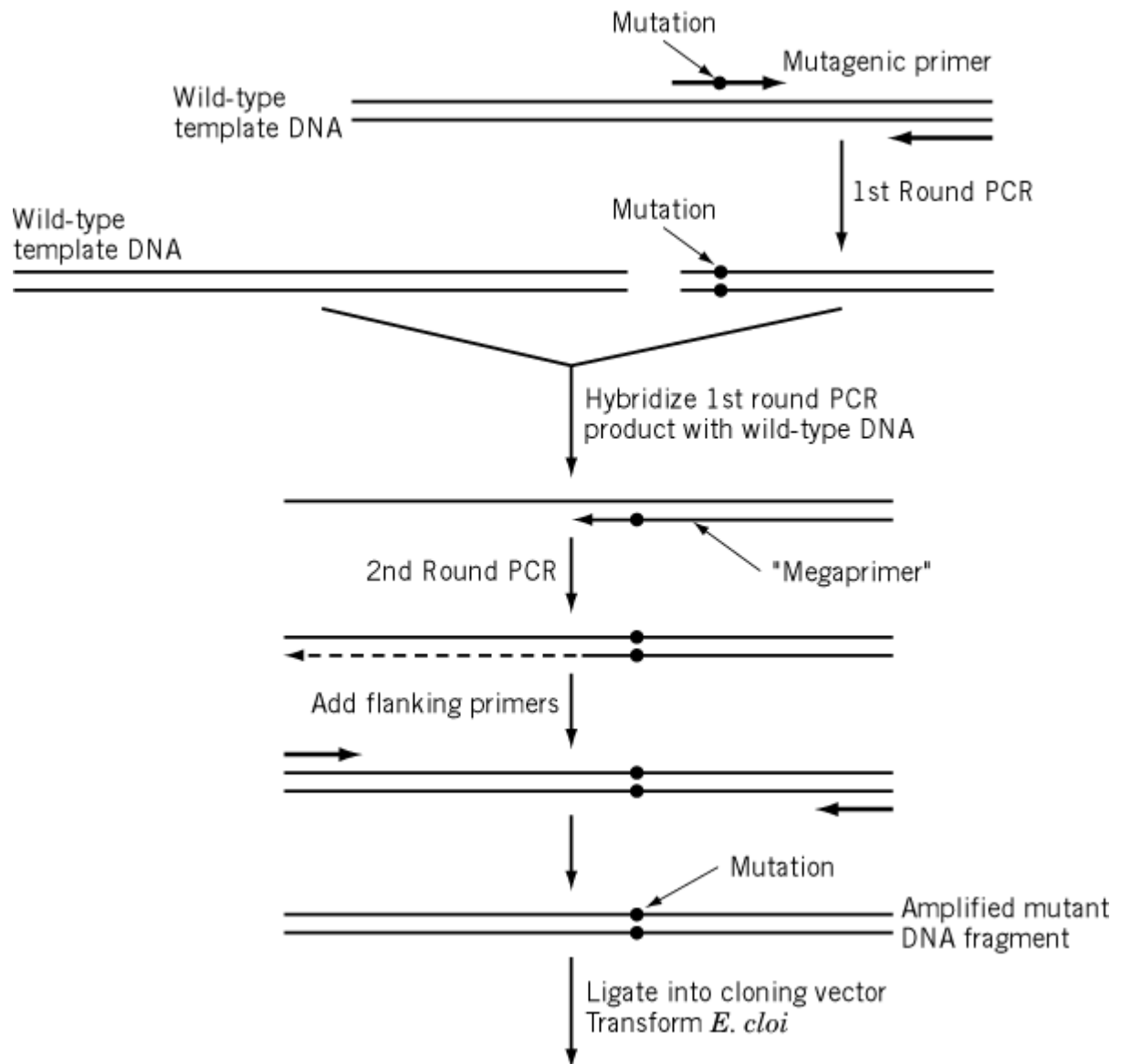
**Figure 3.** Site-directed mutagenesis by overlap-extension PCR. The two first rounds of PCR produced two overlapping fragments of the original template, both containing the mutation within the overlap region. These two PCR products are annealed and then subjected to a second round of PCR to generate the entire fragment with the mutation. The flanking primers contain the restriction sites for ligating the fragment back into the original vector.



#### 4.2. Megaprimer Method

The *megaprimer* method uses only a single mutagenic primer to create mutations in the target template (36-38) (Fig. 4). In the first round of amplification, the wild-type template is amplified using either a sense or anti-sense mutagenic primer and an appropriate flanking primer. The amplified product is then used in a second round of PCR with wild-type template and the other flanking primer to create a fragment of the same length as the original target DNA containing the desired mutation. The key to this method is that the amplified product from the first round of PCR is used as a primer in the second round of PCR. Compared to the four-primer method, this procedure requires only a single mutagenic primer and yields more of the full-length product. This is probably due to the instability of the 10- to 20-basepair overlap between the two mutant templates during the second round of amplification when using the four-primer method. In the megaprimer method, the overlap between the template and mutagenic strands is more extensive.

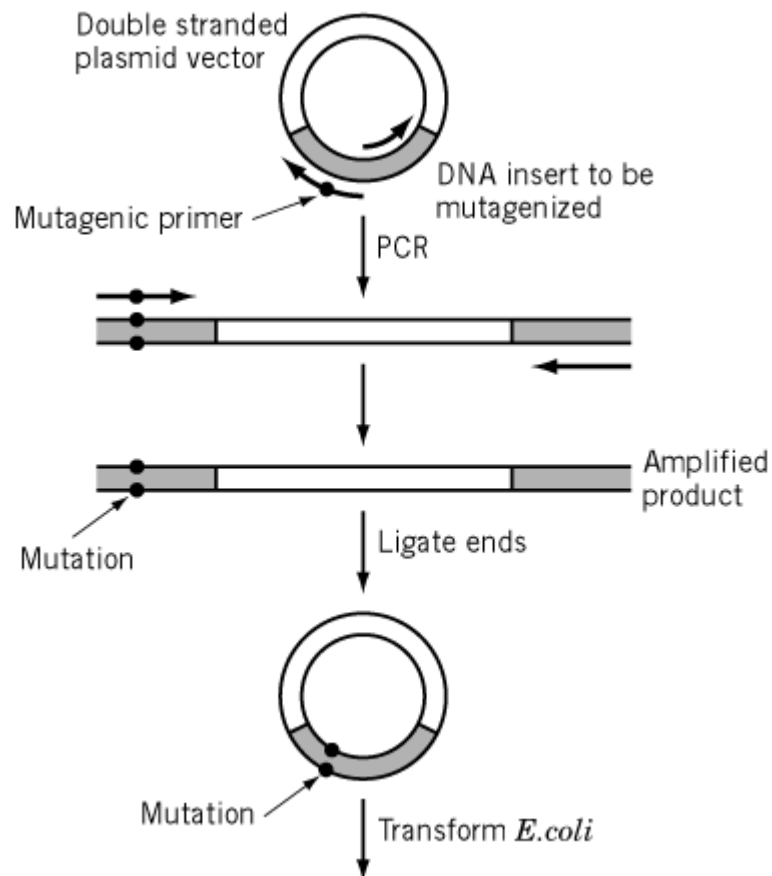
**Figure 4.** Site-directed mutagenesis by the megaprimer PCR method. The first round of PCR is used to make a fragment of the template DNA containing the desired mutation. This “megaprimer” is then hybridized to the wild-type template DNA, and a second round of PCR is carried out, to generate the entire molecule with the mutation. The flanking primers contain the restriction sites for cloning the fragment back into the vector.



### 4.3. Inverse PCR

A third approach, termed *inverse PCR*, uses only two primers to create the desired mutation (39) (Fig. 5). The key feature of this method is that in making the mutation, the entire vector is amplified. The two primers, one containing the desired mutation, extend on the circular template DNA in opposite directions. Amplification ultimately yields a linear, double-stranded DNA molecule containing the mutation at one end. Following amplification, the ends are ligated, and the resulting circular DNA molecule is transformed into *E. coli*. There are a number of variations of this method that improve the efficiency of mutagenesis (reviewed in 33).

**Figure 5.** Site-directed mutagenesis by inverse PCR. The vector is cleaved by a restriction enzyme at a single site, close to the site to be mutated. Two primers from the two ends of the linearized molecule, one containing the desired mutation, are used to amplify the linear molecule and produce molecules with the mutation. After ligating the ends to regenerate a circular molecule, the vector is used to transform *E. coli*.



The major disadvantage of PCR-based *in vitro* mutagenesis methods is the lower fidelity of [Taq DNA polymerase](#) in DNA polymerization than the phage polymerases used in the traditional methods. This can result in undesired mutations being incorporated into the amplified DNA fragment due to nucleotide misincorporation. If this happens in an early round of amplification, the undesired mutation will be present in a majority of cloned fragments. One simple way to alleviate this problem is to use fewer rounds of DNA amplification and conditions for the highest fidelity of DNA polymerization (40). Alternatively, the use of Pfu or Vent thermal stable polymerases, which are less error-prone than Taq polymerase, will result in fewer undesired mutations. Regardless of which method is used, the DNA sequence of the entire amplified fragment should be determined to ensure that only the desired changes have been made.

## 5. Evolution of Site-Directed Mutagenesis Applications

Initially, site-directed mutations were constructed one at a time, using a single oligonucleotide. However, multiple mutations can be generated using a *degenerate* oligonucleotide containing a mixture of related sequences. For example, an oligonucleotide that was degenerate at three adjacent positions was used to generate a series of proteins containing different amino acids at a selected position (11). *Doped* oligonucleotides are synthesized by spiking each nucleotide monomer with a small amount of the other three monomers. The resulting oligonucleotide will direct a variety of mutations over a defined window of the wild-type sequence (41, 42). The use of degenerate oligonucleotides is particularly powerful when used in conjunction with a genetic screen or selection (43, 44). A particularly powerful method for protein evolution, termed *DNA shuffling*, uses the polymerase chain reaction to generate proteins with novel or improved properties (45, 46).

## 6. Summary

In the 20 years since its original report, site-directed mutagenesis has evolved into a standard tool

with which to evaluate gene function. It can be accomplished by a variety of methods, many of which are now commercially available in kit form. Given the ease and speed by which site-directed mutagenesis experiments can be performed, the challenge has shifted from simply being able to construct a mutant to designing good experiments.

## Bibliography

1. C. Montell, E. F. Fisher, M. H. Caruthers, and A. J. Berk (1982) *Nature* **295**, 380–384.
2. B. C. Cunningham and J. A. Wells (1989) *Science* **244**, 1081–1085.
3. G. Winter et al. (1982) *Nature* **299**, 756–758.
4. C. A. Hutchison et al. (1978) *J. Biol. Chem.* **253**, 6551–6560.
5. The Royal Swedish Academy of Sciences (1993). The Nobel Foundation, press release.
6. K. Itakura et al. (1977) *Science* **198**, 1056–1063.
7. M. Goulian and A. Kornberg (1967) *Proc. Natl. Acad. Sci. USA* **58**, 1723–1730.
8. C. A. Hutchison and M. H. Edgell (1971) *J. Virol.* **8**, 181–189.
9. A. Razin, T. Hirose, K. Itakura, and A. D. Riggs (1978) *Proc. Natl. Acad. Sci. USA* **75**, 4268–4270.
10. M. D. Matteucci and H. L. Heyneker (1983) *Nucl. Acids Res.* **11**, 3113–3121.
11. J. A. Wells, M. Vasser, and D. B. Powers (1985) *Gene* **34**, 315–323.
12. R. P. Wharton and M. Ptashne (1985) *Nature* **316**, 601–605.
13. R. B. Wallace et al. (1980) *Science* **209**, 1396–1400.
14. I. Kudo, M. Leineweber, and U. L. Raj Bhandary (1981) *Proc. Natl. Acad. Sci. USA* **78**, 4753–4757.
15. G. F. Simons et al. (1982) *Nucl. Acids Res.* **10**, 821–832.
16. C. Traboni, R. Cortese, G. Ciliberto, and G. Cesareni (1983) *Nucl. Acids Res.* **11**, 4229–4239.
17. M. J. Zoller and M. Smith (1982) *Nucl. Acids Res.* **10**, 6487–6500.
18. W. Kramer, K. Schughart, and H. J. Fritz (1982) *Nucl. Acids Res.* **10**, 6475–6485.
19. W. Kramer et al. (1984) *Nucl. Acids Res.* **12**, 9441–9456.
20. M. Schold, A. Colombero, A. A. Reyes, and R. B. Wallace (1984) *DNA* **3**, 469–477.
21. R. J. Zagursky and M. L. Berman (1984) *Gene* **27**, 183–191.
22. D. A. Mead, E. S. Skorupa, and B. Kemper (1985) *Nucl. Acids Res.* **13**, 1103–1118.
23. M. J. Zoller, and M. Smith (1983) *Methods Enzymol.* **100**, 468–500.
24. R. B. Wallace et al. (1981) *Nucl. Acids Res.* **9**, 3647–3656.
25. K. Norris, F. Norris, L. Christiansen, and N. Fiil (1983) *Nucl. Acids Res.* **11**, 5103–5112.
26. M. J. Zoller, and M. Smith (1984) *DNA* **3**, 479–488.
27. T. A. Kunkel (1985) *Proc. Natl. Acad. Sci. USA* **82**, 488–492.
28. J. W. Taylor, J. Ott, and F. Eckstein (1985) *Nucl. Acids Res.* **13**, 8765–8785.
29. P. Carter, H. Bedouelle, and G. Winter (1985) *Nucl. Acids Res.* **13**, 4431–4443.
30. M. M. Waye, M. E. Verhoeyen, P. T. Jones, and G. Winter (1985) *Nucl. Acids Res.* **13**, 8561–8571.
31. T. Hultman et al. (1990) *Nucl. Acids Res.* **18**, 5107–5112.
32. K. Mullis et al. (1986) *Cold Spring Harbor Symp. Quant. Biol.* **51**, 263–273.
33. M. M. Ling and B. H. Robinson (1997) *Anal. Biochem.* **254**, 157–178.
34. R. Higuchi, B. Krummel, and R. K. Saiki (1988) *Nucl. Acids Res.* **16**, 7351–7367.
35. S. N. Ho et al. (1989) *Gene* **77**, 51–59.
36. R. M. Nelson, and G. L. Long (1989) *Anal. Biochem.* **180**, 147–151.
37. G. Sarkar and S. S. Sommer (1990) *Biotechniques* **8**, 404–407.



38. S. Perrin, and G. Gilliland (1990) Nucl. Acids Res. **18**, 7433–7438.
39. A. Hemsley et al. (1989) Nucl. Acids Res. **17**, 6545–6551.
40. K. A. Eckert and T. A. Kunkel (1991) PCR Methods Appl. **1**, 17–24.
41. J. B. McNeil and M. Smith (1985) Mol. Cell Biol. **5**, 3545–3551.
42. C. A. Hutchison, S. K. Nordeen, K. Vogt, and M. H. Edgell (1986) Proc. Natl. Acad. Sci. USA **83**, 710–714.
43. J. F. Reidhaar-Olson and R. T. Sauer (1991) Proteins **7**, 306–316.
44. H. B. Lowman, S. H. Bass, N. Simpson, and J. A. Wells (1991) Biochemistry **30**, 10832–10838.
45. W. P. Stemmer (1994) Nature **370**, 389–391.
46. A. Cramer, S. A. Raillard, E. Bermudez, and W. P. Stemmer (1998) Nature **391**, 288–291.

### Suggestions for Further Reading

47. D. Botstein and D. Shortle (1985) Strategies and applications of *in vitro* mutagenesis. Science **229**, 1193–1201.
48. M. M. Ling and B. H. Robinson (1997) Approaches to DNA mutagenesis: An overview. Anal. Biochem. **254**, 157–178.
49. M. Smith (1985) *In vitro* mutagenesis. Ann. Rev. Genet. **19**, 423–462.
50. M. Smith (1994) Nobel lecture. Synthetic DNA and biology. Biosci. Rep. **14**, 51–66.
51. J. D. Watson, M. Gilman, J. Witkowski, and M. Zoller (1992) *In vitro* mutagenesis, *Recombinant DNA*, 2nd ed., W. H. Freeman, New York, pp. 191–211.

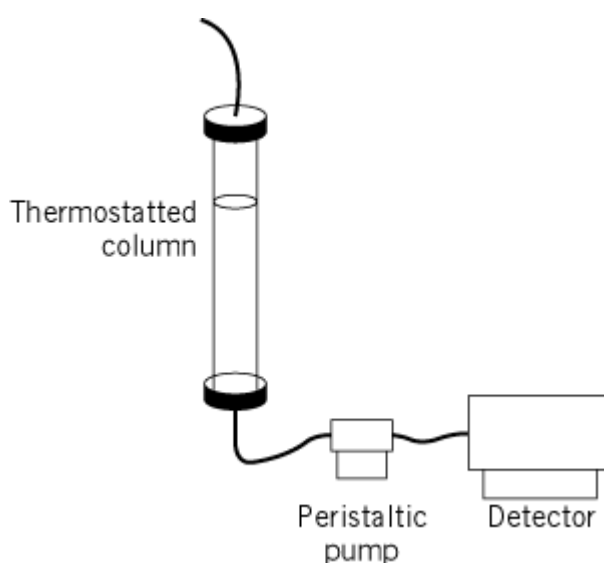
## Size Exclusion Chromatography

Molecular biologists have recently become increasingly interested in purification and biochemical analysis of **gene** products that function in a broad range of biological systems. Size exclusion chromatography (SEC) provides a method for both analytical and preparative-scale separations of biological **macromolecules**. Historically, SEC has been referred to as gel filtration chromatography when applied to the separation of **proteins** in aqueous (native) conditions (1). The most common uses of SEC are for determination of the size of macromolecules and for their purification. A less appreciated, but very powerful, application of the method is for the analysis of the interactions of macromolecules, to determine their stoichiometries and energetics. In SEC, solutes are separated on the basis of their size. In the case of proteins, the molecular size or **hydrodynamic volume** is usually specified in terms of the **Stokes Radius**(1). SEC differs from other types of **chromatography**, such as **ion-exchange** or **affinity chromatography**, in that the analytes are not separated on the basis of their differential interaction with the column resin. In fact, an inherent assumption of gel filtration chromatography is that there are no interactions between the protein solute and the stationary phase. This entry provides descriptions of the SEC technique and of methods of analysis of SEC data.

The basic experimental apparatus for standard SEC is shown in Figure 1 (2). The size exclusion resin is packed into a column, which may be jacketed to maintain temperature *via* connection to a circulating water bath. Detailed discussion of commercially available resins is provided in **molecular sieve resins**. Protein samples to be analyzed by SEC are loaded onto the top of the resin as either a

small or large zone (see discussion below), and allowed to enter the resin at a flow rate determined by a pump connected at the column outlet. Once a sufficient amount of sample has entered the column, the sample is replaced by column buffer, and transport of the macromolecules proceeds. The eluted proteins can be detected most conveniently using an in-line detection system, most commonly a UV–visible **absorbance** detector. However, **fluorescence**, or even in-line scintillation counting of **radioactivity**, can also measure the protein eluted. Alternatively, the eluted material can be characterized after its elution and collection, using a very wide variety of techniques. Although standard low pressure chromatographic equipment has traditionally been used for SEC, the use of FPLC and **HPLC** systems has become increasingly popular for performing small-zone measurements (3). The utility of HPLC for large-zone experiments, however, has been rigorously demonstrated in only one case (4).

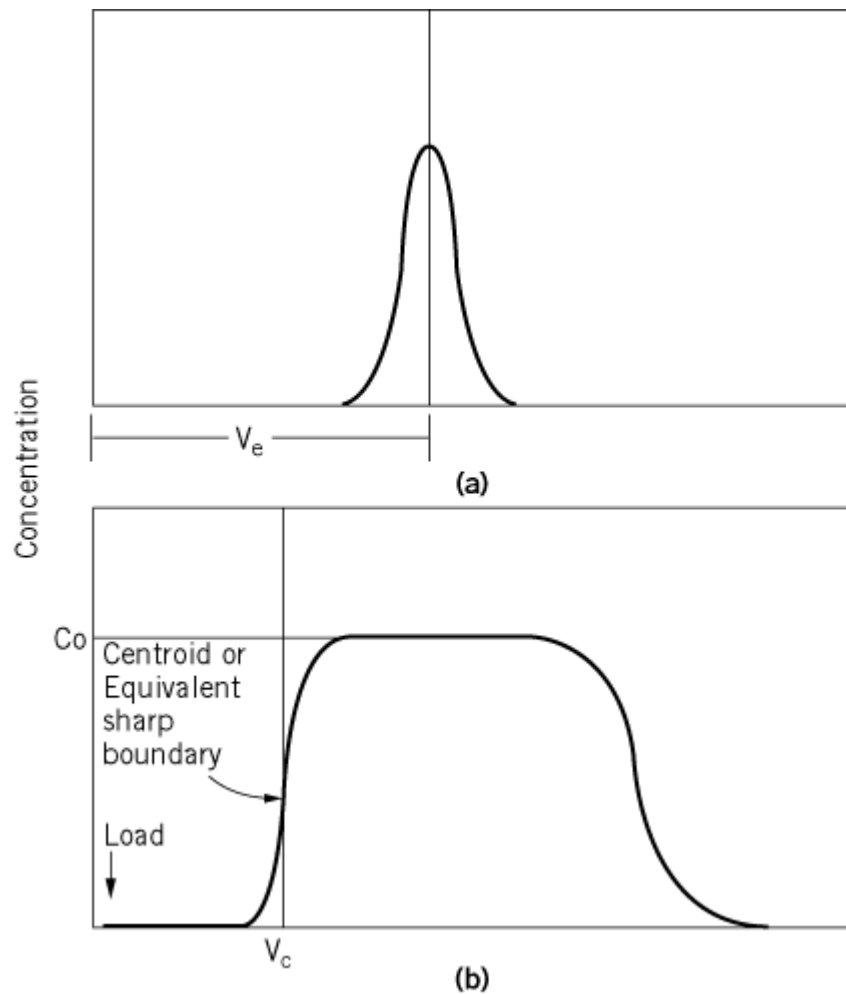
**Figure 1.** Schematic diagram of a standard low pressure gel filtration chromatography experimental setup. For a discussion of the experimental details, the reader is referred to Reference 2.



## 1. Zonal Analysis

The most common type of SEC experiment involves loading a small volume of solute onto the column. This is referred to as a small-zone experiment (Fig. 2a) or zonal analysis, and it may be used for either analytical or preparative separations. The volume of the sample used in this type of experiment is approximately 1–4% of the total column bed volume (1). As the sample flows through the column, its transport is accompanied by spreading of the zone as a result of **diffusion**. A single homogeneous peak detected at the exit port of the column should have a Gaussian shape, the apex of which defines the elution volume of the solute. Applying a mixture of solutes with different elution properties will result in multiple Gaussian-shaped peaks.

**Figure 2.** Schematic representations of results of (a) small-zone and (b) large-zone size exclusion chromatography:  $V_e$ —the elution volume of the small zone,  $C_0$ —the initial concentration of the loaded sample in the large-zone experiment,  $V_c$ —the centroid or equivalent sharp boundary of the leading edge of the large zone. This last term is analogous to the elution volume measured in the small-zone experiment.



## 2. Boundary Analysis of Large Zones

A second type of SEC measurement involves chromatographing a large volume of solute on the resin. This type of experiment, which is referred to as large-zone gel filtration chromatography, is especially useful for measuring [protein–protein interactions](#) (5, 6). One advantage of this technique over small-zone experiments is that dilution of the sample does not have to be considered. As a small zone of protein is transported through a column, the zone is diluted, so the concentration at the column outlet can be as much as 100-fold lower than that of the loaded sample. If dilution of the protein results in its dissociation, the oligomeric state of the protein at the point of elution may differ significantly from that of the loaded sample. This makes it impossible to interpret the results of small-zone experiments in terms of protein assembly equilibrium. In a large-zone experiment, in contrast, a sufficiently large volume of protein is loaded onto the column so that the edges of the zone, leading and trailing, are fed by a plateau of protein at the constant initial loading concentration (Fig. 2b). The elution volume for a large zone is determined from the centroid or equivalent sharp boundary of the leading or trailing edge of the zone. This elution position reflects the average size of the protein at its initial loading concentration. If a protein dissociates, its average size, as determined from its elution volume, will, according to mass action, decrease with zones loaded at decreasing protein concentrations. The data relating the measured partition coefficient to protein concentration are analyzed using appropriate mathematical models, to obtain the stoichiometry and energetics of the protein assembly process (6).

The information obtained from large-zone SEC experiments is thus analogous to that obtained from [sedimentation equilibrium centrifugation](#) measurements. The advantages of the SEC technique over

centrifugation are that the instrumentation is of only moderate cost, and consequently available to a larger number of laboratories, and that the method can be more readily applied to study of very tight macromolecular associations. This has been demonstrated in studies of the [lambda phage cI repressor](#), for which dimerization equilibrium constants in the nanomolar range of concentration were measured by large-zone SEC (7). An excellent discussion of the application of large-zone analytical gel filtration chromatography for measurement of protein-protein interactions is provided in Reference 6.

### 3. Size Measurements by SEC

In order to obtain size information about a solute from an SEC experiment, the column must first be characterized in terms of the volumes accessible to analytes. The volume of an SEC column is divisible into three parts: (1) the volume external to the packing material,  $V_0$ ; (2) the volume contained within the porous resin beads that is accessible to small molecules,  $V_i$ ; and (3) the volume occupied by the packing material itself,  $V_g$  (1, 5). The values of  $V_0$  and  $V_i$  are determined experimentally by measuring the elution volumes of, respectively, a large solute that is totally excluded from the interior of the resin and a small solute that has access to all pores of the resin. The elution volume,  $V_e$ , of a solute that is partially included in the pores of the gel filtration resin can be related to the void and internal volumes of a column by the following equation:

$$V_e = V_0 + \sigma V_i$$

where  $s$  is the [partition coefficient](#) of the solute or the fraction of the internal volume of the resin,  $V_i$ , that is accessible to the analyte. A detailed discussion that relates these measured volumes to the microscopic properties of the sieving resin is contained in [molecular sieve resins](#).

It is the partition coefficient,  $s$ , that, when compared with the values measured for proteins of known size, provides information about the molecular size of an unknown macromolecule. If a series of proteins of known size or Stokes radius are subjected to SEC, a linear relationship between partition coefficient and size is observed (5). For proteins that are well approximated as hydrated spheres, a linear relationship between the elution volume and molecular weight will be observed (1, 5). One approach to analysis of SEC data to obtain information about the size of the protein was developed by Ackers (8). This relationship assumes that penetrable volumes within column packing material are distributed randomly with respect to the size of the protein that can penetrate them. The Stokes radius,  $a$ , for a protein in Å units ( $1.0 \times 10^{-10}$ ) is related to the inverse error function complement ( $\text{erfc}^{-1}$ ) of  $s$  as

$$a = a_0 + b_0 \text{erfc}^{-1} \sigma$$

where  $a_0$  and  $b_0$  are calibration constants for a column and are obtained by chromatographing a series of proteins of known size or Stokes radius. Values of the inverse error function complement are compiled in tables provided by the National Institute of Standards and Technology (NIST), or they may be readily calculated (9). The Stokes radius for a protein obtained from SEC measurements will accurately reflect the radius of a spherically shaped molecule, but not all proteins are well modeled as spheres; approaches to deal with other shapes are given in [Stokes Radius](#).

### Bibliography

1. M. E. Himmel and P. G. Squire (1988) In *Aqueous Size Exclusion Chromatography*, P. L. Dubin, ed., Elsevier, New York, pp. 3–22.
2. *Gel Filtration: Theory and Practice*, 6th ed., Amersham Pharmacia Biotech, Piscataway, NJ.
3. B. Seville, C. Vidal-Madjar, and A. Jaulmes (1991) In *HPLC of Proteins, Peptides and*

*Polynucleotides*, M. T. W. Hearn, ed., VCH, New York, pp. 397–451.

4. E. Nenortas and D. Beckett (1995) *Anal. Biochem.* **222**, 366–373.
5. G. K. Ackers (1975) In *The Proteins*, H. Neurath and R. L. Hill, eds., Academic, New York, pp. 1–94.
6. R. Valdes, Jr. and G. K. Ackers (1979) *Meth. Enzymol.* **61**, 125–143.
7. D. Beckett, K. S. Koblan, and G. K. Ackers (1991) *Anal. Biochem.* **196**, 69–75.
8. G. K. Ackers (1967) *J. Biol. Chem.*, **242**, 3237–3238.
9. A. V. Aston (1954) *Tables of the Error Function and Its Derivative*, N.B.S., Applied Mathematics Series, Vol. **41**, pp. 3829–3838.

### Suggestion for Further Reading

10. H. G. Barth, B. E. Boyes, and C. Jackson (1996) *Anal. Chem.* **68**, 445R–466R. (A comprehensive review of recent advances in SEC methodologies, including recent developments in application of SEC to the study of biological macromolecules.)

## Slime Molds

*Slime molds* is a term commonly used to describe a group of soil amoebas that have the ability to use bacteria as a source of food. They are neither a mold nor a fungus, but rather organisms that share many basic biochemical and genetic mechanisms with higher eukaryotic cells. They can be readily cultivated in the laboratory, and a considerable array of techniques have been developed to manipulate them on a biochemical, genetic, and cellular level. To the molecular biologist, these organisms can present opportunities to examine the details of the regulation and function of systems common to higher organisms, but in a less complicated and more accessible setting.

The two species most studied in the laboratory are *Dictyostelium discoideum* (see [Dictyostelium](#)) and *Physarum polycephalum* (see [Physarum](#)). They both have the ability to change their basic cellular architecture in a programmed, developmental pathway to optimize their growth and survival. They differ in that *Physarum* has a number of vegetative forms, including a flagellated swimming form, a migrating amoeboid form, and a unique single large cell form containing many nuclei, called a plasmodium. When starved, it can develop into a variety of environmentally resistant forms, including spores. The *Dictyostelium* developmental program differs from that of *Physarum* in that the cells aggregate when starved and develop into a multicellular organism, called a slug. The cell walls do not fuse, and cells differentiate into unique types that perform different functions. If this faster moving slug fails to find better conditions, it can develop into a terminal fruiting structure containing resistant spores. The vegetative form of *Dictyostelium* is a single cell that uses amoeboid motility to find food.

Slime molds have been classified in the *Mycetozoa* phylum. There are three recognized subclasses: acellular (myxogastroid), cellular (dictyostelid), and protostelid. It has been proposed that slime molds are related to the multicellular eukaryotes, more closely to fungal and animal cells than to plants ([1](#)). Outside the laboratory, they can be found in forest detritus, where they occupy a unique niche, eating bacteria. In the laboratory, by combining genetic, biochemical, and cell biological techniques, these organisms have given investigators multifaceted approaches to problems that are intractable elsewhere.

## Bibliography

1. S. L. Baldaur and W. F. Doolittle (1997) Proc. Natl. Acad. Sci. USA **94**, 12007–12012.

## Suggestion for Further Reading

2. S. L. Stephenson and H. Stempen (1994) *A Handbook of Slime Molds*, Timber Press, Portland, OR.

## Slow-Binding Enzyme Inhibition

Inhibitors have proved to be useful probes of chemical and [kinetic mechanisms](#) of **enzyme**-catalyzed reactions. The action of inhibitors has also provided background information for the development of specific bioactive compounds to act as chemotherapeutic agents or as herbicides. Most studies have been performed with classical inhibitors. These are substrate analogues that give rise to linear [competitive inhibition](#) with respect to the substrate through the reversible formation of **dead-end** complexes that can only dissociate back to the components from which they were formed. Usually, the [kinetic](#) investigations have been performed under conditions in which the concentrations of substrate and inhibitor are much greater than the concentration of the enzyme, and all the equilibria are set up rapidly, conforming to [Michaelis-Menten kinetics](#). The rapid interaction of an inhibitor (I) at the [active site](#) of an enzyme (E) to form an EI complex is described by equation [1](#).



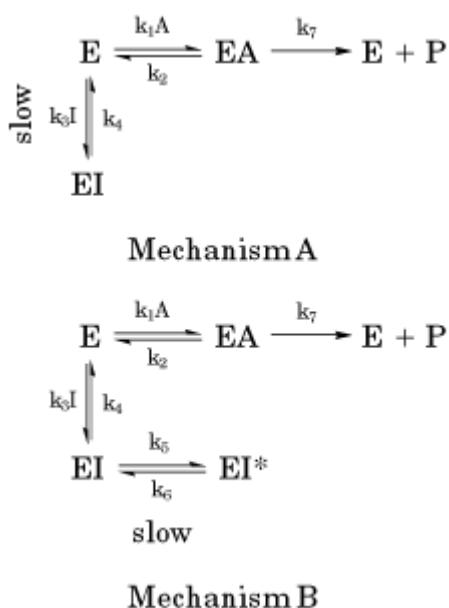
The **dissociation constant** ( $K_i$ ) for the reaction, a thermodynamic quantity, is defined by the relationships given in equation [\(2\)](#).

$$K_i = \frac{k_4}{k_3} = \frac{(E)(I)}{(EI)} \quad (2)$$

Rapid establishment of the equilibrium, on the time scale of a steady-state kinetic experiment of seconds to minutes, requires that the magnitudes of the apparent first-order rate constants for the formation of EI,  $k_3I_p$ , and for the dissociation of EI,  $k_4$ , be relatively high, with values greater than about  $1\text{s}^{-1}$ . If the binding is very tight, however, dissociation will be intrinsically slow; for example, if  $K_i$  is in the micromolar range,  $k_4$  can be no greater than  $10^{-6} k_3$  in units of  $\text{s}^{-1}$ . Consequently, a tight-binding inhibitor will tend to dissociate slowly and might not fulfill the requirements for rapid equilibration.

It has now been established that the equilibrium between enzyme, inhibitor, and the enzyme-inhibitor complex is not always set up rapidly, so that the inhibition becomes time-dependent. Compounds that behave in this manner have been referred to as *slow-binding inhibitors* [\(1\)](#). This term conveys the idea that binding, which is the establishment of all equilibria involving the inhibitor, occurs slowly on the steady-state time scale of seconds to minutes. Slow-binding inhibition resembles transient-state kinetics on a different time scale. The definition is an operational one, as it is always possible that the time to reach equilibrium could vary with the experimental conditions. The simplest scheme to illustrate slow-binding inhibition of an enzyme-catalyzed reaction is shown as mechanism A in Figure [1](#).

**Figure 1.** Kinetic mechanisms for the slow-binding inhibition of enzyme-catalyzed reactions by substrate analogues.



There are two reasons why formation of an EI complex might be slow. If the concentration of  $I$  required to demonstrate the inhibition is relatively low, the rate of forming EI ( $k_3 I_t$ ) would also be low. For example, if the values for  $K_i$ ,  $I_t$  and  $k_3$  are  $10^{-9}$  M,  $10^{-8}$  M and  $10^6 \text{ M}^{-1} \text{ s}^{-1}$ ,  $k_3 I_t$  would be  $10^{-2} \text{ s}^{-1}$  and the [half-time](#) for the forward reaction 69 s. As  $k_4 = K_i k_3 = 10^{-3} \text{ s}^{-1}$ , the half-time for the dissociation reaction will be 11.6 min. It should be noted that, for these calculations, it is assumed that the inhibitor concentration is at least 10 times that of the enzyme, so that tight-binding inhibition does not occur (see [Active-Site Titrants](#)). Analysis of slow-binding data for the mechanism under discussion will yield a true value for  $k_3$ , whose magnitude is comparable to that of the usual second-order rate constants for the interaction of enzymes and substrates. An alternative explanation for slow-binding inhibition conforming to mechanism A is that the inhibitor encounters barriers to its binding at the active site. The value of  $k_3$  might then be sufficiently low to permit the monitoring, on a steady-state time scale, of the time-dependent inhibition.

A second mechanism to account for slow-binding inhibition is mechanism B in [Figure 1](#), which involves the rapid formation of a collision complex (EI), with the inhibitor behaving initially as a classical competitive inhibitor (cf eqs. [1](#) and [2](#)). The EI complex then undergoes a slow conformational change, or isomerization reaction, to form a more stable complex, EI\*. The overall dissociation constant for the reaction,  $K_i^*$ , would be defined as

$$K_i^* = \frac{(E)(I)}{(EI) + (EI^*)} = \frac{K_i k_6}{k_5 + k_6} \quad (3)$$

where  $K_i = k_4/k_3$ . The degree to which the initial binding is enhanced through the isomerization reaction will depend on the ratio  $k_5/k_6$ . The absolute values of these rate constants must be such as to allow observation, on a steady-state time scale, of the isomerization as manifested by the slow increase in inhibition.

Under conditions in which the total concentration of a slow-binding inhibitor is at least 10 times greater than the total enzyme concentration and the reaction is started by the addition of enzyme, the progress curve in the presence of a single inhibitor concentration is described for either mechanism A or mechanism B by an integrated rate equation (eq. 4).

$$P = v_s t + (v_o - v_s)(1 - e^{-kt})/k \quad (4)$$

For this equation, which contains both a linear term and an exponential term,  $P$  represents the concentration of product;  $v_s$  and  $v_o$  denote steady-state and initial velocities, respectively;  $k$  represents an apparent first-order rate constant whose meaning varies with the mechanism. Equation 4 predicts, and curve b of Figure 2 illustrates, that when the reaction is started by adding enzyme, there is an initial burst or transient phase. The velocity then settles down to a slower steady-state rate, as represented by the asymptote of the curve, because of the slow establishment of the equilibrium between enzyme, inhibitor, and the enzyme-inhibitor complex(es). If the enzyme is preincubated with the inhibitor and the reaction is started by the addition of substrate, equation 4 still applies. But there is now a slow decrease in the degree of inhibition, an apparent activation, because of the establishment of new equilibria that involve interactions of both substrate and inhibitor with the enzyme (Fig. 2, curve c); consequently, the initial velocity is less than the steady-state velocity. Whichever procedure is used to study the inhibition, the same steady-state velocity will be obtained. All this discussion assumes that the enzyme is not inactivated or activated by other phenomena.

**Figure 2.** Progress curves for an enzyme-catalyzed reaction in (a) the absence, and (b, c) the presence, of a slow-binding inhibitor. The reactions were started (b) by adding enzyme or (c) by adding substrate after preincubation of enzyme and inhibitor. The dashed lines represent steady-state rates in the presence of inhibitor.

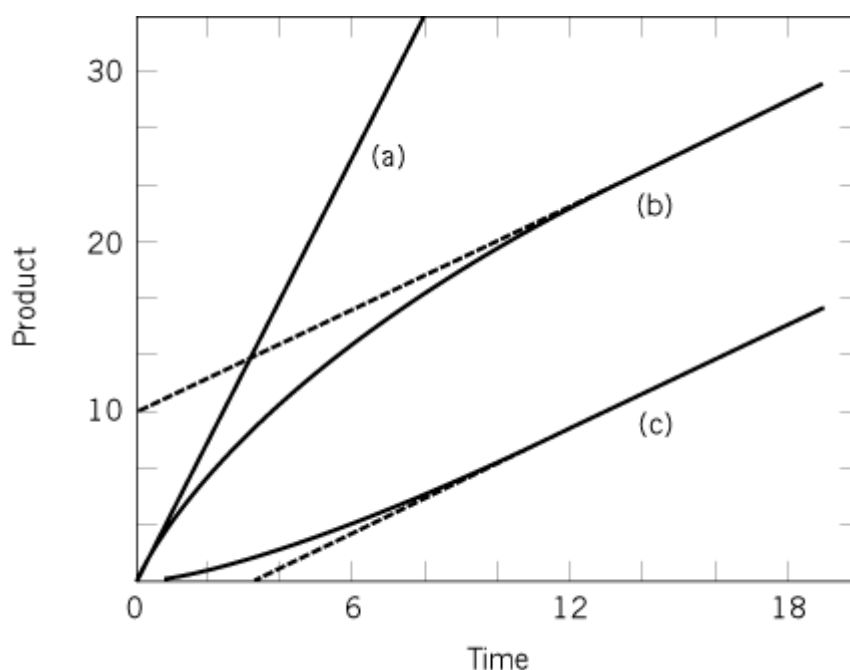
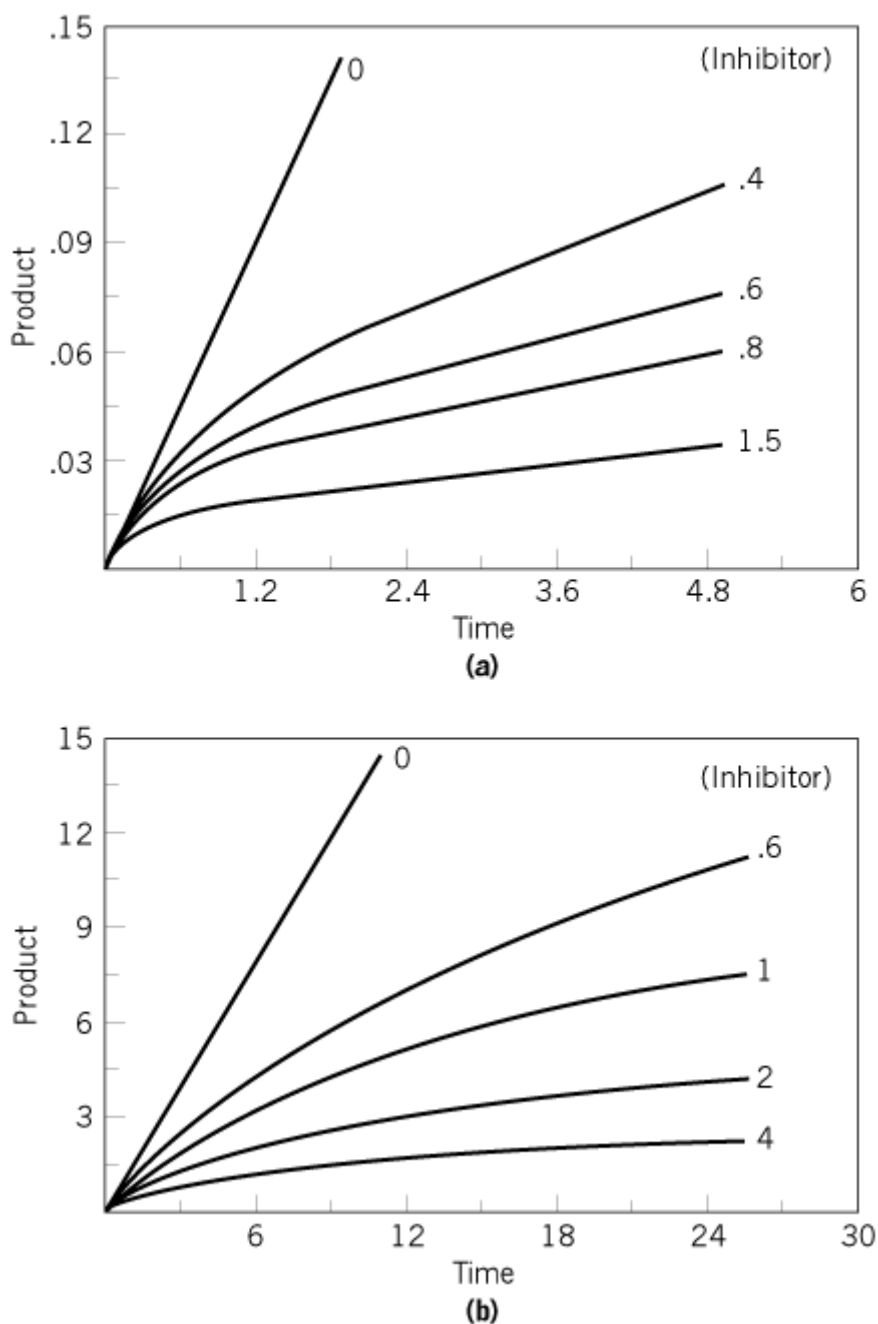


Figure 3 shows the basic types of plots obtained when reactions conforming to either mechanism A or mechanism B are run in the presence of different concentrations of a slow-binding inhibitor, and started by the addition of enzyme. The two families of curves are similar except that, for mechanism A, the initial velocity of the reaction is independent of the inhibitor concentration. By contrast, the initial velocity of a reaction that conforms to mechanism B is described by the equation for linear [competitive inhibition](#) (Table 1). This same general equation describes the variation with inhibitor



concentration of the asymptotes for each family of curves (Table 1). As the concentration of the inhibitor increases, so also does the rate at which the curves for the two mechanisms turn over. This is a function of the value for the apparent first-order rate constant ( $k$ ), which increases with inhibitor concentration (Table 1). It will be noted that, in the presence of a fixed concentration of substrate ( $A$ ), the variation of  $k$  as a function of inhibitor concentration is linear for mechanism A and hyperbolic for mechanism B. The linear plot has an intercept with the vertical ordinate equal to  $k_4$ . The hyperbolic plot has limiting values of  $(k_5 + k_6)$  and  $k_6$  at infinite and zero inhibitor concentrations, respectively.

**Figure 3.** Progress curves for the slow-binding inhibition of enzyme-catalyzed reactions conforming to mechanisms (a) and (b). In each case, the reaction was starting by adding enzyme to a mixture of substrate and varying amounts of the inhibitor.



| Parameter in equation<br>(4) | Mechanism A   | Mechanism B  |
|------------------------------|---|--|
| Initial velocity             | $v_o$   | $v_o = \frac{VA}{K_o(1+I/K_i)+A}$  |
| Steady-state velocity        | $v_s = \frac{vA}{K_o(1+I/K_i)+A}$                         | $v_s = \frac{VA}{K_o(1+I/K_i^*)+A}$  |
| $k$                          | $k = k_4[1 + \frac{i}{K_i(1+A/K_o)}]$<br>$k_4 = kv_s/v_o$ | $k = k_6 \left\{ \frac{1+I/[K_i^*(1+A/K_o)]}{1+I/[K_i(1+A/K_o)]} \right\}$<br>$k_6 = kv_s/v_o$ |

<sup>a</sup>  $K_i = k_4/k_3K_i^* = K_i[k_6(k_5+k_6)]$  (see Fig. 1)

Although values for the kinetic parameters associated with slow-binding inhibition can be determined by graphic procedures, it is greatly preferable to make an overall least-squares fit of the data to equations that describe mechanisms A and B, namely, equation 4, plus those listed in Table 1. This approach allows discrimination between the two mechanisms, as well as determination of values for each kinetic parameter, together with their standard errors. The nonlinear regression analysis involves using about 20 points from each curve and the utilization of a graphic procedure to obtain estimates of the  $k$  value for each curve (2, 3). The same basic approach is used for the study and analysis of slow-binding inhibition when an enzyme has multiple substrates. Analysis of inhibition data is simplified, however, if all substrates other than the one for which the inhibitor is an analogue are present at saturating concentrations.

There are several points to note about the detection and study of slow-binding inhibition (2). It will not be observed, on starting the reaction by adding the enzyme, if the rate of interaction of enzyme and inhibitor (mechanism A) or the rate of isomerization of the enzyme-inhibitor complex (mechanism B) are very slow relative to the rate of the reaction catalyzed. Therefore, the enzyme should also be preincubated with the inhibitor, and inhibition should be studied as a function of enzyme concentration. The ratio of the concentration of substrate to its  $K_m$  (**Michaelis constant**) should not be high; if it is, the concentration of free enzyme is reduced considerably and there is a concomitant reduction in the collision rate for enzyme and inhibitor. Thus, it could appear that the inhibition conforms to mechanism A rather than to mechanism B. With any progress curve investigation, it is important that there is no other inactivation of the enzyme. This can be confirmed by the observation of an extended linear asymptote or by the application of Selwyn's test (4). Reactions in the presence of the inhibitor should be run for a sufficiently long period, but not beyond the point at which the reaction in the absence of inhibitor ceases to be linear.

In the original formulation of mechanism B, consideration was not given to the possibility that the step for the formation of an initial collision complex between E and I was at steady state rather than at thermodynamic equilibrium (Fig. 1). Subsequent theoretical investigations have demonstrated the difficulties associated with the use of steady-state kinetic techniques to distinguish mechanism B from the more general but more complex mechanism B, with both steps involving the formation of enzyme-inhibitor complexes being at steady state equilibrium (5).

When a slow-binding inhibitor causes inhibition at concentrations comparable to that of the enzyme, the inhibitor is described as being of the *slow, tight-binding* type (1). The inhibition can still be described in terms of mechanisms A and B (Fig. 1), but the equations to describe the inhibition are more complex because of the change in the concentration of free inhibitor as the inhibition proceeds. The availability of a computer program that allows straightforward analysis of

slow, tight-binding inhibition data has facilitated the study of this type of inhibition (3). Consideration has been given recently to kinetic schemes for slow- and tight-binding inhibition where the inhibitor can also combine with an enzyme-substrate complex and influence the binding or the rate of product formation, or both (6) (see [Competitive Inhibition](#) and [Noncompetitive Inhibition](#) ).

From a list of enzymes subject to slow-binding and slow, tight-binding inhibition (1), it appears that the predominant mechanism of inhibition is that described by mechanism B.

#### Bibliography

1. J. F. Morrison and C. T. Walsh (1988) *Adv. Enzymol.* **59**, 201–301.
2. J. F. Morrison and S. R. Stone (1985) *Comments Mol. Cell. Biophys.* **2**, 347–368.
3. M. J. Sculley and J. F. Morrison (1986) *Biochim. Biophys. Acta* **874**, 44–53.
4. M. J. Selwyn (1965) *Biochim. Biophys. Acta* **105**, 193–195.
5. M. J. Sculley, J. F. Morrison, and W. W. Cleland (1996) *Biochim. Biophys. Acta* **1298**, 78–86.
6. S. E. Szedlacsek and R. G. Duggleby (1995) *Meth. Enzymol.* **249**, 144–180.

### Small Cytoplasmic RNA (scRNA), Small Cytoplasmic RNP (scRNP)

The nucleoplasm and cytoplasm of eukaryotic cells contain a large number of small ribonucleoprotein (RNP) complexes, which are called [small nuclear RNPs](#) (snRNPs) and small cytoplasmic RNPs (scRNPs), respectively. They are composed of one or more proteins and a small stable RNA molecule (chain length <~300 nucleotides). The high abundance and the evolutionary conservation of RNA sequences within these small RNPs indicate that they have an important biological role. While a large body of information is available on the function of snRNPs in regulating nuclear RNA processing (see [RNA Splicing](#)), relatively little is known about the function of most scRNPs. Nonetheless, their cellular localization suggests a role in regulation of the last step of gene expression—that is, [translation](#) or disposition of newly synthesized proteins [reviewed by Baserga and Steitz (1)].

The discovery of the scRNPs was facilitated by **autoantibodies** present in sera of patients with [autoimmune diseases](#) in which antibodies against scRNPs components were found. The antibodies are usually targeted against a protein [epitope](#), indicating that a common protein is associated with scRNPs of the same class. For example, **antibodies** against the 60-kDa Ro protein, which is a component of several scRNPs, also called Ro RNPs, were found in sera of patients with systemic lupus erythematosus (SLE) or with Sjögren's syndrome (2-6). Sera of such patients frequently contain autoantibodies against the 50-kDa La protein (2, 3). These antigens were termed Ro and La, based on the names of the patients in whom they were first identified (3).

#### 1. Ro RNPs

The Ro RNPs were first identified with the aid of anti-Ro autoantibodies (6). They contain one molecule of small cytoplasmic RNA that belongs to a family called Y RNAs, to indicate their cYtoplasmic location, and are assembled with the 60-kDa Ro protein and the 50-kDa La protein. The Y RNAs sequences, as well as the sequences of the Ro and La proteins, have been conserved during

[evolution](#). This, together with their high abundance ( $\sim 10^5$  copies per cell), suggests an important physiological role. However, their biological function is yet unknown [reviewed by Baserga and Steitz (1) and by Van Venrooij et al. (7)].

The Y RNAs are produced by [transcription](#) by **RNA polymerase III**. Four different Y RNA molecules, ranging in size between 69 to 112 nucleotides, were identified in humans and were termed Y1, Y3, Y4, and Y5 RNAs (Y2 is a truncated form of Y1). In different species, however, different numbers of types of Y RNAs were found associated in Ro RNPs. Thus, while humans and [Xenopus](#) cells contain all four types of Y RNAs, two types are found in mouse cells, and only one type is found in *Caenorhabditis elegans* (8-12).

The predicted secondary structure of sequenced Y RNAs can be drawn as a long base-paired stem composed of the 5' and the 3' ends of the molecule and an internal loop containing a pyrimidine-rich sequence (9, 13). Within the stem region, a highly conserved base-paired sequence with a bulged cytidine was proposed to be the binding site of the 60-kDa Ro protein (14, 15). The La antigen was proposed to bind to the 3' oligo-uridylylate sequence (15-17).

In *Xenopus* oocytes, the 60-kDa Ro protein is found also complexed with a heterogeneous population of 5S ribosomal RNA precursors that contain internal mutations. This binding is attributed to an alternative fold of the RNA in the mutant 5S rRNA. Because the mutant 5S rRNA precursors are processed inefficiently and are finally degraded, it was proposed that the Ro protein may function in a quality-control pathway of 5S rRNA biosynthesis (18, 19).

The 60-kDa Ro protein is an abundant and conserved protein that contains an RNA-recognition motif (RRM) (see [RNA-Binding Proteins](#)) of  $\sim 80$  amino acid residues (10, 20). An additional protein of 52 kDa (Ro52), which is recognized by autoantibodies from patients with SLE and Sjögren's syndrome, has been proposed to be a component of Ro RNPs in human cells (21), but other reports indicate otherwise (22, 23). The 52-kDa protein does not bind directly to the Y RNA, but probably associates with the Ro RNP through [protein-protein interactions](#) (24).

The 50-kDa La autoantigen is a highly abundant phosphoprotein that is found complexed with a variety of small RNAs to form small RNP particles [reviewed by Van Venrooij et al. (7)]. These include precursors to 5S rRNA, [transfer RNA](#), 7S RNA, and Ro cytoplasmic Y RNAs, which are transcribed by RNA polymerase III (8, 25, 26). The La protein has an RNA recognition motif (RRM) (27). Its binding to the above-mentioned small RNAs occur via this RRM and is directed, at least in part, to a short uridylylate sequence at the 3' end of the relevant RNA (15-17).

The La protein has been proposed to be involved in a number of physiological processes within the cell's nucleus. For example, experiments *in vitro* have indicated that the La protein is required as a transcription termination factor for RNA polymerase III transcripts (28-30), which is consistent with the protein's [ATPase](#) activity that can melt DNA/RNA hybrids *in vitro* (31). Nevertheless, its association with scRNPs suggests a cytoplasmic physiological function as well. This view is supported by **immunofluorescence** experiments showing that the La protein is localized in the cytoplasm at sites of translation (32).

## 2. Viral scRNPs

An example of viral scRNPs are complexes of VAI and VAII small [adenovirus](#) RNAs that are transcribed by RNA polymerase III and have the La autoantigen as a common protein partner. These viral scRNPs are highly abundant in infected cells ( $\sim 10^8$  copies per cell) and are found in the cytoplasm. The  $\sim 160$ -nucleotide VAI RNA plays a role in enabling **protein biosynthesis** late in the adenovirus life cycle (33). Due to its high abundance and its binding to the interferon-induced 68-kDa kinase, it competes effectively with the cellular components that are involved in the **interferon**-induced transcription shutdown of translation (34, 35).

### 3. Signal Recognition Particle (SRP)

Another class of scRNPs are the [signal recognition particles](#) (SRPs) that help translocation of secretory proteins during translation. Mammalian SRPs are composed of 7SL RNA of ~300 nucleotides in length, which is complexed with six proteins having molecular weights of 9, 14, 19, 54, 68, and 72 kDa, to form an 11 S RNP particle [reviewed by Walter and Johnson (36) and by Lutcke (37)].

The 7SL SRP RNA is an abundant transcript of RNA polymerase III (~10<sup>6</sup> copies per cell). It is a member of a class of mammalian intermediate-repeat sequences called short interspersed elements ([SINES](#)). The 7SL RNA has a remarkable sequence [homology](#) with the **Alu** family of repetitive sequences in primates and rodents. The “Alu domain” of the SRP consists of the Alu homologous region of 7SL RNA to which SRP9 and SRP14 are bound.

SRPs are required for cotranslational targeting of nascent secretory proteins to the [endoplasmic reticulum](#). Translation of secretory proteins begins on [ribosomes](#) that are free in the cytosol. The SRP binds to the **hydrophobic** amino-terminal nascent peptide (15 to 30 amino acid residues) of the secreted nascent protein while it emerges from the large subunit of the ribosome. The binding arrests or delays the translation of the nascent polypeptide chain. The SRP then targets the nascent polypeptide–ribosome complex to the SRP receptor in the rough endoplasmic reticulum. The release of SRP from the signal sequence/ribosome complex requires GTP hydrolysis [reviewed by Walter and Johnson (36) by Lutcke (37), and by Bovia and Strub (38)]. Specifically, the 54-kDa protein, also called SRP54, binds directly to the signal peptide through a methionine-rich carboxyl-terminal domain. This protein is a [GTP-binding protein](#), and when bound to the 7SL RNA it helps target the nascent chain by interacting with the heterodimeric SRP receptor. The function of the other proteins has not yet been characterized in detail. SRP19 aids the binding of SRP54 to the 7SL RNA, and the heterodimer SRP68/72 was proposed to interact with the  $\alpha$  subunit of the SRP receptor. SRP68/72 and the Alu domain of SRP are required to confer elongation arrest activity of the particles.

The components of the SRP and SRP receptor were highly conserved during evolution. This is in accordance with their function in targeting proteins to the endoplasmic reticulum, which is shared by all eukaryotic organisms examined [reviewed by Wolin (39)].

### 4. Alu-Related scRNPs

Several Alu-related small RNAs were found to accumulate stably in the cytoplasm of rodent and primate cells. They are less abundant (10<sup>3</sup> to 10<sup>4</sup> copies per cell) and less well characterized than the above described scRNPs. Among them are the primate scAlu RNA (120 nucleotides), the rodent scB1 (140 nucleotides), and the BC200 RNA (200 nucleotides). The latter is expressed specifically in nerve cells by RNA polymerase III. These scRNAs have Alu-like sequences at their 5' and 3' ends, similar to the 7SL RNA, and can be drawn in a cruciform secondary structure. These Alu-like scRNAs are found complexed in small cytoplasmic RNPs and have been shown to bind the SRP14/19 heterodimer *in vitro*. Although the role of these scRNPs is yet unknown, a role in translation has been suggested [reviewed by Bovia and Strub (38)].

### Bibliography

1. S. J. Baserga and J. A. Steitz (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 359–381.
2. J. R. Anderson, K. G. Gray, J. S. Beck, W. U. Buchanan, and A. J. McElhinney (1962) *Ann. Rheum. Dis.* **21**, 360–369.
3. G. Clark, M. Reichlin, and T. B. Tomasi, Jr. (1969) *J. Immunol.* **102**, 117–122.
4. M. Mattioli and M. Reichlin (1974) *Arthritis Rheum.* **17**, 421–429.

5. M. A. Alspaugh and E. M. Tan (1975) *J. Clin. Invest.* **55**, 1067–1073.
6. M. R. Lerner, J. A. Boyle, J. A. Hardin, and J. A. Steitz (1981) *Science* **211**, 400–402.
7. W. J. Van Venrooij, R. L. Slobbe, and G. J. Pruijn (1993) *Mol. Biol. Rep.* **18**, 113–119.
8. J. P. Hendrick, S. L. Wolin, J. Rinke, M. R. Lerner, and J. A. Steitz (1981) *Mol. Cell. Biol.* **1**, 1138–1149.
9. S. L. Wolin and J. A. Steitz (1983) *Cell* **32**, 735–744.
10. C. A. O'Brien, K. Margelot, and S. L. Wolin (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7250–7254.
11. D. J. Van Horn, D. Eisenberg, C. A. O'Brien, and S. L. Wolin (1995) *RNA* **1**, 293–303.
12. A. D. Farris, C. A. O'Brien, and J. B. Harley (1995) *Gene* **154**, 193–198.
13. C. W. Van Gelder, J. P. Thijssen, E. C. Klaassen, C. Sturchler, A. Krol, W. J. Van Venrooij, and G. J. Pruijn (1994) *Nucleic Acids Res.* **22**, 2498–2506.
14. S. L. Wolin and J. A. Steitz (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1996–2000.
15. G. J. Pruijn, R. L. Slobbe, and W. J. Van Venrooij (1991) *Nucleic Acids Res.* **19**, 5173–5180.
16. M. B. Mathews and A. M. Francoeur (1984) *Mol. Cell. Biol.* **4**, 1134–1140.
17. J. E. Stefano (1984) *Cell* **36**, 145–154.
18. C. A. O'Brien and S. L. Wolin (1994) *Genes Dev.* **8**, 2891–2903.
19. H. Shi, C. A. O'Brien, D. J. Van Horn, and S. L. Wolin (1996) *RNA* **2**, 769–784.
20. S. L. Deutscher, J. B. Harley, and J. D. Keene (1988) *Proc. Natl. Acad. Sci. USA* **85**, 9479–9483.
21. E. Ben-Chetrit, E. K. Chan, K. F. Sullivan, and E. M. Tan (1988) *J. Exp. Med.* **167**, 1560–1571.
22. A. Kelekar, M. R. Saitta, and J. D. Keene (1994) *J. Clin. Invest.* **93**, 1637–1644.
23. G. Boire, M. Gendron, N. Monast, B. Bastin, and H. A. Menard (1995) *Clin. Exp. Immunol.* **100**, 489–498.
24. R. L. Slobbe, W. Pluk, W. J. Van Venrooij, and G. J. Pruijn (1992) *J. Mol. Biol.* **227**, 361–366.
25. J. Rinke and J. A. Steitz (1982) *Cell* **29**, 149–159.
26. J. C. Chambers, M. G. Kurilla, and J. D. Keene (1983) *J. Biol. Chem.* **258**, 11438–11441.
27. G. J. Pruijn, R. L. Slobbe, and W. J. Van Venrooij (1990) *Mol. Biol. Rep.* **14**, 43–48.
28. E. Gottlieb and J. A. Steitz (1989) *EMBO J.* **8**, 851–861.
29. R. J. Maraia, D. J. Kenan, and J. D. Keene (1994) *Mol. Cell. Biol.* **14**, 2147–2158.
30. H. Fan, A. L. Sakulich, J. L. Goodier, X. Zhang, J. Qin, and R. J. Maraia (1997) *Cell* **88**, 707–715.
31. M. Bachmann, K. Pfeifer, H. C. Schroder, and W. E. Muller (1990) *Cell* **60**, 85–93.
32. M. Bachmann, H. C. Schroder, D. Falke, and W. E. Muller (1988) *Cell Biol. Int. Rep.* **12**, 765–789.
33. M. B. Mathews and T. Shenk (1991) *J. Virol.* **65**, 5657–5662.
34. K. H. Mellits, M. Kostura, and M. B. Mathews (1990) *Cell* **61**, 843–852.
35. G. D. Ghadge, S. Swaminathan, M. G. Katze, and B. Thimmapaya (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7140–7144.
36. P. Walter and A. E. Johnson (1994) *Annu. Rev. Cell Biol.* **10**, 87–119.
37. H. Lutcke (1995) *Eur. J. Biochem.* **228**, 531–550.
38. F. Bovia and K. Strub (1996) *J. Cell. Sci.* 2601–2608.
39. S. L. Wolin (1994) *Cell* **77**, 787–790.

## Small Nuclear Rnps (Snrnps)

In addition to the short conserved sequences at the branch site and the 5' and 3' [splice sites](#), [RNA splicing](#) of nuclear pre-mRNA is dependent on the activity of a large number of [trans-acting](#) factors. These can be divided into two classes: the snRNPs, evolutionarily conserved ribonucleoprotein particles, and non-snRNP protein splicing factors (reviewed in Refs. [1](#) and [2](#)). These *trans*-acting factors assemble in an ordered fashion onto the pre-mRNA substrate to form the [spliceosome](#), wherein the catalysis of splicing occurs. The spliceosomal snRNPs play a central role in the recognition and alignment of the pre-mRNA splice sites during spliceosome formation and are also involved in the catalytic steps of splicing (see [Spliceosome](#)). The major spliceosomal snRNPs, which are involved in the removal of U2-dependent introns (the most abundant type of nuclear pre-mRNA intron), include U1, U2, U5, and U4/U6, and are named according to their snRNA component(s). Due to their greater abundance ( $2 \times 10^5$  to  $10^6$  particles per cell, compared with 50 to 100 in yeast), the major spliceosomal snRNPs have been best characterized at the biochemical and structural level in higher eukaryotes and will thus be emphasized in the following paragraphs. In addition to U1, U2, U5, and U4/U6, a minor class of spliceosomal snRNPs involved in the excision of U12-dependent introns (which represent a small minority of nuclear pre-mRNA introns) has also been recently identified. These include the U11, U12, and U4atac/U6atac snRNPs, which, unlike the ubiquitous major snRNPs, appear to be lacking in some eukaryotes (eg, in the budding yeast *S. cerevisiae*) ([3-5](#)). As a consequence of their low abundance (estimated at  $10^3$  to  $10^4$  particles per cell) and more recent identification, the minor spliceosomal snRNPs are presently not well-characterized and thus will not be discussed in detail here.

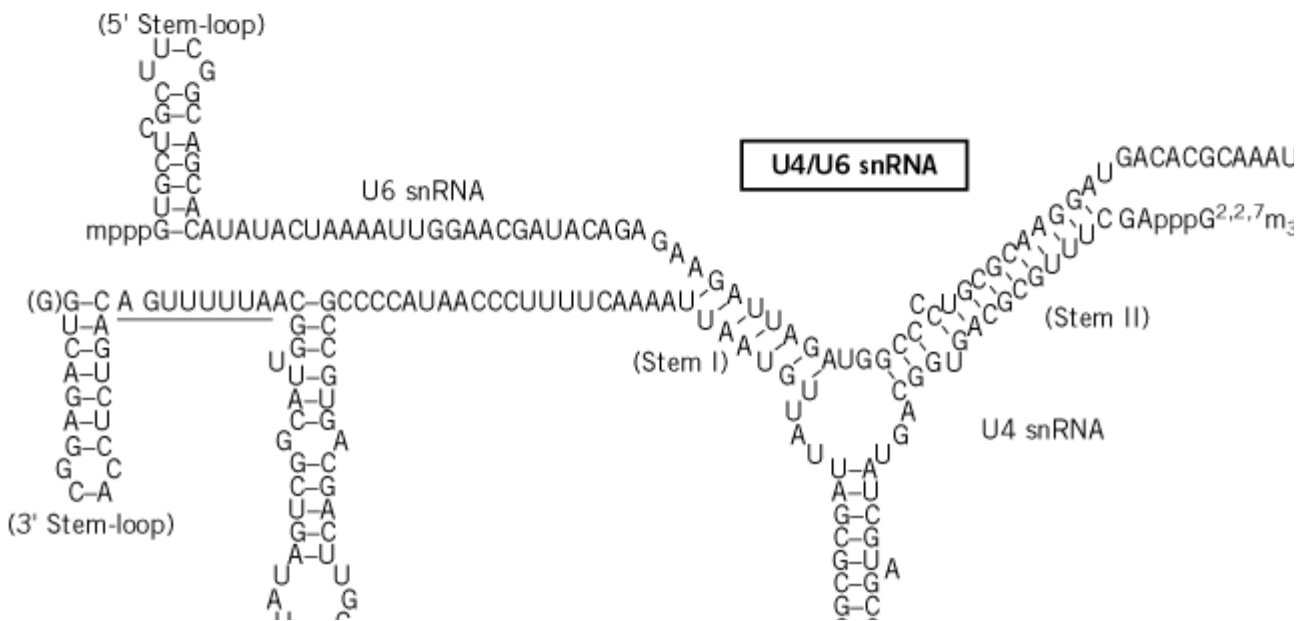
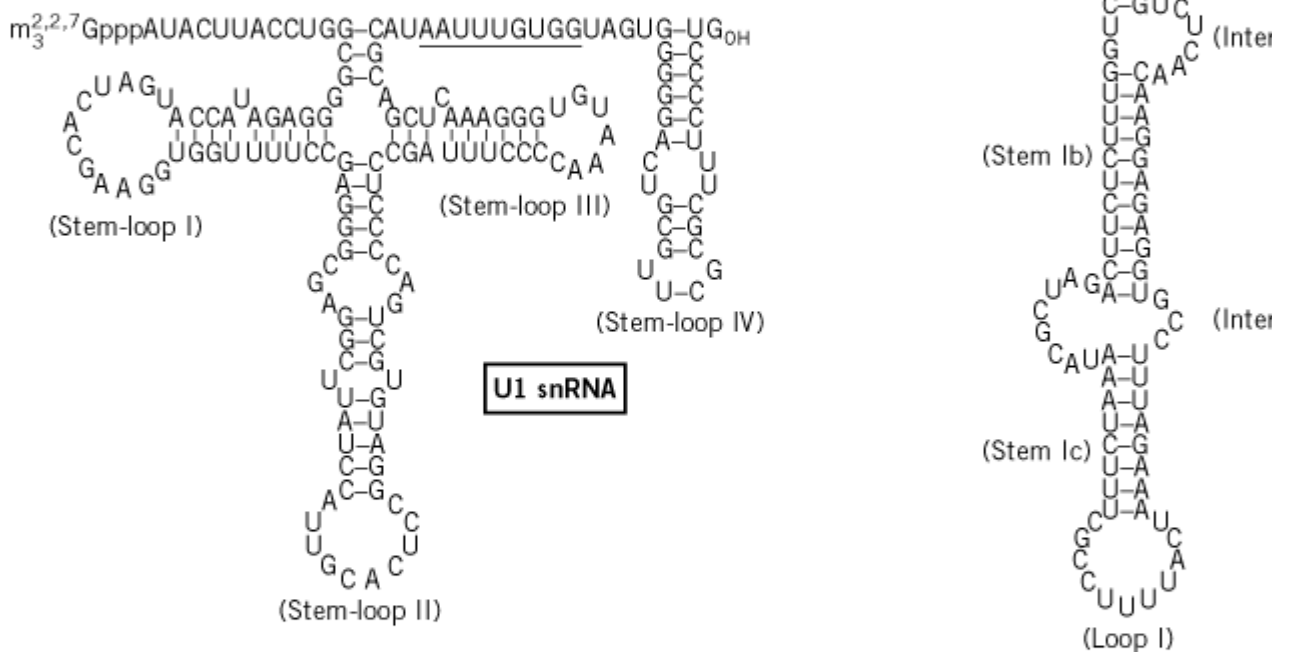
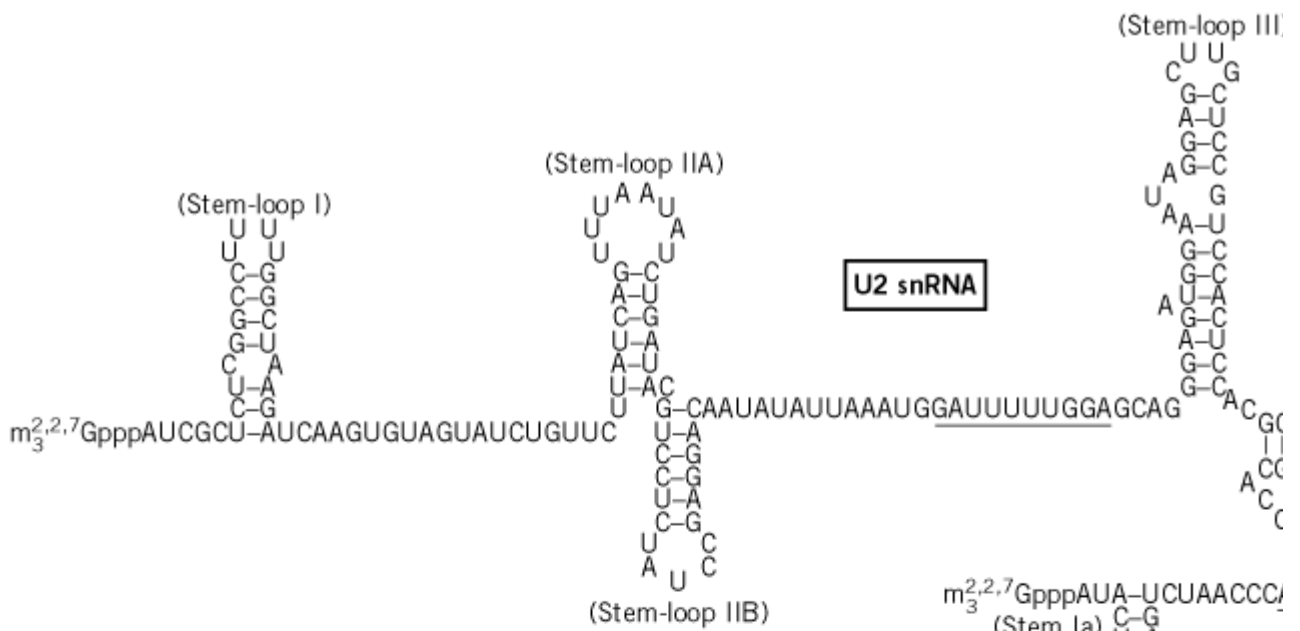
### 1. snRNA

The major spliceosomal snRNPs consist of one snRNA molecule (or two in the case of U4/U6) complexed with a number of proteins. The U1, U2, U4, U5, and U6 snRNAs are characterized by their small size (164, 187, 145, 116, and 106 nucleotides, respectively in humans), metabolic stability, and high degree of sequence conservation. They also contain a large number of modified nucleotides, such as pseudouridine and 2'-*O*-methylated nucleosides, and, with the exception of U6, possess a unique 2,2,7-trimethylguanosine [-cap](#) structure. In addition to their sequence conservation, the secondary structures of the metazoan snRNAs are also highly conserved such that a consensus secondary structure model can be generated for each of them. The sequence and most probable secondary structures of the human U1, U2, U4, U5, and U6 snRNAs are shown schematically in [Figure 1](#). The majority of the U4 and U6 snRNAs interact by extensive base pairing to form a characteristic Y-shaped structure and are present in this form within the U4/U6 snRNP. The most highly conserved sequences of these snRNAs are typically single-stranded regions that either base-pair with the pre-mRNA or other snRNAs during spliceosome assembly (eg, the 5' end of U1) or serve as binding sites for snRNP proteins (eg, the so-called Sm site; see text below). With the exception of U4 and U6, the major spliceosomal snRNAs are somewhat larger in yeast (eg, U1 and U2 are 568 and 1175 nucleotides long, respectively), and they therefore exhibit different secondary structures. It should be noted that the structures of most of the spliceosomal snRNAs are not fixed, but rather undergo conformational changes during splicing that reflect the dynamic nature of the splicing process. A prime example is the U6 snRNA, which first base-pairs with U4 in the U4/U6 snRNP and then, after the incorporation of the U4/U6 snRNP into the spliceosome, dissociates from this complex and basepairs with the U2 snRNA and the pre-mRNA (see [Spliceosome](#)). The U4 snRNA and the 5' half of the U2 snRNA also undergo significant conformational rearrangements during splicing ([6](#)). During splicing, the spliceosomal snRNAs engage in multiple base pairing interactions with other spliceosomal snRNAs and with the pre-mRNA. These RNA-RNA interactions play essential roles in the selection and favorable alignment of the 5' and 3' splice sites for splicing catalysis (see [Spliceosome](#)). Some of the structural features of the complex RNA network that is formed during spliceosome assembly mimic RNA structures that are essential for the

autocatalytic splicing of group II pre-mRNA [self-splicing introns](#). The spliceosomal snRNAs—in particular, U2, U5, and U6—thus appear to form, at least in part, the active sites of the spliceosome and take part directly in the catalysis of nuclear pre-mRNA splicing.

**Figure 1.** Sequence and secondary structure models of the human U1, U2, U5, and U4/U6 snRNAs. The consensus secondary structures shown are the models of Guthrie and Patterson ([15](#)), and that of U2 was proposed by Ares and Igel ([16](#)). The conserved sequence elements (conserved sequence elements of the common snRNP proteins) is underlined. The 5' ends of the snRNAs possess a cap structure (2,2,7-trimethyl-guanine cap) in spliceosomal snRNAs (with the exception of U6, which possesses a monomethyl phosphate cap).





## 2. snRNP Proteins

The spliceosomal snRNAs are present in cells as discrete ribonucleoprotein complexes. In human cells, the U1 and U2 snRNAs are organized as a 12 S and 17 S particle, respectively. U4/U6 and U5 snRNAs can be isolated as individual 12 S and 20 S particles, respectively, but they also associate with one another to generate a 25 S [U5.U4/U6] tri-snRNP complex, and they are integrated into the spliceosome as such. The protein composition of the major spliceosomal snRNPs has been best characterized in [HeLa Cells](#) and is summarized in Table 1. Proteins associated with the U1, U2, U5, and U4/U6 snRNPs fall into two classes. The first class consists of the so-called common or Sm proteins, denoted B', B, D1, D2, D3, E, F, and G, which are tightly associated with all snRNP particles. The Sm proteins play an important role in the biogenesis of the snRNPs and are essential for their import into the [nucleus](#) (see text below). The second class is comprised of the particle-specific proteins, which associate with a particular snRNP particle or complex. These proteins exhibit a wide range of binding affinities, and their association with an snRNP particle is thus dependent on the ionic strength of the particle's environment. In addition to the eight common snRNP proteins, the metazoan U1 snRNA is associated with three proteins, designated 70K, A, and C. Under physiological conditions, the U2 snRNA is complexed with 12 U2-specific proteins, whose molecular weights range from 33 to 160 kDa (Table 1). With 16 specific proteins and two sets of common proteins, the 25 S [U4/U6.U5] tri-snRNP possesses the most complex protein composition of the snRNPs (Table 1). The vast majority of proteins identified in isolated HeLa snRNPs are also present in the spliceosome (7), and thus most spliceosomal proteins have an snRNP origin. The majority of the human snRNP proteins have now been **cloned** and sequenced, and several have been shown to perform essential functions during spliceosome assembly and the catalytic steps of splicing (see [Spliceosome](#)). For example, snRNP proteins (eg, U1-70K, U5-220 kD, and the U2-specific proteins comprising SF3a and SF3b) are involved in [protein–protein](#) interactions and protein–RNA interactions that facilitate the interactions of the U1, U2, and U4/U6.U5 tri-snRNPs with the pre-mRNA during spliceosome assembly. Furthermore, several snRNP proteins (eg, the 20 kD, 100 kD, 116 kD, and 200 kD U4/U6.U5 tri-snRNP proteins) appear to possess enzymatic activities, such as RNA duplex unwinding or protein isomerization activity, that potentially drive the many conformational changes that occur in the spliceosomal RNA and/or protein networks. In this way, snRNP proteins make important contributions to the assembly of the active sites responsible for the catalysis of splicing. However, whether snRNP proteins contribute directly to catalysis is currently an open question.

**Table 1. Protein Composition of Human snRNPs<sup>a</sup>**

| Name | Approximate $M_r$<br>(kDa) | 12 S<br>U1 | 17 S<br>U2 | 25 S<br>U4/U6.U5 | <i>S. cerevisiae</i><br>Homologue |
|------|----------------------------|------------|------------|------------------|-----------------------------------|
| G    | 9                          | ^          | ^          | ^                | G                                 |
| F    | 11                         | ^          | ^          | ^                | F                                 |
| E    | 12                         | ^          | ^          | ^                | E                                 |
| D1   | 16                         | ^          | ^          | ^                | D1                                |
| D2   | 16.5                       | ^          | ^          | ^                | D2                                |
| D3   | 18                         | ^          | ^          | ^                | D3                                |
| B    | 28                         | ^          | ^          | ^                | B                                 |
| B'   | 29                         | ^          | ^          | ^                |                                   |

|      |      |          |   |         |
|------|------|----------|---|---------|
| C    | 22   | <i>f</i> |   | y U1-C  |
| A    | 34   | <i>f</i> |   | Mud1p   |
| 70K  | 70   | <i>f</i> |   | Snplp   |
| B''  | 28.5 |          | ○ |         |
| A'   | 31   |          | ○ |         |
|      | 33   |          | ○ |         |
|      | 35   |          | ○ |         |
|      | 92   |          | ○ |         |
| ^    | 60   |          | ○ | Prp9p   |
| SF3a | 66   |          | ○ | Prp11p  |
| ^    | 110  |          | ^ | Prp21p  |
| ^    | 53   |          | ○ | Hsh49p  |
| SF3b | 120  |          | ○ |         |
|      | 150  |          | ○ | Cus1p   |
| ○    | 160  |          | ○ |         |
|      | 15   |          |   |         |
|      | 40   |          |   |         |
|      | 65   |          |   |         |
|      | 100  |          |   | Prp28p  |
|      | 102  |          |   |         |
|      | 110  |          |   |         |
|      | 116  |          |   | Snu114p |
|      | 200  |          |   | Snu246p |
|      | 220  |          |   | Prp8p   |
|      | 15.5 |          |   | +       |
|      | 20   |          |   | +       |
|      | 60   |          |   | +       |
|      | 90   |          |   | +       |
|      | 27   |          |   |         |
|      | 61   |          |   |         |
|      | 63   |          |   |         |

<sup>a</sup> The presence of a given protein in a particular snRNP is indicated by the various symbols. The 25 S [U4/U6.U5] tri-snRNP complex contains two sets of common proteins (E, F, G, D1, D2, D3, B, B'). Tri-snRNP proteins that also associate with U4/U6 snRNPs are indicated by a ( ), and those also present in U5 snRNPs are marked with a solid diamond. Proteins associated solely with the 25 S [U4/U6.U5] are marked with a solid square. SF3a is composed of the 60, 66, and 110 kDa U2-specific proteins, and SF3b consists of the 53, 120, 150, and 160 kDa U2-specific proteins. Additional information (including references) about the yeast *S. cerevisiae* common proteins and U1 snRNP specific proteins are reported in Ref. [8](#) and for the Prp proteins in Ref. [2](#). Snu246 is also referred to as Brp2p ([17](#)), Slt22p ([18](#)), and Rss1p ([19](#)). The identification of Prp3p and Prp4p as homologues of the U4/U6 90 and 60 kDa proteins is described in Refs. [12](#) and [20](#). Hsh49p and Cus1p are described in Refs. [21](#) and [22](#), respectively. (Courtesy of Claudia Schneider.)

**Homologues** of some of the HeLa U1, U2, U5, and U4/U6 snRNP proteins have also been identified

either genetically or biochemically in yeast (summarized in Table 1). Significant progress has recently been made in the biochemical characterization of snRNPs from the yeast *S. cerevisiae*. This has complemented, as well as extended, the somewhat limited information obtained through yeast genetic techniques. Biochemical characterization of the yeast U1 snRNP has demonstrated that it possesses a significantly more complex protein composition than the metazoan U1 snRNP, containing not only homologues of the U1-70K, A and C proteins, but also six additional yeast-specific proteins (8). Some of these additional proteins appear to mediate the bridging of the 5' and 3' splice sites during the earliest stages of spliceosome formation in yeast. Although the number of snRNP protein homologues that have been characterized at the molecular level in yeast is presently limited, there are striking examples of both structural and functional conservation between yeast snRNP proteins and their counterparts in evolutionarily distant organisms. This not only underscores the functional importance of the snRNP proteins in nuclear pre-mRNA splicing, but also confirms that many aspects of the splicing process are conserved between higher and lower eukaryotes.

### 3. snRNP Biosynthesis and Structure

The snRNPs undergo a complex pathway of biogenesis that involves their shuttling between the nucleus and cytoplasm. Subsequent to their [transcription](#) in the nucleus, the major spliceosomal snRNAs, with the exception of U6, migrate to the cytoplasm. In the cytoplasm, some of the nucleosides of the snRNAs are modified (primarily 2' OH methylation and pseudouridylation) and their 5' m<sup>7</sup>G cap is hypermethylated to a 2,2,7-trimethylguanosine (m<sub>3</sub>G) cap. Cap hypermethylation is dependent upon the formation of an snRNP core structure. Core snRNPs are formed by the association of the common snRNP proteins (B', B, D3, D2, D1, E, F, G) with the Sm site, an evolutionarily conserved structural motif found in the U1, U2, U4, and U5 snRNAs (see Fig. 1). The Sm proteins form distinct heteromeric complexes (eg, E/F/G or B/D3) that interact with the Sm site in an ordered manner. Subsequent to core snRNP formation and cap hypermethylation, the snRNPs are translocated to the nucleus by an active, receptor-mediated process. The [nuclear import](#) pathway of the spliceosomal snRNPs is distinct from that of karyophilic proteins. snRNPs possess a bipartite nuclear localization signal that is comprised of the m<sub>3</sub>G cap and the snRNP core structure. At least one of the receptors responsible for snRNP import, Snurportin1, which specifically recognizes the m<sub>3</sub>G cap, has recently been described (9). SnRNP assembly is completed in the nucleus, where most of the particle-specific proteins are thought to associate with the snRNPs.

Although the vast majority of spliceosomal snRNP proteins have already been identified in humans, and a growing number are rapidly being identified in yeast, significantly less information about structural aspects (ie, RNA–protein and protein–protein interactions) of the spliceosomal snRNPs is currently available. Because the snRNPs serve as spliceosomal subunits, information regarding their higher-order structure may provide much needed information about the three-dimensional structure of the spliceosome. Information about the general morphology of both the mammalian and yeast spliceosomal snRNPs has been obtained through [electron microscopy](#). More detailed information about RNA–protein and protein–protein interactions within the snRNPs is, however, limited. The molecular characterization of snRNP proteins has revealed that they possess a variety of interesting structural motifs, which include, among others, [RNA-binding](#) and protein-interaction **domains** (eg, RRM, RS, and WD-40 domains). The best-characterized interactions within the snRNPs are those formed between snRNP proteins containing an RNA-binding domain (eg, U1-70K, U1-A, U2-B'') and their cognate snRNA. The atomic structure of interactions between the U1-A protein and U1 snRNA, and between the U2-B''/U2-A' proteins and the U2 snRNA, have been obtained by [X-ray crystallography](#) (10, 11). Well-characterized protein–protein interactions include those formed between U2-snRNP-specific proteins comprising the splicing factors SF3a and SF3b, as well as those among the Sm proteins. More recently, a number of novel snRNP protein–protein interactions have been reported that involve proteins associated with the [U4/U6.U5] tri-snRNP complex (12-14).

### Bibliography

1. J. F. Cáceres and A. R. Krainer (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford U.K., pp. 174–212.
2. P. E. Hodges, M. Plumpton, and J. D. Beggs (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford U.K., pp. 213–241.
3. S. L. Hall and R. A. Padgett (1996) *Science* **271**, 1716–1718.
4. W.-Y. Tarn and J. A. Steitz (1996) *Cell* **84**, 801–811.
5. W.-Y. Tarn and J. A. Steitz (1996) *Science* **273**, 1824–1832.
6. H. D. Madhani and C. Guthrie (1994) *Annu. Rev. Genet.* **28**, 1–26.
7. R. Reed and L. Palandjian (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, U.K. pp. 103–129.
8. A. Gottschalk, J. Tang, O. Puig, et al., (1998) *RNA* **4**, 374–393.
9. J. Huber, U. Cronshagen, M. Kadokura, C. Marshallsay, T. Wada, M. Sekine, and R. Lührmann (1998) *EMBO J.* **17**, 4114–4126.
10. C. Oubridge, N. Ito, P. R. Evans, C. H. Teo, and K. Nagai (1994) *Nature* **372**, 432–438.
11. S. R. Price, P. R. Evans, and K. Nagai (1998) *Nature* **394**, 645–650.
12. D. S. Horowitz, R. Kobayashi, and A. R. Krainer (1997) *RNA* **3**, 1374–1387.
13. S. Teigelkamp, T. Achsel, C. Mundt, S.-F. Göthel, U. Cronshagen, W. S. Lane, M. Marahiel, and R. Lührmann (1998) *RNA* **4**, 127–141.
14. T. Achsel, K. Ahrens, H. Brahms, S. Teigelkamp, and R. Lührmann (1998) *Mol. Cell. Biol.* **18**, 6756–6766.
15. C. Guthrie and B. Patterson (1988) *Annu. Rev. Genet.* **22**, 387–419.
16. M. Ares and A. H. Igel (1990) *Genes Dev.* **4**, 2132–2145.
17. S. M. Noble and C. Guthrie (1996) *Genetics* **143**, 67–80.
18. D. Xu, S. Nouraini, D. Field, S.-J. Tang, and J. D. Friesen (1996) *Nature* **381**, 709–713.
19. J. Lin and J. J. Rossi (1996) *RNA* **2**, 835–848.
20. J. Lauber, G. Plessel, S. Prehn et al., (1997) *RNA* **3**, 926–941.
21. H. Igel, S. Wells, R. Perriman and M. Ares (1998) *RNA* **4**, 1–10.
22. S. E. Wells, M. Neville, M. Haynes, J. Wang, H. Igel, and M. Ares (1996) *Genes Dev.* **10**, 220–232.

### **Suggestions for Further Reading**

23. C. L. Will and R. Lührmann (1997) "snRNP Structure and Function. In" *Eukaryotic mRNA Processing* (A. R. Krainer, ed.) IRL Press, Oxford, U.K., pp. 130–173.
24. A. Krämer (1996) The structure and function of proteins involved in mammalian pre-mRNA splicing. *Annu. Rev. Biochem.* **65**, 367–409.
25. C. L. Will and R. Lührmann (1997) Protein functions in pre-mRNA splicing. *Curr. Opin. Cell Biol.* **9**, 320–328.
26. C. Guthrie and B. Patterson (1988) Spliceosomal snRNAs. *Annu Rev. Genet.* **22**, 387–419.

### **Small-Angle Scattering**

Analysis of the small-angle scattering of X-rays, neutrons, or light is a valuable tool for studying the

structures and interactions of biological [macromolecules](#) in solution. Importantly, small-angle scattering experiments can be done on proteins and polynucleotides in a variety of solution conditions that mimic different physiological factors. Studies can be done on biological macromolecules and the assemblies they form over a very wide molecular weight range that encompasses the smallest proteins to the large animal [viruses](#). The information obtained is inherently low-resolution, generally confined to molecular shapes. When neutrons are used with [contrast variation](#), the shapes and dispositions of component structures in complexes can be obtained (see [Neutron Diffraction And Scattering](#) and [Contrast Variation](#)). [Light scattering](#) is fundamentally the same as X-ray scattering, except that the wavelengths are longer and there are differences in the theoretical formalisms (see [Light Scattering](#)).

## 1. Theory of Small-Angle Scattering

Small-angle scattering (also known as low-angle scattering, or sometimes solution scattering) results from the constructive interference of secondary waves that are scattered when a plane wave interacts with matter. The small-angle [scattering intensity distribution](#) from a particle in solution is maximum at zero scattering angle and falls off with a rate that depends upon the size and shape of the scattering particle. The larger the particle, the faster the falloff. [Small-angle scattering](#) generally refers to the elastic, coherent scattering of X-rays or neutrons, both of which have the properties of plane waves. Scattering theory assumes the wavelength of the radiation being used as a structural probe is smaller than the dimensions of the object being studied. X-rays and neutrons can be produced with wavelengths from 0.1 to 10's of Å (1 Å = 10<sup>-10</sup>m), and hence they are useful for studies of biological molecules with dimensions from 10 to 1000's of Å.

The total [scattering intensity distribution](#) for a homogeneous solution of monodisperse, randomly oriented, noninteracting particles in a solution can be expressed as

$$I(Q) = \int \int \Delta\rho(\mathbf{r}_1)\Delta\rho(\mathbf{r}_2) \frac{\sin Q|\mathbf{r}_1 - \mathbf{r}_2|}{Q|\mathbf{r}_1 - \mathbf{r}_2|} d\mathbf{r}_1 d\mathbf{r}_2 \quad (1)$$

where  $Q$  is the momentum transfer or scattering vector amplitude and is equal to  $4\pi(\sin q)/\lambda$  ( $q$  is half the scattering angle and  $\lambda$  is the wavelength of the radiation; note that  $Q$  is also called  $K$ ,  $h$ , or  $q$ , and some authors prefer to use  $s = Q/2\pi$ ).  $\Delta\rho(\mathbf{r}) = \rho(\mathbf{r}) - \rho_s$ , and is the “contrast,” or scattering density difference between the particle and solvent. A very useful function in small-angle scattering is the pair distance, or vector length distribution function  $P(r)$ , which is related to  $I(Q)$  by a Fourier transformation:

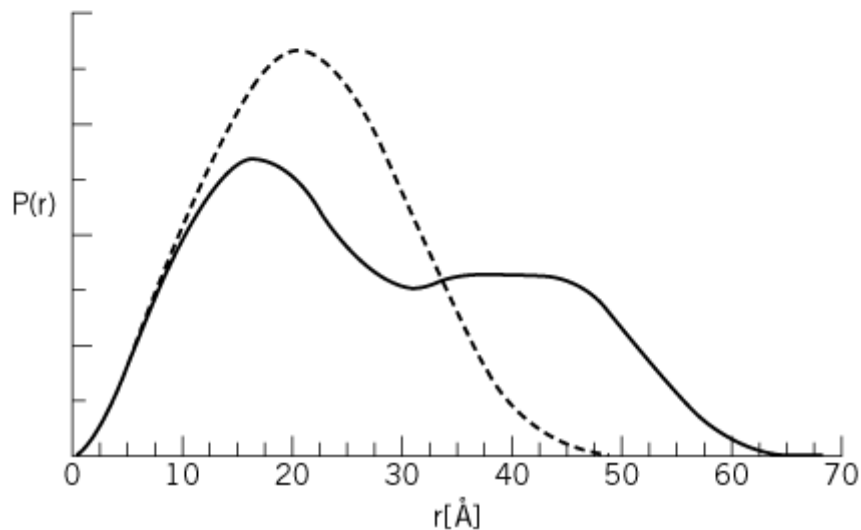
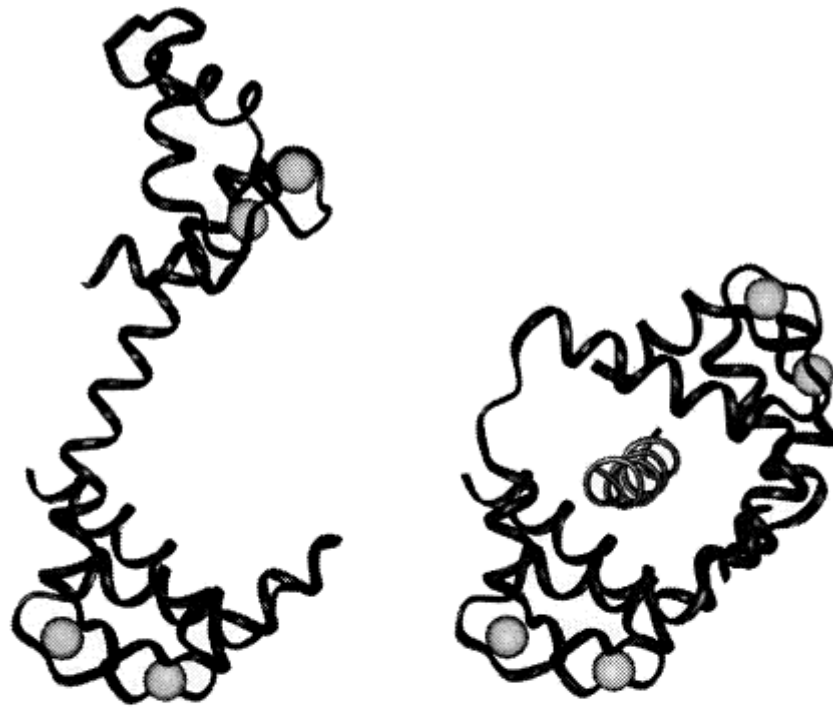
$$I(Q) = 4\pi \int P(r) \frac{\sin Qr}{Qr} dr \quad (2)$$

$$P(r) = \frac{1}{2\pi^2} \int I(Q) Q \cdot r \sin(Q \cdot r) dQ \quad (3)$$

$P(r)$  is the frequency of vector lengths connecting small-volume elements within the scattering particle, weighted by their scattering densities.  $P(r)$  goes to zero at a value corresponding to the maximum dimension of the particle,  $d_{\max}$ . Because  $I(Q)$  can only be measured over a finite range of  $Q$  values,  $P(r)$  is generally calculated using an indirect Fourier transform method in which the coefficients of the Fourier transform of a set of functions in real space are fit to the scattering data. Several such methods have been developed using different basis sets (1, 2); alternatively, a regularization method can be used (3).  $P(r)$  is extremely sensitive to the symmetry of the scattering particle and to the relationships between domains or repeating structures. This effect is demonstrated in Figure 1, which shows the structures and corresponding model  $P(r)$  functions for the calcium-

binding protein [calmodulin](#) and its complex with the peptide sequence of its binding domain in myosin light-chain kinase. The two globular domains of the uncomplexed calmodulin give rise to a  $P(r)$  function with a maximum whose position is dominated by the most frequently occurring vector lengths within the individual domains, and a shoulder whose position is influenced by vector lengths between domains. The calmodulin complex is more compact and symmetric and has a correspondingly more symmetric  $P(r)$ , and its  $d_{\max}$  is smaller. Small-angle scattering experiments showed this dramatic effect in 1989 (4), which has since been confirmed by high-resolution solution [NMR](#) studies (5) and has become a key element in our understanding of calcium/calmodulin enzyme activation (6).

**Figure 1.** Ribbon representations of the structures of the crystal form of calmodulin [*upper left* (22)], its complex with its binding domain from myosin light chain kinase (MLCK) determined by NMR [*upper right* (5)], and the corresponding  $P(r)$  functions represented by solid and dashed lines, respectively. Calcium atoms are represented as spheres. Small-angle scattering experiments (4) detected this dramatic conformational collapse by calmodulin that occurs upon binding the helical MLCK binding domain and is effected by flexibility in the helix connecting the two lobes of calmodulin. This figure is based on work done at Los Alamos National Laboratory.



A parameter often used in the interpretation of small-angle scattering data is the [radius of gyration](#),  $R_g$ , which can be calculated as the second moment of  $P(r)$ . The zeroth moment of  $P(r)$  gives the forward scatter,  $I(0)$ , which, when normalized to the protein concentration expressed in units of mg/mL, is directly proportional to the molecular weight of the particle for particles of equal scattering density (7). If the  $I(0)$  values are normalized to molar concentrations, they will be proportional to the square of the molecular weight. Because most proteins have similar X-ray scattering densities, protein standards known to be monodisperse in solution can be used to determine the molecular weight of the scattering particle and to evaluate possible sample aggregation. For neutron scattering, the data can be scaled absolutely and precise molecular weights calculated (8). For a monodisperse solution of particles, the volume,  $V$ , of the scattering particle can also be calculated using the Porod relation:



$$V = \frac{2\pi^2 I(0)}{Q_i} \quad (4)$$

where  $Q_i$  is the scattering invariant:

$$Q_i = \int_0^\infty Q^2 dQ I(Q) \quad (5)$$

## 2. Interparticle Interference and Solvent Effects

Small-angle scattering theory assumes there is no correlation in positions or orientation of the particles in solution. For a dilute solution (<10mg/ml of a small protein ~50 kDa) the average distance between particles is ~200 Å and long-range [electrostatic interactions](#) are generally canceled by the distribution of positive and negative charges on the surfaces of the protein molecules. For highly charged macromolecules, or for higher concentrations, this approximation may not hold. In these cases, the scattering function will be convoluted with an “interparticle” structure factor  $S(Q)$  which is less than 1 at  $S(0)$ , oscillates around unity, and is strongly damped with increasing  $Q$  (9).  $S(Q)$  has a first maximum at  $\sim 2\pi/r_0$ , where  $r_0$  is the average distance between neighboring objects.

The contribution of the interparticle structure factor to the total scattering will be diminution of the lowest- $Q$  data, giving rise to an apparent reduction in the structural parameters derived from the scattering data. Often there is a linear dependence on particle concentration,  $c$ , to these interparticle interference effects, and they can be eliminated by extrapolating  $I(Q)/c$  versus  $c$  to infinite dilution ( $c = 0$ ). Alternatively, one can minimize interparticle interference by adjusting the solution conditions, such as the pH (moving closer to the [isoelectric point](#) for the molecule can reduce the net charge on the molecule) or ionic strength (counterions can effectively screen charges on the surface of the molecule). If attractive forces exist between particles, the interparticle structure factor  $S(0)$  is greater than 1, and the scattering particles will appear larger than they are.

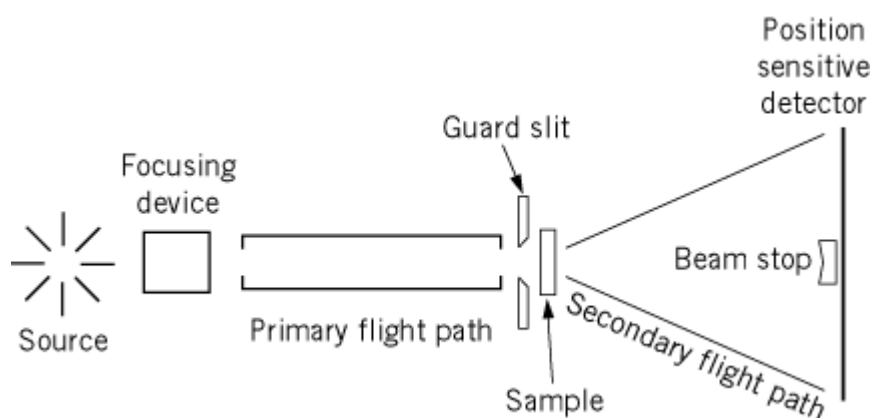
It is in general a good approximation to assume that a soluble protein in dilute aqueous buffer (<0.1 M salts) can be considered as a two-component solution of solvent and solute. In the case of highly charged molecules (such as DNA and very acidic or basic proteins in neutral pH solutions) with relatively high concentrations of salts, however, there can be preferential binding or repulsion of counterions to the surface of the protein that will effect scattering measurements (10). Such systems effectively have three components: solute, bulk solvent, and solvent at the surface of the solute. The effect on the scattering data is to make the scattering particles appear larger than they are.

## 3. The Experimental Setup

A small-angle scattering experiment is conceptually simple (Fig. 2). A source of radiation is required, followed by a collimated path (the primary flight path). In general the radiation is monochromatized, although instruments on pulsed neutron spallation sources have increased intensity and  $Q$  range using a “white” beam and time-of-flight methods to determine the wavelength of the neutrons arriving at the detector (11). In the case of X-rays, a focusing device is installed between the source and the primary flight path. Immediately after the exit aperture of the primary flight path, guard slits can be placed so as not to touch the beam, but to eliminate parasitic scatter from the experiment. The sample is placed immediately after the guard slits. Typically, sample volumes for X-ray experiments are quite small (~10 to 50 mL), while neutron experiments use larger beams and hence require larger sample volumes (~200 to 600 mL). The solute concentrations required for both X-ray and neutron experiments are similar, 1 to 10 mg/mL, depending on the molecular weight of the solute. Because of the molecular weight dependence of the scattering intensity, larger molecules can be measured using samples with lower solute concentrations. A position-sensitive detector is placed at some distance from the sample. This distance is

approximately equal to the primary flight path in order to minimize beam divergence. The path after the sample to the detector is called the secondary flight path. Detectors may be one- or two-dimensional position-sensitive devices. Because small-angle scattering from randomly oriented particles in solution gives rise to circularly symmetric scattering, two-dimensional detectors offer the opportunity for circular averaging in order to increase signal-to-noise ratio. The total flight path of the instrument is designed so that ideally, in the absence of a sample, the radiation enters the collimator and passes undeflected through the instrument, to be detected at a single point on the detector. In practice, this point is a small area due to beam divergence. For X-rays, highly polished mirrors with a slight bend can be used as focusing devices to increase the intensity of X-rays throughout the systems. There are no simple focusing devices for neutrons, and hence the collimation is designed simply to select neutrons whose path is already directed to a point on the detector. When a sample is placed in position, it is bathed in radiation, and a fraction of the radiation is deflected to give rise to the small-angle scattering pattern. Because the intensity of the scattered radiation is small, very sensitive detectors are used that cannot tolerate the high intensity of the undeflected beam, so a small beam stop is placed in front of the detector to absorb the direct beam. The sample-to-detector distance needs to be large enough to resolve very small angles of fractions of degrees. The greater this angular resolution, the larger the objects that can be studied. For a small-angle X-ray instrument with copper  $K_{\alpha}$  radiation ( $\lambda = 1.54 \text{ \AA}$ ,  $1 \text{ \AA} = 10^{-10} \text{ m}$ ), a 65-cm secondary path length, and an ionizing gas chamber linear wire detector,  $Q$  values as low as  $0.006 \text{ \AA}^{-1}$  can be measured, which allows study of objects with dimensions of 100's of  $\text{\AA}$ . Larger molecular dimensions can be measured using instruments with longer flight paths and/or longer wavelength radiation. Neutron small-angle scattering instruments typically can select wavelengths over a fairly wide range (4 to 10  $\text{\AA}$  at reactors, 0.1 to 20  $\text{\AA}$  at pulsed neutron spallation sources) and can have variable sample-to-detector distances, ranging from 1 to tens of meters.

**Figure 2.** Schematic diagram of a small-angle scattering experiment geometry (not to scale).  $2\theta$  is the scattering angle.



The theory of small-angle scattering assumes point sources of radiation and monochromatic (single wavelength) radiation. This is not true, in practice and in fact some instruments are designed to use a slit beam geometry and/or a range of wavelengths to increase the flux on the sample, which results in a “smearing” of the scattering pattern that affects both the amplitude and shape of the scattering profile. Corrections can be made for both factors. In general it is easier to “smear” the models used to interpret the scattering data rather than to desmear the data. In addition, error propagation can be done more rigorously in the former case.

#### 4. Examples of Applications of Small-Angle Scattering

Significant advances have been made in a broad range of areas in structural biology based on X-ray or neutron small-angle scattering analyses. Small-angle [X-ray scattering](#) has been used to characterize molecular hinges in bilobal proteins that allow for cleft closure to initiate either ligand binding and/or substrate binding and catalysis (12, 13) and the conformational changes involved in biochemical regulation by calmodulin (6). Small-angle [X-ray scattering](#) has also contributed to our understanding of the degree of compaction in **protein folding** intermediates, such as [molten globules](#) (14-16). Small-angle neutron scattering and [contrast variation](#) has yielded a model structure showing the shapes and interactions of the muscle proteins troponin C and troponin I (17) that has implications for the calcium-sensitive “switch” mechanism that controls the contractile apparatus. In the late 1970s, the solution structure of the [nucleosome](#) core particle was determined using neutron [small-angle scattering](#) and [contrast variation](#) (18, 19). The important result of this landmark study was that the DNA was found to wind around the outside of the histone protein core, which has important implications for DNA accessibility in [DNA replication](#) and [transcription](#) mechanisms. Neutron [small-angle scattering](#) with deuterium labeling and [contrast variation](#) has also revealed the dispositions of the component proteins and/or nucleotides in the [ribosomes](#) (20, 21) that are the **protein biosynthesis** factories in cells (see [Neutron Diffraction And Scattering](#)).

### Bibliography

1. P. B. Moore (1980) *J. Appl. Crystallogr.* **13**, 168–175.
2. O. Glatter (1977) *J. Appl. Crystallogr.* **10**, 415–421.
3. D. I. Svergun, A. V. Semenyuk, and L. A. Feigen (1988) *Acta Crystallogr.* **A44**, 244–250.
4. D. B. Heidorn et al. (1989) *Biochemistry* **28**, 6757–6764.
5. M. Ikura et al. (1992) *Science* **256**, 632–638.
6. J. K. Krueger et al. (1998) *Biochemistry* **37**, 13997–14004.
7. W. R. Kringbaum and F. R. Kugler (1970) *Biochemistry* **9**, 1216–1223.
8. B. Jacrot and G. Zaccai (1981) *Biopolymers* **20**, 2413–2426.
9. S.-H. Chen and D. Bendedouch (1986) *Methods Enzymol.* **130**, 79–116.
10. G. Zaccai (1986) *Methods Enzymol.* **127**, 619–629.
11. P. A. Seeger and R. P. Hjelm (1994) *Transactions of the American Crystallography Association, Proc. Symp. Time-of-Flight Diffraction at Pulsed Neutron Sources*, **29**, 63–77.
12. G. A. Olah et al. (1993) *Biochemistry* **32**, 3649–3657.
13. C. A. Pickover, D. B. McKay, D. M. Engelman, and T. A. Steitz (1979) *J. Biol. Chem.* **254**, 11323–11329.
14. G. V. Semisotnov et al. (1996) *J. Mol. Biol.* **262**, 559–574.
15. T. R. Sosnick and J. Trewhella (1992) *Biochemistry* **31**, 8329–8335.
16. L. Chen, K. O. Hodgson, and S. Doniach (1996) *J. Mol. Biol.* **261**, 658–671.
17. G. A. Olah and J. Trewhella (1994) *Biochemistry* **33**, 12800–12806.
18. J. F. Pardon et al. (1975) *Nucleic Acids Res.* **2**, 2163–2176.
19. P. Suau et al. (1977) *Nucleic Acids Res.* **4**, 3769–3786.
20. M. S. Capel et al. (1987) *Science* **238**, 1403–1406.
21. R. P. May et al. (1992) *EMBO J.* **11**, 373–378.
22. Y. S. Babu, C. E. Bugg, and W. J. Cook (1988) *J. Mol. Biol.* **204**, 191–204.

### Suggestions for Further Reading

23. L. A. Feigen and D. I. Svergun (1987) *Structure Analysis by Small-Angle X-ray and Neutron Scattering*, Plenum Press, New York.
24. O. Glatter and O. Kratky (1982) *Small-Angle X-Ray Scattering*, Academic Press, New York.
25. P. B. Moore (1982) "Small-Angle Scattering Techniques for the Study of Biological Macromolecules and Macromolecular Aggregates. In" *Methods of Experimental Physics*, Vol.

## Sodium Dodecyl Sulfate (SDS)

Of all the [detergents](#) used, SDS is probably the strongest **denaturing** one known. In addition, earlier studies have also shown that SDS is efficient in solubilizing both [proteins](#) and [lipids](#) present in cell [membranes](#) (1). However, the highly denaturing property of SDS almost precludes its use in obtaining solubilized [membrane proteins](#) in their active conformation. In a few cases, SDS has been used to solubilize amphiphilic peptides and small membrane proteins for solution-state nuclear magnetic resonance ([NMR](#)) analysis (2). SDS is widely used in SDS–polyacrylamide gel electrophoresis ([SDS–PAGE](#)) for analyzing denatured proteins based on their molecular weights. SDS–PAGE has also been used on a preparative scale in the purification of proteins for N-terminal protein **sequencing**, raising **antibodies** and diagnostic studies of inborn neurological disorders (3, 4). Following separation, gel slices that contain the protein bands (or spots, when the proteins are resolved by [two-dimensional gel electrophoresis](#)) are excised from the [polyacrylamide](#) gel, and the proteins are then separated from the gel by [electroelution](#) .

The essential principle of SDS–PAGE is that protein molecules are denatured in SDS such that the **hydrophobic** portions of the long-chain SDS molecules envelope the denatured [polypeptide chain](#) . Consequently, the native charge of any protein is completely masked by the negative charge on the SDS molecules, and all proteins adopt similar denatured conformations. Thus, once the charge and shape differences among all the different protein molecules are eliminated, the proteins move (under an electric field) through the polyacrylamide gel matrix according to their molecular weights, that is, the high molecular weight proteins move more slowly than the low molecular weight proteins.

### Bibliography

1. P. Banerjee, J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
2. G. D. Henry and B. D. Sykes (1994) *Methods Enzymol.* **239**, 515–535.
3. H. Schagger and G. von Jagow (1991) *Anal. Biochem.* **199**, 223–231.
4. J. Schupbach, R. W. Ammann, and A. U. Freiburghaus (1991) *Anal. Biochem.* **196**, 337–343.

## Solvent Flattening

The result of a structure determination by [X-ray crystallography](#) is the calculation of an electron density map and the construction of a molecular model of the structure in this map. Protein crystals contain large amounts of solvent that occupy the voids between the protein molecules. This solvent has a dynamic character. Therefore, in a high-quality electron density map, the solvent region is flat and has a relatively low density, whereas the protein region has a higher density and is certainly not flat. However, if the map is of poor quality, the solvent region will contain peaks and valleys of density. In solvent flattening they are simply removed by modifying the density to a low, constant value. An important condition is that the regions occupied by the protein molecules can be identified

as regions with relatively high density and that the boundary between the protein and the solvent region can be traced ([1](#), [2](#)). This modified map is closer to the correct map than the first, preliminary electron density map.

From this improved map, [structure factors](#) including phase angles are calculated. Then, the calculated phase angles are combined with the observed structure factor amplitudes to calculate a new electron density map. This solvent flattening procedure is repeated until no further improvement is obtained.

#### Bibliography

1. B.-C. Wang (1985) *Methods Enzymol.* **115**, 90–112.
2. A. G. W. Leslie (1987) *Acta Crystallogr.* **A43**, 134–136.

#### Suggestion for Further Reading

3. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

### Solvent Perturbation Spectroscopy

The absorption spectra of the aromatic amino acids ([phenylalanine](#), [tyrosine](#), and [tryptophan](#)) depend on their molecular environments. Alterations in the solvent can thus be used to change the environment of those aromatic side-chains in a folded protein that are exposed to the solvent (see [Protein Structure](#)). In principle, this provides a method to identify residues that are located at the protein surface.

In general, small shifts in the wavelength of absorption occur as a response to changes in the solvent: a blue-shift of the spectra to shorter wavelengths is observed when the polarity of the solvent increases and a red-shift when the polarity decreases. For example, the maximum of the absorption of the phenol chromophore is blue-shifted by about 3 nm when the solvent is changed from CCl<sub>4</sub> to water ([1](#)). The solvent perturbation method has, however, found only limited applications. First, the changes in absorbance are generally small; secondly, the changes in the nature of the solvent have to remain very small in order not to perturb the conformation of the protein under investigation. The exposure of [tryptophan](#) residues to the solvent is much more easily investigated by fluorescence quenching (see [Fluorescence Spectroscopy](#)).

#### Bibliography

1. S. Yanari and F. A. Bovey (1960) *J. Biol. Chem.* **235**, 2818–2826.

#### Suggestion for Further Reading

2. J. W. Donovan (1973) Ultraviolet difference spectroscopy—new techniques and applications. *Methods Enzymol.* **27**, 497–525.

### Somatic Hypermutations

Somatic hypermutations are a unique mechanism used by [B cells](#) to amplify [immunoglobulin](#) (Ig) diversity and permit [antigen](#) to select clones that express high-affinity **antibodies**. The rate of such mutations is very high, generally estimated to be 1 per 1000 base pairs per cell division. At the onset of a primary [immune response](#), antibodies are of the [IgM](#) isotype and are generally of moderate affinity, because they essentially result from the expression of germline genes that have rearranged their mosaic of VDJ and VJ gene segments for the heavy and light chains, respectively, and that have no or very few mutations. As a rule, **T-cell-dependent** antigens will induce a second wave of antibody production, now of the [IgG](#) isotype. This takes place in secondary lymphoid organs (such as spleen, tonsils, or lymph nodes) that have been colonized by some stimulated B cells, generating a germinal center, where active cellular interactions take place, including cooperation between T and B cells, which is an absolute prerequisite for the initiation of hypermutations. B cells divide rapidly, giving rise to a clonal amplification *in situ*. They present a blastic morphology and are known as centroblasts, which define the so-called dark zone of the germinal center. As the centroblasts actively divide, they accumulate somatic mutations that are localized mostly in the regions of the VH and VL segments that do interact with the antigen. As a result of these mutations, the array of affinities will be largely increased, and the organization of the germinal center will permit the antigen to select positively the B cells that express at their surface the immunoglobulins with the greatest affinities. At this stage, centroblasts cease to divide, and they become small nondividing centrocytes that accumulate in the light zone of the germinal center, where they actively interact with numerous follicular dendritic cells. As a result of this interaction, the selected centrocytes can now be directed toward either plasma cell differentiation or become long-lived B [memory cells](#).

The molecular mechanisms that drive somatic hypermutations are still only partially understood. Very clearly, this is a specific feature of B cells; although T cells share many common features with B cells, they make their **T-cell-receptor repertoire** by **gene rearrangements** alone and do not mutate. Somatic mutations have long been suspected on mouse I chains, because it was known that only three VI genes were present in this species. Comparing their genomic structure with that of [complementary DNAs](#) prepared from I $\mu$ -expressing plasmacytomas indicated that additional mechanisms had to operate to amplify the diversity. This was also extensively studied with large collections of [monoclonal antibodies](#) derived at various times after immunization of inbred mice against the nitro-iodo phenyl arsonate hapten. The number of nucleotide substitutions clearly increased with time. They were scattered all along the variable regions, with a clear accumulation in the hypervariable regions, which suggested that a selection mechanism had taken place. A puzzling observation was that mutations were limited to the **V region**. Compared with the rearranged genomic gene, it was apparent that mutations were found not only in the coding region, but also in the 5'- and 3'-untranslated flanking regions. More precise information was obtained by Neuberger and Milstein at Cambridge (1). First, it was shown with the mouse k chain model that mutations were strictly dependent on the presence of the **enhancer**-containing region, located at the 3' position of the C $\kappa$  gene. Second, it was demonstrated that if the V $\kappa$  gene region was replaced in a construct by another gene, such as the **b-globin** gene, mutations still occurred in the inserted sequence, clearly demonstrating that the mutational events were not dependent upon the V region structure itself. The information that is lacking to date is the nature of the enzymatic mechanism that must generate these mutations. Despite intense efforts from several laboratories, this problem has still not been solved.

See also entries [Affinity Maturation](#), [Antibody](#), [Class Switching](#), and [Immune Response](#).

#### Bibliography

1. C. J. Jolly, S. D. Wagner, C. Rada, N. Klix, C. Milstein, and M. S. Neuberger (1996) The targeting of somatic hypermutation. *Semin. Immunol.* **8**, 159–168.

#### Suggestions for Further Reading

2. N. S. Green, M. M. Lin, and M. D. Scharff (1998) Somatic hypermutation of antibody genes: a hot spot warms up. *Bioessays* **20**, 227–234.
3. A. Ehlich, V. Martin, W. Muller, and K. Rajewsky (1994) Analysis of the B-cell compartment at the level of single cells. *Curr. Biol.* **4**, 573–583.

## Somatic Mutation

Somatic mutations are heritable changes in the [genome](#) of the non **germ-line**, **somatic cells** of an individual. When cells carrying a somatic [mutation](#) divide, they give rise to **clones** that are also [mutant](#). The proportion of mutant cells in the population depends on the timing of the mutation. If the mutations occur early in [development](#), the ratio of mutants to **wild-type** in the somatic tissue population is high, whereas if the mutation arises late in development, the proportion of mutants is low. In addition, somatic mutations can give rise to cancer (see [Neoplastic Transformation](#)). Somatic mutations occur at high frequency in **antibody**-producing cells. Some of these mutations result in antibodies with higher affinity for the [antigen](#), thereby improving the specificity of the antibody response (see [Affinity Maturation](#)).

## Sonic Hedgehog

The sonic hedgehog (SHH) protein belongs to a family of intercellular signaling molecules involved in patterning, growth, and cell-fate specification during vertebrate embryogenesis. SHH derives its name from the **hedgehog** (*hh*) gene of *Drosophila*, a segment polarity gene that was identified in a large screen for [mutations](#) that disrupt embryonic development. The bristle pattern on the cuticle of hh mutant flies is aberrant, which gives them a rough or hedgehog-like appearance. Several **Hedgehog** (*HH*) **genes** have been **cloned** from chick, rodent, human, frog, and [zebrafish](#), and are expressed in tissues that have patterning activity and at sites of epithelial-mesenchymal interactions all over the developing embryo. The HH [signal transduction](#) pathway is highly conserved from *Drosophila* to vertebrates, and [mutations](#) in the *SHH* gene and components of its signaling pathway are responsible for developmental defects and tumorigenesis in humans.

### 1. SHH Protein Biochemistry

The SHH gene encodes for a secreted protein that consists of a [signal peptide](#), a highly conserved N-terminal region and a more divergent C-terminal **domain**. SHH is produced as a 45-kDa precursor protein that cleaves itself **proteolytically** to generate two fragments, a 25-kDa carboxy-terminal fragment and a 19-kDa amino-terminal, biologically active fragment (SHH-N). The proteolysis, which is **catalyzed** by conserved residues within the carboxy-terminal domain of the precursor protein, also results in the covalent addition of a cholesterol molecule to the C-terminus of the SHH-N fragment (1). The lipid modification of SHH-N has profound consequences for localization of the protein. Cholesterol-modified SHH-N remains associated with the cell that produces it, whereas unmodified SHH-N produced from truncated [complementary DNA](#) is freely diffusible. Cell association is important for the patterning functions of HH in development. It may create a high local concentration of the protein and may also restrict its domain of activity within a tissue (2).

## 2. SHH and Patterning in the Vertebrate Embryo

HH proteins play key roles in patterning many tissues in the developing vertebrate embryo. Indian HH and Desert HH, two other mammalian HH **homologues**, regulate bone development and spermatogenesis, respectively. SHH is expressed at numerous sites of epithelial-mesenchymal interactions in the developing embryo, including the gut, lung, hair, whisker, and tooth. It also plays a role in establishing some of the major axes in early embryonic development: left-right asymmetry of the body, dorsal-ventral polarity in the central nervous system, and anterior-posterior polarity in the limb. SHH also induces specific neuronal cell fates in the central nervous system and promotes proliferation in a number of tissues.

### 2.1. SHH and Left-Right Asymmetry of the Body

The vertebrate body plan is patterned along three axes, anterior-posterior, dorsal-ventral, and left-right. The establishment of left-right asymmetry is essential for normal development and the correct positions of the internal organs. In vertebrates the heart is normally on the left and the liver on the right. At early developmental stages in the chick embryo, SHH expression is restricted to the left side of Hensen's node (a signaling center), well before the appearance of morphological asymmetries. The left-sided SHH expression restricts the expression of nodal, a secreted molecule of the [transforming growth factor](#) b (TGFb) superfamily, to the left side of the node. Ectopic expression of SHH on the right side of the embryo induces nodal expression on the right side and alters the position of the heart. These results suggest that SHH is part of the molecular cascade that controls left-right asymmetry in the developing chick embryo.

### 2.2. SHH and Neural Tube Patterning

The embryonic nervous system consists of a tube of proliferating cells that extends along the anterior-posterior axis of the embryo. The anterior region of the neural tube gives rise to the brain, and more posterior regions give rise to the spinal cord. The neural tube is patterned along its dorsal-ventral axis by signals from the overlying surface ectoderm and from the underlying mesodermally derived notochord. Signals from the notochord ventralize the neural tube and induce three major cell types: floor plate cells at the ventral midline, motoneurons at the ventrolateral position, and interneurons at more dorsal locations. SHH is expressed in the notochord at the stage in which it mediates its inductive effects on the neural tube, and treatment with SHH [recombinant protein](#) mimics the effects of the notochord on floor plate and motoneuron induction in neural tube explants cultured *in vitro*. The floor plate and motoneurons do not develop in SHH-null animals, suggesting that SHH signaling is required to induce midline structures in the neural tube (3).

The effect of SHH on cell fate induction in the neural tube is concentration-dependent. High concentrations induce floor plate, intermediate concentrations induce motoneurons, and low concentrations induce interneurons. Graded SHH signaling also regulates the expression of [transcription factors](#) in the neural tube. Low concentrations of SHH initially repress the expression of paired domain transcription factors PAX-3 and PAX-7 in the ventral neural tube, generating a population of ventralized progenitor cells (see [Pax Genes](#)). Later, graded SHH signaling establishes regional differences in the expression of two additional transcription factors, Nkx2.2 and PAX-6. Nkx2.2 is induced ventrally, and PAX-6 is expressed in a gradient dorsal to the Nkx2.2 expression domain. The position of the different SHH-induced neuron subclasses correlates with the regional differences in transcription-factor expression along the dorsal-ventral axis of the neural tube. SHH may establish neuronal identity in the neural tube by differentially regulating the expression of transcription factors, thereby establishing different progenitor cell populations along the dorsal-ventral axis of the neural tube. The response of progenitor cells in the neural tube to SHH signaling also depends on their position along the anterior-posterior axis of the neural tube; SHH induces the development of cholinergic neurons in the forebrain, dopaminergic neurons in the midbrain, and motoneurons in the spinal cord. The different outcomes of SHH signaling may be due to differences in progenitor cell identity established by signals that pattern the neural tube along its anterior-posterior axis.



### 2.3. SHH and Developmental Defects in Humans

Disruption of the normal midline patterning in the forebrain in humans results in a condition known as holoprosencephaly (HPE). In alobar HPE, the most severe form of HPE, the forebrain fails to separate into left and right hemispheres, and these patients have cyclopia and severe facial anomalies. The developmental defects observed in humans who have HPE are similar to those observed in SHH-null mice, raising the possibility that some cases of HPE may result from genetic alterations that disrupt SHH signaling. Subsequently, it has been shown that mutations in the SHH gene itself are responsible for some cases of autosomal-dominant HPE. Most of the mutations are point mutations that result in premature truncation of the protein or alteration of highly conserved residues, suggesting that this form of HPE results from haploinsufficiency of the SHH gene.

### 2.4. SHH and Somite Patterning

The skeletal muscle of the body is derived from segmented arrays of mesodermal structures called somites. Signals from the notochord and neural tube pattern each somite into two compartments: the ventral sclerotome, which gives rise to the vertebrae and ribs, and the dorsal dermomyotome, which gives rise to the dermis and the myotome. SHH can substitute for notochord-derived signals in promoting the expression of sclerotome markers in explanted somites, and in combination with WNT (a secreted factor that is present in the dorsal neural tube, see [Wingless Signaling](#)), it also promotes myogenesis in explanted somites. Three lines of evidence suggest that, instead of acting as an inductive signal for sclerotome and myoblast development in the somite, SHH is required for the survival and expansion of cells of the somitic lineages. First, cells in the somites undergo [programmed cell death](#) when the notochord and neural tube are surgically ablated (4). Secondly, small numbers of myoblasts and sclerotome cells still develop in SHH-null mice (3). Thirdly, SHH stimulates the proliferation of sclerotome cells and myoblasts (5, 6).

### 2.5. SHH and Limb Development

Vertebrate limbs develop from the limb buds, which are outgrowths of the body wall. Each limb bud is patterned along three axes, proximal-distal, dorsal-ventral, and anterior-posterior. Patterning along the anterior-posterior axis is controlled by the zone of polarizing activity (ZPA), a group of mesenchymal cells located in the posterior limb bud. When the ZPA is grafted onto the anterior side of a limb bud, it induces a mirror-image duplication of the digits. SHH is expressed in the cells of the ZPA, and SHH protein mimics the limb-duplicating activity of the ZPA. The use of SHH to pattern the vertebrate limb is a striking example of conservation in evolution—hh also controls anterior-posterior patterning in the developing wing and leg in *Drosophila*.

### 2.6. SHH and Long-Range Signaling

Induction of the floor plate by SHH is a short-range interaction that requires contact with the SHH-expressing cells of the notochord. The induction of motoneurons and sclerotome, however, occurs over a longer distance and does not require contact with the notochord. Because SHH is tethered to the membrane via its cholesterol anchor, it is not clear how it mediates these long-range effects. One possibility is that SHH induces the production of a secondary signaling molecule, which then diffuses and acts at a distance to pattern tissues. In *Drosophila*, the long-range patterning effects of the hh protein are mediated by the induction of *decapentaplegic*, a fly TGF $\beta$  homologue. Although there is no evidence that motoneuron induction in the neural tube occurs via the production of a secondary signaling molecule, there is a correlation in many other tissues of the developing embryo between the expression of SHH and the expression of bone morphogenetic proteins (BMPs), secreted signaling molecules that belong to the TGF $\beta$  superfamily. SHH and BMPs are expressed in adjacent domains in the lung, gut, tooth, hair, and other tissues, and overexpression of SHH induces BMP expression in the developing gut. This suggests that some of the patterning activities of SHH may be mediated by BMPs.

## 3. SHH Signal Transduction and Cancer

Although most of the information about the signaling pathway activated by SHH has come from genetic analysis of *Drosophila* development, the fundamental aspects of the pathway are conserved

in vertebrates. The SHH receptor is Patched (PTC), a transmembrane protein with 12 transmembrane segments that sits in a complex at the cell membrane with Smoothed (SMO), a seven-pass transmembrane protein. According to the current model, PTC is required to inhibit the activity of SMO, thereby repressing the expression of HH-induced target genes. SHH binding to PTC relieves the inhibition of SMO and allows it to activate downstream components of the SHH signal cascade, which include members of the GLI family of **zinc-finger** transcription factors. A paradoxical feature of SHH signaling is that it also up-regulates PTC expression. One explanation for this cellular response is that, in addition to regulating the activity of SMO, PTC may also sequester SHH and, therefore, limit its range of activity. Animals that have targeted disruptions of the PTC gene (7) have severe neural tube defects that include an expansion of the floor plate. This suggests that PTC activity is required to regulate SHH signaling during vertebrate development.

In addition to its patterning activities during embryogenesis, SHH also plays a role in growth control in a number of tissues, including the lung, sclerotome, retina, myoblasts, and skin. Several lines of evidence suggest that inappropriate activation of the SHH signaling pathway plays a role in human tumorigenesis. First, SHH target genes GLI and PTC are expressed constitutively in a number of human basal cell carcinomas (BCCs), the most common tumor type in humans, suggesting that the SHH signal cascade is activated in these tumor cells. Secondly, overexpression of SHH in the skin of transgenic mice is sufficient to induce BCCs. Thirdly, mutations in the PTC gene are associated with the autosomal-dominant condition basal cell nevus syndrome (BCNS): in addition to developmental defects, patients with BCNS have a high incidence of BCCs and medulloblastomas. Fourth, mutations in PTC and SMO (8) have been found in a high proportion of spontaneous BCCs and medulloblastomas. Taken together, these data suggest that components of the SHH signaling cascade are important targets for genetic alterations in human cancer.

#### Bibliography

1. J.A. Porter, K.E. Young, and P.A. Beachy (1996) *Science* **274**, 255–259.
2. J.A. Porter et al. (1996) *Cell* **86**, 21–34.
3. C. Chiang, Y. Litngtung, E. Lee, K.E. Young, J.L. Corden, H. Westphal, and P.A. Beachy (1996) *Nature* **383**, 407–413.
4. M.A. Teillet, Y. Watanabe, P. Jeffs, D. Duprez, F. Lapointe, and N.M. Le Douarin (1998) *Development* **125**, 2019–2030.
5. C.-M. Fan, J.A. Porter, C. Chiang, D.T. Chiang, P.A. Beachy, and M. Tessier-Lavigne (1995) *Cell* **81**, 457–465.
6. D. Duprez, C. Fournier-Thibault, and N. Le Douarin (1998) *Development* **125**, 495–505.
7. L. Goodrich, L. Milenkovic, K. Higgins, and M.P. Scott (1997) *Science* **277**, 1109–1113.
8. J. Xie, M. Murone, S-M. Luoh, A. Ryan, Q. Gu, C. Zhang, J.M. Bonifas, C.-W. Lam, M. Hynes, A. Goddard, A. Rosenthal, E.H. Epstein, Jr., and F.J. de Sauvage (1998) *Nature* **391**, 90–92.

#### Suggestions for Further Reading

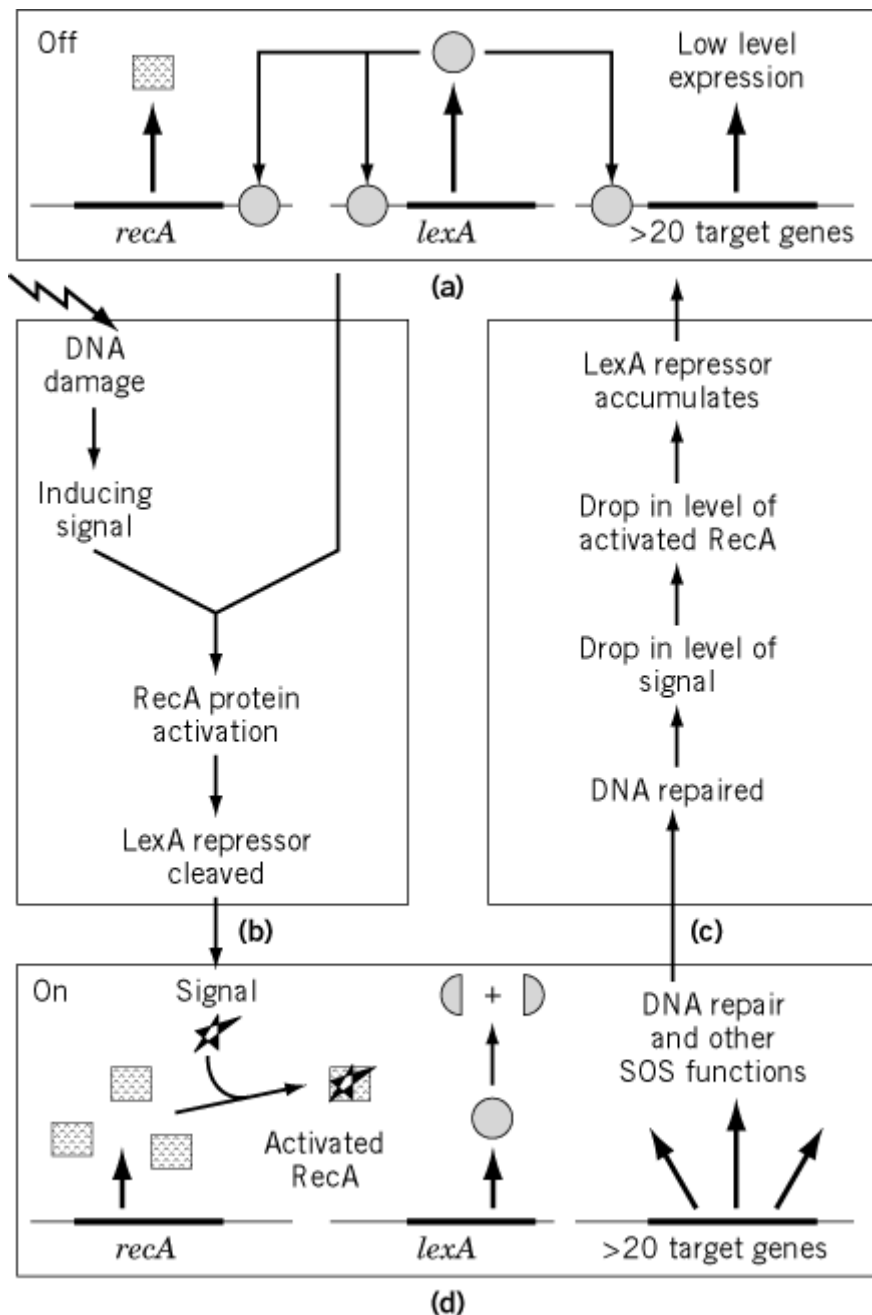
9. M. Hammerschmidt, A. Brook, and A.P. McMahon (1996) The world according to hedgehog, *Trends Genet.* **13**, 14–21
10. Y. Tanabe and T.M. Jessell (1996) Diversity and pattern in the developing spinal cord, *Science* **274**, 1115–1123.
11. R.L. Johnson and C.J. Tabin (1997) Molecular models for vertebrate limb development, *Cell* **90**, 979–990.
12. E.J. Robertson (1997) Left-right asymmetry, *Science* **275**, 1280.
13. P.D. Currie, and P.W. Ingham (1998) The generation and interpretation of positional information within the vertebrate myotome, *Mech. Dev.* **73**, 3–21
14. P.W. Ingham (1998) The patched gene in development and cancer, *Curr. Opin Genet. Dev.* **8**, 89–94

15. J.E. Ming and M. Muenke (1998) Holoprosencephaly: from Homer to Hedgehog, *Clin. Genet.* **53**, 155–163.
16. S. Pfaff and C. Kintner (1998) Neuronal diversification: Development of motor neuron subtypes, *Curr. Opin. Neurobiol.* **8**, 27–36.

## SOS Response

The SOS response is a global **gene** regulatory system used by **bacteria**, such as *Escherichia coli*, to respond to treatments that damage **DNA** or inhibit [DNA replication](#). During this response, about twenty genes are expressed at increased rates. These genes are individually termed *SOS genes* and collectively the *SOS regulon*. Their expression is controlled by the interplay of two regulatory proteins: the [LexA repressor](#), which inhibits expression of the SOS genes during normal cell growth (see also [Repressors](#)), and the **RecA** protein, which is activated by treatments that turn on the SOS response (Fig. 1). Activation of RecA leads in turn to inactivation of LexA by specific **proteolytic** cleavage (see [LExA Repressor](#)). It is thought that expression of the SOS genes helps the cell counteract the damage to the DNA, leading to repair of the damage and resumption of normal cell growth.

**Figure 1.** The SOS regulatory system. **(a)** State of the system during normal cell growth. The LexA protein is active and represses synthesis of RecA (left) and the SOS proteins (right). **(b)** Induction and transition to the induced state caused by DNA damage. The activated RecA protein causes the LexA protein to cleave itself into two pieces, which inactivates it as a repressor. **(c)** Induced SOS state. In the absence of active LexA, the *recA* and SOS genes are expressed in large amounts. If the cell contains a  $\lambda$  prophage, a prolonged stay in this state results in prophage induction. **(d)** Transition to the normal growth state. Activation of RecA is reversible, so that DNA repair leads to loss of RecA activation. The state of the system and the transitions between the two states are controlled by the level of RecA coprotease activity. Modified from Little and Mount (1982) with permission.



Although DNA is chemically stable, it can be modified by various agents in ways that alter its base-pairing properties (see [Mutagenesis](#)). Cells have a variety of mechanisms for repairing DNA (see [DNA Repair](#)). These mechanisms are expressed in *E. coli* at low levels that suffice to repair sporadic damage, but when the level of damage becomes high enough to threaten cell survival, the SOS system is turned on. Induction occurs within a few minutes, and the cells begin to recover from small amounts of DNA damage within half an hour. Many of the SOS genes play direct roles in DNA repair. Others, it is thought, aid cell survival in more indirect ways. For still others, no function is known, and they may play roles in cell survival under conditions different from those normally used to study *E. coli* in the laboratory.

The SOS response can be regarded as a [signal transduction](#) pathway, in which an ordered set of events occurs in response to a signal. In this particular system, the central part of the pathway is the best understood aspect. The nature of the signals that activate RecA and the functions of the SOS gene products are less clear. For this reason, the regulatory proteins are discussed first.

## 0.1. Regulatory Proteins

RecA is a complex protein that plays central roles in genetic [recombination](#) and DNA repair in addition to its regulatory role. From *in vitro* studies it is believed that the activated form of RecA is a helical filament containing many RecA molecules polymerized on single-stranded DNA and containing ATP or dATP. This filament also plays crucial roles in recombination and repair.

LexA is a repressor similar to phage  $\lambda$  repressor. It contains an N-terminal DNA-binding domain and a C-terminal domain involved in forming dimers. The protein binds to sites on DNA with dyad symmetry, and the bound form is a dimer. The specific proteolytic cleavage site lies between the two domains of the protein. Therefore, cleavage separates the DNA-binding function from the dimerization function, and DNA binding is inhibited because the protein cannot dimerize. The proteolytic reaction is unusual biochemically in that cleavage is a self-processing reaction carried out by groups in LexA, not in RecA. Activated RecA serves as a “co-protease” to stimulate this reaction (see [LExA Repressor](#)).

If the cell contains a **prophage** for phage  $\lambda$  or a related bacterial virus (see [Lambda Phage](#)), the CI repressor of the prophage can undergo the same type of cleavage reaction, again stimulated by activated RecA. The result of this reaction is termed prophage induction. The viral lytic genes, normally controlled by CI, are then derepressed, and the virus proceeds to grow lytically, releasing a burst of progeny virus. It is believed that the responsiveness of the CI repressor to activated RecA is an adaptation of the virus to a signal that the cell is in trouble and may not survive. In this view, the cell has not developed a regulatory mechanism for the benefit of the virus. Rather, the virus puts an existing cellular mechanism to its own uses. Cleavage of the  $\lambda$  repressor occurs at a much slower rate than that of LexA. Accordingly, if the DNA damage is not severe, cleavage of CI does not proceed to completion before the DNA damage is repaired (see Fig. 1), and prophage induction does not occur. This differential sensitivity of LexA and phage repressors reflects their differing regulatory niches. LexA has evolved to respond rapidly even to small inducing treatments, whereas  $\lambda$  has evolved to respond slowly and only to levels of damage that threaten cell survival. Indeed, prophage induction becomes efficient only at doses of DNA damage that begin killing the cells. Another difference between the two systems is that SOS induction is reversible, in that the cells return to the normal growth state. By contrast, prophage induction is irreversible because the complicated regulatory circuitry of the  $\lambda$  genetic switch prevents further CI expression, even if cleavage ceases.

## 0.2. Inducing Treatments and Inducing Signals

A variety of treatments induce the SOS response. DNA-damaging treatments include irradiation with ultraviolet light and treatment with reactive chemicals, such as [mitomycin C](#) and methyl methane sulfonate. These agents have in common their ability to modify the DNA bases so that they cannot form **Watson–Crick base pairs**. Other treatments that interfere with DNA replication are also effective inducing treatments. Examples are **nalidixic acid**, which inhibits DNA gyrase (see [DNA Topology](#)), starvation of a thymine auxotroph for thymine, and exposure of many conditional replication mutants to restrictive conditions. All of these inducing treatments converge on RecA protein, which can be regarded as a sensor of the state of the DNA.

It is believed that one important (but perhaps not the only) way that RecA is activated *in vivo* is a consequence of the mechanism of DNA replication. When a [replication fork](#) encounters an abnormal, noncoding base in the template strand, further extension of the growing strand is blocked. A new [Okazaki fragment](#) begins at a distance, typically thought to be a few hundred nucleotides, downstream of the damage, leaving a gap of single-stranded DNA, termed a daughter-strand gap. One likely way that RecA is activated then is by binding to this gap. One prediction of this model is that if replication were inhibited in a way that does not produce gaps, DNA damage would not induce the SOS response. Various lines of *in vivo* evidence are consistent with this prediction (1). Moreover, it is thought that inducing treatments that do not damage DNA directly (such as nalidixic acid treatment) also generate single-stranded DNA. However, it remains possible that other types of inducing signals also activate RecA.

### 0.3. Functions of SOS Gene Products

Many SOS gene products are directly involved in [DNA repair](#) . These include the excision repair proteins UvrA and UvrB and proteins involved in daughter-strand gap repair, including RuvA, RuvB, and RecA itself. Other SOS gene products play less direct roles in DNA repair. The Sula protein is a cell division inhibitor. While the SOS response is activated, cells cannot divide. It is likely that this helps daughter-strand gap repair by keeping several copies of the chromosome in the same cell, thereby providing templates for damaged molecules.

The SOS gene products UmuD and UmuC play central roles in the process of SOS mutagenesis. This process is believed to be a “court of last resort,” in which irreparable damage is converted to a readable sequence, usually producing errors. UmuD protein is processed after synthesis to its active form, termed UmuD'. Once again, this processing reaction is a specific proteolytic cleavage reaction entirely parallel to those that inactivate LexA and  $\lambda$  repressor. It differs, however, in that the cleavage product is activated, rather than inactivated, by the cleavage. UmuD' and UmuC work together with RecA, acting in still another role in DNA metabolism. It is postulated that these proteins act to reduce the proofreading activity of **DNA polymerase III**.

Still other SOS genes code for functions that are not understood. This may be a consequence of the fact that *E. coli* is not typically studied in its natural environment. Laboratory conditions (exponential growth in rich media) are chosen for their reproducibility and convenience, but in nature the organism is generally found in the gut or in dilute aqueous environments, such as streams. It is likely that some functions that promote cell survival after DNA damage in these environments differ from each other and from laboratory conditions.

### Bibliography

1. M. Sassanfar and J. W. Roberts (1990) *J. Mol. Biol.* **212**, 79–96.

### Suggestions for Further Reading

2. E. C. Friedberg, G. C. Walker, and W. Siede (1995) *DNA Repair and Mutagenesis*, ASM Press, Washington, D.C., Chapters "10" and "12". Excellent and comprehensive.
3. J. W. Little and D. W. Mount (1982) The SOS regulatory system of *Escherichia coli*. *Cell* **29**, 11–22. Still a clear and concise summary of the development of the SOS model.

### Southern Blots (DNA Blots)

In 1975, Edwin Southern published the method for transferring electrophoretically resolved DNA [restriction fragments](#) from [agarose](#) gels to [nitrocellulose](#) membrane filters to be subsequently **hybridized** with DNA probes (1). This procedure has since been termed Southern blotting (or DNA blotting) and is a milestone in molecular biology (2). Southern blotting provides the connection between a specific nucleotide sequence (in a known or synthesized DNA fragment) and a resolved DNA molecule, a possibility that has been central to the discovery and molecular analysis of genes and their **transcripts** and the structural elements that regulate them (cross reference).

Typically, DNA fragments are resolved in agarose gels, which are then treated with alkaline buffers to **denature** the DNA prior to transfer to the [blotting matrix](#). Then the gel is neutralized in mildly acidic buffers and is prepared for blotting. Traditionally the transfer is accomplished via convection (see [Blotting](#)), where the gel is placed on a paper wick that draws buffer from a reservoir (3, 4). A

blotting matrix (originally a nitrocellulose membrane filter, although nylon membranes are often a preferred alternative, see [Blotting Matrices](#)) is applied to the surface of the gel and a stack of absorbent paper (paper towels) is placed on top of the matrix. A weight is added to ensure uniform pressure over the surface of the blot. Upon completion of transfer (commonly performed overnight), the DNA is fixed to the nitrocellulose membrane by baking at 80°C, which should be conducted in a vacuum oven to prevent igniting the matrix (with nylon, UV [cross-linking](#) is often employed). Then the blot is blocked, ie, prehybridized, originally with a solution containing a mixture of irrelevant DNA (eg, fragmented salmon sperm), polyvinylpyrrolidone, bovine [serum albumin](#), Ficoll and a [detergent](#), although other mixtures have since been found equally effective, especially with nylon matrices. The purpose of this step is blocking all unoccupied areas of the matrix to minimize nonspecific interactions with the probe. Probing the blot is achieved by incubating it in a solution containing the radioactively labeled single-stranded DNA or RNA of interest, typically at 65°C for 6 to 12 h. After incubation, the blot is washed extensively and ultimately autoradiographed.

Numerous modifications and improvements to the traditional procedure described have been introduced over the years (5). The time and temperature of probing, the salt conditions, and the presence of various reagents, such as **formamide**, all affect the stability of the hybrid formed on the surface of the blot. Thus the stringency of the assay can be varied, from which the degree of [homology](#) between the probe and bound DNA is inferred. The blotting matrices used and the methods of transfer have also been modified, for example, electrotransfer of DNA from **polyacrylamide gels** to nylon membrane filters in low salt buffers, to be cross-linked to the matrix with UV light (as opposed to baking at 80°C) is still considered Southern blotting.

Characteristically, the blot can be erased after having been analyzed by boiling the filter in an appropriate solution to remove the bound probe and thus enable a second round of reaction with a different probe. This can be repeated a number of times, which conserves DNA samples and allows efficient comparisons of probes and repeated examination of rare samples.

#### Bibliography

1. E. M. Southern (1975) *J. Mol. Biol.* **98**, 503–517.
2. J. Meinkoth and G. Wahl (1984) *Anal. Biochem.* **138**, 267–284.
3. L. G. Davis, M. D. Dibner, and J. F. Battey (1986) *Basic Methods in Molecular Biology*, Elsevier, New York, pp. 62–65.
4. (1988) *Nucleic Acid Hybridization: A Practical Approach* (B. D. Hames and S. J. Higgins, eds.), IRL Press, Oxford, UK.
5. M. R. Evans, A. L. Bertera, and D. W. Harris (1994) *Mol. Biotechnol.* **1**, 1–12.

#### **Soybean Trypsin Inhibitor (Kunitz), STI**

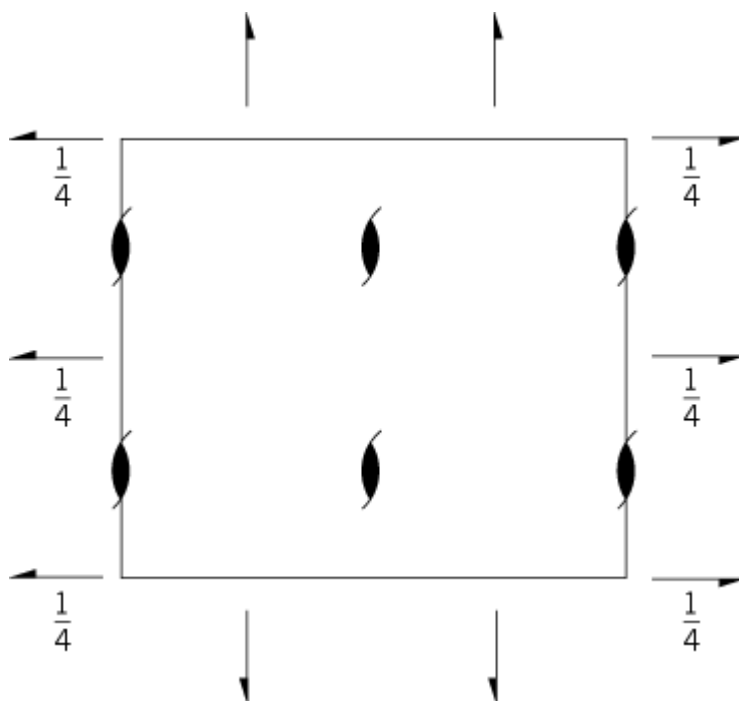
Soybean trypsin inhibitor (Kunitz) (see **Serine proteinase inhibitors** and **Proteinase inhibitor, protein**) is a storage protein in soybeans. It is a polypeptide of 181 residues crosslinked by two disulfide bridges. The Arg<sup>63</sup>-Ile bond serves as the reactive site for this powerful ( $K_1 = 1 \times 10^{-11}$  M) inhibitor of trypsin. Chymotrypsin is also inhibited at this reactive site, albeit more weakly ( $K_1 = 1 \times 10^{-5}$  M). Soybean trypsin inhibitor (Kunitz) gave its name to a large family of proteins. Many, but not all, of these are proteinase inhibitors. The use of the word *Kunitz* for a family of predominantly plant inhibitors, as well as its use for the important family of predominantly animal

inhibitors [see **BPTI (bovine pancreatic trypsin inhibitor) (Kunitz)**], attests to Moses Kunitz's great achievement. However, it also causes massive confusion. Plant scientists think of STI, but medical and animal scientists of BPTI. Omitting (Kunitz) from STI was suggested, but that would also lead to confusion as soybeans contain many other trypsin inhibitors, among them several members of the Bowman–Birk family (see **Bowman–Birk proteinase inhibitors**).

## Space Group

[Crystallography](#) is very important in molecular biology because [X-ray crystallography](#) can determine detailed and accurate structures of macromolecules. The molecules in a crystal are packed regularly throughout the crystal. This often leads to a symmetrical relationship among the molecules. Several types of symmetry elements occur: rotation axes, mirror planes, and centers of symmetry. In addition, the axes and mirror planes can be combined with a translation. The symmetry elements can be grouped in 230, and not more than 230, different ways, called space groups. The 230 space groups and their properties are tabulated in Ref. 1. As an example, one projection of the [unit cell](#) in space group  $P2_12_12_1$  is given in Fig. 1.

**Figure 1.** One projection of the unit cell in space group  $P2_12_12_1$ . The graphic symbols indicate twofold screw axes. Molecules are related to each other by a  $180^\circ$  rotation around the axis plus a translation along the axis over half of the unit cell. The horizontal screw axes are above the plane of the drawing at one-quarter of the unit cell.



## Bibliography

1. The International Union of Crystallography (1992) International Tables for Crystallography, Vol. A, Kluwer Academic, Dordrecht, Boston, London.



## Suggestions for Further Reading

2. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.
3. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists*, VCH, New York, Weinheim, Cambridge, Chap. "4".

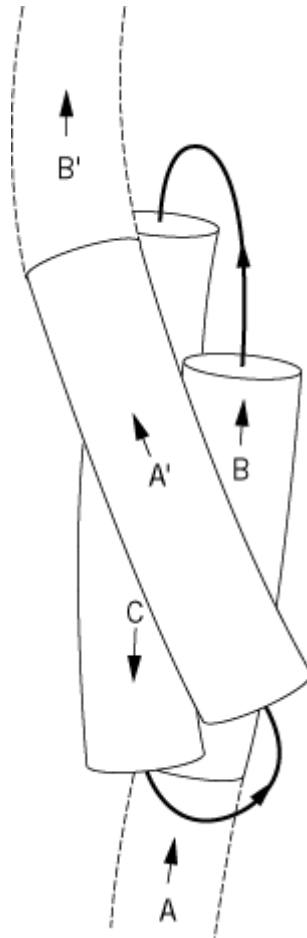
## Spectrin

Spectrin is a flexible and elongated intracellular structural protein, originally isolated from **erythrocytes**, but now known to be a characteristic of many other cell types too. It is a major component of a supportive meshwork lying in the cytoplasmic face of the plasma membrane and in all cases plays an important mechanical role in helping to specify cell shape and in providing [membrane](#) stability. Spectrin associates indirectly with the membrane through the interaction of [ankyrin](#) with Band 3 protein (the **anion exchange protein**); it also interacts with a variety of other proteins, including [actin](#). There are two types of spectrin molecule: (i) a-spectrin, which consists of a single a-chain of molecular weight about 240 kDa, and (ii) b-spectrin, which consists of a single b-chain of molecular weight about 230 kDa. Each has a length of about 100 nm and each can be conveniently divided into *N*-terminal, rod, and *C*-terminal **domains**. The *N*-terminal domain of b-spectrin (but not a-spectrin) contains an actin-binding site; the *C*-terminal domain in a-spectrin contains putative **calcium-binding** sites, and that in b-spectrin has an ankyrin-binding site. In between these end domains, the amino acid sequence has about 18 to 20 tandem repeats, each of length 106 residues. The repeat is predicted to be highly **a-helical**, an observation consistent with the overall observed a-helicity of the spectrin molecule (65% to 70%). It is these repeats that form the rod structure and account for the bulk of the observed molecular length. Approximately at the center of the rod domain is a region of sequence corresponding to part of the tenth repeat that is homologous to the *src*-family of **protein kinases** and to [phospholipase C](#). a- and b-spectrin form an antiparallel dimer *in vivo* in which the molecules are almost perfectly overlapped. These in turn aggregate via their end domains (the *N*-terminal domains of two a-spectrin molecules and two *C*-terminal domains of b-spectrin molecules) to generate a structure with actin-binding sites at each end separated by about 200 nm. The rod domain acts as a spacer between functionally important parts of the molecule (actin-binding site, Ca<sup>2+</sup> binding site) and appears to have much potential for flexing and kinking.

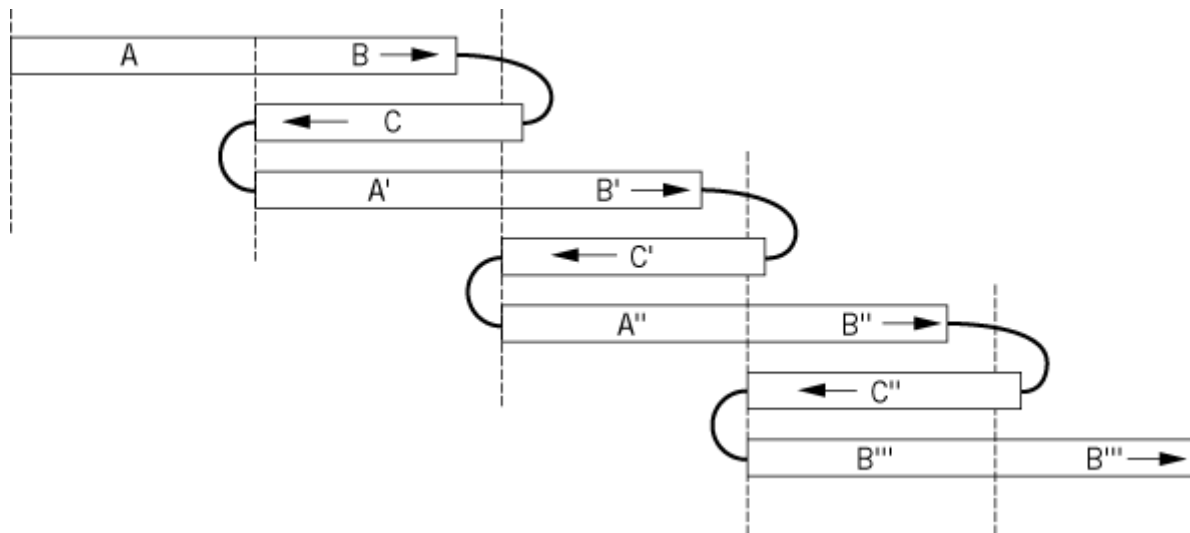
If the entire sequence of the rod domain were a-helical, the predicted length would be close to three times that observed. Furthermore, there is evidence to suggest that a structure comprising a long single a-helix would be unlikely to be stable in an aqueous environment. Consequently it was proposed that short non-a-helical regions (perhaps b- [turns](#)) occurred in each repeat and that these would allow the a-helices to form an antiparallel bundle from a single spectrin chain, thus giving the structure its stability (1, 2). In such a bundle, two of the a-helical segments would lie parallel to one another but antiparallel to the third. This structure is known as the *three-a-helix motif* (3). An important development occurred when Winograd et al. (4) expressed appropriately phased fragments of *Drosophila* a-spectrin from which they were able to isolate a highly a-helical segment of length 109 residues. This folded into a stable structure, with properties very similar to those of native a-spectrin, and was proposed to correspond to the structural repeat. Its structure determined by X-ray crystallography (5) and nmr solution studies (6) confirmed the model of Parry et al. (3) in which the three a-helices were arranged in a left-handed bundle with left-handed chain connectivity (Fig. 1). The axial repeat was close to 5.0 nm as predicted, and the motif was stabilized by both nonpolar and interchain [electrostatic interactions](#), again as proposed. Perfect [coiled-coil](#) packing, however, was not

maintained between two of the constituent  $\alpha$ -helices, even though the sequences of all three  $\alpha$ -helices exhibited the [heptad repeat](#) characteristic of such a conformation (5). Although there are elements of three antiparallel  $\alpha$ -helices within an axial distance of about 5 nm along the rod domain, the structural unit can be considered as one long  $\alpha$ -helix followed by a b-turn, followed by a short antiparallel  $\alpha$ -helix of about half the length and another b-turn (Fig. 2). A very similar conformation is shared by both [dystrophin](#) and  $\alpha$ -actinin.

**Figure 1.** Schematic diagram of a three- $\alpha$ -helix motif, showing the A, B, and C  $\alpha$ -helices forming part of a left-handed coiled-coil structure. The connectivity was deduced as being left-handed on the basis of the number of potential interhelix interactions. The subsequent structure of a dimer determined by X-ray crystallography and nmr solution studies confirmed all the main points of this model. (From Ref. 3, with permission.)



**Figure 2.** Schematic diagram of a portion of the rod domain of the spectrin superfamily of proteins. The repeat can be considered either as one  $\alpha$ -helix about 9 nm in length (A and B) followed by a turn and then a second shorter  $\alpha$ -helix about 5 nm long (C) followed by another turn, or as the three- $\alpha$ -helix motif illustrated in Fig. 1. It is predicted that the axial repeats will be different in spectrin and  $\alpha$ -actinin. (From Ref. 3, with permission.)



## Bibliography

1. M. Koenig, A. P. Monaco, and L. M. Kunkel (1988) The complete sequence of dystrophin predicts a rod-shaped cytoskeletal protein. *Cell* **53**, 219–228.
2. R. R. Dubreuil, T. J. Byers, A. L. Sillman, D. Bar-Zvi, L. S. B. Goldstein, and D. Branton (1989) The complete sequence of *Drosophila* alpha-spectrin: conservation of structural domains between alpha-spectrins and alpha-actinin. *J. Cell Biol.* **109**, 2197–2205.
3. D. A. D. Parry, T. W. Dixon, and C. Cohen (1992) Analysis of the three  $\alpha$ -helix motif in the spectrin superfamily of proteins. *Biophys. J.* **61**, 858–867.
4. E. Winograd, D. Hume, and D. Branton (1991) Phasing the conformational unit of spectrin. *Proc. Natl. Acad. Sci. USA* **88**, 10788–10791.
5. Y. Yan, E. Winograd, A. Viel, T. Cronin, S. C. Harrison, and D. Branton (1993) Crystal structure of the repetitive segments of spectrin. *Science* **262**, 2027–2030.
6. J. Pascual, M. Pfuhl, D. Walther, M. Saraste, and M. Nilges (1997) Solution structure of the spectrin repeat—a left-handed anti-parallel triple-helical coiled-coil. *J. Mol. Biol.* **273**, 740–751.

## Suggestion for Further Reading

7. V. Bennett and D. M. Gilligan (1993) The spectrin-based membrane skeleton and micron-scale organization of the plasma membrane. *Annu. Rev. Cell Biol.* **9**, 27–66.

## Spectroscopy

Spectroscopy is the collective name for a diverse group of experimental techniques in which the interactions of electromagnetic radiation with a sample of interest are studied. Spectroscopies are grouped into different categories depending on:

1. the energies that are involved,
2. which type of transitions between energy levels are triggered in a molecule by the incident radiation, and

3. whether the molecules of interest absorb the energy or emit it (usually after prior excitation). The energy  $E$  of an electromagnetic wave is proportional to its frequency and thus inversely proportional to its wavelength:

$$E = h\nu = hc/\lambda \quad (1)$$

where  $h$  is Planck's constant and  $c$  the velocity of light.

Most widely used are spectroscopic techniques in which the interaction of molecules with ultraviolet (UV) or visible (vis) light is studied. In this range, the energy of the incident radiation is equivalent to 150 to 400 kJ/mol, which is sufficient to promote electrons from the ground state to the first excited state. The absorbance of energy from the incident light as a function of its wavelength is measured in [absorption spectroscopy](#). In fluorescence and phosphorescence spectroscopy (see [Fluorescence Spectroscopy](#)), the return from the activated to the ground state is studied. This deactivation is not always accompanied by the emission of light. Other pathways for deactivation can exist in which energy can be exchanged with the surrounding solvent, in a radiationless process.

Optically active substances can absorb left- and right-handed circularly polarized light to different extents. This phenomenon is called [circular dichroism](#) (CD). CD spectroscopy is especially well suited to study the conformations and conformational transitions of proteins and nucleic acids.

Transitions between different vibrational states of a molecule require energies in the range of about 50 kJ/mol (see [Vibrational Spectroscopy](#)). They lead to absorption of radiation in the infrared (IR) region of the spectrum, which is measured in infrared (IR) spectroscopy. Information about vibrations in molecules can also be obtained by Raman spectroscopy. The analysis of Raman spectra can complement that from IR spectra. To produce a Raman band, the polarizability of the molecule must change, rather than the dipole moment in IR spectroscopy, so the magnitudes of the IR and Raman signals differ.

#### Suggestions for Further Reading

S. B. Brown (ed.) (1980) *An Introduction to Spectroscopy for Biochemists*, Academic Press, London. An excellent compilation of basic articles on the various spectroscopies.

C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part II, W. H. Freeman and Co., San Francisco, CA. Chapters 7–9 of this textbook cover advanced and theoretical aspects of spectroscopy.

## Sperm

The spermatozoa in a [male](#) ejaculate are called sperm. A spermatozoon is the [haploid](#), motile **gamete** of the male that is specialized in the transfer of genetic information and induction of the [egg](#) at [fertilization](#). Spermatozoa are the smallest cells in the animal kingdom. They are produced in the testis from precursor **diploid** spermatogonia. These cells are located close to the basal lamina and are supported by sertoli cells. During proliferation and growth of spermatogonia, they undergo meiosis, develop to spermatocytes, and migrate from the basal lamina to the inner lumen of the testis. Differentiated haploid spermatozoa are transported to the epididymis, where they undergo further maturation to gain full motility and fertilization capacity.

The mature sperm cell consists of three distinct parts: a head, a midpiece, and a tail. The head consists of **acrosomal** and postacrosomal regions. The acrosomal cap covers the acrosomal vesicle. The head has dimensions  $10 \times 3 \mu\text{m}$  in mammals and contains mostly DNA, which is heavily condensed. This is due to the replacement of [histones](#) by [protamines](#) during spermatogenesis (1). **Gene expression** from the haploid [genome](#) in the fully differentiated sperm head has not been reported. The DNA in the sperm head is present solely for the transport of genetic information, along with a minimum of other material.

The midpiece of the sperm contains numerous [mitochondria](#), which are closely attached to the dense fibers to provide energy for the movement of the **flagella**. These mitochondria do not become part of the new [zygote](#) and, consequently, mitochondria are inherited maternally exclusively. No [ribosomes](#), [endoplasmic reticulum](#), or [Golgi apparatus](#) are present in the midpiece or in the head. The proteins involved in sperm metabolism are the diazepam-binding **receptor** (2) and [dynein](#) (3). The sperm metabolism is also affected when deletions of the mitochondrial DNA are present, which appear to be associated with dead sperm in subfertile man.

The tail is important for the motility of the sperm. The flagella contain  $9+9+2$  dense fibers, of which the outer dense fibers are coded by *odf* genes (4-7). The *Fsc1* gene codes for a major **cytoskeletal** structure of the mammalian sperm flagellum (8). The ATP for motility is provided by the mitochondria in the midpiece.

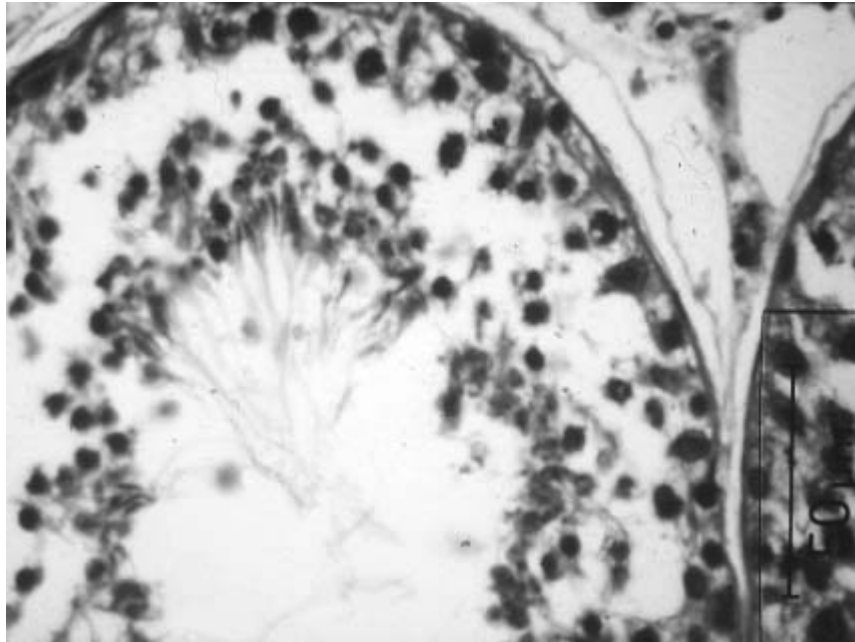
As spermatozoa are specialized to transfer genetic information and to induce the egg to promote the development of the new individual, they must have the ability to recognize specifically the **oocytes** to be fertilized. Adhesion molecules located on the surface of the acrosomal cap are called ligands and are responsible for species-specific interactions (see [Acrosome](#) and [Fertilization](#)). A plethora of adhesive proteins were found to interact with the zona pellucida of the oocyte. For instance, AQN and AWN, the so-called sperm adhesins (9), **beta-galactosidase**, a **trypsin**-sensitive site, and a number of proteins characterized immunochemically by their interaction with [monoclonal antibodies](#) were identified to be involved in binding to the egg. These are primary ligands (see [Acrosome](#)), most of which are expressed in the epididymis and then are associated with the acrosomal cap or the sperm head.

On contact with the zona pellucida of the egg, spermatozoa undergo an [exocytosis](#), which is called the acrosome reaction. This essential reaction during fertilization leads to the exposure of those proteins that are otherwise protected by the acrosomal cap. The main component of the acrosome is *proacrosin*, which exhibits high affinity for zona pellucida **glycoproteins** and, after activation to *acrosin*, acts as a [proteinase](#) to digest zona pellucida glycoproteins locally (see [Acrosome](#)).

## 1. Spermatogenesis

During embryogenesis, primordial gonads develop that harbor the **germ cells**. The primordial urogenital tract develops either to a testis or an ovum. The primordial germ cells, the ancestors of which are the cells located in the [polar plasma](#), increase in number by mitosis but, shortly after birth, mitosis comes to a standstill. From puberty on, because of an increase in testosterone levels, the type A stem spermatogonia start to divide again (Fig. 1). Diploid type A1 spermatogonia are produced, which undergo further rounds of mitotic divisions. After several rounds of mitosis, primary spermatocytes are formed, which enter a **tetraploid** phase. Subsequently, two cycles of meiotic division take place that are directed by the dose-dependent gene *roughex* (10) to form secondary spermatocytes and, subsequently, spermatids with a haploid genome. At the beginning of cell differentiation, the spermatids still have the shape of diploid cells. The process of histogenesis reshapes the spermatids. The change in morphology serves to produce cells that are specifically adapted to gene transfer. One feature of this specialization is the compaction and condensation of the [nucleus](#), which is directed by transition proteins and protamines (11).

**Figure 1.** Development of sperm cells. Type A spermatogonia are located adjacent to the interstitial cells. They divide mitotically until primary spermatocysts are formed. The subsequent reduction of the genome occurs by two divisions during meiosis. Spermatids are formed, which then differentiate to spermatozoa. During this process, developing sperm cells migrate to the lumen. In the center of the tubuli, spermatozoa are visible. Further maturation of spermatozoa occurs in the epididymis. (Slide, with friendly permission from A. Bielefeld.)



Gene expression during the entire process of spermatogenesis can be divided into those genes that are expressed at the diploid stage (eg, *proacrosin*) or at the haploid stage (transition protein 2 gene) (12). Cytoplasmic bridges connect cells with each other so that an interchange of all necessary RNA, proteins, and so forth, is possible. During histogenesis, a tail is formed from a distal **centriole**, and the Golgi apparatus converts to the acrosomal vesicle and to the acrosomal cap. Mitochondria of the cell are arranged close to the distal centriole to form the midpiece later. The amount of **cytoplasm** is constantly reduced, and fully differentiated spermatozoa do not have any cytoplasm at all. Sometimes a droplet of cytoplasm that migrated to the distal end of the sperm tail is visible in ejaculated spermatozoa that have not matured completely.

#### Bibliography

1. S. K. Choudhary, S. M. Wykes, J. A. Kramer, A. N. Mohamed, F. Koppitch, J. E. Nelson, and S. A. Krawetz (1995) *J. Biol. Chem.* **270**, 8755–8762.
2. M. Kolmer, M. Pelto-Huikko, M. Parvinen, C. Hoog, and H. Alho (1997) *DNA Cell Biol.* **16**, 59–72.
3. S. M. King, S. P. Marchese-Ragona, S. K. Parker, and H. W. Detrich 3rd (1997) *Biochemistry* **36**, 1306–1314.
4. Y. Kim, I. M. Adham, T. Haack, H. Kremling, and W. Engel (1995) *Chem. Hoppe. Seyler.* **376**, 431–435.
5. S. Burmester, and S. J. Hoyer-Fender (1996) *Mol. Reprod. Dev.* **45**, 10–20.
6. H. Brohmann, S. Pinnecke, and S. Hoyer-Fender (1997) *J. Biol. Chem.* **272**, 10327–10332.
7. X. Shao, H. A. Tarnasky, U. Schalles, R. Oko, and F. A. van der Hoorn (1997) *J. Biol. Chem.* **272**, 6105–6113.
8. K. D. Fulcher, C. Mori, J. E. Welch, D. A. O'Brien, D. G. Klapper, and E. M. Eddy (1995) *Biol. Reprod.* **52**, 41–49.

9. E. Töpfer-Petersen, and J. J. Calvete (1996) *J. Reprod. Fertil. Suppl.* **50**, 55–61.
10. P. Gonczy, B. J. Thomas, and S. DiNardo (1994) *Cell* **77**, 1015–1025.
11. S. M. Wykes, J. E. Nelson, D. W. Visscher, D. Djakiew, and S. A. Krawetz (1995) *DNA Cell Biol.* **14**, 155–161.
12. K. Nayernia, I. Adham, H. Kremling, K. Reim, M. Schlicker, G. Schlüter, and W. Engel (1996) *Int. J. Dev. Biol.* **40**, 379–383.

### Suggestion for Further Reading

13. J. P. Dadoune (1994) The cellular biology of mammalian spermatids: A review. *Bull. Assoc. Anat.* **78**, 33–40.

## Spermine, Spermidine

Spermine and spermidine are polyamines bound to **DNA** in semen and in other body tissues. Spermine is a tetramine,  $\text{H}_2\text{N}-[(\text{CH}_2)_3\text{-NH}]_3\text{-H}$ , and spermidine is a triamine,  $\text{H}_2\text{N}-[(\text{CH}_2)_3\text{-NH}]_2\text{-H}$ . The concentration of polyamines in the [nucleus](#) may modulate the binding of other proteins, such as [histones](#), to DNA and facilitates the progression of enzymes like **RNA polymerase** along the DNA **double-helix** in a [chromatin](#) environment.

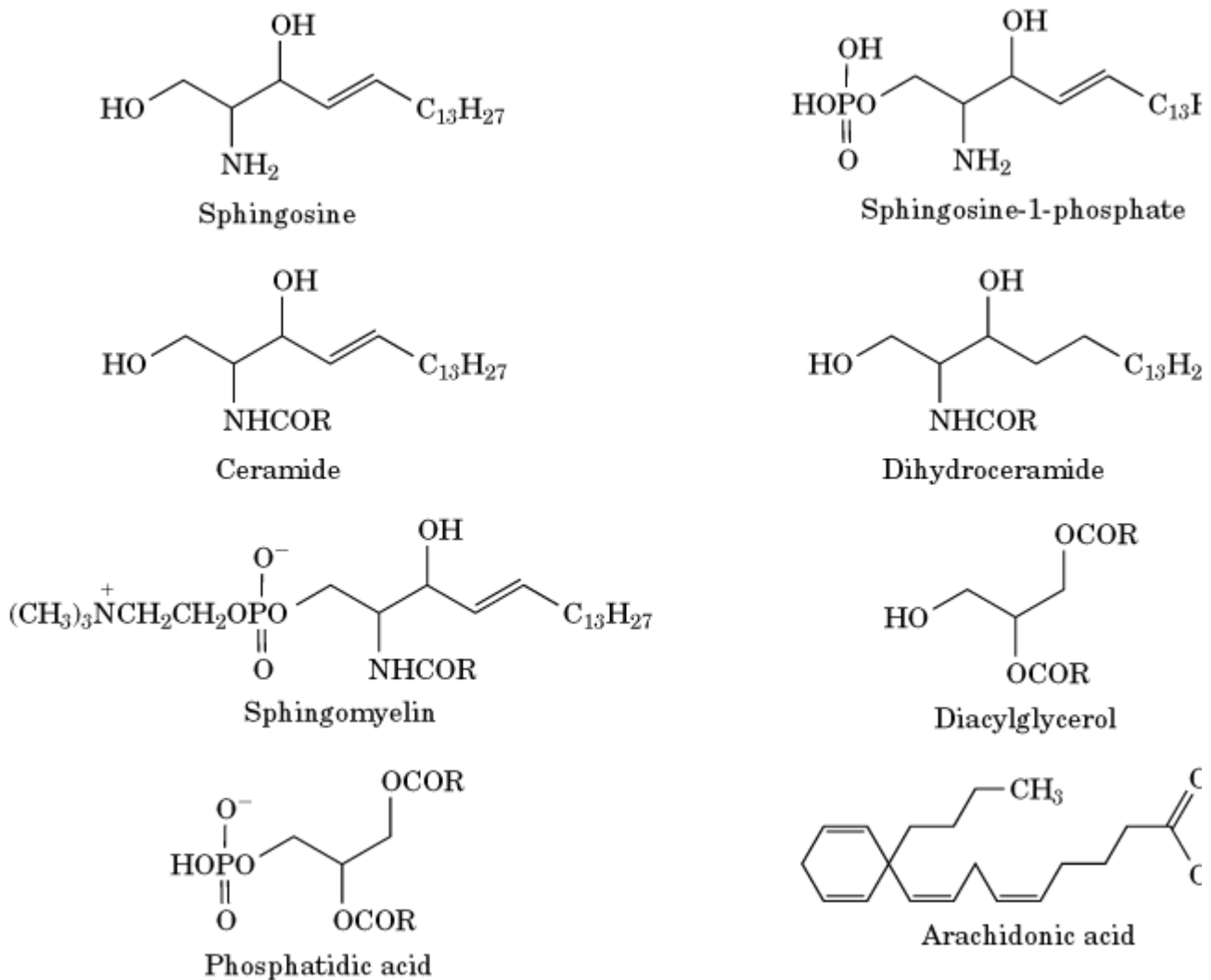
## Sphingomyelinases

Sphingolipids constitute a large family of molecules based on a sphingoid backbone (such as sphingosine). Sphingolipids include sphingomyelin, ceramide, glycolipids, sulfatides, gangliosides, and other derivatives of sphingosine and ceramide.

### 1. Sphingolipid Structure and Function

The typical sphingolipid is composed of a sphingoid base, an acyl chain in amide linkage, and a substituent at the 1-hydroxyl position (see Fig. [1](#) for some representative structures). Sphingolipids are known to participate in several functions of the cell. For example, certain glycolipids can function as **receptors** or components of receptors. Gangliosides and glycolipids play a role in cell recognition and cell–cell contact response. Many cell types display specific patterns of glycolipids and gangliosides, and these have been associated with specific stages of [development](#), [differentiation](#), or tumor progression ([1](#), [2](#)). More recent studies have identified emerging roles for sphingolipids in [signal transduction](#) in a manner analogous to the role of glycerolipids in the phosphatidylinositol (PI) cycle (see [Phospholipases C](#)) and other pathways of signal transduction ([3-5](#)).

**Figure 1.** Structures of selected sphingolipids. The symbol R represents long-chain alkyl groups that can be saturated or unsaturated.



## 2. Sphingomyelinases and Regulation of Ceramide Formation

Many extracellular agents and stimuli cause activation of sphingomyelinases and/or accumulation of endogenous levels of ceramide. These include cytokines (such as **tumor necrosis factor- $\alpha$** , the Fas ligands, and **nerve growth factor**), agents of differentiation or growth suppression (such as [glucocorticoids](#) and 1-25-dihydroxyvitamin D<sub>3</sub>), chemotherapeutic agents (such as cytosine arabinoside, vincristine, and daunorubicin), and other agents of stress or injury (such as heat and irradiation) (3).

At least two distinct sphingomyelinases have been implicated in response to one or more of the above agents. Acid sphingomyelinase (with a pH optimum in the acidic range) has been cloned and is known to be deficient in the sphingolipid storage disorder, Niemann–Pick disease (6, 7). A secreted form of the enzyme is zinc-dependent and has been suggested to play a role in atherosclerosis. A neutral and magnesium-dependent sphingomyelinase has been implicated in the response to most of the above inducers of sphingomyelin hydrolysis and ceramide formation (3).

## 3. Ceramide and Sphingolipids in the Eukaryotic Stress Response

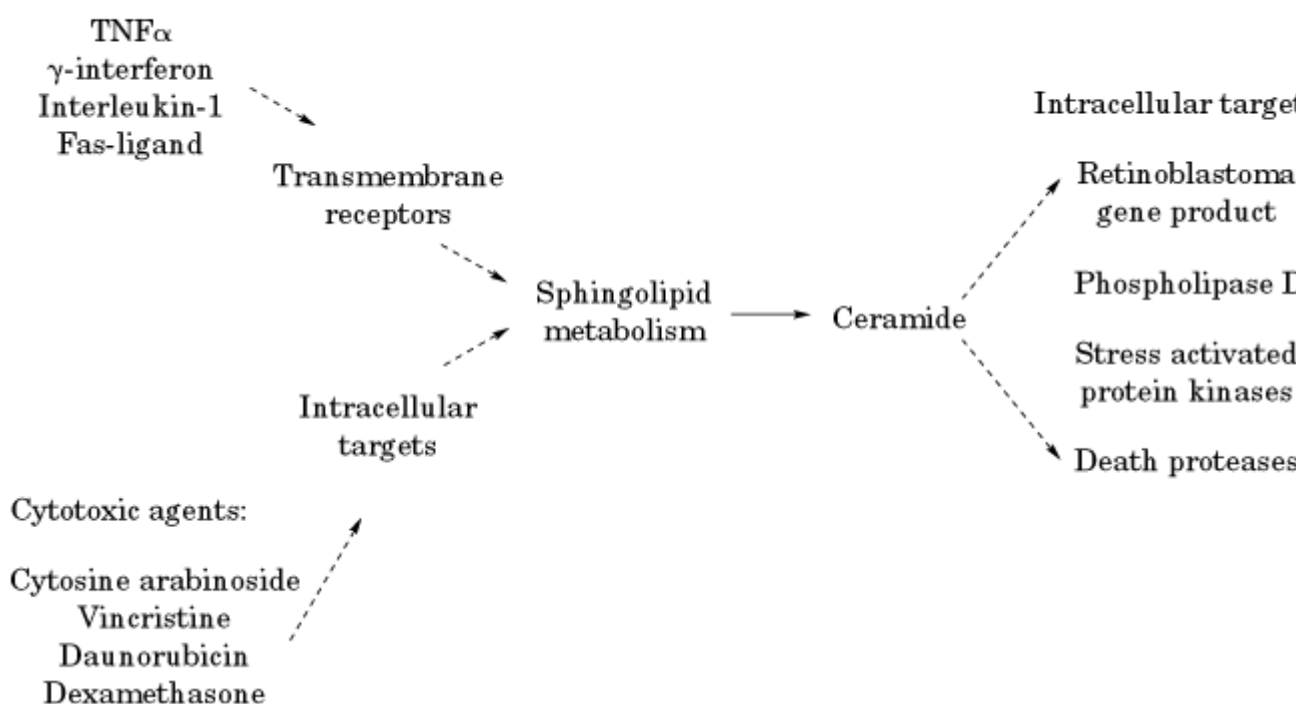
Ceramide formed in response to these agents participates in the **stress response** of cells (Fig. 2). The addition of cell-permeable analogs of ceramide or the modulation of endogenous ceramide



metabolism has been shown to regulate programmed cell death ([apoptosis](#)), [cell cycle](#) arrest, terminal cell [differentiation](#), cell senescence, and inflammatory responses ([8-11](#)). In the cell, ceramide can regulate protein **phosphatases** and [kinases](#), cell death [proteinases](#), the **retinoblastoma** gene product, phospholipase D, and other elements critical in the response of cells to **mitogenic** or growth suppressor stimuli ([12](#)).

**Figure 2.** Scheme of the involvement of sphingolipids in stress responses. Various cytokines and cytotoxic agents modulate and cause the accumulation of ceramide, which in turn serves to regulate a number of cellular processes related to growth responses.

**Stress cytokines:**



Sphingolipids are essential for the viability of **yeast**, as determined by the deletion of the first and rate-limiting enzyme of yeast sphingolipid biosynthesis, serine–palmitoyl transferase. Suppressors of this mutation result in viable cells that are defective in their response to environmental stresses such as heat, osmolarity, and acid pH ([13](#)). This can be overcome by the addition of specific sphingolipids. These results have been taken to suggest an important role for yeast sphingolipids in the response of yeast to various stress stimuli.

These studies demonstrate the significance of sphingolipids in stress responses of both “lower” and “higher” eukaryotic cells.

**Bibliography**

1. H. Wiegandt (1985) in *Glycolipids*, H. Wiegandt, ed., Elsevier, New York, pp. 199–259.
2. S. Hakomori (1981) *Annu. Rev. Biochem.* **50**, 733–764.
3. Y. A. Hannun (1994) *J. Biol. Chem.* **269**, 3125–3128.
4. S. Spiegel and A. H. Merrill, Jr. (1996) *FASEB J.* **10**, 1388–1397.
5. Y. H. Zhang and R. Kolesnick (1995) *Endocrinology* **136**, 4157–4160.
6. R. O. Brady (1983) in *The Metabolic Basis of Inherited Disease*, J. B. Stanbury, J. B.

Wyngaarden, D. S. Fredrickson, J. L. Goldstein, and M. S. Brown, eds., McGraw-Hill, New York, pp. 831–841.

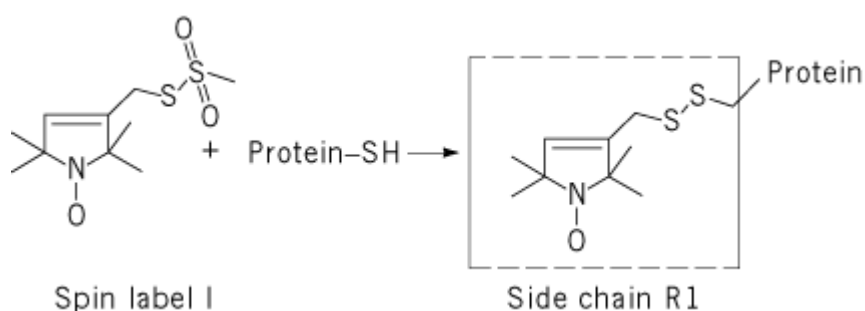
7. L. Quintern, E. H. Schuchman, O. Levrán, M. Suchi, K. Ferlinz, H. Reinke, K. Sandhoff, and R. J. Desnick (1989) *EMBO J.* **8**, 2469–2473.
8. L. M. Obeid, C. M. Linardic, L. A. Karolak, and Y. A. Hannun (1993) *Science* **259**, 1769–1771.
9. T. Okazaki, R. M. Bell, and Y. A. Hannun (1989) *J. Biol. Chem.* **264**, 19076–19080.
10. M. E. Venable, J. Y. Lee, M. J. Smyth, A. Bielawska, and L. M. Obeid (1995) *J. Biol. Chem.* **270**, 30701–30708.
11. L. R. Ballou, C. P. Chao, M. A. Holness, S. C. Barker, and R. Raghov (1992) *J. Biol. Chem.* **267**, 20044–20050.
12. Y. A. Hannun (1996) *Science* **274**, 1855–1859.
13. J. L. Patton, B. Srinivasan, R. C. Dickson, and R. L. Lester (1992) *J. Bacteriol.* **174**, 7180–7184.

## Spin Labeling

### 1. Spin Labeling

Unpaired electrons in molecules can be probed using **electron paramagnetic resonance (EPR)** techniques, but nature does not always provide the necessary free-radical at a site of interest. (1) Spin labels are stable, paramagnetic molecules that can be easily tailor-made by organic synthesis to provide selective structural and dynamic information via EPR spectroscopy (2, 3). An advantage of spin-label EPR studies is that they can be carried out at physiologically relevant temperatures. Nitroxide-labeled **nucleic acids** (4), **lipids** (5) and enzyme substrates (6) are among the varieties of spin label analogs that have been made. A molecular biology-based approach, called site-directed spin labeling (SDSL), for studying the dynamics, folding, and structure of **proteins**, has been developed (2, 3). In SDSL, a reactive nitroxide (for example, a methanethiosulfonate) reacts with a **cysteine** side chain that has been placed at one or more selected sites in a protein by **site-directed mutagenesis**. Figure 1 shows the reaction scheme most often used in the SDSL approach. Spin labels are also used in nuclear magnetic resonance (NMR) studies for drug discovery (7).

**Figure 1.** Introduction of a site-directed spin label. A methanethiosulfonate spin label reacts with the natural, or genetically engineered, sulfhydryl side chain of a protein to give a site-directed, spin-labeled protein. (Redrawn with permission from reference 2, Fig. 1.)



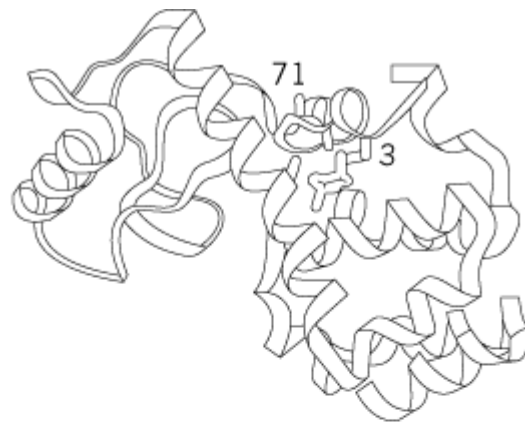
### 2. Dynamics of Proteins, DNA, and Membranes

The effects of motion of a spin labeled macromolecule can be observed indirectly as averaging of features in the EPR spectrum or directly by time domain EPR experiments. The effects of motion on a spin label EPR spectrum arise because the unpaired electron in the nitroxide molecule is in an asymmetric environment that renders its interaction with an external magnetic field sensitive to their relative orientations. When the spin label is rigidly attached to a larger molecule, the EPR spectra will reflect the motion of the macromolecule. The time-scale of motions that can be detected using spin-label EPR is determined by the frequency of the particular spectrometer used and the details of the detection scheme. Motions from milliseconds to fractions of a nanosecond have been measured. Applications include resolving spectra from distinct structural states of myosin in muscle (8), detecting time-dependent structural changes in bacteriorhodopsin (9), studies of the persistence length for DNA bending (4), and measurements of dynamics of lipids in [membranes](#) (5).

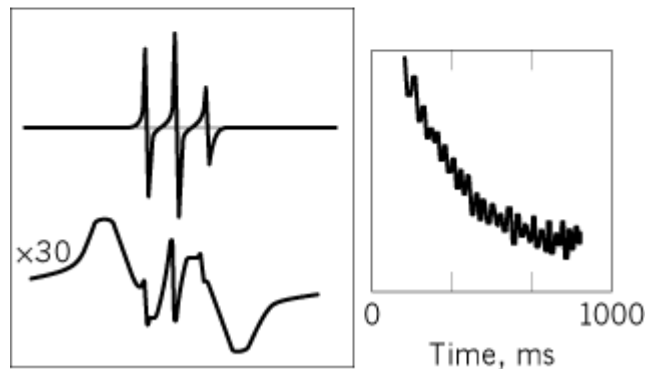
### 3. SDSL

This topic has been touched on in the preceding paragraphs (2, 3). The varied applications include determination of protein **secondary structure**, in which stepwise substitution of nitroxides on a protein structural element and magnetic effects of oxygen or paramagnetic ions in solution are used to determine surface exposure, and measurement of distances between two amino acid residues, in which two spin label sites are probed by EPR for magnetic interactions between them. Several experimental approaches can be taken to determine distances between spin labels on a macromolecule (2, 10). A measurement of inter-residue distance is illustrated in Figure 2. The inset (b) shows EPR spectra of folded, di-spin-labeled T4 [lysozyme](#) (bottom) and of the same protein unfolded in urea (top). The overall structure of the protein shown in (a) indicates the close proximity of the two spin label side chains. Spin labels separated by 7 to 24Å give broadened spectra from which the inter-residue distance can be determined. Once a protein derivative has been prepared with two interacting spin labels, spectral changes may be employed to measure the kinetics of refolding, as illustrated in (c). Inter-residue distances in [membrane proteins](#), such as the lactose **permease** (11), can be estimated by the complementary techniques of **cysteine**-scanning mutagenesis and paired site-directed spin labeling.

**Figure 2.** Spin labeling of phage T4 lysozyme. The structure of a mutant form of T4 lysozyme in which two free cysteine residues were introduced (Ile3Cys and Val71Cys) and allowed to react with spin labels after reaction with spin labels is shown Figure 2a. The proximity of the spin labels leads to the broad electron paramagnetic resonance (EPR) spectrum shown in the inset (b, bottom). The EPR spectrum of the same sample unfolded in urea is sharp and shows little interaction between spin labels (b, top). The time-resolved rate of refolding can be followed by EPR (c). (Reproduced with permission from reference 2, Fig. 5.)



(a)



(b)

(c)

## Bibliography

1. H. M. McConnell (1971) *Annu. Rev. Biochem.* **40**, 227–236.
2. W. L. Hubbell, H. S. Mchaorab, and C. Altenbach (1996) *Structure* **4**, 779–783.
3. W. L. Hubbell, D. S. Cafiso, and C. Altenbach (2000) *Nature Struct. Biol.* **7**, 735–739.
4. B. H. Robinson, C. Mailer, and G. Drobny (1997) *Annu. Rev. Biophys. Biomol. Struct.* **26**, 629–658.
5. D. Marsh and L. I. Horváth (1998) *Biochim. Biophys. Acta* **1376**, 267–296.
6. T. Palm, C. Coan, and W. E. Trommer (2001) *Biol. Chem.* **382**, 417–423.
7. W. Jahnke, L. B. Perez, C. G. Paris, A. Strauss, G. Fendrich, and C. M. Nalin (2000) *J. Am. Chem. Soc.* **122**, 7394–7395.
8. J. E. Baker, I. Brust-Mascher, S. Ramachandran, L. E. W. LaConte, and D. D. Thomas (1998) *Proc. Natl. Acad. Sci.* **95**, 2944–2949.
9. R. Mollaaghababa, H. J. Steinhoff, W. L. Hubbell, and H. G. Khorana (2000) *Biochemistry* **39**, 1120–1127.
10. M. Persson, J. R. Harbridge, P. Hammarström, R. Mitri, L.-G. Mårtensson, U. Carlsson, G. R. Eaton, and S. S. Eaton (2001) *Biophys. J.* **80**, 2886–2897.
11. J. Wu, J. Voss, W. L. Hubbell, and H. R. Kaback (1996) *Proc. Natl. Acad. Sci.* **93**, 10123–10127.

## Suggestions for Further Reading

12. *Biological Magnetic Resonance, Vol. 14, Spin Labeling. The Next Millennium* (L. J. Berliner,

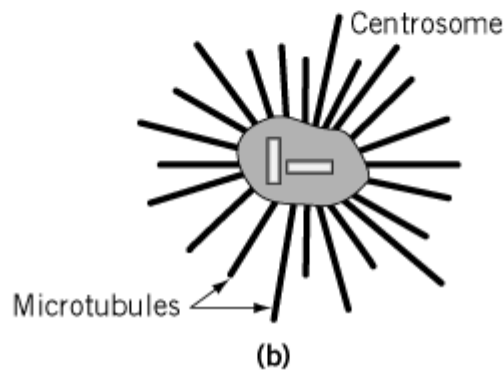
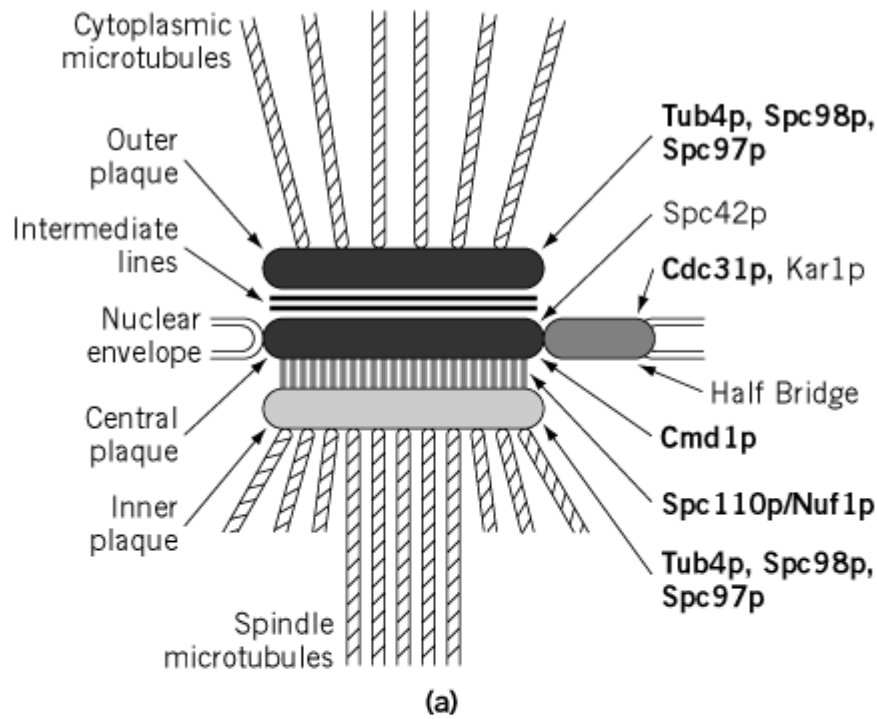
ed.), Plenum Press, New York, 1998.

## Spindle Pole Body

In yeast, [microtubules](#) are organized by the spindle pole body (SPB). The SPB is the fungal equivalent of the **centrosome**, and both are microtubule organizing centers (MTOC). Microtubules provide the physical means by which [chromosomes](#) are segregated during mitosis, and by which nuclei fuse during mating. SPBs organize the nuclear microtubules of the **mitotic spindle** and the cytoplasmic microtubules.

Since the baker's yeast *Saccharomyces cerevisiae* has been developed into a model genetic system, much has been learned about the SPB. The SPB is physically quite distinct from the animal centrosome, but the functional conservation is intriguing (Fig. 1). The centrosome consists of two centrioles surrounded by pericentriolar material, whereas the SPB is a multilayered cylindrical structure. Nevertheless, as more is learned about both MTOCs, some general themes emerge; there may be many more similarities than might be inferred from electron micrographs of these structures. Because *S. cerevisiae* is the model system in which the SPB has been most extensively studied, this article will emphasize what is known about the SPB in this yeast.

**Figure 1.** A diagram of the SPB of *S. cerevisiae* (**a**) with an inset diagram of the centrosome for comparison (**b**). For simplicity, the only proteins included are those both shown to localize specifically to the SPB (and not to the spindle microtubules near the SPB) and mentioned in the text. SPB proteins for which there exists a homologue in the centrosome are in **bold** (reviewed in (3, 17, 18)). Although the SPB has been studied by electron microscopy for many years, recent work has clarified some elements of SPB structure (the intermediate lines) (19). In the centrosome diagram (**b**), the gray represents the pericentriolar material; the inner rectangles are centrioles.



## 1. Structure of the Spindle Pole Body

The *S. cerevisiae* SPB is an electron-dense, short, cylindrical structure composed of several layers (Fig. 1). It is a permanent structure: it is always embedded in the [nuclear envelope](#), and no nuclear envelope breakdown occurs during mitosis in *S. cerevisiae*. This is in contrast with most other eukaryotes, where nuclear breakdown is concomitant with mitosis. Although permanence of the nuclear envelope is a common feature of yeast mitosis, a permanent association of the SPB with the nuclear envelope is not. In *Schizosaccharomyces pombe*, for example, the SPB is not permanently associated with the nuclear envelope.

The central plaque of the SPB is co-planar with the nuclear envelope and is interposed between the inner and outer plaques, from which the nuclear and cytoplasmic microtubules, respectively, emanate throughout the [cell cycle](#) (see Fig. 1). Separated by the nuclear envelope, the cytoplasmic and nuclear arrays are always separated compartmentally. Early in the cell cycle, a single SPB exists within the nuclear envelope.

Components of the SPB have been identified through biochemical, genetic, and direct cytologic approaches. Through biochemical approaches, some of the proteins that comprise the layers of the SPB have been identified (for example, Spc98p, Spc42p, and Spc110p/Nuf1p) (1-3). Some proteins

are thought to play a structural role, whereas others are more actively involved in association with microtubules. The essential phosphoprotein Spc42p, for example, is present at the central plaque. In [conditional-lethal mutant](#) strains containing a defective Spc42p, loss of viability occurs at the restrictive temperature as SPB duplication is attempted. SPB duplication is unsuccessful, and dense material collects above the half-bridge, suggesting that Spc42p is essential during SPB duplication, possibly for attachment of the SPB to the nuclear envelope. Overexpression of Spc42p results in a greatly lengthened central plaque, due to formation of a polymer of Spc42p (4).

Fibers connecting the inner and central plaques are composed of homodimers of Spc110p/Nuf1p, a [coiled-coil](#) protein. The coiled-coil encoding region of the SPC110/NUF1 gene is largely dispensable, and deletions of it lessen the distance between the central and inner plaques in a manner proportional to the length of the deletion. This distance between the inner and central plaques, although not essential for cell growth, was suggested to lend flexibility to the SPB, which might be important for microtubule array orientation after SPB duplication (5). Alternatively, such an element could allow regulatory molecules to access the dense SPB during SPB duplication or microtubule nucleation.

## 2. g-Tubulin and Microtubule Nucleation and Organization at the Spindle Pole Body

The major function of an MTOC is the nucleation and organization of microtubules. Microtubules do not polymerize spontaneously in vivo: The levels of [tubulin](#) within the cell are below the critical concentration necessary for spontaneous microtubule nucleation and assembly. Microtubules are nucleated in vivo at MTOCs, which obviate this requirement for a critical tubulin concentration. g-Tubulin, termed Tub4p in *S. cerevisiae*, functions specifically in microtubule nucleation. First discovered in *Aspergillus nidulans* (6), g-tubulin has been shown to localize to the MTOC and to be central to microtubule assembly. g-Tubulin is **homologous** to both a- and b-tubulin, but it is not incorporated into the length of the microtubule. Instead, it is a specialized tubulin that binds to the minus end of microtubules and is important for microtubule nucleation. In *S. cerevisiae*, Tub4p interacts with Spc98p and Spc97p, proteins that have homologues in higher eukaryotes (7-9). Together, these proteins form a complex in the cytoplasm, which then assembles at the inner and outer plaques (10, 11).

The organization of microtubules at the SPB is highly regulated. Because tubulin exchange occurs at both ends of the microtubule, the SPB presumably controls the activity at the minus end (the end proximal to the SPB). It has been suggested that g-tubulin has a role in depolymerization at the minus end of the microtubule, as a mutant g-tubulin gene in *A. nidulans* was able to suppress a microtubule hyperstability phenotype conferred by certain b-tubulin mutations. These g-tubulin mutants were proposed to do this by increasing the rate of depolymerization at the minus end of microtubules as a consequence of tubulin heterodimer removal (12).

## 3. SPB Duplication

Because SPB duplication marks the beginning of mitotic spindle formation, the timing of mitosis depends on the timing of this duplication. Although several events in SPB duplication have been identified by **electron microscopic** analysis (13, 14), and proteins that function in this process have been determined, the mechanism of this duplication is still unclear. It is known that early in the G1 stage of mitosis the SPB is not duplicated and is flanked by a “satellite” that is thought to be the precursor of the new SPB. This satellite abuts the cytoplasmic side of the half-bridge, a thick portion of the nuclear envelope next to the SPB (Fig. 1). After passage through Start (see [Cell Cycle](#)), two SPBs appear, separated by a full bridge. How this occurs is an active area of research. Several proteins that localize by **immunofluorescence microscopy** to the half-bridge cause a block in SPB duplication when defective, supporting the hypothesis that the half-bridge is important for duplication (reviewed in (14)).

SPB duplication is known to be a highly regulated process. SPB duplication must occur exactly once

per cell division cycle. Although many SPB proteins appear to be regulated by **phosphorylation**, the exact functions of these phosphorylations are a continuing area of study. Mps1p is a [kinase](#) essential for SPB duplication (15). There are other likely means of regulation of SPB duplication. For example, Pcs1p, a [proteasome](#) component, is thought to be important for degrading some component prior to SPB duplication (16).

#### 4. The SPB and the Centrosome

Despite the morphologic differences between the SPB and centrosome, research reveals an increasing number of protein homologues that are shared between these two organelles (Fig. 1). It appears that elements of the SPB that are closest to the minus ends of the microtubules are those most likely to have animal cell counterparts. Tub4p, Spc97p, Spc98p, and Spc110p/Nuf1p, among other proteins, all have homologues in centrosomes. In addition, regulatory proteins such as centrin and [calmodulin](#) exist in both the yeast SPB and the centrosome. It is unknown how extensive the analogy between SPBs and centrosomes will be, but it is remarkable that so many elements of the MTOC are so highly conserved in organisms as disparate as fungi and animals.

#### Bibliography

1. M. P. Rout and J. V. Kilmartin (1990) *J. Cell Biol.* **111**, 1913–1927.
2. M. P. Rout and J. V. Kilmartin (1991) *Cold Spring Harbor Symp. Quant. Biol.* **61**, 687–692.
3. P. A. Wigge, O. N. Jensen, S. Holmes, S. Soues, M. Mann and J. V. Kilmartin (1998) *J. Cell Biol.* **141**, 967–977.
4. A. D. Donaldson and J. V. Kilmartin (1996) *J. Cell Biol.* **5**, 887–901.
5. J. V. Kilmartin, S. L. Dyos, D. Kershaw and J. T. Finch (1993) *J. Cell Biol.* **123**, 1175–1184.
6. C. E. Oakley and B. R. Oakley (1989) *Nature* **338**, 662–664.
7. O. C. Martin, R. N. Gunawardane, A. Iwamatsu and Y. Zheng (1998) *J. Cell Biol.* **141**, 675–687.
8. S. M. Murphy, L. Urbani and T. Stearns (1998) *J. Cell Biol.* **141**, 663–674.
9. A.-M. Tassin, C. Celati, M. Moudjou and M. Bornens (1998) *J. Cell Biol.* **141**, 689–701.
10. S. Geissler, G. Pereira, A. Spang, M. Knop, S. Soues, J. Kilmartin and E. Schiebel (1996) *EMBO J.* **15**, 3899–3911.
11. M. Knop, G. Pereira, S. Geissler, K. Grein and E. Schiebel (1997) *EMBO J.* **16**, 1550–1564.
12. B. R. Oakley (1994) In *Microtubules* (J. S. Hyams and C. W. Lloyd, eds.), Wiley-Liss, Inc., New York, pp. 33–45.
13. B. Byers and L. Goetsch (1974) *Cold Spring Harbor Symp. Quant. Biol.* **38**, 123–131.
14. M. Winey and B. Byers (1993) *Trends Genetics* **9**, 300–304.
15. E. Weiss and M. Winey (1996) *J. Cell Biol.* **132**, 111–123.
16. H. B. McDonald and B. Byers (1997) *J. Cell Biol.* **137**, 539–553.
17. T. Stearns and M. Winey (1997) *Cell* **91**, 303–309.
18. S. G. Sobel (1997) *Journal of Experimental Zoology* **277**, 120–138.
19. E. Bullitt, R. P. Rout, J. V. Kilmartin and C. W. Akey (1997) *Cell* **89**, 1077–1086.

#### Suggestions for Further Reading

20. D. Botstein, D. Amberg, J. Mulholland, T. Huffaker, A. Adams, D. Drubin and T. Stearns (1997) "The yeast cytoskeleton". In *The Molecular and Cellular Biology of the Yeast Saccharomyces: cell cycle and cell biology* (J. R. Pringle, J. R. Broach and E. W. Jones, eds.), Cold Spring Harbor, New York, pp. 72–78. This text has a comprehensive section on the cytoskeleton of yeast.
21. A. D. Donaldson and J. V. Kilmartin (1996) Spc42p: a phosphorylated component of the *S. cerevisiae* spindle pole body (SPB) with an essential function during SPB duplication. *J. Cell Biol.* **5**, 887–901. This work is an elegant example of how a SPB gene is studied and how its



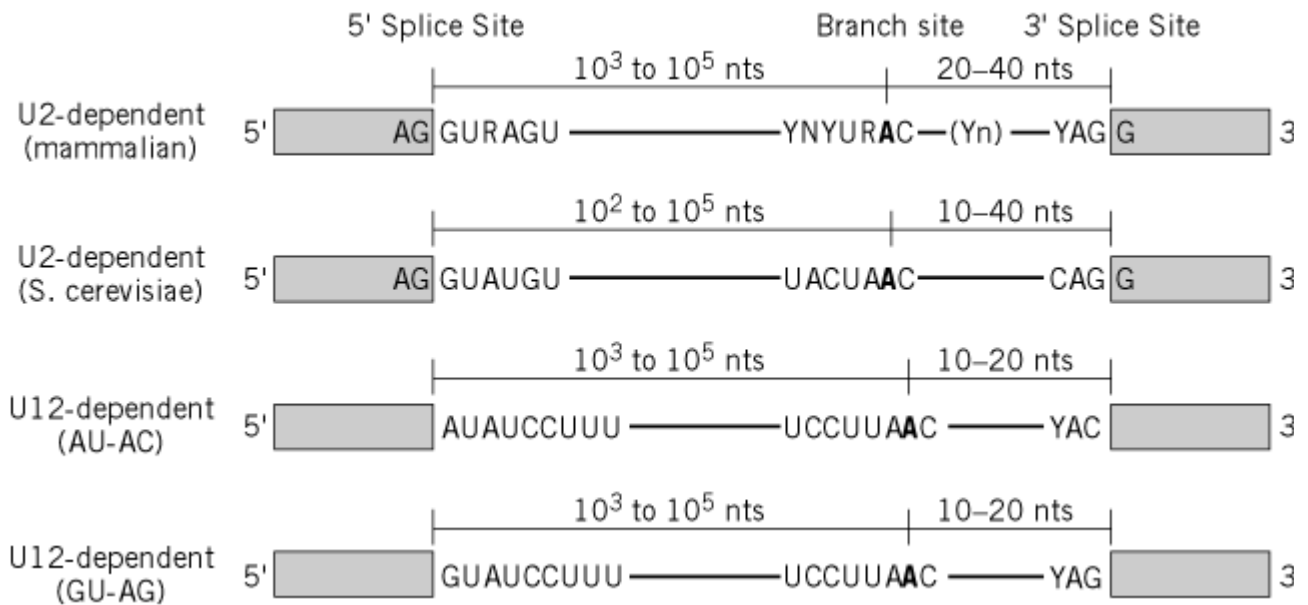
function is elucidated.

22. S. G. Sobel (1997) Mitosis and the spindle pole body in *Saccharomyces cerevisiae*. *Journal of Experimental Zoology* **277**, 120–138. A general review containing detailed information on the *S. cerevisiae* SPB as well as related topics such as mitosis and microtubule structure.
23. T. Stearns and M. Winey (1997) The cell center at 100. *Cell* **91**, 303–309. This is an up-to-date view of SPB and centrosome research that contrasts these two structures.
24. P. A. Wigge, O. N. Jensen, S. Holmes, S. Soues, M. Mann and J. V. Kilmartin (1998) Analysis of the *Saccharomyces* spindle pole by matrix-assisted laser desorption/ionization (MALDI) mass spectrometry. *J. Cell Biol.* **141**, 967–977. This exciting work utilizes MALDI for the identification of SPB proteins and shows EM images of the SPB.
25. M. Winey and B. Byers (1993) Assembly and function of the spindle pole body in budding yeast. *Trends Genetics* **9**, 300–304. This paper reviews SPB duplication in *S. cerevisiae*.

## Splice Sites

Most eukaryotic **genes** are **transcribed** as [messenger RNA](#) precursors that contain noncoding regions called **introns**. Nuclear pre-mRNA introns are excised and the coding regions (exons) are ligated together by a process called [RNA splicing](#). The exact phosphodiester bonds to be cleaved and ligated (ie, the so-called 5' and 3' [splice sites](#), which are also referred to as the 5' donor or 3' acceptor sites, or left and right splicing junctions, respectively) are determined primarily by short conserved sequences at the exon–intron junctions (1). By comparing the sequences of nuclear pre-mRNA introns, **consensus sequences** surrounding the 5' and 3' splice sites have been identified. Based on the nature of these consensus sequences and the type of [spliceosome](#) that assembles on them, nuclear pre-mRNA introns have been divided into two classes: U2-dependent, which represent the vast majority of introns and are found in all eukaryotes, and U12-dependent, which are less abundant and absent from some eukaryotes (eg, the yeast *Saccharomyces cerevisiae*) (2). In the case of U2-dependent introns, mammalian 5' splice sites are marked by the consensus sequence AG/GURAGU while those in the yeast *S. cerevisiae* are marked by AG/GUAUGU (where/ indicates the splice site, R = purine, Y = pyrimidine, and N = any nucleotide) (see Fig. 1). The 3' splice site can be divided into three distinct sequence elements that are typically found within the first 40 nucleotides upstream of the actual exon–intron junction (see Fig. 1). These include the so-called branch site, which is characterized by the consensus sequence YNYUAC in mammals and UACUAAC in yeast (where A represents the site of branch formation during pre-mRNA splicing), and the actual 3' splice junction, which conforms to the consensus sequences YAG/G and CAG/G in mammals and yeast, respectively. In addition, many vertebrate U2-dependent introns contain a 10- to 20-nucleotide stretch of pyrimidines (designated the polypyrimidine tract) that is found just upstream of the 3' splice junction. These splice-site consensus sequences are highly conserved in yeast, whereas in mammals the consensus represents the most prevalent sequence, but it often fluctuates considerably from one intron to the next. However, the terminal intron dinucleotides (GU and AG) are almost absolutely conserved in both higher and lower eukaryotes; greater than 99.9% of U2-dependent introns follow this so-called GU–AG rule. Mutations in these and other highly conserved nucleotides of both splice site sequences often lead to the activation of nearby sites (so-called cryptic splice sites) that are only used in the absence of the authentic site. Furthermore, splice site mutations account for a significant portion (up to 15%) of human genetic diseases that result from single point mutations (3).

**Figure 1.** Consensus sequences of nuclear pre-mRNA introns. 5' and 3' splice-site and branch-site consensus sequences are shown for U2-dependent introns in mammals and in the yeast *S. cerevisiae* and for U12-dependent introns, which are divided into (AU-AC) and (GU-AG) subtypes. R indicates a purine, Y a pyrimidine, and N any nucleotide. The branch point adenosine is indicated by bold lettering. In mammalian U2-dependent introns, a weakly conserved, 10 to 20-nucleotide-long pyrimidine tract (Yn) is found between the branch site and 3' splice site. Exons are represented by shaded boxes and introns by lettering or a solid line. The intron lengths and distances from the branch point to the 3' splice site, in nucleotides (nts), represent typical values for the majority of the specified introns. Due to limited information, the indicated length of U12-dependent introns is at present only an estimated value.



The recently identified minor class of nuclear pre-mRNA introns (ie, U12-dependent introns) contain distinct 5' and 3' splice-site consensus sequences, which are also summarized in Fig. 1. U12-dependent introns consist of two subtypes, AU-AC and GU-AG, which differ from one another solely in their terminal dinucleotides (4, 5). Both subtypes, which are estimated to comprise less than 0.2% of all nuclear pre-mRNA introns, possess strongly conserved consensus sequences at their 5' splice site, branch site, and 3' splice site and appear to lack a polypyrimidine tract (2, 6). The fact that these sequences are more highly conserved than the consensus sequences present in mammalian U2-dependent introns suggests that they play a greater role in directing the sites of [spliceosome](#) formation and thus ultimately the precise sites of cleavage and ligation during splicing.

The conserved sequences that define the 5' and 3' splice sites are recognized in a highly ordered manner by multiple factors during spliceosome assembly and the catalytic steps of splicing. For example, the 5' splice-site consensus is bound initially by the U1 snRNP and subsequently by the U5 and U6 snRNPs, whereas the branch site is initially recognized by the splicing factor SF1/BBP and then interacts with the U2 snRNP. These and subsequent interactions ultimately lead to the formation of a catalytically active spliceosome that is responsible for the cleavage/ligation reactions of splicing. Due to their limited length and high degree of degeneracy in higher eukaryotes, however, the 5' and 3' consensus sequences do not contain enough information to distinguish authentic splice sites from other pre-mRNA sequences that mimic them. This is particularly a problem in higher eukaryotes where intron lengths can exceed  $10^5$  nucleotides and most pre-mRNAs contain multiple introns; yeast pre-mRNAs, on the other hand, contain at most a single intron whose length is typically much less than  $10^3$  nucleotides (7). How the correct 5' and 3' splice sites are first selected from a large number of potential candidates and then correctly paired is currently poorly understood (see also [Alternative Splicing](#)). In single-intron pre-mRNAs where the intron length is relatively short (ie, the

majority of pre-mRNAs used in *in vitro* splicing systems), splice-site recognition and pairing are thought to be mediated by interactions occurring at the early stages of spliceosome assembly between factors bound at the 5' splice site and the downstream 3' splice site (ie, across an intron) (see [Spliceosome](#)). In pre-mRNAs that possess multiple and relatively long introns, 5' and 3' splice-site recognition has been proposed initially to involve interactions that occur across an exon (reviewed in Ref. 8). This so-called exon definition model stems from the observation that the vast majority of exons in vertebrates are fewer than 300 nucleotides long. This would limit the number of competing nonauthentic splice sites located between the 3' splice site and 5' splice site flanking the exon and would ensure a relatively constant, short distance between them. Furthermore, this model is consistent with the observation that factors bound at a 5' splice site can enhance the interaction of those binding to an upstream 3' splice site (ie, across an exon). Because the splice sites within the same intron are not initially paired in this model, interactions between the proper pair of 5' and 3' splice sites have been proposed to occur at a later stage of spliceosome assembly, for example during B complex formation when the U6 snRNA interacts both with the 5' splice site and with factors bound at the branch site (ie, the U2 snRNA) (see [Spliceosome](#)). While the exon definition model provides a general mechanism for this process, many aspects of splice-site recognition and pairing remain unclear and are thus the focus of much discussion and investigation in the nuclear pre-mRNA splicing field.

### Bibliography

1. M. J. Moore, C. C. Query, and P. A. Sharp (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 303–358.
2. P. A. Sharp and C. B. Burge (1997) *Cell* **91**, 875–879.
3. M. Krawczak, J. Reiss, and D. N. Cooper (1992) *Hum. Genet.* **92**, 41–54.
4. Q. Wu and A. R. Krainer (1997) *RNA* **3**, 586–601.
5. R. C. Dietrich, R. Incorvaia, and R. A. Padgett (1997) *Mol. Cell* **1**, 151–160.
6. S. L. Hall and R. A. Padgett (1994) *J. Mol. Biol.* **239**, 357–365.
7. P. E. Hodges, M. Plumpton, and J. D. Beggs (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, U.K., pp. 213–241.
8. S. M. Berget (1995) *J. Biol. Chem.* **270**, 2411–2414.

### Suggestions for Further Reading

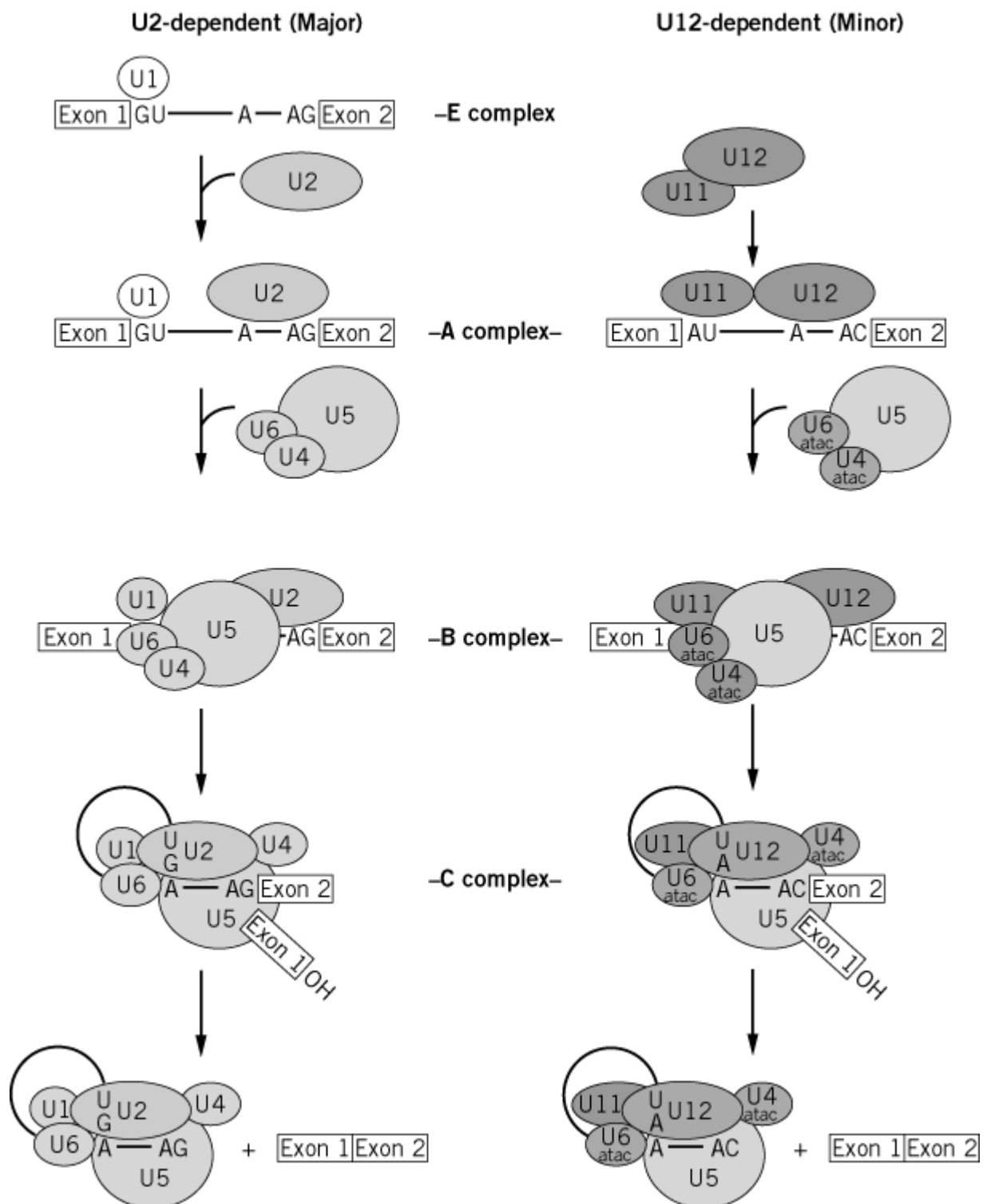
9. C. B. Burge, T. H. Tuschl, and P. A. Sharp (1998) "Splicing of Precursors to mRNAs by the Spliceosomes". In *The RNA World II* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, in press.
10. R. Reed (1996) Initial splice site recognition and pairing during pre-mRNA splicing. *Curr. Opin. Genet. Dev.* **6**, 215–220.

## Spliceosome

[RNA splicing](#) of nuclear pre-mRNA is catalyzed by a large (50 S to 60 S) multicomponent ribonucleoprotein (RNP) complex termed the *spliceosome*. Spliceosomes are assembled stepwise on the pre-mRNA and consist of the UsnRNPs [see [Small Nuclear Rnps \(Snrmps\)](#)] and numerous non-snRNP proteins. During spliceosome assembly, the pre-mRNA is folded so that its reactive groups (eg, the 5' splice site and branch site adenosine) are brought into close proximity. In this way, the

[active site](#)(s) responsible for the catalysis of splicing are formed. The majority of nuclear pre-mRNA introns are removed by the so-called major or U2-dependent spliceosome, which is comprised of the U1, U2, U5, and U4/U6 snRNPs. U12-dependent introns, which represent less than 0.2% of all nuclear pre-mRNA introns and have to date been identified in only a limited number of eukaryotes, are excised by the minor or U12-dependent spliceosome (1-3). The latter contains the U11, U12, U5, and U4atac/U6atac snRNPs. Spliceosome formation is a highly ordered process that requires ATP at all but the earliest stages of assembly. The assembly pathways of the major and minor spliceosomes are shown in Figure 1. The general assembly pathway of the U2-dependent spliceosome is identical in metazoans and yeast; due to the apparent lack of U12-dependent introns in the budding yeast *S. cerevisiae*, however, the latter does not possess a U12-dependent spliceosome (3). Because little is currently known about the U12-dependent spliceosome, we focus primarily on the mammalian U2-dependent spliceosome and refer to yeast only on those occasions where the two systems differ significantly.

**Figure 1.** Stepwise assembly of the U2- and U12-dependent spliceosomes. Only those steps that can be resolved by biochemical methods (eg, native gel electrophoresis or gel filtration) under normal conditions with mammalian splicing extracts are shown. For the sake of simplicity, the ordered interactions of the snRNPs (indicated by ellipses), but not those of non-snRNP proteins, are shown. The various spliceosomal complexes are named according to the metazoan nomenclature. Although not yet identified, a post-splicing complex (containing only the excised intron) similar to that formed upon dissociation of the U2-dependent spliceosome is depicted for the U12-dependent spliceosome. Exon and intron sequences are indicated by boxes and lines, respectively. The first two and last two intron nucleotides, as well as the branch site adenosine, are also shown.



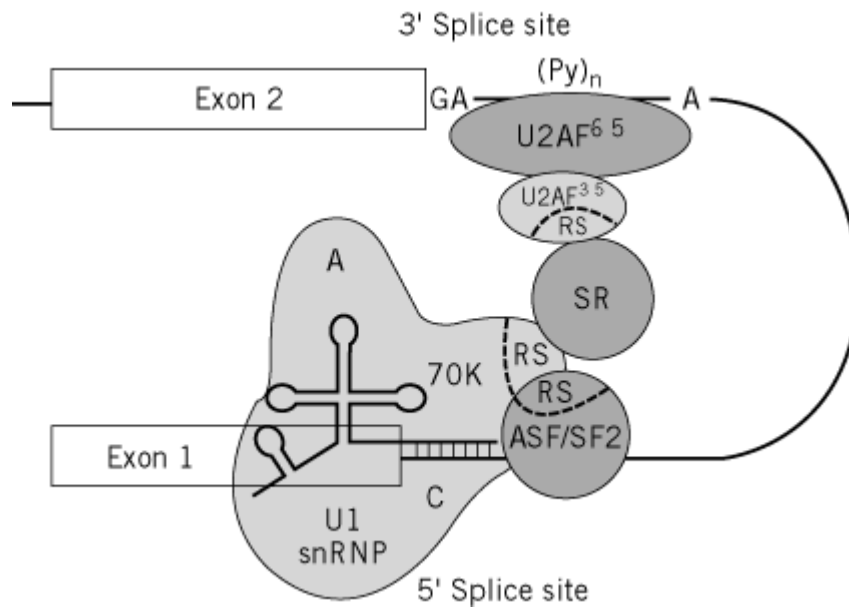
Although intermediate complexes exist (as evidenced, for example, by genetic techniques in yeast), under normal splicing conditions four distinct spliceosomal complexes can be readily detected by native [gel electrophoresis](#) and/or gel filtration. Assembly of the major spliceosome is initiated by the interaction of the U1 snRNP with the 5' splice site of the pre-mRNA, which leads to the formation of spliceosomal complex E. The U2 snRNP then binds to the branch site, giving rise to the so-called pre-spliceosome or A complex. U5 and U4/U6 snRNPs, in the form of a preassembled [U4/U6.U5] tri-snRNP complex, subsequently interact to generate the spliceosome (B complex). After a conformational rearrangement, a catalytically active spliceosome (C complex) is ultimately formed.

Numerous non-snRNP splicing factors also interact at defined steps of the spliceosome assembly pathway and during catalysis (reviewed in Refs. [4](#) and [5](#)). A similar assembly pathway has also been deduced recently for the U12-dependent spliceosome. In this case, initially the U11 and U12 snRNPs interact with the 5' splice site and branch site, respectively ([1](#), [2](#), [6](#), [7](#)). Because the U11 and U12 snRNP form a stable 18 S di-snRNP particle, they, unlike U1 and U2, are thought to interact with the pre-mRNA together as a preformed complex ([8](#)). In the subsequent step, the U5 and U4atac/U6atac snRNPs also appear to interact as a preassembled [U4atac/U6atac.U5] tri-snRNP complex to form mature spliceosomes ([9](#)). Non-snRNP splicing factors present in the U12-dependent spliceosome have not yet been identified, but are also thought to associate at various stages of the assembly process. Although its characterization has just begun, initial studies indicate that an RNA–RNA network similar to that present in the U2-dependent spliceosome (see below) is also formed in the U12-dependent spliceosome. Based on these observations, it has been suggested that similar active sites are formed in both types of spliceosome.

## 1. E Complex

Spliceosome assembly is initiated by the formation of the commitment complex (yeast designation) or the early (E) complex (metazoan designation), in which the 5' and 3' splice sites are initially recognized by the U1 snRNP and the splicing factor U2AF (U2 snRNP auxiliary factor; Mud2 in yeast), respectively. E complex formation does not require ATP and may be considered a nucleation event for spliceosome assembly. That is, the pre-mRNA in these complexes is “committed” to spliceosome assembly, as evidenced by the inability of an excess of noncomplexed competitor pre-mRNA to inhibit the subsequent conversion of E complexes to active spliceosomes. In higher eukaryotes, E complex formation is dependent upon a complex set of molecular interactions involving primarily members of the SR protein superfamily. The latter proteins, which are not found in the yeast *S. cerevisiae*, are characterized by a domain rich in Arg–Ser dipeptide sequences (RS domain) and often one or more **RNA-binding** domains (RBDs) (reviewed in Refs. [10](#) and [11](#)). Due to their modular organization, SR proteins can simultaneously bind RNA (eg, the pre-mRNA) via their RBD and other factors through **protein-protein interactions** involving their RS domain. Association of the U1 snRNP with the pre-mRNA involves base pairing between the 5' end of the U1 snRNA and the 5' splice site (Fig. [2](#)) and is facilitated by several proteins, including the SR protein ASF/SF2 and the U1 snRNP 70K and C proteins (see [Small Nuclear Rnps \(Snrnps\)](#)). U2AF consists of a 35- and a 65-kDa subunit and interacts, via an RNA-binding domain in the latter subunit, directly with the pre-mRNA's polypyrimidine tract just upstream of the 3' splice site (see [Splice Sites](#)). During E complex formation, a network of protein–protein interactions involving SR proteins is formed between the U1-70K protein and U2AF, both of which also contain an RS domain (Fig. [2](#)). In this way, the 5' and 3' splice sites are apparently brought into close proximity with one another at the earliest stages of spliceosome assembly. In yeast, this bridging of the 5' and 3' splice sites may be carried out by proteins that are uniquely associated with the yeast U1 snRNP ([12](#)) (see [Small Nuclear Rnps \(Snrnps\)](#)).

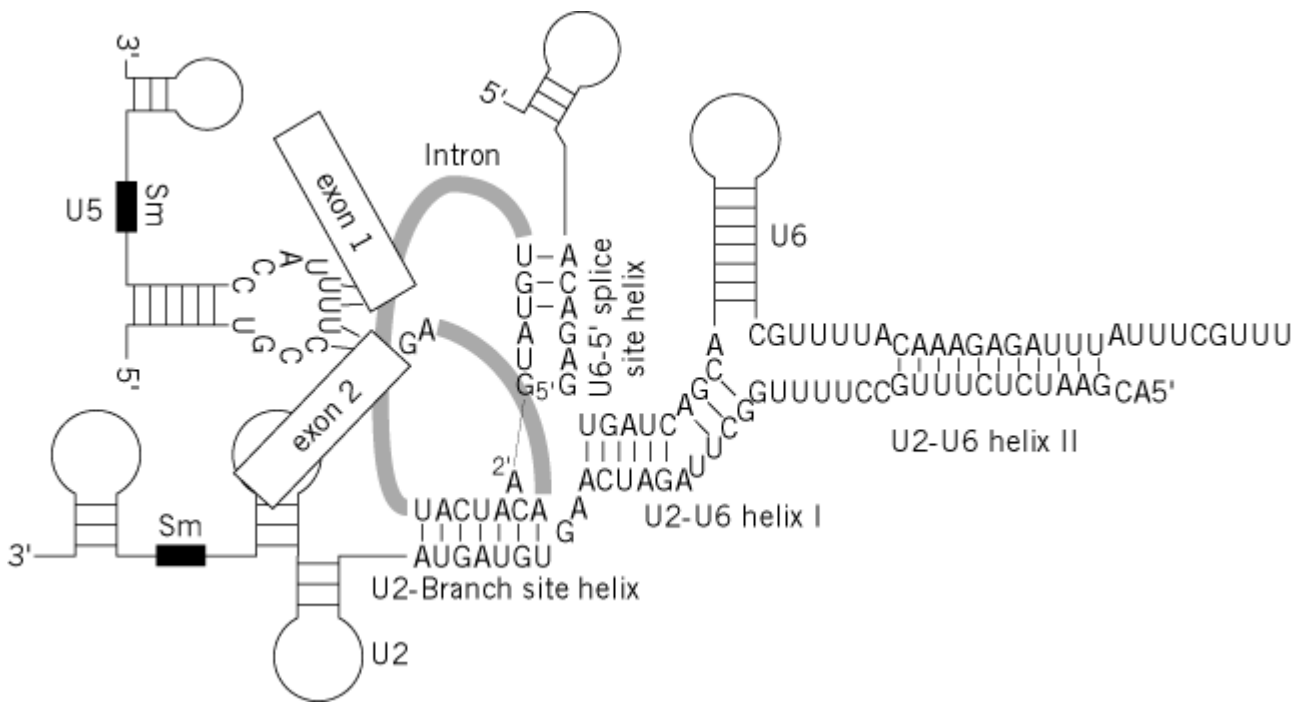
**Figure 2.** Network of protein–protein and protein–RNA interactions that occur during spliceosomal E complex formation. The U1 snRNP consists of the U1 snRNA (thick solid line), whose 5' end base-pairs with the pre-mRNA's 5' splice site, the snRNP core proteins (not indicated), and the U1-specific 70K, A, and C proteins. The U1 snRNA/5' splice site interaction is facilitated by the interaction of the RS domain of the SR protein ASF/SF2 with that of the U1-70K protein. SR proteins are thought to bridge the 5' and 3' splice sites by interacting simultaneously with U1-70K and the 35 kDa subunit of U2AF, which also contains an RS domain. Exon and intron sequences are indicated by boxes and thin solid lines, respectively. RS domains are delineated by a dotted line. The branch point adenosine (A), the polypyrimidine tract (Py)<sub>n</sub>, and the conserved dinucleotide at the 3' splice site (GA) are indicated.



## 2. A Complex

E complexes are converted to the pre-spliceosome or A complex by the ATP-dependent interaction of the U2 snRNP with the branch site. U2 snRNP binding involves base pairing between nucleotides near the 5' end of the U2 snRNA and the conserved branch site sequence (see Fig. 3). This results in the bulging out of the branch site adenosine, which is thought to position it favorably for its nucleophilic attack at the 5' splice site (see [RNA Splicing](#)). The interaction of U2 with the branch site is dependent on the presence of U2AF, as well as the splicing factor SF1/mBBP (BBP in yeast), which interacts with the branch site prior to A complex formation (13, 14). A subset of U2 snRNP proteins also plays a critical role in stabilizing the U2 snRNP interaction with the pre-mRNA. These proteins include the heteromeric splicing factors SF3a and SF3b, which are essential for A complex assembly and are part of the functional 17 S form of the U2 snRNP (see [Small Nuclear Rnps \(Snrnps\)](#)). SF3a and SF3b components interact with the pre-mRNA near the branch site and are thus thought to stabilize the U2 snRNP/pre-mRNA interaction (15). Most aspects of A complex formation appear to be conserved in yeast, including an essential role for the yeast **homologues** of most of the proteins comprising SF3a and SF3b (16, 17).

**Figure 3.** Network of RNA interactions occurring among the U2, U5, and U6 snRNAs and the pre-mRNA within the pre-spliceosome. In this model, the pre-mRNA has undergone the first transesterification reaction. The cleaved exon 1 is held in place by interactions with nucleotides of the U5 loop 1 sequence. Exon sequences are indicated by open boxes; and intron sequences, other than conserved nucleotides at the 5' splice site (GUAUGU), branch site (UACUAACA), and the 3' splice site (AG), are depicted by a thick shaded line. The 2', 5' phosphodiester bond between the first intron nucleotide (G) and branch point adenosine is depicted by a thin dark line. Cartoon structural models of U2, U5, and U6 snRNAs are according to the metazoan consensus structure, whereas all nucleotide sequences are from the yeast *S. cerevisiae*. Courtesy of Patrizia Fabrizio.



### 3. B and C Complexes

Pre-spliceosomes are converted to spliceosomes by the association of a preassembled tri-snRNP complex that contains the U4/U6 and U5 snRNPs. B complex assembly is an ATP-dependent process that results in the formation of a complex pre-mRNA–snRNA network. In this network, the pre-mRNA is folded so that its chemically reactive groups (eg, the 5' and 3' splice sites and the branch site) are juxtaposed (Fig. 3) (reviewed in Ref. 18). Subsequent to tri-snRNP association, the U6 snRNA dissociates from the U4 snRNA and forms new RNA duplexes with the U2 snRNA and intron nucleotides of the 5' splice site (Fig. 3). These interactions appear to juxtapose the 5' splice site with the branch site adenosine, thereby facilitating the first transesterification reaction. Furthermore, the U6 snRNA/5' splice site interaction appears to specify the precise site of nucleophilic attack at the 5' splice junction. During B complex formation, the U1 / 5' splice-site base-pairing interaction is disrupted, and the U5 snRNA, via its highly conserved loop 1, interacts with a small number of nonconserved exon nucleotides at the 5' splice site (Fig. 3). This interaction is thought to be one of the factors responsible for retaining the cleaved 5' exon in the spliceosome after the first transesterification reaction. Subsequent to this reaction, U5 loop 1 nucleotides also interact with exon nucleotides at the 3' splice site, consistent with the idea that U5 also aligns the 3' OH of the 5' exon and the 3' splice site for the second catalytic step of splicing. Although poorly understood when compared with the E and A complexes, B complex formation also undoubtedly involves the formation of multiple protein–protein and protein–RNA interactions. One well-established interaction is that between the U5 snRNP 220 kDa protein (Prp8p in yeast) and the pre-mRNA. This protein interacts with sequences surrounding both the 5' and 3' splice sites and is thus thought to stabilize U5 snRNA interactions with the pre-mRNA (19). In mammals, SR proteins have also been implicated in this step of spliceosome assembly, although the exact nature of their presumed interactions remains to be established (20). Recently, two proteins containing an RS domain, and therefore potential SR protein interaction partners, have been identified in the tri-snRNP complex (the U5-100 kDa protein and tri-snRNP associated 27 kDa protein) (21, 22). Thus, a network of SR protein interactions has also been proposed to help stabilize the interaction of the tri-snRNP complex with the pre-mRNA during B complex formation.

Just prior to or concomitant with the first transesterification reaction, major structural changes occur



in the spliceosome, generating complex C. This complex represents a catalytically active spliceosome that has undergone the first step of splicing and thus contains the pre-mRNA splicing intermediates (cleaved 5' exon and the intron-3' exon lariat; see [RNA Splicing](#)). The U1 and U4 snRNPs appear to be only loosely associated with the spliceosome at these later stages of the splicing reaction. Upon completion of both catalytic steps, the mRNA is released from the spliceosome in an ATP-dependent process, and the excised intron remains in a complex that is poorly characterized but contains at least the U2, U5, and U6 snRNPs. Ultimately, this post-splicing complex dissociates, and the intron lariat is debranched and degraded. The released snRNPs are believed to reassemble and participate repeatedly in new rounds of spliceosome assembly.

#### 4. Dynamic Nature of the Spliceosome

Understanding the dynamics of RNA and protein interactions during spliceosome assembly and the catalytic steps of splicing is currently one of the major focuses of the nuclear pre-mRNA splicing field. Spliceosome assembly, as well as disassembly, is accompanied by changes in many of the interactions formed between the pre-mRNA and other spliceosomal components (snRNPs and non-snRNP proteins) and among the latter spliceosomal components themselves. The formation of some of these interactions is dependent on the disruption of others (ie, they are mutually exclusive), a phenomenon that appears to contribute to the highly ordered nature of spliceosome assembly. The best, albeit still poorly, understood conformational changes occur in the spliceosomal RNA–RNA network. For example, as eluded to above, major rearrangements occur at the 5' splice site, which first interacts with the U1 snRNA and then, after U1 dissociation, with the U5 and U6 snRNAs. These multiple recognition events have been proposed to act as a proofreading mechanism that decreases aberrant cleavage at the 5' splice site. Furthermore, the U6 snRNA exchanges its interaction partners, first base-pairing with the U4 snRNA in the tri-snRNP complex and later during spliceosome assembly with the U2 snRNA and the 5' splice site. These and other RNA rearrangements appear to be catalysed by spliceosomal proteins belonging to the **DEAD/H-box** superfamily of [ATPases](#) and [RNA Helicases](#). Members of this protein family share eight highly conserved sequence motifs, among them ATP-binding and hydrolysis motifs. At present, seven putative ATP-dependent RNA helicases have been detected in yeast and four in mammals (reviewed in Ref. [24](#)). Some of these proteins (eg, Prp5p, UAP56, and Prp28p/U5-100 kD) perform essential functions during spliceosome assembly, facilitating U2 snRNP or [U4/U6.U5] tri-snRNP association. Others are required for the first or second catalytic steps of splicing (eg, Prp2p, Snu246/U5-200 kD, and Prp16p) or for the release of the splicing products from the spliceosome (eg, Prp22/HRH1 and Prp43p). Most of these proteins are known to exhibit RNA-stimulated ATPase activity. Thus, the requirement for ATP at various steps of the splicing process is due, at least in part, to the activity of these proteins. Only very recently has it been possible to demonstrate *in vitro* RNA duplex-unwinding activity for three of these spliceosomal proteins, namely Prp16p ([25](#)), Snu246/U5-200 kD ([26](#), [27](#)), and Prp22p ([28](#), [29](#)). Potential RNA substrates for each of these proteins have been proposed based on results obtained from *in vitro* studies and what is known about the step at which each of them is required for splicing. However, their exact substrates and functions within the spliceosome remain to be established. Spliceosomal rearrangements could also be facilitated by the U5 snRNP 116 kDa protein (Snu114 in yeast), which shares significant homology with the ribosomal [GTPase](#) EF2 ([30](#)). The latter protein facilitates translocation of the [ribosome](#) during [translation](#), and the U5-116 kD protein has thus been proposed to function in an analogous manner in the spliceosome.

Although currently not well understood, changes in spliceosome structure also appear to involve the ordered formation and disruption of both protein–protein and protein–RNA interactions. For example, several spliceosomal proteins are known to interact transiently with the pre-mRNA and/or the spliceosome ([5](#), [31](#)). Changes in those interactions involving proteins may arise indirectly through the activity of RNA helicases or may be directly regulated by the activity of protein isomerases or **molecular chaperones**. Interestingly, a component of the [U4/U6.U5] tri-snRNP complex, the 20 kDa protein, shares extensive homology with known [peptidyl-prolyl cis/trans](#)

[isomerases](#) of the [cyclophilin](#) family, suggesting that it may promote protein conformational changes during splicing and thereby contribute to the dynamics of the spliceosome ([32](#), [33](#)). Protein–protein and protein-RNA interactions within the spliceosome also appear to be regulated by protein [kinases](#) and [phosphatases](#). For example, SR proteins are known to undergo phosphorylation/dephosphorylation cycles during splicing, and their **phosphorylation** state has been shown to affect both spliceosome assembly and splicing catalysis (reviewed in Refs. [10](#) and [11](#)). Several SR protein kinases (eg, SRPK1 and Clk/Sty) have been identified, and recently it has become clear that RS-domain phosphorylation can enhance the interaction of an SR protein with other RS domain-containing proteins ([34](#)). This suggests that SR protein interactions within the spliceosome are dynamic in nature. Moreover, at least two serine/threonine protein phosphatases, namely PP1 and PP2a, are required for both catalytic steps of splicing, and the former appears to be responsible for the dephosphorylation of one or more SR proteins ([35](#)). Despite recent advances, however, the regulation of nuclear pre-mRNA splicing by protein kinases and phosphatases remains a relatively poorly understood area of the splicing field.

### Bibliography

1. S. L. Hall and R. A. Padgett (1996) *Science* **271**, 1716–1718.
2. W.-Y. Tarn and J. A. Steitz (1996) *Cell* **84**, 801–811.
3. P. A. Sharp and C. B. Burge (1997) *Cell* **91**, 875–879.
4. J. F. Cáceres and A. R. Krainer (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, U.K. pp. 174–212.
5. P. E. Hodges, M. Plumpton, and J. D. Beggs (1997) In *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, U.K., pp. 213–241.
6. Y.-T. Yu and J. A. Steitz (1997) *Proc. Natl. Acad. Sci. USA* **94**, 6030–6035.
7. I. Kolossova and R. A. Padgett (1997) *RNA* **3**, 227–233.
8. K. M. Wassarman and J. A. Steitz (1992) *Mol. Cell. Biol.* **12**, 1276–1285.
9. W.-Y. Tarn and J. A. Steitz (1996) *Science* **273**, 1824–1832.
10. X.-D. Fu (1995) *RNA* **1**, 663–680.
11. J. L. Manley and R. Tacke (1996) *Genes Dev.* **10**, 1569–1579.
12. A. Gottschalk, J. Tang, O. Puig et al., (1998) *RNA* **4**, 374–393.
13. N. Abovich and M. Rosbash (1997) *Cell* **89**, 403–412.
14. J. A. Berglund, K. Chua, N. Abovich, R. Reed, and M. Rosbash (1997) *Cell* **89**, 781–787.
15. O. Gozani, R. Feld, and R. Reed (1996) *Genes Dev.* **10**, 233–243.
16. P. E. Hodges and J. D. Beggs (1994) *Curr. Biol.* **4**, 264–267.
17. H. Igel, S. Wells, R. Perriman, and M. Ares (1998) *RNA* **4**, 1–10.
18. H. D. Madhani and C. Guthrie (1994) *Annu. Rev. Genet.* **28**, 1–26.
19. S. Teigelkamp, A. J. Newman and J. D. Beggs (1995) *EMBO J.* **14**, 2602–2612.
20. R. F. Roscigno and M. A. Garcia-Blanco (1995) *RNA* **1**, 692–706.
21. S. Fetzner, J. Lauber, C. L. Will, and R. Lührmann (1997) *RNA* **3**, 344–355.
22. S. Teigelkamp, C. Mundt, T. Achsel, C. L. Will, and R. Lührmann (1997) *RNA* **3**, 1313–1326.
23. F. V. Fuller-Pace (1994) *Trends Cell Biol.* **4**, 271–274.
24. J. P. Staley and C. Guthrie (1998) *Cell* **92**, 315–326.
25. Y. Wang, J. D. Wagner, and C. Guthrie (1998) *Curr. Biol.* **8**, 441–451.
26. B. Lagerbauer, T. Achsel, and R. Lührmann (1998) *Proc. Natl. Acad. Sci. USA* **95**, 4188–4192.
27. P. Raghunathan and C. Guthrie (1998) *Cur. Biol.* **8**, 847–855.
28. J. D. O. Wagner, E. Jankowsky, M. Company, A. M. Pyle, and J. N. Abelson (1998) *EMBO J.* **17**, 2926–2937.
29. B. Schwer and C. H. Gross (1998) *EMBO J.* **17**, 2086–2094.

30. P. Fabrizio, B. Laggerbauer, J. Lauber, W. S. Lane, and R. Lührmann (1997) *EMBO J.* **16**, 4092–4106.
31. M. D. Chiara, O. Gozani, M. Bennett, P. Champion-Arnaud, L. Palandjian, and R. Reed (1996) *Mol. Cell. Biol.* **16**, 3317–3326.
32. D. S. Horowitz, R. Kobayashi, and A. R. Krainer (1997) *RNA* **3**, 1374–1387.
33. S. Teigelkamp, T. Achsel, C. Mundt, S.-F. Göthel, U. Cronshagen, W. S. Lane, M. Marahiel, and R. Lührmann (1998) *RNA* **4**, 127–141.
34. S.-H. Xiao and J. L. Manley (1997) *Genes Dev.* **11**, 334–344.
35. J. E. Mermoud, P. T. W. Cohen, and A. I. Lamond (1994) *EMBO J.* **13**, 5679–5688.

### Suggestions for Further Reading

36. R. Reed and L. Palandjian (1997) "Spliceosome Assembly. In" *Eukaryotic mRNA Processing* (A. R. Krainer, ed.), IRL Press, Oxford, U.K., pp. 103–129.
37. A. Krämer (1996) The structure and function of proteins involved in mammalian pre-mRNA splicing. *Annu. Rev. Biochem.* **65**, 367–409.
38. C. L. Will and R. Lührmann (1997) Protein functions in pre-mRNA splicing. *Curr. Opin. Cell Biol.* **9**, 320–328.
39. T. W. Nilsen (1998) "RNA–RNA Interactions in Nuclear Pre-mRNA Splicing. In" *RNA Structure and Function* (R. W. Simons and M. Grunberg-Manago, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 279–307.
40. J. P. Staley and C. Guthrie (1998) Mechanical devices of the spliceosome: motors, clocks, springs and things. *Cell* **92**, 315–326.
41. W.-Y. Tarn and J. A. Steitz (1997) Pre-mRNA splicing: the discovery of a new spliceosome doubles the challenge. *Trends Biochem. Sci.* **22**, 132–137.

## Spore, Sporulation

Two of the major challenges facing all life forms are propagation of the species and survival under adverse conditions. One mechanism used by many organisms to meet these challenges is to encapsulate their [genome](#) in a specialized cell structure, known as a spore, that is metabolically inactive and often highly resistant to environmental assault. The spore protects the organism's genome until an appropriate time when the instructions encoded on the **DNA** can be used to produce a new individual. Because the morphogenic process of spore formation (sporulation) is often associated with extracellular production in many prokaryotic species of commercially valuable commodities (**antibiotics**, [enzymes](#), bioinsecticides, and more) and is related to pathogenicity and toxicity of others, spore research remains an active and important scientific area. Additionally, sporulation cycles in some organisms serve as relatively simple examples of cellular [differentiation](#) events that can be analyzed at the molecular level.

Spores can be grouped into two broad types that are not necessarily mutually exclusive: reproductive spores and resting spores. The former, made by various **plant** species and some **fungi**, can be either single-celled or multicellular. Produced in enormous amounts and passively disseminated by wind, water, or animals after detachment from the parental organism, reproductive spores can each give rise to a new individual and thus serve for rapid increases in the population of the species. Resting, or dormant, spores are produced by many prokaryotic species (notably in the genera *Bacillus*,

*Clostridium*, *Streptomyces*, *Myxococcus*) and some lower eukaryotes (eg, *Dictyostelium*) as a survival mechanism when the organism encounters an unfavorable environment. These spores can remain in a quiescent state until favorable growth conditions return, at which time they undergo a differentiation process, known as germination and outgrowth, to regenerate actively growing and dividing cells (vegetative cells). There is good evidence that dormant bacterial spores can remain viable for up to a century (1), and some recent observations suggest that survival for millions of years may also be possible (2). Although considerable progress continues to be made in unraveling the details of spore formation in *Myxococcus*, *Streptomyces* and *Dictyostelium* (for reviews see Refs. 3-5), the molecular biology of sporulation has been best characterized for *Bacillus subtilis*.

Extensive research on the structure and properties of *Bacillus* spores has served as a paradigm in the field. The two most important factors responsible for the longevity of these bacterial spores are:

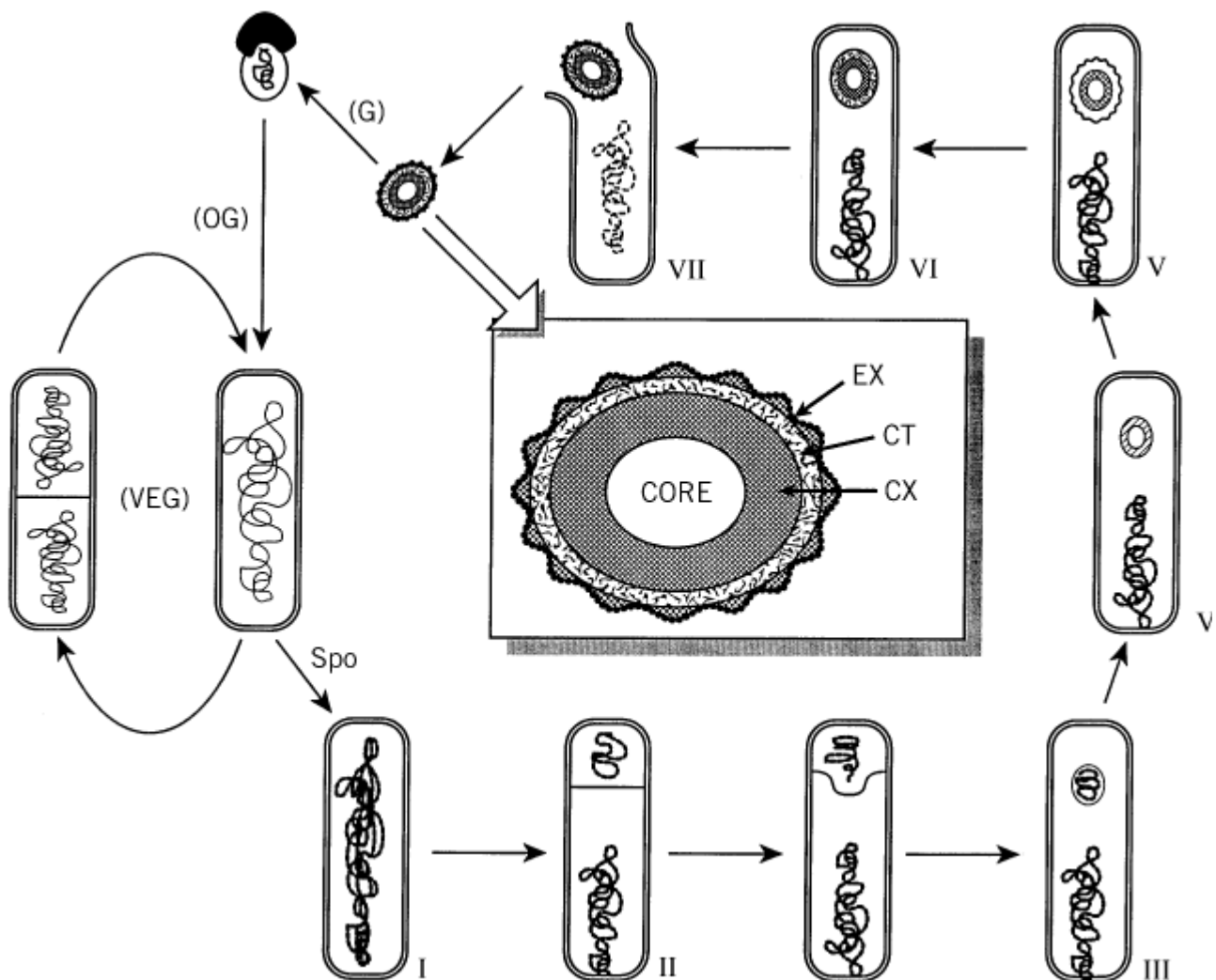
1. encasement of inactivated enzymes, [ribosomes](#), and DNA in a dehydrated central core surrounded by a multilayered protective shell; and
2. protection of the DNA present in the core from environmental damage.

The factors responsible for the inactivity of spore enzymes are not entirely understood, but the lowered [water](#) content of the core (about 5- to 10-fold less than that present in the cytoplasm of growing cells) is undoubtedly a major contributor to both this property and to overall heat resistance (6). In spores of *Bacillus* and related genera, a class of fascinating proteins, known as a/b SASP (small acid-soluble proteins), coat the spore DNA and protect it from damaging assaults caused by desiccation, heat, oxidation and ultraviolet radiation (7). The multilayered integument (8) surrounding the core provides a permeability barrier preventing access of some damaging agents and functions to maintain physicochemical integrity of the core itself. From the inside out, the protective layers of a mature *Bacillus* spore generally include:

1. a cytoplasmic membrane enveloping the core;
2. a thin layer of peptidoglycan known as the germ cell wall that is very similar, or identical, to vegetative cell peptidoglycan;
3. a thick layer called the cortex that is comprised of a different form of peptidoglycan;
4. a second membrane;
5. a thick layer of proteinaceous coat that itself can be subdivided into layers of distinct appearance; and
6. an outer membranous layer known as the exosporium.

The morphological landmarks of *B. subtilis* sporulation are illustrated schematically in Figure 1. [Mutation](#) analysis has identified over 100 **genes** required for normal spore formation. Many of these genes are designated by a specialized nomenclature indicating at which morphological stage mutations in the gene halt development. For example, *spoII* mutants arrest at stage II (asymmetric septum formation) and do not proceed to stage III (forespore engulfment). Different [operons](#) are distinguished by a single capital letter (*spoIIA*, *spoIIB*, etc.), and different cistrons within the same operon are denoted by additional letters (*spoIIAA*, *spoIIAB*, etc.). Not all sporulation-associated genes use this nomenclature and, in many instances, elucidation of the biochemical function of the gene product has led to renaming of the gene (eg, *spoIIAC* was discovered to encode a [s factor](#) of **RNA polymerase**, and the gene is now called *sigF*).

**Figure 1.** Sporulation of *Bacillus subtilis*. Landmark morphological stages are: I, axial filament of chromosome; II, polar septation dividing the sporangium into smaller forespore and larger mother cell compartments; III, engulfment; IV, cortex formation; V, coat formation; VI, maturation; VII, lysis of mother cell and release of spore. Under laboratory conditions 37%, the process takes about 8 h. Germination (G), outgrowth (OG), and vegetative growth (VEG) complete the life cycle. The center inset is a diagram illustrating major features of a mature spore. The layers are not drawn to strict scale; spore membranes and distinctions in the coat layer are not shown. CX, cortex; CT, protein coat; EX, exosporium.



Sporulation is but one of the developmental options available to *B. subtilis* cells when they are deprived of nutrients and enter a stationary phase of growth. Numerous interconnected regulatory circuits sense the exact environmental and metabolic conditions and dictate which option is chosen. Initiation of the sporulation pathway is absolutely dependent upon the level of phosphate flux through a [signal transduction](#) system known as the phosphorelay (9). The phosphorelay is a complex assemblage of protein [kinases](#), protein **phosphatases** and phosphorylatable proteins, with each component apparently being a focal point for input of signals indicative of nutritional and physiological status. There is only sketchy evidence suggesting how the activities of the kinases are regulated (10, 11) and how intracellular GTP levels might affect the phosphoprotein phosphotransferase encoded by *spo0B* (12). Most current excitement concerns the signaling roles played by small [peptides](#) that affect the activity of the phosphatases (13-15). The signals, whatever their exact nature, regulate the level of the ultimate output of the phosphorelay: Spo0A protein in its phosphorylated form (Spo0A~P<sub>4</sub>). Spo0A~P<sub>4</sub> is a transcriptional regulatory factor that is responsible for activating expression of three operons (*spoIIA*, *spoIIG*, *spoIIE*) required for transition from stage II to stage III of sporulation. Additionally, Spo0A~P<sub>4</sub> represses transcription of the *abrB* gene, and this leads to expression of many functions that had been kept silent during vegetative growth due to the regulatory actions of AbrB protein (16). Spo0A~P<sub>4</sub> is also required for the hallmark of stage II: formation of an asymmetrically-positioned division septum. Although the exact mechanistic role of Spo0A~P<sub>4</sub> in this event is unclear, it is clear that commitment to sporulate occurs soon after formation of this septum (17). Up until this point, the sporulation pathway could

have been aborted and the cell channeled either back into vegetative growth or else into an alternative developmental state of stationary phase.

Completion of the sporulation septum divides the cell (now termed the sporangium) into two compartments: the forespore (prespore) and the mother cell. [Because the mother cell subsequently engulfs the forespore compartment (stage III) and further development and maturation occurs in this configuration, the spore produced is more properly called an endospore.] Each compartment proceeds to express different sets of genes necessary for endospore formation. This compartmentalization is primarily due to differential activity of alternative RNA polymerase sigma factors (18), but it is likely that the differences in [chromosome](#) condensation in each compartment also play an important role (19). A fascinating aspect of sigma factor compartmentalization is that intercompartmental communication mechanisms temporally link the activity occurring in one compartment to the differential activity that occurs within the other. The general scheme of this crisscross regulation (20) is that  $s^F$  activity in the forespore leads to activation of  $s^E$  in the mother cell;  $s^E$  activity in the mother cell is necessary for  $s^G$  activation in the forespore;  $s^G$ -directed transcription in the forespore is required for  $s^K$  activity to follow  $s^E$  activity in the mother cell. Much remains to be elucidated concerning how these complex interplays occur, but many of the details now known provide illustrative examples as to how cellular morphology can affect **gene expression**, and *vice versa*.

Spo0A~P<sub>04</sub>-activated expression of the operon (*spoIIA*) encoding  $s^F$  occurs prior to completion of the sporulation septum, and so all three protein products of the operon (SpoIIAA, SpoIIAB,  $s^F$ ) are present in both compartments of the sporangium. However,  $s^F$  activity occurs only in the forespore due to a regulatory mechanism termed “partner switching” (21). The SpoIIAB protein is an anti-sigma factor that can bind  $s^F$  and sequester it in an inactive form. SpoIIAB can also bind SpoIIAA (an anti-anti- $s$  factor), and its relative affinity for each mutually exclusive partner is influenced by the intracompartamental ATP/ADP ratio (see [Adenylate Charge](#)). Additionally, SpoIIAB is a kinase able to phosphorylate SpoIIAA, and SpoIIAA~P<sub>04</sub> cannot bind SpoIIAB (22). Although many questions remain unanswered, the current model is that in the forespore the ATP/ADP ratio is low; this condition favors binding of SpoIIAB–SpoIIAA and release of free, active  $s^F$ . Furthermore, a membrane-associated phosphatase encoded by *spoIIE* (whose [transcription](#) is activated by Spo0A~P<sub>04</sub>) becomes localized to the septum but becomes active only on the forespore side, thus counteracting SpoIIAB kinase action upon SpoIIAA in that compartment. In the mother cell, SpoIIAB- $s^F$  binding predominates, due apparently due to a higher ATP/ADP ratio, SpoIIAB phosphorylation of SpoIIAA, and absence of SpoIIE phosphatase activity (23).

Two genes expressed in the forespore due to  $s^F$ -directed transcription are *sigG* and *spoIIR*. Gene *sigG* encodes the  $s$  factor responsible for late forespore gene expression, and *spoIIR* encodes a protein necessary for  $s^E$  activation in the mother cell. As in the case of *spoIIA*, the operon (*spoIIG*) encoding  $s^E$  is expressed prior to septum completion as a result of Spo0A~P<sub>04</sub>-mediated transcriptional activation. The initial proteins produced are a membrane-bound putative [proteinase](#) (SpoIIGA) and an inactive precursor of  $s^E$  (pro- $s^E$ ). The production of SpoIIR in the forespore compartment is the signal responsible for initial processing of pro- $s^E$  to  $s^E$  in the mother cell (24). SpoIIR is secreted out of the forespore, and this process somehow activates SpoIIGA-dependent proteinase activity on the mother cell side of the septal membrane (25). It is not clear whether SpoIIR provides an absolute directionality that prevents pro- $s^E$  to  $s^E$  processing from occurring in the forespore or if additional mechanisms that selectively inactivate  $s^E$  in that compartment also exist (26).

Relatively little is currently understood concerning regulation of  $s^G$  activation in the forespore, but it

is apparent that  $s^E$ -controlled transcription of the eight-**cistron** *spoIIIA* operon in the mother cell is a crucial event (27). It also is clear that  $s^G$  activity in the forespore is required for  $s^K$  activity to appear in the mother cell. Synthesis and activation of  $s^K$  involves a number of mechanisms, including one that could occur only in a terminally differentiated cell destined to lyse. An intact, contiguous gene (*sigK*) encoding  $s^K$  is not normally present on the chromosome: The *sigK* **open reading frame** is interrupted by a 48-kb element referred to as *skin*. During sporulation,  $s^E$  and a transcriptional regulatory protein, SpoIIID, activate transcription of a gene (*spoIVCA*) in the *skin* element. The SpoIVCA protein is a site-specific **recombinase** that catalyzes excision of a circular DNA molecule in which the truncated portions of *sigK* become fused in-frame (see [Translation](#)). The resultant intact copy of *sigK* is transcribed by  $s^E$ /SpoIIID, and the **messenger RNA** is translated to produce an inactive pro- $s^K$  protein. This precursor is then processed to active  $s^K$  by a proteinase whose activity is dependent upon  $s^G$ -directed expression of genes in the forespore (28). Once activated in their respective compartments, both  $s^K$  and  $s^G$  direct further transcription of their own genes via positive autoregulatory feedback loops and are responsible for expression of gene products necessary for the final steps of endospore development and maturation.

The cascade of *s* factors during *B. subtilis* sporulation is obviously a primary means of achieving temporally regulated expression of genes required for the ordered assembly of the morphological structure of the endospore. Yet many details concerning this process remain to be elucidated, and it seems highly probable that other, perhaps novel, mechanisms await discovery. Although many genes responsible for the physical components and assembly of the spore structures (eg, the peptidoglycan cortex, protein coat) have been identified, very little is actually known about the biochemistry and enzymology underlying these assemblages. Additionally, much work remains to be done in order to understand the process of spore germination and outgrowth, an area that has received far too little attention in the past. Nevertheless, the continuing rapid pace of discovery has shown that sporulation is an excellent experimental system for examining fundamental aspects of cellular differentiation.

## Bibliography

1. M. J. Kennedy, S. L. Reader, and L. M. Swierczynski (1994) *Microbiol.* **140**, 2513–2529.
2. R. J. Cano and M. K. Borucki (1995) *Science* **268**, 1060–1064.
3. M. Dworkin (1996) *Microbiol. Rev.* **60**, 70–102.
4. W. F. Loomis (1996) *Microbiol. Rev.* **60**, 135–150.
5. W. C. Champness and K. F. Chater (1994) In *Regulation of Bacterial Differentiation* (P. Piggot, C. P. Moran, and P. Youngman, eds.), American Society for Microbiology, Washington, D.C., pp. 61–94.
6. P. Gerhardt and R. E. Marquis (1989) In *Regulation of Prokaryotic Development* (I. Smith, R. Slepecky, and P. Setlow, eds.), American Society for Microbiology, Washington, D.C., pp. 17–63.
7. P. Setlow (1995) *Annu. Rev. Microbiol.* **49**, 29–54.
8. D. J. Tipper and J. J. Gauthier (1972) In *Spores V* (H. O. Halvorson, R. Hanson, and L. L. Campbell, eds.), American Society for Microbiology, Washington, D.C., pp. 3–12.
9. J. A. Hoch (1993) *Annu. Rev. Microbiol.* **47**, 441–465.
10. M. A. Strauch, D. de Mendoza, and J. A. Hoch (1992) *Mol. Microbiol.* **6**, 2909–2917.
11. V. Dartois, T. Djavakhishvili, and J. A. Hoch (1996) *J. Bacteriol.* **178**, 1178–1186.
12. J. Kok, K. A. Trach, and J. A. Hoch (1994) *J. Bacteriol.* **176**, 7155–7160.
13. M. Perego et al. (1994) *Cell* **79**, 1047–1055.
14. M. Perego and J. A. Hoch (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1549–1553.
15. M. Perego, P. Glaser, and J. A. Hoch (1996) *Mol. Microbiol.* **19**, 1151–1157.
16. M. A. Strauch and J. A. Hoch (1993) *Mol. Microbiol.* **7**, 337–342.

17. G. F. Parker, R. A. Daniel, and J. Errington (1996) *Microbiology* **142**, 3445–3452.
18. W. G. Haldenwang (1995) *Microbiol. Rev.* **59**, 1–30.
19. B. Setlow et al. (1991) *J. Bacteriol.* **173**, 6270–6278.
20. R. Losick and P. Stragier (1992) *Nature* **355**, 601–604.
21. L. Duncan, S. Alper, and R. Losick (1994) *Curr. Opin. Gen. Dev.* **4**, 630–636.
22. T. Magnin, M. Lord, J. Errington, and M. D. Yudkin (1996) *Mol. Microbiol.* **19**, 901–907.
23. F. Arigoni et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 3238–3242.
24. M. L. Karow, P. Glaser, and P. J. Piggot (1995) *Proc. Natl. Acad. Sci. USA* **92**, 2012–2016.
25. A. E. M. Hofmeister et al. (1995) *Cell* **83**, 219–226.
26. P. J. Piggot (1996) *Curr. Opin. Gen. Dev.* **6**, 531–537.
27. E. M. Kellner, A. Decatur, and C. P. Moran (1996) *Mol. Microbiol.* **21**, 913–924.
28. E. Ricca, S. Cutting, and R. Losick (1992) *J. Bacteriol.* **174**, 3177–3184.

### Suggestions for Further Reading

29. J. A. Hoch (1995) "Control of cellular development in sporulating bacteria by the phosphorelay two-component signal transduction system", In *Two-Component Signal Transduction* (J. A. Hoch and T. J. Silhavy, eds.), ASM Press, Washington, D.C, pp. 129–144.
30. P. Stragier and R. Losick. (1996) Molecular genetics of sporulation in *Bacillus subtilis*. *Annu. Rev. Genet.* **30**, 297–341.
31. J. Errington (1993) *Bacillus subtilis*: Regulation of gene expression and control of morphogenesis. *Microbiol. Rev.* **57**, 1–33.
32. M. A. Strauch (1993) Regulation of *Bacillus subtilis* gene expression during the transition from exponential growth to stationary phase. *Prog. Nucleic Acids Res. Mol. Biol.* **46**, 121–153.
33. A. D. Grossman (1995) Genetic networks controlling the initiation of sporulation and the development of genetic competence in *Bacillus subtilis*. *Annu. Rev. Genet.* **29**: 477–508.

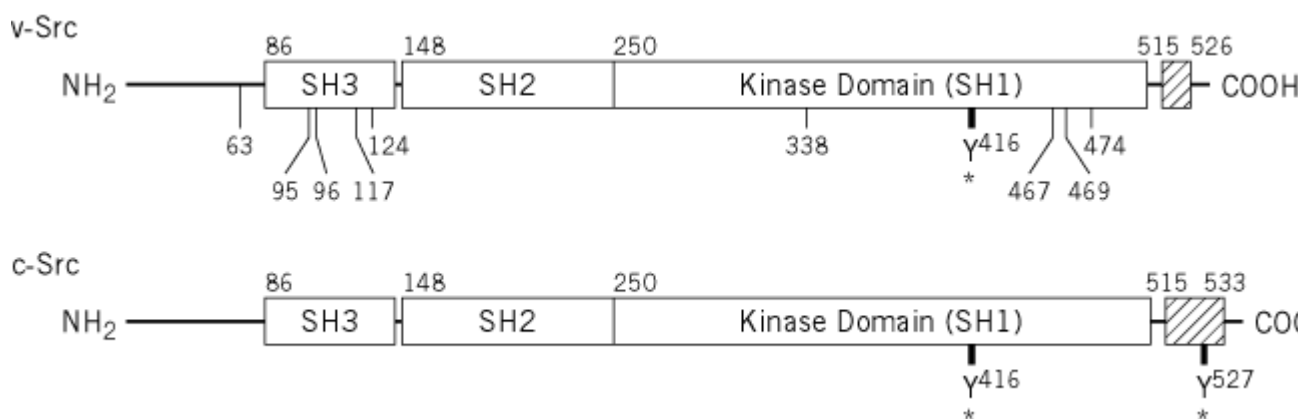
### SRC Genes

Viral (v-) and cellular (c-) *src* genes were the first identified oncogene and proto-oncogene, respectively. The v-*src* gene was identified as the transforming gene of Rous sarcoma virus (RSV), a typical RNA tumor virus that causes neoplastic transformation of chicken cells (1). Hence, the gene name *src* (*sarcoma*). The translational product of v-*src* is a phosphoprotein of 60 kDa, designated as v-Src or pp60<sup>v-*src*</sup>, which was shown to be the first example of tyrosine kinase, a novel family of protein kinases (2). The transforming ability of v-*src* is dependent on the tyrosine kinase activity of the gene product. On the other hand, the c-*src* gene was originally found as a cellular homologue of v-*src* in normal chicken cells and is now known to be widespread in higher eukaryotes. There are 12 exons in the chicken c-*src* (3), whereas v-*src* has no introns. It is believed that v-*src* has been derived from a spliced messenger RNA of c-*src*, and, thus, the c-*src* gene is called a proto-oncogene. The product of c-*src*, termed c-Src or pp60<sup>c-*src*</sup>, is also a 60-kDa phosphoprotein and has an intrinsic tyrosine kinase activity. The gene expression of c-*src* is essentially ubiquitous. However, the expression level of c-Src is usually very low, in contrast to v-Src, which is highly expressed from the viral genome. In addition, the activity of c-Src is strictly regulated and much lower than that of v-Src. The c-*src* gene has been shown not to be oncogenic, even if it is artificially overexpressed.



The major structural difference between the v-Src and c-Src proteins resides at their carboxy termini (Fig. 1). The 19-amino acid residue tail of c-Src is replaced by an unrelated sequence of 12 amino acids in v-Src. Biochemical and molecular biological analyses have revealed that this difference is responsible for the elevated tyrosine kinase activity of v-Src and its transforming ability. Some other point mutations are also present in v-Src. These amino acid replacements may also contribute to the kinase activation. Several artificial mutations have been shown to activate c-Src. For example, the single amino acid change of Tyr527 to Phe can cause activation of the kinase and elevation of the oncogenic potential of the mutant gene product, although the product is still less potent as a transforming protein than v-Src.

**Figure 1.** Schematic structures of v-Src and c-Src. Both v-Src (from Rous sarcoma virus SR-A strain) and chicken c-Src are shown. The kinase domain (SH1) and two domains (SH2 and SH3) that are conserved among the Src family kinases are shown by **open bars**. The carboxyl terminal sequences (shown by **hatched bars**) differ from each other in both sequence and length. Nine single-amino acid substitutions are also known and indicated by **vertical lines** below v-Src with their residue numbers. v-Src or activated c-Src undergoes autophosphorylation on Tyr416 (Y<sup>416</sup>), and inactive c-Src is phosphorylated on Tyr527 (Y<sup>527</sup>). These sites are shown by **asterisks**.



After the discovery of *src* genes, several closely related genes have been identified as oncogenes of viral origin or proto-oncogenes, and these are now included in the Src family. The Src family members identified thus far are Src, Fyn, Yes, Fgr, Lck, Lyn, Hck, Blk, and Yrk (1-6). The Src family genes are evolutionary conserved in higher eukaryotes, and all encode nonreceptor tyrosine kinases of around 60 kDa. Structural comparison of the Src family kinases has revealed a common feature, designated as the Src homology domain, which consists of three conserved domains: Src homology 1 (SH1), Src homology 2 (SH2), and Src homology 3 (SH3). The SH1 domain corresponds to the carboxy-terminal kinase (catalytic) domain itself. The SH2 and SH3 domains reside in the amino-terminal noncatalytic region (Fig. 1). The amino-terminal sequence is divergent among the family members, except for the several N-terminal residues that contain the invariant Gly2. This glycine residue is shown to be necessary for the attachment of myristoyl group by which Src family kinases bind to the plasma membrane.

The inactive forms of Src family kinases have been shown to be phosphorylated on the carboxy-terminal tyrosine residue (Fig. 1). by another tyrosine kinase, named Csk (7). The SH2 domain binds to this terminal phosphotyrosine residue through an intramolecular interaction, while the SH3 domain reinforces this interaction by binding to a proline-rich motif in the same molecule. Therefore, one of the functions of SH2 and SH3 domains in the Src family proteins is to maintain the structure of the proto-oncogene product in an inactive state. This fact was demonstrated by crystallographic analysis of c-Src and Hck (8, 9). Because the carboxy-terminal tyrosine is absent from the v-Src sequence (Fig. 1), the SH2 domain in v-Src cannot form the intramolecular interaction. This is

thought to be the why v-Src kinase is deregulated and constitutively activated. The activation mechanism under physiological conditions of the normal Src family kinases is not well understood. However, it is conceivable that any kind of signal that leads to the dissociation of either SH2- or SH3-mediated interaction will activate the Src family kinases.

As to their physiological function, the Src family kinases are shown to be physically and functionally coupled with various surface receptors of hematopoietic cells. For example, Fyn is associated with the T cell receptor, Lyn with the B cell receptor, and Lck with CD antigens of T cells. Some receptor tyrosine kinases, such as platelet-derived growth factor receptor, epidermal growth factor receptor, and CSF-1 receptor are also known to be coupled with the Src family members ([10](#)).

### Bibliography

1. R. Jove and H. Hanafusa (1987) *Annu. Rev. Cell Biol.* **3**, 31–56.
2. T. Hunter and J. A. Cooper (1985) *Annu. Rev. Biochem.* **54**, 897–930.
3. T. Takeya and H. Hanafusa (1983) *Cell* **32**, 881–890.
4. J. A. Cooper (1990) "The src family of protein-tyrosine kinases", In *Peptides and Protein Phosphorylation*, (B. E. Kemp, ed.), CRC Press, Boca Raton, FL, pp. 85–113.
5. S. M. Dymecki, J. E. Niederhuber, and S. V. Desiderio (1990) *Science* **247**, 332–336.
6. M. Sudol, H. Greulich, L. Newman, A. Sarkar, J. Sukegawa, and T. Yamamoto (1993) *Oncogene* **8**, 823–831.
7. M. Okada and H. Nakagawa (1989) *J. Biol. Chem.* **264**, 20886–20893
8. W. Wu, S. C. Harrison, and M. J. Eck (1997) *Nature* **385**, 595–602.
9. F. Sicheri, I. Moarefi, and J. Kuriyan (1997) *Nature* **385**, 602–609.
10. S. M. Thomas and J. S. Brugge (1997) *Annu. Rev. Cell Dev. Biol.* **13**, 513–609.

### Src Homology Domain

The Src homology (SH) domain is a general term for the three conserved domains (SH1, SH2, and SH3) originally identified by amino acid sequence analysis of protein products of *src* genes and related genes. The SH1 domain represents the catalytic tyrosine kinase domain. The SH2 and SH3 domains are modules that bind to phosphotyrosine-containing and proline-rich sequence motifs, respectively. Both SH2 and SH3 domains are now found in a wide variety of biosignaling molecules and have been shown to play important roles in the mechanism of signal transduction through protein-protein interactions.

### Suggestion for Further Reading

- T. Pawson (1995) Protein modules and signaling networks. *Nature* **373**, 573–580.

### Stabilization And Destabilization By Co-Solvents

The effect of a co-solvent on the stability of a protein (or other macromolecule) is determined strictly by the difference between the interactions of the co-solvent with the protein in the unfolded and native states (see [Protein Stability](#)). The essential interactions that must be known are [preferential binding](#) (or **preferential hydration**) and the [transfer free energy](#), ie, **free energy** of interaction with ([binding](#) to) the protein of the structure stabilizer or destabilizer (1).

Let the unfolding (stability) process be represented by a simple equilibrium between the native  $N$  and denatured  $D$  states, and let this equilibrium be affected by a ligand (stabilizing or destabilizing co-solvent):



The effect of the co-solvent can be expressed then relative to either of two reference states:

1. In the first reference state the given co-solvent concentration is given. The question is: Does the co-solvent at this concentration stabilize the protein relative to an infinitesimally smaller concentration? The effect the co-solvent has then is expressed by the *Wyman linkage relation*:

$$\frac{d \log K}{d \log a_L} = \nu^D - \nu^N = \Delta \left( \frac{\partial m_L}{\partial m_{pr}} \right)_{T,P,\mu_L} \quad (2)$$

where  $n^D$  and  $n^N$  stand for the preferential binding of the co-solvent to the unfolded and native protein ( $a_L$  is the thermodynamic activity of the co-solvent, which frequently may be approximated by the concentration). Therefore, the slope of a log-log plot of the equilibrium constant versus the ligand concentration gives the difference in preferential binding between the two end states.

Conversely, knowledge of  $n^D$  and  $n^N$  indicates whether a ligand (co-solvent) will be a stabilizer or destabilizer of the structure: Greater [equilibrium dialysis](#) binding to the unfolded form enhances the unfolding reaction; greater binding to the native form results in stabilization of the native structure.

(2) In the second reference state, the effect is measured relative to [water](#) as solvent. It is expressed by the difference between the binding free energies to the two end states:

$$\delta \Delta G^{0(N \rightarrow D)} = \Delta G_L^0 - \Delta G_w^0 = \Delta \mu_{pr,tr}^D - \Delta \mu_{pr,tr}^N = \delta \Delta \mu_{pr,tr}^{N \rightarrow D} \quad (3)$$

where  $DG_L^0$  is the standard free energy of the unfolding equilibrium at ligand concentration  $L$  and  $DG_w^0$  is that in water (this is related to the equilibrium constant by  $DG^0 = -2.303RT \log K$ ).  $Dm_{pr,tr}$  is the transfer free energy of the protein from water (dilute buffer) to the solvent of the given composition in the denatured  $D$  and native  $N$  states; it, in fact, is the free energy of binding of the ligand (co-solvent) to the protein ( $DG_b = Dm_{pr,tr}$ ) in its two states.

#### Bibliography

1. S. N. Timasheff (1993) *Ann. Rev. Biophys. Biomol. Struct.* **22**, 67–97.

#### Suggestion for Further Reading

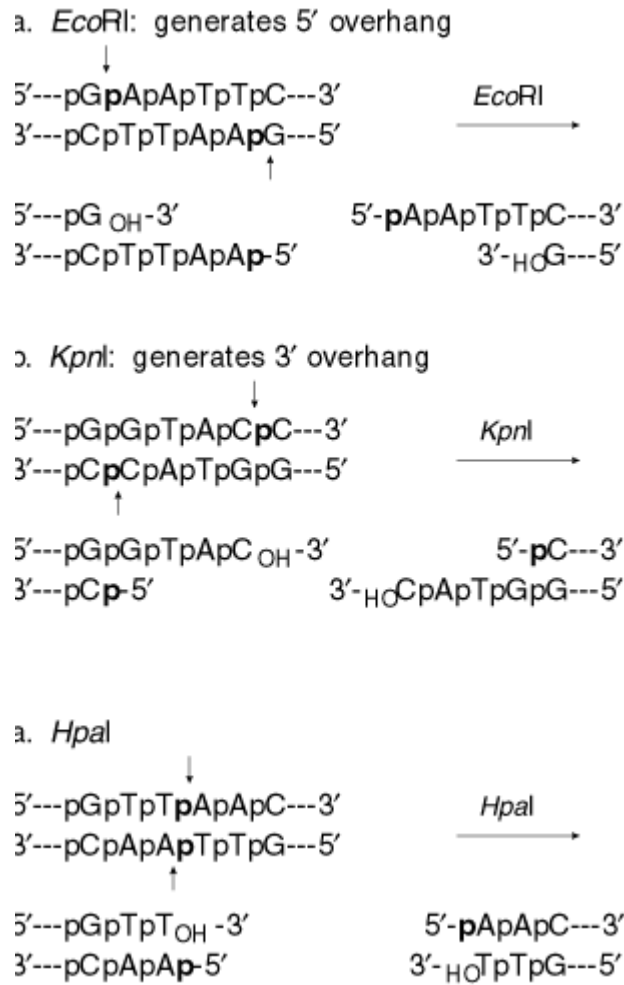
2. S. N. Timasheff (1995) In *Protein-Solvent Interactions* (R. B. Gregory, ed.), Marcel Dekker, New York, Chap. "11".

## Staggered Cut

[Restriction Enzymes](#) are bacterial endonucleases that recognize defined sequences of double-stranded DNA and catalyze phosphodiester bond hydrolysis of both strands at positions within or adjacent to the recognition sequence (see [Restriction Enzymes](#)). This article focuses on the type II restriction endonucleases, which represent greater than 90% of the known restriction enzymes and are those most used in [recombinant DNA](#) manipulations (see [Cloning](#)). For more information on types I, II, and III restriction enzymes, see the article on [Restriction–Modification Systems](#).

The sequences recognized by restriction enzymes are 4–8 base pairs in length and are called *restriction sites*. The hydrolytic scission occurs 5' to the phosphate in the phosphodiester bond of each strand, so that 5' phosphoryl and 3' hydroxyl termini are produced. The hydrolytic cuts can be made either directly opposite one another or offset from one another along the helix, generating either blunt or cohesive ends, respectively, in the resultant DNA fragments (see [Restriction Fragment](#)) (Fig. 1). Within the type II family of enzymes, there are several classes (types IIp, IIw, II<sub>n</sub>, and II<sub>t</sub>), which are distinguished by the nature of the sequence recognized and by the position of cleavage within or outside of the site (1). Representative examples are given in Table 1.

**Figure 1.** Types of cleavage by restriction endonucleases and resulting DNA termini: (a) staggered cut resulting in cohesive ends; (b) blunt cut resulting in blunt ends. Arrows indicate position of cleavage, and boldface **p** indicates phosphodiester bond cleaved.



**Table 1. Classes of Type II Restriction Enzymes**

| Type II        | Recognition  |  |  |
|----------------|--|--|--|
| Class          | Specificity  | Position of Cleavage                                   | Example(s)   |
| I <sub>p</sub> | Palindromic sequence, with or without degeneracy <sup>a</sup> with even number of base pairs | Within recognition sequence                            | <i>Bam</i> HI: G ^ GATCC<br><i>Ban</i> I: G ^ G PyPu CC <sup>b</sup><br><i>Sma</i> I: CCC ^ GGG<br><br><i>Pst</i> I: CTGCA ^ G |
| I <sub>s</sub> | Asymmetric   | Outside recognition sequence (up to 20base pairs away) | <i>Hga</i> I: GACGC (5/10) <sup>c</sup><br><i>Fok</i> I: GGATG (9 / 13)  |

|     |   |   |  |
|-----|---|---|--|
| Iiw | Palindromic with degeneracy and odd number of base pairs    | Within recognition sequence                             | <i>HinFI</i> : G ^ ANTC <sup>d</sup><br><i>Sau96I</i> : G ^ GNCC |
| IIn | Degenerate spacer segment interrupting palindromic sequence | Within spacer region                                    | <i>PshAI</i> : GACNN ^ NNGTC<br><i>PflMI</i> : CCANNNN ^ NTGG    |
| IIt | Asymmetric  | Similar to IIs with one cut inside recognition sequence | <i>BsrI</i> : 5'-ACTGGN ^ N-3'<br>3'-TGAC ^ CNN-5'               |

<sup>a</sup> The allowance of one or more nucleotide substitutions at a specific position in the restriction site.

<sup>b</sup> Py is Pyrimidine containing nucleotide residue (C,T), and Pu is purine containing nucleotide residue (A,G).

<sup>c</sup> Indicates the position of the site of cleavage of the “top” and “bottom” strands 3' to the recognition site.

<sup>d</sup> N is G,A,T, or C.

The cuts produced by restriction enzymes yield defined fragments of DNA that can be exploited in several recombinant DNA procedures. Appropriately cleaved fragments of duplex DNA produced from restriction digestion can be ligated (joined) to other fragments using a variety of methods (Table 2) (see also [Ligation](#)). For fragments containing compatible cohesive ends (3' or 5' overhangs) resulting from digests with the same enzyme, such as *EcoRI* or *Pst I*, ligation of both fragments with [DNA Ligase](#) regenerate the same restriction site. Fragments produced by different enzymes but having compatible termini can generate novel cleavage specificities on ligation. For example, *SalI* (G ^ TCGAC)-digested and *AvaI* (C ^ TCGAG)-digested DNA fragments contain 5' overhangs that are compatible with one another. On ligation of the fragments, both the *SalI* and *AvaI* sites are destroyed, and only a *TaqI* (^ TCGA) restriction site remains in the product DNA.

**Table 2. Manipulations of DNA Fragment Ends by DNA Ligation**

| DNA fragments | DNA restriction    | DNA restriction    | Resultant, new        |
|---------------|--------------------|--------------------|-----------------------|
| containing:   | fragment 1, with a | fragment 2, with a | restriction site upon |
|               | given Terminus:    | given Terminus:    | ligation of           |
|               |                    |                    | fragments 1 & 2:      |

|   |   |   |  |
|---|---|---|--|
| Compatible cohesive ends—<br>same restriction enzyme                          | 5'-G <sub>OH</sub><br>3'-CTTAA <sub>p</sub><br><i>EcoRI</i> (G ^ AATTC)   | pAATTC-3'<br>OHG-5'<br><i>EcoRI</i> (G ^ AATTC)           | 5'-GAATTC-3'<br>3'-CTTAAG-5'<br><i>EcoRI</i> (G ^ AATTC)   |
| Compatible cohesive ends—<br>restriction enzymes with different specificities | 5'-G <sub>OH</sub><br>3'-CAGCT <sub>p</sub><br><i>SalI</i> (G TCAGC)  | pTCGAG-3'<br>OH <sup>C</sup> -5'<br><i>AvaI</i> (C TCGAG) | 5'-GTCGAG-3'<br>3'-CAGCTG-5'<br><i>TaqI</i> (5- ^ TCGA-3') |
| Different blunt-ended termini   | 5'-GAT <sub>OH</sub><br>3'-CTA <sub>p</sub><br><i>EcoRV</i> (GAT ^ ATC)   | pCCT-3'<br>OHGGA-5'<br><i>StuI</i> (AGG ^ CCT)            | 5'-GATCCT-3'<br>3'-CTAGGA-5'<br><i>MboI</i> (^ GATC)       |
| Incompatible ends   | 5'-GAATT <sub>OH</sub> <sup>a</sup><br>3'-CTTAA <sub>p</sub><br><i>EcoRI</i> treated with Klenow fragment of DNA polymerase I | pCC-3'<br>OHGG-5'<br><i>HaeIII</i> (GG ^ CC)              | 5'-GAATTCC-3'<br>3'-CTTAAGG-5'<br>5'-GAATTCC-3'            |

---

<sup>a</sup> Italics represent nucleotides added by Klenow fragment of DNA polymerase I.

New cleavage sites can also be generated by ligating blunt-ended termini together (see [Blunt-End Ligation](#)). For example, a *MboI* site (^ GATC) can be generated by ligating two blunt-ended fragments: one generated by the blunt-end producer *EcoRV* (GAT ^ ATC) and the other by *StuI* (AGG ^ CCT). Fragments with incompatible cohesive 5' ends can be converted to blunt-ended, double-strand molecules with the **Klenow fragment** of **DNA polymerase I**. For example, fragments containing the 5' overhang of an *EcoRI* site (G ^ AATTC) can be filled in with Klenow fragment to form blunt-ended molecules with the 5' sequence AATT and ligated to blunt-ended fragments produced by *HaeIII* (GG ^ CC); the product DNA then contains a regenerated *EcoRI* site (GGAATTC). Fragments containing either 5' or 3' overhanging ends can also be made blunt-ended by removing the single-stranded segment with either [S1 nuclease](#) or mung bean nuclease and ligated to other blunt-ended termini. Fragments produced by digestion with two different restriction enzymes, for example, with single *PstI* and *BamHI* sites, will have different termini (*PstI* or *BamHI*) on each end that can be ligated into cut DNA with similar, compatible termini.

Orientation problems are usually found with blunt-end ligations. For example, a gene borne on a blunt-ended DNA fragment ligated into a vector downstream of a promoter has only a 50% chance of being ligated in the proper orientation for expression, because blunt-end fragments can be joined in either orientation. This can be overcome with ligations using cohesive-end or double-digested fragments, which assure correct orientation of the gene on ligation, using restriction enzymes of the type I<sub>p</sub>, I<sub>w</sub>, I<sub>n</sub>, and I<sub>t</sub> classes (Table 1) (1).

The type II<sub>s</sub> restriction enzymes (see [Restriction–Modification Systems](#)) offer unique opportunities for molecular cloning experiments because they cleave outside, but 3' adjacent to, the sequence recognized, producing long staggered ends (2). One application involves the design of an adapter with a type II<sub>s</sub> site within an engineered single-stranded region that anneals to a complementary single-stranded DNA template. This adapter can first be utilized as a primer for DNA polymerase to synthesize a double-stranded template, and this product can subsequently be hydrolyzed with the type II<sub>s</sub> restriction enzyme to achieve the desired cut. Another application is for the precise excision of DNA fragments and for gene assembly using a DNA vector containing a unique restriction site immediately between two inward-facing type II<sub>s</sub> sequences (2, 3). These applications illustrate why type II<sub>s</sub> enzymes have been referred to as “universal” restriction enzymes.

Restriction enzymes can be used in combination with methyltransferases to generate new or rare cleavage specificities (see [Methylation, DNA](#) and [Methyltransferase](#)). For example, premethylation of (GGATCCGG) sequences with the methylase M•*Msp*I to form mCCGG inhibits methylation by M•*Bam*HI (GGATmCC), but not cleavage by the restriction enzyme R•*Bam*HI. Therefore, all *Bam*HI sequences not adjacent to methylated *Msp*I sites will be rendered insensitive to *Bam*HI cleavage after methylation by M•*Bam*HI (4). Another application, methylase-limited partial digestion, involves incubation of DNA for limited times with methyltransferases in order to methylate some but not all canonical sequences. Therefore, some of the canonical sequences remain susceptible to cleavage by the partner endonuclease, while others do not. This method is convenient when large restriction fragments are desired (5).

Given the variety of sequence specificities, cleavage and methylation properties of available restriction and methylation enzymes, and the numerous ways in which they can be used in DNA manipulation and analysis, it is clear why the discovery of these enzymes revolutionized molecular cloning. Other applications of these enzymes are discussed in the articles on [restriction map](#) and [Restriction Fragment](#).

#### Bibliography

1. Kessler and V. Manta (1990) *Gene* **92**, 1–248.
2. W. Szybalski, S. C. Kim, N. Hasan, and A. J. Podhajski (1991) *Gene* **100**, 13–26.
3. H.-Y. Eun (1996) in *Enzymology Primer for Recombinant DNA*, Academic Press, San Diego, pp. 233–306.
4. A. Pingould, J. Alves, and R. Geiger (1993) in *Methods in Molecular Biology Vol. 16, Enzymes of Molecular Biology*, M. M. Burrell, ed., Humana Press, Inc., Totowa, N.J., p. 167.
5. M. Nelson, E. Raschke, and M. McClelland (1993) *Nucl. Acids Res.* **21**, 3139–3154.

#### Star Activity



*Star activity* refers to the aberrant DNA cleavage by type II restriction endonucleases at sites other than their canonical recognition sequences (see [Restriction–Modification Systems](#) and [Restriction Enzymes](#)). The star activity of an enzyme is usually denoted with an asterisk following the name of the enzyme, such as *EcoRI\**. Ordinarily, the recognition specificity of restriction endonucleases for the canonical sequence is very high. The canonical site is generally cleaved  $10^6$  times more frequently than other DNA sequences under optimal conditions (1). With high enzyme concentration or with long incubation periods, however, cleavage can also occur at sites differing by one or more base pairs (referred to as *star sites*), although the canonical site continues to be selected preferentially. Presumably star sites are cleaved much more slowly because of the absence of critical interactions of the enzyme with the bases or the phosphate backbone that are needed for optimal binding or catalysis (2).

Specificity can also be altered as a result of suboptimal reaction conditions: high pH, low ionic strength,  $Mn^{2+}$  instead of  $Mg^{2+}$ , or the presence of cosolvents (glycerol, DMSO, or ethanol). Under these conditions, selectivity for the canonical restriction site is reduced as the enzyme relaxes its sequence specificity, and cleavage at degenerate or shortened versions of the canonical sequence occurs more often. The allowable degeneracy may range from (1) a single-base-pair change, (2) loss of recognition of distal base pairs in the canonical segment, to (3) several possible base pair substitutions in the primary recognition site. In the presence of organic solvents, *EcoRI*, which normally cleaves at the sequence GAATTC, has a star activity, *EcoRI\**, toward sequences containing the core tetranucleotide sequence AATT (1). *XbaI*, which cleaves at the canonical sequence TCTAGA, is an example of an extreme star activity; it becomes a nonspecific nuclease under conditions of elevated pH and in the presence of glycerol and DMSO (1). Isoschizomeric enzymes (see [Isoschizomer](#)) have been shown to have different cleavage specificities under star reaction conditions (3), so they respond differently to conditions. Star activity can be a problem when restriction enzymes are used as tools for precisely fragmenting DNA. On the other hand, the characteristic changes in specificity patterns of restriction enzymes under conditions inducing star activity has been useful in mechanistic studies (4).

### Bibliography

1. A. Pingoud, J. Alves, and R. Geiger (1993) in *Methods in Molecular Biology*, Vol. 16, Enzymes of Molecular Biology, M. M. Burrill, ed., Humana Press, Totowa, N.J. p. 122.
2. A. Jeltsch, J. Alves, H. Wolfes, G. Maass, and A. Pingoud (1994) *Biochemistry* 33, 10215–10219.
3. H.-Y. Eun (1996) in *Enzymology Primer for Recombinant DNA*, Academic Press, San Diego, pp. 233–306.
4. C. R. Robinson and S. G. Sligar (1995) *Proc. Natl. Acad. Sci. USA* 92, 3444–3448.

### Suggestion for Further Reading

5. E. G. Malygin and V. V. Zinoviev (1989) Studies on the role of symmetry in the specific recognition of natural and synthetic DNA by type II restriction and modification enzymes. *Sov. Sci. Rev. Devel. Physiochem. Biol.* 9, 87–139.

### Start Codons

In *Escherichia coli*, the **codon** AUG of the [genetic code](#) is an efficient initiator of [translation](#) in **protein biosynthesis** and is by far the most common one. GUG and UUG are also efficient initiation

codons and are not uncommon. AUU is only rarely used, most notably in the gene for translation [initiation factor 3](#).

At initiation of translation, the anticodon of a special initiator transfer RNA,  $\text{tRNA}^{\text{Met}}_{\text{f}}$  pairs with the initiation codon bound at the P site of a [ribosome](#). In **eubacteria**, [chloroplasts](#) and [mitochondria](#), the initiating amino acid is formylmethionine (1). The formyl group permits binding by initiation factor 2. This, together with special features of  $\text{tRNA}^{\text{Met}}_{\text{f}}$  ensures that formylmethionine is delivered to the ribosomal P site and becomes the first amino-terminal amino acid to be incorporated into a polypeptide chain. The formyl group is removed subsequently, however, so that methionine is the first encoded amino acid. Formylmethionine  $\text{tRNA}^{\text{Met}}_{\text{f}}$  is not recognized by the [elongation factor](#) EF-Tu, which delivers aminoacyl tRNA to the ribosomal A site for incorporation at internal positions in polypeptide proteins. Instead, the “elongator” that decode AUG, GUG, UUG, and AUU at internal positions in coding regions are not acylated with formylmethionine, but rather with methionine, valine, leucine, and isoleucine, respectively, as in the normal genetic code.

There is a special context for initiator codons that gives them their distinctive function and meaning. In the overwhelming majority of cases in eubacteria, this special context is a short sequence of bases in the [messenger RNA](#) known as a **Shine–Dalgarno** sequence after its discoverers. It is located approximately 5 or more nucleotides 5' of the initiator codon. The Shine–Dalgarno sequence of the messenger RNA pairs with a complementary sequence very close to the 3' end of 16 S rRNA in the small (30 S) ribosomal subunit. The interaction serves to position the initiation codon for interaction with initiator tRNA at the ribosomal P site. In a small number of mRNAs from eubacterial bacteriophage or extreme thermophilic **archae**, there is no Shine–Dalgarno sequence, and occasionally no 5' sequence at all. In these cases, some sequence 3' to the start codon that is only partially defined serves to facilitate initiation.

In **eukaryotes**, AUG is again by far the most common initiator codon, but CUG, ACG, and GUG are also found. The special initiator tRNA,  $\text{tRNA}^{\text{Met}}_{\text{i}}$ , inserts methionine and not formylmethionine. The discriminatory features used to specify when these codons are to function as initiators are different from that utilized in eubacteria. With the aid of initiation factors, the small, 40 S subunit of eukaryotic ribosomes generally associates with mRNA via a [-cap](#) structure and, at least in yeast, the [poly A](#) tail at the 3' end of the mRNA also plays a role (2, 3). The recruited subunit then scans until a potential initiator codon flanked by a consensus sequence is encountered (4). The large, 60 S ribosomal subunit subsequently joins to complete assembly of the intact, 80 S ribosome. The great majority of eukaryotic start codons are recognized by ribosomes recruited in this manner, but there are two interesting exceptions.

1. The mRNA for a number of **proto-oncogenes** and [transcription factors](#) have one or more short open [reading frames](#) (ORF) preceding the main protein coding sequence. Ribosomes that have translated a particular short ORF may have the potential to reinitiate translation at the main ORF start codon location some distance in the 3' direction. This mode of initiation has been studied in great depth in yeast GCN4 (5).
2. Special structures known as IRES (internal ribosome entry segments) permit ribosomes to make initial contact with internal regions in mRNA and initiate translation at appropriately positioned start codons (6). IRES-mediated internal initiation was initially discovered in translation of the picornaviruses, [polio virus](#) and encephalomyocarditis, where it enables viral protein synthesis to escape the consequences of viral-induced inhibition of cap-mediated initiation (7, 8). This inhibition serves to inactivate host translation. However, the start codons of a small number of chromosomally encoded genes are also recognized by IRES-mediated initiation.

## Bibliography

1. D. Mazel, E. Coïc, S. Blanchard, W. Saurin, and P. Marlière (1997) *J. Mol. Biol.* **266**, 939–949.
2. N. Iizuka, L. Najita, A. Franzusoff, and P. Sarnow (1994) *Mol. Cell Biol.* **14**, 7322–7330.
3. S. Z. Tarun and A. B. Sachs (1996) *EMBO J.* **15**, 7168–7177.
4. M. Kozak (1997) *EMBO J.* **16**, 2482–2492.
5. A. G. Hinnebusch (1996) In *Translational Control* (J. W. B. Hershey, M. B. Mathews, and N. Sonenberg, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 199–244.
6. R. J. Jackson and A. Kaminski (1995) *RNA* **1**, 985–1000.
7. J. Pelletier and N. Sonenberg (1988) *Nature* **334**, 320–325.
8. S. K. Jang, H. G. Krausslich, M. J. H. Nicklin, G. M. Duke, A. C. Palmenberg, and E. Wimmer (1988) *J. Virol.* **62**, 2636–2643.

## Suggestions for Further Reading

9. J. W. B. Hershey, M. B. Mathews, and N. Sonenberg (1996) *Translational Control*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp. 1–794.
10. A. B. Sachs, P. Sarnow, and M. W. Hentze (1997) Starting at the beginning, middle, and end: Translation initiation in eukaryotes. *Cell* **89**, 831–838.

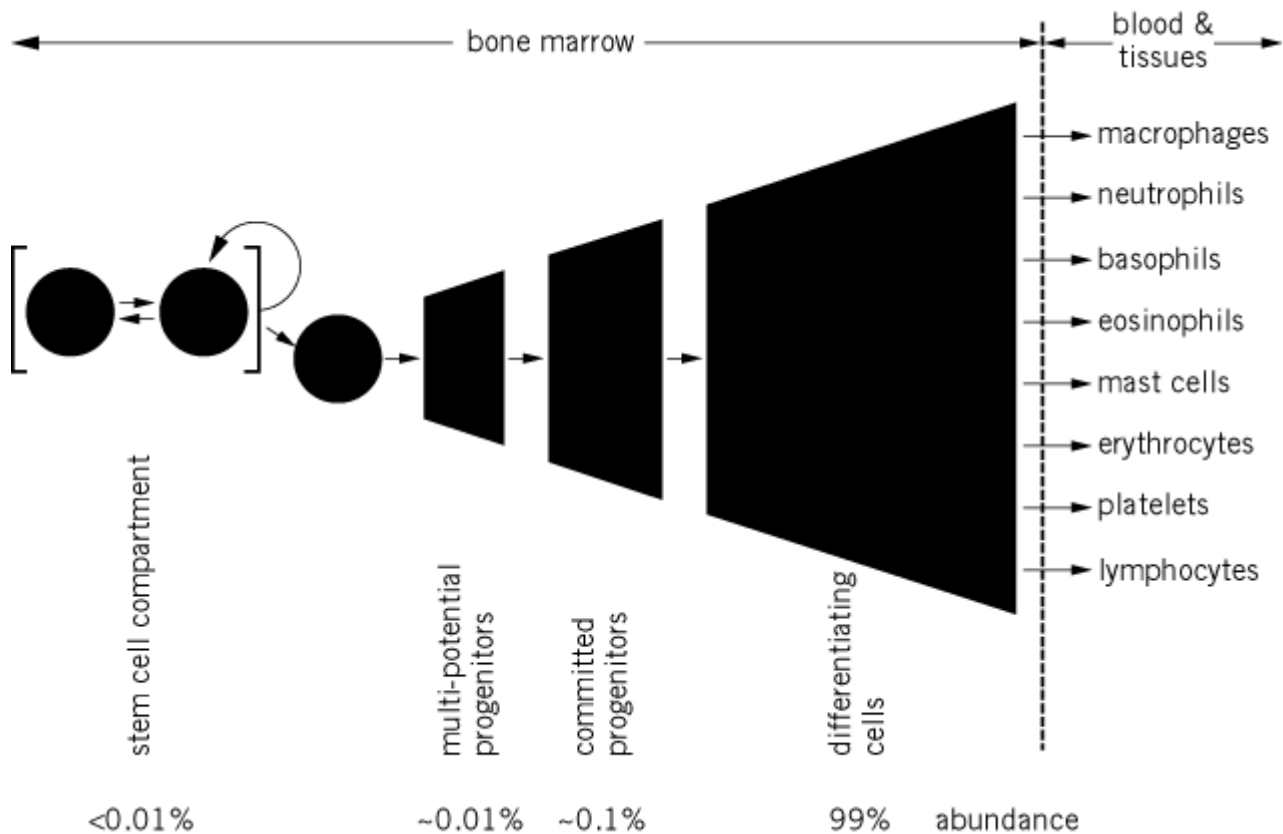
## Stem Cells

In multicellular organisms, certain tissues, such as those of the **hematopoietic** system and epithelia, are characterized by high cell turnover. Throughout life, **senescent** mature cells are shed or removed by phagocytosis or [apoptosis](#). These cells are continuously replaced by the proliferation and differentiation of immature cells that, in turn, arise continuously from small numbers of “stem cells” in the tissue. For example, it has been calculated that in steady-state hematopoiesis in the mouse,  $2.6 \times 10^8$  cells must be replaced each day. Transplantation experiments with purified stem cells in lethally irradiated mice have indicated that engraftment of as few as three such cells can bring about long-term hematopoietic reconstitution ([1](#)), illustrating their enormous potential for giving rise to differentiated progeny. Experiments in mice with cytotoxic agents (such as 5-fluorouracil) that selectively kill actively dividing cells have shown that, during steady-state hematopoiesis, the majority of primitive hematopoietic stem cells (HSCs) are quiescent, but are rapidly recruited into the [cell cycle](#) following depletion of more mature cells ([2](#)). More recent observations on *in vivo*-labeled stem cells purified by procedures based on their small size, lack of lineage markers, and low staining with vital dyes such as rhodamine 123 or Hoechst 33342 indicate that the earliest identifiable murine HSCs actually cycle at a low rate ([3](#)). The relative quiescence of early HSC accounts, at least in part, for the selectivity of chemotherapeutic regimens used in the treatment of hematological malignancies.

As stated above, stem cells have great proliferative capacity that is used to maintain differentiated cell populations throughout the life of the animal. Intrinsic to this process is the maintenance of the stem-cell pool, a process termed “self-renewal.” That is, primitive stem cells are believed to undergo asymmetric divisions, giving rise to a proportion of progeny with stem-cell characteristics and a proportion with the capacity to differentiate. Furthermore, the hematopoietic stem cell, and probably other stem cells, has the capacity to give rise to multiple lineages of differentiated cells (Fig. [1](#)); that is, it is “multipotential.” True “stem” cells are rare; in the case of human bone marrow, they

represent fewer than  $1/10^4$  of the total cells. More abundant are their immediate progeny, the hematopoietic progenitor cells, which display active cell cycling, although lower proliferative potential, little or no self-renewal capacity, and greater lineage restriction; they account for approximately  $1/10^3$  of normal bone marrow cells. The bulk of the cell population in normal bone marrow consists of postproliferative maturing cells, especially those of the neutrophil lineage.

**Figure 1.** Hierarchy of the hematopoietic system.



The following discussion will mostly use the hematopoietic system to illustrate the principles of stem cell behavior. An attempt will be made to illustrate the generality of the concepts by reference to keratinocyte differentiation. For more comprehensive treatment of a variety of other stem cell systems, including plants and invertebrates, see “Suggestions for Further Reading.”

## 1. Assays for Stem Cells

Because of their rarity, stem cells have generally been studied indirectly by their ability to give rise to colonies and/or differentiated progeny *in vivo* or *in vitro*. The “gold standard” for assay of HSCs is their ability, at limiting numbers, to give rise to multilineage engraftment of lethally irradiated animals, usually syngeneic mice. When highly purified primitive stem cells are used for marrow repopulation, co-infusion of more mature cells is necessary to ensure survival until long-term reconstitution occurs (4). Congenic markers, such as Ly5 **allotypes**, or sex-mismatched donors are used to distinguish between reconstitution by donor stem cells and that due to short-term repopulating cells or autologous recovery. Self-renewal capability is assessed by the ability of donor-derived cells to engraft in multiple serial transplants (5, 6). Indefinite self-renewal may not occur, due to [telomere](#) shortening (1, 7). The fate of individual stem cells and their progeny, including the ability to contribute to multiple hematopoietic lineages, has been monitored by following retrovirus-

marked stem cells that have unique insertion sites detectable by [restriction fragment length polymorphism \(RFLP\)](#) analysis (8). Another *in vivo* clonal assay, the colony-forming unit–spleen (CFU–S) assay, which depends on the ability of infused bone marrow cells to give rise to clones of maturing hematopoietic cells in the spleens of irradiated mice after 8 to 12 days, was used extensively in early studies, but is now considered to measure more mature progenitor or “transit” cells (9, 10).

*In vivo* assays with limiting numbers of cells are clearly not applicable in studies on human HSCs, which have generally been dependent on *in vitro* assays where cells are seeded onto supportive stroma (11), in semisolid medium (eg, Ref 12), or, more recently, in liquid culture supplemented with defined **cytokine** combinations (13). The nature of the cells that proliferate and differentiate in these assays has been reviewed (10). It is unclear to what extent they measure HSCs capable of long-term marrow repopulation, although quiescence, high proliferative potential, and the ability to give rise to large numbers of committed progenitors and/or differentiated cells of multiple lineages are taken as indicators of “stemness.” Recently, human stem and progenitor cell populations have been shown to be capable of giving multilineage engraftment of irradiated, immunocompromized mice (14). This may provide a useful model for studying early human HSCs *in vivo*.

## 2. Purification of HSCs

As described above, HSCs and progenitor cells are rare in hematopoietic tissues. Resistance to drugs that are selectively cytotoxic for cycling cells has long been used for enrichment of HSCs and very early progenitor cells (2). With the advent of fast and accurate multiparameter fluorescence-activated cell sorting (FACS; see [Flow Cytometry](#)), methods based on [light-scattering](#) characteristics, stains that correlate with quiescence/cycling, and multiple cell surface “differentiation” markers, detected by [monoclonal antibodies](#), have been developed (15, 16). This is an area of great interest and potential clinical application in separation of normal and malignant cells for autografting in cancer therapy, as well as for identification of appropriate targets for gene therapy, and a number of procedures for large-scale antibody-based cell selection have been developed (17). In the mouse, the HSC is generally considered to have the phenotype lineage marker  $(lin)^-$ ,  $thy-1^{low}$ ,  $Ly6A/E$  (Sca-1) $^+$ ,  $c-kit^+$  (15, 16), although recent data suggest that the earliest HSCs lack a detectable cell-surface  $c-kit$  (6). In humans, the most definitive phenotype for the HSC at this time appears to be  $CD34^+$ ,  $CD38^-$ , rhodamine<sup>DULL</sup>,  $c-KIT^+$  (14, 18, 19).

## 3. Regulation of Stem Cell Growth and Differentiation

The process of production of mature hematopoietic cells requires at least three sets of decisions regulating (a) entry of quiescent stem cells into the cell cycle, (b) whether a daughter cell resulting from stem cell division retains stem cell characteristics or differentiates, and (c) lineage commitment of differentiating cells. These processes clearly involve the actions of both positive and negative regulators produced by the hematopoietic microenvironment, but the basis of the action of such factors remains controversial. Two schools of thought exist. The first of these, the “stochastic” model (20), proposes that decisions relating to self-renewal versus differentiation and to lineage commitment are probabilistic and intrinsic to the stem cell. The hematopoietic microenvironment is proposed to provide a permissive milieu that enables cells to survive and fulfill their intrinsic potential (21). In contrast, the “deterministic” model proposes that external influences, such as cytokines, instruct the differentiation of hematopoietic stem cells (22, 23). The lack of specificity of the intracellular [signal transduction](#) pathways induced by receptors for different “lineage specific” cytokines, along with a lack of evidence for “deterministic” effects of cytokines on primary cells, has been taken as supporting the stochastic model (24). Furthermore, some evidence suggests that the developmental program of hemopoietic cells may be intrinsically driven by **homeobox** genes (7). Thus, the function of external regulators such as cytokines is proposed to be to ensure the survival of cells so that they can fulfill their intrinsic developmental potential. On the other hand, there is direct evidence for the influence of the nature and level of the cytokine on differentiation of individual

daughter cells from divisions of granulocyte-macrophage progenitor cells (25, 26). It is not possible at the present time to resolve this issue.

The process of induction of quiescent stem cells into cell cycle is also largely unknown and is of importance for stem cell expansion and gene therapy (1). Positive and negative regulators may be involved. For example, **transforming growth factor b** (TGFb), produced constitutively by bone marrow stroma and in an autocrine fashion by hematopoietic stem and/or progenitor cells, can inhibit proliferation *in vitro*. The levels of this cytokine are, however, unaltered in mice undergoing hematopoietic regeneration following irradiation (27). Recently, Haylock et al. (28), using a limiting dilution serum-free culture system, demonstrated that Flt3 ligand was capable of recruiting quiescent primitive human HSC into cycle. Similarly, regulation of self-renewal versus differentiation on HSC division is not understood, and argument continues as to whether it is a stochastic process or regulated by environmental factors (eg, see Ref. 7). So far, factors capable of supporting self-renewal of HSCs *in vitro* have not been identified.

Cytokines clearly play a crucial role in hematopoiesis, although, as described above, the precise mechanism is unknown. Positive regulators include (a) **interleukins-3** and -6, granulocyte-macrophage **colony-stimulating factor**, granulocyte colony-stimulating factor, and erythropoietin, which are ligands for members of the cytokine receptor family, and (b) stem cell factor (SCF), macrophage colony-stimulating factor (M-CSF), and Flt3 ligand, which activate members of the **receptor tyrosine kinase** family. Originally it was thought that distinct **growth factors** acted at different stages of hematopoiesis, but it is now known that individual factors act at multiple levels. For example, thrombopoietin and c-mpl ligand, which were thought to be specific megakaryocyte differentiation factors, are now known also to promote the growth of multipotential progenitor cells (29). Furthermore, cytokines act in a combinatorial fashion (30). A characteristic of primitive HSC is their requirement for multiple cytokines to promote their proliferation and differentiation *in vitro*, whereas committed progenitor cells require fewer factors (10). The actions of growth factors are antagonized by negative regulators, such as macrophage inflammatory protein (MIP)-1a and TGF-b; the latter may act in part by down-regulating the expression of cytokine receptors and/or by inducing apoptosis (27).

#### 4. Role of the Hematopoietic Microenvironment

In humans, hematopoiesis occurs exclusively in the supportive environment of the bone marrow, and HSCs infused intravenously home selectively to the marrow. The bone marrow is architecturally complex, with different microenvironments apparently supporting the development of different lineages. At least two classes of molecules produced by the marrow stroma are involved: (a) growth regulatory factors as described above and (b) molecules mediating the adhesion between stroma and hematopoietic stem and progenitor cells (31). These sets of molecules are interactive and, in some cases, multifunctional. For example, certain growth factors, such as (SCF) and Flt-3 ligand, exist as **transmembrane** proteins, as well as in soluble form, and may directly mediate adhesion by binding to their receptors on HSCs. Furthermore, binding of SCF to its receptor c-KIT was shown to up-regulate **integrin**-mediated adhesion of HSCs to the **extracellular matrix** protein, **fibronectin** (32). Indirect evidence that c-KIT may be involved in retaining stem and progenitor cells in the bone marrow comes from the observations that SCF administered systemically is a potent agent for mobilizing HSCs from the marrow to the peripheral blood, and, irrespective of the mobilizing agent, peripheral blood HSCs show apparent down-regulation of c-KIT (33). It is clear that both positive and negative regulatory influences on HSC development are mediated by contact with the stroma. Output of both total cells and committed progenitor cells in cytokine-supplemented stroma-free long-term HSC cultures far exceeds that observed in the presence of stroma (19). This may in part be due to the production of negative regulatory factors by the stroma, but it almost certainly also involves contact-mediated interactions. Novel molecules implicated in this process are the mucins expressed by the HSCs. Interaction of these molecules with their stromal ligands may negatively regulate HSC proliferation and/or differentiation and, under some circumstances, appears to trigger apoptosis (34, 35).

## 5. Generalization to Other Stem Cell Systems

Hematopoiesis is arguably the best understood system of somatic cell differentiation in mammals. Another system that has been extensively studied, probably because of the relative ease of culturing nontransformed cells, is keratinocyte differentiation. Like hematopoietic cells, epidermal cells are characterized by rapid turnover throughout the life of the animal, with replenishment from a small population of stem cells. These cells, which are located in the basal layer of the epidermis, have been identified *in vivo* based on their relative quiescence, in that they display long-term retention of radiolabeled DNA after infusion of  $^3\text{H}$ -thymidine, and on their great proliferative potential *in vivo* and *in vitro*. They comprise between 1% and 10% of the basal keratinocyte population in normal skin (36, 37). Intermediate between stem cells and their differentiating progeny are “transit amplifying” cells, which are actively cycling, have more limited proliferative potential, and may be the equivalent of the committed progenitor cells of the hematopoietic system. Until recently, characterization of epidermal stem cells has been hampered by the lack of markers to enable their isolation. Limited purification was accomplished based on their high expression of b1 integrins and rapid adhesion to type IV collagen (38), but the frequency of the cells isolated on the basis of these criteria suggested that the population also contained transit amplifying cells. More recently, a population of putative keratinocyte stem cells was isolated based on a high level of a6 integrin expression, combined with low expression of a cell-surface proliferation-associated marker (39). These cells, which represented approximately 10% of the a6 integrin-positive basal keratinocyte population, were quiescent based on cell cycle analysis, and they displayed extensive proliferative capacity *in vitro*, resulting in  $10^6$ -fold or more expansion on long-term culture. While cells with characteristics of quiescence and capacity for extensive proliferation, both *in vitro* and *in vivo*, have been demonstrated in interfollicular epidermis, other stem-cell attributes of multipotentiality and self-renewal capacity (based on serial transplantation) have not been demonstrated thus far. Cells with the capacity to give rise to keratinocytes of epidermal phenotype in raft cultures have been identified in the sheath of the hair follicle and are probably derived from the bulge region, but it is not known whether they are precursors of epidermal stem cells in normal skin. These cells may also serve as stem cells for sebaceous glands (36).

## 6. Embryonic Stem Cells

This article has focused on stem cells that persist throughout the life of mammals. The most primitive stem cell in these animals is the embryonic stem (ES) cell, which has unlimited proliferative potential under appropriate culture conditions *in vitro*, but retains the ability to give rise to all tissues *in vivo* and multiple cell types *in vitro* (eg, Ref. 40). Because they can be genetically manipulated *in vitro*, then introduced into blastocysts and implanted into surrogate mothers, ES cells have provided the means for generating “gene-knockout” mice that are playing a vital role understanding the function of genes and proteins *in vivo*. Furthermore, their ability to undergo multilineage differentiation *in vitro* provides an important system for molecular analysis of early events in development.

### Bibliography

1. M. A. S. Moore (1997) Stem Cells **15** (Suppl. 1), 239–251.
2. G. S. Hodgson and T. R. Bradley (1979) Nature **281**, 381–382.
3. G. B. Bradford, B. Williams, R. Rossi, and I. Bertonecello (1997) Exp. Hematol. **25**, 445–453.
4. S. J. Szilvassy, R. K. Humphries, P. M. Lansdorp, A. C. Eaves, and C. J. Eaves (1990) Proc. Natl. Acad. Sci. USA **87**, 8736–8740.
5. X. Q. Yan, Y. Chen, C. Hartley, P. McElroy, F. Fletcher, and I. K. McNiece (1998) Bone Marrow Transplant. **21**, 975–981.
6. H. Doi et al. (1997) Proc. Natl. Acad. Sci. USA **94**, 2513–2517.
7. P. M. Lansdorp (1997) Biol. Blood and Marrow Transplant. **3**, 171–178.

8. G. Keller, C. Paige, E. Gilboa, and E. F. Wagner (1985) *Nature* **318**, 149–154.
9. B. I. Lord (1997) In *Stem Cells* (C. S. Potten, ed.), Academic Press, London, pp. 401–422.
10. M. A. S. Moore (1991) *Blood* **78**, 1–19.
11. H. J. Sutherland, C. J. Eaves, A. C. Eaves, W. Dragowska, and P. M. Lansdorp (1989) *Blood* **74**, 1563–1570.
12. A. G. Leary and M. Ogawa (1987) *Blood* **69**, 953–956.
13. J. Brandt, E. F. Srour, K. van Besien, R. A. Bridell, and R. Hoffman (1990) *J. Clin. Invest.* **86**, 932–941.
14. J. E. Dick, M. Bhatia, O. Gan, U. Kapp, and J. C. Wang (1997) *Stem Cells* **15**(Suppl. 1), 199–203.
15. C. L. Li and G. R. Johnson (1995) *Blood* **85**, 1472–1479.
16. N. Uchida, A. S. Tsukamoto, D. He, A. M. Frieria, R. Scollay, and I. L. Weissman (1998) *J. Clin. Invest.* **101**, 961–966.
17. E. Wunder, H. Slovalat, P. R. Hénon, and S. Serke (1994) *Hematopoietic Stem Cells. The Mulhouse Manual*, AlphaMed Press, Dayton, OH, pp. 123–288.
18. L. W. Terstappen, S. Huang, M. Safford, P. M. Lansdorp, and M. R. Loken (1991) *Blood* **77**, 1218–1227.
19. P. J. Simmons, G. W. Aylett, S. Niutta, L. B. To, C. A. Juttner, and L. K. Ashman (1994) *Exp. Hematol.* **22**, 157–165.
20. J. E. Till, E. A. McCulloch, and L. Siminovitch (1964) *Proc. Natl. Acad. Sci. USA* **51**, 29–36.
21. C. J. Eaves, R. K. Humphries, and A. C. Eaves (1981) In *Hemoglobins in Development and Differentiation* (G. Stamatoyannopoulos and A. W. Nienhuis, eds.), A. R. Liss, New York, pp. 35–44.
22. J. J. Trentin (1971) *Am. J. Pathol.* **65**, 621–628.
23. D. Metcalf (1991) *Proc. Natl. Acad. Sci. USA* **88**, 11310–11314.
24. M. Socolovsky, H. F. Lodish, and G. Q. Daley (1998) *Proc. Natl. Acad. Sci. USA* **95**, 6573–6575.
25. D. Metcalf (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5327–5330.
26. D. Metcalf and A. W. Burgess (1982) *J. Cell. Physiol.* **111**, 275–283.
27. G. J. Graham (1997) *Bailliere's Clin. Haematol.* **10**, 539–559.
28. D. N. Haylock et al. (1997) *Blood* **90**, 2260–2272.
29. K. Kaushansky (1998) *Blood* **92**, 1–3.
30. D. Metcalf (1993) *Blood* **82**, 3515–3523.
31. P. J. Simmons, A. C. W. Zannettino, S. Gronthos, and D. I. Leavesley (1994) *Leuk. Lymphoma* **12**, 353–363.
32. J.-P. Levesque, D. I. Leavesley, S. Niutta, M. Vadas, and P. J. Simmons (1995) *J. Exp. Med.* **181**, 1805–1815.
33. L. B. To, D. N. Haylock, T. Dowse, P. J. Simmons, S. Trimboli, L. K. Ashman, and C. A. Juttner (1994) *Blood* **84**, 2930–2939.
34. V. Bazil, J. Brandt, S. Chen, M. Roeding, K. Luens, A. Tsukamoto, and R. Hoffman (1996) *Blood* **87**, 1272–1281.
35. A. C. W. Zannettino, H.-J. Bühring, S. Niutta, S. M. Watt, M. A. Benton and P. J. Simmons (1998) *Blood* **92**, 1–18.
36. S. J. Miller, R. M. Lavker, and T.-T. Sun (1997) In *Stem Cells* (C. S. Potten, ed.), Academic Press, London, pp. 331–362.
37. P. H. Jones (1997) *Bioessays* **19**, 683–690.
38. P. H. Jones, S. Harper, and F. M. Watt (1995) *Cell* **80**, 83–93.



39. A. Li, P. J. Simmons, and P. Kaur (1998) *Proc. Natl. Acad. Sci. USA* **95**, 3902–3907.
40. G. M. Keller (1995) *Curr. Opin. Cell Biol.* **7**, 862–869.

### Suggestions for Further Reading

41. C. S. Potten (1997) *Stem Cells*, Academic Press, London.
42. S. J. Morrison, N. M. Shah and D. J. Anderson (1997) Regulatory mechanisms in stem cell biology. *Cell* **88**, 287–298.

## Stereoisomers

Two molecules are stereoisomers if they differ in the spatial orientation of those atoms that cannot be rapidly interconverted by rotation about single bonds. The ability of two molecules to assume a common geometry “rapidly” differentiates conformers (see [Conformation](#)) from stereoisomers. Stereoisomers contain the same number and type of bonds and have the same chemical name, except for a prefix that is sometimes used to discriminate between the different stereoisomers. Stereoisomers are divided into [enantiomers](#) and [diastereomers](#). Enantiomers are molecules with nonsuperimposable mirror images, whereas diastereomers comprise all other stereoisomers ([1](#)).

### Bibliography

1. E. L. Eliel, (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York.

### Suggestions for Further Reading

2. J. March (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, Chap. "4".
3. K. Mislow (1966) *Introduction to Stereochemistry*, W. A. Benjamin, New York.
4. B. Testa (1982) "The geometry of molecules: Basic principles and nomenclatures". In *Stereochemistry* (C. Tamm, ed.) New Comprehensive Biochemistry, Vol. **3**, Elsevier, Amsterdam, pp. 1–48.

## Steroid Hormone Receptors

The gonads and adrenal gland produce five major groups of [steroid hormones](#): [estrogens](#), progestins, androgens, [glucocorticoids](#), and mineralocorticoids. Clinical symptoms associated with a reduced level of these hormones were known long before a concept of [hormones](#) was developed. Aristotle, for example, identified the effects of castration in humans and birds, and in 1855 Addison described the symptoms of chronic adrenal insufficiency. In 1896, Beatson observed, that female sex steroid hormones play a pivotal role in promoting mammary tumor growth. The same was found by Huggins and Hodges in 1941 for androgens and the growth of prostatic neoplasms. In the first half of the twentieth century, it became possible to isolate individual steroids from extracts derived from source tissues by employing biological complementation assays. In 1929 Butenandt was able to crystallize estrone from the urine of pregnant women, and in the 1930s Reichstein identified corticosteroids,

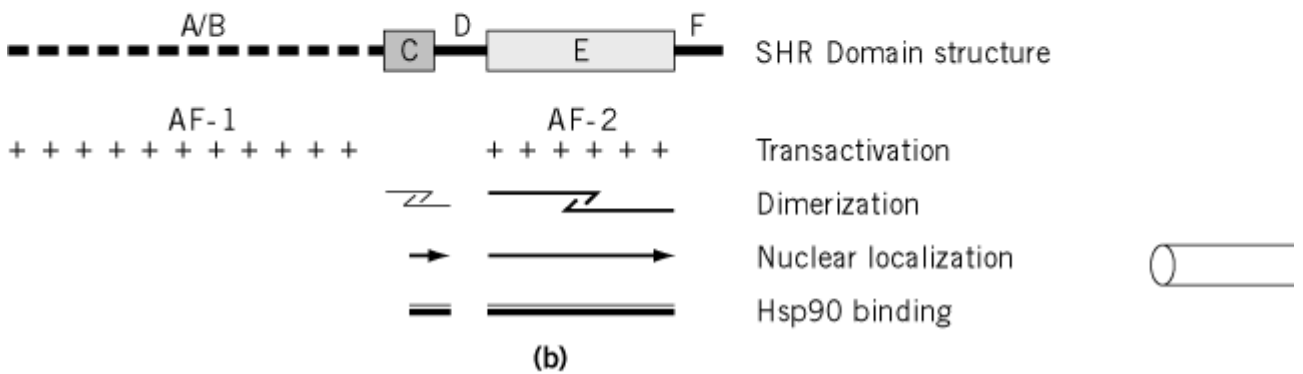
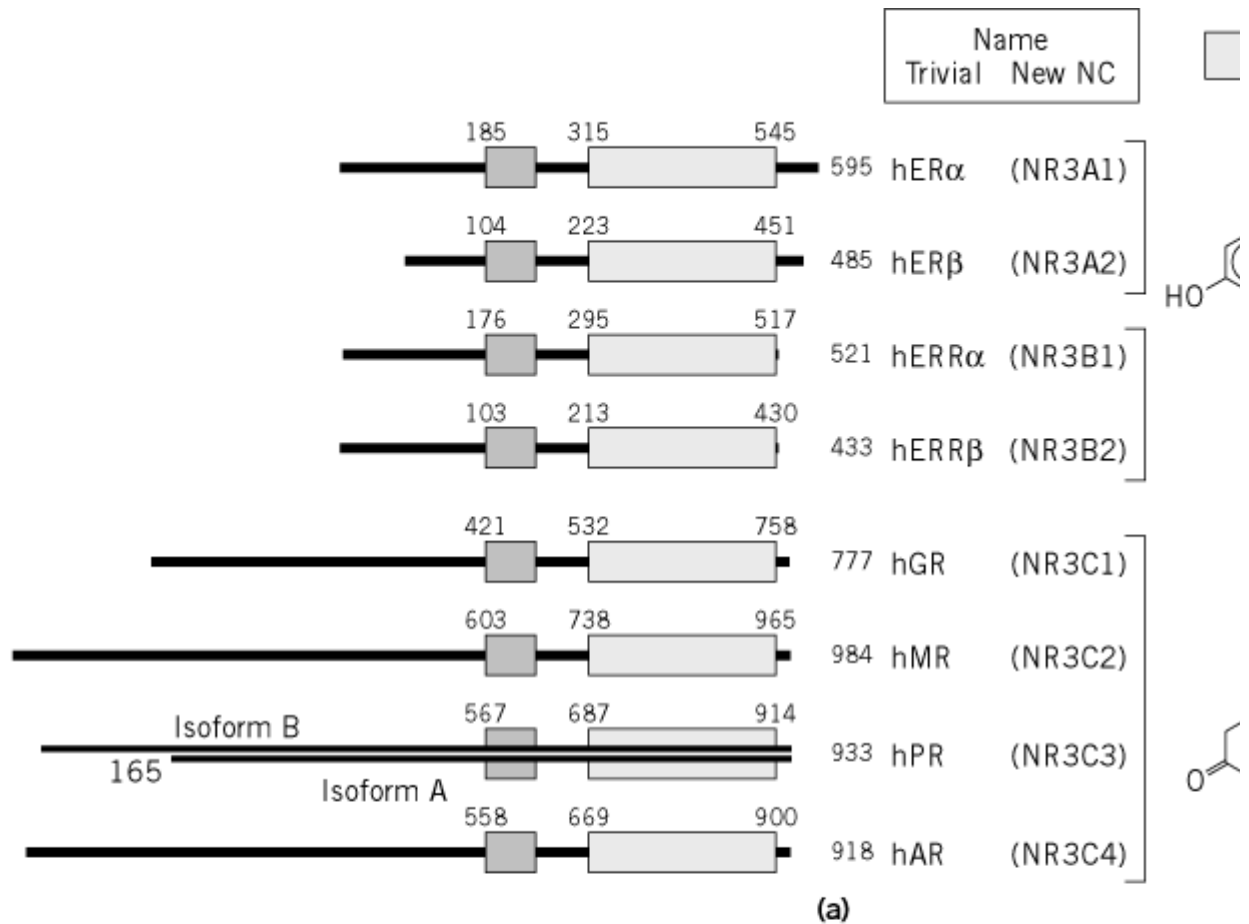
which gained popularity as anti-inflammatory and immunosuppressive therapeutic agents soon after (see [Dexamethasone](#)). The end of this period was marked by the isolation of aldosterone by Simpson and Tait in 1953. When <sup>3</sup>H- and <sup>14</sup>C-labeled hormone compounds became available in the late 1950s, it was possible to follow the fate of the steroid hormones up to their target tissues. Thus it could be shown that the hormone itself, not a metabolite, produced the response via the activation of **gene expression** mechanisms. The concept that steroid hormones are involved in control of **gene transcription** was also triggered by the observation of **ecdysone**-induced giant chromosome **puffs** by Clever and Karlson in 1960. About 10 years later, the groups of Jensen and Gorski (1) had established a two-step model that involved specific high-affinity steroid hormone receptors (SHRs) within the target cells that had to be activated by hormone in order to induce expression of hormone responsive genes (see [Glucocorticoid Response Element](#) Fig. 1). A further 10 years later, the **genes** for the hormone receptors for glucocorticoid and estrogen were **cloned**, and these receptors became the first known cellular **transcription factors** for **RNA polymerase II**. At around the same time, cloning of some SHR targets, such as the human [metallothionein](#) gene and the mouse mammary tumor virus [genome](#), led to the identification of specific binding sites for SHRs termed [hormone response elements](#) (**HREs**; see also [Glucocorticoid Response Element](#)) serving as inducible [enhancers](#). Within the past decade, many more facets of steroid receptors as hormone-activated transcriptional regulators have emerged and will be discussed briefly.

## 1. Nuclear Receptor Superfamily

The homology of the SHRs to the *v-erbA* **oncogene** led to the discovery of the *c-erbA* locus as the receptor for [thyroid hormones](#). The identification of another structurally related receptor for the vitamin A metabolite retinoic acid solidified the concept of a nuclear receptor [superfamily](#). All nuclear receptors are characterized by a central DNA-binding **domain** (see text below), which targets the receptor to the HRE. Among other approaches, use of this conserved DNA-binding domain as a screening probe under low stringency conditions proved very successful in the discovery of related receptors. Today, more than 60 different nuclear receptors have been described, including receptors for the known nuclear hormones and a vast number of so-called orphan receptors, with no, or unknown, and possibly novel ligands. With so many different receptors, nomenclature has become a real problem. Therefore, the majority of leading researchers in the field have now agreed to use a **unified nomenclature** (see text below) based on the system developed for the [cytochrome P450](#) gene superfamily (2). Most nuclear receptors bind as homo- or heterodimers to an HRE that is often present in two copies (half-sites), which can be arranged as direct, inverted (**palindromic**), or everted repeats; the half-sites are separated by one to eight nucleotides. On the basis of their dimerization behavior, the nuclear receptors can be operationally placed into four subgroups: (1) the homodimerizing steroid receptors, (2) the receptors forming heterodimers with the retinoid-X receptor (RXR), and the (3) dimeric and (4) monomeric orphan receptors. This grouping is also supported by **phylogenetic** analyses (2) showing that the ligand-binding ability appears to be acquired by orphan receptors during the course of evolution, just as the ability to heterodimerize appears to be acquired by homodimerizing receptors. Nuclear receptor diversification probably occurred in two steps: one that took place very early, during the multicellularization event leading to the metazoan phyla, and resulting in the receptor subgroups we know today, and a second step occurring later on, corresponding to the advent of vertebrates and leading to the formation of multiple **paralogous** receptors within the existing subgroups. Therefore, the steroid and steroid-like receptors (nuclear receptor subfamily 3, NR3, with its three groups A, B, and C according to the new nomenclature system) depicted in Figure 1a must be viewed as a late “invention” of [evolution](#). They include two receptors for estrogens (hERa and hERb or NR3A1 and NR3A2), two estrogen-related orphan receptors (hERRa and hERRb, or NR3B1 and NR3B2), and one receptor each for the other major steroid hormone classes: glucocorticoids (hGR or NR3C1), mineralocorticoids (hMR or NR3C2), progestins (hPR or NR3C3), and androgens (hAR or NR3C4). hERRa and b are the only receptors of the subfamily that bind to DNA as monomers rather than as homodimers. Notably, the receptors for secosteroids (vitamin D) and ecdysteroids (see [Ecdysone](#)) do not belong to this group. They form heterodimers with RXR or RXR-like proteins, such as the very recently identified

pregnane X receptor (PXR), which binds synthetic glucocorticoids and pregnenolone derivatives.

**Figure 1.** The human steroid hormone receptor family. (a) Illustrated primary structures of all family members, cognate numbers refer to amino acid residue positions. Highlighted are the DNA-binding domain (DBD) and the ligand-binding domain (LBD). Extensions are drawn to scale. The name of each receptor is indicated (trivial name and the abbreviation according to the human; ER, estrogen receptor; ERR, estrogen-related receptor; GR, glucocorticoid receptor; MR, mineralocorticoid receptor; AR, androgen receptor). (b) SHR domain structure and structure–function relationships. The A/B, C, D, E, and F segments w sequence conservation amongst different SHRs. AF-1 and -2 are the transcription activation functions 1 and 2. (c) Illustration of cognate hormone response element (HRE). The HRE is of an inverted repeat (IR) type (palindrome), and the half sites (a (N) nucleotides (IR3-HRE). The hormone ligand is indicated with a sterane formula.



## 2. Domain Structure

All steroid hormone receptors are modular proteins composed of distinct regions, as shown in Figure 1b, top panel (3). When the chicken estrogen receptor cDNA became available, comparison with the human estrogen and glucocorticoid receptor sequences led to the nomenclature based on regions A–F, which is used for all members of the nuclear receptor superfamily. From numerous functional and structural analyses, it became clear that these distinct regions correspond to functional and structural units called **domains**. Region C, the **DNA-binding** domain, and region E, the **ligand-binding** domain, display a high degree of sequence conservation, whereas no significant conservation is detected between the paralogous SHRs for the regions A/B, D, and F. Region D is considered as a flexible hinge region between the DNA- and ligand-binding domains. Its amino terminus is an integral part of the DNA-binding domain and is involved in its dimerization. The F region, now defined as the C-terminus beyond [alpha-helix 12](#) of the ligand-binding domain (see text below), is essential for hormone binding in the receptors for progestins (PR), glucocorticoids (GR), and androgens (AR), but not in ERα. Recent work indicates that this domain is also important for the discrimination between agonistic and antagonistic hormone ligands. Region A is highly conserved only between chicken and human estrogen receptors, but this distinction is much less clear in the other steroid receptors. Therefore, regions A and B are combined to an A/B region in most cases.

Regions C and E are not only responsible for DNA and ligand binding, respectively, and they encode other functions as well (see Fig. 1b, lower panel). The intracellular distribution of steroid receptors is the result of nuclear–cytoplasmic diffusion and cytoplasmic–nuclear shuttling. At equilibrium, the majority of steroid receptors is in the nucleus due to the presence of so-called nuclear localization signals (NLSs) that are believed to be required for nuclear pore recognition (see [Nuclear Import, Export](#)). A constitutive NLS is located at the border of region C/D, whereas a second NLS in the ligand-binding domain is ligand-dependent. Intracellular localization is less clear for the glucocorticoid and mineralocorticoid receptors because hormone-induced nuclear translocation has been reported in these cases.

All unliganded SHRs are associated with a large multiprotein complex of **molecular chaperones**, including Hsp90, which maintains the receptor in an inactive state but keeps it well prepared for hormone binding. Again, the region at the C/D border together with the ligand-binding domain are required for an SHR–Hsp90 interaction to take place. Most likely these chaperones play an active role in keeping the SHRs functional.

Initially all SHRs were thought to bind as homodimers to their cognate response elements within the **promoters** of target genes, which is still a specific characteristic of this subfamily within the nuclear receptor superfamily. But with the identification of the estrogen-related receptors a and b, the family has now two orphan members that do bind to DNA as monomers. Moreover, it is well known that the predominant form of the glucocorticoid receptor is monomeric in solution and that dimerization occurs only after binding to an HRE, accounting for the **cooperativity** observed during DNA-binding. There are also reports showing that the glucocorticoid receptor can bind to an HRE half-site as a monomer, although with lower affinity. A weak dimerization domain encoded within the DNA binding domain and a strong dimerization interface provided by the ligand-binding domain are responsible for dimerization. It is generally believed that the glucocorticoid receptor only contains the weak dimerization interface within the DNA-binding domain.

Ligand binding confers transcriptional competence onto SHRs, which is exerted in most receptors by two independent **transactivation** functions: (1) a constitutively active one in the A/B region located close to the DNA-binding domain, referred to as *activation function 1* (AF-1; also known as q-1 or enh-1) and (2) a ligand-inducible activation function in the ligand-binding domain: AF-2. The long progesterone receptor isoform, PR-B (see text below), harbors a third activation function at its amino terminus, AF-3, which is specific for the promoter and the cell. The AFs connect the receptor to the transcription apparatus via direct interactions with basal transcription factors, sequence-specific transcription factors, and/or transcription intermediary factors (see text below).

### 3. Family Members

The human estrogen and glucocorticoid receptors were the first steroid hormone receptors which genes were cloned. The estrogen receptor (now called *hERa*; Fig. 1a) was cloned from a [cDNA library](#) produced from the breast cancer cell line MCF-7. Recently, a second estrogen receptor, ERb, was cloned which is very similar to ERa in terms of structure and function but also shows subtle and important functional differences (see [Estrogen Receptors](#)). Both receptors bind the ligands estradiol, diethylstilbestrol, estriol, and estrone with high affinity. cDNA clones for the estrogen-related receptors a and b (*hERRa* and b) were isolated from a human testis cDNA library, using the *hERa* DNA-binding domain as a probe.

For the glucocorticoid receptor, two different classes of cDNA have been described, *hGRa* and b, which are the result of [alternative splicing](#) from a single gene transcript. Both **isoforms** are identical up to amino acid residue 727 and then diverge; *hGRa* is slightly larger (777 amino acids) than *hGRb* (742 amino acids). In contrast to *hGRa*, which is commonly considered as the bona fide hGR, *hGRb* was long dismissed as a cloning artifact, but it is now shown to be expressed at modest but varying levels in a range of tissues. *hGRb* does not bind hormone, is transcriptionally inactive, and, therefore, acts as a ligand-independent negative regulator of glucocorticoid action in [transfection](#) experiments. However, the b-isoform is not conserved across species, and the relative expression of both isoforms is not known. *hGRa* shows high affinity for the artificial glucocorticoid dexamethasone, moderate affinity for the physiologic steroids cortisol and corticosterone, and low affinity for mineralocorticoids and progesterone.

The progesterone receptor was the first SHR shown to exist in two common isoforms generated by differential use of alternative promoters. One promoter initiates transcription at positions +1 and +15 of the gene, which gives rise to the longer isoform B (see Fig. 1a). The second promoter initiates *hPR* transcripts between nucleotides +737 and +842, to encode the 164-residue shorter *hPR* form A. Both isoforms display different target gene specificities because they differ in their A/B domains and, hence, in their transactivation functions (see text above). First regarded as an exception in the family, differential promoter usage and alternative splicing are now the rule for all members. Both progesterone receptor isoforms show high affinity for the natural ligand progesterone and the synthetic agonist R5020.

The human mineralocorticoid (*hMR*) and androgen receptors (*hAR*) were cloned 2–3 years after the glucocorticoid receptor had been described. The androgen receptor exists in two isoforms, *hAR-A* and *hAR-B*, which are structurally analogous to the two *hPR* isoforms. In contrast to *hPR-A/B*, however, *hAR-A* is expressed at substantially lower levels than the B form, and its contribution to androgen action is not known. The androgen receptor binds the two naturally occurring ligands dihydrotestosterone and testosterone with high affinity, whereas the mineralocorticoid receptor shows high and equivalent affinity for aldosterone, physiologic corticosteroids, and progesterone.

#### 4. Genomic Organization

All the genes encoding SHRs have multiple **introns** and exons and range in length from at least 40 kbp for the human ERb gene to more than 140 kbp for the human ERa gene. Most genes contain 8 exons, with the positions of the introns being strictly conserved, and are scattered over the human genome (see Table 1). Only *ERRa* and progesterone receptor, and ERb and *ERRb*, colocalize to overlapping regions of chromosomes 11 and 14, respectively. Compared to the wealth of information on SHR structure–function relationships, relatively little is known about the regulation of the SHR genes themselves. Of the *hERa* gene, it is known that it is regulated by three tissue-specific promoters. As a result, *hERa* transcripts from various tissues can differ in the proximal part of their 5'-untranslated regions, depending on which promoter is used.

**Table 1. Chromosomal Location of Human Steroid Hormone Receptor Genes**

---

| Gene <sup>a</sup> | New Nomenclature | Chromosome Location            |
|-------------------|------------------|--------------------------------|
| ERa               | (NR3A1)          | 6q25.1                         |
| ERb               | (NR3A2)          | 14q22-q24                      |
| ERRa              | (NR3B1)          | 11q12-q13                      |
|                   |                  | 13q12.1 (processed pseudogene) |
| ERRb              | (NR3B2)          | 14q24.3                        |
| GR                | (NR3C1)          | 5q31-q32                       |
| MR                | (NR3C2)          | 4q31.1                         |
| PR                | (NR3C3)          | 11q13-q22                      |
| AR                | (NR3C4)          | Xq11.2-q12                     |

---

<sup>a</sup> ER, estrogen receptor; ERR, estrogen-related receptor; GR, glucocorticoid receptor; MR, mineralocorticoid receptor; PR, progesterone receptor; AR, androgen receptor.

## 5. DNA Binding

The receptors for progestins, glucocorticoids, mineralocorticoids, and androgens recognize the same DNA sequence (half-site: AGAACA). This sequence deviates in two nucleotides from the DNA-sequence bound by both estrogen receptors (AGGTCA). The nonestrogen receptors recognize their half-site mainly by a **hydrophobic** interaction with the methyl group of the thymine in position 4 of the complementary strand (AGAACA) that is not present in the DNA-sequence recognized by the estrogen receptor (4). The half-site spacing of 3 base pairs is such that an SHR homodimer binds both half-sites on the same face of the DNA (see Fig. 1c). The DNA-binding domains of SHRs comprise approximately 80 residues encoded by region C, plus some 14 N-terminal residues of region D. The domain contains two patterns that are reminiscent of, but clearly distinguishable from, the **zinc finger** motifs first observed in the *Xenopus laevis* transcription factor IIIA. The SHR zinc fingers are also able to coordinate a zinc atom tetrahedrally, but are of the type (C<sub>2</sub>-C<sub>2</sub>), with four **cysteine residues** as ligands (5). Only very few amino acid residues within the SHR zinc fingers are responsible for specific recognition of the cognate HRE (see [Estrogen Receptors](#) for details).

## 6. Ligand Binding

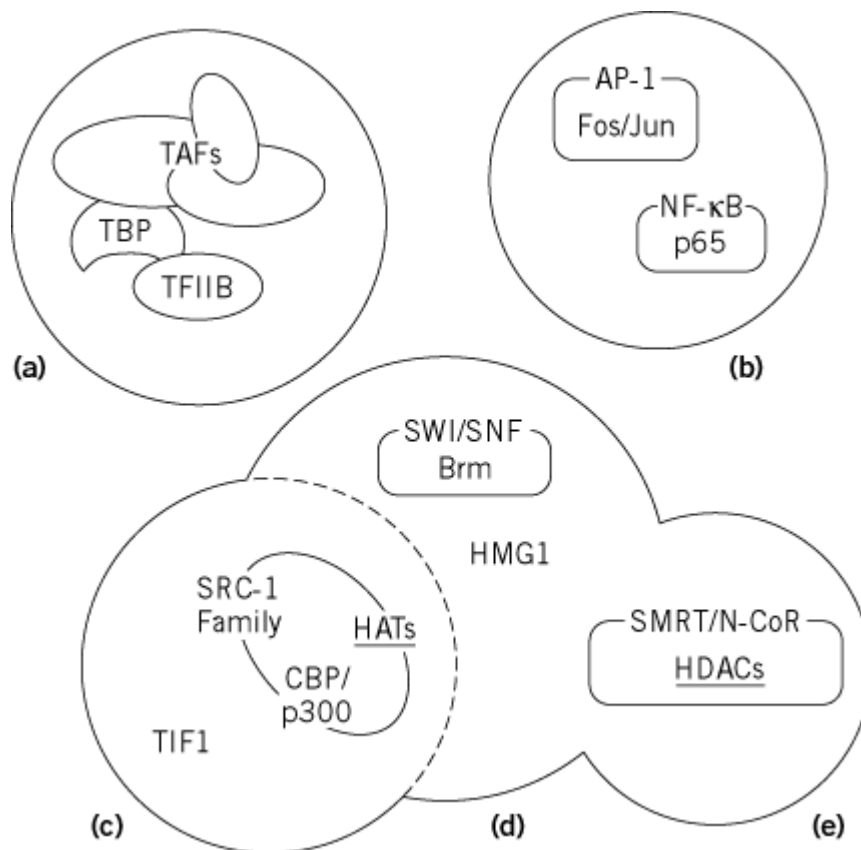
Although not proved, it is believed that the lipophilic steroid hormones, as well as synthetic antihormones, enter the target cell by simple [diffusion](#) and bind, within the cytoplasm, to a multiprotein complex of chaperones and SHR. Hormone binding induces a transformation of the SHR complexes that is associated with an increase in their affinity for DNA and a decrease in size of the complex. For example, it is well known that the unliganded hERa receptor sediments as an 8S complex on low salt sucrose [density gradient centrifugation](#), but it can be activated with high salt, temperature, or hormone ligand, to yield a more compact and **proteinase**-resistant 4S form. Hormone-induced transformation of the 8S receptor complex appears to reflect loss of the associated chaperone Hsp90, followed by a tight association with the nuclear compartment. For hERa (PDB accession nos. 1ERE and 1ERR; see [Structure Databases](#)) and hPR (PDB accession no. 1a28), the three-dimensional structures of their ligand binding domains complexed with 17 $\beta$ -estradiol and the antihormone raloxifene or with progesterone, respectively, have been determined. Like other known ligand-binding domains of nuclear receptors, those of SHRs are folded into a three-layered, antiparallel  $\alpha$ -helical sandwich that creates a wedge-shaped molecular scaffold, with the ligand-binding cavity at the narrower end of the domain. This cavity is completely partitioned from the

external environment and is closed by  $\alpha$ -helix 12, operating as a “lid” after hormone has entered the binding pocket. The relocation of this [amphipathic](#) helix over the hormone binding site generates a new surface that allows coactivators to bind to the ligand-binding domain, thereby mediating the activity of activation function 2, which forms the core of helix 12. The ligand-binding domain dimerization interfaces are formed by  $\alpha$ -helices that line up and/or intertwine and differ slightly between ER $\alpha$  and progesterone receptor.

## 7. Activation of Transcription by Steroid Hormone Receptors

To regulate transcription, agonist-liganded SHRs must communicate with the general transcription machinery (6). This can be achieved either by a direct contact between SHRs and the general transcription factors (GTFs; see [Transcription](#)), which are building up the transcription preinitiation complex at the transcription starting point, or by means of transcriptional intermediary factors (TIFs; also called *coactivators*, *mediators*, or *bridging factors*), other sequence-specific transcription factors ([Transcription Factors](#)) or **chromatin factors** that function as interpreters between SHRs and the general transcription machinery (Fig. 2).

**Figure 2.** Nuclear partners of steroid hormone receptors. Typical examples are shown for each class. In the case of the general transcription factors (a), the TATA-box-binding protein (TBP), the general transcription factor IIB (TFIIB), and some of the TBP-associated factors (TAFs), all parts of the transcription preinitiation complex, were shown to interact with SHRs. Of the interacting sequence-specific transcription factors (b), listed are the activator protein-1 (AP1) with its main subunits Fos (named after the oncogene identified in the FBJ murine osteogenic sarcoma virus) and Jun (named after the avian sarcoma virus 17; *jun* is the Japanese word for 17), and the nuclear factor  $\kappa$ B (NF- $\kappa$ B) with its p65 subunit. Both factors are key mediators of inflammatory and immune responses. Important coactivators (c) are the transcription intermediary factor 1 (TIF1), the steroid receptor coactivator (SRC-1) family of proteins, and the general coactivators cAMP response element binding protein-binding protein (CBP), and a closely related factor, p300. Of the chromatin factors (d), brahma (Brm), a component of the very large SWI/SNF-complex (SWI stands for mating-type switching and SNF for sucrose nonfermenting) and the HMG (high-mobility group) protein 1 (HMG1) are mentioned. The silencing mediator for retinoid and thyroid hormone receptors (SMRT) and the nuclear receptor corepressor (N-CoR), (e) as well as the listed coactivators, are all involved in chromatin remodeling, through either histone acetyltransferase (HAT, underlined) or histone deacetylase (HDAC, underlined) activities, or direct contacts to structural chromatin components (TIF1). Therefore, chromatin remodeling seems to be a key process during steroid hormone receptor action.



For example, hER $\alpha$  has been shown to interact with the **TATA-box**-binding protein (TBP), the TBP-associated factor hTAF<sub>II</sub>30, and the general transcription factor TFIIB. Likewise, the human progesterone receptor was shown to interact with *Drosophila* TAF<sub>II</sub>110 and TFIIB. Because an interaction with the same GTFs has also been documented for many other sequence-specific transcription factors, the specificity of an individual interaction remains to be determined.

### 7.1. Transcription Intermediary Factors (TIFs)

TIFs are thought to bridge between DNA-bound sequence-specific transcription factors and GTFs. The suspicion that SHRs may require such TIFs for activation was derived from the squelching phenomenon. “Squelching” refers to the observation that excess receptor can inhibit its own transactivation, as well as transcription by other transactivators. One possible explanation for this behavior is that additional factors, present in limited amounts and required for transactivation, are trapped by excess SHRs in unproductive complexes.

Today many of these transcription intermediary factors are known (see Fig. 2). One of the best characterized examples is the steroid receptor coactivator-1 (SRC-1) (7). SRC1 interacts with the human progesterone, glucocorticoid, and estrogen- $\alpha$  receptors and enhances their hormone-dependent transcriptional activities, without altering the basal activity of a target promoter. Moreover, by coexpressing SRC-1, it is possible to reactivate a target promoter that is squelched by excess SHR. Another important and general coactivator shown to interact with SHRs and SRC-1 is the cAMP response-element-binding protein, (CREB)-binding protein (CBP), and the related protein p300 (see text below).

For the binding of at least a subset of these TIFs to agonist liganded SHRs, the presence of one or more TIF signature motifs, Leu-X-X-Leu-Leu, are necessary and sufficient; in all known TIFs, at least one of these signature motifs is present in those parts of the proteins that were shown to interact with the SHRs. **Secondary structure prediction** analysis shows that the Leu-X-X-Leu-Leu



signature motif is part of an  $\alpha$ -helix. After binding of an agonistic ligand to the SHR, helix 12, containing the AF-2 core of the receptor, is relocated (see text above), thus creating a new interface that is poised for interaction with the signature motif of a TIF.

## 7.2. Corepressors

For some members of the nuclear receptor family, such as the thyroid hormone receptor (TR) and the retinoid acid receptors (RARs), which form heterodimers with the retinoid X receptor (RXR), it is well known that they repress basal transcription in the unliganded state. Recently, it could be shown that this silencing effect is mediated by [corepressors](#) that were called nuclear receptor corepressor (N-CoR) and silencing mediator for retinoid and thyroid hormone receptors (SMRT). N-CoR and SMRT associate with DNA-bound TR/RXR or RAR/RXR heterodimers in the absence of hormone, but not in its presence. SMRT and N-CoR are multiprotein complexes that exhibit [histone](#) deacetylase activity. The same proteins are also implicated in steroid hormone action, as it was demonstrated that overexpression of N-CoR and SMRT represses the partial agonistic activity of E $\alpha$  bound to the antihormone tamoxifen and of progesterone receptor bound to the antagonist RU486. Because all previous investigations in this field have been done by using artificial assays, including yeast [two-hybrid systems](#) and transfection experiments, these results must be considered with caution in terms of their *in vivo* relevance.

In summary, SHRs can exist in a multitude of conformations depending on the nature of their bound ligand. The various interaction surfaces are determined by  $\alpha$ -helix 12, and most probably also the following C-terminal residues, the F region. Transactivation by SHRs appears not to be a simple net activation, but the sum of the relief from repression by corepressors and activation by coactivators. The switch from the repressed to the activated state is promoted by the hormone ligand through an **allosteric** change in the SHR structure.

## 7.3. Interactions with Sequence-Specific Transcription Factors

SHRs also control the activities of natural promoters through positive and negative interactions with sequence-specific transcription factors. Particularly interesting is an interaction that represses the activity of both partners, as exemplified in the case of human glucocorticoid receptor and the heterodimeric transcription factor AP-1. Typically, this mutually inhibitory interaction occurs with composite binding sites, specifically, overlapping or closely spaced hormone response elements and AP1 sites, or on separate HRE or AP1 sites, and depends on the precise composition of the AP1 dimer. Similar interactions have been described between glucocorticoid receptor, the p65 subunit of the transcription factor NF- $\kappa$ B, and the transcription factor GATA-1. From experiments with numerous SHR mutants, it can be concluded that repression is most probably mediated by a [protein–protein interaction](#) with an SHR monomer that requires the N-terminal portion of the receptor at least (the N-terminal portion of an SHR monomer). Because most immunoregulatory genes, as well as genes involved in inflammation, are positively regulated by the transcription factors AP1 and NF- $\kappa$ B, it is conceivable that the immunosuppressive and anti-inflammatory activities of glucocorticoids are mediated through inhibition of AP1- and NF- $\kappa$ B-mediated transactivation by glucocorticoid receptor.

## 8. Crosstalk with Other Signal Transduction Pathways

Steroid hormone receptors are phosphoproteins, and their functions are modified not only by hormone, but also by **phosphorylation**. Most of the identified phosphorylation sites are [serine](#) and [threonine](#) residues, but some of the family members are also phosphorylated on [tyrosine](#). In case of the chicken progesterone receptor, for example, four phosphorylated serine residues have been identified that are common to both isoforms PR-A and PR-B. The two N-terminal sites are only moderately phosphorylated in the absence of hormone, whereas an increase in phosphorylation of these sites, and the appearance of two new phosphorylation sites, are observed after hormone treatment. Mutation of these serine residues results in a cell- and promoter-specific variation of receptor activity when tested in transfection experiments.

Modulation of [kinase](#) activity can also cause activation of some SHRs, even in the absence of hormone. [Epidermal growth factor](#) (EGF) treatment of ovariectomized mice, for example, results in nuclear translocation and an altered phosphorylation state of the estrogen receptor ER $\alpha$ . In transient transfection studies, EGF can activate the ER $\alpha$ . In contrast, it was found in ER-containing MCF-7 mammary tumor cells that estradiol stimulates within minutes the c-Src kinase and mitogen-activated protein (MAP) kinase signal transduction pathways. Here, ER $\alpha$  was shown to interact directly with c-Src. Likewise, in the progesterone receptor-positive T47D mammary tumor cell line, progestins rapidly and reversibly stimulate the c-Src/p21ras/Erk-2 pathway. This activation requires not only PR-B but also ligand-free ER $\alpha$  (8). In contrast to ER $\alpha$ , PR-B does not interact with c-Src, but with ER $\alpha$  via its amino terminal part, which is not present in PR-A. The transactivation function of PR-B is not required to activate this signal transduction pathway. Therefore, extensive crosstalk takes place between peptide [growth factors](#) and SHRs, which could explain the mitogenic effects of steroid hormones and the steroidlike effects of growth factors.

## 9. Interaction with Chromatin

The interaction of SHRs with DNA, the general transcription machinery, corepressors and coactivators, as well as sequence-specific transcription factors, takes place in the nucleus with its DNA compacted into [chromatin](#) (see [Nucleosome](#)). Genetic analyses have demonstrated a widespread involvement of chromatin structure in gene regulation in general. For the glucocorticoid receptor, it was shown that components of the so-called SWI/SNF complex, a set of pleiotropic transactivators that counteract repressive functions of chromatin, are required for transactivation in yeast. In human cells lacking brahma (hBrm), the homologue of yeast SWI2, transactivation by glucocorticoid receptor is weak and can be selectively enhanced by ectopic expression of hBrm. Similar to SWI2 in yeast, hBrm is part of a large multiprotein complex that mediates ATP-dependent disruption of a nucleosome. These results document the link between SHRs and the cellular machinery involved in chromatin dynamics.

One well-documented example for SHR-chromatin interactions is the mouse mammary tumor virus (MMTV) promoter. The MMTV promoter is transcriptionally controlled by steroid hormones, in particular glucocorticoids and progestins. The cognate hormone receptors bind to a cluster of hormone response elements (HREs) and facilitate the binding of other transcription factors, including nuclear factor 1 (NF1). Many investigators find nucleosomes in defined positions over the MMTV promoter. One dominant nucleosome phase found both in mammalian and yeast cells carrying the MMTV promoter permits SHR binding to the HREs, while precluding binding of NF1. This difference between SHRs and NF1 seems to reflect the way in which a [DNA-binding protein](#) recognizes the major groove of its cognate DNA site. Whereas SHR–DNA interactions take place largely on one face of the DNA double helix, and within a small sector of the helix circumference, many sequence-specific transcription factors, including NF1, exhibit more contacts around the helix circumference and almost completely embrace the double helix. This means that an HRE positioned on the surface of a nucleosome is still accessible to SHRs, so long as the major groove recognized by the receptor is pointing away from the nucleosome. But in the case of NF1, no matter what the rotational orientation of its binding site on the nucleosome surface might be, a fraction of the required protein–DNA contacts is always occluded by the histones, and binding is not possible. Therefore, the SHRs might represent a class of transcription factors whose cognate binding sites are still accessible within chromatin, if the rotational HRE orientation is favorable, and thus can act as initiators of chromatin remodeling so that other transcription factors, exemplified by NF1, can bind in a second step.

Apart from the SWI/SNF complex and related ATP-dependent chromatin remodeling machines, increasing evidence has accumulated in recent years implicating the modification of [histones](#) in the regulation of gene transcription. Histones are composed of a histone-fold domain, involved in wrapping the DNA in nucleosomes, and an *N*-terminal tail rich in [lysine](#) residues that is protruding out of the nucleosome. **Acetylation** of the tail lysine residues greatly reduces the affinity of the histone tails for DNA, and it is believed to render DNA more accessible to transcription factors,

while still maintaining a nucleosomal architecture. Moreover, some of the identified coactivators for the SHRs such as CBP/p300 and SRC-1 (see Fig. 2), as well as the largest subunit of the pivotal general transcription factor TFIID, TAF<sub>II</sub>230 / 250, turned out to be histone acetyl transferases. All these observations lead to the **current two-step model** for transcriptional activation by SHRs that is built on the recruitment of coactivators and other transcription factors with histone acetyl transferase activity, resulting in the local destabilization of repressive histone–DNA interactions, followed by direct or indirect interactions with the basal transcription machinery.

### Bibliography

1. J. Gorski and F. Gannon (1976) *Annu. Rev. Physiol.* **38**, 425–450.
2. V. Laudet, J. Auwerx, E. Baulieu, M. Beato, M. Becker-Andre, et al. *EMBO J.* in press.
3. R. M. Evans (1988) *Science* **240**, 889–895.
4. M. Beato (1989) *Cell* **56**, 335–344.
5. J. W. Schwabe and D. Rhodes (1991) *TIBS* **16**, 291–296.
6. M. Beato and A. Sanchez Pacheco (1996) *Endocr. Rev.* **17**, 587–609.
7. S. A. Onate, S. Y. Tsai, M. J. Tsai, and B. W. O'Malley (1995) *Science* **270**, 1354–1357.
8. A. Migliaccio, D. Piccolo, G. Castoria, M. Di Domenico, A. Bilancio, et al. (1998) *EMBO J.* **17**, 2008–2018.

### Suggestions for Further Reading

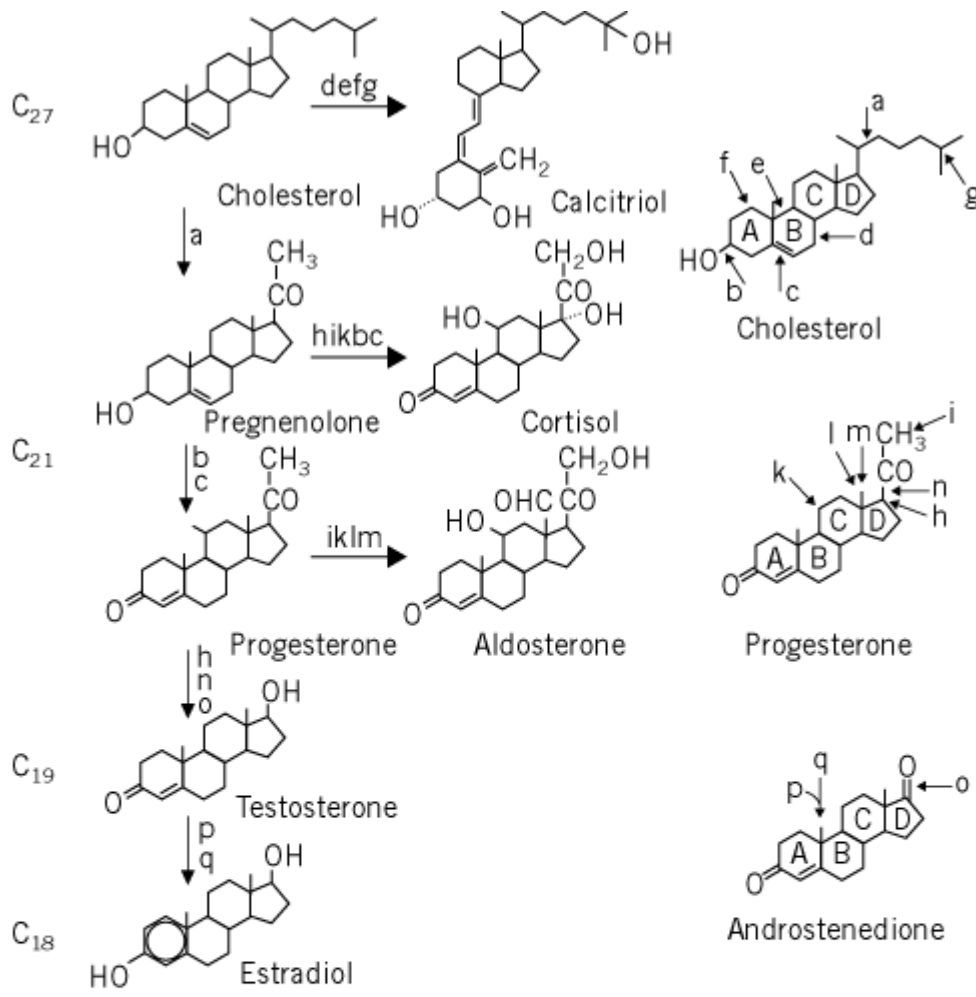
9. E.-E. Baulieu and P. A. Kelly (1990) *Hormones. From Molecules to Disease*, Chapman and Hall, New York and London.
10. M. G. Parker (1991) *Nuclear Hormone Receptors. Molecular Mechanisms, Cellular Functions, Clinical Abnormalities*, Academic Press, London.
11. M. Beato, P. Herrlich, and G. Schütz (1995) Steroid hormone receptors: many actors in search of a plot. *Cell* **83**, 851–857.
12. D. J. Mangelsdorf, C. Thummel, M. Beato, P. Herrlich, G. Schütz, et al. (1995) The nuclear receptor superfamily: the second decade. *Cell* **83**, 835–839.
13. H. Gronemeyer and V. Laudet (1995) Transcription factors 3: Nuclear receptors. *Protein Profile* **2**, 1173–1308.
14. The Nuclear Receptor Resource: <http://nrr.georgetown.edu/nrr/nrr.html>

## Steroid Hormones

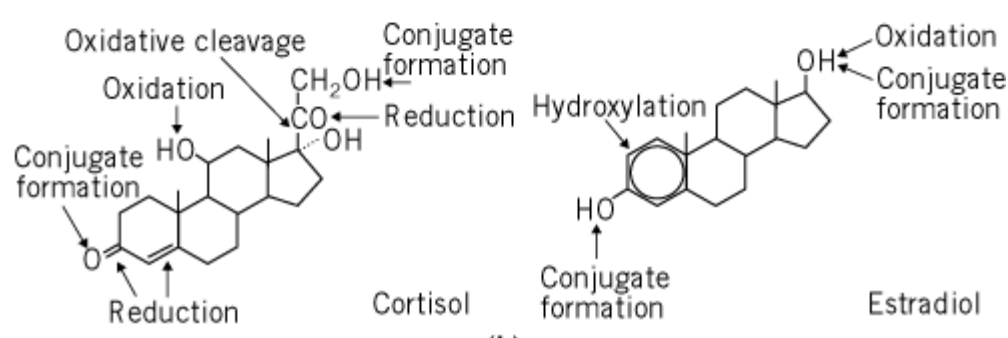
Steroid hormones comprise a major group of [hormones](#) that is characterized by the presence of the cyclopentanperhydrophenanthrene ring. The main steroid hormones and their major biosynthetic pathways are shown in Fig. 1 a. Steroid hormones are derived from blood cholesterol or from cholesterol synthesized in the gland cell. The first step of the biosynthetic pathway is cleavage of the cholesterol side chain; the only steroid hormones that retain the side chain are the ecdysteroids. Further enzymatic steps (eg, cleavages, isomerization, and aromatization) lead to the various steroid hormones.

**Figure 1.** Biosynthesis and inactivation of steroid hormones. ( a ) Biosynthesis. The enzymes involved catalyze the

following reactions: hydroxylations (a, f, g, h, i, k, l, p), dehydrogenations b, d, m), isomerization (c), hydrogenation (o), cleavages (a, e, n), aromatization (q). To illustrate the sites of action of each type of enzyme, the three steroids cholesterol, progesterone, and androstenedione (an intermediate in testosterone biosynthesis) are shown. (b) Inactivation. (From Ref. 11, with permission.)



(a)



(b)

Steroids released in the blood circulation are bound to [albumin](#). [Glucocorticoids](#) and the sex hormones are also bound specifically to corticosteroid-binding globulin and sex-hormone-binding globulin, respectively. Inactivation of steroid hormones (Fig. 1b) takes place primarily in the liver, but also in other organs and cells.

Steroid hormones act in target cells by regulating gene [transcription](#) (see [Hormone Receptors](#)). They

bind to intracellular **receptors** (see [Estrogen Receptors](#)) in a highly specific manner (1, 2). Interaction of the ligand with the receptor leads to dimerization of the protein and binding to specific nucleotide sequences, the [hormone response elements](#). These are [enhancer](#) elements, placed at various distances from the **promoter** of the hormonally regulated gene. Binding of the receptor dimer to these sequences, and its interaction with other [transcription factors](#) and regulatory proteins, results in modulation of gene transcription. In addition, some steroid hormones also exert rapid, nongenomic effects on intracellular  $\text{Ca}^{2+}$  and other ion concentrations, acting at the level of the cell membrane.

### 1. Glucocorticoids

The main representative [glucocorticoids](#) are cortisol (in man) and corticosterone (in rodents). They are synthesized in the zona fasciculata of the adrenal cortex, and they regulate carbohydrate, lipid, and protein metabolism. In the liver, glucocorticoids stimulate protein synthesis, whereas they stimulate **protein degradation** in muscle, bone, and lymph organs. Some of the [amino acids](#) released by protein degradation are metabolized in the liver to glucose precursors and subsequently to glucose and glycogen (gluconeogenesis). Glucocorticoids induce the biosynthesis of the crucial enzymes that are involved in gluconeogenesis (3).

### 2. Mineralocorticoids

The main representative of this class is aldosterone (11b,21-dihydroxy-3,20-dioxopregn-4-en-18-al), which is synthesized in the zona glomerulosa of the adrenal cortex. Aldosterone regulates  $\text{Na}^+$  reabsorption and  $\text{K}^+$  excretion, by inducing the biosynthesis of the  $\text{Na}^+/\text{K}^+$ -ATPase in the epithelial cells of the kidney involuted tubules. Together with  $\text{Na}^+$ ,  $\text{Cl}^-$  ions are also reabsorbed, as well as water, which then returns to the blood circulation (4).

### 3. Androgens

Androgens are the male sex hormones. The effects of castration were recognized very early, and the amelioration of these effects by transplantation of testis can be regarded as one of the first experiments in the field of Endocrinology. The major representatives of this group of hormones are testosterone (D4-androsten-(7b-ol-3-on) and its D4-reduced derivative, 5a-dihydrotestosterone (DHT). Testosterone is synthesized in the Leydig cells of the testis, whereas 5-DHT is produced within the target cells, by the action of 5a-reductase, an enzyme whose synthesis is induced by testosterone. Steroids with androgen activity are also synthesized in the adrenals (androstenedione and 11-hydroxyandrostenedione) (5). In adrenal tumors, the concentrations of these steroids rise, causing virilism in women.

Testosterone acts during embryogenesis, causing the sexual differentiation of the male (differentiation of the Wolffian ducts to epididymis, sperm duct, and sperm cyst), whereas 5-DHT, which is formed in a later period, is responsible for the development of the prostate, the external genitals (penis and scrotum), and growth of male hair (6). In adults, the continuous production of androgens is necessary for sperm maturation and for the activity of the accessory glands of the genital tract. The action of androgens on spermatogenesis is a paracrine effect, as testosterone synthesized in the Leydig cells diffuses directly to the seminiferous tubule and acts there. Androgens also stimulate **erythropoiesis**, affect the immune system, and increase muscle mass. This last effect is due to increased protein synthesis and increased nitrogen retention (“anabolic action”). Androgen derivatives lacking significant androgen action but exerting the anabolic effects (anabolic steroids) are used in therapy, but are also misused by athletes. Steroids with antiandrogen effects have been synthesized, acting through competitive inhibition of the androgen at the level of the androgen receptor (cyproterone).

### 4. Estrogens

The main representative of [estrogens](#) is 17 $\beta$ -estradiol (1,3,5-estratriene-3,17 $\beta$ -diol). It is synthesized in the ovary and in the placenta. Estrogens are responsible for the development of the secondary female sex characteristics (breast, hair, skin, and fat distribution) (7). They have proliferative effects on the endometrial mucosa, acting during the first half of the menstrual cycle. Estrogens induce the biosynthesis of a series of proteins involved in cell proliferation and the [cell cycle](#), such as [growth factors](#), [transcription factors](#), and [cyclins](#).

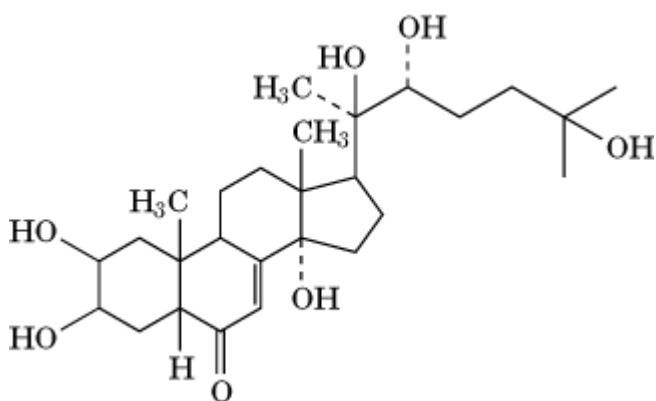
## 5. Progesterone

Progesterone (pregn-4-ene-3,20-dione) is synthesized in the ovary (corpus luteum). It acts during the second half of the menstrual cycle, inducing the secretory phase of the endometrial mucosa and thereby preparing the endometrium for possible embryo implantation (8). Progesterone acts antagonistically to estrogens with regard to cell proliferation, inhibiting transcription of genes encoding growth factors. Furthermore, progesterone induces the biosynthesis of a series of uterus-specific proteins. One important effect of progesterone is the triggering of the acrosome reaction in human sperm cells, which is accomplished by a nongenomic, direct membrane effect of the hormone, activating a signal transduction pathway.

## 6. Ecdysteroids

[Ecdysone](#)((22R)-20-hydroxy-20-(2-hydroxy-2-methylpropyl)cholesta-7,24-dien-3-one) and its 20-hydroxylated derivative, ecdysterone, are the major members of this class and the steroid hormones of insects and crustaceans (Fig. 2). Ecdysone is produced in the prothoracic glands, transported by the hemolymph to target cells, and there it is hydroxylated to the hormonally active ecdysterone. One other hormonally active ecdysone derivative is ponasterone (25-deoxy-ecdysterone). A major function of ecdysteroids is the regulation of ecdysis and metamorphosis. In addition, ecdysteroids affect many other processes, namely, embryogenesis, ovarian maturation, vitellogenesis, ovulation, spermiogenesis, and behavior. In these processes, cell proliferation and [apoptosis](#) are involved and are modulated by ecdysteroids. Ecdysteroids act by regulating gene transcription. This was first shown by Clever and Karlson (9) in the midge [Chironomus tentans by injecting ecdysone into larvae and observing rapid puff induction in salivary gland polytene chromosomes](#). On the basis of this observation, the hormone-gene activation hypothesis was formulated; it has subsequently been corroborated in several vertebrate systems and is still acquiring general validity. The receptor for ecdysteroids has been characterized in several insect species, and three isoforms have been detected in *Drosophila melanogaster*. The isoforms bind to DNA as heterodimers with the **orphan receptor** Ultraspiracle (see [Hormone Receptors](#)), a homologue of the vertebrate **retinoid X-receptor** (10).

**Figure 2.** Structure of ecdysterone (20-hydroxy-ecdysone).



## Bibliography

1. D. J. Mangelsdorf et al. (1995) *Cell* **a83**, 835–839.
2. M. Beato, P. Herrlich, and G. Schutz (1995) *Cell* **83**, 851–857.
3. S. J. Pilgis and D. K. Granner (1992) *Annu. Rev. Physiol.* **54**, 885–909.
4. D. Marver (1985) "The Mineralocorticoid Receptor". In *Biochemical Actions of Hormones*, Vol. **VII** (G. Litwack, ed.), Academic Press, New York, pp. 385–431.
5. F. F. G. Rommerts and F. G. Focko (1990) In *Testosterone Action Deficiency Substitution* (E. Nieschlag, H. M. Behre, eds), Springer, Berlin, 1990, pp. 1–22.
6. J. D. Wilson, F. W. George, and J. E. Griffin (1981) *Science* **211**, 1278–1284.
7. S. S. C. Yen (1991) In *Reproductive Endocrinology*, 13th ed. (R. B. Jaffe, eds.), Saunders, Philadelphia, 1991, pp. 273–308.
8. E. Y. Adashi and P. C. K. Leung (1993) *The Ovary*, Raven Press, New York.
9. U. Clever and P. Karlson (1960) *Exp. Cell Res.* **20**, 623–626.
10. C. S. Thummel (1995) *Cell* **83**, 871–877.
11. J. Koolmann and K.-H. Rohm (1996) *Color Atlas of Biochemistry*, Thieme, Stuttgart.

## Suggestions for Further Reading

12. D. N. Orth, W. J. Kovacs, and C. R. DeBold (1992) "The Adrenal Cortex". In *Williams Textbook of Endocrinology*, 8th ed. (J. D. Wilson and D. W. Foster, eds.), Saunders, London.
13. E. E. Baulieu and P. A. Kelly (eds.) (1990) *Hormones: From Molecules to Disease*, Hermann, Paris.
14. M. G. Parker (1993) *Steroid Hormone Action*, IRL Press, Oxford, U.K.

## Sterol Response Element

The sterol response element is a 10-bp [response element](#) that controls the activation of specific **genes** in response to a fall in the levels of cholesterol ([1](#)). The regulation of specific genes in response to the levels of cholesterol is necessary because cholesterol can be obtained by mammalian cells either by uptake of cholesterol-containing **lipoprotein** particles or by *de novo* biosynthesis. Mammalian cells control these two pathways so as to ensure an appropriate supply of cholesterol by feedback **repression** of several key genes involved in cholesterol metabolism.

In response to lowered levels of cholesterol, the supply is increased both by activating the expression of the gene encoding the **low-density lipoprotein** (LDL) receptor, which plays a central role in cellular uptake of cholesterol, and by activating the genes encoding various enzymes involved in the *de novo* biosynthesis of cholesterol, such as HMG-CoA reductase and HMG-CoA synthase. The genes encoding these proteins contain a sterol response element (SRE) that is able to bind two SRE-binding proteins (SREBP-1 and SREBP-2). These two proteins are synthesized as high-molecular-weight molecules that are bound to cellular [membranes](#); hence, they cannot enter the nucleus and activate transcription by binding to the SRE. Following a decrease in cholesterol levels, SREBP-1 and SREBP-2 are cleaved **proteolytically**, so that a smaller active fragment of each is released from the membrane; it enters the nucleus, where it binds to the SRE and activates transcription ([2](#), [3](#)).

In this manner, these [transcription factors](#) are able to activate specific genes in response to decreased

levels of cholesterol by binding to the SRE and activating the transcription of genes involved in the uptake or synthesis of cholesterol. This simple homeostatic mechanism thus ensures the precise regulation of cholesterol levels within these cells.

### Bibliography

1. J. L. Goldstein and M. S. Brown (1990) *Nature* **343**, 425–430.
2. X. Wang, R. Sato, M. S. Brown, X. Hua, and G. L. Goldstein (1994) *Cell* **77**, 53–62.
3. J. Saki, E. A. Duncan, R. B. Rawson, X. Hua, M. S. Brown, and J. L. Goldstein (1996) *Cell* **85**, 1037–1046.

### Suggestions for Further Reading

4. G. P. Gasic (1994) Basic helix–loop–helix transcription factor and sterol sensor in a single membrane bound molecule. *Cell* **77**, 17–19.
5. H. C. Towle (1995) Metabolic regulation of gene transcription in mammals. *J. Biol. Chem.* **270**, 23235–23238.

## Stokes Radius

In examining the function of biological [macromolecules](#), it is useful to have information about the size and shape of each relevant macromolecule. In the absence of high resolution structural information, like that obtained from [X-ray crystallography](#) or [NMR](#) spectroscopic methods, transport studies can provide low resolution information about these structural parameters and the [hydrodynamic volume](#). Examples of transport techniques include [size exclusion chromatography](#), [sedimentation velocity centrifugation](#), and [electrophoresis](#). The rate of transport of any molecule through a medium using any of these techniques is dictated in part by the frictional resistance of the molecule to motion. This frictional resistance, which is expressed as a **frictional coefficient**, is, in turn, determined in part by the size and shape of the macromolecule. Here the size of a macromolecule will be discussed in terms of the Stokes radius. Use of the Stokes radius to estimate the **molecular weight** will also be presented, as well as the complications arising from the effects of [hydration](#) and molecular shape.

For a spherical molecule, the frictional coefficient is related to its size by

$$f = 6\pi\eta R \quad (1)$$

where  $\eta$  is the viscosity of the fluid through which the molecule is transported and  $R$  is the radius of the sphere ([1-3](#)).  $R$  is referred to as the Stokes radius, and Equation ([1](#)) is known as Stokes' law. This law was originally proposed in 1851 by the British scientist Sir George Stokes from consideration of the forces acting on a particle as it sinks through a liquid column under the influence of gravity ([4](#)). For a given total particle volume, a sphere has the minimum possible frictional coefficient, which is referred to as  $f_0$ .

In using the Stokes radius to obtain information about the size and shape of a biological macromolecule, it is important to consider two factors: (1) biological macromolecules are highly solvated and, (2) many are nonspherical. The solvation factor can be readily taken into account. The radius of a sphere is related to its volume by  $r = (3\text{volume}/4\pi)^{1/3}$ . For a solvated sphere, however,



the volume is expressed as follows:

$$V_h = \frac{M}{N_0(\bar{V}_2 + \delta_1 \bar{V}_1)} \quad (2)$$

where  $M$  is the molecular weight of the unsolvated macromolecule,  $N_0$  is Avogadro's number,  $\bar{V}_2$  is the **partial specific volume** of the macromolecule in  $\text{cm}^3/\text{g}$ ,  $\delta_1$  is the hydration defined as the grams of water per gram of protein, and  $\bar{V}_1$  is the partial specific volume of the aqueous solvent or water in  $\text{cm}^3/\text{g}$ . This last term is simply the inverse of the solvent density. For a more detailed discussion of macromolecular hydration and the measurement or estimates of these terms, refer to the [Hydrodynamic Volume](#) article. Using Equation (2), the radius of a solvated spherical molecule is related to the molecular weight of the molecule by the following equation:

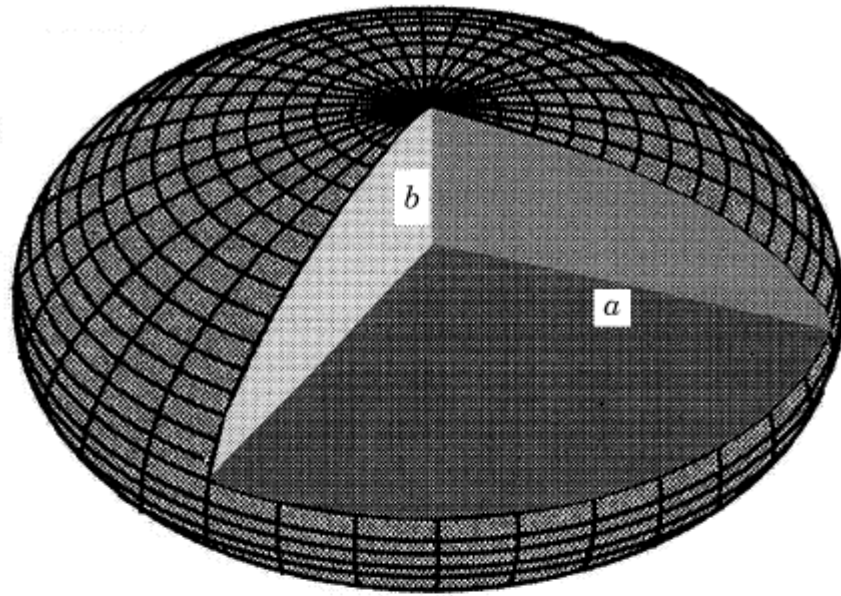
$$R = \left[ \frac{3M_w \bar{V}_2}{4\pi N_0} \left( 1 + \frac{\delta_1 \bar{V}_1}{\bar{V}_2} \right) \right]^{1/3} \quad (3)$$

As indicated above, this equation will apply only to spherical molecules. For those that are nonspherical, use of a Stokes radius determined by, for example, size exclusion chromatography to calculate the molecular weight will lead to an erroneously high value. Nonspherical molecules can, however, be modeled using different assumed shapes, such as ellipsoids or rods, on the basis of information obtained from transport measurements. For such molecules, a term  $R_e$  is defined as the radius of a sphere that is equivalent in volume to that of the nonspherical particle and

$$f_0 = 6\pi\eta R_e \quad (4)$$

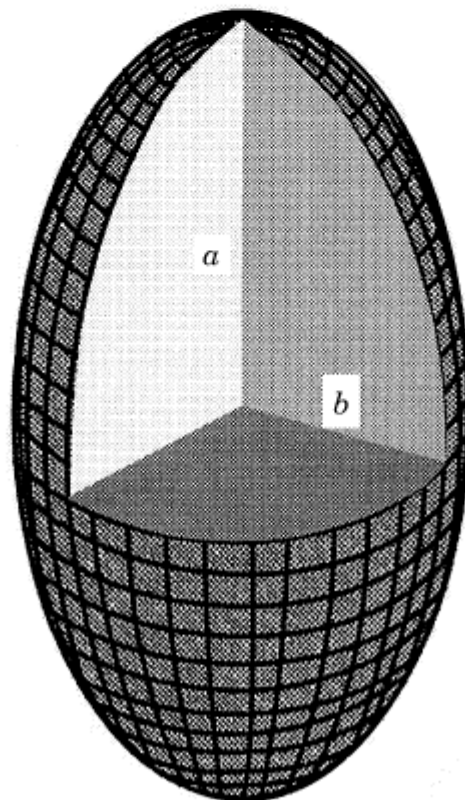
$R_e$  is referred to as the equivalent radius; it can be calculated using Equation (3), provided that the molecular weight of the macromolecule is known. In Equation (4),  $f_0$  is the frictional coefficient that would be measured for a sphere of volume equivalent to that of the nonspherical particle. It is a “minimal” frictional coefficient. The ratio of the true frictional coefficient of the nonspherical particle,  $f$ , to this minimal frictional coefficient of a hypothetical sphere of equivalent volume is useful in obtaining information about the shape of a macromolecule. Mathematical relationships have been derived for prolate and oblate ellipsoids (Fig. 1), as well as for macromolecules that are best described as rods (5). The equations pertaining to ellipsoids are shown in Table 1. From the known frictional coefficient ratio,  $f/f_0$ , one can estimate the shape of the molecule and its asymmetry as expressed by the axial ratio. These relationships are shown graphically in Figure 2.

**Figure 1.** Diagrams of perfect ellipsoidal shapes: (a). Oblate ellipsoids; (b) Prolate ellipsoids. A prolate ellipsoid can be considered to be an elongated sphere, and an oblate ellipsoid is generated by flattening a sphere. For both shapes,  $a$  is the semimajor axis length and  $b$  is the semiminor axis length. The frictional ratios,  $f/f_0$ , for each shape were calculated using the equations shown in Table 1.



$$a/b = 3.0 \quad f/f_0 = 1.105$$

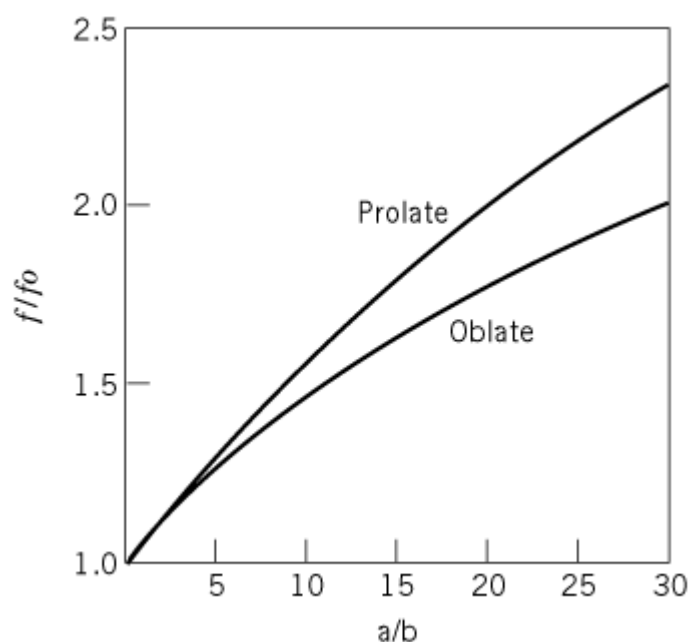
(a)



$$a/b = 2.0 \quad f/f_0 = 1.044$$

(b)

**Figure 2.** Graphical representations of the relationship between the fictional ratio ( $f/f_0$ ) and the axial ratio ( $a/b$ ) for oblate and prolate ellipsoids. The equations relating these two quantities are shown in Table [1](#).



**Table 1. Relationships between Frictional and Axial Ratios for Ellipsoidal Shapes**

| Shape             | $f/f_0^a$  | $R_e^b$        |
|-------------------|--|----------------|
| Prolate ellipsoid | $\frac{(\frac{a}{b})^{-1/3}((\frac{a}{b})^2-1)^{1/2}}{\ln[(\frac{a}{b})+(\frac{a}{b})^2-1]^{1/2}}$ | $(ab^2)^{1/3}$ |
| Oblate ellipsoid  | $\frac{((\frac{a}{b})^2-1)^{1/2}}{(\frac{a}{b})^{2/3}\tan^{-1}[(\frac{a}{b})^2-1]^{1/2}}$          | $(a^2b)^{1/3}$ |

<sup>a</sup> The axial ratio,  $a/b$ , is the ratio of the semimajor axis length,  $a$ , to the minor axis length,  $b$ . These terms are shown pictorially in Figure 1. Remember that the inverse tangent should be taken for a number expressed in radians.

<sup>b</sup> The parameter  $R_e$  is the equivalent radius or the radius of the sphere that is equivalent in volume to that of the nonspherical particle.

#### Bibliography

1. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, Part II, pp. 557–562. W. H. Freeman and Co., San Francisco.
2. K. E. van Holde, W. C. Johnson, and P. S. Ho (1998) *Principles of Biophysical Chemistry*, Prentice-Hall, Englewood Cliffs, NJ, pp. 192–194.
3. I. Tinoco, Jr., K. Sauer, and J. C. Wang (1995) *Physical Chemistry: Principle and Applications in Biological Sciences*, Prentice-Hall, Englewood Cliffs, NJ, pp. 277–281.
4. “Stokes’ law,” *Britannica Online*, <http://www.eb.com>.
5. F. Perrin (1936) *J. Phys. Radium*, **7**, 1–11.

## Stop Codons

In the standard [genetic code](#), three **codons**, UAG, UAA, and UGA, are assigned as *stop codons* to terminate [translation](#) of the [messenger RNA](#) in **protein biosynthesis**. For historical genetic reasons, they are often referred to as the *amber*, *ochre*, and *opal* codons, respectively. To some extent, however, the termination signals are effectively quadruplet since the identity of the 3' base immediately following substantially influences the efficiency of termination (1). The triplets themselves do not have the same efficiency; UGA is somewhat less efficient than the other two but still functions with at least 97% efficiency (ribosomes that do not terminate readthrough by inserting an amino acid, often Trp, when the stop codon is UGA). In **eubacteria**, the relative frequency of use of the termination triplets is UAA > UGA >> UAG. Healthy mutants, especially of eubacteria and **yeast**, have been isolated with the [anticodon](#) of a [transfer RNA](#) species altered so that it can decode a stop codon. Such “**suppressor**” mutants, especially those for UAG, have been very useful for genetic analyses. Normally, no tRNA corresponds to the termination codons. Instead, protein [termination factors](#) mediate release of the polypeptide chain from the [ribosome](#). In *Escherichia coli*, there are two codon-specific release factors, whereas **eukaryotes** have only one (2). In both cases, there is an additional codon-independent protein factor involved.

In eubacteria, there is not always independent ribosome entry to different coding regions contained on the same mRNA. In several cases, the stop codon for one **cistron** overlaps with, or is very close to, the start codon for the next coding region. Although the polypeptide chain is released at the end of the first coding region, the 30S ribosomal subunit, or perhaps the entire 70S ribosome, need not dissociate from the mRNA but immediately initiates decoding of the next sequence. In the absence of such translational coupling, the ribosome recycling factor functions to generate separate ribosomal subunits for independent initiation elsewhere (3).

Stop codons are essential for release of the nascent polypeptide from ribosomes and, consequently, release of ribosomes from mRNA. On occasion, [transcription](#) to synthesize mRNA stops prematurely, or an mRNA is cleaved internally by a **ribonuclease**, so that a coding region without an in-frame stop codon results. When a ribosome comes to the 3' end of such a terminatorless mRNA, it, along with the incomplete nascent peptide, gets stuck. Recently, it has been found that eubacteria have a remarkable rescue mechanism involving an RNA, known as tmRNA, which has both tRNA and mRNA-like features. The *E. coli* tmRNA consists of 363 nucleotides. Its 5' and 3' ends base-pair to form a partial tRNA-like structure that is aminoacylated. The current model indicated by the experimental observations is that tmRNA enters the A site of ribosomes stalled at the end of terminatorless mRNA, donates its amino acid to the carboxy terminus of the stalled nascent peptide, and then functions as an mRNA. The mRNA portion is internal to the molecule and encodes 10 amino acid residues before the ribosome encounters a stop codon and terminates normally, with recycling of ribosomes (4-6). The carboxy-terminal extension of the nascent polypeptide encoded by tmRNA constitutes a recognition site for a [proteinase](#), which destroys the aberrant protein (5).

### Bibliography

1. E. S. Poole, C. M. Brown, and W. P. Tate (1995) EMBO J. **14**, 151–158.
2. L. Frolova et al (1994) Nature **372**, 701–703.
3. L. Janosi, R. Ricker, and A. Kaji (1996) Biochimie **78**, 959–969.
4. G.-F. Tu, G. E. Reid, J.-G. Zhang, R. L. Moritz, and R. J. Simpson (1995) J. Biol. Chem. **270**, 9322–9326.
5. K. C. Keiler, P. R. H. Waller, and R. T. Sauer (1996) Science **271**, 990–993.

6. A. Muto, M. Sato, T. Tadaki, M. Fukushima, C. Ushida, and H. Himeno (1996) *Biochimie* **78**, 985–991.

### Suggestion for Further Reading

7. Y. Nakamura, K. Ito, and L. Isaksson (1996) Emerging understanding of translation termination. *Cell* **87**, 147–150.

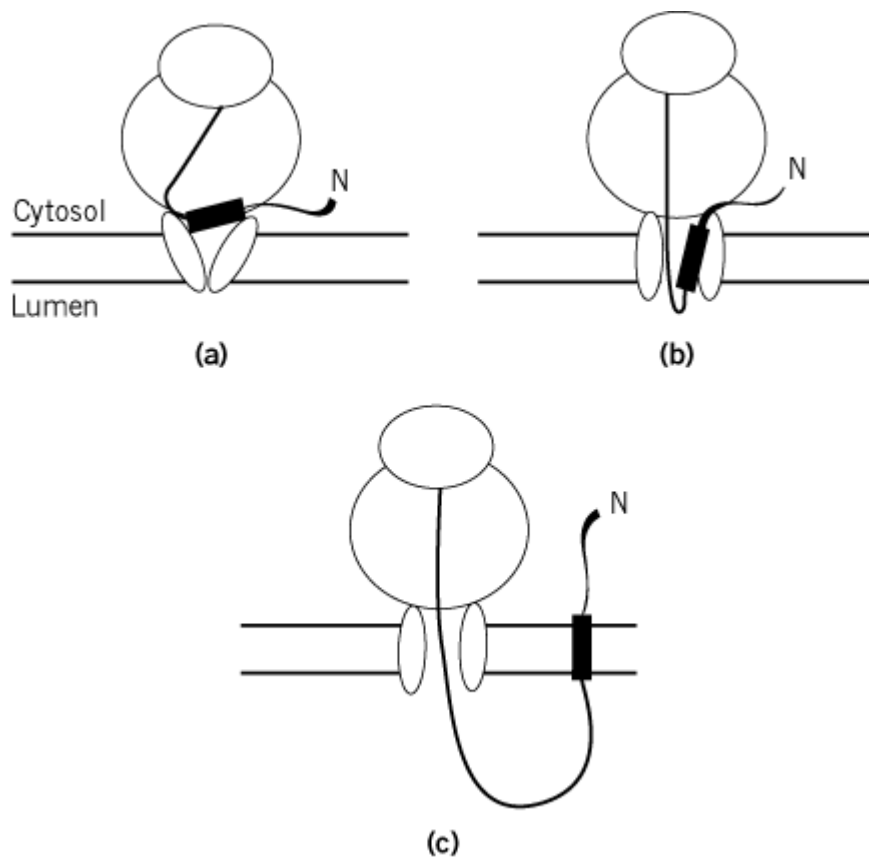
## Stop-Transfer Sequence

During the biosynthesis of an integral [membrane protein](#), the nascent [polypeptide chain](#) passes through a transmembrane channel system known as the translocon. When a peptide containing a stop-transfer sequence [a concept originally suggested by Blobel ([1](#))] inserts into the channel, the translocation process is interrupted. Subsequently, the translocation channel releases the polypeptide chain carrying the stop-transfer sequence laterally to the lipid bilayer. The stop-transfer sequence thus becomes a **transmembrane** segment ([2](#)).

Especially in higher eukaryotes, the translocation of most integral membrane proteins across the [endoplasmic reticulum](#) (ER) membrane occurs cotranslationally and with the help of metabolic energy via the same route as the translocation of secreted soluble proteins—that is, through the Sec61p-complex ([2](#), [3](#)) (see [Protein Secretion](#)). The latter is a heterotrimeric transmembrane protein complex that forms an aqueous channel across the ER membrane and interacts with a number of other components of the protein translocation machinery. In yeast and bacteria, most membrane proteins seem to be inserted to the membrane posttranslationally rather than cotranslationally. Despite this, the majority of them use the Sec-complex for integration into the membrane. However, a Sec-independent route may also exist.

The translocation process is initiated when one of the **Sec proteins** recognizes the complex formed by the **signal sequence**, the [signal recognition particle](#) (SRP), and the SRP receptor. This causes the [ribosome](#) from which the signal sequence is emerging to bind to the translocation channel, sealing its entrance to the cytoplasm (Fig. [1](#)). It is thought that the signal sequence then interacts with the channel, which opens, allowing the nascent polypeptide to insert into it. The polypeptide is translocated across the membrane until the process is halted by a stop-transfer sequence, which interacts with the channel. Although the *N*-terminus of the stop-transfer sequence enters the channel first, the orientation of the sequence relative to the membrane can change during its stay in the channel. Finally, the translocation channel releases the transmembrane segment with the desired orientation laterally to the lipid bilayer. The ribosome remains bound to the Sec61p complex, and the synthesis and translocation of the polypeptide proceeds.

**Figure 1.** Cotranslational integration of a membrane protein into the ER membrane. **(a)** A ribosome that is synthesizing a single-span membrane protein binds to the translocation channel. The stop-transfer sequence (which in this case is also an *N*-terminal signal-anchor sequence) interacts with the channel and opens it. **(b)** The stop-transfer sequence adopts the orientation in which the *N*-terminus is in the cytosol and the *C*-terminus is in the ER lumen. **(c)** The stop-transfer sequence has been released to the membrane, and the synthesis and translocation of the luminal domain of the protein continue. For details and more complex situations, see Ref. [2](#).



During synthesis of polytopic membrane proteins, with multiple transmembrane segments, the translocation channel may store more than one transmembrane segment at a time until these reach the correct topology in the bilayer and are released simultaneously. Such a cooperativity could allow utilization of the favorable **free energy** of membrane insertion of one transmembrane helix to be coupled to the insertion of another helix with less favorable insertion energetics, perhaps due to the presence of charged residues.

The general features of a stop-transfer sequence include [hydrophobicity](#) and positively charged residues in the nontranslocated flanking sequence (the “positive-inside rule”; see [Hydropathy](#)). A stop-transfer sequence typically comprises about 20 hydrophobic residues, as expected for a transmembrane  $\alpha$ -helix. However, a systematic analysis of the requirements for the stop-transfer function has shown that a stretch of only nine [leucine](#) residues is enough for stable insertion into the membrane. In contrast, more than 19 [alanine](#) residues are required for the same function.

A concept related to the stop-transfer sequence is the stop-transfer effector. In the [prion](#) protein, it is the [hydrophilic](#) sequence that precedes a sequence with a conditional stop-transfer function (4). Mutations in both stop-transfer sequence and effector lead to a change in the transmembrane orientation in a population of the prion protein. It appears that the Gerstmann–Straussler–Scheinker-type prion disease, caused by an Ala117Val mutation in the stop-transfer sequence, is linked to the generation of a misoriented transmembrane form of the protein. In general, not many examples of membrane proteins with more than one topology are known. However, ductin and the proteolipid of a V-type ATPase appear represent a case in which the same polypeptide is inserted in the membrane in two different topologies (5).

#### Bibliography

1. G. Blobel (1980) Intracellular protein topogenesis. Proc. Natl. Acad. Sci. USA **77**, 1496–1500.
2. K. E. S. Matlack, W. Mothes, and T. A. Rapoport (1998) Protein translocation: tunnel vision.

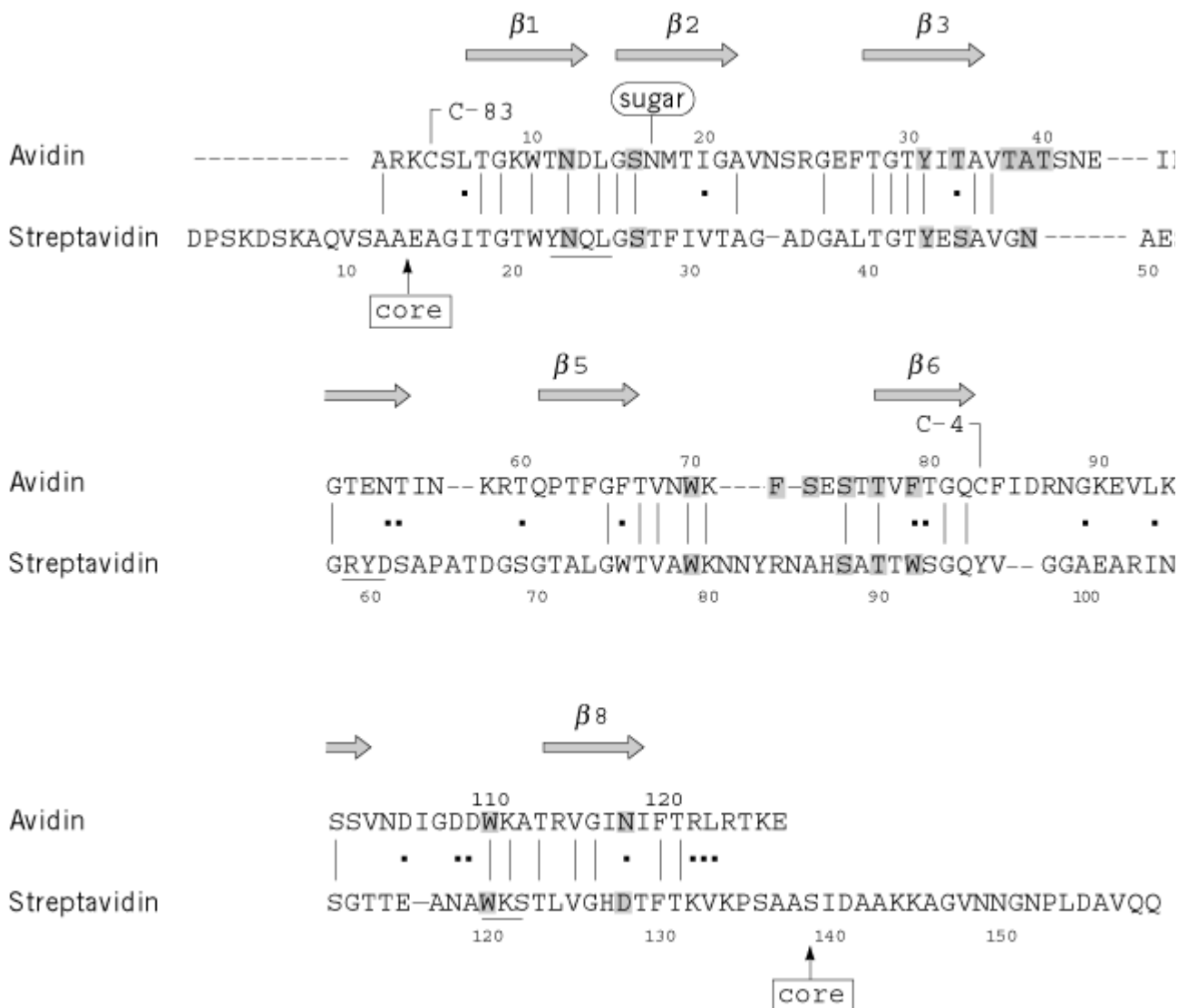
Cell **92**, 381–390.

3. G. Von Heijne (1997) Getting greasy: how transmembrane polypeptide segments integrate into the lipid bilayer. *Mol. Microbiol.* **24**, 249–253.
4. R. S. Hedge, J. A. Mastrianni, M. R. Scott, K. A. Defea, P. Tremblay, M. Torchia, S. J. Dearmond, S. B. Prusiner, and V. R. Lingappa (1998) A transmembrane form of the prion protein in neurodegenerative disease. *Science* **279**, 827–834.
5. J. Dunlop, P. C. Jones, and M. E. Finbow (1995) Membrane insertion and assembly of ductin: a polytopic channel with dual orientations. *EMBO J.* **14**, 3609–3616.

## Streptavidin

Streptavidin, a bacterial protein produced by various strains of the genus *Streptomyces* ([1](#), [2](#)), is a distant relative of egg-white [avidin](#). Similar to avidin, the streptavidin monomer associates into a tetramer and binds the vitamin [biotin](#) with a similarly high **association constant** [ $K_a$  (streptavidin) =  $2.5 \times 10^{13} \text{ M}^{-1}$  versus  $K_a$ (avidin) =  $1.7 \times 10^{15} \text{ M}^{-1}$ ]. Twelve of the 16 biotin-binding residues that make up the binding site of avidin are conserved in streptavidin (Fig. [1](#)). Unlike avidin, which is a positively charged **glycoprotein**, streptavidin is unglycosylated in the native state and bears an overall charge close to neutral.

**Figure 1.** Structure-based sequence alignment of avidin and streptavidin. Identical residues are indicated by vertical line similar residues by dots. The positions of the  $\beta$ -strands are shown by arrows. Residues that participate in binding biotin a The glycosylation site of avidin is shown, as are the N- and C-terminal cleavage sites of streptavidin that are sensitive to digestion. Cellular recognition motifs on streptavidin that cause erratic targeting are underscored.



Streptavidin is the protein component of a two-component antibiotic complex that acts mainly on **gram-negative** bacteria (3). It binds to a group of low molecular weight peptides (termed “stravidins”), each of which contains an unusual amino acid called amicloenomycin (Acm) (4). Acm is a biotin antimetabolite that inhibits one of the enzymes of the biotin biosynthetic pathway. It is interesting to note that avidin can replace streptavidin as the protein component of the antibiotic, even though such a function has not been described in the egg.

Like avidin, streptavidin binds to the dye 4'-hydroxyazobenzene-2-carboxylic acid (HABA) but with a significantly lower association constant. Streptavidin also binds peptides that contain the consensus sequence His-Pro-Gln (HPQ), which has been termed *Strep-tag* when used for affinity-based separations. Egg-white avidin fails to recognize this peptide but binds to a different consensus sequence, His-Pro-Tyr-Pro (HPYP), with a much lower affinity.

Streptavidin gained prominence due to its suitability for use in avidin-biotin technology (see [Avidin-Biotin System](#)). Because of its nearly neutral charge and lack of oligosaccharide moiety, streptavidin is often preferred over egg-white avidin in many applications. Nevertheless, a certain level of familiarity with its molecular properties is advised, because, like avidin, the molecule contains a variety of undesirable characteristics that may cause “nonspecific” or “unwanted” binding with extraneous substances in the experimental system under study.



Streptavidin is secreted as a tetramer in the native state, has a molecular weight of about 66,500 and a subunit molecular weight of about 16,600. The protein undergoes postsecretory **proteolysis** to a molecular weight of 52,800 (13,200 for the subunit). The cleavage site occurs near a Ser-Ala-Ala sequence that is located at both the *N*- and *C*-termini of the protein (5). This truncated protein was termed “core” streptavidin (6) and is the commercially available form generally used in the various applications of the avidin-biotin system.

Core streptavidin has been crystallized (7, 8), and the structure of the binding site is very similar to that of avidin (9), which suggests that nature likes to conserve the prime binding sites. In contrast, despite the identical fold (see [Avidin](#)) and remarkable superimposition of the two proteins, there is only about 30% similarity in the residues that are not part of the binding site (Fig. 1). The dissimilarity is strongest in the loops that interconnect the  $\beta$ -strands, and this dissimilarity is so strong that antibodies raised against one protein fail to interact with the other.

Streptavidin contains several intracellular and extracellular motifs that avidin lacks. This includes the Arg-Tyr-Asp (RYD) sequence, which mimics the Arg-Gly-Asp (RGD) motif—the universal cell recognition domain present in a series of **cell-adhesion**-related proteins (10). In fact, streptavidin interacts specifically with cell-surface [integrins](#) in an RGD-dependent manner, similar to that of [fibronectin](#), vitronectin, and other adhesion molecules (11).

Streptavidin also contains a Trp-Lys-Ser (WKS) sequence. This motif is very rare among known proteins and it was postulated that it has a significant role in the high-affinity interaction of human tissue factor with its ligand. The protein also contains another consensus sequence, Tyr-Xaa-Xaa-Leu (YXXL), known as the immunoreceptor tyrosine-based inhibitory motif (ITIM), which is present in many cytoplasmic protein **domains** from **B** and [T cells](#).

Finally, it seems that the “alanine-scanning information principle” (12), which is currently popular in the development of second-generation drugs, was used by nature in the case of streptavidin. Streptavidin contains 25 [alanine](#) residues (!) whereas avidin contains only five. Different kinds of amino acids were apparently replaced by alanine. Perhaps streptavidin is a later version of avidin in evolution?

Streptavidin may be a “second-generation” protein, but it is certainly the preferred choice, at least currently, for use in the avidin-biotin system. Nevertheless, because of the different intrinsic binding motifs that characterize the streptavidin molecule, it should be used with care in avidin-biotin technology.

## Bibliography

1. E.O. Stapley, J.M. Mata, I.M. Miller, T.C. Demny, and H.B. Woodruff (1963) *Antimicrob. Agents Chemother.* **3**, 20–27.
2. E.A. Bayer, T. Kulik, R. Adar, and M. Wilchek (1995) *Biochim. Biophys. Acta* **1263**, 60–66.
3. L. Chalet and F.J. Wolf (1964) *Arch. Biochem. Biophys.* **106**, 1–5.
4. K.H. Baggaley, B. Blessington, C.P. Falshaw, W.D. Ollis, L. Chalet, and F.J. Wolf (1969) *Chem. Commun.* **1969**, 101–102.
5. E.A. Bayer, H. Ben-Hur, Y. Hiller, and M. Wilchek (1989) *Biochem. J.* **259**, 369–376.
6. C.E. Argaraña, I.D. Kuntz, S. Birken, R. Axel, and C.R. Cantor (1986) *Nucleic Acids Res.* **14**, 1871–1882.
7. W.A. Hendrickson, A. Pähler, J.L. Smith, Y. Satow, E.A. Merritt, and R.P. Phizackerley (1989) *Proc. Natl. Acad. Sci. USA* **86**, 2190–2194.
8. P.C. Weber, D.H. Ohlendorf, J.J. Wendoloski, and F.R. Salemme (1989) *Science* **243**, 85–88.
9. O. Livnah, E.A. Bayer, M. Wilchek, and J.L. Sussman (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5076–5080.

10. R. Alon, E.A. Bayer, and M. Wilchek (1990) *Biochem. Biophys. Res. Commun.* **170**, 1236–1241.
11. R. Alon, E.A. Bayer, and M. Wilchek (1993) *Eur. J. Cell Biol.* **60**, 1–11.
12. H.B. Lowman and J.A. Wells (1993) *J. Mol. Biol.* **234**, 564–578.

### Suggestions for Further Reading

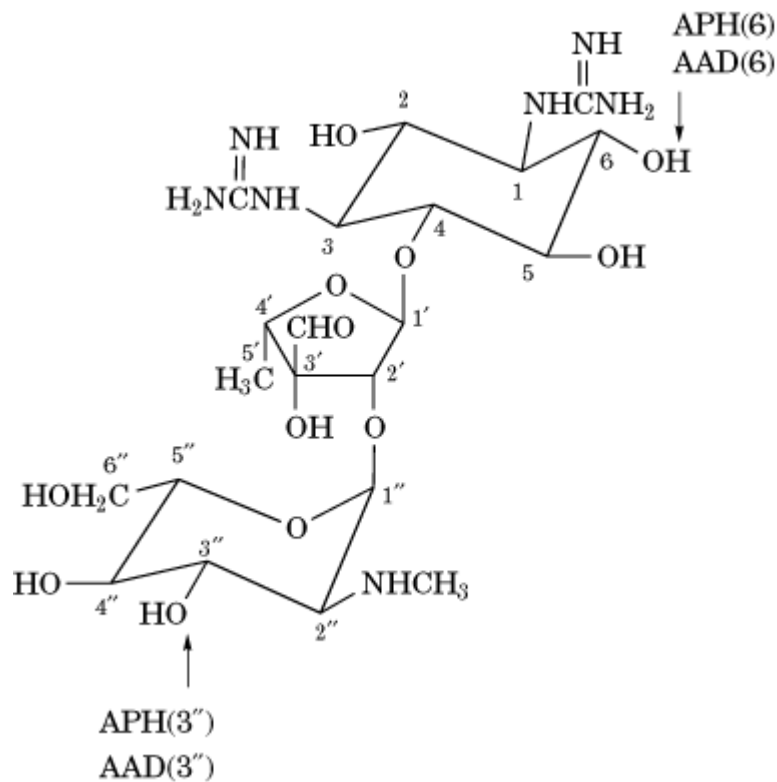
13. E.A. Bayer, and M. Wilchek (1994) "Modified avidins for application in avidin-biotin technology: An improvement on nature", In *Egg Uses and Processing Technologies* (J.S. Sim and S. Nakai, eds.), CAB International, Wallingford, UK, pp. 158–176.
14. N.M. Green (1975) Avidin, *Adv. Protein Chem.* **29**, 85–133.
15. N.M. Green (1990) Avidin and streptavidin, *Methods Enzymol.* **184**, 51–67.
16. M. Wilchek and E.A. Bayer, eds. (1990) *Avidin-Biotin Technology*, *Methods Enzymol.* Vol. 184, Academic Press, San Diego.

## Streptomycin

Streptomycin (SM) was the first aminoglycoside antibiotic agent, introduced in 1944 by Schatz, Bugie, and Waksman (1). It has been used as an effective antimicrobial substance for the treatment of infectious diseases caused by aerobic **Gram-negative** and **Gram-positive** organisms, especially as a first-line drug for the therapy for tuberculosis. However, the appearance of resistant microorganisms (2), SM's ototoxicity, and the introduction of the less toxic broad-spectrum b-lactam antibiotics decreased SM's clinical usefulness in recent years. It has also been used as a tool for analysis of the mechanism of [translation](#) during **protein biosynthesis**, because it binds to the 30 S subunit of [ribosomes](#) and causes misreading of the [messenger RNA](#).

SM was discovered from culture supernatants of *Streptomyces griseus*, although it is now known also to be produced by *S. bikiniensis*, *S. olivaceus*, *S. mashuensis*, *S. galbus*, *S. rameus*, and *S. glovisporus streptomicine*. SM is one of the aminoglycosidic aminocyclitols and belongs to the oligosaccharide group of basic water-soluble antibiotics (Fig. 1). It consists of streptidine, L-streptose, and *N*-methyl-L-glucosamine (C<sub>21</sub>H<sub>39</sub>N<sub>7</sub>O<sub>12</sub>, molecular weight 581.58) (3).

**Figure 1.** The structure of streptomycin. The sites that are modified by inactivating enzymes are indicated. Key: APH, *O*-phosphotransferase; AAD, *O*-adenyltransferase.



## 1. Mechanism of Action

SM exhibits bactericidal activity (killing bacteria) by inhibiting prokaryotic protein biosynthesis, although most such inhibitors have only bacteriostatic activity, ie, inhibiting bacteria's growth but not killing them. Its bactericidal activity can be explained by four processes: (i) SM associates with the cell surface by ionic binding and diffuses through the outer membrane into the periplasmic space, in an energy-independent process (4, 5). This diffusion process may be owing to a “self-promoted uptake” mechanism, in which SM displaces divalent cations that cross-bridge adjacent lipopolysaccharide molecules, permeabilizing the outer membrane (6);

(ii) SM enters the cytoplasm, using the electron transport system to cross the inner membrane (energy-dependent phase I) (6);

(iii) SM binds to the 30 S ribosomal subunit at nucleotides 911–915 of the 16 S ribosomal RNA in *Escherichia coli* (7, 8). The 900 stem and the 530 stem/loop regions of the 16 S rRNA interact with ribosomal proteins S4 and S12 (9-12). Proteins S4, S5, and S12 have been reported to be close neighbors in the protein map of the 30 S ribosomal subunit (13), and these proteins and their associated RNA segments constitute an area of the 30 S ribosomal subunit known as the accuracy region (13). SM binds irreversibly to the 30 S ribosomal subunit and subsequently forms abnormal initiation complexes, so-called streptomycin monosomes, by fixing the 30 S/50 S ribosomal complex at the start codon of the mRNA (14). Accumulation of the abnormal initiation complex blocks further translation of the mRNA, elicits premature termination, and incorporates incorrect amino acids, producing abnormal proteins; and

(iv) The abnormal proteins are incorporated into the cell membrane, and this alters its permeability. The alteration of the cell membrane stimulates the further entry of SM through the energy-dependent phase II process (15) and causes its bactericidal activity (16).

The energy-dependent transport across the inner membrane is inhibited by divalent cations (eg,  $\text{Ca}^{2+}$

and  $Mg^{2+}$ ) and by acidic and anaerobic conditions (6). Therefore, SM shows a reduced level of bactericidal activity in the presence of  $Ca^{2+}$  and  $Mg^{2+}$ , or in acidic conditions, and does not have bactericidal activity against anaerobic organisms and facultative bacteria cultured anaerobically (17). The bactericidal activity of SM competes with other inhibitors of protein bio-synthesis, such as [chloramphenicol](#), [erythromycin](#), and [tetracycline](#), but not with [puromycin](#) (18).

It has been reported that organisms with mutations at the *rpsL* gene, encoding the S12 protein, exhibit SM-resistant or SM-dependent phenotypes. Furthermore, mutations in the genes encoding the S4 or S5 protein—*ram* (ribosomal ambiguity) mutations—of an SM-dependent mutant that cannot grow in the absence of SM cause it to revert to SM independence and to grow in the absence of SM (19).

## 2. Resistance

Two resistance mechanisms to SM have been reported. One is the modification of SM by modifying [enzymes](#), and the other is mutation of the genes encoding the target sites (20, 21).

The modifying enzymes are usually encoded by resistance **plasmids**, but in *Serratia marcescens* and *Providencia stuartii*, they are encoded by chromosomal genes (22). Two modifying enzymes have been elucidated: *O*-phosphotransferase, encoded by the *aph* gene, which phosphorylates at the 3''-OH or 6-OH groups of SM, and *O*-adenyltransferase, encoded by the *aad* gene, which adenylates at the 3''-OH or 6-OH groups (Figure 1). Cross resistance between SM and [kanamycin](#) has been reported. However, because kanamycin, [neomycin](#), paromomycin, and gentamicin are not affected by these modifying enzymes, SM-resistant strains possessing the modifying enzymes are sensitive to these aminoglycosides. Resistance plasmids encoding the modifying enzymes have been found in *Enterobacteriaceae*, *Pseudomonas aeruginosae*, and *Staphylococcus*, but not in mycobacteria.

The other resistance mechanism is the alteration of target sites. This occurs in the S12 protein (the *rpsL* gene) and in the region of the 16 S rRNA (the *rrs* gene) that interacts with the S12 protein. In *E. coli*, substitution mutations on the *rpsL* gene are replacement of Lys43 by arginine, isoleucine, or asparagine and replacement of Lys88 by arginine (23). Similar mutations of the *rpsL* gene are found in mycobacteria (24-30), with Lys43 replacement by arginine or threonine and Lys88 replacement by arginine or glutamine. In addition, Arg9 replacement by histidine and Val93 replacement by methionine have been reported. Resistance mutations in the *rrs* gene have occurred in the regions of the 900 stem and the 530 stem/loop, which interact with the S12 protein as described in the text above: C→T [transition mutations](#) at positions 491, 512, or 516, A→C or T [transversion mutations](#) at position 513, C→A or G transversions at position 903, or A→G transition at position 904. Meyer et al (26) have shown that approximately 30% of resistant *M. tuberculosis* strains have missense mutations at Lys43 or Lys88 of the S12 protein. They have also reported that the resistance level to SM is dependent on the mutation site; mutations at Lys43 or Lys88 of the S12 protein confer high-level resistance, whereas mutations at the 530 stem/loop of 16 S rRNA give intermediate-level resistance. In mycobacteria, approximately 80% of SM-resistant strains carry mutations in either the *rpsL* or the *rrs* gene (28-30), but low-level resistant strains possess no such mutations (25, 26). It has been suggested that alteration of the permeability of the cell wall is a third resistance mechanism, because the addition of the membrane-active [detergent](#) Tween 80 lowered the resistance level (26). Resistance mutations in S12 protein and 16 S rRNA are phenotypically recessive. Therefore, mutation of a single gene cannot generate the resistant **phenotype** in organisms that have multiple copies of the *rrs* gene, such as *E. coli* (with seven copies) and *M. smegmatis* (with two). Because the slowly growing mycobacteria, eg, *M. tuberculosis* and *M. leprae*, have a single rRNA gene copy (31, 32), mutation of the *rrs* gene can confer resistant-phenotype organisms.

SM-dependent mutant strains that require SM for growth possess mutations of the S12 protein: Pro42 altered to leucine, Lys43 altered to glutamic acid, deletion of Lys88, Pro91 altered to leucine or arginine, Gly92 altered to aspartic acid, and deletion of Arg94 (23); plus a mutation of the S4

protein (the *rpsD* gene). Furthermore, a revertant mutation from SM dependence to SM independence has been reported for Gln73 (to proline) in S4 protein (33).

### 3. Clinical Uses

SM is active against mycobacteria and a number of aerobic Gram-negative rods (*Enterobacteriaceae*) and Gram-positive cocci (*Enterococcaceae* and *Staphylococcus*) (2). It is not active, however, against *Streptococcus pyogenes*, *S. pneumoniae*, fungi, and anaerobic bacteria. SM is now used rarely except for the treatment of tuberculosis and some unusual infections, for example, tularemia (*Francisella tularensis*), plague (*Yersinia pestis*), and Weil's disease (*Leptospira interrogans* serovar. *icterohaemorrhagiae*). With the emergence of multiple **drug-resistant** (MDR) strains of *M. tuberculosis*, SM has been renewed as the first-choice drug for tuberculosis. However, SM cannot enter into living cells, so it cannot kill and eradicate intracellular microbes such as mycobacteria.

As SM is a polar cation, it is poorly absorbed from the gastrointestinal tract. It is primarily excreted via the kidneys, essentially unchanged. Therefore, oral administration cannot be expected to be effective against systemic infection.

Considerable intrinsic toxicity, mainly in the form of nephrotoxicity and vestibular or auditory toxicity, is a characteristic of all the aminoglycosides. The toxicities are caused by damage to the eighth cranial nerve and to hair cells in the cochlea. This side effect is irreversible, even after discontinuance of the drug. The nephrotoxic potential varies among the aminoglycosides and is usually reversible when the drug is discontinued.

### Bibliography

1. A. Schatz, E. Bugie and S. A. Waksman (1944) Proc. Soc. Exp. Biol. Med. **55**, 66.
2. G. L. Mandel and W. A. Petri, Jr. (1996) In *Goodman & Gilman's The Pharmacological Basis of Therapeutics* (J. G. Hardman and L. E. Limbird, eds.), McGraw-Hill, New York, pp. 1159–1161.
3. S. Budavari (1996) *The Merck Index: An Encyclopedia of Chemicals, Drugs, and Biologicals*, Merck Research Laboratories, Whitehouse Station, N.J., p. 1507.
4. R. Nakae and T. Nakae (1982) Antimicrob. Agents Chemother. **22**, 554–559.
5. R. E. Hancock, S. W. Farmer, Z. S. Li and K. Poole (1991) Antimicrob. Agents Chemother. **35**, 1309–1314.
6. L. E. Bryan and S. Kwan (1983) Antimicrob. Agents Chemother. **23**, 835–845.
7. D. Moazed and H. F. Noller (1986) Cell **47**, 985–994.
8. G. Schreiner and K. H. Nierhaus (1973) J. Mol. Biol. **81**, 71–82.
9. O. Pongs and V. A. Erdmann (1973) FEBS Lett. **37**, 47–50.
10. S. Stern, T. Powers, L.-M. Changchien and H. F. Noller (1988) J. Mol. Biol. **201**, 683–695.
11. P. N. Allen and H. F. Noller (1989) J. Mol. Biol. **208**, 457–468.
12. T. Powers and H. F. Noller (1994) J. Mol. Biol. **235**, 156–172.
13. M. S. Capel, D. M. Engelman, B. R. Freeborn, M. Kjeldgaard, J. A. Langer, V. Ramakrishnan, D. G. Schindler, D. K. Schneider, B. P. Schoenborn, I. Y. Sillers, S. Yabuki and P. B. Moore (1987) Science **238**, 1403–1406.
14. L. Luzzatto, D. Apirion and D. Schlessinger (1969) J. Mol. Biol. **42**, 315–335.
15. H. J. Busse, C. Wostmann and E. P. Bakker (1992) J. Gen. Microb. **138**, 551–561.
16. B. D. Davis (1988) J. Antimicrob. Chem. **22**, 1–3.
17. L. E. Bryan and S. Kwan (1981) J. Antimicrob. Chemother. **8** (suppl.D), 1–8.
18. H. Yamaki and N. Tanaka (1963) J. Antibiot. **16**, 222.

19. H. G. Wittmann and D. Apirion (1975) *Mol. Gen. Genet.* **141**, 331–341.
20. B. G. Spratt (1994) *Science* **264**, 388–393.
21. J. Davies (1994) *Science* **264**, 375–382.
22. P. N. Rather, E. Orosz, K. J. Shaw, R. Hare and G. Miller (1993) *J. Bacteriol.* **175**, 6492–6498.
23. A. R. Timms, H. Steingrimsdottir, A. R. Lehmann and B. A. Bridges (1992) *Mol. Gen. Genet.* **232**, 89–96.
24. M. Finken, P. Kirschner, A. Meier, A. Wrede and E. C. Bottger (1993) *Mol. Microbiol.* **9**, 1239–1246.
25. J. Nair, D. A. Rouse, G.-H. Bai and S. L. Morris (1993) *Mol. Microbiol.* **10**, 521–527.
26. A. Meier, P. Sander, K.-J. Schaper, M. Scholz and E. C. Bottger (1996) *Antimicrob. Agents Chemother.* **40**, 2452–2454.
27. J. M. Musser (1995) *Clin. Microbiol. Rev.* **8**, 496–514.
28. N. Honore and S. T. Cole (1994) *Antimicrob. Agents Chemother.* **38**, 238–242.
29. C. Katsukawa, A. Tamura, Y. Miyata, C. Abe, M. Makino and Y. Suzuki (1997) *J. Appl. Microbiol.* **83**, 634–640.
30. S. Morris, G. H. Bai, P. Suffys, L. Portillo-Gomez, M. Fairchok and D. Rouse (1995) *J. Infect. Dis.* **171**, 954–960.
31. H. Bercovier, O. Kafri and S. Sela (1986) *Biochem. Biophys. Res. Commun.* **136**, 1136–1141.
32. Y. Suzuki, K. Yoshinaga, Y. Ono, A. Nagata and T. Yamada (1987) *J. Bacteriol.* **169**, 839–843.
33. B. L. Randolph-Anderson, J. E. Boynton, N. W. Gillham, C. Huang and X. Q. Liu (1995) *Mol. Gen. Genet.* **247**, 295–305.

## Stringency

The results of **hybridization** reactions of nucleic acids depend on the solution conditions under which they are performed. Some conditions will favor hybridization of the probe nucleic acid to the target nucleic acid. Low stringency conditions favor duplex formation. High stringency conditions destabilize the duplex, thereby selecting only the most stable duplexes. Under high stringency solution conditions, probes will bind only to regions of high sequence [homology](#). Under low stringency conditions, the requirement for sequence specificity is relaxed, and significant numbers of mismatches, bulges, or hairpins may be tolerated. Regions of the target nucleic acid can be discriminated as a function of sequence homology by adjusting the stringency of the hybridization reaction.

Typically temperature, salt concentration, and chemical **denaturant** concentration are adjusted to optimize the stability of probe–target nucleic acid duplexes. Duplex stability decreases with decreasing salt concentration, with increasing temperature, and with increasing denaturant concentration. Selection of the appropriate level of stringency requires consideration of the influence of the various adjustable conditions on duplex stability.

Several factors, in addition to the base composition, sequence, and length, affect the stability of nucleic acid complexes. These factors include environmental variables such as temperature, the nature and concentration of salts, concentration of added denaturant, and in some cases pH. The concentration of nucleic acid strands is a significant determinant of thermal stability when oligonucleotides are involved in the equilibrium.

The influence of cations on nucleic acid complex stability may be due to specific interactions or to nonspecific binding. Nonspecific, or “territorial”, binding is characterized by an elevated cation concentration over that in the bulk solution in the vicinity of the nucleic acid. This type of binding, which is also called “counterion condensation,” occurs when the linear charge density of a polymer exceeds a critical value.

Both single-stranded and higher order complexes of nucleic acids possess condensed counterions. The linear charge density of the isolated strands is less than  $1/n$  (where  $n$  is the molecularity) of the linear charge density of the complex, due to elongation of the single strands. Therefore, the complex binds more cations than do the isolated strands. When the strands separate, the excess cations are released into solution. As the bulk cation concentration increases, this cation release is increasingly unfavorable. Thus, more energy is required to effect the transition, and the complex is stabilized. This effect depends on the counterion valence, with higher valence counterions promoting greater stabilization than monovalent ions at the same concentration.

Because of the large negative charge of nucleic acid duplexes, anions seldom bind strongly to these molecules. The influence of anions is exerted via effects on cations and on [water](#) structure. Anion effects are usually small; at very high concentration, however, chaotropic anions can reduce the thermal stability of DNA.

Chemical denaturants are used in hybridization reactions to reduce the  $T_m$  into an experimentally convenient range. This is particularly important for experiments using RNA, which is degraded readily at high temperature. Formamide is used frequently as a chemical denaturant because it does not react with the DNA to make or break chemical bonds, it reduces the  $T_m$  of duplexes significantly at reasonable concentrations, it is readily miscible in water, and it has desirable optical properties (low absorbance in the UV region relative to alternatives).

Because neither the nucleic bases nor the sugar–phosphate backbone have titratable groups near neutral pH, the stabilities of nucleic acid complexes are typically insensitive to pH over a wide range. The primary exception is the pyrimidine–purine–pyrimidine **triple helix**. Formation of this complex is accompanied by protonation of the third-strand cytosine residues. This protonation equilibrium is coupled to triplex formation and therefore influences the stability of the complex; it can increase the apparent  $pK_a$  of the cytosine by more than 2 pH units.

Optimization of the hybridization conditions requires consideration of the thermal stability and the rate of formation of the hybrid duplex. The rate depends on temperature and usually reaches a maximum at about  $T_m - 25^\circ\text{C}$  ([1](#)). Selection of the solution conditions typically is made with the aid of empirical equations that predict the  $T_m$  of the fully formed duplexes (p. 31 of Ref. [1](#)). The estimates of  $T_m$  are rather crude, but high precision is usually not required. Different equations are used for DNA–DNA, RNA–RNA, and DNA–RNA duplexes. These equations contain terms for sodium ion concentration, %GC, length ( $L$ ), and % formamide, but are valid only for duplexes of length  $L \geq 50$  base pairs. DNA–DNA duplexes:

$$T_m(^{\circ}\text{C}) \approx 81.5 + 16.6 \log\left(\frac{[\text{Na}^+]}{1 + 0.7[\text{Na}^+]}\right) \\ + 0.41(\%GC) - 500/L - 0.63(\%\text{formamide})$$

RNA–RNA duplexes:

$$T_m(^{\circ}\text{C}) \approx 78 + 16.6 \log\left(\frac{[\text{Na}^+]}{1 + 0.7[\text{Na}^+]}\right) \\ + 0.7(\%GC) - 500/L - 0.35(\%\text{formamide})$$

DNA–RNA duplexes:

$$T_m(^{\circ}\text{C}) \approx 67 + 16.6 \log([\text{Na}^+]/1 + 0.7[\text{Na}^+]) + 0.8(\%GC) \\ - 500/L - 1.0(\% \text{formamide} : 0-20\%) \\ - 0.3(\% \text{additional formamide})$$

If additional components are included in the solutions, the equations may fail to provide useful estimates. A pH value near neutrality is normally used, so a pH term is not included in the empirical equations, but significant deviation from neutral pH will compromise use of these empirical equations. Examination of the equations shows that the %GC and  $L$ -dependent terms are determined by the construction of the molecule. The  $T_m$  can be adjusted by manipulating the concentrations of  $\text{Na}^+$  and formamide. A given  $T_m$  can be obtained by infinitely many different combinations of  $\text{Na}^+$  and formamide concentration. This gives the experimenter flexibility in the design of hybridization conditions.

The  $T_m$  of nucleic acid polymers depends linearly on the GC content, a property that is exploited in the empirical equations shown above. This linearity is due to the sequence effects averaging out over a polymer. For oligonucleotides, the sequence effects become important in determining the thermal stability. Therefore, explicit account of sequence must be made when estimating the  $T_m$  of oligonucleotides. Another feature that differs from the polymer case is the dependence of  $T_m$  on  $[\text{Na}^+]$ , which depends on the length for oligonucleotides. Reliable estimates for the enthalpy,  $DH^{\circ}$ , and free-energy,  $DG^{\circ}$ , changes associated with formation of DNA (2) and RNA (3) duplexes and for  $T_m$  can be computed with knowledge of the base sequence and solution conditions. Thermodynamic data for RNA–DNA duplexes and for higher molecular structures, such as **triple helices**, are sparse relative to the databases for DNA–DNA and RNA–RNA duplexes. Data for these systems are being accumulated, and reliable predictive models should be forthcoming.

Adjustments in stringency can be used to discriminate among nucleic acids with different backbones, thereby favoring binding of a desired hybridization probe. Use of high formamide concentration (80%) selects for RNA–DNA hybrids over DNA–DNA duplexes, so it is useful when RNA probes are employed. The physiochemical properties of the various backbone-modified nonnatural nucleic acid analogues can be exploited by adjustment of the conditions of the binding reaction. Discrimination also can be made between nucleic acid complexes of different numbers of strands. When specific base composition requirements are met, three-stranded nucleic acid complexes can form. Cytosine protonation is required for the formation of the pyrimidine–purine–pyrimidine triple helix, thereby making its formation strongly dependent on pH. Triple helices also exhibit different salt concentration dependencies than do duplexes. Thus, solution conditions can be adjusted to favor or disfavor formation of three-stranded complexes relative to duplexes.

#### Bibliography

1. P. Tijssen (1993) "Hybridization with Nucleic Acid Probes, Part 1", *Theory and Nucleic Acid Preparation*, Elsevier, Amsterdam, p. 44.
2. K. J. Breslauer, R. Frank, H. Blöcker, and L. A. Marky (1986) *Proc. Natl. Acad. Sci. USA* **83**, 3746–3750.
3. S. M. Freier, R. Kierzek, J. A. Jaeger, N. Sugimoto, M. H. Caruthers, T. Neilson, and D. H. Turner (1986) *Proc. Natl. Acad. Sci. USA* **83**, 9373–9377.

#### Suggestion for Further Reading

4. V. A. Bloomfield, D. M. Crothers, and I. Tinoco (1974) *Physical Chemistry of Nucleic Acids*,



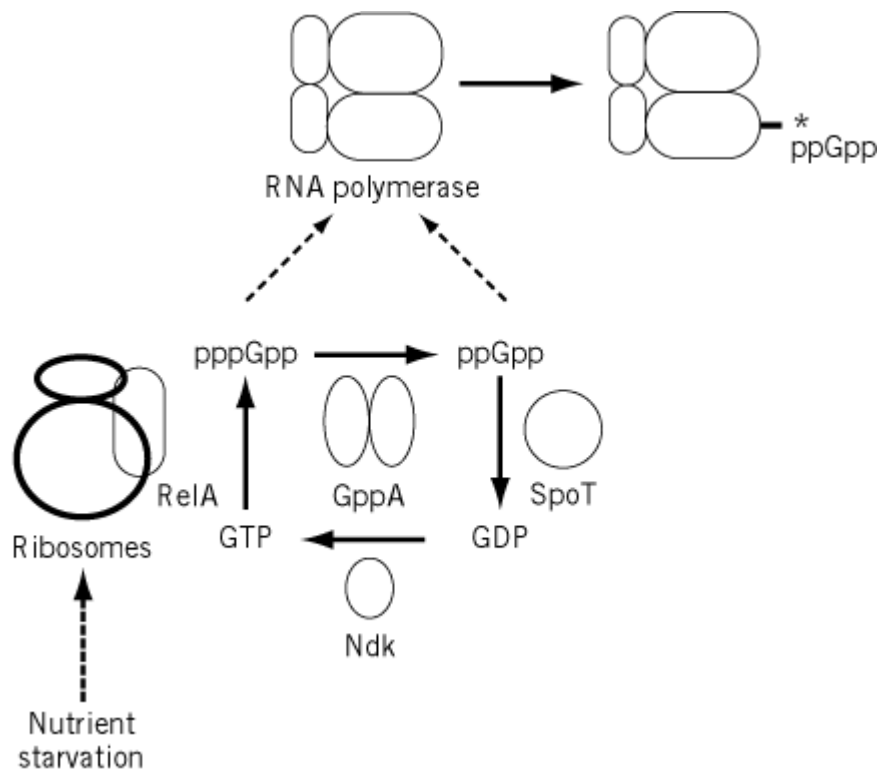
## Stringent Control

On [amino acid](#) starvation, bacterial cells exhibit an adaptive response, known as the stringent response, for repression of production of the machinery for cell growth, such as the [translation](#) and [transcription](#) apparatus ([1](#), [2](#)). The requirement for amino acids is, however, relaxed by [mutations](#) in a single RNA control locus called *relA* ([3](#)). The search for a regulatory molecule under the control of *relA* led to the discovery of accumulation under starved conditions of unusual guanine **nucleotides**, called magic spots I and II ([4](#)), which were later identified as ppGpp and pppGpp, respectively ([5](#), [6](#)).

### 1. Synthesis and Degradation of ppGpp

The *relA1* gene encodes the RelA [enzyme](#) (p)ppGpp synthetase I, originally called stringent factor, that catalyzes the synthesis of pppGpp and ppGpp ([7-9](#)) [Fig. [1](#)]. This RelA protein is associated with [ribosomes](#), but it can be released by high-salt washes. The (p)ppGpp synthesis reaction requires that the ribosomes have bound [messenger RNA](#) and codon-specified uncharged [transfer RNA](#) bound at the acceptor (A) site ([10](#), [11](#)). Thus, the synthesis of (p)ppGpp takes place at an elongation step in protein biosynthesis without ribosomal movement ([10](#)). The purified RelA enzyme, free of ribosomes, is virtually inactive, but it can be activated by inducing conformational changes with methanol and [detergents](#).

**Figure 1.** Regulatory circuit of the stringent response. Ribosome-associated RelA catalyzes the synthesis of both pppGpp and ppGpp from GTP and GDP, respectively, but, from the high  $K_m$  value, GTP is considered to be the natural substrate. The conversion of pppGpp to ppGpp is catalyzed by several different enzymes but, from its low  $K_m$  value, the *gppG* gene product is considered to play a major role *in vivo*. SpoT catalyzes the breakdown of (p)ppGpp to GTP or GDP; in the presence of high concentrations of GTP or GDP, however, it acts as (p)ppGpp synthetase II. The most probable target of (p)ppGpp action is the  $\beta$  subunit of RNA polymerase.



The synthetic reaction of (p)ppGpp is a pyrophosphoryl group transfer from ATP to the 3' hydroxyl of GDP or GTP acceptor nucleotides (7, 12). The high  $K_m$  (Michaelis constant) values for GTP and GDP (0.5 mM) are within the physiological concentration range for GTP but not for GDP, suggesting that pppGpp is the most likely primary product *in vivo* (7, 9). Conversion of pppGpp to ppGpp can be catalyzed by a variety of enzymes *in vitro*, including two exopolyphosphatases, products of the *gppG* and *ppx* genes (13-15) [Fig. 1].

The *spoT* gene encodes an enzyme for ppGpp degradation, ie, removal of the 3'-pyrophosphate residue from either ppGpp or pppGpp (16) [Fig. 1]. The protein SpoT is also associated with ribosomes but is active even after release by high-salt treatment (17, 18). The *spoT* operon includes five genes in the order: *gmk* (guanylate kinase), *rpoZ* (RNA polymerase w subunit), *spoT* (ppGpp 3'-pyrophosphatase), *spoU* (RNA 2'-O-methyltransferase) and *recG* (junction-specific RNA helicase) (19). At least three gene products are involved in guanine nucleotide metabolism. RNA polymerase-associated w protein was suggested to control the ppGpp sensitivity of RNA polymerase (20), but this concept has been challenged (21). SpoT catalyzes the *relA*-independent synthesis of (p) ppGpp, the reverse reaction of (p)ppGpp degradation and, therefore, SpoT is now recognized as (p) ppGpp synthetase II (22, 23).

## 2. Target of ppGpp Action

The concept that the target of ppGpp action is **RNA polymerase** is proposed on the basis of observations that *in vitro* transcription of stringently controlled genes by purified RNA polymerase is inhibited by ppGpp in the absence of additional factors, even though the step of transcription affected by ppGpp seems to vary, depending on which template is used (24-28). Early *in vitro* studies, however, gave conflicting results regarding the effect of ppGpp on transcription of the stringently controlled genes. This disagreement has now been attributed to various factors, including the heterogeneity in RNA polymerase preparations used in early studies, the involvement of several additional factors for maximum transcription of the stringent **promoters**, the presence of multiple promoters with different specificities in the same and different genes, nonlinear DNA dependency of

the stringent promoters, and inhibition by ppGpp of not only initiation but also transcription elongation.

In parallel with in vitro transcription studies, genetic studies have reached the same conclusion. Some **rifampicin**-resistant *rpoB* mutants of an *Escherichia coli* B *relA* parental strain were found to be hypersensitive to ppGpp levels (29). Among point mutations of the  $\beta$  subunit of RNA polymerase generated by the **suppression** of [amber mutations](#) in the *rpoB* gene, three amino acid substitutions were found to result in a relaxed RNA-control **phenotype** (30). This result has been confirmed by in vitro transcription studies using RNA polymerase purified from these *rpoB* mutants with relaxed phenotype (31).

The outcome of these biochemical and genetic observations supported the simple model that ppGpp is one of the [transcription factors](#) that interacts with RNA polymerase and modulates its promoter selectivity (32-34). Evidence for direct contact between ppGpp and the RNA polymerase  $\beta$  subunit has accumulated [Fig. 1]. Owens et al. (35) demonstrated [crosslinking](#) of azido-ppGp to RNA polymerase. Reddy et al (36) showed stoichiometric binding of a fluorescence-labeled ppGpp to RNA polymerase and estimated its distance from the rifampicin-binding domain of RNA polymerase. Chatterji et al (37) synthesized azido-ppGpp and, using a **radiolabeled** derivative, identified the sites of ppGpp crosslinking on the  $\beta$  subunit of RNA polymerase, which are close to those of *rpoB* mutations conferring relaxed phenotype (31). Spontaneous **missense** suppressor alleles, which are able to confer complete **prototrophy** to *E. coli* mutants lacking both of the (p) ppGpp synthases, RelA and SpoT, were localized in RNA polymerase subunit genes, 97% in *rpoB* and *rpoC* and 3% in *rpoD*. Mutations in conserved region 3 of [sigma factor](#)  $\sigma^{70}$  increases the ppGpp sensitivity of RNA polymerase (38). One possibility is that these mutant  $\sigma^{70}$  induce conformational changes of the  $\beta$  subunit, leading to alteration in its affinity for ppGpp.

### 3. Inhibition of RNA Synthesis by ppGpp

*Escherichia coli* possesses seven ribosomal RNA (rRNA), *rrn* operons, which all contain a *tRNA* gene in the spacer region between the 16S and 23S *rRNA* genes and the 5S *rRNA* gene, in the region distal to the 23S gene. The promoter region of these long transcripts is sufficient to determine their regulation during the stringent response, as well as during growth-rate control. Two promoters have been identified, upstream P1 and downstream P2, which together display complex transcriptional regulation. Upstream of the P1 promoter are regions implicated in transcript activation, ie, the UP element (or prokaryotic [enhancer](#)) and the Fis protein-binding site (39). Activity of the *rrn* promoters *in vitro* is strongly dependent on DNA **supercoiling** in purified assays (40). The C-terminal domain of RNA polymerase alpha-subunit is involved in recognition of both the DNA UP element and the Fis protein (33, 34). From the three-dimensional structure of this domain (41), the same protein surface seems to be involved in contact with both the DNA and protein factors, each leading to modulation of the promoter-recognition properties of RNA polymerase. Direct measurement of P1 and P2 activities *in vivo* indicated that the strong promoter P1 plays a major role in rRNA synthesis in rapidly growing cells and is subject to inhibition during the stringent response, but that P2, a weaker constitutive promoter, is insensitive to amino acid starvation (42, 43). Using an *in vitro* mixed transcription system, Kajitani and Ishihama (25) revealed that the upstream P1 promoter was specifically repressed by ppGpp, whereas the downstream P2 promoter was virtually unaffected.

Several lines of evidence indicate that transcription of stringently controlled genes, including *tRNA* genes and some ribosomal protein genes, by purified RNA polymerase holoenzyme  $E_s^{70}$  is also inhibited in the presence of ppGpp (25, 44). Taken together, these observations indicate that these stringently controlled genes share a common regulatory DNA sequence that is specifically recognized by ppGpp-associated RNA polymerase.

### 4. Stringent Signal

Sequence comparisons suggested the presence between the -10 [Pribnow Box](#) and the +1 RNA transcription start site of a GC-rich region that has been called a “discriminator” and has been proposed as a common feature of all promoters negatively regulated during the stringent response ([32](#), [45](#)). When the GC-rich discriminator of the *tyrT* gene was replaced with an AT-rich sequence, the promoter became more active than the wild-type one. The mutant promoter also became resistant to negative stringent control *in vivo* ([44](#)).

Transcripts of the *tufB* operon contain four *tRNA* genes upstream of the *tufB* gene and display negative stringent control *in vivo* and ppGpp inhibition *in vitro* ([46](#)). Studies of ppGpp inhibition of promoter mutants implicated the critical role of base-pair positions -7 to -4, which normally have the sequence GCGC ([47](#)). Changing each of the GC pairs at the -7 to -4 positions individually to AT pairs led to a reduction in the degree of ppGpp inhibition *in vitro*. The upstream and strong *rrn* promoter P1 is responsible for both stringent control and growth-rate-dependent control, whereas the activity of P2 is weak and is involved in a low level of constitutive expression of the *rrn* genes. At slow growth rates, the P2 activity predominates and is responsible for the persistence of rRNA synthesis, because it is resistant to ppGpp inhibition. The core promoter region of P1, from nucleotides -41 to +1, is sufficient for growth-rate-dependent control ([48](#)), which overlaps the discriminator signal for stringent control.

## 5. Transcription Activation by ppGpp

Venetianer ([49](#), [50](#)) observed that the [his operon mRNA](#) was abundant during the stringent response. This finding suggested that the expression of certain genes, including those of amino acid biosynthetic [operons](#), is activated during the stringent response. It is currently thought that many genes are subject to positive ppGpp control. [Two-Dimensional Gel Electrophoresis](#) analysis of proteins synthesized during the stringent response reveals that the fraction of proteins showing positive control is approximately equal to those showing inhibited synthesis ([51](#), [52](#)).

As with negative control, ppGpp has been implicated as the regulatory signal mediating positive control. The *relA* gene product is required for maximal expression *in vivo* of the *his* operon ([53](#), [54](#)). Again, these effects are promoter-specific, for they occur even when the *his* **attenuator** is deleted ([54](#)). Positive regulation of *his* operon expression has been verified as occurring at the promoter ([55](#), [56](#)). Enhancement of transcription of the *arg* and [trp operons](#) *in vitro* was also found in the presence of ppGpp ([25](#), [57](#)).

## 6. Role of ppGpp in Growth Control

Early studies suggested that, during steady-state growth, both relaxed and stringent strains displayed an inverse correlation between basal ppGpp levels and both growth rate and RNA accumulation levels ([58](#), [59](#)). The relationship between ppGpp levels and rRNA accumulation during very slow steady-state growth is, however, complicated by a number of features: (1) rRNA is synthesized but is not assembled into ribosomes, as a result of unbalanced synthesis of ribosomal proteins and/or rapid degradation of newly synthesized rRNA ([59-61](#)) and (2) ppGpp affects the expression of many genes in different ways, either repression or activation, and to various extents. Thus, the hierarchy of **gene expression** is markedly influenced, affecting directly or indirectly the rate of rRNA synthesis.

Growth-rate control is observed both in the presence and in the absence of ppGpp ([39](#)) suggesting that ppGpp does not play a major role in growth-rate control. Instead, Gaal et al. ([62](#)) provided conclusive evidence that the concentration of substrate nucleoside triphosphates determines the growth-rate-dependent control of rRNA transcription. Thus, the core promoter sequence may specify the concentration of nucleoside triphosphates required for efficient initiation of transcription.

Although ppGpp is not directly involved in growth-rate-dependent control, it is involved in growth-phase control by inducing the synthesis of RNA polymerase  $\sigma^S$  subunit for transcription of stationary

phase-specific or **stress response**–specific genes (63). In the stationary phase of *E. coli* growth, ppGpp also induces the production of polyphosphate (64), which ultimately leads to alteration in the gene expression pattern by binding to the RNA polymerase and modulating its promoter selectivity (65).

## Bibliography

1. J. A. Gallant (1979) Stringent control in *E. coli*. *Annu. Rev. Genet.* **13**, 393–415.
2. M. Cachel et al. (1996) "The stringent response". In *Escherichia coli and Salmonella* (F. C. Neidhardt, ed.), ASM Press, Washington, DC, 2nd ed., pp. 1458–1496.
3. L. Alfoldi, G. S. Stent, and R. C. Clowes (1962) The chromosomal site for the RNA control (RC) locus in *Escherichia coli*. *J. Mol. Biol.* **5**, 348–355.
4. M. Cashel, and J. Gallant (1969) Two compounds implicated in the function of the RC gene of *Escherichia coli*. *Nature (London)* **22**, 838–841.
5. M. Cashel and B. Kalbacher (1970) The control of ribonucleic acid synthesis in *Escherichia coli* V. Characterization of a nucleotide associated with the stringent response. *J. Biol. Chem.* **245**, 2309–2318.
6. L. Que et al. (1973) Guanosine 5 $\epsilon$ -diphosphate, 3 $\epsilon$ -diphosphate: Assignment of structure by <sup>13</sup>C nuclear magnetic resonance spectroscopy. *Proc. Natl. Acad. Sci. USA* **70**, 2563–2566.
7. J. W. Cochran and R. W. Byrne (1974) Isolation and properties of a ribosome-bound factor required for ppGpp and pppGpp synthesis in *Escherichia coli*. *J. Biol. Chem.* **249**, 353–360.
8. W. A. Haseltine et al. (1972) MS1 and MS2 are made on ribosomes in an idling step of protein synthesis. *Nature* **238**, 381–384.
9. F. S. Pedersen and N. O. Kjeldgaard (1977) Analysis of the relA gene product of *Escherichia coli*. *Eur. J. Biochem.* **76**, 91–97.
10. W. A. Haseltine and R. Block (1973) Synthesis of guanosine tetre- and pentaphosphate requires the presence of a codon-specific, uncharged transfer ribonucleic acid in the acceptor site of ribosomes. *Proc. Natl. Acad. Sci. USA* **70**, 1564–1568.
11. F. S. Pedersen, E. Lund, and N. O. Kjeldgaard (1973) Codon specific, tRNA-dependent in vitro synthesis of ppGpp and pppGpp. *Nature* **243**, 12–15.
12. J. Sy, Y. Ogawa, and F. Lipmann (1973) Nonribosomal synthesis of guanosine 5 $\epsilon$ ,3 $\epsilon$ -polyphosphates by the ribosomal wash of stringent *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **70**, 2145–2148.
13. M. Akiyama, E. Crooke, and A. Kornberg (1993) An exopolyphosphatase of *Escherichia coli*. The enzyme and its ppx gene in a polyphosphate operon. *J. Biol. Chem.* **268**, 633–639.
14. E. Hamel and M. Cashel (1973) Role of guanine nucleotides in protein synthesis: Elongation factor G and guanosine 5 $\epsilon$ -triphosphate, 3 $\epsilon$ -diphosphate. *Proc. Natl. Acad. Sci. USA* **70**, 3250–3254.
15. C. R. Somerville and A. Ahmed (1979) Mutants of *Escherichia coli* defective in the degradation of guanosine 5 $\epsilon$ -triphosphate, 3 $\epsilon$ -diphosphate (pppGpp). *Mol. Gen. Genet.* **169**, 315–323.
16. E. A. Heinemeyer and D. Richter (1978) Mechanism of the in vitro breakdown of guanosine 5 $\epsilon$ -diphosphate 3 $\epsilon$ -diphosphate in *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **75**, 4180–4183.
17. E. A. Heinemeyer and D. Richter (1977) *In vitro* degradation of guanosine tetraphosphate (ppGpp) by an enzyme associated with the ribosomal fraction from *Escherichia coli*. *FEBS Lett.* **84**, 357–361.
18. I. Sy (1977) In vitro degradation of guanosine 5 $\epsilon$ -diphosphate, 3 $\epsilon$ -diphosphate. *Proc. Natl. Acad. Sci. USA* **74**, 5529–5533.
19. B. C. Persson, G. Jaeger, and C. Gustafsson (1997) The spoU gene of *Escherichia coli*, the fourth gene of the *spoT* operon, is essential for tRNA (Gm18) 2 $\epsilon$ -O-methyltransferase activity. *Nucleic Acids Res.* **25**, 4093–4097.
20. K. Igarashi, N. Fujita, and A. Ishihama, (1989) Promoter selectivity of *Escherichia coli* RNA

- polymerase: Omega factor is responsible for ppGpp sensitivity. *Nucleic Acids Res.* **17**, 8755–8765.
21. D. Gentry et al. (1991) The omega subunit of *Escherichia coli* K-12 RNA polymerase is not required for stringent RNA control in vivo. *J. Bacteriol.* **173**, 3901–3903.
  22. V. J. Hernandez and H. Bremer *Escherichia coli* ppGpp synthetase II activity requires *spoT*. *J. Biol. Chem.* **266**, 5991–5999.
  23. H. Xiao et al. (1991) Residual guanosine 3 $\epsilon$ ,5 $\epsilon$ -bispyrophosphate synthetic activity of *relA* null mutants can be eliminated by *spoT* null mutations. *J. Biol. Chem.* **266**, 5980–5990.
  24. J. Hamming, A. B. Geert, and M. Gruber (1980) *E. coli* RNA polymerase–rDNA promoter interaction and the effect of ppGpp. *Nucl. Acids Res.* **8**, 3947–3963.
  25. M. Kajitani and A. Ishihama (1984) Promoter selectivity of *Escherichia coli* RNA polymerase. Differential stringent control of the multiple promoters from ribosomal RNA and protein operons. *J. Biol. Chem.* **259**, 1951–1957.
  26. R. E. Kingston and M. J. Chamberlin (1981) Pausing and attenuation of in vitro transcription in the *rrnB* operon of *E. coli*. *Cell* **27**, 523–531.
  27. C. C. Pao, P. P. Dennis, and J. A. Gallant (1980) Regulation of ribosomal and transfer RNA synthesis by guanosine 5 $\epsilon$ -diphosphate-3 $\epsilon$ -monophosphate. *J. Biol. Chem.* **255**, 1830–1833.
  28. A. J. J. Van Ooyen, M. Gruber, and P. Jorgenson (1976) The mechanism of action of ppGpp on rRNA synthesis in vitro. *Cell* **8**, 123–128.
  29. R. Little, J. Ryals, and H. Bremer, (1983) *rpoB* mutation in *Escherichia coli* alters control of ribosome synthesis by guanosine tetraphosphate. *J. Bacteriol.* **154**, 787–792.
  30. V. Nene and R. E. Glass (1983) Relaxed mutants of *Escherichia coli* RNA polymerase. *FEBS Lett.* **153**, 307–310.
  31. R. E. Glass, S. T. Jones, and A. Ishihama (1986) Genetic studies on the subunit of *Escherichia coli* RNA polymerase. VII. RNA polymerase is a target for ppGpp. *Mol. Gen. Genet.* **203**, 265–268.
  32. A. Travers (1980) Modulation of RNA polymerase specificity by ppGpp. *Mol. Gen. Genet.* **147**, 225–232.
  33. A. Ishihama (1993) Protein-protein communication within the transcription apparatus. *J. Bacteriol.* **175**, 2483–2489.
  34. A. Ishihama (1997) "Promoter selectivity control of *Escherichia coli* RNA polymerase". *Nucleic Acids and Molecular Biology*, Vol. 11, *Mechanism of Transcription* (F. Eckstein and D. Lilley, eds.) Springer-Verlag, Heidelberg, Germany, pp. 53–70.
  35. J. R. Owens, A.-Y. M. Woody, and B. E. Haley (1987) Characterization of the guanosine-3 $\epsilon$ -diphosphate-5 $\epsilon$ -diphosphate binding site on *E. coli* RNA polymerase using a photoprobe, 8-azido guanosine-3 $\epsilon$ ,5 $\epsilon$ -bisphosphate. *Biochem. Biophys. Res. Commun.*, **142**, 964–971.
  36. P. S. Reddy, A. Raghavan, and D. Chatterji (1995) Evidence for a ppGpp-binding site on *Escherichia coli* RNA polymerase: Proximity relationship with the rifampicin-binding domain. *Mol. Microbiol.* **15**, 225–256.
  37. D. Chatterji, N. Fujita, and A. Ishihama (1998) The mediator for stringent control, ppGpp, binds to the  $\omega$ -subunit of *Escherichia coli* RNA polymerase. *Genes to Cells* **3**, 279–287.
  38. V. J. Hernandez and M. Cashel (1995) Changes in conserved region 3 of *Escherichia coli* sigma 70 mediate ppGpp-dependent functions in vivo. *J. Mol. Biol.* **252**, 536–549.
  39. R. L. Gourse, H. A. deBoer, and M. Nomura (1986) DNA determinants of rRNA synthesis in *E. coli*: Growth rate–dependent regulation, feedback inhibition, up-stream activation, anti-termination. *Cell* **44**, 197–205.
  40. G. Glaser, P. Sarmientos, and M. Cashel (1983) Functional interrelationship between two tandem *E. coli* ribosomal RNA promoters. *Nature* **302**, 74–76.
  41. Y. H. Jeon et al. (1995) Solution structure of the C-terminal domain of RNA polymerase  $\omega$  subunit responsible for the contact with transacting transcription factors and *cis*-acting UP

- element of promoter. *Science* **270**, 1495–1997.
42. R. L. Gourse et al. (1996) rRNA transcription and growth rate-dependent regulation of ribosome synthesis in *Escherichia coli*. *Annu. Rev. Microbiol.* **50**, 645–677.
  43. P. Sarmientos and M. Cashel (1983) Carbon starvation and growth rate-dependent regulation of the *Escherichia coli* ribosomal RNA promoters: Differential control of dual promoters. *Proc. Natl. Acad. Sci. USA* **80**, 7010–7013.
  44. A. I. Lamond and A. A. Travers (1985) Genetically separable functional elements mediate the optimal expression and stringent regulation of a bacterial tRNA gene. *Cell* **40**, 319–326.
  45. A. Travers (1984) Conserved features of coordinately regulated *E. coli* promoters. *Nucleic Acids Res.* **12**, 2605–2618.
  46. J. Mizushima-Sugano, J. A. Miyajima, and Y. Kaziro (1983) Selective inhibition of transcription of the *E. coli* operon by guanosine-5'-diphosphate-3'-diphosphate. *Mol. Gen. Genet.* **189**, 185–192.
  47. J. Mizushima-Sugano and Y. Kaziro (1985) Regulation of the expression of the *tufB* operon: DNA sequences directly involved in the stringent control. *EMBO J.* **4**, 1053–1058.
  48. M. S. Barlett, and R. L. Gourse (1994) Growth rate-dependent control of the *rrnB* P1 core promoter in *Escherichia coli*. *J. Bacteriol.* **176**, 5560–5564.
  49. P. Venetianer (1968) Preferential synthesis of the messenger RNA of the histidine operon during histidine starvation. *Biochem. Biophys. Res. Commun.* **33**, 959–963.
  50. P. Venetianer (1969) Level of messenger RNA transcribed from the histidine operon in repressed, derepressed, and histidine-starved *Salmonella typhimurium*. *J. Mol. Biol.* **45**, 375–384.
  51. A. V. Furano and E. P. Witel (1976) Effect of the *relA* gene on the synthesis of individual proteins in vivo. *Cell* **8**, 115–122.
  52. P. H. O'Farrell (1978) The suppression of defective translation by ppGpp and its role in the stringent response. *Cell* **14**, 545–557.
  53. S. W. Artz, and J. R. Broach (1975) Histidine regulation in *Salmonella typhimurium*: An activator-attenuator model of gene regulation. *Proc. Natl. Acad. Sci. USA* **72**, 3453–3457.
  54. J. C. Stephens, S. W. Artz, and B. N. Ames (1975) Guanosine 5'-diphosphate 3'-diphosphate (ppGpp): Positive effector for histidine operon transcription and general signal for amino-acid deficiency. *Proc. Natl. Acad. Sci. USA* **72**, 4389–4393.
  55. K. E. Rudd et al. (1985) Mutations in the *spoT* gene of *Salmonella typhimurium*: Effects of *his* operon expression. *J. Bacteriol.* **163**, 534–542.
  56. M. E. Winkler, D. J. Roth, and P. E. Hartman (1978) Promoter- and attenuator-related metabolic regulation of the *Salmonella typhimurium* histidine operon. *J. Bacteriol.* **133**, 993–1000.
  57. M. J. Zidwick, J. Korshus, and P. Rogers (1984) Positive control of expression of the *argECBH* gene cluster in vitro by guanosine 5'-diphosphate 3'-diphosphate. *J. Bacteriol.* **159**, 647–651.
  58. R. A. Lazzarini, M. Cashel, and J. Gallant (1971) On the regulation of guanosine tetraphosphate levels in stringent and relaxed strains of *Escherichia coli*. *J. Biol. Chem.* **246**, 4381–4385.
  59. J. Ryals, R. Little, and H. Bremer (1982) Control of rRNA and tRNA synthesis in *Escherichia coli* by guanosine tetraphosphate. *J. Bacteriol.* **151**, 1261–1268.
  60. K. Gausing (1977) Regulation of ribosome production in *Escherichia coli*: Synthesis and stability of ribosomal RNA and of ribosomal protein messenger RNA at different growth rates. *J. Mol. Biol.* **115**, 335–354.
  61. T. E. Norris and A. Koch (1972) Effect of growth rate on the relative rates of messenger, ribosomal, and transfer RNA in *Escherichia coli*. *J. Mol. Biol.* **64**, 633–649.
  62. T. Gaal et al. (1997) Transcription regulation by initiating NTP concentration: rRNA synthesis in bacteria. *Science* **278**, 2092–2097.
  63. D. R. Gentry et al. (1993) Synthesis of the stationary phase specific sigma factor  $\sigma^S$  is positively

regulated by ppGpp. *J. Bacteriol.* **175**, 7982–7989.

64. A. Kuroda et al. (1997) Guanosine tetra- and pentaphosphate promote accumulation of inorganic polyphosphate in *Escherichia coli*. *J. Biol. Chem.* **272**, 21240–21243.
65. S. Kusano and A. Ishihama (1997) Functional interaction of *Escherichia coli* RNA polymerase with inorganic polyphosphate. *Genes to Cells* **2**, 433–441.

## Structure Databases

Structure databases contain the three-dimensional structures of biological [macromolecules](#) determined by [X-ray crystallography](#) and by [NMR](#). The primary structure database is Protein Data Bank (PDB), which is an archive of all publicly available three-dimensional structures of [proteins](#), [nucleic acids](#), carbohydrates, [viruses](#), and biomolecular complexes. Unlike the [sequence databases](#), which are organized by using database management systems, PDB is basically a collection of text files, each of which contains a structure entry deposited by the authors. An entry file consists of the textual information of definition, source, references, and comments, the sequence information, the **secondary structure** information, the three-dimensional information of atomic coordinates, crystallographic [structure factors](#), and NMR experimental data. Other useful structure databases include Nucleic Acid DataBase (NDB) for nucleic acids and Cambridge Structural Database (CSD) for organic and metal-organic compounds. The WWW addresses for these databases are given in [Table 1](#).

**Table 1. WWW Addresses for the Structure Databases**

| Database | Address  |
|----------|--|
| PDB      | <a href="http://www.pdb.bnl.gov">www.pdb.bnl.gov</a>             |
| NDB      | <a href="http://ndbserver.rutgers.edu">ndbserver.rutgers.edu</a> |
| CSD      | (site currently unavailable)                                     |

The first three-dimensional [protein structure](#) was elucidated for [myoglobin](#) in 1960 by John C. Kendrew (1). The need to computerize X-ray crystallographic data was immediately apparent, and the PDB was established in 1971 (2). Compared to one-dimensional sequence information, the three-dimensional structure information is far more difficult to analyze and comprehend. The PDB has stimulated computational research in this area, especially of the **protein folding** problem of predicting the three-dimensional structure from the amino acid sequence (see [Protein Structure Prediction](#)). PDB is the basis of our current knowledge of macromolecular structures, which is used, for example, in [threading protein sequences](#), where an amino acid sequence is matched against a library of representative three-dimensional folds of the polypeptide backbone. The information in PDB has been used to develop a number of value-added databases. For example, SCOP (3) provides a hierarchical classification of protein folds. Individual entries of the PDB also comprise a useful resource that can be examined in more detail, perhaps by using computer graphics and **computer simulation**, for basic understanding of structure–function relationships and applied research in drug



design and other areas. Although the majority of the databases in molecular biology are publicly available and freely accessible in the WWW (see [Databases](#)), some structure databases that inherited the tradition of chemical information are not in the public domain.

### Bibliography

1. J. C. Kendrew et al. (1960) *Nature* **185**, 422–427.
2. F. C. Bernstein et al. (1977) *J. Mol. Biol.* **112**, 535–542.
3. A. G. Murzin, S. E. Brenner, T. Hubbard, and C. Chothia (1995) *J. Mol. Biol.* **247**, 536–540.

### Structure Factor

In [X-ray crystallography](#), the structure factor  $\mathbf{F}(h\ k\ l)$  for a reflection  $(h\ k\ l)$  is related to the electron density distribution in the [unit cell](#).

$$\mathbf{F}(hk\ell) = V \int_{x=0}^1 \int_{y=0}^1 \int_{z=0}^1 \rho(xyz) \times \exp[2\pi i(hx + ky + \ell z)] dx dy dz \quad (1)$$

where  $x$ ,  $y$ , and  $z$  are fractional coordinates along the unit cell axes  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ ;  $V$  is the volume of the unit cell, and  $\rho$  the electron density distribution in the unit cell.  $\mathbf{F}$  can be regarded as a vector whose amplitude is  $|F|$  and whose phase angle is  $\alpha$ :  $F = |F| \exp(i\alpha)$ . The crystal scatters (reflects) only in specific directions given by the indices  $h$ ,  $k$ , and  $\ell$ . The reason is that the beams originating from the individual unit cells positively interfere only in these directions. In other directions they extinguish each other. The crystal “amplifies” the scattering by one unit cell and the scattered beam, or reflection  $(h\ k\ \ell)$ , is characterized by the structure factor  $\mathbf{F}(h\ k\ \ell)$  as calculated for one unit cell. Its amplitude  $|F(h\ k\ \ell)|$  is, apart from correction factors, equal to  $\sqrt{I(hk\ell)}$ , where  $I(h\ k\ \ell)$  is the intensity of the scattered beam. The phase angles cannot be derived straightforwardly from the X-ray diffraction pattern (see [Phase Problem](#)).

### Suggestion for Further Reading

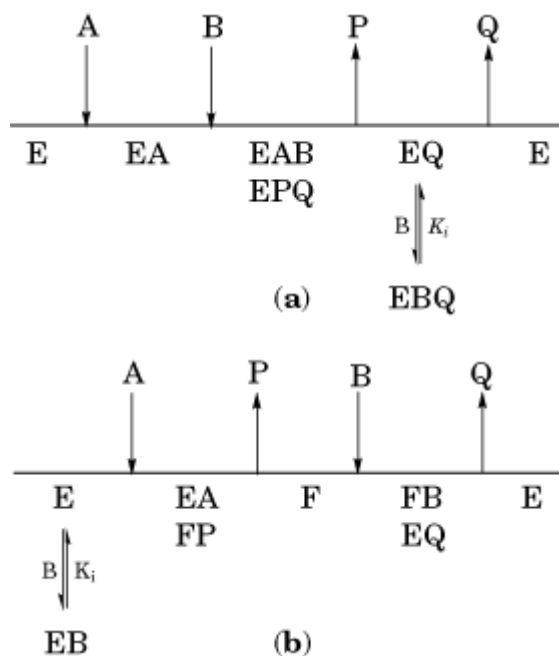
J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

### Substrate Inhibition

Inhibition of an [enzyme](#) by its substrate occurs whenever a **dead-end** enzyme-substrate complex forms. Such a reaction usually occurs at high substrate concentrations and is due to that substrate binding to a form of enzyme with which the product of the substrate normally combines. Two common **kinetic mechanisms** for substrate inhibition are illustrated in [Figure 1](#). For the ordered mechanism of [Figure 1a](#), substrate B combines with EQ to form a dead-end EBQ complex. B is behaving like P in combining with the EQ complex. The substrate inhibition by B with respect to A

would be [uncompetitive inhibition](#). Pyruvate behaves in this way for the [lactate dehydrogenase](#) reaction, which conforms to an ordered mechanism, with NADH being the first substrate to add and NAD the last substrate to leave the enzyme ([1](#)). The dead-end complex is E-NAD-pyruvate.

**Figure 1.** The two most common mechanisms of substrate inhibition: (a) Uncompetitive inhibition by substrate B with respect to the variable substrate A for an ordered Bi-Bi kinetic mechanism; (b) Competitive inhibition by substrate B with respect to the variable substrate A for a Ping-Pong kinetic mechanism.



Ping-Pong mechanisms are subject to substrate inhibition because of the structural similarity of the free enzyme (E) and the covalently modified enzyme (F). In Figure [1b](#), substrate B, which normally combines with F, reacts with E also. Since B, as an inhibitor, and A, as a varied substrate, combine with the same form of enzyme, the inhibition by B would be linear competitive with respect to A (see [Competitive Inhibition](#)). Double substrate inhibition can occur with Ping-Pong mechanisms because of the formation of both EB and FA complexes. When such inhibition occurs, the double-reciprocal [Lineweaver-Burk plot](#) looks quite complex. At high concentrations of the variable substrate, there is upward curvature of the plot, and the slope of the line increases as the concentration of the other, fixed substrate moves into the inhibitory range (see below).

The general equation for substrate inhibition is shown in equation [1](#).

$$v = \frac{VA}{K_a + A + A^2/K_i} \tag{1}$$

where  $V$  is the maximum velocity of the enzyme-catalyzed reaction,  $A$  the concentration of the substrate,  $K_a$  its  $K_m$  (Michaelis constant), and  $K_i$  its inhibition constant. This equation indicates that the reaction velocity would increase initially with the concentration of  $A$  and then decrease with further increases in substrate concentration (Fig. [2a](#)). Equation [1](#) can be rearranged in reciprocal form as equation [2](#):

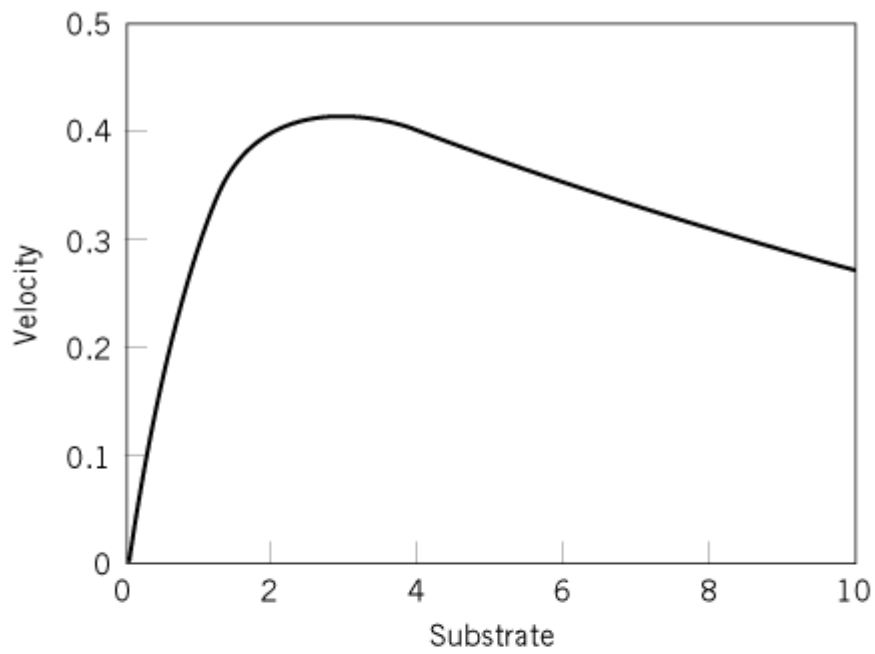
$$\frac{1}{v} = \frac{K_a}{V} \frac{1}{A} + \frac{1}{V} + \frac{A}{VK_i} \tag{2}$$

It should be noted that  $A$  appears in the denominator of the equation as a substrate and in the numerator as an inhibitor. The result is that an upward curvature is observed in double-reciprocal plots at higher concentrations of substrate (Fig. **2b**). Over the high range of substrate concentrations, equation **2** reduces to equation **3**:

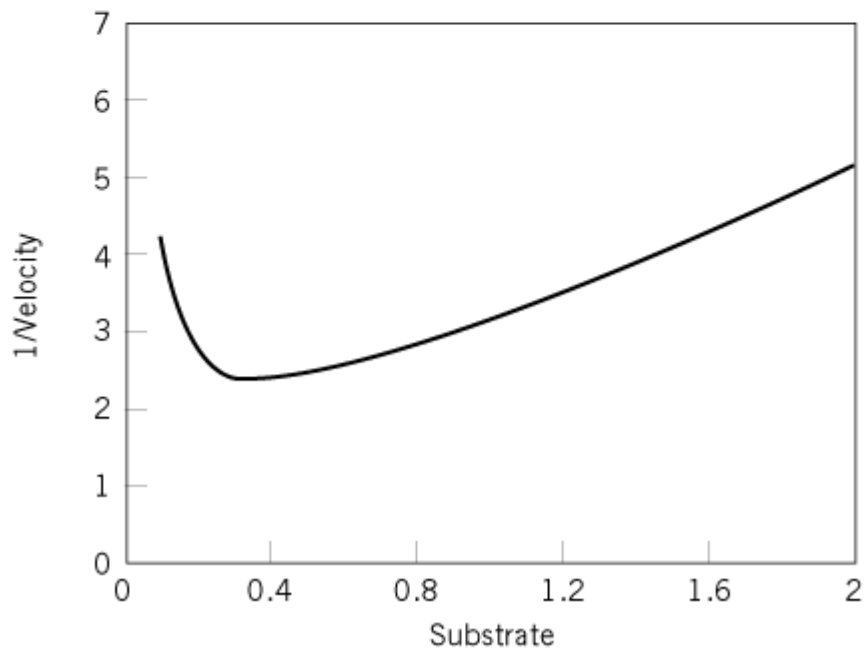
$$\frac{1}{v} = \frac{1}{VK_i}A + \frac{1}{V} \quad (3)$$

A plot of data from the curved portion of Figure **2b** would yield a straight line that intersects the abscissa at the point where  $-A = K_i$ .

**Figure 2.** Kinetic plots for substrate inhibition: (a) variation of initial velocity as a function of substrate concentration, and (b) double-reciprocal plot of the variation of initial velocity with substrate concentration.



(a)



(b)

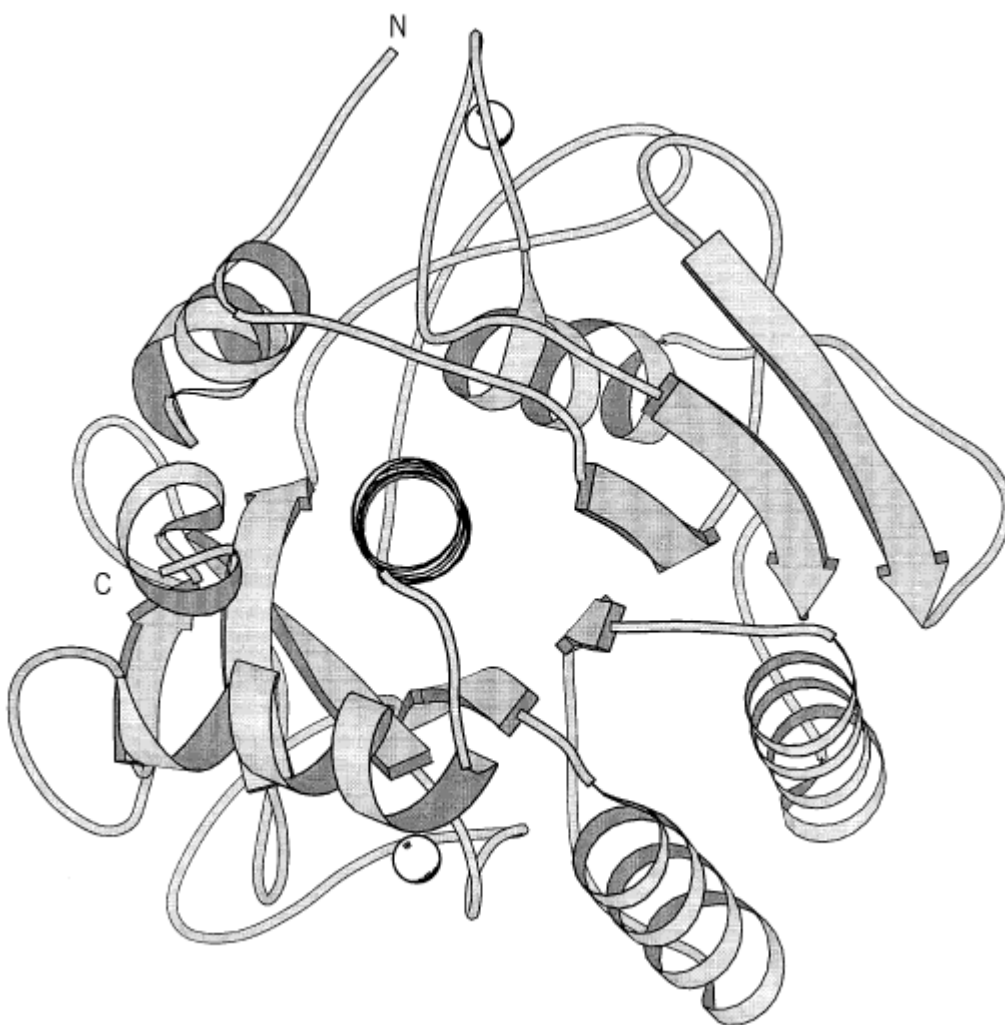
### Bibliography

1. E. H. Braswell (1975) In *Isozymes I Molecular Structure* (C. L. Markert, ed.) Academic Press, New York, p. 119.

### Subtilisin

Subtilisin is a member of the subtilase family of [serine proteinases](#), of which more than 50 have been identified (1) (Fig. 1). Family members are recognized on the basis of amino acid sequence [homology](#). The name subtilisin was originally given to an enzyme isolated from a culture of *Bacillus subtilis*. Later this culture was reclassified as *B. amyloliquefaciens*, and the enzyme is now known as subtilisin BPN' (2). A second subtilisin has been obtained from *B. licheniformis* and is known as subtilisin Carlsberg. In an attempt to avoid confusion, an alternative nomenclature has been introduced. Subtilisin Carlsberg is called subtilopeptidase A or subtilisin A, and BPN' is subtilopeptidase B or subtilisin B. The latter has also been referred to as Nagarse in earlier literature. The Enzyme Commission of the International Union of Biochemistry and Molecular Biology has assigned them to subsection E.C. 3.4.21.62 (serine proteinases are E.C. 3.4.21.). It should be noted that some suppliers refer to subtilisins as bacterial proteinases. Other members of the family have been given quite different trivial names such as thermitase and [proteinase K](#).

**Figure 1.** The three-dimensional structure of subtilisin. Only the polypeptide chain is depicted as a ribbon, with arrows for b-strands and coils for a-helices.



The active site of each subtilisin contains a [catalytic triad](#) typical of serine proteinases, and the catalytic mechanism similarly involves an acylenzyme intermediate. These enzymes have broad substrate specificity, although there is some preference for cleavage adjacent to neutral and acidic

amino acid residues. They can bring about extensive degradation of denatured proteins and can even degrade some native proteins. This ability to act on native proteins has been used to advantage to explore structure–function relationships of various proteins (3).

Subtilisin BPN' is such an effective proteinase that it has acquired commercial applications as a component of cleaning solutions. It is stable as a dry powder but undergoes autolytic degradation in solution. Stock solutions prepared in glycerol are stable for many months in the frozen state. Subtilisins are inactivated by **DIPF** and **PMSE**, as well as by oxidants. An enormous number of mutants have been prepared in efforts to alter virtually all aspects of the activity and stability of the enzyme (4).

#### Bibliography

1. R. J. Siezen, W. M. de Vos, J. A. Leunissen, and B. W. Dijkstra (1991) *Protein Eng.* **4**, 719–737.
2. A. J. Russell and A. R. Fersht (1986) *Nature* **321**, 733–737.
3. F. M. Richards and H. W. Wyckoff (1973) In *Atlas of Molecular Structures in Biology*, Vol. 1 (D. C. Phillips and F. M. Richards, eds.), Oxford University Press, Oxford, U.K., pp. 1–75.
4. J. A. Wells and D. E. Estell (1988) *Trends Biochem. Sci.* **13**, 291–297.

### Subtractive Hybridization

Subtractive hybridization is a technique for identifying and characterizing differences between two populations of **nucleic acids**. It detects differences between the **RNA** in different cells, tissues, organisms, or sexes under normal conditions, or during different growth phases, after various treatments (ie, **hormone** application, **heat shock**) or in diseased (or mutant) versus healthy (or wild-type) cells. Subtractive hybridization also detects **DNA** differences between different **genomes** or between cell types where deletions or certain types of genomic rearrangements have occurred. Subtractive hybridization techniques have identified many differentially expressed sequences from a wide variety of organisms. For example, recent studies using subtractive hybridization have identified mouse **complementary DNA** (cDNA) associated with **retinoic acid**-induced growth arrest (1), bacterial DNA differentiating clinical from nonclinical populations of *Staphylococcus aureus* (2), and plant **senescence**-associated genes from *Brassica napus* (3).

Subtractive hybridization requires two populations of nucleic acids; the tester (or tracer) contains the target nucleic acid (the DNA or RNA differences that one wants to identify), and the driver lacks the target sequences. The two populations are **hybridized** with a driver to tester ratio of at least 10:1. Because of the large excess of driver molecules, tester sequences are more likely to form driver–tester hybrids than double-stranded tester. Only the sequences in common between the tester and the driver hybridize, however, leaving the remaining tester sequences either single-stranded or forming tester–tester pairs. The driver–tester, double-stranded driver and any single-stranded driver molecules are subsequently removed (the “subtractive” step), leaving only tester molecules enriched for sequences not found in the driver. Usually multiple rounds of subtractive hybridization are necessary to identify truly tester-specific nucleic acid sequences.

There are five basic steps to subtractive hybridization: (1) choosing material for isolating tester and driver nucleic acids; (2) producing tester and driver; (3) hybridizing; (4) removing driver–tester hybrids and excess driver (subtraction); and (5) isolating of the complete sequence of the remaining target nucleic acid. Variations are possible at each step, and the materials used and methods chosen

depend on the desired results.

When choosing appropriate sources for driver and tester, it must be kept in mind that the less complex the source of tester and driver and the more sequences they have in common, the easier it is to isolate specific target sequence differences. For example, it is easier to identify RNA differences between cell types than it is to identify differences between tissues because fewer genes are expressed in single cells.

## 1. Preparation of Driver and Tester

In principle, both tester and driver samples can be either DNA or RNA, but it is often most practical for the tester to be DNA (because the tester is present in a low concentration, and DNA is more stable than RNA), and for the driver to be RNA (after hybridization, excess driver RNA can be eliminated enzymatically or by alkali degradation). In the basic subtractive hybridization protocol, RNA from the tester source is **reverse transcribed** into [complementary DNA](#) (cDNA) and hybridized to poly A<sup>+</sup> driver RNA. The tester–driver hybrids are removed, excess fresh driver is added, and the hybridization is repeated once. The remaining “target” cDNA is either **cloned** or used to make a probe. This basic procedure is useful if the starting material is not very complex and is easy to isolate.

If little starting tissue is available or if the starting material is complex, multiple rounds of hybridization-subtraction are needed, and it is necessary to use a **library-** or a **PCR-**based technique. Tester and driver are prepared from [cDNA libraries](#) as **phagemids** or as library inserts amplified by PCR or *in vitro* transcription. Alternatively, cDNA from tester and driver sources is ligated to different primers, amplified by PCR, and hybridized. The steps are repeated as needed.

## 2. Hybridization

When single-stranded nucleic acids are hybridized to each other, more abundant sequences anneal more rapidly because they encounter each other more frequently (see **C<sub>0</sub>t curve**). During subtractive hybridization, the hybridization step is driven by the excess driver sequences, so tester sequences that have complementary sequences in the driver population rapidly form driver–tester hybrids, whereas sequences unique to the tester population remain single-stranded or form tester–tester pairs more slowly. Rare sequences from either population take longer to pair up than abundant sequences. The ratio of driver to tester, the overall concentration of driver, the temperature, and the length of hybridization should be chosen based on the complexity of the driver and tester, the abundance class of the target nucleic acids, and the length of the driver and tester sequences used.

### 2.1. Subtraction

The purpose of the subtraction step is to remove driver–tester hybrids formed during the hybridization step, leaving behind tester enriched for the target sequences. Many different methods are used for subtraction, depending on the nature of the driver and the tester. A few possibilities are mentioned.

[Hydroxyapatite chromatography](#) is used to bind double-stranded driver and driver–tester hybrids, leaving single-stranded nucleic acids behind. This is a good choice if the driver is RNA because single-stranded RNA can be removed chemically or enzymatically, leaving only single-stranded cDNA tester after the subtraction.

If the tester is a single-stranded phagemid library and the driver is first-strand cDNA, after hybridization the double-stranded driver–tester hybrids can be digested with a frequent-cutting [restriction enzyme](#), and the hybridization mixture used to infect **bacteria**. Only the single-stranded tester phagemids infect, and they can thus be isolated.

A common procedure is to use [biotin](#)–streptavidin binding to separate nucleic acids. [Streptavidin](#) binds to biotinylated driver sequences, and phenol extraction is used to remove the streptavidin protein and the bound driver and driver–tester hybrids. Streptavidin can also be attached to beads or to a column and used to remove excess driver and driver–tester hybrids.

The effectiveness of the subtraction is monitored by using **radiolabeled** tester and determining whether the levels of single-stranded tester decrease after subtraction. Alternatively, enrichment for target sequences is monitored. If there are known genes common to the driver and tester and one or more specific to the tester, it can be determined, after each round of hybridization and subtraction, whether the tester-specific gene is becoming more abundant compared with the common genes.

## 2.2. Isolation of Target Sequences

After one or more hybridization and subtraction steps, the resulting tester nucleic acids should be greatly enriched for target sequences. However, it is still possible that rare sequences common to both the driver and the tester remain, and in many cases the sequences isolated are only partial gene sequences. The remaining tester sequences are isolated and analyzed in a variety of ways. Tester can be made into an enriched library and probed with driver and tester sequences to look for tester-specific clones, or the tester is labeled and used to probe tester and driver libraries and to isolate full-length clones. It is necessary to further analyze isolated tester sequences by **Northern blotting**, [in situ hybridization](#) or **PCR** methods to determine whether the sequences are truly tester-specific.

## 2.3. Alternatives to Standard Subtractive Hybridization Techniques

### 2.3.1. Positive Selection

An important alternative to subtractive hybridization is positive selection. Hybridization of tester and driver are still carried out but, rather than removing unwanted driver–tester and driver sequences by subtraction during step 4, double-stranded tester sequences are positively selected for selective cloning or selective amplification. Again, various methods are employed to carry out positive selection. A simple method is to digest tester with a restriction enzyme producing **cohesive ends**, while using sonication to shear the driver DNA randomly. After hybridization, [DNA Ligase](#) and **vector** DNA are added. Only double-stranded tester is cloned into the vector, and then it can be used to **transform** bacteria.

### 2.3.2. Representational Difference Analysis (RDA)

RDA is a positive selection technique employing PCR. RDA was originally used to identify differences between complex genomes, such as those caused by chromosomal rearrangements or losses due to cancer, infections with pathogens, and polymorphisms between individuals (4), and it was later adapted to analyze differences in gene expression (5). In both cases, tester and driver are ligated to adapters, amplified by PCR, the original adapters are removed, and new adapters (T2) are ligated only to the tester. After hybridization, only tester–tester DNA is amplified by using primers specific for the T2 adapter. The amplified tester is used again in further rounds of hybridization.

### 2.3.3. Suppression Subtractive Hybridization (6)

In this positive selection technique, both driver and tester are digested with a frequent-cutting restriction enzyme to give **blunt ends**. Tester is divided into two samples, which are ligated to different adapters, P1 and P2, and then hybridized to excess driver. Then the two tester populations are mixed, and additional driver is added. Hybrids formed between members of the two subtracted tester populations are selectively amplified by PCR using primers specific to P1 and P2. Molecules that have either P1 or P2 adapters at both ends form “panhandles” as the adapters hybridize to each other. and these molecules are not amplified by PCR (this results in the “suppression”).

### 2.3.4. Differential Display (7)

Differential display is a PCR-based technique that uses random amplification of cDNAs in different populations to identify differences between the populations. For each reaction, one primer is 5'-T<sub>11</sub>NN, where NN are any two specific nucleotides. This primer binds to a subset of cDNAs containing the two nucleotides complementary to NN immediately adjacent to the poly A tail. The



second primer is an arbitrary 10-mer. PCR using these two primers amplifies the same subset of expressed messenger RNA in each sample. Differences in amplification products between different samples are visualized by running the products on a sequencing gel. Band differences on the gel are cut out, and the DNA is eluted and cloned for further analysis. Differential display has the advantage that, for each primer set, a large number of different populations are comparable side by side on one sequencing gel. However, one disadvantage is that a large number of primers and reactions are needed to survey all of the expressed genes in a population.

#### 2.3.5. Serial Analysis of Gene Expression (SAGE) (8)

SAGE is based on the fact that nine bp of sequence located at its 3'-end is all the sequence information needed to identify a gene unambiguously. The first step in SAGE involves generating a 9-bp cDNA tag for each of the mRNAs in a population. Then many unrelated tags are concatenated, the concatenated tags are cloned, and random clones are sequenced. These sequences give a spectrum of the genes expressed in the tissue and indicate their relative abundance. Many genes can be analyzed at once, because only nine bp of each gene are sequenced and many tags are sequenced in a single reaction. To be useful for most purposes, however, full-length genes corresponding to the tags must be subsequently identified and isolated.

#### 2.3.6. Microarrays

cDNAs representing either known or unknown genes can be spotted onto glass and probed with different sources of fluorescently-labeled mRNA (9). Alternatively, 20-mer oligonucleotides (oligos) can be synthesized *in situ* in high-density arrays (10). Oligo sequences are derived from known gene sequences, and a number of oligos are prepared for each gene, so that there are many internal controls. Arrays can be synthesized that contain hundreds of thousands of oligos and can thus simultaneously monitor differences in the expression of tens of thousands of different genes corresponding to the cDNAs. A number of different probes can be analyzed at once because differently colored fluorescent labels are available. As the full repertoire of expressed sequences becomes available for different organisms, this technique will be extremely useful in monitoring genomewide changes in gene expression in different tissues, developmental states, and mutant backgrounds.

### 3. Conclusion

Many different methods have been used to analyze DNA differences between genomes and differences in gene expression (RNA differences). The use of subtractive hybridization has resulted in the isolation of many useful tissue-specific markers and interesting genes, but it is quite labor-intensive. Various positive selection techniques are more rapid, but they may not result in identifying all differentially expressed genes. Both subtractive hybridization and positive selection only compare two nucleic acid populations at once. The final products are used to make a subtracted probe or to produce a subtracted library, which then are used to identify and isolate full-length target sequences. If full-length cDNA libraries are used, it is not necessary to include this final step. Techniques, such as differential display, SAGE, and the analysis of microarrays, allow comparing a number of different mRNA populations, but differentially expressed genes are not specifically selected. Target genes are identified as a band on a gel (differential display), a very short sequence (SAGE), or a cDNA sequence (microarrays), so that a further step is needed to identify and/or clone full-length sequences. All of the techniques described are used with success, so the choice of technique for detecting DNA and RNA differences depends on the materials available to the investigator and the end results desired.

In the future, there will be so much sequence information available from genome sequencing and expressed sequence tag (see [Expressed Sequence Tag](#)) projects that, for investigators studying certain organisms, these types of differential cloning experiments will be superseded by sequence [database](#) searches. The potential for this type of database search has been demonstrated with the isolation of three prostate-specific genes identified by the comparison of human ESTs isolated from different tissue-specific libraries (11).

## Bibliography

1. R. A. Spanjaard, P. J. Lee, S. Sarkar, P. S. Goedegebuure, and T.J. Eberlein (1997) *Cancer Res.* **57**, 5122–5128.
2. W. A. el-Adhami, P. R. Stewart, and K. I. Matthaai (1997) *J. Med. Microbiol.* **46**, 987–997.
3. V. Buchanan-Wollaston and C. Ainsworth (1997) *Plant Mol. Biol.* **33**, 821–834.
4. N. Lisitsyn, N. Lisitsyn, and M. Wigler (1993) *Science* **259**, 946–951.
5. M. Hubank and D. G. Schatz (1994) *Nucleic Acids Res.* **22**, 5640–5648.
6. L. Diatchenko et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6025–6030.
7. P. Liang and A. B. Pardee (1992) *Science* **257**, 967–971.
8. V. E. Velculescu, L. Zhang, B. Vogelstein, and K. W. Kinzler (1995) *Science* **270**, 484–487.
9. M. Schena, D. Shalon, R. W. Davis, and P. O. Brown (1995) *Science* **270**, 467–470.
10. D.J. Lockhart et. al (1996) *Nature Biotechnol.* **14**, 1675–1680.
11. G. Vasmatazis, M. Essand, U. Brinkmann, B. Lee, and I. Pastan (1998) *Proc. Natl. Acad. Sci. USA* **95**, 300–304.

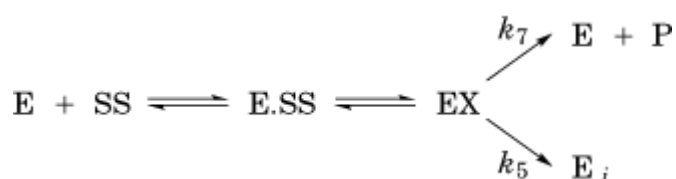
## Suggestion for Further Reading

12. C.G. Sagerstrom, B.I. Sun, and H.L. Sive (1997) *Subtractive cloning: Past, present, and future*, *Ann. Rev. Biochem.* **66**, 751–783. Excellent comprehensive review, many useful references for specific techniques.

## Suicide Inhibitor

A suicide inhibitor, also known as a *suicide substrate*, is a structural analogue of a substrate for an [enzyme](#) that contains a latent reactive group. Following reversible binding of the inhibitor at the enzyme's [active site](#), this group is activated by the enzyme during the course of its catalysis of the normal reaction to release a chemically reactive group within the active site. This reactive group can react covalently with an appropriately located active-site residue, inactivating the enzyme. That is, the enzyme catalyzes its own suicide (1). A wide range of functional groups can be catalytically unmasked by their target enzymes to produce electrophiles that cause inactivation. These include acetylenes, olefins, and beta-substituted amino acids. The chemical mechanism of their action has been outlined (2).

Suicide substrates are characterized kinetically by the concomitant formation of product and inactivation of the enzyme. The reactions can be formulated as



After the reversible formation of the enzyme–suicide substrate complex (E.SS), the enzyme catalyzes the reversible formation of an EX complex, where X is an activated intermediate that can either (1) be released as product P or (2) modify the enzyme. The partition ratio that describes the relative rates of the two reactions,  $k_7/k_5$ , remains constant over the time course of the reaction and is,

in effect, a measure of enzyme turnover relative to enzyme inactivation.

With a suicide inhibitor, a plot of the concentration of P as a function of time would show a time-dependent decrease in the rate of product formation, long before the equilibrium of the reaction was reached. Further, at infinite time, the enzyme would be completely inactivated, and the kinetic plot would exhibit a horizontal asymptote. In this respect, the behavior of a suicide substrate is similar to that of a **slow-binding enzyme inhibitor** whose action is, essentially irreversible. At infinite time, with classical, non-**allosteric** enzymes, a stoichiometric amount of intermediate should be covalently bound to all the enzyme molecules so that  $P_{\infty}/E_p$ , the ratio of the total amount of product formed,  $P_{\infty}$ , to the total amount of enzyme present,  $E_p$ , equals  $k_7/k_5$ , the partition ratio (3). Partition ratios vary from one for the inactivation of GABA aminotransferase by gabaculine and L-aspartate aminotransferase by vinylglycine, to values in excess of 1000 (1). The lower the partition ratio, the greater is the effectiveness of a suicide inhibitor as an enzyme inactivator.

#### Bibliography

1. C. Walsh (1982) *Tetrahedron* **38**, 871–909.
2. C. T. Walsh (1984) *Ann. Rev. Biochem.* **53**, 493–535.
3. S. G. Waley (1980) *Biochem. J.* **185**, 771–773.

#### Sulfate Salts

Ammonium, sodium, and magnesium sulfate have been used as agents for [precipitation](#) and [crystallization](#) of [proteins](#) for more than half a century.  $\text{Na}_2\text{SO}_4$  and  $(\text{NH}_4)_2\text{SO}_4$ , being on the **salting out** end of the [Hofmeister series](#), are good protein precipitants;  $\text{MgSO}_4$  is a weaker precipitant but is useful when  $\text{Mg}^{2+}$  ions interact with the protein.  $(\text{NH}_4)_2\text{SO}_4$  is the most frequently used salt, due to its great solubility.

The usual procedure in the use of these salts in the purification of a protein is that of fractional precipitation by a step-wise increase in salt concentration (1). To a crude protein extract, one adds  $(\text{NH}_4)_2\text{SO}_4$  up to a certain salt concentration. The precipitated material is centrifuged off, and more salt is added to the supernatant, precipitating more protein. In such manner, a step is reached that precipitates the protein of interest, usually detected by a specific bioassay. The fractionation is then refined. The active precipitate is dissolved in dilute buffer, and  $(\text{NH}_4)_2\text{SO}_4$  is added in much narrower steps, with the isolation of the active fraction.

In protein fractionation, ammonium sulfate is used most frequently because of its high solubility in water ( $\sim 4 M$ ). It would be preferable to use  $\text{Na}_2\text{SO}_4$  because of possible problems with subsequent nitrogen analysis and interference of the  $\text{NH}_4^+$  ion with some biological assays. The drawback to  $\text{Na}_2\text{SO}_4$  is its lower solubility ( $\sim 2 M$ ) that frequently is insufficient for the desired fractionation.

#### Bibliography

1. J. B. Sumner and G. F. Somers (1947) *Chemistry and Methods of Enzymes*, Academic Press, New York.

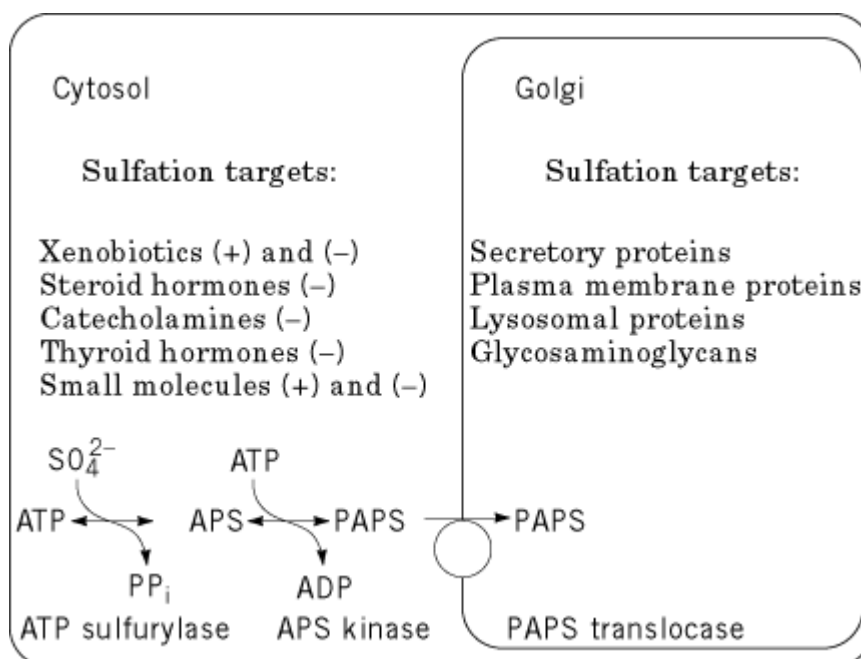
## Sulfation

Biotransformations that involve sulfate conjugation were first described more than 120 years ago. Originally, sulfation was connected mainly with drug inactivation and metabolism, which is initiated by an oxidative reaction that inactivates the target (phase I), followed by subsequent conjugation of, for example, sulfate or [glutathione](#) (phase II). Such conjugation increases the drug's water solubility and thereby facilitates excretion of the end product. Since then, many additional targets and functions of sulfate conjugation has been revealed: (i) inactivation of hormones and catecholamines, (ii) activation and inactivation of xenobiotics, (iii) removal of endogenous products (eg, bile salts), and (iv) modulation of [protein structure](#) and function.

The sulfation reactions are catalyzed by members of a family of enzymes designated sulfotransferases (ST). They use 3'-phosphoadenosine 5'-phosphosulfate (PAPS) as a co-substrate and sulfate donor. Thus, sulfation denotes conjugation of a given substrate with a sulfuryl group ( $-\text{SO}_3^-$ ); accordingly, the reaction should more correctly be designated as *sulfonation*. The functional groups of the affected substrates are typically either hydroxyl groups, resulting in formation of an ester, or unprotonated amines, leading to sulfamates.

Sulfation is widespread in nature, and sulfotransferases are found in both plants and animals. In mammals, sulfation appears to be ubiquitous in all tissues, but the type of sulfation and the target specificity may well be tissue-specific. The targets range from various xenobiotics, neurotransmitters, [steroid hormones](#), and [thyroid hormones](#), to [proteins](#), [glycoproteins](#), and glycosaminoglycans. The targets are localized both in the cytosol and in the [protein secretion](#) pathway (Fig. 1). These intracellular localizations define two subsets of sulfotransferases: (i) the soluble sulfotransferases in the cytosol and (ii) the membrane-associated sulfotransferases in the [Golgi apparatus](#). Along with their differences in enzymatic properties and substrate specificities, each subset will be described according to their intracellular localization.

**Figure 1.** Schematic diagram of the intracellular localization of sulfation reactions and targets in a typical mammalian cell. In the cytosol, brackets indicate whether sulfation leads to activation (+) or inactivation (-). An overview of PAPS synthesis and localization is also shown at the bottom. Xenobiotics refer to exogenic products administered to the cell, such as drugs.



## 1. PAPS Biosynthesis

The co-substrate PAPS is produced in the cytosol in a two-step enzymatic reaction from ATP and inorganic sulfate (Fig. 1). For sulfation reactions in the Golgi, PAPS is translocated to the lumen of the Golgi by a specific translocase (1). The rate-limiting step in PAPS formation is the formation of APS, and its steady-state concentrations are relatively low and display tissue-specific variations. Even in the liver, where PAPS concentrations are highest, stores can be depleted within minutes during maximal sulfation rates. Under normal conditions, PAPS concentrations are independent of the availability of inorganic sulfate. In rat liver, however, challenge of the conjugation reaction with high doses of exogenous substrate decreases the PAPS concentrations and depletes sulfate in the serum. This is in contrast to mice, in which the sulfotransferase activity limits the sulfation rate under corresponding conditions (2). Thus, sulfate availability, which changes with age and during physiological changes and diseases, can be affected by drug intake. Such depletion of sulfate stores can introduce serious toxic effects. Hence, regulation of PAPS synthesis is equally important for sulfation, as is the regulation of expression of the sulfotransferases.

## 2. Sulfation in the Cytosol

Sulfation in the cytosol involves a number of different functions and targets. The reactions have been studied both by endocrinologists and pharmacologists, and enzymes have been named after the substrate used for investigation. The overlapping substrate specificities has led to a confusing nomenclature of the enzymes, which is being unraveled by the [cloning](#) of the genes for the sulfotransferases. At present, five cytosolic sulfotransferases have been characterized in humans, but more will probably be identified in coming years (see Table 1).

**Table 1. Overview of the Best Characterized Human Cytosolic Sulfotransferases** <sup>a</sup>

| Enzyme | Endogenous Substrate | Xenobiotic Substrates |
|--------|----------------------|-----------------------|
|--------|----------------------|-----------------------|

|                   |  |                                       |
|-------------------|--|---------------------------------------|
| TS PST1 (P-PST)   | Iodothyronine, estrogen                      | Hydroxyarylamines, simple phenols     |
| TS PST2 (SULT1A2) | Unknown                                      | Hydroxyarylamines, simple phenols     |
| TL PST (M-PST)    | Catecholamines, iodothyronine                | 1-Naphthol, minoxidil, salbutamol     |
| EST               | Estrogen                                     | Ethinylastradiol, equilenin           |
| DHEAST (HST)      | Androgens, cholesterol, bile salts, estrogen | Aliphatic alcohols, benzylic alcohols |

<sup>a</sup> The most common names of the enzymes are given in abbreviations: TS PST-1 (–2), thermostable phenol ST-1 (–2); P-PST, phenol-preferring phenol ST; SULT1A2, sulfotransferase 1A2; TL PST, thermolabile phenol sulfotransferase; M-PST, mono-amine preferring phenol ST; EST, estrogen ST; DHEAST, dehydroepiandrosterone ST; HST, hydroxysteroid ST.

Sulfation in the cytosol has a number of endogenous targets. For example, sulfation plays a role in the inactivation of thyroid hormones. The prohormone thyroxine (T4) is converted into T3, the active component, in peripheral tissue. Inactivation occurs by deiodination or by sulfation by thermostable phenol sulfotransferase 1 (TS PST1). Sulfated T3 does not bind to the thyroid receptor, and sulfated T3 and T4 are rapidly degraded. Similarly, sulfation of catecholamines such as dopamine by TL PST (thermolabile phenol sulfotransferase) renders dopamine inactive as a neurotransmitter. In fact, about 90% of the dopamine circulating in plasma is found as sulfate conjugates, but the exact role of this derivative is at present unclear.

In addition, sulfation plays a key role in homeostasis and metabolism of another class of important hormones: the steroids. Sulfation of steroids excludes binding of the steroids to their receptors and thereby eliminates their effect. However, steroid sulfation may also be a way of transporting more soluble steroids, and the biological activity can be regained by the aid of tissue-specific sulfatases leading to specific receptor binding with subsequent signaling. The most abundant steroid in plasma is sulfated dehydroepiandrosterone (DHEA). The sulfated form of DHEA is metabolized 100-fold more slowly than nonsulfated DHEA, but the biological significance of its high plasma concentrations, which decrease with age, is not yet understood. However, DHEA serves as a precursor of both estrogen and androgen. DHEA is the target of DHEA sulfotransferase, which is highly abundant in the liver. This enzyme is also responsible for the sulfation of cholesterol and secondary bile salts. Secondary bile salts are generated from salts that are secreted by the liver and by bacteria in the intestine. These salts have a toxic effect on hepatocytes, but sulfation facilitates the excretion of the bile salts through the urine.

Sulfation also plays an important role in the biochemistry of estrogen. Sulfated DHEA is the major precursor of estrogens and androgens. Also, sulfated estrone is the most abundant estrogen in the plasma. Like other steroid hormones, estrogens are inactivated by sulfation. Hence, sulfotransferases serve a regulatory role of estrogen activity. For example, estrogen sulfotransferase is expressed in a cyclic manner in the endometrium and is induced by progesterone. Thus, this sulfotransferase may regulate stimulation by estrogen during the menstrual cycle. However, three sulfotransferases are capable of sulfating estrogens, albeit with different affinities, and this makes characterization of the reactions difficult.

Besides inactivation, sulfation can activate a number of compounds. One example is minoxidil, an antihypertensive agent, which is active only after endogenous sulfation. However, sulfation is also engaged in a less fortunate type of bioactivation: generation of [mutagens](#). Many [carcinogens](#), including dietary procarcinogens such as safrole and estragole, only obtain their mutagenic

properties after metabolic modification, such as by sulfation. This effect of sulfation was first demonstrated for hydroxamic acids and hydroxylamines, but at present the list of affected compounds is long, and most cytosolic sulfotransferases are involved in the process.

In recent years, much has been learned about the cytosolic transferases and their activities. Recently, the first three-dimensional [protein structure](#) of estrogen sulfotransferase was presented (3), and new sulfotransferase variants will probably be identified. The availability of sulfotransferases as [recombinant proteins](#) will also facilitate future studies of the bioactivity of the sulfotransferases.

### 3. Sulfation in the Golgi

In the Golgi, two classes of molecules are sulfated: proteins, which are sulfated on [tyrosine](#) residues, and glycosaminoglycans. The types of sulfotransferases involved are different and do not overlap in substrate specificity.

Sulfation of tyrosine residues was first discovered as a modification of [fibrinogen](#) in the mid-1950s. In the 1960s, the modification was found in a few peptide hormones, but only in the beginning of the 1980s was the truly widespread nature of tyrosine sulfation recognized (4). Since then, sulfation of a number of proteins has been demonstrated. Tyrosine sulfation occurs in the *trans* Golgi network; as a consequence, the proteins affected are secretory proteins, [membrane proteins](#), and presumably lysosomal proteins (5). The modification is the most common side-chain modification of tyrosine residues, and up to 1% of the tyrosine in a given cell may be sulfated. Tyrosine sulfation appears to be ubiquitous in all tissues and has been traced back to *Caenorhabditis elegans*.

Tyrosine sulfate appears to be involved in various types of [protein–protein interactions](#). The following effects are well established: (i) Tyrosine sulfation of the peptide hormone cholecystokinin is a prerequisite for receptor binding; (ii) sulfation of the leech anticoagulant protein [hirudin](#) increases the affinity toward [thrombin](#); (iii) sulfation increases the affinity between factor VIII and von Willebrand factor, (iv) sulfation is necessary for the binding of P-selectin to P-selectin glycoprotein ligand 1; (v) sulfation increases the intracellular processing of the peptide hormone gastrin; and (vi) sulfation alters the secretion kinetics of Drosophila yolk protein 2. For many of the proteins that are known to be sulfated (eg, the b-[amyloid precursor protein](#)), the function is unknown. This is in part due to the difficulties in prediction of tyrosine sulfation sites. Furthermore, the high frequency of tyrosine sulfation suggests that more proteins than those known today will be found to be sulfated. Based on the structures surrounding known tyrosine sulfation sites, and *in vitro* experiments using synthetic peptides, a **consensus sequence** for prediction of sulfation has been suggested. The main feature is the presence of an acidic residue on the amino-terminal side of the tyrosine and additional acidic residues neighboring the sulfation site. At least partial sulfation can be obtained, however, without an acidic residue on the amino-terminal side of the tyrosine residue, as found in the peptide hormone gastrin. It is at present unclear how common partial sulfation is and what are the precise structural requirements for sulfation. With the recent identification and cloning of two tyrosylprotein sulfotransferases [TPST-1 (6) and TPST-2 (7, 8)], new tools are available for precise definition of the requirements. It is also possible that additional TPST family members will be identified.

The two TPST [complementary DNAs](#) each encode a type II integral [membrane protein](#) that displays an overall identity of 64% in amino acid sequence. Both enzymes are **glycoproteins**, of 370 and 377 amino acid residues, respectively, and contain a short (7-residue) cytosolic tail that is probably involved in intracellular transport and localization of the enzymes. Little is known about the regulation of the enzymes, but, with the recent cloning, this will be a topic for further investigation.

Another class of molecules sulfated in the Golgi is the glycosaminoglycans (GAGs). These molecules are disaccharides and constituents of the proteoglycans, proteins that are also glycosylated to an extreme degree, ending up with molecular weights of up to several million. The proteoglycans consists of a complex protein core of varying composition conjugated with between 1 and 100 GAG

chains. There are seven classical GAGs, six of which are sulfated in at least one position: chondroitin-4-sulfate, chondroitin-6-sulfate, dermatan sulfate, keratan sulfate, heparin, and heparan sulfate. Accordingly, large proteoglycans can contain a high number of negatively charged groups. These charges are central for the physiological functions of glycosaminoglycans and proteoglycans in interacting with other proteins, similar to the biological role of tyrosine sulfation. Glycosaminoglycans were traditionally considered structural components of the [extracellular matrix](#) in connective tissue. However, recent years have disclosed much more active roles of glycosaminoglycans.

The best-characterized effect of sulfation is in the anticoagulant effect of heparin. Heparin binds antithrombin III in a reaction that is strictly dependent on sulfation of a specific position of heparin and alters the conformations of both molecules. Heparin also modulates the effect of basic [fibroblast growth factor](#) (FGF), and the receptor binding of this growth factor is affected by heparan sulfate. Recent studies suggest that selection between FGF-1 and FGF-2 by the activating proteoglycan during development is mediated by changes in heparan sulfate side chains, and possibly by a change in sulfation pattern during differentiation and growth arrest (9). Similarly, the chondroitin sulfates, which are involved in cell adhesion, cell migration, and possibly also neural development, show differences in sulfation pattern during development and tumor progression that are likely to influence the physiological function of the proteoglycans carrying these sulfates. Considering the increased interest in the functions of proteoglycans in [signal transduction](#) and other types of cell signaling, along with the developmental changes in compositions of glycosaminoglycans, additional examples of functional significance of sulfate in proteoglycans can be expected.

#### 4. Homology of Sulfotransferases

The sulfotransferases with all the different targets described here have surprisingly little [homology](#), suggesting that these enzymes may have evolved independently. This was disputed recently, when homology between two protein regions involved in PAPS binding from cytosolic transferases and from several glucosaminoglycan/heparan sulfotransferases was demonstrated (10). Following the cloning of the genes for the two TPSTs, a general comparison of all the types of sulfotransferases is possible. This shows that the TPSTs differ from the other membrane-associated sulfotransferases and from the cytosolic enzymes, even in the most conserved regions (Fig. 2). Hence, if the sulfotransferase have evolved from a common ancestor, they have diverged sufficiently to have very different [primary structures](#).

**Figure 2.** Alignment of partial protein sequences of sulfotransferases. Two spatially separated regions that are important for PAPS binding are the primary homologous regions of sulfotransferases. They are denoted the PSB loop and the 3'-binding site. All sulfotransferases were aligned as described (8) except for the TPSTs. The abbreviations of the protein sequences are: mEST, mouse estrogen sulfotransferase; mPST, mouse phenol sulfotransferase; mHST, mouse hydroxysteroid sulfotransferase; FST, flavonol sulfotransferase (plants); hTPST, human tyrosylprotein sulfotransferase; GalCerST, human 3'-phosphoadenylylsulfate galactosylceramide sulfotransferase; hHSNST, human heparan sulfate *N*-deacetylase/*N*-sulfotransferase; hHS2OST, human heparan sulfate 2-sulfotransferase.



|          |     | 5' PSB   |     | 3' PB |                  |      |
|----------|-----|----------|-----|-------|------------------|------|
| mEST     | ATY | PKSGTTWI | SEV | CKM   | IYLCRNAKDVAVSYYY | FLLM |
| mPST     | STY | PKSGTNWM | SEI | IKV   | IYVARNAKDVVVSYYN | FYKM |
| mHST     | LTY | PKSGTNWL | NEI | AKA   | IYLMRNPRDILVSGYF | FWGN |
| FST      | ASY | PKSGTTWL | KAL | CKI   | VYIYRNMKDVIVSYYH | FLRQ |
|          |     |          |     |       |                  |      |
| hTPST-1  | GGV | PRSGTTLM | RAM | AKF   | LLMVRDGRASVHSMIS | RKVT |
| hTPST-2  | GGV | PRSGTTLM | RAM | SKF   | LLMVRDGRASVHSMIT | RKVT |
|          |     |          |     |       |                  |      |
| GalCerST | LKT | HKTASSTL | LNI | AIF   | ITVLRDPARLFESSFH | YFGP |
| hHSNST   | IGP | QKTGTTAL | YLF | AKV   | LTILINPADRAYSWYQ | HQRA |
| hHS2OST  | NRV | PKTASTSF | TNI | PIY   | INVIRDPIERLVSYYY | FLRF |

## Bibliography

1. C. Abeijon, E. C. Mandon, and C. B. Hirschberg (1997) *TIBS* **22**, 203–207.
2. L. Liu and C. D. Klaassen (1996) *Toxicol. Appl. Pharmacol.* **139**, 128–134.
3. Y. Kakuta, L. G. Pedersen, C. W. Carter, M. Negishi, and L. C. Pedersen (1997) *Nature Struct. Biol.* **4**, 904–908.
4. W. B. Huttner (1982) *Nature* **299**, 273–276.
5. Hille, T. Braulke, K. V. Figura, and W. B. Huttner (1990) *Eur. J. Biochem.* **188**, 557–586.
6. Y.-B. Ouyang, W. S. Lane, and K. L. Moore (1998) *Proc. Natl. Acad. Sci USA* **95**, 2896–2901.
7. R. Beisswanger, D. Corbeil, C. Vannier, C. Thiele, U. Dohrmann, R. Kellner, K. Ashman, C. Niehrs, and W. B. Huttner (1998) *Proc. Natl. Acad. Sci. USA* **95**, 11134–11139.
8. Y.-B. Oyang and K. L. Moore (1998) *J. Biol. Chem.* **273**, 24770–24774.
9. Y. G. Brinckman, V. Nurcombe, M. D. Ford, J. T. Gallagher, P. F. Bartlett, and J. T. Turnbull (1998) *Glycobiology* **8**, 463–471.
10. Y. Kakuta, L. G. Pedersen, L. C. Pedersen, and M. Negishi (1998) *TIBS* **23**, 129–130.

## Suggestions for Further Reading

11.

Papers presented at the Third International Sulfation Workshop in Drymen, Scotland in 1996 are in a special issue of *Chem.-Biol. Interact* (1998) volume 109. Additional reading:

### Sulfation in the Cytosol:

12. R. M. Weinshilboum et al. (1997) *FASEB J.* **11**, 3–14.
13. M. A. Runge-Morris (1997) *FASEB J.* **11**, 109–117.
14. C. N. Falany (1997) *FASEB J.* **11**, 206–216.

### PAPS biosynthesis:

15. C. D. Klaassen and J. W. Boles (1997) *FASEB J.* **11**, 404–418.

### Sulfation in the Golgi Network:

16. C. Niehrs, R. Beisswanger, and W. B. Huttner (1994) *Chem.-Biol. Interact.* **91**, 257–271.

17. P. Jollés (ed.) (1994) *Proteoglycans*, Birkhäuser Verlag, Basel.

## Sulfur Isotopes

Sulfur (or sulphur) (1) is element number 16 in the periodic table. Thirteen isotopes of sulfur have been studied (2), ranging in atomic mass from  $^{28}\text{S}$  (half-life = 0.12s) to  $^{40}\text{S}$  (half-life = 9s) (see [Radioactivity](#) and [Radioisotopes](#)). Four stable isotopes are found in nature:  $^{32}\text{S}$  at 95.02%,  $^{33}\text{S}$  at 0.75%,  $^{34}\text{S}$  at 4.21%, and  $^{36}\text{S}$  at 0.02% abundance.

The most important radioactive isotope of sulfur is  $^{35}\text{S}$  (half-life = 87.44days). It decays by beta-minus emission to chlorine-35, which is stable. Sulfur-35 yields one beta particle per decay, with an energy of 0.167 MeV maximum, 0.0488 MeV on average.

Sulfur-35 is prepared in an accelerator by bombarding a stable chlorine-37 target with deuterons according to the reaction  $^{37}\text{Cl}(d,a)^{35}\text{S}$ . Carrier-free sulfur-35 may be prepared in a reactor according to the reaction  $^{35}\text{Cl}(n,p)^{35}\text{S}$ . Sulfur-35 may also be prepared by activation of natural sulfur according to the reaction  $^{34}\text{S}(n,g)^{35}\text{S}$ .

Sulfur-35 was first used in biological studies in the 1930s to study sulfur metabolism. Penicillin was labeled with sulfur-35 in 1940. Plasma proteins were first labeled with  $^{35}\text{S}$ -**cysteine** in 1943 using cyclotron-produced sulfur-35 (3). Ceccaldi and co-workers in 1953 used sulfur-35 to identify the origin of enterocytes and mucocytes and their migration along villi walls (4).

Sulfur-35-labeled [cysteine](#) and [methionine](#) are incorporated into proteins and used to monitor protein biosynthesis. One of the advantages of  $^{35}\text{S}$ -labeling is that high specific activities can be achieved (eg, about 200 times greater than that possible with  $^{14}\text{C}$  labels). Consequently, sulfur-35 is used in molecular biology whenever possible. Sulfur-35 is detected by beta-particle liquid scintillation counting and by [autoradiography](#) and [fluorography](#).

## Bibliography

1. D. R. Lide and H. Pr. Frederikse, eds. (1995) *CRC Handbook of Chemistry and Physics*, CRC Press, Boca Raton, Fla.
2. Knolls Atomic Power Laboratory (1966) *Chart of the Nuclides*, 15th ed., available from General Electric Company, San Jose, Calif.
3. M. Brucer (1990) *A Chronology of Nuclear Medicine*, Heritage Publications, Inc., St. Louis, Mo.
4. P. F. Ceccaldi, J. Verne, J. Biez-Chareton, and R. Wegmann (1953) in *Proc. French Society of Histochemistry*, Dec. 1953 abstract book, p. 21.
5. *The Radiochemical Manual*, 2nd ed. (1966) Amersham, Bucks, England.

## Suggestion for Further Reading

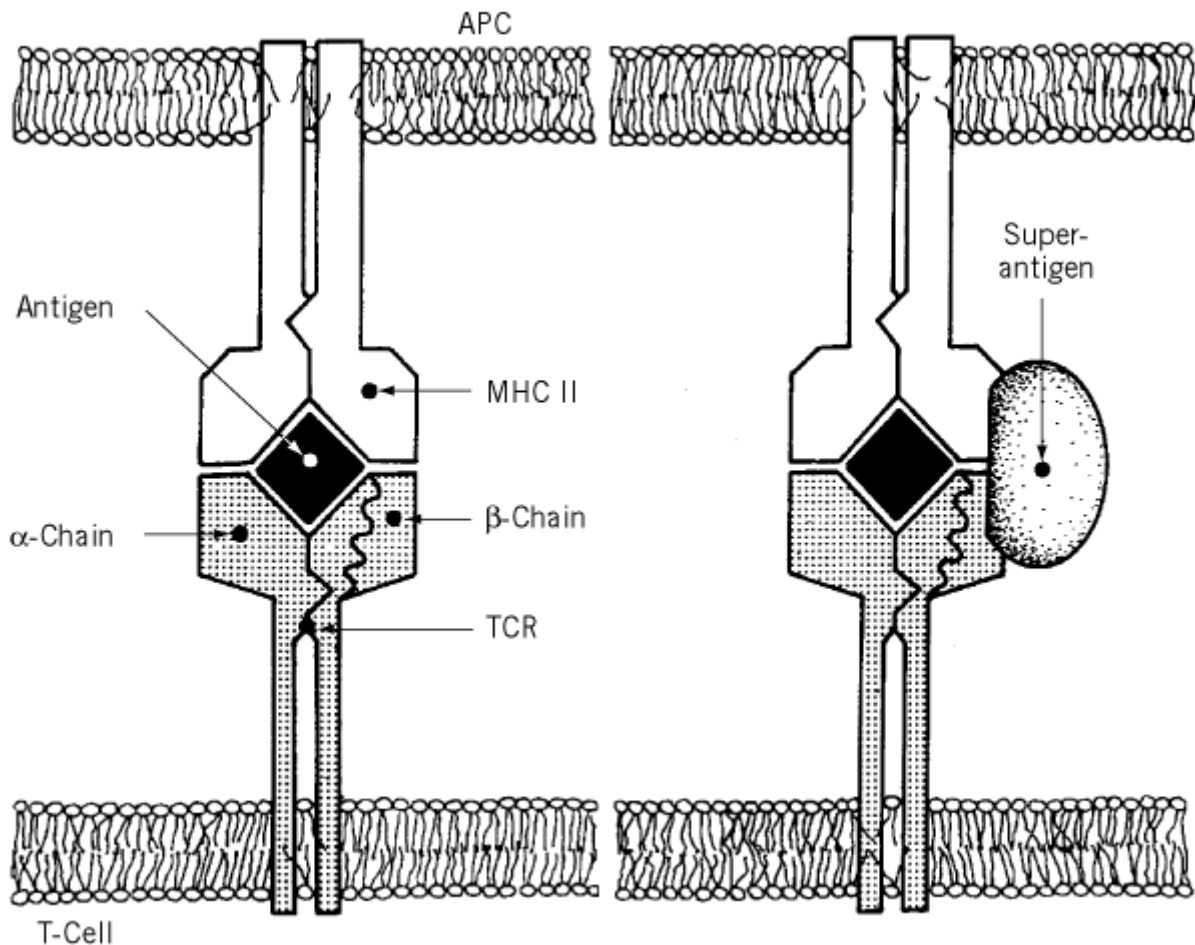
6. D. D. Dziewiatkowski (1962) "Sulfur", In *Mineral Metabolism*, Vol. 2, part B, C. L. Comar and F. Bronner, eds., Academic Press, New York.

## Superantigen Toxins

The vital importance of the immune and inflammatory responses against microbial infections is underlined by the many bacterial [toxins](#) active against cells involved in body defense. Several toxins, such as **mellitin and membrane-perturbing toxins** and [pertussis toxin](#), manifest their effects primarily, but not uniquely, on inflammatory and immune cells. More specific are the 22- to 28-kDa protein superantigenic toxins involved in staphylococcal food poisoning and toxic shock syndromes, which are produced by different species of *Streptococci* and *Staphylococci* (1, 2). It is likely that similar protein toxins are produced by several other bacteria. These toxins share various degrees of [primary structure](#) similarities, but are structurally and functionally very similar (3). The protein molecules have two lobes: One is organized as a [beta-sheet](#) and the other as a b-barrel, with hydrophobic residues lining the inner surface of the cylinder.  $\alpha$ -Helices are concentrated in the region connecting the two **domains**, with a key role played by a central long  $\alpha$ -helix. Two major grooves are present on the opposite sides of this helix, one of them appears to be involved in binding [major histocompatibility complex](#) (MHC) class II (MHC-II) and is adjacent to the [T-cell receptor](#) (TCR) binding site.

The most relevant aspect of their superantigen toxin action appears to be the stimulation of T-cell proliferation, which has led to their general classification as [superantigens](#). This stimulation is quite different from that of most other mitogens like [lectins](#) (eg, concanavalin A). Superantigen toxins are capable of a double interaction: with MHC-II molecules and with the TCR (Fig. 1). A ternary complex forms upon antigen presentation to T-cells by an **antigen-presenting** cell in the presence of a superantigen toxin, which stimulates T-cell proliferation and the release of **cytokines**. However, at variance from normal [antigens](#), which are processed intracellularly and presented on a groove on the MHC-II surface, superantigen toxins bind to a different site of the MHC-II molecule, with **dissociation constants** in the  $10^{-5}$  to  $10^{-8}$  M range. One group of superantigen toxins binds MHC-II molecules via an atom of  $Zn^{2+}$ , which bridges three toxin histidine residues, one located on the  $\beta$  chain of MHC-II; other superantigen toxins bind directly to different sites on the class II molecule.

**Figure 1.** Stimulation of T-cell proliferation by conventional antigens (left) and bacterial superantigens (right). Both antigens are recognized only when presented to the T-cell antigen receptor (TCR) by a class II molecule of the major histocompatibility complex (MHC-II). Antigens are processed intracellularly by antigen-presenting cells (APC) and presented to the T cells as small peptides on a groove on the surface of MHC-II. In contrast, superantigens bind a different site on MHC-II, leaving the groove for peptide presentation free and functional. The normal antigen–MHC-II complex interacts with TCR using multiple contacts involving both variable (V) and constant segments of the TCR  $\alpha$  and  $\beta$  chains, and this ensures stimulation of the few T cells bearing the appropriate receptor. The superantigen–MHCII complex interacts with the  $V\beta$  region of TCR, leading to an unspecific stimulation of a large proportion of the total T-cell population.



The MHC-II-complexed T-cell [epitopes](#) interact with the TCR via multiple contacts to both its  $\alpha$  and  $\beta$  chains (Fig. 1). In contrast, the superantigen toxin interacts merely with the  $V\beta$  region of the TCR, whereas it interacts with the  $V\gamma$  region of  $gd$ -T-cell receptors. The  $V\beta$  region is well conserved, and it might be involved in the crucial elimination of antiself T-cells. This reduces the possibility that the host mutates its T-cell receptor to delete the superantigen-binding region.

Superantigen toxins induce an unspecific stimulation of a large proportion of the total T-cell population of the host (from 8 to 40%, depending on the type of toxin) with a massive release of cytokines, which cause a variety of pathological consequences, including fever, malaise, nausea, vomiting, and diarrhea (common symptoms of food poisoning) and systemic shock reaction. Superantigen toxins can also potentiate the effect of lipopolysaccharide endotoxins, and they may activate other cells, such as mast cells. T-cell clones amplified by superantigens often disappear, or become inactive after being stimulated, and this might lead to immune depression.

#### Bibliography

1. G. Menestrina, G. Schiavo, and C. Montecucco (1994) *Mol. Aspects Med.* **15**, 81–193.
2. B. Fleischer (1995) *Rev. Med. Microbiol.* **6**, 49–57.
3. T. S. Jardetzky et al. (1994) *Nature* **368**, 711–718.

## Superantigens

Superantigens are a special type of [antigen](#) that bind to **multiple histocompatibility complex** (MHC) class II molecules and to polymorphic regions of the b chain of the **ab T-cell receptor** (TCR), without involving interacting structures that condition normal MHC–peptide complex recognition by the TCRab specific combining site. All **T cells** that express a TCR with a b chain bound by a superantigen will be activated, which may represent an important fraction of the T-cell population.

Extensive analysis of this phenomenon provided understanding of the nature of Mls (minor lymphocyte stimulating) antigens, described long ago by Festenstein. Mls antigens were detected *in vitro* using lymphocytes from strains of mice having identical MHC, but different Mls. When they were mixed, the lymphocytes proliferated, which was called *mixed-lymphocyte reaction* (MLR). It turned out that these Mls antigens were encoded by mouse mammary tumor viruses (MMTVs) that had been integrated into the mouse [genome](#) and genetically transmitted. An interesting phenomenon that is associated with these viral superantigens is that they induce selection against T cells that express the TCR b chain that interact with the endogeneously expressed MMTV. This is the result of clonal deletion that occurs at the double-positive (DP) stage, when thymocytes are about to leave the thymic cortex. Entire Vb families may thus be eliminated by one given Mls/MMTV antigen and are therefore absent from the corresponding mouse strains. It has been shown that MMTV viral antigen binds to the b chain of the MHC class II molecules, at the external side of the **a-helix**, in a region that is not involved in recognition by the TCR combining site.

Another well-studied example of superantigens is provided by bacterial toxins [such as staphylococcal enterotoxins (Se A-E)] that cause food poisoning in humans or toxic shock syndrome (TSST-1) (see [Cholera Toxin and Enterotoxins](#)). [X-ray crystallography](#) structures of complexes of MHC class II molecules with these enterotoxins indicate that the bacterial superantigen binds to the a chain.

**B-cell** superantigens have more recently been described as proteins that also bind to **variable regions of immunoglobulins**, leading to activation of naive B cells at a much greater frequency than when stimulated in the conventional way.

More and more examples of superantigens are being described and found to be responsible for the induction of immunopathological disorders; this is undoubtedly an expanding chapter of pathology, as recently illustrated by the observation of superantigen activity associated with human **herpes virus**, [cytomegalovirus](#), and **Epstein–Barr virus**.

See also entries [Antigen](#), [T Cell](#), **T-cell receptor**.

### Suggestions for Further Reading

P. Marrack, E. Kushnir, and J. Kappler (1991) A maternally inherited superantigen encoded by a mammary tumour virus. *Nature* **349**, 524–526.

H. Acha-Orbea and E. Palmer (1991) Mls—a retrovirus exploits the immune system. *Immunol. Today*, **12**, 356–361.

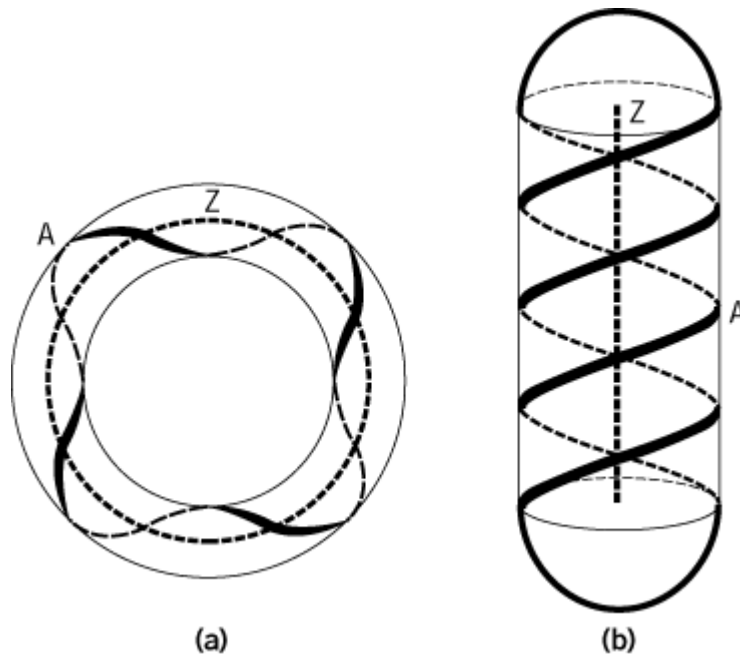
B. T. Huber, P. N. Hsu, and N. Sutkowski (1996) Virus-encoded superantigens. *Microbiol. Rev.* **60**, 473–482.

G. J. Silverman, J. V. Nayak, and A. La Cava (1997) B-cell superantigens: molecular and cellular implications. *Int. Rev. Immunol.* **14**, 259–290.

## Supercoiling of DNA

The geometry of closed circular DNA (see [DNA Topology](#)) includes the trajectory of the axis and the twist of the duplex taken to be a continuous variable along the duplex axis. For this purpose, the DNA is usually modeled either as an elastic rod or as a full or half ribbon ([1](#)) (see [DNA Structure](#)). There are two classes of superhelical structure, which can be reversibly interconverted without strand scission. The class characteristic of purified, monomeric closed circular DNA is the plectonemic or interwound superhelix, in which the DNA axis winds solenoidally up and down a virtual cylinder, forming a double helix, with smoothly curved connection at the ends of the cylinder. The type often encountered in DNA either associated with other DNA molecules (catenated DNA) or wrapped on a protein surface is the toroidal superhelix. Here the DNA axis winds solenoidally around the surface of a torus or a segment of a torus, forming a single helix. The two types of superhelix are illustrated in Figure [1](#).

**Figure 1.** The two types of superhelix winding in superhelical DNA. **(a)** Model of a toroidal superhelix. The figure shows the DNA axis, A, winding four times in a right-handed sense about the superhelical axis, Z. The superhelical axis is circular, and the DNA axis appears to be wrapped on the surface of a circular torus. **(b)** Model of a plectonemic or interwound superhelix. The figure shows the DNA axis, A, winding two times up and two times down a linear superhelical axis, Z. The axis lies at the center of a cylinder, on the surface of which the DNA appears to be winding. The DNA axis is joined to itself at the cylinder caps along a great circle.



Concentration affects both factors; the stiffness (bending modulus) decreases with increasing ionic strength ([2](#)), but so do the intraduplex repulsions ([3](#)). The response of superhelical structure to changing ionic conditions, as well as to the temperature, is therefore expected to be complex.

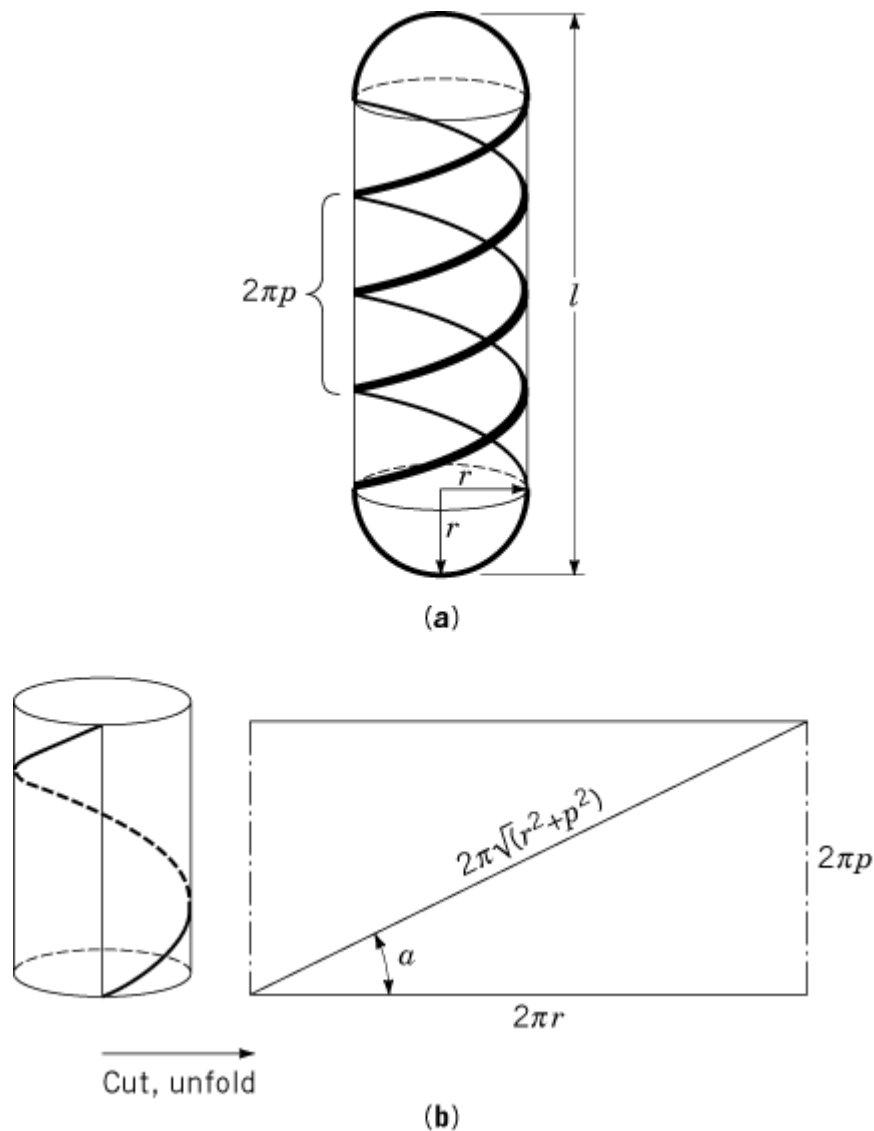
DNA behaves like a chain polymer in solution and is consequently flexible. This applies to superhelical DNA as well, leading to the expectation that the structure of the superhelix fluctuates significantly about some sort of average. The detailed structure will be determined by a balance among forces. One of the more important of these forces results from the tendency to minimize the bending distortion, which encourages the opposing duplex strands to approach nodal self-contact. Actual self-passage would, of course, nearly eliminate the bending distortion. This factor tends to reduce the superhelix radius,  $r$ . Polyanionic repulsion forces, on the other hand, tend to repel the opposing duplex strands and to increase  $r$ . Any change in the salt composition and

The geometry of plectonemic superhelical DNAs has been investigated experimentally within a fairly narrow range of conditions. The important tertiary structural properties are the number of supercoils ( $n$ ); the superhelix pitch ( $2pp$ ) and radius ( $r$ ), the number of branches, and the average twist and writhe. These properties are illustrated schematically in Figure 2. The geometric quantities have been estimated for closed circular DNAs of 3.5 and 7.0 kbp by a combination of conventional [electron microscopy](#) and site-specific [recombination](#) (4). The results apply to a buffer containing 20 mM Tris-HCl, 50 mM NaCl, and 20 mM MgCl<sub>2</sub>, although some experiments in the electron microscope used much lower salt (10 mM Tris-HCl); no major differences were observed between these two solvents by electron microscopy. Under these conditions, the superhelix pitch was found to be nearly constant, independent of  $s$  and  $N$ , at 54–56°. The superhelix radius depends on the superhelix density according to

$$\frac{1}{r} = 0.00153 - 0.268\sigma \quad (1)$$

where  $r$  is measured in Å. The length of the superhelix axis is independent of  $s$  but proportional to the contour length of the DNA. The ratio of the superhelix axis length to the DNA length is 0.41. The number of branch points is independent of  $s$  but approximately proportional to the DNA length. The number of supercoils is proportional to  $DLk$ , and  $n = 0.89vbm0DLk vbm0$ . The ratio between the twist increment and the writhe is constant, and  $DTw/Wr \approx 3$ . Individually,  $DTw = 0.28DLk$  and  $Wr = 0.72DLk$ .

**Figure 2.** Geometry of an interwound (plectonemic) DNA superhelix. **(a)** In this schematic model the DNA axis winds up the surface of a virtual cylinder, crosses the top in a hemispheric circle, winds back down the cylinder with the same handedness, and finally completes the path as a second hemispheric circle at the bottom. The height of the capped cylinder is  $l$ , the radius is  $r$ , and the height of the cylindrical portion is  $l-2r$ . In this model, the total number of superhelical turns,  $n$ , is 4, with 2 winding up and 2 winding down the surface of the cylinder. The pitch,  $2pp$ , is shown for a representative superhelical turn. **(b)** Here the surface has been cut open and unfolded into a plane. The DNA axis forms a diagonal across the resulting rectangle. The height of the rectangle is the superhelical pitch  $2pp$ , and the base of the rectangle is the circumference of the cylinder,  $2\pi r$ . The superhelix winding angle,  $\alpha$ , is given by  $\sin(\alpha) = p/\sqrt{r^2 + p^2}$ .



The salt dependence of the superhelix radius is uncertain and somewhat controversial. It has been reported, based upon experiments employing [cryoelectron microscopy](#), that the superhelix radius is highly salt-dependent, with the opposing duplex segments approaching self-contact at moderate counterion concentrations ([5](#), [6](#)). On the other hand, the results of static and dynamic [light scattering](#), time-resolved fluorescence polarization anisotropy of intercalated [ethidium bromide](#), and [circular dichroism](#) indicate that the structure remains much more relaxed at high salt, with no evidence of lateral intraduplex contacts ([7](#)). Yet another study using [atomic force microscopy](#) found somewhat intermediate results ([8](#)). Superhelical DNA was found to form an increasingly tight plectonemic superhelix with increasing salt concentration, including regions of close DNA–DNA contacts. The structures observed were, however, relatively looser than those observed with cryoelectron microscopy.

#### Bibliography

1. W. R. Bauer, F. H. C. Crick, and J. H. White (1980) *Sci. Am.* **243**, 118–122.
2. P. J. Hagerman (1988) *Ann. Rev. Biophys. Biophys. Chem.* **17**, 265–286.
3. V. V. Rybenkov, N. R. Cozzarelli, and A. V. Vologodskii (1993) *Proc. Natl. Acad. Sci. USA* **90**, 5307–5311.
4. T. C. Boles, J. H. White, and N. R. Cozzarelli (1990) *J. Mol. Biol.* **213**, 931–951.



5. M. Adrian, B. ten Heggeler-Bordier, W. Wahli, A. Z. Stasiak, A. Stasiak, and J. Dubochet (1990) *EMBO J.* **9**, 4551–4554.
6. J. Bednar, P. Furrer, A. Stasiak, J. Dubochet, E. H. Egelman, and A. D. Bates (1994) *J. Mol. Biol.* **235**, 825–847.
7. J. A. Gebe, J. J. Delrow, P. J. Heath, B. S. Fujimoto, D. W. Stewart, and J. M. Schurr (1996) *J. Mol. Biol.* **262**, 105–128.
8. Y. L. Lyubchenko and L. S. Shlyakthenko (1997) *Proc. Natl. Acad. Sci. USA* **94**, 496–501.

### Suggestion for Further Reading

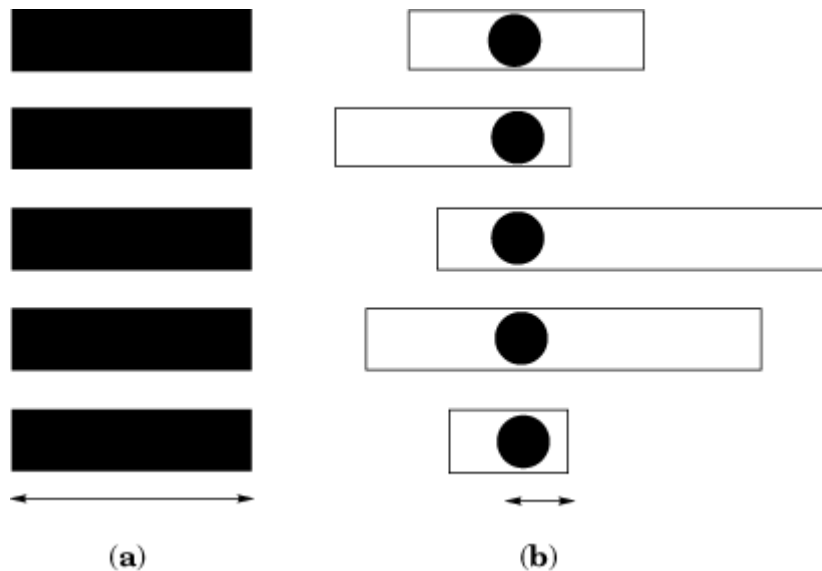
9. N. R. Cozzarelli, T. C. Boles, and J. H. White (1990) Primer on the topology and geometry of DNA supercoiling, in *DNA Topology and Its Biological Effects* N. R. Cozzarelli and J. C. Wang, eds., Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 139–184. (This review article contains a clear and detailed description of all aspects of the geometry of superhelical DNA.)

## Superfamily

A superfamily is a set of several that have similar amino acid sequences or **tertiary structures** of the they encode, even though each member of the family has a different function. For example, **immunoglobulins** and **T-cell receptors** are direct counterparts, each of which is produced by its own type of lymphocyte, **B cells** or **T cells**. Their **protein structures** are related to each other, and their gene organizations are also similar. Moreover, the mechanisms of generating their variability also share common features. By including the **major histocompatibility complex** (MHC) proteins, which also have common features in various aspects, the immunoglobulins, T-cell receptor, and MHC constitute a superfamily (1). Indeed, it is quite possible that this superfamily originated during from a common ancestor.

Superfamily is a classification of gene families on the basis of comparisons of the entire gene region as a unit. It has become apparent recently, however, that there are many kinds of **mosaic proteins**, which contain more than one different functional domain (see also ). It follows that many genes may share only a part of their gene regions, and that their entire gene regions are not necessarily similar. This suggests that a superfamily, on the basis of the entire gene regions, may not be sufficient to reflect the variety of protein-coding genes (Fig. 1). Thus, a superfamily must be redefined on the basis of their common domains.

**Figure 1.** Superfamilies defined on the basis of their entire gene regions (a) or with only a common domain being shared among different mosaic proteins (b). The solid boxes and circles indicate the homologous regions.



## Bibliography

“Superfamily” in , Vol. 4, p. 2468, by T. Gojobori; “Superfamily” in (online), posting date: January 15, 2002, by T. Gojobori, National Institute of Genetics, Shizuoka, Japan.

1. B. Lewin, (1994) *Genes V*, Oxford Univ. Press, Oxford.

## Superhelical DNA Energetics

### 1. Topology-Dependent Energetics

The constancy of the **linking number** for a [DNA topology](#) domain places an additional constraint upon all chemical reactions and other changes of state of the [DNA structure](#). This is reflected in significantly altered thermodynamic quantities, thermal **denaturation** patterns, transitions to non-**B DNA** duplex structures, and binding of various reagents. Generally speaking, two opposing tendencies must be balanced in analyzing superhelical DNA energetics. The first is minimization of the deviation of the **twist**,  $T_w$ , of closed duplex DNA (cdDNA ) from that of the corresponding nicked duplex DNA. The second is minimization of the bending distortion that accompanies superhelix formation. [DNA Ligase](#) The net free-energy change due to a change in linking difference,  $\Delta Lk$ , results from some combination of these two factors and is known as the free energy of superhelix formation. In terms of the superhelix density,  $s$  the free energy is

$$\Delta G(\sigma, T) = \frac{NRT}{h_0^2} q(T) \sigma^2 \quad (1)$$

where  $N$  is the number of base pairs in the DNA, and  $h_0$  is the helical repeat of nicked circular or linear duplex DNA, 10.5 bp/turn. The superhelix density is  $s = (Lk - Lk_0)/Lk_0$ , where  $Lk$  is the cdDNA linking number and  $Lk_0 = N/h_0$ . [Gel Electrophoresis](#)

#### 1.1. Experimental Determination of Free Energy

The free energy that appears in Equation (1) has been determined by three independent methods. The first depends on the perturbation of the binding constant of intercalating drugs with changes in  $s$  (1). This method determines  $DG$  over the entire accessible range of  $\pm s$ . The second method is based on the width of the distribution of topoisomers when a closed circular DNA is equilibrated in the presence of a **topoisomerase** (2-4). This method determines  $DG$  over a narrow range of  $\pm s$  in the vicinity of  $s = 0$ . The third method relates the extent of the linking difference to the local denaturation (5). This method determines  $DG$  over a wide range of values of  $s < 0$ . All methods give comparable results under similar environmental conditions.

The experimental results show that the free-energy coefficient  $q(T)$  in Equation (1) is independent of superhelix density,  $s$ , and of the length of the DNA,  $N$ , for  $N > 2400$ bp. For shorter DNAs,  $q$  increases as  $N$  decreases (6, 7) according to the relationship (8)

$$q(37^\circ, N) = 3939 - 1.1N \text{ for } N < 2400 \text{bp} \quad (2)$$

No data have been reported to date for the temperature dependence of  $q$  for DNAs in this size range. The value of  $q(T)$  for  $N > 2400$ bp is 1160 bp at 37°C and decreases with temperature (9, 10) according to the relationship

$$q(T) = \frac{0.968 \times 10^6}{T} - 1.92 \times 10^3 \text{bp for } N > 2400 \text{bp} \quad (3)$$

Both the entropy and enthalpy of superhelix formation are quadratic functions of the superhelix density but are independent of temperature. Numerical values have been determined for pBR322 (4, 9), pSM1 (10), and ColE1 *amp* plasmid (11, 12). Experiments with pSM1 ( $N = 5804$ bp) and pBR322 ( $N = 4363$ bp) DNAs, indicate that the enthalpy change per mole base pair is  $DH/N = 17.3s^2$  kcal/mol, and the entropy change per mole base pair is  $DS/N = 35$  cal/mol/deg. At 37°C, the free-energy change per mole base pair is  $DG/N = 6.4s^2$  kcal/mol.

## 1.2. Boltzmann Distribution of Topoisomers

Perhaps the clearest manifestation of the free energy of superhelix formation is in the appearance of the thermal distribution of topoisomers that results from equilibration in the presence of a DNA topoisomerase (2, 3). In this type of experiment, a cdDNA is incubated with topoisomerase (or the appropriate nicked circular DNA is incubated in the presence of [Linking Number Of DNA](#)), and thermal equilibrium established. The enzyme is then inactivated, and the resulting Boltzmann distribution of topoisomers is fractionated by [gel electrophoresis](#), as shown in Figure 1. The Boltzmann distribution is determined by the free energy difference between topoisomer  $i$  and the hypothetical completely relaxed species at the center of the distribution, with  $DLk = \epsilon$ . The distribution is then

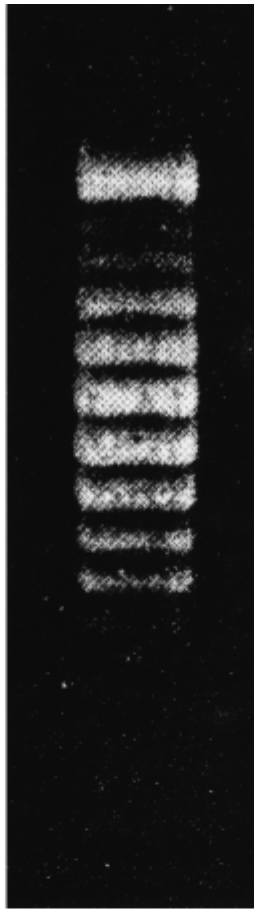
$$c_i = \left( \frac{1}{Q} \right) \exp \left[ \frac{\Delta G(i + \epsilon)}{RT} \right] \quad (4)$$

where the factor  $Q$  is the partition function, here simply a constant. Incorporating Equation (1), using the definition of  $s$  and simplifying, the Boltzmann distribution has the form

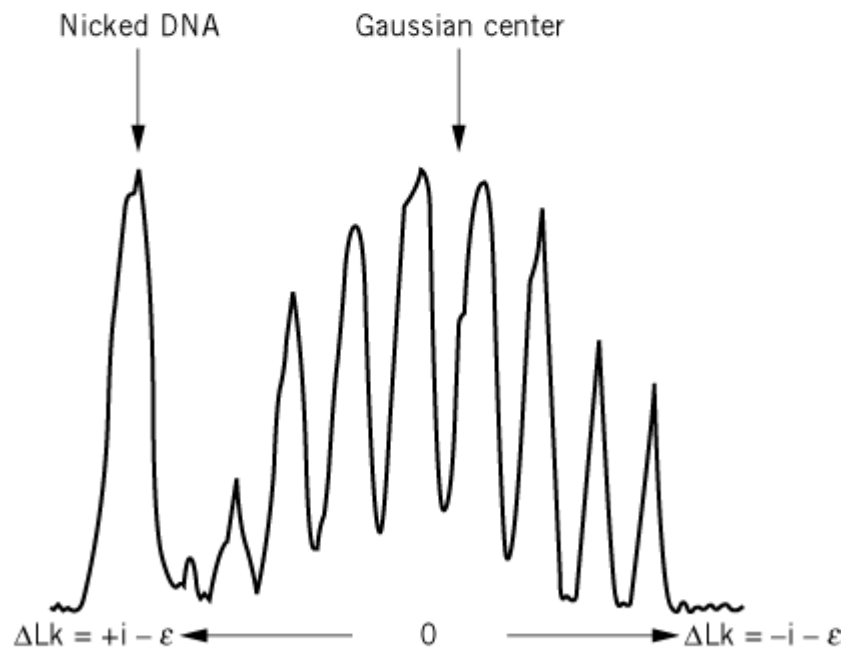
$$c_i = \left( \frac{1}{Q} \right) \exp \left[ \left( \frac{q}{N} (i + \epsilon)^2 \right) \right] \quad (5)$$

**Figure 1.** Experimental measurement of DNA superhelical energetics. (a) Fractionation of a distribution of topoisomers by gel electrophoresis. The topoisomers are offset from the nicked circular DNA because of a difference

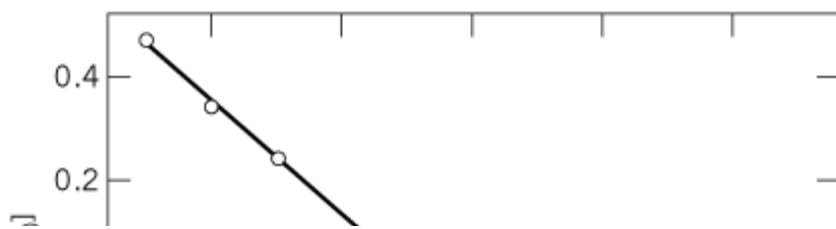
in temperature between the topoisomeric equilibration conditions and the gel electrophoresis conditions. (Modified with permission from Ref. [2](#); copyright 1975, National Academy of Sciences.) **(b)** (See next page) Densitometric tracing of the same gel, showing the intensity profile as a function of distance. The location of the nicked DNA is indicated, as is the approximate location of the Gaussian center of the distribution. **(c)** Analytical plot of the topoisomer distribution. The  $D Lk$  values of the topoisomers are  $i + \epsilon$ , where  $i = 0, \pm 1, \pm 2$ , and so on.



(a)



(b)



For analytical purposes, it is convenient to recast the result in Equation (5) as a ratio of the concentration of topoisomer  $i$  relative to that of the most prominent topoisomer, taken to lie at  $i = 0$ . Writing Equation (5) for each and taking the ratio, we obtain

$$\frac{c_i}{c_{i=0}} = \exp \left[ \left( \frac{q}{N} \right) i(i + 2\varepsilon) \right] \quad (6)$$

If the relative concentrations of the various topoisomers are plotted on a linear scale, the resulting distribution is the expected Gaussian distribution. The separation between the maximum of this distribution and the location of the nearest topoisomer is equal to  $\varepsilon$ , the difference in winding between a nicked circular DNA and the nearest relaxed, closed DNA (see [Linking Number Of DNA](#)). Taking the logarithms of both sides of Equation (6) gives a straight line representation that can be applied to the data

$$\left( \frac{1}{i} \right) \ln \left( \frac{c_i}{c_{\max}} \right) = \left( \frac{q}{N} \right) (i + 2\varepsilon) \quad (7)$$

This is the equation that is plotted in Figure 1(c), with the slope giving  $q/N$  and the intercept at  $i = 0$  yielding  $\varepsilon$ .

## 2. Topology-Dependent Chemical Reactivity

The free energy of superhelix formation enters into all chemical reactions that involve a change in the average duplex rotation. The classical example of such a reaction is the intercalation of [ethidium bromide](#), which unwinds the DNA duplex by  $26^\circ$ , or 0.072 duplex turn. In general, the free energy for the binding of  $n$  ligands per base pair to a closed circular DNA,  $DG(n,s)$ , is the sum of two contributions: the intrinsic binding free energy,  $DG^\circ(n)$ , and the associated change in the free energy of superhelix formation,  $DDG(s)$ . The binding constants are related to the intrinsic and total binding free energies by the usual relationship  $RT \ln(K) = -DG$  and  $RT \ln(K^\circ) = -DG^\circ$ , where  $K^\circ$  is the intrinsic binding constant for a nicked circular or linear duplex DNA. The change in the superhelix free-energy term  $DG(s)$  is  $[dD G(s)/ds]Ds$ . If a bound ligand alters the duplex rotation by  $f$  turns, then  $Ds = h_0 f/N$ . Combining these terms, the binding constant is dependent on the superhelix density and temperature according to

$$K(\sigma, T) = K^\circ(T) \exp \left[ -\frac{2f}{h_0} q(T) \sigma^2 \right] \quad (8)$$

The effect on  $K$  can be significant under biologically relevant conditions. For example, if a protein binds to a DNA of  $s = -0.07$  and removes one duplex turn in the process,  $K(s = -0.07, 37^\circ\text{C}) = 3K^\circ$ . This means that the avidity of such a closed circular DNA for the protein binding is approximately three times that of a comparable nicked circular or linear DNA under the same environmental conditions. This effect can be even more pronounced for a **supercoiled** DNA that has been prepared with a very high (negative) superhelix density. Thus, if  $s = -0.17$ , with one duplex turn removed in the binding process,  $K(s = -0.17, 37^\circ\text{C}) = 600K^\circ$ . Consequently, if two DNAs of low and high (negative) superhelix density compete for binding of the same protein, the closed DNA of high  $|s|$  has a significant advantage. In the event that  $s > 0$ , of course, the effect is reversed and  $K(s > 0) < K^\circ$ .

## Bibliography

1. W. Bauer and J. Vinograd (1970) *J. Mol. Biol.* **47**, 419–435.
2. R. E. Depew and J. C. Wang (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4275–4279.
3. D. E. Pulleyblank, M. Shure, D. Tang, J. Vinograd, and H. P. Vosberg (1975) *Proc. Natl. Acad. Sci. USA* **72**, 4280–4284.
4. M. Duguet (1993) *Nucleic Acids Res.* **21**, 463–468.
5. C. J. Benham (1992) *J. Mol. Biol.* **225**, 835–847.
6. D. Shore and R. L. Baldwin (1983) *J. Mol. Biol.* **170**, 983–1007.
7. D. S. Horowitz and J. C. Wang (1984) *J. Mol. Biol.* **173**, 75–91.
8. W. R. Bauer and R. Gallo (1989) in *Chromosomes: Eukaryotic, Prokaryotic, and Viral*, K. W. Adolph, ed., CRC Press, Boca Raton, FL, Vol. I, pp. 87–126.
9. W. R. Bauer and C. J. Benham (1993) *J. Mol. Biol.* **234**, 1184–1196.
10. W. R. Bauer, H. Ohtsubo, E. Ohtsubo, and C. J. Benham (1995) *J. Mol. Biol.* **253**, 438–452.
11. A. Seidl and H.-J. Hinz (1984) *Proc. Natl. Acad. Sci. USA* **81**, 1312–1316.
12. C. H. Lee, H. Mizusawa, and T. Kakefuda (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2838–2842.

## Suggestion for Further Reading

13. C. J. Benham (1990) "Free energies associated with superhelical deformations of DNA, in" *Structure and Methods: DNA and RNA*, R. H. Sarma and M. H. Sarma, eds., Adenine Press, New York, Vol. 3, pp. 255–266.

## Supersecondary Structure

Supersecondary structure is an alternative term for [protein motif](#). The term is used to describe preferred packing arrangements of combinations of **secondary structure** elements in **protein structures**.

[See [Protein Structure](#).]

## Suppressor Mutation

Suppressor mutations relieve the phenotypic effects of preexisting mutations. These “second-site” mutations have been extremely useful in genetic and biochemical analyses. The suppressor mutation identifies a new site affecting the phenotype of interest, and suppression is strong evidence that the two genetic sites functionally interact. Suppressor mutations can occur in the same gene as the initial mutation (intragenic suppression), or the two mutations may be in separate genes (intergenic suppression). Intragenic suppression may occur if the amino acids specified by the two sites must interact to give the normal function of the gene product. Intragenic suppression may also occur by the relief of a polar effect on expression. Intergenic suppression has been extensively used to characterize interactions between gene products. This analysis has been used to identify proteins that

interact in various biochemical, gene expression, and signal transduction pathways. A special kind of intergenic suppression occurs when mutations occur in transfer RNAs so that they translates stop (nonsense) codons as sense codons. These tRNAs (nonsense suppressors) can allow expression of genes that have suffered mutations that create stop codons within coding sequences. Another special kind of suppression occurs when mutants are “rescued” by normal genes expressed from high-copy plasmids. See [Nonsense Suppression](#) for more complete discussions and references.

## Suppressor tRNA

Suppressor transfer RNAs (tRNAs) are mutant tRNAs that insert a suitable amino acid at a mutant site in protein encoding genes. The term “suppressor” is used because the mutant tRNA “suppresses” the phenotypic effect of the coding mutation. Most suppressor tRNAs contain a mutation in either the anticodon, changing codon specificity, or at some position that alters the aminoacylation identity of the tRNA. Suppressor tRNAs have been isolated that decode each of the three termination (nonsense) codons. Suppressors of missense and frameshift mutations are also known. Suppressor tRNAs have been used extensively to study tRNA functions. Suppressors have also been used to insert specific amino acids at particular positions within proteins to study effects of specific amino acid substitutions. For other discussions of suppressor tRNAs, see [Nonsense Suppression](#).

Because suppressor tRNAs can potentially affect the production of a large fraction of the cellular protein, many suppressors have relatively low translational activities, so that normal translation predominates. Exceptions include the UAG termination codon suppressors (amber suppressors) in *Escherichia coli*. The UAG termination codon is rare in *E. coli* and, apparently, frequent translation as a sense does not severely affect this organism. Because the normal translational functions of tRNAs are essential, most suppressor tRNA mutations occur in one of a family of duplicate tRNA genes, and the nonmutated copies continue to perform normally. Suppressor alleles for single copy tRNA genes can be generated using cloned copies of the gene.

A great deal of information about tRNA function has been gained by measuring the efficiency with which suppressor tRNAs decode. Suppressor tRNA genes are easy to manipulate, and suppression efficiencies are readily quantified; together, these properties make suppression the assay of choice for probing translational mechanisms. Suppressor tRNAs must compete with normal molecules for decoding. Nonsense suppressors, for example, must compete with the peptide release factors that normally decode nonsense codons. Suppressors that are good competitors allow for efficient suppression of the nonsense mutation. To study tRNA structure and function, many workers have created additional mutations in nonsense suppressors to determine how those secondary mutations affect suppression efficiency. In an extensive mutational analysis of the anticodon arm, for example, Yarus et al. (1) showed that the native sequence was optimal for suppression, and that mutations reduce suppression efficiency by various degrees. It was concluded that tRNAs utilize an “extended anticodon” in which the nucleotides near the anticodon contribute to its ability to translate codons.

Murgola and co-workers have isolated a large number of missense suppressors in *E. coli* (2). These tRNAs insert suitable amino acids at missense mutations. Many of these mutant tRNAs have anticodon mutations that change codon specificity, changing a glycine tRNA so that it reads the tryptophan codon, for example. Other suppressors are aminoacylated with a different amino acid. These tRNAs may thus insert a suitable amino acid at missense sites that contain their cognate codons. Missense suppressors are generally inefficient, probably because efficient translation of codons with a different meaning is toxic.

Frameshift suppressor tRNAs are also known (3, 4). In the classic examples the tRNAs have an extra



base inserted into the anticodon loop such that they may decode with a four-base anticodon. Certain of these tRNAs may in fact decode this way. For example, tRNA<sup>Trp</sup> derivatives that have four-base anticodons cognate to UAGU, UAGC, or UAGU sites decode those sites as four-base codons. Interestingly, they can also decode these sites as three-base codons. Together, these data have been used to argue that tRNA anticodon loop conformation determines the number of bases decoded (5). Not all frameshift suppressors may act in this way. There is strong evidence that some frameshift suppressor tRNAs that have an extra large anticodon loop probably do not decode at all. The insertion of the extra base into the “suppressor” tRNA may prevent it from functioning and thus “starve” its codon. Then, frameshifting may occur when the starved codon is read by a normal but noncognate tRNA such that the reading frame is somehow perturbed, thus linking misreading and frameshifting (6). There is at least one example in which a frameshift suppressor with a normal sized anticodon loop directly causes frameshifting by misreading (7). Frameshift suppression by mutations in tRNA nucleoside modification enzymes may function in a similar way. The undermodified tRNAs may act poorly in translation, thus starving their codons and allowing for misreading/frameshifting events (6). These studies have been instrumental in our understanding of both reading frame maintenance and translational accuracy.

McClain and co-workers (8) have used suppressor tRNAs to demonstrate changes in the aminoacylation specificities of several tRNAs and thus map the determinants for tRNA amino acid “identity.” They engineered an amber termination site near the 5' end of the gene for the easily isolable dihydrofolate reductase enzyme. They then mutagenized tRNA suppressors and determined which amino acids were inserted by the mutant tRNAs by sequencing the amino termini of isolated enzymes. They found that a small number of nucleotides are mostly responsible for defining tRNA identity. The use of suppressor tRNAs was critical to this work because there was no ambiguity about which tRNA actually decodes the amber site.

Studies of the roles of nucleoside modifications are greatly facilitated by the use of suppressor tRNAs. Modifications are ubiquitous in RNA, and work thus far shows that they are important for the functions of tRNAs (9, 10). Comparisons of the activities of suppressor tRNAs that either do or do not contain certain modifications have shown, for example, that modifications near the anticodon increase translational efficiency. Mechanisms are not always clear, but a common theme is that bulky modification 3' to the anticodon increases translational efficiency, as if increased base stacking stabilizes anticodon:codon complexes.

Suppressor tRNAs are also used to show that codon translation is affected by neighboring nucleotides (codon context). Many studies show that nonsense codons are more readily suppressed if the 3' neighbor nucleotide is a purine. This context effect is likely to have two sources, base stacking with the 3' purine stabilizes the anticodon:codon complex (11, 12), and the release factors are highly dependent on the neighbor for termination and a 3' U appears to be optimal (13, 14).

Transfer RNA Suppressors will also be useful for protein engineering. In one approach, Abelson and collaborators have constructed an extensive set of tRNA suppressors that can be used to insert a wide range of amino acids at amber mutations (15). These tRNAs may allow systematic tests of the effects of various amino acids on protein structure and function. Transfer RNA suppression systems for labeling proteins with novel amino acids at specific positions are also in development. Amber suppressor tRNAs are being aminoacylated with nonstandard amino acids. Then these tRNAs direct the incorporation of the novel amino acid into a nascent protein whose gene has been mutated to contain an amber codon at a specific site within its coding sequence. At this time, only a few nonstandard amino acids may be incorporated in this way, but it is anticipated that it will soon be possible to label proteins with a large variety of amino acids containing affinity tags, fluors, or reactive groups (16).

## Bibliography

1. M. Yarus, S. W. Cline, P. Weir, L. Breeden, and R. C. Thompson (1986) *J. Mol. Biol.* **192**,

235–255.

2. E. J. Murgola (1994) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, D. C., pp. 491–509.
3. M. R. Culbertson, P. Leeds, M. G. Sandbaken, and P. G. Wilson (1990) In *The Ribosome: Structure, Function, & Evolution* (W. E. Hill, A. Dahlberg, R. A. Garrett, P. B. Moore, D. Schlessinger, and J. R. Warner, eds.), American Society for Microbiology Press, Washington, D. C., pp. 559–570.
4. J. F. Atkins, R. B. Weiss, S. Thompson, and R. F. Gesteland (1991) *Annu. Rev. Genet.* **25**, 210–228 (1991).
5. J. F. Curran and M. Yarus (1987) *Science* **238**, 1545–1550.
6. Q. Qian, J.-N. Li, T. G. Hagervall, P. J. Farabaugh and G. R. Björk (1998) *Mol. Cell* **1**, 471–482.
7. F. T. Pagel, T. M. F. Tuohy, J. F. Atkins, and E. J. Murgola (1992) *J. Bacteriol.* **174**, 4179–4182.
8. W. H. McClain (1994) In *tRNA: Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), American Society for Microbiology Press, Washington, D. C., pp. 335–347.
9. J. F. Curran (1998) In *Modification and Editing of RNA* (H. Grosjean and R. Benne, eds.), American Society for Microbiology Press, Washington, D. C., pp. 493–516.
10. G. R. Björk (1996) In *Escherichia coli and Salmonella: Cellular and Molecular Biology* (F. C. Neidhardt, R. Curtis III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology Press, Washington, D. C., pp. 861–886.
11. M. Yarus and J. F. Curran (1992) In *Transfer RNA in Protein Synthesis* (D. L. Hatfield, B. J. Lee and R. M. Pirtle, eds.), CRC Press, Boca Raton, FL, pp. 319–365.
12. R. H. Buckingham (1994) *Biochimie* **76**, 351–354.
13. W. T. Pedersen and J. F. Curran (1991) *J. Mol. Biol.* **219**, 231–241.
14. E. S. Poole, C. M. Brown, and W. P. Tate (1995) *EMBO J.* **14**, 151–158.
15. L. G. Kleina, J.-M. Masson, J. Normanly, J. Abelson, and J. H. Miller (1990) *J. Mol. Biol.* **213**, 704–717.
16. L. Wang, A. Brock, B. Heberich, and P. G. Schultz (2001) *Science* **292**, 498–500.

## Surface Wrapping Of DNA

Closed duplex DNA (cdDNA) is often found wrapped on a protein surface, either entirely or partially. Such wrapping provides a frame of reference in which to analyze the topological and geometric properties of the DNA (see [DNA Topology](#)), thereby greatly simplifying many of the associated relationships. For this purpose, both the protein surface and the DNA are modeled as continuous (smooth) objects. The protein surface is defined by moving a **water**-sized spherical probe around the van der Waals surface of each external atom. The solvent-accessible surface is the continuous sheet defined by the locus of the center of the probe (**1**) (see [Accessible Surface](#)). The surface on which the DNA axis lies is then defined by extending outward by 1 nm from the protein surface in order to account for the thickness of the [DNA structure](#).

Even more generally, a closed duplex DNA not actually wrapped on a protein often behaves in some respects as if it were bound to a surface of a particular shape. Some of the general properties of DNA

wrapped on a physical surface apply to wrapping on such a virtual surface as well. The relationships appropriate to surface wrapping begin with the fundamental relationship for a topological domain (see [Linking Number Of DNA](#))

$$Lk = Tw + Wr \quad (1)$$

which relates the two geometric quantities: **twist**,  $Tw$ , and **writhe**,  $Wr$ , to the linking number,  $Lk$ . This equation applies to all closed duplex DNAs, regardless of whether they are wrapped on a surface. Unfortunately, except in a few special cases, neither  $Tw$  nor  $Wr$  is directly accessible experimentally or readily interpretable in structural terms. Equation (1) is analogous to the combined first and second laws of thermodynamics, in that  $Lk$  is a state function, while both  $Tw$  and  $Wr$  are dependent on the path of any deformation.

## 1. Surface Linking Number

The presence of a surface reference frame permits the reformulation of the above fundamental relationship into a much more tractable form. For DNA associated with a surface, the twist divides into two parts, one determined by the contour of the surface at the loci of attachment of the DNA axis ( $STw$ ), and the other by the local winding of the DNA about its axis ( $F$ ). The total twist is then the sum of two independent contributions (see [Twist, DNA](#))

$$Tw = STw + \Phi \quad (2)$$

Both  $STw$  and  $F$  are experimentally accessible. This expression for the twist is combined with the fundamental relationship, Equation (1), to obtain

$$Lk = STw + Wr + \Phi \quad (3)$$

Equation (3) now contains two quantities that are determined by the shape of the surface/axis ( $STw$  and  $Wr$ ) and one that is intrinsic to the DNA ( $F$ ).

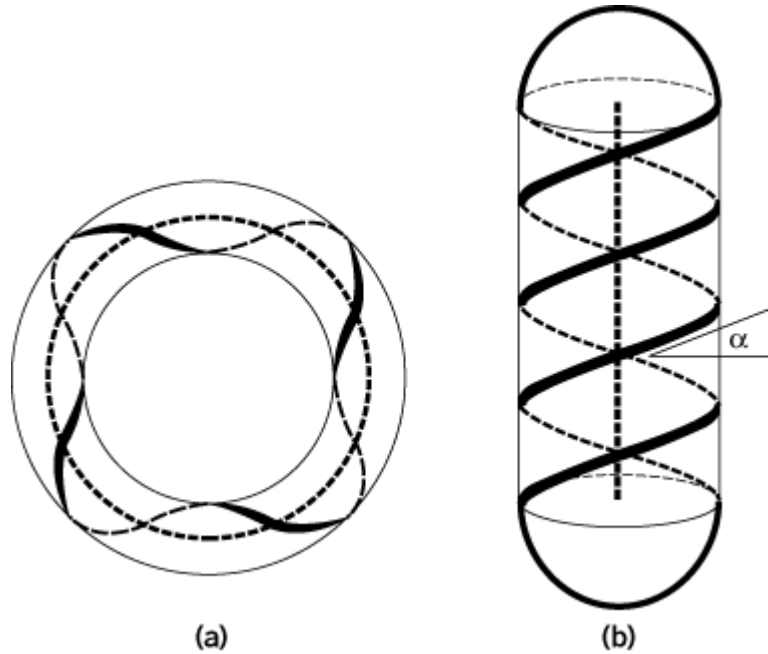
The surface-dependent quantities are combined into the surface linking number,  $SLk$  (2):

$$SLk = STw + Wr \quad (4)$$

This equation is another fundamental relationship, analogous to the more general Equation (1), which applies to surface-wrapped closed duplex DNA. The surface linking number, like the linking number, is an integer for a superhelical DNA. The surface linking number is a topological quantity and is therefore unaltered if the surface is distorted smoothly. In addition, for any surface distortion that does not break the surface near the DNA axis, any change in  $STw$  is offset by an equal and opposite change in  $Wr$ . For example, a plectonemic superhelix can undergo very large distortions, including the formation of one or more branches, without transformation to a toroidal superhelix. In all such distortions,  $SLk$  is unchanged; hence changes in surface twist are equal in magnitude and opposite in sign to changes in writhe. Since calculation of  $STw$  is usually much easier than direct calculation of  $Wr$ , the latter is usually unnecessary in this very common case (see [Writhe, DNA](#)).

The most important types of surfaces from the practical view are the spheroid and the circular torus. The spheroid applies to DNA wrapped entirely on the surface of a globular protein; it also applies, as a virtual surface, to all interwound cdDNAs (3). The circular torus applies to DNA wrapped on a cylinder, such as a [nucleosome](#) core. For all spheroidal surfaces,  $SLk = 0$ . For wrapping  $n$  times around a toroid,  $SLk = +n$  if the wrapping is right-handed and  $SLk = -n$  if the wrapping is left-handed. Examples of the two types of surface wrapping for DNA are illustrated in Figure 1.

**Figure 1.** Wrapping of a closed circular DNA onto two different types of surfaces. (a) The DNA is wrapped in a right-handed sense on the surface of a spherical torus, where the number of superhelical turns  $n = +4$ . (b) The DNA is wrapped on the surface of a capped cylinder. Here the wrapping is also right-handed, and  $n = +2$  up the cylinder and  $+2$  down, for a total of  $+4$ . The DNA axis is connected at the tops by a perfect half circle.



Combining Equations (3) and (4), an alternative expression for the linking number is obtained:

$$Lk = SLk + \Phi \quad (5)$$

The linking number of superhelical DNA wrapped on a protein surface or a well-defined virtual surface is thus the sum of two integers, each of which can be independently and experimentally determined. This expression for the linking number in Equation (5) is much more straightforward than that of Equation (1). Here the component expressing the axis trajectory,  $SLk$ , is rigorously separated from the winding of either strand about the axis,  $F$ . For the common case of a DNA wrapped entirely as plectonemic superhelix, Equation (5) simplifies even more to  $Lk = f$ . For a DNA wrapped entirely as a toroidal superhelix,  $Lk = \pm n + f$ .  $Lk$ ,  $SLk$ , and  $F$  have been independently measured (4) and found to be in agreement with Equation (5).

## 2. DNA Helical Repeat

The second independent quantity in Equation (5) is  $F$ , the winding number of the DNA axis about the normal to the protein surface.  $F$  is  $\pm \frac{1}{2}$  the number of times the backbone curve  $C$  intersects the surface as measured by chemical or **nuclease** probes. The helical repeat ( $h$ ) of surface wrapped DNA is directly calculated from  $F$  and the number of base pairs ( $N$ ) by

$$h = \frac{N}{\Phi} \quad (6)$$

The related quantity,  $N/Tw$ , differs from  $h$ , usually by a few percent, depending on the details of the surface geometry (3). Specifically

$$\frac{N}{Tw} = \frac{h}{1 + (STw/\Phi)} \quad (7)$$

For a regular, right-handed plectonemic superhelix,  $STw = \pm n \sin(\alpha)$ , where  $\alpha$  is the superhelix pitch angle. Then for a plectonemic superhelix, Equation (7) simplifies to

$$\frac{N}{Tw} = \frac{h}{1 + (n/N)h\sin(\alpha)} \quad (8)$$

where  $n$  is the number of plectonemic superhelical turns. These geometric quantities are specified in Figure 1.

For the case of a nicked circular DNA, or of a closed DNA with  $DLk = 0$ , the helical repeat is that of DNA with a linear axis,  $h_0 = N/f_0$ , where  $f_0 = Lk_0$ , the pseudo-linking number of open circular DNA [see [Linking Number Of DNA](#)]. These considerations lead to a general expression for the helical repeat of a closed circular DNA. Equations (5) and (6) are combined, along with the definition of the superhelix density,  $s = (Lk - Lk_0)/Lk_0$ , to yield

$$h = \frac{h_0}{\sigma - (SLk/\Phi_0) + 1} = \frac{h_0}{(1 + \sigma) - (SLk/N)h_0} \quad (9)$$

### 2.1. Interwound Superhelix

The result here applies to purified, noncatenated DNA. Since  $SLk = 0$  for any plectonemic superhelix, the helical repeat in this case is a simple function of the linking number or superhelix density. Equation (8) then simplifies to

$$h = \frac{N}{\Phi} = \frac{h_0}{1 + \sigma} \quad (10)$$

The relationship in this equation has been confirmed for several cDNAs. For a typical naturally occurring cDNA, with  $s = -0.06$ , the helical repeat is increased from the linear axis value of 10.5 bp/turn to the significantly higher value of 11.2 bp/turn. Superhelix densities as great as  $-0.17$  have been reported, corresponding to  $h = 13.1$  bp/turn (5). Positively supercoiled DNA, having  $s > 0$ , can be generated with the DNA reverse gyrase (6) (see **Topoisomerases**). For example, if  $s = +0.10$ ,  $h = 9.55$  bp/turn.

### 2.2. Toroidal Superhelix

The result here applies to DNA wrapped on a protein surface and to the closed circular submolecules of catenated DNA. In the latter case, either submolecule of the catenated pair can be taken as defining a virtual toroidal surface about which the other winds. (The former must have a circular axis but need not contain its own topological domain.) Now  $SLk = \pm n$  for any toroidal superhelix, with  $n > 0$  for right-handed winding and  $n < 0$  for left-handed winding. The helical repeat for a DNA topological domain wrapped on either type of surface is then

$$\frac{1}{h} = \frac{1}{h_0}(1 + \sigma) \pm \frac{n}{N} \quad (11)$$

For example, the variation of  $n$  with  $s$  has been determined for catenanes of DNAs of various lengths (7). In one such case, identical submolecules of length 3.5 kbp were wound right-handed  $n$  times, with the average value of  $n$  varying from 2.5 to 5.5. It was found experimentally that  $s = (0.32n -$

0.30) $Lk_0$ . Applying Equation (9), we see that the helical repeat in this catenane is increased from 10.50 bp/turn ( $n = 0$ ) to 10.56 bp/turn ( $n = 2.5$ ) to 10.63 bp/turn ( $n = 5.5$ ). Another example is the **SV40** minichromosome (8). Here the DNA is wrapped left-handed 1.8 times about each of 26 **nucleosomes**, hence  $n = SLk = -46.8$ . (Since  $SLk$  must be an integer for the DNA as a whole, an extra contribution of  $-0.2$  comes from the closure of the axis around a virtual connecting torus, giving  $SLk = -47$  in total.) Applying Equation (9), we note that the helical repeat of DNA on the octamers is reduced from 10.50 to 10.17 bp/turn, a result consistent with nuclease digestion results (9).

### Bibliography

1. F. M. Richards (1977) *Annu. Rev. Biophys. Bioeng.* **6**, 151–176.
2. J. H. White, N. R. Cozzarelli, and W. R. Bauer (1988) *Science* **241**, 323–327.
3. W. R. Bauer and J. H. White (1990) in *Nucleic Acids and Molecular Biology*, F. Epstein and D. M. J. Lilley, eds., Springer-Verlag, Berlin Heidelberg, Vol. **4**, pp. 39–54.
4. W. R. Bauer, J. J. Hayes, J. H. White, and A. P. Wolffe (1994) *J. Mol. Biol.* **236**, 685–690.
5. C. K. Singleton and R. D. Wells (1982) *Anal. Biochem.* **122**, 253–257.
6. A. Kikuchi and K. Asai (1984) *Nature* **309**, 677–681.
7. S. A. Wasserman, J. H. White, and N. R. Cozzarelli (1988) *Nature* **334**, 448–450.
8. J. H. White and W. R. Bauer (1989) *Cell* **58**, 9–10.
9. H. R. Drew and A. A. Travers (1985) *J. Mol. Biol.* **186**, 773–790.

### Suggestions for Further Reading

10. J. H. White, R. M. Gallo, and W. R. Bauer (1992). Closed circular DNA as a probe for protein-induced structural changes, *Trends Biochem. Sci.* **17**, 7–12. (This article presents practical examples of how to apply the surface wrapping quantities to biological situations.)
11. D. R. Bates and A. Maxwell (1993) *DNA Topology*, Oxford Univ. Press, London. (Chapter 3 of this monograph contains an excellent and highly readable discussion of DNA surface wrapping.)

### SV40 (Simian Virus 40)

The virion of the simian **virus** (SV40) has a diameter of about 45 nm and a **sedimentation coefficient** of 240 S (see [Papovavirus](#)). SV40 grows permissively and produces thousands of progeny viruses in monkey cells, its natural host. Cells from rodent and other mammalian origins are abortively infected by SV40. Some, but not all, of such nonpermissive primary-culture cells are immortalized. Some immortalized cells are released from the tightly controlled cellular growth regulatory system and become transformed (see [Neoplastic Transformation](#)).

The **genome** of SV40 is double-stranded circular DNA, 5243 bp in length, and encodes both early and late genes. Two proteins with a molecular weight of 90 kDa and 17 kDa, named large and small **T antigens**, respectively, are **translated** from the early **messenger RNAs** that are transcribed primarily before the onset of viral DNA replication. The large T antigen is localized exclusively in the **nucleus**. The small T antigen exists in all fractions as a soluble form. Viral capsid proteins—VP1, VP2, and VP3—are encoded from the late gene.

By binding to many cellular proteins that can regulate the [cell cycle](#) at the G1 phase, such as **p53**, **retinoblastoma** protein, p107, and p130, the large T antigen inactivates them and makes the infected cells enter S phase. Moreover, the large T antigen participates directly in viral [DNA replication](#). The large T antigen binds to the GGAGGC sequence at the [replication origin](#) and unwinds the adjacent region by its [ATPase](#) and [DNA helicase](#) activities (the preinitiation complex). By recruitment of **DNA polymerase**  $\alpha$ -[primase](#) and RPA (replication protein A), the preinitiation complex is converted to the initiation complex. The large T antigen functions as a DNA helicase to unwind the two strands of preforked parental DNA. The replicative chains are elongated by the function of DNA polymerases  $\delta$  and  $\epsilon$ , replication factor C, and **proliferating cell nuclear antigen** (PCNA). **Phosphorylation** of residue Thr124 of the large T antigen by cyclin-dependent kinase is required for viral DNA replication. At the late stage of DNA replication, the large T antigen autoregulates its own mRNA synthesis and transactivates the late-gene transcription by bridging [transcription factors](#) TEF-1 (transcription enhancer factor 1) and TBP (TATA box-binding protein), which are associated with the late **promoter**. Newly synthesized SV40 DNA and capsid proteins are assembled in the nucleus to produce virus particles. Eventually, the infected cells die and lyse, to release progeny viruses.

The small T antigen forms a complex with protein **phosphatase** 2A and inhibits its enzymatic activity. Thereby, the enhanced mitogen-activated protein (MAP) kinase pathway stimulates cell proliferation. This function is dispensable for viral DNA replication, however, because the mutant lacking small T antigen still produces progeny viruses at the same efficiency as does wild type.

In SV40-infected nonpermissive cells, viral DNA replication is not induced due to the absence of a permissive factor, which is considered to be the nature of DNA polymerase  $\alpha$ -primase complex, even though the large T antigen is fully expressed. Continuous expression of the large T antigen inactivates p53 and retinoblastoma family proteins and causes the infected cells to be immortalized and/or transformed by the deregulation of the cell growth control. Such transformed cells can grow in soft agarose media, but they hardly induce tumors in animals upon injection.

Human primary-culture cells expressing the large T antigen, prepared by DNA [transfection](#), also become transformed, but not immortalized. In addition to the deregulation of the functions controlled by p53 and the retinoblastoma family proteins, inactivation of the putative [senescence](#) genes, which seem to act against immortalization, is also required for the immortalization of human cells. However, the efficiency of spontaneous immortalization of T-antigen-expressing human cells is increased about  $10^5$ -fold over that of control cells not expressing T antigen. This is presumably caused by the increased **mutation** frequency in the senescence genes, due to the enhanced [recombination](#) and rearrangements of [chromosomes](#) induced by the large T antigen. Therefore, the large T antigen is often used to immortalize human primary-culture cells prepared from rare tissues or biopsy materials.

#### Suggestions for Further Reading

- C. N. Cole (1996) *Polyomavirinae: The Viruses and Their Replication*. In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott -Raven, Philadelphia, pp.1997–2025.
- E. Fanning and R. Knippers (1992) Structure and function of simian virus 40 large tumor antigen. *Annu. Rev. Biochem.* **61**, 55–85.
- J. J. Manfredi and C. Prives (1994) The transforming activity of simian virus 40 large T antigen. *Biochim. Biophys. Acta* **1198**, 65–83.
- E. Moran (1993) DNA tumor virus transforming proteins and the cell cycle. *Curr. Opin. Genet. Dev.* **3**, 63–70.

## Svedberg Unit(S)

The *Svedberg unit* is the unit in which **sedimentation coefficients**, or *s*-values, are normally expressed. It is a unit of time defined as  $1 \times 10^{-13}$  s. Its symbol is S. Thus  $1 \times 10^{-13}$  sec is equal to 1 Svedberg unit or 1S. See [Analytical Ultracentrifugation](#) and [Sedimentation Velocity Centrifugation](#).

## Switch Region

[Isotype](#) or [class switching](#) occurs during the primary [immune response](#) upon stimulation with **T-cell-dependent antigens**. It happens only after the [immunoglobulin \(Ig\) gene rearrangement](#) has taken place for the heavy chain, so that the V–D–J encoded region of the heavy chain will be conserved after switching. Because the light chain remains unmodified, this implies that isotype switching does not change the [antibody](#) specificity. Switching may take place between m and any other H-chain isotype, to change from [IgM](#) to another isotype, although the most frequent event is  $m \rightarrow g$  ( $\text{IgM} \rightarrow \text{bIgGb}$ ). The choice of the acceptor isotype is controlled by cytokines synthesized by T cells. The distinct  $T_{H1}$  and  $T_{H2}$  populations drive different switches, because they produce different cytokines. Thus in the mouse  $T_{H1}$  will induce B cells to switch to [IgG2a](#), whereas  $T_{H2}$  will drive a change to [IgG1](#) or [IgE](#). In humans, where the  $T_{H1}/T_{H2}$  dichotomy is not as clear, **interleukin-4** (IL-4), a typical  $T_{H2}$  cytokine in mice, also induces a  $m \rightarrow \epsilon$  switch.

Switch regions are located at the 5' position of each set of heavy-chain constant genes, except for the d gene, which accounts for the simultaneous expression of [IgM](#) and [IgD](#) on mature B cells. They are designated as  $S_m$ ,  $S_{g1}$  to  $S_{g4}$ ,  $S_{a1}$ ,  $S_{a2}$ , and  $S_{\epsilon}$ , according to the flanking isotype. Their length varies between 2 and 10 kbp, depending upon the isotype. They contain highly repetitive sequences in which GAGCT and GGGT motifs are frequently present. In the mouse, for example, the  $S_g$  regions contain repeated sequences of 49 nucleotides. These regions provide the basis for the switch mechanism, because this mechanism occurs by nonhomologous [recombination](#). The enzymatic process, which is not dependent on [recombinase](#) activating genes, has not been elucidated thus far. Class switching most usually involves deletion of the DNA section located between the two switch regions. This piece of DNA is circularized, and such DNA circles have indeed been isolated and sequenced. Switching may also occur by inversion, which leaves the possibility of having revertants, although this appears rather exceptional.

Switch regions are preceded 5' by a **promoter**-like region that is involved in synthesis of I transcripts, which initiate just before the switch occurs. For example,  $I_m \rightarrow I_{g1}$  transcripts will be produced before  $m \rightarrow g1$  switching. These transcripts contain the I region, plus the constant region of the 3' isotype. Many speculations have been made regarding a possible role for these transcripts in switch control, including the possible formation of a transient triplex. To date, no convincing data have supported any hypothesis. These transcripts may simply be a passive witness of DNA accessibility that would condition binding of the still elusive enzymatic complex. Isotype switching is controlled by a regulatory gene located at the 3' position of the entire IGCH locus.

See also entries [Antibody](#), [Immune Response](#), [Immunoglobulin](#), [IgA](#), [IgE](#), [IgG](#), [IgM](#), and [Isotype](#).

Suggestions for Further Reading



- C. M. Snapper, K. B. Marcu, and P. Zelazowski (1997) *Immunity* **6**, 217–223.  
C. Esser and A. Radbruch (1990) *Annu. Rev. Immunol.* **8**, 717–735.  
M. Cogné et al. (1994) *Cell*, **77**, 737–747.

## Synaptonemal Complexes

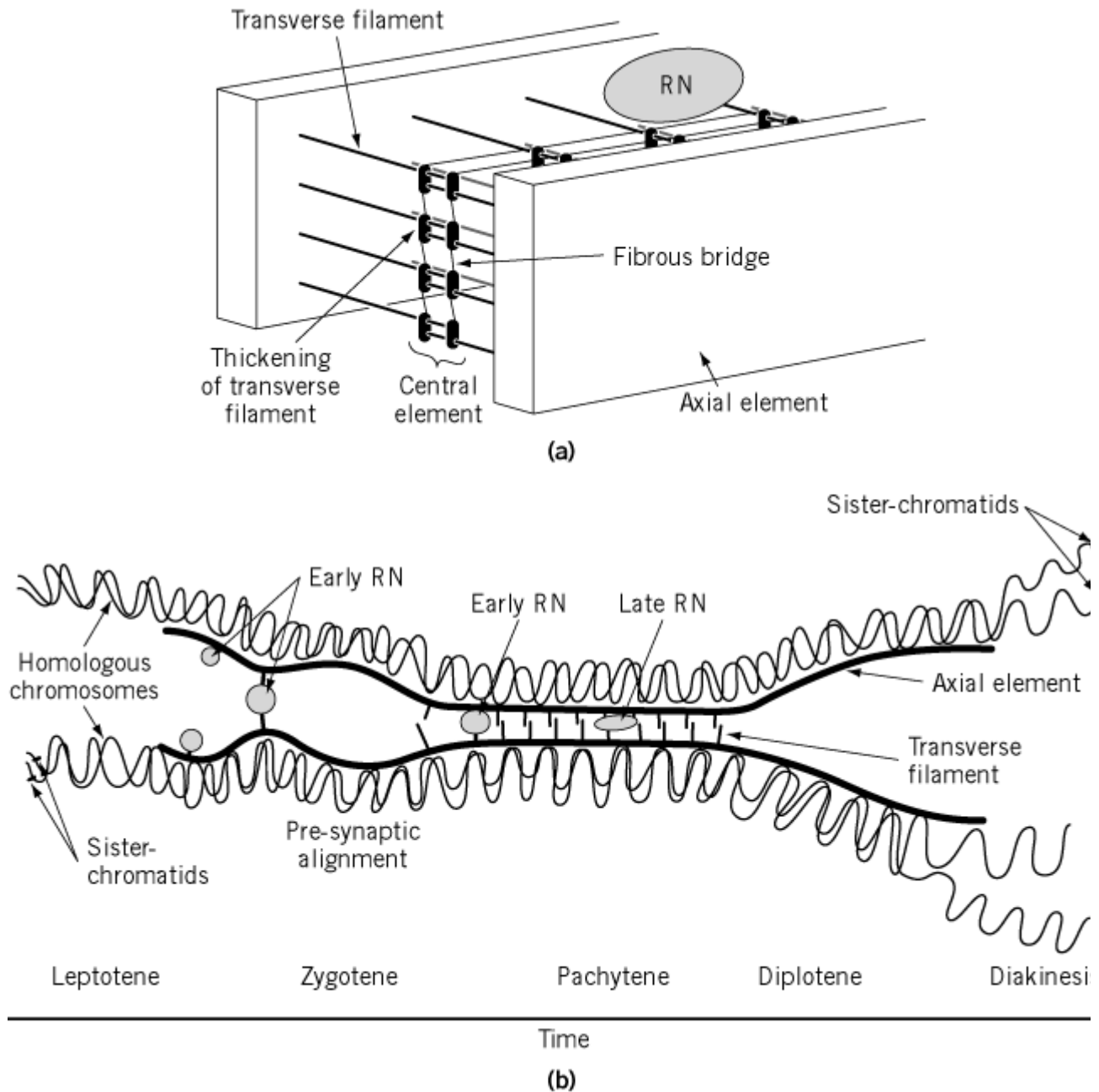
Synaptonemal complexes (SCs) are proteinaceous, zipperlike structures that are formed between homologous [chromosomes](#) during the prophase of the first meiotic division ([1](#)). SCs were first described in 1956 ([2](#), [3](#)), and they have since been found in almost all sexually reproducing eukaryotes analyzed thus far ([1](#)). Their assembly and disassembly accompanies the successive chromosomal alterations of meiotic prophase, namely pairing, [recombination](#), and condensation of homologous chromosomes. Originally, it was thought that SCs created the preconditions for meiotic **crossing over** by bringing corresponding parts of homologous chromosomes in close apposition ([1](#)). However, this view was gradually undermined during recent years, and SCs are now considered as structures that control both the number and distribution of reciprocal exchanges (crossovers) and convert crossovers into functional chiasmata ([4](#)). How SCs fulfill these roles remains to be elucidated.

### 1. Assembly and Disassembly of SCs

Meiosis consists of two successive divisions, meiosis I and meiosis II, which follow a single round of DNA replication. After the premeiotic S-phase, during the prophase of meiosis I, homologous chromosomes pair, non-sister chromatids of homologous chromosomes recombine with each other, the chromosomes condense while they are still paired, and then they disjoin at meiosis I. As a result, **diploid** nuclei produce [haploid](#) nuclei with novel assortments of genes. The assembly and disassembly of synaptonemal complexes accompany these chromosomal alterations. After premeiotic [DNA replication](#), in the leptotene stage of meiotic prophase, the two sister chromatids of each chromosome develop a single proteinaceous axis, called axial element or axial core. Axial elements of SCs differ in several respects from the axes (also called scaffolds) of mitotic metaphase chromosomes: (i) Axial elements of SCs are much longer; (ii) they contain several meiosis-specific proteins (see below and Table [1](#)), and (iii) they are shared by the two sister chromatids of a chromosome, whereas in a mitotic chromosome there is a separate axis (scaffold) for each sister chromatid. Usually, the assembly of axial elements starts at the [nuclear envelope](#) ([1](#), [5](#)). Subsequently, in the zygotene stage, the axial elements of homologous chromosomes are connected (“synapsed”) along their length by numerous transverse filaments to form the structure of an SC (Fig. [1a](#)); the connecting process is called synapsis. In most species, synapsis initiates preferentially, but not exclusively, in the (sub)**telomeric** regions ([1](#), [5](#)). In the pachytene stage of meiotic prophase, synapsis is complete, and the axial elements are connected along their length by transverse filaments. Axial elements within the SC structure are referred to as lateral elements. Between the two lateral elements of an SC, a third longitudinal structure, the central element, is formed on the transverse filaments. Both lateral elements, together with the central element, make up the tripartite structure of the SC (Fig. [1a](#)). In some organisms, complete axial elements are assembled and aligned along their entire length before synapsis ([6](#)). This phenomenon of presynaptic alignment is probably a common feature of meiosis, although it may be difficult to observe in species in which the assembly and alignment of axial elements is rapidly followed by synapsis ([6](#), [7](#)). Aligned axial elements are further apart ( $\geq 300\text{nm}$ ) than synapsed lateral elements (about 100 nm). Fibrous connections appear to pull aligned axial elements together at irregular intervals ([6](#), [7](#)), and synapsis starts where axial elements have approached each other within 300 nm, possibly because transverse filament components

protruding from opposite axial elements can then interact.

**Figure 1.** Structure, assembly, and disassembly of the SC. (a) Overview of the organization of the central region. RN, recombination nodule. This overview is mainly based on Refs. 8 and 9. (b) Assembly and disassembly of SCs. The successive stages of meiotic prophase are indicated from early (left) to late (right). RN, recombination nodule.



**Table 1. Protein Components and Candidate Components of Synaptonemal Complexes**

| Component | First Identified in | Gene or cDNA | Molecular Weight <sup>a</sup> | Localization within SC | Remark <sup>a</sup> TF, transverse filament. | Reference |
|-----------|---------------------|--------------|-------------------------------|------------------------|--|-----------|
|-----------|---------------------|--------------|-------------------------------|------------------------|--|-----------|

|         |        | Cloned |              |                |                      |           |
|---------|--------|--------|--------------|----------------|----------------------|-----------|
| SCP1    | Rat    | Yes    | 120,000 Da   | Central region | TF component         | (18)      |
| SCP2    | Rat    | Yes    | 170,000 Da   | Axial element  | Synapsed+unsynapsed  | (21)      |
| SCP3    | Rat    | Yes    | 30,000 Da    | Axial element  | Synapsed+unsynapsed  | (20)      |
| SC65    | Rat    | Yes    | 50,000 Da    | Central region | Not meiosis-specific | (22)      |
| Unnamed | Rat    | No     | 48,000 $M_r$ | Central region |                      | (19)      |
| Unnamed | Lilium | No     | 52,000 $M_r$ | Axial element  | Synapsed+            | (23)      |
|         |        |        | -70,000      |                | unsynapsed           |           |
| Zip1p   | Yeast  | Yes    | 110,000 Da   | Central region | TF component         | (7),(27)  |
| Hop1p   | Yeast  | Yes    | 70,000 Da    | Axial element  | Unsynapsed           | (24),(26) |
| Red1p   | Yeast  | Yes    | 90,000 Da    | Axial element  | Synapsed+unsynapsed  | (25)      |

<sup>a</sup> Da, molecular weight predicted from nucleotide sequence of gene or cDNA;  $M_r$ , relative electrophoretic mobility.

At the ultrastructural level, the organization of the region between the lateral elements (the central region) is essentially the same in different species (8-10). The structural unit of the central region consists of a single transverse filament, which spans the distance between the axial elements, and carries two symmetrically placed thickenings, at fixed positions (Fig. 1b). Fibrous bridges connect the thickenings with each other in two longitudinal rows, which together constitute the central element. Three to five layers of transverse filaments occur on top of each other within the central region; these layers are also kept in register by fibrous bridges between the thickenings (9). In most species, this tripartite structure of SCs has been conserved. There is, however, much variation in the detailed structure, in particular with respect to the degree of order in the central region, the structure of the central element, and the appearance of the lateral elements (1, 10); in some fungi, the lateral elements display a regular, transverse striping (11). At least two species with an otherwise normal meiosis do not have detectable SCs, namely *Schizosaccharomyces pombe* (fission yeast) (12) and *Aspergillus nidulans* (a filamentous fungus) (13). Meiotic prophase cells of *S. pombe* contain so-called linear elements, which possibly correspond to the lateral elements of SCs (12, 14).

At the end of meiotic prophase, the SC structure is disassembled. In most species, the transverse filaments disappear first and then the axial elements (diplotene stage). Until their disassembly, the axial elements remain attached through their ends to the nuclear envelope; some species have special thickenings, called attachment plaques, by which the axial/lateral elements are attached to the nuclear envelope (1). The disassembly of SCs is in many species followed by the diffuse diplotene stage, in which no axial elements or chromosomal axes can be distinguished (15). In postdiffuse diplotene and diakinesis, the chromosomes recondense in preparation of the first meiotic division; the two sister chromatids then become clearly distinguishable, and each displays their own longitudinally supporting structure, the chromatid scaffold, instead of the single, common axial element of the SC.

## 2. Recombination Nodules

The assembly of SCs is accompanied by the appearance of spherical or ellipsoidal electron-dense structures called recombination nodules (RNs) (reviewed in Ref. [16](#)). In several species, two types of RNs can be distinguished on the basis of their morphology and time of appearance: early and late RNs ([16](#)). Late RNs are found on top of the central region of pachytene SCs (Fig. [1a](#)). Their frequency and distribution along the bivalents in mid-pachytene correlates with the frequency and distribution of chiasmata and crossovers ([16](#)). Late RNs are therefore probably protein complexes that are involved in crossover formation.

Early RNs are more numerous than late RNs; they occur throughout the nucleus in leptotene. During zygotene and early pachytene, they concentrate along unpaired axial elements and on (part of) the fibrous connections between the axial elements of aligned chromosomes. In late zygotene and early pachytene, they are also found on top of the central region of the SC, but their number decreases gradually, until in mid-pachytene all early RNs have disappeared ([16](#), [17](#)). The distribution of early RNs along the bivalents does not correlate with the distribution of crossovers. It has been suggested that early RNs are involved in homology search and that they leave **gene conversions** as footprints of their activity ([16](#)). It is possible that some of the early RNs develop into late RNs (reviewed in Ref. [17](#)).

## 3. SC Components

Protein components can only be assigned with certainty to SCs by ultrastructural studies. Table [1](#) lists the protein components that have thus been assigned to SCs, along with a number of strong candidate SC components. In rodents ([18-22](#)) and *Lilium* ([23](#)), SC proteins were identified in preparations of isolated SCs, and they were localized within the SC structure by immunolabeling. The Hop1 ([24](#)) and Red1 ([25](#)) proteins of budding yeast (*Saccharomyces cerevisiae*) were identified as candidate SC components on the basis of **immunofluorescence** or immunogold studies and cytological analyses of hop1 and red1 mutants ([24](#), [26](#)). Zip1p of yeast has been assigned to the central region on the basis of immunofluorescence studies and a systematic ultrastructural analysis of SCs in a series of mutants carrying deletions or insertions in the zip1 gene ([27](#)). No amino acid sequence similarity has been detected between the SC-proteins of yeast and mammals. However, mammalian SCP1 ([18](#)) and yeast Zip1 ([7](#)) have similar predicted **secondary structures** and are incorporated in a similar way into the SC ([27](#), [28](#)). Most probably, SCP1 and Zip1 are components of the transverse filaments.

In the [mouse](#), an increasing number of proteins that are known or suspected to play a role in recombination and/or **DNA repair** (RAD51, RPA, MLH1, POLb), or in signaling pathways that link DNA-repair to progression of the **cell cycle** (ATM, ATR, BRCA), have been localized on or along SCs by immunofluorescence (reviewed in Ref. [17](#)). In yeast, comparable studies have been performed ([29](#), [30](#)). Most probably, these recombination proteins are not part of SCs. They associate with SCs in a characteristic order and pattern (reviewed in Ref. [17](#) for the mouse), presumably to perform or monitor certain steps in meiotic recombination. Of these proteins, only RAD51 has been localized ultrastructurally ([31](#), [32](#)). It is part of early RNs and localizes along (not on) the axial elements of SCs in zygotene and early pachytene. The immunofluorescence pattern of MLH1 strongly suggests that it is localized in late RNs ([17](#)). Furthermore, in male but not female mice, Hsp70-2 ([33](#)) has been localized along the axial elements, whereas **topoisomerase II** gradually congregates onto the axial elements toward the end of meiotic prophase ([34](#)).

The axial elements/lateral elements of SCs are rich in DNA, whereas very little DNA occurs in the central region ([35](#)). At irregular intervals, threads of DNA or **chromatin** cross the central region; some of these threads are associated with RNs, which appear to be enriched in DNA ([35](#)).

## 4. Possible Functions of SCs

Originally, SCs were considered as structures that bring corresponding parts of homologous chromosomes in close apposition to enable correct meiotic genetic exchange (1). However, this view has been abandoned for various reasons (reviewed in Ref. 4). One obvious reason is provided by *S. pombe* and *A. nidulans*, which lack SCs but perform correct meiotic recombination (12, 13). In budding yeast (*Saccharomyces cerevisiae*), the assembly of SCs has been analyzed systematically in relation to meiotic chromosome pairing and recombination, by the determination of the order of events at the DNA and chromosomal levels and by the analysis of SC assembly and recombination in meiotic mutants (reviewed in Refs. 36 and 37). In yeast, SCs are not required for high levels of meiotic recombination (36); on the contrary, successive steps in SC assembly appear to depend on successive steps in recombination (36, 37). In other organisms, however, apparently normal SCs are assembled in the total absence of meiotic recombination (38, 39). The following possible roles of SCs have been proposed on the basis of the studies in yeast and in various other organisms: (i) SCs are required for crossover interference (36), which is the phenomenon that the presence of a crossover reduces the probability of another crossover nearby on the same chromosome pair (bivalent). Interference can be considered as a manifestation of a regulatory mechanism that ensures at least one crossover per bivalent, even if the average number of crossovers per bivalent is close to one. At least one crossover per bivalent is required for the proper segregation of chromosomes at meiosis I. The intact, tripartite SC structure is required for crossover interference (36), but it is not at all clear how the SC is involved. *S. pombe* and *A. nidulans*, which lack SCs, do not display crossover interference. The average number of crossovers per chromosome pair is so high in these organisms that a regulatory mechanism that ensures at least one crossover per bivalent is not required (12, 13). (ii) The axial elements of SCs may have a role in monitoring the recombination process, so that it can be related to progression of the meiotic cell cycle (40). (iii) The axial elements direct initiated recombination events to the homologous chromosome rather than to the sister chromatid (41). (iv) Axial element components are required to link recombination at the DNA level to crossing over at the level of chromosomal axes (36, 37); this is essential for the formation of functional chiasmata—that is, structures that can ensure proper segregation of homologous chromosomes at meiosis I.

To summarize, SCs appear to steer recombination in the right direction, and they also appear to link recombination events to the meiotic cell cycle and to the mechanisms of meiotic chromosome segregation (4, 42). It remains to be elucidated how SC components are involved in these possible functions.

## Bibliography

1. D. Von Wettstein, S. W. Rasmussen, and P. B. Holm (1984) *Annu. Rev. Genet.* **18**, 331–413.
2. M. J. Moses (1956) *J. Biophys. Biochem. Cytol.* **2**, 215–218.
3. D. W. Fawcett (1956) *J. Biophys. Biochem. Cytol.* **2**, 403–406.
4. R. S. Hawley and T. Arbel (1993) *Cell* **72**, 301–303.
5. H. Scherthan et al. (1996) *J. Cell Biol.* **134**, 1109–1125.
6. J. Loidl (1994) *Experientia* **50**, 285–294.
7. M. Sym, J. Engebrecht, and G. S. Roeder (1993) *Cell* **72**, 365–378.
8. K. Schmekel, U. Skoglund, and B. Daneholt (1993) *Chromosoma* **102**, 682–692.
9. K. Schmekel, J. Wahrman, U. Skoglund, and B. Daneholt (1993) *Chromosoma* **102**, 669–681.
10. K. Schmekel and B. Daneholt (1995) *Trends Cell Biol.* **5**, 239–242.
11. D. Zickler (1973) *Chromosoma* **40**, 401–416.
12. J. Kohli and J. Baehler (1994) *Experientia* **50**, 295–306.
13. M. Egel-Mitani, L. W. Olson, and R. Egel (1982) *Hereditas* **97**, 179–187.
14. J. Baehler, T. Wyler, J. Loidl, and J. Kohli (1993) *J. Cell Biol.* **121**, 241–256.
15. A. J. J. Dietrich and R. J. P. Mulder (1981) *Chromosoma* **83**, 409–418.
16. A. T. C. Carpenter (1994) *Bioessays* **16**, 69–74.

17. T. Ashley and A. Plug (1997) In *Meiosis and Gametogenesis* (M. A. Handel, ed.), Academic Press, San Diego, pp. 201–239.
18. R. L. J. Meuwissen, H. H. Offenberg, A. J. J. Dietrich, A. Riesewijk, M. van Iersel, and C. Heyting (1992) *EMBO J.* **11**, 5091–5100.
19. A. Smith and R. Benavente (1992) *Exp. Cell Res.* **198**, 291–297.
20. J. H. M. Lammers et al. (1994) *Mol. Cell. Biol.* **14**, 1137–1146.
21. H. H. Offenberg et al. (1998) *Nucleic Acids Res.* **26**, 2572–2579.
22. Q. Chen, R. E. Pearlman, and P. B. Moens (1992) *Biochem. Cell Biol.* **70**, 1030–1038.
23. M. K. Anderson, S. M. Stack, R. J. Todd, and R. P. Ellis (1994) *Chromosoma* **103**, 357–367.
24. N. M. Hollingsworth, L. Goetsch, and B. Byers (1990) *Cell* **61**, 73–84.
25. E. A. Thompson and G. S. Roeder (1989) *Mol. Gen. Genet.* **218**, 293–301.
26. A. V. Smith and G. S. Roeder (1997) *J. Cell Biol.* **136**, 957–967.
27. K. S. Tung and G. S. Roeder (1998) *Genetics* **149**, 817–832.
28. K. Schmekel et al. (1996) *Exp. Cell Res.* **226**, 20–30.
29. D. Bishop (1994) *Cell* **79**, 1081–1092.
30. P. Ross-McDonald and G. S. Roeder (1995) *Cell* **79**, 1069–1080.
31. L. K. Anderson, H. H. Offenberg, W. H. M. Verkuijlen, and C. Heyting (1997) *Proc. Natl. Acad. Sci. USA* **94**, 6868–6873.
32. P. B. Moens et al. (1997) *Chromosoma* **106**, 207–215.
33. J. W. Allen et al. (1996) *Chromosoma* **104**, 414–421.
34. P. B. Moens and W. Earnshaw (1989) *Chromosoma* **98**, 317–322.
35. G. H. Vázquez Nin et al. (1993) *Chromosoma* **102**, 457–463.
36. G. S. Roeder (1995) *Proc. Natl. Acad. Sci. USA* **92**, 10450–10456.
37. G. S. Roeder (1997) *Genes & Dev.* **11**, 2600–2621.
38. A. F. Dernburg et al. (1998) *Cell* **94**, 387–398.
39. K. S. McKim et al. (1998) *Science* **279**, 876–878.
40. L. Xu, B. M. Weiner, and N. Kleckner (1997) *Genes & Dev.* **11**, 106–118.
41. A. Kleckner and A. Schwacha (1997) *Cell* **90**, 1123–1135.
42. M. P. Maguire (1995) *J. Heredity* **86**, 330–340.

### Suggestions for Further Reading

43.

References [1](#), [6](#), [16](#), [36](#), and [42](#) of the above list are suitable for further reading. Reference [1](#) is a comprehensive review of the cytology of SCs. Reference [6](#) is more up-to-date than Ref. [1](#) and relates the cytology of SCs to *in situ* hybridization studies and analyses of meiotic recombination at the DNA level. Reference [16](#) is focused on recombination nodules. Reference [36](#) is the best start to get familiar with the complex literature on meiosis and synaptonemal complexes in yeast. Reference [42](#) is a good and up-to-date introductory text.

44. N. Kleckner (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8167–8174. This review considers various possible roles of SCs in meiotic chromosome pairing and recombination.

### Syncytium

A syncytium is a cell that contains multiple **nuclei** in the same cytoplasm. Some organisms are completely syncytial, such as the slime mold [Physarum](#) *Physarum*. In other organisms, only some tissues are a syncytium. For example, human muscles are syncytia, formed by the fusions of individual cells called myoblasts. A syncytium may also be present at only some stages in [development](#). Early *Drosophila* [embryos](#) (prior to [blastoderm](#) formation) are a syncytium (1). The first nuclear divisions after fertilization take place without any cellular division. After 13 divisions, the germ cells have formed separate cells, but the somatic nuclei remain in one cytoplasm. Most of the somatic nuclei are at the surface and form a stage called the syncytial blastoderm. Only after 13 divisions does cytoplasmic division follow the somatic nuclear divisions. Cellular membranes form around all of the somatic nuclei at the periphery of the blastoderm, to form the cellular blastoderm stage. A few nuclei remain in the interior of the blastoderm in a syncytium with the yolk.

#### Bibliography

1. F. R. Turner and A. P. Mahowald (1976) *Dev. Biol.* **50**, 95–108.

#### Suggestion for Further Reading

2. S. J. Counce (1961) The analysis of insect embryogenesis. *Annu. Rev. Entomol.* **6**, 295–312.

### Synonymous Substitution

Synonymous substitution is an **evolutionary** term meaning the fixation of a [silent mutation](#) in a population or subpopulation. Because the [genetic code](#) is redundant, particularly in the third position (see [Degeneracy of the Genetic Code](#)), silent mutations are frequent. If codons for the 20 amino acids were distributed equally in **genomes** and if all types of base substitutions were equally likely, 19% of all [base-pair substitution](#) mutations would be silent. The rate of synonymous substitution during evolution, however, is far higher than this. Indeed, synonymous substitutions can be 70 to 90% of all substitutions (1). The obvious reason is that most amino acid changes have deleterious effects, whereas silent mutations can be **neutral** or nearly so. Thus, they may not be **selected** against and come to fixation by [genetic drift](#). But organisms can have preferred codons for certain amino acids (see [Codon Usage and Bias](#)), in which case particular silent mutations may be selected for or against, especially when they occur in highly expressed genes.

#### Bibliography

1. T. S. Whittam (1996) In *Escherichia coli and Salmonella; Cellular and Molecular Biology*, 2nd ed. (F. C. Neidhardt, R. Curtiss, III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaechter, and H. E. Umbarger, eds.), American Society for Microbiology Press, Washington, DC, pp. 2708–2720.

### Syntenic Genes

*Syntenic genes* are **genes** situated on the same [chromosome](#). Genes that are syntenic on one chromosome of a species are often also syntenic on a chromosome of another species. A few examples are given below. For example, syntenic relationships of genes for *galactokinase* and thymidine [kinase](#) are conserved between humans (chromosome 17) and chimpanzees (chromosome 19). Other syntenic relationships and gene assignments are consistent with proposed homologies between chimpanzee and human chromosomes (1). In the same vein, homologous genes for enolase, phosphogluconate dehydrogenase, phosphoglucomutase, and adenylate kinase are syntenic on mouse chromosome 4 and human chromosome 1 (2). Many syntenic groups have been found between the baboon *Papio papio* and *Homo sapiens* by cosegregation analysis of [hybrid cells](#) obtained between *Papio* fibroblasts and a mouse [cell line](#) deficient in thymidine kinase 3, 4. Mapping genes spanning the chromosomes in humans, cattle, and mice permits syntenic comparisons between prototypic genomes of three mammalian orders and provides insight into the evolutionary history of the ancestral mammalian genome (5).

### Bibliography

1. S. Chen et al. (1976) *Somatic Cell Genet.* **2**, 205–213.
2. P. A. Lalley, U. Francke, and J. D. Minna (1978) *Proc. Natl. Acad. Sci. USA* **75**, 2382–2386.
3. N. Creau-Goldberg et al. (1982) *Ann. Genet.* **25**, 14–18.
4. N. Creau-Goldberg et al. (1983) *Ann. Genet.* **26**, 75–78.
5. N. Zhang and J. E. Womack (1992) *Genomics* **14**, 126–130.

### Synthase/Lyase

Lyases are [enzymes](#) that catalyze the cleavage or formation of C—C, C—O and C—N bonds by means other than hydrolysis or oxidation. They have a single substrate in one direction and two substrates in the other. The carbon-carbon lyases include those that bring about carboxylation reactions, reverse aldol condensations, or the reversible cleavage of a 3-hydroxy acid. The action of carbon-oxygen lyases leads to the formation of unsaturated products, and the recommended names for specific enzymes belonging to this group are dehydratases, hydratases, and synthases. The carbon-nitrogen lyases catalyze the elimination of ammonia with the formation of a double bond and are either ammonia or amidine lyases. Carbon-sulfur lyases also belong to this group. It should be mentioned that, under conditions in which the synthesis of the single substrate is considered to be more important, the alternative name of synthase is used. A synthase differs from a **synthetase** in not having a requirement for ATP or another nucleoside triphosphate.

### Synthetases/Ligases

Ligases are a group of [enzymes](#) that are involved with linking together two molecules covalently, with the concomitant hydrolysis of a phosphoryl or pyrophosphoryl group of a nucleoside triphosphate that is generally **ATP**. The bonds formed in the joining of two molecules are C—O, C—S, and C—C. Enzymes catalyzing reactions of this type were referred to originally as *synthetases*, but this term proved to be unsatisfactory because of the confusion with **synthase**. Nevertheless,



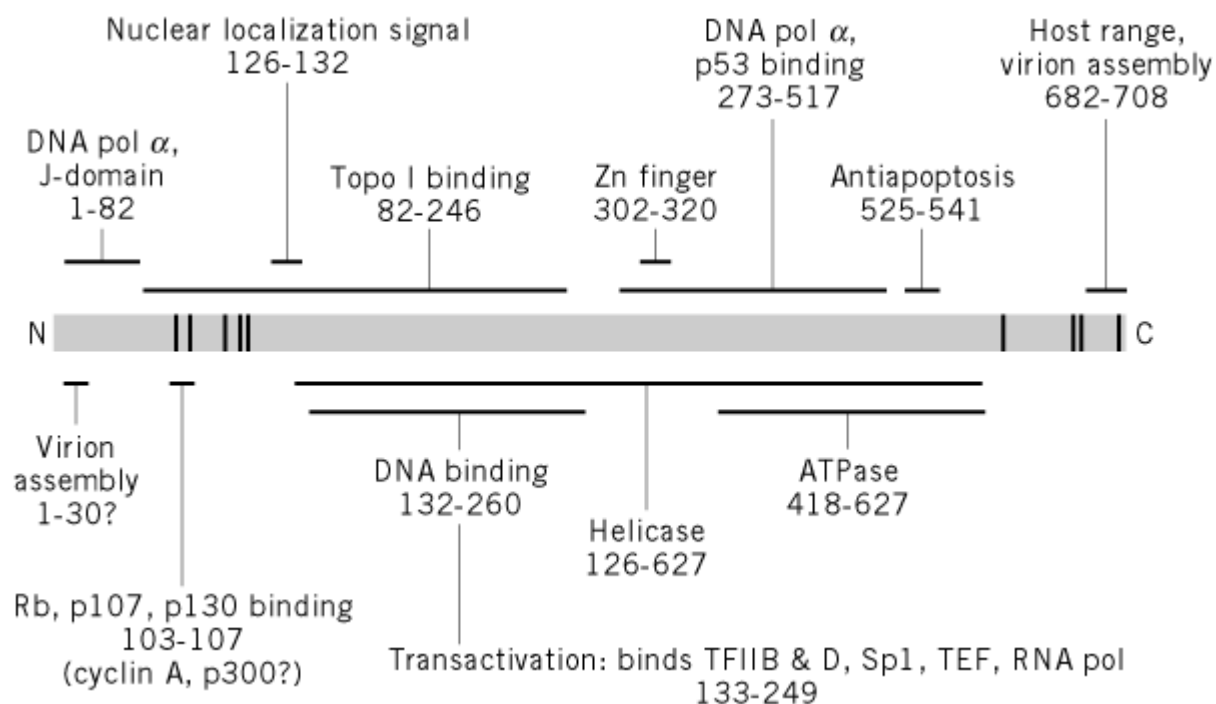
*synthetase* is still being used as a trivial name. It should also be noted that older names are being retained. For example, *pyruvate carboxylase* is the name given to the enzyme responsible for the conversion of pyruvate to oxaloacetate at the expense of ATP, and *synthase* is used for enzymes that restore phosphodiester bonds in nucleic acids.

## T Antigen

T antigens are multifunctional [proteins](#) encoded by **papovaviruses**, such as simian virus 40 (**SV40**), mouse **polyoma virus**, and human BK and JC viruses. Under certain circumstances, animals infected with a papovavirus can develop tumors and will produce **antibodies** against the early gene products of the virus, which are called tumor or “T” antigens. T antigens regulate expression of both viral and host genes, modulate cellular proliferation, and participate in viral DNA synthesis and virion assembly.

The most intensively studied T antigen is that of SV40 virus. Although there are some differences, the T antigens of other papovaviruses function quite similarly to that of SV40. Pre- [messenger RNA](#) transcribed from the A gene of the SV40 virus is differentially spliced (see **Splicing, RNA**) to produce two mRNAs that yield the SV40 large T and small t antigens. The 708 amino acid-residue large T antigen has a molecular weight of about 82,500 (Fig. 1). The 17,000-Da small t antigen comprises the amino-terminal 82 amino acid residues of the large T antigen plus an additional 97 residues encoded by the **intron** of the A gene.

**Figure 1.** Distribution of functional domains on SV40 large T antigen. The polypeptide chain of T antigen is represented by the stippled bar with *N* at the amino terminus and *C* at the carboxy terminus. Vertical lines within the bar represent phosphorylation sites. Numbers indicate the residues that constitute each domain, which is indicated by horizontal lines. (Modified, with permission, from the *Annual Review of Biochemistry*, Vol. 61, © 1992 by Annual Reviews Inc.)

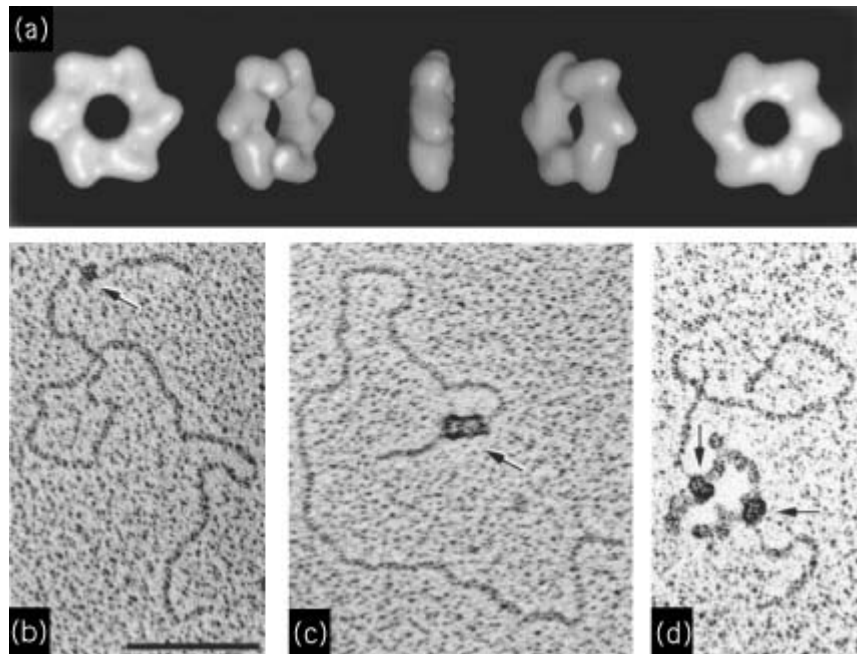


Large T antigen undergoes numerous types of [post-translational modifications](#) that include acylation, glycosylation, ribosylation, adenylation, and phosphorylation. These modifications modulate the activity and distribution of the protein within the cell. The phosphorylation sites are clustered near the amino and carboxy ends of the polypeptide chain (Fig. 1). Phosphorylation of the various sites within the clusters modulates both positive and negative control of viral [DNA replication](#) .

T antigen is composed of several functional **domains** that have a variety of functions (1-9) (Fig. 1). One of these domains binds DNA and consists structurally of a central five-stranded antiparallel **b-sheet** that is flanked by a pair of **a-helices** on one side and an a-helix and a  $3_{10}$  helix on the other side (see [Protein Structure](#)). Residues 152–155 and 203–207 are juxtaposed to form a surface that contacts the recognition sequence of the DNA (10). This domain has a greater affinity for single-stranded DNA than for double-stranded, unless the duplex DNA contains the pentanucleotide sequence GAGGC, which is specifically recognized and bound. Clusters of this recognition sequence are located near and within the SV40 **origin of replication**, which separates the **promoters** for the early and the late genes. The spacing of these pentanucleotides within the clusters promotes protein–protein contacts that further strengthen the affinity of T antigen for duplex DNA. The binding of T antigen to these clusters governs the transcription of viral genes and the initiation of viral DNA replication.

Another domain possesses an [ATPase](#) activity that is strongly stimulated by the binding of single-stranded DNA. The binding of either ADP or ATP to T antigen induces its oligomerization into a hexamer, which is propeller-shaped with a central channel (11) (Fig. 2a). The hexameric form of T antigen functions as a DNA helicase, which translocates in a 3'–5' direction along one strand of the duplex as it unwinds the DNA. The helicase domain spans the ATPase and DNA-binding domains and contains a **zinc finger**.

**Figure 2.** Reconstructed image and electron micrographs of T antigen. (a) Reconstructed three-dimensional images of hexameric T antigen assembled in the presence of ADP. (Image panel provided by M. C. San Martin and J. M. Carazo; reprinted with permission from reference 11). (b) T antigen bound to a linearized **plasmid** containing the SV40 origin of replication near one end of the DNA molecule. Arrow points to the bound T antigen. (c) Double hexamer of T antigen bound to a linearized plasmid containing the SV40 origin near one end of the DNA molecule. Arrow points to the double hexamer. (d) Unwound structure formed by mixing linear plasmid DNA containing the SV40 origin with T antigen, single-stranded DNA binding protein, topoisomerase, and ATP. Arrows point to hexameric T antigen at the unwinding forks. (Bar = 0.1  $\mu\text{m}$ ) (Electron micrographs provided by Dr. F. B. Dean.)



At the amino terminus of T antigen is a domain that is functionally homologous with the J-domain of DnaJ/hsp40 **molecular chaperones** (see [DnaK/DnaJ Proteins](#)). This domain interacts with the host hsc70 protein (see **Hsc, hsp proteins**) and is also required for efficient synthesis of viral DNA. The extreme carboxy terminus determines the host range of the virus, and in conjunction with the J domain, also governs virion maturation. Other domains modulate host [transcription](#), [cell cycle](#) progression, and defense responses by interacting with the cellular factors that regulate these processes.

T antigen regulates both viral and host transcription. After the initial synthesis of T antigen, transcription of early viral genes is down-regulated by the binding of monomeric or dimeric T antigen to one of the clusters of pentanucleotide recognition sites at the origin (Fig. [2b,c](#)). Positive regulation of late viral transcripts by T antigen is partially achieved by binding to TEF-1, which is a cellular [transcription factor](#) that represses late transcription of viral genes ([12](#)).

T antigen transactivates simple host promoters that contain a [TATA box](#) and binding sites for the transcription factors SP-1, ATF, AP1, or TEF-1. T antigen binds SP1, TFIIB, TATA-binding protein, and other components of transcription factor TFIID, and the 140-kDa subunit of **RNA polymerase II** ([3](#)). T antigen also transactivates **ribosomal** RNA genes by binding to SL1 protein, which is required for efficient transcription of this class of genes ([13](#)). The mechanism by which the binding of T antigen to these factors transactivates host genes is not completely known. T antigen may simply function as a TATA-binding protein-associated factor in a complex with TFIID. Alternatively, the interaction of T antigen with host transcription factors may block the binding of [repressors](#) that inhibit the formation of preinitiation complexes.

T antigen can also repress the expression of certain host genes by sequestering transcriptional activators such as p300, which modulates the expression of genes controlled by **cAMP-response elements** (CRE). T antigen binds specific **isoforms** of p300, which results in repression of CRE promoters ([8](#)).

T antigen interacts with several cellular factors that govern cell growth, such as the **retinoblastoma** tumor suppressor (Rb) and [p53](#). T antigen contains domains that bind to each of these cellular proteins (Fig. [1](#)). Rb is a growth suppressor that normally sequesters the transcription factor E2F. By sequestering Rb, T antigen frees E2F, which then promotes the expression of factors that advance the

entry of the cell into the **S phase** of the cell cycle. The freed E2F also promotes the expression of factors required for nucleotide synthesis and for DNA replication. The Rb-related factors p107 and p130, which interact with E2F-like transcription factors, are also bound and inhibited by T antigen (14). p53 is a damage response factor that indirectly stimulates the activity of Rb and also participates in the induction of [apoptosis](#). By sequestering p53, T antigen further inhibits the suppression of E2F by Rb and impedes the induction of apoptosis. T antigen also induces the expression of several factors that are required for progression of the cell cycle, including [cyclin A](#), cyclin B, p33cdk2, and p34cdc2 (15). Two of these factors, cyclin A and p33cdk2, form a complex with T antigen at replication initiation sites on the viral DNA (16).

T antigen can induce the formation of tumors in rodents that have been inoculated with SV40. Rodents are not a natural host for SV40, and the virus is unable to complete its life cycle. However, the viral DNA may be integrated into the [genome](#) of the rodent cells, which leads to expression of T antigen. Because T antigen affects the progression of the cell cycle in rodents just as it does in simians, critical cell-cycle checkpoints are disabled (15). [Karyotype](#) instability eventually ensues, which may then lead to oncogenic [transformation](#) of the infected cells. SV40 may also induce tumors in its natural host if the virus contains a mutation that blocks completion of its life cycle but does not alter the overall structure or function of T antigen itself.

T antigen participates in viral DNA replication, where it functions as a DNA helicase and binds at the origin of replication. On binding ATP, T antigen assembles into two hexamers that encircle the DNA at the replication origin (17) (Figure 2c). The hexamers then recruit host-encoded replication factors to the origin via specific protein-protein interactions. The host replication factors that specifically interact with T antigen include single-stranded-DNA binding protein (SSB), topoisomerase I (see [DNA Topology](#)), and DNA polymerase  $\alpha$  (2, 18). Upon association with the SSB and topoisomerase, each hexamer begins unwinding the origin DNA (Fig. 2 d). DNA replication is then initiated at each unwound fork by DNA polymerase  $\alpha$ . Additional replication factors are then recruited so that the synthesis of viral DNA can be completed.

## 1. Small t Antigen

Although not required for viral growth, small t antigen stimulates viral DNA synthesis (19). Small t antigen is required for efficient transformation of growth-arrested cells, where it promotes the transition to the S phase of the cell cycle. The promotion of this transition is probably a consequence of small t antigen's ability to induce the expression of [cyclin D1](#) and to bind and block the activity of protein **phosphatase 2A**, which dephosphorylates threonine residues (4, 20). Inhibition of this phosphatase may activate **mitogen**-activated protein kinases that induce signaling pathways that initiate the transition to S phase. Since phosphatase 2A can dephosphorylate large T antigen, small t may also facilitate the activation of large T antigen by sequestering this phosphatase. Small t antigen may also activate and repress the expression of various host genes (21).

## Bibliography

1. S. L. Spence and J. M. Pipas (1994) *Virology* **204**, 200–209.
2. D. T. Simmons, T. Melendy, D. Usher, and B. Stillman (1996) *Virology* **222**, 365–374.
3. S. D. Johnston, X. M. Yu, and J. E. Mertz (1996) *J. Virol.* **70**, 1191–1202.
4. E. Fanning (1992) *J. Virol.* **66**, 1289–1293.
5. B. Damania and J. C. Alwine (1996) *Genes Develop.* **10**, 1369–1381.
6. S. D. Conzen, C. A. Snay, and C. N. Cole (1997) *J. Virol.* **71**, 4536–4543.
7. K. S. Campbell, K. P. Mullane, I. A. Aksoy, H. Stubdal, J. Zalvide, J. M. Pipas, A. Silver, T. M. Roberts, B. S. Schaffhausen, and J. A. DeCaprio (1997) *Genes Develop.* **11**, 1098–1110.
8. M. L. Avantaggiati, M. Carbone, A. Graessmann, Y. Nakatani, B. Howard, and A. S. Levine (1996) *EMBO J.* **15**, 2236–2248.
9. J. P. Adamczewski, J. V. Gannon, and T. Hunt (1993) *J. Virol.* **67**, 6551–6557.

10. X. Luo, D. G. Sanford, P. A. Bullock, and W. W. Bachovchin (1996) *Nat. Struct. Biol.* **3**, 1034–1039.
11. M. C. San Martin, C. Gruss, and M. J. Carazo (1997) *J. Mol. Biol.* **268**, 15–20.
12. L. C. Berger, D. B. Smith, I. Davidson, J. J. Hwang, E. Fanning, and A. G. Wildeman (1996) *J. Virol.* **70**, 1203–1212.
13. W. Zhai, J. A. Tuan, and L. Comai (1997) *Genes Devel.* **11**, 1605–1617.
14. J. Zalvide and J. A. DeCaprio (1995) *Mol. Cell. Biol.* **15**, 5800–5810.
15. T. H. Chang, F. A. Ray, and D. A. Thompson (1997) *Oncogene* **14**, 2383–2393.
16. D. Cannella, J. M. Roberts, and R. Fotedar (1997) *Chromosoma* **105**, 349–359.
17. J. A. Borowiec, F. B. Dean, P. A. Bullock, and J. Hurwitz (1990) *Cell* **60**, 181–184.
18. S. Waga and B. Stillman (1994) *Nature* **369**, 207–212.
19. C. Cicala, M. L. Avantaggiati, A. Graessmann, K. Rundell, A. S. Levine, and M. Carbone (1994) *J. Virol.* **68**, 3138–3144.
20. G. Watanabe, A. Howe, R. J. Lee, C. Albanese, I. W. Shu, A. N. Karnezis, L. Zon, J. Kyriakis, K. Rundell, and R. G. Pestell (1996) *Proc. Natl. Acad. Sci. USA* **93**, 12861–12866.
21. W. B. Wang, I. Bikel, E. Marsilio, D. Newsome, and D. M. Livingston (1994) *J. Virol.* **68**, 6180–6187.

### Suggestions for Further Reading

22. E. Fanning, and R. Knippers (1992) Structure and function of simian virus 40 large tumor antigen. *Ann. Rev. Biochem.* **61**, 55–85. (This review is quite comprehensive and contains an extensive bibliography.)
23. J. W. Ludlow (1993) Interactions between SV40 large-tumor antigen and the growth suppressor proteins pRB and p53. *FASEB J.* **7**, 866–871.
24. J. M. Pipas (1992) Common and unique features of T antigens encoded by the polyomavirus group. *J. Virol.* **66**, 3979–3985. (This review discusses the properties of T antigens encoded by several other important papovaviruses.)
25. J. Tooze, ed. (1981) *Molecular Biology of Tumor Viruses*, part 2, rev. ed. Cold Spring Harbor Laboratory, Cold Spring Harbor, N. Y. (Although somewhat dated, this reference is still quite useful and extensively describes the biology of papovaviruses and the T antigens that they encode.)

## T Cell

T lymphocytes represent an essential component of the immune system, because they are responsible for cell-mediated immunity, such as allograft rejection or delayed hypersensitivity. They are also necessary for most humoral responses—that is, production of circulating **antibodies** by **B cells** to the so-called **T-cell-dependent antigens**, which are mainly **proteins**. T cells that are involved in cell-mediated immunity are called *effector T cells*, such as **cytotoxic T lymphocytes**, whereas those that control B (and effector T) expression are designated as helper T cells, or Th.

T cells recognize antigens in a manner that is totally different from that of B cells, because they do not interact with native **epitopes** but with **peptides** derived from **antigen processing** and presented in association with **major histocompatibility complex** (MHC) molecules at the cell surface of antigen-presenting cells. T-cell recognition is mediated by the **T-cell receptor** (TCR), which is

expressed exclusively at the cell surface, as part of a membrane receptor complex that contains also a signaling module, termed CD3. TCR is a heterodimer composed either of  $\alpha$  and  $\beta$  or  $\gamma$  and  $\delta$  polypeptide chains, which are generated upon [gene rearrangement](#) of discrete mosaic genes organized in a fashion (ie, V–D–J and V–J) similar to that of [immunoglobulin](#) (Ig) genes. In addition to the TCR–CD3 complex, T cells express many other surface components, among which CD4 and CD8 are of prime importance, because they condition interactions with MHC class II and class I molecules, respectively.

Differentiation of T cells takes place in the thymus, which is first colonized by early progenitors that emigrate from the bone marrow. The discrete steps that lead to fully differentiated T cells have been dissected with a large collection of mice inactivated for most genes implicated in this complex cascade of events. Critical markers are the sequential rearrangements of the TCR  $\alpha$  and  $\beta$  genes and expression of the CD4 and CD8 molecules. Early progenitors first migrate in the thymic cortex with a TCR<sup>-</sup>, CD8<sup>-</sup> CD4<sup>lo</sup> phenotype. In the subcortical epithelium, progenitors start to divide and are designated as triple-negative (TN) because they express neither CD4, CD8, nor TCR molecules. The first gene rearrangements (V–D–J of the  $\beta$  locus) take place. The  $\beta$  chain associates with preT $\alpha$  and becomes expressed as a heterodimer at the cell surface. This molecule controls [allelic exclusion](#) of the  $\beta$  locus, thus ensuring monoclonal expression of the  $\beta$  chain. As thymocytes migrate deeper into the cortical epithelium, the  $\alpha$  locus V–J gene rearranges, and the TCR $\alpha\beta$  may now be expressed at the cell surface. Both the CD4 and CD8 molecules become expressed, so that the double-positive (DP) CD4<sup>+</sup>CD8<sup>+</sup> stage is now reached. The following steps will be critical, because they will drive the expressed [repertoire](#) of circulating T cells. At the DP stage, when still in the cortical epithelium, cells will interact with MHC molecules that have bound a huge variety of endogenous peptides. T cells that bind with a high affinity will be eliminated, as dangerously autoreactive. T cells that do not encounter any MHC peptide complex will die by “neglect.” Those that interact with a moderate affinity will be positively selected and enter terminal maturation in the thymic medulla. The DP cells now interact with either MHC class I or MHC class II, resulting in the occurrence of single-positive (SP) cells that are CD8 or CD4, respectively. The “choice” of interacting with either MHC class is more or less predicted by the sequence of the TCR molecules, especially in some complementarity determining regions, CDR1 and CDR2 regions (see [Immunoglobulin Structure](#)), that may have coevolved with monomorphic interacting regions of the MHC molecules. Recent **X-ray crystallographic** structures have suggested a structural basis for a better understanding of the positive selection.

Further functional specialization of T cells is dependent on antigenic stimulation and a complex interplay with a variety of cytokines. Helper T cells (Th) are generated from CD4 cell Th precursors (Thp) upon antigenic stimulation that takes place through interaction with an antigen-presenting cell. Thp will mature to Th0 cells that, depending on the nature of the antigenic stimulus and/or the cytokines present in the microenvironment (the two later being closely linked), will finally acquire the distinctive properties of Th1 or Th2 cells. Th1, stimulated by **interleukin-12** (IL-12), are involved in cell-mediated immune responses and secrete IL-2, [tumor necrosis factor](#), and [interferon  \$\gamma\$](#)  (IFN $\gamma$ ). In contrast, Th2 are reactive to IL-4 and are implicated in the control of T-dependent antibody responses, including immediate hypersensitivity mediated by [IgE](#); they synthesize IL-4, IL-5, IL-6, IL-10, and IL-13. Th1 and Th2 down-regulate each other, mostly through IFN $\gamma$  that down-regulates Th2, and IL-4 that represses Th1. The distinction between Th1 and Th2, described mostly in mice, may not be so clear-cut, and it may be better to consider instead the cytokines of type I or type II that are organized as a complex interacting network.

T cells may also behave as direct effectors for cell-mediated responses. This is the case for CD4 T cells that are responsible for delayed-type hypersensitivity and of CD8 [cytotoxic T lymphocytes](#). The latter are derived from cytotoxic T precursors (Tcp) that interact with helper T cells to become potent effectors of T cell cytotoxicity, which may target a large variety of cells that they are induced to kill by [apoptosis](#) after mobilization of the Fas system.

Finally, a word should be mentioned about suppressor T cells (Ts), once extremely popular but now less so. The existence of Ts cells was postulated when it was reported that certain **tolerance** states could be transferred from mouse to mouse, thus excluding the simple explanation of tolerance as merely due to clonal deletion. The suppression problem has not really been solved, especially with the realization of dominant tolerance, but it might also be understood as a complex regulatory network involving subsets of T cells.

See also entries **Antigen presentation**, [Epitope](#), [Gene Rearrangement](#), [Recombinase](#), and [T-Cell Receptor \(TCR\)](#).

#### Suggestions for Further Reading

J. Fraser, D. Strauss, and A. Weiss (1993) Signal transduction events leading to T cell lymphokine gene expression. *Immunol. Today* **14**, 357–362.

G. Anderson, N. C. Moore, J. J. T. Owen, and E. J. Jenkinson (1996) Cellular interactions in thymocyte development. *Annu. Rev. Immunol.* **14**, 73–99.

H. J. Fehling and H. von Boehmer (1997) Early T cell development in the thymus of normal and genetically altered mice. *Curr. Opin. Immunol.* **9**, 263–275.

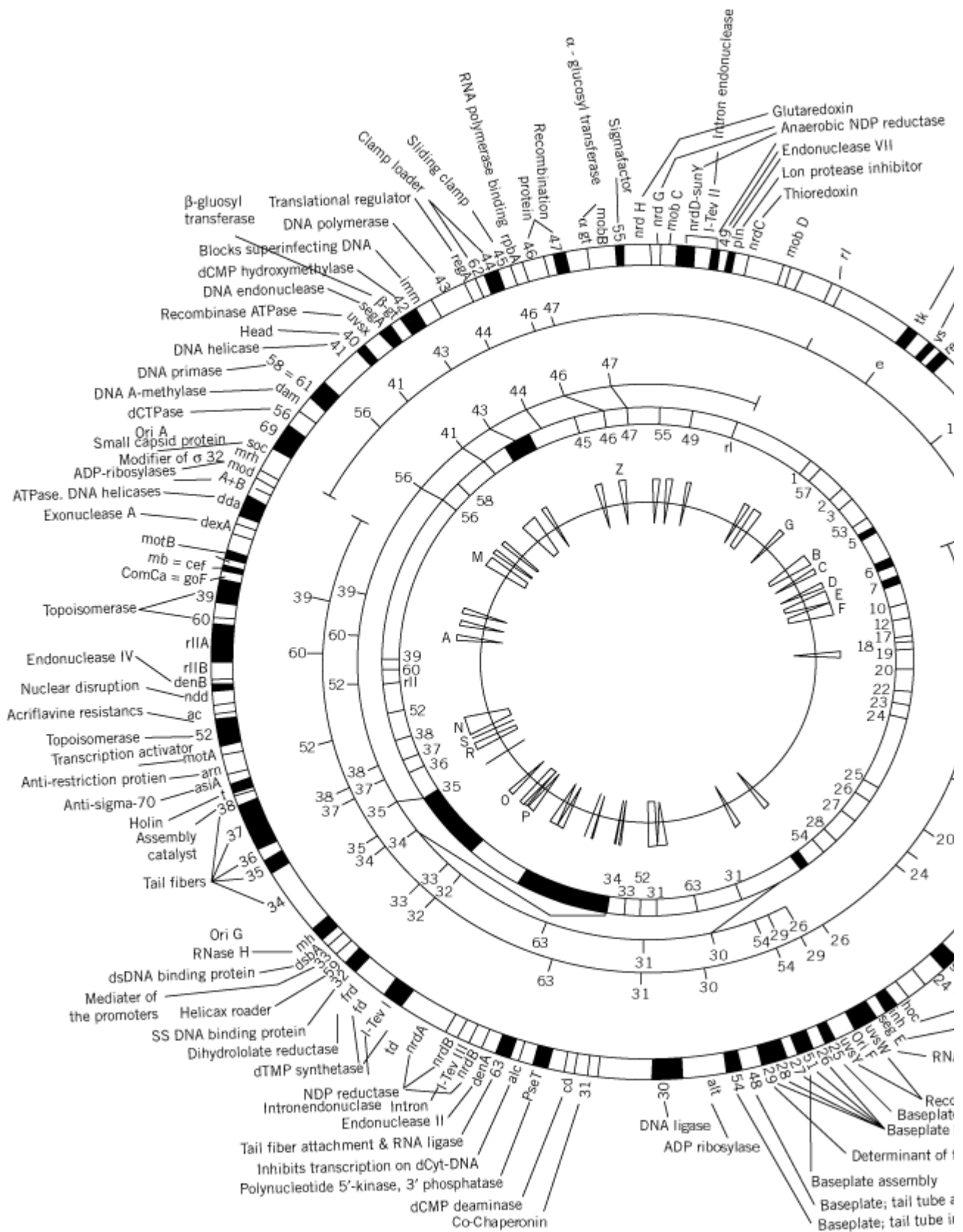
P. Kisielow and H. von Boehmer (1995) Development and selection of T cells: facts and puzzles. *Adv. Immunol.* **58**, 87–209.

## T4 and T-Even Bacteriophages

### 1. History and Overview

Bacteriophages T2, T4 and T6 were among seven *Escherichia coli* phages selected by Max Delbrück to study the fundamentals of viral replication. They are serologically and genetically related and are collectively called the “T-even” phages (1). Based on similar genomic organization, regulatory patterns, virion structure and sequence similarity of their “essential” genes, many other phages from different parts of the world belong to this family (2). More distantly related phages have been called pseudo T-even phages (3) and schizo T-even phages (4). The genomes of these phages are contained in large ( $\approx 174$ -kbp) linear, double-stranded DNA molecules, whose termini contain repetitions of 3 to 5% of the 168,903 base pair genome (5). The termini are randomly permuted over circular maps (Fig. 1). The DNA contains hydroxymethylcytosine (HMC) instead of cytosine. In most members of the family, the HMC residues are further glycosylated to different extents. These modifications allow escape from host [restriction enzymes](#) and are important for the developmental strategy of these phages.

**Figure 1.** Comparison of several maps of T4 genes. The outermost circle shows the map of known genes and origins of 1 (5). The next three overlapping circular segments drawn as thin lines show the positions of the indicated genes based on packaging between these genes and reference markers in *rI*, *rII*, and *rIII* (76). The next circle shows the position of these frequencies (107). The innermost circle represents the heteroduplex loops between T2 and T4 DNA (90). The substitutio gene 69 (40).

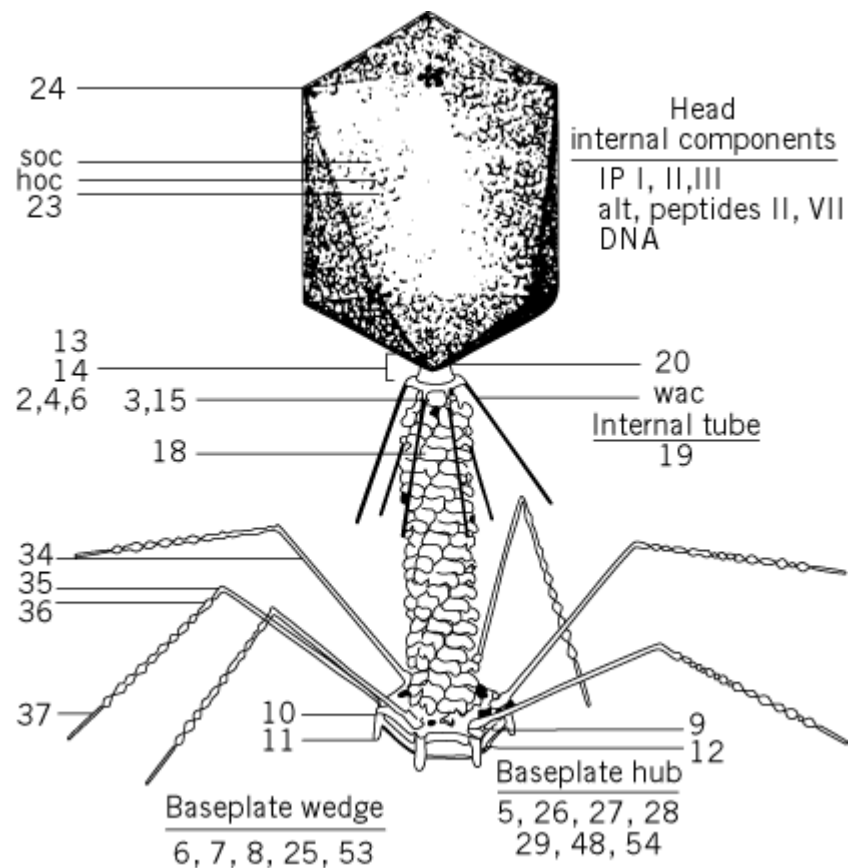


DNA is packaged in elongated “heads” of quasi-**icosahedral symmetry**. Heads are joined to tails whose baseplates and fibers (Fig. 2) are instrumental for recognition, adsorption, and injection of the DNA into host bacteria. Different phages of the T-even family recognize different receptors in different host strains. Apparently, illegitimate recombination with genomes of plasmids and other



phage genomes in the same host cell facilitates rapid evolution of T-even tail fiber genes and adaptation to different hosts (6).

**Figure 2.** Diagram of the T4 virion, based on electron microscopy at 2 to 3 nm resolution. Locations of proteins are indicated by the corresponding gene numbers (see (Fig. 1)). The portal vertex (gp20) is attached to the upper rings of the neck structure inside the head itself. An internal tail tube lies inside the tail sheath. The short fibers of gp12 are shown in stored, folded conformation. (Reproduced from (Ref. 80)).



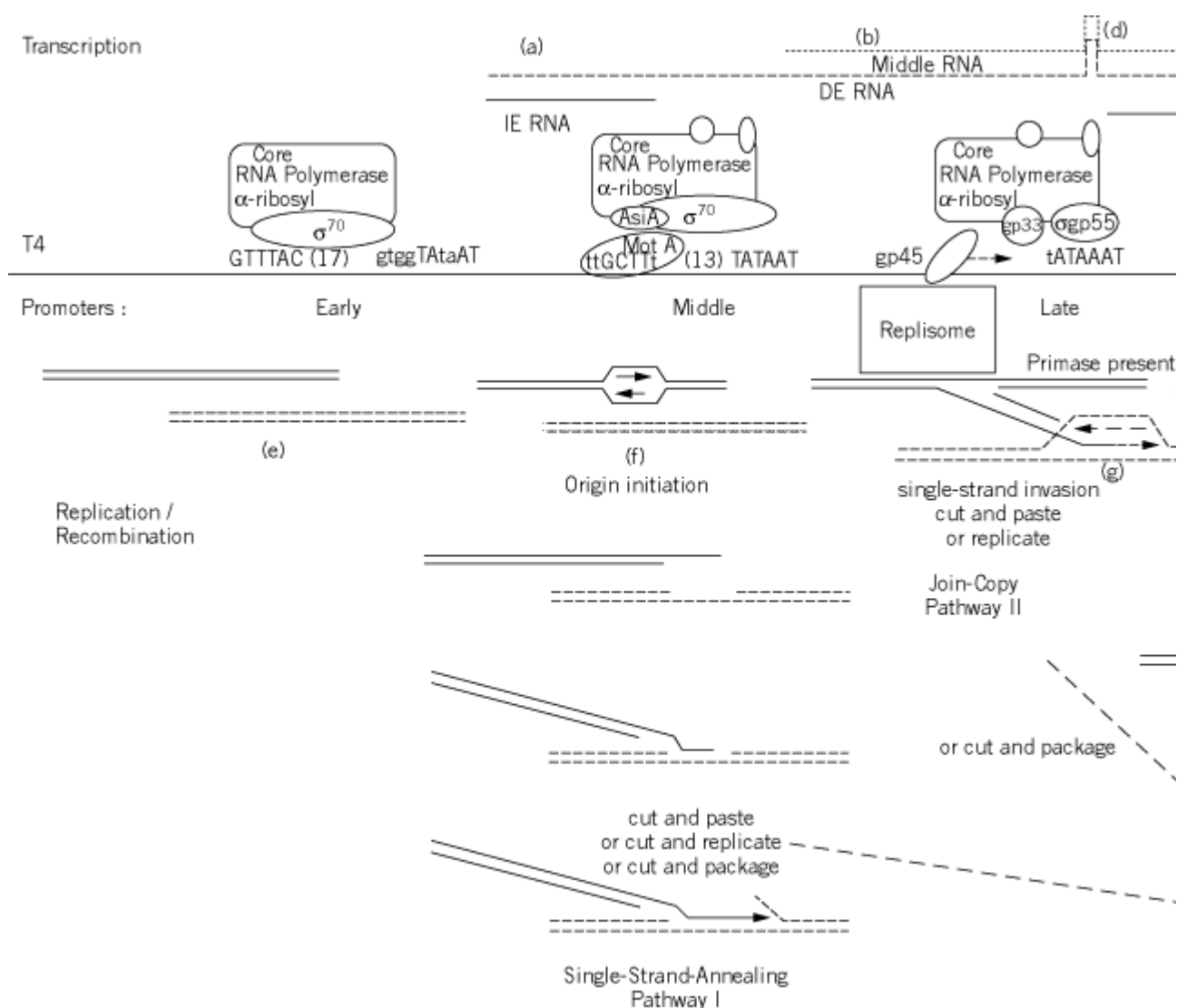
T-even phages are some of the most successful molecular parasites. Their genomes code for phage-specific [DNA replication](#), [recombination](#), and [DNA repair](#) functions and provide well-characterized genes and proteins for detailed analyses of these fundamental biological processes.

Nevertheless, like all viruses, T-even phages depend for their propagation on many vital structures and functions of their hosts, for example, [membranes](#), energy metabolism, [transcription](#), and [translation](#). They manage to usurp host structures and functions in an exquisite choreography that allows adaptations to different environmental conditions including different physiological states of the host.

Most of the recent work on T-even phages has been done with T4, mainly because the isolation of a large collection of [conditional lethal mutants](#) provided a powerful impetus for molecular analyses by biochemical and biophysical methods. The results of the combined efforts of many research groups that worked on T4 and related phages are summarized in a monograph (1), and comparisons with other members of this phage family have been reviewed (2-4, 7). Because of space limitations, the references cited here are restricted mainly to chapters in Ref. 1 and subsequent summarizing publications. The recipe for the success of T4 as a molecular parasite is based on multiple redundant pathways for most, if not all, physiologically important DNA transactions, which are interconnected

at many levels (Fig. 3). These redundancies and cross-connections allow flexibilities during development and evolution.

**Figure 3.** A diagram of the relationships of different transcription modes and pathways of DNA replication, recombination (progressing in time after infection from left to right). The upper panel (a–d) shows a stylized example of overlapping ea region of the T4 genome. Ribosome-binding sites for late proteins in prereplicative transcripts are sequestered in hairpin: panel indicates modifications of host RNA polymerase after infection and the consensus sequences of early, middle, and different forms of RNA polymerase. The lower panel shows different stages of DNA replication and recombination. (e) (f) Bidirectional origin initiation in one of them. Leading-strand synthesis is primed by RNA polymerase-generated transcripts by short RNAs synthesized by primase. As soon as the first growing point reaches an end (only one is shown), the partial homologous region of another chromosome (g) or the terminal redundancy of the same chromosome (not shown). “Join-involving 3’ DNA end. When an endonuclease cuts the invaded DNA at the backward recombinational junction at the bla recombination can be initiated from the cut 3’ end to allow copying of single-stranded segments of an invading DNA. (A arrowheads does not allow direct initiation of DNA replication but breaks the DNA.) Together, both processes generate l increasingly more complex by reiteration. This DNA is processed into mature, unbranched chromosomes during packagi lines, filled with different patterns. Newly synthesized DNA is drawn as thin lines. Discontinuous synthesis of Okazaki fragments continuous synthesis by solid lines. Arrows indicate the directions of RNA or DNA synthesis.



The gradual subversion of host functions to support different aspects of phage propagation is

achieved in many small steps:

1. A cascade of phage-induced proteins modifies the host **RNA polymerase** and its accessory [transcription factors](#) covalently and noncovalently, allowing sequential recognition of different classes of **promoters** and modulating the processivity of RNA polymerase to transcribe HMC-containing DNA selectively. Thereby, all host transcription is inactivated, and timed functioning of different classes of phage promoters is achieved. At least one of these proteins, an enzyme that ADP-ribosylates one a subunit of host RNA polymerase, is packaged into virions and injected with the phage DNA into the next host bacterium.
2. RNA processing by phage and host enzymes, [translational repressors](#), and still poorly understood modifications of [ribosomes](#), all modulate T4 gene expression. No transcriptional repressors are known, and it is surmised that posttranscriptional modulations can better adjust to rapid physiological changes during the short T-even development: One growth cycle is finished in less than 30 minutes at 37°C.
3. The onset of the first phage DNA replication from specific origins requires host RNA polymerase to generate primers and is thereby influenced by the physiological regulatory processes of the host. Most subsequent T4 replication depends on phage-encoded replication and recombination proteins and on DNA primers that are intermediates of homologous [recombination](#). This subsequent replication is entirely phage controlled.
4. The requirement of late transcription for the sliding clamp protein of replisomes, couples DNA replication and transcription.
5. During the later stages of development, proteins involved in packaging DNA also become important for DNA replication and repair and thereby coordinate these processes.

## 2. Genome Structure and Genetic Map

The genome of T4 resides in 168,903 bp of double-stranded DNA containing glucosylated HMC residues. The complex **modification** and **restriction** of T4 DNA can best be rationalized as the result of an ongoing evolutionary process that includes exchanges between the phage, its host, and prophages resident in the host.

T-even phages destroy dCTP, synthesize dHMCTP, and use the latter for DNA synthesis. This modification protects the T4 DNA against T4-encoded restriction endonuclease II that together with endonuclease IV and the gene 46/47-controlled **nuclease** degrades the host DNA as part of the parasitic strategy to usurp the host (8-10). It also makes phage DNA resistant to most type I, type II, and type III restriction enzymes. Although, HMC residues confer resistance to most I, type II, and type III restriction systems, they render DNA susceptible to the Mcr restriction systems of the host. These host functions were the first restriction systems (then called Rgl) discovered. They are now called McrA and McrBC, because they restrict DNA that contains methylcytosine or hydroxymethylcytosine. These enzymes are inactive, when the HMC residues are glycosylated. In T4 DNA, all HMC residues are modified; 70% with a- and 30% with b-glycosyl linkages. There are no a-glycosyltransferases, in T2 and T6 DNA and 25% of the HMC residues remain unglycosylated. T6 contains many diglycosylated residues. In addition, an early T4 antirestriction protein (Arn) protects nonglycosylated T4 DNA against McrA but not McrBC.

Phages T2 and T4, but not T6, encode a Dam **methylase** that methylates 0.5 to 1.5% of the adenine residues at the N<sup>6</sup> positions, mostly but not exclusively at GATC sequence. These enzymes exhibit patches of similarity at the protein level to the *E. coli* Dam methylase and the *DpnII* methylase of *Diplococcus pneumoniae*. The only proven physiological role of adenine methylation is protection against the phage P1 restriction system, when there is no HMC glycosylation (8).

Using genetic tricks, T4 mutants that contain unmodified cytosines or adenines in their DNA have

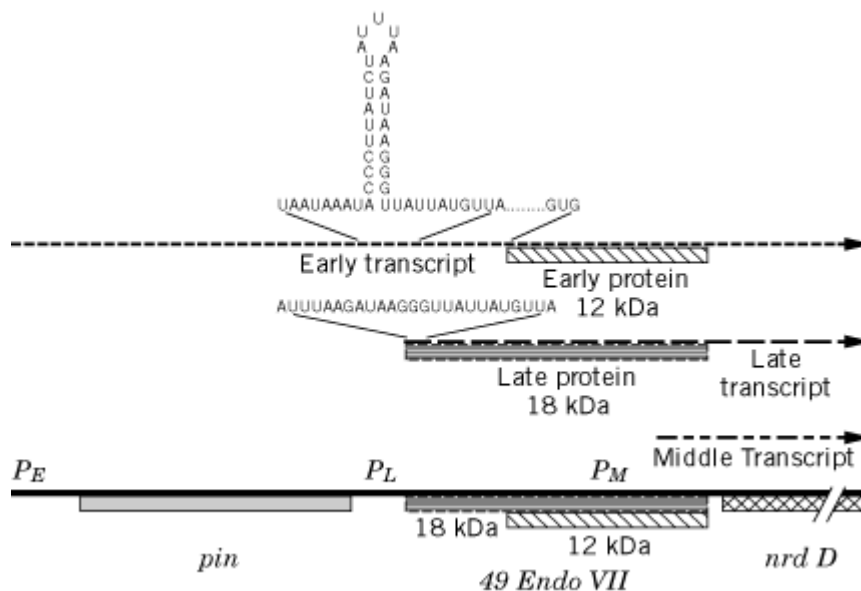
been isolated (8, 10). The cytosine-containing mutants have limited variability and host range. Their DNA has been instrumental in mapping and sequencing the T4 genome, but the lack of modification affects regulation of transcription termination (10, 11), and origin-dependent DNA replication (12, 13), and initiation of packaging (14) appear to be different from wild-type T4.

Mature T4 DNA molecules (chromosomes) packaged into virions are linear. In contrast, intracellular replicating DNA contains multiple covalently linked copies of the genome, called “concatemers,” that are highly branched as a result of recombination and replication. These are debranched and cut during packaging (see below) at nearly random (i.e., circularly permuted) map positions (1) and generate “terminal redundancies” at the ends so that 3 to 5% of the genome is diploid.

Numerous mutations and their assignments to complementation groups and open reading frames (ORFs) have defined 161 proteins (counting proteins initiated from gene-internal ribosome binding sites as separate proteins) that have known or deduced functions. Approximately 127 additional ORFs were revealed by DNA sequence analyses. Most, if not all of these are transcribed, and their protein products have been detected by acrylamide gel electrophoresis. However, the many functions of many of them are still unknown. There are many overlapping coding regions; several genes direct synthesis of more than one protein from in-frame internal start codons; many ORFs are very small; and in at least two regions, both complementary DNA strands encode proteins. In spite of the large genome size, only a few short regions have no coding capacity. These and other observations suggest that apparently redundant and “nonessential” genes confer no selective disadvantages.

In contrast to transcription units of many other (e. g., the **lambdoid phages** and **T7-related phages**, T-even early, middle and late transcription units are interdigitated. Most genes are transcribed from multiple promoters, gene for related or interacting proteins are not necessarily clustered, and in many cases early and late genes are cotranscribed (15) (see Figs. 1 and 4). In this and other respects, the T-even phages more closely resemble the **Herpes viruses** than certain other phages. Intriguingly, some T4 proteins (e.g., **DNA polymerase** and terminase) also resemble the corresponding proteins of Herpes viruses.

**Figure 4.** An example of interwoven transcriptional and posttranscriptional controls in the T4 gene *49* (Endo VII) region (see (Fig. 1). The locations of early ( $P_E$ ), middle ( $P_M$ ), and late ( $P_L$ ) promoters and the protein-encoding segments of genes *pin*, *49* and *nrdD* are marked. Overlapping transcripts are distinguished by differently patterned horizontal lines with arrows. The *nrdD* transcripts encode two proteins; one from spliced and the other (a truncated peptide) from unspliced RNA. The late gene *49* transcript is predominantly translated into the 18-kDa Endo VII. In the larger, early transcript, the Shine–Dalgarno sequence is sequestered in a hairpin. A shorter 12-kDa protein can be initiated in frame from an internal GUG. (Modified from Ref. 15.)



### 3. Temporally Controlled Gene Expression During T-Even Development

The classical Hershey–Chase experiment established unambiguously that T-even DNA is sufficient to establish viral growth (16). T4 adsorption and infection (17, 18) occur in steps, described below. Complete injection of phage DNA and a few packaged proteins through the membranes of host cells requires membrane potential. Inactivation of host functions and temporal regulation of expression of different classes of phage genes from the injected DNA and its copies are then exerted at several levels: **transcript** initiation, elongation, termination, and stability, **translational** controls, and combinations thereof (15). In contrast to T7-related phages, the T-evens use the host's core RNA polymerase throughout their life cycle. Differential transcription initiation from host or T4 early, middle, or late promoters is due, in part, to a cascade of RNA polymerase modifications: several accessory T4 proteins (1) associate noncovalently with the core RNA polymerase, and the  $\alpha$ -subunits are covalently **ADP-ribosylated**.

Different classes of phage genes can be distinguished in terms of timing as early (also called immediate early, IE), middle (also called delayed early, DE), or late. Another classification criterion distinguishes all genes that are expressed before the onset of DNA replication as “prereplicative” or “early” from “postreplicative” or “late” genes, whose expression depends on DNA replication. However, the regulation of T4 gene expression resembles more a web of interacting regulatory networks than a cascade of linear timed pathways.

Operationally, early (IE) genes are defined as transcribed by host RNA polymerase when phage gene expression is prevented by inhibitors of protein synthesis (eg, [chloramphenicol](#)). Under these conditions, the host RNA polymerase uses unmodified  $s^{70}$  to recognize early T4 promoters, but the core RNA polymerase is altered by T4 Alt protein that is packaged and coinjected with the phage DNA. Alt protein ADP ribosylates one  $\alpha$ -subunit at Arg 265, a position that is important for recognition of so-called UP elements in promoter DNA (19) and for the dimer configuration of the carboxy-terminal domains of RNA polymerase  $\alpha$ -subunits (R. Ebright, personal communication). This ADP-ribosylation is beneficial, but not essential for T4 growth.

Expression of all other T-even genes requires phage-directed protein synthesis for several reasons discussed below. The distinction between different classes is, however, blurred because most T4 genes are under dual or multiple transcriptional controls and many late T4 genes are transcribed from late promoters, but are also cotranscribed with early genes from early or middle promoters. For these

late genes, posttranscriptional mechanisms (discussed below) prevent early expression.

Collectively, prereplicative genes encode i) **nucleases** that degrade the host DNA; ii) enzymes of the deoxyribonucleotide biosynthesis complex (see [Ribonucleotide Reductase](#)); iii) proteins of the DNA replication, recombination, and repair machines; iv) proteins that modify the T4 DNA to protect it from degradation by its own nuclease and from other restriction enzymes; v) several [transfer RNAs](#) that are processed from precursor RNAs and supplement host tRNAs during translation; vi) proteins that modify the structure and function of the host RNA polymerase; vii) at least one ribonuclease (RegB protein) that selectively destroys certain early transcripts; and viii) a [translational repressor](#) (RegA protein). In addition, some prereplicative transcripts serve as primers for leading strand DNA synthesis in **origins of replication** (see text below).

The late genes code for all virion components, for some DNA recombination and DNA repair proteins (defining a late T4 “join-cut-copy” DNA replication, recombination and repair pathway described below), and for proteins that cut and package the complex vegetative DNA into preformed heads. A soluble lysozyme (gp e) lyses the cell wall of host bacteria to release the progeny phage particles. It is paralogous to a baseplate lysozyme (gp5) that attacks the cell wall from the outside during infection (20). The gp e lysozyme gains access to its substrate in the outer cell wall through a protein channel that forms from the holin (T protein) that causes the membrane potential to break down. In T4, this process is regulated, in part, in response to superinfection sensed by the RI protein (17, 21, 22).

T4 early promoters are recognized by the *E. coli* RNA polymerase that contains the major sigma factor  $s^{70}$ , when the host DNA is still largely intact. These promoters resemble the **consensus sequence** of *E. coli* promoters, but most of them also have some differences and some additional information content (23) (Fig. 3a). Many of them have upstream poly A or poly T tracts that may enhance transcription as bendable sequences or as UP elements or both.

Several factors, none of them essential for T4 growth, have been proposed as contributors to the preferential transcription from T4’s early promoters compared with *E. coli* promoters (24):

1. Most importantly, once the early T4 Alc protein is made, it inhibits all transcript elongation on (cytosine-containing) host DNA.
2. Transcription from various T4 or *E. coli* promoters is affected differentially by ADP ribosylation of Arg<sup>265</sup> (in the C-terminal activation and dimerization domain) of one  $\alpha$ -subunit of the host’s RNA polymerase by the T4 Alt protein, as mentioned above, and subsequent ADP ribosylation of Arg<sup>265</sup> of the other  $\alpha$ -subunit by the early T4 ModA protein (25, 26).
3. When infection occurs, the host DNA is associated with nonspecific (eg, HU, H-NS) or semispecific (eg, IHF, FIS) DNA binding proteins. In contrast, the infecting phage DNA is at first largely free of proteins and may be more readily accessible to the host’s RNA polymerase (24).
4. The T4 Ndd protein, a DNA-binding protein, disrupts the *E. coli* nucleoid (27, 28).
5. The early T4 AsiA protein binds to the C-terminal segment of  $s^{70}$  (29) and interferes with transcription from host promoters that contain consensus –35 regions (30, 31).
6. The host DNA is relaxed and degraded by early T4 proteins.

Both early and middle promoters require host  $s^{70}$ . Binding of AsiA protein to  $s^{70}$  prevents transcription of early T4 promoters with “standard” –35 regions, but infection can proceed because most prereplicative genes are also transcribed from middle promoters (11) (Fig. 3b). At middle promoters, the RNA polymerase, whose  $s^{70}$ -subunit has acquired AsiA protein, recognizes consensus –10 regions and T4 MotA protein bound to consensus DNA sequences called *motA* boxes (11, 32, 33) (Fig. 3). Transcriptions from early and middle promoters can overlap in time, another

example why linear classification schemes are oversimplified.

RNA polymerase is directed to late promoters by a phage-encoded sigma factor, gp55. Initiation from late promoters (34) (Fig. 3), requires, in addition, the adapter protein gp33 and the **sliding clamp** of the replisome, gp45, which is loaded onto the DNA at single-strand interruptions by a gp44–gp62 complex (see [DNA Replication](#)). This sliding clamp tracks along DNA and allows RNA polymerase, bound at late promoters, to form open complexes. The requirement for the sliding clamp couples late T4 transcription to DNA replication. In certain recombination — and repair — defective mutants, late transcription can occur without DNA replication because single-strand interruptions in the DNA provide entry sites for the sliding clamp (34).

Even when some genes, are located downstream of a promoter, they may be poorly expressed because of i) premature transcriptional termination, ii) RNA processing or degradation, or iii) [ribosome](#) binding sites are sequestered by secondary structures (see Attenuation). Certain mutations that affect the transcriptional termination factor Rho (also called host defective, *hdF* or *nusD*) of *E. coli* prevent growth of wild-type T4 by causing premature termination of many T4 transcripts (11). Mutations that allow growth in these *rho* mutants have been found in three nonessential T4 genes. It was initially thought that these encode transcriptional **antiterminators**, but current evidence suggests that at least one of them, *goF*, allows T4 growth in the (*nusD*) *rho* mutants, because it stabilizes the few functional transcripts that have not been prematurely terminated. Plasmid-encoded Rop protein has similar stabilizing effects, and pBR322 plasmids allow T4 growth in these **rho** mutant hosts (35, 36).

Several additional nonessential T4 proteins also associate with or modify host RNA polymerase subunits (Table 1). Some of them affect the binding of minor host  $\sigma$  factors to core RNA polymerase and prevent competition with  $\sigma^{70}$  that is essential for T4 development. For example, the Mrh protein modulates phosphorylation of the host's  $\sigma^{32}$ , which affects T4 growth at 42°C (37, 38).

**Table 1. Phage T4-Encoded Proteins That Modify Structure or Function of the Host RNA Polymerase<sup>a</sup>**

| Modification   | Product of Gene | Molecular Mass (kDa) | Time of Synthesis   |
|--|-----------------|----------------------|---|
| ADP ribosylation of one $\alpha$ -subunit at Arg <sup>265</sup>                        | <i>alt</i>      | 76.8–75.9            | Late, packaged into heads, cleaved in the head, and injected with DNA |
| ADP ribosylation of second or both $\alpha$ subunits at Arg <sup>265</sup>             | <i>modA</i>     | 23.3                 | Early   |
| ADP ribosylation of ribosomal protein S1 and yet undetermined proteins                 | <i>modB</i>     | 24.2                 | Early   |
| Association with ternary transcribing complex, causing termination on C-containing DNA | <i>alc</i>      | 19.1                 | Early   |
| Modulation of host $\sigma^{70}$   | <i>asiA</i>     | 10.6                 | Early   |
| Activation of middle   | <i>motA</i>     | 23.6                 | Early   |

|  |             |      |                  |
|--|-------------|------|------------------|
| promoters, binds to DNA and to host s <sup>70</sup>          |             |      |                  |
| Modulation of phosphorylation of host s <sup>32</sup>        | <i>mrh</i>  | 18.5 | Early            |
| Similarity to (decoy of?) host s <sup>32</sup>               | <i>srh</i>  | 8.1  | Early            |
| Similarity to host s <sup>70</sup>                           | <i>srd</i>  | 29.0 | Early            |
| Binding to RNA polymerase s factor for late promoters        | <i>rpbA</i> | 15.5 | Early and middle |
|  | 55          | 21.5 | Middle           |
| Bridge between $\sigma^{8P55}$ and sliding clamp             | 33          | 12.8 | Middle           |
| Binding to some late promoters                               | <i>dsbA</i> | 10.4 | Middle           |
| Activator of late promoters, sliding clamp of DNA polymerase | 45          | 24.8 | Early and middle |
| Clamp loader   | 44          | 35.8 | Early and middle |
| Clamp loader   | 62          | 21.4 | Early and middle |

<sup>a</sup> Refs. [15](#) and [38](#).

#### 4. Posttranscriptional Controls

Due to the interspersion of early and late genes on the T4 map (Fig. 1), many early and middle transcripts are extended into late genes (Figs. 3 and 4). Nevertheless, few corresponding late proteins are synthesized early below 37°C. There is some early expression of these late genes *in vivo* at higher temperatures. It is remarkable that expression of at least ten such late genes investigated so far follows similar pattern: in the long early transcripts, a hairpin sequesters the translation initiation region, either the [Shine–Dalgarno sequence](#) or the [start codon](#), or both. A late promoter, immediately upstream of these late genes, directs synthesis of transcripts that cannot from the hairpin; these late transcripts are efficiently translated (Figs. 3d and 4). As an example, expression of gene 49 coding for endonuclease VII is depicted in Figure 4. In this case, a shorter peptide, initiated from an internal GUG, plays a role in join-cut-copy replication/recombination ([39](#)), discussed below. The late T4 lysozyme gene *e* is similarly transcribed early and late, but mainly translated late. These late genes can be translated early *in vivo* at high temperatures, or *in vitro*, when RNA is partially degraded, because under these conditions the ribosome binding sites become accessible. I speculate that the corresponding position of late promoters provides some selective advantage because if RNA polymerase that enter at the early or middle promoters, pauses and terminates near the hairpins formed in early transcripts, it might facilitate acquisition of T4-encoded s<sup>gp55</sup> and transcription initiation from the late promoters, which reside in each case in the corresponding palindromic DNA sequence.

Three additional regulatory systems of T4 depend on translational controls ([40a](#)). Two genes are autogenously regulated by translational repression: gene 32, which codes for the major single stranded DNA binding protein (SSB) involved in DNA replication, recombination and repair, and gene 43, encoding DNA polymerase. A more general, nonessential translational repressor (RegA protein), which binds to several specific transcripts, reduces translation of several T4 replication proteins and of some host proteins. This repression is most apparent under experimental conditions



that prolong early transcription.

Translation of certain early transcripts is prevented at later times by a T4 *regB*-encoded ribonuclease, which selectively cleaves these transcripts at the ribosome binding sites.

## 5. DNA Replication and Recombination

The combination of genetic experimentation and virtuoso biochemical and biophysical characterization of replicative proteins has led to an understanding of the functions and interactions of replicative proteins in the basic replisome, a biological machine that moves the replicative fork (or through which replicating DNA passes) (see [DNA replication](#)).

The proteins that assemble into basic replisomes in all organisms share sequence and structural similarities, and some of the T4 replication proteins can function partially in eukaryotic *in vitro* systems ([41](#)).

Seven T4 proteins, corresponding to genes *43* (DNA polymerase), *44* and *62* (sliding clamp loader), *45* (sliding clamp), *41* (DNA helicase), *61* (primase to synthesize primers for Okazaki fragments), and *32* (single-stranded DNA binding protein) together replicate model templates at *in vivo* speeds. In wild-type T4, leading and lagging strand synthesis are coupled ([42](#)), but they can be uncoupled *in vivo* or *in vitro* ([43](#)), for example, in primase- or **topoisomerase**-deficient mutants ([39](#)). T4 DNA polymerase has been well characterized by exquisite mutational, biochemical, and biophysical methods. Mutator and antimutator mutations have delineated domains that are important for polymerizing and for proofreading activities that contribute to the exceptionally high fidelity of this enzyme ([44-46](#)). Mutated sites in replisome components can now be correlated with the 3-D structure of the T4 and related RB69 DNA polymerase, the sliding clamp (gp45), and a large segment of the single-stranded DNA binding protein gp32 ([47-50](#)). Dynamic interactions between these proteins have been dissected in exquisite detail (see Refs. [42](#) and [51](#) for recent reviews).

Gene expression, DNA replication, and recombination are tightly interwoven. Replication and recombination follow several different redundant pathways, shown diagrammatically in Fig. [3](#) ([39](#), [52-54](#)). The functioning of each pathway in time depends in part on timed expression of the genes required for each pathway. Redundant alternative modes of replication and recombination ensure that both processes work under many different conditions and during different stages of development.

The first round of DNA replication is initiated from one of several potential origins. Because of the circular permutation of chromosomal ends (indicated in Fig. [3e](#)), any origin in each individual chromosome is located at a different distance from the DNA termini. Only one origin is used, in most chromosomes perhaps because at first, there are limited supplies of replisome components. Each of the four origins that have been closely investigated has a different sequence. Three origins (*A*, *F*, and *G*) require transcription from middle promoters ([13](#)), whereas *ori E* uses an early promoter ([55](#), [56](#)). Figure [3f](#) depicts initiation from an arbitrary generic origin. Transcripts initiated from these promoters serve as primers for leading-strand DNA synthesis. (They also encode proteins, but priming and directing protein synthesis are mutually exclusive.)

Primase then synthesizes primers for Okazaki fragments ([57](#)). Transition from RNA to DNA synthesis occurs at several sites within a region approximately 1 kb downstream of each origin promoter ([55](#), [56](#), [58](#)).

The transition from  $s^{70}$ -dependent prereplicative transcription to T4  $s^{gp55}$ -dependent late transcription inhibits initiation from these replication origins, because the modified RNA polymerase no longer recognizes origin promoters ([56](#), [59](#)), and/or because the late UvsW protein, an RNA–DNA helicase, actively unwinds potential primers from the DNA template ([54](#), [60](#)).

Subsequent DNA replication is initiated from intermediates of homologous recombination, by a mechanism first compellingly demonstrated in phage T4 (59), (Fig. 3g and i). Therefore, the T4 recombination genes (Table 2) are also important for DNA replication in vivo, and recombination-deficient T4 mutants arrest DNA replication prematurely (13, 53). It is now apparent that there are several ways by which recombination intermediates can prime DNA synthesis (39, 52). Such intermediates can be formed by annealing complementary single strands or by invasion of a single-stranded terminus into the homologous region of another molecule (or at the other end of the same molecule). Both structures can initiate replication using the 3' end of the invading or annealing single strand (join-copy recombination) (Fig. 3g). In addition, invasion intermediates can initiate replication from a 3' end in the invaded strand after an endonuclease or a damaging agent has broken the invaded strand at the branch-migrating junction (join-cut-copy recombination) (Fig. 3i). The join-copy mode can start as soon as a growing point has reached an end. The join-cut-copy mode occurs usually with some delay, because in the absence of DNA nicks imposed by damaging agents, it depends on prior synthesis of late endonucleases. Because only the join-cut-copy mode can bypass the requirement for primase or topoisomerase in T4 DNA replication, mutants defective in these enzymes appear to have a delay in recombination-dependent DNA replication. Radiation and other agents that cause single-strand nicks or double-strand breaks can enhance recombination and reduce the delay in DNA synthesis of primase or topoisomerase mutants.

**Table 2. T4 DNA Replication and Recombination Proteins**

| <b>Protein</b>                                  | <b>Gene</b>  | <b>DNA phenotype of Mutants<sup>a</sup></b> | <b>Time of Expression</b> |
|---|--------------|---|---------------------------|
| DNA polymerase                                  | 43           | D0  | Early and middle          |
| Sliding clamp                                   | 45           | D0  | Early and middle          |
| Clamp loader                                    | 44 & 62      | D0  | Early and middle          |
| Primase   | 61( = 58)    | DD  | Early and middle          |
| DNA helicase                                    | 41           | DA  | Early and middle          |
| DNA helicase                                    | <i>dda</i>   | DD  | Early                     |
| Loader of gene 41 helicase                      | 59           | DA  | Early and middle          |
| Single-stranded DNA binding protein             | 32           | DA  | Early, middle, and late   |
| dNMP kinase                                     | 1            | DS  | Middle and late           |
| dCMP hydroxymethylase                           | 42           | DS  | Early                     |
| dCTPase, dUTPase, dCDPase, dUDPase              | 56           | DS  | Early and middle          |
| 5' to 3'Exonuclease, RNase H                    | <i>rnh</i>   | WT  | Early                     |
| DNA topoisomerase                               | 39 & 52 & 60 | DD  | Early and middle          |
| DNA ligase                                      | 30           | DA  | Early                     |
| Recombination protein and nuclease <sup>b</sup> | 46 & 47      | DA  | Early and middle          |
| RecA-like recombination protein                 | <i>uvsX</i>  | DA  | Early and middle          |

|                              |                        |                 |                                   |
|------------------------------|------------------------|-----------------|-----------------------------------|
| Helper of UvsX protein       | <i>uvsY</i>            | DA              | Middle and late                   |
| DNA and RNA-helicase         | <i>uvsW</i>            | WT <sup>b</sup> | Late                              |
| Endonuclease VII             | <i>49</i>              | WT <sup>b</sup> | Early (weak) and<br>Late (strong) |
| Packaging terminase          | <i>16 &amp; 17</i>     | WT <sup>b</sup> | Late                              |
| Membrane-associated proteins | <i>rIIA &amp; rIIB</i> | WT <sup>b</sup> | Early and middle                  |
| Endonuclease II              | <i>denA</i>            | WT <sup>b</sup> | Early and middle                  |
| Exonuclease A                | <i>dexA</i>            | WT <sup>b</sup> | Early                             |

<sup>a</sup> D0, no DNA synthesis; DA, DNA arrest; DD, recombination-dependent DNA arrest; DD, recombination-dependent DNA replication delayed; DS, slow DNA synthesis; WT, normal DNA synthesis, except under special conditions.

<sup>b</sup> DNA synthesis is altered in combination with other phage or host mutations.

Ultimately, reiteration of recombination-dependent DNA replication generates a highly branched concatemeric DNA network in which no individual chromosomes can be distinguished.

Of course, not all recombination junctions need to be converted to replication forks. T4 recombination can occur, albeit later than usual, when DNA replication is inhibited. Electron micrographs of such T4 DNA intermediates provided the first compelling evidence for the importance of branch migration in homologous recombination (61). Under these conditions, no viable progeny particles are produced, because no packageable concatemers are formed, there is little late transcription (which depends on DNA replication), and thus there are no heads that can be filled.

Recombination-dependent DNA replication *in vitro* has been achieved for some time (13, 62). Consistent with genetic analyses, reactions require several recombination proteins in addition to the seven basic replicative proteins just mentioned. T4 UvsX (a **RecA** homologue) and T4 UvsY (a mediator, that binds to single- and to double-stranded DNA and is needed only at low UvsX concentrations) promote strand invasion. Gp59 facilitates loading of the replicative gp41 helicase, which is also important in driving branch migration of recombining DNA. These proteins also interact with the major T4 ssDNA binding protein, gp32. *In vitro* recombination-dependent replicative reactions are primed by 3' ends of single-stranded DNA invading homologous duplex DNA. With an excess of UvsX protein, and in the absence of primase or topoisomerase, they result in conservative DNA replication called “bubble migration,” in which the nascent DNA strand is extruded from the DNA template, much like nascent transcripts are extruded by RNA polymerase (62). This process has been postulated to provide the single-stranded DNA for homing of **introns** by a recombination model dubbed the synthesis-dependent strand-annealing model (63, 64). There is no direct evidence that bubble migration occurs *in vivo* at physiological concentrations of UvsX protein, primase, and topoisomerase (reviewed in 13). Moreover, these *in vitro* systems do not depend on gps 46 and 47, which are required *in vivo* for recombination-dependent DNA replication (59).

Recently, *in vitro* initiation has been accomplished from a plasmid-borne T4 *ori F* (*uvsY*) (58). A short transcript, hybridized to supercoiled plasmid DNA *in vitro*, primed leading-strand DNA synthesis in this system. This priming occurred at a position different from the RNA–DNA transition observed at this origin *in vivo* (55, 56).

Many of the T4-encoded DNA enzymes, most importantly DNA ligase, polynucleotide kinase, DNA polymerase, and ssDNA binding protein, are now standard components of cloning and sequencing

procedures and kits.

## 6. DNA Packaging

Branched, concatemeric, intracellular T4 DNA can be debranched by T4 endonuclease VII (gp49) and by a terminase, a heteromeric protein encoded by genes *16* and *17*. These proteins associate with DNA and with gp20 at the portal vertex of the head to form a “packasome” that drives DNA into performed heads (65, 66) by an ATP-dependent process. Gene *17* produces several proteins of different sizes by initiation from in-frame internal initiation codons (67). At least two of these have nuclease activities that can be abolished by mutations (68, 69); the largest also binds to folded single-stranded DNA segments (68). Packaging of nonmodified dC-containing DNA starts preferentially at apparent *pac* sequences located in the overlap of genes *16* and *17* and duplicated in gene *19* (66). Other evidence indicates that T4 initiates packaging of HMC-modified DNA at random positions of the genome at folded single-stranded DNA segments (68), suggesting that packaging, like DNA replication, is also initiated by redundant mechanisms. Processive packaging of 103 to 105% genome lengths to fill the preformed heads generates the random circular permutation of the ends in mature virion DNA. Endonuclease VII, which cuts **Holliday junctions**, Y-junctions, and mismatched base pairs *in vitro*, can trim the branches of vegetative DNA *in vivo* and *in vitro* (70). This enzyme binds *in vitro* to the portal protein gp20 (71). DNA ligase, endonuclease V, and topoisomerase are also required for packaging, presumably to ascertain that any lesions and single-stranded segments in DNA are repaired before the DNA is packaged (65).

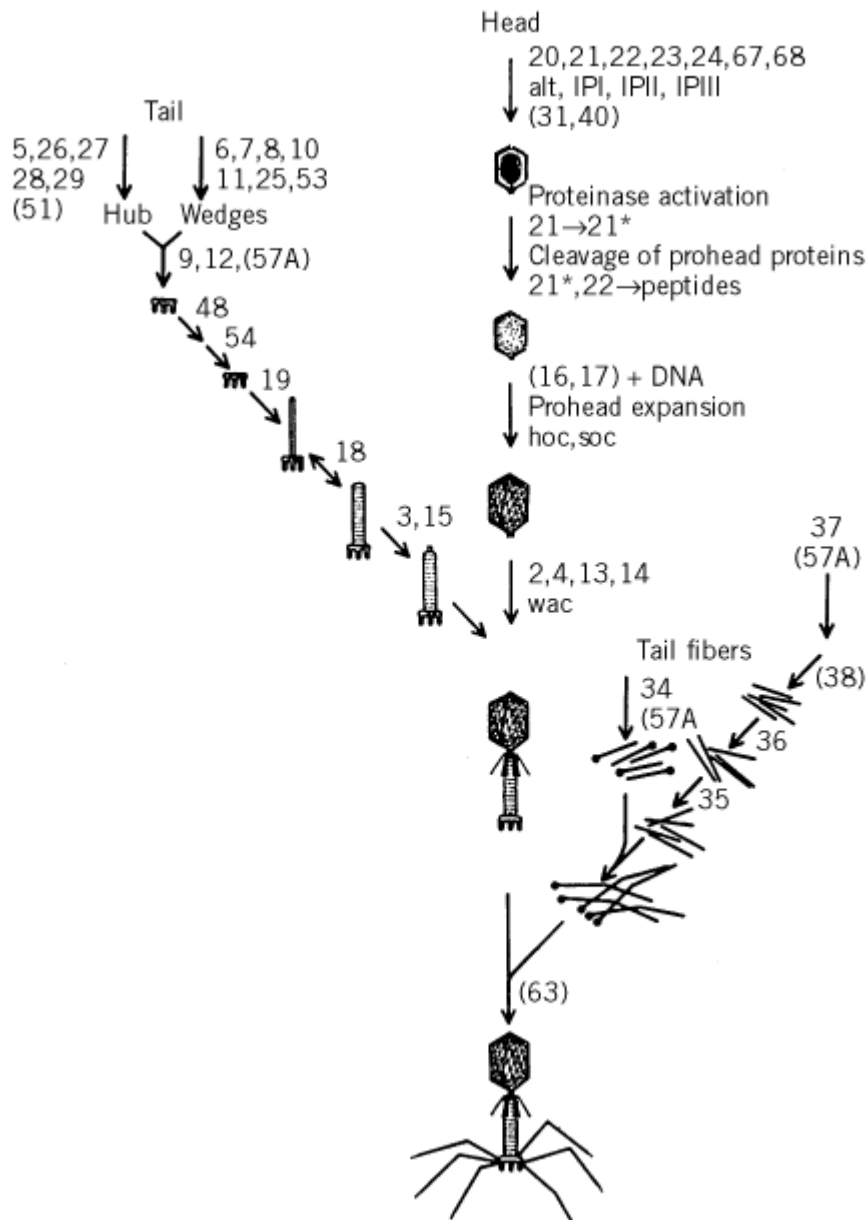
Both endonuclease VII and the nuclease activity of gene *17* products play important roles in the late, join-cut-copy, recombination-dependent DNA replication described above (39, 52). Thus, these activities can link DNA replication and packaging at a physiological level.

## 7. Virion Structure and Assembly

T-even phages build some of the most complex virus particles, which resemble lunar landing modules (Fig. 2). They devote more than 40% of their genetic information to synthesizing and assembling the protein components of these particles. The importance of host functions for assembly became apparent from “host defective” (*hd*) mutations and compensating phage mutations and led to the discovery of the first **chaperonin**, now called GroEL. The small subunit of this chaperonin in *E. coli*, GroES, is not necessary for T4 growth, because T4 gene *31* encodes a similar protein that is better adapted to interact with the major capsid protein (72).

Exquisite mutational and biochemical analyses revealed the complex assembly pathways (Fig. 5) that accommodate efficient packaging of DNA and efficient genome transmission (65). This system can also be used to package and deliver large pieces of foreign DNA (73).

**Figure 5.** Pathways of T4 virion assembly. Gene numbers in parentheses correspond to assembly factors that are not part of the final structure. The identities of gp26, gp28, and gp51 are uncertain. (Reproduced from Ref. 80 with permission.)



Twenty-four genes are involved in head morphogenesis, more than 25 encode structural components of the tail and tail fibers, and five more are needed as morphogenetic catalysts. Heads, tails, and tail fibers are assembled independently and put together after the heads are filled with DNA (Fig. 5).

Protein components of heads, tails, and baseplates and shapes of the protein subassemblies are similar among the T-even phages, except that one or two nonessential proteins (Hoc, Soc) that decorate and stabilize T4 heads are missing in some members of the T-even family.

*In vivo* assembly of heads is initiated on the bacterial membrane from the portal protein (gp20). The major head subunits (gp23) assemble around a scaffolding core, and gp 24, which is related to gp 23, forms pentamers at the vertices. Both gp23 and gp24 have to be converted into an assembly-competent state by the chaperonins mentioned above. Exquisite cooperative conformational changes lead to changes in the relative positions of proteins while the head precursors are being assembled. Subsequently, a phage-encoded **proteinase**, which is derived from gp21, one of the scaffolding proteins, cleaves each head subunit, the scaffolding proteins, and other packaged proteins, eg, the ADP-ribosylating Alt protein (see Gene Expression). Thus, head assembly is irreversible. Proteolysis is probably triggered *in vivo* by closure of the head vertices and perhaps by entry of some DNA into

the head; it can be achieved *in vitro* by other conditions. In any case, subsequent rearrangement of the subunits that lead to expansion of the head volume is triggered by the entry of DNA (74). The T4 heads are then stabilized by decoration with the Hoc and Soc proteins mentioned above. The complete assembly pathway has resisted *in vitro* reconstitution (65).

The normal, elongated T4 heads can be described by skewed **icosahedral symmetry**, in which the bottom and top triangles have a triangulation number  $T = 13$ , and the faces are elongated with a Q number  $Q = 21$ . These heads can accommodate full-length chromosomes with terminal redundancies. Anomalous heads of different sizes and shapes are made in low proportions (less than 5%) by wild-type T4, and in higher proportions by certain mutants or in the presence of an arginine analog, canavanine. Head lengths and shapes are determined at an early assembly stage, before formation of the unprocessed proheads.

Small heads (sometimes called “petites”) contain incomplete genomes, which represent nearly random permutations of the genetic map. Most of these can be depicted as isometric (not elongated as the normal head) icosahedra and have a predicted volume that matches the measured DNA content (two-thirds of normal length). After infection with two such incomplete genomes, recombination reconstitutes complete genomes in one third of the doubly infected bacteria at the frequency expected if they represent random circular permutations of the map (75, 76). Heads intermediate in size between the isometrics and normal elongated ones contain intermediate sized DNA molecules (77, 78). Larger heads, called “giants,” contain oversized chromosomes that represent linear concatemers (79). A small proportion of all heads contains more than one DNA molecule; some heads are asymmetrical and contain more than one portal (80). The anomalous head sizes and shapes occur because assembly depends on numerous interactions of several proteins, each of which can exist in multiple conformational states and in different concentrations.

The base plates are built like exquisite cocked mechanical devices that lock the DNA into the virion and serve as valves that allow DNA to exit during infection. They consist of a central hub, six outer wedges, and six tail spikes; each of these substructures is assembled in independent pathways. The tails have a remarkably uniform length, determined by a subunit of the base plate that acts like a tape measure. Tails have a tubular inner core, through which DNA passes from the head to the base plate, and an outer sheath that contracts during infection by concerted conformational changes of individual subunits. Six bent tail fibers are attached to each tail. Each fiber consists of two half-fibers whose proteins are joined at an angle. The inner (proximal) half-fibers are attached to the base plate. In many newly formed phage particles, the other end is attached to short fibers that emanate from the junctions of heads and tails.

These complex assembly pathways, it is thought, facilitate the stability of the free virion while allowing subsequent conformational changes in response to contacts with the next host bacterium, as a prelude to DNA injection. The outer (distal) half-fibers contact phage-specific receptors on the surface of the bacterial cell wall during the first reversible step of adsorption. These contacts facilitate a second irreversible step in which the base plate contacts different receptors. Then, a concerted conformational change of all base plate components from a hexagon to a star configuration opens a hole in the tail to allow DNA to exit from the particles. This transition also activates contraction of the tail sheath and cleavage of gp5, a base plate protein, to release a lysozyme that actively digests host cell wall components from the outside.

Heads, tails, and tail fibers, made in different bacteria after infection with appropriate mutants, can be assembled *in vitro* to infectious particles, into which DNA can be packaged *in vitro*, allowing packaging of foreign DNA.

## 8. Restriction and Modification

In the following discussion, the term “restriction” is used in its broadest meaning and is not limited to type II [restriction enzymes](#). The complex modifications and restrictions of T4 DNA and restriction

of other DNA by T4 can be best rationalized as a result of an ongoing evolution that is accelerated by strong selective forces and exchanges between the genomes of T4 and its host, especially plasmid and prophage genomes that reside in the host.

As mentioned before, HMC-containing DNA is protected from the T4-encoded nucleases that degrade host DNA. However, it is susceptible to McrA and McrBC restriction systems, unless it is glycosylated (81). McrBC is a GTP-dependent enzyme that moves along DNA (82), and assembles as rings (83). McrA of *E. coli* K12 resides in a cryptic prophage-like element, *e14*, that is not present in all *E. coli* strains (84). It is counteracted by the T4 protein Arn. Such restriction systems provide powerful selective advantages for the different DNA glycosylating enzymes of various T-even phages.

Lit, another protein also found in bacteria containing *e14*, cleaves the host's [elongation factor](#) EF Tu, if and only if it is directed to its target by a short internal peptide (*gol*) of the major T4 head protein, gp23. Thus, in *e14*-containing bacteria, all late T4 protein synthesis is inhibited (85, 86).

The classical example of phage exclusions by genes of resident prophages is that of T4 *rII* mutants by the *rexA* and *rexB* genes that are expressed in [lambda phage](#) lysogens. This exclusion was elegantly used in Benzer's classical analyses of structure and function of a gene. It occurs at the time of transition from join-copy to join-cut-copy recombination discussed under DNA replication and recombination, and it involves the late T4 endo VII that is important in the join-cut-copy pathway (52). The molecular mechanism of this exclusion is still unknown. It has been proposed, without direct experimental evidence, that the RexB proteins form an ion channel that, when opened after infection with T4 *rII* mutants, leads to cessation of host metabolism (87).

Another cryptic DNA element of certain *E. coli* strains, *prc*, encodes a PrrC protein that excludes mutants deficient in T4 [RNA ligase](#) or **polynucleotide kinase**. PrrC protein is a cryptic RNA endonuclease that is activated by the small (26 residue) T4 Stp protein to cleave the anticodon loop of an essential host tRNA<sup>Lys</sup>. The T4 RNA ligase can repair this damage, but in the absence of RNA ligase, the cleavage of this tRNA is lethal to T4 protein synthesis. Intriguingly, the *E. coli prcC* gene is located between three genes of a type IC restriction cassette. The corresponding proteins are thought to inhibit PrrC RNase activity in uninfected cells.

Phage P2 lysogens exclude T4 by two mechanisms: the Tin protein poisons the single-stranded DNA binding protein gp32 (except in certain T4 gene 32 mutants) that is essential for all T4 DNA replication and recombination (88), and the P2 Old protein can degrade T4 and lambda DNA from ends, nicks, and gaps (89), unless they are protected by certain proteins.

## 9. Evolution

The sequences and map positions of homologous "essential" genes of most T-even phages are similar. In contrast, different members of the family have different "nonessential" genes interspersed between the homologous essential genes (Fig. 1). The differences contribute to apparent exclusions of alleles of one phage by another, and to the species barriers between different members of the family. The different "nonessential" genes became evident first as insertion or substitution loops in electron micrographs of heteroduplex DNA prepared *in vitro* by annealing single strands of T2, T4, and T6 DNA (90) and have been confirmed by sequence comparisons in many cases (2, 3, 39, 91). Although inactivating these nonessential genes has little or no consequence for phage growth in the laboratory, these genes are known or suspected to be important under different physiological conditions, in different hosts, and in the face of various restrictive systems imposed by different hosts.

In some cases, the sequence divergence reflects gene amplifications, inversions, and/or permutations of duplicated sequences (7, 92). In other cases, illegitimate pairing of partially homologous

sequences, join-copy and join-cut-copy recombination, and partial repair of mismatched heteroduplexes (discussed above) are proposed to lead to horizontal gene transfer. For example, in different T-even phages, different nonessential genes were acquired adjacent to the essential dCTPase gene, generating different apparent mutations in the corresponding dCTPase genes (39, 91). This mechanism can explain variabilities of other phage genes as well, for example, differences in the DNA polymerase genes of T4 and RB69, differences in the *e* lysozyme and the base plate lysozyme genes of phage T4, and similarities of base plate gene of T4 to phage P2 (93). Moreover, substitutions of sequence blocks of individual tail fiber genes by foreign sequences can account for the differences in tail fibers and host range of different members of the family and for the remarkably rapid coevolution of viral and host genomes (6).

## 10. Perspectives

The T-even phages have been instrumental in the initial formulations of many fundamental biological concepts: the unambiguous recognition of nucleic acids as genetic material; the operational distinctions in defining the gene by mutational, recombinational, or functional analyses; the demonstration of messenger RNA; the nature of the genetic code; light-dependent and light-independent DNA repair mechanisms; restriction and modification of DNA; the presence of introns in prokaryotes, and the facility for self-splicing or skipping of introns during translation; and the importance of an rules governing macromolecular assemblies during morphogenesis and DNA metabolism.

The essential role of homologous recombination in initiating DNA replication was first demonstrated in T4. Recombination-dependent DNA replication has also been observed in other phages, bacteria, chloroplasts, mitochondria, and yeast. It is a most important component of DNA repair (53, 94-99). It is important for survival of the T-even phages under all conditions because origin-dependent DNA replication is inactivated as a consequence of the developmental program. RNA-polymerase-generated transcripts are also used in other systems to initiate leading-strand DNA synthesis at origins (99, 100). A bypass of primase deficiencies by recombination enzymes has so far only been demonstrated in T4, but is also likely to operate in ColE1-like plasmids.

Other phages in which recombination deficiencies appear to result in premature arrest of DNA replication, include T1 (101, 102), in which a similar reason as in T4 is suspected, P22 (103), and a family of *Bacillus subtilis* phages (104, 105).

One of the more general aspects of T4 biology is the remarkable redundancy of pathways and proteins to initiate fundamental processes, such as DNA replication and recombination. The importance of redundant pathways for global evolution of novel development circuits has been discussed (106). To maintain redundant pathways during evolution, it is important that they are based on different principles and subject to different pressures. The persistence of mechanically different redundant pathways in T4 can be readily rationalized because they are best adapted to different stages of transcription, gene expression in general, and DNA packaging during phage development (39, 52, 55) (Fig. 3). Therefore, these redundancies are ready to bypass certain lesions in other important genes, (eg, the T4 primase gene), they facilitate cross-talk, coordination and optimization of transcription, translation, DNA replication, recombination, repair, and packaging, and they confer tolerance to changes that facilitate evolution.

## 11. Acknowledgment

Supported by National Science Foundation grant MCB 9983568.

## Bibliography

1. J. D. Karam et al. (eds.) (1994) *Molecular Biology of Bacteriophage T4*, American Society for Microbiology, Washington, DC.



2. E. Kutter et al (1996) *Virus Genes* **11**, 213–225.
3. C. Monod, F. Repoila, M. Kutateladze, and H. M. Krisch (1997) *J. Mol. Biol.* **267**, 237–249.
4. F. Tetart et al (2001) *J. Bacteriol.* **183**, 358–366.
5. E. Kutter et al (2001) manuscript in preparation; GenBank Accession no. AF158101.
6. U. Henning and S. Hashemol-Hosseini (1994) In Ref. 1, pp. 291–298.
7. F. Tetart, C. Desplats, and H. M. Krisch (1998) *J Mol Biol.* **282**, 543–556.
8. K. Carlson, E. A. Raleigh, and S. Hattman (1994) In Ref. 1, pp. 369–381.
9. K. Carlson and L. D. Kosturko (1998) *Mol. Microbiol.* **27**, 671–676.
10. E. Kutter et al (1994) In Ref. 1, pp. 491–519.
11. B. Stitt, and D. Hinton (1994) In Ref. 1, pp. 142–160.
12. J. K. Yee, and R. C. Marsh (1985) *J. Virol.* **54**, 271–277.
13. K. N. Kreuzer, and S. W. Morrical (1994) in Ref. 1, pp. 28–42.
14. H. Lin, and L. W. Black (1998) *Virology* **242**, 118–127.
15. G. Mosig, and D. H. Hall (1994) In Ref. 1, pp. 127–131.
16. A. D. Hershey, and M. Chase (1952) *J. Gen. Physiol.* **36**, 39–56.
17. S. T. Abedon (1994) in Ref. 1, pp. 397–405.
18. E. Goldberg, L. Grinius, and L. Letellier (1994) In Ref. 1, pp. 347–356.
19. W. Ross et al (1993) *Science* **262**, 1407–1413.
20. G. Mosig, G. W. Lin, J. Franklin, and W. H. Fan (1989) *New Biol.* **1**, 171–179.
21. P. Paddison et al. (1998) *Genetics* **148**, 1539–1550.
22. H. K. Dressman, and J. W. Drake (1999) *J. Bacteriol.* **181**, 4391–4396.
23. K. Wilkens, and W. Ruger (1994) In Ref. 1, pp. 132–141.
24. E. Kutter et al (1994) In Ref. 1, pp. 357–368.
25. K. Wilkens, and W. Ruger (1996) *Plasmid* **35**, 108–120.
26. K. Wilkens, B. Tiemann, F. Bazan, and W. Ruger (1997) In *ADP-Ribosylation in Animal Tissue* Haag, and Koch-Nolte, eds.), Plenum Press, New York, pp. 71–82.
27. J.-Y. Bouet, H. M. Krisch, and J.-M. Louarn (1998) *J. Bacteriol.* **180**, 5227–5230.
28. J. Y. Bouet, N. J. Campo, H. M. Krisch, and J. M. Louarn (1996) *Mol. Microbiol.* **20**, 519–528.
29. D. M. Hinton, R. March-Amegadzie, J. S. Gerber, and M. Sharma (1996) *J. Mol. Biol.* **256**, 235–248.
30. E. N. Brody et al (1995) *FEMS Microbiol. Lett.* **128**, 1–8.
31. F. Colland, G. Orsini, E. N. Brody, H. Buc, and A. Kolb (1998) *Mol. Microbiol.* **27**, 819–829.
32. R. March-Amegadzie, and D. M. Hinton (1995) *Mol. Microbiol.* **15**, 649–660.
33. J. C. Hinton (1997) *Mol. Microbiol.* **26**, 417–422.
34. E. P. Geiduschek (1995) *Semin. Virol.* **6**, 25–33.
35. R. S. Washburn and B. L. Stitt (1996) *J. Mol. Biol.* **260**, 332–346.
36. S. Sozhamannan and B. L. Stitt (1997) *J. Mol. Biol.* **268**, 689–703.
37. M. W. Frazier and G. Mosig (1990) *Gene* **88**, 7–14.
38. G. Mosig, N. E. Colowick and B. C. Pietz (1998) *Gene* **223**, 143–155.
39. G. Mosig, J. Gewin, A. Luder, N. Colowick and D. Vo (2001) *Proc. Natl. Acad. Sci. U.S.A.* **98**, 8306–8311.
40. E. S. Miller, J. D. Karam and E. Spicer (1994) In Ref. 1, pp. 193–205.
41. N. G. Nossal (1994) In Ref. 1, pp. 43–53.
42. S. J. Benkovic, A. M. Valentine and F. Salinas (2001) *Annu. Rev. Biochem.* **70**, 181–208.

43. F. A. Kadyrov and J. W. Drake (2001) *J. Biol. Chem.* **4**, 4.
44. L. J. Reha-Krantz (1998) *Genetics* **148**, 1551–1557.
45. L. J. Reha-Krantz et al (1998) *J. Biol. Chem.* **273**, 22969–22976.
46. N. G. Nossal (1998) *Genetics* **148**, 1535–1538.
47. J. Wang et al (1997) *Cell* **89**, 1087–1099.
48. Y. Shamoo, A. M. Friedman, M. R. Parsons, W. H. Konigsberg and T. A. Steitz (1995) *Nature* **376**, 362–366.
49. Y. Shamoo, A. M. Friedman, M. R. Parsons, W. H. Konigsberg and T. A. Steitz (1995) *Nature* **376**, 616.
50. Y. Shamoo and T. A. Steitz (1999) *Cell* **99**, 155–166.
51. M. A. Trakselis, M. U. Mayer, F. T. Ishmael, R. M. Roccasecca and S. J. Benkovic (2001) *Trends Biochem. Sci.* **26**, 566–572.
52. G. Mosig (1998) *Annu. Rev. Genet.* **32**, 379–413.
53. K. N. Kreuzer (2000) *Trends Biochem. Sci.* **25**, 165–173.
54. J. W. George, B. A. Stohr, D. J. Tomso and K. N. Kreuzer (2001) *Proc. Natl. Acad. Sci. U.S.A.* **98**, 8290–8297.
55. G. Mosig, N. Colowick, M. E. Gruidl A. Chang and A. J. Harvey (1995) *FEMS Microbiol. Rev.* **17**, 83–98.
56. R. Vaiskunaite, A. Miller, L. Davenport and G. Mosig (1999). *J. Bacteriol.* **181**, 7115–7125.
57. K. G. Belanger and K. N. Kreuzer (1998) *Mol. Cell* **2**, 693–701.
58. N. G. Nossal K. C. Dudas and K. N. Kreuzer (2001) *Mol. Cell* **7**, 31–41.
59. A. Luder and G. Mosig (1982) *Proc. Natl. Acad. Sci. U.S.A.* **79**, 1101–1105.
60. K. Carles-Kinch, J. W. George and K. N. Kreuzer (1997) *EMBO J.* **16**, 4142–4151.
61. T. R. Broker (1973) *J. Mol. Biol.* **81**, 1–16.
62. T. Formosa and B. M. Alberts (1986) *Cell* **47**, 793–806.
63. M. Belfort and P. S. Perlman (1995) *J. Biol. Chem.* **270**, 30237–30240.
64. J. E. Mueller, J. Clyman, Y.-J. Huang, M. M. Parker and M. Belfort (1996) *Genes Dev.* **10**, 351–364.
65. L. W. Black, M. K. Showe and A. C. Steven (1994) In Ref. 1, pp. 218–258.
66. L. W. Black (1995) *Bio essays* **17**, 1025–1030.
67. J. G. Franklin and G. Mosig (1996) *Gene* **177**, 179–189.
68. J. L. Franklin, D. Haseltine, L. Davenport and G. Mosig (1998) *J. Mol. Biol.* **277**, 541–557.
69. S. P. Bhattacharyya and V. B. Rao (1994) *Gene* **146**, 67–72.
70. B. Kemper (1998) In *DNA Damage and Repair: DNA Repair in Prokaryotes and Lower Eukaryotes* (J. A. Nickoloff and M. Hoekstra, eds.), Humana Press, Totowa, NJ, Vol. **1**, pp. 179–204.
71. S. Golz and B. Kemper (1999) *J. Mol. Biol.* **285**, 1131–1144.
72. D. Ang, F. Keppel, G. Klein, A. Richardson and C. Georgopoulos (2000) *Annu. Rev. Genet.* **34**, 439–456.
73. Y. R. Hong and L. W. Black (1993) *Gene* **136**, 193–198.
74. P. J. Jardine and D. H. Coombs (1998) *J. Mol. Biol.* **284**, 661–672.
75. G. Mosig (1963) *Cold Spring Harbor Symp. Quant. Biol.* **28**, 35–42.
76. G. Mosig (1968) *Genetics* **59**, 137–151.
77. G. Mosig et al (1972) *J. Virol.* **9**, 857–871.
78. J. A. Haynes, F. A. Eiserling (1996) *Virology* **221**, 67–77.
79. A. W. Kozinski and L. D. Kosturko (1976) *J. Virol.* **17**, 801–804.

80. F. A. Eiserling and L. W. Black (1994) In Ref. 1, pp. 209–217.
81. K. Carlson, E. Raleigh, and A. S. Hattman (1994) In Ref. 1, pp. 369–381.
82. F. J. Stewart, D. Panne, T. A. Bickle, and E. A. Raleigh (2000) *J. Mol. Biol.* **298**, 611–622.
83. D. Panne, S. A. Müller, S. Wirtz, A. Engel, and T. A. Bickle (2001) *EMBO J.* **20**, 3210–3217.
84. V. A. Barcus, A. J. B. Titheradge, and N. E. Murray (1995) *Genetics* **140**, 1187–1197.
85. T. Georgiou et al (1998) *Proc. Natl. Acad. Sci. U.S.A.* **95**, 2891–2895.
86. L. Snyder, and G. Kaufmann (1994) In Ref. 1, pp. 391–396.
87. D. H. Parma et al (1992) *Genes Dev.* **6**, 497–510.
88. G. Mosig et al (1997) *Virology* **230**, 72–81.
89. R. Calendar et al. (1998) in *Horizontal Gene Transfer* (M. Syvanen, ed.), Chapman and Hall, pp. 241–252.
90. J.-S. Kim, and N. Davidson (1974) *Virology* **57**, 93–111.
91. T. P. Gary, N. E. Colowick, and G. Mosig (1998) *Genetics* **148**, 1461–1473.
92. F. Tetart, C. Monod, and H. M. Krisch (1996) *J. Mol. Biol.* **258**, 726–731.
93. G. Mosig, and R. Calendar (2001) In *Horizontal Gene Transfer II* (M. Syvanen, ed., Academic Press, London, Chap. "13", in press.
94. M. M. Cox et al. (2000) *Nature*, 37–41.
95. M. M. Cox (2001) *Proc. Natl. Acad. Sci. USA* **98**, 8173–8180.
96. S. C. Kowalczykowski (2000) *Trends Biochem. Sci.* **25**, 156–165.
97. B. Michel et al (2001) *Proc. Natl. Acad. Sci. USA* **98**, 8181–8188.
98. E. Kraus, W. Y. Leung, and J. E. Haber (2001) *Proc. Natl. Acad. Sci. USA* **98**, 8255–8262.
99. T. Kogoma (1997) *Microbiol. Mol. Biol. Rev.* **61**, 212–248.
100. A. Kornberg, T. A. Baker (1992) *DNA Replication*. W. H. Freeman, New York, 2 ed.
101. J. Liebeschuetz, R. D. Harris, and D. A. Ritchie (1987) *J. Gen. Virol.* **68**, 2049–2062.
102. J. R. Christensen (1994) in *Encyclopedia of Virology* (R. G. Webster, and A. Granoff, eds.), Academic Press, San Diego, CA Vol. **3**, pp. 1371–1376.
103. D. Botstein, and M. J. Mat (1970) *J. Mol. Biol.* **96**, 87.
104. A. Bravo, and J. C. Alonso (1990) *Nucl. Acids Res.* **18**, 4651–4657.
105. C. R. Steward (1994) In *Encyclopedia of Virology* (R. G. Webster, and A. Granoff, eds.), Academic Press, San Diego, Vol. **3**, pp. 1352–1356.
106. M. Kirschner, and J. Gerhart (1998) *Proc. Natl. Acad. Sci. U.S.A.* **95**, 8420–8427.
107. R. S. Edgar, and W. B. Wood (1966) *Proc. Natl. Acad. Sci. U.S.A.* **55**, 498–505.

### **Suggestions for Further Reading**

108. J. D. Karam et al (eds.) (1994) *Molecular Biology of Bacteriophage T4*. American Society for Microbiology, Washington, DC.
109. E. Kutter et al (1996) *Virus Genes* **11**, 213–225.
110. E. P. Geiduschek (1995) *Semin. Virol.* **6**, 25–33.
111. G. Mosig et al (1995) *FEMS Microbiol. Rev.* **17**, 83–98.
112. G. Mosig (1998) *Annu. Rev. Genet.* **32**, 379–413.
113. E. Kutter et al (2001) manuscript in preparation; GenBank Accession no. AF158101

## 1. General Properties

One of the original Type phages isolated by Demerec and Fano, T7 is now considered the prototype of a group of virulent phages that have an icosahedral head and a short, stubby tail, to which are attached six tail fibers (1). A major feature distinguishing it from other members of the *Podoviridae* is the synthesis, early after infection, of a new **RNA polymerase** that is resistant to the antibiotic [rifampicin](#) and highly specific for **promoters** on the phage [genome](#). T7 was isolated as a phage that grows on *Escherichia coli* B, but T7 and its close relatives grow equally well on *E. coli* K-12 and C strains and on some of the Shigellae. It does not form plaques on fresh natural isolates of *E. coli* because it does not adsorb well (if at all) to smooth or capsule-containing bacteria, and it can be propagated on *Salmonella typhimurium* strains only if the latter are rough (and lack the *hsdLT* restriction system).

The T7 group of phages has been subdivided into three on the basis of the promoter specificity of the phage RNA polymerase: (i) Of the more than 60 coliphage isolates, T3 appears unique and is more closely related to the *Yersinia enterocolitica* phage YeO3-12 than other coliphages. (ii) There are three members of a second group (BA14, BA127, and BA156). (iii) The remainder are like T7 (2). Members of a given subdivision undergo [recombination](#) efficiently with each other, but recombinants between phages in different subdivisions are rare. In part this is due to divergence of the [terminal repeat](#) sequences that are essential for growth, and in part to the specificity of the phage RNA polymerase for its cognate promoters. In addition, T3 and T7 (the BA phages have not been tested) also exhibit mutual exclusion in coinfections, with most such coinfecting cells producing either T3 only or T7 only. Nevertheless, the ancestors of T7 and T3 did recombine, because the gene *I7* promoter in T3 is actually of T7 specificity and is not used during a normal T3 infection.

[Heteroduplex](#) analyses of a few T7-like genomes have shown that they exhibit varying degrees of [homology](#), containing regions of >90% sequence identity and others with no apparent similarity, this has been confirmed by DNA sequence data for T7 and T3. T7-like phages that infect enteric bacteria other than *E. coli* are significantly less closely related than the coliphages; for example, the *S. typhimurium* phage SP6 has not even maintained the same gene order as T7. As may then be expected, phages growing in nonenteric bacteria have diverged even further.

## 2. Genetic Structure

The genetic map of T7 is based on the DNA sequence of 39,937 bp (3; Genbank Accession V01146) and contains 56 known or potential genes (Table 1). The genes are all **transcribed** from the same strand of DNA and are ordered by number sequentially from the genetic left end, which is the first end of the genome to enter the infected cell. Thus gene *0.3* is the first gene to be expressed, and *19.5* is the last. Integral numbered genes are unconditionally essential (except that gene 2 [amber mutants](#) grow on *E. coli* C), whereas most noninteger genes are nonessential or conditionally essential, exceptions being 2.5, 6.7, and 7.3. Coding sequences occupy almost 92% of the genome; most of the remainder contains the terminal repeats, **origins of replication**, **promoters**, and ribonuclease III processing sites, or other recognizable genetic signals. The T7 genes are described as close-packed, but there are five instances of potential overlapping genes that are read in a different [reading frame](#): genes *4.1* and *4.2* are almost entirely within gene 4 coding sequences, gene *18.7* lies within gene *18.5*, and both genes *19.2* and *19.3* overlap the gene *19* sequence. In addition, gene 4 codes for two proteins, gp4A and gp4B, due to the use of an in-frame internal [initiation codon](#). Furthermore, programmed [frameshifting](#) of the [ribosome](#) to the +1 frame during [translation](#) of gene *0.6* and to the -1 frame during translation of genes *5.5* and *10* affords, respectively, gp0.6B, a 168-residue *5.5-5.7* fusion protein, and gp10B (3, 4). No biological function for these frame-shifted products is known, but T3 contains a different, yet comparable, shift sequence at the same relative position in gene *10* that also leads to a gp10B product (5). The gene *10* homologues of other T7-like phages are also thought to make two products; gp10B may therefore be important under some conditions of

infection.

**Table 1. Genetic Map of T7**

| T7 Gene  | $M_r$ product       | Function or Comment  |
|----------|---------------------|--|
| Class-I  |                     |  |
| 0.3      | 13,678              | Nonessential; inactivates Type I restriction enzymes                                 |
| 0.4      | 5,621               | Nonessential   |
| 0.5      | 4,744               | Nonessential   |
| 0.6A     | 6,201               | Nonessential   |
| 0.6B     | 13,250 <sup>a</sup> | Nonessential   |
| 0.7      | 41,124              | Nonessential; protein kinase, inactivates host transcription. Col Ib exclusion       |
| 1        | 98,092              | RNA polymerase   |
| 1.1      | 5,180               | Nonessential   |
| 1.2      | 10,059              | Nonessential; inhibits <i>E.coli</i> deoxyguanosine triphosphohydrolase, F exclusion |
| 1.3      | 41,133              | DNA ligase; exclusion by <i>Shigella sonnei</i> D <sub>2</sub> 371-48                |
| Class II |                     |  |
| 1.4      | 5,446               | Nonessential   |
| 1.5      | 3,174               | Nonessential   |
| 1.6      | 9,946               | Nonessential   |
| 1.7      | 22,053              | Nonessential   |
| 1.8      | 5,781               | Nonessential   |
| 2        | 7,043               | Inactivates <i>E. coli</i> RNA polymerase  |
| 2.5      | 25,562              | Single-stranded DNA-binding protein  |
| 2.8      | 15,617              | Nonessential; derived from group I intron? Absent in T3 genome                       |
| 3        | 17,040              | Endonuclease, Holliday junction resolvase  |
| 3.5      | 16,806              | Amidase (lysozyme), regulates T7 RNA polymerase activity                             |
| 3.8      | 14,329              | Nonessential; derived from group I intron? Absent in T3 genome                       |
| 4A       | 62,656              | Primase/helicase   |
| 4B       | 55,743              | Helicase   |
| 4.1      | 4265                | Overlapping gene, expression not determined  |
| 4.2      | 12,653              | Overlapping gene, expression not determined  |
| 4.3      | 7,927               | Nonessential   |
| 4.5      | 9,960               | Nonessential   |
| 4.7      | 15,208              | Nonessential. Absent in T3 genome  |
| 5        | 79,692              | DNA polymerase   |
| 5.3      | 13,067              | Nonessential; derived from group I intron?   |

|           |         |   |
|-----------|---------|---|
| 5.5       | 7,280   | Nonessential; binds to H-NS, <i>l rex</i> exclusion. 5.5–5.7 fusion protein also made |
| 5.7       | 7,280   | Nonessential  |
| 5.9       | 5,913   | Nonessential, inactivates RecBCD  |
| 6         | 34,371  | 5' → 3' double-stranded exonuclease   |
| 6.3       | 4,088   | Nonessential  |
| Class III |         |   |
| 6.5       | 9,474   | Nonessential  |
| 6.7       | 9,207   | Adsorption  |
| 7         | 15,303  | Nonessential; host range. Absent in T3 genome   |
| 7.3       | 9,937   | Initiation of infection   |
| 7.7       | 14,737  | Nonessential; derived from group I intron? Absent in T3 genome                        |
| 8         | 58,989  | Head–tail connector   |
| 9         | 33,766  | Scaffolding protein   |
| 10A       | 36,414  | Major capsid protein, <i>Shigella sonnei</i> D <sub>2</sub> 371-48 and F exclusion    |
| 10B       | 41,700  | Minor capsid protein, <i>Shigella sonnei</i> D <sub>2</sub> 371-48 and F exclusion    |
| 11        | 22,289  | Tail protein  |
| 12        | 89,265  | Tail protein  |
| 13        | 15,852  | Internal head protein; initiation of infection  |
| 14        | 20,836  | Internal core protein ejected from particle at initiation of infection                |
| 15        | 84,210  | Internal core protein ejected from particle at initiation of infection                |
| 16        | 143,840 | Internal core protein ejected from particle at initiation of infection                |
| 17        | 61,411  | Adsorption, tail fiber protein  |
| 17.5      | 7,391   | Holin for cell lysis  |
| 18        | 10,145  | DNA packaging, small subunit  |
| 18.5      | 16,243  | Cell lysis; <i>l Rz</i> homologue   |
| 18.7      | 9,195   | Overlapping gene; cell lysis; <i>l RzI</i> homologue                                  |
| 19        | 66,130  | DNA packaging, large subunit  |
| 19.2      | 9,264   | Overlapping gene, expression not determined   |
| 19.3      | 6,429   | Overlapping gene, expression not determined   |
| 19.5      | 5,434   | Nonessential  |

---

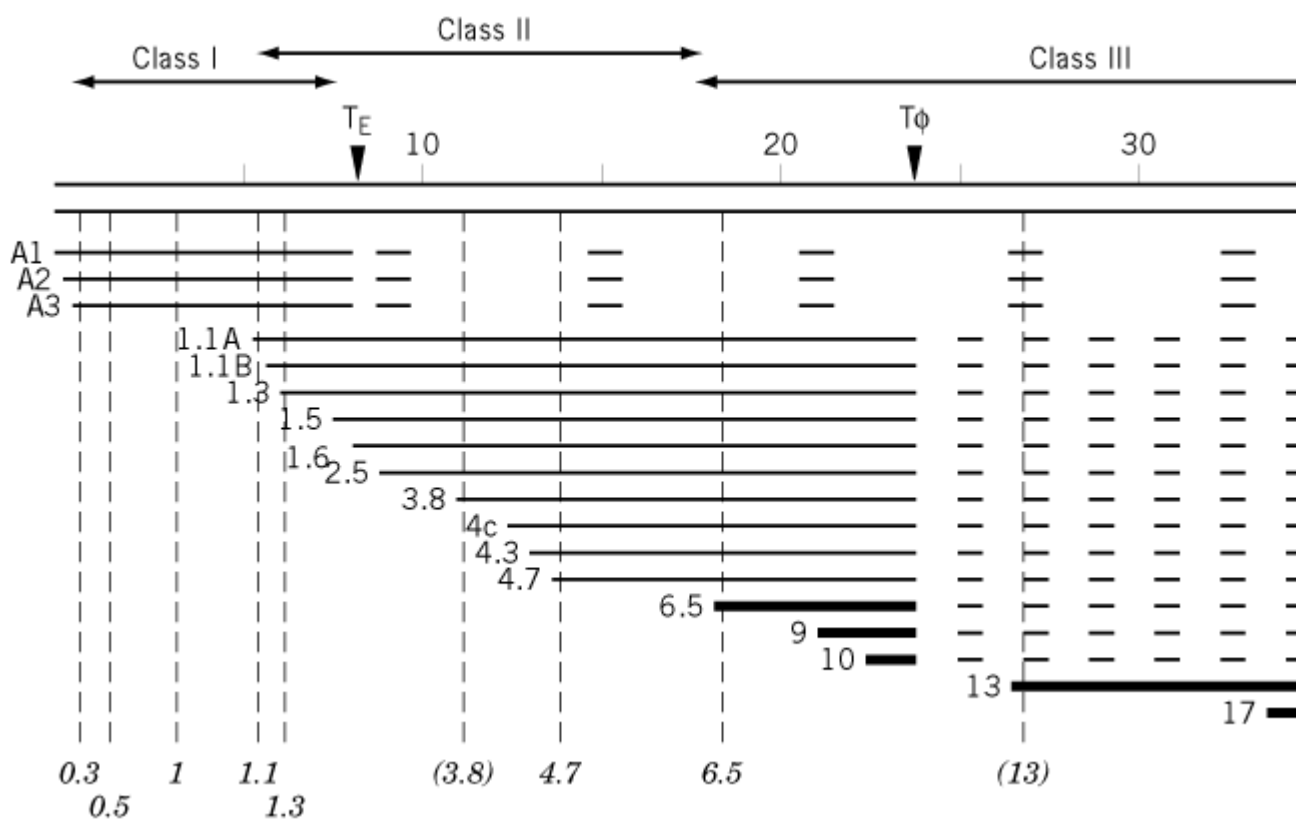
<sup>a</sup> Predicted only, site of frameshift not determined experimentally.

T7 DNA contains only the four normal bases; a limited amount of adenine **methylation** occurs from the Type I **restriction–modification** enzyme and from Dam methylase, but not all potential sites are fully methylated, because the rate of T7 DNA replication and packaging is so high. Some of the

other T7-like phage DNAs, including that of T3, are completely unmethylated, because they code for an [S-adenosylmethionine](#) hydrolase, which removes the donor of the methyl groups. The sequence GATC, as well as sequences corresponding to several common [restriction enzymes](#), is grossly underrepresented in T7 DNA. For example, although expected 156 times on statistical grounds, the sequence GATC is present only six times in T7 DNA, ten times in T3 DNA, and not once in BA14 DNA.

Three classes of T7 genes have been recognized (Fig. 1): class I genes are those expressed from about 2 min to 8 min after infection at 30°C; class II genes from about 6 min to 15 min; and class III genes from 8 min until lysis, which occurs after ~25min (6). The ten class I genes are transcribed by *E. coli* RNA polymerase; their functions are mainly to subvert the bacterium into a phage-producing factory. The only essential class I gene is gene *1*, coding for the T7 RNA polymerase. The latter enzyme transcribes the class II genes, which are involved mainly in phage DNA metabolism, and the class III genes, whose functions are predominantly morphogenetic.

**Figure 1.** Schematic depicting the RNAs of phage T7. The three classes of genes are shown above the T7 genome, and terminators are indicated.  $T_E$  only affects *E. coli* RNA polymerase.  $T_\phi$  is the major T7 RNA polymerase terminator; on r dependent terminator, designated  $T_{pac}$ , is important for DNA packaging. The promoter for each transcript is indicated at first gene transcribed; dashed lines indicate that readthrough of the terminators occurs. Vertical dashed lines indicate site by the first gene following the site of cleavage; sites where cleavage is inefficient are in parentheses.



All but five T7 proteins are initiated with an AUG codon; the remaining five use GUG. All coding sequences are preceded by a recognizable [Shine-Dalgarno sequence](#) that can potentially form at least four base pairs with the T7 mRNA. Some of the [messenger RNAs](#) are more efficiently translated than others, suggestive of some level of translational control. The Class III gene *10* mRNA is one of the most actively translated RNAs known. More than  $10^5$  copies of gp10 are made in about 20 min at 30°C. In addition to the abundance of the gene *10* mRNA, its leader sequence and the initial codons

appear designed to maximize ribosome binding and translation initiation.

### 3. The T7 Particle

The phage capsid (**icosahedral symmetry**,  $T = 7$ ) contains 415 molecules of either gp10A or gp10B (7). About 10% of the capsid is normally gp10B, but phages making only gp10A or gp10B are viable; gp10A-only phages grow like the wild type, but gp10B-only phages produce distinctly smaller plaques. At one vertex of the icosahedron, 12 molecules of gp8 act as the head–tail connector, whose structure is known from [cryoelectron microscopy](#) (8, 9). A central hole extends the length of the connector and is collinear with the internal channel of the ~20–nm tail; the latter consists of 6 molecules of gp12 and 12 copies of gp11. Six tail fibers, each containing trimers of gp17, are attached to the tail just below its junction with the connector protein. The 87-residue gp6.7 should also be part of the tail or tail fiber, because 6.7 null mutants fail to adsorb to cells. In addition to DNA, the head also contains a number of internal proteins, most prominent of which are gp14 (18 copies), gp15 (12 copies), and gp16 (3 copies), which form a hollow cylindrical structure, coaxial with the tail; this internal core is attached to the inner surface of the capsid and the connector (7, 10). The core is essential both for DNA ejection and for morphogenesis of new phage particles. The double-stranded DNA is spooled in about six layers around the axis of the core–connector–tail structure (11).

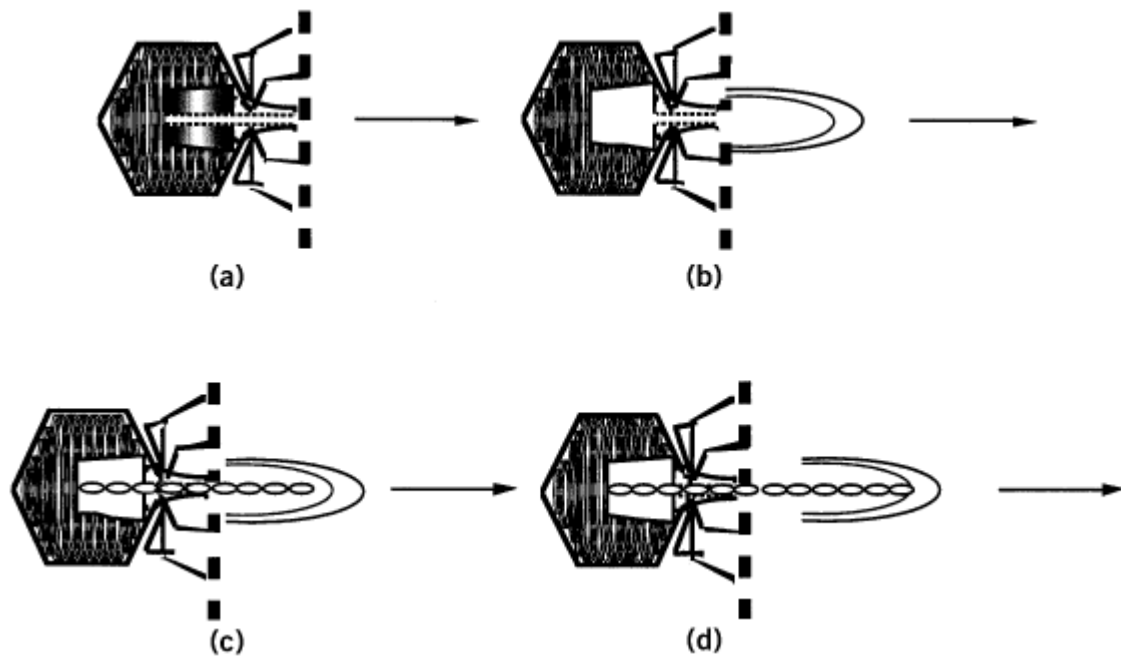
### 4. Infection Cycle

#### 4.1. Adsorption, DNA Penetration, and Transcription

The tail fibers attach to the lipopolysaccharide of the *E. coli* outer membrane to initiate adsorption. A second receptor probably exists, and probably interacts with a tail component, but it has not yet been identified. Following adsorption, the infection is made irreversible by the degradation of the internal head proteins gp7.3 and gp13; the core then disaggregates, and its components (gp14, gp15, and gp16) pass through the connector and tail and enter the cell (Fig. 2). Both gp15 and gp16 are large proteins (84.2 and 143.8 kDa, respectively), and they may have to unfold to pass through the connector–tail channel, whose diameter is only ~3.7nm.

**Figure 2.** Conceptual model for entry of the T7 genome into the cell. **(a)** Proteolysis allows disaggregation and ejection of the internal core structure into the cell. Gp14 localizes to the outer membrane, possibly forming a channel. Gp16 inserts into the inner membrane; gp15 is also in the cell. **(b)** Probably with the help of gp15, gp16 establishes a DNA translocase. **(c)** About 850 bp of the genome is ratcheted into the cell, causing separation of the translocase from its energy source; DNA translocation stops. **(d)** Transcription, initially by *E. coli* RNA polymerase and then by T7 RNA polymerase, brings the remainder of the 40-kbp genome into the cell.





Gp14 localizes to the outer membrane, presumably forming a channel. The 143-kDa gp16 can then enter the **periplasm**, from where it inserts into the cytoplasmic membrane, thereby completing a protein channel across the cell envelope for DNA translocation. Gp15 is thought to pass through this channel into the cell cytoplasm, although it cannot yet be excluded that the protein has localized to an osmotic shock-resistant part of the periplasm. The natural cross-linking of peptidoglycan in the outer membrane provides a barrier to molecules greater than  $\sim 50\text{kDa}$  (12), but the *N*-terminal portion of gp16 is homologous to the active site region of *E. coli* *sltY*, the soluble lytic transglycosylase Slt70, and this part of the protein may enlarge a hole through the peptidoglycan layer. This presumed activity of gp16 is important only when the conditions of infection favor increased peptidoglycan cross-linking (low temperature, high cell density).

Only when the transmembrane channel has been formed does DNA leave the phage head, but even then only  $\sim 850\text{bp}$  of the 40-kb genome enter the cell immediately (13). Thus the often espoused idea that phages are like syringes with their contents under pressure does not apply to T7 (and probably not to any phage). The remainder of the genome is normally translocated into the cell as a result of [transcription](#) (see text below). The [proton motive force](#) is necessary for the leading part of the genome to enter the cell, energy seemingly being required both for the insertion of gp16 into the cytoplasmic membrane and to power the motor that ratchets T7 DNA into the cell. The motor is currently thought to consist of gp16 and gp15: certain gene 16 missense mutants uncouple genome entry from transcription and cause the entire 40-kb genome to enter the cell at a constant  $\sim 70\text{bp/s}$  at  $30^\circ\text{C}$  (14); gene 15 mutants have also been shown to affect the kinetics of genome entry.

The leading 850 bp of T7 DNA that enter the cell contain three strong *E. coli* promoters and a T7 promoter. Transcription from any of the *E. coli* promoters not only makes RNAs for class I gene expression, but it also pulls the first 19% of the genome (to the Rho-independent terminator TE) into the cell. The rate of genome entry at this stage is  $\sim 40\text{bp/s}$  at  $30^\circ\text{C}$ , the same rate that *E. coli* RNA polymerase transcribes DNA *in vitro* under optimal conditions (13). Again via transcription, T7 RNA polymerase pulls the remainder of the phage genome into the cell at  $\sim 250\text{bp/s}$ , simultaneously making mRNAs for class II and class III gene expression. About 10 min, one-third of the latent period, is occupied with the complete internalization of the T7 genome. The slow entry may play a role in control of gene expression, but the most compelling explanation for the coupling of genome entry to transcription is that the first gene transcribed is *o.3*, whose product binds and inactivates Type I restriction enzymes before most of the genome enters the cell. T7 normally is fully resistant

to Type I restriction, but its genome is rapidly degraded if introduced into the cell from a [lambda phage](#) particle (13).

T7 has no defense, other than lacking the pertinent sequences, against Type II restriction enzymes, although some members of the T7 family do grow on hosts that contain Type II enzymes, even though their genomes do contain cognate recognition sites. T7 is also sensitive to the **P1 phage**-encoded Type III enzyme, but not to that of *EcoP15*, whose methylase recognition sequence is 5'-CAGCAG/5'-CTGCTG. The *EcoP15* restriction enzyme requires two copies of this sequence in opposite orientations, but the 36 copies of this hexamer are all oriented the same way in the T7 genome (15). T3 is restricted by *EcoP15*, because its genome contains the recognition sequence in both orientations.

*In vitro*, the linear T7 genome is susceptible to degradation by the **RecBCD** nuclease. Gp5.9 protein binds stoichiometrically to RecBCD holoenzyme, but not to its isolated subunits, and inhibits all its **nuclease** and [ATPase](#) activities. Gp5.9 is a class II protein and is not found in the phage particle, so how the entering T7 genome is protected from RecBCD at the beginning of the infection is unknown. Furthermore, gene 5.9 appears to be nonessential for growth in RecBCD-containing cells (16); either the ends of T7 DNA are protected from the RecBCD, or a second but unknown T7 gene may have an activity comparable to that of gp5.9.

T7 primary transcripts are processed by RNase III at ten sites, although the enzyme is not essential for phage development (3). Except for the gene 1 RNA, the processed RNAs have multiple **cistrons**, but there is no evidence for polarity effects due to amber mutations. In the late region of the genome, this can be easily explained by the fast rate of transcription by T7 RNA polymerase. T7 RNAs are stable throughout the infection, possibly due to the formation of base-paired structures at the 3' ends of the processed primary transcripts. The 0.3 mRNA, which is used *in vivo* for only about 8 min of infection, is said to be competent for *in vitro* translation experiments even when isolated late after infection. How the translational shut-off of class I and class II proteins is achieved is unknown.

In addition to gp0.3 and gp1, the phage RNA polymerase, other early gene products include gp0.7, a **serine/threonine kinase** that **phosphorylates** several host proteins involved in translation and RNA metabolism, including RNase III (17). Distinct from its kinase activity, gp0.7 also causes the abrupt shut-off of host-catalyzed transcription (18). Despite its multiple functions, mutants lacking 0.7 grow better than wild-type in rich medium in most laboratory strains. However, 0.7 is necessary for phage growth under nutrient-poor conditions. Gp1.2 inhibits *E. coli* dGTPase (*dgt*), and gp1.3 is a [DNA Ligase](#); these proteins are essential only if the host overexpresses *dgt* or is ligase-defective, respectively.

The [X-ray crystallography](#) structure of the 883-residue T7 RNA polymerase has been determined at 3.3 Å resolution (19). About 45% of the molecule shows extensive structural homology to the **Klenow fragment** of *E. coli* **DNA polymerase I** and resembles a cupped right hand, with finger, palm, and thumb domains that form a cleft. Amino acid residues in the cleft can be **cross-linked** to the promoter. The thumb subdomain wraps around the template strand and probably confers processive synthesis of RNA. A loop of residues 743 to 773 project out from the fingers domain of the cleft and contain promoter-specificity determinants.

T7 promoters consist of a highly conserved 23-bp sequence that runs from nucleotides -17 to +6 relative to the transcription start site. Two functional domains have been recognized, an initiation domain that spans positions -5 to +6, and a binding domain between -17 and -6 (20). The Asn748Asp mutant form of the T7 enzyme recognizes T3 (and not T7) promoters; conversely, changing positions -10 and -11 of a T7 promoter to those found in T3 promoters abolishes T7 promoter activity, but allows full recognition by T3 polymerase (21, 22). Similarly, positions -8 and -9 are the primary determinants of SP6 versus T7 promoters (23). A second specificity determinant that has been recognized is residue 758 of the T7 enzyme, which interacts with the promoter at

position -8 (24).

There are 17 T7 promoters in the phage genome, ten upstream of class II genes, five upstream of class III genes, and two,  $\phi OL$  and  $\phi OR$ , near the genome ends.  $\phi OL$  is used only weakly, while  $\phi OR$  is a strong, consensus promoter that can provide an mRNA only for gene 19.5; it has been proposed that  $\phi OL$  and  $\phi OR$  are replication promoters, but neither one is essential for viability. There is a major terminator  $T\phi$  for T7 RNA polymerase just downstream of gene 10. All transcripts from class II promoters and those from the class III promoters  $\phi 6.5$ ,  $\phi 9$ , and  $\phi 10$  terminate at  $T\phi$  with about 85% efficiency (25). This results in a nested set of mRNAs, all of which contain gene 10 near their 3' end, thereby providing the necessary high dosage of the major capsid protein mRNA. Readthrough of the terminator  $T\phi$  is programmed into the T7 developmental cycle; the approximately 15% of transcripts that fail to terminate at  $T\phi$  include RNA for the essential genes 11 and 12, which are not expressed at high levels (3).

Class III promoters all have the consensus sequence, whereas class II promoters vary at one to three positions. Class II promoters produce less mRNA, mainly because of an increased frequency of abortive initiation, which is further enhanced by formation of a complex of T7 RNA polymerase (gp1) with lysozyme (gp3.5). Gp3.5 exerts its effects on T7 RNA polymerase by preventing its isomerization from the initiation to elongation complex (26). Abortive initiation is the main reason why class II promoters appear to be weaker than class III promoters, and why class II gene expression is selectively shut off during infection. The inhibitory gp1:gp3.5 interaction constitutes a feedback loop because gene 3.5 is a class II gene that is expressed by gp1. This interaction is important, although not essential, in the T7 life cycle: Its absence reduces RNA polymerase function in both initiating DNA replication and in the maturation/packaging steps of replicated DNA.

The crystal structure of gp3.5 has been determined to 2.7 Å and that of a mutant protein lacking residues 2 to 5 to 2.2 Å (27). The N-terminal eight amino acid residues, together with a region containing residues 30 to 42, form a surface that interacts with T7 RNA polymerase. Gp3.5 is a single-domain globular protein first characterized as possessing a **lysozyme**-like activity, shown later to be a zinc-containing amidase. Zinc is necessary for catalytic activity, but not for polymerase binding. Amidase activity is important, but not essential, for cell lysis by T7; amidase mutants only delay lysis, and the major role of the amidase activity of gp3.5 is to release phage from cell debris. The lysis defect observed with gene 3.5 mutants that fail to interact with gene 1 is due to defects in [DNA replication](#) and maturation; a comparable lysis defect is also observed with mutants defective in other replication functions.

#### 4.2. DNA Replication

DNA replication *in vivo* requires T7 RNA polymerase, DNA polymerase, [primase/helicase](#), [single-stranded DNA binding protein](#), endonuclease, and exonuclease. Two other phage gene products are less directly required: The importance of the gp1:gp3.5 interaction mentioned above may be to stall RNA polymerase near the promoter(s) in an origin, causing the origin DNA to remain strand-separated and to allow assembly of the replication machinery. The gene 2 protein is also normally required in DNA replication to inhibit *E. coli* RNA polymerase, which, although transcriptionally inactive by the time DNA synthesis begins, can interfere with the late steps of replication. Phage replication is independent of all *E. coli* DNA replication genes, including topoisomerases; the latter is due to the T7 genome not being circularized. However, *E. coli* [thioredoxin](#) is absolutely required for T7 replication; the protein makes a 1:1 complex with T7 gp5 to form a highly processive DNA polymerase. Purified gp5 has extremely limited DNA polymerase and 3' → 5' exonuclease activities. The presence of thioredoxin increases the processivity of the polymerase 1000-fold (to a level equivalent to more complicated DNA polymerases that require a **sliding clamp** for processive synthesis), thereby also increasing the overall rate of synthesis (28). Stimulation by thioredoxin is independent of its **redox** capacity, although oxidized thioredoxin does not bind to gp5. Replacing either active-site [cysteine](#) residue in thioredoxin with Ser or Ala reduces its affinity for gp5 but does not reduce the maximum level of processivity increase (29). Addition of gp2.5 to T7 DNA polymerase further increases the processivity and also increases the specific rate of DNA synthesis *in*

*vitro* (30). A remarkable crystal structure of T7 DNA polymerase, complexed with a primer-template and nucleoside triphosphate in the polymerase active site, has recently been generated at 2.2 Å resolution (31). Like the RNA polymerase, T7 DNA polymerase shares extensive sequence homology with *E. coli* DNA polymerase I, and both are members of the Pol I structural family. The addition of the 176 residue thioredoxin-binding domain of gp5 to the comparable region of Pol I makes a protein that forms a 1:1 complex with thioredoxin. Interestingly, and consistent with the function of thioredoxin in T7 DNA polymerase, both Pol I and the modified enzyme synthesize DNA distributively, but addition of thioredoxin results in a large increase in processivity by the modified enzyme (32).

The primary origin of replication is an A + T-rich region about 15% from the genetic left end, upstream of and overlapping the 5' end of gene *I.1*. This origin is not essential, because deletion mutants lacking this region grow well. Several secondary origins have been identified using [electron microscopy](#) of replicating molecules and by both *in vivo* and *in vitro* plasmid replication assays (33, 34). The *oOR* secondary origin can also be deleted from a primary origin mutant; it is likely that several sequences can be used to initiate replication under normal laboratory conditions, and perhaps all potential origins can be active under different growth conditions. Replication from the primary origin is bidirectional on the linear T7 genome, although the leftward moving fork is assembled later than the rightward fork. Electron microscopy has revealed replication bubbles and Y-shaped molecules (the origin is closer to the left end), and multiple initiation events on the same molecule have been observed.

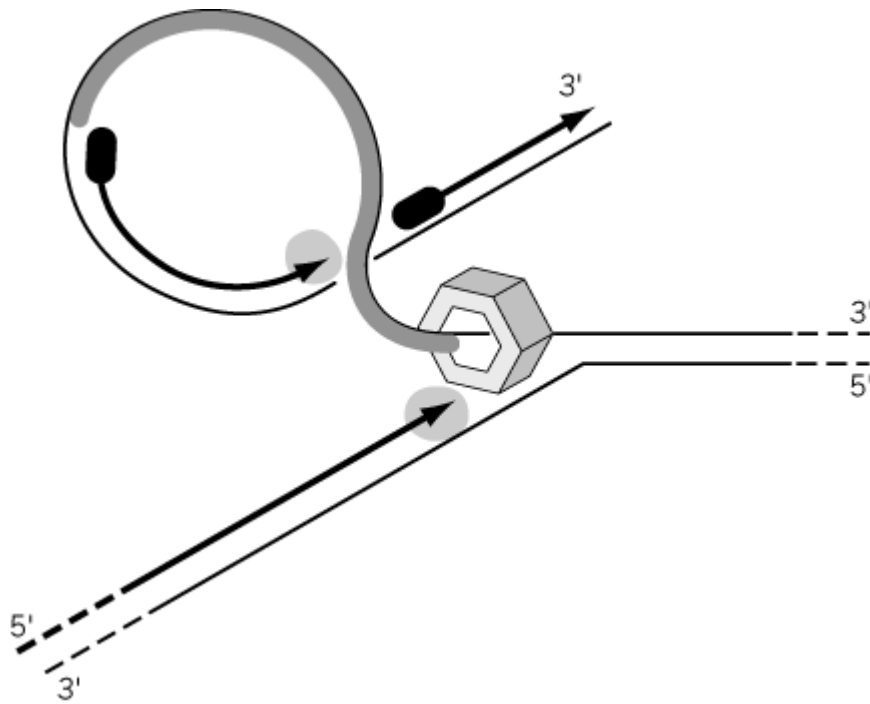
The replicated monomer genome necessarily has single-stranded 3' ends; degradation of the RNA primer and complete erosion of the terminal repeat into unique sequences on the newly synthesized strand is normally catalyzed by gp6, a 5' → 3' double-stranded exonuclease. The direct repeat allows the formation of linear concatemers; although exonuclease-treated genomes can be circularized *in vitro*, circular forms of replicating T7 DNA are extremely rare *in vivo*, and linear molecules of three to ten genome equivalents, with only a single copy of the 160-bp terminal repeat at each genome junction, are found. As replication proceeds, the linear concatemers are converted into highly branched recombining and replicating structures that contain >100 genome equivalents. Recombination is independent of the host recombination functions and, in addition to replication functions, requires the single-strand DNA binding and strand-exchange properties of gp2.5, plus the gp4 helicase, the gp6 exonuclease, and the **Holliday junction**-resolving activity of the gp3 endonuclease (35-37). The latter ultimately resolves the recombining DNA into linear concatemers that are the substrate for the packaging machinery. Gp2.5 interacts physically with T7 DNA polymerase, and both recombination and bidirectional replication requires gp2.5. In addition to their role in replication and recombination, the gp3 and gp6 nucleases completely degrade the host chromosome; ~80% of the nucleotides necessary for replication are derived from *E. coli* DNA. Thus T7 DNA replication and recombination are inseparable. The overall rate of T7 DNA synthesis is very fast, and more than 200 genome equivalents are made in less than 10 min at 30°C. This high rate suggests that numerous [replication forks](#) are active, consistent with the idea that recombination intermediates can be converted into replication forks, as in **T4 phage**-infected cells. However, unlike T4 (and many other phages), T7 DNA replication is not necessary for late transcription (6). Mutants that fail to replicate their DNA have normal patterns of transcription; in fact, transcription of T7 DNA is almost over before DNA replication is initiated. The many nicks that accompany DNA replication and recombination are sealed by DNA ligase, either of the host or of T7. T7 ligase mutants plate normally, although their burst size is reduced on most *E. coli* strains. When *E. coli* ligase is defective, however, T7 ligase (gp1.3) becomes an essential protein (38).

T7 replication is extremely fast and, unless it is blocked, *E. coli* chromosome-derived nucleotides cannot be detected as acid-soluble material. Thymidine incorporation into DNA is a misleading measure of T7 replication rates because the nucleotide pools vary enormously throughout an infection. If the onset of replication is delayed, as after infection by 3.5 mutants, thymidine-uptake measurements underestimate the rate of DNA synthesis. Conversely, when nucleotide pools are

nearly exhausted near the time of cell lysis, thymidine uptake measurements grossly overestimate the replication rate.

Much of our understanding of T7 DNA replication has come from *in vitro* studies (Fig. 3). At the primary origin of replication, T7 RNA polymerase synthesizes primers of 10 to 60 nucleotides from both promoters upstream of the origin; there is no defined site where RNA synthesis stops and DNA synthesis begins. A hexamer of the gp4 primase/helicase, bound to the single-stranded DNA through its central hole, translocates 5' → 3' on the [Okazaki fragment](#) template strand (39). Hydrolysis of ribo- or deoxynucleoside triphosphates (dTTP being preferred *in vitro*) can be used as the energy source to unwind unreplicated DNA. Based upon the similarity of the nucleotide-binding change mechanism of the helicase with that of the F<sub>1</sub>F<sub>0</sub>-ATPase (see [ATP Synthase](#)), it has been suggested that rotational movement around the single-stranded DNA leads to unidirectional translocation along single-stranded DNA and consequent unwinding of double-stranded DNA (40). The helicase also interacts physically with both gp2.5 and DNA polymerase, to help make DNA synthesis more processive. Relative to the gp4B helicase, gp4A contains an extra 63 amino acid residues at the *N*-terminus that confer primase activity. The *N*-terminal extension contains a **zinc finger** that interacts with the sequences 3'-CTGG(G/T) or 3'-CTGTG in single-stranded DNA and synthesizes ribonucleotide primers 5'-ACC(A/C) or 5'-ACAC, which serve to initiate synthesis of Okazaki fragments (41). Mixed hexamers of gp4A and gp4B are formed *in vitro*, and phage mutants that make only gp4A are viable, suggesting that helicase unwinding is not constant, periodically pausing to allow primer synthesis. Synthesis on the leading and lagging DNA strands is tightly coupled. A replisomal complex consisting of T7 DNA polymerase, primase/helicase, and single-stranded DNA-binding protein simultaneously synthesizes both strands of DNA (42). Measurements of leading and lagging strand synthesis showed that they are coupled and are synthesized by distinct DNA polymerase molecules. The lagging strand polymerase cycles from one Okazaki fragment to the next; the length of the Okazaki fragments, which are initiated by [primase](#), are thought to be controlled by single-stranded DNA-binding protein. Furthermore, the DNA loop predicted by coordinated synthesis on leading and lagging strands was observed using electron microscopy.

**Figure 3.** Schematic of a T7 DNA replication fork. Two molecules of DNA polymerase (gray circles) synthesize processively the leading and lagging strands (thick lines). A hexamer of the primase-helicase (nut-shaped) unwinds the template duplex ahead of the leading strand polymerase and synthesizes the RNA primers (black ovals). Single-stranded DNA is coated with single-stranded DNA binding protein (gray stripe). Note that all three components, DNA polymerase, primase-helicase, and single-stranded DNA binding protein, must interact for coupling of leading and lagging strand synthesis (42).



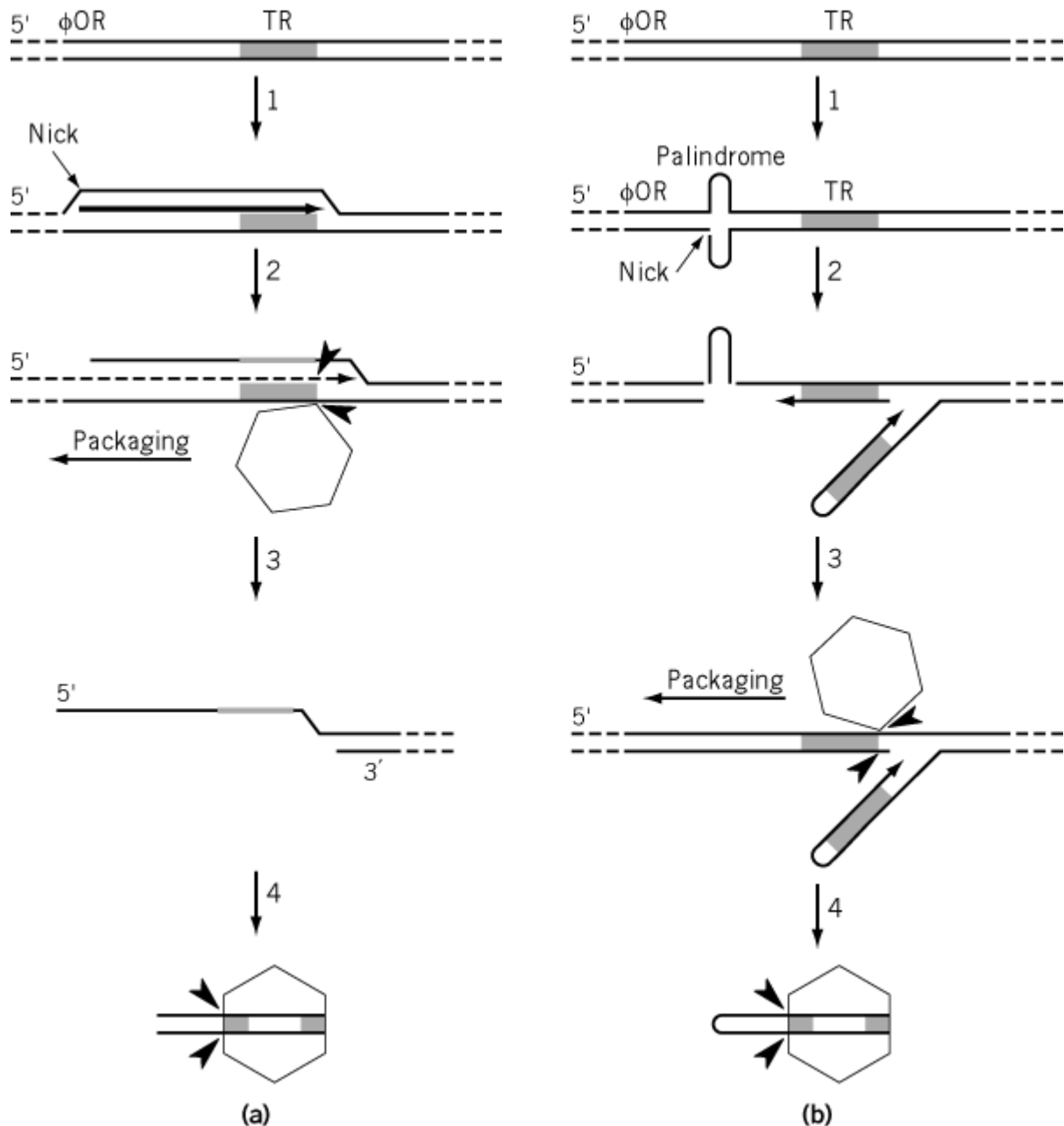
#### 4.3. DNA Maturation and Phage Morphogenesis

Packaging of T7 DNA occurs using preformed proheads. Assembly of these may be initiated using gp10 around a scaffold of gp9, generating an incomplete prohead shell, a process that would be followed by insertion of the connector–core complex (gp8-[gp14+gp15+gp16]). Alternatively, the connector–core complex may serve to nucleate the assembly of the scaffold and the capsid protein. Unlike several other phages, T7 is known to grow in *groE* mutants, but mutants defective in other **chaperonins** do not appear to have been tested and could facilitate the formation of the connector. The T3 connector (which can substitute for T7 gp8 in an otherwise T7 particle), when isolated from phage particles or filled heads, exhibits 12-fold rotational symmetry, but T7 12-mers and 13-mers are found when gene 8 is expressed from a clone (43). It has been suggested that the 13-mer is incorporated into proheads and converted into the 12-mer during DNA packaging, perhaps triggering conversion of the prohead into the mature capsid, with the concomitant loss of the scaffolding protein to allow DNA packaging (44). How the multitiered core structure is assembled is not yet known.

The substrate for DNA packaging is a concatemer consisting of about 5 to 10 genome equivalents, but each genome requires a second copy of the terminal repeat to be synthesized to generate the mature, full-length genome (Fig. 4). Transcription from the  $\phi$  OR promoter is known to cross the terminal repeat, but it pauses and terminates shortly thereafter. Termination at this site requires gp3.5 lysozyme, the absence of this protein leads to a defect in DNA maturation and packaging (45, 46). Transcription presumably originates from the upstream promoters  $\phi$  17 and/or  $\phi$  13 in  $\phi$  OR deletion mutants. Using an *in vitro* system, packaging was shown to be dependent on transcription, and the demonstrated specificity for packaging concatemeric DNA resides in the transcribing polymerase. Furthermore, an interaction between T7 RNA polymerase and the packaging protein gp19 has been shown by genetic **suppressor** analysis (47).

**Figure 4.** Two models for duplicating the terminal repeat and packaging T7 DNA (a) 1. Transcription from  $\phi$ OR allows a nick to be made in the displaced non-template strand, providing a 3'-OH. 2. Strand-displacement DNA synthesis crosses the terminal repeat. 3. DNA packaging is initiated by a prohead and terminase proteins (gp18 and gp19), which create a double-strand break, forming the mature right end of the genome. Removal of the left-hand

genome (via packaging) reveals a template-primer on the genome to the right for DNA synthesis to complete duplication of the terminal repeat. 4. Packaging is terminated by a double-strand cleavage forming the mature left end of the genome and release of a small duplex DNA. (Adapted from Ref. 48.) (b) 1. The base of a palindrome in genomic DNA is nicked, allowing a hairpin with a 3'-OH to form; the position of the nick is about 350 bp downstream of  $\phi OR$ . 2. DNA synthesis extends the hairpin and crosses the terminal repeat. Primase-initiated DNA synthesis on the other strand completes duplication of the terminal repeat. 3. DNA packaging is initiated by a prohead and terminase proteins (gp18 and gp19), which create a double-strand break, forming the mature right end of the genome. 4. Packaging is terminated by a double-strand cleavage, forming the mature left end of the genome, and release of a small single-stranded hairpin DNA. (Adapted from Ref. 49.)



One model posits that a nick is made in the single-stranded, nontemplate DNA in a transcription complex (48). This nick can be extended by strand-displacement DNA synthesis to cross the terminal repeat. Packaging is then initiated by a double-strand break to form the mature, genetic right end of the DNA, and it proceeds leftward through that genome. The 3'-OH on a genome not yet packaged can now serve to prime DNA synthesis on the displaced strand, to complete the duplication of the terminal repeat; the mature genetic left end of the genome is formed by a second double-strand

cleavage event. This idea does not explain the presence of a novel hairpin structure that is found in replicating DNA and appears to be a kinetic intermediate in the formation of the genetic left end (49). A second model suggests that a nick is created in the template strand, creating a 3'-OH to serve as a primer for DNA synthesis. Synthesis can then duplicate the terminal repeat and terminate whenever the replication fork encounters a prohead that is in the process of encapsidating DNA. This DNA synthesis generates a single strand that can be a template for primase-primed DNA synthesis, enabling the second strand of the terminal repeat to be copied. Double-strand breaks form both ends of the mature genome. This idea (49) accounts for the existence of the hairpin structure that is found attached to an immature genome genetic left end, but it does not account for the role of transcription and gp3.5-mediated transcription termination in DNA packaging. The model is also not supported by the fact that the site of nicking (and thus hairpin formation) can be deleted in phage that exhibit no severe growth defect. Perhaps elements of both models are correct, but it is likely that T7 has evolved more than one mechanism to duplicate the terminal repeat and generate a mature genetic left end. Secondary mechanisms certainly exist, because deletion mutants that remove more than about 15% of the T7 genome (the minimum size that is packagable) do not recognize the normal genetic left end in concatemeric DNA and package a genome whose left end therefore lies within normally internal sequences.

Gp18 and gp19 form the packaging enzyme; six molecules of gp19 associate with the connector of the prohead with the help of ATP as an **allosteric** effector (50). The association may be with the inner face of the connector, which has six domains surrounding the channel (9). The small subunit gp18 associates with the DNA, which then combines with the prohead gp19 complex. A site-specific, double-strand break on a DNA concatemer then initiates the packaging process. This endonuclease activity requires proheads for specificity; in their absence, gp19 alone possesses a nonspecific endonuclease activity that is inhibited by added gp18. Translocation of DNA into the prohead is driven by ATP hydrolysis at an estimated rate of 350 bp/s at 30°C, and 1.8 bp of DNA is packaged per ATP molecule hydrolyzed (50-52).

Despite the fact that lysis of the T7-infected host cell is more abrupt and complete than cells infected by most other phages, little is known of the process. Genes 3.5 and 17.5 are known to be required, perhaps together with 18.5 and the hypothetical 18.7, the latter potentially being the T7 homologues of *lRz* and *Rz1*. Gp17.5 appears to be the T7 holin, and it may disrupt the membrane to form nonrefractile cell ghosts, allowing access of gp3.5 to the peptidoglycan to effect lysis and release of phage from cell debris (53). Lysis is only slightly delayed by 3.5 deletion mutants, however, and mutants defective in DNA replication are equally delayed for lysis. These observations suggest that some other factor, perhaps a product of DNA replication or maturation, is responsible for initiating cell lysis.

Aside from lysis, the entire T7 growth cycle has been modeled as a coupled set of 93 ordinary differential equations, using parameters derived directly from experimental data (54). The computer simulation treats the T7 genome as containing 122 distinct genetic elements and accurately portrays the various stages of T7 development, including genome entry, gene expression, the switch from *E. coli* to T7 RNA polymerase-mediated transcription, and DNA replication. The latest version of the model allows a predictive assessment of the fitness (ie, phage production per unit time) of any mutant T7, including those whose genetic elements have been scrambled relative to the wild type.

## 5. Host-Parasite Interactions

A number of natural *E. coli* strains, as well as other enterobacterial hosts, are nonpermissive for T7. In many cases, the nonpermissive host carries a prophage or plasmid, but in other situations the host gene(s) involved have not been identified. The most interesting cases are those where inhibition of phage growth is not due to DNA restriction or defects in adsorption, but where the initial stages of infection appear normal but later events are grossly aberrant (55).

The *l rex* genes inhibit the growth of a 5.5 missense mutant of T7, perhaps by a similar mechanism



to that of *rex*-exclusion of T4 *rII* mutants (3, 16, 56). Gp5.5 also binds to the small DNA-binding protein H-NS; although the latter is not required for T7 development, an H-NS mutant has been isolated that then requires T7 to be 5.5<sup>+</sup> for growth (56). It is not known whether the above observations are related beyond the commonality of gp5.5. The plasmid Col Ib contains a gene that inhibits the growth of T7 0.7 mutants (57). These mutants also fail to grow on certain *rpoC* mutants of *E. coli*, and it seems likely that the defect in both cases involves *E. coli* RNA polymerase. To explain why the gp3 endonuclease does not degrade T7 DNA when it is degrading the bacterial chromosome, it was originally suggested that RNA polymerase, when bound to DNA, provides some signal to trigger degradation (58). Perhaps gene 2 fails to inactivate the mutant or Col Ib-modified RNA polymerase completely and causes an infection to abort.

Gene 1.2 mutants fail to grow on cells that overexpress *dgt*, a dGTPase; the enzyme is inhibited by gp1.2 (59). The lower dGTP pool size in *dgt* overexpressing mutants (*optA1*) aborts phage DNA replication and triggers breakdown of the intracellular DNA. Suppressors of a gene 1.2 null mutation have been isolated that allow phage growth in *optA1* mutants; these appear to be altered in gene 5, T7 DNA polymerase. Consistent with this idea is that several T4 gene 43 (DNA polymerase) mutator mutants that have increased 3' → 5' exonuclease activity also fail to grow on *optA1* mutants (60).

T7 and most of its relatives, T3 being an exception, fail to grow in cells harboring the F plasmid (55). The only F function required for exclusion is *pifA*, a gene that is nonessential for conjugation or plasmid replication or maintenance. The F *pifA* gene interferes with the normal function of T7 genes 1.2 and 10 and causes inhibition of all macromolecular synthesis in abortively infected cells (61). The infection is aborted before about half of the infecting phages have translocated their entire genome into the cell (61); in the remaining half of the population, transcription over the class III region of the genome is grossly aberrant (62). T7 requires three mutations to grow in F-containing cells, a missense or null mutation in gene 1.2 plus two missense mutations in gene 10 (63). T3 gene 1.2 is responsible for that phage growing in the presence of F, and T3 1.2 mutants respond like wild-type T7 when infecting F-containing strains. However, as in the case of T7, two missense mutations in T3 gene 10 can restore the ability of a 1.2 mutant phage to grow (64). The known functions of genes 1.2 and 10 do not yet provide a clear explanation for the inhibition of phage growth by F *pifA*. The proteins specified by gene 10 are not known to have additional roles beyond capsid shell formation, but normal capsid assembly by gp10 is not required for the infection to abort. The anti-dGTPase activity of gp1.2 also seems unrelated to exclusion by F, and T3 gene 1.2 mutants have been isolated that are excluded by F but retain the ability to grow on *optA1* hosts, and vice versa (65).

In the absence of other phage genes, plasmids expressing either T7 gene 1.2 or gene 10 kill cells containing PifA (66); the dysfunctions that occur as cells are dying are the same as those following an abortive phage infection, suggesting that the phage and F proteins interact with a host function. This cellular function has not been identified; it appears to be the same regardless of whether gp1.2 or gp10 interacts with PifA, and it can be protected by overexpression of the integral inner membrane FxsA protein (67). FxsA has no known independent function, and insertional inactivation or moderate overexpression of *fxsA* provide no obvious cellular phenotype. Nevertheless, if overexpression of an integral membrane protein suppresses F exclusion of T7, it is likely that the cellular location where T7 gp1.2 or gp10 interact with PifA is also the cell membrane.

Another abortive infection that involves expression of gene 10 is growth on *Shigella sonnei* D<sub>2</sub> 371-48 (68). T7 and most members of the family fail to grow on this strain; midway through the infection cycle, there is a massive breakdown of the phage DNA, and the infection is aborted. A Q → P change at codon 62 of T7 gene 10 prevents DNA degradation and allows phage growth. As in F exclusion of T7, the trigger for the abortive infection appears not to require capsid assembly. T3 grows normally on the *Shigella sonnei* host, even though its capsid protein has the same sequence around codon 62. The difference between T3 and T7 in this context appears to reside in gene 1.3, DNA ligase. T7-infected *Shigella sonnei* appear to be deficient in ligase activity, and gene 1.3 is

essential for T7 growth even when the gene *10-Q62P* mutation is present. Furthermore, a hybrid T7–T3 phage that contains T3 gene *1.3* productively infects the *Shigella sonnei* strain.

## 6. T7 as a Tool in Molecular Biology

### 6.1. Expression Vectors

Many RNAs and proteins can be produced in large quantities in *E. coli* using a gene expression system based on T7 RNA polymerase (69, 70). *E. coli* RNA polymerase does not recognize T7 promoters, and T7 RNA polymerase does not recognize *E. coli* or eukaryotic RNA polymerase promoters, nor does it recognize many sequences as transcription termination sites. T7 RNA polymerase is a single polypeptide chain enzyme that transcribes DNA at ~250bp/s at 30°C and is naturally resistant to rifampicin, an inhibitor of the *E. coli* enzyme. A wide variety of plasmid vectors have been developed for expressing genes under T7 RNA polymerase control, and they are available commercially, currently from Novagen, Stratagene, or Promega. The pET vectors, in particular, have features for a wide variety of cloning and expression projects, and gene *1* has also been cloned under arabinose promoter control in a pBAD plasmid. A typical vector places T7 gene *1* under *lac* promoter control, with the normal T7 promoter for the target gene modified by the addition of a *lac* operator sequence; thus both gene *1* and target gene expression are inducible by the addition of IPTG. If necessary, basal level synthesis of T7 RNA polymerase can be inhibited by the presence of the gene 3.5 plasmids pLysS or pLysE (S and E refer to the orientation [silent or expressed] of gene 3.5 relative to the vector tetracycline promoter). pLysS and pLysE are pACYC184-based plasmids that are compatible with the majority of cloning vector plasmids; pLysE is usually only necessary when the target gene is very toxic to *E. coli*. Addition of **IPTG** induces gene *1* to levels that overcome the inhibitory activity of the pLys plasmids. Use of the pLys plasmids often provides a secondary benefit, in that freeze-thawing cells, or treating them with the nonionic detergent Triton X-100, allows the gp3.5 lysozyme access to the peptidoglycan and promotes efficient cell lysis. In cases where the target gene is extremely toxic, gene *1* can also be delivered via infection by  $\lambda$  CE6, a  $\lambda$  phage derivative that expresses gene *1* from the  $\lambda$  p<sub>L</sub> promoter. The most commonly used host for overexpression is BL21(DE3), which contains gene *1* on a  $\lambda$  prophage. BL21(DE3) lacks a Type I restriction system and the [proteinases](#) Lon and OmpT, features that can help cloning and prevent **protein degradation** during purification. It is not always appreciated that toxic proteins often turn on the *E. coli* [SOS response](#) system and induce the prophage from BL21(DE3); in such cases, the *E. coli* *recA* strain HMS174(DE3) can be a useful host. Other frequent causes of the failure to express at high levels include attempts to induce at saturation density, where the cells are no longer making much protein, and loss of the plasmid itself, or mutation of the target gene, either of which can occur if the gene product is toxic to the cell. The [ampicillin](#) used with many vectors can be degraded by any secreted **b-lactamase**, leading to the loss of selection for plasmid maintenance. Vectors that employ other antibiotics are less of a problem, but even these can be abused by serially propagating strains.

### 6.2. Display Vectors

T7 has also been modified as an [epitope](#) or protein display vector that is currently available as part of a kit from Novagen. The phage particle is stable to 1% [SDS](#), 4 M [urea](#) or 2 M **guanidinium** chloride, 10 mM [EDTA](#), 0.1 M [dithiothreitol](#), and a pH range from 4 to 10. Note that T7 is not stable to the concentrations of acetic acid commonly used with the filamentous phage display systems, but the procedures for selecting T7 from a display library are otherwise comparable. The stability of the phage particle, plus the fact that the capsid protein gp10 exists in two forms, gp10A and gp10B, which are both incorporated into the phage particle, combine to make T7 a useful display tool. The frameshifted, C-terminal extension of gp10B is exposed on the outer faces of the icosahedral capsid, but it is not essential for phage growth; viable phage mutants exist that contain only gp10A or only gp10B. The usual display vectors make no gp10A and a truncated form of gp10B. Cloning sites have been introduced at the 3' end of the modified gene *10B*, and a very high efficiency *in vitro* packaging extract is available. Foreign C-terminal extensions of up to about 50 amino acid residues have been engineered to replace the natural C-terminal extension of gp10B, and the resulting phages then

contain 415 copies of the fusion protein on the capsid surface. Extensions longer than about 50 residues can also be made, but these require use of a complementing plasmid that supplies gp10A; using current protocols, a 1200-amino-acid residue extension has been displayed on the phage surface at an average copy number of 0.1 to 1 per phage particle. Note that there is a size limit on the phage genome for efficient packaging into viable particles, and the T7 strains used to display short or long extensions are different.

### 6.3. A Vector for Studying Codon Usage

The gene 9 mRNA has been manipulated to provide vectors for examining the effects of rare, or low-usage, codons. Insertions can be made after codon 13, 223, or 307 of a 313-codon test gene (71). Insertion of nine consecutive CUA [leucine](#) codons after codon 13, but not after codons 223 or 307, strongly inhibited translation of the modified RNA, without concomitant inhibition of translation of other CUA-containing mRNAs. Surprisingly, this elegant system has not been widely employed to examine the effects of other low-usage codons.

### Bibliography

1. K. H. Korsten, C. Tomkiewicz, and R. Hausmann (1979) The strategy of infection as a criterion for phylogenetic relationships of non-coli phages morphologically similar to phage T7. *J. Gen. Virol.* **43**, 57–73.
2. R. Hausmann (1988) "The T7 Group. In" *The Bacteriophages*, Vol. 1, R. Calendar, ed., Plenum Press, New York, pp. 259–289.
3. J. J. Dunn and F. W. Studier (1983) Complete nucleotide sequence of bacteriophage T7 DNA and the locations of genetic elements. *J. Mol. Biol.* **166**, 477–535.
4. B. G. Condrón, R. F. Gesteland, and J. F. Atkins (1991) An analysis of sequences stimulating frameshifting in the decoding of gene 10 of bacteriophage T7. *Nucleic Acids Res.* **19**, 5607–5612.
5. J. P. Condreay, S. E. Wright, and I. J. Molineux (1989) Nucleotide sequence and complementation studies of the gene 10 region of bacteriophage T3. *J. Mol. Biol.* **207**, 555–561.
6. F. W. Studier (1972) Bacteriophage T7. *Science* **176**, 367–376.
7. A. C. Steven and B. L. Trus (1986) In *Electron Microscopy of Proteins*, Vol. 5: *Viral Structure* (J. R. Harris and R. W. Horne, eds.), Academic Press, New York, pp. 1–35.
8. E. Kocsis, M. E. Cerritelli, B. L. Trus, N. Cheng, and A. C. Steven (1995) Improved methods for determination of rotational symmetries in macromolecules. *Ultramicroscopy* **60**, 219–228.
9. J. M. Valpuesta and J. L. Carrascosa (1994) Structure of viral connectors and their function in bacteriophage assembly and DNA packaging. *Q. J. Biophys.* **27**, 107–155.
10. P. Serwer (1976) Internal proteins of bacteriophage T7. *J. Mol. Biol.* **107**, 271–291.
11. M. E. Cerritelli, N. Cheng, A. H. Rosenberg, C. E. McPherson, F. P. Booy, and A. C. Steven (1997) Encapsidated conformation of bacteriophage T7. *Cell* **91**, 271–280.
12. P. Demchick and A. L. Koch (1996) The permeability of the wall fabric of *Escherichia coli* and *Bacillus subtilis*. *J. Bacteriol.* **178**, 768–773.
13. L. R. García and I. J. Molineux (1995) Rate of translocation of bacteriophage T7 DNA across the membranes of *Escherichia coli*. *J. Bacteriol.* **177**, 4066–4076.
14. L. R. García and I. J. Molineux (1995) Incomplete entry of bacteriophage T7 DNA into F plasmid-containing *Escherichia coli* strains. *J. Bacteriol.* **177**, 4077–4083.
15. A. Meisel, T. A. Bickle, D. H. Kruger, and C. Schroeder (1992) Type III restriction enzymes need two inversely oriented recognition sites for DNA cleavage. *Nature* **355**, 467–469.
16. L. Lin (1992) Ph.D. dissertation, SUNY, New York.
17. E. S. Robertson, L. A. Aggison, and F. W. Studier (1994) Phosphorylation of elongation factor G and ribosomal protein S6 in bacteriophage T7-infected *Escherichia coli*. *Mol. Micro.* **11**, 1045–1057.
18. L. B. Rothman-Denes, S. Muthukrishnan, R. Haselkorn, and F. W. Studier (1973) "A T7 Gene

Function Required for Shut-off of Host and Early T7 Transcription. In" *Virus Research* (C. F. Fox and W. S. Robinson, eds.), Academic Press, New York. pp. 227–239.

19. R. Sousa, Y. J. Chung, J. P. Rose, and B. C. Wang (1993) Crystal structure of bacteriophage T7 RNA polymerase at 3.3 Å resolution. *Nature* **364**, 593–599.
20. K. A. Chapman, S. I. Gunderson, M. Anello, R. D. Wells, and R. R. Burgess (1988) Bacteriophage T7 late promoters with point mutations: quantitative footprinting and *in vivo* expression. *Nucleic Acids Res.* **16**, 4511–4524.
21. J. F. Klement, M. B. Moorefield, E. Jorgensen, J. E. Brown, S. Risman, and W. T. McAllister (1990) Discrimination between bacteriophage T3 and T7 promoters by the T3 and T7 RNA polymerases depends primarily upon a three base-pair region located 10 to 12 base-pairs upstream from the start site. *J. Mol. Biol.* **215**, 21–29.
22. C. A. Raskin, G. A. Diaz, and W. T. McAllister (1993) T7 RNA polymerase mutants with altered promoter specificities. *Proc. Natl. Acad. Sci. USA.* **90**, 3147–3151.
23. S. S. Lee and K. C. Kang (1992) A two-base-pair substitution in T7 promoter by SP6 promoter-specific base pairs alone abolishes T7 promoter activity but reveals SP6 promoter activity. *Biochem. Int.* **26**, 1–5.
24. M. Rong, B. He, W. T. McAllister, and R. K. Durbin (1998) Promoter specificity determinants of T7 RNA polymerase. *Proc. Natl. Acad. Sci. USA.* **95**, 515–519.
25. L. E. Macdonald, R. K. Durbin, J. J. Dunn, and W. T. McAllister. 1994. Characterization of two types of termination signal for bacteriophage T7 RNA polymerase. *J. Mol. Biol.* **238**, 145–158.
26. X. Zhang and F. W. Studier (1997) Mechanism of inhibition of T7 RNA polymerase by T7 lysozyme. *J. Mol. Biol.* **269**, 964–981.
27. X. Cheng, X. Zhang, J. W. Pflugrath, and F. W. Studier (1994) The structure of bacteriophage T7 lysozyme, a zinc amidase and an inhibitor of T7 RNA polymerase. *Proc. Natl. Acad. Sci. USA.* **91**, 4034–4038.
28. S. Tabor, H. E. Huber, and C. C. Richardson (1987) *Escherichia coli* thioredoxin confers processivity on the DNA polymerase activity of the gene 5 protein of bacteriophage T7. *J. Biol. Chem.* **262**, 16212–16223.
29. H. E. Huber, M. Russel, P. Model, and C. C. Richardson (1986) Interaction of mutant thioredoxins of *Escherichia coli* with the gene 5 protein of phage T7. The redox capacity of thioredoxin is not required for stimulation of DNA polymerase activity. *J. Biol. Chem.* **261**, 15006–15012.
30. Y. T. Kim, S. Tabor, J. E. Churchich, and C. C. Richardson (1992) Interactions of gene 2.5 protein and DNA polymerase of bacteriophage T7. *J. Biol. Chem.* **267**, 15032–15040.
31. S. Doublé, S. Tabor, A. M. Long, C. C. Richardson, and T. Ellenberger (1998) Crystal structure of a bacteriophage T7 DNA replication complex at 2.2 Å resolution. *Nature* **391**, 251–258.
32. E. Bedford, S. Tabor, and C. C. Richardson (1997) The thioredoxin binding domain of bacteriophage T7 DNA polymerase confers processivity on *Escherichia coli* DNA polymerase I. *Proc. Natl. Acad. Sci. USA.* **94**, 479–484.
33. S. D. Rabkin and C. C. Richardson. 1988. Initiation of DNA replication at cloned origins of bacteriophage T7. *J. Mol. Biol.* **204**, 903–916.
34. S. D. Rabkin and C. C. Richardson (1990) *In vivo* analysis of the initiation of bacteriophage T7 DNA replication. *Virology* **174**, 585–592.
35. D. Kong, N. G. Nossal, and C. C. Richardson (1997) Role of the bacteriophage T7 and T4 single-stranded DNA-binding proteins in the formation of joint molecules and DNA helicase-catalyzed polar branch migration. *J. Biol. Chem.* **272**, 8380–8387.
36. D. Kong and C. C. Richardson (1996) Single-stranded DNA binding protein and DNA helicase of bacteriophage T7 mediate homologous DNA strand exchange. *EMBO J.* **15**, 2010–2019.
37. B. De Massy, R. A. Weisberg, and F. W. Studier (1987) Gene 3 endonuclease of bacteriophage T7 resolves conformationally branched structures in double-stranded DNA. *J. Mol. Biol.* **193**,

38. F. W. Studier (1973) Genetic analysis of non-essential bacteriophage T7 genes. *J. Mol. Biol.* **79**, 227–236.
39. E. H. Egelman, X. Yu, R. Wild, M. M. Hingorani, and S. S. Patel (1995) Bacteriophage T7 helicase/primase forms rings around single-stranded DNA that suggest a general structure for hexameric helicases. *Proc. Natl. Acad. Sci. USA* **92**, 3869–3873.
40. M. M. Hingorani, M. T. Washington, K. C. Moore, and S. S. Patel (1997) The dTTPase mechanism of T7 DNA helicase resembles the binding change mechanism of the F<sub>1</sub>-ATPase. *Proc. Natl. Acad. Sci. USA* **94**, 5012–5017.
41. H. Nakai and C. C. Richardson (1986) Interactions of the DNA polymerase and gene 4 protein of bacteriophage T7: protein–protein and protein-DNA interactions involved in RNA-primed DNA synthesis. *J. Biol. Chem.* **266**, 9818–9830.
42. J. Lee, P. D. Chastain II, T. Kusakabe, J. D. Griffith, and C. C. Richardson (1998) Coordinated leading and lagging strand synthesis on a mini-circular template. *Mol. Cell* **1**, 1001–1010.
43. E. Kocsis, M. E. Cerritelli, B. L. Trus, N. Cheng, and A. C. Steven (1995) Improved methods for determination of rotational symmetries in macromolecules. *Ultramicroscopy* **60**, 219–28.
44. M. E. Cerritelli and F. W. Studier (1996) Purification and characterization of T7 head-tail connectors expressed from the cloned gene. *J. Mol. Biol.* **258**, 299–307.
45. D. L. Lyakhov, B. He, X. Zhang, F. W. Studier, J. J. Dunn, and W. T. McAllister (1997) Mutant bacteriophage T7 RNA polymerases with altered termination properties. *J. Mol. Biol.* **26**, 28–40.
46. X. Zhang and F. W. Studier (1997) Mechanism of inhibition of T7 RNA polymerase by T7 lysozyme. *J. Mol. Biol.* **269**, 964–981.
47. X. Zhang and F. W. Studier (1995) Isolation of transcriptionally active mutants of T7 RNA polymerase that do not support phage growth. *J. Mol. Biol.* **250**, 156–168.
48. H. Fujisawa and M. Morita (1997) Phage DNA Packaging. *Genes to Cells* **2**, 537–545.
49. Y.-B. Chung and D. C. Hinkle (1990) Bacteriophage T7 DNA packaging. III A “hairpin” end formed on T7 DNA concatemers may be an intermediate in the processing reaction. *J. Mol. Biol.* **216**, 939–948.
50. H. Shibata, H. Fujisawa, and T. Minagawa (1987) Characterization of the bacteriophage T3 DNA packaging reaction in vitro in a defined system. *J. Mol. Biol.* **196**, 845–851.
51. M. Morita, M. Tasaka, and H. Fujisawa (1993) DNA packaging ATPase of bacteriophage T3. *Virology* **193**, 748–752.
52. M. Morita, M. Tasaka, and H. Fujisawa (1995) Structural and functional domains of the DNA packaging protein of bacteriophage T3: importance of the C-terminal region of the large subunit in prohead binding. *J. Mol. Biol.* **245**, 635–644.
53. R. Young (1992) Bacteriophage lysis: mechanism and regulation. *Microbiol. Rev.* **56**, 430–481.
54. D. Endy, D. Kong, and J. Yin (1997) Intracellular kinetics of a growing virus: a genetically structured simulation for bacteriophage T7. *Biotechnol. Bioeng.* **55**, 375–389.
55. I. J. Molineux (1991) Host–parasite interactions: recent developments in the genetics of abortive phage infections. *New Biol.* **3**, 230–236.
56. Q. Liu and C. C. Richardson (1993) Gene 5.5 protein of bacteriophage T7 inhibits the nucleoid protein H-NS of *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **90**, 1761–1765.
57. B. Gomez and L. Nualart (1977) Requirement of the bacteriophage T7 *0.7* gene for phage growth in the presence of the Col 1b factor. *J. Gen. Virol.* **35**, 99–106.
58. P. Q. Mooney, R. North, and I. J. Molineux (1980) The Role of Bacteriophage T7 Gene 2 Protein in DNA Replication. *Nucleic Acids Res.* **8**, 3043–3053.
59. S. M. Wurgler and C. C. Richardson (1990) Structure and regulation of the dGTP triphosphohydrolase from *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **87**, 2740–2744.

60. P. Gauss, D. H. Doherty, and L. Gold (1983) Bacterial and phage mutations that reveal helix unwinding activities required for bacteriophage T4 DNA replication. *Proc. Natl. Acad. Sci. USA* **80**, 1669–1673.
61. L. R. García and I. J. Molineux (1995) Incomplete entry of bacteriophage T7 DNA into F plasmid-containing *Escherichia coli* strains. *J. Bacteriol.* **177**, 4077–4083.
62. P. J. Beck and I. J. Molineux (1991) Defective transcription of the right end of bacteriophage T7 DNA during an abortive infection of F plasmid-containing *Escherichia coli*. *J. Bacteriol.* **173**, 947–954.
63. I. J. Molineux, C. K. Schmitt, and J. P. Condreay (1989) Mutants of bacteriophage T7 that escape F restriction. *J. Mol. Biol.* **207**, 563–574.
64. J. P. Condreay and I. J. Molineux (1989) Synthesis of the capsid protein inhibits development of bacteriophage T3 mutants that abortively infect F plasmid-containing strains. *J. Mol. Biol.* **207**, 543–554.
65. M. P. Schmitt, P. J. Beck, C. A. Kearney, J. L. Spence, D. DiGiovanni, J. P. Condreay, and I. J. Molineux (1987) Sequence of a conditionally essential region of bacteriophage T3, including the primary origin of DNA replication. *J. Mol. Biol.* **193**, 479–495.
66. C. K. Schmitt and I. J. Molineux (1991) Expression of gene *1.2* and gene *10* of bacteriophage T7 is lethal to F plasmid-containing *Escherichia coli*. *J. Bacteriol.* **173**, 1536–1543.
67. W.-F. Wang (1997) Ph.D. dissertation, University of Texas at Austin.
68. P. J. Beck, J. P. Condreay, and I. J. Molineux (1986) Expression of the unassembled capsid protein during infection of *Shigella sonnei* by bacteriophage T7 results in DNA damage that is repairable by bacteriophage T3, but not T7, DNA ligase. *J. Bacteriol.* **167**, 251–256.
69. F. W. Studier, A. H. Rosenberg, J. J. Dunn, and J. W. Dubendorff (1990) Use of T7 RNA polymerase to direct expression of cloned genes. *Methods Enzymol.* **185**, 60–89.
70. A. H. Rosenberg and F. W. Studier (1998) "Producing Proteins in a T7 Expression System. In" *Cells: A Laboratory Manual* (D. L. Spector, R. Goldman, and L. Leinwand, eds.), Cold Spring Harbor Press, Cold Spring Harbor, NY.
71. E. Goldman, A. H. Rosenberg, G. Zubay, and F. W. Studier (1995) Consecutive low-usage leucine codons block translation only when near the 5' end of a message in *Escherichia coli*. *J. Mol. Biol.* **245**, 467–473.

### **Suggestions for Further Reading**

72. F. W. Studier (1969) The genetics and physiology of bacteriophage T7. *Virology* **39**, 562–574. (The introduction to laboratory experiments using T7.)
73. R. Hausmann (1973) The genetics of T-odd phages. *Annu. Rev. Microbiol.* **27**, 51–67.
74. R. Hausmann (1976) Bacteriophage T7 genetics. *Curr. Top. Microbiol. Immunol.* **75**, 77–110. (Two thoughtful reviews of T7 biology.)
75. D. H. Krüger and C. Schroeder (1981) Bacteriophage T3 and bacteriophage T7 virus–host cell interactions. *Microbiol. Rev.* **45**, 9–51. (A very comprehensive review of T7.)
76. D. H. Duckworth, J. Glenn, and D. J. McCorquodale (1981) Inhibition of bacteriophage replication by extrachromosomal genetic elements. *Microbiol. Rev.* **45**, 52–71. (A review on host–parasite interactions.)

### **T-Box Genes**

T-box genes are members of a rapidly growing and highly conserved family of [transcription factors](#) that share a region of [homology](#) to the **DNA-binding** domain (T-box) of the mouse *Brachyury* (or *T*) gene product ([1](#), [2](#)). Most T-box genes have been found through experiments designed to identify genes that have specific activities in embryonic development or that cause developmental defects. T-box genes have been isolated from a wide variety of vertebrates and invertebrates, including *Caenorhabditis elegans* (roundworm), *Drosophila melanogaster* (fruitfly), ascidians, echinoderms, amphioxus, *Xenopus laevis* (frog), *Gallus gallus* (chick), *Danio rerio* ([zebrafish](#)), *Mus musculus* ([mouse](#)), and humans. T-box genes are expressed during gastrulation and organogenesis in derivatives of all three germ layers (ie, ectoderm, mesoderm, and endoderm). During embryogenesis, the spatial and temporal expression patterns of T-box genes are unique, although overlapping, indicating that they are differentially regulated during development ([3](#)). Mutations in T-box genes have been associated with defects of development in *Drosophila* ([4](#)), zebrafish ([5](#)), mouse ([6](#)), and humans ([7-9](#)). **Phylogenetic** analyses suggest that the origin of the T-box gene family predates the divergence of arthropods and chordates from a common ancestor more than 600 million years ago ([10](#)). Consequently, it is likely that T-box genes play critical roles in the development of all animal species. Identifying the elements that control T-box gene expression and those that are activated by T-box gene products are major areas of investigation.

## 1. The Brachyury (T) Gene

In 1927, Dobrovolskĭa-Zavadskĭa first described the phenotype of a semidominant mouse mutant with a short, often slightly kinked tail ([11](#)). He named this locus *Brachyury* (Greek, *short-tail*) or *T* (*tail*). Mice that were heterozygous for a mutation at the *T* locus had only short tails, while homozygous *Brachyury* mutants were much more severely affected. They had a thickened primitive streak and lacked a notochord, and the mesoderm posterior to somite 7 was abnormal or absent. Additionally, the allantois of homozygotes failed to fuse with the chorion, causing embryonic death at 10.5 days postcoitum. Although *Brachyury* was one of the earliest well-documented mouse mutants, its molecular basis was not characterized for more than 60 years.

The mouse *T* gene, **cloned** in 1990, is the best studied of the T-box genes. It encodes a protein of 436 amino acid residues that localizes to the nucleus and activates [transcription](#) ([6](#)). The DNA-binding domain (T-box) of T is composed of 180 residues located in the amino-terminal portion of the protein. Analysis of deletion mutants indicates there are two activation domains and two repressor domains located toward the carboxyl terminus. [X-ray crystallography](#) demonstrates that the T-box of the frog homologue of *T*, *Xbra*, binds as a dimer to a 24-nucleotide **palindromic** DNA duplex, interacting with the major and minor grooves of DNA ([12](#)). The mechanism by which the T protein binds DNA appears to be unique among known transcription factors ([12](#)).

*Brachyury* is caused by a 200-kbp deletion that removes the entire *T* gene and effectively produces a [null mutation](#). Consequently, the observed defects are thought to result from haploinsufficiency of T protein. In contrast, other *T* alleles (eg,  $T^{wis}$ ,  $T^c$ ,  $T^{c-2H}$ ) are due to insertions and deletions that cause frameshift mutations and truncated T proteins ([6](#), [13](#)). These alleles produce more severe developmental defects and may act as dominant negative mutations (ie, the mutant protein perturbs the function of the normal protein). Homologues of the *T* gene have been isolated in human (*T*), frog (*Xbra*), chick (*T*), and zebrafish (*no tail* (*ntl*) or *zf-T*).

The *T* gene is located in a region of mouse chromosome 17 defined by four inversions that are called the ***t* haplotype**. The *t* haplotype contains at least two additional genes, *small t* (*t*) and *Brachyury the second* (*T2*), alleles of which produce phenotypes similar to *Brachyury*. Until about 20 years ago, *t* haplotypes were thought to be mutant alleles of the *T* locus. However, genetic and molecular tools have shown that the *t* haplotype is a chromosomal region containing many different genes that affect development, male fertility, and recombination ([14](#)). Thus, the *T* gene is only one gene within the *t* haplotype.

The main defect in mice lacking T protein is a failure to undergo appropriate gastrulation movements. Studies in the mouse have demonstrated that these aberrant movements during gastrulation may be related to the lack of formation of particular regions of mesoderm (eg, posterior mesoderm). This is consistent with the [messenger RNA](#) expression pattern of *T*: transcripts are first detected in the primitive streak at the onset of gastrulation and persist in this area as gastrulation proceeds, and *T* is expressed longer in cells fated to become axial and posterior mesoderm, the same regions absent in *Brachyury* (6, 13).

The actions of *T* homologues in zebrafish (*zf-T*) and frog (*Xbra*) appear to be similar to the roles of *T* in mouse, suggesting that *T* has been conserved functionally. Mutations in *zf-T* produce a phenotype similar to *Brachyury* in which caudal mesoderm is deficient (5). Overexpression of a dominant-negative *Xbra* construct produces a mutant phenotype in *Xenopus* that resembles homozygous *T* mutants in mouse and zebrafish (15). Mutations in the *dm-Trg*, the *Drosophila* ortholog of *T* ([orthologous genes](#) are genes that are direct descendants of a common ancestor), affect posterior structures that may be analogous to the posterior structures of vertebrates (4).

The *T* gene is an immediate-early response gene (ie, transcription is directly activated by a [signal transduction](#) pathway). In frog, *Xbra* can be activated by mesoderm-inducing factors, such as members of the [transforming growth factor](#) b family (eg, activin) and the [fibroblast growth factor](#) family. Furthermore, *Xbra* can alter the fate cells in a dose-dependent fashion (16). Thus, *T* genes appear to control the specification and differentiation of mesoderm subpopulations in the developing embryo.

## 2. The T-Box Gene Family

For a brief period of time, the *T* gene and its orthologs were thought to be the only members of a newly described family of transcription factors. However, further characterization of the *T* gene stimulated a search for, and the eventual isolation of many additional genes sharing a region of homology with the DNA-binding domain of the T protein. These genes compose the T-box gene family (Table 1). The largest number of known T-box genes have been identified in the mouse. However, the first *T* homologue to be identified, *optomotor blind* (*OMB*), was isolated in a *Drosophila* mutant with abnormal optomotor-turning behavior and an inversion on the [X chromosome](#). The original variant allele is produced by a deletion in the **promoter** region that results in a dramatic reduction in *omb* expression. Other *omb* mutants, such as *bifid* and *Quadroon*, exhibit defects of wings and tergites, respectively. Thus, mutations in T-box genes may disrupt development of only specific subpopulations of mesoderm.

**Table 1. Characterized T-Box Genes in Vertebrates and Invertebrates**

| T-Box Gene   | C. <i>elegans</i> | <i>Drosophila</i> | Sea Urchin   | Ascidian                     | Amphioxus                        | <i>Xenopus</i> | Newt              | Zebrafish   | Chick                         | Mouse                   |
|--------------|-------------------|-------------------|--------------|------------------------------|----------------------------------|----------------|-------------------|-------------|-------------------------------|-------------------------|
|              |                   | <i>dm-Trg</i>     | <i>su-Ta</i> | <i>as-T</i><br><i>ci-bra</i> | <i>amBra-1</i><br><i>amBra-2</i> | <i>Xbra</i>    | <i>cp-T</i>       | <i>zf-T</i> | <i>ch-T</i><br><i>ch-TbxT</i> | <i>mu</i>               |
| <i>Tbr-1</i> |                   |                   |              |                              |                                  | <i>x-Eomes</i> |                   |             |                               | <i>mu</i><br><i>Tbi</i> |
| <i>Tbx1</i>  |                   | <i>org-1</i>      |              |                              |                                  |                | <i>nv-Tbox2-3</i> |             |                               | <i>mu</i><br><i>Tb:</i> |



|              |                |                |               |              |                 |            |
|--------------|----------------|----------------|---------------|--------------|-----------------|------------|
| <i>Tbx2</i>  | <i>ce-tbx2</i> | <i>dm-omb</i>  |               |              | <i>ch-</i>      | <i>mu</i>  |
|              |                |                |               |              | <i>Tbx2</i>     | <i>Tb:</i> |
| <i>Tbx3</i>  |                |                | <i>x-ET</i>   |              | <i>ch-</i>      | <i>mu</i>  |
|              |                |                |               |              | <i>Tbx3</i>     | <i>Tb:</i> |
| <i>Tbx4</i>  |                |                |               |              | <i>ch-</i>      | <i>mu</i>  |
|              |                |                |               |              | <i>Tbx4</i>     | <i>Tb:</i> |
| <i>Tbx5</i>  |                |                |               | <i>nv-</i>   | <i>ch-</i>      | <i>mu</i>  |
|              |                |                |               | <i>Tbox1</i> | <i>Tbx5</i>     | <i>Tb:</i> |
| <i>Tbx6</i>  |                |                | <i>x-VegT</i> |              | <i>zf-tbx6</i>  | <i>ch-</i> |
|              |                |                |               |              |                 | <i>mu</i>  |
|              |                |                |               |              | <i>Tbx6L</i>    | <i>Tb:</i> |
| <i>Tbx7</i>  | <i>ce-tbx7</i> |                |               |              |                 |            |
| <i>Tbx8</i>  | <i>ce-tbx8</i> |                |               |              |                 |            |
| <i>Tbx9</i>  | <i>ce-tbx9</i> |                |               |              |                 |            |
| <i>Tbx10</i> |                |                |               |              |                 | <i>mu</i>  |
|              |                |                |               |              |                 | <i>Tb:</i> |
| <i>Tbx11</i> | <i>ce-</i>     |                |               |              |                 | <i>mu</i>  |
|              | <i>tbx11</i>   |                |               |              |                 | <i>Tb:</i> |
| <i>Tbx12</i> | <i>ce-</i>     | <i>dm-HT15</i> |               |              |                 | <i>mu</i>  |
|              | <i>tbx12</i>   |                |               |              |                 | <i>Tb:</i> |
| <i>Tbx15</i> |                |                |               | <i>nv-</i>   |                 | <i>mu</i>  |
|              |                |                |               | <i>Tbox4</i> |                 | <i>Tb:</i> |
| <i>Tbx16</i> |                |                |               |              | <i>zf-tbx16</i> |            |
| <i>Tbx17</i> | <i>ce-</i>     |                |               |              |                 |            |
|              | <i>tbx17</i>   |                |               |              |                 |            |
| <i>T2</i>    |                | <i>as-T2</i>   |               |              |                 |            |

---

Phylogenetic analyses have been useful for reconstructing the evolutionary history and functional relationships among T-box genes. It appears that T-box genes arose from a common progenitor present in the most recent common ancestor of extant metazoans. Among different family members, the size of the T-box domain ranges from approximately 180–190 residues, and the amino acid sequence is highly conserved. However, the T-box may be located at various positions within each protein, and sequence similarities outside of the T-box are observed only for closely related family members. Thus, it appears that the DNA-binding activity between T-box family members has been retained, while the domains that may be involved in [protein–protein interactions](#) have diversified.

There are substantial data suggesting the existence of distinct subfamilies (eg, *T*, *Tbx2 / 3 / 4 / 5*, *Tbx6*, *Tbx1*, and *Tbr* subfamilies) of more closely related T-box genes that exhibit overlapping expression patterns and possibly functional similarities. Classification of specific T-box genes into subfamilies differs, depending on whether the genomic structure or T-box amino acid sequences are assessed (17). Most discrepancies arise, however, when distantly related genes are compared. For more closely related genes, different classification schemes produce similar results. As additional T-box genes are isolated and characterized, the complexity of subfamilies will likely increase.

### 3. The Role of T-Box Genes in Development

T-box genes have been implicated in a wide variety of developmental processes. The T-box genes *Tbx2*, *Tbx3*, *Tbx4*, and *Tbx5*, have well-defined orthologs in chick, mouse, and human. They appear to have arisen from an ancestral gene, *TBX2 / 3 / 4 / 5*, common to bony fishes and tetrapods.

Unequal **crossing over** produced the orthologous and tightly linked pairs *Tbx2 / 3* and *Tbx4 / 5*. Subsequent duplication of the chromosomal region generated *Tbx2/Tbx4* and *Tbx3/Tbx5* (10).

Expression patterns of these T-box genes have been well characterized in chick and mouse<sup>3</sup>, and some data are available for humans (7-9). Although there is moderate overlap between the expression patterns of nonorthologous genes, the expression patterns of orthologous genes are very similar between species. However, there are important differences between orthologous genes that may reflect functional recruitment of T-box genes for novel functions during evolution. For example, *Tbx4* and *Tbx5* are not expressed in chick allantois, but are expressed in the allantois of the mouse, suggesting that they may have been recruited in the evolution of the mammalian placental unit.

*Tbx5* and *Tbx4* are broadly expressed in the developing embryo, although a major focus of investigation has been the role of *Tbx5* and *Tbx4* in vertebrate limb development (18, 19). *Tbx5* and *Tbx4* are expressed in the lateral plate mesoderm prior to limb bud initiation. In the chick and mouse, *Tbx5* is expressed exclusively in the forelimb, while *Tbx4* expression is limited to the hindlimb. In ectopic limbs induced in chick flank, *Tbx5* and *Tbx4* are differentially expressed in the anterior and posterior halves, respectively (18). These data suggest that forelimb and hindlimb identity may be specified by *Tbx5* and *Tbx4*, respectively (19).

*Tbx6* is expressed in the paraxial mesoderm that will eventually form somite precursors. In the absence of *Tbx6*, cells destined to form posterior somites differentiate along a neural pathway (20). Thus, *Tbx6* appears to be required for specification of posterior paraxial mesoderm. Furthermore, the spatial expression pattern of *Tbx6* partially overlaps the expression pattern of *T*, suggesting that the two proteins may interact. This is further supported by the observation that *Tbx6* expression is downregulated in homozygous *Brachyury* mutants.

*Tbx1* expression is first detectable in the anterior embryonic mesoderm. It is expressed in the mesenchyme of the first, second, and third pharyngeal arches and the epithelium of the first, second, and third pharyngeal pouches. The pharyngeal pouches form the tonsils, parathyroids, and thymus.

*Tbr1* is expressed only in the postmitotic cells of the developing and adult telencephalon, and it may play a role in the differentiation of cells in the neocortex. Thus, it is somewhat distinct from other T-box genes (21).

#### 4. T-Box Genes and Birth Defects

The medical importance of T-box genes as a cause of human birth defects is well established. Mutations in T-box genes have been associated with two autosomal dominant multiple birth-defect syndromes. Mutations in *TBX3* cause limb, apocrine, genital, and dental defects in ulnar–mammary syndrome (UMS) (7), and mutations in *TBX5* cause limb and heart defects in Holt–Oram syndrome (HOS) (8, 9). Most of the mutations observed in *TBX3* and *TBX5* are small insertions or deletions that produce frameshifts and a truncated protein product. It has therefore been suggested that UMS and HOS are caused by haploinsufficiency of *TBX3* and *TBX5* proteins, respectively. Missense mutations have also been identified in *TBX3* and *TBX5*, although their effects on transcription and/or protein function are unknown.

The tissues affected in individuals with UMS and HOS are a subset of tissues in which *Tbx3* and *Tbx5* are expressed in chick and mouse. In other words, development of most tissues in which *Tbx3* and *Tbx5* are expressed is normal in individuals with these conditions. This suggests that different tissues may require quantitatively different levels of normal *TBX3* or *TBX5* protein. Alternatively, other genes including T-box genes with expression domains overlapping *TBX3* and *TBX5* may be able to compensate for reduced levels of normal *TBX3* and *TBX5* protein. Such functional redundancy is also observed for other families of transcription factors that play prominent roles in development (eg, HOX genes).

The limb defects observed in individuals with UMS are characterized by deletions or duplications of the distal posterior elements of the upper limb. In contrast, the distal anterior skeletal elements of the upper limb are disrupted in individuals with HOS. This suggests that *TBX3* and *TBX5* may partly control patterning of the anteroposterior axis of the vertebrate upper limb. Thus, analogous to the role of *T* in specifying axial mesodermal subpopulations, *TBX3* and *TBX5* may specify mesodermal subpopulations in the upper limb.

The human homologue of *Tbx1* maps to a region of chromosome 22 that is commonly deleted in individuals with velocardiofacial (VCF) syndrome (22). VCF syndrome is characterized by defects of the palate, hypoplasia or aplasia of the thymus, and conotruncal heart anomalies. These defects are concordant with the expression pattern of *Tbx1* in mouse. However, many additional genes are deleted in individuals with VCF, and thus the role of *TBX1* haploinsufficiency, if any, is unclear.

Since the identification of *T*, the T-box gene family has grown rapidly. Given the important and diverse roles that T-box genes play in embryogenesis, it will not be surprising if additional human birth-defect syndromes are found to be caused by mutations in T-box genes.

### Bibliography

1. R. J. Bollag, Z. Siegfried, J. A. Cebra-Thomas, N. Garvey, E. M. Davison, and L. M. Silver (1994) *Nat. Genet.* **7**, 383–389.
2. S. I. Agulnik, R. J. Bollag, and L. M. Silver (1995) *Genomics* **25**, 214–219.
3. D. L. Chapman et al. (1996) *Dev. Dyn.* **206**, 379–390.
4. A. Kispert, B. G. Herrmann, M. Leptin, and R. Reuter (1994) *Genes Dev.* **8**, 2137–2150.
5. S. Schulte-Merker, F. M. van Eeden, M. E. Halpern, C. B. Kimmel, and C. Nusslein-Volhard (1994) *Development* **120**, 1009–1015.
6. B. G. Herrmann, S. Labiet, A. Poustka, T. King, and H. Lehrach (1990) *Nature* **343**, 617–622.
7. M. Bamshad et al. (1997) *Nat. Genet.* **16**, 311–315.
8. C. T. Basson et al. (1997) *Nat. Genet.* **15**, 30–34.
9. Q. Y. Li et al. (1997) *Nat. Genet.* **15**, 21–29.
10. I. Ruvinsky and L. M. Silver (1997) *Genomics* **40**, 262–266.
11. N. Dobrovolskĭa-Zavadskaĭa (1927) *C. R. Seanc. Soc. Biol.* **97**, 114–116.
12. C. W. Muller and B. G. Herrmann (1997) *Nature* **389**, 884–888.
13. B. G. Herrmann (1991) *Development* **113**, 913–917.
14. L. M. Silver (1993) *Trends Genet.* **9**, 250–254.
15. F. L. Conlon, S. G. Sedgwick, K. M. Weston, and J. C. Smith (1996) *Development* **122**, 2427–2435.
16. J. C. Smith, B. M. J. Price, J. B. A. Green, D. Weigel, and B. G. Herrmann (1991) *Cell* **67**, 79–87.
17. S. Wattler, A. Russ, M. Evans, and M. Nehls (1998) *Genomics* **48**, 24–33.
18. H. Ohuchi et al. (1998) *Development* **125**, 51–60.
19. M. Logan, H. Simon, and C. Tabin (1998) *Development* **125**, 2825–2835.
20. D. L. Chapman and V. E. Papaioannou (1998) *Nature* **391**, 695–697.
21. C. Chieffo et al. (1997) *Genome* **43**, 267–277.
22. A. Bulfone et al. (1995) *Neuron* **15**, 63–78.

### Suggestions for Further Reading

23. V. E. Papaioannou and L. M. Silver (1998) The T-box gene family, *BioEssays* **20**, 9–19.
24. A. I. Kavka and J. B. A. Green (1997) Tales of tails: Brachyury and the T-box genes, *Biochem. Biophys. Acta* **1333**, F73–F84.

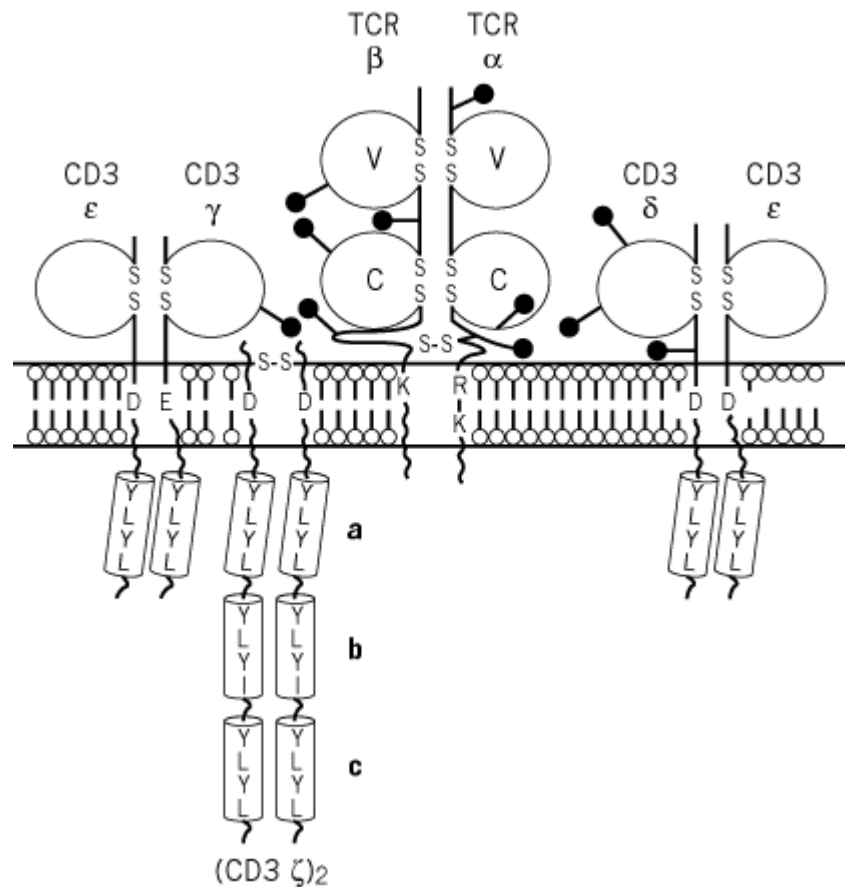
25. B. G. Herrmann and A. Kispert (1994) The *T* genes in embryogenesis, Trends Genet. **10**, 280–286.
26. J. Smith (1997) Brachyury and the T-box genes, Curr. Opin. Genet. Dev. **7**, 4374–4380.

## T-Cell Receptor (TCR)

Recognition of **antigens** by [T cells](#) is mediated by a specific receptor, called T-cell receptor (TCR). Each cell produces a unique TCR, as in the case of [immunoglobulin biosynthesis](#). TCR has a general architecture resembling that of **immunoglobulins** and the [B-cell receptor](#) (BCR), although there are three main differences between them: 1. TCR recognizes **peptides** derived from the protein antigen after it has been processed and presented by antigen-presenting cells in close association with **multiple histocompatibility complex** (MHC) class I or class II molecules. 2. TCRs have no [somatic mutations](#). 3. There is no soluble form of TCR, which always remains expressed at the surface of T cells.

TCR molecules are **disulfide-bonded** heterodimeric structures composed of one  $\alpha$  and one  $\beta$  [polypeptide chain](#), each having an immunoglobulin-like organization—that is, a variable NH<sub>2</sub>-terminal **variable domain**, a constant C-domain, a hydrophobic **transmembrane** region, and a short cytoplasmic COOH-terminus (Fig. 1). These TCR $\alpha\beta$  are found on over 90% of circulating T cells; in addition, another type has been described as TCR $\gamma\delta$ , which has a somewhat similar organization, centered on the corresponding  $\gamma\delta$  heterodimer, that is preferentially expressed in dermal and mucosal tissues.

**Figure 1.** Schematic domain organization of the TCR. The recognition module is that of a TCR $\alpha\beta$ . It is expressed at the cell surface with the CD3 signaling module, composed of  $\gamma$ ,  $\delta$ ,  $\epsilon$ ,  $\zeta$ , and  $\eta$  chains, assembled as indicated. CD3 chains contact protein **tyrosine kinases** by their ITAM (immunoreceptor tyrosine-based activation motif) motifs.



The organization of the TCR chain genes, as well as their controlled [gene rearrangements](#) by the same RAG1/RAG2 [recombinase](#), are very similar to those of immunoglobulins. In humans, the TCRB locus contains 41 to 46 Vb genes, organized in 24 families, that recombine randomly to two sets of D(6 J)C b gene segments; TCRA has 50 Va genes, an unusually large number of Ja gene segments (56), and one Ca gene. Organization of the TCRD locus is peculiar, in that it is entirely contained within the TCRA locus; when the latter has rearranged, the TCRD locus is completely deleted. TCRD has a VDJ organization, with three Vd, four Jd, and one Cd. Finally the TCRG locus contains six Vg, five Jg, and two Cg genes. Similar to the situation described for the Ig heavy chains, the complementarity determining region 3 (CDR3) of the b chain is the most variable area of T-cell receptors (see [Immunoglobulin Structure](#)).

Another similarity to the BCR organization is the close association of the TCR heterodimer with the CD3 signaling module (Fig. 1). CD3 has a complex molecular structure, composed of five discrete polypeptide chains, g, d,  $\epsilon$ , z, and h. Three have an external Ig-like domain (although its similarity with immunoglobulin (Ig) sequences is rather low), a transmembrane region, and a small cytoplasmic region. These chains are organized as two heterodimers, g $\epsilon$  and d $\epsilon$ , and contain on each of their cytoplasmic region one *immunoreceptor tyrosine-based activation motif* (ITAM) that provides interactions with cytosolic protein kinases, which are activated after stimulation of the receptor. Transmembrane regions of the three g, d, and  $\epsilon$  chains contain a negatively charged residue that stabilizes the interaction with the TCR, which has positively charged residues in the corresponding transmembrane area. The z chains have only a very short external segment, a transmembrane region, and a long cytoplasmic region that contains three ITAMs that, upon **phosphorylation** by Fyn or Lck, bind and activate an important partner of the T-cell transduction cascade, ZAP 70. They form homodimers or, alternatively, heterodimers with a spliced form known as the h chain.

Two other molecules are of particular interest at the surface of T cells, namely CD4 and CD8, considered to be co-receptors because they bind MHC class II or class I, respectively, of antigen-presenting cells and thus may play an important role in docking the TCR onto the MHC–peptide complex. CD4 is a monomer consisting of four external Ig-like domains; D1 and D2 form a rigid rod 60 Å long, connected to D3 and D4 by a flexible hinge and one cytoplasmic region, which contacts the Lck tyrosine kinase. CD4 binds to the b2 domain MHC class II molecules, independently of the TCR–MHC–peptide complex, although there may be effects on their binding energies. CD8 is a disulfide-bonded a b heterodimer; each chain has a single Ig-like domain and a long extended segment that terminates with a cytoplasmic tail that also contacts Lck. CD8 binds to the a3 domain of MHC class I molecules.

The general TCR three-dimensional structure also resembles that of Ig domains; it consists of b-sheets and external loops that protrude out of the main structure. These loops make up the Va and Vb domains and are directly involved in recognition of the peptide–MHC complex, as recently confirmed by X-ray crystallographic structures determined independently by several laboratories. These loops contain the CDRs, which are also similar to those of Ig H and L chains (see [Immunoglobulin Structure](#)). Of prime interest is the first high-resolution structure of the complex of a human TCRA b bound to an HLA-A2-nonapeptide; it revealed that the TCR is oriented diagonally across the peptide binding site of the MHC molecule, with CDR1a lying between the MHC a1 and a2 helices, whereas CDR2a interacts with the MHC a2 helix. The CDR3a and CDR3b loops contact both the central and COOH-terminal regions of the bound peptide, but they also interact with the MHC a1 and a2 helices. CDR1a and CDR1b each contribute a single peptide-binding residue. CDR1b and CDR2b have no direct contact with the complex, although they approach the COOH-terminus of the MHC a1 helix. It would now be of great interest to get the complete structure, including the CD4 or CD8 co-receptors. A crystal structure of a complex between a CD8a homodimer and HLA-A2 has recently been reported, making unlikely a previous hypothesis that CD8 might also contact the TCR/CD3 ectodomains. One cannot exclude, however, indirect constraints on the detailed structure of the TCR, through conformations induced as a result of interactions between cytoplasmic partners that initiate the transduction cascade.

Finally, it is likely, in view of other recently published structures, that variations on the same general diagonal orientation provide a general view of TCR–MHC–peptide interactions, especially because there is an obvious degree of degeneracy in overall T-cell recognition, like that already mentioned for antibodies.

See also entries **Antigen presentation**, [Epitope](#), [Gene Rearrangement](#), [Recombinase](#), and [T Cell](#).

#### Suggestions for Further Reading

M. M. Davis and P. J. Bjorkman (1988) T-cell antigen receptor genes and T-cell recognition. *Nature* **334**, 395–401.

D. N. Garboczi, P. Ghosh, U. Utz, Q. R. Fan, W. E. Biddison, and D. C. Wiley (1996) Structure of the complex between human T cell receptor, viral peptide and HLA-A2. *Nature* **384**, 134–141.

G. Mazza, D. Housset, C. Piras, C. Gregoire, S. Y. Lin, J. C. Fontecilla-Camps, and B. Malissen (1998) Glimpses at the recognition of peptide/MHC complexes by T-cell antigen receptors. *Immunol. Rev.* **163**, 187–196.

#### **T-Complex, -DNA, -Region, -Strand**

The T-complex is the **DNA** and associated [proteins](#) that are transferred from [Agrobacterium](#) to the **plant** cell during *Agrobacterium*-mediated cell [transformation](#). This process is a fascinating example of interkingdom **gene** transfer. It is the physical basis of tumor formation by *Agrobacterium* and the basis of **vectors** used in the creating **transgenic** plants by *Agrobacterium*-mediated transformation that is so vital to [plant genetic engineering](#).

The transferred, or T-, DNA of the [Ti plasmids](#) and [Ri plasmids](#) is the region of DNA transferred from the bacterium to the infected plant cell during transformation. The T-DNA is flanked and defined by 25-bp imperfect sequence repeats: TGGCAGGATATATTC(*or* G)XG(*or*A)T(*or*G)TGTA AAA(*or*T)T(*or*C) (1). Different strains of bacteria contain Ti plasmids that have different T-DNA. Nopaline strains (eg, C58, T37) contain a Ti plasmid that has a single T-DNA of approximately 24 kbp. Octopine strains (eg, B6S3, Ach5) contain two regions of T-DNA, a 14-kbp T<sub>L</sub> region and a 7-kbp T<sub>R</sub> region. The T<sub>L</sub> region is normally present in plant cells transformed by octopine strains of bacteria, and this region is the functional equivalent of the right-hand portion of the nopaline T-DNA. Interestingly, the T<sub>R</sub> region is not always present in transformed cells and indeed, if present, may not be linked to the T<sub>L</sub> region, suggesting independent transfer of the two DNAs.

The genes encoded by the T-DNA of Ti plasmids are generally involved in biosynthesizing the **phytohormones**, [auxins](#), and [cytokinins](#), or the **opines**. Opines are a novel source of fixed carbon and nitrogen synthesized by the host plant cell and are catabolized by the infecting bacteria. The Ti plasmid is generally grouped according to the opine whose biosynthetic gene is encoded by the T-DNA. The genes are not transcribed in the bacteria, but they become active once the T-DNA has integrated into the plant [genome](#). Nopaline T-DNA encodes at least 13 transcripts, and the protein products of 7 have been defined: *iaaM* and *iaaH*, which convert L-tryptophan to the auxin indole-3-acetic acid, *via* indole-3-acetamide; *ipt*, that encodes isopentyl transferase, which mediates the synthesis of cytokinin; gene 5 leads to the accumulation of indole-3-lactate, an auxin antagonist; gene 6a is thought to be involved in secreting opines; nopaline synthase (*nos*) and agrocinopine synthase (*acs*) are genes responsible for synthesizing opines. Two genes on T<sub>R</sub> are responsible for the synthesizing mannopine (transcript 1' and 2'), and one gene encodes a product that converts mannopine to agropine (transcript 0') (2).

The Ri plasmids can contain one or two T-DNAs, depending on the strain, analogously to Ti plasmids. The genes encoded by the Ri plasmid T-DNA differ, however, from those of the Ti plasmid. In contrast to the T-DNA genes of the Ti plasmid, the functions of the genes encoded by the Ri plasmid T-DNA are little understood, although the *rolA*, *B*, and *C* gene products disturb the response of the infected cell to phytohormones.

Transfer of the T-DNA to the plant cell is a multistep process that is similar to DNA transfer during bacterial **conjugation** (3). The proteins that mediate the transfer process are encoded by the virulence (*vir*) regions of the Ti and Ri plasmids. The transfer process is initiated by activation of the expression of the genes of the *vir* region on the bacteria that sense a wounded plant cell. *VirD* encodes four peptides, of which virD1 binds to the border repeats. VirD2 forms a complex with virD1 and introduces a single-strand nick in the DNA. Following nicking, a single strand of DNA is released by strand-replacement DNA synthesis, and this is called the T-strand. *virE2* is a **single-stranded**, [DNA-binding protein](#) that it is thought, binds to the T-strand and protects it from **nucleases**.

The T-strand bound with *virE2* and *virD2* forms the T complex, and this is transferred to the plant cell. Transfer is mediated *via* a sexual **pili**-like organ, made up of peptides encoded by *virB*. Both *virE2* and *virD2* contain nuclear targeting signals (see [Nuclear Import, Export](#)), which, it is thought, help in targeting the T-complex to the [nucleus](#). The mechanism of insertion into the genome is not completely understood, but it is considered similar to illegitimate [recombination](#). A single

unrearranged T-DNA can insert into the genome, but multimers of T-DNA can occur at single sites within the plant genome, depending on the experimental transformation, the bacteria, and the T-DNA *vir* region combination used. Once integrated, the T-DNA is generally maintained stably and, if a single insert, is inherited in a Mendelian manner (4).

#### Bibliography

1. K. Wang, L. Herrera-Estrella, M. Van Montagu, and P. Zambryski (1984) *Cell*. **38**, 455–462.
2. P. Zambryski, J. Tempe, and J. Schell (1989) *Cell* **56**, 193–201.
3. M. Lessl and E. Lanka (1994) *Cell* **77**, 321–324.
4. C. Baron and P. Zambryski (1995) *Trends Biotechnol.* **13**, 356–361.

#### Suggestions for Further Reading

5. W. Ream (1989) *Agrobacterium tumefaciens* and interkingdom genetic exchange, *Annu. Rev. Phytopathol.* **27**, 583–618.
6. J. R. Zupan and P. Zambryski (1997) The *Agrobacterium* DNA transfer complex, *Crit. Rev. Plant Sci.* **16**, 279–295.

#### Tac Promoter

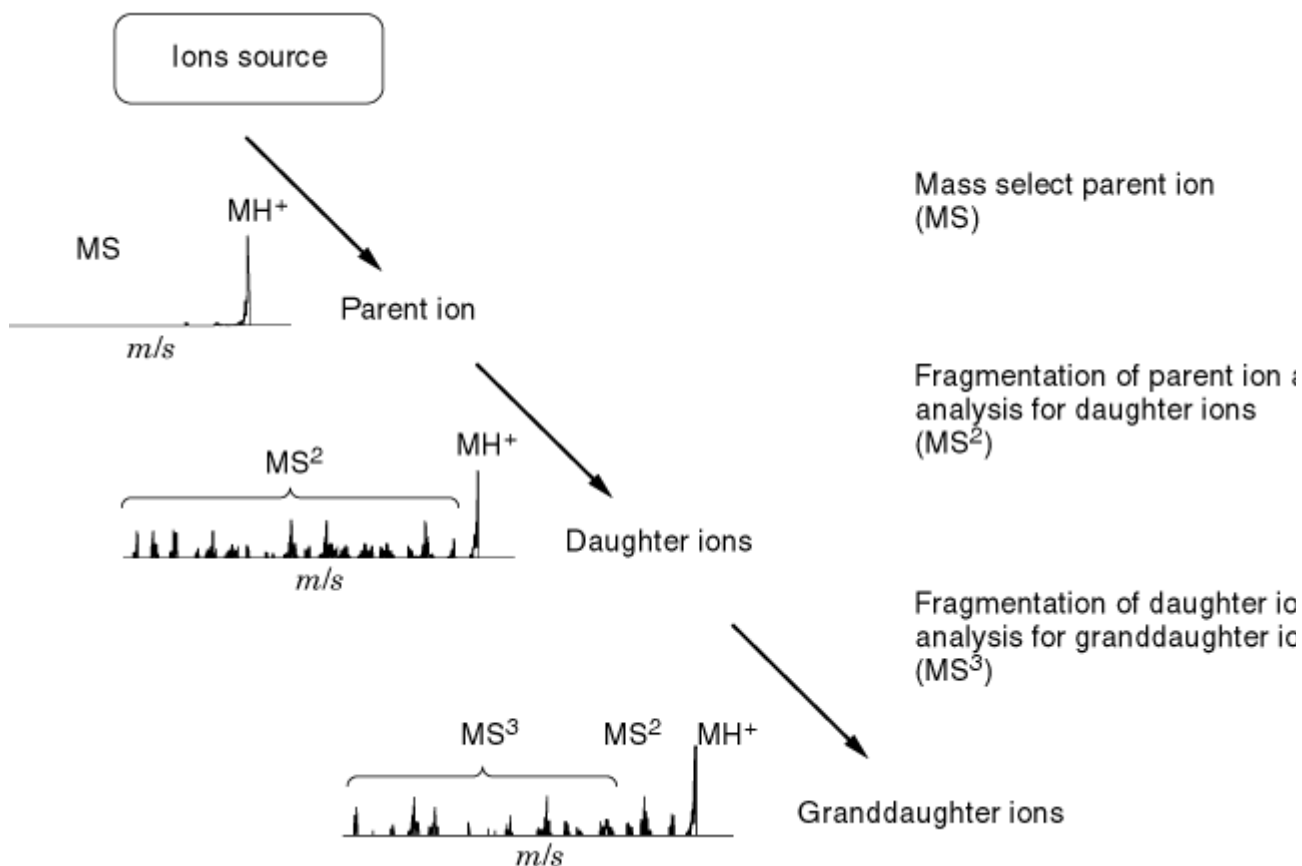
The tac promoter was artificially constructed from the *trp* and *lac* promoters. It is utilized by the prevalent form of the **RNA polymerase** from *E. coli*, the  $E\sigma^{70}$  holoenzyme. It has consensus –10 and –35 regions, each obtained from a different parent promoter. The tac promoter was popular in the early 1980s for driving the overexpression of cloned genes (see **Genetic engineering**), but now it has largely been supplanted by systems based on the use of **bacteriophage** RNA polymerases and their promoters.

#### Tandem Mass Spectrometry

An important goal of [mass spectrometry](#) (MS) is the routine acquisition of complete sequence information for **biopolymers** such as [proteins](#) and oligonucleotides. The realization of this goal has significantly progressed as a direct result of the development of [electrospray ionization](#) (ESI) with tandem mass spectrometry. By typically using an *ion trap* or *triple quadrupole* mass analyzer, ESI serves to generate the intact molecular ion in the gas phase. The ion is then exposed to ion/molecule collisions with neutral atoms such as argon or helium. The resulting *collision-induced dissociation* (CID) results in fragment ions that can be mass analyzed. This ability to induce fragmentation and then to perform successive mass spectrometry experiments on these fragment ions is known as tandem mass spectrometry (abbreviated  $MS^n$ , where  $n$  refers to the number of generations of fragment ions being analyzed) and is illustrated in Fig. 1. Because the collision-induced dissociation behavior of **peptides** is already well characterized, tandem mass spectrometry is used routinely to acquire the partial or total sequence of small peptides (<3 kDa); it has also been used for short oligonucleotides.



**Figure 1.** Generation of fragment ions via collision-induced dissociation (CID) and the mass analysis ( $MS^n$ ) of the progeny fragment ions. The terms parent, daughter, and granddaughter ions are used in this discussion; however, precursor, product, second-generation product ions are also commonly used. In the  $MS^2$  experiment the molecular ion  $MH^+$  can be selected in the mass analyzer and made to undergo CID, which results in its fragmentation and the subsequent mass analysis of the fragment ions. An  $MS^3$  experiment can be performed by selecting a daughter fragment ion, exposing it to CID, and thus generating granddaughter fragment ions.



Another biopolymer sequencing method, *ladder sequencing*, was pioneered (1) for proteins and further developed (2) for oligonucleotides. This method uses [matrix-assisted laser desorption/ionization](#) (MALDI) mass spectrometry (MS) in combination with enzyme digestion or chemical digestion to generate sequence-specific ladders of proteins and oligonucleotides. For instance, protein ladder sequencing involves the simultaneous analysis of a mixture of peptides/proteins that have undergone a stepwise [Edman Degradation](#) (see [Protein Sequencing](#)). Ladder-generating chemistry generates a family of sequence-defining peptide fragments in which each differs from the next by one amino-acid residue. Once the mixture of peptides is obtained, the mixture can be analyzed by MALDI-MS, in order to generate a sequence ladder. The analysis of oligonucleotides by MALDI has also employed ladder sequencing technology using the time-dependence of exonucleases to generate the ladder. Electrospray ionization tandem mass spectrometry and MALDI ladder sequencing have both similarities and distinguishing features, which are summarized in Table 1.

**Table 1. Characteristics of Mass Spectrometry-Based Sequencing Approaches**

| ESI Tandem Mass Spectrometry                  | MALDI-MS Ladder Sequencing                    |
|---|---|
| Picomole sensitivity                          | Picomole sensitivity                          |
| Partial (sometimes full) sequence information | Partial (sometimes full) sequence information |
| Range to about 25 residue peptides            | Range to 150 residue peptides                 |
| Range to about 10 residue oligonucleotides    | Range to about 50 residue oligonucleotides    |
| Interpretation of data is reasonably involved | Interpretation of data is straightforward     |

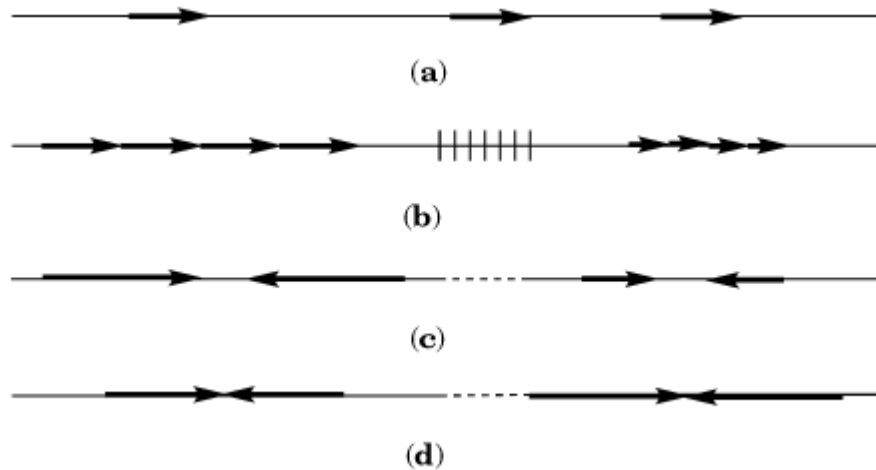
## Bibliography

1. B. T. Chait, R. Wang, R. C. Beavis, and S. B. Kent (1993) *Science* **262**, 89–92.
2. U. Pielers, W. Zurcher, M. Schar, and H. E. Moser (1993) *Nucleic Acids Research* **21**, 3191–3196.

## Tandem Repeats

A repeat, or repeating, **nucleotide sequence** is any of two or more identical segments of a longer sequence or identical segments of DNA. If the repeats are separate from one another, they are called dispersed or interspersed repeats. The repeats may also be contiguous, either head-to-head [[inverted repeats](#)] or head-to-tail, in which case they are called tandem repeats. Various types of repeats are outlined schematically in [Figure 1](#).

**Figure 1.** Traditional schematic representations of various nucleotide sequence repeats. The continuous lines in [direct repeats](#) (a) and tandem repeats (b) may correspond to either single-stranded or double-stranded molecules. In [inverted repeats](#) (c) and [palindromes](#) (d), the continuous lines mean double-stranded molecules. The arrows indicate the direction of reading the sequence in the 5' to 3' direction; left-to-right arrows are read from one strand, whereas the right-to-left arrows are read from the complementary strand. Only double-stranded complementary symmetric palindromes are shown in the scheme. No conventional presentation exists for true palindromes.



The most conspicuous family of tandem repeats, typical of eukaryotes, is the so-called simple repeats, in which the repeating unit consists of only a few bases. These may appear in hundreds of tandem copies in many separate arrays. For example,  $(TG)_n$  is the dominating simple repeat in the human and most eukaryotic **genomes**. Simple repeats are highly polymorphic, with a variable repeat copy number. Many of them are involved in biological functions in a copy number-dependent way, and the hypothesis has been put forward (1) that the tandemly repeating sequences may function to tune the nearby gene expression by changing the repeat copy number. This is consistent with the existence of so-called triplet expansion diseases associated with an excessive increase in copy numbers of certain tandemly repeating triplets, such as  $(CGG)_n$  in the case of fragile X syndrome and  $(CAG)_n$  in Huntington's disease (2). The expansion type diseases are not confined to only triplet repeats. Diseases caused by expansion of tandemly repeating units of other sizes also exist (3, 4).

Telomeric tandem repeats, eg, the human **telomeric** sequence  $(AGGGTT)_n$ , are important in maintaining intact the ends of chromosomes (5). It was already understood in the early 1970s that the absence of such regenerating repeats at the ends of the linear chromosomes would lead to incomplete replication of the chromosomes (6, 7).

Tandemly repeating sequences can amount to an appreciable proportion of the genome. In eukaryotes, most of these are **satellite** DNA sequences. Their G + C composition is frequently different from the bulk of the genomic DNA, so that the satellite DNAs are readily detected as satellite fractions by all techniques that are sensitive to the G+C content. Most satellite DNA is involved in the noncoding **heterochromatin** regions of the chromosomes. The size of the repeat unit of the satellite DNA can be a few to several hundred base pairs.

The coding sequences often form tandem (head-to-tail) arrays as well, both in eukaryotes (8) and in prokaryotes (9), as a result of **gene amplification** (see). For example, this is typical of **ribosomal** RNA genes.

#### Bibliography

1. E. N. Trifonov (1989) Bull. Math. Biol. **51**, 417–432.
2. C. T. Ashley Jr. and S. T. Warren (1995) Annu. Rev. Genet. **29**, 703–728.
3. J. C. T. van Deutekom et al. (1993) Hum. Mol. Genet. **2**, 2037–2042.
4. J.-L. Mandel (1997) Nature **386**, 767–769.
5. E. H. Blackburn (1990) J. Biol. Chem. **265**, 5919–5921.
6. A. M. Olovnikov (1971) Doklady Biochem. **201**, 394–397.

7. J. D. Watson (1972) *Nature New Biol.* **239**, 197–201.
8. J. L. Hamlin et al. (1991) *Prog. Nucl. Acid Res. Mol. Biol.* **41**, 203–239.
9. D. Romero and R. Palacios (1997) *Annu. Rev. Genet.* **31**, 91–111.

### Suggestion for Further Reading

10. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts and J. D. Watson (1994) *Molecular Biology of the Cell*, 3rd ed., Garland Publishing Inc., New York.

## Taq DNA Polymerase

The thermostable [DNA-dependent DNA polymerase I](#), Taq DNA polymerase (Taq pol), is the most extensively used DNA polymerase for amplifying genetic material by the polymerase chain reaction (PCR) (1) and for **sequencing DNA**. Taq pol is obtained from *Thermus aquaticus*, an **archaeobacterium** that thrives at elevated temperatures in deep thermal vents, and the enzyme is highly resistant to **denaturation** at elevated temperatures. Taq DNA polymerase is closely related to DNA polymerase I of *Escherichia coli* in both its primary and three-dimensional, tertiary structures. Both are single subunit enzymes with C-terminal polymerase and N-terminal 5'-3' exonuclease **domains**. However, Taq pol lacks the 3'-5' exonuclease (proofreading) function found in *E. coli* pol I and therefore cannot excise polymerization errors incurred during synthesis (2).

### 1. Properties

Taq pol has been **cloned** and overexpressed in *E. coli*. The purified full-length protein is composed of a single subunit of 94 kDa with two distinct activities: 5'-3' polymerase and 5'-3' exonuclease. Taq DNA polymerase is highly thermostable, loses only 10% of its activity after 30 min incubation at 72 °C (which is the normal growth temperature for *Thermus aquaticus*), and 50% activity after 30 minute incubation at 95 °C (M. Suzuki and L.A. Loeb, unpublished results). Like DNA pol I of *E. coli*, Taq pol is susceptible to **proteolytic** cleavage and yields a small N-terminal fragment that has 5'-3' exonuclease activity and a large C-terminal fragment that has polymerase activity. The large C-terminal fragment is analogous to the Klenow fragment of DNA pol I and is also called the KlenTaq or Stoffel fragment. It is thermostable, loses only 10% activity after 30 min incubation at 95 °C, and therefore is frequently used in polymerase chain reactions (PCR) and various sequencing protocols that involve prolonged incubations at elevated temperatures.

Taq pol conducts DNA-templated DNA synthesis with moderate accuracy, on average misincorporating only 1 in 9000 nucleotides (2). The predicted fidelity during PCR is one error per 400 bases after 25 cycles. The fidelity during PCR is enhanced twofold by using the truncated KlenTaq fragment instead of full-length Taq pol (3) and up to 10-fold by using other thermostable polymerases that have a 3-5' exonuclease (proofreading) function. Specific methods that enhance the fidelity of Taq pol include incubation at low pH (5 to 6) or with reduced concentrations of MgCl<sub>2</sub>. However, both of these methods decrease Taq pol activity and are generally incompatible with PCR, especially for amplifying long fragments.

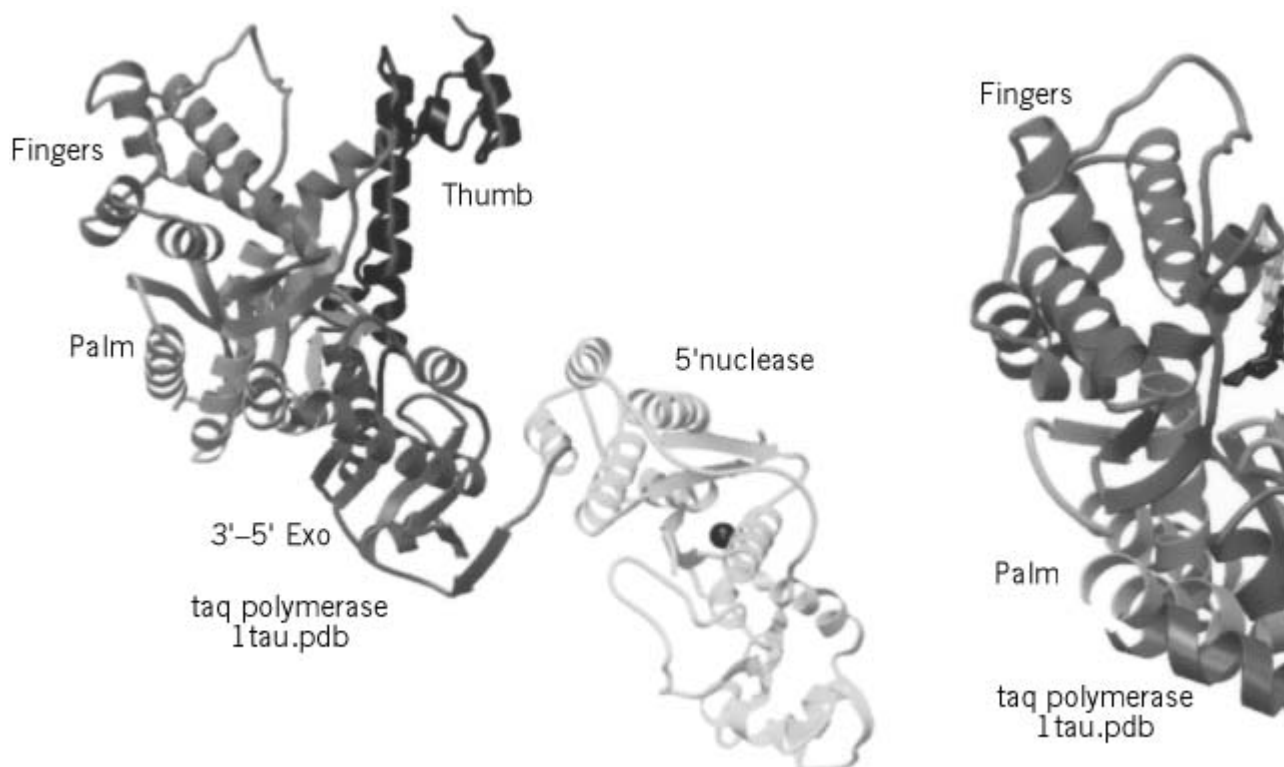
The two components of fidelity include (1) nucleotide misinsertion and (2) extension of the misinserted base. Taq pol extends mispairs at a very low efficiency: 10<sup>-3</sup> for a T–G mispair to 10<sup>-6</sup> for A–A mispair (4). Because of its relative specificity to extend only properly matched primers, Taq

pol is frequently used to detect specific *in vivo* mutations.

## 2. Structure

The **X-ray crystallographic** structures of all polymerases determined to date, including Taq DNA polymerase, resemble in overall morphology a cupped human right hand, complete with subdomains corresponding to the fingers (which bind the single stranded-template), palm (which binds the incoming deoxynucleoside triphosphate, DNTP) and thumb (which binds double-stranded DNA) (5, 6) (Fig. 1). The polymerase subdomains of Taq pol and the Klenow fragment are highly homologous and have 51% amino acid identity. In addition, the overall folding of the two polymerase domains in the three-dimensional structures are virtually identical and have an average (root mean square) deviation in alpha carbons of approximately 1.4 Å. The overall folding and location of the 3'-5' exonuclease domain, which is 30 Å away from the polymerase active site, is also similar in both enzymes. Taq pol, however, lacks four loops found in the Klenow fragment on one side of the 3'-5' exo site (5). In addition, four key amino acid residues believed to bind metal cofactors involved in the exonuclease reaction in the Klenow fragment (Asp424, Asp501, Asp355, and Glu357) are replaced by nonacidic residues in Taq pol (Leu356, Arg405, Gly308, Val310). These changes may explain why the 3'-5' nuclease domain of Taq pol is inactive. A direct comparison of the Klenow and Taq pol structures also shows that Taq contains a more **hydrophobic** core and more favorable **electrostatic interactions**. Both of these properties may contribute to the greater **thermostability** of Taq pol (6).

**Figure 1.** Ribbon diagram of the Taq DNA pol structure in the absence and presence of double-stranded DNA. This poly cupped human right hand. Residues of helix O located in the finger subdomain probably interact with the extended single I interact with the duplex portion of the DNA. The structure contains an inactive, vestigial 3'-5' nuclease proofreading dc catalytic site, and a functioning 5'-nuclease site located 70 Å from the active polymerase site.



Interestingly, the location of the 5'-3' exonuclease active site, which probably functions in **nick translation** during [DNA repair](#) and during removal of [Okazaki Fragments](#), is 70 Å from the polymerase active site. How polymerase and 5'-3' exonuclease sites so far apart work in concert to leave a “repaired” double-stranded DNA with only a nick remains a mystery. The structure of Taq pol bound to blunt-end DNA has been described (7) (Fig. 1). Similar to the complexes with DNA of pol b-DNA and HIV reverse transcriptase, the DNA in the Taq pol active site adopts a structure that is a hybrid between A and B forms. As a consequence, the minor groove, which interacts with amino acid residues in the polymerase active site, is especially wider than in B-form DNA. It is thought that protein side-chains may form [hydrogen bonds](#) with the O2 atom of the pyrimidine ring and with the N3 atoms of purines (7). Because the positions of these two atoms are unchanged in A:T and G:C base pairs, these putative hydrogen-bond interactions between protein side chains and the last base pair may insure proper base pairing before nucleotide incorporation and thus enhance the fidelity of the polymerase.

### 3. Uses

Taq pol is used extensively in polymerase chain reactions to amplify genetic material. **PCR** involves incubation with the **gene** that is to be amplified in the presence of a polymerase, PCR primers that flank the gene of interest, all four dNTP, and magnesium. Briefly, each cycle of PCR involves three steps: (1) incubation at elevated temperature to separate the double-stranded DNA (typically 95°C); (2) incubation at lower temperature to allow primer/template annealing, and (3) incubation at a temperature for polymerization (1). Before the discovery of Taq pol, PCR involved manually adding a small amount of polymerase to the incubation mixture before each polymerization step because the polymerization was inactivated during incubation at elevated temperatures. The thermostability of Taq allows heat denaturation of DNA after each cycle without enzyme inactivation, thus allowing automation of PCR. In general, current PCR technology using full-length Taq pol allows amplification from one molecule to more than  $10^5$  copies of a target sequence, which can be as large as 5 kilobases. Even larger sequences are amplified by combining a truncated KlenTaq fragment with low levels of a thermostable polymerase that contains 3'-5' proofreading activity (eg, Vent DNA polymerase or Pfu polymerase) (8).

PCR is used widely in research and clinical laboratories for diverse procedures including amplifying genes for cloning, diagnosis of diseases, and detecting levels of viruses. Taq DNA polymerase, which possesses weak RNA-templated DNA polymerase activity, is also used to synthesize double-stranded DNA from [messenger RNA](#) templates (9). However, this process, termed RT-PCR, is generally more efficient if a **reverse transcriptase** is used during the first cycle to generate single-stranded DNA and Taq pol is used in subsequent amplification cycles (10). Alternatively, a polymerase from the thermophilic bacteria *Thermus thermophilus* (TTH POL), which contains efficient RNA-templated and DNA-templated DNA polymerase activities in the presence of  $Mn^{2+}$ , can be used for RT-PCR (11). Because of its relative specificity to extend only properly matched primers, Taq pol I is used to detect specific mutations through PCR analysis, simply by choosing primers with a 3' -terminus that is complementary to the mutant. Efficient amplification of the DNA in this procedure suggests the presence of a specific mutation. This protocol has been used successfully to detect mutations in the *ras* **oncogene** and common mutations resulting in inherited disorders.

#### 3.1. Taq Pol in DNA Sequencing

Until recently, enzymatic DNA sequencing using Taq pol was marred because dideoxynucleotides (ddNTP) are not efficiently incorporated by Taq pol (its usage of ddNTP is 1000 times less efficient than that of dNTP in the presence of  $Mg^{2+}$ ; Ref. 12). However, it has been recently shown that the substitution Phe667Tyr of Taq pol increases incorporation of ddNTP relative to dNTP 250 by 8000-fold (13). This mutated Taq (*Thermo sequanase*) enables cycle sequencing, thus producing accurate and analyzable sequences using either **radioactive** or **fluorescent** sequencing technologies (14).

### 3.2. Taq Pol in T/A Cloning

Taq pol shares a characteristic common to all polymerases that lack 3'-5' exonuclease proofreading activity in that it incorporates nontemplated nucleotides (usually adenine) onto blunt-end DNA (15). This property is frequently used in molecular biology to ligate fragments of DNA in T/A cloning protocols. Briefly, PCR amplification of the sequence of interest results in a significant proportion of products containing a nontemplated adenine incorporated onto both 3' ends. Then this PCR product is incubated in the presence of a **ligase** with a linear **vector** that contains 3' thymine residues. The completed reaction results in a circular DNA that contains a cloned insert.

#### Bibliography

1. R. K. Saiki, G. H. Gelfand, S. Stoffel, S. J. Scharf, R. Higuchi, G. T. Horn, K. B. Mullis, and H. A. Erlich (1988) *Science* **239**, 487–91.
2. K. A. Eckert and T. A. Kunkel (1990) *Nucleic Acids Res.* **18**, 3739–3744.
3. W. M. Barnes (1992) *Gene* **112**, 29–35.
4. M. M. Huang, N. Arnheim, and M. F. Goodman (1992) *Nucleic Acids Res.* **20**, 4567–4573.
5. Y. Kim, S. H. Eom, J. Wang, D.-S. Lee, S. W. Suh and T. A. Steitz (1995) *Nature* **376**, 612–616.
6. S. Korolev, M. Nayal, W. M. Barnes, E. Di Cera, and G. Waksman (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9264–9268.
7. S. H. Eom, J. Wang, and T. A. Steitz (1996) *Nature* **382**, 278–281.
8. W. M. Barnes (1994) *Proc. Natl. Acad. Sci. USA* **91**, 2216–2220.
9. M. D. Jones, and N. S. Foulkes (1989) *Nucleic Acids Res.* **17**, 8387–8388.
10. E. S. Kawasaki, S. Clark, M. Y. Coyne, S. D. Smith, R. Champlin, O. N. Witte, and F. McCormick (1988) *Proc. Natl. Acad. Sci. USA* **85**, 5698–5702.
11. T. W. Myers and D. D. Gelfand (1991) *Biochemistry* **30**, 7661–7666.
12. J. W. Brandis, S. G. Edwards, and K. A. Johnson (1996) *Biochemistry* **35**, 2189–2200.
13. S. Tabor, and C. C. Richardson (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6339–6343.
14. M. A. Reeve and C. W. Fuller (1995) *Nature* **376**, 796–797.
15. P. H. Patel and B. D. Preston (1994) *Proc. Natl. Acad. Sci. USA* **91**, 549–553.

#### Suggestions for Further Reading

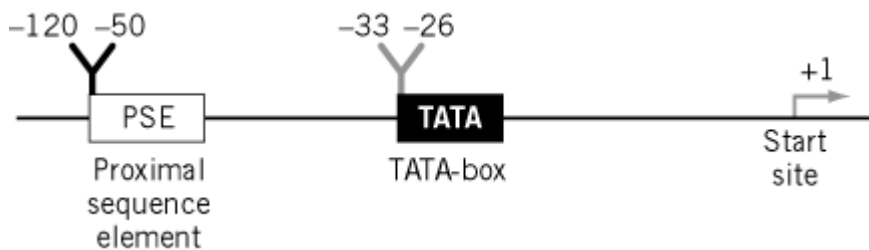
16. A. Kornberg and T. A. Baker (1992) *DNA Replication*, Freeman, New York.
17. K. A. Eckert and T. A. Kunkel (1991) DNA polymerase fidelity and the polymerase chain reaction, *PCR Methods Appl.* **1**, 17–24.

#### TATA Box

In higher eukaryotes, comparison of DNA sequences of a variety of **RNA polymerase-II promoters** revealed an A + T-rich **consensus sequence**, closely resembling the –10 **Pribnow Box** of bacterial promoters and centered approximately 30 base pairs (bp) upstream of the start site of **transcription**. This 8-bp element is called the TATA box because the maximum **homology** covers the first four bases, TATA (Fig. 1). It is also found in lower eukaryotes (such as *Saccharomyces cerevisiae*), although at variable distances from the start site, ranging from positions –30 to –120. *In vivo* as well as *in vitro* studies revealed the importance of this sequence for the initiation of accurate transcription.

Indeed, using *in vitro* transcription assays, it was shown that deletions encompassing part, or all, of the TATA sequence are detrimental for transcription (1, 2). Even a single base change within the TATA box drastically decreases transcription, indicating that the rate of transcription could be a function of the nature of the TATA-box sequence (3). The TATA sequence is involved in the selection of the transcription initiation site; accordingly, deletions of the original start site that leave the TATA box unchanged do not prevent transcription initiation, which still occurs approximately 30 bp downstream of the TATA element.

**Figure 1.** Organization of eukaryotic class II promoters.

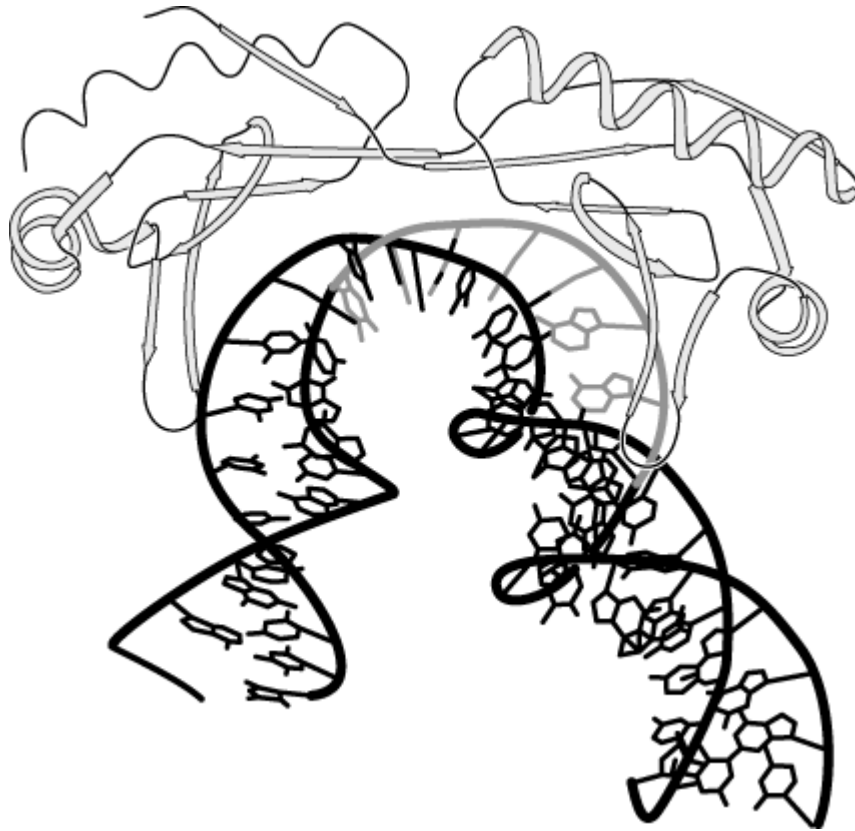


Transcription requires correct positioning of RNA polymerase II on the promoter, in addition to other components called basal (or general) [transcription factors](#). These factors participate in the formation of the preinitiation complex, which is nucleated by the binding of one of them, TFIID, on the TATA box. Targeting of the TATA box by TATA binding protein (TBP), a subunit of TFIID, allows the arrival of the other factors, either in a chronological (sequential) order or in a single step as a large complex, the RNA polymerase II holoenzyme.

[X-ray crystallography](#) studies show that the carboxy-terminal domain of TBP consists of two direct repeats of about 80 amino acid residues and resembles a saddle that sits astride the DNA (4, 5), contacting its minor groove (Fig. 2). Insertion of two phenylalanine residues of TBP between the two last base pairs of the TATA element induces a bend of about 80° in the DNA structure, which widens the minor groove but does not impede the base pairing (6). Bending the DNA around the TATA box allows contacts between regulatory factors bound to their specific sequence, on one side, and the basal transcription machinery on the other side. It could also prevent condensation of DNA in [chromatin](#) structure and repression of transcription. The importance of the DNA kinking varies, depending upon the nature of the TATA box and the surrounding sequences.

**Figure 2.** View of the TBP–TATA complex. The TATA sequence is shown in dark stippling, whereas the surrounding DNA sequences are in bold lines. TBP, which is displayed in light stippling in a ribbon representation, induces the kinking of the DNA (9).





The affinity of TBP for the TATA element probably determines the strength of the corresponding promoter—a strong promoter allowing a high rate of transcription. It is still unclear, however, whether TBP recognizes the primary TATA-box sequence or simply its three-dimensional structure. Such a structural mode of recognition has already been suggested for class I promoters, which lack a canonical TATA box. Even though DNA bending at promoters is induced by the TATA–TBP interaction, recent data demonstrate that TBP is able to bind in a lock-and-key fashion to a similar bent structure provided by DNA damage, such as cisplatin adducts or ultraviolet photoproducts (7). The requirement of TBP for transcription from a TATA-containing promoter has been challenged by the discovery that a complex of TBP-associated factors (TAFs) without TBP (called TFTC for TBP-free TAF<sub>II</sub>-containing complex) could substitute for TFIID (8). TFTC does not bind the TATA-box itself, however, raising the possibility that another component does so during TFTC-mediated transcription.

#### Bibliography

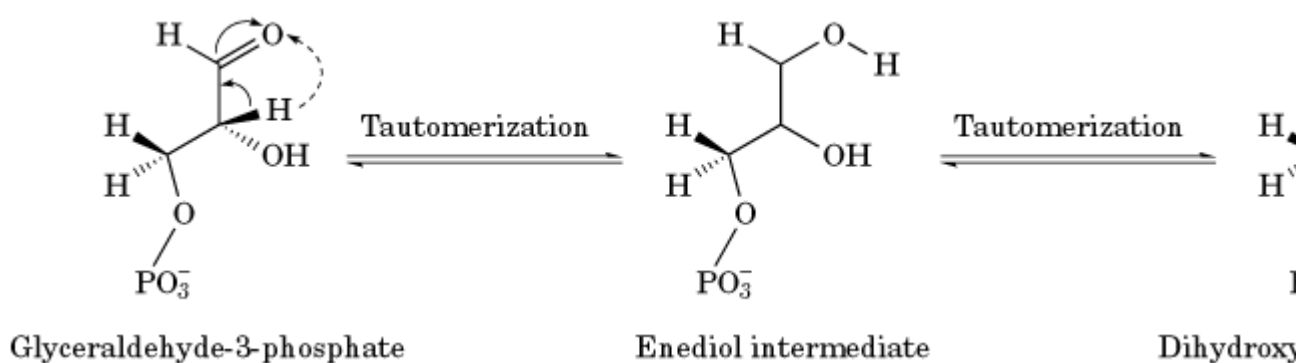
1. F. Gannon, K. O'Hare, F. Perrin, J. P. LePenec, C. Benoist, M. Cochet, R. Breathnach, A. Royal, A. Garapin, B. Cami, and P. Chambon (1979) *Nature* **278**, 428–434.
2. J. Corden, B. Wasylyk, A. Buchwalder, P. Sassone-Corsi, C. Kedinger, and P. Chambon (1980) *Science* **209**, 1406–1414.
3. S. Hahn, S. Buratowski, P. A. Sharp, and L. Guarente (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5718–5722.
4. B. Cavallini, J. Huet, J. L. Plassat, A. Sentenac, J. M. Egly, and P. Chambon (1988) *Nature* **334**, 77–80.
5. S. Buratowski, S. Hahn, P. A. Sharp, and L. Guarente (1988) *Nature* **334**, 37–42.
6. S. K. Burley (1996) *Curr. Opin. Struct. Biol.* **6**, 69–75.
7. P. Vichi, F. Coin, J. P. Renaud, W. Vermeulen, J. H. Hoeijmakers, D. Moras, and J. M. Egly (1997) *EMBO J.* **16**, 7444–7456.

8. E. Wieczorek, M. Brand, X. Jacq, and L. Tora (1998) *Nature* **393**, 187–191.  
 9. D. B. Nikolov, H. Chen, E. D. Halay, A. A. Usheva, K. Hisatake, D. K. Lee, R. G. Roeder, and S. K. Burkley (1995) *Nature* **377**, 119–128.

## Tautomers

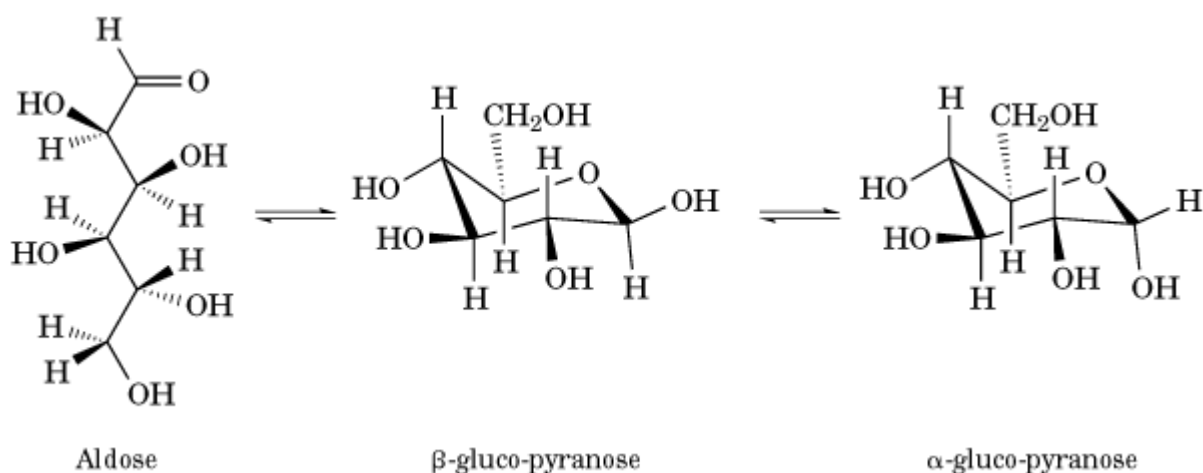
Tautomers are structural **isomers** that interconvert rapidly (1), so they will all exist in a solution of the compound at equilibrium at ambient temperature. The concept of tautomerism was introduced by Laar (2). In most cases, tautomers are generated by the structural and electronic rearrangements induced by moving a single proton (3). In the reaction catalyzed by triosephosphate isomerase, the substrate, glyceraldehyde-3-phosphate, the enediol intermediate, and the product, dihydroxyacetone-phosphate, comprise two tautomeric pairs. As shown in Figure 1, glyceraldehyde-3-phosphate and the enediol are structural isomers that result when a proton is transferred from C2 to the aldehyde oxygen, as indicated by the dashed arrow. Transfer of the proton is accompanied by a rearrangement of the electrons, as indicated by the solid curved arrows. Similarly, dihydroxyacetone-phosphate and the enediol intermediate are tautomers, as a single-proton transfer also permits their interconversion. Because of the imprecision of the phrase “rapidly,” the property of being tautomers may not be transitive. For example, glyceraldehyde-3-phosphate and dihydroxyacetone-phosphate are not considered to be tautomers, although they are both tautomers of the enediol intermediate.

**Figure 1.** Two examples of keto-enol tautomerization. A proton from the  $\alpha$ -methylene of either glyceraldehyde-3-phosphate or dihydroxyacetone-phosphate is transferred to the carbonyl; after the depicted electron rearrangement, a structural isomer



A second form of tautomerism important in biochemistry is the existence of both open-chain and ring forms of monosaccharides. Thus, the three common forms of glucose shown in Figure 2 are all tautomers, because they will interconvert rapidly in aqueous solution. The two cyclic pyranose forms are not only tautomers, but also **epimers** and **diastereomers**.

**Figure 2.** The three tautomeric forms of D-glucose. These are rapidly interconverting structural isomers. Note that the tautomeric forms differ by more than the rearrangement of a proton, consistent with Ingold's expanded definition (1).



The existence of tautomers can be detected spectroscopically, in that all the tautomeric forms may give rise to different spectral features. For glucose, it is possible to observe  $^1\text{H-NMR}$  resonances from all three tautomeric forms. It is also possible to monitor their “rapid” equilibration in aqueous solution by following the difference in optical rotation after dissolving in water a pure crystal of either of the pyranose forms.

The equilibrium constant between tautomeric forms often is very far from 1, but the less favored tautomer may be the biologically active form. In **nucleic acid** bases, keto-enol tautomerism can change the [hydrogen bond](#) acceptor carbonyl to a hydrogen bond donor hydroxyl (4). Tautomeric forms of all four nucleic acid bases exist, but the rare tautomer occurs less than 0.01% of the time (5).

#### Bibliography

1. C. K. Ingold (1969) *Structure and Mechanism of Organic Chemistry*, 2nd ed., Cornell University Press, Ithaca.
2. C. Laar (1885) *Ber. Dtsch. Chem. Ges.* **XVIII**, 652.
3. T. M. Lowry (1927) *Chem. Rev.* **4**, 233.
4. P. Beak (1977) *Acc. Chem. Res.* **10**, 186–192.
5. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer Advanced Texts in Chemistry (C. R. Cantor, ed.), Springer-Verlag, New York, Chap. 5.

#### Suggestions for Further Reading

6. J. Kyte (1995) *Structure in Protein Chemistry*, Garland, New York, pp. 54–58.
7. J. March (1985) *Advanced Organic Chemistry*, Wiley-Interscience, New York, pp. 66–70.

#### Telocentric Chromosome

Telocentric chromosomes appear by cytogenetic techniques to have their [centromeres](#) at the very end of the [chromosome](#). True telocentric chromosomes do not exist, because the **telomeres** are at the tips

of the chromosomes. In contrast to the centromeres, there is little to distinguish telomeres from [euchromatin](#), but molecular biological and genetic analysis suggest that the telomere has a unique structure and function distinct from the rest of the chromosome.

## Telomerase

There is a high level of similarity in the sequence organization of the **telomeres** in eukaryotic [chromosomes](#). Telomeres contain repeat sequences of six to eleven deoxyribonucleotides (see [Terminal Repeats](#)). All of these repeat sequences in vertebrates and in most other organisms end with guanines. Animal and **plant** telomeres contain the repeats: 5'-T<sub>2-4</sub>-A<sub>0-1</sub>-G<sub>3-4</sub>-3'. The sequence 5'-TTAGGG-3' is at the end of all vertebrate chromosomes (1). The [enzyme](#) telomerase is a terminal [transferase](#) that elongates the G-rich strand of telomeres in the absence of a DNA [template](#) (2). A common feature of telomerases is their capacity to recognize single-stranded oligonucleotides ending with G at the 3'-terminus. First defined in *Tetrahymena*, the enzyme is a [ribonucleoprotein](#) containing an **RNA** template for the **DNA** sequence in the telomere. This templating activity enables eukaryotic chromosomes to maintain their length in spite of the incapacity of the normal chromosomal **DNA polymerases** to replicate the linear chromosomal DNA completely. This occurs because DNA replicates in a 5' to 3' direction so that one strand of the **double-helix** can be synthesized to the end, and the other strand cannot because DNA synthesis has to be primed with ribonucleotides. In the absence of telomerase activity, the progressively shortened chromosome would eventually lose genes.

In human **somatic cells**, telomeres become shorter with age, whereas their lengths are maintained in **germ cells**. This has led to the suggestion of reduced telomerase activity in somatic cells, which may contribute to the process of aging. Consistent with this hypothesis, mutations in *Saccharomyces cerevisiae* that lead to the progressive loss of telomeres eventually cause [aneuploidy](#) and [senescence](#).

### Bibliography

1. E. H. Blackburn (1990) *J. Biol. Chem.* **265**, 5919–5923.
2. C. W. Grieder and E. H. Blackburn (1987) *Cell* **51**, 887–897.
3. D. Broccoli and H. Cooke (1993) *Am. J. Hum. Genet.* **52**, 657–660.

### Suggestion for Further Reading

4. D. Kipling (1995) *The Telomere*, Oxford University Press, Oxford, UK.

## Telomere

Telomeres are specialized **domains** found at the ends of eukaryotic [chromosomes](#) (1). They have specialized structures that protect the chromosomes from attack by **exonucleases**, prevent end-to-end fusion of the chromosomes, and promote complete replication of the ends of the linear **DNA** molecules present in chromosomes. Telomeric DNA has been studied in many organisms including

*Tetrahymena* and yeast. Chromosomes of *Saccharomyces cerevisiae* terminate in 250 to 650 bp of the simple repetitive DNA sequence  $C_{2-3}-A_{1-6}$ . These are binding sites for a non-histone protein known as repressor activator protein (RAP) 1. Although some lower eukaryotic telomeres are organized in [nucleosomes](#), those of a number of species, including *Tetrahymena* and yeast, display a protein-dependent protection from **nucleases** in which the size of the protected structure (<140 bp) differs from that anticipated for nucleosomes. This suggests either that non-histone nucleoprotein complexes exist in a regular array or that modified nucleosomes exist at the telomere. The RAP1 protein may be a participant in the non-histone protein-DNA complexes at the chromosomal ends of yeast. RAP1 interacts with yeast telomeres *in vivo* and is important for maintaining telomeric length. The protein is abundant (<4000 copies per cell) and fractionates with the nuclear scaffold. More recently, RAP1 has been used as a marker to localize yeast telomeres to the nuclear periphery in yeast (2). This localization depends on proteins SIR3 and SIR4, one of which, SIR3, is believed to interact with the tails of core histone H4. This observation suggests that RAP1, SIR3, SIR4, and the core histones might be involved in assembling a specialized nucleoprotein complex at the nuclear periphery in yeast. Proximity to the telomeres provides examples of the effect of chromosomal context on **gene expression** (see [Position Effect](#)).

Mammalian telomeres are composed of a tandem array of TTAGGG repeats. The length of the telomeric repeat sequences varies among **germ cells**, tumor cells, and **somatic cells**, indicating that the telomeres are a specialized variable component of the genome. Sperm telomeres are 10 to 14 kbp long, whereas telomeres in somatic cells are several kilobase pairs shorter and very heterogeneous in length. Telomere length in tumor cells is even shorter than in normal somatic cells. Specialized proteins also recognize mammalian telomeric DNA. One such protein, TTAGGG repeat factor (TRF), binds along the entire length of telomeric DNA. Short human telomeres (2 to 7 kbp long) have a very unusual [chromatin](#) structure, characterized by diffuse **micrococcal nuclease** digestion patterns. In contrast, longer telomeres (14 to 150 kbp in humans, mouse, and rat) have a more typical chromatin structure consisting of extensive arrays of close-packed nucleosomes (one every 150 to 165 bp). Human telomeres also fractionate in the nuclear scaffold fraction, consistent with the concept that telomeres assemble into an extended nucleoprotein complex with a specialized functional role (3).

#### Bibliography

1. V. A. Zakian (1989) *Ann. Rev. Genet.* **23**, 579–604.
2. F. Palladino et al. (1993) *Cell* **75**, 543–555.
3. T. DeLange (1992) *EMBO J.* **11**, 717–724.

#### Suggestion for Further Reading

4. D. Kipling (1995) *The Telomere*, Oxford University Press, Oxford, UK.

#### Temperature Factor

The atoms of a molecule are not at rest but vibrate around an equilibrium position (see [Molecular Dynamics](#)). The higher the temperature, the stronger the vibration. In [X-ray crystallography](#), the vibration affects the X-ray scattering by each atom. During an X-ray exposure, scattering is by a “time-averaged” atom, averaged over a great many positions of the atom around its equilibrium. Therefore, the X-rays see an atom that is larger than it really is. The larger the area over which the electron cloud of the atom is distributed, the weaker the atomic scattering. This is expressed in the

temperature factor. If the vibration is equally strong in all directions (spherically symmetrical), it is called isotropic vibration. Then the correction factor to be applied to the atomic scattering factor is given by

$$\text{Temp factor(iso)} = \exp \left[ -B \frac{\sin^2 \theta}{\lambda^2} \right]$$

where  $B$  is the thermal parameter of the temperature factor,  $\theta$  is the reflection angle, and  $\lambda$  the X-ray wavelength.  $B$  is related to the mean square amplitude,  $\overline{u^2}$ , of the atomic vibration by  $B = 8\pi^2 \overline{u^2}$ .

If the vibration is anisotropic, it is usually represented by an ellipsoid (replacing the sphere for isotropic vibration). Now six parameters are involved, three to define the length of the ellipsoid axes and three for their orientation. In elucidating a structure by X-ray crystallography, every atom requires determining three parameters for its position in the unit cell, plus one parameter ( $B$ ) if the vibration is isotropic and six parameters if it is anisotropic. For macromolecular crystals, the amount of X-ray data is usually not sufficient to allow determining anisotropic temperature factors. However, it can be done for relatively small molecules or a very high resolution X-ray pattern (1).

Actual temperature factors are affected by the vibrations within a molecule and also by imperfections in the lattice of the crystal, over which all X-ray diffraction measurements are averaged.

#### Bibliography

1. C. Frazao et al. (1995) *Structure* **3**, 1159–1169.

#### Suggestions for Further Reading

2. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.
3. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists*, New York, Weinheim, Cambridge, Chap. 13.

## Temperature Gradient Gel Electrophoresis

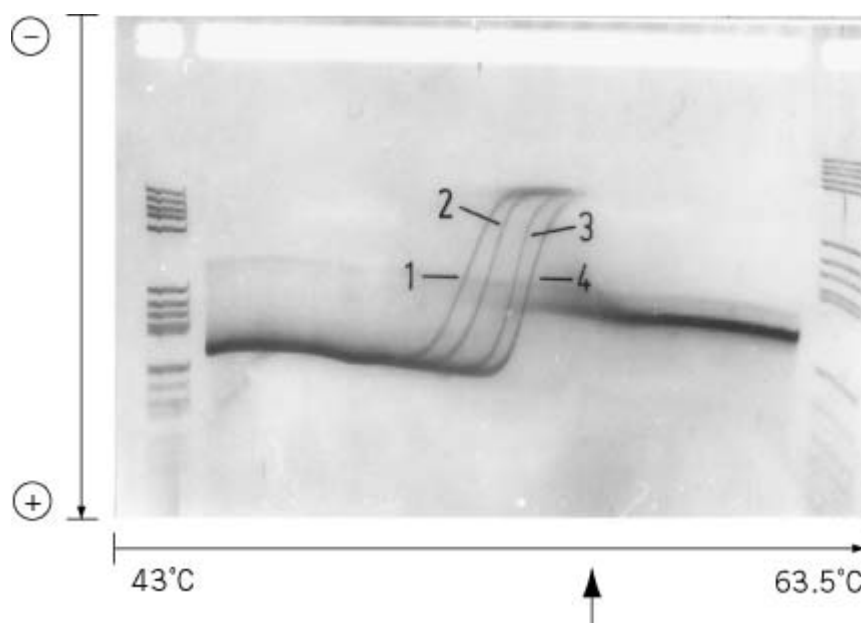
Temperature is a fundamental parameter for the structures of [proteins](#) and **nucleic acids**, which are usually unfolded by high temperatures and sometimes at low temperatures also (see [Protein Stability](#)). Varying the temperature can be combined with **PAGE** and used as the variable in [transverse gradient gel electrophoresis](#) to monitor the structures of proteins and nucleic acids and their thermal stabilities. In this technique, the sample is subjected to electrophoresis across the width of a slab [polyacrylamide](#) gel in which there is a gradient of temperature perpendicular to the direction of migration. Each molecule migrates at a constant temperature, but that temperature varies across the gel. In general, electrophoretic mobility increases with increasing temperature, but superimposed on this are any changes in mobility caused by a change in the macromolecular structure. If such changes are rapid and reversible, a continuous band is produced across the gel and through any of the temperature regions of such changes (see [Transverse Gradient Gel Electrophoresis \(Tgge\)](#)).

Proteins dissociate into subunits, if multimeric, and unfold at high temperatures, when they also tend to precipitate. Dissociation into folded subunits increases the PAGE mobility, whereas unfolding

causes a decrease. Temperature-induced unfolding is similar to that observed with [urea](#) (see [Urea Gradient Gel Electrophoresis](#)).

Temperature gradient gel electrophoresis is most useful with nucleic acids, because temperature has well-documented effects on the melting of **double helices**. The melting of an internal segment of a double-stranded DNA helix converts it into a circular [random coil](#) with a dramatic increase in the dimensions of the molecule. This greatly decreases the electrophoretic mobility through sieving polyacrylamide gels. Such melting is also rapidly reversed because the melted strands are held together by the residual double helix at either or both ends, so the two chains can rapidly reassociate. Consequently, this initial partial melting is observed as a continuous and abrupt decrease in mobility of the band (Fig. 1). At higher temperatures, the double helix is totally melted, and the two single strands dissociate. This is effectively irreversible during electrophoretic separation, so there is a discontinuity in the electrophoretic band at this temperature. The original band is replaced by a band that represents the individual single strands, which have increased PAGE mobilities because of their smaller size. But their mobility is less than that of the original double helix. If the dissociation is prevented by **cross-linking** the ends of the two chains, the original band remains continuous at high temperatures.

**Figure 1.** Temperature gradient gel electrophoresis of double-stranded DNA. The slab gel of 8% (w/v) polyacrylamide contained Tris-acetate buffer at pH 8.4, 0.2 mM EDTA, plus 8 M urea to lower the melting temperature to a practical range. As indicated, a linear temperature gradient was established across the gel during the electrophoretic separation. The sample was a mixture of two 138 base pair, DNA double helices, differing at only one base pair (C:G versus T:A), that were dissociated and then reannealed together, so that both homoduplexes (bands 3 and 4) and both heteroduplexes (bands 1 and 2) are present. The heteroduplexes have a single base mismatch and melt at temperatures lower than the two homoduplexes. From Ref. 1.



[Mutations](#), in which one base pair is replaced by another, affect the stability of the DNA double helix and alter the temperature at which the mobility of the double-stranded DNA decreases. For example, bands 3 and 4 of Figure 1 differ in only one of 128 base pairs, yet they give distinct melting curves. Detecting mutations is more sensitive if there is a mismatch between the two strands. For example, the two heteroduplexes of the two homoduplexes of bands 3 and 4 of Figure 1 are present in bands 1 and 2, which have significantly lower stabilities.

Single-stranded nucleic acids can have complementary segments within the same chain, which

associate and cause the molecule to become more compact. An extreme example is the folded structure of [transfer RNA](#). The melting of such segments is also observed by transverse gradient gel electrophoresis.

#### Bibliography

1. D. Riesner, G. Steger, U. Wiese, M. Wulfert, M. Heibey, and K. Henco (1992) *Electrophoresis* **13**, 632–636.

#### Suggestion for Further Reading

2. K. Henco, J. Harders, U. Wiese, and D. Riesner (1994) Temperature gradient gel electrophoresis for the detection of polymorphic DNA and RNA, *Methods Mol. Biol.* **31**, 211–228.

### Temperature-Sensitive Mutation

Among the [mutations](#) that affect the function of a [protein](#), some allow the protein to be active at the organism's normal temperature but inactive at either higher or lower temperatures. The former are temperature-sensitive (Ts) mutations and the latter are **cold-sensitive mutations**. Both are types of [conditional lethal mutations](#). Classically, Ts mutations have been used to identify proteins encoded by specific genes. Once a Ts mutation is located on a gene, its encoded protein is identified by showing that it is temperature-sensitive in an *in vitro* assay.

Ts mutants are also a powerful genetic way to identify interacting proteins (see [Protein–Protein Interactions](#)) because a protein destabilized at one temperature may be stabilized by mutant versions of an interacting protein. For example, if a Ts **allele** of a given gene is identified, **suppressors** can be selected at the nonpermissive temperature. Most such suppressors have no **phenotype** in the absence of the first mutant allele, but some have a cold-sensitive phenotype. Then the gene responsible for the cold-sensitive phenotype can be identified, and it is a candidate for a protein that interacts functionally with the first protein, and suppressors of the cold-sensitive phenotype can be sought. Like all such suppressor analysis, these suppressors have to be allele-specific to be useful.

#### Suggestion for Further Reading

- E. M. Phizicky and S. Fields. (1995) Protein-protein interactions: Methods for detection and analysis, *Microbiol. Rev.* **59**, 94–123.

### Template

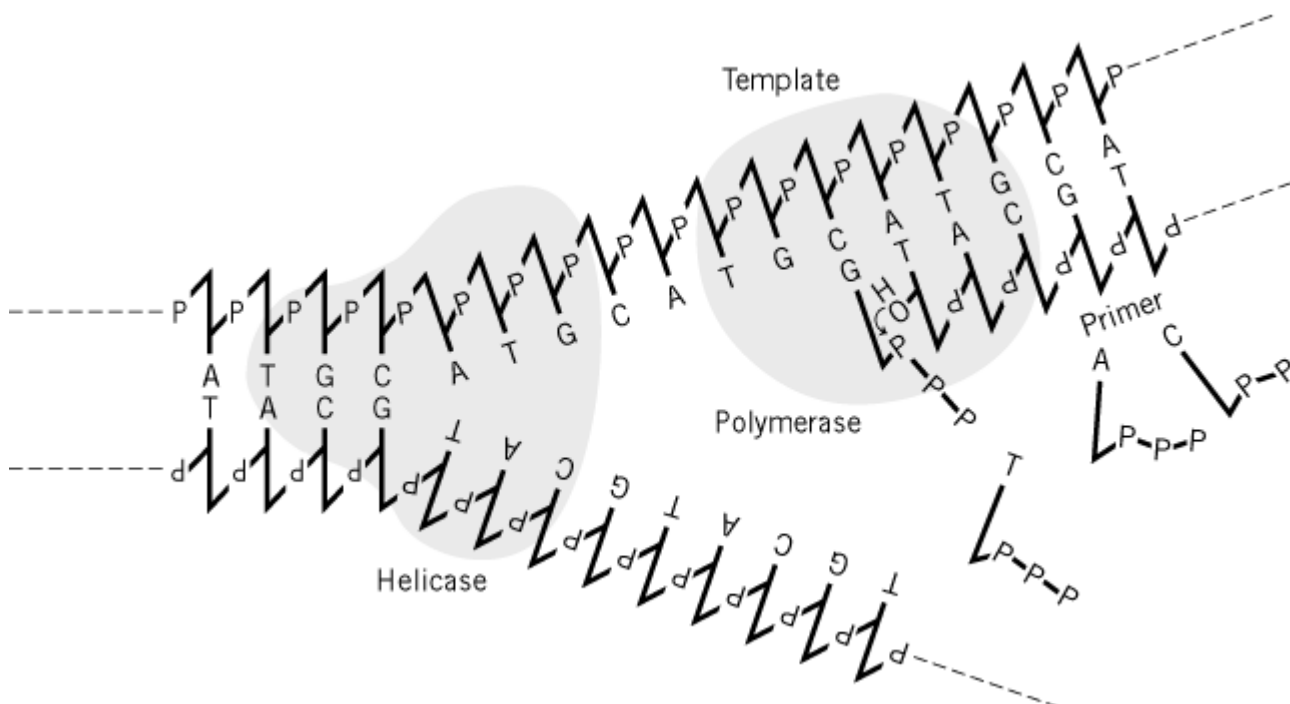
A template is defined in the 1978 *Webster's New Collegiate Dictionary* as a molecule (such as **RNA**) in a biological system that carries the genetic code for another molecule. Also, the 1995 *Concise Oxford Dictionary*, Ninth Edition describes a template as the molecular pattern governing the assembly of a [protein](#), etc. According to these broad definitions, **DNA** is the template for DNA (in [DNA replication](#)) and for RNA (in [transcription](#)), and RNA is the template for protein (in



[translation](#) ). In other words, a template is a molecular entity essential for the transfer of genetic information from DNA to DNA to RNA to protein (the central dogma). More precisely, the word “template” is used when a molecule transferring information (the template) is in direct contact with the recipient molecule (the product). Therefore, in this sense, RNA is not a template for protein because no direct interaction between [messenger RNA](#) and the protein occurs during translation (see **Protein biosynthesis**). In cases which are the exception to the central dogma, viral RNA is synthesized from RNA using RNA as a template (see RNA viruses) and, moreover, DNA is synthesized using RNA as a template during the infection of [retroviruses](#) (see also **Reverse transcription**).

In DNA replication, the double helix is unwound, and each single-stranded DNA molecule is used as a template to synthesize a complementary strand. Since DNA undergoes **semi-conservative replication**, parental DNA serves as a template and becomes a component of the daughter DNA molecule. In transcription, in contrast, a DNA strand serves as a template for the synthesis of RNA that then dissociates from the template (see [Transcription](#)). **DNA polymerases** that assemble deoxyribonucleotides on the template DNA strand according to the Adenine (A):Thymine (T), Guanine(G):Cytosine (C) base complementary rule (Watson–Crick base pairing) require both a template and primer for the reaction to take place (Fig. 1). Consequently, DNA synthesis is preceded by the synthesis of small RNA molecules of fewer than 10 bp by RNA [primase](#), which can initiate the synthesis of RNA complementary to the template DNA (see [Okazaki Fragments](#)). The primer RNA is eventually replaced by DNA through the lagging DNA strand extended from the region preceding the synthesis of RNA primer (see [Replication Fork \(Y-Fork Intermediate\)](#)). The template–product relationship through the base complementarity of A to T and G to C is not 100% accurate, but results in mis-pair formation once every  $10^3$  to  $10^4$  times. On the other hand, mis-pairs generally occur only once in  $10^8$  to  $10^9$  when genomes are replicated within the cell [(1)]. This is achieved by the proofreading activity of DNA polymerases themselves and by other [DNA repair](#) mechanisms.

**Figure 1.** DNA polymerase requires template and primer. At a growing replication fork, single strands of DNA provided by [DNA helicase](#) serve as templates for DNA polymerases. The enzymes synthesize phosphodiester bonds between the 3' end of the primer, newly synthesized strand, and the deoxyribonucleotide-triphosphate that is base-paired correctly with template strand. Only synthesis of the [leading strand](#) is shown schematically here.



The template–product relationship is very important in the transfer of genetic information in biological systems. However, the physico-chemical nature of the relationship itself is not sufficiently accurate to perform the precisely regulated biological processes, and many proteins have been evolved to repair the errors that inevitably result.

#### Bibliography

1. M. D. Topal and J. R. Fresco (1976) *Nature* **263**, 285–293

### Terminal Redundancy

Terminal redundancy is the existence of similar base-pair sequences at both ends of a linear **DNA** molecule. This is a common feature of many retroviral **genomes**. Many [retroviruses](#) and [retrotransposons](#) in the **eukaryotic** genome have identical or almost identical [long terminal repeats](#) (LTRs) at their ends. The length of the LTR is usually between 250 and 600 bp. There is less sequence conservation between the LTRs than within the body of the retrovirus or retrotransposon. The conserved sequences in the LTR required for retrotransposon function are the **promoter** sequences, sites of [transcription](#) and [poly A](#) addition, and short [inverted repeats](#) at the ends of each LTR. The activity of promoter sequences in LTRs is a potential source of genomic instability and inappropriate transcription within the [chromosome](#) (1). It has been proposed that a major function of the **methylation** of **CpG** sequences in vertebrates is to silence such transcriptional activity and to prevent transposition events. It is sometimes difficult to discriminate between *bona fide* promoters in chromosomal DNA and the evolutionary relics of LTRs (2).

#### Bibliography

1. J. A. Yoder, C. P. Walsh and T. H. Bestor (1997) *Trends Genet.* **13**, 335–340.
2. J. A. Yoder, C. P. Walsh and T. H. Bestor (1997) *Trends Genet.* **13**, 470–472.

### Suggestion for Further Reading

3. I. R. Arkhipora, N. V. Lymbomirskaya and Y. V. Ilyin (1995) *Drosophila Retrotransposons*, R. G. Landes, Austin TX.

### Terminal Repeats

**Telomeres** consist of repeated **DNA** sequences, and these terminal repeats are synthesized by [telomerase](#). A number of such distinct DNA sequences (Fig. 1) are utilized in **fungi**, protozoa, **plants** and animals. The existence of terminal repeats is important for replicating DNA at the [chromosome](#) ends and for assembling specialized [nucleoprotein](#) structures. General features of telomeric repeats include strand asymmetry in base composition in which the G-rich strand is oriented 5' to 3' toward the end of the chromosome. The G-rich strand can have a single-stranded terminus for at least some

of the [cell cycle](#), which in turn may facilitate the formation of unusual DNA structures. In humans, the size of the repeat array varies from cell to cell or from tissue to tissue, dependent on [telomerase](#) activity. Heterogeneity in the length of the terminal repeats is caused by variation in telomerase activity and by loss of terminal sequences due to incomplete replication.

**Figure 1.** Terminal repeat sequences at the telomeres of different organisms.

| <i>Sequence</i>                     | <i>Organism</i>                 |
|-------------------------------------|---------------------------------|
| TTTTGGGG                            | <i>Euplotes, Oxytricha</i>      |
| TTGGGG                              | <i>Tetrahymena</i>              |
| TG <sub>1-6</sub> TG <sub>2-3</sub> | <i>Saccharomyces cerevisiae</i> |
| TTAGGG                              | <i>Neurospora crassa</i>        |
| TTAGG                               | <i>Bombyx mori</i>              |
| TTTAGGG                             | <i>Arabidopsis thaliana</i>     |
| TTAGGG                              | <i>Homo sapiens</i>             |

#### Suggestion for Further Reading

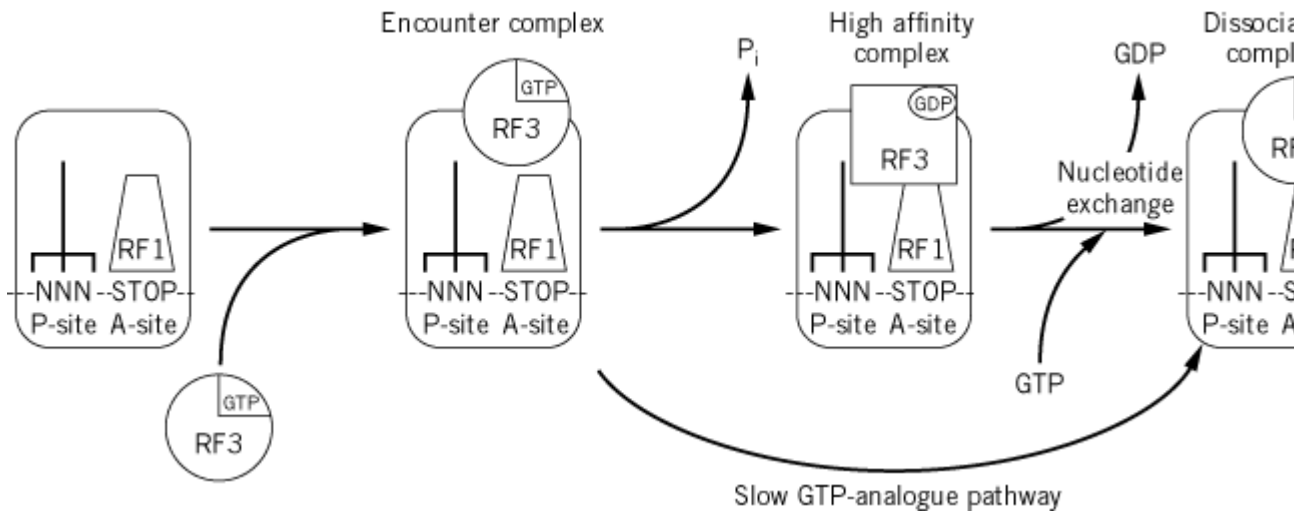
D. Kipling (1995) *The Telomere*, Oxford University Press, Oxford, UK.

### Termination Factor

Most textbooks end the description of protein biosynthesis with the release-factor-mediated release of the completed [polypeptide chain](#) from the peptidyl- transfer RNA. This is a gross oversight, because there is an additional crucial step in protein synthesis: recycling of ribosomes through dissociation of the termination complex (1). In bacteria, this process requires a ribosome-recycling factor (RRF, originally called ribosome-releasing factor). This process must be fundamental, because the gene for RRF is essential for cell growth, and the living cell must reuse the ribosome, release factors, and tRNA for the next round of protein synthesis. Upon release of the polypeptide chain, the ribosomal P and A sites remain occupied with a deacylated tRNA and a tRNA-mimicking release factor, respectively. A translocase is probably required to forward deacylated tRNA and release factor to the E and P sites of the ribosome, respectively. It is speculated that RF-3 or EF-G may catalyze this final [translocation](#) reaction. Alternatively, this post-release ribosomal complex may be dissociated directly by the action of RRF and RF-3 or EF-G (Fig. 1). Bacterial release factor RF-3 accelerates the dissociation of RF-1 and RF-2 from the ribosome in a GTP-dependent manner, and fast recycling of ribosomes requires both RF-3 and RRF (2).

**Figure 1.** An RF-3 recycle factor model (2). This explains how RF-3 accelerates the dissociation of RF-1 / 2 from the rit

tRNA.



## Bibliography

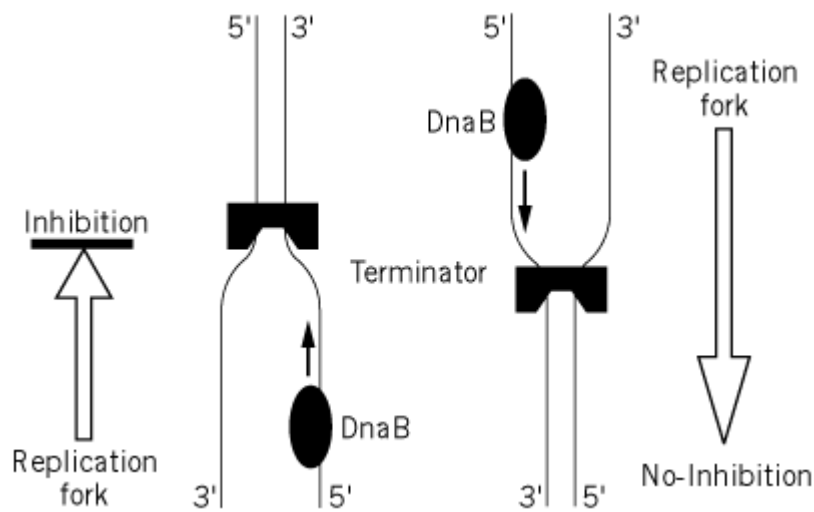
1. L. Janosi, H. Hara, S. Zhang, and A. Kaji (1996) *Adv. Biophys.* **32**, 121–201.
2. M. Y. Pavlov, D. V. Freistoffer, J. MacDougall, R. H. Buckingham, and M. Ehrenberg (1997) *EMBO J.* **16**, 4134–4141.

## Termination Of DNA Replication

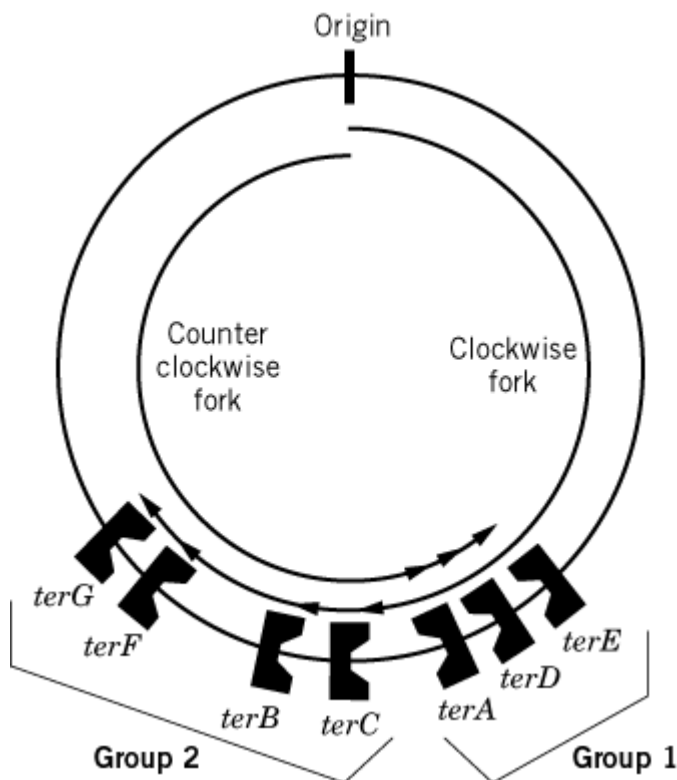
[DNA replication](#) can be divided into three distinct steps: initiation, elongation, and termination. The bidirectional replication of a [circular chromosome](#) of bacteria terminates at a position where the two [replication forks](#) meet. Bacteria have developed a system that ensures termination will occur within a restricted terminus region. This is achieved by a combination of a DNA motif of 20 to 30 bp, called the *ter* sequence, and a cognate termination protein that recognizes *ter* sites and binds to them tightly.

Such systems were discovered in *Escherichia coli* and *Bacillus subtilis*, through the identification of the accumulation of Y-shaped replication intermediates at specific sites in the terminus regions (1), and they have been extensively characterized genetically and biochemically (2-4). The *E. coli* termination protein is encoded by the *tus* (terminus utilization substance) gene and has a molecular weight of 36 kDa (309 amino acid residues). The Tus protein specifically binds to the *ter* sites containing the **consensus sequence** of about 20 bp, and the Tus–Ter complex arrests a replication fork approaching from one direction but not from the other. This arrest is thought to be due to the orientation-dependent inhibition of unwinding of the DNA duplex by the DnaB [DNA helicase](#) at the apex of the replication fork (Fig. 1). Seven *ter* sites have been identified in the terminus region of the *E. coli* [chromosome](#), as shown schematically in Figure 2. The clockwise replication fork can pass through the group1 *ter* sites, but if it arrives at the group2 *ter* sites before it meets the counterclockwise replication fork, it will stall there. Similarly, the counterclockwise replication fork will stall at the group1 *ter* sites if it has not met the clockwise fork. Thus, the termination event is regulated to occur in the terminus region opposite the *oriC* [replication origin](#).

**Figure 1.** Polar inhibition of the DNA duplex unwinding by DNA helicases.



**Figure 2.** Location and orientation of seven *ter* sites on the *E. coli* chromosome (6).



A similar system to that of *E. coli* is also present in *B. subtilis*. The *B. subtilis* termination protein, encoded by the *rtp* (replication termination protein) gene, binds to several *ter* sites on the chromosome arranged similarly to those found in *E. coli* (5). The Rtp–Ter complex can either arrest or permit the passage of the replication fork, depending on the direction of its approach. Furthermore, the *B. subtilis* Rtp–Ter system efficiently arrests the progression of the replication fork

in *E. coli*. However, the size (14.5 kDa, 122 residues) and amino acid sequence of Rtp protein are quite different from those of the *E. coli* Tus. Also, the functional *B. subtilis* arrest complex requires Rtp dimers to be bound to a longer *ter* sequence of about 30 bp in length.

The most intriguing aspect of the terminator–Ter complex is its polarity of the replication fork arrest, and two models have been proposed (6). One involves specific [protein–protein interactions](#) between the terminator and DNA helicase, and the other suggests a simple physical block of the complex (polar clump) against the action of the helicase. The findings that the Tus–Ter complex acts as a polar barrier to helicases of both prokaryotic and eukaryotic origins, and that the *B. subtilis* Rtp–Ter system efficiently arrests the progression of the replication fork in *E. coli*, support the latter model. The [X-ray crystallography](#) structures of the Tus and Rtp proteins have provided information about the molecular mechanism of the polar arrest (7, 8).

From an **evolutionary** point of view, *E. coli* and *B. subtilis* might have independently developed similar replication termination systems by [convergent evolution](#), suggesting the advantage of the system for cell growth. However, the biological significance of the polar termination system is unclear. The terminator gene can be deleted without any apparent effect on cell growth in both *B. subtilis* and *E. coli*, indicating that chromosome replication can terminate wherever the two replication forks happen to meet. A similar system of polar replication fork block is characterized in the autonomously replicating plasmid R6K (9). Furthermore, fork-blocking sites have been reported in the **yeast**, pea, frog, and human genomes (10–13).

#### Bibliography

1. A. S. Weiss and R. G. Wake (1984) *Cell* **39**, 683–689.
2. H. Yoshikawa and R. G. Wake (1993) In *Bacillus subtilis and Other Gram-Positive Bacteria* (A. L. Sonenshein, J. A. Hoch, and R. Losick, eds.), American Society for Microbiology, Washington, DC, pp. 507–528.
3. T. A. Baker (1995) *Cell* **80**, 521–524.
4. T. M. Hill (1996) In *Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 1602–1614.
5. A. A. Griffiths and R. G. Wake (1997) *J. Bacteriol.* **179**, 3358–3361.
6. R. G. Wake (1997) *Nature* **383**, 582–583.
7. D. E. Bussier, D. Bastia, and S. W. White (1995) *Cell* **80**, 651–660.
8. K. Kamada, T. Horiuchi, K. Ohsumi, N. Shimamoto, and K. Morikawa (1997) *Nature* **383**, 598–603.
9. D. Bastia, J. Germino, J. H. Crosa, and J. Ram (1981) *Proc. Natl. Acad. Sci. USA* **78**, 2095–2099.
10. B. J. Brewer and W. L. Fangmann (1988) *Cell* **55**, 637–643.
11. S. A. Greenfeder and C. S. Newlon (1992) *Mol. Cell. Biol.* **12**, 4056–4066.
12. P. Hernandez, S. S. Lamm, C. A. Bjerknes, and J. Van't Hof (1988) *EMBO J.* **7**, 303–308.
13. R. D. Little, T. H. K. Platt, and C. L. Schildkraut (1993) *Mol. Cell. Biol.* **13**, 6600–6613.

#### Suggestions for further reading

14. T. M. Hill (1996) "Features of the chromosome terminus region". In *Escherichia coli and Salmonella* (F. C. Neidhard, ed.), American Society for Microbiology, Washington, DC, pp. 1602–1614.
15. K. Kamada, T. Horiuchi, K. Ohsumi, N. Shimamoto, and K. Morikawa (1997) Structure of a replication-terminator protein complexed with DNA. *Nature* **383**, 598–603.

## Tertiary Structure

Protein structure is classified in a hierarchical manner into [primary structure](#), **secondary structure**, [tertiary structure](#), and [quaternary structure](#). The tertiary structure refers to the overall three-dimensional fold of the [polypeptide chain](#) of the protein. Some [fibrous proteins](#) (eg, [collagen](#)) adopt a regular repeating three-dimensional structure. In contrast, the tertiary structure of most proteins is much more complex and is formed by packing of the protein's secondary structural elements into one or more compact globular units (usually called **domains**). The tertiary structure provides information on the three-dimensional structure of each of the domains and on how the domains pack together. The complexity and diversity of protein tertiary structure gives rise to the complexity and diversity of protein function. Defining what a protein looks like, by determining its tertiary structure, is a significant step in understanding the biological function of that protein.

The tertiary structure of a protein is difficult to predict (see [Protein Structure Prediction](#)) but can be determined experimentally by protein [X-ray crystallography](#), **nuclear magnetic resonance (NMR)** or [cryoelectron microscopy](#). The stable folded tertiary structure of a protein conforms to certain rules of protein structure. For example, the [side chains](#) of most globular proteins are distributed in a nonrandom manner. Most **hydrophobic** residues are located in the inner core of the structure, and charged side chains are generally found on the surface. In addition, more than 90% of amino acid residues in a protein usually adopt backbone conformations that correspond to the **a-helices**, [b-sheets](#), or [turns](#) type of secondary structure. These and other rules are used to assess the quality of experimentally determined structures and can be used to both assist and validate protein structure prediction.

[See also [Protein Structure](#) and [Domain, Protein](#).]

### Suggestions for Further Reading

C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, W. H. Freeman, New York.

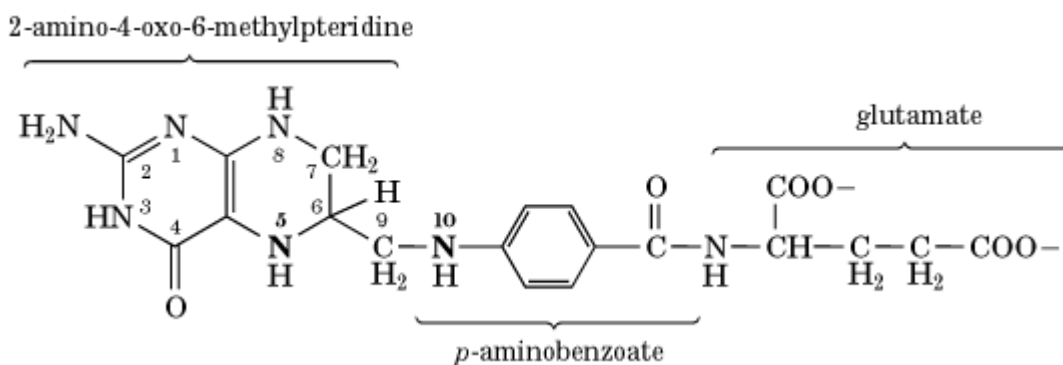
N. J. Darby and T. E. Creighton (1993) *Protein Structure*, IRL Press, Oxford, U.K.

## Tetrahydrofolate

The importance of folate compounds in metabolism has been established for over 50 years. Folate derivatives participate in a myriad of biosynthetic reactions involving transfers of groups containing a single carbon atom. For example, these functional units are essential components in the metabolism of the [amino acids](#) glycine, serine, methionine, and histidine, and the biosynthesis of purines and pyrimidines. Tetrahydrofolate (THF), the reduced form of folic acid, is the primary carrier of key components of the cellular machinery involving the mobilization and utilization of single-carbon functional groups at the methyl, methylene and formyl oxidation levels. Tetrahydrofolate derivatives serve as useful storage units of these functional groups that are not detrimental to the host organism, but are poised for reactivity when recruited by biosynthetic [enzymes](#). Bacterial cells can synthesize THF, while it can be formed in mammalian cells either by intestinal microorganisms or by the

successive reduction of dietary folate by the enzyme [dihydrofolate reductase](#) to first 7,8-dihydrofolate and finally to 5,6,7,8- THF. THF consists of three functional groups: a bicyclic substituted pteridine, *p*-aminobenzoate, and a glutamic acid tail (Fig. 1). The glutamate tail can exist either as a single residue or as multiple copies linked in amide bonds through its  $\gamma$ -carboxylate group. The most active forms of THF contain polyglutamyl tails.

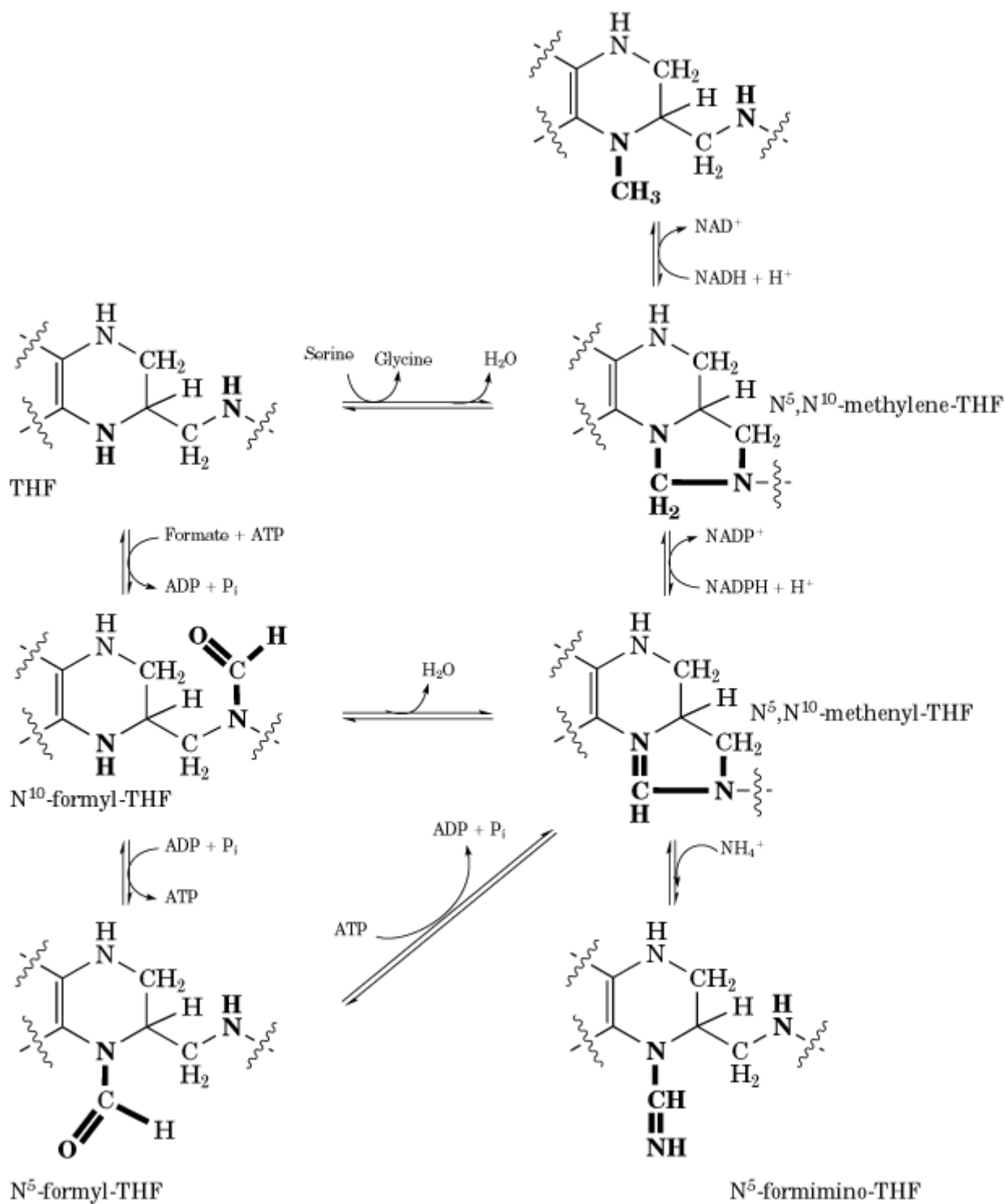
**Figure 1.** Chemical structure of THF. The N<sup>5</sup> and N<sup>10</sup>-nitrogen atoms that can carry one-carbon functional groups are in boldface.



One-carbon functional groups can be appended at position N<sup>5</sup>, N<sup>10</sup> or bridged between the N<sup>5</sup> and N<sup>10</sup> positions of THF (Fig. 2). The most reduced form of THF carries a methyl group at the N<sup>5</sup> position, the intermediate form carries a methylene group that bridges N<sup>5</sup> and N<sup>10</sup>, and the most oxidized forms carry a methenyl, formyl, or formamino group at either the N<sup>5</sup> or the N<sup>10</sup> position. Several functional groups, such as the methenyl, formyl, formimino, methylene, and methyl groups, can be interconverted while the carbon moiety is attached to THF. This feature allows a carbon group that was removed from one molecule in one oxidation state to be transferred to a second molecule in a different oxidation state. For example, N<sup>10</sup>-formyl THF (the most oxidized) can be synthesized from THF, formate and ATP by the enzyme N<sup>10</sup>-formyl THF synthetase. The newly generated N<sup>10</sup>-formyl-THF can be dehydrated to form N<sup>5</sup>, N<sup>10</sup>-methenyl-THF (most oxidized), which can be reduced to form N<sup>5</sup>, N<sup>10</sup>-methylene-THF (intermediate oxidation state), interconverted to N<sup>5</sup>-formyl THF, or deamidated to form N<sup>5</sup>-formimino-THF (most oxidized). The reduced N<sup>5</sup>, N<sup>10</sup>-methylene-THF can be further reduced by the enzyme N<sup>5</sup>, N<sup>10</sup>-methylene-THF reductase to form the THF derivative N<sup>5</sup>-methyl-THF (most reduced) which is an essential component in the conversion of homocysteine to methionine. N<sup>5</sup>, N<sup>10</sup>-methylene-THF is also generated from THF by the enzyme serine hydroxymethyl transferase, in the conversion of serine to glycine. All these THF derivatives can serve as either cofactor in a variety of biosynthetic and metabolic pathways, or as a means to neutralize and store activated one-carbon units that are byproducts of degradation reactions. Because THF and its derivatives are essential components in a number of biosynthetic reactions, folate-dependent enzymes are good targets for both antimicrobial and antineoplastic drugs.

**Figure 2.** The interconversions of one-carbon functional groups on THF. The N<sup>5</sup> and N<sup>10</sup>-nitrogens and the one-carbon functional groups are in boldface.





### Suggestions for Further Reading

R. L. Blakely and S. J. Benkovic (1984) *Folates and Pterins*, Vol. 1, Wiley, New York.

R. L. Blakely and S. J. Benkovic (1985) *Folates and Pterins*, Vol. 2, Wiley, New York.

R. L. Blakely and S. J. Benkovic (1986) *Folates and Pterins*, Vol. 3, Wiley, New York. These three reviews cover all the enzyme reactions known to involve THF derivatives.

J. E. Ayling, M. Gopal Nair, and C. M. Baugh (1993) *Chemistry and Biology of Pteridines and Folates*, Plenum Press, New York. An up-to-date compendium of the enzymes involved in folate metabolism.

## Tetranitromethane

Tetranitromethane (TNM) is commonly used for the [nitration](#) of [tyrosine](#) residues in [proteins](#). It has the structure  $C(NO_2)_4$  and a molecular weight of 196.04. It is normally a pale yellow liquid with a melting point of +13.8°C and a boiling point of 126°C. TNM is prepared by nitrating acetic anhydride with anhydrous nitric acid. TNM is also used to detect double bonds in organic compounds.

## Tetraploidy

Tetraploidy is common in **plants**. This type of **polyploidy** is characterized by four basic [genomes](#) in the cell nucleus with four copies of each of the [chromosomes](#). If the chromosomes are identical (AAAA), the organism is called *autotetraploid*. If the cell has arisen by hybridization between related species, followed by production of unreduced gametes, it is called *allotetraploid* (AABB). Because of their even number of chromosomes, both types of tetraploids undergo regular **meiosis**, and the tetraploidy is passed on to future generations.

## Thalassemia

The term thalassemia comes from the Greek *thalassa*, the word for the Mediterranean sea. An earlier designation for the disorder was Mediterranean anemia, as the earliest cases were described in populations in that area, particularly in Italy and Greece. Subsequently, similar conditions occurred extensively in Africa and in Asia. The term Cooley's anemia is also used.

Thalassemia has proved to be a problem of [hemoglobin](#) (Hb) production. All Hb's are tetrameric proteins, composed of four polypeptide chains of two types, a and b. The a-type globin chains are coded in humans by a complex of **genes** on [chromosome](#) 16, and the b-type chains by a complex of genes on chromosome 11. The molecular formula for Hb is written  $a_2b_2$ ,  $a_2d_2$ , etc., depending upon which combination of a-like and b-like polypeptides are involved.

Persons who have the classic Mediterranean version of thalassemia (thalassemia major) are severely anemic in childhood and have poor survival. Others in the same families may have mild anemia (thalassemia minor) that is not life threatening. Genetic studies indicated that the severe cases are caused by **homozygosity** for a gene that, in **heterozygous** combination with its normal **allele**, produces mild anemia. Studies of families in which hemoglobin variants were also segregating confirmed that most thalassemia mutations in the Mediterranean area and in Africa map to the b-globin gene complex on chromosome 11. Those that did not map to the b-globin complex, and this

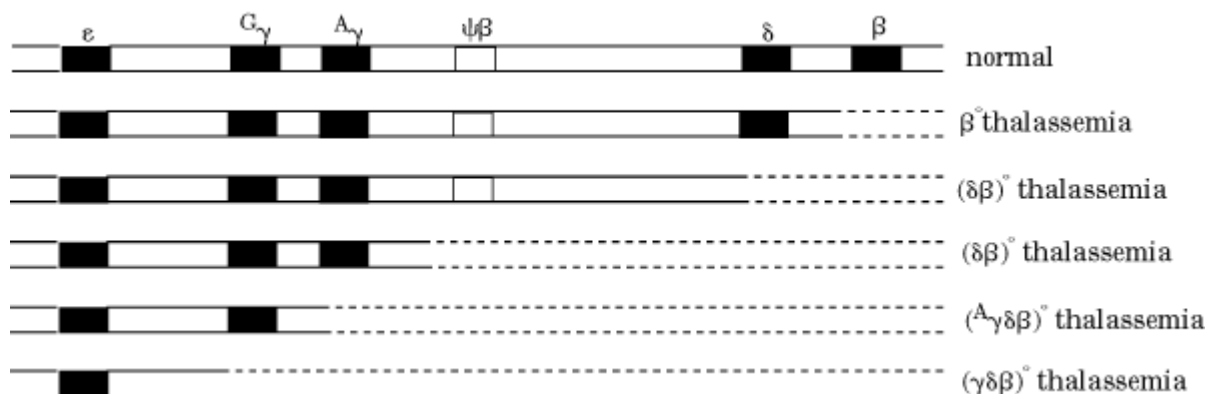
included most of the Asian examples, mapped to the  $\alpha$ -globin gene complex on chromosome 16.

Because the hemoglobins looked structurally normal in thalassemia, it was at first thought that thalassemia is due to regulatory mutations in the globin complexes. Subsequently, it was observed that there are structural modifications of globins in the case of some, but not all, thalassemia mutations. Therefore, the definition of thalassemia has been expanded to include all mutations in the globin gene complexes that significantly reduce the production of functional globin.

### 1. Beta-thalassemia

Numerous mutations in the  $\beta$ -globin complex on chromosome 11 result in thalassemia. Some are single [base-pair substitutions](#) that may be in coding regions (**exons**), in **introns**, or in regulatory regions. Many gene deletions also have been identified. Figure 1, a diagram of the  $\beta$ -globin complex, shows examples of deletions associated with thalassemia. If any part of the  $\beta$ -globin gene is deleted, no  $\beta$ -globin is produced. If the deletion affects only the  $\beta$ -globin gene, the mutation is described as a  $\beta^{\circ}$ -thalassemia mutation. Other globin genes would not be affected, but a child who is homozygous for a  $\beta^{\circ}$  mutation cannot make Hb A ( $\alpha_2\beta_2$ ), the major Hb in children and adults. The  $\gamma$  chains that are part of fetal Hb (Hb F;  $\alpha_2\gamma_2$ ) function normally in the fetus and would do so in an adult were the  $\gamma$  genes to remain active beyond the fetal period, as they do in the condition known as hereditary persistence of fetal hemoglobin. The  $\delta$  gene is active in  $\beta^{\circ}$ -thalassemia, but the amount of Hb A<sub>2</sub> ( $\alpha_2\delta_2$ ) produced is insufficient to meet the requirement for oxygen transport.

**Figure 1.** Genetic structure of the human  $\beta$ -globin gene complex on chromosome 11. The  $\psi\beta$  locus is a nonfunctional [pseudogene](#). Examples of deletion alleles associated with  $\beta$ -thalassemia are shown by the dotted line. Homozygous and heterozygous combinations of these alleles produce various forms of thalassemia, depending on the specific loci deleted. Gene positions and sizes are not drawn to scale.



Some deletions involve both the  $\beta$  and  $\delta$  genes, homozygosity for which would be designated  $(\beta\delta)^{\circ}$ -thalassemia. Clinically, such patients are not different from those with  $\beta^{\circ}$  thalassemia. Deletions that extend from the  $\beta$  gene into or beyond the  $\gamma$  genes are described as  $(\beta\gamma\delta)^{\circ}$  mutations. In this case, fetal Hb and adult Hb synthesis is impaired, and fetuses homozygous for such mutations would be aborted. If the entire  $\beta$  complex is deleted, even embryonic Hb would not be made, and abortion of the embryo would occur. Several examples of  $\beta^{\circ}$ -thalassemia are due to insertions or deletions of one or two nucleotides in the coding regions that produce [frameshift mutations](#).

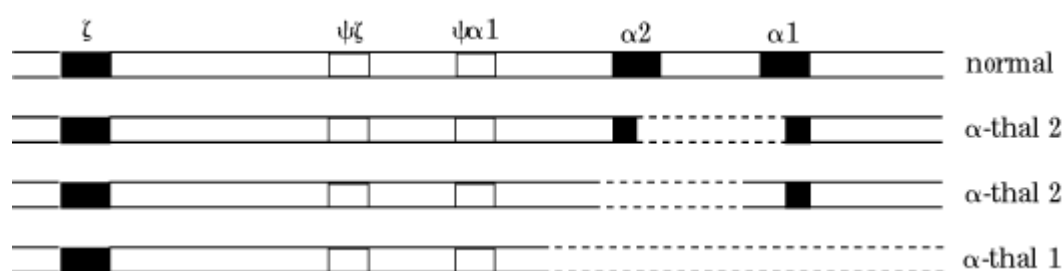
Some mutations greatly diminish [transcription](#) of the  $\beta$ -globin gene but do not abolish it completely. These are typically nucleotide substitutions, within coding regions or in introns, in untranslated regions of the transcript, or in regulatory regions. If the substitution is in the coding region, an

altered globin may be produced. In the case of Hb E, for example, a substitution in codon 26 of the b-globin gene changes the amino acid from glutamic acid to lysine and also creates a splice site that competes with the normal **splicing** position in codon 30. When normal splicing of the RNA occurs, Hb E is produced. Abnormal splicing yields no recognizable product. Mutations that lead to diminished synthesis rather than absence of b-globin are described as b<sup>+</sup>-thalassemia mutations.

## 2. Alpha-thalassemia

Much of what has been said about b-thalassemias applies to the a-globin gene complex on chromosome 16. One difference, however, is the need for a-globin chains during the fetal period. The predominant Hb in the fetus is fetal Hb (Hb F,  $\alpha_2\gamma_2$ ). Any major disruption of a-globin synthesis is detrimental to fetuses and to postnatal individuals. There is also a difference in structure between the a- and b-globin loci, as illustrated in Fig. 2. Normal chromosomes have two identical copies of the a-globin gene. Thus, a normal diploid genomic complement has four functional a-globin genes.

**Figure 2.** Genetic structure of the human a-globin gene complex on chromosome 11. There are two nonfunctional pseudogenes ( $\psi$ ). Examples of deletion alleles associated with a-thalassemia are illustrated by the dotted lines. The deletion break points in the upper a-thal 2 allele are at homologous positions in the two a genes, creating a single fused functional a gene. Gene positions and sizes are not drawn to scale.



Mutations of the a-globin genes occur widely but are particularly common in Asia. There are two common mutant alleles. One has a deletion of one locus (called the a-thalassemia 2 allele), and the other has deletion of both loci (called the a-thalassemia 1 allele). This generates diploid combinations of chromosomes with 0, 1, 2, 3, or 4 a-globin genes. Phenotypically, it makes no difference how the sum of the genes is achieved, that is, a person who has two functional genes could either be homozygous for chromosomes with one gene each (homozygous a-thalassemia 2) or heterozygous for chromosomes with 0 and 2 genes (heterozygous a-thalassemia 1). Persons who have three functional a-globin genes are phenotypically normal. If only two functional genes are present, Hb levels are slightly lower than normal, and there is mild microcytosis. If only one functional gene is present, moderate hemolytic anemia occurs, along with other hematological symptoms. Hb H ( $\beta_4$ ) precipitates in the red cells, and this form of thalassemia is also known as Hb H disease. A fetus homozygous for an a-thalassemia 1 chromosome has no functional a-globin genes. Some synthesis of z-globin chains occurs, forming Hb Portland ( $z_2g_2$ ), but most of the Hb is in the form of Hb Bart's ( $g_4$ ). Such fetuses have hydrops fetalis, severe anemia, and hepatosplenomegaly. Birth is premature, and the child is stillborn or dies shortly afterward. This condition is known almost solely in Asia, where the a-thalassemia 1 mutation is sufficiently common to produce homozygotes.

## 3. Combinations of thalassemia mutations with structural variants

One of the first indications that thalassemia is directly related to the globin genes came from observations on compound heterozygotes (double heterozygotes for a Hb variant and thalassemia).

Because  $\beta$ -thalassemia occurs in Africa, as does **sickle cell anemia** (Hb S) and Hb C, most of the early studies involved combinations of these variants. People who have a Hb S allele from one parent and a thalassemia allele from the other produce Hb S that has little or no Hb A. The red cells would sickle, although the condition is milder than typical sickle cell anemia, particularly if some Hb A is produced. Based on the proof that thalassemia mutations are allelic to  $\beta$ -globin loci in most Africans and the distinction between  $\beta^+$  and  $\beta^\circ$  mutations, the phenotypic effects became readily understood. The combination Hb-S/ $\beta^\circ$ -thal has no potential to produce normal Hb A. Combinations with  $\beta^+$  mutations produce reduced amounts of Hb A. Similar studies with the  $\beta$ -globin mutation that produces Hb C gave the same results, although the C/thal heterozygote has only mild anemia.

Combinations of structural mutations of  $\alpha$ -globin loci that have  $\alpha$ -thalassemia mutations are typically benign, unless the structural mutation is associated with depressed globin production. In that case, it would also qualify as a thalassemia mutation. An example is Hb Constant Spring, a widely occurring  $\alpha$ -globin mutation in which the termination codon is replaced with a sense codon, causing [translation](#) to continue beyond the usual stopping point. The amount of this globin is greatly reduced. When it is combined with other  $\alpha$ -thalassemia mutations, the clinical picture is much like it would be if the Hb Constant Spring gene were deleted.

The combination of  $\alpha$ -thalassemia mutations with  $\beta$ -globin mutations ameliorates the impact of the  $\beta$ -globin mutations somewhat. For example, if a person is heterozygous for sickle-cell anemia and  $\alpha$ -thalassemia 1, the ratio of Hb S to Hb A drops, apparently because the  $\beta^S$ -globin chains do not compete as effectively as  $\beta^A$ -globin chains for the limited supply of  $\alpha$ -chains.

#### 4. Population Genetics of Thalassemia

Thalassemia occurs predominantly in hot and tropical areas and in descendants of persons whose ancestral origins are from those areas. In different areas, however, quite different spectra of thalassemia mutations are prevalent. In the Mediterranean area and in Africa, most thalassemia is due to  $\beta$ -globin gene mutations. In Asia, mutations are predominantly in the  $\alpha$ -globin genes. The most common mutation in the Mediterranean is a G to A substitution at position 110 of intron 1. This creates an alternative splice site at the 3'-end of the intron. Because some of the pre-mRNA is spliced correctly, this is a  $\beta^+$ -thalassemia. In the so-called Ferrara-type thalassemia, which occurs primarily in Italy, codon 39 of the  $\beta$ -globin gene is changed to a termination codon, producing a  $\beta^\circ$ -thalassemia. In China and India, a common  $\beta^\circ$ -thalassemia mutation involves a deletion of four nucleotides in codons 41 and 42. The Hb E mutation is a common cause of  $\beta^+$ -thalassemia in Asia.

The predominant  $\alpha$ -thalassemia mutations are deletions. In Asia, either one or both  $\alpha$ -globin genes on a chromosome may be deleted, generating the large number of combinations discussed earlier. In Africa,  $\alpha$ -thalassemia also occurs, but only the chromosome from which one  $\alpha$ -globin gene is deleted, is common. Thus the Hb H disease and hydrops fetalis observed in China are rare in Africa.

The frequent occurrence of thalassemia in Africa, Southeast Asia, and the Mediterranean area, but not in more temperate regions, raised the possibility that these detrimental genes have some counterbalancing selection. This was supported by the recognition that a different spectrum of mutations is involved in thalassemia in different areas. The areas in which thalassemia genes are common correspond generally to the occurrence of falciparum malaria, and it is now thought that persons with milder versions of thalassemia have increased resistance to falciparum malaria, as is the case with sickle cell anemia. See also [Hemoglobin Mutations](#).

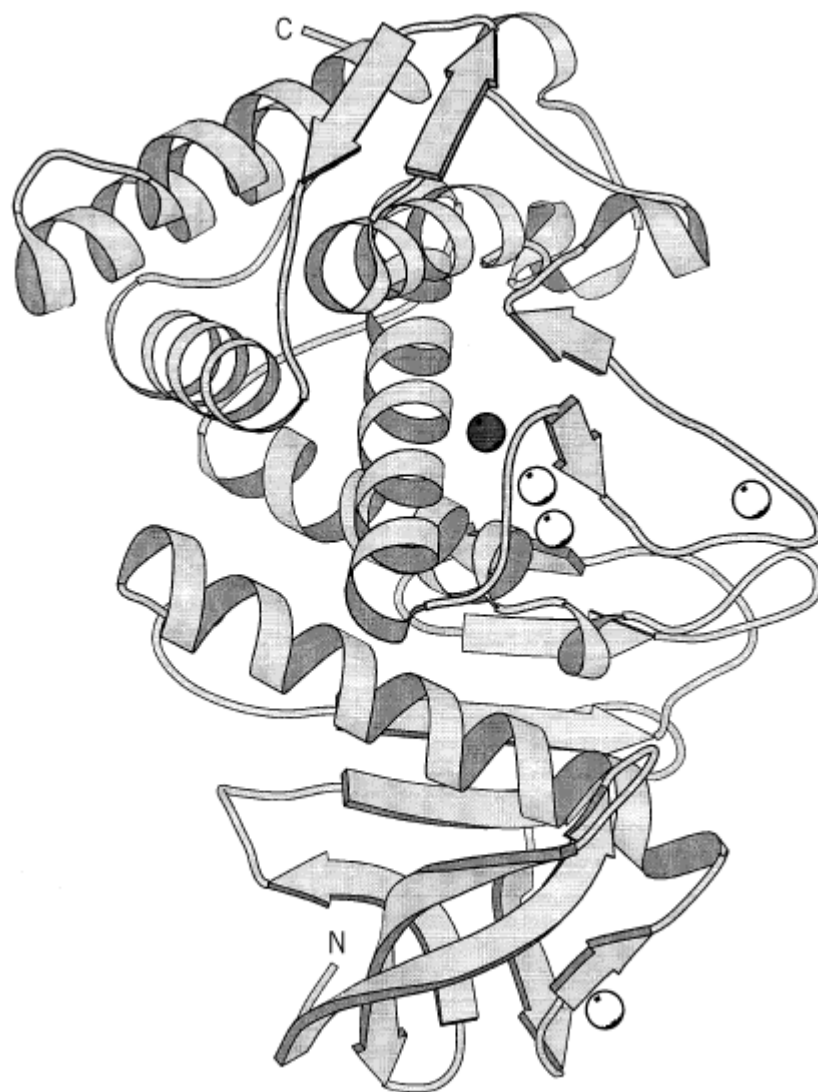
#### Suggestion for Further Reading

H. F. Bunn and B. G. Forget (1986) *Hemoglobin: Molecular, Genetic and Clinical Aspects*, W. B. Saunders, Philadelphia.

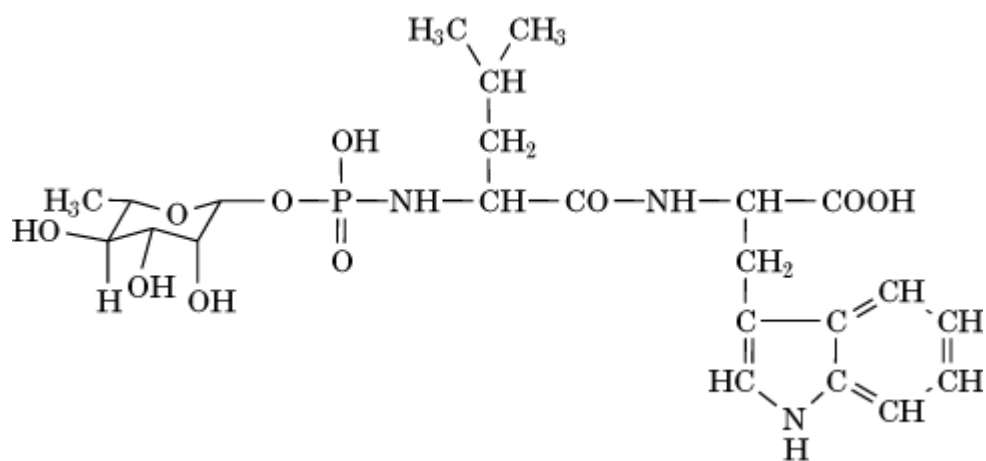
## Thermolysin

Thermolysin is a zinc **metalloenzyme** that is secreted into the culture medium by the thermostable **microorganism**, *Bacillus thermoproteolyticus* (Fig. 1). It is an **endoproteinase** (E.C. 3.4.24.27) that catalyzes the hydrolysis of **peptide bonds** in which the NH- group is contributed by a **hydrophobic** amino acid residue (1). Thermolysin, as its name implies, is a useful enzyme for degrading proteins that are resistant to other proteinases, because it is remarkably stable in solution at fairly high temperatures (ie, 80°C). It can be used, therefore, under conditions where most other proteins are unfolded (denatured by heat), thus exposing susceptible peptide bonds to the enzyme. It does not contain any stabilizing **disulfide bonds**, but it contains four calcium ions that are important for stability (Fig. 1). Removal of calcium by treatment with **EDTA** or by mild acidic conditions rapidly abolishes activity, as a result of autolytic degradation. Thermolysin is typical of a family of bacterial neutral proteinases that have similar neutral pH optima, amino acid sequences, metal contents, and mechanisms of action. A natural product synthesized by actinomyces and secreted into the culture medium is a potent inhibitor of thermolysin and related neutral proteinases. It is known as phosphoramidon, and chemically it is *N*-( $\alpha$ -L-rhamnopyranosyloxyphosphinyl)- L-leucyl-L-tryptophan (2) (Fig. 2). The sugar part of this inhibitor is not essential for good inhibition, and in fact phosphoryl-leucyl-tryptophan is even more effective than phosphoramidon ( $K_i = 2\text{nM}$  versus 28 nM).

**Figure 1.** The three-dimensional structure of thermolysin. Only the backbone is depicted schematically as a ribbon, with arrows for  $\beta$ -strands and coils for  $\alpha$ -helices. The dark sphere is the zinc ion at the active site, whereas the open spheres are stabilizing calcium ions.



**Figure 2.** Chemical structure of phosphoramidon (rhamnose-phosphoryl-leucyl-tryptophan), an inhibitor of many bacterial metalloendopeptidases, such as thermolysin. Interaction of the phosphoryl group with the active site zinc contributes to the specificity of inhibition for this class of enzymes.



## Bibliography

1. R. L. Heinrikson (1977) *Methods Enzymol.* **47**, 175–189.
2. T. Aoyagi and H. Umezawa (1975) In *Proteases and Biological Control* (E. Reich, D. B. Rifkin, and E. Shaw, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 429–454.

## Thin Filament

The thin filament is the **actin**-based protein polymer that interacts with myosin to generate muscle contraction. The thin filament contains different actin-associated proteins in striated muscle and in smooth muscle. In striated (skeletal or cardiac) muscle it is interdigitated between the myosin-based thick filaments in a highly ordered lattice. It is approximately 100 Å in diameter (see [Microfilament](#)) and can have a length of ~110 μm. The striated muscle thin filament contains the regulatory proteins tropomyosin and troponin in addition to actin. X-ray fiber diffraction has been used to generate an atomic model of the actin-tropomyosin component of the thin filament (1). Skeletal striated muscle thin filaments also contain the giant molecule nebulin, which is believed to function as a regulator of thin filament length (2). Most thin filaments have a defined stoichiometry of one tropomyosin [coiled-coil](#) dimer and one troponin complex for every seven actin subunits. The troponin complex contains three subunits: troponin C (for calcium-binding), troponin I (for inhibitory), and troponin T (for tropomyosin-binding). The release of Ca<sup>2+</sup> by the sarcoplasmic reticulum of striated muscle leads to the binding of Ca<sup>2+</sup> by troponin T, and then to an azimuthal rotation of the tropomyosin–troponin complex about the actin filament. This unblocks the myosin-binding site on the surface of the actin filament, as proposed in the “steric-blocking” model for thin filament regulation (3). Recent evidence has suggested that the initial binding of myosin to F-actin is required for an additional shift of the tropomyosin–troponin complex that leads to full activation of the thin filament (4).

Within smooth muscle, the thin filament contains the actin-binding proteins caldesmon and calponin, in addition to tropomyosin. The regulatory role of these accessory proteins in smooth muscle contraction (where the main regulation resides on the myosin filaments) is much less understood than it is for striated muscle thin filaments, and structural studies have suggested that the roles of calponin (5) and caldesmon (6) in smooth muscle cannot be the same as that played by troponin in striated muscle.

## Bibliography

1. M. Lorenz, K. J. V. Poole, D. Popp, G. Rosenbaum, and K. C. Holmes (1995) *J. Mol. Biol.* **246**, 108–119.
2. S. Labeit, T. Gibson, A. Lakey, K. Leonard, M. Zeviani, P. Knight, J. Wardale and J. Trinick (1991) *FEBS Lett.* **282**, 313–316.
3. H. E. Huxley (1972) *Cold Spring Harb. Symp. Quant. Biol.* **37**, 361–376.
4. P. Vibert, R. Craig, and W. Lehman (1997) *J. Mol. Biol.* **266**, 8–14.
5. J. L. Hodgkinson, M. el-Mezgueldi, R. Craig, P. Vibert, S. B. Marston, and W. Lehman (1997) *J. Mol. Biol.* **273**, 150–159.
6. W. Lehman, P. Vibert, and R. Craig (1997) *J. Mol. Biol.* **274**, 310–317.

## Suggestion for Further Reading

7. L. S. Tobacman (1996) Thin filament-mediated regulation of cardiac contraction. *Ann. Rev. Physiol.* **58**, 447–481.



## Thin-Layer Chromatography/Electrophoresis

### 1. Chromatography

Thin-layer chromatography (TLC), in which the adsorbent is in a thin and uniform layer fixed on a suitable supporting plate of a material, such as glass or plastic, was first suggested by Izmailov and Shraiber in 1938 (1). This type of [chromatography](#) was not widely used until the late 1950s, probably because the development of paper and **gas-liquid** chromatography was proceeding rapidly at that time. Its rapid development began about 1956, mainly because of the work of Stahl (2, 3), who devised convenient methods for preparing plates and showed that TLC is applicable to a wide variety of separations.

In thin-layer chromatography, a solution of the sample in a volatile solvent is applied to the bottom of a TLC plate. When the spot has dried, the plate is placed vertically in a suitable tank with its lower edge immersed in the selected mobile phase. The mobile phase rises by capillary action, producing an ascending chromatographic separation and resolving the various components of the sample mixture into discrete spots. At the end of the run, the mobile phase is allowed to evaporate from the plate, and the separated spots are located and identified by physical and/or chemical methods.

Thin-layer chromatography has the advantage of performing a separation easily in a minimum of time and with a minimum of chemicals and instrumentation. TLC also offers good selectivity and a wide variety of possible chromatographic interactions. TLC run in a two-dimensional mode has more separating power than the well-developed **HPLC** method. Recent developments in TLC instrumentation make the field more and more exciting and provide increased possibilities for performing this technique with high accuracy (4).

### 2. Electrophoresis

The basis of [electrophoresis](#) is the differential migration rate of ionic molecules in an electrolyte solution under the influence of an applied electric field. Although electrophoresis is not in principle a chromatographic technique, it is used in conjunction with paper chromatography, and it provides an extremely useful method for separating charged substances, such as [proteins](#) and **nucleic acids**. Many forms of electrophoresis are carried out in gel media. Electrophoresis carried out on cellulose or paper strips is called zone electrophoresis. The capillary walls provide mechanical support for the carrier electrolyte in [capillary zone electrophoresis](#). Detailed practice in electrophoresis is beyond the scope of this article, so interested readers are directed to a number of excellent monographs (5, 6).

### Bibliography

1. N. A. Izmailov and M. S. Shraiber (1938) *Farmaciya* **3**, 1; (1940) *Chem. Abstr.* **34**, 855.
2. E. Stahl et al. (1956) *Pharmazie* **11**, 633.
3. E. Stahl (ed.) (1962) *Thin Layer Chromatography*, Academic Press, New York.
4. N. Grinberg (ed.) (1990) *Modern Thin-Layer Chromatography (Chromatographic Science Series Vol. 52)*, Marcel Dekker, New York.
5. Z. Deyl (ed.) (1979) *Electrophoresis: A Survey of Techniques and Applications (Journal of Chromatography Library, Vol. 18)*, Elsevier, Amsterdam.
6. R. Weinberger (1993) *Practical Capillary Electrophoresis* Academic Press, Boston.

## Thiol Groups

Thiol groups are encountered in biological systems in [cysteine](#) residues and in cofactors such as lipoamide and lipoic acid. They are also called *sulfhydryl* and *mercapto* groups. When Zeise discovered  $C_2H_5SH$  in 1834, he called it “mercaptan” (corpus mercurium captans) because the formation of mercury derivatives was a striking characteristic. The thiol group is the most chemically reactive group that is normally encountered in biological systems. It is a powerful nucleophile that undergoes a wide variety of chemical reactions, many of which are exploited in its biological functions. A most important property is its tendency to ionize,



to the thiolate anion, which is usually the reactive species. The nonionized thiol group is usually unreactive. Most alkyl thiol groups, such as those of cysteine residues, have  $pK_a$  values close to 9. Therefore, they are reactive only at alkaline pH values, where significant amounts of thiolate anion are present. The thiolate anion is known as a soft nucleophile, poorly solvated, highly polarizable, with vacant *d*-orbitals and nucleophilic power much greater than would be predicted from its basicity. Thiol groups have only very weak **hydrogen-bonding** capabilities.

Thiol groups under too many chemical reactions to catalogue completely. Only the most important are described here.

### 1. Oxidation

Thiol groups are readily oxidized by oxygen, especially in the presence of trace amounts of metal ions, such as  $Cu^{2+}$ ,  $Fe^{2+}$ ,  $Co^{2+}$ , and  $Mn^{2+}$ ; it is likely that the complex of metal and thiol (see below) is the actual reactant with oxygen. Thiol groups may be oxidized to various oxidation states, but some of them are intrinsically unstable. In addition to the thiol form, only two oxidation states are generally encountered, the disulfide and the sulfonic acid. The disulfide is usually the end product of air oxidation:

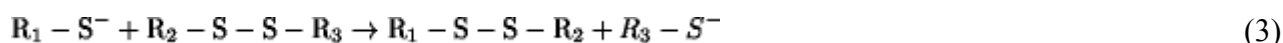


[Disulfide bonds](#) between cysteine residues are often found in proteins, especially those secreted from cells. These disulfide bonds are not produced by air oxidation, however, but are incorporated by the catalyst [protein disulfide isomerase](#).

The sulfonic acids are produced by more potent oxidizing agents. For example, performic acid oxidizes both thiol and disulfide forms of Cys residues of proteins to cysteic acid, with the  $CH_2-SO_3^-$  side chain. The intermediate sulfenate ( $-SO^-$ ) and sulfinic acid ( $-SO_2^-$ ) oxidation states are generally unstable and not normally present. They have been identified, however, when specifically stabilized in a protein structure, as in the case of the cysteine sulfenic acid that is involved in the catalytic mechanism of NADH peroxidase (1).

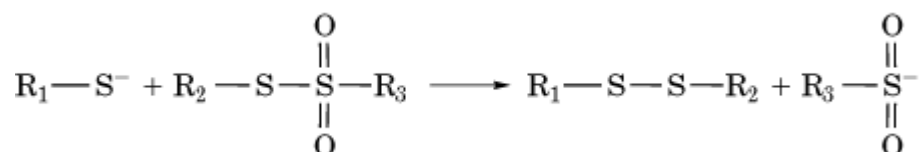
### 2. Thiol–disulfide exchange

Thiolate anions react rapidly with disulfide bonds, displacing one sulfur atom of the disulfide bond and taking its place:

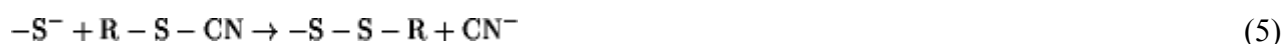


Such [thiol–disulfide exchange](#) reactions will continue among all the thiol and disulfide species present until an equilibrium mixture is generated. Thiol–disulfide exchange reactions are extremely specific, and that with 5,5'-dithio-bisnitrobenzoic acid (DTNB), or **Ellman's reagent**, is the most convenient and accurate method of assaying thiol groups quantitatively.

If one of the sulfur atoms of the disulfide is oxidized to the sulfonate, it does not take part in thiol–disulfide exchange, and only a single reaction will take place if the reagent is present in great excess over the thiol group:

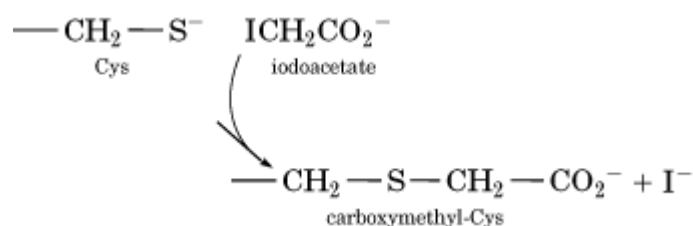


Such reactions are very useful for making specific mixed disulfide species in stoichiometric quantities. The same can be accomplished with thiocyanate derivatives in which the thiol displaces the cyanide ion:



### 3. Alkylation

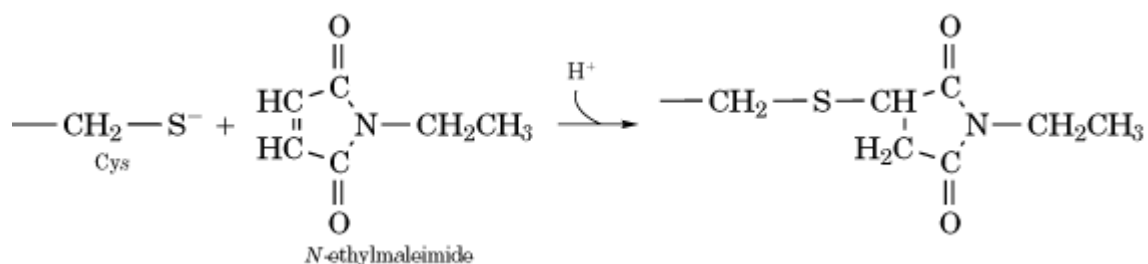
Thiolate anions react rapidly with alkyl halides, such as **iodoacetamide**, **iodoacetate**, and methyl iodide:



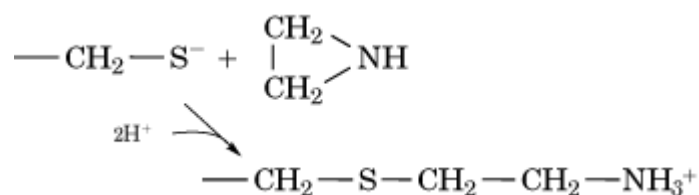
Such reactions are irreversible and the adduct generated is very stable.

### 4. Addition across double bonds

Thiolate anions are sufficiently nucleophilic to add across C=C double bonds, as in maleic anhydride. *N*-Ethylmaleimide is the classic reagent that is used most frequently to modify thiol groups:



Thiol groups can also open the ring of ethylenimine:



With cysteine residues in proteins, the resulting side-chain is now positively charged and is analogous to a [lysine](#) or [arginine](#) side-chain. Consequently, **proteolytic** enzymes such as **trypsin** will cleave the peptide bond after such a residue, so this modification is used frequently in [peptide mapping](#) studies of proteins.

## 5. Metal ions

Thiol groups form complexes of varying stabilities with a variety of metal ions. The most stable are those with divalent mercury,  $\text{Hg}^{2+}$ , but its divalency means that complexes with a variety of stoichiometries are formed. Consequently, univalent organic [mercurials](#) of the type  $\text{R-Hg}^+$  tend to be used instead, because they more reproducibly form one-to-one complexes with thiol groups. Such reactions with mercurials are the most obvious and one of the more useful ways to make heavy-atom derivatives for [X-ray crystallography](#) determination of protein structure (see [Isomorphous Replacement](#)).

Thiol complexes with silver are less stable than those of mercury, but univalent  $\text{Ag}^+$  reacts stoichiometrically and can be used to titrate thiol groups. Copper, iron, zinc, cobalt, molybdenum, manganese, and cadmium ions all form various complexes with thiol groups.

## Bibliography

1. J. I. Yeh, A. Claiborne, and W. G. J. Hol (1996) *Biochemistry* **35**, 9951–9957.

### Suggestions for Further Reading

2. P. C. Jocelyn (1972) *Biochemistry of the SH Group*, Academic Press, London.
3. Yu. M. Torchinskii (1974) *Sulphydryl and Disulphide Groups of Proteins*, Plenum Press, New York.
4. R. J. Huxtable (1986) *Biochemistry of Sulfur*, Plenum Press, New York.
5. G. L. Kenyon and T. W. Bruice (1977) Novel sulphydryl reagents. *Methods Enzymol.* **47**, 407–430.
6. A. Russo and E. A. Bump (1988) Detection and quantitation of biological sulphydryls. *Methods Biochem. Anal.* **33**, 165–241.
7. S. Patai, ed. (1990) *The Chemistry of Sulphenic Acids and Their Derivatives*, Wiley, New York.

## Thiol Proteinase

This class of [proteinase](#) (E.C. 3.4.22) catalyzes [peptide bond](#) hydrolysis by a mechanism that involves the [thiol group](#) of the side-chain of a [cysteine](#) residue. These enzymes are also known as *cysteine* or *sulphydryl proteinases*. The thiol group interacts with the carbonyl group of the

substrate's peptide bond that is to be hydrolyzed and forms a thiol ester intermediate. A histidine side chain enhances the reactivity of the thiol group in much the same way that a serine hydroxyl group is activated in [serine proteinases](#). Thiol proteinases tend to act at moderately acidic pH, about 5, and are the predominant hydrolases found in **lysosomes** and [endosomes](#). A well-known thiol proteinase derived from the papaya plant is [papain](#), the key ingredient in a commonly used meat tenderizer.

Any reagent that will react chemically with a sulfhydryl group will inactivate a thiol proteinase. **Iodoacetate** is such a nonspecific inactivator, whereas E64 [L-transepoxy succinyl-leucyl-amido (4-guanidino) butane] is a more useful general inhibitor of thiol proteinases. Protein inhibitors present in blood plasma, and known as [cystatins](#), block activity by covering the [active site](#) of the enzyme. [Leupeptin](#), a peptide aldehyde, is another general thiol proteinase inhibitor that binds tightly to the catalytic center.

A thiol proteinase appears to be responsible for the activation of the precursor peptide that gives rise to the [enkephalins](#), and E64 has been proposed as an inhibitor of prohormone processing (1). A thiol proteinase has been shown to be involved in the conversion of pre-renin to *renin*, which is the enzyme catalyzing the rate-determining step in the production of the hypertensive peptide, *angiotensin II* (2). *Cruzipain*, a thiol proteinase, has a critical function in the infectious disease trypanosomiasis, which affects millions of Central and South Americans and, hence, is a target for selective inhibition (3).

An important subclass of the cysteine proteinases includes the **interleukin-1 $\beta$**  converting enzyme, which converts an inactive precursor to a proinflammatory **cytokine** (4). These enzymes specifically cleave peptide bonds after [aspartic acid](#) residues and, hence, are called *caspases*. One member of this class, caspase-3, cleaves poly (ADP-ribose) polymerase and hence plays a key role in programmed cell death (see [Apoptosis](#)) (5), a process of fundamental importance to biology.

#### Bibliography

1. V. Y. Hook, A. V. Azaryan, S. R. Hwang, and N. Tezapsidis (1994) *FASEB J.* **8**, 1269–1278.
2. W. A. Hsueh and J. D. Baxter (1991) *Hypertension* **17**, 469–477.
3. M. E. McGrath et al. (1995) *J. Mol. Biol.* **247**, 251–259.
4. Y. Gu et al. (1997) *Science* **275**, 206–209.
5. M. Tewari et al. (1995) *Cell* **81**, 801–809.

#### Suggestion for Further Reading

6. N. A. Thornberry and S. M. Molineaux (1995) Interleukin-1 beta converting enzyme: a novel cysteine protease required for IL-1 beta production and implicated in programmed cell death. *Protein Sci.* **4**, 3–12.

#### Thiol–Disulfide Exchange

[Thiol Groups](#) and [disulfide bonds](#) undergo a spontaneous chemical reaction, in which the thiol group displaces one sulfur atom of the disulfide bond in an  $S_N^2$  type of reaction:

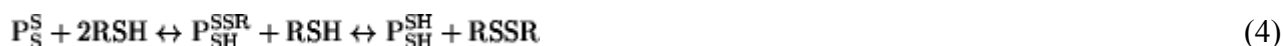


group, to give the largest rate, is the same as the pH of the reaction. A thiol with a greater  $pK_a$  value will have less ionized form; with a lower  $pK_a$  value, the thiol is ionized, but its nucleophilicity is lower.

The above considerations also apply equally to the reverse of the reaction, and the expected equilibrium constant can be estimated from Eq. 4 for the reaction in both directions. The equilibrium constant is pH-dependent if the attacking and leaving sulfur atoms have different  $pK_a$  values; it varies over the pH interval between the two thiol  $pK_a$  values. The equilibrium favors the thiol group with the lower  $pK_a$  value; just the opposite would have been predicted considering only the effect of ionization of the two thiol groups on mass action.

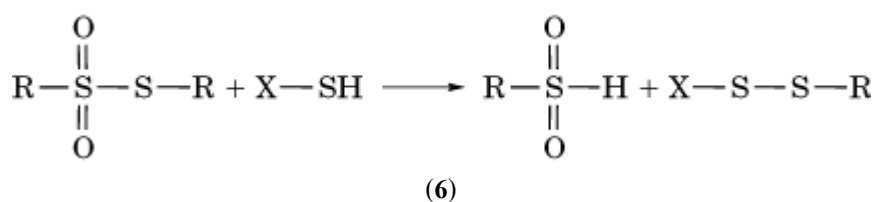
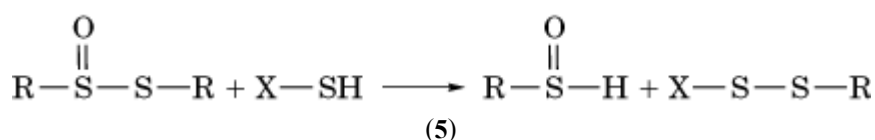
Of course, the reaction occurs more rapidly than predicted if there are positive charges near the disulfide bond, to attract the attacking thiolate anion, or if the disulfide bond is strained; for example, the strained disulfide bond of lipoic acid reacts considerably more rapidly than otherwise expected (1), which is accounted for about 3.8 kcal/mol of conformational strain in the 5-membered ring (the CSSC dihedral angle is only about  $30^\circ$  and not the favored  $90^\circ$ ). Conversely, the reaction is inhibited by the presence of negative charges near the disulfide bond, if either the thiol or disulfide group is buried and inaccessible, or if there are bulky substituents adjacent to the sulfur atoms. Electrostatic effects on the equilibria of thiol–disulfide exchange reactions are small in magnitude, but occur in the expected direction. Formation of a mixed disulfide with unlike charges close on the two moieties is favored, while that with like charges is disfavored. Such effects are substantial only when the charged groups are on adjacent residues, and their magnitude can be decreased in magnitude by electrostatic screening with high salt concentrations.

The thiol–disulfide exchange reaction is routinely used in molecular biology when protein disulfide bonds ( $P^S_S$ ) are reduced by reagents such as dithiothreitol or **b-mercaptoethanol**. In this case, two sequential thiol–disulfide exchange reactions are necessary, proceeding through a mixed disulfide between the reagent (RSH) and the protein:



Of course, thiol–disulfide exchange reactions are readily reversible, and the above reactions can be used to add disulfide bonds to a protein (see [Protein Folding In Vitro](#)).

The thiol–disulfide exchange reaction can be simplified if a mono- or dioxide form of the disulfide is used. The oxidation increases the reactivity, but only the nonoxidized sulfur atom undergoes thiol–disulfide exchange:



Consequently, the reaction stops at this stage and will be stoichiometric.

Some proteins, such as [thioredoxin](#), [glutaredoxin](#), and **protein disulfide isomerase (PDI)**,

participate in thiol–disulfide exchange reactions much more rapidly than expected from the  $pK_a$  values of their [active-site](#) thiol groups (Eq.(3)). This is believed to be due to these proteins tending to bind noncovalently the molecules that they react with, as an [enzyme](#) binds its substrate, and to them actually stabilizing the [transition state](#) for the reaction, again like an enzyme. When present in small quantities, such proteins appear to catalyze thiol–disulfide exchange reactions between other thiol and disulfide compounds. They do this by reacting more rapidly with both the thiol and disulfide compounds, so they do not catalyze the direct reaction between them but instead provide a more rapid alternative reaction pathway (3).

### Bibliography

1. T. E. Creighton (1975) *J. Mol. Biol.* **96**, 767–776.
2. R. P. Szajewski and G. M. Whitesides (1980). *J. Amer. Chem. Soc.* **102**, 2011–2026.
3. N. J. Darby, S. Raina, and T. E. Creighton (1998) *Biochemistry* **37**, 783–791.

### Suggestion for Further Reading

4. R. Singh and G. M. Whitesides (1993). "Thiol-disulfide interchange". In *The Chemistry of Sulphur-containing Functional Groups* (S. Patai and Z. Rappoport, eds.), Wiley, New York, pp. 633–658.

## Thioredoxin

Thioredoxin (Trx) is a 12-kDa **redox protein** that has a redox-active **disulfide bond** in its [active site](#) with the conserved sequence -Cys-Gly-Pro-Cys- (1, 2). The reduced form of Trx, Trx-(SH)<sub>2</sub>, is a powerful protein disulfide reductase, and the disulfide in oxidized thioredoxin (Trx-S<sub>2</sub>) is reduced by NADPH and the flavoenzyme [thioredoxin reductase](#) (3, 4). Thioredoxin, NADPH, and thioredoxin reductase, the thioredoxin system, exists in all living cells and is a hydrogen donor for synthesizing deoxyribonucleotides by **ribonucleotide reductase**, which is essential for [DNA replication](#). The thioredoxin system plays a major role in maintaining a reducing environment in the cytosol. Thioredoxin catalyzes dithiol-disulfide oxidoreductions and has a large number of functions (reviewed in Ref. 5), for example, as a hydrogen donor for reductive [enzymes](#), in phage T7 DNA replication, in filamentous phage assembly, in chloroplast **photosynthetic** enzyme regulation by light, in redox regulation of enzymes and [transcription factors](#) by thiol/disulfide redox control, in defense against oxidative stress, and as a secreted **cocytokine** in mammalian cells.

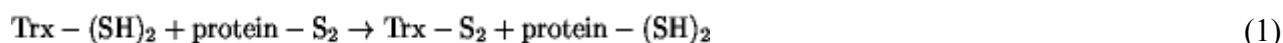
### 1. Structure

The three-dimensional **protein structures** of both oxidized and reduced thioredoxins have been determined in solution by both [X-ray crystallography](#) and **NMR**. They have the same globular “a/b sandwich” structure that is now called the thioredoxin fold. The active-site cysteine residues are located on a protrusion at the C-terminus of a [beta-strand](#) and at the beginning of an **alpha-helix**. The sulfur atom of the N-terminal Cys residue is solvent-exposed and its **thiol group** is the nucleophile that has a low  $pK_a$  value in Trx-(SH)<sub>2</sub>. A unique conserved *cis* proline residue (Pro76) is adjacent to the active site, which is surrounded by a flat **hydrophobic** surface used in binding interactions with other proteins. Thioredoxin-fold proteins include a growing [superfamily](#). Some have several thioredoxin **domains**, as in [protein disulfide isomerase](#) (PDI), which catalyzes formation of native disulfides in secreted proteins.



## 2. Thioredoxin Functions

Thioredoxin was originally defined as a hydrogen donor for *Escherichia coli* ribonucleotide reductase (1). This enzyme contains a disulfide bond after each catalytic cycle that is reduced by thioredoxin via its general disulfide reductase activity:



The enzymes that reduce sulfate (PAPS reductase) and methionine sulfoxide also use thioredoxin as an electron donor (5). Trx-(SH)<sub>2</sub> in *E. coli* is essential for phage T7 DNA replication in its role as a subunit of T7 **DNA polymerase** (6). This enzyme is a 1:1 complex of the virally coded polymerase (gene 5 protein of 80 kDa) and the host Trx-(SH)<sub>2</sub> (7). The role of the bound Trx is to give the enzyme high processivity in DNA synthesis by forming a complex and does not involve electron transport. **Recombinant** exonuclease-deficient T7 DNA polymerase is widely used in [DNA sequencing](#). *E. coli* cells that lack Trx (*trxA*<sup>-</sup>) are viable because [glutaredoxin](#) and [glutathione](#) act as a hydrogen donor for ribonucleotide reductase (8). Thioredoxin, an abundant protein in *E. coli*, has about 10,000 copies per cell and increases in amount in stationary-phase cells. It is located at the plasma membrane at the junctions of the inner and outer membrane (Bayer's adhesion sites) (9). This may explain its role in assembly of filamentous phage, where Trx-(SH)<sub>2</sub> is essential (10), but Trx-S<sub>2</sub> is not. This also gives the phage assembly process a requirement for thioredoxin reductase. The localization of Trx at the plasma membrane in *E. coli* and its release from cells by osmotic shock is a basis for a Trx-fusion system for overexpressing proteins. Soluble fusion proteins are selectively released from cells by osmotic shock (11). Recently, the **gene** for a second *E. coli* thioredoxin, called Trx2 (*trxC*), has been **cloned** (12), and those **genomes** sequenced generally encode two or more thioredoxin genes.

The tryptophan **fluorescence** of *E. coli* Trx is strongly quenched by the disulfide bond in Trx-S<sub>2</sub> (13). Reduction to Trx-(SH)<sub>2</sub> results in three-fold increased fluorescence caused by a localized conformational change that affects Trp28, whereas the conserved Trp31 still has a low quantum yield (13). The kinetics of oxidation or reduction of thioredoxin are readily measured by fluorescence. Such experiments demonstrated that Trx-S<sub>2</sub> is reduced by **dithiothreitol** (DTT) two orders of magnitude faster than the disulfide bonds of insulin and that Trx-(SH)<sub>2</sub> reacts four orders of magnitude more rapidly with the insulin disulfide bonds than DTT (4). This led to the realization that thioredoxin catalyzes the reduction of insulin disulfides by DTT. This can be used as a simple Trx assay by following the turbidity that results from insulin B-chain precipitation (4). The [oxidation/reduction potential](#) ( $E_o'$ ) of Trx-S<sub>2</sub>/Trx-(SH)<sub>2</sub> is -270mV, and a Pro34His mutant has an  $E_o'$  of -235mV (14).

Unique chloroplast thioredoxins in photosynthetic organisms regulate several photosynthetic enzymes by light *via* electrons from **ferredoxin**-thioredoxin reductase (15). Trx<sub>f</sub> regulates fructose-bis-phosphatase and Trx<sub>m</sub> regulates malate dehydrogenase (16). In each case the enzyme inactive in the dark contains a disulfide bond that is reduced to activate the enzyme in light. Trx<sub>h</sub> in the cytosol is reduced by NAPH and thioredoxin reductase.

Thioredoxin and thioredoxin reductase from mammals differ in some important respects from the prokaryotic proteins, and they were first purified to homogeneity from rat liver (17) (Table 1). Thioredoxin has roles in regulating many transcription factors by thiol/disulfide redox control that

keep critical thiol groups reduced (18). Trx is also a secreted cocytokine from both virally infected and normal cells (19). Secretion occurs by a mechanism that does not require a [signal peptide](#) (20). Measurements of human plasma levels of thioredoxin show a correlation with AIDS progression (21). Thioredoxin protects against oxidative stress as an electron donor to thioredoxin peroxidases (peroxiredoxins) (22).

**Table 1. Properties of Thioredoxin Systems**

| Organism                       | Thioredoxin   | Thioredoxin Reductase   |
|--------------------------------|---|---|
| <i>E. coli</i> , yeast, plants | 12 kDa  | M <sub>r</sub> 70,000 (2 subunits)<br><i>Highly specific</i>  |
| Mammalian cells                | 12 kDa<br><i>Inactivated by oxidation of structural SH-groups</i> | M <sub>r</sub> 116,000 (twosubunits)<br><i>Broad substrate specificity, selenoprotein homologous to glutathione reductase</i> |

### 3. Sequence Comparisons

Thioredoxins from archaebacteria to humans have 27 to 69% sequence identity to the well-studied *E. coli* thioredoxin, which indicates that all thioredoxins have similar three-dimensional structures (23). The amino acid chain length in all organisms studied thus far is approximately the same (see Fig. 1), and has an unusually low frequency of gaps. When gaps exist, they are only one residue long. In addition to the active site sequence Cys32 – Gly33 – Pro34 – Cys35, a number of residue positions are highly conserved (using the *E. coli* numbering): Asp26, Ala29, Trp31, Asp61, Pro76 (*cis*), and Gly92. Avian and mammalian thioredoxins contain additional nonredox active cysteine residues that are not involved in redox function but are implicated in regulation. Human thioredoxin forms an inactive intermolecular homodimer *via* a disulfide bond between Cys73 in each monomer (24). In general, residue substitutions among the thioredoxins occur predominantly on the surface of the molecule and distant from the active site.

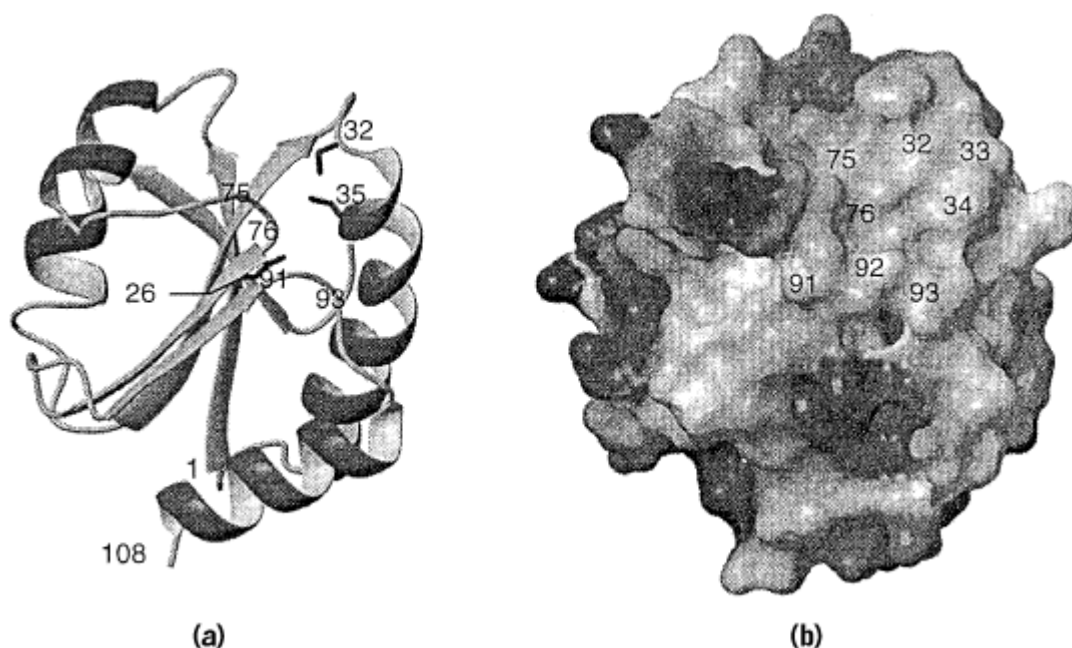
**Figure 1.** Amino acid sequence alignment of selected thioredoxins. Original references to the sequences are given in Ref sequences were aligned on the basis of the three-dimensional structures of *E. coli* (29) and human thioredoxins (30).

|                            |  |    |    |    |     |
|----------------------------|--|----|----|----|-----|
|                            | 10   | 20 | 30 | 40 | 5   |
| <i>E.coli</i> <sup>a</sup> | SDKIIHLTDDSFDTDLKADG-AILVDFWAEWCGPCKMIAPILDEIADEY        |    |    |    |     |
| yeast <sup>b</sup>         | MVTQLKSASEYDSALASGDK-LVVVDFFATWCGPCKMIAPMIEKFAEQY        |    |    |    |     |
| spinach <sup>c</sup>       | MEAIVGKVTEVNKDTFWPIVKAAGDKPVVLDMFTQWCGPCKAMAPKYEKLAEEY   |    |    |    |     |
| human <sup>d</sup>         | VKQIESK-TAQEALDAAGDKLVVDFSATWCGPCKMIKPFHSLSEKY           |    |    |    |     |
| structure                  | β  | α  | β  | α  |     |
|                            | 60   | 70 | 80 | 90 | 100 |
| <i>E.coli</i> <sup>a</sup> | GKLTVAKLNIDQNPGTAPKYGIRGIPTLLLFKNGEVAATKVGALSKGQLKEFLD   |    |    |    |     |
| yeast <sup>b</sup>         | DAAFYKLDV-DEVSDVAQKAEVSSMPTLIFYKGGKEVTRVVGANPAAIKQAIAS   |    |    |    |     |
| spinach <sup>c</sup>       | DVIIFLKLDNCNENKTLAKELGIRVVPTFKILKENSVVGEVTGAKYDKLLEA IQA |    |    |    |     |
| human <sup>d</sup>         | -NVIFLEVDVDDCQDVASECEVKCMPTFOFFKKGQKVGFEFSGANKEKLEATINE  |    |    |    |     |
| structure                  | β  | α  | α  | β  | β   |

#### 4. Thioredoxin Three-Dimensional Structures

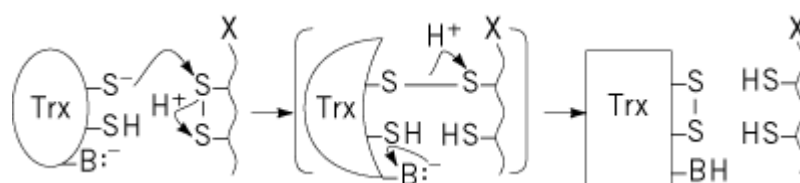
The first three-dimensional structure of a thioredoxin determined by X-ray crystallography techniques (25) was that of the oxidized *E. coli* protein. It consists of a five b-strands and four  $\alpha$ -helices. The molecule is organized into an N-terminal bab and a C-terminal bba motif connected by an  $\alpha$  helix (Fig. 2). *E. coli* Trx can be cleaved proteolytically into two fragments by **trypsin** at Arg73, roughly between these two motifs. Although neither fragment alone is active enzymatically, they can reconstitute a native-like molecule that has enzymatic activity (26). The central five-stranded mixed b sheet lies between the second and last helices on one face and the connecting helix on the other.

**Figure 2.** Representative three-dimensional NMR solution structure of reduced *E. coli* thioredoxin, illustrating the active site cysteines (Cys32 and Cys35) and the hydrophobic interaction surface area. (a) Cartoon of the polypeptide backbone where the side chains of residues 26, 32, and 35 are displayed as sticks. (b) The molecular surface is colored according to the electrostatic potential (dark gray: positive, medium gray: negative, and light gray: uncharged). Residues are labeled at their approximate positions.



Several high-resolution thioredoxin structures have been determined by using both X-ray crystallography, for example, *E. coli* (27) and human (28), and NMR spectroscopy, *E. coli* (29) or human (30). Comparison of the available three-dimensional structures confirms the close structural similarity suggested by their sequence [homology](#). Despite the biochemical evidence of functional differences between the oxidized and reduced forms of the *E. coli* protein (eg, with T7 gene 5 protein), structural differences between the two forms are subtle and localized to the active-site region (31). Differentiation between the two forms could be the consequence of increased accessibility of conformational states other than the ground state in the reduced form, relative to the oxidized forms. The N-terminal redox-active sulfur atom (Cys32) is located in a protuberance of the molecule and is accessible to solvent from one side in both oxidized and reduced forms of the protein. The  $pK_a$  of this thiol group is believed to be  $\sim 7$  which is 1 to 2 pH units lower than normal. Adjacent to the accessible face of nucleophilic Cys32 is a large, relatively flat hydrophobic patch on the protein surface composed of atoms from residues Gly33, Pro34, Ile75, Pro76, Val91, Gly92, and Ala93 (Fig. 2). The existence of a single hydrophobic interaction surface is consistent with the broad specificity and covalent disulfide intermediate in the mechanism of thioredoxin action (Fig. 3). A number of charged side chains, for example, Glu30, Lys36 and Lys57, are on the opposite face of the protuberance, directly behind the active site cysteine residues, and are of possible mechanistic importance. Largely buried between these residues is the highly conserved Asp 26 which is believed to have an abnormally high  $pK_a$  of 7.5 (32) and has been proposed as the general base for deprotonation of the C-terminal (Cys 35) cysteine in the mixed-disulfide intermediate.

**Figure 3.** Mechanism of reduction of the disulfide bond of a protein substrate (X) by thioredoxin (Trx). Reduced thioredoxin (left) binds to a target protein via the hydrophobic interaction surface, followed by nucleophilic attack of the N-terminal active site thiolate on the target disulfide in a **thiol-disulfide exchange** reaction, resulting in a transient protein–protein mixed disulfide (center). Intramolecular attack by the second thiol group of Trx results in oxidized thioredoxin and reduced target protein.



## 5. Intermolecular Complexes

Structures of intermolecular complexes between thioredoxins and the proteins and fragments with which they interact have helped to elucidate the nature of the [protein–protein interactions](#) involved (7, 33, 34). The interface between Trx and its substrate protein includes the aforementioned hydrophobic interaction surface. Interestingly, human Trx binds substrate peptides in two essentially antiparallel orientations, with hallmark **hydrogen bonds** between the backbone of the substrate Cys residue and backbone atoms centered on the residue of Trx preceding the *cis*-Pro76. Thus, Trx can utilize the hydrophobic interaction surface to accommodate substrates that have little sequence homology and might also involve **molecular chaperone**-like conformational changes.

## 6. Thioredoxin Superfamily

A growing number of proteins contain thioredoxin-like structural motifs, and these proteins are now

called members of the thioredoxin superfamily. Proteins that contain the thioredoxin fold can be grouped into at least six classes: thioredoxins, glutaredoxins, DsbA, protein disulfide isomerase, **glutathione transferases**, and glutathione peroxidases (35). The thioredoxin fold common to each of these protein families consists of the bab and bba motifs, that have insertions or extensions of polypeptide chain. Although the sequence homology among these six classes is limited and no function or activity is common to all of them, there is a functional similarity common to four of these members. Thioredoxins, glutaredoxins, DsbA, and protein disulfide isomerase are redox-active proteins that contain a -Cys-X<sub>1</sub>-X<sub>2</sub>-Cys- active-site motif (where X represents any of the 19 commonly-occurring non-cysteine amino acids). Furthermore, although the remaining two proteins, the glutathione transferases and glutathione peroxidases, lack the -Cys-X<sub>1</sub>-X<sub>2</sub>-Cys- active-site motif, they share with the glutaredoxins a specific interaction with the ubiquitous cysteine-containing peptide glutathione. Interestingly, the orientation of the glutathione in the glutaredoxin and glutathione transferases is similar to that in the cysteine-containing segment of the REF-1 peptide in the complex with human Trx, possibly indicating an overall conserved mode of interaction.

### Bibliography

1. T. C. Laurent, E. C. Moore, and P. Reichard (1964) *J. Biol. Chem.* **239**, 3436–3444.
2. A. Holmgren (1968) *Eur. J. Biochem.* **6**, 475–484.
3. A. Holmgren (1979) *J. Biol. Chem.* **254**, 9113–9119.
4. A. Holmgren (1979) *J. Biol. Chem.* **254**, 9627–9632.
5. A. Holmgren (1989) *J. Biol. Chem.* **264**, 13963–13966.
6. D. F. Mark and C. C. Richardson (1976) *Proc. Natl. Acad. Sci. USA* **73**, 780–784.
7. S. Doublíé, S. Tabor, A. M. Long, C. C. Richardson, and T. Ellenberger (1998) *Nature* **391**, 251–258.
8. A. Holmgren (1976) *Proc. Natl. Acad. Sci. USA* **73**, 2275–2279.
9. C. A. Lunn and V. Pigiet (1986) In *Thioredoxin and Glutaredoxin Systems: Structure and Function*. (A. Holmgren, C. I. Brändén, H. Jörnvall, and B.-M. Sjöberg, eds.), Raven Press, New York, pp. 165–176.
10. M. Russel and P. Model (1985) *Proc. Natl. Acad. Sci. USA* **82**, 29–33.
11. E. R. LaVallie, E. A. DiBlasio, S. Kovacic, K. L. Grant, P. F. Scheudel, and J. M. McCoy (1993) *Biotechnology* **11**, 187–193.
12. A. Miranda-Vizuete, A. E. Damdimopoulos, J.-Å. Gustafsson, and G. Spyrou (1997) *J. Biol. Chem.* **272**, 30841–30847.
13. G. Krause and A. Holmgren (1991) *J. Biol. Chem.* **266**, 4056–4066.
14. F. Åslund, K. D. Berndt, and A. Holmgren (1997) *J. Biol. Chem.* **272**, 30780–30786.
15. B. B. Buchanan (1980) *Annu. Rev. Plant Physiol.* **31**, 341–374.
16. P. Schürmann (1995) *Methods Enzymol.* **252**, 274–283.
17. M. Luthman and A. Holmgren (1982) *Biochemistry* **21**, 4600–4606.
18. K. Schulze-Osthoff, H. Schenk, and W. Dröge (1995) *Methods Enzymol.* **252**, 253–264.
19. Y. Tagaya, Y. Maeda, A. Mitsui, N. Kondo, H. Matsui, Y. Hamuro, N. Brown, K. Arai, T. Yokata, H. Wakasugi, and J. Yodoi (1989) *EMBO J.* **8**, 757–764.
20. A. Rubartelli, A. Bajetto, G. Allavena, E. Wollmann, and R. Sitia (1992) *J. Biol. Chem.* **267**, 24161–24164.
21. H. Nakamura, S. De Rosa, M. Roederer, M. T., Anderson, J. G. Dubs, J. Yodoi, A. Holmgren, L. A., Herzenberg, and L. Herzenberg (1996) *Int. Immunol.* **8**, 603–611.
22. H. Z. Chae, S. J. Chung, and S. G. Rhee (1994) *J. Biol. Chem.* **269**, 27670–27678.
23. H. Eklund, F. K. Gleason, and A. Holmgren (1991) *Proteins* **11**, 13–28.
24. X. Ren, M. Björnstedt, B. Shen, M. Ericson, and A. Holmgren (1993) *Biochemistry* **32**, 9701–9708.

25. A. Holmgren, B.-O. Söderberg, H. Eklund, and C.-I. Brändén (1975) *Proc. Natl. Acad. Sci. USA* **72**, 2305–2309.
26. I. Slaby and A. Holmgren (1979) *Biochemistry* **18**, 5584–5591.
27. S. K. Katti, D. M. LeMaster, and H. Eklund (1990) *J. Mol. Biol.* **212**, 167–184.
28. A. Weichsel, J. R. Gasdska, G. Powis, and W. R. Montfort (1996) *Structure* **4**, 735–751.
29. M.-F. Jeng, A. P. Campbell, T. Begley, A. Holmgren, D. A. Case, P. E. Wright, and H. J. Dyson (1994) *Structure* **2**, 853–868.
30. J. Qin, G. M. Clore, and A. M. Gronenborn (1994) *Structure* **2**, 503–522.
31. A. Holmgren (1995) *Structure* **3**, 239–243.
32. H. J. Dyson, M.-F. Jeng, L. L. Tennant, I. Slaby, M. Lindell, D. S. Cui, S. Kuprin, and A. Holmgren (1997) *Biochemistry* **36**, 2622–2636.
33. J. Qin, G. M. Clore, W. M. P. Kennedy, J. R. Huth, and A. M. Gronenborn (1995) *Structure* **3**, 289–297.
34. J. Qin, G. M. Clore, W. M. P. Kennedy, J. Kuszewski, and A. M. Gronenborn (1996) *Structure* **4**, 613–620.
35. J. L. Martin (1995) *Structure* **3**, 245–250.

## Thioredoxin Reductase

Thioredoxin reductase (TrxR) contains the cofactor FAD and catalyzes the reduction of the disulfide form of [thioredoxin](#) (TRX -S<sub>2</sub>) to the dithiol form (TRX -(SH)<sub>2</sub>):



### 1. Thioredoxin Reductase in Prokaryotes, Yeast, and Plants

The dimeric enzyme from *Escherichia coli* is highly specific for NADPH. It has a [K<sub>m</sub> \(Michaelis constant\)](#) of 1.2 μM, and a *K<sub>m</sub>* of 2.8 μM for Trx-S<sub>2</sub>. It acts by a ping-pong mechanism and has a turnover rate of 2000 min<sup>-1</sup> per FAD (see [Kinetic Mechanisms, Enzyme](#)). NADP<sup>+</sup> is a **competitive inhibitor** and has a *K<sub>i</sub>* of 15 μM (1). The [protein structure](#) of the enzyme determined by [X-ray crystallography](#) shows surprisingly large differences from other members of the pyridine-nucleotide disulfide oxidoreductase family (2). The [active-site](#) sequence -Cys-Ala-Thr-Cys- is located in the N-terminal FAD **domain**, suggesting [convergent evolution](#). An interesting exception to the rule of separate thioredoxin and thioredoxin reductase proteins exists in *Mycobacterium leprae*, where a functional hybrid protein is found (3).

### 2. Thioredoxin Reductase in Mammalian Cells

Mammalian thioredoxin reductase has properties strikingly different from the enzymes from *E. coli* or yeast (see Table 1 of [Thioredoxin](#)). The enzyme from rat liver, the first purified to homogeneity, is a dimer that has two FAD molecules and subunits with a molecular mass of 58 kDa (4). All mammalian thioredoxin reductases have a surprisingly wide substrate specificity (Table 1). In

particular, the enzyme reduces selenite directly and is a lipid hydroperoxide reductase (5). Recently, it was shown that human adrenocarcinoma thioredoxin reductase contains a [selenocysteine](#) (SEC) residue (6). The enzyme is **homologous** to glutathione reductase, and has an elongation that contains the conserved C-terminal sequence -Gly-Cys-Sec-Gly in human, rat, and bovine enzymes (Fig. 1). Insertion of Sec requires a selenocysteine-insertion sequence (SECIS), present in the 3'-untranslated region, to decode the UGA normal [stop codon](#) as selenocysteine (7). The -Sec-Gly sequence is required for catalytic activity, because its removal by [carboxypeptidase](#) digestion inactivates the enzyme (7). The selenocysteine in mammalian thioredoxin reductase explains its wide substrate specificity and inhibition by 1-chloro-dinitrobenzene (DCNB) (8). The latter compound selectively modifies both the Sec and Cys residues in the enzyme and induces a 30-fold increase in NADPH-oxidase activity in the alkylated enzyme (9). The requirement for selenium in the enzyme may explain why selenite is required to grow some cells in synthetic tissue-culture media. Tumor cells often have a strong up-regulation of TrxR (10), which may protect the cells. The selenocysteine nature of the enzyme also explains why it is inhibited by some drugs used clinically, such as gold thioglucose or nitrosourea compounds (11).

**Figure 1.** The primary structure of one mammalian thioredoxin reductase subunit from the rat that has 498 residues (7). The conserved N-terminal active -site disulfide region, which is identical to that of glutathione reductase (Cys-Val-Asn-Val-Gly-Cys, or CVNVGC), is shown, plus the conserved C-terminal (Gly-Cys-Sec-Gly, or GCUG) selenocysteine-containing sequence in the human, rat, and bovine enzymes. The positions of the FAD-binding, NADPH-binding, and interface domains are indicated.



**Table 1. Substrates for Mammalian Thioredoxin Reductase**

---

|  |
|--|
| Trx-S <sub>2</sub> from many species                             |
| Vitamin K  |
| Alloxan  |
| Selenite and selenodiglutathione                                 |
| 5,5'-Dithiobis-(2-nitrobenzoic acid) ( <b>Ellman's reagent</b> ) |
| Protein disulfide isomerase (PDI)                                |
| Cu <sup>2+</sup> ions  |
| Lipid hydroperoxides and H <sub>2</sub> O <sub>2</sub>           |
| S-Nitrosoglutathione(GSNO)                                       |

---

## Bibliography

1. C. H. Williams, Jr. (1992) In *Chemistry and Biochemistry of Flavoenzymes* (F. Müller, ed.), FL, CRC Press, Boca Raton, pp. 121–211.

2. G. Waksman, T. S. Krishna, C. H. Williams, Jr., and J. Kuriyan (1994) *J. Mol. Biol.* **236**, 800–816.
3. B. Wieles, J. van Noort, J. W. Drijfout, R. Offringa, A. Holmgren, and T. Ottenhof (1995) *J. Biol. Chem.* **270**, 25604–25606.
4. M. Luthman and A. Holmgren (1992) *Biochemistry* **21**, 6628–6633.
5. M. Björnstedt, M. Hamberg, S. Kumar, J. Xue, and A. Holmgren (1995) *J. Biol. Chem.* **270**, 11761–11764.
6. T. Tamura and T. C. Stadtman (1996) *Proc. Natl. Acad. Sci. USA* **93**, 1006–1011.
7. L. Zhong, E. S. J. Arnér, J. Ljung, F. Åslund, and A. Holmgren (1998) *J. Biol. Chem.* **273**, 8581–8591.
8. E. S. J. Arnér, M. Björnstedt, and A. Holmgren (1995) *J. Biol. Chem.* **270**, 3479–3482.
9. J. Nordberg, L. Zhong, A. Holmgren, and E. S. J. Arnér (1998) *J. Biol. Chem.* **273**, 10835–10842.
10. G. Powis, J. R. Gasdaska, and A. Baker (1997) *Adv. Pharmacol.* **38**, 329–359.
11. K. E. Hill, G. W. McCoolum, and R. F. Bruk (1997) *Anal. Biochem.* **253**, 123–125.

## Thiotemplate Mechanism Of Peptide Antibiotic Synthesis

**Peptide antibiotics** are unusual in that they contain many unusual and modified amino acid residues, even D-amino acid residues, and often have cyclic **polypeptide chains**. Although some are synthesized on [ribosomes](#), in the usual manner of **protein biosynthesis**, and are then **post-translationally modified**, many have unusual, nonribosomal mechanisms of synthesis. The thiotemplate mechanism refers to this type of mechanism of forming [peptide bonds](#) sequentially to generate a polypeptide precursor of peptide antibiotics. The name was coined when thioester intermediates and the presence of 4'-phosphopantetheine, derived from coenzyme A, were discovered in an enzyme system catalyzing the synthesis of the cyclic decapeptide **gramicidin S** (1). It was suggested that amino acid residues were polymerized sequentially on enzyme subunits as thioesters attached to the 4'-phosphopantetheine swinging arm of a carrier protein. It transformed the fatty acid elongation cycle proposed by Lynen (2) to enzymatic peptide formation and introduced the concept of a protein [template](#) for assembly of a peptide sequence (3).

The thioester type of assembly mechanism has been established for a wide variety of peptides, including [tyrocidine](#), bacitracin, polymyxin, enniatin, beauvericin, **actinomycin**, alamethicin, ferrichrome, [cyclosporin](#), and the [penicillin](#) precursor peptide d-(L-a-aminoadipoyl)-L-cysteinyl-D-valine. It soon became apparent that it was the principal method of forming peptide bonds in the absence of ribosomes, and it has been shown to operate for at least 25 steps, polymerizing that number of amino acids specifically (4).

When it became evident that the protein template enzymes were unusually large structures, with masses often exceeding 500 kDa, Kurahashi introduced the term *multienzyme thiotemplate mechanism* (5). It was later pointed out, however, that the principle of the 4'-phosphopantetheine swinging arm did not account for all of the experimental observations, including (i) the accumulation of all intermediates preceding a certain step when the process was interrupted there, (ii) the loss of thioester formation, with retainment of adenylate activation, and (iii) constraints on the access of a single cofactor on an integrated carrier domain to 11 amino acid intermediates in a system like cyclosporin synthetase.



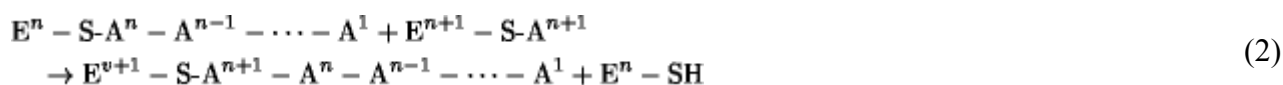
These issues were resolved after the amino acid sequences of several peptide synthetases were determined from the sequences of their genes (6). This revealed that the peptide synthetases had modular structures composed of (i) adenylate-forming activation domains, (ii) carrier proteins containing sites for attachment of the respective cofactor, (iii) condensation domains, and (iv) a variety of modification functions responsible for epimerization, methylation, cyclization and hydrolytic release of the precursor peptide. Most striking was the absence of conserved [cysteine](#) residues that had been implied by the thioester concept. It became obvious that each activation domain was associated with a carrier protein **homologous** to the *acyl carrier protein* of the fatty acid synthesis systems, and the subsequent addition of 4'-phosphopantetheine to the carrier proteins has been demonstrated by [mass spectrometry](#). These findings have led to a further refinement of the process, which is now called the *multiple carrier thiotemplate mechanism* (7).

## 1. Multiple Carrier Thiotemplate Mechanism

To illustrate the multiple carrier thiotemplate mechanism, a scheme for the formation of a tripeptide is illustrated in Figure 1. Each amino acid,  $A^i$ , is adenylated and activated by a specific adenylation domain of the enzyme and then transferred to the [thiol group](#) of the pantetheine cofactor residing on its respective carrier domain,  $E^i$ -SH, as a thioester,  $E^i$ -S- $A^i$ . Peptide bond formation is initiated with the amino acid that will be C-terminal,  $A^1$ , becoming linked to the next amino acid,  $A^2$ ,

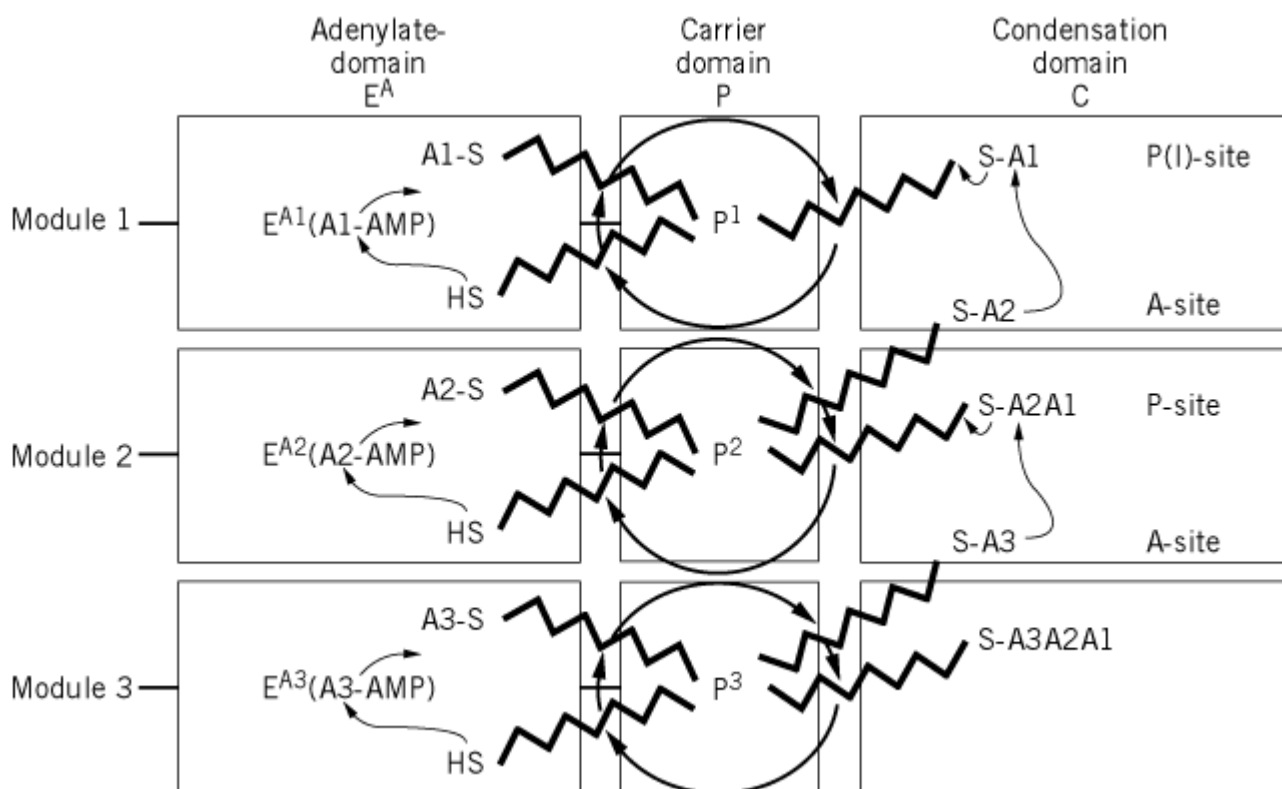


Elongation proceeds similarly and sequentially, adding a third and further amino acid thioesters to the peptide:



In this way, a polypeptide chain is synthesized starting from the carboxyl terminus. At each step, the assembly process is specific for the amino acid to be incorporated and for the appropriate [enantiomer](#), either D or L. Hydroxy acids can also be incorporated, and the corresponding ester bonds are formed similarly, using acyl intermediates. The process thus resembles in principle the ribosomal mechanism, in that there is transfer of an activated peptidyl intermediate to an aminoacyl intermediate (8). Release of the complete peptide chain occurs on the appropriate termination domain, either as the linear chain, which is produced by hydrolysis or aminolysis, or as the cyclic peptide or as peptidolactone.

**Figure 1.** Illustration of the multiple carrier thiotemplate mechanism, as elucidated for synthesis of the d-(L-aminoadipoyl)-L-cysteinyl-D-valine tripeptide  $\beta$ -lactam precursor (11). Each of the three amino acids is activated at a specific adenylation domain,  $E^A$ . Each adenylate domain is linked to a specific carrier domain,  $P^i$ . The aminoacyl moiety  $A^i$  of the  $A^i$ -AMP adenylates are transferred to the respective thiol groups of the 4'-phosphopantetheine cofactor, illustrated by the thick zigzag line. Peptide bond formation occurs, in analogy to the ribosomal cycle, from aminoacyl (A) and peptidyl (P) sites on specific condensation domains. In the case of the initiation reaction, the P-site is called the initiation site, or I site. The tripeptide formed,  $A^3$ - $A^2$ - $A^1$ , is thioester bound and may be transferred for further elongation or released by a thioesterase, as in penicillin formation.



The reactions linking the successive amino acids and analogues are thought to be catalyzed by the condensation domains; they are essentially irreversible. In addition, catalytic sites may be involved in the modification of the thioester intermediates, such as their *N*-methylation or epimerization and, especially in the case of polyketide biosynthesis, in the reduction and dehydration of the respective keto intermediate.

The key component of the thiotemplate mechanism is the acyl-, aminoacyl-, or peptidyl-carrier protein or domain, to which the respective intermediate remains attached covalently during the elongation cycle. Carrier domains can be identified readily by the conserved elements of their [protein structure](#), including a cavity formed by three  **$\alpha$ -helices**, an additional  $\alpha$ -helix, and a variable loop containing the 4'-phosphopantetheine cofactor binding site (9). Addition of the cofactor to the carrier domain is catalyzed by a family of 4'-phosphopantetheine transferases, using coenzyme A (10). Condensation processes require at least one carrier protein, and stepwise processes involve the intra- or intermolecular transfer of acyl intermediates between two carrier protein proteins. The largest known thiotemplates contain 11 carrier domains in peptide biosynthesis, for cyclosporin, and six carrier domains in polyketide biosynthesis, for rapamycin.

The thiotemplate mechanism contrasts with the mechanism of formation of other small peptides, such as [glutathione](#) or peptidoglycan components, which are readily distinguished by (i) their use of phosphate to activate the amino acids, rather than adenylate, (ii) the absence of covalently attached intermediates, and (iii) the activation of peptide carboxyl groups, rather than amino acid carboxyls.

#### Bibliography

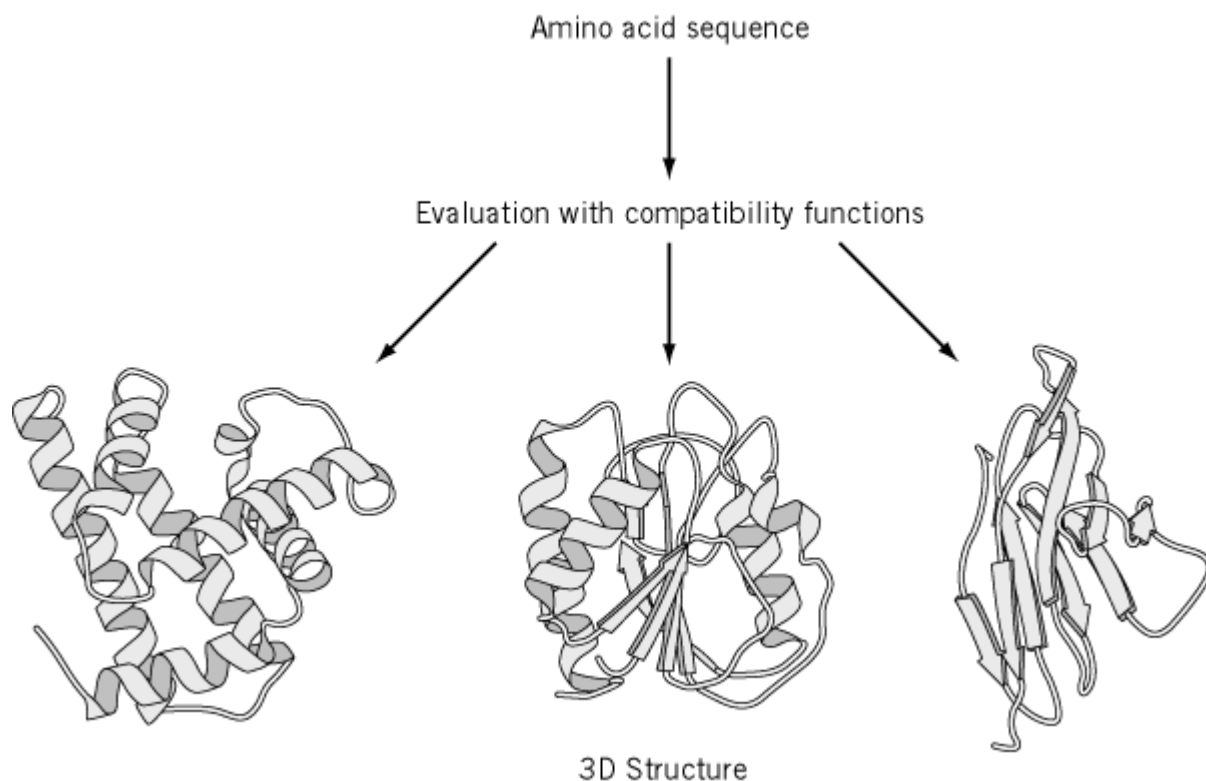
1. S. G. Laland and T.-L. Zimmer (1973) *Essays Biochem.* **9**, 31–57.
2. F. Lynen (1967) *Biochem. J.* **102**, 381–400.
3. F. Lipmann (1973) *Science* **173**, 875–884.
4. H. Kleinkauf and H. von Döhren (1990) *Eur. J. Biochem.* **192**, 1–15.

5. K. Kurahashi (1981) "Peptide Antibiotics". In *Biosynthesis*, Vol. IV (J. W. Corcoran, ed.), Springer, Berlin, pp. 325–352.
6. H. Kleinkauf and H. von Döhren (1996) *Eur. J. Biochem.* **236**, 335–351.
7. T. Stein, J. Vater, V. Kruft, B. Wittmann-Liebold, P. Franke, M. Panico, R. McDowell, and H. R. Morris (1996) *J. Biol. Chem.* **271**, 15428–15435.
8. D. A. Hopwood (1997) *Chem. Rev.* **97**, 2465–2497.
9. M. P. Crump, J. Crosby, C. E. Dempsey, J. A. Parkinson, M. Murray, D. A. Hopwood, and T. J. Simpson (1997) *Biochemistry* **36**, 6000–6008.
10. R. H. Lambalot et al. (1996) *Chem. Biol.* **3**, 923–936.
11. H. von Döhren, U. Keller, J. Vater, and R. Zocher (1997) *Chem. Rev.* **97**, 2675–2705.

## Threading Protein Sequences

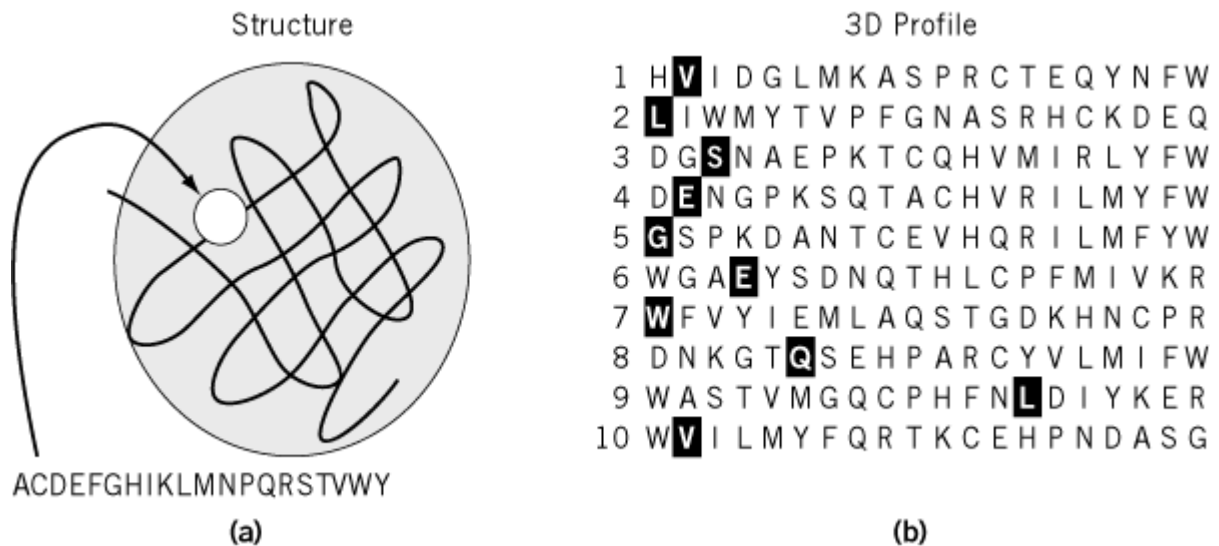
Formulation of the [inverse folding problem](#), in which the amino acid sequence, or [primary structure](#), most compatible with a given three-dimensional [protein structure](#) is sought, provided a clue to the problem of **protein structural prediction**, in which the three-dimensional structure is to be predicted from just the primary structure. Given a query sequence  $a$  of a protein of unknown structure (A), mount the sequence onto a number of known protein structures (B, C, D, etc.) one-by-one, using all feasible alignments of sequence to structure. Then, evaluate the fitness of sequence  $a$  to each one of the structures to find the structure most compatible with the sequence  $a$  (Fig. 1). If the structure so found is B, it is inferred that the structure A to be predicted is probably similar to B. This is called the “3D–1D compatibility method,” “fold recognition,” or “prediction by threading.” It is an extension of the conventional comparison of sequences alone (1D–1D), but in principle is more effective in detecting similarities between proteins because the 3-D structure is more conserved during [evolution](#) than the sequence (see [Homological Modeling](#)). It is observed empirically that **homologous** natural proteins that arose by divergent evolution from a common ancestor and belong to a **gene family** share a common polypeptide fold. Moreover, the number of folds found in natural proteins is believed to be fairly limited (1). Because of the rapid growth of the protein structural [database](#), it was logical for the threading approach to emerge. The logic of supposing that a known structure can be a model for a new protein arose from the inverse folding problem, but the difference lies in that the threading method considers a certain (query) sequence in relation to all possible model structures, whereas inverse folding considers a model structure in relation to all the possible sequences.

**Figure 1.** Protein structure prediction by threading. A query sequence is to be threaded through each of template structures taken from the structural database, and the sequence/structure fitness is quantitatively evaluated with compatibility functions to find a structure most compatible with the sequence.



In threading a given 1-D sequence through a given 3-D structure, there are a large number of possibilities, which is known as the alignment problem. Bowie (2, 3) was the first to solve this problem by applying the dynamic programming technique that had been fully developed in the field of [aligning sequences](#), so-called homological searches (see [Sequence Analysis](#)). To reduce the magnitude of the computational problem, Bowie et al. (3) introduced the 3-D profile table, which was constructed from the model structure (Fig. 2). This is a  $(20 \times n)$  table, where columns of 20 amino acid residues are arrayed along the  $n$  residues of the structure. Each number in the table gives the fitness of the respective amino acid residue for a given residue site, which depends on its **secondary structure, hydrophilic/hydrophobic** environment, etc. Given such a 3-D profile table, it is straightforward to compare it with any sequence by using the dynamic programming algorithm and to obtain the optimum path, including gaps (i.e., the best 3-D–1-D alignment), plus the alignment score. In actual predictive procedures, a structural library that has been converted to a set of 3-D profiles is scanned with a query sequence to seek the model structure and alignment that gives the highest score (or lowest energy). If the score obtained is sufficiently high, a convincing prediction has been obtained. The computational time is no greater than that for the usual sequence homological search methods.

**Figure 2.** Construction of the 3-D Profile. (a) Each of 20 amino acids is placed one-by-one at a site of a given 3-D structure of a protein to evaluate its compatibility with the structural environment. A, C, D . . in the figure are one-letter codes of amino acids standing for Ala, Cys, Asp. ., respectively. The results are tabulated in a profile table (b). The example in this table was constructed from sperm whale myoglobin structure (PDB code: 1mbd). The profile table for the first 10 residue sites (one row corresponding to one site) shown consists of columns of the length equal to the total number of residue sites in the structure. In the profile table of this illustration, the 20 amino acids are sorted from left to right at each site, according to the compatibility score. Those amino acids in the native sequence that are highlighted sit on the left-hand side of the table, implying that they are energetically favorable to the structure.



In addition to finding the best 3D–1D alignment, it is important to evaluate the fitness between a sequence and a structure. Bowie et al. (3) used a simple measure, called the “single-body approximation,” which considers a single amino acid residue in its entire structural environment. The Sippl potential (see [Inverse Folding Problem](#)) has a more advanced form of the two-body function defined for two interacting residues. Interactions between a central residue and all surrounding residues would give the total interaction energy for the central residue. The two-body function is not logically compatible with the 3-D profile, however, because the surrounding residues will not be known until after the alignment is fixed. This dilemma has been solved by the “frozen approximation” (4), in which it is assumed that the native amino acids of the model structure are in the surrounding sites. In this way, a combination of the Sippl-type potential (or others) with the 3-D profile treatment provides rapid and more reliable predictive methods (5).

Other types of methods can be envisaged for directly threading the entire sequence through a structure without using the 3-D profile, and they could be compatible with any type of potential functions. The problem, however, is the large number of alignments possible. Jones et al. (6) managed to solve the problem by introducing the double dynamic programming algorithm, which had been developed in pairwise comparisons of 3-D structures by Taylor (7). Alternatively, **Monte Carlo**-type optimizations may be used (8). The computations for these methods of direct threading are very time-consuming.

Many investigators have contributed to the development of various threading methods. To compare the effectiveness of each, a worldwide prediction contest, named (Critical Assessment of Techniques for Protein Structure Prediction CASP) has been held twice thus far, in 1994 and 1996 (9). A completely blind test was arranged so that all predictor would make their own predictions for target proteins, whose 3D structures were not yet known but were in the process of being determined. Then, an appointed assessor evaluated all of the predictions by comparing them with the structures that had subsequently been determined and presented the results at a joint meeting of all of the predictors. The results indicated that several predictions, but not all, actually inferred the correct folds of target proteins from their sequences alone (10). The sequence identity between the known and predicted proteins was around 20% or less, too low for any methods other than threading to identify the structural relationship. Therefore, the prediction contests clearly demonstrated the applicability of the threading method. At the same time, however, some drawbacks were revealed. One of them is the ambiguity often seen in 3-D–1-D alignments obtained by direct comparisons of structures (11). There are still problems to be overcome, and the method should be refined further. In any case, several successful examples uncovered by CASP were the first clear verification of the correct prediction of 3-D protein structures after a number of unsuccessful trials over many years.

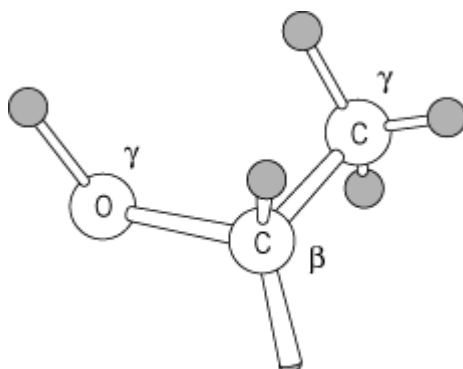
## Bibliography

1. C. Chothia (1992) *Nature* **357**, 543–544.
2. J. U. Bowie, N. D. Clarke, C. O. Pabo, and R. T. Sauer (1990) *Proteins: Struct. Function Genet.* **7**, 257–264.
3. J. U. Bowie, R. Luthy, and D. Eisenberg (1991) *Science* **253**, 164–170.
4. A. Godzik, A. Kolinski, and J. Skolnick (1992) *J. Mol. Biol.* **227**, 227–238.
5. Y. Matsuo and K. Nishikawa (1994) *Protein Sci.* **3**, 2055–2063.
6. D. T. Jones, W. R. Taylor, and J. M. Thornton (1992) *Nature* **358**, 86–89.
7. W. R. Taylor and C. A. Orengo (1989) *J. Mol. Biol.* **208**, 1–22.
8. S. H. Bryant and C. E. Lawrence (1993) *Proteins: Struct. Function Genet.* **16**, 92–112.
9. J. Moulton, J. T. Pedersen, R. Judson, and K. Fidelis (1995) *Proteins: Struct. Function Genet.* **23**, 2–4.
10. C. M.-R. Lemer, M. J. Rooman, and S. J. Wodak (1995) *Proteins: Struct. Function Genet.* **23**, 337–355.
11. A. Marchler-Bauer and S. H. Bryant (1997) *Proteins: Struct. Function Genet. (Suppl. 1)*, 74–82.

## Threonine (Thr, T)

The [amino acid](#) threonine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to four **codons**—ACU, ACC, ACA, and ACG—and represents approximately 5.8% of the residues of the proteins that have been characterized. The threonyl residue incorporated has a mass of 101.11 Da, a **van der Waals volume** of  $93 \text{ \AA}^3$ , and an [accessible surface](#) area of  $146 \text{ \AA}^2$ . Thr residues are changed moderately frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [serine](#), [alanine](#), [valine](#), and [lysine](#) residues.

The side chain of Thr is dominated by its hydroxyl group:



This hydroxyl group is normally no more reactive chemically than is that of ethanol, so there are few chemical reactions that can modify Thr residues in a protein specifically and readily. The only reaction that is generally useful is acetylation with acetyl chloride in aqueous trifluoroacetic acid. Note that the side chain has a center of asymmetry at  $C_{\beta}$  and that only the one isomer occurs naturally.

The hydroxyl group is very polar and [hydrophilic](#) and can function as either a donor or acceptor in [hydrogen bonds](#). Consequently, Thr residues can be fully buried in protein structures, and about 23% are, especially if their hydroxyl group is paired in a hydrogen bond. The hydroxyl group is situated sterically to interact with polar groups of the polypeptide backbone, which affects its conformation and reactivity. For example, peptide bonds adjacent to Thr residues are especially susceptible to acid hydrolysis. Thr residues do not favor the  $\alpha$ -helical conformation, and they occur in folded proteins most frequently in [beta-sheets](#).

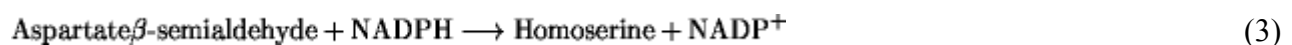
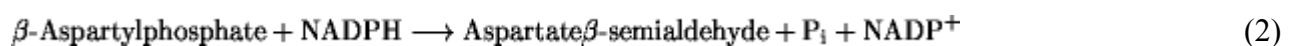
In native proteins, Thr residues are frequently modified by the addition of oligosaccharides (see [O-Glycosylation](#)) and phosphate groups.

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Threonine Operon

The four carbon atoms of aspartic acid in *Escherichia coli* are the precursors of the four carbon atoms of threonine. In addition, with pyruvic acid they contribute to the biosynthesis of lysine and diaminopimelate and, with the  $\beta$ -carbon of serine, to the biosynthesis of methionine (1). The reactions leading from aspartate to threonine are as follows:



Reaction 1 is catalyzed by the [enzyme](#) aspartate kinase I, encoded by the **gene** *thrA*; reaction 2 by aspartate semialdehyde dehydrogenase (*asd*); reaction 3 by homoserine dehydrogenase, encoded by *thrA*; reaction 4 by homoserine kinase (*thrB*); and reaction 5 by threonine synthase (*thrC*).

Note that both reactions 1 and 3 are catalyzed by a single bifunctional protein, aspartate kinase I-homoserine dehydrogenase I, encoded by *thrA*. Another bifunctional enzyme, aspartate kinase II-homoserine dehydrogenase II, is encoded by *metL*, whereas aspartate kinase III is encoded by *lys C*. The three aspartate kinases catalyze the same reaction and differ in the mode of regulating their activity and their synthesis (1).

The genes *thrA*, *thrB*, *thrC* are organized in this order in an [operon](#), localized at min 0 on the *E. coli*

[chromosome](#) (2, 3). The full-length [messenger RNA](#) is 4860 nucleotides long (4, 6). The coding sequences of *thrA* and *thrB* are separated by only one nucleotide, whereas those of *thrB* and *thrC* are contiguous (5). Therefore the [Shine–Dalgarno sequence](#) of *thrB* is within the coding sequence of *thrA*, and that of *thrC* is within the coding sequence of *thrB*.

A strong **promoter** is located 190 bp upstream of the [translation](#) start of *thrA*. In addition to the transcripts initiated at this promoter, an internal promoter at the 3' end of *thrA* allows the transcription of *thrB*, although with much lower efficiency (7). Translational coupling between *thrA* and *thrB* has been demonstrated (8). A detailed genetic map of *thrA* shows that the gene is indeed composed of two segments corresponding, respectively, to aspartate kinase I and homoserine dehydrogenase I (9).

Expression of the threonine operon depends on the intracellular concentrations of both threonine and isoleucine (10) by a mechanism called *multivalent repression* (11). In diploids, some derepressed mutations were cis-acting, while other mutations acted in trans. The trans-acting mutations affected either the threonyl-tRNA synthetase (12, 13) or the isoleucyl-tRNA synthetase (14) (see [Aminoacyl tRNA Synthetases](#)). The cis-acting mutations were first thought to be operator mutants (15, 16) but were localized in a region of the gene that had all the characteristics of an **attenuation** mechanism: (1) the mRNA leader sequence shows the possibility of several mutually exclusive secondary structures; (2) there is a r-independent signal for transcription termination; (3) finally, the leader peptide sequence contains numerous threonine and isoleucine codons: Thr-Thr-Ile-Thr-Thr-Thr-Ile-Thr-Ile-Thr-Thr-Thr (17). Further study of derepressed mutants supports the attenuation mechanism (18). In particular, a mutant carrying a deletion of the leader sequence is derepressed (19).

The organization of the three genes is different in other species, such as *Bacillus subtilis* (20), *Corynebacterium glutamicum* (21), and *Pseudomonas aeruginosa* (22). In some of these organisms, there is a single aspartate kinase that is not covalently linked with homoserine dehydrogenase in a multifunctional protein (23).

## Bibliography

1. G. N. Cohen (1983) In K. M. Herrmann and R. L. Somerville, eds., *Amino acids: Biosynthesis and Genetic Regulation*, Addison–Wesley, London, U.K., pp. 147–171.
2. J. Thèze et al. (1974) *J. Bacteriol.* **117**, 133–143.
3. J. Thèze and I. Saint Girons (1974) *J. Bacteriol.* **118**, 990–998.
4. M. Katinka et al. (1980) *Proc. Natl. Acad. Sci. USA* **77**, 5730–5733.
5. P. Cossart, M. Katinka, and M. Yaniv (1981) *Nucleic Acids Res.* **9**, 339–347.
6. C. Parsot et al. (1983) *Nucleic Acids Res.* **11**, 7331–7345.
7. I. Saint Girons and D. Margarita (1985) *J. Bacteriol.* **161**, 461–462.
8. S. Little et al. (1989) *J. Bacteriol.* **171**, 3518–3522.
9. I. Saint Girons and D. Margarita (1978) *Mol. Gen. Genetics* **162**, 101–107.
10. G. N. Cohen and J-C Patte (1963) *Cold Spring Harbor Symp. Quant. Biol.* **28**, 513–516.
11. M. Freundlich (1963) *Biochem. Biophys. Res. Comm.* **10**, 277–282.
12. G. Nass and J. Thomale (1974) *FEBS Lett.* **39**, 182–186.
13. E. J. Johnson, I. Saint Girons, and G. N. Cohen (1975) *J. Bacteriol.* **129**, 66–70.
14. J. M. Blatt and H. E. Umbarger (1972) *Biochem. Genet.* **6**, 99–118.
15. J. F. Gardner and O. H. Smith (1975) *J. Bacteriol.* **124**, 161–166.
16. I. Saint Girons and D. Margarita (1975) *J. Bacteriol.* **124**, 1137–1141.
17. J. F. Gardner (1979) *Proc. Natl. Acad. Sci. USA* **76**, 1706–1710.
18. J. F. Gardner (1982) *J. Biol. Chem.* **257**, 3896–3904.
19. C. Parsot, I. Saint Girons, and P. Cossart (1982) *Mol. Gen. Genet.* **188**, 455–458.



20. C. Parsot and G. N. Cohen (1988) *J. Biol. Chem.* **263**, 14654–14660.
21. M. T. Follettie, H. K. Shin, and A. J. Sinskey (1988) *Mol. Microbiol.* **2**, 53–62.
22. C. Clepet et al. (1992) *Mol. Microbiol.* **6**, 3109–3119.
23. G. N. Cohen (1994) *Biosyntheses*, Chapman and Hall, New York, pp. 378–386.

## Thrombin

Thrombin is the pivotal enzyme in the process of blood clotting or coagulation. It is the final proteinase of the proteolytic cascade, catalyzing the cleavage of two peptide bonds in its primary substrate fibrinogen, to generate fibrin, which subsequently polymerizes to form an insoluble clot. Thrombin has a multiplicity of additional functions, however, that contribute to both the amplification of blood clotting and its regulation. Thrombin promotes feedback activation of the coagulation cascade by activating the cofactors, factors V and VIII, as well as activating the proteinase zymogens, factors XI and XIII. Thrombin also stimulates the aggregation of blood platelets by activation of the thrombin receptor or PAR-1 (protease-activated receptor-1), one of four members of a sub-family of seven-transmembrane G-protein-coupled receptors, by a novel proteolytic “tethered ligand” mechanism (1). Thrombin also has anticoagulant properties and functions, because binding to the endothelial cell membrane protein thrombomodulin alters the specificity of thrombin, enabling it to activate another proteinase zymogen, protein C, which can, in turn, proteolytically inactivate factors V and VIII. Through activation of PAR-1, PAR-3, and PAR-4 thrombin can also mediate responses in a variety of cell types other than platelets, including endothelial cells, vascular smooth muscle cells and fibroblasts. Thus, despite its unique role in the generation of fibrin, thrombin can be viewed as a pleiotropic molecule, exhibiting both enzymatic and hormone-like properties.

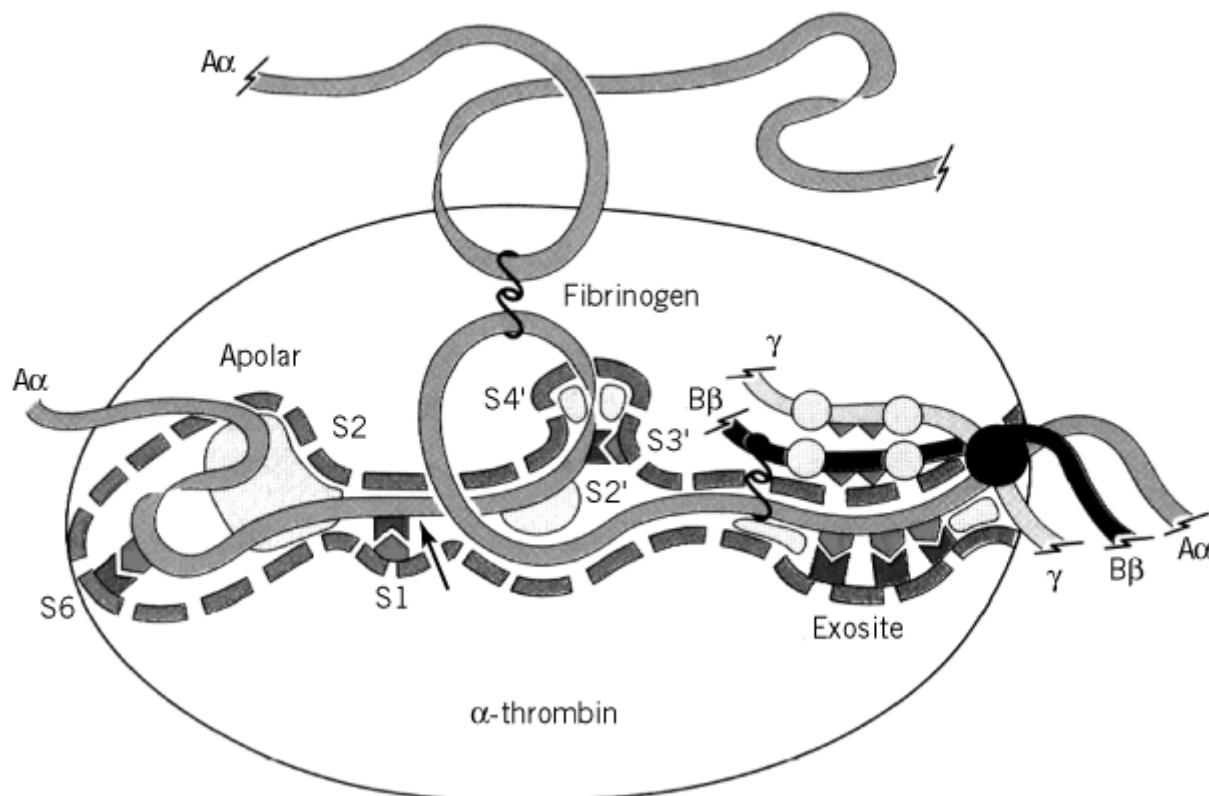
Thrombin itself is generated by limited proteolysis of prothrombin, a 579 residue single-chain glycoprotein ( $M_r$  71,600) expressed in the liver as the product of a 19 kb gene comprising 13 introns and 14 exons located on chromosome 11. Prothrombin undergoes important post-translational modifications in addition to N-linked glycosylation at Asn 78, 100 and 373. The first 10 N-terminal glutamic acid residues, contained within a single 32-residue domain (known as the *Gla domain*), are modified to  $\gamma$ -carboxyglutamic acid (Gla) by a vitamin K-dependent carboxylase, a modification common to all of the “vitamin K-dependent” coagulation factors, namely, factors VII, IX, X, protein C, and prothrombin. This modification is signaled for by residues –17 to –1 of the preprothrombin sequence, adjacent to the secretion signal peptide at –43 to –18. These Gla residues chelate  $Ca^{2+}$ , leading to an ordering of the structure of the N-terminal structural domain, which is then able to bind to negatively-charged phospholipids. Prothrombin also contains two kringle domains that lack the lysine-binding properties of other kringle modules [see [Plasminogen](#)] and appear to have no discrete biological function, and the C-terminal serine proteinase domain. Structures for all the domains of prothrombin – specifically, fragment 1 (Gla-kringle), fragment 2 (kringle 2), and the serine proteinase domain – have been solved by X-ray crystallography and/or NMR. The structure of a-thrombin has been studied particularly extensively, in order to elucidate the molecular basis of its highly restricted substrate specificity.

The serine proteinase factor  $X_a$  in the prothrombinase complex [see [Blood Clotting](#)], activates prothrombin by cleavage at two positions (2). The first cleavage is at the Arg320-Ile321 peptide bond, generating “meizothrombin”, a disulfide bond-linked two-chain proteinase that has a competent active-site, although lacking the specificity of mature thrombin. Subsequent cleavage at

Arg271-Thr272 liberates the N-terminal fragment 1·2, which retains the phospholipid binding characteristics of prothrombin, plus fully mature  $\alpha$ -thrombin. In the absence of cofactors, factor  $X_a$  cleaves the same two peptide bonds in the reverse order, first generating fragment 1·2 and the catalytically inactive prethrombin 2, which is subsequently activated to  $\alpha$ -thrombin. The loss of the N-terminal domains on activation is unique among the blood coagulation proteinases and results in thrombin having full activity in the solution-phase, with no cofactor requirements.

Thrombin consists of two disulfide-bridged polypeptide chains. The shorter A chain forms an integral part of the structure but not involved in catalysis. The B chain exhibits the characteristic fold of the trypsin-like serine proteinases, having two 6-stranded  $\beta$ -barrels with the equivalents of the catalytic residues Ser195, His57, and Asp102 (chymotrypsin numbering, corresponding to Ser525, His363, and Asp419 of thrombin, respectively) located at their interface. As with trypsin, substrate recognition involves a basic P1 residue, but thrombin has a much more extensive substrate recognition surface, as shown in Figure 1. Substrate recognition, and its very restricted specificity, is dictated largely by a number of unique insertion loops that border the active-site cleft, rendering it both deeper and narrower than that of chymotrypsin. These loops limit the access of many potential substrates and inhibitors to the active site, as well as imposing constraints on individual subsites. An example of this is the “60-loop” (a four-residue insertion at a position corresponding to Val60 of chymotrypsin), which restricts substrate P2 utilization to small hydrophobic residues, and forms part of the “apolar binding site” (3).

**Figure 1. The active-site of thrombin.** Schematic diagram showing the major interactions of the fibrinogen A $\alpha$  chain with thrombin and highlighting extensive area involved in substrate recognition. Reprinted from *Trends Biochem. Sci.* **20**, The clot thickens: clues provided by thrombin structure, 23–28, Copyright (1995), with permission from Elsevier Science.



Substrate recognition by thrombin also involves the use of “exosites” to an unusually high degree.

These are regions of the active-site cleft far removed from the catalytic residues, the most significant of which is the “fibrinogen recognition exosite” (4). This is a cationic region to the “east” of the active-site (in the “standard orientation” shown in Figure 1). It is involved in catalytic interactions with fibrinogen and the thrombin receptor and non-catalytic interactions with thrombomodulin and the inhibitor hirudin (derived from the blood sucking leech *Hirudo medicinalis*). The latter is of particular interest as both the remarkable specificity and affinity ( $K_d \sim 10$  fM) of this inhibitor are primarily dictated by this exosite interaction. The other major exosite is the heparin-binding site lying to the north-west of the active site (not shown in Figure 1), which is involved in the physiological inhibitory reaction with the serpin antithrombin bound to heparin or other sulfated glycosaminoglycans (5).

Many of the principles of substrate recognition elucidated in the extensive studies of thrombin will also apply to other serine proteinases with restricted substrate specificity; high resolution structures of them are also being solved.

## Bibliography

“Thrombin” in , Vol. 4, pp. 2544–2546, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ; “Thrombin” in (online), posting date: January 15, 2002, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ.

1. T.K. Vu, D.T. Hung, V.I. Wheaton and S.R. Coughlin, (1991) *Cell*, **64**, 1057–1068.
2. S. Krishnaswamy, K.G. Mann and M.E. Nesheim, (1986) *J. Biol. Chem.* **621**, 8977–8983.
3. W. Bode, I. Mayr, U. Baumann, R. Huber, S.R. Stone and J. Hofsteenge, (1989) *EMBO J.* **8**, 3467–3475.
4. M.T. Stubbs, H. Oschkinat, I. Mayr, R. Huber, H. Angliker, S.R. Stone and W. Bode, (1992) *Eur. J. Biochem.* **206**, 187–195.
5. Z.R. Gan, Y. Li, Z. Chen, S.D. Lewis and J.A. Shafer, (1994) *J. Biol. Chem.* **269**, 1301–1305.

## Suggestions for Further Reading

6. S.R. Coughlin (2000) Thrombin signalling and protease-activated receptors (Review). *Nature* **407**, 258–264.
7. E. Di Cera and A.M. Cantwell (2001) Determinants of thrombin specificity. *Ann. N Y Acad. Sci.* **936**, 133–146.
8. M.T. Stubbs and W. Bode (1995) The clot thickens: clues provided by thrombin structure. *Trends Biochem Sci.* **20**, 23–28.

## Thymidylate Synthase

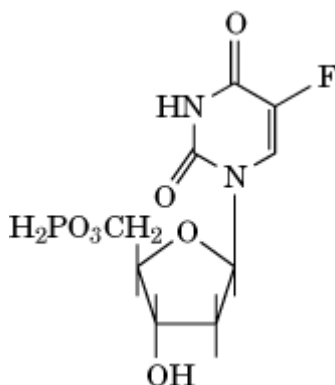
Thymidylate synthase (TS) shares with [ribonucleotide reductase](#) the distinction of being responsible for one of the two chemical differences between DNA and RNA. Just as ribonucleotide reductase carries out the production of deoxyribose at the nucleoside di- or triphosphate level, TS generates the methyl group of thymine, acting at the deoxyribonucleoside 5'-monophosphate level:



This is the only known biological [methylation](#) reaction not involving [S-adenosylmethionine](#) or the synthesis of methionine itself. Aside from this distinctive role, the enzyme is of interest from several standpoints—(1) mechanistically, (2) as the first **enzyme** identified as an anticancer drug target, (3) as one of the most active current targets for computer-assisted drug design, (4) as the first **prokaryotic** protein shown to result from processing of a split gene, (5) as an agent in translational autoregulation, and (6) as a component of [multifunctional proteins](#) and multienzyme complexes.

### 1. A Target for Fluorinated Pyrimidines

In 1957, Friedkin and Kornberg (1) described a tetrahydrofolate-dependent enzyme that converts deoxyuridine monophosphate (dUMP) to thymidine monophosphate (dTMP) (Eq. (1)). The methyl group of dTMP was shown to originate as the methylene carbon of 5,10-methylenetetrahydrofolate ( $\text{CH}_2 = \text{THF}$ ). This discovery of thymidylate synthase was intertwined with the early work of Heidelberger, who synthesized 5-fluorouracil and [5-fluorodeoxyuridine](#) as potential anticancer drugs (2). Heidelberger had noted that tumor cells metabolize uracil much more rapidly than normal cells do, and he expected that substitution of fluorine for a similarly sized hydrogen atom might inhibit the cellular uptake and usage of uracil selectively in tumor cells. In 1958, the action of both fluorinated pyrimidines was shown by Cohen et al (3) to result from their conversion *in vivo* to the deoxyribonucleoside monophosphate 5-fluorodeoxyuridylate (FdUMP), which they



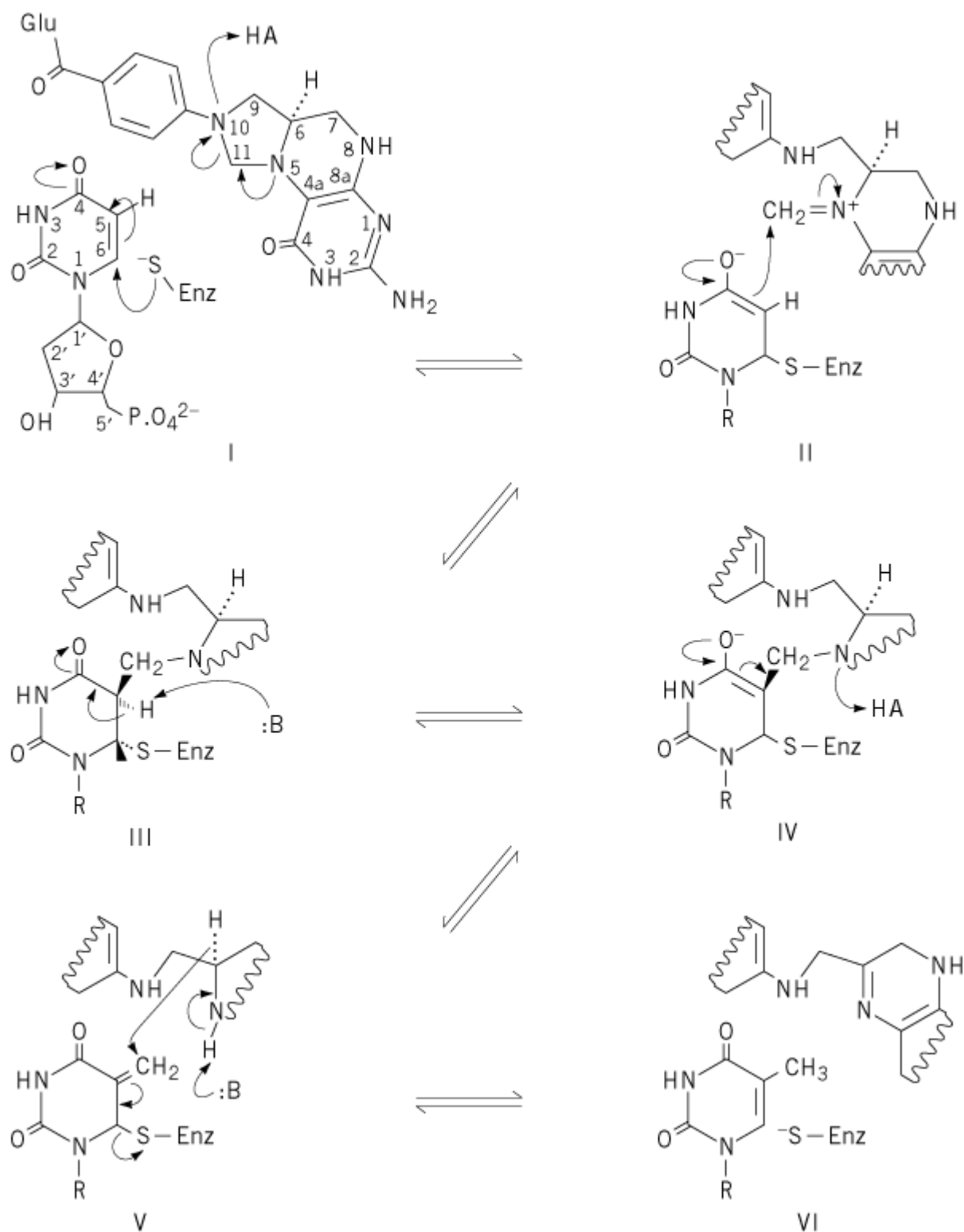
showed to act as a potent inhibitor of TS. Fluorine has a **van der Waals radius** close to that of hydrogen, for which it is substituted in the analogues. This substitution creates a kinetically irreversible inhibition, which eventually was shown to result from formation of a covalently bonded enzyme-inhibitor complex. This is a ternary complex, whose formation requires  $\text{CH}_2 = \text{THF}$  as well. These factors suggested that the structure of this complex would yield clues to the mechanism of action of the enzyme.

Another important early development was the demonstration by Friedkin and colleagues that the folate coenzyme serves as the agent, not only for transfer of a single-carbon functional group but also for reduction of that group from the methylene to the methyl level (4). [Radioisotope](#) labeling studies showed that both the methylene group and a hydride ion—the hydrogen linked to C6 of  $\text{CH}_2 = \text{THF}$ —were transferred to dTMP essentially without dilution. This is the only known reaction in which tetrahydrofolate serves as a **redox** cofactor. These findings suggested a bridged intermediate in which the methylene carbon of  $\text{CH}_2 = \text{THF}$  is linked transiently to C5 of the pyrimidine ring. Inhibition by FdUMP was postulated to result from formation of a similar complex, with the fluorine atom acting to prevent cleavage of the N5-methylene carbon bond.

But how does FdUMP bind covalently to the enzyme? Degradation of the ternary complex revealed a [cysteine](#) sulfur atom linked to C6 of the pyrimidine ring (reviewed in 5). These findings suggested

a mechanism for the reaction in which the cysteine [thiol group](#) initiates nucleophilic attack on C6, converting C5 into a nucleophile that attacks C11, the methylene carbon of 5,10-methylenetetrahydrofolate. Much evidence now supports the mechanism that is outlined in Figure 1.

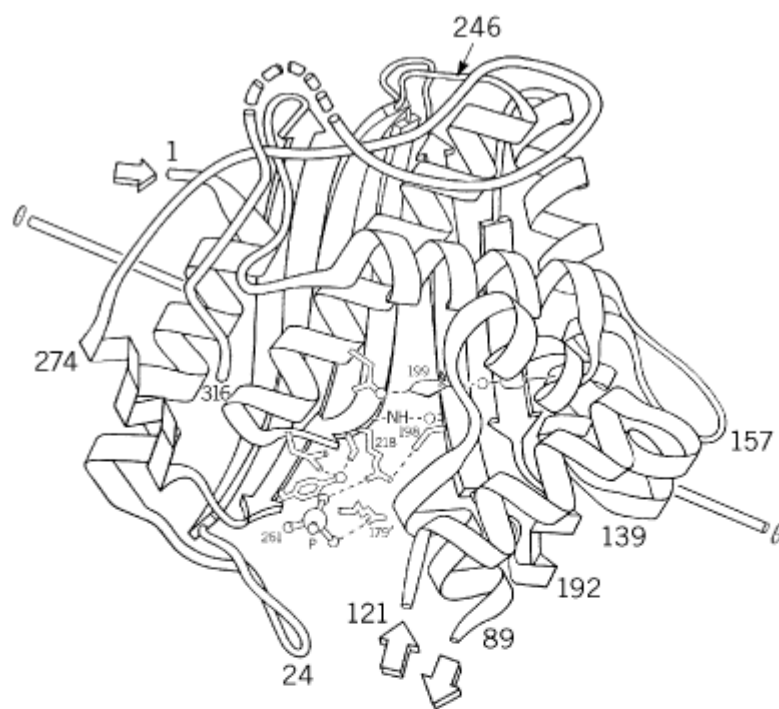
**Figure 1.** A plausible mechanism for the thymidylate synthase reaction, based on radiolabeling studies and the structure of the FdUMP—CH<sub>2</sub> = THF—TS ternary complex. In this complex, the C—F bond cannot be broken, which blocks the reaction at the stage of conversion of III to IV. Reproduced with permission from D. C. Hyatt, F. Maley, and W. R. Montfort (1997) *Biochemistry* **36**, 4585–4594.



## 2. Structure of Thymidylate Synthase

TS from nearly all sources investigated is a homodimeric protein with a protomer size of 30 to 40 kDa. The enzyme is one of the most highly conserved proteins known, with some 18% of its residues invariant among the two dozen known sequences. The first [X-ray crystallography](#) structure reported was of the *Lactobacillus casei* TS (6). The structure shows two substrate-binding clefts, each comprising residues from both polypeptide chains. One such cleft is shown in Figure 2a, which also depicts conserved residues found within the cleft. Cys198 is the initiating nucleophile, and the nearby Arg218 is thought to lower the  $pK_a$  of Cys198, converting the thiol group to the reactive thiolate ion. Arg179', from the other polypeptide chain, is thought to interact with the phosphate group on the substrate. Several nearby [lysine](#) residues (not shown) are thought to interact with the polyglutamate tail of the folate cofactor, which explains the observation that thymidylate synthase binds folate polyglutamates about 100-fold more tightly than the monoglutamate.

**Figure 2.** The crystal structure of thymidylate synthase. (a) One subunit of the *L. casei* TS dimer, showing (with numbers) conserved residues in the active site. Reproduced with permission from L. W. Hardy, J. S. Finer-Moore, W. R. Montfort, M. O. Jones, D. V. Santi, and R. M. Stroud (1987) *Science* **235**, 448–455. (b) One subunit of the ternary complex between *E. coli* TS, FdUMP, and  $CH_2 = THF$ . The ligands are bold, with the folate cofactor above and dUMP below. Reproduced with permission from D. C. Hyatt, F. Maley, and W. R. Montfort (1997) *Biochemistry* **36**, 4585–4594.



(a)

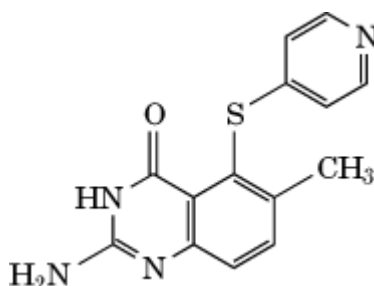


(b)

Several other crystal structures are now known for thymidylate synthase. Figure 2b shows the structure of the FdUMP—CH<sub>2</sub> = THF—enzyme ternary complex for one subunit of *Escherichia coli* thymidylate synthase. Detailed analysis of this structure (7) and others confirmed the interactions postulated from studies of the free enzyme, and it also helped to resolve some stereochemical anomalies of the reaction that had been revealed by radioisotope labeling experiments.

The availability of several TS structures, including that of the human enzyme, plus the attractiveness

of TS as a chemotherapeutic target, has made this enzyme one of the most active subjects of computer-assisted drug discovery research; compounds are being sought that bind as tightly as FdUMP but with greater selectivity and/or lower toxicity. In one approach (8), a docking program was used to identify known compounds that bind tightly in the TS active site. Phenolphthalein, a compound with no evident similarity to substrate or cofactor, was found to bind with an affinity in the micromolar range. Another approach (9) used computational analysis of the binding site to predict structures that would bind tightly in the cofactor site. One such compound, shown below, has no evident structural relationship to methylenetetrahydrofolate,



but it binds the enzyme with an inhibition constant ( $K_i$ ) of around 30 nM.

Although the thymidylate synthase molecule contains two structurally identical substrate binding sites, considerable evidence indicates that TS displays **half-of-the-sites reactivity**, ie, that just one of the two sites is catalytically active. Consistent with this, Maley et al (10) found that simply mixing two inactive mutant forms of *E. coli* TS led to restoration of full wild-type activity, suggesting that mutant subunits had exchanged. Because one of the mutations affected the critical active-site cysteine residue (Cys146 in *E. coli*), only one functional active site per enzyme molecule was evidently enough to give full activity.

### 3. Thymidylate Synthase and Prokaryotic Introns

In bacteriophage T4, a virus-specific TS is encoded by gene *td*. In 1984, biologists were startled to learn that split genes were not confined to the eukaryotic kingdom, when Chu et al (11) showed that the T4 *td* gene contains a 1-kbp **intron**. Further analysis showed this and a handful of other prokaryotic introns to be of the Group I, self-splicing type. These and subsequent findings (12) have contributed to a lively debate over whether **introns** are evolutionarily ancient or recent.

### 4. Thymidylate Synthase and Translational Autoregulation

In eukaryotic cells, thymidylate synthase gene expression is closely aligned with the cell cycle, a fact that has led to detailed investigations of TS **promoter** structure and gene regulation. What is most novel about TS gene regulation, however, is that it is one of the first eukaryotic genes shown to be autoregulated at the level of **translation** (see **Translation Repressors**). TS protein has been shown *in vitro* to bind specifically to its own mRNA and to prevent translation of that mRNA (13). The ability of TS protein to inhibit its own synthesis was blocked by TS substrates or by FdUMP. The physiological significance of this autoregulation is not yet known.

### 5. Multienzyme Complexes and Multifunctional Proteins

Thymidylate synthase is associated structurally with functionally related enzymes. This was first shown in bacteriophage T4, where at least 10 enzymes of DNA precursor biosynthesis interact to form a complex, called dNTP synthetase, which evidently facilitates the flow of DNA precursors into DNA. Within this complex, an especially strong interaction links thymidylate synthase with the phage-coded **deoxycytidylate deaminase**, the source of most of the dUMP that is used by TS (14). Similar complexes have been described in eukaryotic cells although a role in coupling dNTP synthesis to DNA replication has not yet been established. Surprisingly, of those TS molecules that



have been shown to be localized within the nucleus (a prerequisite for the coupling of dNTP synthesis to replication), most are located in the [nucleolus](#) (15), a finding that suggests possible additional metabolic roles for the TS protein.

Quite different is the situation in protozoa and some **plants**, where TS exists as a bifunctional protein also containing [dihydrofolate reductase](#) activity. Kinetic analysis of the *Leishmania* DHFR-TS (16) shows that the dihydrofolate released by the TS reaction is channeled to the DHFR active site—transported from site to site without release from the enzyme surface—even though the active sites are 40 Å apart. Elcock et al (16) have ascribed this unusual behavior to the charge distribution on the enzyme surface. Because of its unusual structure and function, the bifunctional DHFR-TS is being actively studied as a target for antiparasitic drugs.

## Bibliography

1. M. Friedkin and A. Kornberg (1957) In W. D. McElroy and B. Glass, eds., *The Chemical Basis of Heredity*, Johns Hopkins Press, Baltimore, MD, pp. 609–620.
2. C. Heidelberger, N. K. Chaudhuri, P. Danenberg, D. Mooren, L. Griesbach, R. Duschinsky, R. J. Schnitzer, E. Plevin, and T. Scheiner (1957) *Nature* **179**, 663–666.
3. S. S. Cohen, J. G. Flaks, H. D. Barner, M. R. Loeb, and J. Lichtenstein (1958) *Proc. Natl. Acad. Sci. USA* **44**, 1004–1012.
4. M. Friedkin (1959) *Federation Proc.* **18**, 230.
5. P. V. Danenberg (1977) *Biochim. Biophys. Acta* **473**, 73–92.
6. L. W. Hardy, J. S. Finer-Moore, W. R. Montfort, M. O. Jones, D. V. Santi, and R. M. Stroud (1987) *Science* **235**, 448–455.
7. D. C. Hyatt, F. Maley, and W. R. Montfort (1997) *Biochemistry* **36**, 4585–4594.
8. B. K. Shoichet, R. M. Stroud, D. V. Santi, I. D. Kuntz, and K. M. Perry (1993) *Science* **259**, 1445–1450.
9. K. Appelt et al (191) *J. Med. Chem.* **34**, 1925–1934.
10. F. Maley, J. Pedersen-Lane, and L. Changchien (1995) *Biochemistry* **34**, 1469–1474.
11. F. K. Chu, G. F. Maley, F. Maley, and M. Belfort (1984) *Proc. Natl. Acad. Sci. USA* **81**, 3049–3053.
12. M. Belfort, M. E. Reaban, T. Coetzee, and J. Z. Dalgaard (1995) *J. Bacteriol.* **177**, 3987–3903.
13. E. Chu, D. M. Koeller, J. L. Casey, J. C. Drake, B. A. Chabner, P. C. Elwood, S. Zinn, and C. J. Allegra (1991) *Proc. Natl. Acad. Sci. USA* **88**, 8977–8981.
14. K. M. McGaughey, L. J. Wheeler, J. T. Moore, G. F. Maley, F. Maley, and C. K. Mathews (1996) *J. Biol. Chem.* **271**, 23037–23042.
15. W. A. Samsonoff, J. Reston, M. McKee, B. O'Connor, J. Galivan, G. Maley, and F. Maley (1997) *J. Biol. Chem.* **272**, 13281–13285.
16. A. H. Elcock, M. J. Potter, D. A. Matthews, D. R. Knighton, and J. A. McCammon (1996) *J. Mol. Biol.* **262**, 370–374.

## Suggestions for Further Reading

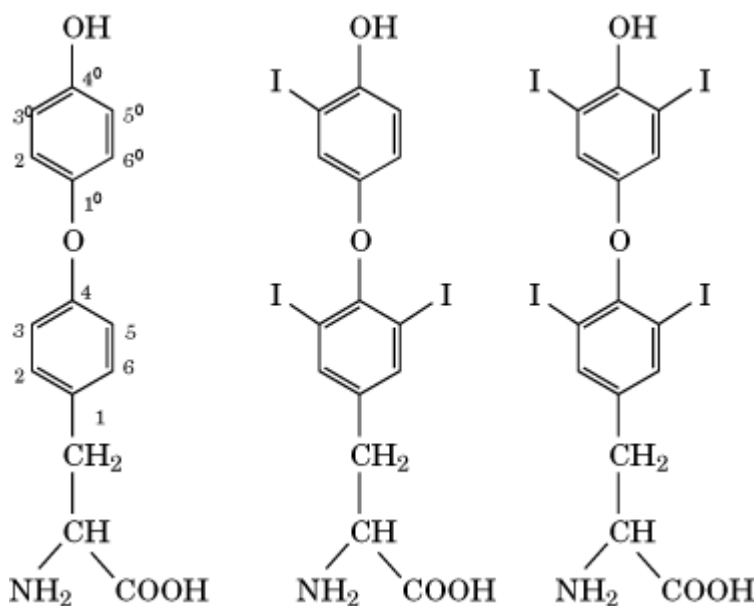
17. P. V. Danenberg (1977) Thymidylate synthase, a target for cancer chemotherapy. *Biochim. Biophys. Acta* **473**, 73–92. An excellent review of the older literature.
18. R. J. Cisneros, L. A. Silks, and R. B. Dunlap (1988) Mechanistic aspects of thymidylate synthase: molecular basis for drug design. *Drugs of the Future* **13**, 859–881. A well-referenced review of TS mechanism and inhibitor binding.
19. K. M. Perry, E. B. Fauman, J. S. Finer-Moore, W. R. Montfort, G. F. Maley, F. Maley, and R. M. Stroud (1990) Plastic adaptation toward mutations in proteins: structural comparison of thymidylate synthases. *Proteins: Structure, Function, Genetics* **8**, 315–333. A detailed discussion of the TS crystal structure.

20. M. P. Costi (1998) Thymidylate synthase inhibition: a structure-based rationale for drug design. *Med. Res. Rev.* **18**, 21–42. A very recent review of thymidylate synthase as a target for structure-based drug design.
21. E. Chu and C. J. Allegra (1996) The role of thymidylate synthase as an RNA binding protein. *BioEssays* **18**, 191–198. A recent review of translational autoregulation of thymidylate synthase.
22. Mathews, C. K. (1993) Enzyme organization in DNA precursor biosynthesis. *Prog in Nucleic Acid Res. Mol. Biol.*, **44**, 167–203. A review of multienzyme complexes in DNA precursor biosynthesis.

## Thyroid Hormones

Thyroid hormones are iodinated derivatives of [tyrosine](#) that are synthesized in the follicular cells of the thyroid gland. The noniodinated thyronine can be regarded as the basic structure of these [hormones](#) (Fig. 1). 3,5,3'-Triiodothyronine (T3) is the major hormonally active derivative, whereas thyroxine (3,5,3',5'-tetraiodothyronine, T4) has one-fifth of its biological activity. Biosynthesis of the hormones involves thyroglobulin, a tyrosine-rich protein stored in the follicle lumen. Iodine is oxidized to I<sub>2</sub> by a [peroxidase](#) of the follicular cells and then coupled covalently to tyrosine residues of thyroglobulin. The side chains of other tyrosine residues of thyroglobulin are subsequently coupled to the iodinated tyrosine residues, with the formation of diphenylether covalent bonds and completion of the basic thyronin structure. In this form, the thyroid hormones are still part of the thyroglobulin molecule. Under the stimulatory influence of thyrotropic hormone (TSH, thyrotropin), by way of a membrane receptor and activation of adenylyl cyclase ([1](#)), the iodinated thyroglobulin is taken up by the thyrocytes and degraded by [proteolysis](#), and the freed T4 and T3 are released into the circulation. The bulk of T4 and T3 in the circulation is protein-bound. Only the free hormones can penetrate cells. In target cells, T4 is 5'-deiodinated to T3, which, at the cellular level, is the active hormone.

**Figure 1.** Structures of thyronine, T3, and T4.



Thyroid hormones exert a broad range of effects on cell growth, development, and metabolism. A classical effect of thyroxine, instrumental in delineating the mode of action of the hormone, is the triggering of amphibian metamorphosis (2). This change from larva to the adult animal is accompanied biochemically by a change in the mode of amino acid nitrogen excretion. Larvae excrete nitrogen in the form of  $\text{NH}_4^+$ , whereas frogs excrete it as urea. Thyroxine induces the biosynthesis of key enzymes involved in the production of urea by molecular mechanisms similar to those of [steroid hormones](#). These hormones bind to intracellular **receptors**, belonging to the superfamily of nuclear receptors (3). Two T3 receptor genes have been found, *a* and *b*, located on human chromosomes 17 and 3, respectively. The respective receptors (TRa and TRb) have been identified. TRa is the *c-erbA protooncogene*. Two different isoforms of *b* (TRb1 and TRb2) have been isolated and exhibit tissue-specific expression. The receptors recognize the respective [hormone response elements](#) of DNA, to which they bind either as homodimers or as heterodimers with other nuclear proteins, such as RXR (retinoid X-receptor). In some cases, the receptor binds in the absence of hormone and acts as a [repressor](#). Derepression is accomplished by its ligand binding to the receptor (4). In other cases, the hormone-receptor complex acts as a positive regulator of transcription.

Thyroid hormone deficiency during embryogenesis results in mental and somatic developmental defects (cretinism). Experimental extirpation of the thyroid gland leads to delayed growth and sexual maturation. In the fully developed organism, thyroid hormones show general metabolic effects, increasing metabolic rates,  $\text{O}_2$  consumption, and thermogenesis. These effects result from the induction of the  $\text{Na}^+/\text{K}^+$  ATPase. The synthesis of RNA and proteins are also stimulated. Furthermore, the biosynthesis of other enzymes (eg, the mitochondrial glycerol-1-phosphate dehydrogenase) are also induced by the thyroid hormones. Thyroid hormones stimulate both lipogenesis and lipolysis. Lipogenesis is the result of induction of biosynthesis of enzymes of the lipogenic pathway (malate dehydrogenase, glucose 6-phosphate dehydrogenase, and fatty acid synthetase). In hypothyroidism, **lipoprotein** metabolism is altered, serum concentrations of LDL-cholesterol are increased, and hepatic [lipase](#) activity is decreased. Thyroid hormones reduce systemic vascular resistance, enhance cardiac contractility, and have a chronotropic effect. T3 stimulates transcription of myosin heavy-chain<sub>a</sub> and inhibits expression of heavy-chain b-genes.

#### Bibliography

1. G. A. Brent (1994) N. Engl. J. Med. **331**, 847–853.

2. J. R. Tata (1993) *BioEssays* **15**, 239–248.
3. D. J. Mangelsdorf et al. (1995) *Cell* **83**, 835–839.
4. F. J. Piedrafita et al. (1995) *Mol. Endocrinol.* **9**, 563–578.

### Suggestions for Further Reading

5. M. A. Lazar (1993) Thyroid hormone receptors: multiple forms, multiple possibilities. *Endocr. Rev.* **14**, 184–193.
6. A. Munoz and J. Bernal (1997) Biological activities of thyroid hormone receptors. *Eur. J. Endocrinol.* **137**, 433–445.

### Ti Plasmid

The Ti plasmid of *Agrobacterium tumefaciens* is responsible for producing tumors on infected dicotyledonous **plants**. During the infection process, a defined region of **DNA**, the transferred or T-DNA, is transferred from the **bacteria** to the infected plant cell and integrated into the plant **genome** (see **T-Complex, -DNA, -Region, -Strand**). This natural form of genetic **transformation** is the basis of **vectors** used in **plant genetic engineering** for producing **transgenic** plants.

The Ti plasmid is a megaplasmid (approximately 220 kbp) present in a single copy in *A. tumefaciens*. Three types of Ti plasmid have been defined, octopine, nopaline, or agropine, depending on which **opine** is synthesized by the plant cells transformed by the corresponding T-DNA. Genetically the functions of the Ti plasmid can be defined as (1) oncogenicity, (2) nature of opines synthesized, (3) utilization of opines, (4) conjugative transfer of the Ti plasmid, (5) sensitivity to the **bacteriocin** agrocin 84, (6) host range, and (7) exclusion of **bacteriophage** AP-1 (1). Octopine and nopaline Ti plasmids share four regions of significant homology: the T-DNA; the virulence (*vir*) region, containing genes encoding proteins responsible for transfer of the T-DNA to the plant cell; the **origin of replication**; and a region responsible for conjugative transfer of the plasmid (2). The Ti plasmid is not stable in bacteria cultured above 30°C.

Because of their large size, the use of Ti plasmids in constructing plant transformation vectors was initially problematic. This was overcome, however, by creating **cointegrative vectors**. These are based on Ti plasmids into which foreign DNA to be transferred to the plant cell can be engineered by recombination.

### Bibliography

1. K. Kersters and J. De Ley (1984) "Genus III *Agrobacterium* Conn 1943", In *Bergey's Manual of Systematic Bacteriology*, 1 Vol 1 (N. R. Krieg and J. G. Holt, eds.), Williams and Wilkins, Baltimore, pp. 244–254.
2. A. Depicker, M. Van Montagu, and J. Schell (1978) *Nature* **275**, 150–153.

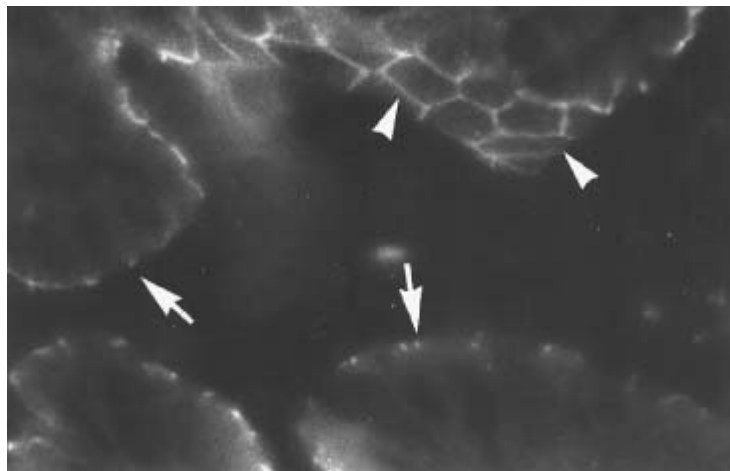
### Suggestion for Further Reading

3. G. Kahl and J. Schell (1982) *Molecular Biology of Plant Tumors*, Academic Press, London.

## Tight Junction

The tight or occludens junction is a [cell junction](#) that occludes or regulates the permeability of paracellular channels. These junctions are widespread in simple epithelia and some endothelia, cell layers that function to divide physiological compartments within the body—for example, the epithelium of the intestine and the endothelium of the blood–brain barrier. In transmission [electron microscopy](#), tight junctions appear as punctate areas where adjacent plasma membranes approach each other so closely that their outer leaflets seem to touch or fuse. By [freeze fracture](#), they appear as continuous strands of particles located in the outer or P-face of the plasma membrane, with complementary grooves located on the inner or E-face. The strands may be single or may form a continuous anastomosing network (1). The junctions are zonular in nature, encircling the entire perimeters of cells within an epithelium or endothelium—hence their alternative name, zonula occludens (Fig. 1).

**Figure 1.** Tight junctions in small intestinal epithelium stained with antibody to ZO-1. Vertical sections through the cells show the concentration of staining at the extreme apicolateral borders (arrows), and transverse sections show the zonular nature of the junctions (arrowheads). (From Ref. 29, with permission from Chapman and Hall.)

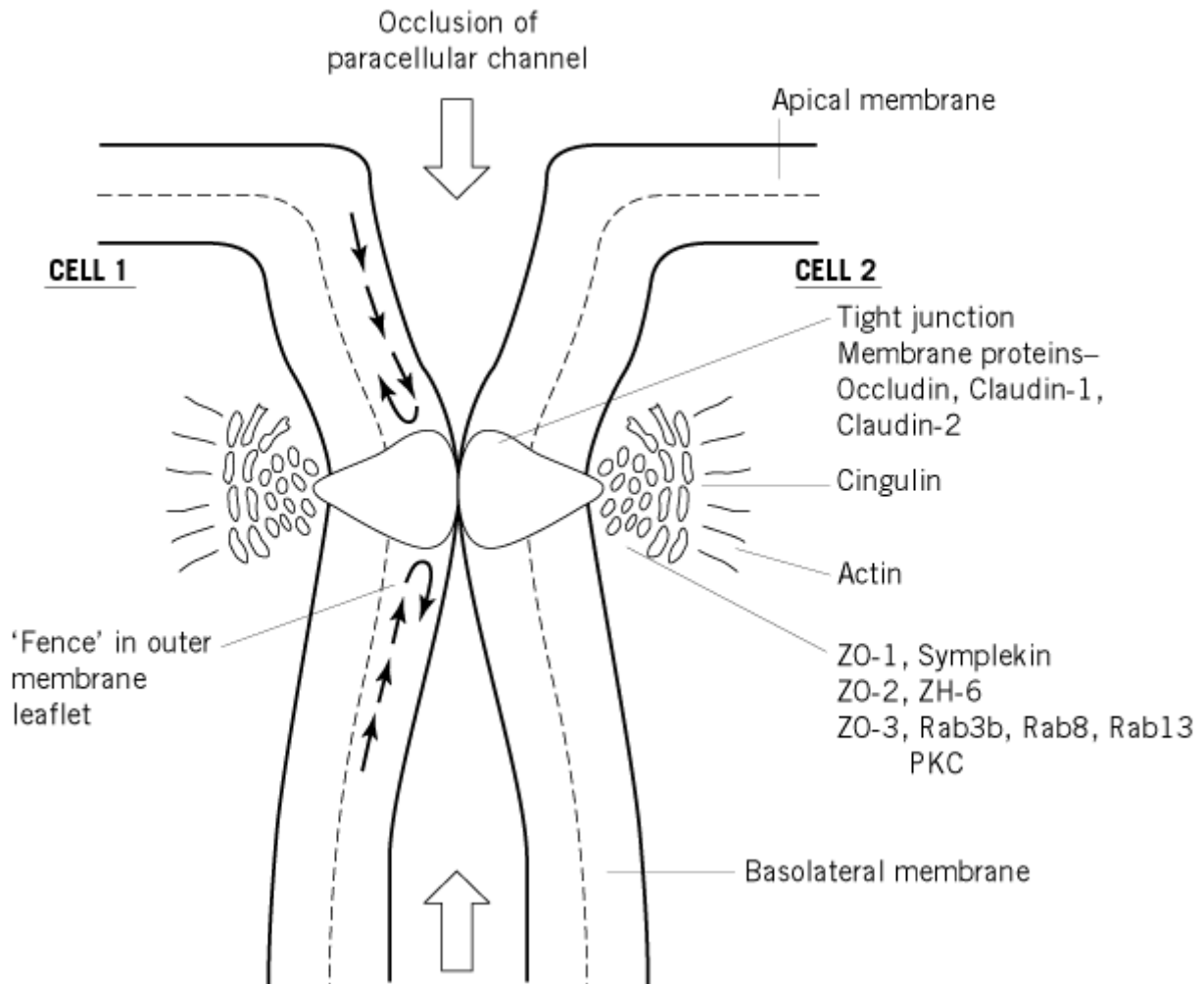


Tight junctions are largely impermeable to molecules with radii greater than 1.5 nm but are selectively permeable to ions (1, 2). Junctional permeability may be measured by determining the flux of cell-impermeable radioactive or labeled molecules across an epithelium. Restriction of the passage of ions results in an electrical resistance between one side of an epithelium and the other, called the transepithelial resistance (TER). This may be measured as an indication of tight junction function. Junctions that give rise to low TER are said to be “leaky,” and junctions of high resistance are said to be “tight.” Extremes of TER lie at 6 to 7 W/cm<sup>2</sup> for rat proximal tubule and 6×10<sup>4</sup> to 3×10<sup>5</sup>W/cm<sup>2</sup> for rat urinary bladder (2). Generally, the more strands seen in a junction by freeze fracture and electron microscopy, the higher the TER (1, 3).

Tight junctions also make an important contribution to the regulation of epithelial cell polarity (4). In simple epithelia, the protein and phospholipid compositions of the apical and basolateral plasma membranes differ. Molecules that are not anchored to the [cytoskeleton](#) are normally freely diffusible within the plane of the plasma membrane. Tight junctions have been shown to restrict the diffusion of molecules in the outer leaflet of the plasma membrane, thus confining them to either the apical or the basolateral domain. The regulation of paracellular channels is sometimes referred to as the “gate”

function of tight junctions, and the restriction of diffusion of membrane molecules is sometimes called the “fence” function (Fig. 2).

**Figure 2.** Diagrammatic representation of the structure of the tight junction, indicating its principal functions: (a) the “gate” function regulating paracellular permeability and (b) the “fence” function restricting diffusion of membrane components between the different membrane domains of the cell. (From Ref. 29, with permission from Chapman and Hall.)



Several molecular components of tight junctions have been identified. A major integral membrane component is the 60-kDa protein occludin, which has four **transmembrane** domains (5). This is directly involved in the formation of tight junction strands (6) and forms tight junction-like vesicular structures when overexpressed in insect Sf9 cells (7). Occludin also raises TER when expressed in epithelial cells (8), decreases TER when expressed in a C-terminally truncated form (9), and shows **cell adhesion** activity (9). Peptides corresponding to the second extracellular loop of occludin cause tight junction disassembly when applied to cells (10). However, ES cells in which the occludin gene has been deleted form good tight junction strands (11), indicating the existence of other integral **membrane proteins**. This led to the discovery of two more four transmembrane domain proteins, claudin-1 and claudin-2 (22.9 and 24.5 kDa), that localize to junctional strands (12).

The cytoplasmic components of tight junctions include three members of a larger family of proteins called membrane-associated guanylate kinase (MAGUK) proteins. These are ZO-1, ZO-2, and ZO-3 (13-15), with ZO-1 being the first discovered tight junction protein. These proteins characteristically

have DHR (discs-large homology region), **SH3** (src homology 3), and GuK (guanylate kinase) **domains** and are suggested to have membrane receptor clustering, cytoskeletal interaction, and signaling functions (16). Cingulin is a tight junction protein with a [coiled-coil](#) rod domain, showing some [homology](#) to the myosin heavy chain and suggesting the possibility of dimerization (17). It may form a filamentous network and mediate cytoskeletal associations. A number of potential regulatory proteins associated with tight junctions are three small [GTP-binding proteins](#) [rab3B, rab8, and rab13 (18, 19)], G-proteins, and [protein kinase C](#) (20). Tight junctions also contain a number of other proteins of unknown function, including 7H6 (21) and symplekin (22).

Numerous factors have been found to modulate tight junction permeability. Many may function through an effect on the [actin](#) cytoskeleton, which is associated either with the tight junction directly or the closely associated zonula adherens (see [Intermediate Junction](#)) (23). **Phosphorylation** of tight junction components, such as ZO-1 and cingulin, is also seen as important in regulation of both permeability and assembly (24, 25).

Maintenance of epithelial permeability barriers is vital to normal physiological function; and compromise of barrier function may be found in a variety of diseases, including inflammatory bowel disease and asthma (26, 27). Tight junction proteins are expressed very early in normal development, where they appear to be involved in polarity determination in the first epithelium, the trophectoderm (28).

#### Bibliography

1. J. L. Madara and K. Dharmasathaphorn (1985) *J. Cell Biol.* **101**, 2124–2136.
2. D. W. Powell (1981) *Am. J. Physiol.* **250**, G275–G288.
3. P. Claude (1978) *J. Membr. Biol.* **39**, 219–232.
4. K. Simons (1990) In *Morphoregulatory Molecules* (G. M. Edelman, B. A. Cunningham, and J-P. Thiery, eds.), Wiley, New York, pp. 341–356.
5. M. Furuse et al. (1993) *J. Cell Biol.* **123**, 1777–1788.
6. K. Fujimoto (1995) *J. Cell Sci.* **108**, 3443–3449.
7. M. Furuse et al. (1996) *J. Cell Sci.* **109**, 429–435.
8. M. S. Balda et al. (1996) *J. Cell Biol.* **134**, 1031–1049.
9. C. M. Van Itallie and J. M. Anderson (1997) *J. Cell Sci.* **110**, 1113–1121.
10. V. Wong and B. M. Gumbiner (1997) *J. Cell Biol.* **136**, 399–409.
11. M. Saitou et al. (1998) *J. Cell Biol.* **141**, 397–408.
12. M. Furuse et al. (1998) *J. Cell Biol.* **141**, 1539–1550.
13. B. R. Stevenson (1986) *J. Cell Biol.* **103**, 755–766.
14. B. Gumbiner et al. (1991) *Proc. Natl. Acad. Sci. USA.* **88**, 3460–3464.
15. J. Haskins et al. (1998) *J. Cell Biol.* **141**, 199–208.
16. D. F. Woods and P. J. Bryant (1993) *Mech. Dev.* **44**, 85–89.
17. S. Citi et al. (1990) *J. Cell Biol.* **111**, 409a.
18. E. Weber et al. (1994) *J. Cell Biol.* **125**, 583–594.
19. A. Zahraoui et al. (1994) *J. Cell Biol.* **124**, 101–115.
20. V. Dodane and B. Kachar (1996) *J. Membr. Biol.* **149**, 199–209.
21. Y. Zhong et al. (1993) *J. Cell Biol.* **120**, 477–483.
22. B. H. Keon et al. (1996) *J. Cell Biol.* **134**, 1003–1018.
23. J. L. Madara (1989) *J. Clin. Invest.* **83**, 1089–1094.
24. B. Stevenson (1989) *Biochem. J.* **263**, 597–599.
25. S. Citi and N. Denisenko (1995) *J. Cell Sci.* **108**, 2917–2926.
26. L. Lora et al. (1997) *Gastroenterology* **113**, 1347–1354.

27. H. Wan et al. (1998)
28. T. P. Fleming et al. (1993) *Development* **117**, 1135–1144.
29. D. R. Garrod and J. E. Collins (1992) "Intercellular Junctions and Cell Adhesion in Epithelial Cells". In *Epithelial Organization and Development* (T. P. Fleming, ed.), Chapman and Hall, London, pp. 1–52.

### Suggestions for Further Reading

30. S. Citi and N. Cordenosi (1999) "The Molecular Basis for the Structure, Function and Regulation of Tight Junctions". In *Adhesive Interactions of Cells* (D. R. Garrod, M. A. J. Chidgey, A. J. North, eds.), JAI Press, Greenwich, CT, pp. 202–232. (An excellent review covering all aspects of tight junctions.)
31. J. M. Anderson and C. M. Van Itallie (1995) Tight junctions and the molecular basis for the regulation of paracellular permeability. *Am. J. Physiol.* **269**, 467–475. (Another excellent review.)

## TIM Barrel

The TIM barrel is named after the [protein](#) triose phosphate isomerase, or TIM, in which this type of **domain** of a [protein structure](#) was first observed. The TIM barrel has a central cylinder or barrel of [b-sheet](#) formed from eight parallel [b-strands](#), with strands 1 and 8 hydrogen bonded to each other to close the b-barrel. Each b-strand is connected to the next by a linking [a-helix](#); the eight a-helices form a concentric layer surrounding the central parallel b-barrel. Both the strands and helices of the TIM barrel have a pronounced right-hand twist (Fig. [1](#)). The topology of the TIM barrel gives rise to its other names, the (ba)<sub>8</sub> or (a/b)<sub>8</sub> barrel.

**Figure 1.** A TIM barrel. **(Top)** Schematic representation of the topology of the TIM barrel, with b-strands depicted as purple arrows and a-helices depicted as green cylinders. The connecting regions between the helices and cylinders are shown in yellow, and the *N*- and *C*-terminal ends of the motif are labeled. **(Bottom)** Schematic representation of the backbone of triose phosphate isomerase ([1](#)) showing a typical TIM barrel structure. The center of the barrel is formed from the b-strands forming a parallel b-sheet (shown as purple arrows), and the outer layer is formed from a-helices (shown as green coils) that connect the b-strands. The connecting regions between helices and strands are shown in yellow. This figure was generated using Molscript ([2](#)) and Raster3D ([3, 4](#)). See color insert.





This very common fold is found in many proteins with diverse functions and no detectable sequence identity, so their similar structures are thought to be an example of [convergent evolution](#). Although the activities and sequences of these proteins vary, their [active sites](#) are all formed by loop regions at the carboxyl ends of the b-strands that connect to the a-helices. The particular enzymatic activity of the [enzyme](#) is dependent on the amino acid residues in these regions. The TIM barrel forms one of two major classes of a/b domains, and the other class is typified by the [nucleotide-binding motif](#).

[See [Domain, Protein](#) and compare with [Antiparallel beta-barrel motifs](#).]

#### Bibliography

1. E. Lolis et al. (1990) *Biochemistry* **29**, 6609–6618.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.
3. E. A. Merritt and M. E. P. Murphy (1994) *Acta Crystallogr.* **D50**, 869–873.
4. D. J. Bacon and W. F. Anderson (1988) *J. Mol. Graphics* **6**, 219–222.

#### Suggestions for Further Reading

5. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
6. J. S. Richardson (1981) The anatomy and taxonomy of protein structure. *Adv. Protein Chem.* **34**, 167–339.

#### Tissue Culture

Tissue culture was originally understood to mean the maintenance and, in some cases, growth *in vitro* of small explanted fragments of tissue from an animal or plant, (1, 2), but it has become a generic term covering a wide range of *in vitro* cultivation techniques. Cell culture, where the tissue is dissociated into individual cells, which may be propagated adherent to glass or plastic, or in suspension, is the technique most widely used. Organ culture, where the tissue is not disaggregated, but maintained at the air–liquid interface as a complete fragment, is a method that retains the architecture of the tissue, but does not allow propagation. Maintenance of small fragments of tissue at the solid–liquid interface, that is, attached to the glass or plastic substrate (primary explant culture), permits the outgrowth of cells on to the substrate by migration and cell proliferation, and can give rise to a propagated [cell line](#).

Tissue culture, defined in its generic sense, is used when it is important to regulate the physical and physiological microenvironment of the cells, maintain a consistent, validated stock of characterized cells, and/or minimize the use of laboratory animals. It also allows the introduction of exogenous genomic DNA with a high efficiency (see [Transfection](#)), the study of signaling mechanisms under controlled conditions with a purified population of cells, the measurement of cytotoxic and genotoxic responses, and the production of [proteins](#) with appropriate [post-translational modifications](#) (see [Protein Engineering](#)). Tissue culture, however, will not mimic *in vivo* conditions accurately, because of the difference in the cellular microenvironment, the absence of metabolizing enzymes normally found in the liver, and the difficulty of reproducing the fully differentiated cell **phenotype** *in vitro*. Continuous cell lines are also genetically unstable and often genotypically heterogeneous.

## 1. History

Opinions vary as to when tissue culture actually commenced, as there was a gradual transition from the earlier studies of pathologists with isolated fragments of tissue to the demonstration of actual growth (3). Tissue culture implies that the explanted tissue or cells can be maintained for more than a few hours and will display some aspect of growth. The first success in this area is usually attributed to Harrison (1), who was able to maintain spinal cord from frog embryo in a natural medium of clotted frog lymph for up to 4 weeks and observed outgrowth of nerve fibers from the explant. Others were able to repeat Harrison's observations and made attempts to passage the resulting cultures (4). Propagated cell lines (see [Cell Line](#)) were derived from the outgrowth of several types of tissue, but it was not until after World War II that extensive use was made of **trypsin** to passage cells, producing a single-cell suspension that could be divided accurately and even cloned (5).

Numerous attempts were made to generate new cell lines (6), assisted by the development of antibiotics, which were incorporated into culture media to prolong culture life without contamination. Subsequently, the development of laminar flow cabinets in the late 1960s and early 1970s provided a protected environment, facilitating aseptic technique without the necessity of continued culture in antibiotics. Around this same period, it was observed that cells could be stored frozen in culture medium, with a preservative added, for extended periods (7-9), and this gave further protection against loss by contamination or incubator failure. This led to a rapid proliferation of cell lines, and it was not until Gartler (10) initiated studies with enzyme polymorphisms, which could distinguish some cell lines from others by the **isoenzymes** characteristic of the original tissue, that a more serious contamination problem was discovered, namely, cross-contamination. This was confirmed by [karyotype](#) analysis (11), and by the mid-1970s, it was estimated that around 30% of all cell lines in current use in North America were cross-infected with [HeLa Cells](#) (12). Many of these, such as KB, Hep-2, Girardi heart, and Chang liver, are still in current use, often not identified by the user as HeLa-contaminated. This problem has been recognized by the major cell banks, who now label stocks suspected as carrying HeLa markers (13).

The continued use of antibiotics, combined with natural products, such as serum and trypsin, brought in another crisis when Barile (14) and others reported that large numbers of propagated cell lines were contaminated with [mycoplasma](#), too small to detect readily without special staining, and not always as catastrophic to the culture as are most bacterial infections. These microorganisms were too

small to be filtered out by the sterilization procedures of the day and, once contaminating a cell line, were very difficult to remove. The development of a **fluorescent** staining technique (15) and other assays for mycoplasma detection, together with a reduced porosity of sterilization filters used in serum manufacture, means that contamination is readily detected, and further contamination preventable, by constant screening and use of materials from reputable sources. Nevertheless, mycoplasma contamination, introduced from infected cell lines, primary cultures, or the operator, is still widespread and requires constant vigilance to eliminate.

Modern tissue culture, carried out in high grade facilities with ultrapure water and media components and proper validation and authentication procedures, is now a highly reproducible technique employed in academic research laboratories, diagnostic laboratories, and industrial-scale production of biopharmaceuticals.

## 2. Types of Culture

### 2.1. Cell Culture

This is the maintenance and propagation of cells as a uniform monolayer or suspension (16). A cell culture may be derived from a tissue by three main methods of primary culture.

*1. Primary explant technique.* A small fragment of tissue is placed on the base of a culture vessel so that it adheres, either spontaneously, facilitated by scratching the plastic surface, trapping the tissue under a coverslip, or use of surface tension or clotted plasma. A primary culture arises by outgrowth from the explant across the surface of the dish.

*2. Mechanical disaggregation.* The tissue may be disaggregated by chopping with scalpels or scissors or by forcing the tissue through a mesh screen or syringe needle. The resultant suspension of cells and small fragments is allowed to settle and form an adherent cell monolayer on a glass substrate or the correct grade of plastic (see below).

*3. Enzymatic disaggregation.* Tissue may be disaggregated by incubation in a number of different [proteinases](#); trypsin, collagenase, and Dispase are the most common. Trypsin requires the absence of serum, which contains [trypsin inhibitors](#), but collagenase and Dispase are insensitive to [proteinase inhibitors](#) in serum. As for mechanical disaggregation, the resultant cell suspension, washed free of enzyme by [centrifugation](#) and resuspension, is allowed to settle to the base of the dish to form an adherent monolayer (17).

The outgrowth from a primary explant culture, or the monolayer generated from mechanically or enzymatically disaggregated cells, may be subcultured (see [Cell Line](#)) and transferred to fresh culture vessels (17). This process, also known as *passage*, is usually achieved by treating the monolayer with trypsin to dissociate the cells from each other and the substrate, resuspending the cells in fresh medium with serum (to inhibit the residual trypsin), and diluting into fresh culture vessels (17). After the first subculture the culture becomes a cell line (see [Cell Line](#)).

## 3. Organ Culture

Organ cultures are not as widely used as cell cultures, as they require considerably more effort to initiate and cannot subsequently be propagated. They do, however, provide a means to retain at least some of the histological structure of the tissue, and with it some of its phenotypic characteristics. Organ cultures can be maintained for up to 3 weeks (although longer periods have been recorded) and may contain areas of localized cell division, but growth is limited. Tissue, or preferably a whole organ, from the embryo survives better, and may show some net growth in culture, but the growth and cell survival are limited by the dimensions of the tissue explant. Under normal conditions (20% oxygen at atmospheric pressure), the radius of a spherical organ culture is limited to a maximum of 500  $\mu\text{m}$ . The organ culture may extend in shape in two dimensions, as long as one dimension

remains at 500  $\mu\text{m}$  or less. The inward diffusion of oxygen, and the outward diffusion of  $\text{CO}_2$ , are optimized by growing the tissue at the air–liquid interface; access to the nutrient medium is provided by a permeable support, usually a porous membrane on a stainless steel mesh grid. Organ cultures have been found particularly useful for cultivation of skin (18), fetal bone (20), and various embryonic organs during organogenesis (21).

#### 4. Histotypic Culture

It is possible to create the cell–cell interactions provided by tissue-like densities, while still having the convenience and reproducibility of cell lines, by growing cell cultures to high cell densities. This can be achieved in several ways: (1) By growing cells in a filter well insert, where the cells are crowded, but adequate medium is provided by access to a relatively large reservoir (22); (2) growing cells as stirred aggregates, called *spheroids*; generated by growing the cells at high concentrations on agar or [agarose](#) in a multiwell plate, so that the cells form aggregates in the bottom of the meniscus generated by the agar in the well—these are then stirred at low speed in suspension; similar aggregates may be generated by growing cells in suspension in zero gravity in a slowly rotating chamber (23); and (3) Growing cells on the outer surface of perfused microcapillary bundles, where the cells are seeded in the outer chamber holding the hollow fibers, and medium is pumped through the fibers from a reservoir (24).

#### 5. Organotypic Culture

Histotypic culture can be adapted to cocultivation of two, or more, different cell types, in an attempt to simulate heterotypic cell interactions, in addition to the homotypic cell interactions achieved in histotypic culture. Skin culture has become a classic example, where epidermal epithelium maintained in coculture with dermal fibroblasts embedded in [collagen](#) will generate well-differentiated keratinocytes expressing involucrin, filaggrin, and cross-linked [keratin](#) (25). The requirement for collagen is similar to many specialized cultures, which require elements of the [extracellular matrix](#), such as collagen, [laminin](#), and [fibronectin](#), for proliferation and differentiation.

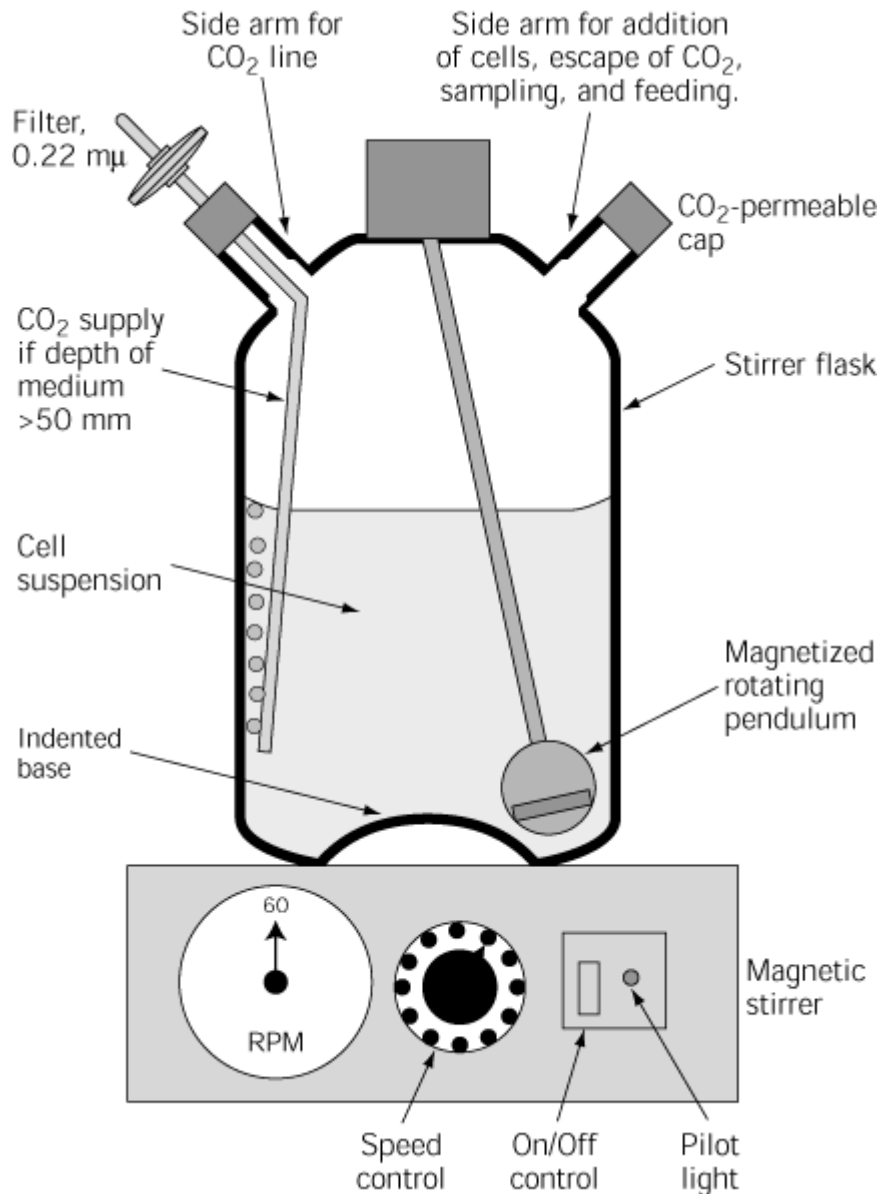
#### 6. Monolayer and Suspension Cultures

Most normal cells attach to plastic and require to spread out on the plastic before cell division will commence (see [Contact Inhibition](#)). Attachment is mediated by transmembrane [integrin](#) receptors (see [Cell Adhesion Molecules](#)), which bind components of the extracellular matrix deposited on the substrate by the cells. Hence, coating the substrate with matrix molecules, such as fibronectin or collagen (26), or conditioning the substrate by previously culturing cells on it and removing them with [detergent](#) in water (27), can often facilitate cell attachment and subsequent proliferation.

Some cells will grow readily in suspension, either spontaneously (eg, murine ascites tumors) or by mechanical stirring (eg, HeLa-S<sub>3</sub>) (Fig. 1). These cells are said to be *anchorage independent*, and the cultures are usually derived from the transformed cells of a tumor or from normal cells that have transformed *in vitro*, either spontaneously or following [mutagenesis](#) (see [Neoplastic Transformation](#)). Suspension cultures can be propagated in large bulk without elaborate mechanisms for increasing the surface area. Gas exchange is improved when the depth of medium exceeds 50 mm by sparging the medium with a mixture of 5%  $\text{CO}_2$  in air.

**Figure 1.** Cross-sectional diagram of a stirrer vessel on a magnetic stirrer. The side arm at top right is used to add cell suspension and collect samples; it has a permeable cap to allow escape of  $\text{CO}_2$ . Filtered  $\text{CO}_2$  in air is supplied via a filter to the side arm at top left; it is required if the depth of the culture medium exceeds 50 mm. Agitation is achieved by a magnet, enclosed in a glass pendulum, and is driven by the magnetic stirrer at ~60 rpm. (After design of stirrer vessel marketed by Techne, Cambridge, U.K.). Modified from Freshney's *Culture of Animal Cells, A Multimedia*

Guide, 1999, Wiley-Liss, New York.



Suspension cultures can be scaled up easily, deliver a large bulk of cells at one time, and do not need trypsin for harvesting the cells from the culture. They can also be maintained in a steady state of growth by regulating the rate at which medium is added to balance the rate of cell proliferation and by withdrawal of surplus cells. Such cultures are called biostats or vivostats and enable the concentrations of cells, nutrients, and products, and the pH, osmolality, and gas tension, to be kept constant over prolonged periods (28).

## 7. Culture Vessels

Cell cultures are usually grown in disposable plastic flasks, petri dishes, or multiwell plates, that have been treated with plasma discharge, or some similar process, to create a net negative charge on the surface of the flask. Flasks are preferable for long-term propagation, dishes for cloning or where direct access to the growth surface is required (and are cheaper), and multiwell plates for replicate sampling. Where a large number of cells is required ( $>1 \times 10^9$ ), roller bottles, or multisurface propagators are required for attached cells, and large fermentors for suspension-grown cells.

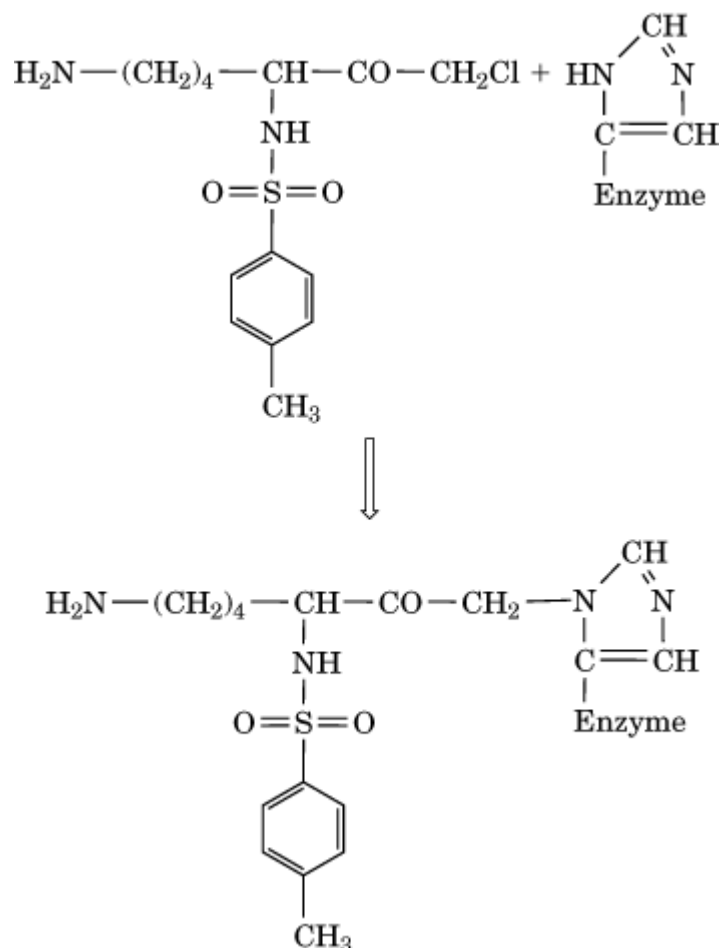
## Bibliography

1. R. G. Harrison (1907) *Proc. Soc. Exp. Biol. Med.* **4**, 140–143.
2. A. Carrel (1912) *J. Exp. Med.* **15**, 516–528.
3. A. Fischer (1925) *Tissue Culture; Studies in Experimental Morphology and General Physiology of Tissue Cells in Vitro*, Heineman, London.
4. A. Carrel (1914) *J. Exp. Med.* **20**, 1
5. K. K. Sanford, W. R. Earle, and G. D. Likely (1948) *J. Natl. Cancer Inst.* **9**, 229.
6. G. O. Gey, W. D. Coffman, M. T. Kubicek (1952) *Cancer Res.* **12**, 364–365.
7. S. P. Leibo, P. Mazur (1971) *Cryobiology* **8**, 447–452.
8. L. W. Harris, J. B. Griffiths (1977) *Cryobiology* **14**, 662–669.
9. J. E. Lovelock, and M. W. H. Bishop (1959) *Nature* **183**, 1394–1395.
10. S. M. Gartler (1967) *2nd Bicennial Review Conf. Cell, Tissue and Organ Culture*, NCI Monographs, pp. 167–195.
11. W. A. Nelson Rees, D. W. Daniels, R. R. Flandermeyer (1981) *Science* **212**, 446–452.
12. C. S. Stulberg, W. D. Peterson, Jr., W. F. Simpson (1976) *Am. J. Hematol.* **1**, 237–242.
13. R. J. Hay (1991) *Dev. Biol. Stand.* **75**, 193–204.
14. M. F. Barile (1977) in R. T. Acton, J. D. Lynn, eds., *Cell Culture and Its Applications*, Academic Press, New York, p. 291.
15. T. R. Chen (1977) *Exp. Cell. Res.* **104**, 255.
16. W. I. Schaeffer (1990) *In Vitro Cell Dev. Biol.* **26**, 97–101.
17. R. I. Freshney (2000) *Culture of Animal Cells, a Manual of Basic Technique*, Wiley-Liss, New York.
18. P. B. Medawar (1948) *Quart. J. Micr. Sci.*, **89**, 187.
19. I. Lasnitzki (1992) in *Animal Cell Culture, a Practical Approach*, R. I. Freshney, ed., IRL Press at Oxford Univ. Press, Oxford, pp. 213–261.
20. J. Bornstein, P. L. Schwartz, R. E. H. Wettenhall (1973) in *Tissue Culture, Methods and Applications*, P. F. Kruse, Jr. and M. K. Patterson, Jr., (eds.), Academic Press, New York, pp. 331–333.
21. R. Auerbach, C. Grobstein (1958) *Exp. Cell. Res.* **15**, 384–397.
22. M. Chambard, B. Vemer, J. Gabrion, J. Mauchamp, J. C. Bugeia, C. Pelassy, B. Mercier (1983) *J. Cell Biol.* **96**, 1172–1177.
23. H. E. Zhau, T. J. Goodwin S.-M. Chang T. L. Baker L. W. K. Chang (1997) *In Vitro Cell. Dev. Biol., Animal* **33**, 375.
24. R. A. Knazek, P. Gullino, P. O. Kohler, R. Dedrick (1972) *Science* **178**, 65–67.
25. A. Bohnert, J. Hornung, I. C. Mackenzie, and N. E. Fusenig (1986) *Cell Tissue Res.* **244**, 413–429.
26. K. A. Elliget, J. F. Lechner (1992) in *Culture of Epithelial Cells*, R. I. Freshney, ed., Wiley-Liss, New York, pp. 181–196.
27. D. Gospodarowicz, D. Delgado, I. Vlodavsky (1980) *Proc. Natl. Acad. Sci. USA* **77**, 4094–4098.
28. A. Kadouri and R. E. Spier (1997) *Cytotechnology* **24**, 89–98.

**TLCK (*N-P*-Tosyl-Lysine Chloromethyl Ketone)**

This reagent was developed by Shaw and co-workers as a specific inhibitor of **trypsin** (1). Its design was based on the principle of [affinity labeling](#), which couples substrate specificity with chemical reactivity to generate a covalent bond between the reagent and a functional group at the [active site](#) of an enzyme. In this case, the [lysine](#) part of the reagent confers specificity for trypsin and the chloromethyl ketone provides chemical reactivity (Fig. 1). Although trypsin is a [serine proteinase](#), the reagent does not interact with the active site serine residue but instead alkylates the [histidine](#) residue that is part of the [catalytic triad](#). The reagent has essentially no effect on the activity of **chymotrypsin** or other serine proteinases and is often used to inactivate the small amounts of trypsin present in purified chymotrypsin preparations. It is also often combined with inhibitors of other types of proteinases to prevent unwanted proteolysis that might occur during the course of protein isolation from cell or tissue homogenates (2).

**Figure 1.** Interaction of TLCK (tosyl-lysyl-chloromethyl ketone, ie, L-1-chloro-3[4-tosylamido]-7-amino-2-heptanone) with the active-site histidine side chain of a trypsin-like serine proteinase. The reaction product is catalytically inactive.



#### Bibliography

1. E. Shaw, M. Mares-Guia, and W. Cohen (1965) *Biochemistry* **4**, 2219 ff.
2. R. J. Beynon (1989) In *Protein Purification Methods: A Practical Approach* (E. L. V. Harris and S. Angal, eds.), IRL Press, Oxford, U.K., pp. 40–49.

## Tobacco Mosaic Virus


*Tobacco mosaic virus* (TMV) was the first **virus** to be recognized as a disease entity different from **bacteria** when Beijerinck (1) concluded that it was a new type of infectious agent, “Contagium vivum fluidum.” Since then, it has spurred the development of various concepts of viruses, being, among other aspects, the first virus to be purified and crystallized and shown to be composed of **protein** and **RNA** (2). Details of the history of TMV are given in Fraenkel-Conrat (3).

TMV is found worldwide and locally can cause an important disease in tobacco. It has a large host range, infecting at least 200 species from 30 families (4). The virus occurs in very high concentrations in infected leaves ( $> 10^6$  particles per cell) and is very stable, even retaining its infectivity in nonsterile extracts at room temperature for more than 50 years (5). These properties make it very readily transmitted by contact between plants. Man is the main vector, either directly by handling infected and then healthy plants or indirectly by farm machinery. The virus survives in dead crop debris and thus can be transmitted to the next season's crop at planting time.

TMV is the type member of the tobamovirus genus, which contains 13 species and two other possible members.

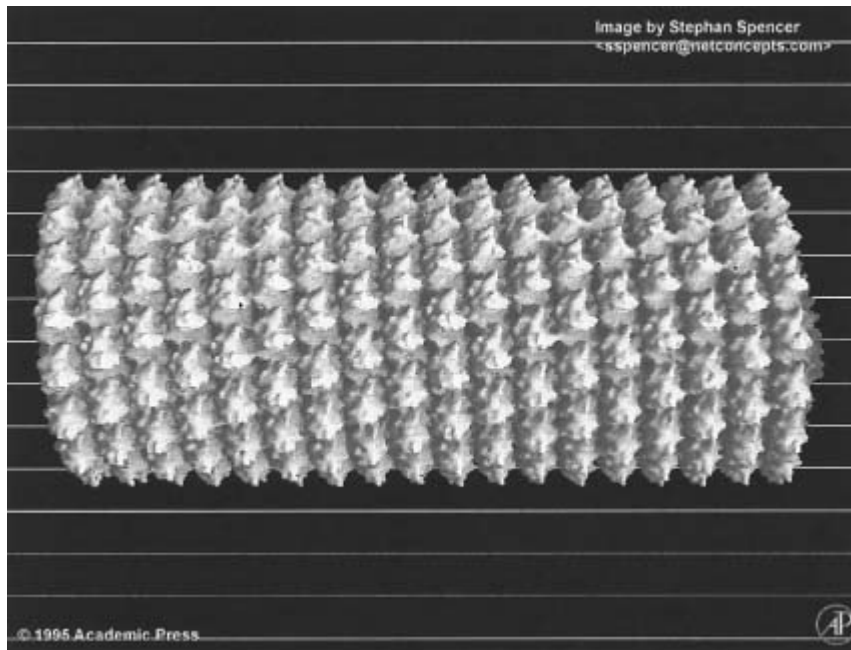
### 1. Virus Structure

The **virions** of TMV are rod-shaped, 300 nm long, and 19 nm in diameter, and composed of a single-coat protein species encapsidating a single (+)-strand linear RNA molecule of about 6.4 kb. The detailed structures of the virions and coat protein subunits have been determined by **electron microscopy** and **X-ray crystallography**. In a virion, approximately 2100 subunits are closely packed

in a single right-handed helix of pitch 2.3 nm and with 16  coat protein molecules per turn (Fig. 1). The RNA binds at a radius of about 4 nm, with three nucleotides per protein subunit. There is a central canal of about 2-nm radius. The length of the rod is determined by the length of the RNA, which is fully encapsidated.

**Figure 1.** The structure of part of TMV. The portion indicated includes 18 turns of the helical array, with  $16\frac{1}{3}$  subunits/turn. The entire virus consists of about 128 turns.





Much is known about the detailed structure of the protein subunits and their interactions in forming the virion (6). The protein monomer aggregates in solution in various ways, depending on ionic strength, pH, and temperature (7, 8). The structure of one of these forms, a double disk of 17 subunits per disk that is an intermediate in virus assembly, has been determined by X-ray crystallography to 2.8 Å resolution. Up to a radius of 4 nm, no structure is resolved, which suggests a disordered state. Much of the rest of each protein subunit is made up of four **a-helices**, with their distal ends bound together by regions of **b-sheet**. Both the N- and C- termini of the polypeptide chain occur at the circumference of the disk. The subunits have **polar** and **hydrophobic** interactions with each other. In the virion, there are also **electrostatic interactions**, between both the RNA and protein one involving pairs of carboxyl groups with anomalous  $pK_a$  values (near pH 7.0) on adjacent subunits. It is thought that these carboxyl groups play an important role in particle assembly and disassembly.

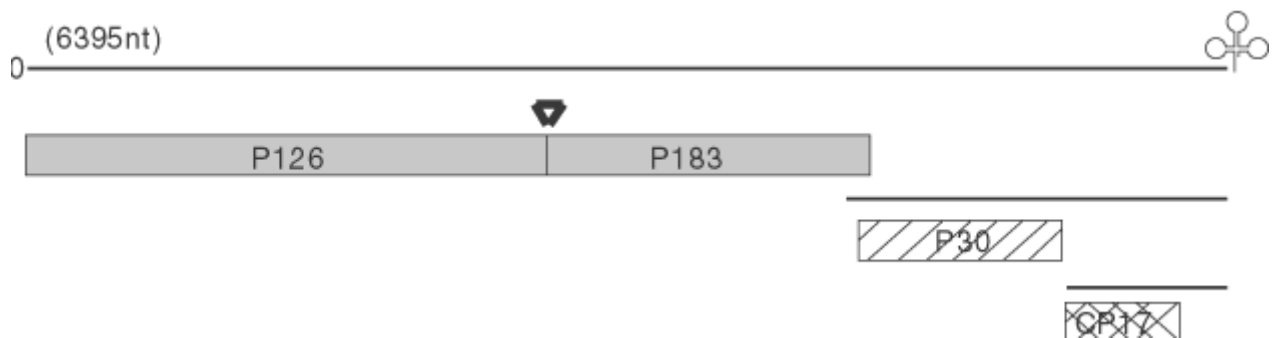
The **self-assembly** of TMV coat protein subunits has been studied extensively *in vitro* by the addition of virion RNA to coat protein preparations, which leads to the reassembly of virus particles. Assembly starts at a specific site on the RNA, termed the *origin of assembly*, which is a region of extensive stem-loop structures about 1100 to 900 nucleotides from the 3' end. In the initial initiation event, the hairpin loop of the origin of assembly interacts with a double disk of coat protein subunits so that both ends of the RNA protrude from the same side of the disk. The hairpin loop opens up as the RNA intercalates between the two layers of the double disk, which changes to a helical (double-lockwasher) form. A second double disk adds to the first on the side away from the protruding RNA tails and, as it switches to a helical form, it interacts with the RNA 5' to the origin of assembly. The helical rod continues to grow in the 5' direction by the further addition of double disks, effectively pulling the RNA through the axial hole until all the RNA 5' of the origin of assembly is encapsidated. The assembly in the 3' direction is much slower than that in the 5' direction and occurs by the addition of small aggregates of coat protein (A protein form) to the helical rod.

The *in vivo* disassembly of TMV involving the cotranslational disassembly mechanism is discussed under [Virus infection, plant](#).

## 2. Virus Genome

The [genome](#) of TMV comprises 6395 nucleotides and codes for four, or possibly five, proteins (Fig. 2). The first **open reading frame**, encoding a protein of 126 kDa, ends in an **amber** stop codon, but this is occasionally *read through* to also give a protein of 183 kDa. These two proteins are involved in TMV replication and make up the viral **RNA-dependent RNA polymerase** of the “Sindbis virus supergroup.” There are amino acid sequence motifs suggestive of a **methyl transferase** and **helicase** domain in the 126-kDa protein and of an RNA polymerase in the readthrough portion of the 183-kDa protein. There is a suggestion that the 54-kDa readthrough part of this protein might also be expressed independently from a subgenomic RNA. The other two downstream proteins, one at 30 kDa and the other at 17.6 kDa, are expressed from subgenomic [messenger RNAs](#). The 30-kDa protein is involved in cell-to-cell movement (see [Virus Infection, Plant](#)), and the 17.6-kDa protein is the viral coat protein.

**Figure 2.** Genome organization of TMV. The upper line represents the single-stranded RNA genome; the circle at the 5' end indicates the cap, and the tRNA-like structure at the 3' end is also illustrated. The other two lines are the subgenomic RNAs. The open reading frames are shown by the boxes, and their position on the diagram indicates from which RNA they are expressed. The proteins are distinguished by their approximate molecular weights, in kDa. The stippled box (P183) indicates the RNA polymerase, that striped (P30) is the protein believed to be involved in movement, whereas that cross-hatched (P17) is the coat protein; ▼ indicates readthrough.



The viral RNA has a methyl guanosine [cap](#) and leader sequence of 67 nucleotides that is very AU-rich. This leader sequence, termed the W sequence, appears to have no secondary structure and enhances [translation](#) (9). The 3'-end folds into a **transfer RNA**-like structure and accepts histidine.

## Bibliography

1. M. W. Beijerinck (1898) *Versl. Gewone Vergad. Wis- Natuurkd. Afd., K. Akad. Wet. Amsterdam* **7**, 229–235.
2. F. C. Bawden and N. W. Pirie (1937) *Proc. R. Soc. Lond. Ser. B* **123**, 274–320.
3. H. Fraenkel-Conrat (1986) In *The Plant Viruses*, Vol. **2**: The Rod-shaped Plant Viruses (M. H. V. van Regenmortel and H. Fraenkel-Conrat, eds.), Plenum Press, New York, pp. 5–17.
4. M. Zaitlin and H. W. Israel (1975) *Commonwealth Mycological Institute/Association of Applied Biologists*, Descriptions of Plant Viruses No. 151.
5. G. Silber and L. G. Burk (1965) *Nature (Lond.)* **206**, 740–741.
6. K. E. Richards and R. C. Williams (1976) *Compr. Virol.* **6**, 1–37.
7. A. C. H. Durham, J. T. Finch, and A. Klug (1971) *Nature (Lond.) New Biol.* **229**, 37–42.
8. R. E. F. Matthews (1991) *Plant Virology*, Academic Press, San Diego, CA, pp. 106–113 and 217–223.
9. D. R. Gallie, D. E. Sleat, J. W. Watts, P. C. Turner, and T. M. A. Wilson (1987) *Nucl. Acids Res.* **15**, 8693–8711.

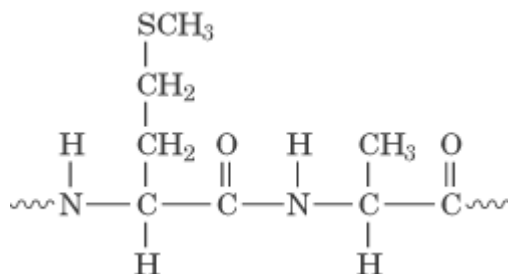
### Suggestions for Further Reading

10. D. J. Lewandowski and W. O. Dawson (1994) "Tobamoviruses". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1436–1441.
11. M. H. V van Regenmortel and H. Fraenkel-Conrat, eds. (1986) *The Plant Viruses*, Vol.2: *The Rod-shaped Plant Viruses*, Plenum Press, New York.

### TOCSY Spectrum

Total correlation spectroscopy (TOCSY) is a type of two-dimensional nuclear magnetic resonance (NMR) spectrum in which cross peaks result from coherence transfers. For studies of biological macromolecules, TOCSY normally monitors hydrogen atoms (protons), with cross peaks present because of coherence transfer from hydrogen to hydrogen of the macromolecule. TOCSY is a powerful tool for making assignments of proton NMR signals to specific protons of the molecule being examined. The information provided in a TOCSY spectrum contains the data present in a correlation spectroscopy (COSY)-type experiment but includes additional cross peaks that arise because of coherence transfers to or from all protons in a network of spins.

The requirement for the appearance of cross peaks in a proton TOCSY spectrum is the presence of a collection of mutually spin-coupled protons. Consider the alanine residue in the segment of polypeptide shown below:



Protons of the methyl group are J-coupled to the proton attached to the  $\alpha$ -carbon ( $J \sim 7\text{Hz}$ ), and the proton on the  $\alpha$ -carbon is spin-coupled to the amide N-H proton ( $J \sim 2\text{--}10\text{Hz}$ ). The methyl protons are separated from the amide N-H proton by four chemical bonds; thus, the coupling constant between the methyl spins and the peptide N-H is too small to be resolved under typical experimental conditions. Because a resolved coupling constant is necessary for detectable coherence transfer, no such transfer between the methyl protons and the N-H proton would be expected in a standard COSY experiment. A pathway exists for such transfer, however, involving transfer of methyl proton coherence to the  $\alpha$ -hydrogen and then transfer from the  $\alpha$ -hydrogen to the peptide N-H. The results of these coherence transfers are seen in a TOCSY spectrum as cross peaks that are characterized by two [chemical shift](#) coordinates. The coordinate in one dimension corresponds to the precessional frequency of the coherence before transfer, whereas the other chemical shift coordinate corresponds to the chemical shift of the coherence after the coherence transfer process has taken place. In a TOCSY spectrum, an alanine residue will be represented by three diagonal peaks, corresponding to the shifts of the  $\text{CH}_3$ ,  $\text{C}_\alpha\text{H}$ , and N-H protons, and by six cross peaks, corresponding to all possible origins of coherence and all possible destinations for coherence transfer.

The usefulness of the TOCSY experiment becomes apparent on recognition that the spin-coupling constants are zero between protons on adjacent residues of a polypeptide, so that coherence transfers *cannot* take place between different amino acid residues of the polypeptide chain. A network of cross peaks in TOCSY can only arise from protons within the *same* amino acid. If coherence transfer takes place between all possible partners within a residue, the chemical shifts of all spins in that residue can be identified.

Some of the cross peaks present in TOCSY will also be present in a COSY spectrum of the same sample. A significant advantage of TOCSY over COSY is that all components of a TOCSY cross peak tend to be of the same sign, whereas COSY cross peaks have both positive and negative components. COSY cross peaks can therefore become self-canceling when the separation of the cross peak components is inadequate.

The critical element of a TOCSY experiment is the portion known as the isotropic mixing period, during which transfers of coherence take place. The length of this mixing time typically ranges from 20 to 70 ms. For values near the low end of the range, only coherence transfer to nearby spins takes place, and the TOCSY spectrum is similar to a COSY-type spectrum. For longer mixing times, coherence transfers over the entire spin-coupled network are possible. The extent of coherence transfer (reflected in the intensity of a cross peak) is a complex function of the spin-coupling constants in the network, the length of the mixing time, and the method used to achieve isotropic mixing. No assurance exists that an expected cross peak will indeed appear in the TOCSY spectrum for a specific mixing time; thus, for unambiguous identification of all possible coherence transfers in a particular spin system, it may be advantageous to have TOCSY spectra recorded using several mixing times.

The elements of the TOCSY experiment can be built into experiments that produce three-dimensional or higher NMR spectra. (See also **NMR, COSY spectrum, Scalar coupling**.)

#### Suggestions for Further Reading

L. Braunschweiler and R. R. Ernst (1983) *J. Magn. Reson.* **53**, 521–533.

J. Cavanagh, W. J. Fairbrother, A. G. Palmer III and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.

## Tolerance, Immunological

Tolerance is a central notion in immunology, because the [immune response](#) has *a priori* the potential power of self-recognition, which implies the possibility of fatal aggression against self constituents (see [Autoantibody](#), [Autoimmunity](#), and [Autoimmune Diseases](#)). This is already present in the famous formula of *horror autotoxicus* put forward by Ehrlich in 1901. As the cells endowed with specific recognition, [T cells](#) and [B cells](#), are permanently synthesized, an *ad hoc* mechanism must operate to prevent the emergence and expansion of potentially autoaggressive clones. Natural autoantibodies permanently circulate in the body, but these are not harmful because of their low concentration, low affinity, polyspecificity, connectivity, and germline-encoded nature.

The first known example of tolerance, given in all immunology textbooks, was the observation by Owen in the 1940 that, as a consequence of a natural cross-circulation of blood during fetal life, dizygotic twin calves have a mixed population of red cells and can accept grafts from each other. In 1953, Brent, Billingham, and Medawar showed that injection of **allogenic** lymphoid cells at birth

induced a tolerant state that allowed the treated newborn, once adult, to accept skin grafts from an animal of the same strain as the initial donor. This indicated that tolerance was an acquired phenomenon. It was also shown to be specific, because the pretreated animal rejected the graft of a third strain. This experiment was crucial for Burnet's proposals of the clonal theory that explained tolerance by elimination of the autoreactive lymphocytes (the “forbidden clones”).

Some years later, tolerance induction was extended to soluble [proteins](#), when it was observed that the nature of the immune response was very much dependent on the dose of [immunogen](#) administered. An antibody response is obtained only in a relatively narrow range, whereas low and high doses induce unresponsiveness, called low-dose and high-dose tolerance. The two phenomena appear to be different, because the high-dose tolerance tends to revert spontaneously as the [antigen](#) is eliminated catabolically, whereas the low-dose tolerance is long lasting. The mechanisms of tolerance are still subject to debate, because clonal deletion, anergy, and suppression have been proposed, and all are supported by a large body of data.

The acquisition of central tolerance takes place in the primary lymphoid organs—that is, in the thymus and in the bone marrow for T and B cells, respectively, where a screen for autoreactivity of newly emerging cells operates before they may enter the periphery. For T cells, this happens when  $CD4^+CD8^+$  double-positive (DP) thymocytes that have just rearranged their [T-cell receptor](#) genes are subjected to negative selection in the thymic cortex (see [Gene Rearrangement](#)). It is accepted that thymocytes having a high affinity for major histocompatibility complex (MHC) peptides of the thymic environment are deleted. In fact, this is quite difficult to prove, and mechanisms to account for negative selection are still awaited.

B lymphocytes are screened for autoreactivity before leaving the bone marrow, at the immature B-cell stage, when expressing only [IgM](#) at their surface. When exposed to antibody directed against their heavy chain, m, these cells die, suggesting that their elimination in the bone marrow results from interaction with a multivalent antigen. This view is also supported by experiments with transgenic mice that express H and L genes encoding an antibody directed against H-2<sup>k</sup>, an MHC class I molecule. In non H-2<sup>k</sup> mice, the B cells develop normally and express the transgenic antibody. In H-2<sup>k</sup> mice, a severe blockage of differentiation occurs at the preB stage, indicating that most B cells die *in situ*. This is a typical example of clonal deletion, as postulated by Burnet. It is of interest to note that blockage in this transgenic models may be leaky, because some level of endogenous gene rearrangement of light chain genes may occur, implying a transient reactivation of the [recombinase](#). These new light chains may now associate with the transgenic m chain, resulting in the occurrence of antibodies endowed with new specificity. This mechanism, known as *receptor editing*, is both an additional mechanism to generate diversity and a way to rescue clones that otherwise would have been deleted. More sophisticated models of double transgenic mice have been studied, using transgenes for both the antigen and the corresponding antibody. In that case, tolerance is set up in two distinct ways. When the transgenic antigen is expressed as a membrane-bound molecule, B cells are deleted. When it is expressed as a soluble form, B cells are present, but anergized. No surface IgM is found, although cytoplasmic IgM is present. Furthermore, surface [IgD](#) is present that might also be linked to an alteration of the transduction cascade, which questions whether IgD might have a possible role in this type of tolerance.

Finally, it should be mentioned that tolerance may also be seen as the result of a network of regulation, possibly involving cytokines and an equilibrium between Th cells that might have an active suppressive effect. In this regard, the possibility of inducing dominant tolerance by grafting thymic cells from an exclusively epithelial origin is compatible with such a view, because this tolerant state can be transferred by lymphoid cells. This does not necessarily imply the existence of suppressor T cells as such, but could result from a systemic generalized regulation.

Experimentally induced tolerance may be broken, for example, by attempting stimulation with cross-reactive antigens. This opens the way to an approach to a better understanding of **autoimmune**

diseases, which certainly could be considered to be the consequence of a dysbalance in the regulation of the immune system, inducing a break in the physiological tolerance to self components.

See also entries [Autoantibody](#), [Autoimmunity](#), [Autoimmune diseases](#), [Gene rearrangement](#), [Recombinases](#), [Clonal selection theory](#), and [Immunogen](#).

#### Suggestions for Further Reading

C. C. Goodnow (1996) Balancing immunity and tolerance: deleting and tuning lymphocyte repertoires. *Proc. Natl. Acad. Sci. USA* **93**, 2264–2271.

N. Le Douarin et al. Evidence for a thymus-dependent form of tolerance that is not based on elimination or anergy of reactive T cells. *Immunol. Rev.* **149**, 35–53.

E. L. Prak, M. Trounstein, D. Huszar, and M. Weigert (1994) Light chain editing in kappa-deficient animals: a potential mechanisms of B cell tolerance. *J. Exp. Med.* **180**, 1805–1815.

G. J. V. Nossal (1991) B-cell selection and tolerance. *Curr. Opin. Immunol.* **3**, 193–198.

## Tomato Bushy Stunt Virus

*Tomato bushy stunt virus* (TBSV) was first described in England in 1935 (1) and has now been recognized in many countries in Europe, the Mediterranean regions, and North and South America. It causes economic problems in some countries, with tomato crops but is not generally recognized as a major problem. It has been found to infect several other plant species naturally, and its artificial host range is wide. Although the virus spreads readily through crops, no specific vector has been identified. It is transmitted by contact and through the soil and is found in rivers, which may account for its long-distance contamination. The virus occurs in high concentrations in infected plants, is stable, and is easily purified. Thus, it has been a good model for studies on virus structure and basic molecular biology.

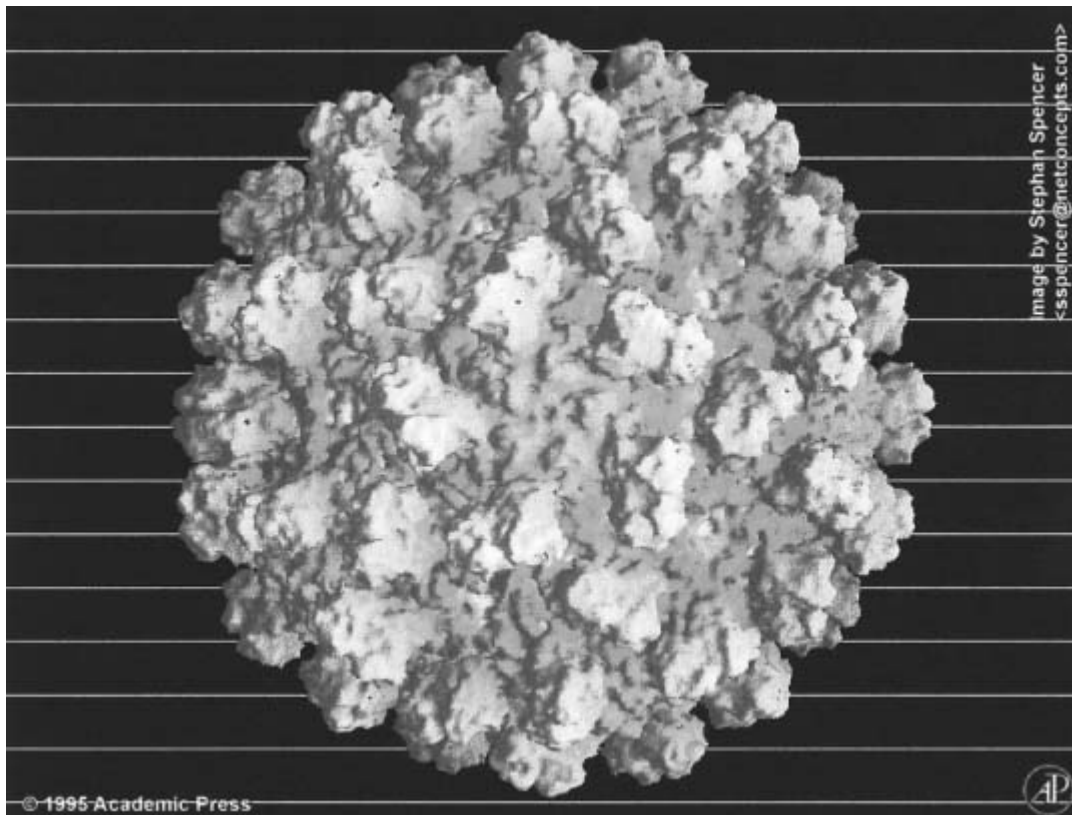
TBSV is the type member of the Tombusvirus genus, which currently contains 13 species.

### 1. Virus Structure

The **virion** of TBSV is isometric, with a diameter of about 30 nm (Fig. 1). It is composed of a single type of coat protein encapsidating a single (+)-strand linear **RNA** molecule of about 4.7 kb. The detailed structures of the virions and coat protein subunits have been determined by [electron microscopy](#) and [X-ray crystallography](#). In fact, TBSV was the first virus shown to have **icosahedral symmetry** (2) and thus has formed the basis of our understanding of isometric virus structure. The 180 coat protein subunits of the virus particle are arranged in  $T=3$  symmetry (see [Virus Structure](#)), and the virion structure has been determined at 2.9 Å resolution. Each coat protein subunit comprises 387 amino acids, which are folded into three distinct **domains**. The N-terminal domain (the R domain) is composed of 66 residues, with a high proportion of basic amino acids, and interacts internally with the RNA. This is connected by an arm of 35 residues to the rest of the molecule, which is made up of the globular shell (*S*) domain (167 residues) and surface-protruding (*P*) domain (114 amino acids); the *S* and *P* domains are connected by a flexible hinge of five residues. The *P* domains are clustered as 90 dimers, which are important in virion assembly and stability. The *S* domain forms a barrel structure composed of eight [Beta-Strand](#), which is now recognized as being characteristic of the coat proteins of many isometric viruses. The 180 coat protein subunits of the virus particles have three different **quasi-equivalent** conformations, which are necessary for icosahedral symmetry packing. The A and B conformations are very similar, with disordered N-

terminal arms. The C conformation has an ordered N-terminal arm that interdigitates with two other C subunit arms around the icosahedral three-fold axis.

**Figure 1.** The external surface of tomato bushy stunt virus, viewed down a local five-fold axis.



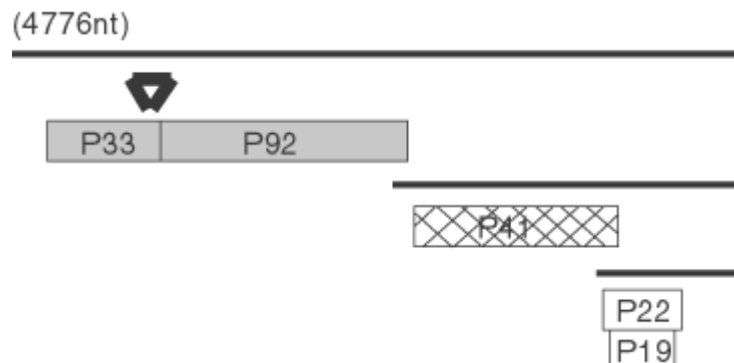
Three major interactions stabilize the virus particle: (1) [electrostatic interactions](#) between the *R* domain and genomic RNA, (2) a pH-dependent interaction between protein subunits that protonates above pH 7, and (3) an interaction involving  $\text{Ca}^{++}$  ions. The  $\text{Ca}^{++}$  links adjacent subunits via the carboxylate groups from the side-chains of **glutamate** and **aspartate** residues in the *S* domain.

## 2. Virus genome

The [genome](#) of TBSV comprises 4776 nucleotides and codes for five proteins (Fig. 2). The first **open reading frame**, encoding a protein of 33 kDa, terminates in an **amber** stop codon, but this is occasionally *read through* to yield a protein of 92 kDa (see [Viruses, Plant](#)). These two proteins are essential for virus replication, and it is thought that they make up the viral **RNA-dependent RNA polymerase**. The readthrough portion of the 92-kDa protein contains the motif of Gly-Asp-Asp (GDD) typical of such RNA polymerases. The third open reading frame encodes the viral coat protein and is expressed from a subgenomic 2.1-kb [messenger RNA](#). The fourth and fifth open reading frames overlap each other, and both are expressed from another subgenomic RNA (0.9 kb). [Mutagenesis](#) indicates that both of these proteins are essential for full infection of the host, and there is a suggestion that the 22-kDa polypeptide is involved in cell-to-cell movement (3).

**Figure 2.** Genome organization of TBSV. The upper line represents the single-stranded RNA genome, and the other two lines are the subgenomic RNAs. The five open reading frames are shown by the boxes, and their position on the

diagram indicates from which RNA they are expressed. The proteins are distinguished by their approximate molecular weights, in kDa. The stippled box indicates the RNA polymerase, whereas the cross-hatched box is the coat protein; ▽ indicates readthrough.



### Bibliography

1. K. M. Smith (1935) *Ann. Appl. Biol.* **22**, 731–741.
2. D. L. D. Caspar (1956). *Nature* **177**, 475–476.
3. H. B. Scholthof, T. J. Morris, and A. O. Jackson (1993) *Mol. Plant-Microb. Interact.* **6**, 309–322.

### Suggestions for Further Reading

4. G. P. Martelli, D. Gallitelli, and M. Russo (1988) "Tombusviruses". In *The Plant Viruses, Vol.3: Polyhedral Virions with Monopartite RNA Genomes* (R. Koenig, ed.), Plenum Press, New York, pp. 13–72.
5. D. M. Rochon (1994) "Tombusviruses". In *Encyclopedia of Virology* (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 1447–1452.
6. M. Russo, J. Burgyan, and P. Martelli (1994) *Molecular biology of Tombusviridae*. *Adv. Virus Res.* **44**, 381–428.

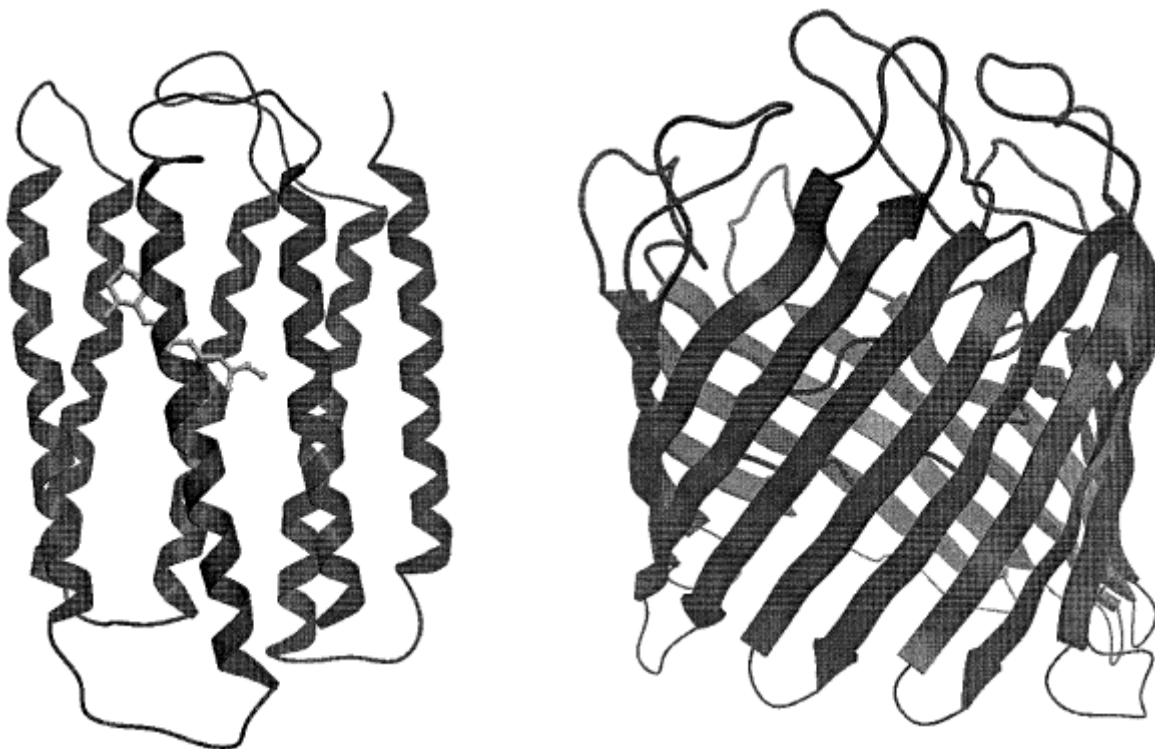
### Topogenesis

Integral [membrane proteins](#) are embedded in the lipid bilayer either by membrane-spanning **a-helices**, which for multispanning or oligomeric proteins form a helix bundle, or by **b-strands** forming a b-barrel. This is attested by the few solved three-dimensional structures of membrane proteins. As examples, Figure 1 shows the backbone structures of [bacteriorhodopsin](#) (1) composed of seven transmembrane a-helices, and of the matrix [porin](#) OmpF (2), a barrel of 16 antiparallel b-strands. b-Barrel proteins are found exclusively in the outer membranes of bacteria, [mitochondria](#), and [chloroplasts](#). Little is known about the process of their integration into the membrane. Because the **hydrophobic** surface is generated only when the b-sheet has assembled, it is likely that **protein folding** precedes insertion. In contrast, individual hydrophobic helices can fold and insert independently, and they can assemble into the bundle in a separate second step (3). This model is supported by refolding experiments and by the functional assembly of membrane proteins from separately expressed fragments (eg, Ref. 4). The basic topology of helical membrane proteins is thus defined in the first step, the insertion of membrane-spanning segments into the bilayer. The



determinants that direct topogenesis and the machinery that decodes them are largely conserved between prokaryotes and eukaryotes. In the best-studied systems, the cytoplasmic membrane of bacteria and the [endoplasmic reticulum](#) (ER) of eukaryotic cells, a hydrophobic signal sequence is recognized by a targeting system and guided to a gated hetero-oligomeric transmembrane channel. Exoplasmic [hydrophilic](#) portions of the [polypeptide chain](#) are translocated across the membrane, whereas apolar sequences of ~15 to 25 residues are released into the lipid bilayer as transmembrane  $\alpha$ -helices. Their orientation is determined mainly by flanking charged residues, but also by their hydrophobicity and length, by their position within the polypeptide, and by the folding properties of adjacent polar domains. Our current understanding of topogenic determinants allows reasonably efficient prediction of protein topology. However, the molecular mechanisms by which the determinants are recognized have not been established yet.

**Figure 1.** Helix bundles versus  $\beta$ -barrels. The backbone structures of bacteriorhodopsin (1) of *Halobacterium halobium* (a) and matrix porin OmpF (2) of *Escherichia coli* (b). The extracellular surfaces are facing up. (Courtesy of Ansgar Philippsen, Biozentrum, University of Basel.)



## 1. Machinery for Membrane Targeting and Insertion

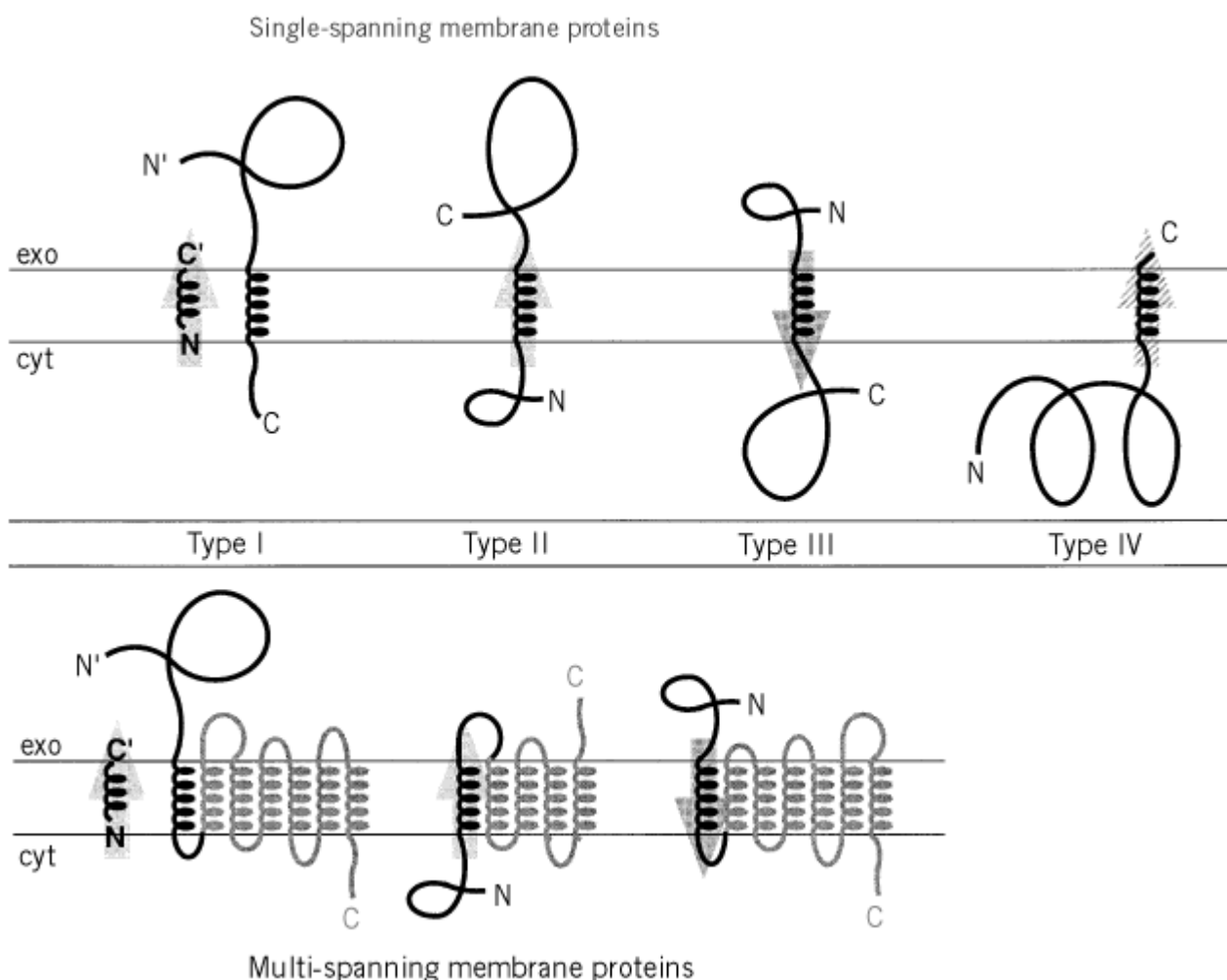
The machinery for targeting, translocation, and membrane insertion of proteins is largely conserved between prokaryotic and eukaryotic cells (5, 6). Secretory and membrane proteins are targeted to the membrane by [signal peptides](#) near the amino terminus of the polypeptide in a cotranslational or posttranslational process. In cotranslational targeting (which appears to be the predominant mechanism in mammalian cells, but not in yeast and bacteria), the signal is recognized by the [signal recognition particle](#) (SRP; Ffh-RNP in *Escherichia coli*) and is delivered to the translocation machinery by interaction with the SRP receptor (FtsY in *E. coli*). In posttranslational targeting, cytosolic **molecular chaperones** (eg, SecB in *E. coli*, Hsp70p in yeast) are required to keep the proteins in a translocation-competent state. In both cases, translocation and membrane insertion are mediated by a heterotrimeric complex, the Sec61 complex in eukaryotes and SecYEG in prokaryotes, which forms a proteinaceous channel. Reconstitution experiments showed that the

Sec61 complex in combination with SRP, SRP receptor, and (depending on the signal sequence) TRAM (translocating chain-associating membrane protein) is sufficient for cotranslational translocation (7, 8). Posttranslational translocation in yeast requires the heterotetrameric Sec62 / 63 complex, in addition to the Sec61 complex, as well as Kar2p as a luminal **translocation motor** (9). The minimal machinery in bacteria consists of the SecYEG complex and SecA, which serves to recruit secretory and membrane proteins to the membrane and to drive translocation by its **ATPase** activity (10, 11).

## 2. Classification of Membrane Protein Topologies

The determinants within a polypeptide that direct it to the membrane and that define its topology are largely conserved and interchangeable between prokaryotes and eukaryotes. Targeting and membrane insertion are initiated by signal sequences. Three types of signals can be distinguished that cooperate with the SecYEG or Sec61 translocation machinery (Fig. 2): cleaved signals, signal-anchor, and reverse signal-anchor sequences. Common to all of them is a stretch of apolar residues.

**Figure 2.** Classification of membrane protein topologies. The classification is based on the characteristics of the initial signal sequence (cleaved signal, signal-anchor, reverse signal-anchor, *C*-terminal signal), not just the final disposition of *N*- and *C*-termini. The arrows indicate the *N*-to-*C* direction in the signals to highlight the difference between signals translocating the carboxy-terminal or the amino-terminal sequence. (Modified after Ref. 20.)



Cleaved signals are typically composed of a positively charged amino-terminal segment of 1 to 5

residues, a central hydrophobic domain of 7 to 15 apolar residues, and a more polar carboxy-terminal segment of 3 to 7 residues at the site of cleavage by [signal peptidase](#). Like signal sequences as a whole, signal cleavage sites are degenerate. The main characteristics are a preferred proximity to the hydrophobic domain and a strong preference for small uncharged residues at positions -1 and -3 with respect to the cleaved peptide bond (12). The signal is inserted into the translocon as a loop (13, 14), exposing the amino terminus to the cytosol and the carboxy-terminal cleavage site to the exoplasmic surface. The carboxy-terminal sequence is then translocated, either completely, resulting in a secretory protein, or up to a second topogenic element, a hydrophobic stop-transfer sequence. In the latter case, translocation is stopped, and the final protein is anchored in the membrane by the stop-transfer sequence. Such single-spanning membrane proteins with a processed exoplasmic *N*-terminus and a cytoplasmic *C*-terminus (ie, an  $N_{\text{exo}}/C_{\text{cyt}}$  topology) are classified as type I membrane proteins. Examples are glycophorin A and filamentous phage protein gpIII.

In contrast to cleaved signals, signal-anchor sequences are not necessarily located at the very amino terminus, and they lack a signal cleavage site. Their hydrophobic domains are longer, approximately 20 residues, which is sufficient to span the core of the lipid bilayer in an  $\alpha$ -helical conformation and to anchor the final protein in the membrane. As for cleaved signals, the hydrophobic cores of signal-anchor sequences are flanked by predominantly positive charges toward the amino terminus and by negative or fewer positive charges toward the carboxy terminus. They also promote translocation of the carboxy-terminal portion of the polypeptide. The resulting type II proteins thus have an  $N_{\text{cyt}}/C_{\text{exo}}$  topology, as exemplified by the asialoglycoprotein receptor and the bacterial [penicillin-binding protein Ib](#).

Reverse signal-anchor sequences, in contrast, promote translocation of the amino-terminal portion of the polypeptide, resulting in an  $N_{\text{exo}}/C_{\text{cyt}}$  topology without proteolytic processing. They differ from signal-anchors mainly by a reversed charge distribution in the flanking segments. Single-spanning proteins inserted by reverse signal-anchors are classified as type III membrane proteins (or, because of the identical final orientation, as “type I without cleaved signal” or “type Ib”). Examples are **cytochrome P-450** and bacteriophage Pf3 coat protein.

Multispanning proteins may be classified according to the most amino-terminal signal that initiates targeting and membrane insertion. For example, the [thrombin](#) receptor with an amino-terminal cleaved signal and seven-transmembrane domains (15) is a type I multispanning protein. Other members of the seven-transmembrane receptor family, such as the  $\beta$ -adrenergic receptor (16), are of type III, because they acquire the same topology with a reverse signal-anchor sequence. The erythroid anion transporter band 3 is an example of a multispanning type II protein: It consists of a 40-kDa amino-terminal, cytoplasmic domain, followed by a portion with 12 potential transmembrane segments (17), the first of which is likely to function as a signal-anchor sequence for targeting and insertion. Maltose transporter MalF and [bacteriorhodopsin](#) are prokaryotic examples of type II and type III multispanning membrane proteins, respectively.

Some membrane proteins do not span the membrane but are tethered to the exoplasmic leaflet of the bilayer by a lipid anchor. Bacterial lipoproteins are basically secreted proteins made with a cleavable signal sequence that is recognized by a specific lipoprotein signal peptidase (LspA) (18). Cleavage occurs at a [cysteine](#) residue that is first modified to diacylglycerylcysteine. In addition, the new amino terminus is also fatty acylated. In eukaryotes, a group of proteins is anchored in the membrane by glycosyl phosphatidylinositol ([GPI ANCHORS](#)) (19-21). These proteins are synthesized initially as type I single-spanning membrane proteins, generally without a hydrophilic cytoplasmic domain. In a transamidation reaction, the stop-transfer sequence is replaced, still in the ER, by a precursor GPI anchor consisting of phosphatidylinositol, a linear tetrasaccharide core ( $\text{Man}_3\text{GlcN}$ ), and a phosphoethanolamine group.

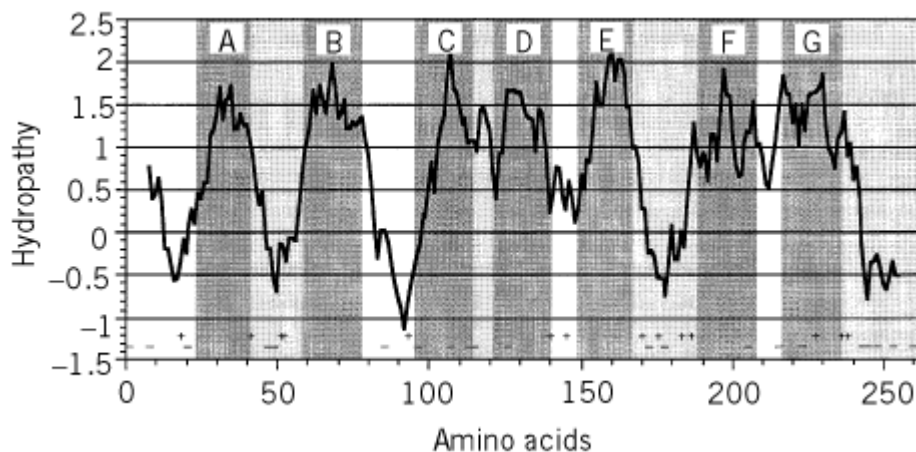
Another rapidly growing family of eukaryotic proteins, including cytochrome  $b_5$  and the SNARE

proteins, lacks classical signal sequences. Instead, these proteins are inserted into the ER, and perhaps other membranes, by a carboxy-terminal hydrophobic domain that spans the lipid bilayer (Fig. 2, type IV proteins). Insertion is independent of the Sec61 machinery and is not well-characterized at present (22-24).

### 3. Prediction of Protein Topology

The first step in predicting the topology of a protein from its sequence is to identify the transmembrane segments. Most membrane-spanning helical segments can be identified by a hydropathy analysis that plots the average hydrophobicity of amino acid side chains for a window of 11 to 21 residues along the sequence (25, 26), as illustrated in Figure 3 for bacteriorhodopsin. Ambiguities arise because transmembrane segments of multispinning or oligomeric proteins may contain hydrophilic or even charged residues, thus resulting in a lower average hydrophobicity than expected for a helix exposed to the lipid layer. In addition, adjacent relatively short transmembrane helices without a clearly hydrophilic spacer may be difficult to distinguish from a single relatively long helix.

**Figure 3.** Hydropathy plot of bacteriorhodopsin. The average hydrophobicity of a sliding window of 15 residues was calculated according to Kyte and Doolittle (25). Positive and negative residues are indicated with + and -. Exoplasmic, transmembrane, and cytoplasmic localization is indicated by white, dark gray, and light gray background, respectively. The transmembrane helices are labeled A-G.



The hydrophobic sequences of transmembrane segments do not provide useful clues as to their orientation in the bilayer, but the hydrophilic flanking sequences do. Statistical analysis revealed an enrichment of positive charges in cytoplasmic sequences and their depletion from exoplasmic ones (27). This applies to single- and multispinning proteins of eubacteria, archaebacteria, and eukaryotes (21), particularly also to signal sequences of any type. In prokaryotes, this observation was initially formulated as the “positive-inside rule” (27). For eukaryotic signal-anchor sequences, the charge difference between the two flanking segments, rather than the positive charge per se, correlates with the orientation: The cytoplasmic sequence generally carries a more positive charge than the exoplasmic one (28). The best correlation was found for the 15 flanking residues on either side of the apolar domain. The charge rule is also reflected in the structure of cleaved signals: The positively charged amino-terminus stays in the cytoplasm, whereas the carboxy-terminal end, which is depleted of positive charges, is translocated (29). The positive-inside and charge-difference rules are not restricted to the initial signal sequence, but they apply also to stop-transfer sequences and later transmembrane segments in multispinning proteins, although to a lesser degree in eukaryotic proteins than in prokaryotic (30, 31).

A combination of a hydropathy plot with the analysis of the charge distribution within the hydrophilic portions of a protein sequence yielded a useful algorithm to predict the topology [TopPred (30); also available on the internet: (site currently unavailable)]. If several related sequences are available, their alignments can be used for an improved topology prediction (32, 33) (site currently unavailable). An alternative algorithm is available for prokaryotic proteins that is independent of the charge distribution and uses the so-called dense alignment surface method (based on low-stringency [dot-plots](#) of the query sequence against a collection of non-homologous membrane proteins; (site currently unavailable)) (34).

The hallmark of b-barrel proteins are sequences of 7 to 9 residues, comprising one transmembrane strand of the b-sheet, in which every second side chain is apolar, to form the hydrophobic surface. Algorithms for the identification and topology prediction of b-barrel proteins have been developed exploiting, in addition, the typical accumulation of aromatic residues in flanking positions in porins (35) or employing a neural network (36) ([http://strucbio.biologie.uni-konstanz.de/~kay/om\\_topo\\_predict.html](http://strucbio.biologie.uni-konstanz.de/~kay/om_topo_predict.html)).

To be certain about the topology of a protein, it has to be verified experimentally. A number of approaches to define the membrane sidedness of individual hydrophilic segments have been developed for this purpose (37). Exposure to the cell exterior or to the cytosol can be assayed by sensitivity to [proteinases](#) added to intact cells or to microsomes, respectively. Similarly, the accessibility of different segments to sequence-specific **antibodies** can be tested. In multispanning proteins, cleavage patterns may be difficult to interpret, and efficient antibodies against short loops may be hard to generate. Insertion of specific cleavage sites [eg, for factor Xa (37)], or of known antigenic [epitopes](#) (38) at various positions throughout the protein have been successfully used to map protein topology. In eukaryotes, **N-glycosylation** at a natural or an engineered site is conclusive evidence for luminal exposure of the sequence. Another approach involves the fusion of truncated portions of the protein to a reporter sequence whose localization can be easily determined (eg, by enzymatic activity). The prototype of this method is the use of PhoA fusions in bacteria as a sensor for its location (39). With all assays involving mutagenesis, it has to be kept in mind that the mutation may affect the protein's topology.

#### 4. Topogenic Determinants of Single-Spanning Proteins

##### 4.1. Charged Residues

Charged residues flanking the hydrophobic segment of signal sequences were confirmed by [site-directed mutagenesis](#) to function as topogenic determinants. Insertion of positive amino acids within 30 residues carboxy-terminal of a cleaved signal inhibited its function (40, 41), more effectively so in *E. coli* than in the mammalian ER (42). The reverse signal-anchor of cytochrome P450, a type III ( $N_{\text{exo}}/C_{\text{cyt}}$ ) single-spanning protein, was converted to a type II ( $N_{\text{cyt}}/C_{\text{exo}}$ ) signal-anchor by insertion of positively charged residues into the short polar amino-terminal domain preceding the transmembrane sequence (43-45). Mutation of flanking charges in the asialoglycoprotein receptor H1 and in the paramyxovirus hemagglutinin-neuraminidase, two type II proteins, caused only a fraction of the polypeptides to insert with the opposite type III ( $N_{\text{exo}}/C_{\text{cyt}}$ ) topology (46-48).

Positive charges had a stronger effect on topogenesis than negative ones, and they were more effective the closer they were to the hydrophobic segment. In these and other studies (eg, Ref. 49), however, the asymmetric distribution of flanking charges in mutant proteins was not sufficient to generate a unique topology, indicating that additional factors influence topogenesis.

For bacteria, the simplest explanation for the topogenic role of charged residues is an involvement of the electrochemical potential across the membrane in the insertion mechanism. Indeed, the potential stimulates the translocation of chain segments containing negatively charged residues and inhibits translocation of positively charged segments in mutant forms of leader peptidase and M13 procoat, suggesting an **electrophoresis**-like membrane transfer mechanism (50, 51). In addition, the content of acidic phospholipids in the membrane was found to affect charge-dependent topogenesis (52). It

was suggested that anionic lipids prevent translocation of positive flanking charges, possibly by a direct, [electrostatic interaction](#). No general membrane potential exists across the ER membrane. It appears likely that the topogenically active charges interact with components of the translocation machinery.

#### 4.2. Folding of Hydrophilic Domains

The polypeptide needs to be largely unfolded for translocation across the membrane. For example, it was observed that the folding of a [ubiquitin](#) domain in a chimeric protein prevented its (posttranslational) translocation in yeast ([53](#)). In cotranslational translocation of type I or type II proteins, the docking of the [ribosome](#) to the Sec61 complex prevents premature protein folding ([54](#), [55](#)). The amino-terminal domain of type III proteins, however, is completed and potentially folded in the cytoplasm before the signal sequence emerges from the ribosome. Folding of this segment and its size are thus potential obstacles for transfer across the membrane. Truncation of the amino-terminal domain facilitated type III insertion of charge mutants of H1 and hemagglutinin-neuraminidase ([46](#), [47](#)). In contrast, extension of the amino-terminal domain of H1 by the complete coding sequence of [dihydrofolate reductase](#) (197 residues) blocked amino-terminal translocation and forced the protein to insert with type II topology, despite the charge distribution of a reverse signal-anchor ([56](#)). Destabilizing point mutations in the reductase largely restored its transport through the membrane. Similarly, a 30-residue **zinc finger** domain, but not a defective mutant sequence, hindered insertion as a type III protein ([56](#)). The folding properties of the amino-terminal domain are thus a limiting factor for type III insertion, whereas the size (at least up to ~230 residues) does not appear to be a significant obstacle. The amino-terminal domains of natural type III proteins are relatively short (generally less than 100 residues). They probably evolved to be devoid of rapidly and stably folding sequences to ensure a unique topology.

Whereas in eukaryotes proteins of type I, II, and III all depend on the Sec61 machinery for insertion, bacterial translocation of amino-terminal tails can be independent of functional SecA and SecY ([57-59](#)) and sensitive to the size of the amino-terminal domain ([57](#), [58](#)).

#### 4.3. Hydrophobicity

An influence of the hydrophobic segment of signal sequences on topogenesis was suggested by deletion experiments ([60](#)). Systematic analysis of signals with artificial hydrophobic domains consisting of oligoleucine sequences of different lengths revealed that translocation of the amino terminus was favored by long, hydrophobic sequences and that translocation of the carboxy terminus was favored by short ones ([61](#), [62](#)). The topogenic contributions of the hydrophobic sequence, the flanking charges, and the hydrophilic amino-terminal domain were additive. In combination, these determinants were sufficient to achieve unique membrane insertion in either orientation ([62](#)).

The topogenic contribution of the hydrophobic signal domain was also tested for natural signal sequences that did not conform with the charge difference rule ([63](#)). The results showed that the hydrophobicity of the signal is important for the correct and uniform orientation of proteins, particularly of those with unusual flanking charges.

### 5. Topogenesis of Multispanning Proteins

Extended stretches of apolar residues cannot be translocated across the membrane and are eventually released into the lipid bilayer as membrane-spanning segments. The length of the apolar segment and its hydrophobicity are the main criteria for a stop-transfer sequence ([64](#), [65](#)). For example, 19 successive alanine residues, but only 9 leucines, were required for efficient stop-transfer activity in the ER ([66](#), [67](#)). Subsequent charged residues, particularly positive ones, may increase the stop-transfer efficiency of shorter apolar sequences, but are not essential. Natural stop-transfer sequences are generally about 20 residues long to span the membrane comfortably and thus can also accommodate a few polar amino acids, which may be necessary for the protein's function.

It is an attractive model that multispanning proteins are generated by successive signal and stop-

transfer sequences, initiating with the most amino-terminal signal and proceeding sequentially toward the carboxy terminus (68). Each topogenic element might function independently in the order of appearance from the ribosome. Such linear insertion of successive transmembrane domains into the ER membrane has been demonstrated for artificial proteins composed of up to four repeats of type II signal-anchor sequences separated by relatively long hydrophilic spacers of 50 to 100 residues (69, 70). The first and third repeats inserted with a  $N_{\text{cyt}}/C_{\text{exo}}$  orientation, the second and fourth with the opposite  $N_{\text{exo}}/C_{\text{cyt}}$  orientation. The second insertion occurred at a time when the first exoplasmic loop was already translocated and could also be mediated by a mutant signal unable to interact with SRP (69). SRP is thus required only to target the first signal domain to the ER membrane, but not for insertion of subsequent segments.

However, model proteins with shorter hydrophilic spacer sequences of ~25 residues separating the transmembrane helices did not follow a strict sequential start–stop model of insertion (71). Positive charges introduced into exoplasmic loops inhibited their translocation, resulting in “frustrated” molecules in which potential transmembrane segments did not span the membrane. This demonstrated that sequences other than the initial signal may exert topogenic influence.

For bacteria, there is extensive evidence against a strictly linear insertion process in which the topology is determined by the first transmembrane domain followed by passive serpentine insertion of the subsequent helices. Deletion of individual membrane-spanning segments did not affect the topology of the following transmembrane domains toward the carboxy terminus (72–74). Positively charged residues play an important role and, inserted into normally exoplasmic loops, caused local disruption of insertion (75, 76). This reflects competition of conflicting signals within the polypeptide.

## 6. Conclusion

So far, topogenic determinants have been studied mainly by challenging the translocation machinery *in vitro* or *in vivo* with mutant proteins. Because the components of the machinery have been identified in both prokaryotes and eukaryotes, the molecular mechanisms of the insertion process can be analyzed more directly. In [crosslinking](#) experiments, it has been shown that signal-anchor sequences are already in contact with lipids early in the insertion process, suggesting lateral opening of the translocation channel (77), and how transmembrane helices exit into the lipid environment during the translocation process has been analyzed (78, 79). Further studies will be directed to determine in molecular detail how topogenic signals are decoded by the translocation machinery.

## Bibliography

1. N. Grigorieff, T. A. Ceska, K. H. Downing, J. M. Baldwin, and R. Henderson (1996) *J. Mol. Biol.* **259**, 393–421.
2. S. W. Cowan, T. Schirmer, G. Rummel, M. Steiert, R. Ghosh, R. A. Pauptit, J. N. Jansonius, and J. P. Rosenbusch (1992) *Nature* **358**, 727–733.
3. J. L. Popot and D. M. Engelman (1990) *Biochemistry* **29**, 4031–4037.
4. T. W. Kahn and D. M. Engelman (1992) *Biochemistry* **31**, 6144–6151.
5. G. Schatz and B. Dobberstein (1996) *Science* **271**, 1519–1526.
6. M. Pohlschröder, W. A. Prinz, E. Hartmann, and J. Beckwith (1997) *Cell* **91**, 563–566.
7. D. Görlich and T. A. Rapoport (1993) *Cell* **75**, 615–630.
8. S. Voigt, B. Jungnickel, E. Hartmann, and T. A. Rapoport (1996) *J. Cell Biol.* **134**, 25–35.
9. J. L. Brodsky, J. Goeckeler, and R. Schekman (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9643–9646.
10. L. Brundage, J. P. Hendrick, E. Schiebel, A. Driessen, and W. Wickner (1990) *Cell* **62**, 649–657.

11. M. Hanada, K. I. Nishiyama, S. Mizushima, and H. Tokuda (1994) *J. Biol. Chem.* **269**, 23625–23631.
12. G. von Heijne (1986) *Nucleic Acids Res.* **14**, 4683–4690.
13. A. Kuhn (1987) *Science* **238**, 1413–1415.
14. A. S. Shaw, P. J. M. Rottier, and J. K. Rose (1988) *Proc. Natl. Acad. Sci. USA* **85**, 7592–7596.
15. T. K. Vu, D. T. Hung, V. I. Wheaton, and S. R. Coughlin (1991) *Cell* **64**, 1057–1068.
16. E. Rands, M. R. Candelore, A. H. Cheung, W. S. Hill, C. D. Strader, and R. A. Dixon (1990) *J. Biol. Chem.* **265**, 10759–10764.
17. R. R. Kopito and H. F. Lodish (1985) *Nature* **316**, 234–238.
18. V. Braun and H. C. Wu (1994) In *New Comprehensive Biochemistry: Bacterial Cell Wall*, Vol. **27** (J. M. Ghuyssen and R. Hakenbeck, eds.), Elsevier Science, Amsterdam, pp. 319–341.
19. P. T. Englund (1993) *Annu. Rev. Biochem.* **62**, 121–138.
20. M. J. McConville and M. A. Ferguson (1993) *Biochem. J.* **294**, 305–324.
21. T. Kinoshita, K. Ohishi, and J. Takeda (1997) *J. Biochem. (Tokyo)* **122**, 251–257.
22. U. Kutay, E. Hartmann, and T. A. Rapoport (1993) *Trends Cell Biol.* **3**, 72–75.
23. U. Kutay, G. Ahnert-Hilger, E. Hartmann, B. Wiedenmann, and T. A. Rapoport (1995) *EMBO J.* **14**, 217–223.
24. P. Whitley, E. Grahn, U. Kutay, T. A. Rapoport, and G. von Heijne (1996) *J. Biol. Chem.* **271**, 7583–7586.
25. J. Kyte and R. F. Doolittle (1982) *J. Mol. Biol.* **157**, 105–132.
26. D. Eisenberg, E. Schwarz, M. Komaromy, and R. Wall (1984) *J. Mol. Biol.* **179**, 125–142.
27. G. von Heijne (1986) *EMBO J.* **5**, 3021–3027.
28. E. Hartmann, T. A. Rapoport, and H. F. Lodish (1989) *Proc. Natl. Acad. Sci. USA* **86**, 5786–5790.
29. G. von Heijne (1990) *J. Mem. Biol.* **115**, 195–201.
30. G. von Heijne (1992) *J. Mol. Biol.* **225**, 487–494.
31. L. Sipos and G. von Heijne (1993) *Eur. J. Biochem.* **213**, 1333–1340.
32. B. Persson and P. Argos (1994) *J. Mol. Biol.* **237**, 182–192.
33. B. Persson and P. Argos (1997) *J. Protein Chem.* **16**, 453–457.
34. M. Cserzo, E. Wallin, I. Simon, G. von Heijne, and A. Elofsson (1997) *Protein Eng.* **10**, 673–676.
35. T. Schirmer and S. W. Cowan (1993) *Protein Sci.* **2**, 1361–1363.
36. K. Diederichs, J. Freigang, S. Umhau, K. Zeth and, J. Breed (1998) *Protein Sci.* **7**, 2413–2420.
37. H. P. Wessels, J. P. Beltzer, and M. Spiess (1991) *Methods Cell Biol.* **34**, 287–302.
38. C. Kast, V. Canfield, R. Levenson, and P. Gros (1996) *J. Biol. Chem.* **271**, 9240–9248.
39. C. Manoil, J. J. Mekalanos, and J. Beckwith (1990) *J. Bacteriol.* **172**, 515–518.
40. K. Yamane, S. Matsuyama, and S. Mizushima (1988) *J. Biol. Chem.* **263**, 5368–5372.
41. H. Andersson and G. von Heijne (1991) *Proc. Natl. Acad. Sci. USA* **88**, 9751–9754.
42. M. Johansson, I. Nilsson, and G. von Heijne (1993) *Mol. Gen. Genet.* **239**, 251–256.
43. S. Monier, P. Van Luc, G. Kreibich, D. D. Sabatini, and M. Adesnik (1988) *J. Cell Biol.* **107**, 457–470.
44. E. Szczesna-Skorupa, N. Browne, D. A. Mead, and B. Kemper (1988) *Proc. Natl. Acad. Sci. USA* **85**, 738–742.
45. E. Szczesna-Skorupa and B. Kemper (1989) *J. Cell Biol.* **108**, 1237–1243.
46. J. P. Beltzer, K. Fiedler, C. Fuhrer, I. Geffen, C. Handschin, H. P. Wessels, and M. Spiess (1991) *J. Biol. Chem.* **266**, 973–978.



47. G. D. Parks and R. A. Lamb (1991) *Cell* **64**, 777–787.
48. G. D. Parks and R. A. Lamb (1993) *J. Biol. Chem.* **268**, 19101–19109.
49. D. W. Andrews, J. C. Young, L. F. Mirels, and G. J. Czarnota (1992) *J. Biol. Chem.* **267**, 7761–7769.
50. H. Andersson and G. von Heijne (1994) *EMBO J.* **13**, 2267–2272.
51. G. Cao, A. Kuhn, and R. E. Dalbey (1995) *EMBO J.* **14**, 866–875.
52. W. van Klompenburg, I. Nilsson, G. von Heijne, and B. de Kruijff (1997) *EMBO J.* **16**, 4261–4266.
53. N. Johnsson and A. Varshavsky (1994) *EMBO J.* **13**, 2686–2698.
54. K. U. Kalies, D. Gorlich, and T. A. Rapoport (1994) *J. Cell Biol.* **126**, 925–934.
55. P. Whitley, I. M. Nilsson, and G. von Heijne (1996) *J. Biol. Chem.* **271**, 6241–6244.
56. A. J. Denzer, C. E. Nabholz, and M. Spiess (1995) *EMBO J.* **14**, 6311–6317.
57. H. Andersson and G. von Heijne (1993) *EMBO J.* **12**, 683–691.
58. G. Q. Cao and R. E. Dalbey (1994) *EMBO J.* **13**, 4662–4669.
59. P. Whitley, T. Zander, M. Ehrmann, M. Haardt, E. Bremer, and G. von Heijne (1994) *EMBO J.* **13**, 4653–4661.
60. T. Sato, M. Sakaguchi, K. Mihara, and T. Omura (1990) *EMBO J.* **9**, 2391–2397.
61. M. Sakaguchi, R. Tomiyoshi, T. Kuroiwa, K. Mihara, and T. Omura (1992) *Proc. Natl. Acad. Sci. USA* **89**, 16–19.
62. J. M. Wahlberg and M. Spiess (1997) *J. Cell Biol.* **137**, 555–562.
63. A. Eusebio, T. Friedberg, and M. Spiess (1998) *Exp. Cell Res.* **241**, 181–185.
64. N. G. Davis and P. Model (1985) *Cell* **41**, 607–614.
65. A. Saaf, E. Wallin, and G. von Heijne (1998) *Eur. J. Biochem.* **251**, 821–829.
66. T. Kuroiwa, M. Sakaguchi, K. Mihara, and T. Omura (1990) *J. Biochem. Tokyo* **108**, 829–834.
67. T. Kuroiwa, M. Sakaguchi, K. Mihara, and T. Omura (1991) *J. Biol. Chem.* **266**, 9251–9255.
68. G. Blobel (1980) *Proc. Natl. Acad. Sci. USA* **77**, 1496–1500.
69. H. P. Wessels and M. Spiess (1988) *Cell* **55**, 61–70.
70. J. Lipp, N. Flint, M. T. Haeuptle, and B. Dobberstein (1989) *J. Cell Biol.* **109**, 2013–2022.
71. G. Gafvelin, M. Sakaguchi, H. Andersson, and G. von Heijne (1997) *J. Biol. Chem.* **272**, 6119–6127.
72. E. Bibi, G. Verner, C. Y. Chang, and H. R. Kaback (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7271–7275.
73. M. Ehrmann and J. Beckwith (1991) *J. Biol. Chem.* **266**, 16530–16533.
74. K. McGovern and J. Beckwith (1991) *J. Biol. Chem.* **266**, 20870–20876.
75. K. Yamane, Y. Akiyama, K. Ito, and S. Mizushima (1990) *J. Biol. Chem.* **265**, 21166–21171.
76. G. Gafvelin and G. von Heijne (1994) *Cell* **77**, 401–412.
77. B. Martoglio, M. W. Hofmann, J. Brunner, and B. Dobberstein (1995) *Cell* **81**, 207–214.
78. H. Do, D. Falcone, J. Lin, D. W. Andrews, and A. E. Johnson (1996) *Cell* **85**, 369–378.
79. W. Mothes, S. U. Heinrich, R. Graf, I. Nilsson, G. von Heijne, J. Brunner, and T. A. Rapoport (1997) *Cell* **89**, 523–533.

## Topoisomer

The discovery of ring-shaped DNA molecules in the 1960s pointed to a connection between the topology of links and knots, subjects that had long been cherished in pure mathematics, and the molecular biology of DNA, by then a blossoming discipline following the unveiling of the **double-helix** structure in 1953. It became clear that the topology of DNA is deeply rooted in its double-helix structure. The discovery of **supercoiled** DNA in 1965 (1) laid the cornerstone for this realization; and the discovery of the DNA **topoisomerases**, [enzymes](#) that solve the topological problems of DNA (2), showed the biological importance of DNA topology. Topological problems arise in nearly all transactions of intracellular DNA, and Nature invented not one, but three or more, distinct subfamilies of enzymes to solve these problems. It also became clear that the topological problems of intracellular DNA did not originate from its circularity in some organisms, but from its extreme length in nearly all organisms.

## 1. Topoisomers (Topological Isomers)

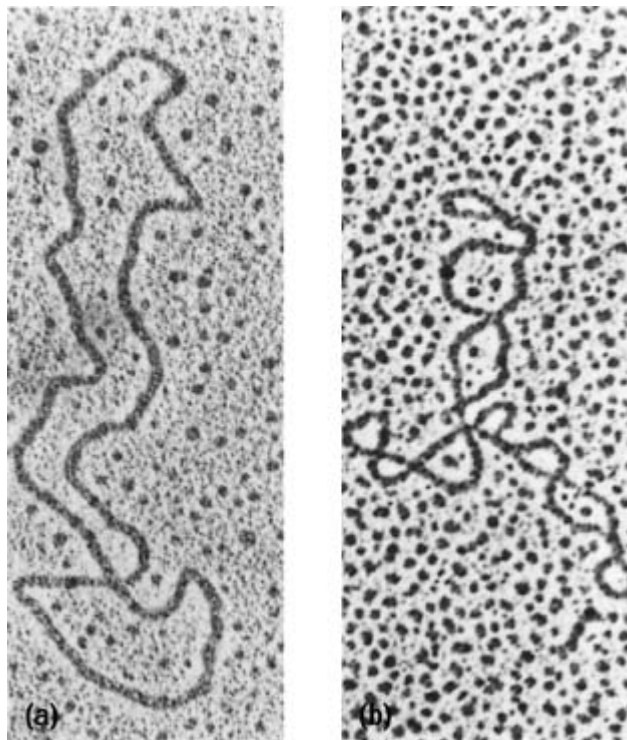
### 1.1. Supercoiled DNA

Several novel forms of DNA rings were identified following the discovery of circular DNA, of which the finding of supertwisted or supercoiled DNA in 1965 (1) is of particular importance. The supercoiling of a double-stranded DNA ring is best described in terms of a topological quantity called the **linking number**. Because of the double-helical structure of DNA, a DNA ring with two intact strands is comprised of a pair of multiply intertwined single-stranded rings. The pair of strands revolve around each other, and they cannot be separated without a break in at least one of them. The circular strands are thus topologically linked, and the order or degree of this topological linkage is specified by the linking number,  $Lk$ . The linking number can be defined as the minimal number of strand passage events that are needed to unlink the two strands of a DNA ring. In each event, one strand of the duplex DNA is transiently broken, the other strand is passed through the break, and the broken strand is then resealed. For a typical DNA in its most stable structure, called a *relaxed* DNA, the two complementary strands make one right-handed helical turn every 10.5 base pairs under physiological conditions; the linking number  $Lk^0$  of a relaxed DNA ring  $N$  base pairs in length is thus  $N/10.5$  under physiological conditions (3). In general, the helical structure of a DNA, and therefore the value of  $Lk^0$  of a DNA of a given size, is dependent on its nucleotide sequence, as well as ambient conditions. A DNA segment of alternating cytosines and guanines, for example, may assume a left-handed double-helix structure under certain conditions, and such a segment would contribute a negative quantity to the overall  $Lk^0$  of the DNA ring under these conditions.

The difference between  $Lk$  and  $Lk^0$ , called the *linking difference*, gives a measure of how far the topological state of a DNA ring deviates from its relaxed state. If the magnitude of  $Lk - Lk^0$  differs significantly from zero, the DNA ring would be torsionally and flexurally strained. Such a strained molecule would assume a distorted shape, in much the same way as a torsionally unbalanced multistranded rope (Fig. 1). The descriptions “supercoiled,” “superhelical,” or “supertwisted” have been used interchangeably to describe such a molecule. Because both  $Lk$  and  $Lk^0$  are dependent on the size of a DNA ring, the normalized quantity  $s = (Lk - Lk^0)/Lk^0$ , called the *specific linking difference*, is a convenient measure of the extents of supercoiling of DNA rings of different sizes. DNA rings that differ only in their linking numbers are topological isomers, or topoisomers. A topoisomer is negatively supercoiled if the pair of strands in it are significantly underwound ( $Lk < Lk^0$  and  $Lk - Lk^0 < 0$ ) and positively supercoiled if they are overwound ( $Lk > Lk^0$  and  $Lk - Lk^0 > 0$ ). DNA rings purified from bacteria and eukarya are usually negatively supercoiled, with values of  $s$  around  $-0.06$ .

**Figure 1.** Electron micrographs showing a relaxed (a) and a supercoiled (b) DNA molecule. The size of the DNA is

about 10,000 base pairs. To facilitate viewing, the apparent diameter of the double helix was made several times larger by coating the molecules with a protein. (From Ref. [43](#), with permission.)

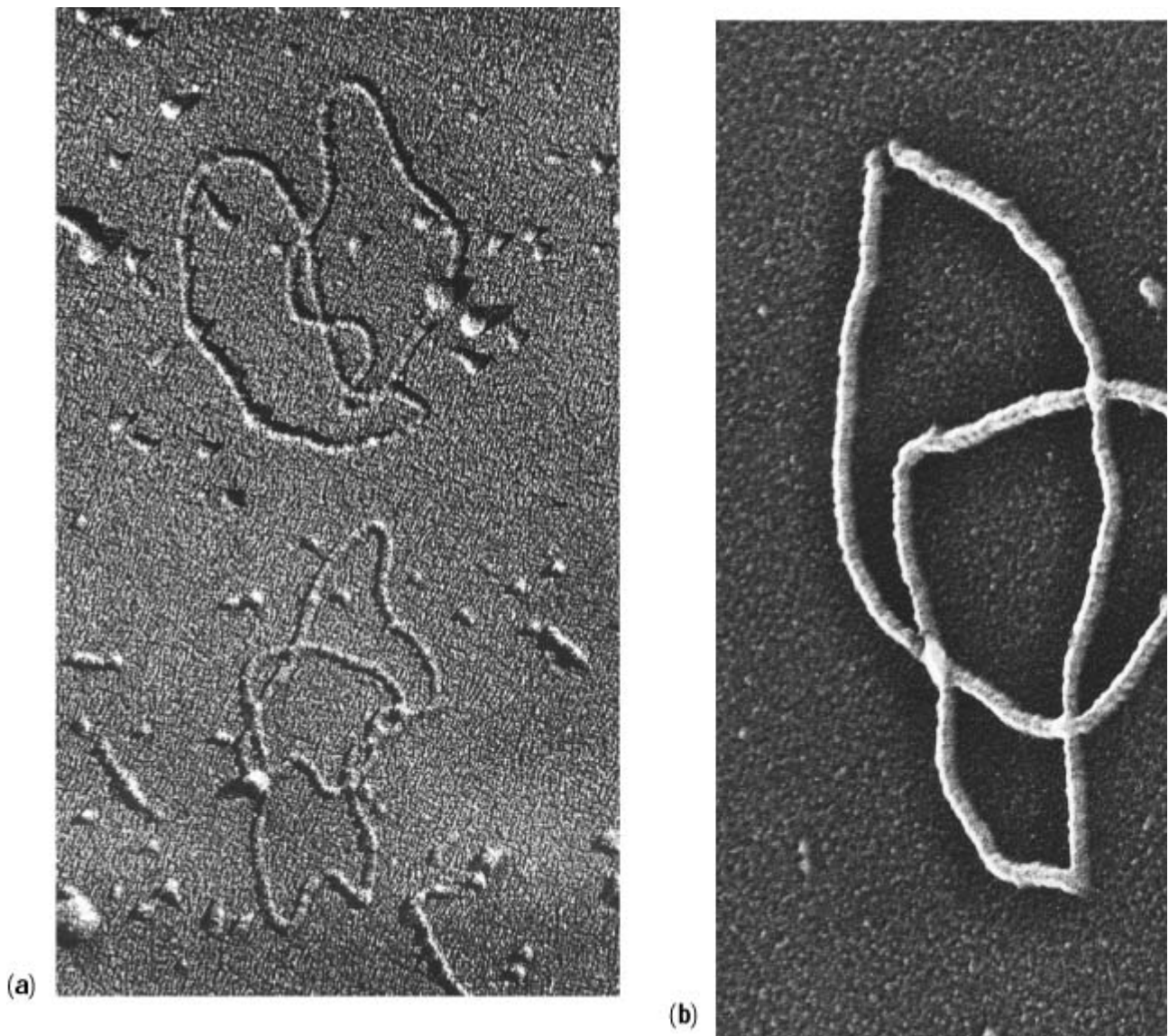


Supercoiled DNA is distorted in two major ways. First, the helical twist of the DNA is altered because of the topological strain. Second, there is a distortion of the overall shape of the ring, with the helical axis of the DNA taking on a spatial **writhe** (see Fig. [1](#)). These two types of distortions are coupled, and their sum is determined by the linking difference ([3](#), [4](#)). Because of these distortions, a supercoiled DNA has unique properties relative to its relaxed topoisomer ([4-8](#)). A local unwinding of two strands of a DNA, for example, can be accomplished much more easily in a negatively supercoiled than in a relaxed DNA. Negative supercoiling of a DNA would, for example, facilitate the binding of proteins that preferentially bind to single-stranded or underwound duplex DNA, proteins such as the [histone](#) octamer complex that require left-handed wrapping of DNA around them, and planar aromatic molecules that **intercalate** in between DNA base pairs. Negative supercoiling of a DNA could also drive changes in the DNA helical structure that are otherwise thermodynamically unfavorable under physiological conditions, such as a change of an alternating CG sequence from the right-handed B-helical form to the left-handed Z-helical form (see [Z-DNA](#)), and the extrusion of a pair of hairpin structures from a **palindromic** sequence ([cruciform](#) formation). Positive supercoiling of a DNA, on the other hand, would disfavor processes that are favored by negative supercoiling. Both negative and positive supercoiling may also facilitate the binding of proteins to DNA crossovers or to bent DNA segments.

## 1.2. Catenanes and Knots

Catenanes and knots are also topoisomers of simple rings (Figs. [2a](#) and [2b](#)). Catenanes, or interlocked DNA rings, were prepared *in vitro* and found *in vivo* in 1967 ([9](#), [10](#)). Inside the cells, they can be formed by [DNA replication](#), if the strands of the final stretch of unreplicated parental DNA come apart before their complete unraveling ([11](#)), or by [recombination](#), if the joining of two distant sites in a DNA ring traps some of the DNA crossovers in between ([12](#)). Knotted DNA rings were discovered as the products of certain reactions *in vitro* in the 1970s ([13](#), [14](#)) and subsequently found in **plasmid** preparations ([15](#)).

**Figure 2.** (a) An electron micrograph showing two dimeric catenanes of simian virus 40 DNA. The viral DNA was isolated from cells grown under conditions that inhibit DNA topoisomerase II. (From Ref. 5, with permission; micrograph was originally from A. Varshavsky.) (b) An electron micrograph of a knotted DNA ring. (From Ref. 5, with permission; micrograph was originally from A. Stasiak.)



The discovery of DNA supercoils, catenanes, and knots helped to bring into focus the topological problems of DNA. How are the parental strands of a DNA ring unlinked during replication? Why are DNA rings isolated from cells negatively supercoiled? If catenanes can form by replication or recombination, how do they come apart? These topological problems are not limited to ring-shaped DNA molecules. [Chromosomes](#) in eukaryotic cells, for example, are linear in the topological sense because the ends of a chromosome are not joined. The chromosomes are probably organized into many loops inside the cells, however, and are therefore subject to topological constraints in much the same way as DNA rings. The extreme length of a chromosome also makes it difficult to disentangle the molecules through movements of their ends. The discovery of DNA rings in various topological forms greatly simplified the study of DNA topology, however, and provided the necessary tool in the search of enzymes, called the DNA *topoisomerases*, that can interconvert DNA topoisomers and solve all topological problems of intracellular DNA (see text below).

## 2. DNA Topoisomerases

### 2.1. Classification

DNA topoisomerases are [enzymes](#) that catalyze the interconversion of DNA topoisomers. In their presence, the DNA strands in a duplex DNA, or different duplex DNA segments in the same or different molecules, can move through each other. Although the discovery of this family of enzymes started around 1970 with the partial purification of an activity from *Escherichia coli* cell extracts that can remove negative supercoils ([16](#)), the term topoisomerase was not coined until nearly a decade later ([17](#)).

The DNA topoisomerases can be divided into two types ([2](#)). The type I enzymes catalyze linking number changes of topoisomers by introducing a transient break in one strand, passing the other strand through this break, and then rejoining the interrupted strand. They can also change the topology of double-stranded DNA catenanes or knots, provided that at least one of the strands contains a gap or nick. The type II enzymes catalyze the transport of one DNA double helix through another, using the transient breakage of both strands in one of them. These enzymes can interconvert linking number topoisomers, catenanes, or knots, independent of the presence of a gap or nick. In the interconversion between linking number topoisomers, a consequence of passing one duplex DNA segment through a transient break in another segment of the same DNA is that the linking number of the DNA is altered by two ([18](#), [19](#)). In linking number changes mediated by type I topoisomerases, each DNA strand passage event alters  $Lk$  by unity. The change of  $Lk$  in steps of two is a signature of type II DNA topoisomerases.

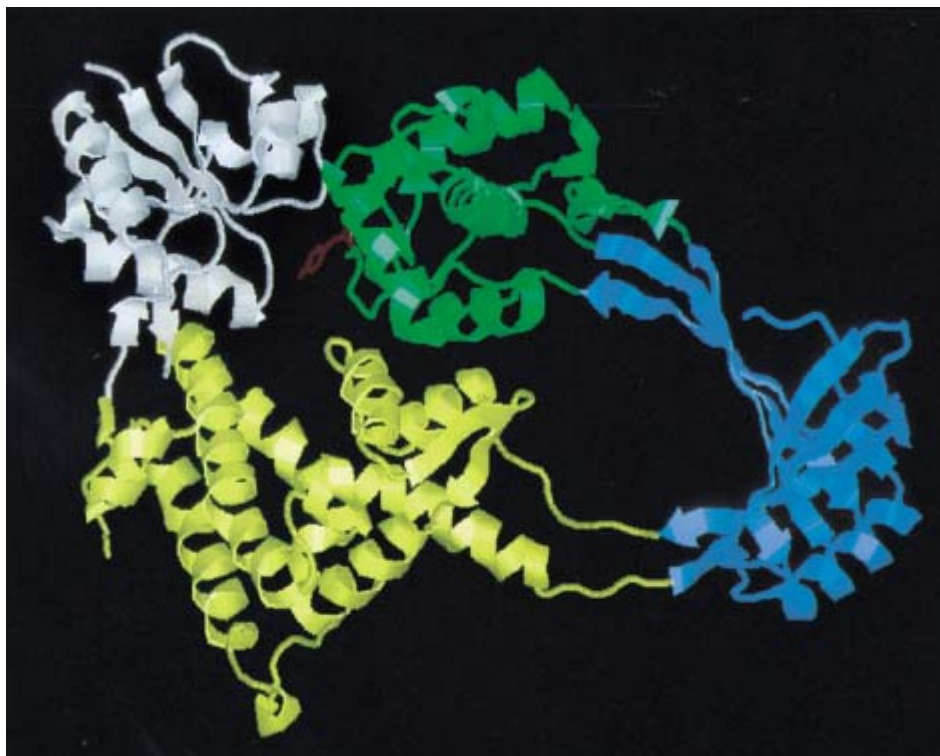
The type I enzymes can be further divided into the IA and IB subfamilies ([2](#)). The former includes bacterial DNA topoisomerases I and III, eukaryotic DNA topoisomerase III, and an enzyme “reverse gyrase” found in archaea; the latter includes eukaryotic DNA topoisomerase I and pox virus topoisomerase. The amino acid sequences and the biochemical properties of the type IA enzymes are very different from those of the IB enzymes. The type II DNA topoisomerases are found in bacteria, eukarya, and archaea. They were thought to form a single homologous subfamily, but recent work suggests that, in addition to the archetype type II enzymes, there is a new subfamily of enzymes in archaea, and perhaps in eukarya as well ([20](#)). These subfamilies are referred to as type IIA and IIB, respectively. *E. coli* has four topoisomerases: two type I enzymes DNA topoisomerases I and III, and two type II enzymes DNA gyrase (DNA topoisomerase II) and DNA topoisomerase IV. Bacterial gyrase is unique among type II DNA topoisomerases in that it can catalyze the negative supercoiling of a relaxed DNA ring; other type II enzymes can relax positively and negatively supercoiled DNA, but cannot supercoil a relaxed DNA. The budding yeast *S. cerevisiae* has three topoisomerases excluding any plausible type IIB enzymes: two type I enzymes DNA topoisomerases I and III, and one type IIA enzyme DNA topoisomerase II. Other eukaryotes are similar to yeast, but there may be more than one member in each subfamily. In humans, two type IA, two type IB, and two type IIA **isozymes** are known. Genes encoding six human enzymes, DNA topoisomerases I, IIA, IIB, DNA topoisomerases IIIA, IIIB, have been located to chromosome regions 20q12-13.2, 17q21-22, 3p24, and 17p11.2-12 22q11-12, respectively ([2](#), [21](#), [22](#)).

### 2.2. Mechanisms

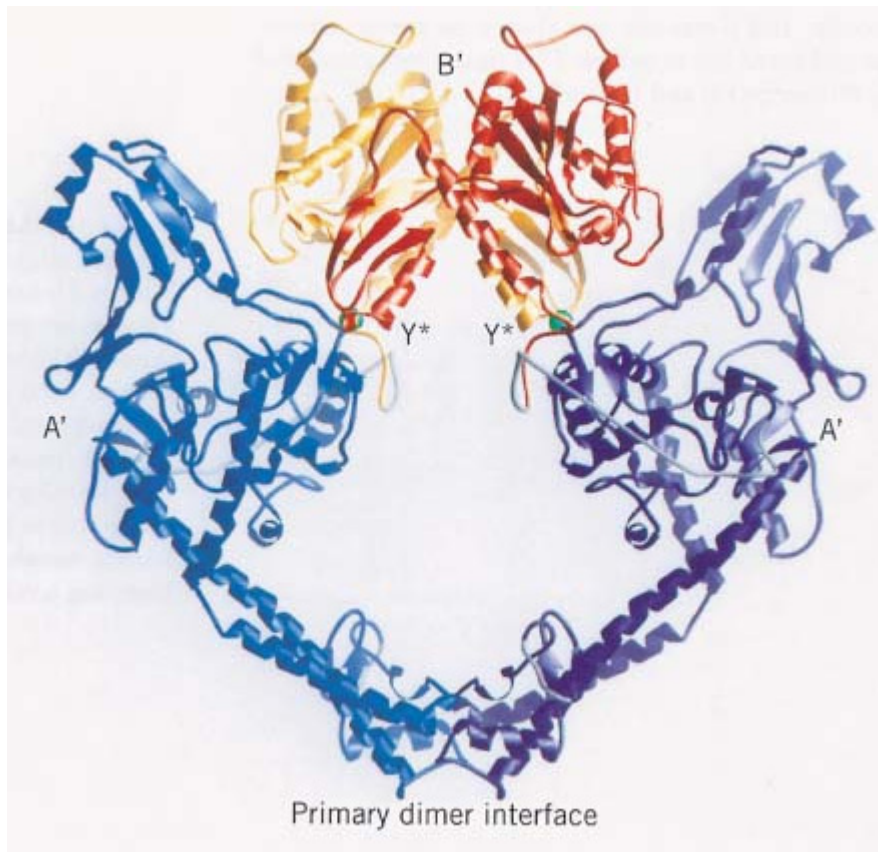
How do the DNA topoisomerases perform their magic? First, they must catalyze the transient breakage of DNA strands. It was postulated in 1971 that this transient breakage involved transesterification between an enzyme hydroxyl group and a DNA internucleotide phosphodiester bond ([16](#)). Nucleophilic attack of the enzyme hydroxyl on a DNA phosphorous would form a covalent enzyme–DNA bond and break a DNA phosphodiester bond, generating a deoxyribose hydroxyl group in the process. Rejoining of the broken DNA would then ensue by a second transesterification reaction, which could be the microscopic reversal of the first: The deoxyribose hydroxyl group would attack the phosphorous in the covalent enzyme–DNA link, rejoining the DNA strand and breaking the enzyme–DNA link. This postulate has been confirmed, and DNA cleavage by all DNA topoisomerases is found to go through a covalent enzyme–DNA intermediate in which a DNA phosphoryl group is linked to the phenolic oxygen of an enzyme [tyrosine](#) residue. For the type

IA and type II enzymes, in the covalent intermediate the tyrosyl residue is linked to a DNA 5'-phosphoryl group; for the type IB enzymes, a DNA 3'-phosphoryl group. The particular tyrosines in a number of DNA topoisomerases have been identified (2); and the positions of the tyrosines in several enzymes, such as in the type IA enzyme *Escherichia coli* DNA topoisomerase I (23) and a type II enzyme *Saccharomyces cerevisiae* DNA topoisomerase II (24), can be seen in their three-dimensional structures determined by [X-ray crystallography](#) (Figs. 3 and 4; see also (28, 29). Side chains in addition to the [active-site](#) tyrosines are also involved in stabilizing the attacking and leaving groups during covalent catalysis.

**Figure 3.** A ribbon illustration of the crystal structure of a 67-kDa *N*-terminal fragment of *E. coli* DNA topoisomerase I. Drawing is based on the coordinates of the structure reported in Ref. 21. This fragment can bind and cleave single-stranded DNA, but cannot remove supercoils from a duplex DNA ring (21). The four domains I to IV are colored in white, blue, green, and yellow, respectively; the active site tyrosine, Tyr319, is located in domain III and shown in red. See color insert.

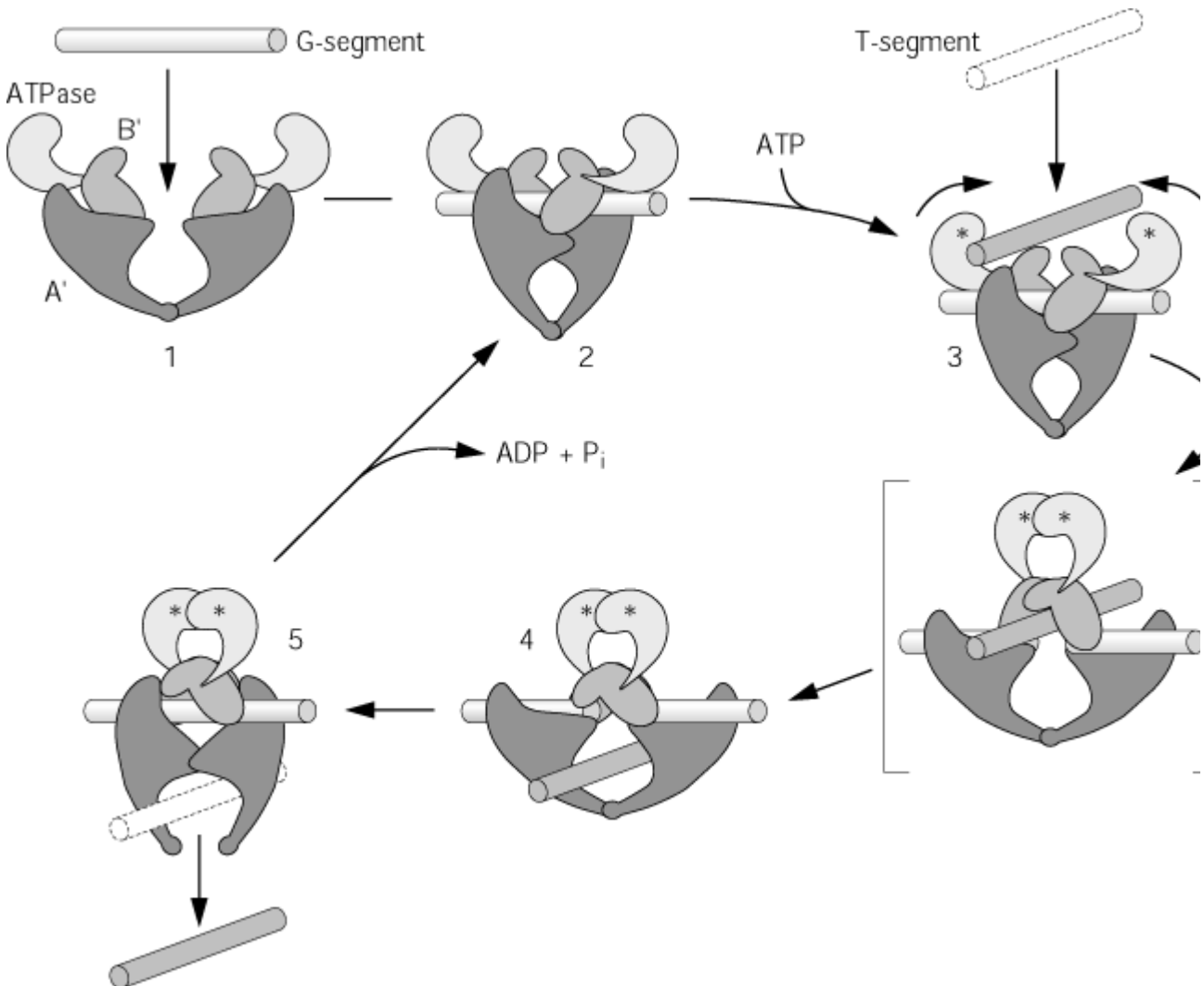


**Figure 4.** A ribbon illustration of the crystal structure of a 92-kDa fragment of yeast DNA topoisomerase II. This fragment is comprised of amino acid residues 410 to 1202 of the intact enzyme. It can bind and cleave double-stranded DNA, but lacks the ATPase domain located in the *N*-terminal region of the intact enzyme and cannot catalyze DNA transport (23). The pair of active-site tyrosines in the dimeric enzyme are represented by two green spheres labeled Y\*. A' and B' refer to two regions of the polypeptide, from residues around 660 to 1202 and from 410 to around 660, respectively; the A' fragment corresponds to the *N*-terminal two-thirds of the A-subunit of bacterial gyrase, and the B' fragment corresponds to the *C*-terminal half of the B-subunit of bacterial gyrase. (From Fig. 3a of Ref. 24, with permission.) See color insert.



Second, each enzyme must mediate the passage of DNA strands. Recent biochemical and three-dimensional structural studies of yeast DNA topoisomerase II have led to a molecular model of how a single type II enzyme molecule couples ATP binding and hydrolysis to the transport of one DNA double helix through another (24, 25). This is illustrated in Fig. 5. The enzyme acts as a dimeric protein clamp, with each half composed of several fragments denoted by ATPase, B' and A' (drawing 1 of Fig. 5; the ATPase fragment is composed of the *N*-terminal portion of the polypeptide from amino acids 1 to 409, the B' fragment from 410 to about 660, and the A' fragment from around 660 to 1202). Either in the free form (1), or when bound to a DNA segment called the G-segment (2), the clamp can close and open through formation and disruption, respectively, of contacts between the pair of ATPase domains. The opening and closing of this protein gate, termed the *N*-gate because of the proximity of the ATPase domain to the *N*-terminus of the polypeptide, are modulated by ATP binding and hydrolysis: ATP-binding closes the *N*-gate, which reopens after ATP hydrolysis and the release of the hydrolyzed products (ADP and orthophosphate). When the *N*-gate is in the open state, a second DNA-segment, called the T-segment, can enter the gate (3). This entrance is probably facilitated by weak but favorable enzyme-DNA interaction. When the clamp closes upon ATP-binding, a T-segment that has entered the clamp is trapped. Steric repulsion between the trapped T-segment and the closed ATPase domains is believed to force open the enzyme-mediated DNA-gate in the G-segment, as well as expel the T-segment into the large cavity below the T-segment (3 and 4). This expulsion relieves the steric repulsion, and the DNA gate returns to the more stable closed state. The closing of the DNA-gate contracts the size of the hole containing the T-segment, however, and thus forces the T-segment to exit through a second protein gate, called the C-gate, formed by amino acid residues 1031–1046 and 1114–1130 of the yeast enzyme (5). Following ATP hydrolysis and product release, the closed enzyme clamp reopens (2) to complete a cycle of DNA transport. All type IIA enzymes are likely to catalyze DNA transport in a similar way.

**Figure 5.** The transport of one DNA double helix through another by a type IIA DNA topoisomerase. The macromolecules involved are drawn to scale in this illustration. Each half of the dimeric enzyme is composed of three parts denoted by ATPase, which represents the *N*-terminal 409 amino acid stretch containing the ATPase domain, and B' and A' are the fragments defined in the legend to Figure 4. The DNA G-segment, within which a transient break is to be introduced to open a gate in it, and T-segment, which is to be transported through the DNA gate, are shown as rods. The asterisk (\*) represents a bound ATP. See text for a description of the model. (From Fig. 5 of Ref. 24, with permission.)



The unique ability of bacterial gyrase to catalyze negative supercoiling can be attributed to its ability to orient the T-segment relative to the G-segment: In the complex between gyrase and DNA, a 140-k long DNA segment is wrapped right-handedly around the enzyme (26), which positions the T-segment for transport through the G-segment in a directional way. For the other members of the type IIA subfamily, the relative orientation of the G- and T-segment is determined by the DNA conformation, which always favors the removal of supercoils (2, 27).

How the type IA enzymes mediate DNA strand passage is less clear. In the crystal structure of a large fragment of the type IA enzyme *E. coli* DNA topoisomerase I (Fig. 3), the polypeptide is found to fold into four distinct domains. Domains I and IV form the “base” of the structure, which is connected through a pair of long strands to the “lid” formed by domains II and III (23). The active site tyrosine, Tyr319, is located in domain III at the junction between this domain and the base. This strategic location suggests that there are likely large movements between the lid and the base during the various steps of strand cleavage and passage (23). Following the cleavage of a DNA strand, for example, domain III with the covalently attached 5'-end of the broken DNA would presumably be lifted away from the base to generate a gap between the domains sufficiently large for the passage of another strand or perhaps even a double-helical DNA segment. After strand passage, this gap would



have to close again to allow rejoining of the cleaved strand, and further movements of the domains would be needed for the enzyme to return to a state ready for the next cycle of reaction. The molecular details of these movements are unknown. Unlike the reaction catalyzed by the type II enzymes, in which DNA transport is coupled to ATP binding and hydrolysis, DNA strand transport in reactions catalyzed by the type IA enzymes does not involve a high-energy cofactor (with the exception of “reverse gyrase”). Therefore for the majority of the type IA enzymes, changes in enzyme–DNA interactions in the various steps must play crucial roles in the ordered processing of the reaction (23). Conceptually, instead of viewing the enzyme as the mover and shaker and the DNA as the passive substrate, both should be thought of as integral parts in the orchestration of the intricate reaction steps.

For the type IB enzymes, the solution of the crystal structure of a large fragment of a pox virus topoisomerase (28), and several crystal structures of catalytically active fragments of human DNA topoisomerase I in complex with DNA (29), has provided much insight into their mechanism. The type IB enzymes are structurally similar to a family of tyrosine recombinases including phage I integrase and a phage P1 site-specific recombinase called Cre. In the human DNA topoisomerase I–DNA complexes, the bilobed protein clamps around the duplex DNA and the contacts mostly DNA phosphate groups (29). Based on these structures and earlier biochemical studies, a detailed mechanism of the type IB enzymes has emerged (28-31).

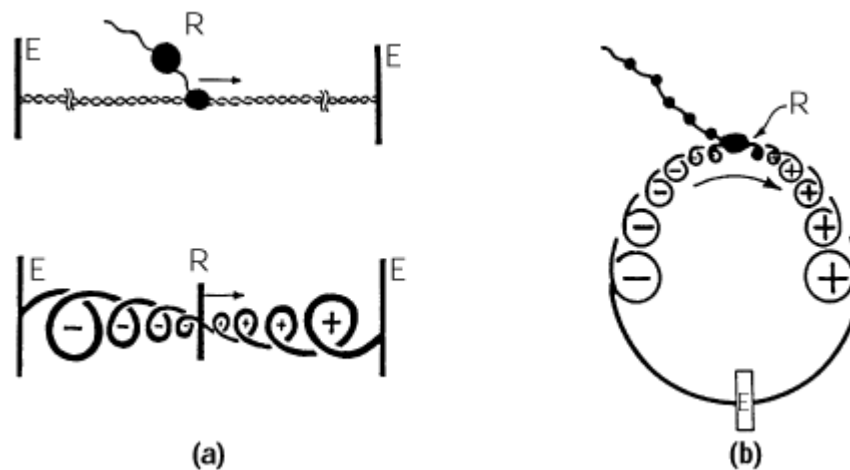
An interesting aspect of DNA strand passage by a type I enzyme is the number of passages that can occur in each strand breakage and rejoining event (18, 32, 33). This issue is closely related to the state of the transiently cleaved DNA strand: One end of the broken strand is attached covalently to the active-site tyrosine, but what about the other end? In the enzyme-bridging model, the other end is also postulated to be enzyme-bound through noncovalent interactions. In this model, each transient cleavage would allow only one strand passage event. Strand passage would be mediated by movements of the enzyme–DNA complex, and the cleaved strand would have to be rejoined before either end can dissociate from the enzyme. In the DNA rotation model, in the covalent intermediate the enzyme is closely contacting only the side of DNA containing the covalently attached DNA end. The other side of the DNA could therefore rotate around a single bond in the uncleaved strand across the cleaved bond. In the latter model, the enzyme would have a passive role in strand passage, and multiple strand passage events could be achieved by multiple rotations of one side of the cleaved DNA relative to the other side. To relax a supercoiled DNA with a linking difference of  $Lk - Lk^0$ , the enzyme-bridging model would require  $Lk - Lk^0$  events of strand breakage and rejoining, but the DNA rotation model could require as few as one single event. Biochemical and kinetic experiments suggest that the type IA enzymes are likely to follow the enzyme-bridging mechanism (34) and that the type IB enzymes are likely to follow the DNA rotation mechanism (35). The crystal structures of human DNA topoisomerase I in complex with DNA are also suggestive that the type IB enzymes can allow multiple strand passage events during each round of DNA breakage and rejoining (30).

### 3. Biological Ramifications of DNA Topology

The topological properties of DNA rings are rooted in the double-helix structure of DNA. One major cause of catenane formation, for example, is the incomplete removal of helical intertwinings between the parental strands of a replicating DNA (11). Supercoiling of intracellular DNA is also closely tied to the translocation of macromolecular assemblies along a helical DNA template (36). In Fig. 6, the movement of a transcribing **RNA polymerase** along its template is illustrated. The helical geometry of the template would require that the transcriptional ensemble, including the RNA polymerase, the nascent transcript attached to it, and proteins bound to the nascent transcript, be rotated around the DNA as transcription proceeds. But what if a transcription ensemble is prevented from rotating around its template? For example, in bacteria, where [transcription](#) and [translation](#) are coupled, nascent polypeptides of a membrane protein would be made before the termination of the [messenger RNA](#). These mRNA-attached polypeptides could become attached to the cell membrane and would thus prevent the entire transcriptional assembly from rotating around the DNA template. Interaction

between RNA polymerase and template-bound regulatory proteins could also prevent the rotation of the transcription ensemble around the DNA. In either case, the DNA would be forced to rotate relative to the transcription ensemble. The DNA in front of the moving polymerase would become overwound or positively supercoiled, and the DNA behind the polymerase would become underwound or negatively supercoiled (36).

**Figure 6.** Supercoiling of the DNA template by transcription (36). (a) Because of the double-helix geometry of the DNA, a transcriptional ensemble R (including the RNA polymerase, the nascent transcript attached to the polymerase, and proteins associated with the transcript) would rotate around the DNA as transcription proceeds (upper drawing). If the ensemble is prevented from turning around the DNA, owing to, for example, the anchoring of the ensemble through the association of the polymerase with a template-bound regulatory protein or through insertion of mRNA-associated polypeptides into membranes, then the DNA would be forced to turn around its helical axis. In such a case, positive supercoils (+ signs) would be generated ahead of the ensemble and negative supercoils (– signs) behind it (lower drawing). In (a), the ends of the DNA segment are shown to be anchored to hypothetical elements denoted by E. In general, E represents a barrier to the rotation of DNA around its helical axis, as illustrated in (b) for a circular DNA template. In bacteria, gyrase appears to be primarily responsible for the removal of the positive supercoils and DNA topoisomerase I the negative supercoils; in eukaryotes, DNA topoisomerase I or II can remove both. (From Ref. 48, with permission.)



The close link between DNA topology and the double-helix structure of DNA also foretells the influence of DNA topology, and thus the involvement of DNA topoisomerases, in nearly all aspects of cellular transactions of DNA. In replication, for example, these enzymes are needed for the separation of the parental strands. In *E. coli*, the action of gyrase, perhaps augmented by the other topoisomerases, is necessary for the unlinking of the strands of the ring-shaped chromosome. In yeast propagation of the replication forks requires DNA topoisomerase I or II. A type II enzyme is also needed to unlink DNA catenanes or pairs of intertwined chromosomes at the time of their segregation (2). In *E. coli*, this role is mainly fulfilled by DNA topoisomerase IV; in yeast, by DNA topoisomerase II.

In transcription, one or more DNA topoisomerases is involved. In bacteria, gyrase is believed to be the enzyme primarily responsible for the removal of the positive supercoils ahead of the transcription apparatus, and DNA topoisomerase I is believed to be responsible for the removal of the negative supercoils behind (36). Excessive supercoiling is likely to pose problems for cellular processes. In *E. coli*, inactivation of DNA topoisomerase I stimulates hybrid formation between the nascent RNA and the template strand of DNA, which is presumably a consequence of excessive negative supercoiling (37). This kind of **R-loop** formation is detrimental to the cell, especially in strains lacking **ribonuclease H**, which hydrolyzes the RNA strand in the hybrid. Participation of eukaryotic DNA topoisomerases in **gene regulation**, through their association with other proteins involved in the

initiation of transcription, has also been suggested (2).

Replication and transcription are two of the better-studied cellular processes involving one or more DNA topoisomerases. These enzymes have also been implicated in chromosome condensation and decondensation, recombination, the maintenance of [genome](#) stability, and [DNA repair](#). Few other enzymes are involved in so many important and different cellular processes.

#### 4. Topoisomerase-Targeting Toxins, Antibiotics, and Anticancer Agents

Shortly after the discovery of bacterial gyrase, two different classes of **antibiotics** represented by **nalidixic acid** and coumermycin were found to target the enzyme (38-40). Extensive studies showed that nalidixate, a member of the fluoroquinolone class of antibiotics that includes ciprofloxacin and norfloxacin, acts by trapping the covalent complex between DNA and gyrase, and coumermycin acts by inhibiting the gyrase [ATPase](#) activity. These drugs have similar effects on bacterial DNA topoisomerase IV, and their primary target in a bacterium could be either one or both of the two type IIA DNA topoisomerases, depending on the particular bacterium and the particular drug (41, 42). The quinolone class of drugs are widely used. The trapping of the covalent enzyme–DNA complex appears to lead to cell killing through arrest of cell division. In addition to these antibiotics, two very different plasmid-encoded toxins were also found to target bacterial gyrase. The *ccdB* or *letD* gene product of the *E. coli* F plasmid appears to interact with the A-subunit of gyrase to prevent DNA from rejoining after the formation of the gyrase–DNA covalent intermediate (43). This cytotoxicity facilitates the stable inheritance of the plasmid, because it also encodes a labile antidote of the toxin. Another gyrase-targeting toxin, microcin B17, has been identified in strains of *E. coli* bearing one of group of naturally occurring plasmids. This toxin is first expressed as a 69-residue peptide, which is then extensively processed to give a 43-residue product with four thiazole and four oxazole rings (44).

In the 1980s, a number of antitumor agents were shown to target eukaryotic DNA topoisomerase II by trapping the covalent DNA–enzyme complex (45, 46). These agents include the clinically useful drugs doxorubicin, etoposide, and mitoxantrone. A plant alkaloid, camptothecin, was found in 1985 to trap the covalent complex between DNA and eukaryotic DNA topoisomerase I, and several of the camptothecin derivatives have since been developed for cancer treatment. Topoisomerase-targeting cytotoxic agents that act by inhibiting one or more topoisomerases, rather than trapping a covalent intermediate, such as the gyrase inhibitor coumermycin and the eukaryotic DNA topoisomerase II inhibitors bisdioxopiperazines, are presently of limited clinical use but are of potential in future developments of anticancer, antibacterial, antifungal, and antiparasitic therapeutics.

#### Bibliography

1. J. Vinograd, J. Lebowitz, R. Radloff, R. Watson, and P. Lapis (1965) *Proc. Natl. Acad. Sci. USA* **53**, 1104–1111.
2. J. C. Wang (1996) *Annu. Rev. Biochem.* **65**, 635–692.
3. J. C. Wang (1980) *Trends Biochem. Sci.* **5**, 219–221.
4. N. R. Cozzarelli, T. C. Boles, and J. H. White (1990) In *DNA Topology and Its Biological Effects* (N. R. Cozzarelli and J. C. Wang, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 139–184.
5. J. C. Wang (1986) In *Cyclic Polymers* (J. A. Semlyen, ed.), Elsevier, New York, pp. 225–260.
6. M. D. Frank-Kamenetskii (1990) In *DNA Topology and Its Biological Effects* (N. R. Cozzarelli and J. C. Wang, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 185–215.
7. A. V. Vologodskii and N. R. Cozzarelli (1994) *Annu. Rev. Biophys. Biomol. Struct.* **23**, 609–643.
8. A. V. Vologodskii and N. R. Cozzarelli (1994) *Curr. Opin. Struct. Biol.* **4**, 372–375.
9. J. C. Wang and H. Schwartz (1967) *Biopolymers* **5**, 953–966.

10. B. Hudson and J. Vinograd (1967) *Nature* **216**, 647–652.
11. O. Sundin and A. Varshavsky (1980) *Cell* **21**, 103–114.
12. S. A. Wasserman and N. R. Cozzarelli (1986) *Science* **232**, 951–960.
13. L. F. Liu, R. E. Depew, and J. C. Wang (1976) *J. Mol. Biol.* **106**, 439–452.
14. L. F. Liu, L. Perkocho, R. Calendar, and J. C. Wang (1981) *Proc. Natl. Acad. Sci. USA* **78**, 5498–5502.
15. K. Shishido, N. Komiyama, and S. Ikawa (1987) *J. Mol. Biol.* **196**, 215–218.
16. J. C. Wang (1971) *J. Mol. Biol.* **55**, 523–533.
17. J. C. Wang and L. F. Liu (1979) In *Molecular Genetics*, Part III (J. H. Taylor, ed.), Academic Press, New York, pp. 65–88.
18. P. O. Brown and N. R. Cozzarelli (1979) *Science* **206**, 1081–1083.
19. L. F. Liu, C. C. Liu, and B. M. Alberts (1980) *Cell* **19**, 697–707.
20. A. Bergerat, B. de Massy, P. C. Varoutas, A. Nicolas, and P. Forterre (1997) *Nature (London)* **386**, 414–417.
21. K. Kawasaki, S. Minoshima, E. Nakato, K. Shibuya, A. Shintani, J. L. Schmeits, J. Wang, and N. T. Shimizu (1997) *Genome Res.* **7**, 250–261.
22. H. Zhang J. M. Barcelo B. Lee G. Kohlhagen D. B. Zimonjic N. C Popsu Y. Pommier (2001) *Proc. Natl. Acad. Sci. USA* **98**, 10608–10613.
23. C. Lima, J. C. Wang, and A. Mondragón (1994) *Nature* **367**, 138–146.
24. J. M. Berger, S. J. Gamblin, S. C. Harrison, and J. C. Wang (1996) *Nature* **379**, 225–232.
25. J. Roca, J. M. Berger, S. C. Harrison, and J. C. Wang (1996) *Proc. Natl. Acad. Sci. USA* **93**, 4057–4062.
26. L. F. Liu and J. C. Wang (1978) *Cell* **15**, 979–984.
27. J. Roca and J. C. Wang (1996) *Genes to Cells* **1**, 17–27.
28. C. Cheng, P. Kussie, N. Pavleitch, and S. Schman (1998) *Cell* **92**, 841–850.
29. M. R. Redinbo, L. Stewart, P. Kuhn, J. J. Champoux, and W. G. Hol (1998) *Science* **279**, 1504–1513.
30. J. J. Champoux (2001) *Annu. Rev. Biochem.* **70**, 369–413.
31. B. O. Krogh and S. Shuman (2000) *Mol. Cell.* **5**, 1035–1041.
32. J. C. Wang (1982) In *Nucleases* (R. Roberts and S. Linn, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 41–57.
33. J. J. Champoux (1990) In *DNA Topology and Its Biological Effects* (N. R. Cozzarelli and J. C. Wang, eds.), Cold Spring Harbor Laboratory, Cold Spring Harbor, New York, pp. 217–242.
34. K. A. Kirkegaard and J. C. Wang (1985) *J. Mol. Biol.* **185**, 625–637.
35. J. T. Stivers, S. Shuman, and A. S. Mildvan (1994) *Biochemistry* **33**, 327–339.
36. L. F. Liu and J. C. Wang (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7024–7027.
37. M. Drolet, X. Bi, and L. F. Liu (1994) *J. Biol. Chem.* **269**, 2068–2074.
38. M. Gellert, M. H. O'Dea, T. Itoh, and J.-I. Tomizawa (1976) *Proc. Natl. Acad. Sci. USA* **73**, 4474–4478.
39. A. Sugino, C. L. Peebles, K. N. Kreuzer, and N. R. Cozzarelli (1977) *Proc. Natl. Acad. Sci. USA* **74**, 4767–4771.
40. M. Gellert, K. Mizuuchi, M. H. O'Dea, T. Itoh, and J.-I. Tomizawa (1977) *Proc. Natl. Acad. Sci. USA* **74**, 4772–4776.
41. L. Ferrero, B. Cameron, B. Manse, D. Lagneaux, J. Crouzet, A. Famechon, and F. Blanche (1994) *Mol. Microbiol.* **13**, 641–653.
42. A. B. Khodursky, E. L. Zechiedrich, and N. R. Cozzarelli (1995) *Proc. Natl. Acad. Sci. USA* **92**,

11801–11805.

43. P. Bernard, K. E. Kezdy, L. van Melderen, J. Stayaert, L. Wyns, et al. (1993) *J. Mol. Biol.* **234**, 534–541.
44. P. Yorgey, J. Lee, J. Kordel, E. Vivas, P. Warner, et al. (1994) *Proc. Natl. Acad. Sci. USA* **91**, 4519–4523.
45. L. F. Liu, ed. (1994) *DNA Topoisomerases: Biochemistry and Molecular Biology, Advances in Pharmacology*, Vol. **29A**, Academic Press, San Diego.
46. L. F. Liu, ed. (1994) *DNA Topoisomerase and Their Applications in Pharmacology, Advances in Pharmacology*, Vol. **29B**, Academic Press, San Diego.
47. J. C. Wang (1982) *Sci. Am.* **247**, 94–109.
48. J. C. Wang (1991) *J. Biol. Chem.* **266**, 6659–6662.

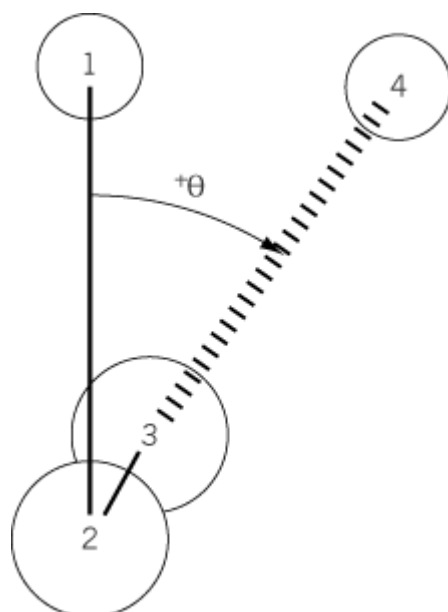
### Suggestions for Further Reading

49. For introductory accounts of DNA topology and DNA topoisomerases, see W. Bauer, F. H. C. Crick, and J. H. White (1980) *Sci. Am.* **243**, 118–133 (1980); J. C. Wang (1982), Ref. [47](#); and J. C. Wang (1994), in Ref. [46](#). A 1976 article by F. H. C. Crick (*Proc. Natl. Acad. Sci. USA* **73**, 2639–2643) also touched upon several interesting aspects of DNA topology. For more detailed and more recent studies, see the various chapters in the monographs given in Refs. [4](#), [45](#), and [46](#) and the reviews on DNA topoisomerases ([2](#), [30](#)), and references therein. References [45](#) and [46](#) contain a number of chapters on therapeutics targeting various DNA topoisomerases.

### Torsion Angle

Torsion angles within a molecule are defined by four atoms connected by three bonds. To determine the torsion angle of the bond connecting atoms 2 and 3, an observer sights down this and measures the angle that the bond between atoms 1 and 2 must be rotated through to become eclipsed with the bond between atoms 3 and 4 (Fig. [1](#)). The torsion angle is taken as positive if the front bond is rotated in a clockwise direction ([1](#)).

**Figure 1.** The torsion angle defined by four atoms connected by three covalent bonds. The clockwise rotation of the bond in front so as to be superimposed on the bond in back defines a positive torsion angle  $\varphi$ . For biological macromolecules, the backbone torsion angles are defined by four contiguous backbone atoms.

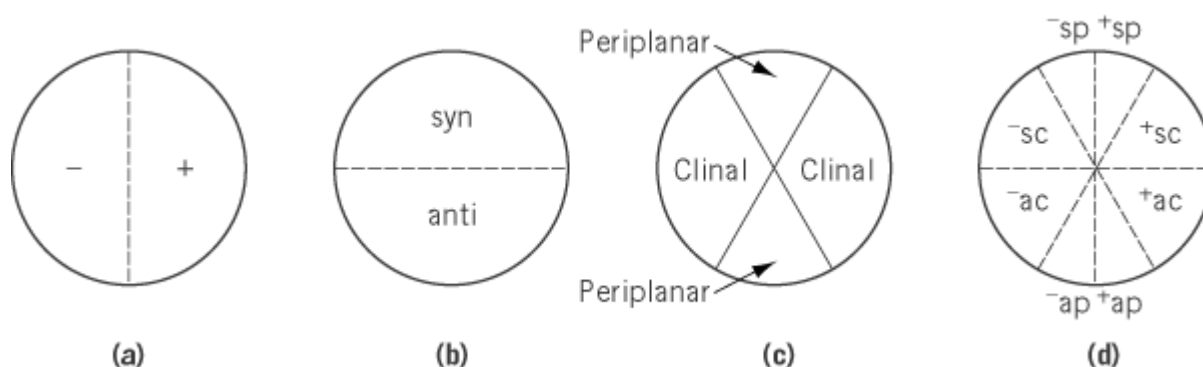


Torsion angles are preferentially given between  $-180$  and  $+180^\circ$  (2). This convention results in [polypeptide](#) torsion angles being  $180^\circ$  less than originally defined by Ramachandran (3) (see [Ramachandran Plot](#)). For polynucleotides, the torsion angles are defined by the conventions adapted by the IUPAC (4).

There is a technical difference between torsion angles and **dihedral angles**. A dihedral angle is defined by two planes containing four distinct atoms that are not necessarily sequentially covalently bonded. This is not possible with the definition of torsion angle.

Molecules that differ in torsion angles, ie, rotation about a single bond, are rotamers (see [Stereoisomers](#)). Torsion angles about single bonds may be described qualitatively, as well as quantitatively. The torsion angles are divided by three criteria: as positive and negative, as *syn* and *anti*, and as *periplanar* and *clinal* (Fig. 2, see bottom of page) (5). To give a qualitative description of the torsion angle, the abbreviation for all three divisions is given in Figure 2(d). The designation of torsion angles is of some importance, since this designation defines the difference between most [conformations](#).

**Figure 2.** Classifications of torsion angles. They are grouped by three criteria: (a) positive and negative rotation, (b) *syn* and *anti*, and (c) *clinal* and *periplanar*. This gives rise to the eight different rotamers, shown with their two-letter abbreviations in Figure 2(d).



## Bibliography

1. W. Klyne and V. Prelog (1960) *Experientia* **16**, 521.
2. IUPAC-IUB Commission on Biochemical Nomenclature (1970) *Biochemistry* **9**, 3471–3480.
3. G. N. Ramachandran, C. Ramakrishnan, and V. Sasisekharan (1963) *J. Mol. Biol.* **7**, 95–99.
4. IUPAC-IUB Commission on Biochemical Nomenclature (1983) *Eur. J. Biochem.* **17**, 193–201.
5. D. J. Millen (1962) In *Progress in Stereochemistry* (P. B. D. de la Mare, and W. Klyne, eds.), Vol.3, Butterworths, Washington, DC, pp. 138–164.

## Suggestions for Further Reading

6. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry, Part I: The Conformation of Biological Macromolecules*, W. H. Freeman, San Francisco, CA.
7. W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer-Verlag, New York, pp. 13–24.

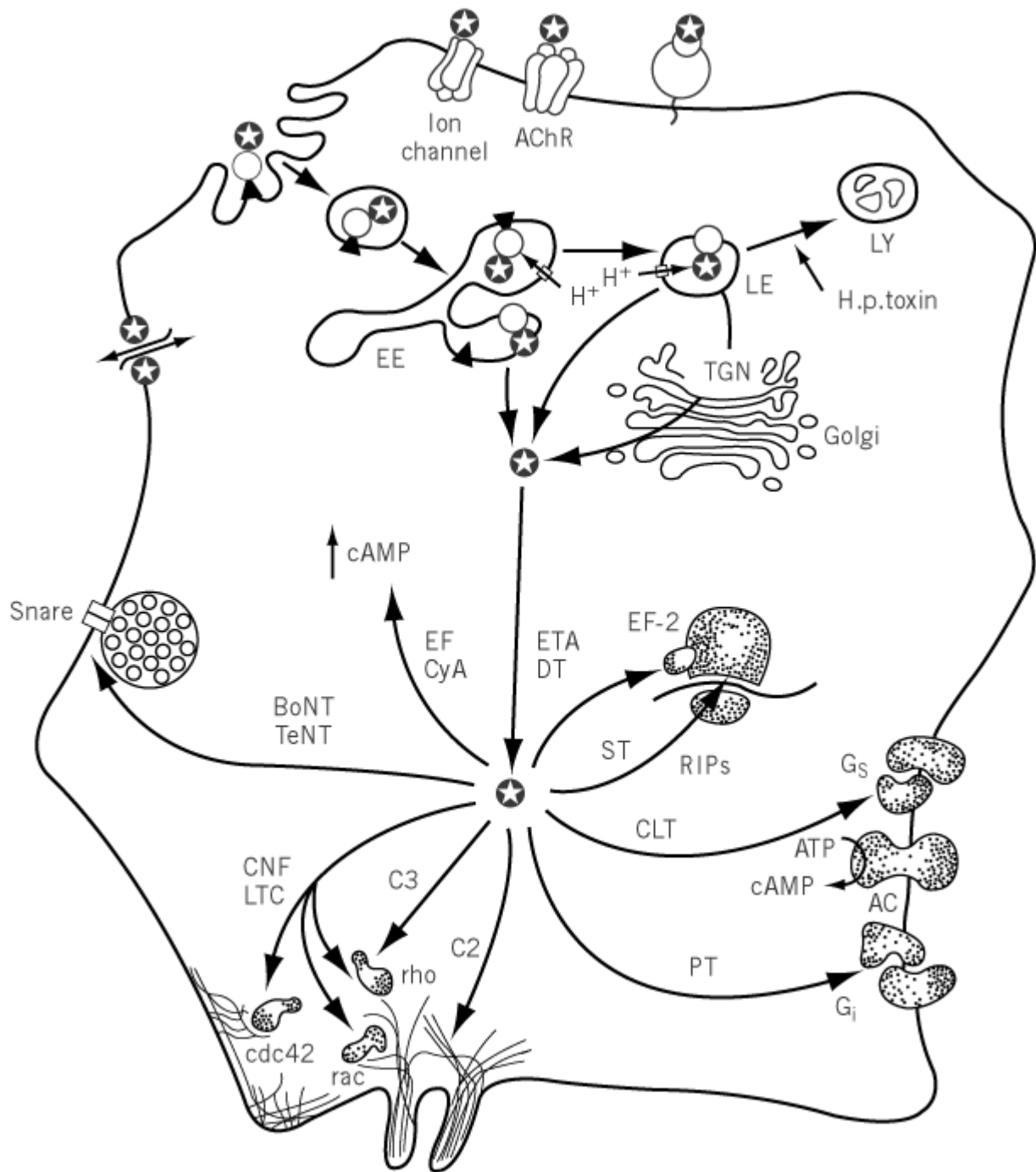
## Toxins

Thousands and thousands of living organisms of different complexity, from **prokaryotes** to **plants** and animals, produce and release a variety of toxins (1-5). A toxin can be generally defined as any substance that causes damage to the normal physiology of living organisms. Most toxins are produced with the objective of increasing the chance of survival of the toxin-producing species. Toxins can be of any chemical complexity, from the very simple formic acid to the long polypeptide chains of clostridial toxins. They are the end-result of a long-term evolution of the toxin-producing organism in its environment and hence reflect its particular life style. Thus, by learning about the mechanism of action of a toxin, we can, at the same time, learn about important physiological aspects of the living species interacting with the venomous organism. This concept has been summarized by Claude Bernard, after his classical studies on the mechanism of action of curare, “Poisons are chemical scalpels for the dissection of physiological processes.”

Most, if not all, animal protein toxins play a role in the prey–predator relationship and are delivered by the venomous animal with the intention of immobilizing a prey or of warning a predator (2, 3). The effects of their action have to develop very rapidly. Given these premises, it is no surprise that many animal toxins act on cell surface **receptors** of neuronal or muscle cells and their privileged site of action is the neuromuscular junction. The preferred target molecules are the nicotinic or muscarinic **acetylcholine receptors** and protein **channels** specific for **sodium**, **potassium**, or **calcium** ions (Fig. 1). Depending on the particular binding site, they can cause inhibitory paralysis (eg, by blocking the sodium or calcium channel function) or excitatory paralysis (eg, by prolonging the opening time of the sodium channel). An inhibitory action is frequently used to immobilize a prey, particularly when it has dimensions larger than those of the venomous animal itself, whereas an excitatory action is more frequently used to “warn” predators, because an excitatory stimulus is memorized much better than an inhibitory one. Most venoms contain a variety of toxins, with different targets and mechanisms of action (1-5). Several snakes produce, in addition to **neurotoxins**, hemorrhage factors that interfere with the coagulation cascade or degrade **extracellular matrix** proteins (3). Other preferred targets are **membrane** phospholipids, which are hydrolyzed to smaller components, thereby altering the integrity and permeability barrier of the plasma membrane (6, 7).

**Figure 1.** Cellular interactions of protein toxins. This schematic picture summarizes the large array of interactions of known toxins (indicated by a star) with cells. A very large number of toxins bind to ion channels or cell surface receptors, causing their inactivation or activation (top part). Other toxins form pores or channels across the plasma membrane (left). Many toxins display their activity inside cells. They bind to a receptor (o) and enter the endocytic pathway through early endosomes (EE) and late endosomes (LE). They may be transported retrogradely up to the **Golgi** cisternae and [endoplasmic reticulum](#) or proceed to **lysosomes** (LY), where they are degraded. At some stage along these intracellular pathways, the catalytic toxin subunit (indicated by a star) crosses the intracellular membrane and translocates into the cytosol, where it recognizes and modifies specifically a target molecule. The *Helicobacter pylori* toxin (H.p. toxin) causes a swelling of a hybrid LE-LY compartment. The edema factor of *Bacillus anthracis* (EF) and [adenylate cyclase](#) of *Bordetella pertussis* (CyA) elevate the cellular **cAMP** level. Diphtheria toxin (DT) and exotoxin A from *Pseudomonas aeruginosa* (ETA) **ADP-ribosylate** the [elongation factor 2](#) and block protein synthesis. Shiga toxins (ST) and the ribosomal-inactivating toxins (RIPs) remove a single adenine residue from the 28S ribosomal RNA and stop protein synthesis. Cholera toxin (CLT) and related toxins ADP-ribosylate the  $G_s$  subunit of hormone-coupled adenylate cyclase complex and elevate the levels of cAMP. Pertussis toxin (PT) ADP-ribosylates the  $G_i$  subunit of hormone-coupled adenylate cyclase complex, as well as other **G proteins**, and causes an increase of cAMP. The colony necrotic factor of *Escherichia coli* (CNF), large clostridial toxins (LCT), and small clostridial C3 affect [actin](#) polymerization-controlling proteins of the **ras** superfamily, whereas the clostridial C2 and related toxins ADP-ribosylate directly G-actin. Tetanus neurotoxin (TeNT) and botulinum neurotoxins (BoNT) are zinc-endopeptidases, which cleave specifically three protein components of the neuroexocytosis apparatus and block [exocytosis](#).





Plants produce a variety of toxins including proteins that kill **eukaryotic** cells by inactivating their [ribosomes](#), thereby blocking protein synthesis (8). These toxins act slowly, in contrast to animal toxins, whose rapidity of action is essential. They exert their activity inside cells and exploit [endocytosis](#) to enter into the **cytosol**. The role of such toxins in plant survival is not clear, but they have been suggested to protect plant tissues and seeds from parasites (8). As well as plants, microorganisms do not depend on speed of action, and many bacterial toxins act slowly, but effectively, inside cells or on the plasma membrane. Microorganisms that live on or inside animals coevolve with their host much more rapidly than animals with their prey. Thus, they have developed, and still do, a larger variety of strategies to modify the host physiology in such a way as to promote their multiplication and/or spread, or to counteract the inflammatory and immune defense responses (9). Microbial toxins are part of these strategies, and they hit on vital physiological functions of animals by modifying selected cell targets (Fig. 1) (10). Many of them alter cell responses by

binding to particular receptors as [superantigen](#) toxins or endotoxins do. [Endotoxins](#) are lipopolysaccharides of the outer membrane of many **gram-negative** bacteria, which induce the release of **cytokines** upon binding to inflammatory cells. Other toxins alter the plasma membrane permeability barrier, and cells die by colloid-osmotic lysis and/or by loss of cytosolic components. Several bacterial protein toxins are [enzymes](#) acting inside cells, and their structural organization reflects the need to enter into cells to perform their activity (10, 11). Together with the plant toxins mentioned above, these protein toxins are known as A-B toxins, because they consist of two parts: A (active) is an enzyme connected via a peptide linker and a single [disulfide bond](#) to polypeptide B (binding), which is responsible for cell binding and internalization. Figure 1 depicts schematically the route of entry of A-B toxins into cells, a four-step process that begins with binding to a cell surface receptor. These proteins are activated by specific cleavage of the peptide linker by bacterial or tissue [proteinases](#) or by cell surface-associated proteinases, such as furin, which cleaves between positively charged residues (12). Binding is rapidly followed by internalization of the toxin–receptor complex inside **vesicles** (second step), whose nature depends on the toxin receptor: Different toxins travel different vesicular trafficking pathways inside the cell. In any case, in order to display its toxic activity, the A subunit has to leave the vesicle lumen and translocate across the membrane into the cytosol (third step), where it modifies its target (fourth step) and intoxicates the cell. The yield of the process of toxin cell entry may be low, but this is compensated by a catalytic type of activity that leads to the inactivation, one after the other, of all target molecules present in the cell. At variance from most animal toxins, which act on a one-to-one basis, at least in principle, one single copy of A is sufficient to intoxicate a cell completely. So far, five different enzymic activities have been found associated with proteins that act intracellularly (Table 1). The cellular and systemic effect of such activities depends on the role of the target molecule in cell function and of the intoxicated cell in the physiology of the poisoned organism.

**Table 1. Catalytic Activities, Targets, and Cell Effects of Bacterial Protein Toxins with Intracellular Targets**

| Toxin             | Activity               | Target                        | Effect                                       |
|-------------------|------------------------|-------------------------------|--|
| DT                | ADP-ribosyltransferase | EF-2                          | Blockade of protein synthesis and cell death |
| ETA               | ADP-ribosyltransferase | EF-2                          | Blockade of protein synthesis and cell death |
| CLT               | ADP-ribosyltransferase | $G_s$ , $G_p$<br>$G_{olfs}$   | Increase cAMP                                |
| LT                | ADP-ribosyltransferase | $G_s$ , $G_p$<br>$G_{olfs}$   | Increase cAMP                                |
| PT                | ADP-ribosyltransferase | $G_i$ , $G_o$ ,<br>$G_{gust}$ | Increase cAMP                                |
| $C_2$             | ADP-ribosyltransferase | Actin                         | Inhibition of actin                          |
| $C_2$ -liketoxins | ADP-ribosyltransferase | Actin                         | Polymerization, depolymerization of F-actin  |
| $C_3$             | ADP-ribosyltransferase | Rho                           | Depolymerization of F-actin                  |
| <i>C.d.</i> LCT   | Glucose-transferase    | Rho, Rac,                     | Depolymerization of F-                       |

|                 |   |                           |  |
|-----------------|---|---------------------------|--|
| <i>C.s.</i> LCT | Glucose-transferase                           | Cdc42<br>Rac, Ras,<br>Rap | actin<br>Breakdown of stress<br>fibers         |
| <i>C.n.</i>     | Toxin N-acetyl-<br>glucosamine<br>transferase | Rho, Rac,<br>Cdc42        | Depolymerization of F-<br>actin                |
| STs             | Adenine-<br>glycohydrolase                    | r-RNA<br>28S              | Blockade of protein<br>synthesis and celldeath |
| EF              | Adenylcyclase                                 | None                      | Increase cAMP                                  |
| CyA             | Adenylcyclase                                 | None                      | Increase cAMP                                  |
| TeNT, BoNT/B    |   |                           |  |
| BoNT/D, /F, /G  | Metallo-protease                              | VAMP                      | Inhibition of exocytosis                       |
| BoNT/A, /C, /E  | Metallo-protease                              | SNAP-25                   | Inhibition of exocytosis                       |
| BoNT/C          | Metallo-protease                              | Syntaxin                  | Inhibition of exocytosis                       |

A-B toxins can be divided into two groups on the basis of the structural organization of the B portion: (1) oligomeric B toxins (cholera toxin, shiga toxins, etc.) and (2) three-domain toxins ([diphtheria toxin](#), clostridial neurotoxins, etc.) that are composed of a two-domain B polypeptide chain plus an A domain. These latter toxins use their carboxyl-terminal B domain to bind to a protein receptor and enter the lumen of intracellular compartments that are acidic. Low pH induces the transition from a neutral conformation to an acidic one, capable of inserting both the A and B subunits into the lipid bilayer. There is evidence that the B amino-terminal domain forms a transmembrane channel, open laterally to lipids, that mediates the membrane translocation of A ([10](#), [11](#)). This event and the release of A in the cytosol are coupled to the reduction of the intersubunit disulfide bond, which appears to be the rate-limiting step of the entire cell entry process. Oligomeric B toxins bind to the saccharide moiety of glycolipids or **glycoproteins**, localized on the apical membrane of epithelial cells. Their mode of cell entry is less understood than that of the three-domain bacterial toxins, but they do not appear to depend on low luminal pH to translocate A into the cytosol.

As depicted in Figure [1](#), the various toxins attack different intracellular targets with a variety of consequences. Plant toxins, diphtheria, shiga, and pseudomonas toxins block protein synthesis and cause cell death. However, the majority of protein toxins with intracellular targets do not kill the cell; instead, they alter cell structure and and/or physiology in such a way as to promote bacterial proliferation and/or spread. Nonetheless, cellular alterations can be such as to be incompatible with the survival of the tissue or organism.

#### Bibliography

1. E. Habermann (1972) *Science* **177**, 314–322.
2. P. N. Strong (1990) *Pharmacol. Ther.* **46**, 137–162.
3. A. L. Harvey, ed. (1991) *Snake Toxins*, Pergamon Press, New York.
4. B. M. Olivera et al. (1990) *Science* **249**, 257–263.
5. R. Rappuoli and C. Montecucco (1997) *Guidebook to Protein Toxins and Their Use in Cell Biology*, Sambrook and Tooze, Oxford University Press, Oxford, UK.
6. H. L. Harvey (1990) In *Handbook of Toxicology* (W. T. Shier and D. Mebs, eds.), Marcel Dekker, New York, pp. 1–66.
7. R. M. Kini, ed. (1997) *Venom Phospholipase A2 Enzymes: Structure, Function and Mechanism*,

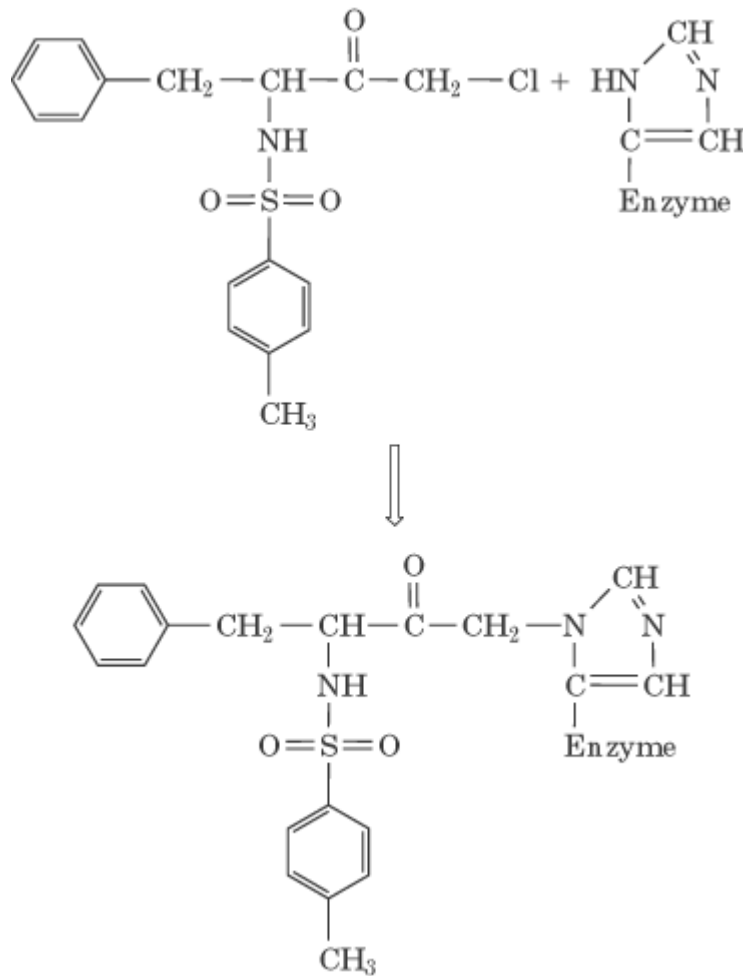
Wiley, New York.

8. L. Barbieri, M. G. Battelli, and F. Stirpe (1993) *Biochim. Biophys. Acta* **1154**, 237–282.
9. C. E. Mims (1995) *The Pathogenesis of Infectious Diseases*, Academic Press, London.
10. G. Menestrina, G. Schiavo, and C. Montecucco (1994) *Mol. Aspects Med.* **15**, 81–193.
11. J. E. Alouf and J. H. Freer, eds. (1991) *Sourcebook of Bacterial Protein Toxins*, Academic Press, San Diego, CA.
12. M. Tsuneoka et al. (1993) *J. Biol. Chem.* **268**, 26461–26465.

### TPCK (*N-P-Tosyl-Phenylalanine Chloromethyl Ketone*)

This reagent was developed by Schoellman and Shaw as a specific inhibitor of **chymotrypsin** (1). Its design was based on the principle of [affinity labeling](#), which couples substrate specificity with chemical reactivity to generate a covalent bond between the reagent and a functional group at the [active site](#) of the enzyme. In this case, the [phenylalanine](#) part of the reagent confers specificity for chymotrypsin, while the chloromethyl ketone provides chemical reactivity (Figure 1). Although chymotrypsin is a [serine proteinase](#), the reagent does not interact with the active site serine residue but instead alkylates the [histidine](#) residue that is part of the [catalytic triad](#). The reagent has essentially no effect on the activity of **trypsin** or other serine proteinases and is often used to inactivate the small amounts of chymotrypsin present in purified trypsin preparations. It is also often combined with inhibitors of other types of proteinases to prevent unwanted **proteolysis** that might occur during the course of protein isolation from cell or tissue homogenates (2).

**Figure 1.** Reaction of TPCK (tosyl-phenylalanyl chloromethyl ketone, ie, L-1-chloro-3-[4-tosylamido]-4-phenyl-2-butanone) with the active-site histidine side chain of a chymotrypsin-like serine proteinase. The reaction product is catalytically inactive.



## Bibliography

1. G. Schoellmann and E. Shaw (1963) *Biochemistry* **2**, 252 ff.
2. R. J. Beynon (1989) in *Protein Purification Methods: A Practical Approach* (E. L. V. Harris and S. Angal, eds.), IRL Press, Oxford, U.K., pp. 40–49.

## Tracking Dyes

In [gel electrophoresis](#) using discontinuous [buffer](#) systems, [disc electrophoresis](#), dyes migrating with a net mobility intermediate between that of the leading and trailing ions of a moving boundary are used to mark the moving boundary (“tracking dyes”). They are used to monitor the extent of electrophoresis, and the mobilities of electrophoretic bands within samples are often expressed relative to those of the tracking dye.

Bromphenol blue is the most popular example of a tracking dye, being used generally in [SDS-PAGE](#), but it, like any other dye, is limited in this function by its net mobility. Thus, it may mark a chloride/glycinate moving boundary (stack) in SDS-PAGE incorrectly. This happens when it is displaced by the faster-migrating micellar **SDS**, as occurs in a typical polyacrylamide gel of

moderate or low concentration. Only after molecular sieving effects start to retard micellar SDS (eg, in a 12% (w/v) gel) can bromphenol blue occupy and correctly mark the boundary. By contrast, the dye pyroninY, when complexed with SDS, has a mobility greater than that of micellar SDS and marks the boundary correctly at low gel concentrations. As a result of its binding of SDS, however, it is large and is retarded by molecular sieving at gel concentrations greater than about 12% (w/v), so as to migrate behind the moving boundary front. At even higher gel concentrations, even a small dye like bromphenol blue is retarded by sieving behind the moving boundary (1). Thus, tracking dye stacking needs to be tested by comparing its position on the gel with that of the boundary determined more directly, such as by precipitation or [radioactivity](#) of the leading or trailing ions. Since a moving boundary front (stack) marked by a tracking dye is frequently used to define the characteristic relative mobility ( $R_f$ ) of a band, that  $R_f$  is valid only in the gel concentration range in which the tracking dye migrates with the moving boundary front. Beyond that gel concentration range, the dye cannot be used to define the  $R_f$ .

In gel electrophoresis in a continuous buffer, any dye can be used for the calculation of the  $R_f$  of a band at any gel concentration. In contrast to the stacked dye in a discontinuous buffer system, however, the relatively small molecular weight dye band spreads rapidly by [diffusion](#) in proportion to the migration distance, so that its usefulness as a reference band for  $R_f$  measurement progressively diminishes with the time of electrophoresis.

#### Bibliography

1. M. Wyckoff, D. Rodbard, and A. Chrambach (1977) *Anal. Biochem.* **78**, 459–482.

#### Trans Configuration

*Trans* is a prefix that is used to describe the [configuration](#) of molecules, indicating that two substituents are on opposite sides of a bond or ring (1). The alternative, when the substituents are on the same side, is denoted by [cis configuration](#). In specifying the *trans* configuration about a double bond, the prefix (*E*), for the German *entgegen*, is preferred (2). The term *trans* is sometimes used in specifying a [conformation](#) when the [torsion angle](#) approaches 180°. To avoid confusion, it is preferable to use *anti* for describing conformations.

#### Bibliography

1. A. Baeyer (1888) *Liebig's Annalen der Chemie* **CCXLV**, 137.
2. IUPAC (1974) *Pure Appl. Chem.* **45**, 13–30.

#### Suggestions for Further Reading

3. E. L. Eliel (1962) *Stereochemistry of Carbon Compounds*, McGraw-Hill, New York, pp. 318–371.
4. J. March (1985) *Advanced Organic Chemistry*, 3rd ed., Wiley-Interscience, New York, pp. 109–115.

## Trans-Acting

*Trans*-acting is the effect of one **gene** on the activity of another regardless of whether or not the two are physically linked. The term is used in contradistinction to [cis-acting](#) which describes the situation where physical linkage (usually close) of genetic elements is necessary for interaction.

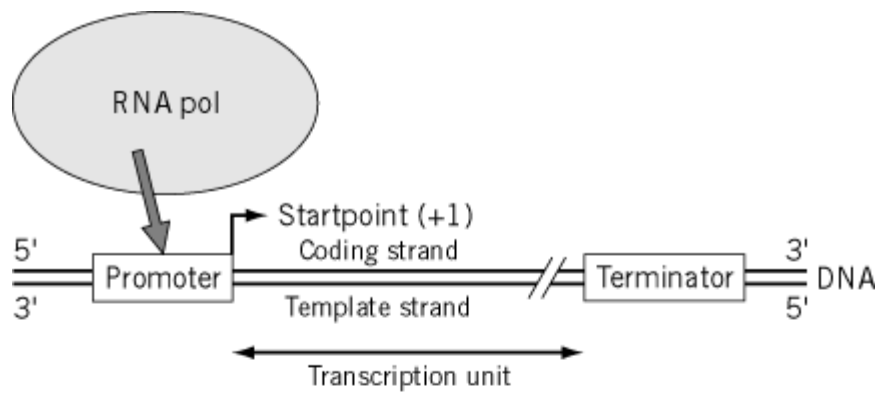
In the networks of gene interactions that play such a major part in the controls of metabolism and development, *trans*-activity of genes, through their protein products, is the general rule. Among *trans*-acting genes are all those that encode **transcriptional factors** and proteins whose binding is necessary for the function of [enhancer](#) sequences, proteins involved in intron splicing (see *introns*, *spliceosomes*), and proteins of [chromatin](#) complexes.

In general, proteins that function by binding to DNA or (for splicing control) to RNA are all *trans*-acting and the nucleic acid sequences to which they bind are *cis*-acting. However, in the fly *Drosophila melanogaster* there are some examples of DNA sequences that act in *trans*, enhancing the activity of a *cis*-positioned gene and also, to some extent, of the *trans* homologue, if the latter lacks an enhancer of its own. This unusual phenomenon, called *transvection*, depends on the close pairing of homologous [chromosomes](#) in somatic cells of *Drosophila*.

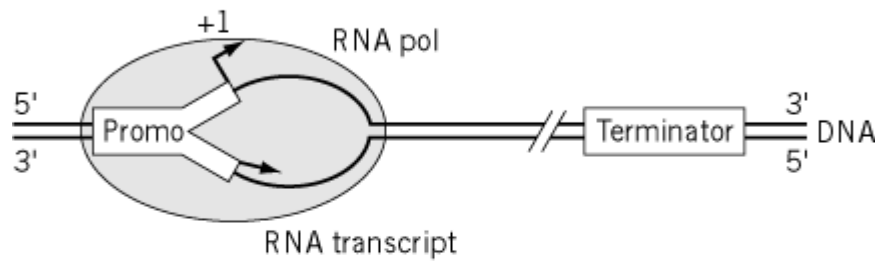
## Transcription

Transcription transfers the genetic instructions from the DNA template to RNA and is the first step in expression of the [genome](#). This highly regulated process consists of the synthesis of an RNA molecule by an **RNA polymerase** (RNA pol) that reads the template (or noncoding) strand of the DNA; the primary transcript will therefore carry the same genetic message as the other DNA strand, namely the coding strand. The synthesis of RNA involves the formation of a partial heteroduplex between the DNA template strand and the growing RNA transcript, leading to a bubble structure. RNA pol catalyzes the formation of successive phosphodiester bonds between the 5'-triphosphate group of nucleotide precursors and the 3'-OH group of the last nucleotide of the nascent RNA polynucleotide chain; the growth of the RNA chain proceeds in the 5' to 3' direction. Transcription is generally divided into four steps (Fig. 1): (i) Preinitiation involves the recognition of particular sequences constituting the [promoter](#) by the transcriptional apparatus, including the binding of RNA pol and opening of the double-stranded DNA around the initiation site; (ii) initiation is usually considered to be the formation of the first phosphodiester bond of the transcript; (iii) elongation begins with clearance of the promoter, which refers to the escape of RNA pol from its anchorage site on the promoter; the enzyme then progresses along the gene to elongate the RNA transcript, leading to concomitant movement of the bubble structure by successive unwinding and rewinding of the DNA duplex; and (iv) finally, termination occurs when the elongation complex encounters a terminator sequence, with subsequent release of both RNA pol and the RNA product. Depending on the type of gene from which it derives, the primary transcript will undergo differential processing to give rise to various types of RNAs: **ribosomal RNA**, [transfer RNA](#), **small nuclear RNA**, or [messenger RNA](#) that are further **translated** into proteins.

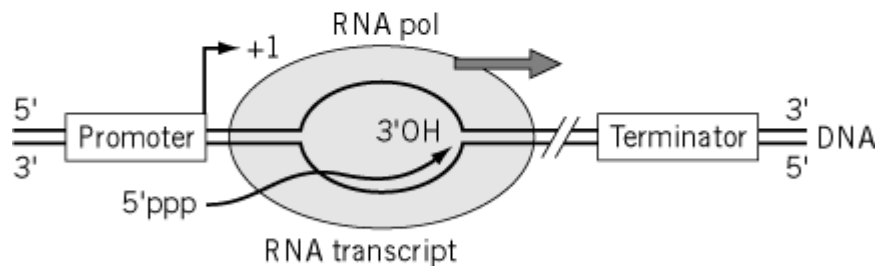
**Figure 1.** Major steps of the transcription process. (a) Promoter recognition; (b) initiation; (c) elongation; (d) termination.



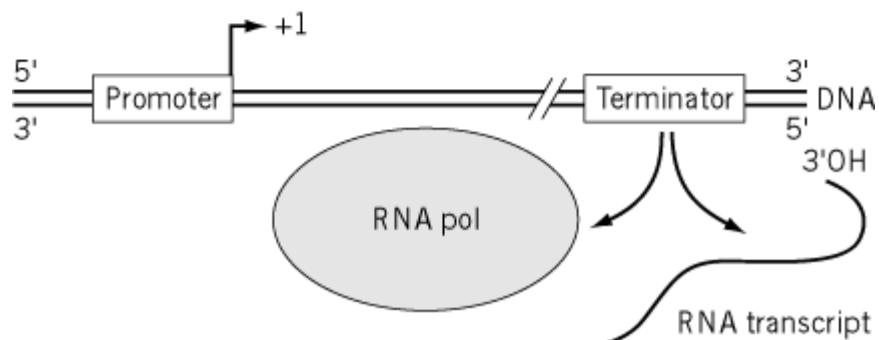
(a)



(b)



(c)



(d)

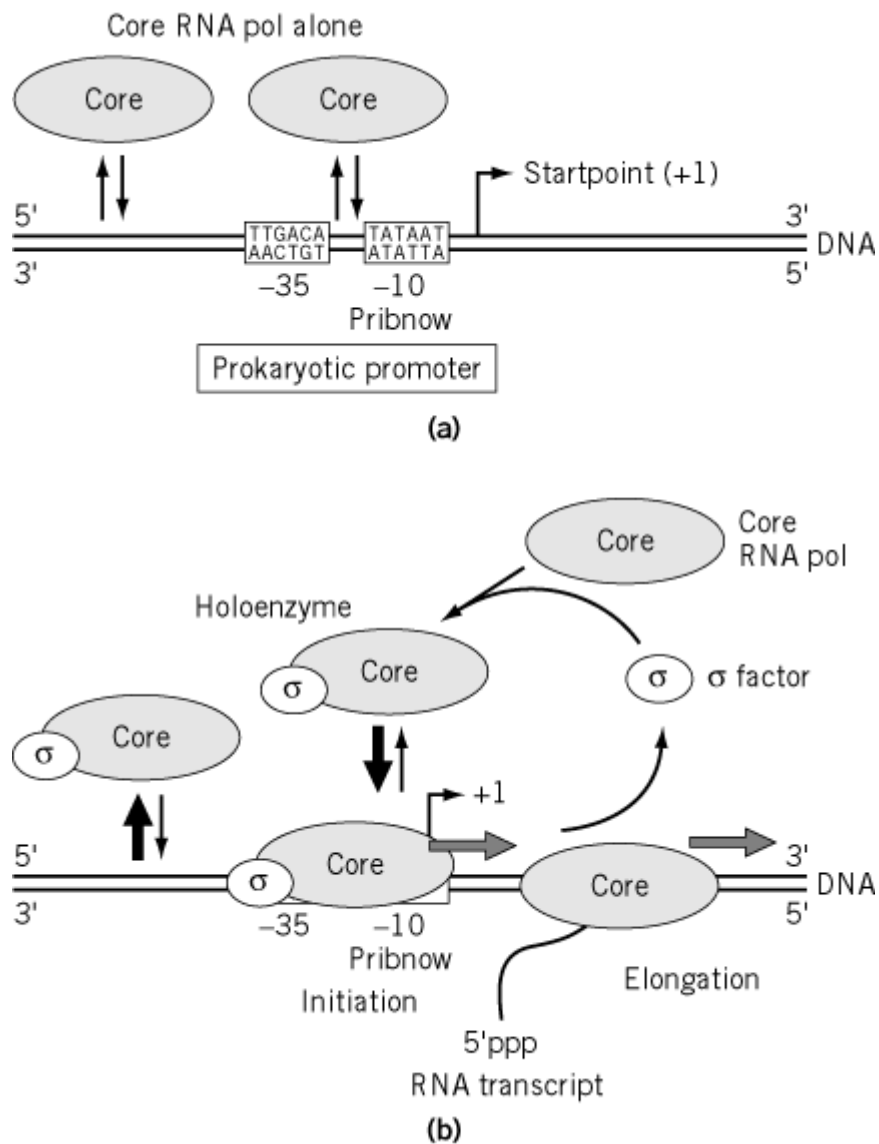
Although the general scheme of transcription has been conserved during [evolution](#), some noteworthy differences exist between the prokaryotic and the eukaryotic systems; these include the promoter structure, the components of the transcription machinery, the function of accessory factors, and the mode of regulation.

## 1. Transcription in Prokaryotes



Prokaryotes possess a single RNA pol to synthesize all RNAs. RNA pol is composed of only five subunits: four of them constitute the core enzyme ( $\alpha_2\beta\beta'$ ) and the last, called the [sigma factor](#), associates with the core to form the holoenzyme. The  $\beta$  and  $\beta'$  subunits together form the catalytic center, whereas the two  $\alpha$  subunits maintain the structural conformation. Accurate initiation of transcription occurs only in the presence of holoenzyme. The  $\sigma$  factor confers promoter specificity in two ways: It destabilizes nonspecific RNA pol–DNA complexes, while at the same time promoting formation of a stable RNA pol–promoter complex. As soon as short transcripts (10 to 15 nucleotides) have been synthesized,  $\sigma$  dissociates from the transcription complex, and the core enzyme continues elongation (Fig. 2).

**Figure 2.** Accurate initiation from a prokaryotic promoter requires RNA polymerase holoenzyme. Core polymerase forms non-productive complexes with either nonspecific or promoter DNA sequences (a). Association of  $\sigma$  factor with core polymerase, to generate holo-polymerase, confers promoter specificity and allows accurate initiation of transcription. When the RNA product reaches a length of 10 to 15 nucleotides,  $\sigma$  is recycled in other initiation events (b).

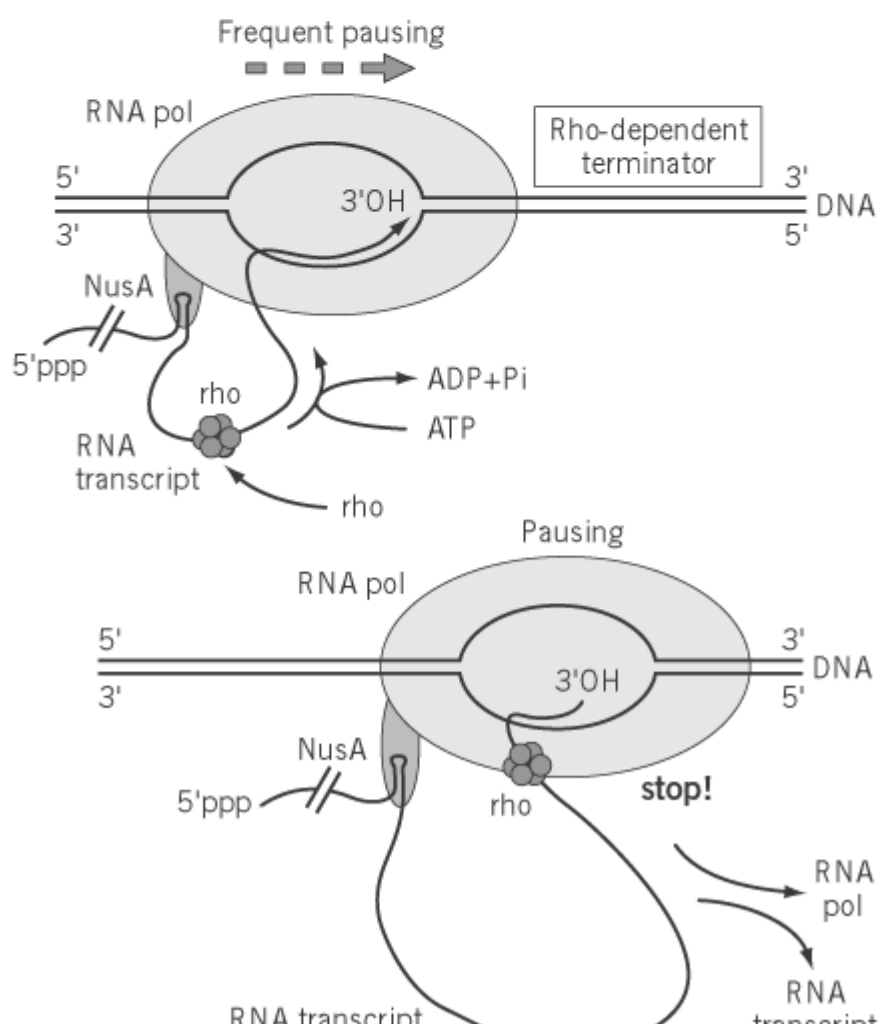
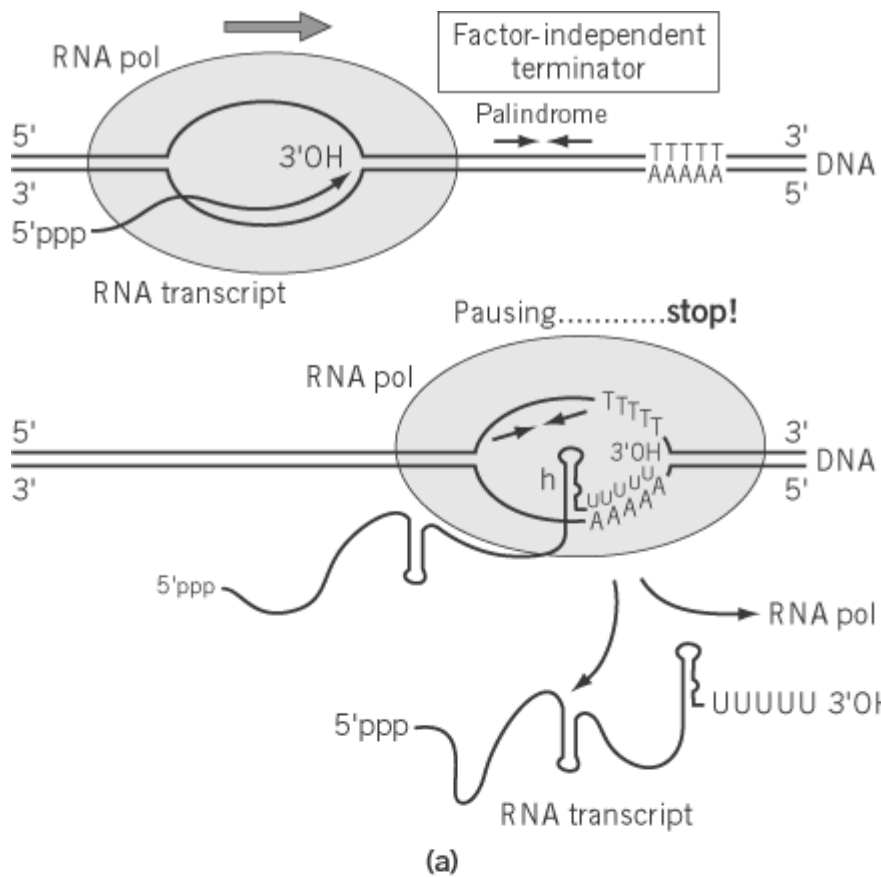


Prokaryotic promoters are characterized by four conserved features: (i) the startpoint, also called the

initiation or start site (+1), usually a purine; (ii) the [Poly A](#), a short A-T-rich sequence located at the -10 position; (iii) an element located around the -35 position; and (iv) the fixed distance that separates both of these elements (16 to 18 bp). During transcription initiation, RNA pol covers a region extending from -50 to +20 around the start site by interacting with the Pribnow box and the -35 element. The -35 sequence is recognized by the RNA pol, whereas subsequent opening of the promoter from nucleotides -10 to +3 depends on the Pribnow box.

During elongation, the core polymerase may pause following the formation of secondary structure within the nascent RNA. The presence of pausing sites is also linked to termination events. Termination occurs through two mechanisms: (i) The signal of arrest at the end of the gene is provided by a G-C-rich **palindromic** sequence, followed by a [poly-A](#) sequence; transcription of these sequences leads to the formation of a stable hairpin directly followed by a poly-U sequence (Fig. [3a](#)). Once the hairpin structure has induced pausing of the polymerase, the poly-U/poly-A heteroduplex allows further release of the transcript and the enzyme. (ii) Termination may require the presence of the hexameric releasing factor rho, which exhibits an RNA-dependent [Operons](#) and an **RNA-DNA helicase** activity (Fig. [3b](#)). Rho binds the RNA upstream of the rho-dependent terminator and subsequently translocates along the transcript. Termination most likely occurs upon interaction with the RNA pol stalled at the rho-dependent terminator. Consistently, rho is able to interact with the  $\beta$  subunit of the enzyme. Similarly, Nus A factor regulates transcription by increasing the tendency of RNA pol to pause.

**Figure 3.** Models for termination of prokaryotic transcription. ( **a** ) Factor-independent termination: Transcription of palindromic and polyA sequences from the factor-independent terminator leads to the formation of a hairpin structure (*h*) and to a poly-U sequence, respectively. The hairpin induces pausing of the polymerase; together with the adjacent polyA-polyU hybrid, this promotes transcription arrest and release of the RNA transcript and the RNA polymerase. ( **b** ) Rho-dependent termination: Termination factor rho binds the RNA transcript upstream of the rho-dependent terminator and migrates toward the elongating polymerase. Termination occurs when rho catches up the polymerase pausing at the terminator. The factor NusA facilitates termination by increasing the frequency of pausing.



Prokaryotic genes are clustered in [operons](#); an operon is controlled by a single promoter and contains several genes all devoted to the same metabolic pathway, such as lactose catabolism in the case of the *lac operon*. Regulation of transcription is performed essentially at two stages: (i) At the initiation level, variation in the pattern of transcribed operons can result from the use of different  $\sigma$  factors, each specific to a type of promoter. In this respect, events such as **sporulation** or response to **heat shock** are driven by a cascade of  $\sigma$  factors. The host RNA pol can be substituted by a newly synthesized viral RNA pol upon bacteriophage infection. An alternative pathway requires **trans-acting** factors that bind to the **operator** element, a specific sequence generally located nearby the promoter, to enhance (activators) or repress (repressors) the initiation of transcription. (ii) When the elongating RNA pol encounters a first termination sequence, it may either stop or continue to transcribe adjacent genes. This will depend on factors that fasten on the elongating polymerase after having bound to a specific site on the DNA or on the transcript. In the case of **antitermination** factors (or antiterminators), RNA pol will read through the stop signal, but some other factors exist that will increase its propensity to stop, probably by inducing frequent pausing. At some operons, the phenomenon of **attenuation** can also occur, when transcription is intimately coupled with translation: progression of the ribosome along the nascent RNA modulates the formation of secondary structures, enabling RNA pol to read through intrinsic terminators (or attenuators).

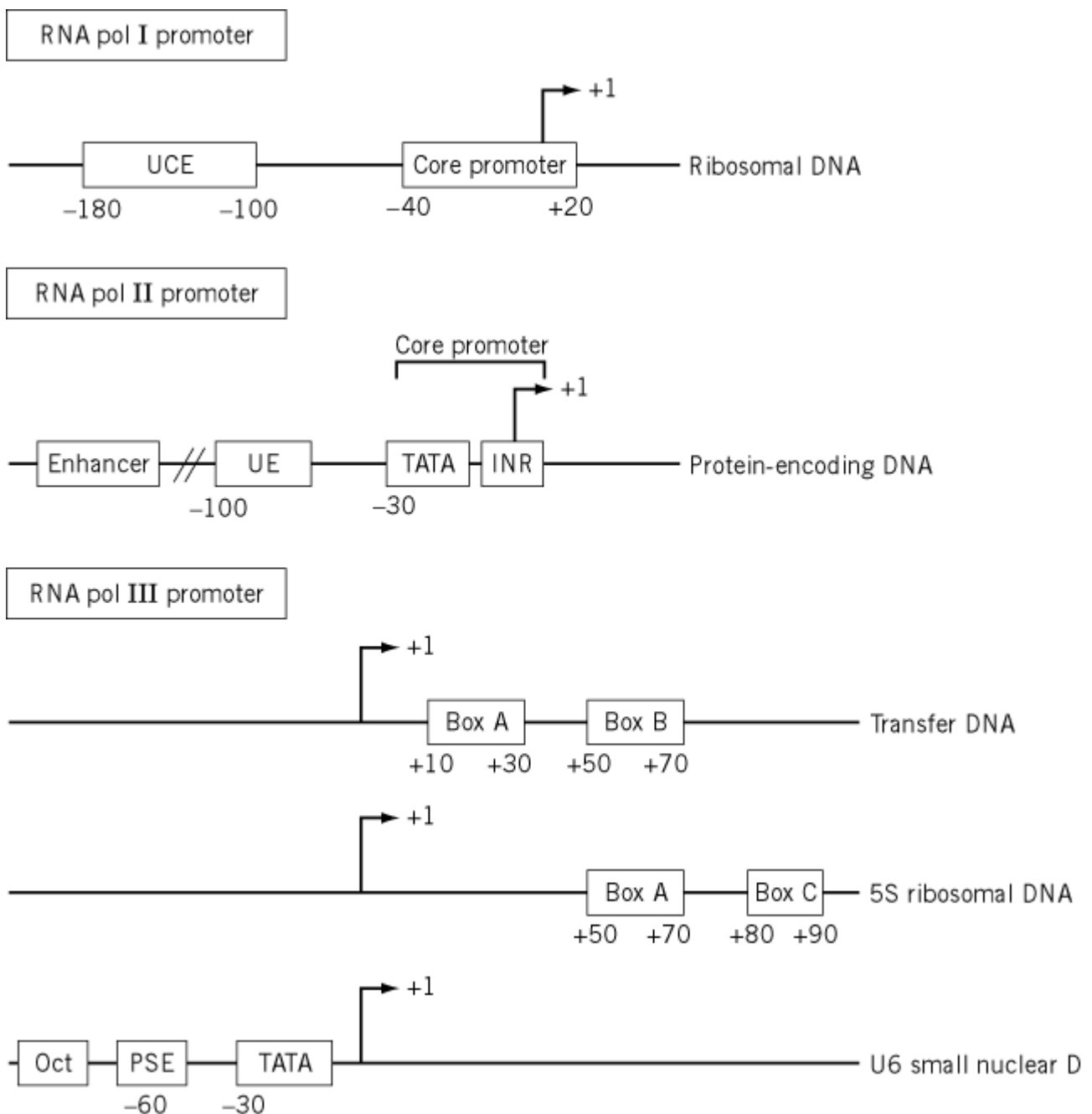
## 2. Transcription in Eukaryotes

Consistent with a higher complexity in the genome organization, the transcription process in eukaryotes involves a much larger number of proteins than in prokaryotes. There are three different RNA polymerases, each specialized in transcription of a particular class of genes:

1. RNA pol I transcribes ribosomal RNA genes.
2. RNA pol II transcribes protein-encoding genes and certain small nuclear RNA (sn RNA) genes.
3. RNA pol III gives rise to small RNAs such as transfer RNAs, 5 S rRNA, and snRNAs.

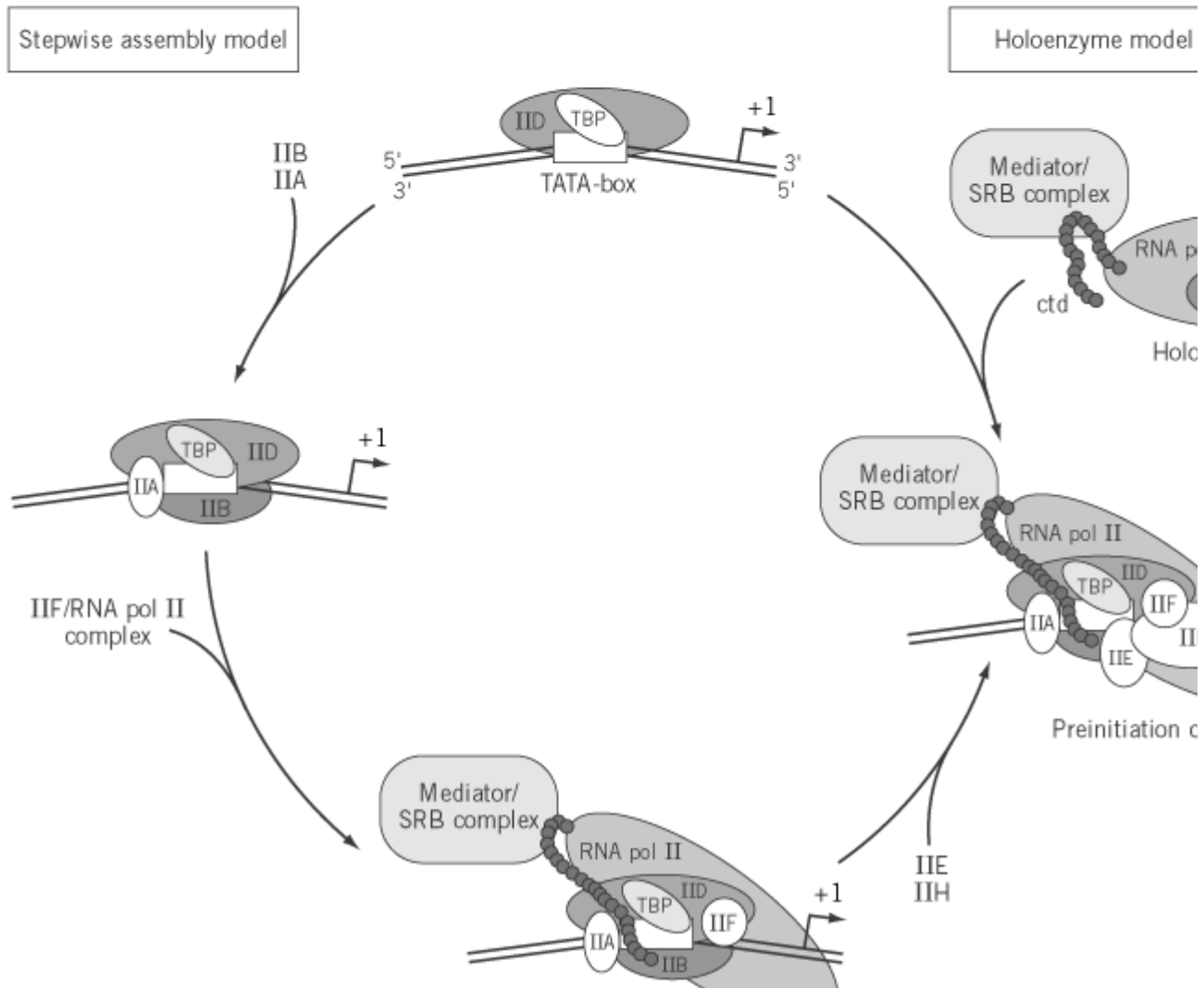
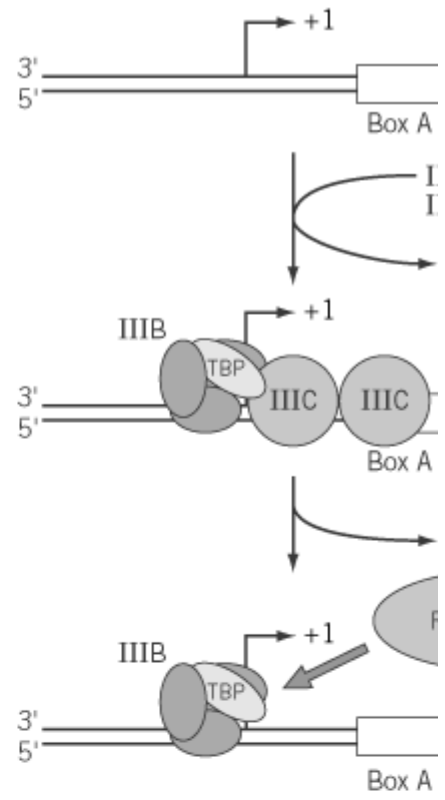
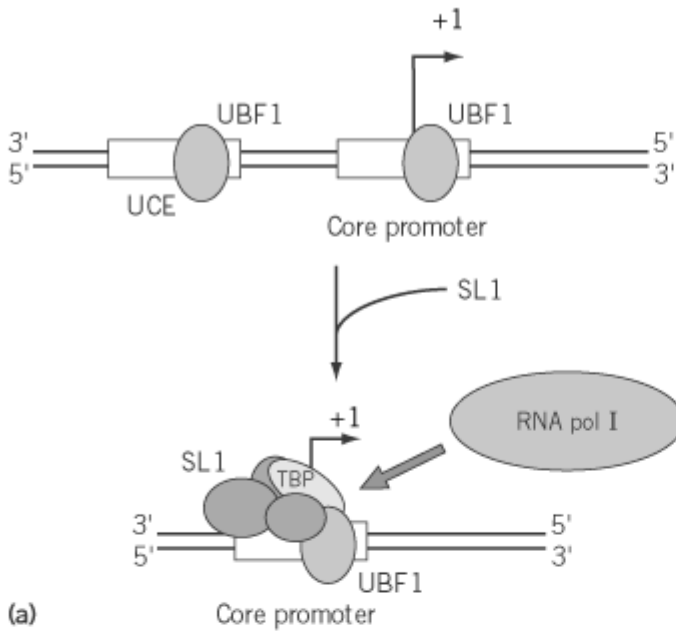
Each enzyme is composed of at least 12 subunits, three of which are common to all enzymes. They initiate RNA synthesis from three different types of promoters (Fig. 4). Moreover, promoter recognition by any RNA pol is absolutely dependent upon preliminary binding of accessory factors.

**Figure 4.** Different classes of eukaryotic promoters. TATA box (TATA) and initiator element (INR) constitute the RNA polymerase II core promoter. UCE, upstream control element; UE, upstream element; Oct, octamer sequence; PSE, proximal sequence element. +1 indicates the startpoint.



RNA pol I promoters include a core promoter (or selector) and a regulatory upstream control element (UCE). Both elements contain a G-C-rich sequence that is recognized by the ancillary factor UBF1. UBF1 allows the binding of the SL1 factor on the promoter, before RNA pol I finally joins the initiation complex (Fig. 5a). RNA pol III promoters present more diversity in their structure (Fig. 4). Downstream (or internal) promoters contain two conserved sequences (the box A sequence associated with either box B or box C), whereas upstream promoters possess a [TATA-box](#) element, in addition to Oct (octamer) and PSE (proximal sequence) regulatory elements. Positioning of RNA pol III on these promoters is mediated by the [transcription factor](#) TFIIB. TFIIB binds to the promoter either independently via the TBP subunit (TATA-binding protein) when a TATA-box sequence is present, or through the assistance of TFIIA and TFIIC in the absence of any TATA-box (Fig. 5b).

**Figure 5.** Initiation from different eukaryotic promoters. Positioning of RNA pol I (**a**) and RNA pol III (**b**) on class I and promoters, respectively. UBF1 and SL1 refer to class I transcription factors, and IIIA, IIIB, and IIIC to class III transcription subunit common to the “positioning” factors SL1 and IIIB. (**c**) Stepwise assembly model and holoenzyme model for transcription complex on RNA pol II core promoter. IIA, IIB, IID, IIE, IIF, and IIH are the class II basal transcription factors. TBP is The C-terminal domain (ctd) of the largest subunit of RNA pol II interacts with the mediator complex. +1 indicates the start site.



Transcription of protein-encoding genes is much more complex and requires two sets of factors (see [Enhancer](#)): (i) The basal or general transcription factors (GTFs) are necessary for transcription from a minimal promoter; they assemble onto the promoter to form, together with the RNA pol II, the preinitiation complex (PIC). (ii) The sequence-specific factors bind their cognate upstream elements, or [enhancers](#), and then interact with the GTFs in order to regulate transcription. These interactions can be mediated by certain intermediary factors or coregulators.

An RNA pol II promoter (or core promoter) is generally composed of two conserved elements, either alone or functioning in synergy: a TATA-box element, usually located 30 nucleotides upstream of the start site in higher eukaryotes, and an initiator element (Inr) encompassing the start site. Promoter recognition is mediated by the binding of the basal factor TFIID through its TBP subunit sitting astride the TATA-box sequence. The role of TFIID is dual: First, by anchoring the PIC to the promoter, it will position RNA pol II and select the startsite; second, it will induce a bend in the DNA structure that reduces [nucleosome](#) packaging from the [chromatin](#) and brings regulatory proteins, already bound to the DNA, closer to the promoter. Transcription from TATA-less promoters also uses TBP; in that case, binding of TFIID results from recognition of Inr by its TAF subunits ([TBP-associated factors](#)) or additional factors.

The association of TFIID with the promoter nucleates the assembly of the PIC, allowing the arrival of the other five GTFs in addition to RNA pol II in a sequential order. Alternatively, the TFIID/promoter complex can be targeted in a single step by the RNA pol II “holoenzyme” complex (Fig. 5c). In addition to RNA pol II, this holoenzyme includes a subset of GTFs, as well as additional factors, such as SRB proteins (for [suppressor of RNA pol B](#)), that constitute the mediator or SRB complex (in yeast) and mediate transcriptional activation. In the stepwise assembly model, TFIIB joins the TFIID-promoter complex, contacting bent DNA both upstream and downstream of the TATA box, thus stabilizing the complex. A specific TFIIB recognition element (BRE), located just upstream of the TATA box, exists at certain promoters. The binding of TBP can also be stabilized by TFIIA, which may relieve inactivation by certain factors, such as MOT1, DR, or NC. The RNA pol II/TFIIF complex then targets the promoter. Completion of the PIC assembly occurs after the successive binding of TFIIE and TFIIH. Opening of the DNA region encompassing the startpoint (nucleotides -9 to +2), which gives the RNA pol II access to the template strand, requires both of these factors, in addition to ATP hydrolysis, except in the case of **supercoiled** DNA. The multisubunit complex TFIIH is the only GTF that possesses several known enzymatic activities: ATP-dependent [DNA helicase](#) activity probably separates the two DNA strands, while its DNA-dependent [Polyadenylation](#) activity could provide part of the energy. In addition, the cdk7 subunit is able to **phosphorylate** the C-terminal domain (CTD) of the largest subunit of RNA pol II. TFIIE not only recruits TFIIH in the PIC but also stimulates TFIIH kinase and ATPase activities. CTD-phosphorylation occurs concomitantly with initiation, converting the unphosphorylated form of RNA pol (IIa) into the hyperphosphorylated elongating form (IIo). Phosphorylation facilitates escape of the elongation complex during promoter clearance by severing contacts between CTD and components of the initiation complex.

TFIIF remains associated with RNA pol II during elongation, together with four additional elongation factors: TFIIIS, p-TEFb, elongin SIII, and ELL. These factors control the rate of mRNA synthesis by regulating the processivity of RNA pol. TFIIIS and p-TEFb allow the arrested complexes to recover elongation activity. The p-TEFb factor, which also exhibits a CTD-kinase activity, is thought to regulate the phosphorylation state of RNA pol II, thus stabilizing the elongation-proficient conformation. Termination requires the presence of two sequences downstream of the transcribed gene: the signal for [polyadenylation](#) of the transcript, followed by a pausing site. Cleavage of the 3' end of the RNA, mediated by three factors (CPSF, CstF, CF), is necessary for termination to occur.



Transcription regulation involves interaction of the sequence-specific factors, bound to their cognate DNA elements, with the basal transcription machinery. This interaction can be either direct or mediated by cofactors like TAFs, SRB proteins, or TIFs (transcription intermediary factors). The specific elements are located at variable distances from the core promoter: Upstream elements (UEs) such as the [CAAT box](#) or [GC box](#) are usually found between positions –50 and –100, whereas enhancers or **silencers** act at distances up to tens of thousands of base pairs, either upstream or downstream of the startpoint. Activation factors usually possess a **DNA-binding** domain (DBD) and an activation domain (AD) (see [-cyclic AMP, cAMP](#)). Some of them are inducible, either newly synthesized or activated, before they stimulate transcription via their response element. Activation includes various mechanisms, such as the recruitment of GTFs onto the promoter or removal of repressor proteins.

Repression of the transcription is usually coupled to DNA compaction into nucleosomes; derepression thus requires chromatin remodeling, most probably involving complexes such as SWI/SNF, which is part of the RNA pol II holoenzyme in yeast, or NURF (nucleosome remodeling factor). Some proteins also act as repressors of transcription by either counteracting activators or impeding the assembly of the PIC.

Over the past decade, the combination of biochemistry and molecular biology has shed some light on our understanding of the molecular bases of human diseases. Unexpectedly, a variety of human genetic disorders are linked to defects in RNA pol II transcription machinery. For example, defects in TFIIF result in xeroderma pigmentosum, Cockayne syndrome, and trichothiodystrophy. Mutation of the VHL **tumor suppressor** gene, the product of which regulates elongin SIII, affects the elongation process and results in the von Hippel–Lindau disorder; and mutations in the **cyclic AMP** response element-binding protein (CREB) binding protein (CBP), which is involved in coactivation process, are found associated with the Rubenstein–Taybi syndrome.

## Bibliography

### Suggestions for Further Reading

S. M. Uptain, C. M. Kane, and M. J. Chamberlin (1997) Basic mechanisms of transcript elongation and its regulation. *Annu. Rev. Biochem.* **66**, 117–172.

### Prokaryotes

C. Yanofsky, K. V. Konan, and J. P. Sarsero (1996) Some novel transcription attenuation mechanisms used by bacteria. *Biochimie* **78**, 1017–1024.

W. G. Haldenwang (1995) The sigma factors of *Bacillus subtilis*. *Microbiol. Rev.* **59**, 1–30.

P. Stragier and R. Losick (1990) Cascades of sigma factors revisited. *Mol. Microbiol.* **4**, 1801–1806.

L. V. Richardson and J. P. Richardson (1996) Rho-dependent termination of transcription is governed primarily by the upstream Rho utilization (rut) sequences of a terminator. *J. Biol. Chem.* **271**, 21597–21603.

J. Greenblatt, J. R. Nodwell, and S. W. Mason (1993) Transcriptional antitermination. *Nature* **364**, 401–406.

T. M. Henkin (1996) Control of transcription termination in prokaryotes. *Annu. Rev. Genet.* **30**, 35–57.

### RNA pol I and III Transcription

E. P. Geiduschek and G. A. Kassavetis (1995) Comparing transcriptional initiation by RNA polymerases I and III. *Curr. Opin. Cell. Biol.* **7**, 344–351.

R. H. Reeder and W. H. Lang (1997) Terminating transcription in eukaryotes: lessons learned

from RNA polymerase I. *Trends Biochem. Sci.* **22**, 473–477.

B. S. Shastry (1993) Gene expression: surprises from the class III side. *Mol. Cell. Biochem.* **124**, 85–89.

### **RNA pol II Transcription**

M. Hampsey (1998) Molecular genetics of the RNA polymerase II general transcriptional machinery. *Microbiol. Mol. Biol. Rev.* **62**, 465–503.

G. Orphanides, T. Lagrange, and D. Reinberg (1996) The general transcription factors of RNA polymerase II. *Genes Dev.* **10**, 2657–2683.

T. Lagrange, A. N. Kapanidis, H. Tang, D. Reinberg, and R. H. Ebright (1998) New core promoter element in RNA polymerase II-dependent transcription: sequence-specific DNA binding by transcription factor IIB. *Genes Dev.* **12**, 34–44.

D. B. Nikolov and S. K. Burley (1997) RNA polymerase II transcription initiation: a structural view. *Proc. Natl. Acad. Sci. USA* **94**, 15–22.

F. Coin and J. M. Egly (1998) Ten years of TFIIH. *Cold Spring Harbor Symp. Quant. Biol.*, Volume **63**.

S. McCracken, N. Fong, K. Yankulov, S. Ballantyne, G. Pan, J. Greenblatt, S. D. Patterson, M. Wickens, and D. L. Bentley (1997) The C-terminal domain of RNA polymerase II couples mRNA processing to transcription. *Nature* **385**, 357–361.

A. S. Neish, S. F. Anderson, B. P. Schlegel, W. Wei, and J. D. Parvin (1998) Factors associated with the mammalian RNA polymerase II holoenzyme. *Nucleic Acids Res.* **26**, 847–853.

V. Ossipow, J. P. Tassan, E. A. Nigg, and U. Schibler (1995) A mammalian RNA polymerase II holoenzyme containing all components required for promoter-specific transcription initiation. *Cell* **83**, 137–146.

R. D. Kornberg (1996) Special Issue: The RNA polymerase II transcriptional machinery. *Trends Biochem. Sci.* **9**, 325–356.

### **Chromatin and Transcription**

L. G. Burns and C. L. Peterson (1997) The yeast SWI–SNF complex facilitates binding of a transcriptional activator to nucleosomal sites *in vivo*. *Mol. Cell. Biol.* **17**, 4811–4819.

G. Mizuguchi, T. Tsukiyama, J. Wisniewski, and C. Wu (1997) Role of nucleosome remodeling factor NURF in transcriptional activation of chromatin. *Mol. Cell* **1**, 141–150.

A. P. Wolffe, J. Wong, and D. Pruss (1997) Activators and repressors: making use of chromatin to regulate transcription. *Genes Cells* **2**, 291–302.

### **Regulation in Eukaryotes**

M. Keaveney and K. Struhl (1998) Activator-mediated recruitment of the RNA polymerase II machinery is the predominant mechanism for transcriptional activation in yeast. *Mol. Cell* **1**, 917–924.

T. I. Lee and R. A. Young (1998) Regulation of gene expression by TBP-associated proteins. *Genes Dev.* **12**, 1398–1408.

K. Struhl and Z. Moqtaderi (1998) The TAFs in the HAT [comment]. *Cell* **94**, 1–4.

### **Transcription Syndromes**

T. Aso, A. Shilatifard, J. W. Conaway, and R. C. Conaway (1996) Transcription syndromes and the role of RNA polymerase II general transcription factors in human disease. *J. Clin. Invest.* **97**, 1561–1569.

## Transcription Factors

The expression of eukaryotic gene is controlled tightly with respect to both time and space. While there are many control points during gene expression, the control of initiation of transcription is, for most genes, the most prominent among them.

The production of RNA from a DNA template is catalysed by RNA polymerases (RNAP). The RNA polymerase holoenzyme of *Escherichia coli*, a protein of molecular weight 465 kD, is composed of a core enzyme with the subunit composition  $\alpha_2\beta\beta'$ , which is responsible for RNA polymerisation (1, 2), and a sigma factor that directs the holoenzyme to the promoters of genes (3-5).

The most simple types of RNA polymerases are found in bacteriophages. The RNAP of T7 bacteriophage is a single subunit enzyme of molecular weight 99 kD (6), but it is composed of domains with specific functions (7). For example, the relative locations of the active site and the domain for promoter recognition allow the orientation of the polymerase on the template to be defined. The specificity of T7 RNAP is controlled by a loop consisting of approximately 40 amino acid residues. A single point mutation within this loop is sufficient to switch the specificity of the enzyme from the T7 promoter sequence to that of the T3 promoter (8).

While all prokaryotic RNA polymerases are capable of recognizing gene promoters and initiating transcription from the proper start point, the efficiency of transcriptional initiation is further modulated by trans-acting factors known as transcription factors. Transcription factors are called repressors or activators, depending on whether their effect is to decrease or increase the level of transcription. Because of relatively high efficiency of RNA synthesis by prokaryotic RNAPs, the main theme for regulation in prokaryotes is repression.

Eukaryotes have three distinct RNA polymerases that catalyze nuclear gene transcription (9). RNAP-I and -III catalyze the production of ribosomal RNA and transfer RNA, respectively, while RNA-II [the structure of this multisubunit protein has recently been solved at a resolution of 2.8 Å (10)] is responsible for the synthesis of all messenger RNA. The eukaryotic RNAPs have between 8 and 14 subunits and a molecular weight of approximately 500 kD. The purified enzymes can catalyse low-level template-dependent transcription of RNA, but, despite their structural complexity, these enzymes are unable to recognize the promoter of genes and depend fully on a set of auxiliary proteins known as general transcription factors to initiate transcription from the corresponding class I, II, and III gene promoters (11).

In the case of transcription by RNAP-II, the basal transcription factors necessary for low-level transcriptional initiation from the proper start point have been separated into biochemically defined fraction referred to as TFIIA, -B, -D, -E, -F, H, and -J (12-15). Subsequent cloning and biochemical analyses have shown that many of these fractions consist of multiple proteins. For example, TFIID consists of the TATA box binding protein (TBP) and 8 to 12 tightly bound TAFs (TBP-associated factors) (16, 17). A low resolution structure of TFIID was obtained by electron microscopy and single particle analysis (18, 19). TFIID appears to adopt the structure of a horseshoe with a deep groove that most likely should accommodate the DNA.

The stupefying complexity of biological processes like embryonal development make it obvious that gene expression must be exceedingly tightly controlled (e.g., see Refs 20-22). For the initiation of RNA synthesis, this is achieved through the action of transcription factors. These proteins bind with high affinity to their target DNA sequences and regulate the rate of transcription through interaction with the basal transcriptional machinery. Because of the low efficiency of the basal transcriptional

machinery, most eukaryotic transcription factors are activators.

The specificity of transcriptional initiation is controlled through the interaction of transcription factors with their DNA-binding sites in the promoter and enhancer regions of genes, so it is vital to understand the forces that govern the recognition of DNA by proteins. Nature has used a modular approach to assemble transcription factors. The individual domains within a transcription factor, such as, for example, the DNA-binding domain and the transcriptional activation domain, are thereby, to a first approximation, independent of each other. As a consequence, the properties of the DNA-binding domains can be studied in the absence of the activation domains. The limited amount of structural data available for activation domains is in sharp contrast with the wealth of information about the DNA-binding domains of transcription factors.

## 1. DNA-Binding Motifs Observed in Eukaryotic Transcription Factors

One of the central observations that emerged from X-ray crystallographic and NMR-spectroscopic studies and sequence comparison is that most DNA-binding proteins can be grouped into families based on the structural motif that they rely on for sequence-specific DNA recognition (23). Each of these motifs contains a simple element of secondary structure complementary to the structure of the DNA (24).

## 2. Helix-Turn-Helix Motif

The helix-turn-helix (HTH) motif was the first DNA-binding motif to be discovered, and extensive structural and functional studies have been performed with many HTH proteins (25, 26), including  $\lambda$ -phage and 434-phage Cro repressors (27, 28), the DNA-binding domains of  $\lambda$ - and 434-phage repressors (29-32), Lac repressor (33-36), that of the *trp* operon (37-40), the *E. coli* transcriptional activator cyclic AMP receptor protein (CRP)/catabolite gene activator protein (CAP) (41), *E. coli* inversion stimulation factor (FIS) (42), *Salmonella* *hin* recombinase (43), and the *E. coli* biotin repressor (44). Comparison of the structures of these proteins revealed a shared DNA recognition motif consisting of an  $\alpha$ -helix, a turn of four amino acid residues, and a second helix (Fig. 1) (26) (see Helix-Turn-Helix Motif). The structural architecture of the HTH-motif relies on a hydrophobic core through which the two  $\alpha$ -helices pack together. The second helix (recognition helix) interacts with the major groove of the DNA. Interestingly, the isolated HTH motif is not a stable protein domain and, unlike most other DNA binding motifs, relies on other parts of the proteins for stability.

**Figure 1.** The global structure of the DNA complex of the N-terminal domain of  $\lambda$ -repressor (29). The protein binds to the DNA as a dimer with one monomer contacting each half site of the  $O_L1$  operator. The recognition helix fits snugly into the minor groove of the DNA, while the second helix of the HTH-motif binds across the major groove. Note that the HTH motifs are stabilized through interactions with the remainder of the repressor subunits. The N-terminal extension of one subunit wraps around the DNA to contact the minor groove. The HTH-motifs of each monomer are a darker shade of grey.  $\alpha$ -Helical regions of the protein are represented as cylinders.



The crystal structure of MarA, a member of the AraC prokaryotic transcriptional activator family, shows a bipartite HTH motif (45). This motif is composed of seven  $\alpha$ -helices and folds into two structurally similar subdomains, which both contain an HTH motif and which are connected by helix-4. The recognition helices of either HTH-motif are inserted into adjacent major groove segments on the same face of the DNA. Because of the rather low separation of 27Å, the DNA is bent by 35°. The DNA complex of the *E. coli* transcription factor Rob displays a similar architecture, but only the N-terminal HTH domain engages the major groove of the unbent DNA (46).

### 3. Homeodomain

If the definition of the HTH motif is relaxed to also allow for longer loops, then many eukaryotic transcriptional regulators can be included. For example, the structures of the yeast mating type protein  $\alpha$ -MAT (47, 48), as well as those of the *Drosophila* homeotic proteins Engrailed (49) and Antennapedia (50-52) revealed that they contained an HTH motif. Unlike the prokaryotic HTH domains, isolated homeodomains can fold correctly and bind to DNA with specificity similar to that of the parent protein (24, 53). The overall structure of these homeodomains is visualized rather easily (Fig. 2) (24, 54, 55). Helices 1 and 2 pack against one another in an antiparallel fashion. Helix 3 is amphiphatic, and the hydrophobic face packs against the first two  $\alpha$ -helices to which it is oriented perpendicular. Helices 2 and 3 form the HTH motif. The main contacts between homeodomains and DNA are made by residues of helix 3, which lies in the major groove of the DNA and by amino acids in an N-terminal extension.

**Figure 2.** The overall structure of the DNA complex of the homeodomain of the *Drosophila* protein engrailed (49). The amino-terminal arm of the homeodomain contacts the minor groove of the DNA. Helix 3 contacts the major groove of the DNA and, together with helix 2, constitutes the HTH-motif. The  $\alpha$ -helical regions of the protein are represented as cylinders.



Winged HTH-proteins include two or more b-sheets packed against the core of the three a-helices that form the HTH-motif. An extended loop that connects the two b-sheets adjacent to the recognition helix is involved in DNA-binding through contacts to the phosphate backbone. Members of this class include Elk-1 (56), SAP-1 (57), and PU.1 (58) of the ETS family of transcription factors and IRF-1 of the interferon regulatory factors (59).

Besides proteins containing a homeodomain, the eukaryotic HTH family of transcriptional regulators also includes the POU domain proteins (60), yeast heat shock regulatory factor (61), histone H5 (62), rat hepatocyte nuclear factor 3g (63), and the product of the murine c-Myb proto-oncogene (64, 65). The DNA binding motif of Oct-1 contains two HTH-motifs, namely a canonical homeodomain and a POU-specific domain, that, despite poor sequence homology, resembles the HTH motif of the repressors of phage 1 and 434, except for an extended loop-structure connecting the two helices (66, 67).

The paired domain, which is found in transcription factors such as paired and gooseberry from *Drosophila* and the mammalian Pax1 to Pax9, contains two structurally independent domains (68, 69). The N-terminal domain comprises a short region of antiparallel b-sheet followed by a type II b-turn, three a-helices with a fold resembling the classical homeodomain, and an extended C-terminal tail. The C-terminal domain contains a fold that is similar to a homeodomain.

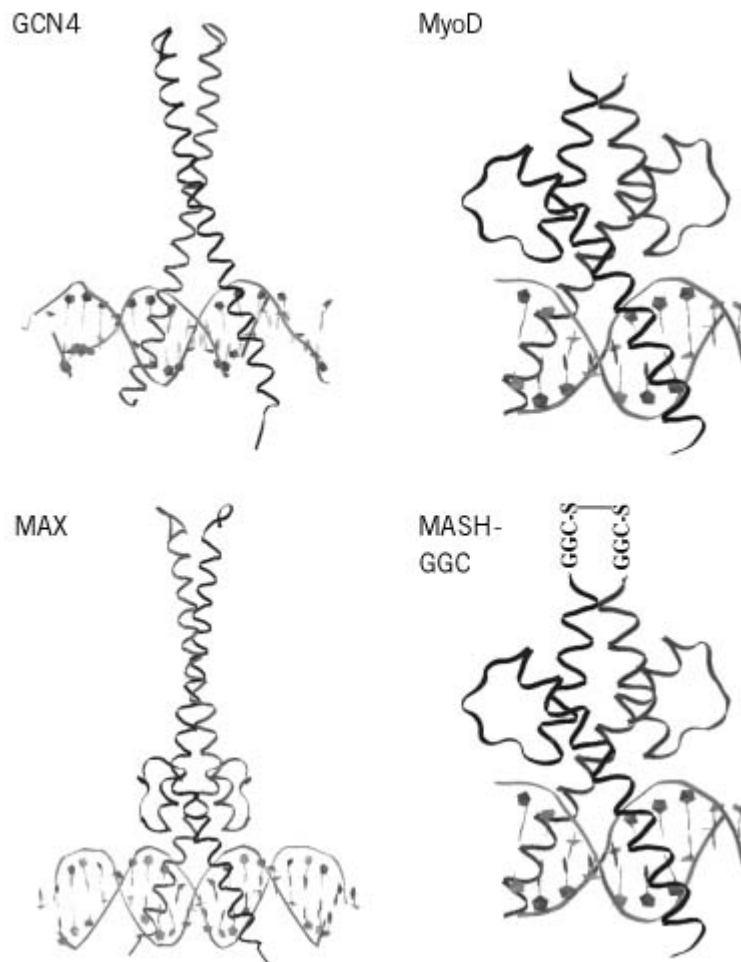
#### 4. MADS-Box Proteins

These proteins have the function of DNA-binding and dimerization integrated into a single protein design. They bind to DNA as a homodimer, and the primary DNA-binding element is an antiparallel coiled coil of the two amphiphatic a-helices (helix-1), one from each subunit. The basic residues of these two a-helices fit into the major groove of the DNA double helix (see MADS-box proteins).

## 5. Basic Leucine Zipper, Basic Helix-Loop-Helix, and Basic Helix-Loop-Helix Zipper Motifs

The proteins of the basic leucine zipper (BZ), the basic helix-loop-helix (BHLH), and the basic helix-loop-helix-zipper (BHLHZ) families of transcription factors use similar strategies for DNA binding. They all bind to DNA as dimers through a N-terminal recognition helix consisting of clusters of basic residues. The basic region is directly adjacent to the dimerization motif. X-ray structures of GCN4 (70) and Pap1 (71) show that dimerization is mediated by the zipper domain, which contains a heptad repeat of leucine residues (Fig. 3) (70-73). The leucine residues from two zipper domains pack tightly against each other, thereby stabilizing the overall structure of a coiled coil.

**Figure 3.** The structures of the DNA complexes of the DNA-binding domains of the BZ protein GCN4 (70), the BHLH protein MyoD (74), the BHLHZ protein Max (76), and the engineered MASH-GGC (139). For MASH-GGC, the tripeptide Gly-Gly-Cys was added to the C-terminal end of the BHLH subunit.



Structural of the BHLH proteins MyoD (Fig. 3) (74) and E47 (75) complexed to DNA showed that, in this class, dimerization occurs through a left-handed parallel four-helix bundle. The helices are connected by extended loops.

Combination of the BZ motifs and BHLH motifs led to the BHLHZ family, characterized by a

leucine zipper domain that extends from the C-terminus of the BHLH domain (Fig. 3) (76, 77). The zipper domain adds substantial buried surface area to the dimerization domain, thereby increasing both the stability and the DNA binding specificity of the protein (*vide infra*) (see Leucine zipper and helix-loop-helix motif).

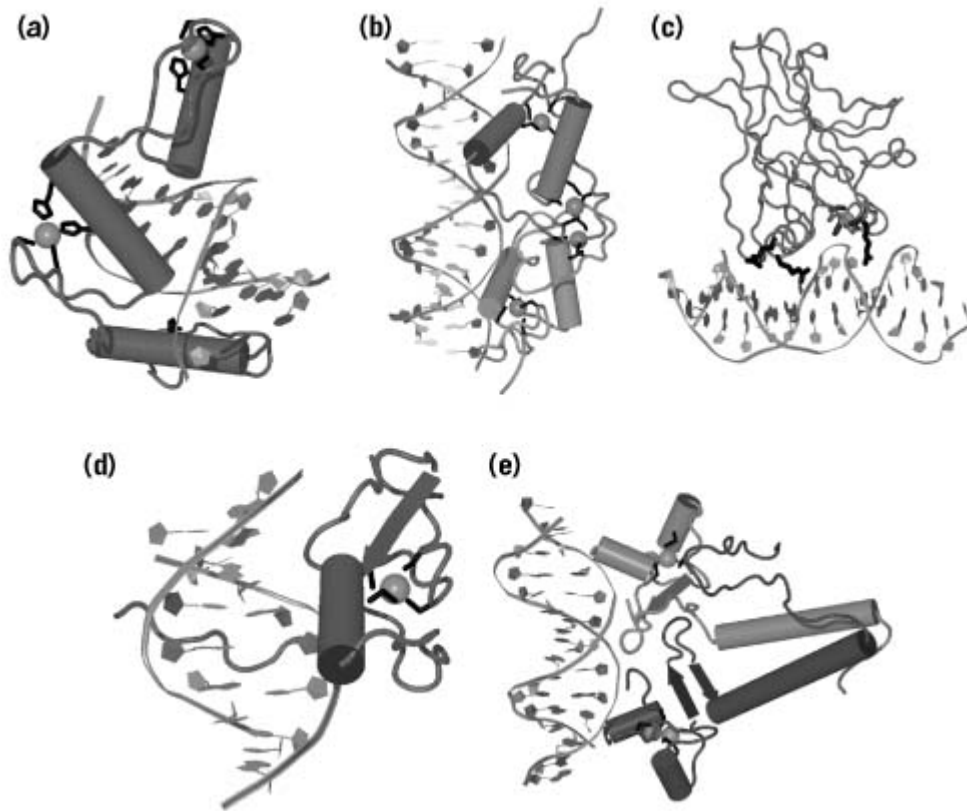
## 6. Zinc Containing DNA-Binding Domains

The family of zinc-bearing DNA-binding motifs is the largest group of eukaryotic transcription factors, and it has been estimated that perhaps 1% of all mammalian genes encode zinc-finger proteins (78, 79). Their abundance might be a consequence of the structural economy of stabilizing modular protein folds through metal ions (80). In sharp contrast to the structural similarities observed within the DNA binding families described above, diversity appears to be the hallmark of the zinc-containing DNA-binding domains. To date, seven subfamilies of Zn-containing proteins have been characterized:

1. Zinc fingers have the same canonical bba structure, stabilized through a hydrophobic core and a tightly bound zinc ion that interacts with two cysteine residues of the second strand and two histidine residues that are part of the  $\alpha$ -helix (Fig. 4a). Finger specificity seems to depend on three amino acid residues just before and within the  $\alpha$ -helix that make specific contacts to the DNA (81-87). Zinc fingers occur in multiples connected by flexible linkers and can be viewed as independent DNA-binding units. The structures solved comprise ADR1 (88) and SWI5 (89) from yeast, murine Zif 268 (81, 90), the human factors MBP-1 (91), and GLI (92), *Xenopus* Xfin (93), and the *Drosophila* tramtrack protein (94).

**Figure 4.** Structures of zinc-containing transcription factors. The  $\alpha$ -helices and  $\beta$ -sheets are represented as **cylinders** and **arrows**, respectively. (a) The global structure of the DNA complex of the three zinc finger peptide of murine Zif268 (81, 90). The three zinc fingers track along the major groove of the DNA, which is contacted through residues of the recognition  $\alpha$ -helix, which is part of the bba structure. This domain is held together through a tightly packed hydrophobic core and a Zn ion coordinated to two cysteines and two histidines. (b) Overall structure of the dimeric DNA-binding domain of estrogen receptor complexed with DNA (98, 99). In each monomer, one Zn ion stabilizes the conformation of the recognition helix, which fits into the major groove of the DNA, while the second Zn helps stabilize the dimer interface. Each of the Zn ions is coordinated by the side chains of four cysteine residues. (c) Ribbon drawing of the DNA complex of the p53 tumor suppressor core domain (108). The core domain structure consists of a  $\beta$ -sandwich that serves as a scaffold to orient a loop-sheet-helix motif and two large loops, the conformation of which is stabilized by a tetrahedrally coordinated zinc. The zinc is bound by the side chains of three cysteines and one histidine residue. The Zn domain orients the side chain of an arginine residue that contacts the minor groove of the DNA, while two arginine residues of the loop-sheet-helix motif mediate contacts with the major groove. (d) Illustration of the DNA complex of the DNA binding domain of the chicken erythroid transcription factor GATA-1, (103). A single Zn ion is used to stabilize the core domain of GATA-1, which consists of two irregular antiparallel  $\beta$ -helix. The  $\alpha$ -helix and the loop connecting the two antiparallel sheets interact with the major groove of the DNA, while the carboxy-terminal tail wraps around the DNA to contact the minor groove. The  $\beta$ -sheets and the  $\alpha$ -helices of GATA-1 are represented as **arrows** and **cylinders**, respectively. (e) The structure of the DNA complex of the yeast transcriptional activator PPR1 (106). Two zinc ions are coordinated by six cysteines in each of the two PPR1 subunits. The two central cysteine side chains, thereby, contact both zincs. The two zinc domains lie above the major groove of the DNA, and the DNA contacts are mediated through the residues of the  $\alpha$ -helical recognition helices.



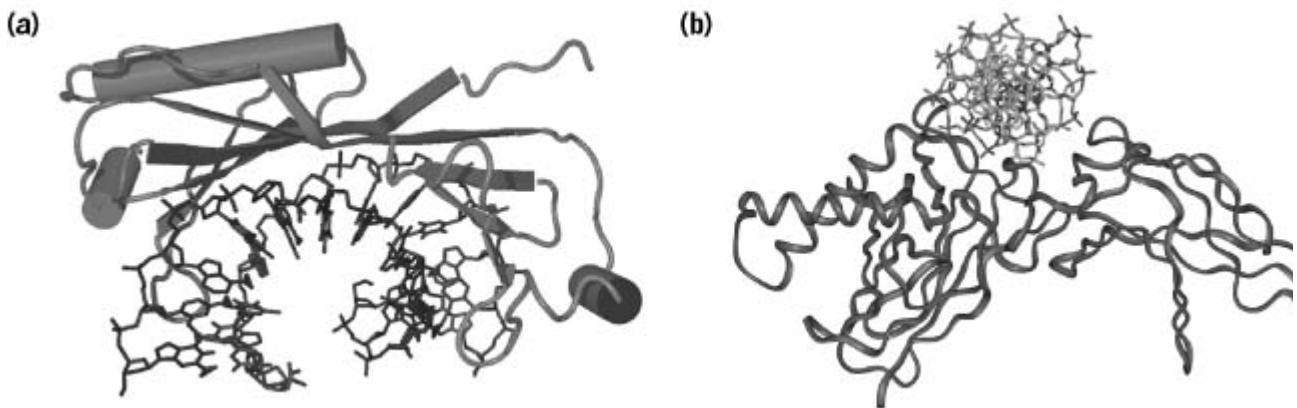


2. The structures of the DNA-binding domains of the glucocorticoid (95-97) and estrogen (98, 99) receptors, the retinoic acid receptor-b (100), and the retinoid X receptor-a (101) revealed a second type of Zn-containing DNA binding motif (Fig. 4b). It has eight cysteine residues; four from each of two peptide loops bind to two zinc ions (78, 95-102). The cysteine residues are, thereby, located C-terminal to an  $\alpha$ -helix that fits into the major groove of the substrate.
3. An  $\alpha$ -helix located just N-terminal of the cluster of four cysteine residues in the GATA-1 group of Zn-containing DNA-binding proteins also interacts with the major groove (Fig. 4d) (103). The GATA-1 domain is the only known Zn-containing DNA-binding protein that relies on only a single metal core to bind to DNA with high affinity.
4. The yeast proteins GAL4 (104, 105) and PPR1 (Fig. 4e) (106), as well as the fungal protein PUT3 (107), also use an  $\alpha$ -helix for interaction with the major groove of the DNA; they contain a unique Zn-binding motif in which six cysteine residues are coordinated to two spatially close Zn-ions (104-106). Two of the cysteines bind simultaneously to both zinc atoms.
5. The tumor repressor p53 (108) (Fig. 4c). The core domain structure of the tumor suppressor p53 (108) consists of a  $\beta$ -sandwich that serves as a scaffold for two large loops, which are held together, in part, by a tetrahedrally coordinated zinc and which interact with the major groove of the DNA and for a loop-sheet-helix motif (Fig. 4c).
6. The human transcriptional elongation factor TFIIS (109, 110). NMR studies revealed that the Cys4 nucleic acid binding domain of the human transcriptional elongation factor TFIIS consists of a three-stranded  $\beta$ -sheet ( $\beta$ -ribbon) (109, 110). The X-ray structure of yeast RNA polymerase II revealed six zinc-binding proteins. The C-terminal Zn-binding module of Rpb9 and Rpb12 fold into a ribbon motif similar to that described for TFIIS (111, 112).
7. Subunit Rpb10 of RNA polymerase from *Methanobacterium thermoautotrophicum* contains an unusual CysX<sub>2</sub>CysX<sub>N</sub>CysCys Zn binding motif. The Zn atom is located at one end of an unusual cluster of three  $\alpha$ -helices (113).

## 7. DNA Recognition via $\beta$ -Sheets

Whereas most of the proteins described above rely on an  $\alpha$ -helix for the formation of contacts with the major groove of the DNA, the prokaryotic repressors P22 Arc (114-116), P22 Mnt (117, 118), and MetJ (119, 120) use a pair of anti-parallel  $\beta$ -sheets to contact the major groove. No eukaryotic protein is known that uses a  $\beta$ -sheet to contact the DNA's major groove. However, the X-ray structures of DNA complexes of the TATA box-binding proteins (TBPs) (121-124) showed that DNA binding is mediated by TBP's curved eight-stranded antiparallel  $\beta$ -sheet, which provides a concave surface for contacts with the bases in the minor groove and with the phosphate backbone (Fig. 5A). The chromosomal protein Sac7d of the hyperthermophile *Sulfolobus acidocaldarius* uses the triple stranded  $\beta$ -sheet ( $\beta_3$ - $\beta_4$ - $\beta_5$ ), which is part of a larger five-stranded incomplete  $\beta$ -barrel, capped at the opening by a C-terminal  $\alpha$ -helix (125).

**Figure 5.** **Figure (a)** Global conformation of the TBP-DNA complex (122). The protein binds the TATA box from the minor groove side and bends the DNA into the major groove. The  $\alpha$ -helices are represented as **cylinders**, and the  $\beta$ -strands are represented as **arrows** **(b)** Structure of the human NF- $\kappa$ B p52 homodimer-DNA complex (128). DNA binding relies on a  $\beta$ -barrel with an s-type immunoglobulin fold and is mediated by protein loops.



Runx1-CBF $\beta$  (126), NFAT (127), NF- $\kappa$ B (128), p53 (129), STAT (130), and the Brachyury T-Box family of proteins (131) rely on a  $\beta$ -barrel with an s-type immunoglobulin fold for DNA binding. DNA binding is mediated by loops that extend from one end of the Ig-motif (Fig. 5B).

## 8. Interactions between Transcription Factors

The observation that eukaryotic transcriptional regulators bind to relatively short DNA sequences with modest DNA-binding specificity indicates that nature has chosen a combinatorial approach to regulate transcriptional initiation in eukaryotes. Thus, many genes can be differentially regulated through the use of various combinations of relatively few transcription factors. If the binding reactions of individual transcription factors to their respective sites are independent of one another, the specificities multiply (corresponding to the sum of the free energies of the individual binding reactions). For example, the binding reaction of a MEF-2C dimer and a MyoD/E12 heterodimer to a myogenin minimal promoter containing only an E-box and a MEF-2C binding site (22) would be characterized by a DNA-binding specificity of  $1.2 \times 10^3$ .

DNA-binding reactions, however, often show a high degree of cooperation. A particularly well-characterized case is the cI repressor of bacteriophage  $\lambda$ . It binds to the three 17 bp binding sites  $O_R1$ ,  $O_R2$ , and  $O_R3$  of the right  $\lambda$ -operator with different affinities and specificities. Therefore, if cI

repressor binds to an artificial DNA sequence containing  $O_R1$ ,  $O_R2$ , and  $O_R3$  separated by linkers long enough to “isolate” the binding sites from one another, the affinity and specificity observed are simply the product of the individual binding reactions. However, the binding reaction of cI repressor to the wild-type operator  $O_R$  is highly cooperative, and the quaternary complex with the wild-type sequences is stabilized by 4.8 kcal/mol relative to the case where the binding sites are isolated (132). Of course, cooperative binding relies on a properly folded repressor, and, therefore, it is likely that the specificity of binding is also increased approximately 4,000-fold because the non-specific complex is simply a loose association between the proteins and the DNA. Similarly, two Arc repressors bind cooperatively to the 21 bp Arc operator. When the left or the right operator half-site is occupied by an Arc dimer, the affinity of the second dimer is increased approximately 5,900-fold (133).

The eukaryotic transcriptional activator MyoD binds cooperatively to two sites in muscle-specific enhancer sequences (134). It is known, however, that the specific and non-specific DNA complexes of BHLH-proteins like MyoD are rather similar (135-137). The increases in the DNA-binding affinity might, therefore, not lead to an increase in the DNA-binding specificity. Interestingly, the DNA-binding specificity of basic helix-loop-helix (BHLH) proteins can be affected by alterations of the dimerization motif. The DNA-binding specificity of the bHLHZ transcription factor USF, for example, is dependent on the presence of the coiled-coil domain because its removal led to a dramatic reduction of the DNA binding specificity (138). The BHLH-protein MASH-1 binds to DNA with low sequence specificity (136). When the tripeptide Gly-Gly-Cys was added to the C-terminal end of the BHLH domain, however, a disulfide bond could be formed between the Cys residue of each monomer (Fig 3), and the affinity for E-box containing DNA was increased, while the stability of the complex with heterologous DNA was reduced (139). A similar increase in the DNA binding specificity of MASH-1 was obtained when the C-terminus of the first BHLH subunit was covalently linked to the N-terminus of the second through the introduction of a peptide linker, to generate a “single-chain dimer” (140). Circular dichroism spectroscopy revealed that, like wild type MASH-BHLH, the “single-chain dimers” of MASH-1 adopted only partly folded structures in the absence of DNA, but they underwent a folding transition to mainly  $\alpha$ -helical conformations on DNA binding. The affinity of the “single-chain dimers” for E-box containing DNA sequences was increased almost 30-fold when compared with wild-type MASH-BHLH, while the stability of the complex with heterologous DNA was only 10 times greater. The free energy of transferring a protein molecule from a nonspecific to a specific site was, therefore, increased more than threefold through the introduction of the peptide linker.

Covalently linking the subunits of MASH-1, either through a disulphide bond or through a peptide linker, might reduce the number of conformations accessible to the protein in the disordered state and, therefore, diminish the entropic penalty that accompanies folding and DNA-binding. Limiting the number of accessible conformations of the basic region of BHLH-proteins could stabilize the complex with specific DNA and destabilize the complex with nonspecific DNA sequences, thereby increasing the overall DNA-binding specificity. (NMR) spectroscopic studies showed that, while the flexibility of side chains at the protein-DNA interface is restrained when compared with the uncomplexed state, it is still greater than the flexibility of side chains in the protein core (141). A precedent for such a mechanism is provided by the tryptophane repressor of *E. coli*. Replacing Ala77 with valine results in a local stabilization of the DNA-binding domain of the repressor, reduced conformational flexibility (142), and increased specificity through reduced affinity for nonspecific DNA (143). Similarly, metallopeptide complexes containing two copies of a peptide comprising the basic and the spacer region of the bZ protein GCN4 assembled into a dimer through a bis(terpyridyl) iron(II) complex displayed increased site specificity when compared with native GCN4 (144). A similar observation was made with a single-chain variant of the P22 Arc repressor, in which the subunits were connected through a 15-residue peptide linker (145).

Substitution of Cys for Val152 of the Lac repressor, which permits disulfide formation between the two N-terminal DNA binding domains, resulted in an increase in the stability of the operator

complex by approximately six-fold relative to the wild-type protein (146). This increased affinity derives from the decreased entropic cost in properly orientating the two N-terminal domains for operator binding.

Within an eukaryotic cell, similar increases of the DNA binding affinity and specificity might result from the interaction with other components of the transcriptional machinery. While one side of helix-1 of BHLH proteins directly contacts the major groove of DNA, the other side is exposed to the solvent providing an ideal contact surface for additional proteins. The significantly reduced DNA binding specificity of BHLH proteins when compared with proteins containing an HTH-motif might, in part, be the consequence of the fact that the conformation of the recognition helix in the latter is stabilized through interactions with other parts of the protein (26). Such interactions should reduce the conformational flexibility of the recognition helix, thereby increasing the specificity of the proteins.

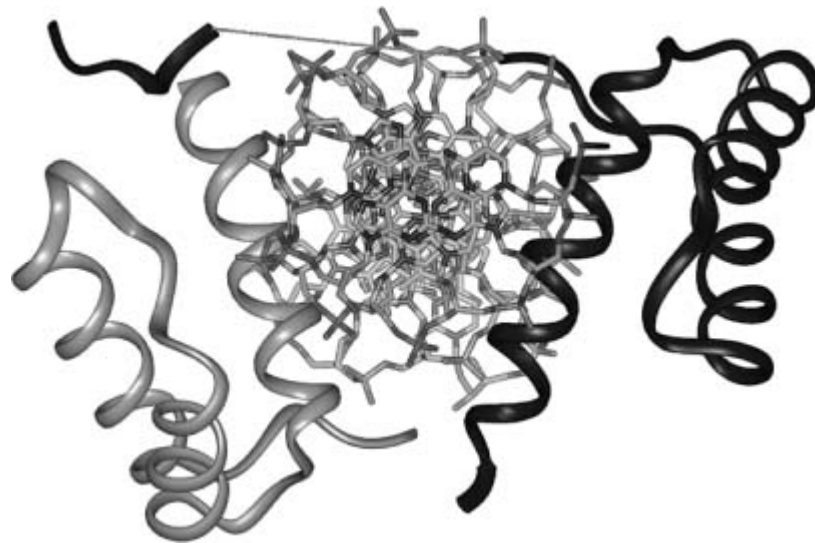
MASH-1, a BHLH protein that promotes the differentiation of neuronal precursor cells, could be converted into a protein capable of inducing myogenesis in fibroblasts by replacing Leu130 of MASH-1 with lysine and introducing an additional turn into its basic recognition helix (147). These changes did not significantly alter the DNA binding properties of the proteins in cell free conditions. Crystallographic data for the DNA complexes of MyoD suggested that Leu130 pointed away from the DNA into the solvent. The identity of the amino acid in position 130 appeared to be important for protein-protein interactions that might affect the DNA binding specificities displayed by BHLH-proteins in vivo and could, thereby, form the molecular basis of the different physiological properties of the myogenic and neurogenic BHLH proteins.

Similarly, both the human T-cell leukemia virus type-1 Tax protein (148-150) and the hepatitis B virus protein pX (151) affect the DNA binding properties of their target proteins. The nonhomologous proteins Tax and pX interact with the basic region of bZ proteins to enhance their DNA-binding affinity and discriminate between peptide-DNA complexes based on the identity of the DNA target site (148, 152, 153). Consistent with the model presented above, binding of Tax and pX to various bZ proteins also changed their relative affinities for different DNA-binding sites (150, 151), thereby providing a molecular rationale for the altered transcriptional patterns observed in infected cells (154-157).

The Hox-gene products are a set of transcription factors that are important for the embryonic development of multicellular organisms (158). Surprisingly and contrary to the high specificity of their biological function, many Hox proteins, such as *Drosophila* Ultrabithorax (Ubx), display only modest DNA binding specificity on their own. However, together with the homeoprotein Extradenticle (Exd) (or Pbx1, which is the mammalian homologue of Exd), they bind DNA as heterodimers with significantly increased affinity and specificity (159, 160). The homeodomains of the heterodimer bind in a tandem manner to opposite faces of the DNA. The co-operative interactions arise from the YPWM amino-acid motif of Ubx that forms a reverse turn and inserts into a hydrophobic pocket on the Exd homeodomain surface (161) (Fig. 6). Similarly, the homeodomain protein MATA2 and the MADS box protein MCM1 exemplify combinatorial control in determining the mating type in *Saccharomyces cerevisiae* (162). The otherwise flexible amino-terminal extension of the MATA2 homeodomain forms a b-hairpin in the complex that grips the MCM1 surface through parallel b-strand hydrogen bonds and hydrophobic interactions (163). This interaction, together with MCM1 induced DNA bending, which brings the two proteins into closer contact, forms the basis of the observed cooperation. Similar cooperative binding was observed for the ternary DNA complexes of SAP-1 and SRF (164), MATA2 and Matal (165), an OCA-B (OBF-1, Bob1) peptide and the Oct-1 POU domain (166), and NFAT1 and the heterodimer of Fos and Jun (127).

**Figure 6.** The DNA complex of Ubx and Exd (161). The Ubx homeodomain is shown in **black** and the Exd homeodomain in **grey**. The **dissolved line** represents the disordered linker, which connects the YWPM motif to the

Ubx homeodomain. This motif contacts a hydrophobic pocket of the Exd homeodomain, which is the source of the cooperative interaction between the two proteins.



## 9. Summary

A thorough chemical and physical understanding of DNA recognition by transcription factors requires a combination of structural and thermodynamic studies. The precise recognition of a defined DNA sequence by a given transcription factor necessitates an optimal shape complementarity between the interacting species. The binding reactions are, thereby, characterized through the formation of large complementary surfaces at the interface between the DNA and the protein. While rigid body association has been observed, significant conformational changes of both the DNA and the protein often occur on complex formation to ensure an optimal fit. It is, therefore, important to understand both the intrinsic and the induced properties of the reactants.

While most prokaryotic transcriptional regulators bind to DNA with high specificity, many eukaryotic transcription factors display only modest DNA-binding specificity. The amount of DNA-binding specificity of a protein depends on the conformational flexibility of its DNA recognition element. The recognition helix of the prokaryotic HTH motifs, for instance, is stabilized through the interaction with other parts of the protein (indeed, the HTH motif adopts a stable conformation only in the structural context of the whole protein). The recognition helices of the eukaryotic BHLH and BZ motifs, on the other hand, adopt well-defined structures only upon binding to DNA. However, even in the complexes, large parts of the recognition helices face the solvent. The conformational flexibility of such open recognition elements can be reduced and their DNA-binding specificity altered through the interaction with other components of the transcriptional machinery. Therefore, transcriptional regulation in eukaryotes relies to a large part on multiprotein complexes with the potential for combinatorial interactions.

## Bibliography

1. A. Polyakov, E. Severinova, and S. A. Darst (1995) *Cell* **83**, 365–373.
2. S. A. Darst, J. W. Roberts, A. Malhotra, M. Marr, K. Severinov, and E. Severinov (1997) In *Nucleic Acids and Molecular Biology* (D. M. J. Lilley and F. Eckstein eds.), Springer, Berlin, Heidelberg, New York, pp. 27–40.
3. J. D. Helmann and M. J. Chamberlin (1988) *Annu. Rev. Biochem.* **57**, 839–872.
4. A. Ishihama (1998) *Trend Genet* **4**, 282–286.
5. A. Ishihama (1997) In *Nucleic Acids and Molecular Biology* (D. M. J. Lilley and F. Eckstein,

- eds.), Springer Verlag, Berlin, Heidelberg, New York, pp. 53–70.
6. P. Davanloo, A. H. Rosenberg, J. J. Dunn, and F. W. Studier (1984) *Proc. Natl. Acad. Sci. U.S.A.* **81**, 2035–2039.
  7. R. Sousa, Y. Je Chung, J. P. Rose and B.-C. Wang (1993) *Nature* **364**, 593–599.
  8. C. A. Raskin, G. Diaz, K. Joho and W. T. McAllister (1992) *J. Mol. Biol.* **228**, 506–515.
  9. A. Sentenac (1985) *CRC Crit. Rev. Biochem.* **18**, 31–90.
  10. P. Cramer, D. A. Bushnell, and R. D. Kornberg (2001) *Science* **292**, 1863–1876.
  11. J. J. Eloranta and S. Goodbourn (1996) "Frontiers in molecular biology" In *Eukaryotic Gene Transcription* (S. Goodbourn, ed.) IRL Press at Oxford University Press, Oxford, pp. 1–33.
  12. R. Weinmann (1992) *Gene Expression* **2**, 81–91.
  13. R. C. Conaway and J. W. Conaway (1993) *Annu. Rev. Biochem.* **62**, 161–190.
  14. P. W. J. Rigby (1993) *Cell* **72**, 7–10.
  15. L. Zawel and D. Reinberg (1993) *Prog. Nucl. Acid. Res. Mol. Biol.* **44**, 67–108.
  16. W. P. Tansey and W. Herr (1997) *Cell* **88**, 729–732.
  17. S. K. Burley and R. G. Roeder (1996) *Annu. Rev. Biochem.* **65**, 769–799.
  18. M. Brand, C. Leurent, V. Mallouh, L. Tora, and P. Schultz (1999) *Science* **286**, 2151–2153.
  19. F. Andel III, A. G. Ladurner, C. Inouye, R. Tijan, and E. Nogales (1999) *Science* **286**, 2153–2156.
  20. J. Whiting, H. Marshall, M. Cook, R. Krumlauf, P. W. J. Rigby, D. Stott, and R. K. Allemann (1991) *Genes Dev.* **5**, 2048.
  21. P. A. Lawrence (1992) *The Making of a Fly: The Genetics of Animal Design*, Blackwell Scientific Publications, Oxford.
  22. S.-P. Yee and P. W. J. Rigby (1993) *Genes Dev.* **7**, 1277–1289.
  23. S. C. Harrison (1991) *Nature* **353**, 715–719.
  24. C. O. Pabo and R. T. Sauer (1992) *Annu. Rev. Biochem.* **61**, 1053–1095.
  25. R. G. Brennan and B. W. Matthews (1989) *J. Biol. Chem.* **264**, 1903–1906.
  26. S. C. Harrison and A. K. Aggarwal (1990) *Annu. Rev. Biochem.* **59**, 933–969.
  27. R. G. Brennan, S. L. Roderick, Y. Takeda, B. W. Matthews (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8165–8169.
  28. A. Mondragon, C. Wolberger, and S. C. Harrison (1989) *J. Mol. Biol.* **205**, 179–188.
  29. L. J. Beamer and C. O. Pabo (1992) *J. Mol. Biol.* **227**, 177–196.
  30. J. E. Anderson, M. Ptashne, S. C. Harrison (1987) *Nature* **326**, 846–852.
  31. A. K. Aggarwal, D. W. Rodgers, M. Drottar, M. Ptashne, and S. C. Harrison (1988) *Science* **242**, 899–907.
  32. L. J. W. Shimon and S. C. Harrison (1993) *J. Mol. Biol.* **232**, 826–838.
  33. A. M. Friedman, T. O. Fischmann, T. A. Steitz (1995) *Science* **268**, 1721–1727.
  34. M. Lewis, G. Chang, N. C. Hortin, M. A. Kercher, H. C. Pace, M. A. Schumacher, R. G. Brennan, and P. Lu (1996) *Science* **271**, 1247–1254.
  35. C. A. E. M. Spronk, A. M. J. J. Bonvin, P. K. Radha, G. Melacini, R. Boelens, and R. Kaptein (1999) *Structure* **7**, 1483–1492.
  36. C. E. Bell and M. Lewis (2000) *Nat. Struct. Biol.* **7**, 209–214.
  37. R. Kaptein, E. R. P. Zuiderweg, R. M. Scheek, R. Boelens, and W. F. van Gunsteren (1985) *J. Mol. Biol.* **182**, 179–182.
  38. R. Boelens, R. M. Scheek, J. H. van Boons, and R. Kaptein (1987) *J. Mol. Biol.* **193**, 213–216.
  39. Z. Otwinowski, R. W. Schevitz, R.-G. Zhang, C. L. Lawson, A. Joachimiak, R. Q. Marmorstein, B. F. Luisi, and P. B. Sigler (1988) *Nature* **335**, 321–329.

40. C. L. Lawson and J. Carey (1993) *Nature* **366**, 178–182.
41. S. C. Schultz, G. C. Shields, and T. A. Steitz (1991) *Science* **253**, 1001–1007.
42. D. Kostrewa, J. Granzin, C. Koch, H.-W. Choe, S. Raghunathan, W. Wolf, J. Labahn, R. Kahmann, and W. Saenger (1991) *Nature* **349**, 178–180.
43. J.-A. Feng, R. C. Johnson, and R. E. Dickerson (1994) *Science* **263**, 348–355.
44. K. Wilson, L. M. Shewchuk, R. G. Brennan, A. J. Otsuka, and B. W. Matthews (1992) *Proc. Natl. Acad. Sci. USA* **89**, 9257–9261.
45. S. Rhee, R. G. Martin, J. L. Rosner, and D. R. Davies (1998) *Proc. Natl. Acad. Sci. U.S.A.* **95**, 10413–10418.
46. H. J. Kwoon, M. H. J. Bennik, B. Demple, and T. Ellenberger (2000) *Nature Struct. Biol.* **7**, 424–430.
47. C. Wolberger, A. K. Vershon, B. Liu, A. D. Johnson, and C. O. Pabo (1991) *Cell* **67**, 517–528.
48. C. L. Phillips, A. K. Vershon, A. D. Johnson, and F. W. Dahlquist (1991) *Genes Dev.* **5**, 764–772.
49. C. R. Kissinger, B. Liu, E. Martin-Bianco, T. B. Kornberg, and C. O. Pabo (1990) *Cell* **63**, 579–590.
50. Y. Q. Qian, M. Billeter, G. Otting, M. Müller, W. J. Gehring, and K. Wüthrich (1989) *Cell* **59**, 573–580.
51. G. Otting, Y. Q. Qian, M. Billeter, M. Müller, M. Affolter, W. J. Gehring, and K. Wüthrich (1990) *EMBO J.* 3085–3092.
52. M. Billeter, Y. Q. Qian, G. Otting, M. Müller, W. Gehring, and K. Wüthrich (1993) *J. Mol. Biol.* **234**, 1084–1097.
53. M. Müller, M. Affolter, M. W. Leupin, G. Otting, K. Wüthrich, and W. J. Gehring (1988) *EMBO J.* 4299–4304.
54. W. J. Gehring, M. Müller, M. Affolter, A. Percival-Smith, M. Billeter, Y. Q. Qian, G. Otting, and K. Wüthrich (1990) *Trends Genet.* **6**, 323–329.
55. W. J. Gehring, Y. Q. Qian, M. Billeter, K. Furukubo-Tokunaga, A. F. Schier, D. Resendez-Perez, M. Affolter, G. Otting, and K. Wüthrich (1994) *Cell* **78**, 211–223.
56. Y. Mo, B. Vaessen, K. Johnston, and R. Marmorstein (2000) *Nat. Struct. Biol.* **7**, 292–297.
57. Y. Mo, B. Vaessen, K. Johnston, and R. Marmorstein (1998) *Mol. Cell* **2**, 201–212.
58. R. Kodandapani, F. Pio, C.-Z. Ni, G. Picialli, M. Klemsz, S. McKercher, R. A. Maki, and K. R. Ely (1996) *Nature* **380**, 456–460.
59. C. R. Escalante, J. Yie, D. Thanos, and A. K. Aggarwal (1998) *Nature* **391**, 103–106.
60. P. E. Wright (1994) *Curr. Opin. Struct. Biol.* **4**, 22–27.
61. C. J. Harrison, A. A. Bohm, H. C. M. Nelson (1994) *Science* **263**, 224–227.
62. V. Ramakrishnan, J. T. Finch, V. Graziano, P. L. Lee, R. M. Sweet (1993) *Nature* **362**, 219–223.
63. K. L. Clark, E. D. Halay, E. Lai, and S. K. Burley (1993) *Nature* **364**, 412–420.
64. K. Ogata, H. Hojo, S. Aimoto, T. Nakai, H. Nakamura, A. Sarai, S. Ishii, and Y. Nishimura (1992) *Proc. Natl. Acad. Sci. USA* **89**, 6428–6432.
65. K. Ogata, S. Morikawa, H. Nakamura, H. Hojo, S. Yoshimura, R. Zhang, S. Aimoto, Y. Ametani, Z. Hirata, A. Sarai, S. Ishii, and Y. Nishimura (1995) *Nat. Struct. Biol.* **2**, 309–320.
66. N. Assa-Munt, R. J. Mortishire-Smith, R. Aurora, W. Herr, and P. E. Wright (1993) *Cell* **73**, 193–205.
67. N. Dekker, M. Cox, R. Boelens, C. P. Verrijetzer, P. C. van der Vliet, and Kaptein (1993) *Nature* **362**, 852–855.
68. W. Xu, M. A. Rould, S. Jun, C. Desplan, and C. O. Pabo (1995) *Cell* **80**, 639–650.
69. H. E. Xu, M. A. Rould, W. Xu, J. A. Epstein, R. L. Maas, and C. O. Pabo (1999) *Genes Dev.*

- 13, 1263–1275.
70. T. E. Ellenberger, C. J. Brandl, K. Struhl, and S. C. Harrison (1992) *Cell* **71**, 1223–1237.
  71. Y. Fujii, T. Shimizu, T. Toda, M. Yanagida, and T. Hakoshima (2000) *Nature Struct. Biol.* **7**, 889–893.
  72. P. König and T. J. Richmond (1993) *J. Mol. Biol.* **233**, 139–154.
  73. W. Keller, P. König, and T. J. Richmond (1995) *J. Mol. Biol.* **254**, 657–667.
  74. P. C. M. Ma, M. A. Rould, H. Weintraub, and C. O. Pabo (1994) *Cell*, **77**, 451–459.
  75. T. Ellenberger, D. Fass, M. Arnaud, and S. C. Harrison (1994) *Genes & Development* **8**, 970–980.
  76. A. R. Ferr<sup>o</sup>-D'Amaré, G. C. Prendergast, E. B. Ziff, and S. K. Burley (1993) *Nature* **363**, 38–45.
  77. P. Brownlie, T. A. Ceska, M. Lamers, C. Romier, G. Stier, H. Teo, and D. Suck (1997) *Structure* **5**, 509–520.
  78. M. Schmiedeskamp and R. E. Klevit (1994) *Curr. Opin. Struct. Biol.* **4**, 28–35.
  79. J. P. MacKay and M. Crossley (1998) *Trends Biochem. Sci.* **23**, 1–4.
  80. B. F. Luisi (1992) *Nature* **356**, 379–380.
  81. N. P. Pavletich and C. O. Pabo (1991) *Science* **252**, 809–817.
  82. J. M. Berg (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 11109–11110.
  83. G. H. Jacobs (1992), *EMBO J.* **11**, 4507–4517.
  84. R. W. Kriwacki, S. C. Schultz, T. A. Steitz, and J. P. Caradonna (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 9759–9763.
  85. J. R. Desjarlasi and J. M. Berg (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 2256–2260.
  86. R. C. Hoffman, S. J. Horvath, and R. E. Klevit (1993) *Protein Sci.* **2**, 951–965.
  87. J. Kuwahara, A. Yonezawa, M. Futamura, and Y. Sugiura (1993) *Biochemistry* **32**, 5994–6001.
  88. R. E. Klevit, J. R. Herriott, and S. J. Horvath (1990) *Proteins* **7**, 215–226.
  89. D. Neuhaus, Y. Nakaseko, J. W. Schwabe, A. Klug (1992) *J. Mol. Biol.* **228**, 637–651.
  90. M. Elrod-Erickson, M. A. Rould, L. Nekludova, and C. O. Pabo (1996) *Structure* **4**, 1171–1180.
  91. J. G. Omichinski, G. M. Clore, M. Robien, K. Sakaguchi, E. Appella, and A. M. Gronenborn (1992) *Biochemistry* **31**, 3907–3917.
  92. N. P. Pavletich and C. O. Pabo (1993) *Science* **261**, 1701–1707.
  93. M. S. Lee, G. P. Gippert, K. V. Soman, D. A. Case, and P. E. Wright (1989) *Science* **245**, 635–637.
  94. L. Fairall, J. W. R. Schwabe, L. Chapman, J. T. Finch, and D. Rhodes (1993) *Nature* **366**, 483–487.
  95. T. Härd, E. Kellenbach, R. Boelens, B. A. Maler, K. Dahlman, L. P. Fredman, J. Carlstedt-Duke, K. R. Yamamoto, J. A. Gustafsson, and R. Kaptein (1990) *Science* **249**, 157–160.
  96. B. F. Luisi, W. Xu, Z. Otwinowski, L. P. Freedman, and K. R. Yamamoto (1991) *Nature*, **352**, 497–505.
  97. H. Baumann, K. Paulsen, H. Kovács, H. Berglund, A. P. H. Wright, J.-A. Gustafsson, and T. Härd (1993) *Biochemistry* **32**, 13463–13471.
  98. J. W. R. Schwabe, D. Neuhaus, and D. Rhodes (1990) *Nature* **348**, 458–461.
  99. J. W. R. Schwabe, L. Chapman, J. T. Finch, and D. Rhodes (1993) *Cell* **75**, 567–578.
  100. R. M. A. Knegt, M. Kathaira, J. G. Schilthuis, A. M. J. J. Bonvin, R. Boelens, D. Eib P. T. Saag, and Kaptein (1993) *J Biomolec NMR* **3**, 1–17.
  101. M. S. Lee, S. A. Kliewer, J. Provencal, P. E. Wright, and R. M. Evans (1993) *Science* **260**,



1117–1121.

102. L. P. Freeman and B. F. Luisi (1993) *J. Cell. Biochem.* **51**, 140–150
103. J. G. Omichinski, M. G. Clore, O. Schaad, G. Felsenfeld, C. Trainor, E. Appella, S. J. Stahl, and A. M. Gronenborn (1993) *Science* **261**, 438–446.
104. R. Marmorstein, M. Carey, M. Ptashne, and S. C. Harrison (1992) *Nature* **356**, 408–414.
105. P. J. Kraulis, A. R. C. Raine, P. L. Gadhavi, and E. D. Laue (1992) *Nature* **356**, 448–450.
106. R. Marmorstein and S. C. Harrison (1994) *Genes Dev.* **8**, 2504–2512.
107. K. Swaminathan, P. Flynn, R. J. Reece, and R. Marmorstein (1997) *Nat. Struct. Biol.* **4**, 751–759.
108. Y. Cho, S. Gorina, P. D. Jeffrey, and N. P. Pavletich (1994) *Science* **265**, 346–355.
109. X. Qian, C. Jeon, H. Yoon, K. Agarwal, and M. A. Weiss (1993) *Nature* **365**, 277–279.
110. X. Qian, S. N. Gozani, H. Yoon, C. J. Jeon, K. Agarwal, and M. A. Weiss (1993) *Biochemistry* **32**, 9944–9959.
111. P. Cramer, D. A. Bushnell, J. Fu, A. L. Gnat, B. Maier-Davis, N. E. Thompson, R. R. Burgess, A. M. Edwards, P. R. David, R. D. Kornberg (1999) *Science* **288**, 640–649.
112. B. Wang, D. N. Jones, B. P. Kaine, M. A. Weiss (1998) *Structure* **6**, 555–569.
113. C. D. Mackereth, C. H. Arrowsmith, A. M. Edwards, and L. P. McIntosh (2000) *Proc. Natl. Acad. Sci. USA* **97**, 6316–6321.
114. J. N. Breg, J. H. J. van Opheusden, M. J. M. Burgering, R. Boelens, and R. Kaptein (1990) *Nature* **346**, 586–589.
115. B. E. Raumann, M. A. Rould, C. O. Pabo, and R. T. Sauer (1994) *Nature* **367**, 754–757.
116. B. E. Raumann, B. E. Brown, and R. T. Sauer (1994) *Curr. Opin. Struct. Biol.* **4**, 36–43.
117. M. J. Burgering, R. Boelens, D. E. Gilbert, J. N. Breg, K. L. Knight, R. T. Sauer, R. Kaptein (1994) *Biochemistry* **33**, 15036–15045.
118. C. D. Waldburger and R. T. Sauer (1995) *Biochemistry* **34**, 13109–13116.
119. J. B. Rafferty, W. S. Somers, I. Saint-Girons, and S. E. V. Phillips (1989) *Nature* **341**, 705–710.
120. W. S. Somers and S. E. V. Phillips (1992) *Nature* **359**, 387–393.
121. Z. S. Juo, T. K. Chiu, P. M. Leibermann, I. Baikalov, A. J. Berk, and R. E. Dickerson (1996) *J. Mol. Biol.* **261**, 239–254.
122. Y. Kim, J. H. Geiger, S. Hahn, and P. B. Sigler (1993) *Nature* **365**, 512–529.
123. J. L. Kim, D. B. Nikolov, and S. K. Burley (1993) *Nature* **365**, 520–527.
124. J. L. Kim and S. K. Burley (1994) *Nat. Struct. Biol.* **1**, 638–653.
125. H. Robinson, Y.-G. Gao, B. S. McCrary, S. P. Edmondson, J. W. Shriver, and A. H.-J. Wang (1998) *Nature* **392**, 202–205.
126. J. Bravo, Z. Li, N. A. Speck, and A. J. Warren (2001) *Nature Struct. Biol.* **8**, 371–378.
127. L. Chen, M. Glover, P. G. Hogan, Rao, and S. C. Harrison (1998) *Nature* **392**, 42–48.
128. P. Cramer, C.J. Larson, G. L. Verdine, and C. W. Müller (1997) *EMBO J.* **16**, 7078–7090.
129. Y. Cho, S. Gorina, P. D. Jeffrey, and N. P. Pavletic (1994) *Science* **265**, 346–355.
130. S. Becker, B. Groner, and C. W. Müller (1998) *Nature* **394**, 145–151.
131. C.W. Müller and B.G. Herrmann (1997) *Nature* **389**, 884–888.
132. D. F. Senear and R. Batey (1991) *Biochemistry* **30**, 6677–6688.
133. B. M. Brown and R. T. sauer (1993) *Biochemistry* **32**, 1354–1363.
134. H. Weintraub, R. Davis, D. Lockshon, and A. Lassar (1990) *Proc. Natl. Acad. Sci. USA* **87**, 5623–5627.
135. A. G. E. Künne, M. Sieber, D. Meierhans, and R. K. Allemann (1998) *Biochemistry* **37**, 4217–4223.

136. D. Meierhans, C. el Ariss, M. Neuenschwander, M. Sieber, J. F. Stackhouse, and R. K. Allemann (1995) *Biochemistry* **34**, 11026–11036.
137. A. G. E. Künne, D. Meierhans, and R. K. Allemann (1996) *FEBS Lett.* **391**, 79–83.
138. M. Sha, A. R. Ferré-D'Amaré, S. K. Burley, and D. J. Goss (1995) *J. Biol. Chem.* **270**, 19325–19329.
139. A. G. E. Künne and R. K. Allemann (1997) *Biochemistry* **36**, 1085–1091.
140. M. Sieber and R. K. Allemann (1998) *Nucleic Acids Res.* **26**, 1408–1413.
141. H. Berglund, H. Baumann, S. Knapp, R. Ladenstein, and T. Härd (1995) *J. Am. Chem. Soc.* **117**, 12883–12884.
142. M. R. Gryk and O. Jardetzky (1996) *J. Mol. Biol.* **255**, 204–214.
143. D. N. Arvidson, J. Pfau, J. K. Hatt, M. Shapiro, F. S. Pecoraro, and P. Youderian (1993) *J. Biol. Chem.* **268**, 4362–4639.
144. B. Cuenoud and A. Schepartz (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 1154–1159.
145. C. R. Robinson and R. T. Sauer (1996) *Biochemistry* **35**, 109–116.
146. C. M. Falcon, L. Swint-Kruse, and K. S. Matthews (1997) *J. Biol. Chem.* **272**, 26818–26821.
147. C. Dezan, D. Meierhans, A. G. E. Kuenne, and R. K. Allemann (1999) *Biol. Chem.* **380**, 705–710.
148. S. Wagner and M. R. Green (1993) *Science* **262**, 395–399.
149. A. M. Baranger, C. R. Palmer, M. K. Hamm, H. A. Giebler, A. Brauweiler, J. K. Nyborg, and A. Schepartz (1995) *Nature* **376**, 606–608.
150. G. Perini, S. Wagner, and M. R. Green (1995) *Nature* **376**, 602–605.
151. C. R. Palmer, L. D. Gagnas, and A. Schepartz (1997) *Biochemistry* **36**, 15349–15355.
152. A. P. Armstrong, A. A. Franklin, M. N. Uittenbogaard, H. A. Giebler, and J. K. Nyborg (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 7303–7307.
153. J. M. Cox, L. S. Sloan, and A. Schepartz (1995) *Chem. Biol.* **2**, 819–826.
154. R. Grassmann, C. Dengler, I. Müller-Fleckenstein, B. Fleckenstein, K. McGuire, M.-C. Dokhelar, J. G. Sodroski, and W. A. Haseltine (1989) *Proc. Natl. Acad. Sci. USA* **86**, 3351–3355.
155. T. Oka, H. Sonobe, J. Iwata, I. Kubonishi, H. Satoh, M. Takata, Y. Tanaka, M. Tateno, H. Tozawa, S. Mori, T. Yoshiki, and Y. Ohtsuki (1992) *J. Virol.* **66**, 6686–6694.
156. G. Natoli, M. L. Avantaggiati, P. Chirillo, A. Costanzo, M. Artini, C. Balsano, and M. Levrero (1994) *Mol. Cell. Biol.* **14**, 989–998.
157. C. Balsano, O. Billet, M. Bennoun, C. Cavard, A. Zider, G. Grimber, G. Natoli, P. Briand, and M. Levrero (1994) *J. Hepatol.* **21**, 103–109.
158. W. McGinnis and R. Krumlauf (1992) *Cell* **68**, 283–302.
159. R. S. Mann (1995) *Bioassays* **17**, 855–863.
160. R. S. Mann and S.-K. Chan (1996) *Trends Genet.* **12**, 258–262.
161. J. M. Passner, H. D. Ryoo, L. Shen, R. S. Mann, and A. K. Aggarwal (1999) *Nature* **397**, 714–719.
162. A. D. Johnson (1995) *Curr. Opin. Genet. Dev.* **5**, 552–558.
163. S. Tan and T. J. Richmond (1998) *Nature* **391**, 660–666.
164. M. Hassler and T. J. Richmond (2001) *EMBO J.* **20**, 3018–3028.
165. T. Li, M. R. Stark, A. D. Johnson, and C. Wolberger (1995) *Science* **270**, 262–269.
166. D. Chasman, K. Cepek, P. A. Sharp, and C. O. Pabo (1999) *Genes Dev.* **13**, 2650–2657.

### **Suggestions for Further Reading**

167. C. O. Pabo and R. T. Sauer (1992) *Transcription factors: structural families and principles of*

DNA recognition. *Annu. Rev. Biochem.* **61**, 1053–1095.

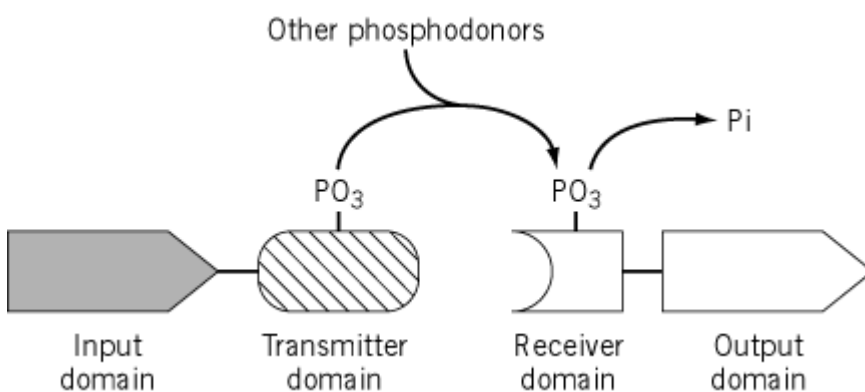
168. B. Luisi (1995) "Frontiers in molecular biology". In *DNA-Protein: Structural Interactions* (D. M. J. Lilley ed.,) IRL Press at Oxford University Press, Oxford, pp. 1–48.
169. D. S. Latchman (1998) *Eukaryotic Transcription Factors*, 3 ed., Academic Press, ISBN: 0124371779.

## Transducer Proteins

The transducer proteins are [proteins](#) that receive and help process information so that proper action by the cell ensues. This article emphasizes transducer proteins that bind ligands on the outside of **bacterial cells** and, consequently, cause changes within the cell to optimize cell survival in the particular environment. Two particular types of transducers are described, the sensor kinase of the two-component [signal transduction](#) system and the methylated receptors of [chemotaxis](#) in bacteria. Both have ectodomains responsible for interacting with environmental factors, which induce conformational changes in the receptor to regulate the cytoplasmic coupled kinase activity.

An important means by which bacteria sense changes in their environment is by a signal transduction strategy known as the two-component system. It consists of a membrane-spanning sensor kinase and a cognate cytoplasmic response regulator. A sensor [kinase](#) with two **domains**, an input domain usually outside the plasma membrane and a cytoplasmic transmitter domain, autophosphorylates its transmitter domain on a histidine residue at a rate that depends on the input domain's conformation. The response regulator usually also has two domains, a receiver domain and an output domain. Using the phosphorylated sensor kinase as a substrate, the receiver domain catalyzes the transfer of the phosphate to an aspartate residue within itself, so that inhibition is relieved on the output domain (Fig. 1). The output domain becomes altered, usually resulting in activation of the [transcription](#) of an appropriate [operon](#).

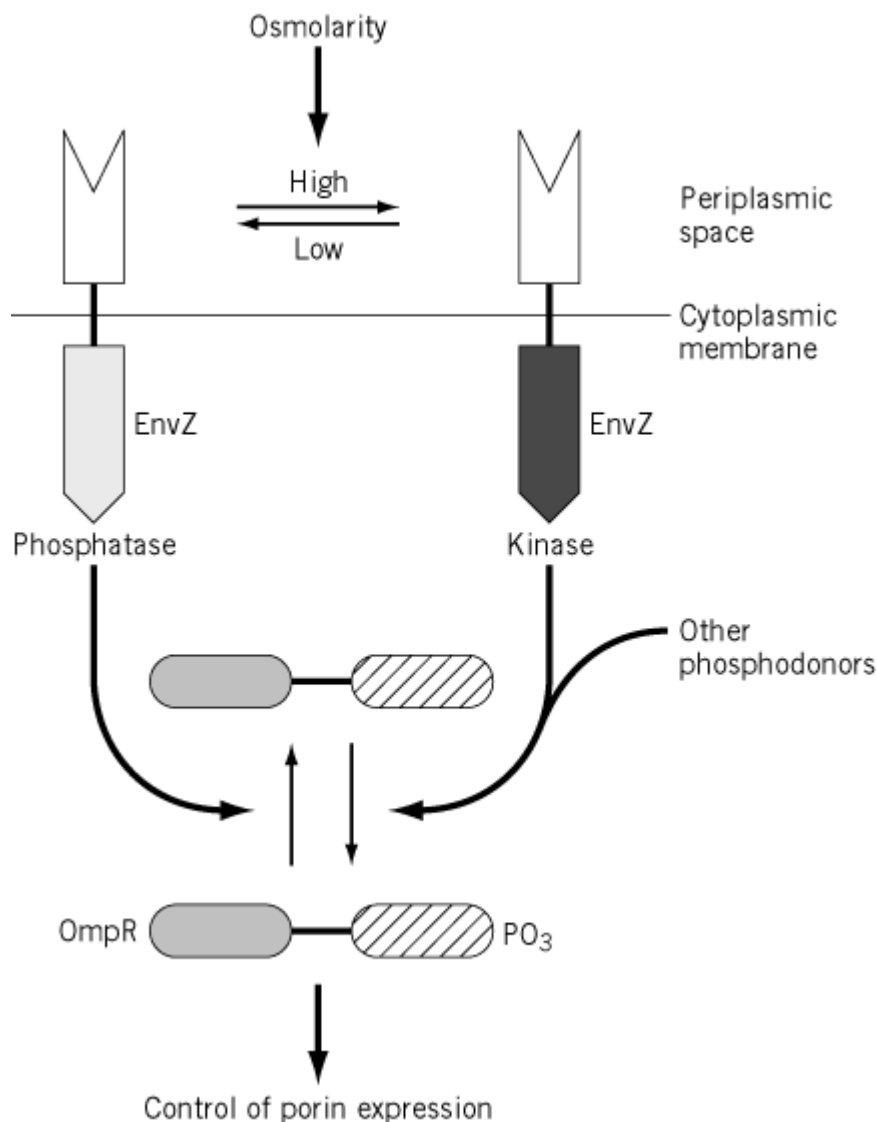
**Figure 1.** Prototypical two-component system in bacteria.



One particularly interesting instance of this is from the osmoregulation circuit in *Escherichia coli*. The sensor kinase EnvZ has an input domain in the **periplasmic** space beyond the membrane and a cytoplasmic output domain that phosphorylates OmpR, the controller of [porin](#) expression. In an

unknown way, high osmolarity makes EnvZ act as a kinase, leading to high concentrations of phosphorylated OmpR (OmpR-P), and low osmolarity makes EnvZ act as a phosphatase, leading to low concentrations of OmpR-P (1) (Fig. 2).

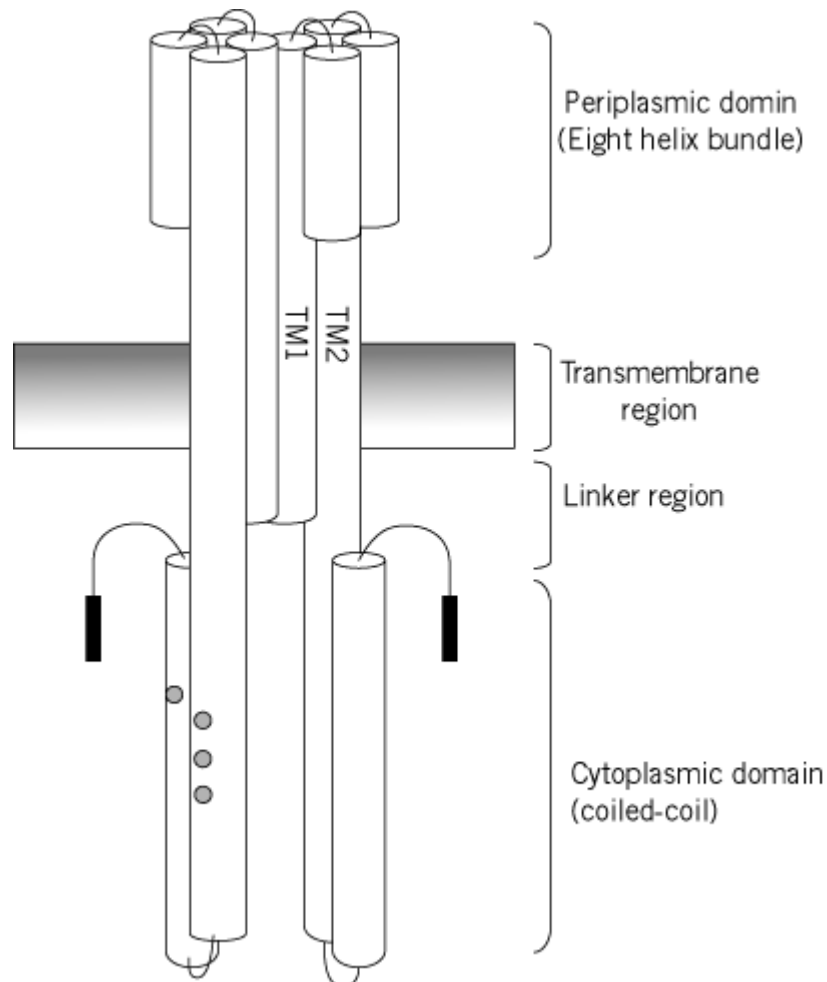
**Figure 2.** Sensor kinase of the two-component system that senses osmolarity in bacteria. Open arrows indicate control of the designated signaling step. EnvZ is transducer protein that senses the osmolarity of the external medium, with effects on its opposite activities as a kinase or a phosphatase of OmpR. In turn, the phosphorylation state of OmpR controls the expression of the gene for porin.



The primary means by which bacteria sense changes in concentration of chemotactic attractant or repellents is through chemoreceptors. Chemoreceptor proteins exist *in vivo* as stable dimer molecules that are entirely  $\alpha$ -helical throughout. The extracellular domain organization of the *E. coli* Tar dimer consists of an eight-helix bundle, and the cytoplasmic domain consists mostly of a four-helix coiled-coil (Fig. 3)(3). The chemoreceptors may bind attractants and repellents directly, like Tar binding aspartate, or indirectly, by interaction with ligand-binding proteins like Tar-binding maltose-binding protein, associated with maltose. Evidence indicates that Tar-binding aspartate produces a piston-like shift in transmembrane helix 2 toward the cytoplasm. This conformational change results in a reduction in CheA kinase activity, the kinase of the chemotaxis system. Reduction in CheA-P levels in turn reduces levels of CheY-P (see [Chemotaxis](#)), causing a change in direction of flagellar

rotation.

**Figure 3.** Dimeric structure of the chemoreceptors in bacteria. The cylinders represent  $\alpha$ -helical regions of each monomer. Attractants bind between the monomers causing a downward shift of transmembrane helix 2 of one of the monomers.



Chemotaxis transducers are unique, because they contain conserved, reversibly methylated glutamate residues in the cytoplasmic four-helix coiled-coil domain. These conserved glutamate residues are the substrates for the methyltransferase enzyme, CheR, which uses S-adenosyl methionine as a methyl donor (see [Methyltransferase](#)). Methyl glutamate residues are substrates for the methyl-esterase enzyme, CheB. CheB is a response regulator protein like CheY, which also uses CheA as a phospho-donor, and is activated by phosphorylation to demethylate the chemoreceptors. The methylation system serves as an adaptational mechanism to make the chemotactic receptors phasic (on for a short period) rather than tonic (always on when ligand is bound), a unique feature among two component systems. The methylation region is located between the signaling region, where the CheA kinase binds, and the cell exterior, where the chemoeffectors bind, and is thus strategically located to compensate for the movement caused by chemoeffector binding. These taxis transducer proteins are found among both the **eubacteria** and **archaea**, and thus probably antedate the separation of these two domains of life. In all known cases, these transducer proteins control a sensor CheA kinase. In *E. coli*, they operate as chemoreceptors, as described before. In two instances in the archeon *Halobacterium salinarium*, the transducer protein complexes with a sensory [rhodopsin](#), which senses light by using a retinal chromophore and interacts with the methylated transducer protein to control the CheA kinase (4).

Bibliography

1. J. S. Parkinson (1993) *Cell* **73**, 857–871.
2. J. B. Stock and M. G. Surette (1996) In *Escherichia coli and Salmonella, Cellular and Molecular Biology*, 2 ed. (F. C. Neidhardt, ed.), ASM Press, Washington, DC, Vol. I, pp. 1103–1129.
3. J.J. Falke (2001) *Trends Biochem. Sci.* **26**, 257–265.
4. W. D. Hoff, K.H. Jung, and J. L. Spudich (1997) *Annu. Rev. Biophys. Biomol. Struct.* **26**, 221–256.

### Suggestions for Further Reading

5. J.J. Falke and S.H. Kim (2000) Structure of a conserved receptor domain that regulates kinase activity: The cytoplasmic domain of bacterial taxis receptors. *Curr. Opin. Struct. Biol.* **10**, 462–469.
6. T. W. Grebe and J. B. Stock (1999) The histidine protein kinase superfamily. *Adv. Microb. Physiol.* **41**, 139–227.

## Transducin

Transducin, also known as  $G_t$  by the standard [GTP-binding protein](#) nomenclature, is the major [heterotrimeric G protein](#) of photoreceptor cells of the vertebrate retina, comprised of  $\alpha$ ,  $\beta$ , and  $\gamma$  subunits. Activation of  $G_t$  is initiated by the photoactivated form of [rhodopsin](#) (known as “bleached,” meta II), and it in turn activates the effector protein cyclic GMP phosphodiesterase (PDE). By hydrolyzing [cyclic GMP](#) in response to light, this signaling pathway leads to the closure of cyclic GMP-activated cation channels in photoreceptor cells. The resultant hyperpolarization of the photoreceptor cell is conveyed synaptically to organizational interneurons in the retina, which ultimately send the visual information to the brain.

Transducin's mechanism of action is typical for a G protein. Bleached rhodopsin catalyzes the release of GDP from  $G_{\alpha_t}$ , and the binding of GTP in its place, in an extremely rapid and efficient reaction. Several hundred  $G_t$  molecules can be activated in 1 s, and a single bleached rhodopsin molecule can catalyze the activation of 2000  $G_t$  molecules in a mammalian photoreceptor cell (1). Without rhodopsin, the rate of GDP/GTP exchange by  $G_t$  is extremely slow; noncatalyzed stoichiometric exchange has not been observed *in vitro* (see [Guanine Nucleotide Exchange Factors](#)). Such slow spontaneous activation is needed to minimize background signaling in the dark, which is in turn needed to allow photoreceptor signaling to operate over a huge range of light intensities.

Activated  $G_{\alpha_t}$ -GTP binds and sequesters the two inhibitory PDE  $\gamma$  subunits from the inactive  $\alpha\beta\gamma_2$  PDE heterotetramer, thus activating the catalytic PDE  $\alpha\beta$  subunits. Although  $G_{\beta\gamma}$  is presumably released upon  $G_t$  activation, no regulatory role for  $G_{\beta\gamma}$  in photoreceptor cells has been clearly established. Signaling is terminated by hydrolysis of  $G_{\alpha_t}$ -bound GTP and the consequent release of PDE  $\gamma$ . In photoreceptor cells, GTP hydrolysis and its consequent deactivation are much faster than is observed with isolated  $G_t$  ( $\sim 1s^{-1}$  rather than  $0.02s^{-1}$ ). Deactivation is accelerated in photoreceptors by RGS9, a GTPase-activating protein (GAP) whose activity is potentiated by PDE $\gamma$

subunits (2) (see [RGS Proteins](#)).

There are two  $G_{\alpha_t}$  polypeptide chains,  $G_{\alpha_{t1}}$  and  $G_{\alpha_{t2}}$ , encoded by two closely related **genes** in the two classes of photoreceptor cells. Both are members of the  $G_{\alpha_i}$  family and are typically sensitive to [pertussis toxin](#). In rod cells, which are responsible for vision in dim light,  $G_t$  is composed of  $G_{\alpha_{t1}}$ ,  $\beta_1$ , and  $\gamma_1$  subunits. Its great abundance, ~15% of total protein in the photoreceptor membranes of rod outer segments (3), made rod transducin one of the two experimental prototypes in the elucidation of G-protein signaling mechanisms. Transducin in cone photoreceptor cells (there is one cone cell each for red-, green-, and blue-sensitive opsins) contains  $G_{\alpha_{t2}}$ . These two  $G_{\alpha_t}$  polypeptide chains are also found at low abundance in a few other places in the brain.

Transducin is unusual among G proteins in several respects. First, it is relatively more water-soluble than other G proteins, despite the addition of a **farnesyl group** to  $G_{\beta_1}$  and a **myristoyl group** to  $G_{\alpha_t}$  (which has no **palmitoyl groups**). It can be solubilized from photoreceptor membranes by high ionic strength in the absence of [detergent](#), and  $G_{\alpha_t}$  can be solubilized quantitatively by dissociating it from Gbg after activation with a nonhydrolyzable GTP analog. Its mode of regulating its effector by chelation of inhibitory subunits is also unique. Last, basal GDP/GTP exchange on  $G_{\alpha_t}$  is at least 10-fold slower than that on other G proteins.

#### Bibliography

1. P. A. Liebman, K. R. Parker, and E. A. Dratz (1987) *Annu. Rev. Physiol.* **49**, 765–791.
2. W. He, C. W. Cowan, and T. G. Wensel (1998) *Neuron* **20**, 95–102.
3. H. E. Hamm and M. D. Bownds (1986) *Biochemistry* **25**, 4512–4523.

#### Suggestions for Further Reading

4. P. A. Liebman, K. R. Parker, and E. A. Dratz (1987) The molecular mechanism of visual excitation and its relation to the structure and composition of the rod outer segment. *Annu. Rev. Physiol.* **49**, 765–791.
5. L. Stryer (1986) Cyclic GMP cascade of vision. *Annu. Rev. Neurosci.* **9**, 87–119.

## Transfection

Transfection is the introduction of foreign **DNA** into the **nuclei** of **eukaryotic** cells, a fundamental technology in molecular and cellular biology. There are many reasons for transfecting cells (Table 1). Most often, cells are transfected to characterize the function and control of **genes**. More recently, economic (see [Transgenic Technology](#)) or therapeutic (human gene therapy) purposes have become increasingly important. After transfection, the foreign DNA, which is also referred to as the *transgene*, is usually neither integrated into cellular [chromosomes](#) nor replicating as an **episome**. In mitotically active cells, the transgene persists for only a few days, with a peak of expression observed within the first 24–48 h after exposure to DNA (transient transfection). Stable transfection results only if the transgene is integrated into the cellular [genome](#). This is a random event, facilitated by chromosomal strand breaks and involving cellular [recombinases](#) and [DNA repair](#) pathways in an as yet poorly defined manner. The frequency of stable integration after transient transfection is dependent on the target cell population and on the method and protocol used; it usually occurs in  $10^{-}$

$2$  to  $10^{-6}$  of transiently transfected cells. Therefore, dominant selectable marker genes encoding drug resistance or cell surface or cytoplasmic markers are used for the identification of stable clones. Regulatory elements of stably inserted foreign DNA are subject to environmental influences from neighboring chromosomal sequences. This can result in either support or extinction (silencing) of transgene expression, leading to a broad clonal variability, depending on the integration site. Semistable transfection can be accomplished by introducing episomally replicating **plasmids** that are inherited by random distribution in mitosis. Usually, foreign DNA is more autonomously expressed in the episomal form, but can easily be lost if not conferring a selective advantage.

**Table 1. Reasons for Transfecting Eukaryotic Cells**

---

|   |
|---|
| Study <i>cis</i> -active elements of DNA (control regions of genes)               |
| Study coding sequences of DNA (function of genes)                                 |
| Overexpress genes for protein purification (including industrial production)      |
| Alter cellular genotype or phenotype for fate-mapping in complex organisms        |
| Alter cellular genotype or phenotype for diagnosis or therapy in humans           |
| Alter characteristics of cells or organisms for scientific or economical purposes |

---

Starting from the original description of calcium phosphate precipitation as the first transfection method in 1973 (1), a great variety of methods and protocols have been developed to date that, in principle, allow transfection of any cell. However, with the increasing sophistication of the scientific goals, efficacy is no longer the only parameter to be considered (Table 2). Equally important criteria for choosing a method are the impact of the transfection procedure on cellular **stress responses** and viability, plus the fate of DNA and transfection vehicles inside the cell. This includes the potential degradation of DNA in **lysosomes**, the copy number of transgenes in the nucleus, the type and loci of transgene integration, and the potential of rearranging the cellular genome during or after transfection.

**Table 2. Criteria for Evaluating Transfection Methods**

---

|  |
|--|
| Stress and toxicity of method (side effects for cellular genotype or phenotype)                |
| General applicability (cell numbers and amount of DNA required, accessibility of target cells) |
| Selectivity (possibility of targeted delivery when approaching complex cell systems or organs) |
| Potential of lysosomal degradation (might generate unwanted side products)                     |
| Duration of transfection (transient, stable, semistable)                                       |
| Quality of stable transfection (nuclear copy number and rearrangement potential)               |
| Efficacy (percentage of transfected cells, level of gene expression)                           |

---



These parameters are well understood in the natural ways of introducing genetic information into cells, such as **virus infection** pathways, some of which were exploited for developing very efficient, yet often laborious, viral transfection systems. The description of viral methods for introducing DNA into eukaryotic cells is covered elsewhere (see [Adenovirus](#), [Baculovirus](#), [Polyomavirus](#), [Retroviruses](#), [Rous Sarcoma Virus \(RSV\)](#), and [Vaccinia Virus](#) ). As is evident from these articles, viruses have evolved using diverse strategies to overcome the most important hurdles to introducing foreign DNA successfully into cells. These are the cellular membrane barrier, lysosomal degradation of endocytosed nucleic acids, nuclear transport of DNA, and, most importantly, replication and persistence of DNA in episomal moieties or integrated into the cellular genome. Some properties of viral replication pathways have been exploited in the design of physicochemical transfection methods. These have the general advantage of using biochemically defined material for transfection of cells. In many physicochemical strategies, however, the quality parameters discussed above and listed in Table 2 are still not entirely investigated. Despite this lack of basic knowledge, one encounters an ever increasing diversity of protocols and recommendations, whose variety is due to the use of different types of transfection materials or instruments and varying DNA, targets (cells or even organs), culture conditions, experimental scales and scientific goals.

This article provides an overview of the physicochemical techniques most commonly used for transfecting cells. Some transfection methods rely on physical forces to enhance the uptake of DNA into cells; these include injection, biolistics, and electroporation . Chemical and biological carrier-mediated transfection methods include calcium phosphate precipitation, complexation with DEAE-dextran, polybrene or activated dendrimers, lipofection , and, finally, receptor-mediated methods. Description of protocols is beyond the scope of this *Encyclopedia*. For this, the reader should refer to the literature cited and, where possible, to experienced scientists.

## 1. Transfection by Feeding

The simplest way of offering DNA to cells is feeding. In hematopoietic cells, a lineage-restricted cellular **receptor** for nucleic acids has been described (2); this is a heparin-binding [integrin](#). Similar receptors that trigger uptake of nucleic acids by [endocytosis](#) might also exist in other tissues. However, the biological significance of DNA receptor-mediated endocytosis of nucleic acids remains to be investigated. The endosome eventually fuses with lysosomes, after which the DNA is subjected to an acidic milieu and cellular **nucleases**. A surprising result has been obtained by feeding mice with not more than 50 µg of bacteriophage M13 DNA: fragments of M13 of size of up to 1.6 kbp have been detected by **PCR** in blood cells after feeding (3). This indicates that feeding can be quite successful for transferring small pieces of nucleic acids, such as oligonucleotides. However, it is—fortunately—not applicable for transfecting cells with genetically defined, biologically active genes.

## 2. Injection of DNA in Organs

To bypass the membrane barrier, physical methods have been developed. The simplest procedure is injection of DNA. Discovered surprisingly late, injection of DNA into postmitotic striated muscle tissue results in persistence and expression over prolonged periods of time (weeks to months) (4). In contrast to [microinjection](#) employed in the process of developing transgenic animals, injection does not directly target the nuclei. Because no carrier is used for the DNA, this implies that uptake and nuclear transport of foreign DNA occurs solely by cellular factors, thus representing a natural pathway. The preferred application of injecting DNA into tissues is DNA vaccination, which has turned out to be a most potent, versatile, and cost-effective way of immunization that is applicable to a large number of diseases, including viral and bacterial infections and possibly also cancer (4). More recently, instillation of DNA in hypertonic fluids has led to successful transfection of hepatocytes, indicating that hypertonic exposure to naked DNA after intravascular delivery might be applicable to other parenchymal organs (5).

### 3. Biolistics

More elaborate than injection using syringes is the method of biolistics. The principle of this method is to coat heavy-metal (preferably gold) particles (size about 1  $\mu\text{m}$ ) with DNA which thereafter are accelerated in a high voltage electric discharge or gas-pressure-driven apparatus (the so-called gene gun). Biolistics, also known as particle bombardment, is sufficiently versatile, efficient and flexible to be applied successfully in a wide range of cell types and cell environments, importantly also including organ slices *in vitro* and many types of solid tissues or organs *in vivo* (6-9). The method requires a fairly small amount of DNA and, compared to other transfection procedures, a small number of cells. Transfection results in high nuclear copy numbers (usually  $>20$ ). When two or more different DNA are simultaneously coated on the metallic bead, almost 100% cotransfection is observed. Stable transduction is reported to occur at relatively high frequency (in the range of 1% of transiently transfected cells) and apparently results in integration of multiple arrays of transgenes. This may lead to inaccuracies when analyzing regulatory elements of transgenes and, moreover, might render the cellular genome unstable. Biolistics has found wide acceptance for applications in almost any eukaryotic cell and in cell environments not accessible by other transfection procedures. Unfortunately, the potential drawbacks have not been systematically addressed, although some enthusiasm has already been raised regarding therapeutic applications (preferentially for DNA vaccination) in humans.

### 4. Electroporation and Iontophoresis

Electroporation also employs physical forces to bypass the membrane barrier. Cells suspended at high densities ( $2 \times 10^6$ – $10^7 \text{mL}^{-1}$ ) in a buffer are subject to a short electric pulse lasting microseconds to a few milliseconds (10-12). The pulse is generated by discharging an electrical field of high voltage through the cell suspension kept in a cuvette. Alternatively, cells can be grown on microporous membranes for transfection *in situ*. Electroporation has also been proposed for improving delivery of drugs or nucleic acids by transdermal *iontophoresis* (13, 14), a technique using an electrical potential gradient for facilitating the movement of solute ions across membranes. Electroporation is highly stressful, with best transfection efficiencies obtained under conditions resulting in more than 50% mortality. Surviving cells are expected to reorganize the membrane immediately after being pulsed. Careful examination, however, reveals membrane instabilities even hours after transfection. Parameters influencing the transfection efficacy and mortality are those defining the pulse length (condenser capacity and voltage; electrode gap distance; volume, resistance, and temperature of the cell suspension), the amount of DNA and, generally, the fitness of cells. Further variations result from the use of machines generating different types of pulse waveforms (exponential decay, rectangular plateau, and radiofrequency) or pulse repetitions (single pulse vs double pulse). Thus, almost any cell can be transfected using electroporation, but the definition of a suitable protocol can be cumbersome. Fortunately, some conditions seem to be more widely applicable (15, 16). An advantage of electroporation is that stable integrations of transfected DNA often occur as single copies, rendering electroporation the method of choice for gene trap or gene knockout studies in developmental biology (see [Gene Targeting](#)). Not sufficiently analyzed is whether the procedure of electroporation can induce chromosomal rearrangements. Moreover, it can be envisaged that the stress exerted by electroporation interferes with some transient transfection studies where the results of the experiment are monitored within 24–48 h after transfection.

### 5. Calcium Phosphate Transfection

Calcium phosphate transfection was developed in 1973 to introduce Adenovirus DNA into mammalian cells (1). This actually represents the first description of a transfection method, made possible by carefully evaluating the parameters leading to a fine DNA–salt precipitate that can be applied to cells for endocytotic uptake. The method works well for most adherent cells *in vitro* and has the general advantage of being exceptionally cheap (17). However, the efficacy of the method is

heavily dependent on the pH of the media and calcium phosphate buffers, explaining the low reproducibility in the hands of researchers not precisely controlling this factor. Stable transformants occur at a frequency of  $\sim 10^{-5}$ , often containing multiple copies of transfected DNA. This might be a result of previous complexation, as in biolistics, which would imply that the complete DNA–salt precipitate is transported into the nucleus. Using a special buffer allowing complex formation on cells in culture dishes, stable transfection frequencies can be improved 10–100-fold (18). Calcium phosphate transfection is not applicable *in vivo* or to suspension cells *in vitro*.

## 6. DEAE–Dextran Transfection

DEAE–dextran transfection is another simple and cheap method based on chemical complexation of DNA, which works well in transient expression experiments and can also be applied to some suspension cells (19). The mechanisms involved in DNA uptake and nuclear delivery in this method are poorly understood. The transfected DNA is present as minichromosomes and somehow excluded from stable integration. Parameters that need to be optimized for each given cell type are the concentrations and ratio of DEAE–dextran to DNA, the duration of transfection, and the time course of expression of the transfected gene. A DMSO (dimethylsulfoxide) or glycerol shock of cells after exposure to DEAE–dextran/DNA complexes can improve their uptake by chemical permeabilization of the membrane.

## 7. Polybrene/DMSO-Assisted Transfection

In this method, DNA is bound to the amorphous polycation polybrene, which is adsorbed electrostatically by negatively charged residues of **glycoproteins** located on the cell surface. Uptake is accomplished by chemically disrupting the membrane with a short exposure (in the region of 5 min) of cells to DMSO (20, 21). Toxicity is low to moderate, once the amount of polybrene and, more importantly, the conditions for the DMSO shock have been adjusted. Polybrene might protect the DNA from nuclease attack, perhaps similar to the lysosomal neutralization and disruption pathway described for activated dendrimers (see discussion below). Inside the nucleus, the DNA presumably dissociates from polybrene to facilitate [transcription](#) and integration. The method works best in adherent cell layers where high transient, as well as stable, transfection rates can be achieved; in suspension cells, it is not so efficacious. Integration can occur both as multiple copies arranged in tandem links at single sites and as single copies at distinct sites.

## 8. Activated Dendrimer Technology

More recently, so called activated dendrimers (22, 23) have been proposed for complex-mediated transfection. In this technique, DNA is bound electrostatically to tiny spherical cationic polyamidoamine polymers (dendrimers). Activated dendrimers with higher transfection efficacy are produced by partially degrading dendrimers on exposure to heat and solvents before binding of DNA. Uptake is mediated by endocytosis. No chemical permeabilization is required. Degradation of DNA is presumably blocked by neutralizing the acidic pH of lysosomes with an excess of positive charges of the dendrimer. Lysosomal escape might also involve osmotic swelling and subsequent disruption of the organelle. The subsequent fate of the dendrimer is yet unknown, as is the mechanism for nuclear translocation of transfected DNA. This relatively nontoxic method is applicable with good reproducibility for transient and stable transfection of adherent and suspension cells *in vitro*. *In vivo* applications have not been reported thus far. The copy number and integrity of stably transfected DNA still need to be investigated.

## 9. Lipofection

A method offering Trojan Horses for transfection is lipofection. The DNA is coated in surrogate membranes known as [liposomes](#), which are then added to cells in order to fuse with the cellular membrane for releasing their content into the cytoplasm (24, 25). Fusion of liposomes with cellular

membranes is a rare event, however; uptake occurs preferentially by endocytosis. Liposomes are artificial spherical or cochleate-like vesicles consisting of lipid bilayers covering small volumes of aqueous solution. Many different neutral or cationic [lipids](#) are utilized by many different companies for the commercial formulation of liposomes. The rationale for using cationic lipids is an improved interaction with negatively charged DNA and cellular membranes. Almost any formulation is reported to yield excellent transient transfection rates in selected target cells. In fact, different cell types need different lipid formulations to achieve good transfection rates. Therefore, several formulations need to be tested in a given experimental setting. Once the lipid of choice and the optimal DNA:lipid ratio is defined, lipofection works in a highly reproducible manner, for both transient and stable transfection. Of special interest is the reproducibly good stable transfection efficiency achieved with lipofection of large pieces of DNA, such as **yeast artificial chromosomes** ([26](#)). Lipofection is also employed for transfecting cells in complex organisms *in vivo*, including applications in gene therapy. However, not every formulation withstands inactivation by serum or membrane surfactants ([27](#)). Moreover, in poorly adjusted experimental conditions, lipofection can be hampered by a profound toxicity. Interesting prospects of lipofection are the introduction of membrane components (glycoproteins) of viruses in order to increase the potential for membrane fusion or lysosomal escape ([28](#)) and the combination with specific ligands for targeted delivery in specialized cell types (see the following paragraph). Furthermore, proteins can be codelivered with DNA, offering novel perspectives for nuclear and genomic targeting.

## 10. Receptor-Mediated Transfection

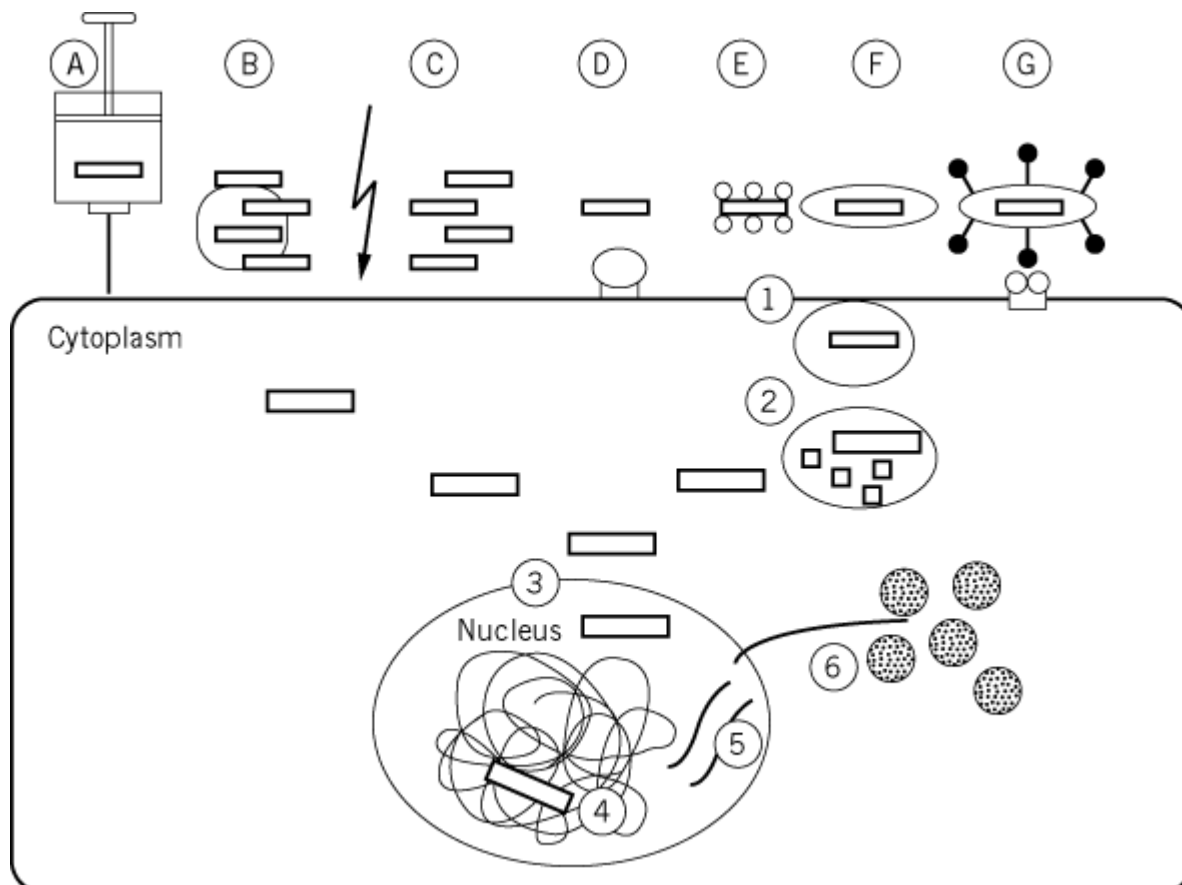
All methods described thus far are nonspecific with respect to the target cell population. Receptor-mediated transfection has been developed to overcome this limitation. In the archetypal form known as transferrin infection, DNA is electrostatically bound to polylysine that is coupled to [transferrin](#) ([29](#)). Uptake by endocytosis predominantly occurs in cells expressing the cognate transferrin receptor. High transfection efficiencies were reported when lysosomal degradation was blocked by treatment with **chloroquine**, or, even better, when the complexes had been coupled with pH-controlled membrane-disruption activities derived from [adenovirus](#) or [influenza virus](#) ([30](#)). This observation has led researchers to design multiple types of artificial DNA shuttles, consisting of a ligand for cell-type specific uptake (eg, for the asialoglycoprotein receptor exclusively expressed in liver epithelium ([30](#))), a DNA-binding moiety (polylysine, DNA-binding domains of [transcription factors](#), liposomes, possibly dendrimers in the future) and, optionally, the lysosome-disruption activity. The primary goal is to design biochemically defined pseudoviruses for applications in gene therapy. Criteria for the usefulness of these systems include (1) the specificity and efficacy of transduction, (2) the resistance to serum or other natural barriers such as membrane surfactants, and (3) the potential mutagenicity of randomly integrated DNA.

## 11. Summary and Outlook

In summary, enormous progress has been made in developing defined physicochemical procedures for transfecting cells *in vitro* and *in vivo* (summarized in Fig. [1](#) and Table [3](#)). To date, the challenge is no longer whether a cell can be transfected, but rather what the consequences of the transfection method are for the genotype and phenotype of the transfected cells. The primary choice to be made is the method, a decision that should be influenced by a basic understanding of its technological background and its possible consequences for the efficiency and quality of the experiment. Important parameters in this regard are the stress (for the cell, of course, not so much for the researcher) and toxicity of the method, the copy number of uptaken DNA and its fate inside the cell, and the consequences of the newly-introduced DNA for the cellular phenotype. The secondary but still critical choice is that of the protocol, which has profound impact on the outcome of the experiment and often will have to be adapted according to the individual setting under investigation.

**Figure 1.** Schematic overview of physicochemical transfection methods. The transgene is represented as a horizontal

bar. Injection (A); biolistics (B); electroporation (C); feeding—possibly involving a cellular receptor for nucleic acids (D); chemical complex-mediated methods including calcium phosphate transfection, DEAE/dextran-transfection, activated dendrimers, and polybrene/DMSO transfection, the latter requiring chemical permeabilization of the membrane (E); lipofection (F); receptor-mediated transfection (G). The first step in transfection is to bypass the membrane barrier, requiring either endocytosis, fusion with, or forced disruption of the membrane (1). Endocytosed nucleic acids can be subject to nuclease attack in lysosomes (2); lysosomal escape is a prerequisite for successful transfection with defined transgenes. After nuclear transport, which is mediated by yet unknown cellular cofactors, the transgene is present in episomal form (3). Integration in the cellular genome occurs randomly, and rarely (4). The transgene is transcribed to RNA from its episomal or integrated form (5). The cellular phenotype is altered according to the coding sequence and expression levels of the transgene (6).



**Table 3. Overview of Transfection Methods**

| Method    | Principle                           | Uptake               | Transient Transfection | Stable Transfection |
|-----------|-------------------------------------|----------------------|------------------------|---------------------|
| Feeding   | Exposure to DNA in aqueous solution | Endocytosis          | Yes (degraded DNA)     | Unclear             |
| Injection | DNA in aqueous                      | Unknown <sup>a</sup> | <i>In vivo</i> only    | Not carefully       |

|                                | solution injected in tissue                          |  |                                    | analyzed  |
|--------------------------------|--|--|------------------------------------|---|
| Biolistics                     | Bombardment with DNA coated on gold particles        | Forced disruption of membrane            | <i>In vitro</i> and <i>in vivo</i> | In multiple copies, rate $\approx 10^{-2}$        |
| Electroporation                | Disruption of membrane barrier in electric field     | Forced disruption of membrane            | Preferentially <i>in vitro</i>     | In single copies, rate $\approx 10^{-5}$          |
| Calcium phosphate transfection | Complexation of DNA by salt precipitation            | Endocytosis                              | <i>In vitro</i> only               | In multiple copies, rate $\approx 10^{-2}$        |
| DEAE-dextran transfection      | Complexation of DNA with large sugars                | Endocytosis or chemical permeabilization | <i>In vitro</i> only               | No  |
| Polybrene/DMSO transfection    | Complexation of DNA on cationic amorphous salt       | Chemical permeabilization                | <i>In vitro</i> only               | Single or multiple copies, rate $\approx 10^{-4}$ |
| Activated dendrimers           | Complexation of DNA on cationic particles            | Endocytosis                              | <i>In vitro</i>                    | Possibly in multiple copies                       |
| Lipofection                    | Coating of DNA in aqueous solution in lipid bilayers | Membrane fusion or endocytosis           | <i>In vitro</i> and <i>in vivo</i> | Single copy possible, rate $\approx 10^{-4}$      |
| Receptor-mediated transfection | DNA bound to ligand-displaying carrier               | Endocytosis                              | <i>In vitro</i> and <i>in vivo</i> | Yes, copy number presumably carrier-dependent     |

---

<sup>a</sup> Possibly involving forced membrane disruption by locally increased pressure or nucleic acid receptor-mediated endocytosis.

In future, we will encounter an increasing diversity of materials and methods proposed for transfection. Special attention will be paid to developing selective and efficient systems that allow safe and controlled transduction of primary cells *in vitro* and *in vivo*. With an improved knowledge of the underlying mechanisms, continuous elaboration of transfection methods will thus raise the quality of basic research, as well as the economical or even therapeutic applications of molecular genetics.

Bibliography

1. F. L. Graham and A. J. van der Eb (1973) *Virology* **52**, 456.
2. L. Benimetskaya et al. (1997) *Nature Med.* **3**, 414–420.
3. R. Schubbert, C. M. Lettmann, and W. Doerfler (1994) *Mol. Gen. Genet.* **242**, 495–504.
4. A. Wolff et al. (1990) *Science* **247**, 1465–1468.
5. V. Budker et al. (1996) *Gene Ther.* **3**, 593–598.
6. N.-S. Yang et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 9568–9572.
7. N.-S. Yang (1992) *Crit. Rev. Biotechnol.* **12**, 335–356.
8. C. Sanford, F. D. Smith, and J. A. Russell (1993) *Meth. Enzymol.* **217**, 483–509.
9. D. Arnold, L. Feng, and N. Heintz (1994) *Proc. Natl. Acad. Sci. USA* **91**, 9970–9974.
10. H. Potter, L. Weis, and P. Leder (1984) *Proc. Natl. Acad. Sci. USA* **81**, 7161–7165.
11. H. Potter (1988) *Anal. Biochem.* **174**, 361–373.
12. C. Weaver (1993) *J. Cell. Biochem.* **51**, 426–435.
13. P. Singh and H. I. Maibach (1994) *Crit. Rev. Ther. Drug Carrier Syst.* **11**, 161–213.
14. B. Bommaman et al. (1994) *Pharm. Res.* **11**, 1809–1814.
15. G. Chu, H. Hayakawa, and P. Berg (1987) *Nucleic Acids Res.* **15**, 1311–1326.
16. C. Baum et al. (1994) *BioTechniques* **17**, 1059–1062.
17. C. Chen and H. Okayama (1988) *BioTechniques* **6**, 632–638.
18. M. Ishiura et al. (1982) *Mol. Cell. Biol.* **2**, 607–616.
19. J. Sussman and G. Milman (1984) *Mol. Cell. Biol.* **4**, 1641.
20. S. Kawai and M. Nishizawa (1984) *Mol. Cell. Biol.* **4**, 1172–1174.
21. R. Aubin, M. Weinfeld, and M. C. Paterson (1988) *Somatic Cell Mol. Genet.* **14**, 155–167.
22. J. Haensler and F. C. Szoka (1993) *Bioconjug. Chem.* **4**, 372–379.
23. X. Tang, C. T. Redemann, and F. C. Szoka (1996) *Bioconjug. Chem.* **7**, 703–714.
24. L. Felgner et al. (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7413–7417.
25. M. Strauss (1996) *Meth. Mol. Biol.* **54**, 307–327.
26. T. Lee and R. Jaenisch (1996) *Nucleic Acids Res.* **24**, 5054–5055.
27. M.-F. Tsan, G. L. Tsan, and J. E. White (1997) *Human Gene Ther.* **8**, 817–825.
28. I. Yanagihara et al. (1995) in *Molecular and Cell Biology of Human Gene Therapeutics*, G. Dickson, ed., Chapman & Hall, London, pp. 64–82.
29. E. Wagner et al. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 3410–3414.
30. C. Plank et al. (1992) *Bioconjug. Chem.* **3**, 533–539.

### **Suggestions for Further Reading**

31. Background information on transfection methods including key papers and concrete protocols are included in the references listed above. Moreover, two comprehensive textbooks are recommended:
32. I. Ausubel, R. Brent, R. E. Kingston, D. D. Moore, J. G. Seidman, J. A. Smith, K. Struhl, and M. Frederick. (1994) *Current Protocols in Molecular Biology*, Wiley, New York (this standard textbook of molecular biology protocols also is a rich source of detailed protocols and helpful guidelines for troubleshooting in basic transfection methods).
33. R. S. Tuan (1997) *Recombinant Gene Expression Protocols*, Humana Press Totowa, NJ (this book contains a number of protocols and background information on transfection and selection techniques employed in yeast, plants as well as nonmammalian and mammalian higher eukaryotes).

## Transfer Free Energy

The transfer free energy  $Dm_{pr,tr}$  of a protein is the increment in **free energy** of transferring the protein from water into a solvent of a given composition (eg, 1M sucrose, 8M [urea](#), 6M **guanidinium chloride**). This means that it is the free energy of interaction of the protein with the solvent of the given composition, in other words, the free energy of [binding](#) of the cosolvent to the protein,  $DG_b$ .

The transfer free energy defines quantitatively the effect of a cosolvent on [protein stability](#) (see [Stabilization And Destabilization By Co-Solvents](#), solubility), and [self-assembly](#) (such as the formation of **organelles**, eg, [microtubules](#)).

The transfer free energy is obtained directly by integration over ligand concentration of the [preferential binding](#) measured by [equilibrium dialysis](#), when that parameter is expressed by the perturbation of the chemical potential of the protein,  $m_{pr}$ , by the ligand (1):

$$(\partial\mu_{pr}/\partial m_L)_{T,P,m_{pr}} = -(\partial m_L/\partial m_{pr})_{T,P,\mu_L} (\partial\mu_L/\partial m_L)_{T,P,m_{pr}} \quad (1)$$

and

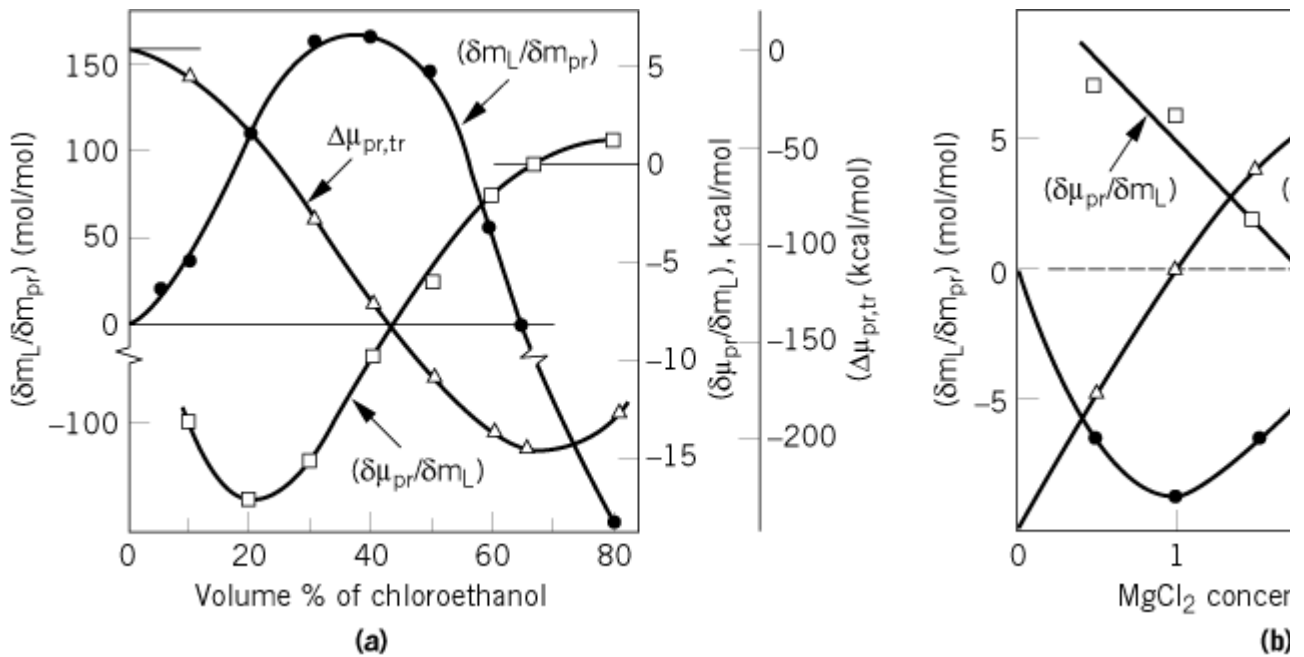
$$\Delta\mu_{pr,tr} = \int_0^L (\partial\mu_{pr}/\partial m_L)_{T,P,m_{pr}} dm_L = \Delta G_b \quad (2)$$

The term  $(\partial\mu_L/\partial m_L)_{T,P,m_{pr}}$  is the expression of the nonideality of the ligand, ie, cosolvent, and is equal to  $RT(1/m_L + \partial \ln g_L/\partial m_L)$  where  $R$  is the gas constant,  $T$  the temperature, and  $g_L$  is the *activity coefficient* of the ligand, eg, sucrose, [urea](#), the values of which can be found in many standard tables.

The variation with the cosolvent concentration of the three fundamental interaction parameters, namely, preferential binding, chemical potential perturbation, and transfer free energy, is illustrated in Figure 1 for two cosolvent systems. It is clear that their mutual variation is complex and a single point does not give much information. It is also clear that positive values of preferential binding correspond to negative values of the chemical potential perturbation, ie, the interactions are favorable, and vice versa. There is no such direct relation with the transfer free energy (free energy of binding). For example, in the 2-chlorethanol system,  $Dm_{pr,tr}$  is negative at all cosolvent concentrations, even though the preferential binding goes through a maximum and becomes negative. Just the converse is true for the b-lactoglobulin-MgCl<sub>2</sub> system, where  $Dm_{pr,tr}$  is always positive (increasingly unfavorable interactions relative to water), whereas preferential binding starts with negative values, passes through a minimum, and finally becomes positive. This reflects the difference in the reference states and illustrates the vast amount of information that may be obtained from equilibrium dialysis measurements, if these are fully interpreted.

**Figure 1.** Variation of the the thermodynamic interaction parameters with cosolvent concentration: preferential binding (interaction parameter ( $f$ )), and transfer free energy  $Dm_{pr,tr}$  ( $\Delta$ ). (a) **b-Lactoglobulin** in aqueous 2-chloroethanol; (b) b-la pH 3.0. Note that  $Dm_{pr,tr}$  is at its peak (minimal for the favorably interacting 2-chloroethanol and maximal for the unfav where the binding measured by equilibrium dialysis is zero. (Reproduced, with permission, from (2).)





## Bibliography

1. S. N. Timasheff (1995) In *Protein-Solvent Interactions* (R. B. Gregory, ed.), Marcel Dekker, New York, Chap. "11".
2. S. N. Timasheff (1993) *Ann. Rev. Biophys. Biomol. Struct.* **22**, 67–97.

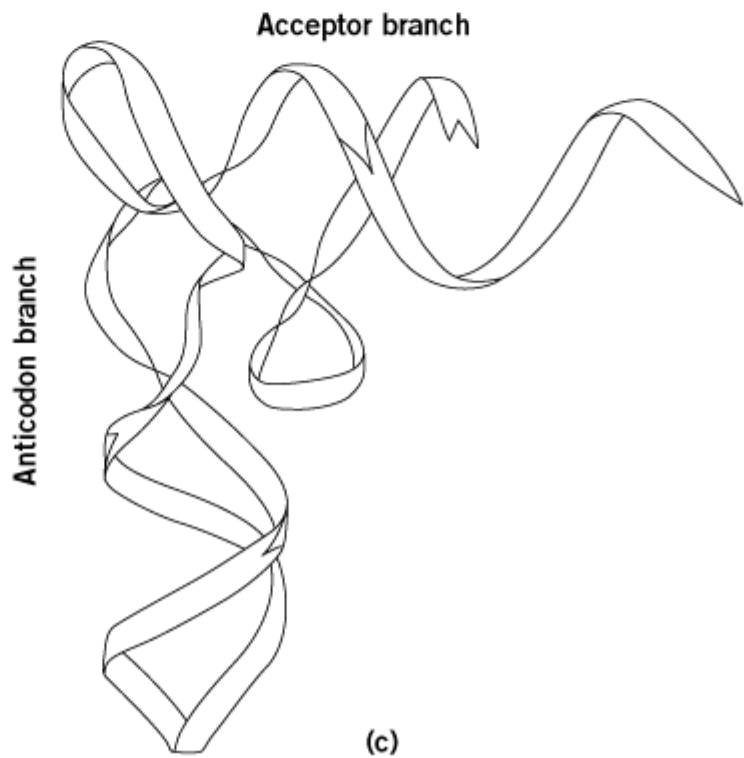
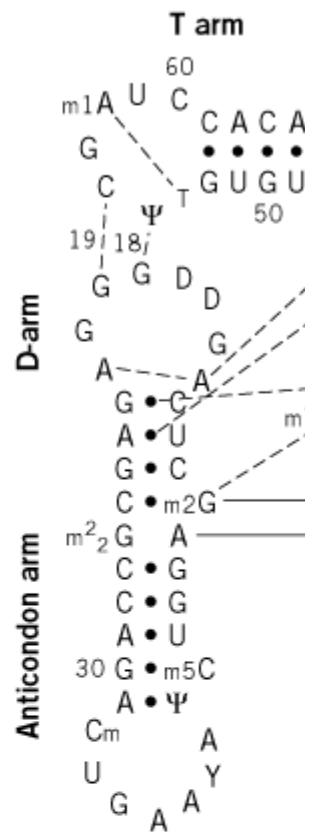
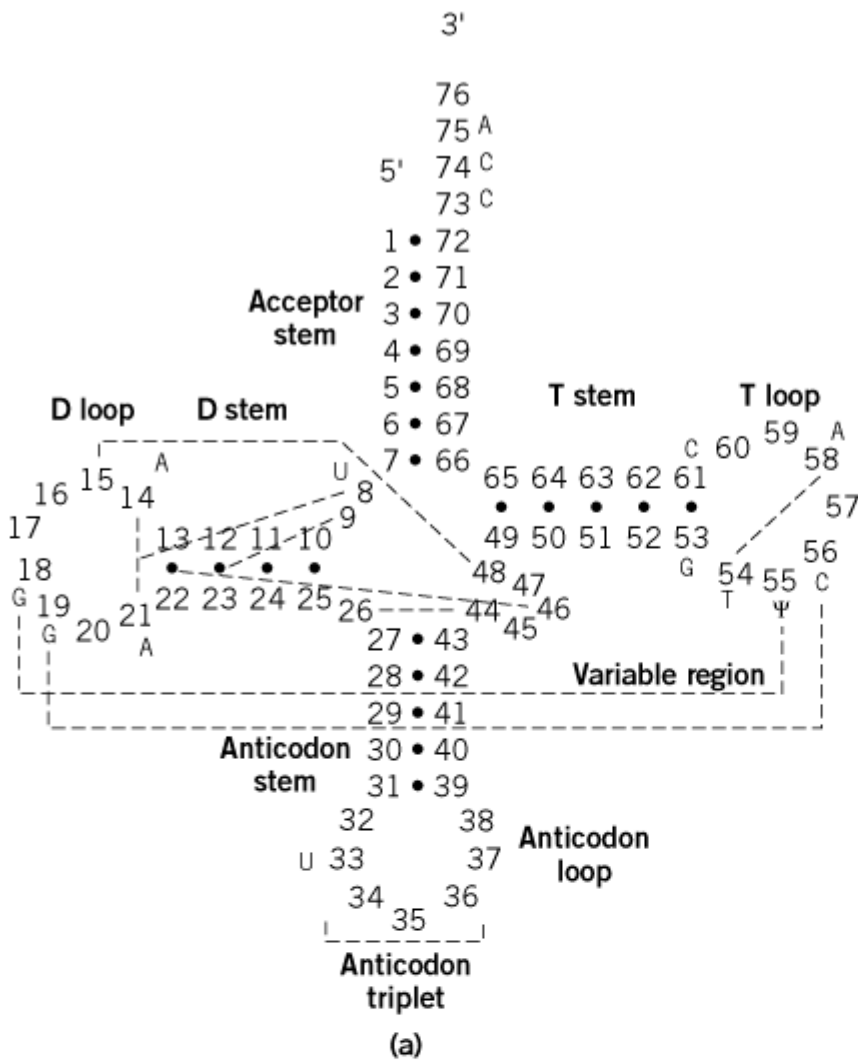
## Transfer RNA

Transfer RNAs, usually abbreviated tRNAs, are small, well-characterized RNA molecules with a key role in **protein biosynthesis**. They must interact with a number of proteins, including specific [aminoacyl-tRNA synthetases](#), **initiation** and [elongation factors](#), and components of [ribosomes](#). Specific interaction between a region of the tRNA called the [anticodon](#) and **codons** of the [messenger RNA](#) (mRNA) guarantees that the nucleotide information in the mRNA is translated correctly into protein. The concept of an adapter to provide the interface between nucleic acid language and protein language was introduced by Crick in 1955 (1). tRNA (the abbreviation for transfer ribonucleic acid) was, however, not discovered until 3 years later (2). tRNAs also participate in nonprotein synthetic capacities. They prime **reverse transcription** in [retrovirus](#) life cycles and serve as amino acid donors during synthesis of the cell wall, chlorophyll, and heme. Structural variants of tRNA, the tRNA-like molecules, are probably of distinct evolutionary origin but share some tRNA properties. These tRNA mimics are present in some messenger RNA and in a number of viral RNA, where they contribute to **translational** control and viral replication, respectively. One other small RNA, 10 S RNA, will also be considered here because it has aminoacylation properties like tRNA. Interestingly, it also has mRNA-like properties and is an important component in a process that eliminates defective proteins from **prokaryotic** cells.

### 1. Structural Properties

tRNAs are single RNA chains of 75–93 nucleotides, present in the cytosol and organelles of all living cells. Holley and his co-workers determined the first tRNA sequence in 1965 (3); today more than 2700 sequences are known either from direct sequencing of tRNA or by sequencing of their genes (4). Direct sequencing of tRNA reveals a high level of modified versions of the classic nucleotides. More than 80 different modifications have been characterized (5). All occur posttranscriptionally and may be modifications of the base (methylation, thiolations, etc) or of the ribose (methylations). Dictated by their primary sequence, tRNA fold into cloverleaf-like secondary structures with well-defined stems and loops that make up the acceptor arm, D arm and loop, anticodon arm and loop, and the T arm and loop (Fig. 1a). Regardless of the length of the tRNA, the numbering of conserved nucleotides remains constant. The acceptor stem always has 7 base pairs and 4 single-stranded nucleotides, including an absolutely conserved CCA sequence. The D stem and loop are of variable length, whereas the anticodon stem has 5 nucleotides and the anticodon loop has 7 nucleotides. The variable region usually has 4–5 nucleotides but can contain up to 24 nucleotides. Finally, the T stem always has 5 base pairs and the T loop has 7 nucleotides. High-resolution crystal structures of yeast tRNA<sup>Phe</sup> and yeast tRNA<sup>Asp</sup> revealed a three-dimensional (3D) L-shaped structure. The two domains of this structure correspond to the acceptor branch (stacking of the acceptor arm and the T arm) and the anticodon branch (anticodon arm, D arm) (Fig. 1). The angle between both branches is about 90°, and the distance between the extremities of the branches, that is, between the acceptor end and the anticodon triplet, is 75–80 Å. This 3D structure forms through a number of interactions between conserved or semiconserved nucleotides distant in the covalent structure, but close in the folded conformation, especially those in the D and T loops and within the core of the molecule (Fig. 1). Among the nine tertiary interactions present in classic tRNA, three involve triple interactions between three nucleotides (between residues 8-14-21, 9-12-23, 10-25-45, and 13-22-46).

**Figure 1.** Secondary and tertiary structural models of transfer RNA. Common subdomains are highlighted by the same color. (a) Secondary structure of tRNA. The network of long-range interactions between nucleotides responsible for the tertiary structure is shown by black lines. Nucleotides are indicated by number; conserved nucleotides are shown next to their numbered position. (b) L-shaped structure of yeast tRNA<sup>Phe</sup> represented in two dimensions. (c) Global three-dimensional folding of yeast tRNA<sup>Phe</sup> as deduced by X-ray crystallography. Modified nucleosides are I pseudouridine, m1A 1-methyladenosine, m7G 7-methylguanosine, m2G 2-methylguanosine, methylcytidine, T 5-methyluridine, and Y wybutosine.



tRNA are classified according to the length of their variable region, which can be 4 or 5 nucleotides (class I tRNA) or 10–24 nucleotides (class II tRNA). Only tRNA<sup>Leu</sup>, tRNA<sup>Ser</sup>, and eubacterial—as well as some organelle—tRNA<sup>Tyr</sup> are class II tRNA. A new class II tRNA, discovered more recently in many organisms, serves as the adapter for [selenocysteine](#), the “21st amino acid” (6). Some tRNA from mitochondria present peculiar structural features, such as the absence of one or more stems and loops within the secondary structure (7) and/or the absence of conserved nucleotides. Nonetheless, these tRNA are thought to fold into the L-shaped 3D structures necessary to fulfill their functions (8).

In each cell or organelle, there are 20 families of tRNA, one per each amino acid. Isoacceptor tRNAs belong to the same family and are charged with the same amino acid during protein synthesis. It is critical that the amino acid be appropriate to the sequence of the anticodon triplet, which, in turn, is complementary to a specific codon. Charging of tRNA with an amino acid is governed by specific aminoacyl-tRNA synthetases (aatRS). These enzymes esterify an amino acid to either the 2'- or 3'-hydroxyl group of the 3'-terminal ribose, a mechanistic aspect of the aminoacylation reaction that is used to classify aatRS into two classes (9, 10). The reaction occurs in two steps: (1) activation of the amino acid to an adenylate and (2) transfer of the activated amino acid to the tRNA. Both reactions are magnesium-dependent. Some aatRS require the presence of the homologous tRNA for recognition of the cognate amino acid (11).

### 1.1. Charging with Amino Acids

The accuracy of protein synthesis is dependent on the accuracy of aminoacylation, so the elements that dictate the recognition of tRNA by aatRS have been referred to as the second genetic code (12). Specific aminoacylation of tRNA by their cognate aatRS is governed by identity determinants present within the tRNA. These include both positive signals allowing recognition by the cognate aatRS and negative signals hindering recognition by noncognate enzymes (13-16). Positive elements continue to be identified using two similar approaches. The *in vivo* approach exploits [suppressor tRNA](#), tRNA with anticodon triplets that are altered for complementarity to codons that do not correspond to an amino acid and are referred to as [stop](#) codons. Mutations that alter the specificity of aminoacylation identify “identity elements”. The *in vitro* approach examines aminoacylation properties of purified tRNAs produced by *in vitro* transcription. Alterations in the nucleotide sequence of these tRNAs that lead to significant decreases in aminoacylation efficiency are referred to as recognition elements. Both approaches reveal a similar set of nucleotides important for aminoacylation identity, which are referred to as identity determinants .

Positive identity elements are limited in number (a single nucleotide and up to 13 nucleotides). Thus, histidine identity requires a single base added at the 5' end of the tRNA, (G1); specific alanylation is governed by a single base pair (G3-U70) in the acceptor stem; only five nucleotides within yeast tRNA<sup>Phe</sup> dictate specific phenylalanylation (A20, G34, A35, A36, and A72). Identity elements are concentrated mostly in the anticodon loop and the 3' end of the acceptor stem, including the discriminator base, first recognized in the early 1970s. Identity elements have, however, also been found in the D loop or the anticodon stem. The preferential location of identity elements at the extremities of tRNA led investigators to design minihelices , small RNAs that are missing substantial portions of the original tRNA. These minihelices contain complete or partial identity sets (17-19). Minihelices as small as the acceptor stem can be specific substrates for aminoacyl-tRNA synthetases. Anticodon-like minihelices, containing major identity elements, stimulate the aminoacylation capacity of the partner acceptor stem minihelix in yeast valine and *Escherichia coli* isoleucine systems (20, 21).

Identity sets transplanted into an unrelated host tRNA confer a new identity on the host. A simple anticodon switch changes the identity of methionine and valine tRNA (22). Insertion of the 5-nucleotide phenylalanine identity set into yeast tRNA<sup>Arg</sup>, tRNA<sup>Tyr</sup>, and tRNA<sup>Met</sup> converts these tRNA to very efficient phenylalanine acceptors. A leucine tRNA can be transformed into a serine

tRNA by changing 12 nucleotides (23). Interestingly, nonoverlapping identity sets allow the design of some tRNAs with multiple identities (24). Identity sets have been investigated for all *E. coli* tRNAs, as well as for a number of eukaryotic tRNAs. Although there are some exceptions, identity sets are usually conserved throughout different kingdoms (16).

That identity elements usually in direct contact with the aaRS is revealed by X-ray structures of complexes of the *E. coli* glutamine (25) and yeast aspartate (26, 27) tRNA and their respective synthetases. All chemical groups of an identity nucleotide are not involved in specificity. Specific recognition of alanine tRNA by its cognate synthetase is dictated by the exocyclic NH<sub>2</sub> group of guanosine G3 (28). Specificity of isoleucylation is governed by modification of the 5'-anticodon nucleotide, C34. When modified with a lysine, to become "lysidine" 34, the tRNA is efficiently recognized by isoleucyl-tRNA synthetase. In the absence of the lysidine modification, the isoleucine identity is lost and the tRNA becomes a substrate for methionyl-tRNA synthetase (29).

Optimal expression of aminoacylation identity also depends on higher order structural features, and disruption of tertiary interactions leads to dramatic losses in aminoacylation efficiencies. Some aminoacyl-tRNA synthetases are sensitive to local subtle conformational changes in their cognate tRNA. Contribution of the structure to efficient aminoacylation has been demonstrated for alanine-specific tRNA. A local distortion in the acceptor stem helix caused by the G-U base-pair identity element contributes to specificity (30). In *E. coli* tRNA<sup>Arg</sup>, identity is dependent on the presence of a "variable" pocket, formed by nucleotides from the D and T loops, which are in close vicinity in the 3D structure (31). Similarly, *E. coli* tRNA<sup>Glu</sup> has major structural identity determinants within the "augmented D helix" (32). In *E. coli* tRNA<sup>Cys</sup>, the tertiary G15-C48 pair is a structural element required for specific aminoacylation (33).

The higher-order structure necessary during the recognition step is not maintained following the initial interaction between synthetase and tRNA. Significant conformational changes do occur after synthetases bind to tRNA. During the aspartate tRNA synthetase interaction, the tRNA anticodon loop nucleotides are spread out to reach pockets in the complementary protein environment of the aspartyl-tRNA synthetase (27). Dramatic changes are also observed in the anticodon loop during the glutamine tRNA synthetase interaction, as well as at the acceptor stem, where the CCA end folds back on the acceptor stem helix (25).

Negative elements also contribute to specificity by hindering tRNA recognition by noncognate aminoacyl-tRNA synthetases. Thus, lysidine, present at position 34 of *E. coli* tRNA<sup>Ile</sup> (see discussion above), prevents recognition of this tRNA by methionyl-tRNA synthetase (29). A single methyl group present on G37 of yeast tRNA<sup>Asp</sup> hinders its recognition by the noncognate arginyl-tRNA synthetase (34).

In three cases (Arg, Glu, Gln), the tRNA not only is a substrate in the aminoacylation reaction but also plays an important role as a cofactor in allowing correct recognition of the amino acid by aaRS. Interaction between tRNA identity nucleotides and their recognition sites in glutaminyl-tRNA synthetase determines the affinity of the enzyme for the cognate amino acid (11, 35). This participation of tRNA in optimizing amino acid recognition reveals a novel mechanism for maintaining translational fidelity and also provides a strong basis for the coevolution of tRNA and their cognate synthetases.

In rare instances, production of a correctly charged tRNA depends on more than a correct aminoacyl-tRNA synthetase interaction. To produce Gln-tRNA<sup>Gln</sup> in *E. coli* and in yeast mitochondria, the tRNA is first "mischarged" with glutamic acid (Glu-tRNA<sup>Gln</sup>) and then the esterified glutamic acid is converted into glutamine. Synthesis of Asn-tRNA<sup>Asn</sup> in *Halofax volcanii* occurs in a similar fashion; aspartate is charged on tRNA<sup>Asn</sup> and subsequently converted to asparagine by

transamidation (36). Synthesis of selenocysteine (Sec)-charged tRNA<sup>Sec</sup> also requires two steps. The tRNA is aminoacylated by seryl-tRNA synthetase with serine, and in a second step, selenium is inserted into the serine by selenocysteine synthase to produce selenocysteine (6, 37).

## 2. Amino Acid Transfer

After specific aminoacylation, tRNAs arrive at the ribosome where their amino acids are connected together as [peptide bonds](#) are formed by peptidyl transferase. A particular tRNA<sup>Met</sup>, tRNA<sup>iMet</sup>, the initiator tRNA, recognizes the initiation codon AUG. Methionine residues located internally are provided by the elongator tRNA<sup>Met</sup>, tRNA<sup>eMet</sup>. Both tRNA are charged with methionine; in prokaryotes and in mitochondria, the methionine charged to the initiator tRNA is subsequently modified by formylation. Initiator tRNA bear particular structural features that allow their discrimination from elongator tRNA by initiation factors. These include the formylated amino acid, if present (38), and three consecutive G:C base pairs, conserved in the anticodon stem of virtually all initiator tRNA from eubacteria, eukaryotes, and archaeobacteria (39).

All aminoacylated elongator tRNA form ternary complexes with a GTP-activated elongation factor, which is necessary for binding to the ribosomal A site. The major elements required for recognition by prokaryotic elongation factor EF-Tu-GTP are the  $\alpha$ -NH<sub>2</sub> group of the bound amino acid and nucleotide A76. Additional direct contacts between the tRNA and the protein occur via riboses and phosphates contained within the amino acid acceptor stem, as revealed by the [X-ray crystallography](#) structure of the ternary complex formed by yeast tRNA<sup>Phe</sup>, elongator factor EF-Tu from *Thermus aquaticus*, and a structural analog of GTP (40). Neither conserved nucleotides nor modified bases within this domain are necessary.

Negative elements allow elongation factors to discriminate initiator from elongator tRNA. For prokaryotic tRNA<sup>fMet</sup>, the formyl group present specifically on the  $\alpha$ -amino group of methionine is the negative element. A modification at the 2' carbon of the ribose at position 64 in plant and fungi tRNA<sup>iMet</sup> provides this function. An additional signal that distinguishes initiator from elongator tRNA is the absence of a base pair between residues 1 and 72.

Interaction of aminoacylated tRNA with the ribosome depends on initiation or elongation factors. Initiator tRNA enters the ribosomal P (peptidyl-) site. Elongator tRNAs enter the A (aminoacyl) site, matching the anticodon with the mRNA codon, and acquire the growing peptide chain attached to the tRNA in the P site. The peptidyl-transferase activity of the large ribosomal subunit catalyzes a nucleophilic attack by the amino acid of the A site tRNA on the carboxy-terminal amino acid of the growing peptide chain. The tRNA carrying the peptide chain then moves to the P site by [translocation](#). Uncharged tRNAs leave the ribosome via the E (exit) site. Many studies, including [footprinting](#) and [crosslinking](#), seek to define detailed information about the interactions between tRNA and the ribosome (41-44). It is known that the anticodon of the tRNA interacts with the small ribosomal subunit, whereas the acceptor stem, especially the -CCA-amino acid end, interacts with the large ribosomal subunit. Direct contacts with *E. coli* 23 S RNA have been observed. Recently, three tRNAs bound to *E. coli* 70 S ribosomes were directly visualized with [cryoelectron microscopy](#) (45). The detailed arrangement of A- and P-site tRNA inferred from this study by 3D reconstruction allowed localization of the ribosome sites for anticodon interaction and peptide bond formation.

## 3. Other Functions

tRNAs have several functions in addition to ribosomal-dependent protein biosynthesis. Specific host tRNA are used to prime minus-strand DNA synthesis during the replication cycle of retroviruses. tRNA<sup>Trp</sup>, tRNA<sup>Pro</sup>, tRNA<sup>Lys</sup>, and tRNA<sup>iMet</sup> are most widely used as primers. They are specifically encapsidated in viral particles. Reverse transcription of the genomic RNA of human immunodeficiency type 1 virus (HIV-1) is primed by tRNA<sup>3Lys</sup>. Primer tRNA bind to the primer

binding site (PBS) of the viral genome via hybridization of 18 or more (31 nucleotides of tRNA<sup>3<sup>Lys</sup></sup>) 3'-terminal nucleotides, leaving the 3'-OH accessible for initiation of reverse transcription. In some cases, additional interactions take place between the primer tRNA and the retro-RNA outside the PBS. Modified nucleosides of the tRNA contribute to the efficiency and accuracy of the different steps of the process. Since tRNA secondary and tertiary structures are very stable, hybridization to the viral RNA requires assistance from nucleocapsid protein and reverse transcriptase (46).

tRNA<sup>Glu</sup> is involved in the formation of d-aminolevulinic acid, a precursor for heme and chlorophyll synthesis. In this process, tRNA<sup>Glu</sup>, charged with glutamate, is the substrate for Glu-tRNA reductase, which reduces Glu-tRNA<sup>Glu</sup> to glutamate 1-semi-aldehyde and releases free tRNA<sup>Glu</sup>. Glutamate semialdehyde is further converted to d-aminolevulinic acid by glutamate-1-semialdehyde-2, 1-aminomutase in a reaction that is independent of tRNA (47).

tRNA-like domains are present at the 3' end of some RNA genomes of some plant viruses (48). These domains fold into simple L shapes or more complex conformations that include an L shape (49). These structures are recognized by **nucleotidyltransferases**, are efficiently aminoacylated by host or heterologous aminoacyl-tRNA synthetases, and interact with elongation factors. Tymoviruses are typically charged with tRNA for valine, tobamoviruses with histidine, and bromoviruses with tyrosine. Rules for aminoacylation are the same as those for classic tRNA. Despite their functional and structural similarities to classic tRNA, these tRNA-like domains are not involved in protein synthesis. They serve as initiation sites for viral replication, and also play a role similar to that of telomeres by protecting the viral RNA against degradation (50). In bromoviruses, aminoacylation is not required for replication. In tymoviruses, the presence of an amino acid is necessary for replication (51). tRNA-like structures are present in the upstream leader sequences of some *E. coli* mRNA that code for aminoacyl-tRNA synthetases and are involved in translational regulation of aatRS synthesis (52). The tRNA-like structure of the threonyl-tRNA synthetase mRNA is the best characterized. The leader folds into two anticodon-like arms. When the level of cellular tRNA<sup>Thr</sup> is low, free ThrRS binds to the leader sequence instead. This hinders binding of the ribosome to the [Shine–Dalgarno](#) sequence and prevents use of the initiation codon. Alternatively, high levels of tRNA<sup>Thr</sup> engage all the available ThrRS, and translation from the mRNA proceeds (53).

#### 4. 10SaRNA/tmRNA

10SaRNA or tmRNA, a 350-nucleotide RNA, has a dual role as both tRNA and mRNA and is present in a wide variety of prokaryotes (54). It has features within its secondary structure that mimic canonical tRNA arms (55). Aminoacylated with alanine, it enters the A site of ribosomes that are stalled on mRNA and thus cannot terminate protein synthesis properly. Peptidyl transferase catalyzes the transfer of the nascent peptide to the alanine, and the ribosome transfers to the 10Sa/tmRNA; it continues to translate the 10Sa/tmRNA, until a stop codon is reached and the ribosome, the protein, and the last tRNA are released normally. All proteins released by this process have the same carboxy-terminal amino acid sequence, since 10Sa/tmRNA directs its synthesis. This carboxy-terminal tag is a signal for directing the truncated, defective proteins to specific cellular [proteinases](#). Thus, 10Sa/tmRNA protects the cell from defective proteins and allows ribosome cycling by directing ribosome switching from one mRNA to another (56).

#### 5. Neuromuscular Diseases

A large number of human neuromuscular genetic diseases resulting in severe oxidative phosphorylation deficiencies correlate with the occurrence of point mutations in mitochondrial tRNA genes (57, 58). Mitochondrial DNA in individuals suffering from these diseases is heteroplasmic, and the severity of the symptoms generally relates to the ratio of wild-type to mutant DNA. To date, 32 independent mutations have been found in 13 tRNA genes. Among the most widely studied are mutation 3243 in tRNA<sup>Leu</sup> (UUR) gene, responsible for mitochondrial myopathy with lactic acidosis

and stroke-like episodes (MELAS) and for diabetes mellitus, and mutation 8344 in the gene of tRNA<sup>Lys</sup>, responsible for myoclonic epilepsy with ragged red fibers (MERRF). The molecular mechanisms by which mutations in the genes cause the associated phenotypes are not yet clear.

### Bibliography

1. F. H. C. Crick (1955) in E. Chargaff, and J. N. Davidson, eds., *The Nucleic Acids*. Academic Press, New York, Vol. **3**, p. 349.
2. M. B. Hoagland, M. L. Stephenson, J. F. Scott, L. I. Hecht, and P. C. Zamecnik (1958) *J. Biol. Chem.* **231**, 241–257.
3. R. W. Holley, J. Apgar, G. A. Everett, J. T. Madison, M. Marquise, S. H. Merrill, J. R. Penswick, and R. Zamir (1965) *Science* **147**, 1462–1465.
4. M. Sprinzl, C. Steegborn, F. Hubel, and S. Steinberg (1996) *Nucleic Acids Res.* **24**(1), 68–72.
5. P. A. Limbach, P. F. Crain, and J. A. McCloskey (1994) *Nucleic Acids Res.* **22**, 2183–2196.
6. C. Baron, and A. Böck (1995) in D. Söll, and U. L. RajBhandary, eds., *tRNA: Structure, Biosynthesis, and Function*, American Society for Microbiology Press, Washington, D.C., pp. 529–544.
7. R. Okimoto, and D. R. Wolstenholme (1990) *EMBO J.* **9**, 3405–3411.
8. D. R. Wolstenholme, R. Okimoto, and J. L. McFarlane (1994) *Nucleic Acids Res.* **22**, 4300–4306.
9. G. Eriani, M. Delarue, O. Poch, J. Gangloff, and D. Moras (1990) *Nature* **347**, 203–206.
10. S. Cusack, C. Berthet-Colominas, M. Hartlein, N. Nassar, and R. Leberman (1990) *Nature* **347**, 249–255.
11. K.-W. Hong, M. Ibba, I. Weygand-Durasevic, M. J. Rogers, H.-U. Thomann, and D. Söll (1996) *EMBO J.* **15**(8), 1983–1991.
12. C. de Duve (1988) *Nature* **333** (May), 117–118.
13. J. Cavarelli, and D. Moras (1993) *FASEB J.* **7**, 79–86.
14. R. Giegé, J. D. Puglisi, and C. Florentz (1993) *Prog. Nucleic Acid Res. Mol. Biol.*, **45**, 129–206.
15. W. H. McClain (1993) *J. Mol. Biol.* **234**, 257–280.
16. M. E. Saks, and J. R. Sampson (1995) *J. Mol. Evol.* **40**, 509–518.
17. C. Francklyn, K. Musier-Forsyth, and P. Schimmel (1992) *Eur. J. Biochem.* **206**, 315–321.
18. C. Francklyn and P. Schimmel (1990) *Proc. Natl. Acad. Sci. USA* **87**, 8655–8659.
19. S. A. Martinis, and P. Schimmel (1995) In D. Söll, and U. L. RajBhandary, eds., *tRNA: Structure, Biosynthesis, and Function*, American Society for Microbiology Press, Washington, D.C., pp. 349–370.
20. M. Frugier, C. Florentz, and R. Giegé (1992) *Proc. Natl. Acad. Sci. USA* **89**(9), 3990–3994.
21. O. Nureki, T. Niimi, T. Muramatsu, H. Kanno, T. Kohno, C. Florentz, R. Giegé, and S. Yokoyama (1994) *J. Mol. Biol.* **236**, 710–724.
22. L. H. Schulman, and H. Pelka (1988) *Science* **242**, 765–768.
23. J. Normanly, R. C. Ogden, S. J. Horvath, and J. Abelson (1986) *Nature* **321**, 213–219.
24. M. Frugier, C. Florentz, P. Schimmel, and R. Giegé (1993) *Biochemistry* **32**, 14053–14061.
25. M. A. Rould, J. J. Perona, D. Söll, and T. A. Steitz (1989) *Science* **246**, 1135–1142.
26. M. Ruff, S. Krishnaswamy, M. Boeglin, A. Poterszman, A. Mitschler, A. Podjarny, B. Rees, J.-C. Thierry, and D. Moras (1991) *Science* **252**, 1682–1689.
27. J. Cavarelli, B. Rees, M. Ruff, J.-C. Thierry, and D. Moras (1993) *Nature* **362**, 181–184.
28. K. Musier-Forsyth, and P. Schimmel (1992) *Nature* **357**, 513–515.
29. T. Muramatsu, K. Nishikawa, F. Nemoto, Y. Kuchino, S. Nishimura, T. Miyazawa, and S. Yokoyama (1988) *Nature* **336**, 179–181.
30. W. H. McClain, K. Gabriel, and J. Scheinder (1996) *RNA* **2**, 105–109.



31. W. H. McClain, and K. Foss (1988) *Science* **241**, 1804–1807.
32. S. Sekine, O. Nureki, K. Sakamoto, T. Niimi, M. Tateno, M. Go, T. Kohno, A. Brisson, J. Lapointe, and S. Yokoyama (1996) *J. Mol. Biol.* **256**, 685–700.
33. Y.-M. Hou, E. Westhof, and R. Giegé (1993) *Proc. Natl. Acad. Sci. USA* **90**, 6776–6780.
34. J. Putz, C. Florentz, F. Benseler, and R. Giegé (1994) *Nature Struct. Biol.* **1**, 580–582.
35. M. Ibba, K.-W. Hong, J. M. Sherman, and D. Söll (1996) *Proc. Natl. Acad. Sci. USA* **93**, 6953–6958.
36. A. W. Curnow, M. Ibba, and D. Söll (1996) *Nature* **382**, 589–590.
37. N. Hubert, R. Walczak, C. Sturchler, E. Myslinski, C. Schuster, E. Westhof, P. Carbon, and A. Krol (1996) *Biochimie* **78**, 590–596.
38. S. Li, N. V. Kumar, U. Varshney, and U. L. RajBhandary (1996) *J. Biol. Chem.* **271**(2), 1022–1028.
39. N. Mandal, D. Mangroo, J. J. Dalluge, J. A. McCloskey, and U. L. Rajbhandary (1996) *RNA* **2**, 473–482.
40. P. Nissen, M. Kjeldgaard, S. Thirup, G. Polekhina, L. Reshetnikova, B. F. C. Clark, and J. Nyborg (1995) *Science* **270**, 1464–1472.
41. D. Moazed and H. F. Noller (1990) *J. Mol. Biol.* **211**(1), 135–145.
42. A. Huttenhofer, and H. F. Noller (1992) *Proc. Natl. Acad. Sci. USA* **89**, 7851–7855.
43. J. Rinke-Appel, N. Junke, M. Osswald, and R. Brimacombe (1995) *RNA* **1**, 1018–1028.
44. U. von Ahsen and H. F. Noller (1995) *Science* **267**, 234–237.
45. R. K. Agrawal, P. Penczek, R. A. Grassucci, Y. Li, A. Leith, K. H. Nierhaus, and J. Frank (1996) *Science* **271**, 1000–1002.
46. C. Isel, C. Ehresmann, G. Keith, B. Ehresmann, and R. Marquet (1995) *J. Mol. Biol.* **247**, 236–250.
47. E. Verkamp, A. M. Kumar, A. Lloyd, O. Martins, N. Stange-Thomann, and D. Söll (1995) In D. Söll, and U. L. RajBhandary, eds., *tRNA: Structure, Biosynthesis and Function*, American Society for Microbiology Press, Washington, D.C., pp. 545–550.
48. C. Florentz and R. Giegé (1995) In D. Söll, and U. L. RajBhandary, eds., *tRNA: Structure, Biosynthesis, and Function*, American Society for Microbiology Press, Washington, D.C., pp. 141–163.
49. C. W. A. Pleij (1994) *Curr. Opin. Struct. Biol.* **4**, 337–344.
50. A. M. Weiner and N. Maizels (1987) *Proc. Natl. Acad. Sci. USA* **84**, 7383–7387.
51. T. W. Dreher, C.-H. Tsai, and J. M. Skuzeski (1996) *Proc. Natl. Acad. Sci. USA* **93**, 12212–12216.
52. M. Springer, M. Graffe, J. Dondon, and M. Grunberg-Manago (1989) *EMBO J.* **8**, 2417–2424.
53. P. Romby, J. Caillet, C. Ebel, C. Sacerdot, M. Graffe, F. Eyermann, C. Brunel, H. Moine, C. Ehresmann, B. Ehressman, and M. Springer (1996) *EMBO J.* **15**, 5976–5987.
54. Y. Komine, M. Kitabatake, T. Yokogawa, and K. Nishikawa (1994) *Proc. Natl. Acad. Sci. USA* **91**, 9223–9277.
55. B. Felden, H. Himeno, A. Muto, J. P. McCutcheon, J. F. Atkins, and R. F. Gesteland (1997) *RNA* **3**, 89–103.
56. K. C. Keiler, R. H. Waller, and R. T. Sauer (1996) *Science* **271**, 990–993.
57. D. C. Wallace (1992) *Annu. Rev. Biochem.* **61**, 1175–1212.
58. N.-G. Larsson, and D. A. Clayton (1995) *Annu. Rev. Genet.* **29**, 151–178.

### Suggestions for Further Reading

59. D. Söll and U. L. Rajbhandary (1995) *tRNA Structure, Biosynthesis, and Function*, American Society for Microbiology, Washington, D.C.

60. K. H. Nierhaus, ed. (1993) *The Translational Apparatus*, Plenum Press, New York.
61. D. Söll (1993) "Transfer RNA: An RNA for all seasons," in *The RNA World*, R. Gesteland and J. F. Atkins, eds., Cold Spring Harbor Laboratory Press., New York.
62. W. H. McClain (1993) Transfer RNA Identity, *FASEB J.* **7**, 72–78.

## Transferase

Transferases are a diverse group of [enzymes](#) with respect to the identity of the groups that they transfer and the structure of the acceptor molecules.

1. The transfer of one-carbon groups is catalyzed by methyltransferases; hydroxymethyl-, formyl-, and related transferases; carboxyl- and carbamoyltransferases; and amidinotransferases.
2. Acyl groups are transferred by the acyltransferases, with the formation of esters or amides, and the donor is usually the corresponding coenzyme A (CoA) derivative. Aminoacyl group transfer also occurs.
3. Glycosyltransferases form a major subdivision within the transferase group of enzymes; the general reaction is the transfer of a sugar moiety from oligosaccharides or a high-energy compound, such as UDP-galactose, to another carbohydrate molecule. They can catalyze hydrolytic and phosphorolytic reactions by transfer of a glycosyl group to water and inorganic phosphate, respectively.
4. *Aminotransferases* catalyze the transfer of an amino group from an amino acid to a keto acid through the involvement of [pyridoxal phosphate](#) (see [Coenzyme](#), [Cofactor](#)).
5. The enzymes concerned with the transfer of phosphorus-containing groups include the phosphotransferases (see above), diphosphotransferases, nucleotidyltransferases, and the transferases that utilize CDP, GDP, and UDP derivatives.
6. There are also the sulfurtransferases, the sulfotransferases, and the CoA transferases.

## Transferrin

No other metal surpasses iron in versatility of uses to which it is put by nature. Iron is essential for the transport of oxygen by hemoglobin and hemerythrin; it has critical functions in the energy-transducing pathways of electron transport; it is indispensable in the DNA-synthesizing activity of mammalian ribonucleotide reductase; it participates in the catalysis of biological oxidations and oxygenations by dioxygen and hydrogen peroxide, as well as in protection against oxidative stress induced by these sometimes noxious species; it has a central function in nitrogen fixation; it engages in a rich variety of reactions entailing hydration-dehydration or hydrolytic cleavage; it is a magnetosensor for orienting organisms as diverse as magnetotactic bacteria, honey bees, and pigeons; and it is even used for mineralization of teeth by invertebrates of the limpet family.

The need for iron is accompanied by a corresponding need for iron carriers. In an aerobic environment at physiological pH the stable state of iron is the ferric,  $\text{Fe}^{3+}$ , but hydrolysis of ferric

iron to form insoluble polynuclear complexes of ferric hydroxide limits the concentration of simple aquated  $\text{Fe}^{3+}$  to  $10^{-17}$  M. Iron carriers are needed to overcome the solubility problem and to keep iron from catalyzing the generation of noxious and reactive oxygen species — an event prone to occur when iron is surrounded by reductants of metabolism — the pyridine nucleotides or reduced flavins, for example. Iron carriers are as indispensable to life as iron. Derangements of iron metabolism, whether hereditary as in hemochromatosis, or acquired as in transfusion-dependent or nutritional iron-deficiency anemias, are among the most common, ubiquitous and dangerous of human maladies.

## 1. Transferrins: General Descriptions

In the universe of multicellular organisms with circulatory systems, including species at least as ancient as the cockroach and tobacco hornworm (1-3), iron transport is often managed by the transferrins, a class of single-chain, two-sited, iron-binding proteins. The human transferrin molecule is a bilobar structure of 769 amino acids and two multiantennary glycans, with a molecular weight of near 80,000 depending on the branching of its two glycans (4, 5). A 45% identity in amino acid sequences of the two lobes, with conservative equivalence of most of the nonidentical residues, and a nearly exact homology in intron/exon organization of corresponding halves of the human transferrin gene (6), establish that the modern transferrin structure arose from duplication and fusion of an ancestral gene specifying a single-sited transferrin precursor. No species is known in which a verified single-sited half-transferrin persists, nor is the evolutionary advantage enjoyed by a two-sited protein clearly understood.

The transferrins are customarily divided into three major subclasses: serum transferrin, ovotransferrin (formerly called conalbumin), and lactoferrin or lactotransferrin, depending on the source from which each is isolated. Iron-binding ligands are identical among these subclasses, as is the general organization of the protein molecules. Each lobe of a transferrin molecule is divided into two domains, designated N1 or C1 and N2 or C2, depending on location in B- or C-lobes, respectively. Thus, the protein molecule consists of two lobes and four domains. The lobes are joined by a short connecting  $\beta$ -strand in human transferrin, by a loop structure in hen ovotransferrin (7), and by a three-turn helix in human lactoferrin (8). Slight differences in arrangements of lobes and domains help account for the sometimes striking differences in iron-binding and functional properties among the transferrins.

Serum transferrin is, under normal circumstances, the major or sole source of iron for the metabolic needs of most cells. In man, the half-life of circulating transferrin is about 7.5 days, whereas that of transferrin-borne iron is about 1.7 h (9), implying that a transferrin molecule experiences over 100 cycles of iron transport during its lifetime. Understanding mechanisms of iron binding and release is therefore central to understanding the biological function of transferrin. Like mammalian transferrin, ovotransferrin, the product of the same gene in the chicken as serum transferrin and differing from the serum protein only in its glycosylation, also can serve as an iron source for developing chick embryo red cells (10). Its primary function, however, may be in sequestering iron, thereby defending the egg against microbial invasion. In contrast to these proteins, the principal role of lactoferrin appears to be in conferring bacteriostatic activity to the milk and other physiological fluids in which it is found. Such activity originates not only from the strong iron-binding properties of the protein, but in a basic sequence of its N-lobe that has been termed lactoferricin (11). Lactoferrin sequesters iron much more tenaciously than do serum transferrin and ovotransferrin, so that release of iron to cells probably entails intracellular degradation of the protein. The concentration of lactoferrin in blood plasma is less than 1 mg/ml, but its turnover, primarily the result of uptake by hepatocytes (12), is sufficiently rapid so that lactoferrin may serve a true transport function for iron, complementing that of serum transferrin.

## 2. Regulation of Transferrin Expression

The human transferrin gene lies on chromosome 3 (4), as do the genes for the transferrin receptor and ceruloplasmin. Transferrin is expressed by many cells largely in a constitutive manner, the liver being by far the most important source of the circulating protein (13). Secretion is promoted by cytokines (14); expression is enhanced by estrogen, hypoxia, iron deficiency anemia, inflammation, and other factors, but seldom to a pronounced degree. Regulation may be at transcriptional and post-transcriptional levels, with tissue-specific transcriptional factors identified in Sertoli cells and hepatocytes (15).

### 3. Structure of Serum Transferrin

Human serum transferrin consists of 679 amino acids disposed in a single polypeptide chain divided into two highly homologous lobes (46% identity of amino acid residues, and a much higher degree of conservative structure) roughly equal in size. Domain 1 of each lobe is formed by the first 90 or so amino acids of that lobe. The main chain then crosses to form the full second domain comprising some 150 residues, before returning to complete the first domain. Thus, the two domains are connected by a pair of antiparallel strands; the two lobes are joined by a single short  $\beta$ -strand of some seven residues. Each domain and each interdomain connecting strand provides a ligand for the iron-binding site: Asp 63 (392) from domain 1 of N-lobe (or C-lobe), Tyr 95 (426) from the first strand connecting domain 1 to domain 2, Tyr 188 (517) from domain 2, and His 249 (485) from the connecting strand returning from domain 2 to complete domain 1. Thus, the protein contributes four ligands to each of its specific sites; the remaining two coordination requirements of  $\text{Fe}^{3+}$  are satisfied by a carbonate anion bound in bidentate manner to the metal ion and further joined to the protein in a complex network of hydrogen and electrostatic bonds.

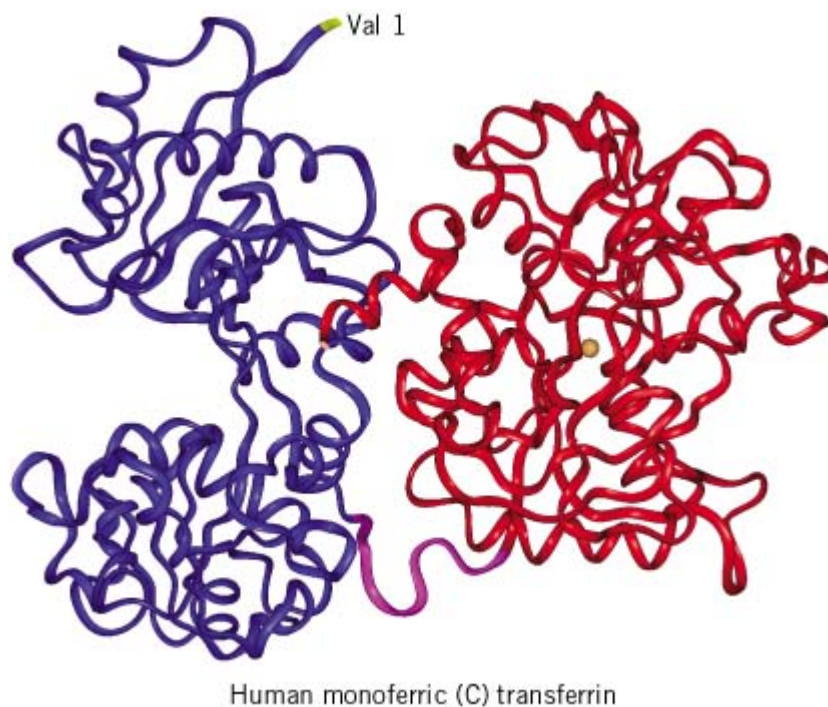
### 4. Carbonate and the Binding of Iron to Transferrin

Binding of iron and carbonate are highly cooperative, neither being strongly bound in the absence of the other. This interdependence of iron and carbonate binding may be taken as the unique and defining characteristic of all transferrins. A variety of studies suggest that carbonate binding precedes iron binding, perhaps by introducing negative charge at the binding site in preparation for receiving positively charged iron (16, 17). For this reason, carbonate (or, in its absence, other bifunctional anion capable of occupying the anion-binding site) is often referred to as the synergistic anion (18). Iron is always bound in the ferric state,  $\text{Fe}^{3+}$ , the binding of the ferrous form,  $\text{Fe}^{2+}$ , being much too weak to be of practical significance.

### 5. Conformational Changes Associated with Binding and Release of Iron

Each lobe of transferrin, when devoid of iron, assumes an open-jaw configuration of its two component domains. Binding of iron is accompanied by a rotation of the domains with respect to each other, thereby closing the cleft bearing the iron-binding site to secure the metal and shield it from hydrolysis and formation of insoluble ferric hydroxide (Fig. 1). Interdomain interactions, and to a lesser extent interlobe interactions, contribute to stabilization of bound iron (19-21).

**Figure 1.** Crystal structure at 3.3 Å resolution of monoferric transferrin bearing iron in the C-terminal lobe. Residue Val 1 in the empty N-lobe is labeled, and the iron atom in the C-lobe is depicted as a pale pink ball. The N-lobe is in blue, the C-lobe in red, and the interlobe connecting strand in violet.



## 6. Binding and Release of Iron by Transferrin

Transferrin will accept iron in a variety of forms, ranging from free  $\text{Fe}^{2+}$  to  $\text{Fe}^{3+}$  bound to a large variety of chelators. A recently described iron transporter, IREG1 (also known as ferroportin1 and MTP1) (22-24), probably acting in concert with hephaestin, a multicopper oxidase (25), has been implicated in the export of iron from duodenum and reticuloendothelial cells. Very likely, therefore, iron experiences a sequence of intracellular redox reactions before it is presented to circulating transferrin as  $\text{Fe}^{3+}$ .

Under physiological conditions, the effective stability constants for the binding of iron by transferrin exceed  $10^{21} \text{ M}^{-1}$ . Despite this great affinity for  $\text{Fe}^{3+}$ , the transferrin molecule can deliver its iron to erythroid cells within a minute or two without damage to its structure. How this is accomplished is still only dimly understood, but recent studies, reinforcing old conjectures, have offered insights into the mechanisms of reversible iron release. Protonation of the synergistic carbonate anion as pH is lowered, with consequent disruption of its coordination to iron, undoubtedly contributes to release of iron from transferrin within acidified endosomes (26-28). A further effect of low pH is seen in the N-lobe, where two lysine residues on opposite domains lie in close apposition to each other (29). At the pH of the circulation, 7.4, both lysines cannot be positively charged by protonation. As pH is lowered in the endosome, however, both lysines become protonated and mutually repulsive, thereby helping force the domains into the open configuration from which release is facile (30). Other protonations, of iron ligands or remote residues, as well as ligand exchange reactions, are likely to be involved in iron release as well (31, 32). Binding of transferrin to its receptor also modulates iron release, impeding at extracellular pH, 7.4, but facilitating at endosomal pH, 5.6 (33). Even at endosomal pH, iron sequestering agents are required to achieve release of iron in physiologically reasonable times. The chemical nature and concentration of such agents and the ionic strength of the solution in which release occurs markedly affects release rates (34), with differing effects on the two lobes of transferrin, but the physiological roles of these variables are not known. At the cellular level, binding of the hemochromatosis protein, HFE, to receptor interferes with concomitant binding of transferrin and uptake of its iron (35, 36).

Because the endosomal iron transporter DMT1 requires  $\text{Fe}^{2+}$  (37), a reductive event during iron release from transferrin and its export to the cytoplasm must occur, but where and how this takes place is unknown. Intestinal absorption of iron also involves reduction of  $\text{Fe}^{3+}$ , again emphasizing the role of redox reactions in iron metabolism (38). The many concerted interactions among transferrin and its partner proteins of iron metabolism are still incompletely understood and subjects of active investigation by many laboratories.

## 7. Acknowledgments

Preparation of this article was supported, in part, by Grant DK15056 from the National Institutes of Health, U.S. Public Health Service.

The author thanks Dr. Harmon Zuccola for the crystal structure coordinates of human monoferric transferrin bearing iron in the C-terminal lobe.

## Bibliography

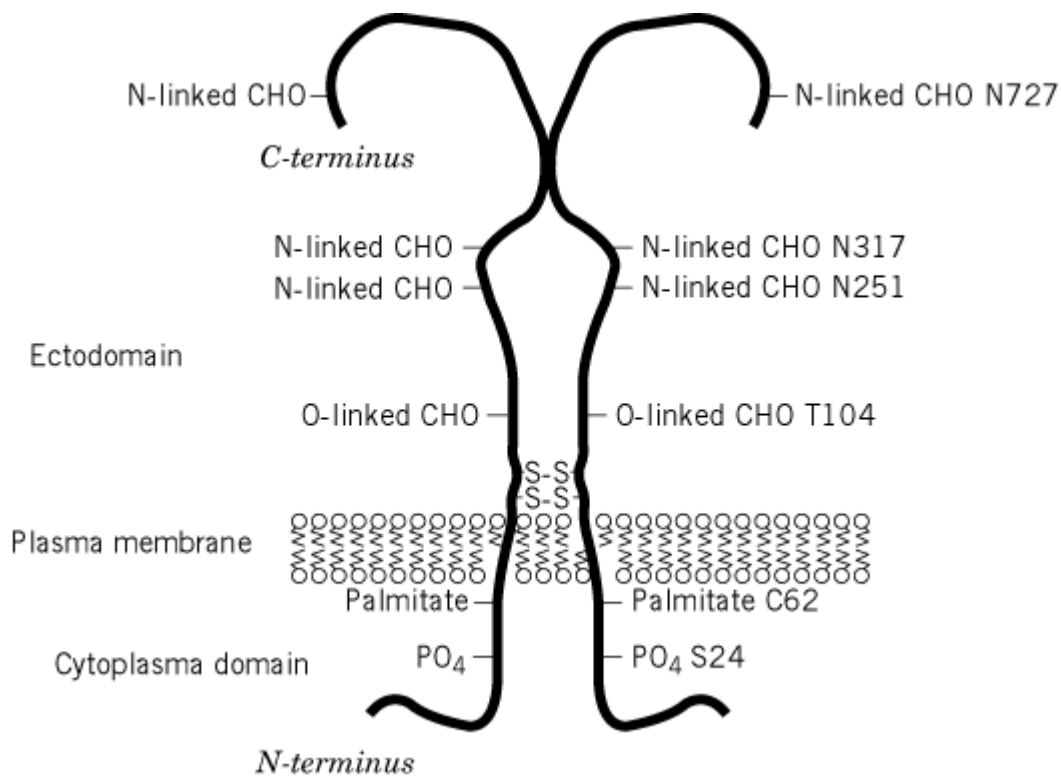
1. R. C. Jamroz, J. R. Gasdaska, J. Y. Bradfield, and J. H. Law (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 1320–1324.
2. J. R. Gasdaska, J. H. Law, C. J. Bender, and P. Aisen (1996) *J. Inorg. Biochem.* **64**, 247–258.
3. N. S. Bartfeld and J. H. Law (1990) *J. Biol. Chem.* **265**, 21684–21691.
4. F. Yang, J. B. Lum, J. R. McGill, C. M. Moore, S. L. Naylor, P. H. van Bragt, W. D. Baldwin, and B. H. Bowman (1984) *Proc. Natl. Acad. Sci. U.S.A.* **81**, 2752–2756.
5. A. B. Mason, M. K. Miller, W. D. Funk, D. K. Banfield, K. J. Savage, R. W. A. Oliver, B. N. Green, R. T. A. MacGillivray, and R. C. Woodworth (1993) *Biochemistry* **32**, 5472–5479.
6. I. Park, E. Schaeffer, A. Sidoli, F. E. Baralle, G. N. Cohen, and M. M. Zakin (1985) *Proc. Natl. Acad. Sci. U.S.A.* **82**, 3149–3153.
7. H. Kurokawa, B. Mikami, and M. Hirose (1995) *J. Mol. Biol.* **254**, 196–207.
8. E. N. Baker (1994) *Adv. Inorg. Chem.* **41**, 389–461.
9. J. H. Katz (1961) *J. Clin. Invest.* **40**, 2143–2152.
10. S. C. Williams and R. C. Woodworth (1973) *J. Biol. Chem.* **248**, 5845–5853.
11. K. S. Hoek, J. M. Milne, P. A. Grieve, D. A. Dionysius, and R. Smith (1997) *Antimicrob. Agents Chemother.* **41**, 54–59.
12. D. D. McAbee (1995) *Biochem. J.* **311**, 603–609.
13. E. H. Morgan (1981) *Mol. Aspects Med.* **4**, 3–123.
14. E. Hoeben, J. Van Damme, W. Put, J. V. Swinnen, and G. Verhoeven (1996) *Endocrinology* **137**, 514–521.
15. E. Schaeffer, F. Guillou, D. Part, M. M. Zakin (1993) *J. Biol. Chem.* **268**, 23399–23408.
16. J. L. Zweier, J. B. Wooten, and J. S. Cohen (1981) *J. Biol. Chem.* **20**, 3505–3510.
17. O. W. Nadeau, A. M. Falick, and R. C. Woodworth (1996) *Biochemistry* **35**, 14294–14303.
18. M. R. Schlabach and G. W. Bates (1975) *J. Biol. Chem.* **250**, 2182–2188.
19. P. Aisen, A. Leibman, and J. Zweier (1978) *J. Biol. Chem.* **253**, 1930–1937.
20. C. L. Day, K. M. Stowell, E. N. Baker, and J. W. Tweedie (1992) *J. Biol. Chem.* **267**, 13857–13862.
21. O. Zak, P. Aisen, J. B. Crawley, C. L. Joannou, K. J. Patel, M. Rafiq, and R. W. Evans (1995) *Biochemistry* **34**, 14428–14434.
22. A. T. McKie, P. Marciani, A. Rolfs, K. Brennan, K. Wehr, D. Barrow, S. Miret, A. Bomford, T. J. Peters, F. Farzaneh, M. A. Hedger, M. W. Hentze, and R. J. Simpson (2000) *Mol. Cell* **5**, 299–309.
23. A. Donovan, A. Brownlie, Y. Zhou, J. Shepard, S. J. Pratt, J. Moynihan, B. H. Paw, A. Drejer,

- B. Barut, A. Zapata, T. C. Law, C. Brugnara, S. E. Lux, G. S. Pinkus, J. L. Pinkus, P. D. Kingsley, J. Palis, M. D. Fleming, N. C. Andrews, and L. I. Zon (2000) *Nature* **403**, 776–781.
24. S. Abboud and D. J. Haile (2000) *J. Biol. Chem.* **275**, 19906–19912.
  25. C. D. Vulpe, Y. M. Kuo, T. L. Murphy, L. Cowley, C. Askwith, N. Libina, J. Gitschier, and G. J. Anderson (1999) *Nat. Genet.* **21**, 195–199.
  26. J. H. Jandl, J. K. Inman, R. L. Simmons, and D. W. Allen (1959) *J. Clin. Invest.* **38**, 161–185.
  27. P. Aisen and A. Leibman (1973) *Biochim. Biophys. Acta* **304**, 797–804.
  28. J.-M. El Hage Chahine and R. Pakdaman (1995) *Eur. J. Biochem.* **230**, 1102–1110.
  29. J. Dewan, B. Mikami, M. Hirose, and J. C. Sacchettini (1993) *Biochemistry* **32**, 11963–11968.
  30. L. M. Steinlein, C. M. Ligman, S. Kessler, and R. A. Ikeda (1998) *Biochemistry* **37**, 13696–13703.
  31. R. Pakdaman, F. B. Abdallah, and J. M. E. Chahine (1999) *J. Mol. Biol.* **293**, 1273–1284.
  32. F. B. Abdallah and J. M. E. Chahine (2000) *J. Mol. Biol.* **303**, 255–266.
  33. P. K. Bali, O. Zak, P. Aisen (1991) *Biochemistry* **30**, 324–328.
  34. S. A. Kretchmar and K. N. Raymond (1988) *Inorg. Chem.* **27**, 1436–1441.
  35. J. A. Lebrón, A. G. West Jr., and P. J. Bjorkman (1999) *J. Mol. Biol.* **294**, 239–245.
  36. C. N. Roy, D. M. Penny, J. N. Feder, and C. A. Enns (1999) *J. Biol. Chem.* **274**, 9022–9028.
  37. N. C. Andrews (1999) *Int. J. Biochem. Cell Biol.* **31**, 991–994.
  38. A. T. McKie, D. Barrow, G. O. Latunde-Dada, A. Rolfs, G. Sager, E. Mudaly, M. Mudaly, C. Richardson, D. Barlow, A. Bomford, T. J. Peters, K. B. Raja, S. Shirali, M. A. Hediger, F. Farzaneh, and R. J. Simpson (2001) *Science* **291**, 1755–1759.

## Transferrin Receptor

Most cells acquire iron from [transferrin](#) by a **receptor**-mediated process, and the transferrin receptor controls iron uptake. The transferrin receptor is a glycosylated homodimer [membrane protein](#) of 90-kDa subunits, linked by two intramembranous [disulfide bonds](#) (1). The receptor spans the plasma [membrane](#) once (Fig. 1), with the transmembrane portion of the receptor serving as its own **signal sequence** for membrane insertion (2). Insertion is stabilized by covalently bound [fatty acid](#) in or near the intramembranous portion of the protein (1). Most of the receptor, including the transferrin-recognition sequence in its C-terminal region (3), lies outside the membrane. The extracellular portion of the receptor contains three sites for [N-glycosylation](#) and one for [O-glycosylation](#); **glycosylation** is critical in the proper folding and function of the molecule (4, 5). A soluble 85-kDa form of the receptor, resulting from cleavage at an extracellular site adjacent to the region of the membrane insertion site, is found in the circulation and correlates with iron deficiency and enhanced erythropoiesis (6).

**Figure 1.** Model of the transferrin receptor showing post-translational modification sites. (Figure generously provided by Dr. Caroline Enns.)



## 1. Cell Biology

Distribution of the transferrin receptor between plasma membrane and intracellular locations is variable among different cell types, with as much as 75% located within cells and therefore inaccessible to extracellular transferrin (7, 8). [Endocytosis](#) of receptor is triggered by binding of transferrin, but the receptor is continually shuttled between plasma membrane and intracellular sites even in the absence of its ligand (9). **Phosphorylation** of the receptor enhances its rate of internalization, without altering its cellular distribution (10). [Microtubule](#) inhibitors retard internalization but exert little effect on [exocytosis](#) (11, 12); conversely, inhibitors of phosphoinositide [kinase](#) inhibit exocytosis without affecting endocytosis (13). A tetrapeptide tight-turn sequence, -Tyr-X-Arg-Phe- (or YXRF), within the cytoplasmic domain of the receptor has been incriminated as the recognition signal for high efficiency endocytosis (14).

## 2. Regulation of Transferrin Receptor Gene Expression

Expression of the transferrin receptor gene is enhanced by iron demand or deficiency and suppressed, even to the point of nondetectability, by iron overload (15). Regulation of expression is primarily at the stage of [translation](#). Hairpin **stem-loop structures** in the 3'-untranslated region of the receptor [messenger RNA](#), termed iron-responsive elements, (IRE) bind [trans-acting](#) proteins known as iron regulatory proteins (IRPs) (formerly designated IRE-binding proteins). Five IRE motifs are found in the receptor mRNA. Occupancy of these regulatory sequences by IRPs stabilizes and substantially prolongs the lifetime of the mRNA. Two classes of IRP have been identified. The first of these, IRP-1, displays high homology with mitochondrial *aconitase*, an **iron-sulfur** enzyme with a 4Fe-4S cluster. When depleted of iron by cellular iron deficiency, iron chelators, or reducing agents, IRP-1 binds avidly to the IREs, protecting their mRNA from endonucleolytic attack. Conversely, IRP-1 loses IRE-binding activity when loaded with iron, with consequent loss of such protection. In contrast, IRP-2 is not definitely known to bind iron; instead, iron effects a threefold increase in its degradation rate by [proteasomes](#) (16) as well as increased *de novo* synthesis (17). The IRE targets of the two IRP may not be identical (17, 18), but each responds to the iron status of the



cell, increasing receptor mRNA levels in iron deficiency and decreasing mRNA in iron-replete states.

### 3. Receptor-Dependent Uptake of Iron from Transferrin by Cells

At the pH of the extracellular space, 7.4, where transferrin and the transferrin receptor meet, iron-bearing transferrin is preferentially bound by receptor (19). The effect of this preference is to keep the cell from useless internalization of transferrin molecules devoid of iron, which predominate in iron-deficient and normal states, when transferrin saturation is 30% or less. Binding of transferrin to its receptor is rapidly followed by internalization of the transferrin–transferrin receptor complex to a membrane-bounded **clathrin**-coated vesicle, which soon matures to an ATP-driven proton-pumping **endosome**. As the pH is lowered to about 5.6 (the final value varies in different cell types), iron is lost from transferrin, a reaction driven in part by the receptor, which now binds apotransferrin more tightly than iron-bearing transferrin. In addition to the lowered pH of the endosome and the accelerating effect of receptor on iron release (20), an iron-sequestering agent is required to effect release; this requirement may be met by physiological iron chelators such as ATP, citrate, and oxalate.

Released iron must now escape the endosome; details of how this is accomplished remain elusive. One attractive suggestion is that reduction of released iron to the ferrous state enables it to traverse the endosomal membrane (21). In **yeast**, reduction of iron is an essential step for uptake (22), and ferrous iron is capable of crossing reticulocyte membranes for incorporation into **hemoglobin** (23). The endosome bearing iron-depleted transferrin still bound to its receptor is now returned to the cell surface, where fusion with the plasma membrane again exposes transferrin to the extracellular pH, 7.4. At this pH, apotransferrin dissociates from receptor to reenter the circulation for another cycle of iron transport ( (24, 25)

### 4. Receptor-Independent Pathways of Iron Uptake from Transferrin

In addition to the receptor-mediated pathway for iron uptake from transferrin, at least two other routes are available for acquisition of transferrin-borne iron by cells. The hepatocyte may secure as much or more iron from transferrin by a receptor-independent pathway as by the receptor-mediated process (26, 27). The ultimate fate of iron acquired by the receptor-independent path is indistinguishable from that of iron taken up by the receptor-dependent mechanism, but the details of receptor-independent are still to be elucidated.

Melanoma cells can manage uptake of iron from transferrin by yet another mechanism. These malignant cells express a membrane transferrin appropriately named *melanotransferrin*, which can function as a shuttle for transferring iron from transferrin to the cell interiors (28). The mechanism of iron exchange between proteins, and the physiological importance of melanotransferrin as a vehicle for iron uptake, like so many other intricacies of iron metabolism, are not known.

### 5. Acknowledgment

Preparation of this manuscript was supported, in part, by Grant DK15056 from the National Institutes of Health, U.S. Public Health Service.

### Bibliography

1. S. Jing and I. S. Trowbridge (1987) *EMBO J.* **6**, 327–331.
2. M. Zerial, P. Melancon, C. Schneider, and H. Garoff (1997) *EMBO J.* **5**, 1543–1550.
3. F. Buchegger, I. S. Trowbridge, L. F. Liu, S. White, and J. F. Collawn (1996) *Eur. J. Biochem.* **235**, 9–17.
4. G. R. Hayes, C. A. Enns, and J. J. Lucas (1992) *Glycobiology* **2**, 355–359.

5. G. R. Hayes, A. Williams, C. E. Costello, C. A. Enns, and J. J. Lucas (1995) *Glycobiology* **5**, 227–232.
6. J. D. Cook, R. D. Baynes, and B. S. Skikne (1994) *Adv. Exp. Med. Biol.* **352**, 119–126.
7. J. L. Frazier, J. H. Caskey, M. Yoffe, and P. A. Seligman (1982) *J. Clin. Invest.* **69**, 853–865.
8. J. Sainte-Marie, M. Vidal, P. Bette-Bobillo, J. R. Philippot, and A. Bienvenüe (1991) *Eur. J. Biochem.* **201**, 295–302.
9. N. Gironès and R. J. Davis (1989) *Biochem. J.* **264**, 35–46.
10. T. Eichholtz, P. Vossebeld, M. Van Overveld, and H. Ploegh (1992) *J. Biol. Chem.* **267**, 22490–22495.
11. M. Jin and M. D. Snider (1993) *J. Biol. Chem.* **268**, 18390–18397.
12. H. S. Thatte, K. R. Bridges, and D. E. Golan (1994) *J. Cell. Physiol.* **160**, 345–357.
13. P. R. Shepherd, M. A. Soos, and K. Siddle (1995) *Biochem. Biophys. Res. Commun.* **211**, 535–539.
14. J. F. Collawn, M. Stangel, L. A. Kuhn, V. Esekogwu, S. Q. Jing, I. S. Trowbridge, and J. A. Tainer (1990) *Cell* **63**, 1061–1072.
15. M. Lombard, A. Bomford, M. Hynes, N. V. Naoumov, S. Roberts, J. Crowe, and R. Williams (1989) *Hepatology* **9**, 1–5.
16. B. Guo, J. D. Phillips, Y. Yu, and E. A. Leibold (1995) *J. Biol. Chem.* **270**, 21645–21651.
17. B. R. Henderson and L. C. Kühn (1995) *J. Biol. Chem.* **270**, 20509–20515.
18. E. R. Henderson, E. Menotti, and L. C. Kühn (1996) *J. Biol. Chem.* **271**, 4900–4908.
19. S. P. Young, A. Bomford, and R. Williams (1984) *Biochem. J.* **219**, 505–510.
20. P. K. Bali, O. Zak, and P. Aisen (1991) *Biochemistry* **30**, 324–328.
21. J. A. Watkins, J. D. Altazan, P. Elder, C.-Y. Li, M.-T. Nunez, X.-X. Cui, and J. Glass (1992) *Biochemistry* **31**, 5820–5830.
22. G. J. Anderson, A. Dancis, D. G. Roman, and R. D. Klausner (1994) *Adv. Exp. Med. Biol.* **356**, 81–90.
23. A. Egyed (1988) *Br. J. Haematol.* **68**, 483–486.
24. R. D. Klausner, J. V. Ashwell, J. B. VanRenswoude, J. Harford, and K. Bridges (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2263–2266.
25. A. Dautry-Varsat, A. Ciechanover, and H. F. Lodish (1983) *Proc. Natl. Acad. Sci. USA* **80**, 2258–2262.
26. D. Trinder, E. H. Morgan, and E. Baker (1988) *Biochim. Biophys. Acta* **943**, 440–446.
27. D. Trinder, O. Zak, and P. Aisen (1996) *Hepatology* **23**, 1512–1520.
28. D. R. Richardson and E. Baker (1990) *Biochim. Biophys. Acta Mol. Cell Res.* **1053**, 1–12.

### **Suggestion for Further Reading**

29. M. W. Hentze and L. C. Kühn (1996) Molecular control of vertebrate iron metabolism: mRNA-based regulatory circuits operated by iron, nitric oxide, and oxidative stress, *Proc. Natl. Acad. Sci. USA* **93**, 8175–8182.

### **Transformation**

Many bacteria can bind naked **DNA** to their surface, transport it through the cell wall, and

incorporate it into their [genomes](#). This phenomenon is called *transformation*. Typically, transforming DNA substitutes for a **homologous** segment of the host [chromosome](#) by [recombination](#), although **plasmids** usually resume their autonomous existence. The term *transformation* originally referred only to the alteration in behavior and appearance of one pneumococcal strain upon mixing with heat-killed cells of another (1), but it has since come to signify the entire DNA-uptake process.

To import DNA, bacteria must be in a particular physiological state, usually transitory, termed [competence](#). Only then are the **genes** that code for the DNA uptake apparatus turned on. Competence occurs naturally in many species, but far from all. However, the capacity for DNA uptake can often be induced by laboratory manipulation of bacteria that do not display natural competence. Two examples of competence, one natural and one artificial, transformed biology itself: Natural transformation of *Pneumococcus* (now *Streptococcus pneumoniae*) was the experimental procedure used by Avery and his colleagues to demonstrate that DNA is the genetic material (2), and calcium treatment of normally noncompetent *Escherichia coli* was crucial to the emergence of [recombinant DNA](#) technology (3, 4). Our intention here is less ambitious: to show how transformation takes place and to consider its function in nature.

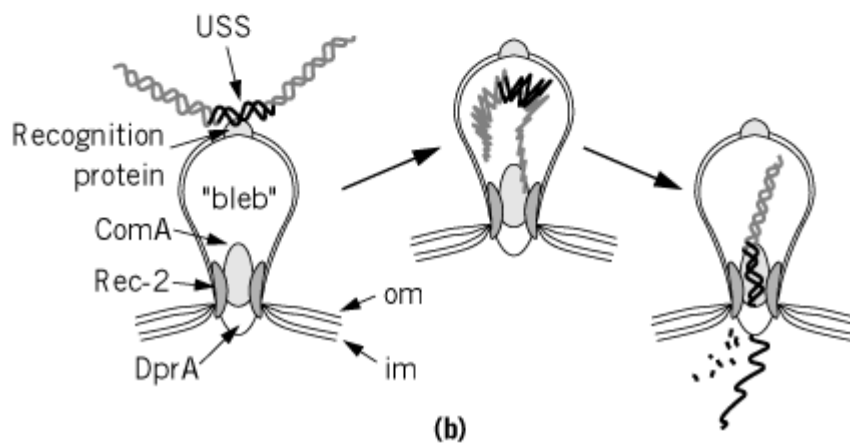
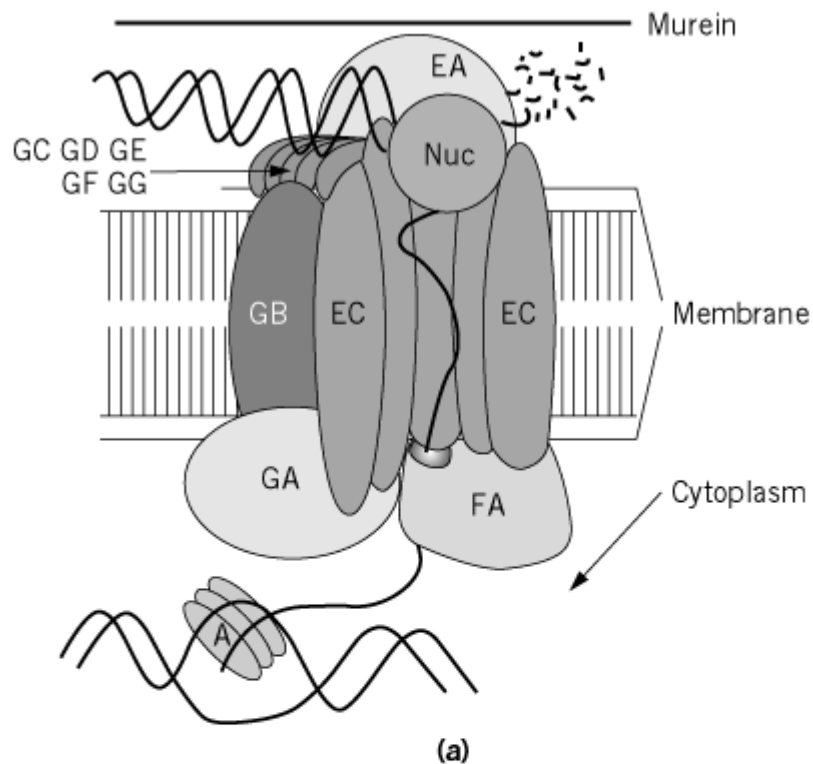
## 1. The Basic Process

Nearly everything we know about transformation and competence comes from experiments carried out on four species: *S. pneumoniae*, *Bacillus subtilis*, *Haemophilus influenzae*, and *Neisseria gonorrhoeae*. These studies suggest that whereas induction of competence is idiosyncratic, possibly reflecting the diversity of natural habitats in which naked DNA is encountered (see [Competence](#)), the basic mechanisms of transformation are similar in all naturally transformable bacteria. Thus we can view the process of transformation as passing through four stages: (i) Competent cells sporting a limited number of DNA-binding receptors attach duplex DNA to the cell surface. (ii) The DNA is processed by endonucleolytic cleavage and, usually, by exonucleolytic removal of one of the strands. (iii) The other strand is drawn through a pore into the cytoplasm, possibly in concert with the exonucleolytic degradation. (iv) The entering single strand of DNA undergoes homologous recombination with the chromosome, creating a new recombinant if the recombining partners carry different **alleles**. We shall look at how each of the four bacterial species works its variations on this theme.

### 1.1. *B. Subtilis*

As cells become competent, [transcription](#) of the “late” competence genes begins. One of these transcription units, *comG*, specifies five small proteins (ComGC to GG) that appear to form external **receptors** for DNA, as well as two larger proteins (ComGA, GB) that probably anchor the receptors to the membrane (Fig. 1). Each competent cell has about 50 of these receptors, as measured by saturation of binding capacity for radioactively labeled DNA (5, 6). When transforming DNA is added, the receptors bind it rapidly and without regard to sequence (unlike *H. influenzae* and *N. gonorrhoeae*; see text below), and for about 30 seconds the DNA remains on the surface, susceptible to added **nucleases** but not removable by washing.

**Figure 1.** The apparatus of DNA uptake during transformation of (a) *B. subtilis* (modified from Ref. 8, with permission of the publisher) and (b) *H. influenzae*. The locations and presumed roles of several of the proteins shown are partly speculative, but the overall mechanisms depicted are based on good evidence. The *B. subtilis* Com proteins are described in the text; Nuc denotes nucleases that cut bound DNA and degrade the 5'-ended strand; A shows RecA protein, mediating invasion of the resident chromosome. The appearance of the *H. influenzae* bleb is based on the electron micrograph shown in Ref. 34. The pore location of the *H. influenzae* proteins is based on their presumed function (see text).



Processing of the DNA then begins, initiated by double-strand cleavage. The length distribution of the resulting fragments is consistent with the DNA being cut once at each of the sites to which it is bound (7). This cleavage provides the free ends used for transport across the membrane. Within 1–2 minutes, one strand of each molecule becomes immune to added nuclease, reflecting transport, while the other is degraded. Of the known competence functions, the ComE proteins are the most intimately involved in transport. ComEC is needed only for transport, not for DNA binding, and its properties suggest that it may form a transmembrane pore through which the DNA passes (8). ComEA appears to be involved in both binding and transport: Its C-terminus, which projects from the membrane to the exterior, carries **DNA-binding** motifs, suggesting a role in processing of DNA (9). However, neither ComEA nor any of the other proteins has been proven to exhibit the nuclease activities associated with processing. The 17-kDa ComI protein, which forms part of a complex obtained from competent cells, was found to have endonuclease activity (10), but a *comI* null mutant is fully transformable (11). Perhaps there is more than one transformation endonuclease.

The entering single strand appears to be driven into the cytoplasm by the inner-membrane protein, ComFA, using the energy of ATP hydrolysis to translocate the DNA, rather like the PriA protein of

*E. coli* that it resembles (12, 13). ComFA mutations that abolish ATP binding diminish rather than eliminate transformation, suggesting that the DNA can sometimes find its way in by simple diffusion.

Upon entering the cell, the single-stranded DNA encounters a recombination system primed by induction of competence: an increased frequency of single-strand gaps in the recipient chromosome (14) and major increases in the concentrations of **RecA** and the nuclease/**DNA helicase**, AddAB (15, 16). Incorporation of DNA into the homologous chromosomal region is thus very efficient; about 70% of it is integrated (11), and the overall probability of integration is greater than 0.5 per molecule initially bound at the cell surface. Although sequence differences between the added and host DNA cause the transformation frequency to drop (17), apparently because of reduced duplex stability, the DNA metabolism of *B. subtilis* is otherwise well-adapted for incorporation of new alleles. **Restriction-modification systems** do not seem to limit transformation in interspecific crosses (18), and **mismatch repair** following transformation has not been reported, although genes analogous to the *hex* mismatch correction system of *S. pneumoniae* (see text below) appear in the *B. subtilis* genome sequence.

### 1.2. *S. Pneumoniae*

The individual components of the DNA binding and uptake process are not as well known as for *B. subtilis*, but the process appears to be very similar. There are 30–80 duplex-DNA binding sites per cell, at which DNA is cut to initiate transport of a single strand into the cytoplasm. In this bacterium, however, more is known of the nuclease activities. A specific membrane-bound endonuclease, EndA, introduces double-strand breaks into the DNA, thus providing free ends for entry (19, 20). Although introduction of an *endA* null mutation reduces transformation frequency over a thousandfold, transformants are still obtained in such an *endA* strain (21), possibly because a single-strand nicking activity that accompanies binding can also create free ends (19). The mere presence of ends is insufficient, however; ongoing transport is curtailed when the mechanism encounters a nick in the entering strand (22). A particular geometry of the DNA end and the transport initiator seems to be needed.

Unlike *B. subtilis*, *S. pneumoniae* transports DNA in an oriented manner. The entering strand goes in 3' first, while the other strand is degraded from its 5' end to acid-soluble oligonucleotides (22, 23). The identical kinetics of uptake and degradation suggest that these processes are concerted (22), although there is as yet no direct evidence for such coupling.

Other proteins involved in DNA processing or the assembly of transmembrane pores have not been identified, although it is possible that a 19.5-kDa peptide found associated with the single strand in the cytoplasm is one such example (24), because by complexing with the DNA it confers on it resistance to nucleases. Evidence that DNA uptake is energy-dependent has been presented, with the suggestion that an ATP-driven import protein translocates the entering single strand (25).

The DNA is released from the nuclease-resistant complex to recombine with the homologous chromosomal region via RecA-mediated strand invasion. A function that probably catalyzes branch migration, MmsA, is also needed to complete transformation (26). Genetic transformation must run the gauntlet of the Hex repair system, which recognizes and removes mismatched bases, reducing transformation frequencies (27). The Hex system is of limited capacity, however, being readily overwhelmed by related but divergent DNA (28). The way is therefore open to recombination with DNA of other *Streptococci*, with obvious evolutionary consequences.

### 1.3. *H. Influenzae*

Initial DNA binding by this **Gram-negative** bacterium contrasts sharply with that of *B. subtilis* and *S. pneumoniae*, in that only the DNA of *H. influenzae* itself can participate (29). The binding apparatus specifically recognizes a 9-bp sequence (5'-AAGTGCGGT), termed an *uptake signal sequence* (USS), within an extended consensus region of 29 bp that is rich in adenine and thymine. *H. influenzae* contains 1465 of these USS's, nearly 200 times more than expected by chance, and they

are distributed more or less randomly within its genome (30, 31), accounting for the species-specific binding observed. A further difference is the appearance on the surface of competent cells of about four to eight physically discrete membranous vesicles, termed *transformasomes* (or “blebs”), which bind DNA and protect it from nucleases (Fig. 1b) (32-34). The blebs are necessary for transformation; certain transformation-deficient mutants form the blebs, but shed them prematurely (32, 35).

Transformation in *H. influenzae* begins when a transformasome packages a double-stranded USS-containing DNA molecule. Both linear and circular (or closed-end) DNA is readily taken up by transformasomes, but only linear DNA is conducted out again (36); there appears to be no endonuclease to create ends in bound DNA. The transformasome extrudes DNA with a free end through a pore into the cytoplasm; [electron microscopic](#) images suggest that the two cell-wall membranes are fused where the bleb attaches (34). As with *S. pneumoniae*, the 5'-ended strand is degraded, while the 3'-ended strand enters the cell (36), apparently by a translocation process that involves the products of the following genes: *rec-2*, which is probably a pore protein (37), *dprA*, which is possibly a translocase (38), and *comA*. Unlike *S. pneumoniae*, however, the entering strand is not protected and suffers partial degradation, prior to being incorporated by recombination into the chromosome (36). Rec-1 protein mediates the recombination, with a recipient chromosome rendered 5% single-stranded during competence development (39).

Although the catalogue of proteins that participate in *H. influenzae* transformation is far from complete, the availability of the entire *H. influenzae* sequence (40) will accelerate identification of proteins with functions similar to those in other transformable species.

#### 1.4. *N. Gonorrhoeae*

This gram-negative species inhabits the human genital tract, attaching itself to the epithelia by filamentous [pili](#). It is another bacterium that selects for transformation by its own DNA, or that of close relatives, through USS recognition. In this case, the USS is the 10-bp sequence 5'-GCCGTCTGAA, of which there are several hundred examples in the genome (41, 42).

DNA uptake depends strictly on the major component of the pilus, PilE, as well as a protein involved in pilus assembly, PilC (43-45), but not on the pilus itself. Nonpiliated mutants can still be transformed (46), as can strains that make normal PilE in quantities insufficient for pilus formation (45). *Pil C*<sup>-</sup> strains that have reacquired the ability to make pili are normally not transformable, but become so upon addition of a PilC protein fraction to the cells, indicating that PilC and PilE act at the cell exterior. Nonpiliated *pilC*<sup>-</sup> strains cannot be so complemented, suggesting that PilC and PilE act as a complex to promote DNA uptake (45). They do not, however, specifically recognize DNA containing the *Neisseria* USS (47); the protein that does so remains to be identified. Deficiencies in other proteins involved in pilus formation, such as PilT, can impair transformation (48), leading to the suggestion that these proteins also mediate the reverse reaction, pilus depolymerization, thus pulling bound DNA across the outer membrane (47).

Bound DNA is subjected to limited double-strand cleavage (48), but enters the **periplasm** otherwise intact. The available data (49) suggest that it then crosses the murein layer with the aid of two proteins involved in cell shape and division, ComL and Tpc, and finally passes through the cytoplasmic membrane aided by ComA, a protein with similarities to the pore-forming proteins of *B. subtilis* and *H. influenzae*. Uniquely among the organisms studied, *Neisseria* does not generate single-stranded DNA during uptake. Even the further processing that occurs in the cytoplasm appears to consist of double-strand exonucleolysis (49). Strand separation is presumably concerted with the other steps of homologous recombination mediated by RecA (50). Nevertheless, efficient transformation with single-stranded DNA has been observed (51), indicating a degree of versatility in the transport machinery.

Transformation of *Neisseria* has revealed some striking relationships. First is the intimacy of the

connection between DNA transformation and the pili. The pilus is the point of both attachment to the host and attack by it; it expands the former and limits the latter by binding related DNA containing alternative pilus gene alleles. Second is the economy of the 10-bp USS's, which serve not only as the recognition elements for DNA binding but also, in the form of pairs of inverted sequences, as transcriptional terminators (41). [Although the examples of *H. influenzae* and *N. gonorrhoeae* had suggested that sequence-specific binding was a feature of all gram-negative bacteria, this is not the case; *Acinetobacter calcoaceticus* binds any DNA (52).] Third is the extensive homology of the PilE protein and proteins involved in pilus assembly with a large family of proteins responsible for transporting macromolecules across membranes: the DNA-binding proteins of competent *B. subtilis* and *H. influenzae*, the type IV pilus proteins of a variety of pathogenic bacteria, and proteins involved in filamentous phage assembly, DNA transfer during conjugation in *E. coli* and *Agrobacterium tumefaciens*, and protein secretion (53).

## 2. Artificial Transformation

The discovery that *E. coli* could be induced to take up DNA by treatment with  $\text{Ca}^{2+}$  ions (3) led to the realization that transformation need not be limited to species that attain competence naturally. The need for genetic manipulation of an expanding range of bacterial species motivated the search for new methods of transformation. These fall into three general classes: chemical treatment, protoplasting, and electroporation. These methods nearly always have as their goal the uptake of intact double-stranded DNA, and the physical mechanisms, insofar as they are known, appear to bear little relation to those of natural transformation.

1. Chemical treatment is typified by the suspension of *E. coli* cells in cold  $\text{CaCl}_2$  solution, and the ensuing membrane changes that might be responsible for allowing DNA uptake are discussed in the entry [Competence](#). This type of method has been most successful with gram-negative bacteria, especially when variations such as the use of other divalent cations or the addition of a freeze-thaw step, which presumably makes the membrane fragile, are included (54, 55). Nevertheless, addition of **polyethylene glycol** to the cell-DNA mixture has permitted transformation of certain gram-positive species (56).
2. Protoplast transformation involves mild digestion of the murein layer with lysozyme, incubation of the resulting protoplasts with DNA in the presence of polyethylene glycol, and regeneration of the treated cells to restore the cell wall, followed by selection of transformants (57). Although more laborious than other methods, protoplast transformation has proven particularly useful for strains that cannot be rendered competent otherwise. Most of these are gram-positive bacteria, such as *Streptomyces*, *Corynebacterium*, and *Lactobacillus*, suggesting that the main barrier to DNA entry is the thick murein layer, which is removed during formation of protoplasts. For some species, deoxyribonucleases associated with the protoplasts hinder transformation, but this problem has been circumvented in certain cases by enclosing the DNA in phospholipid vesicles called liposomes (58). These presumably fuse with protoplast membranes to deliver the DNA, rather like the blebs of competent *H. influenzae*.
3. Electroporation works by delivering a brief pulse of current to cells in the presence of DNA, producing temporary holes in the membranes that allow DNA uptake (59, 60). It is fast and simple and is probably the method of choice for a bacterium of unknown transformability. It does not work with every species, however: Some (eg, *S. pneumoniae*), although transformable naturally, are not transformed by electroporation, because only cells in which competence has developed contain adequate levels of recombination enzymes. For this method, as for the others, the mechanisms underlying DNA uptake are not fully understood.

## 3. Biological Significance

Until recently, it was not unreasonable to question whether transformability is biologically important

or just a by-product of the altered metabolism with which it is associated. It has, however, become increasingly difficult to sustain this skepticism. Perhaps the strongest argument for the biological significance of transformation is the existence of mosaic genes in transformable species isolated from nature after exposure to antibiotics or host defenses. This phenomenon implies a role for transformation in generating the variability that has allowed fitter populations to emerge from threatened ones. Nevertheless, other functions for uptake of DNA have been advanced.

### 3.1. Restoration?

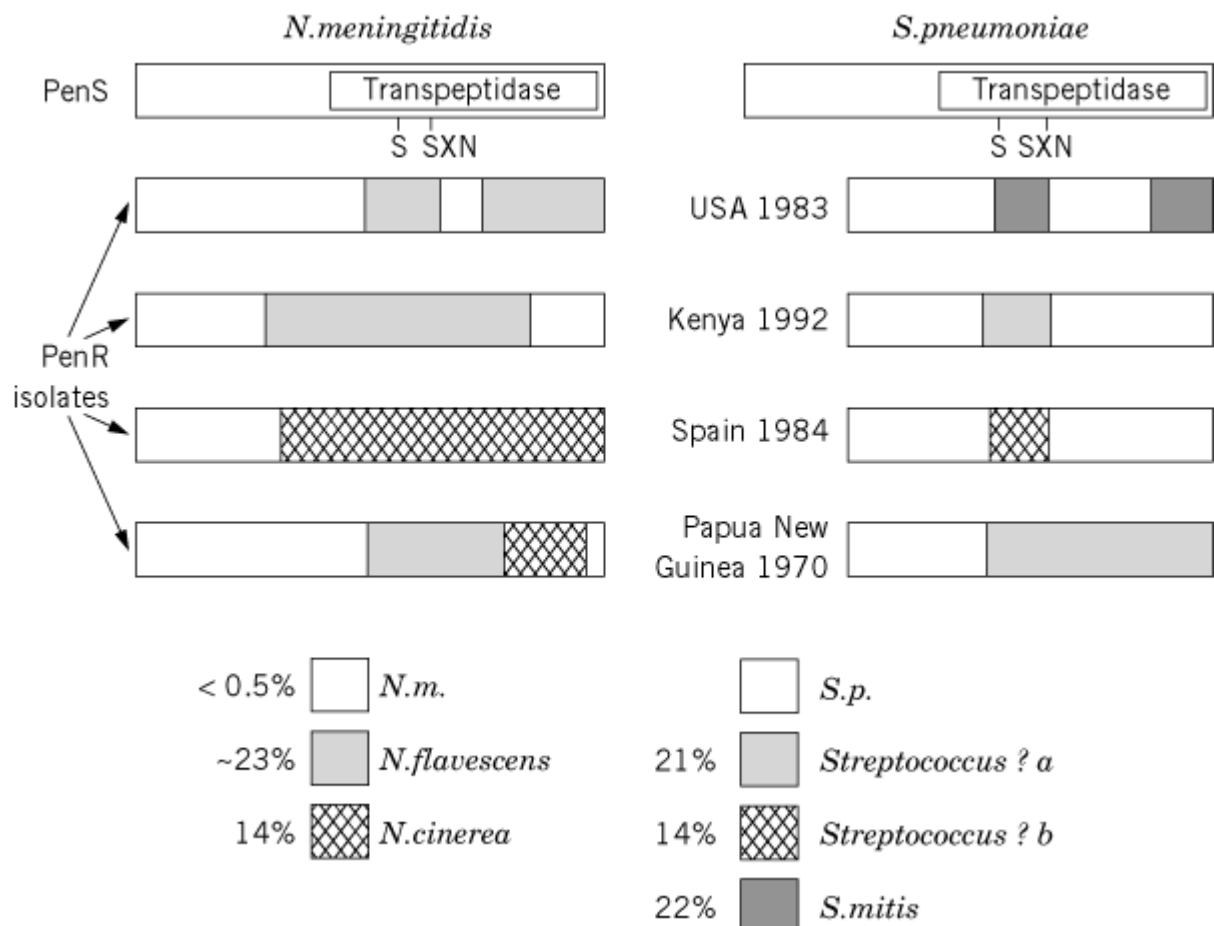
The induction of competence by nutrient deprivation suggested that transformation might be a form of scavenging, in which cells increased their chances of survival by eating DNA released from their dead cohabitants (61, 62). There is, however, no direct evidence that transformed cells use DNA in this way. Moreover, the wastage of one strand by degradation during transport, the lack of extensive degradation of the incoming strand, and the elaborate nature of the uptake apparatus render a general nutritional role unlikely. Another suggested function for transformation is the use of the imported DNA as a template to repair damaged chromosomes: Addition of homologous DNA increased the survival of UV-irradiated *B. subtilis* and *H. influenzae* cells (63, 64). However, in *H. influenzae* a sequence constituting only 1% of the chromosome was just as effective as the whole genome, indicating that transforming DNA enhances survival indirectly, rather than by simply providing undamaged copies (64). Furthermore, DNA damage failed to enhance competence in either species (65). Thus, although transformation might serve in nature as one means of correcting errors, [DNA repair](#) does not appear to be its major role.

### 3.2. Mosaic Genes

The clearest indication that transformation might increase the fitness of bacterial populations comes from the analysis of **penicillin**-resistant isolates of *Neisseria* and *Streptococcus* (66). Penicillin acts as a substrate analogue, binding to the transpeptidases that cross-link murein chains. Resistance arises by alterations in these penicillin-binding proteins (PBPs) that reduce their affinity for penicillin. Comparison of the PBP2 genes in various *Neisseria* species revealed that those in penicillin-sensitive isolates of each species have a very uniform sequence, whereas those in resistant isolates of the same species have a composite structure (Fig. 2). The resistant PBP2s of *N. gonorrhoeae* and *N. meningitidis* consist of regions nearly identical to those in sensitive strains, interrupted by islands of divergent sequence corresponding to the PBP2 genes of *N. flavescens* and *N. cinerea* (67). Similar gene structures are found in the PBP2B genes of diverse isolates of *S. pneumoniae*, where it has also been possible in some cases to trace the origin of the divergent regions (68, 69). It is unlikely that these mosaic genes result from successive point mutations. The subtle changes in [protein structure](#) needed to reduce affinity for penicillin, while continuing to allow recognition of substrate, can in general only be achieved through multiple amino acid changes. Because the probability that any single one of these changes provides selective advantage is low, all or most of the changes are most likely to have been introduced simultaneously. Only recombination can bring this about. Because **phage**-mediated and **plasmid** mobilization-mediated mobilization of chromosomal DNA are unknown in *Neisseria* and *Streptococcus*, transformation by DNA released from related species almost certainly generates the resistant PBPs.

**Figure 2.** Mosaic composition of the PBP2 and PBP2B genes in penicillin-resistant isolates of *N. meningitidis* and *S. pneumoniae*, respectively. The figure is an abbreviated composite of Figures 1 and 2 in Ref. 66. The domain responsible for the transpeptidase activity of these homologous proteins is shown, together with the active site serine (S) and a shared sequence motif, -Ser- X-Asn- (SXN), that is probably important for activity. The percentages indicate the degree of sequence dissimilarity of each divergent region with the equivalent sequence in the penicillin-sensitive strains. For *S. pneumoniae*, only the sequenced portion (~70%) is shown. The source of two of the divergent PBP2B regions is unknown (?a and?b). (Modified from Ref. 66, with permission from the author.)





Alteration of the major *N. gonorrhoeae* pilus protein, Pile, provided the first evidence for the involvement of transformation in the creation of mosaic genes (70). The pili are a target of immune attack, but frequent switching of the expressed *pilE* allele helps the bacterium evade it. This switching occurs in cultured cells, either by recombination with endogenous silent *pil* genes or by transformation with DNA from cells expressing a different allele. Because *Neisseria* are constitutively competent and continually release DNA by autolysis, it is likely that recombination with transforming DNA is also a major source of pilus switching *in vivo*. Evidence for mosaic gene formation in nature comes from the structure of IgA1 proteinase (*iga*) genes, virulence determinants thought to act by cleaving the antibody and using the Fab fragment to mask the pilus. Genes from independent *N. gonorrhoeae* isolates contained a pattern of related segments (71). A more extensive analysis has uncovered evidence of interspecies recombination between *iga* genes of *N. gonorrhoeae* and *N. meningitidis*, as well as indications of composite gene formation in *H. influenzae* (72).

### 3.3. Microcosms

Transformation has often been demonstrated in pseudonatural habitats ("microcosms"). *H. influenzae* cells contained in diffusion chambers and implanted in rat peritoneum quickly acquired competence and were transformed with coimplanted DNA (73); differently marked *B. subtilis* strains inoculated into sterile soil were replaced by populations of cells carrying markers from both parents (74); and transformation of *Pseudomonas stutzeri* has been observed in a sterile marine water and sediment environment (75). Even with this evidence that transformation can play an important part in adapting natural populations to environmental challenge, we cannot yet assess its extent in nature. Future studies aimed at determining how transformation contributes to the evolution of natural populations and the virulence of pathogenic bacteria will be of particular interest. Indications that naked DNA can be fed to an animal with whose genome it shares no sequence homology, and yet become integrated in the chromosomes of several tissues of the animal (76), provocatively suggest

that we might have underestimated the wider significance of transformation.

## Bibliography

1. F. Griffith (1928) *J. Hygiene* **27**, 113–159.
2. O. T. Avery, C. M. MacLeod, and M. McCarty (1944) *J. Expl. Med.* **79**, 137–159.
3. M. Mandel and A. Higa (1970) *J. Mol. Biol.* **53**, 159–162.
4. S. N. Cohen, A. C. Y. Chang, H. W. Boyer, and R. B. Helling (1973) *Proc. Natl. Acad. Sci. USA* **70**, 3240–3244.
5. R. M. Singh (1972) *J. Bacteriol.* **110**, 266–272.
6. D. Dubnau and C. Cirigliano (1972) *J. Mol. Biol.* **64**, 31–46.
7. D. Dubnau (1976) In *Microbiology 1976* (D. Schlessinger, ed.), ASM, Washington D.C., pp. 14–27.
8. D. Dubnau (1997) *Gene* **192**, 191–198.
9. G. S. Inamine and D. Dubnau (1995) *J. Bacteriol.* **177**, 3045–3051.
10. H. Smith, K. Wiersma, G. Venema, and S. Bron (1984) *J. Bacteriol.* **157**, 733–738.
11. D. Dubnau (1991) *Microbiol. Rev.* **55**, 395–424.
12. J. A. Londonyo-Vallejo and D. Dubnau (1994) *Mol. Microbiol.* **13**, 197–205.
13. J. A. Londonyo-Vallejo and D. Dubnau (1994) *J. Bacteriol.* **176**, 4642–4645.
14. W. J. Harris and G. C. Barr (1971) *Mol. Gen. Genet.* **113**, 316–330.
15. C. M. Lovett, P. M. Love, and R. E. Yasbin (1989) *J. Bacteriol.* **171**, 2318–2322.
16. B. J. Haijema, L. W. Hamoen, J. Kooistra, G. Venema, and D. van Sinderen (1995) *Mol. Microbiol.* **15**, 203–211.
17. M. S. Roberts and F. M. Cohan (1993) *Genetics* **134**, 401–408.
18. R. M. Harris-Warrick and J. Lederberg (1978) *J. Bacteriol.* **133**, 1237–1245.
19. S. A. Lacks, B. Greenberg, and M. Neuberger (1975) *J. Bacteriol.* **123**, 222–232.
20. A. L. Rosenthal and S. A. Lacks (1980) *J. Mol. Biol.* **141**, 133–146.
21. A. Puyet, B. Greenberg, and S. A. Lacks (1990) *J. Mol. Biol.* **213**, 727–738.
22. V. Méjean and J.-P. Claverys (1993) *J. Biol. Chem.* **268**, 5594–5599.
23. V. Méjean and J.-P. Claverys (1988) *Mol. Gen. Genet.* **213**, 444–418.
24. D. A. Morrison and B. Mannarelli (1979) *J. Bacteriol.* **140**, 655–665.
25. C. Clavè, F. Martin, and M.-C. Trombe (1989) In *Genetic Transformation and Expression* (L. O. Butler et al., eds.), Intercept, Andover, pp. 27–40.
26. B. Martin, G. J. Sharples, O. Humbert, R. G. Lloyd, and J.-P. Claverys (1996) *Mol. Microbiol.* **19**, 1035–1045.
27. J.-P. Claverys and S. A. Lacks (1986) *Microbiol. Rev.* **50**, 133–165.
28. O. Humbert, M. Prudhomme, R. Hakenbeck, C. G. Dowson, and J.-P. Claverys (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9052–9056.
29. K. L. Sisco and H. O. Smith (1979) *Proc. Natl. Acad. Sci. USA* **76**, 972–976.
30. S. H. Goodgal and M. A. Mitchell (1990) *J. Bacteriol.* **172**, 5924–5928.
31. H. O. Smith, J.-F. Tomb, B. A. Dougherty, R. D. Fleischmann, and J. G. Ventner (1995) *Science* **269**, 538–540.
32. M. Kahn, M. Concino, R. Gromkova, and S. H. Goodgal (1979) *Biochem. Biophys. Res. Commun.* **87**, 764–772.
33. R. A. Deich and H. O. Smith (1980) *Mol. Gen. Genet.* **177**, 369–374.
34. M. E. Kahn, P. Barany, and H. O. Smith (1983) *Proc. Natl. Acad. Sci. USA* **80**, 6927–6931.
35. M. F. Concino and S. H. Goodgal (1982) *J. Bacteriol.* **152**, 441–450.

36. F. Barany, M. E. Kahn, and H. O. Smith (1983) *Proc. Natl. Acad. Sci. USA* **80**, 7274–7278.
37. S. W. Clifton, D. McCarthy, and B. A. Roe (1994) *Gene* **146**, 95–100.
38. S. Karudapuram, X. Zhao, and G. J. Barcak (1995) *J. Bacteriol.* **177**, 3235–3240.
39. D. McCarthy and D. M. Kupfer (1987) *J. Bacteriol.* **169**, 565–571.
40. R. D. Fleischmann et al (1995) *Science* **269**, 496–512.
41. S. D. Goodman and J. J. Scocca (1988) *Proc. Natl. Acad. Sci. USA* **85**, 6982–6986.
42. C. Elkins, C. E. Thomas, H. S. Seifert, and P. F. Sparling (1991) *J. Bacteriol.* **173**, 3911–3913.
43. T. F. Meyer, E. Billyard, R. Haas, S. Storzbach, and M. So (1984) *Proc. Natl. Acad. Sci. USA* **81**, 6110–6114.
44. C. P. Gibbs, B.-Y. Reinmann, E. Schultz, A. Kaufmann, R. Haas, and T. F. Meyer (1989) *Nature* **338**, 651–652.
45. T. Rudel, D. Facius, R. Barten, E. Nonnenbacher, and T. F. Meyer (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7986–7990.
46. Q. Y. Zhang, D. DeRyckere, P. Lauer, and J. M. Kooimey (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5366–5370.
47. M. Fussenegger, T. Rudel, R. Barten, R. Ryll, and T. F. Meyer (1997) *Gene* **192**, 125–134.
48. G. D. Biswas, K. L. Burnstein, and P. F. Sparling (1986) *J. Bacteriol.* **168**, 756–761.
49. D. Facius, M. Fussenegger, and T. F. Meyer (1996) *FEMS Microbiol. Lett.* **137**, 159–164.
50. J. M. Kooimey, E. C. Gotschlich, K. Robbins, S. Bergstrum, and J. Swanson (1987) *Genetics* **117**, 391–398.
51. D. C. Stein (1991) *Can. J. Microbiol.* **37**, 345–349.
52. R. Palmen, B. Vosman, P. Buijsman, C. K. Breck, and K. J. Hellingwerf (1993) *J. Gen. Microbiol.* **139**, 295–305.
53. M. Hobbs and J. S. Mattick (1993) *Mol. Microbiol.* **10**, 233–243.
54. J. R. Saunders, A. Docherty, and G. O. Humphreys (1984) *Methods Microbiol.* **17**, 61–95.
55. M. J. Merrick, J. R. Gibbons, and J. R. Postgate (1987) *J. Gen. Microbiol.* **133**, 2053–2057.
56. M. E. Sanders and M. A. Nicholson (1987) *Appl. Environ. Microbiol.* **53**, 1730–1736.
57. D. A. Hopwood (1981) *Annu. Rev. Microbiol.* **35**, 237–272.
58. R. T. Fraley, C. S. Fornari, and S. Kaplan (1979) *Proc. Natl. Acad. Sci. USA* **76**, 3348–3352.
59. S. K. Harlander (1986) In *Streptococcal Genetics* (J. J. Feretti and R. Curtiss III, eds.), ASM, Washington D.C., pp. 229–233.
60. B. M. Chassy and J. L. Flickinger (1987) *FEMS Microbiol. Lett.* **44**, 173–177.
61. G. J. Stewart and C. A. Carlson (1986) *Annu. Rev. Microbiol.* **40**, 211–235.
62. R. J. Redfield (1993) *J. Hered.* **84**, 400–404.
63. M. A. Hoelzer and R. E. Michod (1991) *Genetics* **128**, 215–223.
64. J. A. Mongold (1992) *Genetics* **132**, 893–898.
65. R. J. Redfield (1993) *Genetics* **133**, 755–761.
66. B. G. Spratt (1994) *Science* **264**, 388–393.
67. R. Hakenbeck, H. Ellerbrock, T. Briese, S. Handwerger, and A. Tomasz (1986) *Antimicrob. Agents Chemother.* **30**, 553–558.
68. B. G. Spratt, Q. Y. Zhang, D. M. Jones, A. Hutchison, J. A. Brannigan, and C. G. Dowson (1989) *Proc. Natl. Acad. Sci. USA* **86**, 8988–8992.
69. T. J. Coffey, C. G. Dowson, M. Daniels, J. Zhou, C. Martin, B. G. Spratt, and J. M. Musser (1991) *Mol. Microbiol.* **5**, 2255–2260.
70. H. S. Seifert, R. S. Ajioka, C. Marchal, P. F. Sparling, and M. So (1988) *Nature* **336**, 392–395.
71. R. Halter, J. Pohlner, and T. F. Meyer (1989) *EMBO J.* **8**, 2737–2744.

72. H. Lomholt, K. Poulsen, and M. Kilian (1995) *Mol. Microbiol.* **15**, 495–506.
73. M. Dargis, P. Gourde, D. Beauchamp, B. Foiry, M. Jacques, and F. Malouin (1992) *Infect. Immun.* **60**, 4024–4031.
74. J. B. Graham and C. A. Istock (1979) *Science* **204**, 637–639.
75. G. J. Stewart and C. D. Sinigalliano (1990) *Appl. Environ. Microbiol.* **56**, 1818–1824.
76. R. Schubbert, D. Renz, B. Schmitz, and W. Doefler (1997) *Proc. Natl. Acad. Sci. USA* **94**, 961–966.

### Suggestions for Further Reading

77. M. G. Lorenz and W. Wackernagel (1994) Bacterial gene transfer by natural genetic transformation in the environment. *Microbiol. Rev.* **58**, 563–602. Especially good for information on less-studied species and on transformation in natural environments.
78. B. Dreiseikelmann (1994) Translocation of DNA across bacterial membranes. *Microbiol. Rev.* **58**, 293–316. Highlights the parallels with other DNA transfer systems and homologies among transport proteins.
79. *Gene* **192**, Pt. 1 (1997). Report of a workshop on Type-IV pilins and the protein and DNA traffic they mediate; reviews of transformation of *B. subtilis*, *Neisseria*, and *A. calcoaceticus*.
80. B. Spratt (1994) Resistance to antibiotics mediated by target alterations. *Science* **264**, 388–393. An authoritative and accessible review of mosaic genes in penicillin-resistant isolates.
81. A. Mercenier and B. M. Chassy (1988) Strategies for the development of bacterial transformation systems. *Biochimie* **70**, 503–517. A slightly dated, but still relevant, review of approaches to artificial transformation of bacteria.

## Transformation Of Fungi

Transformation systems for filamentous [fungi](#) came into their own during the 1980s. The first DNA-mediated transformation into a filamentous fungus was demonstrated for *Neurospora crassa* ([1](#)). Since then, procedures have been developed for all of the major fungal classes ([2-8](#)); similar transformation frequencies, stabilities of transformants, and multiple heterologous integration events have been observed for most of the species commonly studied (see Refs. [9](#) and [10](#) for review).

### 1. Getting DNA into Cells

Most fungal transformation techniques involve either the formation of protoplasts or some method to remove cell-wall components. Cells without walls are further incubated with DNA in the presence of calcium ions and **polyethylene glycol**, followed by regeneration of the cell wall and appropriate selection. Protoplasts are usually formed using young cultures and may be stabilized with either sorbitol or sucrose. In certain instances, transformation efficiencies were reported to be increased with a brief **heat shock** step before recovery and plating. The following partial list of fungi have been successfully transformed by protoplasts: *Candida*, *Aspergillus*, *Podospora*, *Ascobolus*, *Coprinus* ([5](#), [11](#), [12](#)). Intact *Neurospora* has been reported to be transformed with lithium acetate ions ([13](#)), but this methodology does not appear to be in general use for other fungi. In a few instances, fungal cells have been reported to be transformed using electroporation techniques, although efficiencies are considerably lower than those using protoplasts ([14](#), [15](#)).

The efficiency of transformation varies with the type of fungal cell being manipulated, the age of the

culture, and the quality of the DNA preparation. In addition, the efficiency varies with the type of transformation that occurs—replicative versus integrative (discussed in greater detail below).

### 1.1. Vectors

Vectors for transformation of filamentous fungi are not that different from those used for transformation of the yeasts. They contain a selectable marker and bacterial sequences for propagation and selection (usually in *E. coli*). They may also contain multiple cloning sites (MCS) for easy directed cloning at unique sites within the vector. In some instances, DNA fragments are included in the vector that permit autonomous replication in the host fungus (16, 17). Nuclear plasmids analogous to 2  $\mu$  in *Saccharomyces* have not been found in other filamentous fungi, and the dearth of naturally occurring plasmids has made constructing replicative vectors in the filamentous fungi difficult. More recently, cosmid vectors have been described for *M. grisea* and *U. maydis* (18).

Recently, a sequence of DNA from *A. nidulans*, *AMA1*, has been found that substantially increases transformation frequencies concomitant with a decrease in stability of the vector and an elevation in copy number (10–30 copies/haploid genome in *A. nidulans*); all properties of [autonomously replicating sequences](#) (ARS). From a practical standpoint, this same sequence confers similar properties on other species (19, 20), but it has been observed that these plasmids are lost during conidiation (N. Coleman, personal communication). In another series of experiments, a portion of the 2  $\mu$  of *Saccharomyces* was shown both to stabilize and to decrease the copy number of two *Hansenula polymorpha* plasmids (21). The mechanism by which this occurs is thought to be through mitotic partitioning of proteins bound to the 62-bp repeated sequence of the STB locus on 2  $\mu$ . Whether or not these effects can be extended to other model fungal systems remains to be determined.

## 2. Selectable Markers

For many of the fungi, auxotrophic mutants exist that permit the isolation of wild-type genes using a simple selection scheme (ability to grow on medium lacking a particular amino acid or nutrient). For example, mutants in orotidylic acid pyrophosphorylase (OMPppase) have been isolated from nearly every microorganism in which they have been screened for; and either the host gene has been isolated and cloned, or sufficient functional homology exists between the OMPppase genes of different organisms that one OMPppase has been used successfully to transform another species. For example, the first successful transformation in *A. nidulans* used the *pyr-4* gene of *N. crassa* (11). In fact, this gene is so conserved that some fungal genes have been isolated by complementing the defects in an *E. coli pyrE* mutant (22).

Other auxotrophic markers include *trp-1* and *arg12* of *N. crassa* [encoding a trifunctional enzyme for tryptophan biosynthesis (see [TRP Operon](#)) and ornithine carbamoyl transferase, respectively]; and *amdS* and *niaD* of *A. nidulans* (encoding acetamidase and nitrate reductase, respectively). In many instances, these genes are conserved so that their selective use in heterologous hosts is widely known. For example, *trpC* gene and *pyrG* gene of *A. nidulans* correspond to the *trp-1* and *pyr-4* of *N. crassa*. Many, many more auxotrophic markers are known than are listed here and have been used in a wide variety of fungi.

Auxotrophic markers are an excellent source of selectivity in fungi because a single copy of the gene is usually suitable for complementation. Second, the marker itself may be used to direct a chromosomal integration event at its homologous site. And third, selections against the auxotrophy usually give low backgrounds so that false positives are kept to a minimum. However, the disadvantage of using an auxotrophy as a selectable marker is that one requires a starting mutant strain. While some mutations may be directly selected for, such as *ura3* using 5-FOA (5-fluoroorotic acid), or screened visually using color pigmentation, such as *ade2* (23), most auxotrophies require screening thousands of mutants.

If a suitable auxotrophic marker does not exist, a number of “[antibiotic](#)”-resistance genes have been

cloned from a variety of organisms that function in heterologous hosts and have been used successfully in cloning strategies. For example, resistance to the anti-[tubulin](#) drug benomyl is provided by the *Bml* gene of *Neurospora* and the *BenAR* gene of *Aspergillus nidulans* (24, 25). Many antibiotic resistance genes from bacteria also function suitably as selective markers in fungal transformations: phleomycin, hygromycinB, and methotrexate (7, 26-28). **Drug-resistance** markers have the advantage that one does not require a recessive mutant strain as a host. In addition, drug resistance markers provide a simple means of multiply transforming a strain without the necessity of genetic crosses or other manipulations required to make a strain multiply mutant. On the other hand, unlike markers where there is homology between the gene present on the plasmid and the gene in the chromosome, these heterologous markers cannot be used to direct integration events within the host genome. That very property maybe used to advantage when one desires to lower background integration events by using heterologous marker selection. One caution is noted here: Selection for drug resistance tends to cause high backgrounds, because single-gene mutations (common during transformation procedures) often lead to elevated levels of resistance.

### 3. Integrative Transformation

The characteristics of transformation systems developed for the majority of filamentous fungi are very similar to one another. In most instances, for example, transformation with circular double-stranded plasmids leads to integrative transformation at both nonhomologous and homologous sites, although nonhomologous or ectopic integration events exceed homologous events in most systems (12). The integrations may be single or multiple events and are often accompanied by plasmid rearrangement. This is in contrast to events that most often occur in yeasts, in which homologous integrations are favored and plasmid rearrangements are rare. Even in mammals, where nonhomologous integrations far exceed homologous events, plasmid integrity is favored. In both *Podospora anserina* and *Ascobolus immersus*, nonhomologous integrations are more frequent than homologous ones, although during a homologous integration in *Podopsora* an integration/duplication event is the more likely outcome, whereas in *Ascobolus* gene substitution is more likely to occur (29).

Increasing relative rates of homologous versus nonhomologous integrations has been achieved in several ways: First, increasing the length of homology increased homologous integration events by up to 50–95% in *P. anserina* (30). However, the extent of the homology in these experiments was generated by using cosmids and reached 30 to 50 kbp in some instances, generally impractical for routine laboratory use. Using single-stranded DNA or double-strand cutting within the region of homology was observed to increase the proportion of homologous gene substitutions in *A. immersus* up to 65–70% of the events (31) and for the phytopathogenic filamentous fungus, *Ashbya gossypii* (32). Both of these methodologies have been well documented in *S. cerevisiae* system as well. A further modification of this procedure called restriction enzyme-mediated integration (REMI) has been highly successful in increasing transformation frequencies in most filamentous fungi (see text below).

### 4. Nonintegrative or Replicative Transformation

This is the exception in most filamentous fungi (12). Most filamentous fungal vectors are not maintained extrachromosomally. While this may be due to a failure to obtain host autonomously replicative sequences (ARSS), it is difficult to reconcile that in those instances where large fragments of chromosomes have been cloned successfully, origins of replication are not present. It is possible that [centromeres](#) or [telomeres](#) or both are needed to stabilize plasmids in filamentous fungi. In *Podospora*, linear plasmids containing *Tetrahymena* telomeres were shown to be self-replicating; however, they were observed to be unstable when placed on nonselective media (33).

The advantage of having a nonintegrative system available is that it often imparts a very high transformation efficiency upon the host strain. In addition, because extrachromosomally maintained plasmids usually have elevated copy numbers, questions concerning function of extra copies of

proteins produced in a transformed strain may be addressed without prior knowledge of inducible promoters. However, a certain level of risk also occurs in that overexpression of certain proteins is known to cause lethality in some instances.

## 5. Methods of Gene Cloning and Manipulation

The tools that one uses can provide a wealth of information about a wide variety of fungi—even those with intractable genetic systems. Clearly, isolating fungal genes by [complementation](#) of previously characterized mutations is a highly successful strategy. However, it may be impractical in some systems to generate suitable mutations. In those instances, one can take advantage of the ectopic nature of fungal transformation in the filamentous fungi in order to introduce mutations that may be in the form of duplication, disruption, or activation of a gene. Less obvious is the ability to clone genes by insertional mutagenesis, in which the transformation results in the creation of a mutant **allele** that is marked by a molecular tag. Vectors containing bacterial [transposons](#) have been constructed that have been used for gene disruption in *A. nidulans* or *M. grisea* strains auxotrophic for arginine or sensitive to hygromycin, respectively (34). The use of transposons has proven useful in many systems, but it has not been heavily pursued in the filamentous fungi. Nevertheless, putative transposons have been identified from *A. fumigatus* (35), *N. crassa* (36), and *M. grisea* (37) that may be adaptable to this problem.

**Restriction enzyme**-mediated integration (REMI) is a recent development in which the addition of a restriction enzyme to the DNA solution is used to enhance significantly the frequency of integrative events. It has been widely used in the fungal field (38-49). A modification of this approach, in which the transforming DNA solution also contains an inducible promoter along with a selectable marker, permits the random insertion of DNA sequences upstream from open reading frames, which may then be controlled by the ectopic promoter (T. Adams, personal communication). Identification of genes by activation may, of course, be achieved by creating transcription or translation fusion libraries *in vitro* and then introducing them into suitable hosts.

Other strategies for gene identification and cloning in fungi are widely known and practiced in other systems (42). Genes may be cloned by heterologous complementation in a related mutant host [such as the *pyr4* gene of *N. crassa*, or the *cdc10* gene of *Candida albicans* (43)]; or they may be identified by homology during sequencing projects (see the WEB sites for the various genome projects described below) and identified either by degenerate **PCR** or short homologous regions in cloned [libraries](#). Furthermore, [reporter genes](#) such as *lacZ* (encoding **b-galactosidase**), *uidA* (encoding **b-glucuronidase** of GUS), and [green fluorescent protein](#) from *Aequorea victoria* (44-49) have all been shown to be functional in one or more of the fungi. Thus, the fundamental tools available for molecular analysis of the yeasts are, by and large, available for the filamentous fungi. These tools have bridged the gap between the few fungi studied by geneticists for decades and those fungi whose genetic systems had remained largely undeveloped, but whose biology nevertheless compels us to understand their unique place in the biological world.

## Bibliography

1. M. E. Case, M. Schweizer, S. R. Kushner, and N. H. Giles (1980) Proc. Natl. Acad. Sci USA **76**, 5259–5263.
2. R. K. Beri and G. Turner (1987) Curr. Genet. **11**, 639–641.
3. F. Sanchez, M. Lozano, V. Rubio, and M. A. Penalva (1987) Gene **51**, 97–102.
4. L. LeChevanton and G. Leblon (1989) Gene **77**, 39–49.
5. D. M. Binninger, C. Skrzynia, P. J. Pukkila, and L. A. Casselton (1987) EMBO J **6**, 835–340.
6. F. G. Chumley, K. A. Parsons, and B. Valent (1985) J. Cell. Biochem. Suppl. **9C**, 197.
7. N. Durand, P. Reymond, and M. Fevre (1990) Transformations of penicillium to hygromycin B and phleomycin resistance. *British Mycological Society, 8th General Meeting*, Nottingham, England.

8. J. C. Edman and K. J. Kwon-Chung (1990) *Mol. Cell Biol.* **10**, 4538–4544.
9. J. Rambosek and J. Leach (1987) *CRC Crit. Rev. Biotechnol.* **6**, 357–393.
10. M. J. Hynes (1997) *J. Genet.* **75**, 297–311.
11. D. J. Ballance, F. P. Buxton, and G. Turner (1983) *Biochem Biophys. Res. Commun.* **112**, 284–289.
12. J. R. S. Fincham (1989) *Microbiol. Rev.* **53**, 148–170.
13. S. S. Dhawale, J. V. Paietta, and G. A. Marzluf (1984) *Curr. Genet.* **8**, 77–79.
14. K. N. Faber, W. Harder, G. Ab, and M. Veenhuis (1995) *Curr. Genet.* **25**, 305–310.
15. M. Ward, K. H. Kodama, and L. J. Wilson (1989) *Exp. Mycol.* **13**, 289–293.
16. M. B. Kurtz, D. R. Kirsch, and R. Kelly (1988) *Microbiol. Sci.* **5**, 58–63.
17. A. Aleksenko and A. J. Clutterbuck (1997) *Fungal Genet. Biol.* **21**, 373–387.
18. Z. An, M. L. Farman, A. Budde, S. Taura, and S. A. Leong (1996) *Gene* **176**, 93–96.
19. A. Y. Aleksenko and A. J. Clutterbuck (1995) *Curr. Genet.* **28**, 87–93.
20. A. Aleksenko, I. Nikolaev, Y. Vinetski, and A. J. Clutterbuck (1996) *Mol. Gen. Genet.* **253**, 242–246.
21. A. I. Bogdanova, O. S. Kustikova, M. O. Agaphonov, and M. D. Ter-Avanesyanyan (1998) *Yeast* **14**, 1–9.
22. J. Begueret, V. Razanamparany, M. Perrot, and C. Barreau (1984) *Gene* **32**, 487–492.
23. M. B. Kurtz, M. W. Cortelyou, and D. R. Kirsch (1986) *Mol. Cell. Biol.* **6**, 142–149.
24. M. J. Orbach, E. B. Porro, and C. Yanofsky (1986) *Mol. Cell. Biol.* **6**, 2452–2461.
25. P. W. Dunne and B. R. Oakley (1988) *Mol. Gen. Genet.* **213**, 339–345.
26. M. Kolar, P. J. Punt, C. A. M. J. J. van den Hondel, and H. Schwab (1988) *Gene* **62**, 127–134.
27. B. Austin, R. M. Hall, and B. M. Tyler (1990) *Gene* **93**, 157–162.
28. O. C. Yoder, K. Weltring, B. G. Turgeon, and H. D. Van Etten (1986) "Technology for Molecular Cloning of Fungal Virulence Genes". In *Biology and Molecular Biology of Plant-Pathogen Interactions*, (J. A. Bailey, ed.), Springer-Verlag, Berlin.
29. J.-L. Rossignol and M. Picard (1991) "*Ascobolus immersus* and *Podospora anserina*: Sex, Recombination, Silencing and Death". In *More Gene Manipulations in Fungi*, (J. W. Bennett, L. L. Lasure, eds.), Academic Press, New York, pp. 267–290.
30. M. Picard, R. Debuchy, J. Julien, and Y. Brygoo (1987) *Mol. Gen. Genet.* **210**, 129–134.
31. C. Goyon and G. Faugeron (1989) *Mol. Cell. Biol.* **9**, 2818–2827.
32. S. Stenier, J. Wendland, M. C. Wright, and P. Philippsen (1995) *Genetics* **140**, 973–987.
33. M. Perrot, C. Barreau, and J. Begueret (1987) *Mol. Cell. Biol.* **7**, 1725–1730.
34. L. Hamer and S. Gilger (1998) proteal found on web (site currently unavailable).
35. C. Neuveglise, J. Sarfati, J. P. Latge, and S. Paris (1996) *Nucleic Acids Res.* **24**, 1428–1434.
36. E. B. Cambareri, J. Helber, and J. A. Kinsey (1994) *Mol. Gen. Genet.* **242**, 658–665.
37. P. Kachroo, S. A. Leong, and B. B. Chatoos (1994) *Mol. Gen. Genet.* **245**, 339–348.
38. M. Bolker, H. U. Bohnert, K. H. Braun, J. Gohl, and R. Kahmann (1995) *Mol. Gen. Genet.* **248**, 547–552.
39. R. S. Redman and R. J. Rodriguez (1994) *Exp. Mycol.* **18**, 230–246.
40. Z. Shi, D. Christian, and H. Leung (1995) *Phytopathology* **85**, 329–333.
41. J. A. Sweigard, A. M. Carroll, L. Farrall, F. G. Chumley, and B. Valent (1998) *Mol. Plant Microbe Interact.* **11**, 404–412.
42. J. Agnan, C. Korch, and C. Selitrennikoff (1997) *Fungal Genet. Biol.* **21**, 292–301.
43. B. J. DiDomenico et al. (1994) *Mol. Gen. Genet.* **242**, 689–698.
44. I. N. Roberts, R. P. Oliver, P. J. Punt, and C. A. M. J. J. Van denHondel (1989) *Curr. Genet.* **15**,



177–180.

45. R. F. M. Van Gorcom, P. J. Punt, P. H. Pouwels, and C. A. M. J. Vanden Hondel (1986) *Gene* **48**, 211–217.
46. Y. Couteaudier, M. J. Daboussi, A. Eparvier, T. Langin, and J. Orcival (1993) *Appl. Environ. Microbiol.* **59**, 1767–1773.
47. D. C. Prasher (1995) *Trends Genet.* **11**, 320–323.
48. R. Kahmann (1996) *Mol. Gen. Genet.* **252**, 503–509.
49. J. Morschhauser, S. Michel, and J. Hacker (1998) *Mol. Gen. Genet.* **257**, 412–420.

## Transforming Growth Factors

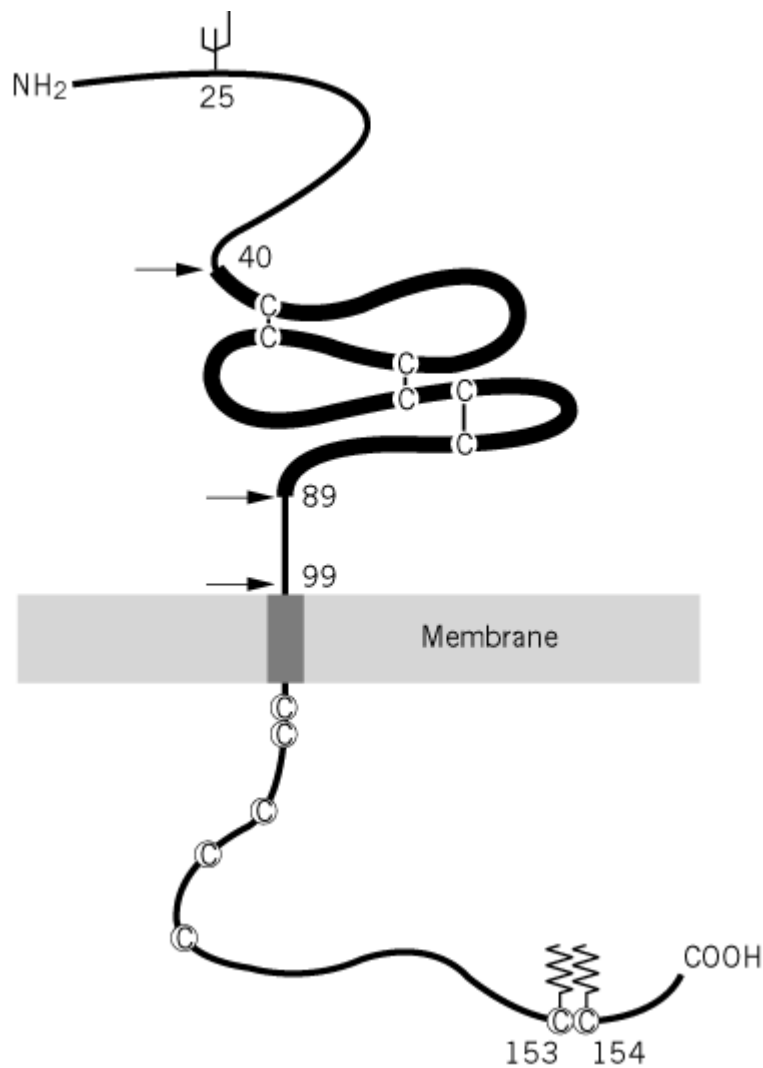
The biological activity of transforming growth factor (TGF) was first described in 1978 by De Larco and Todaro (1). The addition of supernatant of Rous-sarcoma virus-transformed fibroblasts to normal rat kidney cells or mouse fibroblasts enabled these cells to form colonies under anchorage-independent growth conditions in agarose. This growth property is limited to transformed and neoplastic cells and therefore the activity was called TGF. Further studies have shown that the transforming capacity was mediated by two unrelated proteins that have been termed TGF $\alpha$  and TGF $\beta$ , respectively. Their coidentification in this particular assay was rather accidentally as both have little sequence homology and display distinct biological activities. TGFs are now known to regulate a number of processes during normal development. The activity originally described as TGF $\beta$  has developed into a superfamily of numerous related molecules. They bind to a heterocomplex of receptor molecules with Ser/Thr kinase activity. Only recently a new family of molecules has been identified that serve as intracellular targets for these kinases and function as signal transducers to the nucleus.

### 1. Transforming Growth Factor $\alpha$ (TGF $\alpha$ )

#### 1.1. Structural Properties of the TGF $\alpha$ Protein

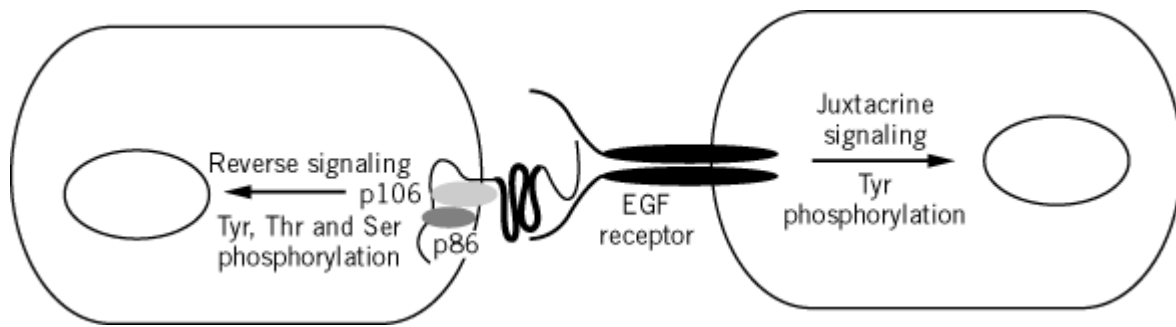
The TGF $\alpha$  gene spans a region of about 100 kb and includes 6 exons that encode a membrane-anchored precursor molecule (proTGF $\alpha$ ) of 160 amino acids (2) (Fig. 1). The mature 50–amino acid polypeptide is released from the precursor by cleavage of Ala/Val bonds at both termini by an elastase-like protease. Cleavage at the N-terminus occurs rapidly while the release of mature TGF $\alpha$  by proteolysis of the C-terminal site is slow and tightly regulated (3). Many cell types secrete larger TGF $\alpha$  species due to differential cleavage at the C-terminus and heterogeneous glycosylation of the extracellular domain.

**Figure 1.** Schematic representation of the TGF $\alpha$  precursor molecule. The bold line indicates the mature TGF $\alpha$  molecule. Proteolytic cleavage sites are indicated by arrows. N-linked carbohydrates (position Asn 25) and palmitoylation sites (positions Cys 153 and Cys 154) are indicated. Cys residues are marked by ©.



The highly conserved epidermal growth factor (EGF)-domain, characterized by six Cys residues, is essential for proper interaction of soluble TGF $\alpha$  with the EGF receptor. This binding results in receptor dimerization. In addition to autocrine and paracrine stimulation of EGF receptors by soluble TGF $\alpha$ , activation of EGF receptors may also occur through proTGF $\alpha$ . By means of site-directed mutagenesis the cleavage sites were altered yielding noncleavable proTGF $\alpha$  precursors. Coculture experiments confirmed that the membrane-bound proTGF $\alpha$  was biologically active and could activate the EGF receptor in a similar way as the soluble ligand (4, 5), a phenomenon called juxtacrine signaling (Fig. 2). Two cellular proteins (p106 and p86) associate with the transmembrane and cytoplasmic part of proTGF $\alpha$  and this complex can phosphorylate exogenous substrates on Tyr, Ser and Thr residues. The 39-amino acid long and highly conserved cytoplasmic part itself does not possess any enzymatic activity. However, experiments that introduced deletions and mutations revealed its critical role for complex formation and kinase activity. Two palmitoylated Cys residues at the C-terminus are necessary for association with p86 and truncation of the cytoplasmic tail almost completely abrogated kinase activity (6). Complex formation between the transmembrane growth factor and a kinase activity could therefore lead to a two-directional signaling creating a situation in which proTGF $\alpha$  would serve both as a ligand and a receptor: the extracellular domain would elicit a signal in neighboring cells through juxtacrine mechanisms while the complex between cytoplasmic region and kinase would signal into the opposite direction (Fig. 2).

**Figure 2.** Two-directional signal transduction of membrane TGF $\alpha$ . The extracellular domain of proTGF $\alpha$  can trigger signal transduction after interaction with EGF receptors on neighboring cells. Cellular proteins that associate with the intracellular domain attribute kinase activity to proTGF $\alpha$  and allow reverse signaling in the TGF $\alpha$ -producing cell.



## 1.2. Biological Activities of TGF $\alpha$

TGF $\alpha$  is widely expressed in developing embryos and in a number of adult mammalian tissues. Important functions have been attributed to TGF $\alpha$  in wound healing and the homeostasis of several tissues. In addition, TGF $\alpha$  is thought to effectively drive the proliferation of different cell types, especially those of epithelial origin and to regulate normal development at various stages (7). Therefore, it was rather surprising that mice in which both TGF $\alpha$  alleles had been inactivated developed normally and were generally healthy and fertile. Such mice show two major phenotypic changes: curly hair as a consequence of a dramatic disorganization of hair follicles and abnormal eye development of variable incidence and severity (8, 9). The lack of more dramatic changes might, however, reflect a compensation of TGF $\alpha$  by other members of the EGF-superfamily that also can serve as a ligand for the EGF receptor.

Probably all epithelial cells synthesize TGF $\alpha$  under normal conditions. They coexpress the EGF receptor thereby being able to respond to TGF $\alpha$  in an autocrine fashion. The overexpression of TGF $\alpha$  in mice has underlined its mitogenic activity for epithelial and mesenchymal cells *in vivo*. Deregulation of the ligand might, therefore, lead to impaired growth properties. In fact, the pathological hyperproliferation of skin keratinocytes that is characteristic for psoriasis could be explained by locally increased production of TGF $\alpha$  (10). It is well established that numerous neoplastic and transformed cell lines express TGF $\alpha$ , quite often to a higher extent than their normal counterparts. Transfection of a weakly tumorigenic carcinoma cell line with TGF $\alpha$  enhanced its tumorigenicity in nude mice (11) suggesting a role of TGF $\alpha$  during certain stages of carcinogenesis. This question has been investigated in transgenic mice where TGF $\alpha$  was overexpressed under the control of different promoters. The results suggest that the effects of TGF $\alpha$  overexpression are pleiotropic and tissue-specific. Several tissues and organs displayed hyperplasia of epithelial cells while the tissue architecture remained largely intact. However, in certain tissues such as liver and mammary glands TGF $\alpha$  overexpression induced neoplastic transformation (12-14). This process may occur preferentially in predisposed cells that have undergone other changes such as overexpression of the EGF receptor or dysregulated oncogenes (7) though the overexpression of TGF $\alpha$  might substitute for some of these events (15). Thus, besides being a physiological ligand that is important for cell proliferation and development, TGF $\alpha$  might also act as an oncoprotein *in vivo* under certain circumstances.

## 2. Transforming Growth Factor b (TGFb)

### 2.1. The TGFb Superfamily

The TGFb superfamily is a group of polypeptide molecules that affect cell proliferation and differentiation in numerous ways. The first report describing a biological activity of a member of the TGFb superfamily appeared in 1981 (16). The prototype molecule of the family is TGFb1 which shows sequence homology to all family members. In vertebrates, five genes encoding distinct TGFb-

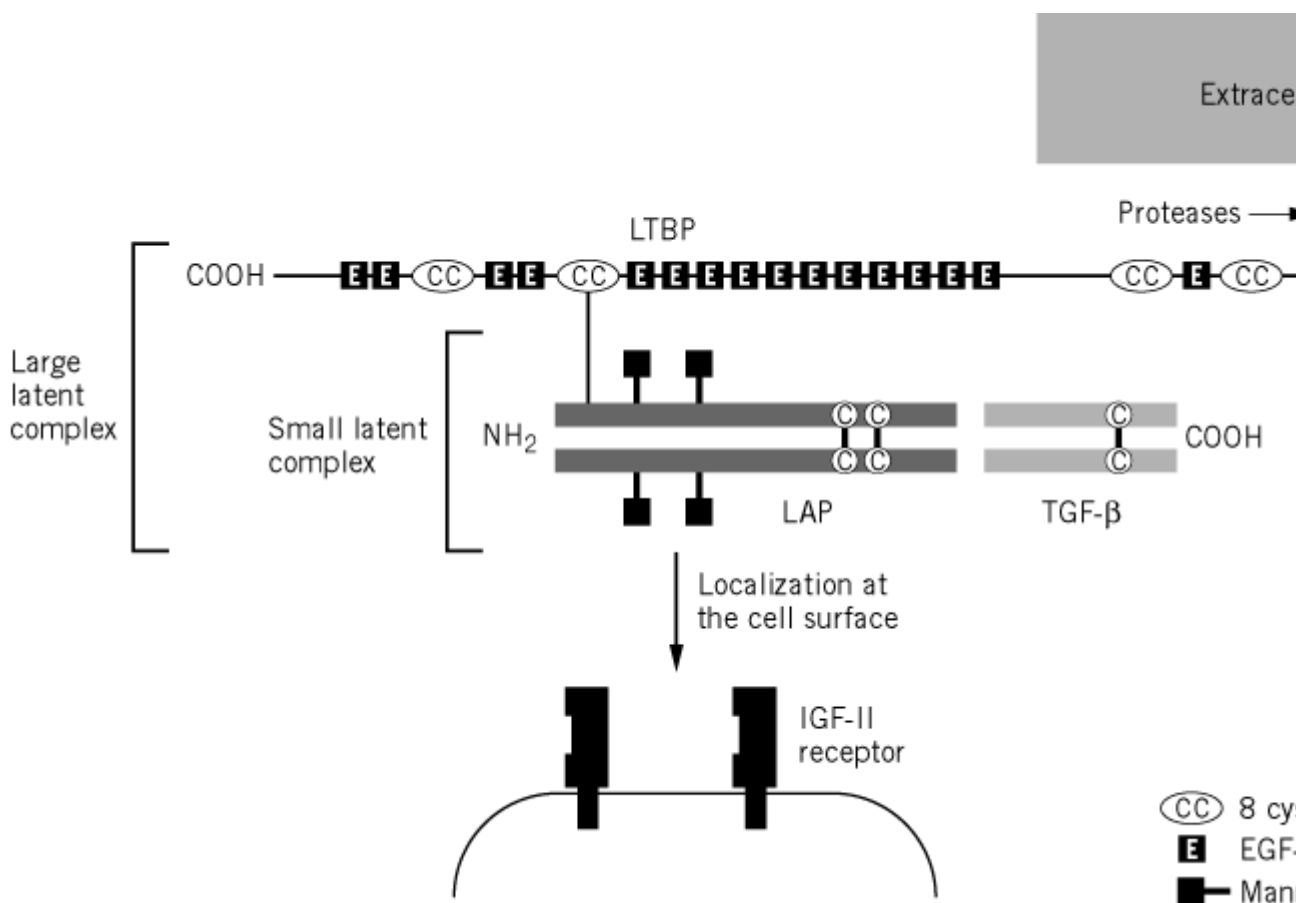
isoforms have been identified. Each gene spans more than 100 kb and contains seven exons. Mammalian cells express three differentially regulated isoforms (TGFb1, TGFb2 and TGFb3). Additional isoforms have been identified from chicken (TGFb4) and *Xenopus laevis* (TGFb5). The peptides are synthesized as large precursor molecules. The amino-terminal part contains a leader sequence (15–25 amino acids) and a prodomain of variable size (50–375 amino acids) which is poorly conserved. The biological activity is located in the C-terminal region which shares 70 to 80% sequence homology between different isoforms. The interspecies homology of TGFb isoforms is extremely high: mature human TGFb1 is identical to porcine, simian and bovine TGFb1 while the murine isoform differs by a single amino acid. The most characteristic structure of mature proteins is the so-called cysteine knot. Crystallographic studies of TGFb2 revealed that this motif is a ring of eight amino acids that are held together by two disulfide bonds while a third one passes through the middle of it (17). A similar three-dimensional structure has also been identified in members of the nerve growth factor and platelet-derived growth factor (PDGF) families. The seventh conserved cysteine forms an intermolecular disulfide bond between two monomers to form a dimer. Increasing information on the structural properties of TGFb-isoforms has enabled the identification of further related molecules. In humans and mice these new members include activins, Müllerian inhibitory substance, bone morphogenic proteins, nodal, growth differentiation factors and others. Further members of the TGFb-superfamily have been identified in nonmammalian organisms: the decapentaplegic gene (*dpp*) plays an essential role during several stages of development in *Drosophila*, *dorsalin* appears to be important for the neural patterning in chicken while the *Vgl* gene participates in embryonic axis formation of *Xenopus laevis* (18, 19). Presently, the TGFb-superfamily comprises more than 30 members.

## 2.2. Activation of Latent TGFb

The bioactive form of TGFb1 is a 25-kDa molecule that is usually composed of two identical disulfide-linked monomers though heterodimeric forms have also been described (20, 21). However, TGFb is not directly secreted as a bioactive molecule but is released in a latent form (Fig. 3) that is biologically inactive and cannot interact with cell surface receptors (22). The precursor is cleaved into the N-terminal latency associated protein (LAP) and the mature TGFb peptide. Both molecules form homodimers that interact via electrostatic interactions to form the tetrameric small latent complex (SLC). LAP is thought to be essential for the proper folding and transport of the complex. In addition, it contains mannose-6-phosphate (M6P) residues that allow interaction with the M6P/IGF-II receptor. Thereby TGFb could be localized at the cell surface and this interaction has been shown to be necessary for the activation of latent TGFb (23). Additional proteins are frequently associated with the SLC to yield the large latent complex. The latent TGFb-binding protein (LTBP) is a glycoprotein of 125 to 190 kDa that is covalently bound to LAP through its third eight cysteine region. LTBP appears to have a function in determining TGFb-bioavailability through targeting the latent complex to the extracellular matrix (ECM). LTBP apparently is covalently linked to the ECM as only proteases release it. Experimentally, latent TGFb can be activated by denaturing conditions (heat, acidification, chaotropes, detergents) and by radiation though such mechanisms are unlikely to be of general physiological relevance. Mechanisms that are more likely to represent physiological mechanisms include nonproteolytic activation, enzymatic deglycosylation of LAP and activation by different proteases (22). A central role in TGFb activation has been attributed to the thrombospondins (TSPs), a family comprising five extracellular proteins (24). The best characterized member is TSP-1 which interacts with a number of membrane-bound molecules but also with growth factors such as PDGF and TGFb. TSP-1 is composed of globular domains at the N- and C-terminus, a domain homologous to procollagen and three repetitive motifs, designated TSP type 1, 2, and 3 repeats (TSRs). The WSHWSPW sequence in the second TSR interacts with TGFb while an <sup>412</sup>RFK<sup>415</sup> motif located at the C-terminus of TSR-1 mediates the activation of TGFb1 upon interaction with the <sup>54</sup>SKL<sup>57</sup> sequence at the amino terminus of LAP. This mechanism of activation appears to be conserved among TGFb-isoforms as TGFb2 can be activated in a similar way (25). There is strong evidence that activation of TGFb by TSP-1 also occurs in vivo: TSP-1 deficient mice first develop normally but acquire histologic abnormalities in several organs similar to those recorded in TGFb1 knockout mice. Treatment of these mice with synthetic KRFK peptides resulted

in a marked improvement of lung abnormalities (26). In striking contrast are results which indicate that TSP-1 fails to activate TGF $\beta$  produced by platelets (27). Apparently further requirements such as posttranslational modifications of TSP-1 or the presence of additional cofactors are needed in these cells suggesting the existence of tissue- and/or cell type-specific mechanisms for the activation of TGF $\beta$ .

**Figure 3.** Schematic representation of latent TGF $\beta$ . The disulfide bond between the latent TGF $\beta$ -binding protein (LTBP) protein (LAP) is indicated by a thin line. Covalent bonds link latent TGF $\beta$  through LTBP to the extracellular matrix while mannose-6-phosphate and the IGF-II receptor allow localization at the cell surface. Reproduced with modifications from



### 2.3. TGF $\beta$ Receptors

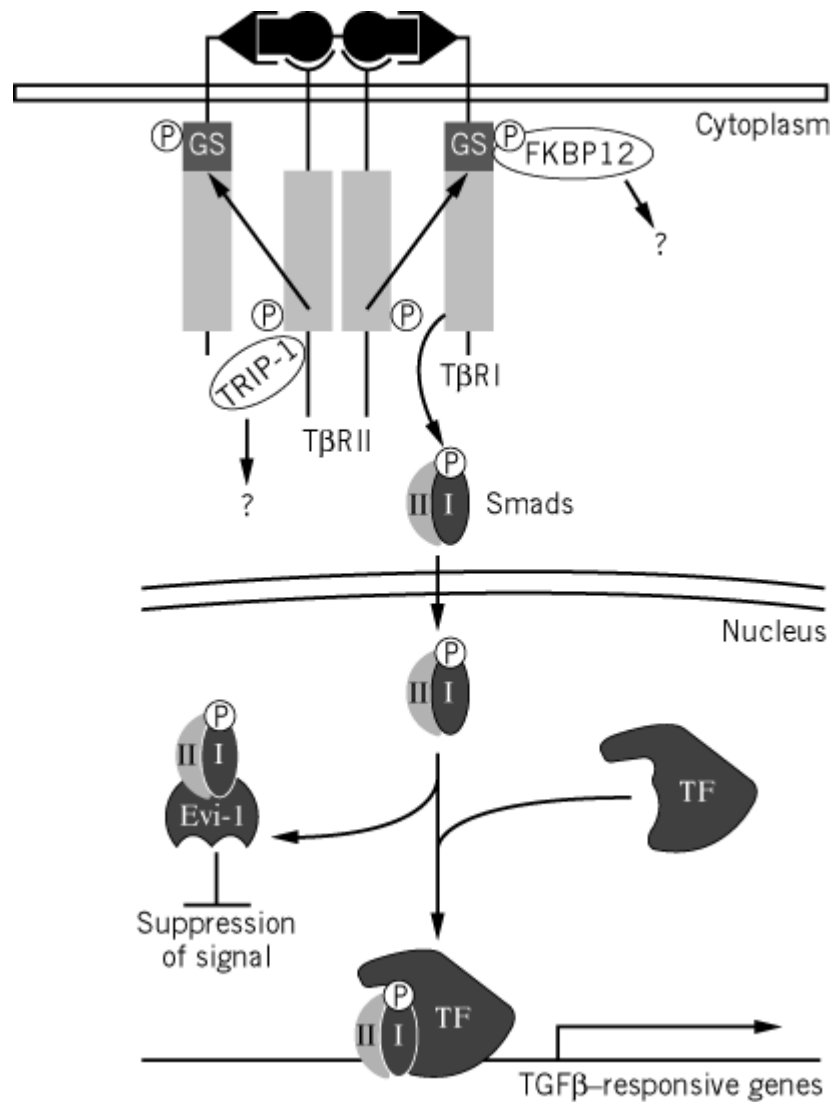
Affinity crosslinking studies using iodinated TGF $\beta$  have identified at least nine different surface molecules that can interact with TGF $\beta$ . However, subsequent studies revealed that sole interaction does not necessarily trigger a biological signal in every case. The major TGF $\beta$ -binding proteins are molecules of 53, 70 to 85, and 200 to 400 kDa, which are termed TGF $\beta$ -receptors type I, II and III (TbRI, TbRII, TbRIII), respectively. Meanwhile it is clear that only TbRI and TbRII are directly involved in signaling. The importance of both receptors for signaling was shown by chemically induced TGF $\beta$ -resistant mutants of epithelial cells which either failed to bind TGF $\beta$  to TbRI (R mutants) or to both TbRI and TbRII (DR mutants). Fusion of R and DR cells resulted in normal expression of TbRI and TbRII and a completely restored TGF $\beta$ -responsiveness of these cells (28). The puzzling fact that TbRI were present in hybrids despite its absence in both R and DR mutants was explained by a model in which TbRI required TbRII for ligand binding and the necessity of both receptors for signaling. This model has been confirmed upon cloning of the receptors.

The cloning of TbRI (29) and TbRII (30) revealed that both receptors lack a tyrosine kinase domain

but instead are transmembrane Ser/Thr kinases. The overall structure of both receptors is similar, comprising cysteine-rich extracellular domains, a single transmembrane region and the intracellular tail containing the kinase domain. However, each receptor has characteristic and unique structural features. TβRII possess a C-terminal tail rich in Ser and Thr and bind TGFβ when expressed alone possibly through a site formed by several Cys residues. Besides this cysteine box ligand specificity appears to be determined by the distribution of Cys residues in the extracellular domain. The TβRIII (beta-glycan) appears to modulate TGFβ-binding in positive and negative ways. In the absence of TβRIII relatively few TβRII bind TGFβ isoforms, especially TGFβ2, with high affinity. However, TGFβ presented by TβRIII results in high affinity binding of all three isoforms to TβRII (31). In contrast, in its soluble form TβRIII can effectively antagonize the action of TGFβ by sequestration of free ligand (32). The TβRI does not bind ligand when expressed alone and contains a characteristic conserved 29–amino acid sequence (GS domain) located at the N-terminal of the kinase domain.

The kinase of TβRII is constitutively active and autophosphorylates Ser and Thr residues. Neither its activity nor the phosphorylation sites are altered upon binding of TGFβ. However, upon ligand-binding a heterocomplex is formed with TβRI which is phosphorylated on Ser and Thr in the GS domain. This transphosphorylation activates the kinase of TβRI which then in turn delivers the signal to further cellular targets (33). Recent evidence suggests that the functional receptor complex is a heterotetramer that consists of two TβRI and TβRII molecules (34) (Fig. 4).

**Figure 4.** Signal transduction by TGFβ receptors. Binding of TGFβ to TβRII leads to formation of a heterocomplex with TβRI. The kinase (indicated in gray) of TβRII transphosphorylates the GS-domain of TβRI which activates the kinase of TβRI. The signal is transmitted to the nucleus by Smad proteins that associate with transcription factors (TF) and finally activate TGFβ-responsive genes. R = R-Smads, Co = Co-Smads.



#### 2.4. Smad Proteins—Key Molecules of TGF $\beta$ Signaling

Using the yeast two hybrid system several proteins have been identified that interact with the cytoplasmic part of TbRI or TbRII. FKBP12 interferes with TGF $\beta$  signaling by inhibition of TbRI phosphorylation (35). TRIP-1 associates with the TbRII and decreases TGF $\beta$  mediated induction of transcription from the plasminogen activator inhibitor-1 promoter (36). Recently, a closely related homolog has been identified in breast cancer cells suggesting a role for these molecules in tumor progression (37).

The fact that TbR possess Ser/Thr kinase activity lead to the speculation that the intracellular downstream targets would be different than those used by tyrosine kinases. The first protein identified was the *Mad* (*Mothers against dpp*) gene product of *Drosophila* (38) followed by the discovery of three genes (*sma-2*, *sma-3*, *sma-4*) that share mutant phenotypes with *daf-4*, a TbRII in *C. elegans* (39). Meanwhile nine related vertebrate genes have been described and this novel family has been termed Smads, for *sma*- and *Mad*-related proteins. Except for Smad6 these proteins contain an N-terminal DNA-binding domain termed MH-1 (*Mad*-homology domain 1) and a C-terminal MH-2 domain involved in transcriptional activation, Smad-oligomerization and interaction with receptor serine kinases. According to their structural and functional properties three classes of Smad proteins have been defined (40). R-Smads ('receptor-regulated' Smads) are phosphorylated at the C-terminal SSXS motif by the serine kinase of TbRI and activin type I (Smads-2 and -3) or BMP type I receptors (Smads-1, -5 and -8). Upon phosphorylation, R-Smads physically associate with a second

group termed Co-Smads ('common' Smads). Co-Smads which include Smad-4 and -4b appear to be required for the formation of a functional transcriptional complex. The third group of so-called Anti-Smads (Smads-6 and -7) function as inhibitors of TGF $\beta$  signaling by interference with the phosphorylation of R-Smads.

A network of regulators influences activation of T $\beta$ R and Smad proteins (40). Antagonists of activation include the immunophilin FKBP12 (35), the pseudoreceptor BAMBI (41) and Anti-Smads (42, 43). SARA (Smad anchor for receptor activation) binds to R-Smads thereby directing them to T $\beta$ Rs (44). Phosphorylation of R-Smads reduces their affinity to SARA. Phosphorylated R-Smads associate with Co-Smads in the cytoplasm and move to the nucleus. Both R- and Co-Smads are able to directly activate transcription via their MH-2 domain or they act in concert with other transcription factors (40) (Fig. 4). In addition, R-Smads may initiate transcription by interference with transcriptional repressors (45).

## 2.5. Biological Activities of TGF $\beta$ s

A prominent biological activity of TGF $\beta$ s is growth inhibition of most normal cells, especially those of epithelial origin. TGF $\beta$  prevents the progression into the S phase of the cell cycle by maintaining the retinoblastoma (Rb) protein in a hypophosphorylated form that arrests the cells in late G<sub>1</sub> (46).

Complexes of cyclins and cyclin-dependent kinases (cdk) are involved in hyperphosphorylation of Rb. TGF $\beta$  can inhibit the kinase activity of this complex by downregulation of cdk or induction of cdk inhibitors (47). In some cells of mesenchymal origin TGF $\beta$  induces growth stimulation. This activity has been reported to be mediated by PDGF which is induced by TGF $\beta$  (48). TGF $\beta$ s are also important molecules in wound healing. At the site of tissue injury platelets release high amounts of TGF $\beta$  which acts as a chemoattractant for monocytes, neutrophils and T lymphocytes. After the early phase TGF $\beta$  exerts immunosuppressive effects that counteract the inflammatory process. T lymphocytes acquire TGF $\beta$ -responsiveness and are growth-inhibited primarily by the interference of TGF $\beta$  with interleukin-2-mediated growth stimulation. The potential physiological immunosuppressive role of TGF $\beta$  is strengthened by results obtained with TGF $\beta$ 1-knockout mice which are born normally but die within 20 days of a wasting syndrome that is accompanied by multifocal inflammatory responses in several organs and elevated levels of inflammatory cells and cytokines (49). In contrast, homozygous TGF $\beta$ 3 null mutants are not viable and die within 24 hours after birth. They display a defective palatogenesis and exhibit histopathological abnormalities in the pulmonary system (50, 51). Certain developmental processes appear to be strictly controlled by one isoform of TGF $\beta$ . Its absence can not be compensated by other isoforms though the biological activity of TGF $\beta$  peptides is often quite similar in vitro. Excess of TGF $\beta$  may equally have dramatic consequences. Animals in which TGF $\beta$ 1 was overexpressed under the control of an epidermal-specific promoter displayed a heavily disturbed skin development and died within 24 hours (52). Targeted expression of TGF $\beta$ 2 in osteoblasts resulted in a dramatic loss of bone mass and the appearance of an osteoporosis-like phenotype (53). Inactivation of components of the T $\beta$ RI signaling cascade caused severe early developmental abnormalities and embryonic lethality in homozygous animals (54, 55). As a single Smad protein is used by different ligands these dramatic effects might be explained by the simultaneous inactivation of multiple signal pathways of TGF $\beta$ -related factors.

## 2.6. TGF $\beta$ and Cancer

The majority of neoplastic cells is resistant to the growth inhibitory activity of TGF $\beta$  suggesting that TGF $\beta$ -unresponsiveness may attribute a growth advantage to the malignant cell. However, only recently some of the responsible mechanisms have been unraveled on a molecular level. In colorectal carcinoma frameshift mutations have been detected that introduce stop codons on the N-terminal site of the transmembrane domain of T $\beta$ RII. This mutation leads to the synthesis of truncated receptor proteins that are functionally inactive. Unexpectedly, these lesions occur exclusively in a subpopulation of colon cancer cells that are characterized by mismatch repair deficiencies and demonstrate microsatellite instability (56). Meanwhile the same mutational inactivation of the T $\beta$ RII has been detected in a variety of human tumors. Transfection of wild type T $\beta$ RII into tumor cells that express no or only mutated T $\beta$ RII suppressed anchorage-independent growth in soft agar and abolished tumor formation in nude mice (57, 58). These experiments provide a direct proof that



TbRII can function as a tumor suppressor gene. Some of the intracellular target molecules may have similar activities. Functional mutations and loss of heterozygosity have been discovered in DPC4 (Smad4) and Smad2 (59, 60) which are both essential for the signaling of TGFb to the nucleus. Finally, the ligand itself and molecules upstream of TGFb such as the IGF-II receptor are candidate tumor suppressor genes (61). Future experiments will have to confirm the functional activity of these components to suppress tumor growth. Furthermore, mutational inactivation of the TGFb signaling pathway entitles neoplastic cells to use endogenously produced TGFb as an effective weapon against attacking immune cells. TGFb is certainly the most potent immunosuppressive cytokine that can influence almost every cell type and reaction of the immune system (62, 63). Through TGFb-mediated suppression of immune responses tumor cells may escape from immune surveillance and gain further growth advantage.

### 3. Conclusion

Since the discovery of TGFs more than twenty years ago our knowledge on their molecular structure and their biological activities has tremendously increased. While this progress has answered many questions it has also created new ones. It even holds true for the ligands that are well-known for quite some time on the molecular level. Though the activation of mature TGFb has been unraveled in some tissues it remains still unclear how the cytokine is activated in other cell types. To date neither the physiological importance of membrane-bound TGFa nor the reverse signaling pathway are fully understood. The Smad family has been identified as a key player in TGFb-mediated signal transduction and a number of cellular proteins have been discovered which modulate their activity. Additional partners may reveal completely new functions of these molecules and open yet unidentified connections to other pathways. Likewise, the understanding of the transcriptional control of Smad proteins is just at the beginning. The investigation of these complex interactions may open novel possibilities in different fields of research enabling us to generate the next round of interesting questions.

### Bibliography

1. J.E. De Larco and G.J. Todaro, Proc. Natl Acad. Sci. U.S.A. **75**, 4001–4005 (1978).
2. R. Derynck et al., Cell **38**, 287–297 (1984).
3. A. Pandiella and J. Massagué, Proc. Natl. Acad. Sci. U.S.A. **88**, 1726–1730 (1991).
4. S.T. Wong et al., Cell **56**, 495–506 (1989).
5. R. Brachmann et al., Cell **56**, 691–700 (1989).
6. L. Shum, C.W. Turck, and R. Derynck, J. Biol. Chem. **271**, 28502–28508 (1996).
7. R. Derynck, Adv. Cancer Res. **58**, 27–52 (1992).
8. N.C. Luetkeke et al., Cell **73**, 263–278 (1993).
9. G.B. Mann et al., Cell **73**, 249–261 (1993).
10. J.T. Elder et al., Science **243**, 811–814 (1989).
11. B.L. Ziober et al., J. Biol. Chem. **268**, 691–698 (1993).
12. E.P. Sandgren et al., Cell **61**, 1121–1135 (1990).
13. C. Jhappan et al., Cell **61**, 1137–1146 (1990).
14. Y. Matsui et al., Cell **61**, 1147–1155 (1990).
15. R. Vassar, M.E. Hutton, and E. Fuchs, Mol. Cell. Biol. **12**, 4643–4653 (1992).
16. A.B. Roberts et al., Proc. Natl. Acad. Sci. U.S.A. **78**, 5339–5343 (1981).
17. M.P. Schlunegger and M. Grütter, Nature **358**, 430–434 (1992).
18. D.M. Kingsley, Genes Dev. **8**, 133–146 (1994).
19. J. Massagué, L. Attisano, and J.L. Wrana, Trends Cell Biol. **4**, 172–178 (1994).
20. S. Cheifetz et al., Cell **48**, 409–415 (1987).
21. Y. Ogawa et al., J. Biol. Chem. **267**, 2325–2328 (1992).

22. J.S. Munger et al., *Kidney Interna.* **3**, 1376–1382 (1997).
23. P.A. Dennis and D.B. Rifkin, *Proc. Natl. Acad. Sci. U.S.A.* **88**, 580–584 (1991).
24. J. Lawler, *Curr. Opin. Cell Biol.* **12**, 634–640 (2000).
25. S.M.F. Ribeiro et al., *J. Biol. Chem.* **274**, 13586–13593 (1999).
26. S.E. Crawford et al., *Cell* **93**, 1159–1170 (1998).
27. M. Abdelouahed, A. Ludlow, G. Brunner, and J. Lawler, *J. Biol. Chem.* **275**, 17933–17936 (2000).
28. M. Laiho et al., *J. Biol. Chem.* **266**, 9108–9112 (1991).
29. P. Franzén et al., *Cell* **75**, 681–692 (1993).
30. H.Y. Lin et al., *Cell* **68**, 775–785 (1992).
31. F. Lopéz-Casillas, J.L. Wrana, and J. Massagué, *Cell* **73**, 1435–1444 (1993).
32. F. Lopéz-Casillas, H.M. Payne, J.L. Andres, and J. Massagué, *J. Cell Biol.* **124**, 557–568 (1994).
33. J.L. Wrana et al., *Nature* **370**, 341–347 (1994).
34. F. Weis-Garcia and J. Massagué, *EMBO J.* **15**, 276–289 (1996).
35. T. Wang et al., *Cell* **86**, 435–444 (1996).
36. L. Choy and R. Derynck, *J. Biol. Chem.* **273**, 31455–31462 (1998).
37. S. Matsuda et al., *Cancer Res.* **60**, 13–17 (2000).
38. J.J. Sekelsky et al., *Genetics* **139**, 1347–1358 (1995).
39. C. Savage et al., *Proc. Natl. Acad. Sci. U.S.A.* **93**, 790–794 (1996).
40. C.M. Zimmerman and R.W. Padgett, *Gene* **249**, 17–30 (2000).
41. D. Onichtchouk et al., *Nature* **401**, 480–485 (1999).
42. A. Hata, G. Lagna, J. Massagué, and A. Hammati-Brivanlou, *Genes Dev.* **12**, 186–197 (1998).
43. H. Hayashi et al., *Cell* **89**, 1165–1173 (1997).
44. T. Tsukazagi et al., *Cell* **95**, 779–791 (1998).
45. X. Shi et al., *J. Biol. Chem.* **274**, 13711–13717 (1999).
46. M. Laiho et al., *Cell* **62**, 175–185 (1990).
47. J. Massagué and K. Polyak, *Curr. Opin. Genet. Dev.* **5**, 91–96 (1995).
48. E.B. Leof et al., *Proc. Natl. Acad. Sci. U.S.A.* **83**, 2453–2457 (1986).
49. M.M. Shull et al., *Nature* **359**, 693–699 (1992).
50. G. Proetzl et al., *Nature Genetics* **11**, 409–414 (1995).
51. V. Kaartinen et al., *Nature Genetics* **11**, 415–421 (1995).
52. K. Sellheyer et al., *Proc. Natl. Acad. Sci. U.S.A.* **90**, 5237–5241 (1993).
53. A. Erlebacher and R. Derynck, *J. Cell Biol.* **132**, 195–210 (1996).
54. M. Nomura and E. Li, *Nature* **393**, 786–790 (1998).
55. C. Sirard et al., *Genes Dev.* **12**, 107–119 (1998).
56. S. Markowitz et al., *Science* **268**, 1336–1338 (1995).
57. J. Wang et al., *J. Biol. Chem.* **270**, 22044–22049 (1995).
58. L. Sun et al., *J. Biol. Chem.* **269**, 26449–26455 (1994).
59. S.A. Hahn et al., *Science* **271**, 350–353 (1996).
60. K. Eppert et al., *Cell* **86**, 543–552 (1996).
61. S.A. Markowitz and A.B. Roberts, *Cytokine Growth Factor Rev.* **7**, 93–102 (1996).
62. S. Chouaib et al., *Immunol. Today* **18**, 493–497 (1997).
63. J.J. Letterio and A.B. Roberts, *Clin. Immunol. Immunopath.* **84**, 244–250 (1997).

## Additional Reading

64. Böttinger E.P., Letterio J.J., and Roberts A.B., Biology of TGF- in knockout and transgenic mouse models, *Kidney Intern.* **51**, 1355–1360 (1997).
65. Lee D.C., Fenton S.E., Berkowitz E.A., and Hissong M.A., Transforming growth factor : expression, regulation, and biological activities, *Pharmacol. Rev.* **47**, 51–85 (1995). Covers all aspects of TGF biology, more than 300 references.
66. Massagué J. and Wotton D., Transcriptional control by the TGF- /Smad signaling system, *EMBO J.* **19**, 1745–1754 (2000).
67. Murphy-Ullrich J.E. and Poczatek M. Activation of latent TGF-B by thrombospondin-1: mechanisms and physiology, *Cytokine Growth Fact. Rev.* **11**, 59–69 (2000).
68. Roberts A.B. and Sporn M.B., "The transforming growth factor s", in M.B. Sporn and A.B. Roberts, eds., *Peptide growth factors and their receptor: Handbook of Experimental Pharmacology*, Springer Verlag, Heidelberg, Germany 1990, pp. 419–472.

## Transgenic Technology

Transgenic technology consists of the introduction of a defined genetic material into the **germline** of mice or other animals. This technique has attracted considerable attention, because it has been found to be a very powerful approach for different experimental fields, including molecular and **developmental** biology as well as medical research. Using different experimental strategies, it allows the *in vivo* analysis of **gene** function, the mechanisms controlling **gene expression**, and the generation of specific animal models for human diseases. This has resulted in new insights into normal and pathological processes, thus improving our knowledge and tools in diagnostic and therapeutic fields.

This review focuses on the use of transgenic mice for different purposes including a) **promoter** analysis, and gain- and loss-of-function [mutations](#): and b) their application in molecular and developmental biology. Further aspects will be discussed in the generation of animal models and their use in diagnostic and therapeutic approaches.

### 1. Experimental Approaches

Three basic methods are widely used for the generation of transgenic animals:

1. [Microinjection](#) of DNA into the male **pronucleus** of one-cell embryo.
2. Infection of pre- and postimplantation [embryos](#) by [retroviruses](#) (which will not be discussed here)
3. Transfer of the desired genetic material using embryonic [stem cells](#).

### 2. Pronuclear Injection of Genetic Material

The pronuclear injection of linearized DNA is the technique most used for generating transgenic animals including transgenic mice ([1](#), [2](#)). Female mice are superovulated by hormonal treatment and mated to produce one-cell stage embryos in which the two **pronuclei** are still visible ([2](#)). The linearized DNA devoid of vector sequences, known as the *transgene*, is injected by micromanipulation into the male pronucleus, which is larger ([2](#)). The following day, the two-cell stage embryos are transferred to a pseudopregnant female. Pseudopregnancy is achieved by mating

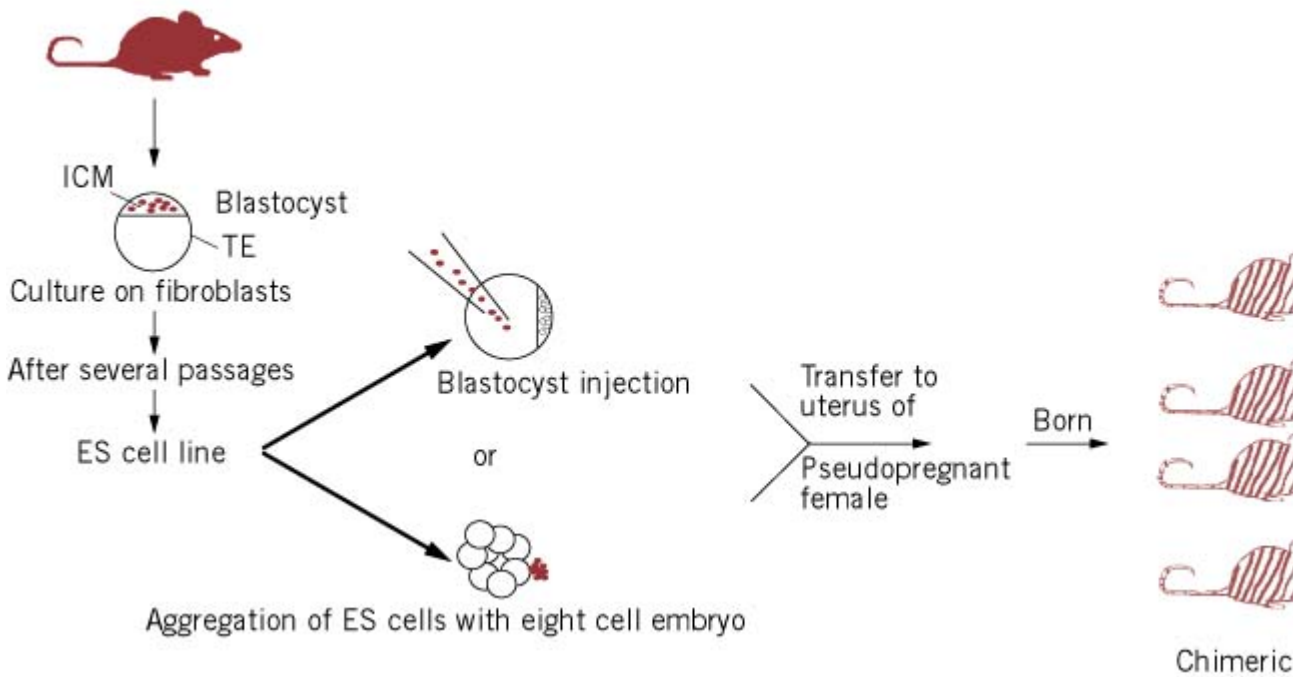
females with vasectomized males. Integration of the introduced DNA normally occurs randomly at the one-cell stage, resulting in a multicopy tandem of head to tail array of the transgene into the host [genome](#). After the manipulations, the newborn animals are tested using standard procedures, such as [Southern blot](#) or polymerase chain reaction (PCR) for transgenesis, using DNA derived from tail biopsies. Animals that carry the transgene are called *founders*. Founder mice are crossed with wild-type animals to establish a transgenic line, which carried a stably integrated transgene in the germline. Sometimes the integration of the introduced DNA occurs at later embryonic stages, rather than at the one-cell stage, resulting in a mosaic integration of the transgene into the animals. Consequently, further mating may not transmit the transgene to the germline.

### 3. Embryonic Stem Cell Techniques

Embryonic stem (ES) cell technology is usually used for the establishment of loss-of-function mutations (3). As it becomes more effective, it may in the future replace the pronuclear injection technique.

Embryonic stem cells are derived from the inner cell mass of the mouse blastocyst (Fig. 1) (4, 5). The blastocyst is cultured on a feeder layer of primary embryonic fibroblasts. After a few days, the blastocyst attaches to the feeder and the inner cell mass starts to grow. From an inner cell mass that is picked selectively, a homogenous population of embryonic stem cells can be established (Fig. 1). The ES cells, when cultured on embryonic fibroblasts and/or in the presence of leukemia inhibitory factor (LIF), can be propagated *in vitro* and retain their undifferentiated state and pluripotency. When ES cells are introduced into the mouse embryo, they contribute to all tissues, including the germline (5, 6). They have been shown to be the ideal vehicle for the introduction of exogenous DNA sequences into the mouse genome, because they can be genetically manipulated *in vitro* before starting *in vivo* procedures (7). Thus, more precise DNA manipulation such as site-specific integration by homologous [recombination](#) can be carried out. ES cells with stably integrated exogenous DNA are introduced into the embryo by blastocyst injections or by aggregation with eight cell embryos (5), and these manipulated embryos are then transferred to foster mothers. Animals derived from both host and ES cells are termed *chimeric*. The extent of ES cell contribution to the chimera can be assessed by coat color (Fig. 1). These chimeras can transmit the introduced genetic sequences to the mouse germline by mating with wild-type animals, resulting in the establishment of transgenic lines. The advantage of ES technology over pronuclear injection lies in the more controllable introduction of the DNA in ES cells. Stable ES cell lines may be analyzed for expression, copy number of the transgene, and so on, prior to proceeding to the mouse germline. However, the ES cell procedure is more time-consuming. Exogenous DNA, when introduced into cells, normally integrates randomly within the genome. However, specific targeting of exogenous DNA to a defined locus in the cellular genome is now possible by the means of homologous recombination. Mammalian somatic cells possess the basic machinery to mediate homologous recombination between endogenous chromosomal loci and exogenous DNA molecules (8, 9). In conjunction with ES cell technology, homologous recombination provides a powerful tool to study gene function through the generation of transgenic mice harboring specific mutations for any gene of interest.

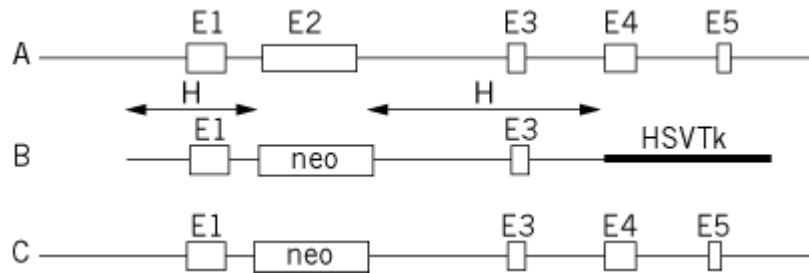
**Figure 1. Establishment of embryonic stem cells (ES).** Blastocysts are collected from mice (with agouti coat) and cultured on a feeder layer of embryonic fibroblasts. After further passages ES cell lines can be established. Introduction of ES cell mouse embryos by blastocyst injection or aggregation to eight cell embryos and subsequent transfer to foster mothers give rise to chimeric mice. Chimeric mice are schematically shown with brown stripes.



There are two types of targeting vectors in common use. Both consist of a selection marker flanked by sequences that are homologous to the gene of interest. Usually, [neomycin](#) (Neo) is the selection marker of choice and confers resistance to Geneticin (G418). In the first type of vector, the targeting construct is linearized in such a way that it remains collinear with the endogenous sequences of the gene to be targeted.

During homologous recombination events, the endogenous sequences are replaced by the targeting vector, and sequences that lie outside the domains of homology are excluded from the integration. This kind of vector is termed *replacement vector* ([10](#), [11](#)). The second type of vector is linearized so that the pairing of the targeting vector with the endogenous sequences of the target gene results in the insertion of the whole targeting construct, including plasmid sequences. The result is a partial duplication of the target gene. This kind of vector is termed *insertion vector* ([12](#)). Both types of vectors have been successfully used, and they exhibit similar targeting frequencies ([3](#)). In most cases, however, replacement vectors are used, due to their easier handling and construction. The targeting vector used for homologous recombination experiments consists of two domains with homology to the gene to be targeted, one at the 5' end and the other on the 3' end of the targeting construct (Fig. [2](#)). The selection marker, usually the Neo gene, or less usually that for hygromycin resistance or the **hypoxanthine-guanine phosphoribosyl transferase** (HPRT) [minigene](#) ([13](#)), is usually inserted between the 5' end and the 3' end homologous segments. The Neo gene is used first to monitor the presence of introduced DNA and second to interrupt, or **knock out**, the function of the gene. Therefore, the Neo gene is introduced into one exon, or it may replace a deleted coding sequence (Fig. [2](#)). A typical knock out experiment will start by **electroporating** the targeting construct into ES cells. Subsequent selection (24 h after electroporation) using geneticin (G418) kills all the cells that have not integrated the construct. After 8 to 10 days of selection, the G418-resistant colonies are screened by PCR or Southern blot for homologous recombination events. Homologous recombinant clones are further used to generate germline chimeras via aggregation or blastocyst injection (Fig. [1](#)). ES cells are in most cases derived from 129Sv agouti strain, and chimerism is identified by the coat color of agouti in contrast to the black or albino, of the host embryo (C57B1/6, NMRI, or CD1 strains).

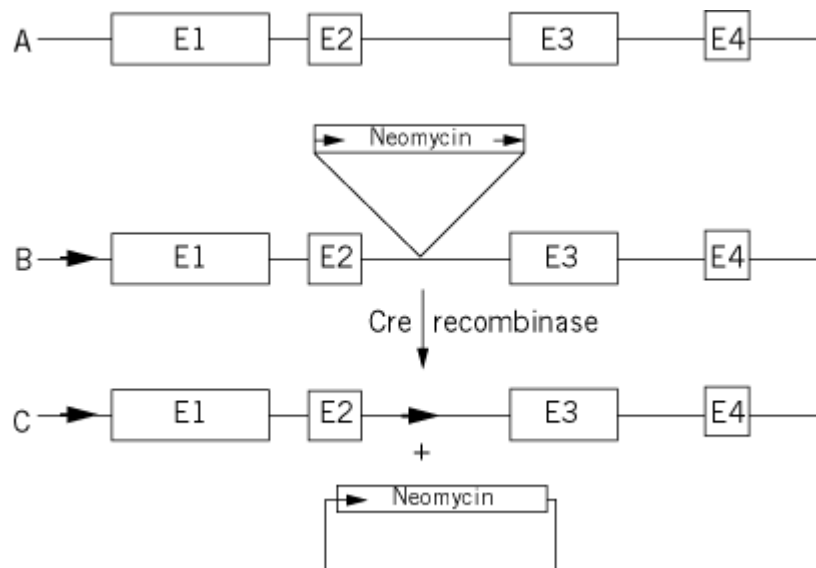
**Figure 2. Constructs for homologous recombination.** (a) The genomic organization normal gene to be targeted. Boxes represent exons (coding sequences E). (b) The targeting construct using the positive–negative selection procedure. The second exon of the gene was replaced by the PGKNeo. HSV-Tk is the thymidine kinase gene for the selection with gancyclovir, and lies outside of the segments of homologous DNA (indicated with H). (c) genomic organization of the targeted allele after homologous recombination.



The frequency of homologous recombination in mammalian cells is normally a rare event, in contrast to **yeast**. Two factors influence the frequency of homologous recombination: 1) the length of the homologous segments used for the construct and 2) the origin of the genomic DNA (14, 15). Increased targeting frequencies are obtained if using isogenic DNA (from the strain of mice from which the ES cells are derived) and if more than 10 kb of homologous DNA is included in the construct. One strategy to enhance the targeting frequency is the promoter-less construct. If the gene of interest is expressed in undifferentiated ES cells, a promoter-less Neo gene can be introduced into the construct. Basically, the Neo gene coding sequence is fused in frame to one exon or may replace the first exon containing the start site for [translation](#). Upon introduction of the construct into ES cells, the Neo gene is only active if it integrates in the vicinity of a promoter or is inserted into the target locus (16). The number of clones to be screened is low, so consequently the targeting frequency is highly enriched (about 100-fold). Most of the genes are not expressed in undifferentiated ES cells, however, and the targeting frequency for most genes remains very low. Capecchi and colleagues (17) have developed a technique that allows enrichment of homologous recombinant events when targeting genes that are not expressed in ES cells. In addition to the insertion of the Neo into one exon to disrupt gene function, a **thymidine kinase** (TK) gene from the Herpes simplex virus is placed outside of the homologous regions. Two selection drugs are now applied to the medium. The G418 kills all the cells that have not integrated the Neo gene (random and homologous recombinant clones), whereas Gancyclovir kills all the cells that have retained the TK gene (random integrations); the TK gene is lost during homologous recombination events. This method is called the *positive–negative selection* (PNS) procedure (17) and provides an enrichment factor of between four- and 10-fold (Fig. 2). The PNS procedure can also be used to achieve point mutations in the *hit-and-run* technique, in which a variation of homologous recombination is used to introduce deletions or point mutations in two steps (18). HPRT has been successfully introduced as a selection marker. The *in–out* targeting strategy is similar, but uses in the first step and *hpprt* minigene for selection, rather than neomycin resistance and thymidine kinase (13). It is also suitable for the generation of point mutation or subtle deletions. The selection system requires HPRT-negative ES cells. Recently, more sophisticated strategies have been developed, based on the Cre-Lox recombination system of **bacteriophage** P1 or Flp recombinase system from yeast (19). The Cre recombinase performs recombination at *Lox* sites (locus of X-ing-over) in **bacteria** and eukaryotic cells and is the most widely used system (20, 21). DNA that is flanked by 34-bp *LoxP* sequences is excised by the Cre recombinase. This system allows the removal of specific DNA segments, such as selection markers, and only the desired point mutations or deletions will be retained (Fig. 3). This strategy can be used to perform so-called *knock-in* experiments and *temporal* knockouts (22). Genes are then altered in certain tissues and at determined times of embryonic development (23). In fact, because the inactivation of several genes led to early lethal phenotype, the Cre-*LoxP* system is now used to design knockouts in such a way that tissue-specific inactivation of the targeted locus is

performed. In principle, in a first targeting experiment the gene of interest is specifically modified so that some of the coding sequences (exons) are flanked by two *LoxP*-sites in cis (Fig. 3). The *LoxP* sites, which are introduced into introns, should not affect the function of the gene. The thus modified allele is termed floxed allele. The *LoxP* sites are now recognized by Cre-recombinase to perform the excision of the coding sequences flanked by them (Fig. 3). Practically, floxed mice are generated from ES cells, which carry the floxed allele. The floxed mice are then crossed with transgenic mice expressing the Cre-recombinase under a specific promoter in a defined tissue or cell type. There are now several transgenic mice expressing the Cre-recombinase in specific tissues (so-called Cre-mice).

**Figure 3. *LoxP*-Cre recombinase system.** Scheme for the site-specific recombination of the *LoxP*-Cre system of the phage P1 in mammalian cells. A typical construct for the generation of so-called floxed allele is schematically shown. The Neo gene inserted into intron 2 is flanked by two *LoxP* sites shown as arrowheads. A third *LoxP* site is introduced in front of exon 1. In a first knockout experiment the gene shown in (a) is targeted and replaced by the targeting construct (b). After addition of the Cre recombinase, the neo gene is excised and from Neo only one *LoxP* site is left. The floxed allele (b) exhibits two *LoxP* sites flanking exon 1 and 2. The floxed allele is not the only product, which appears after addition of the Cre recombinase. The addition of the Cre recombinase may be performed after homologous recombination by transfecting ES cells (carrying allele B) with Cre-expressing DNA or by crossing transgenic mice expressing Cre with mice having the knockout allele B, where the neomycin gene is still present. (a) Normal allele; (b) targeting construct; (c) floxed allele. Exons (coding sequences are labeled with E).



During the generation of the floxed allele the selection marker (usually neomycin also flanked by *LoxP* sites in cis) is removed by Cre. This is of special interest, because it has been reported that the neomycin promoter/enhancer may affect the expression of neighboring genes. In gene clusters, where this may occur, a bias of the phenotype is induced. The first knockout mice of the *Myf5* gene, which show a rib cage defect, illustrate clearly this phenomenon. In the newly targeted allele, however, where Cre excised the neomycin gene, *Myf5*-deficient embryos exhibit a normal rib cage (24). One of the experiments illustrating the use of *LoxP* sites is the tissue-specific inactivation of the *Pax6* gene in the eye. Mutation of the *Pax6* normally induces several defects in the eye, nose, spinal cord, and brain tissues. The tissue-specific inactivation of *Pax6* in the lens primordium indicates that it is required for lens formation and for correct placement of a single retina in the eye (25).

In addition, the *LoxP* system is now widely used to generate so called temporal knockouts in mice by using inducible site-specific recombination. For example, the engineering of a fusion protein between the Cre recombinase and the ligand-binding domain of a modified steroid receptor provides a functional inducible site-specific recombination system in mice (26, 27). The application of the

drug tamoxifen, which is still recognized by the mutant form of the estrogen receptor, efficiently induces the activation of the Cre fusion protein at the desired time of development. An alternative strategy is using a tetracycline controlled gene activation or repression protocol.

Finally, the Cre-recombinase system may be used to generate transgenic mice expressing a certain protein in a certain tissue and at a defined time of development. This is in fact necessary, when the expression of the desired gene may induce lethal defects at early stages of development.

#### 4. Potential of Transgenic Technologies

In recent years a large number of genes and gene families have been cloned. Most of them exhibit a restricted temporal and spatial expression pattern during murine development. To study the function of these genes during development, *in vitro* experiments are not sufficient, because they cannot reproduce the complex cell–cell communications and interactions of the intact organism. Transgenic technology is therefore the method of choice to analyze the various mechanisms of gene regulation and gene function. Functional analysis of the gene can be carried out in the intact animal, and new insights may be gained into such processes as tumorigenesis and disorders of cardiovascular, immune, inflammatory, and neurodegenerative natures. Furthermore, animal models may be generated for each gene suspected of having a role in disease.

##### 4.1. Gene Regulation Analysis

The analysis of gene regulation and the search for elements necessary and sufficient to reproduce the *in vivo* expression pattern of genes has profited from the pronuclear injection technology. Generally, the bacterial **b-galactosidase** gene (*LacZ*; see *lac* Operon) is inserted in the construct under the control of different DNA elements of the gene to be analyzed, and transgenic embryos are generated. The activity of the elements is followed by staining the embryos for the *LacZ* gene product using **X-gal** as a substrate, to produce a blue color. Normally, several deletion constructs and corresponding transgenic lines are necessary to determine the regulatory elements of a gene (28, 29). *LacZ* is a valuable marker that can also be inserted in the knockout construct, so that the expression of the targeted gene can be followed during the analysis of the phenotype (30)(Fig. 4).

**Figure 4. Expression pattern of the *Pax7* gene at E11 of development.** The *Pax7* gene has been inactivated by homologous recombination by inserting the b-galactosidase gene in frame to the first exon of the paired box. Mice carrying the mutation express the b-galactosidase gene under the promoter of *Pax7*, and the blue staining reflects the expression of the gene during development.





A recent technique takes the advantage of *LacZ* as a [reporter gene](#) to search for new genes involved in development. This approach has been first described using pronuclear injection. The *LacZ* gene lacking active [cis-acting](#) elements is microinjected into mouse embryos. If the *LacZ* happens to integrate at a site controlled by a strong promoter or [enhancer](#) that is active in specific cells, it will be also expressed in these cells (31). A modified technique has been designed for ES cells. Again *LacZ* is inserted between a splice acceptor site at the 5' end, and a neomycin-resistance gene at the 3' end. Upon electroporation into ES cells, the construct may happen to integrate into a gene, when  $\beta$ -galactosidase may be fused to the endogenous protein. This enzyme is the expressed under the trapped promoter, and reflects the expression pattern of the endogenous gene. The corresponding ES cells are introduced into the embryo, and the expression pattern can be followed in chimeras or in mouse transgenic lines. At the same time, the trapped gene can cloned using RACE-technique (**R**apid **A**mplification of **c**DNA **E**nds using **P**CR). In this way, a new gene is identified by its expression pattern. The insertion of the *LacZ* leads to the inactivation of the trapped gene. Thus this method, known as *gene trap*, allows the cloning of new genes and their targeted inactivation in ES cells and mice (32, 33).

#### 4.2. Functional Analysis

One of the classical methods to analyze the function of a gene is to modify its expression pattern *in vivo*, by increasing the activity of the gene product by its overexpression in transgenic mice and thus disturbing its physiological activity. The addition of genetic information, from mouse, rat, or human may result in alterations in the **phenotype**, and consequently help to define the role of the gene of interest. It may also lead to the generation of animal models for human disease, such as hypertension, arthritis, atherosclerosis, inflammatory disease, and neurodegenerative disorders (34, 35). One of the earliest transgenic mice generated by microinjection arose from the classical experiment of Palmiter et al (36). They used the [metallothionein](#) promoter to achieve widespread expression of the human [growth hormone](#) in mice. The transgenic mice were significantly larger than normal mice, confirming the role of the hormone in body growth. A related type of experiment was the use of the  $\beta$ -actin promoter upstream of the *HoxA7* coding region to express *HoxA7* at ectopic sites. Several founder animals were generated. The transgenic offspring died soon after birth and exhibited craniofacial abnormalities reminiscent of the effect of [retinoic acid](#) application during pregnancy. Also a posterior transformation of the cranial vertebrae has been observed, thus consistent with a homeotic function for this gene (37).

Expression of a certain gene in a target tissue requires the regulatory elements specific for the desired site. Promoter analysis is very time-consuming, however, and only a few genes have been analyzed in detail. The recent use of large DNA fragments to generate transgenic mice is of general interest. **Yeast artificial chromosomes** (YACs) are now used for the analysis of gene regulation and function. Such large DNA fragments are able to contain all regulatory elements, thereby ensuring expression in an integration-site-independent manner (38). The overexpression of *Pax6* using YAC transgenic technology led to a phenotype similar to that observed in a loss of function mutant and confirmed the dosage effect of this **paired-box** gene during eye development (39). The use of YAC transgenics provides a unique opportunity to analyze several basic questions in biological research. In a number of diseases where large DNA fragments are involved, YAC transgenics may now be generated. This may be used to test the hypothesis of the involvement of the **amyloid** precursor protein in Alzheimer's disease. Additionally, the molecular bases of diseases (such as Down syndrome) resulting from trisomy or other chromosomal translocations are now able to be dissected. This would be done by inserting large consecutive fragments of DNA from regions near the site of translocation to form transgenic mouse lines. Furthermore, mouse lines can be developed to express human **antibodies** for therapeutic and diagnostic purposes (40).

Before the development of homologous recombination in ES cells, dominant-negative mutations have been achieved by means of classical transgenic mice. The purpose is to alter one of the subunits by mutation so as to disrupt the function of the entire protein. One of the earliest experiments was the introduction of a transgene containing a point mutation into the pro- $\alpha$ 1(I) **collagen** gene. The mutation is identical to the one found in human corresponding gene. The transgenic mice developed a dominant-lethal phenotype similar to the human disease osteogenesis imperfecta (41).

Some investigators have utilized some toxic genes and performed some cell ablation experiments *in vivo*. The toxic genes generate a toxic product in the cell where they are expressed, or they may produce a substance that first becomes toxic after it is metabolized (42). Examples are the A-chain of the **diphtheria toxin** (DT-A) (43) and the Herpes simplex thymidine kinase gene TK mentioned above (17). In this way, cell-specific regulatory elements like the d-crystalline promoter can be used for the ablation of a specific cell type or cell population.

Homologous recombination remains the method of choice, however, to generate loss of function mutations. Several genes have been inactivated and many of them have been found to be necessary for embryonic development (44). In **gene families**, redundant functions have been confirmed by single and double mutants. This was clearly illustrated with the MyoD family of **transcription factors**. The inactivation of MyoD, Myf5, and myogenin revealed that although muscle is normally produced in mice lacking MyoD or Myf5 (45, 46), mice with null alleles for both genes produce no muscle, and no myogenin is transcribed (47). In further experiments the myogenin **complementary DNA** was inserted into the Myf5 locus by homologous recombination, disrupting its function. This knock-in of myogenin into the Myf5 locus leads to viable mice with normal rib cage. Myf5-deficient mice die prenatally and exhibit a rib cage defect (22). The knock-in experiment clearly demonstrates the functional redundancy of Myf5 and myogenin for rib formation. This kind of redundancy is also found for the **engrailed** protein (48).

Among other gene families involved in embryonic development, the *Pax* (49) and the *Hox* gene families (50) are notable. The generation of loss of function for these genes confirmed their crucial role in embryonic development. Analysis of some of the *Hox* mutants indicated that the expression of these genes is required to establish the anterior posterior axis (51-53). Analysis of the *Pax* loss of function mutants revealed not only their important function in establishing certain brain territories of the nervous system, but also that they have an essential role for the generation of different organs of the body (54). These experiments indicate that *Pax* genes are necessary for the development of the eye, the nose, the inner ear, the heart, the kidney, the pancreas, and the thyroid gland (54).

#### 4.3. Animal Models for Human Diseases

A number of spontaneous mouse mutants have been described that can be used as animal models for

human genetic disorders (55). Because the number of these mutants is limited, the homologous recombination technique provides a powerful tool to generate animal models for known human genetic disease, where it allows the inactivation of each gene in the mouse. Although transgenic mouse models are not always suitable (some neurological and vascular disorders are better studied in the rat), some already existing models are encouraging. The first animal model engineered in mice was the HPRT-negative mouse as an animal model for the Lesh–Nyhan syndrome. Unfortunately, these mice exhibit none of the neurological symptoms characteristic of that syndrome (56, 57). However, the administration of an adenine phosphoribosyltransferase (APRT) inhibitor to the HPRT-deficient mice leads to the development of the typical symptoms of the syndrome, including compulsive self-injurious behavior (58).

Cystic fibrosis (CF) is an autosomal recessive disorder affecting about one in every 2500 newborn individuals in Caucasian populations, and is caused by a defect in the chloride transport in epithelial cells. The gene encodes a protein called *Cystic fibrosis transmembrane conductance regulator* (CFTR). Several animal models have been generated in mice. Two generated by replacement vectors have identical phenotypes, and mice die perinatally of intestinal obstructions (59, 60). In contrast, a much greater percentage of the mice generated by an insertional type vector survive, and they developed pathological changes in the lung that are characteristic for human CF (61). Whitsett and colleagues (62) prevented premature death of CF mice by introducing a transgene encoding human CFTR expressed under a gut-specific promoter. Now the CF mice survive and do not exhibit intestinal obstructions, but they develop similar symptoms to the human CF. These models are useful to establish proper conditions for aerosol-based gene therapy.

The HD gene responsible for the development of the neurodegenerative disorder Huntington's disease has been knocked out to generate an animal model. The loss of function mutation is lethal, however, and the phenotype observed in heterozygotes does not reflect the expected animal model (63, 64). However, a valuable animal model was generated by the microinjection of constructs containing a portion of the mutated HD gene in humans carrying variable CAG repeat expansions into the first exon (65). This suggests that if the real cause of the disorder is unknown, it may be suitable in some cases to generate transgenic mice by the classical procedure.

In the field of cancer research, homologous recombination has permitted the generation of mice lacking the p53 protein. Mice deficient in this protein are normal but develop cancer within 6 months after birth; consequently, they represent an animal model of carcinogenesis (66). Analysis of the mutated mice revealed that the p53 gene is necessary for the **programmed cell death (apoptosis)** of cells that are proliferating abnormally. Mice deficient for the **proto-oncogene** Mdm2 gene die early in development (67, 68). The Mdm2 oncoprotein forms a complex with the p53 tumor-suppressor protein, and inhibits p53-mediated transregulation of expression. Mice deficient for both the Mdm2 and the p53 proteins develop normally and are viable, supporting the hypothesis that the critical role of Mdm2 during development is regulation of p53 function (67, 68). Mdm2-deficient mice die because the p53 function cannot be down-regulated in early stages of development. Furthermore, some of the p53-negative embryos die *in utero*, exhibiting an exencephaly. These neural tube defects are associated predominantly with female mice, and a similar sex bias is also observed in human embryos with neural tube defects (69-71). This indicates that the p53 may also have some crucial role in the closure of the neural tube. Two other **tumor suppressor genes** have been inactivated in mice, but the phenotype was different from that observed in humans; mice with the inactivated **retinoblastoma** or Wilms' tumor gene do not develop tumors (72). These results may emphasize some differences in the pathogenesis of tumors in mice and man.

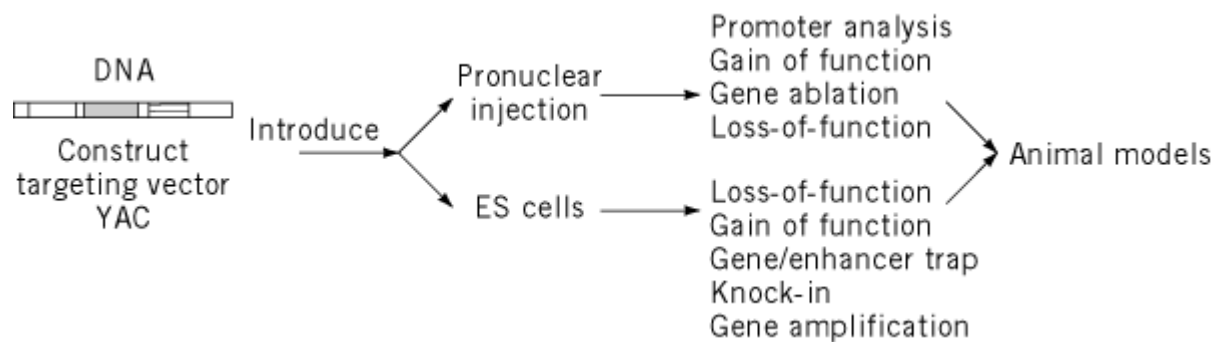
It has long been debated whether DNA **methylation** plays an important role during carcinogenesis. The generation of mice deficient in the gene for DNA **methyltransferase**, Dnmt (which is responsible for methylation of newly synthesized CpG dinucleotides) (73) can be used to test this possibility. In fact, APC<sup>Min</sup> mice carrying a point mutation for multiple intestinal neoplasia (Min) in the mouse homologue of the APC gene have been used for this purpose. Inheritance of mutations of the APC

gene in man is the cause of familial adenomatous polyposis coli. A small proportion of cases develops to a malignant state. Mice heterozygous for the APC<sup>Min</sup> allele exhibit a similar phenotype, so this represents a suitable animal model (72). Mice heterozygous for both APC<sup>Min</sup> and Dnmt deficiency show a significant reduction in the number of intestinal polyps when compared to those heterozygous for solely APC<sup>Min</sup> (74). Furthermore, the use of 5-aza-2' deoxycytidine, an inhibitor of DNA methylation, leads to a more drastic reduction of the polyp number in these double heterozygous mice (72). These results indicate how useful transgenic mice may be for the elucidation of debated hypotheses, and they may also give us new insights for the development of new concepts for cancer therapy. The efficacy of the homologous recombination technique has been recently shown by the production of a **fusion protein** that normally occurs by chromosomal translocation in humans, to produce an oncogenic protein Mll-AF9 causing acute leukemia (75).

#### 4.4. Future Perspectives

Transgenic technology has opened new possibilities in studying gene regulation and function (Fig. 5). It provides us with animal models for human diseases. Furthermore, it may give us new insights into the molecular basis of known human disorders. One such a disease is inflammatory bowel disease, called *ulcerative colitis*, which has previously been proposed to originate from an *autoimmune* response. Mice lacking interleukin-2 function develop a similar bowel inflammation, suggesting that interleukin-2 deficiency may be involved in this process (76). The site-specific recombination system and the use of YAC transgenics have opened new possibilities. Animal models can now be produced using the Cre/lox system to design subtle mutations, deletions, or even translocations (77). Temporal knockout experiments are opening new avenues in the analysis of mouse development.

**Figure 5. General scheme for the potential of the transgenic mice.** Using pronuclear injection or ES cell technology, transgenic mice can be generated. Constructs are variable and may include also YACs. Gain and loss of function mutants help to analyze gene regulation and function and provide animal models for human diseases.



Finally, the establishment of embryonic stem cells from different species including livestock has been reported. However, so far none could pass the germline, as in the mouse. The recent production of a gene-targeted sheep by nuclear transfer from cultured somatic cells (fibroblasts) provides a new tool for the introduction of a defined genetic material to create transgenic animals (78). It has opened a new field for the generation of transgenic animals from different species to use for transplantation and useful biofactories.

#### 5. Acknowledgments

The author thanks Prof. Peter Gruss for constant support and encouragement, Jens Krull for excellent technical assistance, and Prof. Peter Gruss, Dr. Kenneth Ewan and Dieter Treichel for critically reading the manuscript. I apologize that it was not possible to cite a number of colleagues due to the

limited scope of this review. This work is supported by the Max-Society.

## Bibliography

1. R. E. Hammer, V. G. Pursel, C. E. Rexroad Jr, R. J. Wall, D. J. Bolt, K. M. Ebert, R. D. Palmiter, and R. L. Brinster (1985) *Nature* **315**, 680–683.
2. B. Hogan, R. Beddington, F. Constantini, and E. Lacy (1994) *Manipulating the Mouse Embryo*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 218–252.
3. E. P. Brandon, R. L. Idzerda, and G. S. McKnight (1995) *Curr Biol* **5**, 625–634, 758–765, and 873–881.
4. M. J. Evans and M. H. Kaufman (1981) *Nature* **292**, 154–156.
5. E. J. Robertson (1987) *Teratocarcinomas and Embryonic Stem Cells*, IRL Press, Oxford, Washington, pp. 71–152.
6. A. Bradley, M. Evans, M. H. Kaufman, and E. J. Robertson (1984) *Nature* **309**, 255–256.
7. A. Gossler, T. C. Doetschman, R. Korn, E. Serfling, and R. Kemler (1986) *Proc. Natl. Acad. Sci. U.S.A.* **83**, 9065–9069.
8. K. R. Folger, E. A. Wong, G. Wahl, and M. R. Capecchi (1982) *Mol. Cell Biol.* **2**, 1372–1378.
9. O. Smithies, R. G. Gregg, S. S. Boggs, M. A. Kuralewski, and R. S. Kucherlapati (1985) *Nature* **317**, 1230–1234.
10. K. R. Thomas and M. R. Capecchi (1987) *Cell* **51**, 503–512.
11. T. C. Doetschman, R. G. Gregg, N. Maeda, M. L. Hooper, D. W. Melton, S. Thomson, and O. Smithies (1987) *Nature* **330**, 576–578.
12. M. R. Capecchi (1989) *Science* **244**, 1288–1292.
13. V. Valancius and O. Smithies (1991) *Mol. Cell Biol.* **11**, 1402–1408.
14. K. R. Thomas, C. Dend, and M. R. Capecchi (1992) *Mol. Cell Biol.* **12**, 2919–2923.
15. H. T. Riele, E. R. Maandag, and A. Berns (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 5128–5132.
16. P. L. Schwartzberg, S. P. Goff, and E. J. Robertson (1989) *Science* **246**, 799–803.
17. S. L. Mansour, K. R. Thomas, and M. R. Capecchi (1988) *Nature* **336**, 348–352.
18. R. Ramirez-Solis, H. Zheng, J. Whiting, R. Krumlauf, and A. Bradley (1993) *Cell* **73**, 279–294.
19. B. Sauer (1993) *Methods Enzymol.* **225**, 890–900.
20. B. Sauer and N. Henderson (1990) *New Biol.* **2**, 441–449.
21. H. Gu, Y-R. Zou, and K. Rajewski (1993) *Cell* **73**, 1155–1164.
22. Y. Wang, P. N. J. Schnegelsberg, J. Dansman, and R. Jaenisch (1996) *Nature* **379**, 823–825.
23. H. Gu, J. D. Marth, P. C. Orban, H. Mossman, and K. Rajewski (1994) *Science* **265**, 103–106.
24. A. Kaul, M. Küster, N. Neuhaus, and T. Braun (2000) *Cell* **102**, 17–19.
25. R. Ashery-Padan, T. Marquardt, X. Zhou, and P. Gruss (2000) *Genes Dev.* **14**, 2701–2711.
26. D. Metzger, J. Clifford, H. Chiba, and P. Chambon (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 6991–6995.
27. K. Kellendonk, F. Tronche, A-P. Monaghan, P-O. Angrand, F. Stewart, and G. Schütz (1996) *Nucleic Acids Res.*, 1404–1411.
28. A. W. Püschel, R. Balling, and P. Gruss (1990) *Development* **108**, 435–442.
29. T. Knittel, M. Kessel, M. H. Kim, and P. Gruss (1995) *Development* **121**, 1077–1088.
30. A. Mansouri, A. Stoykova, M. Torres, and P. Gruss (1996) *Development* **122**, 831–838.
31. N. D. Allen, D. G. Cran, S. C. Barton, S. Hettle, W. Reik, and A. Zurani (1988) *Nature* **333**, 852–855.
32. A. Gossler, A. L. Joyner, J. Rossant, and W. C. Skarnes (1989) *Science* **244**, 463–465.
33. W. C. Skarnes, B. A. Auerbach, and A. L. Joyner (1992) *Genes Dev.* **6**, 903–918.
34. R. Lathe and J. Mullins (1993) *Transgen. Res.* **2**, 286–299.

35. E. F. Wagner and F. Theuring (1993) *Transgenic Animals as Model Systems for Human Diseases*. E. Schering Foundation Workshop 6, Springer, Heidelberg.
36. R. D. Palmiter, R. L. Brinster, R. E. Hammer, M. E. Trumbauer, M. G. Rosenfeld, N. C. Birnberg, and R. M. Evans (1982) *Nature* **300**, 611–615.
37. R. Balling, G. Mutter, P. Gruss, and M. Kessel (1989) *Cell* **58**, 337–347.
38. A. Schedl, L. Montoliu, G. Kelsey, and G. Schüz (1993) *Nature* **362**, 258–261.
39. A. Schedl, A. Ross, M. Lee, D. Engelkamp, P. Rashbass, V. Van Heyningen, and N. D. Hastie (1996) *Cell* **86**, 71–82.
40. B. T. Lamb and J. D. Gearhart (1995) *Curr. Opin. Gen. Dev.* **5**, 342–348.
41. A. Stacy, J. Bateman, T. Choi, T. Mascara, W. Cole, and R. Jaenisch (1998) *Nature* **332**, 131–136.
42. G. A. Evans (1989) *Genes Dev.* **3**, 259–263.
43. M. L. Breitman, H. Rombola, I. H. Maxwell, G. K. K. Klintworth, and A. Bernstein (1990) *Mol. Cell Biol.* **10**, 474–479.
44. A. J. Copp (1995) *TIG* **11**, 87–93.
45. T. Braun, M. A. Rudnicki, H. H. Arnold, and R. Jaenisch (1992) *Cell* **71**, 369–382.
46. M. A. Rudnicki, T. Braun, S. Hinuma, and R. Jaenisch (1992) *Cell* **71**, 383–390.
47. M. A. Rudnicki, P. N. J. Schnegelsberg, R. H. Stead, T. Braun, H. H. Arnold, and R. Jaenisch (1993) *Cell* **75**, 1351–1359.
48. M. Hanks, W. Wurst, L. Anson-Cartwright, A. B. Auerbach, and A. Joyner (1996) **269**, 679–684.
49. A. Mansouri, A. Stoykova, and P. Gruss (1994) *J. Cell Sc. suppl.* **18**, 35–42.
50. M. Kessel and P. Gruss (1990) *Science* **249**, 374–349.
51. O. Chisaka and M. R. Capecchi (1991) *Nature* **350**, 473–479.
52. T. Lufkin, A. Dierich, M. Lé Meur, M. Mark, and P. Chambon (1991) *Cell* **66**, 1105–1119.
53. H. Le Mouellic, Y. Lallemand, and P. Brûlet (1992) *Cell* **69**, 251–264.
54. A. Mansouri, M. Hallonet, and P. Gruss (1996) *Curr. Opin. Cell Biol.* **8**, 851–857.
55. S. M. Darling and C. M. Abbott (1992) *Bioessays* **14**, 359–366.
56. M. R. Kuehn, A. Bradley, E. J. Robertson, and M. J. Evans. (1987) *Nature* **326**, 295–298.
57. M. Hooper, K. Hardy, A. Handyside, S. Hunter, and M. Monk (1987) *Nature* **326**, 292–295.
58. C. L. Wu and D. W. Melton (1993). *Nat. Genet.* **3**, 235–240.
59. J. N. Snouwaert, K. K. Brigman, A. M. Latour, N. N. Malouf, R. C. Boucher, O. Smithies, and B. H. Koller (1992) *Science* **257**, 1083–1088.
60. W. H. Colledge, R. Radcliff, D. Foster, R. Williamson, and M. J. Evans (1992) *Lancet* **340**, 680.
61. J. R. Dorin, P. Dickinson, E. W. F. W. Alton, S. N. Smith, D. M. Geddes, B. J. Stevenson, W. L. Kimber, S. Fleming, A. R. Clarke, M. L. Hooper, L. Anderson, R. S. P. Beddington, and D. J. Porteous (1992) *Nature* **359**: 211–215.
62. L. Zhou, C. R. Dey, S. E. Wert, M. D. DuVall, R. A. Frizzell, and J. A. Whitsett (1994) *Science* **266**, 1705–1708.
63. N. Jamal, S. B. Floresco, J. R. O’Kusky, V. M. Diewert, J. M. Richman, J. Zeisler, A. Borowski, J. D. Marth, A. G. Phillips, and M. R. Hayden (1995) *Cell* **81**, 811–823.
64. S. Zeitlin, J. P. Liu, D. I. Chapman, V. E. Papaioannou, and A. Efstratiadis (1995) *Na. Genet.* **11**, 155–164.
65. I. Mangiarini, K. Sathasivam, M. Seller, B. Cozens, A. Harper, C. Hetherington, M. Lawton, Y. Trotter, H. Lehrach, S. W. Davies, and G. P. Bates (1996) *Cell* **87**, 493–506.
66. L. A. Donehower, M. Harvey, B. L. Slagle, M. J. McArthur, C. A. Montgomery Jr, J. S. Butel, and A. Bradley (1992) *Nature* **356**, 215–221.

67. R. M. de Oca Luna, D. S. Wagner, and G. Lozano (1995) *Nature* **378**, 203–206.
68. S. N. Jones, A. E. Roe, L. A. Donehower, and A. Bradley (1995) **378**, 206–208.
69. J. F. Armstrong, M. H. Kaufman, D. J. Harrison, and A. R. Clarke (1995) *Curr Biol* **5**, 931–936.
70. V. P. Sah, L. D. Attardi, G. J. Mulligan, B. O. Williamson, R. T. Bronson, and T. Jacks (1995) *Nat. Genet.* **10**, 175–180.
71. S. Darling (1996) *Curr. Opin. Genet. Dev.* **6**, 289–294.
72. B. O. Williams and T. Jacks (1996) *Curr. Opin. Genet. Dev.* **6**, 65–70.
73. W. Li, T. H. Bestor, and R. Jaenisch (1992) *Cell* **69**, 915–926.
74. P. W. Laird, L. Jackson-Grusby, A. Frazeli, S. L. Dickinson, W. E. Jung, E. Li, R. A. Weinberg, and R. Jaenisch (1995) *Cell* **81**, 197–205.
75. J. Corral, I. Lavenir, H. Impey, A. J. Warren, A. Forster, T. A. Larson, S. Bell, A. N. McKenzie, G. King, and T. H. Rabbitts (1996) **85**, 853–861.
76. B. Sadlack, H. Nerz, H. Schorle, A. Schimpl, A. C. Feller, and I. Horak (1993) *Cell* **75**, 253–261.
77. R. Ramirez-Solis, P. Liu, and A. Bradley (1995) *Nature* **378**, 720–724.
78. K. J. McCreath, J. Howcroft, K. H. S. Campbell, A. Collman, A. E. Schnieke, and A. J. Kind (2000) *Nature* **405**, 1066–1069.

### Suggestions for further Reading

79. P. M. Wassarman and M. L. DePamphilis (ed.) (1993) *Guide to Techniques in Mouse Development, Section IX-XI*, *Methods in Enzymology* Vol. **25**, Academic Press, New York, pp. 719–930.
80. R. M. Torres and R. Kühn (1997) *Laboratory Protocols for Conditional Targeting*. Oxford University Press.

## Transit Peptides

Cells contain distinct subcellular compartments, such as the [nucleus](#), the [endoplasmic reticulum](#) (ER), the [Golgi apparatus](#), [peroxisomes](#), [mitochondria](#), and, in plants, [chloroplasts](#). Each of these compartments has a well-defined morphology and is surrounded by one or more limiting [membranes](#). The various functions of the different compartments depend critically on an array of diverse proteins. Because the majority of the organellar proteins are encoded in nuclear DNA and synthesized on cytosolic [ribosomes](#), efficient mechanisms must exist to transport proteins to their final destinations.

### 1. Export and Import Systems Operate in Intracellular Protein Trafficking

Two types of protein transport systems have been classified: export and import systems ([1](#)). Export systems operate in the transport of proteins from the cytosol into extracytosolic compartments. *Extracytosolic* means that the final destination of a given protein is a noncytoplasmic space, such as the lumen of the ER or the chloroplast thylakoid lumen, or a distinctive membrane, such as the peroxisomal membrane or the inner membranes of mitochondria and chloroplasts, respectively. Import systems, by contrast, transport proteins from the cytosol into an environment that is either functionally equivalent to, or evolutionarily derived from, the cytosol. Among such cytosolic-equivalent compartments are the internal compartments of mitochondria (called the *mitochondrial*

*matrix*), chloroplasts (stroma), and peroxisomes (peroxisomal matrix). The nucleoplasm of the eukaryotic cell nucleus is cytosolic as well, but transport of proteins through the nuclear pores does not obey the common principles governing protein translocation in export and import systems. Consequently, nuclear protein transport will not be discussed here (but see Ref. [2](#) for review and also see [Nuclear Import, Export](#)).

## 2. The Principal Features of transit Peptides

Export and import systems operate by similar mechanisms ([1](#)). A given polypeptide chain has to contain a molecular “stamp,” a topogenic signal ([3](#)) that dictates its proper compartmentation. Topogenic signals can reside either in the polypeptide itself (internal targeting signals, as found in almost all nuclear and peroxisomal proteins and some proteins destined to the inner membranes of mitochondria and chloroplasts) or in NH<sub>2</sub>-terminal cleavable extensions known as **signal sequences**, leader sequences, or presequences ([1](#), [3](#)) (see [Signal Peptide](#)). These sequences, which are rather short and therefore often collectively referred to as transit peptides, display a remarkable variability in terms of structure and function. However, powerful methods have been developed to identify and characterize transit peptides and to predict to which compartment a transit peptide will direct a protein ([4](#), [5](#)).

Proteins destined for the ER contain **hydrophobic** signal peptides (see Ref. [1](#) and references cited therein). These are remarkably different in length and overall amino acid composition, as found in different secretory proteins. Present evidence suggests, however, that there might be common **secondary structure** motifs that are essential for proper signal peptide function. As shown for the prepro- $\alpha$ -factor, the signal peptide consists of a hydrophobic core and two to three turns of  $\alpha$ -**helix** ([6](#)). Other signal peptides might contain similar helical structures, in which the minimal length of the hydrophobic core could be determined by the requirement to build subsequent turns of the helix (6 to 7 residues)([6](#)).

Proteins destined to the matrix of mitochondria contain mitochondrial transit peptides of about 20 to 60 amino acid residues in length (see Ref. [7](#)). Mitochondrial transit peptides are rich in positively charged amino acids but usually lack negative charges. Furthermore, mitochondrial transit peptides have the potential to form amphipathic  $\alpha$ -helices in nonaqueous environments, such as membranes. Both the basic and the amphipathic characters of mitochondrial transit peptides are essential for their function.

Proteins destined to chloroplasts contain chloroplast transit peptides. These are variable in length and amino acid composition ([8](#)). They are characterized by an uncharged NH<sub>2</sub>-terminus and a high content of hydroxylated amino acids (Ser, Thr). Chloroplast transit peptides are neither strongly basic nor amphiphilic, and their nature remains mysterious ([1](#)).

All presequences have in common that they are removed **proteolytically** during or upon transport of the higher-molecular-weight precursors to their final destinations (see Ref. [1](#) and literature cited therein). The actual cleavage sites are conserved in proteins destined to the same compartment, but they differ for secretory proteins and mitochondrial and chloroplast precursors. This reflects the different specificities of the signal, mitochondrial matrix, and chloroplast stroma processing peptidases.

## 3. Transit Sequences Mediate Protein Export and Import, Sometimes Both

Transit peptides are operative in the transport of proteins along the export pathway (eg, precursors destined to the ER or the chloroplast thylakoid lumen) and the import pathway (eg, precursors destined to the mitochondrial matrix and the chloroplast stroma). Import signals are generally [hydrophilic](#), whereas export signals appear to be hydrophobic ([1](#)).



Most transit peptides are simple in terms of structure and function: They just direct the precursor across one or two lipid bilayers (1). Some transit peptides, however, display a more complex molecular architecture. They represent bipartite presequences that contain both *import signals*, directing the precursor across the limiting outer and inner membranes of mitochondria and chloroplasts, and *export signals*, which determine the final intraorganellar protein sorting (1). For example, the import domains of mitochondrial and chloroplast precursors reside in the NH<sub>2</sub>-terminal halves of their transit peptides and are removed by the matrix and stroma processing peptidases, respectively (1). The export domains located in the carboxy-terminal halves of the transit peptides then direct the intermediate-sized proteins to their final destinations—for example, the mitochondrial inner membrane and the thylakoid lumen (1). At either place, the export domains are removed proteolytically, giving rise to the mature, properly localized protein.

#### 4. Transit Peptides Interact with Cytosolic Targeting Factors and Membrane-Bound Receptor/Translocase Complexes During Transport

##### 4.1. Cytosolic Factors with Presumed Targeting Functions

How can the targeting information, specifying to which compartment a given precursor must be transported, be decoded? There are at least two crucial elements in this decoding process. These are provided by precursor-specific targeting factors and presequence-specific receptors in the target membranes. A nascent polypeptide chain just emerging from the ribosome binds sequentially to these proteins. However, the very first partner of a nascent chain is a heteromeric protein that forms a nascent-chain-associated complex (NAC) (9-11). As shown for mammals, NAC sits on the surface of the ribosome just at the site where the nascent polypeptide chain appears. NAC binds to the chain and thereby prevents the [signal recognition particle](#), a supramolecular protein/RNA complex involved in cotranslational transport of proteins into the ER (see text below), from associating with polypeptide chains that do not contain ER-specific signal sequences. If such proteins do, however, contain such NH<sub>2</sub>-terminal signals, SRP displaces NAC from the nascent polypeptide chain and initiates co-translational protein secretion into the ER. In addition, NAC prevents spontaneous contacts between the ER and vacant ribosomes (9-11).

Remarkably, NAC of yeast has recently been demonstrated to operate also in protein transport to mitochondria (12). NAC was shown to protect nascent polypeptide chains from all cytosolic proteins until 30 to 50 amino acids have been translated. Subsequently, the mitochondrial presequence was found to become accessible to a cytosolic protein, Mft52 (13), and to the **molecular chaperone** Hsp70 that, together with DnaJ (a co-chaperone; see [DnaK/DnaJ Proteins](#)), delivers the precursor to the cytosolic face of the mitochondrial outer membrane (12).

From rabbit reticulocyte lysates, two other targeting factors have been purified. These are the presequence binding factor (14) and the targeting factor (15). They interact with typical mitochondrial transit peptides (see text above). Rat liver cytosol contains another such factor, the mitochondrial import stimulating factor (MSF) (16, 17), that recognizes and targets mitochondrial precursors to their final destination (18, 19). In addition to the aforementioned factors, a cytosolic Hsp70 molecular chaperone was suggested to operate in mitochondrial protein import (20, 21).

The relative degrees to which the different cytosolic factors determine the various targeting pathways in a cell has not been investigated. This is due in part to the fact that not all factors have been purified to homogeneity from one and the same source. On the other hand, complex regulatory circuits might exist that regulate intracellular protein trafficking at the level of targeting factor expression and interaction (see Ref. 22 for review). A general involvement of NAC in intracellular protein trafficking is likely but has not been proven in all cases (eg, for chloroplasts).

##### 4.2. Membrane-Bound Receptor/Translocase Complexes Decode Export and Import Signals

Although transit peptides per se are able to interact with their respective target membranes, the aforementioned factors promote proper binding and/or translocation. In the case of the ER (see Refs.

[23](#) and [24](#) for review), signal recognition particle (SRP) binding to the signal peptide of nascent secretory proteins causes an arrest of elongation of translation. Simultaneously, SRP pilots the ribosome-nascent chain–SRP complex to the ER membrane. At this place, a heterodimeric SRP receptor (SR) mediates docking, after which SRP transfers the signal sequence to the translocon and then cycles back to the cytosol. GTP hydrolysis appears to be involved in the regulation of targeting: SRP and SR reciprocally stimulate GTP hydrolysis. As a final step, the signal sequence inserts into the actual translocation channel formed by the Sec61p complex so that co-translational transport can begin. Recent studies have shown that Sec61p constitutes both the signal peptide receptor and the actual translocation channel ([25](#)).

The post-translational transport of secretory proteins into the ER, found in yeast, is independent of SRP but requires a signal sequence ([6](#)) and involves the same Sec61p complex as reported before, as well as an additional complex, named Sec62 / 63p ([26](#)).

Co- and post-translational protein transport into the ER differ functionally in their driving forces ([26](#), [27](#)). In the co-translational pathway, the major energy source appears to be provided by the ribosome itself. In the post-translational pathway, a luminal member of the Hsp70 family, called BiP, operates as an ATP-powered ratchet and/or import motor and pulls the precursor in Ref. [27](#). Despite this major difference, co- and post-translational translocation converge at one important point; in each case the signal sequence is presumed to cause opening of the translocation channel ([6](#), [25](#)).

In the case of mitochondria, two different receptor subcomplexes have been characterized on the outer mitochondrial membrane (see Refs. [7](#) and [28](#) for review). Translocon outer membrane protein Tom70 (the number indicates the apparent molecular mass in kilodaltons) and Tom37 form one subcomplex ([29](#)) that is involved in binding of MSF-complexed precursors, as mentioned before. Upon binding MSFs, an [ATPase](#) is activated and, after ATP hydrolysis, MSF is released. The precursor bound to Tom70/Tom37 is then transferred via another Tom, Tom5 ([30](#)), to the Tom20/Tom22 subcomplex ([31](#), [32](#)). In turn, the precursor interacts with Tom40 that, presumably together with three small Tom components, forms the actual translocation site ([33](#)). Precursors that do not need MSF for binding bypass Tom70/Tom37 and interact either directly or in a Hsp70-complexed form with the Tom22/Tom20 subcomplex ([19](#), [21](#)).

Transport of cytosolic precursors across the mitochondrial outer membrane is indeed a remarkable process. The positive charges in the basic transit peptides appear to interact sequentially with the negatively charged cytosolic domains of the different Tom proteins ([34](#)). Presumably by [electrostatic interactions](#), the precursors pass Tom20, perhaps Tom5, and Tom22. After translocation across the outer membrane, the transit peptide reaches Tim23 ([35](#)), which represents a key component of one of the two principal protein import machineries in the mitochondrial inner membrane (see Ref. [7](#) for further details). Whether the intermembrane space-exposed domain of Tom22 associates with Tim23 or forms a binding site for the positively charged or other, yet to be identified, amino acid residues of the transit peptide has not been ultimately answered ([36-39](#)). Nevertheless, once the intermembrane space-exposed domain of Tim23 has been approached, the transit peptide binds to and triggers dissociation of the Tim23 dimer ([35](#)). As a result of a presumed conformational change, the preprotein-conducting channel, which is constituted by Tim23 and Tim17 and presumably other, as-yet-undetermined components, opens ([40](#)). Then, the presequence is transferred through this Tim complex into the matrix so that it may interact with mitochondrial Hsp70 (mHsp70)([41](#), [42](#)). To accomplish its function, Hsp70 interacts with Tim44, which is associated with the Tim23/Tim17 channel (see Ref. [40](#) for details), and the soluble nucleotide exchange factor MGE ([43](#)). Initial precursor translocation requires a membrane potential  $\Delta\psi$ , but in an equilibrium reaction that is later biased toward inward movement by the binding of mHsp70 ([44](#), [45](#)). Subsequent rounds of mHsp70 binding require hydrolysis of matrix ATP and pull the precursor in.

The proper function of mitochondrial import (matrix-targeting) signals can be counteracted by hydrophobic regions immediately following the transit peptide ([46](#)). The precursor then arrests its movement and may leave the translocation channel of the inner membrane laterally into the

membrane. In addition to such a stop-transfer mechanism, reexport from the matrix has been reported for some inner membrane proteins (47). Furthermore, numerous other mechanisms have been identified in recent years that give rise to an abundance of different protein topologies (48-53). Remarkably, these various mechanisms involve transport across the common Tom machinery but diverge at the intermembrane space and/or inner membrane (48-51). While some mitochondrial inner membrane proteins are synthesized with NH<sub>2</sub>-terminal cleavable transit peptides [eg, D-lactate dehydrogenase (52)], others are targeted via internal signals [eg, carrier proteins (48-50) and the BCS1 protein (53)].

## 5. A Common Scheme of Transit Peptide Function

Transit peptides are crucial elements in the targeting pathways of cytosolic precursors to their final destinations. They can operate as either export or import signals. In either case, they interact sequentially with cytosolic proteins with presumed targeting function, cytosolically exposed receptors on their respective target membranes, and translocation channel components embedded into these membranes. The presequence initiates translocation. The driving force for the actual translocation across the membrane may be provided by the ribosome itself (co-translational import of secretory proteins into the ER) or ATP-powered molecular chaperones on the *trans* side of the target membrane (post-translational transport into the ER, mitochondria, and chloroplasts). During or upon completion of translocation, the NH<sub>2</sub>-terminal presequences are cleaved proteolytically. In case of bipartite transit peptides, the remaining part directs the intermediate-sized protein along its final targeting pathway. Recent studies have shown that thylakoid luminal proteins of chloroplasts can follow two different pathways: a DpH-dependent route (54) and a chloroplast SecA-dependent route (55); both display amazing similarities to the protein export pathways operating in bacteria (see Refs. 56 and 57 for review). Pathway specificity is determined by information residing in the carboxy-terminal halves of the transit peptides, which function as export signals (58, 59).

## 6. Exceptions to the Rules: Transit Peptides with Regulatory Functions

Another interesting facet of transit peptide function was recently discovered: Porphyrins such as heme, which is found ubiquitously in all living organisms, can regulate mitochondrial protein import (60). A conserved heme-binding site was identified in the presequence of 5-aminolevulinic acid synthase, which represents a key enzyme of heme biosynthesis. Binding of heme to this motif blocked import of the cytosolic precursor of this enzyme, and this is assumed to be a key element through which heme biosynthesis is controlled in mitochondria (60). Just the opposite mode of regulation was discovered in case of a key enzyme of chlorophyll biosynthesis of higher plants, the NADPH:protochlorophyllide oxidoreductase (PORA). One of the enzyme's two substrates, protochlorophyllide, bound to a yet unidentified porphyrin binding site in the transit peptide and triggered the import of the cytosolic precursor into chloroplasts (61).

## Bibliography

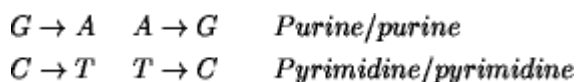
1. G. Schatz and B. Dobberstein (1996) *Science* **271**, 1519–1526.
2. D. Goerlich and I. W. Mattaj (1996) *Science* **271**, 1513–1518.
3. G. Blobel (1980) *Proc. Natl. Acad. Sci. USA* **77**, 1496–1500.
4. H. Nielsen, J. Engelbrecht, S. Brunak, and G. von Heijne (1997) *Protein Eng.* **10**, 1–6.
5. M. C. Claros, S. Brunak, and G. von Heijne (1997) *Curr. Opin. Struct. Biol.* **7**, 394–398.
6. K. Plath et al., submitted for publication.
7. W. Neupert (1997) Protein import into mitochondria. *Annu. Rev. Biochem.* **66**, 863–917.
8. von Heijne et al. (1991) *Plant Mol. Biol. Rep.* **9**, 104–126.
9. B. Wiedmann, H. Sakai, T. A. Davis, and M. Wiedmann (1994) *Nature* **370**, 434–440.
10. B. Lauring, H. Sakai, G. Kreibich, and M. Wiedmann (1995) *Proc. Natl. Acad. Sci. USA* **92**, 5411–5415.

11. Proc. Natl. Acad. Sci. USA **92**, 9435–9439.
12. R. George, T. Beddoe, K. Landl, and T. Lithgow (1998) Proc. Natl. Acad. Sci. USA **95**, 2296–2301.
13. P. Cartwright, T. Beilharz, P. Hansen, J. Garrett, and T. Lithgow (1997) J. Biol. Chem. **272**, 5320–5325.
14. K. Murakami and M. Mori (1990) EMBO J. **9**, 3201–3208.
15. H. Ono and S. Tuboi (1990) Arch. Biochem. Biophys. **280**, 299–304.
16. N. Hachiya et al. (1993) EMBO J. **12**, 1579–1586.
17. N. Hachiya et al. (1994) EMBO J. **13**, 5146–5154.
18. T. Komiya et al. (1994) J. Biol. Chem. **269**, 30893–30897.
19. N. Hachiya et al. (1995) Nature **376**, 705–709.
20. T. Komiya, M. Sakaguchi, and K. Mihara (1996) EMBO J. **15**, 399–407.
21. T. Komiya, S. Rospert, G. Schatz, and K. Mihara (1997) EMBO J. **16**, 4267–4275.
22. J. Nunnari and P. Walter (1996) Regulation of organelle biogenesis. Cell **84**, 389–394.
23. T. A. Rapoport, B. Jungnickel, and U. Kutay (1996) Annu. Rev. Biochem. **65**, 271–303.
24. A. K. Corsi and R. Schekman (1996) J. Biol. Chem. **271**, 30299–30302.
25. W. Mothes, B. Jungnickel, J. Brunner, and T. A. Rapoport (1999) J. Biol. Chem. **142**, in press.
26. K. E. Matlack, W. Mothes, and T. A. Rapoport (1998) Cell **92**, 381–390.
27. K. E. Matlack, B. Misselwitz, and T. A. Rapoport, submitted for publication.
28. N. Pfanner, E. A. Craig, and A. Hoenlinger (1997) Mitochondrial preprotein translocase. Annu. Rev. Cell Dev. Biol. **13**, 25–51.
29. S. Gratzler et al. (1995) J. Cell Biol. **129**, 25–34.
30. K. Dietmeier et al. (1997) Nature **388**, 195–200.
31. A. Hoenlinger et al. (1995) Mol. Cell. Biol. **15**, 3382–3389.
32. A. Mayer, F. E. Nargang, W. Neupert, and R. Lill (1995) EMBO J. **14**, 4204–4211.
33. K.-P. Kuenkele et al. (1998) Cell **93**, 1009–1019.
34. T. Komiya et al. (1999) EMBO J. **17**, in press.
35. M. F. Bauer, C. Sirrenberg, W. Neupert, and M. Brunner (1996) Cell **87**, 33–41.
36. D. A. Court et al. (1996) Mol. Cell. Biol. **16**, 4035–4042.
37. F. E. Nargang et al. (1998) Mol. Cell. Biol. **18**, 3173–3181.
38. J. Brix, K. Dietmeier, and N. Pfanner (1997) J. Biol. Chem. **272**, 20730–20735.
39. M. Moczko et al. (1997) Mol. Cell. Biol. **17**, 6574–6584.
40. U. Boemer et al. (1997) EMBO J. **16**, 2205–2216.
41. P. J. Kang et al. (1990) Nature **348**, 137–143.
42. P. E. Scherrer et al. (1990) EMBO J. **9**, 4315–4322.
43. H.-C. Schneider, B. Westermann, W. Neupert, and M. Brunner (1996) EMBO J. **15**, 5976–5803.
44. J. Martin, K. Mahlke, and N. Pfanner (1991) J. Biol. Chem. **266**, 18051–18057.
45. C. Ungermann, B. Guiard, W. Neupert, and D. M. Cyr (1996) EMBO J. **15**, 735–744.
46. A. Gruhler et al. (1997) J. Biol. Chem. **272**, 17410–17415.
47. K. Hell et al. (1998) Proc. Natl. Acad. Sci. USA **95**, 2250–2255.
48. C. Sirrenberg et al. (1996) Nature **384**, 582–585.
49. C. Sirrenberg et al. (1998) Nature **391**, 912–915.
50. C. M. Koehler et al. (1998) Science **279**, 369–373
51. O. Kerscher et al. (1997) J. Cell Biol. **139**, 1663–1675.
52. E. E. Rojo, B. Guiard, W. Neupert, and R. A. Stuart (1998) J. Biol. Chem. **273**, 8040–8047.

53. H. Foelsch, B. Guiard, W. Neupert, and R. A. Stuart (1996) *EMBO J.* **15**, 479–487.
54. A.M. Settles et al. (1997) *Science* **278**, 1467–1470.
55. J. Yuan, R. Henry, M. McCaffery, and K. Cline (1994) *Science* **266**, 796–798.
56. K. Cline and R. Henry (1996) Import and routing of nucleus–encoded chloroplast proteins. *Annu. Rev. Cell Dev. Biol.* **12**, 1–26.
57. C. Robinson and A. Mant (1997) Targeting of proteins into and across the thylakoid membrane. *Trends Plant Sci.* **2**, 431–437.
58. E. Bogsch, S. Brink, and C. Robinson (1997) *EMBO J.* **16**, 3851–3859.
59. R. Henry et al. (1997) *J. Cell Biol.* **136**, 823–832.
60. J. T. Lathrop and M. P. Timko (1993) *Science* **259**, 522–525.
61. C. Reinbothe, N. Lebedev, K. Apel, and S. Reinbothe (1997) *Proc. Natl. Acad. Sci USA* **94**, 8890–8894.

## Transition Mutation

A transition mutation is a [base-pair-substitution](#) point [mutation](#) that substitutes a **purine** for the other purine or a **pyrimidine** for the other pyrimidine, so that the purine/pyrimidine axis of the DNA molecule is maintained. The four possible transitions are



The term transition was first defined by Ernst Freese in 1959 ([1](#)). By studying mutations induced by the **mutagenic** base analog [5-bromouracil](#) or [2-aminopurine](#), he showed that there are two types of [mutagenesis](#), those that undergo **reversion** by the base analogues, called transitions, and those that do not, called **transversions**. Transitions result from base mispairing due to the ambiguous base-pairing properties of base analogues. 2-Aminopurine mispaired with a thymine results in the insertion of an adenine in the next round of [DNA replication](#), creating a G to A transition. Because these mutations were due to mispairing by the base analog, it was reasoned that the mutations would be revertible with another base analogue.

Transitions are the errors that occur most frequently during DNA replication. **DNA polymerases** are far more likely to misinsert the wrong purine or pyrimidine than to insert a purine instead of a pyrimidine or vice versa. Transitions are also induced by a number of exogenous and endogenous DNA-damaging agents that behave like base analog *in vivo*. Such mutagens include alkylating and some oxidative agents.

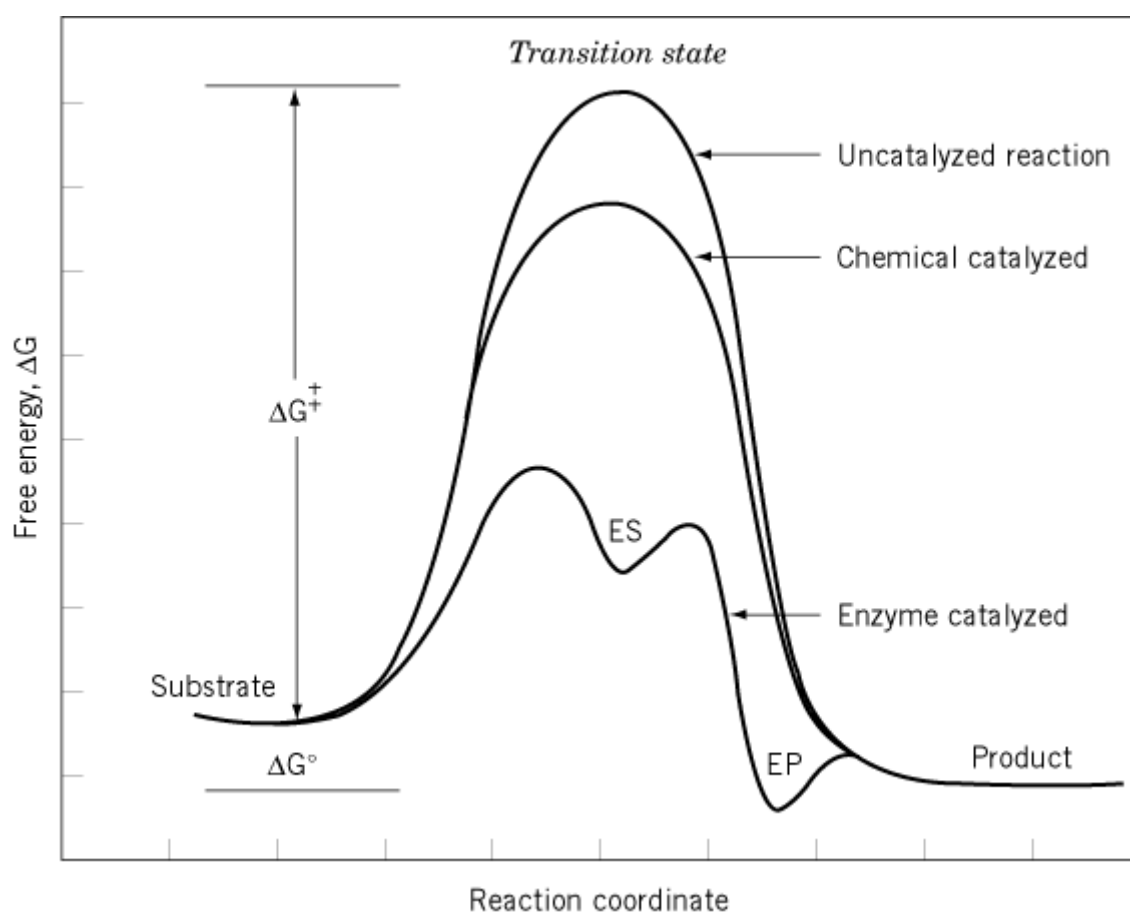
## Bibliography

1. E. Freese (1959) *Proc. Natl. Acad. Sci. USA* **45**, 622–633.

## Transition State

There is no necessary connection between how much and how fast a reaction occurs (see [Kinetics](#)). The more stable the products are relative to the starting reactants, the more complete is the reaction. A **free energy** diagram of the reaction indicates why this is so (Fig. 1). The difference in free energy between the reactant and the product determines the equilibrium constant for the reaction. A large and favorable free energy difference does not, however, ensure that the reaction will occur rapidly. The rate of a reaction is determined by an independent variable, the free energy of the transition state, which is the least stable species occurring during the reaction. If the transition state has a low free energy, the reaction occurs rapidly. The higher the free energy, the slower the reaction. Within the transition state for chemical reactions, covalent bonds are being broken and made.

**Figure 1.** An example of a reaction coordinate diagram. The reaction coordinate is a measure of the extent to which a chemical reaction has occurred in a molecule. The starting molecule is on the left, the product on the right. The free energy of the molecule is given by the solid line. The species with the highest free energy is the transition state for the reaction. The higher its free energy, the slower the reaction. Chemical catalysis is depicted as occurring because the transition state is stabilized. Such stabilization shown for enzymatic catalysis is caused by binding of the reactants to the enzyme.



According to the simplest form of transition-state theory, the rate constant  $k_r$  for a reaction is defined by the free energy  $DG^\ddagger$  of the transition state, in the following equation:

$$k_r = \left( \frac{k_B T}{h} \right) \exp(-\Delta G^\ddagger / RT) \quad (1)$$

where  $k_B$  is the Boltzmann constant,  $T$  is the absolute temperature,  $h$  is Planck's constant, and  $R$  is the gas constant. This equation was derived by assuming that the transition state is in equilibrium with the reactants and that molecules reaching the transition state are converted to products by the most rapid process in which covalent bonds can be broken and made, given by  $k_B T/h$ ; it has a value of  $6 \times 10^{12} \text{s}^{-1}$  at  $25^\circ\text{C}$ .

A related concept is the [activation energy](#), which corresponds to the enthalpy of the transition state. [Enzymes](#) and other catalysts can be imagined to increase the rates of reactions by lowering the free energy of the transition state, generally by interacting with it and stabilizing it.

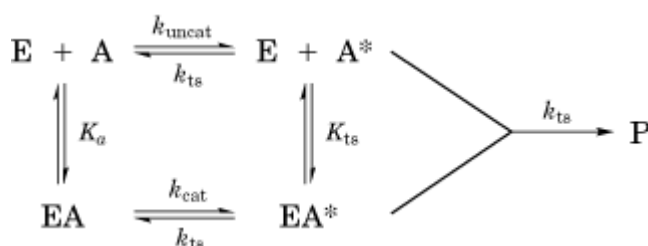
Because they are the least stable species during a reaction, transition states are populated for the least amount of time and are not detectable directly. Their natures can be inferred only by studying the effects of systematically varying the conditions of a reaction or the reactants and examining their effects on the rate of the reaction, i.e., the free energy of the transition state (see [Free Energy Relationships](#)).

#### Suggestion for Further Reading

K. J. Laidler and P. S. Bunting (1973) *The Chemical Kinetics of Enzyme Action* 2nd ed., Oxford University Press, London, pp. 40–50.

### Transition State Analogue

Like all **catalysts**, [enzymes](#) increase the rate of chemical reactions by lowering the [activation energy](#). The result is that the substrate molecules pass through the activated or [transition state](#) at a greater rate than would be the case if the reaction were uncatalyzed (see [Catalysis](#)). The features of an enzyme-catalyzed reaction and the corresponding nonenzymic reaction can be described in terms of transition kinetic theory ([1](#)) by the following scheme:



In this model, the enzyme (E) has only a single substrate, which is converted to a single product (P), and the release of P from the enzyme is not rate-limiting. The rates of formation of the transition state, designated as  $\text{A}^*$ , when bound to the enzyme and when in solution, are given by  $k_{\text{cat}}$  and  $k_{\text{uncat}}$ , respectively. A transition state breaks down to the original ground state or to the product at the same intrinsic rate, independent of the enzyme. In classical transition state theory, these rates are represented by the rate constant  $k_{\text{ts}}$ , which is equal to  $kT/h$ , where  $k$  is Boltzmann's constant,  $h$  Planck's constant, and  $T$  the absolute temperature.  $K_a$  and  $K_{\text{ts}}$  denote **dissociation constants** for the EA and EA\* complexes, respectively. From the kinetic scheme, it follows that the rate of the reaction on the enzyme, relative to the rate of the reaction in solution, is given by  $k_{\text{cat}}/K_a$ , which must be equal to  $k_{\text{uncat}}/K_{\text{ts}}$ , as these are **linked functions** and the **free energy** must not change in a

cyclic series of reactions. To the extent that  $k_{\text{cat}}$  is greater than  $k_{\text{uncat}}$ , the value of  $K_{\text{ts}}$  must be less than  $K_a$ ; ie, the transition state should have a correspondingly higher affinity for the enzyme than the substrate. As enzymes speed up the rate of reactions by as much as  $10^{12}$ -fold, the binding of the transition state by the enzyme should be very much tighter than the binding of the substrate in the ground state. The above relationship can be expressed as  $K_{\text{ts}} = (k_{\text{uncat}}/k_{\text{cat}}) K_a$ , so that, for an enzyme that increases  $10^{12}$ -fold the rate of reaction of a substrate with a  $K_a$  value of  $10^{-5}\text{M}$ , the value for  $K_{\text{ts}}$  would be  $10^{-17}\text{M}$ . Such a low value for  $K_{\text{ts}}$  would ensure that the transition state is not released from the enzyme. The enhanced binding of the transition state would be due to the additional and favorable interactions between it and the enzyme that are not available to the ground state substrate.

The idea of tight-binding of the transition state prompted the development of *transition state analogues* with structural features approaching those of the transition state. It was expected that such compounds would be highly specific and potent inhibitors of enzymes. But it would be difficult, if not impossible, to prepare a good stable analogue of an extremely unstable transition state, which must have an unusual covalent structure with partial and extended bonds and unusual bond angles. Most of the compounds referred to as transition state analogues are, in effect, *intermediate state analogues*, as they are analogues of actual intermediates, such as carbanions and tetrahedral adducts (2, 3). An enzyme would not be expected to bind an intermediate state as tightly as it would a transition state. Examples of intermediate state analogues that bind strongly to enzymes are listed in Table 1.

**Table 1. Selected Examples of Intermediate State Analogues**

| Enzyme                                   | Inhibitor                                  | $K$     | Reference |
|--|--|---------|-----------|
| Aconitase                                | 1-hydroxy-2-nitro-1,3-propanedicarboxylate | 0.7 nM  | (5)       |
| Adenosine deaminase                      | 2'-deoxycoformycin                         | 0.22 nM | (6)       |
| Adenylate deaminase                      | Coformycin-5'-phosphate                    | 55 pM   | (6)       |
| Cytidine deaminase                       | Phosphapyrimidine nucleoside               | 0.9 nM  | (7)       |
| Enolase                                  | (3-hydroxy-2-nitropropyl)-phosphonate      | 6.0 nM  | (8)       |
| Isocitrate lyase                         | 3-nitropropionate                          | 1.5 nM  | (2)       |
| <b>Pepsin</b>                            | Pepstatin                                  | 46 pM   | (9)       |
| <b>Ribulose bisphosphate carboxylase</b> | 4-carboxy-D-arabinitol-1,5-bisphosphate    | 5 pM    | (10)      |
| <u>Thermolysin</u>                       | Phosphonoamidate peptide                   | 9.1 nM  | (11)      |

Finding an intermediate state analogue that functions as a potent inhibitor and has been developed on the basis of the postulated chemical mechanism for the reaction could be taken as evidence for the proposed mechanism. It is possible however, that the strong interaction is due simply to fortuitous interactions. Methotrexate, an analogue of dihydrofolate, is a potent inhibitor of [dihydrofolate reductase](#) from *Escherichia coli* (see [Aminopterin](#), [Methotrexate](#), [Trimethoprim](#), and [Folic Acid](#)) and shows other characteristics of the behavior of intermediate state analogues (see below). But the strong inhibition of dihydrofolate reductase by methotrexate occurs not because it behaves as an



intermediate state analogue, for it is bound upside down relative to the binding of dihydrofolate, but because there are more interactions between the inhibitor and the enzyme than between the enzyme and the substrate.

A feature of intermediate state analogues is that they often exhibit slow-binding characteristics (4) (see [Slow-Binding Enzyme Inhibition](#)). This time-dependent behavior is often due to the rapid formation of an enzyme-inhibitor complex that subsequently undergoes a slow conformational change, leading to marked enhancement of the binding of the inhibitor. This behavior could well be analogous to the conformational change that an enzyme-substrate complex is believed to undergo in the formation of the transition state complex. It could be considered that, in both instances, the conformational changes are responsible for the stronger interactions and tighter binding.

## Bibliography

1. R. Wolfenden (1976) *Ann. Rev. Biophys. Bioeng.* **5**, 271–306.
2. J. V. Schloss and W. W. Cleland (1982) *Biochemistry* **21**, 4420–4427.
3. J. F. Morrison and W. W. Cleland (1983) *Biochemistry* **22**, 5507–5513.
4. J. V. Schloss (1988) *Acc. Chem. Res.* **21**, 348–353.
5. J. V. Schloss, D. J. T. Porter, H. J. Bright, and W. W. Cleland (1980) *Biochemistry* **19**, 2358–2362.
6. C. Frieden, L. C. Kurtz, and H. R. Gilbert (1980) *Biochemistry* **19**, 5303–5309.
7. G. Ashley and P. A. Bartlett (1984) *J. Biol. Chem.* **259**, 13621–13627.
8. V. E. Anderson, P. M. Weiss, and W. W. Cleland (1984) *Biochemistry* **23**, 2779–2786.
9. D. H. Rich and E. T. O. Sun (1980) *Biochem. Pharmacol.* **29**, 2205–2212.
10. J. Pierce, N. E. Tolbert, and R. Barker (1980) *Biochemistry* **19**, 934–942.
11. P. A. Bartlett and C. K. Marlowe (1983) *Biochemistry* **22**, 4618–4624.

## Translation

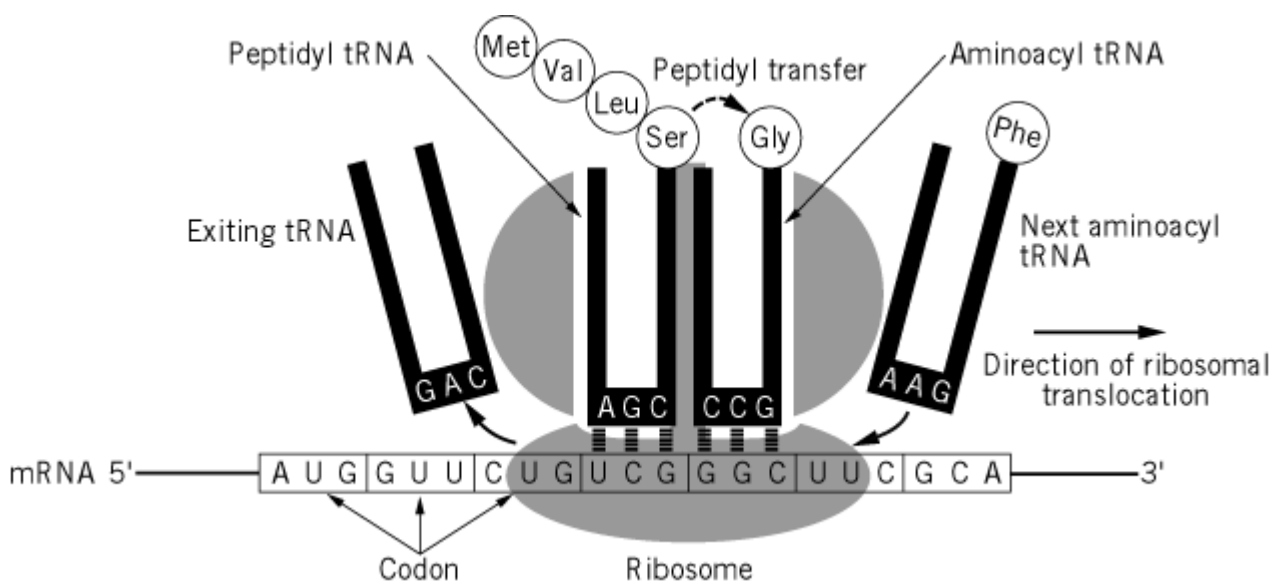
Translation is the synthesis of a [protein](#) according to the genetic information. It is the second step of **gene** expression, following [transcription](#), and is a universal and essential step for life. The process is called translation because it converts the genetic information of nucleic acid sequences, composed of four distinct nucleotides, to the polypeptide sequence composed of 20 distinct amino acids. A unit of the genetic code consisting of three nucleotides (triplet), named a **codon**, corresponds to one amino acid, except for the final [stop codon](#). The conversion of information is often referred to as “decoding.” Being highly accurate (with an error rate of only  $5 \times 10^{-4}$  per amino acid residue) and fast (15 amino acids incorporated per second at 37 °C), translation is one of the most complex biological processes; it involves a number of factors in a coordinate manner and consumes a vast proportion (about 80%) of energy that cells produce. That translation of all the contemporary [genomes](#) uses a common basic mechanism is a strong indication that all organisms originated from a single evolutionary ancestor.

The template for translation is a [messenger RNA](#) that encodes the genetic code as a sequence of 64 (=  $4^3$ ) different codons. The linear order of codons aligned on the mRNA from the 5' end to the 3' end in a nonoverlapping and nongapped manner specifies the amino acid sequence of the polypeptide chain from the *N*-terminus to the *C*-terminus. The principal reaction of translation is the polymerization (by [peptide bonds](#)) of amino acids arrayed according to the codons. The amino acids

are not lined up by direct interaction with the codons, however, but each of them is first bound to an adapter molecule, [transfer RNA](#), and aligned by the interaction between the anticodon of the tRNA and the codon of the mRNA. [Aminoacyl-tRNA synthetases](#) catalyze the covalent attachment of the **α-carboxyl group** of the amino acid to the 3' terminus of the appropriate tRNA, to form aminoacyl tRNA. This reaction couples the hydrolysis of ATP to AMP to charge the amino acid with the chemical energy for peptide bond formation. The accuracy of the aminoacyl-tRNA synthetase is one of the determinant factors for the translational fidelity. The polymerization of amino acid proceeds from the *N*-terminus to the *C*-terminus. For each step in the reaction, the substrates are the elongating peptide and a new amino acid, both of which are attached to tRNAs. The actual polymerization is, therefore, transfer of the elongating peptide from one tRNA to the amino group of the new amino acid-tRNA complex, referred to as **peptidyl transfer**.

Translation is carried out on the intracellular particles called [ribosomes](#) (Fig. 1). Each ribosome consists of three distinct molecules of ribosomal RNAs and many proteins and is made up of two subunits, namely the large and small subunits (see [Ribosomes](#)). The ribosome has two major functions. It binds to mRNA and tRNAs, providing the sites of specific codon-anticodon interaction. It also catalyzes the peptidyl transfer. There are three tRNA-binding sites on the ribosome; two of them, called the A site and the P site, are directly involved in decoding of mRNA, while the remaining one, called the E site, is used as an exit site after the decoding and peptidyl transfer reaction. The P site of the ribosome accommodates the tRNA with the elongating peptide (peptidyl-tRNA), and the A site accepts the aminoacyl-tRNA that correctly base-pairs with the corresponding codon. Binding of the aminoacyl-tRNA is followed by peptidyl transfer. The catalytic center of the reaction is in the large subunit of the ribosome and likely to be a **ribozyme** encoded by domain V of the large rRNA (see [Peptidyl Transferase](#)). The 3' termini of the two tRNAs at the P- and A-sites come to face the catalytic center. The peptidyl transfer reaction leaves deacylated-tRNA at the P site and leaves the newly elongated peptidyl-tRNA at the A site. The ribosome then proceeds toward the 3' end of the mRNA by 3 nucleotides. This [translocation](#) is concomitant with the simultaneous shifts of the peptidyl-tRNA from the A site to the P site and of the deacylated-tRNA from the P site to the E site (see [Translocation](#)). On binding of a new aminoacyl-tRNA to the A site, the deacylated-tRNA dissociates from the E site of the ribosome.

**Figure 1.** Schematic representation of protein synthesis on the ribosome.



The translation process is divided into three steps, namely initiation, elongation, and termination. Initiation not only refers to decoding of the [initiation codon](#) (which is normally the AUG codon of mRNA), but also includes all the processes prerequisite for the first peptide bond formation. It is more diverse between prokaryotes and eukaryotes than the other two steps. The small subunit of the ribosome is bound to mRNA with the assistance of protein factors named [initiation factors](#). In bacteria, binding occurs near the initiation codon through interaction between the small subunit (16 S) rRNA and a specific short sequence of mRNA, the [Shine–Dalgarno sequence](#). For most of the eukaryotic mRNAs, however, the initial binding site of the translational apparatus is the Cap structure of the 5' terminus of mRNA (see [Cap](#)), and the whole 5'-untranslated region is scanned by the small subunit (see [Scanning Hypothesis](#)). In both cases, the small subunit is associated with the large subunit at the initiation codon (see [Initiation Complex](#)). The tRNA that decodes the initiation codon is not used for elongation. Bacterial initiator tRNA is tRNA<sup>f</sup> bound to *N*-formyl methionine, and eukaryotic initiator tRNA is referred as tRNA<sup>iMet</sup> (see [fMET \(N-Formyl Methionine\)](#)). The *N*-formyl groups of mature bacterial peptides, as well as the *N*-terminal methionine residues of many bacterial and eukaryotic peptides, are removed co-translationally. Initiation of protein synthesis constitutes a rate-limiting step in translation and hence provides a main target of regulation in many genes of prokaryotes and eukaryotes (see [Initiation Complex](#)).

Elongation is a reaction cycle that repeats the three processes for each amino acid: (i) decoding of aminoacyl-tRNA at the A site of ribosome, (ii) peptidyl transfer, and (iii) translocation of the ribosome. Two indispensable [GTP-binding proteins](#), called [elongation factors](#), are involved in these processes. Only the aminoacyl-tRNA associated with EF-Tu (in bacteria) or EF-1a (in eukaryotes) and with GTP, in a ternary complex, can bind to the A site of the ribosome. When aminoacyl-tRNA correctly base-pairs with the codon, hydrolysis of GTP to GDP occurs, and the GDP form of EF-Tu and EF-1a dissociates the ribosome, leaving the aminoacyl-tRNA at the A site for the subsequent peptidyl transfer reaction. After peptidyl transfer, the GTP form of EF-G (for bacteria) or EF-2 (for eukaryotes) then binds to the ribosome and translocates the peptidyl-tRNA from the A site to the P site, making the A site vacant for the next round of protein elongation. Apparently, hydrolysis of GTP to GDP serves as an energy source for translocation. Mutually exclusive binding of EF-Tu (EF-1a) and EF-G (EF-2) ensures orderliness of the elongation cycle. Recycling of the GDP forms of EF-Tu and EF-1a are catalyzed by the GDP-GTP exchange factors, EF-Ts (in bacteria) and EF-1b EF-1g (in eukaryotes). They are also classified as elongation factors.

Termination of protein synthesis takes place on the ribosome in response to a stop, rather than a sense, codon in the “decoding” site. Translation termination requires two classes of polypeptide release factors: (i) a class I factor, codon-specific release factors (RF-1 and RF-2 in prokaryotes; eRF-1 in eukaryotes), and (ii) a class II factor, nonspecific release factors (RF-3 in prokaryotes; eRF-3 in eukaryotes), that binds guanine nucleotides (see [Release Factor](#)). Release factors seem to mimic tRNA and recognize directly the stop codon. When release factors correctly recognize the stop codon on the ribosome, the ribosomal peptidyl transferase center (domain V of the large rRNA) changes the catalytic mode and hydrolyzes the peptidyl-tRNA, releasing the polypeptide chain. After polypeptide release, the final dissociation of the termination complex takes place by the action of the ribosome-recycling factor (called RRF in bacteria) (see [Termination Factor](#)), allowing reuse of the ribosome, release factors, and tRNAs for the next round of protein synthesis.

Translation can be studied *in vitro*. Supernatants of low-speed centrifugation of cell lysates show translation activity dependent on exogenously added mRNAs and upon energy supplementation. Products from endogenous mRNAs can be reduced by a ribonuclease treatment. Widely used *in vitro* translation system are rabbit reticulocyte lysates, wheat germ extracts, and the S30 fraction of *Escherichia coli* extracts. Individual steps of translation may be examined in a reconstitution system with purified ribosomes and other components.

#### Suggestions for Further Reading

M. B. Mathews, N. Sonenberg, and J. W. B. Hershey (1996) In *Translational Control* (J. W. B.

Hershey, M. B. Mathews, and N. Sonenberg, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, pp. 1–29.

D. E. Draper (1996) In *Escherichia coli and Salmonella*, 2nd ed. (R. Curtiss III et al., eds.), American Society for Microbiology Press, Washington, D. C., pp. 902–908.

C. G. Kurland, D. Hughes, and M. Ehrenberg (1996) In *Escherichia coli and Salmonella*, 2nd ed. (R. Curtiss III et al., eds.), American Society for Microbiology Press, Washington, D.C., pp. 979–1004.

Y. Nakamura, K. Ito, and L. A. Isaksson (1996) *Cell* **87**, 147–150.

Y. Nakamura and K. Ito (1998) *Genes Cells* **3**, 265–278.

K. S. Wilson and H. F. Noller (1998) *Cell* **92**, 337–349.

## Translation Repressors

[Translation](#) of [messenger RNA](#) into [protein](#) represents an important regulatory step in **gene** expression. **Protein biosynthesis** involves three distinct stages: initiation, elongation, and termination. In eukaryotes, formation of the mRNA–[ribosome](#) initiation complex represents the most critical step for control of translation ([1](#)). The translation machinery is responsive to a wide range of intracellular and extracellular signals. For example, under certain conditions, such as **virus infection** or **heat shock**, there is preferential translation of viral and heat-shock mRNAs, respectively, while translation of most **5'-capped**, cellular mRNAs is inhibited ([2](#), [3](#)). Maternal mRNA during early development, and certain newly synthesized mRNA transcripts in somatic cells, are transported to the cytoplasm, where they are sequestered from the translational machinery and remain untranslated until the encoded proteins are needed ([4](#)).

Translation repressors encompass both [cis-acting](#) mRNA regulatory sequences and [trans-acting](#) protein factors. The best characterized repressors act at the level of translation initiation and fall into four classes: (1) [repressor](#) molecules that target the eukaryotic [initiation factors](#) (eIFs) required to assemble mRNA–ribosome initiation complexes, (2) specific ribonucleotide sequences and structural motifs within the mRNA template that interfere with the translational machinery, (3) specific and nonspecific mRNA-binding proteins, and (4) complementary RNAs that inhibit translation through the formation of heteroduplexes with the mRNA template.

### 1. Inhibition of Translation Initiation Factors

Translation is dependent on the productive association of mRNA with ribosomes. Recruitment of the 40 S ribosomal subunit to mRNA requires ATP and involves the participation of several translation initiation factors ([1](#)). Functional inactivation of these initiation factors disrupts this process. This can occur by a variety of mechanisms, including **phosphorylation**, interaction with inhibitory polypeptides, and **proteolytic** cleavage.

eIF-2 participates in the formation of the ternary complex involving the initiator methionine [transfer RNA](#) and GTP. eIF2 is required for the delivery of the ternary complex eIF2:tRNA(Met):GTP to the 40 S ribosome subunit before mRNA is bound. eIF-2 contains three subunits. Phosphorylation of the  $\alpha$ -subunit of eIF-2 blocks the exchange of GDP for GTP required for the catalytic utilization of eIF2. Three protein [kinases](#), hemin-regulated initiation factor-2 kinase (HRI), double-stranded RNA-dependent kinase (PKR), and yeast protein kinase GCN2, specifically phosphorylate eIF-2a under certain physiological conditions ([5-7](#)). Viral infection induces PKR autophosphorylation and kinase

activation, which lead to phosphorylation of eIF-2a and subsequent repression of host cell protein synthesis (8). On heat shock, eIF-2a is initially phosphorylated by HRI, resulting in a general inhibition of protein synthesis and preferential translation of heat-shock protein mRNAs. During heat shock recovery, the level of eIF-2a phosphorylation decreases as cellular protein synthesis is restored (3). In yeast, GCN2 phosphorylates eIF-2a in response to amino acid starvation, resulting in a decrease in the synthesis of total cellular proteins, but a specific increase in the translation of GCN4 mRNA (9).

eIF-4E and eIF-4G are important components of the eIF-4F complex. eIF-4E binds directly to the 5' cap and plays a critical role in the formation of 48 S mRNA–ribosome preinitiation complex. eIF-4E is present in the cell in limiting amounts relative to other initiation factors. The phosphorylation state of eIF-4E correlates well with the rate of translation in the cell, and overexpression of eIF-4E has been shown to affect cell growth (10). A family of inhibitory proteins have been shown to inhibit translation by interacting with eIF-4E to prevent assembly of the eIF-4F complex. These proteins are called eIF-4E-binding proteins (4E-BPs) (11). The activity of 4E-BPs are also regulated by phosphorylation. Viral infection induces the dephosphorylation of 4E-BPs, resulting in the enhanced association of 4E-BPs with eIF-4E and inhibition of host mRNA translation. In contrast, exposure to [growth factors](#) promotes phosphorylation of 4E-BPs, resulting in dissociation from eIF-4E and translational activation (11).

Following [poliovirus](#) infection, eIF-4G is proteolytically cleaved by the virus-encoded proteinase 2A<sup>pro</sup>. As a result, translation of most capped cellular mRNA is inhibited, while translation of uncapped viral mRNAs is stimulated (12).

## 2. *Cis*-Acting Regulatory RNA Sequences

Most eukaryotic mRNAs are composed of a 5' cap, a 5'-untranslated region (5'-UTR), a central amino acid coding region, and a 3' untranslated region (3'-UTR), followed by a **poly(A)** tail. Translational initiation of most capped mRNAs follows the [scanning hypothesis](#) that was originally proposed by Kozak and Shatkin (13). In this model, the translation initiation complex forms at or near the 5' cap and scans the mRNA toward the 3' end in search of the first [initiation codon](#) that is situated within an appropriate sequence context. The efficiency of translation varies markedly for different mRNAs and is influenced in large part by specific sequences and structural motifs in the mRNA. The presence of upstream translation initiation codons and associated open reading frames (uORFs) in the 5'-UTR, plus specific sequences and secondary structures in the two untranslated regions, have been shown to repress translation of the protein coding region.

The inhibitory effect of uORFs can be explained by the inability of ribosomes to reinitiate translation at the downstream initiator AUG after completion of translation of the uORF (1, 13). The suppressive effect of uORFs can be overcome by several mechanisms, including the removal of upstream AUG codons by [alternative splicing](#) or through the use of alternative **promoters** to initiate [transcription](#). An uORF can be bypassed either through a “leaky scanning” mechanism (13), or through the use of an [internal ribosome entry site](#) to initiate translation (14). In addition, under certain conditions, ribosomes may acquire competence to bind to the ternary complex to reinitiate translation after having completed translation of the uORF (9).

During cap-dependent translation initiation, the formation of the 48 S preinitiation complex near the 5'-cap structure and the subsequent ribosome scanning process can be inhibited by the presence of mRNA secondary structures in the 5'-untranslated regions. The effect of an mRNA 5'-UTR secondary structure on the rate of translational initiation is determined by both the stability and the position of the secondary structure (15, 16). Secondary structures located close to the 5' end of the mRNA may prevent assembly of the preinitiation complex; when distant from the cap, a stem-loop structure with moderate stability may be easily unwound by the scanning 40 S ribosomal subunit and its associated initiation factors. For example, a hairpin structure with a stability of  $-30\text{kcal/mol}$

located 12 nucleotides downstream of the 5' cap completely inhibits the translation of preproinsulin mRNA, but has no effect on translation when it is located 52 nucleotides downstream of the cap structure (16).

### 3. RNA-Binding Proteins

[RNA-binding proteins](#) can be recruited to specific sites or to secondary structures on the mRNA to repress translation. The activity and/or expression of these RNA-binding proteins may determine the translational efficiency of specific mRNAs. The [iron-response element](#) (IRE) located at the cap-proximal regions of [ferritin](#) mRNA contains a 35-nucleotide sequence that folds into a specific stem-loop structure. The IRE is recognized by a repressor protein, IRE-BP. Under iron-starvation conditions, the IRE-BP binds to the IRE with high affinity, blocking recruitment of the 40 S ribosome subunit to ferritin mRNA and preventing its translation. As iron levels increase, IRE-BP dissociates from the IRE, and the translation repression of ferritin mRNA is relieved. The inhibitory effect of IRE-BP on ferritin mRNA translation is also position-dependent, requiring that the IRE be located within 40 nucleotides of the 5' cap (17). This suggests that formation of the cap-proximal IRE:IRE-BP complex may disrupt the interaction between the cap structure and the 40 S ribosome subunit.

In recent years, a growing number of *cis*-acting sequence elements have been identified in the 3' UTRs of mRNAs from various species. These 3'-UTR regulatory sequences have been shown to influence mRNA translation, stability, and localization (4, 18). In general, the 3' *cis*-acting regulatory elements involved in translational regulation interact with specific binding proteins to repress translation. One intriguing proposition is that RNA-binding proteins acting on the 3' end of mRNA disrupt the association between the 3' and 5' ends of mRNA that normally facilitate translation initiation.

During mouse spermatogenesis, protamine mRNAs are stored in the cytoplasm for 7 days prior to translation. A conserved sequence in the 3'-UTR of protamine mRNA is responsible for the appropriate temporal control of translation. The inhibitory function of this conserved sequence element is dependent on binding to an 18-kDa protein; translation is facilitated on dissociation of the factor (19). Erythroid 15-lipoxygenase (LOX) is a key enzyme during erythroid cell differentiation and is required for the breakdown of internal membranes during the late stages of erythrocyte maturation. LOX mRNA is synthesized in bone marrow erythroid precursor cells, but is not translated until the final steps of reticulocyte maturation. The 3'-UTR of LOX mRNA contains a 10-fold repeat of a CU-rich 19 nucleotide sequence that has been shown to mediate translational repression. Two proteins, hnRNP E1 and hnRNP K, bind to the CU-rich motif and mediate repression by disrupting the association of ribosomes with LOX mRNA (20).

Y-box proteins are [transcription factors](#) that bind specifically to Y-box sequences on DNA. In addition, these proteins are components of RNPs that bind non-specifically to stored mRNA molecules in germline cells. Binding of Y-box proteins to mRNA sequences, both *in vitro* and *in vivo*, has been shown to repress translation (21, 22).

### 4. Complementary Antisense RNAs

Small, stable RNA species have been shown to repress translation by pairing with specific sequences on target mRNA molecules. One of the best characterized examples occurs during early development in *Caenorhabditis elegans*, where the small, noncoding *lin-4* RNA binds to a complementary and repeated sequence element present in the 3'-UTR of *lin-14* mRNA. *lin-14* gene mutations, in which the complementary sequences in the 3'-UTR are partially or completely deleted, result in inappropriately high levels of LIN-14 protein expression at late stages of development. It is believed that *lin-4* regulates *lin-14* translation via an **antisense** RNA–RNA interaction (23). Alternatively, a repressor protein may bind to the bulged cytosine residue in the RNA/RNA hybrid structure and

block translation (24).

## Bibliography

1. W. C. Merrick (1992) *Microbiol. Rev.* **56**, 291–315.
2. K. Meerovitch, N. Sonenberg, and J. Pelletier (1990) in *Translation in Eukaryotes*, H. Trachsel, ed., CRC Press, Berne, Switzerland, pp. 273–292.
3. R. Duncan and J. W. B. Hershey (1984) *J. Biol. Chem.* **259**, 11882–11889.
4. D. Curtis, R. Lehmann, and P. D. Zamore (1995) *Cell* **81**, 171–178.
5. C. E. Samuel (1993) *J. Biol. Chem.* **268**, 7603–7606.
6. J.-J. Chen, M. S. Throop, L. Gehrke, I. Kuo, K. Pal, M. Brodesky, and I. M. London (1991) *Proc. Natl. Acad. Sci. USA* **88**, 7729–7733.
7. S. Legon, R. J. Jackson, and T. Hunt (1973) *Nature* **241**, 150–152.
8. M. B. Mathews and J. Shenk (1991) *J. Virol.* **65**, 5657–5662.
9. A. P. Geballe and D. Morris (1994) *Trends Biochem. Sci.* **19**, 159–164.
10. A. Lazaris-Karatzas, K. S. Montine, and N. Sonenberg (1990) *Nature* **345**, 544–547.
11. A. Pause, G. J. Belsham, A. C. Gingras, O. Donze, T.-A. Lin, J. C. Lawrence, and N. Sonenberg (1994) *Nature* **371**, 762–767.
12. N. Sonenberg (1987) *Adv. Virus Res.* **33**, 175–204.
13. M. Kozak and A. J. Shatkin (1978) *Cell* **13**, 201–212.
14. R. J. Jackson and A. Kaminski (1995) *RNA* **1**, 985–100.
15. J. Pelletier and N. Sonenberg (1985) *Mol. Cell. Biol.* **5**, 3222–3230.
16. M. Kozak (1986) *Cell* **44**, 283–292.
17. M. W. Hentze and L. C. Kühn (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8175–8182.
18. R. J. Jackson (1993) *Cell* **74**, 9–14.
19. Y. K. Kwon and N. B. Hecht (1991) *Proc. Natl. Acad. Sci. USA* **88**, 3584–3588.
20. D. H. Ostareck, A. Ostareck-Lederer, M. Wilm, B. J. Thiele, M. Mann, and M. W. Hentze (1997) *Cell* **89**, 597–606.
21. J. D. Richter and L. D. Smith (1984) *Nature* **309**, 378–380.
22. A. P. Wolffe (1994) *Bioassays* **16**, 245–251.
23. R. C. Lee, R. L. Feinbaum, and V. Ambros (1993) *Cell* **75**, 843–854.
24. I. Ha, B. Wightman, and G. Ruvkun (1996) *Genes Dev.* **10**, 3041–3050.

## Suggestions for Further Reading

25. J. W. B. Hershey, M. B. Mathews, and N. Sonenberg (1996) *Translational Control*, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
26. N. Standart and R. J. Jackson (1994) Regulation of translation by specific protein/mRNA interactions, *Biochimie* **76**, 867–879.
27. A. B. Sachs, P. Sarnow, and M. W. Hentze (1997) Starting at the beginning, middle and end: translation initiation in eukaryotes, *Cell* **89**, 831–838.
28. M. Wickens, P. Anderson, and R. J. Jackson (1997) Life and death in the cytoplasm: messages from the 3' end, *Curr. Opin. Genet. Dev.* **7**, 220–232.
29. M. Muckenthal and M. W. Hentze (1997) Mechanisms for post-transcriptional regulation by iron-responsive elements and iron regulatory proteins, *Prog. Mol. Subcell. Biol.* **18**, 93–115.

## Translational Editing

The [genetic code](#) is an algorithm that relates 20 specific [amino acids](#) used in protein synthesis to individual trinucleotide sequences, the **codons**. Because of degeneracy in the code, there are 61 triplet codons for the 20 amino acids. The relationship between these 20 amino acids and the 61 triplets is established in the aminoacylation reactions catalyzed by [aminoacyl tRNA synthetases](#). In these reactions, an amino acid is covalently attached to the 3' end of the [transfer RNA](#) (tRNA) that contains the [anticodon](#) triplet that is **cognate** to that amino acid.

Once an amino acid is attached to a tRNA, the position of the amino acid in a growing [polypeptide chain](#) is determined by the interaction between the anticodon of the tRNA and the matching, complementary codon of mRNA. This conclusion followed from early experiments that showed that, when alanine was attached to a cysteinyl-tRNA (tRNA<sup>Cys</sup>), alanine was incorporated into positions in a growing polypeptide that normally were occupied by cysteine (1) (this experiment, by Lipmann, Benzer, and co-workers, was done by reaction of Raney nickel with Cys-tRNA<sup>Cys</sup>, which converted Cys to Ala *in situ*). The conclusion from these experiments has since been confirmed by other means (2).

Errors in protein synthesis are relatively rare, having a frequency of perhaps  $10^{-4}$  *in vivo* (3-5). They would be much higher, except for translational editing mechanisms that specifically correct particular errors that occur at a much higher frequency. These errors occur because of limitations to molecular recognition at two critical points in translation. First, certain amino acids are difficult for synthetases to discriminate from each other and, for that reason, can be misactivated and potentially attached to the tRNA cognate to the synthetase but noncognate to the amino acid. Second, the anticodon-codon interaction on the [ribosome](#) is determined by base pairing between the complementary triplets, and this interaction is prone to mistakes because, for example, pairing of two rather than three bases, although less probable, is still possible.

In the case of misactivation of amino acids by synthetases, the energy cost for proofreading is in the form of extra hydrolysis of ATP. For the ribosome-dependent system, the cost of proofreading is in an additional hydrolysis of GTP. Errors also occur in [DNA replication](#), and here, as in translation, nucleoside triphosphate hydrolysis is required to correct errors of replication. The general requirement for triphosphate hydrolysis in error-correcting mechanisms was recognized early by Hopfield (6) and Nino (7).

Translational editing by tRNA synthetases requires the presence of the cognate tRNA. This system of proofreading is an example of RNA-dependent amino acid recognition (8). The tRNA synthetases were perhaps among the first proteins to emerge from an [RNA world](#) in connection with the development of the genetic code. Possibly, the present-day RNA-dependent recognition of amino acids in translational editing developed out of an early system in which RNA played a more prominent role in aminoacylation activity in general.

### 1. RNA-Dependent Amino Acid Discrimination in Translational Editing

#### 1.1. Misactivation of Amino Acids

Most tRNA synthetases can activate amino acids in the absence or presence of tRNA. The activated amino acid then reacts with the 3' end of the tRNA cognate to the amino acid:



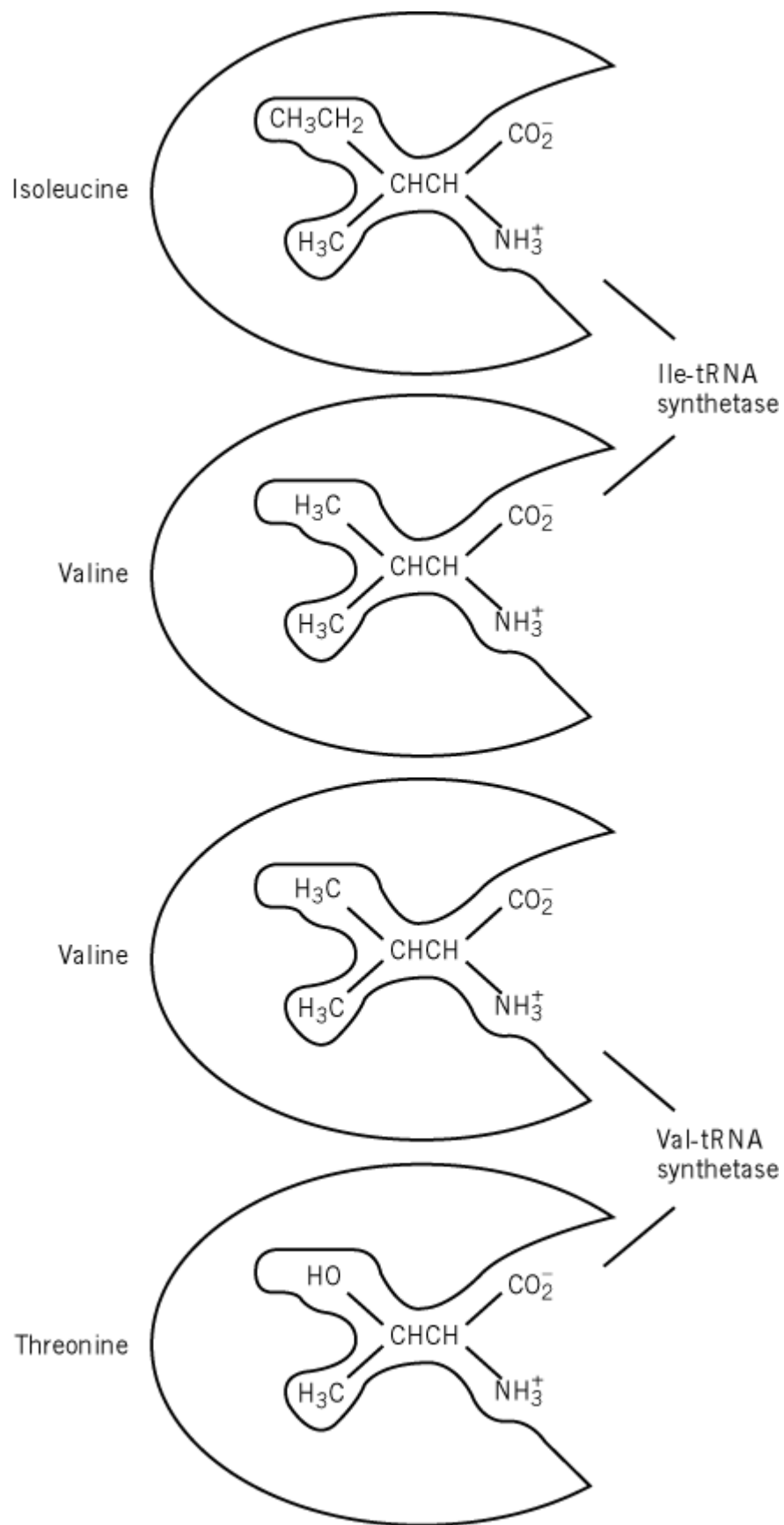




In reaction (1), E, AA, AA-AMP, and  $\text{PP}_i$  denote enzyme, amino acid, aminoacyladenylate, and inorganic pyrophosphate, respectively. In reaction (2), AA-tRNA represents aminoacyl-tRNA. The first step is referred to as amino acid activation, while the second step is a transfer reaction where the aminoacyl moiety is transferred from the adenylate to the 3' end of the tRNA (see [Aminoacyl tRNA Synthetases](#)). In general, the transfer reaction is the rate-determining step in aminoacylation.

Potential errors of aminoacylation occur because of inherent difficulties in molecular recognition associated with step 1 (5, 9, 10). In particular, Pauling pointed out that closely similar amino acids such as [valine](#) and [isoleucine](#), which are both branched at the  $\beta$ -carbon and differ by one methylene group, are difficult for isoleucyl-tRNA synthetase to discriminate on the basis of binding interactions alone (11) (Fig. 1). For example, a pocket that can accommodate the isobutyl side chain of isoleucine should also be able to accommodate the isopropyl group of valine. The loss of [van der Waals interactions](#) associated with the missing methylene group of valine should be calculable from measurements of the transfer of a methylene group from a hydrophobic solvent to water (see [Hydrophobicity](#)). Similarly, valyl-tRNA synthetase is confronted with distinguishing between valine and the nearly isosteric threonine. These expectations are fulfilled, because valine is misactivated by isoleucyl-tRNA synthetase and threonine is misactivated by valyl-tRNA synthetase.

**Figure 1.** Schematic illustration of the binding of isoleucine and valine to IleRS and of valine and threonine to ValRS. The isopropyl side chain of valine can fit easily into the pocket in IleRS that binds the isobutyl side chain of isoleucine. In the case of threonine and valine, the side chains are also isosteric, thus making it possible for threonine to fit into the pocket in ValRS for valine. (Adapted from Ref. 9.)



### 1.2. Translational Editing as an ATPase Activity

The overall translational editing activity was discovered with the isoleucine system (12). Valine was

observed to be misactivated by isoleucyl-tRNA synthetase (IleRS). When challenged with tRNA<sup>Ile</sup>, the misactivated adenylate was hydrolyzed to AMP and valine in the overall reactions:



The sum of reactions (3) and (4) is simple ATP hydrolysis:



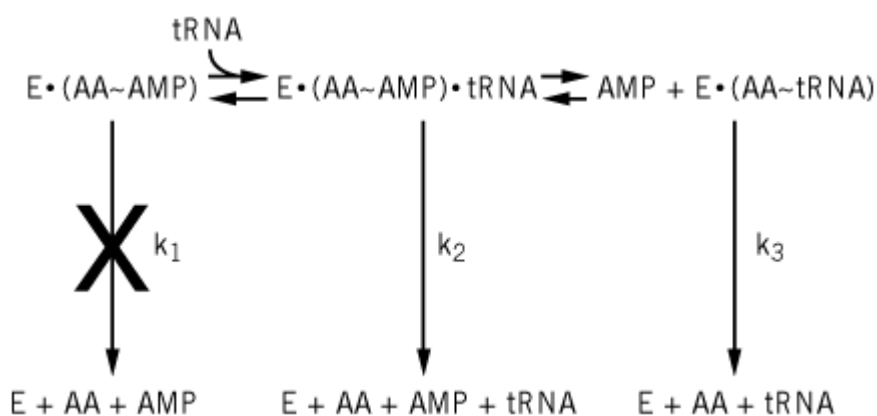
Thus, translational editing can be measured as a tRNA<sup>Ile</sup>-dependent [ATPase](#) activity of IleRS, where the enzyme discriminates between isoleucine and valine in a reaction that depends on tRNA<sup>Ile</sup>.

### 1.3. Two Steps to Translational Editing by tRNA Synthetases

When the ATPase activity associated with RNA-dependent amino acid discrimination was discovered, no details of the mechanism were understood. Now it is known that, when valine is misactivated by isoleucyl-tRNA synthetase, the misactivated amino acid can, in principle, be attached in step 2 to tRNA<sup>Ile</sup> to yield Val-tRNA<sup>Ile</sup>. The species Thr-tRNA<sup>Val</sup> can similarly be generated when threonine is misactivated by valyl-tRNA synthetase. At this point, three potential reactions can occur (Fig. 2). First, the misactivated amino acid can be spontaneously hydrolyzed while bound to the enzyme ( $k_1$  step). This reaction is inconsequential because, in the absence of added tRNA, even the misactivated amino acid is stably bound to the enzyme. Second, the misactivated species can be “induced” to hydrolyze in the presence of the tRNA ( $k_2$  step), or it can be transferred onto the tRNA to give a mischarged species. This mischarged amino acid is subsequently hydrolyzed from the tRNA by the action of an esterase activity of the synthetase that is directed toward the mischarged tRNA ( $k_3$  step) ([13-15](#)).

**Figure 2.** Schematic of editing pathways. In this illustration, a misactivated aminoacyl adenylate (AA-AMP) is stably bound to a tRNA synthetase enzyme, E. No hydrolysis to AA + AMP occurs ( $k_1$  step) in the absence of tRNA.

Addition of tRNA causes breakdown of the adenylate by one of two steps. In the  $k_2$  step, the adenylate is hydrolyzed before the amino acid is transferred to the tRNA (pretransfer editing). In the  $k_3$  step, the mischarged aminoacyl tRNA species is formed and an esterase activity of the enzyme catalyzes removal of the misactivated amino acid from the tRNA.



The  $k_2$  and  $k_3$  steps represent tRNA-dependent amino acid discrimination in translational editing. They occur only when the mischarged amino acid is first bound and then activated to an adenylate that is located at the active site of the enzyme (9, 16). In contrast, activation of the correct amino acid leads to stable attachment of the amino acid to its cognate tRNA. Thus, in the context of its tRNA, a synthetase can distinguish two closely similar amino acids through these hydrolytic editing reactions.

#### 1.4. Elucidation of the Two Steps in RNA-Dependent Amino Acid Discrimination

The  $k_3$  step of translational editing was demonstrated directly. The mischarged Val-tRNA<sup>Ile</sup> was generated by a special procedure. When IleRS was added, an esterase activity of the enzyme removed valine from tRNA<sup>Ile</sup>:



This esterase activity was first seen as a weak activity that synthetases had toward their cognate, properly charged, aminoacyl tRNA, such as Ile-tRNA<sup>Ile</sup> (13). This activity was then observed to be greatly enhanced when the mischarged species was presented to the enzyme (14).

The  $k_2$  step indicates that the misactivated adenylate is hydrolytically removed before transfer to the tRNA. This step was first inferred from kinetic studies that were able to distinguish two steps to the overall editing reaction, one of which occurred before transfer of the aminoacyl group to tRNA (9). Thus, the  $k_2$  and  $k_3$  steps can be referred to as pretransfer and posttransfer editing reactions, respectively. In the pretransfer reaction, the tRNA is imagined to act as an effector and not as an amino acid-accepting substrate.

#### 1.5. DNA-Aptamer-Induced Translational Editing

If the amino acid acceptor function of tRNA<sup>Ile</sup> were still necessary for the  $k_2$  step, then that step might be thought of as closely related or similar to the  $k_3$  step, where a mischarged species is deacylated. For example, the acceptor hydroxyl groups of the tRNA itself might participate in the hydrolytic reaction. Independent demonstration of the pretransfer step was accomplished by searching for a surrogate effector that was incapable of accepting an amino acid. A DNA **aptamer** was selected from a pool of more than  $10^{14}$  molecules. DNA lacks a 2'-OH, which is the initial site of amino acid attachment by the isoleucine enzyme (17). The selected aptamer replaces tRNA<sup>Ile</sup> in reaction (4), inducing the hydrolysis of Val-AMP bound to IleRS. In contrast, the aptamer had no effect on the stability of Ile-AMP bound to the isoleucine enzyme. As expected, the aptamer was not an acceptor for isoleucine. These data show that a constellation of nucleotides per se can act as a surrogate effector for RNA-induced translational editing, without requirement for amino acid acceptance activity. They provide strong support for the existence of the  $k_2$  step of translational editing.

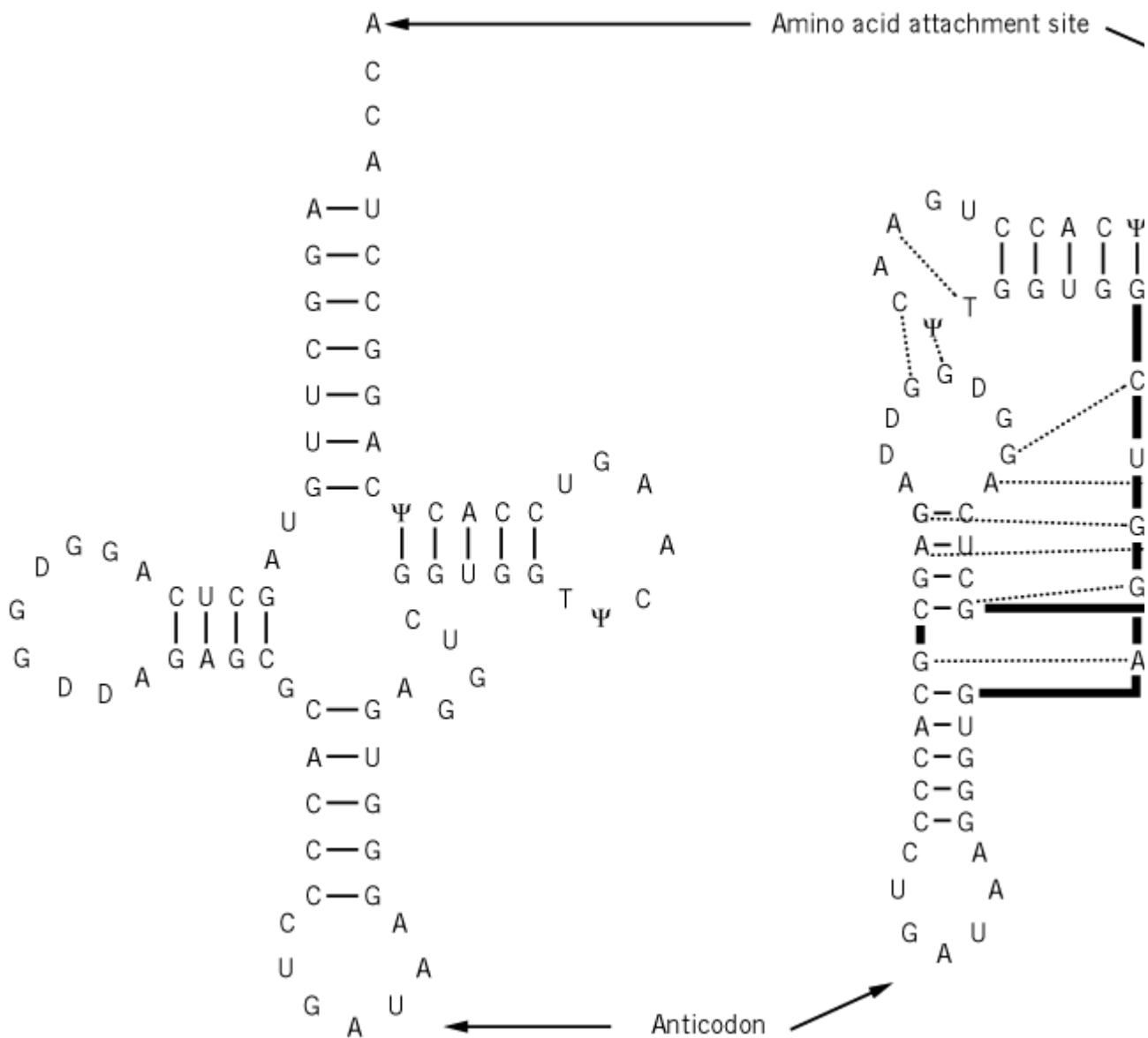
#### 1.6. tRNA Nucleotides for Translational Editing Are Distinct from Those for Aminoacylation

The idea that a constellation of nucleotides per se can elicit the ATPase activity (Eqs. 3-5) motivated experiments to identify the location of those nucleotides in a tRNA. Actually, tRNA<sup>Val</sup> does not substitute for tRNA<sup>Ile</sup> as the effector in equation (4); that is, addition of tRNA<sup>Val</sup> to IleRS(Val-AMP) results in no hydrolytic breakdown of Val-AMP. Thus, some of the nucleotides that distinguish tRNA<sup>Ile</sup> from tRNA<sup>Val</sup> must account for the difference in translational editing activity. Because nucleotide differences also must explain why one is charged with isoleucine and the other is not, the possibility that nucleotides needed for charging are also needed for editing has to be considered.

Nucleotides needed for aminoacylation of tRNA are typically located within the [acceptor stems](#) and, in many instances, within the anticodon trinucleotides as well (see [Aminoacyl tRNA Synthetases](#)). Transplantation of the GAU anticodon of tRNA<sup>Ile</sup> into tRNA<sup>Val</sup> converts the latter to an isoleucine acceptor that is designated as tRNA<sup>Val/GAU</sup>. Isoleucine acceptance is further enhanced if the acceptor stem of tRNA<sup>Ile</sup> is introduced into tRNA<sup>Val/GAU</sup>. Yet neither of these “chimerized” valine tRNAs is active in stimulating the ATPase activity that is characteristic of the editing reactions ([eqs. 3-5](#)) ([18](#)). These results showed that nucleotides sufficient for charging a tRNA do not in themselves confer editing.

Instead, nucleotides at the corner of the L-shaped tRNA molecule trigger the editing response ([Fig. 3](#)). The editing function of tRNA<sup>Ile</sup> or of tRNA<sup>Val/GAU</sup> can be turned off or on by simple manipulation of just three nucleotides. These nucleotides are part of a highly differentiated region of the molecule where two loops of the tRNA cloverleaf are brought together. When the three-dimensional structure of tRNA was first determined, this highly differentiated corner was suggested to be a potential protein discrimination site ([19](#)). Possibly it is a discrimination site used for triggering the editing response. Additionally or alternatively, this part of the molecule may be involved in transducing a signal between the two domains of the tRNA structure. This signal would be part of a system to ensure that the trinucleotide anticodon and the attached amino acid were matched according to the rules of the genetic code.

**Figure 3.** Nucleotides in a tRNA that trigger the editing response. The L shape of the tRNA structure is depicted. Nucleic tertiary interactions are joined by dotted lines. The nucleotides that are needed for the editing response are shown in bold



### 1.7. Editing Site on a Synthetase Is Distinct from That for Aminoacylation

The editing activity is directed toward an adenylate and an aminoacyl-tRNA, species that are also for cognate amino acid. Early work suggested, however, that the editing activity occurred at a site that was distinct from the active site (13). To address the question of whether the two sites can be separated, mutational analysis and [directed mutagenesis](#) of IleRS found several mutations that severely reduced the catalytic activity by more than two orders of magnitude. One mutation was obtained that eliminated valine and isoleucine in the initial amino acid activation step (20) (Eq. 1). In contrast, when challenged with valine and isoleucine, the enzyme efficiently discriminated against misactivated Val-AMP in the hydrolytic editing reaction (E was mutationally isolated from the active site).

In another set of experiments, mutations were made in ValRS that severely affected amino acid activation for valine. In an analysis of aminoacylation error correction, the deacylation activity of two of the mutant Thr-tRNA<sup>Val</sup> was studied. In each instance, this activity was unaffected by the mutations (21). This suggests a design for at least some tRNA synthetases that separates the catalytic site for aminoacylation from the site for amino acid discrimination.

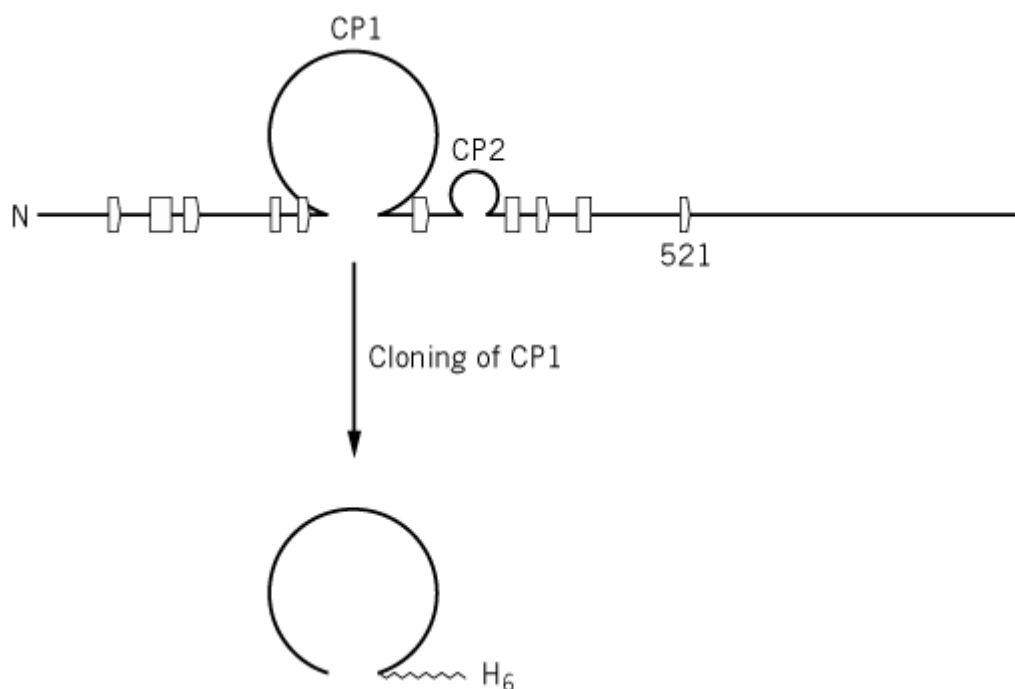
### 1.8. Cloned Insertions from Two tRNA Synthetases Destroy Errors of Aminoacylation

Isoleucyl- and valyl-tRNA synthetases are examples of class I enzymes. These enzymes have a core

constructed from alternating b-strands and a-helices in an alternating b–a–b pattern (see [Aminoacyl tRNA synthetase](#) divided into two  $b_3a_2$  halves that are separated by an insertion into the active site known as connective loop 1 (CP1). An experiment designed to determine whether the aminoacyl moiety of Ile-tRNA<sup>Ile</sup> and Val-tRNA<sup>Ile</sup> both require a special modification of the aminoacyl moiety was created. This modification enabled covalent labeling of each of these two aminoacyl tRNAs (22). While each labeled the same site in the nucleotide-binding pocket, the modified Val-tRNA<sup>Ile</sup> species also labeled a site within CP1. This result raised the possibility that CP1 is involved in the deacylation of Val-tRNA<sup>Ile</sup>.

This possibility was tested by [cloning](#) the CP1 insertion from IleRS and from ValRS (Fig. 4) (23). The cloned CP1 proteins can be used directly in an assay for deacylation of mischarged tRNA. The CP1 insertion from ValRS (CP1<sup>Ile</sup>) was shown to have an esterase activity that catalyzes hydrolysis of Val-tRNA<sup>Ile</sup>, but not of Ile-tRNA<sup>Ile</sup>. The CP1 insertion from IleRS (CP1<sup>Val</sup>) catalyzes hydrolysis of Thr-tRNA<sup>Val</sup> but does not destroy Val-tRNA<sup>Val</sup>. Thus, the CP1 insertions are responsible for destroying errors of aminoacylation (the  $k_3$  step of Fig. 2).

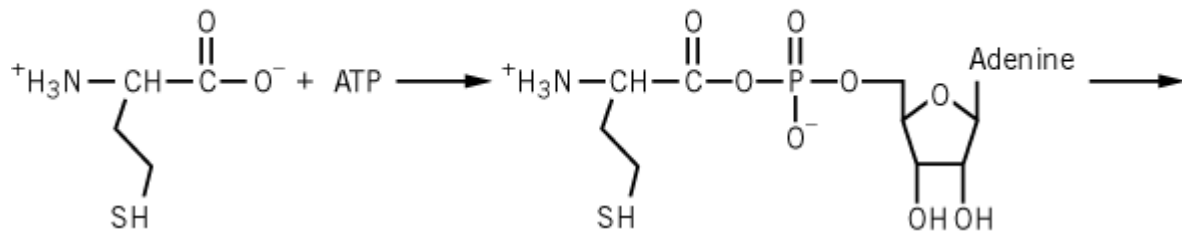
**Figure 4.** Cloning of the CP1 insertion. A schematic diagram of the *Bacillus stearothermophilus* ValRS is shown. The protein is composed of alternating b-strands (arrows) and a-helices (rectangles). The fold is split by the CP1 insertions, which can be cloned as a separate domain.



## 2. RNA-Independent Amino Acid Recognition in Translational Editing

Methionyl-tRNA synthetase is a special example that uses a tRNA-independent pathway to reject the *homocysteine* (5). This amino acid is generated by transfer of the active methyl group from *S*-adenosylmethionine to give *S*-adenosylhomocysteine, followed by hydrolysis of *S*-adenosylhomocysteine to adenosylhomocysteine and pyruvate. Homocysteine is also generated by the conversion of cystathionine to homocysteine and pyruvate. Homocysteine is then activated by MetRS (Eq. 1) to form homocysteinyladenylate. In the absence of tRNA<sup>Met</sup>, the misactivated adenylylated homocysteine thiolactone, which is released from the enzyme (Fig. 5). By conversion to the thiolactone, homocysteine is prevented from being incorporated into proteins (homocysteine is also activated by IleRS and, similarly, converted to methionine by MetRS). The physiological consequences of the production of homocysteine thiolactone is a question of great interest.

**Figure 5.** Formation of homocysteine thiolactone. Homocysteiny adenylate is formed by reaction of homocysteine with methionyl-tRNA synthetase. The activated homocysteiny adenylate is subsequently converted to the thiolactone and re consequence, homocysteine is not attached to methionyl tRNA and is not incorporated into proteins.



### 3. Translational Editing on the Ribosome

Aminoacylated tRNA binds to the A site on the [ribosome](#); specificity is determined by the codon–an [Translation](#)). The specificity of the codon-anticodon interaction is enhanced by an editing mechanism: GTP to GDP and  $P_i$ . The source of the GTP in this hydrolytic reaction is that which is bound to [elong](#); member of the G-factor superfamily, EF-Tu has two forms—one that is “active” (GTP-bound) and an inactive form (GDP-bound). The active form binds aminoacyl-tRNA while the inactive form does not.

After aminoacylation, the charged tRNA is bound to the GTP complex of EF-Tu and carried to the ribosome and bound to the A site. With a perfect codon–anticodon match, the GTP is hydrolyzed and the aminoacyl moiety is transferred to the A site, where the aminoacyl moiety can react with the peptide group of the peptidyl tRNA in the A site. In the case of an imperfect anticodon match, some data showed that two molecules of GTP are hydrolyzed per aminoacyl tRNA interaction. If the interaction is imperfect (eg, two base pairs rather than three are complementary), then the imperfectly bound aminoacyl tRNA is rejected from the A site in a reaction that also results in hydrolysis of the GTP that is bound to EF-Tu. Thus, in order for a mismatched aminoacyl tRNA to bind again at the A site, it must rebind to a new molecule of EF-Tu-GTP and repeat the cycle. In consequence, an abortive cycle is established in which GTP hydrolysis accompanies the successive binding and rejection of mismatched aminoacyl-tRNA.

### 4. Conclusions

The translational editing reactions described above share in common the use of nucleoside triphosphates to prevent errors in protein synthesis. In both the RNA-dependent amino acid discrimination by tRNA and the editing on the ribosome, RNA plays a critical role as an effector in the editing mechanism. The effector function is not clear, but the tRNA synthetase system shows clearly that specific nucleotides in the tRNA actually affect the editing response. Whether specific nucleotides (other than the anticodon triplet) in a particular part of the tRNA structure affect the editing response is also a question that requires more investigation. These RNA-dependent events may have had even more roles in protein synthesis than it does in contemporary life forms, when RNA may have had even more roles in protein synthesis than it does in contemporary life forms.

### Bibliography

1. F. Chapeville, F. Lipmann, G. von Ehrenstein, B. Weisblum, W. J. Ray, Jr., and S. Benzer (1962) *Proc. Natl. Acad. Sci. USA* **48**, 1086–1092.
2. J. Normanly and J. Abelson (1989) *Annu. Rev. Biochem.* **58**, 1029–1049.



3. R. B. Loftfield and D. Vanderjagt (1972) *Biochem. J.* **128**, 1353–1356.
4. R. B. Loftfield (1963) *Biochem. J.* **89**, 82–92.
5. H. Jakubowski and E. Goldman (1992) *Microbiol. Rev.* **56**, 412–429.
6. J. J. Hopfield (1974) *Proc. Natl. Acad. Sci. USA* **71**, 4135–4139.
7. J. Nino (1975) *Biochimie* **57**, 587–595.
8. P. Schimmel and E. Schmidt (1995) *Trends Biochem. Sci.* **20**, 1–2.
9. A. Fersht (1985) *Enzyme Structure and Mechanism*, W. H. Freeman and Company, New York, p. 102.
10. F. Cramer, U. Englisch, W. Freist, and H. Sternbach (1991) *Biochimie*. **73**, 1027–1035.
11. L. Pauling (1957) in *Festschrift für Prof. Dr. Arthur Stoll*, Birkhauser-Verlag, Basel, pp. 597–600.
12. A. N. Baldwin and P. Berg (1966) *J. Biol. Chem.* **241**, 839–845.
13. A. A. Schreier and P. R. Schimmel (1972) *Biochemistry* **11**, 1582–1589.
14. E. W. Eldred and P. R. Schimmel (1972) *J. Biol. Chem.* **247**, 2961–2964.
15. M. Yarus (1972) *Proc. Natl. Acad. Sci. USA* **69**, 1915–1919.
16. W. Freist (1989) *Biochemistry* **28**, 6787–6795.
17. S. Hale and P. Schimmel (1996) *Proc. Natl. Acad. Sci. USA* **93**, 2755–2758.
18. S. P. Hale, D. S. Auld, E. Schmidt, and P. Schimmel (1997) *Science* **275**, 1250–1252.
19. J. E. Ladner, A. Jack, J. D. Robertus, R. S. Brown, D. Rhodes, B. F. C. Clark, and A. Klug (1977) *Proc. Natl. Acad. Sci. USA* **74**, 4414–4418.
20. E. Schmidt and P. Schimmel (1994) *Science* **264**, 265–267.
21. L. Lin and P. Schimmel (1996) *Biochemistry* **35**, 5596–5601.
22. E. Schmidt and P. Schimmel (1995) *Biochemistry* **34**, 11204–11210.
23. L. Lin, S. P. Hale, and P. Schimmel (1996) *Nature* **384**, 33–34.
24. A. Weijland and A. Parmeggiani (1993) *Science* **259**, 1311–1314.

## Translocation

[Translation](#) of a [messenger RNA](#) during protein biosynthesis requires a coupled movement of mRNA and transfer RNAs throughout the elongation stage. Each new amino acid is recruited to the [ribosome](#) as an aminoacyl-tRNA:EF-Tu:GTP ternary complex. Elongation factor EF-G catalyzes a precise movement of the tRNA–mRNA complex within the ribosome after **peptidyl transfer** to empty the A site for the next aminoacyl-tRNA (see [Elongation Factors \(EFs\)](#)). During each translocation step, the elbow of tRNA moves on the order of 50 Å. Although the elongation cycle under physiological conditions requires the two G proteins EF-Tu and EF-G, it has been shown that polypeptide synthesis can be carried out by the ribosome itself, in the absence of factors or GTP under certain *in vitro* conditions [such as poly(U)-dependent polyphenylalanine synthesis]. Thus, the ability to move mRNA and tRNA is an inherent property of the ribosome; the factors serve to increase the speed and accuracy of elongation in a GTP-dependent manner.

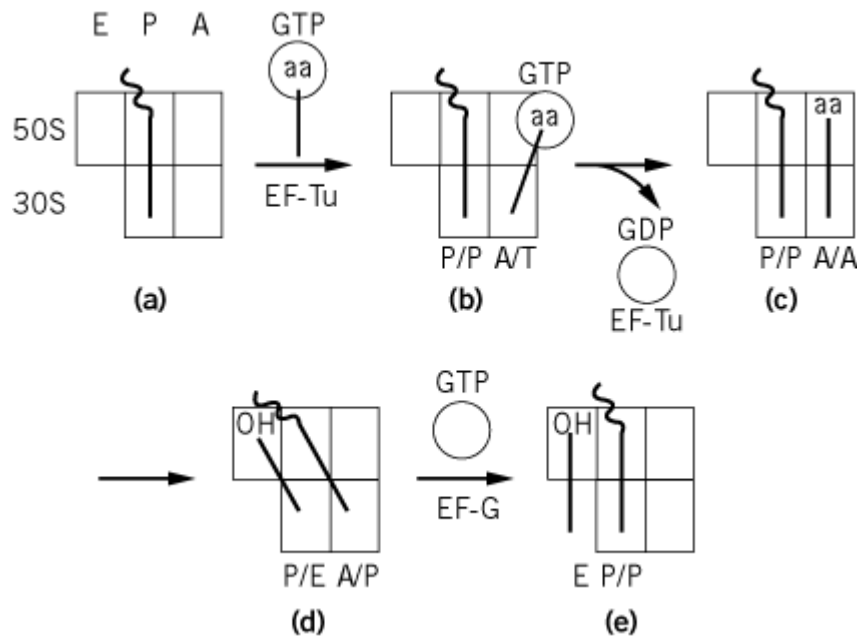
Eukaryotic and prokaryotic ribosomes are composed of one large subunit and one small subunit that fit together to form a complete ribosome with mass of several millions daltons (see [Ribosomes](#)). The small subunit matches the tRNA to the **codons** of the mRNA, while the large subunit catalyzes formation of the [peptide bonds](#) that link the amino acids together into a polypeptide chain. The

ribosome contains more than 50 different proteins and three rRNA molecules, provides specific binding sites for mRNA, tRNA, and the various initiation, elongation, and termination factors, and catalyzes peptide bond formation. The ribosome also stabilizes codon–[anticodon](#) interactions and preserves the [reading frame](#).

How does the ribosome choreograph the movements required for translation? There are two major tRNA-binding sites on the ribosome: one incoming site, called the A site for aminoacyl–tRNA, and the P site for peptidyl–tRNA, where new peptide bonds are formed. A third site, the E site, is transiently occupied by the deacylated end of the tRNA that has completed amino acid transfer and is about to be ejected. According to the three-site version of the classical model, beginning with an initiator or peptidyl–tRNA in the P site, a new aminoacyl tRNA, with an anticodon complementary to the mRNA codon positioned at the A site, is introduced as an aminoacyl–tRNA:EF-Tu:GTP ternary complex. Upon correct binding of the aminoacyl–tRNA, GTP is hydrolyzed, and the GDP form of EF-Tu is released from the ribosome, along with inorganic phosphate. The anticodon ends of both tRNAs interact with the small 30 S subunit, and their acceptor ends interact with the large 50 S subunit. The ribosome is now ready for formation of the next peptide bond. This step is catalyzed by the peptidyl transferase activity of the large ribosomal subunit, and the growing peptide chain carried by the tRNA at the P site links to the amino acid carried by the tRNA at the A site. Movement of the newly created peptidyl–tRNA from the A to the P site is accomplished by EF-G in a GTP-dependent reaction. At the same time, the deacylated tRNA moves to the E site. The E site is most likely located on the 50 S subunit; this would mean that the elongation cycle is really a two-and-a-half-site, rather than a three-site, model. Deacylated tRNA, bound weakly to the E site, dissociates from the ribosome to complete the cycle of elongation.

Chemical [footprinting](#) studies in certain intermediate states of elongation showed that the two ends of a tRNA could be in different states on the two ribosomal subunits; for example, a tRNA could simultaneously occupy the 30 S A site and the 50 S P site. Interpretation of these experiments resulted in the hybrid states model for the elongation cycle ([1](#), [2](#)). According to this model, translocation of tRNAs on the ribosome occurs in two separable steps (Fig. [1](#)). In the first step, the movement of the acceptor ends of tRNAs relative to the 50 S ribosomal subunit occurs spontaneously following peptide bond formation, independent of EF-G and GTP (Fig. [1](#)). The peptidyl–tRNA is in the A/P state, in which its anticodon is in the 30 S A site, while its acceptor end is in the 50 S P site. The deacylated tRNA is in the P/E site, in which its anticodon is in the 30 S P site and its acceptor end in the 50 S E site. The second step, which is EF-G and GTP-dependent, results in translocation of the anticodon arms of the two tRNAs, together with their associated mRNA, relative to the 30 S subunit. The peptidyl–tRNA moves into the P/P state and deacylated tRNA into the E site (Fig. [1](#)). Independent movement of tRNAs relative to the two ribosomal subunits suggests that tRNA translocation could involve relative movement of the subunits, providing an explanation for the universal two-subunit structure of ribosomes.

**Figure 1.** Hybrid states model for the translational elongation cycle ([2](#)). The tRNA binding sites on the 50 S and 30 S subunits are represented schematically by the upper and lower rectangles, respectively. The 50 S subunit is subdivided into A, P, and E sites; the 30 S subunit is subdivided into A and P sites. The tRNAs are represented by vertical bars, and the nascent polypeptide chain is represented by a wavy line; circles represent EF-Tu and EF-G, respectively. Binding states, indicated at the bottom of each panel, indicate the state of each tRNA relative to each ribosomal subunit; A/P indicates interaction of the anticodon end to a tRNA with the 30 S A site and its acceptor end to the 50 S P site.



The discovery of factor-independent translocation established that it is fundamentally a ribosomal mechanism. Equally significant is the finding that translocation can take place in the absence of mRNA. Under certain *in vitro* conditions, lysyl-tRNA can be polymerized by the ribosome into polylysine. Thus, the translocation machine acts directly on tRNA and the movement of mRNA is passive, driven by its association with tRNA.

Recent studies have provided evidence for further insight into the hybrid states model. First, initial binding of the EF-Tu ternary complex is rapid and reversible and does not involve codon recognition (3). This initial complex can form even if the A site is blocked with another aminoacyl-tRNA. A series of conformational changes follow the initial binding to lead to codon recognition and GTP hydrolysis. Upon dissociation of EF-Tu:GDP, the aminoacyl-tRNA enters the 50 S site, leading to peptide bond formation. Second, until recently, binding of EF-G:GTP to the ribosome was thought to induce translocation, followed by GTP hydrolysis and release of EF-G:GDP from the ribosome. However, recent pre-steady-state kinetic experiments show that GTP hydrolysis occurs before translocation and accelerates translocation more than 50-fold relative to that observed with nonhydrolyzable GTP analogues. It reevaluates the role of GTP hydrolysis in translocation, suggesting that a conformational transition in EF-G itself is coupled to translocation (4).

The structures of EF-G and its GDP-complex that have been solved by [X-ray crystallography](#) resemble the aminoacyl-tRNA:EF-Tu:GTP ternary complex (see [Elongation Factors \(EFs\)](#)). The tRNA anticodon stem mimicry by domains IV of EF-G is essential for its translocase activity to fit the A site of the ribosome. Superficially, EF-Tu and EF-G appear to be involved in translational steps that seem mechanistically unrelated; EF-G catalyzes translocation, while EF-Tu introduces aminoacyl-tRNA to the ribosome. Their fundamental mechanism, however, may be to catalyze the two very similar sets of rotational steps, which may themselves be innate properties of ribosomal mechanics.

#### Bibliography

1. Moazared and H. F. Noller (1989) *Nature* **342**, 142–148.
2. K. S. Wilson and H. F. Noller (1998) *Cell* **92**, 131–139.
3. M. V. Rodnina, T. Pape, R. Fricke, and W. Wintermeyer (1995) *Biochem. Cell. Biol.* **73**, 1221–1227.

4. M. V. Rodnina, A. Savelsbergh, V. I. Katunin, and W. Wintermeyer (1997) *Nature* **385**, 37–41.

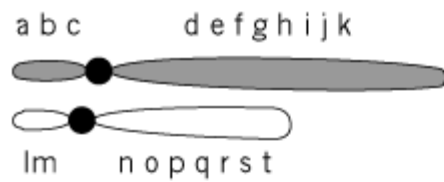
### **Suggestion for Further Reading**

5. H. Stark et al. (1997) *Cell* **88**, 19–28; (1997) *Nature* **389**, 403–406.

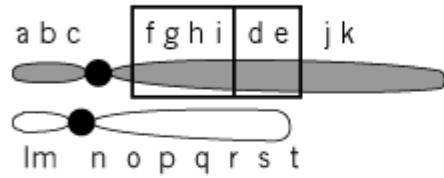
## **Translocation, Chromosomal**

Chromosomal translocations occur when segments of one [chromosome](#) are relocated to another. They result from chromosome breakage and fusion (Fig. 1). Three breaks in a single chromosome lead to an intrachromosomal translocation, in which the translocated segment might be inverted. Single breaks lead to other chromosomal changes (Fig. 1). [Reciprocal translocations](#) occur when two non[homologous chromosomes](#) exchange segments. Translocations of chromosomal segments to the end of the chromosome also occur. This type of terminal translocation might account for the existence of **telomeric** repeats within the chromosomal arms. Finally, Robertsonian events are an extreme form of translocation when two [acrocentric chromosomes](#) fuse to generate a **metacentric**.

**Figure 1.** Examples of chromosomal translocations. **(a)** Chromosomal translocation leading to the insertion (intercalation) of one chromosomal segment into a distinct position in the same chromosome. **(b)** Chromosomal translocation leading to the intercalation of one chromosomal segment into a non homologous chromosome. **(c)** Reciprocal chromosomal translocation between nonhomologous chromosomes. Two possible outcomes are shown: a balanced pair of chromosomes or a combination of a dicentric and acentric chromosome.

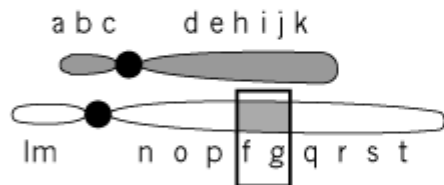


Interchromosomal transposition

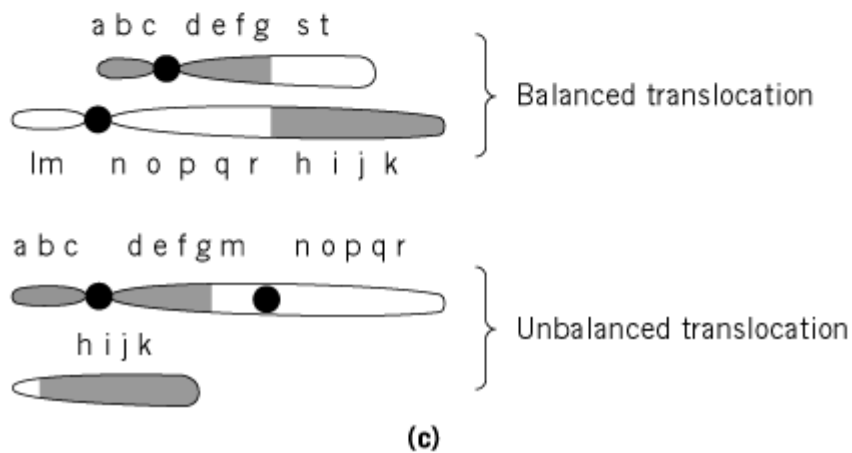


(a)

Interchromosomal transposition



(b)



(c)

### Suggestion for Further Reading

R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes: a synthesis*, Wiley-Liss, New York.

### Transmembrane $\alpha$ -Helix

In most [membrane proteins](#), the transmembrane regions are composed of  **$\alpha$ -helices**. Each helix

typically consists of 20 to 30 consecutive **hydrophobic** amino acid residues, although single [polar](#) or potentially charged residues are occasionally inserted into the membrane-spanning sequence. Helical **secondary structure** is a natural way of embedding a [polypeptide chain](#) into lipid. [Hydrogen bonding](#) between the peptide groups is satisfied by the network that forms along the helical axis between residues that are close in the sequence.

Structural analysis of the [photosynthetic reaction center](#) (1) and the mitochondrial cytochrome *c* oxidase (2), which contain a total of 39 transmembrane  $\alpha$ -helices, suggests the following general features of transmembrane helices:

1. The helices are about 40 Å long when they penetrate the bilayer as straight rods. In most cases, however, they are slightly tilted relative to the membrane normal.
2. The ends of the  $\alpha$ -helices have polar residues.
3. The hydrophobic core of a helix spans about 20 Å and is composed of hydrophobic residues, such as [phenylalanine](#), [valine](#), [leucine](#), [methionine](#), and [isoleucine](#).
4. Polar aromatic residues ([tyrosine](#), [tryptophan](#)) are found on both sides of the core. Close to the hydrophobic core, aromatic residues are exposed to the lipid headgroup region, whereas the buried tyrosine and tryptophan residues are slightly further from the core.
5. The latter region is also populated by [asparagine](#) and [glutamine](#) residues that are either exposed to the surface or buried in the protein structure.
6. The occurrence of [proline](#) residues peaks around 25 Å from the core, just outside the ends of the helices. Proline and [glycine](#) residues are buried in the structure when they occur in the hydrophobic core.

#### Bibliography

1. D. C. Rees, H. Komiya, T. O. Yeates, J. P. Allen, and G. Feher (1989) The bacterial photosynthetic reaction center as a model for membrane proteins. *Annu. Rev. Biochem.* **58**, 607–633.
2. E. Wallin, T. Tsukihara, S. Yoshikawa, G. Von Heijne, and A. Elofsson (1997) Architecture of helix bundle membrane proteins: an analysis of cytochrome *c* oxidase from bovine mitochondria. *Protein Sci.* **6**, 808–815.

## Transposable Elements

Transposable elements are discrete pieces of DNA that can move between nonhomologous positions in the [genome](#). Within a single cell, they may translocate between positions within a [chromosome](#), between chromosomes, or between chromosomes and extrachromosomal DNAs such as plasmids. Such elements have been found in virtually every organism examined and can take a variety of forms. In some cases, they encode only functions to promote their own movement—that is, a [recombinase](#) called a [transposase](#) (or [integrase](#), when the transposable element is a **virus** that integrates into the host genome) and special recombination sequences at the tips of the element that are the recombination substrates on which the transposase acts. In other cases, transposable elements additionally encode other information, such as [antibiotic-resistance](#) genes or viral genomes.

All transposable elements, whether they encode only transposase or also additional genes, have an “activity”: Their movement to new insertion sites is very often accompanied by changes in host gene expression. If the element inserts into a **gene**, that gene is disrupted and usually inactivated; such

inactivation could, for example, lead to cell death if the gene product is essential. Alternatively, the expression of a host gene adjacent to the point of element insertion may come under the influence of the transposon, for example by element-embedded [enhancers](#), thereby disrupting the usual host-specified expression pattern. Thus transposable elements can influence their host by altering genomic DNA upon insertion; that is, the translocation and insertion of transposable elements can cause [mutations](#). Indeed, in some organisms they are the predominant form of spontaneous mutagenesis. Probably because of this potential for lethal mutagenesis of the host, transposition frequency is generally very low, often less than one translocation per thousands of cell divisions. One strategy for controlling transposition is to highly regulate and limit the expression of the transposase; indeed, in some elements, transposase expression is limited at virtually every level of gene expression. Host factors that participate in or influence transposition may also play roles in controlling transposition.

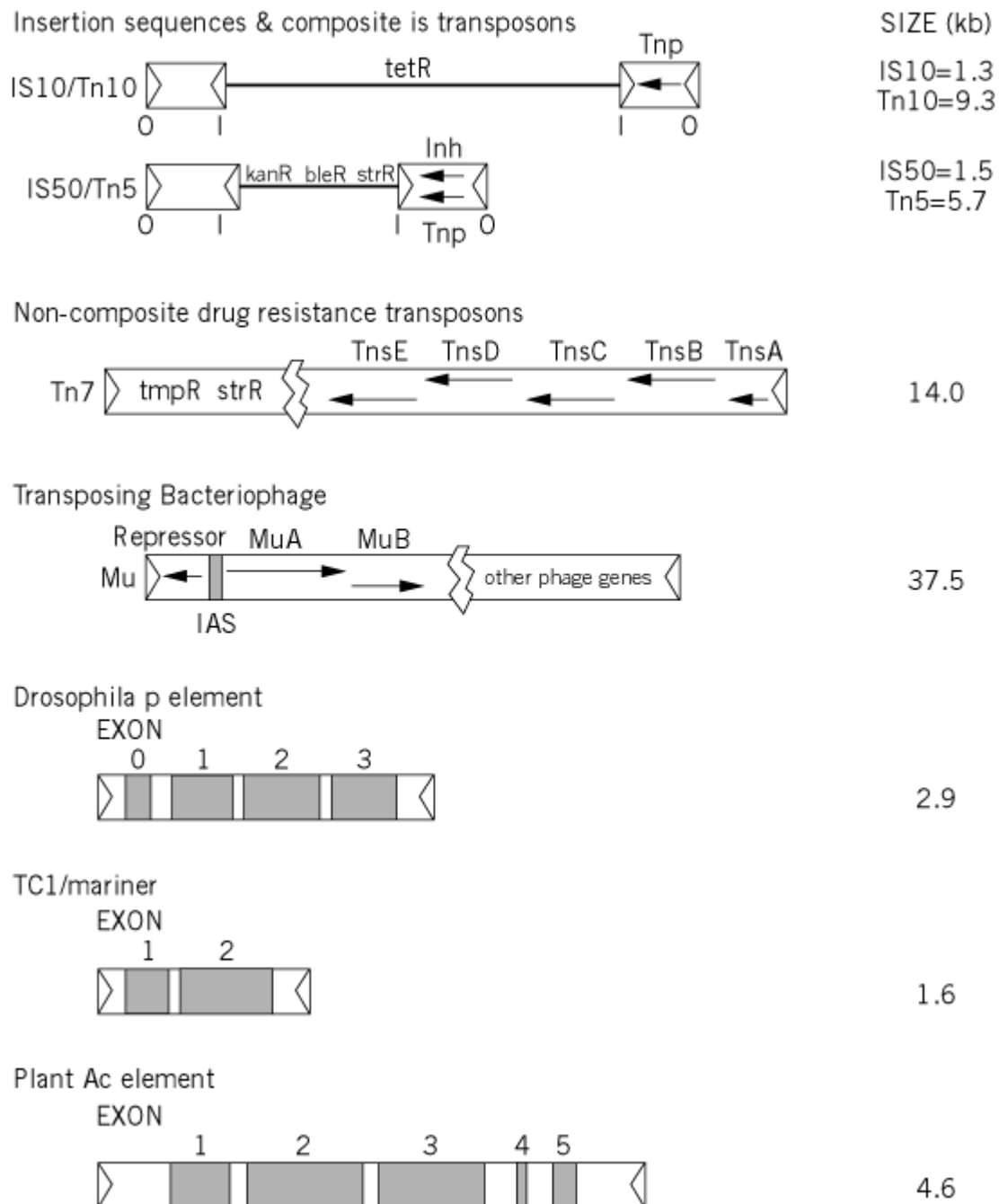
Transposable elements were discovered in maize in the 1940s by Barbara McClintock ([1](#)). While looking for deletions of certain regions of genomic DNA, she observed a large number of “unstable mutations”; that is, although the mutant phenotype was present in early cells of developing kernels, the mutation “reverted” to wild-type during the later stages of growth. We now know that this instability results from the insertion of a transposable element to cause the mutation and then the excision of the element to restore it to wild-type; chromosome breakage can also result from the excision of the element.

McClintock's great insight was to realize that there were discrete genetic elements that could both influence the expression of host genes, including causing unstable mutations, and move from place to place in the genome. She found such elements at the sites of certain unstable mutations she examined. One such element, which she called Ds (for dissociation, because she also found it to be associated with chromosome breaks), we now know is a deleted transposable element that has recombination sequences at its tips but does not encode a transposase. She deduced that this Ds element was activated by another discrete element elsewhere in the genome that she called Ac (for activator); we now know that activator is an intact transposable element that makes a transposase that can act both on both Ac and Ds because they have highly related terminal recombination sequences. She determined that these elements could both affect gene expression and move from place to place; both were novel and remarkable properties. To emphasize the effect that these elements could have on host gene expression, she called them “[controlling elements](#).” She also believed that the effects of controlling elements on gene expression were critical to development ([2](#)). These mobile controlling elements are now generally referred to as *transposable elements*, emphasizing their ability to translocate from place to place.

McClintock discovered and described these [mobile elements](#) in genetic terms. It was not until the late 1960s, when transposable elements were discovered in bacteria, that mobile elements acquired a molecular description. Several groups isolated unusual mutations in bacterial [operons](#) that were highly polar on downstream gene expression ([3](#), [4](#)). These mutations could be analyzed physically because they could be located on bacteriophage DNA, which is far, far smaller than the maize genome. Such molecular analysis revealed that these polar mutations resulted from the insertion of small discrete pieces of DNA that contained [transcription](#) terminators. Mobile DNA in bacteria was also discovered by the analysis of antibiotic resistance genes that could translocate between plasmids. Subsequently, transposable elements have been discovered and described in many organisms; particularly well-studied elements are the [P element](#) of *Drosophila* and the Tc1 element of [nematodes](#) and its close relative mariner, found in *Drosophila* ([5](#)) (Fig. [1](#)). Another important step in the study of transposable elements was the realization that [retroviruses](#), such as human immunodeficiency virus ([HIV](#)), and some nonviral elements similar to retroviruses, such as the [Ty elements](#) of yeast, are also transposable elements ([6](#)).

**Figure 1.** Examples of transposable elements. Well-studied transposable elements from a variety of organisms are

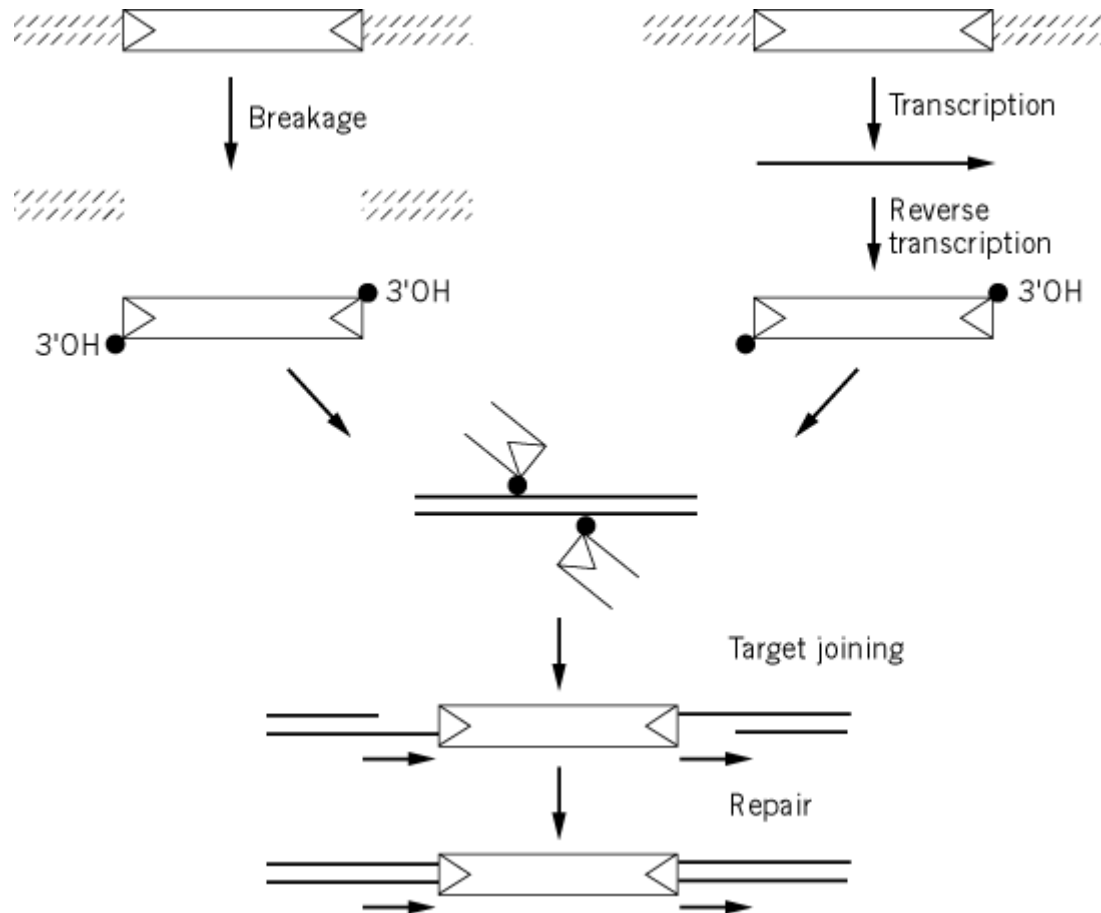
shown. The special recombination sequences present as inverted at the ends of the elements are indicated by triangles. All the elements encode proteins for transposition as indicated by arrows (top three sections) or gray boxes (exons of eukaryotic elements). Some elements encode other determinants; for example, IS50 encodes a transposition inhibitor, many elements encode antibiotic resistance determinants, and Mu encodes a viral genome.



The mechanism of transposition has been studied intensively, facilitated by the development of *in vitro* transposition systems for a number of elements. Many elements transpose via a “cut and paste” mechanism, in which the transposase binds specifically to the ends of the element and introduces double-strand breaks that excise the element from the donor backbone (7); a key feature of this step is that the 3’OH ends of the element are exposed (Fig. 2). The 3’OH ends of the element then join to the target DNA at staggered positions, one end on each strand; these positions have a spacing characteristic of each element. Thus the newly inserted element is flanked by short gaps that are repaired by the host [DNA repair](#) machinery; this repair results in target sequence duplications that flank the newly inserted element.



**Figure 2.** The mechanism of transposition. Reaction pathways for both DNA-only elements (left path) and RNA + DNA elements (right path) are shown. The key first step is to expose to the 3'OH ends of the element, usually by transposase cleavage, although in some retrotransposons the 3'OH ends are generated by cDNA synthesis. The 3' OHs at each end of the element join to staggered positions of the target, often 5 or 9 bp apart. In the strand transfer product, the transposon and target DNAs are covalently linked through the 3' ends of the element. The 5' ends of the transposon are flanked by short gaps (determined by the spacing of the staggered positions of joining). These gaps are filled in by the host repair machinery.



The transposition of a retrovirus involves the same DNA breakage and joining steps (7). Although one form of the retrovirus is RNA, the actual substrate for transposition is a DNA copy of this RNA generated by **reverse transcription**. In some cases, reverse transcription leaves the 3'OH ends of the element exposed, and in other cases DNA cleavages are necessary to expose them. As with DNA-only elements, the 3'OH ends of the retroviral DNA attack the target DNA at staggered positions, generating an insertion product flanked by small gaps that are subsequently repaired by the host.

All transposable elements that have been examined use this basic chemical scheme of exposing 3' OHs and using them to attack the target DNA. It is also likely that the transposases that execute these reactions from bacteria to yeast to *Drosophila* to humans are fundamentally related.

#### Bibliography

1. B. McClintock (1948) Carnegie Inst. Washington Yearbook **47**, 155–169.
2. B. McClintock (1956) Cold Spring Harbor Symp. Quant. Biol. **21**, 197–216.
3. E. Jordan, H. Saedler, and P. Starlinger (1968) Mol. Gen. Genet. **102**, 353–365.

4. J. A. Shapiro (1969) *J. Mol. Biol.* **40**, 93–105.
5. H. Saedler and A. Gierl (1996) *Curr. Top. Microbiol. Immunol.* **204**, 27–48.
6. J. D. Boeke and J. P. Stoye (1997) In *Retroviruses* (H. Varmus, S. Hughes, and J. Coffin, ed.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 343–435.
7. K. Mizuuchi (1992) *J. Biol Chem.* **267**, 21273–21276.

### Suggestions for Further Reading

8. *The Dynamic Genome: Barbara McClintock's Ideas in the Centure of Genetics* (1992) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
9. B. B. Finlay and S. Falkow (1997) Common themes in microbial pathogenicity revisited (review). *Microbiol. Mol. Biol. Rev.* **61**(2), 136–139.
10. J. D. Boeke and S. E. Devine (1998) Yeast retrotransposons: finding a nice quiet neighborhood. *Cell* **93**, 1087–1089.
11. M. L. Pardue, O. N. Danilevskaya, K. Kowenhaupt, F. Slot, and K. L. Traverse (1996) *Drosophila* telomeres: new views on chromosome evolution. *Trends Genet.* **12**(2), 48–52.
12. D. B. Roth and N. L. Craig (1998) VDJ recombination: a transposase goes to work. *Cell* **94**, 1–20.

## Transposase

Transposase is the [recombinase](#) that promotes the movement of its cognate [transposable element](#) from place to place in a [genome](#). The transposase binds to special specific [inverted terminal repeat](#) sequences at the ends of a transposable element and executes the DNA breakage and joining reactions that underlie [transposition](#). Recognition of these specific sequences is the basis of recognition of a cognate element by a transposase.

Transposase executes multiple steps in [recombination](#). It specifically recognizes an element by binding to particular sequences at the element termini. These special sequences are usually arranged as inverted repeats, so the interaction of transposase with these sequences positions a transposase molecule at each end of the element in identical fashion. For the elements in which mutagenesis has been used to evaluate the functions of these inverted repeats ([1](#), [2](#)), it has been determined that the internal side of an inverted repeat consists of specific sequences that are the actual transposase binding site. This binding sequence itself is separated from specific sequences at the extreme tips of the element by “spacer” base pairs, whose actual sequence is less critical to recombination. The tip base pairs are the key in specifying the positions for cleavage and joining. In some cases, there are other transposase binding sites internal to the inverted repeats that assist the assembly of appropriate transposase–transposon end complexes and in mediating the DNA breakage and joining reactions. The active form of transposase in the systems that have been examined is a multimer; for example, IS10 transposase functions as a dimer, with one transposase acting at each element end ([3](#)).

Transposase also interacts with the target DNA and is often intimately involved in choosing a target site ([4](#)). Some elements recognize specific features in the target DNA [for example, retroviral transposases ([integrases](#)) interact preferentially with bent DNA ([5](#))], while other transposases recognize particular sequences in the target DNA [for example, the bacterial element Tn10 preferentially recognizes sequences related to a 9-bp **consensus sequence** ([6](#), [7](#))]. Following end recognition, transposase then mediates the synapsis of the element ends with each other, thus

identifying a transposable element with two ends and juxtaposing the ends so they can act coordinately.

The transposase then executes the breakage reactions at the ends of the element that disconnect it from the donor backbone (8). In many cases, these are double-strand breaks that completely sever the element from the donor DNA; in some other cases, breaks happen only at the 3' ends. In some elements, these cleavages occur only upon end synapsis; for other elements, interaction of the ends with the target DNA is also required for the initiation of DNA breakage (9, 10). Once the 3' ends of the transposon are exposed by these transposase cleavage reactions, transposase then promotes the attack of the 3'OH ends on the target DNA at staggered positions, often separated by 5 or 9 bp. Thus the transposon becomes covalently linked to the target DNA.

The *in vitro* analysis of transposable elements from bacteria to yeast to [nematodes](#) to humans has revealed that all of these reactions occur by a fundamentally similar mechanism (8): in all cases, recombination is Mg<sup>2+</sup>-dependent and involves cleavages to expose the 3' ends of the elements and the attack of these 3'OH ends at staggered positions on the target DNA; for some elements, cleavages at the 5' ends of the transposon can also occur, completely disconnecting the element from the donor backbone. Recent **X-ray crystallographic** studies of transposases from bacteria [the [Mu phage](#) transposase (11)] and humans [the HIV integrase (12)] have revealed that although there is little sequence homology between these proteins, the same global three-dimensional structure exists in the catalytic regions. Most strikingly, in each recombinase there is a cluster of acidic residues (the D,D,E motif) that forms a Mg<sup>2+</sup> binding site that is likely the key to the active site for DNA breakage and joining. Although displaced from each other in the primary sequence, the folding of both recombinases closely juxtaposes these acidic residues in space. This D,D,E motif has been observed in many transposases (13, 14), and it is likely that many other transposases are also members of the retroviral integrase superfamily.

The activity of the transposase may also require or be influenced by other proteins. For example, some systems require accessory host DNA-bending proteins that also interact with the element termini to facilitate the assembly of particular transposase–transposon end structures (15-17). Proteins other than transposase can also play a key role. These other regulatory proteins can also play a key role in target-site selection (4).

## Bibliography

1. D. Haniford and N. Kleckner (1994) *EMBO J.* **13**, 3401–3411.
2. J. Sakai, R. M. Chalmers, and N. Kleckner (1995) *EMBO J.* **14**, 4374–4383.
3. S. Bolland and N. Kleckner (1996) *Cell* **84**, 223–233.
4. N. L. Craig (1997) In *Annual Review of Biochemistry*, Annual Reviews, Inc., Palo Alto, CA, pp. 437–474.
5. P. M. Pryciak and H. E. Varmus (1992) *Cell* **69**, 769–780.
6. J. Bender and N. Kleckner (1992) *EMBO J.* **11**, 741–750.
7. J. Bender and N. Kleckner (1992) *Proc. Natl. Acad. Sci. USA* **89**, 7996–8000.
8. K. Mizuuchi (1992) *J. Biol. Chem.* **267**, 21273–21276.
9. R. J. Bainton, K. M. Kubo, J.-N. Feng, and N. L. Craig (1993) *Cell* **72**, 931–943.
10. M. Mizuuchi, T. A. Baker, and K. M. Mizuuchi (1995) *Cell* **83**, 375–385.
11. P. Rice and K. Mizuuchi (1995) *Cell* **82**, 209–220.
12. F. Dyda, A. B. Hickman, T. M. Jenkins, A. Engelman, R. Craigie, and D. R. Davies (1994) *Science* **266**, 1981–1986.
13. J. Kulkosky, K. S. Jones, R. S. Katz, J. P. G. Mack, and A. M. Skalka (1992) *Mol. Cell Biol.* **12**, 2331–2338.

14. P. Polard and M. Chandler (1995) *Mol Microbiol.* **15**, 1–23.
15. R. Chalmers, A. Guhathakurta, H. Benjamin, and N. Kleckner (1998) *Cell* **93**, 897–908.
16. B. D. Lavoie and G. Chaconas (1993) *Genes Dev.* **7**, 2510–2519.
17. B. D. Lavoie and G. Chaconas (1996) *Curr. Top. Microbiol. Immunol.* **204**, 83–102.

## Transposition

[Transposable elements](#) are discrete DNA segments that can move between nonhomologous positions within a [genome](#) and have been found in virtually all organisms examined. The [recombination](#) pathway by which such elements move is called *transposition*. Most elements encode a [transposase](#), that is, the [recombinase](#) that executes the DNA breakage and joining reactions that underlie transposition, as well as special recombination sequences at the ends of the transposon arranged as [inverted terminal repeats](#) that include transposase binding sites; elements lacking a transposase can often be mobilized by the transposase from another cognate element. The insertion of a transposable element into a new insertion site alters the host DNA at that point and often results in a [mutation](#) through gene disruption.

Transposition requires the binding of transposase to both ends of the transposable element and subsequent synapsis of the ends. This requirement for synapsis prior to DNA breakage and joining ensures that an intact two-ended element is present. Many elements—from bacteria to *Drosophila* to fish—transpose by a “cut and paste” mechanism, in which  $Mg^{2+}$ -dependent double-strand breaks at the ends of the element separate the [transposon](#) from the donor backbone; that is, the transposable element is excised from the donor site (1) (see Fig. 2 in [Transposable Elements](#)). The cleavages that expose the 3'OH ends are key, because these reactions expose the transposon ends that transposase will join to the target DNA. Cleavage of the other strand at the 5' ends of the element can occur by a variety of means, including simple endonucleolytic cleavage by the transposase or using an alternative transposase subunit specific for 5' ends cleavage (2).

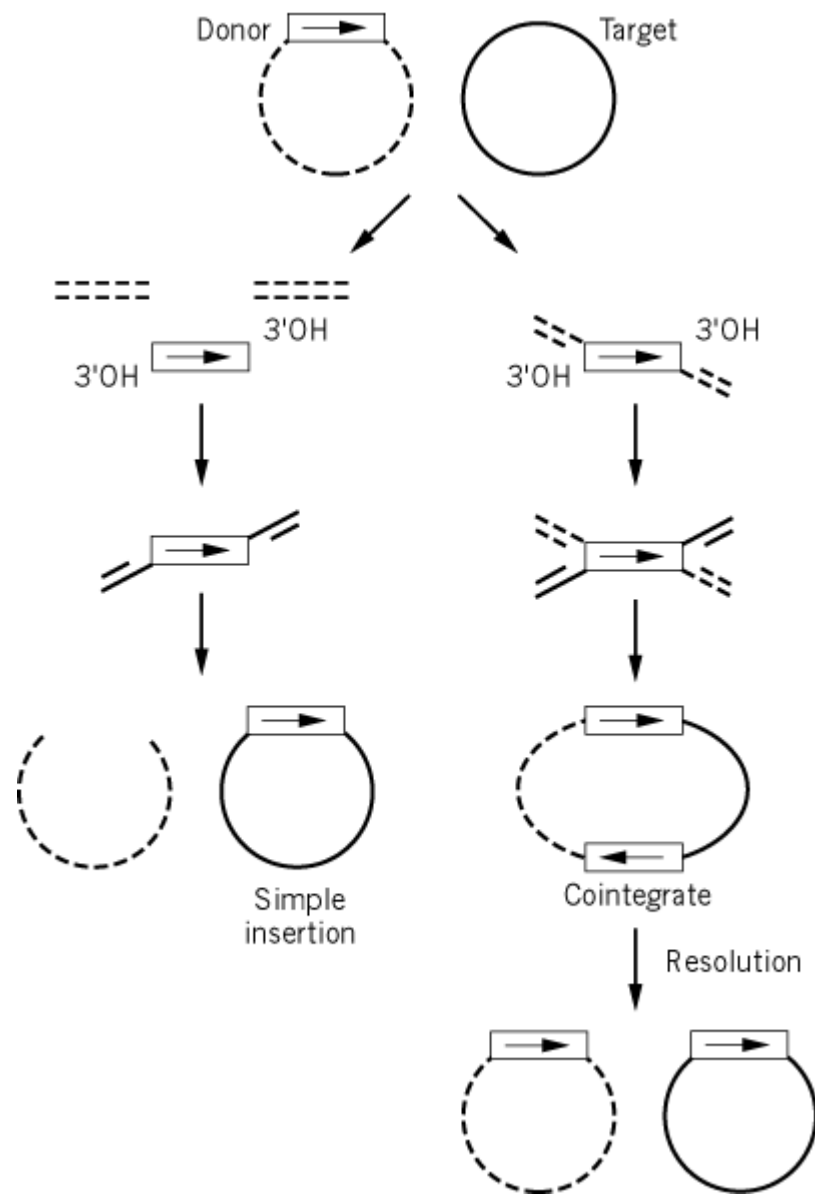
In the subsequent target-joining step, the exposed 3'OH ends of the transposon directly attack the target DNA. It is important to appreciate that this excision and insertion cycle occurs in the absence of any covalent protein–DNA intermediates; instead, these steps occur via simple one-step transesterification reactions in which  $H_2O$  attacks the transposon end to disconnect it from the donor DNA, and the transposon end joins to the target DNA by the direct attack of the 3'OH transposon end (1). This attack of the transposon ends occurs at staggered positions on the target DNA—that is, the positions of end joining displaced from each other by several nucleotides. Because of these staggered positions of attack, a small gap the size of the stagger flanks each end of the newly inserted element. Repair of this gap occurs by the host [DNA repair](#) functions to generate intact duplex DNA, in which the newly inserted transposon is flanked by a target sequence duplication.

The transposition of retroviruses and retroviral-like transposons occurs by making a [messenger RNA](#) copy from the DNA provirus present in a genome, which is then turned into a DNA copy of the element by **reverse transcriptase** (3). The actual DNA breakage and joining steps that underlie the insertion of this DNA copy into a target DNA are very similar to the mechanisms described above for bacterial elements and for elements from *C. elegans* and *Drosophila*. The ends of the viral DNA are often trimmed to expose the actual 3'OH ends of the transposon, and the resulting 3'OH ends the attack the target DNA at staggered positions. As with the excision–insertion elements described

above, the newly inserted transposon is flanked by target-sequence duplications resulting from the repair of these gaps.

For some bacterial elements, in particular [Mu phage](#) and probably also Tn3, transposition involves 3' end cleavage and joining of the ends to the target DNA, but no 5' end cleavage. The products of these transposition reactions are different than the simple insertions that arise by the cut-and-paste pathway (Fig. 1). Transposon end cleavage occurs only at the 3' ends of Mu; the 5' ends of Mu remain attached to the donor DNA. Thus, when the 3' ends attack the target DNA, the transposon is still linked to the donor DNA via its 5' ends and is now also linked to the target DNA. This structure is variably called a *fusion product*, a *strand transfer product*, or a *Shapiro intermediate*. The exposed target ends that flank the newly inserted transposon have 3'OH ends that can serve as primers for [DNA replication](#). Such replication results in a structure called a *cointegrate*, which contains two transposon copies linked by the donor backbone and the target DNA. In this reaction, as opposed to a simple cut-and-paste reaction, the transposon is copied during recombination; thus this type of reaction is called *replicative transposition*.

**Figure 1.** Simple insertion and cointegrate products of transposition. Some elements transpose through a cut-and-paste mechanism in which the transposon is completely excised from the donor site and is then inserted into the target DNA. Other elements, such as phage Mu and Tn3-like elements, carry out replicative transposition in which an additional copy of the element is made by DNA replication. In this pathway, recombination begins with cleavage to expose the 3' ends of the element and the resulting 3'OH termini attack the target DNA. Replication of the resulting joint structure from 3'OHs in the target DNA that flank the newly inserted element results in a molecule called a cointegrate in which two copies of the transposon link the donor and target replicons. Recombination between the directly repeated copies of the element can generate two species, one which looks like an intact donor molecule and the other which looks like a simple insertion. Tn3-like elements encode a special resolvase enzyme and recombination site within the element that very efficiently promote resolution.



For Mu, the final product is a cointegrate; Tn3 cointegrates are, however, processed further. Tn3 encodes both a transposase and a **resolvase**, an enzyme that can act at special sites to exchange DNA duplexes. Resolvase action on the cointegrates results in a target molecule that contains a copy of the transposon and regeneration of a donor molecule containing the transposon. It should be noted that different transposition reactions can both give rise to simple insertion products; for example, both cut-and-paste elements and Tn3-like elements yield similar insertions as their final product, but are generated by distinct mechanisms. It is very difficult to infer the recombination mechanism from *in vivo* studies, because the products observed may have undergone other recombination reactions prior to or after the actual transposition event.

Another feature that affects whether transposition appears replicative or not is how the gapped donor DNA is dealt with after translocation of an element to a new site. In some cases, the gapped backbone is repaired by double-strand gap repair using a sister chromosome as a template (4-6); in this case, the donor can be restored to its transposon-containing state. When such repair occurs, transposition appears replicative; that is, there is one transposon copy at the donor site and one at the target site, although transposition itself occurred by a nonreplicative cut-and-paste mechanism, and the transposition copy at the donor site was generated by homologous recombination. In other cases, the gapped donor is repaired by an end-joining reaction, generally without restoration of the donor to

the original pre-transposon state. Thus transposition can leave “footprints” that may alter gene expression at the insertion site, even though the transposable element itself is no longer present (4, 6, 7). With retroviruses and retrotransposons, the “donor” provirus site is not altered, because the translocation substrate is a DNA copy of the element made by reverse transcription of an RNA copy.

## Bibliography

1. K. Mizuuchi (1992) *J. Biol. Chem.* **267**, 21273–21276.
2. E. W. May and N. L. Craig (1996) *Science* **272**, 401–404.
3. J. D. Boeke and J. P. Stoye (1997) In *Retroviruses* (H. Varmus, S. Hughes, and J. Coffin, ed.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 343–435.
4. W. R. Engels, D. M. Johnson-Schlitz, W. B. Eggleston, and J. Sved (1990) *Cell* **62**, 515–525.
5. A. T. Hagemann and N. L. Craig (1993) *Genetics* **133**, 9–16.
6. R. H. Plasterk (1991) *EMBO J.* **10**, 1919–1925.
7. E. S. Coen, T. P. Robbins, J. Almeida, A. Hudson, and R. Carpenter (1989) In *Mobile DNA* (D. A. Berg and M. M. Howe, ed.), American Society for Microbiology, Washington, D.C., pp. 413–436.

## Transposon

Transposon is another word for a [transposable element](#)—that is, a DNA segment that can move between nonhomologous positions in a [genome](#) via a type of [recombination](#) called [transposition](#). A transposable element generally encodes a [recombinase](#), a [transposase](#) that executes the DNA breakage and joining reactions that underlie transposition, by acting on special recombination sequences, [inverted terminal repeats](#), at the ends of the element. Transposable elements are present in Eubacteria, Archea, and Eucaryotes—indeed, in virtually every organism that has been examined. The term “transposon” is often, but certainly not exclusively, used in conjunction with bacterial transposable elements that encode, in addition to their transposition functions, other determinants such as [antibiotic-resistance](#) genes. The structure of a number of bacterial transposons comprises special recombination sequences up to perhaps 150 bp in length at the termini of the element, and the interior of the element encodes a transposase (and possibly other transposition proteins) and a gene (s) for another determinant, for example, an [antibiotic-resistance](#) gene. Small elements that lack determinants other than transposase are usually called **insertion sequences**.

Some transposons are actually composite elements in which two insertion sequence (IS) elements flank a determinant such as an antibiotic-resistance gene. One IS element may move independently of the other, and both IS elements can collaborate to move the entire IS–drug resistance–IS segment. For example, the transposon Tn10 is composed to two IS10 elements flanking a segment encoding a tetracycline-resistance determinant. The transposon Tn5 consists of 2 IS50 elements flanking a **kanamycin**-resistance determinant. Transposase supplied by both (or sometimes one) of the IS elements acts on the outside ends of the composite element to move them from place to place.

Despite the custom of using the terms *insertion sequences* and *transposons* to describe certain bacterial elements, it is important to remember that both types of DNA segments are actually transposable elements, highly related to many eukaryotic transposable elements in structure and transposition mechanism. However, some types of transposable elements have, to date, been identified in bacteria. For example, composite transposons with two internal individually transposable elements (IS elements) have not been identified to date in eukaryotes. Furthermore,

elements that execute replicative transposition such as elements bacteriophage Mu and Tn3 have also not been identified in eukaryotes.

## Transposon Tagging

**Transposons** are mobile genetic elements that move (transpose) from one region of the **genome** to another. Transposons are used as tools for gene **cloning** because insertion of a transposon into a **gene** disrupts its function, often producing a visible mutant **phenotype**. When the **DNA** sequence of the transposon is known, it is possible to clone the disrupted gene by using the transposable element as a “tag” to identify the segment of DNA harboring the element. Transposon tagging involves inducing transposition, screening for **mutations** caused by transposon insertion, identifying the element causing the mutation, and cloning the tagged gene.

Transposons are found in almost all organisms where they have been looked for, including **bacteria**, **yeast**, **plants**, and mammals. Transposons from one organism often also transpose in heterologous systems. Therefore it is possible to use transposon-tagging methods to clone genes in a wide variety of organisms, including those where there are few other tools available for gene isolation.

### 1. Characteristics of Transposons Relevant to Tagging

There are many families of transposons, and they are generally divided into two classes. Class I is the **retrotransposons**, and Class II is the DNA-based transposons. Class I transposons are used for tagging in mammals and yeast, and class II transposons are those used most frequently in tagging schemes for bacteria, plants, *Drosophila*, and **Caenorhabditis elegans**. Class II elements have inverted repeats at their ends and produce the products needed for their own excision and insertion, termed **transposase**.

#### 1.1. Regulation of Transposition

It is important to have some method for regulating the frequency of transposition. Certain families of transposable elements contain-“**nonautonomous**” **controlling elements** that cannot produce transposase but move in the presence of the transposase produced by an “autonomous” family member. A similar “two-element” system is often employed in transposon-tagging schemes where a stable transposase source (eg, a transposon immobile due to deletion of one of its **direct repeats**) is used to mobilize nonautonomous elements. This system has the advantage that, once new insertions have been generated by combining the two types of elements, the transposase source can be segregated. This stabilizes the new insertions and thus stabilizes any mutant phenotypes generated. Transposase can be reintroduced later to remobilize the element to revert the mutant phenotype or to create new alleles of the tagged gene.

In some organisms, the frequency of transposition is regulated by environmental factors, such as temperature. For example, in *Antirrhinum*, the transposition frequency of Tam3 is increased 1000-fold in plants grown at 15°C rather than 25°C (1). In *Saccharomyces cerevisiae*, transposition of Ty1 is also fairly high at low temperatures but decreases 20-fold at 30°C (2).

#### 1.2. Obtaining the Sequence of the Transposon “Tag”

To use a transposon as a molecular “tag” for gene cloning, it is necessary to use a transposon with a known DNA sequence. Many transposons have already been well characterized. New transposons can be “trapped” by first mobilizing transposition and then identifying transposon insertions into known genes (generally by looking for unstable **alleles** of the gene). Then the mutant allele is sequenced, thus determining the sequence of the inserted transposon.



### 1.3. Integration Site Preference

The efficiency of transposon tagging any particular gene depends on a number of factors. The perfect transposon for tagging purposes would have no preference for the integration site in the host genome, but many, if not most, transposons show some type of site preference. Preferences range from that of the *C. elegans* Tc1 element, which has a TA dinucleotide target site (3), to that of the *Drosophila* P element which preferentially inserts at the 5'-ends of transcription units (4). Such target site preferences explain why transposon-tagged alleles are never recovered for some genes, whereas many tagged alleles have been discovered for other genes.

### 1.4. Transposition to Linked Sites

A feature of many transposons is their tendency to transpose to linked sites. This is useful for some types of screens (see “Directed tagging” later) but makes random transposon mutagenesis more difficult. In some cases, methods have been devised to select for transposition to unlinked sites. An *Arabidopsis* Ds element carrying a positive selectable marker was linked to IAAH (a marker that can be selected against). The Ds element was mobilized, and by selecting for the Ds marker and selecting against IAAH, lines were identified in which the transposon had moved to a position unlinked to IAAH (5).

### 1.5. Endogenous Versus Heterologous Transposons

Another consideration in transposon-tagging experiments is whether to use endogenous unmarked transposons, which are plentiful in some backgrounds, or whether to construct transposons specifically for transposon tagging and to mobilize a single transposon or a few transposons per genome. The greater the number of mobile transposons and the higher the transposition frequency, the greater the chance of tagging any particular gene. A large number of background transposons, however, complicates the linkage studies that must be done to identify the specific transposon inserted into the gene of interest. High rates of transposition also result in unstable mutations, making it difficult to link a transposon with a mutant phenotype. The use of only one or a few transposons per genome greatly facilitates identification of the transposon inserted into the gene of interest and thus speeds up cloning of the gene. The trade-off is that more individuals must be screened to obtain insertions into the gene of interest. Using a heterologous two-element, transposon-tagging system is one way of ensuring that only a few elements are mobilized in the genome in the presence of the transposase and allows removing the transposase by segregation to prevent further transposition.

### 1.6. Transposon and Transposase Engineering

It is useful to have transposons designed so that (1) they carry selectable markers and (1) they carry part of a **plasmid** that can be selected for in *E. coli* to facilitate cloning of flanking DNA by plasmid rescue. A marker linked to the transposase source facilitates removal of the transposase later by segregation from the tagged gene. It is also sometimes possible to increase transposition frequencies by deleting or altering part of the transposase gene, by using a powerful constitutive **promoter** to drive transposase, or by altering transposon size (6-8).

## 2. Transposon Tagging Strategies: Obtaining Insertions into Genes

### 2.1. Random Mutagenesis

In random mutagenesis, transposons are mobilized to create a **library** of individuals with different transposon insertions. Then the library of insertions is screened for mutant phenotypes of interest. In a two-element system, transposons are mobilized by putting stable transposase into the background of the nonautonomous transposons, either by genetic crosses or, in the case of the *Drosophila* P element, by **microinjecting** the transposase into an embryo containing nonautonomous elements (9). In other organisms, growth at low temperature mobilizes transposition. Dominant visible mutations are seen in the M1 generation, and recessive visible mutations are seen segregating in the M2 progeny of M1 individuals. Individuals containing a transposed element are often termed “transposants”.

## 2.2. Directed Mutagenesis

There are two types of directed tagging schemes. In both cases, the target gene has been identified previously by its mutant phenotype. In the first type of directed mutagenetic scheme, insertion into a specific gene is selected for by crossing an individual homozygous for the previously identified recessive mutation with a wild-type individual carrying an autonomous transposon (alternatively, an individual carrying both a stable transposase source and the mutation of interest is crossed with a wild-type individual carrying a nonautonomous transposon). The next generation, the M1, is screened for individuals exhibiting the mutant phenotype. These individuals carry both the original mutation and a new transposon-induced allele of the same gene.

The second type of directed mutagenetic experiment is useful in systems where the transposon is known to move preferentially to linked sites on the [chromosome](#). A mapped transposable element that is known to be linked to the gene of interest is mobilized. Insertions into the gene of interest are identified in the M1 by the strategy outlined previously, or the linked transposon is mobilized in a wild-type background, and the M2 is screened for mutations. In either case, if there is a strong preference for transposition to linked sites, the frequency of mutations induced in the specific gene should increase relative to a random tagging procedure ([10](#)).

## 2.3. Enhancer and Gene Traps

Enhancer traps are transposons designed to identify genes on the basis of their expression patterns. The transposon carries a [reporter gene](#) under the control of a weak or minimal promoter. Insertion of the transposon near an [enhancer](#) causes the reporter gene to be expressed in a pattern governed by that enhancer. In this way, it is possible to identify genes having interesting tissue-specific or developmentally regulated patterns. Enhancer-trap screens should identify genes not usually found in regular tagging screens, such as those having no obvious phenotype when disrupted. They can also be used to suggest new roles for genes. For example, nearby insertion of an enhancer trap could reveal interesting adult expression patterns in a gene that, when disrupted, is lethal in an embryo. Enhancer-trap lines also provide useful tissue, cell-type, and developmentally specific markers that can be analyzed in different mutant backgrounds or under different conditions. Enhancer-trap screens are successfully used to identify interesting gene expression patterns in a number of different organisms, including *Arabidopsis* ([5](#)) and *Drosophila* ([11](#)).

Gene traps are similar to enhancer traps, except that the transposon used for tagging carries a reporter gene lacking a promoter. Thus, for reporter gene expression, the transposon must insert downstream of the promoter of an expressed gene. RNA **splicing** -acceptor sequences are introduced upstream of the coding region of the reporter gene to permit fusion of the reporter gene to the tagged gene if the transposon inserts into an **intron** ([5](#)).

Cloning genes identified by a gene-trap transposon is often easier than cloning enhancer-trapped genes, because gene traps must insert in the correct orientation downstream of a promoter to allow reporter gene expression. Enhancers act at a distance, and it is possible that the gene normally regulated by the “trapped” enhancer is at some distance from the enhancer-trap insertion. To clone the endogenous enhancer-regulated gene, it is sometimes necessary to initiate a **chromosome walk** from the site of the transposon insertion.

## 3. Cloning tagged genes: from insertion to sequence

The first step in cloning a tagged gene is to identify a transposon linked to the mutation. This is done by [Southern blotting](#) DNA from progeny obtained from the putatively tagged mutant. A transposon-specific probe probes the blot to identify a band present in the homozygous mutant progeny and absent in wild-type progeny that do not segregate for the mutant phenotype. If many transposons segregate in the background, finding a transposon linked to the mutant phenotype is difficult. A large number of progeny need to be examined, or the mutant must back-crossed to a strain lacking transposons to segregate background transposons. Once a linked transposon is identified, there are a number of different ways to clone the gene in which a transposon has inserted. Three commonly

used methods are described here.

1. **Plasmid rescue.** If the transposon used for tagging contains part of a selectable plasmid (eg, pBR322) at one end of the transposon, plasmid rescue is used to isolate a fragment of flanking host DNA. Genomic DNA from the tagged individual is digested with a [restriction enzyme](#) that releases the selectable plasmid from the transposon but does not cut within the plasmid itself. This creates a linear piece of DNA containing the plasmid sequence and a small fragment of flanking host DNA. The products of the digestion are religated at a high dilution to ensure intramolecular ligation products. The ligation products are **transformed** into *E. coli*, and then the cloned flanking DNA is isolated.
2. **Inverse polymerase chain reaction (IPCR).** IPCR is another method for isolating host sequences flanking a transposon insertion. As for plasmid rescue, genomic DNA from a tagged individual is isolated and digested with a restriction enzyme that releases the end of the transposon and a piece of adjoining host DNA. Ligation is used to circularize the linear host DNA–transposon fragment. PCR using two transposon-specific oligonucleotide primers, each reading outward from the ends of the transposon sequence into the flanking DNA, are used to amplify the flanking host DNA. Then the amplified product is cloned.
3. **Library screening.** A third option is to make a library using genomic DNA from the tagged mutant that has been digested with an enzyme that does not cut within the inserted transposon. The recombinant vector containing the transposon also contains some flanking host DNA. A transposon-specific probe is used to screen the library and identify clones containing the transposon and the flanking DNA.

#### 4. Confirmation

Once the flanking host DNA is cloned, it should be used to probe Southern blots of DNA extracted from homozygous mutant and wild-type plants to look for a band difference, indicating transposon insertion into the complementary sequences in the mutant. This is done to prove that the cloned DNA actually represents the host DNA flanking the transposon and is not an artifact of cloning. Then the cloned flanking sequences are sequenced directly or used as probes to obtain full-length sequences for further analysis. If the transposon has not inserted into the coding sequence of the disrupted gene itself (or in the case of an enhancer trap, if the transposon is not in a gene), it is necessary to use the flanking sequences for initiating chromosome walk to identify the gene of interest.

The most convincing way to prove that a mutant phenotype results from transposon disruption of the cloned gene is to rescue the mutant phenotype by transforming the cloned gene back into the mutant. It is also possible to isolate and analyze different alleles of the gene from independently derived mutants. If all of the mutant alleles harbor mutations in the cloned gene, it is good evidence that the correct gene has been identified. Phenotypic revertants that show a loss of the transposon from the cloned sequence, when they are examined by sequencing or Southern blotting techniques, also confirm the identity of the tagged gene. In the case of enhancer-trapped genes, where there is no mutant phenotype associated with the insertion, expression analysis of the cloned gene should closely mimic that of the enhancer-trap reporter gene.

#### 5. Transposon Tagging Strategies: From Gene to Mutant

##### 5.1. Site-Selected Mutagenesis

Often a gene is identified based on the basis of its expression pattern or its homology to another gene, but its phenotype is unknown. The goal of site-selected mutagenesis is to obtain a transposon insertion into a gene of known sequence to obtain an idea of the null phenotype of the gene.

To obtain a site-selected mutation in a gene of interest, a large population of transposants is generated and divided into pools. Two types of primers are needed for the **PCR** reaction: (1) primers

specific to the gene of interest and (2) transposon-specific primers, one reading “out” from each end of the transposon into the host DNA flanking the insertion. PCR is performed using one gene-specific primer, one transposon-specific primer, and DNA obtained from a pool of transposants. If an individual has an insertion in or near the gene of interest, the gene-specific primer and the transposon-specific primer are positioned so that a PCR product is amplified. Nested primers are often used to ensure specificity of the amplified band. Then pools of transposants from which a product is amplified are subdivided and rescreened until a single individual harboring the insert is obtained. Because the original individuals are often dead by the time tagged individuals are identified, it is necessary to maintain viable siblings of the original transposants for further studies. In the case of plants, this is simple because seeds from the transposants can be maintained for reasonable periods of time. In the case of *C. elegans*, a frozen “transposon insertion mutant bank” has been developed (10).

In some cases, the original insertion does not disrupt gene function sufficiently to generate a visible phenotype. Then, it is often necessary to remobilize the transposon and select for excision/reinsertion events that fully disrupt the gene.

## 6. Conclusion

Transposon tagging methods are valuable for several reasons. They are used to identify and clone numerous genes having visible phenotypes, and enhancer-trap and gene-trap methods are used to identify genes based on their expression patterns. Possession of a tagged allele of a gene is often a shortcut to cloning, compared with traditional methods, such as chromosome walking. Now, development of site-selected tagging screens allows identifying transposon insertions into genes of known sequence. Such insertions result in null alleles or are used to produce null alleles of genes for which no previous mutation is identified. Null alleles are extremely useful for identifying the functional roles of the gene under study. They are used in genetic studies, such as analyzing double-mutant interactions, and they provide null backgrounds for mutational analysis of a gene. Transposon-induced alleles are also used to produce an allelic series by remobilizing the transposon and selecting for reinsertion elsewhere in the gene or coding sequence. Multiple alleles of a gene can provide insights into its function. Somatic excision events are also used to analyze the role of a gene in different tissues or at different times during development. For all of these reasons, transposon tagging is and will continue to be a useful tool for gene cloning and gene analysis.

## Bibliography

1. E.S. Coen, T.P. Robbins, J. Almeida, A. Hudson, and R. Carpenter (1989) In *Mobile DNA* (D.E. Berg and M.M. Howe, eds.), American Society for Microbiology, Washington, D.C., p. 41.
2. J. D. Boeke (1989) In *Mobile DNA* (D.E. Berg and M.M. Howe, eds.), American Society for Microbiology, Washington, D.C., p. 357.
3. D.G. Moerman and R.H. Waterston (1989) In *Mobile DNA* (D.E. Berg and M.M. Howe, eds.), American Society for Microbiology, Washington, D.C., p. 540.
4. W.R. Engels (1989) In *Mobile DNA* (D.E. Berg and M.M. Howe, eds.), American Society for Microbiology, Washington, D.C., p. 451–452.
5. V. Sundaresan, P. Springer, T. Volpe, S. Haward, J.D.G. Jones, C. Dean, H. Ma, and R. Martienssen (1995) *Genes Dev.* **9**, 1797–1810.
6. I. Bancroft, A.M. Bhatt, C. Sjodin, S. Scofield, J.D.G. Jones, and C. Dean (1992) *Mol. Gen. Genet.* **233**, 449–461.
7. J. Swinburne, L. Balcells, S. R. Scofield, J.D.G. Jones, and G. Coupland (1992) *Plant Cell* **4**, 583–595.
8. N. Kleckner (1989) In *Mobile DNA* (D.E. Berg and M.M. Howe, eds.), American Society for Microbiology, Washington, D.C., p. 246.
9. L. Cooley, R. Kelley, and A. Spradling (1988) *Science* **239**, 1121–1128.

10. D. Long, J. Goodrich, K. Wilson, E. Sundberg, M. Martin, P. Puangsomlee, and G. Coupland (1997) *Plant J.* **11**, 145–148.
11. C. Wilson, R.K. Pearson, H.J. Bellen, C.J. O'Kane, U. Grossniklaus, and W.J. Gehring (1989) *Genes Dev.* **3**, 1301–1313.
12. R.R. Zwaal, A. Broeks, J. van Meurs, J.T.M. Groenen, and R.H.A. Plasterk (1993) *Proc. Natl. Acad. Sci. USA* **90**, 7431–7435.

### Suggestions for Further Reading

13. V. Walbot (1992) Strategies for mutagenesis and gene cloning using transposon tagging and T-DNA insertional mutagenesis, *Ann. Rev. Plant Physiol. Plant Mol. Biol.* **43**, 49–82. Useful description of strategies, although discussion of tagging in heterologous systems is now outdated.
14. B. I. Osborne, and B. Baker (1995) Movers and shakers: maize transposons as tools for analyzing other plant genomes. *Curr. Opin. Cell. Biol.* **3**, 406–413. Describes heterologous tagging for cloning genes in tobacco, tomato, and flax.
15. R. Hehl (1994) Transposon tagging in heterologous host plants, *Trends Genet.* **10**, 385–386.
16. A. C. Spradling, D.M. Stern, I. Kiss, J. Roote, T. Laverty, and G.M. Rubin (1995) Gene disruption using P transposable elements: An integral component of the *Drosophila* genome project, *Proc. Natl. Acad. Sci. USA* **92**, 10824–10830.
17. K. Kauser (1990) From gene to phenotype in *Drosophila* and other organisms, *BioEssays* **12**, 297–301.
18. L. Cooley, C. Berg, and A. Spradling (1988) Controlling P element insertional mutagenesis, *Trends Genet.* **4**, 254–258.

### Transverse Gradient Gel Electrophoresis (Tgge)

[Electrophoresis](#) through [polyacrylamide](#) gels (**PAGE**) is one of the most commonly used techniques for characterizing proteins and nucleic acids. It is especially useful to monitor changes in their [conformations](#) because the size and shape of a molecules are two of the main determinants of electrophoretic mobility through sieving gels (the other is the net charge). Electrophoresis, however, is primarily a comparative technique because little information is obtained solely from the mobility of a macromolecule under any single set of conditions. It is most informative when a change in mobility is observed. For example, unfolding of a **protein** or **nucleic acid** is apparent by a dramatic decrease in electrophoretic mobility. Phenomena whose conformation is induced to change under the influence of some agent, such as a **denaturant**, are readily studied by electrophoresis in transverse gradient gels. In such gels, the perturbing agent is included in the gel in a continuous gradient, perpendicular to the direction of electrophoretic mobility. In a sample applied uniformly across the top of such a gel, each molecule migrates under constant conditions but different from those adjacent to it on the gel. The population of molecules migrates under continuously varying conditions. Changes in the macromolecular structure at any point within the gradient are apparent from the change in mobility and give a characteristic shape to the electrophoretic band after electrophoresis. Such changes are most apparent if they occur abruptly, such as a cooperative change from form A to form B over a small part of the transverse gradient, when form A predominates at one extreme of the gradient and form B at the other extreme.

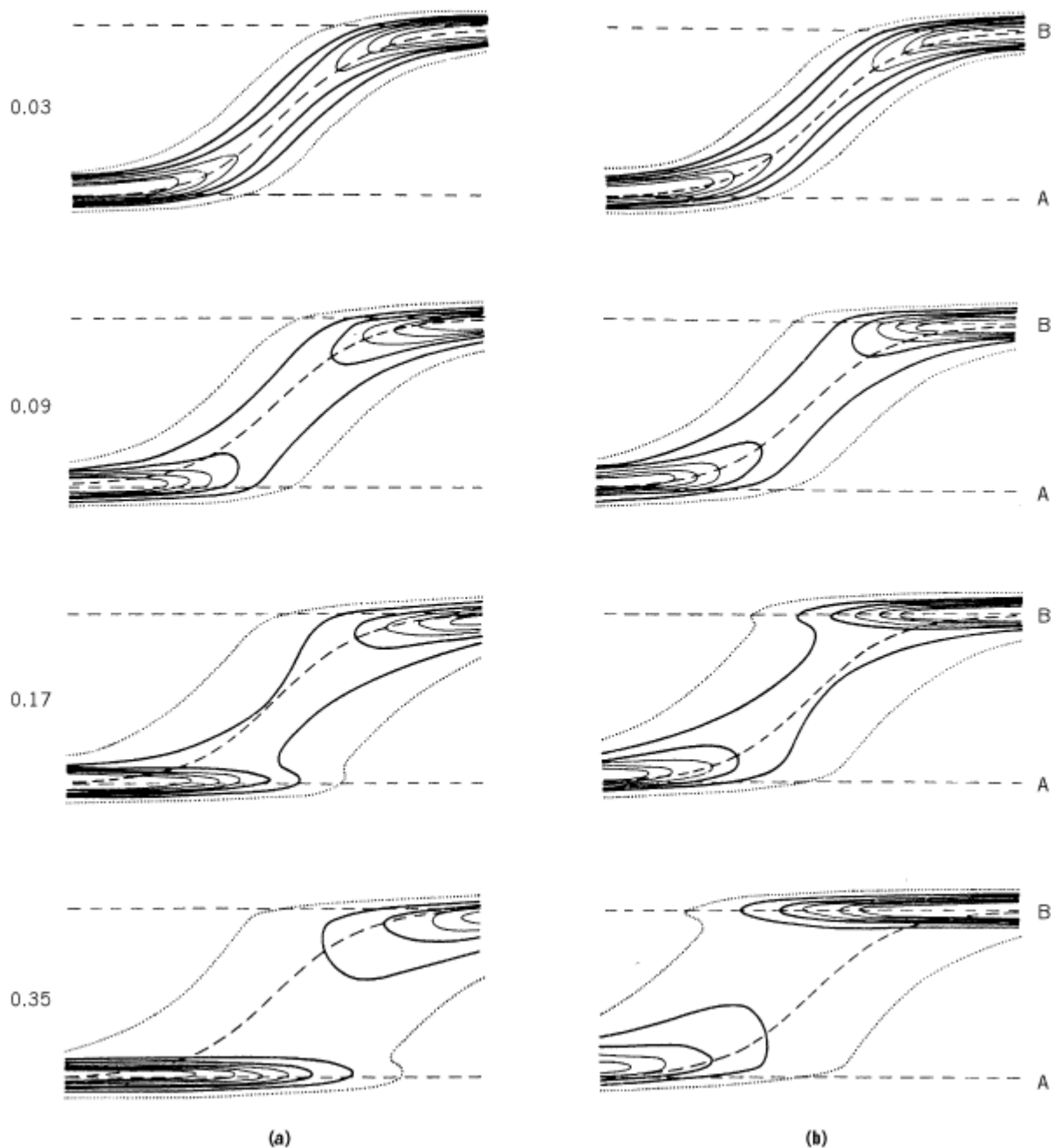
TGGE has the many advantages of electrophoretic techniques. The method is simple, rapid, and

requires no expensive equipment, only small quantities of sample are required, and the sample need not necessarily be homogeneous. Although only semiquantitative, TGGE has the advantage that close comparisons can be made of different macromolecules by running them on the same gel. Even small differences in their electrophoretic mobilities are detected.

When a macromolecule interconverts between two forms, say  $A \leftrightarrow B$ , as a result of the transverse gradient, a major question is whether one or two electrophoretic bands are observed at positions where the two species coexist. Numerical simulations indicate that this depends almost exclusively on the rate at which the two states are interconverted, relative to the duration of the electrophoretic separation (Fig. 1). If their interconversion is very slow, individual bands corresponding to both A and B are present to the extent that they were both present in the original sample. At the other extreme, with very rapid rates of interconversion and a **half-time** less than about 0.1 the duration of the electrophoretic separation, a single band will be present that has an intermediate electrophoretic mobility at the weighted average of the two. In this case, a change from one state to another at some point in a transverse gradient gel is apparent after the electrophoresis as a sigmoidal change in the mobility of a single, continuous band across the gradient. Then the final pattern describes the equilibrium between the two forms as a function of the variable in the gradient. With intermediate rates of interconversion, the electrophoretic patterns are smeared between the two positions of the two states A and B as various molecules convert to the other state at various times during the separation.

**Figure 1.** Simulated TGGE patterns expected for a macromolecule undergoing a two-state transition between forms A and B at four different rates and starting with the sample initially as either all A (left) or all B (right.) The transverse gradient was imagined to favor the A state on the left and the B state on the right, and the apparent equilibrium constant  $[A]/[B]$  varies logarithmically across the gradient from  $10^2$  to  $10^{-2}$ . It was assumed that both the forward and reverse rate constants vary inversely across the gradient by a factor of  $10^2$ . The rate of the interconversion is at a minimum at the middle of the transverse gradient, and the half-time is  $t_{1/2}$ . The rate there is expressed on the left as the ratio of  $t_{1/2}$  to the duration of the electrophoretic separation ( $t_{\text{elect}}$ ). Electrophoretic migration is from top to bottom, and state A has greater mobility. The distribution of the macromolecule in the final gel is given by the contours.

With the most rapid interconversion (top), the final pattern is a single, continuous band that describes the equilibrium between A and B across the gradient (dashed sigmoidal curve). The patterns are the same, irrespective of whether form A or B is applied to the gel. With a slower transition, where the interconversion takes place with a half-time comparable to the duration of the electrophoretic separation, the continuous band is replaced by a smear because the initial molecules convert to the other form only once and at various times during the electrophoresis. The form applied to the gel persists increasingly with a slower transition. In the limit of a very slow transition, only the form applied to the gel would be present after the electrophoretic separation. Adapted from Ref. 1.



If the starting macromolecule is homogeneous and all in one state, only one band should be apparent finally at each point in the gradient. Its mobility should be either that of the original state, if its interconversion is slow, or, if the interconversion is fast, of the final equilibrium mixture. The rates of such transitions can also be investigated more directly by comparing patterns obtained starting with the sample in either of the two states. Identical patterns are observed with rapid interconversions, but the original form persists if the interconversion is slow. If more than two states are possible, the resulting patterns may be more complex, but the same rules should apply to the interconversion of each pair of states.

The transverse gradient can be of any parameter that does not unduly affect the electrophoresis. Unfortunately, high concentrations of ionic species, such as **guanidinium chloride**, have a drastic effect on electrophoretic mobility and therefore cannot be included in electrophoretic gels. Small concentrations of other charged molecules making up the gradient are altered by their electrophoresis within the gel. Therefore reagents with a net charge are not ideal for TGGE, but any other agent that

does not interfere with electrophoresis may be used in the transverse gradient.

For example, the size and shape of a molecule can be inferred from its electrophoretic mobility as a function of acrylamide and bisacrylamide concentrations (see [Ferguson Plot](#)). These electrophoretic measurements are usually carried out in a number of gels with different acrylamide concentrations, but they are more readily performed in a single transverse gradient gel in which the acrylamide concentration varies across the gel.

The most common agents used to alter the structures and conformations of proteins and nucleic acids are the denaturant [urea](#) and temperature. Their use in TGGE is described in the entries **urea gradient electrophoresis** and [temperature gradient gel electrophoresis](#), where they give information about the folding and unfolding of proteins and nucleic acids.

#### Bibliography

1. T. E. Creighton (1979) *J. Mol. Biol.* **129**, 235–264.

#### Suggestion for Further Reading

2. D. P. Goldenberg and T. E. Creighton (1984) Gel electrophoresis in studies of protein conformation and folding, *Anal. Biochem.* **138**, 1–18.

### Transversion Mutation

A transversion mutation is a [mutation](#) that substitutes a **purine** for a **pyrimidine** or a pyrimidine for a purine, inverting the purine/pyrimidine axis of the **DNA** molecule. These are

*G → C   G → T   A → C   A → T   Purine → Pyrimidine*

*C → G   C → A   T → G   T → A   Pyrimidine → Purine*

Transversions may result from the mispairing of a purine with a second purine, then replication of a nascent strand completes the transversion mutation by pairing that purine with a pyrimidine. Transversions were named by Freese ([1](#)) as a class of mutations that are not easily subject to **reversion** by base analogs (see [Transition Mutation](#)). However, transversions typically result when a nucleotide is altered by a bulky DNA adduct so that the base-pairing capacity is lost.

#### Bibliography

1. E. Freese (1959) *Proc. Natl. Acad. Sci. USA* **45**, 622–633.

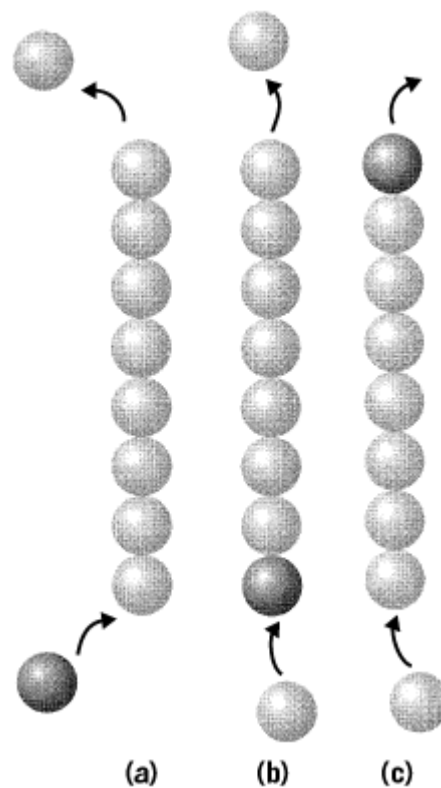
### Treadmilling



In the absence of an energy source, a [polymer](#) will exist in equilibrium with a pool of monomers. In such an equilibrium, monomers will randomly associate and dissociate from the polymer, with no net change in length. Since many protein polymers (such as [actin](#), [tubulin](#), and [RecA](#) ) bind and hydrolyze a nucleotide cofactor, the possibility exists that this energy source can be used by the polymer to drive a flux of subunits through the filament. This will result in a steady state where the polymer is maintaining a constant length but moving subunits vectorially, a process known as treadmilling.

With no energy source, a polymer may have a “fast” growing end and a “slow” growing end, but the equilibrium constant (the ratio of the rate constants for subunit addition and loss) must be the same at both ends of the filament. Thus, if the rate of addition is tenfold greater at one end of the filament than at the other, the rate of subunit loss at this fast-growing end must be tenfold greater than at the slow-growing end. The critical concentration of the monomer (the concentration at which the monomer will be in equilibrium with the polymer, with no net growth or depolymerization) must be the same at both ends of the filament at equilibrium. In the presence of an energy source, it was first realized for actin (1) that ATP hydrolysis could lead to a situation where the critical concentrations are different at the two ends. There could thus exist a state where the monomer concentration leads to depolymerization at one end, but polymerization at the opposite end. Under these steady-state conditions, the filament length could remain constant, but there would be a net flux of subunits travelling through the polymer (Fig. 1).

**Figure 1.** In a polymer where subunits add preferentially at one end (bottom, a) and depolymerize from the opposite end (top, a), a treadmilling phenomenon will exist, with a flux of subunits passing through a polymer that maintains a constant length. In (b) a “labeled” subunit adds at the bottom, and this subunit would travel the length of the polymer before exiting, as in (c). This flux of subunits would require an energy source (such as ATP hydrolysis).



A direct *in vitro* demonstration of the existence of treadmilling was shown for [microtubules](#) (2), but the role of treadmilling *in vivo* for both actin and microtubules has remained controversial for many

years. Recent observations, however, have shown the existence of *in vivo* treadmilling for both actin (3) and microtubules (4, 5). The resulting vectorial movement of subunits through a polymer could be an important mechanism for motility and intracellular transport.

### Bibliography

1. A. Wegner (1976) *J. Mol. Biol.* **108**, 139–150.
2. R. L. Margolis and L. Wilson (1978) *Cell* **13**, 1–8.
3. M. F. Carlier (1998) *Curr. Opin. Cell Biol.* **10**, 45–51.
4. V. I. Rodionov and G. G. Borisy (1997) *Science* **275**, 215–218.
5. C. M. Waterman-Storer and E. D. Salmon (1997) *J. Cell Biol.* **139**, 417–434.

### Suggestions for Further Reading

6. J. M. Neuhaus, M. Wanger, T. Keiser, and A. Wegner (1983) Treadmilling of actin, *J. Muscle Res. Cell Motility* **4**, 507–527.
7. C. M. Waterman-Storer and E. D. Salmon (1997) Microtubule dynamics: Treadmilling comes around again, *Curr. Biol.* **7**, R369–R372.

## Trifluoroethanol

Trifluoroethanol (2,2,2-trifluoroethanol,  $\text{CF}_3\text{CH}_2\text{OH}$ , often abbreviated TFE) is a substituted organic alcohol often used as a co-solvent. It is a liquid at room temperature and miscible with water. The property that makes it of interest in molecular biology is its ability to stabilize noncovalent structure in peptides. The [a-helix](#) is a common **secondary structure** motif in proteins, but it is only slightly stable in most short peptides (see [Alpha-Helix Formation](#)). TFE stabilizes helical structure in some peptides (1-3). Several proteins have been split into peptide fragments, and these peptide fragments have shown a tendency to form helices in TFE (4-6), even if they are unstructured in water. This tendency is particularly strong for peptides corresponding to helical regions in the intact protein, but it has also been present even in fragments that contained no helix in the intact protein. Relatively high concentrations of TFE are needed to get the full stabilization, and the amount of a-helical structure generally reaches a maximum around 40% to 60% (v/v) TFE. TFE has also been shown in certain instances to stabilize [b-sheet](#) secondary structure (3, 7).

Because TFE weakens **hydrophobic** interactions, it is a protein **denaturant**. However, the denatured state induced by TFE is usually structured. For example, the denatured state of the protein hen egg-white [lysozyme](#) in 70% TFE consists of regions of high helical content (8), corresponding to regions that were helical in the native protein, although this denatured state does lack [tertiary structure](#) interactions. Study of proteins denatured in TFE can give an idea of the local structure-forming tendencies of regions of the polypeptide chain and may suggest sites at which the folding of the protein is initiated.

The mechanism by which TFE stabilizes structure is unclear. One suggestion is that TFE interferes with the ability of [water](#) to solvate peptide groups that make up part of the peptide backbone (2, 9). Therefore the unstructured conformation of a peptide, whose backbone peptide groups are exposed to water, is less favored in TFE solution than the same conformation in pure water. Thus, structures featuring these groups [hydrogen-bonded](#) to other groups in the peptide (as in the a-helix) are favored.

## Bibliography

1. J. E. Brown and W. A. Klee (1971) *Biochemistry* **10**, 470–476.
2. J. W. Nelson and N. R. Kallenbach (1986) *Proteins* **1**, 211–217.
3. F. D. Sönnichsen, J. E. Van Eyk, R. S. Hodges, and B. D. Sykes (1992) *Biochemistry* **31**, 8790–8798.
4. J. Kemmink and T. E. Creighton (1995) *Biochemistry* **34**, 12630–12635.
5. D. Hamada, Y. Kuroda, T. Tanaka, and Y. Goto (1995) *J. Mol. Biol.* **254**, 737–746.
6. J. J. Yang, M. Buck, M. Pitkeathly, M. Kotik, D. T. Haynie, C. M. Dobson, and S. E. Radford (1995) *J. Mol. Biol.* **252**, 483–491.
7. E. de Alba, M. Angeles-Jimenez, M. Rico, and J. L. Nieto (1996) *Folding & Design* **1**, 133–144.
8. M. Buck, H. Schwalbe, and C. M. Dobson (1995) *Biochemistry* **34**, 13219–13232.
9. A. Cammers-Goodwin, T. J. Allen, S. L. Oslick, K. F. McClure, J. H. Lee, and D. S. Kemp (1996) *J. Am. Chem. Soc.* **118**, 3082–3090.

## Suggestions for Further Reading

10. A. Janasoff and A. R. Fersht (1994) Quantitative determination of helical propensities from trifluoroethanol titration curves. *Biochemistry* **33**, 2129–2135. An attempt at a quantitative description of the effect of TFE on helix stability, using treatments similar to those given to protein denaturants like urea or guanidinium chloride.
11. N. R. Kallenbach, P. Lyu, and H. Zhou (1996) "CD spectroscopy and the helix-coil transition in peptides and polypeptides". In *Circular Dichroism and the Conformational Analysis of Biomolecules* (G. D. Fasman, ed.), Plenum Press, New York, pp. 201–259. An excellent overall review of helix formation in peptides.
12. J. W. Nelson and N. R. Kallenbach (1989) Persistence of the alpha-helix stop signal in the S-peptide in trifluoroethanol solutions. *Biochemistry* **28**, 5256–5261. A continuation of the studies on the S-peptide of ribonuclease A.
13. K. Shiraki, K. Nishikawa, and Y. Goto (1995) Trifluoroethanol-induced stabilization of the helical structure of  $\beta$ -lactoglobulin: implication for non-hierarchical protein folding. *J. Mol. Biol.* **245**, 180–194. A comparison of the TFE-induced denatured states of 20 proteins with their native structures.

## Trigger Factor

Trigger factor is an abundant **cytosolic** protein of **bacteria**, first discovered in *Escherichia coli* by Wickner and co-workers in 1987 (1). They searched in a biochemical screen for cytosolic components involved in [protein secretion](#) of the precursor of **OmpA protein**, proOmpA, and identified a protein with an apparent molecular weight of 60 kDa that had the ability to form 1–1 stoichiometric complexes with proOmpA; this stabilized the precursor and facilitated its translocation into membrane **vesicles**. The protein was termed “trigger factor” because of its ability to trigger the folding of proOmpA into a [membrane](#) assembly-competent form *in vitro*.

Recent investigations revealed new interesting features of trigger factor. That from *E. coli* has [peptidyl-prolyl-cis/trans isomerase](#) (PPIase) activity and is capable of catalyzing **protein folding** *in vitro* much more efficiently than all other PPIases tested so far (2). In addition, trigger factor binds to the large subunit of [ribosomes](#) (3) and associates with cytosolic and secretory nascent polypeptide

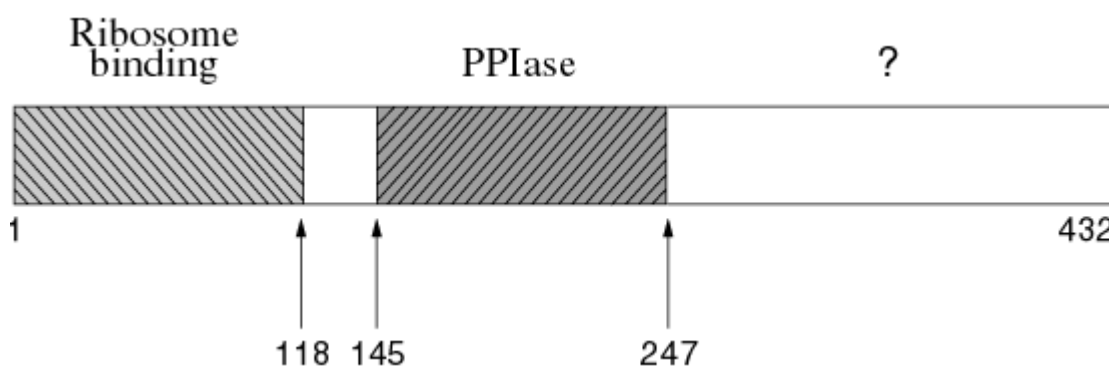
chains *in vitro* (4, 5). These findings led to the hypothesis that trigger factor acts as a cotranslational folding catalyst for newly synthesized polypeptides. Trigger factor was also described to cooperate with the [chaperonin](#) GroEL in stimulating **protein degradation** of an unstable protein, apparently by enhancing the interaction of GroEL with the substrate (6). Together, these findings suggest an involvement of trigger factor in processes related to cellular protein folding and metabolism. However, *in vivo* functions of trigger factor have not yet been demonstrated.

Genes encoding homologous trigger factor proteins were recently identified in other **eubacteria**, including *Bacillus subtilis*, *Campylobacter jejuni*, *Haemophilus influenzae*, *Haemophilus actinomycetemcomitans*, *Mycoplasma pneumoniae*, *Mycoplasma genitalium*, and *Synechocystis PCC6803*, whereas no homologous protein has been described so far in **archaeobacteria** and **eukaryotic** cells. This review will focus mainly on the *E. coli* trigger factor, since most mechanistic analyses were performed with this homologue.

### 1. A peptidyl-prolyl-*cis/trans* Isomerase

Peptidyl-prolyl-*cis/trans* isomerases (PPIases) are enzymes that catalyze the ***cis-trans* isomerization** of **peptide bonds** preceding **proline** residues (Xaa-Pro peptide bonds) (see [Peptidyl Prolyl Cis/Trans Isomerases](#)). This isomerization can be rate-limiting in the folding of a polypeptide chain into its native structure (see **Protein folding**). Trigger factor is proposed to represent a new member of the family of [FK506](#) binding (FKBP)-type PPIases on the basis of the following findings. Stoller et al. (2) identified trigger factor as a ribosome-bound PPIase. Using series of tetrapeptides of the sequence Suc-Ala-Xaa-Pro-Phe-*p*-nitroanilide to monitor PPIase activity, trigger factor resembled in its Xaa-Pro substrate specificity PPIases of the FKBP-family (2). Independently, two groups reported on a significant sequence homology of trigger factor with PPIases of the FKBP-family (5, 7). In particular, **hydrophobic** and **aromatic** residues forming the substrate-binding pocket in the well-known structure of the human FKBP12 homologue are found in trigger factor as well. A puzzling finding, however, is that the *E. coli* trigger factor is not inhibited by FK506 at concentrations of up to 100 mM (2), while the *B. subtilis* homologue is half maximally inhibited by 0.5  $\mu$ M FK506 (8). Confirmation of the FKBP-domain as a structural and functional element of trigger factor was achieved by limited **proteolysis** of the native protein, which generated a stable fragment of 12 kDa constituting the predicted FKBP domain (Fig. 1) When assayed for PPIase activity toward tetrapeptides, this fragment displayed the same specific activity as the full-length protein (9, 10).

**Figure 1.** Schematic representation of the domain structure of *E. coli* trigger factor. Hatched boxes show the domains for which activities have been assigned. Arrows with numbers indicate the residues at the domain boundaries.



Alignment of the sequences of trigger factor proteins from different prokaryotic species revealed that the central parts of the proteins representing the FKBP domain exhibit a strikingly higher degree of conservation than do their *N*- or *C*-terminal parts, suggesting that the PPIase activity is a particularly

conserved feature of the trigger factor proteins (10).

## 2. Chaperone activity

The ability of trigger factor to form 1–1 stoichiometric complexes with unfolded proOmpA and to stabilize the translocation-competent form of this precursor is reminiscent of **molecular chaperones**. A chaperone-like activity is further indicated by the trigger factor ability to associate with high affinity (apparent dissociation constant  $\sim 0.7 \mu\text{M}$ ) with the unfolded form of a mutant **ribonuclease T1** (11), and this high affinity binding is a prerequisite for trigger factor's excellent catalysis of refolding of the same protein *in vitro* (2, 11). Refolding of this protein is rate-limited by the *cis–trans* isomerization of a single proline peptide bond. Because of the high affinity of trigger factor for the unfolded substrate, it outscores all other PPIases in its efficiency of catalyzing the prolyl isomerization, although the catalytic constant is low ( $k_{\text{cat}} \sim 1.3 \mu\text{M}$ ) (11). The activity of the isolated central PPIase **domain** of trigger factor is reduced 800-fold, whereas its PPIase activity toward tetrapeptide substrates is comparable to that of full-length trigger factor (11).

These findings indicate distinct sites on trigger factor for high affinity substrate binding and prolyl-isomerization and suggest that the exceptional catalytic efficiency of trigger factor originates from a cooperation of its PPIase and chaperone-like activities (11). Tight substrate binding may position the stretch of the protein substrate that contains the critical proline residue close to the active site of the PPIase.

## 3. Association with Nascent Polypeptide Chains

Two research groups independently demonstrated that *E. coli* trigger factor can associate with nascent polypeptide chains. Lührink and coworkers discovered that trigger factor can be efficiently **crosslinked** to all nascent chains arrested on the ribosome (4). Trigger factor competes with P48, a component of the *E. coli* **signal recognition particle** (SRP), for crosslinking to nascent chains of precursors of secretory proteins. In contrast to P48, however, trigger factor also crosslinks efficiently to nascent chains of cytoplasmic polypeptides (4). Bukau and co-workers identified trigger factor as the main component that associates with ribosomes translating **b-galactosidase** (5). This association was resistant to high salt but was disrupted by **puromycin** treatment, which leads to premature termination of translation, suggesting an association of trigger factor with b-galactosidase nascent chains (5). Further evidence for an association of trigger factor with nascent chains was provided by its photocrosslinking to nascent chains of secretory as well as nonsecretory derivatives of preprolactin (5).

The ability of trigger factor to associate with nascent chains emerging from the ribosome suggests that this protein has a binding site on the ribosome that positions it near the exit site for nascent chains. Trigger factor is indeed exclusively associated with the large ribosomal subunit, which harbors this site (3), and trigger factor cannot be **crosslinked** to nascent polypeptides after their puromycin-mediated release from the ribosome (4, 5). Interestingly, trigger factor also associates with eukaryotic **wheat germ** ribosomes, suggesting conservation in evolution of its binding site on the ribosome (5).

## 4. Role with GroEL in Binding and Degradation of Proteins

A further role for trigger factor in protein metabolism is indicated by the finding that trigger factor stimulates the degradation of CRAG, an artificial protein composed of fragments of the **l cro repressor**, b-galactosidase, and **protein A** that has been used as a model unfolded substrate for studies of **protein degradation** in *E. coli*. Proteolysis of CRAG by the ClpP **proteinase** depends on the interaction of CRAG with the **chaperonin** GroEL, and formation of this complex was the **rate-limiting step** in degradation. Trigger factor accelerates the chaperone-dependent degradation of CRAG by promoting its binding to GroEL. A ternary complex of CRAG, GroEL, and trigger factor

was demonstrated by [affinity chromatography](#). Addition of ATP causes dissociation of the complex of GroEL and trigger factor from the substrate (6). Trigger factor interacts directly with GroEL prior to association with CRAG. This interaction differs from that of GroEL with denatured substrates in that addition of ATP and GroES, which causes substrate release from GroEL, does not dissociate the complex (12). Trigger factor also enhances the binding of GroEL to unfolded proteins other than CRAG (12). This indicates that trigger factor is involved in both degradation and folding of denatured substrates by stimulating their binding to the chaperonin (12).

## 5. Domain Organization

The *E. coli* trigger factor consists of 432 amino acids with a molecular weight of 48 kDa. Limited proteolysis of the full-length protein with [proteinases](#) revealed a compactly folded proteinase-resistant central domain comprising residues 145–247 (Fig. 1) (9, 10). This domain has homology to FKBP and displays PPIase activity, and thus represents the catalytic core of the protein.

The *N*-terminal part of trigger factor forms another structural and functional module. The *N*-terminal 144 residues are necessary and sufficient for specific binding of trigger factor to the large ribosomal subunit, both *in vitro* and *in vivo* (13). This fragment contains a compactly folded domain comprising the amino-terminal 118 amino acids (Fig. 1). It copurifies with *E. coli* ribosomes from cell extracts and is capable of associating with isolated ribosomes *in vitro*; therefore, it represents the ribosome-binding domain of trigger factor (13).

Intriguingly, the *N*-terminal fragment alone, or with the adjacent PPIase domain, is unable to form complexes with translating ribosomes that are resistant to high salt concentrations. Furthermore, the association of both fragments with ribosomes was not altered by puromycin, suggesting that the sensing of the translational status requires the C-terminus of trigger factor (13). It is conceivable that the C-terminus mediates the tight association with nascent polypeptide chains or with the translating ribosomes, which consequently leads to the formation of high salt-resistant complexes of trigger factor with the translating ribosome.

## 6. Cellular Function

Despite the considerable progress made in the biochemical dissection of trigger factor *in vitro*, the role of this protein in *E. coli* remains unclear. The known properties of trigger factor strongly suggest an important biological function in folding of newly synthesized proteins and raises the attractive hypothesis that trigger factor acts as a cotranslational folding catalyst. Interestingly, expression of the *tig* gene encoding trigger factor is growth-phase-controlled and coregulated with genes encoding ribosomal components (14). A possible role of trigger factor in GroEL-assisted protein metabolism (12) is not yet established.

Although a **knockout** mutation of the *tig* gene is not available yet, the available data indicate that trigger factor is not essential for *E. coli* growth at 37°C. Depletion of trigger factor to less than 5% of the cellular wild-type levels causes formation of short filaments, indicating cell division defects. Trigger factor-depleted cells remain fully viable at 37°C, however, and have no additional cellular defects (14). For example, although trigger factor enhances translocation of proOmpA into vesicles *in vitro* (1), trigger factor-depleted cells have no defect in translocation of the precursor, even in a *secB* mutant background that lacks activity of the secretion-specific chaperone SecB (14). It is unclear whether the activities of other chaperones or PPIases existing in the *E. coli* cytosol can compensate at 37°C for the missing activity of trigger factor in the depleted cells.

The only detected phenotype of trigger factor-depleted cells was observed at low temperature. Depleted cells stored at 4°C die faster than do normal cells, while cells overproducing trigger factor show enhanced viability at this temperature (15). This suggests an involvement of trigger factor in cell survival at temperatures below the growth temperature limit. Trigger factor may be particularly important at low temperatures for folding of proteins requiring prolyl-bond isomerization. This

isomerization is slower at low temperature, so there is an increased requirement for catalysis. Furthermore, the cellular level of trigger factor is about twofold higher at 16°C and 4°C than at 37°C (15).

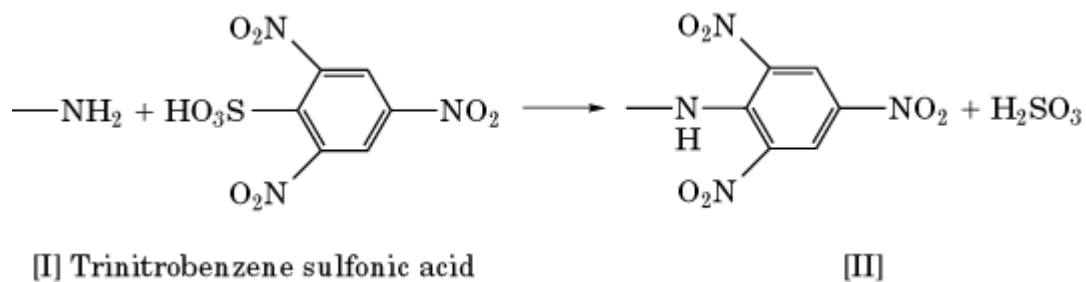
The importance of trigger factor for metabolism of bacteria is also indicated by the discovery of a *tig* gene in *Mycoplasma genitalium*. This bacterium is believed to be free from genetic redundancy and thus contains only the minimal set of genes required for life. Trigger factor appears to be the only PPIase of this organism (16). Identification of the *in vivo* functions of trigger factor thus seems to be of importance for understanding the fundamental process of protein folding in the bacterial cytosol.

## Bibliography

1. E. Crooke and W. Wickner, (1987) Proc. Natl. Acad. Sci. USA **84**, 5216–5220.
2. G. Stoller, K. P. Rücknagel, K. H. Nierhaus, F. X. Schmid, G. Fischer, and J.-U. Rahfeld, (1995) EMBO J. **14**, 4939–4948.
3. R. Lill, E. Crooke, B. Guthrie, and W. Wickner, (1988) Cell **54**, 1013–1018.
4. Q. A. Valent, D. A. Kendall, S. High, R. Kusters, B. Oudega, and J. Luirink, (1995) EMBO J. **14**, 5494–5505.
5. T. Hestekamp, S. Hauser, H. Lütcke, and B. Bukau, (1996) Proc. Natl. Acad. Sci. USA **93**, 4437–4441.
6. O. Kandror, M. Sherman, M. Rhode, and A. L. Goldberg, (1995) EMBO J. **14**, 6021–6027.
7. I. Callebaut, and J.-P. Mornon, (1995) FEBS Lett. **374**, 211–215.
8. S. F. Göthel, R. Schmid, A. Wipat, N. M. Carter, P. T. Emmerson, C. R. Harwood, and M. A. Marahiel, (1997) Eur. J. Biochem. **244**, 59–65.
9. G. Stoller, T. Tradler, K. P. Rücknagel, J.-U. Rahfeld, and G. Fischer, (1996) FEBS Lett. **384**, 117–122.
10. T. Hestekamp and B. Bukau, (1996) FEBS Lett. **385**, 67–71.
11. C. Scholz, G. Stoller, T. Zarnt, G. Fischer, and F. X. Schmid, (1997) EMBO J. **16**, 54–58.
12. O. Kandror, M. Sherman, R. Moerschell, and A. L. Goldberg, (1997) J. Biol. Chem. **272**, 1730–1734.
13. T. Hestekamp, E. Deuerling, and B. Bukau, (1997) J. Biol. Chem. **272**, 21865–21871.
14. B. Guthrie and W. Wickner, (1990) J. Bacteriol. **172**(10), 5555–5562.
15. O. Kandror and A. L. Goldberg, (1997) Proc. Natl. Acad. Sci. **94**, 4978–4981.
16. C. M. Fraser, J. D. Gocayne, O. White, M. D. Adams, R. A. Clayton, R. D. Fleischmann, J., B. C., Kerlavage, A. R., A. R. Sutton, and J. M. E. A. Kelly, (1995) Science **270**, 397–403.

## Trinitrobenzene Sulfonic Acid

2,4,6-Trinitrobenzene sulfonic acid (TNBS) [I], usually obtained as the sodium salt has a molecular weight of 315.15. TNBS reacts specifically with [amino groups](#) of [amino acids](#), [peptides](#), and [proteins](#) under mildly alkaline conditions (Scheme 1) and yields trinitrophenyl derivatives [II] with a yellow color (340 to 350 nm). This reaction is employed to detect or measure of amino groups (1, 2). TNBS also reacts with [thiol groups](#).



## Bibliography

1. R. Haynes, D.T. Osuga and R.E. Feeney (1967) *Biochemistry* **6**, 541–547.
2. R. Fields (1972) *Methods Enzymol.* **25**, 464–468.

## Trinucleotide Repeats

The application of molecular **genetic** approaches to identifying responsible **genes** has advanced our understanding of the molecular mechanisms of a number of hereditary neurological diseases. Among these diseases, a distinct group constitutes a new type of disease mechanism, caused by an unstable expansion of a trinucleotide repeat in the [genome](#), such as a CAG trinucleotide repeat. The expansion of trinucleotide repeats in eleven human gene loci causes the development of at least eight hereditary diseases that affect the nervous system: (1) fragile [X-chromosome](#) syndrome; (2) myotonic dystrophy; (3) spinal and bulbar Kennedy's amyotrophy; (4) Huntington's chorea; (5) type 1 spinocerebellar ataxia; (6) Friedreich's ataxia; (7) dentatorubral-pallidoluysian atrophy; and (8) Machado–Joseph disease.

The triplet repeat may be located in the 5' untranslated region of the gene (for example, an expanded CGG repeat in the case of the fragile-X syndrome) or in the 3' untranslated region (for example, an expanded CTG repeat in the case of myotonic dystrophy). An expanded GAA repeat within the first intron of the gene has been identified in Friedreich's ataxia. Within protein coding regions, an expansion of CAG repeats (and therefore of a corresponding protein with long stretches of glutamine residues) has been identified in spinal and bulbar muscular atrophy, Huntington's disease, spinocerebellar ataxia type 1, dentatorubral-pallidoluysian atrophy, and Machado–Joseph disease. Similar unstable trinucleotide repeat expansions have been suggested, but not unequivocally proven, in many neurological diseases, including dominant ataxia, schizophrenia, and bipolar illness.

Triplet repeats are present in normal individuals. The number of repeat units ranges up to nearly 35. Expanded alleles with the number of trinucleotide repeats exceeding the normal range lead to the expression of disease phenotypes. All the repeats seem to be intrinsically unstable. Unaffected individuals have varying size repeats within the normal range, and most have different numbers of repeats in their maternal and paternal alleles. Importantly, none of the affected genes shows a continuous gradient of repeat sizes from normal to abnormal. Instead, normal genes keep their slightly unstable repeats within the normal range, whereas the repeats in a small subset of abnormal genes jump into the abnormal range. Once abnormal, the repeats become highly unstable and lead to disease. Indeed, all the disorders seem to show [linkage disequilibrium](#). Inheritance of the affected gene is linked to the inheritance of other genetic markers. This indicates that the lineage of most patients with each disease can be traced back to one or a few shared ancestors who acquired the archetype expansion mutation. Thus, it is an excessively rare event for any of the normal frequency



repeats to jump into the pathological range. The increase in the number of repeats appears to account for the phenomenon called *anticipation* in which the clinical severity of the mutation increases from generation to generation.

## 1. Some Examples

In the case of the fragile-X syndrome, the length polymorphism involving the CGG repeat at the 5' end of the FRX gene can give rise to an increased length variant (premutation) that has no clinical phenotype, but is genetically unstable and can generate longer CGG repeats associated with the fragile-X phenotype. The normal length of the repeat expansion  $(CGG)_n$  varies from  $n = 6$  to 52, whereas it increases to  $n = 52$  to 200 in diseased individuals. **PCR** amplification with primers flanking the CGG repeat distinguish the normal length range of polymorphic variants from the longer premutation variants. Some of the unstable variants, however, are sufficiently long to make PCR amplification of the CGG repeat inefficient. Then **Southern blotting** may be required.

In myotonic dystrophy, the CTG repeat, whose length is normally between 5 and 30 repeats, can reach 50 to 2000 repeats in the disease. In Huntington's disease, the respective figures for  $(CAG)_n$  are  $16 < n < 36$  and  $42 < n < 86$ . In spinobulbar muscular atrophy, the respective figures for  $(CAG)_n$  are  $17 < n < 26$  and  $40 < n < 52$ . In Friedreich's ataxia, expansion of the GAA trinucleotide repeat ranges from 200 to 1200 repeat units with no overlap in the normal range.

Although genetic rather than [epigenetic](#) events occur at these loci to produce affected individuals, the characteristics of these diseases often fulfill the definition of genetic [imprinting](#) because the disease locus is affected differently upon passage through female versus male gametogenesis. For example, in congenital myotonic dystrophy and fragile-X syndrome, expansion of the triplet repeat leading to the disease phenotype occurs only through female transmission, whereas in Huntington's disease the most dramatic repeat expansion occurs in the father's germ line.

## 2. Biochemical Data Linked to the Triplet Repeat Expansion

### 2.1. Chromosomal Fragile Sites

Unusual DNA secondary structures have been implicated in the expansion of trinucleotide repeats. Evidence has been presented consistent with folding of the DNA structure of the repeats into hairpin loops at the center of a long **palindromic** segment *in vivo* (1). In two types of fragile-X chromosome syndrome, FRAXA and FRAXE, the trinucleotide repeat expansion may alter the long-range [chromatin](#) structure, which could influence [transcription](#) of nearby gene sequences. FRAXF, another X-chromosomal fragile site that has been cloned, harbors a polymorphic compound triplet array,  $(GCCGTC)_n(GCC)_n$ . Expansion and **methylation** of the GCC repeat and the neighboring CpG-rich region result in chromosomal fragility. It has been suggested that hairpin formation may occur, accompanied by **DNA polymerase** slippage.

High-resolution [NMR](#) spectroscopy and other biophysical methods have given a global picture of the solution behavior of the disease-related CXG ( $X = A, C, G$  or  $T$ ), which show a propensity for folding at lengths as short as 12 residues (2). The fragile FRAXA sites stain poorly on the chromosome, which suggests an altered chromatin structure. Repeating CCG DNA from FraX patients were tested for their ability to assemble into [nucleosomes](#), the basic units of chromatin, using *in vitro* nucleosome reconstitution, [electron microscopy](#), and competitive assembly [gel retardation assays](#). CCG blocks of 50 repeats displayed strong nucleosome exclusion, providing a possible explanation for the nature of the fragile sites (3), although this conclusion is contradicted by others (4).

The CGG/CCG repeats, which unlike the CAG repeats contain a methylatable CpG dinucleotide, acquire considerable methylation when the repeat length exceeds about 200 trinucleotides.

Expansion appears to precede methylation, but how the DNA methylation is acquired remains unclear.

### 2.1.1. CAG and CTG Repeats

The consequences of the CAG repeat expansion seem to be at the protein level, conferring a gain or an altered function rather than a loss of function. The gene responsible for spinal and bulbar Kennedy's amyotrophy has been localized to Xq12 and the androgen receptor gene. This gene contains a trinucleotide repeat within the first **exon** that encodes a polyglutamine tract. The effect on the protein of lengthening this tract is currently unknown. Because the androgen receptor is a [transcription factor](#), one hypothesis is that polyglutamine expansion alters the **gene expression**, leading to neuron degeneration.

The gene responsible for type 1 spinocerebellar ataxia contains an open reading frame that predicts an 87-kDa polypeptide of unknown function (*ataxin1*) containing a polymorphic polyglutamine tract interrupted by two histidine residues. By immunohistochemical methods, both wild-type and mutant ataxin1 have been found in cultured cells, predominantly in the nuclei. Mice containing a human transgene with 82 repeats experience the disease, confirming the dominant nature of the expanded allele.

In the case of Huntington's chorea, the *huntingtin* protein that is enlarged by the expanded polyglutamine tract is predominantly cytoplasmic and of unknown function.

The gene responsible for dentatorubral–pallidolusian atrophy normally encodes a 4.5-kb transcript that is expressed in all tissues. The predicted normal polypeptide chain is 1184 residues long and contains the usual polyglutamine tract corresponding to the CAG repeat that is enlarged in the mutant protein.

In myotonic dystrophy, the aberrantly expanded CTG repeat resides in the 3' untranslated region of a **serine–threonine protein kinase**.

A large scale screening project for CAG and CTG repeats in human reference [cDNA](#) has been undertaken. Nine new cDNA containing polymorphic CAG/CTG repeats have been identified and assigned to chromosomes. Three of them are highly polymorphic and represent the most likely candidate genes for inherited neurodegenerative diseases and, perhaps, neuropsychiatric disorders of multifactorial origin (5).

### Bibliography

1. J. M. Darlow and D. R. Leach (1995) *Genetics* **141**, 825–832.
2. M. Zheng et al. (1996) *J. Mol. Biol.* **263**, 511–516.
3. Y. H. Wang et al. (1996) *J. Mol. Biol.* **263**, 511–516.
4. J. S. Godde and A. P. Wolffe (1996) *J. Biol. Chem.* **271**, 15222–1522.
5. C. Neri et al. (1996) *Hum. Mol. Genet.* **5**, 1001–1009.

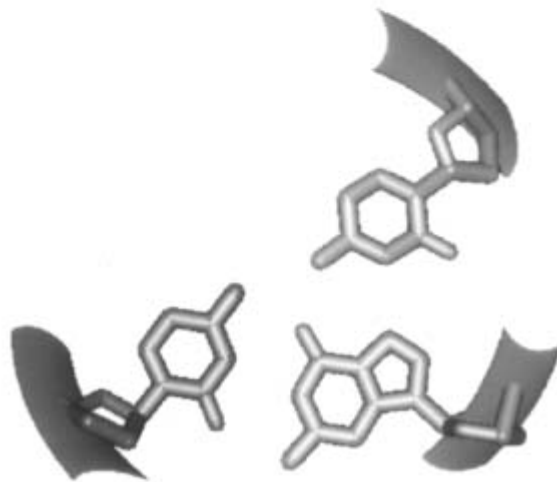
### Suggestion for Further Reading

6. C. T. Ashley, Jr. and S. T. Warren (1995) Trinucleotide repeat expansion and human disease, *Ann. Rev. Genetics* **29**, 703–728.

## Triple Helix, Nucleic Acids

The formation of a triple-stranded nucleic acid structure, specifically, a triplex helix, was demonstrated nearly 40 years ago. It was found that titration of poly(A) with poly(U) resulted in two distinct types of stable structures: that of the 1:1 duplex and the 1:2 triplex. In the latter structure, the third strand of poly(U) is bound to the poly(A) strand of the poly(A).poly(U) duplex using Hoogsteen base pairs. The recent resurgence of interest in the triple helix is due primarily to its possible role in gene [transcription](#) and its possible application in the design of [antisense oligonucleotides](#). It has been shown that stretches of  $(GA)_n:(TC)_n$  sequence in a **plasmid**, which were found to be hypersensitive to **nucleases** during active [transcription](#), underwent a notable change in [DNA structure](#) when the plasmid was negatively **supercoiled** under low-pH conditions. The structural change has been interpreted to be the formation of a triple helix for part of the  $(GA)_n:(TC)_n$  sequence (named H-DNA), rendering the extrusion of a single-stranded loop. The three-dimensional structures of a number of triple-stranded DNA molecules have been studied by [NMR](#). One such example is shown in [Figure 1](#). As noted above, the formation the T.A.T triple base pair requires a Hoogsteen base pair. A similar triple C.G.C base pair is not stable unless the second C is protonated (thus forming a C.G.C<sup>+</sup> triple base pair). This explains the need for a low pH for H-DNA. Thus far, it has been shown that the formation of a stable triplex helix of the (Py).(Pu).(Py) type is relatively straightforward. It remains a challenge, however, to design oligonucleotides that can bind to a specific DNA duplex using all four bases.

**Figure 1.** Side and end views of the structure of a DNA triple helix as determined by NMR (Protein Database 1D3X).



#### Suggestions for Further Reading

W. Saenger (1984) *Principles of Nucleic Acid Structure*, Springer-Verlag, New York.

C. R. Calladine and H. R. Drew (1992) *Understanding DNA*, Academic Press, San Diego.

R. R. Sinden (1994) *DNA Structure and Function*, Academic Press, San Diego.

#### Triple Resonance

A nuclear magnetic resonance ([NMR](#)) experiment involves application of pulses of RF energy to the sample under examination. At least one of these pulses must have a frequency that is close to the Larmor frequency of the spins of the sample that are to be observed (see [NMR](#)). Thus, application of a pulse at 750 MHz is sufficient to produce proton NMR signals from a sample that resides in a magnetic field of 17.62 T. When all of the RF pulses used for an experiment are at the same frequency, the experiment is referred to as a single-resonance NMR experiment. More complex experiments may use RF pulses at two different radiofrequencies; these are double resonance experiments. NMR experiments used to provide information for the determination of the tertiary structures of proteins or nucleic acids, commonly use materials that have been isotopically enriched in  $^{13}\text{C}$  and  $^{15}\text{N}$ . In these experiments, RF pulses near the  $^{13}\text{C}$ ,  $^{15}\text{N}$ , and proton Larmor frequencies are applied to the sample to effect various coherence transfers or to produce spin decoupling; such experiments are triple resonance experiments. Multiple resonance experiments of almost any type are technically feasible, although the expense of these increases substantially as the number of different RFs involved increases.

Multiple resonance NMR experiments are often indicated by means of the notation  $M\{X, Y, \frac{1}{4}\}$ , where M is the chemical symbol for the nucleus whose NMR signals are detected in the experiment, and X, Y,  $\frac{1}{4}$  are the chemical symbols for other spins of the sample that are affected by RF pulses applied at other radiofrequencies. Thus, the triple resonance experiments used to determine biological structures are indicated by  $^1\text{H}\{^{13}\text{C}, ^{15}\text{N}\}$ . (See also [NMR](#) and **Isotope filtering**.)

#### Suggestions for Further Reading

A. Bax and S. Grzesiek (1993) In *NMR of Proteins* (G. M. Clore and A. M. Gronenborn, eds.), C.R.C. Press, Boca Raton, pp. 33–52.

J. Cavanagh, W. J. Fairbrother, A. G. Palmer III and N. J. Skelton (1995) *Protein NMR Spectroscopy*, Academic, San Diego.

B. Whitehead, C. J. Craven, and J. P. Waltho (1997) In *Protein NMR Techniques* (D. G. Reid, ed.), Humana, Totowa, New Jersey, pp. 29–52.

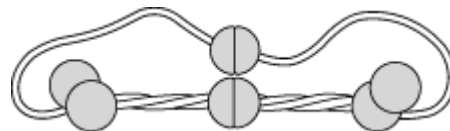
### Triple-Helical Proteins

Reports of proteins with structures containing three interwound  **$\alpha$ -helices** are becoming increasingly common. One of the first to be characterized was [collagen](#), a major component of connective tissue. Type I collagen contains two identical and one different chain, whereas other collagen types may contain either three unique chains or even three identical chains. Hence, both heterotrimers and homotrimers occur naturally *in vivo*.

[Fibrinogen](#), a blood plasma glycoprotein with a molecular weight of about 340 kDa in humans, comprises a pair of heterotrimers, each containing an a, b, and g chain. The *N*-terminal **domains** of all six chains form a central, **disulfide-bonded** globular region, from which two identical three-stranded [coiled-coil](#) rod domains emanate. Each is about 16 nm in length and is stabilized at both ends by a cyclic pattern of interchain disulfide bonds. The rods terminate in globular regions, thus generating a trinodular structure for fibrinogen (Fig 1). These outer nodules consist of the *C*-terminal regions of both the b and g chains. The *C*-terminal region of the a chain, however, loops back and

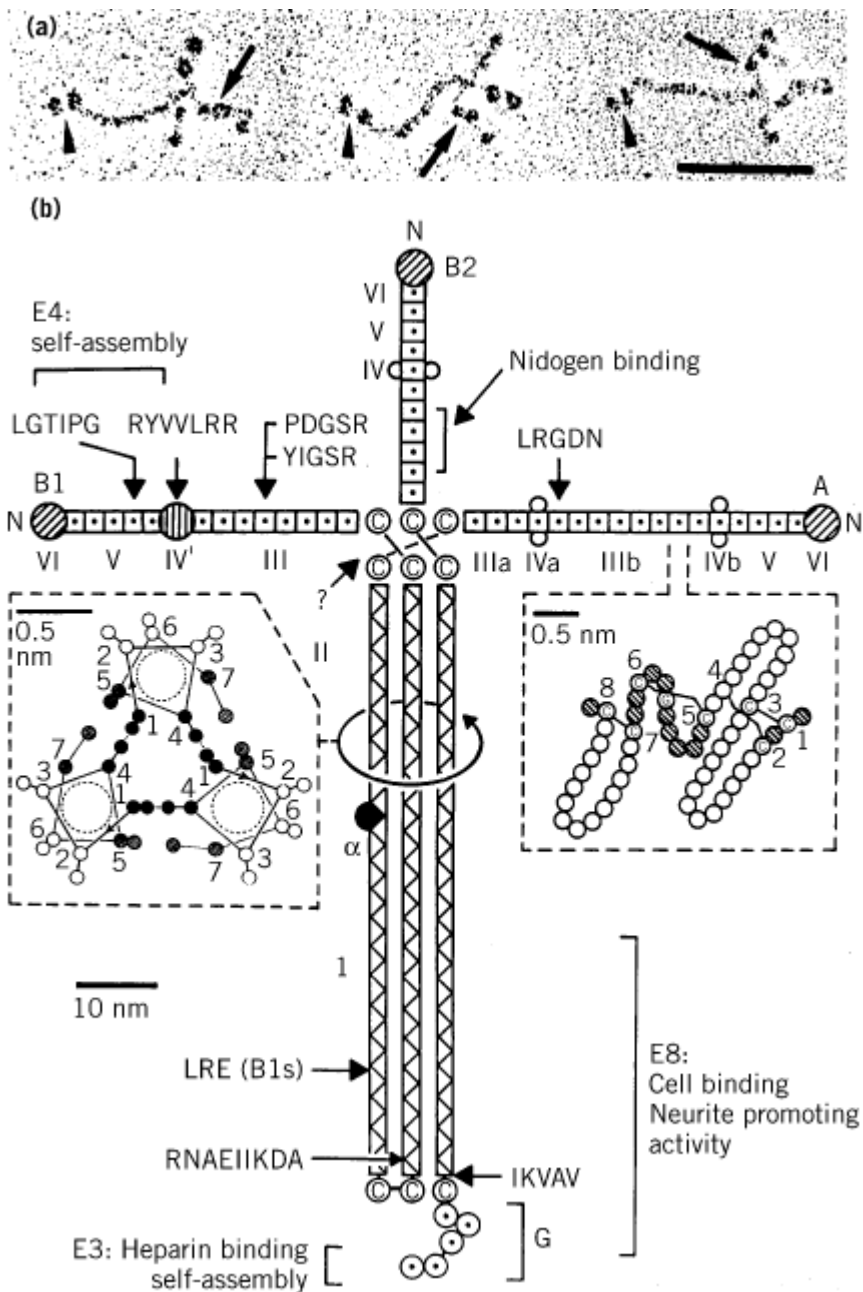
interacts with the central *N*-terminal disulfide-bonded knot. The D fragment from fibrinogen (86 kDa) has been crystallized (1), and its structure is currently being determined at atomic resolution.

**Figure 1.** The trinodular structure of fibrinogen, which consists of a pair of abg heterotrimers. The *N*-terminal regions of the a-, b-, and g-chains are cross-linked by disulfide bonds and together constitute the central nodule. Triple-helical coiled-coil rods extend from the central disulfide-bonded knot and terminate in the two outer nodules, which comprise the *C*-terminal ends of the b- and g-chains. The *C*-terminal end of the a-chain, however, folds back and interacts with the central nodule. (Courtesy of C. Cohen and D. A. D. Parry.)



**Laminins** have many functions in cellular processes and in the organization of the extracellular matrix. Some of the cellular functions include the promotion of growth, differentiation, neurite outgrowth, guidance, mobility, attachment, and mediation of cell communication (2). Laminins are a family of large multidomain **glycoproteins** with molecular weights about 800 to 900 kDa. While there are variations in the structures of different laminin molecules, each consists of three different chains that are disulfide-bonded to one another. They also form a characteristic, asymmetric cross-like structure (Fig 2). The chains aggregate via a long three-stranded coiled-coil domain of somewhat irregular heptad substructure that is located at or relatively close to the *C*-terminal ends of the chains. The three remaining ends of the cross are terminated by small globular domains. Each arm in turn comprises a series of tandem repeats or modules akin to epidermal growth factor motifs.

**Figure 2.** Structures of laminins. (a) Electron micrographs of rotary-shadowed sea urchin laminin showing an asymmetric cross shape. The long arm is 110 nm in length. The arrows indicate one of the short arms with three globular domains (A-chain), and the arrowheads show a short region of extended structure separating globular domains at the *N*-terminal end of the long arm. Bar = 100nm. (b) Model of mouse laminin molecules (A–B1–B2 heterotrimer). Different regions of the molecules are designated by Roman numerals. Squares and circles with a central dot indicate the *eight-cysteine* and the *sex hormone-binding globulin* motifs, respectively. The three-stranded coiled-coil is stabilized at either end by disulfide bonds, indicated by the letter C in a circle. Ref. 11, with permission.)



There are many other triple-helical proteins.

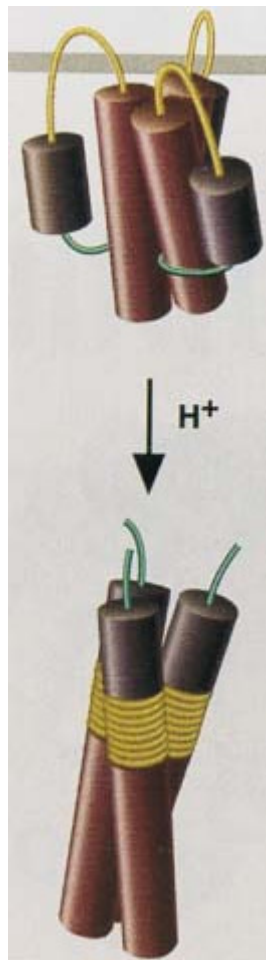
1. These include macrophage scavenger receptor protein (3), a type II **transmembrane protein** in which the sequence contains motifs characteristic of both collagen (a triplet substructure based on glycine) and  $\alpha$ -fibrous proteins (a **heptad repeat** with the **nonpolar** residues alternately three and four apart). The conformation of the extracellular domain in this homotrimeric molecule is believed to contain both a region of collagen-like triple helix and a three-stranded coiled-coil.

2. Mannose-binding protein, a member of the *collectin* family that consists of 18 chains arranged as a hexamer of homotrimers, facilitates the recognition of pathogens by **phagocytes** and plays a key role in first line host defense (4). Each chain has an *N*-terminal region rich in **cysteine** residues, 19 collagen repeats, a heptad-containing substructure, and a *C*-terminal region that is the carbohydrate-recognition domain. Mannose-binding protein, with its three- and seven-residue motifs, has a strong structural similarity with macrophage scavenger receptor protein and lung surfactant D.

3. Another three-stranded structure is found in the gp17 tail fiber of **bacteriophage T7** ([5](#)). It contains a heptad repeat and forms a coiled-coil about 16 nm in length. The 150 residues *N*-terminal to the rod link the fiber to the tail tube of the bacteriophage. The region *C*-terminal to the rod forms four globular domains in a rigid linear array. These form the “foot” of the tail fiber.
4. Evidence has also been presented to show that cartilage matrix protein (CMP) assembles through its *C*-terminal 36 residues into a three-stranded coiled-coil.
5. Tenascin, another major extracellular matrix glycoprotein, has some similarity with fibronectin. Each contains many **epidermal growth factor**-like repeats. Three tenascin chains aggregate via a short region of three to four heptad repeats to initiate the molecular structure.
6. The [spectrin](#) superfamily of proteins also form a triple helical structure, but in these cases a single chain folds to generate a three- $\alpha$ -helix motif.
7. The coiled-coil structure of the fibrous stem of the haemagglutinin in [influenza virus](#) displays a dynamic feature that is quite extraordinary. At neutral pH the stem forms a three-stranded coiled-coil about 8 nm in length. At pH 5, however, the coiled-coil is extended at its *N*-terminal end to generate a much longer coiled-coil (Fig [3](#)). This is formed from a region that was a loop at neutral pH and from another (shorter) piece of coiled-coil that was oppositely orientated at the higher pH ([6](#), [7](#)).

**Figure 3.** The pH-dependent conformational change undergone by influenza virus hemagglutinin. At neutral pH (*top*) the fibrous stem of hemagglutinin consists (mainly) of a three-stranded coiled-coil of approximate length 8 nm, a loop, and a short stretch of  $\alpha$ -helix oppositely oriented to the coiled-coil. At lower pH (*bottom*), however, the loop forms an extension to the coiled-coil at its *N*-terminal end. In turn, the  $\alpha$ -helix becomes oriented in the same direction as the coiled-coil and the loop (now a coiled-coil) and adds even further length to it. The massive spatial change of the  $\alpha$ -helical region promotes membrane fusion and cell entry *in vivo*. (From Ref. [12](#), with permission.) See color insert.





8. Finally, it should be noted that three-stranded protein structures are not restricted to collagen and  $\alpha$ -fibrous proteins. A structure recognized only in recent times is one based on the **b-helix**, as in the crystal structure of the phage P22 tail spike (8). A similar conformation is proposed to occur in the long tail fibers of bacteriophage T4 (9) and in the [adenovirus](#) fiber (10).

#### Bibliography

1. S. J. Everse, H. Pelletier, and R. F. Doolittle, Crystallization of fragment D from human fibrinogen. (1995) *Protein Sci.* **4**, 1013–1016.
2. J. Engel, I. Hunter, T. Schulthess, K. Beck, T. W. Dixon, and D. A. D. Parry (1991) Assembly of laminin isoforms by triple- and double-stranded coiled-coil structures. *Biochem. Soc. Trans.* **19**, 839–843.
3. T. Kodama, M. Freeman, L. Rohrer, J. Zabrecky, P. Matsudaira, and M. Krieger (1990) Macrophage scavenger receptor-type I: a trimer of  $\alpha$ -helical and collagen-like coiled-coils. *Nature (London)* **343**, 531–535.
4. S. Sheriff, C. Y. Chang, and R. A. B. Ezekowitz (1994) Human mannose-binding protein carbohydrate recognition domain trimerizes through a triple-helical  $\alpha$ -helical coiled-coil. *Nature Struct. Biol.* **1**, 789–793.
5. A. C. Steven, B. L. Trus, J. V. Maizel, M. Unser, D. A. D. Parry, J. S. Wall, J. F. Hainfeld and F. W. Studier (1998) Molecular substructure of a viral receptor-recognition protein, The gp17 tail-fiber of bacteriophage T7. *J. Mol. Biol.* **200**, 351–365.
6. C. M. Carr and P. S. Kim (1993) A spring-loaded mechanism for the conformational change of influenza hemagglutinin. *Cell* **73**, 823–832.
7. P. A. Bullough, F. M. Hughson, J. J. Skehel, and D. C. Wiley (1994) Structure of influenza haemagglutinin at the pH of membrane fusion. (1994) *Nature (London)* **371**, 37–43.

8. S. Steinbacher, R. Seckler, S. Miller, B. Steipe, R. Huber, and P. Reinemer (1994) Crystal structure of P22 tailspike protein: interdigitated subunits in a thermostable trimer. *Science* **265**, 383–386.
9. M. E. Cerritelli, J. S. Wall, M. N. Simon, J. F. Conway and A. C. Steven (1996) Stoichiometry and domainal organization of the long tail-fiber of bacteriophage T4: a hinged viral adhesin. *J. Mol. Biol.* **260**, 767–780.
10. P. F. W. Stouten, C. Sander, R. W. H. Ruigrok and C. Cusack (1992) New triple-helical model for the shaft of the adenovirus. *J. Mol. Biol.* **226**, 1073–1084.
11. K. Beck and T. Gruber (1995) "Structure and Assembly of Basement Membrane and Related Extracellular Matrix Proteins". In *Principles of Cell Adhesion* (P. D. Richardson and M. Steiner, eds.), CRC Press, Boca Raton, FL, pp. 219–252.
12. C. Cohen and D. A. D. Parry (1994)  $\alpha$ -Helical coiled coils: More facts and better predictions. *Science* **263**, 488–489.

### Suggestions for Further Reading

13. K. Beck and T. Gruber (1995) "Structure and Assembly of Basement Membrane and Related Extracellular Matrix Proteins". In *Principles of Cell Adhesion* (P. D. Richardson and M. Steiner, eds.), CRC Press, Boca Raton, FL, pp. 219–252.
14. C. Cohen and D. A. D. Parry (1990)  $\alpha$ -Helical coiled-coils and bundles: how to design an  $\alpha$ -helical bundle. *Proteins Struct. Funct. Genet.* **7**, 1–15.
15. A. Lupas (1996) Coiled-coils: new structures and new functions. *Trends Biochem. Sci.* **21**, 375–382.
16. R. A. Kammerer (1997)  $\alpha$ -Helical coiled-coil oligomerization domains in extracellular proteins. *Matrix Biol.* **15**, 555–565.

## Triploidy

In this type of **polyploidy**, somatic **nuclei** possess three copies of the basic [haploid](#) set of [chromosomes](#). As a general rule, triploids undergo irregular **meiosis**, so the triploidy is not normally passed on to succeeding generations.

## Tris Buffer

The primary amine known as Tris,  $(\text{HO}-\text{CH}_2)_3\text{C}-\text{NH}_2$ , is frequently used in buffer solutions (see [Buffers](#)). It has a  $\text{p}K_a$  of 8.1 at 25 °C. The name *Tris* is derived from old forms of the name, such as tris(hydroxymethyl)methylamine; this is unambiguous, even if not quite in accord with modern systematic practice.

The hydroxy groups are electron-withdrawing, so that the lone pair on the nitrogen atom is much less

available than in a typical primary amine. This lowers the  $pK_a$  of Tris to 8.1, so that it is suitable for buffering near neutral pH. This, however, raises the first disadvantage of Tris, that the  $pK$  is not quite low enough for much biological work. If a Tris buffer has a pH of 7.5 or below, less than a fifth of the Tris will be in the  $R-NH_2$  form with an amino group, and over 4/5 as  $R-NH_3^+$  with an ammonio group (since 0.6 pH units below the  $pK_a$  implies a  $[R-NH_3^+]/[R-NH_2]$  ratio of 4). Hence, the buffering will be less for a fixed concentration of such a Tris buffer than for a buffer of the same concentration with a  $pK_a$  closer to the pH to be used.

An advantage of Tris is that both the free base and its hydrochloride are commercially available, so that a buffer can be made reproducibly without depending on pH measurements, by weighing out the required amounts of each form. This is valuable for any buffer, but especially so for Tris; it shares with all amines the property of a large positive  $DH$  for dissociation of the  $R-NH_3^+$  form, so that the  $pK_a$  decreases with heating, by 0.028 per degree Celsius. Hence, accurate temperature control is necessary for pH measurements and adjustments.

The reactivity of Tris as a primary amine is not important for most uses, and its transparency in the UV can be helpful.

## Trisomy

Trisomy is a special case of [aneuploidy](#) where one [chromosome](#) (or part of a chromosome) is represented in three copies instead of the normal two for a **diploid** cell. An extra or missing chromosome is the most frequent cause of mental retardation in humans, in whom partial or total trisomies of all chromosomes have been described. The five chromosomes most involved in this type of human aneuploidy are X, Y, 18, 13 and 21.

Trisomy 21, associated with Down's syndrome, is the most common aneuploidy among humans born alive. The [karyotype](#) of trisomy 21 has 47 chromosomes in all the cells of the organism, or only in certain cells (**mosaic**). A karyotype of 46 chromosomes has, however, been described, in which the supernumerary chromosome 21 is totally translocated onto chromosome 14. Partial trisomies are caused by chromosomal translocations. One example among many is a partial trisomy 21 of the short arm-centromere proximal long arm segment caused by maternal translocation to chromosome 12 (1). Partial trisomy 8 is accompanied by mild mental retardation, strabism, and anomalies of the skeleton. Trisomy 13 is characterized by altered skull and face morphology and often by polydactyly. Partial trisomy 22 is believed to be associated with **retinoblastoma**, etc.

### 1. Fine Localization of Trisomy 21

The preponderant role of the distal third of the long arm in Down's syndrome of trisomy 21 is well established. Trisomy of only band 21q22 results in a state identical to that caused by complete trisomy 21. If the trisomy involves only a part of the above band, the symptoms are diminished, but the appearance of the patient is still reminiscent of Down's syndrome. Monosomy in band 21q22 is lethal. Trisomy of the adjacent region 21q21 is not associated with malformations but leads to mental retardation (2).

### 2. Parental Origin and Cell Stage of Chromosomal Nondisjunction

In a study of 62 aneuploids with total trisomy 18, 56 were maternal in origin and only six paternal. Among the 56 maternally derived trisomics, a postzygotic error in **mitosis** could be excluded in 52 cases. Among those in which chromosomal nondisjunction was attributable to an error in **meiosis**, 11 were the result of meiosis I nondisjunction and 17 were caused by a meiosis II error. This result differs markedly from findings in **acrocentric** chromosomes, where nondisjunction at maternal meiosis I predominates. Among the six paternally derived cases, two originated from a meiotic error, showing that nondisjunction in paternal meiosis is not as rare as previously suggested (3).

Out of 100 cases of trisomy 21, the overwhelming majority of 94 were maternal in origin. This value, obtained with specific DNA markers, is significantly greater than the 75 to 80% nondisjunction rate observed in cytogenetic studies and illustrates the increased accuracy of the molecular approach (4). Among 36 cases of total trisomy 21 in which the supernumerary chromosome was of paternal origin, 15 were consistent with meiosis II errors, 8 with mitotic errors, and only 7 with meiosis I nondisjunction, this contrasts with maternally derived trisomy 21, in which meiosis I errors predominate. Study of the correlation with parental age indicates that the well-known increase in frequency of trisomy 21 with increasing age of the mother is caused by maternal nondisjunction (5).

### 2.1. Biochemical Studies

Patients with Down's syndrome due to trisomy 21 exhibit elevated activity of copper zinc superoxide dismutase, caused by the increased dosage of the **gene** for this **enzyme**, which is encoded at position 21q.22.1 of chromosome 21. Conflicting results have been reported, however, implicating this enzyme in Down's syndrome (6),(7).

### Bibliography

1. H. Chen, M. Tyrkus, and P. V. Woolley (1978) *Ann. Genet.* **21**, 177–180.
2. M. O. Rethore (1981) *Hum. Genet. Suppl.* **2**, 173–182.
3. T. Eggerman et al. (1996) *Hum. Genet.* **97**, 218–223.
4. S. L. Sherman et al. (1991) *Am. J. Hum. Genet.* **49**, 608–620.
5. M. B. Petersen et al. (1993) *Hum. Mol. Genet.* **2**, 1691–1695.
6. D. Minc-Golomb, H. Knobler, and Y. Groner (1991) *EMBO J.* **10**, 2119–2124.
7. R. de la Torre et al. (1996) *Experientia* **52**, 871–873.

## Trithorax Group Genes

The trithorax group (trxG) consists of a genetically defined class of **genes** responsible for maintaining the active expression state of **homeotic genes**. Through the differential expression pattern of homeotic genes, cells become programmed to support specific structures and functions. Mutations in any member of the trxG result in reduced homeotic gene expression and in homeotic body pattern transformations. For the reciprocal process, the genes of the **Polycomb group** (PcG) are necessary for maintaining the repressed state of homeotic gene expression. In *Drosophila*, both groups are part of what is called a “cellular memory” mechanism. Their task is not to initiate the differential expression pattern of homeotic genes, which in *Drosophila* is performed by the early-acting patterning factors encoded by maternal and segmentation genes, but to maintain the genetic expression status stably and heritably over developmental time (for reviews, see Refs. 1, 2). Subsequently, it was found that many other developmentally important genes, whose expression patterns need to be tightly and faithfully maintained, are also targets of this type of regulation in

addition to the homeotic genes.

Molecular analysis of several members of the two groups indicate involvement of their protein products at the level of [chromatin](#) structure. Each group acts in large multiprotein complexes to fulfill its functions. Although PcG proteins cooperate to generate **silenced** chromatin structures, to keep genes inactive, the proteins of the trxB cooperate to counteract the PcG-repressed chromatin by modifying chromatin structures for transcriptional activation.

## 1. The Trithorax protein and its targets

The gene *trithorax* (*trx*), the name-giving member of the group, was identified by its requirement for normal expression of multiple homeotic genes of the **bithorax** and [Antennapedia complex](#). Different homeotic genes and their **promoters** are affected differently in *trx* mutants, suggesting a coordinated spatial and temporal requirement for *trx* to achieve normal expression patterns (3, 4). Several parts of homeotic gene expression are dispensable of *trx* function and probably require other factors of the trxB, like *ash1*, for maintenance (5, 6).

The *trx* gene expresses a complex pattern of [messenger RNA](#) transcripts which encode two large protein isoforms (4, 7, 8). The TRX proteins are characterized by containing several evolutionarily conserved protein motifs, a novel variant of the nuclear receptor-type **DNA-binding domain**, a 130- to 140-residue motif known as the SET domain, and a C<sub>4</sub>HC<sub>3</sub> **zinc-finger** motif, known as the PHD finger. The SET domain is also found in ASH1 and in other chromatin-associated proteins, including Enhancer of Zeste, a member of the PcG (6, 9). The SET domain was functionally analyzed in yeast on a protein having this conserved part. It was demonstrated that the **domain** is necessary for transcriptional silencing and for other cellular processes dependent on defined chromatin structures (10). The PHD finger is found in a large class of nuclear proteins, most of which function as adapters between specific activator proteins and other components of the transcription machinery. A potential role in mediating [protein-protein interactions](#) has been proposed, but the precise molecular function of the domain remains unknown. Although no direct sequence-specific DNA-binding activity has been demonstrated for TRX *in vitro*, it was found by **immunostaining** that the protein binds to approximately 75 chromosomal sites on [polytene chromosomes](#) (7, 11). Interestingly, many TRX binding sites overlap with binding sites of PcG proteins at chromosomal elements known as PcG response elements (PRE), suggesting a functional cross talk between these counteracting factors (11, 12).

## 2. Brahma and the chromatin-remodeling machines

The trxB member *brahma* (*brm*) was isolated in a genetic screen searching for dominant **suppressors** of homeotic transformations produced in Polycomb mutants (13). The idea behind this screen was to identify additional genes involved in inactivating the expression of homeotic genes. Indeed, it was found that *brahma* mutants cause developmental defects similar to other mutants that fail to express homeotic genes adequately (14). Molecular analysis showed that the Brahma protein is homologous to the yeast SWI2 protein, a DNA-dependent [ATPase](#) (15) and also part of a multiprotein complex with similarities to the well-known SWI/SNF complex from yeast.

Identification of the function of Brahma and its relationship to yeast components involved in transcriptional activation was the key to isolating other components associated in what are now called chromatin-remodeling machines (16). ISWI is another SWI2 **homologue** of *Drosophila* and is found in several multiprotein complexes involved in opening chromatin structures in an energy-dependent fashion (17, 18). Thus far, four different machines involved in chromatin-remodeling processes have been isolated in *Drosophila*: (1) the BRM-SWI/SNF complex, whose ATPase activity is DNA-dependent; (2) the NURF complex, whose ATPase activity relies on chromatin; (3) CHRAC, where energy is used to increase DNA accessibility in chromatin; and (4) ACF which, like the previous two, is also an ISWI-containing and ATP-utilizing chromatin-assembly and remodeling

factor (for a review, see Ref. [19](#)). The GAGA factor was originally identified biochemically as a DNA-binding transcriptional activator ([20](#)). GAGA factor cooperates with complexes like NURF to remodel chromatin. The recent isolation of mutations in the GAGA factor-encoding gene revealed a **phenotype** similar to *trx*. Hence, the gene was called Trithorax-like (TRL) and was classified as a *trxG* member ([21](#)).

### 3. Vertebrate homologues

Several members of the *Drosophila trxG* have counterparts in humans and mice. The mixed-lineage leukemia (MLL) gene (also known as ALL-1, HRX, Htrx), a TRX homologue, was found through its involvement in the majority of infantile acute lymphocytic and mixed lineage leukemias (for a review, see Ref. [22](#)). Disrupting the murine *Mll* gene by gene targeting causes homeotic transformations of the vertebrae in heterozygotes and loss of homeotic gene expression, supporting the notion that *Mll* is a functional equivalent of *trx* ([23](#)). Moreover, conserved components of chromatin-remodeling complexes were found in mammalian systems. Two genes, *mbrm* and *brg1*, are homologues of *brm* and also part of a SWI/SNF-like complex ([24-26](#)). Another SWI/SNF component is encoded by the human *HSNF5 (Ini1)* gene, and the putative *Drosophila* homologue *snr1* interacts genetically with *trx* and *brm* ([16, 27](#)).

Initially, the genetic definition of the *trxG* been an entry point to identify genes involved in maintaining homeotic gene expression. It is becoming increasingly clear, however, that several members of the group fulfill much broader tasks and are parts of the general transcriptional activation complexes, thus allowing an assessment of the molecular interactions of these important regulatory factors that control transcriptional and eventually, other chromosomal processes at the level of chromatin structure.

### Bibliography

1. J. A. Kennison (1995) *Ann. Rev. Genet.* **29**, 289–303.
2. R. Paro and P. J. Harte (1996) *Epigenetic Mechanisms of Gene Regulation* 507–528.
3. T. R. Breen and P. J. Harte (1993) *Development* **117**, 119–134.
4. Y. A. Sedkov, S. Tillib, L. I. Mizrokhi, and A. Mazo (1994) *Development* **120**, 1907–1917.
5. D. R. LaJeunesse and A. Shearn (1995) *Mech. Dev.* **53**, 123–139.
6. N. Tripoulas, D. R. LaJeunesse, J. Gildea, and A. Shearn (1996) *Genetics* **143**, 913–928.
7. B. Kuzin, S. Tillib, Y. A. Sedkov, L. I. Mizrokhi, and A. Mazo (1994) *Genes Dev.* **8**, 2478–2490.
8. M. J. Stassen, D. Bailey, S. Nelson, V. Chinwalla, and P. J. Harte (1995) *Mech. Dev.* **52**, 209–223.
9. T. Jenuwein, G. Laible, R. Dorn, and G. Reuter (1998) *Cell. Mol. Life Sci.* **54**, 80–93.
10. C. Nislow, E. Ray, and L. Pillus (1997) *Mol. Cell. Biol.* **8**, 2421–2436.
11. V. Chinwalla, E. P. Jane, and P. J. Harte (1995) *EMBO J.* **14**, 2056–2065.
12. V. Orlando, E. P. Jane, V. Chinwalla, P. J. Harte, and R. Paro (1998) *EMBO J.*, in press.
13. J. A. Kennison and J. W. Tamkun (1988) *Proc. Natl. Acad. Sci. USA.* **85**, 8136–8140.
14. B. J. Brizuela, L. K. Elfring, J. Ballard, J. W. Tamkun, and J. A. Kennison (1994) *Genetics* **137**, 803–813.
15. J. W. Tamkun, R. Deuring, M. P. Scott, M. Kissinger, A. M. Pattatucci, T. C. Kaufman, and J. A. Kennison (1992) *Cell* **68**, 561–572.
16. A. K. Dingwall, S. J. Beek, C. M. McCallum, J. W. Tamkun, G. V. Kalpana, S. P. Goff, and M. P. Scott (1995) *Mol. Cell. Biol.* **6**, 777–791.
17. L. K. Elfring, R. Deuring, C. M. McCallum, C. L. Peterson, and J. W. Tamkun (1994) *Mol. Biol. Cell.* **14**, 2225–2234.
18. T. Tsukiyama, C. Daniel, J. Tamkun, and C. Wu (1995) *Cell* **83**, 1021–1026.

19. P. D. Varga-Weisz and P. B. Becker (1998) *Curr. Opin Cell Biol.* **10**, 346–353.
20. H. Granok, B. A. Leibovitch, C. D. Shaffer, and S. C. R. Elgin (1995) *Curr. Biol.* **5**, 238–241.
21. G. Farkas, J. Gausz, M. Galloni, G. Reuter, H. Gyurkovics, and F. Karch (1994) *Nature* **371**, 806–808.
22. J. E. Rubnitz, F. G. Behm, and J. R. Downing (1996) *Leukemia* **10**, 74–82.
23. B. D. Yu, J. L. Hess, S. E. Horning, G. A. J. Brown, and S. J. Korsmeyer (1995) *Nature* **378**, 505–508.
24. H. Kwon, A. N. Imbalzano, P. A. Khavari, R. E. Kingston, and M. R. Green (1994) *Nature* **370**, 477–481.
25. J. C. Reyes, C. Muchardt, and M. Yaniv (1997) *J. Cell Biol.* **137**, 263–274.
26. W. Wang, Y. Xue, S. Zhou, A. Kuo, B. R. Cairns, and G. R. Crabtree (1996) *Genes Dev.* **10**, 2117–2130.
27. C. Muchardt, C. Sardet, B. Bourachot, C. Onufryk, and M. Yaniv (1995) *Nucleic Acids Res.* **23**, 1127–1132.
28. N. Zak, C. Miller, T. Alon, R. Goldman-Levi, K. Watson, and M. Crosby (1998) *A. Conf. Dros. Res.* **39**, 500B.
29. A. L. Adamson, and A. Shearn (1996) *Genetics*, 621–633.

### Suggestions for Further Reading

30. R. E. Kingston, C. A. Bunker, and A. N. Imbalzano (1996) Repression and activation by protein multiprotein complexes that alter chromatin structure, *Genes Dev.* **10**, 905–920. A comprehensive review of chromatin remodeling.
31. A. Schumacher and T. Magnuson, (1997) Murine Polycomb- and trithorax-group genes regulate homeotic pathways and beyond, *Trends Genet* **13**, 167–170.
32. A. Gould (1997) Functions of mammalian Polycomb group and trithorax group related genes, *Curr. Opin Genet. Dev.* **7**, 488–494.

### Triton X-100 And X-114

Structural variations in **hydrophobic** side chains produce two different **detergents** in the Triton family, Triton X-100 and Triton X-114, that have markedly different physicochemical properties (see **Detergents**). Early studies showed that the polyoxyethylene glycol detergent Triton X-100 has a very low **critical micelle concentration** (cmc of 0.2 mM and is highly efficient in solubilizing **membrane proteins** without dissolving much membrane **lipid** (1). This strongly hydrophobic detergent and the structurally related Nonidet P-40 are widely used in lysing cells and solubilizing membrane proteins. Integral membrane proteins, whose functional activity does not critically depend on membrane lipids, are successfully solubilized by using Triton X-100. Examples include the receptors for **GABA** (2), prostacyclin (3), prolactin (4), **transferrin** (5), and **insulin** (6). Although Triton X-100 is highly efficient in solubilizing membrane proteins, it is relatively mild and **nondenaturing** toward soluble, globular proteins, such as the g-globulins (see **Immunoglobulin**). This is underscored by the large number of current reports, in which cell lysis in a buffer containing 1% Triton X-100 or the closely related detergent Nonidet P-40, is successfully followed by **immunoprecipitation** of specific proteins (7-9).

Triton X-114 is prominent in recent literature for specific reasons. This detergent has a strikingly

greater solubility in water at lower temperatures (0°C), whereas at higher temperatures (30°C) it separates from water to form a separate detergent phase (10). Upon the formation of a discrete layer, the detergent retains specific, detergent-solubilized proteins, so the property of phase separation has been used in protein purification studies. When membrane proteins of adrenal medullary chromaffin granules are subjected to such fractionation in Triton X-114, cholesterol- and phospholipid-associated proteins like **ATPase I** and **glycoprotein IV** are obtained as detergent-insoluble proteins following solubilization at 0°C. Next, the detergent solution is warmed and layered over a cushion of 0.25 M sucrose in buffer containing 0.06% Triton X-114 and then centrifuged (see [Density Gradient Centrifugation](#)). The resulting aqueous phase (top layer, 0.04% Triton X-114) contains a mixture of soluble proteins, chromogranin A, soluble DBH, and membrane glycoproteins III, H, J, and K. Then the glycoproteins are isolated by removing the top layer, followed by exhaustive dialysis at 4°C in a buffer containing 1% Amberlite XAD-2 (11). Based on this phase separation method and the earlier observation that detergents with low cmc values (such as Triton X-100 and Triton X-114) are relatively inefficient in solubilizing the **GPI-anchored** proteins, Hooper and Bashir developed a technique of differential solubilization and temperature-induced phase separation in Triton X-114 to distinguish between the GPI-anchored proteins and those anchored by a simple membrane-spanning polypeptide (12, 13). When this method is applied to pig kidney microvillar membranes, abundant in both GPI-anchored and polypeptide-linked ectoenzymes, Triton X-114 at 0°C solubilizes only the ectoenzymes harboring a polypeptide anchor, whereas the GPI-linked ectoenzymes are sedimented by low-speed centrifugation. Then the detergent-solubilized supernatant is further fractionated by phase separation at 30°C into an aqueous phase and a detergent-rich phase, which contains the polypeptide-linked ectoenzymes.

## Bibliography

1. P. Banerjee, J. B. Joo, J. T. Buse, and G. Dawson (1995) *Chem. Phys. Lipids* **77**, 65–78.
2. T. N. Sato and J. H. Neal (1989) *J. Neurochem.* **52**, 1114–1122.
3. A. K. Dutta-Roy and A. K. Sinha (1987) *J. Biol. Chem.* **262**, 12685–12691.
4. H. Okamura, S. Raguette, A. Bell, J. Gagnon, and P. A. Kelly (1989) *J. Biol. Chem.* **264**, 5904–5911.
5. A. P. Turkewitz et al. (1987) *J. Biol. Chem.* **263**, 8318–8325.
6. Y. F. Yamaguchi and J. T. Harmon (1988) *Biochem.* **27**, 3252–3260.
7. R. M. Kluck, E. Bossy-Wetzel, D. R. Green, and D. D. Newmeyer (1997) *Science* **275**, 1132–1136.
8. P. Erhardt and G. M. Cooper (1996) *J. Biol. Chem.* **271**, 17601–17604.
9. H. Mischak et al. (1996) *Mol. Cell. Biol.* **16**, 5409–5418.
10. A. Sánchez-Ferrer, R. Bru, and F. García-Carmona (1994) *Critical Rev. Biochem. Mol. Biol.* **29**, 275–313.
11. J. G. Pryde and J. H. Phillips (1986) *Biochem. J.* **233**, 525.
12. N. M. Hooper and A. J. Turner (1988) *Biochem. J.* **250**, 865.
13. N. M. Hooper and A. Bashir (1991) *Biochem. J.* **280**, 745.

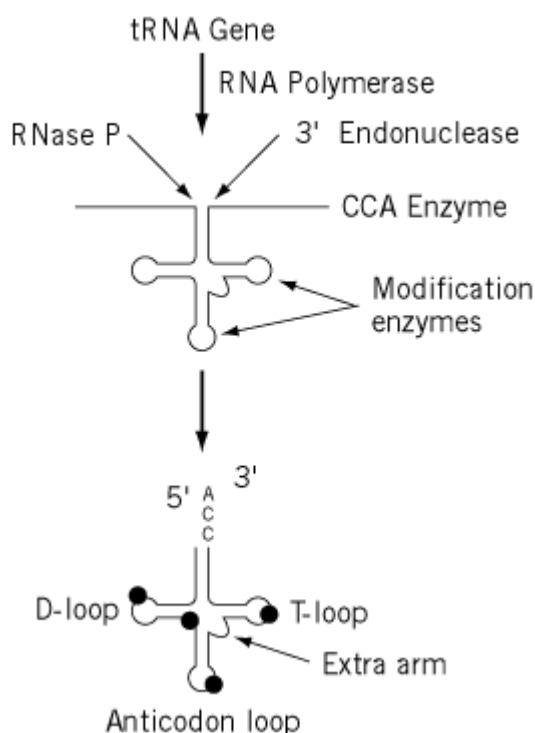
## tRNA Biosynthesis

Depending on the organism and the organelle, any particular [transfer RNA](#) (tRNA) may be **transcribed** from the [genome](#) into an RNA molecule along with other tRNA, with messenger RNA



(mRNA), or with ribosomal RNA (rRNA). Therefore, the first step in tRNA biosynthesis in many cases is the initial release of a precursor tRNA from such a **polycistronic** transcript. tRNA precursors share the need for certain **posttranscriptional processing**, either as part of the releasing process or subsequent to it. These shared needs include removal of the 5'-leader and 3'-trailing sequences and modification of bases (Fig. 1). Addition of the trinucleotide CCA at the 3' end is thought of as a universal requirement, and it is in **eukaryotes**, because the CCA end is not encoded in the gene. In contrast, many prokaryotic tRNA gene transcripts do not require CCA addition as an initial processing event, because it is provided in the primary transcript and becomes the 3'-terminal sequence after other 3'-processing **nucleases** remove downstream nucleotides. Some tRNA require additional posttranscriptional alterations for their maturation. These include 5'-base addition, removal of intervening sequences, and RNA **editing**. While there are critical sequence-specific requirements for some RNA processing enzymes, it is safe to state the general rule that the most critical substrate characteristic for tRNA processing enzymes is retention of the overall tRNA structure in the precursor. It has also been shown that tRNA processing can proceed via alternative routes such that end processing can either precede or follow splicing (1) and splicing can begin with cleavage at either the 5' or 3' end of the intron (2) (see discussion below).

**Figure 1.** Schematic diagram of tRNA biosynthesis. Transcription produces a pre-tRNA that is trimmed on the 5' end by RNase P and on the 3' end by an endonuclease (or exonuclease) in preparation for CCA addition. Base modification reactions occur throughout the processing pathway. Some tRNA genes have intervening sequences which are not shown here but are described in the text.



## 1. Removal of 5' Leader

The 5' leader is removed from tRNA precursors by [ribonuclease P](#) (RNaseP), an enzyme found in all organisms and organelles in which tRNA genes are found. The cleavage reaction is a hydrolysis, requires both monovalent and divalent cations, and leaves a 5'-phosphate on the mature tRNA and a 3'-hydroxyl on the 5' leader. Prokaryotic, eukaryotic and organelle RNaseP enzymes, with some possible exceptions, are [ribonucleoproteins](#) and require both the RNA and protein subunits for

activity *in vivo*. Remarkably, the RNA subunit is responsible for catalysis in the bacterial enzymes (3) and is likely to serve the same function in eukaryotes. All RNaseP RNAs share short regions of sequence similarity and are predicted to take on similar elements of secondary structure. Although they can vary in size, most cluster in the 300–400-nucleotide range. The prokaryotic enzymes are about 10% protein, with a protein subunit of about 14 kDa (kilodaltons). Archaeal and eukaryotic enzymes are 50–70% protein, and the proteins identified so far are larger and share no sequence similarity to each other or to the prokaryotic proteins.

## 2. Removal of 3' Leader

In some organisms, exoribonucleases perform 3'-end tRNA processing, while in others, endonucleases carry out that function. The best-studied exoribonucleases are from *Escherichia coli*. Six exoribonucleases that can trim tRNA precursors *in vitro* have been identified (4). Five can participate in 3'-end formation *in vivo*, but extensive functional overlap is demonstrated by the observation that cells retaining only one of these nucleases are viable. Complete 3' processing *in vitro* required **ribonuclease II** and [polynucleotide phosphorylase](#) to shorten long 3' trailers to intermediates extended by only two to four nucleotides. Final trimming by ribonucleases T and/or PH completes the process (5). The combined genetic and biochemical approaches of Deutscher and his colleagues have revealed that the relative importance of each exoribonuclease is different for different tRNA precursors, and that the combined action of the nucleases is most effective in maturing all tRNA. Whether exonuclease(s) or endonuclease(s) mature the 3' end of nucleus coded tRNA in *S. cerevisiae* is still not established. However, examination of tRNA precursors and 3' trailers from *Drosophila* (6), *Xenopus* (7), and silkworms (8) demonstrate that endonuclease cleavage prepares the 3' end of nucleus-coded tRNA in these eukaryotes (9). Organelle tRNA precursor structure also establishes that 3'-end formation of mitochondrial and chloroplast tRNA is carried out by endoribonucleases. Indeed, these enzymatic activities have been measured in crude extracts and in some cases partially purified (10).

## 3. Addition of CCA 3' End

ATP(CTP)-tRNA-specific nucleotidyltransferase adds the CCA end to tRNA, following 3'-end processing to provide the proper substrate. The enzyme also serves as a repair enzyme to restore the CCA end to tRNA victims of cellular nucleases and is required for this function regardless of whether the CCA end is gene-encoded. Since eukaryotic precursor tRNAs that contain intervening sequences have the CCA end, it is clear that the enzyme can bind partially processed tRNA. This is consistent with a model in which the enzyme interacts with tRNA at the corner of the tRNA where the T and D loops are juxtaposed and extends that interaction across the aminoacyl stem to the 3' end (11).

## 4. Base Modification

Transfer RNA contains the highest concentration of modified bases of any type of RNA. There are 80 different nucleoside modifications known, and many display species and/or phylogenetic specificity (12). Modified bases are introduced after transcription and, depending on the enzyme, either precursor or fully processed tRNA can serve as the preferred substrates. In both prokaryotes and eukaryotes, characterization of precursors demonstrates that modifications do occur in a stepwise fashion (13, 14). Further enforcing an order of addition for nucleus-coded tRNAs, certain tRNA modification enzymes are confined to either the [nucleus](#) or the **cytosol**. Mutants deficient in a particular modification enzyme activity, so that a particular modification is not made, do complete subsequent modifications (10). Thus it is clear that an obligate order of modification per se is not required. Certain [anticodon](#) modifications are not found in unspliced precursors, suggesting that there is a requirement for splicing prior to the synthesis of the modified bases located in the anticodon region (15).

Histidyl tRNAs contain an extra 5'-G nucleotide relative to all other tRNAs. In *E. coli*, the G is derived transcriptionally, and RNaseP does not remove the extra G when it processes the 5' end (16). In yeast, however, the 5'-most G in the histidyl tRNA is added by a histidyl tRNA guanylyltransferase in an ATP-dependent reaction (17). The 5'-G of spinach chloroplast histidyl tRNA is derived by transcription (18), whereas that of animal mitochondria is added posttranscriptionally (19).

## 5. Intron Removal

Since intervening sequences (introns) were discovered in yeast tRNA genes (20), the mechanisms involved in their removal have been sought. Splicing is initiated by an endonuclease that cleaves the pre-tRNA at the splice sites. Biochemical purification of the enzyme has been completed, and it is known to consist of three polypeptide subunits of 31, 42, and 51 kDa (21). Release of the intron leaves a 2',3'-cyclic phosphate on the 3' end of the 5' half of the tRNA (22), and a 5'-OH on the 5' end of the 3'-half of the tRNA (23). The two half-molecules serve as the substrate for a tRNA ligase, which carries out the following steps: (i) creation of a 2'-phosphate by **phosphodiesterase** activity at the 2',3'-cyclic phosphate; (ii) phosphorylation of the 5'-OH on the 3' half by a **polynucleotide kinase** that uses GTP as the phosphate donor (24); (iii) transfer of AMP from the **adenylylated** form of the enzyme to the 5'-P on the 3' half of the tRNA; and (iv) joining of the tRNA halves with the release of AMP (25). At this stage, the product contains a 3',5'-phosphodiester linkage, as well as a 2'-phosphomonoester. The latter is removed by an **NAD-dependent 2'-phosphotransferase** in a reaction in which the 2'-phosphate displaces the nicotinamide part of NAD, to yield ADP-ribose 1''–2'' cyclic phosphate (26). The *RLG1* gene coding for the ligase has been characterized, and deletion analysis supports the idea that various activities can be attributed to separate domains in the protein (27). Similar enzyme activities are present in **HeLa** cell extracts, demonstrating that this is a conserved mechanism (28).

## 6. Editing

Like other RNA, tRNAs are substrates for **editing**, the remarkable process that alters the base sequence of an RNA after synthesis. So far, this type of tRNA processing has been identified only in **mitochondria**. Insertional editing in *Physarum* restores functional structure to mitochondrial tRNA (29). Substitutional editing also usually restores base pairing in stems. In *Acanthamoeba castellanii* A can be changed to G, U to G, and U to A, to correct mismatched base pairs (30). In the land snail, mitochondrial tRNAs sustain changes of C, U, and G to A that in all cases, except one, restore base pairing in the stem (31). This restoration of base pairing has been shown to be required for efficient removal of 5' leaders and 3' trailers in plant mitochondria (32, 33) consistent with the observation that tRNA-processing enzymes require a tRNA structure for optimal function. Finally, a C-to-U change in the second position of the anticodon of marsupial mitochondrial tRNA<sup>Asp</sup> is necessary for proper codon recognition (34). Although the mechanisms are not clear, there is no reason to expect that tRNA transcript editing is different from those editing processes operating on other RNAs.

## Bibliography

1. J. P. O'Connor and C. L. Peebles (1991) *Mol. Cell. Biol.* **11**, 425–439.
2. F. Miao and J. Abelson (1993) *J. Biol. Chem.* **268**, 672–677.
3. C. Guerrier-Takada, K. Gardiner, T. Marsh, N. R. Pace, and S. Altman (1983) *Cell* **35**, 849–857.
4. N. B. Reuven and M. P. Deutscher (1993) *FASEB J.* **7**, 143–148.
5. Z. Li and M. P. Deutscher (1994) *J. Biol. Chem.* **269**, 6064–6071.
6. D. Friendewey, T. Dingermann, L. Cooley, and D. Söll (1985) *J. Biol. Chem.* **260**, 449–454.
7. J. G. Castano, J. A. Tobias, and M. Zasloff (1985) *J. Biol. Chem.* **260**, 9002–9008.

8. R. L. Garber and L. P. Gage (1979) *Cell* **18**, 817–829.
9. O. Hagenbuchle, P. Larson, G. I. Hall, and K. U. Sprague (1979) *Cell* **18**, 1217–1229.
10. A. K. Hopper and N. C. Martin (1992) *The Molecular and Cellular Biology of the Yeast Saccharomyces: Gene Expression*, Vol. **II**, Cold Spring Harbor Press, Cold Spring Harbor, N. Y., pp. 99–141.
11. P. Spacciapoli, L. Doviken, J. J. Mulero, and D. L. Turlow (1989) *J. Biol. Chem.* **264**, 3799.
12. P. F. Crain and J. A. McCloskey (1997) *Nucleic Acids Res.* **25**, 126–127.
13. H. Sakano, S. Yamada, T. Ikemura, Y. Shimura, and H. Ozeki (1974) *Nucleic Acids Res.* **1**, 355.
14. K. Nishikura and E. M. DeRobertis (1981) *J. Mol. Biol.* **145**, 405.
15. G. Knapp, J. S. Beckmann, P. F. Johnson, S. A. Fuhrman, and J. Abelson (1978) *Cell* **14**, 221.
16. O. Orellana, L. Cooley, and D. Söll (1986) *Mol. Cell. Biol.* **6**, 525.
17. S. Pande, D. Jahn, and D. Söll (1991) *J. Biol. Chem.* **266**, 22826.
18. U. Burkard and D. Söll (1988) *J. Biol. Chem.* **263**, 9578–9581.
19. D. L'Abbe, B. F. Lang, P. Desjardins, and R. Morais (1990) *J. Biol. Chem.* **265**, 2988–2992.
20. H. M. Goodman, M. V. Olson, and B. D. Hall (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5453–5457.
21. R. Rauhut, P. R. Green and J. Abelson (1990) *J. Biol. Chem.* **265**, 18180–18184.
22. C. L. Peebles, P. Gegenheimer, and J. Abelson (1983) *Cell* **32**, 525–536.
23. G. Knapp, R. C. Ogden, C. L. Peebles, and J. Abelson (1979) *Cell* **18**, 37–45.
24. H. G. Belford, S. K. Westaway, J. Abelson, and C. L. Greer (1993) *J. Biol. Chem.* **268**, 2444–2450.
25. C. L. Greer, C. L. Peebles, P. Gegenheimer, and J. Abelson (1983) *Cell* **32**, 537–546.
26. G. M. Culver, S. M. McCraith, M. Zillmann, R. Kierzek, N. Michaud, R. D. LaReau, D. H. Turner, and E. M. Phizicky (1993) *Science* **261**, 206–208.
27. B. L. Apostol, S. K. Westaway, J. Abelson, and C. L. Greer (1991) *J. Biol. Chem.* **266**, 7445–7455.
28. M. Zillmann, M. A. Gorovsky, and E. M. Phizicky (1992) *J. Biol. Chem.* **267**, 10289–10294.
29. D. Miller, R. Mahendran, M. Spottswood, H. Costandy, S. Wang, M. L. Ling, and N. Yang (1993) *Semin. Cell. Biol.* **4**, 261–266.
30. K. M. Lonergan and M. W. Gray (1993) *Science* **259**, 812–816.
31. S. Yokobori and S. Paabo (1995) *Proc. Natl. Acad. Sci. USA* **92**, 10432–10435.
32. A. Marchfelder, A. Brennicke, and S. Binder (1996) *J. Biol. Chem.* **271**, 1898–1903.
33. L. Marechal-Drouard, A. Cosset, C. Remacle, D. Ramamonjisoa, and A. Dietrich (1996) *Mol. Cell. Biol.* **16**, 3504–3510.
34. A. Janke and S. Raabo (1993) *Nucleic Acids Res.* **21**, 1523–1525.

## tRNA Ligase

tRNA ligase is a form of [RNA ligase](#) that is thought to function in all eukaryotic cells. Its biological role is to catalyze the splicing of the precursor to [transfer RNA](#). In contrast to the [RNA splicing](#) that occurs with [messenger RNA](#), which proceeds via transesterification, without a protein to catalyze the

cleaving and ligating reactions, tRNA splicing requires both latter enzyme activities. The tRNA ligase of *Saccharomyces cerevisiae*, which is the most-studied (1), is a single [polypeptide chain](#) of 95.4 kDa comprised of **domains** that correspond to cyclic phosphodiesterase, kinase, and adenylation activities. The overall reaction is quite similar to those of T4 **polynucleotide kinase** plus T4 RNA ligase, the only difference being that a 2'-phosphate group occurs with tRNA ligase, rather than a 2'-hydroxyl (see Fig. 2 of [RNA Ligases](#)). The mechanism of the reaction involves (i) opening of the 2',3'-cyclic phosphate to a 2'-phosphomonoester (phosphodiesterase activity), (ii) addition of a phosphate group to the 5'-hydroxyl of the donor (kinase activity), (iii) transfer of an adenylyl group from an adenylylated enzyme, which is produced in advance by another kinase activity of the same enzyme, to the 5'-phosphate of the donor (adenylylation activity), and (iv) attack of the acceptor 3'-hydroxyl on the activated donor phosphoanhydride, to form a 3',5'-phosphodiester bond (bond formation activity). The mature tRNA is then generated by the action of a distinct enzyme, 2'-phosphotransferase, which removes the 2'-phosphate from the product.

Yeast tRNA ligase has a rather strong substrate-specificity, while wheat germ tRNA ligase, another well-studied tRNA ligase, has less specificity and can also ligate oligo(A) and oligo(U<sub>n</sub>G), in addition to tRNA precursors. The wheat germ and other plant (such as tobacco and *Chlamydomonas*) tRNA ligases are assumed to be exploited by [viroids](#) and [virusoids](#) for circularizing their **genomic RNAs** after **rolling-circle replication** and self-cleavage.

In addition to its catalytic function, tRNA ligase may play a role in transport within the [nucleus](#) of the precursor molecule that is to be spliced; this process seems to be performed by a group of PRP (precursor RNA processing) proteins in the case of mRNA splicing (2). Another novel function of yeast tRNA ligase is that it is involved in the splicing of the mRNA of Hac1 protein (a [transcription factor](#) required for the unfolded protein response) (3). This might unveil another important function of tRNA ligase, having a role in a crucial regulatory process.

Structural analysis of tRNA ligase is underway. In yeast tRNA ligase, Lys114 is responsible for accepting an adenylyl group by a phosphoamide bond. The region of the [primary structure](#) surrounding this lysine is well-conserved among tRNA ligases and T4 RNA ligase (see [RNA Ligases](#)) (Fig. 1). The region involved in nucleoside triphosphate-binding is also fairly well conserved among tRNA ligases, T4 polynucleotide kinase, and human 3'-cyclic nucleotide 3'-phosphodiesterase (4), implying that they have an evolutionary relationship.

**Figure 1.** Optimal alignment of a conserved region of RNA ligases. *S. cer.*, *Saccharomyces cerevisiae*; *S. pom.*, *Schizosaccharomyces pombe*; lysine (K) residue that accepts an adenylyl group is indicated by an arrow. Conserved amino acid residues are highlighted.

|                |      |    |         |         |      |         |        |           |           |          |      |      |          |      |   |
|----------------|------|----|---------|---------|------|---------|--------|-----------|-----------|----------|------|------|----------|------|---|
| <i>S. cer.</i> | tRNA | 59 | ITLPCN  | ARGLF   | ISDD | .       | TTNPV  | IV        | ARGYDKFFN | VGEV     |      |      |          |      |   |
| <i>S. pom.</i> | tRNA | 57 | ITLPTN  | ARGLF   | TGYD | YESKRHR | IV     | IRGYDKFFN | IDEV      |          |      |      |          |      |   |
| T4             | RNA  | 47 | LPDALEC | RGIMFEM | DGE  | .       | KPVR   | IAS       | RPMEKFFN  | LNEM     |      |      |          |      |   |
| ↓              |      |    |         |         |      |         |        |           |           |          |      |      |          |      |   |
| <i>S. cer.</i> | tRNA |    | ..CT    | GPYDVT  | I    | K       | ANGCI  | IF        | ISGL      | EDGTLVVC | SKHS | 1    |          |      |   |
| <i>S. pom.</i> | tRNA |    | ..TK    | GPYEL   | I    | V       | KENGCI | IF        | IAAL      | PDGQI    | I    | VS   | SKHS     | 1    |   |
| T4             | RNA  |    | ...     | VDY     | I    | L       | I      | .         | KED       | CSLV     | .    | STYL | DGDEILFK | SKGS | 1 |

1. B. Apostol and C. L. Greer (1991) *Nucleic Acids Res.* **19**, 1853–1860.
2. M. J. Moore, C. C. Query, and P. A. Sharp (1993) In *The RNA World* (R. F. Gesteland and J. F. Atkins, eds.), Cold Spring Harbor Laboratory Press, Cold Spring Harbor N.Y., pp. 303–357.
3. C. Sidrauski, J. S. Cox, and P. Walter (1996) *Cell* **87**, 405–413.
4. E. V. Koonin and A. E. Gorbalenya (1990) *FEBS Lett.* **268**(1), 231–234.

## TRP Operon

The **aromatic amino acid** tryptophan is synthesized by **plants** and by many **microorganisms**, but not by animals. The **genes** and **proteins** required for tryptophan biosynthesis were partially characterized very early in the history of modern molecular biology. Analysis and identification of all of the intermediates in the tryptophan pathway provided the background information and material that fostered numerous subsequent investigations. Studies were performed on the relationship between gene, **enzyme**, and biochemical reaction, enzyme structure and function, enzyme reaction mechanisms, and **operon** organization and regulation. Seven enzymatic functions are required to synthesize this amino acid from its common precursor, chorismate. Chorismate is also the source of the other amino acids, phenylalanine and tyrosine, as well as additional aromatic metabolites. Therefore, tryptophan formation must compete with other essential biosynthetic processes. The products of other pathways are also required for tryptophan formation. These include three co-substrates, L-glutamine, 5-phosphoribosyl-1-pyrophosphate, and L-serine, and the coenzyme, **pyridoxal phosphate**.

In this article we summarize our current understanding of the many features of the organization and functions of the *trp* operon of *Escherichia coli* and *Salmonella enterica*. We describe the operon's polypeptide products, the enzyme complexes they form, and how these enzymes catalyze the individual reactions of tryptophan formation. We also describe how operon expression is regulated and discuss the significance of the various regulatory mechanisms that these organisms use.

### 1. The *TRP* Operon of *E. coli* and *S. enterica*

The *trp* operon of *E. coli* and *S. enterica* is a single transcriptional unit that has five major structural genes (Fig. 1). The regulatory region has features that allow these organisms to sense and respond to changes in the availability of both tryptophan and tryptophan-charged tRNA<sup>Trp</sup>, Trp-tRNA<sup>Trp</sup> (2, 3). Transcription initiation at the principal **promoter**, *trpP1*, is regulated by the tryptophan-activated **trp repressor** (2). Continuation of transcription beyond the regulatory region into the structural genes of the operon is regulated by transcription **attenuation**. This mechanism permits these bacteria also to regulate operon expression in response to changes in the intracellular level of Trp-tRNA<sup>Trp</sup> (3). An inefficient internal promoter, *trpP2*, is located at the distal end of *trpD* (Fig. 1) (1). *trpP2* directs the synthesis of a minor transcript containing *trpC*, *trpB*, and *trpA* genetic information (1). The operon's terminus is defined by two adjacent sites of transcription termination, designated *t* and *t'* (1)(Fig. 1). The first of these, *t*, is a relatively inefficient factor-independent termination site, whereas the second, *t'*, is a site of efficient **Rho**-dependent termination.

**Figure 1.** The *trp* operon, and the intermediates, reactions, and enzymes of tryptophan biosynthesis. The *trp* operon of *E. coli* consists of five major structural genes (shaded rectangles) that encode the five polypeptides that catalyze tryptophan

formation. Two of the structural genes, *trpD* and *trpC*, have two genetic segments and encode bifunctional proteins. The structural genes are preceded by a transcribed leader region that contains *trpL*, a coding region for a 14-residue leader peptide. Transcription of the operon is initiated at the principal promoter, *trpP1*, and at an internal promoter, *trpP2*. (The transcription initiation sites are marked by horizontal arrows below the operon map) Transcription of the operon can be terminated at the attenuator (*attn.*) in the regulatory region or can continue to the end of the operon and stop at either of the tandem terminators, marked *t* and *t'*. Tryptophan biosynthesis from chorismate requires catalysis of seven reactions by enzyme domains designated TrpA through TrpG. The first two reactions are performed by a complex that contains the TrpE and TrpG domains. The next three reactions are catalyzed by the independent functional domains, TrpD, TrpF, and TrpC; however TrpG and TrpD are covalently joined, as are TrpC and TrpF. The last two reactions proceed at the separate active sites of TrpA and TrpB polypeptides of the tryptophan synthase complex, the product of the TrpA reaction, indole, being channeled to the active site of TrpB. The features of the *trp* operon of *S. enterica* and its products are virtually identical to those of *E. coli*. PRPP = 5-phosphoribosyl-1-pyrophosphate. G3P = D-glyceraldehyde-3-phosphate.





## 2. Reactions and Enzymes of Tryptophan Biosynthesis

Biosynthesis of tryptophan proceeds from chorismate, the precursor of most aromatic compounds in the cell, by the pathway shown in Fig. 1 (1, 2). The aromatic ring of tryptophan is derived from chorismate, the carbons and the nitrogen of the indolyl ring from 5-phosphoribosyl-1-pyrophosphate (PRPP) and L-glutamine, respectively, and the alanyl side chain from L-serine. In *E. coli* and *S. enterica* the seven reactions of the tryptophan pathway are catalyzed by three **multifunctional enzymes** made up of the five polypeptide products of the *trp* operon. Two of the polypeptides are bifunctional; each consists of two **domains**. Thus, the pathway relies on the sequential action of seven functional domains (Fig. 1).

The first three reactions of the pathway—production of  $\text{NH}_3$  from glutamine, synthesis of anthranilate from chorismate, and phosphoribosylation of anthranilate with PRPP—are carried out by an enzyme complex consisting of the TrpE and TrpG-D polypeptides, which have anthranilate synthase and glutamine amidotransferase-anthranilate phosphoribosyl (APR) transferase activities, respectively. The fourth and fifth reactions, the conversion of phosphoribosyl anthranilate (PRA) to 1-(*o*-carboxyphenylamino)-1-deoxyribulose-5-phosphate (CdRP) and of CdRP to indoleglycerol phosphate (IGP), are catalyzed by the bifunctional IGP synthase-PRA isomerase polypeptide (TrpC-F). The sixth and seventh reactions, the formation of indole from IGP, and L-tryptophan from indole and serine, are performed by the IGP aldolase (TrpA) and L-serine hydro-lyase (adding indole) (TrpB) activities of the tryptophan synthase complex ( $\text{TrpA}_2\text{TrpB}_2$ , or TSase  $a_2b_2$ ).

The enzymes of *E. coli* and *S. enterica* are very similar in structure and function. Active heterologous anthranilate synthase and tryptophan synthase complexes that contain one subunit from *E. coli* and one from *S. enterica* readily assemble and function *in vivo* and *in vitro*. Nevertheless, there is considerable variation among other microorganisms with respect to the molecular organization of the seven protein functional domains. In some bacteria, all seven domains exist independently, whereas in other bacteria and some fungi, a variety of domain fusions are found, including TrpE-TrpD, TrpG-TrpC, TrpG-TrpC-TrpF, and TrpA-TrpB.

## 3. Structure and Activity of the Tryptophan Biosynthetic Enzymes

### 3.1. Anthranilate Synthase and Anthranilate Phosphoribosyl Transferase

Anthranilate synthase, APR transferase of *E. coli* and *S. enterica* is a multifunctional, heterotetrameric, multienzyme complex ( $M_r = 228,200$ ) composed of two TrpE and two TrpG-D polypeptides (Fig. 1) (4, 5). Assembly of the subunits into the complex occurs spontaneously *in vitro* and depends solely on interactions between the TrpE and the TrpG domains. The tetrameric complex is very stable and has no detectable dissociation or subunit exchange in solution. Anthranilate synthase activity, the synthesis of anthranilate (*o*-aminobenzoate) from chorismate and glutamine (Fig. 1), is catalyzed by the TrpE and TrpG domains of the complex. The APR transferase reaction, the formation of *N*-(5'-phosphoribosyl)-anthranilate from anthranilate and 5-phosphoribosyl-1-pyrophosphate (Fig. 1), is catalyzed by the TrpD domain of the complex.

In the anthranilate synthase reaction, the TrpG domain acts as a classic glutamine amidotransferase, hydrolyzing L-glutamine via a g-glutamyl-*S*-cysteinyl enzyme intermediate, and delivering its amide group to the TrpE subunit as nascent  $\text{NH}_3$  (6). TrpG must be complexed with TrpE to elaborate its glutaminase activity, which is activated 20- to 30-fold by binding of chorismate at the active site of TrpE. The TrpE subunit uses the nascent  $\text{NH}_3$  to aminate chorismate, forming an aminocyclohexadiene reaction intermediate that it then aromatizes to anthranilate with the release of

pyruvate (7). Complexed TrpE can use exogenous  $\text{NH}_3$  in place of glutamine. The uncomplexed TrpE subunit, which lacks glutamine-dependent activity, can catalyze the  $\text{NH}_3$ -dependent anthranilate synthase reaction, albeit at a fivefold reduced rate. In contrast, the APR transferase activities of uncomplexed and complexed TrpG–TrpD are equivalent.

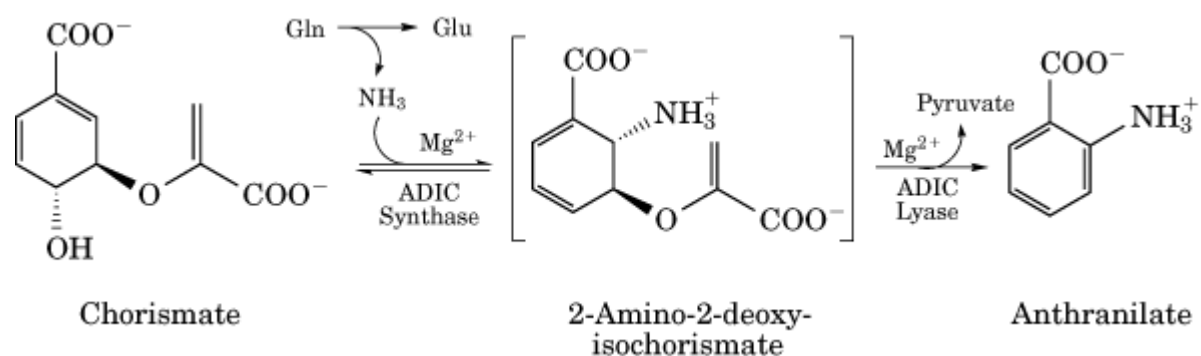
The anthranilate synthase and APR transferase reactions have an absolute requirement for a divalent metal cofactor, either  $\text{Mg}^{2+}$  or  $\text{Mn}^{2+}$ , whereas the glutaminase reaction is metal-independent.

**Electron paramagnetic resonance**, **NMR**, and **kinetic** studies have demonstrated binding interactions between metal, chorismate, and the TrpE subunit, and between metal and the TrpD domain (8). It is believed that the role of the metal ion in the anthranilate synthase reaction is to serve as a general acid in catalyzing the elimination of the C4 hydroxyl group of chorismate. All of the activities of the complex are sensitive to **feedback inhibition** by L-tryptophan, mediated by its binding at a regulatory site on the TrpE subunit.

The TrpE subunit contains 520 residues ( $M_r = 57,100$ ). Random mutational analysis and proteolytic probing has indicated that the polypeptide is composed of an amino-terminal regulatory domain of about 310 residues and a carboxy-terminal catalytic domain of about 200 residues (6, 9). Two separate clusters of residues have been identified by **mutagenesis** as essential for tryptophan binding and feedback inhibition: Glu39-Ser40-Ala41 and Asn288-Pro289-/Met293-Phe294-/Gly305. These residues are highly conserved in all known bacterial TrpE sequences. It is envisioned that these two polypeptide regions interact to form the tryptophan binding site in the folded polypeptide. Photoaffinity labeling of the TrpE subunit with 6-azidotryptophan, a false feedback inhibitor, has confirmed the role of the latter cluster as a structural component of the feedback site. Similar mutational analysis has identified five residues in the carboxy-terminal domain that are essential for catalytic activity: Thr329, His398, Thr425, Gly485, and Glu495. In addition, two arginine residues, Arg370 and Arg382, have been characterized as essential contact residues for assembling TrpE with TrpG in the formation of the intact complex.

The anthranilate synthase reaction proceeds as two sequential partial reactions (Fig. 2) (7). The first step, which is **rate-limiting**, is amination of chorismate, leading to the formation of the reaction intermediate, 2-amino-2-deoxyisochorismate (ADIC). ADIC remains enzyme-bound and is rapidly converted to anthranilate with the release of pyruvate. The two sequential reactions are referred to as the ADIC synthase and ADIC lyase reactions. Both the ADIC synthase and ADIC lyase activities of TrpE require  $\text{Mg}^{2+}$  and are feedback-inhibited by tryptophan. Mutant complexes that are substituted at His398 of the TrpE subunit have impaired ADIC lyase activity and therefore accumulate the intermediate. This made possible its isolation, identification, and use as a reagent in kinetic studies. It is likely that the role of His398 in the ADIC lyase reaction is either the abstraction of the proton from C-2 of ADIC or protonation of the leaving enol-pyruvyl group (Fig. 2).

**Figure 2.** The anthranilate synthase reaction. The reaction catalyzed by anthranilate synthase is a composite of two sequential subreactions, the conversion of chorismic acid to 2-amino-2-deoxy-isochorismate (ADIC) (ADIC synthase activity) and the conversion of ADIC to anthranilic acid (ADIC lyase activity). ADIC normally exists as an enzyme-bound reaction intermediate.



The TrpG-D subunit contains 531 residues ( $M_r = 57,000$ ). The TrpG domain is formed from the amino-terminal 193 residues of the polypeptide ( $M_r = 20,800$ ) and the TrpD domain consists of the carboxy-terminal 331 residues ( $M_r = 35,500$ ). The two domains are held together by a short peptide hinge. The bifunctional protein undoubtedly arose during evolution by the fusion of independent *trpG* and *trpD* genes (10). The TrpG and TrpD domains can be expressed as active, independent proteins, such as exist in other microorganisms, upon appropriate engineering of the *trpD* gene. The TrpD domain can also be removed from both the intact complex and the uncomplexed TrpG–TrpD subunit by *in vitro* trypsinolysis (5). Digestion by trypsin proceeds in three sequential steps from the carboxyl-terminal end of TrpD, suggesting the existence of three structural subdomains within TrpD. The limit cleavage event is at Lys193, the carboxy-terminal residue of TrpG. The resulting TrpE–TrpG “partial complex” is fully active in the anthranilate synthase reaction but lacks APR transferase activity. The independent glutamine amidotransferase domain assembles normally with TrpE, forming the tetrameric partial complex.

Anthranilate synthase is classified as a member of the family of Triad (formerly type G) glutamine amidotransferases (11). These are enzymes characterized by a glutamine amidotransferase domain that has active site Cys, His, and Glu residues arranged in a classic [catalytic triad](#), similar to the arrangement of the Cys/His/Glu catalytic triad of cysteinyl proteases. The relevant catalytic residues in TrpG are Cys84, His170, and Glu172. Cys84 is the nucleophilic residue involved in forming the *g*-glutamyl-*S*-cysteinyl enzyme intermediate. His170 and Glu172 play essential roles as general acid–base catalysts, first in the formation of the glutamyl intermediate and then in the formation and release of products of the reaction,  $\text{NH}_3$  and glutamate.

The three-dimensional structures of several Triad glutamine amidotransferases have been determined, including *E. coli* guanosine monophosphate synthetase and *E. coli* carbamoyl phosphate synthetase. Extrapolating from these structures, it has been postulated (11) that in the anthranilate synthase complex, activation of the glutaminase of the TrpG domain by the binding of chorismate to the active site of TrpE is the result of induced conformational changes that create both a “specificity pocket” for glutamine binding and a hydrophobic tunnel for the sequestered transfer of  $\text{NH}_3$  between the active sites.

Recently, the three-dimensional structures of three different microbial anthranilate synthases have been solved by X-ray crystallography. These include the native TrpE–TrpG complexes from *Sulfolobus solfataricus* (12) and *Serratia marcescens* (13) and the genetically engineered TrpE–TrpG partial complex from *Salmonella typhimurium* (14). The respective TrpE and TrpG folds of the three enzymes are similar, and the subunit structures can be readily compared by superpositioning. Each of the heterotetrameric enzymes is organized as a dimer of TrpE–TrpG dimers related by the crystallographic twofold rotational axis. The heterodimer interface is the same in the three enzymes; however, the heterotetramer interface of the *Salmonella* and *Serratia* enzymes involves TrpE–TrpE interactions, whereas that of the *Sulfolobus* enzyme consists mainly of TrpG–TrpG contacts. It has

been postulated that this difference may be related to the lack of cooperativity in reaction kinetics in the latter complex.

The TrpE subunit has a novel fold that has a complicated a/b topology featuring two structural subdomains. The core of the subunit is made up of two orthogonal, antiparallel  $\beta$ -sheets arranged to form a  $\beta$ -sandwich. The sandwich interface, formed from the hydrophobic surfaces of the  $\beta$ -sheets of the two subdomains, contains two cavities located  $\sim 18\text{\AA}$  apart. One cavity serves as the active site of the subunit that contains the binding sites for chorismate and metal cofactor. This site, which is made up of residues from each of the two subdomains of the sandwich, is appropriately oriented for communication with the GAT active site of the TrpG subunit of the heterodimer.

The coordination sphere of the chorismate molecule bound at the active site still remains to be determined directly. However, the crystal structure of the *Serratia* enzyme complexed with a benzoate derivative and pyruvate, products that presumably arose *in situ* either from some aberrant enzymatic conversion of chorismate or from the spontaneous or radiation-induced degradation of bound substrate or product, has been solved and can be used to infer the details of chorismate binding. Accordingly the carboxyl group of the planar ring of chorismate is coordinated with the main-chain nitrogens of Gly328 and Gly485, as well as with the bound metal cofactor. The carboxyl group of the substrate's enolpyruvyl side chain is hydrogen bonded with the hydroxyl of Tyr449, the main-chain nitrogen of Gly483 and the guanidinium nitrogen of Arg469. Besides its interaction with chorismate, the metal cofactor is coordinated by a charged pocket provided by Glu358 and Glu361, with which it directly interacts, and two water molecules held in position by Glu495 and Glu498. The putative catalytic residues, Thr329 and His398, are appropriately located close to the planar ring of chorismate.

The second cavity within the TrpE  $\beta$ -sandwich contains the binding site for the feedback inhibitor, L-tryptophan (13, 14). Tryptophan is positioned in the feedback site by multiple hydrophobic and hydrogen-bonding interactions with residues contributed by both subdomains of the  $\beta$ -sandwich. Its indolyl ring is anchored by hydrophobic interactions with the side chains of Leu38, Glu39, Tyr292, Gly454, and Cys465. The amino and carboxyl groups of tryptophan are hydrogen-bonded with the side chains of residues Ser40, Lys50, and Pro291, as well as with two fixed water molecules; in addition, the nitrogen of the indolyl ring of tryptophan is coordinated by the main-chain nitrogen of Met293. The large distance that separates the active site and the feedback site in the TrpE subunit requires that the competitive inhibition kinetics displayed by the enzyme are a manifestation of alternate conformations that accompany the binding of substrate and inhibitor.

The topology of the TrpG subunit of the anthranilate synthase complex is an a/b-structure that has the fold typical of the triad-type amidotransferase family (12-14). The core of the structure is an open, seven-stranded, mixed  $\beta$ -sheet. The residues of the catalytic triad, Cys84, His170, and Glu172, are appropriately poised for catalysis; Cys84 is strained in a characteristically unfavorable backbone conformation. Glutamate is bound in the active site via a covalent thioester linkage to Cys84 and additional hydrogen-bonding interactions with Ser135, Ser136, Gln89, Gly58, and Leu86 (13). The bound glutamyl moiety is positioned about  $25\text{\AA}$  from the chorismate-binding site on the opposing TrpE subunit of the complex. However, the anticipated hydrophobic tunnel for transporting nascent ammonia from the TrpG active site to the TrpE active site is not apparent in the crystal structure of the enzyme bound with the benzoate derivative and pyruvate (13). It has been suggested that, instead of a closed tunnel, ammonia may have a directed passage down the TrpE crevice to the chorismate-binding site. It is also possible that because of the lability of chorismate, the structure of the fully activated complex has yet to be characterized.

### 3.2. PRA Isomerase and IGP Synthase

In both *E. coli* and *S. enterica* the fourth and fifth reactions on the tryptophan biosynthetic pathway (Fig. 1) are catalyzed by a monomeric bifunctional protein encoded by the *trpC* (*trpC-trpF*) gene (15). Some microorganisms have a monofunctional PRA isomerase, and some have a monofunctional IGP synthase. In the enzyme pair from enteric bacteria, the IGP synthase reaction is

catalyzed by the N-terminal domain ( $M_r = 28,800$ ) and the PRA isomerase reaction by the C-terminal domain ( $M_r = 21,200$ ). The bifunctional protein probably arose during [evolution](#) by in-frame fusion of adjacent *trpC* and *trpF* genes that originally encoded two independently active proteins (see [Gene Fusion](#)). The advantage of gene fusion appears to be the mutual stabilization of the joined proteins against denaturation and degradation.

Both of the PRA isomerase and the IGP synthase domains possess the symmetrical barrel fold, known as the [TIM barrel](#) (Fig. 3) (16). It was observed first in triose phosphate isomerase and is currently the most frequently observed protein domain fold. Whereas PRA isomerase displays the canonical eightfold b/a-barrel (ie, no extensions at the termini and no major insertions in surface loops), IGP synthase has an extension of 49 amino acid residues at its N-terminus. The active sites in both cases are located at the C-terminal ends of the central b-barrels. The two functional domains are connected by a glycine residue between helix a8 of IGP synthase and strand b1 of PRA isomerase. Judged by the independent binding of the substrate/product analog *N*-5'-phosphoribityl anthranilate (rCdRP, CdRP where the keto group at C2' is reduced to an alcohol group) to both PRA isomerase and IGP synthase, as well as by unchanged catalytic activities of the domains after their separation by protein engineering (17), the domains do not communicate reciprocally with each other as observed in the bienzyme complexes of anthranilate synthase (TrpE and TrpG–TrpD) and tryptophan synthase (TrpA and TrpB). Moreover, the product of the PRA isomerase reaction, CdRP, is not channeled to the active site of IGP synthase because the two active sites are both accessible to solvent and point away from one another (Fig. 3). The genetically separated domains of PRA isomerase (TrpF) and IGP synthase (TrpC), as well as of the  $\alpha$ -subunit of tryptophan synthase (TrpA), unfold at equilibrium with increasing concentrations of **denaturants** ([urea](#) or **guanidinium chloride**), via a partially folded intermediate, I, ( $N \leftrightarrow I \leftrightarrow D$ ), where N is the native and D the denatured state of the monomeric protein (18). In each of the three proteins, the intermediate contains the first 6 ba-modules [ba(1–6)] folded, and the last two ba modules [ba(7–8)] are disordered. Adherence to the same “6 + 2” folding transition argues for a common evolutionary origin of these three eightfold b/a-barrel enzymes of tryptophan biosynthesis. PRA isomerase catalyzes the practically irreversible conversion of a ribosylamine, PRA, to an  $\alpha$ -amino ketone, CdRP (Fig. 1). This so-called *Amadori rearrangement* occurs spontaneously, albeit slowly in absence of the enzyme. In contrast, IGP synthase catalyzes a reaction that does not proceed in the absence of the enzyme, the irreversible ring closure of CdRP to IGP and concomitant loss of CO<sub>2</sub> and water.

**Figure 3.** Ribbon diagram of the bifunctional enzyme IGP synthase-PRA isomerase from *E. coli* (16). N- and C-termini are labeled. The filled circle marks the fusion site between helix a8 of IGP synthase and strand b1 of PRA isomerase at residue Gly255. Bound phosphate is shown as a ball-and-stick model, corresponding to the phosphate groups of the substrates (see Figs. 4a and 5a).

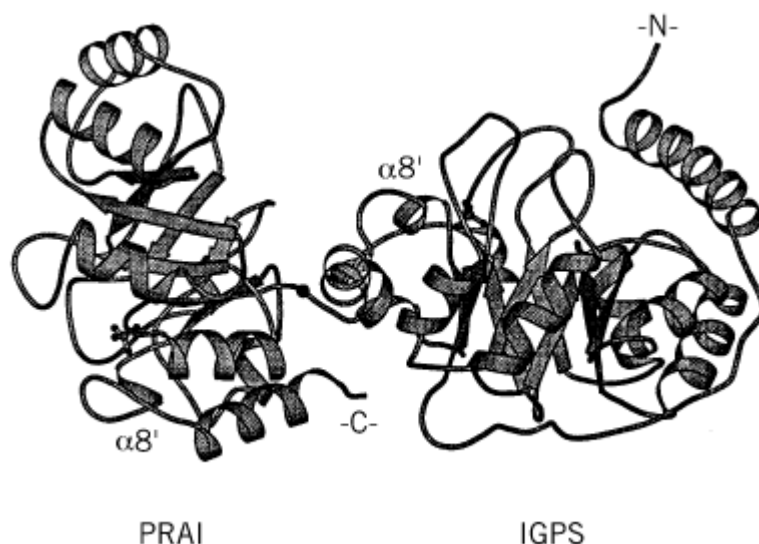
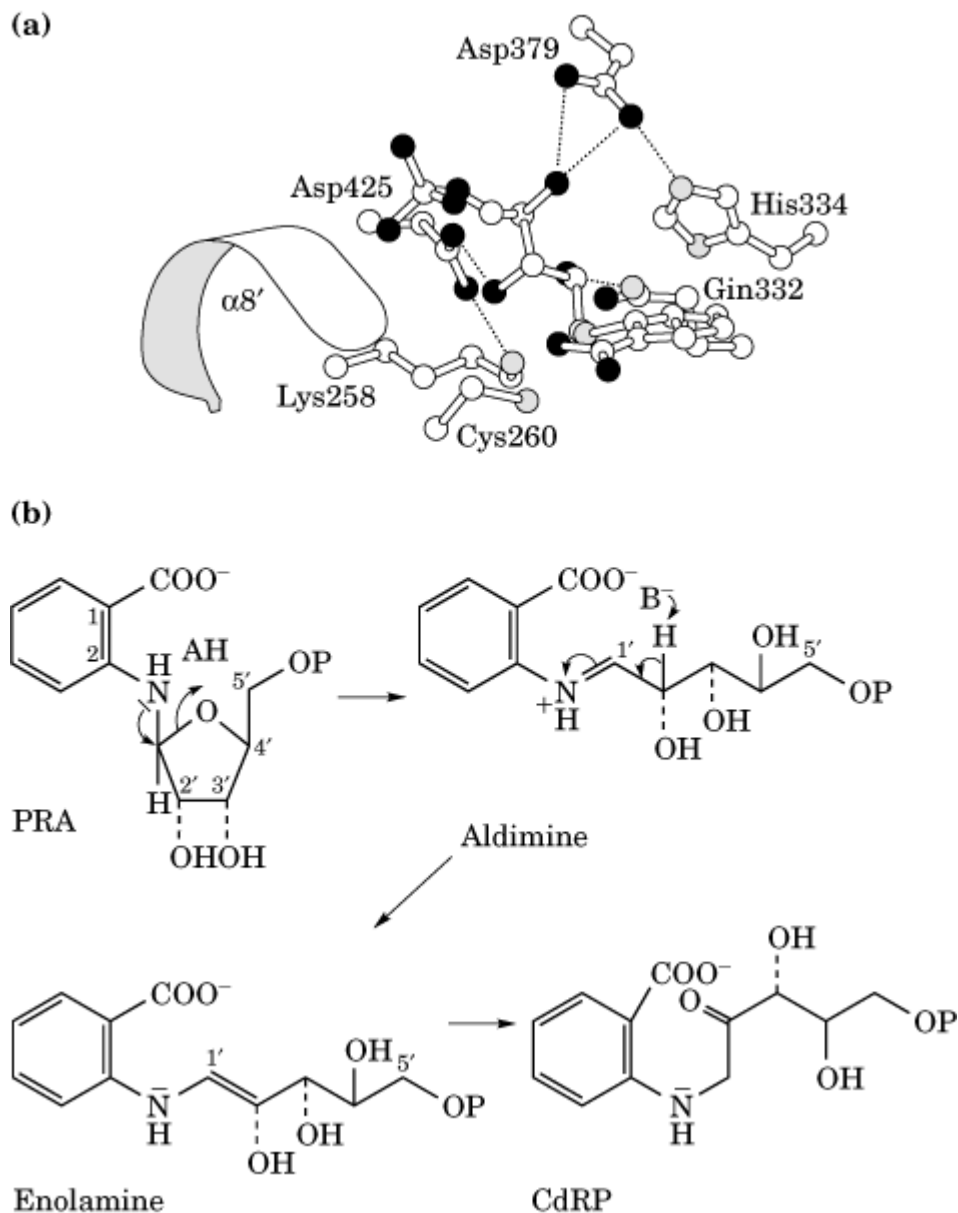


Figure 4a shows the product analog rCdRP fixed at the active site of PRA isomerase (Hennig and Kirschner, unpublished). The phosphate group is bound to the ba(7–8)-subdomain via [salt bridges](#) and [hydrogen bonds](#) and is stabilized further by the dipole of helix  $\alpha 8'$ . This binding mode is also observed in IGP synthase, the  $\alpha$  subunit of tryptophan synthase, and in triosephosphate isomerase, further supporting the notion of a common evolutionary origin for these polypeptides (19). The anthranilate group is bound to a hydrophobic pocket at the opposite end of the active site; the catalytic residues are in between. Among the invariant, polar residues in this region, Cys260 and Asp397 have been identified as catalytically important. Figure 4b indicates how Asp397 as the general acid AH could protonate the furanose ring oxygen and lead to the aldimine. The anion of Cys260 (as  $B^-$ ) could abstract a proton from C2' of the aldimine, leading to the enolamine. The subsequent isomerization to the final product CdRP occurs spontaneously without enzymic catalysis (20).

**Figure 4.** Catalysis by PRA isomerase. **(a)** The active site of PRA isomerase from *Thermotoga maritima* showing invariant residues and the bound product analog rCdRP. Numbering is as for PRA isomerase from *E. coli*.  $\alpha 8'$  is the helix that stabilizes the phosphate group (Fig. 3). **(b)** Tentative intermediates (aldimine and enolamine) of the catalyzed reaction, where Asp397 is the general acid (AH), and Cys260 the general base ( $B^-$ ).



3). Removal of helix  $\alpha 8'$  by deletion mutagenesis shows that it is important for binding the substrate (ground state) but not for stabilizing the [transition state](#) with respect to the bound substrate. Figure [5a](#) presents the product IGP bound at the active site (Hennig and Kirschner, unpublished). The phosphate group is fixed to the ba(7–8)-subdomain, as in the complex of the product analog rCdRP with PRA isomerase (Figure [4a](#)). The indole group is firmly anchored in a hydrophobic pocket at the other end of the active site. Several invariant, polar side chains are located in the region between the phosphate and indole group-binding sites. Among these, Lys114 and Glu163 are essential as catalysis ([19](#)), whereas the invariant residues Glu53 and Asn184 stabilize the active site and position the catalytic residues by a salt bridge and hydrogen bonding, respectively. Figure [5b](#) indicates how bond formation between C1 of the anthranilate group and C2' of the ribulose chain (step 1) could be catalyzed by Lys114 as a general acid (proton donor AH) and by Glu163 as a general base (proton acceptor  $B^-$ ). Arg186 perhaps compensates for the negatively charged carboxylate group of anthranilate, thus stabilizing the sterically strained transition state. The spontaneous decarboxylation of  $I_1$  to  $I_2$  (step 2) is followed by removal of a proton and an hydroxyl ion from  $I_2$  in step 3 and is probably catalyzed by the same general acid–base side chains as in step 1.



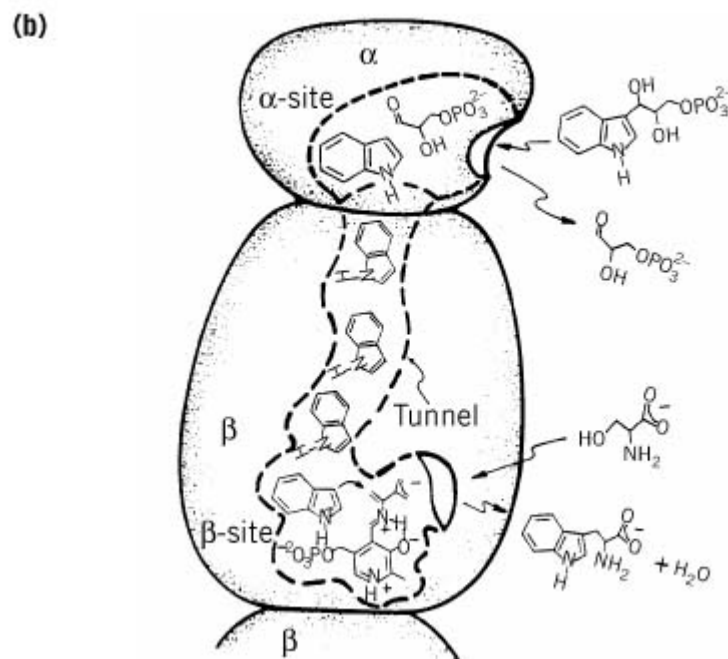
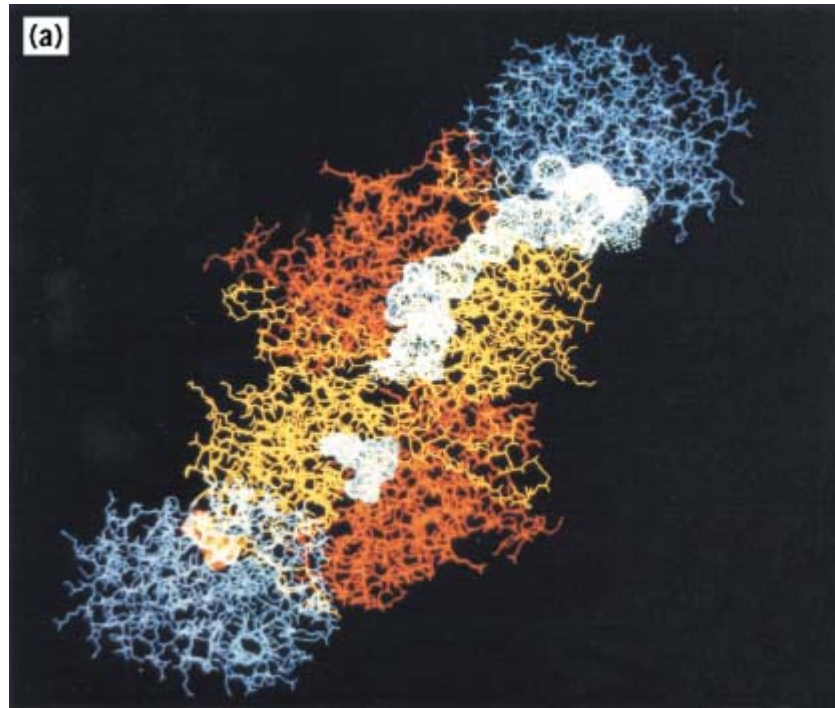


( $M_r = 86,000$ ), hereafter called the b-subunit. The a-subunit, which is encoded by *trpA*, is termed IGP aldolase; it catalyzes the reversible aldolytic cleavage of IGP to indole and D-glyceraldehyde 3-phosphate (G3P). The b-subunit, which is encoded by *trpB*, is termed L-serine hydro-lyase (adding indole); it catalyzes the pyridoxal phosphate-dependent b-elimination of water from L-serine followed by the b-addition of indole to form L-tryptophan. Although each reaction is catalyzed by the individual subunit, the rate is increased by one to two orders of magnitude when the a- and b-subunits associate to form an  $a_2b_2$ -complex. Thus, complex formation results in mutual activation of the a- and b-subunits. Because indole, the intermediate in the overall conversion of IGP to L-tryptophan, is not released free into solution, it must be a **channeled** intermediate that is transferred from one active site to the other within the enzyme complex.

### 3.3.1. Indole Channeling

The physical basis of indole channeling is evident in the [X-ray crystallography](#) structure of the tryptophan synthase  $a_2b_2$ -complex from *S. enterica* (21). The complex has an extended, nearly linear abba arrangement  $\sim 150$  Å long (Fig. 6a). The active sites of neighboring ab-pairs are  $\sim 25$  Å apart and are connected by a buried **hydrophobic** tunnel that has a diameter close to that of indole. This unique tunnel, it is believed, provides a passageway for diffusion of the indole generated at the active site of the a-subunit to the active site of the b subunit, as illustrated by the cartoon in Fig. 6b. Rapid kinetic studies have established that indole is a channeled intermediate [see (22, 23) for reviews].

**Figure 6.** Architecture of the tryptophan synthase  $a_2b_2$ -complex of *Salmonella enterica*. **(a)** The overall structure determined by X-ray crystallography (21). The a-subunits (blue) are on opposite ends of the b subunit dimer. The b-subunit residues 1–204 are in yellow; residues 205–397 are in red. Dot surfaces show the locations of indole-3-propanol phosphate at the active site of the a-subunit and of pyridoxal phosphate at the active site of the b-subunit in the lower ab pair. A tunnel, that connects the two active sites is shown in the upper ab pair. **(b)** Diagram based on the crystal structure showing the reaction of IGP at the active site of the a-subunit, passage of the intermediate indole through the tunnel, and the pyridoxal phosphate-dependent reaction of L-serine and indole at the active site of the b-subunit. See **(a)**.

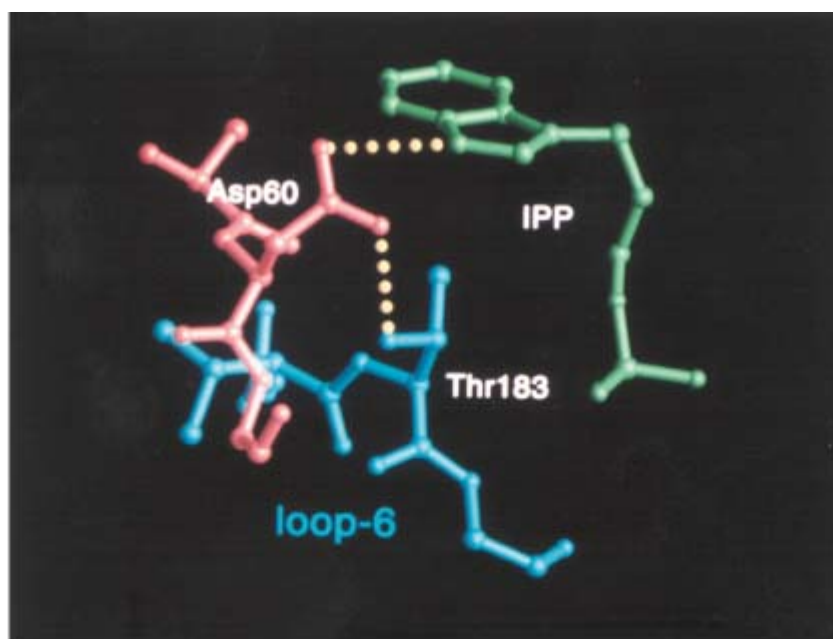


### 3.3.2. Structure and Function of the Tryptophan Synthase $\alpha$ -Subunit

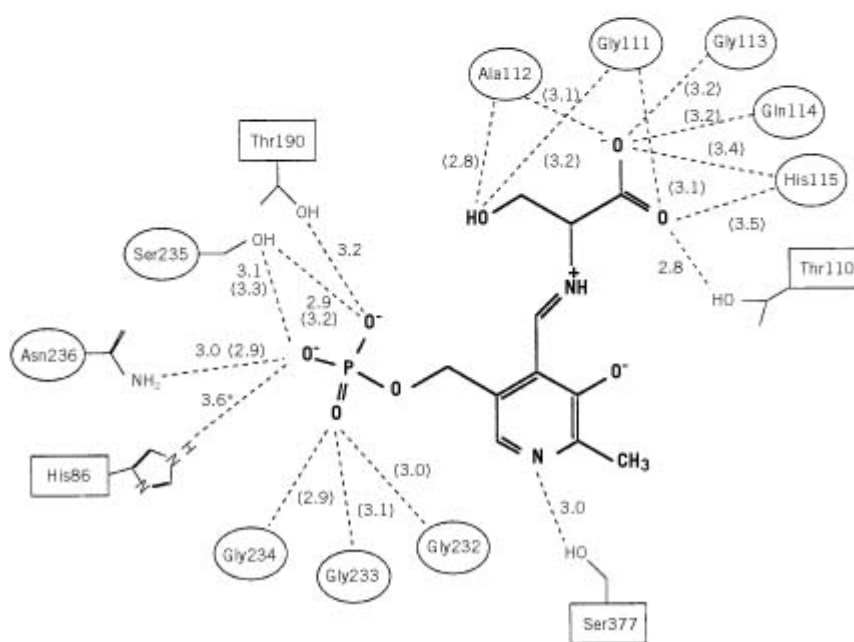
Each  $\alpha$ -subunit in the complex has a  $(\text{ba})_8$ -barrel fold. The structure is similar to that of PRA isomerase and IGP synthase, described earlier. The three enzymes contain a common phosphate group binding site created by residues of the loop between strand-7, helix-7, and the N-terminus of helix 8'. The phosphate of the substrate analog, indole-3-propanol phosphate, binds at this position in the active center of the  $\alpha$ -subunit. Most of the residue positions at which changes inactivate the  $\alpha$ -subunit are located near the bound analog. Mutagenesis studies provide evidence that two of these residues, Asp60 and Glu49, catalyze the aldolytic cleavage of IGP (22) (see [Protein Engineering](#) and [Site-Directed Mutagenesis](#)). Thr183, another position at which there are inactivating changes, cannot be seen in the unliganded structure (21) because it is located in a

flexible loop, loop-6, between strand-6 and helix-6 that has weak electron density. Recent crystal structures provide direct evidence for ligand-induced closure of loop-6 (24). Loop-6 closes down over the bound inhibitor and interacts with residues in Loop-2, including the catalytic residue Asp60. The hydroxyl group of Thr183 interacts with Od of Asp60 and may play a role in positioning the carboxylate group for interaction with the substrate (Fig. 7a).

**Figure 7.** Residues at the active sites of the a-subunit (**a**) and the b-subunit (**b**) from X-ray crystallographic structures of the  $a_2b_2$ -complex that had an inactivating mutation in the b-subunit (24). (**a**) Hydrogen-bonding interactions (yellow dotted lines) between a-subunit residues Thr183 and Asp60 and indole-3-propanol phosphate (IPP). (**b**) Interactions between polar atoms of the external aldimine between pyridoxal phosphate and L-serine and protein residues in an X-ray structure are indicated by dashed lines. The numbers correspond to interatomic distances in angstroms.



(a)



(b)

### 3.3.3. Structure and Function of the Tryptophan Synthase b-Subunit

The larger b-subunit in the complex is composed of two domains of similar size (Fig. 6a). These domains are largely derived from N-terminal residues 1–204 (yellow residues) and C-terminal residues 205–397 (red residues). The pyridoxal phosphate coenzyme is buried in the interface between these two domains at one end of the indole tunnel. The substrate and product binding sites in the b-subunit have been identified in crystal structures of  $\alpha_2\beta_2$ -complexes that contain an inactive, mutant form of the b-subunit in which the active site Lys87 is replaced by threonine (bLys87Thr) (24). The mutant b-subunits bind the substrate (L-serine) or the product (L-tryptophan) as stable external aldimine, Schiff Base intermediates, as illustrated in Fig. 7b for the external aldimine with L-serine. The structures identify residues of the b-subunit that interact with the coenzyme and substrate or product in these intermediates.

Monovalent cations activate the tryptophan synthase b-subunit and  $a_2b_2$ -complex.

Crystallographic studies have identified a binding site for  $\text{Na}^+$  about 8 Å from the phosphate of the bound pyridoxal phosphate. Exchange of  $\text{K}^+$  or  $\text{Cs}^+$  for  $\text{Na}^+$  induces local and long-range changes in the structure of the  $a_2b_2$ -complex [reviewed in (23)]. For example, the side chain of Lys167 forms a salt bridge with a-subunit Asp56 in the presence of  $\text{K}^+$ , but shifts ~7 Å to form a salt bridge with b-subunit Asp305 in the presence of  $\text{Na}^+$ . The interaction between bLys167 and aAsp56 may play a key role in allosteric communication between the a- and b-subunits (see below).

The reaction of L-serine at the active site of the b-subunit proceeds through a series of pyridoxal phosphate intermediates. Pyridoxal phosphate acts as an electron sink to facilitate electron withdrawal from the  $\alpha$ -carbon of the L-serine aldimine shown in Fig. 7b. b-Subunit Lys87 serves a catalytic role as the acceptor of the  $\alpha$ -proton of L-serine in converting of L-serine to aminoacrylate. Indole reacts as a nucleophile on the  $\beta$ -carbon of the aldimine of aminoacrylate as shown in Fig. 6b. Lys87 then protonates the resulting carbanion intermediate and facilitates the release of L-tryptophan.

#### 3.3.4. Reciprocal Communication Between the a- and b-Subunits

Ligands or substrates that bind or react at the active site of either the a-subunit or the b-subunit alter the reaction kinetics of the partner subunit ~25 Å distant. For example, the reaction of L-serine to form aminoacrylate at the active site of the b-subunit (Fig. 6b) results in a 20-fold increase in the rate of turnover of IGP at the active site of the a-subunit. This **allosteric** activation coordinates the activities at the two active sites. The observed kinetic changes result from ligand-induced conformational changes in the protein that are communicated between the active sites of the two subunits. These conformational changes, it is proposed, switch the enzyme between an open, low-activity state, to which ligands bind, and a closed, high-activity state, which prevents the escape of indole (23). Mutagenesis studies provide evidence that a-subunit loops-2 and -6 play important roles in forming of the closed conformation.

Crystal structures of a mutant tryptophan synthase  $a_2b_2$ -complex that has ligands bound to the active sites of the a- and b-subunits provide direct evidence for ligand-induced conformational changes (24). Changes in the a-subunit include closure of loop-6 and the interaction of the closed loop-6 with loop-2 (see before and Fig. 7a). Because closure of loop-6 isolates the active site from solvent, the finding supports the proposal that closure of loop-6 prevents the escape of indole from the enzyme (23). Other conformational differences between the wild type and the two mutant structures include a rigid-body rotation of the a-subunit of ~5 Å relative to the b-subunit and large movements of part of the b-subunit (residues 93–189) toward the rest of the b-subunit (24). The movement of part of the b-subunit creates new interactions between certain residues in the two parts of the b-subunit. The observed conformational changes may contribute to allosteric communication between the a- and b-subunits. Residues at the interface between the a- and b-subunits may also be important in intersubunit communication. Although the subunit interface residues and interactions are not greatly altered in the two mutant structures that have ligands bound at the active sites of the a- and b-subunits, the side chains of two b-subunit residues, Lys167 and Arg175, do exhibit large changes in their interactions with other residues in the two mutant structures. A role for b-subunit Lys167 in allosteric communication between the a- and b-subunits is supported by mutagenesis studies (23). Several residues in the closed a-subunit loop-6 also interact with residues in the b-subunit and may be important for intersubunit communication. Thus, the tryptophan synthase  $a_2b_2$ -complex appears to have evolved so that it can carry out efficient synthesis of L-tryptophan by channeling an indole intermediate between the active sites of the a- and b-subunits. The rates of the reactions are regulated by reciprocal communication between the two subunits.

Recent solution studies provide additional insights into the linkage between the conformation of tryptophan synthase and the regulation of catalysis. Temperature, allosteric ligands, monovalent cations, and mutations affect the equilibrium between a low-activity, “open” conformation and a high-activity, “closed conformation” (25, 26). Mutation of either bLys167 or aAsp56 produces deleterious effects on activity that can be repaired by increased temperatures in combination with CsCl or with NaCl plus an a-subunit ligand (26). The results provide evidence that interaction between aAsp56 and bLys167 is important in stabilizing of the active, closed form of the  $\alpha_2\beta_2$ -complex (26). Spectroscopic and kinetic studies demonstrate that both binding of monovalent cations and formation of the bLys167–aAsp56 salt bridge are important in transmitting allosteric signals between the a- and b-sites (27).

#### 4. Regulation of Gene Transcription

Transcription of the five major structural genes of the *trp* operon of *E. coli* and *S. enterica* is regulated by both repression and transcription attenuation. As previously mentioned, these regulatory mechanisms allow a response, respectively, to changes in the intracellular concentrations of tryptophan and tryptophan-charged tRNA<sup>Trp</sup>, Trp-tRNA<sup>Trp</sup>. The intracellular tryptophan concentration is determined by several events: tryptophan import from the cell's environment, tryptophan produced internally by biosynthesis, and the rate of use of tryptophan during **protein synthesis**. The concentration of charged tRNA<sup>Trp</sup> also depends on several factors: the intracellular concentrations of tryptophan, tRNA<sup>Trp</sup>, and tryptophanyl-tRNA synthetase (see [Aminoacyl tRNA Synthetases](#)), and the overall rate of protein synthesis. It is apparent that the regulatory strategies that control transcription of the *trp* operon of *E. coli* and *S. enterica* were designed to adjust the rate of synthesis of this amino acid in response to all extracellular and intracellular events that alter the availability of tryptophan and Trp-tRNA<sup>Trp</sup> for protein synthesis. **Feedback inhibition** of anthranilate synthase provides an additional important level of regulation by tryptophan by controlling entry of chorismate into the biosynthetic pathway.

Repression of the *trp* operon is sensitive to slight to moderate decreases in intracellular tryptophan concentration. In contrast, attenuation is designed to allow transcription of the structural genes of the operon only when there is nearly complete depletion of the pool of charged tRNA<sup>Trp</sup>. This occurs under two conditions; when the supply of intracellular tryptophan is grossly insufficient and/or when the demand for Trp-tRNA<sup>Trp</sup> greatly exceeds its availability. As we shall describe, the leader peptide coding region, *trpL*, of the *trp* operon of these two bacterial species has only two tryptophan codons, which are adjacent. This design allows relief of termination only when the intracellular concentration of Trp-tRNA<sup>Trp</sup> is so low that rapid translation of both of these tryptophan codons cannot occur.

When *E. coli* or *S. enterica* is grown in the presence of excess tryptophan, repression reduces the rate of transcription initiation at the *trp* promoter to about 1/80th the rate observed in the absence of repression (2). The extent of repression reflects several factors: the intracellular concentrations of tryptophan and repressor, the rate of repressor synthesis, and the fraction of repressor molecules that is available for *trp* operator binding. Repressor synthesis is autoregulated; thus the repressor level drops to one third when cells are grown with excess tryptophan. From the finding that increased production of repressor specified by a **multicopy plasmid** decreases operon expression appreciably, it may be inferred that the *trp* repressor concentration is normally limiting for repression.

Feedback inhibition of anthranilate synthase can have a significant impact on the extent of repression. Thus, when *E. coli* is synthesizing the tryptophan it needs for growth and thus,

repression is incomplete, feedback inhibition by the available tryptophan partially inhibits anthranilate synthase activity. This results in the production of slightly higher levels of the tryptophan biosynthetic enzymes, which is achieved by partial relief of repression. Feedback inhibition is a particularly effective regulatory mechanism because it reduces tryptophan synthesis instantaneously, whereas repression and attenuation have delayed effects. By employing complementary regulatory mechanisms, *E. coli* and *S. enterica* exploit the advantages that can be derived from different means of sensing and responding to tryptophan availability for protein synthesis.

Transcription termination at the *trp* attenuator is relieved completely only when cells are severely deficient in charged tRNA<sup>Trp</sup> (3). Under these conditions, there is about a six- to eight fold increase in the level of *trp* operon structural gene mRNA above the level observed when there is an adequate level of Trp-tRNA<sup>Trp</sup>. Starvation for other amino acids contained in the leader peptide does not have a comparable regulatory effect, except for arginine. Severe arginine starvation also eliminates transcription termination at the *trp* attenuator. The arginine starvation effect is due to the presence of a single arginine codon immediately following the two tryptophan codons in the leader peptide coding segment of the *trp* transcript. Upon arginine starvation, the ribosome that translates the leader peptide coding region presumably also stalls on the critical segment of the *trp* transcript, that is, the segment that contains the tryptophan codons.

#### 4.1. The Operator Region of the *trp* Operon and the Mechanism of Repression

The promoter region of the *trp* operon of *E. coli* and *S. enterica* contains multiple repressor binding sites, or **operators** (Fig. 8). The tryptophan-activated *trp* repressor can bind at these sites and inhibit transcription initiation at the major *trp* **promoter**. Three-dimensional structures have been determined by X-ray crystallography and by **NMR** of the tryptophan-free *trp* aporepressor, the tryptophan-activated *trp* repressor, and the repressor-operator complex. Both the aporepressor and repressor exist as dimers composed of interlocking identical polypeptide chains. Each *trp* aporepressor dimer has two tryptophan binding sites. Bound tryptophan displaces helix E of the two **helix-turn-helix** domains of the dimer and positions these domains so that the dimer can effectively bind the symmetrical operator. Binding of tryptophan at the two sites of the aporepressor is cooperative. The structure of the repressor-operator complex shows that the repressor makes multiple contacts with target operator sequences and that some important specific contacts are water-mediated. This mechanism of recognition has been termed *indirect readout* (28). Although crystallographic analyses suggest that the preferred operator binding sequence is 5'ACTAGT3' (28), other studies indicate that the consensus operator sequence is 5'GNACT3' and that there are three operators in the *trp* promoter/operator region (Fig. 8) (29). The presence of multiple bound repressor dimers in the operator region could facilitate inhibition of transcription initiation of the operon by reducing the likelihood that the promoter/operator will be repressor-free. Repressor dimers, it has been observed, associate in solution, but the significance of this association is not yet known. The *trp* repressor also regulates initiation of the transcription of several other operons concerned with tryptophan metabolism.

**Figure 8.** The *trp* operon promoter/operator region. The -35, -10, and +1 sequences of the promoter are marked. The sequences in the two DNA strands that are assumed to be recognized specifically by each *trp* repressor dimer are indicated in outline form. Note that the diagram implies that three repressor dimers can be bound per operator region. Based on the studies of Yang et al. (29).



#### 4.2. The Mechanism of Transcription Attenuation

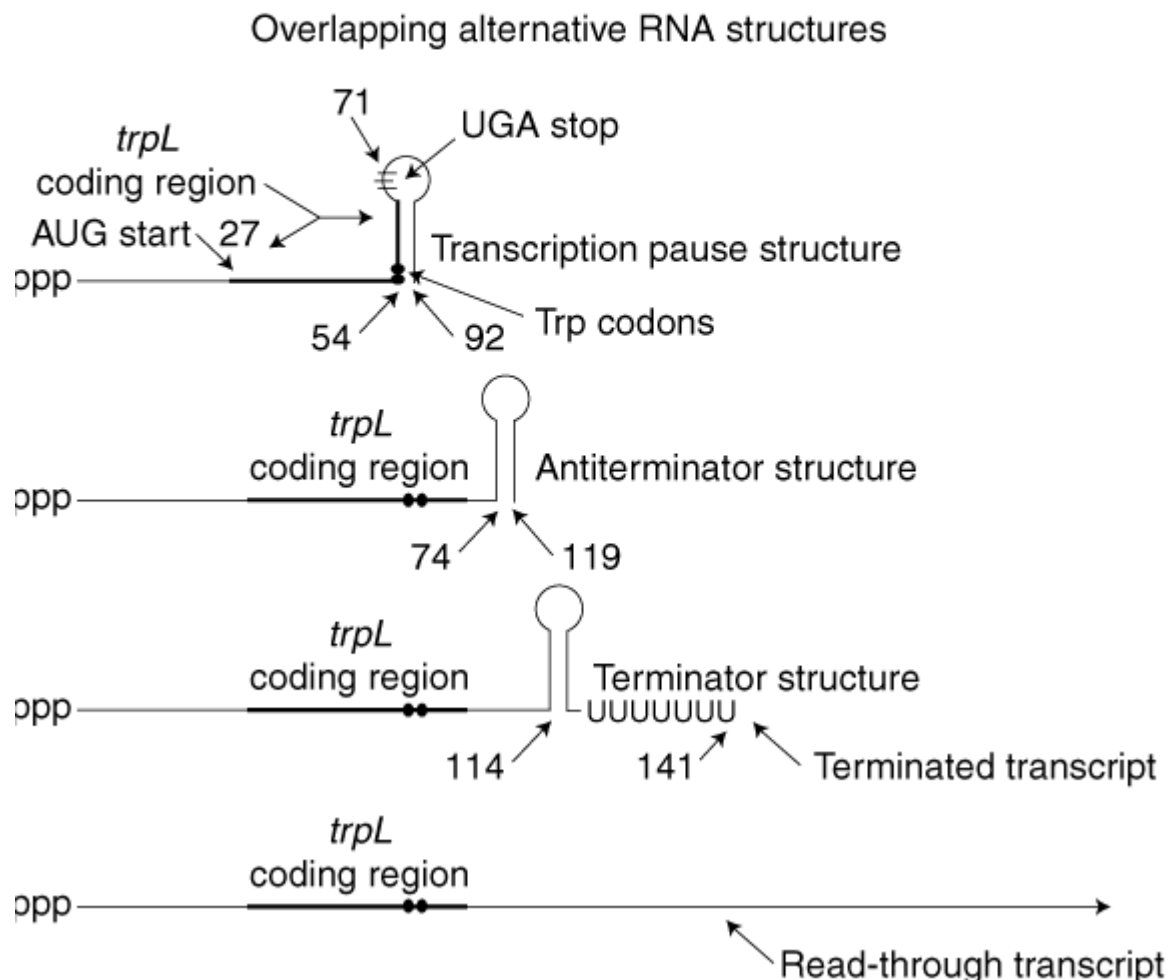
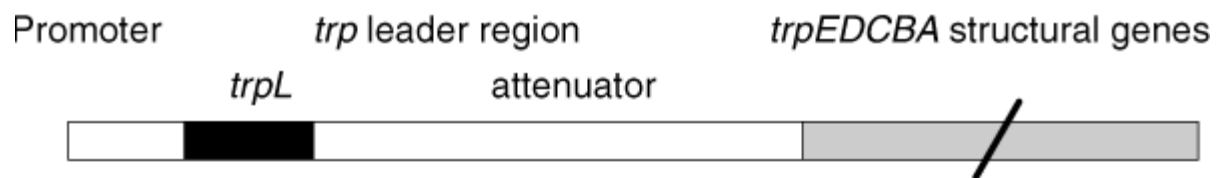
The principal *trp* promoter, *trpP1*, is a relatively strong promoter; in tryptophan-deficient cells, **RNA polymerase** can initiate transcription at this promoter every 4–5 seconds (3). Inasmuch as a transcribing polymerase molecule could reach the attenuator region within a few seconds, it is essential that the decision be made quickly whether or not to terminate transcription at the attenuator. Because this decision is based on the transcript position of the ribosome that is attempting synthesis of the leader peptide, it is essential to synchronize [translation](#) of the leader peptide coding region with transcription of the leader region. Synchronization is achieved by features of the nucleotide sequence of the leader transcript; other features of the transcript determine the appropriate transcription termination response, which depends on the cellular level of charged tRNA<sup>Trp</sup> (3).

#### 4.3. Transcription Pausing, the Initial Important Event in Transcription Attenuation

After initiation of transcription at the *trp* promoter, RNA polymerase synthesizes an RNA segment that forms a **hairpin** structure that causes the transcribing polymerase to pause in the leader region of the operon (Fig. 9). The paused polymerase probably remains at the pause site for less than 1 second (3). The sequence of the pause RNA hairpin, particularly its upper base-paired half, it is believed, is the predominant pause signal (3). However, the DNA region immediately downstream from the pause RNA-encoding region also contributes to pause stability, as does the RNA polymerase accessory factor, NusA. In the paused complex, the 5' end of the nascent transcript is exposed and is available for ribosome loading and initiation of leader peptide synthesis. In fact it is the act of synthesis of the initial portion of the leader peptide that disrupts the transcription pause complex and reactivates transcript elongation. Thus, transcription pausing is the event that is essential to regulation by this form of transcription attenuation.

**Figure 9.** Organization and features of the *trp* leader region and the *trp* leader transcript. The operon map at the top expands the *trp* leader region and shows the locations of *trpL*, the leader peptide coding region, the attenuator, and the five structural genes. Below the operon map are transcripts that have the three alternative hairpin structures, the pause, antiterminator, and terminator. The terminated and the read-through transcripts are marked. Part of the leader peptide coding region, *trpL*, overlaps the pause structure, as indicated. The start and stop codons of *trpL* and the positions of its adjacent Trp codons are shown. Nucleotide positions that define the various transcript segments are indicated.





#### 4.4. Transcription Beyond the Pause Site

When the paused RNA polymerase resumes transcription, the position assumed by the ribosome engaged in synthesizing the leader peptide dictates whether the RNA antiterminator or terminator structure will form (Fig. 9). If a cell has adequate levels of Trp-tRNA<sup>Trp</sup> (and other required charged tRNAs), the translating ribosome will reach the leader peptide stop codon while transcription of the leader region of the operon is proceeding. A ribosome at this position would block formation of the antiterminator structure, thereby allowing what would be the 3' strand of the RNA antiterminator to pair with the immediately following RNA segment to form the terminator (Fig. 9). Formation of the terminator within the transcription complex would signal the transcribing polymerase to terminate transcription. If the Trp-tRNA<sup>Trp</sup> concentration was very low, the translating ribosome would stall during attempted translation of the adjacent tryptophan codons. Such stalling would promote formation of the RNA antiterminator structure. The antiterminator would prevent the subsequent formation of the terminator because the 3' strand of the antiterminator must be free for it to serve as the 5' strand of the terminator (Fig. 9). Thus the crucial location of the ribosome engaged in synthesizing the leader peptide dictates

which RNA hairpin structure forms, and this, in turn, determines whether termination or read-through will occur.

#### 4.5. Basal Level Control

When cells have adequate levels of Trp-tRNA<sup>Trp</sup> to satisfy their needs for protein synthesis and synthesis of the *trp* leader peptide can be completed, only about 90% of *trp* transcripts are terminated in the leader region, there is about 10% read-through transcription. These read-through transcripts provide a “basal level” of all five *trp* polypeptides. Read-through is principally due to ribosome dissociation from the leader peptide stop codon when the synthesis of the leader peptide is completed. Ribosome release allows the leader segment of the transcript to form an antiterminator structure occasionally, which results in read-through. This read-through depends on the ability of cells to initiate and complete synthesis of the leader peptide. When initiation of synthesis of the *trp* leader peptide is defective, transcription termination at the *trp* attenuator increases from about 90% to 98%. This increased termination is due to increased formation of the terminator and decreased formation of the antiterminator. The mechanism responsible for this change is termed basal level control. Although basal level control is a natural consequence of the sequence and arrangement of the *trp* transcript, it is conceivable that this regulatory feature was introduced to reduce *trp* operon transcription further when a severe amino acid deficiency interfered with general protein synthesis.

#### 4.6. Cessation of Leader Peptide Synthesis

Whether transcription stops at the attenuator or continues into the structural genes of the operon, subsequent translation of the leader peptide coding region would no longer be necessary. In fact, after the attenuation decision has been made, leader peptide synthesis is shut down; both the terminated and read-through transcripts contain leader nucleotide sequences that can pair with the ribosome binding site of the leader peptide coding region, and block further peptide synthesis (3).

### 5. Feedback Inhibition

Regulation of carbon flow in the tryptophan pathway is mediated by classical negative **feedback inhibition** of the various activities of the anthranilate synthase-APR transferase complex by L-tryptophan. Feedback inhibition results from the binding of one molecule of tryptophan to the feedback site of the TrpE subunit of the complex. Complete inhibition of anthranilate synthase activity can be attained with the intact TrpE-TrpG-TrpD complex, the TrpE-TrpG partial complex, and the uncomplexed TrpE subunit. This is also true for the individual ADIC synthase, ADIC lyase, and glutaminase activities of the complexes (Fig. 2). In contrast, inhibition of APR transferase activity is partial and never exceeds a maximum of about 60%. Tryptophan inhibition of APR transferase occurs only when TrpG-TrpD is complexed with TrpE, consistent with the location of the tryptophan binding site on the TrpE subunit. In all cases, inhibition is competitive with respect to chorismate and noncompetitive with respect to the other substrates, indicating allosteric communication between the inhibitor site and the active sites of the TrpE and TrpG-TrpD subunits.

Significant conformational effects accompany tryptophan binding to the complex. These are manifest in the cooperative kinetics of tryptophan binding, in the tryptophan-induced cooperative kinetics for chorismate utilization, and in the tryptophan-induced alterations in the chromatographic behavior of the enzyme. The global nature of the conformational effects of tryptophan binding are perhaps best demonstrated in the properties of a doubly mutant hybrid complex, assembled *in vitro*, that contains one catalytically active, feedback-insensitive TrpE subunit, and one catalytically inactive, feedback-sensitive TrpE subunit (30). The feedback sensitivity of the anthranilate synthase activity of the hybrid complex demonstrates that the binding of a single tryptophan molecule to one TrpE subunit elicits conformational and kinetic

effects that are propagated to the active site of the unliganded companion TrpE subunit.

Comparison of the three-dimensional structures of the active and the inhibited enzyme forms suggests a likely molecular mechanism of feedback inhibition in the AS-APR complex (13). In the tryptophan-bound enzyme, a structural motif of the TrpE subunit, formed by residues 327–363 and 387–403, is shifted about 7Å away from the chorismate-binding pocket relative to its position in the chorismate-bound enzyme. Significantly, this movement displaces catalytic residues Thr329 and His398 as well as two metal-binding glutamate residues from the vicinity of the active site and results in the loss of catalytic activity. This rearrangement is the result of the ordering of a disordered loop made up of residues 42–49 in response to tryptophan binding (13, 14). In contrast, the binding of chorismate at the active site moves this same motif in the opposite direction toward the active site crevice and leads to activation of the enzyme. These reciprocal effects explain the competitive behavior of chorismate and tryptophan because binding of one of these molecules to TrpE necessarily precludes binding of the other. Also, the fact that the structural loops involved are situated at the heterotetramer interface explains how a single tryptophan molecule can inhibit both TrpE subunits; in view of the molecular twofold axis of the complex, a movement in one heterodimer must be accompanied by a corresponding movement in the other.

Binding of tryptophan by the TrpE subunit also induces relevant conformational changes in the TrpG subunit of the complex. Two  $\beta$ -strands of TrpG, made up of residues 103–109 and 130–140, become disordered in the tryptophan-bound enzyme (13) in response to the tryptophan-induced movement of an interacting TrpE  $\alpha$ -helix at the heterodimer interface. These changes disrupt the glutamine binding site of the TrpG subunit. Thus, a unitary mechanism for both tryptophan inhibition and chorismate activation of the glutaminase activity of the complexed TrpG subunit is manifest. The mechanism of the partial tryptophan inhibition of the APR activity of the complexed TrpG-APR subunit remains unknown.

Clearly, the unique allosteric properties of the anthranilate synthase-APR transferase complex, coupled with its critical location at the beginning of the pathway, provide a rapid and efficient mechanism of metabolic regulation that modulates and complements the elegant regulatory circuitry at work at the transcriptional and translational levels.

## 6. Other Features That Influence *trp* Operon Expression

### 6.1. Internal Promoter *trpP2*

A secondary, internal promoter, *trpP2*, exists near the 3' end of the *trpD* gene of *E. coli* and *S. enterica* (Fig. 1) (31). *trpP2* is a low-efficiency, **constitutive** promoter, that initiates transcription about 150 nucleotides upstream of the *trpD* stop codon (32). It is responsible for about 60% of the fully repressed expression of the *trpC*, *trpD*, and *trpA* genes. The *trpP2* sequence lacks a conserved nucleotide in both the -35 and -10 regions of the canonical promoter sequence, accounting for its poor efficiency. The *in vivo* role of *trpP2*, it is thought, is to aid the cell in responding to an environment that has extreme variations in the supply of exogenous tryptophan. The elevated basal level of the three-terminal pathway enzymes resulting from *trpP2* function might provide an advantage to the cell in adapting rapidly to conditions of tryptophan starvation. The conservation of this secondary transcriptional element in numerous enteric species indicates that it is physiologically significant, not merely a vestige of an earlier regulatory structure.

### 6.2. Translational Coupling

The multigene *trp* operon of *E. coli* and *S. enterica* is designed to permit regulated synthesis of a polycistronic mRNA that encodes the five polypeptides needed for tryptophan formation. The required concentration of each polypeptide would depend on its inherent catalytic activity, its affinity for its substrates, and its stability. Other potential variables that could influence the rate

of tryptophan synthesis are the efficiency of translation of each mRNA coding segment and the mRNA segments' functional half-lives. As mentioned, two pairs of *trp* polypeptides, TrpE and TrpD, and TrpB and TrpA, form enzyme complexes in which separate catalytic events are coordinated. To optimize catalysis, equimolar levels of each member of these complexes are achieved by translational coupling (2). Thus, the *trpD* and *trpA* ribosome binding sites are designed so they are inefficiently used unless translation of *trpE* and *trpB*, respectively, proceeds to the end of its coding region. Coupling ensures equimolar polypeptide synthesis.

### 6.3. Preferential Synthesis upon Tryptophan Starvation

TrpE and TrpA of *E. coli* and *S. enterica* lack tryptophan. A consequence of this deficiency is that upon severe tryptophan starvation—when synthesis of most tryptophan-containing proteins is reduced or prevented—synthesis of TrpE and TrpA continues, and the relative levels of these proteins increase (2). The disproportionate elevated level of TrpE would result in preferential sequestration of chorismate into the tryptophan pathway. This would occur despite the lack of a comparable increase in the formation of TrpG–D, which normally contributes the glutamine amidotransferase activity used for anthranilate formation, because anthranilate formation by anthranilate synthase does not absolutely depend on glutamine amidotransferase activity. The TrpE polypeptide acting alone can convert chorismate to anthranilate, using free ammonia as the amino source. Inasmuch as anthranilate synthesis proceeds by an irreversible reaction, the conversion of chorismate to anthranilate *in vivo* would lock chorismate into the tryptophan biosynthetic pathway.

### 6.4. *trp* mRNA Degradation

The half-life of *trp*mRNA is comparable to that of many *E. coli* mRNAs (see [Messenger RNA](#)). There is a modest gradient of decay along its length, *trpEm*mRNA is about half as stable as *trpAm*mRNA (2). As in other mRNAs, **hairpin** structures at the 5' and 3' ends of *trp*mRNA probably contribute to its stability.

### 6.5. *trp* Protein Turnover

During log phase growth in a minimal or rich medium, the five *trp* biosynthetic proteins of *E. coli* are relatively stable. Only upon prolonged starvation is there significant turnover of one of the *trp* proteins, the product of *trpC*. Apparently, *E. coli* adjusts the levels of the *trp* biosynthetic enzymes predominantly by regulating synthesis and by growth-dependent dilution.

## 7. Evolutionary Considerations

Although the same seven polypeptide domains appear to be responsible for tryptophan biosynthesis in all organisms that can synthesize this amino acid (33), in many organisms these domains are fused in different arrangements than in *E. coli* and *S. enterica* (15). In addition, chromosomal order and groupings of *trp* genes, as well as the mechanisms used to regulate *trp* gene/operon expression, vary in different species (15, 34). Providing explanations for these differences will require information that is difficult to obtain, such as understanding the functional demands that were imposed on each species during its evolution. There are species in which as many as five of the *trp* enzymatic functions reside in a single multifunctional polypeptide (35). In some species, novel metabolic reactions occur that use tryptophan pathway intermediates for other purposes; understanding these accessory reactions may provide explanations for the different gene arrangements and regulatory mechanisms that are observed (34).

The evolutionary origin of each of the *trp* polypeptide domains is also of considerable interest. The 3-D structures of six of the seven *trp* enzymatic domains are known now. As mentioned, three of these domains, TrpA, TrpF, and TrpC, form barrel structures that have similar characteristics; thus each could have evolved from one of the others or from other preexisting members of this family. In fact, the IGP synthase domain of *E. coli* has been mutationally altered

to yield a polypeptide that can catalyze the PRA isomerase reaction (36). In addition, the HisA barrel protein of *Thermotoga maritima*, an enzyme that catalyzes an Amadori rearrangement of a phosphoribosylamine similar to the PRA isomerase reaction (Fig. 4), has been mutationally converted into an enzyme that catalyzes the PRA isomerase reaction (37). These findings illustrate how simple it would be for a protein to evolve with a new catalytic activity, starting with a copy of a gene encoding an enzyme that could catalyze a similar reaction. Alternatively, both enzymes may have evolved from an ancestral enzyme that had broad substrate specificity (37). The TrpB structure is homologous to proteins in a family of pyridoxal phosphate-dependent enzymes, many of which, like TrpB, can catalyze an amino acid dehydratase reaction. The TrpG domain is highly homologous to members of the family of Triad glutamine amidotransferases. A close homolog of TrpE participates in synthesizing an analog of anthranilate, *p*-aminobenzoate (15). Identifying the true origins of the *trp* genes and tryptophan proteins probably will be difficult because there are numerous ancestral possibilities, and the true evolutionary intermediates may no longer exist.

## 8. Dedication

We dedicate this article to the fond memory of Irving P. Crawford, who also loved the genes and enzymes of tryptophan biosynthesis.

## Bibliography

1. C. Yanofsky, T. Platt, I. P. Crawford, B. P. Nichols, G. E. Christie, H. Horowitz, M. van Cleemput and A. M. Wu (1981) *Nucl. Acids. Res.* **9**, 6647–6668.
2. C. Yanofsky and I. P. Crawford (1989) In *Cellular and Molecular Biology*, Vol. **2** (F. C. Neidhardt, J. L. Ingraham, K. B. Low, B. Magasanik, M. Schaecter, and H. E. Umbarger, eds.), ASM Press, Washington, DC, pp.1453–1472.
3. R. Landick, C. L. Turnbough, Jr., and C. Yanofsky (1996) In *Cellular and Molecular Biology*, Vol. **1** (F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaecter, and H. E. Umbarger, eds.) ASM Press, Washington, DC., pp. 1263–1286.
4. H. Zalkin (1973) *Adv. Enzymol. Relat. Areas Mol. Biol.* **38**, 1–39.
5. R. Bauerle, J. Hess, and S. French (1987) *Methods Enzymol.* **142**, 366–386.
6. H. Zalkin (1993) *Adv. Enzymol. Relat. Areas Mol. Biol.* **66**, 203–309.
7. A. Morollo and R. Bauerle (1993) *Proc. Natl. Acad. Sci. USA* **90**, 9983–9987.
8. A. Summerfield, R. Bauerle, and C. M. Grisham (1988) *J. Biol. Chem.* **263**, 18793–18801.
9. M. G. Caligiuri and R. Bauerle (1991) *J. Biol. Chem.* **266**, 8328–8335.
10. G. F. Miozzari and C. Yanofsky (1979) *Nature* **277**, 486–489.
11. H. Zalkin and J. L. Smith (1998) *Adv. Enzymol. Relat. Areas Mol. Biol.* **72**, 87–144.
12. T. Knochel, A. Ivens, G. Hester, A. Gonzalez, R. Bauerle, M. Wilmanns, K. Kirschner, and J. N. Jansonius (1999) *Proc. Natl. Acad. Sci. USA* **96**, 9479–9484.
13. G. Spraggon, C. Kim, X. Nguyen-Huu, M-C. Yee, C. Yanofsky, and S. Mills (2001) *Proc. Natl. Acad. Sci. USA* **98**, 6021–6026.
14. A. A. Morollo and M. J. Eck (2001) *Nature Struct. Biol.* **8**, 243–247.
15. B. P. Nichols (1996) In *Cellular and Molecular Biology*, Vol. **1** (F. C. Neidhardt, R. Curtiss III, J. L. Ingraham, E. C. C. Lin, K. B. Low, B. Magasanik, W. S. Reznikoff, M. Riley, M. Schaecter, and H. E. Umbarger, eds.), ASM Press, Washington, DC, pp. 2638–2648.
16. M. Wilmanns, J. P. Priestle, T. Niermann, and J. N. Jansonius (1992) *J. Mol. Biol.* **223**, 477–507.
17. M. Eberhard, M. Tsai-Pflugfelder, K. Bolewska, U. Hommel, and K. Kirschner (1995)

- Biochemistry **34**, 5419–5428.
18. M. M. Sanchez del Pino and A. R. Fersht (1997) *Biochemistry* **36**, 5560–5563.
  19. B. Darimont, C. Stehlin, H. Szadkowski, and K. Kirschner (1998) *Protein Sci.* **7**, 1221–1232.
  20. U. Hommel, M. Eberhard, and K. Kirschner (1995) *Biochemistry* **34**, 5429–5439.
  21. C. C. Hyde, S. A. Ahmed, E. A. Padlan, E. W. Miles, and D. R. Davies (1988) *J. Biol. Chem.* **263**, 17857–17871.
  22. E. W. Miles (1995) In *Subcellular Biochemistry, Proteins: Structure, Function, and Protein Engineering* (B. B. Biswas and S. Roy, eds.), Plenum Press, New York, **24**, pp. 207–254.
  23. P. Pan, E. Woehl, and M. F. Dunn (1997) *Trends Biochem. Sci.* **22**, 22–27.
  24. S. Rhee, K. D. Parris, C. C. Hyde, S. A. Ahmed, E. W. Miles, and D. R. Davies (1997) *Biochemistry* **36**, 7664–7680.
  25. Y. X. Fan, P. McPhie, and E. W. Miles (2000) *Biochemistry* **39**, 4692–4703.
  26. Y. X. Fan, P. McPhie, and E. W. Miles (2000) *J. Biol. Chem.* **275**, 20302–20307.
  27. E. Weber-Ban, O. Hur, C. Bagwell, U. Banik, L.-H. Yang, E. W. Miles, and M. F. Dunn (2001) *Biochemistry* **40**, 3497–3511.
  28. Z. Shakked, G. Guzikovich-Guerstein, F. Frolow, D. Rabinovich, A. Joachimiak, and P. B. Sigler (1994) *Nature* **368**, 469–473.
  29. J. Yang, A. Gunasekera, T. A. Lavoie, L. Jin, D. E. A. Lewis, and J. Carey (1996) *J. Mol. Biol.* **258**, 37–52.
  30. M. Caliguri and R. Bauerle (1991) *Science* **252**, 1845–1888.
  31. R. Bauerle and P. Margolin (1967) *J. Mol. Biol.* **26**, 423–436.
  32. H. Horowitz, J. VanArsdell, and T. Platt (1983) *J. Mol. Biol.* **169**, 775–797.
  33. I. P. Crawford (1989) *Annu. Rev. Microbiol.* **43**, 567–600.
  34. M. Chang, D. E. Essar, and I. P. Crawford (1990) In *Biotransformations, Pathogenesis, and Evolving Biotechnology* (S. Silver, A. M. Chakarabarty, B. Iglewski, and S. Kaplan eds.), ASM Press, Washington, DC, pp. 292–302.
  35. T. Schwarz, K. Uthoff, C. Klinger, H. E. Meyer, P. Bartholmes, and M. Kaufmann (1997) *J. Biol. Chem.* **272**, 10616–10623.
  36. M. M. Altamirano, J. M. Blackburn, C. Aguayo, and A. Fersht (2000) *Nature* **403**, 617–622.
  37. C. Jurgens, A. Strom, D. Wegener, S. Hettwer, M. Wilmanns, and R. Sterner (2000) *Proc. Natl. Acad. Sci. USA* **97**, 9925–9930.

### **Suggestions for Further Reading**

38. I. P. Crawford (1987) Synthesis of tryptophan from chorismate: comparative aspects. *Methods Enzymol.* **142**, 293–301.
39. C. Yanofsky (1984) Comparison of regulatory and structural regions of genes of tryptophan metabolism. *Mol. Biol. Evol.* **1**, 143–161.
40. C. Yanofsky (1987) Tryptophan synthetase: Its charmed history. *Bioessays* **6**, 133–137.
41. E. W. Miles (2001) Tryptophan synthase: A multienzyme complex with an intramolecular tunnel. *Chem. Rec.* **1**, 140–151.

## **TRP Repressor**

For many helix-turn-helix motifs the selectivity of the protein for a particular DNA sequence depends both on direct contacts between the protein and the bases and also on local conformational constraints on the DNA structure. However in the case of one hth protein, the tryptophan repressor, substantial selectivity is achieved in the absence of such direct contacts. In the presence of the ligand L-tryptophan, this [repressor](#) binds to an **operator** site within the promoter region of the *trp* operon and blocks transcription initiation. This operon encodes the enzymes required for the biosynthesis of tryptophan which on binding to the trp repressor completes a negative feedback loop.

The *trp* repressor has been crystallized as a complex with a symmetrized 18 bp operator, TGTACTAGTTAAGTAGTAC. In this sequence, the base pairs whose function *in vivo* is most sensitive to [mutation](#) are the two ACTAG sequences centered 4–5 bp from the dyad. In the cocrystal there are no direct hydrogen bonded or nonpolar contacts between the helix-turn-helix motif and the DNA that can account for the selectivity of the interaction (The one direct contact between a base and an amino acid side-chain lies outside the mutationally sensitive region). Instead, direct hydrogen bonded contacts are made principally to the phosphate groups in the DNA backbone. These contacts serve to constrain the backbone in a defined configuration. However, three well-ordered and directed water molecules lying in the major groove between the helix-turn-helix motif and the CTAG sequence make hydrogen bonded contacts with the base pairs in this sequence. Notably recognition of the adenine residue immediately preceding the CTAG sequence by a water-mediated hydrogen bond is sequence specific. Mutations in the amino acids that bind to these water molecules result both in a reduced affinity, and in some examples an altered sequence selectivity, of the repressor for the operator. This structure thus suggests that a degree of sequence selectivity could be achieved by water molecules forming an essential component of the recognition surface of the protein and thus mediating direct readout of the operator sequence acting, as it were, by proxy for the protein.

## **True Breeding**

A homozygote is a true breeding individual. A group of individuals with the same phenotype are said to breed true when their phenotype is maintained for any number of generations in the offspring of genetic **crosses** within the group. The group may consist of a single individual capable of self-fertilization. True breeding occurs when all members of the group are homozygous for the same **alleles**. A trait is said to breed true when the members of the group are homozygous for the same alleles of the genes responsible for the trait. A group of true-breeding individuals is called a *pure line*.

## **Trypsin Inhibitors**

Moses Kunitz characterized, purified, and crystallized two very different trypsin inhibitors: bovine pancreatic trypsin inhibitor (Kunitz) and soybean trypsin inhibitor (Kunitz). He also showed that each of these inhibitors formed a strong but reversible stoichiometric complex with trypsin. The

complexes were totally devoid of trypsin activity. For several years after Kunitz's second achievement, only two well-characterized protein inhibitors were known. This gave rise to the incorrect surmise that most protein inhibitors are trypsin inhibitors. Many laboratories tested numerous biological materials for the inhibition of trypsin and frequently found it. This allowed for the isolation of many more trypsin inhibitors, largely reinforcing the mistaken belief about them. The standard mechanism of interaction between serine proteinases and their inhibitors was established by Laskowski by studying the interaction of trypsin with soybean trypsin inhibitor (Kunitz). It was shown that soybean trypsin inhibitor (Kunitz) is a single polypeptide chain of 181 amino acid residues, crosslinked by two (39–86, 136–145) intramolecular disulfide bridges. Within the first disulfide loop lies the Arg<sup>63</sup>-Ile peptide bond that serves the inhibitor as the reactive site for interaction with serine proteinases. In bovine pancreatic trypsin inhibitor (Kunitz), the reactive site peptide bond is Lys<sup>15</sup>-Ala. Thus, there are two types of trypsin inhibitors—arginine and lysine. The latter but not the former lose their activity by acetylation. The replacement of the reactive site Arg or Lys by Trp by enzymatic semisynthesis converted soybean trypsin inhibitor (Kunitz) from a strong trypsin to a strong chymotrypsin inhibitor. Similar experiments are now done by site-specific mutagenesis, and many conversions of a trypsin inhibitor to a chymotrypsin inhibitor, elastase inhibitor, glutamic acid-specific proteinase inhibitor, or furin inhibitor have been reported. More importantly, workers who searched for natural inhibitors for the above enzymes generally found them. Often, they are homologous to known trypsin inhibitors. The differences are primarily in the enzyme inhibitor contact region.

#### Suggestions for Further Reading

M. Laskowski Sr. and M. Laskowski Jr. (1954) Naturally occurring trypsin inhibitors. *Adv. Protein Chem.* **9**, 203–242.

W. Lu et al. (1993) Arg<sup>15</sup>-Lys<sup>17</sup>-Arg<sup>18</sup> turkey ovomucoid third domain inhibits human furin. *J. Biol. Chem.* **268**, 14583–14585.

K. Ozawa and M. Laskowski Jr. (1966) The reactive site of trypsin inhibitors. *J. Biol. Chem.* **241**, 3955–3961.

R. W. Sealock and M. Laskowski Jr. (1969) Enzymatic replacement of the Arginyl by a Lysyl residue in the reactive site of soybean trypsin inhibitor. *Biochemistry* **8**, 3703–3710.

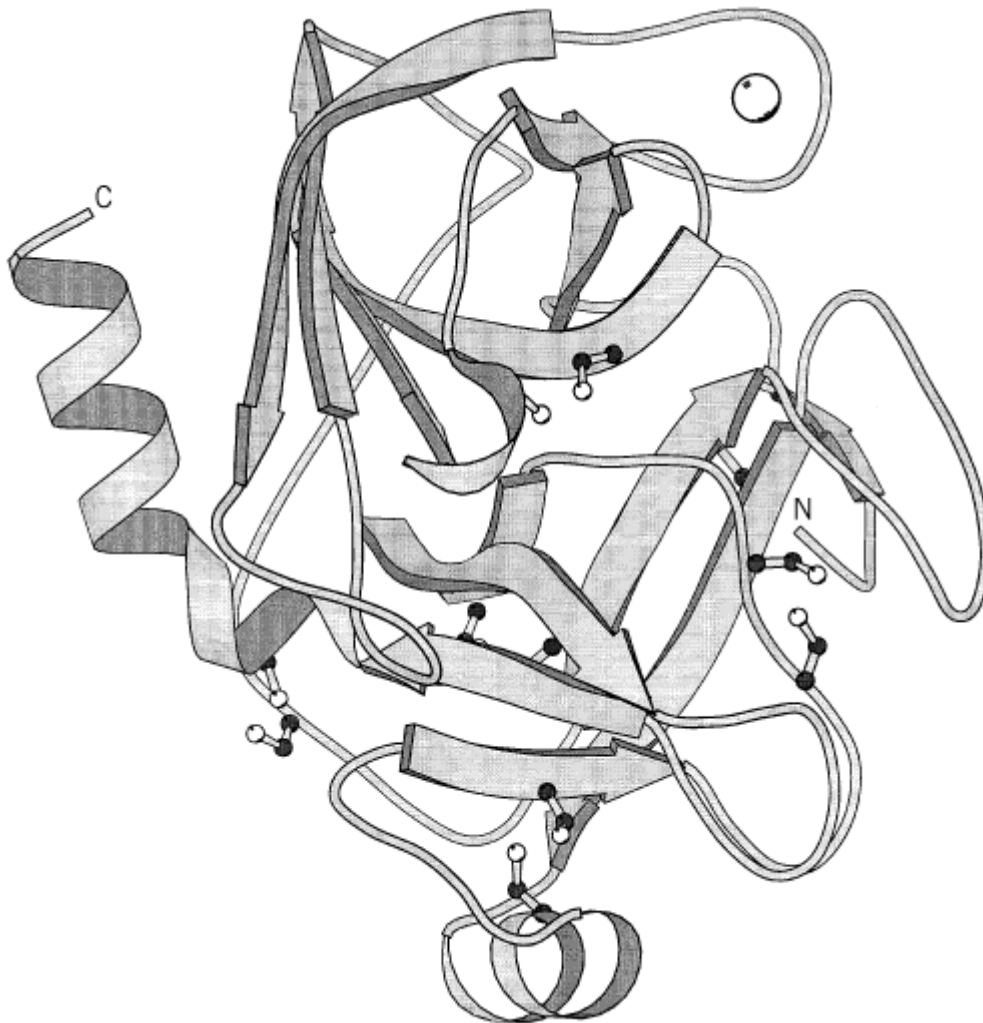
## Trypsin, Trypsinogen

Trypsin is a [serine proteinase](#) of the **chymotrypsin** evolutionary family that is synthesized in the pancreas in the form of an inactive [zymogen](#) precursor, trypsinogen (**1**). When trypsinogen is secreted from the pancreas into the small intestine, it is activated by another serine proteinase, [enterokinase](#), which is attached to the surface of the epithelial cells that line the intestine. Activation is accomplished by removal of a hexapeptide from the amino terminus of the precursor molecule, and this allows the remaining 223-residue polypeptide chain to change its conformation slightly (Fig. **1**) and generate a [catalytic triad](#) at its [active site](#); substrate hydrolysis can then occur. Trypsin is an important proteinase because, once it becomes active, it then activates all the other pancreatic zymogens—for example, chymotrypsinogen, pro-[elastase](#), pro-[carboxypeptidase](#), and so on. It is specific for cleaving peptide bonds in which the carbonyl group is contributed by [arginine](#) or [lysine](#) residues. It does not have much preference for the amino acid that contributes the NH-group to the peptide bond to be cleaved, except that it will be slower if that residue is acidic and will not work at all with [proline](#). This specificity has made trypsin a favorite enzyme for sequence analysis of proteins (see [Peptide Mapping](#)). Arginine and lysine are not particularly abundant constituents of



most proteins. Hence, cleavage of a protein with trypsin usually generates a limited, (ie, experimentally manageable) number of peptides and in good yield. Trypsin preparations for use in sequencing usually have to be treated with [TPCK](#), a potent inhibitor of chymotrypsin, which is the contaminant most frequently encountered in such preparations.

**Figure 1.** The three-dimensional structure of trypsin. The polypeptide backbone is depicted schematically as a ribbon, with arrows for b-strands and coils for a-helices. The sphere is a stabilizing calcium ion.



Trypsin is active at neutral to slightly alkaline pH, is stabilized by calcium ions, undergoes self-digestion fairly quickly, and is reversibly denatured in dilute acid. It is inhibited by, among others, **DIFP**, **TLCK**, **PMSF**, **leupeptin**, **soybean trypsin inhibitor**,  $\alpha_1$ -*antitrypsin*, and  $\alpha_2$ -[macroglobulin](#). The last two normally circulate in blood plasma and thereby prevent inappropriate [blood clotting](#).

#### Bibliography

1. R. Huber and W. Bode (1978) *Acc. Chem. Res.* **11**, 114–122.

## Tryptase

Mast cells are white blood cells found in connective tissue that play a role in inflammation. They contain **secretory granules**—and thus are considered *granulocytes*—which contain substantial amounts of **proteolytic** enzymes known as *granzymes* (1). One of these is tryptase, a [serine proteinase](#) (E.C. 3.4.21.59) that is specific for the catalysis of the hydrolysis of [peptide bonds](#) in which the carbonyl group comes from [arginine](#) or [lysine](#) residues, like **trypsin**. It is secreted from the granules during the immediate reaction that occurs shortly after an individual is exposed to an allergen or an [immunogen](#) to which that person has been sensitized. In fact, it can be quantified in biological fluids by [immunoassay](#) and thereby has diagnostic significance (2). The exact functional role played by the granzymes is not well known, but tryptase can inactivate [fibrinogen](#) and prevent the coagulation that could otherwise occur when blood plasma diffuses into sites of inflammation. It also can lead to the activation of the matrix [metalloproteinases](#) that degrade interstitial tissue and regulate neuropeptide activity.

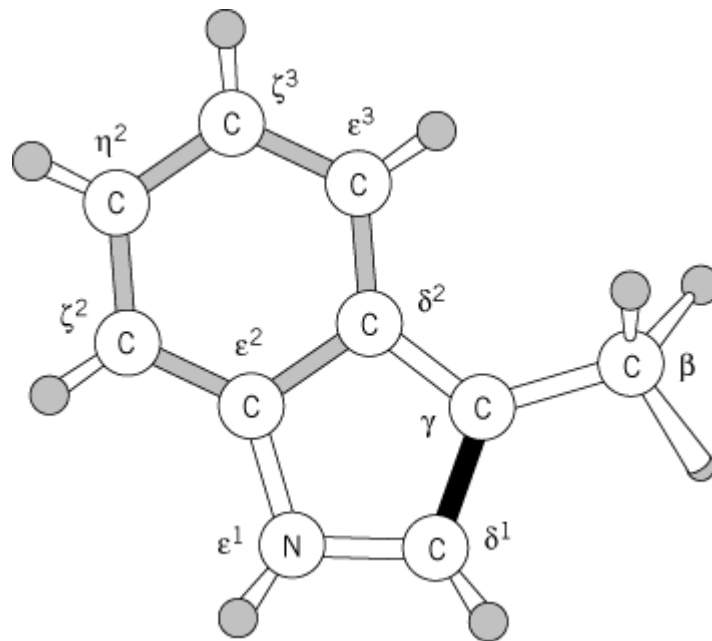
### Bibliography

1. M. J. Smyth, M. D. O'Connor, and J. A. Trapani (1996) *J. Leukoc. Biol.* **60**, 555–562.
2. L. B. Schwartz (1990) *J. Allergy Clin. Immunol.* **86**, 594–598.

## Tryptophan (TRP, W)

The [amino acid](#) tryptophan is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to only one **codon**—UGG. It is the least common amino acid residue and represents only 1.3% of the residues in the proteins that have been characterized. The tryptophanyl residue incorporated is the largest of the twenty normal residues, with a mass of 186.21 Da, a **van der Waals volume** of 163 Å<sup>3</sup>, and an [accessible surface](#) area of 259 Å<sup>2</sup>. Trp residues are those changed least frequently during [divergent evolution](#), and are not replaced by any particular type of residue.

The indole-containing Trp side chain is the largest and most complex of all the normal amino acids:



(1)

The five-membered ring is susceptible to oxidation and various other chemical reactions. Ozone opens the indole ring to that of *N*-formylkynurenine. [Peptide bonds](#) following Trp residues can be cleaved with varying efficiencies by several chemical reagents, especially iodosobenzoic acid and BNPA-skatole. Most of these procedures have the disadvantages of side-reactions and of oxidizing **Cys** and **Met** residues.

The aromatic side chain is largely nonpolar (except for the tendency of the electrons of the peripheral H atoms to be drawn into the aromatic ring), and only the N atom is usually involved in **hydrogen bonding**, although only weakly. Consequently, Trp is considered the most **hydrophobic** of all the amino acid residues, but only about 27% of Trp residues in folded [protein structures](#) are completely buried. In contrast to **Tyr** and **Phe** residues, buried Trp side chains do not normally undergo flipping of the aromatic ring. This is due to the larger size of the indole ring and the absence of symmetry, so that full 360° flips would be required. Trp favors moderately the **alpha-helical** conformation in model peptides and occurs in this type of **secondary structure** in folded proteins, but more frequently in [beta-sheets](#).

Trp residues occur so infrequently that many proteins have only one or a few Trp residues. This and the **absorbance** and **fluorescence** properties of Trp residues make them very useful in characterizing protein structure. Their fluorescence predominates in proteins and is especially sensitive to the environment of the side-chain, but in largely unpredictable ways. This is the only side chain capable of participating in charge-transfer complexes with pyridinium compounds and other electrophiles.

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Tubulin

Tubulins are the protein building blocks of [microtubules](#). They have been isolated and characterized

from many eukaryotic cells and organisms, including vertebrate brain, a favorite source because of the high abundance of microtubules in neurons, unicellular organisms such as the yeast *Saccharomyces cerevisiae*, invertebrate sources such as sea urchin eggs and sperm, plant cells, and cultured mammalian cells. Tubulin is a highly acidic protein, and early methods for its isolation by [ion-exchange chromatography](#) took advantage of its acidic nature. Most often, however, tubulin has been purified by cycles of polymerization at warm temperatures and depolymerization at cold temperatures, taking advantage of the assembly–disassembly properties of microtubules.

It has not yet proven possible to obtain tubulin crystals suitable for [X-ray crystallography](#), but a molecular structure of  $\beta$ -tubulin and the  $\alpha\beta$ -tubulin heterodimer at a resolution of  $\sim 3.7 \text{ \AA}$  have been constructed recently by [electron crystallography](#) of two-dimensional brain-tubulin sheet structures (1). This remarkable work has revealed that each tubulin monomer is a compact globular structure formed by a core of two  [\$\beta\$ -sheets](#) surrounded by 12  [\$\alpha\$ -helices](#). The  $\alpha$ - and  $\beta$ -tubulins are distinct but **homologous** proteins; as a rule, their amino acid sequences are quite similar and they have many structural features in common, even though there is significant **divergence** among tubulins across species. Indeed, the density maps for  $\alpha$  and  $\beta$  tubulin from vertebrate brain are essentially superimposable. There are two guanine nucleotide-binding sites per tubulin dimer, one per subunit (see [GTP-Binding Proteins](#)). With brain tubulin, GTP binds nonexchangeably to  $\alpha$ -tubulin, and the nucleotide remains in the GTP form at all times. GTP binds exchangeably to the  $\beta$ -subunit when the tubulin is in solution and becomes hydrolyzed to GDP plus orthophosphate (Pi) as, or shortly after, the tubulin adds to the growing end of a microtubule. Once incorporated into the microtubule, the  $\beta$ -tubulin-liganded nucleotide remains nonexchangeably bound as GDP until the subunit dissociates from the microtubule.

The  $\alpha$ - and  $\beta$ -tubulins are members of [gene families](#) (2). The  $\beta$  tubulin gene family in vertebrates consists of seven or more genes that encode at least six distinct tubulin isotypes. The  $\alpha$ -gene family has not been as extensively studied, but in the mouse at least seven genes encode six protein isotypes. There is extensive diversity in tubulin genes across species. A major variable region that distinguishes the  $\beta$ -tubulin isotypes is at the carboxyl terminus. This region, which is believed to be a major binding domain for MAPs, appears to be important in the regulation by MAPs of microtubule polymerization and dynamics.

The functional significance of tubulin diversity at both the tubulin gene and protein levels has remained unclear. Evidence primarily with  $\beta$ -tubulin isotypes has indicated that microtubules polymerized *in vitro* are copolymers of all available isotypes. These and related studies initially led to the idea that the multiple tubulin isotypes were not functionally important at the protein level. Recent evidence obtained in cells and *in vitro*, however, now indicates that the different tubulin isotypes—and, perhaps more interestingly, the specific isotype composition of a microtubule population—may indeed be important in specifying microtubule function.

The tubulins are also targets for multiple forms of [post-translational modification](#) (3). These include removal and replacement of the tyrosine residue at the C-terminus of  $\alpha$ -tubulin in many species, polyglutamylation both of  $\alpha$ - and  $\beta$ -tubulins near the C-termini, and **acetylation** of Lys40 of  $\alpha$ -tubulin in many species, **phosphorylation**, and glycation. The functional significance of these diverse post-translational modifications remains poorly understood. Some modifications, such as the detyrosination and acetylation of  $\alpha$  tubulin, are associated with formation of relatively stable microtubules. Others, such as polyglutamylation, which changes the charge properties in the MAP-binding C-terminal region, may regulate the binding of MAPs to the microtubule surface.

## Bibliography

1. E. Nogales, S. G. Wolf, and K. H. Downing (1998) *Nature* **391**, 199–203.
2. K. F. Sullivan (1988) *Annu. Rev. Cell Biol.* **4**, 687–716.
3. T. H. MacRae (1997) *Eur. J. Biochem.* **244**, 265–278.

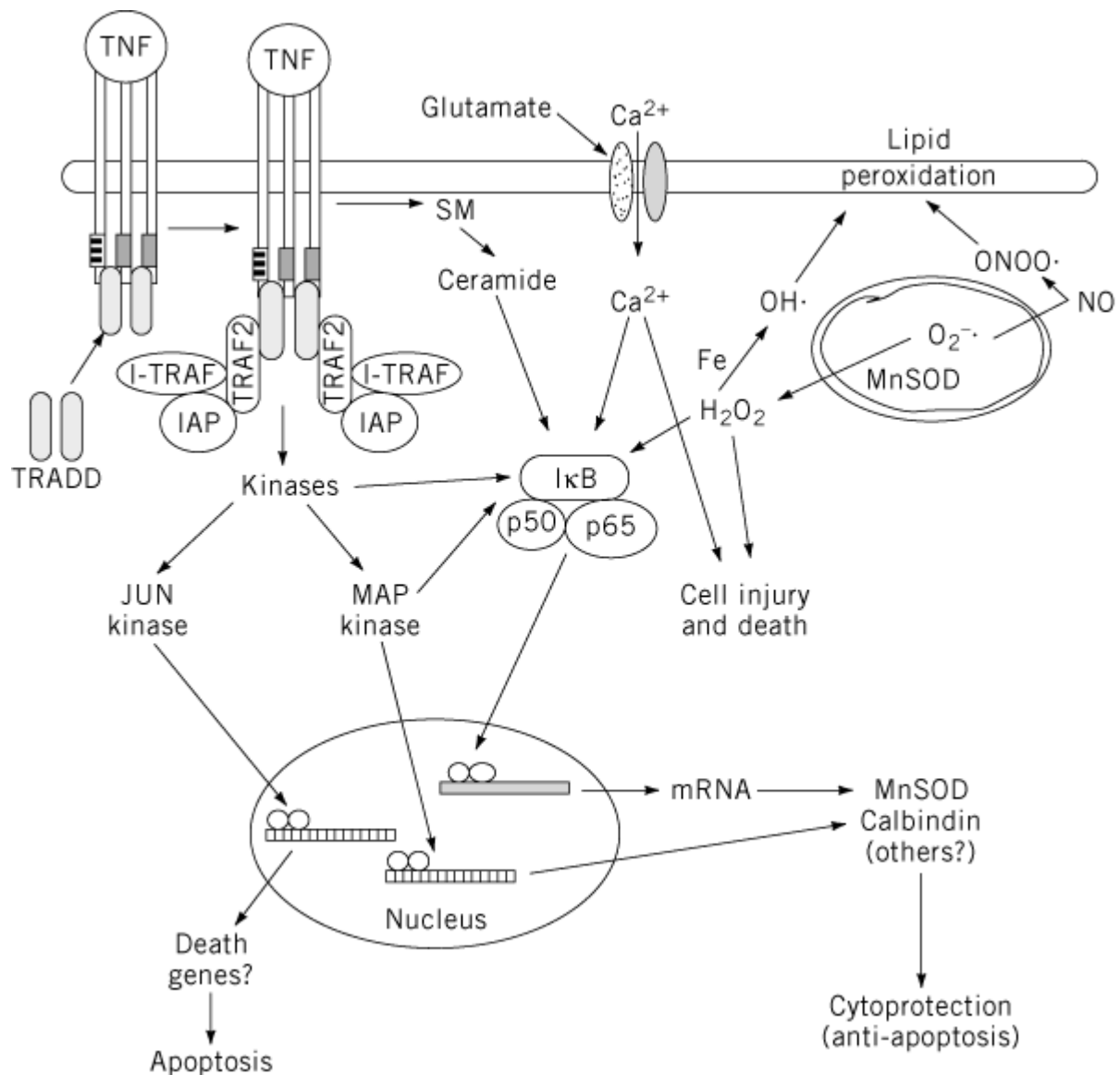
## Suggestion for Further Reading

4. R. F. Luduena (1998) The multiple forms of tubulin: different gene products and covalent modifications. *Int. Rev. Cytol.* **178**, 207–275.

## Tumor Necrosis Factor

Tumor necrosis factor- $\alpha$  (TNF $\alpha$ ) is a 17-kDa protein that is a member of a **cytokine**/trophic factor family that includes TNF $\beta$  (lymphotoxin- $\alpha$ ), CD40, CD30, and ligands of CD27. TNF $\alpha$  is normally expressed at low levels in most tissues, but its expression is rapidly increased in many different organ systems following injury or infection. Many different cell types express receptors for, and respond to, TNF $\alpha$ . Two distinct TNF $\alpha$  receptors p55 and p75, have been identified; engagement of the receptors involves binding of TNF $\alpha$  multimers and receptor trimerization. TNF $\alpha$  binds both receptors with equal affinity, but the receptors appear to be linked to different [signal transduction](#) pathways. Binding of TNF $\alpha$  to the p55 receptor results in recruitment of a panel of cytosolic proteins to the cytoplasmic domains of the receptor, including proteins called TRADD, TRAF2, and IAPs (Fig. 1). This signaling pathway ultimately leads to activation of [kinases](#), including a kinase that **phosphorylates** the I- $\kappa$ B subunit of the [transcription factor](#) NF- $\kappa$ B, and JUN (jun N-terminal) kinase and MAP (mitogen-activated protein) kinase. In addition, TNF $\alpha$  binding to p55 can activate a membrane-associated sphingomyelinase that liberates ceramide from sphingomyelin (Fig. 1). Ceramide may be another TNF $\alpha$ -induced [second messenger](#) that leads to activation of the transcription factor NF- $\kappa$ B. The signaling pathway of the p75 receptor is less well defined, but recent data suggest that this receptor may also be linked to NF- $\kappa$ B activation (1).

**Figure 1.** Signal transduction pathways activated by TNF $\alpha$  acting at the p55 receptor. Activation of the TNF p55 receptor occurs when TNF binds and induces receptor trimerization (upper left). The protein TRADD (TNF receptor-associated death domain) then associates with a “death domain” in the cytoplasmic portion of the TNF receptor which, in turn, causes association of TRAF2 (TNF receptor-associated protein 2) and IAP (inhibitor of apoptosis) with TRADD. Additional proteins, such as I-TRAF, may associate with TRAF2 and modulate signaling. This cascade of protein interactions results in activation of protein kinases and engagement of one or more downstream transcription-regulating pathways. The transcription factor NF- $\kappa$ B exists in an inactive form in the cytoplasm as a complex of three proteins, p50, p65, and I $\kappa$ B. Activation of a kinase by the TNF signaling cascade results in phosphorylation of I $\kappa$ B, liberating the p50/p65 dimer, which then translocates to the nucleus and binds to specific  $\kappa$ B binding sites in the promoters of  $\kappa$ B-responsive genes. An alternate TNF signaling pathway involving activation of sphingomyelinase (SM) and generation of ceramide may also result in NF- $\kappa$ B activation. NF- $\kappa$ B induces the expression of antiapoptotic genes, including those encoding the antioxidant enzyme manganese superoxide dismutase (Mn-SOD) and the calcium-binding protein calbindin. In many cell types, TNF also activates JUN kinase, which may induce the expression of proapoptotic death genes, and MAP kinase, which may regulate genes involved in cellular plasticity. Signaling pathways involving NF- $\kappa$ B, JUN kinase, and MAP kinases are sensitive to reactive oxygen species (eg, H<sub>2</sub>O<sub>2</sub>) and elevation of intracellular calcium levels.



In the inactive state, NF-κB consists of three proteins located in the cell cytosol: p50, p65, and I-κB. Signals that activate NF-κB do so by causing **ubiquitin-mediated protein degradation** of I-κB; the p50-p65 dimer then translocates to the nucleus and binds to DNA **enhancer elements** that contain specific κB-binding sequences. Genes responsive to NF-κB include those encoding various cytokines and [cell-adhesion molecules](#) (2). In addition, NF-κB induces the expression of cytoprotective gene products, including the antioxidant enzyme manganese superoxide dismutase and the [calcium-binding protein](#) calbindin D28k (3-6).

TNFα plays a variety of roles, prominent among them the propagation of cytokine cascades involved in cellular responses to injury and inflammatory responses. In tissues throughout the body, TNFα levels increase rapidly and dramatically following cell injury. For example, TNFα levels increase greatly in the brain following epileptic seizures and cerebral ischemia (7, 8). Studies have shown that TNFα can prevent [apoptosis](#) in cultured neurons by a mechanism involving activation of NF-κB (5, 9). The latter studies showed that activation of NF-κB with ceramide or I-κB **antisense** treatment suppressed apoptosis, whereas blockade of NF-κB activity using decoy DNA abolished the antiapoptotic action of TNFα. Although TNFα prevents death of postmitotic cells, such as neurons, it can directly kill many types of mitotic normal and transformed cells. In the latter cell types, TNFα apparently activates divergent pathways—a death pathway that may involve JUN/SAP kinases and an antideath pathway involving NF-κB (10). The ability of NF-κB to prevent apoptosis has important

implications for therapeutic approaches for a variety of major diseases, including cancer and neurodegenerative disorders. Thus, agents that block NF- $\kappa$ B activation may prove useful in treating certain types of cancer, whereas agents that activate NF- $\kappa$ B in neurons may prove beneficial in the many neurodegenerative conditions that involve apoptosis (eg, stroke, Alzheimer's disease, and Parkinson's disease). It should be noted that, whereas TNF $\alpha$  can act directly on some cell types (eg, neurons) to prevent or promote their death, it also activates macrophages and microglia, cells that produce cytotoxic substances, including [nitric oxide](#) and excitotoxins.

Targeted gene disruption methods have been used to generate mice lacking either p55 or p75, or both receptors. Remarkably, mice lacking one or both TNF $\alpha$  receptors exhibit no overt phenotypes and reproduce normally. However, challenges of the immune systems of these mice have revealed important roles for p55 in [T-cell](#) responses. Thus, mice lacking p55 are resistant to endotoxic shock but exhibit increased sensitivity to *L. monocytogenes* infection ([13](#), [14](#)). Responses of neurons and microglia to brain injury are also altered in mice lacking p55. Thus, the extent of neuronal degeneration following severe seizure activity and focal cerebral ischemia is increased in p55 knockout mice, whereas the microglial response to injury is suppressed ([7](#), [15](#)). Lack of TNF $\alpha$ -mediated upregulation of manganese superoxide dismutase may play a role in the altered cellular responses to brain injury in mice lacking p55 ([7](#), [16](#)).

Beyond its critical roles in cell injury responses, TNF $\alpha$  may serve important functions under normal physiological conditions. For example, recent studies have shown that TNF $\alpha$  can modulate neuronal excitability and plasticity in the brain. Thus, exposure of cultured rat hippocampal neurons to TNF $\alpha$  resulted in an increase in whole-cell voltage-dependent calcium currents and a decrease in currents through ionotropic glutamate receptors ([11](#)). Measurements of intracellular calcium levels in cultured neurons showed that glutamate-induced calcium influx was attenuated in neurons pretreated with TNF $\alpha$  ([3](#)). TNF $\alpha$  may also modify learning and memory processes, as suggested by recent studies showing that long-term depression of synaptic transmission in the hippocampus (a cellular correlate of learning and memory) is altered in mice lacking TNF $\alpha$  receptors. Interestingly, recent data suggest a role for TNF $\alpha$  in regulating sleep and circadian rhythms ([12](#)). These types of data suggest that TNF $\alpha$  may have important functions beyond its quite dramatic roles in tissue injury and repair.

## Bibliography

1. M. Rothe, M. G. Pan, W. J. Henzel, T. M. Ayres, and D. V. Goeddel (1995) *Cell* **83**, 1243–1252.
2. E. N. Benveniste and D. J. Benos (1995) *FASEB J.* **9**, 1577–1584.
3. B. Cheng, S. Christakos, and M. P. Mattson (1994) *Neuron* **12**, 139–153.
4. M. P. Mattson, B. Cheng, S. Baldwin, V. L. Smith-Swintosky, J. Keller, J. W. Geddes, S. W. Scheff, and S. Christakos (1995) *J. Neurosci. Res.* **42**, 357–370.
5. M. P. Mattson, Y. Goodman, H. Luo, W. Fu, and K. Furukawa (1997) *J. Neurosci. Res.* **49**, 681–697.
6. B. B. Warner, M. S. Burhans, J. C. Clark, and J. R. Wispe (1991) *Am. J. Physiol.* **260**, 296–301.
7. A. J. Bruce, W. Boling, M. S. Kindy, J. Peshon, P. J. Kraemer, M. K. Carpenter, F. W. Holtsberg, and M. P. Mattson (1996) *Nature Med.* **2**, 788–794.
8. T. Liu, R. K. Clark, P. C. McDonnell, P. R. Young, M. S. White, F. C. Barone, and G. Z. Feuerstein (1994) *Stroke* **25**, 1481–1488.
9. S. W. Barger, D. Horster, K. Furukawa, Y. Goodman, J. Krieglstein, and M. P. Mattson (1995) *Proc. Natl. Acad. Sci. USA* **92**, 9328–9332.
10. A. A. Beg and D. Baltimore (1996) *Science* **274**, 782–784.
11. K. Furukawa and M. P. Mattson (1998) *J. Neurochem.* **70**, 1876–1886.
12. J. Y. Fang and J. Krueger (1997) *J. Neurosci.* **17**, 5949–5955.
13. K. Pfeffer, T. Matsuyama, T. M. Kundig, A. Wakeham, K. Kishihara, A. Shahinian, K. Wiegmann, P. S. Ohashi, M. Kronke, and T. W. Mak (1993) *Cell* **73**, 457–467.

14. L. Zheng, G. Fisher, R. E. Miller, J. Peschon, D. H. Lynch, and M. J. Lenardo (1995) *Nature* **377**, 348–351.
15. D. S. Gary, A. J. Bruce-Keller, M. S. Kindy, and M. P. Mattson (1998) *J. Cereb. Blood Flow Metab.* **18**, 1283–1287.
16. J. N. Keller, M. S. Kindy, F. W. Holtsberg, D. St. Clair, H. C. Yen, A. Germeyer, S. M. Steiner, A. J. Bruce-Keller, J. B. Hutchins, and M. P. Mattson (1998) *J. Neurosci.* **18**, 687–697.

### Suggestions for Further Reading

17. C. A. Smith, T. Farrah, and R. G. Goodwin (1994) The TNF receptor superfamily of cellular and viral proteins: activation, costimulation, and death. *Cell* **76**, 959–962.
18. B. G. Darnay and B. B. Aggarwal (1997) Early events in TNF signaling: a story of associations and dissociations. *J. Leukocyte Biol.* **61**, 559–566.
19. M. P. Mattson and O. Lindvall (1997) "Neurotrophic factor and cytokine signaling in the aging brain". In *The Aging Brain* (M. P. Mattson and J. W. Geddes, eds.), *Adv. Cell Aging Gerontol.* **2**, 299–345.
20. M. P. Mattson (1998) Free radicals, calcium, and the synaptic plasticity—cell death continuum: emerging roles of the transcription factor NF- $\kappa$ B. *Int. Rev. Neurobiol.* **42**, 103–168.

## Tumor Promoters

### 1. Definition

A tumor promoter is generally defined as a chemical, a complex of chemicals, or a biological agent that promotes a later stage of carcinogenesis, called *tumor promotion*, by altering expression of the genetic information, rather than altering the structure of DNA.

### 2. Tumor Promotion

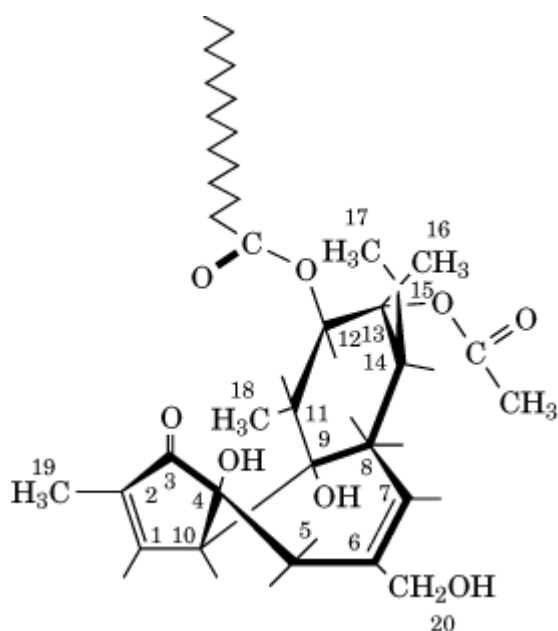
**Carcinogenesis** in skin and other epithelia can be divided into three distinct stages: (i) initiation, (ii) tumor promotion, and (iii) tumor progression. Historically, Berenblum and Shubik (1) first demonstrated in 1947 the initiation and promotion stages in skin carcinogenesis of mice. Initiation resulted from the administration of a single dose of a known carcinogen, called the *initiator*, which initiates carcinogenesis but is not sufficient to develop tumors. The stage of promotion requires repeated applications of a second compound, called the *promoter*, which by itself is incapable of inducing carcinogenesis. Reversed order of application of these compounds—that is, first promoter, then initiator—does not produce tumors. Tumors developed by the initiation-promotion protocol (two-stage carcinogenesis protocol) are largely of benign nature, but they progress further to malignant tumors, either spontaneously or by treatment with carcinogens. This stage is called *progression*.

### 3. Tumor-Promoting Agents

The classical experiment of Berenblum and Shubik used croton oil, the seed oil of Euphorbiacea, *Croton tiglium*, as a tumor promoter. Hecker and coworkers (2) demonstrated that the 12,13 diesters of phorbol are the irritant and promoting principles of croton oil, of which the best known and most potent is 12-*O*-tetradecanoylphorbol 13-acetate (TPA), also known as phorbol myristate acetate (PMA) (Fig. 1). TPA is widely used not only as a tumor promoter but also as a modulator of biological responses in a wide variety of experiments.



**Figure 1.** Structure of 12-*O*-tetradecanoylphorbol 13-acetate (TPA), also known as phorbol myristate acetate (PMA).



Very many diverse chemicals, complexes of chemicals, or biological agents are listed as being tumor promoters with a significant degree of tissue specificity. For example, phenobarbital, 2,3,7,8-tetrachlorodibenzo-*p*-dioxin (TCDD), mycotoxins, and hormones have all been shown to promote liver carcinogenesis in rats and mice. Saccharin acts as a promoter in bladder carcinogenesis of rats. Human carcinogenesis may be promoted by chemicals, environmental factors, and life styles. These include dietary fat in breast and colon cancer, sodium chloride in stomach cancer, cigarette smoke in lung cancer, asbestos in lung cancer and mesothelioma, synthetic estrogens in liver adenoma, and alcoholic beverages in liver and esophageal cancer.

#### 4. Biological Responses

Phorbol ester tumor promoters evoke a wide variety of biological responses in cells, including stimulation or inhibition of terminal [differentiation](#), stimulation of [DNA replication](#) and cellular proliferation, modulation of cellular responses to external stimuli (eg, **cytokines**), and disturbance of cell-cell communication. Topical application of TPA to dorsal skin of mice results in epidermal hyperplasia. In skin carcinogenesis models, tumor promoters induce clonal expansion of the initiated cells as a consequence of these biological responses.

#### 5. Activation of Protein Kinase C

The pleiotropic cellular effects of phorbol esters were hardly understood before the revolutionary work of Nishizuka's laboratory in 1982 (3). They demonstrated that phorbol ester tumor promoters activated protein kinase C (PKC), a **serine-threonine kinase** normally activated by metabolic modification of polar head groups of membrane [lipids](#). Physiologically, PKC is normally activated by **diacylglycerol** (DG) in the presence of phosphatidylserine (PS) and Ca<sup>2+</sup> (for conventional PKC). DG is formed through [phosphatidylinositol](#) turnover by phospholipase C. Phorbol ester tumor promoters activate PKC in place of DG. While DG is formed and further metabolized within a short period of time, phorbol esters bind to and activate PKC strongly, resulting in sustained down-regulation.

PKC exists as a family containing at least 10 members classified into three groups: (i) Ca<sup>2+</sup>, PS, and DG-dependent conventional PKC (α, βI, βII, and γ **isoforms**); (ii) Ca<sup>2+</sup>-independent novel PKC (δ, ε, ζ, and η isoforms); and (iii) Ca<sup>2+</sup>- and DG-independent atypical PKC (θ and ι isoforms). It is not known which isoforms are involved in promotion, although the η isoform was found to regulate negatively the tumor promotion process by induction of terminal differentiation of keratinocytes (4).

#### 6. Induction of Ornithine Decarboxylase

[Ornithine decarboxylase](#) (ODC) is a key enzyme regulating the polyamine biosynthetic pathway. Topical application of phorbol esters causes rapid and dramatic induction of ODC (5). In tumors produced by the two-stage skin carcinogenesis protocol, the levels of ODC and polyamines are elevated constitutively. Of particular interest is that tumor promoters are no longer required for two-stage carcinogenesis when ODC is overexpressed in hair follicle keratinocytes, implying a crucial role of ODC induction in carcinogenesis (6).

#### 7. Gene Expression

In contrast to initiators that interact with cellular genes, most tumor promoters alter **gene expression** in an [epigenetic](#) manner. There is convincing evidence that phorbol ester tumor promoters exert their biological effects by induction of the [transcription factor](#) complex AP-1, which consists of members of the *fos* and *jun* nuclear **oncogene** families. The AP-1 complex activates [transcription](#) of the genes having the TPA-responsive element (TRE) in their **promoter** regions (7, 8) (see [Response Element](#)). Thus, the transcription regulated by AP-1/TRE plays a key role in the early events triggered by phorbol ester tumor promoters.

#### 8. Genotoxicity

There is some evidence that phorbol ester tumor promoters lead to genotoxicity at the [chromosome](#) level, including [aneuploidy](#), chromosomal aberrations, and [gene amplification](#). However, these effects might be indirect, resulting from diverse cellular responses.

#### Bibliography

1. I. Berenblum and P. Shubik (1947) Br. J. Cancer **1**, 379–382.
2. E. Hecker (1968) Cancer Res. **28**, 2338–2349.
3. M. Castagna et al. (1982) J. Biol. Chem. **257**, 7847–7851.
4. K. Chida et al. (1995) Cancer Res. **55**, 4865–4869.
5. T. G. O'Brien, R. C. Simisiman, and R. K. Boutwell (1975) Cancer Res. **35**, 1662–1670.
6. T. G. O'Brien, L. C. Megosh, G. Gilliard and A. P. Soler (1997) Cancer Res. **57**, 2630–2637.
7. P. Angel, T. Smeal, and M. Karin (1988) Cell **55**, 875–885.
8. W. Lamph, P. Wamsley, P. Sasson-Corsi, and I. Verma (1988) Nature, **334**, 629–631.

### Tumor Suppressor Genes

**Oncogenes** are usually dominant, mutated forms of proto-oncogenes. In several instances, however, cell [transformation](#) is also induced by recessive mutations of certain **genes**, termed antioncogenes or

tumor suppressor genes. As the name indicates, the loss of the functional activity of these genes, by either point mutations or deletions, leads to oncogenesis, suggesting that the products of these genes suppress tumorigenesis (1). Such loss of tumor suppressor gene function causes retinoblastoma, Wilm's tumor, colorectal carcinoma, neurofibromatosis type 1, and familial adenomatous polyposis. Although several tumor suppressor genes have been identified, many remain unknown. Two well-characterized tumor suppressor genes, the **retinoblastoma** gene and the [p53](#) gene, are discussed here and in their individual articles.

The protein encoded by the retinoblastoma (Rb) gene is a 105-kDa polypeptide chain, also known as p105RB, that can be **phosphorylated**. The phosphorylation state of p105RB is tightly regulated during the cell cycle, is maximal during the S-phase, and minimal soon after mitosis. Stimulation of quiescent T cells leads to hyperphosphorylation of p105RB, whereas differentiation of myeloid cells is associated with very low levels of phosphorylated p105RB. Interestingly, only the hypophosphorylated form of RB is tumor suppressive. p105RB is phosphorylated in the late G<sub>1</sub>/S-phase and dephosphorylated in the late M-phase. In its dephosphorylated state, the Rb gene product binds to a group of transcription factors, termed the E2F class of proteins, which are required for [DNA synthesis](#) (1, 2). Growth factors stimulate phosphorylation of the Rb protein by activating a group of kinases known as CDC kinases, which phosphorylate Rb. The p34CDC2 kinase (2) is a candidate kinase involved in this process. In turn, this process reduces the affinity of p105RB for E2F, allowing these transcription factors to bind to DNA and to activate the transcription of genes required for DNA synthesis. Interestingly, growth inhibitors, such as [transforming growth factor b](#) (TGF-b) which inhibits cell proliferation, prevent phosphorylation of p105RB, even in the G<sub>1</sub>/S-phase. These observations suggest that the phosphorylation of p105RB is a critical regulatory event in cell proliferation.

p105RB binds to the [adenovirus](#) E1A protein, and this association is essential for adenoviral transformation. Hence, transforming DNA viruses induce transformation by sequestering the hypophosphorylated form of RB105. The biochemical mechanism of tumor suppression by p105RB is at the level of expression of the proto-oncogenes *c-myc* and *c-fos* because RB expression suppresses the expression of both genes.

p53 is another well-characterized tumor suppressor gene, which encodes a polypeptide chain of 375 amino acids. As with RB, loss of p53 function results in the onset of oncogenesis. p53 plays a pivotal role during the execution of **apoptotic** ([programmed cell death](#)) pathways that are induced as a result of **DNA damage** caused by agents, such as ultraviolet or ionizing radiation. Current models suggest that DNA damage results in up-regulating p53 gene transcription and accumulation of this protein in the cell (3). Accumulation of p53 results in enhanced transcription of a group of genes, such as the cell cycle-dependent kinase inhibitor p21/WAF-1/cip-1, that are responsible for stalling the cell cycle in either the G<sub>1</sub>- or G<sub>2</sub>/M-phases. This mechanism provides the cell with valuable time to make critical repairs in its genetic material. If the necessary repairs cannot be made, the p53 protein initiates a "suicidal" apoptotic program that results in cell death to prevent clonal expansion of a cell that has a mutated genome. Malignant cells have evolved novel strategies to destroy p53 function. One mechanism is through mutating the p53 gene and protein. Normal p53 protein binds to DNA and transactivates transcription of a distinct set of target genes. Mutated forms of p53, however, have lost the ability to bind DNA in this fashion and therefore behave abnormally. A second mechanism by which normal p53 function is subverted in a cancerous cell is through the MDM2 protein. The MDM2 protein binds to p53 and is amplified in many tumor cells. When these two proteins are bound, p53 cannot mediate growth arrest at the cell cycle checkpoints (as described previously), which leads to rapid clonal expansion of cells that have unstable, mutated genomes.

## Bibliography

1. R. A. Weinberg (1991) *Science* **254**, 1138–1145.
2. X. Graña and E. P. Reddy (1996) *Oncogene* **11**, 211–220.

3. L. J. Ko and C. Prives (1996) *Genes Dev.* **10**, 1054–1072.

## Turnover Number

The turnover number for an **enzyme**-catalyzed reaction is defined as the number of moles of substrate converted to product per mole of enzyme per second under conditions at which the concentration of all substrates is saturating. Turnover numbers are expressed in units of  $s^{-1}$ . The turnover number can also be expressed as moles of substrate converted to product per mole of catalytic subunit per second; this is the most useful definition for the enzymologist, but such a definition can apply only when the number of catalytically competent subunits is known. The relationship between the rate constants associated with a reaction and its turnover number can be illustrated by reference to an ordered Uni–Bi reaction (Fig. 1). Such a **kinetic mechanism** might be catalyzed by a **phosphatase**, with inorganic phosphate as the last product to be released. For the most general case in which catalysis and the release of both products limit the maximum velocity of the reaction  $V$ , the turnover number would be given by the expression

$$\frac{V}{E_t} = \frac{k_3 k_5 k_7}{k_3(k_5 + k_7) + k_7(k_4 + k_5)} \quad (1)$$

where  $E_t$  is the total amount of enzyme present. Equation (1) simplifies to

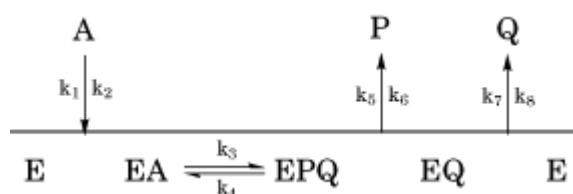
$$\frac{V}{E_t} = \frac{k_3 k_5}{(k_3 + k_4 + k_5)}$$

if the release of Q from the EQ complex is fast compared with catalysis and the release of P. The equation would simplify to

$$\frac{V}{E_t} = \frac{k_5 k_7}{(k_5 + k_7)}$$

if the maximum velocity was limited by the rate of release of the products, which is common for enzyme-catalyzed reactions (1). If the rate of release of the products was much faster than catalysis,  $V/E_t$  would simply be equal to  $k_3$ , which is the relationship that is derived from [Michaelis-Menten kinetics](#). In connection with the aforementioned comments, it should be noted that values for  $k_3$  and  $k_4$  cannot be obtained directly by the use of steady-state kinetic techniques although they might well be determined through the application of transient-state kinetic analysis (2).

**Figure 1.** Ordered Uni–Bi reaction kinetic mechanism to illustrate the unimolecular steps that can limit the maximum velocity of the reaction.



The turnover numbers for enzymes vary considerably from 2 to  $5 \text{ s}^{-1}$  for **ribulose biphosphate carboxylase** (3), which is a very sluggish enzyme, to 1 million  $\text{s}^{-1}$  for [carbonic anhydrase](#) (see [Michaelis–Menten Kinetics](#)).

### Bibliography

1. W. W. Cleland (1975) *Acc. Chem. Res.* **8**, 145–151.
2. K. A. Johnson (1992) *Enzymes* **20**, 1–61.
3. F. C. Hartman and M. R. Harpel (1993) *Adv. Enzymol.* **67**, 1–75.

## Turns

Turn is a general term that describes several types of nonregular [secondary structure](#) in **protein structures** that cause a change in the direction of the [polypeptide chain](#). Turns are also known as reverse turns (since they reverse the chain direction). They include [b-turns](#), [g-turns](#), [hairpins](#), [W-loops](#), and [b-bulges](#) (although for b-bulges, the change in direction is much less than that of other types of turn).

The b-turn defines a four-residue turn stabilized by a backbone [hydrogen bond](#) between the first and fourth residues. Similarly, a g-turn defines a three-residue turn stabilized by a backbone hydrogen bond between the first and third residues. Different types of b-turn and g-turn exist and are classified by their backbone conformation. Hairpin defines a b-strand-turn-b-strand motif, where the turn is often a b-turn or a g-turn. The W-loop defines a loop-shaped [backbone](#) conformation of six or more residues that resembles the Greek letter omega (W); b-bulge defines a disruption to the regular structure of a [b-sheet](#) caused by an additional residue in an edge b-strand.

Turns are distinguished from the regular types of secondary structure ([b-helix](#) and b-sheet) in that they do not have repetitive backbone conformations and hydrogen bonding patterns. Most residues in turns, however, adopt backbone conformations close to the a-helical or b-strand conformations. The difference is that the backbone conformations of successive residues in a turn vary, whereas in b-strands and a-helices the backbone conformations of consecutive residues are usually similar. [Glycine](#) or [proline](#) residues are often found in turns, especially where unusual backbone conformations are required. Turns are usually found at the surface of proteins, where the peptide bonds that do not interact with protein atoms can form hydrogen bonds with the solvent.

[See also [Beta-Turns](#), [Gamma-Turns](#), [Beta-Bulge](#), [Secondary Structure, Protein](#), and [Omega Loop](#).]

### Suggestions for Further Reading

- C. M. Wilmot and J. M. Thornton (1988) Analysis and prediction of the different types of -turn in

proteins. *J. Mol. Biol.* **203**, 221–232.

G. Rose, L. M. Gierasch, and J. A. Smith (1985) Turns in peptides and proteins. *Adv. Protein Chem.* **37**, 1–109.

A. V. Efimov (1993) Patterns of loop regions in proteins. *Curr. Opin. Struct. Biol.* **3**, 379–384.

## Twist, DNA

Twist is an aspect of [DNA structure](#) that concerns [DNA topology](#).

### 1. Total Twist

#### 1.1. General Considerations

The twist is a property of two curves that are connected by a one-to-one correspondence surface. Generally, the twist  $Tw(C', C'')$  of one curve ( $C'$ ) about another ( $C''$ ) is a measure of how often  $C'$  spins about  $C''$  as one advances along  $C''$ . In contrast to the **linking number**, the twist depends on the ordering of the two curves; hence, except for certain special cases (1),  $Tw(C', C'') \neq Tw(C'', C')$ . In considering the twist of DNA, one possible choice of  $C'$  and  $C''$  is to use the two backbone chains,  $C_1$  and  $C_2$ , giving  $Tw(C_1, C_2)$  or, alternatively,  $Tw(C_2, C_1)$ , as the DNA twist. Although this is useful in certain special cases (2), the choice is unnatural for closed duplex DNA (cdDNA). This is because the fundamental relationship (3) for cdDNA,  $Lk = Tw + Wr$ , would then involve the **writhe** of either of the backbone chains, which is seldom desirable.

#### 1.2. The Strand-Axis Twist in DNA

As applied to DNA, the choice of  $C'$  and  $C''$  is dictated by the natural duplex helical geometry: one of the curves is customarily chosen to be the duplex axis,  $A$ ; and the other curve,  $C$ , is taken to be a line along either of the two strands. The symbol  $Tw$  alone is understood to mean  $Tw(C, A)$ . The question then arises as to whether  $Tw(C_1, A) = Tw(C_2, A)$ ; that is, does it matter which strand is chosen? This is answered by writing the fundamental equation for both choices:  $Lk(C_1, A) = Tw(C_1, A) + Wr(A)$  and  $Lk(C_2, A) = Tw(C_2, A) + Wr(A)$ . Now  $Lk(C_1, A) = Lk(C_2, A)$ , except for special cases in which the DNA contains asymmetric local substructures (as, eg, the extrusion of imperfect [palindromes](#) into cruciforms). Such special cases require use of the intersection number, a separate topological quantity (2), in the fundamental equation. For most cases of practical interest, in which the duplex axis is uninterrupted,  $Tw(C_1, A) = Tw(C_2, A)$ . If the DNA axis is a straight line, or if the axis is closed and lies entirely in a plane, the twist is simply the number of revolutions made by  $C$  as it winds around  $A$ . Here both  $C$  and  $A$  are taken to be oriented in the same direction. For example, a [B-DNA](#) containing 15 bp has a twist of  $Tw_0 = N/h_0 = 15/10.5 = 1.429$ , where  $N$  is the number of base pairs and  $h_0$  is the number of base pairs per turn of the helical repeat. DNA is, however, rarely encountered with a truly linear axis. Any out-of-plane motion of axis  $A$  leads to values of  $Tw \neq Tw_0$ , with the extent of the deviation dependent on the amount of axial torsion. A more general way of calculating  $Tw$  is therefore needed.

Since the twist depends on the joint trajectory of  $A$  and  $C$ , it is most straightforward to use vector notation, as described in Figure 1. Let the DNA axis,  $A$ , lie on a surface as shown.  $\mathbf{T}_A$  is the unit tangent vector to  $A$  at a particular point  $a$ . One of the two DNA backbone curves, here  $C$ , winds

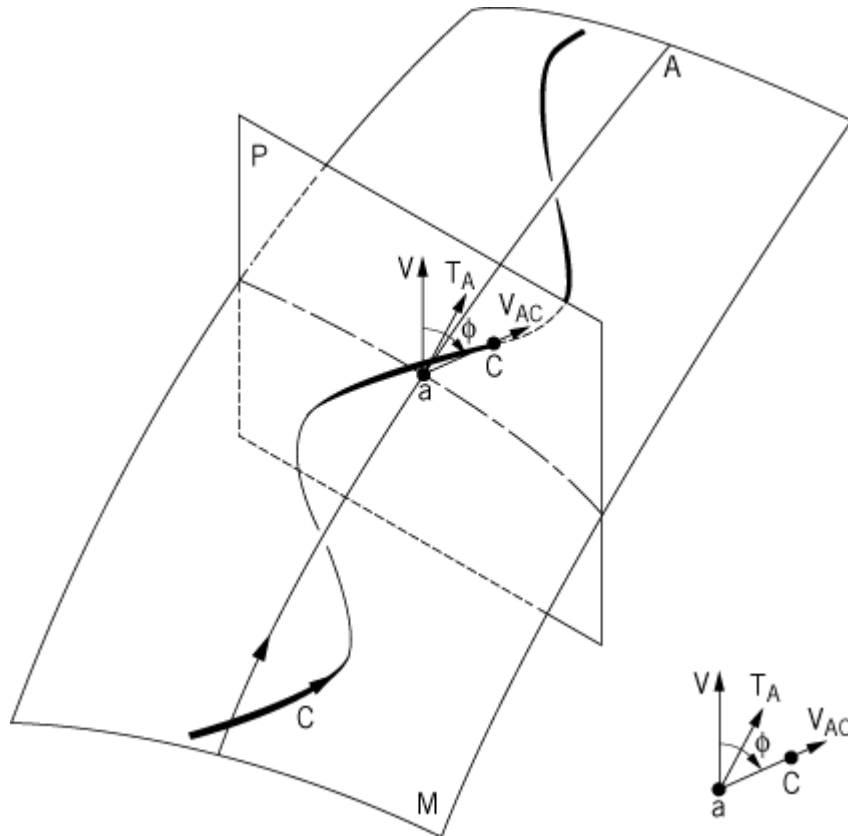
around  $A$  and cuts the surface twice per revolution. The vector  $\mathbf{v}$  is the surface normal at any point on  $A$ . The unit vector  $\mathbf{v}_{AC}$  is perpendicular to  $A$  and connects  $A$  to  $C$ . Then the twist of either of the two DNA strands about the duplex axis is

$$Tw \equiv Tw(C, A) = \frac{1}{2\pi} \int_A [\mathbf{T}_A \times \mathbf{v}_{AC}] \cdot d\mathbf{v}_{AC} \quad (1)$$

This equation can be used to calculate the twist of duplex DNA having any particular axis and backbone geometry. In the case of DNA, the twist is a measure of the spinning of either backbone phosphodiester–sugar chain about the (imaginary) duplex axis. Conceptually,  $Tw$  is the total change in  $\mathbf{v}_{AC}$  in the direction of the normal vector,  $\mathbf{T}_A \times \mathbf{v}_{AC}$ , as the entire DNA segment is traversed.

Physically,  $Tw$  provides the torsional contribution to the elastic energy of the DNA segment modeled as an elastic rod.

**Figure 1.** Definition of the vectors used in calculation of the twist and the surface twist. The DNA axis,  $A$ , lies on a surface  $M$  as shown. One of the two backbone curves,  $C$ , winds around  $A$  and cuts the surface twice per revolution. The vector  $\mathbf{v}$  is the surface normal to  $M$  at any point on  $A$ . An imaginary plane ( $P$ ), is perpendicular to the axis curve and moves along  $A$ . The plane intersects backbone curve  $C$  at successive points  $c$ . The unit vector  $\mathbf{v}_{AC}$  is perpendicular to  $A$  and connects  $A$  to  $C$  at point  $c$ . The backbone chain passes, in a right-hand sense, alternatively above and below the surface.  $\mathbf{T}_A$  is the unit tangent vector to  $A$  at a particular point,  $a$ . The winding number,  $F(A)$ , of  $C$  about  $A$  is given by the number of revolutions of  $\mathbf{v}_{AC}$  about  $\mathbf{v}$  as  $P$  advances along curve  $A$ . This is shown separately in the inset. For the calculation of the total twist, Equation (1), the surface  $M$  is not needed, nor is the vector  $\mathbf{v}$ .



## 2. Twist of DNA Wrapped on a Protein Surface

DNA is commonly associated with proteins; the best documented example is the [nucleosome](#). In

most cases, this interaction brings about a nonlinearity in the DNA axis, meaning that the twist differs from that for a straight-line axis. In general, the twist for DNA that is wrapped on a protein surface may be divided into two parts (4):

$$Tw = \Phi + STw \quad (2)$$

where the two independent components of the twist are  $F(A)$ , the winding number of the DNA axis,  $A$ , and  $STw$ , the surface twist.

### 2.1. Winding Number

The winding number,  $F(A)$ , is a measure of the number of times that  $\mathbf{v}_{AC}$  rotates about  $\mathbf{v}$  as the axis curve  $A$  is traversed:

$$\Phi(A) = \left(\frac{1}{2\pi}\right) \int_A d\phi \quad (3)$$

where  $\phi$  is the angle between  $\mathbf{v}$  and  $\mathbf{v}_{AC}$  (Fig. 1, inset). The winding number has also been called (5) the “local reference frame twist,” in which context  $Tw$  is called the “laboratory frame twist.” The winding number can be measured directly by any technique that is differentially sensitive to either DNA chain when it is near or removed from the protein surface. Examples include **nuclease** digestion (6) and chemical reagents, such as hydroxyl radical **footprinting** (7, 8). In both cases, the DNA is preferentially cleaved when away from the surface, leading to a periodicity in the cleavage pattern that is a direct measure of  $F$ .

### 2.2. Surface Twist

The second term that appears on the right in Equation (2) is the surface twist,  $STw$ . This can be calculated for DNA lying on any smooth surface with the integral (1)

$$STw = \frac{1}{2\pi} \int_C \mathbf{T}_A \cdot \mathbf{v} \times d\mathbf{v} \quad (4)$$

Reference should again be made to Figure 1. As before,  $\mathbf{T}_A$  is the unit tangent vector to  $A$  at a particular point  $a$ . The integral for  $STw$  differs from that for  $Tw$  by the use of  $\mathbf{v}$  instead of  $\mathbf{v}_{AC}$ . The vector  $\mathbf{v}$  is the surface normal at any point on the surface touched by the DNA axis and includes the required information about the local surface contour. The correspondence vector  $\mathbf{v}_{AC}$ , in contrast, contains only information about the correspondence between  $A$  and  $C$ , without regard to the surface chosen.

### 2.3. Example of the Nucleosome

Equation (4) can be used to calculate  $STw$  for any DNA wrapping surface of known geometry (9). For example, the nucleosome can be described as a cylindrical surface  $M$  of radius  $R$  on which the DNA axis  $A$  wraps  $n$  times as a left-handed helix of pitch  $2pp$  (where  $R$  is the sum of the radius of the histone octamer core and the radius of the DNA). The central line of the cylinder is labeled  $L$ . The unit normal vector  $\mathbf{v}$  to  $M$  points radially outward from the surface. The application of Equation (4) to this geometry gives the result (9)

$$STw = \frac{-np}{\sqrt{p^2 + r^2}} \quad (5)$$

The minus sign in this equation is a consequence of the left-handed chirality of the winding. In order to calculate the DNA twist from Equation (2), the winding number,  $F$ , must be determined independently. The experimental data are that  $n = 1.85$ ,  $2pp = 2.8$  nm, and  $R$  (the radius of the cylinder plus the radius of the DNA) = 5.3 nm. The calculated value of  $STw$  is then  $-0.16$ . The value



of  $F$  for the 146 bp of DNA wrapped on the nucleosome is 14.6 turns, based on the experimentally determined helical repeat of 10.0 bp/turn (10). The twist of the DNA about the nucleosome core is the sum of  $STw$  and  $F$ , giving  $Tw = 14.44$ . Had the DNA duplex been linearly extended in solution,  $F;_0 = Tw_0 = 146 / 10.54 = 13.85$ . For the nucleosome, therefore, the **surface wrapping** and accompanying curvature of the DNA axis increases the twist by about 4%.

### Bibliography

1. J. H. White and W. R. Bauer (1988) Proc. Natl. Acad. Sci. USA **85**, 772–776.
2. J. H. White and W. R. Bauer (1987) J. Mol. Biol. **195**, 205–213.
3. J. H. White (1969) Am. J. Math. **91**, 693–728.
4. J. H. White, N. R. Cozzarelli, and W. R. Bauer (1988) Science **241**, 323–327.
5. A. A. Travers and A. Klug (1990) in *DNA Topology and Its Biological Effects*, N. R. Cozzarelli and J. C. Wang, eds., Cold Spring Harbor Press, Cold Spring Harbor, NY, pp. 57–106.
6. A. Klug, L. C. Lutter, and D. Rhodes (1983) Cold Spring Harb. Symp. Quant. Biol. **47**, 285–292.
7. W. J. Dixon, J. J. Hayes, J. R. Levin, M. F. Weidner, B. A. Dombroski, and T. D. Tullius (1991) Meth. Enzymol. **208**, 380–413.
8. J. J. Hayes and A. P. Wolffe (1992) Trends Biochem. Sci. **17**, 250.
9. J. H. White and W. R. Bauer (1986) J. Mol. Biol. **189**, 329–341.
10. H. R. Drew and C. R. Calladine (1987) J. Mol. Biol. **195**, 143–173.

### Suggestions for Further Reading

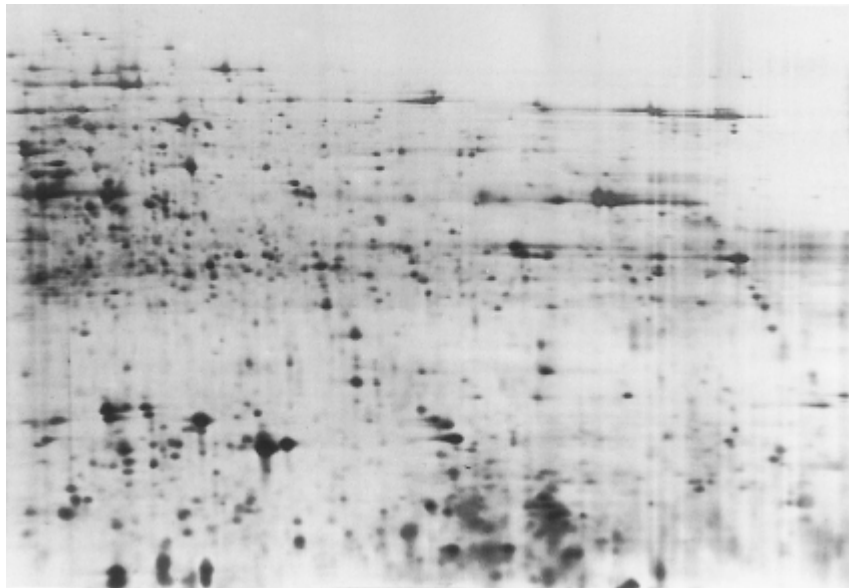
11. J. H. White and W. R. Bauer (1995). Chirality and surface twist of DNA wrapped on protein surfaces, J. Biomol. Struct. Dyn. **12**(4), 815–826. (This article presents sample calculations of the surface twist for several types of protein surfaces.)
12. J. H. White (1989). An introduction to the geometry and topology of DNA structure, in *Mathematical Methods for DNA Sequences*, M. S. Waterman, ed., CRC Press, Inc., Boca Raton, FL, pp. 225–253. (This article describes in detail the various methods for calculation of the twist.)

## Two-Dimensional Gel Electrophoresis

A two-dimensional gel is uniquely suited for analysis by [gel electrophoresis](#) of samples containing many components. The reason is simply geometric: A 10-cm lane covered with evenly spaced idealized bands may be able to resolve 50 such bands; a square of the same dimensions correspondingly can resolve  $50^2$ , or 2500, spots.

Size and net charge are the predominant properties by which proteins differ from one another, so that the first-dimensional gel commonly uses [isoelectric focusing](#) to provide a charge separation whereas the second-dimensional gel provides a size separation by means of [SDS-PAGE](#) (1). Such gels can resolve thousands of proteins, including those from an entire cell (Fig. 1). The position of each protein on the gel is defined by its [isoelectric point](#) ( $pI$ ) and its molecular weight ( $M_r$ ).

**Figure 1.** An example of a two-dimensional gel, separating alkaline yeast cell proteins. First dimension: isoelectric focusing in an Immobiline gel, pH 7–10, containing 0.5% carrier ampholytes, 8-M [urea](#), and 0.5% (v/v) NP-40; electrophoresis was for 90,300 V-h. Second dimension: SDS-PAGE, pore gradient gel of 12% to 15% (w/v) acrylamide, 4% Bisacrylamide, in a Tris-chloride-glycinate discontinuous buffer separating at pH 9.5. The proteins were detected by **silver staining**. [Fig. 4 of A. Goerg et al. (1988), *Electrophoresis* **9**, 37.]



Other combinations of separation techniques may be more suitable for specific applications. For example, [agarose](#) gel electrophoresis of **viruses** employs an initial low-concentration gel for charge separation and a more concentrated second gel for size separation ([2](#)).

With hundreds to thousands of spots, the reproducibility of spot positions on two-dimensional gels is a problem in comparing different gels. It has been solved in part through use of immobilized pH gradients for the isoelectric focusing; these help to produce constant positions for each protein at its *pI* at the steady-state in the first dimension ([3](#)). Densitometry and pattern displacement techniques further assist analysts in comparing different gels. To some degree, immunologic detection of spots after blotting has made it unnecessary to identify each spot from its position alone. The identities of proteins in spots can also be confirmed by their molecular weights, which are determined very accurately by [mass spectrometry](#). Nevertheless, even conventional two-dimensional electrophoresis using isoelectric focusing and SDS-PAGE, along with identification of the spots by their positions alone, has been successful in the construction of ever-growing protein databases ([4](#)).

#### Bibliography

1. M. G. Harrington (ed.) (1991) *Methods*. **3**, 71–141.
2. P. Serwer (1985) *Anal. Biochem.* **144**, 172–178.
3. M. J. Dunn (1987) *Adv. Electrophoresis* **1**, 1–110.
4. J. E. Celis (ed.) (1995) *Electrophoresis* **16**, 2175–2264.

#### Suggestion for Further Reading

5. M. J. Dunn, ed. (1995) Paper symposium: 2D Electrophoresis: From protein maps to genomes. *Electrophoresis* **16**, 1077–1326.

## Two-Hybrid Systems

(*In memoria*)

Maria Dimitrova

Michèle Granger-Schnarr

[Protein–Protein Interactions](#) play a major role in almost all relevant physiological processes occurring in living organisms, including [DNA replication](#) and **transcription, RNA splicing, protein biosynthesis, signal transduction**, and many other processes (1). Since the early 1990s, two-hybrid systems have emerged as a powerful new genetic tool for studying binary protein–protein interactions. These two-hybrid assays are generally conducted in the yeast *Saccharomyces cerevisiae* and use the activation of transcription of one or several [reporter genes](#) to detect the protein–protein interaction in question. In the past, these systems have been used to evaluate the interaction between defined proteins and as a method to identify novel partners for a known protein. As a result of the **sequencing** of several model organisms, however, and the predicted sequencing of the *Homo sapiens* [genome](#) in the near future, the two-hybrid approach will probably occupy an even larger place in the biologist's toolbox. The human genome is likely to encode about 50,000 to 100,000 **genes** (2), giving rise to an even greater number of hypothetical proteins due to [alternative splicing](#). This number is an order of magnitude greater than the 6,000 or so genes in *S. cerevisiae*. Unfortunately, the sequence of a gene does not always reveal its biological function, and even in the laboratory workhorse *S. cerevisiae* fewer than half of the genes could be classified as “functionally characterized” in 1997 (3, 4). It is likely that the situation will be even worse for mouse and human gene products, which frequently lack significant [homology](#) with proteins from simple model organisms and for which gene **knockout** studies are costly and time-consuming or unfeasible for ethical reasons.

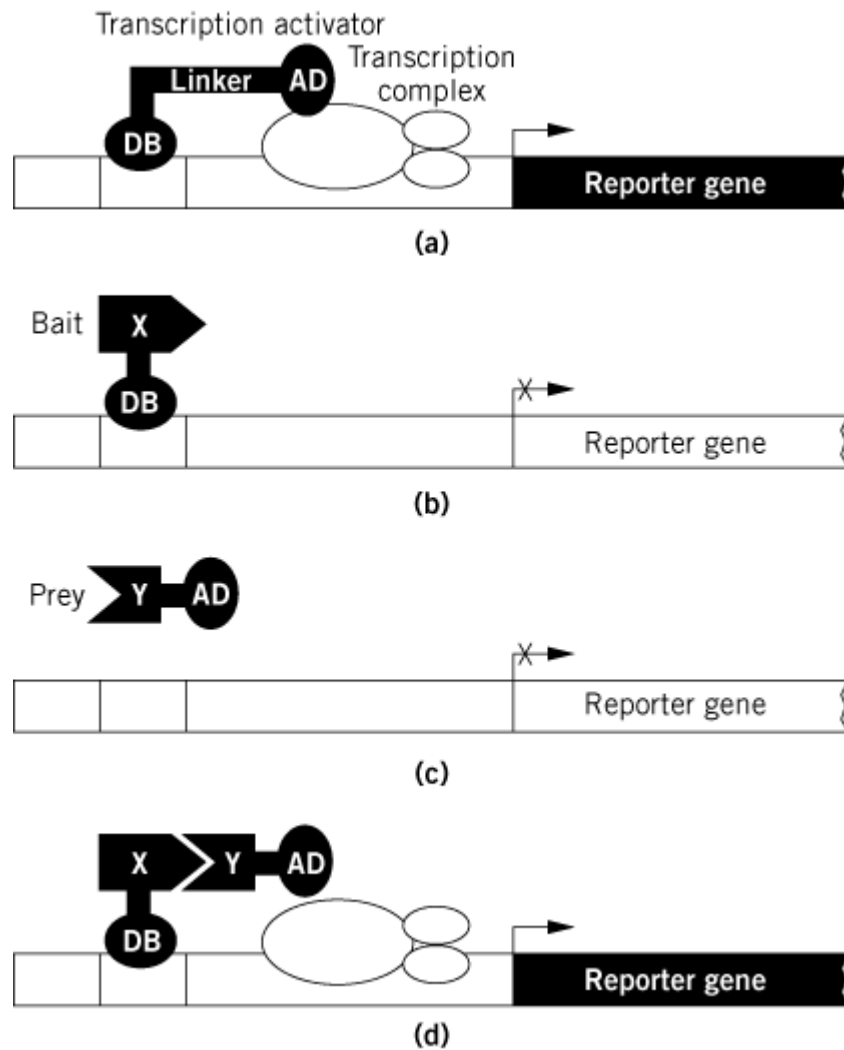
One potential way to classify a protein with unknown function will be to search systematically for its protein partners using two-hybrid methods. If one of these partners has already been functionally characterized, the protein with the yet unknown function is likely to belong to the corresponding regulatory or metabolic pathway.

### 0.1. History of the Yeast Two-Hybrid System Based on Transcriptional Activation

Most **eukaryotic** transcription activators are modular proteins that harbor at least a DNA-binding **domain** and an activation domain held together by a covalent linker (Fig. 1a). Those proteins, it is thought, trigger the assembly of the transcription preinitiation complex (5). The observation that transcription activators act at variable distances from the transcription start point suggested further that transcription activation in eukaryotes occurs via a mechanism with one or more flexible elements. Along these guidelines, Brent and Ptashne (6) showed that a eukaryotic activation domain fused to the DNA-binding domain of the bacterial [LexA repressor](#) functions as a transcriptional activator in yeast. Ma and Ptashne (7) went one step further in showing that transcription is activated even by binary protein complexes in which the DNA-binding and activation domains reside on separate polypeptide chains.

**Figure 1.** Principle of the yeast two-hybrid assay based on transcriptional activation. (a) A classical eukaryotic transcription activator binds upstream of the promoter via its DNA binding domain (DB) and activates transcription via its covalently linked activation domain (AD). (b) A bait protein X fused to DB generally does not activate transcription

because in most cases it lacks a functional activation domain. (c) A prey protein Y fused to AD also does not activate the reporter gene because this protein cannot bind to the promoter. (d) Interaction between X and Y fused respectively, to DB and AD reconstitutes a functional transcription activator and induces transcription of the reporter gene.



Independent work by Fields and Song (8) that gave rise to the first general two-hybrid assay to monitor protein–protein interactions in yeast. In this assay, a first protein of interest (X) is synthesized fused to a DNA-binding domain (DB). This first hybrid protein (DB-X) binds upstream of a reporter gene. If X is devoid of transcriptional activation capacity, binding of DB-X to the promoter does not activate transcription of the reporter gene (Fig. 1b). A second protein (Y) is synthesized fused to an activation domain (AD). Because this second hybrid protein (AD-Y) cannot generally bind to the promoter, it also fails to induce transcription of the reporter gene (Fig. 1c). However, if the two hybrid proteins (DB-X and AD-Y) are coexpressed and if X and Y interact physically, this interaction brings the activation domain close to the promoter and leads to activation of the reporter gene (Fig. 1d). Since 1989, this two-hybrid strategy has been used widely to study the interaction between known proteins and to isolate from activator domain-tagged cDNA libraries new proteins that interact with LexA or Gal4 fusion proteins used as DB-X “baits”. The advantage of this interaction cloning strategy is that the DNA coding for the “prey” protein Y is immediately available. In 1991, Fields and colleagues (9) published the first example of this kind of two-hybrid interaction cloning, using *lacZ* as a screenable reporter gene.

## 0.2. The Yeast Two-Hybrid System: From Screen to Selection

A significant improvement in the initial two-hybrid system was including additional reporter genes

whose expression is required for yeast growth (10-13). Reporters derived from the yeast *Leu2*, *His3*, or *URA3* genes transformed the two-hybrid system from a screen to a selection (see Table 1) allowing relatively rare clones to be isolated. The *URA3* reporter (13) has the additional advantage of using either direct selection for growth on uracil for productive protein–protein interactions, and it also makes possible negative selection by using 5-fluoro-orotic acid, which is toxic for yeast growth if the *URA3* gene is expressed. In this case, the *URA3* reporter may be used to probe the lack of interaction between the proteins X and Y. This is particularly useful for identifying a [mutation](#) in one or the other hybrid protein that disrupts the interaction between the two proteins. This kind of counterselection may also be useful for identifying reagents that inhibit specific protein interactions. One example of this “reverse two-hybrid” approach has been described by Vidal et al. (14).

**Table 1. Different Combinations of DNA Binding (DB) Moieties, Activation Domains (AD), and Reporter Genes Used in Classical Yeast Two-Hybrid Assays**

| DB   | AD   | Screen | Selection                 | Ref. |
|------|------|--------|---------------------------|------|
| GAL4 | GAL4 | lacZ   | –                         | (9)  |
| GAL4 | GAL4 | lacZ   | HIS3 (+3-AT) <sup>a</sup> | (10) |
| LexA | VP16 | lacZ   | HIS3                      | (11) |
| LexA | B42  | lacZ   | LEU2                      | (12) |
| ER   | VP16 | –      | URA3                      | (13) |

<sup>a</sup> 3-AT: 3-aminotriazole.

The various improved two-hybrid variant systems reported in 1993 differ further in the activation domains, using either the strong VP16 activation domain (11), the intermediate Gal4 domain (10), or the weak bacterial B42 domain (12). Weak activation domains potentially increase the spectrum of proteins recovered by avoiding the toxic effects that strong transcriptional activators may have in yeast (15). Most laboratories use either Gal4<sub>1-147</sub> or the full-length LexA repressor as a DNA-binding moiety. Gal4<sub>1-147</sub> corresponds to the amino-terminal domain of the yeast transcriptional activator Gal4. This part of the molecule contains specific DNA-binding and dimerization capacity (16). LexA is a bacterial repressor of 202 amino acid residues that harbors an amino-terminal, DNA-binding, and a carboxyl-terminal dimerization domain (17, 18). The most significant difference in the use of LexA or the Gal4 amino terminus in the yeast two-hybrid assay is that the Gal4 domain contains a nuclear localization signal (see [Lambda Phage](#)), whereas LexA does not (19). This leads to a lower nuclear concentration of LexA, which may be a problem if LexA fusion proteins are expressed from **vectors** that integrate into the genome, leading to a loss of sensitivity not found for integrating Gal4 vectors.

### 0.3. Correlation of Yeast Two-Hybrid Data with *in Vitro* Measurements

Despite the extensive use of the yeast two-hybrid technology by hundreds of different research groups for interaction cloning, little is known about (1) the equilibrium binding constant between proteins X and Y necessary to produce significant reporter gene expression and (2) the extent to which the degree of protein interaction determined by yeast two-hybrid methods correlates with the degree of interaction determined by biochemical techniques.

In regard to the first question, Estojak et al. (20) showed that the yeast two-hybrid assay may detect protein–protein interactions with a **dissociation constant** ( $K_d$ ) as high as ~1  $\mu$ M. These results were

obtained mainly with dimerization-defective versions of the [lambda phage](#) cI repressor and Myc, Max, and Mxi1 **helix-loop-helix** proteins. To extend the range of sensitivity and discrimination of the assay, the authors used different reporter **plasmids** that have one, two, and 8 LexA operators in front of the *lacZ* gene. With the one-operator *lacZ* reporter, the system was able to detect interactions with a predicted  $K_d$  ranging from 20 nM to nearly 1  $\mu$ M, whereas interactions with lower affinities of  $K_d > 1 \mu$ M were not detected. The results with two operators were comparable, but allowed more resolution of differences between weak interactions, plus detection of interactions in the  $K_d$  range of 1  $\mu$ M. Finally, the eight-operator reporter allowed clear detection of interactions in the 1  $\mu$ M  $K_d$  range but had high background levels of activity for many of the baits and showed compression of differences between strong and moderate interactions.

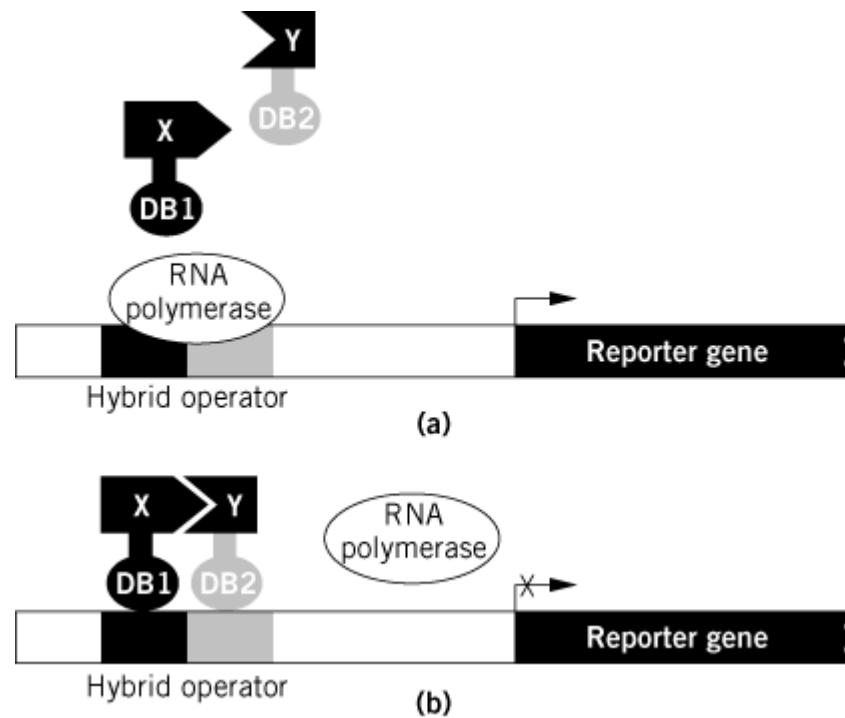
These observations lead immediately to the second question that involves attempts to correlate yeast two-hybrid data more quantitatively with relative binding affinities determined *in vitro*. At least in the case of the phage lambda cI variants, this kind of correlation is difficult to establish (20). A comparison of a lambda wild-type dimerization domain fusion ( $K_d = 20$ nM) with the partially dimerization-deficient Ala152Thr variant ( $K_d \sim 1 \mu$ M) shows that the relative affinities deduced from *in vivo* data would depend strongly on the reporter plasmid. Depending on the reporter used for these studies, formation of a Ala152Thr homodimer would be 200-fold less, seven-fold less, or even two-fold greater than the formation of a cI wild-type repressor homodimer. In this case, the data from the least sensitive *lacZ* reporter (that harbors a single LexA operator) would most closely parallel the *in vitro* data. Surprisingly, the *LEU2* reporters did not detect interactions for this set of lambda cI repressors. This may be due to the fact that in these experiments the *lacZ* reporters were carried on multicopy plasmids, whereas the *LEU2* reporters were integrated into the [chromosome](#).

Given these issues, the data of Estojak et al. (20) suggest that prudent use of a two-hybrid system to study protein interactions should involve assays with low- and high-sensitivity *lacZ* reporters and with at least one *LEU2* reporter. Further, putative interactions should be tested in both orientations, because the two-hybrid assay is subject to directionality for some protein pairs (20).

#### 0.4. Bacterial Two-Hybrid Systems Based on Transcriptional Repression

Given the difficulties encountered with the yeast two-hybrid system, it seemed worthwhile to develop in parallel transcription-based two-hybrid reporter systems in model organisms even simpler than yeast. The eubacterial transcriptional apparatus is much simpler than that of yeast and other eukaryotic organisms, and transcription in eubacteria may be efficiently inhibited by repressor molecules that simply prevent the **RNA polymerase** holoenzyme from interacting with a given promoter. Dimitrova et al. (21) recently developed a LexA-based two-hybrid assay in *Escherichia coli*, whose principle is outlined in Fig. 2. In this assay, a protein Y is fused to the LexA wild-type DNA-binding domain (DBD2) and a protein X is fused to an altered-specificity DNA-binding domain (DBD1) that harbors three point mutations in the LexA helix-turn-helix motif. These two hybrid proteins regulate the expression of a reporter gene placed under the control of a hybrid LexA **operator** that has a wild-type half-site and a mutated half-site recognized by the altered-specificity LexA variant.

**Figure 2.** Principle of a bacterial two-hybrid assay based on transcriptional repression. (a) Protein X is fused to an altered-specificity LexA DNA binding domain and protein Y to a wild-type DNA-binding domain comprising in both cases the 87 first amino acids of LexA. Without a dimerization domain, LexA<sub>1-87</sub> cannot act as a repressor. Thus, if X and Y do not interact, the reporter gene is transcribed. (b) If the proteins X and Y interact, the two hybrid-proteins bind efficiently to the hybrid operator and transcription of the reporter gene is repressed.



In this system, significant DNA binding *in vitro* and transcriptional repression of the *lacZ* reporter *in vivo* are observed only upon coexpression of the two fusion proteins DBD1-X and DBD2-Y (21). This assay allows specific detection of a complex between X and Y even if both X and Y can form homodimers. This possibility was not included in a lambda cI repressor-based system, where X and Y were both fused to a wild-type cI DNA binding domain (22), so that both heterodimerization and homodimerization may give rise to transcriptional repression.

The opposite limitation is inherent in bacterial systems, which use a dominant-negative **phenotype** to monitor heterodimer formation with repressor fusion proteins in *E. coli* (23-26). In this case, a DBD-X homodimer, which confers transcriptional repression, is challenged with a protein expressed either alone (23) or fused to a DNA-binding-deficient domain (24-26). If Y can dissociate an appreciable amount of the DBD-X homodimers, transcriptional repression is relieved, leading to a dominant-negative phenotype. Obviously, this strategy can be applied only if either X or Y can form homodimers. This kind of assay has been used for interaction cloning of a HMG-box protein, which binds to the helix-loop-helix domain of the c-Myc **oncprotein** (25). It remains to be seen if these bacterial two-hybrid systems becomes as popular as the yeast two-hybrid assay for studying known protein interactions and molecular cloning of unknown interaction partners for a known protein.

#### 0.5. Identification of Unknown Interaction Partners with Yeast Two-Hybrid Systems

Using yeast two-hybrid systems, an increasing number of previously unrecognized protein-protein interactions have been detected. Many of these were compiled by Golemis et al. (27). This and other reviews (28-31) describe in much detail the vectors, reporter strains, procedures, and possible problems of two-hybrid interaction cloning experiments. Obviously the bait protein DB-X must have little or no intrinsic ability to activate transcription on its own, must be expressed at reasonably high levels, and must be able to enter the yeast nucleus and bind DNA. Conversely, it has to be verified that the prey protein AD-Y does not activate transcription of the reporter gene on its own. But even in those cases where X and Y definitely interact in the yeast nucleus, this interaction is not necessarily biologically relevant, and its functionality has to be extensively tested *in vivo* and *in vitro*. For example, two heterologous proteins may interact specifically in the yeast nucleus but may not be expressed in the same cell type at the same time. Additionally, several classes of proteins have been isolated by using very different baits, including **heat-shock** proteins, **ribosomal** proteins,

cytochrome oxidase, **mitochondrial** proteins, and [proteasome](#) subunits. Because these proteins generally have no obvious biological relationship to the bait protein, these interactions are generally classified as false positives. Another potential problem of the two-hybrid approach is that a detected interaction supposed to be binary may depend in fact on one or more yeast proteins, which may serve as either bridge or connecting proteins between X and Y, or simply stabilize an existing physical interaction between X and Y.

This potential ability of the two-hybrid assay to detect ternary protein complexes has been intentionally used by several groups. Van Aelst et al. (32) demonstrated, for example, that the signal transduction proteins Ras and MEK interact only in the two-hybrid assay if the Raf protein is overexpressed at the same time. Legrain and Chapon (33) showed that the yeast splicing factors PRP9 and PRP11 do not interact directly but that these two proteins can bind simultaneously to a third protein (SPP91, now called PRP21) to form a three-molecule complex.

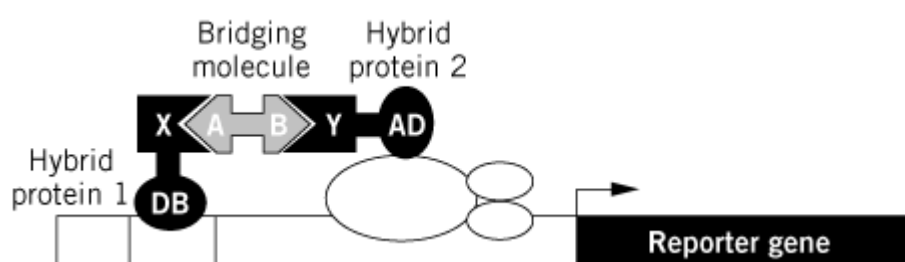
Protein–protein interactions sometimes depend on the **posttranslational modification** of one component of the complex. In particular, protein **phosphorylation** plays a major role in numerous signal transduction pathways. However, yeast cells generally do not employ, for example, tyrosine phosphorylation as a major regulatory modification. To overcome this problem, Osborne et al. (34, 35) coexpressed a **tyrosine kinase**, a cytosolic immunoreceptor domain (that contains phosphorylatable tyrosine residues), and an **SH2 domain**. As expected, the immunoreceptor domain and the SH2 domain interacted in the yeast two-hybrid assay only upon coexpression of a functional tyrosine kinase (34).

To develop a more general system for identifying of a third bridge or connecting protein, Xu et al. (36) developed a two-bait system in which transcription of selectable and counterselectable reporters is directed by different baits. The *URA3* reporter is placed under the control of a tetracycline repressor fusion protein, and the *lacZ* gene is controlled by a LexA repressor fusion. The prey protein (or a prey protein library) is fused as usual to the B42 activation domain.

#### 0.6. Systems that Detect Interactions Dependent on a Third Nonprotein Molecule

Figure 3 shows the rationale of this kind of experiment. The two DB-X and AD-Y fusion proteins do not interact directly but are connected by a covalently linked bifunctional molecule for which the A-moiety interacts with X and the B-moiety with Y. Licitra and Liu (37) used this strategy for a small synthetic ligand, where A was [dexamethasone](#), which readily interacts with the [glucocorticoid](#) receptor (which was fused to LexA), and B was the small molecule [Development](#), which interacts with the FKBP12 protein (fused to the B42 activation domain). As outlined by the authors, the dexamethasone/FK506 hybrid molecule is **hydrophobic** and penetrates readily into yeast cells. This is generally not the case for hydrophilic or charged ligands. To generalize this approach, it is important to develop yeast strains that are more permeable and do not significantly affect yeast viability.

**Figure 3.** Principle of the yeast three-hybrid approach. The two hybrid proteins DB-X and AD-Y do not interact directly as in Fig. 1d, but this interaction is mediated by a third bifunctional bridging molecule. This bridging molecule may be a synthetic organic molecule or a hybrid RNA molecule that has an A moiety that interacts with protein X and a B moiety that is prey for the search of a RNA-binding protein Y.





A–protein interactions are pivotal in fundamental cellular processes, such as **translation**, **messenger RNA** processing, early [development](#), and infection by RNA **viruses**. In spite of the central importance of these interactions, however, few approaches are available to analyze them rapidly *in vivo*. In the three-hybrid approach, the A-moiety in Fig. 3 corresponds to a RNA **stem-loop structure** that interacts with the well-characterized, sequence-specific, RNA-binding coat protein from **bacteriophage MS2**. For the covalently linked B-moiety, the authors used either iron [response element](#) RNA, which interacts with the iron regulatory protein 1, or HIV **transactivation response element** RNA, which interacts with the HIV transactivator protein Tat. Subsequently the same group used this approach to identify an RNA-binding protein that interacts with the 3' untranslated region of a mRNA involved in the regulation of sexual fates in *Caenorhabditis elegans* hermaphrodite germ lines ([39](#)).

#### 0.7. Interaction Mating and Functional Genomics

A major question for most biologists in the approaching post-genome-sequencing era is how to handle the protein sequence data generated by systematic DNA sequencing projects. There will be an enormous amount of protein primary structure information, which must be converted into biologically more relevant information, that is, we have to understand how a given protein is involved in one, or potentially several, regulatory or metabolic pathways. Extrapolating from the experience with the *S. cerevisiae* genome ([3](#), [4](#)), the function of fewer than half of the proteins will be inferred from their primary structure by sequence homology with functionally characterized proteins. For the remaining proteins, new concepts and techniques will have to be developed, especially for those higher eukaryotes without developed genetics. One way to proceed will be establishing genome-wide, two-hybrid information from systematic application of yeast interaction mating experiments ([40-43](#)). [DNA Libraries](#) yeast cells exist in two **mating types** that mate with one another and form diploids if they come into contact. By mating haploid strains of opposite mating types that contain potential interacting proteins, exhaustive and rapid screens of cDNA and genomic [DNA libraries](#) may be conducted.

Fromont-Racine et al. ([43](#)) used this strategy to establish a partial yeast protein interaction map, starting from ten protein baits that are involved in pre-mRNA splicing. Several identified prey proteins were used subsequently as new baits for further rounds of interaction cloning experiments. Conceivably, this strategy may be extended to establish protein linkage maps for whole genomes if most of the work is done by laboratory robots. Based on the nature of the interacting partners, such protein linkage maps will suggest functions for new proteins and lead to hypotheses as to the role of a protein in a special regulatory or metabolic pathway. Given that so many proteins function by interacting with other proteins ([44](#)), in principle two-hybrid approaches can provide useful information for a significant fraction of the **proteome** of any given organism.

#### Bibliography

1. E. M. Phizicky and S. Fields (1995) *Microbiol. Rev.* **59**, 94–123.
2. G. D. Schuler et al. (1996) *Science* **274**, 540–546.
3. A. Goffeau et al. (1996) *Science* **274**, 546–567.
4. B. Dujon (1996) *Trends Genet.* **12**, 263–270.
5. M. Ptashne and A. Gann (1997) *Nature* **386**, 569–577.
6. R. Brent and M. Ptashne (1985) *Cell* **43**, 729–736.
7. J. Ma and M. Ptashne (1988) *Cell* **55**, 443–446.
8. S. Fields and O. K. Song (1989) *Nature* **340**, 245–246.
9. C. T. Chien, P. L. Bartel, R. Sternglanz, and S. Fields (1991) *Proc. Natl. Acad. Sci. USA* **88**, 9578–9582.
10. T. Durfee et al. (1993) *Genes and Dev.* **7**, 555–569.

11. A. B. Vojtek, S. M. Hollenberg, and J. A. Cooper (1993) *Cell* **74**, 205–214.
12. J. Gyuris, E. A. Golemis, H. Chertkov, and R. Brent (1993) *Cell* **75**, 791–803.
13. B. Le Douarin et al. (1995) *Nucleic Acids Res.* **23**, 876–878.
14. M. Vidal et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 10315–10320.
15. G. Gill and M. Ptashne (1988) *Nature* **334**, 721–724.
16. R. Marmorstein, M. Carey, M. Ptashne, and S. C. Harrison (1992) *Nature* **356**, 408–414.
17. M. Schnarr and M. Granger-Schnarr (1993) In *Nucleic Acids and Molecular Biology* (F. Eckstein & D. M. J. Lilley, eds.), Springer-Verlag, Berlin, Vol. **7**, pp. 170–189.
18. P. Dumoulin, R. H. Ebright, R. Knechtel, R. Kaptein, M. Granger-Schnarr, and M. Schnarr (1996) *Biochemistry* **35**, 4279–4286.
19. P. A. Silver, R. Brent, and M. Ptashne (1986) *Mol. Cell. Biol.* **6**, 4763–4766.
20. J. Estojak, R. Brent, and E. A. Golemis (1995) *Mol. Cell. Biol.* **15**, 5820–5829.
21. M. Dimitrova et al. (1998) *Mol. Gen. Genet.* **257**, 205–212.
22. R. Jappelli and S. Brenner (1996) *J. Mol. Biol.* **259**, 575–578.
23. J. K. et al. (1995) *Genes Dev.* **9**, 2986–2996.
24. A. Marchetti, M. Abril-Marti, B. Illi, G. Cesareni, and S. Nasi (1995) *J. Mol. Biol.* **248**, 541–550.
25. C. A. Bunker and R. E. Kingston (1995) *Nucleic Acids Res.* **23**, 269–276.
26. X. Zeng, A. M. Herndon, and J. C. Hu (1997) *Proc. Natl. Acad. Sci. USA* **94**, 3673–3678.
27. E. A. Golemis, I. Serebriiskii and S. F. Law (1997) In *Gene Cloning and Analysis: Current Innovations* (B. C. Schaefer, ed.), pp. 11–28.
28. P. L. Bartel and S. Fields (1995) *Methods Enzymol.* **254**, 241–263.
29. C. Bai and S. J. Elledge (1996) *Methods Enzymol.* **273**, 331–347.
30. E. A. Golemis, I. Serebriiskii, J. Gyuris and R. Brent (1997) In *Current Protocols in Molecular Biology*, (F. Ausubel, R. Brent, R. Kingston, D. Moore, J. Seidmann, J. A. Smith, and K. Struhl, eds.), Wiley, New York, Unit 20.1.1–20.1.35.
31. R. Brent and R. L. Finley (1997) *Annu. Rev. Genet.* **31**, 663–704.
32. L. V. Van Aelst, M. Barr, S. Marcus, A. Polyverino, and M. Wigler (1993). *Proc. Natl. Acad. Sci. USA* **90**, 6213–6217.
33. P. Legrain and C. Chapon (1993) *Science* **262**, 108–110.
34. M. A. Osborne, S. Dalton, and J. P. Kochan (1995) *Bio/Technology* **13**, 1474–1478.
35. M. A. Osborne, M. Lubinus, and J. P. Kochan (1997) In *The Yeast Two-Hybrid System* (S. Fields and P. Bartel, eds.), Oxford University Press, Oxford, pp. 233–258.
36. C. W. Xu, A. R. Mendelsohn, and R. Brent (1997) *Proc. Natl. Acad. Sci. USA* **94**, 12473–12478.
37. E. J. Licitra and J. O. Liu (1996) *Proc. Natl. Acad. Sci. USA* **93**, 12817–12821.
38. D. J. SenGupta et al. (1996) *Proc. Natl. Acad. Sci. USA* **93**, 8496–8501.
39. B. Zhang et al. (1997) *Nature* **390**, 477–484.
40. R. L. Finley and R. Brent (1994) *Proc. Natl. Acad. Sci. USA* **91**, 12980–12984.
41. C. Bendixen, S. Gangloff, and R. Rothstein (1994) *Nucleic Acids Res.* **22**, 1778–1779.
42. P. L. Bartel, J. A. Roecklein, D. SenGupta, and S. Fields (1996) *Nat. Genet.* **12**, 72–77.
43. M. Fromont-Racine, J. C. Rain, and P. Legrain (1997) *Nat. Genet.* **16**, 277–282.
44. B. Alberts (1998) *Cell* **92**, 291–294.

### Suggestions for Further Reading

45. S. Fields and P. Bartel (eds.) (1997) *The Yeast Two-Hybrid System*, Oxford University Press,

Oxford, pp. 1–333.

46. E. A. Golemis, I. Serebriiskii, and S. F. Law (1997) "Adjustment of parameters in the yeast two-hybrid system, In" *Gene Cloning and Analysis: Current Innovations* (B. C. Schaefer, ed.), pp. 11–28.
47. R. Brent and R. L. Finley (1997) Understanding gene and allele function with two-hybrid methods, *Annu. Rev. Genet.* **31**, 663–704.

## Ty Elements

Ty elements are [transposable elements](#) that were originally found in *Saccharomyces cerevisiae* (1). Tys are [retrotransposons](#), resembling [retroviruses](#) in many ways, although Tys are not infectious nor do they have an extracellular phase. Such retrotransposons are widespread and fall into two major classes; both are similar to retroviruses, but differ most obviously in gene order. The two most extensively studied yeast Tys are Ty1 and Ty3, and these elements represent the two classes of retrotransposons. Ty1 and the similar *Drosophila* retrotransposable element copia are the prototypes of one class, while Ty3 and the *Drosophila* retrotransposable element gypsy are prototypes for the other.

As is true with retroviruses, retrotransposons insert into new target sites by making the actual mobile DNA segment by **reverse transcription** of an RNA copy of the element. This is in contrast to other elements that move strictly through DNA intermediates; that is, a DNA segment is excised from the donor chromosome and then inserted into a target site, in steps that involve only DNA breakage and joining events. However, it is important to realize that the chemistry of the DNA processing steps are the same for DNA-only elements and for those that use a reverse transcription-generated copy of the element (2). Indeed, the [transposases](#) that mediate these reactions two types of transposition reactions are structurally related (3, 4).

The Ty1 element was the first nonviral element shown to transpose through an RNA intermediate (5) and has been studied extensively (1). In cells, Ty1 is found in large complex particles that resemble intracellular forms of retroviruses. Isolation of these Ty1 virus-like particles has allowed *in vitro* analysis of transposition. The copy number of Ty1 elements varies from yeast strain to yeast strain. In plants, related elements are often present in very high copy number and are heterogeneous.

An interesting feature of yeast Ty1 is that this element prefers to insert into pol III **promoter** regions (6), which may be a useful strategy to protect the genome from insertions into genes themselves. Although Ty1 transposition can be done *in vitro* using purified virus-like particles, this preference for promoter regions is not observed (7). It is likely that Ty1 is attracted to promoter regions through interaction with some component of the [transcription](#) machinery and that this component is lacking from the *in vitro* system. Another yeast retrotransposon, Ty3, inserts highly specifically into Pol II promoters; in this case, it has been shown that this targeting is mediated in the presence of the [transcription factors](#) that interact with this region of the promoter (8).

## Bibliography

1. J. D. Boeke and J. P. Stoye (1997) In *Retroviruses* (H. Varmus, S. Hughes, and J. Coffin, ed.), Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, pp. 343–435.
2. K. Mizuuchi (1992) *J. Biol. Chem.* **267**, 21273–21276.

3. F. Dyda, A. B. Hickman, T. M. Jenkins, A. Engelman, R. Craigie, and D. R. Davies (1994) *Science* **266**, 1981–1986.
4. P. Rice and K. Mizuuchi (1995) *Cell* **82**, 209–220.
5. J. D. Boeke, D. J. Garfinkel, C. A. Styles, and G. R. Fink (1985) *Cell* **40**, 491–500.
6. S. E. Devine and J. D. Boeke (1996) *Genes Dev.* **10**, 620–633.
7. S. E. Devine and J. D. Boeke (1994) *Nucleic Acids Res.* **22**, 3765–3772.
8. D. L. Chalker and S. B. Sandmeyer (1993) *Proc. Natl. Acad. Sci. USA.* **90**, 4927–4931.

## Type II DNA-Binding Proteins

The bacterial type-II [DNA-binding proteins](#) include the abundant HU and integration host factor (IHF) proteins and the phage-encoded TF1 and are believed to condense their cognate genomes by binding at multiple sites and inducing coherent bends. HU and TF1 possess little, if any, binding sequence specificity, but IHF which is also required for site-specific recombination, [DNA replication](#), and [transcription](#), binds at specific sites characterized by a limited consensus sequence.

All [proteins](#) of this type bind to DNA as dimers and bend the DNA to a substantial extent. In particular IHF induces a bend angle of at least  $160^\circ$  and possibly in excess of  $180^\circ$  within the  $2\frac{1}{2}$  double helical turns that comprises its binding site. IHF is normally a heterodimer; HU can exist either as a homo- or as a hetero-dimer. In both IHF and HU, the two  $\sim 10$  kDa subunits intertwine to form a compact core, from which two long  $\beta$  ribbon arms extend. These arms track along the minor groove from the inside to the outside of the wrapped DNA, where they terminate at the two substantial kinks. In addition to these interactions via the  $\beta$  arms, IHF also clamps the hairpin by minor-groove contacts with a helices from both subunits in the core of the dimer. All the contacts to the DNA are either in the minor groove or part of an extensive network of electrostatic interactions with the phosphate backbone.

The two kinks are induced by the partial intercalation between adjacent basepairs of an absolutely conserved proline side chain located at the tips of the  $\beta$  arms. In the IHF–DNA complex, the DNA bend is maintained by two distinct mechanisms. On the outside of the bend the hydrophobic intercalation stabilizes the opening of the minor groove. On the inside charge neutralization counteracts the enhanced repulsion between the phosphates on opposite sides of the narrowed grooves.

The sequence dependency of IHF binding is determined by the conformation of DNA rather than by base-specific contacts. The “consensus” sequence consists of two short elements separated by approximately half a turn in only one of the two half-sites. The conserved sequence, CAA, at the kink site, can accommodate the severe distortion induced by the protein better than other short sequences, and so is favored.

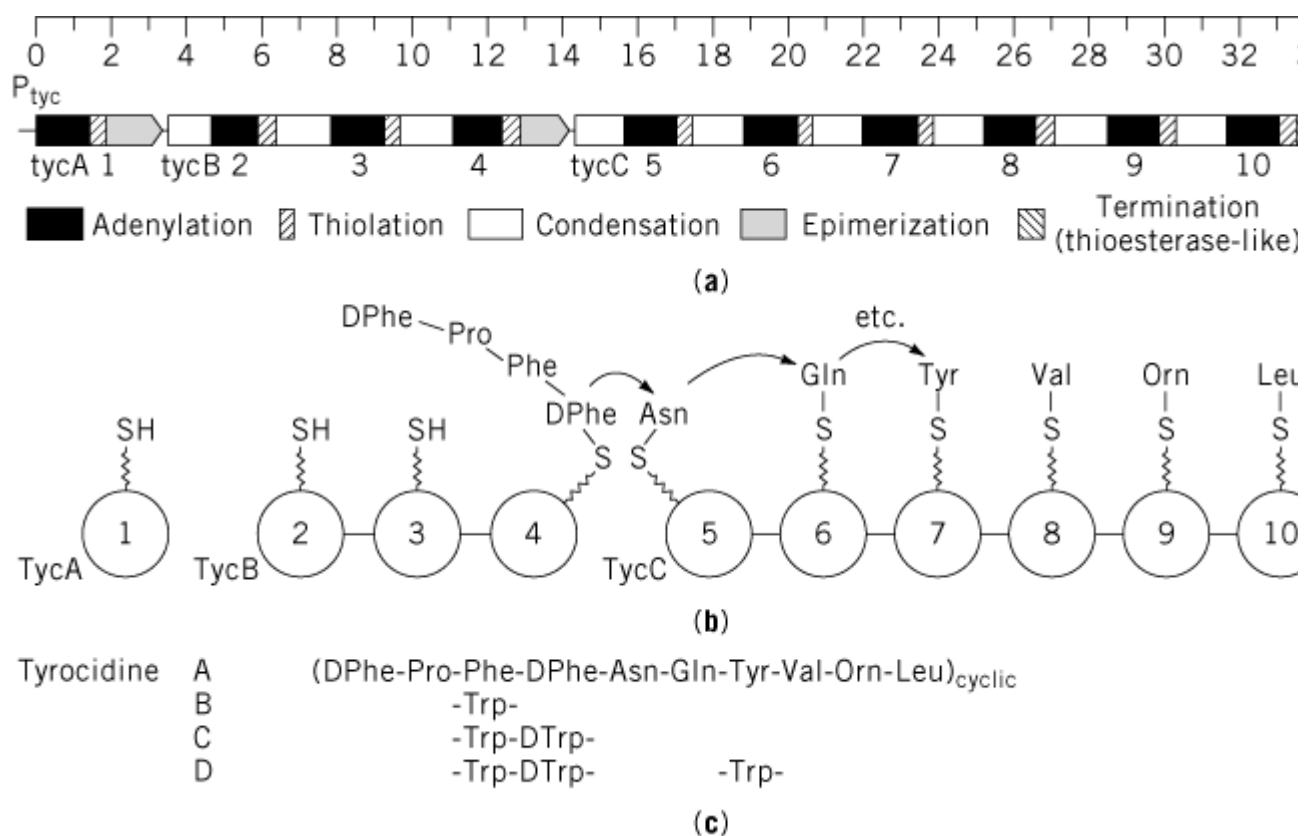
## Tyrosidine

Tyrocidine is cyclic decapeptide antibiotic and secondary metabolite that is synthesized on three nonribosomal peptide synthetases by the so-called thiotemplate mechanism. It was one of the first examples of nonribosomal peptide antibiotic biosynthesis and is best investigated regarding its genetic analysis and possible biological function for the producer *Bacillus brevis*. Tyrocidine possesses antibacterial properties, notably against many Gram-positive bacteria. However, it is considered today to be too toxic for systemic use.

## 1. Structure

Tyrocidine is a mixture of at least four known structural variants (Fig. 1, C). Tyrocidine A, the most prominent of these, is a cyclic decapeptide with the primary structure (-DPhe-Pro-Phe-DPhe-Asn-Gln-Tyr-Val-Orn-Leu-)<sub>cyclic</sub>, where the indicated amino acids are the unusual D-isomer and Orn is the unusual amino acid ornithine. In tyrocidines B, C, and D, the aromatic residues at positions three, four, and seven are gradually replaced by tryptophan (Fig. 1, C). The biosynthesis of a mixture of products is a consequence of the low specificities of the protein templates of the nonribosomal thiotemplate system, compared with the ribosomal system of protein biosynthesis (1, 2). The relative concentrations of the aromatic amino acids in the growth medium affect the ratios of the four variants (3). Although one should refer to tyrocidine as a mixture, the singular will be used here for reasons of simplicity.

**Figure 1.** The genes (a) and proteins (b) for the biosynthesis of the tyrocidines (c) are shown. (a): The biosynthetic gene *tycB*, and *tycC* are transcribed in one operon. The modular organization typical of nonribosomal peptide synthetases is as shown. The modules can be subdivided into domains responsible for the single chemical reactions: adenylation (black), thiolation (stripes), condensation (white), epimerization (grey), and thioesterase domains (coarse stripes), which were defined by analysis of the genes and by biochemical studies of the encoded proteins. (b): Tyrocidine is synthesized on an enzyme called the tyrocidine synthetases TycA, TycB, and TycC. Each module is represented by a sphere, and each incorporates one amino acid into the growing peptide chain. The zigzagged lines represent the prosthetic group 4'-phosphopantetheine. (c): The primary structures of the four known cyclic decapeptide tyrocidines. Amino acid substitutions at positions three, four, and seven are the result of the relatively modest substrate specificities of the corresponding modules.



Tyrocidine is produced by certain strains of the soil bacterium *B. brevis*, such as *B. brevis* ATCC 8185 and ATCC 10068. These strains produce in fact a mixture of two antibiotics, which is called *tyrothricin*, consisting of tyrocidine and the linear pentadecapeptide gramicidin.

## 2. Activity and Mode of Action

The combination of basic groups and nonpolar side chains gives tyrocidine the properties of a cationic detergent and a strong hemolytic activity. Tyrocidine is bactericidal for many Gram-positive organisms, much less so for some Gram-negative bacteria. With its biosurfactant property, it disrupts the integrity of the bacterial cell surface, resulting in the loss of amino acids, pyrimidine and purine bases, and other vital small molecules from the susceptible organism (4). It uncouples oxidative phosphorylation and in eukaryotes it disturbs the function of certain intracellular organelles, such as the mitochondria. All four known variants of tyrocidine display similar activities. Tyrocidine is not absorbed orally and undergoes degradation in the gut, it is primarily utilized as topical agent. The side effects, however, make it today inappropriate for systemic use (5).

## 3. Possible Biological Function

Biosynthesis of the two peptide antibiotics tyrocidine and linear gramicidin is switched on at the onset of stationary growth and shortly preceding the sporulation process. This close correlation in time has always raised speculation about the possible functions of these compounds in the producer cell. The finding that tyrocidine accumulates in the forespore and can also be found in the spore suggested that they have roles in sporulation. Some *in vitro* and *in vivo* experiments indeed pointed to such function, but others were apparently contradictory. Tyrocidine was shown to bind to DNA *in vitro* (6) and to inhibit RNA polymerase isolated from exponentially growing cells. Interestingly, the addition of linear gramicidin had a counteracting effect, suggesting antagonistic regulatory roles of the two antibiotics (7). In contrast to these *in vitro* results, tyrocidine was demonstrated to stimulate RNA synthesis under certain conditions *in vivo*. When exponentially growing cells were shifted to a medium lacking a nitrogen source, the addition of tyrocidine was essential for effective sporulation (8). A mutant strain defective in tyrocidine synthesis was reported to form spores that were less heat-resistant than those formed by wild-type cells. To date, however, the possible roles of tyrocidine remain unclear. No direct target or binding site has been identified. The inaccessibility of the producer strains to genetic manipulation has hindered further investigations.

In general, the biological functions of peptide antibiotics synthesized by the thiotemplate mechanism are obscure (9). Genetic studies on the cascade leading to sporulation in related bacteria, such as *Bacillus subtilis*, do not suggest a necessity for such compounds for regulatory functions. The most obvious explanation for their production is still as a defense weapon in the competition for limited nutrients.

## 4. Biosynthesis

The biosynthesis of tyrocidine proceeds by the thiotemplate mechanism. As one of the best investigated systems of this kind, tyrocidine attracts most of the current interest. According to that mechanism, large multifunctional enzymes, the NRPSs, carry out all steps required to activate, modify and link the amino acids to give the product (1, 2, 10). In some cases other than tyrocidine, further enzymes can be involved in product modification to yield the biologically active molecule. The NRPSs resemble in architecture, an assembly line. They contain one so-called module for each amino acid to be incorporated. It first activates the residue as its aminoacyl adenylate, with the consumption of ATP, then transfers it to the thiol moiety of an individual 4'-phosphopantetheine prosthetic group, which serves as a swinging arm to facilitate transport of the amino acid (or the growing peptide chain) to the next module. The modules can be subdivided into domains that perform the individual chemical reactions (1), for example, the core domains adenylation for

recognition and activation of the substrate monomer thiolation (or peptidyl carrier protein) for binding it as thioester on the 4'-phosphopantetheine, and condensation for formation of the peptide-bond between acyl groups on adjacent modules (Fig. 1). Optional domains like epimerization modify the amino acids upon incorporation into the growing chain.

In the case of tyrocidine, the functional enzyme complex consists of three such NRPSs, namely, TycA, TycB, and TycC (11). TycA is a one-module enzyme (molecular weight 123 kDa) that activates and epimerizes L-Phe. It initiates tyrocidine biosynthesis by transferring the D-Phe onto TycB, where the first peptide bond is formed between D-Phe and L-Pro. TycB (405 kDa) comprises three modules and activates L-Pro, L-Phe, and a second L-Phe, which is also epimerized to the D-enantiomer. TycB catalyzes the stepwise growth of the peptide chain to the tetrapeptide. The tetrapeptide is then transferred to TycC (724 kDa), which incorporates in the same manner the residual six amino acids and then cyclizes head-to-tail the linear decapeptide intermediate, to generate the final product (12) (Fig. 1 B).

The obscure biological function of many antibiotic products of secondary metabolism is in striking contrast to the size of their biosynthetic clusters, which often make up several percent of the genome. The genes *tycA*, *tycB*, and *tycC* encoding the tyrocidine synthetases have been cloned and sequenced (13). They are transcribed in an operon that spans over 35 kbp (Fig. 1 A). The promoter of this operon was shown to be functional in *B. subtilis*, where it is negatively regulated by the DNA-binding protein AbrB during vegetative growth (14). At the end of vegetative growth, AbrB derepresses a number of genes associated with stationary growth and sporulation. Sequence analysis of the tyrocidine NRPSs genes has revealed their typical domain structures (Fig. 1 A). TycC, which activates six amino acids and has six modules, is one of the largest enzymes known. Expression as recombinant protein fragments has allowed many biochemical analyses. It was shown that the adenylation domains that correspond to the variable positions in tyrocidine (Fig. 1 A and C) have only modest substrate specificities. In contrast, the adenylation domains corresponding to the invariant positions have been found to be specific for single amino acids (13). These results explain the naturally occurring variants of tyrocidine. The construction of hybrid NRPSs by fusing gene fragments encoding modules of TycB and TycC in an artificial order allowed the production of the predicted new peptides by the purified enzymes (15). These results indicate the potential of generating new tailored peptide antibiotics by combinatorial biosynthesis. The thioesterase domain (28 kDa) of TycC was excised and shown by using decapeptide-thioesters as substrate substitutes to act as the macrolactamase in the cyclization of the linear decapeptidyl-S-enzyme precursor (16). Directly adjacent to the tyrocidine biosynthesis genes are two genes that encode ABC (ATP-binding cassette) transporters (13). These transporters may confer resistance of the cell against its own product.

## Bibliography

1. M. A. Marahiel, T. Stachelhaus, and H. D. Mootz (1997) *Chem. Rev.* **97**, 2651–2673.
2. H. von Döhren, U. Keller, J. Vater, and R. Zocher (1997) *Chem. Rev.* **97**, 2675–2705.
3. M. A. Ruttenberg and B. Mach (1966) *Biochemistry* **5**, 2864–2869.
4. S. L. Dax (1997) In *Antibacterial Chemotherapeutic Agents*, Blackie Academic & Professional, London, UK, pp. 351–356.
5. F. E. J. Hunter and L. S. Schwartz (1967) In *Antibiotics I (Mode of Action)*, (D. Gottlieb and P. D. Shaw, eds.) Springer-Verlag, New York, NY pp. 142–152.
6. B. Schazschneider, H. Ristow, and H. Kleinkauf (1974) *Nature* **249**, 757–759.
7. H. Ristow, J. Russo, E. Stochaj, and H. Paulus (1982) In *Peptide Antibiotics*, (H. Kleinkauf and H. von Döhren, eds.), W. de Gruyter, Berlin, pp. 381–388.
8. H. Ristow, W. Pschorn, J. Hansen, and U. Winkel (1979) *Nature* **280**, 165–166.
9. E. Katz and A. L. Demain (1977) *Bacteriol. Rev.* **41**, 449–474.
10. D. E. Cane, C. T. Walsh, and C. Khosla (1998) *Science* **282**, 63–68.

11. S. G. Lee and F. Lipmann (1975) *Methods Enzymol.* **43**, 585–602.
12. R. Roskoski, Jr., H. Kleinkauf, W. Gevers, and F. Lipmann (1970) *Biochemistry* **9**, 4846–4851.
13. H. D. Mootz and M. A. Marahiel (1997) *J. Bacteriol.* **179**, 6843–6850.
14. M. A. Marahiel, M. M. Nakano, and P. Zuber (1993) *Mol. Microbiol.* **7**, 631–636.
15. H. D. Mootz, D. Schwarzer, and M. A. Marahiel (2000) *Proc. Natl. Acad. Sci. USA* **97**, 5848–5853.
16. J. W. Trauger, R. M. Kohli, H. D. Mootz, M. A. Marahiel, and C. T. Walsh (2000) *Nature* **407**, 215–218.

### Suggestions for Further Reading

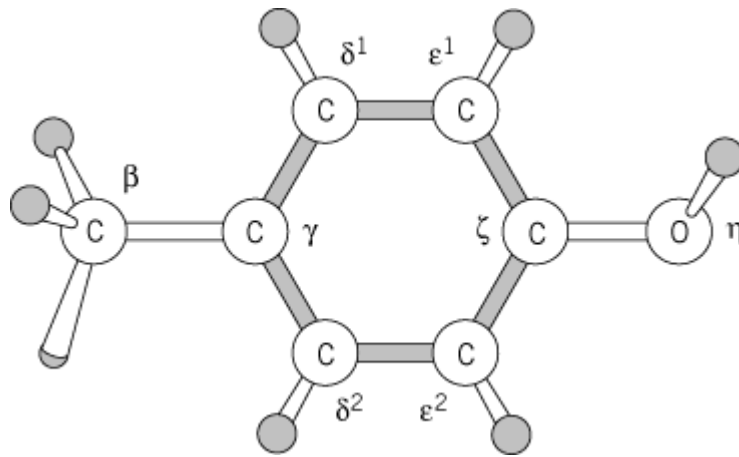
17. F. Lipmann (1971) Attempts to map a process evolution of peptide biosynthesis. *Science* **173**, 875–884.
18. H.-J. Rehm and C. Reed, eds. (1997) *Products of Secondary Metabolism*, Vol. 7 of *Biotechnology* (H. Kleinkauf and H. von Döhren, eds.), VCH, Weinheim.
19. T. Stein, et al. (1996) The multiple carrier model of nonribosomal peptide biosynthesis at modular multienzymatic templates. *J. Biol. Chem.* **271**, 15428–15435.
20. R. H. Lambalot, et al. (1996) A new enzyme superfamily—the phosphopantetheinyl transferases. *Chem. Biol.* **3**, 923–936.
21. E. Conti, T. Stachelhaus, M. A. Marahiel, and P. Brick (1997) Structural basis for the activation of phenylalanine in the nonribosomal biosynthesis of gramicidin S. *EMBO J.* **16**, 4174–4183.
22. T. Stachelhaus, A. Schneider, and M. A. Marahiel (1995) Rational design of peptide antibiotics by targeted replacement of bacterial and fungal domains. *Science* **269**, 69–72.
23. P. J. Belshaw, C. T. Walsh, and T. Stachelhaus (1999) Aminoacyl-CoAs as probes of condensation domain selectivity in nonribosomal peptide synthesis. *Science* **284**, 486–489.
24. T. A. Keating, C. G. Marshall, and C. T. Walsh (2000) Reconstitution and characterization of the *Vibrio cholerae* vibriobactin synthetase from VibB, VibE, VibF, and VibH. *Biochemistry* **39**, 15522–15530.
25. J. W. Trauger, R. M. Kohli, and C. T. Walsh (2001) Cyclization of backbone-substituted peptides catalyzed by the thioesterase domain from the tyrocidine nonribosomal peptide synthetase. *Biochemistry* **40**, 7092–7098.

## Tyrosine (Tyr, Y)

The [amino acid](#) tyrosine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to two **codons**—UAU and UAC—and represents approximately 3.2% of the residues in the proteins that have been characterized. The tyrosinyl residue incorporated has a mass of 163.18 Da, a **van der Waals volume** of 141 Å<sup>3</sup>, and an [accessible surface](#) area of 229 Å<sup>2</sup>. Tyr residues are changed only moderately frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [phenylalanine](#) residues.

The Tyr side chain is predominantly a phenolic group:





The hydroxyl group makes this aromatic ring, in comparison to that of **Phe residues**, relatively reactive in electrophilic substitution reactions. These usually occur at the symmetric positions designated here as  $\epsilon^1$  and  $\epsilon^2$ , which elsewhere are often numbered as 3 and 5, respectively. Consequently, Tyr side chains can be readily **nitrated** and **iodinated**.

The Tyr hydroxyl group ionizes only at very alkaline pH values, with an intrinsic  $pK_a$  of about 11.1. In the nonionized form, it can function as either a donor or acceptor of [hydrogen bonds](#). Consequently, Tyr residues are more [hydrophilic](#) and less **hydrophobic** than the closely similar Phe residues. Tyr residues are usually buried only if the hydroxyl group can participate in hydrogen bonding, and only about 15% of Tyr residues are fully buried in folded [protein structures](#). Nevertheless, even fully buried residues are generally flipping rapidly by  $180^\circ$  rotations about the  $C_b-C_a$  single bond. Tyr residues stabilize the **alpha-helical** conformation only slightly in model peptides and in folded proteins are found twice as frequently in [beta-sheets](#).

The spectral properties of Tyr residues are sensitive to their environment, and both **absorbance** and **fluorescence** are very useful in monitoring changes in the protein structure, for example, on unfolding. Tyr residues are the sites of **phosphorylation** and [adenylation](#) in covalent **enzyme regulation**. They can also be subject to the [post-translational modification](#) of [sulfation](#).

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

## Tyrosine Kinase Receptors

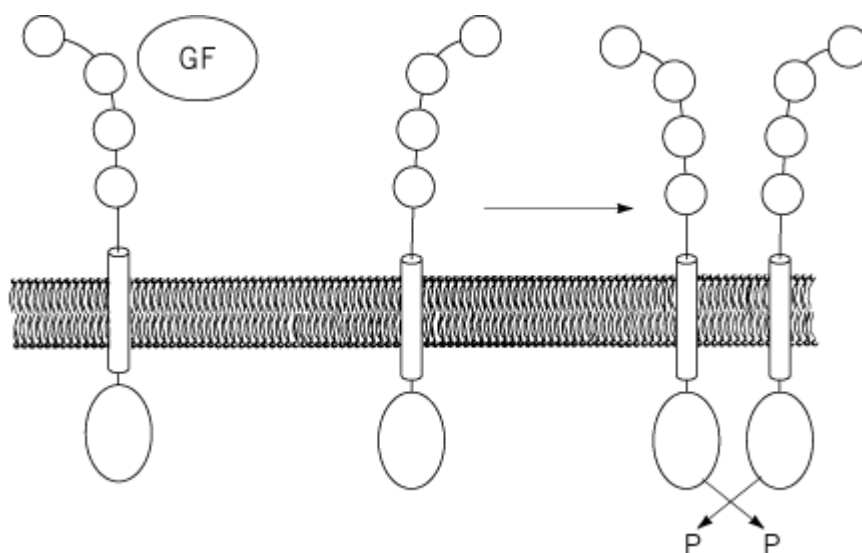
[Growth Factors](#) exert profound effects on the metabolism, survival, and differentiation of numerous cell types. Studies on the genetics and physiology of [development](#) have revealed that the synthesis of these factors is under tight control, and furthermore that their functions are highly specialized; yet the receptors for these molecules appear to be surprisingly similar in structure. Growth-factor receptors have three distinct **domains**: an extracellular ligand-binding region, a **hydrophobic** transmembrane region, and a cytoplasmic region that contains [tyrosine kinase](#) activity. These receptor tyrosine kinases (RTKs) comprise a gene [superfamily](#) of over 30 members that are similarly activated upon ligand binding, undergoing dimerization and autophosphorylation on **tyrosine** residues, which initiates a series of distinct, but often redundant, [signal transduction](#) pathways.

Growth factors induce a variety of changes in transport properties, acute metabolic activities, membrane trafficking, **cytoskeletal** interactions, and a program of gene induction governing cell growth, differentiation, and migration. Some of the factors induce a program of cell growth and viability in some cells, differentiation in others. Interestingly, growth factors with quite different cellular effects induce many of the same signaling pathways upon receptor binding, suggesting that subtle differences result from combinations of divergent and convergent signaling pathways, invoking complex mechanisms that allow cells to respond appropriately.

## 1. Molecular Interactions of Growth Factor Receptors

The gaps in our understanding of signal transduction for RTKs have been narrowed by recent insights into [protein-protein interactions](#) produced as a result of receptor activation. It is generally believed that most of these receptors undergo a similar program of dimerization and transactivation following ligand binding, resulting in tyrosine phosphorylation (Fig. 1) (1, 2). In turn, these tyrosine-phosphorylated proteins serve as binding sites for a number of functionally related signaling molecules. Such interactions are typified by those with proteins containing *src* homology 2 (SH2) domains (3). These are conserved noncatalytic regions of approximately 100 amino acid regions that function as recognition motifs for specific sequences containing phosphotyrosine residues. Some of these proteins, like the pp60<sup>src</sup> tyrosine kinase itself, as well as [phospholipase C](#) (PLC), ras GTPase activating protein (GAP), the protein tyrosine phosphatase (PTPase) SHP2, and the **oncoprotein** vav, have catalytic activities that are involved in signaling pathways. These proteins can be tyrosine-phosphorylated upon binding to RTKs in some cases, although the functional consequences are unclear. For other SH2 proteins, including the 85-kDa regulatory subunit of phosphatidylinositol-(PI)-3 kinase Grb-2 and the oncoproteins crk and nck, catalytic activities have not been detected; this suggests that these molecules serve as adapter proteins, linking receptors or their substrates to other pathways. In all cases, the binding of SH2 proteins is absolutely dependent upon tyrosine phosphorylation of the receptor and is mediated solely through the SH2 domain (4-6).

**Figure 1.** Dimerization of growth factor receptors. These receptors undergo a program of dimerization and autophosphorylation upon binding of growth factor. The tyrosine phosphorylation is thought to occur between receptors brought together by the dimerization.



The structural basis for differential association of RTKs with SH2 domains is now fairly well understood (7). Direct binding measurements have revealed quantitative kinetic differences in these interactions, perhaps accounting for differential activation of pathways that are dictated by receptor

concentration. Analysis of competition binding among SH2 domains to mutated receptors and studies using peptide inhibitors or peptide [libraries](#) have revealed certain preferential [primary structure](#) motifs in receptors for SH2 binding. In general, the amino acid residues contiguous to the carboxy-terminal side of the phosphotyrosine residue prescribe the binding specificity, in which the specific recognition can be dictated by residues 2 to 5 positions toward the *C*-terminus. Additionally, conserved sequences in SH2 domains appear to be critical for phosphotyrosine binding, as demonstrated with the conserved FLVRES (Phe—Leu—Val—Arg—Glu—Ser) sequence. Analysis of these binding data for the **abl** SH2 domain, considered in light of the solution structure, suggested that the two primary [amino groups](#) from the conserved Arg residue undergo bidentate interactions with two of the oxygen atoms of the phosphate, while surrounding Ser hydroxyl groups may [hydrogen bond](#) to the remaining phosphate oxygen.

There are likely to be other protein-interaction domains found in proteins that interact with RTKs. One such domain was identified in the tyrosine kinase substrates Shc and IRS-1, known as the PTB or PI domain (8). These domains are functionally similar to SH2 domains, in that they recognize sequences containing phosphotyrosine, but they appear to have a different sequence specificity, in which residues amino-terminal to the tyrosine are important for binding. Many of the proteins that interact with Shc through this domain have the sequence NPXY (Asn—Pro—X—Tyr), where the Tyr residue may be phosphorylated, although sequence variations on the *N*-terminal side are likely to influence the binding affinity further. Structural studies (9) suggest that these domains may resemble **pleckstrin homology** (PH) domains, another motif commonly found in signaling proteins.

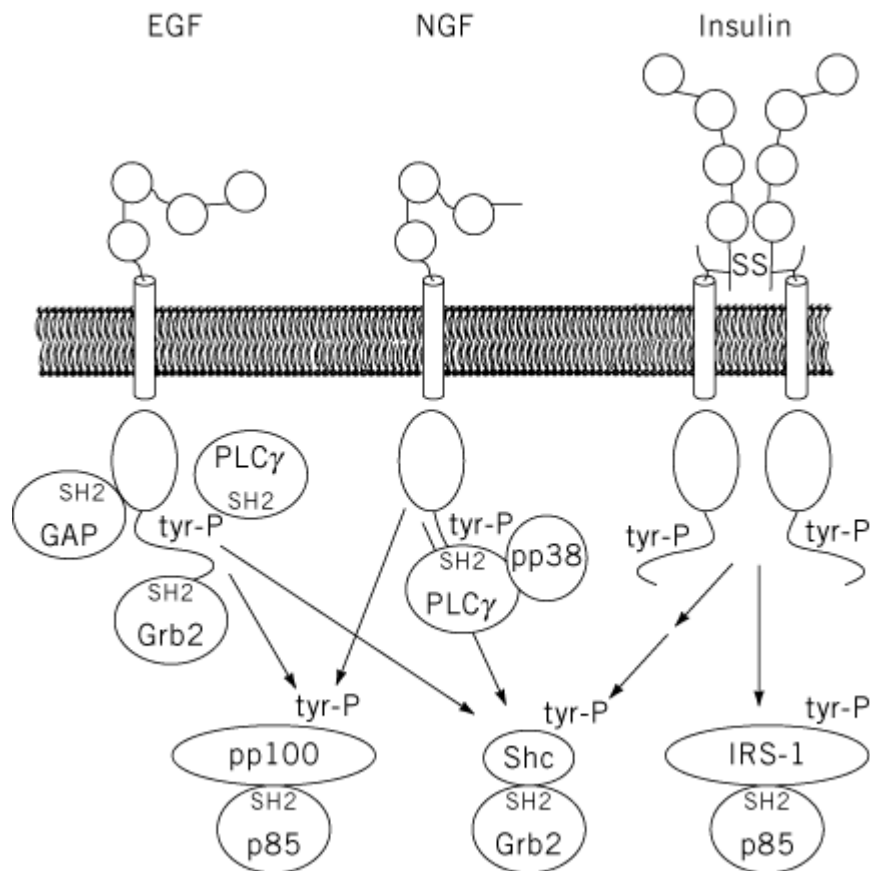
## 2. Signaling Pathways Activated by Receptor Tyrosine Kinases

Although the three-dimensional structures and molecular dynamics for many examples of RTK-SH2 or PTB binding are well understood, the precise functions of these interactions have been more difficult to resolve. In some cases, SH2 domain occupancy directly modulates the activity of the protein, as is the case for the tyrosine phosphatase SHP2, in which the catalytic activity is markedly elevated by the interaction of the phosphatase with tyrosine-phosphorylated receptors or their substrates (10). In other cases, SH2- and PTB-containing proteins are tyrosine-phosphorylated when bound to RTKs, presumably due to the induced proximity. In the case of [phospholipase C \$\gamma\$](#) , tyrosine phosphorylation may be a requirement for activation of the enzyme (11).

The emerging large numbers of RTKs and SH2 proteins begs one central question: where is the specificity? Indeed, most SH2 domains can interact with phosphotyrosine alone, but at significantly higher concentrations than are required for RTK interaction. However, **immunoprecipitation** experiments have indicated that these molecules discriminate in their binding to receptors or substrates, perhaps providing the biochemical basis for initiation of the divergent cellular effects of growth factors. In the case of some receptors, typified by that for insulin and IGF-1, and intracellular kinases, such as the JAK family that are activated by cytokine receptors (see [JAK/STAT Signaling](#)), the tyrosine kinase can phosphorylate endogenous substrates that serve as surrogates for SH2 binding. The phosphorylation of IRS-1 by [insulin](#) or IL-4 induces stable complexes between IRS-1 and SH2 proteins. These interactions can lead to the activation and/or localization of signaling cascades, including **phosphatidylinositol-3 kinase**, [MAP kinase](#), and others. These examples of different pathways of RTK-SH2 protein coupling may reflect a combinatorial diversity that allows for subtle differences in signal discrimination (Fig. 2).

**Figure 2.** Combinatorial diversity in growth factor signaling. The specificity of signal initiation in growth-factor action depends on the molecular interaction of receptors. Interactions of the [epidermal growth factor](#) (EGF), **nerve growth factor** (NGF), and insulin receptors with proteins containing SH2 domains are depicted. These three growth factors produce different effects on cells, but they can activate some of the same pathways. The EGF receptor can form high-affinity complexes with signaling proteins such as PLC $\gamma$ 1, ras GAP, and Grb2, through their respective SH2 domains. The pp140trk NGF receptor can interact with PLC $\gamma$ 1 and an associated 38-kDa tyrosine-phosphorylated protein, the function of which is unknown at present. In contrast, the insulin receptor does not form stable complexes with these

SH2-domain proteins. All three of these growth factors stimulate the activity of phosphatidylinositol-3 kinase, probably by the tyrosine-phosphorylation of proteins that bind to the SH2 domain of the 85-kDa regulatory subunit of the enzyme. In the case of the insulin receptor, this surrogate phosphoprotein, called IRS-1, has 10 known phosphorylation sites within the YMXM (Tyr—Met—X—Met) sequence that is thought to be recognized specifically by p85. The mechanism of activation of phosphatidylinositol-3 kinase by EGF and NGF is less clear, but it may involve phosphorylation of a 100-kDa protein functionally analogous to IRS-1.



### 3. Serine/Threonine Phosphorylation Cascades

In cells treated with growth factors or **transfected** with tyrosine-kinase **oncogenes**, tyrosine phosphorylations are relatively scarce compared to those found on serine and threonine residues. These observations suggested that serine kinases were activated directly or indirectly after tyrosine phosphorylation by growth factor receptors, thus dramatically amplifying the initial signal in the cell. Numerous growth-factor-dependent protein kinases have been described. The best characterized of these is the family of enzymes known collectively as mitogen-activated protein kinases ([MAP kinases](#)).

#### Bibliography

1. A. R. Saltiel, (1998) *Adv. Mol. Cell. Endocrinol.* **2**, 83–98
2. J. Schlessinger and A. Ullrich (1992) *Neuron* **9**, 383–391
3. C. A. Koch, D. Anderson, M. F. Moran, C. Ellis, and T. Pawson (1992) *Science* **252**, 668–674
4. J. R. Downing, B. Margolis, A. Zilberstein, R. A. Ashmum, A. Ullrich, C. J. Sherr, and J. Schlessinger (1989) *EMBO J.* **8**, 3345–3350
5. J. A. Escobido, S. Navankasattusas, W. M. Kavanaugh, D. Milfay, V. A. Fried, and L. T. Williams (1991) *Cell* **65**, 75–82
6. M. Ohmichi, S. J. Decker, and A. R. Saltiel (1992) *Neuron* **9**, 767–777

7. J. Kuriyan and D. Cowburn (1993) *Curr. Opin. Struct. Biol.* **3**, 828–837
8. W. M. Kavanaugh, C. W. Turck, and L. T. Williams (1995) *Science* **268**, 1177–1179
9. M.-M. Zhou, K. S. Ravichandran, E. T. Olejniczak, A. M. Petros, R. P. Meadows, M. Settler, J. E. Harlan, W. S. Wade, S. J. Burakoff, and S. W. Fesik (1995) *Nature* **378**, 584–592
10. S. Ploskey, T. J. Wandless, C. T. Walsh, and S. E. Shoelson (1995) *J. Biol. Chem.* **270**, 2897–3010
11. S. Nishibe, M. I. Wahl, S. M. T. Hernandez-Solomayor, N. K. Tonks, S. G. Rhee, and G. Carpenter (1990) *Science* **250**, 1253–1256

## Tyrosine Kinases and Phosphatases

Phosphorylation of [tyrosine](#) residues occurs at a much lower level than that of [serine](#) and [threonine](#), but it is of critical functional importance in signal transduction. In eukaryotes, phosphorylation of tyrosine residues is the result of the action of either specialized enzymes, which phosphorylate only these residues, or of dual-specificity protein kinases that act on both tyrosine and threonine residues of specific protein substrates. Up to now, specialized tyrosine kinases (Tyr kinases) have been found only in multicellular animals, suggesting that these enzymes have evolved as part of the control of cell–cell interactions. In contrast, dual-specificity protein kinases are found in all eukaryotic species. In addition, phospho-tyrosine residues have been detected in a number of bacterial species belonging to eubacteria and cyanobacteria ([1](#)). However, the function of tyrosine phosphorylation in these organisms, as well as the nature of the protein kinases and phosphatases involved, is not known.

### 1. Tyr Kinases

#### 1.1. Structure and Regulation

Specialized Tyr kinases display strong sequence homologies in their catalytic domains, which are also closely related to those of Ser/Thr kinases ([2](#)). Tyrosine kinases are more closely related to each other than to Ser/Thr kinases, and a number of signature residues are found in all of them, making it possible to predict the specificity of these enzymes on the basis of their amino acid sequences. The three-dimensional structure of the catalytic domain of tyrosine kinases, as determined by [X-ray crystallography](#), is also very similar to that of Ser/Thr kinases ([3-5](#)). Tyrosine kinases can be grouped into two categories (Table [1](#)): those that have a **transmembrane** segment, which are usually receptors for extracellular ligands (receptor tyrosine kinases, RTK), and those which are only intracellular (nonreceptor tyrosine kinases, NRTK) (see [Receptors Linked To Tyrosine Kinases](#) and [Tyrosine Kinase Receptors](#)).

**Table 1. Protein Tyrosine Kinases Families**

---

#### Receptor Tyrosine Kinases

EGF receptor

Eph

Axl

Tie/Tck

PDGF receptor  
FGF receptor Flt, Flk  
Insulin receptor  
LTK/ALK  
Ros/sevenless  
Trk Ror  
DDR/TKT  
HGF receptor  
Nematode Kin15 / 16 family  
Ret

### **Nonreceptor Tyrosine Kinases**

Src  
Tec/Atk Btk  
Csk  
Fes  
Abl  
Syk/Zap70  
Tyk2/JAK  
Ack  
Focal adhesion kinase

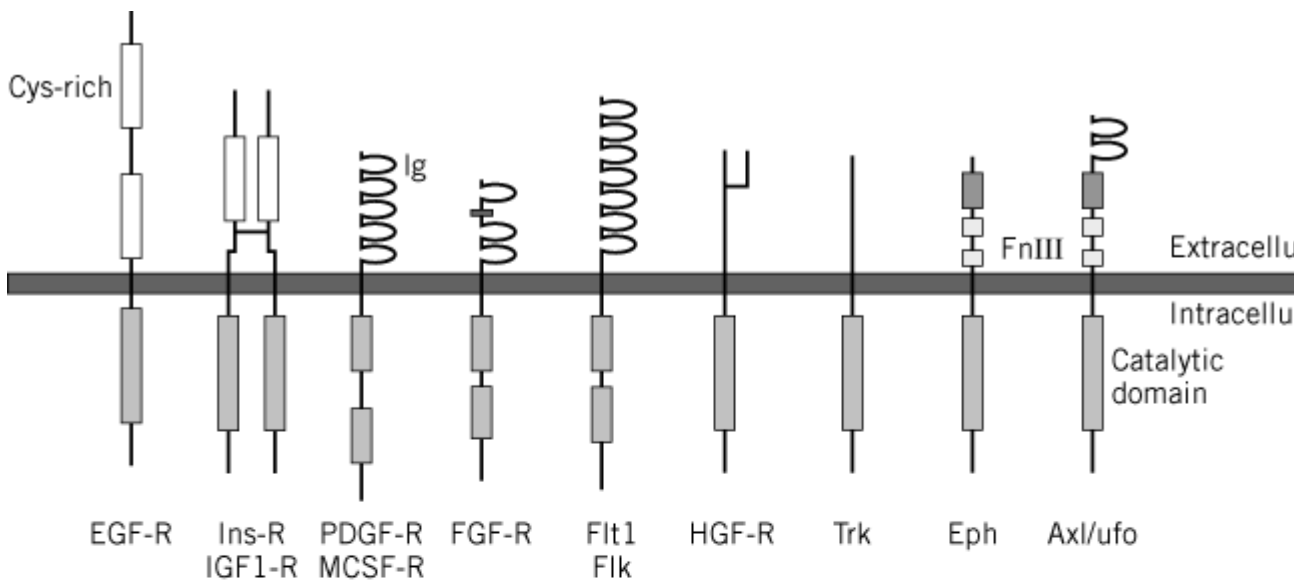
---

#### 1.2. Receptor Tyrosine Kinases

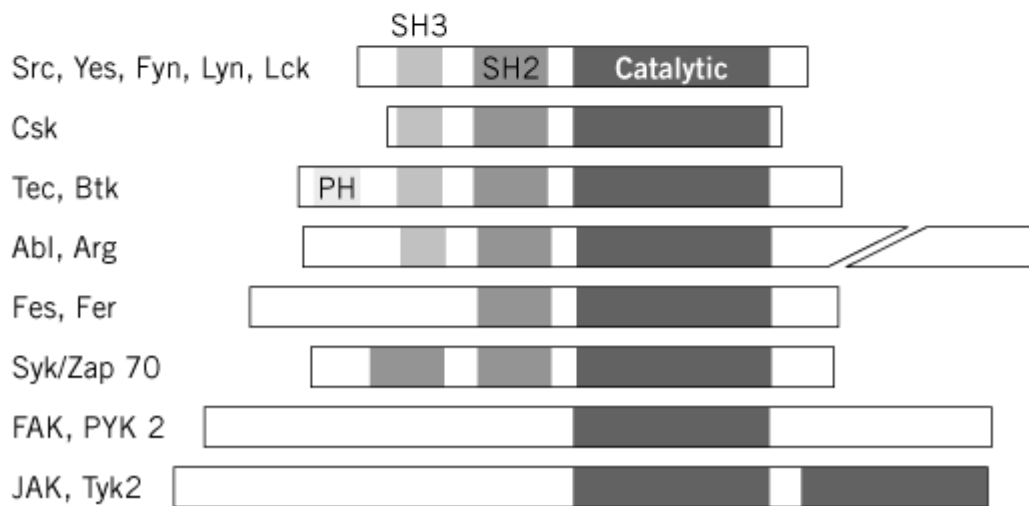
These proteins have a single transmembrane segment, an extracellular **domain** of varying length, and an intracellular catalytic domain (see Fig. 1a and Refs. 6-8). In the absence of ligand, they can be found as monomers or interacting with other proteins. In some cases, the polypeptide chain possesses the transmembrane and catalytic domains, which lacks the receptor domain. The latter is provided by an associated protein, which may itself not be a transmembrane protein but is attached to the transmembrane polypeptide by a **disulfide bond** (eg, insulin receptor) or to the **membrane** by a glycosyl-phosphatidyl-inositol **GPI anchor**, as in the case of the tyrosine kinase Ret, which associates with such a protein that is the receptor for GDNF glial derived neurotrophic factor (9, 10). In all cases that have been well-studied, activation of RTKs requires their dimerization (11, 12). In some cases, the dimer is preformed (eg, the **insulin** receptor) but inactive. In other cases, ligand binding promotes dimerization, which may result from a conformational change of the extracellular domain or from the divalent nature of the ligand itself (11). Upon ligand binding, a *trans*-phosphorylation of the intracellular catalytic domains occurs, leading to the phosphorylation of tyrosine residues located in the activation loop, a critical loop of the catalytic domain, and thereby enhancing its enzymatic activity. Recent work suggests that some receptors have a high rate of basal phosphorylation (ie, in the absence of ligand), compensated by a high rate of dephosphorylation, and that activation could result from a tilting of the balance in favor of the phosphorylation (13). Additional phosphorylated tyrosine residues provide docking sites for protein domains that interact specifically with tyrosine-phosphorylated peptides, the **SH2** (Src-homology 2) and **PTB** (phospho-tyrosine binding) domains (14). In some cases, cytoplasmic proteins that associate with the receptor are also highly phosphorylated on tyrosine residues and provide additional docking sites (eg, the insulin receptor is associated with IRS-1 and 2 proteins). SH2 and PTB domains are found in a variety of proteins that cluster around the activated receptor, forming membrane-associated “signal transfer particles” (11). Thus, a number of enzymes are translocated from the cytosol to the membrane area surrounding the receptor. Several of these enzymes have lipid substrates, such as **phospholipase C<sub>g</sub>** (PLC<sub>g</sub>) and **phosphatidylinositol 3 kinase** (PI3-kinase), and translocation to the membrane enhances their

activity by putting them in contact with their substrates (15). Another class of proteins with SH2 or PTB domains contains adapter proteins that are devoid of enzymatic activity but associate with active proteins (14). One example is Grb2, which possesses a SH2 domain, but interacts also with a [guanine nucleotide exchange factor](#), SOS, via two SH3 domains. Thus, recruitment of Grb2 to a phosphorylated receptor brings SOS in close contact with the membrane to which its substrate, the small [GTP-binding protein](#) Ras, is attached. SOS promotes the GDP/GTP exchange on Ras, which in its GTP-bound form becomes capable of interacting with a Ser/Thr kinase, Raf, that activates the [MAP-kinase](#) cascade (16). Thus, activation of a RTK by its ligand triggers a number of Ser/Thr phosphorylation pathways, since the stimulation of PLC $\beta$  increases the activity of protein kinase C, and the action of PI3-kinases activates additional serine/threonine kinases (17, 18). These various serine/threonine protein kinases mediate many of the effects of [growth factors](#) by altering [transcription](#) and [translation](#). It should be pointed out, however, that each SH2 or PTB domain has its own specificity, in addition to its requirement for a phosphorylated tyrosine residue (14). Therefore, phosphorylation of particular tyrosine residues on a tyrosine kinase receptor or its associated proteins will only lead to the recruitment of specific proteins. Consequently, different receptors may recruit different signaling proteins. Another factor that is important for the specificity of the signaling of each RTK is the rate of its deactivation (19). Some receptors may be rapidly inactivated by dephosphorylation or by endocytosis, whereas others will stay activated for a long time.

**Figure 1.** Domain structures of specific protein tyrosine kinases. (a) Receptor tyrosine kinases. Cys-rich, cysteine rich; I immunoglobulin-like domain, FnIII, fibronectin III-like domain. (b) Nonreceptor tyrosine kinases. SH2, Src-homology domain 2; SH3, Src-homology domain 3; PH, pleckstrin-homology domain.



(a)



(b)

### 1.3. Nonreceptor Tyrosine Kinase

NRTKs form a heterogeneous group of tyrosine kinases in which the catalytic domain is associated with various other domains (Fig. 1b). Some of these NRTKs are functionally equivalent to the catalytic domain of RTKs. For example, the tyrosine kinases of the JAK (Janus kinase) family associate by their amino-terminal domain to the intracellular segment of **cytokine** receptors (20). Activation of the JAKs occurs in response to cytokine binding to its receptor and dimerization of the receptor, and it leads (i) to the phosphorylation of the receptor on tyrosine residues and (ii) the recruitment and tyrosine phosphorylation of a group of [transcription factors](#) containing a SH2 domain, the STATs (signal transducers and activators of transcription) (see [JAK/STAT Signaling](#)). Following tyrosine phosphorylation, STAT dimerizes and translocates to the nucleus, where it activates the transcription of specific genes. Focal adhesion kinase (FAK), which is important in signal transduction in response to the [extracellular matrix](#), also associates with transmembrane proteins. Following engagement of receptors for the extracellular matrix, the [integrins](#), FAK is recruited to focal adhesions and becomes autophosphorylated (21). This allows the binding of Src-family kinase SH2 domains, their activation, and the phosphorylation of a number of associated proteins, leading to the formation of transducing complexes reminiscent of, but distinct from, those



created by RTKs. Src-family kinases possess a SH2 and a SH3 domain, in addition to their tyrosine kinase catalytic domain (22). They can be associated to the membrane by an amino-terminal **myristoyl group**. In the inactive state, they are phosphorylated on a carboxy-terminal tyrosine residue that undergoes an intramolecular interaction with the SH2 domain, blocking access of substrates to the [active site](#). This inactivation is strengthened by the binding of the SH3 domain to a specific peptide sequence located between the SH2 and the catalytic domains (4, 5). Thus, activation of Src-family kinase may result from the dephosphorylation of the carboxy-terminal tyrosine residue and/or from displacement by other ligands of the intramolecular binding by the SH2 or SH3 domains. In many cells, these kinases appear to be involved secondarily, following activation of other tyrosine kinases, including RTKs or FAK (22). They also play a critical role in the activation of lymphocytes, a process in which Src-family kinases and other NRTKs, such as Syk/Zap70 and Btk, interact closely with the antigen receptors. Another NRTK, Csk (carboxy-terminal Src kinase), exerts a potent negative control on the activity of Src-family kinases by phosphorylating their carboxy-terminal tyrosine residue, thereby promoting its inhibitory intramolecular interaction with the SH2 domain, as indicated above. Finally, it should be mentioned that although the majority of tyrosine kinases are localized and activated at the plasma membrane, some tyrosine kinases (eg, Abl, Fer) are located in the nucleus, where they participate in the regulation of [transcription](#) and the [cell cycle](#) (23).

## 2. Tyr Phosphatases

Several groups of enzymes, belonging to different gene families, have the capacity to dephosphorylate tyrosine residues (Table 2). They have in common a signature sequence Cys-X<sub>5</sub>-Arg and a similar configuration of the catalytic site, which also contains a conserved **aspartate** residue (24-26). The catalytic mechanism involves the formation of a covalent phosphate thioester intermediate. Because of the presence of a [cysteine](#) residue in the active site, tyrosine phosphatases are highly sensitive to [thiol-group](#) reagents and oxidizing agents. Tyrosine phosphatases are also inhibited nonspecifically by chemicals that mimic phosphate (eg, sodium orthovanadate, pervanadate) or phosphotyrosine (phenylarsine oxide). The four groups of tyrosine phosphatases include two groups that dephosphorylate exclusively tyrosine residues, the “classical” phosphotyrosine phosphatases (PTP) and the low-molecular-weight PTPs, and two groups of dual-specificity protein phosphatases (25).

**Table 2. Protein Tyrosine Phosphatases and Dual Specificity Phosphatases**

### **Protein Tyrosine Phosphatases**

#### Classical Protein Tyrosine Phosphatases

##### *Transmembrane PTPs*

Classified according to their extracellular domain

Type I (eg, CD45)

Type II (eg, LAR)

Type III (eg, PTPb)

Type IV (eg, PTPa)

Etc.

##### *Nontransmembrane PTPs*

Classified according to their non-catalytic domain(s)

Band 4.1 domain (eg, PTP-H1)

SH2 domain (eg, SH-PTP1)

PEST domain (eg, PTP-PEST)

Etc.

Low  $M_r$  tyrosine phosphatases: Small cytosolic enzymes

### Dual Specificity Phosphatases

VH1 group

Dephosphorylate MAP-kinases (eg, MKP1, CL100)

cdc 25 group

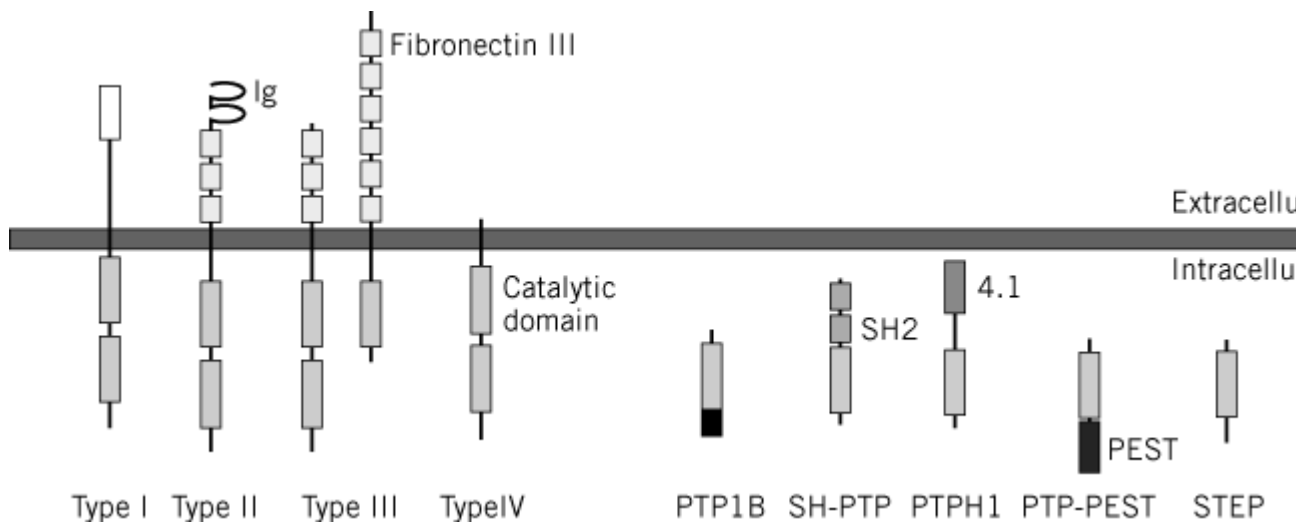
Dephosphorylate cyclin dependent kinases (eg, cdc25)

---

## 2.1. Classical PTPs

Like tyrosine kinases, PTPs may have a transmembrane segment and an extracellular domain, or they can be entirely intracellular (Fig. 2). PTPs from the first group have the organization characteristic of transmembrane receptors and are therefore termed RPTPs (receptor PTPs). Ligands for the extracellular domains of some of these RPTPs have been identified, but no clear regulation of the phosphatase activity by these ligands has yet been demonstrated. The extracellular domains of RPTPs contain modules often found in [cell adhesion molecules](#), and some of them appear to be involved directly in cell–cell interactions (27). Such RPTPs play an important role during [development](#) in axon pathfinding (28). The crystal structure of the intracellular catalytic domain of RPTP-a suggests that it may be inactivated by dimerization (29). An intriguing feature of many RPTPs is the possession of tandem intracellular catalytic domains. However, most or all the phosphatase activity appears associated with only one domain, and the function of the other is not known. PTPs lacking a transmembrane domain nevertheless possess various targeting domains, including SH2 domains and band 4.1 domains (which may target them to membrane-associated proteins) (30). When studied *in vitro*, PTPs have a very high catalytic activity, suggesting that they are tightly regulated in cells. Although these regulations are still poorly understood, they are likely to involve a variety of mechanisms, including phosphorylation, targeting, and ligand binding. PTPs appear to have multiple functions (27, 28, 31-34). On the one hand, they oppose the action of tyrosine kinases, and this role is clearly illustrated in many cells by the sharp increase in tyrosine phosphorylation of many proteins following inhibition of tyrosine phosphatases with nonspecific agents. On the other hand, some PTPs may activate tyrosine phosphorylation pathways—for example, by dephosphorylating specifically the carboxy-terminal inhibitory tyrosine of Src-family kinases (33). Interestingly, PTPs are used by pathogenic microorganisms to deceive cell defenses, as demonstrated in the case of bacteria from the *Yersinia* genus (35). Low-molecular-weight PTPs, which form a separate group of enzymes, are expressed in many cell types, and their regulation and function is not well-characterized (36).

**Figure 2.** Domain structures of classical protein tyrosine phosphatases. Ig, immunoglobulin-like domain; SH2, Src-homology domain 2; 4.1, band 4.1-like domain, PEST, proline, glutamate, serine, threonine-rich domain.



## 2.2. Dual-Specificity Phosphatases

Two groups of protein phosphatases are able to dephosphorylate serine and threonine residues, as well as tyrosines (25). Such phosphatases play a critical role in the control of cell division, because they dephosphorylate tyrosine and threonine residues of cyclin-dependent kinase, cdk2, a step necessary for its activation (37). Another group of enzymes is comprised of relatives of the vaccinia virus gene VH1. Some of these dual-specificity phosphatases dephosphorylate the two residues necessary for the activity of MAP kinases. Interestingly, several of them are immediate-early genes, whose transcription is stimulated by MAP kinase, providing a strong feedback mechanism (38).

### Bibliography

1. P. J. Kennelly and M. Potts (1996) *J. Bacteriol.* **178**, 4759–4764.
2. S. K. Hanks and T. Hunter (1995) *FASEB J.* **9**, 576–596.
3. S. R. Hubbard, L. Wei, L. Ellis, and W. A. Hendrickson (1994) *Nature* **372**, 746–754.
4. W. Q. Xu, S. C. Harrison, and M. J. Eck (1997) *Nature* **385**, 595–602.
5. F. Sicheri, I. Moarefi, and J. Kuriyan (1997) *Nature* **385**, 602–609.
6. W. J. Fantl, D. E. Johnson, and L. T. Williams (1993) *Annu. Rev. Biochem.* **62**, 453–481.
7. J. Schlessinger and A. Ullrich (1992) *Neuron* **9**, 383–391.
8. P. Van der Geer, T. Hunter, and R. A. Lindberg (1994) *Annu. Rev. Cell Biol.* **10**, 251–337.
9. P. Durbec, C. V. Marcos-Gutierrez, C. Kilkenny, et al. (1996) *Nature* **381**, 789–793.
10. S. Q. Jing, D. Z. Wen, Y. B. Yu, et al. (1996) *Cell* **85**, 1113–1124.
11. A. Ullrich and J. Schlessinger (1990) *Cell* **61**, 203–212.
12. C. H. Heldin (1995) *Cell* **80**, 213–223.
13. H. Daub, C. Wallasch, A. Lanckenau, A. Herrlich, and A. Ullrich (1997) *EMBO J.* **16**, 7032–7044.
14. T. Pawson (1995) *Nature* **373**, 573–580.
15. G. Panayotou and M. D. Waterfield (1993) *Bioessays*. **15**, 171–177.
16. R. Marais and C. J. Marshall (1996) *Cancer Surv.* **27**, 101–125.
17. C. G. Proud (1996) *Trends Biochem. Sci.* **21**, 181–185.
18. B. M. Marte and J. Downward (1997) *Trends Biochem. Sci.* **22**, 355–358.
19. C. J. Marshall (1995) *Cell* **80**, 179–185.
20. S. Pellegrini and I. Dusanter-Fourt (1997) *Eur. J. Biochem.* **248**, 615–633.
21. D. D. Schlaepfer and T. Hunter (1998) *Trends Cell Biol.* **8**, 151–157.

22. S. M. Thomas and J. S. Brugge (1997) *Annu. Rev. Cell Dev. Biol.* **13**, 513–609.
23. J. Y. J. Wang (1994) *Trends Biochem. Sci.* **19**, 373–376.
24. R. L. Stone and J. E. Dixon (1994) *J. Biol. Chem.* **269**, 31323–31326.
25. E. B. Fauman and M. A. Saper (1996) *Trends Biochem. Sci.* **21**, 413–417.
26. D. Barford, Z. C. Jia, and N. K. Tonks (1995) *Nature Struct. Biol.* **2**, 1043–1053.
27. S. M. Brady-Kalnay and N. K. Tonks (1995) *Curr. Opin. Cell Biol.* **7**, 650–657.
28. C. B. Chien (1996) *Neuron.* **16**, 1065–1068.
29. A. M. Bilwes, J. Den Hertog, T. Hunter, and J. P. Noel (1996) *Nature* **382**, 555–559.
30. L. J. Mauro and J. E. Dixon (1994) *Trends Biochem. Sci.* **19**, 151–155.
31. M. C. Faux and J. D. Scott (1996) *Cell* **85**, 9–12.
32. T. Matozaki and M. Kasuga (1996) *Cell. Signal* **8**, 13–19.
33. B. G. Neel (1997) *Curr. Opin. Immunol.* **9**, 405–420.
34. N. K. Tonks and B. G. Neel (1996) *Cell* **87**, 365–368.
35. E. G. Ninfa and J. E. Dixon (1994) *Trends Cell Biol.* **4**, 427–430.
36. G. Ramponi and M. Stefani (1997) *Int. J. Biochem. Cell Biol.* **29**, 279–292.
37. B. Sebastian, A. Kakizuka, and T. Hunter (1993) *Proc. Natl. Acad. Sci. USA* **90**, 3521–3524.
38. H. Sun and N. K. Tonks (1994) *Trends Biochem. Sci.* **19**, 480–485.

## Ubiquitin

Ubiquitin (Ub) is a 76 amino acid–residue, heat-stable **protein** that occurs in virtually all cells and is highly conserved in its amino acid sequence. Largely through the work of Hershko and Ciechanover, ubiquitin was first shown in the late 1970s to be an essential cofactor in ATP-dependent **protein degradation** in reticulocyte extracts. Their experiments showed that ubiquitin must first be activated at its C-terminal [glycine](#) residue to be conjugated to a substrate. The activated ubiquitin molecule is then covalently linked through an isopeptide bond to an e-NH<sub>2</sub> group of a lysine residue in the substrate protein. In most Ub-conjugated proteins, the C-terminal glycine residue of one ubiquitin is linked to Lys48 of the next ubiquitin, forming long multiubiquitin chains. These long chains that usually consist of at least five ubiquitin moieties, rather than free ubiquitin, are preferentially bound and degraded by the 26S [proteasome](#).

### 1. Biosynthesis

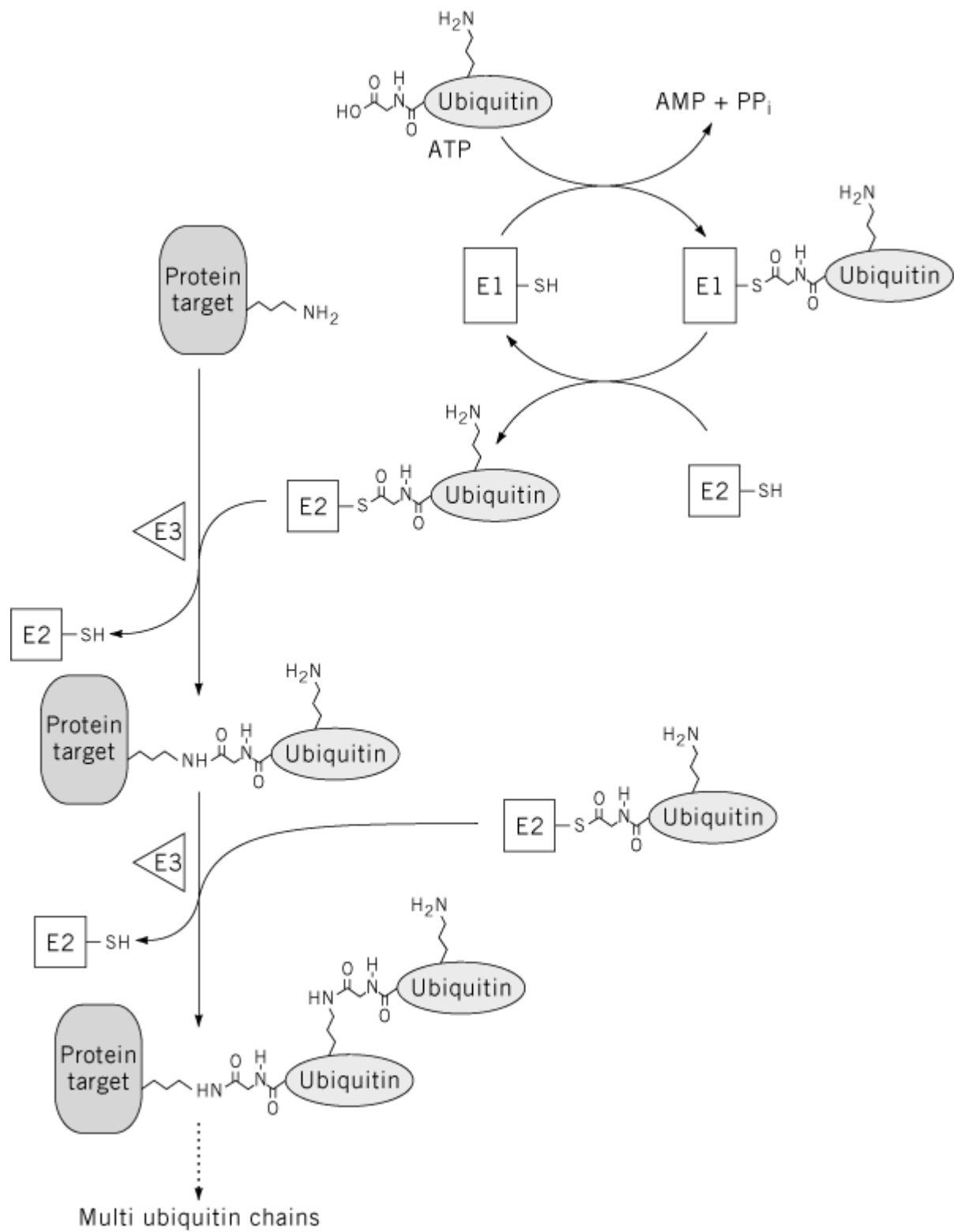
Ubiquitin itself is generated by a novel mechanism. It is found as part of two larger **polypeptide chains** that are cleaved **proteolytically** in a [post-translational modification](#) to release ubiquitin and proteins of [ribosomes](#). However, most of the ubiquitin in stressed cells is generated from a polyprotein containing multiple, tandemly repeated ubiquitin units. After [translation](#), the ubiquitin multimers are rapidly cleaved to monomers by a ubiquitin terminal hydrolase. Unlike ubiquitin-conjugated proteins destined for degradation, in which the ubiquitin is attached by isopeptide linkages to internal lysine residues, the ubiquitin multimers are simple end-to-end chains. A large number of ubiquitin C-terminal hydrolases are found in cells, but their specific functions are still unclear. They are probably important in the regulation of ubiquitin conjugate formation and disassembly. Some of these isopeptidases can be found in association with the 26S proteasome and may be important in the release of free ubiquitin monomers after the ubiquitin chain has been

removed from the protein substrate. In this way, ubiquitin monomers can be reused in degradation of new substrates.

## 2. Linkage of Ubiquitin to Proteins for Degradation

At least three enzymes (Table [1](#)) act in succession to link ubiquitin to protein substrates (Figure [1](#)):

**Figure 1.** Pathway of ubiquitin proteins conjugation. Key: E1, ubiquitin-activating protein; E2, ubiquitin carrier protein; E3, ubiquitin-protein ligase; PP<sub>i</sub>, pyrophosphate.



**Table 1. Ubiquitin-Conjugating Enzymes**

| Enzyme | Name | Abbreviation Family Members |
|--------|------|-----------------------------|
|--------|------|-----------------------------|

|    |                              |     |          |
|----|------------------------------|-----|----------|
| E1 | Ubiquitin activating protein | Uba | 1        |
| E2 | Ubiquitin carrier protein    | Ubc | ~15      |
| E3 | Ubiquitin protein ligase     | Ubr | Numerous |

---

### 2.1. E1 (Ubiquitin-Activating Protein)

The cycle begins with a reaction that generates activated ubiquitin. Hydrolysis of ATP between its  $\alpha$  and  $\beta$  phosphoryl groups drives the formation of a covalently bound C-terminal ubiquitin adenylate; this step is the only one in ubiquitin conjugation that requires energy input. The activated ubiquitin is transferred to a [thiol group](#) of an active-site cysteine residue of E1, liberating AMP, in a reaction similar to the charging of [transfer RNA](#) by [amino acyl-tRNA synthetase](#). The E1-ubiquitin thiolester undergoes another round of AMP-Ub formation, ultimately yielding a ternary complex in which E1 bears two molecules of activated ubiquitin:



This high-energy thiolester form of ubiquitin, generated by E1, is essential for the eventual linkage to substrate in a series of reactions requiring two other ubiquitin-conjugating enzymes, E2 and E3.

E1 is an abundant 110-kDa protein, and it may occur in **phosphorylated** forms that seem to differ in their subcellular localization. In plants, a family of related E1 proteins exists, but in mammalian organisms only a single functional E1 protein has been found to date. E1 is essential for cell viability. **Temperature-sensitive mutants** of E1 have been isolated, originally from a mouse tissue culture cell line, as mutants that block the [cell cycle](#), and they have been very useful in implicating the ubiquitin-proteasome pathway in the degradation of specific protein substrates.

### 2.2. E2s (Ubiquitin Carrier Proteins)

In the second step of the conjugation process, the activated ubiquitin bound to E1 is transferred to a [thiol group](#) of a ubiquitin carrier protein, or E2. At least a dozen such carrier proteins exist in cells. They have a limited ability in vitro to transfer the activated ubiquitin directly to protein substrates, but in vivo they generally serve as intermediary “carriers” in a process catalyzed by E3s. The E2s are small and share a conserved 16-kDa core containing the [cysteine](#) residue that forms a thiolester linkage with the activated ubiquitin. In fact, E2s can be isolated with ubiquitin covalently attached to this residue. A large family of E2 enzymes have been characterized in mammalian cells, plants, and yeast. These genes have been found on the basis of their conserved 16-kDa domain, and the proteins have been isolated biochemically due to their propensity to bind to a ubiquitin [affinity chromatography](#) matrix in the presence of E1 and ATP. The large number of E2s helps to generate the specificity of the ubiquitination system, as selected E2s function in the degradation of different types of substrates.

### 2.3. E3s (Ubiquitin-Protein Ligases)

In the final step of the conjugation process, the activated ubiquitin is transferred from E2 to the protein substrate by E3, the Ub-protein ligase. Most E3s catalyze formation of long ubiquitin chains, which means that the E3 active site is able to transfer ubiquitin directly either to a lysine residue of the substrate or, processively, to the preceding ubiquitin moieties to form a chain. The E3 proteins must contain binding sites for the protein substrate, E2, and at least one, possibly two, ubiquitin molecules. It is believed that E3 alone, or in a complex with an E2, generates the specificity of the ubiquitination process. A variety of structural determinants in the ubiquitinated substrate can be recognized by E3s. These include certain amino-terminal residues, certain **domains** within the substrate (ie, the “destruction box” found within many proteins degraded during mitosis), phosphorylated domains, and, perhaps, having unfolded domains. Only a small fraction of the large number of E3s have been described to date. These fall into several disparate families, with much less structural similarity than among E2s. Several appear to be high molecular-weight multimers or

complexes.

#### 2.4. HECT Family

Proteins with HECT (*homologous to E6-AP carboxyl terminus*) domains form another class of proteins with E3 activity. The prototype of this family, E6-AP, is involved in the degradation of p53 protein in certain cells infected with human papilloma **viruses**. Unlike other E3s, HECT-family E3s actually form a short-lived thiolester with ubiquitin transferred from E2, and then pass the activated Ub on to the substrate. These E3s appear to be important in the turnover of **RNA polymerase**, certain cell cycle regulators, and some membrane-bound **sodium channels**.

#### 2.5. Anaphase-Promoting Complex (APC)/Cyclosome

This ligase is a large multisubunit complex (~1500 kDa) known as the cyclosome or APC and is involved in the conjugation of ubiquitin to mitotic [cyclins](#) and other proteins involved in movement through mitosis. This large E3 complex utilizes a unique E2 enzyme (E2C/E2x/UbcH10) and recognizes a nine-residue motif in the substrate protein, known as the “destruction box.” The destruction box restricts degradation by the E2C/cyclosome pathway to only those proteins that contain that tag.

#### 2.6. SCF

Regulated protein destruction during the cell cycle can also be catalyzed by another set of ubiquitination enzymes. This pathway has been best elucidated in yeast at the transition between the G1 and S phases of the cell cycle. Here, a heterotrimeric E3 complex termed SCF (Skp1/cullin/F-box protein) acts with a specific E2 (Cdc34) to ubiquitinate a cell cycle inhibitor (sic1), allowing progression to S phase. The existence of large families of both cullins and F-box proteins suggests numerous functions in the degradation of other substrates by combining different subunit modules. Members of the SCF family have also been shown to degrade other phosphorylated proteins, such as I $\kappa$ B.

Other E3s undoubtedly exist to ubiquitinate specific cellular substrates. In yeast, for example, a number of E2s serve particular functions (eg, degradation of abnormal proteins by UBC4/5), yet their cognate E3s are unknown.

### 3. Ubiquitin Homologues

A number of proteins with homology to ubiquitin have been discovered recently. Although they have low overall amino acid sequence identities, (perhaps 20%), they share a characteristic ubiquitin-like [tertiary structure](#). Most interestingly, some members (SUMO1 and its yeast homologue Smt3) that contain the same -Gly-Gly C-terminus as ubiquitin have been found conjugated to cellular proteins in a way similar to ubiquitin. This SUMO1 conjugation appears to be involved in targeting RanGAP1 to the [nuclear pore complex](#), where it is required for [nuclear import/export](#). Unlike most cases of ubiquitin conjugation, the conjugation of SUMO1 to cellular protein substrates is not associated with protein degradation. It seems likely that the attachment of ubiquitin and ubiquitin-like proteins to other cellular substrates will prove to be a general mechanism for targeting substrates, to highly regulated cellular processes.

### 4. Degradation

The ubiquitin-conjugated protein is then proteolyzed by the [proteasome](#).

#### Suggestions for Further Reading

M. Glotzer, A. W. Murray, and M. W. Kirschner (1991) Cyclin is degraded by the ubiquitin pathway. *Nature* **349**, 132–138.

V. Chau, J. W. Tobias, A. Bachmair, D. Marriott, D. J. Ecker, D. K. Gonda, and A. Varshavsky (1989) A multiubiquitin chain is confined to specific lysine in a targeted short-lived protein.



Science **243**, 1576–1583.

A. Hershko and A. Ciechanover (1992) The ubiquitin system for protein degradation. *Ann. Rev. Biochem.* **61**, 761–807.

A. Ciechanover (1994) The ubiquitin-proteasome proteolytic pathway. *Cell* **79**, 13–21.

W. E. Mitch and A. L. Goldberg (1996) Mechanisms of muscle wasting. The role of the ubiquitin-proteasome pathway. *New Engl. J. Med.* **335**, 1897–1905.

H. Saitoh, R. T. Pu, and M. Dasso (1997) SUMO-1: wrestling with a new ubiquitin-related modifier. *Trends Biochem. Sci.* **22**, 374–376.

A. M. Weissman (1997) Regulating protein degradation by ubiquitination. *Immunol. Today* **18**, 189–198.

## Ultrabithorax Genes

*Ultrabithorax (Ubx)* is a gene required for the development of the specific morphology of the abdominal and the last thoracic segments of insects. It was discovered in the fruit fly *Drosophila melanogaster*, where it is located in the Bithorax complex, but in other insects it is located in the single cluster that comprises both the genes of the Antennapedia and Bithorax complexes. *Ubx* is a prototypical homeotic gene as mutations of loss or gain of function result in the morphological transformation of one segment into another. *Ubx* encodes a transcription factor with a homeodomain DNA binding motif of the Antennapedia class. The homeodomain of *Ubx* is most closely related to that of *Antennapedia (Antp)* and *abdominal-A (abd-A)*, which evolved from a single ancestral gene that also gave rise to the HOX6-HOX8 orthologs in the vertebrate lineage. Because of its function during development, its homeodomain sequence and its location in a complex with other homeotic genes, *Ubx* is catalogued as a *Hox* gene.

*Drosophila Ubx* is expressed and required in the ectoderm and mesoderm derivatives of the abdominal and the posterior thoracic segments (1, 2). Its expression is very dynamic, not coinciding with segmental landmarks but with parasegments (units shifted by half segment). *Ubx* is the only Hox gene required for the morphological differentiation of parasegments 5 and 6 (that coincide roughly with the third thoracic and first abdominal segment) and contributes to the differentiation of parasegments 7 to 13 in combination with *abd-A* and *Abd-B*. The different morphology of parasegment 5 from parasegment 6 is mainly due to the different temporal and spatial regulation of *Ubx* during the development of these parasegments (3). In the gut, *Ubx* is expressed in the visceral mesoderm in parasegment 7 but is not expressed in the endoderm. Despite this, mutations in *Ubx* also affect the development of the endoderm. This occurs because *Ubx* in the gut mesoderm activates the expression of *decapentaplegic*, a TGF $\beta$  homolog that induces gene expression in the underlying endoderm (4).

Mutations that abolish the function of *Ubx* are lethal. In *Ubx* homozygous larvae, the third thoracic and first abdominal segments are transformed into the second thoracic segment (parasegments 5 and 6 into parasegment 4). These larvae also have some incomplete homeotic transformations in the abdominal segments 2 to 7 (5). *Ubx* function is also required for the adult imaginal structures as shown by mutations in cis-regulatory sequences that remove the expression of *Ubx* in the adult primordia of specific segments. The most spectacular mutants show a complete transformation of one segment into another. For example, the combination of three viable alleles of *Ubx* (*abx*, *bx*, *pbx*, see below) gives rise to flies with two pairs of wings. This is due to the transformation of the metathorax and the halteres typical of dipterans into a mesothorax with a second pair of wings. Flies

mutant for other alleles (*bxd*) develop four pairs of legs due to the formation of an extra pair in the first abdominal segment. Complementary to the loss-of-function mutations, alleles causing expression of *Ubx* in the wing where it is normally absent produce flies in which the wings disappear and are replaced by a pair of halteres.

UBX proteins control the development of segment morphology by transcriptional regulation of downstream target genes in specific regions. It is this battery of genes that is ultimately responsible for the correct morphogenesis of the segment. Several target genes of *Ubx* have been found, and some have been proved to be transcriptionally activated or repressed by UBX. Direct UBX targets include transcription factors (like the homeobox gene *Distalless*); signalling molecules (TGF $\beta$  homolog *decapentaplegic*) and membrane-associated proteins involved in cell adhesion (*connectin*). The number of direct target genes has been estimated to be around 100 to 200 (6, 7).

The spatial regulation of *Ubx* transcription is controlled at several levels (8). Early expression is regulated by a combination of transcriptional activators required for the segmentation of the embryo (homeodomain protein like *fushi-tarazu* and *even-skipped*) and repressors (like the *hunchback* zinc finger protein). These early regulators are only transiently expressed, and other proteins replace them to keep the restricted spatial expression. The *Polycomb* genes maintain the repression of *Ubx*, and the *trithorax* genes maintain its expression. There is a less well characterized group of regulators required for the fine modulation of *Ubx* expression. The homeodomain protein encoded by the *engrailed* gene belongs to this group, but many other transcription factors, like the zinc-finger protein encoded by the *spalt* gene, are likely to be involved in the dynamic pattern of *Ubx* regulation (9).

Mutations in the *Ubx* gene map to a 100 kb region. There is a 40 kb upstream regulatory region, the *bxd* region, that gives rise to several untranslated transcripts, the *bxd*-unit (10). The UBX proteins are encoded in a large transcription unit of 77 kb that generates mRNAs of 3.2 kb and 4.3 kb, the *Ubx* unit. The differential use of two splicing acceptors and the inclusion or not of two microexons results in the translation of six protein isoforms (11, 12). The expression of the different isoforms is tissue-specific but not segment-specific. The longest isoforms are expressed mainly in the ectoderm and mesoderm, while the shortest isoforms lacking both microexons are mainly expressed in the central nervous system (13). The spatial distribution of UBX isoforms has been conserved in *Drosophila* species that diverged 60 million years ago, suggesting a functional requirement for the tissue-specific expression (14). Functional differences between UBX isoforms have been found (15). However, experimentally induced ectopic expression of different protein isoforms shows that they also have equivalent functions (16). These results suggest that isoform diversity may serve to optimise UBX protein function in different contexts, rather than to confer different regulatory properties (15).

The UBX protein has two conserved regions: the homeodomain required for the DNA binding and the YPWM peptide, a stretch of aminoacids N terminal to the homeobox that is required for protein-protein interactions. In vitro, the UBX homeodomain binds to DNA with an ATTA core sequence that is also bound by other homeobox proteins (17). In vivo, the binding specificity of UBX is acquired by interaction with cofactor proteins that increase the sequence specificity. A *Ubx* cofactor protein has been well characterized, *extradenticle* (18). Extradenticle is a homeobox protein, and it could function either by increasing the affinity of UBX for its target DNA sequences or by allowing DNA-bound UBX proteins to act as transcriptional activators (19).

The *cis* regulatory elements of *Ubx* are located in a wide region that includes both the introns and the upstream region. The viable loss-of-function alleles in *cis* regulatory elements do not affect the coding potential but alter the spatial distribution of *Ubx*. According to their phenotypes, mutations in *cis* regulatory regions have been classified in four groups: *anterobithorax* (*abx*), *bithorax* (*bx*), *bithoraxoid* (*bxo*), and *postbithorax* (*pbx*) (20). *abx* and *bx* alleles are homozygous viable and transform, to a variable extent, the anterior third thoracic segment into anterior second segment. The *abx* alleles are deletions while *bx* alleles are caused by insertions of transposable elements in the third intron of the *Ubx* transcription unit (21). The *pbx* and *bxo* alleles map in the upstream region.

The *bx* alleles are caused by insertions of transposable elements; insertions closer to the promoter have a stronger phenotype consisting in the disappearance of the abdominal structures and the formation of an extra pair of legs. The *bx* alleles also have a mild transformation of the posterior metathorax into posterior mesothorax that affects mainly proximal structures. The *pbx* alleles are caused by deletions in a more distal area of the upstream regulatory region. They result in the transformation of the posterior metathorax into posterior mesothorax. All the above viable alleles have mild phenotypes in larvae. Analysis of reporter gene constructs carrying DNA from the *abx*, *bx*, *pbx*, and *bx* *cis* regulatory regions show that these elements are capable of driving expression in the adult as well as in the larva. Analysis of the effect of different breakpoints in the upstream region as well as the work with reporter genes suggest the existence of a degree of redundancy in the *cis* regulatory elements. These elements are required for activation or repression as well as maintenance of the expression (4, 8) and could be acting either as enhancer/silencers or as chromatin isolators that regulate the accessibility to DNA of transcription factors.

The *Contrabithorax* (*Cbx*) alleles are dominant gain-of-function mutations that result in the ectopic expression of UBX proteins in the wing. This ectopic expression results in the transformation of wing to haltere and, in some cases, mesonotum to metanotum. The *Cbx* alleles are mainly caused by breakpoints in the upstream or downstream regulatory regions or in the case of the *Cbx*<sup>1</sup> allele by the transposition of the *pbx cis* regulatory element into the second intron of *Ubx*. Similar phenotypes to the *Cbx* alleles can be obtained in constructs expressing *Ubx* regulated by some *abx*, *pbx*, and *bx* regulatory elements. This minigene is capable of partially rescuing the *Ubx* mutant phenotype, but its lack of restricted expression results in *Cbx* phenotypes (22). *Cbx* phenotypes can also be obtained by mutation of the *Ubx* repressors encoded by the *Polycomb* genes.

The regulatory elements of *Ubx* act mainly in *cis*. However, in certain experimental conditions, it has been shown that they can activate the promoter of the homologous *Ubx* gene. This trans complementation is sensitive to the disruption of pairing between homologous chromosomes and has been termed “transvection” (23). Transvection occurs in other loci of *Drosophila* and happens due to the pairing of the homologous chromosomes during interphase. The chromatin binding protein encoded by the *zeste* gene is required for transvection, and mutations in *Polycomb* affect the degree to which transvection occurs (24, 25).

The Hox genes of *Drosophila* have cross-regulatory interactions. Genes expressed in more posterior segments repress the transcription of more anterior genes. Thus, UBX protein represses the transcription of genes of the Antennapedia Complex, while the ABD-A and ABD-B proteins repress *Ubx* transcription (26). Generalized expression of UBX protein in the embryo results in the transformation of head and thoracic segments into A1, the segment where homogeneous and high UBX protein levels are expressed. In contrast, this ectopic expression has only mild defects in regions posterior to A1. This difference in the effects of ectopic UBX expression is due to the expression of ABD-A and ABD-B proteins in the posterior segments. A similar effect occurs when ANTP and UBX are expressed simultaneously. In this case, expression of both proteins does not give the expected intermediate phenotype but the phenotype of UBX expression alone (27). This epistatic phenomenon at the level of protein expression in which the simultaneous expression of two HOX proteins results in the same phenotypic outcome as expressing the most posterior one has been called in *Drosophila* phenotypic suppression and, in vertebrates, posterior prevalence (28).

From the evolutionary point of view, the *Ultrabithorax* gene predates the evolution of insects and their derived body plan (29). Real *Ubx* homologs (orthologs) have been cloned from several species of arthropods including many insects, crustaceans, and onychophorans (30, 31). In many instances, changes in the spatial expression of *Ubx* orthologs in different species correlate with the morphological differences observed between species (32, 33). In the flour beetle *Tribolium castaneum*, there is strong evidence supporting that mutations in the *Ultrathorax* gene that transform the morphology of abdominal segments towards that of thoracic segments are mutations in the *Ubx* ortholog (34). Experiments expressing the *Ubx* ortholog of *Onychophora* in *Drosophila* suggest that

although many of the functions of the protein have been conserved some have diverged. This divergence is partly due to interaction of the *Ubx* orthologs with different co-factor proteins (35).

## Bibliography

1. M. E. Akam (1983) *EMBO J.* **2**, 2075–2084.
2. R. A. H. White and M. Wilcox (1985) *EMBO J.* **4**, 2035–2043.
3. J. Castelli-Gair and M. Akam (1995) *Development* **121**, 2973–2982.
4. M. Bienz (1994) *Trends Genet.* **10**, 22–26.
5. E. B. Lewis (1978) *Nature* **276**, 565–570.
6. J. Botas and L. Auwers (1996) *Mech. Dev.* **56**, 129–138.
7. G. S. Mastick et al (1995) *Genetics* **139**, 349–363.
8. M. Bienz and J. Müller (1995) *BioEssays* **17**, 775–784.
9. J. Castelli-Gair (1998) *Int. J. Dev. Biol.* **42**, 437–444.
10. H. D. Lipshitz, D. A. Peattie, and D. S. Hogness (1987) *Genes. Dev.* **1**, 307–322.
11. K. Kornfeld et al (1989) *Genes. Dev.* **3**, 243–258.
12. M. B. O'Connor et al (1988) *EMBO J.* **7**, 435–445.
13. A. J. López and D. S. Hogness (1991) *Proc. Natl. Acad. Sci. U.S.A.* **88**, 9924–9928.
14. H. M. Bomze and A. J. López (1994) *Genetics* **136**, 965–977.
15. V. Subramaniam B. H. M. and A. J. López (1994) *Genetics* **136**, 979–991.
16. R. S. Mann and D. S. Hogness (1990) *Cell* **60**, 597–610.
17. S. C. Ekker et al (1994) *EMBO J.* **13**, 3551–3560.
18. R. S. Mann and S.-K. Chan (1996) *TIG* **12**, 258–262.
19. M. D. Biggin and W. McGinnis (1997) *Development* **124**, 4425–4433.
20. I. Duncan (1987) *Annu. Rev. Genet.* **21**, 285–319.
21. M. Peifer and W. Bender (1986) *EMBO J.* **5**, 2293–2303.
22. J. Castelli-Gair, J. Müller, and M. Bienz (1992) *Development* **114**, 877–886.
23. E. B. Lewis (1954) *Am. Nat.* **88**, 225–239.
24. J. A. Kennison (1993) *Trends Genet.* **9**, 75–79.
25. S. Bickel and V. Pirrotta (1990) *EMBO J.* **9**, 2959–2968.
26. G. Struhl and R. A. H. White (1985) *Cell* **43**, 507–519.
27. A. González-Reyes et al (1990) *Nature* **344**, 78–80.
28. D. Duboule and G. Morata (1994) *Trends in Genet.* **10**, 358–364.
29. R. de Rosa et al (1999) *Nature* **399**, 772–776.
30. M. Averof and M. Akam (1993) *Curr. Biol.* **3**, 73–78.
31. J. K. Grenier et al (1997) *Curr. Biol.* **7**, 547–553.
32. M. Averof and N. H. Patel (1997) *Nature* **388**, 682–686.
33. D. L. Stern (1988) *Nature* **396**, 463–466.
34. R. L. Bennett, S. J. Brown, and R. E. Denell (1999) *Dev. Genes Evol.* **209**, 608–619.
35. J. K. Grenier and S. B. Carroll (2000) *Proc. Natl. Acad. Sci. U.S.A.* **97**, 704–709.

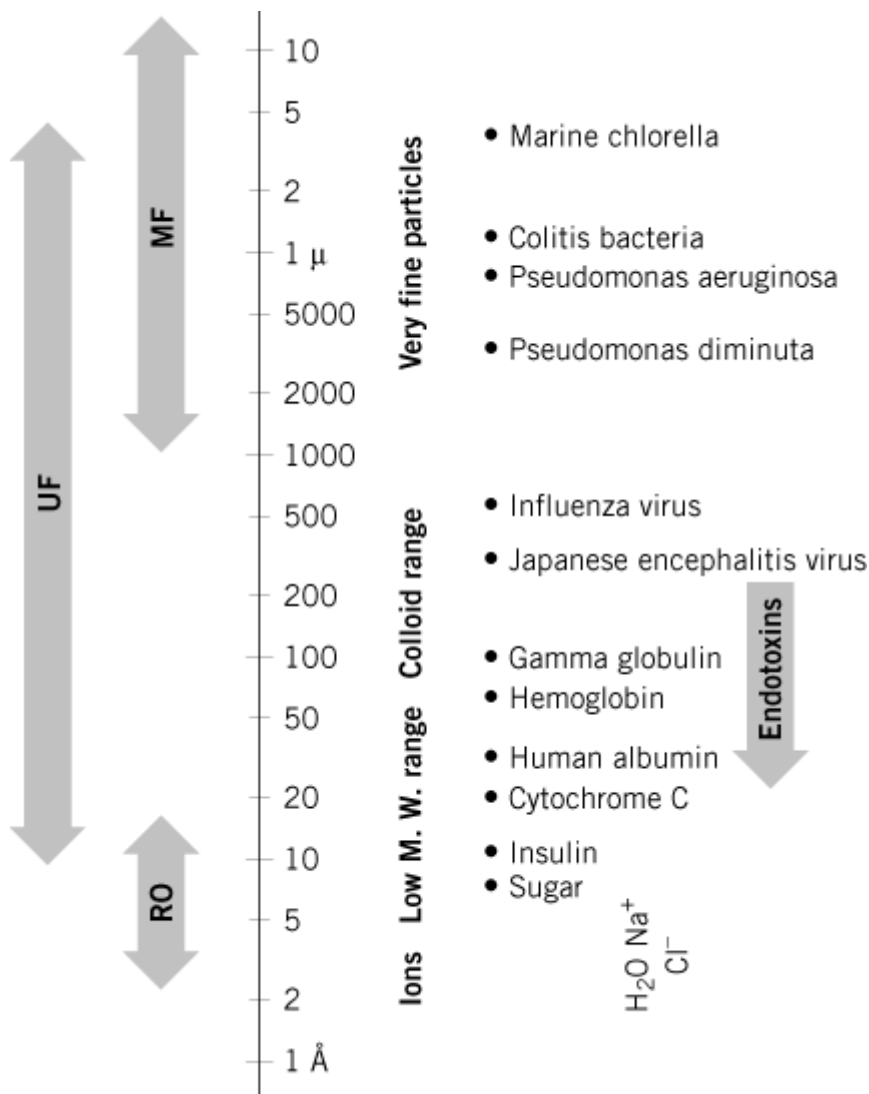
## Suggestions for Further Reading

36. P. A. Beachy (1990) A molecular view of the Ultrabithorax homeotic gene of *Drosophila*. *Trends Genet.* **6**, 46–51.
37. I. Duncan (1987) The bithorax complex. *Ann. Rev. Genet.* **21**, 285–319.
38. R. S. Mann and S.-K. Chan (1996) Extra specificity from extradenticle: the partnership between HOX and PBX/EXD homeodomain proteins. *TIG* **12**, 258–262.

## Ultrafiltration, Microfiltration

*Ultrafiltration* and *microfiltration* are pressure-driven processes that use a semipermeable membrane to separate dissolved and suspended materials on the basis of size; only molecules smaller than the pore size of the membrane may pass through it. The technique of ultrafiltration can deal with molecules with molecular weights in the range of 500 to 1,000,00 da, or approximately 10 to 100 Å (1 to 10 nm) in size. In contrast, microfiltration uses membranes in which the pore sizes range from 0.01 to 10 μm (Fig. 1). Ultrafiltration is suitable for concentrating or washing macromolecules such as [proteins](#), [peptides](#), polysaccharides, [lipids](#), synthetic polymers, **viruses**, and colloids, or the concentration of dilute samples of **nucleic acids**. The range of molecular sizes for which it is used distinguishes the process from *reverse osmosis*, in which even salts or other low molecular solutes are retained.

**Figure 1.** Application range for ultrafiltration, microfiltration, and reverse osmosis.



Ultrafiltration cannot discriminate between macromolecules that have sizes similar to one another, and for this reason the process is unsuited for the separation of soluble molecules whose molecular weights or effective sizes in solution ([Stokes radii](#)) differ by less than a factor of 6. For the separation of such molecules, alternate techniques such as [chromatography](#), [electrophoresis](#), or adsorption should be considered.

### 1. Dead-Ended and Cross-Flow Apparatus

The ultrafiltration process may be either *cross-flow* or *dead-ended*, where the feed stream is forced through the membrane and collected. An example of dead-end filtration is the disposable syringe end filters that use membranes bonded to a sealed plastic support plate. This unit has Luer-type connectors that permit its attachment to a syringe. The fluid to be filtered is placed in the syringe, manually driven through the filter, and collected. A technique particularly suited for the micro-purification of proteins, **antibodies**, and nucleic acids is shown in [Figure 2](#). Samples of up to 500  $\mu\text{l}$  may be placed in the microconcentrator module, which is placed into a centrifuge and spun. The centrifugal force generated forces the filtrate through the membrane and concentrates the molecules too large to pass through it. For larger sample volumes, stirred cells may be used, such as those developed by Amicon ([Fig. 3](#)). In this technique, a flat sheet of membrane is placed on an appropriate support structure and spacer that forms part of the cell assembly. The fluid is added and agitated with a magnetic stirrer. The cell is pressurized, either by air or nitrogen, and the filtrate is forced through the membrane. The use of air pressure avoids the damage of highly labile products that can occur during pumping.

**Figure 2.** Centrifugal concentrator (MicroconR). After centrifugation, the module is inverted and recentrifuged for more complete recovery of the concentrate in the vial. (Courtesy of Millipore Corporation, Bedford, MA.)



**Figure 3.** Stirred cell for ultrafiltration. Pressurized vessel with magnetic stirring bar. The membrane in the form of a disk is situated under the stirring bar. (Courtesy of Millipore Corporation, Bedford, MA.)



The main disadvantage of such techniques is membrane fouling. As the filtrate is forced through the membrane, the retained particles accumulate on the membrane surface. This layer is the *concentration polarization layer*; its development will depend on several factors, such as temperature, solubility, and pH. As the thickness of this layer increases, the efficiency of filtration diminishes. Control of such membrane fouling may be achieved by the use of a thin channel device, in which the feed is forced through a narrow channel that spirals outward over the surface of the membrane. The flow in this narrow channel generates high shear rates, leading to improved solute transfer across the membrane.

For the filtration or concentration of solutions at process scale, hollow fiber or flat sheet cross-flow devices are used (Fig. 4). In cross-flow filtration modules, the feed flow is parallel to the membrane and perpendicular to the filtrate flow. The solvent and solutes smaller than the pore size of the membrane pass through the membrane to form the ultrafiltrate, or permeate. In such devices, the retentate or concentrated solution flows across the membrane when in the form of a flat sheet, or along the fibers. The permeate or filtrate passes across the membrane and flows out of the module for collection. Cross-flow systems are inherently more complex, but they have high rates of mass transfer, due to high shear rate or tangential velocity, as well as turbulence adjacent to the membrane.

**Figure 4.** Hollow fiber ultrafiltration modules for process scale applications. (Courtesy of Millipore Corporation, Bedford, MA.)



## 2. Ultrafiltration Membranes

The membranes that are used in ultrafiltration or related processes may be **hydrophobic** or **hydrophilic**, uniform in their cross section or asymmetric, with a dense thin (0.1 to 1.5  $\mu\text{m}$ ) top layer supported by a porous sublayer containing macrovoids. The dense surface layer determines the membrane's convective transport properties, since it contains the pores, with the substructure providing mechanical strength. The membranes are made from either cellulose (cellulose ester or **nitrocellulose**) or synthetic materials, such as polysulfone and polyacrylonitrile, which exhibit relatively low protein binding.

The basic characteristics of membranes are represented by (1) the flux, (2) rejection of solutes, and molecular-weight cutoff. The flux  $J$  is defined as

$$J = \frac{Q}{A\Delta t} \quad (1)$$

where  $Q$  is the permeated amount,  $A$  the membrane area, and  $Dt$  the sampling time. The rejection  $R$  is calculated from the difference in the concentration of the solute in the feed solution  $C_F$  and the concentration in the permeate  $C_P$ , such that

$$R = 1 - \frac{C_P}{C_F} \quad (2)$$



The cutoff of a membrane is the maximum molecular-weight molecule that the membrane will allow through. For molecules or nucleotides, this parameter is generally defined as the maximum molecular weight, or the number of nucleotides in the DNA fragment, for which 90% of the molecules are retained by the membrane.

### 3. Solute and Solvent Transport Mechanisms During Filtration Processes

The membranes used in ultrafiltration and related processes may be considered a polymeric matrix in which pores are present. Although the actual morphology of the membranes is complex, considerable insights into the behavior of the membrane may be made by the application of a model to describe the transport phenomena. As the mean pore diameters of the membranes suitable for ultrafiltration and microfiltration differ, so do the expressions that are used to determine the flux. For microfiltration membranes, the equations used are based on capillary flow:

$$J = K_P \frac{\Delta P}{L} = K \frac{\Delta P}{\eta L} \quad (3)$$

$$J = \frac{\pi r^4 \Delta P}{8 \eta L} \quad (4)$$

The first of these equations is the D'Arcy flow equation, whereas the second is the Hagen Poiseuille equation. In these equations,  $J$  represents the volume flux,  $DP$  the pressure difference across the two sides of the membrane,  $K$  the permeability or the D'Arcy permeability  $K_p$ ,  $\eta$  the fluid viscosity, and  $p$  the universal constant.  $L$  is the length of the pore capillary, which can be taken to be equal to the thickness of the membrane, since at this level the membrane may be considered a series of parallel cylindrical pores, with a radius of  $r$ , perpendicular to the membrane surface. The equation may also be written to include a tortuosity factor  $T$  in the denominator, which provides a measure of the actual fluid path length through the membrane, relative to the thickness of the membrane.

For ultrafiltration membranes, the equations for the transport of solvent and solute may be expressed as

$$J_V = L_V(\Delta P - \Delta \Pi) \quad (5)$$

$$J_S = L_S(C_M - C_P) + (1 - \sigma)J_V C_S \quad (6)$$

where  $J_V$  and  $J_S$  are the fluxes of the solute and solvent, respectively,  $L_V$  and  $L_S$  their respective hydraulic permeabilities,  $C_M$ ,  $C_P$ , and  $C_S$  the solute concentrations in the membrane, permeate, and solution, respectively,  $DP$  the pressure difference across the membrane,  $\Delta \Pi$  the osmotic pressure, and  $\sigma$  the Staverman reflection coefficient.

Cross-flow filtration techniques are commonly used in which the mixture requiring separation is recirculated over the membrane surface. In this case, the flux rate will depend on the applied or transmembrane pressure  $DP$ , the mass transfer coefficient  $k$  of the device, the solute concentration in the formed gel  $C_G$ , and the bulk solute concentration  $C_B$  such that

$$J = k \log_e \left( \frac{C_G}{C_B} \right) \quad (7)$$

In the filtration of macromolecular solutions of low concentration by ultrafiltration, their osmotic pressures are low compared to the applied pressure and can be neglected. In high concentration solutions, the osmotic pressure effects may be significant, and the volume flux equation (5) requires modification to account for the osmotic pressure of the macromolecular solution. The osmotic pressure exerted by such solution is generally in the form

$$\Delta\Pi = AC + A_1C^2 + A_2C^3 \quad (8)$$

where  $A$ ,  $A_1$ , and  $A_2$  are osmotic virial constants and  $C$  is the concentration of the bulk macromolecular solution. The coefficient  $A$  describes Van't Hoff's limiting law for the osmotic pressure, which is applicable at very dilute concentrations, whereas  $A_1$  can be expressed in terms of the solute molecular weight.

In order to achieve high rates of mass transfer, it is necessary to have a high tangential velocity or shear rate and/or turbulence in the vicinity of the membrane. Theoretical expressions are available for the derivation of the mass transfer coefficients for differing membrane module and flow geometry configurations. Such expressions implicitly assume that the density, viscosity, and solute diffusivity are constant across the boundary layer.

Ultrafiltration and microfiltration are able to prepare reasonably concentrated samples of macromolecules, but they permit only a limited separation of the retained solutes from the smaller or more permeable components present in the solution needing to be filtered. In order to remove such smaller components more effectively from the retained species, the feed solution can be washed in a process known as [diafiltration](#).

#### Suggestions for Further Reading

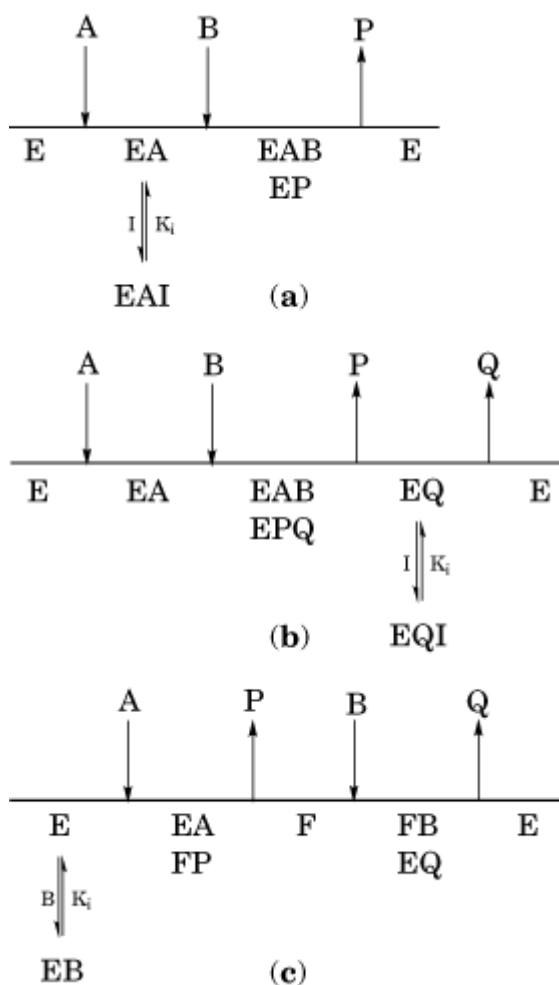
- L. J. Zeman and A. L. Zydney (1996) *Microfiltration and Ultrafiltration Principles and Applications*, Marcel Dekker, New York.
- J. A. Howell, V. Sanchez, and R. W. Field (1993) *Membranes in Bioprocessing: Theory and Applications*, Chapman Hall, London.
- Y. Osada and T. Nakagawa (1992) *Membrane Science and Technology*, Marcel Dekker, New York.
- P. Meares (1976) *Membrane Separation Processes*, Elsevier, Amsterdam.

## Uncompetitive Inhibition

Linear uncompetitive inhibition of an [enzyme](#) occurs when the variable substrate and the inhibitor combine with different enzyme forms and no reversible connection exists between the two points of addition (see [Product Inhibition](#) and [Dead-End Inhibition](#)). This type of inhibition is observed with a two-substrate reaction with an ordered **kinetic mechanism** when an inhibitory analogue of substrate B (I) combines with EA and A is the variable substrate (Fig. [1a](#)). The inhibition by I would be linear **competitive** with respect to B because B and I combine separately with EA. But it would be linear uncompetitive with respect to A, because I and A combine with different enzyme forms and there is no *reversible* connection between the reaction of A with E and of I with EA. Increasing concentrations of A can only increase the concentration of the EA complex with which I reacts. The same type of inhibition would be observed if an inhibitory analogue of substrate B combined only with the EQ complex on the release side of an ordered Bi—Bi mechanism (Fig. [1b](#)). This inhibition

would be linear uncompetitive with respect to both substrates.

**Figure 1.** Enzyme kinetic mechanisms that give rise to linear uncompetitive inhibition.



Linear uncompetitive inhibition is commonly observed with Ping-Pong mechanisms with structural analogues of either of the two substrates (Fig. 1c). When I is an analogue of A, it would compete with A for binding to E and thus cause linear competitive inhibition. The inhibition with respect to B would be linear uncompetitive because B and I combine with different enzyme forms and no reversible connection exists between the combination of I with E and of B with F. The release of P or Q, at what can be considered to be zero product concentration when initial velocities are being measured, breaks the connection between the points of addition of I and B. Similar results would be obtained with an inhibitory analogue of substrate B.

The general form of the equation to describe linear uncompetitive inhibition is shown in equation (1):

$$v = \frac{VA}{K_a + A \left(1 + \frac{I}{K_{ii}}\right)} \quad (1)$$

The reciprocal form of the equation is illustrated in equation (2):

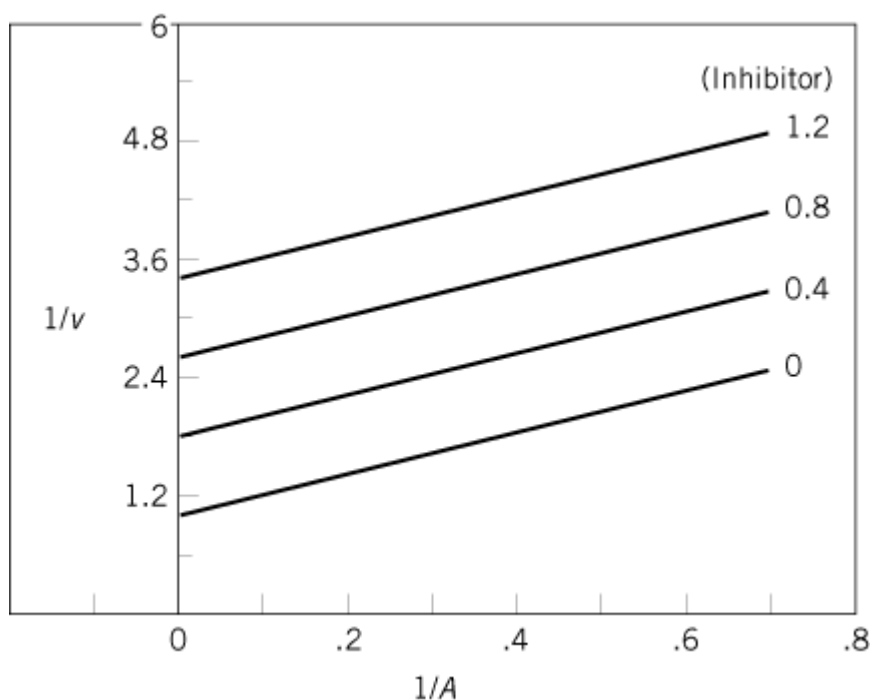
$$\frac{1}{v} = \frac{K_a}{V} \frac{1}{A} + \frac{1}{V} \left( 1 + \frac{I}{K_{ii}} \right) \quad (2)$$

In equations (1) and (2),  $V$  represents the maximum velocity of the reaction,  $K_a$  and  $A$  are the  $K_m$  (Michaelis constant) and concentrations of substrate  $A$ , respectively,  $I$  is the concentration of the inhibitor, and  $K_{ii}$  denotes the inhibition constant associated with the intercepts of double-reciprocal [Lineweaver-Burk plots](#). Because the slope of the curves is independent of the concentration of  $I$  (eq. (2)), a plot of  $1/v$  against  $1/A$  would consist of a family of parallel straight lines (Fig. 2). A replot of the vertical intercepts against  $I$  would yield a straight line according to

$$\text{Intercept} = \frac{1}{VK_{ii}} I + \frac{1}{V}$$

and the horizontal intercept gives the value of  $K_{ii}$ . For the mechanism illustrated in Figure 1a, the  $K_{ii}$  value would not be a true dissociation constant ( $K_i$ ) but rather an apparent value equal to  $K_i(1 + B/K_b)$ . The true value could be calculated using the values of the Michaelis constant and the fixed concentration of  $B$ . The best estimate of  $K_{ii}$  would be obtained by least-squares fitting of the inhibition data to the general equation for linear uncompetitive inhibition (1).

**Figure 2.** Initial velocity pattern for linear uncompetitive inhibition as illustrated by a double reciprocal (Lineweaver-Burk) plot.



## Bibliography

1. W. W. Cleland (1979) *Methods. Enzymol.* **63**, 103–138.

## Unfolded Protein Response

When cells accumulate substantial amounts of [unfolded protein](#) in their [endoplasmic reticulum](#) (ER), as when expressing a mutant form of a secreted protein that cannot fold efficiently or under stress conditions (see **Stress Response**), the ER grows in volume. It also changes from a network surrounding the nuclear envelope to one that fills the peripheral cytoplasm. In yeast, the genes for many ER proteins share a 22-bp transcriptional element that mediates this coordinate unfolded protein response (UPR) (1). A special [transcription factor](#), Hac1p, binds to this element and is required for mounting the UPR. Hac1p is regulated in a novel and complex fashion. Its normal product is synthesized constitutively, but does not induce UPR, either because it is rapidly targeted for **protein degradation** or because its [messenger RNA](#) is sequestered in the [nucleus](#) (2). A more stable form of the factor, Hac1<sup>i</sup>, is produced as a result of nonconventional [RNA splicing](#) of the Hac1 mRNA precursor, in a reaction that does not require the normal RNA processing machinery, but instead requires the [tRNA ligase](#) encoded by the gene RGL1 (3, 4). The endonuclease activity involved in this splicing reaction resides, surprisingly, in the cytosolic domain of the ER transmembrane protein Ire1 (5). This protein is a **serine/threonine kinase**, whose activity is required to induce the UPR in yeast (6, 7). The luminal, N-terminal **domain** of Ire1 can sense, by an unknown mechanism, the folding status in the lumen of the ER and mediate dimerization of the protein in the plane of the ER membrane. Dimerization activates **trans-phosphorylation**, in which each monomer phosphorylates the other; this in turn is thought to activate the endonucleolytic activity and initiate the [alternative splicing](#). Significantly, yeast mutants in the various components of the UPR are also inositol auxotrophs, thus providing a genetic link between the induction of ER **protein biosynthesis** and membrane [lipid](#) synthesis (8). Because ER biogenesis is induced in higher eukaryotes and yeast by similar physiological demands, it is expected that the [signal transduction](#) cascade to the [genome](#) would be similar, yet the mammalian homologues of the yeast machinery are yet to be characterized.

### Bibliography

1. K. Mori, A. Sant, K. Kohno, K. Normington, M. J. Gething, and J. F. Sambrook (1992) *EMBO J.* **11**, 2583–2593.
2. T. Kawahara, H. Yanagi, T. Yura, and K. Mori (1997) *Mol. Biol. Cell* **8**, 1845–1862.
3. J. S. Cox and P. Walter (1996) *Cell* **87**, 391–404.
4. C. Sidrauski, J. S. Cox, and P. Walter (1996) *Cell* **87**, 405–413.
5. C. Sidrauski and P. Walter (1997) *Cell* **90**, 1031–1039.
6. J. S. Cox, C. E. Shamu, and P. Walter (1993) *Cell* **73**, 1197–1206.
7. K. Mori, W. Ma, M. J. Gething, and J. Sambrook (1993) *Cell* **74**, 743–756.
8. J. S. Cox, R. E. Chapman, and P. Walter (1997) *Mol. Biol. Cell* **8**, 1805–1814.

### Suggestions for Further Reading

9. J. S. Cox, C. E. Shamu, and P. Walter (1993) Transcriptional induction of genes encoding endoplasmic reticulum resident proteins requires a transmembrane protein kinase. *Cell* **73**, 1197–1206.
10. K. Mori, W. Ma, M. J. Gething, and J. Sambrook (1993) A transmembrane protein with a cdc2+/CDC28-related kinase activity is required for signaling from the ER to the nucleus. *Cell* **74**, 743–756.
11. C. Sidrauski and P. Walter (1997) The transmembrane kinase Ire1p is a site-specific endonuclease that initiates mRNA splicing in the unfolded protein response. *Cell* **90**, 1031–1039.

## Unfolded Proteins

Disrupting the native [protein structure](#) upon **denaturation** implies disappearance of the specific contacts between the groups of the protein and a decreased compactness of the [polypeptide chain](#) conformation, that is, its unfolding. The completeness of unfolding of the polypeptide chain of proteins upon denaturation is one of the most debated subjects in protein science. It is an important problem because in all theoretical considerations of the mechanism of formation and stabilization of the native protein structure, the completely unfolded polypeptide chain in a [random-coil](#) conformation is a standard reference state. The ideal unfolded state, however, can probably not be achieved in practice. Extreme conditions can change the balance of forces that maintain the native protein structure and can destroy specific interactions between clusters of protein groups, but they certainly cannot eliminate completely all the interactions between the great variety of groups.

It has been believed that proteins are most unfolded in concentrated solutions of [urea](#) or **guanidinium** chloride (GdmCl) (see [Denaturation, Protein](#)), because these solutes significantly increase the intrinsic **viscosity** of proteins at room temperature (an indicator of the hydrodynamic **volume** of the polypeptide chain) and decrease almost to zero the **circular dichroic** (CD) ellipticity at 220 nm (an optical indicator of the content of helical structure). According to these criteria, the least unfolded are thermally denatured proteins, proteins at high temperatures. The intrinsic viscosity of unfolded proteins at high temperatures is significantly smaller than in concentrated solutions of denaturants at room temperature, and thermally unfolded proteins show some residual CD ellipticity (1). From the observed temperature-dependence of the apparent equilibrium constant for unfolding, it has also been assumed that denaturation by denaturants leads to a much greater **heat capacity** increment than denaturation by temperature. Because the denaturation heat capacity increment is regarded as an index of [nonpolar](#) group exposure to [water](#), this fact was considered evidence for more complete unfolding of proteins by denaturants than by temperature (1). However, it was found later that the intrinsic viscosity of unfolded polypeptide chains depends markedly on temperature and decreases significantly with increasing temperature (2, 3). This is the main reason that the intrinsic viscosity of heat-denatured proteins measured at high temperatures is significantly lower than that in concentrated solutions of GdmCl or urea measured at room temperature. The same was found for the CD ellipticity of proteins in concentrated solutions of denaturants. It also depends on temperature (4), and at high temperatures it has the same magnitude as that of the heat-denatured protein in water (3). Furthermore, direct **calorimetric** measurements showed that the heat capacity increase observed upon protein denaturation by urea or GdmCl is greater than that by thermal denaturation only because denaturants interact strongly with proteins and this interaction is enthalpic (5, 6). With increasing temperature, denaturant molecules bound to proteins dissociate gradually, absorbing heat, which increases the value of the apparent heat capacity increment (see [Denaturation, Protein](#)). If the heat effect of gradual dissociation of denaturants on temperature rise is taken into account, the real heat capacity increment of protein denaturation by denaturants and by temperature are similar. With increasing temperature, polypeptide chains in solutions of denaturants come closer to the random-coil conformation because bound denaturant molecules dissociate and because the dissipative forces of thermal motion increase. At sufficiently high temperatures, denaturants do not much affect the state of a polypeptide chain, and it becomes close to that of a heat-denatured protein without denaturants. There may still be some nonspecific, chaotic interactions between the groups of the polypeptide chain, particularly the attractive **hydrophobic** interactions between nonpolar groups, which increase with increasing temperature and squeeze the hydrodynamic volume of the polypeptide chain. But they do not form stable, ordered clusters of groups. Thus, concentrated solutions of urea and GdmCl cannot be regarded as q-solvents, in which the interactions of one group with other groups of a polymer are the same as with the solvent, and the polypeptide chains in these solutions cannot be considered ideal random coils (7, 8). Nevertheless, calorimetric

measurements of the **enthalpies** of unfolding of small globular proteins by various means (GdmCl, acid, and temperature) showed that if the heats of protonation of groups upon unfolding and of solvation by denaturants are properly taken into account, the enthalpy of unfolding is a universal function of temperature (5, 9). This shows that, for small proteins that do not aggregate, the total energy of residual interactions between the groups of denatured proteins is small. Therefore, because there is no way to determine experimentally the energy of these residual interactions in the denatured protein, the denatured state of a protein is usually regarded as the reference unfolded state close to the truly random-coil state.

With increase of molecular weight and length of a polypeptide chain, the local **effective molarities** of groups in the unfolded protein increase, as do their attractive interactions, which significantly complicates the situation. This is especially so when a large protein consists of several domains that unfold under different denaturing conditions. Then, the denatured state of the protein, the state that lacks specific functional activity, does not always correspond to the unfolded state. It might be a partly unfolded state in which some **domains** are unfolded, whereas others preserve their native structure (see [Denaturation, Protein](#)). This partly unfolded state of multidomain proteins is sometimes called a “**molten globule**” state because it is rather compact. However, this designation is hardly justified because the term “molten globule” implies a global change (disruption) of the protein molecular structure, not just of some of its parts (10-12).

### Bibliography

1. C. Tanford and K.C. Aune (1970) *Biochemistry* **9**, 206–211.
2. F. Ahmad and A. Salahuddin (1974) *Biochemistry* **13**, 245–249.
3. P.L. Privalov, et al. (1989) *J. Mol. Biol.* **205**, 737–750.
4. M.L. Tiffany and S. Krimm (1972) *Biopolymers* **11**, 2309–2316.
5. W. Pfeil and P.L. Privalov (1976) *Biophys. Chem.* **4**, 33–40.
6. G.I. Makhatadze and P.L. Privalov (1992) *J. Mol. Biol.* **226**, 491–505.
7. M.L. Tiffany and S. Krimm (1973) *Biopolymers* **12**, 575–587.
8. W. Pfeil and P.L. Privalov (1976) *Biophys. Chem.* **4**, 23–32.
9. K.A. Dill and D. Shortle (1991) *Annu. Rev. Biochem.* **60**, 793–825.
10. T.E. Creighton and D. Shortle (1994) *J. Mol. Biol.* **242**, 670–682.
11. C.M. Dobson, P.A. Evans, and S.E. Radford (1994) *Trends Biochem. Sci.* **19**, 31–37.
12. P.L. Privalov (1995) *J. Mol. Biol.* **258**, 707–725.

### Suggestions for Further Reading

13. K.A. Dill and D. Shortle (1991) Denatured state of proteins, *Annu. Rev. Biochem.* **60**, 795–825.
14. F.M. Richards (1992) "Folded and unfolded proteins", In *Protein Folding* (T.E. Creighton, ed.), Freeman, New York, pp. 1–58.
15. C. Tanford (1968) Protein denaturation. Part A and Part B. *Adv. Protein Chem.* **23**, 121–275.
16. C. Tanford (1970) Protein denaturation. Part C. Theoretical models for the mechanism of denaturation, *Adv. Protein Chem.* **24**, 1–95.

### Unidentified Reading Frame

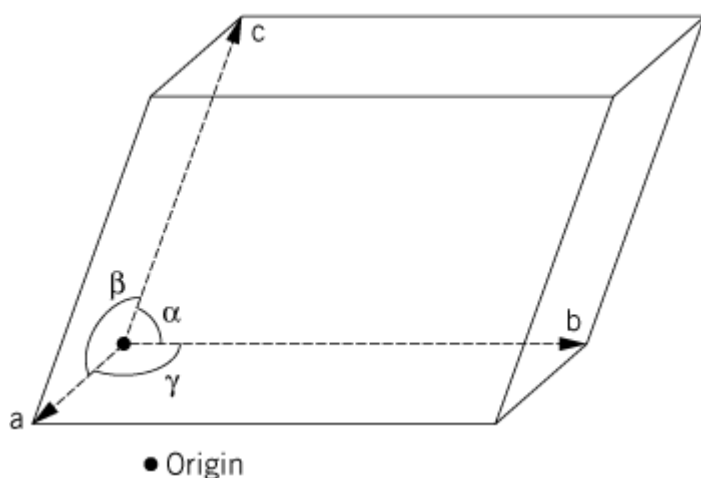
The term *unidentified reading frame* (URF) was previously used by scientists working on the DNA

from [mitochondria](#) and [chloroplasts](#), but it has now been abandoned for the term open reading frame (ORF), which applies to all classes of DNA. By definition, an ORF is unidentified. When it is identified, this means that its product is known and it is now a *bona fide* **gene**.

## Unit Cell

A crystal is characterized by the regular and periodic arrangement of its parts, which are ions, atoms, or molecules (see [Crystallography](#)). In this regular packing, three repeating vectors **a**, **b**, and **c** can be recognized with angles  $\alpha$ ,  $\beta$ , and  $\gamma$  between them. The three vectors define a unit cell in the crystal (Fig. 1). The unit cell can be selected in many different ways, but it is conventional to choose it according to certain rules, given in Ref. 1 (see [Space Group](#)). In protein crystals, the unit cells have dimensions from 30 Å up to several hundred Å. It should be stressed that vectors **a**, **b**, and **c** are imaginary vectors in the crystal and indicate only the periodicity of the crystal lattice. The unit cell in many crystal lattices, is composed of two or more identical structures, related by the symmetry of the lattice and known as the [asymmetric unit](#). The great advantage of the unit cell concept is that only the content of one unit cell or asymmetric unit needs to be described in a report on the crystal structure (see [X-Ray Crystallography](#)).

**Figure 1.** One unit cell in the crystal lattice.



## Bibliography

1. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## Suggestion for Further Reading

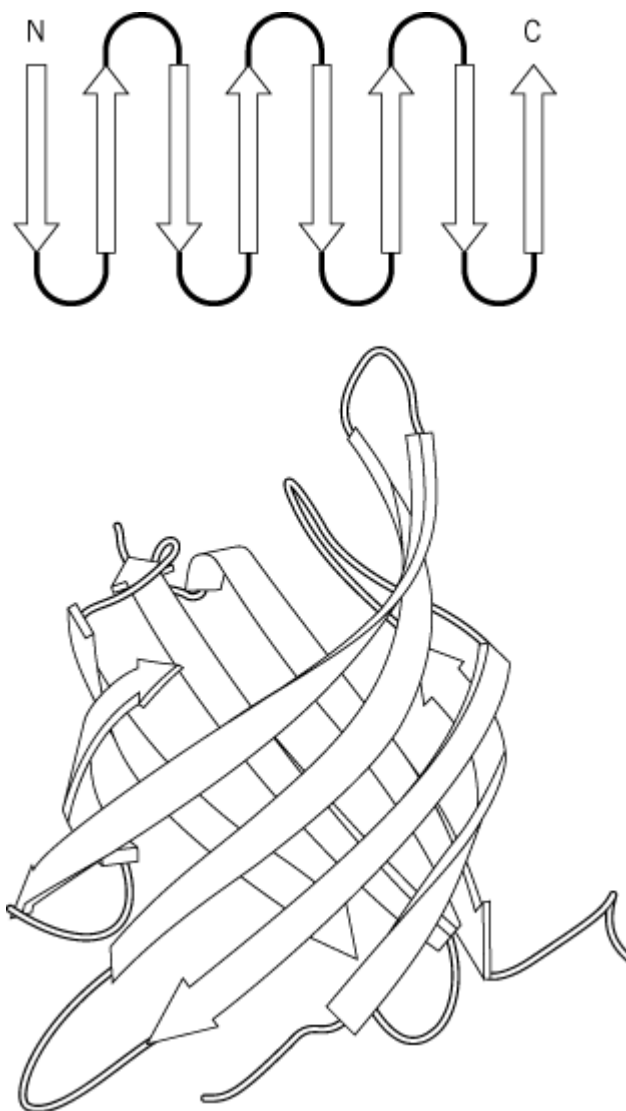
2. The International Union of Crystallography (1992) *International Tables for Crystallography*, Vol. A, Kluwer Academic, Dordrecht, Boston, London. Contains detailed descriptions of the shapes and symmetries of the unit cells in the 230 space groups.



## Up-And-Down b-Barrel

Up-and-down b-barrel describes the topology of a type of [b-barrel motif](#) found in some [protein structures](#), in which the protein [polypeptide chain](#) folds into two [b-sheets](#) that pack against each other to form a cylindrical or barrel shape (Fig. 1, see top of next page). The characteristic feature of up-and-down b-barrels is that all connections between the 8 or 10 [b-strands](#) are **hairpins**, so that b-strands adjacent in sequence are also adjacent in the structure. This topology is equivalent to a [b-meander](#). The barrel is closed by interactions between the first and last b-strands, and a **hydrophobic** cavity is formed in the interior. Proteins having this type of structure often function as **transporters** or solubilizers of hydrophobic ligands (such as retinoids, fatty acids, and bile salts) that bind in the interior of the barrel.

**Figure 1.** Structure of up-and-down b-barrels. **(Top)** Topology of the b-strands in the up-and-down b-barrel, with hairpin connections for all the b-strands. **(Bottom)** Backbone representation of the three-dimensional structure of a typical up-and-down b-barrel protein, retinol binding protein (1). This was generated using Molscrip (2).



[See also [Beta-Sheet](#) and [Antiparallel Beta-Barrel Motifs](#).]

## Bibliography

1. S. W. Cowan, M. E. Newcomer, and T. E. Jones (1990) *Proteins, Struct. Funct. Genet.* **8**, 44–61.
2. P. J. Kraulis (1991) *J. Appl. Crystallogr.* **24**, 946–950.

## Suggestions for Further Reading

3. C. Branden and J. Tooze (1991) *Introduction to Protein Structure*, Garland, New York.
4. J. LaLonde, D. A. Bernlohr and L. J. Banaszak (1994) The up-and-down -barrel. *FASEB J.* **8**, 1240–1247.

## Upstream

**Nucleotide sequence** elements are designated as being *upstream* or [downstream](#), with respect to the direction of [transcription](#) of a given gene. Initiation of transcription of the **gene** and initiation of [translation](#) of [messenger RNA](#) depend on numerous sequence signals located upstream of the [initiation codon](#), usually within the nearest 200–500 bp of the 5' region, but sometimes as far as several kilobase pairs away. Sequence features relevant to translation are located between the transcription and translation starting points, whereas transcription signals reside further upstream. **Promoters** are usually located immediately upstream of the transcription starting points and comprise a separate category of upstream elements.

In prokaryotes, the first upstream element is the [Shine-Dalgarno sequence](#), also called the **ribosome-binding site** (see [Translation](#)). It is located within a few bases from the translation starting point and has the **consensus sequence** GGAGG. This sequence is important for initiation of translation. It base-pairs with a complementary sequence at the 3' end of the small-subunit ribosomal RNA. Upstream of the prokaryotic promoters are **operator** sequences—binding sites for [repressors](#) of transcription. In many instances, the operator sequences are located rather far away—200–300 bp—so that the interaction of the corresponding repressors with the transcription complex requires bending of the DNA, with formation of a loop ([1](#), [2](#)).

In eukaryotes, the sequences upstream of the translation starting points are not necessarily the same as in the **genomic DNA**. The 5' regions of the RNA transcripts before the initiation codon often contain intervening sequences (**introns**), which will be excised during processing of the RNA transcripts, together with any introns that interrupt the protein-coding sequences (see [RNA Splicing](#)).

Nucleotide sequences in eukaryotes immediately upstream of translation starting points are not related to the prokaryotic Shine-Dalgarno sequence. A different type of interaction of the mRNA starting region with the eukaryotic 18 S rRNA is suggested, in which certain trinucleotides that are the most frequent in the starting region of the mRNA make complementary contacts with the rRNA ([3](#)).

The 5' untranslated regions of eukaryotic genes harbor binding sites for an immense variety of species-, tissue-, gene family-, and gene-specific [transcription factors](#) ([4](#)). Many different factors bind to DNA sites with similar **consensus sequences**, thus clustering in groups, or families. The CCAAT group of constitutive transcription factors is one example. Another frequent consensus sequence is GGGCGG, the binding site for Sp1 factor. Both such sequences are usually located in

the region 50–100 bp upstream of the transcription starting point.

Important families of sequences involved in the regulation of transcription in eukaryotes are the [enhancers](#) (5) and [silencers](#) (6). One general feature of these functional sites is the independence of their effect on their location, which can be either upstream or downstream of the protein-coding sequences, at distances from about 100 to several thousand bases away. Changes in their polarity in the sequence also are of no importance for the enhancement or silencing of transcription. The sequences themselves are 100–200 bp. The enhancer sequences harbor a variety of often crowded binding sites for various transcription factors, and they sometimes occur as several copies in tandem.

#### Bibliography

1. M. Irani, L. Orosz and S. Adhya (1983) *Cell* **32**, 783–788.
2. T. M. Dunn et al. (1984) *Proc. Natl. Acad. Sci. USA* **81**, 5017–5020.
3. D. R. Cavener and S. C. Ray (1991) *Nucl. Acids Res.* **19**, 3185–3192.
4. T. Heinemeyer et al. (1998) *Nucl. Acids Res.* **26**, 362–367.
5. H.-P. Müller and W. Schaffner (1990) *Trends Genet. Sc.* **6**, 300–304.
6. P. Laurenson and J. Rine (1992) *Microbiol. Rev.* **56**, 543–560.

#### Suggestions for Further Reading

7. B. Alberts, D. Bray, J. Lewis, M. Raff, K. Roberts and J. D. Watson (1994) *Molecular Biology of the Cell*, rd ed., Garland Publishing Inc., New York.
8. D. S. Latchman (1992) *Eukaryotic Transcription Factors*, Academic Press, London.

## Urea

Urea is one of the most common protein **denaturants**. Classically, it is used at 8 *M* concentration, which unfolds most globular proteins. Its limit of solubility at room temperature is ~10 *M*. As a denaturant, urea has the same efficacy as **guanidinium acetate** and is, thus, a much weaker denaturant than **guanidinium chloride**. It has the advantage of being a neutral molecule, which precludes the various effects of ionization. Urea acts principally by the formation of [hydrogen bonds](#) with protein [peptide](#) groups, although it also has some **hydrophobic** character. At lower concentrations, 0.5 to 1.0 *M*, it is used as a breaker of protein aggregates and solubilizer of proteins synthesized by recombination techniques, following the dissolution of [inclusion bodies](#). See also [Stabilization And Destabilization By Co-Solvents, Denaturants, stabilizers; Guanidinium Salts](#).

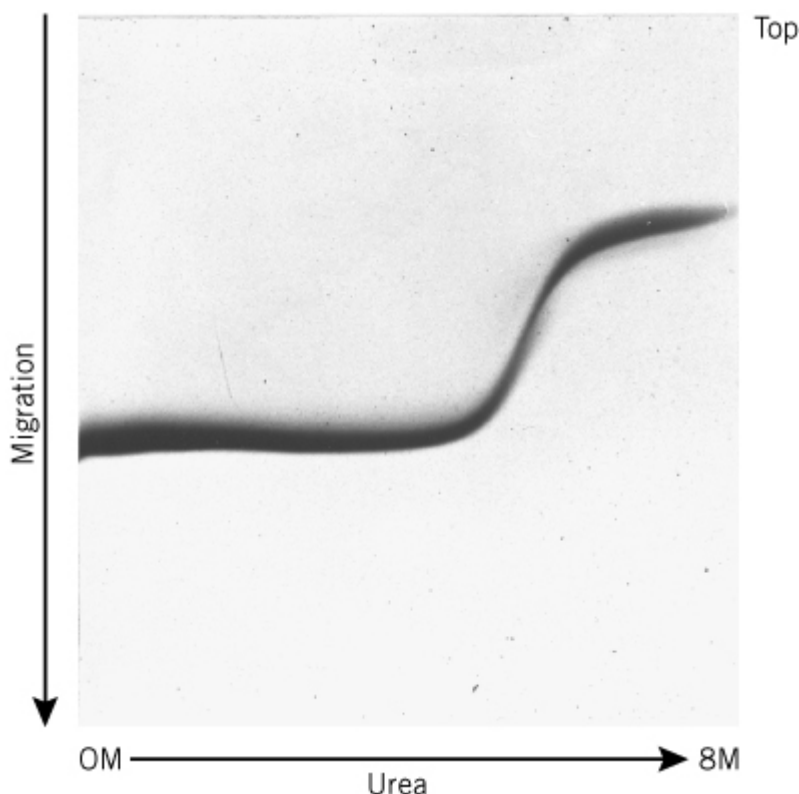
## Urea Gradient Gel Electrophoresis

The denaturant [urea](#) is commonly used in studies of **protein unfolding** and in the denaturation of **nucleic acids**. Because it is nonionic, urea may be incorporated into [polyacrylamide](#) gels for electrophoresis (**PAGE**). When incorporated as a concentration gradient perpendicular to the direction of electrophoretic migration, it is the most commonly used form of [transverse gradient gel](#)

[electrophoresis \(TGGE\)](#). The electrophoretic mobilities of proteins and nucleic acids depend on their conformations, so urea gradient gel electrophoresis is very useful in studying protein and nucleic acid structure and stability.

A concentration gradient of urea is incorporated horizontally into a slab polyacrylamide gel on which a sample of protein or nucleic acid is layered and then subjected to electrophoresis vertically. Each molecule in the sample migrates at a constant denaturant concentration, but the concentration varies continuously across the gel and the sample. The final pattern reveals the electrophoretic mobility of the sample at continuously varying urea concentrations. As with other types of TGGE, the result is a single continuous band across the gel, that has a sigmoidal change in mobility when the structure changes, if the change in structure is reversible and rapid on the timescale of the separation. Slow or irreversible transitions produce discontinuities in the band across the gel. An example of a rapid and reversible unfolding transition of a protein is illustrated in Figure 1. The unfolding transition is apparently two-state,  $N \leftrightarrow U$ , as shown by the single inflection point of the band through the unfolding transition region. The electrophoretic profile graphically displays the relative proportions of the N and U states present at equilibrium as a function of the urea concentration. The results obtained in this way are generally consistent with other studies of the urea-induced unfolding monitored in other ways.

**Figure 1.** Transverse urea gradient gel electrophoresis of cytochrome *c*. A sample of the folded protein was layered on the top of the polyacrylamide slab gel, which contained a linear gradient of urea from left to right. Electrophoresis at pH 4 was from top to bottom. The gel was stained with [Coomassie Brilliant Blue](#). At low urea concentrations, the protein remains folded and migrates relatively rapidly. At high urea concentrations, it unfolds and migrates more slowly. At about 6 M urea, the protein unfolds in a rapid and reversible unfolding transition. The same pattern is obtained starting with unfolded protein. This and the continuous band of protein throughout the unfolding transition region indicate that unfolding and refolding are rapid relative to the electrophoretic time at all urea concentrations. Therefore, the fraction of unfolding at equilibrium determines the migration rate at each urea concentration. The smooth shape of the band through the unfolding transition region and a single inflection point indicate that only two conformational states with different electrophoretic mobilities are present to significant extents. That cytochrome *c* undergoes a rapid, two-state unfolding transition is known from other studies. Taken from Ref. 1.



Urea increases the rate of unfolding ( $N \leftrightarrow U$ ), and decreases that of refolding. In both cases, the urea concentration is proportional to the logarithm of each rate constant (see [Protein Folding In Vitro](#)). Consequently, the logarithm of the equilibrium constant for the reaction  $N \leftrightarrow U$  is proportional to the urea concentration. In this case, the **free energy** difference  $DG_f$  between N and U varies linearly with the urea concentration and is zero at the midpoint of the transition. Its value at other urea concentrations for a two-state unfolding transition are estimated by extrapolating from the transition region (see [Protein Stability](#)). In the case of equilibrium urea gradient gel patterns, like that of Figure 1, the values of  $DG_f$  are extrapolated through the midpoint of the transition defined by the slope of the electrophoretic band at its inflection point. The scale of the free energy difference is defined by the mobilities of the N and U states, which occur at  $DG_f = -2RT$  and  $+2RT$ , respectively.

In a multidomain protein, unfolding of the individual domains are observed as separate transitions if they occur at different urea concentrations. In proteins comprised of multiple identical subunits, the mobility initially increases upon dissociation to folded monomers, and then decreases as they unfold. With different subunits, the dissociation is generally irreversible as the different subunits separate in the gel, unless they just happen to have the same electrophoretic mobility. Similarly, dissociation of a small ligand is generally irreversible. If electrophoresis is carried out rapidly at low temperatures, information about the kinetics of unfolding and dissociation and of refolding are obtained. Multiple unfolded forms that are only slowly interconverted are detected, and the tendency of very slow folding forms to fold is obtained.

Double-stranded nucleic acids lose their structure at high concentrations of urea, which is often augmented with the denaturant **formamide**. The electrophoretic mobility of a double-stranded oligonucleotide usually decreases slightly with increasing urea concentrations, probably because of fraying of the double helix at the ends. Then the mobility decreases dramatically and very abruptly, as the relatively stiff double helix dissociates into two random polynucleotide chains. The decrease in mobility can be very dramatic. Consequently, migration into an increasing urea concentration gradient parallel to the direction of electrophoresis is used to determine at which concentration the mobility change occurs.

[Temperature gradient gel electrophoresis](#), is a related technique, which uses high temperature rather than urea to unfold macromolecules.

#### Bibliography

1. T. E. Creighton 1979 *J. Mol. Biol.* **129**, 235–264.

#### Suggestions for Further Reading

2. D. P. Goldenberg and T. E. Creighton 1984 Gel electrophoresis in studies of protein conformation and folding, *Anal. Biochem.* **138**, 1–18.
3. T. E. Creighton 1986 Detection of folding intermediates using urea-gradient electrophoresis, *Methods Enzymol.* **131**, 156–172.
4. L. S. Lerman, S. G. Fischer, I. Hurley, K. Silverstein, and N. Lumelsky 1984 Sequence-determined DNA separations, *Ann. Rev. Biophys. Bioeng.* **13**, 399–423.

#### Urokinase

Mammalian species have two plasminogen activators, urokinase- or urinary-type plasminogen activator (uPA) and tissue-type plasminogen activator (tPA), that are the products of separate genes and members of the chymotrypsin family of serine proteinases. Both are mosaic proteins with a modular structure similar to that of the blood coagulation proteinases to which they are closely related [see [Blood Clotting](#)]. The two plasminogen activators provide an excellent example of the regulatory potential made available by the modular construction of these proteinases. Although catalyzing the same reaction, *i.e.* specific hydrolysis of Arg560-Val561 of plasminogen, gross differences in the organization of the N-terminal modules of the plasminogen activators, together with subtle differences in the serine proteinase domain, lead to remarkably different functional properties. These are reflected in their biological roles with tPA being predominantly responsible for plasmin-catalyzed fibrin dissolution, whereas uPA is thought to be responsible for generating plasmin activity in the context of extracellular matrix degradation and thus to be involved in tissue remodeling and invasive cell migration (*e.g.* tumor invasion and metastasis; vascular remodeling).

Urokinase, or uPA, is composed of three independent domains; an N-terminal epidermal growth factor-like EGF domain, a kringle module and the serine proteinase domain. It is secreted by a wide variety of cell types as a single-chain glycoprotein of 411 residues ( $M_r$  54,000). This zymogenic form (pro-uPA) is activated by a single cleavage at Lys158-Ile159 giving rise to two disulfide-bridged polypeptides. This cleavage is catalyzed *in vitro* by a variety of trypsin-like proteinases including plasmin, which is the most likely activator *in vivo*. There has been considerable controversy concerning the intrinsic proteolytic activity of pro-uPA, with some investigators reporting an activity similar to that of the two-chain enzyme. However, the absolute level of the activity appears to be less than 0.1% that of the activated enzyme, although there is no consensus yet.

The structure of the catalytic domain of uPA, in complex with the peptide inhibitor Glu-Gly-Arg-CH<sub>2</sub>Cl, has been determined by X-ray crystallography at a resolution of 2.5 Å (1). The solution structure of the N-terminal domains has been solved by heteronuclear NMR (2). No constraints between the domains of uPA have been detected by NMR, demonstrating an exceptionally high degree of interdomain flexibility (3).

In recent years much of the interest surrounding uPA has been focused on the high-affinity interaction between its N-terminal EGF-like module and a specific cell-surface receptor termed uPAR. This is a 55 kD glycosylphosphatidylinositol-anchored (GPI-anchored) membrane protein composed of three domains structurally homologous to snake venom  $\alpha$ -neurotoxins (4). This binding regulates the function of uPA by focusing the enzyme to the pericellular region and providing a mechanism to enhance plasmin generation. The interaction of uPA with uPAR does not alter its activity *per se*, consistent with the observed independence of the EGF-like and serine proteinase domains, but requires the coincident cellular binding of plasminogen (5). This is thought to lead to the formation of ternary or higher order complexes which promote plasminogen activation by appropriate juxtaposition of the reactants and give a large reduction in the apparent  $K_m$  (Michaelis constant) for the reaction (6). The activation of uPAR-bound pro-uPA by cell-associated plasmin is also enhanced, leading to an efficient system of reciprocal zymogen activation and consequently a large amplification of plasmin generation. The precise molecular interactions underlying these kinetic effects have only been partially elucidated.

## Bibliography

“Urokinase” in , Vol. 4, pp. 2728–2729, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ; “Urokinase” in (online), posting date: January 15, 2002, by Vincent Ellis, University of East Anglia, School of Biological Sciences, Norwich, United Kingdom, NR4 7TJ.

1. G. Spraggon, C. Phillips, U.K. Nowak, C.P. Ponting, D. Saunders, C.M. Dobson, D.I. Stuart and E.Y. Jones (1995) *Structure*. **3**, 681–691.
2. A.P. Hansen, A.M. Petros, R.P. Meadows, D.G. Nettlesheim, A.P. Mazar, E.T. Olejniczak, R.X. Xu, T.M. Pederson, J. Henkin and S.W. Fesik (1994) *Biochemistry*, **33**, 4847–4864.
3. A.P. Hansen, A.M. Petros, R.P. Meadows and S.W. Fesik (1994) *Biochemistry*, **33**, 15418–15424.
4. M. Ploug and V. Ellis (1994) *FEBS Letters*. **349**, 163–168.
5. V. Ellis, M.F. Scully and V.V. Kakkar (1989) *J. Biol. Chem.* **264**, 2185–2188.
6. V. Ellis and K. Danø (1993) *J. Biol. Chem.* **268**, 4806–4813.

### Suggestions for Further Reading

7. V. Ellis and K. Danø (1998) "u-Plasminogen activator", in *Handbook of Proteolytic Enzymes*, A.J. Barrett, F. Woessner and N. Rawlings, eds., Academic Press, London, pp 177–184.

## V Genes

The genetic information that encodes the **variable regions** of immunoglobulins, as defined at the protein level, is split into several gene segments; these are not classical exons (see [Introns, Exons](#)), because they have [RNA splicing](#) signals on only one end, whereas the other end contains specific *recombination signal sequences* (RSSs) that allow them to recombine before becoming functional. This organization provides a major combinatorial basis for generating a large number of different immunoglobulin chains from a limited number of genes. Heavy chains are encoded as the result of the [gene rearrangement](#) of three such genes or gene segments, called V (for variable), D (for diversity), and J (for junction); light chains are comprised of only two, V and J. As a consequence, complementarity determining regions CDR 1 and CDR2 are V-gene encoded, whereas CDR3 results from a VDJ or VJ combination (see [Immunoglobulin Structure](#)). This makes the CDR3, especially that of the heavy chain, the most diversified. All these genes are clustered together with the corresponding constant region genes in three loci, IGH, IGK and IGL, located (in humans) at chromosome positions 14q2, 2p12, and 22q11, with sizes of 1350 kbp, 1800 kbp, and 1000 kbp, respectively.

The V gene accounts for most of the variable region, because it encodes for more than 90 amino acid residues, the remaining being contributed by D–J for the heavy chain and J for the light chain. The organization of each V gene has the same pattern, whatever the locus [ie, one leader sequence (L)], encoding a [signal peptide](#) that is cleaved off in the rough [endoplasmic reticulum](#) and is separated from the V gene by a short **intron**. Each L-intron-V unit has its own 5'-**promoter** region that includes several classical sequences, including a [TATA box](#). This organization is quite remarkable, because the V gene that will rearrange will carry over its own promoter region. As a result of this organization, V genes may have individual, very low levels of [transcription](#) before they rearrange, presumably because a transient accessibility of the DNA gives rise to germline V transcription. By contrast, the [enhancer](#) regions are not present at the 5' position of each V gene, but are split into one intronic and one 3' localization. In heavy chains, the enhancer regions are organized so as to be conserved after [class switching](#).

In humans, all three loci have been completely sequenced, so the actual number of V genes is known. There are, however, minor differences between individuals, due to polymorphism. The IGH

locus contains between 46 and 50 VH genes, 27 D, and 6 J gene segments. A significant fraction of the VH genes are [pseudogenes](#) (~40%). The same phenomenon is observed in the mouse. In birds, there is only one potentially functional gene, whereas all others are pseudogenes. In that case, diversity is mostly generated by **gene conversion**, indicating that pseudogenes behave as a reservoir of diversity; this might explain not only why they have been maintained, but also why they have conserved the same general structure, including CDR1 and CDR2. VH genes may be classified in seven families on a homology basis. These families are interspersed, a situation different from that of mice, where VH genes of the same family are clustered. The number of genes that defines a family greatly varies, from a single member to 30, depending on the family.

The IGK locus contains 76 Vk genes (only 34 of which are functional), and 5 Jk genes. The Vk genes may also be grouped in five homology families. The IGVK locus contains a large internal duplication, within which Vk genes have an inverted orientation. As a consequence, rearrangement of these genes is made by inversion instead of the classical deletion of the intervening sequences.

The IGL locus contains 52 Vl genes grouped into 10 families. A little more than 50% are functional. The JI–CI are tandemly organized in seven basic units, three of which are pseudogenes. One region is found duplicated in certain individuals, creating an allelic **polymorphism** in the population.

Chain diversity that results solely from the combinatorial association of the various V–D–J and V–J genes would be sufficient to generate about  $10^7$  distinct antibody molecules if their combination were completely random, but this is apparently not entirely true. For example, careful analysis of the expressed repertoire indicated that some V genes are frequently expressed, whereas some others are rarely used. This probably results from a selective event that takes place early in **B-cell** differentiation. It has been shown that the surrogate light chain, encoded by the I-like (14.1) and the VpreB genes, which associates with the m heavy chain in preB cells, does not combine equally with every heavy chain, so negative selection operates already at the preB stage of differentiation. This may be amplified upon selection of immature B cells before they leave the bone marrow.

See also entries [Antibody](#), [B Cell](#), [Class switch](#), [Gene Rearrangement](#), [Kappa and Lambda chain](#), [Recombinase](#), [Switch Region](#).

#### Suggestions for Further Reading

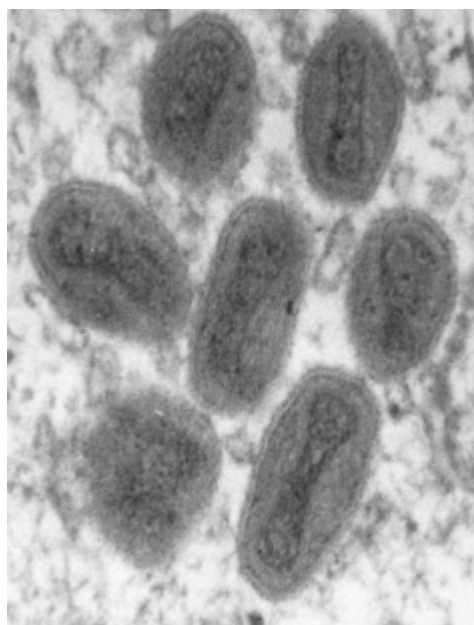
- F. Matsuda et al. (1993) Structure and physical map of 64 variable segments in the 3¢ 0.8-megabase region of the human immunoglobulin heavy-chain locus. *Nature Genet.* **3**, 88–94.
- G. P. Cook et al. (1994) A map of the human immunoglobulin VH locus completed by analysis of the telomeric region of chromosome 14q. *Nature Genet.* **7**, 162–168.
- T. Kirschbaum, R. Jaenichen, and H. G. Zachau (1996) The mouse immunoglobulin kappa locus contains about 140 variable gene segments. *Eur. J. Immunol.* **26**, 1613–1620.
- H. G. Zachau (1993) The immunoglobulin kappa locus-or-what has been learned from looking closely at one-tenth of a percent of the human genome. *Gene* **135**, 167–173.
- J.-P. Fripiat et al. (1995) Organization of the human immunoglobulin lambda light-chain on chromosome 22q11.2. *Hum. Mol Genet.* **4**, 983–991.
- V. Giudicelli et al. (1997) IMGT, the international ImMunoGeneTics database. *Nucleic Acids Res.* **25**, 206–211.
- M. Milili, C. Schiff, M. Fougereau, and C. Tonnelle (1996) The VDJ repertoire expressed in human preB cells reflects the selection of *bona fide* heavy chains. *Eur. J. Immunol.* **26**, 63–69.

## Vaccinia Virus



Vaccinia virus (VV) is a member of the genus *Orthopoxvirus* in the subfamily *Chordopoxvirinae* of the large family *Poxviridae*, comprising a group of complexed animal DNA **viruses**. VV is a “brick-shaped” particle 200 to 400 nm long, the external surface of which is ridged in parallel rows. The detailed structure of the particle is not known. Two kinds of particles are known: extracellular forms containing two **membranes** (EEVs—extracellular enveloped virions) and intracellular particles having only an inner membrane (IMVs—intracellular mature virions). The former is believed to be responsible for dissemination. Thin sections in **electron microscopy** reveal that the outer membrane surrounds the core, which is biconcave (dumbbell-shaped), with two lateral bodies (Fig. 1). The core is composed of a tightly compressed nucleoprotein that contains more than 10 **enzymes**, mostly concerned with **transcription** and **genome** replication.

**Figure 1.** Electron micrographs of thin-sectioned vaccinia virus particles. Magnification is  $\times 96,000$ .



The genome is a linear, double-stranded DNA molecule of approximately 190 kbp. The ends of the genome consist of a terminal hairpin loop with several **tandem repeat** sequences. The ends of the genome form terminal **inverted repeats** (ITRs), identical but oppositely oriented sequences. The majority of the essential **genes** map in the central part of the genome, while the genes nonessential for replication in tissue culture are located at the ends. There are  $\sim 200$  genes in the genome. Virus replication occurs in the cytoplasm, because the virus encodes almost all the functions necessary for genome replication.

The host nucleus contributes to the maturation step in forming virus particles. Receptors for the virus are not known. During penetration, uncoating occurs in two stages: (i) removal of the outer membrane and (ii) further uncoating to produce the naked core that passes into the cytoplasm.

Gene expression is temporarily regulated and divided into three phases: the early genes, intermediate genes, and late genes, each of which have their own characteristic **promoters**. The early genes comprise about half of the genome, are expressed before genome replication, and mainly encode enzymes involved in viral **DNA replication** and **transcription factors** for the intermediate genes.

Expression of the intermediate genes begins immediately after genome replication and produces late-gene transcription factors. The late genes encode major constituents of the virus particles, including early transcription factors, and are expressed after genome replication. Thus, the transcription factors synthesized in each phase govern the gene expression of the next phase.

Replication of the genome involves a self-priming mechanism, leading to the formation of large **concatemeric** branched structures, which are then resolved into unit genomes. Assembly of the progeny virions occurs in inclusions that are formed in the cytoplasm. During the maturation process, various modifications of the virion proteins occur, including **proteolytic** cleavage, acylation, and glycosylation.

IMV is transported on cables of [actin](#) and wrapped with the membranes comprising the **Golgi** cisternae. The virions move further through the cytoplasm to the cell surface and bud into environmental fluid to form EEV. VV produces many secretory proteins that resemble **cytokines** and their receptors, including [interferon](#), **interleukin-1**, [tumor necrosis factors](#), and a component of **complement**. These viral proteins probably compromise the host **immune system**, because mutant viruses lacking these genes have attenuated virulence, although they replicate in tissue culture as well as the wild-type virus.

VV is now used frequently as a live recombinant vector to express heterologous genes, which provides very useful tools for study of immunology and cell biology, as well as for development of novel vaccines.

#### Suggestions for Further Reading

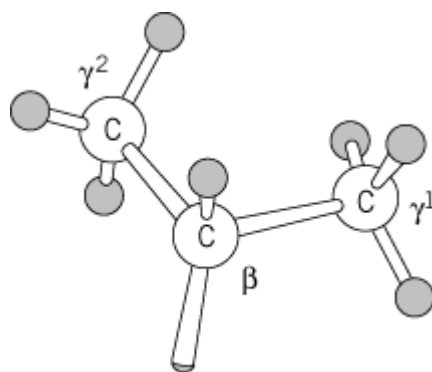
B. Moss (1996) "*Poxviridae: The Viruses and Their Replication*". In *Fields Virology*, 3rd ed. (B. N. Fields et al., eds.), Lippincott-Raven, Philadelphia, pp. 2637–2671.

E. Paoletti (1996) Applications of poxvirus vectors to vaccination: an update. *Proc. Natl. Acad. Sci. USA* **93**, 11349–11353.

## Valine (Val, V)

The [amino acid](#) valine is incorporated into the nascent [polypeptide chain](#) during **protein biosynthesis** in response to four **codons**—GUU, GUC, GUA, and GUG—and represents approximately 6.6% of the residues in the proteins that have been characterized. The valyl residue incorporated has a mass of 99.14 Da, a **van der Waals volume** of  $105 \text{ \AA}^3$ , and an [accessible surface area](#) of  $160 \text{ \AA}^2$ . Val residues are changed moderately frequently during [divergent evolution](#); they are interchanged in **homologous** proteins most frequently with [isoleucine](#), [leucine](#), [alanine](#), and [threonine](#) residues.

The Val side chain is [nonpolar](#) with no functional or reactive groups:



Consequently, Val is one of the most **hydrophobic** residues, and 54% of its residues are completely buried in native protein structures. Val is one of the residues least likely to adopt the **alpha-helical** conformation in model peptides, probably because its branched side chain is constrained to adopt only a limited number of conformations. In folded proteins, it occurs most frequently in [beta-sheet](#) type of **secondary structure**.

#### Suggestion for Further Reading

T. E. Creighton (1993) *Proteins: Structures and Molecular Properties*, 2nd ed., W. H. Freeman, New York.

### van der Waals Interactions

All gaseous atoms and molecules condense to liquid when cooled and then, on further cooling, to the solid state; the single exception is liquid helium. From this observation, we learn that even neutral atoms and molecules attract and repel each other at long and short distances, respectively. Repulsion at short distances gives finite volumes to solid materials. Not surprisingly, [polar](#) molecules attract each other, as they participate in attractive interactions, such as **hydrogen bonding** or dipole interactions, but so do perfectly spherical and [nonpolar](#) molecules, for which neither of the above could be expected. Such interactions became apparent from careful studies of the behavior of real gases. The molecules of an ideal gas occupy no volume and have no interactions, so  $n$  moles of an ideal gas obey the equation of state known as the Boyle–Gay-Lussac’ law:  $PV = nRT$ , where  $P$  equals the pressure,  $V$  the volume,  $R$  the gas constant, and  $T$  the absolute temperature. The behavior of real gases, however, is better described by the equation of state given by J. D. van der Waals:

$$(P + n^2 a/V^2)(V - nb) = nRT \quad (1)$$

where the pressure and volume terms are modified. The parameters  $a$  and  $b$  account, respectively, for the attractive and repulsive intermolecular interactions of the molecules (1). The values of  $a$  and  $b$  determined experimentally for selected gases are given in Table 1. It should be noted that  $a$  and  $b$  have finite values for both polar and nonpolar atoms and molecules.

**Table 1. Values of the Parameters  $a$  and  $b$  of the van der Waals Equation for Various Substances**

---

| Substance                     | $a$   | $b$    |
|-------------------------------|-------|--------|
| H <sub>2</sub>                | 0.248 | 0.0266 |
| He                            | 0.034 | 0.024  |
| H <sub>2</sub> O              | 5.52  | 0.0304 |
| CO                            | 1.47  | 0.0304 |
| HCl                           | 3.72  | 0.0408 |
| CO <sub>2</sub>               | 3.66  | 0.0428 |
| Ar                            | 1.36  | 0.0322 |
| C <sub>6</sub> H <sub>6</sub> | 18.9  | 0.120  |
| Xe                            | 4.17  | 0.0513 |
| O <sub>2</sub>                | 1.38  | 0.0317 |
| N <sub>2</sub>                | 1.37  | 0.0387 |

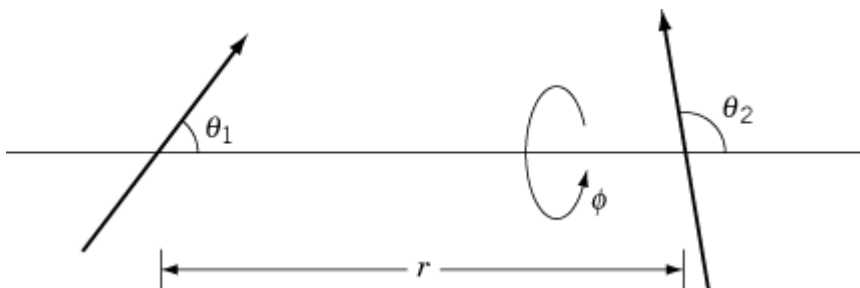
<sup>a</sup> The unit of  $a$  is (bar<sup>2</sup>molar<sup>-2</sup>) and of  $b$  (molar<sup>-1</sup>). Taken from [1](#) with permission.

The attractive interaction between neutral atoms and molecules that is experimentally observed and quantified in the van der Waals equation of state is usually called the *van der Waals interaction* or *van der Waals force*. It is easy to understand that the value of  $a$  is large for polar molecules such as CO<sub>2</sub>, H<sub>2</sub>O, or HCl, where the gas molecules are characterized by finite permanent dipole moments. Two dipole moments  $m_A$  and  $m_B$  with a fixed orientation attract or repel each other when their orientation is antiparallel or parallel, respectively. In both cases, the interaction potential has an inverse 3rd power dependence on the distance between the two dipoles (Fig. [1](#)). When this interaction potential is averaged, with respect to the Boltzmann distribution, over all possible orientations of the two dipoles, and if we assume that the potential is small compared to the thermal energy of  $kT$ , the average potential energy, as obtained by Keesom is [\(2\)](#)

$$U = -2\mu_A^2\mu_B^2/3(4\pi\epsilon_0)^2kTr^6 \quad (2)$$

where  $\epsilon_0$ ,  $k$ , and  $T$  are the dielectric constant of vacuum, the Boltzmann constant, and temperature, respectively.

**Figure 1.** Interaction between two molecules with permanent dipoles.



When one of the two does not have a permanent dipole, but has polarizability  $a$ , it will have an induced dipole in the neighborhood of another molecule with permanent dipole  $m$ . The interaction energy in this case is also dependent on the inverse 3rd power of the distance between the two molecules, and the thermally averaged potential between two molecules A and B was given by Debye (2) as

$$U = -(\alpha_B \mu_A^2 + \alpha_A \mu_B^2) / (4\pi\epsilon_0)^2 r^6 \quad (3)$$

It was difficult initially to explain the nonzero value of  $a$  for nonpolar atoms and molecules, such as Ar, Xe, N<sub>2</sub>, or CH<sub>4</sub>, until 1930, when the newly emerging quantum mechanics were successfully used to explain the interaction between two nonpolar particles. The results may be interpreted in classical terms as follows. The average centers of gravity for the positive and negative charges in a nonpolar atom coincide precisely, so it lacks a permanent dipole moment, but the instantaneous positions of the two centers do not match exactly. Consequently, nonpolar atom and molecules have an instantaneous dipole moment, whose magnitude and direction are constantly changing. When two such particles come into close proximity, the fluctuating dipoles influence each other, and their electron distributions fluctuate so as to produce an attractive interaction between them.

Such an attractive interaction between nonpolar atoms and molecules in the gaseous state is known as the dispersion force of London (3). A similar attractive interaction between such molecules occurs in the liquid state, where molecules are also in constant agitation and rotation, although the influence of the dense medium cannot be neglected. For example, the dispersion force between two nonpolar molecules dissolved in [water](#) is expected to be reduced by at least an order of magnitude compared to their interaction in a vacuum (4). The London dispersion interaction potential as described above has the following dependence on various electronic parameters of the interacting particles:

$$U_L = -\frac{3h\nu_0\alpha^2}{4(4\pi\epsilon_0)^2 r^6} = \frac{3I_A I_B}{2(I_A + I_B)} \frac{\alpha_A \alpha_B}{(4\pi\epsilon_0)^2 r^6} \quad (4)$$

where  $h$  is Planck's constant;  $\nu_0$  is the transition frequency,  $I_A$  and  $I_B$  are the ionization potentials and  $\alpha_A$  and  $\alpha_B$  the polarizabilities of molecules  $A$  and  $B$ . In some textbooks, the term van der Waals interaction is reserved for the London dispersion interactions alone. In conclusion, thermally averaged potentials for all types of interactions, ie, dipole-dipole (Keesom), dipole-induced dipole (Debye), and dispersion interactions between two fluctuating dipoles (London), have inverse 6th power dependence on the distance between them.

The origin of repulsive interactions between closed-shell atoms and molecules is due to Pauli's exclusion principle, which states that two electrons are allowed to occupy the same orbital only if their spins occur in opposite directions. Since each of the electronic orbitals in closed-shell atoms and molecules is already occupied by such paired electrons, they repel other electronic systems that come close enough for their outermost electron clouds to overlap each other.

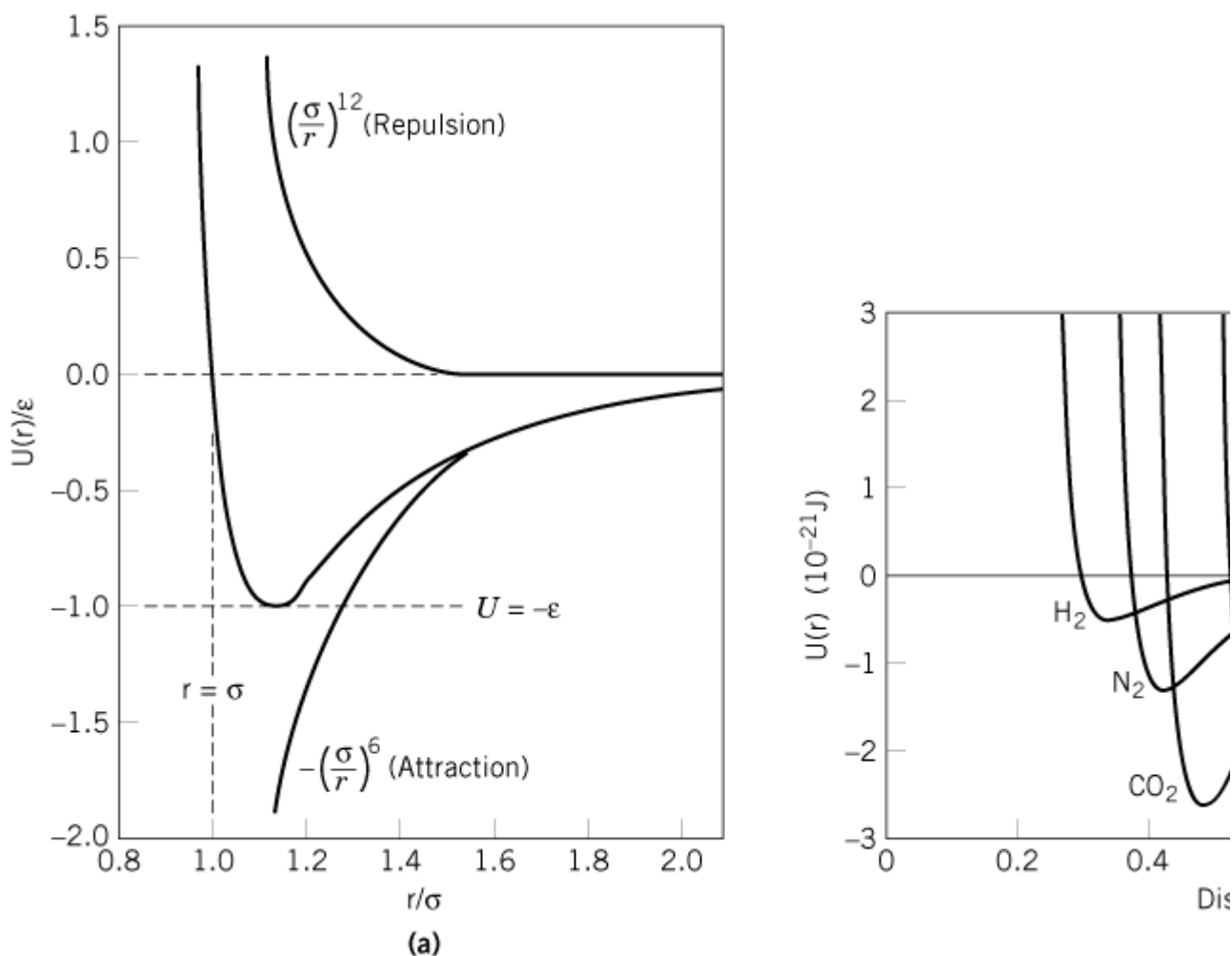
Lennard-Jones presented a concise form of the overall potential between two neutral molecules by expressing it with an equation having two adjustable parameters,  $\epsilon$  and  $s$ :

$$\phi(r) = 4\epsilon \left[ \left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right] \quad (5)$$

The parameters  $\epsilon$  and  $s$  are, respectively, the dissociation energy and inter particle distance where the potential becomes zero, as depicted in Figure 2. The inverse 6th power dependence of attractive interaction becomes inverse 7th power at a large distance, approximately 100 nm, due to the so-

called “retardation effect” resulting from the two molecules exchanging electro-magnetic interactions at a finite speed, that of the speed of light. When the distance is large, this causes a more rapid decrease in the potential with distance. Further refinements of the potential taking into account the interactions involving quadrupole moments will give  $-8$ th and  $-10$ th power dependences on the distance.

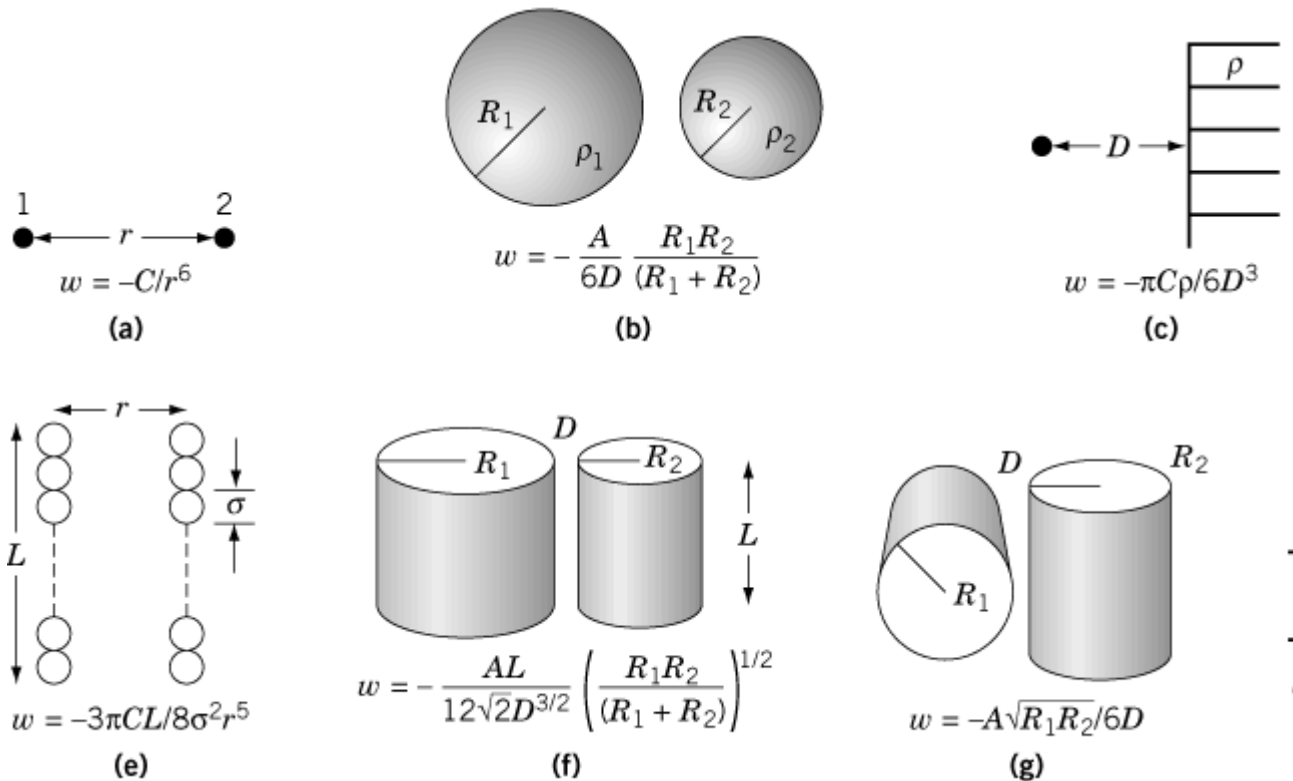
**Figure 2.** The Lennard-Jones potential in terms of two parameters  $\epsilon$  and  $\sigma$  as a function of intermolecular distance  $r$ . The the two terms  $\epsilon$  and  $\sigma$  and breaks down the potential into the attractive and repulsive forces. The graph on the right gives gases. Taken from (1) with permission.



Although its origin from Pauli's exclusion principle is clear, the inverse 12th power dependence of the repulsive potential is not of exact physical derivation. It is sometimes represented by the exponential term  $r \exp(-D_{12})$ , which has a similar steepness against  $r$ . The potential function, when differentiated with respect to  $r$ , gives the force, and therefore the repulsive force has inverse 13th power dependence on distance, and the attractive force has the inverse 7th power. Where the potential is at a minimum, the force is zero.

The inverse 6th power dependence of van der Waals interaction potential on the distance between two point atoms or molecules decreases when applied to spheres of finite sizes and surfaces of infinite extension (Fig. 3).

**Figure 3.** Dependence of van der Waals potential between two objects with finite sizes or infinite extensions.  $D$  is the distance between the objects, where  $r_1$  and  $r_2$  are the number of atoms per unit volume in the two bodies, and  $C$  is the coefficient in the atom–atom pair potential. Other symbols are defined in the figure. (a) Two atoms, (b) two spheres, (c) atom–surface, (d) sphere–surface, (e) two parallel cylinders, (f) two crossed cylinders, and (g) two surfaces. Taken from [5] with permission.



## 1. van der Waals Interactions in Biochemistry and Molecular Biology

An example of the usefulness of the van der Waals potential in biochemistry and molecular biology may be found in computer-aided calculations to explain the behavior of biopolymers, including [proteins](#) and nucleic acids (see [Molecular Dynamics](#)). The majority of biologically important molecules are either ionic or polar molecules, and thus their interaction potentials are dominated by [electrostatic interactions](#), but the van der Waals potential is involved in all kinds of interactions. The importance of van der Waals potential becomes more obvious in molecular model calculations and model building when the participating molecules are not explicitly charged.

## 2. van der Waals Radii

The meaning of the van der Waals radius is closely related to our intuitive notion of the size of atoms and molecules. In the van der Waals equation of state for real gases, the volume of gaseous particles was parameterized in the excluded volume  $b$ , which is expected to be equal to  $4 \times (4/3) \pi r^3$ , where  $r$  is the radius of the gas particle when it is approximated by a hard sphere. Experimentally, the value of  $b$  can be obtained from the measurement of the second virial coefficient. It can be determined more precisely from the radii of atoms and molecules in their closely packed crystals. If one defines such a radius for a single kind of atom in different kinds of molecules, a remarkably constant value can be obtained. The radius is then defined from the balance between the repulsive force due to Pauli's exclusion principle and the van der Waals force. It may thus be related to the parameter  $s$  in the Lennard-Jones potential, for when the distance between two particles is equal to this parameter, the potential is zero. When  $r = 2^{1/6} s$ , the potential energy is at its minimum, and the distance

corresponds to the sum of the two van der Waals radii. The meaning is that when two particles come closer than their van der Waals radii, they repel each other. Table 2 summarizes the currently accepted values of the van der Waals radii for representative atoms and molecules, together with their covalent radii. The van der Waals radius is not a constant parameter for each element but depends on the neighboring atoms with which it is in contact.

**Table 2. Table of Values of the van der Waals Radius ( $r_{\text{vdw}}$ ) and Covalent Radius ( $r_{\text{cov}}$ ) of Common Elements in Biochemistry**

| Elements                      | $r_{\text{vdw}}$ (nm)                          | $r_{\text{cov}}$ (nm) |
|-------------------------------|--|-----------------------|
| H                             | 0.120–0.145                                    | 0.037                 |
| He                            | 0.180 (0.150)                                  | 0.030                 |
| Li                            | 0.180  | 0.134                 |
| C                             | 0.165–0.170                                    | 0.077                 |
| N                             | 0.155  | 0.075                 |
| O                             | 0.150  | 0.073                 |
| F                             | 0.150–0.160 (0.135)                            | 0.071                 |
| Na                            | 0.230  | 0.154                 |
| Mg                            | 0.170  | 0.145                 |
| P                             | 0.185 (0.190)                                  | 0.110                 |
| S                             | 0.180  | 0.120                 |
| Cl                            | 0.170–0.190                                    | 0.099                 |
| K                             | 0.280  | 0.196                 |
| I                             | 0.195–0.212 (0.215)                            |                       |
| <u>Organic groups</u>         |  |                       |
| CH <sub>3</sub>               | 0.200  |                       |
| C <sub>6</sub> H <sub>6</sub> | 0.170 (perpendicular to the ring) <sup>a</sup> |                       |

<sup>a</sup> Values are from (5) (original data are based on (6)). The values in parentheses are from (1), originally based on (7).

#### Bibliography

1. G. M. Barrow (1988) *Physical Chemistry*, McGraw-Hill, New York, Chap. "2", pp. 30–53.
2. W. J. Moore (1972) *Physical Chemistry*, 4th ed., Prentice-Hall, Englewood Cliffs, NJ, Chap. "19", pp. 915–919, and Chap. 4 and 5 of Ref. 5.
3. W. Kauzmann (1957) *Quantum Chemistry: An Introduction*, Academic Press, New York, Chap. "13", pp. 503–517.
4. T. Kihara (1978) *Intermolecular Forces*, John Wiley & Sons, New York (translated from the



Japanese version by S. Ichimaru).

5. J. Israelachvili (1991) *Intermolecular & Surface Forces*, 2nd ed., Academic Press, London, chap. "11", pp. 176–212.
6. J. E. Huheey (1983) *Inorganic Chemistry: Principles of Structure and Reactivity*, 3rd ed., Harper and Row, New York, Chap. "6".
7. A. Bondi (1964) van der Waals volumes and radii, *J. Phys. Chem.* **68**, 441–451.

## van der Waals Surface, Volume

Every atom of a molecule can be regarded as a sphere of the appropriate van der Waals radius, which defines the normal limit of closest approach of other atoms in its environment. The spheres of atoms that are covalently bonded are normally truncated. The surface that results from all these spheres corresponding to all the atoms in a macromolecule is called the van der Waals surface.

The van der Waals volume of a molecule is defined as the space enclosed by the van der Waals spheres of the constituent atoms. The van der Waals volume and surface area of entire small molecules that are not structurally strained may be estimated by simply adding the values for the constituent atoms or chemical groups. In the case of folded, globular proteins, however, the van der Waals surface and volume are not very indicative of the actual surface and volume of the molecule (see [Protein Structure](#)). The interior of a protein molecule is densely packed, with adjacent atoms frequently in van der Waals contact. The average volumes occupied by residues in folded proteins are virtually the same as those they occupy in crystals of the amino acids (Table [1](#)). These volumes are, of course, larger than the van der Waals volume of each amino acid residue.

The van der Waals surface of a folded protein is very complex; much of the surface is in the interior, surrounded by other atoms of the protein and inaccessible to the solvent. Although it has a strictly defined surface area and encloses a defined volume, no chemical procedure can measure this area and volume directly. For that reason, the surface that is generally considered most relevant is the [accessible surface](#), which is defined as that surface that is normally in contact with the solvent.

The packing density of the interior of a macromolecule is defined as the ratio of the van der Waals volume to the volume actually occupied in a solvent. Thus, about 75% on average of the interior volume of a protein is filled with atoms. This is close to the maximum value of 74% possible with close packing of identical spheres and within the range of 70 to 78% found for crystals of small organic molecules (Table [1](#)). This means that the interior of a protein is not an oil drop but resembles instead a molecular crystal. The packing is not necessarily uniform throughout the interior, however, because the atoms are covalently bonded, which exaggerates the close packing. Generally, the atoms in topologically regular parts, such as [alpha-helix](#) and [beta-sheet](#), pack closely, whereas the [active sites](#) of [enzymes](#) are usually packed more loosely. There are small cavities in the packing defects, some of them being large enough to accommodate other molecules. Cavities form only a small percentage of the total volume, but they play an important role in structural flexibility (thermal fluctuation), which is believed to be necessary for most protein functions.

**Table 1. Volumes of Amino Acid Residues**

---

| Residue | Van der Waals volume <sup>a</sup> (Å <sup>3</sup> ) | Residue crystal volume <sup>a</sup> (Å <sup>3</sup> ) | Average volume of buried residues <sup>a</sup> (Å <sup>3</sup> ) |
|---------|---|---|--|
| Ala     | 67  | 96.6  | 91.5   |
| Arg     | 148   |   | 202  |
| Asn     | 96  |   | 135.2  |
| Asp     | 91  | 122.0   | 124.5  |
| Cys     | 86  | 108.7   | 105.6  |
| Gln     | 114   | 148.0   | 161.1  |
| Glu     | 109   | 143.9   | 155.1  |
| Gly     | 48  | 66.5  | 66.4   |
| His     | 118   | 166.3   | 167.3  |
| Ile     | 124   | 169.7   | 168.8  |
| Leu     | 124   |   | 167.9  |
| Lys     | 135   |   | 171.3  |
| Met     | 124   | 176.1   | 170.8  |
| Phe     | 135   |   | 203.4  |
| Pro     | 90  | 124.4   | 129.3  |
| Ser     | 90  | 102.2   | 99.1   |
| Thr     | 93  | 124.3   | 122.1  |
| Trp     | 163   |   | 237.6  |
| Tyr     | 141   | 201.7   | 203.6  |
| Val     | 105   | 143.4   | 141.7  |

<sup>a</sup> Volume enclosed by van der Waals radii of atoms.

<sup>b</sup> Estimated by assuming 11.1Å<sup>3</sup> for the volume lost by an amino acid on becoming a residue.

<sup>c</sup> A residue is defined as buried if 5% or less of its potential accessible surface area is available to solvent contact (1).

## Bibliography

1. C. Chothia (1975) Nature **254**, 304–308.

## Suggestion for Further Reading

2. F. M. Richards (1977) Areas, volumes, packing, and protein structure. Ann. Rev. Biophys. Bioeng. **6**, 151–176.

## Vapor Phase Crystallization

[Crystallization](#) of a **protein** for determination of its three-dimensional structures by [X-ray crystallography](#) requires the gradual development of a state of supersaturation of the protein in the medium that contains the crystallizing agent. One way of approaching that state is by vapor phase

equilibrium. This technique is based on the principle that when two solutions of different concentrations of the same nonvolatile solute are placed in a single sealed chamber, [water](#) molecules will be transported in the vapor phase from the solution of lower concentration (higher vapor pressure) to that of higher concentration until the vapor pressures (solute concentrations) of the two become identical.

This principle is used successfully in the crystallization of proteins. The crystallizing agent used most frequently is *polyethylene glycol* (PEG) of a molecular weight between 1,000 and 20,000. In this procedure, a droplet of the protein solution (10 to 20  $\mu\text{l}$ ) is deposited on a well-siliconized glass plate. This plate is inverted over a well containing 0.5 to 1.0 ml of the solvent at a considerably higher concentration (about twice) of the PEG, and the system is left to equilibrate at constant temperature. When the concentration of the crystallizing agent on the platelet has increased to that in the well by vapor phase transfer of water molecules, the droplet is examined under a microscope for crystal formation. A great advantage of this technique is the ability to examine simultaneously a large number of crystallizing agent concentrations. For example, a 24-well tissue culture tray may be used for the wells, each of which contains a different concentration or size of PEG and above each of which is suspended a droplet from a coverslip.

#### Suggestions for Further Reading

A. McPherson (1985) *Meth. Enzymol.* **114**, 120–125.

A. McPherson (1982) *The Preparation and Analysis of Protein Crystals*, John Wiley & Sons, New York.

### Variable Domain, Region

The variable (V) regions of [immunoglobulins](#) (Igs) are organized as two complementary **domains**, one VH and one VL (see [Immunoglobulin Structure](#)), each encoded by [gene rearrangement](#) of VH–D–JH and VL–JL genes, respectively (see [Immunoglobulin Biosynthesis](#)). VH and VL are associated noncovalently in the so-called Fv fragment that can be released by **proteolytic** digestion of the native immunoglobulin. Because of the symmetrical organization of the molecule, one antibody contains two identical Fv fragments, and each possesses one **antigen-binding** site. Physicochemical measurements have shown that the **dissociation constant** of the isolated Fv is very close, if not identical, to that of the corresponding [antibody](#), indicating that the three-dimensional [protein structure](#) is well-conserved in the isolated fragment. This is in agreement with the domain hypothesis and is a consequence of the rigid structure of the globular organization of Ig domains, which defines the Ig [superfamily](#).

The V domain architecture consists of nine antiparallel **b-strands**, labeled a, b, c, c', c'', d, e, f, and g, that are organized in two **b-pleated sheets** anchored to each other by an intrachain [disulfide bond](#). The structure of a V domain closely resembles that of a C domain, with the notable exception of the additional c' and c'' b strands, which contribute an interacting H–L region in the antibody combining site. [X-ray crystallography](#) structures have demonstrated the complete correspondence between the sequence organization of V regions and the architecture of the combining site. Amino acid sequences of VH and VL regions allowed Kabat to identify within V regions three hypervariable segments, around positions 30, 60, and 90, separated by framework regions that ensure the structure of the domain. Hypervariable regions are designated CDR1, CDR2, and CDR3 (for complementarity determining regions) (see [Immunoglobulin Structure](#)). CDR1 and CDR2 are encoded entirely by V genes, whereas CDR3 have an increased diversity, especially in the heavy chains, because they are

located in the area encoded by the V–D–J joint (and V–J for the L chains). X-ray crystallography structures have clearly demonstrated that all six CDRs from both interacting chains contribute directly the residues involved in binding an antigen.

Fv fragments have retained the complete antibody binding sites, so they may be considered as “mini-antibodies” that might be convenient for therapeutic or laboratory use. This has stimulated their production by several approaches of [protein engineering](#). Phage [combinatorial libraries](#) may be prepared with vectors containing [complementary DNA](#) encoding VH and VL regions; their combinatorial random association will result potentially in an enormous [repertoire](#) of Fv fragments. Alternatively, Fv with preformed specificity may also be constructed and expressed in a filamentous phage as a single chain, called scFv, by introduction of an artificial linker between the VH and the VL cDNAs. Such vectors may be amplified and produced by the classical technologies in phage, bacteria, or Baculovirus. One interesting potential usage is medical imaging to tag a given organ or tumor with a specific scFv bearing a radioactive reporter that may be detected by its radioactivity. Also, “humanized” scFv may be engineered by inserting the CDRs derived from hyperimmunized mice into the framework regions of human origin. Finally, random mutagenesis directed at the CDR sites may amplify the diversity of these new tools almost indefinitely, and all these approaches are being widely explored.

See also entries [Antibody](#), [Idiotypes](#), [Gene Rearrangement](#), and [V Genes](#).

#### Suggestions for Further Reading

- C. Clothia, J. Novotny, R. Brucoleri, and M. Karplus (1985) Domain association in immunoglobulin molecules. The packing of variable domains. *J. Mol. Biol.* **186**, 651–663.
- C. Horne, M. Klein, I. Polidoulis, and K. J. Dorrington (1982) Noncovalent association of heavy and light chains of human immunoglobulins. III. Specific interactions between VH and VL. *J. Immunol.* **129**, 660–664.
- A. F. Williams and A. N. Barclay (1988) The immunoglobulin superfamily domains for surface recognition. *Annu. Rev. Immunol.* **6**, 381–405.
- G. Winter, A. D. Griffiths, R. E. Hawkins, and H. R. Hoogenboom (1994) *Annu. Rev. Immunol.* **12**, 433–455.

## Vertical Gene Transfer

Vertical gene transfer is the mechanism by which a **gene** is transmitted from parents to offspring; this is the normal way to inherit a gene. This term was coined to make a clear distinction in [evolution](#) with [horizontal gene transfer](#), which is a phenomenon where genes are transferred between different species, probably by the mediation of **viruses** or bacteria.

## Vesicular Stomatitis Virus (VSV)

Vesicular stomatitis virus (VSV) is the prototype rhabdovirus and has been used for viral and

biochemical studies because it can be highly purified from infected cells in high yields. This animal pathogen is also frequently used in the laboratory as an efficient inducer of [interferon](#). Like other rhabdovirus, such as rabies virus, the virion has a typical rod-shaped bullet-like morphology, and the VSV genome is a linear single-stranded RNA (11 kb) of negative polarity that serves as the template for both [transcription](#) and replication.

For all negative-strand RNA viruses, including rhabdoviruses, infection begins with the synthesis in the infected cells of the viral [messenger RNA](#). VSV virions contain all the [enzymes](#) (the P and L proteins, see text below) required for the synthesis and processing of five mRNAs, and these virion-associated enzymes are inserted into the cell upon infection, together with the RNA genome; each transcript has a [cap](#) and is **polyadenylated** at its 3' end. Like all other rhabdoviruses, the genomic organization of VSV starts with a small leader sequence at the extreme 3' terminus of the linear single strand RNA and is followed by the template for mRNAs encoding the proteins N (nucleocapsid; 49 kDa), P (phosphoprotein; 29 kDa), M (matrix; 25 kDa), G (glycoprotein; 69 kDa), and then L (large; 241 kDa). Several lines of evidence, such as ultraviolet transcriptional mapping, indicated that the five genes are transcribed sequentially in the same order as they appear in the genome. All of these mRNA synthesis and post-transcriptional modifications can be studied *in vitro* using **detergent**-disrupted purified virions or purified ribonucleoprotein particles (RNPs). Even after many years of study both *in vivo* and *in vitro*, it is not yet fully understood how these five mRNAs can be generated discontinuously from a linear template, in spite of the fact that transcription is apparently sequential. It could be explained by **nuclease** processing of a single precursor RNA molecule or by sequential reinitiation by transcriptase at the 5' end of each mRNA. At the intergenic junctions, polyadenylation and capping should be tightly coupled to VSV RNA synthesis.

The primary transcription causes the accumulation of viral proteins in the infected cells and leads to the onset of genome replication. The **RNA polymerase** switches mode to produce the positive-strand genome-length RNA, which also involves the newly synthesized N-protein, to form the positive-strand RNP. These positive-strand RNPs function as the template for the production of the negative-strand RNPs, which accumulate in the cell and contribute to the synthesis of positive-strand RNPs (secondary transcription). This in turn induces high-level expression of viral proteins in the later stages of infection.

The accumulated negative-strand RNP also serves as the genome of progeny virion, which are assembled at the plasma membrane of the infected cell. When the virions are released from the cells by budding, they have acquired lipid envelopes. Spike-like projections, composed of trimers of G protein, emanate from the envelop, while the M protein lines the inner surface of the envelope or associates with the RNP. So all five viral proteins are found associated with the purified virion. Upon infection, VSV-G protein binds to a cellular receptor, which is known to include such widely present anionic phospholipids as phosphatidylserine, while some protein components might also be necessary for full efficiency. These properties of the G protein partly explain why VSV has a very broad host range; this rhabdovirus not only replicates in several vertebrates, but also can infect insects. The G protein can be also used for the envelope proteins of other distantly related viruses, such as [retrovirus](#) and lentivirus, a phenomenon known as phenotypic mixing. For these reasons, several retroviruses having VSV-G instead of their own envelope protein have been developed for use as stable and high-titer retrovirus vectors with a broad host range.

VSV-G has been also well-studied from the standpoint of **protein targeting** in the cell. When expressed in polarized epithelial cells, VSV-G locates exclusively at the basolateral plasma membrane domain, rather than the apical surface of the same cell. In such a polarized cell, morphogenesis and budding of VSV-G particles occurs specifically in the basolateral plasma membrane.

#### Suggestion for Further Reading

A. K. Banerjee (1987) Transcription and replication of rhabdoviruses. *Microbiol. Rev.* **51**, 66–87.

## Vibrational Spectroscopy

### 1. Infrared, Raman, and Resonance Raman Spectroscopies

The vibrations of molecules provide an avenue to their identification and to the characterization of their structure and environment (see [Spectroscopy](#)). The stretching and bending of bonds occur with characteristic frequencies (energies), depending on force constants and the masses of the connected atoms (1). These frequencies occur in the *infrared* (IR) region of the electromagnetic spectrum, and the vibrations can be detected via absorption of infrared light. The infrared spectrum contains a series of peaks, or bands, whose frequencies are characteristic of the molecule under study. The spectrum is a fingerprint and can be used for identification and quantification. In addition, variations in the band positions and intensities provide information about changes in bonding and conformation that may result from intermolecular interactions.

The same vibrations can be accessed with Raman spectroscopy, a laser [light scattering](#) technique. Most of the laser photons scattering from a sample emerge with unshifted energy (Rayleigh scattering), but a few photons excite vibrations in the molecules and emerge with diminished energy (*Raman scattering*). The shifts in frequency (from the loss of energy) are the same as the infrared frequencies, and the Raman spectrum contains vibrational bands, as does the infrared spectrum. However, the intensities differ in the two spectroscopies: infrared intensity depends on the change in *dipole moment* produced by a vibration, whereas Raman intensity depends on the change in the *polarizability*. Thus, the two techniques provide complementary information.

Raman spectroscopy is advantageous for studies in aqueous solution, because water is a poor Raman scatterer but a strong IR absorber. Also, Raman spectroscopy requires less sample, because the scattering volume of the laser is only microliters, whereas IR spectroscopy requires ~1cm diameter disk of material. However, the use of condensing optics in an IR microscope can overcome this limitation (2, 3). The main obstacle to Raman spectroscopy is interference from **fluorescence**. Because a very small fraction of the laser photons appear at the Raman frequencies, not a great deal of laser-induced fluorescence (whether intrinsic to the molecule under study or from impurities) is necessary to swamp the Raman spectrum. One way to overcome this difficulty is to shift the laser wavelength either into the ultraviolet region (4), as most fluorescence occurs in the visible region, or into the far red region, where excitation of fluorescent species is avoided (5).

Both IR and Raman spectroscopy suffer from low sensitivity and selectivity. Molecular concentrations must be quite high (>0.1M), and spectral crowding is severe for complex molecules because there are  $3N-6$  modes of vibration ( $N$  being the number of atoms). This problem can be overcome by difference spectroscopy, in which subtle changes can be detected by careful subtraction of most of the vibrational spectrum. The advent of Fourier transform (FT) techniques (FTIR and FT-Raman) has made highly accurate subtractions possible, and sensitivity at the single residue level in proteins has been demonstrated (6).

Resonance Raman (RR) spectroscopy provides a powerful way to increase sensitivity and selectivity (7-9). When the laser frequency approaches that of an electronic transition of the molecule, a large increase in scattering power is observed, but this occurs only for those vibrations that carry the molecule into its excited state geometry. Thus, RR scattering is selective for the chromophoric region of the molecule and can be detected at quite low concentrations. It provides a structure probe for [active sites](#) of biomolecules, if their electronic transitions can be accessed by available lasers. Information about nonresonant parts of the molecule, however, are lost. The technique has been

particularly useful in the study of colored cofactors. The advent of reliable ultraviolet lasers has also made it possible to probe nucleic acid bases (10-13), protein aromatic side-chains (14-17), and even the peptide bond (18-20), using RR spectroscopy.

## 2. Proteins

### 2.1. Backbone Vibrations; Amide Bond

The protein backbone has several characteristic vibrations, labeled amide I, II, III, IV, V, VI, VII, A, and B and found at about 1650, 1540, 1300, 650, 725, 300, and 3280, 3090  $\text{cm}^{-1}$ , respectively (21-23). Because all [peptide bonds](#) (OC-NH) are chemically the same, their vibrations overlap in the spectrum. However, the elements of their **secondary structure** (eg,  **$\alpha$ -helix**,  **$\beta$ -sheet**, **turns**) have sufficiently different geometries and interpeptide interactions that their frequencies can be resolved in favorable cases, especially for the amide I and III bands (Table 1). These can be used to evaluate the secondary structure from IR or Raman spectra (21-26).

**Table 1. The Infrared and Raman Frequencies (in  $\text{cm}^{-1}$ ) of the Amide I and III Bands for Various Peptide Secondary Structures**

| Conformation    | Amide I                       |                            | Amide III                  |
|-----------------|-------------------------------|----------------------------|----------------------------|
|                 | Infrared ( $\text{cm}^{-1}$ ) | Raman ( $\text{cm}^{-1}$ ) | Raman ( $\text{cm}^{-1}$ ) |
| $\alpha$ -helix | 1650                          | 1645–1660                  | 1265–1300                  |
| $\beta$ -sheet  | 1632                          | 1665–1680                  | 1230–1240                  |
| $\beta$ -turn   | —                             | 1660–1680                  | 1260–1300                  |
| Unordered       | 1658                          | 1660–1670                  | 1240–1260                  |

The amide II and III bands are strongly enhanced in RR spectra when the laser wavelength approaches the strong amide p-p\* transition at  $\sim 190\text{nm}$  (18-20). An additional band, at  $\sim 1390\text{cm}^{-1}$ , is also enhanced, the intensity of which is sensitive to secondary structure. This mode involves bending of the  $\text{C}_\alpha\text{-H}$  bond and is suppressed in  $\alpha$ -helices because of kinematic effects (20).

The X-Pro peptide bond preceding a [proline](#) residue differs from other peptide bonds because it lacks the N-H bond, an important contributor to the amide II and III modes (21-23, 27, 28). Its amide II frequency is much lower than others ( $\sim 1470\text{cm}^{-1}$ ). In addition, its p-p\* transition is red-shifted, giving it stronger resonance enhancement at wavelengths greater than 200 nm (28). Consequently, the X-Pro bond can be detected in ultraviolet resonance Raman (UVRR) spectra, even for a single proline residue. Its frequency is insensitive to **cis/trans isomerization**, but it does reflect the [hydrogen bond](#) status of the X-Pro carbonyl (29, 30).

### 2.2. Amino Acid Side-Chains

The diversity of amino acids makes it difficult to disentangle side-chain vibrations. However, the aromatic residues—[phenylalanine](#) (Phe), [tyrosine](#) (Tyr), and [tryptophan](#) (Trp), including [histidine](#) (His)—have characteristic ring vibrations that sometimes can be distinguished in IR and Raman spectra. In addition, the p-p\* transitions provide enhancement in UVRR spectra (17), although the transition for histidine is weak (16, 31-33). Except for Phe, the aromatic side-chains have hydrogen bond sites, and the vibrational frequencies and intensity patterns are sensitive to environment. Some

mode frequencies are listed in Table 2.

**Table 2. Selected Vibrational Frequencies of Aromatic Amino Acids<sup>a</sup>**

| Observed Frequency (cm <sup>-1</sup> ) | Assignment   |
|--|--|
| <i>Phenylalanine</i>                   |  |
| 1606                                   | F8a: ring C-C stretch  |
| 1586                                   | F8b: asymmetric ring C-C stretch   |
| 1207                                   | F7a: C-C <sub>ext</sub> stretching   |
| 1182                                   | F9a: ring stretching +C-H bending  |
| 1028                                   | F18a: C-H bending  |
| 1000                                   | F1: ring breathing   |
| <i>Tyrosine</i>                        |  |
| 1617                                   | Y8a: ring C-C stretching   |
| 1601                                   | Y8b: ring C-C stretching   |
| 1519                                   | Y19a: similar to benzene n <sub>19a</sub>  |
| 1443                                   | Y19b: similar to benzene n <sub>19b</sub>  |
| 1263                                   | Y7a': C-O stretching   |
| 1210                                   | Y7a: C <sub>ring</sub> -C <sub>ext</sub> stretching  |
| 1180                                   | Y9a: C-H bending   |
| 853 / 831                              | Fermi resonance pair: Y1 and 2×Y16a interaction  |
| <i>Tryptophan</i>                      |  |
| 1622                                   | W1: benzene-localized mode (n <sub>8a</sub> )  |
| 1578                                   | W2: benzene-localized mode (n <sub>8b</sub> )  |
| 1555                                   | W3: C <sub>2</sub> -C <sub>3</sub> pyrrole-localized mode                                    |
| 1496                                   | W4: similar to benzene n <sub>19b</sub>  |
| 1462                                   | W5: similar to benzene n <sub>19a</sub>  |
| 1434                                   | W6: N <sub>1</sub> -C <sub>2</sub> -C <sub>3</sub> stretch coupled to N <sub>1</sub> -H bend |
| 1361 / 1342                            | W7: Fermi resonance pair   |
| 1305                                   | W8: C <sub>3</sub> -C <sub>9</sub> stretch +N <sub>1</sub> -H bend                           |
| 1238                                   | W10: C-H +C <sub>3</sub> -C <sub>ext</sub>   |
| 1127                                   | W13: similar to benzene n <sub>9b</sub> C-H bend   |
| 1016                                   | W16: benzene C-C stretch   |
| 880                                    | W17: similar to benzene n <sub>12</sub> +N <sub>1</sub> -H                                   |
| 762                                    | W18: indole ring breathing   |



<sup>a</sup> Adapted from R. P. Rava and T. G. Spiro (16).

### 2.2.1. Tyrosine

Several ring modes show sensitivity to hydrogen bonding. These include  $n_{7a}$ ,  $n_{7a'}$ ,  $n_{8a}$ ,  $n_{8b}$ , and  $n_{9a}$  (34, 35) and also the 830 / 850  $\text{cm}^{-1}$  doublet band, which arises from a Fermi resonance between the ring-breathing vibration,  $n_1$ , and an overtone of an out-of-plane ring bending vibration,  $n_{16a}$  (36).

This doublet has been much studied because it is well separated from other bands and is easily distinguishable, even in off-resonance Raman spectra. The intensity ratio of the two components ( $I_{850}:I_{830}$ ) is an environmental indicator, being low for a buried Tyr side-chain and high for a solvent-exposed one (26, 37). In UVRR spectra, the strongly enhanced  $n_{8b}$  is particularly useful because its frequency has been calibrated as a measure of hydrogen bond strength (38, 39).

### 2.2.2. Tryptophan

Like Tyr, Trp has an environmentally sensitive doublet band, at 1360 / 1340  $\text{cm}^{-1}$ ; the intensity ratio is high for buried Trp side-chains and low for solvent-exposed ones (33, 40). The frequencies of the W17, W6, W4, and W2 modes are also environmentally sensitive and have been found in model compounds to correlate well with hydrogen bonding (41-43). W17 should be particularly useful because it is enhanced in UVRR spectra, but some doubt exists about the applicability of the model data to proteins. W17's frequency has been found not to change when [cytochrome c](#) is denatured, even though the single Trp residue is strongly hydrogen-bonded to a heme propionate substituent in the native state (44). W17's frequency is sensitive to deuterium/[hydrogen exchange](#) because it is partially composed of N-H bending. W17 can therefore be used as a monitor of exchange rate (45).

Another useful marker is W3, whose frequency depends on the dihedral angle of the bond connecting the indole ring to the  $C_b$  atom (41). This dependence was very useful in distinguishing the W3 modes of the three Trp residues (per dimer unit) in the UVRR spectrum of [hemoglobin](#) and thereby in discriminating between **tertiary** and **quaternary** interactions (35).

### 2.2.3. Histidine

Like Trp, His has a band,  $n(C_4 = C_5)$ , whose frequency (1575–1600 $\text{cm}^{-1}$ ) is sensitive to the dihedral angle of the bond connecting the ring to the  $C_b$  atom (46). In addition, a marker of its [ionization](#) state (1408 $\text{cm}^{-1}$ ) has been reported in  $D_2O$  (47).

### 2.2.4. Cysteine (Cys)

The S-H stretching vibration ( $\sim 2575\text{cm}^{-1}$ ) is a useful marker for cysteine hydrogen bonding, because, although its intrinsic IR or Raman intensity is low, it is in an isolated region of the vibrational spectrum that has a low background (48).

### 2.2.5. Cystine

The [disulfide bond](#) S-S stretch is readily detected in off-resonance Raman spectra, and its frequency is diagnostic for the disulfide conformation:  $\sim 505\text{cm}^{-1}$  for the *gauche-gauche-gauche* conformation,  $\sim 525\text{cm}^{-1}$  for *gauche-gauche-trans*, and  $\sim 550\text{cm}^{-1}$  for *trans-gauche-trans* (49, 50).

## 2.3. Protein-Bound Chromophores

### 2.3.1. Polyenes

Carotenoids give rise to strong RR scattering (51, 52). The retinal chromophore has been studied extensively in [rhodopsin](#) and [bacteriorhodopsin](#) by RR spectroscopy (52, 53), which has provided structural details of their photocycles.

### 2.3.2. Tetrapyrrole Cofactors

Heme proteins provide rich RR spectra of the heme group that are sensitive to protein-induced distortions (54), and to the oxidation, ligation, and spin states of the central Fe ion (55). The status of heme-bound CO, NO, and O<sub>2</sub> (56) can be assessed through detection of their vibrations, as can the strength of the bond from Fe<sup>2+</sup> to the proximal histidine ligand (57). RR spectroscopy has also been applied to reduced tetrapyrrole chromophores, including chlorins, chlorophylls, bacteriochlorophylls (58), isobacteriochlorins (59), corrins (60), and biliverdin (61).

### 2.3.3. Other Aromatic Chromophores

RR studies have been reported for **flavin** (62), folate (63), molybdopterin (64), pyrroloquinolinequinone (65), **pyridoxal** (66), nicotinamide (67), and **nucleotide** chromophores (68), as well as for tyrosyl radicals (69).

### 2.3.4. Metal Ions

Vibrations of metal-bound ligands can be enhanced in RR spectra excited in ligand-metal charge-transfer absorptions. Studies have been reported for iron-sulfur (70), iron-tyrosinate (71), copper-cysteinate (72), and molybdenum-dithiolene (64) proteins.

## 2.4. Protein Complexes

Protein complexes have been studied by various sensitive difference techniques and artificial labeling. Specific examples include **antibody/antigen**, **hapten/antibody**, **coenzyme/enzyme**, **ligand/hemeprotein**, drug/enzyme, enzyme/ substrate, protein/ **membrane**, and protein/nucleic acid interactions (52, 73-75).

## 3. Nucleic Acids

### 3.1. Backbone Vibrations

The vibrations of the sugar-phosphate backbone are sensitive to conformation. Bands that are diagnostic (Table 3) for secondary structure (A-form, B-form, etc) include the symmetric stretching of the P-O bonds connecting to the sugar rings (~800cm<sup>-1</sup>) and the symmetric and antisymmetric stretching modes of the terminal P-O bonds (~1090 and ~1225cm<sup>-1</sup>) (76, 77).

**Table 3. Conformationally Sensitive Nucleic Acid Vibrations (in cm<sup>-1</sup>)<sup>a</sup>**

| Assignment                   | A form    | B form    | Z form      |
|------------------------------|-----------|-----------|-------------|
| Thy                          | 642       | 665       | —           |
| Gua                          | 640-665   | 682       | 625         |
| backbone/Cyt                 | 783       | 784       | 784         |
| backbone/Thy                 | 779       | 793       | —           |
| O-P-O                        | 800-810   | 820-840   | 740-750     |
| backbone                     | 852       | —         | 855         |
| PO <sub>2</sub> <sup>-</sup> | 1090-1100 | 1085-1090 | 1095        |
| Gua/Cyt                      | 1180      | 1180      | 1180 / 1188 |
| Thy                          | 1239      | 1208      | —           |
| Cyt                          | 1242      | 1240      | 1246        |
| Cyt                          | 1252      | 1260      | 1265        |
| Gua                          | 1314      | 1318      | 1317        |

|     |      |      |      |
|-----|------|------|------|
| Ade | 1334 | 1341 | —    |
| Gua | 1361 | 1362 | 1355 |
| Ade | 1478 | 1483 | —    |
| Gua | 1482 | 1489 | 1486 |

<sup>a</sup> Adapted from W. L. Peticolas and E. Evertsz (76).

### 3.2. Base Vibrations

The purine and pyrimidine bases give rise to many ring vibrations (Table 4), leading to complicated IR and Raman spectra. RR enhancement in the UV absorption bands can be used to discriminate between the different bases, as the excitation profiles differ (13, 78). Base stacking in duplex and triplex structures diminishes RR intensities as it does the absorptivity; RR hypochromism (79-81) can be monitored in a base-selective manner. Frequencies of the purine and pyrimidine ring bending modes at  $\sim 670\text{cm}^{-1}$  and  $\sim 770\text{cm}^{-1}$  are sensitive to the conformation of the sugar ring (76, 82, 83). The carbonyl stretching vibrations ( $\sim 1650\text{cm}^{-1}$ ) are sensitive to hydrogen bonding and dipole coupling effects in the base pairs (83, 84). Secondary structure-sensitive base vibrations are also given in Table 3.

**Table 4. The Vibrations of the Nucleotides<sup>a</sup>**

| Observed Frequency<br>( $\text{cm}^{-1}$ ) | Assignment  |
|--|---|
|  | <i>Adenine (dAMP)</i>                                 |
| 1649                                       | C5C6 (48), dNH <sub>2</sub> (20)                      |
| 1604                                       | dNH <sub>2</sub> (73), C5C6 (15), C6N6' (14)          |
| 1581                                       | C5C4 (48), C4N3 (31)                                  |
| 1509                                       |   |
| 1482                                       | dC2H (29), N9C8 (19), dC8H                            |
| 1423                                       | C4N9 (44), dC8H                                       |
| 1340                                       | N7C5 (39), C8N7 (12)                                  |
| 1309                                       | N9C8 (30), N3C2 (14), dC8H (14), dC2H                 |
| 1254                                       | N1C6 (31), C6N6' (26)                                 |
| 1220                                       |   |
| 1173                                       |   |
| 1012                                       |   |
| 730  | dN7C8N9 (19), N9R (14), dC5N7C8 (12),<br>dC4N9C8 (11) |
|  | <i>Cytosine(dCMP)</i>                                 |
| 1660                                       | C4C5 (38), dNH <sub>2</sub> (29)                      |
| 1652                                       | C2O (67), C2N3 (15)                                   |
| 1605                                       | dNH <sub>2</sub> (58), C4N4' (19)                     |

|      |   |
|------|---|
| 1528 | N3C4 (38), N1C2 (14)                                      |
| 1500 | N1C2 (38), N1C6 (34), C2N3 (27)                           |
| 1412 |   |
| 1374 | C4N4' (29), N1C2 (16)                                     |
| 1294 | N1C6 (28), C5C6 (19)                                      |
| 1250 | dC6H (31), C4N4' (19)                                     |
| 1237 |   |
| 1210 |   |
| 1144 |   |
| 988  | NH <sub>2rock</sub> (36), dC6H (24)                       |
| 782  | N1R (18), C4N (14), dC4C5C6 (13)<br><i>Guanine (dGMP)</i> |
| 1679 | C6O (48), C5C6 (21), C5C4 (11), dN1H (11),<br>N1C6 (10)   |
| 1679 |   |
| 1603 | dNH <sub>2</sub> (83), C2N2' (15)                         |
| 1579 | C4N3 (30), C5C4 (24), N7C5 (16)                           |
| 1539 | C4N9 (33), N7C5 (24)                                      |
| 1489 | dC8H (40), N9C8 (32), C8N7 (21)                           |
| 1419 |   |
| 1364 | C8N7 (26), N1C6 (25), N7C5 (16)                           |
| 1326 | dC8H (25), C8N7 (19)                                      |
| 1253 |   |
| 1217 |   |
| 1179 |   |
| 1082 |   |
| 1034 |   |
| 854  | N9R (16), N3C2 (13)                                       |
| 679  | dN7C8N9 (15), dC5N7C8 (15)<br><i>Uridine (dUMP)</i>       |
| 1686 | C2O (47), C2N3 (24), dN3H (18)                            |
| 1674 | C5C4 (34), C4O (34)                                       |
| 1628 | N3C4 (22), N1C2 (21), C6C5 (20), N1C6 (19)                |
| 1476 | N1C2 (38), C2N3 (17)                                      |
| 1394 | dN3H (48), C4O (27)                                       |
| 1276 |   |
| 1230 | dC6H (23), C2N3 (15)                                      |
| 1000 | dC6H (45), C5C4 (10)                                      |
| 783  | N1C2 (14), N1R (10), C5C4 (10), N1C6 (9), N3C4<br>(8)     |

---

<sup>a</sup> Adapted from S. P. A. Fodor, R. P. Rava, T. R. Hays and T. G. Spiro (13).

### 3.3. Nucleic Acid Complexes

Nucleic acid complexes have also been studied by various difference techniques and artificial labeling. Specific examples include complexes with anticancer drugs, antibiotics, histones, and capsid proteins, DNA/gene regulatory proteins, DNA/RNA hybrids, and 3- and 4-stranded nucleic acids ([85-87](#)).

### 4. Lipids

These molecules do not absorb at accessible laser wavelengths and provide no opportunity for RR enhancement, but because the local molecular concentration is high in membranes, IR and Raman spectra can be obtained for [liposomes](#), planar bilayers, and [membranes](#) using surface-sensitive techniques ([88-91](#)). Lipid vibrational frequencies are given in Table [5](#). Particularly useful are the C-C vibrations of the hydrocarbon tails, which give rise to bands at  $\sim 1065$  and  $\sim 1130\text{cm}^{-1}$  for the all *trans* conformation and at  $\sim 1100\text{cm}^{-1}$  for the *gauche* conformation. The relative intensities of these bands provide a measure of lipid ordering ([93-95](#)).

**Table 5. Vibrational Modes of Lipids and Phospholipids<sup>a</sup>**

| Observed Frequency ( $\text{cm}^{-1}$ ) | Assignment                                      |
|---|---|
| 218 / 224                               | Longitudinal acoustical mode                    |
| 710-720                                 | C-N sym. stretch                                |
| 725                                     | O-P-O sym. stretch                              |
| 860-890                                 | C-C stretch involving acyl carbon               |
| 870-875                                 | C-N stretching                                  |
| 960-1130                                | C-C stretching                                  |
| 967                                     | CH deformation                                  |
| 1065                                    | C-C stretch (all <i>trans</i> )                 |
| 1070-1075                               | C-O stretching                                  |
| 1095-1105                               | O-P-O stretching                                |
| 1100                                    | C-C stretch ( <i>gauche</i> )                   |
| 1130                                    | C-C stretch (all <i>trans</i> )                 |
| 1230                                    | $\text{PO}_2^-$ stretching                      |
| 1200-1470                               | CH / $\text{CH}_2$ / $\text{CH}_3$ deformations |
| 1650-1666                               | CC ( <i>cis</i> )                               |
| 1670-1680                               | C=C ( <i>trans</i> )                            |
| 1730                                    | C-O stretching (ester)                          |
| 1732                                    | CO stretching                                   |
| 2800-3000                               | C-H stretching                                  |

<sup>a</sup> Adapted from A. T. Tu ([92](#)).

The degree of unsaturation can be determined from the ratio of the intensities of the  $\sim 1303\text{cm}^{-1}$  methylene and  $\sim 1265\text{cm}^{-1}$  methine bending vibrations, whereas the position of the C=C stretch is diagnostic for the geometric isomer:  $1675\text{cm}^{-1}$  for *trans* and  $1660\text{cm}^{-1}$  for *cis* (96). The C-H stretches (between  $2800$  and  $3000\text{cm}^{-1}$ ) are sensitive to the lateral packing density of the chains. The ratio of the intensities of the two methylene vibrations at  $\sim 2880$  and  $\sim 2850\text{cm}^{-1}$  is 1 in the crystalline state and decreases as the density decreases (97). Finally, the C-N vibration at  $\sim 715\text{cm}^{-1}$  is sensitive to the O-C-C-N<sup>+</sup> conformation of the head group (98).

## 5. Carbohydrates

As is the case with lipids, these molecules are unsuitable for RR spectroscopy; however, at high concentrations they provide rich IR and Raman spectra (99) (Table 6). Of particular use in probing structural details of these molecules are the ring vibrations at about  $920$  and  $770\text{cm}^{-1}$  and a C-H bending vibration at about  $865\text{cm}^{-1}$ . These modes are sensitive to the a and b conformers of the sugar ring (73, 100, 101). The most sensitive of the three is the C-H bending vibration found at  $\sim 840\text{cm}^{-1}$  in the a conformer and at  $\sim 890\text{cm}^{-1}$  for the b conformer.

**Table 6. Vibrational Frequencies of a-D-Glucose<sup>a</sup>**

| Frequency ( $\text{cm}^{-1}$ ) | Assignment                   |
|--------------------------------|------------------------------|
| 3200-3400                      | O-H stretching               |
| 2800-3000                      | C-H stretching               |
| 1462                           | CH <sub>2</sub> bend         |
| 1442-1360                      | CH bend +OH bend             |
| 1345-1350                      | COH bend                     |
| 1335-1340                      | CH <sub>2</sub>              |
| 1328                           | CH <sub>3</sub>              |
| 1298                           |                              |
| 1270-1280                      | C(6)OH +C (1)OH              |
| 1250                           | C(1)H bend                   |
| 1220-1225                      | CH <sub>2</sub>              |
| 1205                           |                              |
| 1150                           | CO, CC stretch +CH, COH bend |
| 1130                           |                              |
| 1115                           |                              |
| 1070-1076                      | C(1)H, COH bend              |
| 1040-1055                      | C(1)H bend                   |
| 1020-1026                      | COH bend                     |
| 1000-1010                      | CH <sub>3</sub>              |

|         |                      |
|---------|----------------------|
| 913     | C(1)H +COH bend      |
| 890-900 | CH                   |
| 860     | C(2)OH               |
| 720     | CC, CO stretch       |
| 380-600 | Skeletal vibrations  |
| 100-365 | Torsional vibrations |

---

<sup>a</sup> Adapted from J. Twardowski and P. Anzenbacher ([99](#)).

## 6. In vivo Studies

Raman and IR microspectroscopy can provide spatial resolution on the order of  $1\ \mu\text{m}^3$  and have been applied to the study of biological cells and organelles. Examples include bacteria, algae, and cells from plants, blood, sperm, tumors, and photoreceptors ([102](#)).

## 7. Time Resolution: Dynamics

Vibrational spectroscopy can be applied in a time-resolved mode to monitor the evolution of molecular structure following rapid mixing, temperature-jump, or photoexcitation. Examples include enzyme-substrate interactions ([103](#)); protein folding ([104](#)); the photocycles of [bacteriorhodopsin](#) ([105](#)) and of bacterial reaction centers ([106](#)) ligand photodissociation and re-binding in [myoglobin](#) ([107](#)) and cytochrome oxidase ([108](#)); and the **allosteric** transition in hemoglobin ([109](#)).

## 8. Other Vibrational Techniques

### 8.1. Chiroptical Techniques

Vibrational circular dichroism (VCD) ([110](#)) and Raman optical activity (ROA) ([111](#)) spectroscopies afford exquisite sensitivity to stereochemistry. The measurements, which involve the small intensity differences for left- and right-circularly polarized light, are technically demanding, but recent improvements in instrumentation offer considerable promise.

### 8.2. Surface Enhanced Raman Spectroscopy (SERS)

Very large enhancements are observed for Raman signals of molecules adsorbed on small metal particles, (especially silver and gold) through electromagnetic coupling of the vibrations to collective oscillations of the metallic electrons (*surface plasmons*) ([112](#)). Extremely high sensitivities can be achieved, but the preparation of reproducible surfaces has been problematic, although promising new techniques have been developed ([113-115](#)).

### 8.3. Non-linear Optical Techniques

Vibrational spectra can be produced by monitoring the intensities of multiple overlapping laser beams, which can interact through higher-order polarizabilities of the molecule. A variety of configurations are possible, including coherent anti-Stokes Raman spectroscopy (CARS) ([116](#)), Raman gain spectroscopy ([117](#)), and sum-frequency generation spectroscopy ([118](#)). They offer a number of advantages, including fluorescence rejection and sensitivity to molecular environment, but they are technically demanding.

## Bibliography

1. E. B. Wilson, J. C. Decius, and P. C. Cross (1955) *Molecular Vibrations*, McGraw-Hill, New York.
2. R. Barer, R. H. Cole, and H. W. Thompson (1949) *Nature* (London) **163**, 198.

3. (R. G. Messersmidt and M. A. Harthcock, eds.) *Infrared Micro-spectroscopy: Theory and Applications* (1988), Marcel Dekker, New York.
4. W. D. Bowman and T. G. Spiro (1980) *J. Raman Spec.* **9**, 369.
5. T. Hirschfeld and D. B. Chase (1986) *Appl. Spectrosc.* **40**, 133.
6. K. J. Rothschild (1992) *Journal of Bioenergetics and Biomembranes* **24**, 147.
7. T. G. Spiro and P. Stein (1977) *Ann. Rev. Phys. Chem.* **28**, 501.
8. B. B. Johnson and W. L. Peticolas (1976) *Ann. Rev. Phys. Chem.* **27**, 465.
9. P. R. Carey (1982) *Biological Applications of Raman and Resonance Raman Spectroscopies*, Academic Press, New York.
10. D. C. Blazej and W. L. Peticolas (1977) *Proc. Nat. Acad. Sci. USA*, **74**, 2639.
11. Y. Nishimura, A. Y. Hirakawa, and M. Tsuboi (1979) *Adv. Infrared Raman Spectrosc.* **5**, 217.
12. L. D. Ziegler and B. Hudson (1981) *J. Chem. Phys.* **74**, 982.
13. S. P. A. Fodor, R. P. Rava, T. R. Hays, and T. G. Spiro (1985) *J. Am. Chem. Soc.* **107**, 1520.
14. C. R. Johnson, M. Ludwig, S. O'Donnel, and S. A. Asher (1984) *J. Am. Chem. Soc.* **106**, 5008.
15. R. P. Rava and T. G. Spiro (1984) *J. Am. Chem. Soc.* **106**, 4062.
16. R. P. Rava and T. G. Spiro (1985) *J. Phys. Chem.* **89**, 1856.
17. S. P. A. Fodor, R. A. Copeland, C. A. Grygon, and T. G. Spiro (1989) *J. Am. Chem. Soc.* **111**, 5509.
18. L. C. Mayne, L. D. Ziegler, and B. S. Hudson (1985) *J. Phys. Chem.* **89**, 3395.
19. J. M. Dudik, C. R. Johnson, and S. A. Asher (1985) *J. Phys. Chem.* **89**, 3805.
20. Y. Wang, R. Purrello, and T. G. Spiro (1991) *J. Am. Chem. Soc.* **113**, 6359.
21. S. Krimm (1987) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **1**, pp. 1–46.
22. S. Krimm and J. Bandekar (1986) *Adv. Protein. Chem.* **38**, 181.
23. S. Krimm (1987) In *Vibrational Spectra and Structure*, J. R. Durig, ed., vol. **16**, Elsevier: New York.
24. W. K. Surewicz and H. H. Mantsch (1996) In *Spectroscopic Methods for Determining Protein Structure in Solution* (H. A. Havel, ed.), VCH, New York, pp. 135–162.
25. J. C. Austin, T. Jordan, and T. G. Spiro (1993) *Biomolecular Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **20**, pp. 55–127.
26. A. T. Tu (1986) In *Spectroscopy of Biological Systems* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, pp. 72–85.
27. Y. Sugawara, A. Y. Hirakawa, and M. Tsuboi (1984) *J. Mol. Spec.* **108**, 206.
28. D. S. Caswell and T. G. Spiro (1987) *J. Am. Chem. Soc.* **109**, 2796.
29. H. Takeuchi and I. Harada (1990) *J. Raman. Spec.* **21**, 509.
30. T. Jordan, I. Mukerji, Y. Wang, and T. G. Spiro (1996) *J. Mol. Struct.* **379**, 51.
31. B. S. Hudson and L. C. Mayne (1987) In *Biological Applications of Raman Spectroscopy*, (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **2**, pp. 181–210.
32. D. S. Caswell and T. G. Spiro (1986) *J. Am. Chem. Soc.* **108**, 6470.
33. I. Harada and H. Takeuchi (1986) In *Spectroscopy of Biological Systems*, (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **13**, pp. 113–126.
34. H. Takeuchi, N. Watanabe, Y. Satoh, and I. Harada (1989) *J. Raman Spec.* **20**, 233.
35. K. R. Rodgers, C. Su, S. Subramaniam, and T. G. Spiro (1992) *J. Am. Chem. Soc.* **114**, 3697.
36. M. N. Siamwiza, R. C. Lord, M. C. Chen, T. Takamatsu, I. Harada, H. Matsuura, and T. Shimanouchi (1975) *Biochemistry* **14**, 4870.
37. J. Twardowski and P. Anzenbacher (1994) *Raman and IR Spectroscopy in Biology and*



*Biochemistry*, Ellis Horwood, New York, pp. 117–118.

38. P. G. Hildebrandt, R. A. Copeland, T. G. Spiro, J. Otlewski, M. Laskowski Jr., and F. G. Prendergast (1988) *Biochemistry* **27**, 5426.
39. N. T. Yu, B. H. Jo, and D. C. O'Shea (1973) *Arch. Biochem. Biophys.* **156**, 171.
40. N. T. Yu (1974) *J. Am. Chem. Soc.* **96**, 4664.
41. T. Miura, H. Takeuchi, and I. Harada (1989) *J. Raman Spec.* **20**, 667.
42. T. Miura, H. Takeuchi, and I. Harada (1988) *Biochemistry* **27**, 88.
43. T. Maruyama and H. Takeuchi (1995) *J. Raman Spec.* **26**, 319.
44. T. Jordan, J. C. Eads, and T. G. Spiro (1995) *Protein Science* **4**, 716.
45. G. -Y. Liu, C. A. Grygon, and T. G. Spiro (1989) *Biochemistry* **28**, 5046.
46. H. Takeuchi, Y. Kimura, I. Koitabashi, and I. Harada (1991) *J. Raman Spec.* **22**, 233.
47. I. Harada, T. Takamatsu, M. Tasumi, and R. C. Lord (1982) *Biochemistry* **21**, 3674.
48. H. Li and G. J. Thomas Jr. (1991) *J. Am. Chem. Soc.* **113**, 456.
49. H. Sugeta, A. Go, and T. Miyazawa (1972) *Chem. Lett.* **83**.
50. H. Sugeta, A. Go, and T. Miyazawa (1973) *Bull. Chem. Soc. Japan* **46**, 3407.
51. D. Gill, R. G. Kilponen, and L. Rimai (1970) *Nature (London)* **227**, 743.
52. J. Twardowski and P. Anzenbacher (1994) *Raman and IR Spectroscopy in Biology and Biochemistry*, Ellis Horwood, New York, pp. 140–150.
53. R. A. Mathies, S. O. Smith, and I. Palings (1987) In *Biological Applications of Raman Spectroscopy*, (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **2**, pp. 59–108.
54. S. Hu, I. K. Morris, J. P. Singh, K. M. Smith, and T. G. Spiro (1993) *J. Am. Chem. Soc.* **115**, 12446.
55. T. G. Spiro and X.-Y. Li (1988) In *Biological Applications of Raman Spectroscopy*, (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 1–38.
56. E. A. Kerr and N.-T. Yu (1986) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 39–96.
57. T. Kitagawa (1988) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 97–132.
58. M. Lutz and B. Robert (1988) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 347–412.
59. S. Han, J. F. Madden, R. G. Thompson, S. H. Strauss, L. M. Siegel, and T. G. Spiro (1989) *Biochemistry* **28**, 5461.
60. S. Dong, R. Padmakumar, R. Banerjee, and T. G. Spiro (1996) *J. Am. Chem. Soc.* **118**, 9182.
61. J. Matysik, P. Hildebrandt, K. Smit, A. Korkin, F. Mark, W. Gaertner, S. E. Braslavsky, K. Schaffner, and B. Schrader (1995) *J. Mol. Struct.* **348**, 225.
62. J. T. McFarland (1987) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **2**, pp. 211–302.
63. J. C. Austin, A. Fitzhugh, J. E. Villafranca, and T. G. Spiro (1995) *Biochemistry* **34**, 7678.
64. L. Kilpatrick, K. V. Rajagopalan, J. Hilton, N. R. Bastian, R. S. Pilato, and T. G. Spiro (1995) *Biochemistry* **34**, 3032.
65. R. S. Moog, M. A. McGuirl, C. E. Cote, and D. M. Dooley (1986) *Proc. Natl. Acad. Sci. USA* **83**, 8435.
66. M. J. Benecky, R. A. Copeland, R. P. Rava, R. Feldhaus, R. D. Scott, C. M. Metzler, D. E. Metzler, and T. G. Spiro (1985) *J. Biol. Chem.* **260**, 11671.
67. E. G. Rodgers and W. L. Peticolas (1980) *J. Raman Spec.* **9**, 3729.
68. M. Molina, Y. Wang, R. Purrello, and T. G. Spiro (1991) *J. Raman Spectrosc.* **22**, 205.
69. M. L. McGlashen, D. D. Eads, T. G. Spiro, and J. W. Whittaker (1995) *J. Phys. Chem.* **99**,

70. T. G. Spiro, R. S. Czernuszewicz, and S. Han (1988) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 523–554.
71. L. Que Jr. (1988) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **3**, pp. 491–522.
72. C. R. Andrew and J. Sanders-Loehr (1996) *Acc. Chem. Res.* **29**, 365.
73. A. T. Tu (1982) *Raman Spectroscopy in Biology: Principles and Applications*, John Wiley & Sons, New York, pp. 117–130.
74. P. J. Tonge and P. R. Carey (1993) In *Biomolecular Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **20**, pp. 129–161.
75. P. R. Carey (1982) *Biological Applications of Raman and Resonance Raman Spectroscopies*, Academic Press, New York, pp. 154–183.
76. W. L. Peticolas and E. Evertsz (1992) *Methods in Enzymology* **211**, 335.
77. B. Prescott, W. Steinmetz, and G. J. Thomas Jr. (1984) *Biopolymers* **23**, 235.
78. M. Tsuboi, Y. Nishimura, A. Y. Hirakawa, and W. L. Peticolas (1987) *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **2**, pp. 109–180.
79. B. Tomlinson and W. L. Peticolas (1970) *J. Chem. Phys.* **52**, 2154.
80. E. W. Small and W. L. Peticolas (1971) *Biopolymers* **10**, 1377.
81. S. P. A. Fodor and T. G. Spiro (1986) *J. Am. Chem. Soc.* **108**, 3198.
82. W. L. Peticolas, W. L. Kubasek, G. A. Thomas, and M. Tsuboi (1986) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **2**, pp. 81–134.
83. Y. Nishimura and M. Tsuboi (1986) In *Spectroscopy of Biological Systems* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **13**, pp. 177–232.
84. L. Lafleur, J. Rice, and G. J. Thomas Jr. (1972) *Biopolymers* **11**, 2423.
85. P. R. Carey (1982) *Biological Applications of Raman and Resonance Raman Spectroscopies*, Academic Press, New York, pp. 185–207.
86. M. Manfait and T. Theophanides (1986) In *Spectroscopy of Biological Systems* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **13**, pp. 311–346.
87. G. J. Thomas Jr. and M. Tsuboi (1993) *Adv. Biophys. Chem.* JAI Press, Inc., Greenwich, Conn, pp. 1–70.
88. R. N. A. H. Lewis and R. N. McElhaney (1996) In *Infrared Spectroscopy of Biomolecules* (H. H. Mantsch and D. Chapman, eds.), Wiley-Liss, Inc., New York, pp. 159–202.
89. J. C. Conboy, M. C. Messmer, and G. L. Richmond (1996) *J. Phys. Chem.* **100**, 7617.
90. T. M. Cotton, J.-H. Kim, and R. E. Holt (1992) *Adv. Biophys. Chem.* **2**, 115.
91. K. B. Eisenthal (1992) *Ann. Rev. Phys. Chem.* **43**, 627.
92. A. T. Tu (1982) *Raman Spectroscopy in Biology: principles and applications*, John Wiley & Sons, New York, pp. 204–205.
93. J. L. Lippert and W. L. Peticolas (1971) *Proc. Natl. Acad. Sci.* **68**, 1572.
94. R. G. Snyder, D. G. Cameron, H. L. Casal, D. A. C. Compton, and H. H. Mantsch (1982) *Biochim. Biophys. Acta* **684**, 111.
95. P. Yager and B. P. Gaber (1987) In *Biological Applications of Raman Spectroscopy* (T. G. Spiro, ed.), John Wiley & Sons, New York, Vol. **1**, pp. 203–262.
96. G. F. Bailey and R. J. Horvat (1972) *J. Am. Chem. Soc.* **49**, 494.
97. B. P. Gaber and W. L. Peticolas (1977) *Biochim. Biophys. Acta* **465**, 260.
98. H. Akutsu (1981) *Biochemistry* **20**, 7359.
99. J. Twardowski and P. Anzenbacher (1994) *Raman and IR Spectroscopy in Biology and*

*Biochemistry*, Ellis Horwood, New York, pp. 241–242.

100. S. A. Barker, E. J. Bourne, M. Stacey, and D. H. Whiffen (1954) *J. Chem. Soc.* 171.
101. S. A. Barker, E. J. Bourne, R. Stephens, and D. H. Whiffen (1954) *J. Chem. Soc.* 3468.
102. J. Greve and G. J. Puppels (1993) *Biomolecular Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **20**, pp. 231–265.
103. P. R. Carey (1985) *Springer Proc. Phys.* **4**, 233.
104. S. Williams, T. P. Causgrove, R. Gilmanishin, K. S. Fang, R. H. Callender, W. H. Woodruff, and R. B. Dyer (1996) *Biochemistry* **35**, 691.
105. G. G. Kochendoerfer and R. A. Mathies (1995) *Isr. J. Chem.* **35**, 211.
106. K. Wynne, G. Haran, G. D. Reid, C. C. Moser, P. L. Dutton, and R. M. Hochstrasser (1996) *J. Phys. Chem.* **100**, 5140.
107. T. P. Causgrove and R. B. Dyer (1996) *J. Phys. Chem.* **100**, 3273.
108. R. B. Dyer, K. A. Peterson, P. O. Stoutland, and W. H. Woodruff (1994) *Biochemistry* **33**, 500.
109. V. Jayaraman, K. R. Rodgers, I. Mukerji, and T. G. Spiro (1995) *Science* **269**, 1843.
110. T. A. Keiderling and P. Pancoska (1993) In *Biomolecular Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **21**, pp. 267–315.
111. L. D. Barron and L. Hecht (1993) In *Biomolecular Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester, Vol. **21**, pp. 235–266.
112. T. M. Cotton, J.-H. Kim, and G. D. Chumanov (1991) *J. Raman Spec.* **22**, 729.
113. G. Chumanov, K. Sokolov, B. W. Gregory, and T. M. Cotton (1995) *J. Phys. Chem.* **99**, 9466.
114. K. C. Grabar, R. G. Freeman, M. B. Hommer, and M. J. Natan (1995) *Anal. Chem.* **67**, 735.
115. R. P. Van Duyne, J. C. Hulteen, and D. A. Treichel (1993) *J. Chem. Phys.* **99**, 2101.
116. F. W. Schneider (1982) In *Non-Linear Raman Spectroscopy and its Chemical Applications* (W. Kiefer and D. A. Long, eds.), Reidel Publishing Co., Dordrecht, Holland, pp. 445–459.
117. M. D. Morris and R. J. Bienstock (1982) In *Non-Linear Raman Spectroscopy and its Chemical Applications* (W. Kiefer and D. A. Long, eds.), Reidel Publishing Co., Dordrecht, Holland, pp. 543–559.
118. Y. R. Shen (1989) *Nature* **337**, 519.

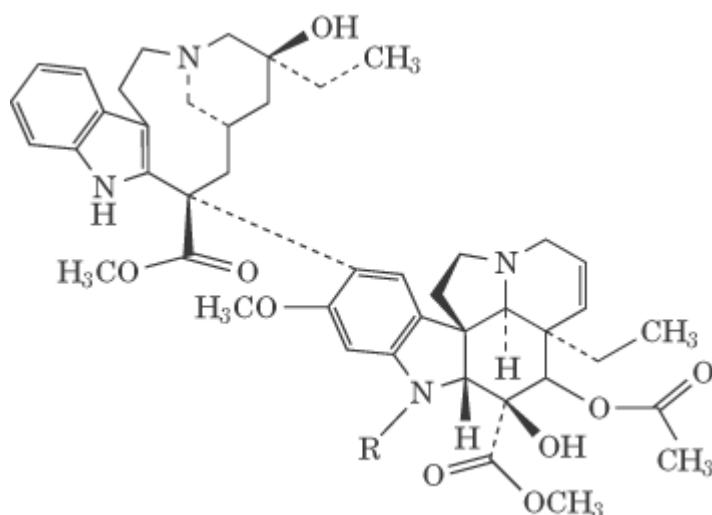
### **Suggestions for Further Reading**

119. P. R. Carey (1982) *Biological Applications of Raman and Resonance Raman Spectroscopies*, Academic Press, New York.
120. *Spectroscopy of Biological Systems* (1986), In *Advances in Spectroscopy* (R. J. H. Clark and R. E. Hester, eds.) John Wiley & Sons, Chichester.
121. *Biomolecular Spectroscopy* (1993), vols. **20** and **21** (R. J. H. Clark and R. E. Hester, eds.), John Wiley & Sons, Chichester.
122. *Spectroscopic Methods for Determining Protein Structure in Solution* (1996), (H. A. Havel, ed.), VCH, New York.
123. *Metallobiochemistry* (1993), In *Methods in Enzymology* (J. F. Riordan and B. L. Vallee, eds.), Academic Press, Inc., San Diego.
124. *Biological Applications of Raman Spectroscopy* (1986) (T. G. Spiro, ed.), John Wiley & Sons, New York.
125. A. T. Tu (1982) *Raman Spectroscopy in Biology: Principles and Applications*, John Wiley & Sons, New York.
126. J. Twardowski and P. Anzenbacher (1994) *Raman and IR Spectroscopy in Biology and Biochemistry*, Ellis Horwood, New York.

## Vinblastine

Vinblastine and its analogue vincristine (Fig. 1), known as “vinca alkaloids,” are natural products obtained from the ornamental flowering plant *Catharanthus roseus*. They were discovered to be potent antimitotic agents in 1959–1960 and since then have become effective chemotherapeutic drugs for treating a number of leukemias and lymphomas, both singly and in combination with other anticancer drugs. Several new vinca alkaloids, including vindesine and vinorelbine, have been introduced into the clinic recently and have shown broad antitumor activity. The vinca alkaloids, and especially vinblastine, have also become important tools in molecular and cell biology for studying the roles of [microtubules](#) and their dynamics in cellular processes. The value of vinblastine as a molecular tool derives from its ability to inhibit microtubule polymerization and, at very low concentrations, to kinetically stabilize microtubule dynamics.

**Figure 1.** Structures of vinblastine (R = CH<sub>3</sub>) and vincristine (R = CHO).



### 1. Vinblastine Binding to Tubulin and Microtubules

Vinblastine binds to tubulin rapidly, reversibly, and independently of temperature between 0°C and 37°C. Vinblastine binding to tubulin induces an unusual conformational change in tubulin that is associated with tubulin self-association (1); this phenomenon appears to be responsible for the vinblastine-tubulin paracrystal formation in cells and *in vitro* that occurs at high vinblastine concentrations ( $\geq 10 \mu\text{M}$ ) (see Ref. 2). The vinblastine-induced increase in the affinity of tubulin for itself also appears to play a critical role in the ability of vinblastine to bind strongly to microtubule ends and to stabilize microtubule dynamics (see text below).

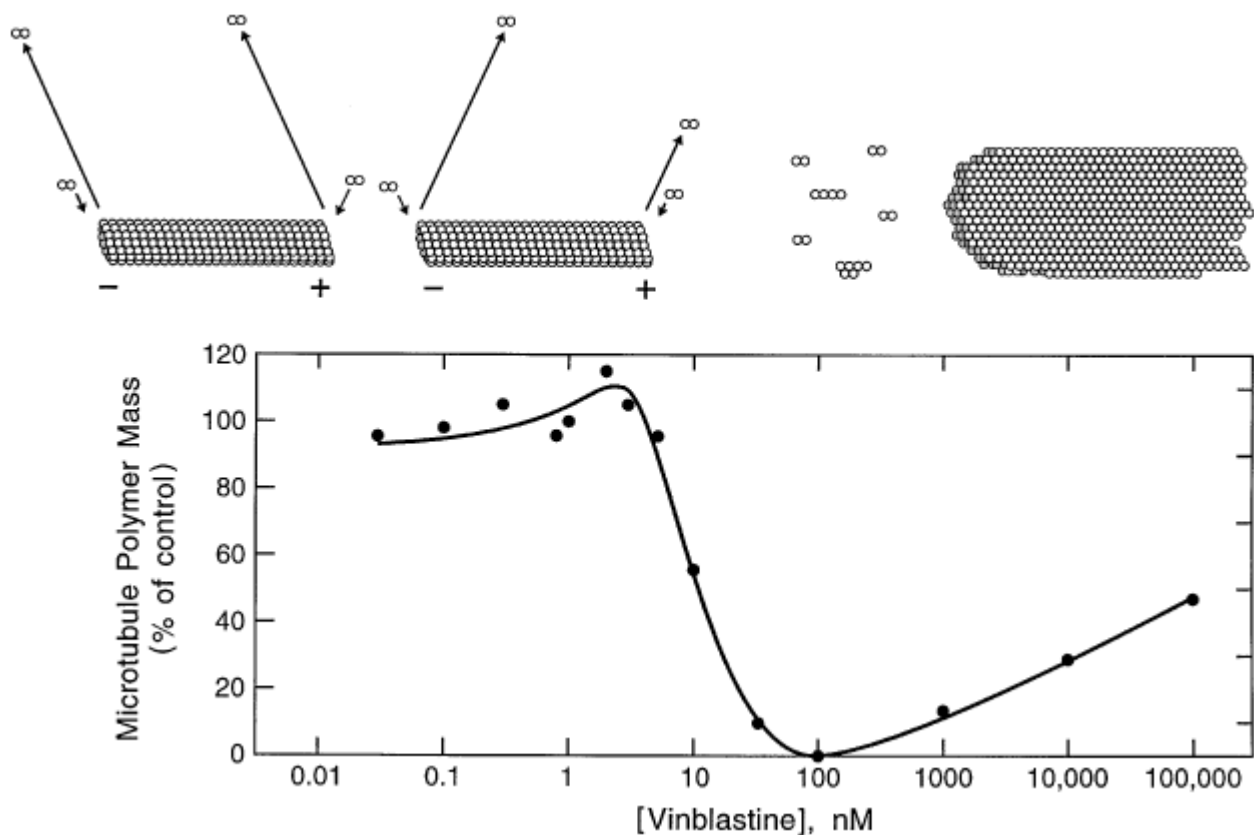
In contrast to the antimitotic drug [colchicine](#), vinblastine can bind directly to microtubules without first forming a complex with soluble tubulin (3). Vinblastine binds to microtubule ends with relatively high affinity ( $K_d$  of 1 to 2  $\mu\text{M}$ ), and it binds with low affinity ( $K_d$  of 0.25 to 0.3 mM) to multiple tubulin sites located along the sides of the microtubule cylinder. There are only 16 to 17

high-affinity binding sites per microtubule end, and the binding of vinblastine to these few sites is responsible for the powerful suppression of tubulin exchange at the ends (described below). The binding of the vinca alkaloids to the low-affinity tubulin binding sites along the microtubule surface is responsible for the ability of relatively high vinca alkaloid concentrations to depolymerize microtubules and to induce tubulin paracrystal formation in cells.

## 2. Effects of Different Vinblastine Concentrations on the Mass, Dynamics, and Organization of Microtubules in Cells

The effects of vinblastine on tubulin and microtubule organization *in vitro* and in cells are strongly concentration-dependent. The cartoon diagrams in Figure 2 (top) show the effects of vinblastine (10 pM to 100  $\mu$ M) on dynamic instability and organization of the tubulin polymers in cells (see Refs. (4-6)). Figure 2 (bottom) shows the effects of vinblastine on the mass of tubulin polymer in [HeLa Cells](#). In the absence of vinblastine, microtubules undergo dynamic instability, episodes of rapid shortening and slow growing at both microtubule ends (top, left-hand microtubule; the rates are arbitrarily represented by the lengths of the arrows; see [Microtubules](#)). Low nanomolar concentrations of vinblastine suppress dynamic instability preferentially at plus ends, without changing the microtubule polymer mass (top, right-hand microtubule). The growing and shortening rates at the plus ends are depicted as being reduced by 50% (shortened arrows), whereas the minus-end rates are unchanged. Intermediate (low micromolar) vinblastine concentrations depolymerize the microtubules. Vinblastine induces self-association of tubulin; thus both tubulin dimers and small oligomers are

**Figure 2.** Effects of vinblastine (10 pM to 100 mM) on microtubules and tubulin. The cartoon drawings (top) show the effects of vinblastine on the organization and dynamics of microtubules and tubulin. The graph (below) shows the effects of vinblastine on the mass of polymerized tubulin. See text for explanation.



Intermediate (low micromolar) vinblastine concentrations depolymerize the microtubules. Vinblastine induces self-association of tubulin; thus both tubulin dimers and small oligomers are

represented at intermediate concentrations [free dimers and small oligomers; Fig. 2 (top); 10 nM to 1000 nM]. High concentrations of vinblastine ( $\geq 10,000$  nM) induce aggregation of tubulin into large paracrystalline arrays (top, far right). Figure 2, (bottom) depicts the polymer mass as a function of vinblastine concentration as occurs in HeLa cells. Low vinblastine concentrations ( $< 1$  to  $5$  nM) do not affect microtubule polymer mass. Intermediate vinblastine concentrations reduce polymer mass in association with increased soluble tubulin and small oligomers, and very high vinblastine concentrations increase polymer mass in association with formation of large tubulin aggregates (paracrystalline arrays).

### 3. Kinetic Suppression of Microtubule Dynamics

Microtubules exhibit two kinds of nonequilibrium dynamics that are important for their functions in cells, [treadmilling](#) and dynamic instability (see [Microtubules](#)). Recent studies have revealed that vinblastine and many other compounds that depolymerize microtubules (see [Colchicine](#)) strongly suppress these dynamics at relatively low concentrations in the absence of appreciable microtubule depolymerization. Vinblastine was found some years ago to inhibit treadmilling *in vitro* and the rate of microtubule disassembly upon dilution of the microtubules (called *kinetic capping*; see Ref. 7). However, its stabilizing effects on dynamics were only fully appreciated with the introduction of differential-interference contrast video microscopy, which enabled one to visualize in real time the stabilizing action of the drug on the growing and shortening dynamics of individual microtubules. Small numbers of vinblastine molecules bound to bovine brain microtubules *in vitro* strongly stabilized microtubule plus ends by suppressing the rate and extent of growing and shortening, decreasing the catastrophe frequency, increasing the rescue frequency, and increasing the percentage of time that the microtubule spent in a state of pause or attenuated (undetectable) dynamic instability (Table 1) (5). In contrast, vinblastine destabilized microtubule minus ends. Vinblastine suppresses plus-end microtubule dynamics in living cells in a similar manner. For example, in BS-C-1 monkey kidney cells, 32 nM vinblastine suppressed the rate of microtubule shortening by 67%, the rate of growing by 20%, and the catastrophe frequency by 67%, and it increased the percentage of time in pause or attenuation by 64% (6) (Table 1). Mitosis was strongly reduced at this vinblastine concentration, and no net microtubule depolymerization accompanied suppression of dynamics. Suppression of mitotic spindle microtubule dynamics appears to be one of the key actions of the vinca alkaloids that result in mitotic block at the metaphase–anaphase transition.

**Table 1. Comparison of Effects of Vinblastine on the Dynamic Instability Parameters of Microtubules in BS-C-1 Cells and Microtubules Assembled from Purified Bovine Brain Tubulin *in Vitro*<sup>a</sup>**

| Parameter                         | Control<br>Microtubules<br>in Cells | Microtubules<br>in<br>Cells, 32 nM<br>Vinblastine<br>(%) | Control<br>Microtubules<br><i>in Vitro</i> | Microtubules<br><i>in Vitro</i> ,<br>0.2 $\mu$ M<br>Vinblastine<br>(%) |
|-----------------------------------|-------------------------------------|--|--|--|
| Rate of<br>growing                | 6.9 $\pm$ 3.9 $\mu$ m/min           | –20  | 1.0 $\pm$ 0.1 $\mu$ m/min                  | –52  |
| Rate of<br>shortening             | 15.5 $\pm$ 10.6 $\mu$ m/min         | –67  | 14.4 $\pm$ 1.2 $\mu$ m/min                 | –56  |
| Growing<br>length <sup>a</sup>    | 1.8 $\pm$ 1.6 $\mu$ m               | –50  | 2.3 $\pm$ 0.2 $\mu$ m                      | –61  |
| Shortening<br>length <sup>a</sup> | 2.4 $\pm$ 2.4 $\mu$ m               | –71  | 5.8 $\pm$ 0.7 $\mu$ m                      | –75  |

|                         |                |           |                 |      |
|-------------------------|----------------|-----------|-----------------|------|
| Time in pause, %        | 36             | +64       | 7.7             | +40  |
| Catastrophe frequency   | 0.03±0.01/sec  | -67       | 0.005±0.001/sec | -45  |
| Rescue frequency        | 0.07±0.01/sec  | No change | 0.013±0.005/sec | +107 |
| Dynamicity <sup>a</sup> | 7.2±1.2 μm/min | -75       | 2.28 μm/min     | -82  |

<sup>a</sup> Data in columns 2 and 3 are from R. Dhamodharan, M. A. Jordan, D. Thrower, L. Wilson, and P. Wadsworth, *Mol. Biol. Cell* **6**, 1215–1229, 1995; data from columns 4 and 5 are from D. Panda, M. A. Jordan, K. C. Chin, and L. Wilson, *J. Biol. Chem.* **271**, 29807–29812, 1996; table is from M. A. Jordan and L. Wilson. The use of drugs to study the role of microtubule assembly dynamics in living cells, *Methods in Enzymology*, **298**, 252–276, 1998.

<sup>b</sup> *Growing length* is the mean microtubule length added during a growing event, *shortening length* is the mean microtubule length lost during a shortening event; *dynamicity* is the total microscopically measurable gain and loss of tubulin subunits at microtubule ends over time; it is a measure of overall dynamic instability.

#### 4. Inhibition of Mitosis by Vinblastine

One of the many microtubule-dependent processes affected by vinblastine is mitosis. Vinblastine totally blocks or slows mitosis in many eukaryotic cells. At high vinblastine concentrations (eg, >10nM in HeLa cells, >300nM in BS-C1 cells), mitotic block is accompanied by, and most likely results from, microtubule depolymerization. Mitotic spindles do not form. [Chromosomes](#) condense normally, but in the absence of microtubules they cannot move. The chromosomes ultimately separate along their lengths but remain attached at the [centromeres](#) (8). Mitotic block in some cells can last for more than 50 h and may result in return to interphase in the absence of anaphase or cytokinesis, accompanied by [chromatin](#) decondensation and reformation of multiple nuclei or multilobed nuclei (see Ref. 8).

At low vinblastine concentrations (<10nM in HeLa cells, <300nM in BSC-1 kidney cells), mitotic block occurs specifically at the transition from prometaphase/metaphase to anaphase. The block occurs in the presence of a full complement of spindle microtubules, and strong evidence indicates that blockage of mitotic progression at low vinblastine concentrations is due to perturbation of spindle microtubule dynamics (4); see also Suggestions for Additional Reading). Low concentrations of vinblastine also inhibit the normal tension-associated fluctuations in separation between sister kinetochores (9), suggesting that microtubule dynamics are required to produce spindle tension.

Low concentrations of vinblastine induce several structural alterations in mitotic spindles in association with mitotic block. Although bipolar spindles form and most chromosomes become attached at both kinetochores and are moved to the metaphase plate, a few to many chromosomes (depending upon the vinblastine concentration), remain near the spindle poles. This most likely occurs because these chromosomes become attached to the spindle at only one kinetochore because of the reduced dynamics. Long-term (20 h) incubation also induces formation of abnormal centrioles, accumulation of membranous vesicles in the centrosome, separation of mother and daughter centrioles by large distances, and a decrease in the number of microtubules attached to kinetochores of congressed chromosomes (10). The significance of these effects on centrosome organization is unknown.

#### Bibliography

1. G. C. Na and S. N. Timasheff (1980) *Biochemistry* **19**, 1347–1354.
2. K. Fujiwara and L. G. Tilney (1975) *Ann. N.Y. Acad. Sci.* **253**, 27–50.
3. L. Wilson, M. A. Jordan, A. Morse, and R. L. Margolis (1982) *J. Mol. Biol.* **159**, 129–149.
4. M. A. Jordan, D. Thrower, and L. Wilson (1991) *Cancer Res.* **51**, 2212–2222.
5. D. Panda, M. A. Jordan, K. Chin, and L. Wilson (1996) *J. Biol. Chem.* **271**, 29807–29812.
6. R. I. Dhamodharan, M. A. Jordan, D. Thrower, L. Wilson, and P. Wadsworth (1995) *Mol. Biol. Cell* **6**, 1215–1229.
7. M. A. Jordan and L. Wilson (1990) *Biochemistry*, **29**, 2730–2739.
8. C. G. Palmer et al. (1960) *Exp. Cell Res.* **20**, 198–265.
9. R. D. Shelby, K. M. Hahn, and K. F. Sullivan (1996) *J. Cell Biol.* **135**, 545–557.
10. K. L. Wendell, L. Wilson, and M. A. Jordan (1993) *J. Cell Sci.* **104**, 261–274.

### Suggestions for Further Reading

11. R. H. Himes (1991) *Pharmacol. Ther.* **51**, 257–267.
12. L. Wilson and M. A. Jordan (1994) "Pharmacological probes of microtubule functions, in" *Microtubules*, J. Hyams and C. Lloyd, eds., Wiley, New York, pp. 59–84.
13. M. A. Jordan and L. Wilson (1998) *Curr. Opin. Cell Biol.* **10**, 123–130.
14. M. A. Jordan and L. Wilson (1998) "The use of drugs to study the role of microtubule assembly dynamics in living cells, in" *Molecular Motors and the Cytoskeleton, Methods in Enzymology*, Vol. **298**, pp. 252–276.

## Viroids

Viroids were discovered in the early 1970s as a result of studies aimed at identifying and characterizing the agents of some **plant** diseases that were originally thought to be **viruses**. Viroids are the only well-defined class of subviral pathogens capable of autonomous **replication** ([1](#)). The 28 known viroid species (Table [1](#)) indicate that they are single-stranded, circular **RNA** molecules of between 246 and 401 bases. Such a **genome** is only about one-tenth the size of that of the smallest known RNA virus. Their other most striking feature is that they are naked RNA molecules; all the available evidence indicates that, in contrast to viruses, viroids do not code for any **protein**. Consequently, viroids are the lowest known step of the biological scale.

**Table 1. Viroid Species with their Abbreviations, Genomic Accession Numbers of Typical Sequence Variants, Sizes, and Genus and Family to which they Belong**

| Viroid species         | Abbreviation | Accession | Size (nt)    | Genus       | Family        |
|------------------------|--------------|-----------|--------------|-------------|---------------|
| Potato spindle tuber   | PSTVd        | V01465    | 356, 359-360 | Pospiviroid | Pospiviroidae |
| Tomato chlorotic dwarf | TCDVd        | AF162131  | 360          | Pospiviroid | Pospiviroidae |



|                            |         |        |                  |              |               |
|----------------------------|---------|--------|------------------|--------------|---------------|
| Mexican papita             | MPVd    | L78454 | 359-360          | Pospiviroid  | Pospiviroidae |
| Tomato planta macho        | TPMVd   | K00817 | 360              | Pospiviroid  | Pospiviroidae |
| Citrus exocortis           | CEVd    | M34917 | 370-375, 463     | Pospiviroid  | Pospiviroidae |
| Chrysanthemum stunt        | CSVd    | V01107 | 354, 356         | Pospiviroid  | Pospiviroidae |
| Tomato apical stunt        | TASVd   | K00818 | 360, 363         | Pospiviroid  | Pospiviroidae |
| Iresine 1                  | IrVd-1  | X95734 | 370              | Pospiviroid  | Pospiviroidae |
| Columnea latent            | CLVd    | X15663 | 370, 372         | Pospiviroid  | Pospiviroidae |
| Hop stunt                  | HSVd    | X00009 | 295-303          | Hostuviroid  | Pospiviroidae |
| Coconut cadang-cadang      | CCCVd   | J02049 | 246-247, 287-301 | Cocadviroid  | Pospiviroidae |
| Coconut tinangaja          | CTiVd   | M20731 | 254              | Cocadviroid  | Pospiviroidae |
| Hop latent                 | HLVd    | X07397 | 256              | Cocadviroid  | Pospiviroidae |
| Citrus IV                  | CVd-IV  | X14638 | 284              | Cocadviroid  | Pospiviroidae |
| Apple scar skin            | ASSVd   | M36646 | 329-330          | Apsacaviroid | Pospiviroidae |
| Citrus III                 | CVd-III | S76454 | 294, 297         | Apsacaviroid | Pospiviroidae |
| Apple dimple fruit         | ADFVd   | X99487 | 306-307          | Apsacaviroid | Pospiviroidae |
| Grapevine yellow speckle 1 | GVYSd-1 | X06904 | 366-368          | Apsacaviroid | Pospiviroidae |
| Grapevine yellow speckle 2 | GVYSd-2 | J04348 | 363              | Apsacaviroid | Pospiviroidae |
| Citrus bent leaf           | CBLVd   | M74065 | 318              | Apsacaviroid | Pospiviroidae |
| Pear blister canker        | PBCVd   | S46812 | 315-316          | Apsacaviroid | Pospiviroidae |
| Australian grapevine       | AGVd    | X17101 | 369              | Apsacaviroid | Pospiviroidae |
| Coleus blumei 1            | CbVd-1  | X52960 | 248, 250-251     | Coleviroid   | Pospiviroidae |
| Coleus blumei 2            | CbVd-2  | X95365 | 301-302          | Coleviroid   | Pospiviroidae |
| Coleus blumei 3            | CbVd-3  | X95364 | 361-362, 364     | Coleviroid   | Pospiviroidae |
| Avocado sunblotch          | ASBVd   | J02020 | 246-250          | Avsunviroid  | Avsunviroidae |

|                                   |        |        |             |                               |
|-----------------------------------|--------|--------|-------------|-------------------------------|
| Peach latent<br>mosaic            | PLMVd  | M83545 | 335-<br>338 | Pelamoviroid<br>Avsunviroidae |
| Chrysanthemum<br>chlorotic mottle | CChMVd | Y14700 | 398-<br>401 | Pelamoviroid<br>Avsunviroidae |

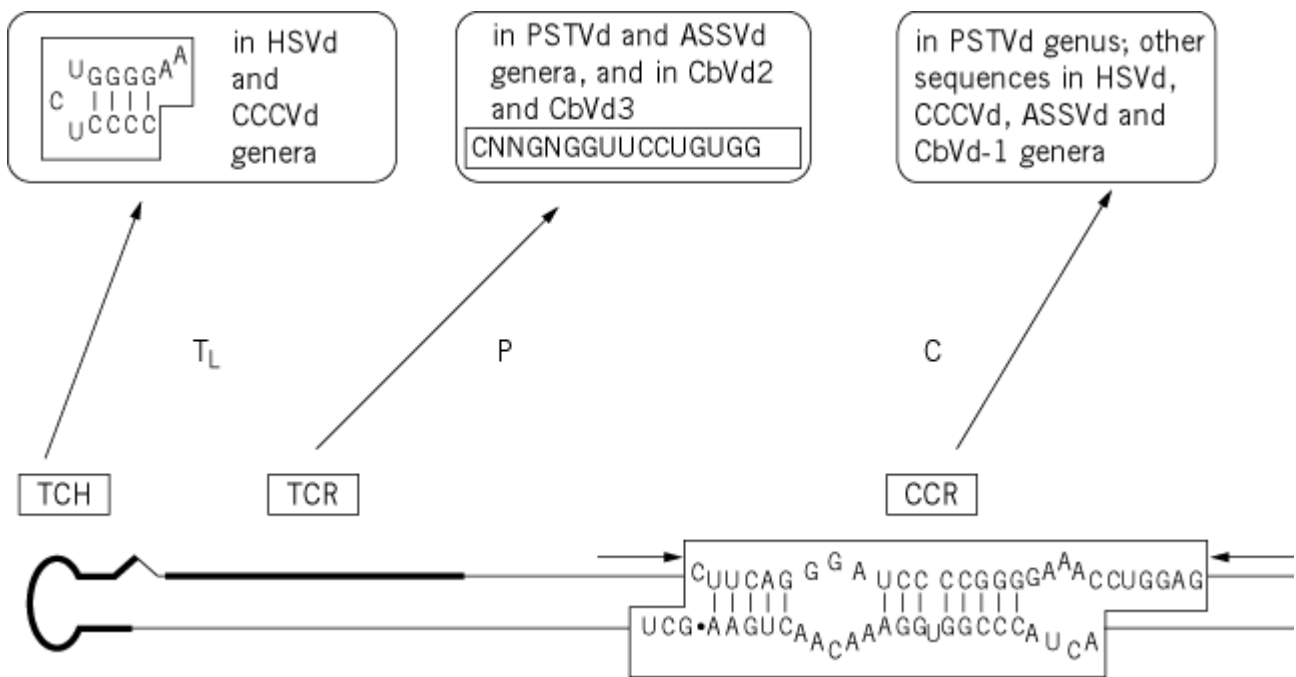
---

Viroids must be replicated by pre-existing host **RNA polymerases**, which are activated, but not coded, by the viroid genome. The viroid RNA itself, or one of its replicative intermediates, must also trigger the pathogenic effects on the host by direct interaction with one or more cellular targets. From the viroid names (Table 1), it is clear that they cause a series of maladies affecting crops of economic importance. Some viroids can also replicate without eliciting any symptom in their hosts. Thus far, no viroid has been found to infect bacterial, fungal or animal cells, although the RNA of human hepatitis delta virus shares some properties with viroid and viroid-like **satellite RNA** (see [Virusoids](#)).

### 1. Conserved Structural Domains and Classification

Viroid RNAs display a high degree of **secondary structure** as a result of extensive intramolecular self-complementarity (2). In most, but not all, cases, they adopt *in vitro* rod-like or quasi-rod-like conformations containing five structural domains that house the various functional determinants (3, 4). Viroids can also adopt other metastable conformations during their replication (5). The structure of viroids *in vivo* is not known, although there are indications of at least a partial rod-like structure in some cases. Within the rod-like structure, several conserved motifs can be distinguished: 1) the central conserved region (CCR) located in the central domain and comprised of two segments of nucleotides flanked by an imperfect inverted repeat in the upper strand, 2) the terminal conserved region (TCR) found in the upper strand of the left terminal domain, and 3) the terminal conserved hairpin (TCH) also located at the end of the left terminal domain and conserved in secondary structure as well as sequence (Fig. 1). The type of CCR as well as the presence or absence of the TRC and TCH motifs are useful criteria for classifying most viroids in different genera within family *Pospiviroidae* (Table 1) (6). Within a second family, *Avsunviroidae*, comprising ASBVd, PLMVd, and CChMVd (6), the situation is very different; they do not have these conserved motifs, but their strands of both polarities are able to adopt hammerhead structures (7) (see [Ribozyme/Catalytic RNA](#)). Moreover, PLMVd and CChMVd adopt *in vitro*, and probably *in vivo*, branched rather than rod-like conformations (7).

**Figure 1.** Schematic rod-like structure model for viroids of the family **Pospiviroidae**. The approximate location of the fi (pathogenic), V (variable) and TL (terminal left) and T<sub>R</sub> (terminal right) are indicated on top of the figure. The core nucleic acid region (CCR), terminal conserved region (TCR) and terminal conserved hairpin (TCH) are shown, as well as the presence of different members of the family. **Arrows** indicate flanking sequences that form, together with the core nucleotides of the inverted repeats.



## 2. Functional and Evolutionary Aspects

Differences between the two viroid families are not restricted to structure. The subcellular replication and accumulation sites of PSTVd and ASBVd, the type members of the two families, are the [nucleus](#) and the [chloroplast](#), respectively. The available evidence indicates that other members of both families behave in this respect as their type species. Viroid replication has been proposed to occur through RNA intermediates (8) and by means of a **rolling circle replication mechanism** (9). In both families, the circular monomeric plus-strand RNA is transcribed into linear multimeric minus strands as a first step. PSTVd and members of the first family subsequently follow an asymmetric pathway in which the multimeric minus strands generate linear multimeric plus strands in a second round of RNA-RNA transcription, which are then processed to unit-length and ligated. ASBVd and other members of the second family follow a symmetric pathway, in which the linear multimeric minus strands of the first step are cleaved and ligated to the monomeric, circular minus RNA strand, which then acts as the initial template for the second half of the cycle. On the basis of sensitivity to **a-amanitin**, nuclear RNA polymerase II is probably involved in the polymerization of the RNA strands of PSTVd and related viroids of family *Pospiviroidae* (10). Additional studies using a monoclonal antibody against a conserved domain in the largest subunit of RNA polymerase II support this view (11). In contrast, an a-amanitin-resistant chloroplastic RNA polymerase appears to catalyze RNA elongation in ASBVd starting from definite sites in both polarity strands (12). One remarkable aspect of the last three viroids is that they self-cleave *in vitro* and almost certainly *in vivo* also, through hammerhead ribozymes to produce unit-length strands. Conversely, members of the PSTVd family may depend on host **ribonucleases** for cleavage, although self-cleavage through a different class of ribozyme structures cannot be excluded. The final circularization step in both families is probably mediated by host [RNA ligases](#) although, at least in the case of PLMVd, this step has been proposed to occur autocatalytically.

An evolutionary origin totally independent from that of RNA viruses has been proposed for viroid and viroid-like satellite RNAs (see [Virusoids](#)). Their circular structure and small size, and especially the frequent occurrence of ribozyme activities in these RNAs, all suggest that they may be molecular fossils of the [RNA world](#) that possibly preceded on Earth the appearance of cellular life based on DNA and proteins. Mutation and recombination events between viroids coinfecting the same host plant would seem to have contributed to their divergent evolution.

There are still fundamental properties of the viroids to be identified: 1) the molecular determinants targeting some viroids to the nucleus and others to the chloroplast; 2) the initiation sites of RNA synthesis in other viroids apart from ASBVd and the nature of promoters directing this synthesis; 3) whether additional RNA polymerases and transcription factors are involved in viroid replication; 4) the exact *in vivo* processing site of members of the PSTVd family; 5) the primary interaction between the viroid RNA, or a replicative intermediate thereof, and a cellular target leading through a signal transduction pathway to the onset of symptoms; and 6) the molecular basis of the interference observed between related viroids co-infecting the same plant.

## Bibliography

1. T. O. Diener (1971) *Virology* **45**, 411–428.
2. H. J. Gross, H. Domdey, and C. Lossow, et al. (1978) *Nature* **273**, 203–208.
3. P. Keese and R. H. Symons (1985) *Proc. Natl. Acad. Sci. USA* **82**, 4582–4586.
4. T. Sano, T. Candresse, R. W. Hammond, T. O. Diener, and R. A. Owens (1992) *Proc. Natl. Acad. Sci. USA* **89**, 10104–10108.
5. F. Qu, C. Heinrich, P. Loss, G. Steger, P. Tien, and D. Riesner (1995) *EMBO J.* **12**, 2129–2139.
6. R. Flores, J. W. Randles, M. Bar-Joseph, and T. O. Diener (2000) In *Virus Taxonomy, Seventh Report of the International Committee on Taxonomy of Viruses*, (M. H. V. Van Regenmortel, C. M. Fauquet, and D. H. L. Bishop, et al., eds.), Academic Press, San Diego, pp. 1009–1024.
7. B. Navarro and R. Flores (1997) *Proc. Natl. Acad. Sci. USA* **94**, 11262–11267.
8. L. K. Grill and J. S. Semancik (1978) *Proc. Natl. Acad. Sci. USA* **75**, 896–900.
9. A. D. Branch and H. D. Robertson (1984) *Science* **223**, 450–454.
10. D. Warrilow and R. H. Symons (1999) *Arch. Virol.* **144** 2367–2375.
11. J. A. Navarro and R. Flores (2000) *EMBO J.* **19**, 2662–2670.

## Suggestions for Further Reading

12. T. O. Diener (1979) *Viroids and Viroid Diseases*, Wiley, New York, NY.
13. T. O. Diener, ed. (1987) *The Viroids (The Viruses)*, Plenum Press, New York.
14. T. O. Diener (1999) Viroids and the nature of viroid diseases. *Arch. Virol. Suppl.* **15**, 203–220.
15. R. Flores, F. Di Serio, and C. Hernández (1997) Viroids: The non-coding genomes. *Semin. Virol.* **8**, 65–73.
16. R. Flores, J. A. Daròs, and C. Hernández (2000) The Avsunviroidae family: Viroids with hammerhead ribozymes. *Adv. Virus. Res.* **55**, 271–323.
17. R. H. Symons (1997) Small catalytic RNAs. *Annu. Rev. Biochem.* **61**, 641–671.
18. J. S. Semancik, ed. (1987) *Viroids and Viroidlike Pathogens*, CRC Press, Boca Raton, FL.

## Virus Infection, Animal

### 1. Definitions

#### 1.1. Viruses

Supramolecular structures made of nucleic acid (the viral genome), proteins, glycosylated or not, and, eventually (in the case of enveloped viruses), lipids. Viruses are obligatory parasites of cells. They use the cell machinery of protein biosynthesis to make enzymes and regulatory and structural

proteins. In some cases, they also use cellular polymerases to synthesize their nucleic acids. The lipid-bilayer, the coat of enveloped viruses, is of cellular origin.

## 1.2. Infection

Invasion of a cell or an organism by an agent that, most of the time, will propagate, kill the host cell, and eventually cause disease.

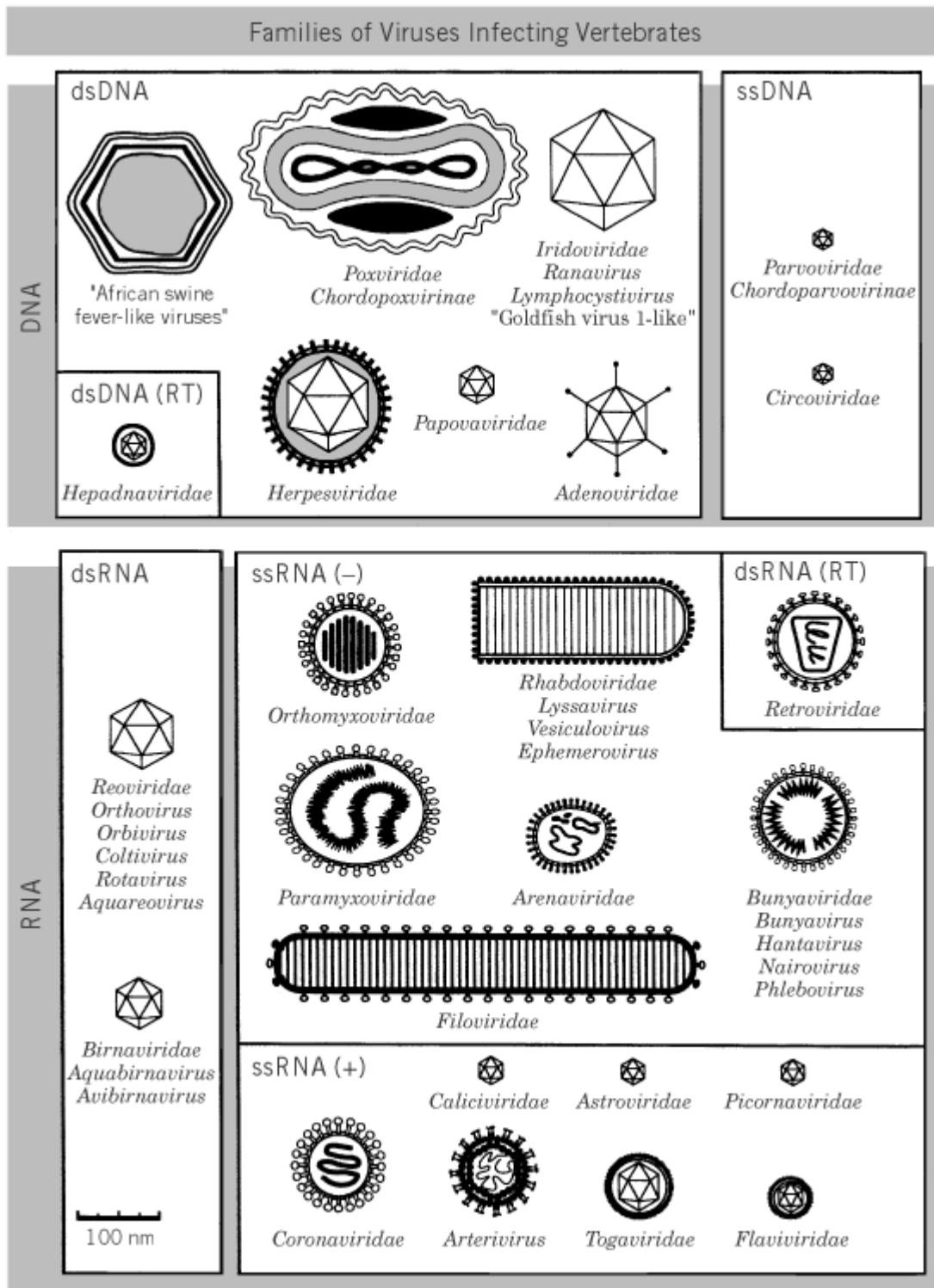
## 2. Classification of Viruses

The International Committee on Taxonomy of Viruses (1) (ICTV) has proposed adopting a universal system of virus taxonomy. The system uses hierarchical levels such as order (suffix -virales), family (suffix -viridae), subfamily (suffix -virinae), genus (suffix -virus) and species ("a virus species is defined as a polythetic class of viruses that constitutes a replicating lineage and occupies a particular ecological niche"). It should be recognized that great difficulties are being encountered in deciding whether a particular virus should be designated as a species or as a subspecies or strain or variant. The properties of viruses used in taxonomy are properties of the virion (morphology, physicochemical, and physical properties) and of the viral proteins, carbohydrates, lipids, and antigenic and biologic properties. The ICTV was recognizing more than 3,600 virus species and hypothesizing that more than 30,000 viruses, virus strains, and subtypes were being tracked in reference centers and culture collections (1). The ultimate goal is to catalogue "data down to subspecies, strain, variant, and isolate levels, that is, levels important in medicine, agriculture, and other scholarly fields." Obviously, the task is gigantic. Its achievement will require major investments in biochemical and biophysical techniques, as well as in bioinformatics. The 1995 state of the art has been updated and amplified (2).

This classification takes no account of the diseases caused by viruses. Indeed, it is more and more often observed that many viruses cause no, or only occasional, disease. They have reached a state of harmony with their host; an example is the case of the African green monkey (AGM) strain of simian immunodeficiency virus (SIV), which is very frequently found in the wild with no induction of disease.

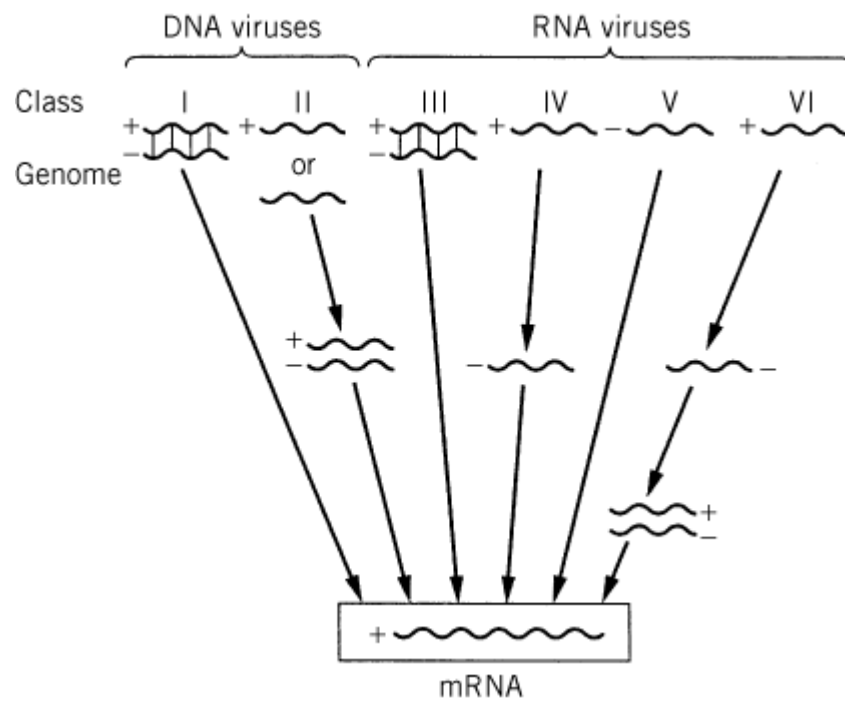
Families of viruses infecting invertebrates or vertebrates are illustrated on Figure 1 as virus diagrams. Taxa containing double-stranded (ds) or single-stranded (ss) DNA or RNA (either the + or – strands) are separated by frames. Nonenveloped viruses are depicted as icosahedra. A dotted or hairy outer layer identifies enveloped viruses. A simple classification of viruses, based on how viral messenger RNA is produced, was proposed by D. Baltimore (3) (Fig. 2).

**Figure 1.** Families of viruses infecting invertebrates and vertebrates. Frames separate taxa containing double-stranded (ds) and single-stranded (ss) genomes. Horizontal blocks separate taxa containing DNA and RNA viruses. Taxa containing reverse transcribing (RT) viruses and the negative (–) and positive (+) ss RNA genomes are also indicated. A **dotted or hairy outer layer** indicates an enveloped virus. Icosahedral structures designate nonenveloped viruses. All diagrams were drawn approximately to the same scale (From Ref. 2, with permission).

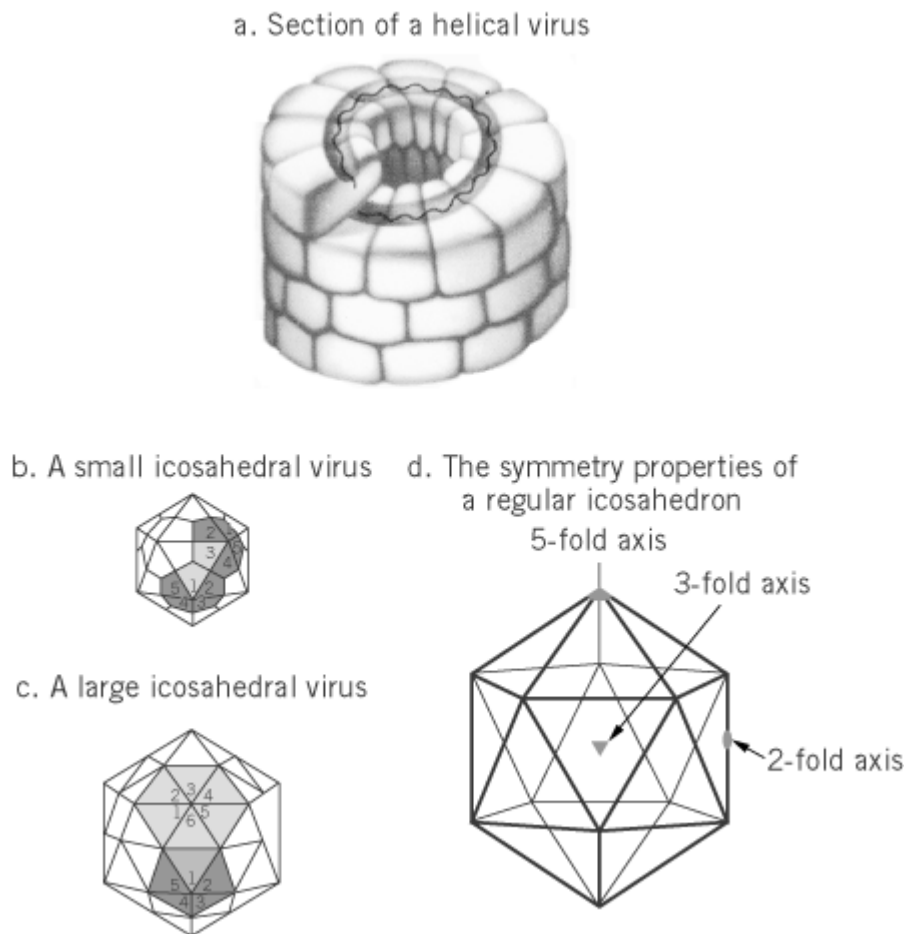


**Figure 2.** Baltimore's simplified classification of viruses. This classification is based on the type of nucleic acid in the genome and the mode of production of mRNA (3). Class I: Viruses with a ds DNA genome; mRNA is transcribed by the classical mechanism of transcription of a ds DNA template. Class II: Viruses with a ss DNA genome; the DNA strand has the same, positive sense as the mRNA, or else the sense of the genomic strand is negative, opposite to the mRNA. Synthesis of viral mRNA does not proceed before the ss DNA genomes of class II viruses are converted to a ds molecule. Class III: Viruses with ds RNA genome; such genomes are segmented. Transcription of mRNA is catalyzed by an RNA-dependent RNA polymerase packaged in the virion. Class IV: Viruses with an ss RNA genome of the same polarity as mRNA; an ss RNA (-) template is made before mRNA synthesis. Class V: Viruses with an ss RNA (-) genome complementary to the mRNA. Virion RNA is transcribed into mRNA by virion-associated enzymes. Class VI: Retroviruses have a diploid RNA genome made of two molecules of RNA (+). The + sense RNA is used as

template by reverse transcriptase for the synthesis of a ds provirus to be integrated in the host DNA; transcription is catalyzed by a cellular DNA-dependent RNA polymerase (From Ref. 3, with permission).



**Figure 3.** Structure of viruses. **(a)** A helical virus. The **wavy red line** represents RNA protected within successive arrays of proteins. **(b, c)** Spherical viruses. Icosahedral capsids have vertices made of symmetrically organized pentons. Hexons make the faces and are arranged along the edges. (From Ref. 19, with permission). **(d)** The fivefold, threefold and twofold symmetry axes of a regular icosahedron (From Ref. 20, with permission).



### 3. Structure of Virions

Figure 1 presents simple diagrams of the structure of virions. Two capsid structures exist; they are either multimers of protein subunits arranged as icosahedra symmetry (quasi-spherical volumes with an external molecular sheet built of 20 identical faces, each of which is an equilateral triangle) or multiple copies of protein subunits arranged as helices. Icosahedral structures can form independently of the presence of nucleic acid, whereas the helical arrangements generally require the presence of nucleic acid to assemble.

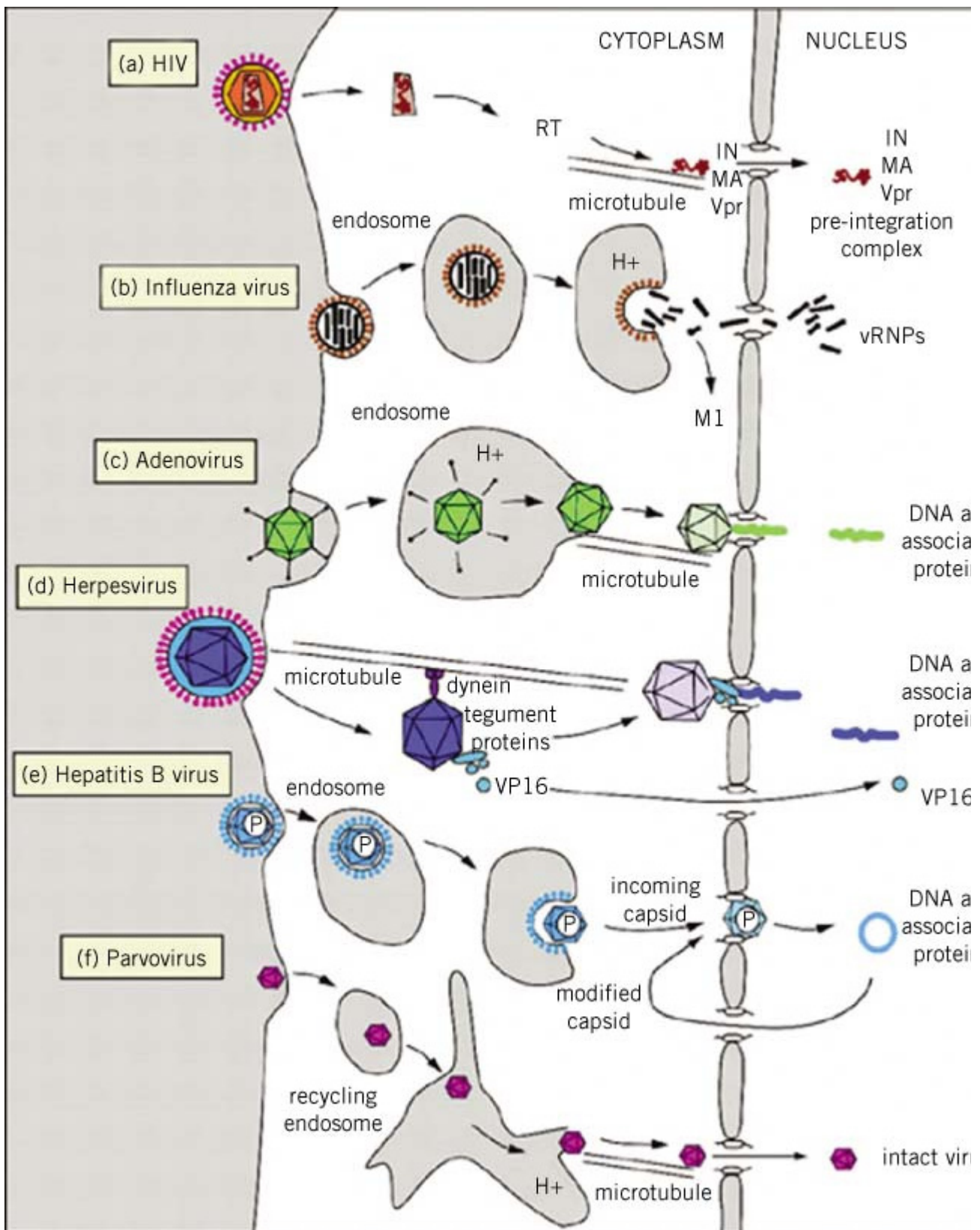
The detailed structure of picornaviridae capsids has been determined by X-ray crystallography. Models of other virus capsids at 28 to 35 Å resolution were provided using cryo-electron microscopy and computer single-particle reconstruction methods (rotaviridae, orbiviridae).

Models of helical and icosahedral viruses are illustrated in Fig. 3. Helical viruses largely expose protein subunits, and nucleic acid is protected within (Fig. 3a). Figure 3b, 3c, and 3d illustrate icosahedral structures with 20 triangular faces (each being made of three identical or equivalent capsid subunits), 12 vertices (where five subunits make contact symmetrically) and 30 edges (along which the protein subunits exhibit 6-fold symmetry, as well as the subunits making up the surfaces in between). Although the shape of the 20 faces is not a flat triangle as illustrated in Figure 3 (b and c), the overall appearance of these solids is that of nearly spherical volumes with triangular faces. At each vertex, edge, or face, symmetry is  $n = 5, 2,$  and  $3,$  respectively (Fig. 3d).

**Figure 4.** Virus import into the nucleus. The pathway from the plasma membrane (PM) to the nuclear membrane is illust



for four families of viruses. Nucleic acid-protein complexes enter the nucleus through nuclear pore complexes. The HIV capsid is released in the cytoplasm and moves to the nucleus. On the way, the capsid protein is disorganized, reverse transcription occurs, and a complex containing the matrix protein, integrase, Vpr, and the nucleocapsid protein (probably cleaved in the nucleus) enters the nucleus via the nuclear pore complex. Influenza virus (an enveloped virus) and adenovirus (a nonenveloped virus) follow similar pathways. They are internalized via the endosomal pathway, and they rupture the endosomal membrane via protein conformational changes induced by the low pH of the late endosome. Ribonucleoproteins (for influenza) and deoxyribonucleoproteins (for adenovirus) enter the nucleus via the nuclear pore. The envelope of herpesviruses fuses with the plasma membrane. The capsid travels along microtubules, moved by dynein. At the nuclear pore, the nucleic acid and its associated proteins are released into the nucleus, whereas the empty capsid remains in the cytosol. Hepatitis B virus may possibly disassemble within the nuclear pore while parvoviruses can enter the nucleus essentially intact (From Ref. with permission). IN = integrase; INM = inner nuclear membrane; MA = matrix; ONM = outer nuclear membrane; RT = reverse transcriptase; Vpr = viral protein R.



The capsid structure of retroviruses corresponds to the above description. It contains an external matrix protein (a key factor in nuclear targeting) underlying the lipid bilayer and an internal capsid protein, limiting the internal cavity, which contains two identical (+) viral RNA molecules associated with the nucleocapsid protein, the reverse transcriptase, and other enzymes of cellular origin. Originally, retroviruses were classified according to their morphological appearance and budding

phenotype as observed by electron microscopy (4) (see virus assembly below).

#### 4. Virus Attachment, Entry, Uncoating, and Transport

Viruses enter cells by interaction with their receptor; human immunodeficiency virus (HIV) and simian immunodeficiency virus (SIV) use two receptors (it is often spoken of as one receptor, CD4, and one coreceptor, which belongs to the family of receptors for chemokines).

Adenoviridae also use two receptors. The presence of the receptor to a given virus does not *ipso facto* mean that the cell will be permissive to that virus. Nonreceptor genetic controls of early events in infection have been amply documented in mice, chickens, hamsters, and so on. Receptors can be proteins, sugars, polysaccharides, ceramides, and so on, present at the cell surface. Virus-receptor attachment rates have been measured on a variety of systems. Some happen to be very fast, with a second-order rate constants of about  $10^{-7} \text{ cm}^3 \text{ min}^{-1}$ , others are much slower,  $10^{-12} \text{ cm}^3 \text{ min}^{-1}$  (from Ref. 7, p. 629).

Two major mechanisms of virus entry have been recognized: virus-cell fusion at the plasma membrane, which is pH-independent, and virus endocytosis, which is pH-dependent. Passage of viral structures through a lipid bilayer probably implies extrusion of hydrophobic peptide sequences, as suggested for poliovirus (5). In the case of fusogenic enveloped viruses, the fusogenic peptide at the amino-terminus of the F protein inserts obliquely in the lipid bilayer of the target cell and destabilizes it (6). Virus-cell fusion occurs at the plasma membrane for pH-independent viruses (eg, retroviruses) or in endosomes for pH-dependent viruses (eg, influenza virus) (7).

Enveloped viruses recognize their target cell via interaction of their external envelope protein with a cell membrane protein acting as receptor. The cases of HIV and SIV are more complex because the virus requires both a receptor and a coreceptor for entry. The most commonly used receptor is CD4; the coreceptor is a chemokine receptor. Induction of a conformational change by interaction of HIV-1 external glycoprotein (gp120) with the CD4 receptor was demonstrated by the crystal structure at 2.5 Å resolution of an HIV-1 gp120 core complexed with a two-domain fragment of human CD4 and an antigen-binding fragment of a neutralizing antibody that blocks chemokine-receptor binding (8). The structure revealed a cavity-laden CD4-gp120 interface and a conserved binding site for the chemokine receptor.

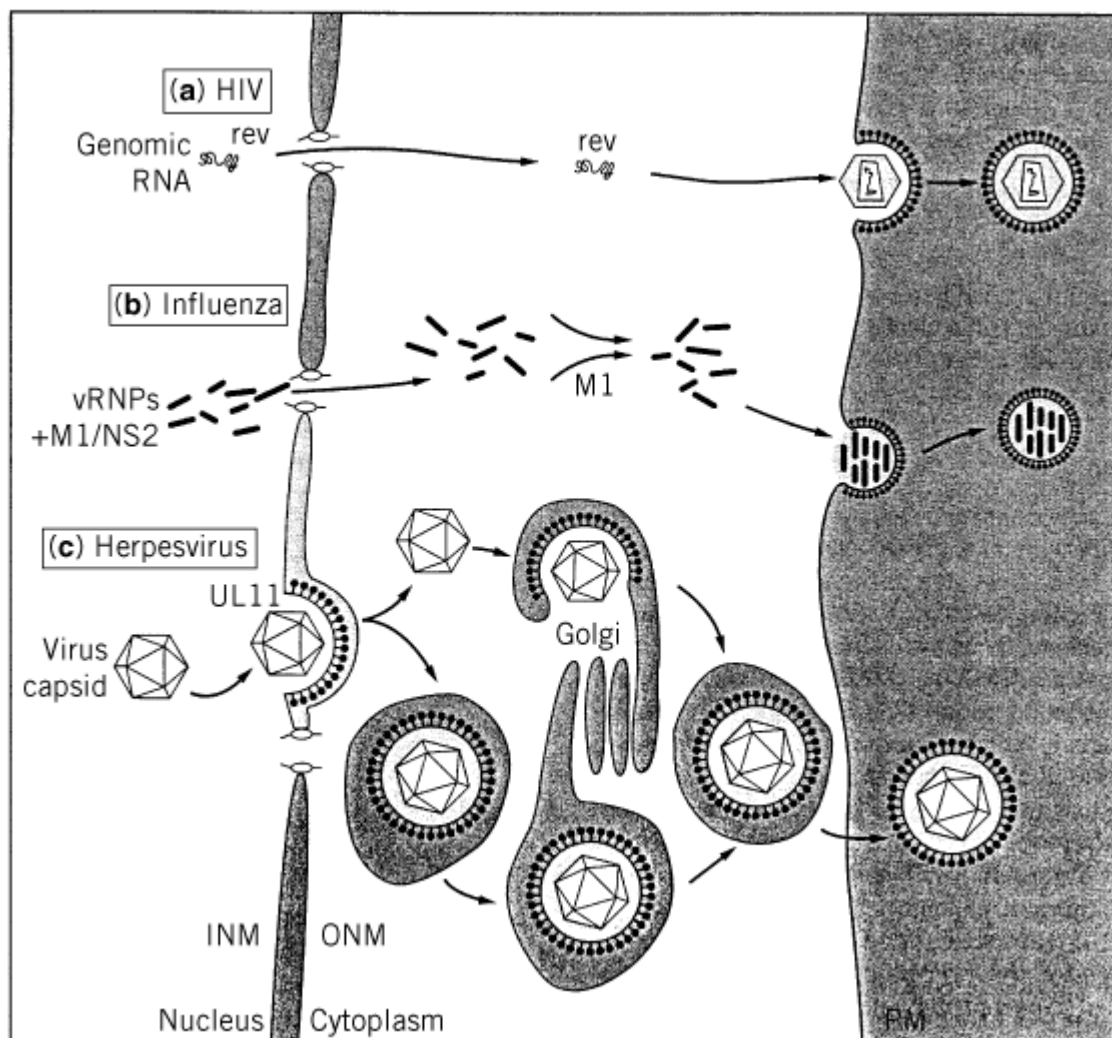
Virus-cell fusion implies a major contribution by the viral transmembrane moiety of the envelope glycoprotein. In pH-dependent viruses (ie, influenza virus), each subunit of the trimeric envelope gets rearranged and reorganized largely as an a helix, the fusion peptide (located at the NH<sub>2</sub> end of the transmembrane protein, HA2 in this case) being cast some 100 Å outwards towards the target membrane (9). X-ray diffraction studies applied to HIV-1 transmembrane protein gp41 and to Moloney murine leukemia virus illustrate the great similarity between pH-dependent and pH-independent envelope transmembrane protein architectures, suggesting a common mechanism for initiating fusion (9-11).

Virion interaction with the virus receptor(s) leads to release in the cytoplasm of the viral capsid or of the viral nucleic acid. Many aspects of uncoating are still undefined, even in thoroughly studied systems such as poliovirus or HIV. In poliovirus, preinfection maturation of the particle involves loss of protein VP4 and swelling of the particle, and, at very low frequency of a few percent, the stage of an infectiousome is supposed to deliver the RNA genome through the membrane and into the cytosol. There is a receptor-mediated extrusion of a hydrophobic structure into the host cell membrane. In this system, the uncoating steps occur outside the cell, with the particle being attached to the receptor (a member of the immunoglobulin family) and in close contact with the cell membrane (Ref. 7, p. 633).

Virus nuclear import (12, 13) is illustrated in Figure 4 and summarized in Table 1. Engulfed in an

endosome, adenovirus capsid external protein causes disruption of the endosomal membrane. Freed in the cytosol, the capsid interacts with the nuclear pore complex through a penton structure at one vertex. Viral DNA and associated proteins penetrate the nucleus, leaving behind the empty capsid. Similar mechanisms mediate the transport of herpesvirus capsids to the nucleus. This enveloped virus fuses at the plasma membrane; the capsid moves to the nucleus along microtubules. The capsid-nuclear pore interaction involves one vertex. DNA and associated proteins penetrate the nucleus.

**Figure 5.** Export of viral material from the nucleus to the cytoplasm and the extracellular compartment. The HIV Rev protein is required for export from the nucleus of full-size and singly spliced viral RNA. An immature capsid gets progressively organized underneath the lipid bilayer of the plasma membrane where viral envelope glycoproteins are concentrated. An analogous process leads to budding of influenza virus particles. Fully assembled herpes virus capsids bud from the inner nuclear membrane and fuse with the outer nuclear membrane, being released naked into the cytoplasm. Naked particles become enwrapped by membranes of the Golgi apparatus and accumulate within a vesicle. Alternatively, the capsid buds through the double membrane of the nucleus, and the outer membrane fuses with a membrane of the Golgi apparatus. The final outcome of this pathway is again an enveloped virus within a cytoplasmic vesicle. Enveloped particles will be released from the cell through fusion of the vesicular membrane with the plasma membrane of the infected cell (From Ref. 14, with permission).



**Table 1. Summary of Viruses That Enter the Nucleus**

| Method of Nuclear import                             | Virus Family                | Selected Example(s)                                     |                   |
|--|-----------------------------|---|-------------------|
| Entry during mitosis                                 | Retroviridae:<br>Oncovirus  | Simian retrovirus type 1                                | RNA<br>(RT)       |
| Uncoating in the cytoplasm + NPC translocation       | Retroviridae:<br>Lentivirus | Human immunodeficiency virus type 1 (HIV-1)             | RNA<br>(RT)       |
|  | Orthomyxoviridae            | Influenza virus   | ss<br>RNA(-<br>)  |
| Uncoating in the cytoplasm                           | Bornaviridae                | Borna disease virus                                     | ss<br>RNA(-<br>)  |
| Docking of virus to the NPC/transit of intact virus? | Hepadnaviridae              | Hepatitis B virus (HBV)                                 | ds<br>DNA<br>(RT) |
| Transit of intact virus                              | Parvoviridae                | Minute virus of mice<br>Parvovirus B19                  | ss<br>DNA         |
|  | Papovaviridae               | Papovavirus Simian virus 40                             | ds<br>DNA         |
| Docking of intact virus to the NPC                   | Adenoviridae                | Adenovirus 2  | ds<br>DNA         |
|  | Herpesviridae               | Herpes simplex virus (HSV)<br>Cytomegalovirus (CMV)     | ds<br>DNA         |
|  | Baculoviridae               | Autographa californica multiple nuclear polyposis virus | ds<br>DNA         |

RT, reverse transcriptase; NPC, nuclear pore complex.

Influenza viruses enter cells by receptor-mediated endocytosis. Uncoating of the virus in endosomes depends on the acidic pH of this compartment. It is supposed that the low pH-activated ion-channel activity of the virion-associated M2 protein permits the flow of ions (essentially protons) inside the virus particle, disrupting of protein-protein interactions and, thus, freeing the ribonucleoproteins (NPs). The low pH also induces a major change in the three-dimensional structure of the transmembrane protein HA2 (see above). The viral envelope fuses with the membrane of the endosome, releasing NPs in the cytosol from where they will move to the nucleus.

Uncoating and transport of retroviral capsids delivered in the cytosol after virus-cell membrane fusion are poorly understood. Reverse transcription, initiated in the producing cell, proceeds during the transport to the nucleus, where integration of the proviral DNA will take place, catalyzed by the viral integrase. Onco-retroviruses depend on the breakdown of the nuclear envelope during mitosis to achieve DNA integration and completion of their life cycle. On the contrary, lentiviruses use the nuclear localization signal of the matrix protein to establish infection in macrophages and quiescent T-lymphocytes (15). In addition to the matrix protein, reverse transcriptase and integrase, the nucleocapsid protein also plays a role by potentiating provirus integration in the presence of  $Mg^{2+}$ .

## 5. Virus Replication

The strategies of viral nucleic acid replication depend on the nature and strandedness of the viral genetic material.

### 5.1. Double-Stranded DNA Viruses

These viruses, Picornaviridae, replicate in the nucleus. DNA duplication can be carried out by either cellular or viral DNA polymerases. Poxviridae are exceptions. Their multiplication is cytoplasmic. The pox genome encodes the viral DNA polymerase. In single-stranded DNA viruses (Parvoviridae), DNA synthesis takes place in the nucleus, uses cellular DNA polymerases, and derives from a self-priming mechanism. The replicative intermediate is a linear duplex molecule. Virions carry either a (+) or a (–) DNA strand. Hepadnaviridae have a peculiar genomic structure. One strand of DNA (negative sense, complementary to the viral messenger RNAs) is full-length; the other is of variable size. Before transcription by the host RNA polymerase II, the DNA is converted into a covalently closed circle. The closed-circular DNA species can serve as template for the synthesis of a full size mRNA that will serve as a template for synthesis of the minus-strand DNA by reverse transcription using a protein primer.

### 5.2. Double-Stranded RNA Viruses

These viruses, Reoviridae, carry a genetic material made of ds RNA segments (10 to 16 segments, according to the genus) and a virion-associated RNA-dependent RNA polymerase (or transcriptase). Transcription products of the double-stranded genome are RNA (+) molecules that are either used as mRNA or encapsulated in nascent virions, where they are used as templates to generate a double-stranded genome.

### 5.3. Single-Stranded RNA (–) Viruses

These viruses, Rhabdoviridae, harbor a negative-sense RNA molecule (11 kb to 15 kb in size) and a virion-associated transcriptase. Replication takes place in the cytoplasm. Host protein factors are required for replication. Replicative intermediates are made of RNA (+) molecules linked to a so-called N protein. In turn, RNA (+) molecules are used as templates for new genomic (–) strands. Segmented negative-strand RNA viruses (ie, orthomyxoviridae [influenza viruses]) have a genome size of 10 kb to 13.6 kb, made of six or seven segments, depending on the genus. At infection, viral nucleocapsids are transported to the cell nucleus, where the virion transcriptase synthesizes mRNA species, which act as templates for new viral RNA synthesis. The latter RNA (–) molecules exist as nucleocapsids in the nucleus of infected cells.

### 5.4. Single-Stranded RNA (+) Viruses

Virions of single-stranded (+) RNA virus (7 kb to 8.5 kb in size for picornaviruses) contain one molecule of infectious, positive-sense ss RNA. Translation of the genomic RNA generates a polyprotein, which is processed by a viral protease, giving rise to structural virus proteins and nonstructural ones, among which are one subunit of replicase (RNA-dependent RNA polymerase). The replicase probably associates with cellular proteins to make up a functional RNA polymerase. This situation is analogous to that of the Q $\beta$  phage, in which the virus-encoded replicase subunit associates with elongation factors Tu and Ts and ribosomal protein S1, to build up a highly active replicase to which the virus contributed minimally. Replication first produces a negative-sense RNA, complementary to the RNA (+) genomic RNA. Due to the higher affinity of replicase for the 3' end sequence of minus strands, many more RNA (+) molecules than RNA (–) molecules are produced, which are either translated or encapsidated in progeny virus.

Retroviridae have a diploid, single-stranded RNA (+) genome and a virion-associated reverse transcriptase (RNA-dependent DNA polymerase). Reverse transcription starts in the virus particle and produces a linear ds DNA copy of the genomic RNA. The linear DNA is integrated into the host cell DNA via the activity of a virus-coded integrase and is transcribed as regular cellular genes. Viruses of that family are oncogenic through their site of integration (activation of proto-oncogenes); other

members of the family have incorporated and mutated proto-oncogenes and, consequently, become acute transforming viruses; others (such as human T-cell leukemia virus) are transforming entities through dysregulation of the cell transcription machinery. The lentivirinae are primarily inducers of immunodeficiency.

## 6. Synthesis of Viral Messenger RNA and Viral Proteins

### 6.1. Double-stranded DNA viruses (ie, Papovaviridae)

The life cycle of many viruses can be divided into an early and a late period. In ds DNA viruses, proteins made in the first period are called early; they play regulatory functions. Proteins made later in the life cycle of the virus are structural molecules. Precursor mRNAs undergo classical post-transcriptional processing, including 5'-capping, polyadenylation, and RNA splicing. Efficient use of the available information involves alternative splicing and use of overlapping open reading frames. A peculiarity seen in Poxviridae, adenoviridae, and some families of parvoviridae is the transcription of early proteins from both strands of the ds genome; this is another strategy to use maximally the available genomic sequences.

### 6.2. Double-Stranded RNA Viruses

Viruses of this family carry a virion-associated transcriptase. Transcription of each segment of the ds segmented genome generates an mRNA coding for a single polypeptide chain, a polyprotein.

### 6.3. Single-Stranded RNA (–) Viruses

The genes of these viruses are transcribed processively, from the 3' to the 5' end of the template virus RNA and in decreasing molar abundance. The mRNAs are 5'-capped, 3'-polyadenylated, and generally monocistronic. In the case of viruses with a segmented genome, the largest genome segments encode one protein each, whereas some of the smaller segments code for additional proteins from spliced or bicistronic mRNAs. Influenza virus mRNA synthesis occurs in the nucleus. It requires initiation by host-cell primers and specifically capped RNA fragments derived from host-cell RNA polymerase II transcripts. In this system, viral mRNAs start with cellular sequences (10 to 13 nucleotides).

### 6.4. Single-Stranded RNA (+) Viruses

A frequent strategy found in this group of viruses is the translation of the + sense genomic RNA into a polyprotein, matured later on by a virus encoded proteinase (ie, picornaviridae, poliovirus). In alphaviridae (Semliki Forest Virus), the full-length genomic RNA codes for nonstructural proteins. The message for structural proteins is subgenomic and corresponds to the 3' region of the genomic RNA.

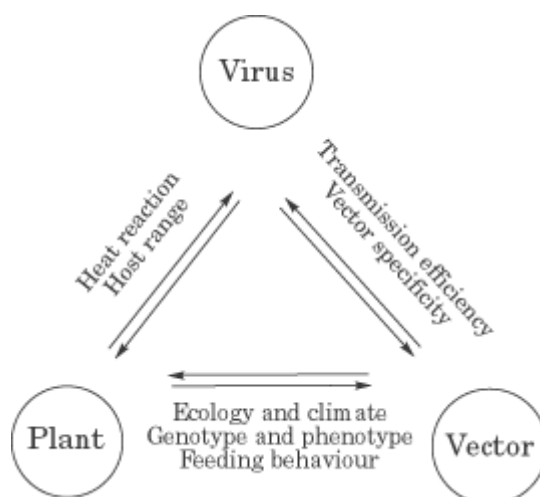
Transcription of retroviral messages occurs from the integrated proviral DNA. It combines all the known mechanisms of eukaryotic transcription. Expression of retroviral information includes expression of regulatory functions from spliced mRNAs, expression of internal structural proteins from genomic, full-length RNA, and expression of envelope proteins from a singly spliced mRNA. Nuclear export of the full-length and singly spliced viral RNA requires interaction with a viral regulatory protein called Rex in Human T-lymphotropic Virus (HTLV) and bovine leukemia virus (BLV). A lentiviral protein with analogous function is called Rev in HIV and SIV.

## **Virus Infection, Plant**

The natural infection of a plant by a virus involves complex interactions between three [genomes](#): that

of the virus, that of the host, and that of the vector (Fig. 1). Superimposed on these interactions are abiotic factors that affect the behavior of the host and vector. The application of molecular techniques is leading to a more detailed understanding of these interactions. However, only a small number of virus–host, virus–vector, and host–vector combinations have so far been studied, and it would be risky to extrapolate from these data to a wide range of viruses. Details of the structures and genome organizations of these viruses are given in [Viruses, plant](#).

**Figure 1.** Genomes that interact in the natural infection of a plant by a virus. Reproduced from (21) by kind permission of W. B. Saunders Co.



## 1. Host Range and Symptoms

Viruses differ in the number of plant species that they can infect systemically. Some viruses have wide host ranges; eg, *cucumber mosaic virus* has been reported to infect more than 800 species in 85 families (1). Other viruses (eg, *apple stem grooving virus*) have narrow host ranges, naturally infecting very few species. Some viruses are restricted to one plant organ, eg, *necroviruses* are usually restricted to the roots, *luteoviruses* are usually restricted to the phloem cells in the vascular bundle. Experimental infections show that some viruses can infect protoplasts of species that they do not infect systemically; this is termed *subliminal infection* and occurs because of the restriction of the cell-to-cell spread of the virus (see below).

Plants react to viruses with a wide range of symptoms, but the symptom type is usually specific to the virus and frequently used in the name of the virus (see Table 1 in [Viruses, Plant](#)). In some virus–host combinations, local symptoms, which are usually *chlorotic* (yellow) or *necrotic* spots, develop on inoculated leaves around the site of virus entry. If the spots become necrotic, the virus often does not spread systemically to the rest of the plant. The number of necrotic spots produced on a leaf after mechanical inoculation can be related to the virus concentration, and a local lesion assay is frequently used for the quantification of an infectious virus. The most common systemic symptom is mottle or mosaic patterns of light and dark green, and sometimes yellow, areas on the leaves. The distribution of light and dark green is usually delimited by the veins on the leaves, giving an irregular mosaic in dicotyledons and striping in monocotyledons. The mosaic expression in flowers is termed *color breaking*, the best-known example being tulip flower color breaking seen in many Dutch old master paintings. Other color symptoms on leaves include yellowing, often associated with phloem-limited viruses, and ringspots formed by concentric rings or irregular linepatterns of yellowed tissue, which may become necrotic. Viruses can also cause growth abnormalities, which include stunting of the plant, reduction of the leaf lamina, enations that are outgrowths from the leaf veins, and



tumorlike growths.

## 2. Plant-to-plant Spread

To infect a plant, a virus has to pass the cuticle and cellulose cell wall, a process that requires mechanical damage. In experimental conditions, many viruses that are not restricted to the phloem can be transmitted by rubbing sap from an infected plant, together with an abrasive, on a leaf of a host species. In nature, some viruses, eg, **tobacco mosaic tobamovirus** (TMV) and potato virus X potexvirus, spread from plant to plant via broken leaf hairs that result from leaves rubbing together. However, for the majority of viruses, natural transmission involves a biological *vector*. These include insects that feed by piercing and sucking (eg, aphids, leafhoppers), [nematodes](#), **fungi** (various lower fungi), and pollen. The interactions between a virus and its vector are very specific, and viruses that are transmitted by one type of insect, say, aphids, are not transmitted by other insect types or by fungi or nematodes. The specificity goes even further, and some viruses are transmitted by only one insect species.

There are two basic forms of virus–vector interactions: one in which the virus is externally borne and does not enter inside the vector and the other, internally borne, in which the virus enters the vector's body. The externally borne interaction is characterized by rapid transmission and short persistence of the virus associated with the vector. Internally borne viruses persist in the vector for much longer periods and, if the virus replicates in the vector, the vector can usually transmit for all its life. In each case, there is a specific interaction. Most is known about details of interactions between viruses and insect vectors.

Externally borne, insect-transmitted viruses interact with specific region(s) of the vectors' anterior food canal. This interaction can be either directly between the virus coat protein and the specific site, as exemplified by aphid transmission of cucumber mosaic cucumovirus (CMV), or involve a virus-encoded nonstructural protein, such as the helper component of potyviruses and the aphid transmission factor (P18 or gene II product) of **cauliflower mosaic caulimovirus** (CaMV). The potyviral helper component is a protein of 53 to 58 kDa (depending on the virus) that is processed proteolytically from the N-terminal part of the viral polyprotein (see Fig. 1 in [Viruses, Plant](#)) and is thought to be a heterodimer in its active form. There are several theories as to how these transmission proteins function (2), the most widely accepted being that they act as a bridge between the viral coat protein and mouthparts of the vector. The predicted **secondary structure** of CaMV P18 indicates that it has an N-terminal **domain**, which is primarily **b-sheet**, separated from the predominantly **a-helical** C-terminal domain by a random structure region. The C-terminal domain of CaMV P18 binds to the viral coat protein, and mutations of the N-terminal domain abolish insect transmissibility (3). The helper component of potyviruses is thought to function in a similar manner, and recently that of tobacco vein mottle potyvirus was shown to bind to the viral coat protein (4). An amino acid residue triplet, Asp-Ala-Gly (DAG), is conserved near the N-terminus of the coat protein of all aphid-transmissible potyviruses and, as this is exposed on the surface, it is thought that this domain interacts with the helper component.

Internally borne viruses have to cross one or more barriers to enter the vector, and further barrier(s) to pass from the vector to the plant. In arthropod vectors, the virus has to cross the gut wall to enter the hemocoel and the basal lamina of the salivary glands to be able to be introduced into the new host together with saliva. The details of the routes by which luteoviruses pass these barriers in their aphid vectors by endocytosis and movement through various vesicles have been detailed by Gildow (5). However, little is known at present about details of the interactions between the virus and vector involved in receptor binding at these barriers; it is these interactions that presumably control vector specificity. Stability within the hemocoel appears to be controlled by *symbionin*, the product of symbiotic bacteria in the vector, which has [chaperonin](#) like properties (6).

Internally-borne viruses that replicate in their vector often have to cross further barriers to enter the insect cells in which they replicate (7). Many of these viruses, which can be regarded as viruses of

insects that are adapting to plants, have more complex structures (such as Reo virus). There are surface proteins on the virus particles that are presumed to interact with cell surface receptors.

### 3. Replication cycle

The details of events leading to the production of new virus particles from incoming virus into an uninfected cell differ according to the type of viral genome, although there are some commonalities in the general strategy. The initial event is the uncoating of the virus particles followed by expression of the gene product(s) involved in genome replication; for DNA viruses and (–)-strand RNA viruses, this will require [transcription](#) of the incoming viral genome. The viral genome is then replicated, and later gene products such as viral coat protein and vector transmission proteins are expressed. The newly synthesised viral genome passes to adjacent uninfected cells (see below) or is encapsidated and accumulates in this initially infected cell. For many viruses with (+)-strand RNA genomes, virus particle uncoating occurs by the process of cotranslational disassembly (8). In this process, the particle structure is relaxed (often by the chelation of divalent cations and/or a pH of about 7), and host cell [ribosomes](#) locate the 5' end of the viral genome, which is an mRNA. The ribosome translocates along the RNA, translating the viral genetic information and, at the same time, displacing the viral coat protein subunits. For TMV, it is thought that cotranslational disassembly slows down or stops at the origin of assembly (see [Tobacco Mosaic Virus](#)); for some isometric viruses, it is thought that the whole genome is disassembled by this process. It is not proven that all (+)-strand RNA viruses employ cotranslational disassembly. Obviously, viruses with DNA genomes and (–)-strand or double-stranded RNA genomes cannot be uncoated by this process, as their genomic nucleic acid is not translated, and the mechanism of their uncoating is unknown. Viruses with (–)-strand and double-stranded RNA genomes have an **RNA-dependent RNA polymerase** incorporated within the virus particle, and thus the genome can be transcribed into a (+)-strand on entry into the cell. In viruses with DNA genomes, the genome is transported to the nucleus, where it is transcribed by the host DNA-dependent RNA polymerase II.

All RNA viruses encode an RNA-dependent RNA polymerase, and most of the viruses with (+)-strand genomes do so at their 5' end (see Fig. 1 in [Viruses, Plant](#), [Tobacco Mosaic Virus](#), and [Tomato Bushy Stunt Virus](#)) (also see [RNA Polymerases, RNA-Dependent](#)); thus, these protein(s) will be expressed early during the cotranslational disassembly process. For (+)-strand RNA viruses, production of the (–)-strand RNA starts at the 3' end of the input (+)-strand, proceeding toward the 5' end and, at least in the case of TMV, completes the disassembly of the virus particles (see [Tobacco Mosaic Virus](#)). The promoters for (–)-strand synthesis are thus located in the 3' region of the viral genome. The sequences at the 3' ends of genomic RNA of several genera of plant viruses can be folded into **transfer RNA**-like structures that can be aminoacylated with histidine (most tobamoviruses), tyrosine (bromoviruses, cucumoviruses), or valine (tymoviruses, one tobamovirus) (see Fig. 1 in [Viruses, Plant](#)). Mutagenesis indicates that these tRNA-like structures are required for RNA replication and especially for the determination of template specificity. Some viruses (see Fig. 1 in [Viruses, Plant](#)) have a **poly(A)** tract at their 3' end and thus resemble the animal-infecting picornaviruses. It is thought that, as with, say, [Polio virus](#), this tract is important in initiating replication. Most other (+)-strand RNA viruses do not have obvious sequence features at their 3' ends, but this region must be involved in initiating replication. The (–)-strand is then the template for the formation of both the full-length and any subgenomic RNA. Similar considerations must be involved in sequences initiating (+)-strand synthesis, but little is known about these. In some cases, however, the sequences required for initiating subgenomic RNA synthesis have been identified [see 9 for a review].

The replication of viruses with (–)-strand or double-stranded RNA genomes is similar to that of (+)-strand RNA viruses except, as described above, the initial event is transcription of the incoming genome by encapsidated RNA-dependent RNA polymerase.

Single-stranded DNA viruses (the geminiviruses) replicate by a **rolling circle** mechanism [see 10, 11

for reviews]. It is thought that the virus gene products interact with the host [cell cycle](#) and host DNA-dependent **DNA polymerase** to overcome the limitation of one round of DNA replication per cell division.

The viruses with double-stranded DNA genomes (badnaviruses and caulimoviruses) are pararetroviruses and replicate by **reverse transcription**. Details of the replication cycle are given in [Cauliflower mosaic virus](#) .

#### 4. Virus Movement through the Plant

To spread from the cell initially infected to an adjacent one, a virus has to pass through the cellulose cell wall. Plant cells communicate with each other by cytoplasmic connections through the cell wall termed *plasmodesmata* [for details of plasmodesmata structure, see ([12](#))] and it is generally accepted that viruses move from cell to cell via them. The effective diameter of plasmodesmata is too small to allow the direct movement of virus particles, and many viruses encode gene product(s) that help with this movement. A further problem is that any permanent increase in plasmodesmatal diameter would affect the control of cell-to-cell communication and thus be detrimental to the plant. Two mechanisms of cell-to-cell movement of viruses are currently recognized, although it is likely that others will be found in the future. In the first, exemplified by TMV, a *movement protein*, termed P30 (see [Tobacco Mosaic Virus](#) for details of genome organization), temporarily increases the exclusion limit of plasmodesmata from 850 Da to more than 10 kDa. P30 binds to the RNA genome of TMV, melting out its secondary and tertiary structure and thus making long narrow nucleoprotein molecules that are thought to pass through the enlarged plasmodesmata. Observations on the intracellular location of P30 indicate that it interacts with [microtubules](#) and to a lesser extent with [actin](#) filaments ([13](#), [14](#)). It is suggested that these interactions might target the movement protein, and hence the movement protein-RNA complex, to the plasmodesmata. After it has functioned, P30 is phosphorylated and “stored” adjacent to plasmodesmata. Cowpea mosaic comovirus (CPMV) and CaMV (see [Cauliflower Mosaic Virus](#) for details of genome organization) are examples of the second mechanism. One or more gene products of these viruses (P58 / 48 of CPMV; gene 1 product of CaMV) is involved in forming permanent tubular structures that pass from cell to cell across the cell wall, most probably through plasmodesmata. Virus particles are observed in these tubules, and it is thought that they are moving from cell to cell. The functional movement proteins have been attributed to gene products of many groups of viruses (see Fig. 1 in [Viruses, Plant](#)). In a few cases, such as those cited above, there is direct evidence for the involvement of these proteins in cell-to-cell movement. In other cases, the evidence is indirect, such as [mutagenesis](#) leading to subliminal infections.

The mechanisms for cell-to-cell movement allow the passage of the virus through the mesophyll tissues to the vascular bundle cell. The mechanisms of movement of viruses across the bundle sheath (cells that surround the vascular bundle) and from cell to cell within the vascular bundle are not clearly understood.

Long-distance movement of viruses causing systemic infection usually occurs in the phloem sieve elements following the source-sink movement of photoassimilates. However, it is not known how viruses enter or exit the sieve element system. For many viruses, the viral coat protein is required for long-distance movement and, most probably, for cell-to-cell movement within the vascular bundle. Long-distance movement in the xylem elements has been suggested for southern bean mosaic sobemovirus.

#### 5. Symptom Production

The ability to manipulate and mutate viral genomes has led to the identification of sequences that influence symptom production, and there are numerous publications on the effects of such changes. For example, a single nucleotide has been identified as being responsible for differences in symptoms of different isolates of the RNA virus TMV ([15](#)) and of the single-stranded DNA virus,

maize streak geminivirus (16); also both gene II and VI products are involved in the symptom and host range determination of CaMV (17, 18). However, little is known about the details of the interactions that these mutations have with the host.

There are many fewer studies on details of how the host reacts to virus infection. A progressive series of metabolic changes occur with a local lesion caused by the infection of *Cucurbita pepo* by CMV (19). At least some Potyviruses transiently inhibit host gene expression in the cells in which they replicate, thus resembling some animal viruses (20).

### Bibliography

1. P. Palukaitis, M. J. Roossinck, R. G. Dietzgen, and R. I. B. Francki (1992) *Adv. Virus Res.* **41**, 281–348.
2. T. P. Pirone and S. Blanc (1996) *Ann. Rev. Phytopathol.* **34**, 227–247.
3. I. Schmidt, S. Blanc, P. Esperandieu, G. Kuhl, G. Devauchelle, C. Louis, and M. Cerutti (1994) *Proc. Natl. Acad. Sci. USA* **91**, 8885–8889.
4. S. Blanc, J.-J. López-Moya, R. Wang, S. García-Lampasona, D. W. Thornbury, and T. P. Pirone (1997) *Virology* **231**, 141–147.
5. F. E. Gildow (1987) *Curr. Topics Vec. Res.* **4**, 93–120.
6. J. F. J. M. van den Heuvel, M. Verbeek, and F. van der Wilk (1994) *J. Gen. Virol.* **75**, 2559–2565.
7. E. D. Ammar (1994) *Adv. Dis. Vec. Res.* **10**, 289–331.
8. T. M. A. Wilson (1985) *J. Gen. Virol.* **66**, 1201–1207.
9. K. W. Buck (1996) *Adv. Virus Res.* **47**, 159–251.
10. D. M. Bisaro (1996) In *DNA Replication in Eukaryotic Cells*, (M. de Pamphelis, ed.) Cold Spring Harbor Laboratory Press, New York, pp. 833–854.
11. S. G. Lazarowitz (1992) *Crit. Rev. Plant Sci.* **11**, 327–349.
12. R. L. Overall and L. M. Blackman (1996) *Trends Plant Sci.* **1**, 307–311.
13. M. Heinlein, B. L. Epel, H. S. Padgett, and R. N. Beachy (1995) *Science* **270**, 1983–1985.
14. B. G. McLean, J. Zupan, and P. C. Zambryski (1995) *Plant Cell* **7**, 2101–2114.
15. N. Banerjee, J.-Y. Wang, and M. Zaitlin (1995) *Virology* **207**, 234–239.
16. M. I. Boulton, D. I. King, J. Donson, and J. W. Davies (1991) *Virology* **183**, 114–121.
17. R. Stratford and S. N. Covey (1989) *Virology* **172**, 451–459.
18. S. G. Qiu and J. E. Schoelz (1992) *Virology* **190**, 773–782.
19. L. I. Técsi, A. M. Smith, A. J. Maule, and R. C. Leegood (1996) *Plant Physiol.* **111**, 975–985.
20. D. Wang and A. J. Maule (1995) *Science* **267**, 229–231.
21. R. Hull (1991) *Seminars Virol.* **2**, 79–80.

### Suggestions for Further Reading

22. L. Bos (1963) *Symptoms of Virus Diseases of Plants*, Centre for Agricultural Publications and Documentation, Wageningen. Common symptoms are listed.
23. R. Hull (1994) Molecular biology of plant virus–vector interactions. *Adv. Dis. Vect. Res.* **10**, 361–386. Addresses movement between plants.
24. R. E. F. Matthews (1991) *Plant Virology*, 3rd ed., Academic Press, San Diego, CA. See for a discussion of general subject.
25. L. A. Mezitt and W. J. Lucas (1996) Plasmodesmatal cell-to-cell transport of proteins and nucleic acids. *Plant Mol. Biol.* **32**, 251–273. See for a discussion of movement within plant.
26. K. Séron and A.-L. Haenni (1966) Vascular movement of plant viruses. *Mol. Plant Microb. Intern.* **9**, 435–442. See for a discussion of movement within plant.

## Virus Structure

**Viruses** spend part of their time inside the living cells they have infected, either in a latent form or, more familiarly, replicating to produce progeny viruses. The rest of the time they are outside the cell in the form of **virions**, the discrete particles that usually come to mind when we think of viruses. Virions are quite varied in shape and size among different known viruses, but detailed examination of their structures reveals that they all share a set of common features (described below) that can be seen to derive from the biological functions that are carried out by the virion. Virions are small, ranging from roughly 30 nm to 300 nm in diameter or longest dimension. The upper end of this range overlaps slightly the lower end of the range of known sizes for cells, but any specific virus is always substantially smaller than the cells it infects.

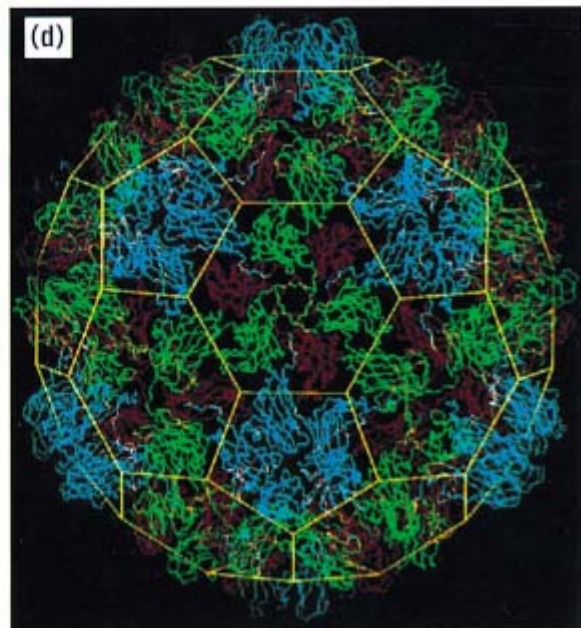
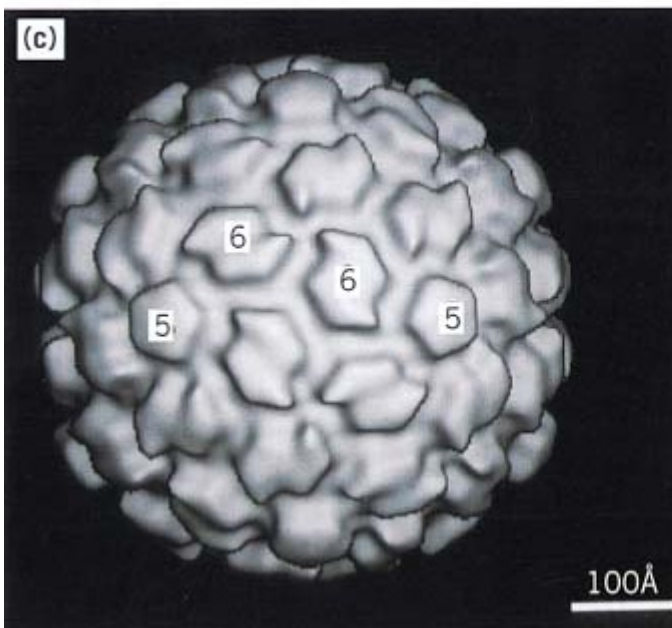
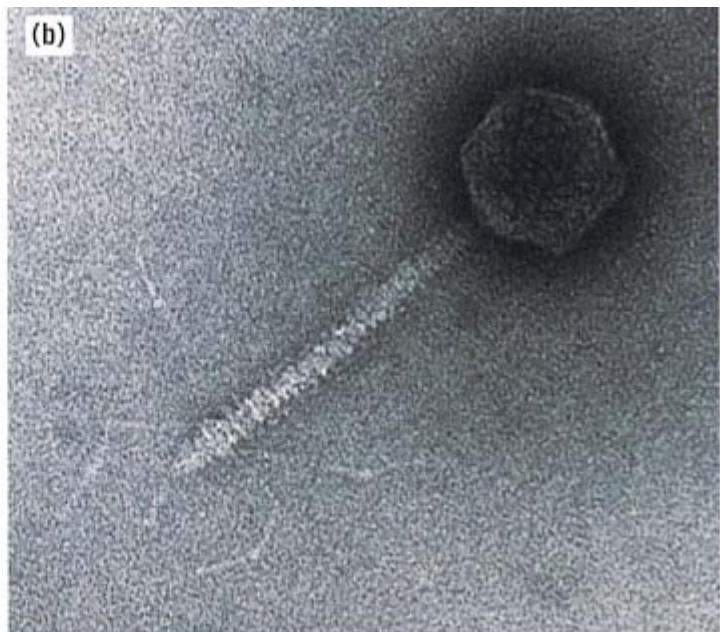
### 1. Methods Used to Determine Virus Structure

Their small size puts virions below the practical resolution limit of the light **microscope**. Because we are therefore denied a direct “look” at the virions, the way we conceptualize them depends on a series of more indirect techniques, which will be reviewed here briefly before describing more specifically what is known about virus structure. When viruses were first discovered a century ago, they were initially defined simply as disease-causing agents that could pass through a filter that retains bacteria. This provides a rough upper limit on the size of virions, and this estimate was refined in subsequent years by using sets of filters with graded series of pore sizes and by biophysical techniques such as **sedimentation velocity** and **sedimentation equilibrium**, which gave estimates of both size and shape for the virus particles under study. However, it was not until the invention of the **electron microscope** that a detailed idea of the shapes of virions was possible. The first experiments with viruses and electron microscopy showed that [tobacco mosaic virus](#) (TMV) is rod-shaped (and of dimensions consistent with the earlier biophysical measurements) and that bacteriophage T2 looks somewhat like a tadpole or a spermatozoan. The quality of the images obtained by electron microscopy improved substantially in the mid-1950s with the advent of the technique of **negative staining**, which provides improved contrast and also stabilizes virions against the distortion produced when they are dried before being inserted into the vacuum of the microscope. Much of what we currently know about virus structure comes from experiments with negative stain electron microscopy (in conjunction with an extensive catalog of genetic, biochemical, and biophysical methods that make it possible to enumerate the component molecules in the virion, dissect the virion, and localize different components in the structure). In recent years, three additional techniques have substantially increased the resolution with which we know the structures of the virions to which they can be applied. These are, first, [cryoelectron microscopy](#), which allows virions to be viewed in the frozen hydrated state, with their structure essentially perfectly preserved. The resulting images have very low contrast and are of limited use by themselves, but the second technique, computerized image processing and filtering, allows a large number of individual images to be combined to produce three-dimensional **reconstructions** of the virions of very good quality. The third technique, [X-ray crystallography](#), has now been applied to a number of viruses and yields three-dimensional structures at near atomic resolution. A particularly interesting recent development is that in favorable cases, the above three techniques can be combined to produce structural results that neither technique alone would be able to produce.

While it is impractical to reproduce here examples of all the structural types that have been determined for different viruses, Figure 1 shows a few representative examples of the different sorts of images that are produced. There are now extensive collections of such images widely available on

the World Wide Web.

**Figure 1.** Different methods of visualizing virion structures. (a) One of the first electron micrographs of a virus. This picture of bacteriophage T2 is from Luria and Anderson (5). (b) Bacteriophage I, negatively stained with uranyl formate (courtesy of Robert L. Duda). (c) Image of the “Prohead II” capsid precursor of bacteriophage HK97, viewed down a twofold symmetry axis at a resolution of 24 Å. This image is reconstructed from images of approximately 100 individual particles, produced by cryoelectron microscopy. The numbers superimposed on the capsid indicate the positions of some of the pentamers and hexamers of the identical protein subunits (courtesy of James F. Conway). (d) Image of Cowpea Chlorotic Mottle Virus at near-atomic resolution, produced by X-ray crystallography and viewed down a three-fold symmetry axis. The red, green, and blue lines are traces of the  $\alpha$ -carbon backbones of the 180 protein subunits. (The positions of the amino acid side chains are known but omitted from this representation for clarity.) Like colors identify protein subunits related to each other by the icosahedral symmetry axes. The yellow lines show the positions of the hexamers and pentamers of the protein subunits (courtesy of John E. Johnson). Approximate magnifications for the four images are: (a), 53,000 $\times$ ; (b) 320,000 $\times$ ; (c) 1,100,000 $\times$ ; (d) 2,000,000 $\times$ . See color insert.



## 2. Composition and Structure of Virions

Virions are the structures by way of which the viral genome travels from one cell to another—from the cell where the virion was produced to the cell it will infect. As such, the virion has two apparent functions. First, it must protect the [genome](#) from environmental damage; and second, it must successfully deliver the genome to the next cell. Viruses with different hosts typically face rather different “problems” in both regards, and these differences, particularly differences in mechanisms of infection, appear to account for many of the observed differences in the structures of virions. Beyond the various structural elaborations that address differences in extracellular environment and differences in the properties of cells to be infected, we will see that the viral coats are constructed according to a rather small number of structural principles.

All virions must of course contain a copy of the nucleic acid that constitutes the viral genome. However, for different viruses the genome can be DNA or RNA, it can be single- or double-stranded, or it can exist as one molecule or several. Whatever its form, virions typically contain only one copy of the genome—that is, they are haploid—but the [retroviruses](#) have two copies per virion of their single-stranded RNA genome. Viral genomes vary in size, from having enough sequence to encode only 3 to 4 proteins (eg, the single-stranded RNA bacteriophages with ~3000 bases/genome) up to genomes of several hundred kilobase pairs. The largest genome of a well-characterized virus is the 330,740-bp dsDNA genome of the algal virus PBCV1. Bacteriophage G has a double-stranded DNA genome measured at ~650,000 bp—bigger than the smallest bacterial genomes—but it is not yet known whether that is all unique sequence DNA. As a rule, RNA genomes cluster at the small end of the genome size scale. DNA genomes span nearly the whole range, and all the biggest of the known viral genomes—over about 20 kbp—are double-stranded DNA. Some genomes are modified by the covalent addition of a protein; examples include [Poliovirus](#), which has a protein called **VpG** attached to the 5'-end of its single-stranded RNA genome with a role in facilitating [translation](#) from the genome when it acts as [messenger RNA](#), and [Adenovirus](#) and some bacteriophages (f29 and PRD1), which have proteins covalently bound to the 5'-ends of their double-stranded DNA genomes to carry out a crucial role in initiation of [DNA replication](#) (1). The genome in many viruses is associated noncovalently with protein, either virus-encoded protein(s) or host [histone](#) proteins, while in other viruses the nucleic acid is packed devoid of protein and accompanied only by water and counterions.

The only other structural feature of virions that is invariant over all known viruses is that a major component of the virion is protein. At a minimum, there is one protein, present in multiple copies (typically several hundred). More often, there are several different proteins present, one or a few in large numbers (“major proteins”) and the rest (“minor proteins”) in smaller numbers (1 to 50 copies/virion). The major protein or proteins make a container for the genome known as a **capsid**. (There is some variability in usage for different virus groups. Here we will take capsid to refer to just the helical or icosahedral proteinaceous container, and not any other nongenome components that might be part of the virion.) As discussed in more detail in the entry **Capsids**, the major protein(s) of the capsid are arranged either with helical symmetry, resulting in a rodlike structure, or with icosahedral symmetry, which produces a roughly spherical protein shell.

For some viruses, there is no more to the structure of the virion than the genome and a simple protein capsid. For many plant viruses, no additional structural elaboration is needed for successful infection because they are delivered directly to their host cells by insects feeding on the plants. Other structurally simple viruses incorporate into their capsid proteins the receptor-binding function they need to attach to cells as the first step of infection. It is also very common, however, for viruses to have additional, more or less elaborate protein structures in addition to the basic capsid, and these usually have a role in getting the virus or its genome into the cell it is infecting. Perhaps the best known of these are the tails and associated tail fibers of the double-stranded DNA bacteriophages (as well as of the morphologically similar examples of archaeal viruses). These six-fold symmetric helical structures attach to one of the five-fold symmetric “corners” of the phage head. Tails can be short structures or extended tubes, in some cases with contractile sheaths. They function to recognize

and attach to the surface of the host cell and to mediate passage of the phage DNA across the complex bacterial cell envelope. Many viruses have fibers or spikes extending out from the capsid shell. A well-studied example is the trimeric fibers of Adenovirus that are found at each of the 12 pentameric “corners” of the virion and which are responsible for the initial attachment of the virus to its host cell. [Reovirus](#) and its relatives, including some bacteriophages, have two concentric icosahedral protein shells surrounding the genome—in this case segmented dsRNA. The outer shell is important for attachment and entry of the virus, but it is removed once the virion is in the cell. The inner shell, in contrast, remains intact with the genome inside and serves as a factory for transcribing the genome. The resulting RNA leaves the inner shell and serves as mRNA and eventually as virion RNA in progeny virions. The transcription into ssRNA and the conversion of ssRNA to genomic dsRNA are both accomplished by a virus-encoded **RNA polymerase** found inside the virion.

Many animal viruses and, much less commonly, other viruses have a [lipid](#) envelope surrounding the protein capsid and genome. The envelope typically has a central role in virus entry through the cellular membrane; this is most often (but not invariably) accomplished by fusion between the viral and cellular [membranes](#), which results in a joining of the contents of the viral envelope with the cytoplasm of the cell. The lipids in the viral envelope are derived from the membrane of the infected cell, and they are arranged in a standard bilayer. The membrane is usually heavily studded with virus-encoded glycoproteins, which have roles in cell entry and in dealing with the host immune system and other aspects of their extracellular environment. In many cases, these viral glycoproteins are responsible for inducing fusion between the viral and cellular membranes. The best-studied such “fusion protein” is the **Influenzavirus hemagglutinin**, which undergoes a dramatic rearrangement of its polypeptide fold as a part of the process of promoting fusion (2). A few viruses have internal membranes that surround the genome but are inside the protein capsid shell [eg, algal virus PBCV1 and bacteriophage PRD1 (1)]. The biological function of these internal membranes is not well understood.

Although virion structures as a class are quite well characterized, there are features of the structures of some virions that are rather poorly understood. Among these we give three examples: (i) [Retroviruses](#) have a protein shell inside the envelope and surrounding the genome, and this has been assumed to have icosahedral symmetry, like all the well-characterized “spherical” capsids. However, a recent study with this structure from the [HIV](#) virion shows that it is not icosahedral but rather has a less well-defined subunit arrangement (3). (ii) Although some of the viruses of the Archaea have the morphology of tailed bacteriophages, others have virions with the shapes of lemons or corn dogs, which have no parallels in the viruses of other hosts (4). The molecular architecture of these virions is very poorly known. (iii) Some viruses with otherwise well-characterized virions have substructures in their virions that are less well understood. For example, [Herpesvirus](#) has an apparently amorphous, thick layer of protein, called the tegument, which surrounds the icosahedrally symmetric capsid, but lies inside the membrane envelope.

See also [Capsids, viral](#).

#### Bibliography

1. J. Caldenty, J. K. Bamford, and D. H. Bamford (1990) *J. Struct. Biol.* **104**, 44–51.
2. P. A. Bullough et al. (1994) *Nature* **371**, 37–43.
3. S. K. Fuller et al. (1997) *Curr. Biol.* **7**, 729–738.
4. P. Stolt and W. Zillig (1994) In *Encyclopedia of Virology*, Vol. **1** (R. G. Webster and A. Granoff, eds.), Academic Press, London, pp. 50–58.
5. S. E. Luria and T. F. Anderson (1942) *Proc. Natl. Acad. Sci. U.S.A.* **28**, 127–130.

#### Suggestions for Further Reading

6. W. Chiu, R. M. Burnett, and R. L. Garcea, (eds.) (1997) *Structural Biology of Viruses* Oxford University Press, New York, 1997.



7. B. N. Fields, D. M. Knipe, and P. M. Howley (eds.) (1996) *Virology* Lippincott-Raven, Philadelphia.

## Viruses, Animal

Viruses have long been invisible partners, because of small size of their particles, of life on this planet, including animals, and the discipline of virology was born only 100 years ago. Compared to cells, virus particles have much simpler structures, with splendid regularity (see Fig. 1 of [Virus Infection, Animal](#)). Therefore, **X-ray crystallographic** studies have been performed successfully for some kinds of viruses and have generated precise structures of the whole virion particles. These structural studies give an insight into molecular basis of interaction of proteins in assembled proteins of virion particles.

Viruses carry their own genomes that encode the genetic information for their replication in infected cells. Because of being obligatory parasites of host cells, viruses largely depend on host components and metabolism for their replication. In general, viruses use macromolecular biosynthesis pathways common to those in host cells. They often have, however, special mechanisms of [DNA replication](#) and gene expression that give advantage to viruses in the competition for cellular resources, resulting in death or immortalization of the host cells. Thus studies of virus replication have frequently led to discoveries of novel mechanisms in infected cells.

A number of molecular interactions between viral and host factors occur in infected hosts. Development of virus diseases must reflect the sum total of the interactions between virus and cells involved in the virus-specific response to replication, dissemination in an entire body, and host [immune responses](#). Virus replications are actually limited to specific species, tissues, and cell types within hosts. These virus responses must be due to the distribution of host factors that support the virus replication. As such host factors, new molecules with unknown natural functions may be discovered, and new research projects with such molecules may open in other fields of the life sciences, such as cell biology. New functions of known molecules may also be found. Thus, identification of host factors affecting the virus replication and dissemination is very important for understanding the development of virus diseases, as well as of host mechanisms at molecular level. Complete understanding of the molecular basis of virus pathogenicity and replication in a whole body would help in the use of viruses as virus-related expression vectors in future gene therapy.

Animal virology has always been one of the leading edges in the field of molecular biology. This volume contains current molecular biological information available from studies on representative animal viruses, and individual entries describe the most recent molecular biological information of a particular virus: [adenovirus](#), [cytomegalovirus](#), [Epstein–Barr virus](#), [helper virus](#), [hepatitis B virus](#), [herpesvirus](#), [HIV](#), [influenza virus](#), [papovavirus](#), [poliovirus](#), [polyomavirus](#), [reovirus](#), [rhinovirus](#), [Rous sarcoma virus](#), [Semliki Forest virus](#), [Sendai virus](#), [SV40](#), [vaccinia virus](#), and [vesicular stomatitis virus](#). The information contained is unlikely to be altered in the future, although new and important concepts of the field of animal virology will continue to contribute to all of the life sciences.

## Viruses, Plant

More than 780 viruses that infect **plants**, some of which are of major agricultural importance, have now been recognized. These have been grouped into 48 genera and 10 families by the International Committee on Taxonomy of Viruses (1) (Table 1). This classification is based on various characters, such as composition and organization of their genomes, particle morphology, and biological properties (see [Virus Infection, Plant](#)). Most of these virus families are unique to plants but some, eg, Reoviridae, Rhabdoviridae, and Bunyaviridae, also have members that infect animals, including invertebrates; it is thought that members of these groups might be viruses of their insect vectors that are becoming adapted to plants. Because plant viruses often occur in relatively high concentrations in plants, many of them can be easily purified and have therefore attracted attention as models for studies on structure, protein:nucleic acid interactions, and molecular biology. Details on the pathology of these viruses are given in [Virus infection, plant](#).

**Table 1. Classification of Plant Viruses**

| Genome <sup>a</sup> | Family         | Genus            | Type member                   | Particle shape | No. genome segments | No. spp. | No. tentative spp. |
|---------------------|----------------|------------------|-------------------------------|----------------|---------------------|----------|--------------------|
| ds DNA (RT)         |                | Badnavirus       | Commelina yellow mottle virus | b <sup>a</sup> | 1                   | 10       | 4                  |
|                     |                | Caulimovirus     | Cauliflower mosaic virus      | i              | 1                   | 11       | 6                  |
| ss DNA              | Geminiviridae  | Subgroup I       | Maize streak virus            | g              | 1                   | 11       | 2                  |
|                     |                | Subgroup II      | Beet curly top virus          | g              | 1                   | 1        | 2                  |
|                     |                | Subgroup III     | Bean golden mosaic virus      | g              | 1 or 2              | 49       | 9                  |
| ds RNA              | Reoviridae     | Fijivirus        | Fiji disease virus            | ic             | 12                  | 5        | 0                  |
|                     |                | Oryzavirus       | Rice ragged stunt virus       | ic             | 10                  | 2        | 0                  |
|                     |                | Phytoreovirus    | Wound tumor virus             | ic             | 12                  | 3        | 0                  |
|                     | Partitiviridae | Alphacryptovirus | White clover cryptic virus 1  | i              | 2                   | 16       | 10                 |
|                     |                | Betacryptovirus  | White clover cryptic virus 2  | 1              | 2                   | 4        | 1                  |

|               |               |                   |                                |     |   |    |    |
|---------------|---------------|-------------------|--------------------------------|-----|---|----|----|
| ss RNA<br>(-) | Rhabdoviridae | Cytorhabdovirus   | Lettuce necrotic yellows virus | bc  | 1 | 8  | 0  |
|               |               | Nucleorhabdovirus | Potato yellow dwarf virus      | bc  | 1 | 6  | 0  |
|               |               | Unassigned        |                                | bc  | 1 | 61 | 0  |
|               | Bunyaviridae  | Tospovirus        | Tomato spotted wilt virus      | ic  | 3 | 2  | 0  |
|               |               | Tenuivirus        | Rice stripe virus              | f   | 4 | 4  | 3  |
| ss RNA<br>(+) | Bromoviridae  | Alfamovirus       | Alfalfa mosaic virus           | b   | 3 | 1  | 0  |
|               |               | Bromovirus        | Brome mosaic virus             | i   | 3 | 6  | 0  |
|               |               | Cucumovirus       | Cucumber mosaic virus          | i   | 3 | 3  | 0  |
|               |               | Iilarvirus        | Tobacco streak virus           | (i) | 3 | 16 | 0  |
|               | Comoviridae   | Comovirus         | Cowpea mosaic virus            | i   | 2 | 15 | 0  |
|               |               | Fabavirus         | Broadbean wilt virus           | i   | 2 | 3  | 0  |
|               |               | Nepovirus         | Tobacco ringspot virus         | i   | 2 | 28 | 8  |
|               | Potyviridae   | Bymovirus         | Barley yellow mosaic virus     | f   | 1 | 5  | 0  |
|               |               | Potyvirus         | Potato virus Y                 | f   | 1 | 75 | 93 |
|               |               | Rymovirus         | Ryegrass mosaic virus          | f   | 1 | 5  | 2  |
|               | Tombusviridae | Carmovirus        | Carnation mottle virus         | i   | 1 | 12 | 7  |
|               |               | Tombusvirus       | Tomato bushy stunt virus       | i   | 1 | 13 | 0  |
|               | Sequiviridae  | Sequivirus        | Parsnip yellow fleck virus     | i   | 1 | 3  | 0  |

|               |  |   |   |    |    |
|---------------|--|---|---|----|----|
| Waikavirus    | Rice tungro<br>spherical<br>virus      | i | 1 | 3  | 0  |
| Carlavirus    | Carnation<br>latent virus              | f | 1 | 29 | 31 |
| Capillovirus  | Apple stem<br>grooving<br>virus        | f | 1 | 3  | 1  |
| Closterovirus | Beet<br>yellow<br>virus                | f | 1 | 6  | 19 |
| Dianthovirus  | Carnation<br>ringspot<br>virus         | i | 1 | 3  | 1  |
| Enamovirus    | Pea enation<br>mosaic<br>virus         | i | 2 | 1  | 0  |
| Furovirus     | Soil-borne<br>wheat<br>mosaic<br>virus | r | 2 | 5  | 5  |
| Hordeivirus   | Barley<br>stripe<br>mosaic<br>virus    | r | 3 | 4  | 0  |
| Idaeovirus    | Raspberry<br>bushy<br>dwarf virus      | i | 2 | 1  | 0  |
| Luteovirus    | Barley<br>yellow<br>dwarf virus        | i | 1 | 16 | 14 |
| Machlomovirus | Maize<br>chlorotic<br>mottle<br>virus  | i | 1 | 1  | 0  |
| Marafivirus   | Maize<br>rayado fino<br>virus          | i | 1 | 3  | 0  |
| Necrovirus    | Tobacco<br>necrosis<br>virus           | i | 1 | 2  | 2  |
| Potexvirus    | Potato<br>virus X                      | f | 1 | 20 | 21 |
| Sobemovirus   | Southern<br>bean<br>mosaic<br>virus    | i | 1 | 10 | 5  |
| Tobamovirus   | Tobacco<br>mosaic<br>virus             | r | 1 | 13 | 2  |

|             |                                |   |   |    |   |
|-------------|--------------------------------|---|---|----|---|
| Tobravirus  | Tobacco rattle virus           | r | 2 | 3  | 0 |
| Trichovirus | Apple chlorotic leafspot virus | f | 1 | 2  | 3 |
| Tymovirus   | Turnip yellow mosaic virus     | i | 1 | 17 | 1 |
| Umbravirus  | Carrot mottle virus            | — | 1 | 5  | 4 |

<sup>a</sup> Genome type: ds DNA(RT) = double-stranded DNA (replicating by reverse transcription); ss DNA = +single-stranded DNA; ds RNA = double-stranded RNA; ss RNA(-) = (-)-sense RNA; ss RNA(+) = (+)-sense RNA(mRNA).

<sup>b</sup> b = bacilliform; bc = complex bacilliform; f = filamentous rod; g = geminate; i = isometric; (i) = possibly isometric; i = complex isometric; r = rigid rod.

## 1. Structure

The particles of most plant viruses have a relatively simple structure, being composed of the viral [genome](#) surrounded by a **capsid** composed of one or more coat [protein](#) species. As well as protecting the genome from degradation by, say, **nucleases**, the capsid presents the external features of the virus to the environment. Thus, it is involved in interactions essential to stages of the virus infection cycle external to the cell. Viruses require such interactions in movement from plant to plant, but not from cell to cell within the plant (see [Virus Infection, Plant](#)). There are two basic forms of structure: *rod-shaped* and *isometric* particles (Table 1).

The rod-shaped particles of viruses range from short rigid rods (eg, [tobacco mosaic virus](#), TMV) to long flexuous rods (eg, beet yellows closterovirus). In these particles, the coat protein subunit is arranged in a helix, with the RNA genome embedded in it (see [Tobacco Mosaic Virus](#) for more detail) and the particles are stabilized by both protein:protein and protein:RNA interactions.

The structure of isometric particles follows the constraints imposed by **icosahedral symmetry** (2), with the **quasi-equivalent** arrangement of coat protein subunits according to strict surface lattice criteria. Most plant isometric particles have a triangulation number  $T = 3$  (180 coat protein subunits), but there are examples of  $T = 1$  (60 subunits; eg, *satellite tobacco necrosis virus*) and  $T = 7$  (420 subunits; eg, [cauliflower mosaic virus](#)). A range of interactions are involved in the stabilization of the isometric particles. For instance, the particles of *alfalfa mosaic alfamovirus* (AMV) and *cucumber mosaic cucumovirus* (CMV) are stabilized primarily by [electrostatic interactions](#) between protein and RNA and dissociated by high salt concentrations. Those of *brome mosaic bromovirus* (BMV) are stabilized by protein:RNA interactions and a protein:protein interaction that involves an unusual carboxylate group which protonates above pH 7. The particles of *southern bean mosaic sobemovirus* (SBMV) have the two interactions found in BMV and also a protein:protein interaction mediated by  $\text{Ca}^{2+}$ . Adjustment of the pH to more than 7 and chelation of  $\text{Ca}^{2+}$  leaves the particles of viruses such as BMV and SBMV stabilized by just protein:RNA interactions; in this situation, the particles appear to swell hydrodynamically. These swollen particles and those of AMV and CMV have been shown to be uncoated by the cotranslational disassembly mechanism (see [Virus Infection, Plant](#)). Other viruses, such as *cowpea mosaic comovirus*, have icosahedral particles stabilized by

strong protein:protein interactions with relatively little involvement of protein:RNA interactions. Preparations of such viruses have a proportion of “empty” particles that do not contain nucleic acid.

Variations on icosahedral symmetry are found among plant viruses. The particles of *geminiviruses* appear to be two  $T = 1$  icosahedra joined together at a five-fold axis (3). Other viruses (eg, alfalfa mosaic virus, *badnaviruses*) have bacilliform particles that can be considered isometric particles cut in two, followed by the insertion of a tubular portion made up of hexamer subunits (4).

There are some groups of plant viruses that have more complex structures. These are viruses such as rhabdoviruses, tospoviruses, and reoviruses that have counterparts among animal viruses and are considered to be insect viruses adapting to plants. Details of these structures can be found in (1).

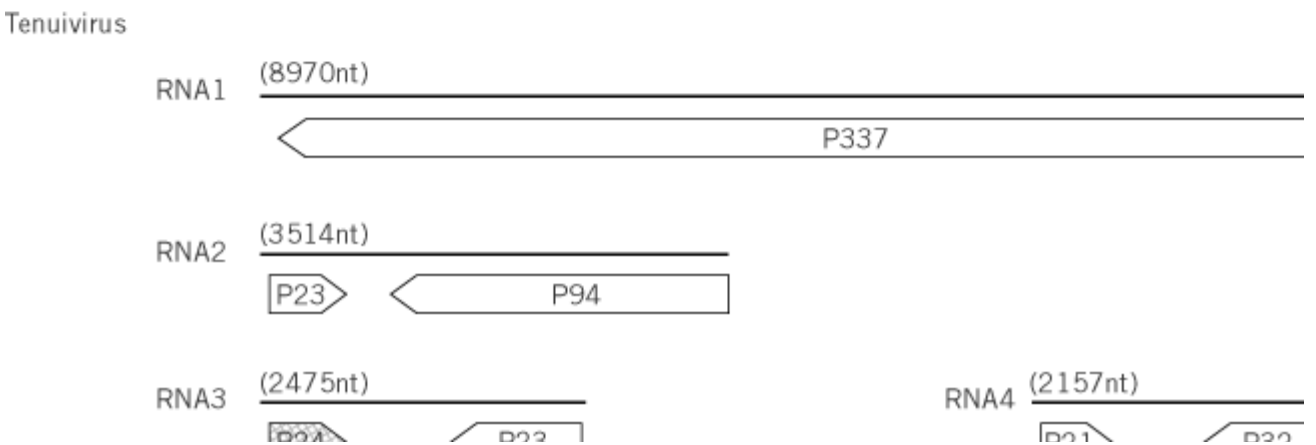
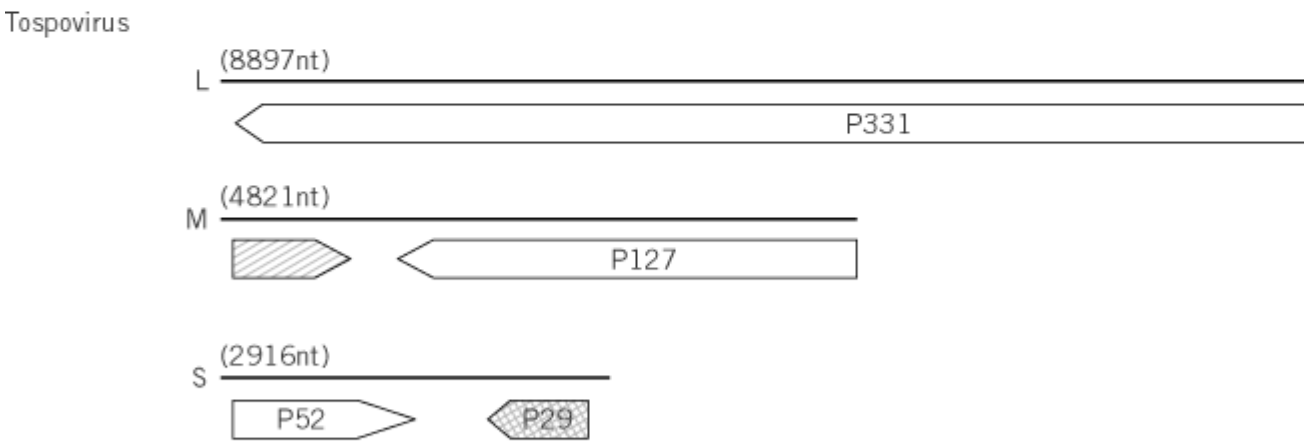
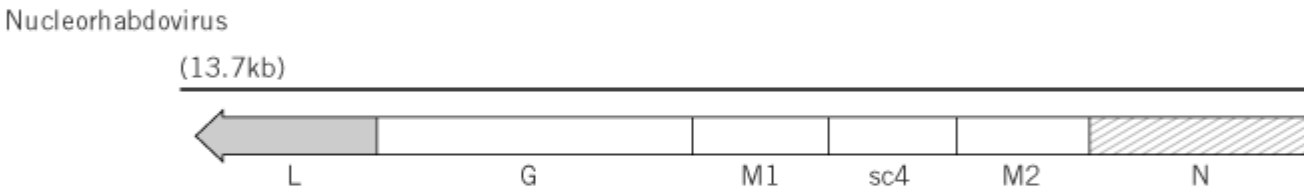
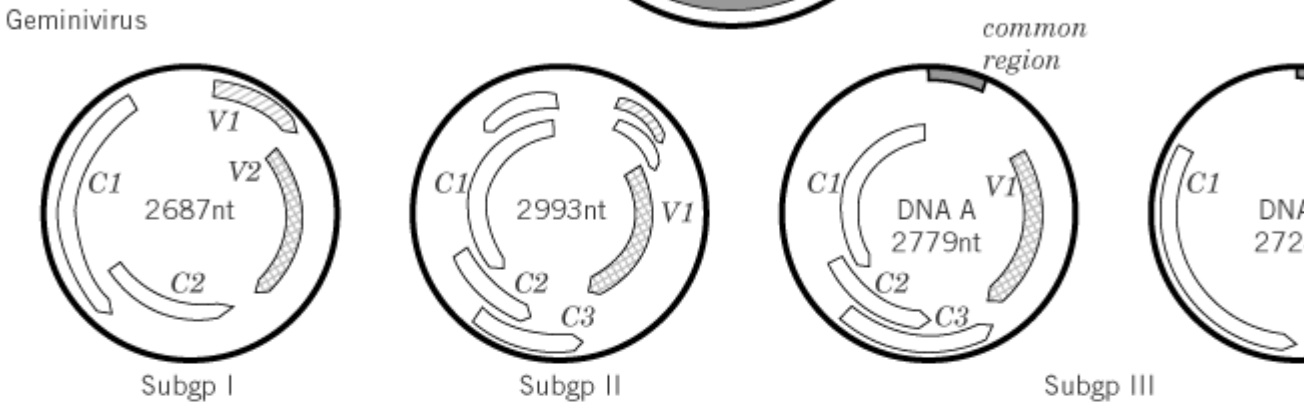
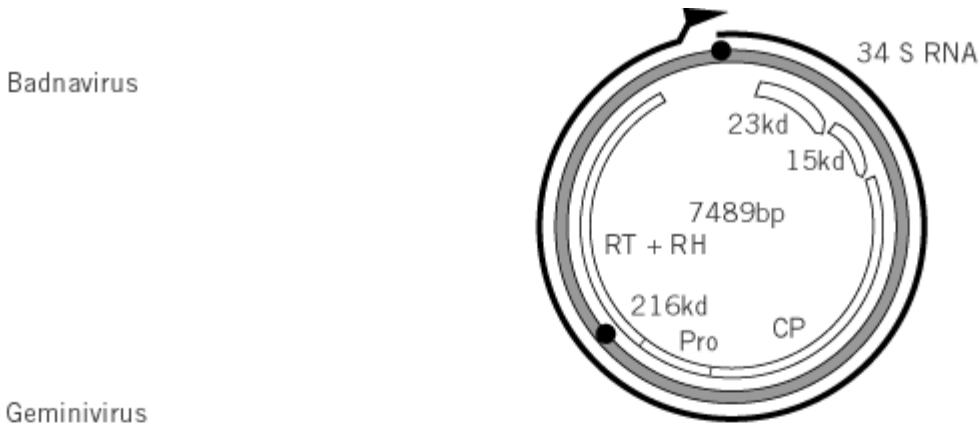
## 2. Genomes

The majority of plant viruses (71%) have (+)-strand (messenger-sense) RNA as their genome; genomes are also found that consist of double-stranded DNA (4%), single-stranded DNA (9.5%), double-stranded RNA (5%), and (–)-strand RNA (10.5%). The double-stranded DNA genome plant viruses replicate by **reverse transcription**, and there are no plant viruses known with double-stranded DNA genomes that replicate the DNA directly.

### 2.1. Genome, Structure, Organization, and Expression

Plant viruses encode products that facilitate their replication, encapsidation, and in most cases, their movement within and between plants (see [Virus Infection, Plant](#)). Expression of these products is constrained by the limitations of eukaryotic [ribosomes](#) being able to **translate** only the 5' **open reading frame** of a [messenger RNA](#) (mRNA). Both plant and animal viruses, have developed various strategies for overcoming this limitation. The genome organizations of typical members of the 48 genera (where sequence data are available) are illustrated in Figure 1.

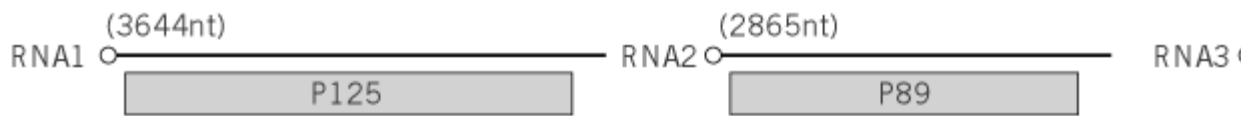
**Figure 1.** Genome organizations of groups of plant viruses. The illustration is of the type member of the group (Table 1) except subgroup III geminiviruses, African cassava mosaic virus; nucleorhabdovirus, sonchus yellow net virus; nepovirus, tomato blackring virus; carlavirus, potato virus M; dianthovirus, red clover necrotic mottle virus; umbravirus, groundnut rosette virus. The genome is represented by two concentric circles for double-stranded circular genomes; one circle for single-stranded circular genomes; a straight line for linear genomes; subgenomic mRNA are shown as thin lines. The size of the genome is given in nucleotides (nt or base pairs, bp). The terminal features of the genomes, when they are known, are listed on the figure. The gene products are indicated by boxes, with the directions of those of ambisense genomes given by an arrow. The molecular size of the gene product is  $\times 10^{-3}$ . The symbols for functions of gene products and expression by read through of a weak stop codon are indicated on the figure.



**Figure 1.** (*continued*) Genome organizations of groups of plant viruses. The illustration is of the type member of the group (Table 1) except subgroup III geminiviruses, African cassava mosaic virus; nucleorhabdovirus, sonchus yellow net virus; nepovirus, tomato blackring virus; carlavirus, potato virus M; dianthovirus, red clover necrotic mottle virus; umbravirus, groundnut rosette virus. The genome is represented by two concentric circles for double-stranded circular genomes; one circle for single-stranded circular genomes; a straight line for linear genomes; subgenomic mRNA are shown as thin lines. The size of the genome is given in nucleotides (nt or base pairs, bp). The terminal features of the genomes, when they are known, are listed on the figure. The gene products are indicated by boxes, with the directions of those of ambisense genomes given by an arrow. The molecular size of the gene product is  $\times 10^{-3}$ . The symbols for functions of gene products and expression by read through of a weak stop codon are indicated on the figure.



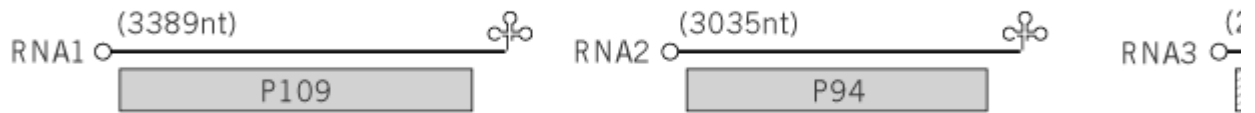
Alfamovirus



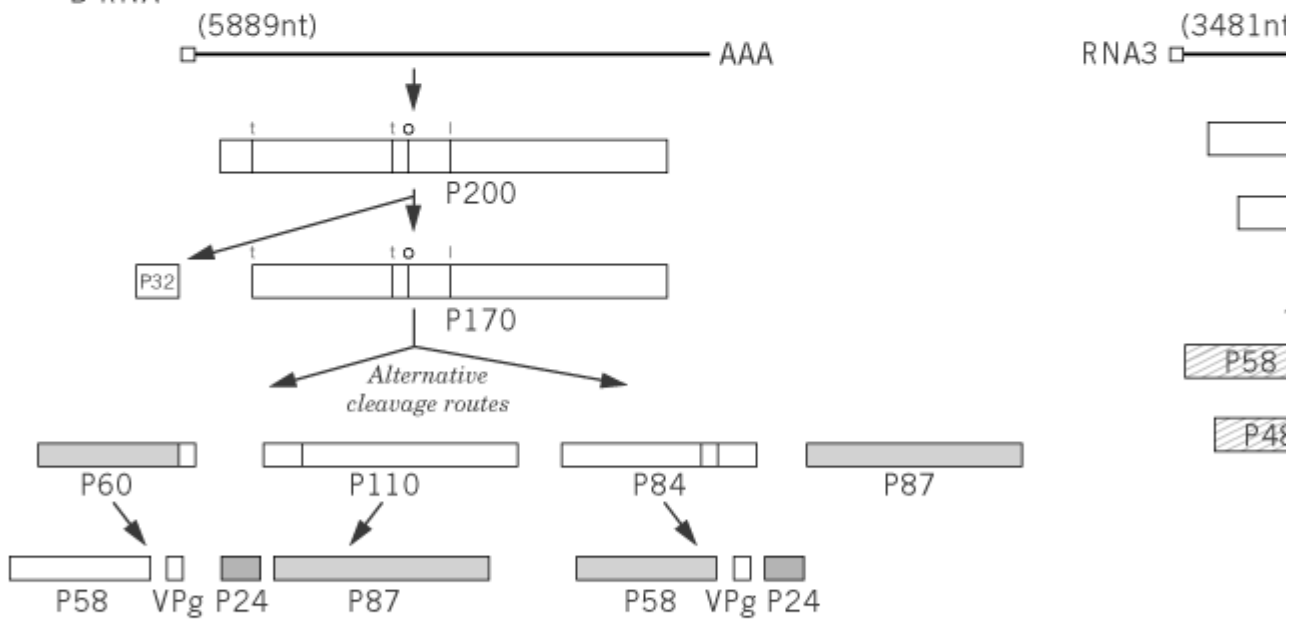
Bromovirus



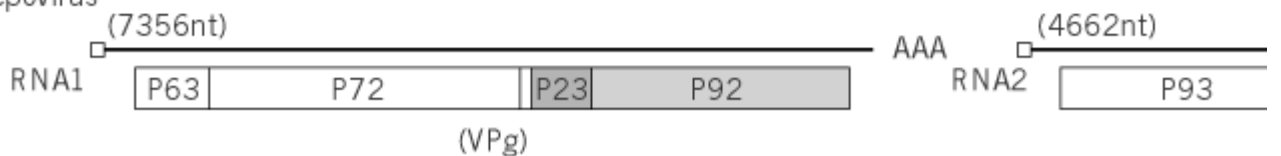
Cucumovirus



Comovirus  
B-RNA



Nepovirus



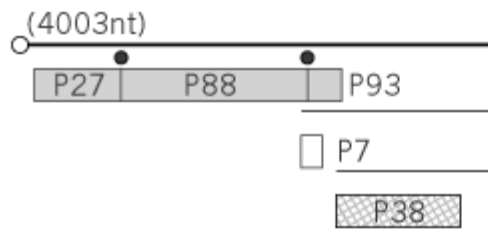
Potyvirus



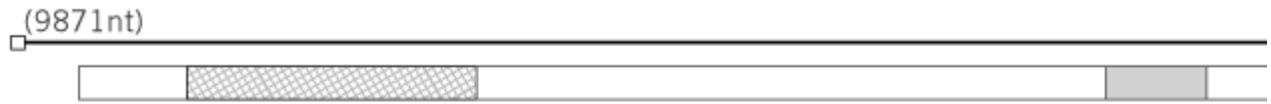
Bymovirus

**Figure 1.** (*continued*) Genome organizations of groups of plant viruses. The illustration is of the type member of the group (Table 1) except subgroup III geminiviruses, African cassava mosaic virus; nucleorhabdovirus, sonchus yellow net virus; nepovirus, tomato blackring virus; carlavirus, potato virus M; dianthovirus, red clover necrotic mottle virus; umbravirus, groundnut rosette virus. The genome is represented by two concentric circles for double-stranded circular genomes; one circle for single-stranded circular genomes; a straight line for linear genomes; subgenomic mRNA are shown as thin lines. The size of the genome is given in nucleotides (nt or base pairs, bp). The terminal features of the genomes, when they are known, are listed on the figure. The gene products are indicated by boxes, with the directions of those of ambisense genomes given by an arrow. The molecular size of the gene product is  $\times 10^{-3}$ . The symbols for functions of gene products and expression by read through of a weak stop codon are indicated on the figure.

Carmovirus



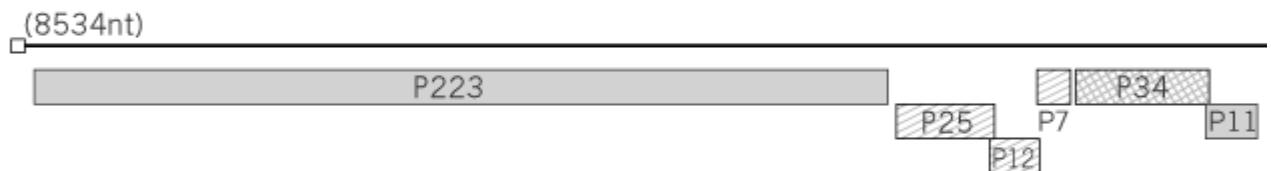
Sequivirus



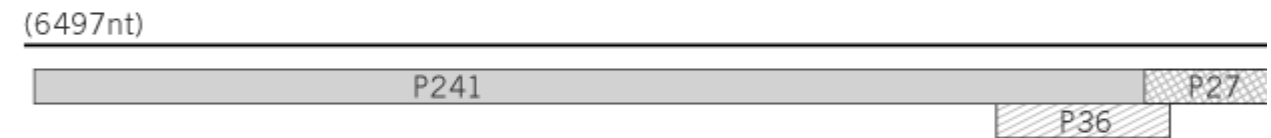
Waikavirus



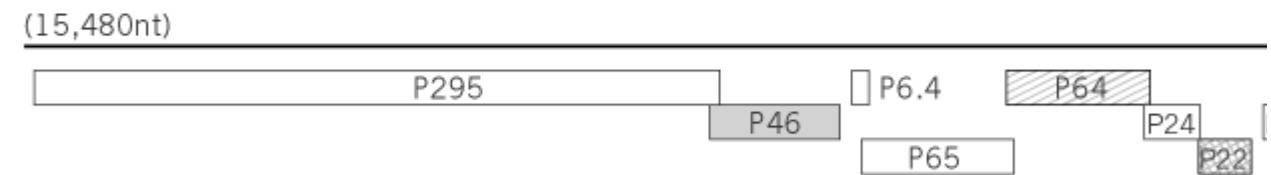
Carlavirus



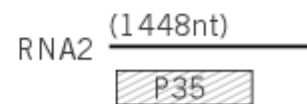
Capillovirus



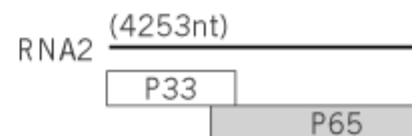
Closterovirus



Dianthovirus



Enamovirus

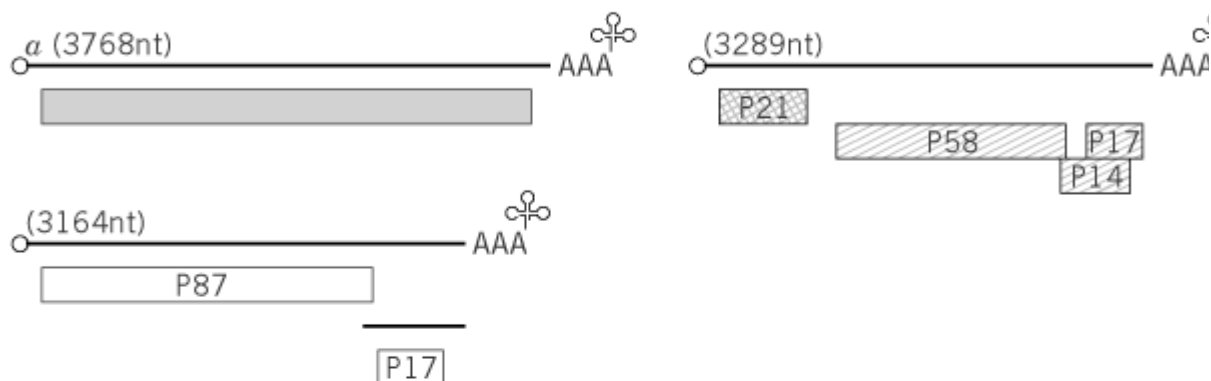


**Figure 1.** (*continued*) Genome organizations of groups of plant viruses. The illustration is of the type member of the group (Table 1) except subgroup III geminiviruses, African cassava mosaic virus; nucleorhabdovirus, sonchus yellow net virus; nepovirus, tomato blackring virus; carlavirus, potato virus M; dianthovirus, red clover necrotic mottle virus; umbravirus, groundnut rosette virus. The genome is represented by two concentric circles for double-stranded circular genomes; one circle for single-stranded circular genomes; a straight line for linear genomes; subgenomic mRNA are shown as thin lines. The size of the genome is given in nucleotides (nt or base pairs, bp). The terminal features of the genomes, when they are known, are listed on the figure. The gene products are indicated by boxes, with the directions of those of ambisense genomes given by an arrow. The molecular size of the gene product is  $\times 10^{-3}$ . The symbols for functions of gene products and expression by read through of a weak stop codon are indicated on the figure.

Furovirus

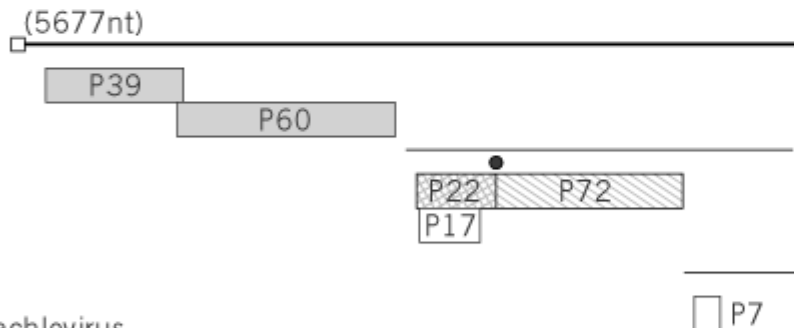


Hordeivirus

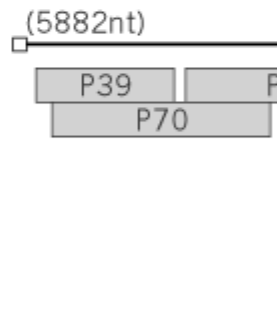


Luteovirus

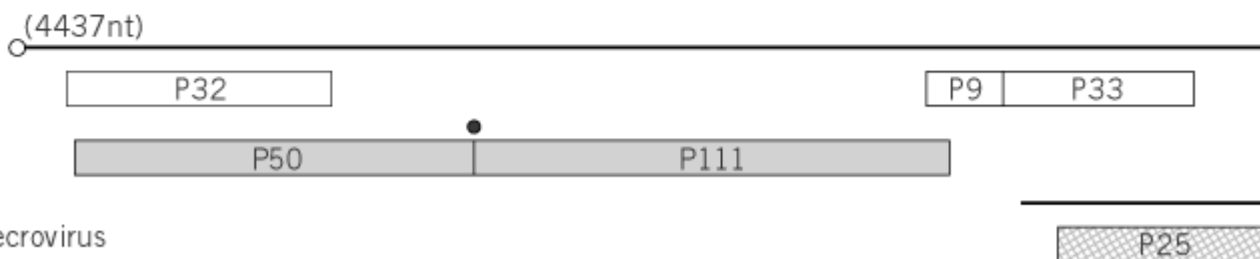
Subgroup I



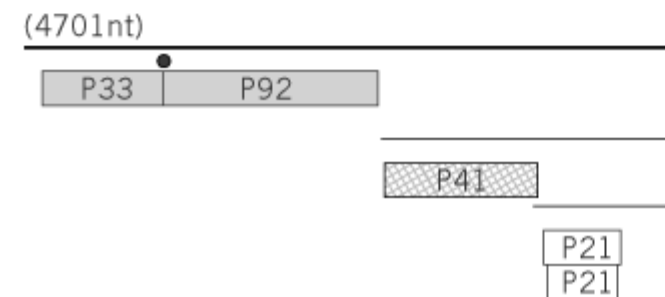
Subgroup II



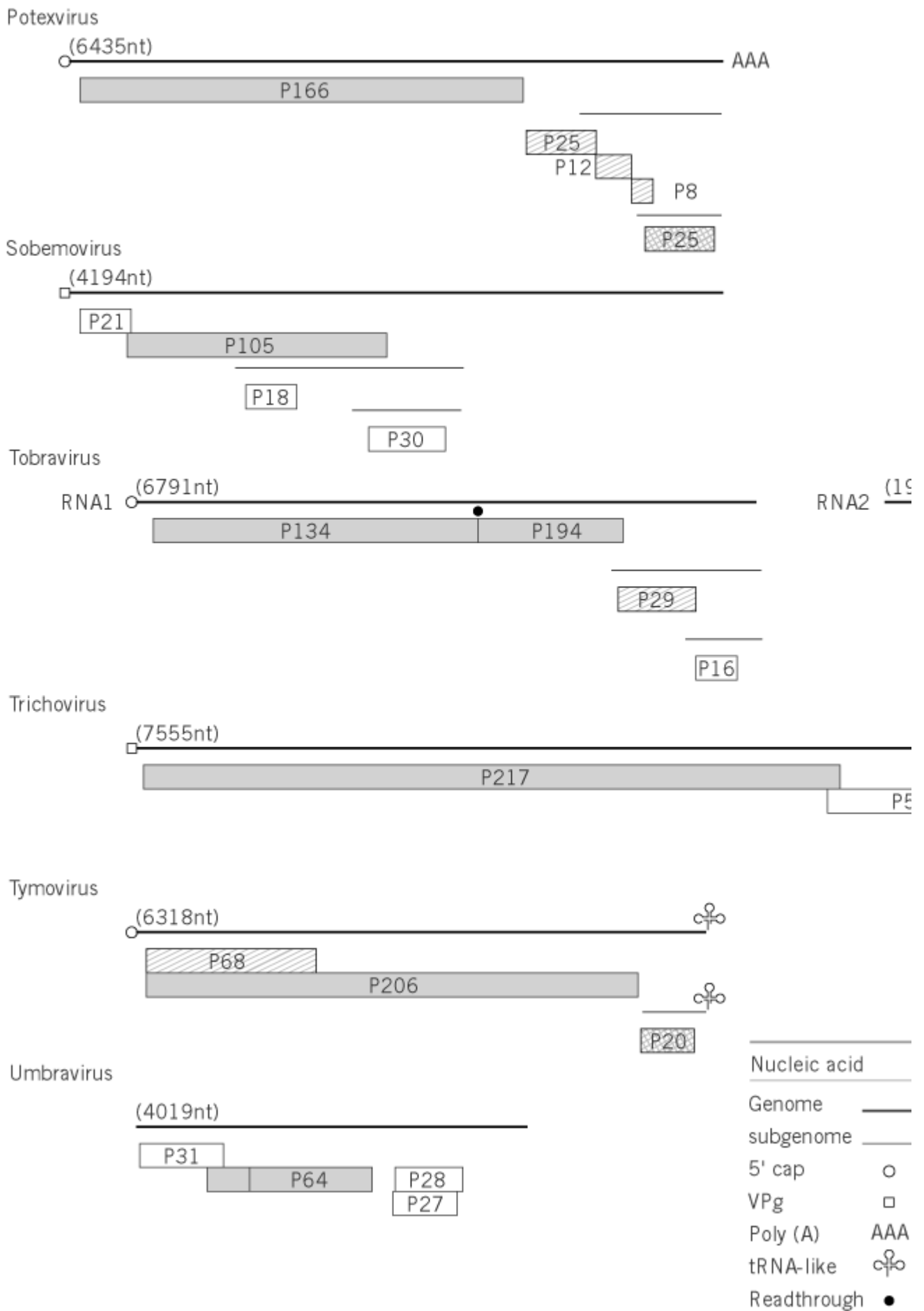
Machlovirus



Necrovirus



**Figure 1.** (*continued*) Genome organizations of groups of plant viruses. The illustration is of the type member of the group (Table 1) except subgroup III geminiviruses, African cassava mosaic virus; nucleorhabdovirus, sonchus yellow net virus; nepovirus, tomato blackring virus; carlavirus, potato virus M; dianthovirus, red clover necrotic mottle virus; umbravirus, groundnut rosette virus. The genome is represented by two concentric circles for double-stranded circular genomes; one circle for single-stranded circular genomes; a straight line for linear genomes; subgenomic mRNA are shown as thin lines. The size of the genome is given in nucleotides (nt or base pairs, bp). The terminal features of the genomes, when they are known, are listed on the figure. The gene products are indicated by boxes, with the directions of those of ambisense genomes given by an arrow. The molecular size of the gene product is  $\times 10^{-3}$ . The symbols for functions of gene products and expression by read through of a weak stop codon are indicated on the figure.



## 2.2. Double-stranded DNA Viruses

The genomes of the two genera of double-stranded DNA viruses are similar in that they are circular molecules of about 7.2 to 8.0 kbp, with one or more discontinuity in each strand; these discontinuities are described in more detail in [Cauliflower mosaic virus](#). The viral DNA of both these genera is transcribed asymmetrically to give a more-than-genome length RNA. The genomes of the badnaviruses encode three or four proteins. The largest badnavirus open reading frame is expressed as a **polyprotein** (Fig. 1), which is processed to yield the viral coat protein (CP), a [carboxyl proteinase](#) (Pro), and reverse transcriptase plus ribonuclease H (RT + RH) (the enzymes involved in reverse transcription). The three open reading frames of most badnaviruses are expressed from the genome-length RNA transcript by unusual mechanisms involving ‘*ribosomal shunting*’ and leaky translation initiation (see (5) for a review).

## 2.3. Single-stranded DNA Viruses

Geminiviruses have circular single-stranded DNA molecules of about 2.7 to 3.0 kb. The three genera of the Geminiviridae have similarities in their genome organizations and expression, but there are differences. Subgroups I and II have monopartite genomes, with members of subgroup I infecting monocotyledonous plants and those of subgroup II dicotyledons. The genome of most subgroup III members (which infect dicotyledonous plants) is divided between two DNA molecules. Members of subgroups I and II are transmitted by leafhoppers, those of subgroup III by whiteflies. In all cases, the genomes are transcribed bilaterally from a characteristic region near the mapping zero point (termed the *common region* for subgroup III members), the transcripts terminating on the opposite side of the genome. Thus, one transcript is from the virus-sense DNA and the other from the complementary-sense DNA. Each of these transcripts expresses one or more proteins (Fig. 1) through mechanisms that involve both different transcript origins and splicing. Molecular aspects of geminiviruses are reviewed in (6, 7).

## 2.4. Double-stranded RNA Viruses

The genome organizations of plant viruses with double-stranded RNA genomes in many cases resemble those of animal-infecting reoviruses. Members of the plant Reoviridae are insect-transmitted, and the virus multiplies in the vector. The genomes of these viruses are divided between 10 or 12 segments of about 2.5 to 0.8 kb that, in most cases, are monocistronic (8). These encode the various proteins that make up the virion structure, together with various nonstructural proteins, including the viral polymerase.

The *cryptic* viruses (Partitiviridae) have no known vectors and are spread by the vegetative propagation of their hosts. Their genomes have only two monocistronic genome segments, one encoding the viral polymerase and the other the coat protein (9).

Each genomic RNA of these double-stranded RNA viruses has a conserved terminal oligonucleotide sequence that is thought to be genus-specific (8).

## 2.5. Single-stranded (–)-Sense RNA Genomes

Two of the three families of plant viruses with (–)-sense RNA genomes resemble viruses that infect animals, and they have many common genome organizational features. They also replicate in their insect vectors. The best characterized rhabdovirus, *Sonchus yellow net virus*, belongs to the *Nucleorhabdovirus* genus and has a genome organization identical to animal rhabdoviruses, except that it has an extra gene product (sc4) (Fig. 1) (10). The genome length (–)-sense RNA is transcribed into monocistronic subgenomic mRNA, molecules.

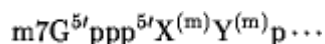
The *Tospoviruses* resemble the animal *Phlebovirus* genus of the *Bunyavirus* family in having three genome segments, the largest (L RNA, Fig. 1) of which is (–) sense and the other two having an *ambisense* expression strategy. In this, the medium and smallest RNA segments of *Tospoviruses* (M and S RNA) each encode one protein from the 5′ end of the viral-sense RNA and another protein from the 5′ end of the complementary RNA strand (ie, the open reading frame is at the 3′ end of the viral-sense RNA). The ambisense strategy is taken even further in the *Tenuiviruses*, which have their



largest genome segment (RNA1, Fig. 1) as (–)-sense RNA and the other three (RNAs 2, 3, and 4) with the ambisense arrangement. In both the Tospoviruses and Tenuiviruses, the largest RNA encodes the RNA-dependent **RNA polymerase**. Molecular aspects of tospoviruses are reviewed in (11) and those of tenuiviruses in (12).

## 2.6. Single-stranded (+)-Sense RNA Genomes

This genome, the most common type found in plant viruses, acts as an mRNA. A range of structures are found at the ends of these RNA (indicated in Fig. 1). Many viruses have a methylated blocked 5' terminal group



where  $X^{(m)}$  and  $Y^{(m)}$  are methylated bases. Others have a virus-encoded protein covalently linked to the 5' end of the genomic RNA. The linkage of this protein, the *Vpg*, has been identified in some cases and, for cowpea mosaic, has been shown to be a phosphodiester bond between the b-hydroxyl group of a [serine](#) residue at the N-terminus of the *Vpg* and the 5'-terminal uridine residue of the two genomic RNA. At the 3' end, some viruses have a **polyadenylate** sequence resembling that of many eukaryotic mRNAs. Other viruses have 3' sequences that fold to resemble [transfer RNA](#) (tRNA). These structures will accept amino acids; eg, that of TMV accepts histidine, BMV and CMV accept tyrosine, and turnip yellow mosaic tymovirus accepts valine. Barley stripe mosaic hordeivirus has a 3'-terminal tRNA-like structure (accepting tyrosine) and an internal poly(A) tract (Fig. 1). The functions of these 3' structures are unknown but thought to be involved with replication.

The 5' noncoding regions of several plant viruses have been shown to have translational enhancing activity. That of TMV, the W sequence, is discussed in [Tobacco mosaic virus](#). The 3' noncoding regions of some viruses fold as **pseudoknots** with unknown function.

An apparently wide range of genome organizations and expression is found in viruses with (+)-sense RNA genomes (Fig. 1), but, if we take into consideration the limitations of eukaryotic ribosome translation, certain basic features emerge. There are two predominant methods by which the four or more gene products are expressed from these viral genomes. In the first, exemplified by carmoviruses, luteoviruses, potexviruses (see Fig. 1), and TMV, the 5' cistron is expressed from the virion RNA, and most of the downstream cistrons are expressed from subgenomic RNAs that are transcribed from the (–)-strand replication intermediate. There are some subtleties in the expression of the 5' genetic information in that, in eg, carmoviruses and TMV, the first two or three open reading frames are separated by weak stop codons that are read through. A variant of this strategy is to **frameshift** between the first and second open reading frame of, eg, luteoviruses, when ribosomes change reading frame from one **cistron** to that of an overlapping cistron. In both read through and frame shift, the larger product is formed to a much lesser extent than the 5' product. The division of the genome into monocistronic units is taken even further in many genera of plant viruses that have multipartite or divided genomes. The genome is divided into two or three segments that are encapsidated separately. For example, members of the Bromoviridae (Alfamovirus, Bromovirus, and Cucumovirus) have their genomes divided into three segments. The larger two segments are monocistronic, whereas the smallest segment is bicistronic but expresses the 3' cistron from a subgenomic RNA.

The second predominant method by which viral gene products are expressed is by having a single open reading frame on the genomic RNA that is translated to give a polyprotein. This is processed to the functional proteins by one or more proteinases encoded within the polyprotein. Potyviruses, Sequiviruses, and Waikaviruses are examples of viruses with a single genome segment that follow this strategy. Comoviruses, Nepoviruses, and Bymoviruses have their genomes divided into two segments, each of which encodes a polyprotein. The processing pathways can be complex, as shown by Comoviruses (Fig. 1).

## 2.7. Genome Replication

The replication of the various type of genome is discussed in [Virus infection, plant](#).

### Bibliography

1. F. A. Murphy, C. M. Fauquet, D. H. L. Bishop, S. A. Ghabrial, A. W. Jarvis, G. P. Martelli, M. A. Mayo, and M. D. Summers (1995) *Virus Taxonomy: Classification and Nomenclature of Viruses*, Springer-Verlag, Wien.
2. D. L. D. Caspar and A. Klug (1962) Cold Spring Harbor Symp. Quant. Biol. **27**, 1–24.
3. R. I. B. Francki, T. Hatta, G. Boccardo, and J. W. Randles (1980) *Virology* **101**, 233–241.
4. R. Hull (1976) *Adv. Virus Res.* **20**, 1–32.
5. R. Hull (1996) *Ann. Rev. Phytopathol.* **34**, 275–297.
6. D. M. Bisaro (1996) In *DNA Replication in Eukaryotic Cells*, (M. de Pamphelis, ed.), Cold Spring Harbor Laboratory Press, New York, pp. 833–854.
7. S. G. Lazarowitz (1992) *Crit. Rev. Plant Sci.* **11**, 327–349.
8. I. Uyeda, I. Kimura, and E. Shikata (1995) *Adv. Virus Res.* **45**, 249–279.
9. G. P. Accotto, C. Marzachi, E. Luisoni, and R. G. Milne (1990) *J. Gen. Virol.* **71**, 433–437.
10. L. A. Heaton, B. I. Hillman, B. G. Hunter, D. Zuidema, and A. O. Jackson (1989) *Proc. Natl. Acad. Sci. USA* **86**, 8665–8668.
11. C. S. Schmaljohn (1996) In *Fundamental Virology*, 3rd ed. (B. N. Fields, P. M. Knipe, and P. M. Howley, eds), Lippincott-Raven, Philadelphia, pp. 649–673.
12. B. C. Ramirez and A.-L. Haenni (1994) *J. Gen. Virol.* **75**, 467–475.

### Suggestion for further reading

13. R. E. F. Matthews (1991) *Plant Virology*, 3rd ed., Academic Press, San Diego, CA.

## Virusoids

Virusoids are a particular class of plant **satellite RNA** (Fig. 1). Satellite RNAs are always functionally dependent on specific [helper viruses](#) and are encapsidated by the coat [protein](#) of these helper viruses. They can be regarded as molecular parasites and in most cases reduce the titer and attenuate the symptom expression of their supporting viruses. Satellite RNAs do not have extensive sequence similarity with the RNA of the helper virus.

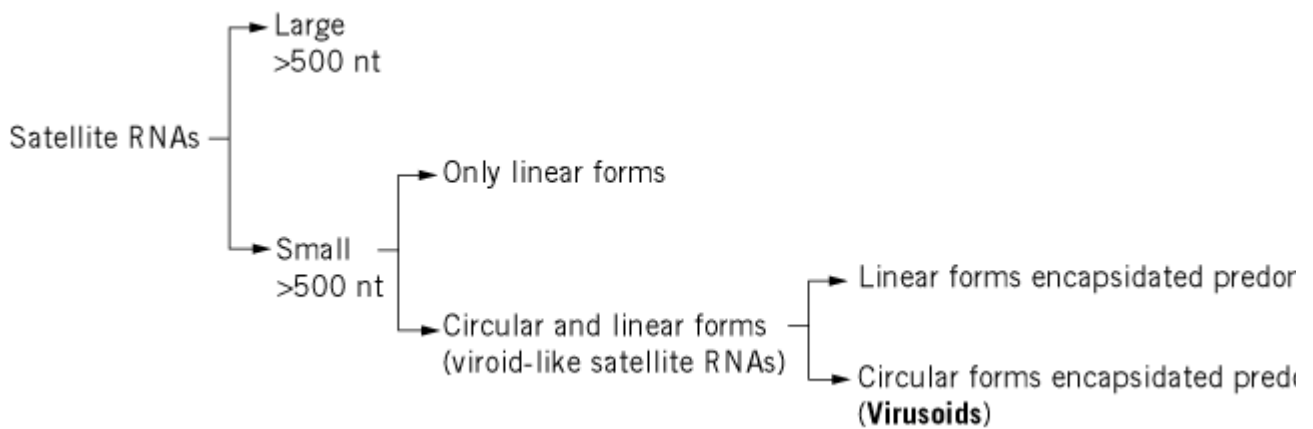


Figure 1 shows the groups into which plant satellite RNAs can be subdivided. Large satellite RNAs can encode proteins, whereas small satellite RNAs do not. Viroid-like satellite RNAs (Table 1) are single-stranded molecules present in infected tissue in both circular and linear forms, are structurally similar to [viroids](#), and most likely replicate by the same type of **rolling circle** mechanism. A property that is unique to the viroid-like satellite RNAs of sobemoviruses is that the circular form is the predominant RNA encapsidated by the coat protein of the helper virus, and on this basis the term virusoid has been proposed for them; in contrast, the linear form is predominantly encapsidated in viroid-like satellite RNAs of nepo- and luteoviruses (1). It is also worth noting that at least one plant viroid-like RNA has been found in association with an **homologous** DNA counterpart forming what has been called a retroviroid-like element (2).

**Table 1. Virusoids and Other Viroid-like Satellite RNAs with Their Abbreviations, Genomic Accession Numbers of Typical Sequence Variants, Sizes, and Group to Which They Belong**

| Viroid-like Satellite RNA                   | Abbreviation  | Accession | Size (nt) | Group       |
|---|---------------|-----------|-----------|-------------|
| Lucerne transient streak virusoid           | vLTSV         | X01984    | 322–324   | Sobemovirus |
| Solanum nodiflorum mottle virusoid          | vSNMoV        | J02386    | 377       | Sobemovirus |
| Subterranean clover mottle virusoid         | vSCMoV        | M33000    | 322, 388  | Sobemovirus |
| Velvet tobacco mottle virusoid              | vVTMoV        | J02439    | 365–366   | Sobemovirus |
| Rice yellow mottle virusoid                 | vRYMV         | AF039909  | 220       | Sobemovirus |
| Cereal yellow dwarf virus RPV satellite RNA | SCYDV-RPV RNA | M63666    | 322       | Luteovirus  |
| Tobacco ringspot virus satellite RNA        | sTRSV RNA     | M14879    | 359–360   | Nepovirus   |
| Arabis mosaic virus satellite RNA           | sArMV RNA     | M21212    | 300       | Nepovirus   |
| Chicory yellow mottle virus satellite RNA   | SCYMoV RNA    | D00721    | 457       | Nepovirus   |

---

All known viroid-like satellite RNAs contain [ribozyme \(catalytic RNAs\)](#) domains in one or in both polarity strands. The ribozymes of virusoids and sCYDV-RPV RNA are of the hammerhead type and can be adopted by either both polarity strands (vLTSV and sCYDV-RPV RNA) or only their plus polarity strands (vSNMoV, vSCMoV, vVTMoV, and vRYMV). Viroid-like satellite RNAs from nepoviruses can form hammerhead and hairpin (or paperclip) ribozymes in their plus and minus polarity strands, respectively. These ribozymes most probably play a role in replication: hairpin structures mediate both RNA self-cleavage and self-ligation ([3](#)), whereas hammerhead structures mediate self-cleavage only ([4-6](#)).

It has been proposed that viroid-like satellite RNAs have a common **phylogenetic** origin with viroids and that they could be molecular fossils of the [RNA world](#) ([7](#)). The RNA of human hepatitis delta virus shares some properties with viroid-like satellite RNAs, including circular structure, functional dependence on a helper virus, and presence of ribozymes in both polarity strands. One key question to be addressed is the characterization of the determinants responsible for the specificity of the interaction between each viroid-like satellite RNA and its helper virus.

### Bibliography

1. G. Bruening, P. A. Feldstein, J. M. Buzayan, H. van Tol, B. K. Passmore, J. deBear, G. R. Gough, P. T. Gilham, and F. Eckstein (1991) In *Viroids and Satellites: Molecular Parasites at the Frontier of Life* K. Maramorosch, ed., CRC Press, Boca Raton, FL, pp. 141–158.
2. J. A. Daròs and R. Flores (1995) *Proc. Natl. Acad. Sci. USA* **92**, 6856–6860.
3. J. M. Buzayan, W. L. Gerlach, and G. Bruening (1986) *Nature* **323**, 349–353.
4. G. A. Prody, J. T. Bakos, J. M. Buzayan, I. R. Schneider, and G. Bruening (1986) *Science* **231**, 1577–1580.
5. A. C. Forster, and R. H. Symons (1987) *Cell* **49**, 211–220.
6. R. Flores, C. Hernández, M. De la Peña, A. Vera, and J. A. Daròs. (2002) *Meth. Enzymol.* **341**, 540–552.
7. S. F. Elena, J. Dopazo, M. De la Peña, R. Flores, T. O. Diener, and A. Moya (2001) *J. Mol. Evol.* **53**, 155–159.

### Suggestions for Further Reading

8. G. Bruening, B. K. Passmore, H. van Tol, J. M. Buzayan, and P. A. Feldstein (1991) Replication of a plant virus satellite RNA: evidence favors transcription of circular templates of both polarities. *Mol Plant-Microbe Interact.* **4**, 219–225.
9. R. H. Symons (1997) Plant pathogenic RNAs and RNA catalysis. *Nucleic Acids Res.* **20**, 2683–2689.
10. R. H. Symons and J. W. Randles (1999) Encapsidated circular viroid-like satellite RNAs (virusoids) of plants. *Curr. Top. Microbiol. Immunol.* **239**, 81–105.
11. M. C. C. Lai (1995) The molecular biology of hepatitis delta virus. *Annu. Rev. Biochem.* **64**, 259–286.
12. J. M. Taylor (1999) Human hepatitis delta virus: an agent with similarities to certain satellite RNAs of plants. *Curr. Top. Microbiol. Immunol.* **239**, 107–122.

## Water

Water is common yet indispensable to life. It is the most abundant constituent of most living organisms, comprising more than 60% of the weight of an adult man. A living creature constructs its own structure, having a definite boundary with its environment but not isolated from it. In physical terms, it is a stationary, open system that maintains its structure and functions by using materials and energy incorporated from the environment. For metabolism to function, a liquid medium is required. Many of the biomolecules involved in metabolism must have [polar](#) groups, and the liquid medium needs to dissolve these substances. Only water satisfies this requirement and exists in large amounts on the earth. Various biological effects of water arise from the characteristics of its interactions with biomolecules, i.e., their [hydration](#). According to their physical nature, hydration phenomena are grouped into two categories: (i) [hydrophilic](#) hydration of ionic and nonionic polar groups, and (ii) **hydrophobic** hydration of [nonpolar](#) groups ([1](#), [2](#)). The structural and physical properties of pure water are described here to help in understanding the role of water in biological phenomena.

Water exhibits three phases—gas, liquid and solid—around its triple point at a pressure of 610.6 pascal (Pa) and a temperature of 273.16 K. Under other conditions, two phases coexist at the phase boundary between them. The gas-liquid coexistence curve ends at the critical point pressure of 22.12 MPa, temperature of 647.3 K, and density of  $0.322 \text{ g cm}^{-3}$ , no liquid phase appears at higher temperatures. Unlike most other liquids, water has a solid-liquid coexistence curve with a negative slope near the normal pressure of 1 atm (0.101325 MPa), at which point an increase in pressure leads to a decrease in the melting (fusion) temperature. At normal pressure, water can be supercooled to  $-39.5^\circ\text{C}$ , but this is only a metastable state, not an equilibrium, and the solution will eventually freeze. Various physical properties of water are listed in [Table 1](#).

**Table 1. Physical Properties of Water<sup>a</sup>**

|   |                         |
|---|-------------------------|
| Molecular weight, $M_w$                         | 18.01529                |
| Bond length, $r_{\text{OH}}$ / pm               | 95.72                   |
| Bond angle, HOH / deg                           | 104.52                  |
| Thermal deBroglie wavelength, / pm              | 23.82                   |
| Effective diameter, $d_e$ / pm                  | 275                     |
| Normal vibration (gas)                          |                         |
| symmetric stretching, $n_1$ / $\text{cm}^{-1}$  | 3656.65                 |
| deformation, $n_2$ / $\text{cm}^{-1}$           | 1594.59                 |
| asymmetric stretching, $n_3$ / $\text{cm}^{-1}$ | 3755.79                 |
| Dipole moment, (gas) $m_g$ / D                  | 1.834                   |
| (liquid) $m_l$ / D                              | 2.45                    |
| Polarizability, $a$ / $\text{m}^3$              | $1.470 \times 10^{-30}$ |
| Density (weight), $r_w$ / $\text{g cm}^{-3}$    | 0.997045                |
| (number), $r_n$ / $\text{nm}^{-3}$              | 33.329                  |
|   | 36.29                   |

|   |                         |
|---|-------------------------|
| Packing density, $h_p$ / %  |                         |
| Coordination number, $n$  | 4.4                     |
| Dielectric constant, $\epsilon$   | 78.54                   |
| Refractive index, $n_D$   | 1.33287                 |
| Triple point pressure, $P_t$ / Pa   | 610.6                   |
| temperature, $T_t$ / K  | 273.16                  |
| Critical point temperature, $T_c$ / K                                       | 647.3                   |
| pressure, $P_c$ / MPa   | 22.12                   |
| density, $r_c$ / g cm <sup>-3</sup>   | 0.322                   |
| Heat of fusion, $DH_f(0^\circ\text{C})$ / kJ mol <sup>-1</sup>              | 6.008                   |
| Heat of vaporization, $DH_v$  |                         |
| (0°C) / kJ mol <sup>-1</sup>  | 45.049                  |
| (25°C) / kJ mol <sup>-1</sup>   | 43.991                  |
| (99.974°C) / kJ mol <sup>-1</sup>   | 40.66                   |
| Specific heat, $C_p$ / J K <sup>-1</sup> g <sup>-1</sup>                    | 4.1796                  |
| Vapor pressure, $P_v$ / kPa   | 3.1675                  |
| Surface tension, $g$ / mJ m <sup>-2</sup>                                   | 71.96                   |
| Thermal expansivity, $a_p$ / K <sup>-1</sup>                                | $2.5721 \times 10^{-4}$ |
| Isothermal compressibility, $k_T$ / GPa <sup>-1</sup>                       | 0.452472                |
| Translational diffusion coefficient, $D_t$ / m <sup>2</sup> s <sup>-1</sup> | $2.14 \times 10^{-9}$   |
| Rotational diffusion coefficient, $D_r$ / s <sup>-1</sup>                   | $6.06 \times 10^{10}$   |
| Viscosity, h/mPa s  | 0.8904                  |
| Ionic product, mol <sup>2</sup> l <sup>-2</sup>                             | $1.0 \times 10^{-14}$   |

---

<sup>a</sup> T = 25°C and normal pressure is assumed unless otherwise noted.

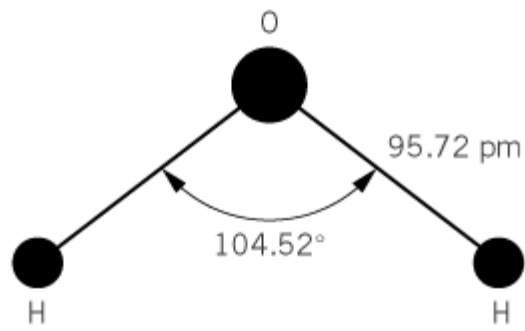
## 1. The Water Molecule

The water molecule is small, composed of one oxygen and two hydrogen atoms, and has an average molecular weight of 18.01529. There are three isotopes for both oxygen and hydrogen atoms: <sup>16</sup>O, <sup>17</sup>O, and <sup>18</sup>O for oxygen, and <sup>1</sup>H (H, or protium), <sup>2</sup>H (D, or deuterium), and <sup>3</sup>H (T, or tritium) for hydrogen. Tritium is a [radioisotope](#) with a half-life of 12.33 years, so there are nine stable isotopic species of water. The natural abundances of the major species are 99.728% <sup>1</sup>H<sub>2</sub><sup>16</sup>O, 0.200% <sup>1</sup>H<sub>2</sub><sup>18</sup>O, 0.040% <sup>1</sup>H<sub>2</sub><sup>17</sup>O, and 0.032% <sup>1</sup>H<sup>2</sup>H <sup>16</sup>O. The molecular weight listed in Table [1](#) is the average weight, with respect to these proportions.

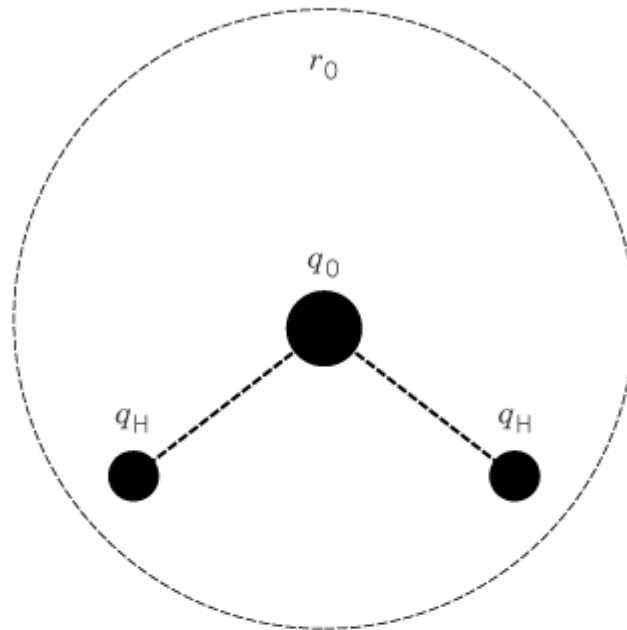
Various unique properties of water originate from the specific structure of the water molecule. It has a bent structure (Fig. [1\(a\)](#)), and the HOH bond angle is 104.52°. The O-H bond length is 95.72 pm

(or 0.9572 Å) for the lowest-energy structure but fluctuates thermally about 10% at room temperature. Two of the six outer-shell electrons of the O atom are involved in the formation of two O-H bonds, and the remaining four electrons make two lone pairs. As the electronegativity of the O atom is significantly greater than that of the H atom, the electron cloud of the O-H bonds is pulled closer to the O atom. As a result, a deficit and an excess of electron density are produced near the H and O atoms, respectively. This is the reason the water molecule has a fairly large electric dipole moment, 1.834 debye (D, or  $6.117 \times 10^{-30}$  Coulomb m), for its small size. In molecular models for computer simulations, the polarized distribution of electrons in the water molecule is approximately taken into account by assuming that the H and O atoms have positive and negative partial point-charges, respectively (Fig. 1(b)) (3).

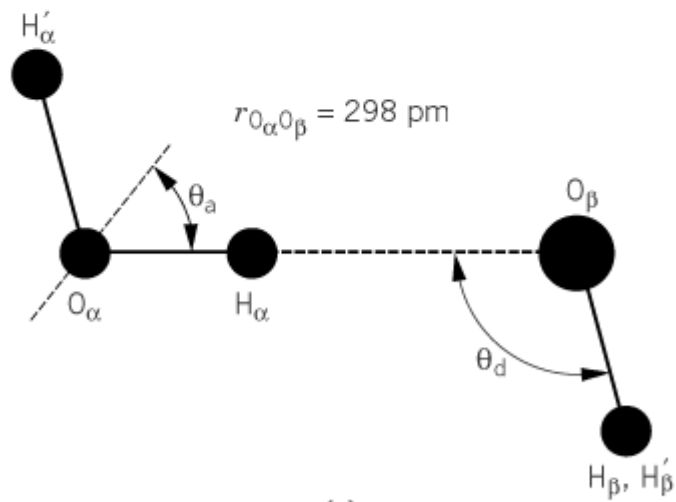
**Figure 1.** The water molecule: (a) the lowest energy structure; (b) a molecular model, TIP3P, where  $r$  and  $\epsilon$  are the position and the depth, respectively, at the minimum of the Lennard-Jones 6-12 potential, and  $q$  is the partial charge in the unit of elementary electric charge (3); (c) the most stable linear conformation of the water dimer, where  $q_d$  and  $q_a$  are the donor and acceptor angles, respectively; (d) the normal modes of vibration:  $n_1$ , symmetric stretching;  $n_2$ , deformation; and  $n_3$ , asymmetric stretching.



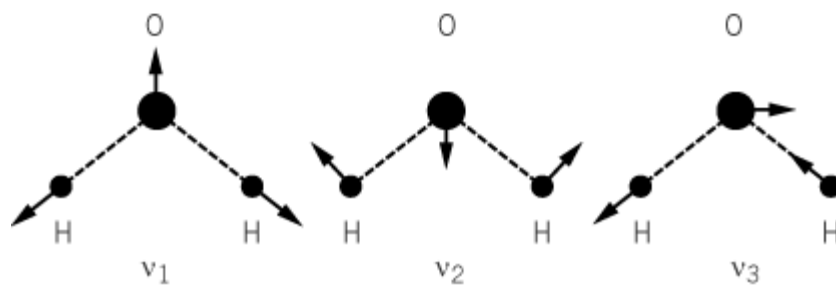
(a)



(b)



(c)





In the gas state, two water molecules form various types of dimers making one or two [hydrogen bonds](#); the most stable dimer is depicted in Figure 1(c) (4). The hydrogen bond involves primarily Coulombic interactions between partial charges on the donor and acceptor molecules. The binding energy of such a hydrogen bond is much greater than thermal energy ( $RT$ ). It is estimated to be  $22.6 \text{ kJ mol}^{-1}$  (5) for the linear conformation in which the three atoms  $\text{O}_a\text{--H}_a\text{--}\frac{1}{2}\text{O}_b$  are in a straight line, as in Figure 1(c). The angles between the bond axis  $\text{O}_a\text{--O}_b$  and the bisectors of the angles  $\text{H}'_a\text{--O}_a\text{--H}_a$  and  $\text{H}_b\text{--O}_b\text{--H}_b$  are called the donor angle  $q_d$  and the acceptor angle  $q_a$ , respectively. Although the hydrogen bond has a directional propensity, allowances for the deviation of  $q_d$  and  $q_a$  from their optimum values are fairly large: about  $25^\circ$  for  $q_d$  and as much as  $100^\circ$  for  $q_a$  (6). Both the directional propensity and the flexibility of the hydrogen bond significantly affect the structural and thermodynamic properties of liquid water and the hydration of biomolecules.

A water molecule has nine degrees of freedom: three translational, three rotational, and three intramolecular vibrational. The three vibrational modes are those of symmetric stretching ( $n_1 = 3656.65 \text{ cm}^{-1}$ ), deformation ( $n_2 = 1594.59 \text{ cm}^{-1}$ ), and asymmetric stretching ( $n_3 = 3755.79 \text{ cm}^{-1}$ ) (Fig. 1(d)). These three frequencies depend strongly on the interactions of the water molecule with neighboring molecules, especially on hydrogen bonding.

## 2. Ice

At normal pressure, water freezes at  $0^\circ\text{C}$  ( $273.15 \text{ K}$ ) to form a crystal called ice- $I_h$  (Fig. 2). On being frozen at a temperature lower than  $-80^\circ\text{C}$  in the gas phase, water forms a crystal named ice- $I_c$ , where  $c$  indicates cubic. The crystal lattice of ice- $I_h$  is hexagonal and has a tridymite-like structure: Each O atom of the water molecules is surrounded tetrahedrally by four O atoms of the nearest neighbor molecules, located at a distance of  $275 \text{ pm}$  ( $2.75 \text{ \AA}$ ) from it. There are twelve second-nearest neighbor molecules, at a distance of  $450 \text{ pm}$  ( $4.50 \text{ \AA}$ ). The number of nearest neighbors in ice- $I_h$ —four—is much smaller than the 12 present in a hexagonal closest-packed structure. The resulting crystal has a large cavity parallel to its  $c$  axis, producing a loosely packed structure. In fact, the density ( $\rho$ ) of ice- $I_h$  at  $0^\circ\text{C}$ , ( $0.91671 \text{ g cm}^{-3}$ ) is only 91.7% of that of liquid water at the same temperature ( $0.99984 \text{ g cm}^{-3}$ ) and is also smaller than the density at the boiling point temperature of  $99.974^\circ\text{C}$  ( $0.95837 \text{ g cm}^{-3}$ ). As the water molecule can be approximated by a sphere with a diameter of  $275 \text{ pm}$  ( $2.75 \text{ \AA}$ ), the packing density of ice at  $0^\circ\text{C}$  can be estimated to be 33.37%, which is only 45% of that for the closest-packed structure (74.05%). Every water molecule in ice- $I_h$  provides two donor sites and two acceptor sites for hydrogen bonds, and the neighboring molecules are hydrogen-bonded to each other. Thus, the hydrogen bond plays an essential role in the tetrahedral structure of ice.

**Figure 2.** Crystal structure of ice- $I_h$ . The filled circles are oxygen atoms and the lines connecting them are hydrogen bonds between water molecules; the hydrogen atoms are not indicated. Thick lines are drawn to indicate hexagonal rings typical of the ice- $I_h$  structure.



Under pressures greater than about 200 MPa, ice exhibits various phases, named II–IX (7). Ice-II to ice-VI and ice-IX have the same number of first-nearest neighbors as ice-I<sub>h</sub>, but the distances are slightly longer. Furthermore, the angles between the two axes of neighboring hydrogen bonds deviate greatly from the tetrahedral angle of 109.47°, and the distances to the second-nearest neighbors are shorter than in ice-I<sub>h</sub>. Ice-VII and ice-VIII, which are made under the high pressure limit, both have a much greater density of  $\approx 1.65 \text{ g cm}^{-3}$ . They have structures in which two identical lattices are overlapped, filling in each other's cavities. Each molecule in the two ices has an equal distance of 286 pm to the first- and second-nearest neighbors.

In crystals of both ice-I<sub>h</sub> and ice-I<sub>c</sub>, all of the O atoms are located regularly, but the position of each H atom is not specified definitely, even at 0°K, which produces a residual **entropy**. There are two potential minima for a proton between two O atoms of water molecules hydrogen-bonded to each other. If a proton moves to another minimum, a pair of ions, HO<sup>-</sup> and H<sub>3</sub>O<sup>+</sup>, are created. Successive transfer of a proton from an H<sub>3</sub>O<sup>+</sup> ion to an HO<sup>-</sup> ion or to neighboring molecules brings about electric conduction. On the other hand, if a proton moves to one of the two hydrogen-bond acceptor sites of the same molecule, it results in an effective rotation of the molecule. Proton jumps of this type in an external electric field contribute to the dielectric permittivity of ice, and its large value of 94 at -2°C in the low frequency region is caused by this mechanism.

### 3. Liquid Water

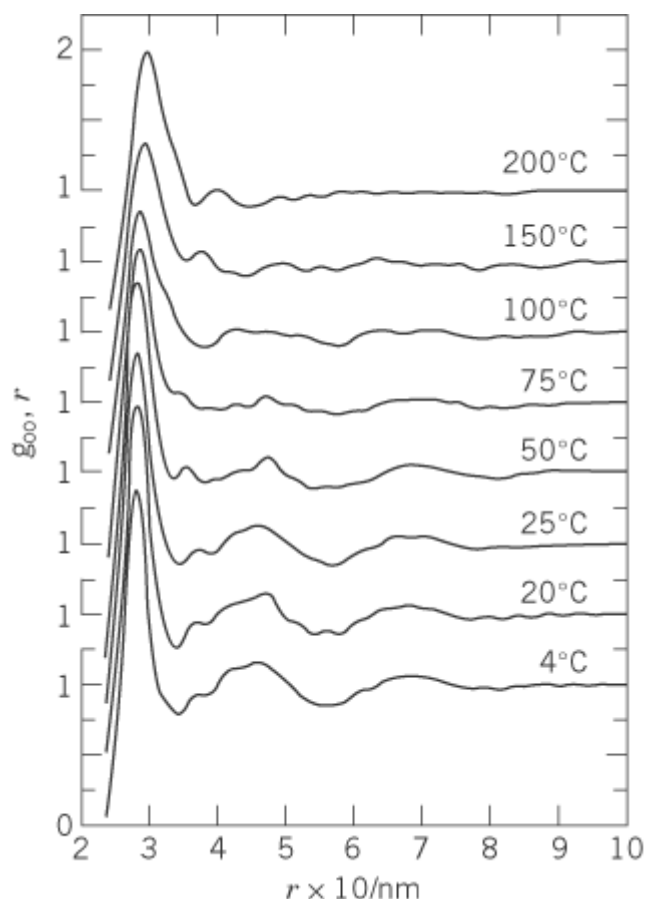
Liquid water has various structural and physical properties that differ from those of other liquids. They have significant effects on the structure and function of proteins and nucleic acids.

#### 3.1. Structure

Using the same approximation made previously for ice, the packing density of water at 25°C is estimated to be only 36.29 %, which is significantly lower than that of nonpolar organic liquids. This suggests that liquid water has a loosely packed structure similar to that of ice. As the liquid molecules continually change their positions, the characteristics of their spatial distribution can be described only statistically. If the ordinary definition of the word “structure” for solid substances is extended to describe the spatial correlation among the molecules or atomic groups in the system

considered, the word can also be applied to liquids. The simplest function to characterize liquid structure in the above sense is the radial distribution function,  $g_{ij}(r)$ . It is defined as the normalized density of atom  $j$  at a distance  $r$  from the center of an atom  $i$ . Experimentally, the radial distribution function of the O atoms of water,  $g_{OO}(r)$ , can be determined from measurement of [X-ray scattering](#) by water (Fig. 3) (8). Furthermore, measurement of the **neutron scattering** profile for mixtures of light and heavy water can determine all of the three radial distribution functions,  $g_{OO}(r)$ ,  $g_{OH}(r)$ , and  $g_{HH}(r)$  (9). The first maximum of  $g_{OO}(r)$  at  $r \approx 280$  pm in Figure 3 is due to the nearest-neighbor molecules. At temperatures lower than about  $50^\circ\text{C}$ , a small but obvious maximum exists at  $r = 450$  pm, corresponding to the distance of the second-nearest neighbors in ice- $I_h$ . Integration of  $g_{OO}(r)$  from  $r = 0$  to the distance of its first minimum yields an estimate of 4.4 for the number of the nearest neighbors,  $n$  (the coordination number), at  $25^\circ\text{C}$ . This is only 10% greater than that for ice and much smaller than the mean value—about 10—for nonpolar liquids. The first peak broadens with increasing temperature, as seen in Figure 3, but  $n$  increases only slightly, to 4.9 at  $83^\circ\text{C}$ . On the other hand, both the minimum and the second peak nearly disappear at  $100^\circ\text{C}$ . These results indicate that a short-range structure like that of ice is nearly preserved in liquid water, but the long-range spatial correlation almost disappears at distances further than those of the second-nearest neighbors.

**Figure 3.** The radial distribution function,  $g_{OO}(r)$ , of liquid water and its temperature dependence (8).



To explain these experimental facts, various models of liquid water have been proposed. The continuum model assumes that the water structure is generated by distorting the tetrahedral structure of ice in various ways (10). The mixture model regards water as an ensemble of water clusters with

different sizes, one example of which is the flickering cluster model (11). In the interstitial model, water is supposed to have a structure in which some of the water molecules enter the ice-like cavities, and the fraction of these molecules increases with increasing temperature (12). Analyses based on these models have shown that every model can reproduce the experiment observations fairly well, which means that it is difficult to define the correct structure of water only by experiments. With advances in the performance of computers, **Monte Carlo** and [molecular dynamics](#) simulations on systems composed of hundreds to thousands of water molecules have provided detailed information about the microscopic structure of water. Molecular models of water, such as ST2, SPC, TIP3P, TIP4P, and so on, are used in these calculations (3). The force field of the water molecule is assumed to be given by the sum of the repulsive forces due to outer-shell electrons, the attractive dispersion forces (see [van der Waals Interactions](#)), and the electrostatic Coulomb forces due to partial point-charges (see [Electrostatic Interactions](#)).

These theoretical studies revealed that in addition to the basic unit of the hydrogen-bond network in ice of a hexagonal ring of water molecules (Fig. 2), a similar number of pentagonal rings also exists in water. Hydrogen-bonded to each other, these rings form water clusters in which formation and rupture of hydrogen bonds continually occur. In most cases, these processes do not occur independently but are strongly correlated in nearby molecules; they occur through rearrangement of hydrogen bonds on a time-scale of picoseconds, and about 10–40 molecules participate in the collective motion (13).

### 3.2. Static and Dynamic Properties

Liquid water has a maximum density of  $0.9997 \text{ g cm}^{-3}$  at a temperature of  $3.98^\circ\text{C}$ . The observed decrease in density with increasing temperature near the freezing point is one of the unique characteristics of water. It is believed to arise from the decrease in volume due to partial collapse of the tetrahedral ice-like structure exceeding the increase in volume due to the general effect of thermal expansion caused by the increase in thermal energy. Compared to other liquids of molecules of a similar size, water has higher temperatures for melting and boiling and larger values for the dielectric permittivity, heat of fusion, heat of vaporization, specific **heat capacity**, and surface tension. All of these come about from water being an associative liquid, where constituent molecules interact with each other through hydrogen bonds. The hydrogen bond is much stronger than the van der Waals interaction which is the major cohesive force of nonpolar liquids. The large dielectric permittivity results from the combined effects of (i) a single water molecule having a large electric dipole for its size, (ii) nearby water molecules combining their dipoles by forming a locally regular structure, and in addition (iii) polarizing each other cooperatively to increase the magnitude of their dipole (14). The average dipole moment estimated for the water molecule in liquid at  $0^\circ\text{C}$ —2.45 D—is significantly greater than that in the gas state—1.834 D. This polar nature of water serves to increase the solubility of polar molecules essential to biological functions. The heat of fusion at the melting point,  $6.008 \text{ kJ mol}^{-1}$ , is only 13.3 % of the heat of vaporization at the melting point,  $45.05 \text{ kJ mol}^{-1}$ , and 14.8 % of that at the boiling point,  $40.66 \text{ kJ mol}^{-1}$ . This clearly shows that the cohesive energy of ice is mostly conserved in liquid water. The large specific heat capacity of water has an important biological significance of minimizing changes in the temperature of living organisms on changes in the ambient temperature. The large interfacial energy between water and a nonpolar medium is closely related to the **hydrophobic interaction** that occurs between nonpolar groups or molecules in water.

The intramolecular and intermolecular vibrational motions of water molecules in the liquid yields several broad bands in the spectra of infrared absorption and Raman scattering (see [Vibrational Spectroscopy](#)) (15). The peak frequency of the band for the stretching mode in the liquid,  $n_s$ , is lower than both  $n_1$  and  $n_3$  for the vapor. In contrast, the peak frequency for the deformation mode in liquid is shifted to a higher frequency. The values of  $n_s$  and  $n_2$  for liquid water are intermediate between their corresponding values for vapor and ice. These changes in frequency can be explained by the effect of hydrogen-bond formation in liquid water. Besides these intramolecular modes,

intermolecular vibrational modes are observed in the frequency range lower than  $1000\text{ cm}^{-1}$ : the constrained rotational vibration, or the libration, whose frequency spreads over the range of  $300\text{--}800\text{ cm}^{-1}$ , and the constrained translational vibrations with frequencies of  $60$  and  $170\text{ cm}^{-1}$ , which originate from the stretching and deformation vibrations of intermolecular hydrogen bonds, respectively. Rotational correlation times of the vector and tensor quantities are obtained from measurements of dielectric dispersion and nuclear magnetic resonance (NMR), respectively. The rotational correlation time for the dipole moment of water molecules in the liquid is estimated to be  $8.25\text{ ps}$  at  $25^\circ\text{C}$ . This agrees fairly well with the theoretical estimate of  $7.06\text{ ps}$  obtained by assuming that the water molecule is approximated by a sphere with a diameter of  $275\text{ pm}$  and applying the Stokes-Einstein formula (see [Hydrodynamic Volume](#)). This indicates that water molecules in the liquid are making random rotational motions on the picosecond time-scale.

### Bibliography

1. A. A. Rashin, ed. (1994) *Biophys. Chem.* **51**, 89-409.
2. K. Soda (1993) *Adv. Biophys.* **29**, 1-54.
3. W. L. Jorgensen and J. Tirado-Rives (1988) *J. Am. Chem. Soc.* **110**, 1657-1666.
4. T. R. Dyke and J. S. Muentzer (1974) *J. Chem. Phys.* **60**, 2929-2930.
5. L. A. Curtiss, D. J. Frurip, and M. Blander (1979) *J. Chem. Phys.* **71**, 2703-2711.
6. J. L. Finney, J. E. Quinn, and J. O. Baum (1990) In *Water Science Reviews* (F. Franks, ed.), Cambridge University Press, Cambridge, U. K., **1**, pp. 93-170.
7. F. Franks (1972) In *Water: A Comprehensive Treatise* (F. Franks, ed.), Plenum Press, New York, **1**, pp. 115-149.
8. A. H. Narten and H. A. Levy (1971) *J. Chem. Phys.* **55**, 2263-2269.
9. W. E. Thiesen and A. H. Narten (1982) *J. Chem. Phys.* **77**, 2656-2662.
10. J. D. Bernal (1964) *Proc. Roy. Soc.* **A280**, 299-322.
11. G. Nemethy and H. A. Scheraga (1962) *J. Chem. Phys.* **36**, 3382.
12. O. Ya. Samoilov (1946) *Zh. Fiz. Khim.* **20**, 1411.
13. I. Ohmine, H. Tanaka, and P. G. Wolynes (1988) *J. Chem. Phys.* **89**, 5852-5860.
14. S. L. Carnie and G. N. Patey (1982) *Mol. Phys.* **47**, 1129-1151.
15. K. M. Mizoguchi, Y. Hori, and Y. Tominaga (1992) *J. Chem. Phys.* **97**, 1961-1968.

### Suggestions for Further Reading

16. *Water: A Comprehensive Treatise* (1972-1982) (F. Franks, ed.), Plenum Press, New York, vols. **1-7**.
17. *Water Science Reviews* (1985-1990) (F. Franks, ed.) Cambridge University Press, Cambridge, U. K., vols. **1-5**.

### Wilson Plot

In [X-ray crystallography](#) the Wilson plot is used to determine the absolute scale of the diffracted intensities and to find the [temperature factor](#). This is based on the equation

$$\ln \frac{\overline{I(hk\ell)}}{\sum (f_i^0)^2} = \ln C - 2B \left( \frac{\sin \theta}{\lambda} \right)^2 \quad (1)$$

The average values of the X-ray intensities,  $\overline{I(hk\ell)}$ , in narrow ranges of  $\frac{\sin \theta}{\lambda}$  are divided by  $\sum (f_i^0)^2$  where the summation is over all atoms in the [unit cell](#) and  $f_i^0$  is the scattering factor of atom  $i$  at rest. The natural logarithm is plotted against  $(\frac{\sin \theta}{\lambda})^2$ , where  $\theta$  is the reflection angle and  $\lambda$  the X-ray wavelength. From the straight line obtained, the thermal parameter  $B$  of the isotropic temperature factor and the constant  $C$  can be derived. This constant  $C$  is a factor that relates the measured intensities to their absolute value:

$$I_{\text{meas}}(hk\ell) = C \times I_{\text{abs}}(hk\ell) \quad (2)$$

Intensities are on the absolute scale when the amplitudes of the [structure factor](#)

$|F(h\ k\ell)| = \sqrt{I(hk\ell)}$  are expressed in electrons.

### Suggestions for Further Reading

J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists* VCH, New York, Weinheim, Cambridge, p. 264.

## Wingless Signaling

The Wnt gene family encodes secreted **glycoproteins** that act as extracellular signals in a variety of biological processes. These proteins are approximately 400 amino acid residues long and are highly conserved at the sequence level within a pattern of 24 [cysteine](#) residues (1). Because of the conservation of Wnt genes and their importance, the elucidation of Wnt gene function has been of paramount interest for diverse disciplines. Wnt genes are important for embryonic [development](#) and cell fate specification in many organisms (2-4). In addition, *Wnt-1*, the first Wnt gene identified, is a common site of insertion for the mouse mammary tumor virus (MMTV) that results in ectopic Wnt expression in the mammary gland (1). Ectopic Wnt also transforms mammary epithelial cells in a paracrine fashion (5). Because the Wnt genes encode ligand signals, the identification of [signal transduction](#) components required for their function has been the focus of much research.

### 0.1. *Drosophila* Wingless as a Genetic Model for Wnt Genes

The *wingless* (*wg*) gene in *Drosophila* is a member of the Wnt gene family, as demonstrated by its ability to mimic *Wnt-1* in transforming mammary epithelial cells (6). *Drosophila* is an excellent model for studying developmental signals because it is possible to combine genetics, developmental biology, and biochemistry. Much of our understanding of Wnt gene function during development has in fact come from studies in *Drosophila*. *wg* acts as a signal to perform a variety of regulatory roles during the development of *Drosophila*. Specifically, *wg* is required for maintaining segment polarity during embryogenesis and in [imaginal disc](#) morphogenesis and visceral mesoderm patterning (reviewed in Ref. 4). However, most of the information about the mechanism of *wg* signaling comes from studies on segment polarity, which will be the focus here.

*wg* is expressed in stripes that are a single cell wide in every segment of the developing *Drosophila* embryo. If *wg* function is absent during embryogenesis, the polarity of segments in the larva is lost and a "lawn" of denticles is all that remains in the ventral region (7, 8). During embryogenesis, *wg* is also required to maintain [transcription](#) of several other segment polarity genes (such as *engrailed*) in adjacent cells (9). These studies have demonstrated that *wg* acts as a signal to instruct adjacent cell

fates. These genetic studies have also been complemented with cell biological analysis that has shown that secreted *wg* protein (WG) travels several cell diameters away from the cells expressing it (10, 11). In addition, the *wg* mutation acts non cell autonomously in genetic mosaic experiments (12). Because Wg functions as a signal, much of the focus on the study of Wg signaling has been the elucidation of the Wg signaling pathway. Two main approaches have led to rapid progress in our understanding. Most of the information has come from (1) genetic analysis, although (2) biochemical approaches in *Drosophila* and in other systems have also led to the identification and analysis of novel Wg signaling components.

## 0.2. Genetic Analysis of the *wg* Signaling Pathway

Because the *wg* phenotype is stereotypical, it has been possible to identify *wg* signaling pathway genes by analyzing mutants whose phenotypes resemble the “lawn” pattern of denticles in *wg* mutant embryos. Historically, this approach was successful in identifying several new components of *wg* signaling in *Drosophila*. Three genes have mutant phenotypes identical to that of *wg*: *dishevelled* (*DSH*), *porcupine* (*PORC*), and *armadillo* (*ARM*) (13, 14). Just as in the case of *wg* mutants, these genes are also required to maintain *en* expression in the embryo and for other *wg*-dependent processes during development (4). The differences in these mutants became apparent in the study of their behavior in genetic mosaics. The *porc* mutation (like *wg*) behaves noncell-autonomously, whereas *dsh* and *arm* behave cell-autonomously (15-17). Simple interpretation of these tests suggested that, as in the case of *wg*, *porc* is required outside the cell receiving the *wg* signal, whereas *dsh* and *arm* are required inside the cell. The molecular characterizations of these genes are consistent with this interpretation. The *porc* gene encodes a protein (Porc) that has eight transmembrane domains and is localized perinuclearly, consistent with its association with the Golgi apparatus (15). The *dsh* gene encodes a pioneer cytoplasmic protein that has been well conserved in evolution (16). *dsh* protein (DSH) also has some [homology](#) to other proteins within three distinct domains. Dsh has a [pleckstrin homology \(PH\) domain](#) that is homologous to Axin (discussed later) and a PDZ (*PSD-95/Dlg/ZO-1*) domain found in a variety of proteins that are targeted to intracellular signaling complexes (17). Finally, the *arm* gene encodes the *Drosophila* homologue of b-catenin, a protein that functions in nuclear signaling pathways (18). Other distinctions in the behavior of the *porc*, *dsh* and *arm* mutants have also contributed to the elucidation of *wg* signaling. Cell biological analysis has shown that *porc* mutants may be defective in the secretion of Wg. In the absence of *porc*, Wg is intracellular and is not detected outside the cells that transcribe *wg* (19). Even though *arm* is expressed ubiquitously in the embryo, *arm* protein (ARM) also accumulates at higher levels in cells that receive Wg during development (20). This accumulation is also relevant at the subcellular level because it is mostly cytoplasmic. This cytoplasmic accumulation is abolished in each of *wg*, *dsh*, and *arm* mutant embryos (20). Therefore, based on these observations, *porc* functions upstream of *wg*, and *arm* functions downstream of all three genes. Because *dsh* acts cell-autonomously and encodes a cytoplasmic protein, it is also logical to assume that *dsh* functions downstream of *wg*. These hypotheses have since been proven more formally.

Because *wg* is expressed in a restricted stripe pattern that is a single cell wide, it was possible to express *wg* ubiquitously to generate a phenotype that can be interpreted as the opposite of *wg* mutants. This was initially done by inducing the expression of *wg* from a heat-inducible transgene (HsWg). *en* expression expands in embryos that express HsWg, and the resulting larva are devoid of denticles in the cuticle's ventral region (21). Instead, these larva have “naked” cuticle that is part of the reiterating pattern of segmentation in wild-type larvae. Simplifying these phenotypes, we can say that the “naked” regions in *wg* mutants are deleted and leave only denticles, and in HsWg embryos, the denticles are deleted and leave only the naked regions. Furthermore, Arm accumulates in the cytoplasm of all cells in the embryo, reflecting the ectopic Wg induction in these embryos (21). Because HsWg creates the phenotype opposite to that of *wg*, it was possible to perform genetic [epistasis](#) tests to determine the *wg* genetic pathway. In the absence of *porc*, HsWg still specifies “naked” cuticle and maintains *en* expression (22). This suggests that *porc* functions upstream of *wg*, as was predicted previously. HsWg is, however, completely ineffective in the absence of either *dsh* or *arm*, suggesting that these genes function downstream of *wg* (22). Another opportunity for epistasis with *wg* pathway genes came with the discovery of the *zeste-white 3* mutant (*ZW3*).

Because *zw3* has a phenotype similar to that of HsWg, it was postulated that the *wg* signaling pathway in *zw3* mutants is constitutively active (23). Therefore, it was possible to make double mutants with *zw3* and each of the *wg* pathway genes. Double-mutant epistatic analysis placed *zw3* function downstream of *wg*, *porc*, and *dsh* but upstream of *arm* (23). The **cloned** *zw3* gene encodes the functional homologue of the **serine/threonine kinase**, glycogen synthase kinase 3b (GSK3B) (24). In *zw3* mutant embryos, Arm accumulates in the cytoplasm of all cells in the embryo, mimicking the effects of HsWg (24, 25). Therefore, it is quite likely that Arm functions downstream of *zw3* protein (ZW3). Other studies have also shown that Arm phosphorylation in the embryo depends on Zw3, suggesting a direct interaction between Zw3 and Arm (26). Therefore, epistatic experiments have ordered the *wg* pathway genes in a purely genetic pathway. Because all identified *wg* signaling components are conserved in other organisms, it was possible to use other systems to confirm and to extend the results of genetic analysis in *Drosophila* by using other approaches (17).

### 0.3. Biochemical Analysis of the Wg/Wnt Signaling Pathway

Biochemical studies in *Drosophila* and in other systems, have also identified *wg*/Wnt signaling components. The transfection of *Drosophila frizzled 2* (*DFZ2*) into cells that are unresponsive to Wg enables these cells to stabilize cytoplasmic Arm in response to Wg (27). Because *Dfz2* encodes a seven-transmembrane receptor-like molecule, it has been postulated that *Dfz2* encodes a receptor for *wg* (27). Although Wg binds *Dfz2*-expressing cells, this binding is general among all *Dfz2* gene family members (17). The identification of a *Dfz2* mutant in *Drosophila* should make it possible to test whether *Dfz2* is a crucial Wg receptor.

Another novel component identified is the product of the *fused* locus in mouse, called Axin. The Axin protein is cytoplasmic and acts as a negative regulator of Wnt signaling (28). Recent studies have provided provocative evidence suggesting that the role of Axin may be as an adaptor protein that allows binding GSK3 to the mammalian Arm homologue, b-catenin (29). This allows the simplest interpretation for the repression of Arm by Zw3, whereby Zw3/GSK3 phosphorylates Arm/b-catenin at the four conserved **phosphorylation** sites of this kinase. The *Drosophila Axin* homologue has not been identified, and genetic tests for this interaction are lacking. Recent studies have shown that the GSK3 phosphorylation targets b-catenin for conjugation with **ubiquitin** and its rapid **protein degradation** via the 26S **proteasome** (30). Genetic evidence for these observations has come from the recent identification of mutations in an enzyme involved in ubiquitin regulation. The E3 class of ubiquitin-conjugating enzymes is responsible for targeting specific proteins for ubiquitination. This gene, named *slimb*, encodes an E3 enzyme, and its mutation results in the accumulation of Arm (31).

In addition to cytoplasmic stabilization, Arm and b-catenin also accumulate in the nucleus and associate physically with **HMG** box **transcription factors** of the Tcf/LEF-1 family (32-34). Although the actual mechanism of Arm/b-catenin nuclear translocation is not known, it is very clear that this translocation is important for Wg/Wnt signaling (32, 35). The identification of a *Drosophila* mutant for Tcf/LEF-1, called *pangolin* (*pan*), has established the validity of these hypotheses. The *pan* gene has a mutant phenotype identical to that of *wg* and also functions downstream of *arm* (36, 37). Physical interaction between Arm and *pan* protein (Pan) also has been demonstrated (36).

It has been known for some time that heparin treatment of Wnt-expressing cells leads to release of Wnt from the extracellular matrix (5). It has also been shown that heparin binds Wg (38). Recent genetic studies have confirmed the biochemical evidence for interaction between heparin-like glycosaminoglycans (GAGs) and Wg /Wnts. *Drosophila* mutants in UDP-glucose dehydrogenase (UDP-GLCDH), which is critical for biosynthesizing GAGs, result in a phenotype that is *wg*-like and block, yet do not completely abolish, the effects of ectopic Wg (39-41). Therefore, GAGs are important for Wg function. The long list of potential GAGs was limited to one because heparanase injections mimic the *wg*-like phenotype, and the injection of only heparan sulfate alleviates the UDP-GlcDH mutant phenotype (39). More biochemical experiments will be required to determine the exact role of heparan sulfate and proteoglycans in *wg* signaling.

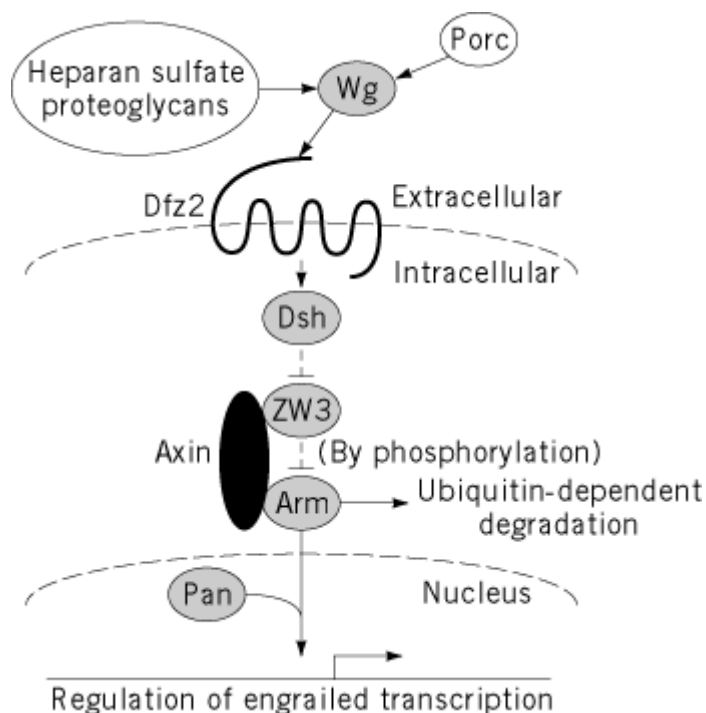


Other potential components of *wg* signaling have been suggested through biochemical analysis. Many of these studies, however, lack other corroborating evidence and await further investigation. Of particular interest is the hypothesis that the *Notch* (*N*) gene may encode a receptor for Wg. Although it was shown that *N* and *wg* share several phenotypic and genetic features during *Drosophila* development (42), the inability of *N* mutants to block the effects of HsWg were used to discount the *N*-*wg* signaling hypothesis (43). However, the lack of biochemical experiments addressing this issue make it difficult to accept or dismiss the idea that *N* protein acts as a receptor for Wg.

#### 0.4. The Mechanism of Wg Signaling

The combined genetic and biochemical approaches discussed before have led to the following scenario for Wg signaling (Fig. 1). Porc is involved in secreting Wg that in turn binds to Dfz2. Heparan sulfate proteoglycans might be involved in all steps leading to this interaction (Dfz2-Wg). Once Wg binds Dfz2, Dsh functions to repress Zw3 function via an unknown mechanism. If Zw3 is not repressed by Dsh, Zw3 enters into a complex with Axin and Arm, leading to the phosphorylation and degradation of Arm (via ubiquitination). When Dsh does repress Zw3, however, Arm is not phosphorylated, cannot be degraded, and is free to accumulate in the cytoplasm. Once Arm accumulates in the cytoplasm, it can enter into a nuclear complex with Pan, resulting in the regulation of genes, such as *en*. The exact mechanism of Arm nuclear translocation is not known.

**Figure 1.** Diagram of the current view of the mechanism of *wg* signaling. Signaling components confirmed via both genetic and biochemical approaches are depicted in gray, and components missing genetic confirmation are depicted in black (eg, Dfz2 and Axin). See text for details.



#### Bibliography

1. R. Nusse and H. E. Varmus (1992) *Cell* **69**, 1073–1087.
2. A. P. McMahon (1992) *Trends Genet.* **8**, 236–242.
3. R. T. Moon, J. D. Brown, and M. Torres (1997) *Trends Genet.* **13**, 157–162.
4. J. Klingensmith and R. Nusse (1994) *Dev. Biol.* **166**, 396–414.

5. S. F. Jue et al. (1992) *Mol. Cell. Biol.* **12**, 321–328.
6. N. R. Ramakrishna and A. M. C. Brown (1993) *Development Suppl.*, 95–103.
7. C. Nüsslein-Volhard and E. Wieschaus (1980) *Nature* **287**, 795–801.
8. A. Martinez-Arias, N. Baker, and P. W. Ingham (1988) *Development* **103**, 157–170.
9. S. DiNardo et al. (1988) *Nature* **332**, 604–609.
10. M. van den Heuvel, R. Nusse, P. Johnston, and P. Lawrence (1989) *Cell* **59**, 739–749.
11. F. Gonzalez et al. (1991) *Mech. Dev.* **35**, 43–54.
12. N. E. Baker (1988) *Dev. Biol.* **125**, 96–108.
13. E. Wieschaus and R. Riggelman (1987) *Cell* **49**, 177–184.
14. N. Perrimon, L. Engstrom, and A. P. Mahowald (1989) *Genetics* **121**, 333–352.
15. T. Kadowaki et al. (1996) *Genes Dev.* **10**, 3116–3128.
16. J. Klingensmith, R. Nusse, and N. Perrimon (1994) *Genes Dev.* **8**, 118–130.
17. K. M. Cadigan and R. Nusse (1997) *Genes Dev.* **11**, 3286–3305.
18. P. D. McCrea, C. W. Turck, and B. Gumbiner (1991) *Science* **254**, 1359–1361.
19. M. van den Heuvel et al. (1993) *EMBO J.* **12**, 5293–5302.
20. R. Riggelman, P. Schedl, and E. Wieschaus (1990) *Cell* **63**, 549–560.
21. J. Noordermeer et al. (1992) *Development* **116**, 711–719.
22. J. Noordermeer, J. Klingensmith, N. Perrimon, and R. Nusse (1994) *Nature* **367**, 80–83.
23. E. Siegfried, E. L. Wilder, and N. Perrimon (1994) *Nature* **367**, 76–80.
24. E. Siegfried, T. B. Chou, and N. Perrimon (1992) *Cell* **71**, 1167–1179.
25. M. Peifer, D. Sweeton, M. Casey, and E. Wieschaus (1994) *Development* **120**, 369–380.
26. M. Peifer, L. M. Pai, and M. Casey (1994) *Dev. Biol.* **166**, 543–566.
27. P. Bhanot et al. (1996) *Nature* **382**, 225–230.
28. L. Zeng et al. (1997) *Cell* **90**, 181–192.
29. S. Ikeda et al. (1998) *EMBO J.* **17**, 1371–1384.
30. H. Aberle et al. (1997) *EMBO J.* **16**, 3797–3804.
31. J. Jiang and G. Struhl (1998) *Nature* **391**, 493–496.
32. M. Molenaar et al. (1996) *Cell* **86**, 391–399.
33. V. Korinek et al. (1997) *Science* **275**, 1784–1787.
34. P. J. Morin et al. (1997) *Science* **275**, 1787–1790.
35. J. Behrens et al. (1996) *Nature* **382**, 638–642.
36. E. Brunner, O. Peter, L. Schweizer, and K. Basler (1997) *Nature* **385**, 829–833.
37. M. Van de Wetering et al. (1997) *Cell* **88**, 789–799.
38. F. Reichsman, L. Smith, and S. Cumberledge (1996) *J. Cell Biol.* **135**, 819–827.
39. R. C. Binari et al. (1997) *Development* **124**, 2623–2632.
40. T. E. Haerry, T. R. Heslip, J. L. Marsh, and M. B. O'Connor (1997) *Development* **124**, 3055–3064.
41. U. Häker, X. Lin, and N. Perrimon (1997) *Development* **124**, 3565–3573.
42. J. P. Couso and A. Martinez-Arias (1994) *Cell* **79**, 259–272.
43. K. M. Cadigan and R. Nusse (1996) *Development* **122**, 2801–2812.

### **Suggestion for Further Reading**

44. Ken M. Cadigan and Roel Nusse (1997). Wnt signaling: A common theme in animal development, *Genes Dev.* **11**, 3286–3305. An exceptionally thorough review of the *wingless/Wnt* signaling field.

## Wobble Pairing

In [translation](#) of the genetic information present in the [messenger RNA](#) using [transfer RNA](#) during **protein biosynthesis**, Crick suggested, in 1966, that base pairing of the first (5') base of the [anticodon](#) of transfer RNA with that at the third (3') position of the **codon** (see [Genetic Code](#)) does not necessarily conform to normal Watson–Crick [base pairing](#) (1). He introduced the term *wobble pairing* for this noncanonical base pairing. Crick pointed out, as a typical example of wobble pairing, that the inosine (I) (a derivative of G) located at the first position of the yeast tRNA<sup>Ala</sup> anticodon can base-pair at the third position of the codon with not only C but also with U and A. Since then, as nucleotide sequences of numerous tRNAs from various sources have been determined (2), knowledge of the extent of wobble pairing has expanded, as shown in Table 1. (See top of next page) The wobble pairings have been verified experimentally both *in vitro* (3) and *in vivo* (4).

**Table 1. Anticodon-Codon Paring Rules So Far Elucidated<sup>a</sup>**

| 5'  | 3'                     | Occurrence                |
|---|------------------------|---------------------------|
| Anticodon <sup>3'</sup>   | Codon <sup>5'</sup>    |                           |
| 1. U or modified U at the wobble position of the anticodon  |                        |                           |
| a. Unmodified U   |                        |                           |
| UNN   | <u>UNN</u>             | <u>UNN</u> <u>UNN</u>     |
| Family boxes in mitochondria, <i>Mycoplasma</i> spp., and chloroplasts  |                        |                           |
| AN'N'   | <u>UN'N'</u>           | <u>CN'N'</u> <u>GN'N'</u> |
| b. U modified to U* (mo <sup>5</sup> U or cmo <sup>5</sup> U)   |                        |                           |
| U*NN  | <u>U*NN</u>            | <u>U*NN</u>               |
| Family boxes in eubacteria  |                        |                           |
| AN'N'   | <u>UN'N'</u>           | <u>GN'N'</u>              |
| c. U modified to U <sup>±</sup> (Um, cmnm <sup>5</sup> U, or mcm <sup>5</sup> U)  |                        |                           |
| U <sup>±</sup> NN   | <u>U<sup>±</sup>NN</u> |                           |
| 2-codon sets in mitochondria, bacteria and eukaryotes   |                        |                           |
| AN'N'   | <u>GN'N'</u>           |                           |
| d. U modified to U <sup>#</sup> (mnm <sup>5</sup> s <sup>2</sup> U, mcnm <sup>5</sup> s <sup>2</sup> U or cmnm <sup>5</sup> s <sup>2</sup> U) |                        |                           |
| U <sup>#</sup> NN   | <u>U<sup>#</sup>NN</u> |                           |
| 2-codon sets in eubacteria and eukaryotes   |                        |                           |
| AN'N'   | <u>GN'N'</u>           |                           |
| 2. G or modified G at the wobble position of the anticodon  |                        |                           |
| a. Case 1 for unmodified G  |                        |                           |

|  |               |               |   |  |
|--|---------------|---------------|---|--|
| GNN  | GNN           |               | 2-codon sets of all organisms, and family boxes in bacteria   |  |
| CN'N'  | <u>U</u> N'N' |               |   |  |
| b. Case 2 for unmodified G                                 |               |               |   |  |
| GNN  | <u>G</u> NN   | <u>G</u> NN   | IleAUN (except AUG) and AsnAAN (except AAG) in starfish mitochondria. SerAGN (except AGG) in <i>Drosophila</i> mitochondria |  |
| CN'N'  | <u>U</u> N'N' | <u>A</u> N'N' |   |  |
| c. G modified to Q   |               |               |   |  |
| QNN  | <u>Q</u> NN   |               | 2-codon sets in eubacteria and eukaryotes   |  |
| CN'N'  | <u>U</u> N'N' |               |   |  |
| d. G modified to I   |               |               |   |  |
| INN  | <u>I</u> NN   | <u>I</u> NN   | ArgCGN in eubacteria, and all family boxes in eukaryotes except GlyGGN  |  |
| CN'N'  | <u>U</u> N'N' | <u>A</u> N'N' |   |  |
| e. G modified to G* (m <sup>7</sup> G)                     |               |               |   |  |
| G*NN   | <u>G</u> *NN  | <u>G</u> *NN  | SerAGN in starfish and squid mitochondria   |  |
| CN'N'  | <u>U</u> N'N' | <u>A</u> N'N' | <u>G</u> N'N'   |  |
| 3. A at the wobble position of the anticodon (rare)        |               |               |   |  |
| ANN  | <u>A</u> NN   | <u>A</u> NN   | <u>A</u> NN   | ThrACU and ArgCGN in mycoplasma spp. and yeast mitochondria. ArgCGN in nematode mitochondria |
| UN'N'  | <u>C</u> N'N' | <u>G</u> N'N' | <u>A</u> N'N'   |  |
| 4. C or modified C at the wobble position of the anticodon |               |               |   |  |
| a. Unmodified C  |               |               |   |  |
| CNN  |               |               | Only 5'N'N'G <sup>3'</sup> codons in all organisms  |  |
| GN'N'  |               |               |   |  |
| b. C modified to C*(f <sup>5</sup> C)                      |               |               |   |  |
| C*AU   | <u>C</u> *AU  |               | MetAUR in animal mitochondria   |  |
| GUA  | <u>A</u> UA   |               |   |  |
| c. C modified to L   |               |               |   |  |
| <u>L</u> AU  |               |               | IleAUA in eubacteria and plant mitochondria   |  |
| <u>A</u> UA  |               |               |   |  |

<sup>a</sup> The anticodon-codon pairing is shown by 5'ANTICODON (NNN)<sup>3'</sup>, 3'CODON(N'N'N')<sup>5'</sup> where N in the anticodon and its corresponding N' in the codon (beneath the N) form Watson-Crick base pairing. The underlined letters are involved in wobble pairing. Family box means that 4 synonymous codons are included in a box in the genetic code table (for example, GUU, GUC, GUA and GUG for Val codons). 2-codon set means that only 2 codons are used for one amino acid (for example, AAA and AAG for Lys codons).

Abbreviations used are: mo<sup>5</sup>U, 5-methoxy U; cmo<sup>5</sup>U, 5-carboxymethoxy U; cmnm<sup>5</sup>U, 5-carboxymethyl-aminomethyl U; mcm<sup>5</sup>U, 5-methylcarboxymethyl U; Um, 2'-O-methyl U; mnm<sup>5</sup>s<sup>2</sup>U, 5-methylaminomethyl-2-thio U, cmnm<sup>5</sup>s<sup>2</sup>U, 5-carboxymethyl-aminomethyl-2-thio U; mcm<sup>5</sup>s<sup>2</sup>U, 5-methylcarboxymethyl-2-thio U; Q, queosine (in eubacteria) or its derivatives, manQ or galQ (in eukaryotes); f<sup>5</sup>C, 5-formyl C; L, lysidine. From Osawa (5) with modifications (6, 7). See also Bjork (8).

The importance of wobble is that there is no need for a unique transfer RNA specific for each codon of the genetic code. Consequently, one transfer RNA can respond to several codons, inserting the same amino acid for each.

### Bibliography

1. F. H. C. Crick (1966) *J. Mol. Biol.* **19**, 548–555.
2. M. Sprinzl, T. Hartmann, J. Weber, J. Blank, and R. Zeidler (1989) *Nucleic Acids Res.* **17** (Suppl.), 1–173.
3. D. Söll, J. Cherayil, D. S. Jones, R. D. Faulkner, A. Hampel, R. M. Bock, and H. G. Khorana (1966) *Cold Spring Harbor Symp. Quant. Biol.* **31**, 51–61.
4. C. Guthrie and J. Abelson (1982) In *The Molecular Biology of the Yeast Saccharomyces: Metabolism and Gene Expression* (J. N. Strathern, E. Jones, and J. Broach, eds.) Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, pp 487–528.
5. S. Osawa (1995) *Evolution of the Genetic Code*, Oxford University Press, Oxford, U.K., pp 1–205.
6. Y. Watanabe, H. Tsurui, T. Ueda, R. Furushima-Shimogawara, S. Takamiya, K. Kita, K. Nishikawa, and K. Watanabe (1997) *Biochim. Biophys. Acta* **1350**, 119–122.
7. S. Matsuyama, T. Ueda, P. F. Clain, J. A. McCloskey, and K. Watanabe (1998) *J. Biol. Chem.* **273**, 3363–3368.
8. G. Björk (1995) In *tRNA, Structure, Biosynthesis, and Function* (D. Söll and U. L. RajBhandary, eds.), ASM Press, Washington D.C. pp. 165–205.

### Writhe, DNA

Writhe is an aspect of [DNA structure](#) that concerns [DNA topology](#).

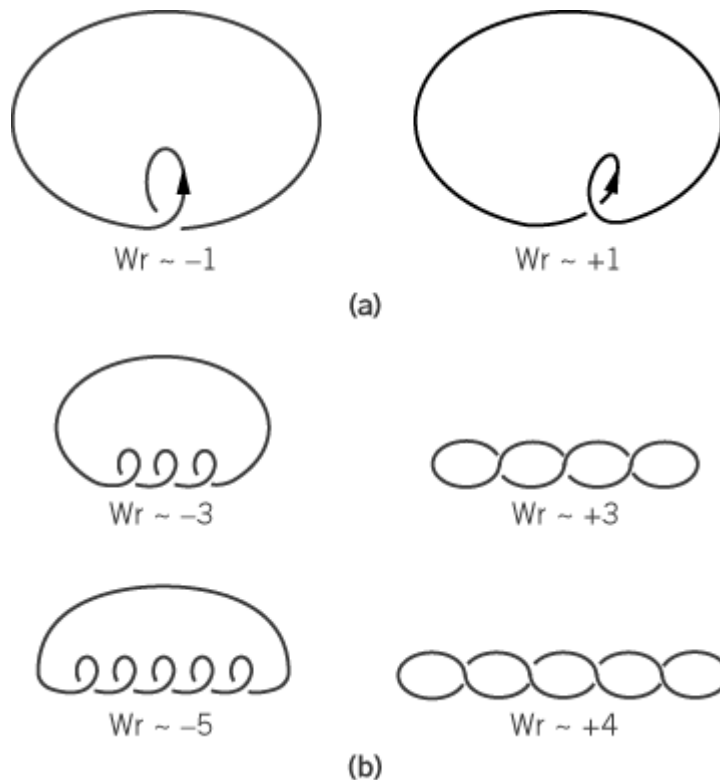
#### 1. Writhing Number

The writhe or writhing number ( $Wr$ ) is a property of a single closed space curve; in application to DNA, the appropriate space curve is taken to be the duplex axis  $A$ . For considerations of DNA structure, the writhe is important insofar as it enters into the fundamental relationship between the **linking number** ( $Lk$ ), the **twist** ( $Tw$ ), and  $Wr$  for a closed circular duplex DNA,  $Lk = Tw + Wr$ . Since  $Wr$  is calculated only for a closed curve, the quantity is not applicable to nicked circular or linear DNAs. The magnitude of the writhe is directly proportional to  $n$ , where  $n$  is the number of superhelical turns (see [Supercoiling of DNA](#)). For a nicked circular duplex DNA, whose axis on the average lies in a plane,  $Wr = 0$ . In all other cases, some of the winding is always converted into changes in  $Tw$ , and  $|Wr| < n$ .

### 1.1. The Index Approach

As with the linking number, conceptually the easiest way to calculate the writhe is with the index approach (see [Linking Number Of DNA](#)). The axis is assigned an orientation and projected onto a plane. Although a specific orientation must be chosen for the axis in order to calculate the directed writhing number, the actual choice of orientation is immaterial, since a reversal of orientation changes the orientations of both curve segments simultaneously. Some examples are shown in Figure 1. For any given projection, an index number of +1 or -1 is assigned to each crossing of one strand by the other. The index number is negative if the tangent vector to the upper portion of the curve must be rotated clockwise to coincide with the tangent vector of the lower part of the curve. The index number is positive in the opposite case. Since the writhe is not a topological quantity, however, its calculation is *not* independent of the projection chosen. This serves to distinguish  $Wr$  clearly from  $Lk$ . The sum of all the index numbers gives the projected writhing number,  $Wr_p$ , for the DNA axis in the projection chosen.

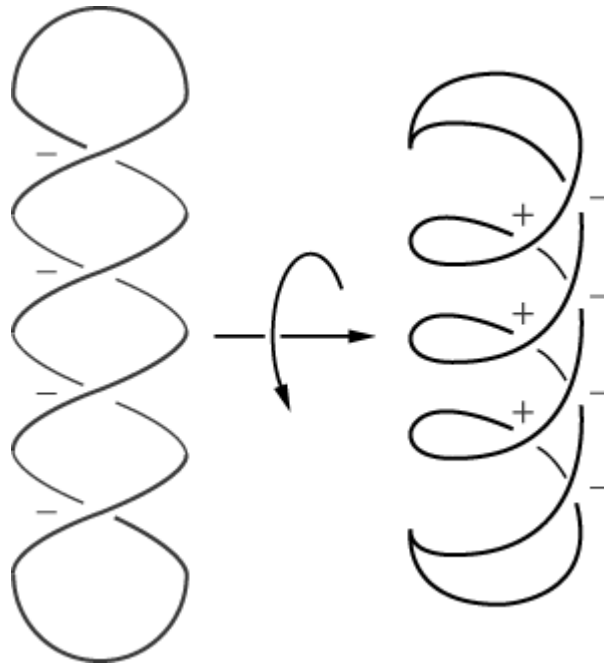
**Figure 1.** Examples of calculation of the writhe for some simple curves using the index approach. (a) The writhing number of curves with one coil. In the projections shown, the projected writhe,  $Wr_p$ , is  $\pm 1$ . However, if the curves are rotated so that other projections are viewed, there are some side views where no crossings are observed. Since  $Wr$  requires that all values of the projected writhe be averaged, the actual values of  $Wr$  will be slightly greater than -1 in one case and slightly less than +1 in the other case. (b) The extension to curves having multiple coils is straightforward. In the views shown here, the projected writhes are the numbers shown under each curve. As with the loops of one coil, however, some projections will have no crossings, giving rise to values of  $Wr$  that are slightly more positive than the negative integers indicated and slightly more negative than the positive integers indicated.



The writhing number,  $Wr$ , is obtained by averaging the projected writhing number over all possible projected views. The values of the directed writhing numbers, and hence of  $Wr$ , are changed by deformation of the axis. As an example, DNA supercoiled into a regular plectonemic superhelix is shown schematically in Figure 2. In most views, the crossings are all negative and equal in number to the number of supercoils. These have a projected writhe of  $-n$ . In some projections, however, positive crossings are also observed. These have a projected writhe more positive than  $-n$ . The

writhe is the sum of all projected writhes, so it is clear that  $|Wr| < n$ . Finally,  $Wr$  changes by 2 if the axis is passed through itself once, as happens by treatment with a type II DNA **topoisomerase**. This is because such a passage converts a +1 crossing to a -1 crossing, a change of -2, or a -1 crossing to a +1 crossing, a change of +2.

**Figure 2.** Analysis of the projected writhe for an interwound DNA superhelix. The solid curve represents the axis of a cdDNA wound into an unbranched, interwound superhelix. Two views of the same structure are depicted. The view on the left shows only the nodes with negative index numbers, and  $Wr_p = -4$ . These result from the crossing of distant DNA segments. The view on the right results from rotation about the horizontal axis so that the top of the molecule is tipped toward the observer. This view contains the same negative contributions but contains additional nodes that have positive index contributions, resulting from the crossing of nearby DNA segments. The value of  $Wr_p$  for this view is thus -1. The  $Wr$  is the average of the  $Wr_p$  values of all views; hence  $|Wr|$  is less than the number of supercoils.



### 1.2. Application of the Gauss Integral

If the trajectory of the DNA axis is known, the writhe may alternatively be calculated from the Gauss integral in a manner analogous to the calculation of  $Lk$  (1) (see [Linking Number Of DNA](#)):

$$Wr = \frac{1}{4\pi} \int_{C \times C} \int \frac{\mathbf{e} \cdot \mathbf{T}_2 \times \mathbf{T}_1}{r^2} ds_1 ds_2 \quad (1)$$

where  $\mathbf{T}_1$  and  $\mathbf{T}_2$  are the unit tangent vectors to the axis curve at any two locations,  $x_1$  and  $x_2$ ;  $r$  is the distance between the points  $r = \|\mathbf{x}_2 - \mathbf{x}_1\|$ ; and  $\hat{\mathbf{e}}$ ; is the unit vector  $(\mathbf{x}_2 - \mathbf{x}_1)/r$ .

### 1.3. Use of the Surface Linking Number

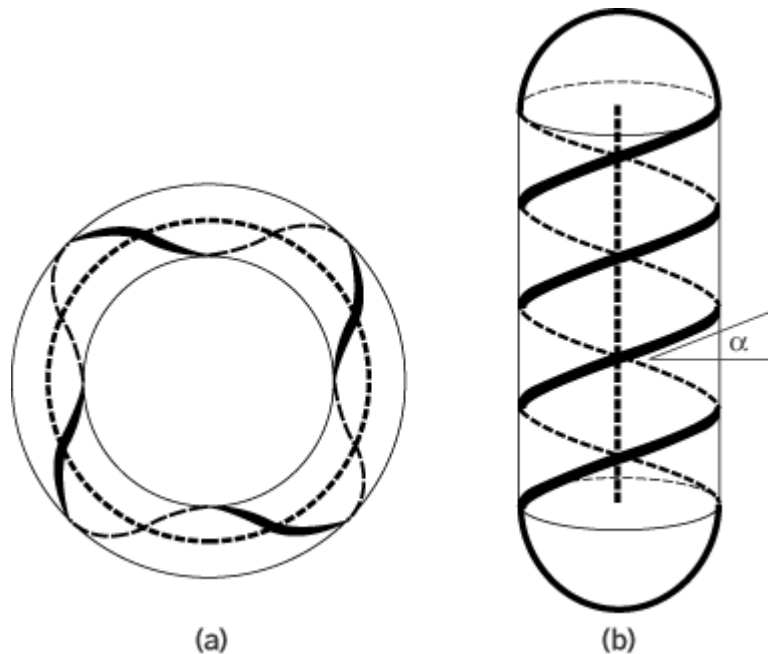
The third and easiest way to calculate the writhe takes advantage of the fundamental relationship for a closed circular DNA (2), in the form  $Wr = Lk - Tw$ . Although this method is indirect, it is considerably easier than the direct methods because both  $Lk$  and  $Tw$  are much more readily obtained than  $Wr$  itself. Even more straightforward is the combination of the fundamental relationship with considerations of surface linking (see [Surface Wrapping Of DNA](#)). The fundamental relationship then becomes

$$Wr = SLk - STw \quad (2)$$

The surface linking number,  $SLk$ , is simply related to the geometry of the real or virtual surface about which the DNA winds. For the two most commonly encountered cases,  $SLk = 0$  for superhelical DNA having the form of any plectonemic superhelix, no matter how distorted;  $SLk = \pm n$  for any toroidal superhelix, where the sign is positive for right-handed winding and negative for left-handed winding. Equation (2) then becomes  $Wr = -STw$  for a plectonemic superhelix and  $Wr = \pm n - STw$  for a toroidal superhelix.

The calculation of the writhe therefore reduces to a problem in the calculation of the surface twist. The calculation of  $STw$  requires knowledge of the shape of the surface and of the associated path of the DNA axis. If the surface is a protein or other real object, the required information can be obtained, at least in principle, by X-ray diffraction. In the case of a free ccdDNA, the axis can sometimes be considered to wrap on a virtual surface. The two most commonly encountered surfaces of this type are the capped cylinder (for plectonemic supercoiling) and the circular torus (for toroidal supercoiling). These are illustrated in Figure 3. For the plectonemic superhelix,  $STw = n \sin(\alpha)$  for DNA whose axis winds as a regular right-handed superhelix of  $n$  turns and superhelix pitch angle  $\alpha$ . It follows immediately that  $Wr = -n \sin(\alpha)$  for a right-handed plectonemic superhelix ( $DLk < 0$ ), and  $Wr = +n \sin(\alpha)$  for a left-handed plectonemic superhelix ( $DLk > 0$ ).

**Figure 3.** Two types of geometry used in the calculation of the writhe by means of the surface linking number (see [Surface Wrapping Of DNA](#)) and surface twist (see [Twist, DNA](#)). (a) Toroidal superhelical DNA wrapped on the surface of a virtual torus. The DNA axis wraps +4 times in a right-handed sense on the surface of a torus.  $SLk = +4$ . (b) Plectonemic interwound DNA, wrapped on the surface of a virtual capped cylinder. The DNA wraps +2 times up and +2 times down a virtual cylinder, with circular closure at the caps.  $SLk = 0$ . The superhelix pitch angle is  $\alpha$ .



If the DNA is wound in a regular right-handed toroidal superhelix,  $SLk = +n$ . The calculation of  $STw$  in this case is more involved than for the plectonemic superhelix and requires knowledge of the dimensions of the virtual toroidal surface on which the DNA is wrapped. If the DNA axis wraps  $n$  times around a circular torus of major radius  $R$  (distance from the center of symmetry to the central axis of the torus) and minor radius  $r$  (radius of the torus), then  $STw$  is given by an elliptical integral that must be solved numerically (3). If the torus is highly extended relative to its width, so that  $r/R \ll 1$ , the result is approximately  $STw \approx nR/\sqrt{R^2 + n^2r^2}$ . In the limiting case



where  $r/R \rightarrow 0$ ,  $STw$  is the same as for the plectonemic superhelix. In the other limiting case that  $r/R \rightarrow 1$ , numerical solutions to the elliptical integral must be obtained.

## 2. Absolute Writhe

The absolute writhe,  $AWr$ , is formally similar to the writhe except that the absolute value of the writhe integrand, instead of the usual signed one, is used in the Gauss integral (3). Although it not part of the basic relationship for a topological domain,  $Lk = Tw + Wr$ ,  $AWr$  has the advantage that it is zero if and only if the curve is planar. Indeed, the closer the absolute writhe is to zero the more nearly planar is the curve. This contrasts with  $Wr$ , which may be zero even though the curve lies out of plane. For example, any curve that lies on the surface of a sphere has  $Wr = 0$ . Such a curve certainly does not have to be planar.

The difference between  $Wr$  and  $AWr$  may also be understood using the process of projections. If a closed space curve is viewed from a distant point, apparent crossings may be seen by the observer. To each of these crossings is assigned a value of +1 or -1, depending on the orientation of the tangents to the crossing segments. The projected writhe is the algebraic sum of these crossing numbers, and the writhe is the average of all projected writhing numbers over all possible projections. The absolute writhe is defined similarly, but with one major difference: instead of attaching a sign to any crossing, the absolute value of the number of crossings is counted. This number is called the *projected absolute writhe*. The absolute writhe is the average of all absolute projected writhing numbers over all possible projections. Thus, for a very tightly wound, almost planar figure eight, the writhe may be +1 or -1 depending on the overlay of the crossing segments; but the absolute writhe is always +1. In terms of the Gauss integral,  $AWr$  is calculated similarly to  $Wr$ :

$$AWr = \frac{1}{4\pi} \int_{C \times C} \int \left| \frac{(\mathbf{y} - \mathbf{x}) \cdot \mathbf{T}_2 \times \mathbf{T}_1}{r^3} \right| ds_1 ds_2 \quad (3)$$

## Bibliography

1. J. H. White (1989) in *Mathematical Methods for DNA Sequences*, M. S. Waterman, ed., CRC Press, Boca Raton, FL, pp. 225–253.
2. J. H. White (1969) *Am. J. Math.* **91**, 693–728.
3. J. H. White and W. R. Bauer (1986) *J. Mol. Biol.* **189**, 329–341.

## Suggestions for Further Reading

4. W. R. Bauer, F. H. C. Crick, and J. H. White (1980) Supercoiled DNA, *Sci. Ame.* **243**, 118–122.
5. F. B. Fuller (1971) The writhing number of a space curve, *Proc. Natl. Acad. Sci USA* **68**, 815–819.

## X-chromosome

In **eukaryotes** that utilize sexual **reproduction**, usually two distinct [chromosomes](#) exist, called [sex chromosomes](#). In the majority of **plants** and animals, the male has one copy of each of these two chromosomes and is called the **heterogametic sex**. The two sex chromosomes (X and Y in humans) have partial homology. They pair during **prophase** of **meiosis** in the male, so that two kinds of **gametes** are produced, each containing just one of the two types of sex chromosomes. In the female, the **homogametic sex**, are two identical sex chromosomes (X in humans), so only one type of gamete results. The presence of different numbers of X-chromosomes in different cells, depending on the sex, relative to the constant number of [autosomes](#) requires that molecular mechanisms exist to compensate for the presence of two X-chromosomes in each cell of a female mammal. The process of **dosage compensation** leads to the inactivation of one of the X-chromosomes in humans (see [X-Chromosome Inactivation](#)). This process is generally random, in that either maternally or paternally derived X-chromosomes are inactivated (see [Random X-Inactivation](#)), and results in the appearance of a [Barr body](#), consisting largely of [facultative heterochromatin](#).

In *Drosophila* the male has a single copy of the X-chromosome, whereas the female is XX with two copies of the X-chromosome. However, the single X-chromosome in males produces the same amount of gene product as the combination of two X-chromosomes in females (1). The male X-chromosome is twice as transcriptionally active as the average female X-chromosome. This is accomplished by the targeting of an [enzyme](#) complex that hyperacetylates the [histone](#) proteins in the [nucleosomes](#) and [chromatin](#) of the male X-chromosome. Histone hyperacetylation promotes the processivity of **RNA polymerase** through nucleosomal arrays and along the chromatin fiber. Thus one gene in the single male X-chromosome in *Drosophila* is as active as the two copies of the gene in the female. An important aspect of the *Drosophila* dosage compensation pathway is the capacity to count the number of sex chromosomes relative to sets of autosomes. If a *Drosophila* egg has equal numbers of X-chromosomes and autosomes, it develops into a female. If there is an extra copy of autosomes relative to the number of X-chromosomes, it does not follow the female differentiation pathway (2).

In the [nematode](#) *Caenorhabditis elegans*, dosage compensation occurs in XX individuals that are [hermaphrodites](#) to normalize gene activity on the X-chromosome with males that are XO. This requires a reduction in [transcription](#) from the two X-chromosomes in hermaphrodites. A chromosomal modification is involved again, but unlike the case in *Drosophila* where gene activation was involved, in *C. elegans* hermaphrodites a specialized (stability and maintenance of chromosome SMC) chromosomal condensation protein is targeted to the X-chromosome to repress its gene activity twofold (3).

#### Bibliography

1. H. J. Muller (1950) Harvey Lect. **43**, 165–175.
2. C. B. Bridges (1921) Science **54**, 252–268.
3. P. T. Chuang, D. G. Albertson, and B. J. Meyer (1994) Cell **79**, 459–474.

#### Suggestion for Further Reading

4. R. P. Wagner, M. P. Maguire, and R. L. Stallings (1993) *Chromosomes. A Synthesis*, Wiley-Liss, New York.

#### X-Chromosome Inactivation

Inactivation of an [X-chromosome](#) in female cells is part of the **dosage compensation** mechanism necessary to allow equivalent expression of X-linked genes in female and male cells, which have two and one X-chromosomes, respectively. The suggestion that X-inactivation might occur was first presented as an hypothesis by Mary Lyon (see [Lyon Hypothesis](#)). The result of the inactivation process is the appearance of the [Barr body](#). In eutherian (placental) mammals, the initial choice between inactivation of maternal or paternal X-chromosomes is random (see [Random X-Inactivation](#)), but once established in a repressed state, the same X-chromosome is inactivated after every cell division. This is an excellent example of establishing and maintaining a chromosomally based state of determination.

Abnormalities in the process of X-chromosome inactivation have allowed the genetic definition of the X-inactivation center required for inactivation to occur (1). The X-inactivation center is actually on the inactive X-chromosome and contains the gene encoding Xist RNA. This site actually remains actively transcribed, whereas the rest of the X-chromosome is silenced. The Xist RNA is the key regulator of inactivation (2, 3). Remarkably the gene does not encode [messenger RNA](#), but it produces instead a long, untranslated RNA that remains associated with the inactive X-chromosome (4, 5). Through some unknown mechanism, the Xist RNA has a causal role in directing the **heterochromatinization** of the inactive X-chromosome. The process of heterochromatinization leads to numerous differences between active and inactive X-chromosomes including DNA **methylation**, histone **acetylation** status, and replication timing (see [Facultative Heterochromatin, Euchromatin](#)).

DNA methylation is characteristic of many inactive **promoters** and genes. Consistent with this observation, many genes in the inactive X-chromosome are heavily methylated in contrast to the active X-chromosome (6). However, the kinetics with which particular sites within the X-linked genes become methylated during the [differentiation](#) of embryonic female **somatic cells** do not always correlate with the timing of transcriptional inactivation (7). Although it remains to be determined if a subset of key sites around regulatory elements is always methylated before transcription is repressed, it seems probable that other mechanisms must supplement any influence of DNA methylation on transcription.

Heterochromatin normally replicates late during **S-phase**. Replication timing has been proposed as a determinant of transcriptional activity. Genes that replicate late during S-phase might do so under conditions of limiting [transcription factors](#), or they might be assembled into a repressive [chromatin](#) structure using components translated only late in S-phase. The active X-chromosome normally replicates early in S-phase, whereas the inactive X-chromosome replicates late (8). However, female lymphoma cell lines have been isolated in which the opposite occurs (9). Thus the inactive X-chromosome does not have to replicate late in S-phase to be transcriptionally quiescent.

The formation of local repressive chromatin structures in which key genetic regulatory elements are rendered inaccessible to transcription factors by inclusion within positioned [nucleosomes](#) is an important mechanism for transcriptional repression. Promoters in the inactive X-chromosome are incorporated into positioned nucleosomes, whereas promoters in the active X-chromosome are free of such structures and have transcription factors bound to them (10). Thus specific local chromatin structures clearly have a role in regulating differential gene activity between the two X-chromosomes. Nevertheless, a causal relationship between chromatin structure and transcription has yet to be established. Nucleosomes on the active X-chromosome contain predominantly acetylated histones, whereas those on the inactive X-chromosome are not acetylated (11). Thus establishing and maintaining specific chromatin structures containing modified histones is an excellent candidate mechanism for establishing and maintaining differential expression of genes between the two X-chromosomes. Differential methylation and replication timing may stabilize these different states of gene activity.

Bibliography

1. H. F. Willard et al. (1993) Cold Spring Harbor Symp. Quant. Biol. **58**, 315–325.
2. N. Brockdorff et al. (1991) Nature **351**, 329–331.
3. C.-J. Brown et al. (1991) Nature **349**, 38–42.
4. N. Brockdorff et al. (1992) Cell **71**, 515–526.
5. C.-J. Brown et al. (1992) Cell **71**, 527–538.
6. S. G. Grant and V. M. Chapman (1988) Annu. Rev. Genet. **22**, 199–233.
7. L. F. Lock, N. Takagi, and G. R. Martin (1987) Cell **48**, 36–46.
8. N. Takagi (1974) Exp. Cell Res. **86**, 127–135.
9. I. Yoshida, N. Kashio, and N. Takagi (1993) EMBO J. **12**, 4397–4405.
10. A. P. Riggs and G. P. Pfeifer (1992) Trends Genet. **8**, 169–174.
11. P. Jeppesen and B. M. Turner (1993) Cell **74**, 281–291.

## X-Ray Crystallography

X-ray crystallography is still the experimental technique that determines the most detailed structures of macromolecules. In 1912 Knipping and Friedrich in Laue's laboratory performed the first X-ray diffraction experiment with a crystal (1). Their result showed at the same time that X rays behave as waves and that crystals display a high degree of order in arranging their atoms, ions, or molecules. The crystal acts as a three-dimensional grating and diffracts X rays just as visible light is diffracted by an optical grating. Since 1912, X-ray crystallography has had an enormous impact on chemistry and biology. At first, great triumphs were obtained in the field of inorganic chemistry. The structures of simple ionic compounds and of complex structures, like the silicates, became understood, and ionic radii could be measured with high accuracy. Unlike inorganic compounds that are mainly ionic, organic compounds consist of molecules. The intensity of the beams diffracted by organic compounds is relatively low and, in the beginning, progress was slow in determining their structures. However, the introduction of improved hardware and Fourier analysis solved the problem, and now the X-ray structure determination of a small or medium-sized organic molecule is an easy job, at least when good quality crystals are grown.

Protein crystallography is the youngest branch of X-ray crystallography. It started in 1934 when Bernal took the first X-ray diffraction picture of a protein crystal (2), first a crystal of pepsin, soon followed by a crystal of insulin. These observations showed that even these large molecules are nicely ordered in their crystals and that their structures might be solved. This was not an easy problem, however, and it took twenty years to find a solution (3). Surprisingly, it could be done by a technique already successfully applied for small organic compounds: the method of [isomorphous replacement](#) in which a heavy atom is attached to the protein structure. It was Perutz' achievement to show that the method could also be used for proteins and that the heavy atoms cause a much larger change in intensity than initially expected. Because of the rather primitive instrumentation, progress was slow at first. With the introduction of more sophisticated instrumentation and computers with suitable software, X-ray crystallography of proteins has grown to a relatively common technique, and structures are being published with ever increasing speed. They are collected in the Protein Data Bank, Brookhaven National Laboratory, Upton, NY 11973-5000, USA (see [Databases](#)).

### 1. X-ray sources

In the home laboratory, X rays are produced in vacuum tubes by bombarding a metal anode with

accelerated electrons. They are either sealed tubes that have a static anode or they have a rotating anode and are continuously pumped. An X-ray tube emits a spectrum of wavelengths whose characteristics depend on the metal. For most diffraction experiments, a single wavelength is selected by a monochromator, 0.71 Å from a molybdenum anode or 1.5418 Å from a copper anode. Diffraction is stronger for the longer wavelength, and this is preferred for proteins. The maximum power of an X-ray tube is limited by heating of the anode caused by the accelerated electrons. Rotating anode tubes have better heat dissipation, produce a higher intensity beam, and therefore, are preferable for protein work.

Much stronger X-ray beams are produced in synchrotrons (4). These devices have electrically charged particles (electrons or positrons) circulating in a vacuum. A continuous spectrum of electromagnetic radiation down to a certain minimum wavelength is emitted when the particles change direction. The minimum depends on the power of the synchrotron. The higher the power, the shorter the minimum wavelength. A wavelength near 1 Å is often chosen for protein crystallography with synchrotron radiation. This has the advantage over the copper wavelength of 1.5418 Å that the absorption of the beam is much lower.

## 2. X-ray detectors

The beams diffracted by crystals of small compounds are usually detected and measured with a diffractometer. Using this instrument, the crystal can be rotated around three or four axes and can reach nearly every orientation. The diffracted beams are measured with a photon counter. Because the beams are measured one after the other, this process is slow and not very suitable for protein crystals where thousands of beams must be measured. Protein crystallographers prefer an instrument that has one or more axes to orient the crystal and is equipped with an area detector on which hundreds or thousands of diffracted beams are registered nearly simultaneously. A popular type is the imaging plate covered with phosphorescent material. After registration, the information collected on the imaging plate is released by a laser and read with a photomultiplier.

## 3. Diffraction theory

The most convenient way to understand the diffraction of X-rays by a crystal was introduced by Bragg (5). It is based on the idea of the reflection of X rays from lattice planes in the crystal (see [Bragg Angle](#)). They are the planes constructed through the lattice points, which are the corners of the [unit cells](#). Within a set, the planes are parallel and equidistant, and have a perpendicular distance  $d$ . Such a set is characterized by three indices:  $h$ ,  $k$ , and  $l$ . The beams reflected at an angle  $\theta$  from the lattice planes reinforce each other only if the path difference for beams from successive planes is equal to a whole number of wavelengths. This occurs when Bragg's law is fulfilled:

$$2d \sin \theta = \lambda \quad (1)$$

The reflected beams have the same indices  $h$ ,  $k$ , and  $l$  as the set of planes from which they are reflected. Bragg's law gives the direction of the diffracted beams, but not their intensity. The latter is determined by the distribution of atoms and molecules in the unit cell of the crystal. It can be shown that Bragg's reflection idea is equivalent to positive interference of the scattering by the individual unit cells in the crystal. The scattering by the unit cells is reinforced in the Bragg directions and extinguished in other directions.

In powder diffraction, the specimen is not a single crystal but consists of a large number of crystal grains. It is exposed to a parallel beam of X rays, and there are always some grains for which the Bragg condition is satisfied. The result is a series of circles centered around the direct beam. The powder method is suitable only for very simple structures, where only a limited number of data are sufficient for the structural determination and for detecting phase transitions. It is completely useless for proteins.

When X rays interact with a crystal, they are scattered by the electrons in the crystal. A unit cell contains a huge number of electrons, and every electron scatters the X-ray beam. These beams interfere with each other, and the net amplitude of the scattering by a unit cell corresponds to a number of electrons  $F$ , which is smaller than the total number of electrons in the unit cell and depends on the electron distribution in the unit cell. Because the crystal scatters only if the beams originating from the individual unit cells positively interfere, the amplitude of the scattered beam  $hk\ell$  is proportional to  $F(hk\ell)$ , the scattering by one unit cell. This is the amplitude of the vector  $\mathbf{F}(hk\ell)$ . It can be derived that

$$\mathbf{F}(hk\ell) = V \int_{x=0}^1 \int_{y=0}^1 \int_{z=0}^1 \rho(xyz) \times \exp[2\pi i(hx + ky + \ell z)] dx dy dz \quad (2)$$

where  $\mathbf{F}(hk\ell)$  is called the [structure factor](#) for reflection  $(hk\ell)$ . It can be regarded as a vector with amplitude  $F$  and phase angle  $\alpha$ :  $\mathbf{F} = |F| \exp[i\alpha]$ ;  $x$ ,  $y$  and  $z$  are fractional coordinates along the unit cell axes  $\mathbf{a}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$ ;  $V$  is the volume of the unit cell and  $\rho$  the electron density distribution in the cell. Using Eq. (2), the scattering by the crystal can be calculated from the electron distribution in the unit cell. In the reverse direction,  $\rho(xyz)$  is derived from the  $\mathbf{F}(hk\ell)$  's of all scattered beams. This is accomplished by applying a mathematical procedure called Fourier transformation. If Eq. (2) is true, then it is also true that

$$\rho(xyz) = \frac{1}{V} \sum_h \sum_k \sum_\ell |F(hk\ell)| \times \exp[-2\pi i(hx + ky + \ell z) + i\alpha(hk\ell)] \quad (3)$$

using Eq. (3), the electron distribution in the unit cell can be calculated provided that the values of  $F(hk\ell)$  and  $\alpha(hk\ell)$  are known for all reflections. The  $F(hk\ell)$  's are easily found because, apart from correction factors, they are proportional to the square root of the scattered intensities  $I(hk\ell)$ . The phase angles  $\alpha(hk\ell)$  cannot be measured, however (see [Phase Problem](#)), but must be derived indirectly by one of these four methods:

1. [direct methods](#)
2. [isomorphous replacement](#)
3. [molecular replacement](#)
4. **MAD**, multiple wavelength anomalous dispersion

Direct methods are the method of choice for small compounds, but they are not yet sufficiently developed for macromolecules. The other three methods are applicable to large molecules.

The degree of detail in the electron density map calculated with Eq. (3) depends on the number of reflections included in the summation. The greater the number, the greater the accuracy. The greater the volume of [reciprocal space](#) included, corresponding to reflections from smaller spacings in the crystal, the greater the atomic detail in the electron density map. This is the "resolution" that is always quoted with every structure determination by X-ray crystallography.

#### Bibliography

1. W. Friedrich, P. Knipping, and M. Laue (1912) Sitzb. kais. Akad. Wiss., München, 303–322.
2. T. L. Blundell and L. N. Johnson (1976) *Protein Crystallography*, Academic Press, New York, London, San Francisco, p. 10.
3. D. W. Green, V. M. Ingram, and M. F. Perutz (1954) Proc. Roy. Soc. A **225**, 287–307.
4. J. R. Helliwell (1992) *Macromolecular Crystallography with Synchrotron Radiation*, Cambridge

University Press, Cambridge.

5. W. L. Bragg (1913) Proc. Cambridge Phil.Soc. **17**, 43–57.

### Suggestions for Further Reading

6. J. P. Glusker, M. Lewis, and M. Rossi (1994) *Crystal Structure Analysis for Chemists and Biologists* VCH, New York, Weinheim, Cambridge.

7. J. Drenth (1999) *Principles of Protein X-ray Crystallography*, Springer, New York.

## X-Ray Scattering

X-ray scattering is an important structural tool for studying biological [macromolecules](#). When X-rays are scattered elastically and coherently by atoms, the scattered waves interfere in a manner related to the spatial distribution of the atoms. The most widespread application of X-ray scattering in structural biology is [crystallography](#), in which the scattering from molecules in three-dimensional crystals is used to determine their structures at high resolution. For biological molecules in unordered or partially ordered systems, lower-resolution structural information is obtained that can also provide critical insights into function.

### 1. Interactions of X-Rays with Matter and the Information Content of Scattering Data

X-rays are photons and can be produced with wavelengths in the range 0.1 to 100 Å (1 Å = 10<sup>-10</sup>m). They are ionizing radiation and therefore can be quite damaging to biomolecules. Like other electromagnetic radiation (such as light), they propagate as waves and are scattered via interactions with the electrons in a sample. The interference of coherent, elastically scattered X-rays interacting with biological samples gives structural information at different levels of resolution, depending on the degree of order in the system. Analysis of the [scattering intensity distribution](#) at small angles of unordered biological macromolecules in solution gives information on their overall shapes and the associations they form. [Small-angle scattering](#) data yield structural parameters such as the [radius of gyration](#),  $R_g$ , molecular weight,  $M$ , and the vector length distribution function,  $P(r)$ , for the scattering particle. These experiments can provide critical information on the conformational changes and associations formed during biological function. For example, small-angle X-ray scattering data from the calcium-binding protein [calmodulin](#) (1) revealed that the [a-helix](#) connecting its two globular domains seen in the crystal form was flexible in solution. This flexibility has proven the key to calmodulin binding and subsequent activation of its regulatory targets (2). X-ray scattering from molecules that are ordered in one, two, or three dimensions is convoluted with the repeating lattice structure, giving rise to diffraction maxima in the scattering pattern. For example, oriented DNA molecules give rise to X-ray diffraction patterns that in 1953 yielded the double helical [DNA structure](#) (3). X-ray diffraction data from natural [membranes](#) were used to establish the lipid bilayer as the predominant structural component of membranes (4). The bilayer structure of membranes give rise to diffraction maxima relating to repeat distances both perpendicular and parallel to the plane of the bilayer (5) that can be important in understanding, for example, the influence of lipid type (6), protein content (7), toxins (8), drugs (9, 10), or molecules like cholesterol (10, 11) on the membrane structure and fluidity. Membranes can also give rise to two-dimensional diffraction patterns from ordered molecules, such as proteins within the planes of the bilayers [eg, [bacteriorhodopsin](#) (7, 12)]. These data can give information on the protein's structure and and/or position within the bilayer. [Crystallography](#) is the interpretation of X-ray diffraction data from three-dimensional crystals of molecules in terms of high-resolution molecular models. In the case of biological macromolecules,

the resolution obtained can be true atomic resolution, although more often it is at the level of individual chemical groups. Three-dimensional crystallographic data on biological macromolecules provide the level of detail that can ultimately reveal the chemistry of biomolecular mechanisms.

## 2. X-Ray Sources

Conventional laboratory X-ray generators work by exciting the characteristic emission lines of metallic elements. An electron beam of a few tens of kilovolts is directed onto a metal target and knocks electrons out of low-lying orbitals, creating “holes.” Because these holes are refilled by electrons from higher orbitals, sharp emission lines are produced. Typically about half the X-ray energy will be converted into these lines. The other half contributes to a white background. The most common metal target used for X-ray production in a laboratory for biological studies is copper, and its  $K_{\alpha}$  emission line ( $\lambda = 1.54 \text{ \AA}$ ) is usually used. The nearby copper  $K_{\beta}$  can be removed using a nickel filter. Synchrotrons provide alternative sources of X-rays in the form of very high-intensity, well-collimated X-ray beams. Synchrotron radiation is emitted in short pulses (a few tenths of nanoseconds) at high frequencies (megahertz) when electrons are accelerated at very high speeds in a circular trajectory using electromagnets. This “white” radiation can be monochromatized for scattering experiments. Synchrotron technology requires large facilities, and there are an increasing number of synchrotron user facilities well-equipped with instrumentation for protein **crystallography, small-angle scattering**, as well as diffraction from partially ordered systems. The increased intensities of the synchrotron sources allow for faster studies on smaller samples, or samples that are inherently weak scatterers. In particular, synchrotron intensities make time-resolved studies possible. For example, they facilitate using small-angle scattering to follow protein conformational changes with time on time scales of milliseconds to seconds (13). The time scales and intensities of synchrotron radiation also facilitate [Laue Diffraction](#) (14) studies that utilize the white spectrum of the source to collect crystal diffraction patterns rapidly. X-ray damage to biological samples can be quite severe with synchrotron intensities. For solution scattering experiments, radiation-induced aggregation can also be a problem. These effects can be minimized by keeping samples at low temperatures, minimizing measurement times, using reducing agents to minimize free-radical concentrations, and/or lowering X-ray intensities by the use of attenuators.

## Bibliography

1. D. B. Heidorn and J. Trehwella (1988) *Biochemistry* **27**, 909–915.
2. J. K. Krueger et al. (1998) *Biochemistry* **37**, 13997–14004.
3. J. D. Watson and F. H. C. Crick (1953) *Nature* **171**, 737–738.
4. M. H. F. Wilkins, A. E. Blaurock, and D. M. Engelman (1971) *Nature New Biology* **230**, 72–76.
5. A. E. Blaurock (1982) *Biochim. Biophys. Acta* **650**, 167–207.
6. B. A. Lewis and D. M. Engelman (1983) *J. Mol. Biol.* **166**, 211–217.
7. A. E. Blaurock (1973) *Nature (London)* **244**, 172–173.
8. A. Colotto, K. Lohner, and P. Laggner (1991) *J. Appl. Crystallogr.* **24**, 846–850.
9. S. V. Balasubramanian and R. M. Straubinger (1994) *Biochemistry* **33**, 8941–8947.
10. P. R. Mason, D. M. Moiesey, and L. Shajenko (1992) *Mol. Pharmacol.* **41**, 315–321.
11. M. J. Janiak, D. M. Small, and G. G. Shipley (1979) *J. Lipid Res.* **20**, 183–199.
12. A. E. Blaurock and W. Stoeckenius (1971) *Nature New Biology* **233**, 152–155.
13. D. Eliezer et al. (1993) *Biophys. J.* **65**, 912–917.
14. A. Cassetta et al. (1993) *Proc. R. Soc. Lond. Ser. A* **442**, 177–192.

## Suggestions for Further Reading

15. C. R. Cantor and P. R. Schimmel (1980) In *Biophysical Chemistry*, Part II: Techniques for the Study of Biological Structure and Function, W. H. Freeman, San Francisco, pp. 687–819.



16. B. Chance et al. (eds.) (1994) *Synchrotron Radiation in the Biosciences*, Oxford University Press, New York.
17. O. Glatter and O. Kratky (1982) *Small-Angle X-ray Scattering*, Academic Press, New York.

## Xenogeneic

Acceptance of an organ graft from one donor to a recipient is strictly conditioned by their respective genetic constitutions. Identity defines *syngeneic* conditions, which occur with monozygous twins or animals from the same inbred strain. Because there is no genetic difference between donor and recipient, the graft is accepted. Whenever the donor and the recipient belong to the same animal species, but differ in their genetic constitution, the *allogeneic* graft is rejected after 10 to 12 days. Allogeneic graft rejection is an immunological reaction, because it is specific, mediated essentially by [cytotoxic T lymphocytes](#), and controlled by molecules of the [major histocompatibility complex](#) (MHC): HLA in humans and H2 in mice. Secondary responses are apparent by an accelerated rejection after a second graft from the same genetic origin as the first one. A *xenogeneic* donor and recipient belong to different animal species. In this case, the rejection is immediate and particularly violent. The mechanisms of xenogeneic graft rejection are numerous and complex. Hyperacute rejection is due to the presence of natural **antibodies** that bind to the endothelium, thereby inducing **complement** fixation, activation of the endothelium, and initiation of [blood clotting](#). Many efforts have been made in the recent past to understand the mechanism of xenogeneic rejection and to define strategies to overcome it, because it might provide a valuable source of organs for human transplantation. So far, significant results have been gained at least for the hyperacute phase, but enormous difficulties are still to be solved. A second immunological barrier is that of the delayed xenograft/acute vascular rejection, which is still not fully understood. Third, rejection mechanisms similar to those encountered in allograft rejection are also operating, with an increased efficiency. All these barriers have to be solved before xenograft really become a reality. To date, no classical immunosuppressive drug can prevent such rejection. The present efforts aim to induce specific unresponsiveness to the most antigenic foreign constituents, but this remains a very difficult task.

See also entries [Clonal Selection Theory](#), [Antibody](#).

### Suggestion for Further Reading

- H. Auchincloss Jr. and D. Sachs (1998) *Annu. Rev. Immunol.* **16**, 433–470.

## Xenopus

More than a thousand different species of frogs have been characterized. The genus *Xenopus* (“strange foot” in Greek) arose more than 120 million years ago and includes 17 species. These aquatic frogs live in ponds in southern (sub-Saharan) Africa. One of these species, *Xenopus laevis*, has become, together with the [mouse](#) and chick, an attractive vertebrate model system for embryologists and developmental biologists. In fact, most of what is known of vertebrate

development comes from the study of the amphibian embryo. Already popular in Europe, *X. laevis* was first introduced in the United States in the early 1950s as a way to test for pregnancy in humans. When injected subcutaneously into female *Xenopus*, the gonadotropin hormone present in pregnant women's urine would induce her to lay eggs the next day. Unlike most other frogs, egg laying in *Xenopus* is not seasonal; these tests could thus be performed all year round.

For embryologists, a major advantage of this frog is the number and size of its eggs. A single female can lay up to several thousand eggs a day. This is particularly important for modern molecular techniques: the source of biological material such as [complementary DNA](#) (cDNA), **RNA**, or [protein](#) is virtually unlimited. In addition, a *X. laevis* egg is approximately 1 mm in diameter, which makes it one of the largest cells in the world, visible to the naked eye, and amenable to microsurgery and microinjection of molecular factors. Moreover, fertilization of the egg is external, so the amphibian embryo can be studied from the beginning, when the embryo is a single cell. Finally, perhaps one of the most important attributes of this system is the amount of knowledge accumulated over more than a century on the [development](#) of the amphibian [embryo](#). There is an amazingly rich body of literature, from descriptive to experimental approaches, on the development of many different types of amphibian embryos, including *Xenopus*. Associated with this literature is a remarkable cast of characters, including Mangold, Spemann, Roux, His, Holtfreter, Hamburger, and Nieuwkoop. Their work and their vision established the foundations of vertebrate experimental embryology, on which modern molecular embryological approaches are built. Although powerful genetic approaches in *Caenorhabditis elegans* and **Drosophila**, and more recently [zebrafish](#) and mouse, have contributed substantially to our knowledge of early embryonic development, experimental embryology in the amphibian had already established many key features of embryonic development. For example, the concept of embryonic **induction**, so widely accepted today as a means to establish embryonic cell fate, was first defined in amphibians (1). This article presents (1) a descriptive overview of *Xenopus laevis* development, including the concepts of [fate maps](#) and gene maps; (2) a review of molecular approaches currently used in the *Xenopus* system in **oocytes** and in embryos; (3) the use of explants and experimentally perturbed embryos in molecular studies of early development; and (4) genomics, nuclear transplantation, transgenesis, and maternal knockout strategies.

## 1. The *Xenopus* Embryo

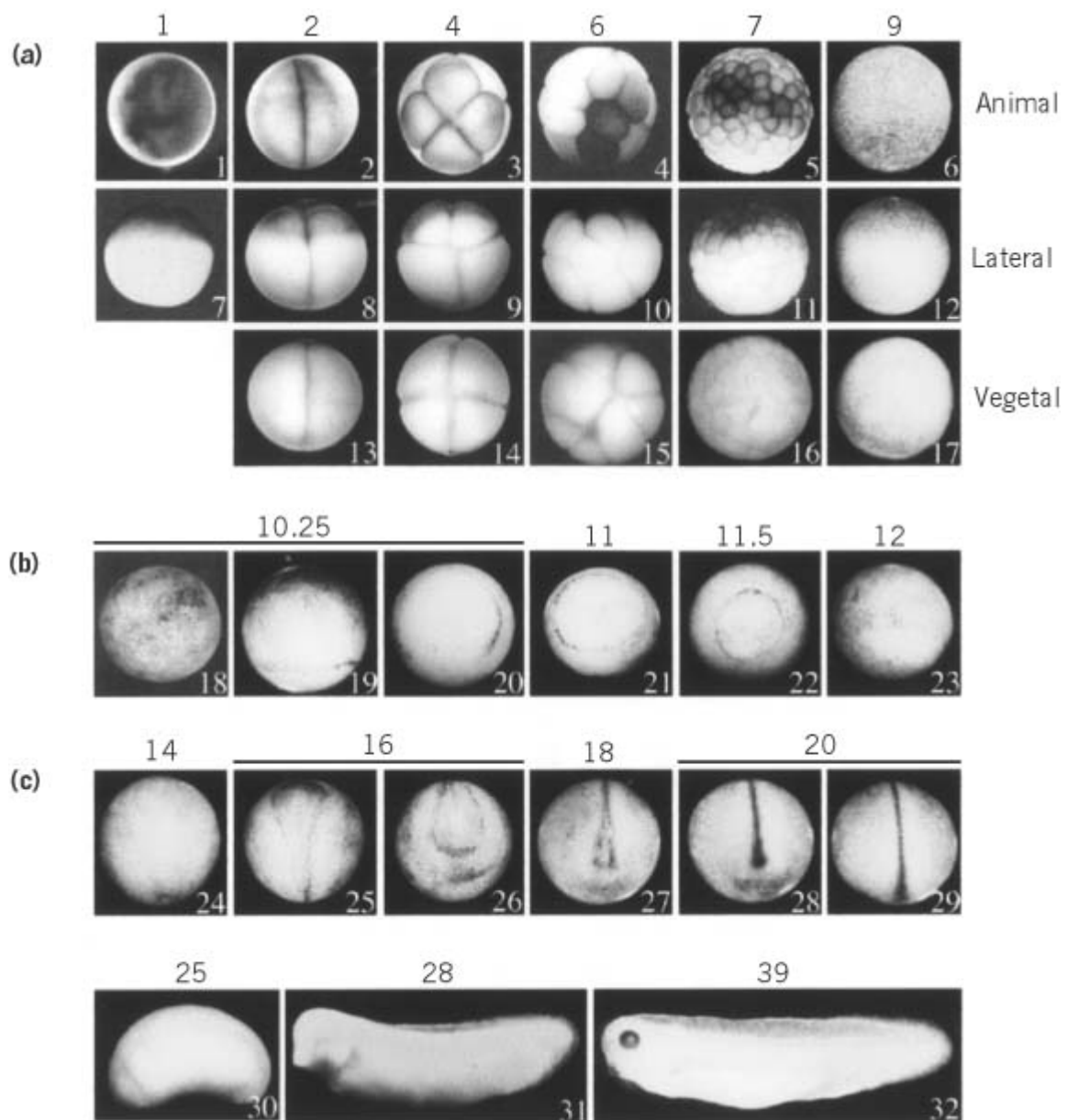
### 1.1. Early Embryonic Development

Large numbers of *Xenopus* embryos can be obtained year-round, either by natural mating or by *in vitro* [fertilization](#). For “natural” mating, both the male and female are injected with human chorionic gonadotropin and left in a quiet dark tank overnight to copulate. A large number of embryos at different stages of development can be harvested the following day. For *in vitro* fertilization, females are injected subcutaneously with the chorionic gonadotropin and left overnight. The next day, a male frog is sacrificed, and fragments of surgically removed testis are used to fertilize eggs gently squeezed from the female. Each squeeze provides a few hundred eggs. Females can be squeezed every hour for several hours. The advantage of *in vitro* fertilization is that all the eggs for a given squeeze develop synchronously; thus, a large number of embryos of a given stage are available. When laid, the eggs have a single axis of cylindrical symmetry running from the pigmented pole at the top, also called the **animal pole**, to the yolky bottom or vegetal pole. This axis of symmetry is imposed maternally during the development of the oocyte (2).

The [sperm](#) can enter at any point around the circumference. The first [cell cycle](#) is unusually long, about 90 minutes at room temperature. About 30 min after sperm entry, a drastic cytoplasmic movement, referred to as *cortical rotation*, takes place: the inside of the egg rotates about 30° relative to the outside. The consequence of cortical rotation is that the original axis of cylindrical symmetry is broken, and one side of the egg is different from the other (2). With 85% accuracy, the side opposite the sperm entry point will become the dorsal side. A series of quick cell divisions follow, without change in the volume of the embryo; thus, the cells become smaller in size following each division (Fig. 1a). During the first 4 hs of development, all cells of the embryo divide synchronously. The length of each cell cycle, after the first, is about 20–30 min at room temperature;

the early [blastomeres](#) of the *Xenopus* embryos have barely enough time to replicate their genome before **cytokinesis** occurs and the next cleavage furrow becomes apparent. The development of the *Xenopus* embryo for the first few hours is entirely under the control of maternal determinants deposited in the egg.

**Figure 1.** Early embryonic stages of *Xenopus laevis* development. (a) The first 6 h of *Xenopus laevis* development. The numbers on top of each panel, or set of panels, represents the developmental stage (45). Each panel is indicated by a number for reference in the text. The first row shows animal pole views; the middle row, lateral views; and the third row, vegetal pole views. The embryo undergoes 13 synchronous cell divisions without an increase in size. (b) Gastrulation. Panel 18 is an animal pole view; 19, a lateral view; and panels 20–23, vegetal posterior views. (c) Neurulation and the formation of a tadpole. Panels 24 and 25 are dorsal views of neural plate and neural groove stage embryos, respectively; panels 26–28 are anterior, head-on views showing the progressive closure of the anterior neural tube; panel 29 is a dorsal view of the neurula with the entire neural tube closed; and panels 30, 31, and 32 are lateral views of late neurula, tailbud, and tadpole, respectively.

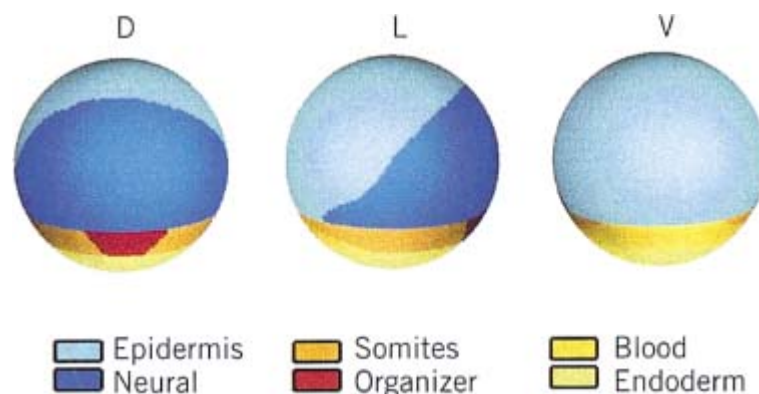


Dramatic changes occur after 13 synchronous cell divisions, when the embryo reaches about 4000 cells (Fig. 1b): (1) the synchronicity of cell division is lost, and embryonic cells divide at different

rhythms in different regions of the embryo; (2) zygotic [transcription](#) begins (3); and (3) finally, cell and tissue movements begin. Among the first of these movements is the spreading of the animal pole region toward the marginal zone, called *epiboly*. Also, the cells of the marginal zone begin to move toward the dorsal side by a process called *convergence and extension*, a process that will continue during **gastrulation** (4). At this stage of development, the three embryonic germ layers—**ectoderm**, **mesoderm**, and **endoderm**—have been specified. These early cell movements set the stage for gastrulation, which begins about 9 h after fertilization.

It is during gastrulation that both the antero-posterior and right-left axes are determined and that cells of the prospective mesoderm and endoderm move to the interior of the embryo. Cells in a subequatorial region on the dorsal side involute first (Fig. 1, panels 19 and 20). This invagination is recognizable morphologically by the formation of a small arc of pigment in the dorsal side, called the *blastopore lip*. The emergence of the blastopore lip defines the dorsal midline. A group of cells located within a 30° angle at each side of the dorsal midline have the capacity to direct development of the entire dorsal axis, including axial mesodermal derivatives, such as somite and notochord, and the entire central nervous system (5). These group of cells are called the “[organizer](#)” (Fig. 2) because classic experiments demonstrated that removing this group of cells from a donor early gastrula embryo and implanting it in the ventral axes (6). One axis is generated by the host organizer; the other one forms under the influence of the ectopic organizer implant.

**Figure 2.** Gastrula fate map. Three views of the superficial gastrula, with the fate map adapted from Keller (1975) (D—is a dorsal view, L—a lateral view, V—a ventral view). The ectoderm is derived from the top of the embryo (animal pole, shades of blue). Mesoderm is derived from the equatorial region (marginal zone, shades of red), and endoderm from the bottom (vegetal pole, yellow). See color insert.



Involution gradually spreads from the dorsal side to the mediolateral and ventral region of the embryo; this is manifested by the expansion of the blastopore lip arc (Fig. 1, panels 20–23). The lip finally becomes a closed circle when cells located in the entire circumference of the embryo have invaginated. The first cells that involute on the dorsal side are the prechordal plate (or head) mesoderm; these cells will contribute to mesodermal components of the craniofacial structures and are therefore anterior in nature. Cells that involute later, on the dorsal side, give rise to somites and notochord (4). As gastrulation proceeds, the blastopore ring becomes smaller and smaller until it closes to engulf the entire vegetal pole (endoderm). The closure of the blastopore demarcates the end of gastrulation and the beginning of neurulation. The blastopore will ultimately become the anus of the tadpole. Thus, at the end of gastrulation the embryo has successfully brought the entire endoderm inside, covered its surface with ectoderm, and placed mesoderm in between.

At the end of gastrulation, the cells of the dorsal ectoderm thicken to form the neural plate (which is two cell layers thick). This is the first morphological sign of the nervous system (Fig. 1c). The neural

plate is flanked laterally by the [neural crest](#) and anteriorly by the sensory placodes and an amphibian-specific organ, the cement gland. Within the neural plate itself, mature neurons emerge in three stripes on each side of the dorsal midline (7). The ones closest to the midline will become motor neurons, those in the intermediate strip will contribute to the interneurons and the most lateral stripe will give rise to sensory neurons. During neurulation, morphogenetic movements elevate the lateral edges of the neural plate to give rise to a neural groove (8). The neural groove will ultimately close to generate the [neural tube](#), with a brain at the anterior and the spinal chord at the posterior end. In *Xenopus laevis*, neural crest migration begins at the neural groove stages and continues following neural tube closure. The eye buds and optic vesicles, as well as the cement gland, become morphologically distinct. Neurulation ends with the closure of the neural tube and the appearance of a tailbud at the posterior end of the embryo. The cells of the heart primordia that flank the organizer on both sides meet in the ventral midline and fuse to form a tube that will later fold to form a functional heart. The progenitors of other organs have begun their differentiation as well. These include the embryonic kidney (pronephros) from lateral mesoderm, as well as the blood islands in the ventral mesoderm that will give rise to cells of the entire hematopoietic pathway.

By the tailbud and tadpole stages, the embryo has a functional nervous system (Fig. 1, panels 31 and 32). Most neuronal axons have reached their targets, and the embryo is motile. A few hours later in the tadpole, the folded heart pumps blood derived from the blood islands into the newly formed arteries and veins. The development of internal organs, such as pancreas, gut, and intestine, is complete. The maternal yolk, which sustained embryonic life up to this point, is exhausted, and the tadpole begins feeding. Concomitant with this is the emergence of specific behavior.

The tadpole continue to grow in size until metamorphosis occurs. During this dramatic stage of the frog's life, major remodeling of almost the entire body changes the morphology and the physiology of the tadpole, to generate a frog. There has been tremendous progress in the understanding of the molecular basis of metamorphosis in amphibians, and excellent reviews are available on this topic (9).

## 1.2. Fate Maps and Developmental Commitment

[Fate maps](#), a fundamental concept in embryology, tell us what regions of the early embryo will contribute to later structures. Of course, the fate map does not tell us about how those structures are formed or their state of commitment, only about their spatial origin. In *Xenopus laevis*, several detailed fate maps of different stages of embryogenesis, as early as the 16- and 32-cell stage, are available (10). The pioneering work of Keller has provided a very high resolution fate map of gastrula and neurula embryos (11). In addition to fate maps of the whole embryo, there are fate maps of the neural plate and sensory placodes (12). At the gastrula stage, the fate maps described in *Xenopus laevis* are very similar to those described in other amphibians, including *X. tropicalis*, *X. borealis*, *Rana pipiens*, and *Triturus* (newt). Fate maps in the amphibian are usually constructed by marking cells with vital dye or by injecting them with other lineage tracers, followed by a statistical assessment of where the labeled cells end up later in the tadpole. Knowledge of the amphibian fate maps has been key in the understanding the molecular processes involved in early development (see text below).

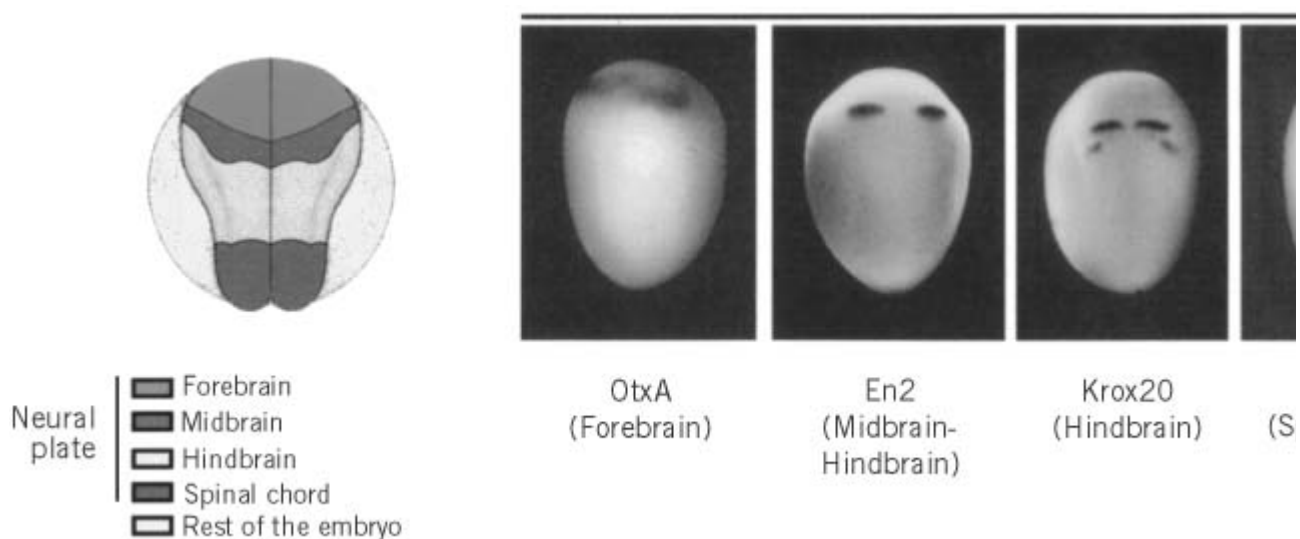
To determine the state of **commitment** of a region of the embryo, approaches such as specification and determination assays can be used (13). In a specification assay, different regions of the early embryo are removed by microsurgery and cultured in a simple saline solution. Because a source of food is contained within each cell, these cultures do not need extrinsic factors to develop. By comparing the behavior of the explants cultured in isolation and the behavior of this region in the context of the embryo, it is possible to assess the state of commitment of that tissue. If the explant in culture gives rise to the same structure as in the embryo, then the explant is **specified**. Additionally, if the explants are derived at different times, it is possible to assess the state of commitment of a tissue at a given time toward a given fate. In determination assays, a cell or group of cells are transferred from a donor embryo into a different environment of a host embryo of the same or different developmental stage, and their development is compared to when they are cultured alone

(as in specification assays). A region of an embryo is **determined** if it can give rise to the same structures regardless of its environment (in culture alone or placed in different region of the embryo). Slack (1984) suggested that analysis of the state of commitment can be expanded to include not only tissue and single cells, but nuclei also. Molecular analyses of these commitment states are used in molecular embryology to assess the molecular pathways involved in cell commitment and cell fate determination (see section on use of *Xenopus* embryos and embryonic explants in molecular biology, below).

### 1.3. Gene Maps

While early embryonic cells can appear morphologically identical, they do not necessarily express the same repertoire of **genes** (Fig. 3). This observation has allowed the unveiling of what Kirschner and Gerhart call the “second anatomy” of the embryo, based on the region-specific distribution of gene products (14). Sometimes the gene map correlates very well with fate maps, and sometimes it does not. When a given gene demarcates the same territory as the fate map for a given time in development, the gene can be considered to be a molecular marker for those cell types (Fig. 3). In molecular embryology, the expression of these genes is used as a testimony for the formation of the cell type in which they are expressed. It is very important to remember that, because some genes have an extremely dynamic pattern of expression, a molecular marker for a given cell type will only attest to the presence or absence of that cell at the appropriate time of expression. For example, the gene brachyury (*Xbra* in Ref. 15) demarcates the entire marginal zone of the early gastrula in *Xenopus*; fate mapping studies have established that the entire mesoderm will derive from this region. A few hours later, at the end of gastrulation and the beginning of neurulation, *Xbra* is expressed in the notochord and as a ring around the blastopore. Thus, *Xbra* is a panmesodermal marker when its expression is measured at gastrula stages; however, at neurula stages *Xbra* is a marker of notochord and blastoporal ring. Other markers are less dynamic in their expression; for example, the neural [cell adhesion molecule](#) (NCAM) (16) is expressed in the neural plate from the onset of its formation and remains entirely neural-specific until at least very late tadpole stages.

**Figure 3.** Neurula fate map and molecular anatomy. The fate map shown in the left panel defines four domains along the axis of the neural plate: forebrain, midbrain, hindbrain, and spinal chord (adapted from Egelson and Harris (12)). The right shows the expression of region-specific molecular markers within the neural plate. From left to right: the homeobox marker of presumptive forebrain; Engrailed-2 (En2) is a homeobox gene that demarcates the midbrain-hindbrain junction; factor Krox20 is a marker of rhombomers 3 and 5 of the hindbrain; the homeobox HoxB9 demarcates the presumptive sp



## 2. Molecular Tools for the Study of Early *Xenopus* Development

The combination of experimental embryological approaches and molecular techniques has armed modern embryologists with powerful tools to dissect the molecular basis of early vertebrate development. Traditional [lambda phage](#)-based gene [libraries](#) are available from almost every stage of embryonic development, as well as from different regions of the embryo at a given stage. More recently, stage- and region-specific [expression libraries](#) have been developed that introduce a very powerful means to **clone** the factors responsible for a given activity. Expression cloning is discussed below (see section on *Xenopus* embryos for molecular studies). As with other animal model systems, screens for gain of function, loss of function, and modifier using expression libraries are available in *Xenopus*. For these approaches, nucleic acids or proteins (wild-type or mutated forms) are **microinjected**, either alone or simultaneously with pools of clones from [cDNA libraries](#), into the oocytes, eggs, or blastulae (17). Injected oocytes and embryos are allowed to develop to the desired stage, after which they are submitted to **phenotypic**, histological, and molecular analysis. While gain-of-function experiments use ectopic expression of wild-type or constitutively active forms of the proteins encoded by the injected DNA or RNA, loss-of-function experiments utilize dominant-negative proteins or injections of **antisense** RNA, antisense DNA oligonucleotides, **ribozymes** or neutralizing **antibodies**. An in-depth review of the use of dominant-negative mutant approaches in embryology has been presented elsewhere (18). It is important to remember that the specificity of effects derived from antisense RNA, ribozymes, or antibody injections need to be interpreted with extreme caution. It is also important to control for the *in vivo* stability of micro-injected reagents.

### 2.1. Microinjection of DNA

Up to 100 pg of supercoiled, nicked, or linear DNA can be microinjected into oocytes or early embryos (17). The [DNA topology](#) has no effect on expression. In oocytes, the size of the nucleus (germinal vesicle) allows direct injection of the DNA into the nucleus. This type of microinjection has been used for transcriptional studies, as well as examination of RNA **splicing** and [DNA recombination](#). In embryos, transcription begins when the blastula reaches about 4000 cells (see text above); consequently, injected DNA is not transcribed for the first few hours of development. Because of this property, injection of DNA, rather than RNA, is sometimes used to study later developmental events. The injected DNA is stable until gastrula-early neurula stages, after which it gradually decays. For reasons that are not well understood, only a fraction of the embryonic blastomeres that inherit the DNA will ultimately express it; thus, expression driven from injected DNA is always **mosaic**. As is the case with genetic studies, mosaic expression has its own advantages when contrasted with uniform expression.

The requirement for **promoters** is not very stringent in *Xenopus laevis*: most vertebrate and **virus** promoters work well in oocytes or embryos. However, with a few exceptions, tissue-specific expression cannot be achieved in this type of transient expression experiment. Another type of DNA injection involves the injection of single-stranded antisense oligonucleotides. In these experiments, regular or chemically modified antisense oligonucleotides are microinjected in the oocyte or embryos with the aim of binding to and eliminating the function of endogenous, sense mRNA.

### 2.2. Microinjection of RNA

The microinjected RNA can be of a biological origin, for example, **poly bAb<sup>+</sup>** RNA isolated from different tissue samples. More frequently, however, single RNA or a population of RNAs are generated synthetically by *in vitro* transcription. In this case, specific **vectors** that contain prokaryotic transcriptional promoters, such as those of SP6, T3, or T7 are used. In addition, these vectors usually contain short stabilizing 5' and 3' untranslated regions flanking the cloning sites. Some vectors also add a polyA-polyC tail to the 3' end of transcripts to increase their stability (17). Synthetic RNAs obtained by *in vitro* transcription of a single clone or pools of clones from expression libraries can be injected into oocytes and embryos. In oocytes, up to 50 ng of RNA can be injected, while embryos can tolerate up to 5 ng of RNA. Injected RNA or DNA may encode either wild-type or mutant versions of proteins. These mutant proteins can include constitutively active forms of known proteins, for example, constitutively active **receptors**, **signal transducers**, or

[transcription factors](#), which will activate signaling pathways or target genes in a ligand-independent fashion. Alternatively, they can encode dominant-negative versions of embryologically active proteins to generate loss of function phenotypes.

### 2.3. Microinjection of Proteins

Finally, injection of proteins, such as [immunoglobulins](#), have also been reported in *Xenopus* embryos (19, 20). In these experiments, neutralizing antibodies against specific proteins are microinjected with the aim of eliminating the function of the endogenous protein. Again, as in the case with antisense techniques, interpretation of **phenotypes** should be done with great care. Rescue experiments, in which the antigen and antibody are coinjected, provide the most stringent control for specificity.

## 3. The *Xenopus* Oocyte

Because of its size, the ease with which it can be obtained, and the relative ease with which it can be cultured, the *Xenopus* oocyte has become popular with researchers with a wide variety of interests. The fully mature *Xenopus laevis* oocyte is a large cell of about 1 mm in diameter. The oocyte is surrounded by a vitelline membrane and several layers of follicular cells. For most experimental conditions, the latter are removed, either mechanically or by enzymatic treatment, for example with collagenase. These defolliculated oocytes still have the vitelline membrane attached, which provides mechanical stability.

Oocytes can be used for a variety of functional studies, a few of which are described below. Up to 50 ng in a total of 50 nL of RNA can be microinjected into the **cytoplasm** of the oocyte. Unlike early embryos, the *Xenopus* oocyte is transcriptionally active. For DNA injections, most, if not all, mammalian promoters work well for expression in the context of the oocyte. Similarly, the [polyadenylation](#) signals from **SV40 virus** or mammals function well in the frog. Up to 10 ng of DNA can be injected directly into the germinal vesicle (nucleus) of the oocyte.

The *Xenopus* oocyte expression system has been used successfully for expression cloning, the study of individual proteins, and large-scale production of secreted proteins. The success of expression cloning depends on both the quality of the library and the specificity of the assay. For example, [membrane proteins](#) such as **ion channels** have been cloned by injecting oocytes with populations of RNA or DNA, followed by electrophysiological analysis to determine whether the desired channel is present in the pool. Further sib-selection of active fractions leads to the isolation of the clone encoding the channel. For isolation of receptors for [hormones](#) or [growth factors](#), similar approaches can be used: binding of a labeled ligand is used to identify positive fractions.

Alternatively, individual cell-surface proteins can also be studied by injection of single RNA or cDNA. This has traditionally been used to study mutant forms of channels and receptors. Finally, although most oocyte experiments involve the analysis of membrane proteins, the *Xenopus* oocyte has been used more recently for the production of large amounts of secreted factors. In fact, one of the most efficient ways to obtain large quantities of highly-purified secreted growth factors, from any source, is to harvest the conditioned medium of oocytes injected with synthetic RNA encoding these factors.

Because of its size, the oocyte is also one of the favorite systems used by cell biologists interested in subcellular trafficking. Thus, trafficking between the [endoplasmic reticulum](#) and the [Golgi apparatus](#) in the secretory pathway, nuclear translocation, and subcellular localization of RNA have all been explored using *Xenopus* oocytes (for a detailed review, see *Methods in Cell Biology*, 1992).

## 4. Use of *Xenopus* embryos and Embryonic Explants in Molecular Biology

Examples of a few assays that are commonly used in *Xenopus* to identify molecular players involved in embryological functions are presented here. It is important to emphasize that, while these techniques are widely used, new approaches emerge continuously that are custom fitted to address



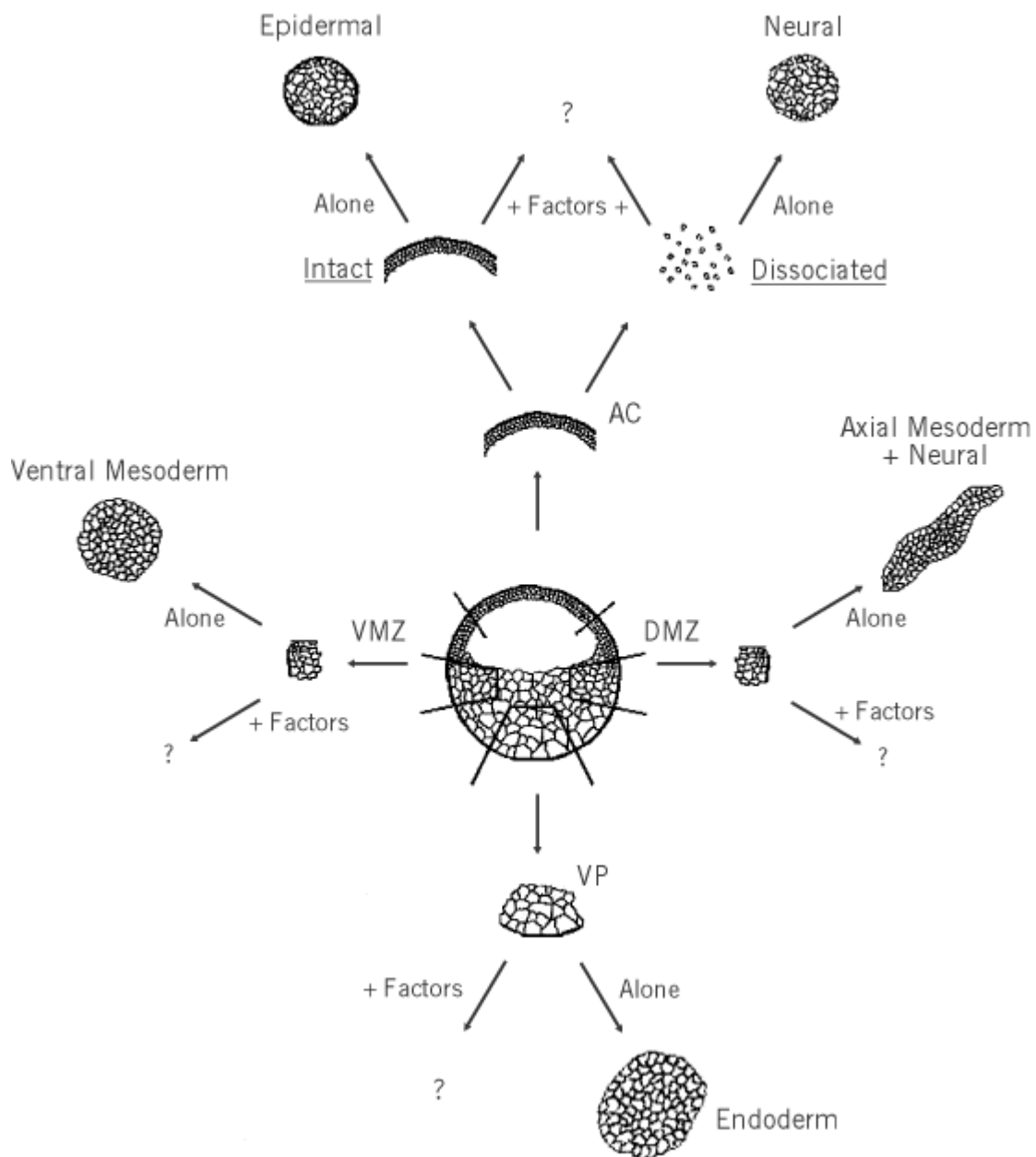
specific questions.

#### 4.1. Embryonic Explants

Specification studies combined with the analysis of cell-type-specific markers have provided a simple but powerful means of identifying factors involved in embryonic development. Explants can be cultured alone or in the presence of factors; alternatively, explants from different regions of the embryo of the same stage or different stages can be recombined to assess their influence on each other, as illustrated by the following examples for each case.

Although many type of explants can be derived from different regions of the amphibian embryo at various times of development, the most frequently used explants are derived from the early blastula or gastrula (Fig. 4). There are three types of explants traditionally used: from the ectoderm (or animal caps), mesoderm (ventral marginal zone, lateral marginal zone, and dorsal marginal zone), and endoderm (or vegetal poles).

**Figure 4.** Embryonic explants. A schematic cross section of a late blastula stage embryo (~4000 cells) is shown in the center. When the animal cap (AC) is removed and kept intact, it will differentiate as epidermis; if dissociated, it will give rise to neural tissue. The explanted dorsal marginal zone (DMZ) cultured in isolation will form axial mesoderm, such as notochord and somites, as well as neural tissue. The vegetal pole (VP) differentiates as endoderm (gut) in culture. The ventral marginal zone (VMZ) will form ventral mesoderm derivatives, such as blood cells, in culture. The behavior of these explants can be challenged in culture by the addition of molecular factors. These factors can be added to the explant culture after its isolation, or they can be presented as injected RNA or DNA at early stages of development by microinjection in the appropriate region of the embryo.



The most popular explant assay is the animal cap assay. When explanted at blastula or early gastrula stages and cultured alone, animal cap ectoderm differentiates as epidermis (Fig. 4). However, these early animal cap cells are pluripotent and are able to change their fate in response to exogenous factors. These factors can be provided as soluble protein directly incubated with the explants (21) or as synthetic RNA injected in the animal hemisphere during the first three embryonic cleavages (22). The injected RNA remains confined to the site of injection, so one can target the RNA to different regions of the embryo. It was the use of *Xenopus laevis* animal caps that first allowed the characterization of molecular candidates involved in embryonic induction. A variant of the same assay involves the dissociation of the animal cap explants in  $\text{Ca}^{2+}$ - and  $\text{Mg}^{2+}$ -free medium. When the cells of the animal cap are dissociated and maintained for several hours in culture, they all adopt a neural fate (Fig. 4 (23)). The type of neural tissue obtained by dissociation of animal cap explants is always anterior in character, including the cement gland (24). Animal caps assays have been used to answer a variety of questions concerning the study of inducers, [morphogens](#) (which are substances that elicit different fates at different concentration), and modifiers (factors that cannot induce specific cell fates on their own but have the ability to modify the character of cell types induced by other

factors) (25). For example, when varying concentrations of the growth factor *activin* (which belongs to the [transforming growth factor](#) (TGF)  $\beta$  ligand superfamily) are applied to dissociated cells, it can induce mesoderm and endoderm in intact animal caps and can display morphogen-like properties. At low concentrations, activin will induce ventral mesoderm; at higher thresholds, axial mesoderm and ultimately endoderm. Growth factors belonging to the Wnt superfamily act instead as modifiers; they cannot induce mesodermal cell fates on their own, but they can change the type of mesoderm induced by activin. Inherent to an understanding of inducers and modifiers is the concept of competence, which relates to the window of time during which a given population of cells can respond to an inductive signal. The animal cap has also provided a means to address competence. Even though our knowledge of competence at the molecular level remains sparse, factors such as [fibroblast growth factor](#) (FGF) have been suggested to regulate the capacity of embryonic cells to respond to mesoderm inducers such as activin (25). It is important to note that animal cap explants have dorsal–ventral polarity (imposed upon them during the first cell cycle by cortical rotation), which may influence the response of these cells to inducers or modifiers. This polarity can be eliminated by using animal caps derived from hyperventralized or hyperdorsalized embryos.

Intact ventral marginal zone explants develop as ventral mesodermal cell types, including those of blood; once dissected, these explants adopt a rough spherical shape in culture. Dorsal marginal zone explants, however, elongate in culture and produce axial mesoderm, such as notochord and somites, as well as neural tissue. Ventral and dorsal marginal zone explants are useful for the assessment of the modifier activity of molecular vis-à-vis mesodermal patterning. Substances that can change the fate of these explants are good candidate modifiers.

Finally, vegetal pole explants of blastula or early gastrula embryo give rise exclusively to endoderm. These explants have been used to assess induction, patterning, and competence in the context of endoderm formation and have been gaining in popularity of late.

#### 4.2. Recombinant Explants

Explants can also be recombined in different ways to evaluate a wide variety of embryological activities, such as induction, patterning, and cell autonomy. One of the most informative recombinant explants was originally performed by Nieuwkoop in 1969 (26). It was known that animal cap explants cultured alone give rise to only ectoderm (epidermis) and that vegetal pole explants give rise solely to endoderm (gut). Nieuwkoop showed that these explants cultured together also form mesoderm. This was the first demonstration of mesoderm induction in the vertebrate. Later experiments using lineage tracers (see text below) established that the vegetal pole was the source of the mesoderm-inducing signal; animal pole cells respond to these signals by changing their fate from ectoderm to mesoderm. Once it was established that the animal pole cells could respond to mesoderm-inducing signals, simple modifications of this assay allowed the molecular characterization of the molecular factors involved. The vegetal pole was replaced by conditioned medium from different cell lines or purified factors; the animal pole cells incubated in the presence of these factors were assayed for the formation of mesoderm.

In general, recombinant explants that generate more tissue types when recombined together than when cultured independently provide the basis for bioassays aimed at isolating the molecules involved. Similar approaches have also been used in the nervous system to recombine anterior and posterior neural plate tissue. Anterior neural plate explants, cultured alone, develop as forebrain, while posterior explants develop as spinal chord; when recombined, intermediate neural tissue also develops, including midbrain and hindbrain (27). Another famous explant system was introduced by Keller and co-workers. Known as “Keller” explants, they include both dorsal mesoderm and ectoderm, contacting each other only at a narrow boundary (28). These explants have been used to study the way planar inductive and patterning signals operate to generate the nervous system (29). It is also possible to recombine tissues of different age. For example, notochord from a neurula or tailbud can be sandwiched together with the animal cap of a blastula stage embryo (like a hotdog, where the bread is the animal cap and the meat is the notochord) to study vertical signals from the axial mesoderm that influence neural patterning in the ectoderm (30).

Heterospecific recombinant explants between *Xenopus* and chick and between *Xenopus* and mouse embryonic tissues have also been described (31, 32). In these cases, particular care should be taken to accommodate the physical requirements for the proper survival of the recombinants, such as temperature of incubation (different between the cold-blooded amphibian and the warm-blooded amniotes) and the ionic strength of the solution (50 mM salt for the amphibian tissue versus 150 mM for amniotes). The use of these heterospecies explants has led to the suggestion, for example, that the molecular natures of neural-inducing signals are conserved from amphibians to birds. Theoretically, therefore, an enormous variety of recombinant explants can be generated from amphibian tissue. The rationale behind their design depends solely on the questions asked and the creativity of the investigator.

#### 4.3. *Xenopus* Embryos for Molecular Studies

In addition to the use of embryonic explants to isolate and characterize embryonic factors, expression cloning strategies have been successfully used in intact *Xenopus* embryos. In these approaches, fractionated cDNA libraries are subjected to *in vitro* transcription to generate synthetic RNA that is microinjected in different regions of either wild-type or perturbed embryos (see text below). Active fractions that generate a specific phenotype can be fractionated further by sib-selection until a single clone is obtained. These approaches have been extremely successful in isolating neural and mesodermal inducers and modifiers (33, 34). Injection in the animal pole, marginal zone, or vegetal pole of a one- or two-cell stage embryo will affect the differentiation of ectoderm, mesoderm and endoderm, respectively.

Expression libraries can also be used to rescue experimentally perturbed *Xenopus* embryos. Three types of perturbed *Xenopus* embryos are traditionally used: ventralized embryos, dorsalized embryos, and exogastrulae. In *Xenopus*, there is a link between the dorsal axis and anterior structures, as well as a link between the ventral and posterior structures (35). Ventralized embryos are obtained by UV irradiation of the vegetal pole during the first cell cycle (Fig. 5). This treatment blocks cortical rotation; thus, the dorsal side is never specified, and the embryo develops without a dorsal axis or head. Screening expression libraries based on their ability to rescue the ventralized phenotype of UV embryos has allowed the characterization of a large number of genes involved in the differentiation of dorsal structures. At the other end of the spectrum, incubation of the early blastula embryo (8- to 16-cell stage) in 300 mM LiCl will hyperdorsalize the embryo so that all ventral and posterior structures are eliminated (Fig. 5). In extreme cases, the embryo develops as a radially symmetrical head with retinal pigmented cells all around the circumference. Rescue of these embryos has been reported by injection of a single RNA, which suggests that these embryos could also be used for expression cloning strategies (33). Finally, exogastrulae represent another form of perturbed embryo in which, instead of invaginating, the cells of the gastrulae move out of the embryo (16, 36). Although these embryos have been used to study ectodermal and mesodermal interactions, they could also theoretically be used in expression cloning experiments, especially for screening genes involved in morphogenetic movements.

**Figure 5.** Experimentally perturbed *Xenopus* embryos. Varying amounts of exposure to UV irradiation during the first cell cycle generates embryos that are gradually ventroposteriorized (left panels), while treatment with LiCl for varying times produces embryos that are dorsoanteriorized. This range of phenotypes represents the dorsoanterior index (DAI) (46) indicated by the number next to each embryo. In this scale, the normal embryo has a DAI of 5 (in the middle), numbers below 5 represent ventralized phenotypes and numbers above 5 dorsalized embryos. At the extremes, an embryo with a DAI of 0 entirely lacks head and dorsal axial structures while an embryo with a DAI of 10 contains only radially symmetrical head structures.



## 5. Genetics

### 5.1. *Xenopus* Genomics

The size of the *Xenopus laevis* genome has been estimated to be approximately  $3 \times 10^9$  bp (37), which is comparable to that of the human genome. The somatic genome has 36 chromosomes (18 pairs). *Xenopus laevis* is a pseudo-**tetraploid** organism: evidence for tetraploidy is based on both comparison of DNA content between *X. laevis* and other Pipidae and the fact that most, but not all, genes examined are duplicated (37). About 20–30% of the entire *X. laevis* genome consists of [repetitive DNA](#), and more than 2500 genes have been cloned from *X. laevis* to date and are available in the [databases](#). In addition, the entire 17,553-bp [mitochondrial DNA](#) genome of *X. laevis* has been sequenced and shown to encode 13 proteins, 22 [transfer RNA](#), and two **ribosomal RNA** (38). There is even a very rough provisional **linkage map** for *X. laevis* (37). Finally, the egg contains 4 ng of total maternal RNA, of which approximately 1% (40 pg/egg) is mRNA.

### 5.2. Genetic Approaches in *Xenopus*

Despite the fact that there are 40 recessive mutants known in *X. laevis* (39), its polyploidy (see text above) and large number of [pseudogenes](#) cause the system to be not accessible presently to classic genetic analysis. Nevertheless, a large number of approaches have recently made progress toward a genetic analysis of *Xenopus* embryonic development. Three such approaches will be described: nuclear transplantation, transgenesis, and maternal knock-outs. Recent work targeting *Xenopus tropicalis* as a bona fide genetic system have been initiated by several groups, however, and classic mutagenesis approaches in progress with this frog should yield their first fruit soon.

#### 5.2.1. Nuclear Transplantation

While nuclear transplantation in mammals, and the generation of the clone Dolly, the lamb have recently received a lot of attention, it is important to remember that the first successful nuclear transplantation was performed in frogs. In 1952, Briggs and King, using the frog *Rana pipiens*, reported the first successful transfer of nuclei from a blastula to an egg and obtained embryos that survived until neurula stages. Shortly after, two groups using genetic markers demonstrated that cloning could work in *X. laevis*. It was, however, John Gurdon and colleagues who first obtained a genetically marked, sexually mature cloned frog (40). Nuclear transplantation requires the isolation of nuclei from donor cells and their subsequent placement into an egg in which the natural [haploid](#) nucleus has been eliminated. The donor cells can be derived from various tissues at various times in development. In *Xenopus*, the younger the nuclei, the better the survival rate. Nuclei isolated from blastula or gastrula stages sustain development beyond metamorphosis with higher efficiency than nuclei derived from the skin or intestine of a tadpole. Unlike what has been reported recently in mammals, however, no nuclei from adult frog tissue have been able to carry development beyond metamorphosis. The recipient cell is always an egg in which the haploid genome has been eliminated

by UV irradiation. The use of genetic markers in the donor and recipient cells, for example nuclei derived from albino embryos transferred into a wild-type pigmented egg, are used to assay for the successful elimination of the maternal nucleus. After nuclear transfer in this case, an albino frog derived from a pigmented egg represents a successful clone.

A modification of the nuclear transplant approach, in which nuclei from a *Xenopus* cell line are used for transplantation, has also been reported. This approach provides the advantage of manipulating the genome of the cultured cell prior to implantation. The number of embryos that develop normally after this procedure is low, and they usually do not survive beyond late neurula stages (41).

### 5.2.2. Transgenesis

Unlike transient expression of RNA or DNA by microinjection, [transgenic technology](#) allows stable and regulated **gene expression**. In *Xenopus*, unlike mouse, DNA injected into the egg or early blastomeres does not integrate into the genome, for reasons that are not clear. Using a technique originally developed in hydra, Kroll and Amaya reported recently the first successful transgenic frogs (42). This technique, called for *restriction-enzyme-mediated integration* (REMI), takes advantage of the nuclear transplantation technique in which nuclei are isolated from mature sperm, rather than somatic cells. The tightly packed sperm [chromatin](#) is decondensed *in vitro* and incubated with DNA that has been digested with [restriction enzymes](#). This DNA is incubated with the decondensed sperm nuclei in the presence of a small amount of the same restriction enzyme and [DNA Ligase](#). During this incubation, it is presumed that the restriction enzyme partially digests the sperm DNA; exogenous DNA that is mixed with the genome then integrates at the digested site(s) with the help of the ligase. Although several such transgenic frogs have been made, and this technique is becoming more widely used, no transgenic F1' progeny have been reported yet; specifically, there have been no reports of germline transmission of transgenic DNA. Most of the transgenic *Xenopus* generated to date do not successfully complete metamorphosis. This is, however, not necessarily a major problem; it is possible to perpetuate a line with nuclear transfer, so there is no requirement for sexual reproduction. Therefore, transgenic lines can, theoretically, be perpetuated indefinitely.

Progress in this field is very rapid, and there are many other attempts in progress. These include the use of [transposable elements](#) and [retroviruses](#) in various species of *Xenopus*. It can be safely predicted that generation of transgenic frogs will become routine before the end of this millennium.

### 5.2.3. Maternal Knockout Strategies

Although conventional **knockout** or knockin strategies, which are commonly used in the mouse, are not presently available in amphibians, it is possible to eliminate maternal transcripts from the egg selectively to assess their role during early embryogenesis. In this approach, antisense oligonucleotides against specific RNA sequences are microinjected into oocytes (43), which are then reimplanted into a female (44). These eggs are marked so that they can be distinguished from controls, either by using the albino/pigmented technique described above or by incubation with lipophilic vital dyes. Later, these eggs are isolated from the surrogate mother and fertilized *in vitro*. Their development is compared to controls that are treated in the exact same way but injected with control sense oligonucleotide. A stringent control for the specificity of the phenotype is provided by rescuing the egg after fertilization by injection of sense RNA.

## 6. Conclusions and Perspectives

In addition to the use of *Xenopus laevis* embryos and oocytes in molecular studies, there are many other features of *Xenopus* that make it amenable to molecular studies, in fields as diverse as metamorphosis, tissue regeneration, behavioral studies, physiology, pharmacology, and neurobiology. Excellent books and reviews are available describing these topics in length and listed under "Suggestions for further reading" (see also Refs. 45 and 46). Finally, many other amazing species of frogs, some now near extinction, have been virtually ignored.

Although there is appreciation of the conservation of molecular factors in the vertical evolutionary scheme, it is important not to undermine horizontal differences in strategies of development among different frog species. Frogs certainly provide a fertile terrain to explore these differences at the molecular level, and it would be no surprise if the differences in molecular strategies turn out to be more impressive than the mechanisms that are conserved.

## 7. Acknowledgments

I would like to thank Curtis Altmann for help in generating Figures [1](#) and [2](#), and Chenbei Chang for help with Figure [5](#). I am indebted to the members of my laboratory for constructive criticism of the manuscript, especially Paul Wilson, Chenbei Chang, Daniel Weinstein, and Bart Eggen. I would also like to thank Shauna Seliy for excellent administrative assistance and Jennifer Marden for technical help. Finally, I thank Margaret Bolce Brivanlou for her love and patience.

## Bibliography

1. V. Hamburger (1988) *The Heritage of Experimental Embryology: Hans Spemann and the Organizer*, Oxford Univ. Press, New York.
2. J. Gerhart et al. (1989) Cortical rotation of the *Xenopus* egg: Consequences for the anteroposterior pattern of embryonic dorsal development. *Development* **37**–51.
3. J. Newport and M. Kirschner (1982) A major developmental transition in early *Xenopus* embryos: 1. Characterization and timing of cellular changes at the midblastula transition. *Cell*, **30**, 675–686.
4. R. E. Keller (1986) "The cellular basis of amphibian gastrulation", in *Developmental Biology: A Comprehensive Synthesis*, Vol. **2**, The Cellular Basis of Morphogenesis, L. Browder, ed., Plenum Press, New York, pp. 241–327.
5. R. M. Stewart and J. C. Gerhart (1990) The anterior extent of dorsal development of the *Xenopus* embryonic axis depends on the quantity of organizer in the late blastula. *Development*, **109**(2), 363–372.
6. H. Spemann and H. Mangold (1924) Uber Induktion von Embryonanlagen durch Implantation artfremder Organisatoren. *Arch. mikr. Anat. EntwMech.* **100**, 599–638.
7. A. Chitnis and C. Kintner (1995) Neural induction and neurogenesis in Amphibian embryos. *Persp. Dev. Biol.* **3**, 3–15.
8. T. E. Schroeder (1970) Neurulation in *Xenopus laevis*: An analysis and model based upon light and electron microscopy. *J. Embryol. Exp. Morphol.* **23**, 427–462.
9. D. Brown et al. (1995) *Amphibian* Metamorphosis: A complex program of gene expression changes controlled by the thyroid hormone. *Recent Prog. Horm. Res.* **50**, 309–315.
10. L. Dale and J. M. W. Slack (1987) Fate map of the 32 cell stage of *Xenopus laevis*. *Development* **99**, 527–551.
11. R. Keller (1991) "Early embryonic development of *Xenopus laevis*", in *Methods in Cell Biology*, B. K. Kay and H. B. Peng, eds., Academic Press, San Diego.
12. G. W. Eagleson and W. A. Harris (1989) Mapping of the presumptive brain regions in the neural plate of *Xenopus laevis*. *J. Neurobiol.* **21**, 427–440.
13. J. M. W. Slack (1991) *From Egg to Embryo: Regional Specification in Early Development*, Cambridge Univ. Press, Cambridge, U.K.
14. J. Gerhart and M. Kirschner (1997) *Cells, Embryos, and Evolution*, Blackwell Science, Malden, MA.
15. J. C. Smith et al. (1991) Expression of a *Xenopus* homolog of brachyury (T) is an immediate–early response to mesoderm induction. *Cell* **67**(1), 79–87.
16. C. R. Kintner and D. A. Melton (1987) Expression of *Xenopus* N-CAM RNA in ectoderm is an early response to neural induction. *Development* **99**, 311–325.
17. P. D. Vize et al. (1991) Assays for gene function in developing *Xenopus* embryos. *Meth. Cell*

Biol. **36**, 367–387.

18. G. Lagna (1997) "Use of dominant negative constructs to modulate gene expression", in *Cellular Techniques in Developmental Biology*, F. de Pablo, A. Ferrus, and C. Stern., eds., Academic Press, pp. 75–96.
19. C. V. E. Wright et al. (1989) Interference with the function of a homeobox gene in *Xenopus* embryos produces malformations of the anterior spinal cord. **59**, 81–93.
20. A. Glinka et al. (1998) Dickkopf-1 is a member of a new family of secreted proteins and functions in head induction. *Nature* **391**, 357–362.
21. J. C. Smith (1987) A mesoderm-inducing factor is produced by a *Xenopus* cell line. *Development* **99**, 3–14.
22. M. Whitman and D. A. Melton (1989) Induction of mesoderm by a viral oncogene in early *Xenopus* embryos. *Science* **244**, 803–806.
23. H. Grunz and L. Tacke (1989) Neural differentiation of *Xenopus laevis* ectoderm takes place after disaggregation and delayed reaggregation without inducer. *Cell Diff. Dev.* **28**, 211–218.
24. P. Wilson and A. Hemmati Brivanlou (1997) Vertebrate neural induction: Inducers, inhibitors, and a new synthesis. *Neuron* **18**, 1–20.
25. P. S. Klein and D. A. Melton (1994) Hormonal regulation of embryogenesis: The formation of mesoderm in *Xenopus laevis*. *Endocrine Rev.* **15**, 326–341.
26. P. D. Nieuwkoop (1969) The formation of mesoderm in urodelean amphibians. I. Induction by the endoderm. *Wilhelm Roux Arch. EntwMech. Org.* **162**, 341–373.
27. W. G. Cox and A. Hemmati-Brivanlou (1995) Caudalization of neural fate by tissue recombination and bFGF. *Development* **121**, 4349–4358.
28. R. E. Keller and M. Danilchik (1988) Regional expression, pattern and timing of convergence and extension during gastrulation of *Xenopus laevis*. *Development* **103**, 193–210.
29. T. Doniach, C. R. Phillips, and J. C. Gerhart (1992) Planar induction of anteroposterior pattern in the developing central nervous system of *Xenopus laevis*. *Science*, **257**, 542–545.
30. A. Hemmati-Brivanlou, R. M. Stewart, and R. M. Harland, Region-specific neural induction of an engrailed protein by anterior notochord in *Xenopus*. *Science* **250** (4982), 800–802.
31. C. R. Kintner and J. Dodd (1991) Hensen's node induces neural tissue in *Xenopus* ectoderm. Implications for the action of the organizer in neural induction. *Development* **113**, 1495–1506.
32. M. Blum et al. (1992) Gastrulation in the mouse: The role of the homeobox gene gooseoid. *Cell* **69**, 1097–1106.
33. W. B. Smith and R. M. Harland (1992) Expression cloning of noggin, a new dorsalizing factor localized to the Spemann organizer in *Xenopus* embryos. *Cell* **70**, 829–840.
34. K. Lustig and M. Kirschner (1995) Use of an oocyte expression assay to reconstitute inductive signaling. *Proc. Natl. Acad. Sci. USA* **92**, 6234–6238.
35. K. R. Kao and R. P. Elinson (1988) The entire mesodermal mantle behaves as Spemann's organizer in dorsoanterior enhanced *Xenopus laevis* embryos. *Dev. Biol.* **127**, 64–77.
36. J. Holtfreter (1933) Die totale Exogastrulation, eine Selbstablosung des Ektoderms von Entomesoderm. *Entwicklung und funktionelles Verhalten nervenloser Organe. Arch. EntwMech. Org.* **129**, 669–793.
37. J.-D. Graf and H. R. Kobel (1991) "Genetics of *Xenopus laevis*", in *Methods in Cell Biology*, B. K. Kay and H. B. Peng, eds., Academic Press, San Diego, pp. 19–34.
38. B. Roe et al. (1985) The complete nucleotide sequence of the *Xenopus laevis* mitochondrial genome. *J. Biol. Chem.* **260**, 9759–9774.
39. A. Droin (1991) Mutants of *Xenopus laevis*. *Meth. Cell Biol.* **36**, 671–673.
40. J. Gurdon, T. R. Elsdale, and M. Fishberg (1958) *Nature* **182**, 64–65.
41. K. Kroll and J. Gerhart (1994) Transgenic *X. laevis* embryos from eggs transplanted with nuclei of transfected cultured cells. *Science* **266**, 650–653.



42. K. Kroll and E. Amaya (1996) Transgenic *Xenopus* embryos from sperm nuclear transplantations reveal FGF signaling requirements during gastrulation. *Development* **122**, 3173–3183.
43. R. M. Harland and H. Weintraub (1985) Translation of mRNA injected into *Xenopus* oocytes is specifically inhibited by antisense RNA. *J. Cell Biol.* **101**, 1094–1099.
44. J. Heasman, S. Holwill, and C. C. Wylie (1991) "Fertilization of cultured oocytes and use in studies of maternally inherited molecules", in *Xenopus laevis: Practical Uses in Cell and Molecular Biology*, Academic Press, San Diego. pp. 213–230.
45. P. D. Nieuwkoop and J. Faber (1967) *Normal Table of Xenopus laevis (Daudin)*, 2nd ed., North-Holland Publishing Company, Amsterdam.
46. K. R. Kao and R. P. Elinson (1989) Dorsalization of mesoderm induction by lithium. *Dev. Biol.* **132**(1), 81–90.

### Suggestions for Further Reading

47. J. Gerhart and M. Kirschner (1997) *Cells, Embryos, and Evolution*, Blackwell Science, Malden, MA.
48. V. Hamburger (1988) *The Heritage of Experimental Embryology, Hans Spemann and the Organizer*, Oxford Univ. Press, New York.
49. P. Hausen, Peter M. Riebesell (1991) *The Early Development of Xenopus Laevis, An Atlas of the Histology*, Springer-Verlag, Berlin.
50. B. K. Kay and B. H. Peng, eds. (1991) "*Xenopus laevis*: Practical uses in cell and molecular biology". *Methods in Cell Biology*, Vol. **36**, Academic Press, San Diego.
51. P. D. Nieuwkoop and J. Faber, eds. (1994) *Normal Table of Venopus Laevis (Daudin)*, Garland Publishing, New York.
52. F. de Pablo, A. Ferrus, and C. D. Stern, eds. (1998) *Cellular and Molecular Procedures in Developmental Biology*, Academic Press, San Diego.
53. R. C. Tinsley and H. R. Kobel, eds. (1996) *The Biology of Xenopus*, Univ. Press New York.

## FX174 Phage

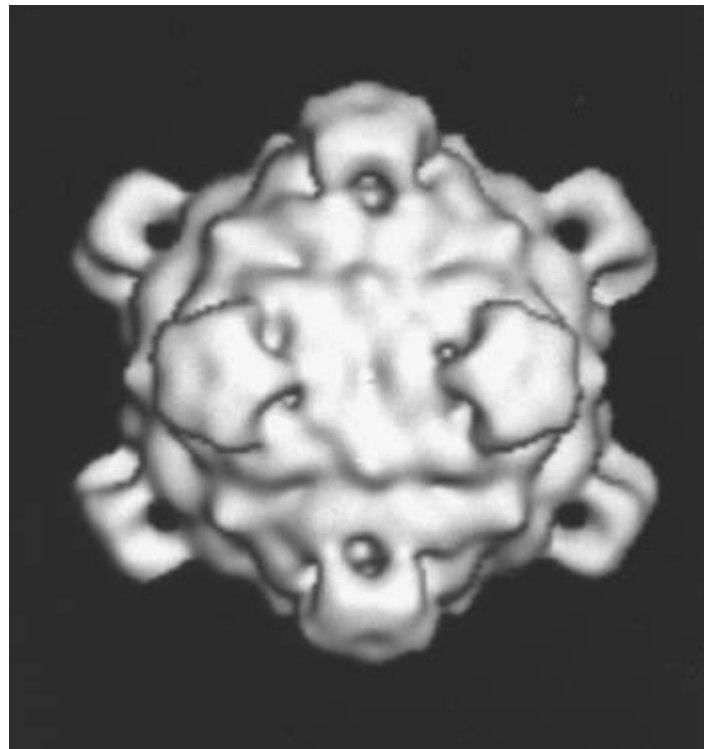
Prior to the development of the **electron microscope**, the size of an object, such as a **virus**, was estimated by filtration through membranes of various pore sizes. By this criterion, the bacterial virus (bacteriophage) fX174, first described in 1935 (1), and its close relative S13 (2) were classified as much smaller than the **T phages** that were under investigation in the early days of molecular biology. Moreover, the fact that their genetic material is **DNA** made these two viruses attractive candidates as front runners in the race to unravel the molecular basis of heredity. The contributions began with the initial characterization of the fX174 DNA as being single stranded, as reported in the first issue of the *Journal of Molecular Biology* (3). The culmination came when fX174 was the first organism to have its **genome** subjected completely to **DNA sequencing** (4), which offered the first direct documentation of the existence of overlapping **genes** in nature.

### 1. Virus structure

The fX174 virion appeared polyhedral with “knobs” or “spikes” in early electron micrographs of metal-shadowed and stained preparations (4). These images led to the proposal that the virion has

**icosahedral symmetry**, a structural feature that was later confirmed at a higher resolution by [cryoelectron microscopy](#) (5) and at atomic resolution by [X-ray crystallography](#) (6). The major contribution of these structural studies was, however, the establishment of the three-dimensional [protein structures](#) of the three major proteins in the crystalline form of the fX174 virion and their [quaternary structure](#) in both the crystalline and aqueous states of the virion. In the virion, the fX174 circular, single-stranded, DNA genome is enclosed within a  $T = 1$  icosahedral protein shell, the simplest of the possible deltahedrons (7). Consequently, the protein shell of the virion contains 60 copies of the proteins encoded in genes F, G, and J, and 12 copies of the gene H protein product (Table 1). Five copies of F capsid protein are clustered around each of the 12 vertices of the icosahedron, forming a shell that has no openings at the center of the 20 triangular faces, as seen in the reconstructed image of the virion exterior in Figure 1. Five copies of G protein form the spikes, or projections, that extend outward from the capsid surface at each vertex of the icosahedron (Fig. 1). The spikes are attached to the capsid outer surface by each G protein interacting with its F protein partner via eight direct [hydrogen bonds](#) and five indirect, water-mediated [hydrogen bonds](#) (8).

**Figure 1.** Three-dimensional reconstructed image of fX174 virion from electron micrographs of unstained frozen-hydrated preparations (courtesy of Norman H. Olson and Timothy S. Baker, Purdue University). The diameter of the virion is 33.5 nm between the exterior edges of the spikes at the fivefold axes (28).



**Table 1. fX174 Structural Proteins**

| Protein   | Mol. Wt.           | No. of Residues  | No. in Virion | No. in Procapsid |
|-----------|--------------------|------------------|---------------|------------------|
| B protein | $13.8 \times 10^3$ | 120 <sup>a</sup> | 0             | 60               |
| D protein | $16.9 \times 10^3$ | 152              | 0             | 240              |
| F protein | $48.4 \times 10^3$ | 426              | 60            | 60               |

|           |                    |                  |    |    |
|-----------|--------------------|------------------|----|----|
| G protein | $19.0 \times 10^3$ | 175              | 60 | 60 |
| H protein | $34.4 \times 10^3$ | 328 <sup>a</sup> | 12 | 12 |
| J protein | $4.2 \times 10^3$  | 36               | 60 | 0  |

<sup>a</sup> Based on DNA sequence (4).

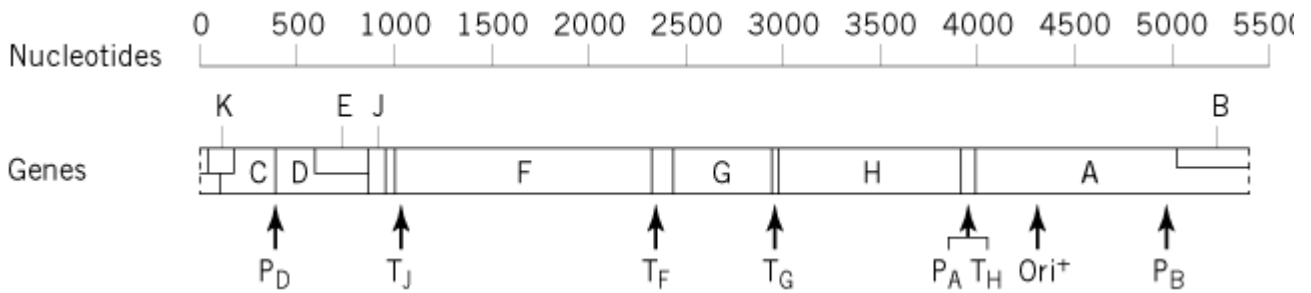
The [single-particle reconstruction](#) (5) and crystallographic data (6) also revealed structural details about the interiors of the capsid layer and the spikes, as well as additional structural components that are interacting with the interior surface of the capsid. For example, both F and G polypeptides contain an eight-stranded antiparallel [beta-sheet](#) motif ([b-barrel](#)) like that found in most viral capsid proteins. In addition, each F pentameric cluster contains a small, central channel that leads into a central channel within each G pentameric cluster. This continuous channel has mainly **hydrophobic** walls and varies between 0.6 and 2.6 nm in diameter (6). The unassigned electron density within the G channel near the junction with the F channel is thought to be the location of a short segment of the H protein that presumably lies within the channel. The remainder of the H polypeptide chain is believed to be located beneath the interior opening of the channel, in contact with both the DNA and the interior surface of the F capsid proteins. The main evidence for this location for each of the 12 H proteins (9) is the fact that [urea](#) treatment of virions removes both the G spikes and the H proteins, leaving a spherical DNA-containing particle (10). Although the H polypeptide chains could not be unambiguously located in the X-ray crystallography electron density map, all 36 residues of the basic J protein could be assigned to an S-shaped path on the inner surface of the capsid, linking two neighboring F proteins across each two-fold axis of symmetry. In addition, a four-nucleotide segment of single-stranded DNA was located in the [difference Fourier](#) map calculated from the differences in the diffraction data on 114 S virions and on 70 S particles that contain less DNA (1). These nucleotides are interacting with four amino-acid side chains of the C-terminal segment of the J protein and with the side chains of seven amino acid residues in the F protein.

## 2. Genome structure

Mutational analysis has identified 11 genes in the circular genome of fX174 (Fig. 2). The A segment of the DNA sequence codes for two gene products, endonuclease A and a truncated version of A that suppresses host DNA synthesis, A\* (not shown in Fig. 2). The A\* protein is **translated** in the same [reading frame](#) as A, but its translation start site is 173 **codons** downstream of the A translation start site. Furthermore, a portion of the DNA sequence in the A segment is read in a different reading frame as gene B, one of the scaffolding proteins required for the assembly of empty procapsids (Table 1). In fact, all three reading frames are used in the DNA sequence that encodes gene K, an unidentified protein that has an effect on progeny virion yield. The protein encoded by gene C is involved in switching DNA synthesis from replication of the double-stranded fX174 replicative form (RF) DNA to the simultaneous synthesis and packaging of fX174 single-stranded progeny DNA into empty procapsids. These packaged DNA strands have the same sequence as fX174 [messenger RNAs](#) that are **transcribed** from the fX RF DNA. Therefore, the packaged strands are labeled positive (+), while the complementary, unpackaged strands are labeled negative (-). The next region of the genome codes for a second set of overlapping genes read in two different reading frames. D is a second scaffolding protein found in procapsids, but not in virions (Table 1), and E induces cell lysis at the bacterial site for elongation and division. Genes J, F, G, and H code for structural proteins found in both virions and procapsids (Table 1).

**Figure 2.** Linear version of the fX174 circular genetic map. The locations of promoters (P), transcription terminations

sites (T), and the origin of replication (*ori*<sup>+</sup>) are shown.



The fX174 genome contains three types of DNA sequences, in addition to those that code for proteins and those that are involved in transcription and translation (**promoters**, terminators, and **ribosome** binding sites; Fig. 2). The first is an **origin of replication** (*ori*<sup>+</sup>) located within gene A, where endonuclease A cleaves between nucleotides 4305 and 4306 (11). The second is an “incompatibility” or “reduction” sequence (12) that interferes with fX174 replication if the sequence is present on a **plasmid** in the uninfected cell. The 85-nucleotide sequence includes parts of genes H and A and the intergenic sequence between them. One suggestion for the function of this sequence is the binding of the incoming single-stranded fX DNA to a site on the bacterial membrane that is essential for fX DNA replication. Presumably, these sites are present in limiting amounts in the uninfected cell. Finally, the third type of DNA sequence includes the intergenic sequences. The H/A untranslated sequence found within the incompatibility sequence shares one feature with two of the three remaining intergenic sequences (J/F, F/G, and G/H): namely, the potential for forming complementary, double-stranded hairpin structures in the single-stranded states of either the DNA or mRNA. However, these potential nucleic acid secondary structures appear to play only subtle roles in fX174 infection, since disruption of the J/F potential hairpin increased the decay rate of only two species of fX mRNA (13).

### 3. Infection cycle

As with most viruses, the fX174 infection cycle begins with the attachment of the virion to receptor molecules on the outer surface of the bacterial host. The core oligosaccharide region of the lipopolysaccharide (LPS) in the outer membrane of *Escherichia coli* C constitutes the initial sites where the fX174 virion binds via several of its spikes (14). Since O antigenic side chains attached to the core monosaccharide units block fX attachment, most **gram-negative bacterial** species are resistant to fX174. Binding is mediated by Ca<sup>+2</sup> cations (15), and there are three separate sites for the cation in the crystal structure of the virion, two in F protein near the threefold axes and one near the outer opening of the channel through the G pentamers in each spike (16). Of interest is the presence of a glucose molecule just below one of the Ca<sup>+2</sup> sites in F. Moreover, another host gene, which lies between *gal* and *aroG* at 17 min in the *E. coli* map, is required for fX binding (17).

Due to the single-strandedness of its DNA genome, fX174 has an additional functional requirement that determines which bacterial species it can infect. The virion must penetrate the outer wall and cytoplasmic membrane to deliver the fX genome to the bacterial DNA replicative complex for synthesis of the complementary (–) strand. This transfer of the viral genome marks the beginning of the eclipse phase, the time period when infectious virions are not present in the infected cell. After binding to the LPS receptor, Ca<sup>+2</sup>-induced structural changes in the fX virion proteins begin the ejection of the single-stranded, circular DNA from within the virion, presumably through an opening in its protein shell (18-20). Most evidence supports a model for ejection through the channel down the center of one “spike” (18, 21). However, a conformational change that would produce an opening where six F proteins meet at the threefold axis cannot be ruled out. Furthermore, it is not clear how

the fX virion encounters the bacterial components involved in DNA replication. Some experiments suggest that these host factors are located where the outer wall adheres to the cytoplasmic membrane (22, 23).

As the ejected fX174 genome enters the cell, the single-stranded DNA is converted to double-stranded RF (Stage I). The host **single-stranded binding protein** (SSB) coats the incoming DNA, and a protein complex (the preprimosome) is assembled on the coated DNA. The preprimosome contains the sequence recognition protein (PriA), PriB, PriC, dnaT, and gyrase (dnaB) (24), and synthesis of several RNA primers begins when **primase** (dnaG) joins the complex. Then **DNA polymerase III** adds deoxyribonucleotides to the primers; DNA polymerase I replaces the RNA primers with DNA; and **DNA Ligase** joins the two ends of the – strand to produce a **supercoiled RF** molecule.

The bacterial transcription system begins to transcribe the fX174 genome, using the RF molecule as the template. Since only two fX genes (A and A<sup>\*</sup>) code for **enzymes**, and the remaining fX genes code for proteins that function as structural components, the virus does not require a sophisticated mechanism for regulating transcription. Except for genes E and K, the fX genes are transcribed and translated at all times during the viral replication cycle. Moreover, the relative amounts of the fX proteins are determined by the relative amounts of the mRNAs that contain the message for each protein (25), Table 2. Thus, the message for endonuclease A (and A<sup>\*</sup>) is present only in the unstable mRNA and the two least abundant transcripts (8 and 9 in Table 2), while the message for the most abundant fX protein, D scaffold (see Table 1), is found in all transcripts (Table 2). Also, the **rho**-independent termination is inefficient, producing two groups of transcripts that begin at the same promoter but end at different termination sites (1, 3, 7 and 2, 4, 5, 6 in Table 2).

**Table 2. fX174 mRNA Transcripts(25)**

| Genes Transcribed   | Promoter                    | Termination    | Order of Abundance |
|---------------------|-----------------------------|----------------|--------------------|
| D, J                | P <sub>D</sub>              | T <sub>J</sub> | 1                  |
| B, C, D, J          | P <sub>B</sub>              | T <sub>J</sub> | 2                  |
| D, J, F             | P <sub>D</sub>              | T <sub>F</sub> | 3                  |
| B, C, D, J, F       | P <sub>B</sub>              | T <sub>F</sub> | 4                  |
| B, C, D, J, F, G    | P <sub>B</sub>              | T <sub>G</sub> | 5                  |
| B, C, D, J, F, G, H | P <sub>B</sub>              | T <sub>H</sub> | 6, 7               |
| D, J, F, G, H       | P <sub>D</sub>              | T <sub>H</sub> | 6, 7               |
| F, G, H             | P <sub>F</sub> <sup>a</sup> | T <sub>J</sub> | 8                  |
| A, B, C, D, J       |                             |                |                    |
| D, J, F, G, H       | P <sub>D</sub>              | T <sub>J</sub> | 9 <sup>a</sup>     |
| A, B, C, D, J       |                             |                |                    |
| A,?                 | P <sub>A</sub>              | ?              | Unstable           |

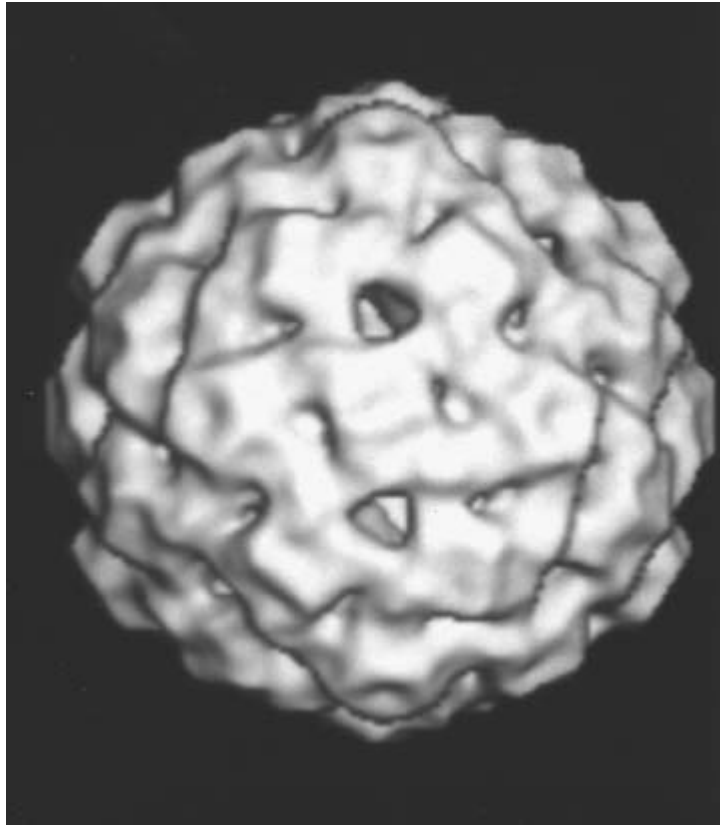
<sup>a</sup> Rare promoter

<sup>b</sup> Single copy

Endonuclease A is the first fX protein synthesized, to create a pool of RF molecules as templates for both transcription and fX genomic DNA replication. This Stage II RF DNA replication begins with endonuclease A cleaving the + strand within the A gene sequence and forming a covalent bond with the 5' phosphate of the cleaved strand. The bacterial [DNA helicase](#) (Rep) unwinds the double-stranded RF, and DNA polymerase III adds deoxyribonucleotides to the 3' OH of the cleaved + strand, until the new + strand is completed. At the same time, SSB coats the newly synthesized single-stranded fX DNA (26). As soon as the primosome assembly site in the F/G intergenic region (Fig. 2) of the new + strand becomes available, the preprimosome assembles on the new + strand and begins synthesis of a new – strand (see discussion of Stage I). To ensure continuous synthesis of RF, endonuclease A cleaves the + strand at the junction between the old and new strands and joins the ends of the old strand as soon as the A gene sequence on the new + strand becomes available. At the same time, A protein forms a covalent bond with the free 5' phosphate group of the new + strand, to begin synthesis of the next RF molecule. Each new RF molecule is released from the replication complex when DNA ligase joins the ends of the completed – strand.

As the pool of RF molecules increases, the level of fX transcripts reaches a point where sufficient amounts of fX capsid (F), spike (G, H), and scaffold (B, D) proteins are available to begin the assembly of 108 S empty procapsids (27). First, 9S complexes are formed of F and 6S complexes of G (both containing five copies of the corresponding fX proteins). These two complexes are joined with H and assembled into a procapsid, with B crosslinking F pentamers on their interior surface across the twofold axes and D forming an exterior scaffold on the surface of the capsid (28, 29). In Figure 3, the pentameric G spikes at each fivefold axis are not as prominent as in the virion (Fig. 1) due to interactions with D on the capsid outer surface. These G/D interactions keep each G b-barrel attached to the capsid surface, but inclined at a 20° angle to the fivefold axis. In the virion, the G b-barrels are roughly parallel to the fivefold axis and held to the capsid surface with the aid of bound water molecules. In the procapsid, the mainly **alpha-helical** D scaffold protein is present in four radically different environments. However, [protein–protein interactions](#) between D molecules form a scaffold that holds each F and G pentameric complex in a shell that has almost three times the internal volume as the protein shell of the mature fX virion. Furthermore, the enlarged procapsid shell has holes with a diameter of 3.0 nm at the threefold axes (see Fig. 3).

**Figure 3.** Three-dimensional reconstructed image of fX174 procapsid from electron micrographs of unstained frozen-hydrated preparations (courtesy of Norman H. Olson and Timothy S. Baker, Purdue University). The diameter of the procapsid is 35.5 nm between the exterior edges of the spikes at the fivefold axes (28).



As procapsids are being assembled, fX C protein begins to compete with SSB for binding to RF at the fX A/helicase complex on the replicating fX DNA. Once bound, C stops Stage II RF replication and provides a binding site for procapsid, presumably at one of the 3.0-nm holes. fX + strand is then simultaneously synthesized and packaged, as 60 molecules of the basic J polypeptide are bound to the single-stranded DNA ([30](#), [31](#)). Each J protein displaces a B polypeptide on the interior surface of two neighboring F proteins as the single-stranded DNA is packaged. The 60 B chains presumably exit from the procapsid through the 19 other holes, forming the 132S provirion.

Meanwhile, small amounts of fX lysis (E) protein are slowly being translated from mRNA containing both D and E messages. The low level of E translation is probably due to a weak ribosome binding site, to rare codons in the E reading frame, and to competition from ribosomes translating D messages. fX E polypeptides begin forming a tunnel through the cell wall at the site of cell elongation and division ([32](#)). Once the cytoplasmic membrane extrudes through the tunnel, lysis occurs, and the fX provirions are exposed to the ionic environment of the medium. The D scaffold proteins dissociate from the provirions, allowing conformational changes in F capsid proteins and minor changes in inclination of G proteins in the spikes, as the protein shell collapses around the fX DNA. The newly formed progeny fX virions then begin the infection cycle when they encounter uninfected cells.

#### 4. Related bacteriophage

Although most of the experimental data have been obtained with fX174, the Seventh Report by the International Committee on the Taxonomy of Viruses includes 38 other phages in the *Microviridae* family. The majority of the members belong to the same genus (*Microvirus*) as fX. The natural hosts for the members of the *Microvirus* genus are various species of *Enterobacteriaceae*. There is, however, at least one member in three additional genera of the *Microviridae* family that infect other eubacterial species. These include members of gram-negative genera *Bdellovibria* and *Chlamydia* and the wall-less genus *Spiroplasma*, which causes several plant diseases.

## Bibliography

1. V. Sertic and N. Bulgakov (1935) *C. R. Soc. Biol. Paris* **119**, 1270–1272.
2. D. E. Lea and M. H. Salaman (1946) *Proc. Roy. Soc. B* **133**, 434–443.
3. R. L. Sinsheimer (1959) *J. Mol. Biol.* **1**, 37–42.
4. C. E. Hall, E. C. Maclean, and I. Tessman (1959) *J. Mol. Biol.* **1**, 192–194.
5. N. H. Olson, T. S. Baker, P. Willingham, and N. L. Incardona (1992) *J. Struct. Biol.* **108**, 168–175.
6. R. McKenna, D. Xia, P. Willingham, L. L. Ilag, S. Krishnaswamy, M. G. Rossmann, N. H. Olson, T. S. Baker, and N. L. Incardona (1992) *Nature* **355**, 137–143.
7. D. L. D. Caspar and A. Klug (1962) *Cold Spring Harbor Symp. Quant. Biol.* **27**, 1–24.
8. R. McKenna, L. L. Ilag, and M. G. Rossmann (1994) *J. Mol. Biol.* **237**, 517–543.
9. A. B. Burgess (1969) *Proc. Natl. Acad. Sci. USA* **64**, 613–617.
10. M. H. Edgell, C. A. Hutchison, III, and R. L. Sinsheimer (1969) *J. Mol. Biol.* **42**, 547–557.
11. S. A. Langeveld, A. D. M. Van Manfeld, P. D. Bass, H. S. Janz, G. A. Van Arkel, and P. J. Weisbeck (1978) *Nature* **271**, 417–420.
12. H. G. Van der Avoort, S. A. Langeveld, A. D. M., G. A. Van Arkel, and P. J. Weisbeck (1984) *J. Virol.* **50**, 533–540.
13. M. N. Hayashi, M. Hayashi, and U. R. Mueller (1983) *J. Virol.* **48**, 186–196.
14. N. L. Incardona and L. Selvidge (1973) *J. Virol.* **11**, 775–782.
15. R. Fujimura and P. Kaesberg (1962) *Biophys. J.* **2**, 433–449.
16. L. L. Ilag, R. McKenna, M. P. Yadav, J. N. BeMiller, N. L. Incardona, and M. G. Rossmann (1994) *J. Mol. Biol.* **244**, 291–300.
17. R. Munikyo, T. Tsuzuki, and M. Sekiguchi (1979) *J. Bacteriol.* **138**, 1038–1040.
18. L. L. Ilag, J. K. Tuech, L. A. Beisner, R. A. Sumrada, and N. L. Incardona (1993) *J. Mol. Biol.* **229**, 671–684.
19. J. E. Newbold and R. L. Sinsheimer (1970) *J. Mol. Biol.* **49**, 49–66.
20. Y. Mano, T. Kawabe, T. Komano, and K. Yazaki (1982) *Agric. Biol. Chem.* **46**, 2041–2049.
21. S. M. Jaswinski, R. Marco, and A. Kornberg (1975) *Virology* **66**, 294–305.
22. M. E. Bayer and T. W. Starkey (1972) *Virology* **49**, 236–256.
23. J. Azuma, J. Morita, and T. Komano (1980) *J. Biochem.* **88**, 525–532.
24. K. Arai, R. Low, J. Kobori, J. Schlomai, and A. Kornberg (1981) *J. Biol. Chem.* **256**, 5273–5280.
25. M. Hayashi, F. K. Fujimura, and M. N. Hayashi (1976) *Proc. Natl. Acad. Sci. USA* **73**, 3519–3523.
26. J. Ikeda, A. Yudelevich, and A. Hurwitz (1976) *Proc. Natl. Acad. Sci. USA* **73**, 2669–2673.
27. R. Mukai, R. K. Hamatake, and M. Hayashi (1979) *Proc. Natl. Acad. Sci. USA* **76**, 4877–4881.
28. L. L. Ilag, N. H. Olson, T. Dokland, C. L. Music, R. H. Cheng, Z. Bowen, R. McKenna, M. G. Rossmann, T. S. Baker, and N. L. Incardona (1995) *Structure* **3**, 353–363.
29. T. Dokland, R. McKenna, L. L. Ilag, B. R. Bowman, N. L. Incardona, B. A. Fane, and M. R. Rossmann (1997) *Nature* **389**, 308–313.
30. A. Aoyama, R. K. Hamatake, and M. Hayashi (1981) *Proc. Natl. Acad. Sci. USA* **78**, 7285–7389.
31. B. Jennings and B. A. Fane (1997) *Virology* **227**, 370–377.
32. W. D. Roof and R. Young (1995) *Fed. Exp. Micro. Soc. Microbiol. Rev.* **17**, 213–219.

## Suggestions for Further Reading



33. D. T. Denhardt, D. Dressler, and D. S. Ray, eds. (1978) *The Single-Stranded DNA Phages*, Cold Spring Harbor Laboratory, Cold Spring Harbor, New York (covers all single-stranded DNA phages).
34. M. Hayashi, A. Aoyama, D. L. Richardson Jr., and M. N. Hayashi, (1988) "Biology of bacteriophage X174. In" *The Bacteriophages*, Vol. 2, (R. Calendar, ed.), Plenum Press, New York and London, pp. 1–71 (contains an exhaustive bibliography).
35. P. D. Bass and H. S. Jansz (1988) Single-stranded DNA phage origins, *Current Topics in Microbiology and Immunology* **136**, 31–37.
36. A. Kornberg and T. Baker (1991) *DNA Replication*, 2nd ed., W. H. Freeman and Company, New York (contains detailed information on X174 DNA replication).

## Y-Chromosome

The Y-chromosome is the sex [chromosome](#) that determines male sexual identity in humans. The Y-chromosome in humans is small, yet contains a pseudoautosomal region in which there are significant homologies to the [X-chromosome](#). This region facilitates chromosome pairing at **meiosis**. Active genes in this region include those for steroid sulfatase and the Xg blood group. The Y-chromosome also contains the sex-determining region Y-gene (SRY gene). The SRY gene is expressed only in the male genital ridge during differentiation of the gonad. The SRY gene product encodes a [DNA-binding protein](#) that is a transcriptional regulator.

It has been proposed that X and Y-chromosomes evolved by **translocating** transcribed regions from the Y-chromosome to the X-chromosome. Then the Y-chromosome would be reduced in size, whereas the X-chromosome would progressively acquire the capacity to inactivate the domains of the chromosome that could not pair in meiosis. The fact that males have only one of each X- and Y-chromosome, whereas females have two X-chromosomes, means that males are more susceptible to genetic mutations on these chromosomes that lead to disease. These include, among very many, Duchenne muscular dystrophy, hemophilia, and fragile X-linked mental retardation.

## Bibliography

1. E. Jablonka and M. J. Lamb (1990) *Biol. Dev.* **65**, 249–276.
2. W. R. Rice (1987) *Genetics* **116**, 161–167.

## Z-DNA

The left-handed, double-helix [DNA structure](#) known as Z-DNA, due to the zigzag path of the phosphate groups along the DNA backbone, was unexpected. It was discovered through [X-ray crystallography](#) analysis of d(CGCGCG)<sub>2</sub> at 0.9 Å atomic resolution ([1](#)). The Z-DNA double helix is tall and thin, relative to B-DNA and A-DNA, with a diameter of ~18 and a helical pitch of ~45 (see Fig. 2 of [DNA Structure](#)). Z-DNA is favored by alternating Py-Pu sequences, especially (dC-dG)<sub>n</sub>.

The repeat unit is a CpG dinucleotide; the glycosyl conformation of cytosine is *anti* and that of guanine is *syn*. The sugar pucker for cytosine is C2'-*endo*, while that for guanine is predominantly C3'-*endo*. The base pairs are slightly inclined by  $-9^\circ$ , and there are 12 base pairs per helical turn. The rise per dinucleotide repeat is  $\sim 7.4$ , and the helical twist angle per dinucleotide repeat is  $-60^\circ$ . The average helical twist angles are about  $-8^\circ$  for the CpG and  $-52^\circ$  for the GpC with a sum of  $-60^\circ$  for the dinucleotide repeat unit (see Fig. 3 of [Base Pairs](#)). Several other oligonucleotides have been crystallized in the Z-DNA conformation, including fragments containing AT base pairs, nonalternating Pu-Py sequences, and guanine-thymine [wobble pairs](#), etc.

The single groove of Z-DNA is narrow and deep. When the **van der Waals** radii of the phosphate group are included, the opening available for access to the groove is 6–7 Å. The major groove becomes very shallow, in fact almost nonexistent, and lies exposed to the solvent. Spermine and hydrated metal ions interact with the phosphate and the N7 of guanine. The interaction of spermine and metal ions with the phosphate groups of the DNA plays a role in charge neutralization. Some spermine molecules that are found located in the deep groove bridge the minor groove by interacting with phosphate groups on both sides, and thus may play a role in stabilizing the Z-DNA helix.

DNA, especially that associated with alternating dC-dG sequences such as poly(dC-dG), can readily undergo the reversible conformational change between B-DNA and Z-DNA. The conformational equilibrium is influenced by negative **supercoiling**, changes in the ionic environment (eg, high salt), and the presence of proteins that bind to Z-DNA. Z-DNA has been demonstrated to exist in [polytene chromosomes](#) of *Drosophila* using fluorescent **antibodies**. The left-handedness of Z-DNA is particularly powerful in relieving the strain of negatively supercoiled DNA (eg, during gene [transcription](#)). More recently, binding proteins specific for Z-DNA have been discovered. It has been established that Z-DNA occurs in metabolically active nuclei and that the level is dependent on DNA torsional strain and increases with transcription. Understanding the definitive biological role of Z-DNA is being actively pursued.

#### Bibliography

1. A. H.-J. Wang et al. (1979) Nature **282**, 680–686.

## Zebrafish

Zebrafish, *Danio (Brachydanio) rerio* (Hamilton-Buchanan), is a member of the family *Cyprinidae*, order *Cypriniformes*. A native of rivers in Southeastern Asia, this small, tropical freshwater fish is now a popular inhabitant of home aquaria throughout the world. A number of features facilitating **embryological** and **genetic** manipulations has made zebrafish the most recent model system in which to study mechanisms of vertebrate embryonic [development](#) ([1](#), [2](#)). Zebrafish are also frequently used in toxicologic studies ([3](#)).

Zebrafish are elongated and flattened laterally, with dark blue and silver-gold stripes running along the longitudinal axis of their body. In their two- to three-year life span, zebrafish grow up to 5 cm in length. Zebrafish achieve sexual maturity at three months of age. They are omnivorous and thus easy to raise using a variety of live and dry foods. Furthermore, they tolerate a relatively wide range of water quality. In captive breeding, zebrafish can produce hundreds of eggs on a weekly basis without seasonal variation. Therefore, large numbers ( $10^5$ ) of zebrafish can be raised and maintained at low cost and labor, facilitating large-scale genetic screens. The mating behavior is photoperiodic; spawning takes place at dawn or, in laboratory conditions, at the beginning of the light cycle. In the

course of mating, the male closely follows the fast-swimming female, culminating in the release of eggs and sperm into the water. Eggs and sperm can be harvested easily from anesthetized fishes and fertilized in vitro. Importantly, the sperm can be preserved by deep freezing, allowing for cost and space-effective maintenance of genotypes (4).

Embryos from a single clutch develop synchronously outside the mother in a range of temperatures, (between 23°C and 33°C, with the optimal temperature being 28.5°C). The embryo and surrounding chorion are translucent, allowing for a detailed microscopic inspection of development. Furthermore, the embryos can survive several days with major developmental defects (eg, without a functional cardiovascular system). The large size of the embryo (0.7 mm in diameter) facilitates [microinjection](#), dissection, transplantation, and other embryological manipulations. All of the above features are critical for effective identification and subsequent analysis of mutant phenotypes.

## 1. Genomics

The haploid genome of zebrafish contains about  $1.7 \times 10^9$ bp (5), organized into 25 approximately metacentric [chromosomes](#) of similar size (6). A genetic linkage map of the zebrafish genome has been developed using PCR (polymerase chain reaction)-based polymorphisms: random amplified polymorphic markers (RAPID) (7) as well as CA dinucleotide repeats and simple sequence length polymorphisms (SSLP) (8). Furthermore, about 500 DNA sequences have been reported for zebrafish, and positions on the genetic map have been determined for 120 of these cloned genes. The current estimate of the size of the map is 3,000 centiMorgan (cM) with an average  $600 \text{ kb cM}^{-1}$ . Taking into account about 1,200 markers for which map positions are known, the average interval between the available zebrafish markers is approximately 1,500 kb. The mapping efforts also revealed a surprising conservation of many large chromosome segments in the genomes of humans and zebrafish (9).

## 2. Embryonic Development

Embryonic development is initiated by the segregation of noncytolytic cytoplasm in the form of a blastodisc situated atop a large sphere of translucent yolk. A series of synchronous and rapid cleavages of the blastodisc leads to the formation of a **blastula** with a mound of [blastomeres](#) on a syncytial yolk cell. The establishment of dorso-ventral asymmetry of the embryo is not strictly correlated with the planes of the first cleavages (10, 11). This process, however, requires the transport of substances from the vegetal hemisphere into marginal blastomeres via cortical [microtubules](#) of the yolk cell (12). An important step in the establishment of the dorsal axis is a transient accumulation of b-catenin in the nuclei on the dorsal side of the zebrafish blastula (13, 14).

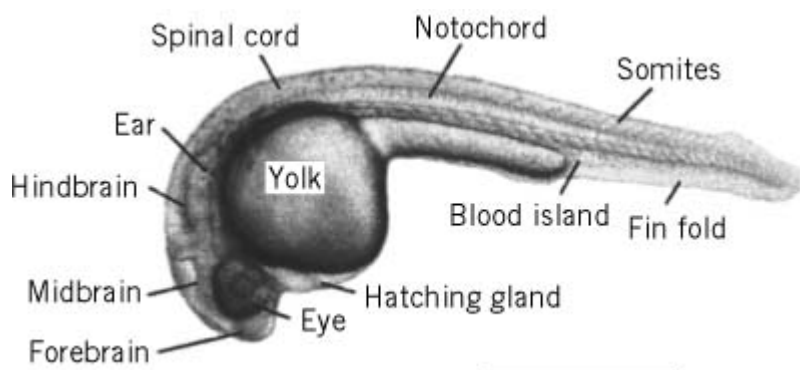
The embryonic body is shaped by three main morphogenetic movements of **gastrulation** (15). During epiboly, the blastoderm becomes thinner as its surface expands to cover the yolk cell completely 10 hours after fertilization. Tissue-restricted lineages can be identified in late blastula, when the blastoderm reaches the equator of the yolk cell. The three classically defined germ layers exhibit an overall organization reminiscent of the early gastrula fate maps of other vertebrates, with progenitors of **ectoderm** arising from near the **animal pole**, and the **mesoderm** and **endoderm** occupying the most marginal positions of the blastoderm (16). The germ layers form by involution/ingression movements at the blastoderm margin that bring prospective mesodermal and endodermal cells underneath the future ectoderm. Concurrent convergent extension movements bring the precursors of various tissues closer to the dorsal midline and are responsible for narrowing of the embryonic axis in the mediolateral direction and for its anterior-posterior elongation. These gastrulation movements create the embryonic shield, a thickening at the dorsal blastoderm margin. The embryonic shield is equivalent to the Spemann gastrula organizer of frog and other vertebrate embryos (17, 18). The dorso-ventral pattern specification involves highly conserved antagonistic interactions between dorsalizing signals emanating from the Spemann gastrula organizer and secreted factors like bone morphogenetic protein-4 (BMP-4) released by the ventral signaling center

(19, 20). Gastrulation movements bring precursors of various organs into their proper positions in the embryo.

During the subsequent segmentation period, organ rudiments form; notochord becomes visible as a prominent dorsal structure and blocks of somites appear on both sides of the notochord sequentially in the trunk and later in the tail. The primordium of the central nervous system becomes visible above the notochord as a thickened neural plate. By the process of epithelial infolding, the neural plate first transforms into a solid neural keel that subsequently hollows into the neural tube (21, 22). The regional morphogenetic processes lead to the formation of distinctive swellings (neuromeres) in the anterior region of the neural plate, corresponding to brain rudiments characteristic of all vertebrates. The most anterior forebrain rudiment consists of telencephalon and ventral diencephalon, from which the eyes develop. The medially located mesencephalon gives rise to the optic tectum dorsally and tegumentum ventrally. The most posterior hindbrain is subdivided into seven rhombomeres. The neural tube posterior to hindbrain gives rise to spinal cord, which develops characteristic dorso-ventral organization, dependent on secreted proteins of the Hedgehog family (23).

During the segmentation period, while the embryo continues to narrow in the mediolateral direction and elongates, the rudiments of kidney, heart, and gut, as well as sensory placodes, also become distinct. The tail, originally closely apposed to the round yolk cell, becomes a prominent structure when it everts from the yolk. At 24 hours, the embryo reaches the pharyngula stage of development. At this phylotypic stage, the zebrafish embryo is most similar to other vertebrate embryos, exhibiting a prominent notochord in the trunk and tail, surrounded by chevron-shaped somites (Fig. 1). The nervous system positioned dorsal to the notochord is hollow and expanded anteriorly, with eyes and brain subdivisions visible in the head. The embryo exhibits circulation and is motile, and rudiments of most organs can be easily identified. The heart tube is located against the ventral body wall below the head and consists of an inner endocardial tube and an outer, myocardial tube. The adult fish heart will consist of a single atrium and a single ventricle, an arrangement characteristic of the early embryonic stages of other vertebrates. Blood cells develop in the intermediate cell mass located in the trunk between the notochord and endoderm, accumulating in a blood island at the trunk/tail border (24). Further development involves the formation of gill arches, including jaws and an outgrowth of pectoral fins. Hatching from the chorion occurs at 2–3 days of development, the swim bladder becomes inflated, and larvae start feeding on day 4 of development.

**Figure 1.** Zebrafish embryo at day 1 of development (pharyngula stage). Scale bar, 0.5 mm.



### 3. Genetics of Zebrafish Development

Efficient methods of induction, recovery and screening for mutations are available in zebrafish. *N-*

ethyl- *N*-nitrosourea is a potent [mutagen](#) of fish spermatogonia and sperm ([25-28](#)). g-ray-irradiation is used to induce mutations in embryos or in sperm of adult males ([29, 30](#)). Furthermore, insertional mutagenesis can be achieved by injecting a pseudotyped retroviral **vector** containing a genome based on the Moloney murine leukemia virus and the envelope **glycoprotein** of the vesicular stomatitis virus into zebrafish embryos ([31, 32](#)). Mutations induced by various mutagens in the G0 generation are transmitted to the F1 generation and are revealed by screening haploid progeny of F1 females (obtained by fertilization with genetically compromised sperm) or gynogenetic diploid progeny of F1 females (obtained by fertilization with genetically compromised sperm followed by inhibition of the second meiotic division). Alternatively, mutations are bred to homozygosity in a classic three-generation screen to be manifest in F3 embryos ([33](#)).

Large-scale genetic screens for mutations that alter embryonic morphology, behavior, or gene expression patterns resulted in the identification of mutations in over 400 genes ([34-36](#)). These screens identified genes involved in practically every aspect of early embryonic and larval development in the zebrafish: early pattern formation, notochord, brain, spinal cord, somites, muscles, heart, circulation, blood, skin, fin, eye, otic vesicle, jaw, and branchial arches ([37](#)). It is estimated that fewer than 50% of genes that could be detected by the screening methods employed have been identified so far. Comparison of the frequency of mutations in known loci with that of embryonic and early larval lethal mutations led to an estimate that 1,500–5,000 genes are essential for early development in zebrafish.

Genes identified by virtue of retroviral insertions are rapidly cloned using the retroviral tag ([38](#)). Genes identified by other types of mutagens are cloned by two main strategies. In the candidate gene approach, genes are identified that have been cloned in other systems and exhibit properties predicted for the mutant locus ([39, 40](#)). The positional cloning approach starts with the identification via genetic mapping of a DNA marker closely linked to a gene to be cloned. Next, this DNA marker is used to isolate chromosomal walking clones encompassing the gene, which is subsequently identified by **gene expression** patterns and [sequence analysis](#), as well as by functional assays ([41](#)).

During the 1990s, zebrafish researchers recovered thousands of mutations that identify hundreds of genes essential for embryonic development. Furthermore, a genetic map of the zebrafish genome has been constructed, and its density is increasing rapidly. Considering the powerful tools available for the analysis of normal and mutant embryos, zebrafish will significantly contribute to our understanding of molecular mechanisms governing early vertebrate development.

## Bibliography

1. G. Streisinger, C. Walker, N. Dower, D. Knauber, and F. Singer (1981) *Nature* **291**, 293–296.
2. C. B. Kimmel (1989) *Trends Genetics* **5**, 283–288.
3. G. Dave and R. Q. Xiu (1991) *Archives of Environmental Contamination & Toxicology* **21**, 126–134.
4. M. Westerfield (1996) *The Zebrafish Book*, University Oregon Press, Eugene, Ore.
5. R. Hinegardner and D. E. Rosen (1972) *American Naturalist* **166**, 621.
6. A. Endo and T. H. Ingalls (1968) *Journal of Heredity* **59**, 382–384.
7. J. H. Postlethwait, S. L. Johnson, C. N. Midson, W. S. Talbot, M. Gates, E. W. Ballinger, D. Africa, R. Andrews, T. Carl, J. S. Eisen, et al (1994) *Science* **264**, 699–703.
8. E. Knapik, A. Goodman, O. S. Atkinson, C. T. Roberts, M. Shiozawa, C. U. Sim, S. Weksler-Zangen, M. R. Trolliet, C. Futrell, B. A. Innes, G. Koike, M. G. McLaughlin, L. Pierre, J. S. Simon, E. Vilalonga, M. Roy, P.-W. Chiang, M. C. Fishman, W. Driever, and H. J. Jacob (1996) *Development* **123**, 451–460.
9. J. H. Postlethwait and W. S. Talbot (1997) *Trends in Genetics* **13**, 183–190.
10. S. Abdelilah, L. Solnica-Krezel, D. Y. Stainier, and W. Driever (1994) *Nature* **370**, 468–471.
11. K. A. Helde, E. T. Wilson, C. J. Cretekos, and D. J. Grunwald (1994) *Science* **265**, 517–520.

12. S. Jesuthasan and U. Strahle (1996) *Current Biology* **7**, 31–42.
13. S. Schneider, H. Steinbesser, R. M. Warga, and P. Hausen (1996) *Mechanisms of Development* **57**, 191–198.
14. G. M. Kelly, D. F. Erezyilmaz, and R. T. Moon (1995) *Mechanisms of Development* **53**, 261–273.
15. R. M. Warga and C. B. Kimmel (1990) *Development* **108**, 569–580.
16. C. B. Kimmel, R. M. Warga, and T. F. Schilling (1990) *Development* **108**, 581–594.
17. J. Shih and S. E. Fraser (1996) *Development* **122**, 1313–1322.
18. E. M. De Robertis, A. Fainsod, L. K. Gont, and H. Steinbeisser (1994) *Development (Supplement)*, 117–124.
19. M. Hammerschmidt, G. N. Serbedzija, and A. McMahon (1996) *Genes and Development* **10**, 2452–2461.
20. S. Fisher, S. L. Amacher, and M. E. Halpern (1997) *Development* **124**, 1301–1311.
21. B. Schmitz, C. Papan, and J. A. Campos-Ortega (1993) *Roux's Arch. of Dev. Biol.* **202**, 250–259.
22. C. Papan and J. A. Camposortega (1994) *Roux's Arch. Devel. Biol.* **203**, 178–186.
23. M. Hammerschmidt, A. Brook, and A. P. McMahon (1997) *Trends in Genetics* **13**, 14–20.
24. H. W. Dietrich, M. W. Kieran, F. Y. Chan, L. M. Barone, K. Yee, J. A. Rundstadler, and L. I. Zon (1995) *Proc. Natl. Acad. Sci. USA* **92**, 37–46.
25. D. J. Grunwald and G. Streisinger (1992) *Genetical Research* **59**, 103–116.
26. M. C. Mullins, M. Hammerschmidt, P. Haffter, and C. Nusslein-Volhard (1994) *Current Biology* **4**, 189–202.
27. L. Solnica-Krezel, A. F. Schier, and W. Driever (1994) *Genetics* **136**, 1401–1420.
28. B. R. Riley and D. L. Grunwald (1995) *Proc. Natl. Acad. Sci. USA* **92**, 5997–6001.
29. C. Walker and G. Streisinger (1983) *Genetics* **103**, 125–136.
30. A. Fritz, M. Rozowski, C. Walker, and M. Westerfield (1996) *Genetics* **144**, 1735–1745.
31. J. C. Burns, T. Friedmann, W. Driever, M. Burrascano, and J. K. Yee (1993) *Proc. Nat. Acad. Sci. USA* **90**, 8033–8037.
32. N. Gaiano, A. Amsterdam, K. Kawakmi, M. Allende, T. Becker, and N. Hopkins (1996) *Nature* **382**, 829–832.
33. W. Driever, D. Stemple, A. Schier, and L. Solnica-Krezel (1994) *Trends in Genetics* **10**, 152–159.
34. W. Driever, L. Solnica-Krezel, A. F. Schier, S. C. F. Neuhaus, J. Malicki, D. L. Stemple, D. Y. R. Stainier, F. Zwartkruis, S. Abdelilah, Z. Rangini, E. Mountcastle-Shah, M. Harvey, J. Belak, and C. Boggs (1996) *Development* **123**, 37–46.
35. P. Hafter, M. Granato, M. Brand, M. C. Mullins, M. Hammerschmidt, D. A. Kane, J. Odenthal, F. J. M. van Eeden, Y.-J. Jiang, C. P. Heisenberg, R. N. Kelsh, M. Furutani-Seiki, E. Vogelsang, U. Beuschle, U. Schach, C. Fabian, and C. Nusslein-Volhard (1996) *Development* **123**, 1–36.
36. P. D. Henion, D. W. Raible, C. E. Beattie, K. L. Stoesser, J. A. Weston, and J. S. Eisen (1996) *Developmental Genetics* **18**, 11–17.
37. (1996) *Zebrafish issue*, The Company of Biologists Limited, Cambridge.
38. M. L. Allende, A. Amsterdam, T. Becker, K. Kawakami, N. Gaiano, and N. Hopkins (1996) *Genes and Development* **10**, 3141–3155.
39. S. Schulte-Merker, F. J. M. VanEeden, M. E. Halpern, C. B. Kimmel, and C. Nusslein-Volhard (1994) *Development* **120**, 1009–1015.
40. W. S. Talbot, W. Trevarrow, M. E. Halpern, A. E. Melby, G. Farr, J. H. Postlethwait, T. Jowett, C. B. Kimmel, and D. Kimelman (1995) *Nature* **378**, 150–157.
41. J. Zhang, W. S. Talbot, and A. F. Schier (1998) Positional cloning identifies zebrafish *one-eyed*

*pinhead* as a permissive EFG-related ligand required during gastrulation . Cell **92**, 241–251.

### Suggestions for Further Reading

42. M. Granato and C. Nusslein-Volhard (1996) Fishing for genes controlling development. Curr. Opinion Genetics Develop. **6**, 461–468.
43. C. B. Kimmel, W. W. Ballard, S. R. Kimmel, B. Ullmann, and T. F. Schilling (1995) Stages of embryonic development of the zebrafish. Developmental Dynamics **203**, 253–310.
44. A. F. Schier (1997) Genetics of neural development in zebrafish. Curr. Opinion Neurobiol. **7**, 119–126.
45. L. Solnica-Krezel, D. L. Stemple, and W. Driever (1995). Transparent things: Cell fates and cell movements during early embryogenesis of zebrafish. BioEssays **17**, 931–939.
46. D. Y. R. Stainier and M. C. Fishman (1994) The zebrafish as a model system to study cardiovascular development. Trends Cardiovasc. Med. **4**, 207–212.
47. <http://zfish.uoregon.edu/>

### Zimm–Bragg Model

The Zimm–Bragg model is a statistical mechanical model that describes the  $\alpha$ -helix to random coil transition in peptides that do not adopt any other folded conformation (1) (see [Alpha-Helix Formation](#)). This and the other popular model for the helix–coil transition, the [Lifson–Roig model](#) (2), are compared and contrasted in [Helix–coil theory](#). In the Zimm–Bragg description of helix formation in a peptide, there are three key parameters:  $n$ , the number of peptide units in the chain;  $s$ , the helix propagation parameter; and  $\sigma$  the helix nucleation parameter. The Zimm–Bragg model defines  $n$  as the number of peptide units ( $\text{—CONH—}$ ) in the chain, which causes some problems in ascribing  $s$  values to particular side chains, because a peptide unit encompasses parts of two amino acid residues. The alternative model, the Lifson–Roig model, does not use this definition and is therefore sometimes more convenient.

There are several implementations of basic Zimm–Bragg theory that can be used to describe helix formation in peptides. In the full treatment, a  $4 \times 4$ -correlation matrix is needed to assign statistical weights to helical or coil residues (see [Helix–Coil Theory](#) for a diagram) in order to calculate the partition function for the peptide system. With this full treatment, residue-specific  $s$  and  $\sigma$  values can be used. A somewhat simplified version of the Zimm–Bragg model employs a  $2 \times 2$ -correlation matrix, where nearest-neighbor interactions can be included. Finally, a useful version of the Zimm–Bragg model for short homopolymers is called the *Zipper Model*. In this approximation, the peptide is treated as a homopolymer, and the single-sequence approximation is used, in which it is assumed that only one stretch of helical conformation exists at one time for a given peptide chain. In other words, because helix nucleation is an unfavorable event, it can be treated as occurring only one time for a given peptide, and the propagation of the helix “zips” up the chain.

This model does not employ matrix methods to calculate the partition function, and the fraction of helical residues can be expressed as

$$f_H = \frac{\sigma s}{(s-1)^3} \left( \frac{ns^{n+2} - (n+2)s^{n+1} + (n+2)s - n}{n\{1 + [\sigma s / (s-1)^2][s^{n+1} + n - (n+1)s]\}} \right) \quad (1)$$

This version of the model is especially useful for investigating the cooperative length-dependence of helix formation by simple repeating homopolymers (3).

### Bibliography

1. B. H. Zimm and J. K. Bragg (1959) *J. Chem. Phys.* **31**, 526–535.
2. S. Lifson and A. Roig (1961) *J. Chem. Phys.* **34**, 1963–1974.
3. J. M. Scholtz, H. Qian, E. J. York, J. M. Stewart, and R. L. Baldwin (1991) *Biopolymers* **31**, 1463–1470.

### Suggestions for Further Reading

4. C. R. Cantor and P. R. Schimmel (1980) *Biophysical Chemistry*, W. H. Freeman, San Francisco. Chapter 20 provides a good introduction to Zimm–Bragg models for helix formation.
5. D. Poland and H. A. Scheraga (1970) *Theory of Helix–Coil Transitions in Biopolymers*, Academic Press, New York. A broad discussion of helix–coil theory. The book also reprints copies of many of the seminal articles on helix-coil theory.
6. H. Qian and J. A. Schellman (1992) Helix–coil theories: a comparative study for finite length polypeptides. *J. Phys. Chem.* **96**, 3987–3994. A comparison of the Zimm–Bragg and Lifson–Roig theories for helix formation in peptides.

## Zinc-Binding Proteins

The 1985 report by Klug and co-workers (1) of a model for the *Xenopus laevis* [transcription factor](#) TFIIIA, which contains a zinc atom coordinated by two [cysteine](#) and two [histidine](#) residues, formally launched the study of metalloregulatory proteins (see [Metalloproteins](#)). A subsequent Zn-EXAFS study confirmed the presence of two N and two S atoms in the inner coordination sphere (2). And complete structures were subsequently determined by [NMR](#). The function of the  $Zn^{2+}$  ion is to organize the [protein structure](#) into one that contains an [alpha-helix](#) suitable for recognizing the major groove of **DNA**. This protein displays the hallmark of a large family of proteins that display the zinc-finger motif and can be identified from the amino acid sequence -Cys- $X_{(2-4)}$ -Cys- $X_{12}$ -His- $X_{(3-5)}$ -His-, where X represents any amino acid. This pattern of cysteine and histidine residues has been used repeatedly to infer the presence of a zinc-finger protein or **domain** from just the amino acid sequence. This protein motif is especially common in eukaryotes, including humans. Some zinc fingers also bind **RNA**, however, complicating the elucidation of a zinc-finger “recognition code.”

Other zinc-binding motifs have been discovered, as shown in Table 1. Zinc-activated transcription factors also include the nuclear hormone receptors, which bind cooperatively to DNA as homodimers (see **Steroid receptors**). The  $Zn^{2+}$  sites in these proteins possess tetrahedral environments of four cysteine residues. Fungal transcription factors include members of the GAL4 family, which also bind as homodimers to DNA. However, these contain a  $(Cys)_3$ -Zn- $(Cys)_2$ -Zn- $(Cys)_3$  dimer that has two cysteine residues bridging the two zinc centers. Structures of the human [estrogen receptor](#) and of GAL4 both show that DNA recognition is also via an  $\alpha$ -helix.

Nucleocapsids of retroviruses (e.g., **HIV-1**) contain large amounts of  $Zn^{2+}$ , probably ligated by three cysteine residues and one histidine (3). The functions of the zinc-binding domains probably involve the recognizing RNA packaging signals. Many other classes of zinc-binding proteins, including those that possess RING and LIM domains, may also function in nucleic acid regulatory capacities.



**Table 1. Major Families of Zinc Metalloregulatory Proteins**

| Family                           | Consensus Sequence <sup>a</sup>  | Recognition Site |
|----------------------------------|--|------------------|
| Zinc finger                      | CX <sub>(2-4)</sub> CX <sub>12</sub> HX <sub>(3-5)</sub> H   | DNA, RNA         |
| Nuclear hormone receptor<br>GAL4 | CX <sub>2</sub> CX <sub>13</sub> CX <sub>2</sub> CX <sub>15</sub> CX <sub>5</sub> CX <sub>12</sub> CX <sub>4</sub> C | DNA              |
| Retroviral nucleocapsid          | CX <sub>2</sub> CX <sub>9-27</sub> CX <sub>2</sub> CX <sub>6</sub> C   | DNA              |
|                                  | CX <sub>2</sub> CX <sub>4</sub> HX <sub>4</sub> C  | RNA <sup>x</sup> |

<sup>a</sup> C, cysteine residue; H, histidine; X, any residue

Why zinc? After iron, zinc is the most plentiful of the trace-metal ions in humans. More importantly, zinc is soluble, exchanges ligands rapidly, does not undergo biological redox reactions, and tolerates a variety of coordination geometries. Sulfur-containing ligands display a marked preference for “soft” metal ions, including Zn<sup>2+</sup>. Hence, it is not surprising that cysteine residues play such a prominent biological role in binding zinc.

#### Bibliography

1. J. Miller, A. D. McLachlan, and A. Klug (1985) *EMBO J.* **4**, 1609–1614.
2. G.P. Diakun, L. Fairall, and A. Klug (1986) *Nature* **324**, 698–699.
3. M.F. Summers, L.E. Henderson, M.R. Chance, J.W. Bess Jr., T.L. South, P.R. Blake, I. Sagi, G. Alvarado Peret, R.C. Sowder, III, D.R. Hare, and L.O. Arthur (1992) *Protein Sci.* **1**, 563–574.

#### Suggestions for Further Reading

4. H. Sigel and A. Sigel (series eds.) (1974–1998) *Metal Ions in Biological Systems*, Dekker, New York, Vols. 1–35.
5. R. J. Cousins (1994) Metal elements and gene expression, *Annu. Rev. Nutr.* **14**, 449–469.

#### Zinc-Containing DNA-Binding Motifs

A number of different DNA-binding motifs contain one or more zinc atoms whose function is to stabilize the modular structure of the domain by coordinating with amino acids, usually [cysteine](#) or histidine, in the appropriate spatial orientation. The only common feature of these domains, of which six have been identified to date, is the presence of zinc. Otherwise they are structurally diverse.

The most ubiquitous of these structures is the “zinc finger” found in many eukaryotic transcription factors. The structure of a typical zinc finger consists of a well-defined  $\alpha$ -helix packed against two  $\beta$ -strands arranged in a hairpin structure. The finger structure is stabilized by the tetrahedral

coordination of zinc to two closely spaced cysteine residues in the b-strands and two closely spaced histidine residues at the C-terminus of the helical region. In addition the structure is further stabilized by a hydrophobic pocket in the interior of the structure.

When bound to DNA, the finger is directed into the major groove such that base-specific contacts are made by the proximal end of the a-helical region with contacting successive amino acids being separated by approximately one helical turn. The sugar-phosphate backbone is contacted by one of the zinc-coordinated histidine residues. A single zinc finger generally contacts three successive base-pairs in DNA, usually in a single-strand of the duplex. However many zinc-finger proteins contain multiple contiguous fingers separated by a short linker. In such examples synergy can occur between adjacent zinc fingers so that each specify four bp subsites overlapping by one base-pair. In proteins with multiple DNA-binding zinc fingers, the structure of the first finger often differs from the canonical structure. For example, the first finger of the SWI5 protein is preceded by a short a-helix and contains an additional [b-strand](#) preceding the normal b hairpin loop. In proteins with multiple zinc fingers, some of the fingers may not be involved in DNA binding but instead can act to mediate protein-protein contacts.

Another DNA-binding motif which contains zinc is exemplified by DNA-binding domain of the steroid receptors. In this fold zinc is coordinated by four cysteine residues. In this case, the path of the polypeptide backbone between the proximal and distal pairs of coordinating cysteines is distinct from that of the TFIIIA type of finger. In the receptor motif, the proximal cysteine residues, as in zinc fingers, reside in the b-sheet structures. By contrast, although the histidine residues involved in zinc coordination are often contained within an a-helix, the corresponding distal cysteine are not. This difference is directly related to DNA binding. The DNA binding domains of the glucocorticoid and oestrogen receptors contain two C-C...C-C coordination domains, the first of which is of similar extent to the TFIIIA fingers whereas the second is significantly shorter. Unlike the TFIIIA fingers both pairs of distal cysteines are immediately followed by an a-helical region. The second a-helix forms direct hydrophobic contacts with the first so that the two fingers and the two a-helices together form a single structural domain, the 'double-loop-zinc-helix'. This domain also includes a region that specifies the dimerization contacts. The determinants for sequence specific recognition lie within the a-helix immediately distal to the first cysteine tetrad. For each receptor the DNA recognition domains bind as dimers to a palindromic sequence containing three conserved bases on each side. But the separation between the conserved trimers differs for different receptors. This means that the rotational orientation of the individual components of the dimer must also be different and be determined by the position of the dimerization contacts relative to the a-helix responsible for sequence recognition.

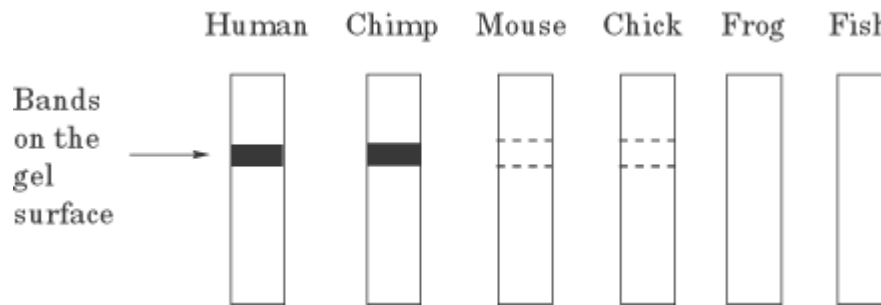
In a third type of zinc-containing DNA binding domain, typified by the GAL4 protein, the zinc-stabilized structure is connected by a flexible region of polypeptide chain to a short amphipathic a-helix which forms hydrophobic contacts with another monomer in dimer formation. In GAL4, as with other zinc-containing motifs, direct contact with the bases is mediated in the major groove by amino acids contained within a short a-helical region. However, most of the contacts to the sugar-phosphate backbone are effected by amino acids in the more extended chain. The GAL4 protein, like the steroid hormone receptors, binds to its DNA recognition site as a dimer with the two monomeric units being held together by a short coiled-coil of two parallel a-helices.

## Zoo Blot

A zoo blot (Fig. 1) is the result of an experiment in which a short DNA fragment from the genomic region of interest is first **cloned** and then used as a radioactive probe to test for related DNA from a

variety of species by [Southern blot](#) analysis. If hybridized fragments related to the probe are found in several species, there is a strong possibility that the probe is an exon of a gene, because of its sequence conservation (1). Such candidates can be further characterized by [DNA sequencing](#) to determine whether they contain open reading frames (ORFs). If so, the identification of the new protein-coding gene is possible by further experiments, such as isolation of [complementary DNA](#) or [messenger RNA](#). Thus, a zoo blot is a powerful tool, not only to see how the probe of interest is conserved among a variety of species, but also to identify a new protein-coding gene.

**Figure 1.** Schematic representation of a zoo blot. Southern blot analysis of genomic DNAs from various species was performed using a particular probe, such as a segment of human DNA.



## Bibliography

1. B. Lewin (1995) *Gene V*, Oxford Univ. Press, Oxford.

## Zwitterion (Dipolar Ion)

A zwitterion is a molecule that simultaneously has one positively charged group and one negatively charged group. The best examples are the individual amino acids that have a cationic amine group and an anionic carboxylate. Proteins usually have both cationic and anionic side chains.

## Zygote

The **diploid** zygote is formed by the fusion of two haploid **gametes** at [fertilization](#). After fertilization, the zygote undergoes cleavage divisions to form a mass of cells called the **blastula**.

## Suggestions for Further Reading

- B. I. Balinsky (1975) *An Introduction to Embryology*, 4th ed., W. B. Saunders, Philadelphia, p. 4.
- J. A. Moore (1972) *Heredity and Development*, Oxford University Press, New York.
- S. Shostak (1991) *Embryology: An Introduction to Developmental Biology*, Harper Collins, New

York.

L. Wolpert et al. (1998) *Principles of Development*, Oxford University Press, Oxford, U.K.

## Zymogen

A zymogen is an enzymatically inactive precursor of an [enzyme](#), often but not always a **proteolytic** enzyme (or [proteinase](#)). Some zymogens are named by adding the suffix *-ogen* to the name of the enzyme itself, as in **trypsinogen** or **pepsinogen**, whereas others are indicated by the prefix *pro-*, as in pro-collagenase or pro-**carboxypeptidase**. No strict rules of nomenclature pertain. Hence, pro-renin rather than reninogen, pre-kallikrein rather than pro-kallikrein, **Factor X** rather than proFactor Xa, etc. Also, some proteins are given the suffix *-ogen* even if they are not a proteinase (or even an enzyme) precursor—for example, [fibrinogen](#) or *angiotensinogen*.

Zymogens are one of Nature's ways of ensuring that proteinase (or other) activity will be generated only at the most appropriate time and in the most appropriate place. Most of the digestive proteinases produced in the pancreas are synthesized as zymogens and stored as such in zymogen granules in the acinar cells. The contents of these granules are only released into the small intestine when food is eaten, and the zymogens are only activated once they enter the small intestine. This latter process is initiated by an intestinal proteinase, [enterokinase](#). Inappropriate activation of these zymogens within the pancreas can cause pancreatitis and have devastating consequences.

Zymogen activation typically involves cleavage of a single [peptide bond](#) in the precursor protein. Sometimes, as in the case of pepsinogen, the proteinase responsible for cleaving that peptide bond is the zymogen itself. Briefly, in the acidic environment of the stomach the pepsinogen molecule unfolds enough to expose the [active site](#), which then catalyzes a process of autoactivation by clipping off a peptide segment from the *N*-terminus of other pepsinogen molecules, thereby activating them and setting off a chain reaction. The pancreatic zymogens are activated by **trypsin**, which is itself generated from trypsinogen by enterokinase. Activation of trypsinogen involves removal of a six-residue peptide from its amino terminus, whereas activation of **chymotrypsinogen** involves cleavage of the peptide bond joining Arg15 and Ile16. The first 15 amino acids are linked to the rest of chymotrypsin by a [disulfide bond](#), so there is no actual removal of a peptide. After activation, chymotrypsin undergoes some proteolytic fine-tuning and loses two internal dipeptides, but this is more for stabilization than for activation.

Many other strategies have evolved for zymogen activation, but basically they all come down to inducing a change in the conformation of the precursor by cleaving a peptide bond and thereby altering the [primary structure](#) (amino acid sequence) of the protein; or by interaction with another molecule, a cell surface, or some other type of structure.

Once a zymogen has been activated, Nature must resort to other means to control its activity, and in the case of proteinases there is no shortage of inhibitors (see [Proteinase Inhibitors](#) and [Proteinase Inhibitors, Proteins](#)). Moreover, proteinases being what they are, they can serve as their own substrates and their activities can be abolished by autolysis (self-digestion).